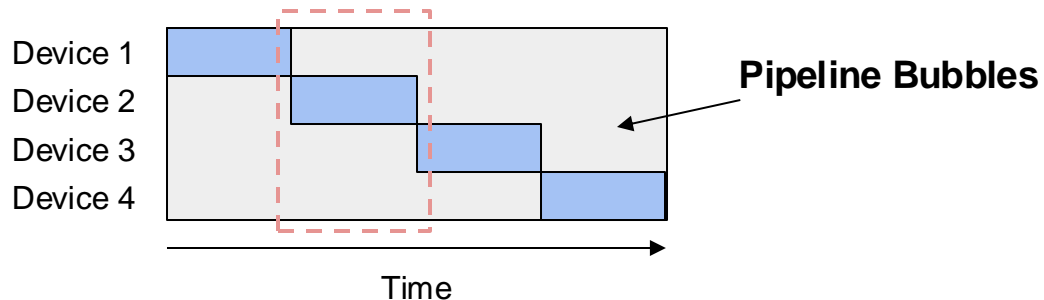



# Where We Are

- Motivation
- History
- Parallelism Overview
- Data parallelism
- Model parallelism
  - **Inter-op parallelism**
  - Intra-op parallelism
- Auto-parallelization

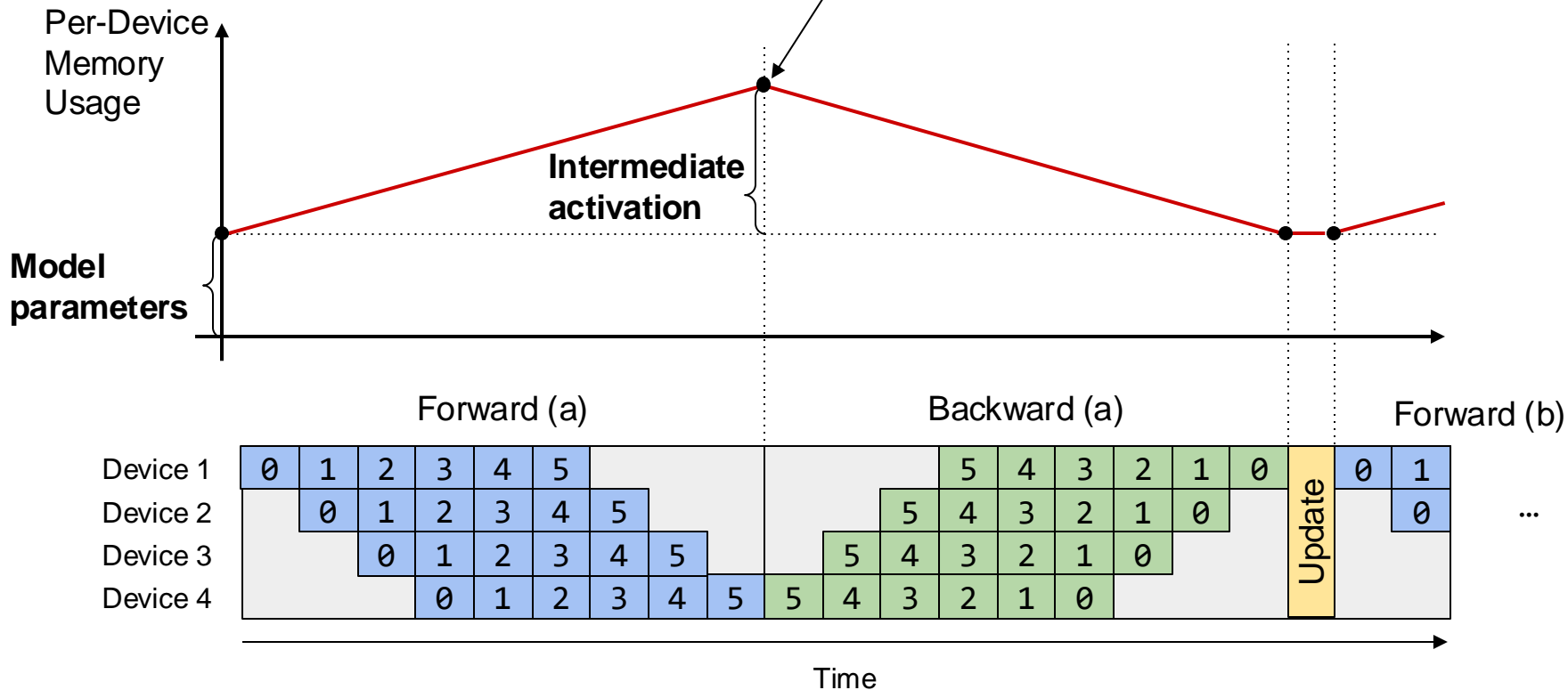
# Recap



- Gray area (  indicates devices being idle (a.k.a. Pipeline bubbles).
- Only 1 device activated at a time.
- **Pipeline bubble percentage** =  $(D - 1) / D$ , assuming  $D$  devices.

# Recap

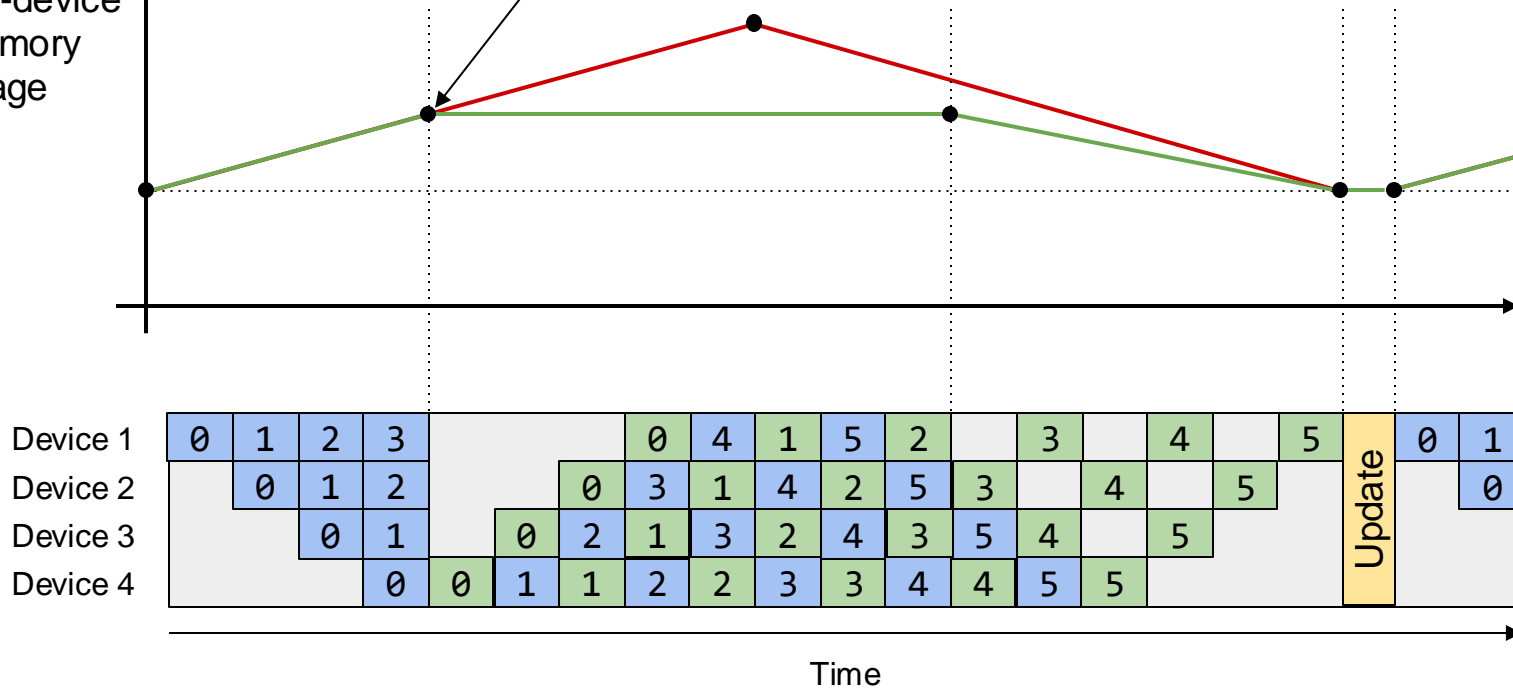
$$= \text{Parameters} + \text{Activation} \times \text{\#Micro-Batches}$$



# Recap

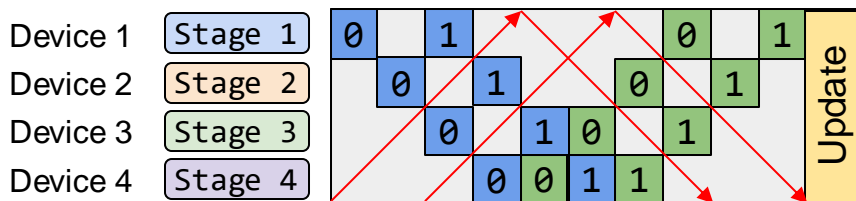
Maximum  
per-device  
memory  
usage

$$= \text{Parameters} + \text{Activation} \times \text{\#Micro-Batches} \times \text{\#Devices}$$

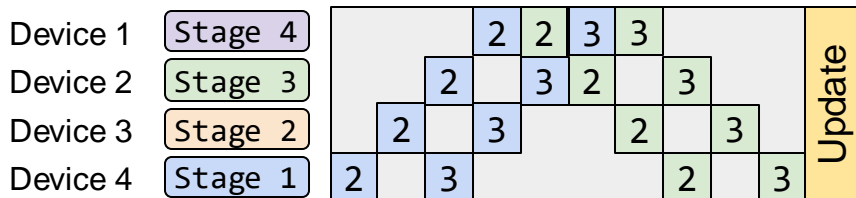


# Recap: Chimera

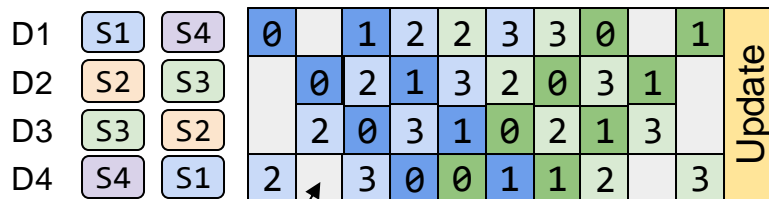
**Idea:** Store bi-directional stages and combine bidirectional pipeline to further reduce pipeline bubbles.



+



Extra copy of parameters & extra synchronization.



Pipeline bubbles percentage  
 $= (D - 2) / (D - 2 + 2N)$   
 with D devices and N micro-batches.

# Synchronous Pipeline Schedule Summary



## Pros:

- Keep the convergence semantics. The training process is exactly the same as training the neural network on a single device.



## Cons:

- Pipeline bubbles.
- Reducing pipeline bubbles typically requires splitting inputs into smaller components, but too small input to the neural network will reduce the hardware efficiency.

# Asynchronous Pipeline Schedules

**Idea:** Start next round of forward pass before backward pass finishes.



## Pros:

- No Pipeline bubbles.



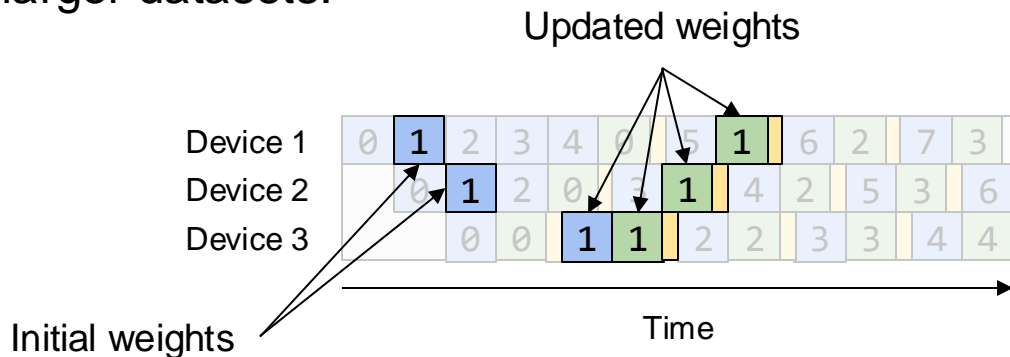
## Cons:

- Break the synchronous training semantics. Now the training will involve stalled gradient.
- Algorithms may store multiple versions of model weights for consistency.

# AMPNet

**Idea:** Fully asynchronous. Each device performs forward pass whenever free and updates the weights after every backward pass.

**Convergence:** Achieve similar accuracy on small datasets (MNIST 97%), hard to generalize to larger datasets.



**PipeMare:** modify the optimizer to improve AMPNet convergence

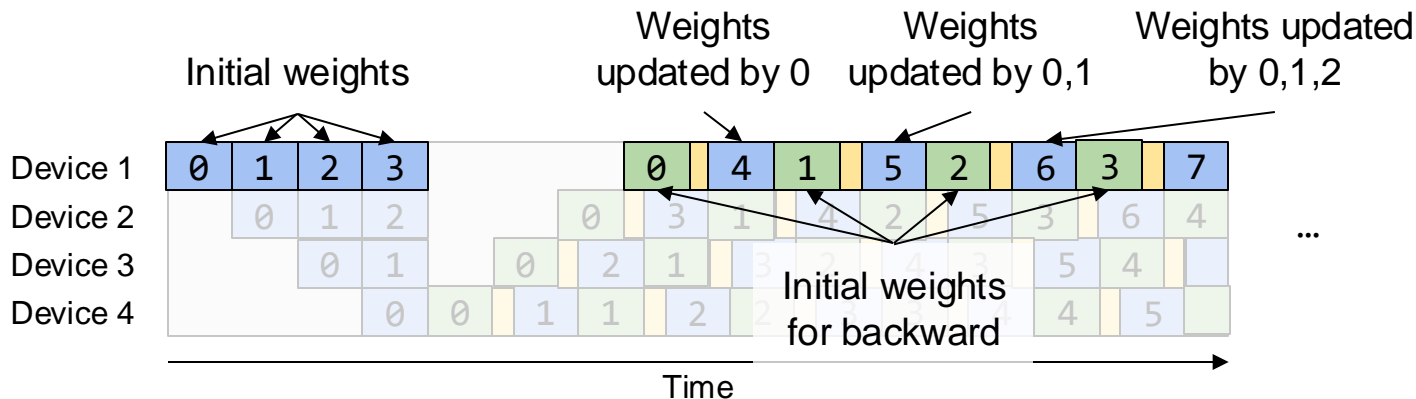


# Pipedream

**Idea:** Enforce the same version of weight for a single input batch by storing multiple weight versions.

**Convergence:** Similar accuracy on ImageNet with a 5x speedup compared to data parallel.

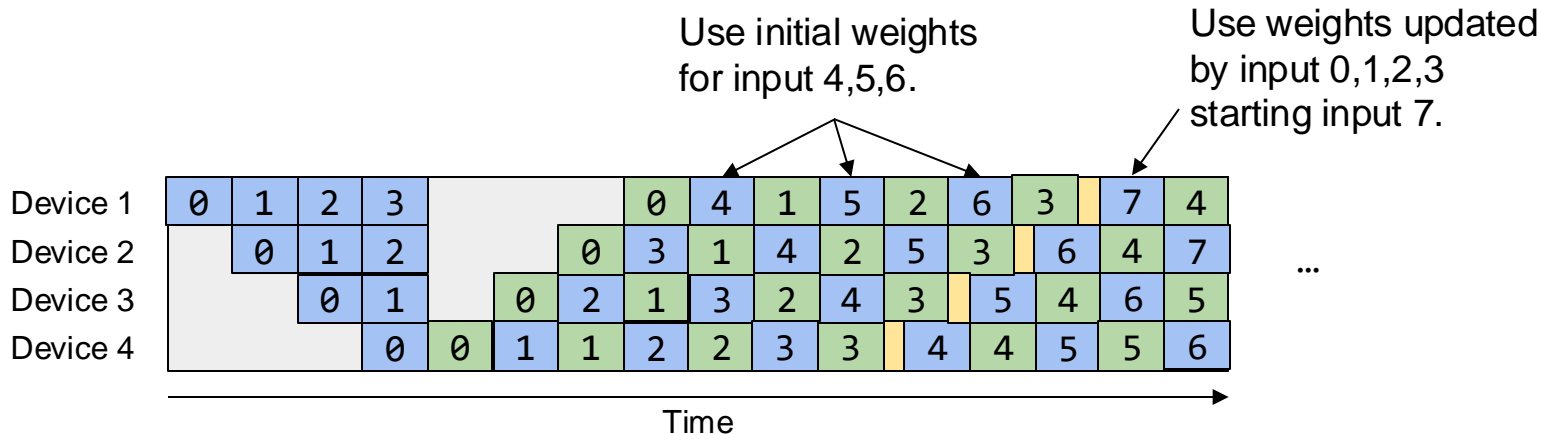
**Con:** No memory saving compared to single device case.



# Pipedream-2BW

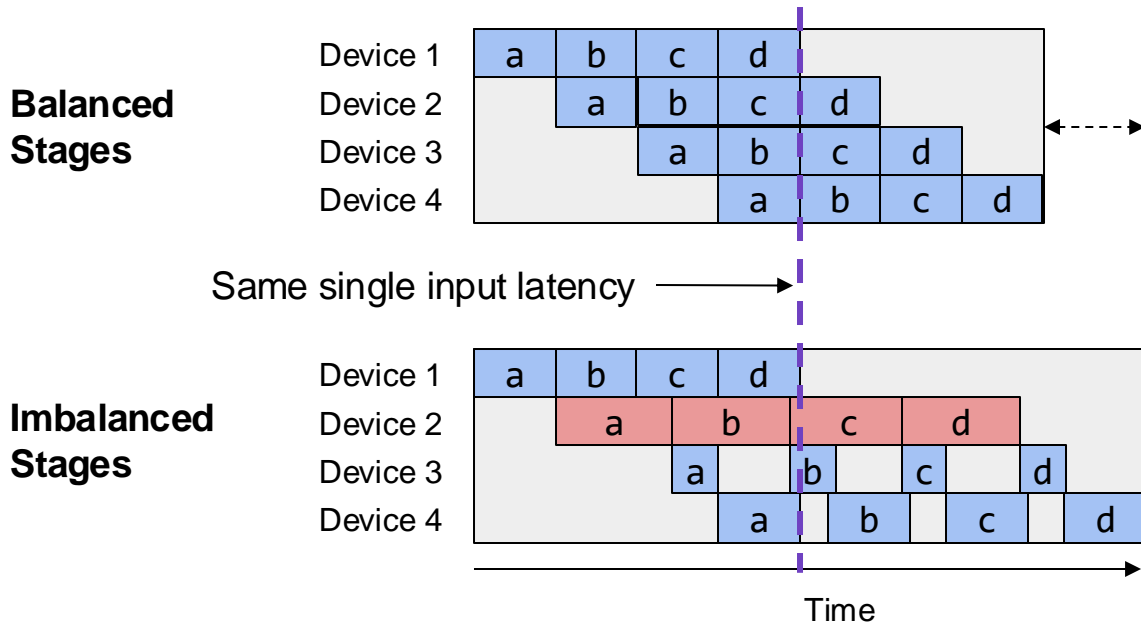
**Idea:** Reduce Pipedream's memory usage (only store 2 copies) by updating weights less frequently. Weights always stalled by 1 update.

**Convergence:** Similar training accuracy on language models (BERT/GPT)



# Imbalanced Pipeline Stages

Pipeline schedules works best with balanced stages:



# Frontier: Automatic Stage Partitioning

**Goal:** Minimize maximum stage latency & maximize parallelization

## Reinforcement Learning Based (mainly for device placement):

1. Mirhoseini, Azalia, et al. "Device placement optimization with reinforcement learning." *ICML 2017*.
2. Gao, Yuanxiang, et al. "Spotlight: Optimizing device placement for training deep neural networks." *ICML 2018*.
3. Mirhoseini, Azalia, et al. "A hierarchical model for device placement." *ICLR 2018*.
4. Addanki, Ravichandra, et al. "Placeto: Learning generalizable device placement algorithms for distributed machine learning." *NeurIPS 2019*.
5. Zhou, Yanqi, et al. "Gdp: Generalized device placement for dataflow graphs." *Arxiv 2019*.
6. Paliwal, Aditya, et al. "Reinforced genetic algorithm learning for optimizing computation graphs." *ICLR 2020*.
7. ...

## Optimization (Dynamic Programming/Linear Programming) Based:

1. Narayanan, Deepak, et al. "PipeDream: generalized pipeline parallelism for DNN training." *SOSP 2019*.
2. Tarnawski, Jakub M., et al. "Efficient algorithms for device placement of dnn graph operators." *NeurIPS 2020*.
3. Fan, Shiqing, et al. "DAPPLE: A pipelined data parallel approach for training large models." *PPoPP 2021*.
4. Tarnawski, Jakub M., Deepak Narayanan, and Amar Phanishayee. "Piper: Multidimensional planner for dnn parallelization." *NeurIPS 2021*.
5. Zheng, Lianmin, et al. "Alpa: Automating Inter-and Intra-Operator Parallelism for Distributed Deep Learning." *OSDI 2022*.
6. ...

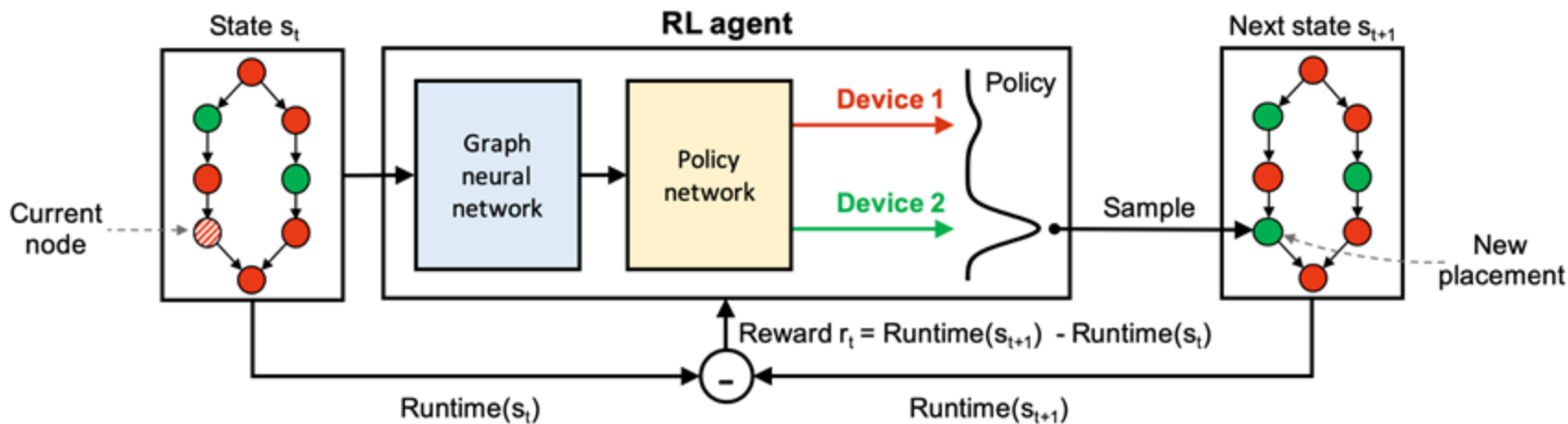
# RL-Based Partitioning Algorithm

**State:** Device assignment plan for a computational graph.

**Action:** Modify the device assignment of a node.

**Reward:** Latency difference between the new and old placements.

Trained with **policy gradient** algorithm.



# Inter-operator Parallelism Summary

**Idea:** Assign different operators of the computational graph to different devices and executed in a pipelined fashion.

Method	General computational graph	No pipeline bubbles	Same convergence as single device
Device Placement	✗	✗	✓
Synchronous Schedule	✓	✗	✓
Asynchronous Schedule	✓	✓	✗

**Stage Partitioning:** Imbalance stage → More pipeline bubble

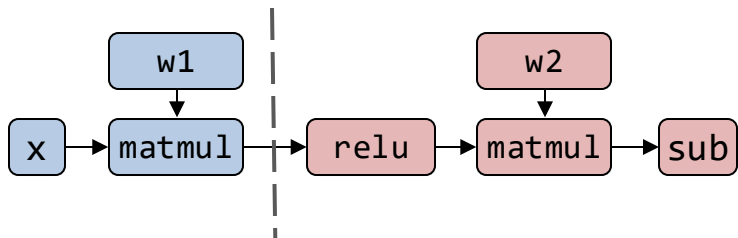
**RL-Based / Optimization-Based** Automatic Stage Partitioning

# Where We Are

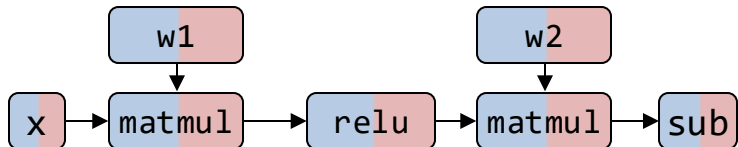
- Motivation
- History
- Parallelism Overview
- Data parallelism
- Model parallelism
  - Inter-op parallelism
  - **Intra-op parallelism**
- Auto-parallelization

# Recap: Intra-op and Inter-op

## Strategy 1: Inter-operator Parallelism



## Strategy 2: Intra-operator Parallelism



**This section:**

1. How to parallelize an **operator** ?
2. How to parallelize a **graph** ?



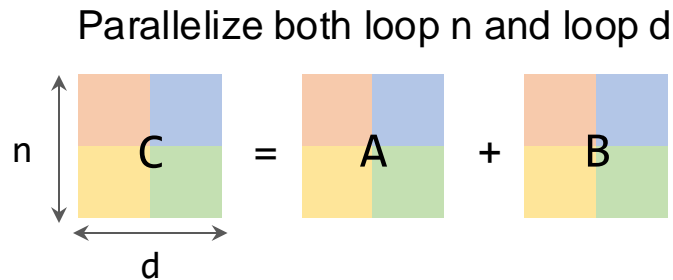
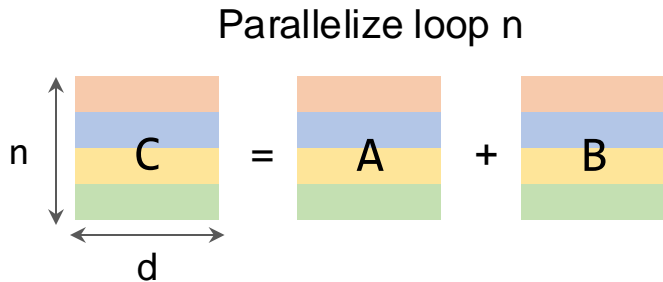
# Parallelize One Operator

## 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



a lot of  
other variants  
...

# Parallelize One Operator

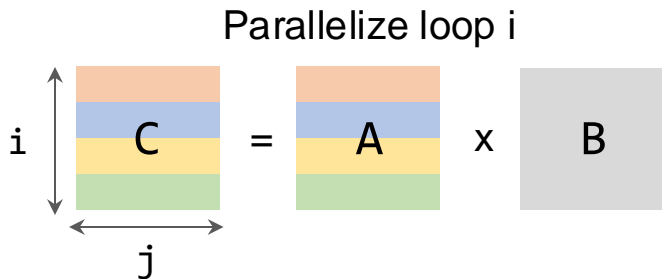
## Matrix multiplication

```
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 replicated



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

# Parallelize One Operator

## Matrix multiplication

```
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]
```

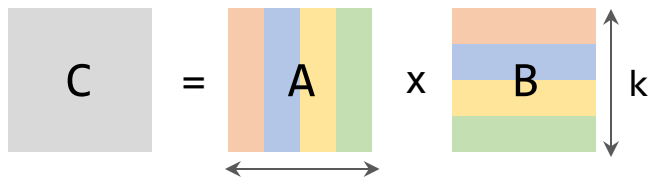
No dependency on the two spatial for-loops.  
Can arbitrarily split the for-loops on different devices.

Accumulation on this reduction loop.

Have to accumulate partial results if we split this for-loop

device 1
  device 2
  device 3
  device 4
  replicated

Parallelize loop k



(got by all-reduce)  $k$

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

# Parallelize One Operator

## Matrix multiplication

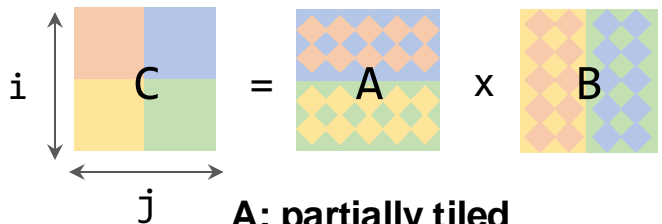
```
for i in range(0, N):  
    for j in range(0, M):  
        for k in range(0, K):  
            C[i,j] = C[i,j] + A[i,k] x B[k,j]
```

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

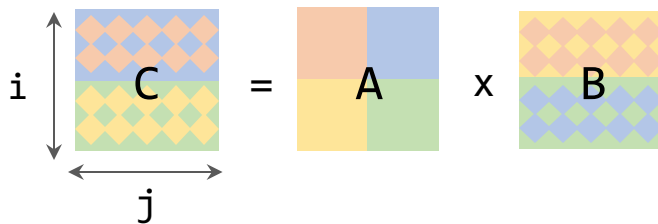
device 1 device 2 device 3 device 4

Parallelize loop i and j



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

Parallelize loop i and k



a lot of  
other variants  
...

# Parallelize One Operator

## 2D Convolution

```
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 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 ( $\leq 5$ ) for parallelization.

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

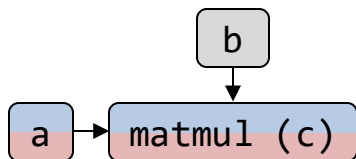
**Complicated case:** Parallelize loop  $h$  and  $w$

# Data Parallelism as A Case of Intra-op Parallelism

 Replicated  Row-partitioned  Column-partitioned

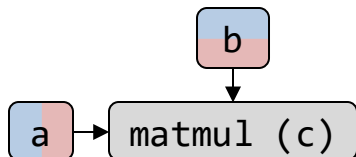
## Matmul Parallelization Type 1

communication cost = 0



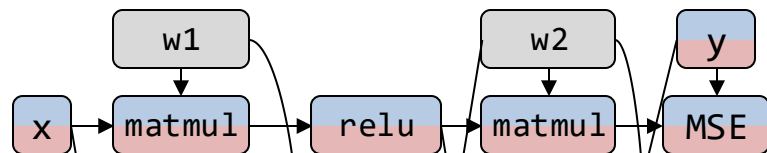
## Matmul Parallelization Type 2

communication cost = all-reduce( $c$ )



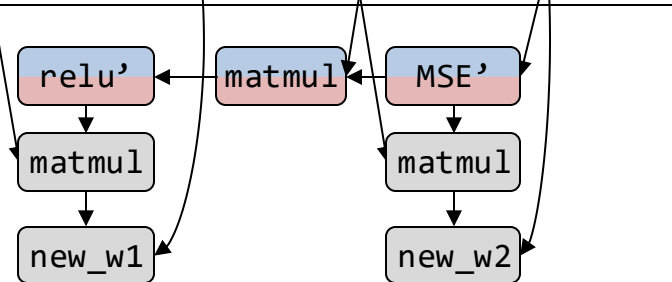
## Forward Pass

Two “Type 1” matmuls: no communication



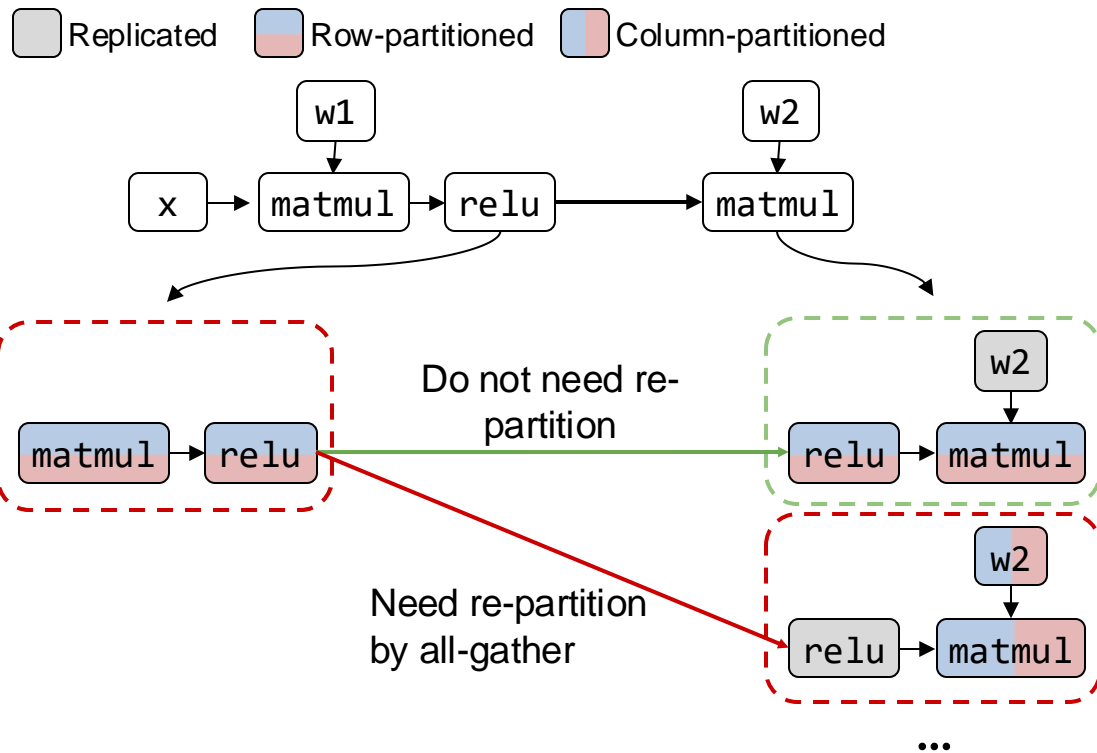
## Backward Pass

One “Type 1” matmul: no communication  
Two “Type 2” matmuls: require all-reduce



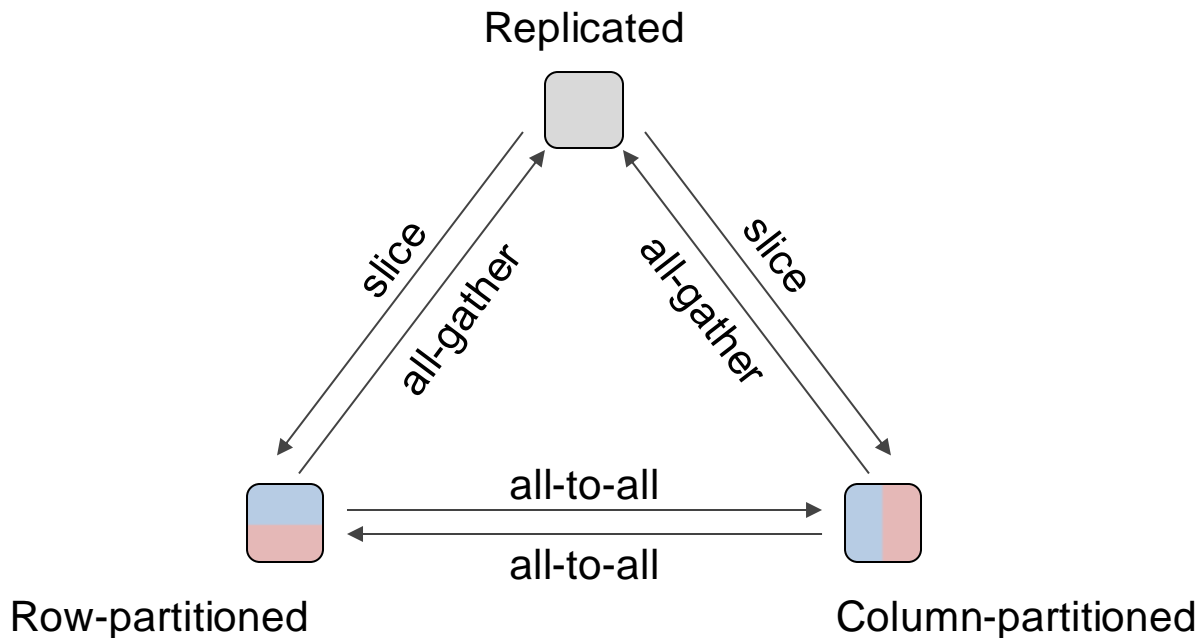
# Re-partition Communication Cost

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

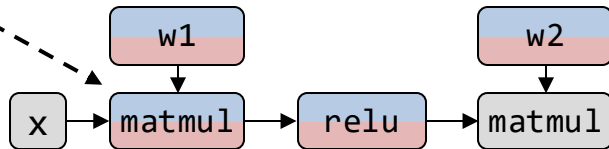




# Parallelize All Operators in a Graph

## 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

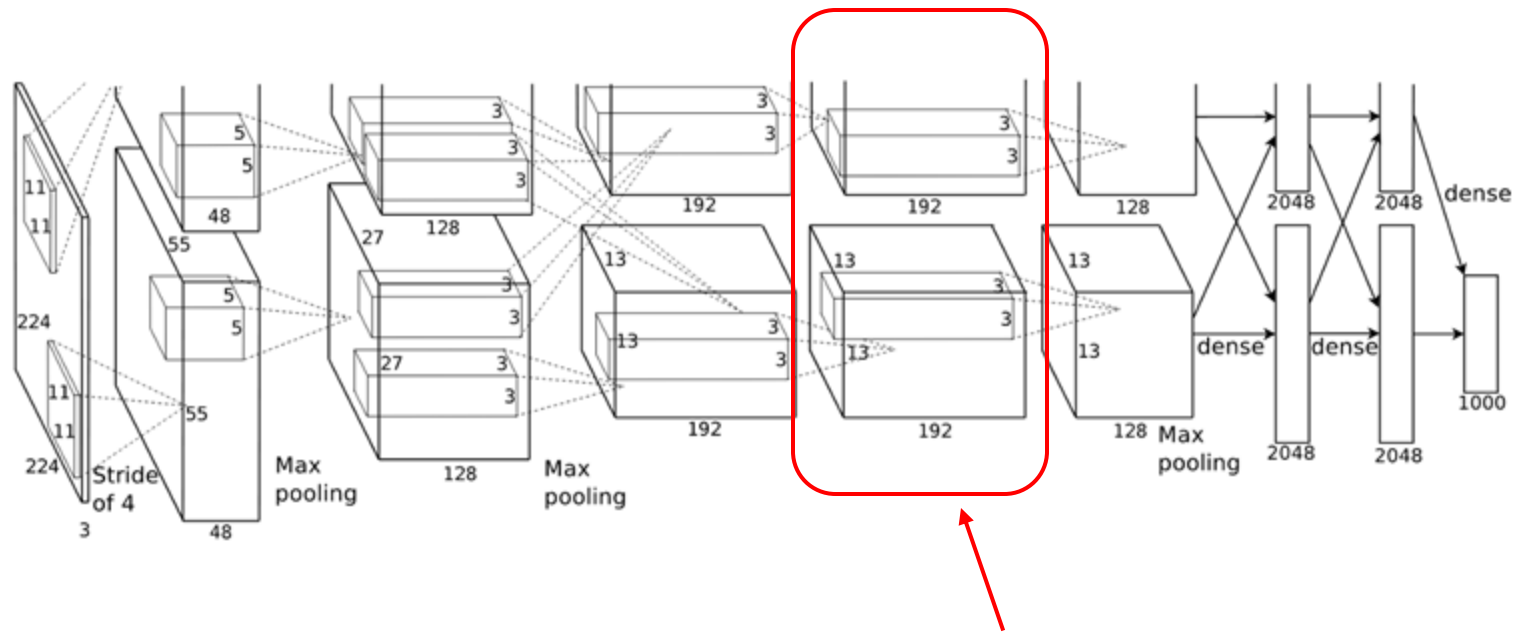
- AlexNet
- Megatron-LM
- GShard MoE

## Systems for Intra-op Parallelism

- ZeRO
- Mesh-Tensorflow
- GSPMD
- Tofu
- FlexFlow

# AlexNet

Result: increase top-1 accuracy by 1.7%

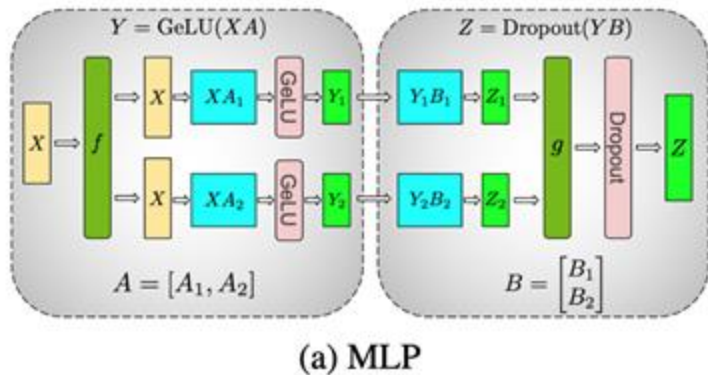


Assign a group convolution layer to 2 GPUs

# Megatron-LM

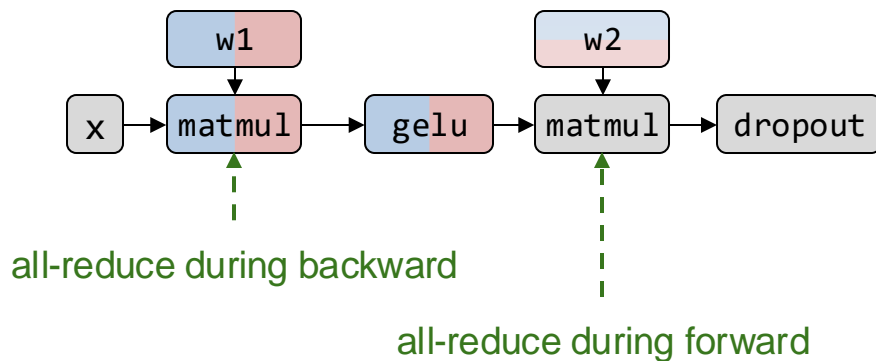
Result: a large language model with 8.3B parameters that outperforms SOTA results

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



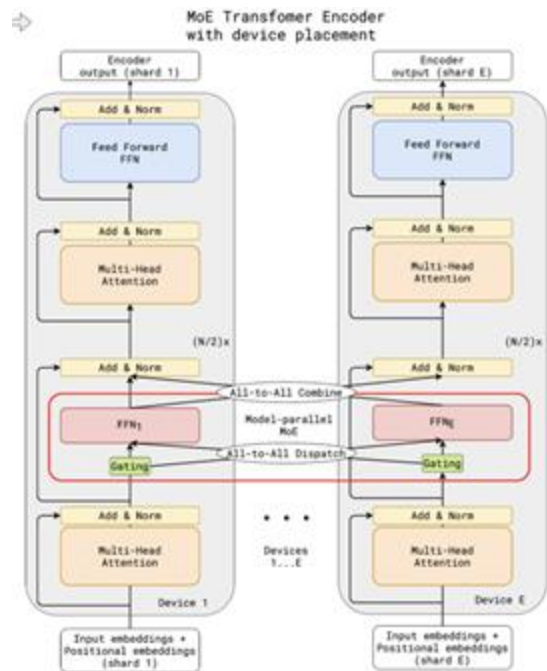
Illustrated with the notations in this tutorial

Replicated Row-partitioned Column-partitioned



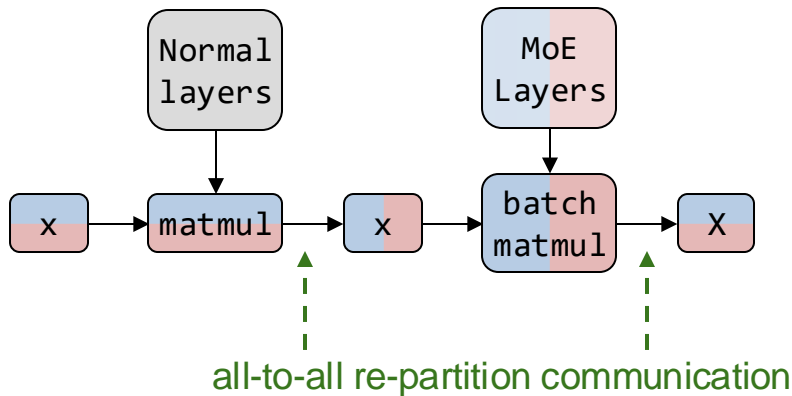
# GShard MoE

Result: a multi-language translation model with 600B parameters that outperforms SOTA



Illustrated with the notations in this class

Replicated Row-partitioned Expert-partitioned



# ZeRO Optimizer

## Problem

Data parallelism replicates gradients, optimizer states and model weights on all devices.

## Idea



Partition gradients, optimizer states and model weights.

M is the number of parameters, N is the number of devices.

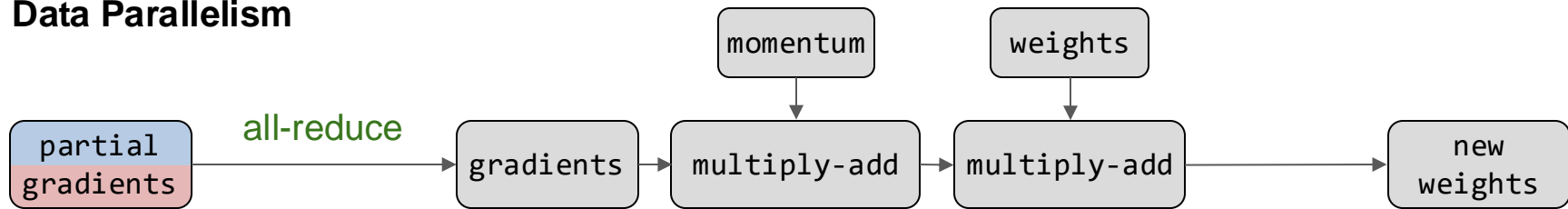
	Optimizer States (12M)	Gradients (2M)	Model Weights (2M)	Memory Cost	Communication Cost
Data Parallelism	Replicated	Replicated	Replicated	$16M$	all-reduce(2M)
ZeRO Stage 1	Partitioned	Replicated	Replicated	$4M + \frac{12M}{N}$	all-reduce(2M)
ZeRO Stage 2	Partitioned	Partitioned	Replicated	$2M + \frac{14M}{N}$	all-reduce(2M)
ZeRO Stage 3	Partitioned	Partitioned	Partitioned	$\frac{16M}{N}$	1.5 all-reduce(2M)

# ZeRO Stage 2

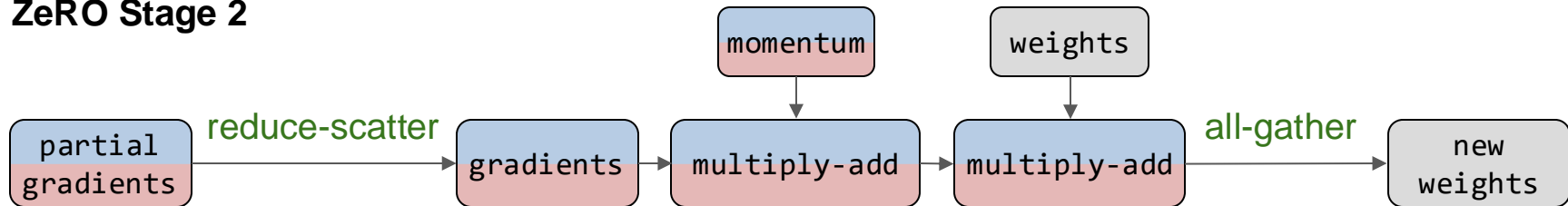
**Key Idea:** all-reduce = reduce-scatter + all-gather

 Replicated  Partitioned

## Data Parallelism



## ZeRO Stage 2



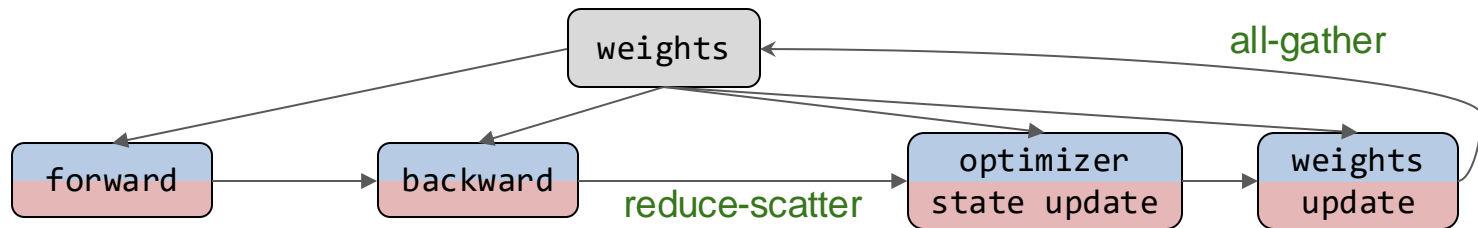
Same communication cost but save memory by partitioning more tensors

# ZeRO Stage 3

Replicated Partitioned

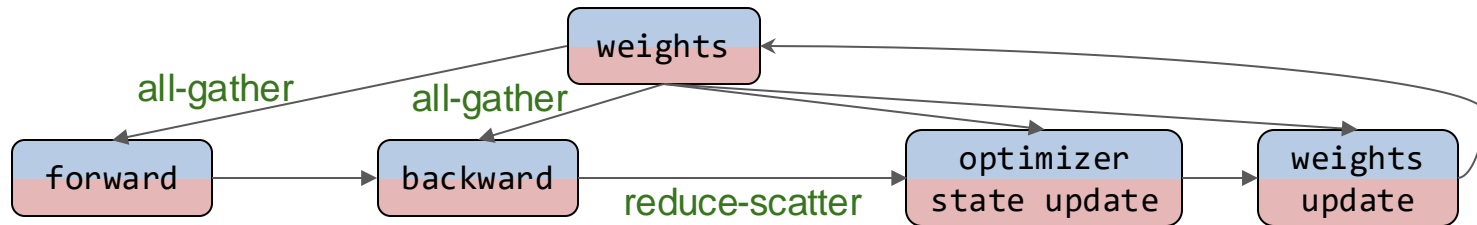
## ZeRO Stage 2

communication cost  
= all-reduce



## ZeRO Stage 3

communication cost  
= 1.5 all-reduce





# Mesh-Tensorflow

Map tensor dimension to mesh dimension for parallelism

```
...  
batch = mtf.Dimension("batch", b)  
io = mtf.Dimension("io", d_io)  
hidden = mtf.Dimension("hidden", d_h)  
# x.shape == [batch, io]  
w = mtf.get_variable("w", shape=[io, hidden])  
bias = mtf.get_variable("bias", shape=[hidden])  
v = mtf.get_variable("v", shape=[hidden, io])  
h = mtf.relu(mtf.einsum(x, w, output_shape=[batch, hidden]) + bias)  
y = mtf.einsum(h, v, output_shape=[batch, io])  
...
```

Tensor dimension

```
mesh_shape = [("rows", r), ("cols", c)]  
computation_layout = [("batch", "rows"), ("hidden", "cols")]
```

Mesh dimension

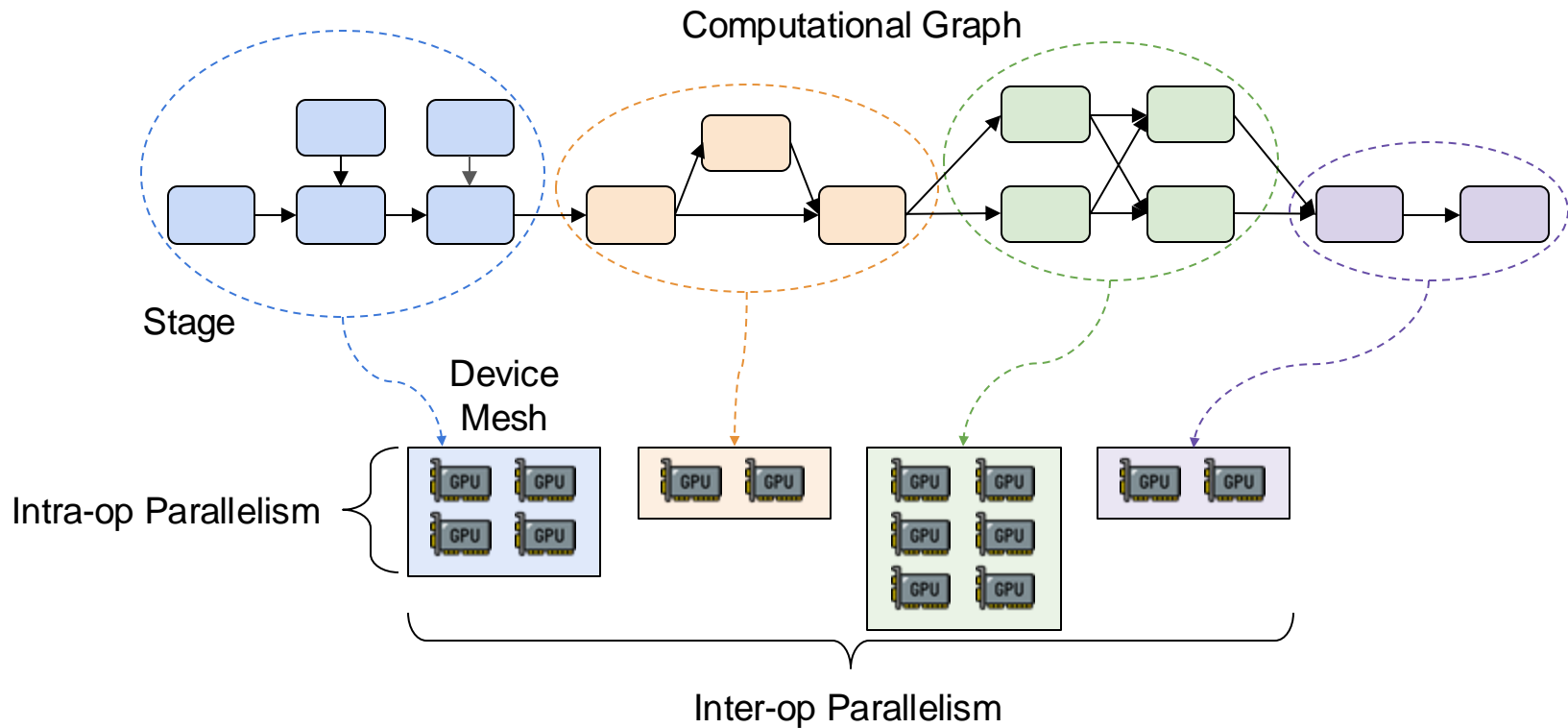
Mapping

# GSPMD

- Use annotations to specify partition strategy
- Propagate the annotations to whole graph
- Use compiler to generate SPMD (Single Program Multiple Data) parallel executables

```
1  # Partition inputs along group (G) dim.
2  + inputs = split(inputs, 0, D)
3  # Replicate the gating weights
4  + wg = replicate(wg)
5  gates = softmax(einsum("GSM,ME->GSE", inputs, wg))
6  combine_weights, dispatch_mask = Top2Gating(gating_logits)
7  dispatched_expert_inputs = einsum(
8    "GSEC,GSM->EGCM", dispatch_mask, reshaped_inputs)
9  # Partition dispatched inputs along expert (E) dim.
10 + dispatched_expert_inputs = split(dispatched_expert_inputs, 0, D)
11 h = einsum("EGCM,EMH->EGCH", dispatched_expert_inputs, wi)
12 ...
```

# Combine Intra-op Parallelism and 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

# Other Techniques for Training Large Models

## System-level Memory Optimizations

- Rematerialization/Gradient Checkpointing
- Swapping

## ML-level Optimizations

- Quantization
- Sparsification
- Low-rank approximation

Chen, Tianqi, et al. "Training deep nets with sublinear memory cost." *arXiv 2016*

Rajbhandari, Samyam, et al. "Zero-infinity: Breaking the gpu memory wall for extreme scale deep learning." *SC 2021*.

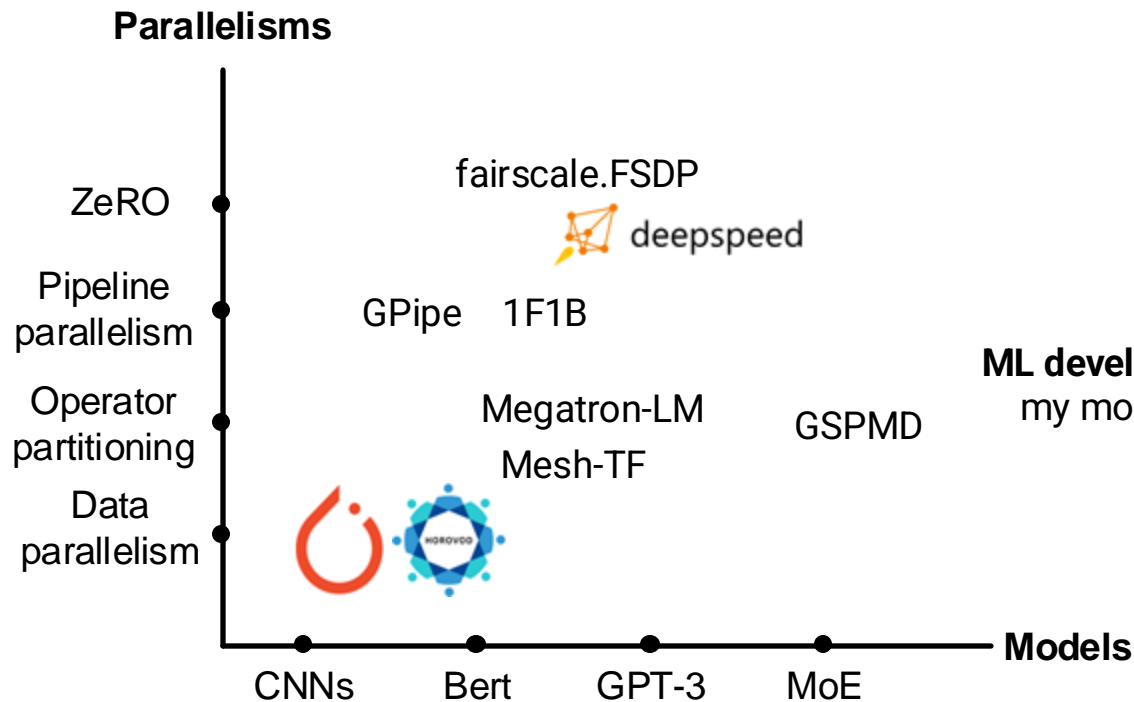
Tang, Hanlin, et al. "1-bit adam: Communication efficient large-scale training with adam's convergence speed." *ICML 2021*.

Shazeer, Noam, and Mitchell Stern. "Adafactor: Adaptive learning rates with sublinear memory cost." *ICML 2018*.

# Where We Are

- Motivation
- History
- Parallelism Overview
- Data parallelism
- Model parallelism
  - Inter-op parallelism
  - Intra-op parallelism
- **Auto-parallelization**

# Auto-parallelization: Motivation



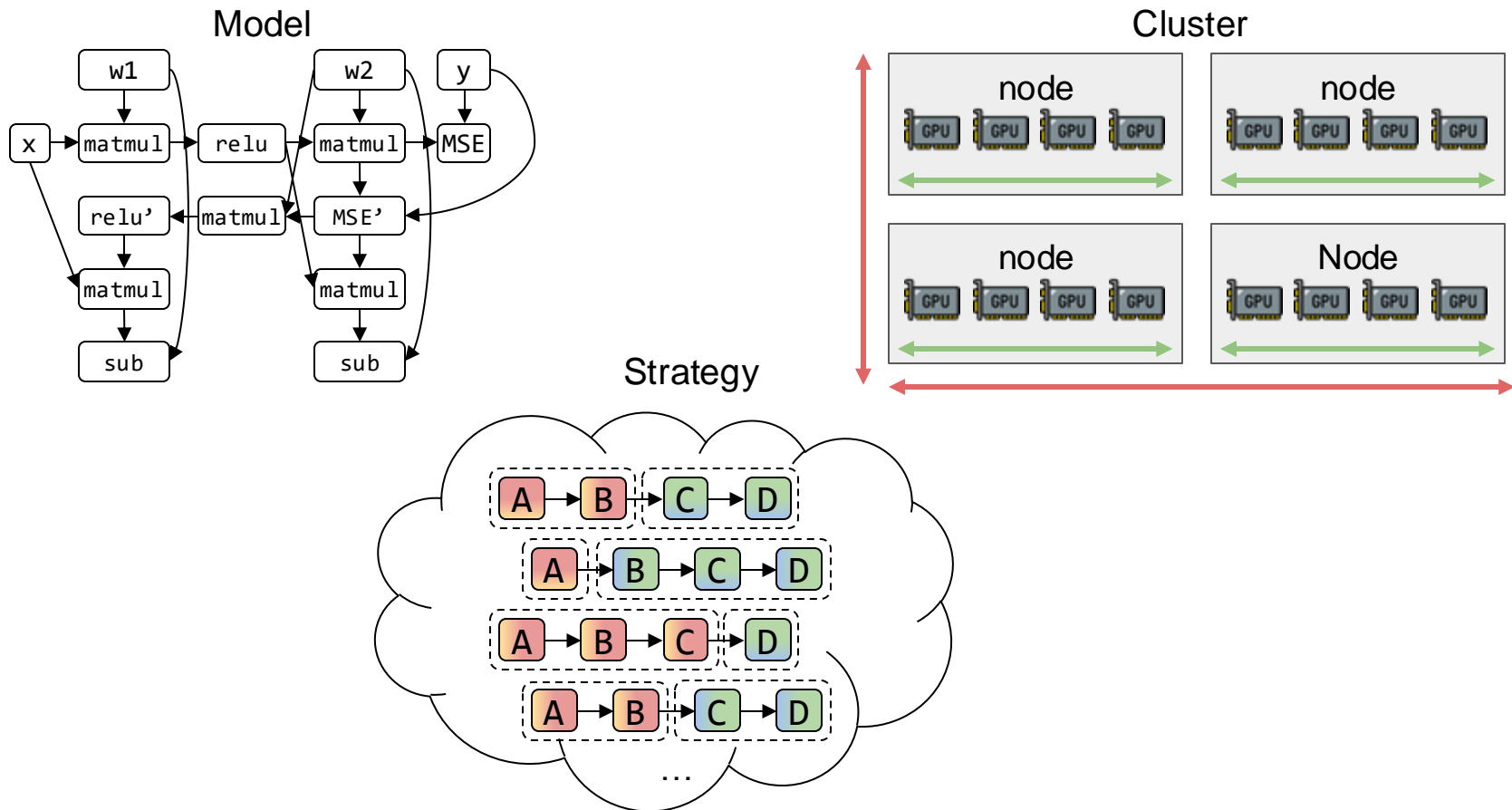
**ML developer:** which one is for my model and my cluster?

# Auto-parallelization: Problem

$$\begin{aligned} & \max_{\text{strategy}} \text{Performance}(\text{Model}, \text{Cluster}) \\ & s.t. \text{ strategy} \in \text{Inter-op} \cup \text{Intra-op} \end{aligned}$$



# Auto-parallelization: Problem



# The Search Space is Huge

**#ops in a real model  
(nodes to color)**

**100 - 10K**

**#op types  
(type of nodes)**

**80 - 200+**

**#devices on a cluster  
(available colors)**

**10s - 1000s**

# Automatic Parallelization Methods

## Search-based methods

- MCMC:
  - [Jia et al., 2018]
  - [Jia et al., 2019]
- Heuristics
  - [Fan et al., 2021]

The complete list of  
references is available  
on the tutorial website

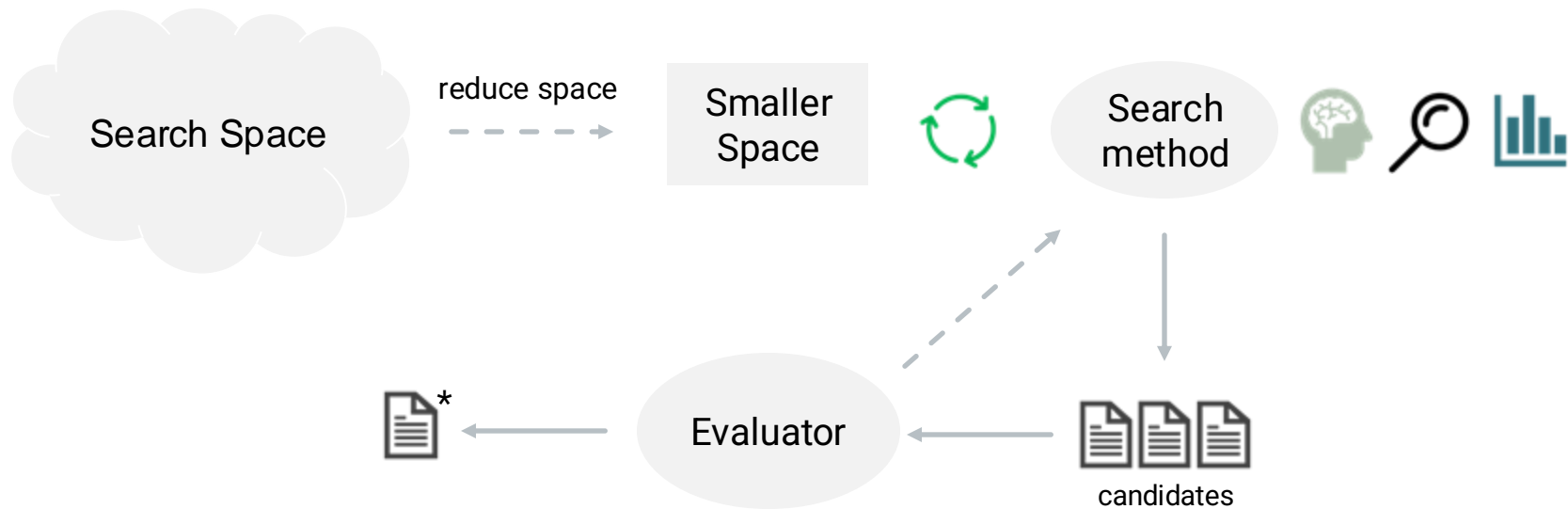
## Learning-based methods

- Reinforcement Learning:
  - [Mirhoseini et al., 2017]
  - [Mirhoseini et al., 2018]
  - [Addanki, et al., 2019]
- ML-based cost model:
  - [Chen et al., 2018],
  - [Zhou et al., 2020],
  - [Zhang, 2020]
- Bayesian optimization:
  - [Sergeev et al., 2018],
  - [Peng et al., 2019]

## Optimization-based methods

- Dynamic programming
  - [Wang, et al., 2018]
  - [Narayanan, et al., 2019]
  - [Li, et al., 2021]
  - [Narayanan, et al., 2012]
  - [Tarnawski, et al., 2020]
  - [Tarnawski, et al., 2021]
- Integer linear programming
  - [Tarnawski, et al., 2020]
- Hierarchical Optimization
  - [Zheng, et al., 2022]

# General Recipe



# Automatic Parallelization Methods

## Search-based methods

- MCMC:
  - [Jia et al., 2018]
  - [Jia et al., 2019]
- Heuristics
  - [Fan et al., 2021]

The complete list of references is available on the tutorial website

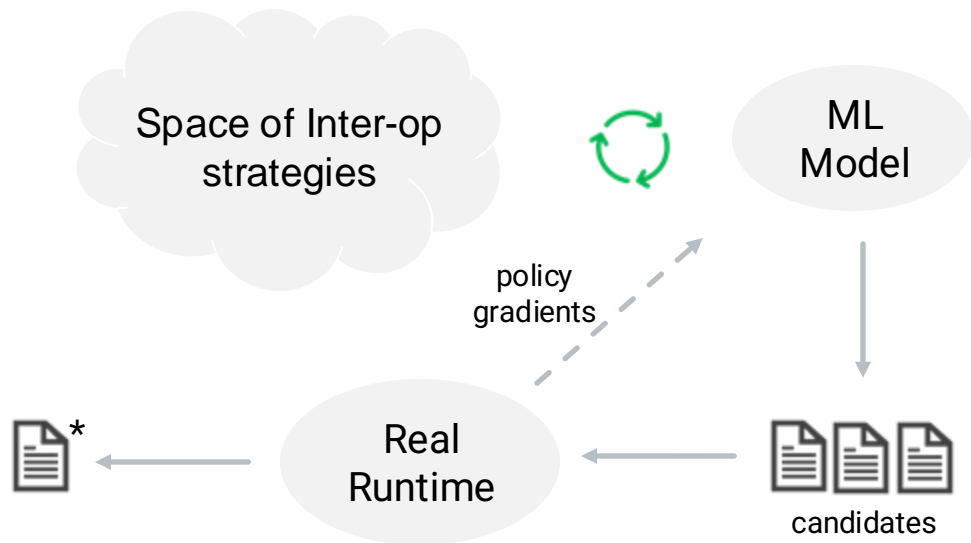
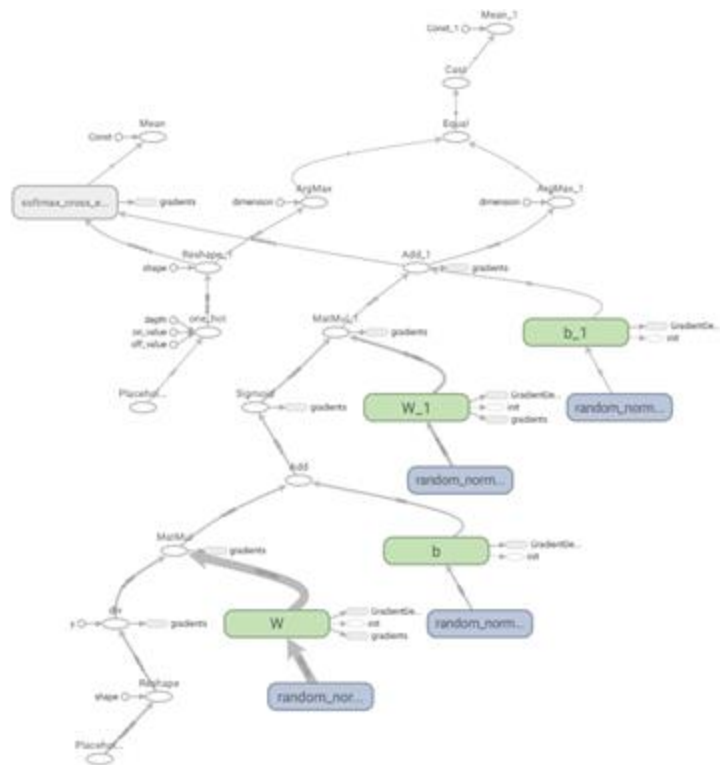
## Learning-based methods

- Reinforcement Learning:
  - **[Mirhoseini et al., 2017]**
  - [Mirhoseini et al., 2018]
  - [Addanki, et al., 2019]
- ML-based cost model:
  - [Chen et al., 2018],
  - [Zhou et al., 2020],
  - [Zhang, 2020]
- Bayesian optimization:
  - [Sergeev et al., 2018],
  - [Peng et al., 2019]

## Optimization-based methods

- Dynamic programming
  - [Wang, et al., 2018]
  - [Narayanan, et al., 2019]
  - [Li, et al., 2021]
  - [Narayanan, et al., 2012]
  - [Tarnawski, et al., 2020]
  - [Tarnawski, et al., 2021]
- Integer linear programming
  - [Tarnawski, et al., 2020]
- Hierarchical optimization
  - **[Zheng, et al., 2022]**

# ColocRL (a.k.a. Device Placement Optimization)



# ColocRL: Model

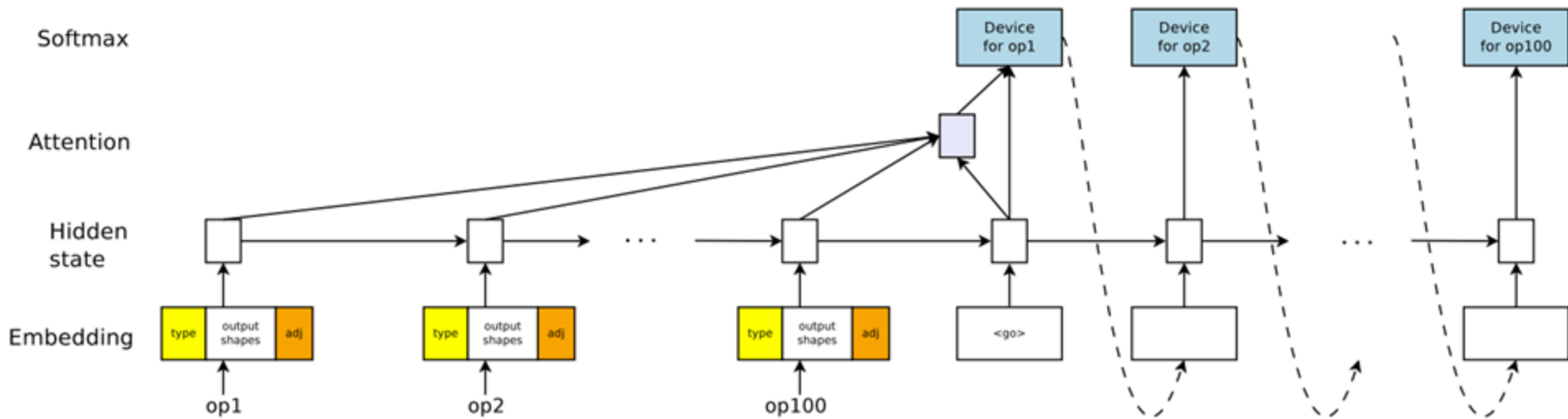


Figure from [Mirhoseini et al., ICML 2017]

# ColocRL: Training

$$\mathbb{E}_{\mathcal{P} \sim \pi(\mathcal{P} \mid \mathcal{G}; \theta)} [R(\mathcal{P}) \mid \mathcal{G}]$$

$\mathcal{G}$ : computational graph

$\mathcal{R}(\mathcal{P})$ : Real runtime of a placement

$\pi(\cdot)$ : output distributed of the RNN



# ColocRL: Other Improvement

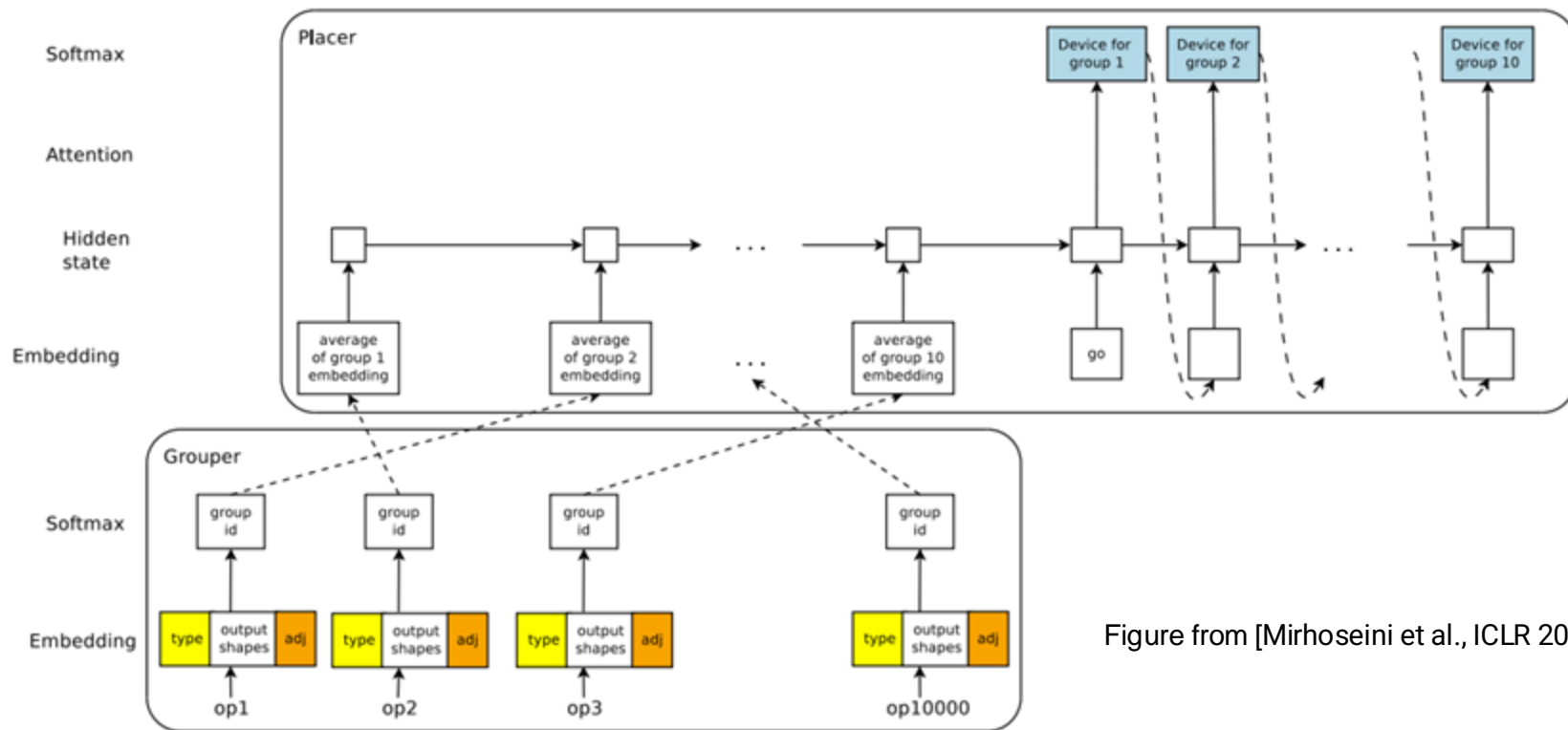
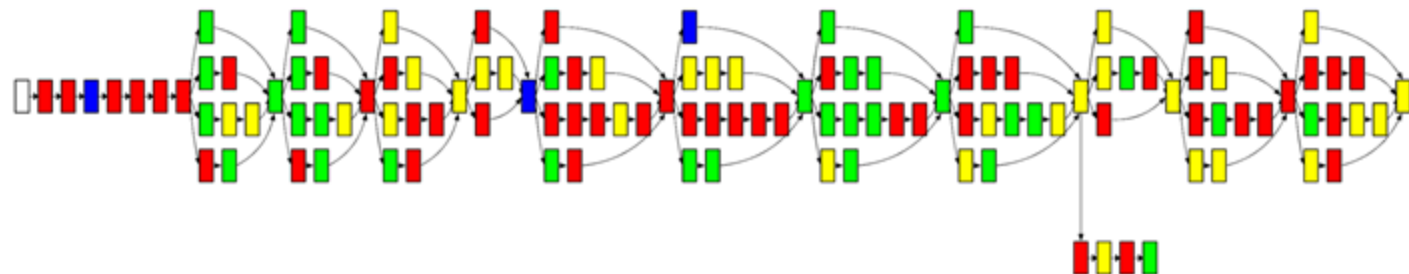


Figure from [Mirhoseini et al., ICLR 2018]

# Results Discussion



Tasks	Single-CPU	Single-GPU	#GPUs	Scotch	MinCut	Expert	RL-based	Speedup
RNNLM (batch 64)	6.89	<b>1.57</b>	2	13.43	11.94	3.81	<b>1.57</b>	0.0%
			4	11.52	10.44	4.46	<b>1.57</b>	0.0%
NMT (batch 64)	10.72	OOM	2	14.19	11.54	4.99	<b>4.04</b>	23.5%
			4	11.23	11.78	4.73	<b>3.92</b>	20.6%
Inception-V3 (batch 32)	26.21	<b>4.60</b>	2	25.24	22.88	11.22	<b>4.60</b>	0.0%
			4	23.41	24.52	10.65	<b>3.85</b>	19.0%

Figure and table from [Mirhoseini et al., ICML 2017]

# Automatic Parallelization Methods

## Search-based methods

- MCMC:
  - [Jia et al., 2018]
  - [Jia et al., 2019]
- Heuristics
  - [Fan et al., 2021]

The complete list of references is available on the tutorial website

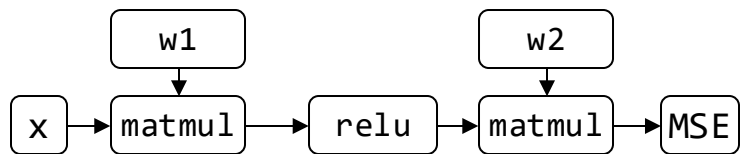
## Learning-based methods

- Reinforcement Learning:
  - [Mirhoseini et al., 2017]
  - [Mirhoseini et al., 2018]
  - [Addanki, et al., 2019]
- ML-based cost model:
  - [Chen et al., 2018],
  - [Zhou et al., 2020],
  - [Zhang, 2020]
- Bayesian optimization:
  - [Sergeev et al., 2018],
  - [Peng et al., 2019]

## Optimization-based methods

- Dynamic programming
  - [Wang, et al., 2018]
  - [Narayanan, et al., 2019]
  - [Li, et al., 2021]
  - [Narayanan, et al., 2012]
  - [Tarnawski, et al., 2020]
  - [Tarnawski, et al., 2021]
- Integer linear programming
  - [Tarnawski, et al., 2020]
- Hierarchical optimization
  - **Alpa [Zheng, et al., 2022]**

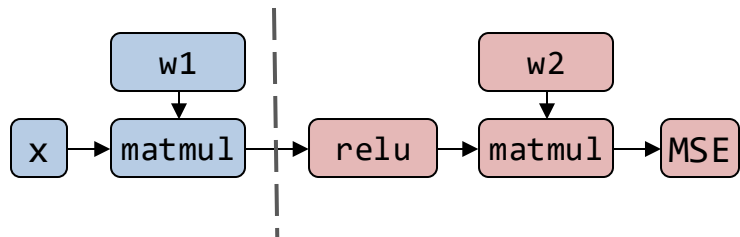
# Optimization-based Method: Alpa



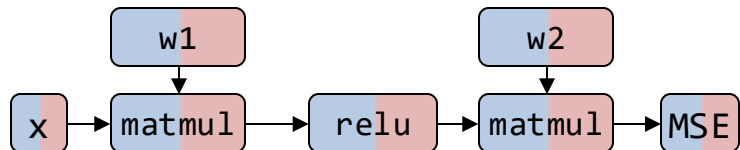
 Device 1

 Device 2

## Inter-op parallelism



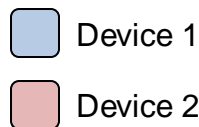
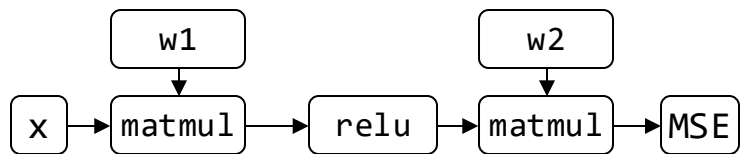
## Intra-op parallelism



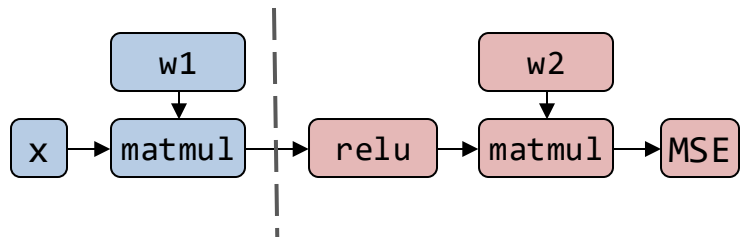
## Trade-off

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

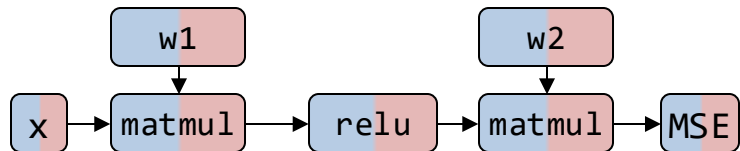
# Alpa Rationale



## Inter-op parallelism

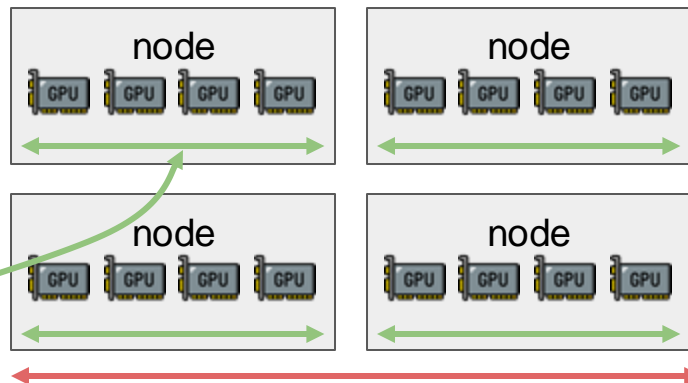


## Intra-op parallelism



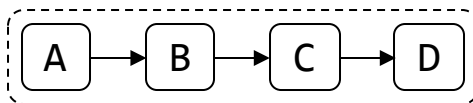
Fast connections (green double arrow)

Slow connections (red double arrow)

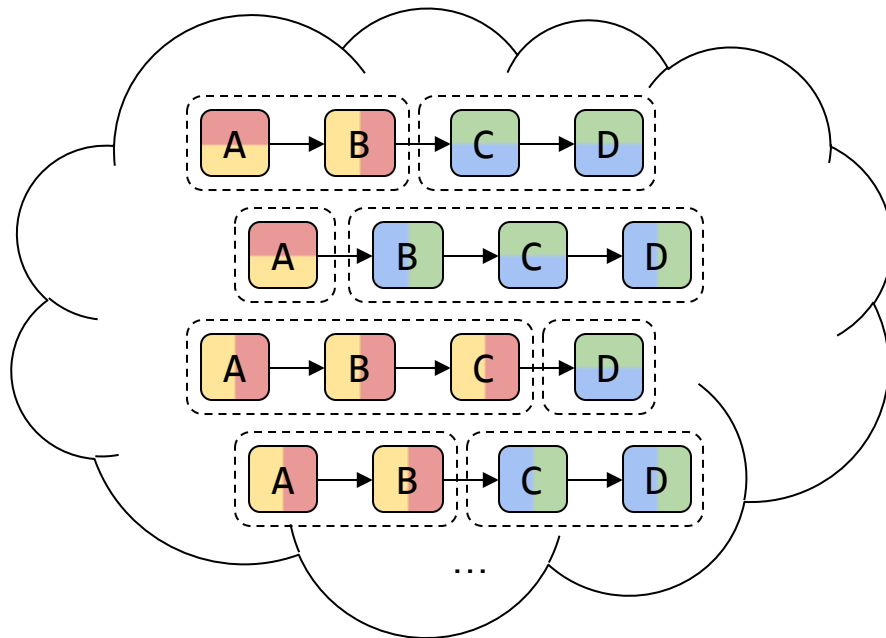


# Search Space

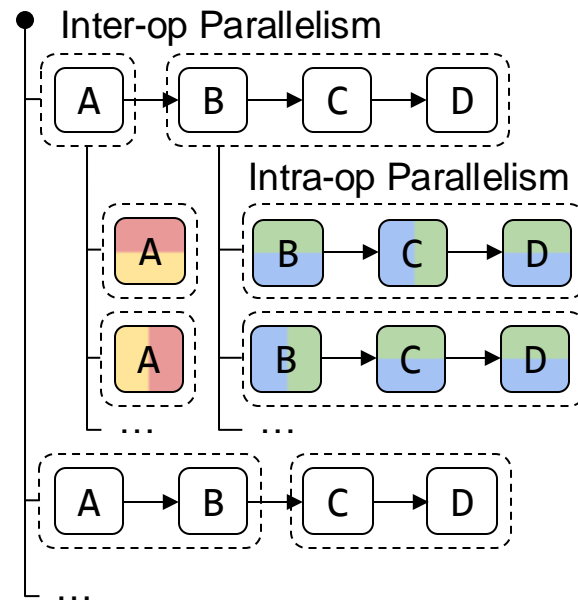
## Computational Graph



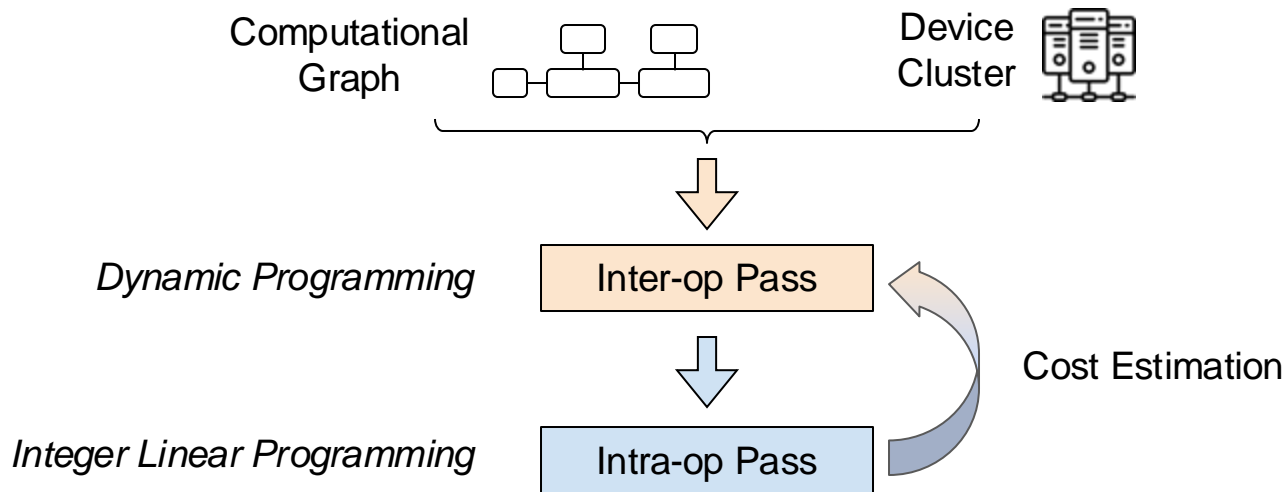
## Whole Search Space



## Alpha Hierarchical Space

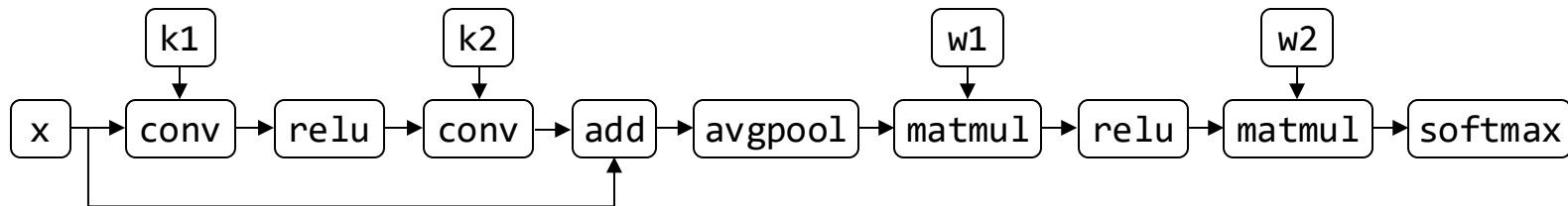


# Alpa Compiler: Hierarchical Optimization



## Inter-op Pass

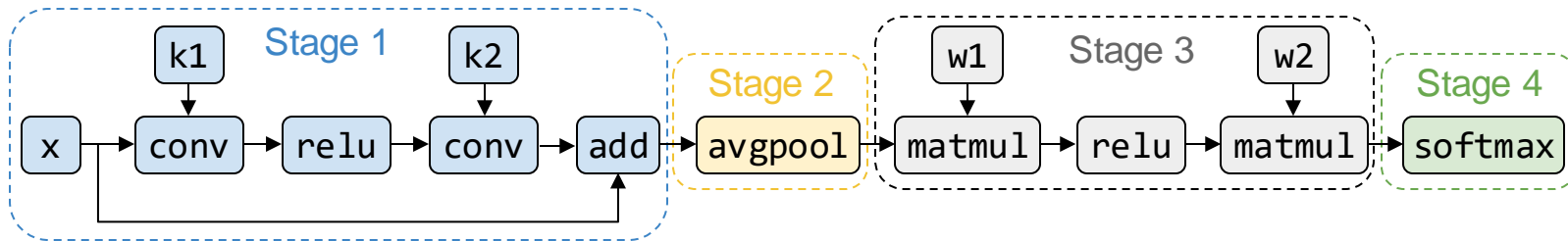
### Computational Graph



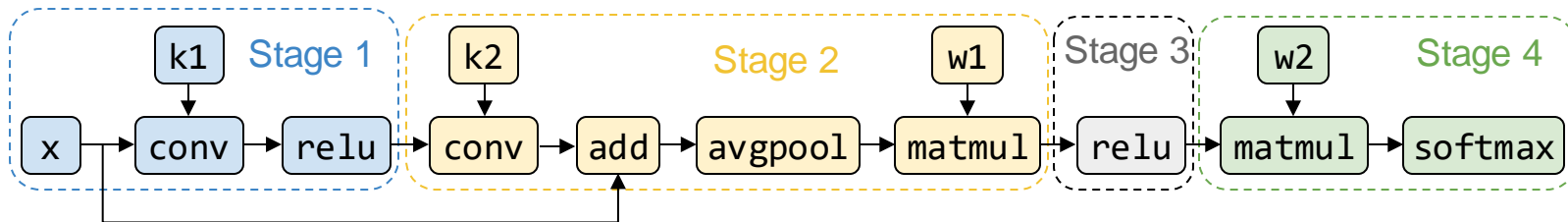


## Inter-op Pass

Graph Partitioning



or

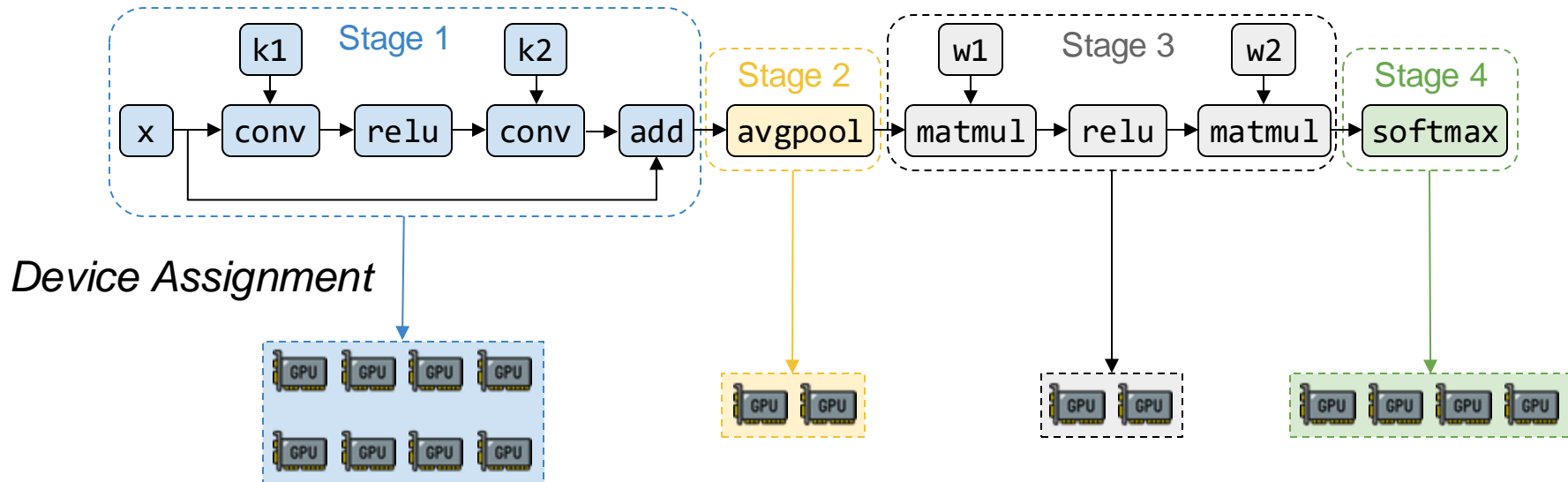


or

...

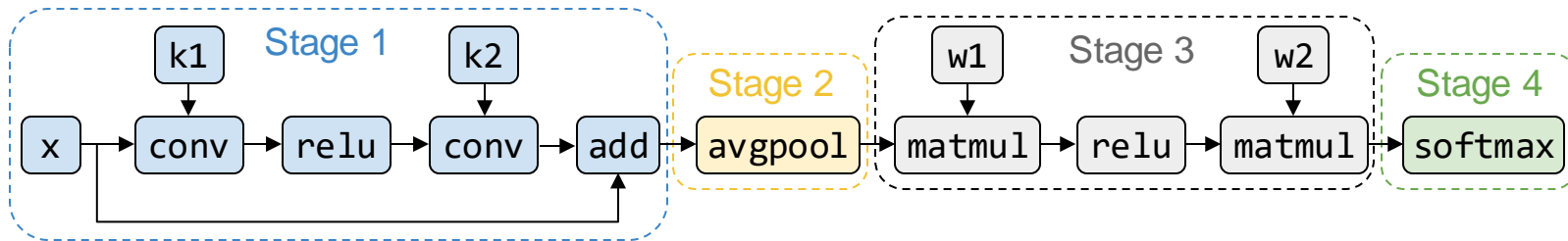
# Inter-op Pass

## Partitioned Computational Graph

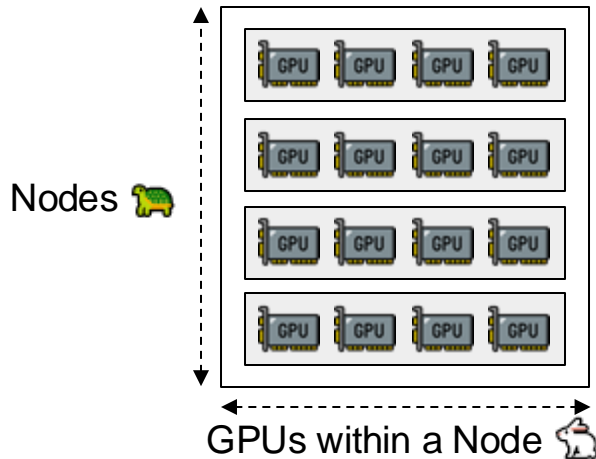


# Inter-op Pass

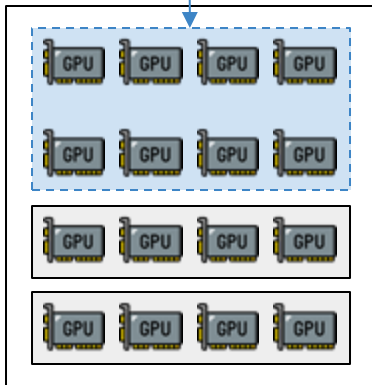
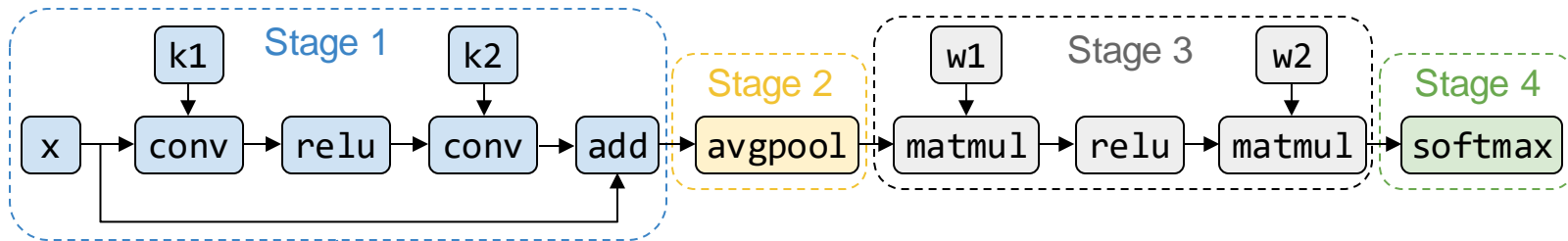
## Partitioned Computational Graph



## Cluster (2D Device Mesh)

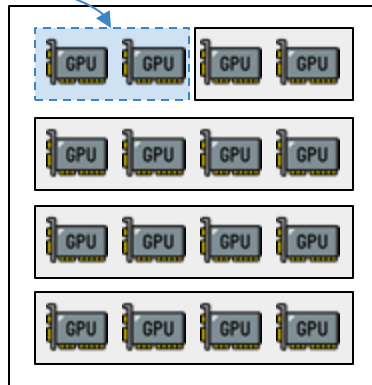


# Inter-op Pass



Submesh Choice 1

or

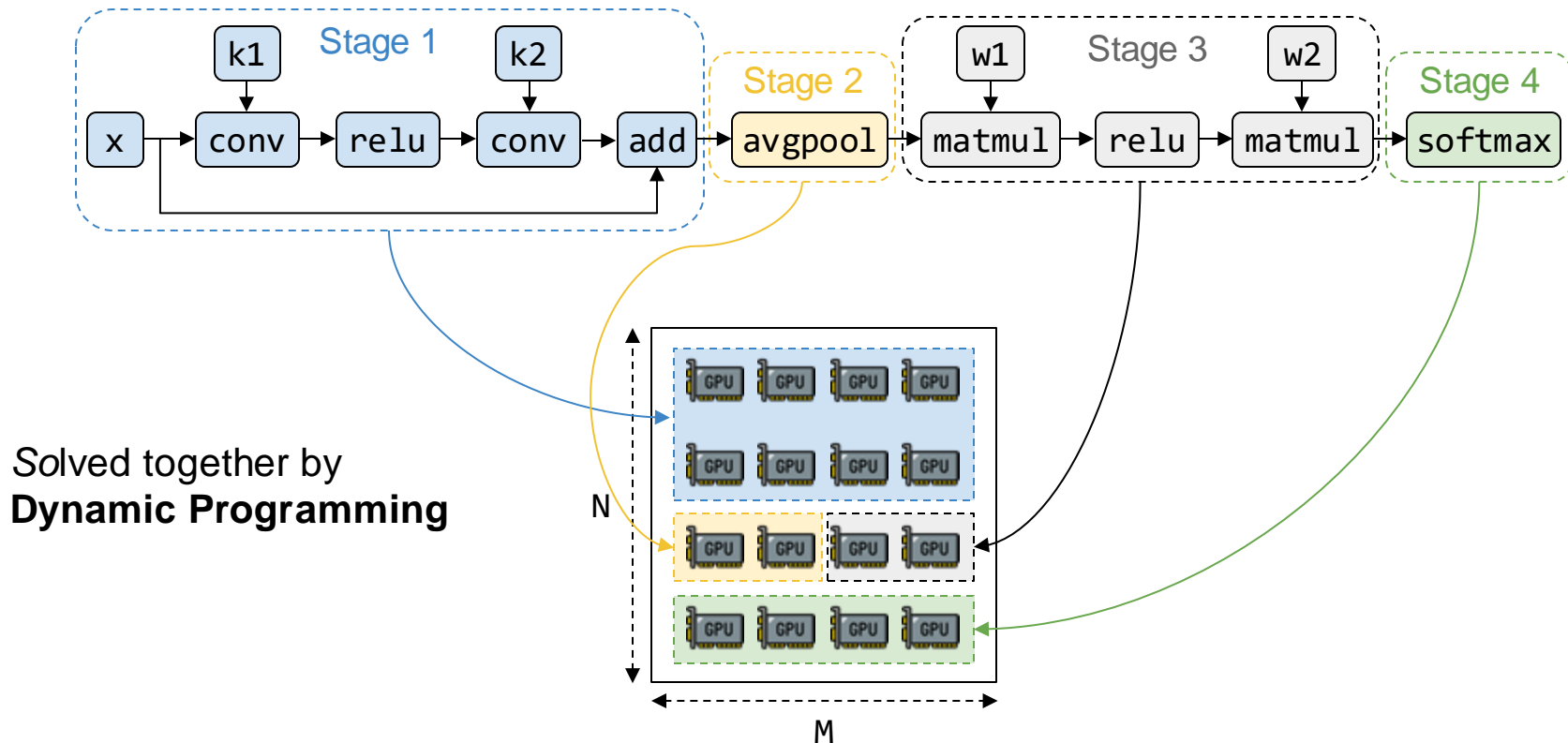


Submesh Choice 2

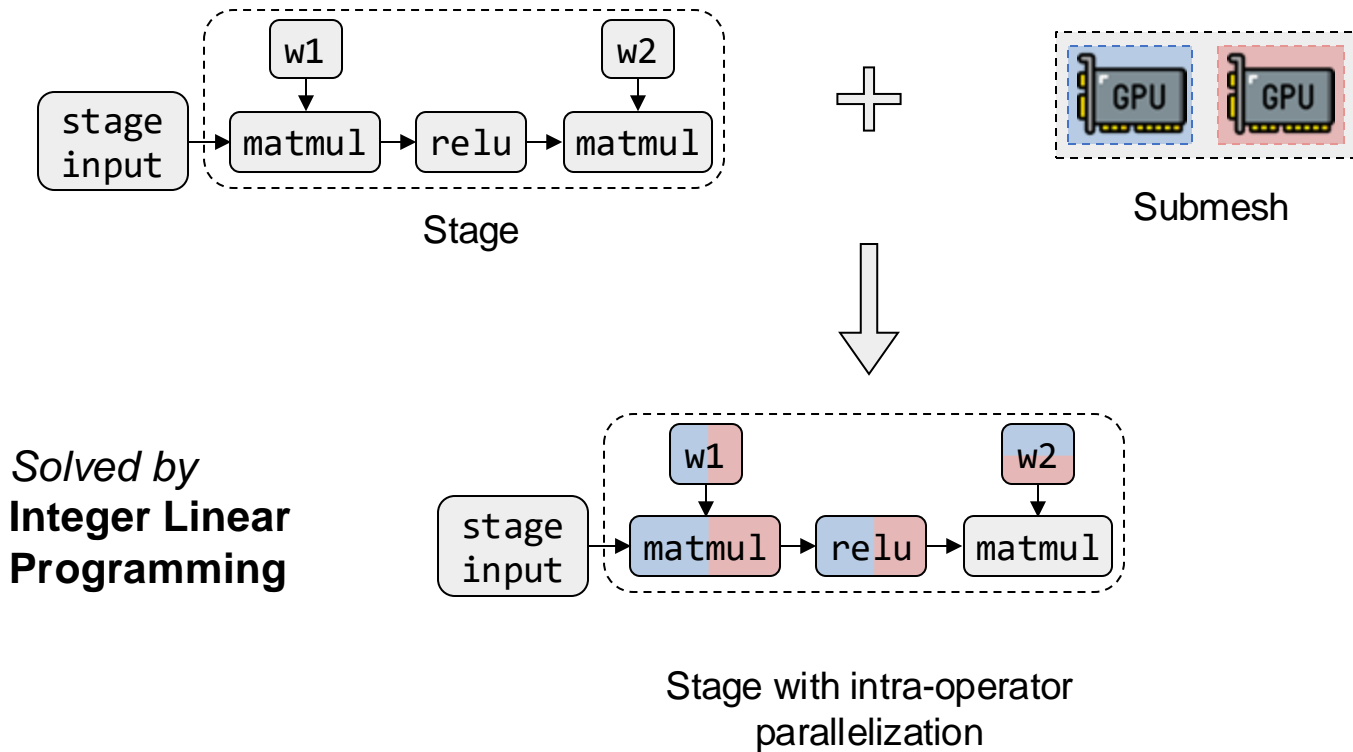
or

...

# Inter-op Pass



# Intra-op Pass

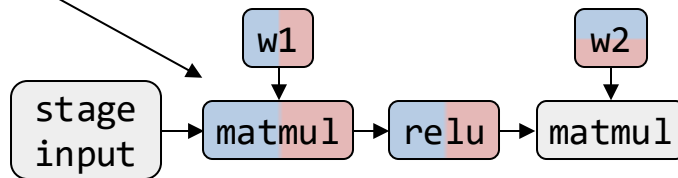


## Intra-op Pass

### Integer Linear Programming Formulation

#### Decision vector

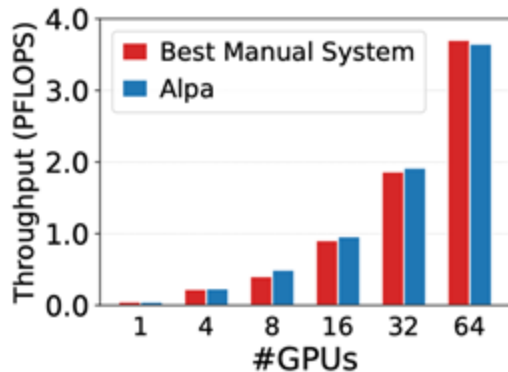
Parallel strategies of each operator



**Minimize** Computation cost + Communication cost

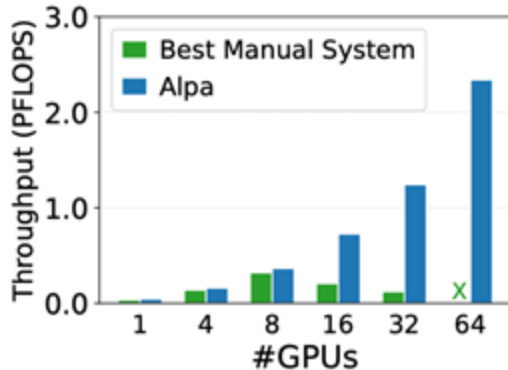
# Evaluation: Comparing with Previous Works

## GPT (up to 39B)



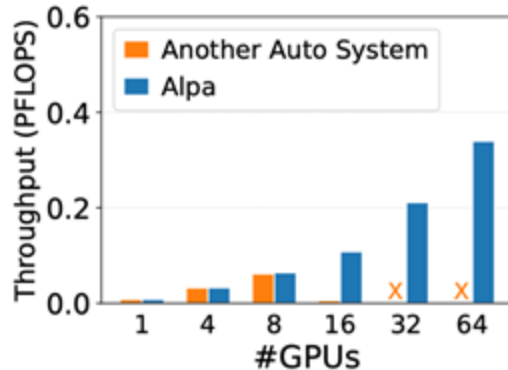
Match specialized manual systems.

## GShard MoE (up to 70B)



Outperform the manual baseline by up to 8x.

## Wide-ResNet (up to 13B)



Generalize to models without manual plans.

*Weak scaling results where the model size grow with #GPUs.*

*Evaluated on 8 AWS EC2 p3.16xlarge nodes with 8 16GB V100s each (64 GPUs in total).*



# Automatic Parallelization Methods

## Search-based methods

- ✓ Easy to extend the search space
- ✓ No training cost
- ✗ High inference cost
- ✗ Not explainable
- ✗ No optimality guarantee

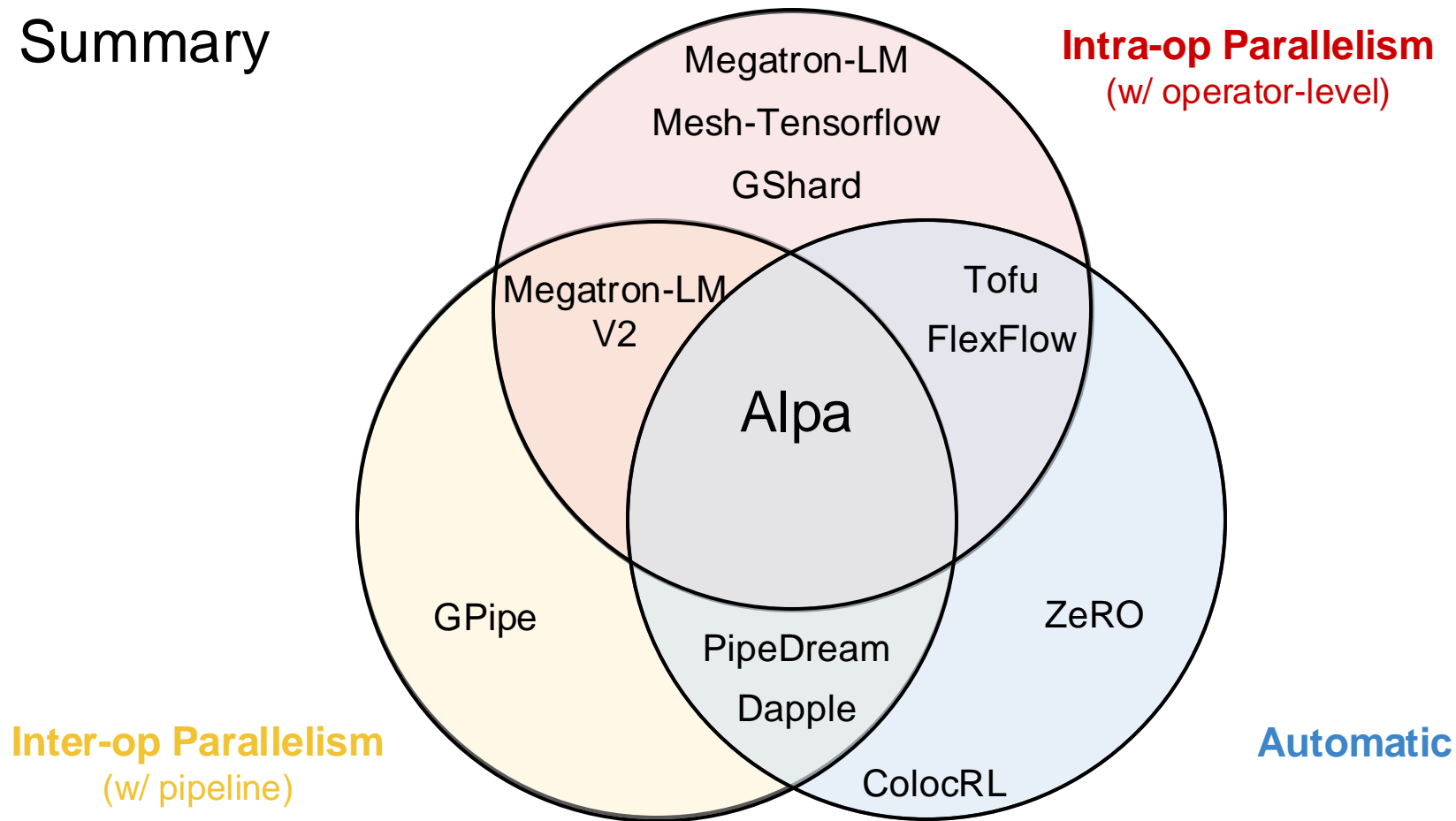
## Learning-based methods

- ✓ Easy to extend the search space
- ✗ High training cost
- ✓ Low inference cost
- ✗ Not explainable
- ✗ No optimality guarantee

## Optimization-based methods

- ✗ Non-trivial to extend the search space
- ✓ No training cost
- ✓ Medium inference cost
- ✓ Explainable
- ✓ Some optimality guarantee

# Summary



# Summary: How to Choose Parallelism

1. Use automatic compiler if not transformer
2. Manual parallelism search for transformers:
  - Factors to consider
    - #GPUs you have
    - Model size
    - JCT (Job completion time)
    - Communication bandwidth
    - etc.

# Hao's Ultimate Guide

