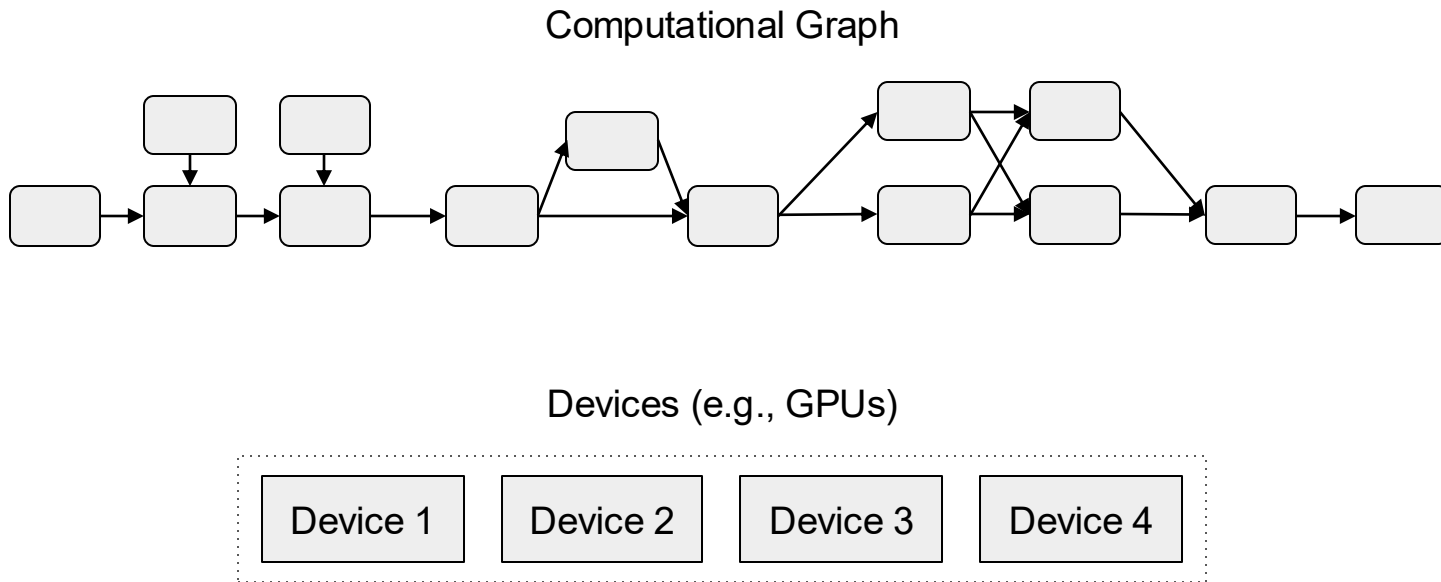


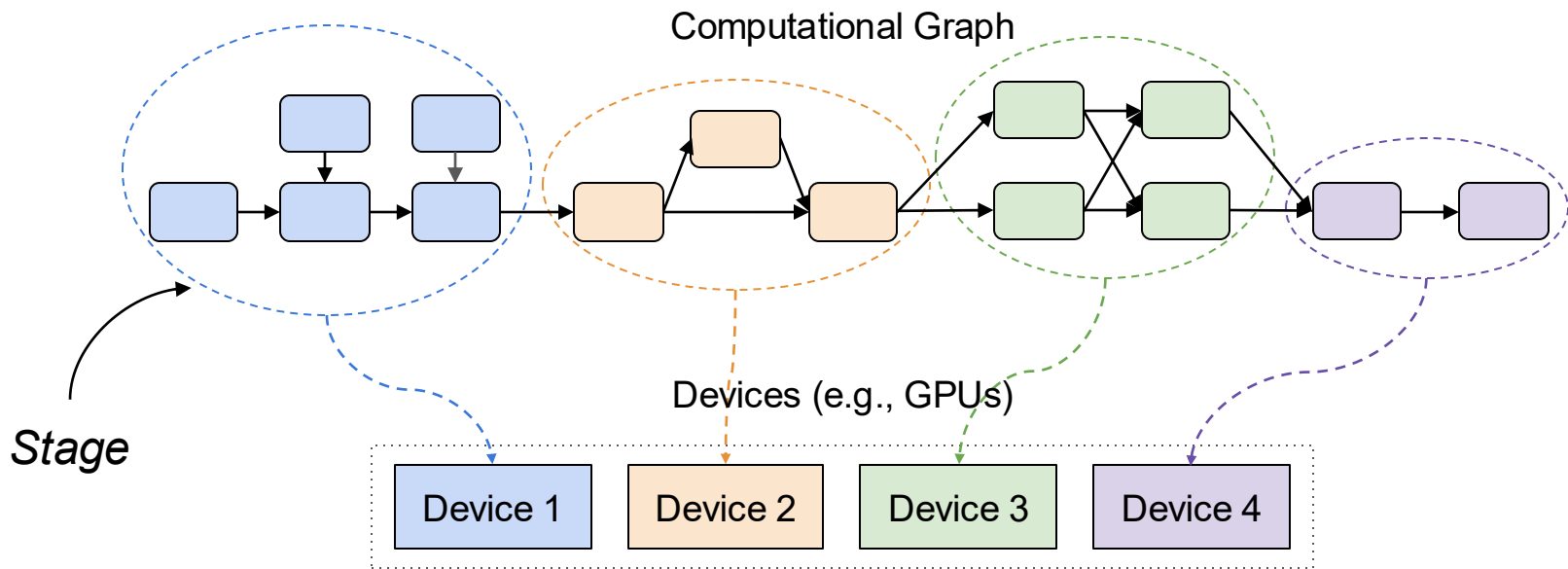
Model Parallelism

- Model parallelism
 - **Inter-op parallelism**
 - Intra-op parallelism

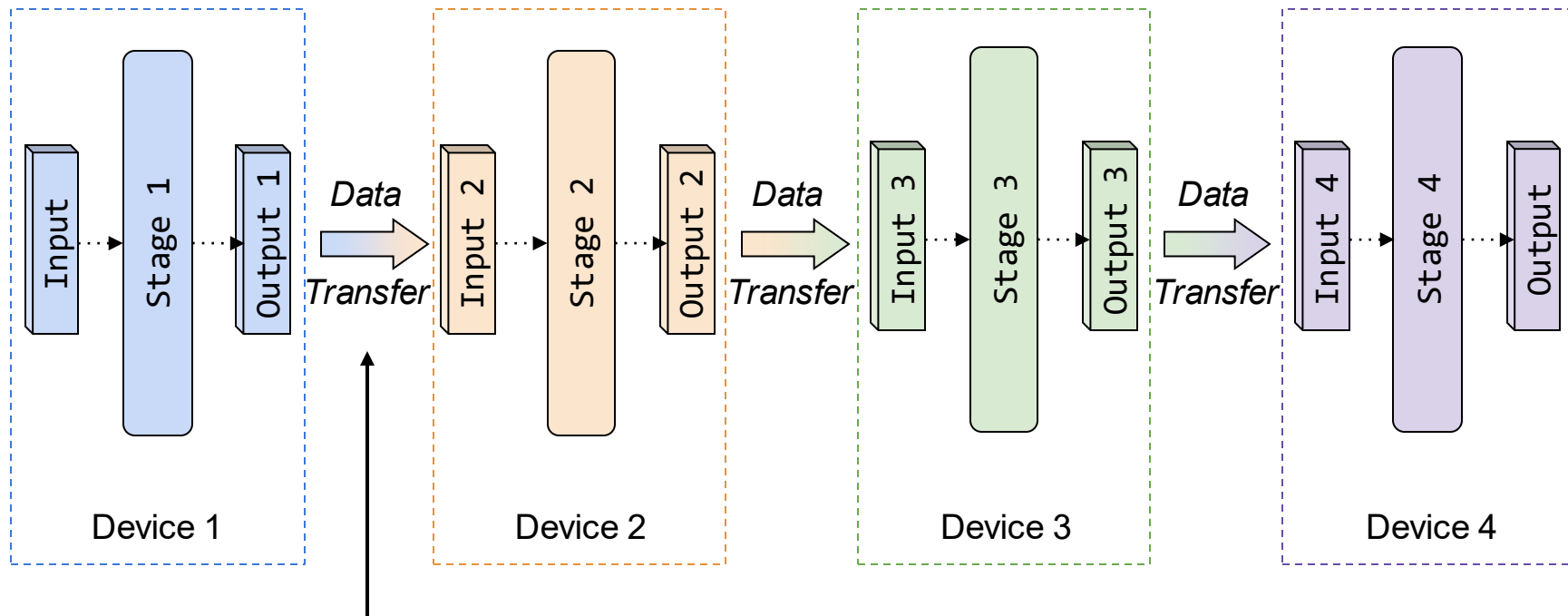
Computational Graph (Neural Networks) → Stages



Computational Graph (Neural Networks) → Stages

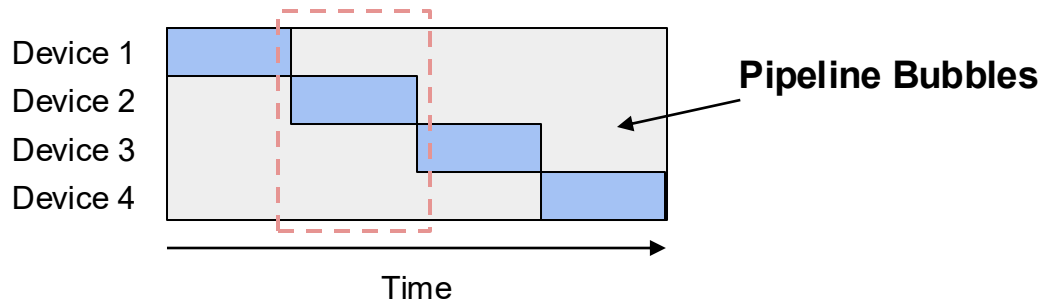



Execution & Data Movement



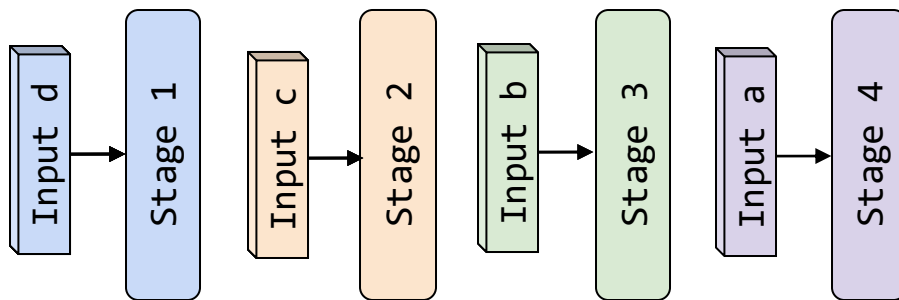
Note: The time spent on data transfer is typically **small**, since we only communicate stage outputs at stage boundaries between two stages.

Timeline: Visualization of Inter-Operator Parallelism

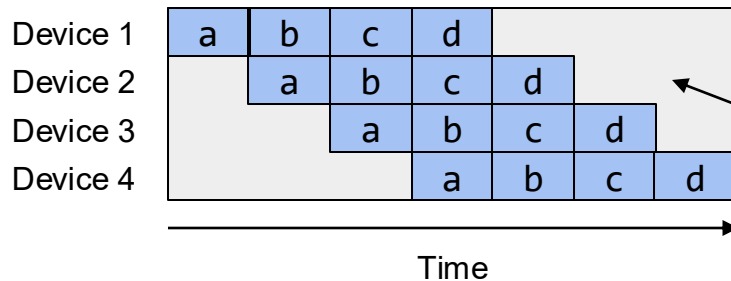


- Gray area ( indicates devices being idle (a.k.a. Pipeline bubbles).
- Only 1 device activated at a time.
- **Pipeline bubble percentage** = $\text{bubble_area} / \text{total_area}$
= $(D - 1) / D$, assuming D devices.

Reduce Pipeline Bubbles via Pipelining Inputs

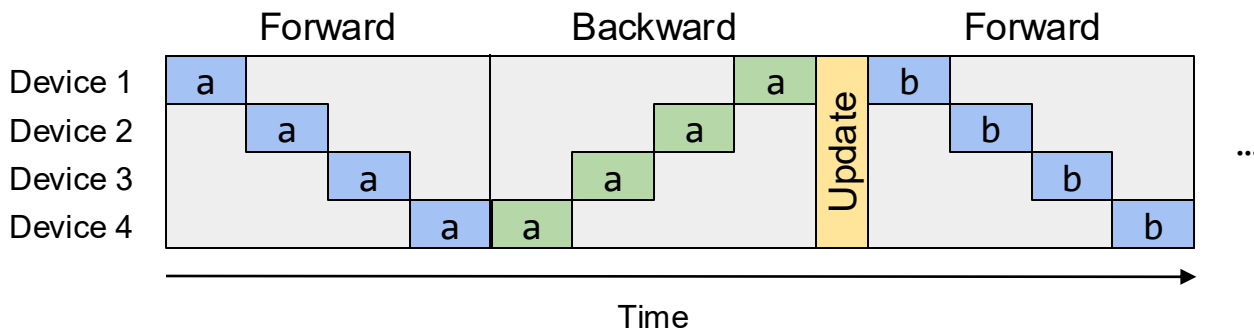
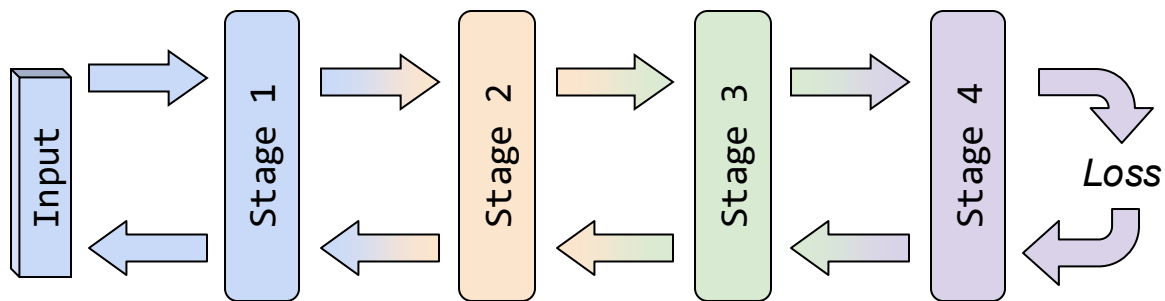


Used in inference.



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

Training: Forward & Backward Dependency

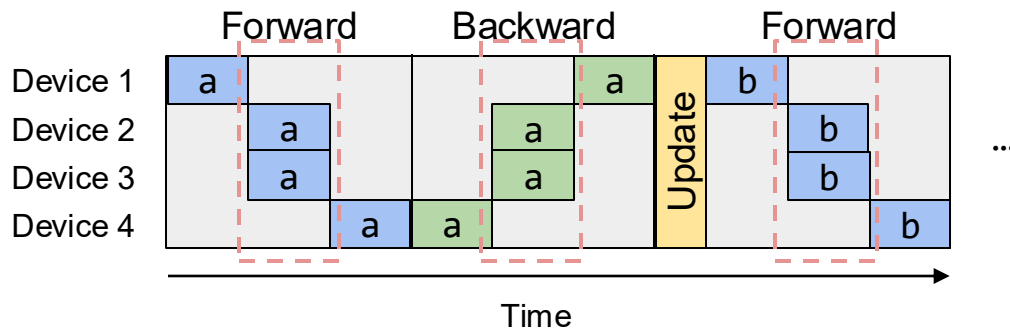
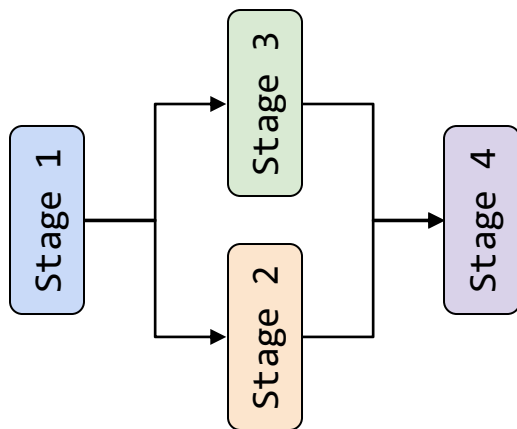


How to Reduce Pipeline Bubbles for Training?

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

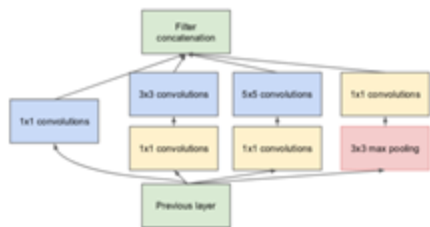
Device Placement

Idea: Slice the branches of a neural network into multiple stages so they can be calculated concurrently.

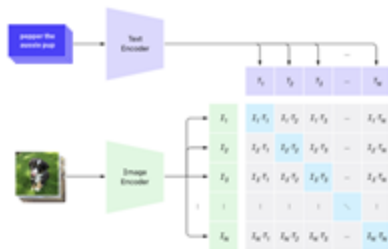


Device Placement: Limitations

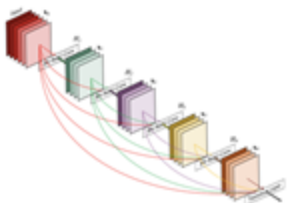
Only works for specific NNs with branches:



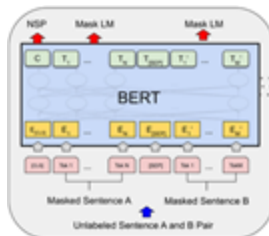
✓ Inception Module



✓ Contrastive Model

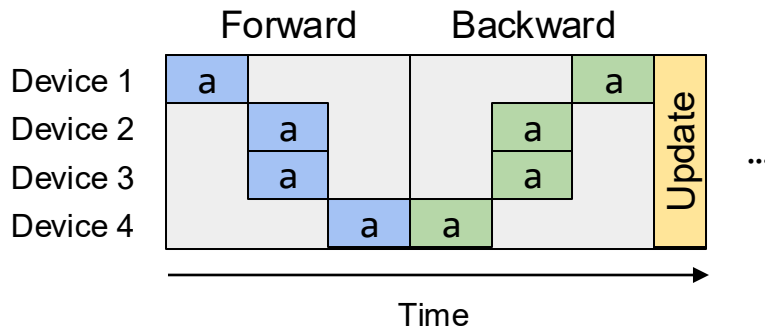


✗ Other ConvNets



✗ Transformers

Device Utilization is still low:



Note: device placement needs to be combined with the other pipeline schedules discussed later to further improve device utilization.

Synchronous Pipeline Parallel Schedule

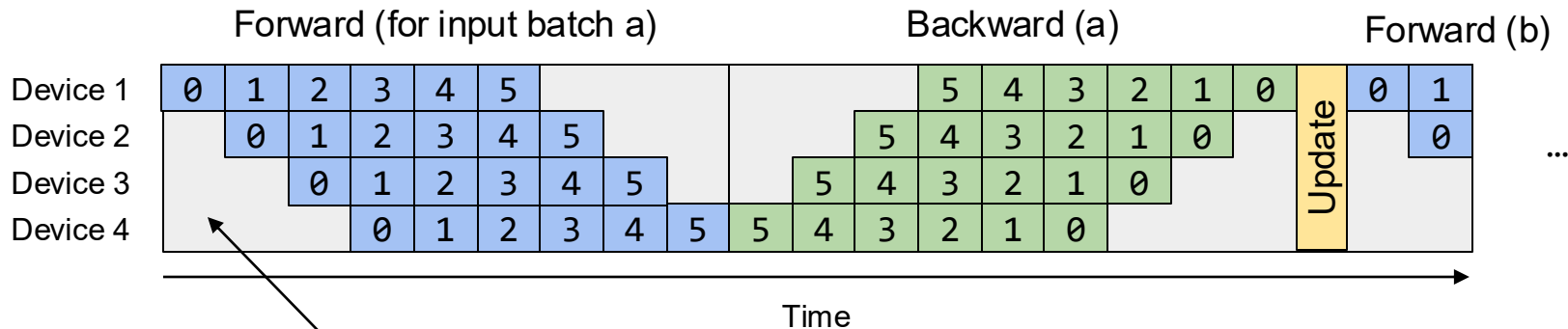
Idea: Modify pipeline schedule to improve efficiency, but keep the computation and convergence semantics exactly the same as if training with a single device.

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.

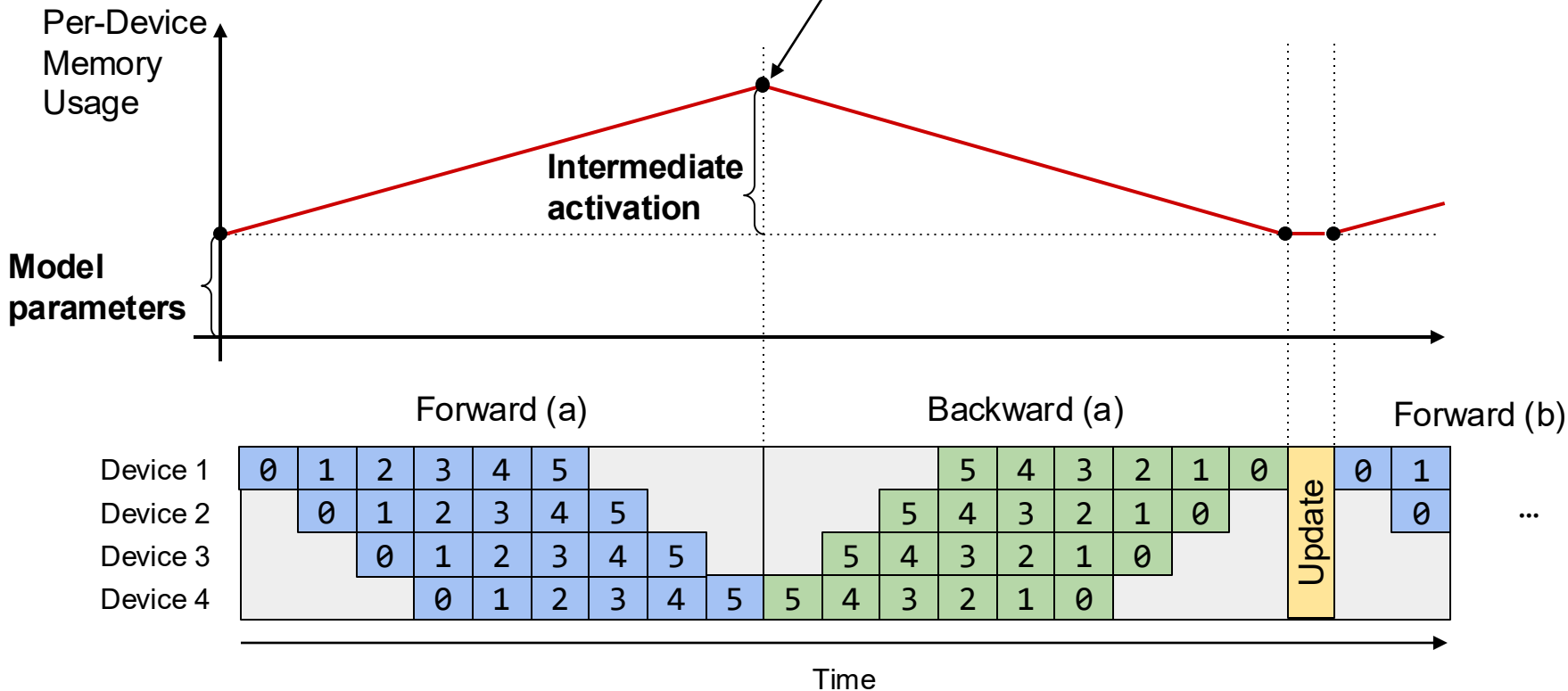
GPipe: Experimental Results

Table: Normalized training throughput using GPipe with different number of devices (stages) and different number of micro-batches M on TPUs.

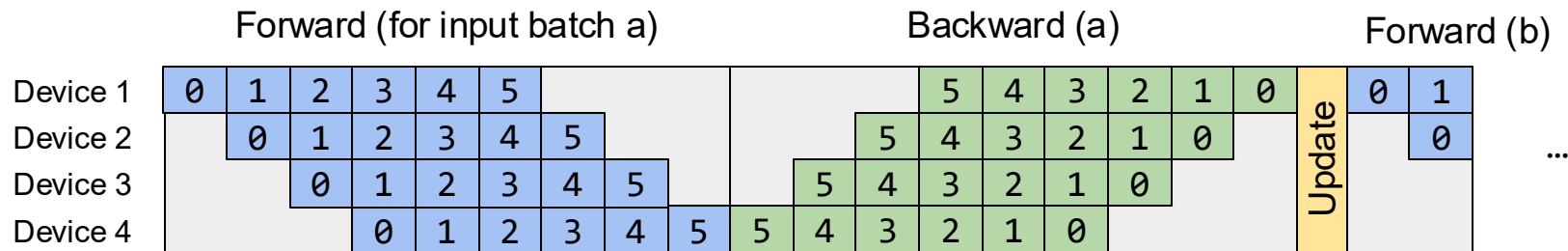
	#TPUs = 2	#TPUs = 4	#TPUs = 8
#Micro-batches = 1	1	1.07	1.3
#Micro-batches = 4	1.7	3.2	4.8
#Micro-batches = 32	1.8	3.4	6.3

GPipe: Memory Usage

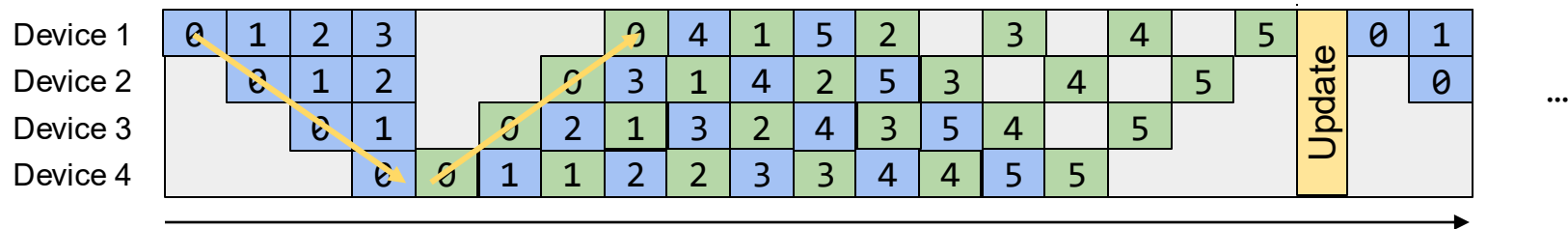
$$= \text{Parameters} + \text{Activation} \times \text{\#Micro-Batches}$$



GPipe Schedule:



1F1B (1 Forward 1 Backward) Schedule:



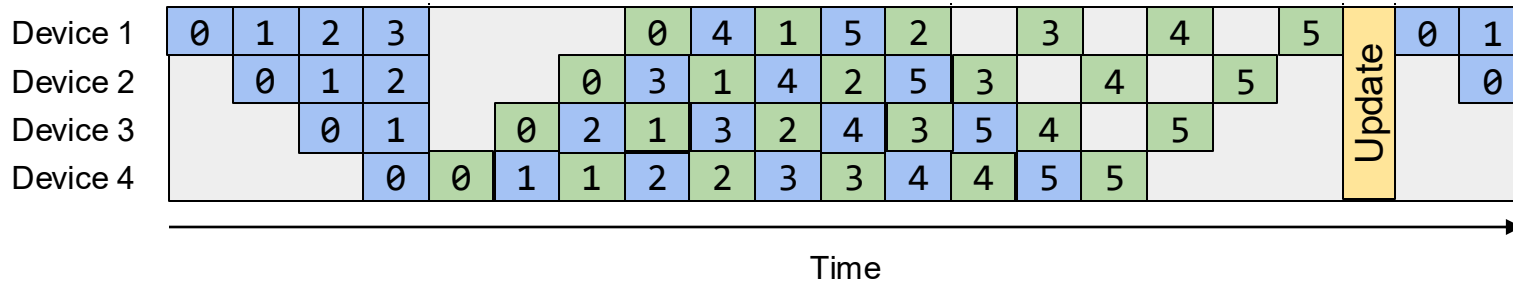
Same Latency

Perform backward as early as possible

1F1B Memory Usage

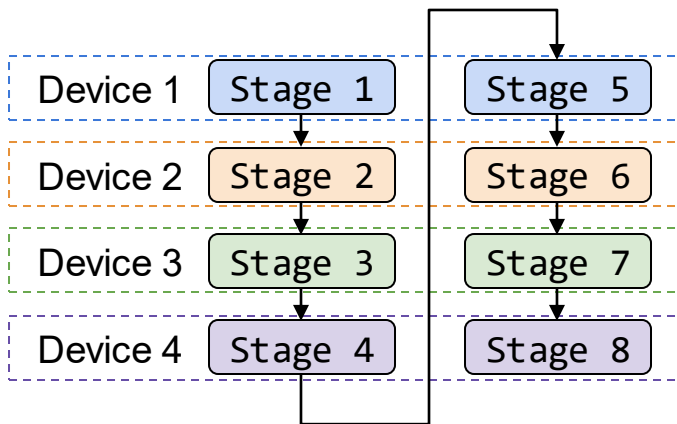
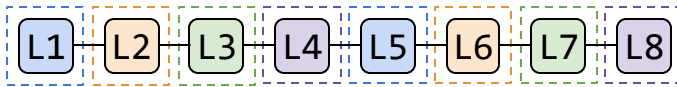
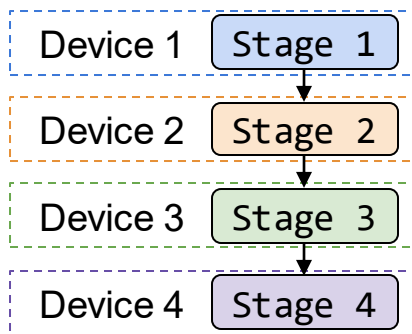
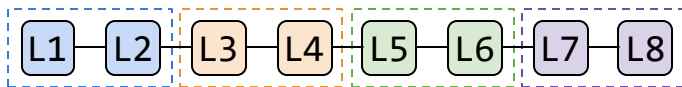
Maximum
per-device
memory
usage

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



Interleaved 1F1B

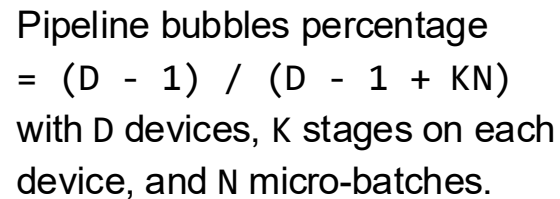
Idea: Slice the neural network into more fine-grained stages and assign multiple stages to reduce pipeline bubble.



Pro:

Con:

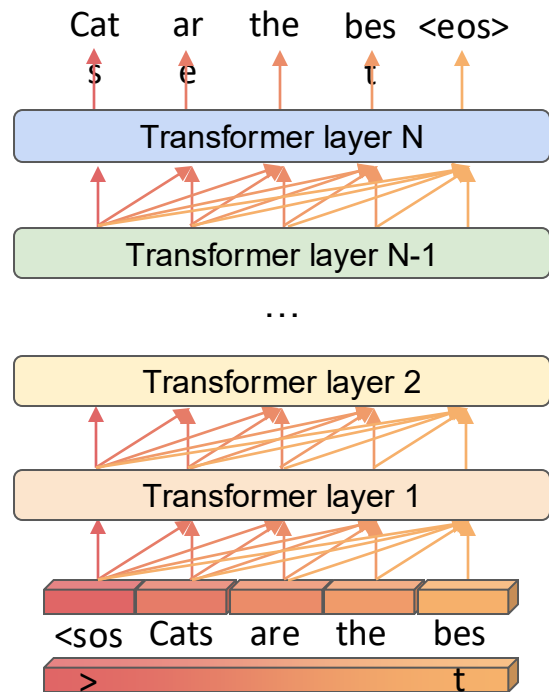
More communication overhead
between stages.



TeraPipe

Idea: The computation of an input token only depends on previous tokens but not future tokens for autoregressive models.

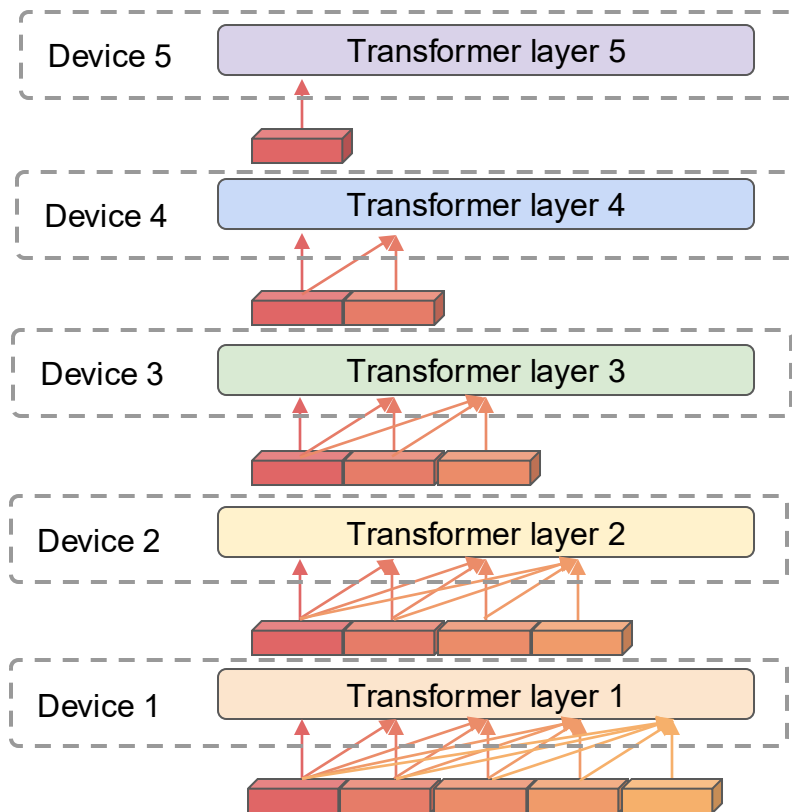
Further reduce the bubble size by pipelining within a sequence.



TeraPipe

Idea: The computation of an input token only depends on previous tokens but not future tokens for autoregressive models.

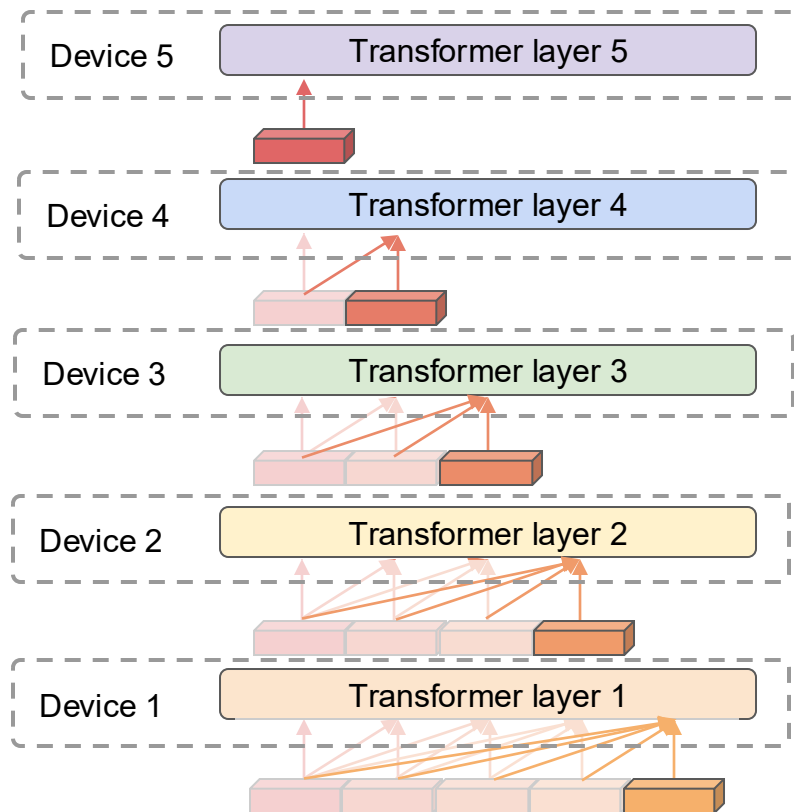
Further reduce the bubble size by pipelining within a sequence.



TeraPipe

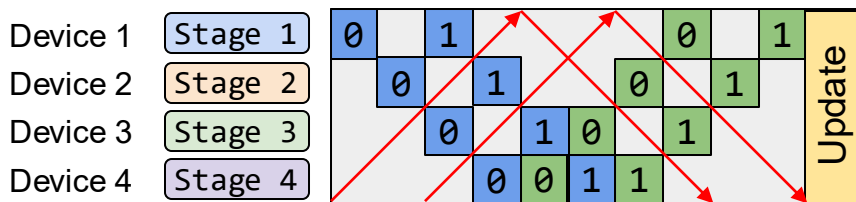
Idea: The computation of an input token only depends on previous tokens but not future tokens for autoregressive models.

Further reduce the bubble size by pipelining within a sequence.

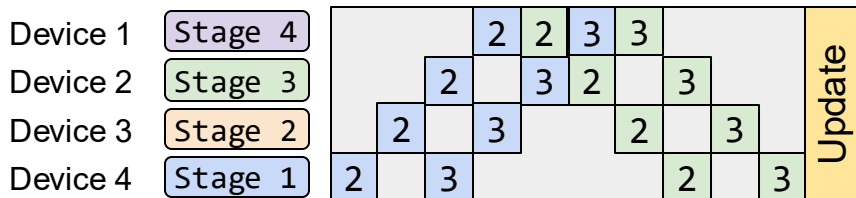


Chimera

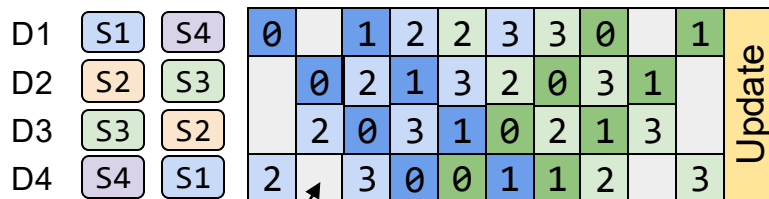
Idea: Store bi-directional stages and combine bidirectional pipeline to further reduce pipeline bubbles.



+



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.



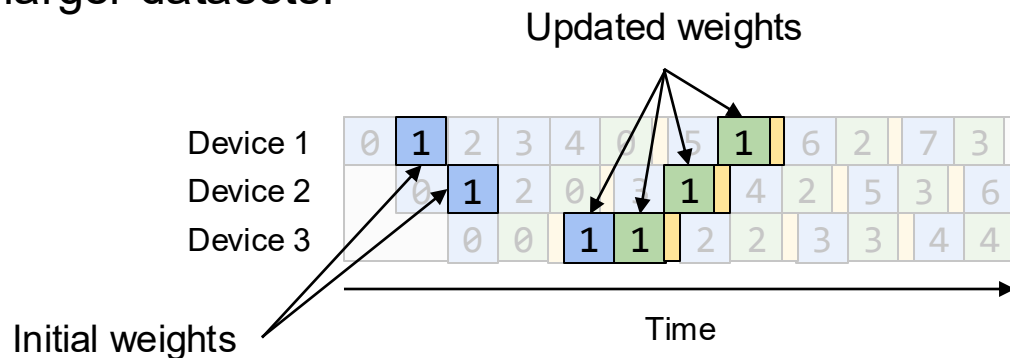
Cons:

- 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.



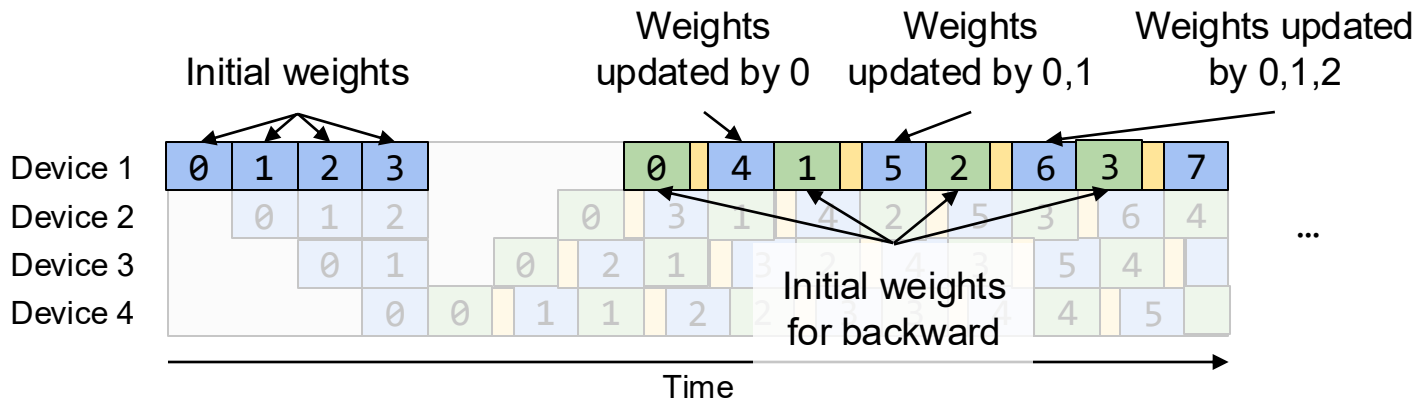
PipeMare: modify the optimizer to improve AMPNet convergence

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.

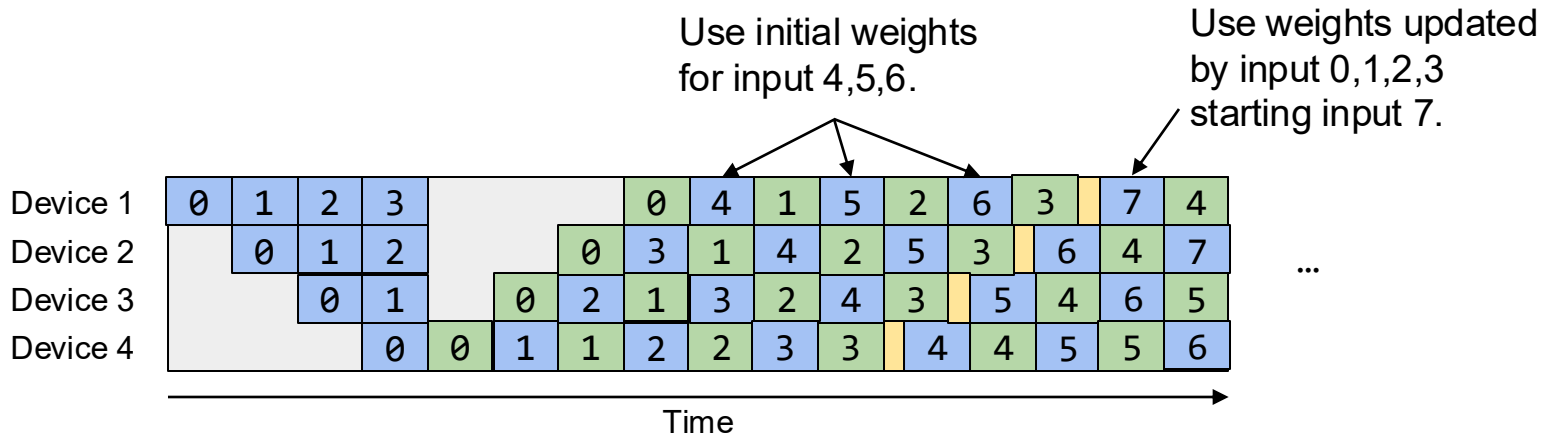
Con: No memory saving compared to single device case.



Pipedream-2BW

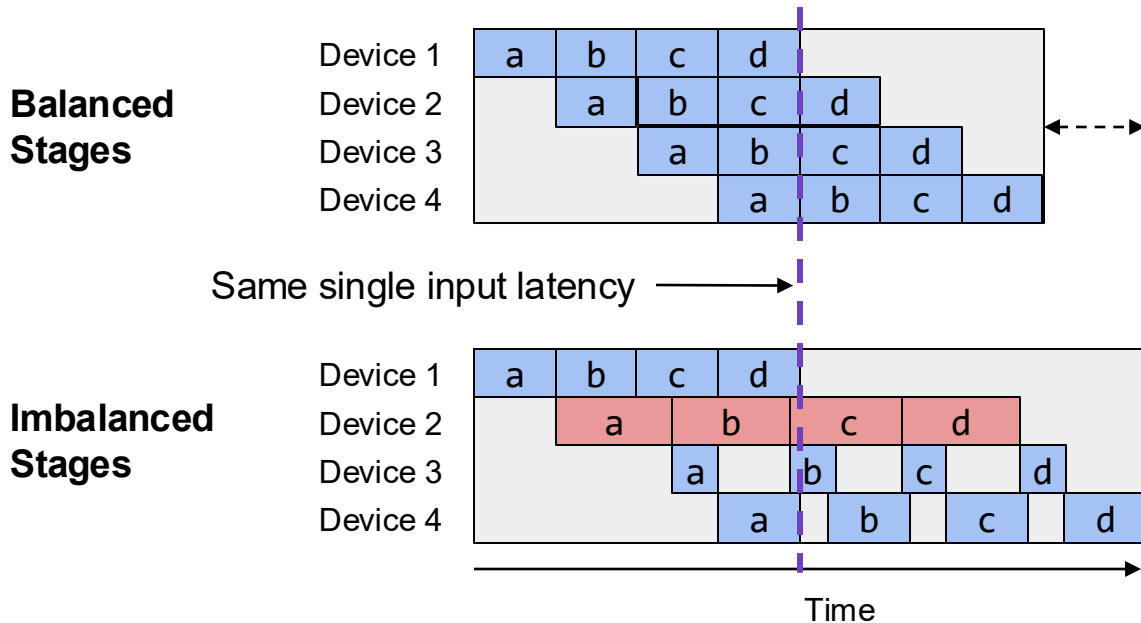
Idea: Reduce Pipedream's memory usage (only store 2 copies) by updating weights less frequently. Weights always stalled by 1 update.

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. ...

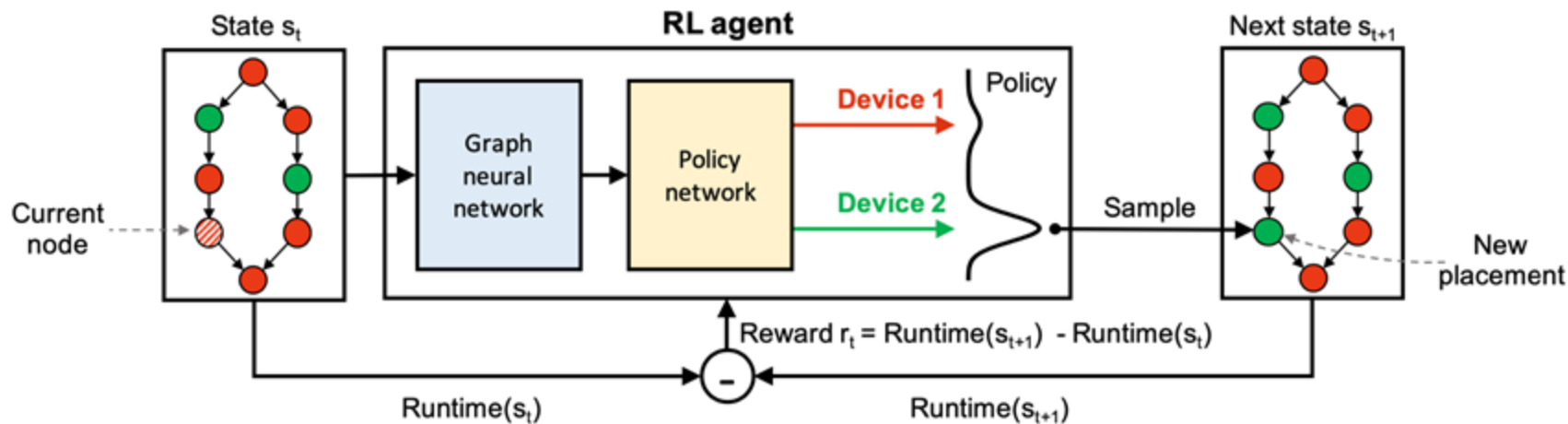
RL-Based Partitioning Algorithm

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

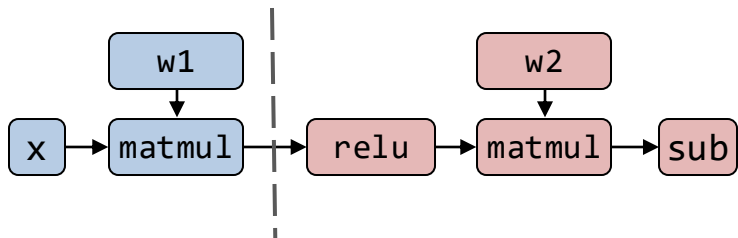
RL-Based / Optimization-Based Automatic Stage Partitioning

Where We Are

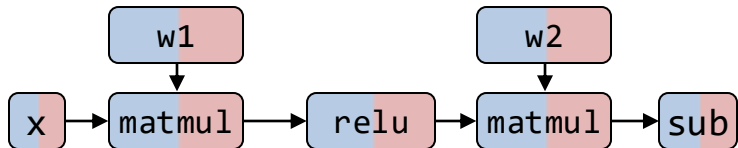
- Model parallelism
 - Inter-op parallelism
 - **Intra-op parallelism**

Recap: Intra-op and Inter-op

Strategy 1: Inter-operator Parallelism



Strategy 2: Intra-operator Parallelism



This section:

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

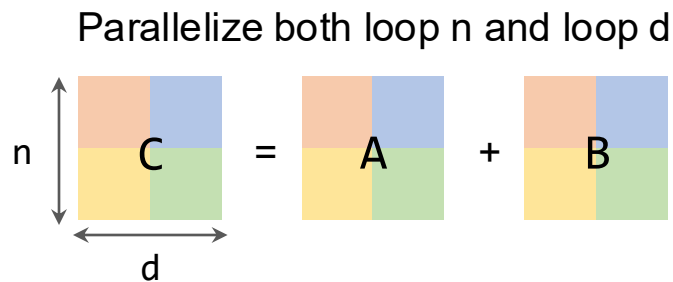
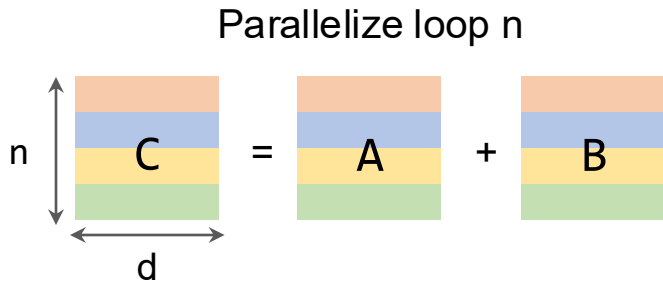
Parallelize One Operator

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

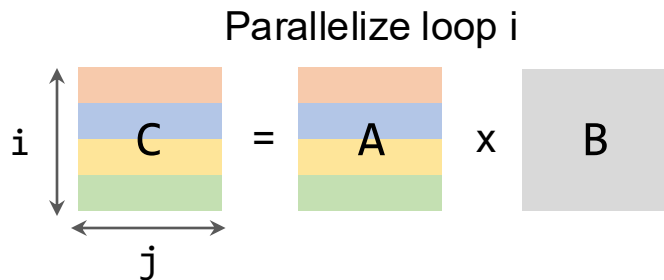
Matrix multiplication

```
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

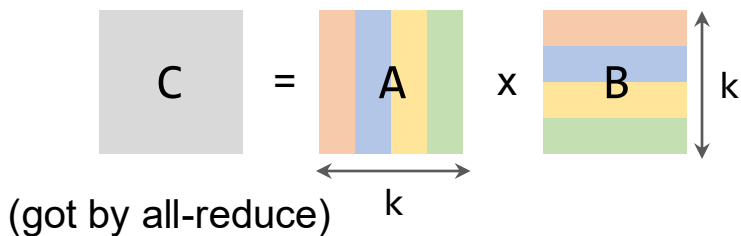
```
for i in range(0, N):  
    for j in range(0, M):  
        for k in range(0, K):  
            C[i,j] = C[i,j] + A[i,k] x B[k,j]
```

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



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

Parallelize One Operator

Matrix multiplication

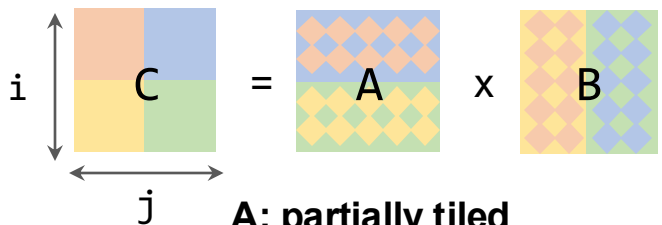
```
for i in range(0, N):  
    for j in range(0, M):  
        for k in range(0, K):  
            C[i,j] = C[i,j] + A[i,k] x B[k,j]
```

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

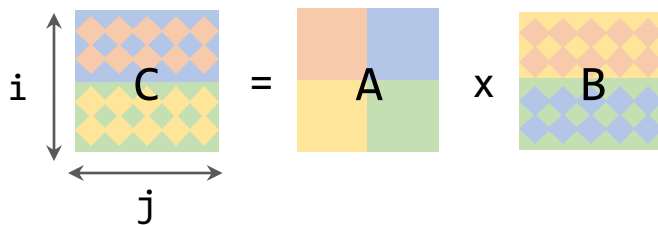
device 1 device 2 device 3 device 4

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 (≤ 5) for parallelization.

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

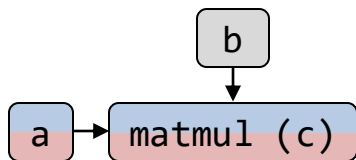
Complicated case: Parallelize loop h and w

Data Parallelism as A Case of Intra-op Parallelism

 Replicated  Row-partitioned  Column-partitioned

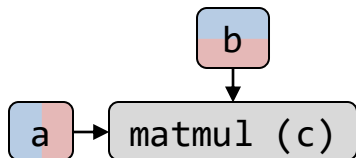
Matmul Parallelization Type 1

communication cost = 0



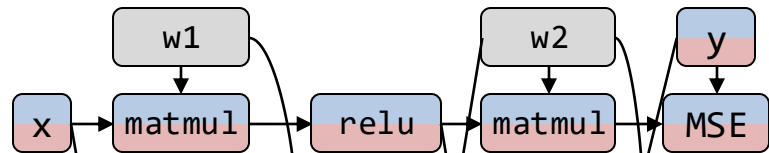
Matmul Parallelization Type 2

communication cost = all-reduce(c)



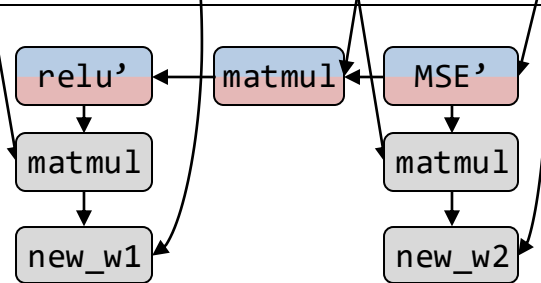
Forward Pass

Two “Type 1” matmuls: no communication



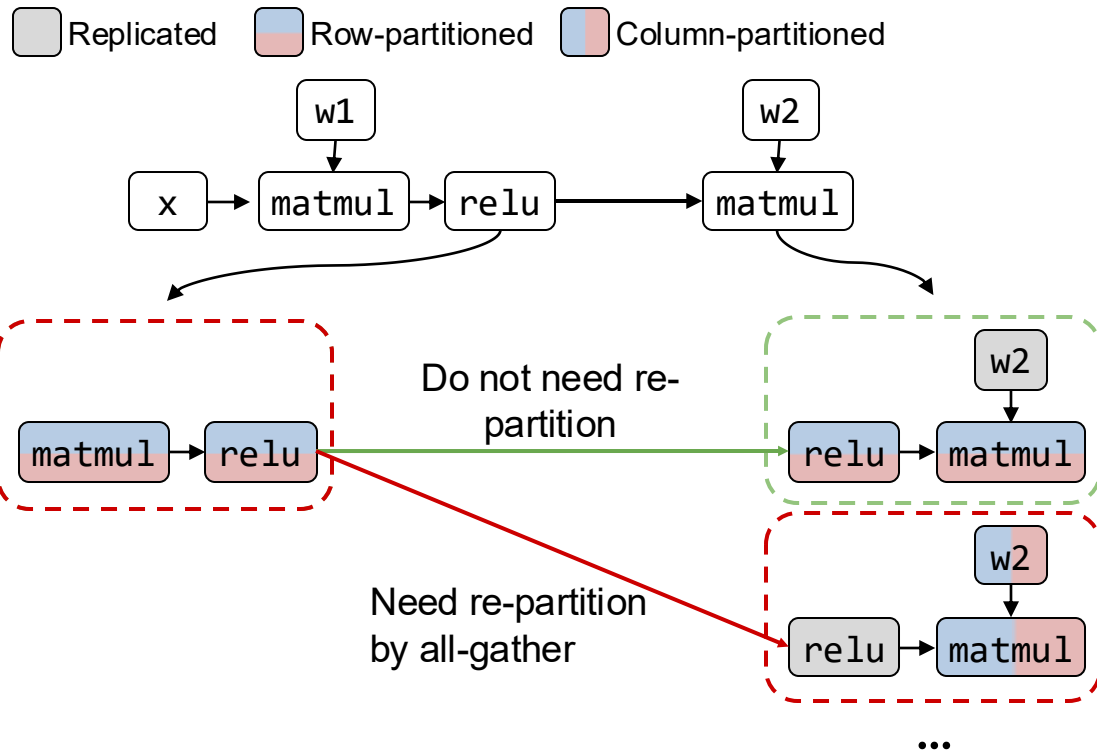
Backward Pass

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



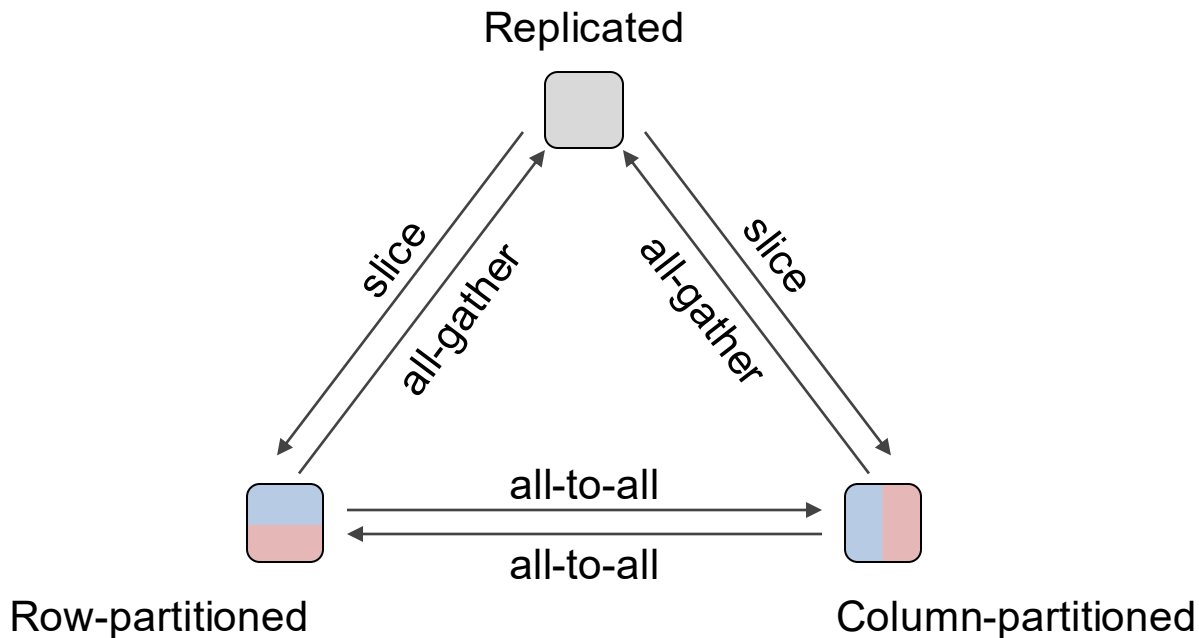
Re-partition Communication Cost

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



Re-partition Communication Cost

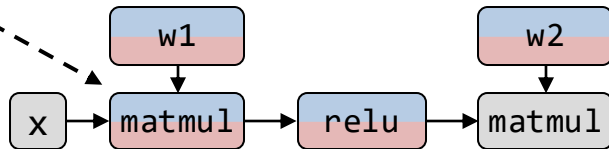
Different operators' parallelization strategies require different partition format of the same tensor



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

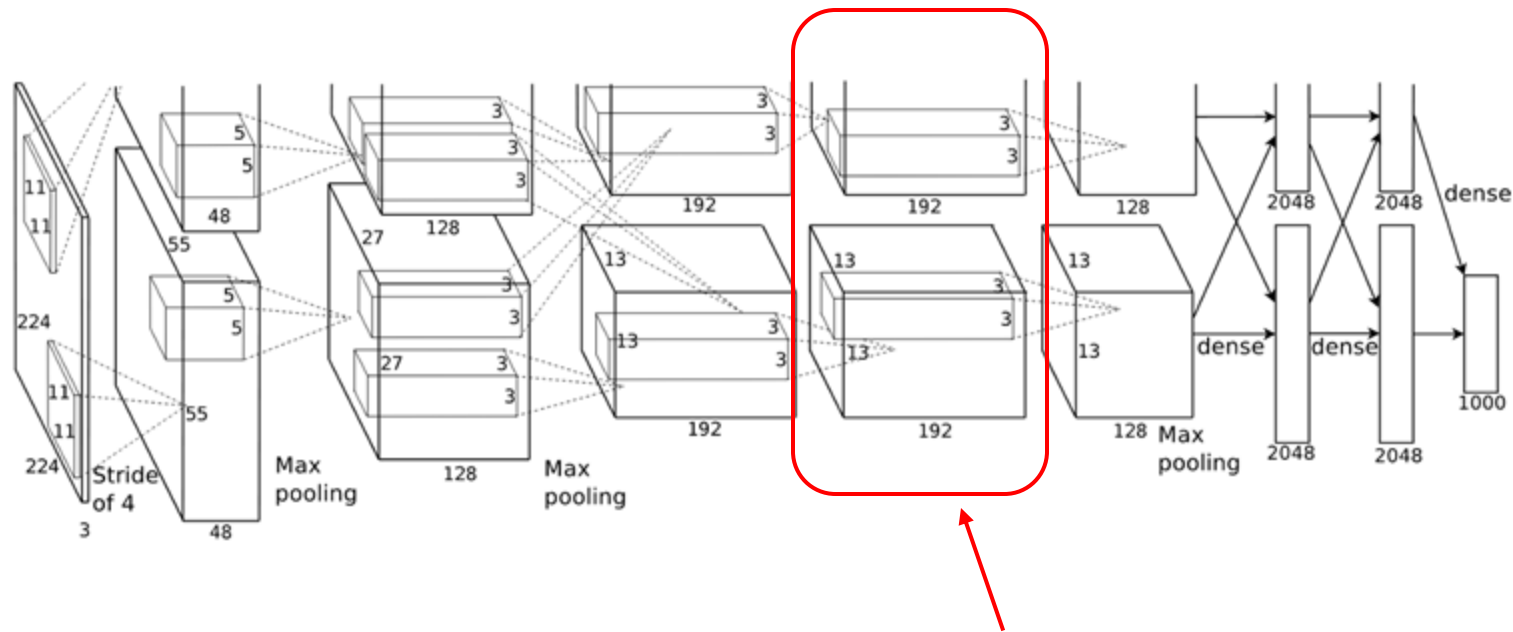
- AlexNet
- Megatron-LM
- GShard MoE

Systems for Intra-op Parallelism

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

AlexNet

Result: increase top-1 accuracy by 1.7%

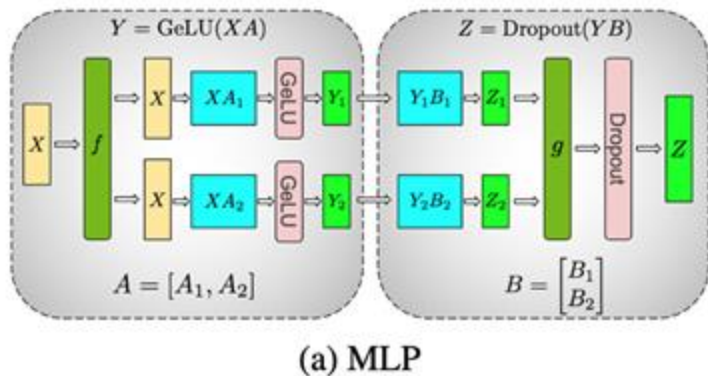


Assign a group convolution layer to 2 GPUs

Megatron-LM

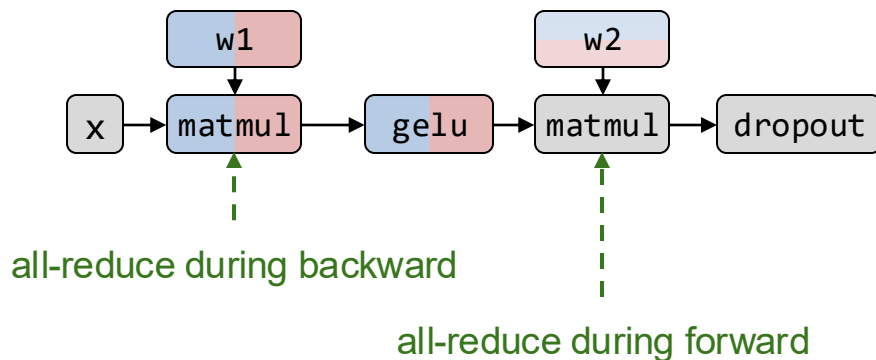
Result: a large language model with 8.3B parameters that outperforms SOTA results

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



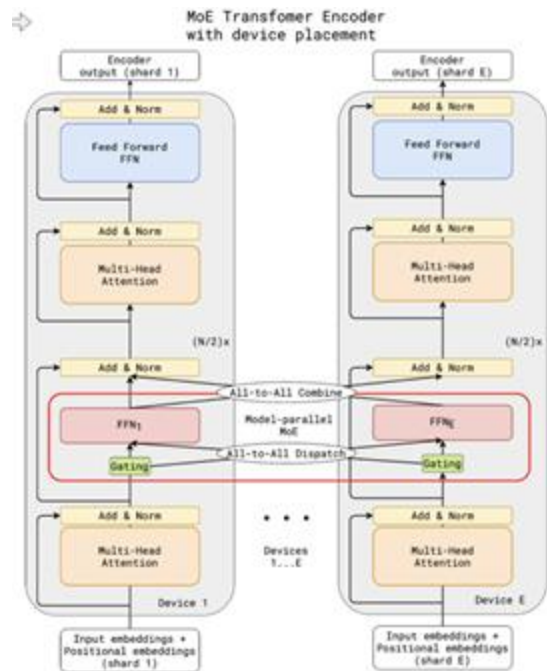
Illustrated with the notations in this tutorial

Replicated Row-partitioned Column-partitioned



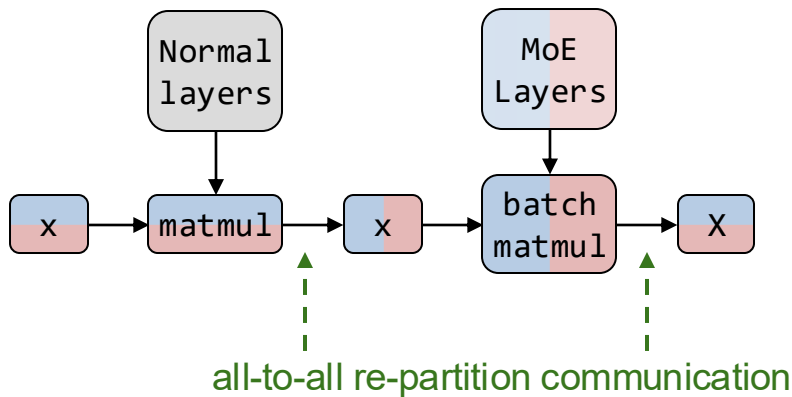
GShard MoE

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.

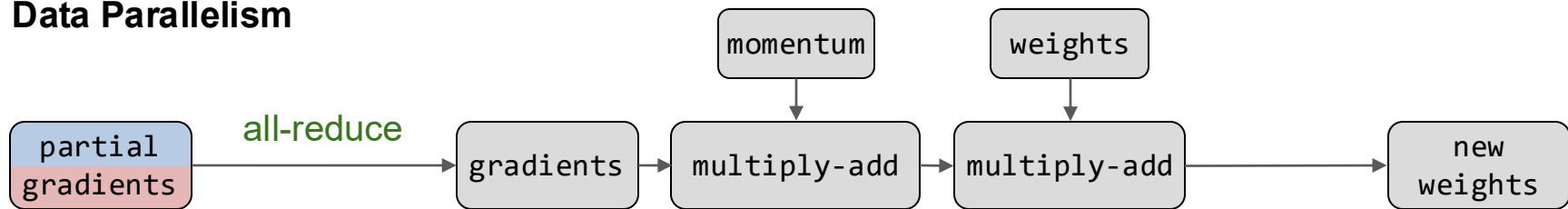
	Optimizer States (12M)	Gradients (2M)	Model Weights (2M)	Memory Cost	Communication Cost
Data Parallelism	Replicated	Replicated	Replicated	$16M$	all-reduce(2M)
ZeRO Stage 1	Partitioned	Replicated	Replicated	$4M + \frac{12M}{N}$	all-reduce(2M)
ZeRO Stage 2	Partitioned	Partitioned	Replicated	$2M + \frac{14M}{N}$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

ZeRO Stage 2

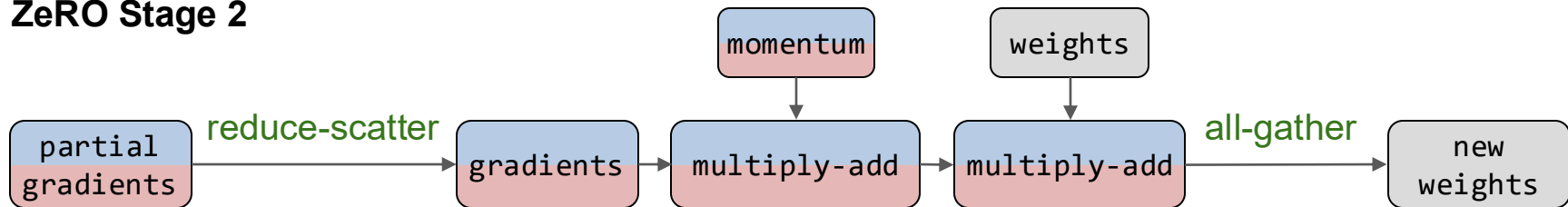
Key Idea: all-reduce = reduce-scatter + all-gather

 Replicated  Partitioned

Data Parallelism



ZeRO Stage 2



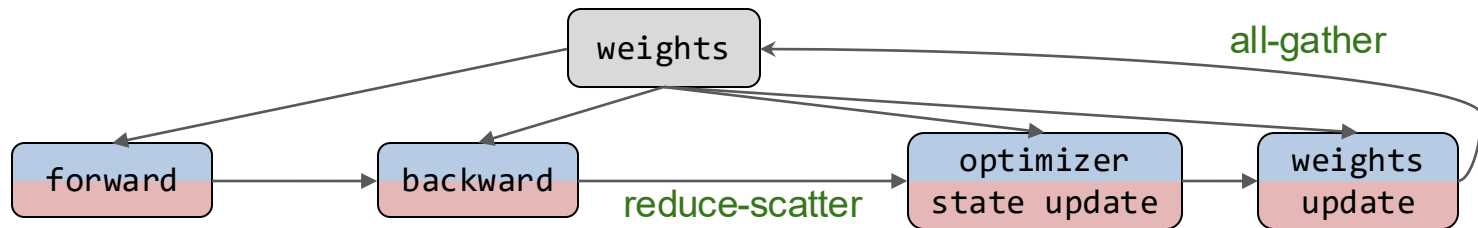
Same communication cost but save memory by partitioning more tensors

ZeRO Stage 3

Replicated Partitioned

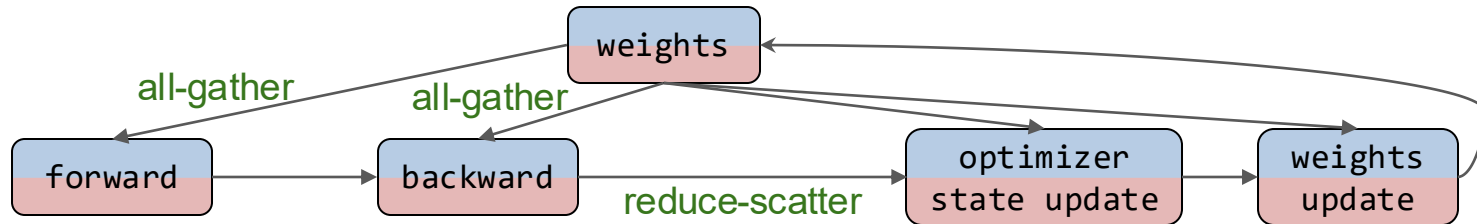
ZeRO Stage 2

communication cost
= all-reduce



ZeRO Stage 3

communication cost
= 1.5 all-reduce



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)

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

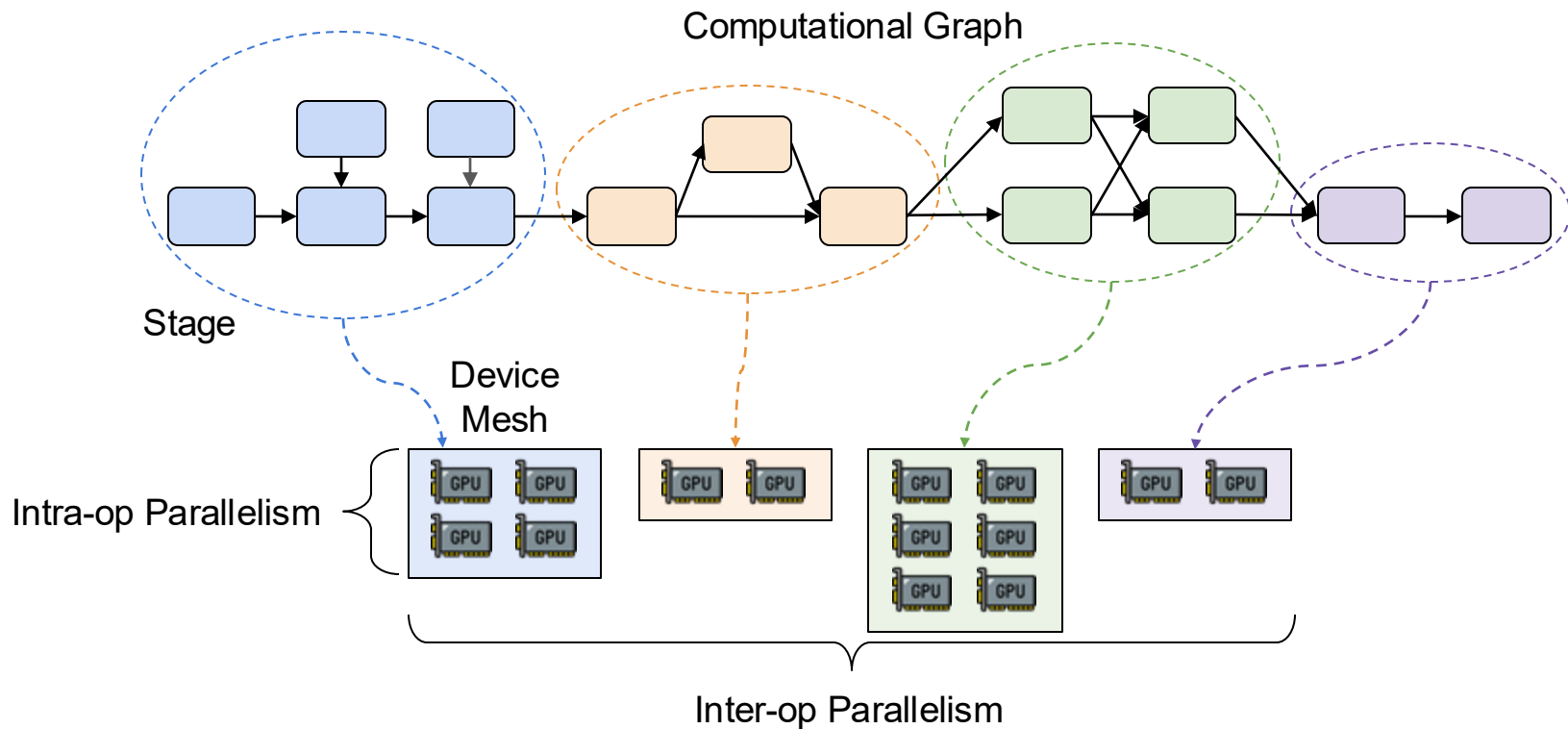
Mapping

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*.

Thanks for Enrolling DSC 204A!