# Where We Are

| | |
|---|---|
| **Machine Learning Systems** | 2012 - Now |
| **Big Data** | 2010 - Now |
| **Cloud** | 2000 - 2016 |
| **Foundations of Data Systems** | 1980 - 2000 |

# Logistics

- Exam date:
  - Final Exam date (tentative): **Friday, March 22, 8 - 11 am, PT**
  - **Decision: In-person Exam**
- Next week:
  - TA will hold multiple hours of Exam review
  - Pay attention to Piazza announcement about scheduling
  - Make sure you attend and get important secret sauces ☺
  - TAs and I will all be available for OH by appointment to help you on exam and wrapping the course!

# ML System history

- ML Systems evolve as more and more ML components (models/optimization algorithms) are unified

Ad-hoc: diverse model family, optimization algos, and data

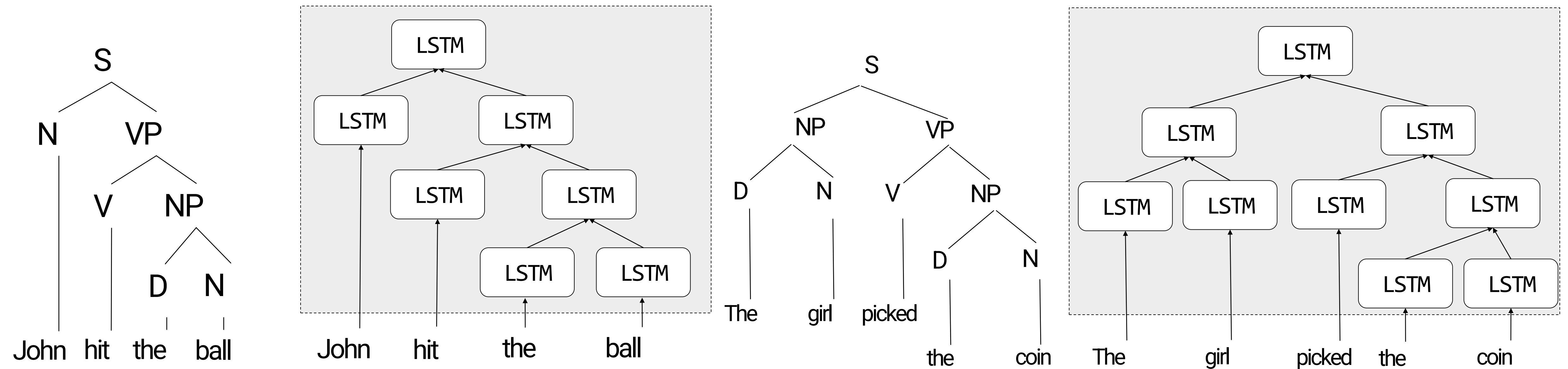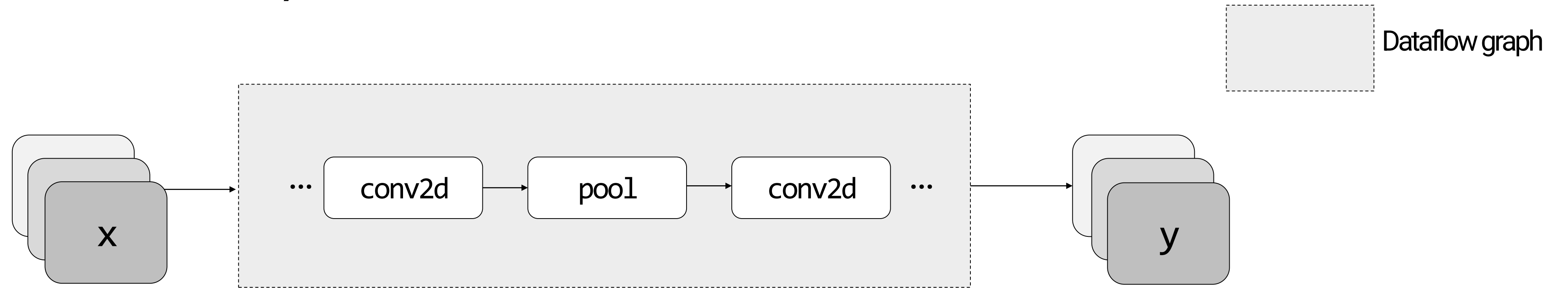Opt algo: iterative-convergent

Model family: neural nets

Model: CNNs/transformers/GNNs

LLMs: transformer decoders

Today: NN, data flow graph, and data parallelism

# Static Models vs. Dynamic Models

# Static vs. Dynamic Dataflow Graphs

- Static Dataflow graphs
  - Define once, execute many times
  - Execution: Once defined, all following computation will **follow** the defined computation
  - Advantages
    - No extra effort for batching optimization, because it can be by nature batched
    - It is always easy to handle a static computational dataflow graphs in all aspects, because of its fixed structure
      - Node placement, distributed runtime, memory management, etc.
  - Benefit the developers

# Static vs. Dynamic Dataflow Graphs

- Can we handle dynamic dataflow graphs?
  - Difficulty in expressing complex flow-control logic
  - Complexity of the computation graph implementation
  - Difficulty in debugging

# How to Handle Dynamic Dataflow Graph?

- In general two ways:

  - Imperative: do not requiring contracting the entire graph before execution

  - Other symbolic representation on top of dataflow graph
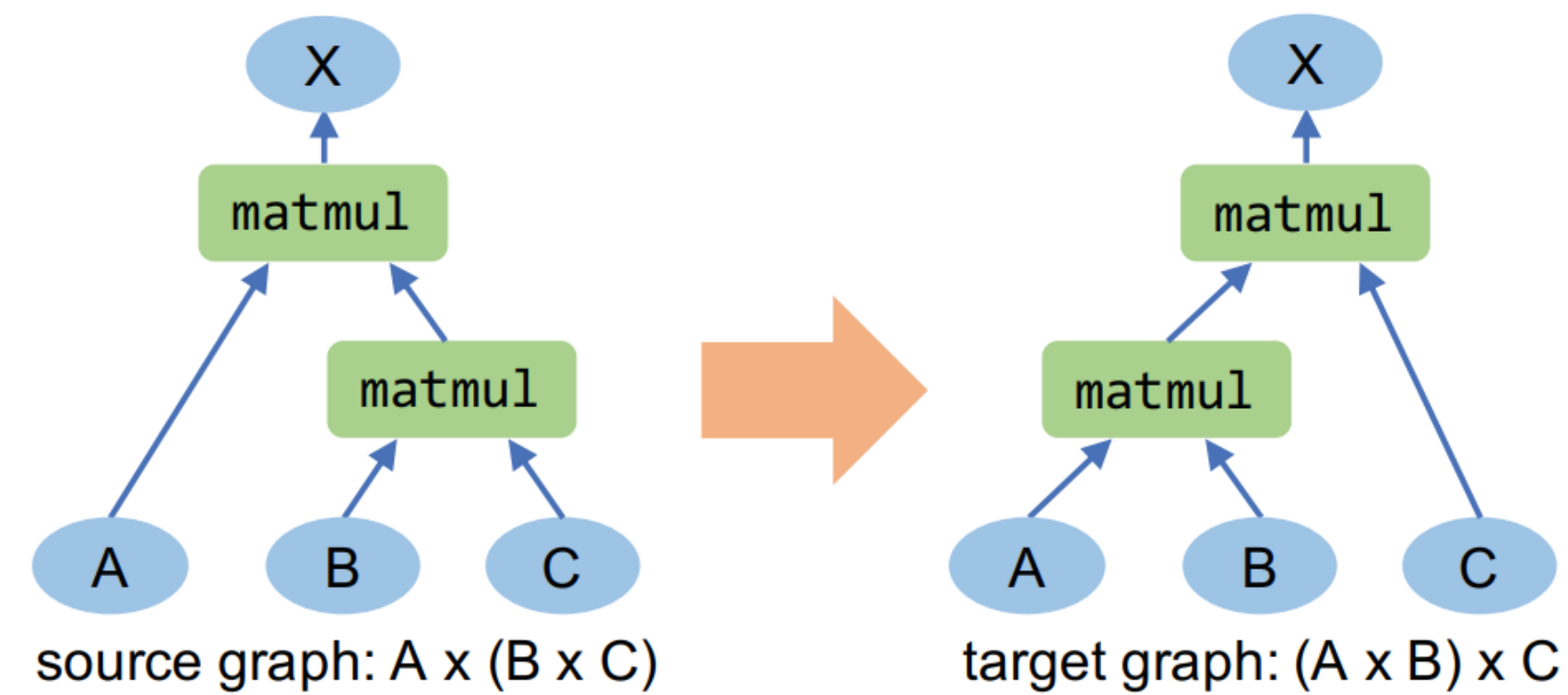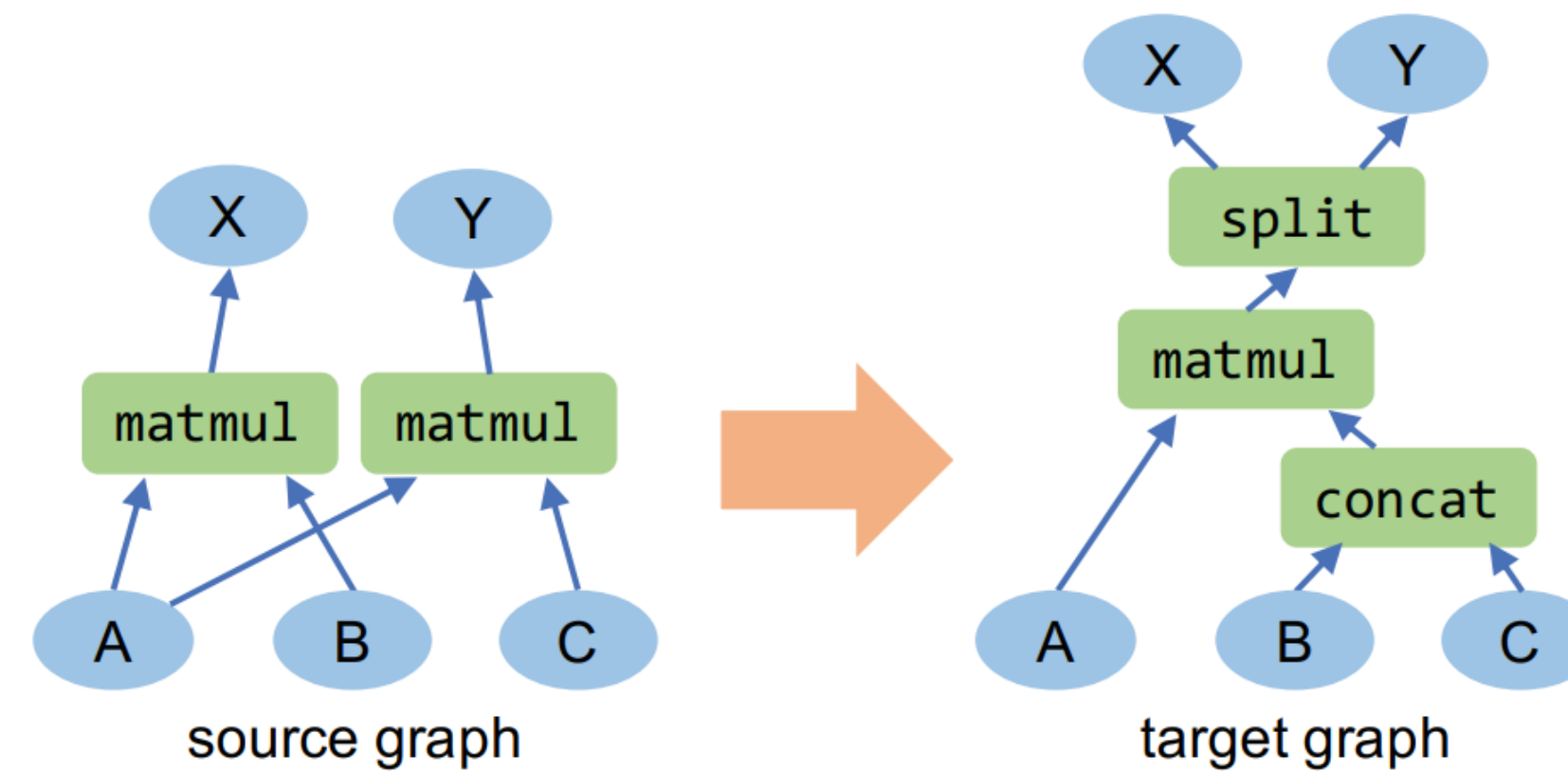
    - vertex-centric representation

Imperative ────────────────────────────►

Imperative

Symbolic

# Questions

- Is **CNN training** static or dynamic graph?

- Is **CNN inference** static or dynamic graph?

- Is **GPT-3 (transformers decoder)** training static graph or dynamic?

- Is **GPT-3 inference with batch size = 1** static or dynamic graph

- Is **GPT-3 serving** static or dynamic graph

# Advanced Topic: DL Dataflow Graph Optimization



**(a)** Associativity of matrix multiplication.



**(b)** Fusing two matrix multiplications using concatenation and split.

# Advanced Topic: DL Graph Compilation
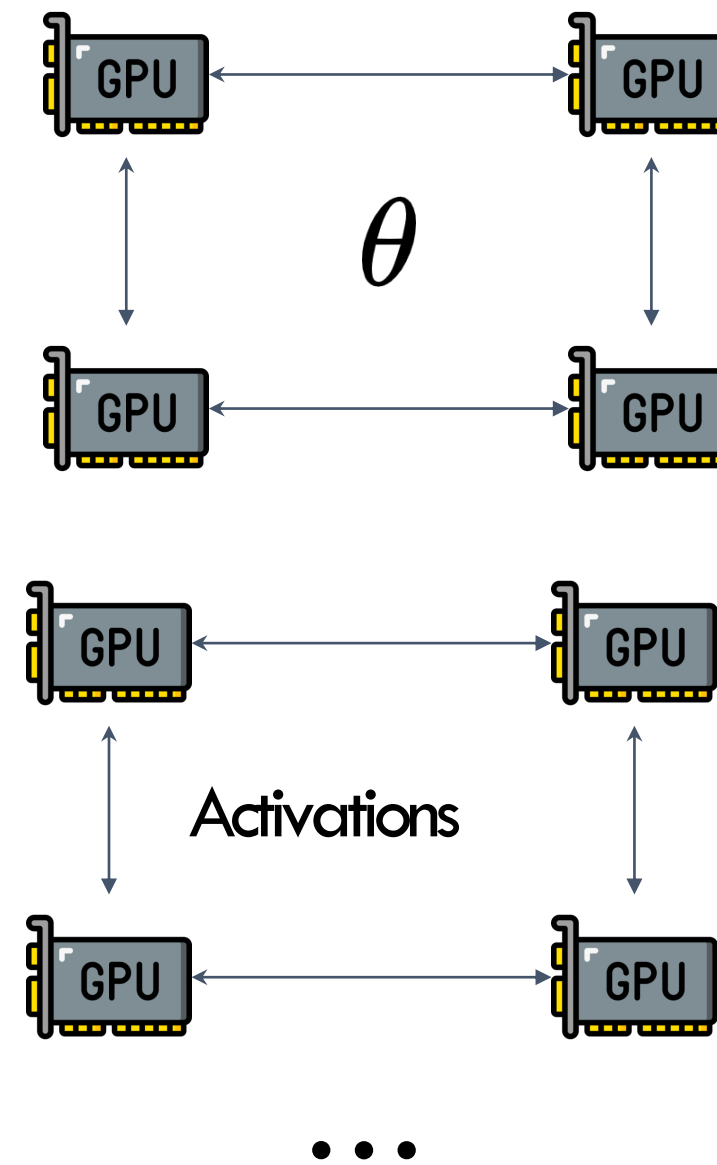
# Where We Are

- Deep Learning as Dataflow Graphs

- Auto-differentiation Libraries

  - Symbolic vs. Imperative

  - Static vs. Dynamic

- **DL Parallelism**

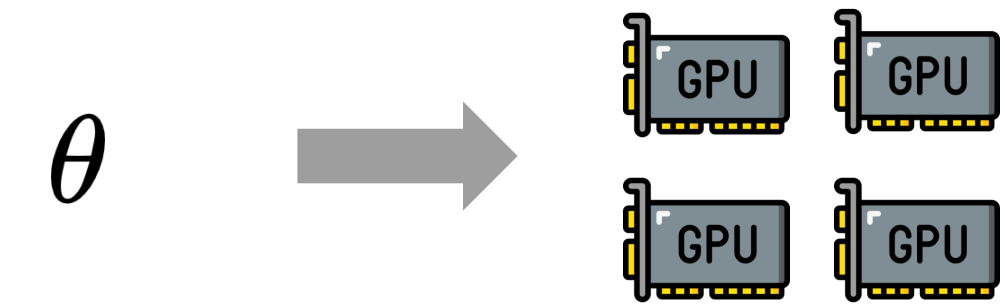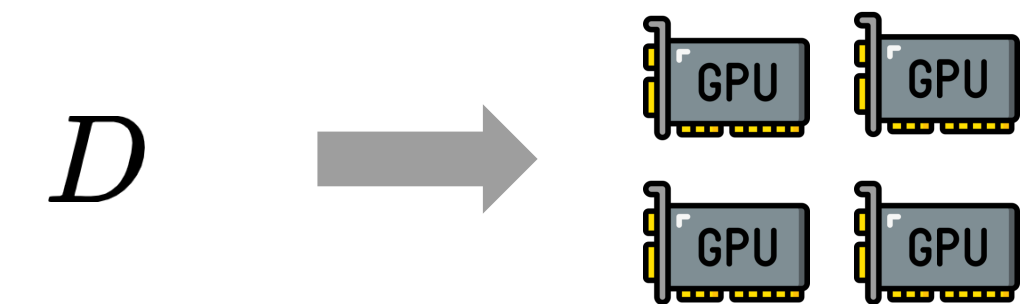# DL Parallelization: 3 Core Problems

**Computing**  **Communication**  **Memory**



$$\nabla_L(\cdot)$$

$$\theta$$

Activations

$$f(\cdot)$$

$$D$$

$$\theta$$

$$\cdots$$

$$\theta^{(t+1)} = f\left(\theta^{(t)}, \nabla_L\left(\theta^{(t)}, D^{(t)}\right)\right)$$

parameter  weight update (sgd, adam, etc.)  model (CNN, GPT, etc.)  data

# Two Views of ML Parallelisms
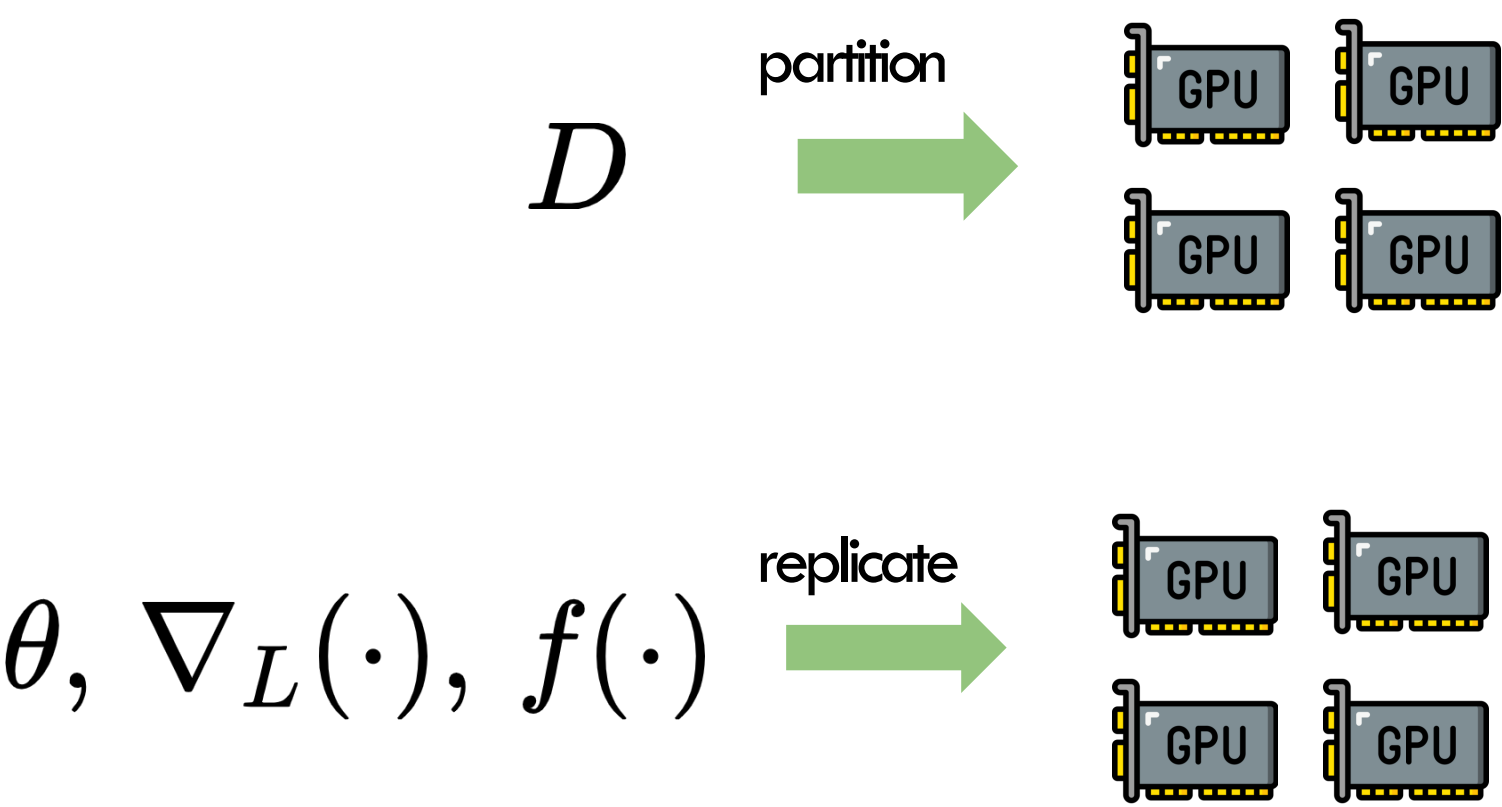
**Classic view**

Data parallelism
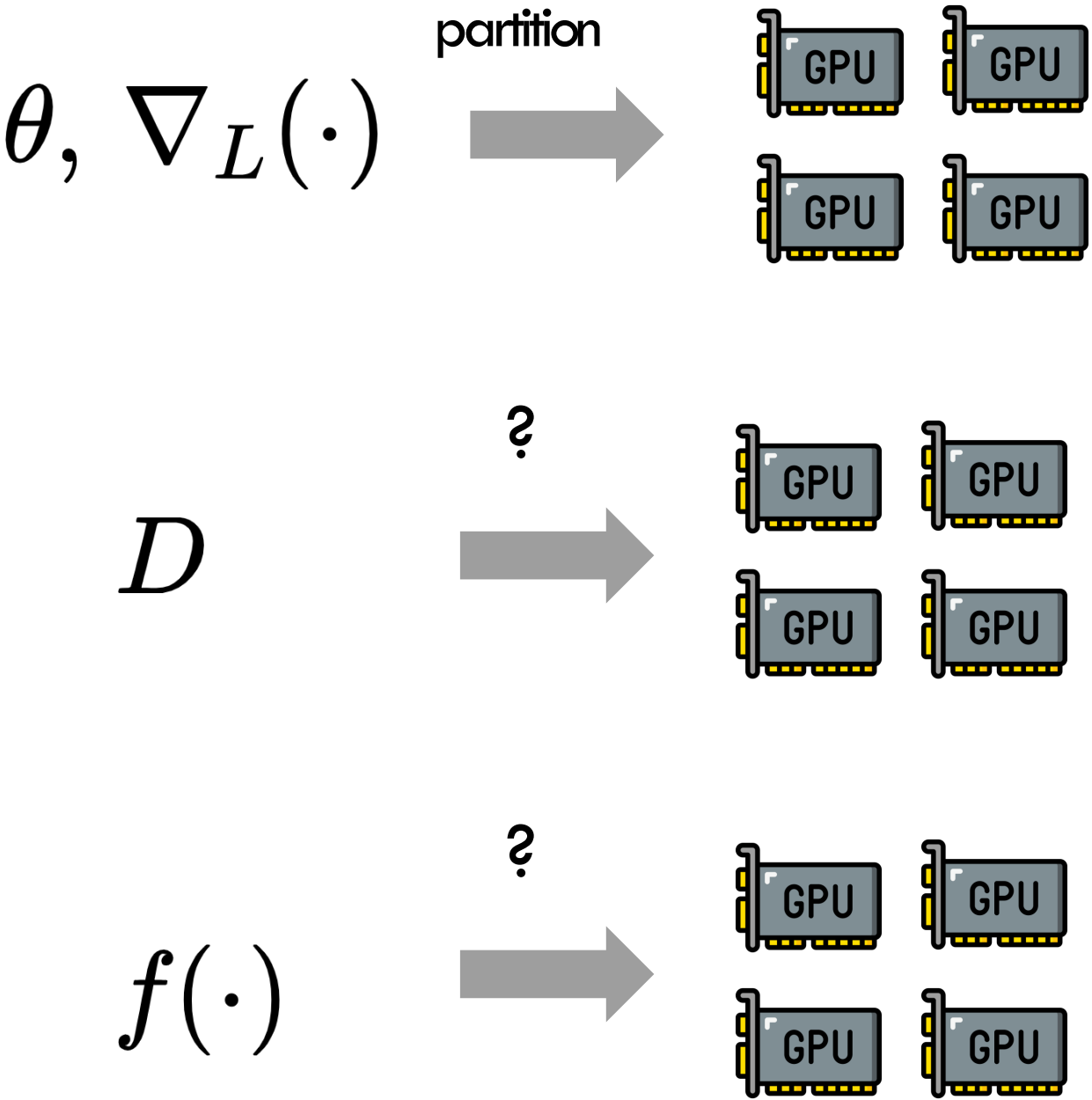
Model parallelism

**New view**

Inter-op parallelism

Intra-op parallelism

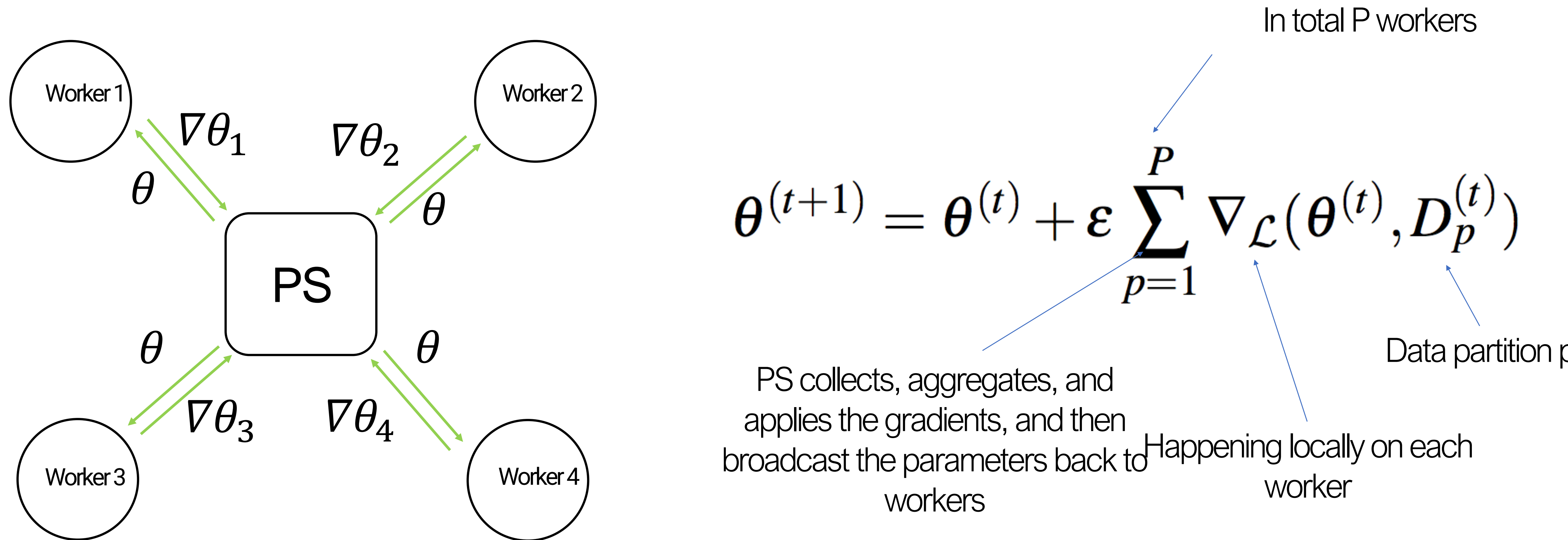# Data and Model Parallelism

**Data parallelism**

**Model parallelism**

$D$    partition  →   [GPU GPU / GPU GPU]

$\theta, \nabla_L(\cdot)$   partition  →   [GPU GPU / GPU GPU]

$\theta, \nabla_L(\cdot), f(\cdot)$   replicate  →   [GPU GPU / GPU GPU]

$D$   ?  →   [GPU GPU / GPU GPU]

$f(\cdot)$   ?  →   [GPU GPU / GPU GPU]

$$\theta^{(t+1)} = f\big(\theta^{(t)}, \nabla_L\big(\theta^{(t)}, D^{(t)}\big)\big)$$

parameter     weight update (sgd, adam, etc.)     model (CNN, GPT, etc.)     data
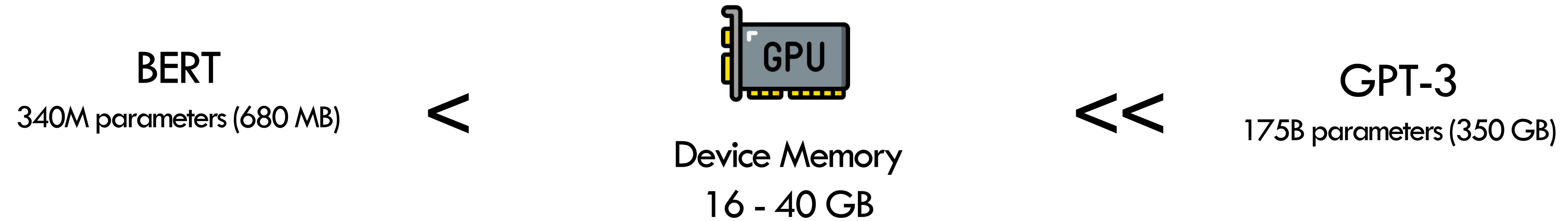
# PS Implements Data Parallelism



In total P workers

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \varepsilon \sum_{p=1}^{P} \nabla_{\mathcal{L}}(\boldsymbol{\theta}^{(t)}, \boldsymbol{D}_p^{(t)})$$

Data partition p

PS collects, aggregates, and applies the gradients, and then broadcast the parameters back to workers

Happening locally on each worker

**Representee Systems: Poseidon, GeePS, BytePS, etc.**

# AllReduce Can Also Handle Data Parallelism Comm



**Representee Systems: Horovod, Torch.DDP**

# Big Model: The Core Computational Challenge

**BERT**

340M parameters (680 MB)

$<$



GPU

Device Memory

16 - 40 GB

$<<$

**GPT-3**

175B parameters (350 GB)

## How to train and serve big models?

## Model Parallelism

# Two Views of ML Parallelisms

Data and model parallelism

- Two pillars: **data** and **model**.

- ✅ "Data parallelism" is general and precise.

- ❓ "Model parallelism" is vague.

- ❓ The view creates ambiguity for methods that neither partitions data nor the model computation.

**New:** Inter-op and Intra-op parallelism.

- Two pillars: **computational graph** and **device cluster**

- ✅ This view is based on their computing characteristics.

- ✅ This view facilitates the development of new parallelism methods.

# Device Cluster
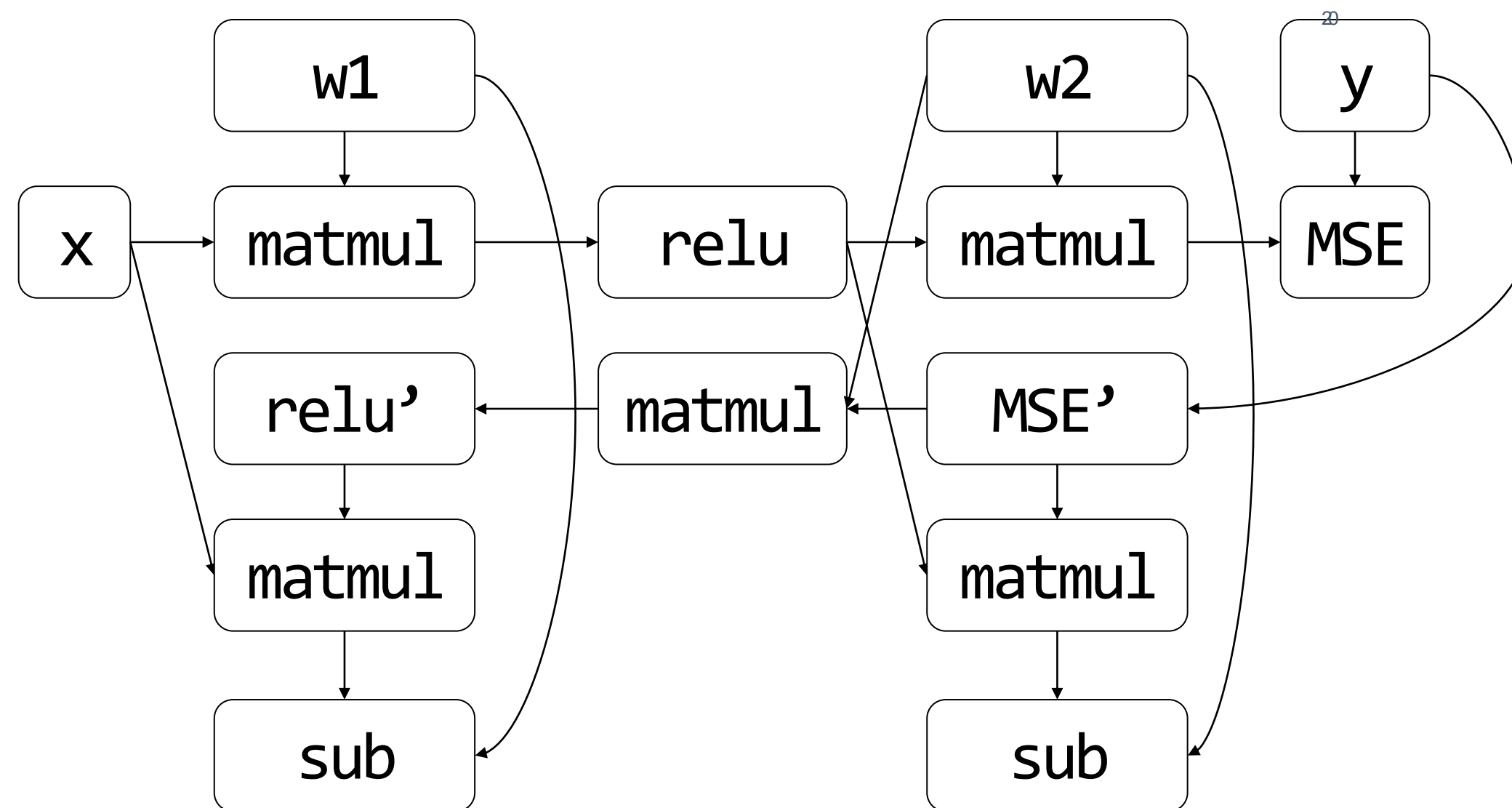
## Nvidia DGX with V100



Figure from NVIDIA

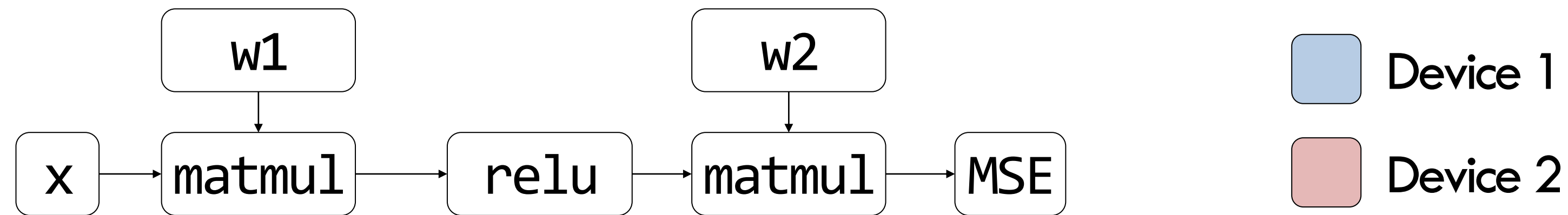## A typical GPU cluster topology
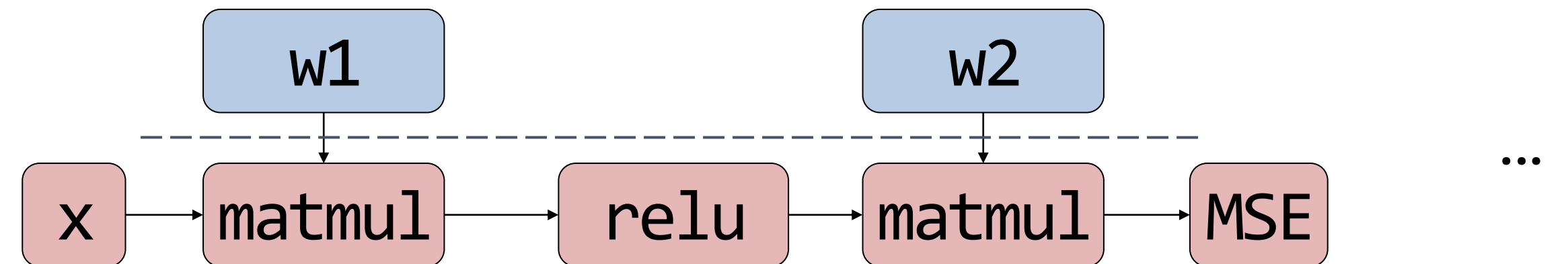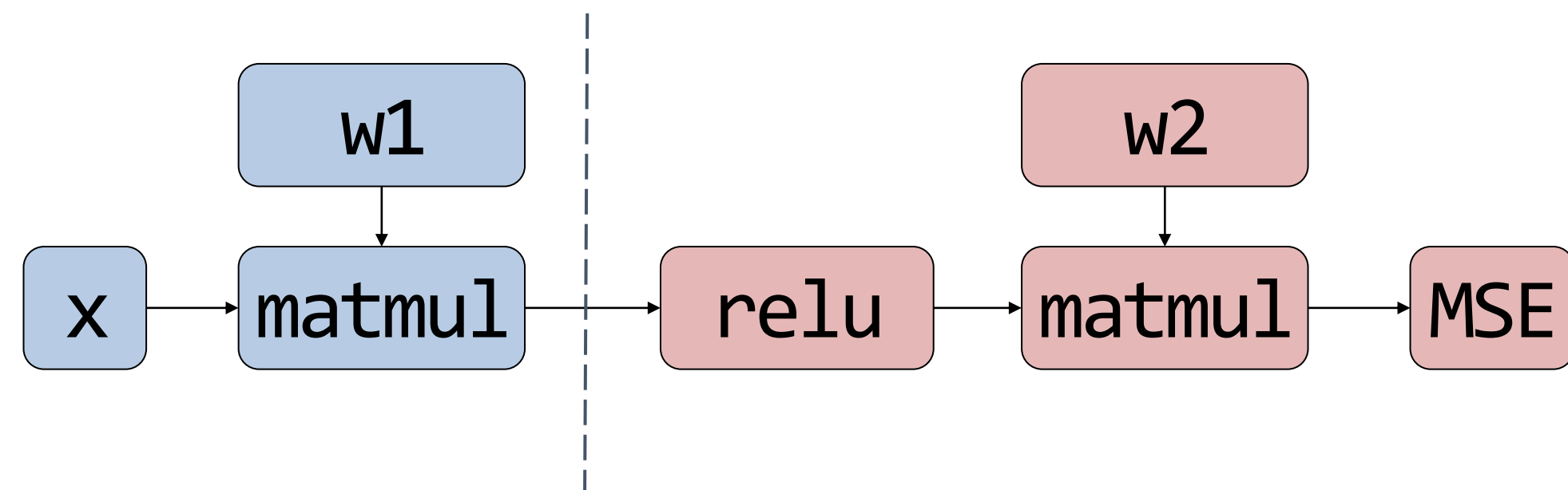
# Partitioning Computation Graph on Device Cluster

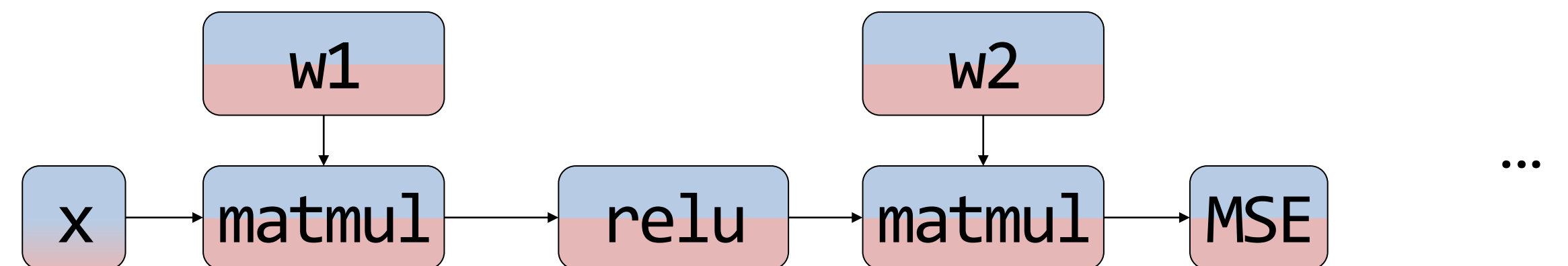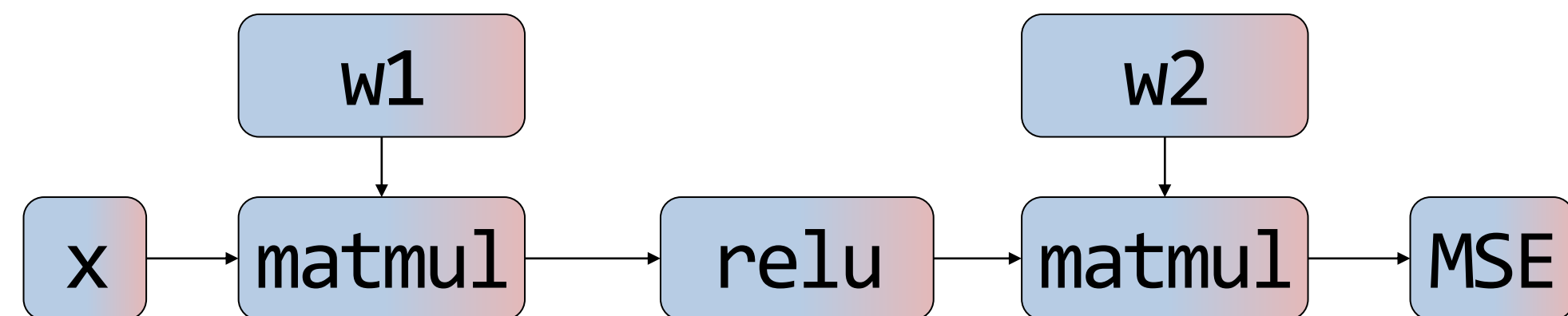How to partition the computational graph
on the device cluster?

# Partitioning Computation Graph

# Partitioning Computation Graph



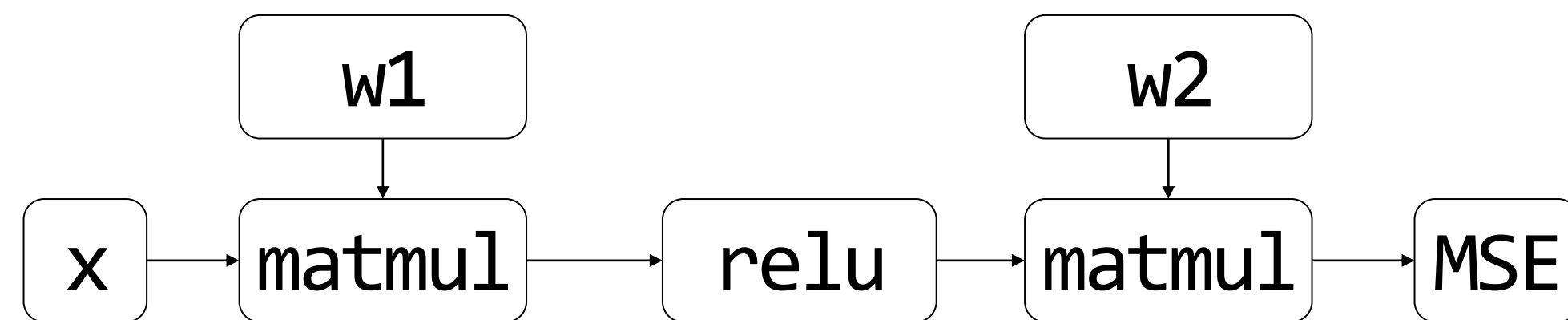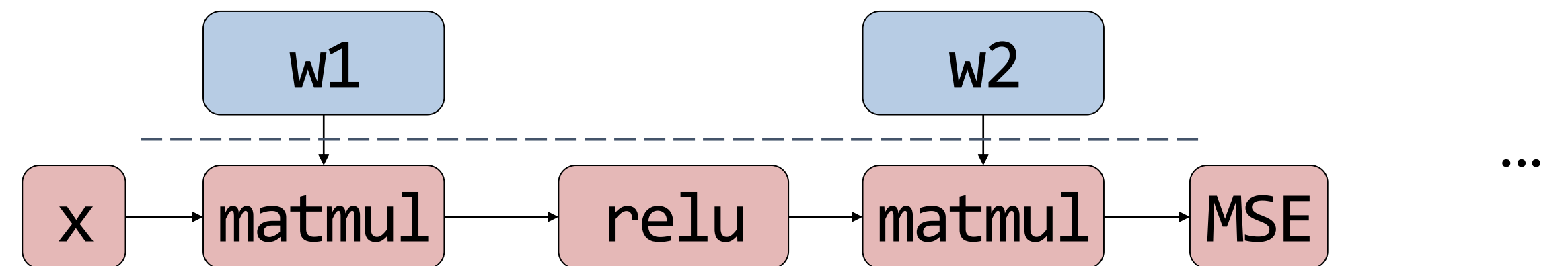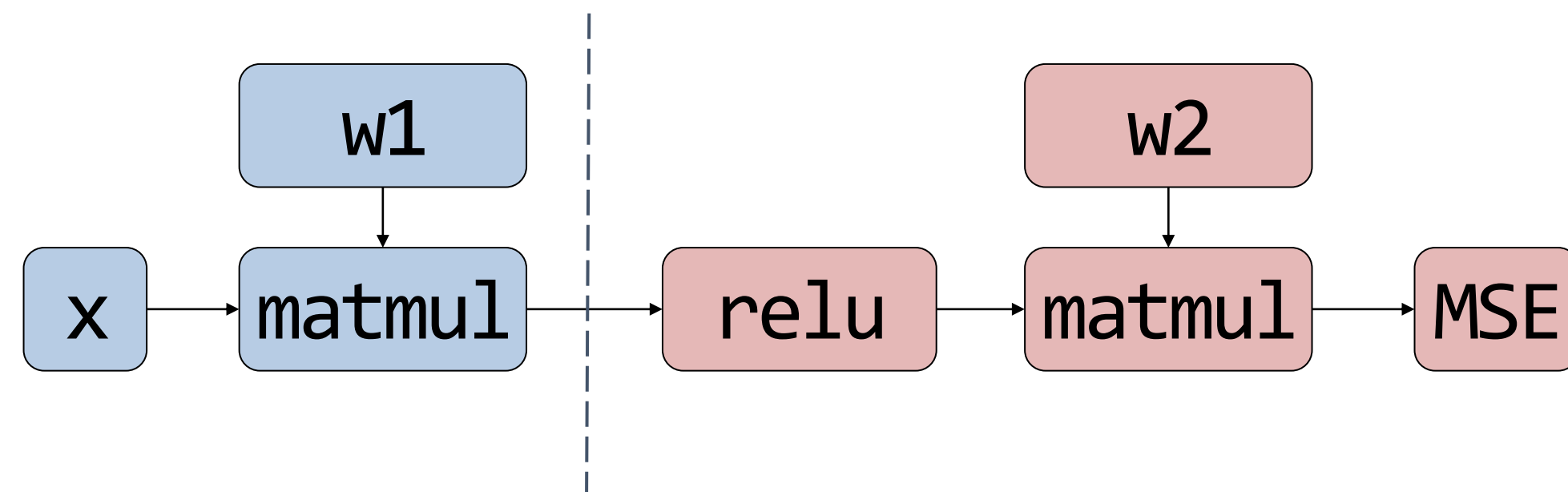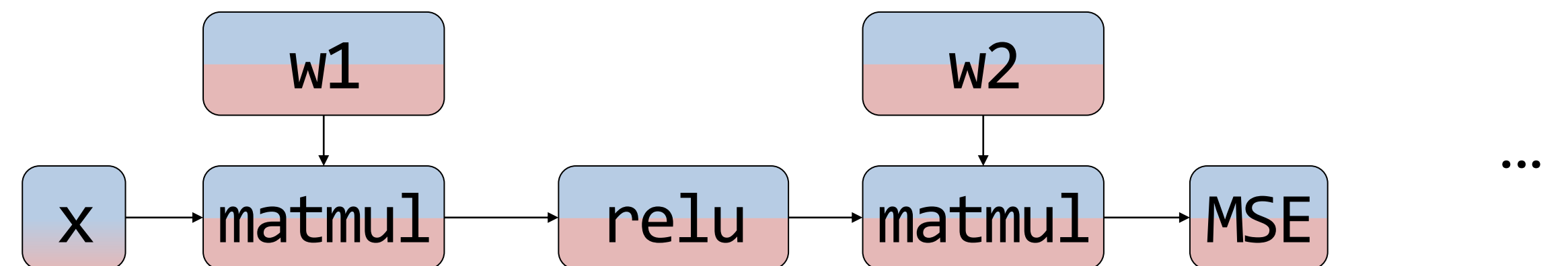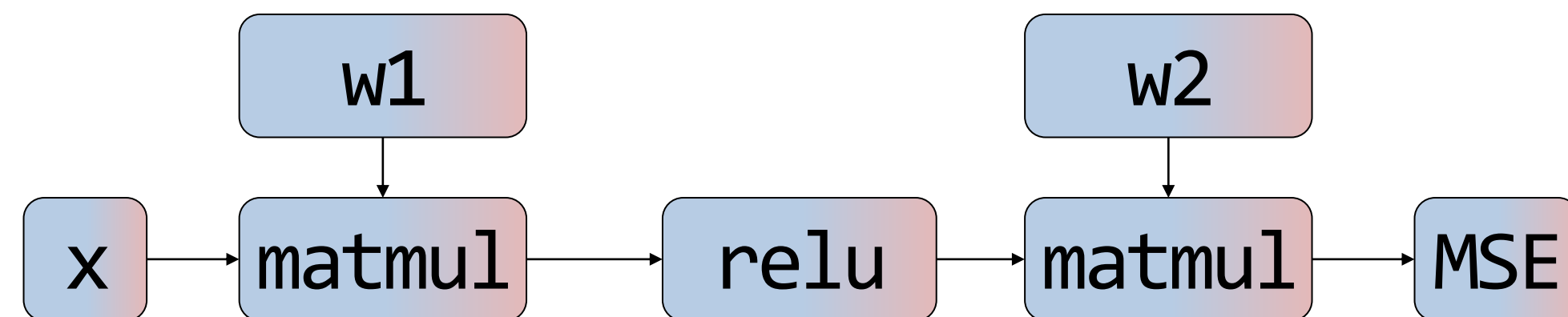## Strategy 1: Inter-operator Parallelism

## Strategy 2: Intra-operator Parallelism

# More Parallelisms...

**Multiple intra-op strategies for a single node**



Row-partitioned   Column-partitioned   Replicated   Device 3   Device 4

**More strategies**

# Summary: Inter-op and Intra-op Parallelisms



**Inter-op parallelism:** Assign different operators to different devices.

**Intra-op parallelism:** Assign different regions of a single operator to different devices.

# Inside Intra- and Inter-op Parallelism

Device 1 (blue)
Device 2 (pink)

Row-partitioned
Column-partitioned
Replicated
all-reduce
P2P

$$Y = X \cdot W_1 \cdot W_2 = X \cdot \begin{bmatrix} W_1^{d1} & W_1^{d2} \end{bmatrix} \cdot \begin{bmatrix} W_2^{d1} \\ W_2^{d2} \end{bmatrix}$$

# Inter-op and Intra-op Parallelism: Characteristics



**Inter-op parallelism:** Requires point-to-point communication but results in device idle

**Intra-op parallelism:** Devices are busy but requires collective communication

# Inter-op and Intra-op Parallelism: Characteristics



**Inter-op parallelism**

**Intra-op parallelism**

**Trade-off**

| | Inter-operator Parallelism | Intra-operator Parallelism |
|---|:---:|:---:|
| Communication | Less | More |
| Device Idle Time | More | Less |

# ML Parallelization under New View



Fast connections
Slow connections

w1
w2
y

x
matmul
r...
...de
GPU GPU

relu'
ma...

**Theme problem:**
What's the best way to execute the graph
subject to memory and communication constraints?

matmul

sub
sub

GPU GPU GPU GPU
GPU GPU GPU GPU

# Where We Are

- Deep Learning as Dataflow Graphs
- Auto-differentiation Libraries
  - Symbolic vs. Imperative
  - Static vs. Dynamic
- DL Parallelism
  - **Inter-op parallelism**
  - Intra-op parallelism

# Computational Graph (Neural Networks) → Stages

## Computational Graph



## Devices (e.g., GPUs)

| Device 1 | Device 2 | Device 3 | Device 4 |

# Computational Graph (Neural Networks) → Stages



Computational Graph

Stage

Devices (e.g., GPUs)

Device 1

Device 2

Device 3

Device 4

# Execution & Data Movement



**Note:** The time spent on data transfer is typically **small,** since we only communicates stage outputs at stage boundaries between two stages.

# Timeline: Visualization of Inter-Operator Parallelism



- Gray area ( ☐ ) indicates devices being idle (a.k.a. Pipeline bubbles).

- Only 1 device activated at a time.

- **Pipeline bubble percentage** = bubble_area / total_area

    = (D - 1) / D, assuming D devices.

# Reduce Pipeline Bubbles via Pipelining Inputs

Input d → Stage 1

Input c → Stage 2

Input b → Stage 3

Input a → Stage 4

Used in inference.

| Device 1 | a | b | c | d |   |   |   |
| Device 2 |   | a | b | c | d |   |   |
| Device 3 |   |   | a | b | c | d |   |
| Device 4 |   |   |   | a | b | c | d |

Time

Pipeline bubbles percentage
$= (D - 1) / (D - 1 + N)$
with D devices and N inputs.

# Training: Forward & Backward Dependency

# How to Reduce Pipeline Bubbles for Training?

- Synchronous Pipeline Parallel Algorithms
  - GPipe
  - 1F1B
  - Interleaved 1F1B
  - TeraPipe
  - Chimera
- Asynchronous Pipeline Parallel Algorithms
  - AMPNet
  - Pipedream/Pipedream-2BW

# How to Reduce Pipeline Bubbles for Training?

- Synchronous Pipeline Parallel Algorithms
  - **GPipe**
  - **1F1B**
  - Interleaved 1F1B
  - TeraPipe
  - Chimera
- Asynchronous Pipeline Parallel Algorithms
  - AMPNet
  - Pipedream/Pipedream-2BW

# GPipe

**Idea:** Partition the input batch into multiple "*micro-batches*". Pipeline the micro-batches.

Accumulate the gradients of the micro-batches:

$$\nabla L_\theta(x) = \frac{1}{N} \sum_{i=1}^{N} \nabla L_\theta(x_i)$$

**Example:** Slice each input batch into 6 micro-batches:



Pipeline bubbles percentage = (D - 1) / (D - 1 + N)
with D devices and N micro-batches.

# GPipe: Memory Usage



= Parameters + Activation × #Micro-Batches

Per-Device Memory Usage

Intermediate activation

Model parameters

Forward (a)    Backward (a)    Forward (b)

|        | 0 | 1 | 2 | 3 | 4 | 5 | | | | | | 5 | 4 | 3 | 2 | 1 | 0 | Update | 0 | 1 |
| Device 1 | 0 | 1 | 2 | 3 | 4 | 5 | | | | | | 5 | 4 | 3 | 2 | 1 | 0 | | 0 | 1 |
| Device 2 | | 0 | 1 | 2 | 3 | 4 | 5 | | | | 5 | 4 | 3 | 2 | 1 | 0 | | | 0 | |
| Device 3 | | | 0 | 1 | 2 | 3 | 4 | 5 | | 5 | 4 | 3 | 2 | 1 | 0 | | | | | |
| Device 4 | | | | 0 | 1 | 2 | 3 | 4 | 5 | 5 | 4 | 3 | 2 | 1 | 0 | | | | | |

...

Time

# GPipe Schedule:



Forward (for input batch a)    Backward (a)    Forward (b)

40

1F1B (1 Forward 1 Backward) Schedule:

Same Latency

Perform backward as early as possible

# 1F1B Memory Usage



= Parameters + Activation × ~~#Micro-Batches~~ #Devices

# Where We Are

- Deep Learning as Dataflow Graphs
- Auto-differentiation Libraries
  - Symbolic vs. Imperative
  - Static vs. Dynamic
- DL Parallelism
  - Inter-op parallelism
  - **Intra-op parallelism**

# Recap: Intra-op and Inter-op

## Strategy 1: Inter-operator Parallelism
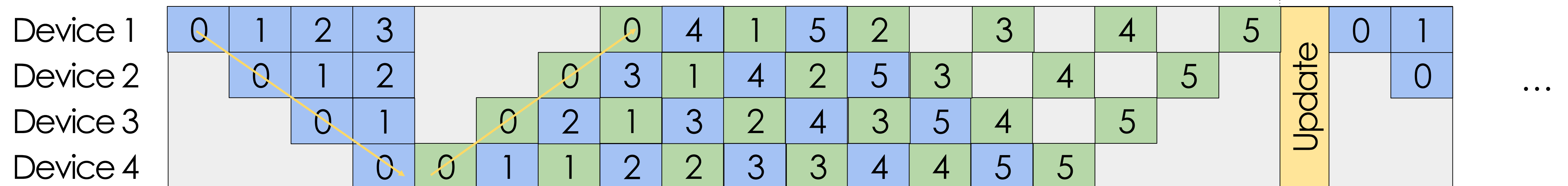


## Strategy 2: Intra-operator Parallelism



**This section:**
How to parallelize an **operator** ?
How to parallelize a **graph** ?

# Parallelize One Operator

Element-wise operators

```
for n in range(0, N):
  for d in range(0, D):
    C[n,d] = A[n,d] + B[n,d]
```
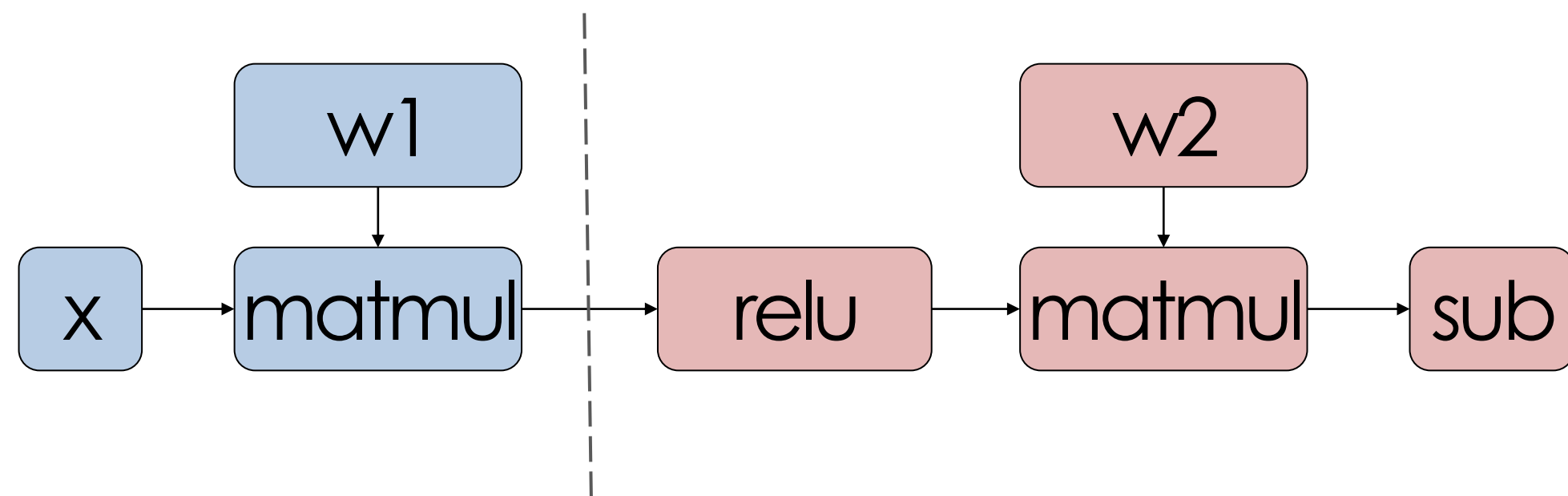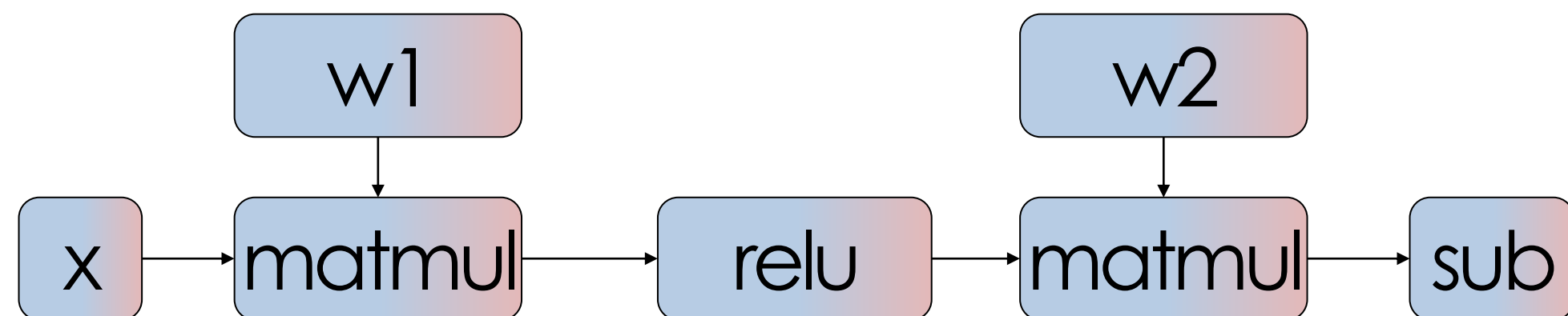
No dependency on the two for-loops.
Can arbitrarily split the for-loops on different devices.

device 1   device 2   device 3   device 4

Parallelize loop n



n    C   =   A   +   B

d

Parallelize both loop n and loop d



n    C   =   A   +   B

d

a lot of other variants
…

# Parallelize One Operator

Matrix multiplication

No dependency on the two spatial for-loops.
Can arbitrarily split the for-loops on different devices.

```
for i in range(0, N):
  for j in range(0, M):
    for k in range(0, K):
      C[i,j] = C[i,j] + A[i,k] x B[k,j]
```
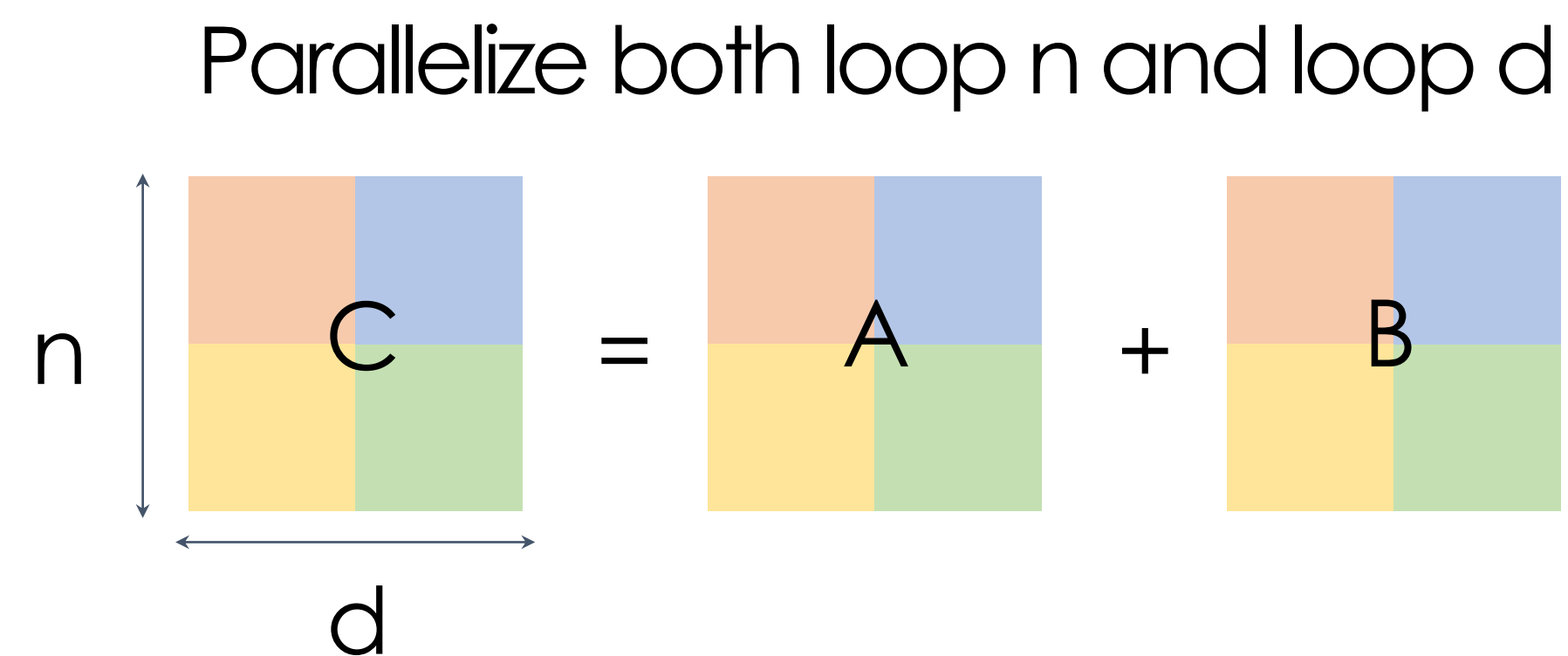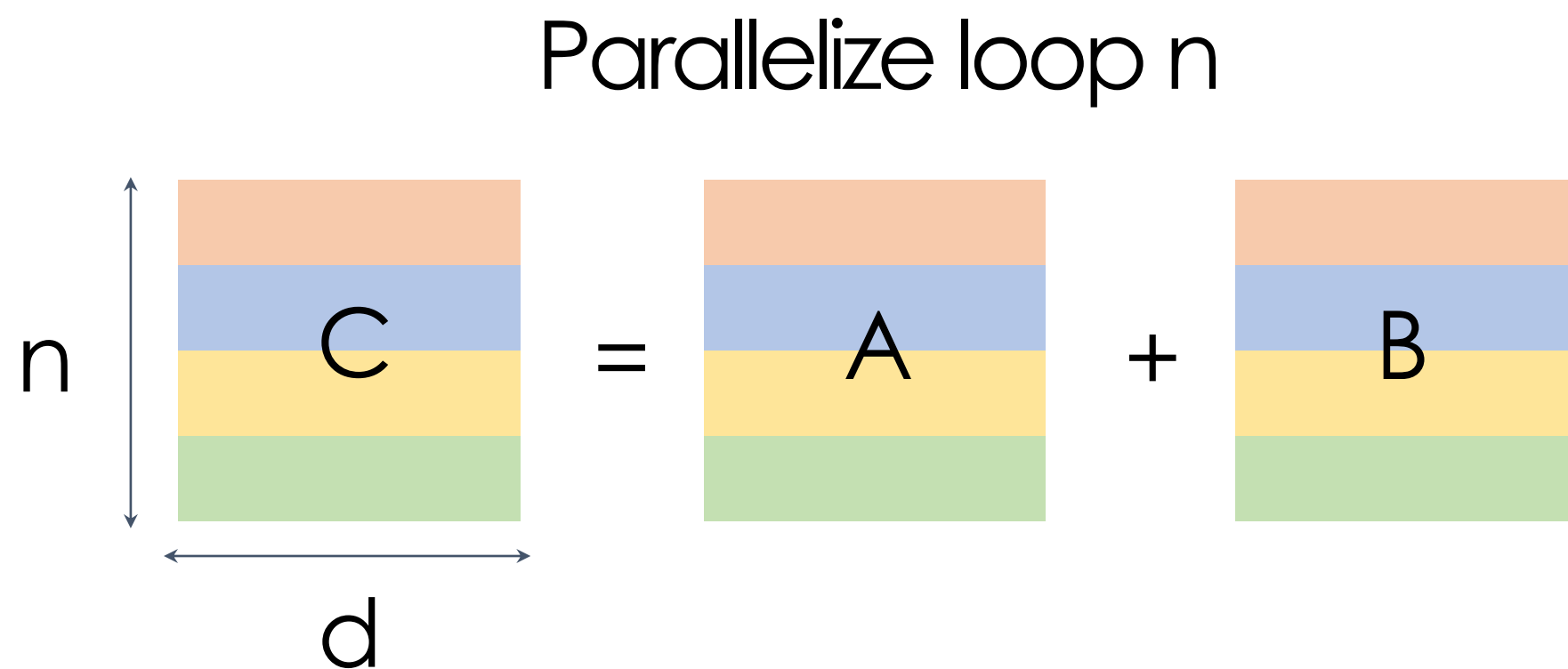
Accumulation on this reduction loop.
Have to accumulate partial results if we split this for-loop

device 1  device 2  device 3  device 4  replicated

Parallelize loop i



$$\begin{bmatrix} C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \\ A_3 \\ A_4 \end{bmatrix} \times B$$

# Parallelize One Operator

Matrix multiplication

No dependency on the two spatial for-loops.
Can arbitrarily split the for-loops on different devices.

```
for i in range(0, N):
 for j in range(0, M):
  for k in range(0, K):
   C[i,j] = C[i,j] + A[i,k] x B[k,j]
```

Accumulation on this reduction loop.
Have to accumulate partial results if we split this for-loop

■ device 1   ■ device 2   ■ device 3   ■ device 4   ■ replicated

Parallelize loop k



$$C = [A_1 \ A_2 \ A_3 \ A_4] \begin{bmatrix} B_1 \\ B_2 \\ B_3 \\ B_4 \end{bmatrix} = A_1 B_1 + A_2 B_2 + A_3 B_3 + A_4 B_4$$

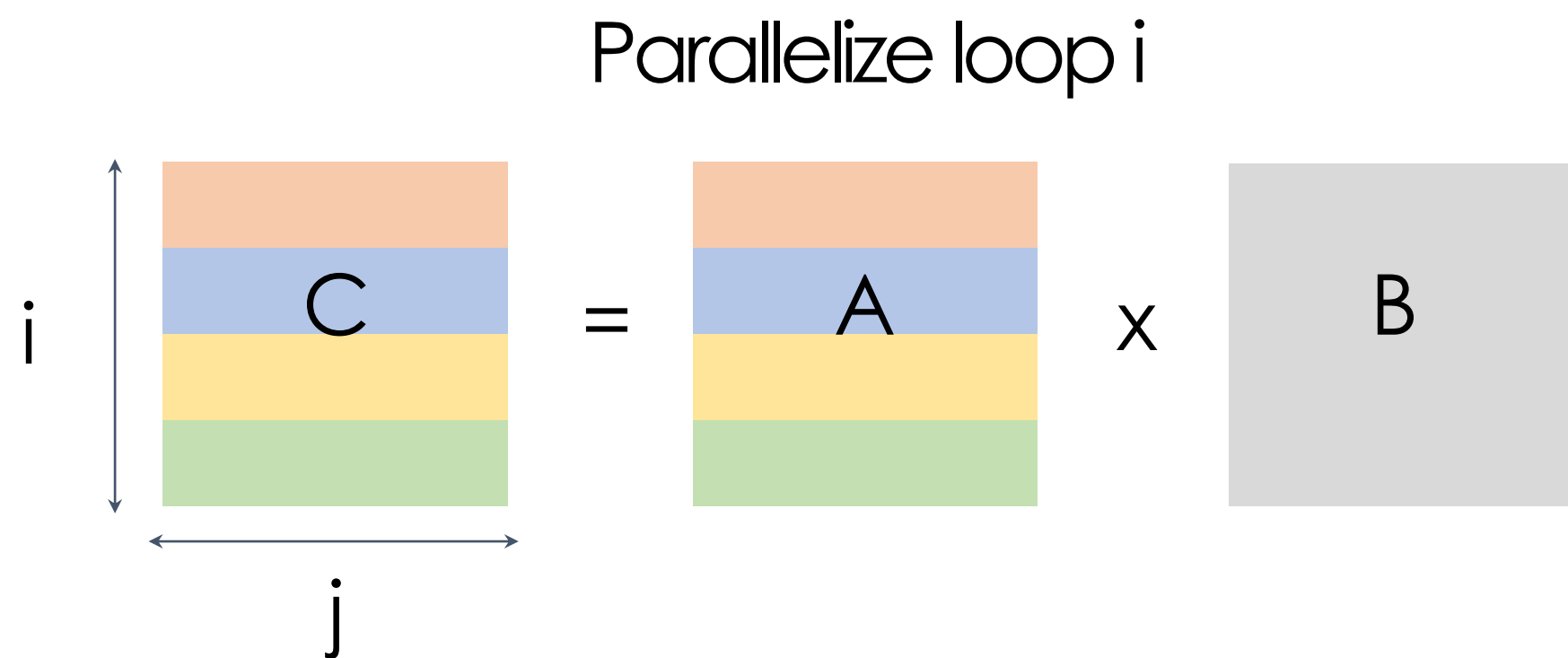(got by all-reduce)

# Parallelize One Operator
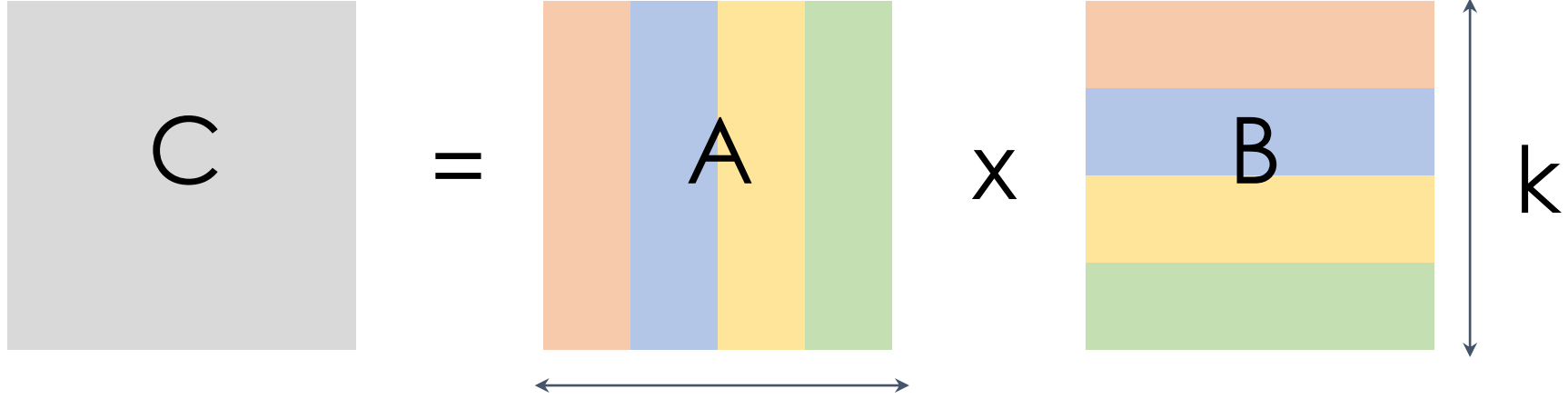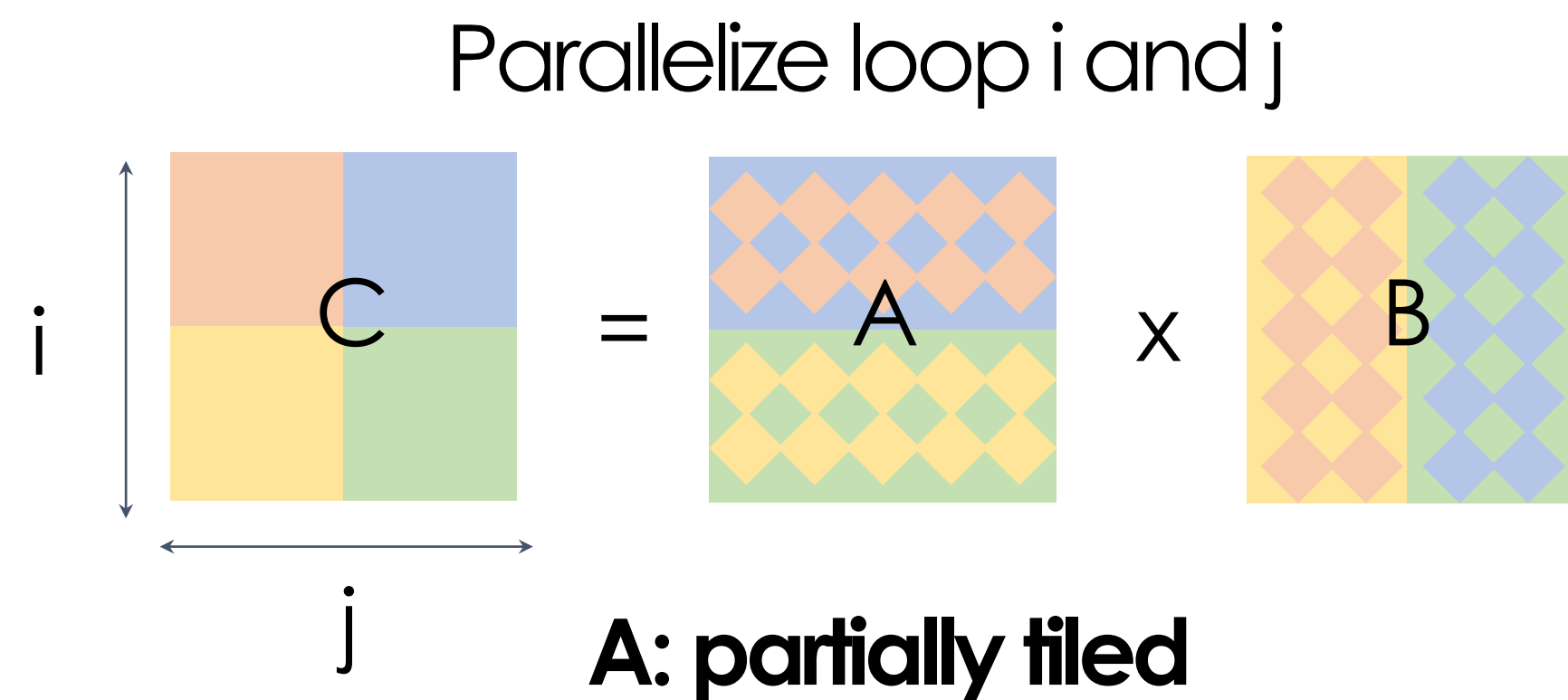
Matrix multiplication

```
for i in range(0, N):
  for j in range(0, M):
    for k in range(0, K):
      C[i,j] = C[i,j] + A[i,k] x B[k,j]
```
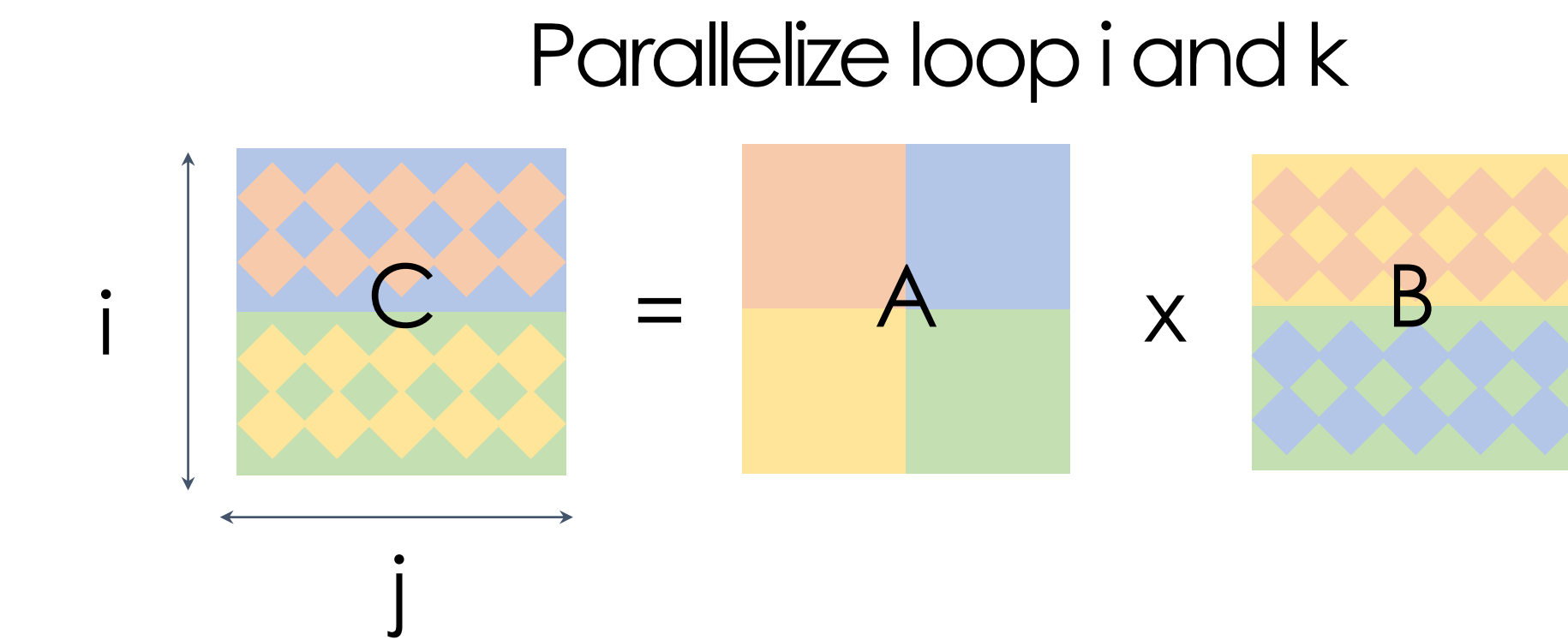
No dependency on the two spatial for-loops.
Can arbitrarily split the for-loops on different devices.

Accumulation on this reduction loop.
Have to accumulate partial results if we split this for-loop

 device 1   device 2   device 3   device 4

Parallelize loop i and j



**A: partially tiled**

Device 1 and 2 hold a replicated tile
Device 3 and 4 hold a replicated tile

Parallelize loop i and k



**C: got by all-reduce**

a lot of other variants

…

# Parallelize One Operator

2D Convolution

Simple spatial loops. Can be arbitrarily split.

Stencil computation loops. Splitting these requires careful boundary handling.

Reduction loop. Need to accumulate partial results.

Reduction loops. But usually too small (<= 5) for parallelization.

```
for n in range(0, N):
  for co in range(0, CO):
    for h in range(0, H):
      for w in range(0, W):
        for ci in range(0, CI):
          for kh in range(0, KH):
            for kw in range(0, KW):
              C[n,co,h,w] += A[n,co,h+kh,w+kw] x B[kh,kw,co,ci]
```

**Simple case:** Parallelize loop n, co, ci, then the parallelization strategies are almost the same as matmul's.
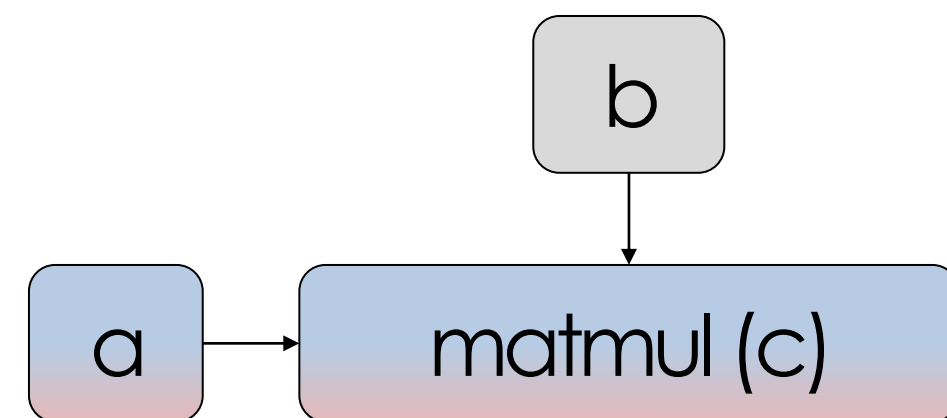
**Complicated case**: Parallelize loop h and w

# Data Parallelism as A Case of Intra-op Parallelism
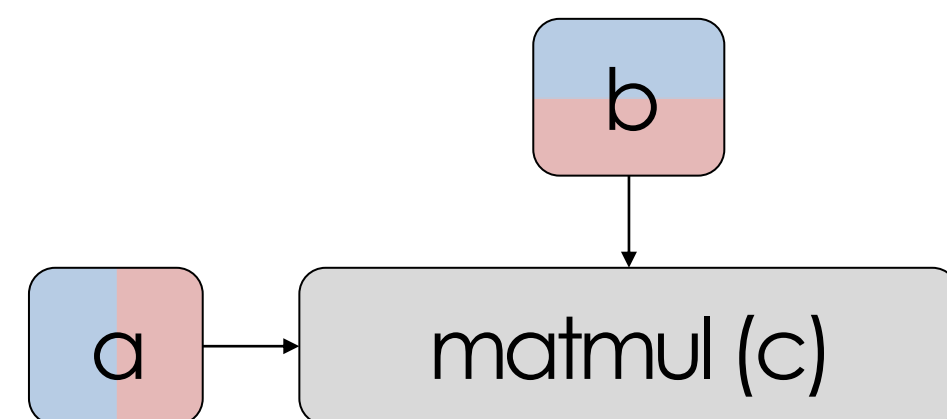
Replicated    Row-partitioned    Column-partitioned

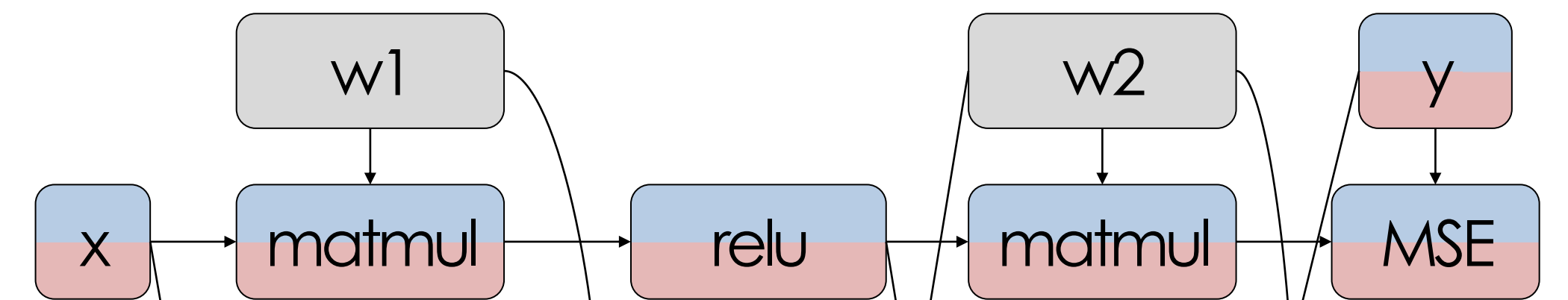## Matmul Parallelization Type 1
communication cost = 0

b → matmul (c) ← a

## Matmul Parallelization Type 2
communication cost = all-reduce(c)

b → matmul (c) ← a

**Forward Pass**
Two "Type 1" matmuls: no communication

w1, x → matmul → relu → w2, matmul → y, MSE

**Backward Pass**
One "Type 1" matmul: no communication
Two "Type 2" matmuls: require all-reduce

relu' ← matmul ← MSE'
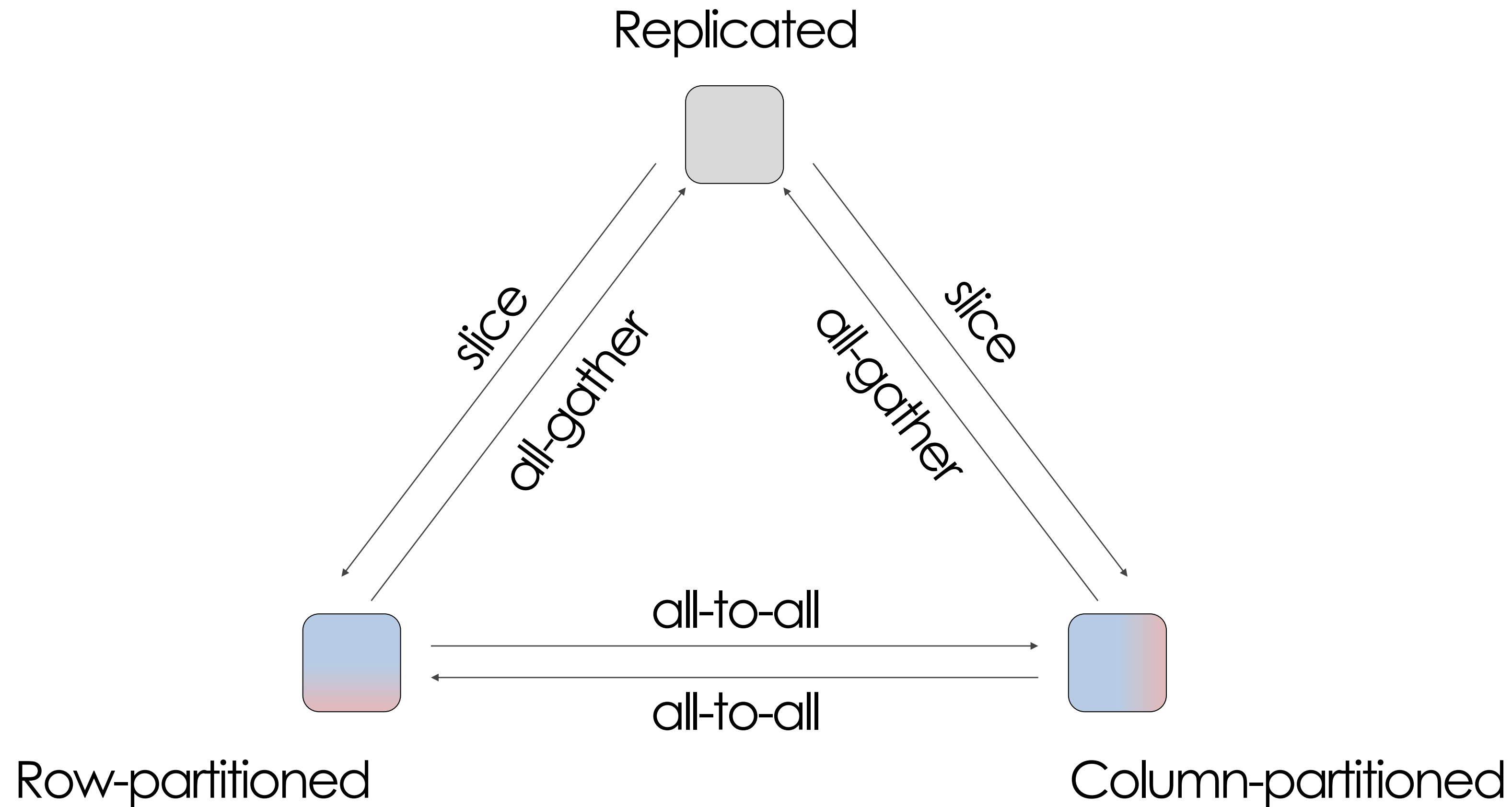relu' → matmul → new_w1
MSE' → matmul → new_w2

# Re-partition Communication Cost

Different operators' parallelization strategies require different partition format of the same tensor
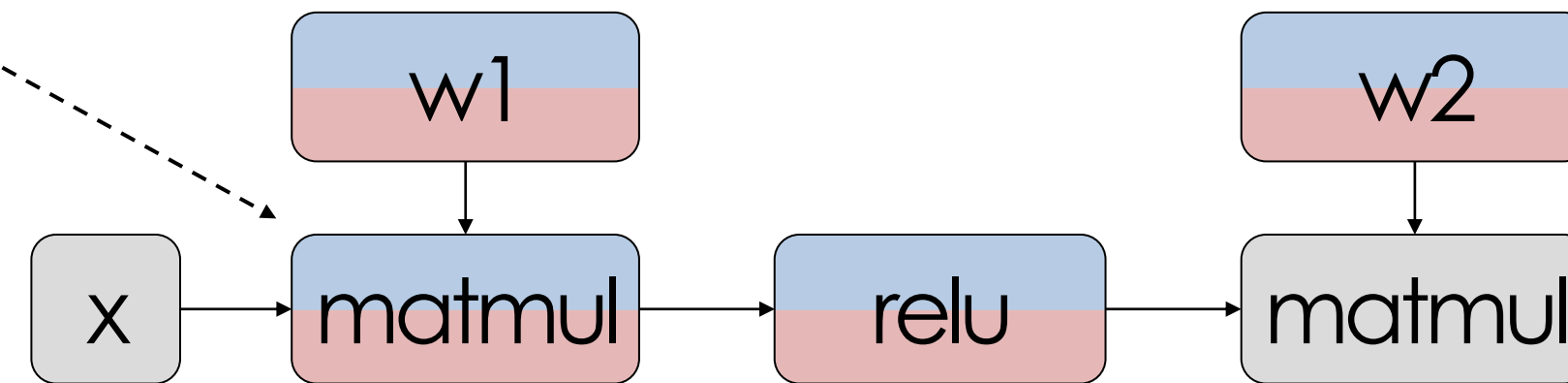
# Re-partition Communication Cost

Different operators' parallelization strategies require different partition format of the same tensor

# Parallelize All Operators in a Graph

**Problem**

Pick a parallel strategy
of each operator



*Minimize*   Node costs (computation + communication) + Edge costs (re-partition communication)
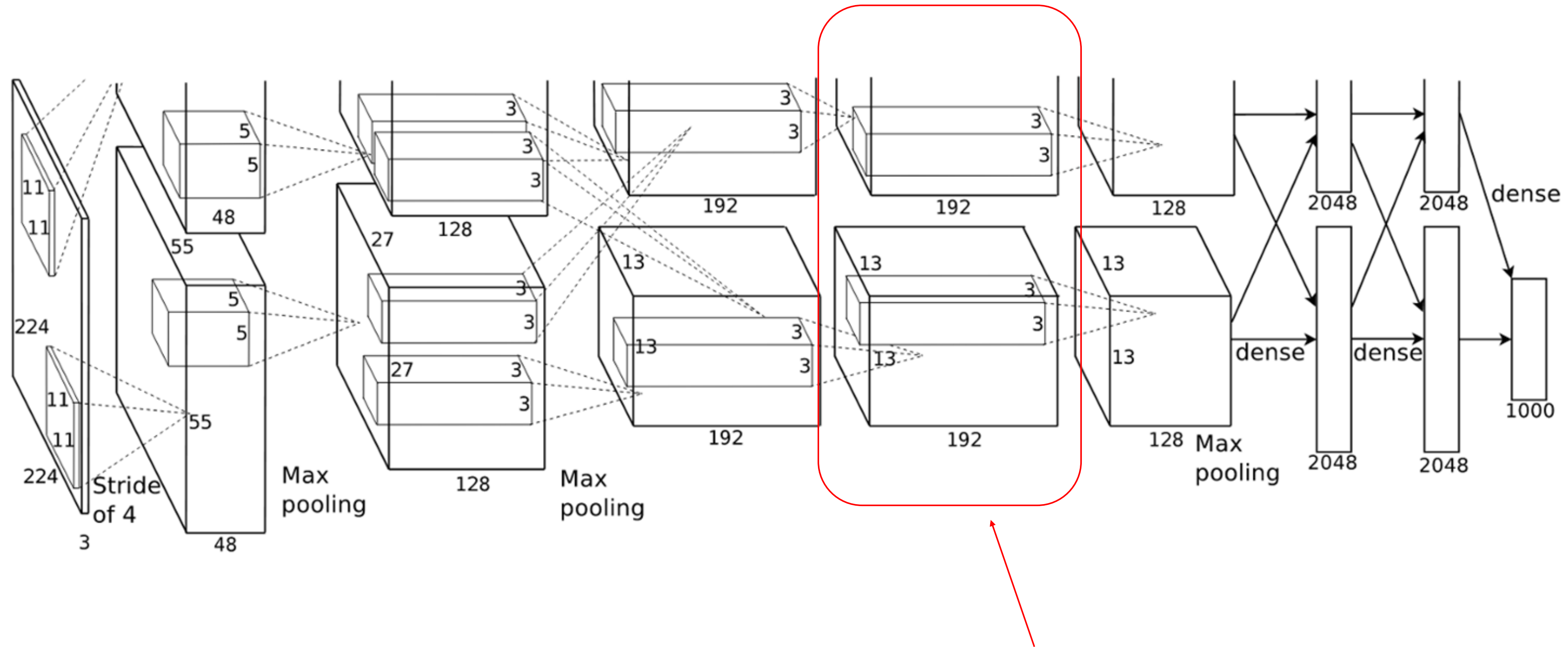
**Solution**

Manual design
Randomized search
Dynamic programming
Integer linear programming

# Important Projects

- Model-specific Intra-op Parallel Strategies
  - **AlexNet**
  - **Megatron-LM**
  - GShard MoE

- Systems for Intra-op Parallelism
  - ZeRO
  - Mesh-Tensorflow
  - GSPMD
  - Tofu
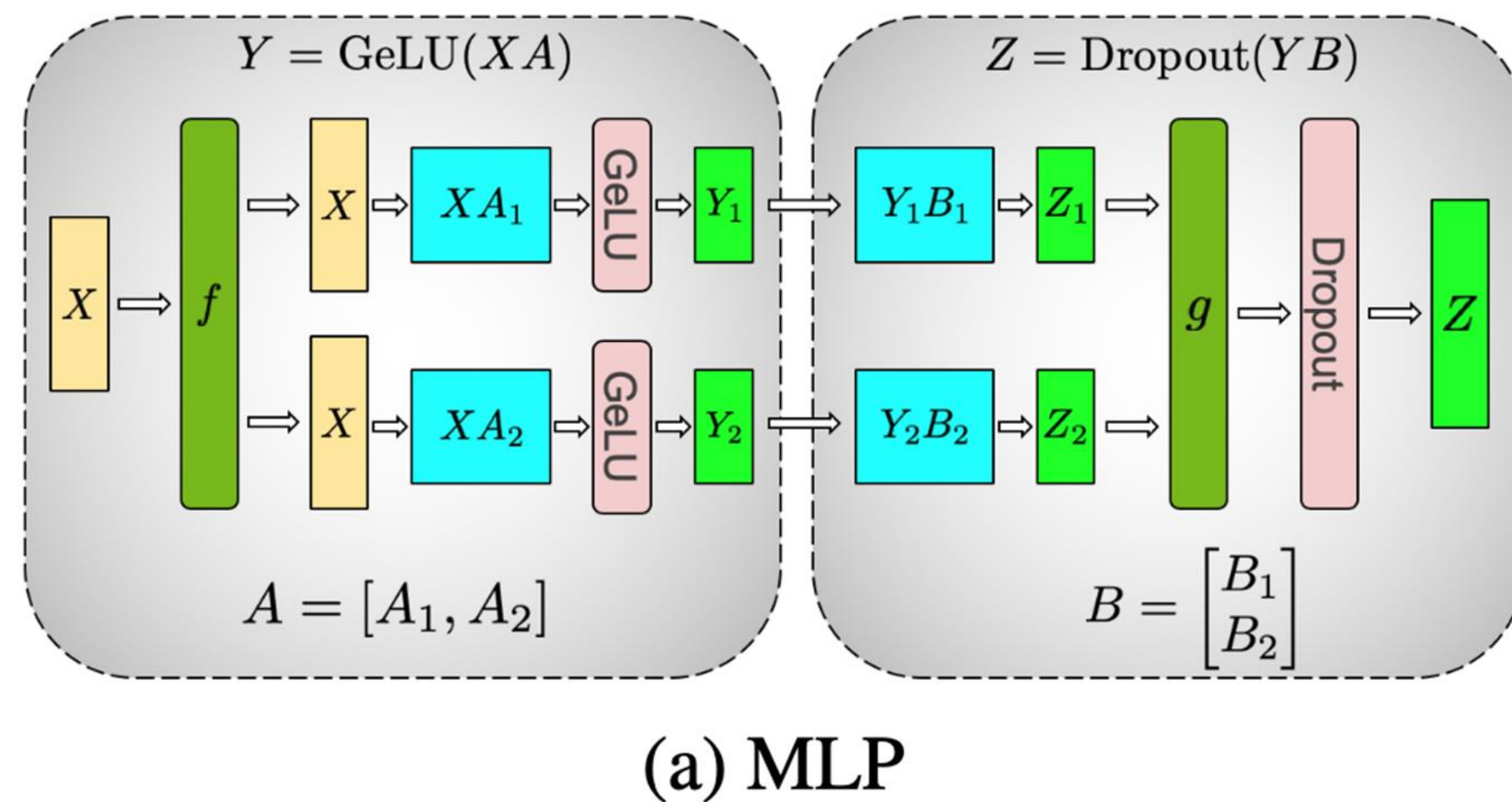  - FlexFlow

# AlexNet

Result: increase top-1 accuracy by 1.7%



Assign a group convolution layer to 2 GPUs

Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional
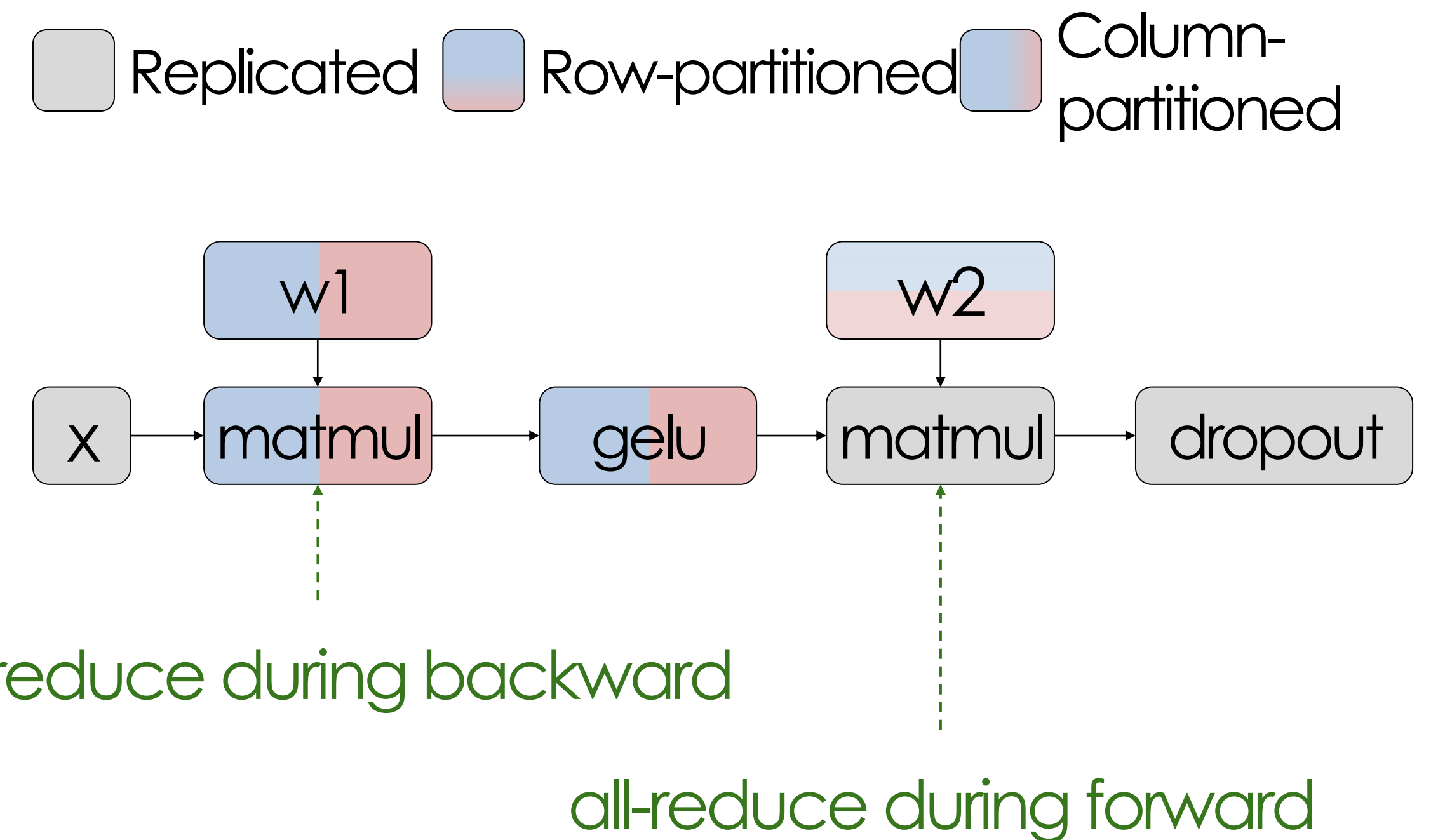
# Megaton-LM

Result: a large language model with 8.3B parameters that outperforms SOTA results

Figure 3 from the paper:
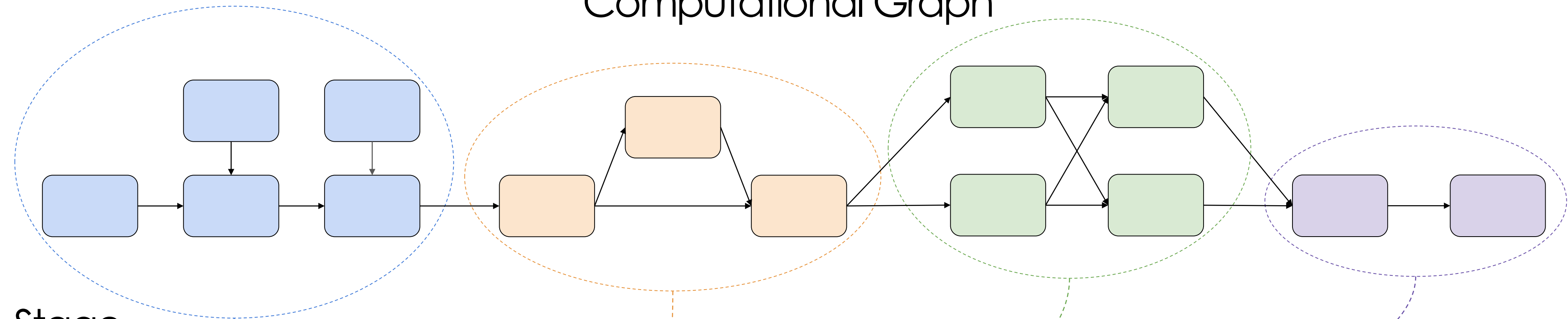How to partition the MLP in the transformer.



$$Y = \text{GeLU}(XA)$$
$$Z = \text{Dropout}(YB)$$
$$A = [A_1, A_2]$$
$$B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

(a) MLP

Illustrated with the notations in this tutorial



Replicated   Row-partitioned   Column-partitioned

all-reduce during backward

all-reduce during forward

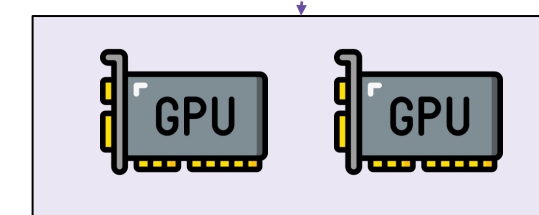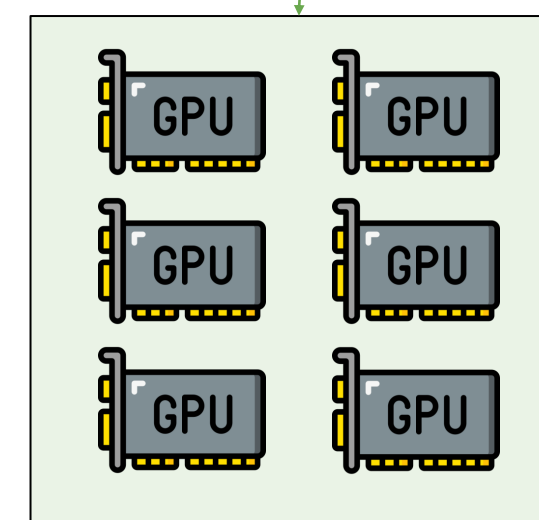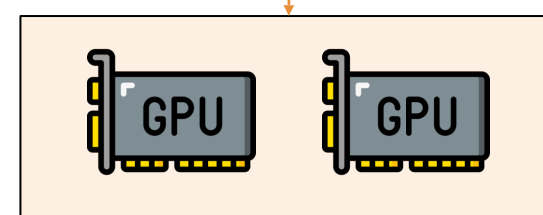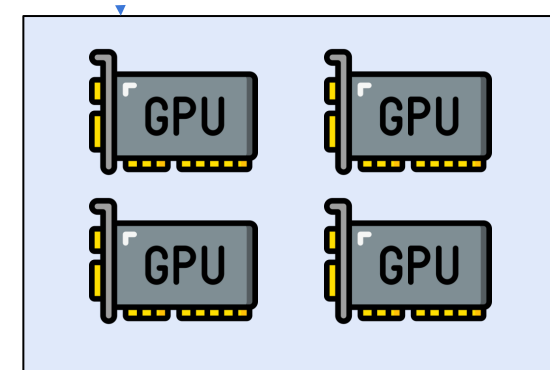# Combine Intra-op Parallelism and Inter-op Parallelism



Computational Graph

Stage

Device Mesh

Intra-op Parallelism

Inter-op Parallelism

# Intra-operator Parallelism Summary

- We can parallelize a single operator by exploiting its internal parallelism

- To do this for a whole computational graph, we need to choose strategies for all nodes in the graph to minimize the communication cost

- Intra-op and inter-op can be combined