

# Unity

Accelerating DNN Training Through Joint Optimization  
of Algebraic Transformations and Parallelization

Colin Unger <sup>\*,1</sup>

Sina Lin <sup>6</sup>

Vinay Ramakrishnaiah <sup>4</sup>

Jamaludin Mohd-Yusof <sup>4</sup>

Jongsoo Park <sup>3</sup>

Zhihao Jia <sup>\*,2,3</sup>

Mandeep Baines <sup>3</sup>

Nirmal Prajapati <sup>4</sup>

Xi Luo <sup>7</sup>

Misha Smelyanskiy <sup>3</sup>

Wei Wu <sup>4,5</sup>

Carlos Efrain Quintero Narvaez <sup>3</sup>

Pat McCormick <sup>4</sup>

Dheevatsa Mudigere <sup>3</sup>

Alex Aiken <sup>1</sup>



1



2



3



4



5



6



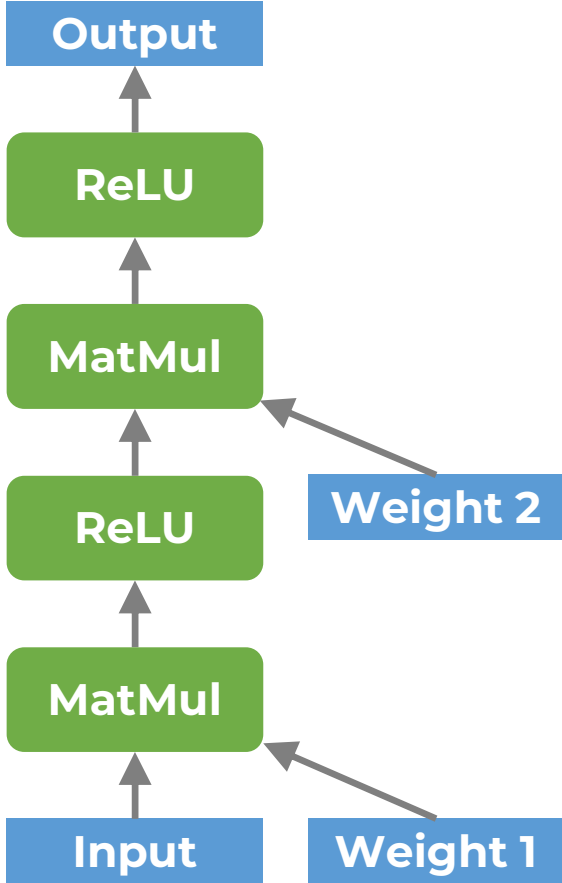
7

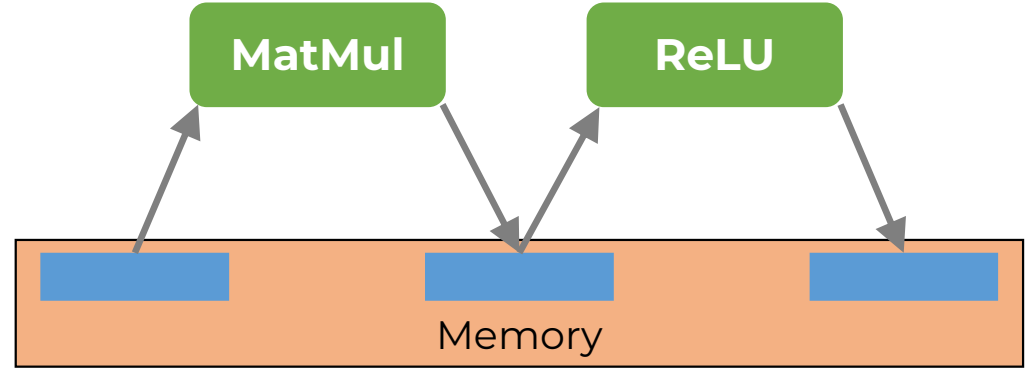
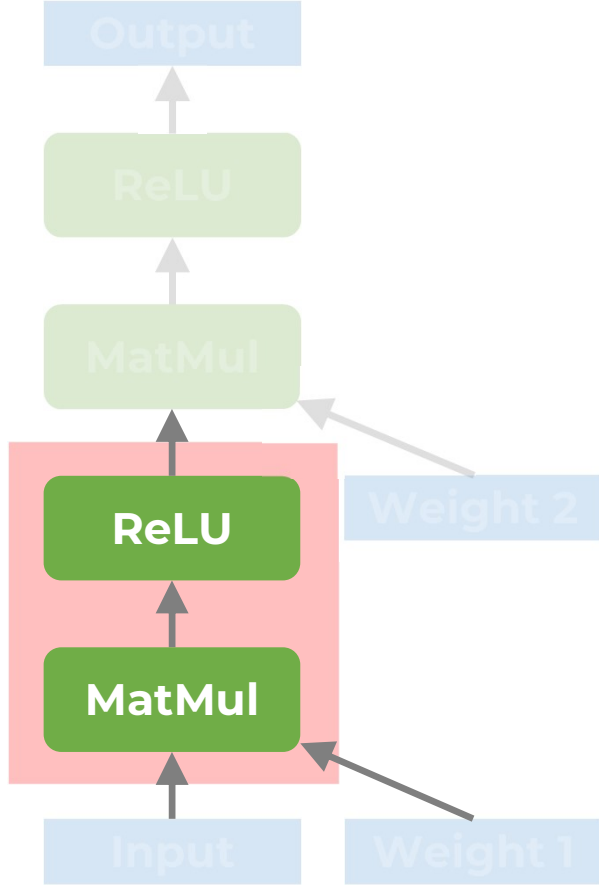
# Unity

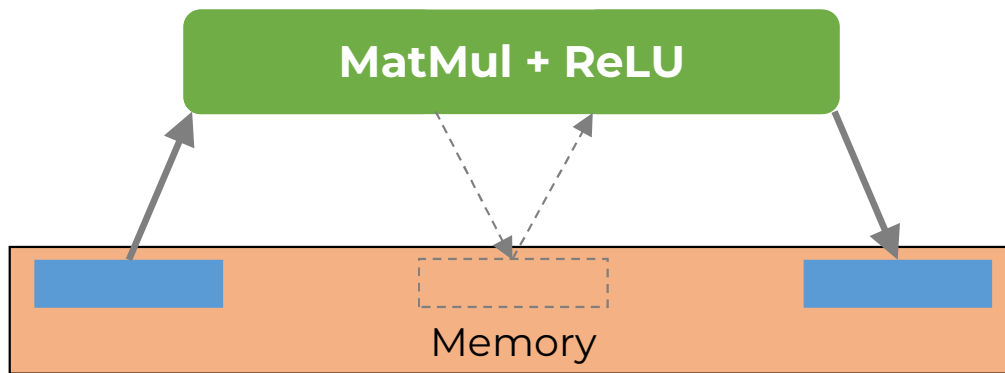
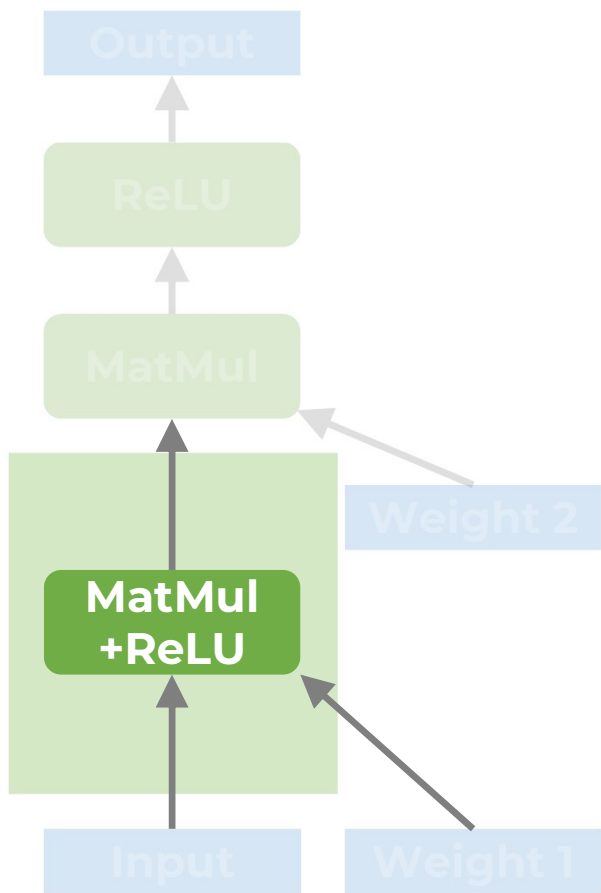
Accelerating DNN Training Through Joint Optimization  
of Algebraic Transformations and Parallelization

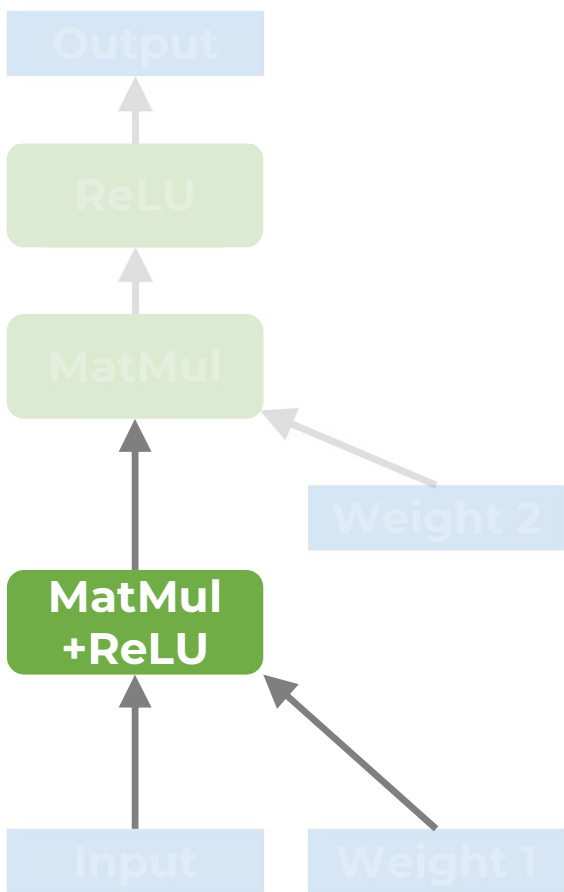
1. Algebraic Transformations
2. Parallelization

# 1. Algebraic Transformations

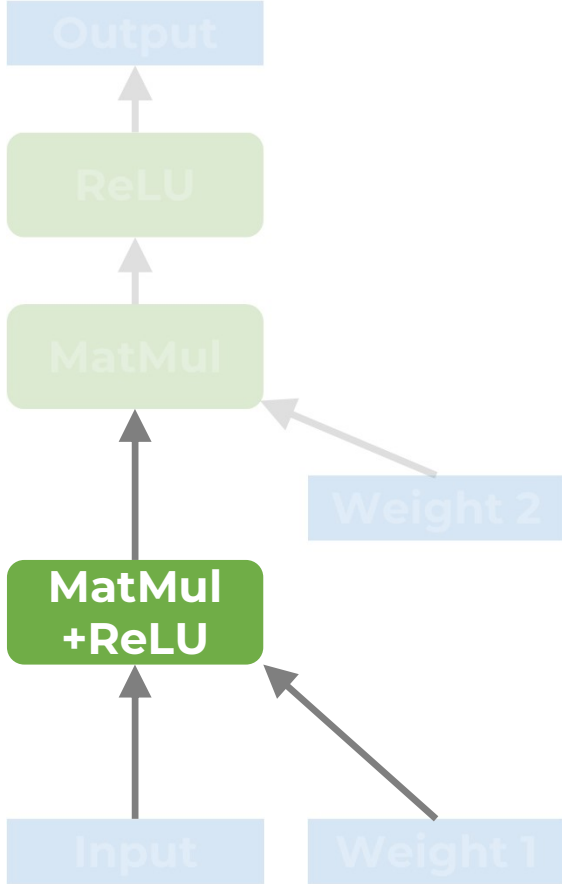






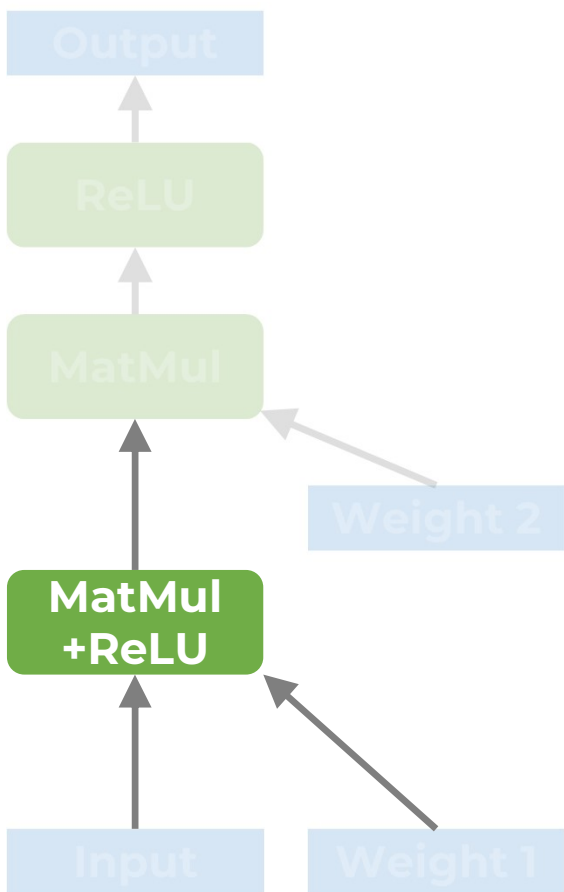


## Operator Fusion

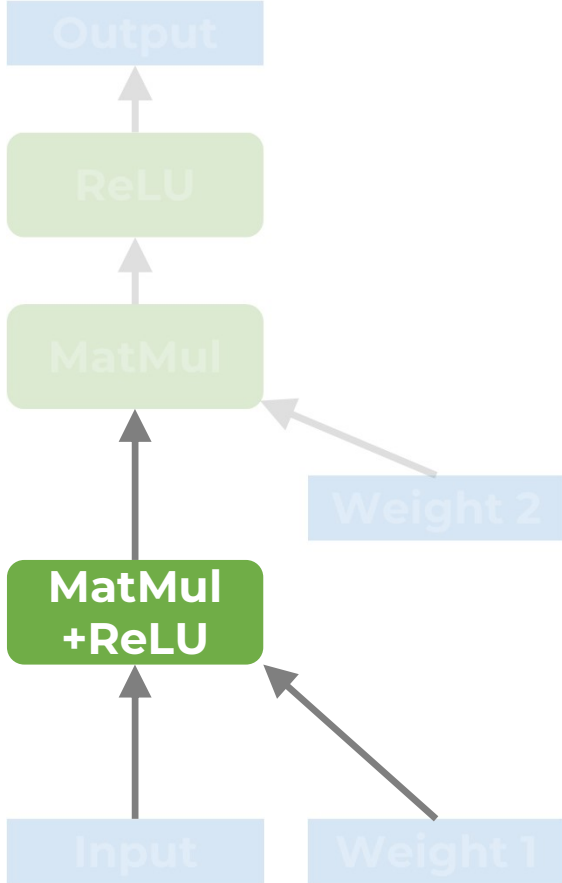


Operator Fusion  
Operator Splitting





Operator Fusion  
Operator Splitting  
Operator Reordering



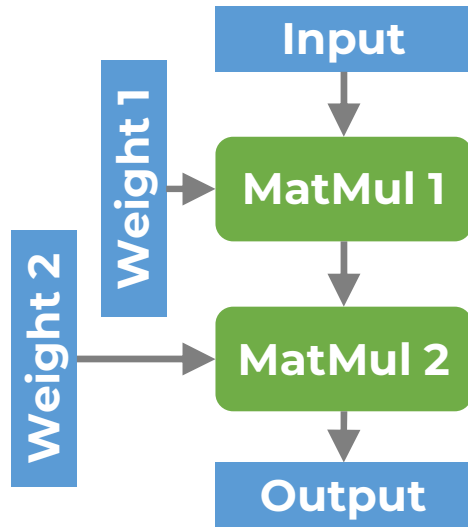
Operator Fusion  
Operator Splitting  
Operator Reordering

...

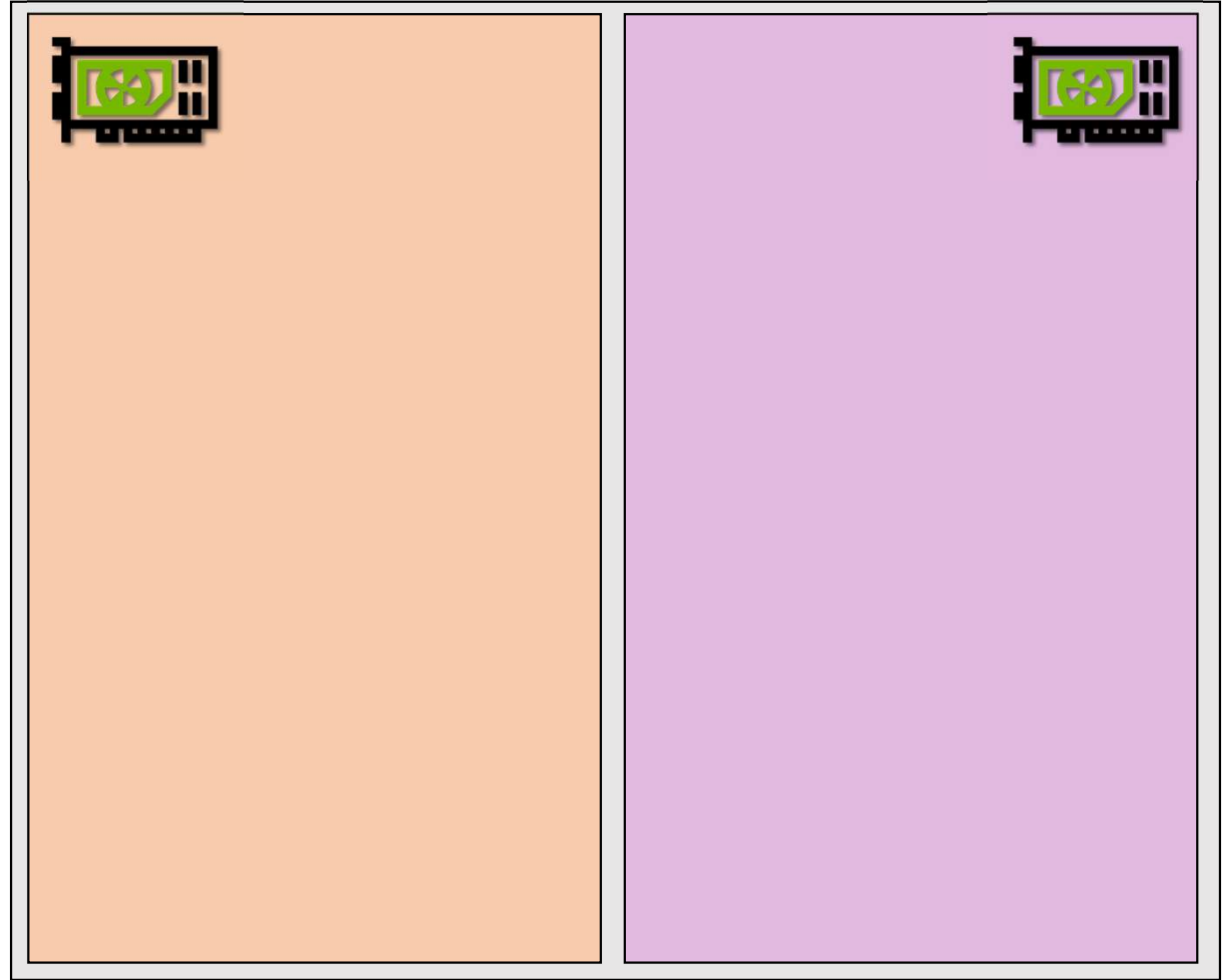
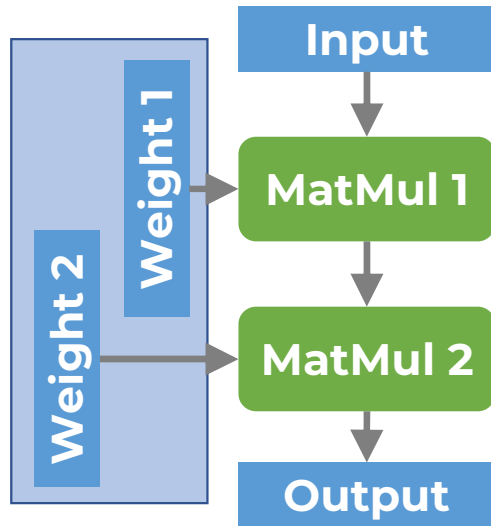
1. Algebraic Transformations
2. Parallelization

# Data Parallel

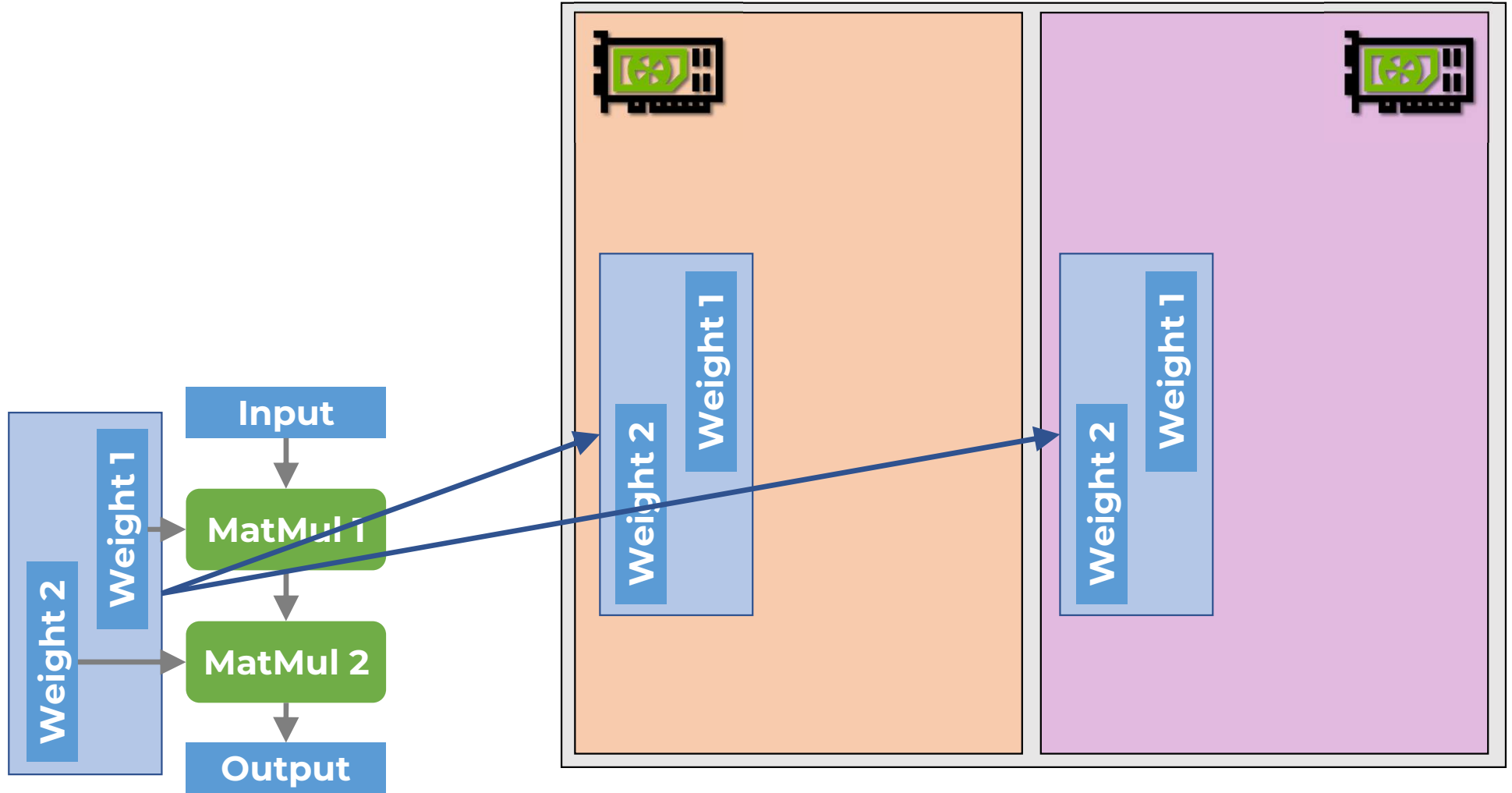
# Data Parallel



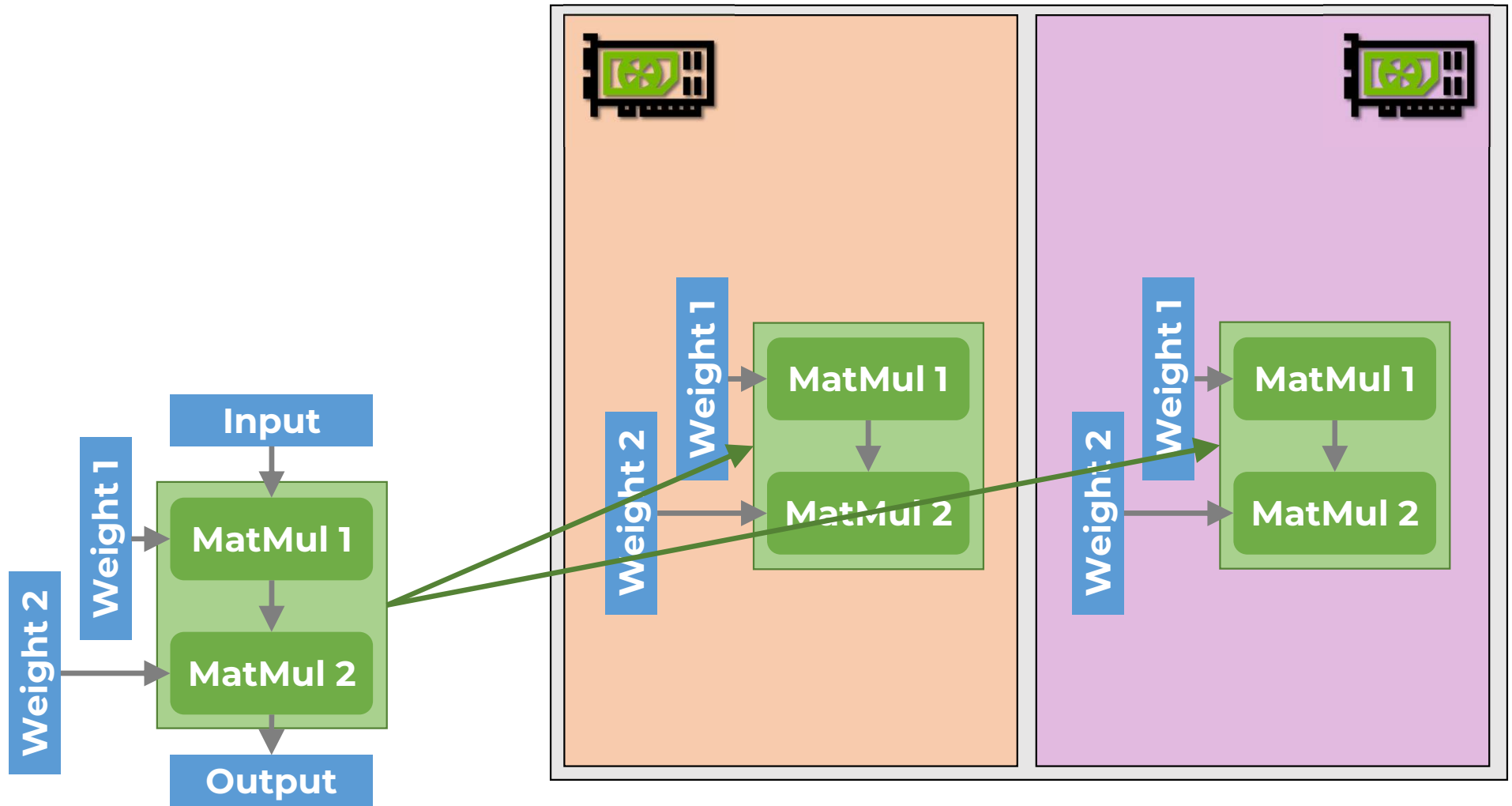
# Data Parallel



# Data Parallel

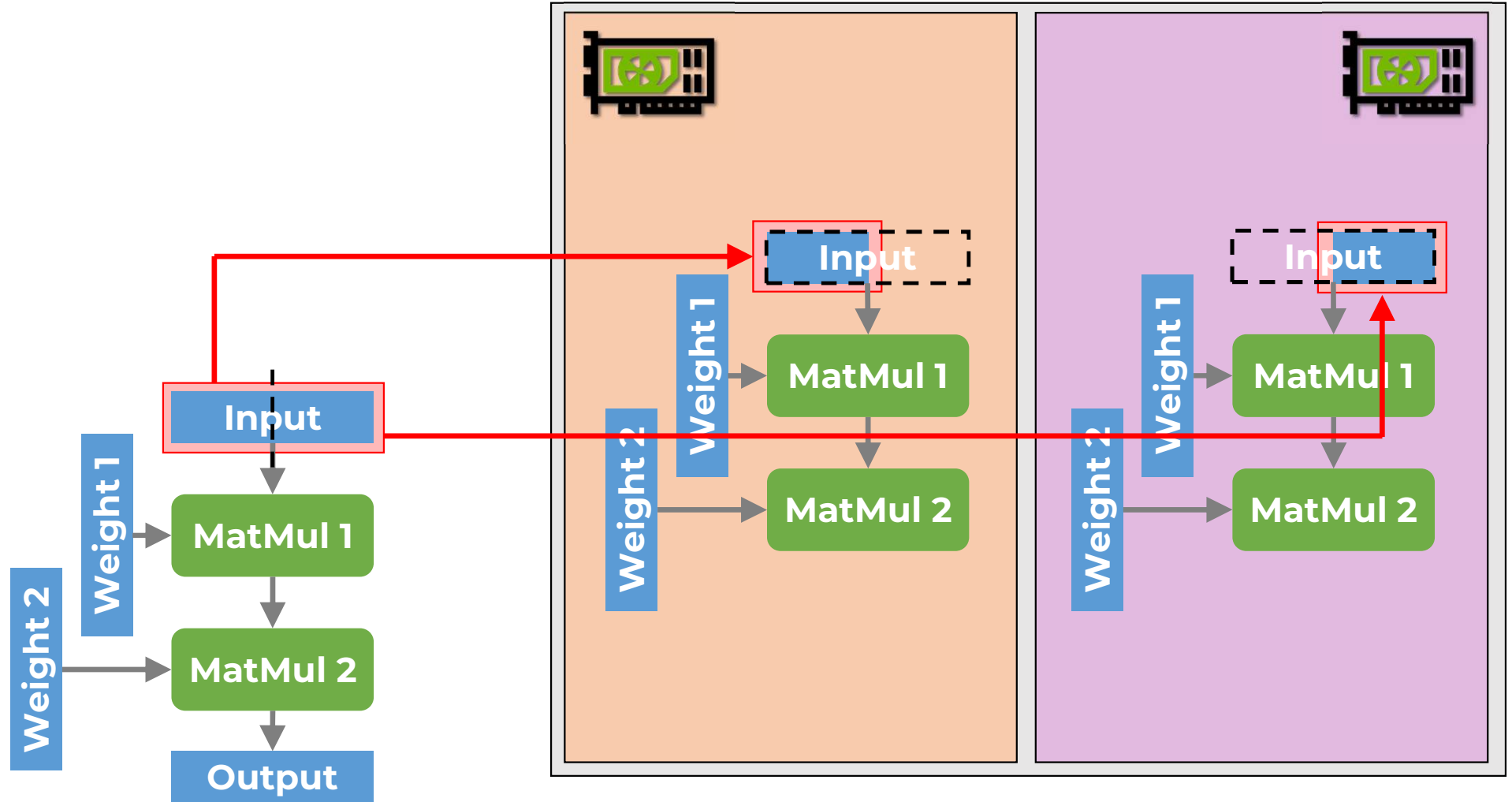


# Data Parallel

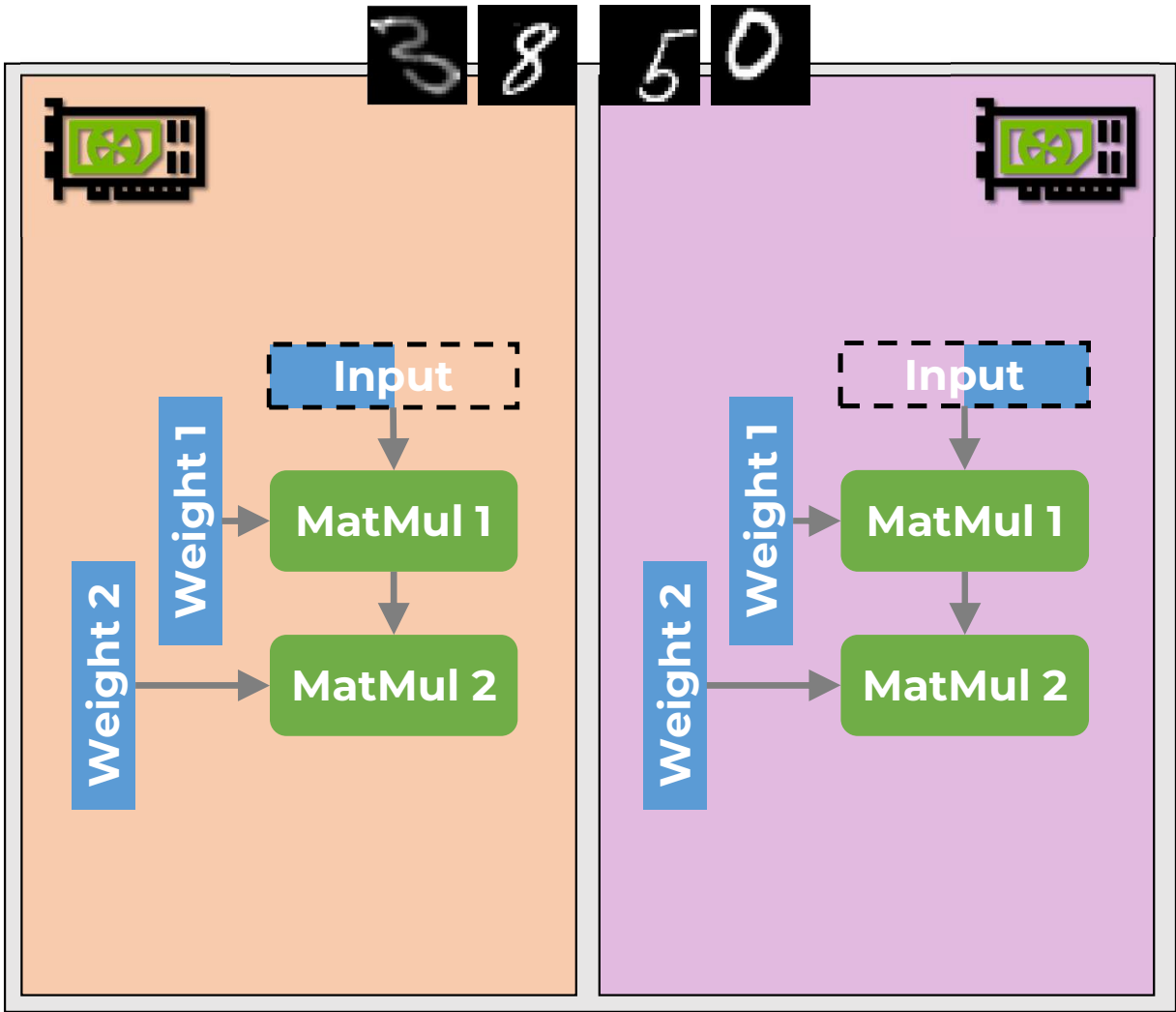
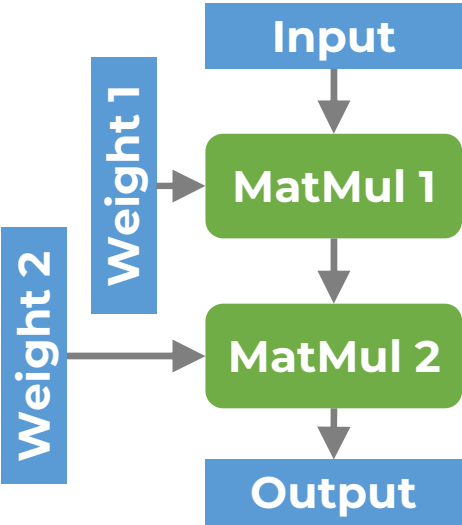




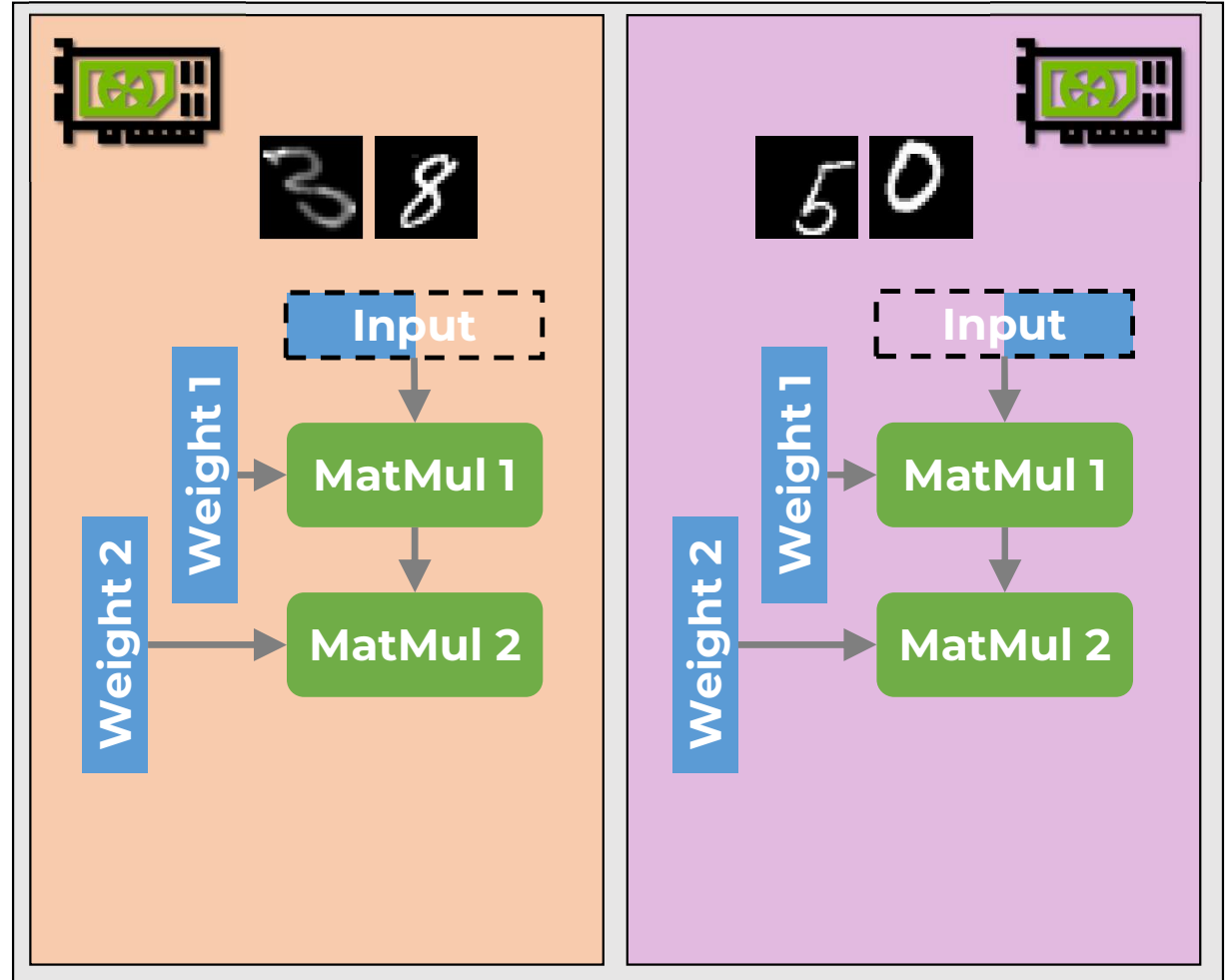
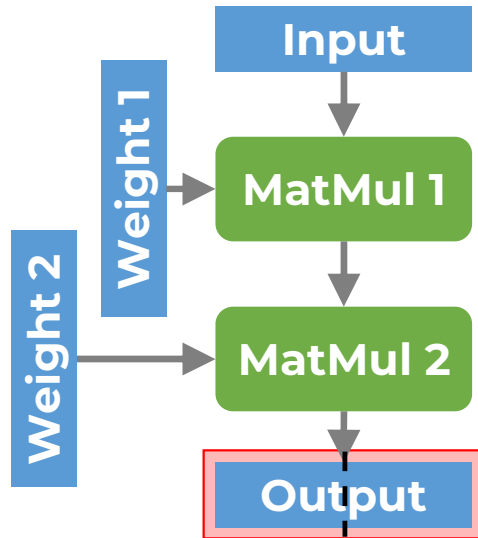
# Data Parallel



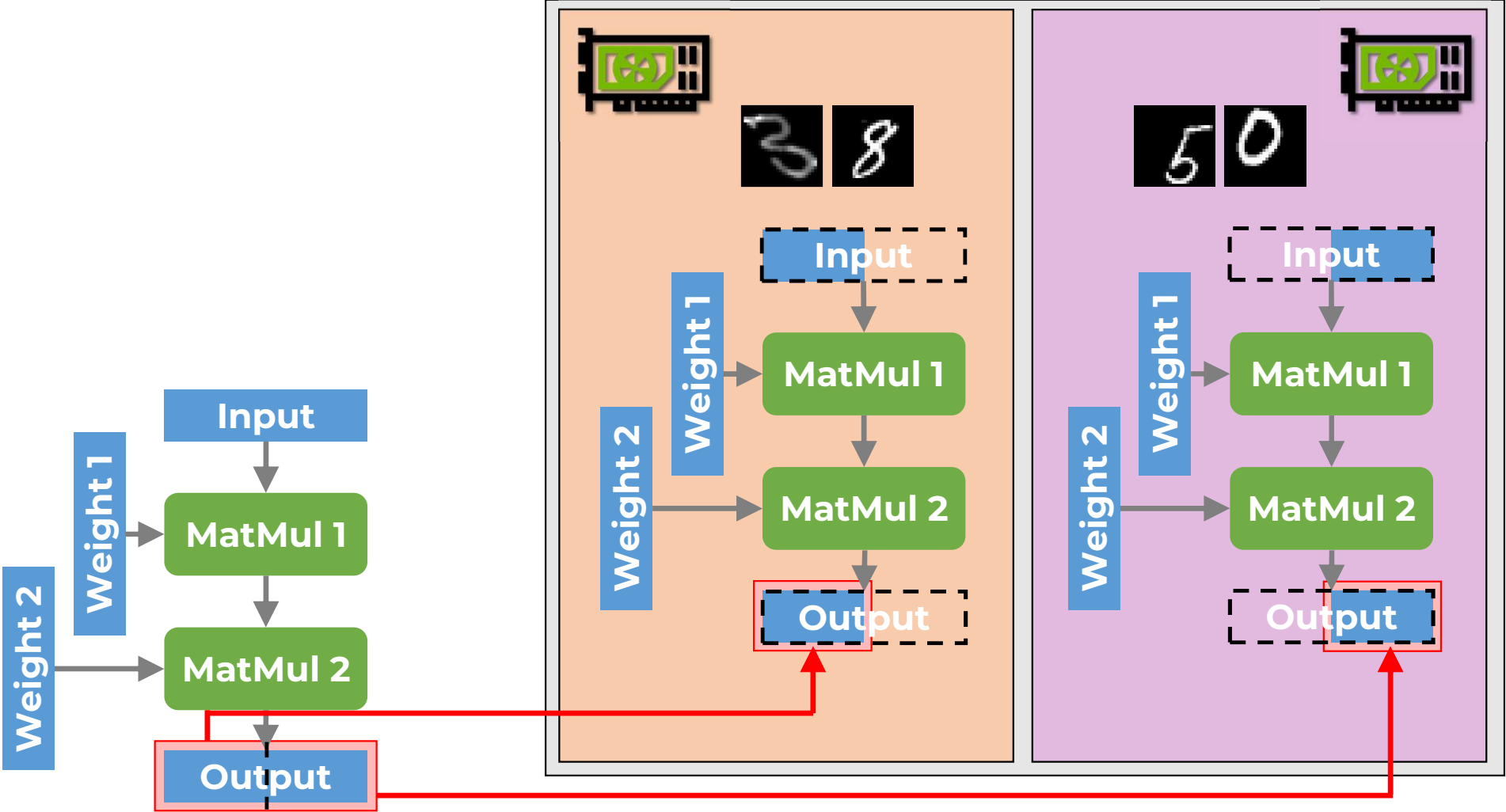
# Data Parallel



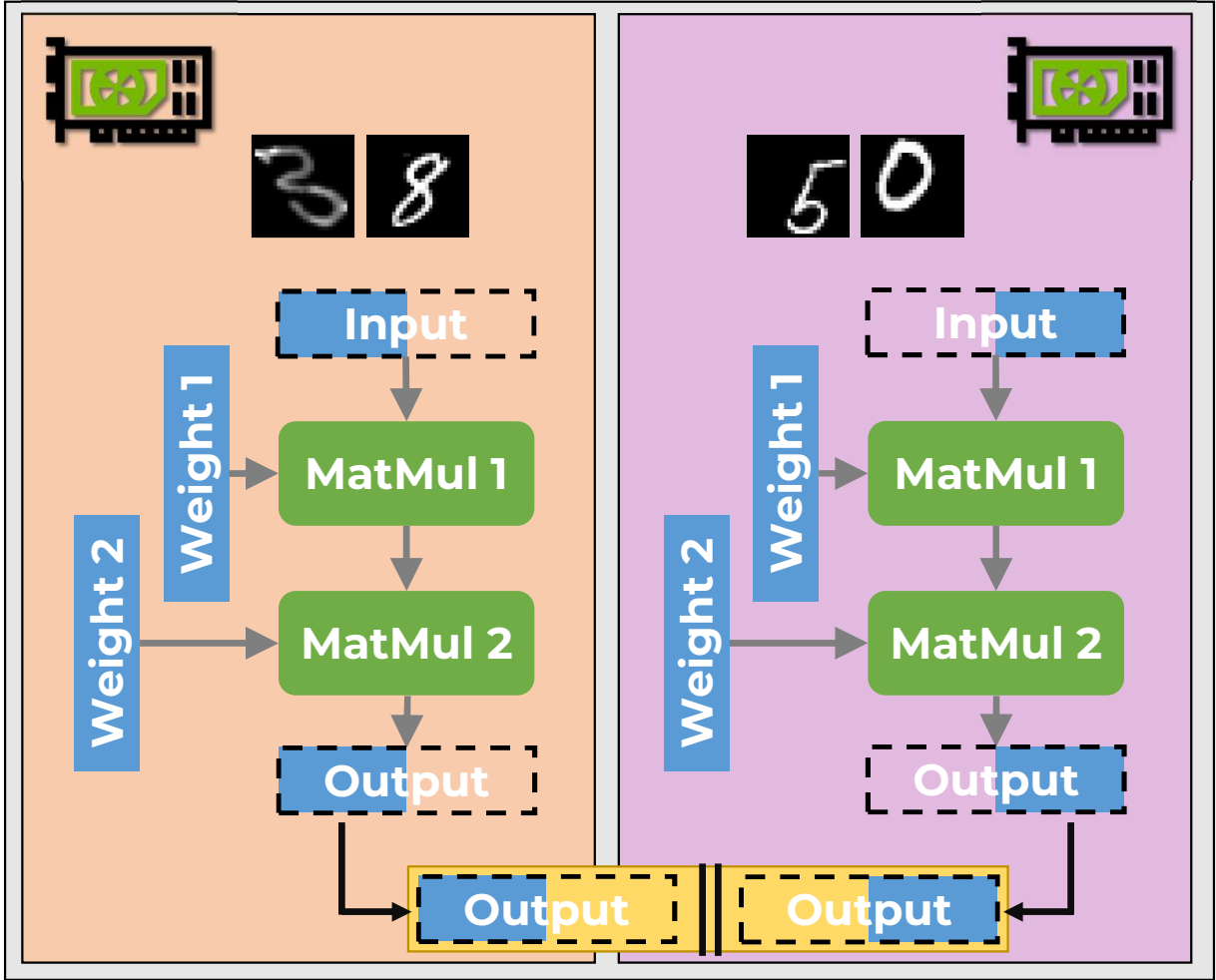
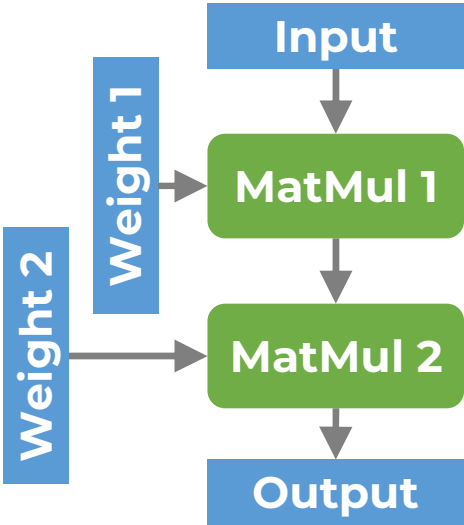
# Data Parallel



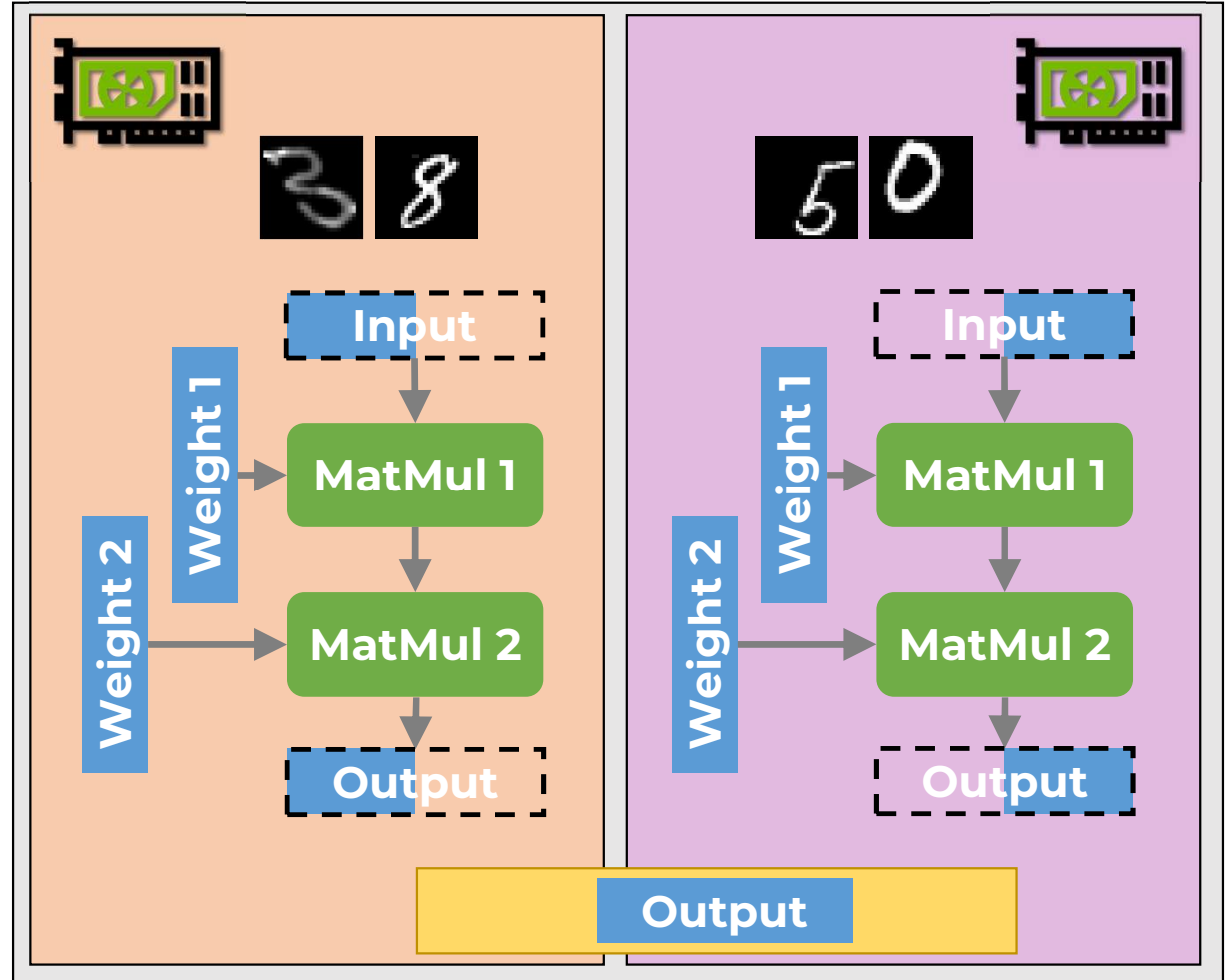
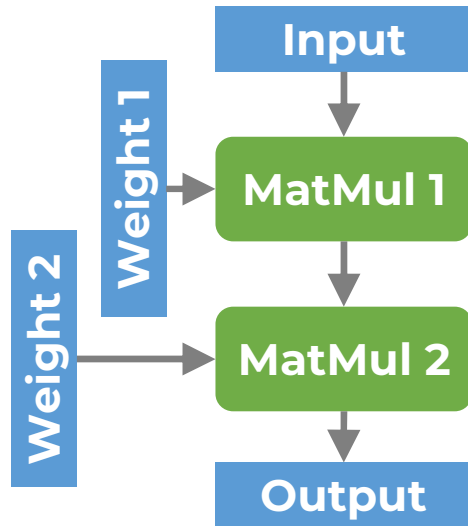
# Data Parallel



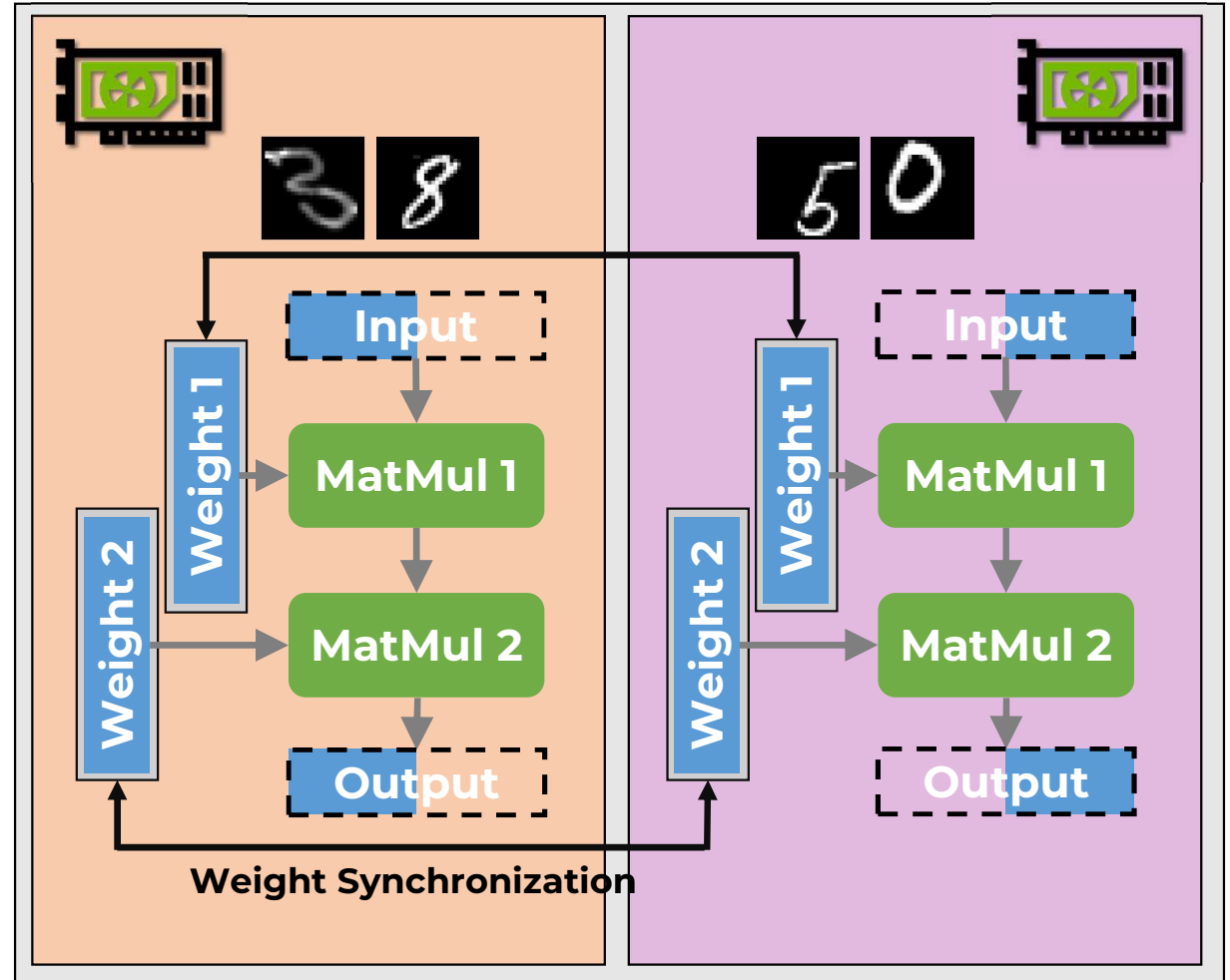
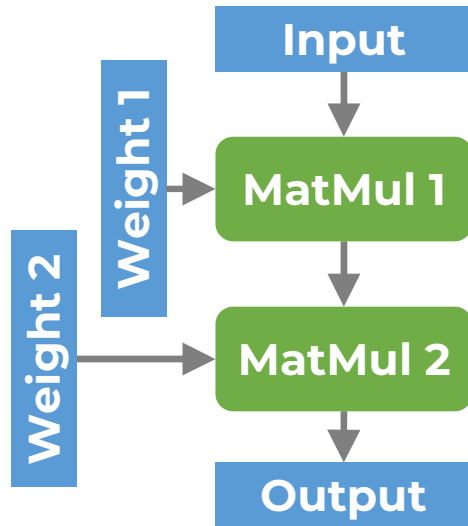
# Data Parallel



# Data Parallel



# Data Parallel



# Data Parallel



Data Parallel  
Model Parallel

Data Parallel

Model Parallel

Attribute Parallel

Operation	Parallelizable Dimensions		
	(S)ample	(A)tttribute	(P)arameter
1D pooling	sample	length, channel	
1D convolution	sample	length	channel
2D convolution	sample	height, width	channel
Matrix multiplication	sample		channel

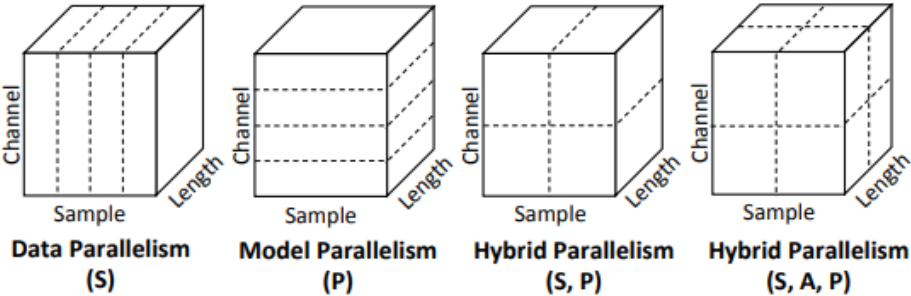


Figure 3: Example parallelization configurations for 1D convolution. Dashed lines show partitioning the tensor.

Data Parallel  
Model Parallel  
Attribute Parallel  
Reduction Parallel

Data Parallel  
Model Parallel  
Attribute Parallel  
Reduction Parallel  
**Parameter Parallel**

Data Parallel  
Model Parallel  
Attribute Parallel  
Reduction Parallel  
Parameter Parallel  
Pipeline Parallel

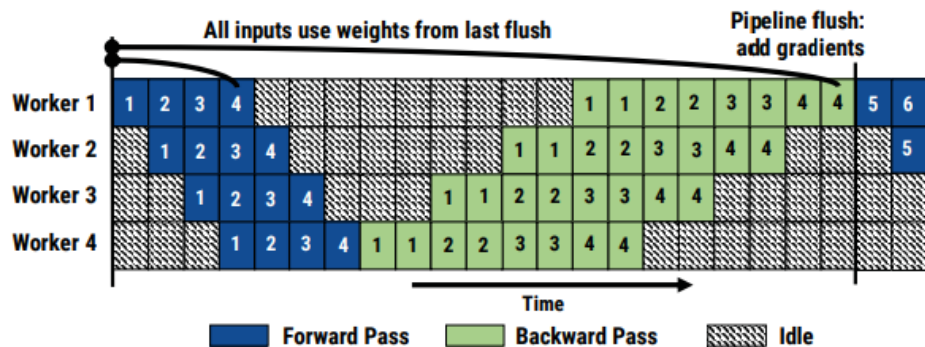


Figure 3: GPipe's inter-batch parallelism approach. Frequent pipeline flushes lead to increased idle time.

Data Parallel  
Model Parallel  
Attribute Parallel  
Reduction Parallel  
Parameter Parallel  
Pipeline Parallel  
...

2 | Parallelization |

## Algebraic Transformations

Operator Fusion  
Operator Splitting  
Operator Reordering  
...

## Auto-Parallelization

FlexFlow [MLSys 19]

Tofu [EuroSys 19]

PipeDream [SOSP 19]

automap [arXiv 19]

Whale [arXiv 21]

Alpa [OSDI 22]

...

## Algebraic Optimizers

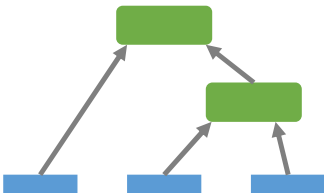
MetaFlow [MLSys 19]

TASO [SOSP 19]

PET [OSDI 21]

Tensat [MLSys 21]

...

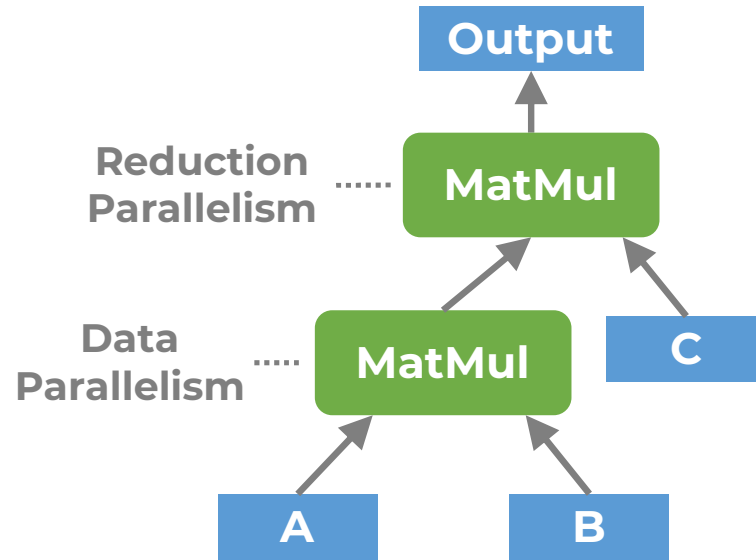


**Auto-  
Parallelization**

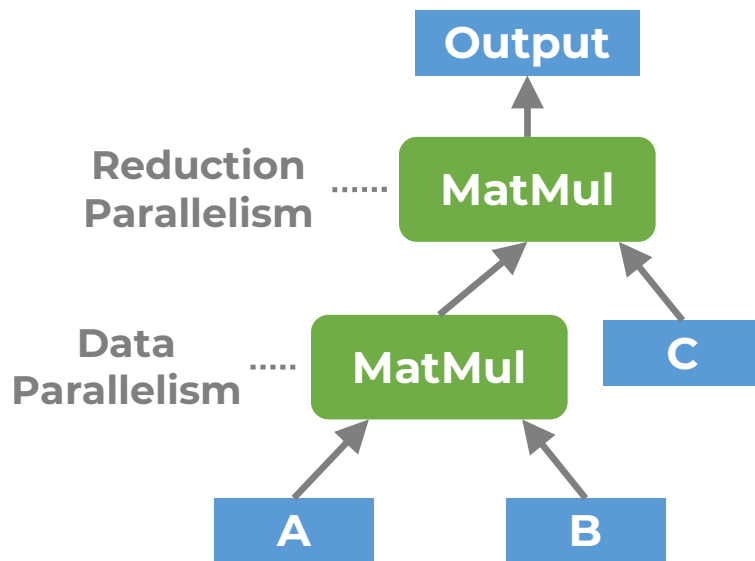
**Algebraic  
Optimizer**

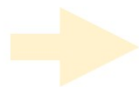
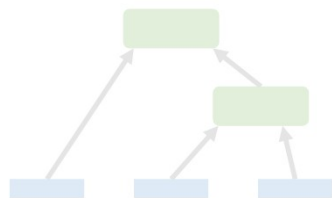
?



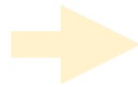


**“annotated computation graph”**

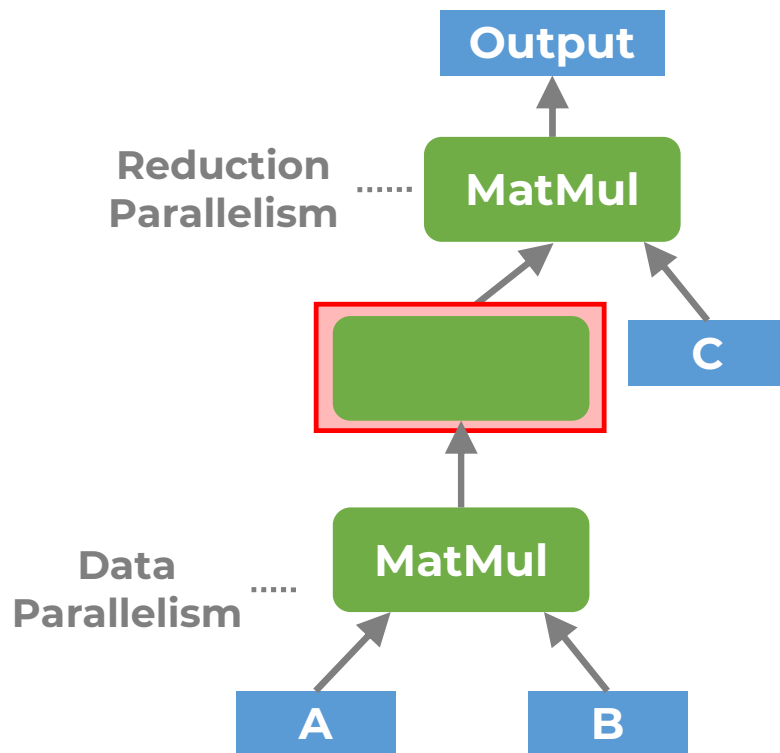
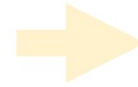


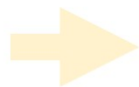
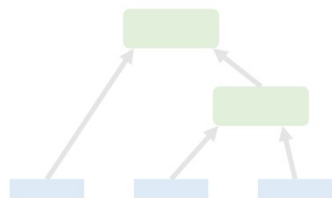


Auto-Parallelization

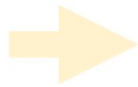


Algebraic Optimizer

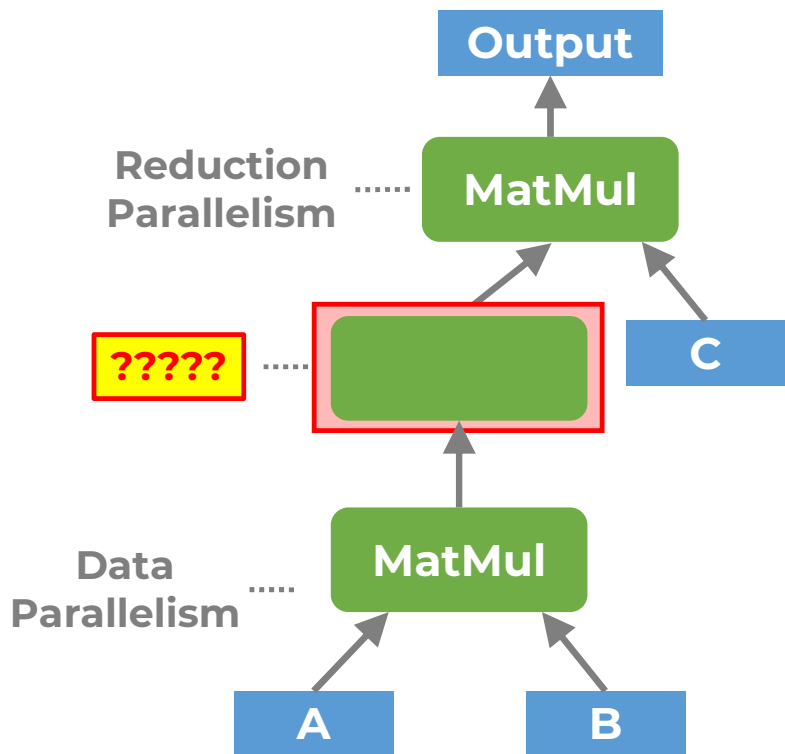
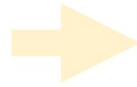


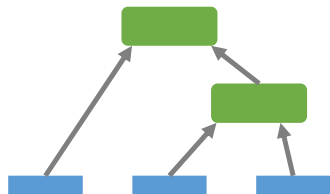


Auto-Parallelization



Algebraic Optimizer



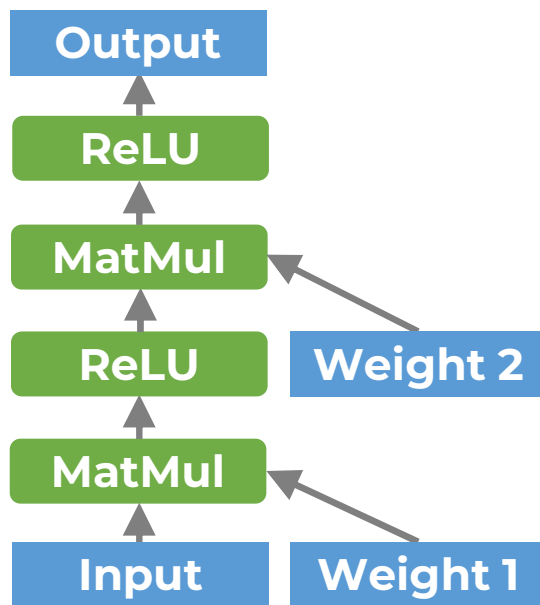


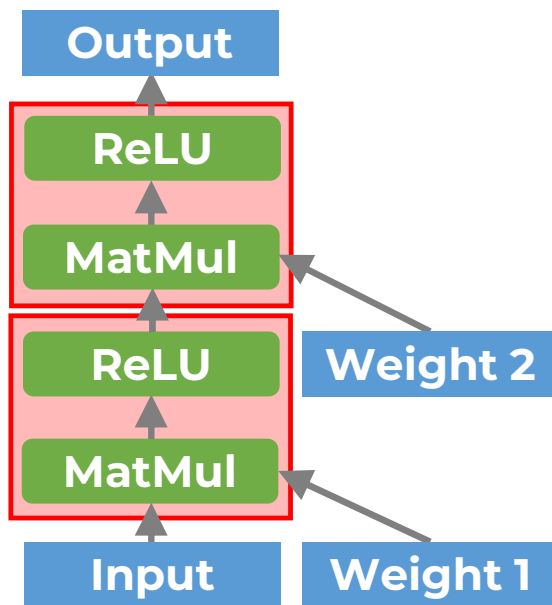
**Algebraic  
Optimizer**

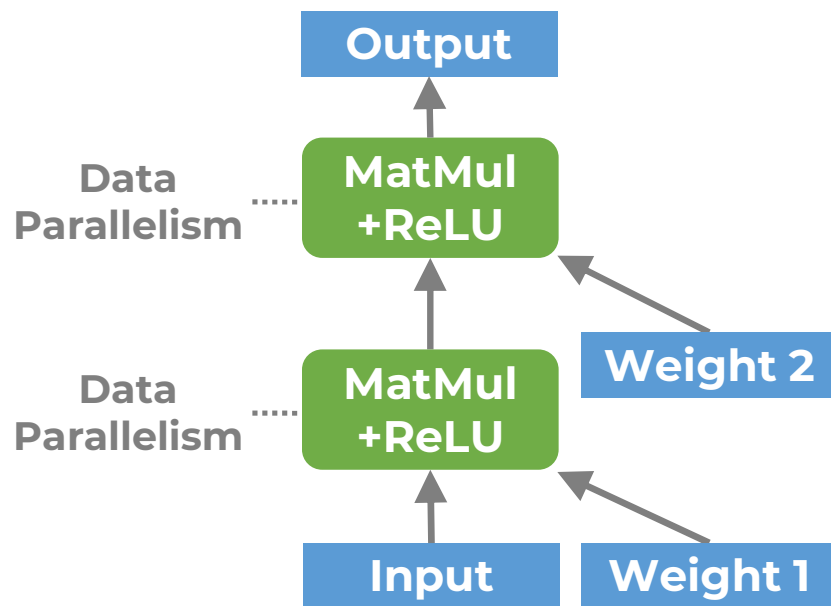


**Auto-  
Parallelization**

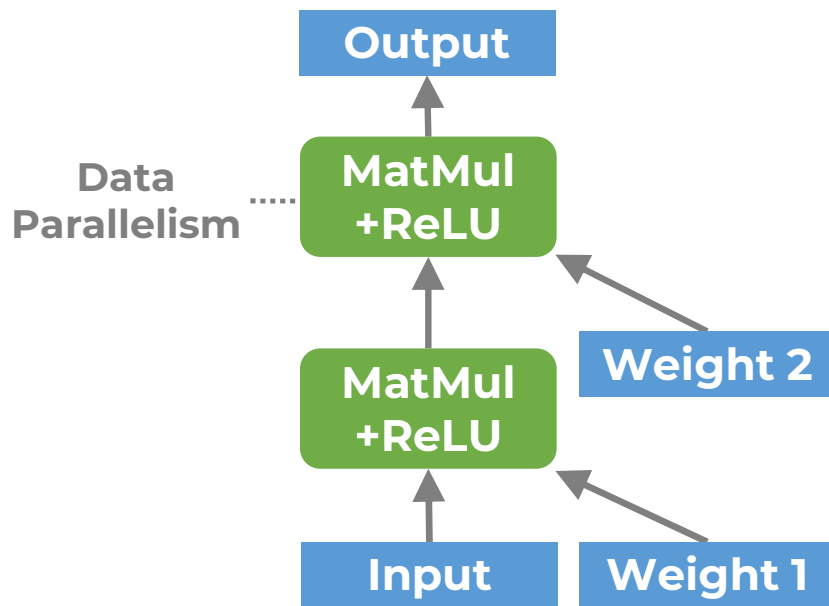


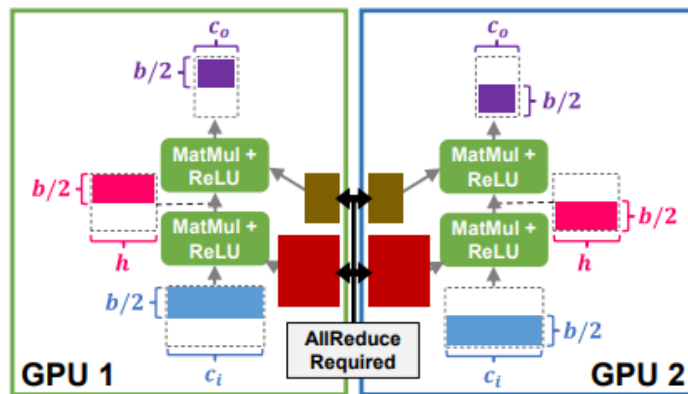
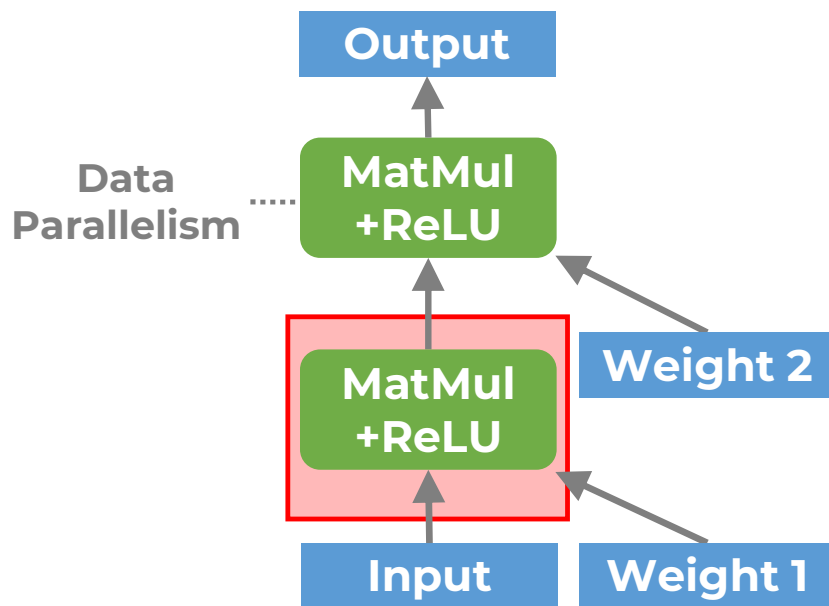


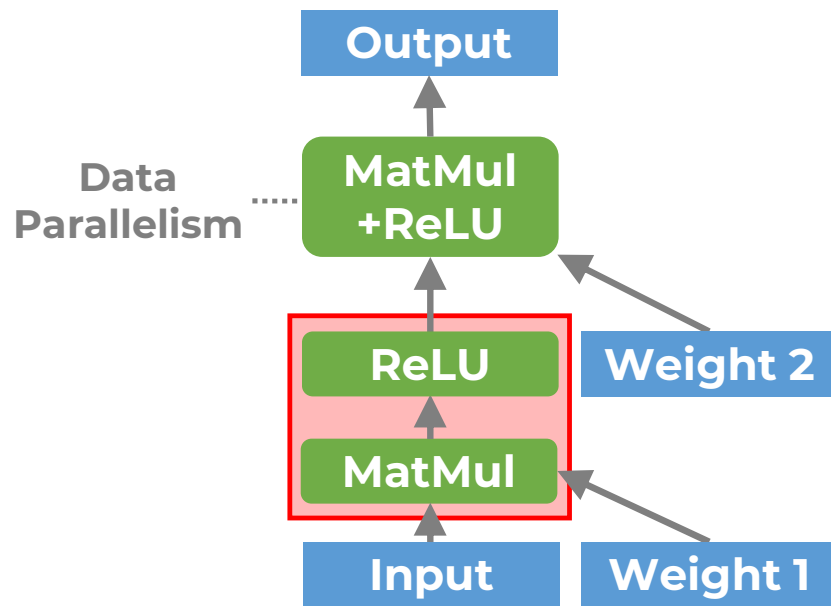


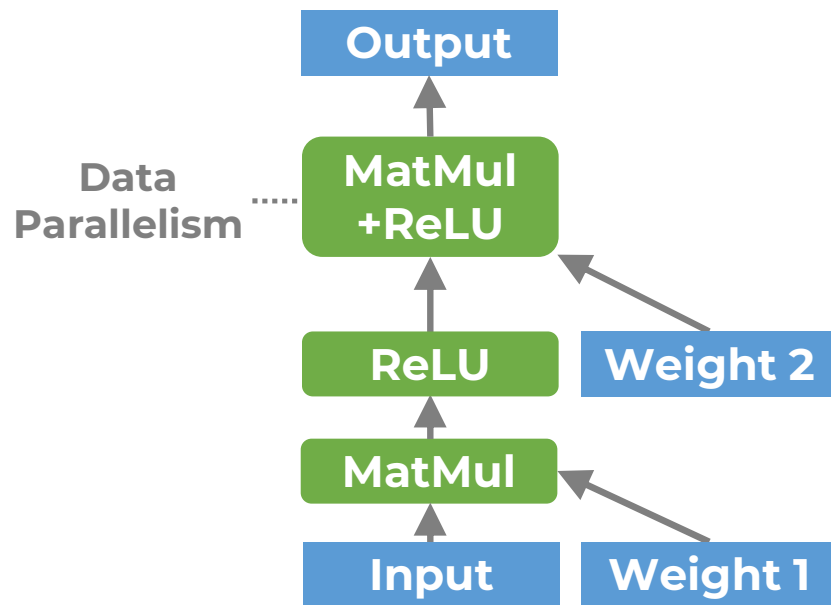


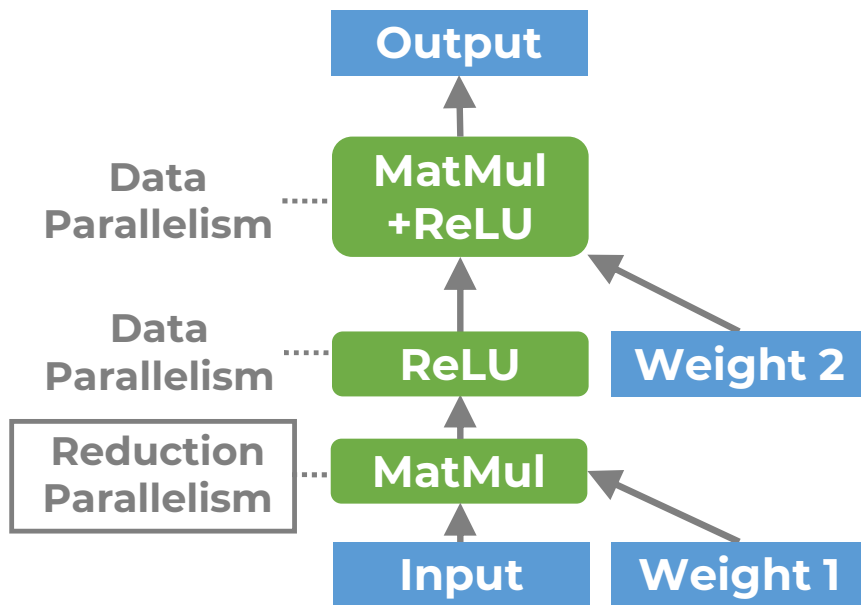


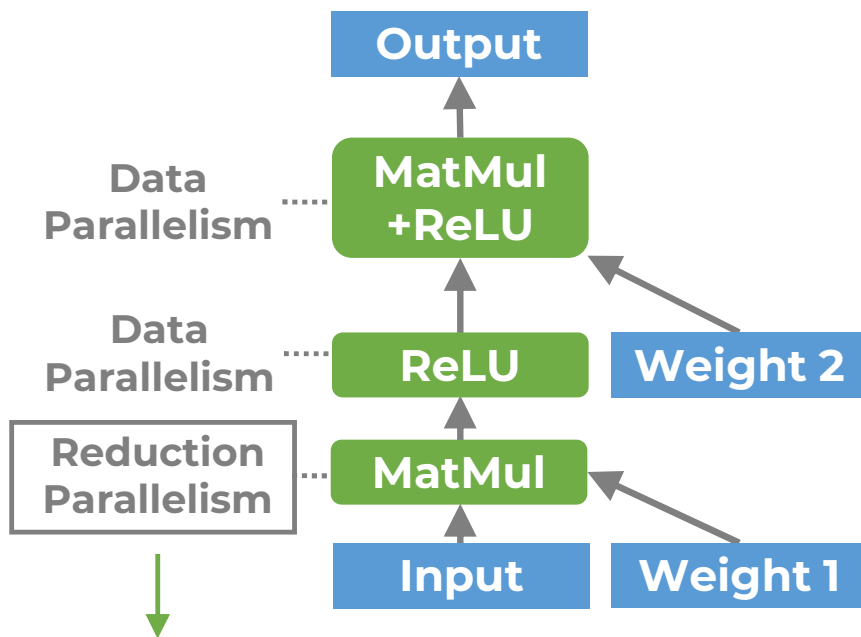
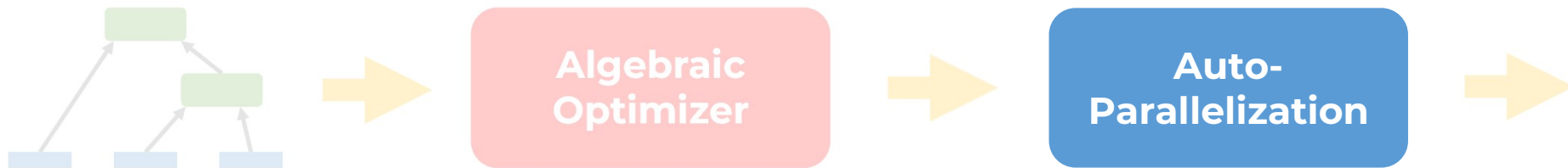




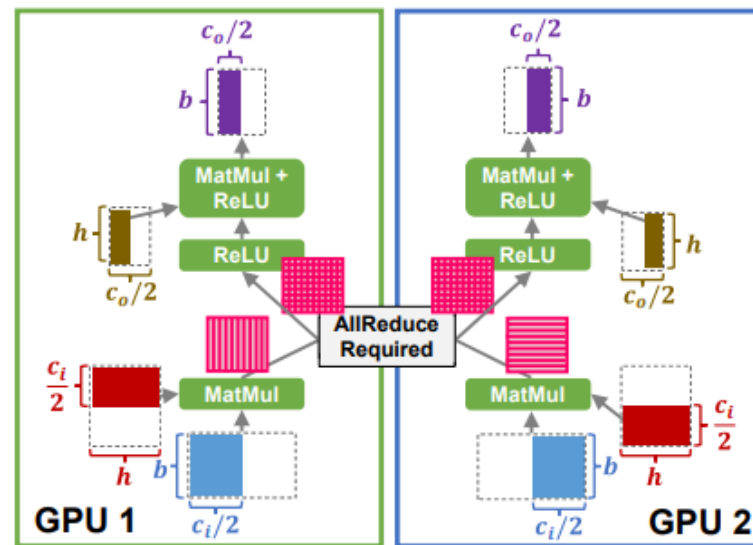








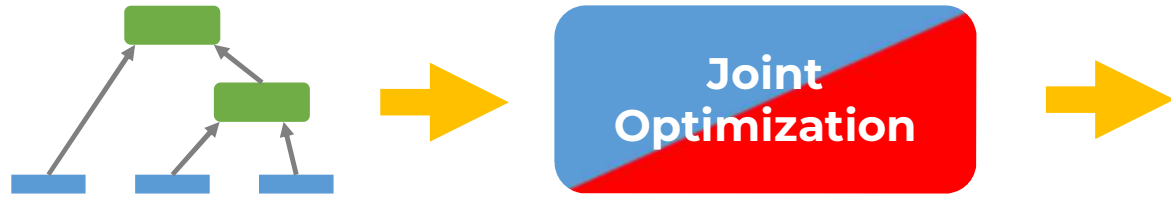
**$\approx 6 \times$  less communication!**





# 1. Representation

## 2.



- 
1. Representation
  2. Scalability



# Unity

Unity — Representation

Representation

Parallel Computation  
Graph (PCG)

Unity

A diagram consisting of a horizontal line extending from the right side of the word 'Unity', followed by a vertical line going upwards, and then a horizontal line extending to the left, connecting to the text 'Parallel Computation Graph (PCG)'.

Representation

Parallel Computation  
Graph (PCG)

Unity

Scalability

Representation

Parallel Computation  
Graph (PCG)

Unity

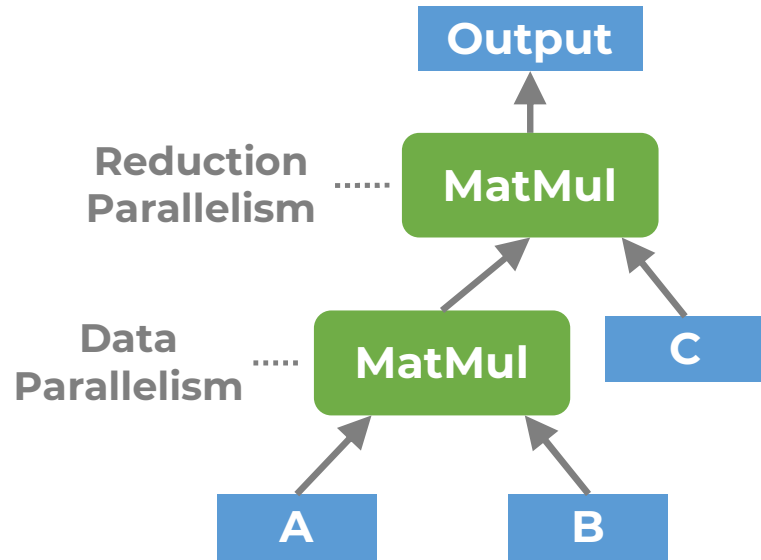
Scalability

Hierarchical Search  
Algorithm

Representation

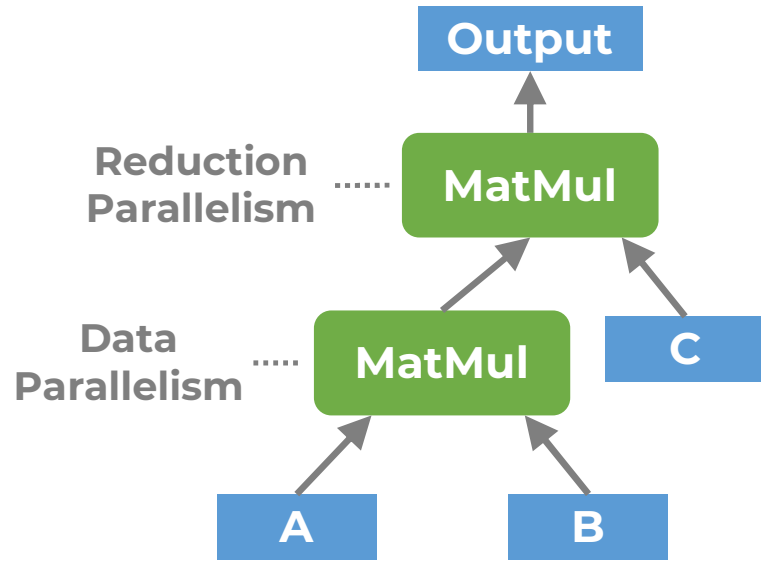
# Parallel Computation Graph (PCG)

# Parallel Computation Graph (PCG)

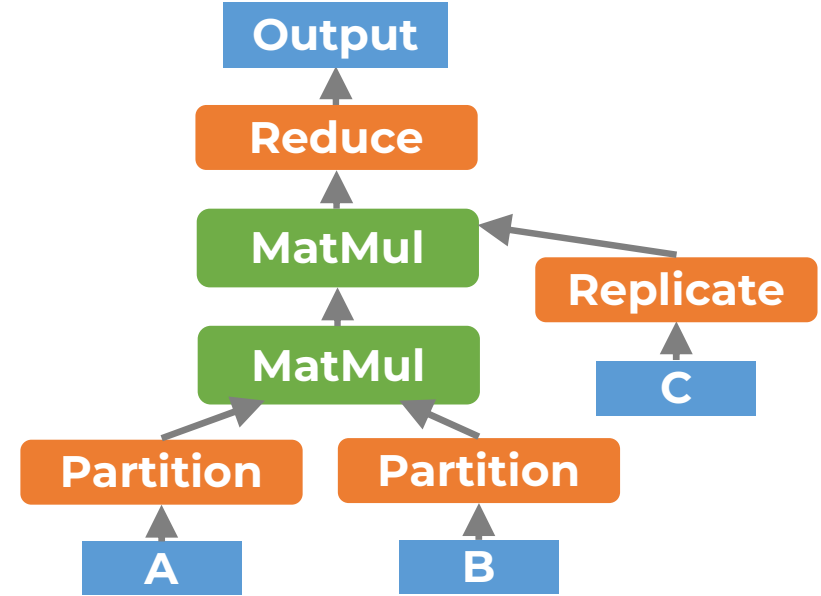


**annotated computation graph**

# Parallel Computation Graph (PCG)



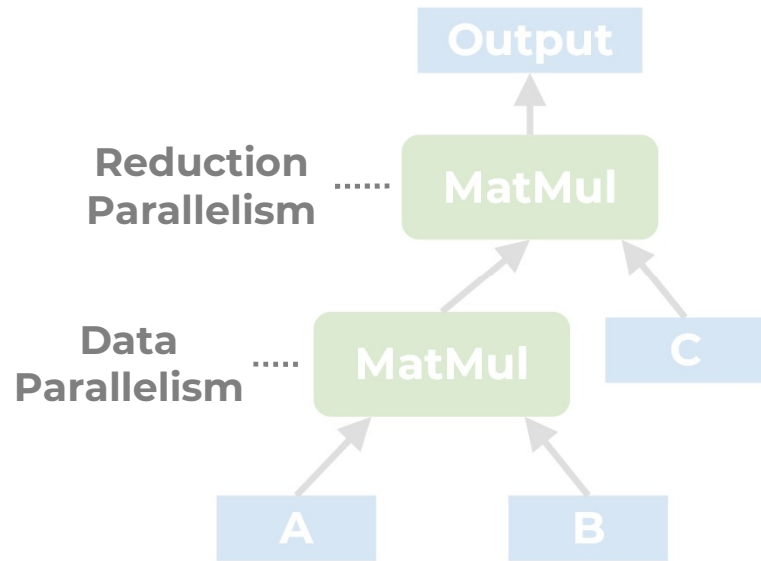
annotated computation graph



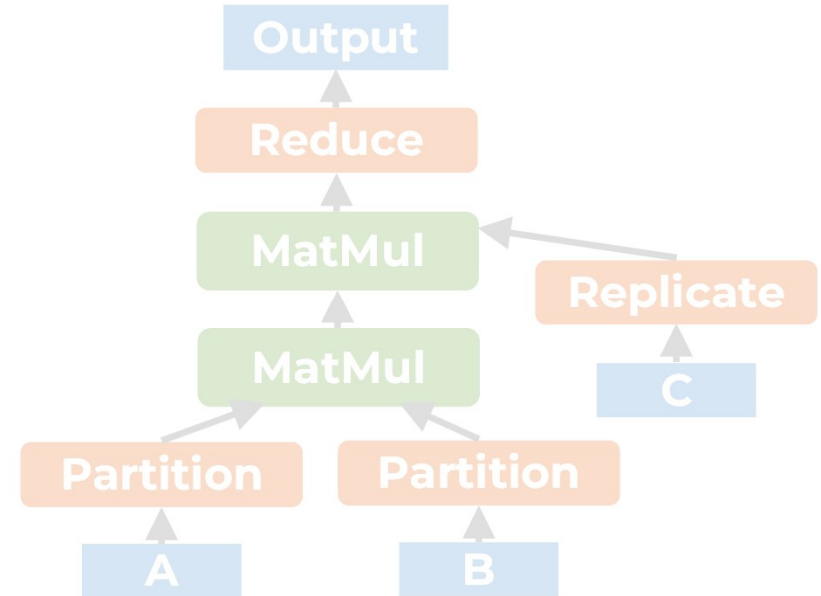
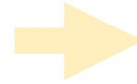
parallel computation graph (PCG)



# Parallel Computation Graph (PCG)

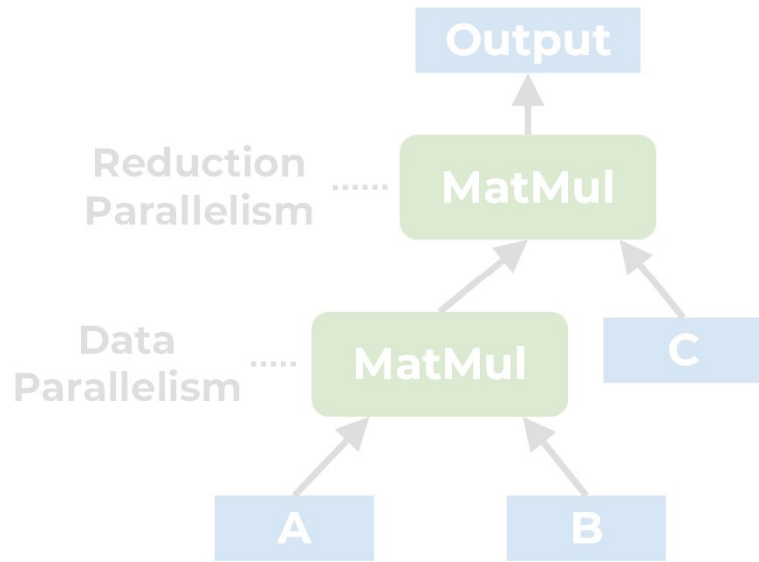


annotated computation graph

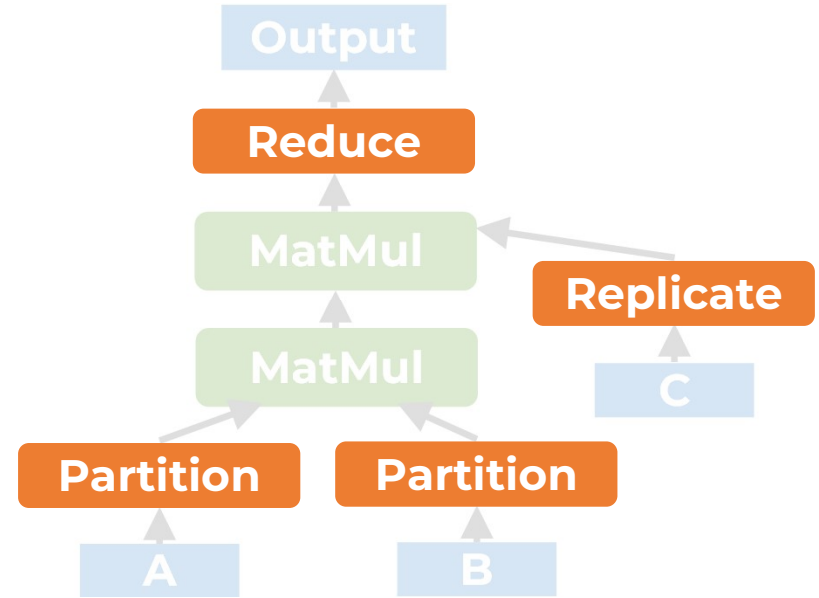
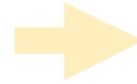


parallel computation graph (PCG)

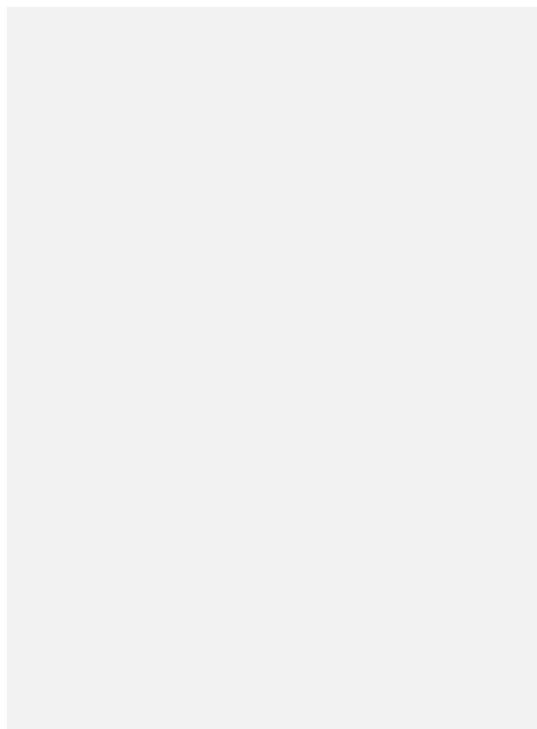
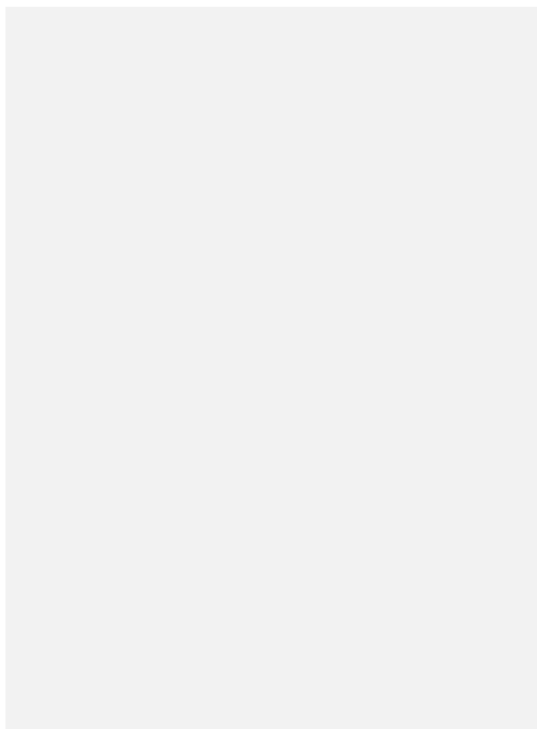
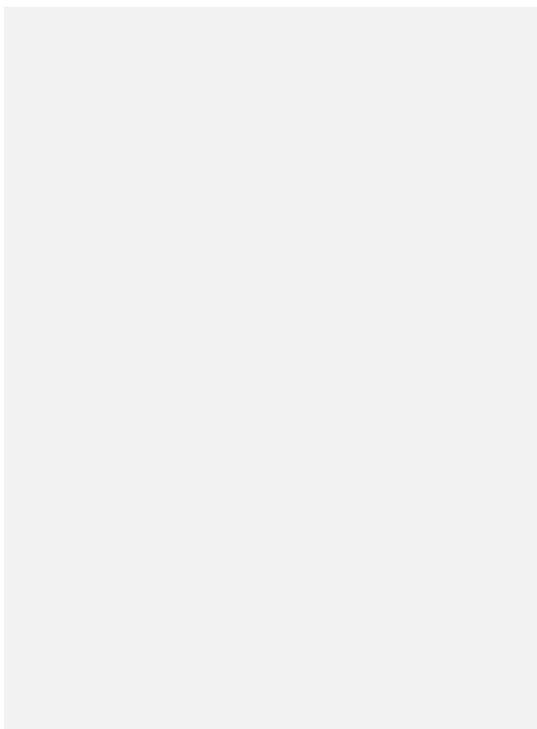
# Parallel Computation Graph (PCG)



annotated computation graph

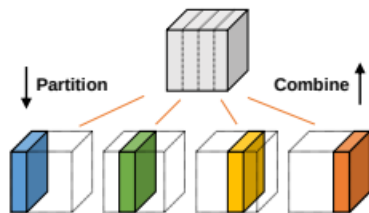


parallel computation graph (PCG)



**Partition**

**Combine**

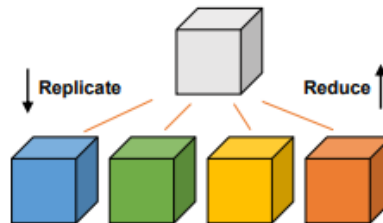
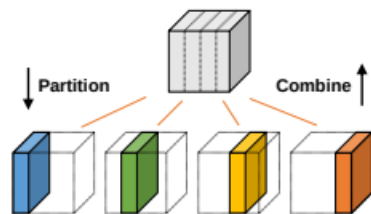


Partition

Replicate

Combine

Reduce



Partition

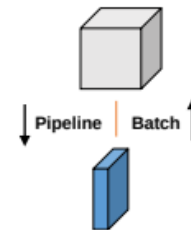
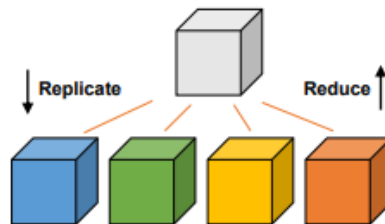
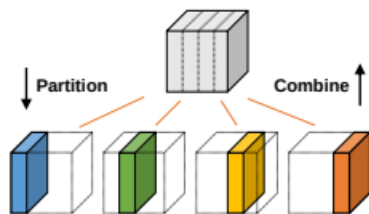
Replicate

Pipeline

Combine

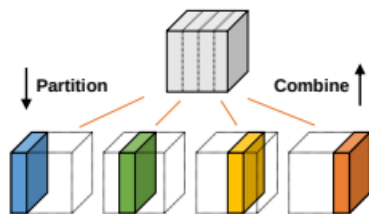
Reduce

Batch



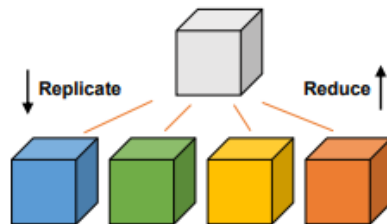
**Partition**

**Combine**



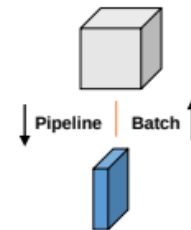
**Replicate**

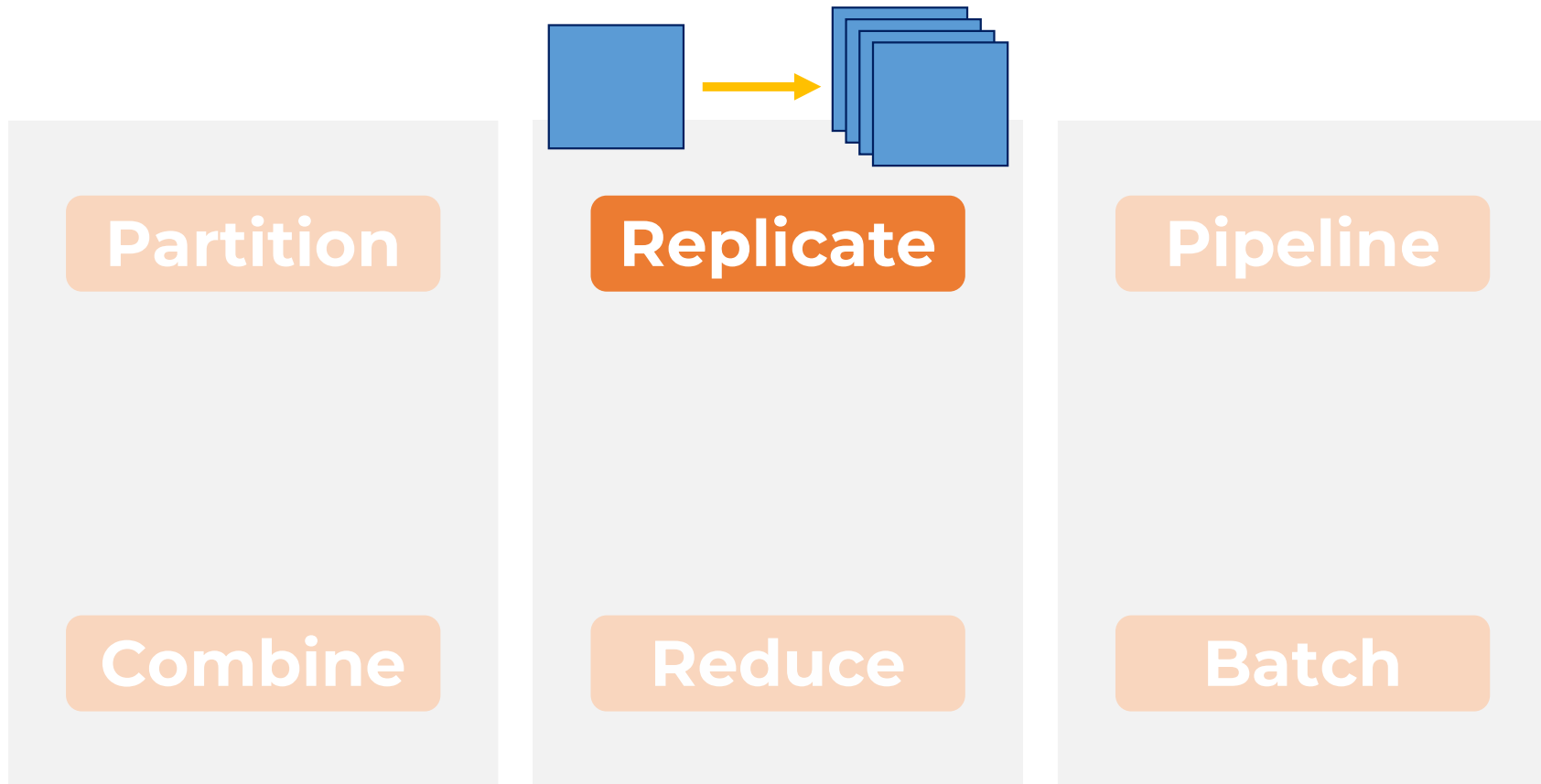
**Reduce**



**Pipeline**

**Batch**







**Partition**

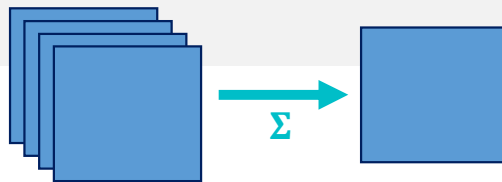
**Replicate**

**Pipeline**

**Combine**

**Reduce**

**Batch**



**Partition**



**Replicate**



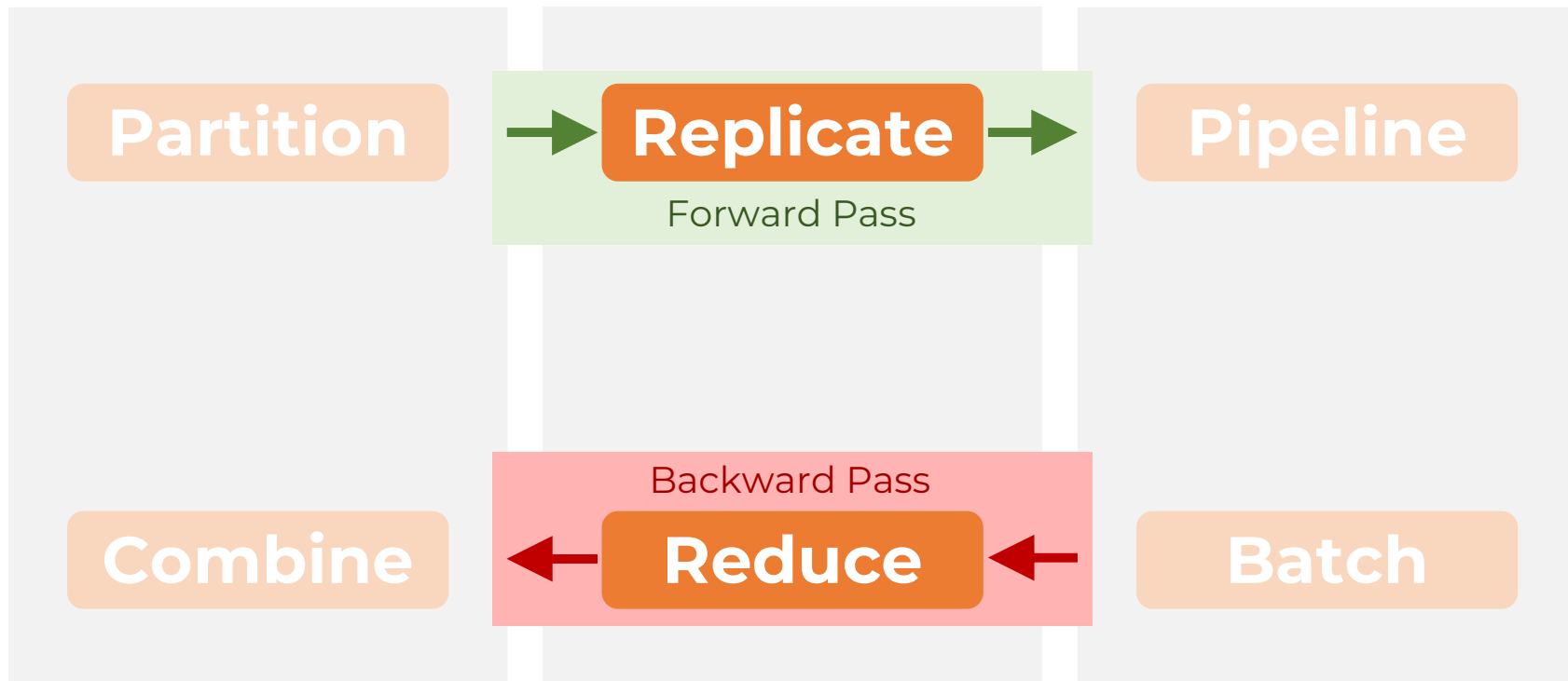
**Pipeline**

Forward Pass

**Combine**

**Reduce**

**Batch**



**Partition**

**Combine**

**Replicate**

**Reduce**

**Pipeline**

**Batch**

**Partition**

**Replicate**

**Pipeline**

**Combine**

Forward Pass



**Reduce**



**Batch**

**Partition**

**Replicate**

Backward Pass

**Pipeline**

**Combine**

Forward Pass

**Reduce**

**Batch**

**Partition**

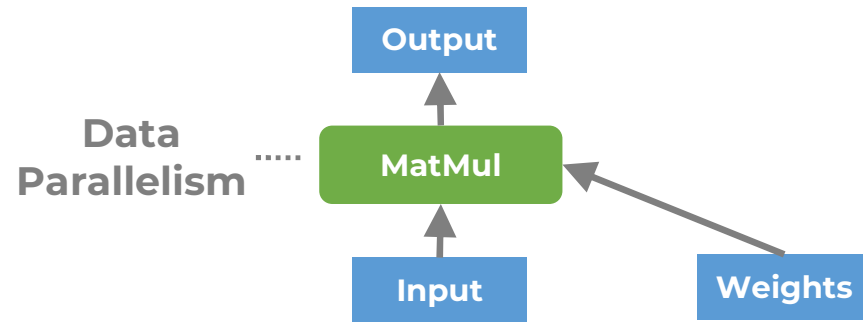
**Combine**

**Replicate**

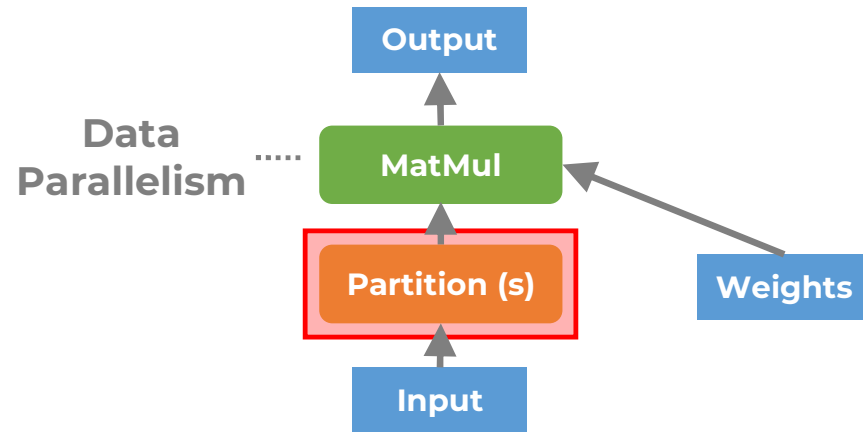
**Reduce**

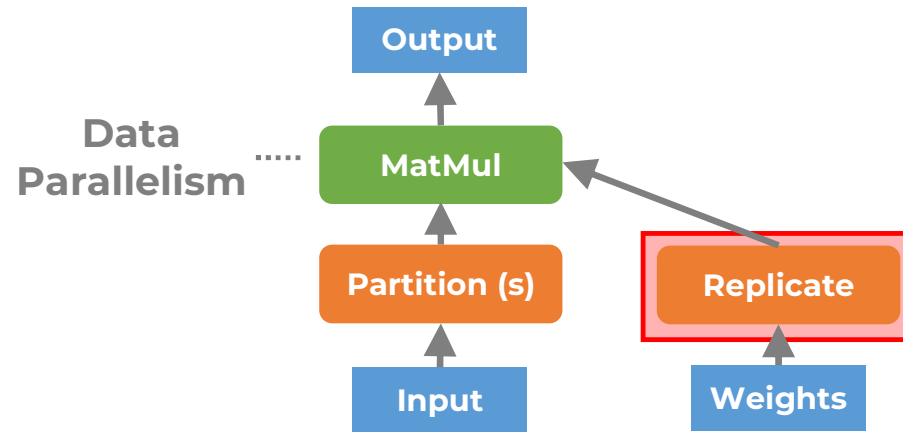
**Pipeline**

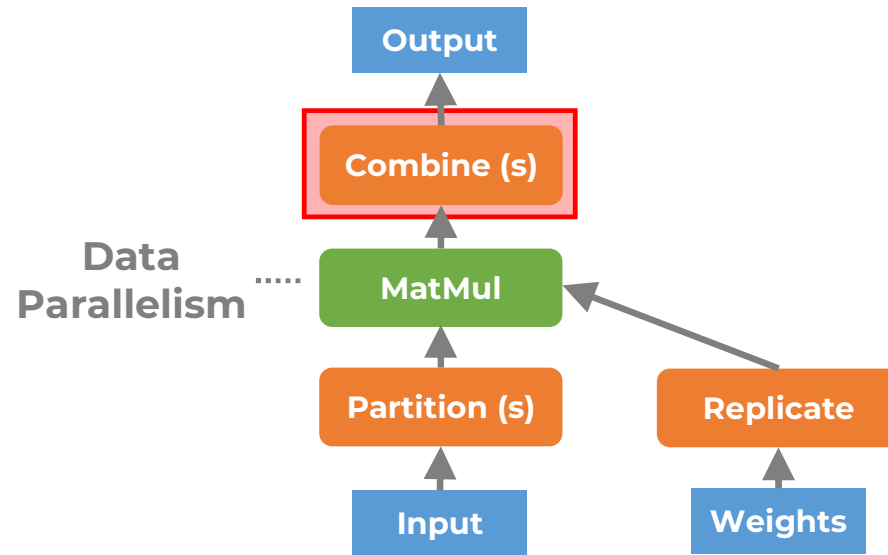
**Batch**

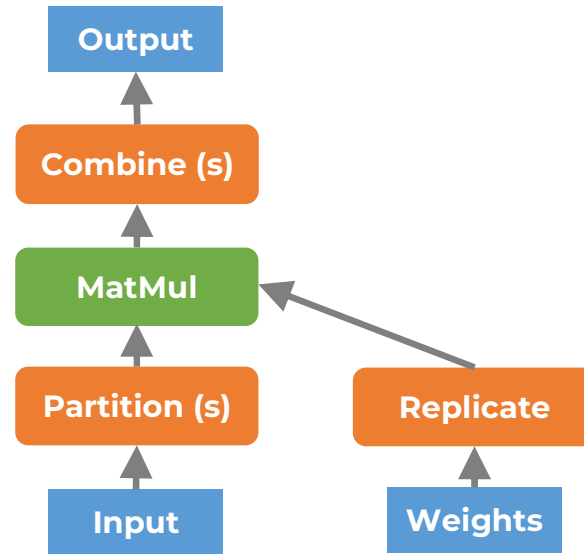


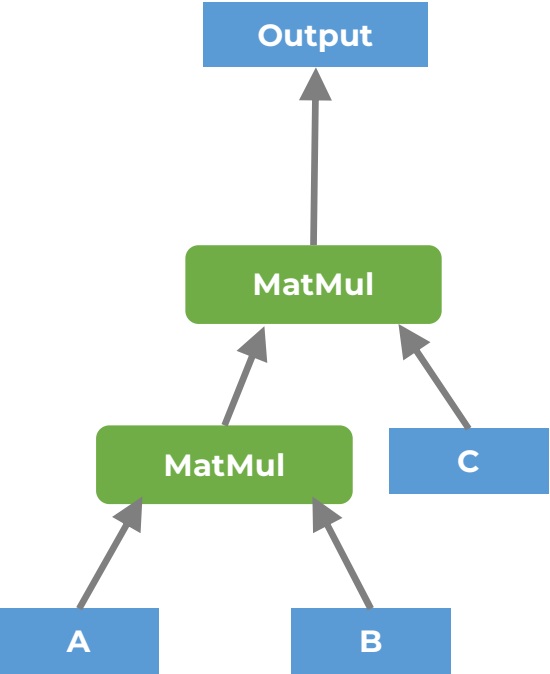


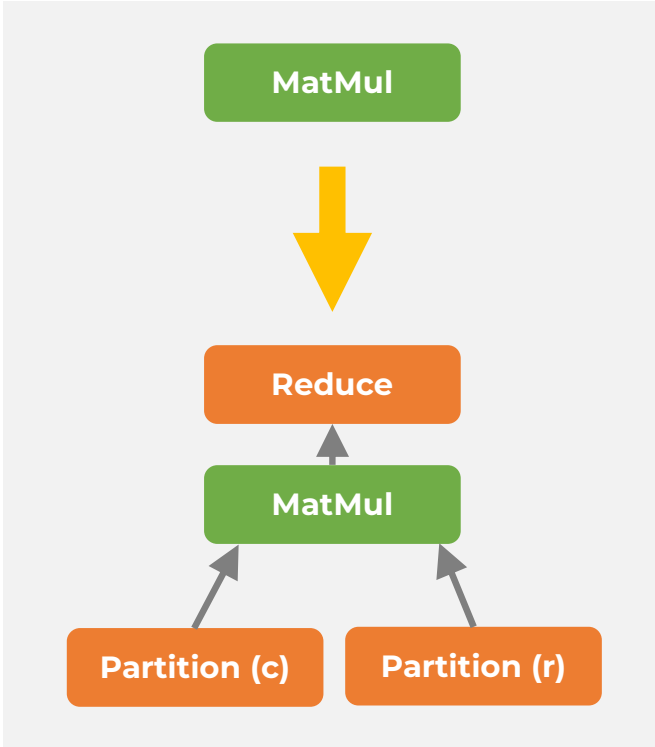
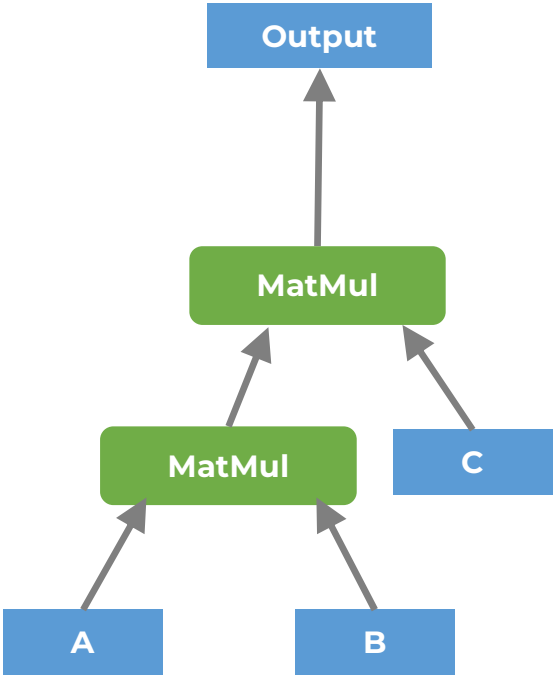


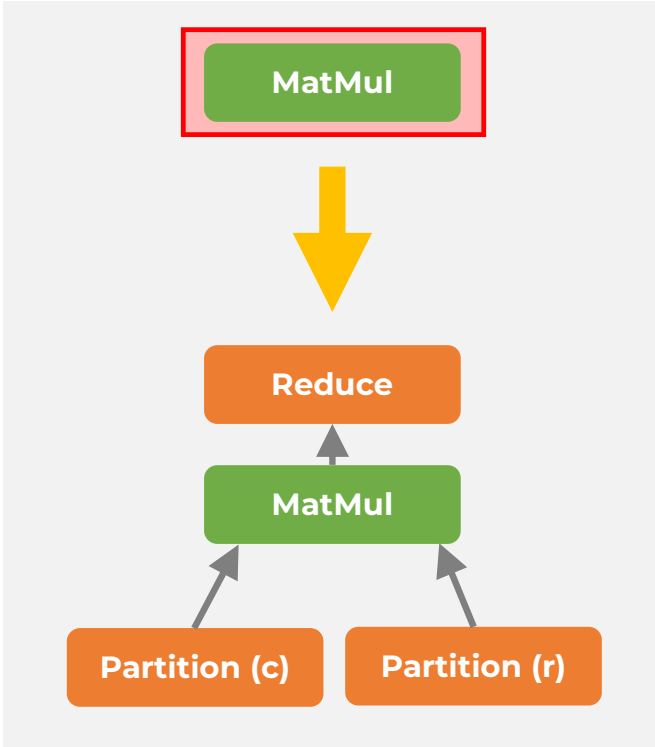
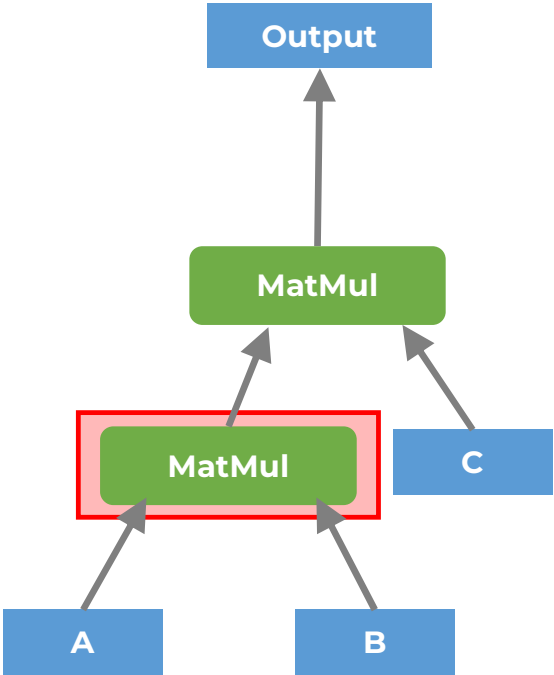


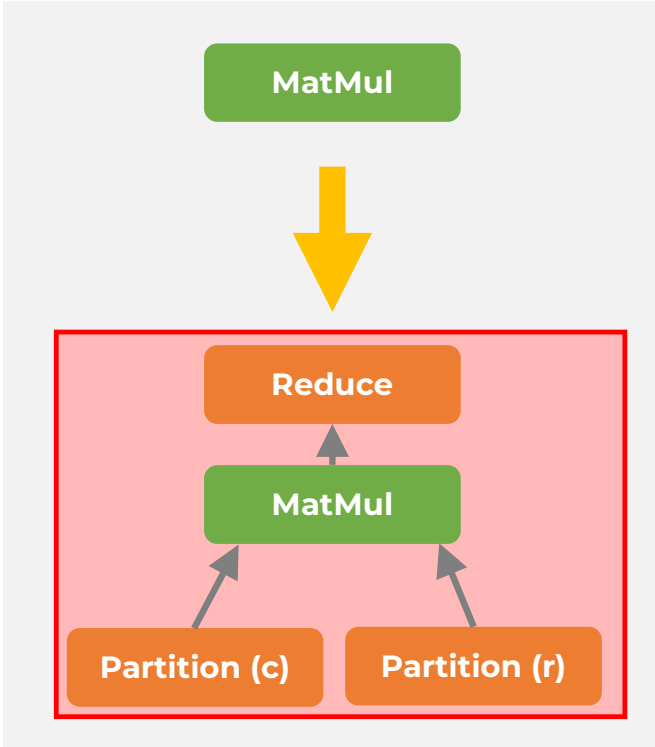
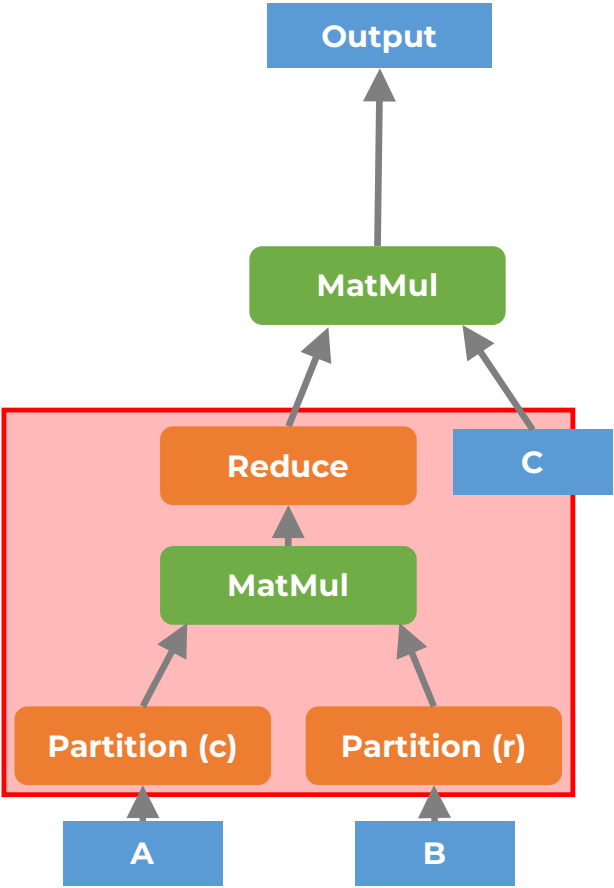




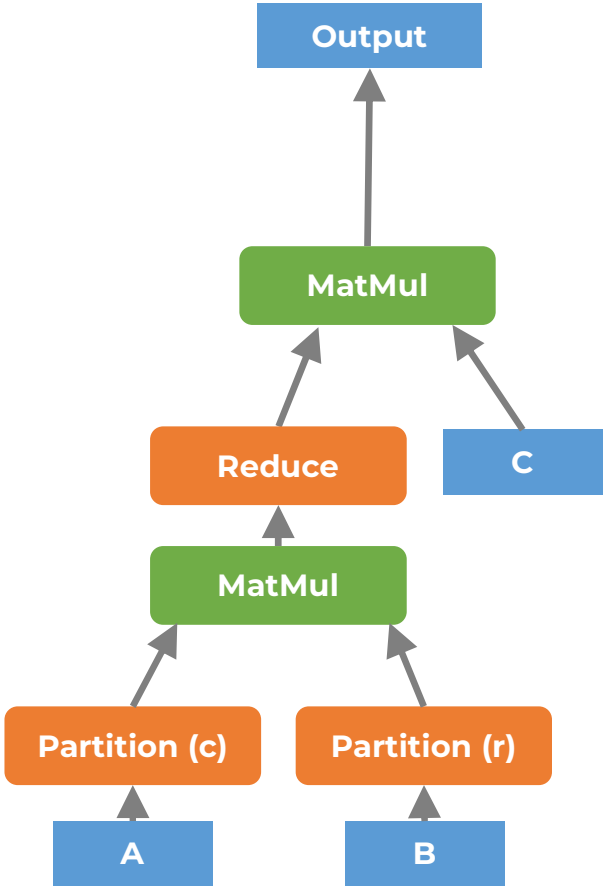


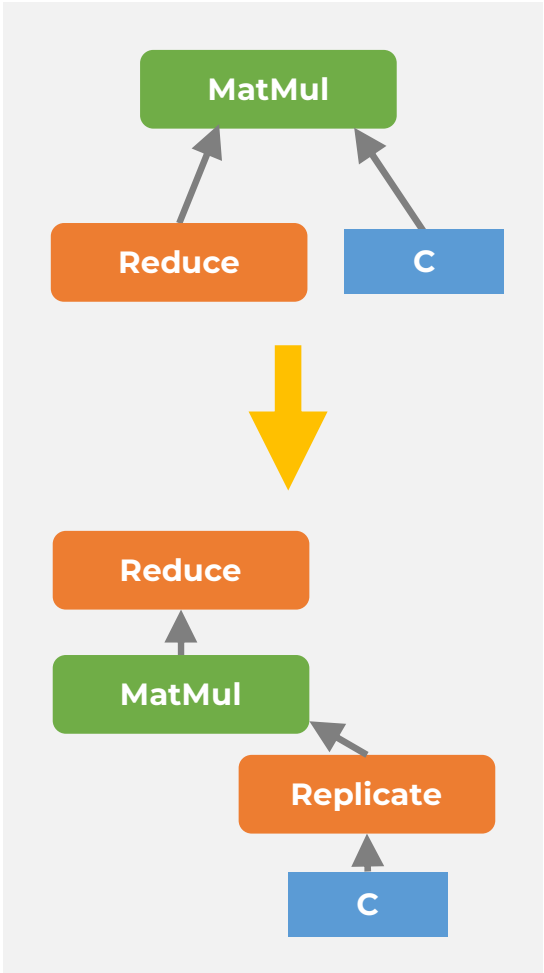
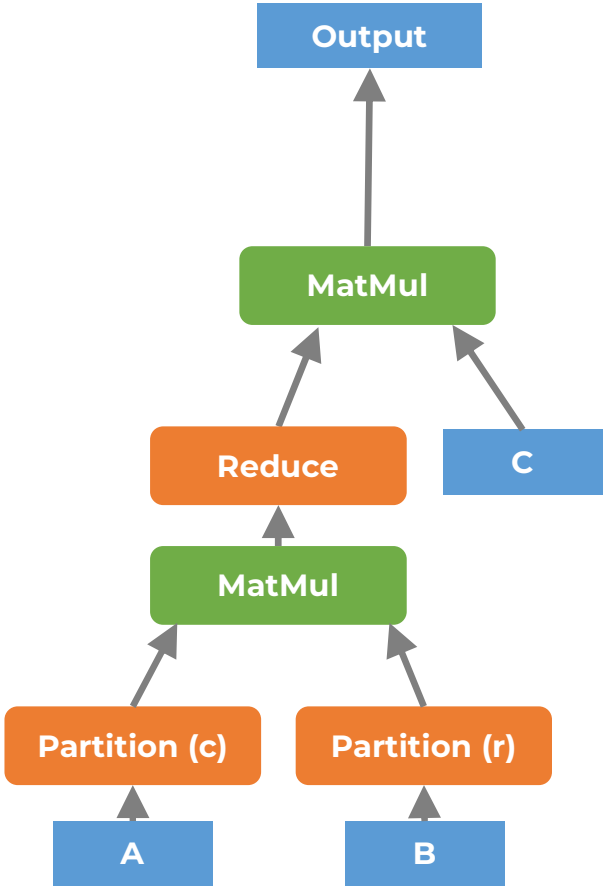


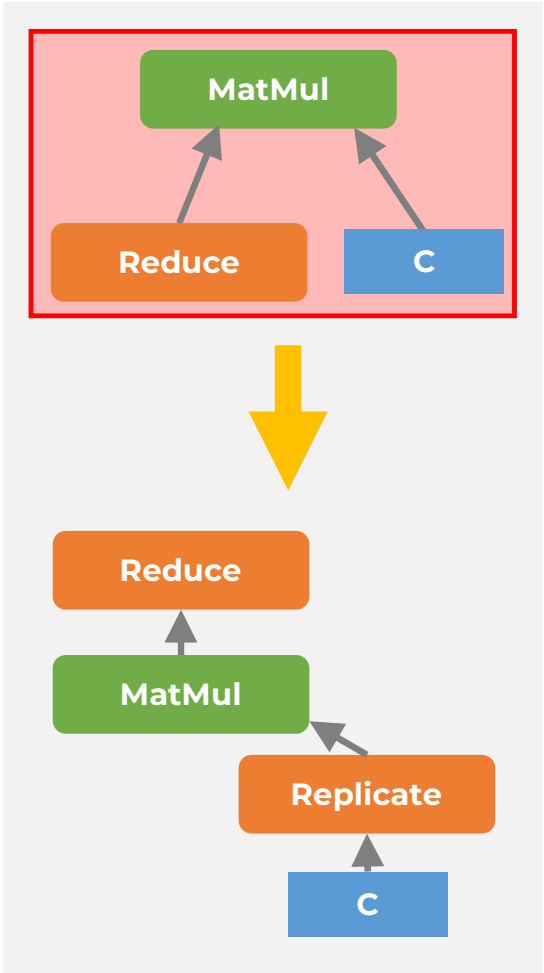
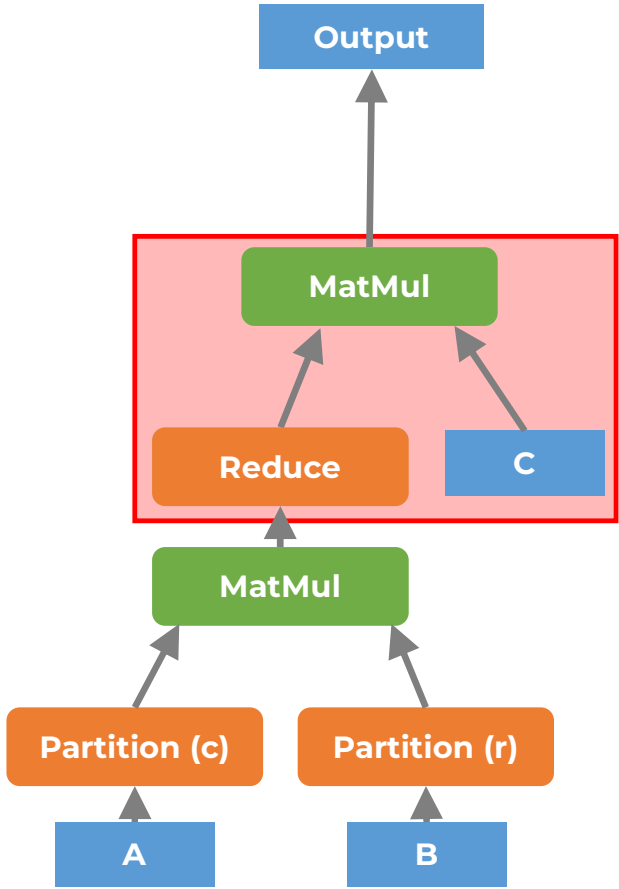


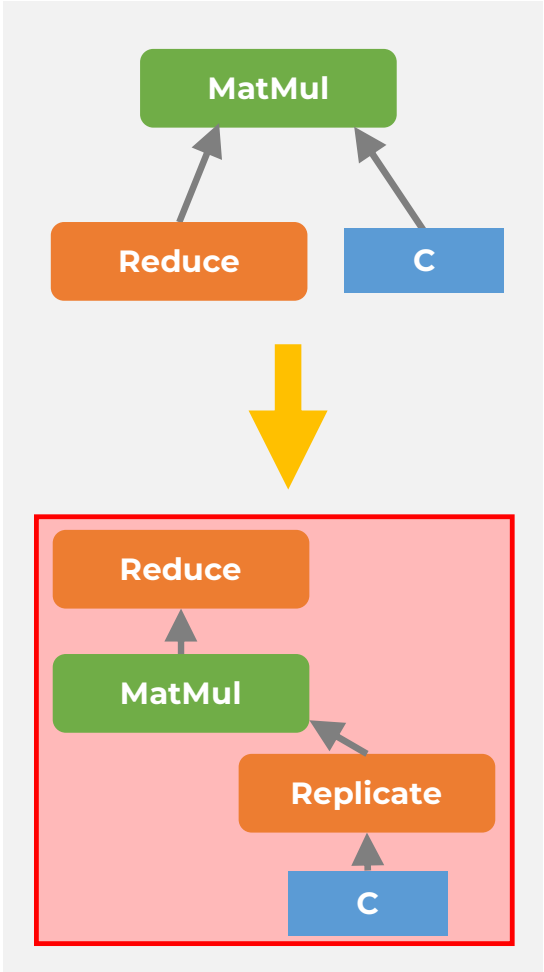
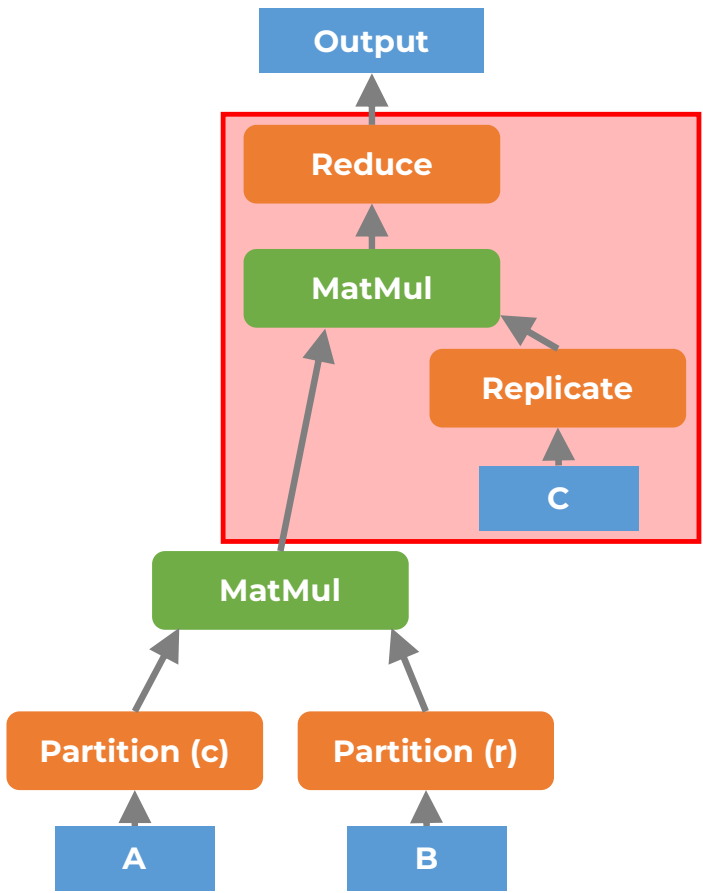


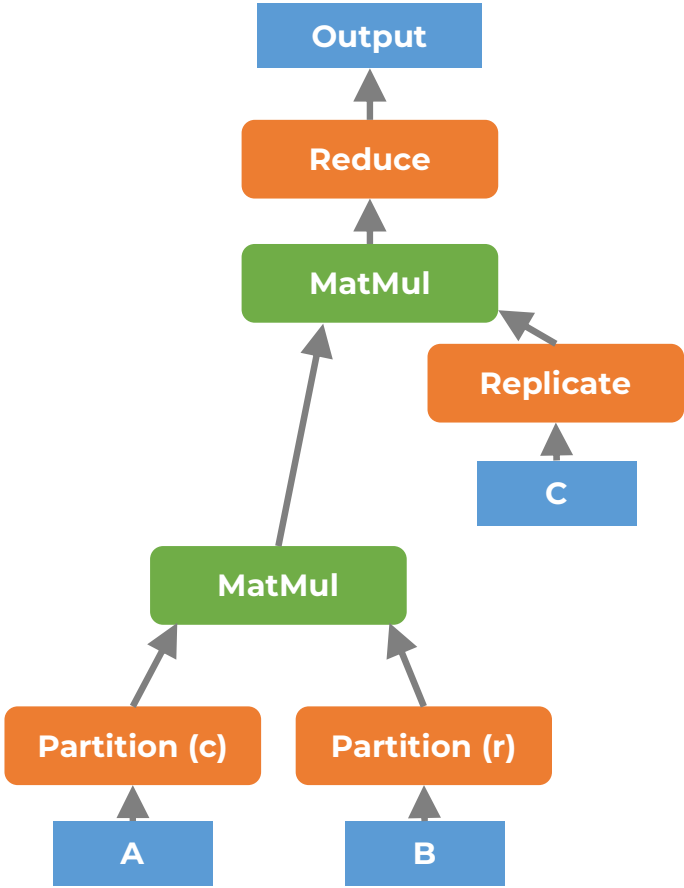














# Separation of concerns

# Separation of concerns

Automatically generate substitutions



# Separation of concerns

Automatically generate substitutions

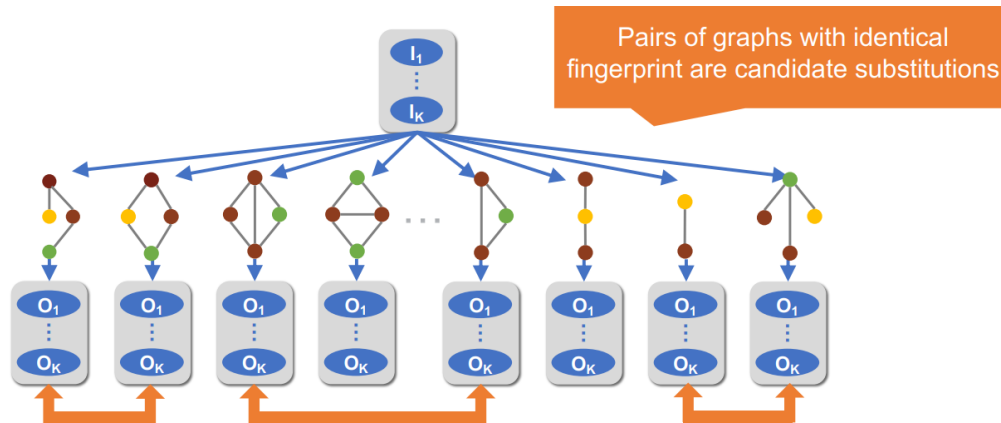
New operators

# Separation of concerns

Automatically generate substitutions

New operators

New forms of parallelism



Separation of concerns

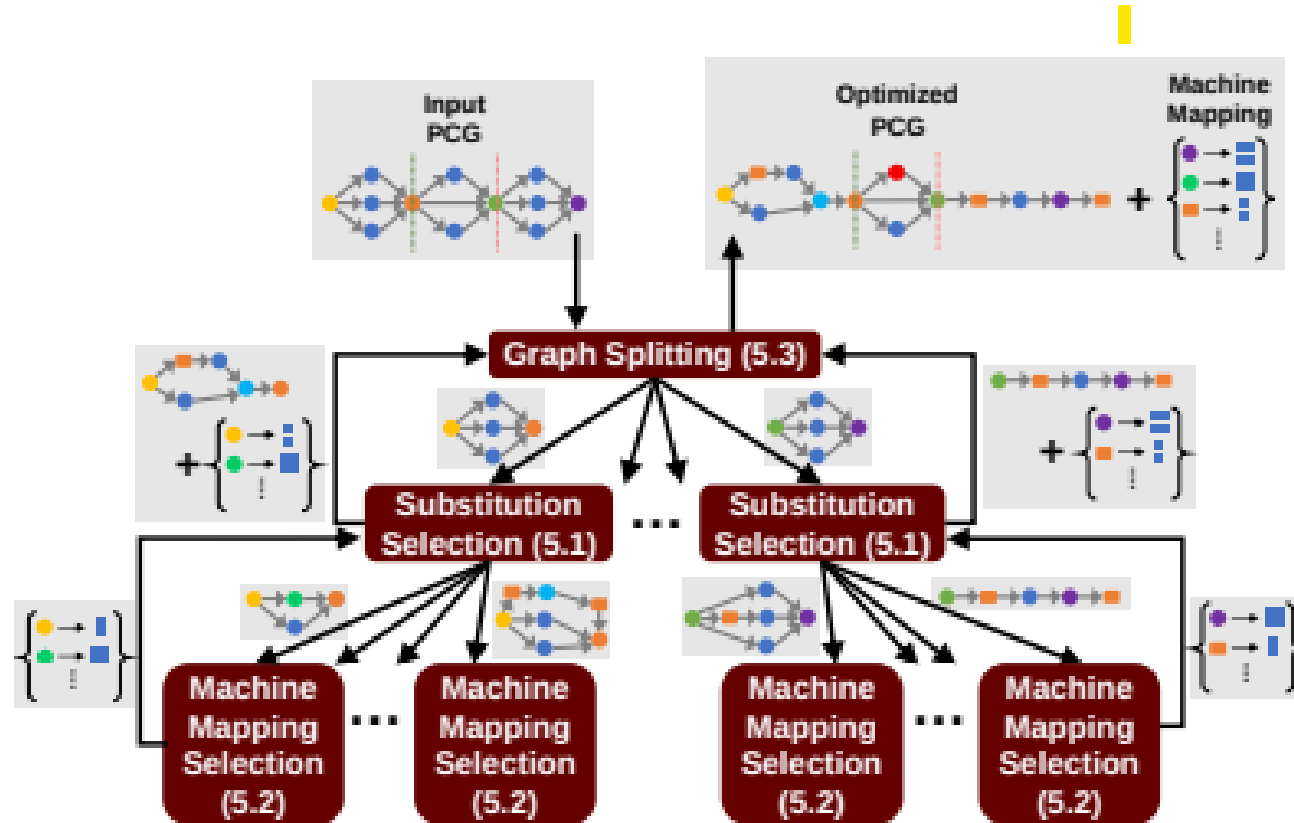
Explicitly represents communication

Separation of concerns

Explicitly represents communication

Concise

# Hierarchical Search Algorithm



# Hierarchical Search Algorithm

Algebraic Transformation

Parallelism Type

Parallelism Degree

Device Mapping

# Hierarchical Search Algorithm

---

Algebraic Transformation

Parallelism Type

Parallelism Degree

---

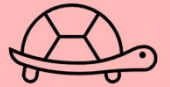
Device Mapping

# Hierarchical Search Algorithm

Algebraic Transformation

Backtracking  
Search

Parallelism Type



Parallelism Degree

Device Mapping



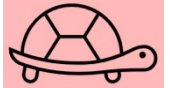
# Hierarchical Search Algorithm

Algebraic Transformation

Parallelism Type

Parallelism Degree

Backtracking  
Search



Device Mapping



Dynamic  
Programming

# Evaluation

# Models

BERT-Large

(Language Modeling)

Candle-UNO

(Precision Medicine)

MLP

(Regression)

DLRM

XDL

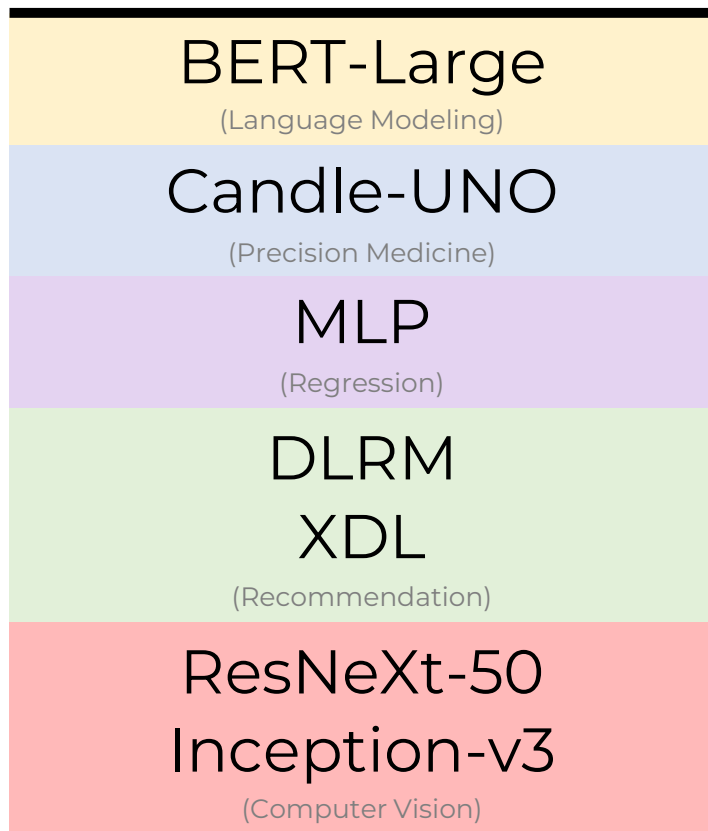
(Recommendation)

ResNeXt-50

Inception-v3

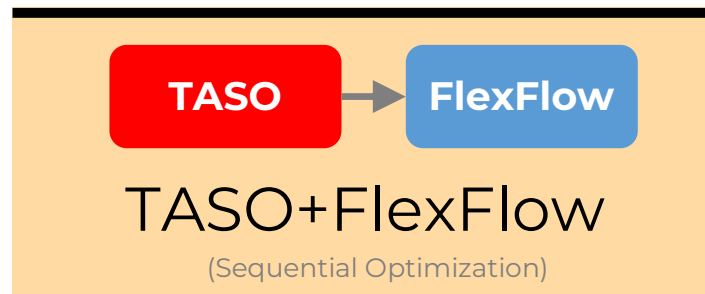
(Computer Vision)

## Models

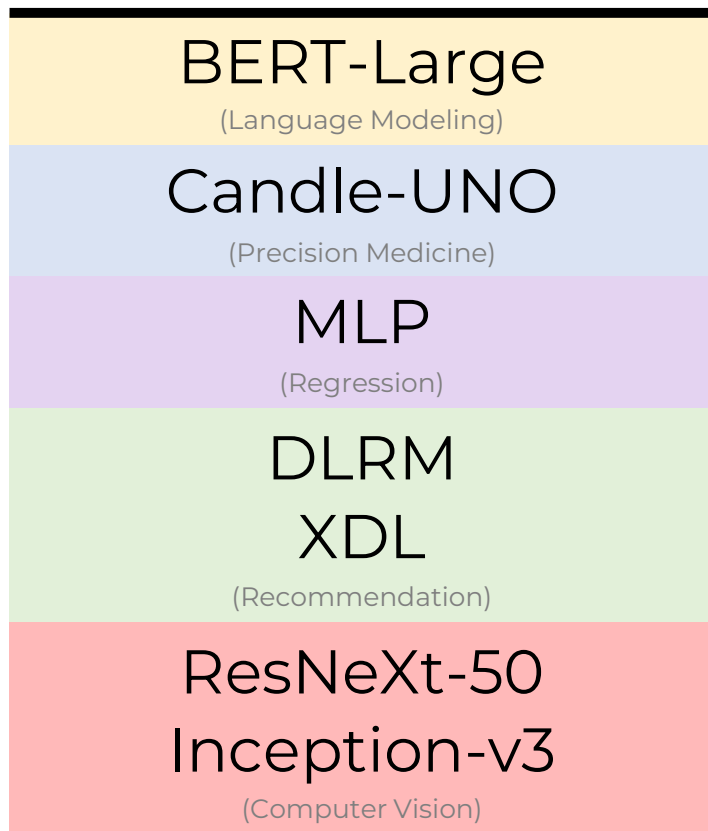


×

## Baselines

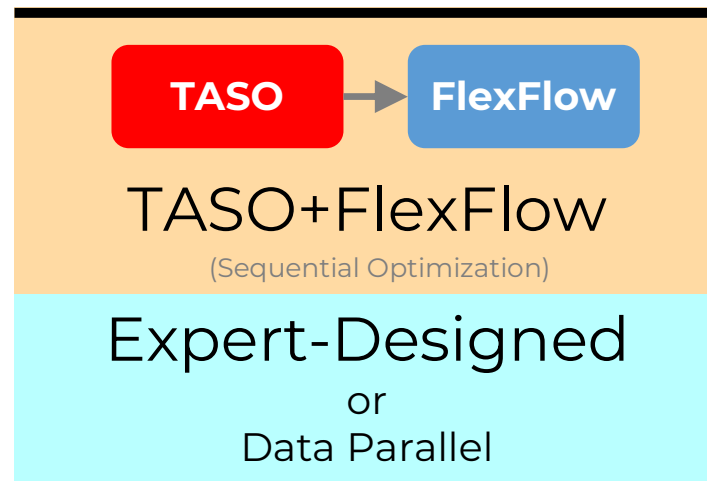


## Models

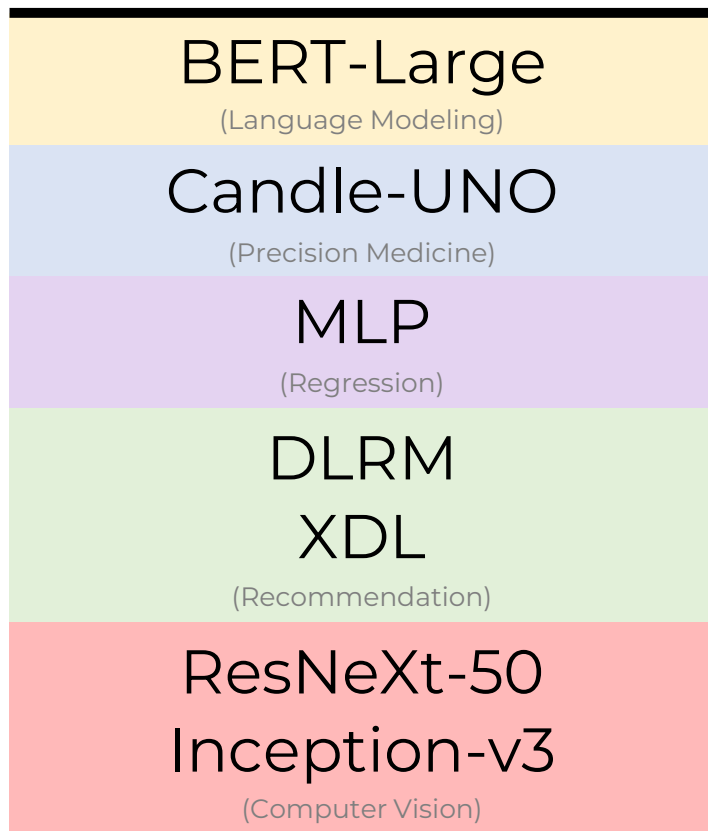


×

## Baselines

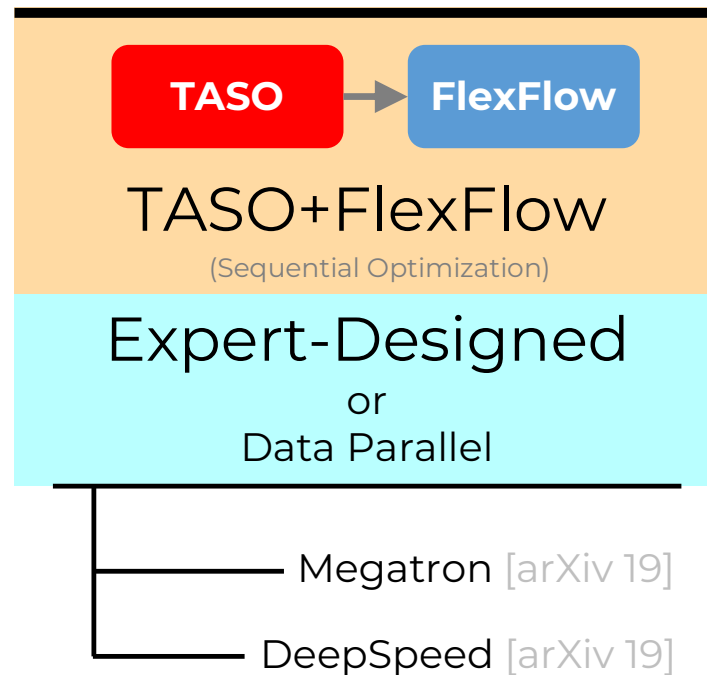


## Models

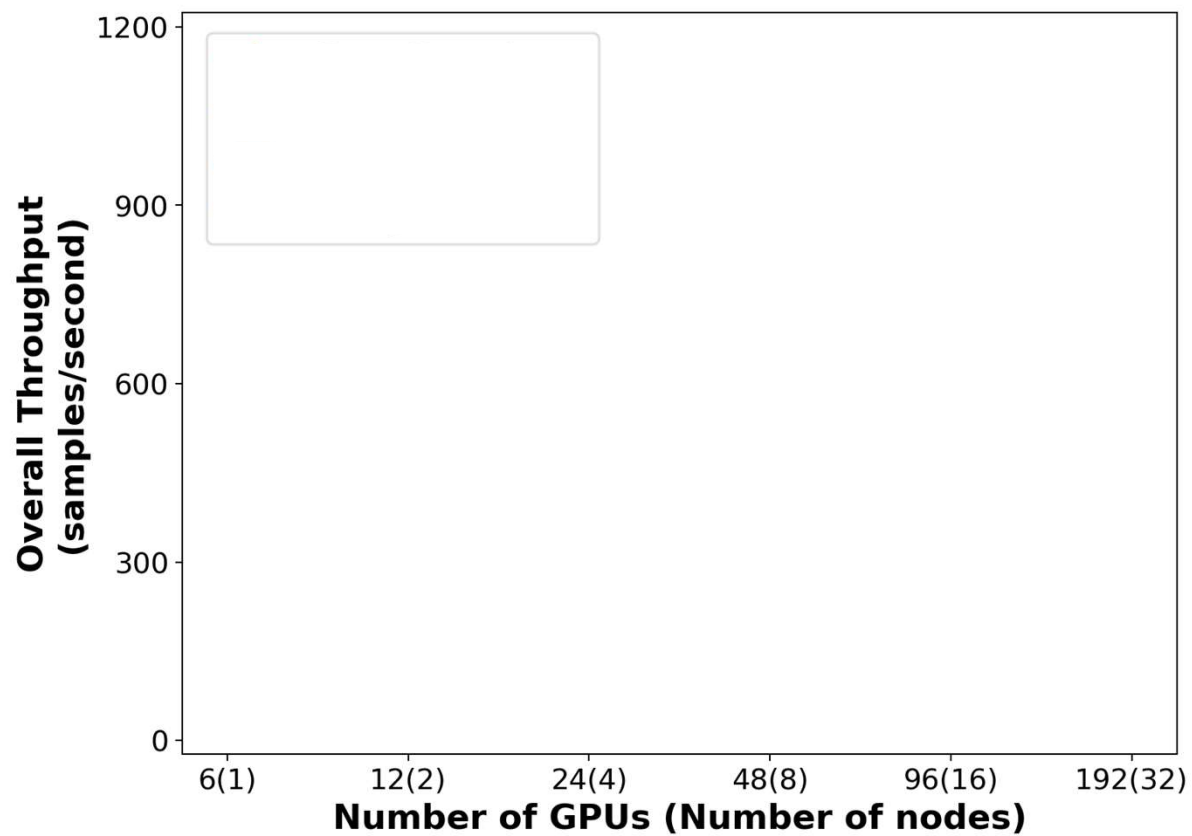


×

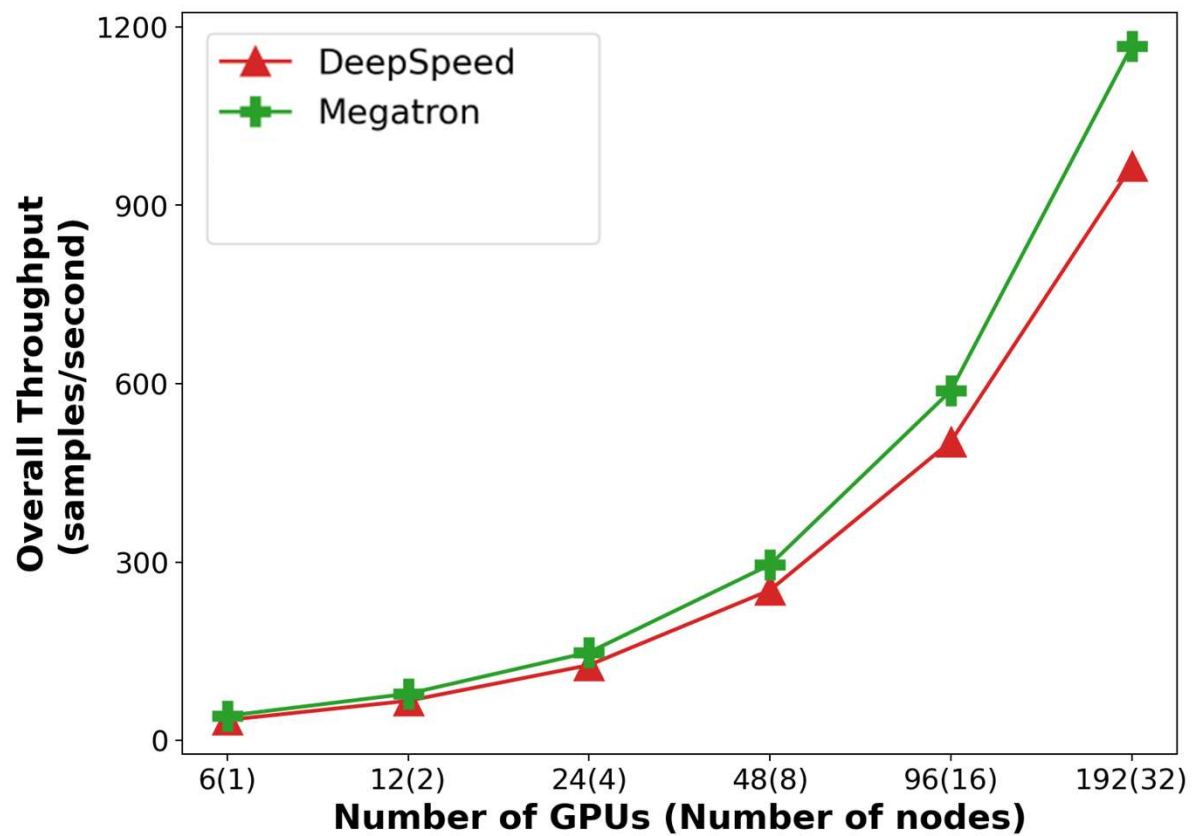
## Baselines



## BERT-Large

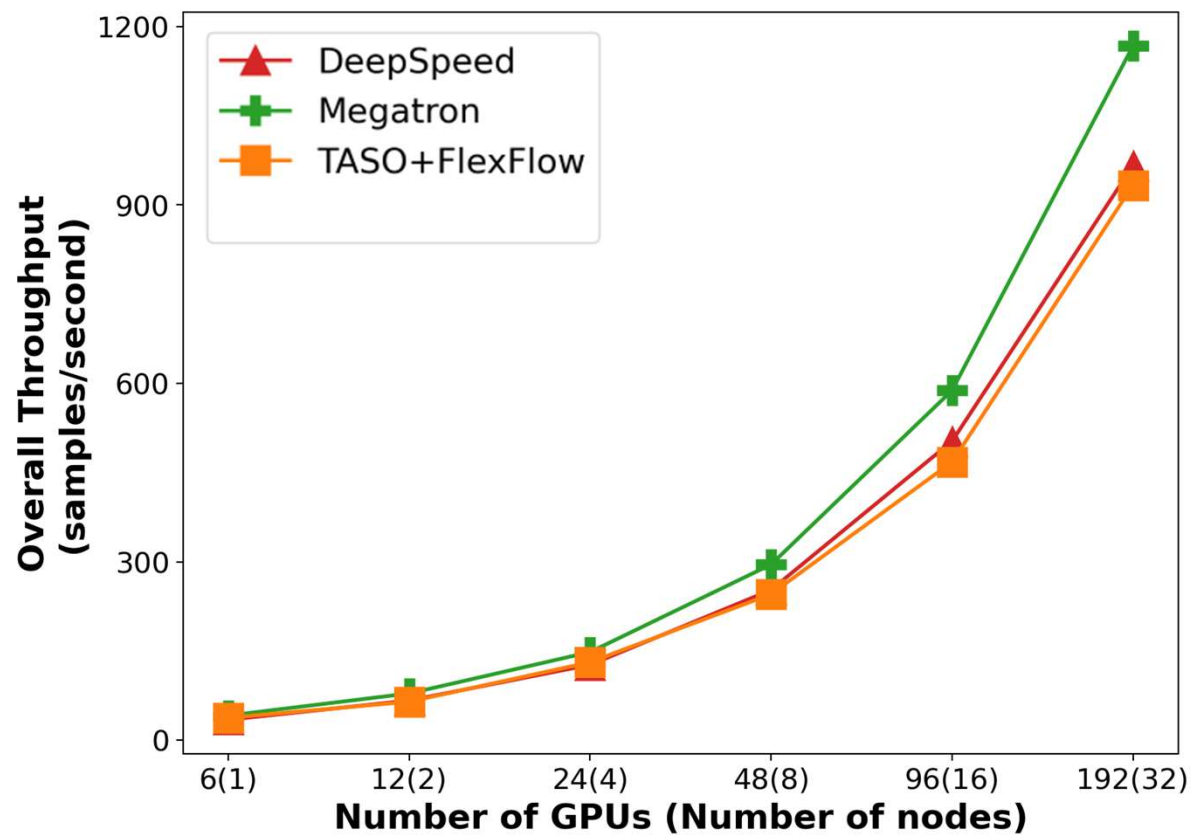


# BERT-Large

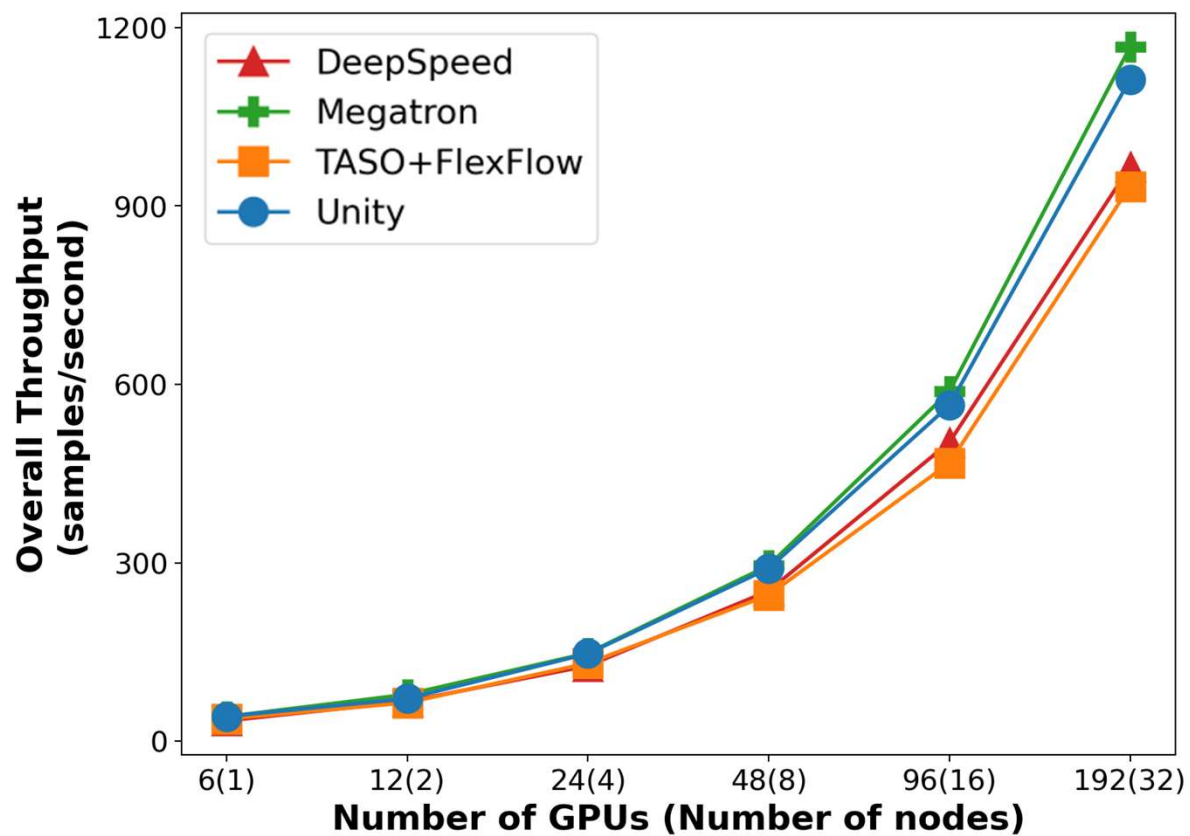




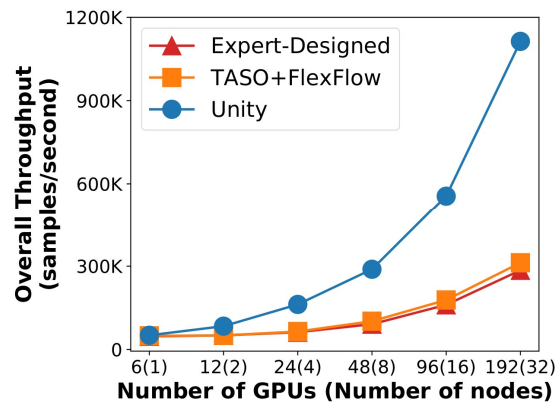
# BERT-Large



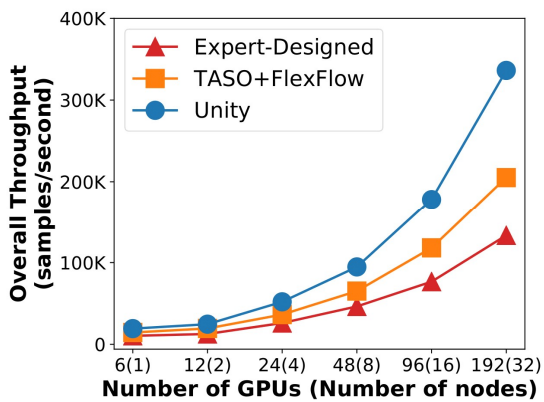
# BERT-Large



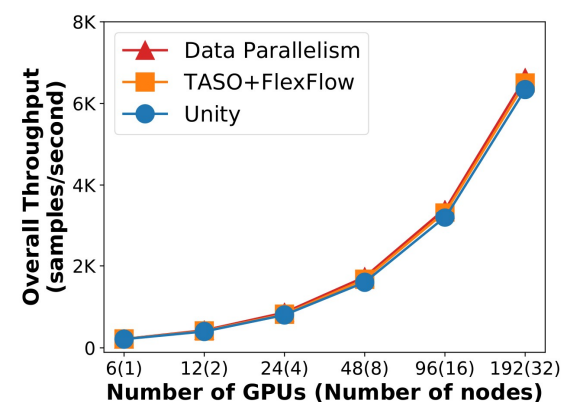
### DLRM



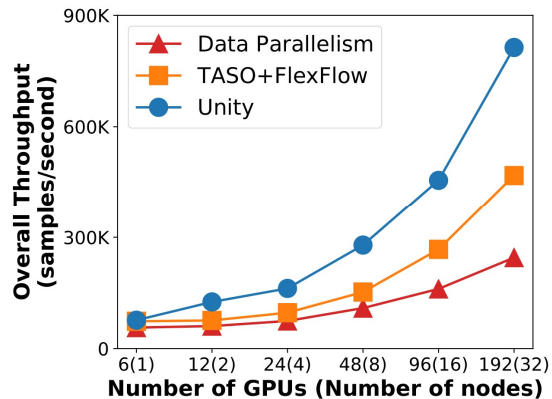
### CANDLE-Uno



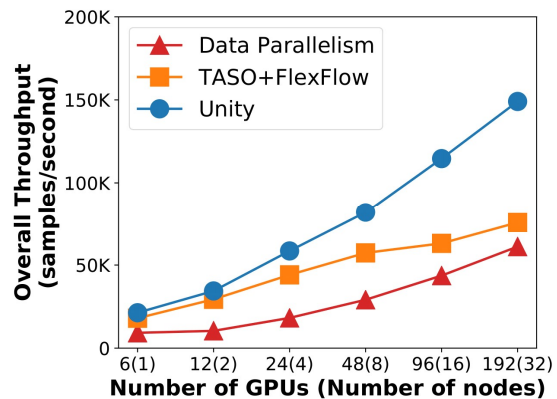
### ResNeXt-50



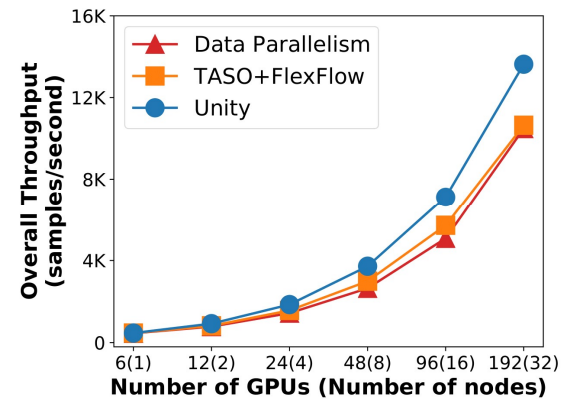
### XDL



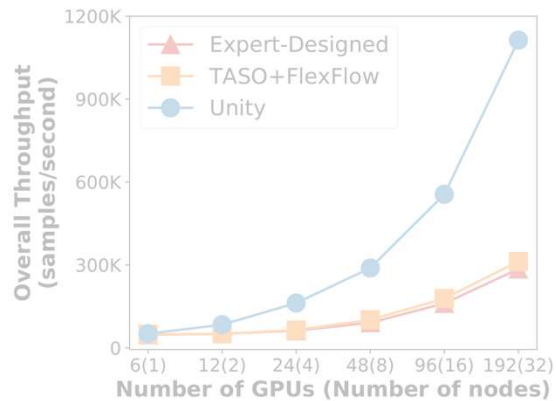
### MLP



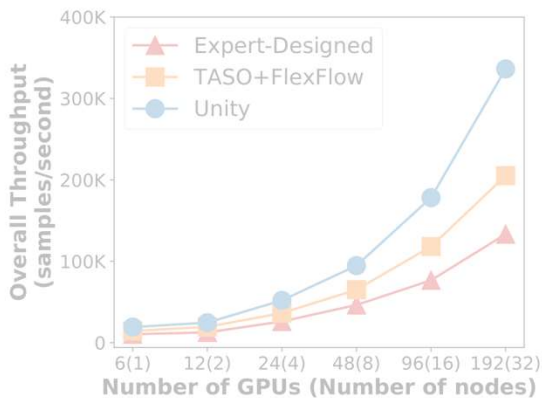
### Inception-v3



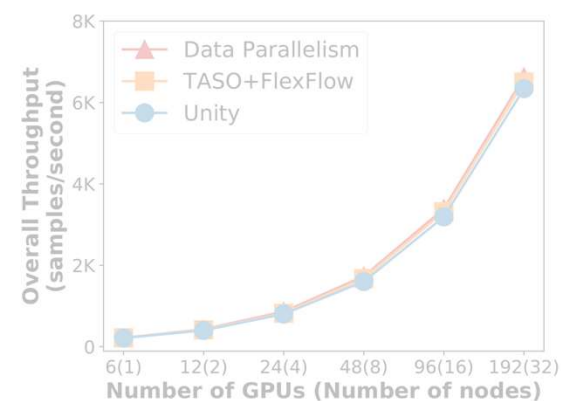
DLRM



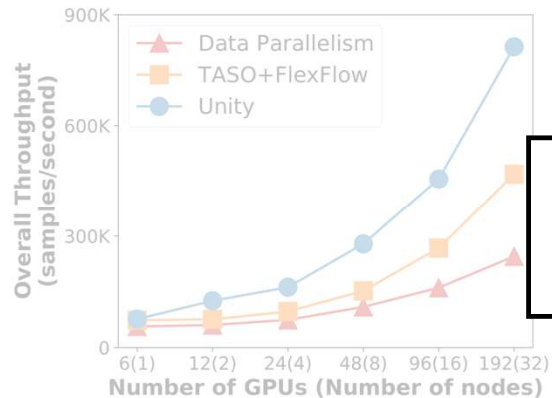
CANDLE-Uno



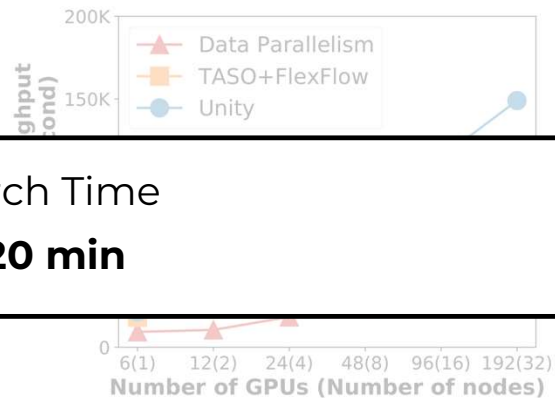
ResNeXt-50



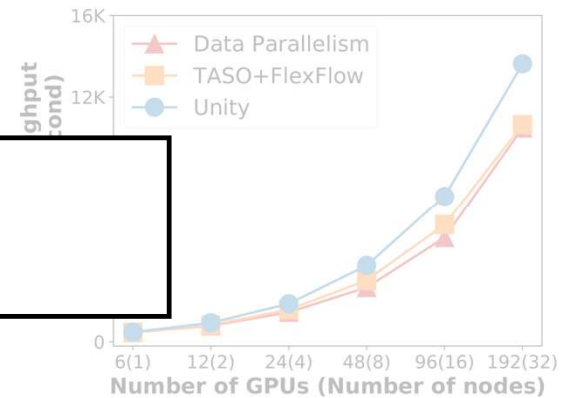
XDL



MLP

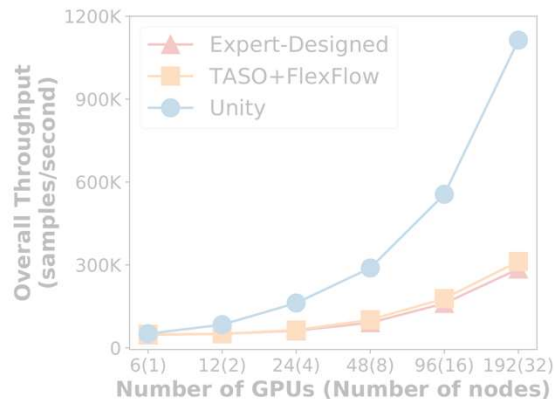


Inception-v3

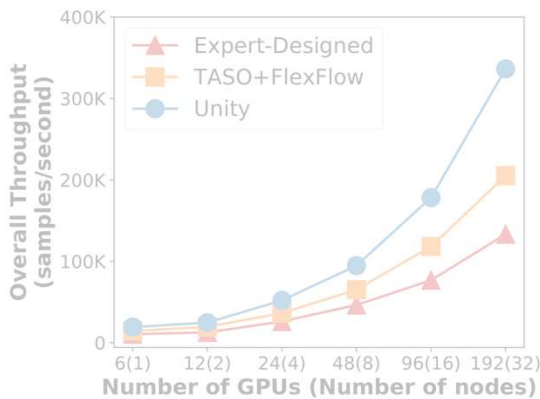


Search Time  
< 20 min

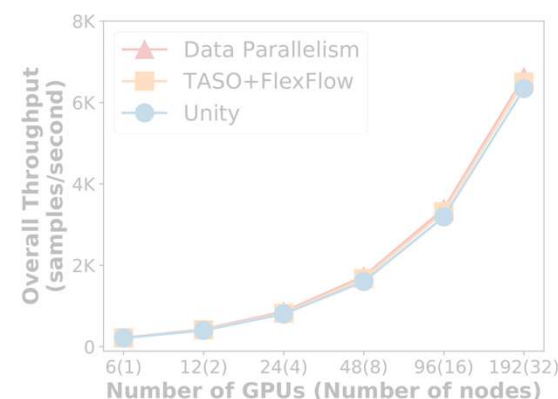
### DLRM



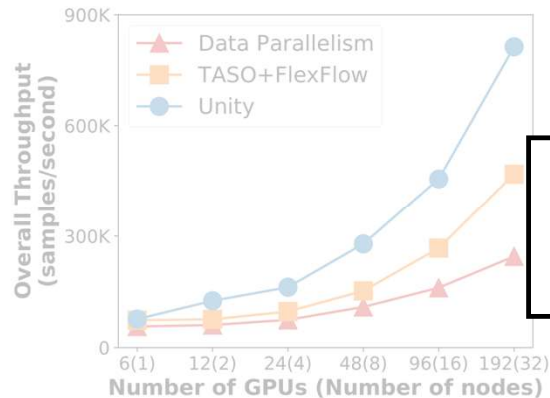
### CANDLE-Uono



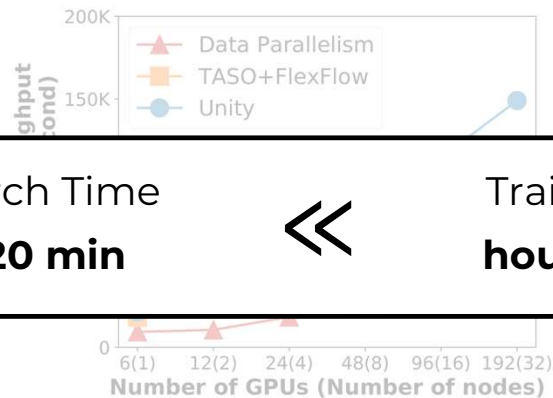
### ResNeXt-50



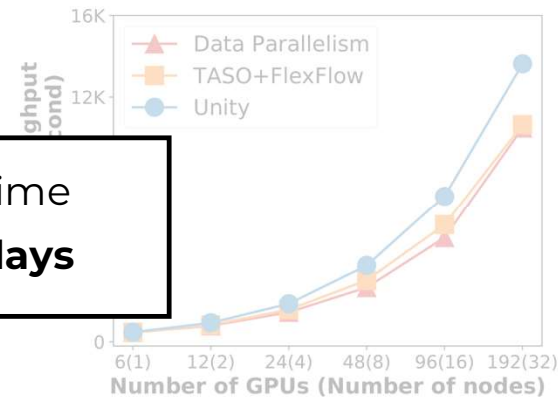
### XDL



### MLP



### Inception-v3



Search Time  
**< 20 min**



Training Time  
**hours or days**

<https://github.com/flexflow/FlexFlow>

kadinzhang Added tests for Linear operator in align/linear (#264)		4884234 · 19 days ago · 1,339 commits
circleci	[CircleCI] update script	12 months ago
align	Added tests for Linear operator in align/linear (#264)	19 days ago
bootcamp_demo	[python] make num_samples a property (#121)	14 months ago
cmake	[CMake] retrieve legion installation path for spack	9 months ago
conda	[Conda] fix conda	17 months ago
config	[Alignment] Add embedding alignment (#255)	3 months ago
deps	[Legion] version update	8 months ago
docker	Adding Dockerfile	2 years ago

A distributed deep learning framework that supports flexible parallelization strategies.

Readme  
Apache-2.0 license  
Code of conduct  
458 stars  
21 watching  
100 forks

Releases  
Release 21.09 (September 30th ...)



Keras



PyTorch



ONNX

python	[Python] fix embedding for pybind11	27 days ago	Publish your first package
scripts	Merging the public branch to the master branch (#74)	18 months ago	



CMU™

