

# 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

# Unity

Accelerating DNN Training Through **Joint Optimization**  
of Algebraic Transformations and Parallelization

1.

2.

# Unity

Accelerating DNN Training Through Joint Optimization  
of Algebraic Transformations and Parallelization

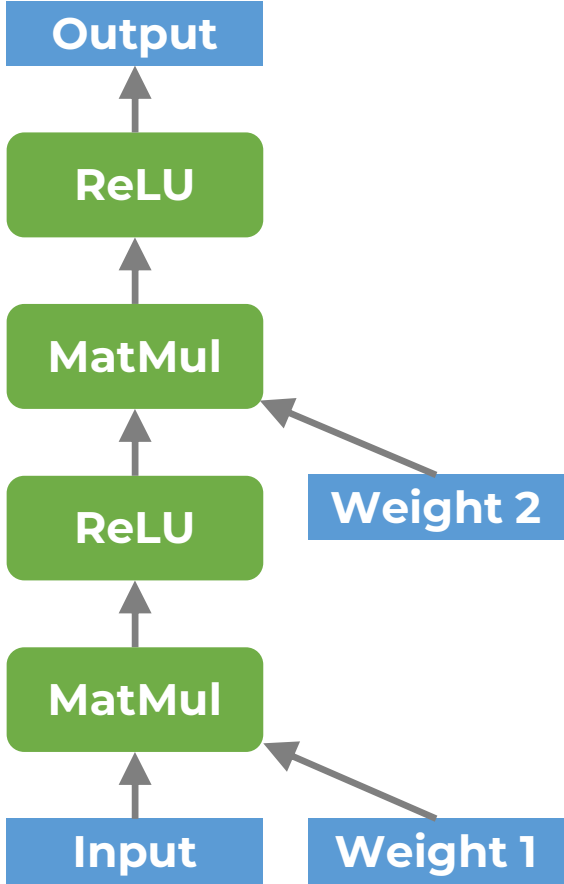
1. Algebraic Transformations
- 2.

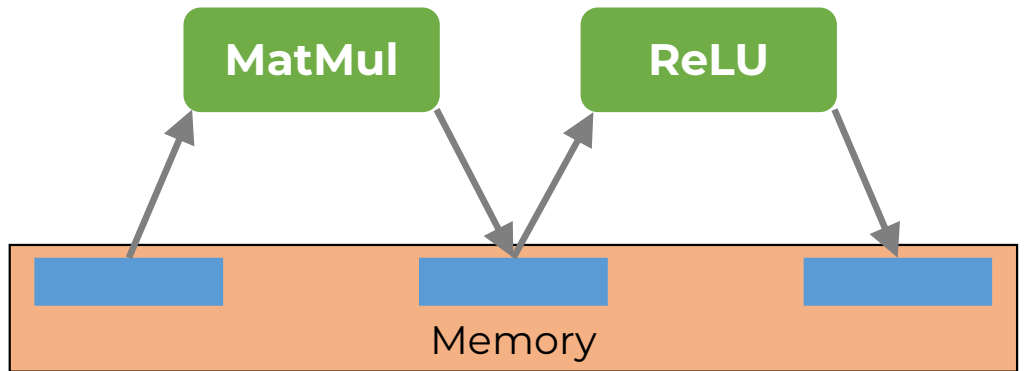
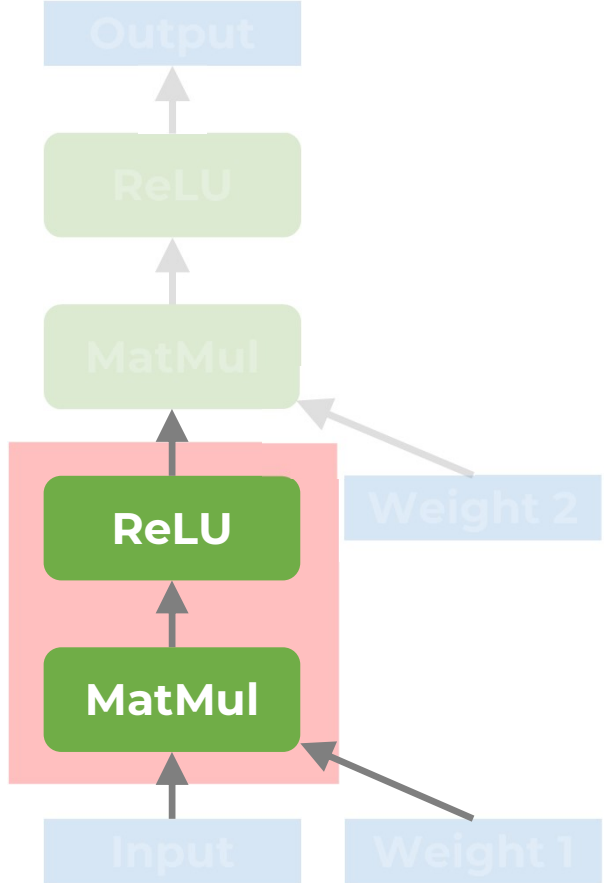
# Unity

Accelerating DNN Training Through Joint Optimization  
of Algebraic Transformations and [Parallelization](#)

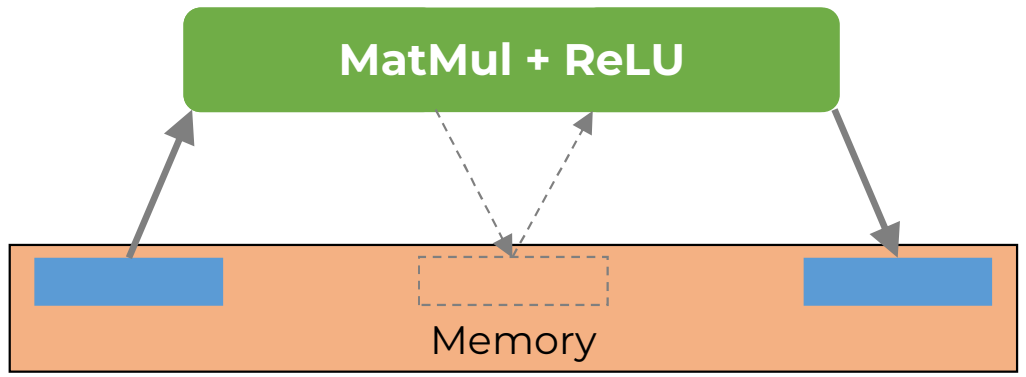
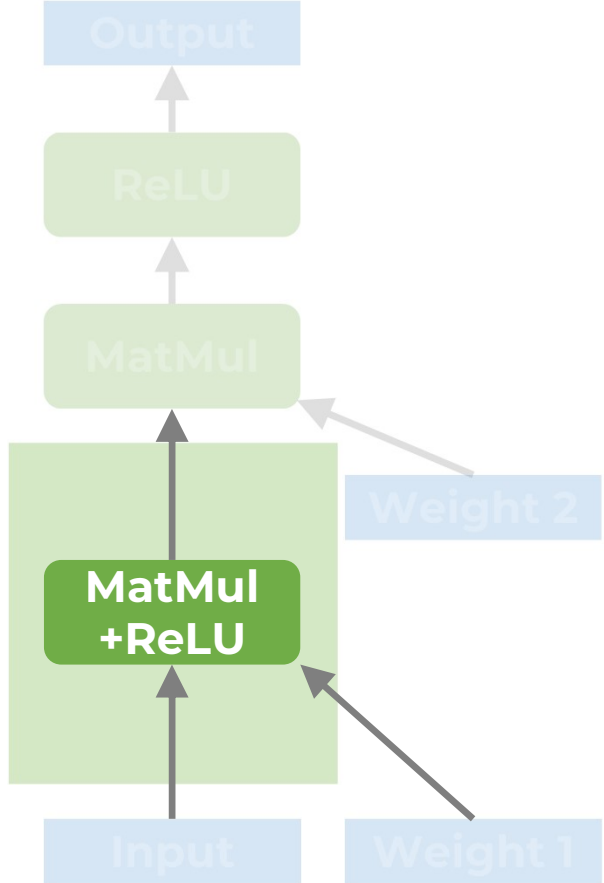
1. Algebraic Transformations
2. Parallelization

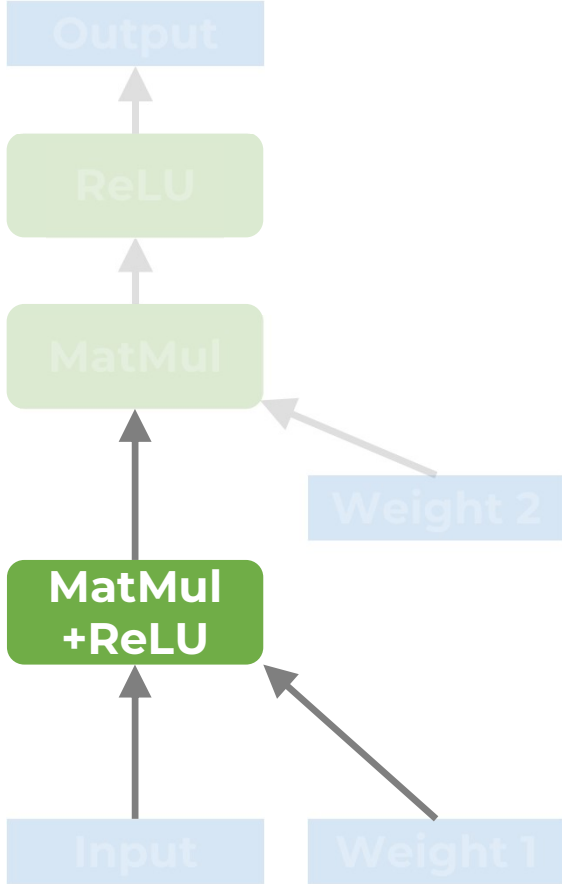
# 1. Algebraic Transformations



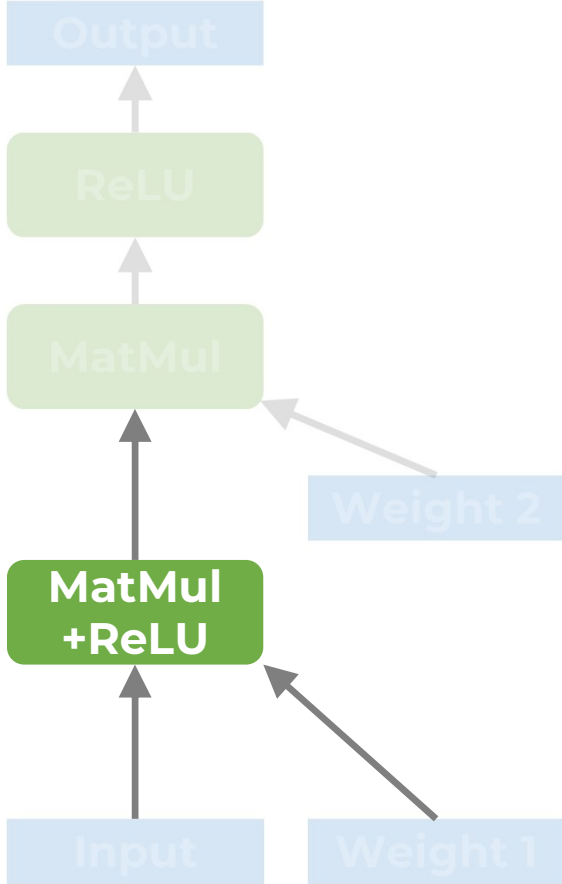




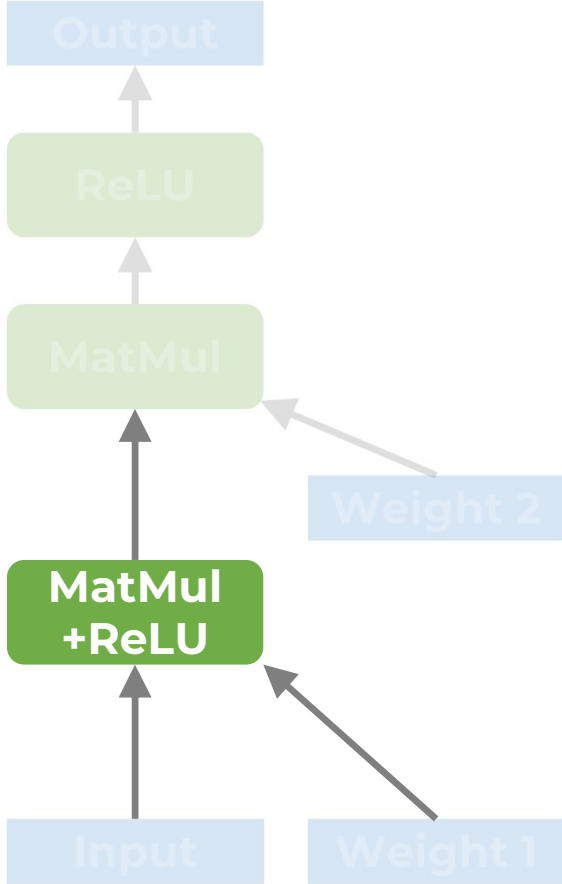




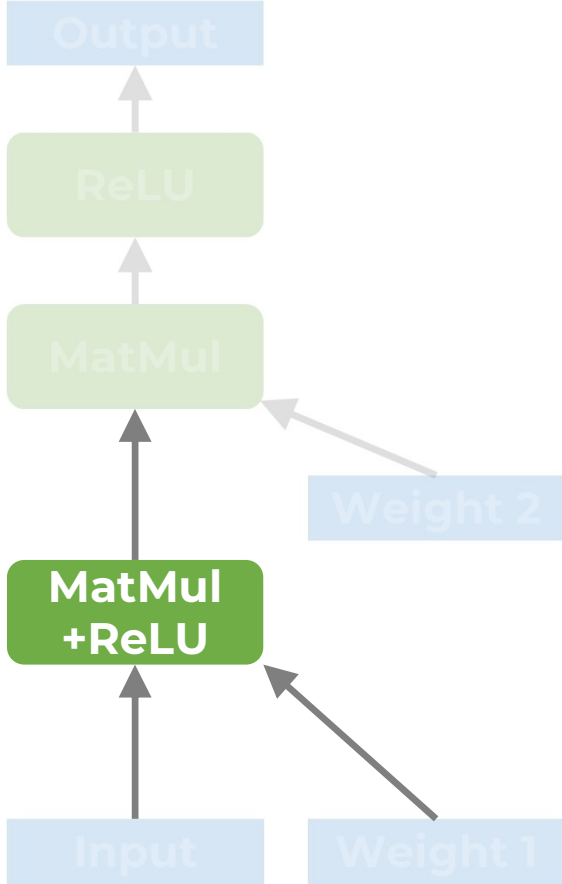
## Operator Fusion



Operator Fusion  
Operator Splitting



Operator Fusion  
Operator Splitting  
Operator Reordering



Operator Fusion  
Operator Splitting  
Operator Reordering

...

Operator Fusion  
Operator Splitting  
Operator Reordering

...

Operator Fusion  
Operator Splitting  
Operator Reordering  
...

## “Algebraic Transformations”

Operator Fusion  
Operator Splitting  
Operator Reordering  
...

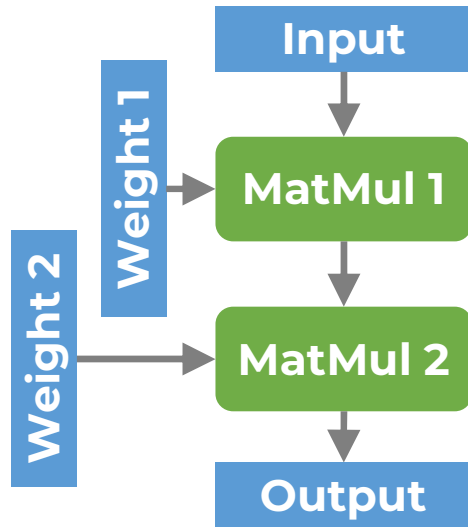


# 1. Algebraic Transformations

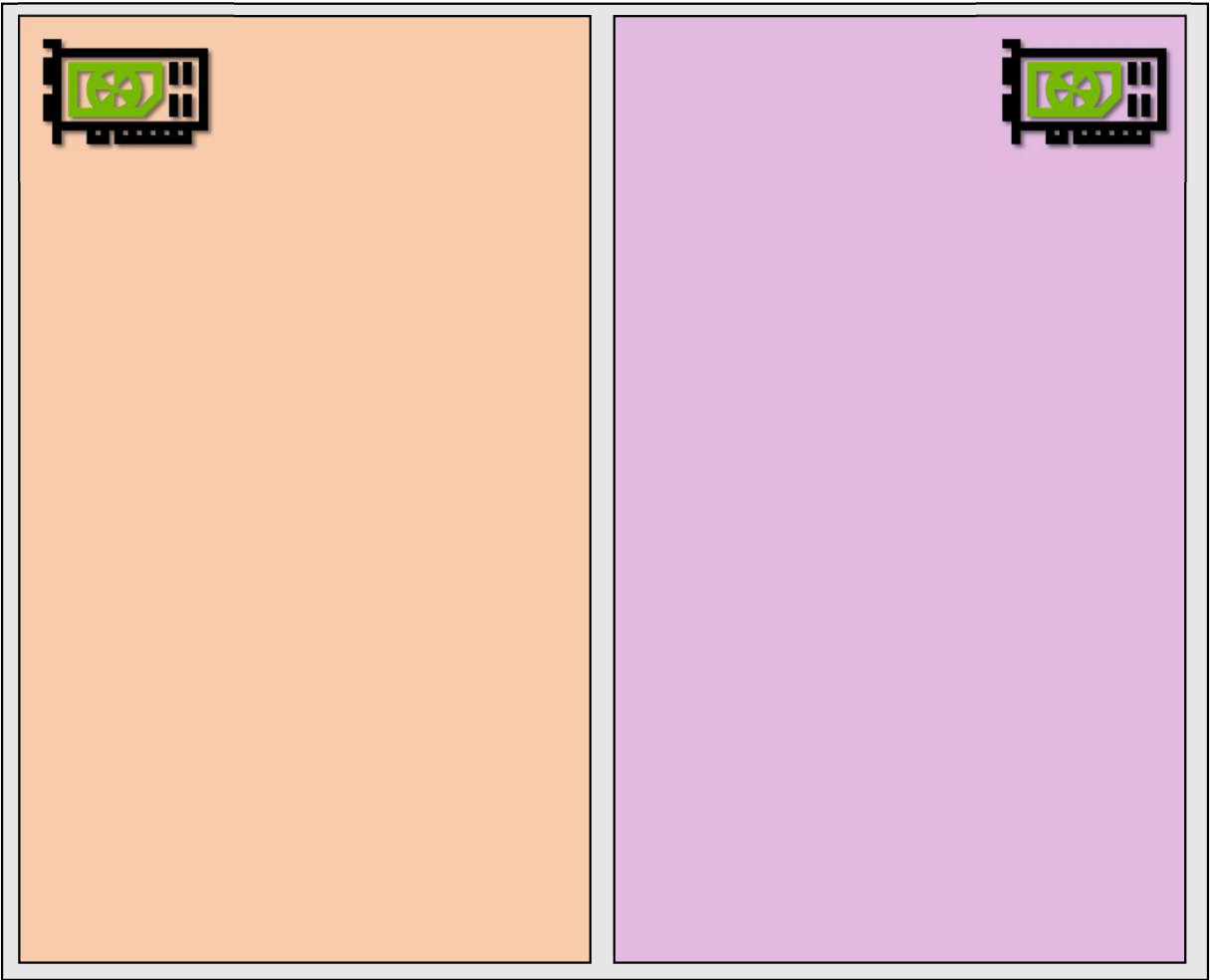
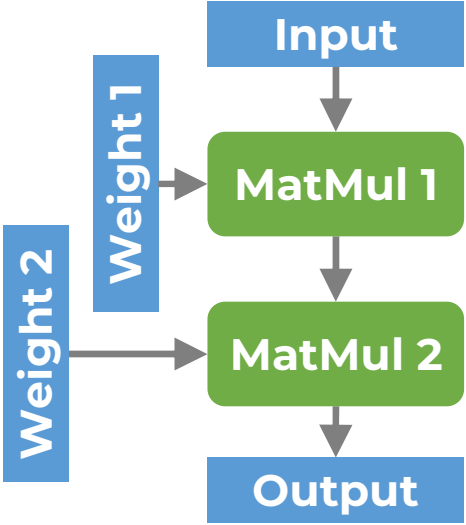
1. Algebraic Transformations
2. Parallelization

# Data Parallel

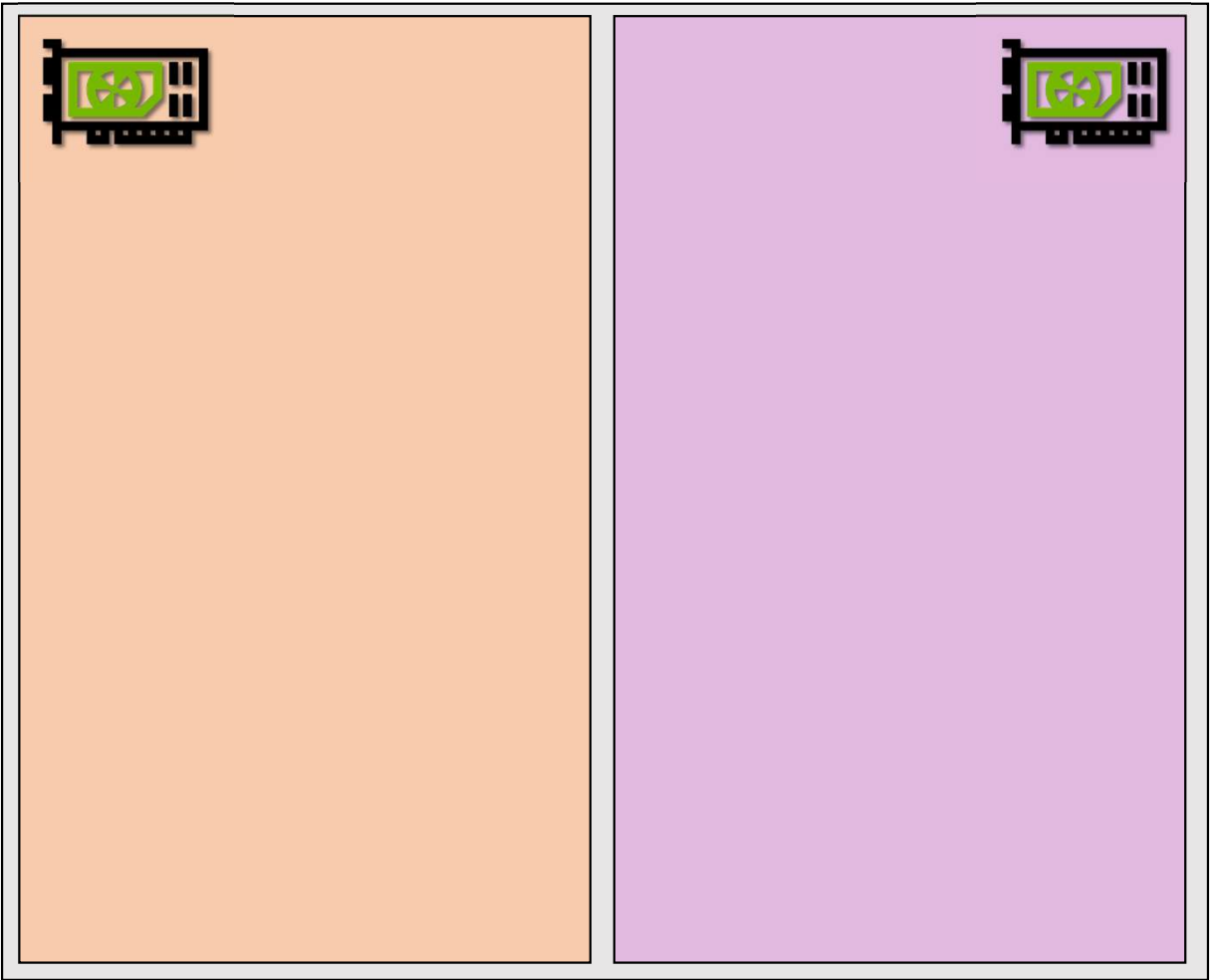
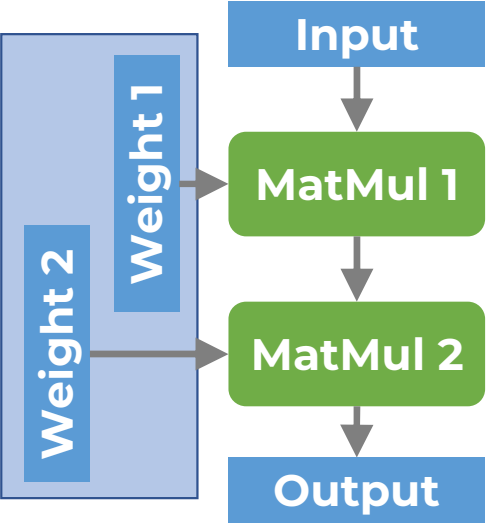
# Data Parallel



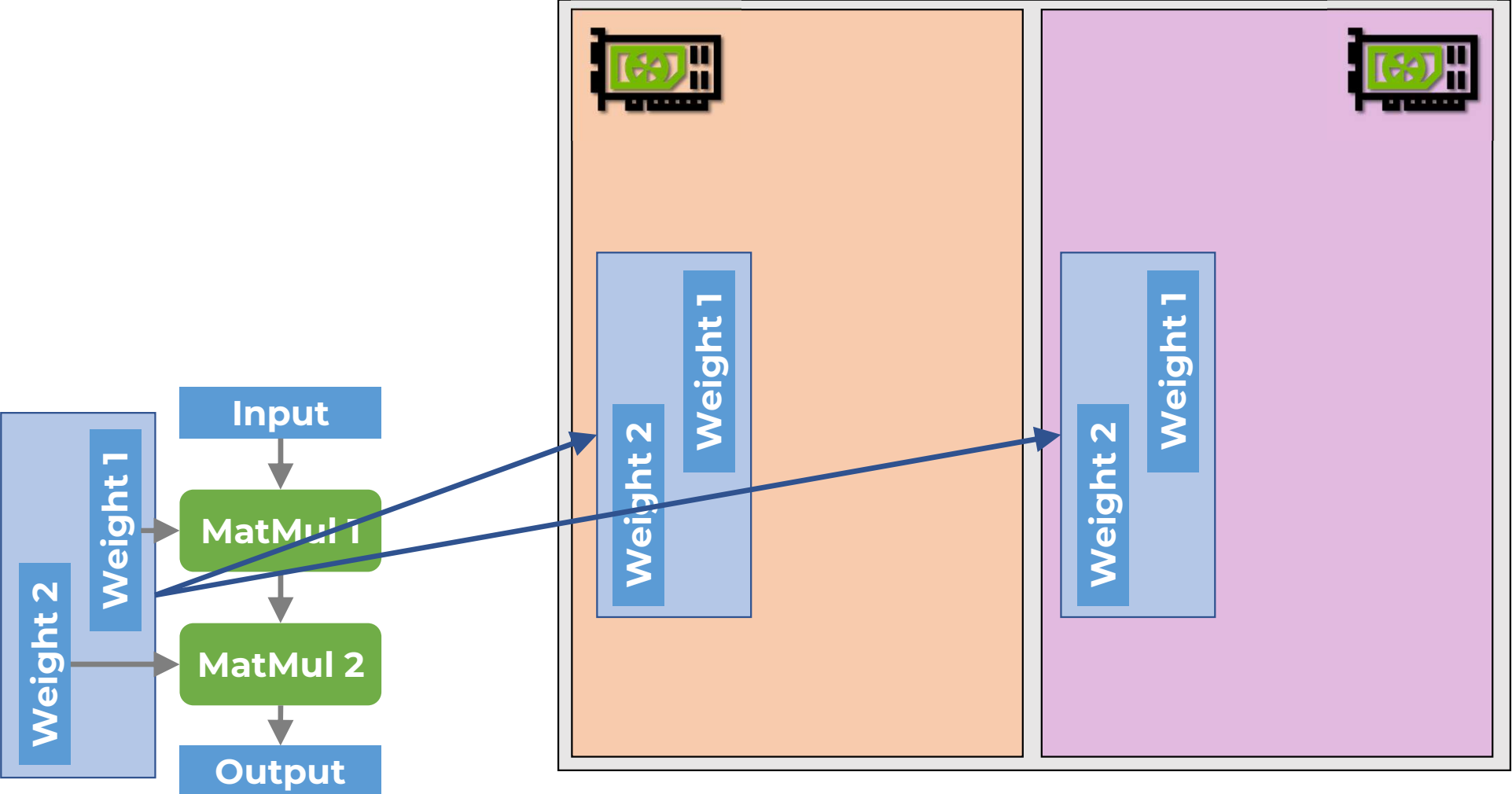
# Data Parallel



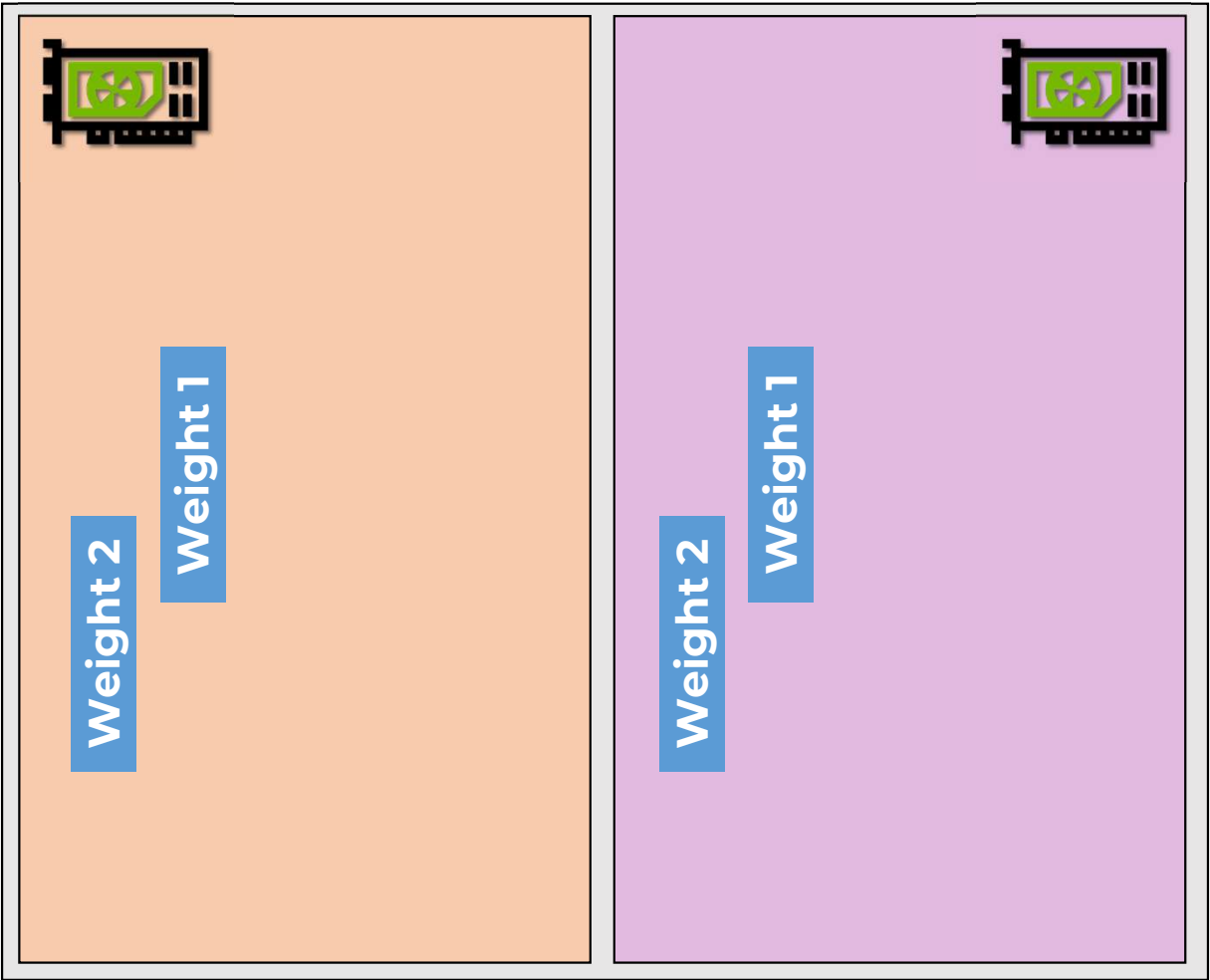
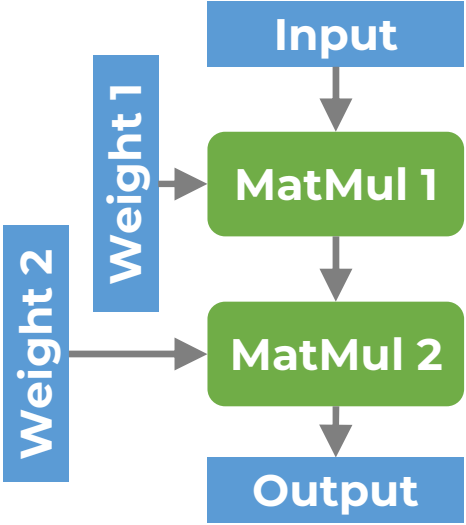
# Data Parallel



# Data Parallel

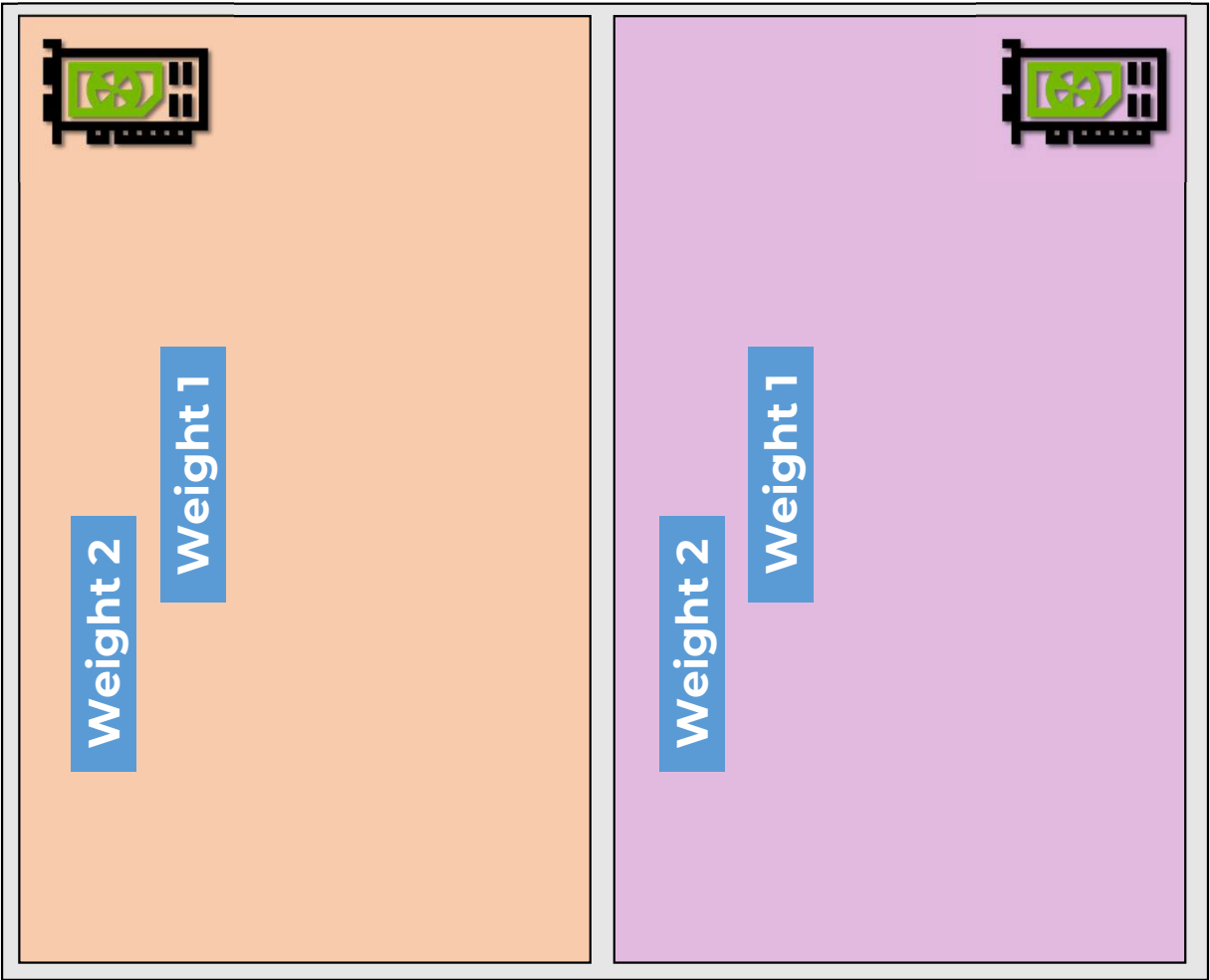
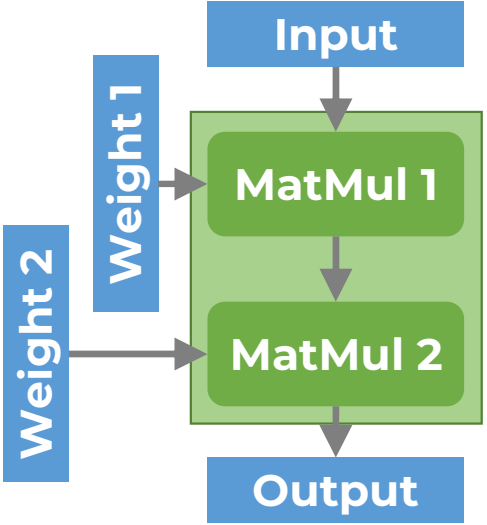


# Data Parallel

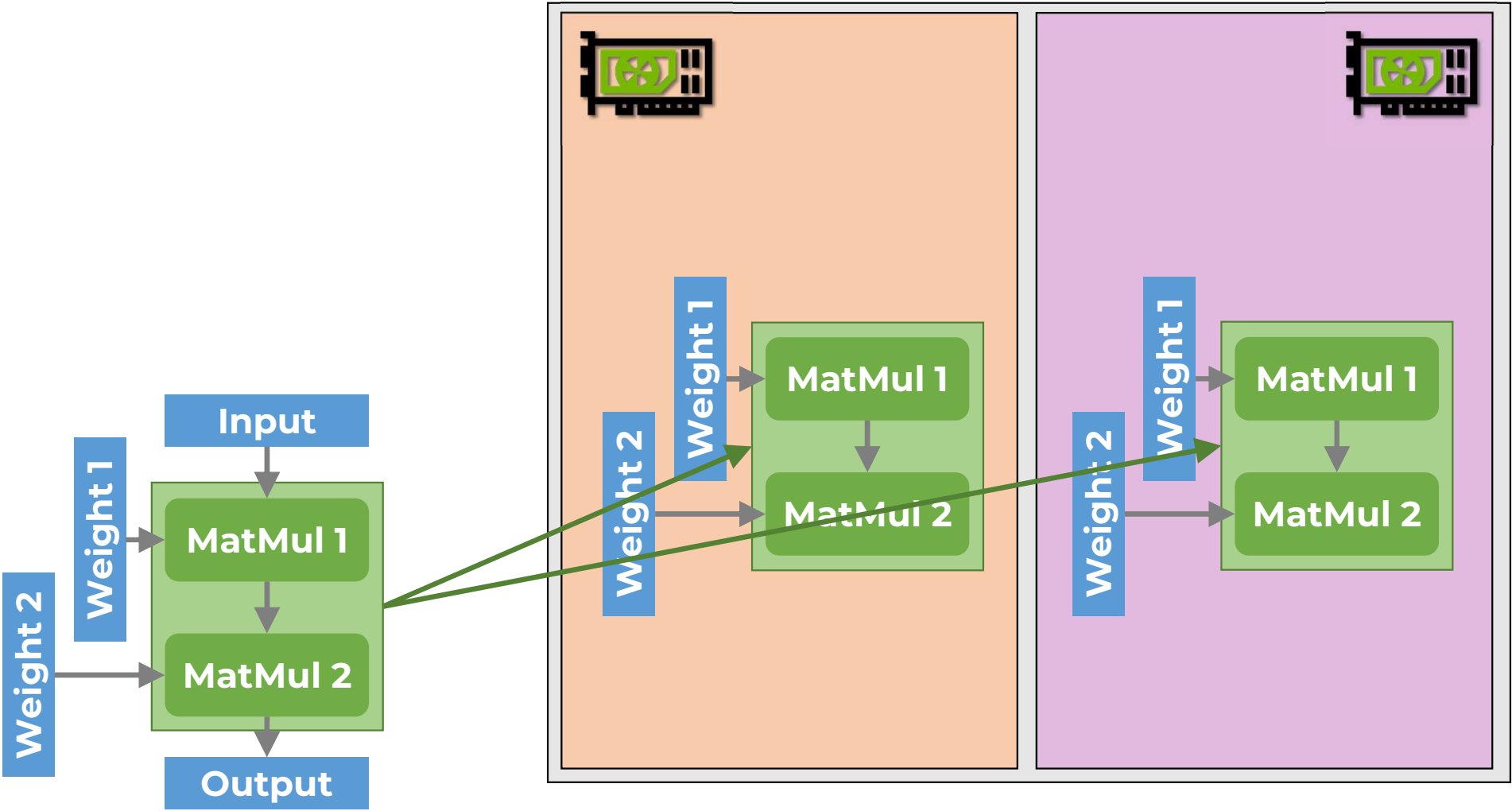




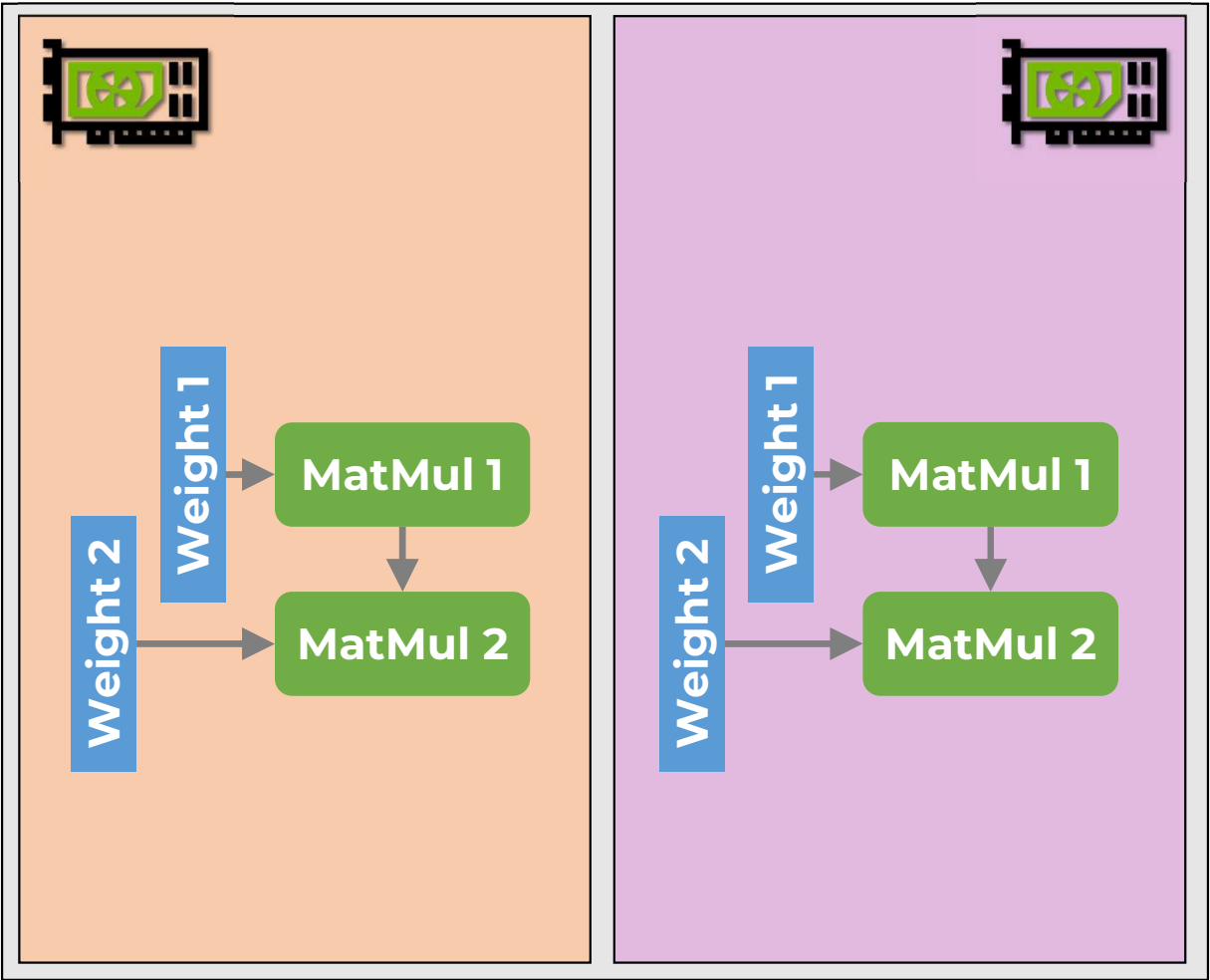
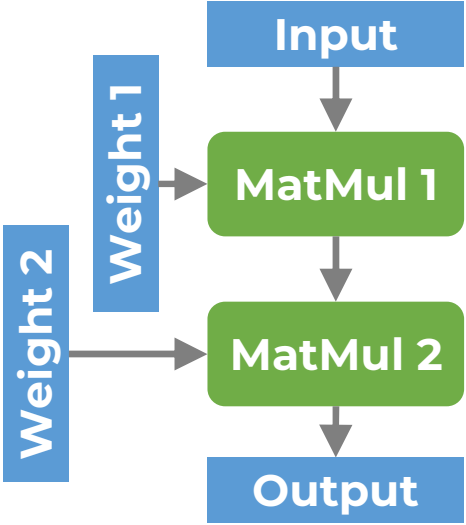
# Data Parallel



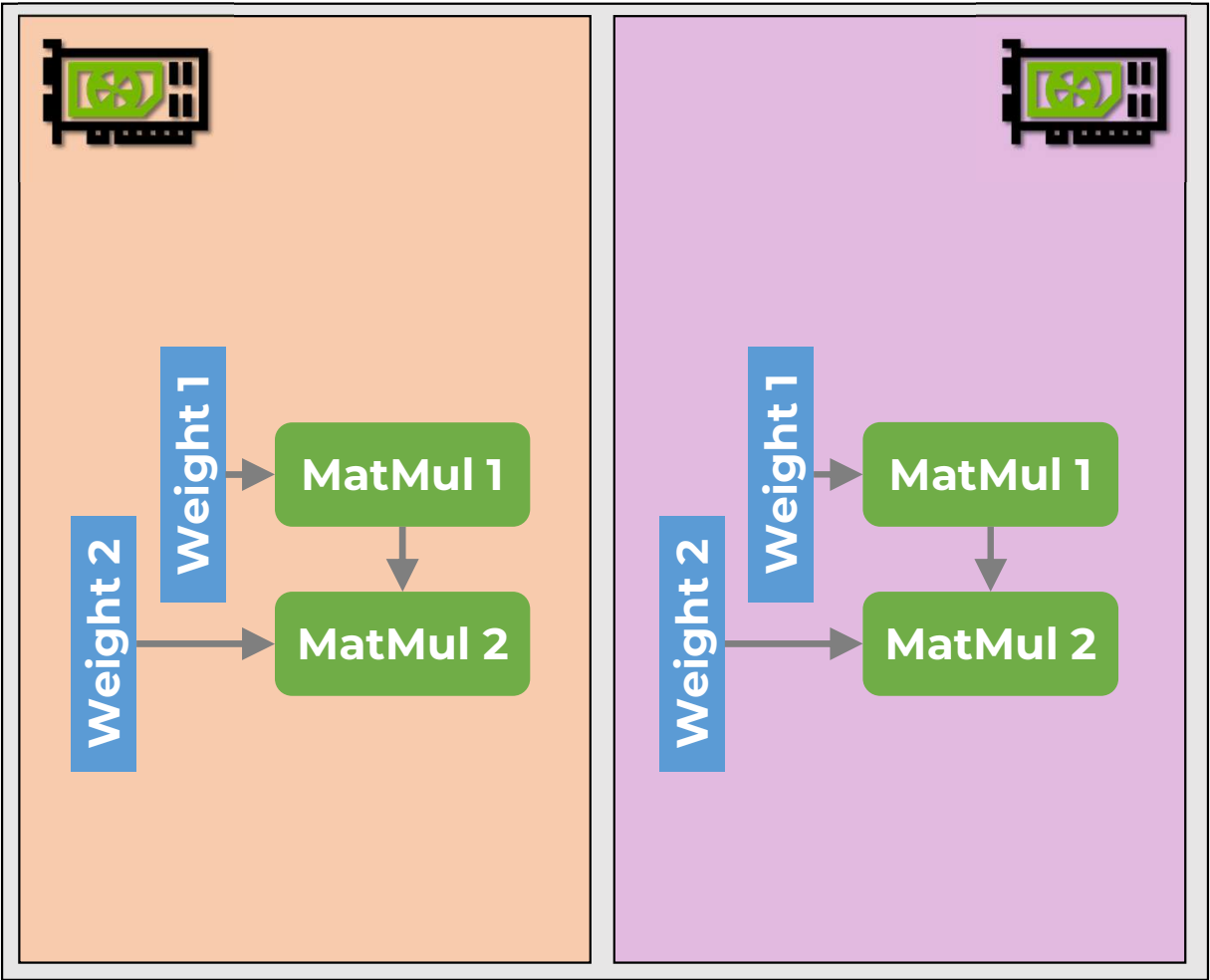
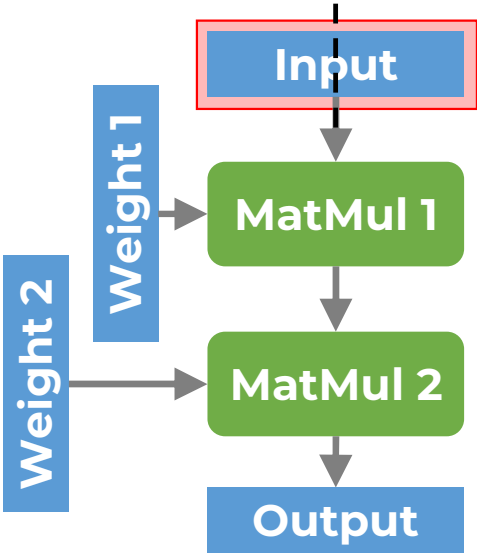
# Data Parallel



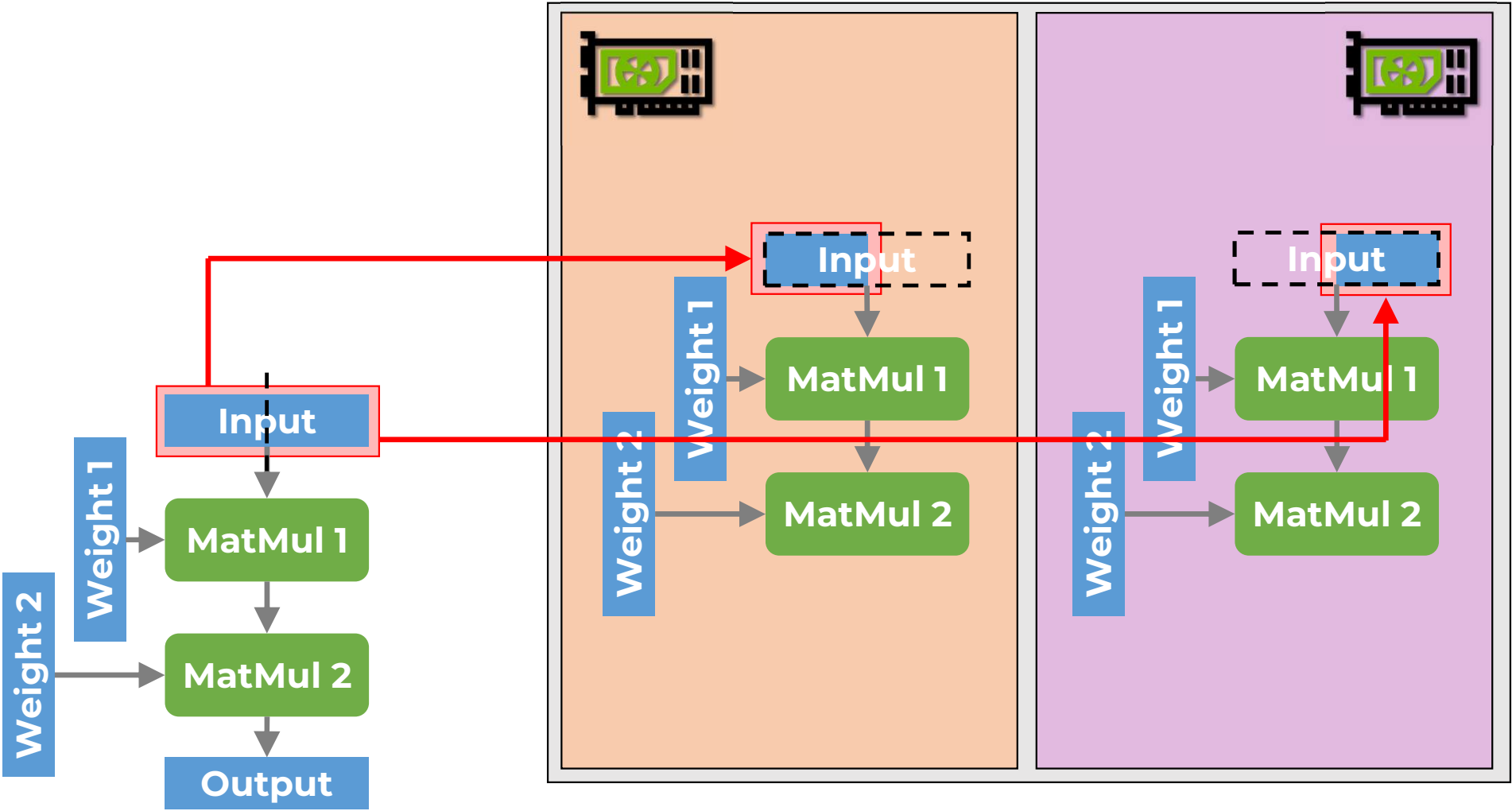
# Data Parallel



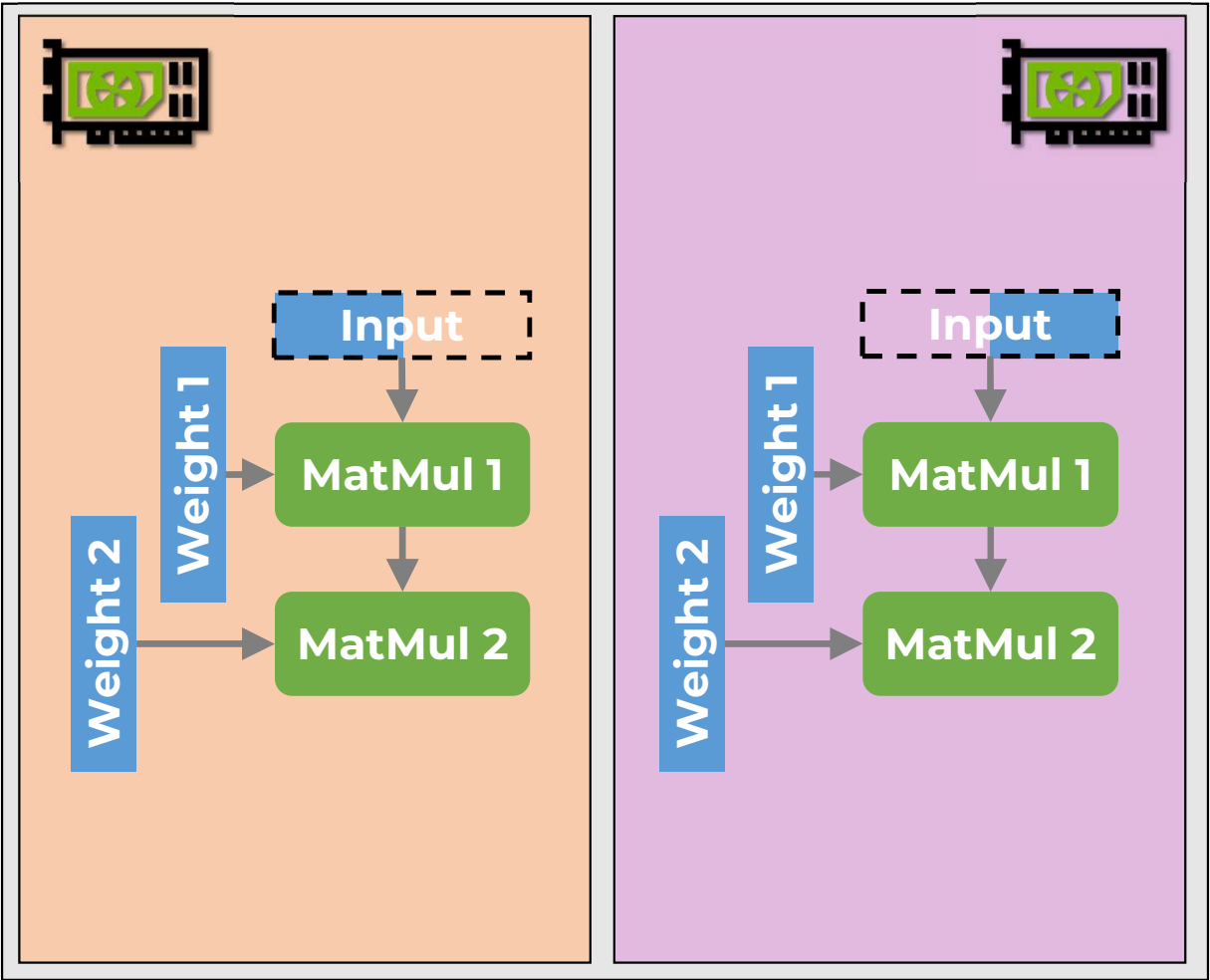
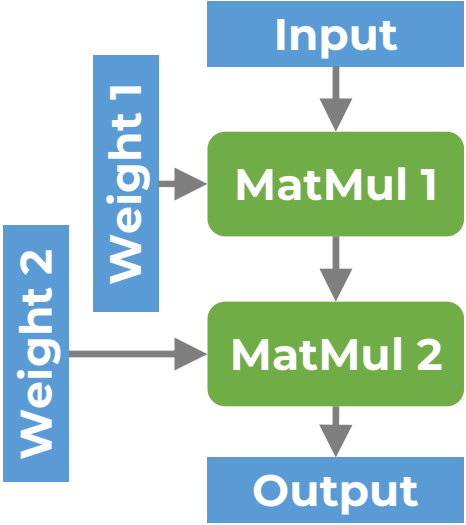
# Data Parallel



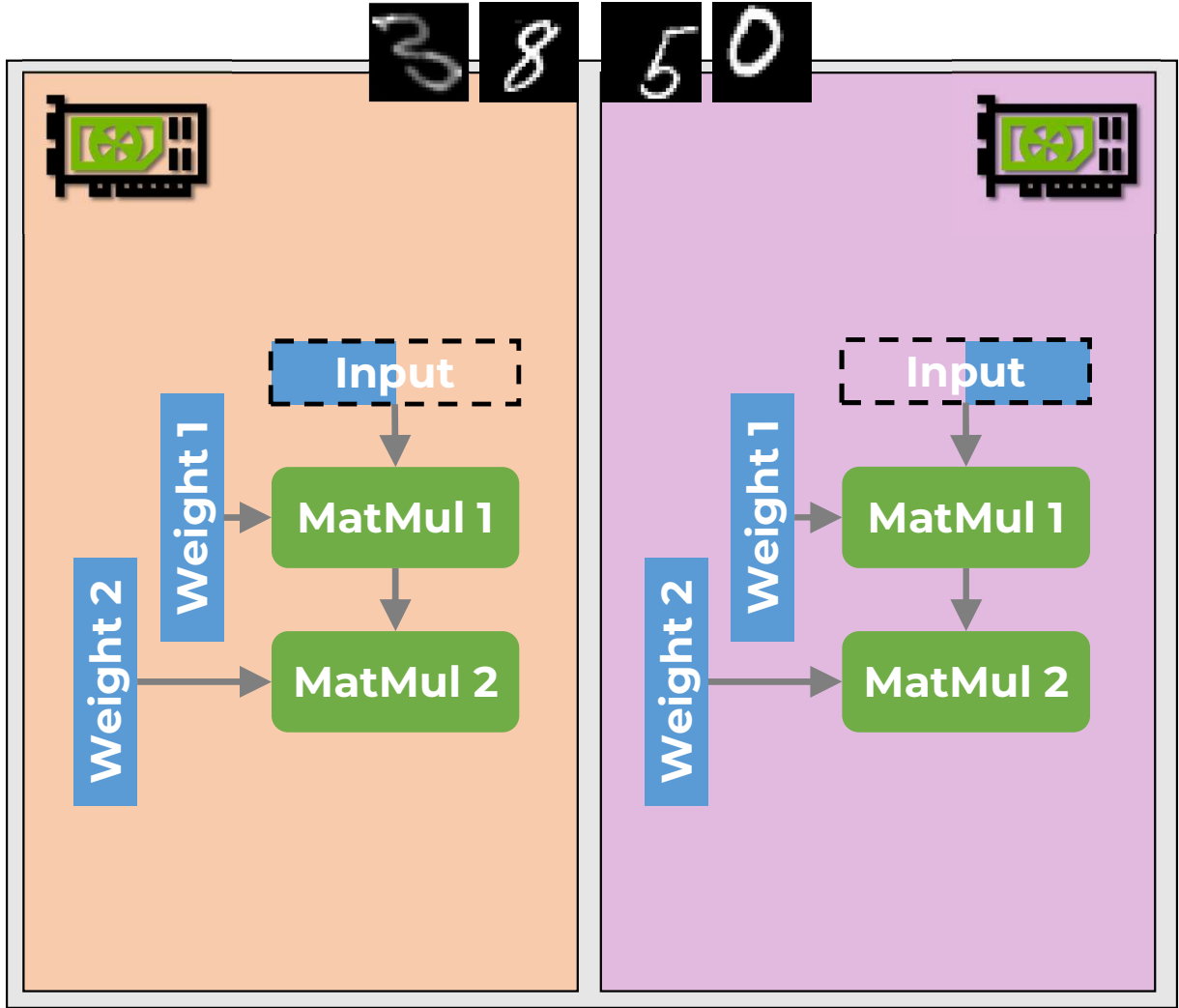
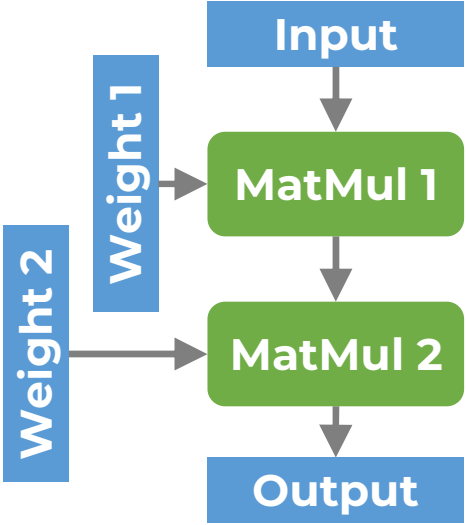
# Data Parallel



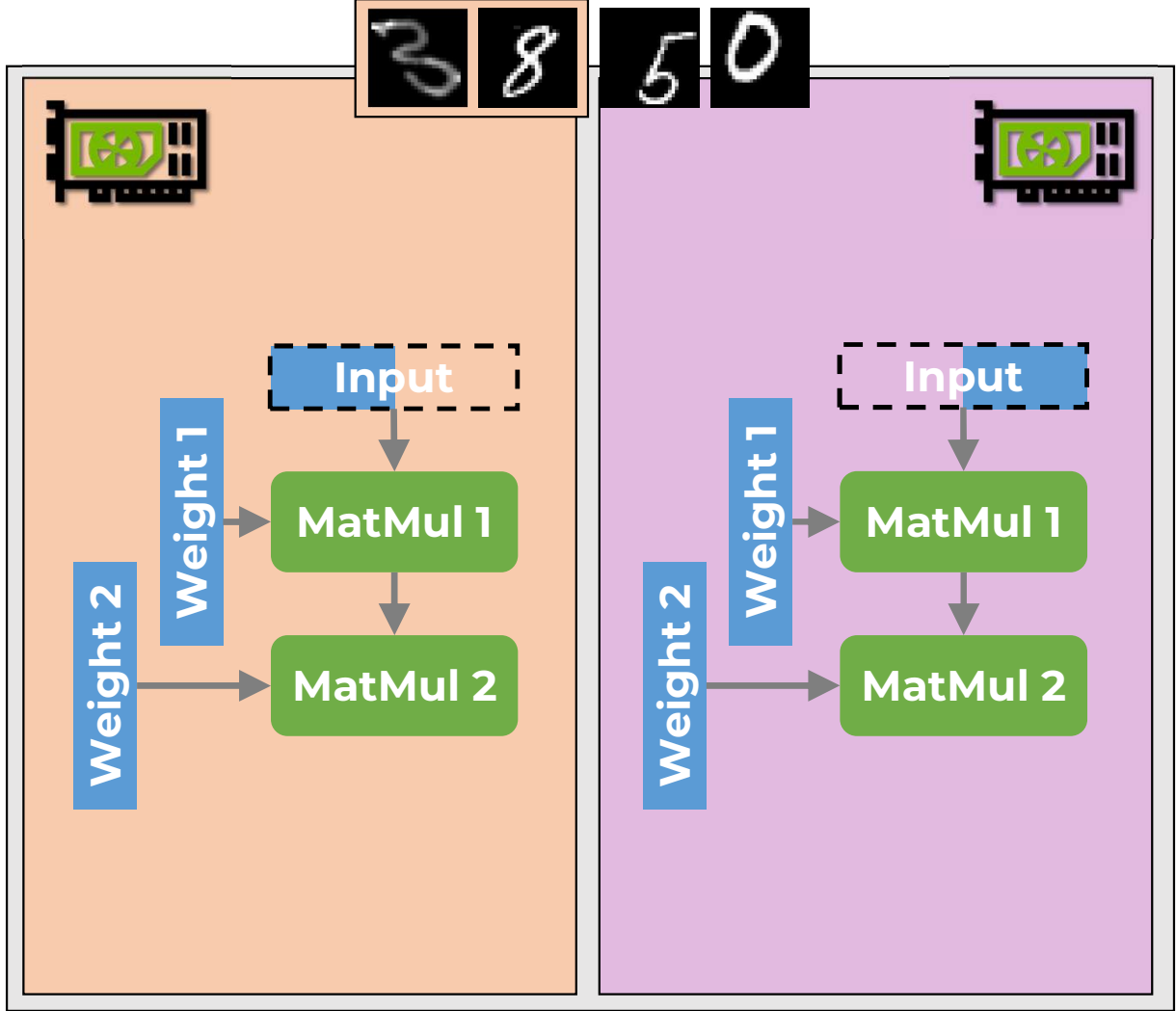
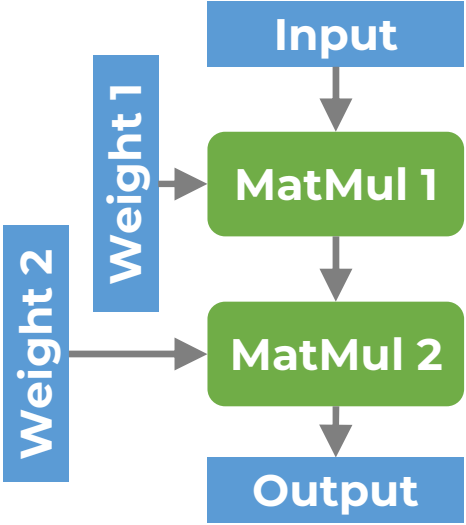
# Data Parallel



# Data Parallel

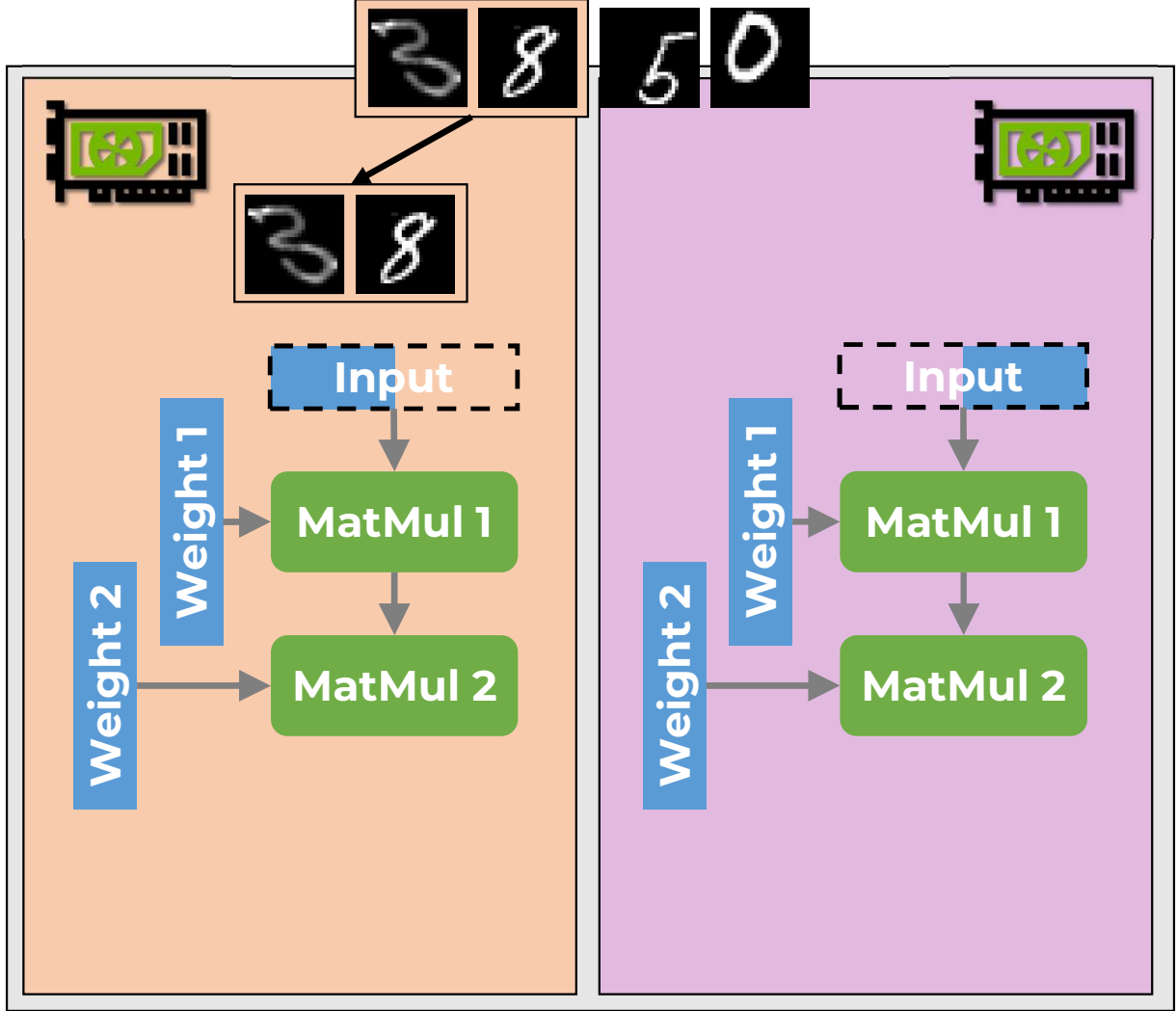
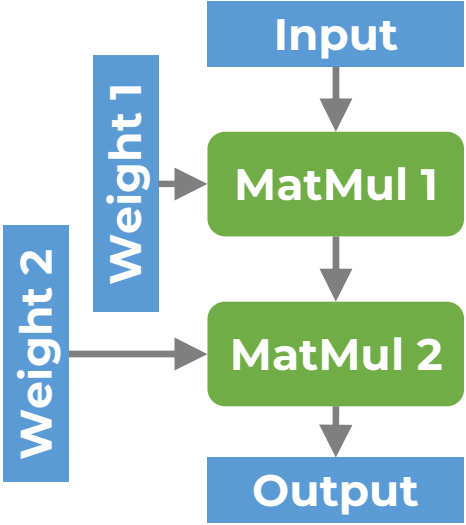


# Data Parallel

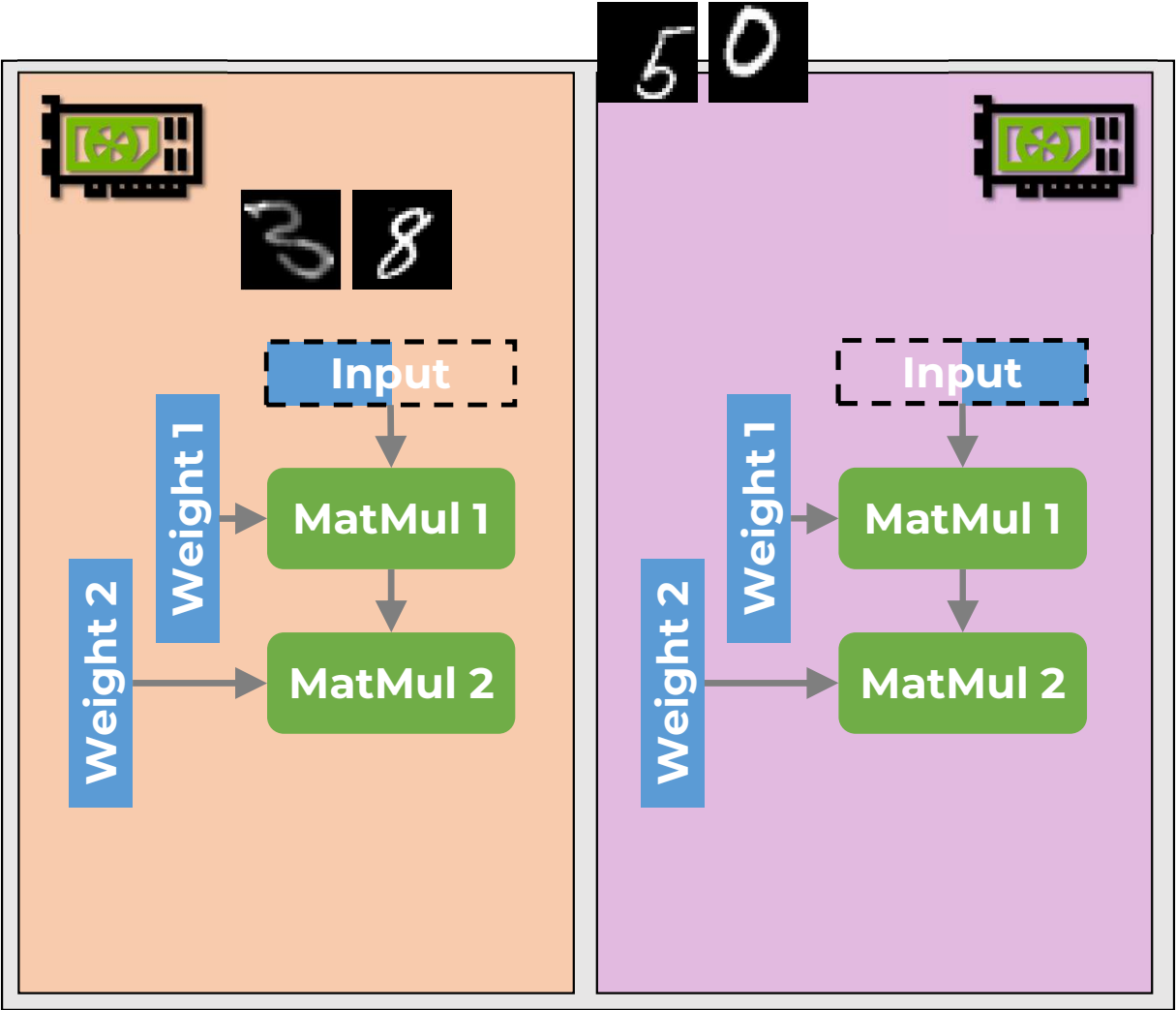
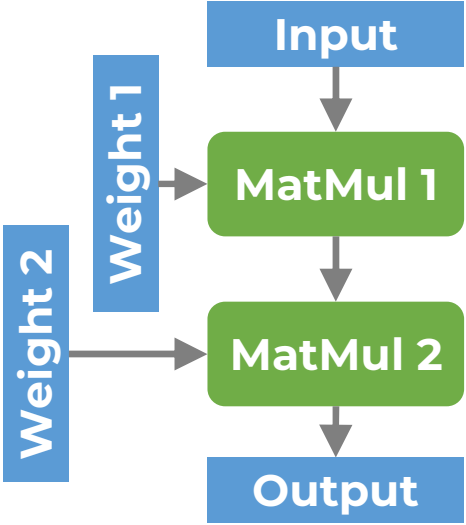




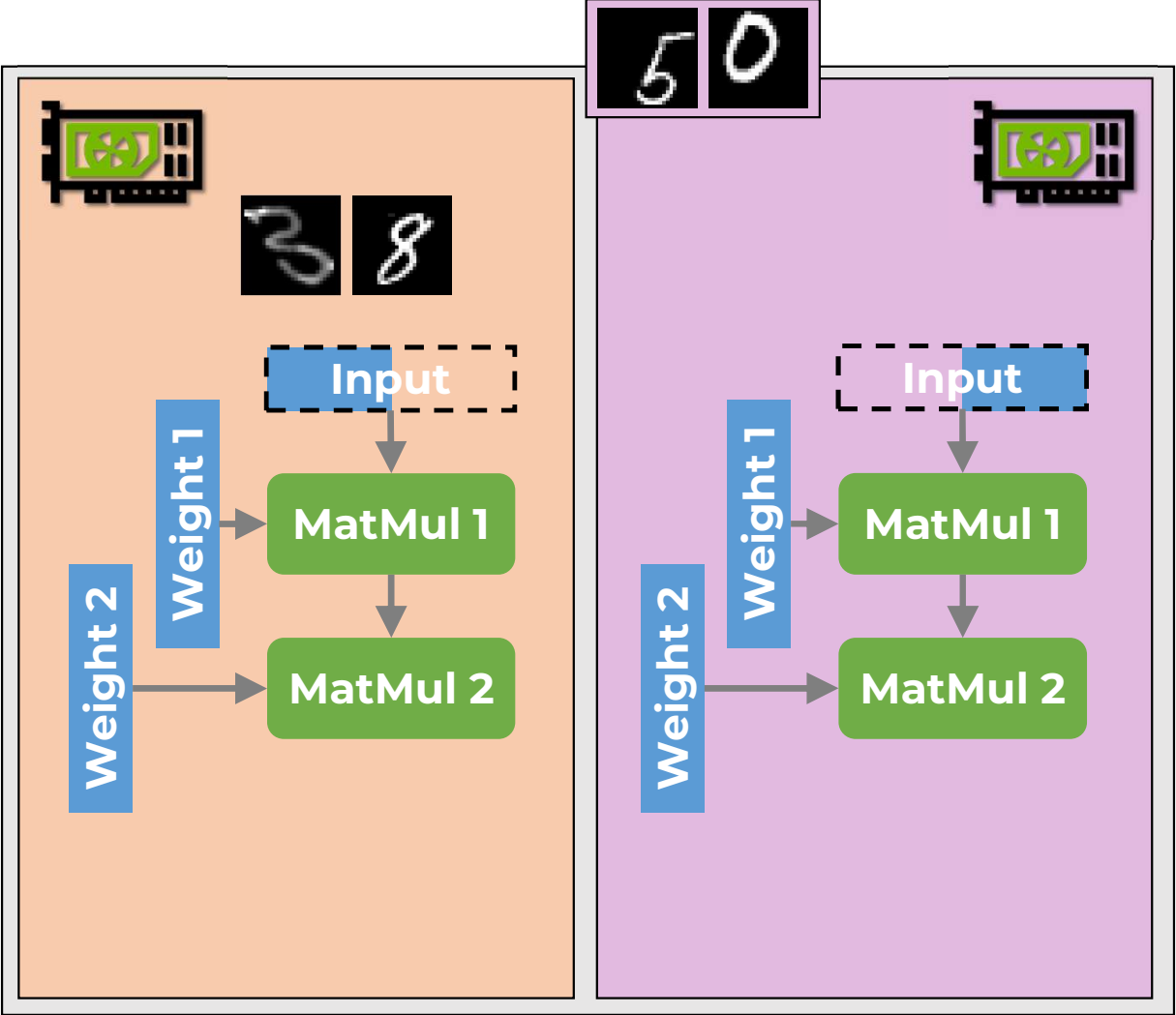
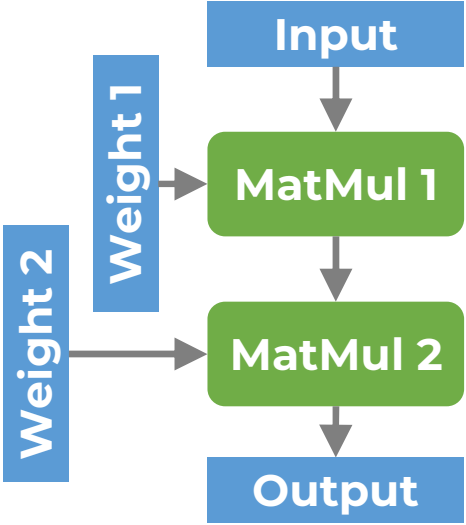
# Data Parallel



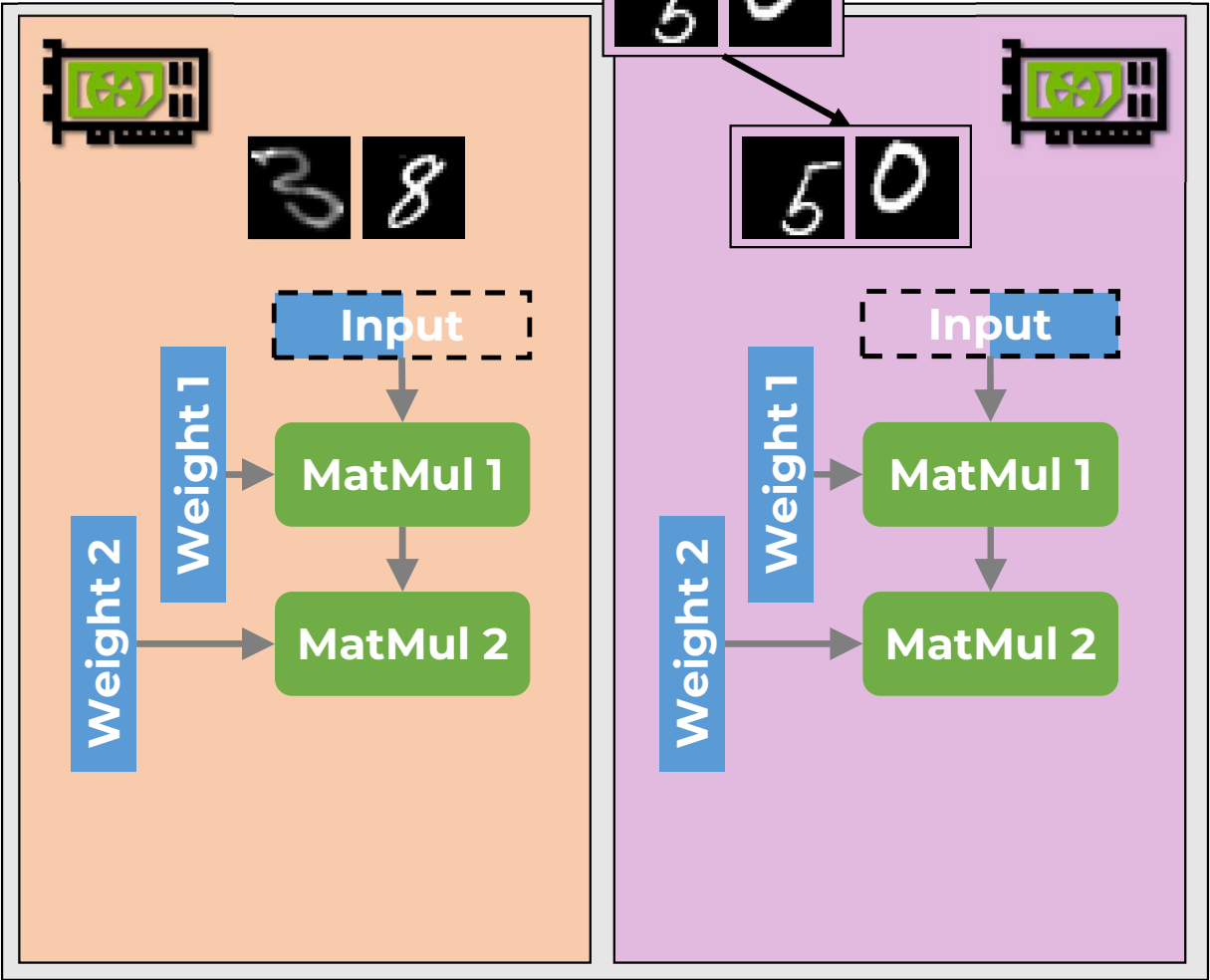
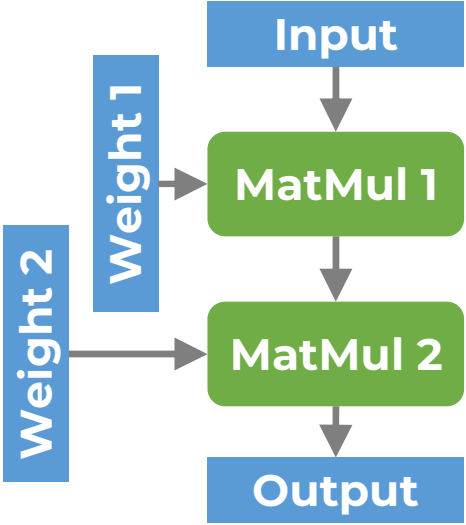
# Data Parallel



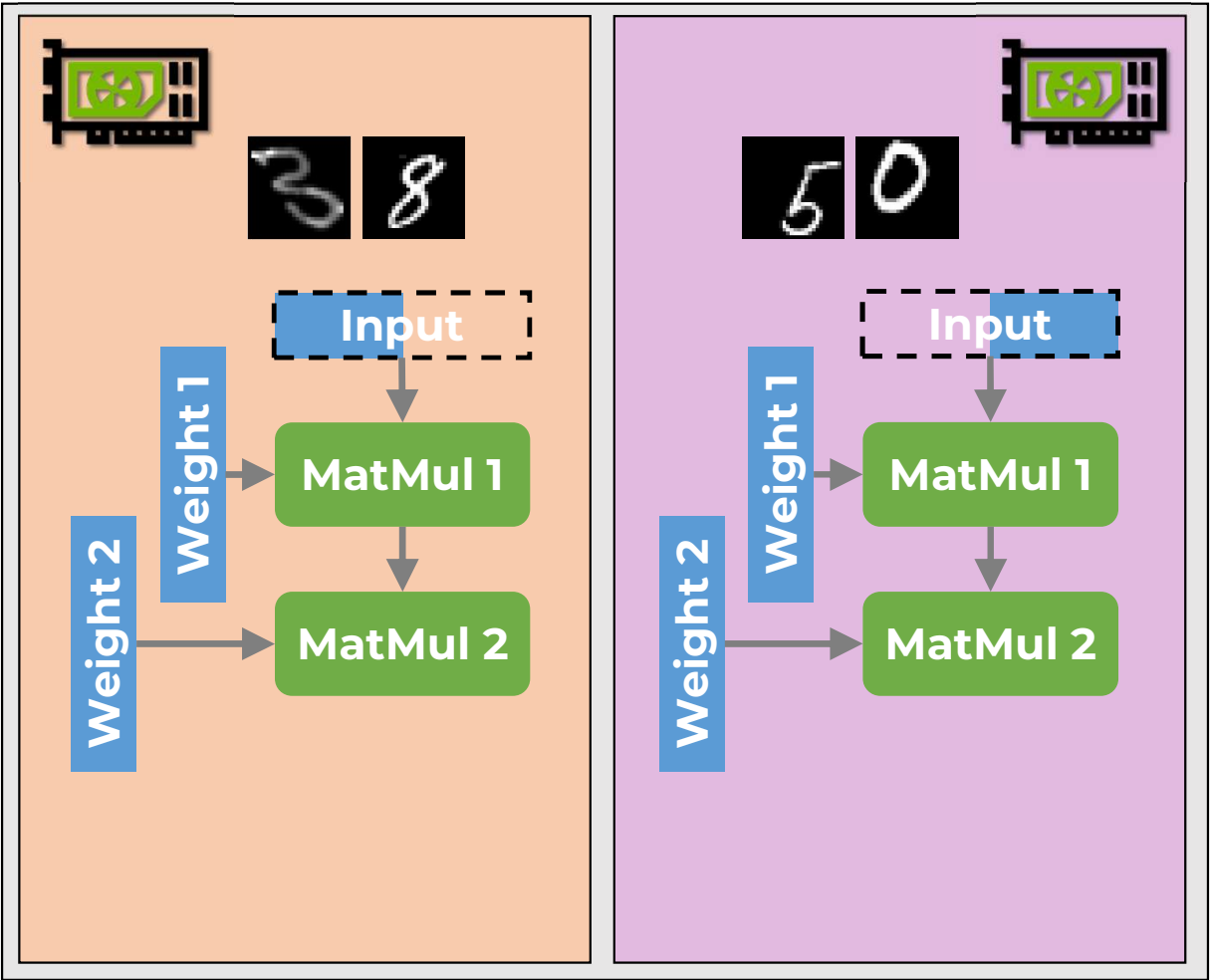
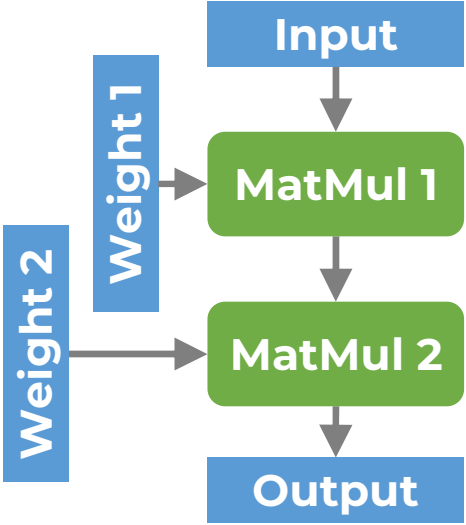
# Data Parallel



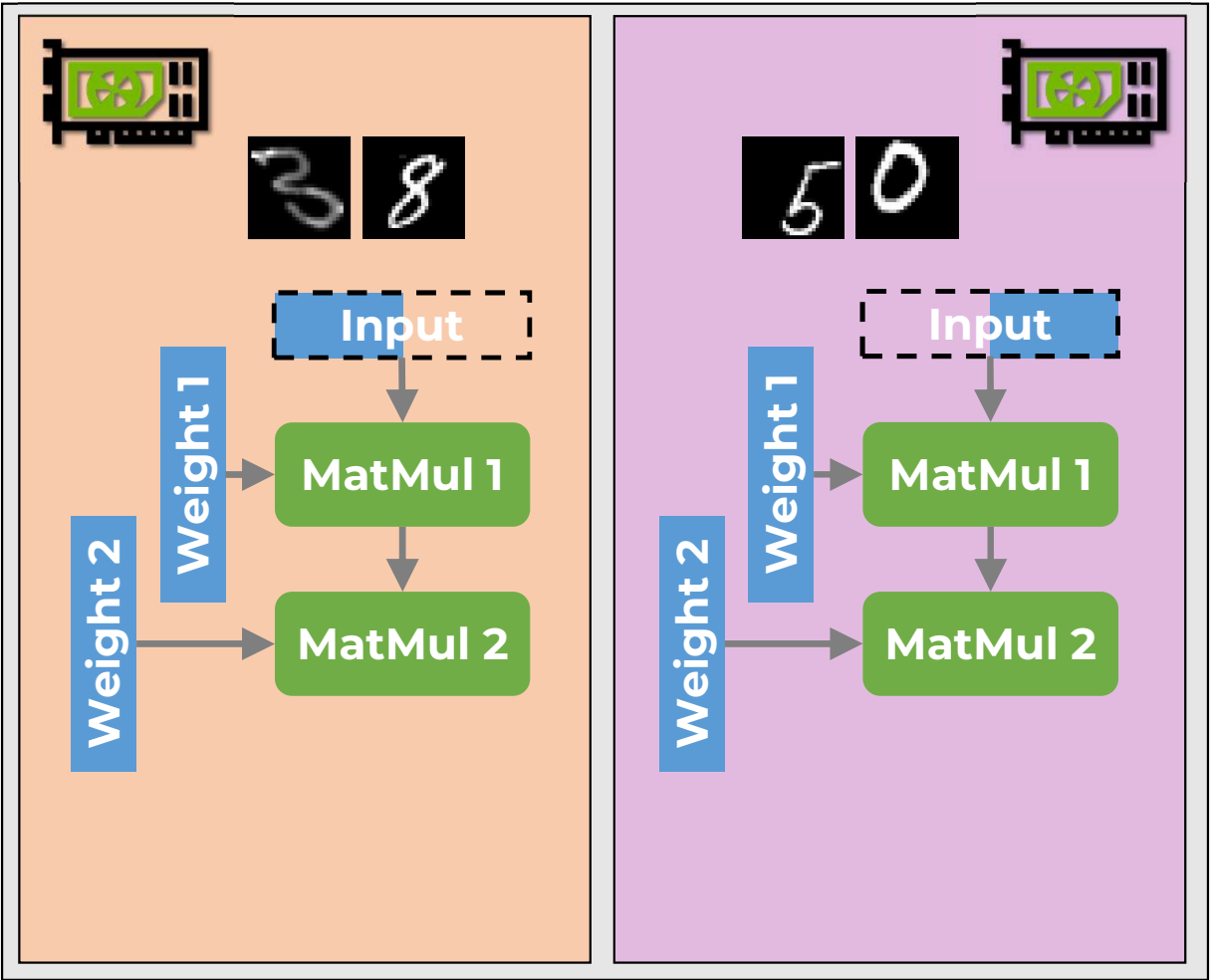
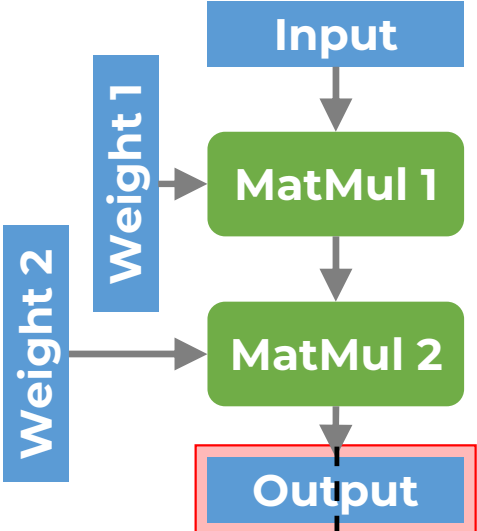
# Data Parallel



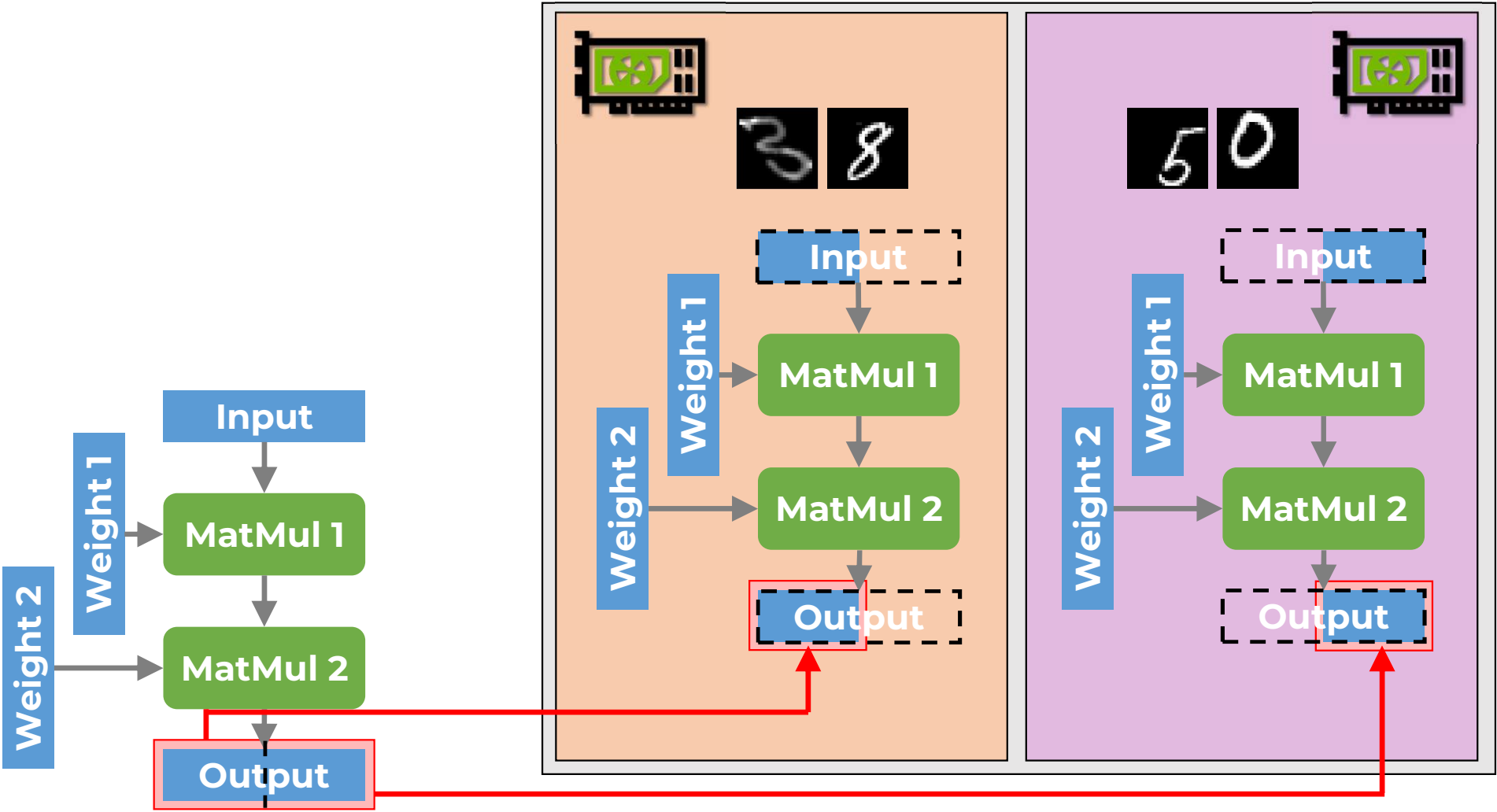
# Data Parallel



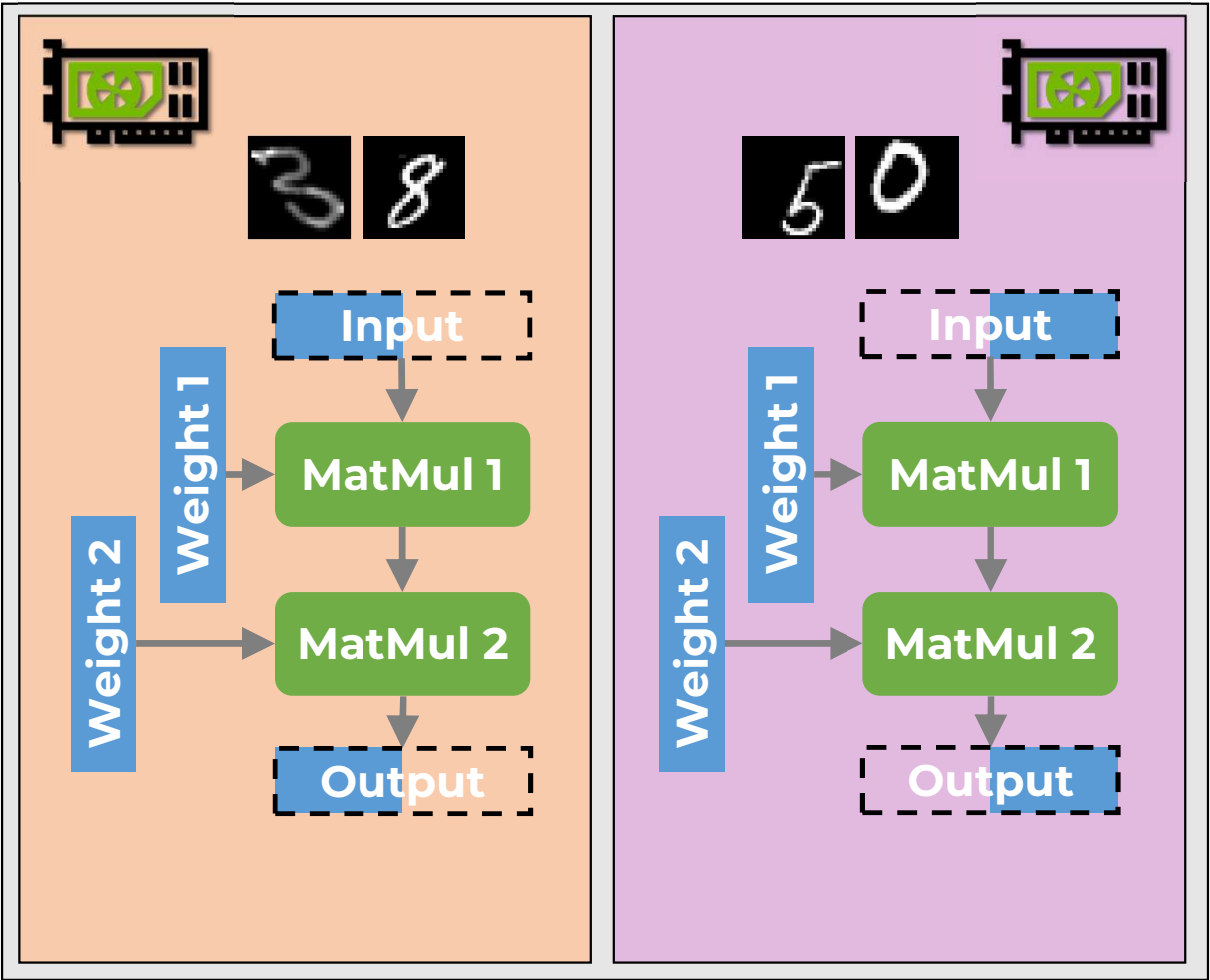
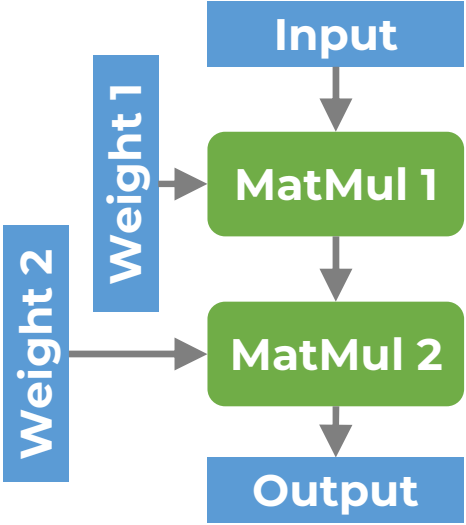
# Data Parallel



# Data Parallel

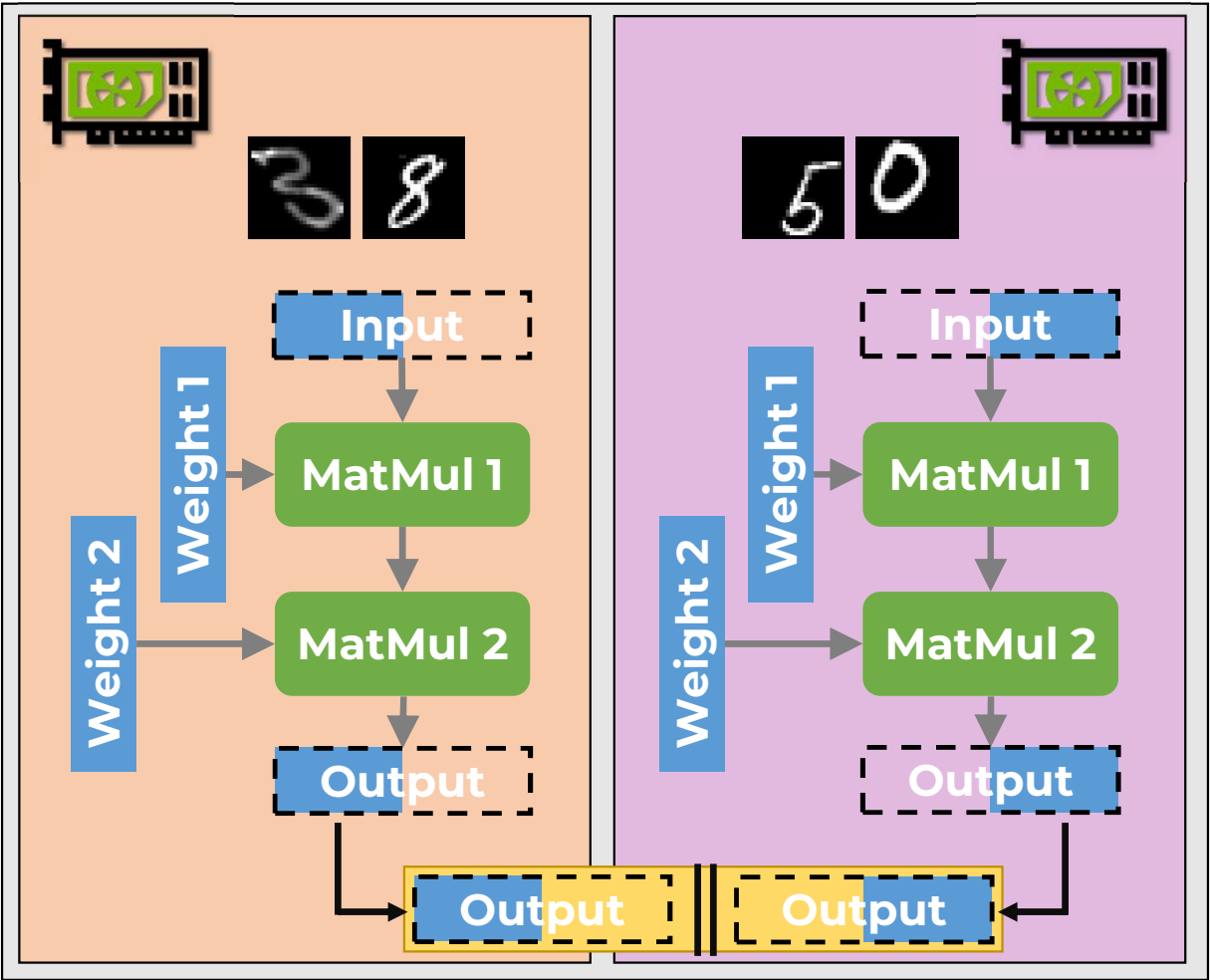
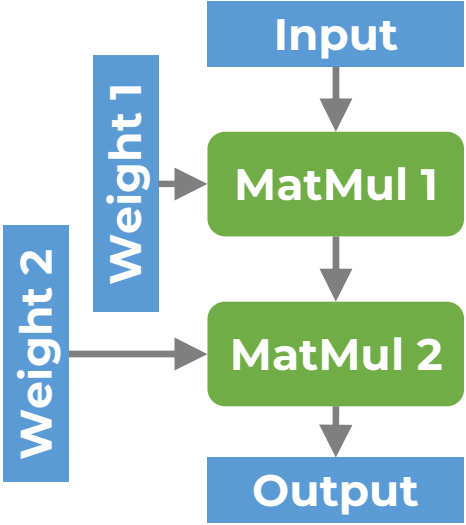


# Data Parallel

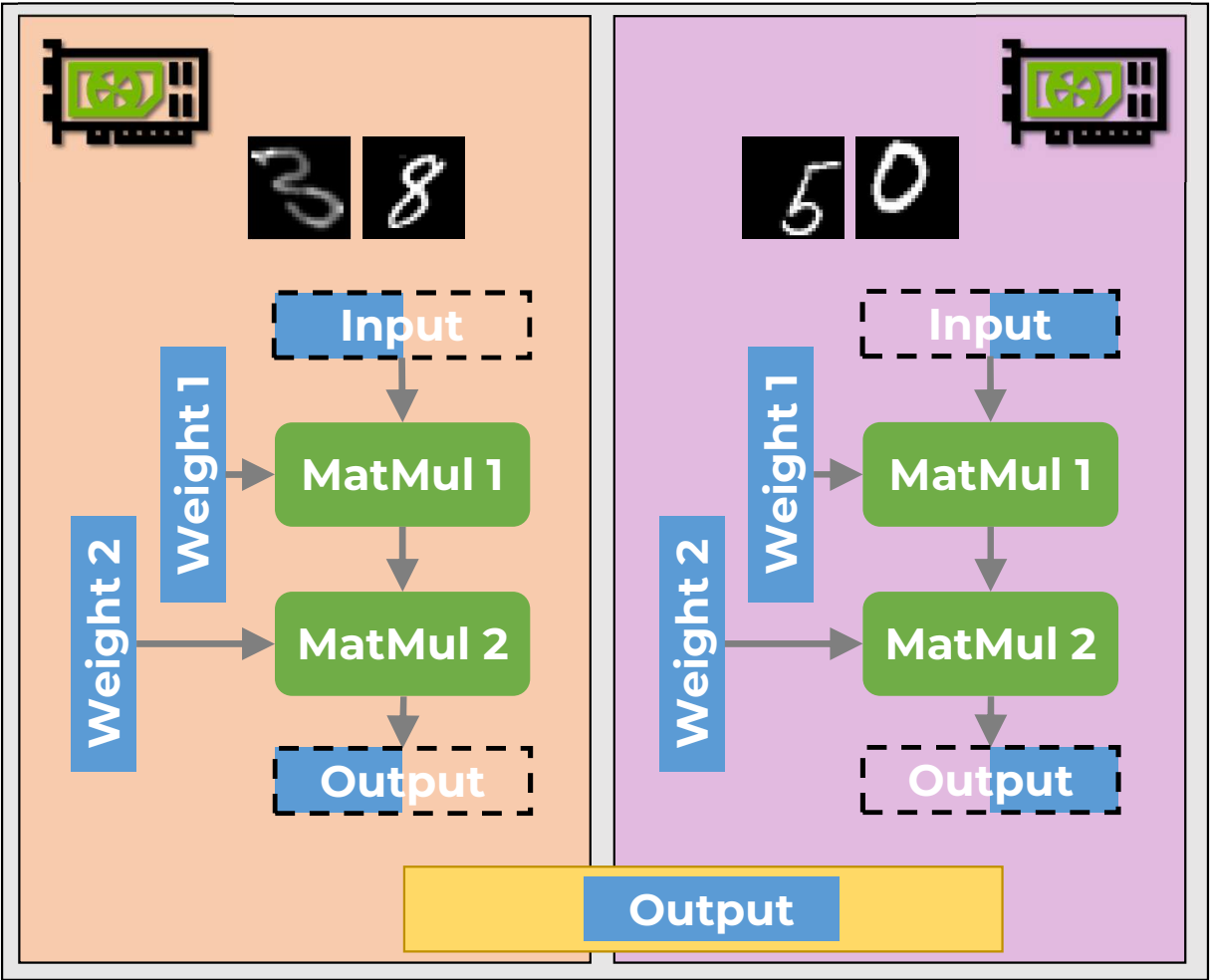
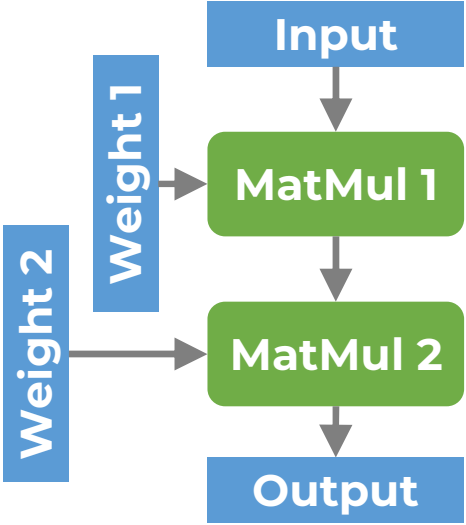




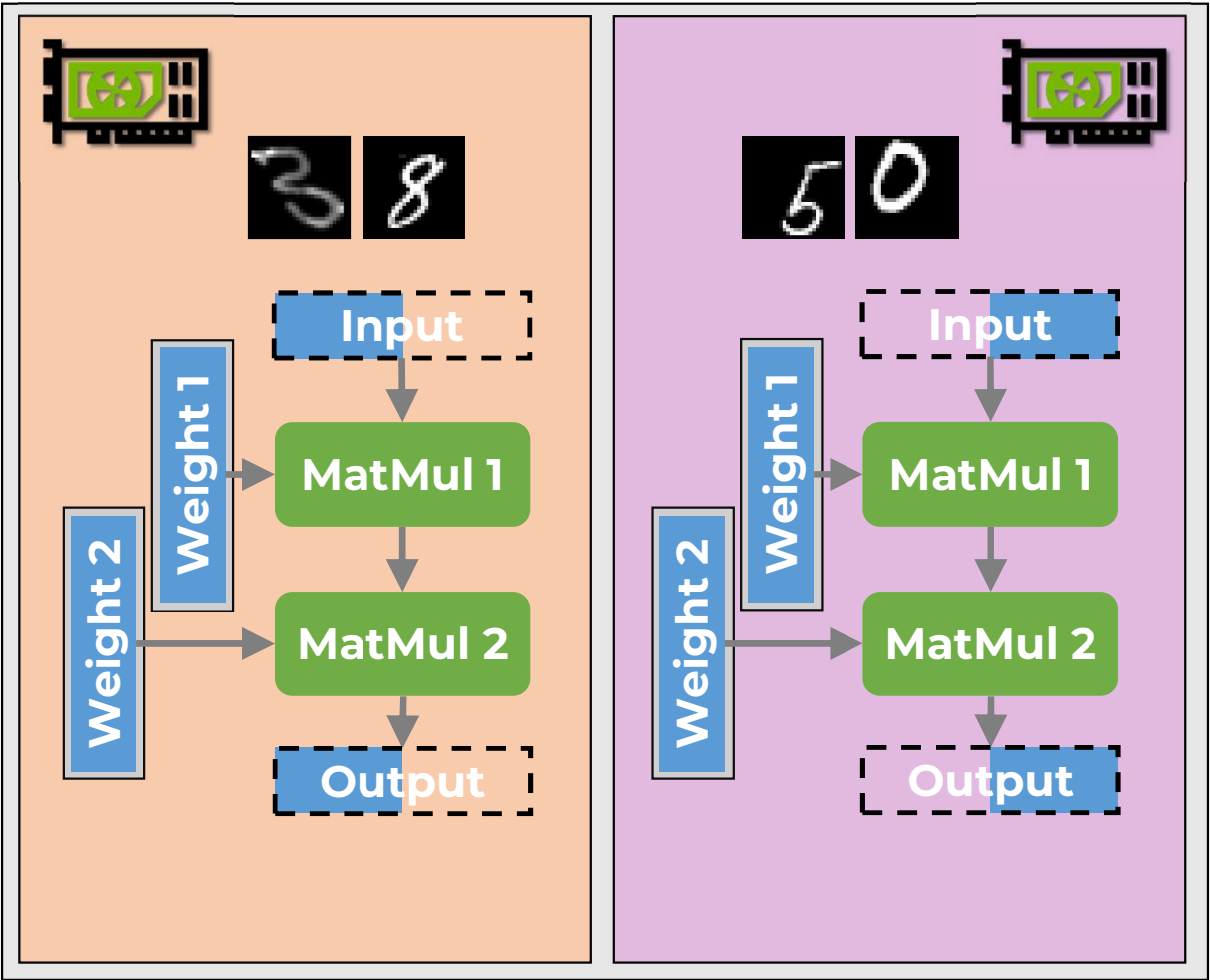
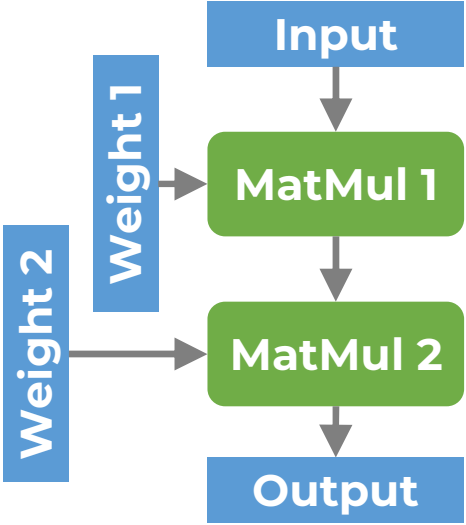
# Data Parallel



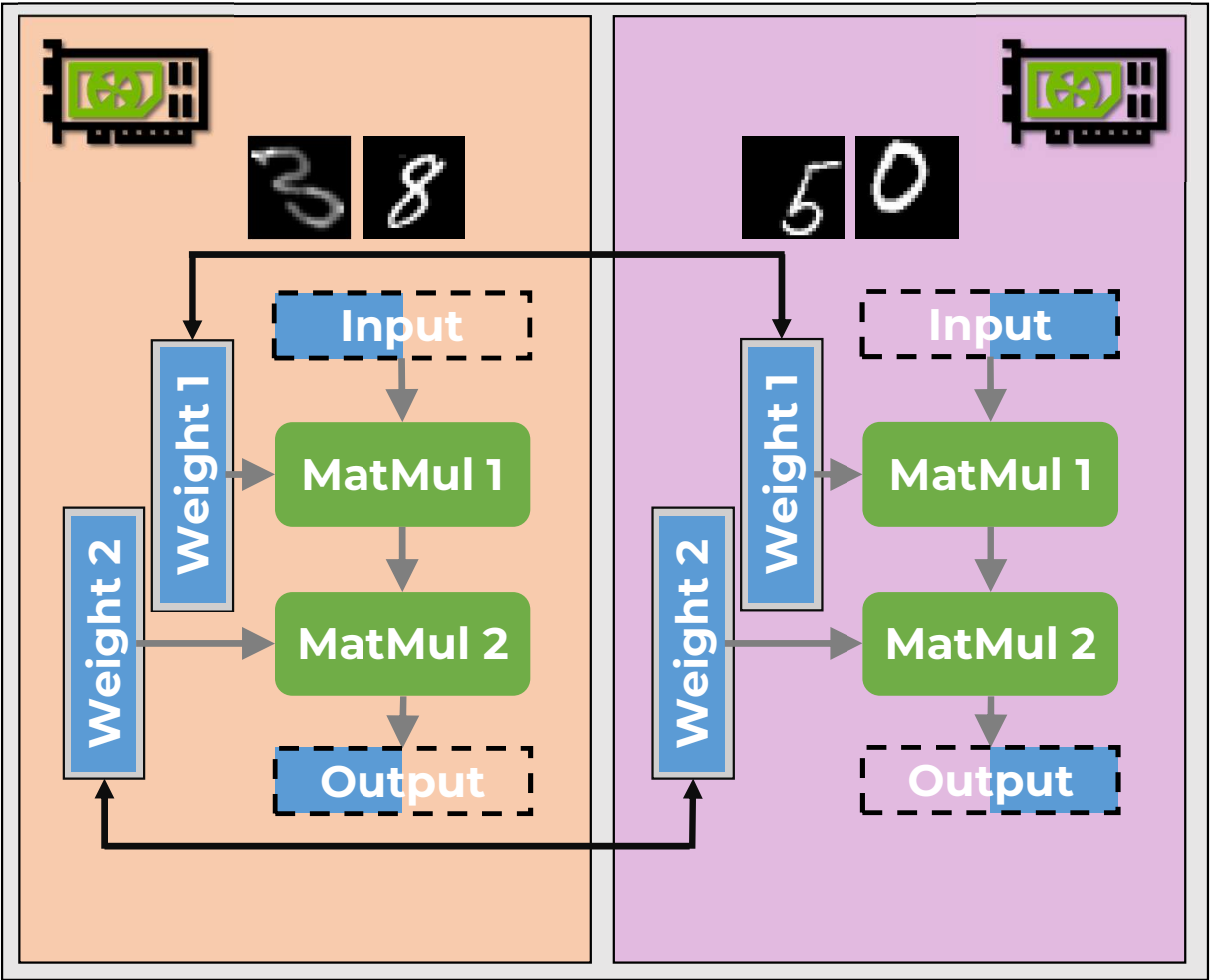
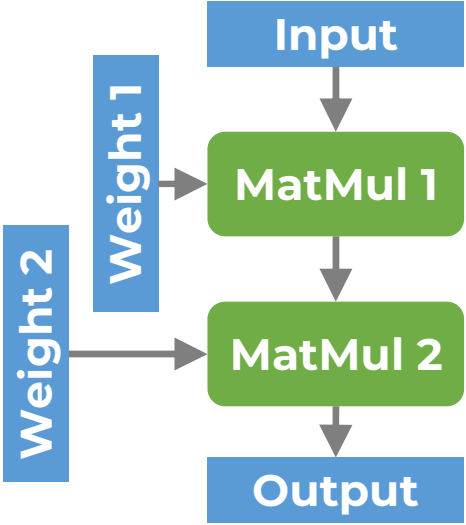
# Data Parallel



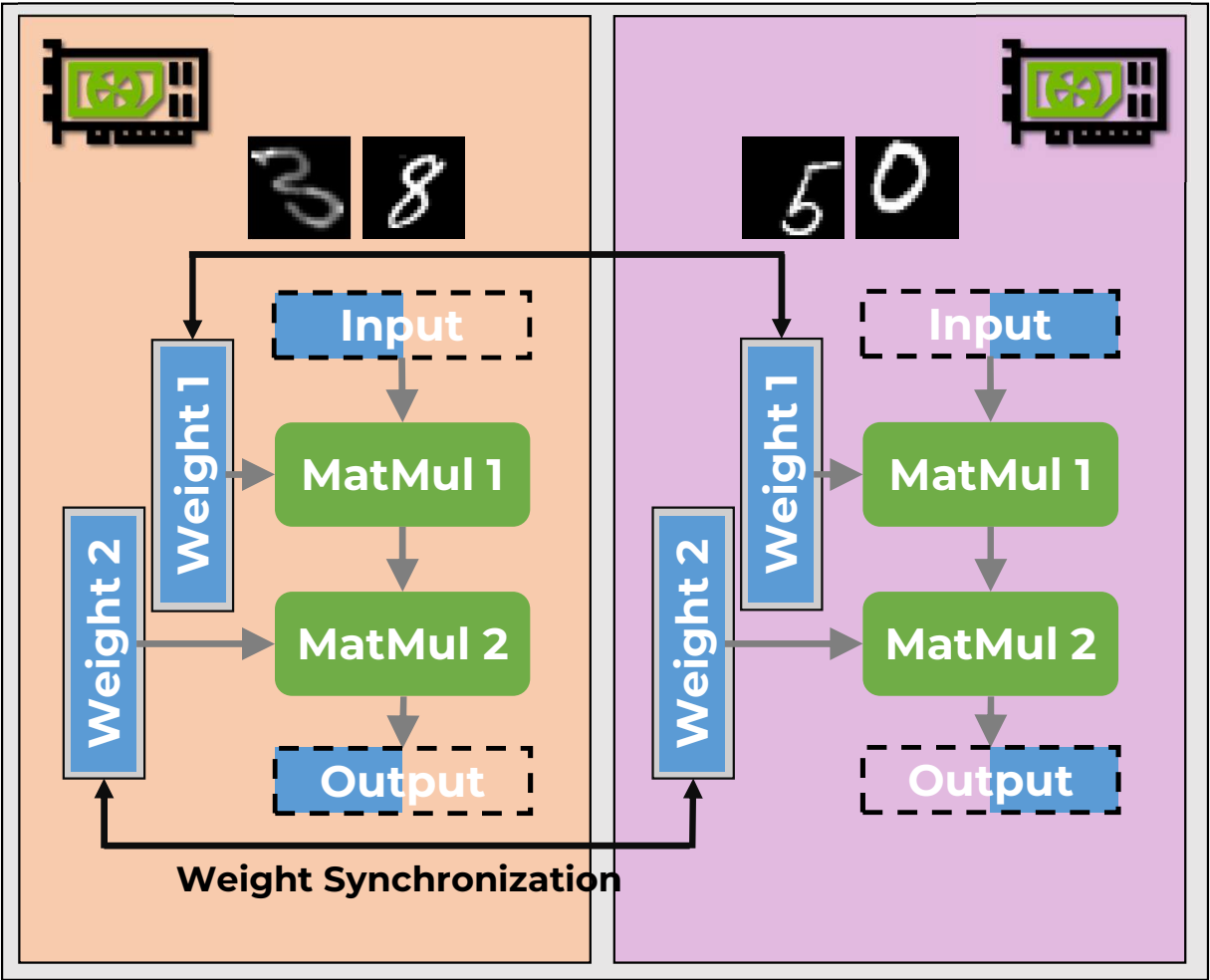
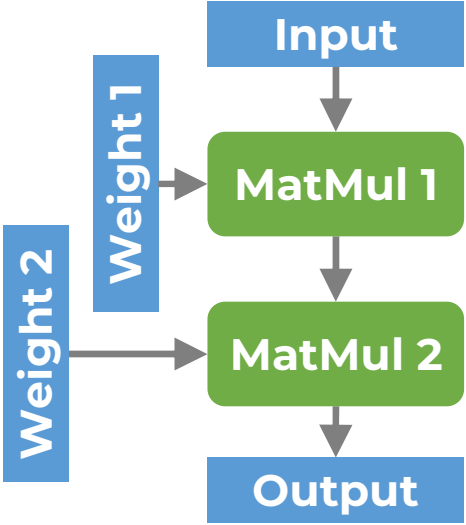
# Data Parallel



# Data Parallel



# Data Parallel



# Data Parallel

# Data Parallel Model Parallel

Data Parallel  
Model Parallel  
Attribute Parallel



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

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

...

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

Parallelization

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

2 | Parallelization |

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

2 | Parallelization |

## Algebraic Transformations

Operator Fusion  
Operator Splitting  
Operator Reordering  
...

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

2 | Parallelization |

2 | 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]

...

## 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

FlexFlow [MLSys 19]

Tofu [EuroSys 19]

PipeDream [SOSP 19]

automap [arXiv 19]

Whale [arXiv 21]

→ Alpa [OSDI 22] ←

...

## Algebraic Optimizers

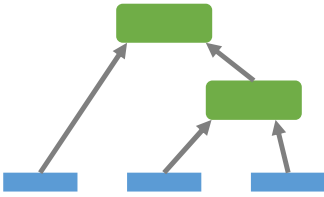
MetaFlow [MLSys 19]

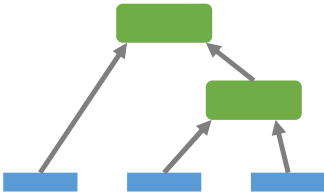
TASO [SOSP 19]

PET [OSDI 21]

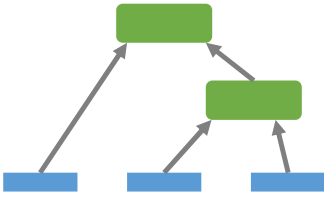
Tensat [MLSys 21]

...



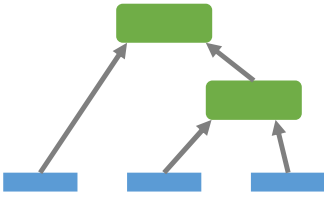


# Auto-Parallelization



**Auto-Parallelization**

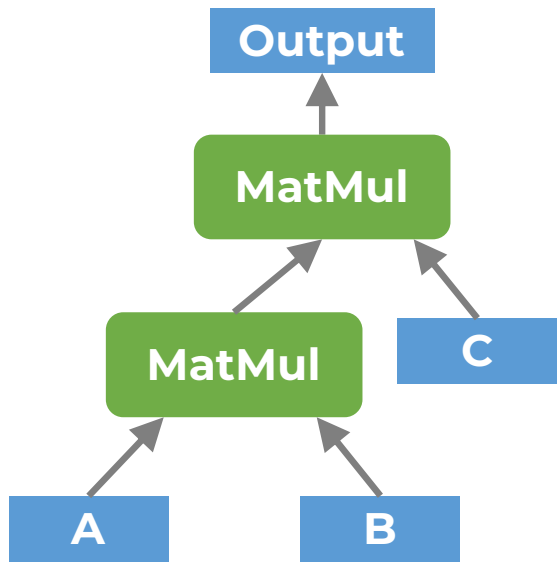
**Algebraic Optimizer**



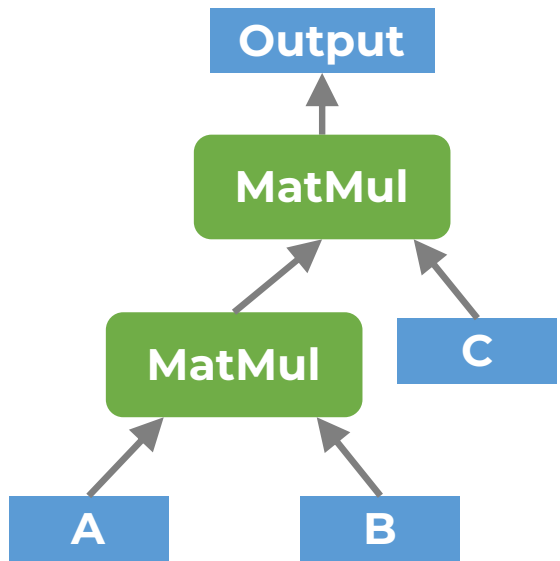
**Auto-Parallelization**

**Algebraic Optimizer**

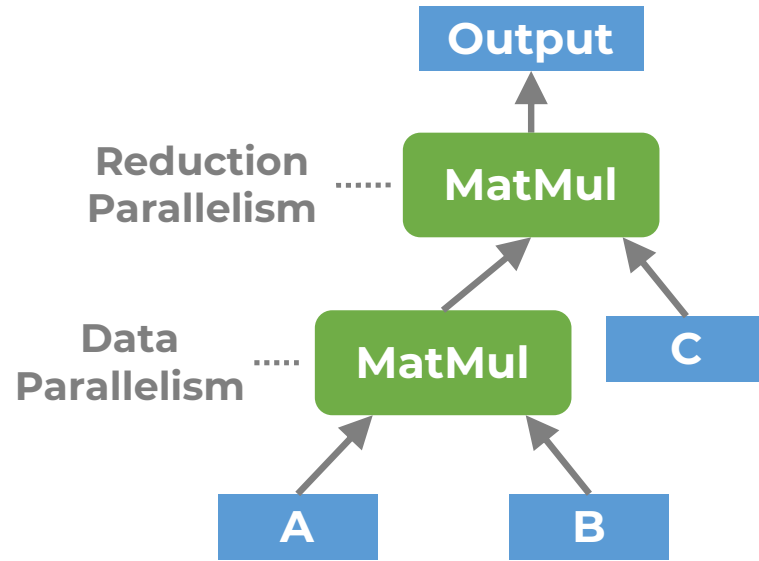
?

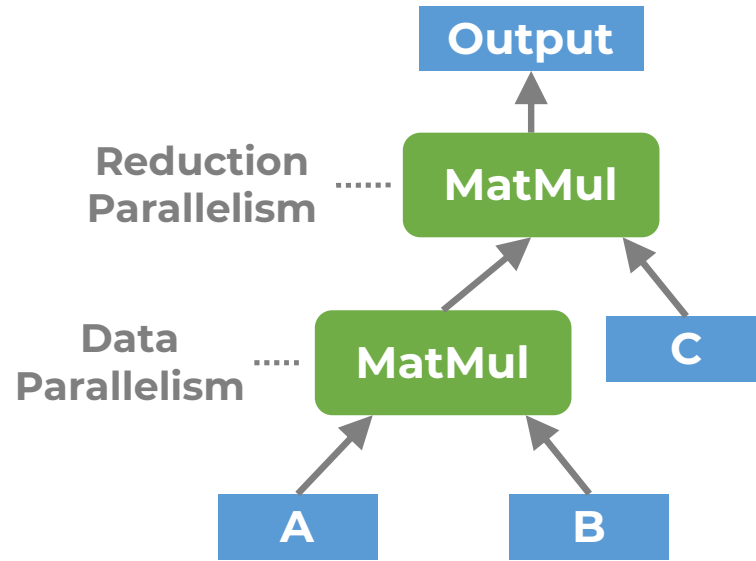




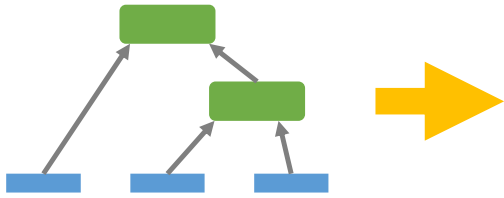


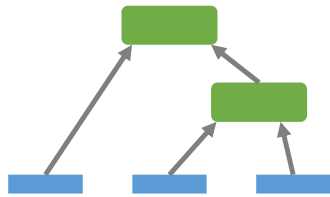
**“computation graph”**





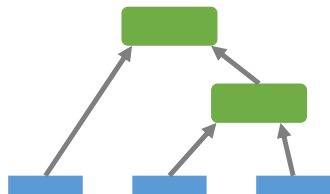
**“annotated computation graph”**





**Auto-Parallelization**



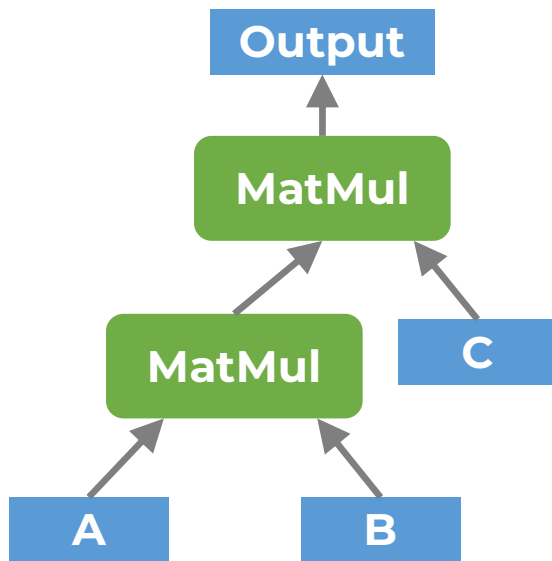


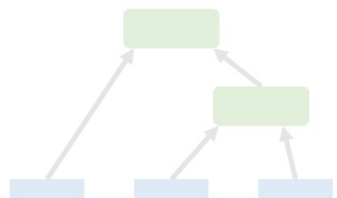
**Auto-Parallelization**



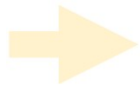
**Algebraic Optimizer**



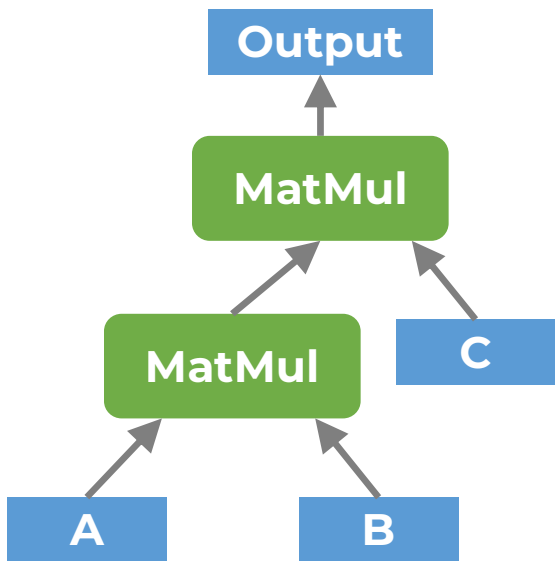
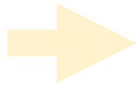




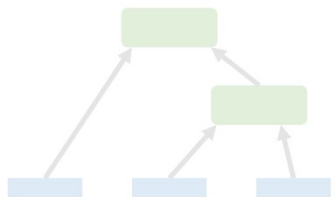
**Auto-Parallelization**



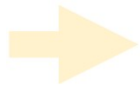
**Algebraic Optimizer**



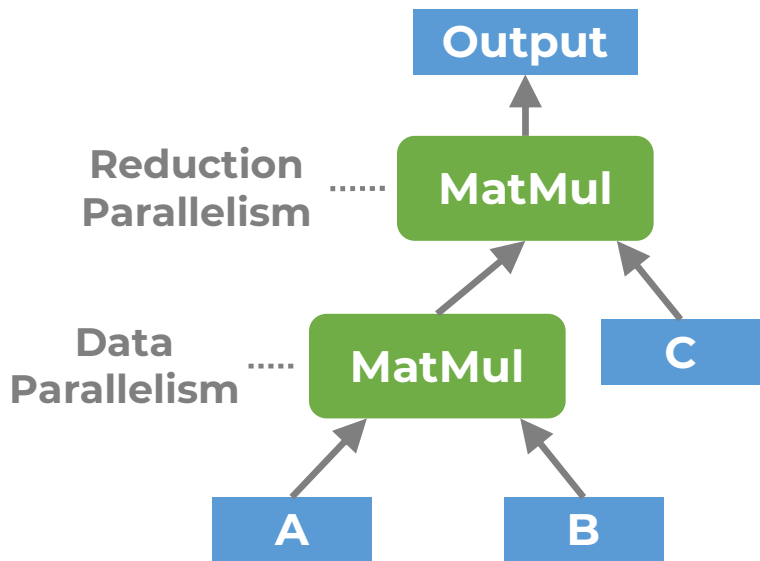
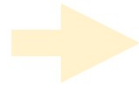


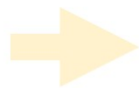
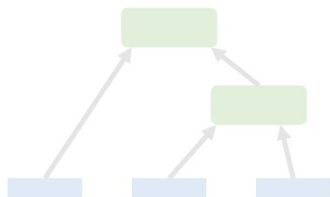


**Auto-Parallelization**

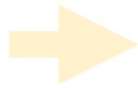


**Algebraic Optimizer**

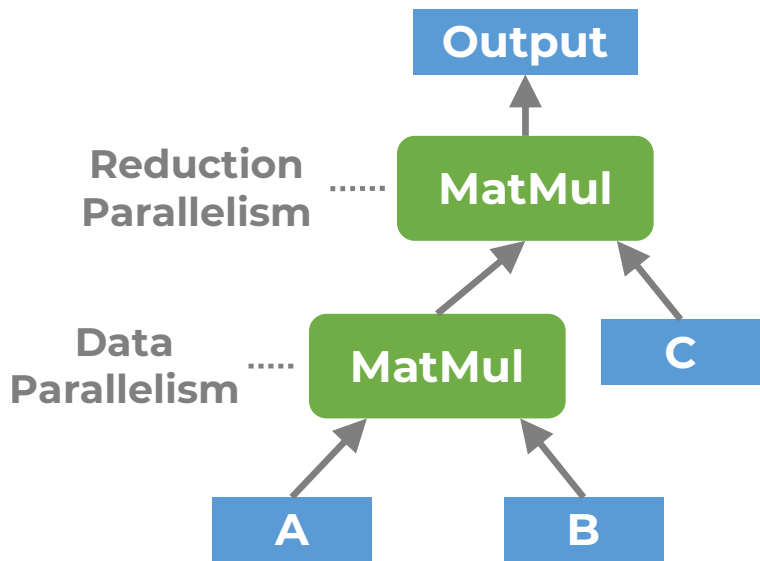
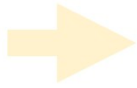


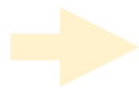
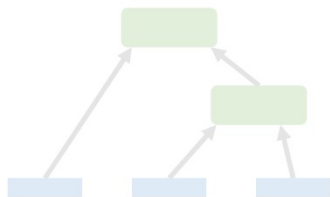


Auto-Parallelization

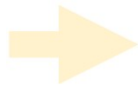


Algebraic Optimizer

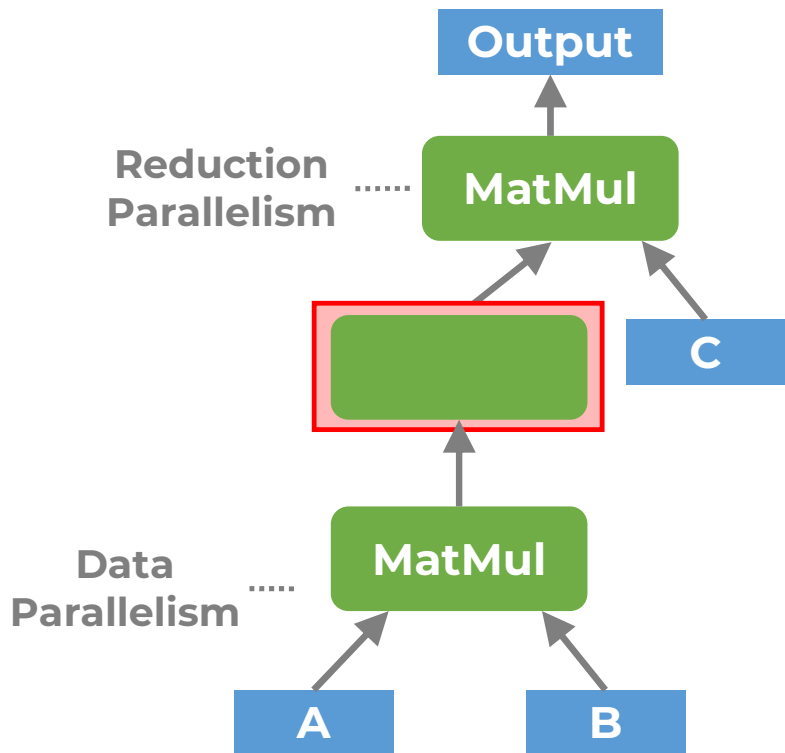
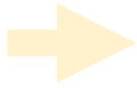


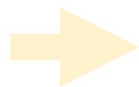
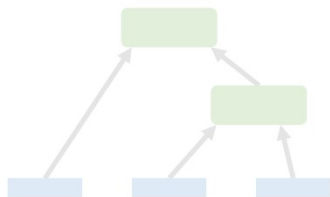


Auto-Parallelization

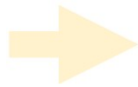


Algebraic Optimizer

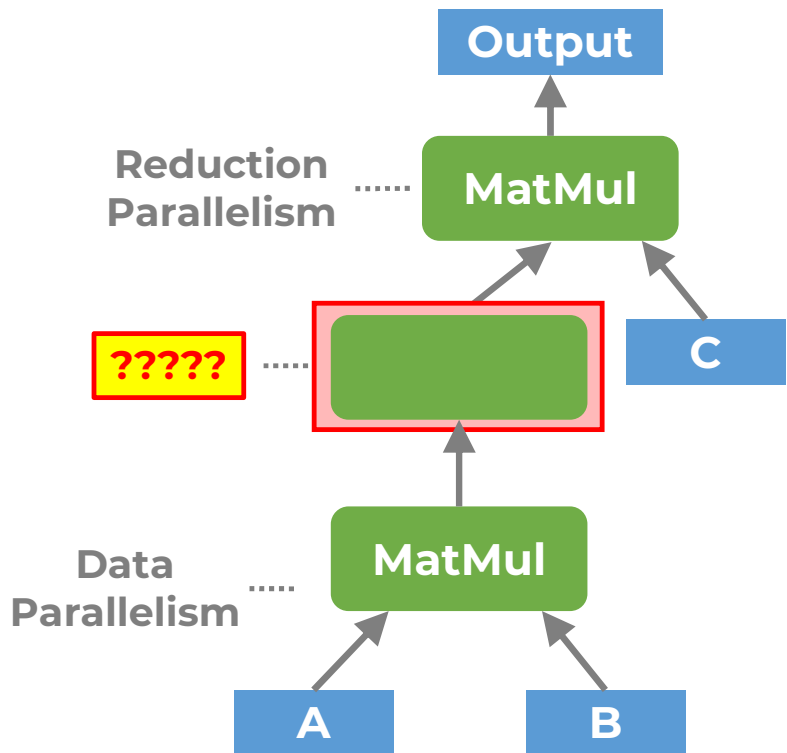
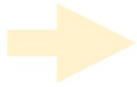


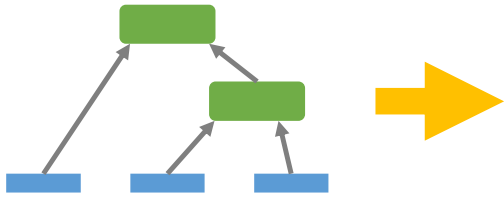


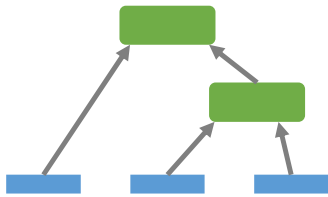
Auto-Parallelization



Algebraic Optimizer

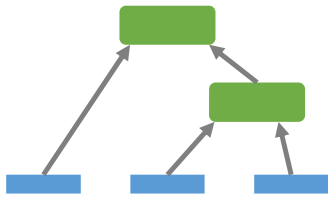






**Algebraic  
Optimizer**



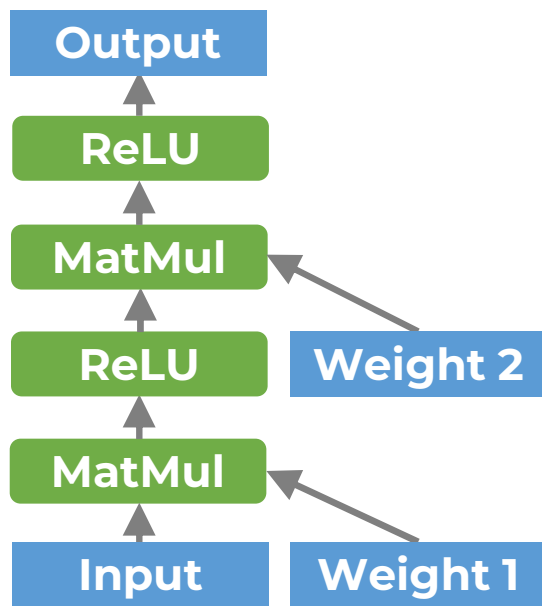


**Algebraic  
Optimizer**

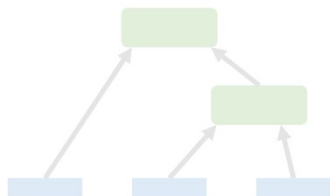


**Auto-  
Parallelization**

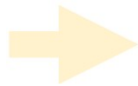




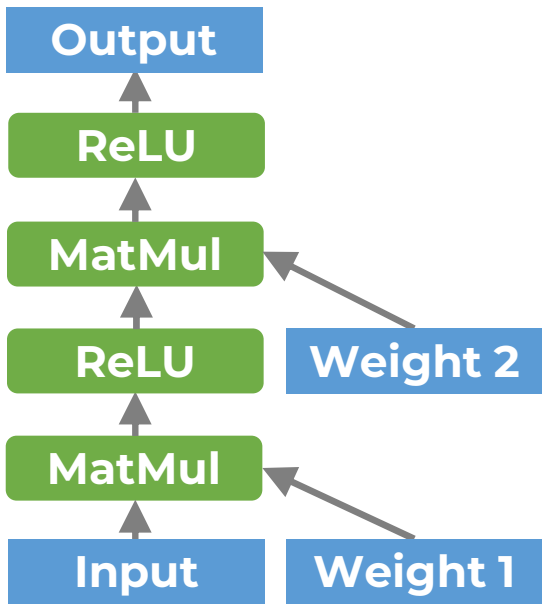
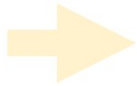


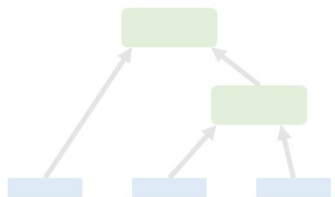


**Algebraic  
Optimizer**

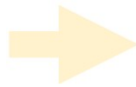


**Auto-  
Parallelization**

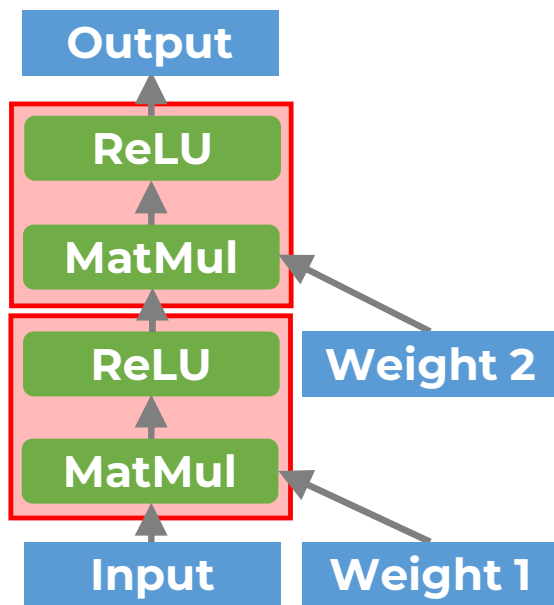
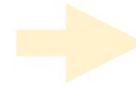


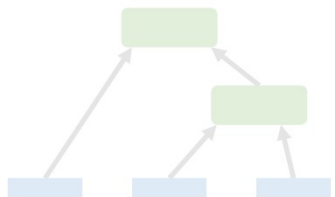


**Algebraic  
Optimizer**

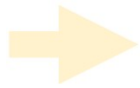


**Auto-  
Parallelization**

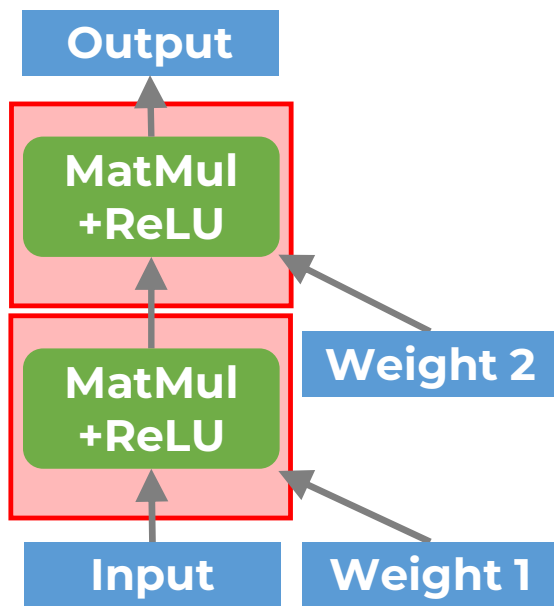
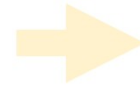


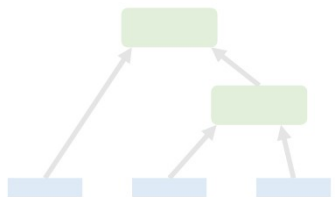


**Algebraic  
Optimizer**

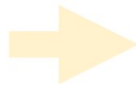


**Auto-  
Parallelization**

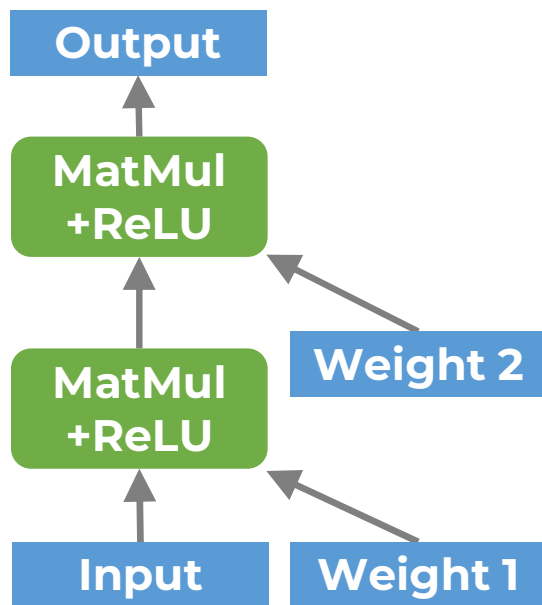
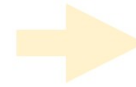


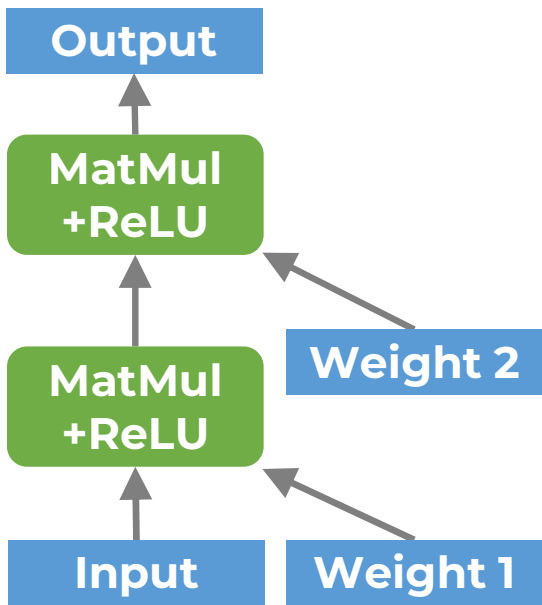
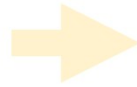
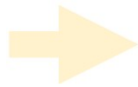
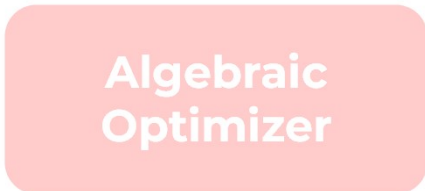
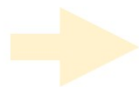
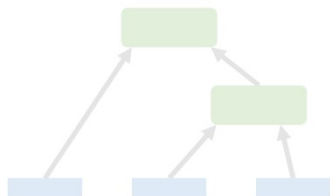


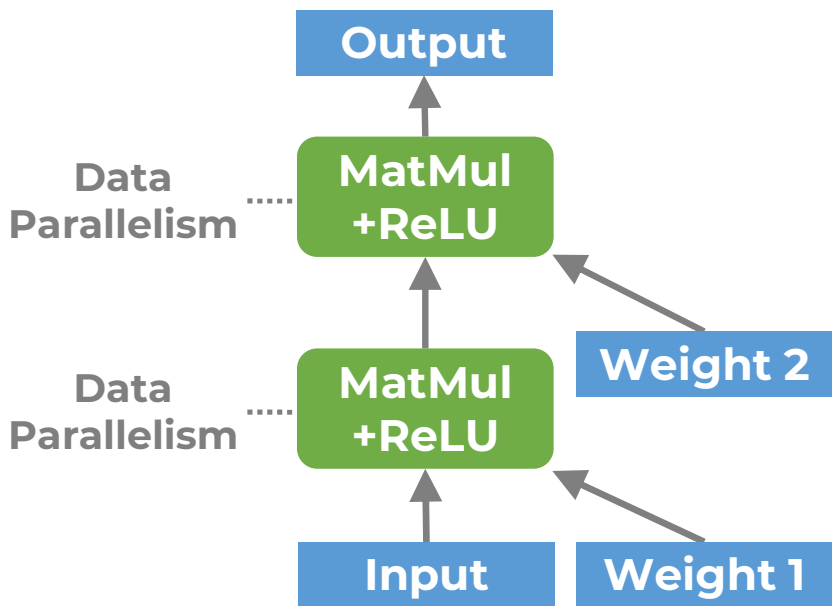
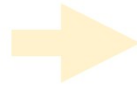
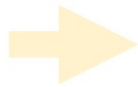
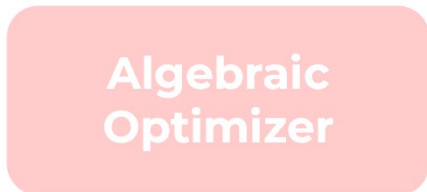
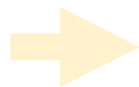
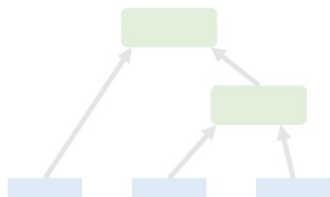
**Algebraic  
Optimizer**

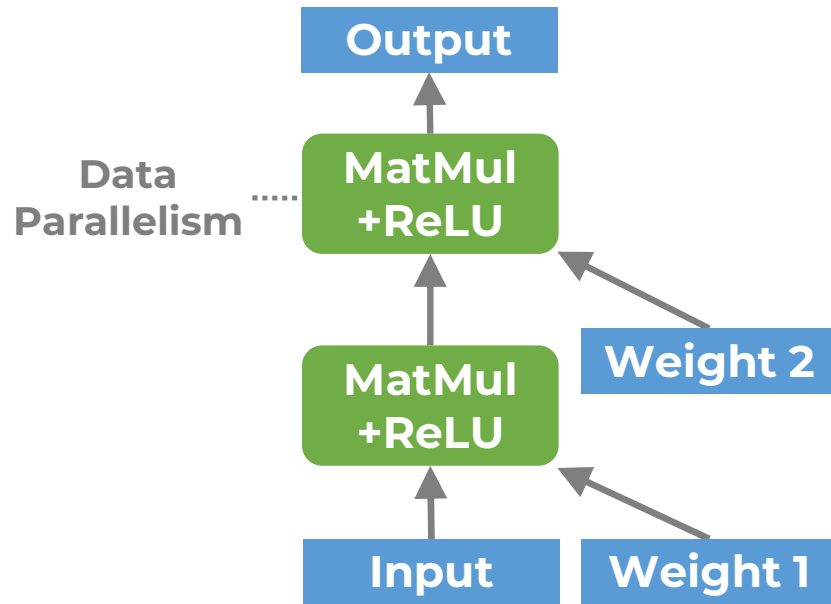


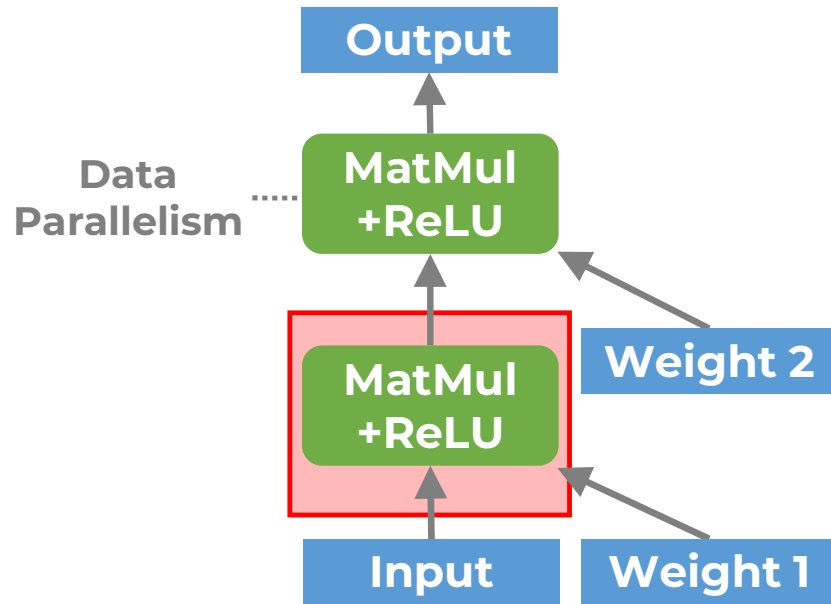
**Auto-  
Parallelization**



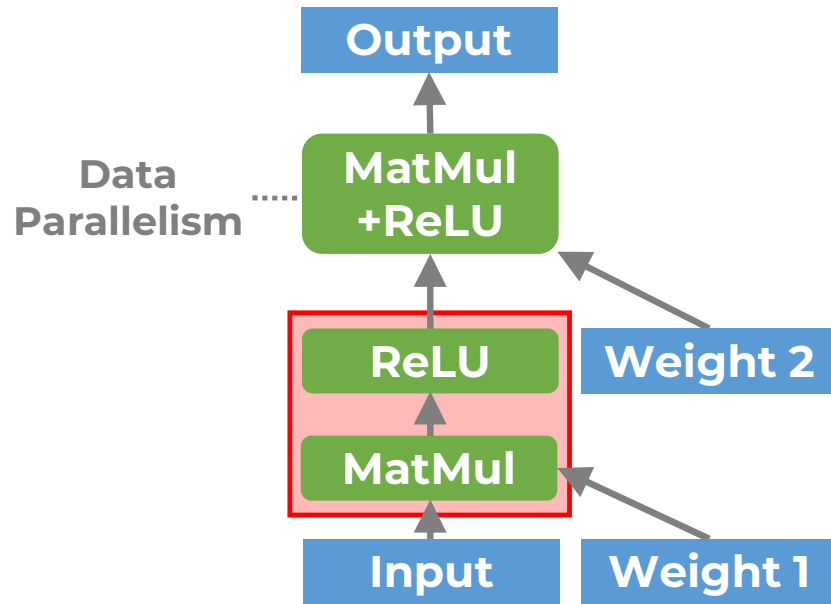


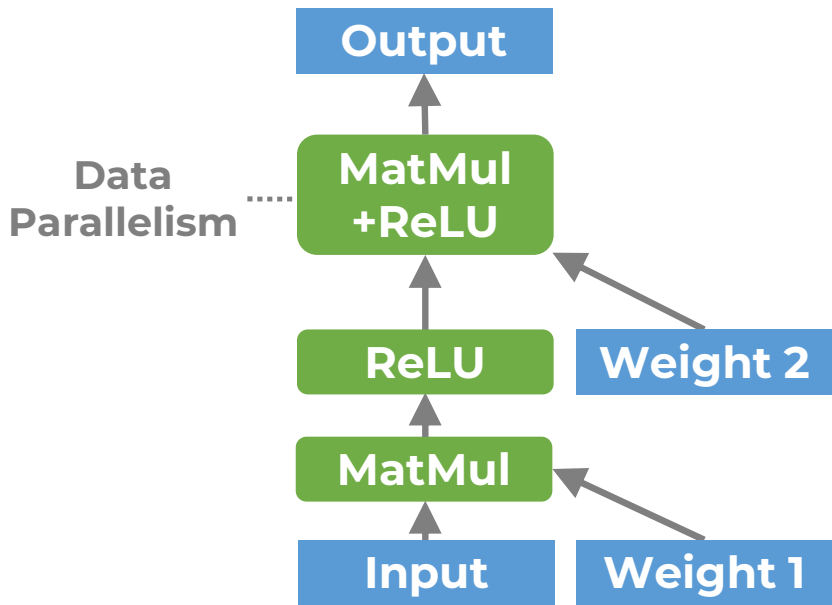
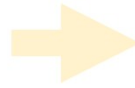
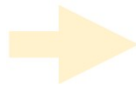
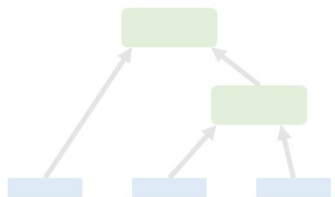


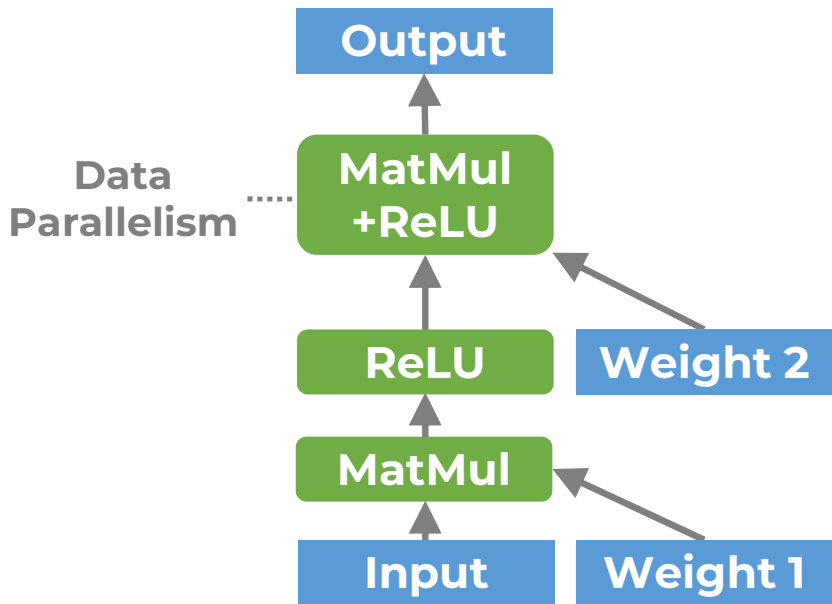
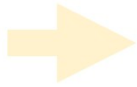
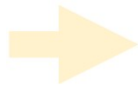
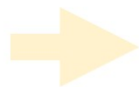
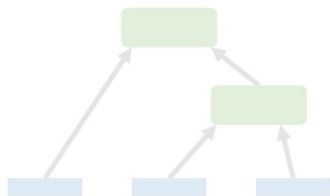


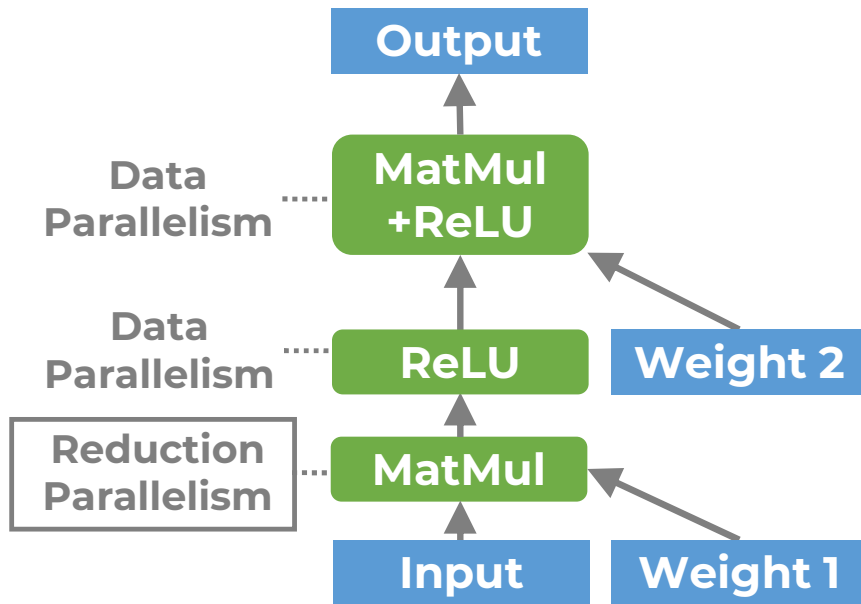


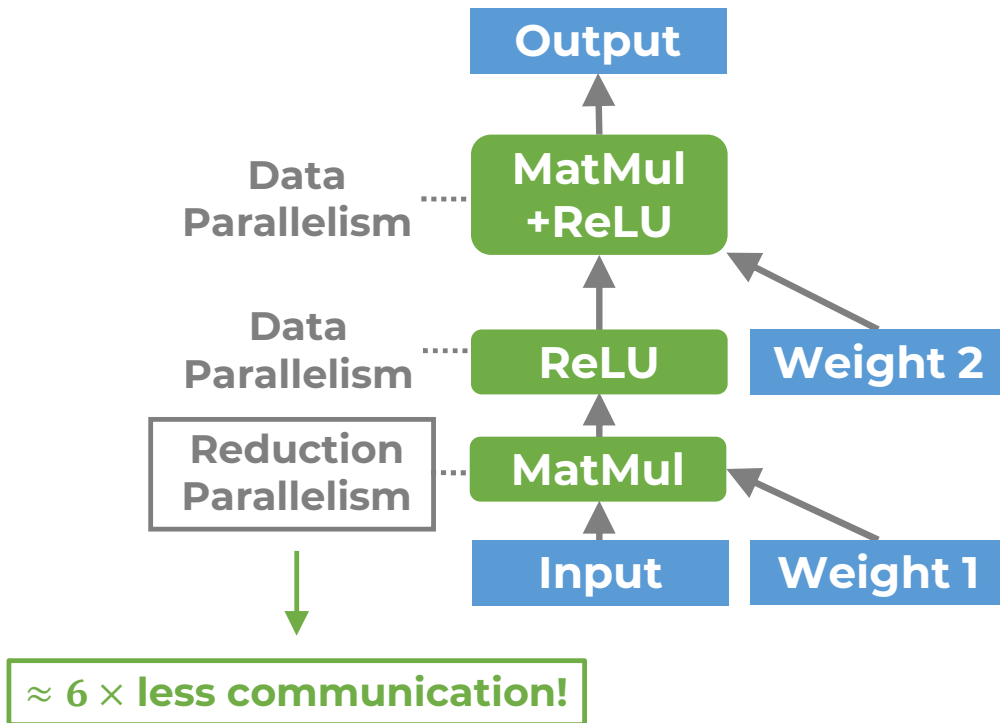










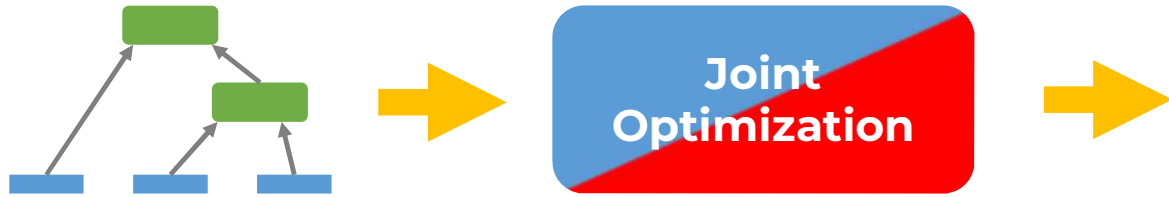






1.

2.

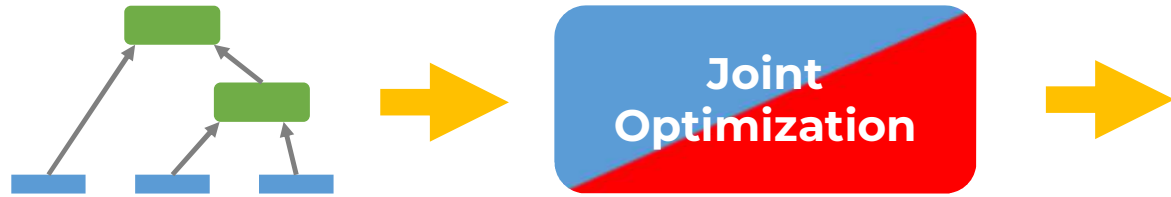


---

1. Representation

2.





1. Representation
2. Scalability

# Unity

Unity — Representation

Representation

Parallel Computation  
Graph (PCG)

Unity



Representation

Parallel Computation  
Graph (PCG)

Unity

Scalability

Representation

Parallel Computation  
Graph (PCG)

Unity

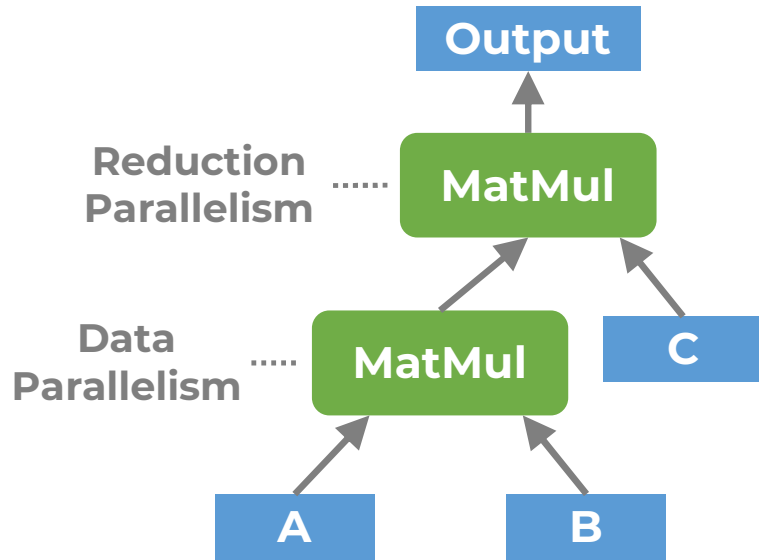
Scalability

Hierarchical Search  
Algorithm

Representation

# Parallel Computation Graph (PCG)

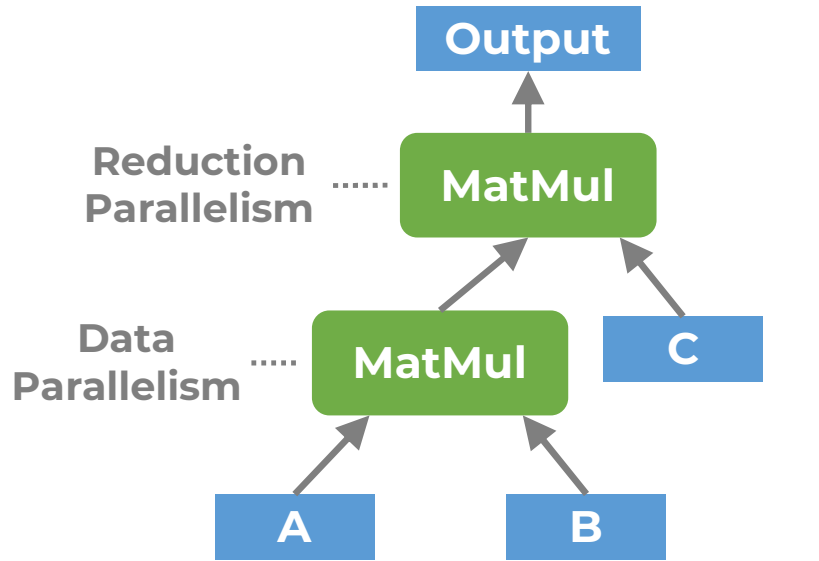
# Parallel Computation Graph (PCG)



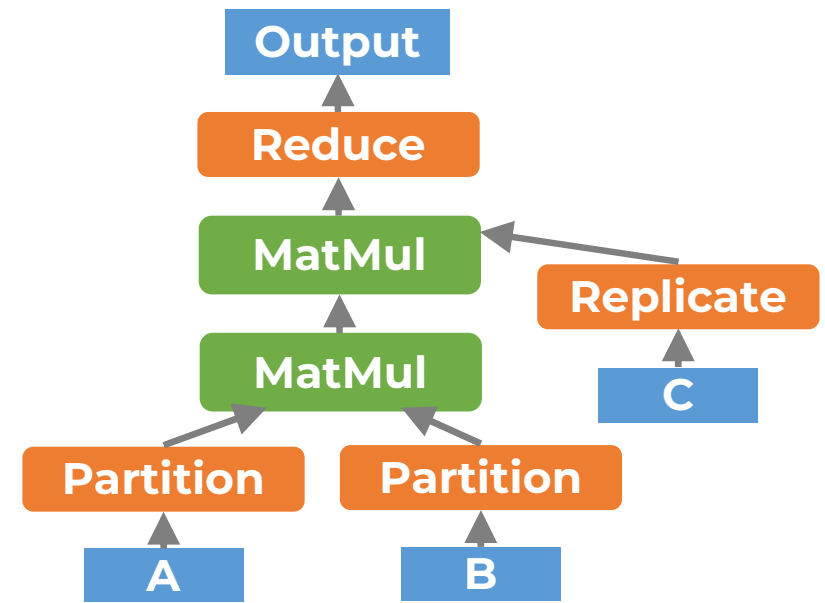
**annotated computation graph**



# Parallel Computation Graph (PCG)

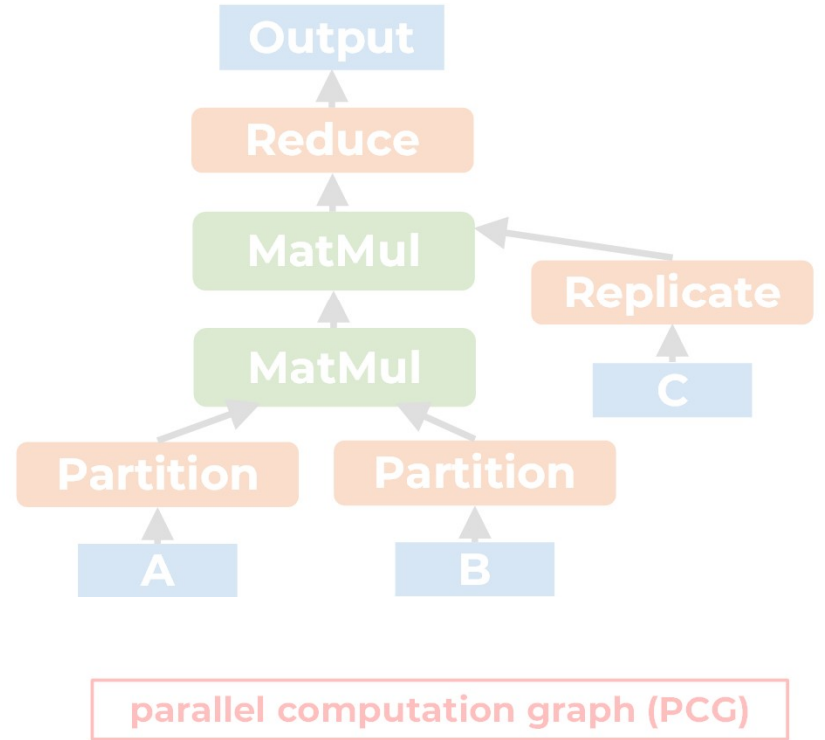
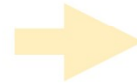
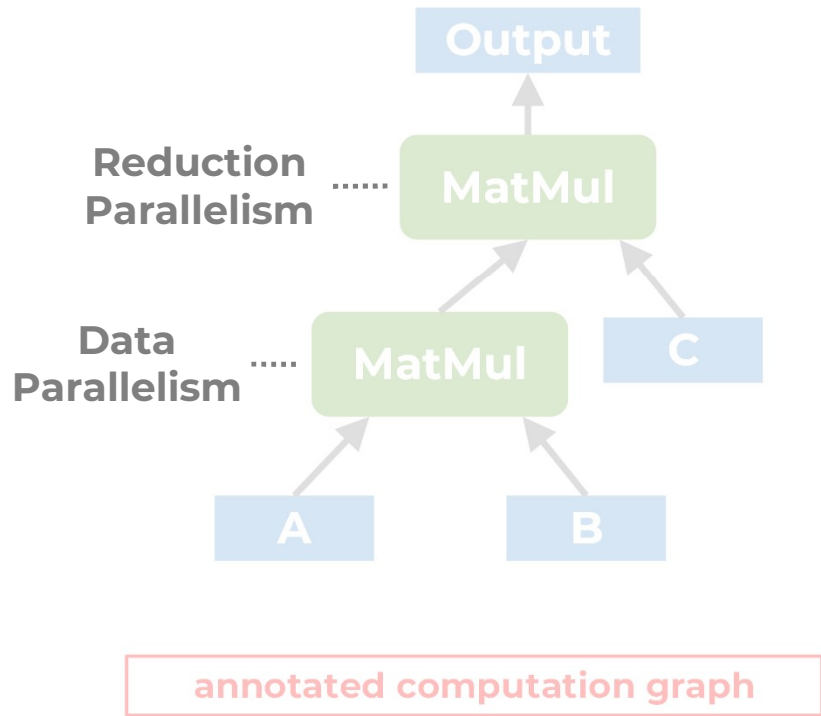


**annotated computation graph**

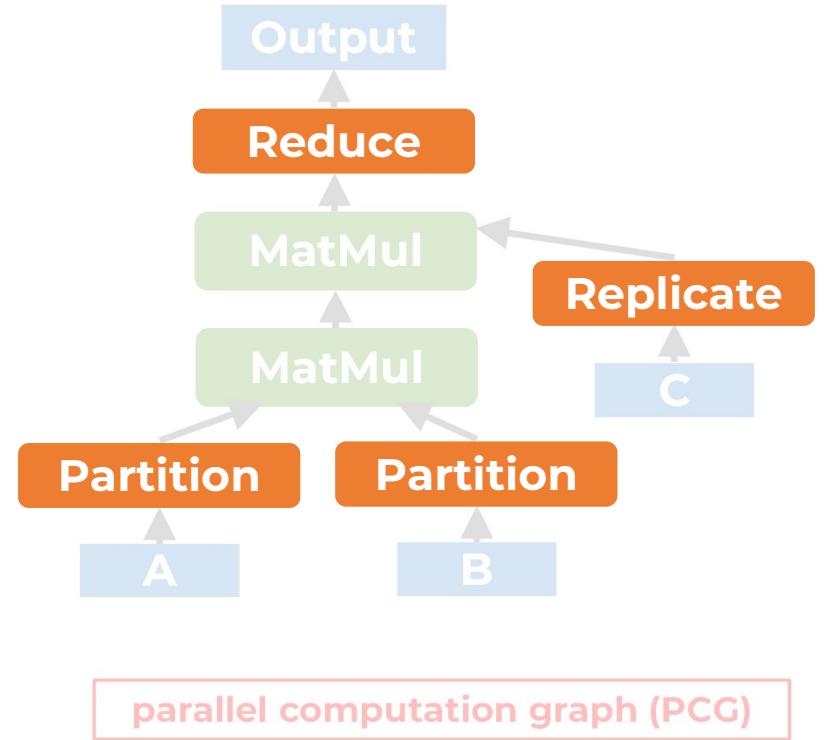
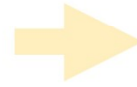
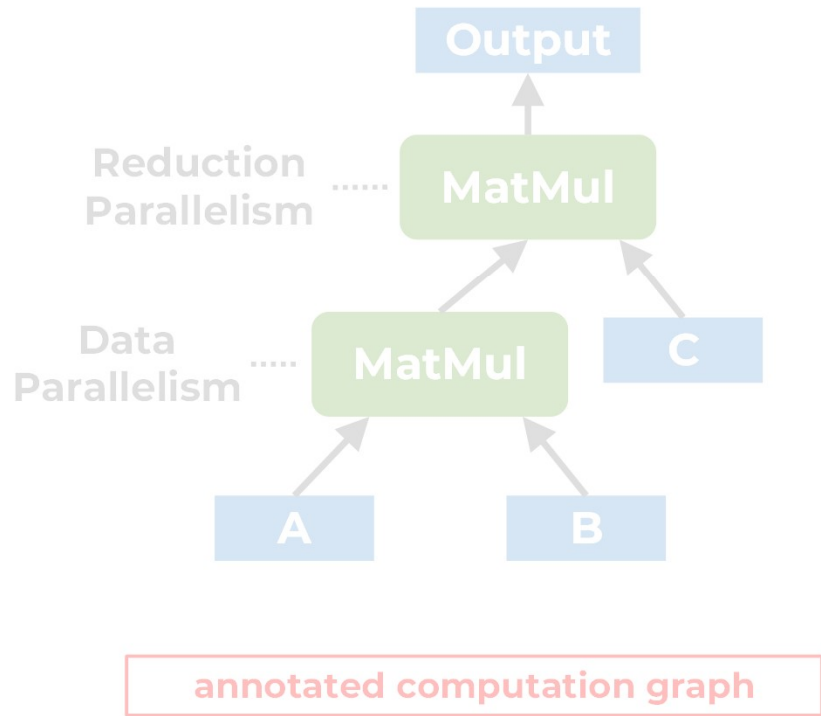


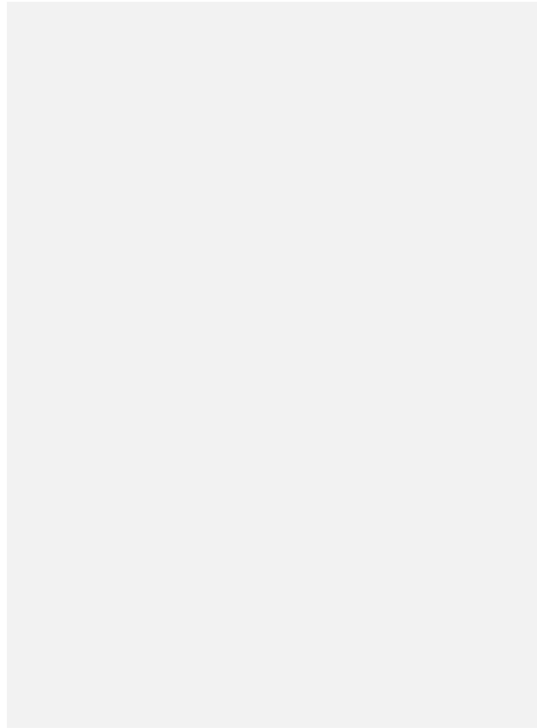
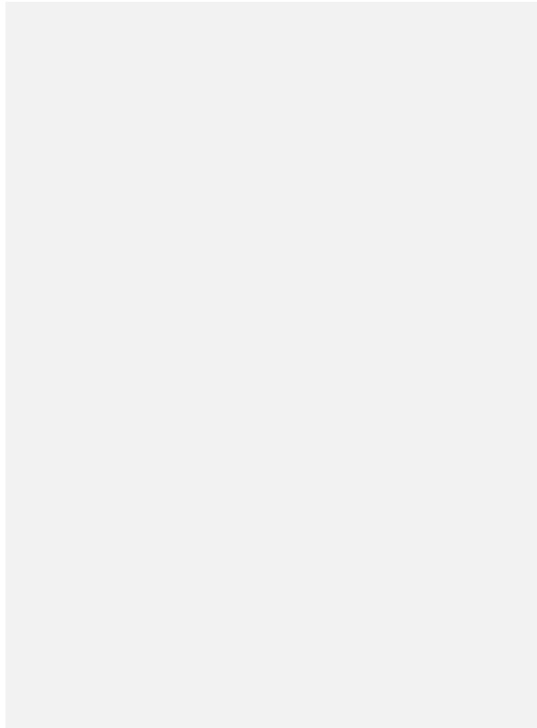
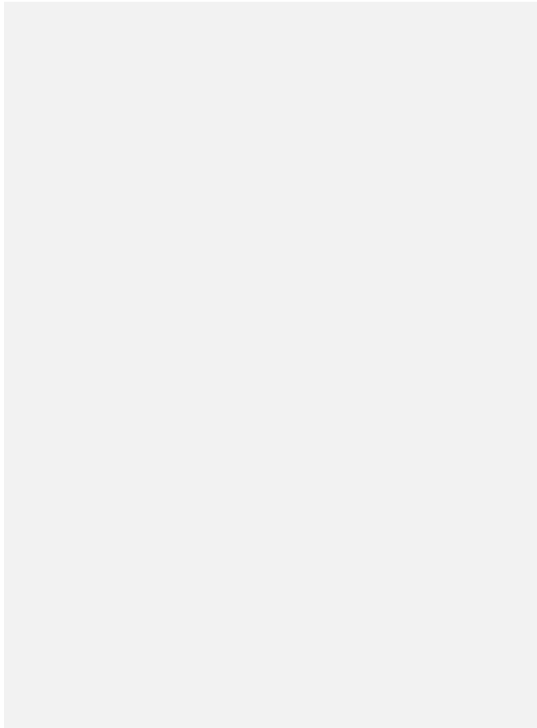
**parallel computation graph (PCG)**

# Parallel Computation Graph (PCG)



# Parallel Computation Graph (PCG)





**Partition**

**Combine**

**Partition**

**Combine**

**Replicate**

**Reduce**

**Partition**

**Combine**

**Replicate**

**Reduce**

**Pipeline**

**Batch**

**Partition**

**Combine**

**Replicate**

**Reduce**

**Pipeline**

**Batch**



**Partition**

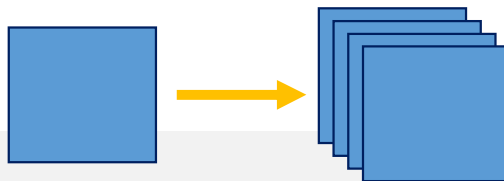
**Combine**

**Replicate**

**Reduce**

**Pipeline**

**Batch**



Partition

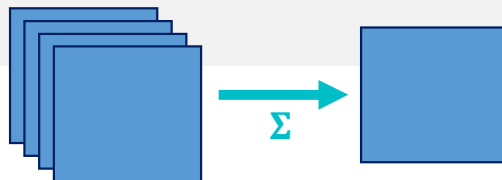
Replicate

Pipeline

Combine

Reduce

Batch



**Partition**



**Replicate**



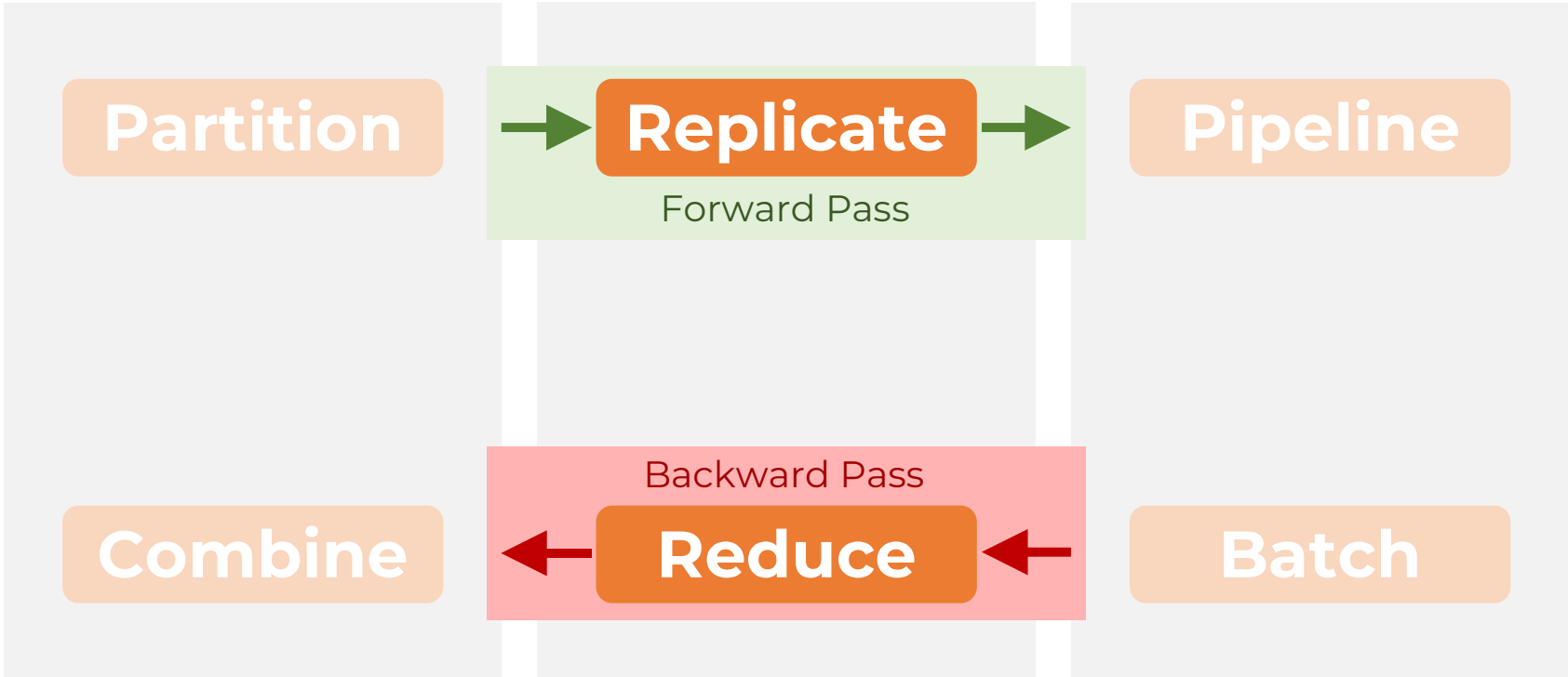
**Pipeline**

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**



**Partition**



**Combine**

**Replicate**

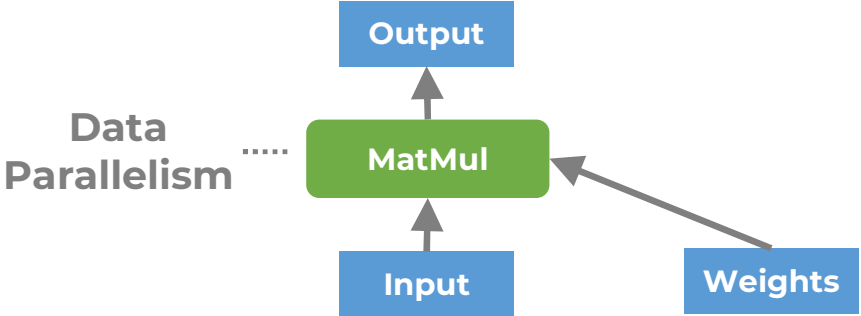


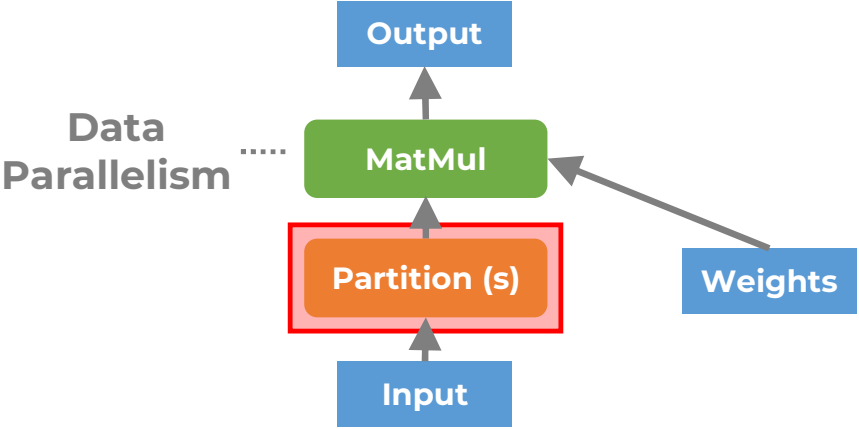
**Reduce**

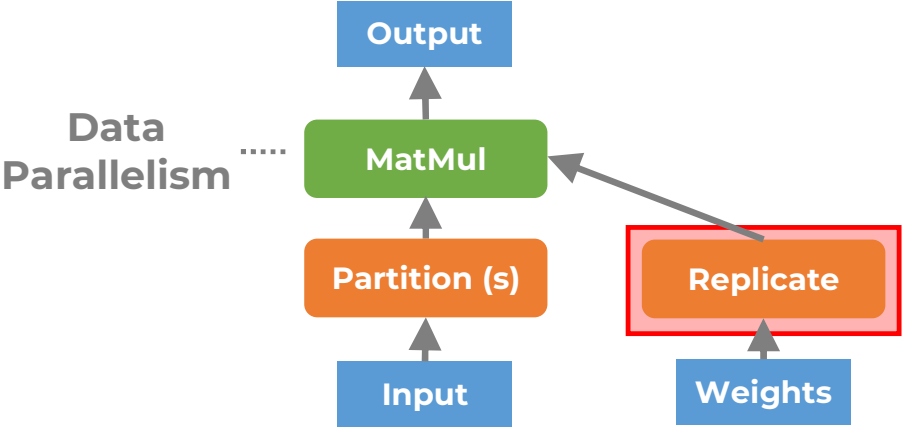
**Pipeline**

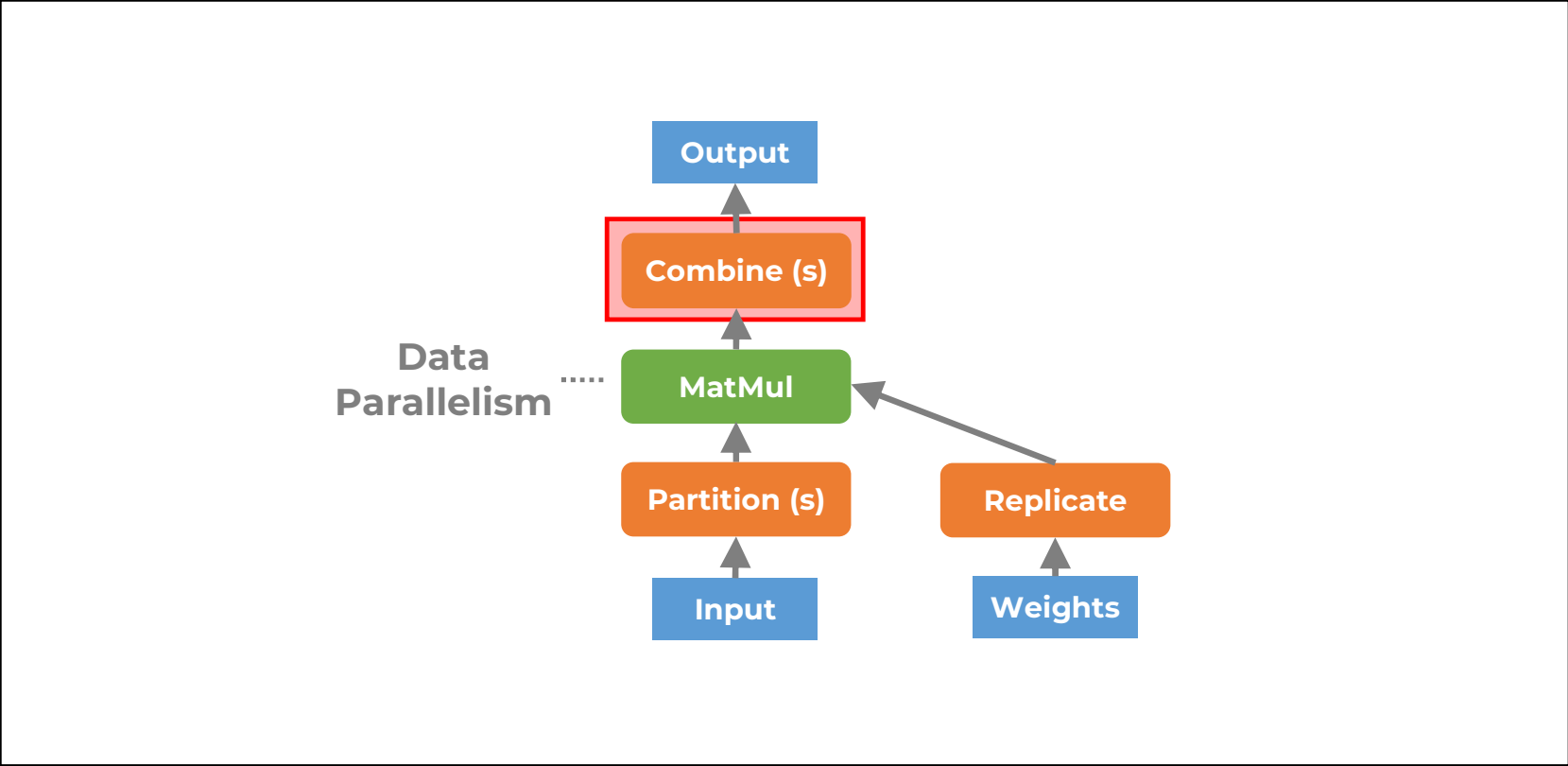


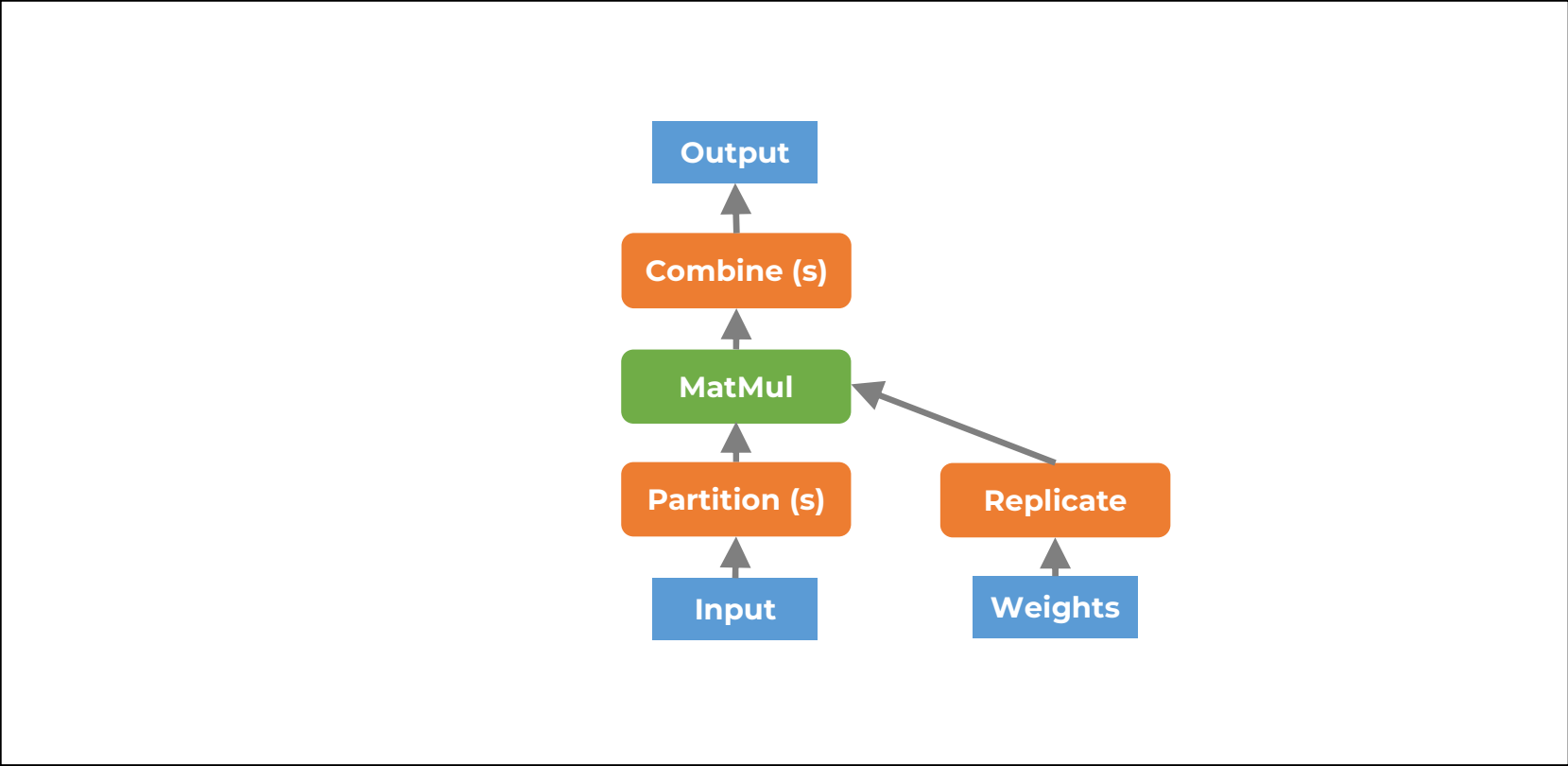
**Batch**

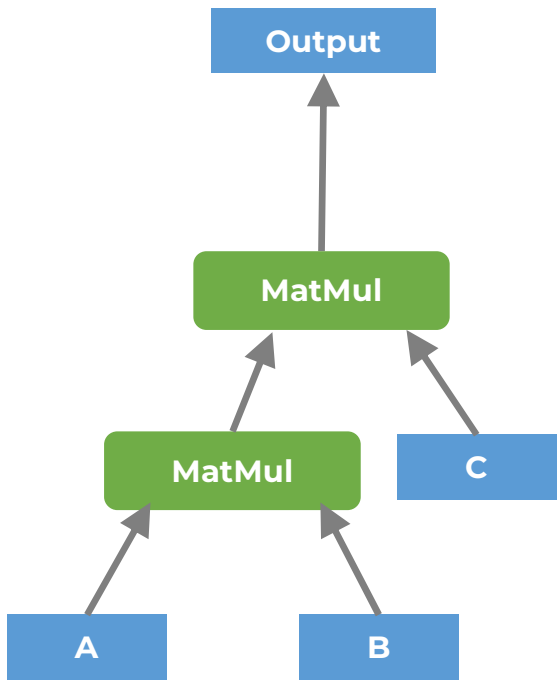


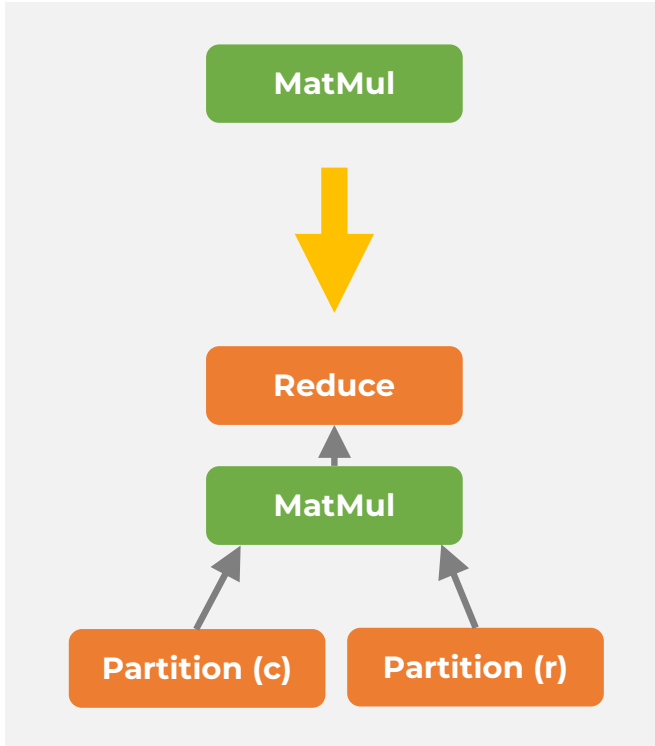
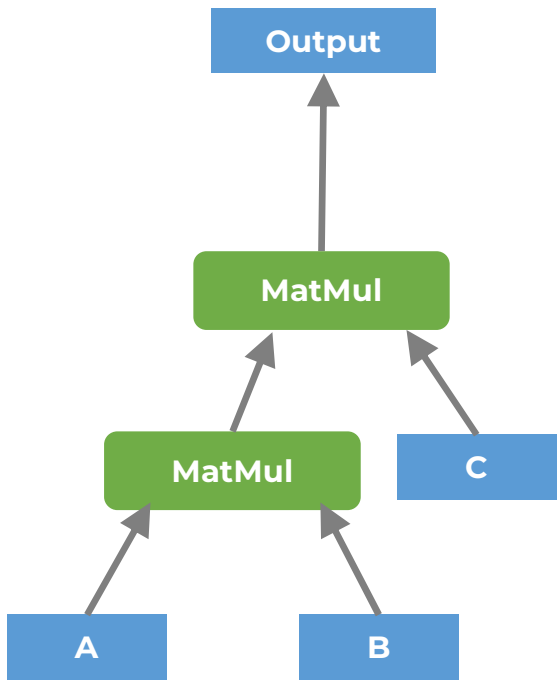




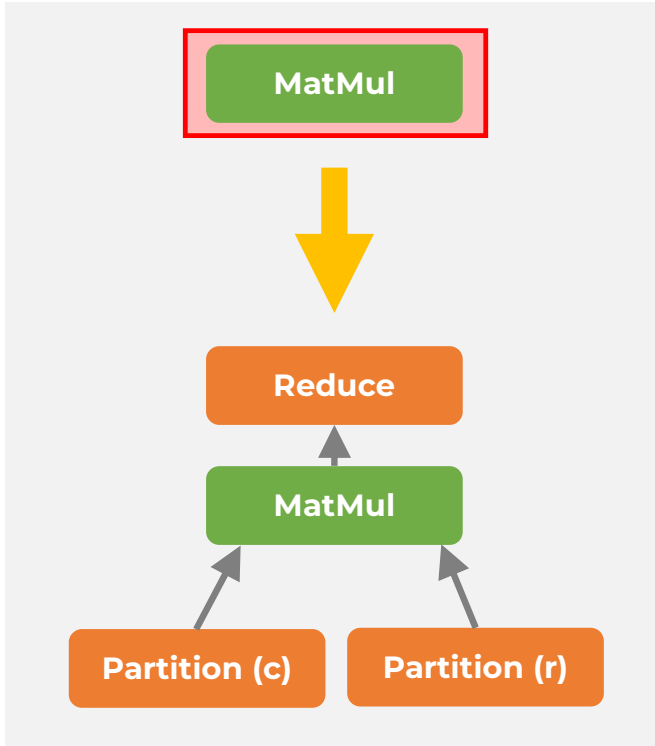
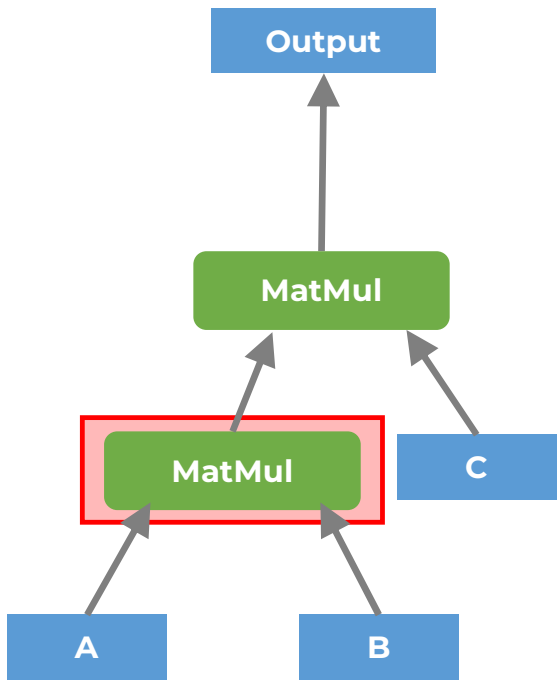


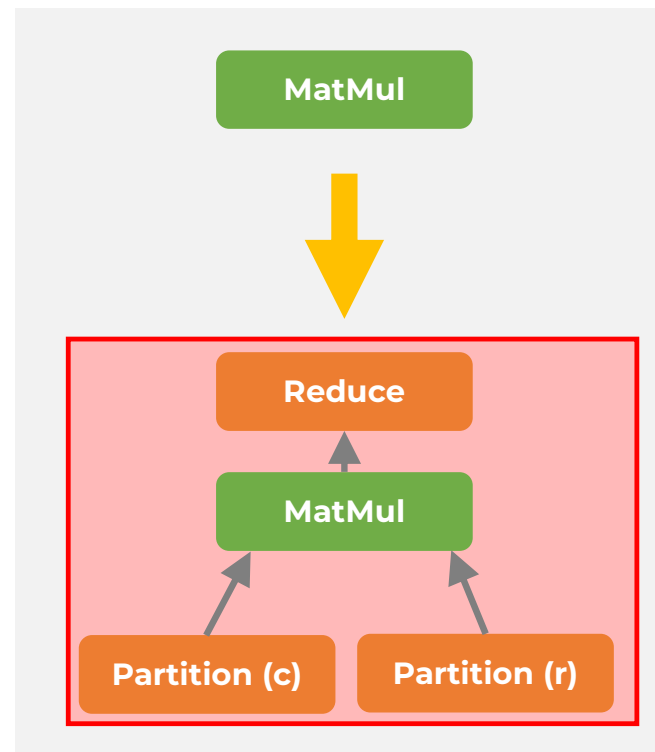
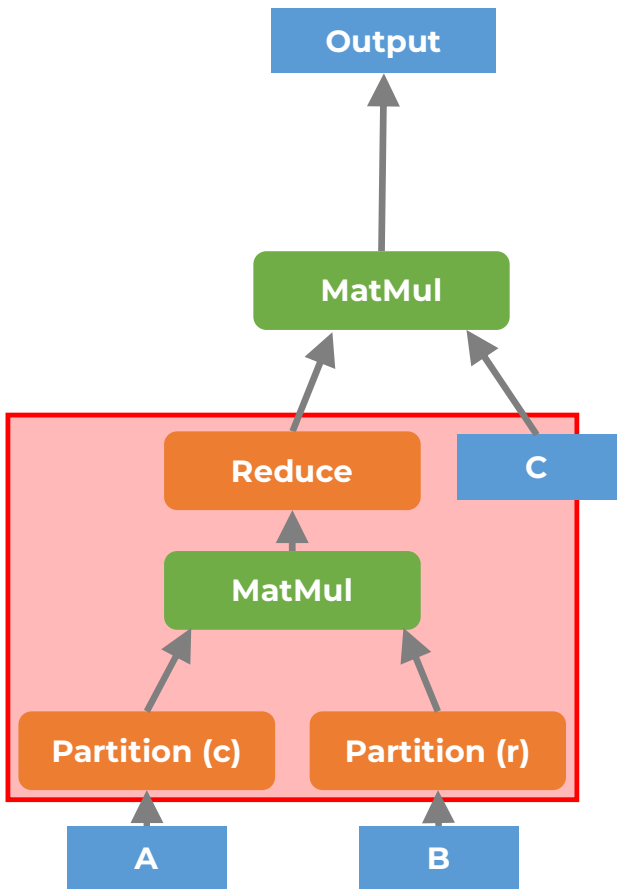


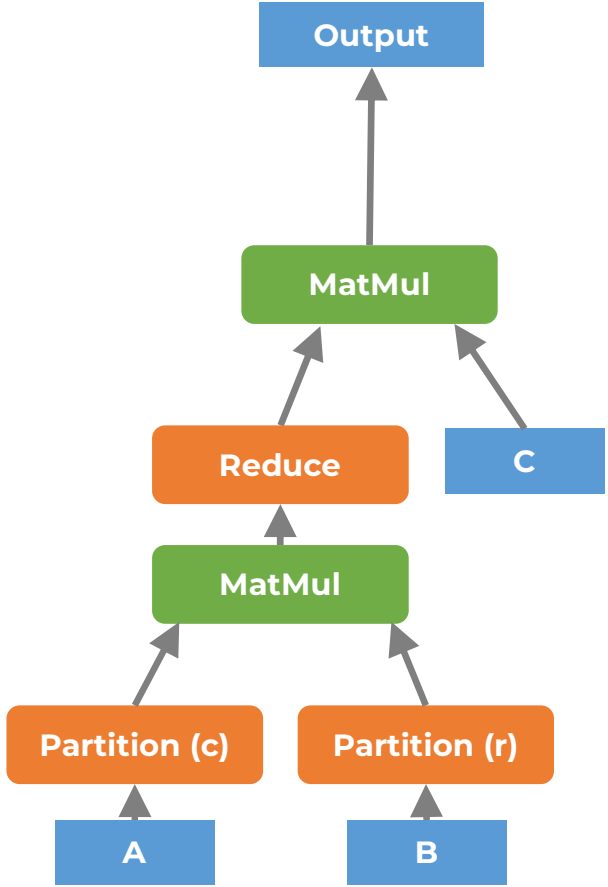




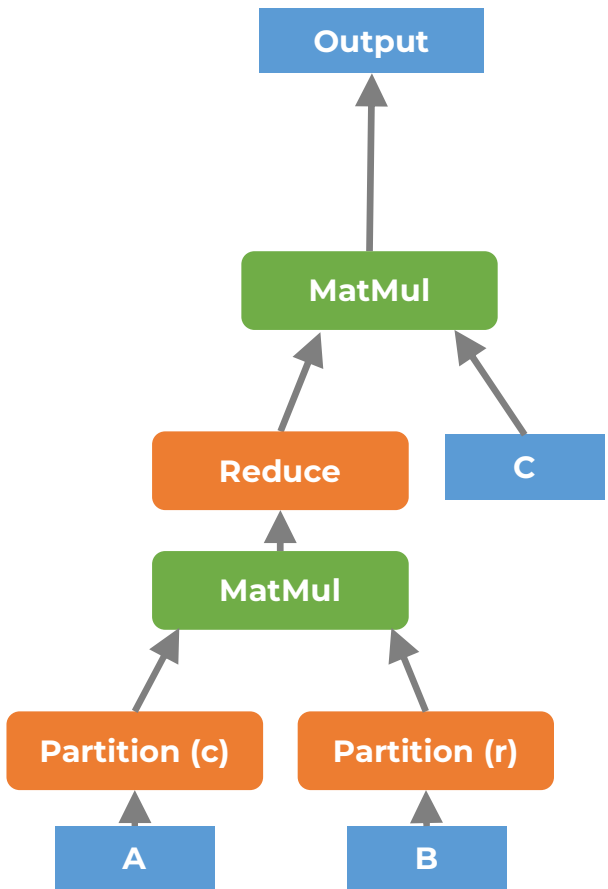




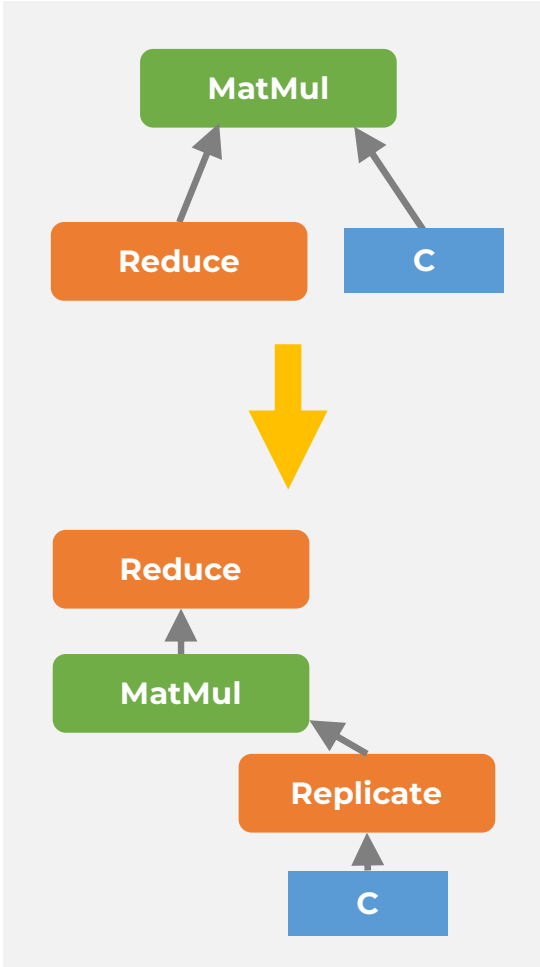




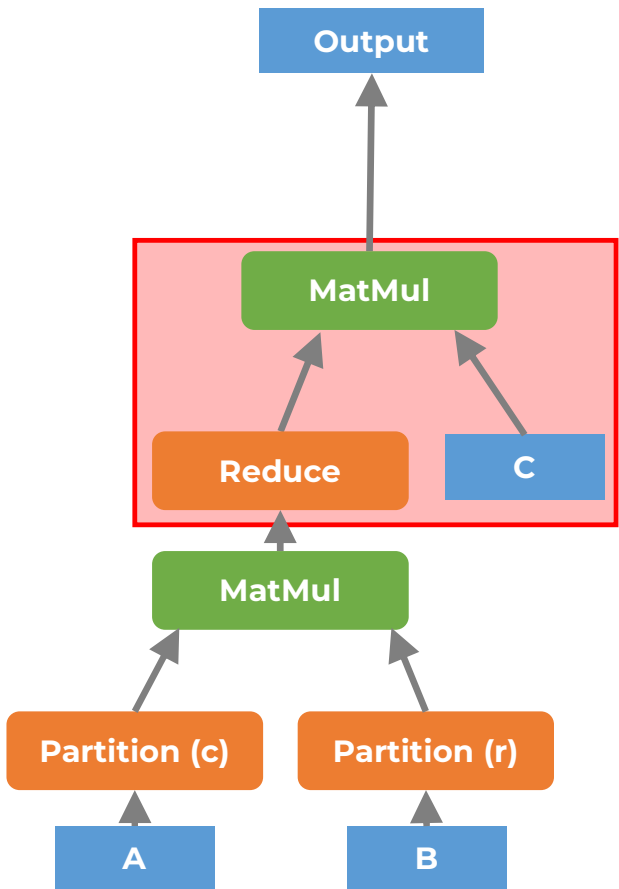
# PCG



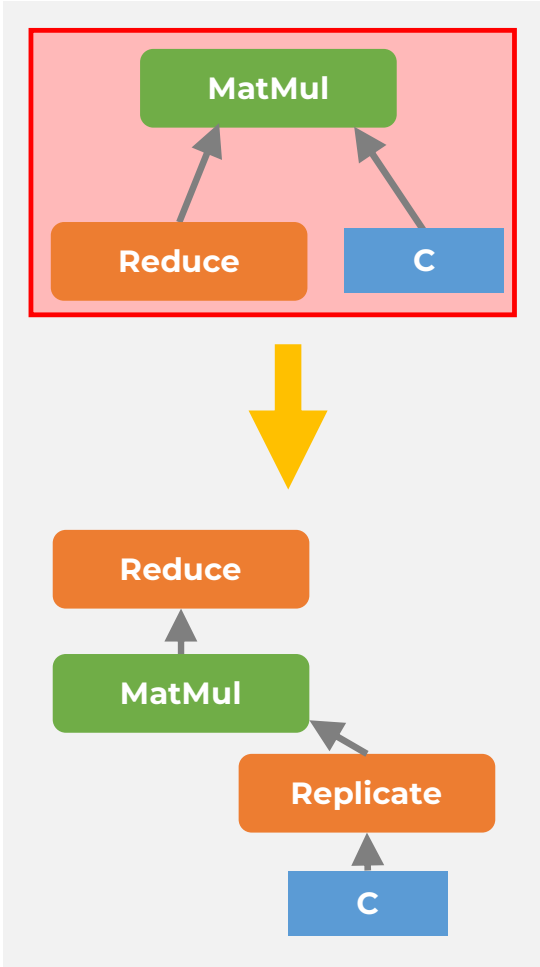
# Substitution

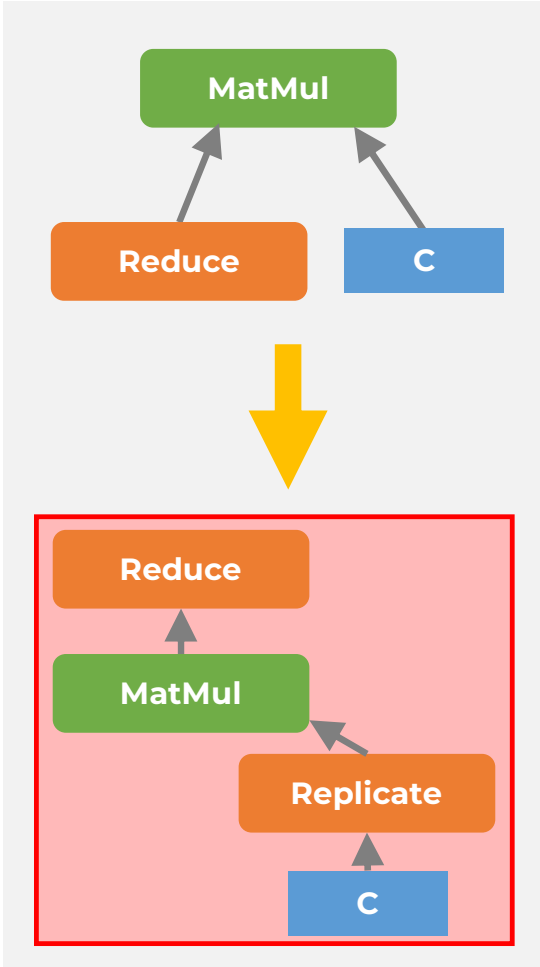
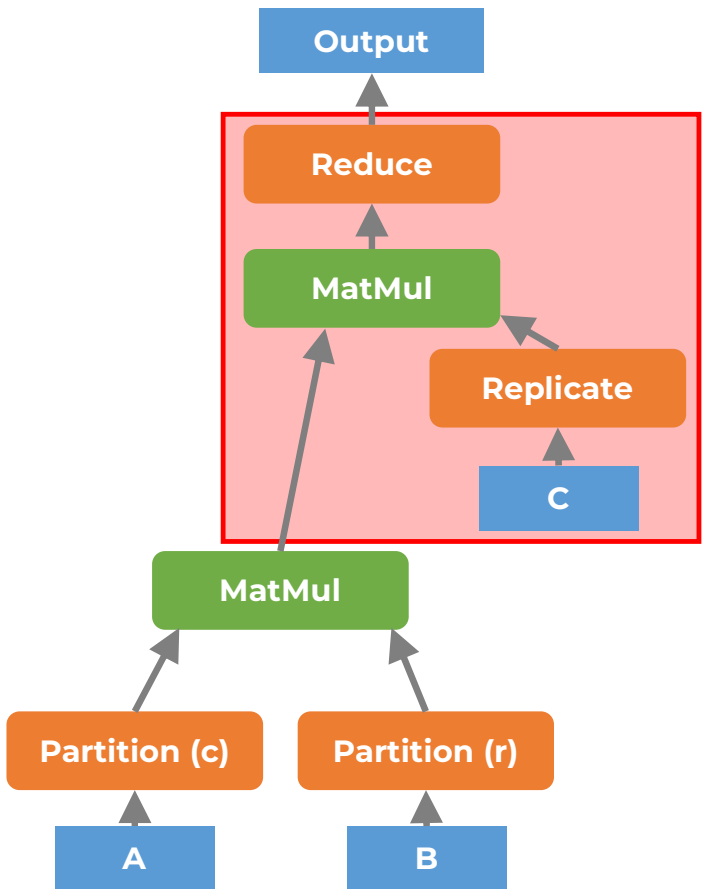


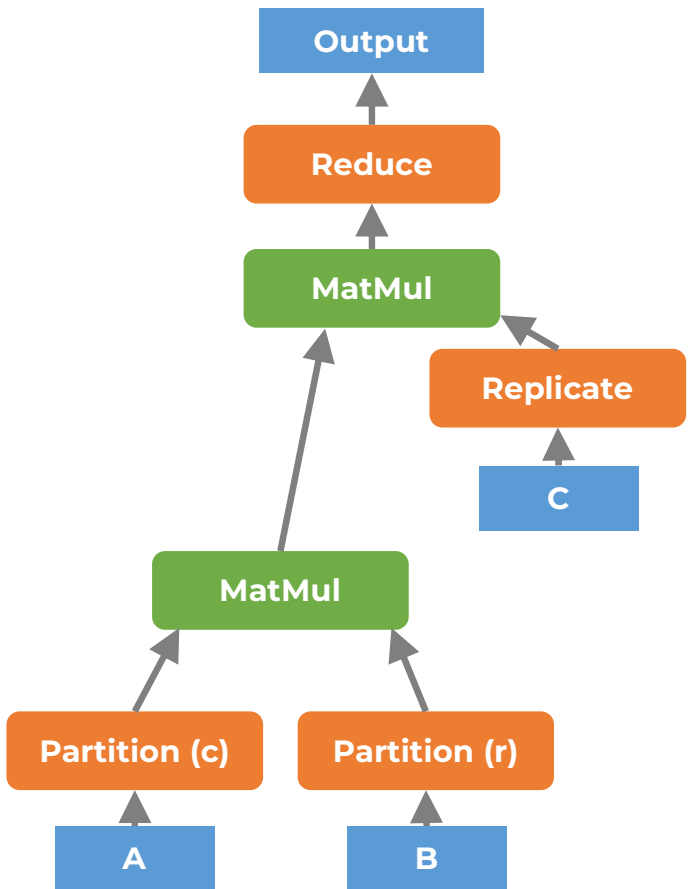
# PCG



# Substitution











# 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



# Hierarchical Search Algorithm

Algebraic Transformation

Parallelism Type

# Hierarchical Search Algorithm

Algebraic Transformation

Parallelism Type

Parallelism Degree

# 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

Backtracking  
Search

Parallelism Type



Parallelism Degree

Device Mapping



Dynamic  
Programming

# Evaluation

# Models

BERT-Large

(Language Modeling)

Candle-UNO

(Precision Medicine)

MLP

(Regression)

DLRM

XDL

(Recommendation)

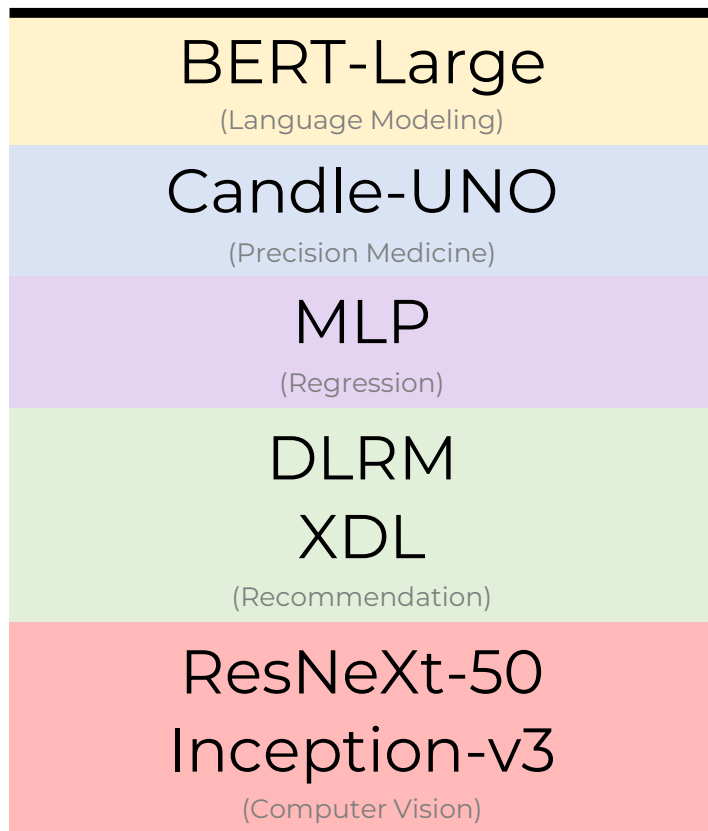
ResNeXt-50

Inception-v3

(Computer Vision)

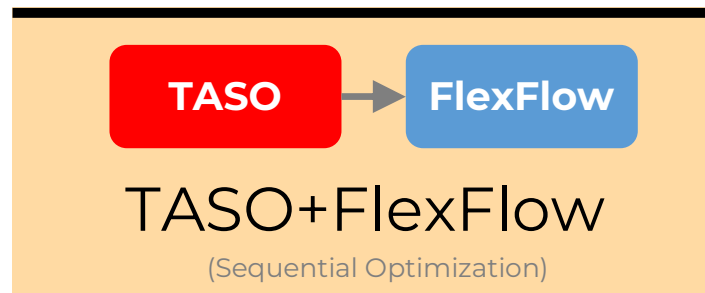


## Models

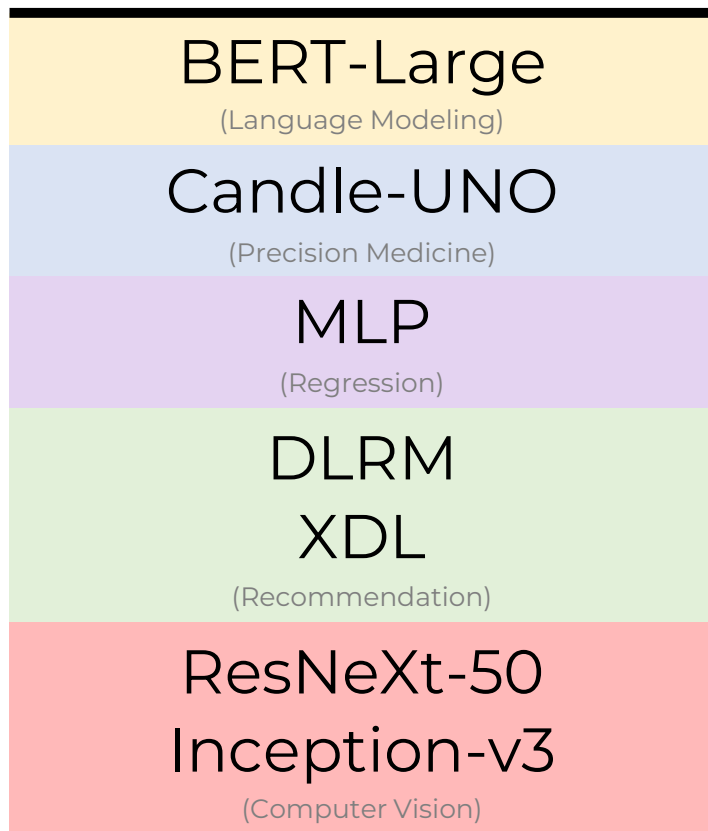


×

## Baselines

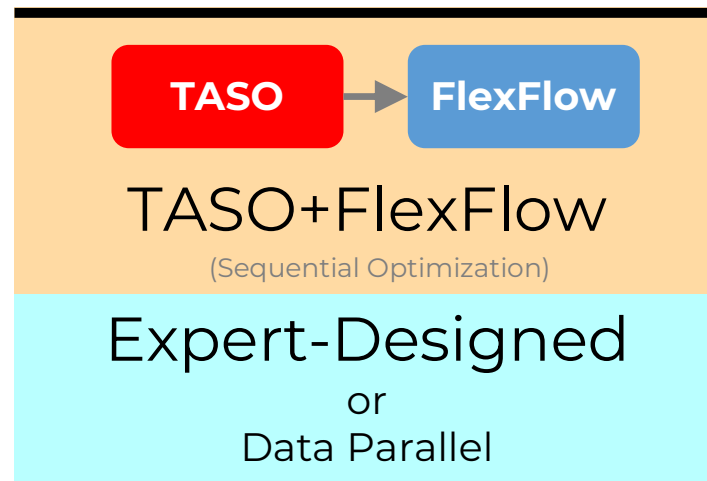


## Models

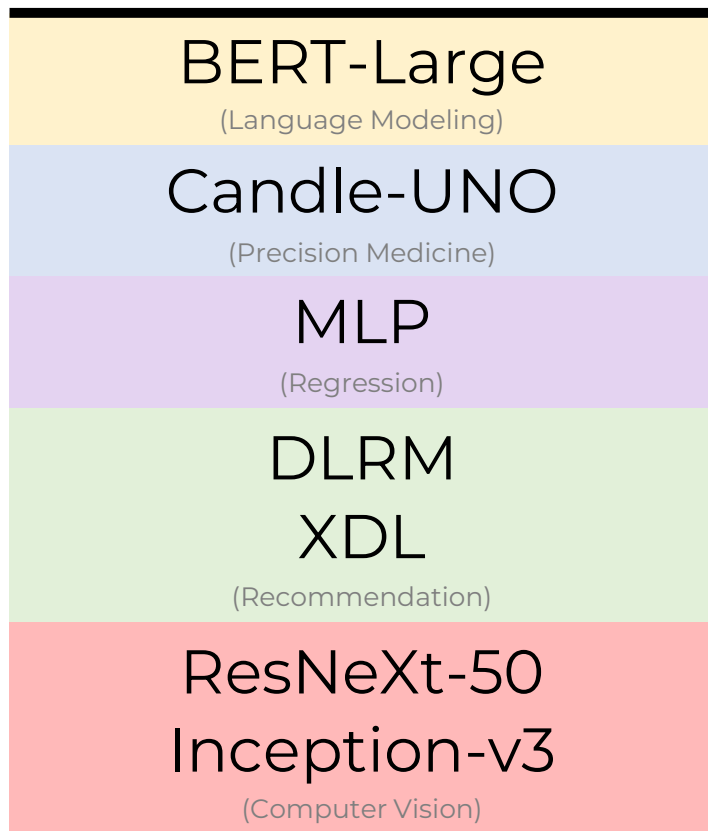


×

## Baselines

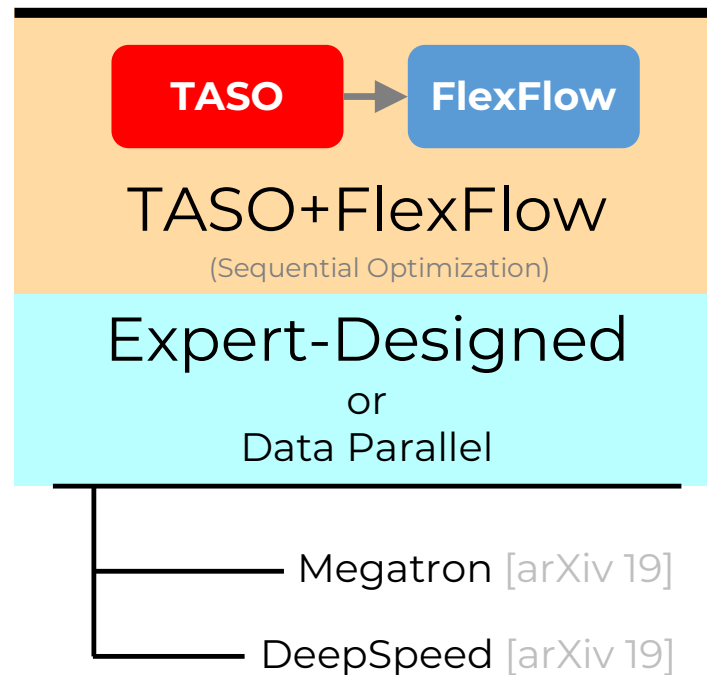


## Models

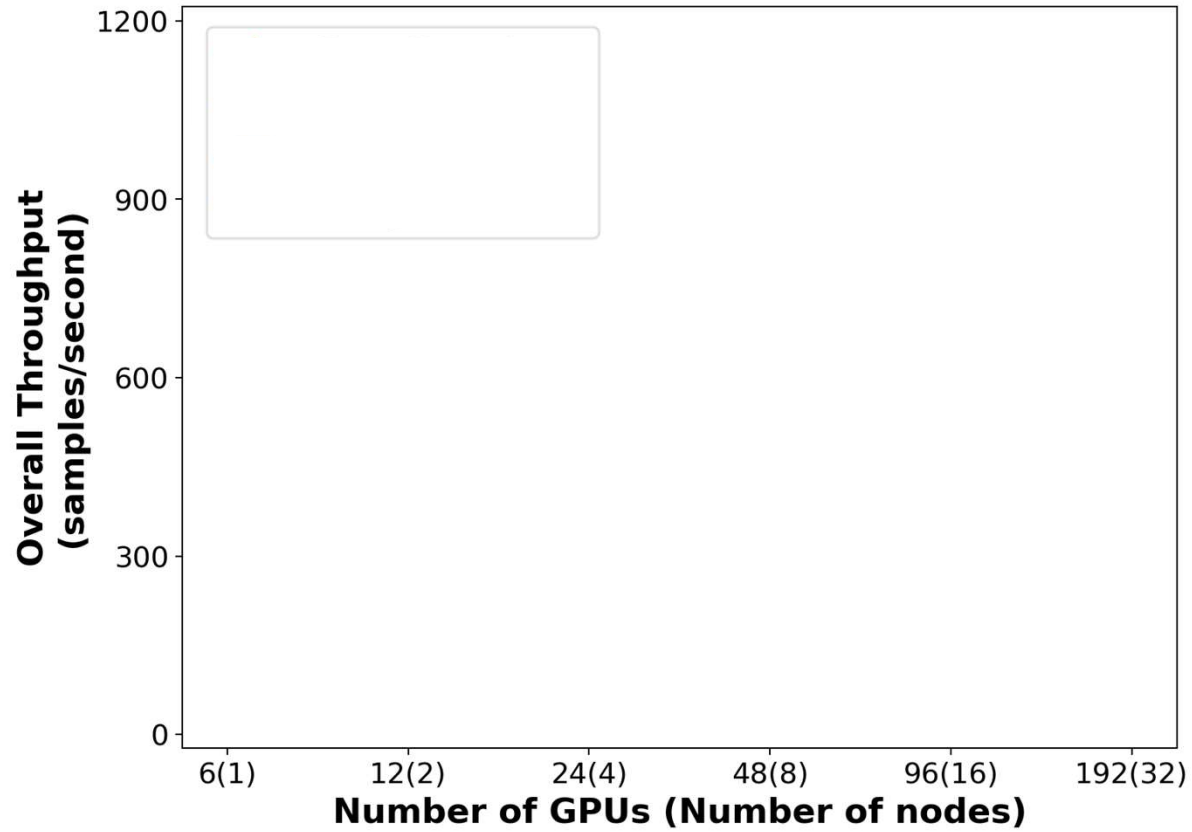


×

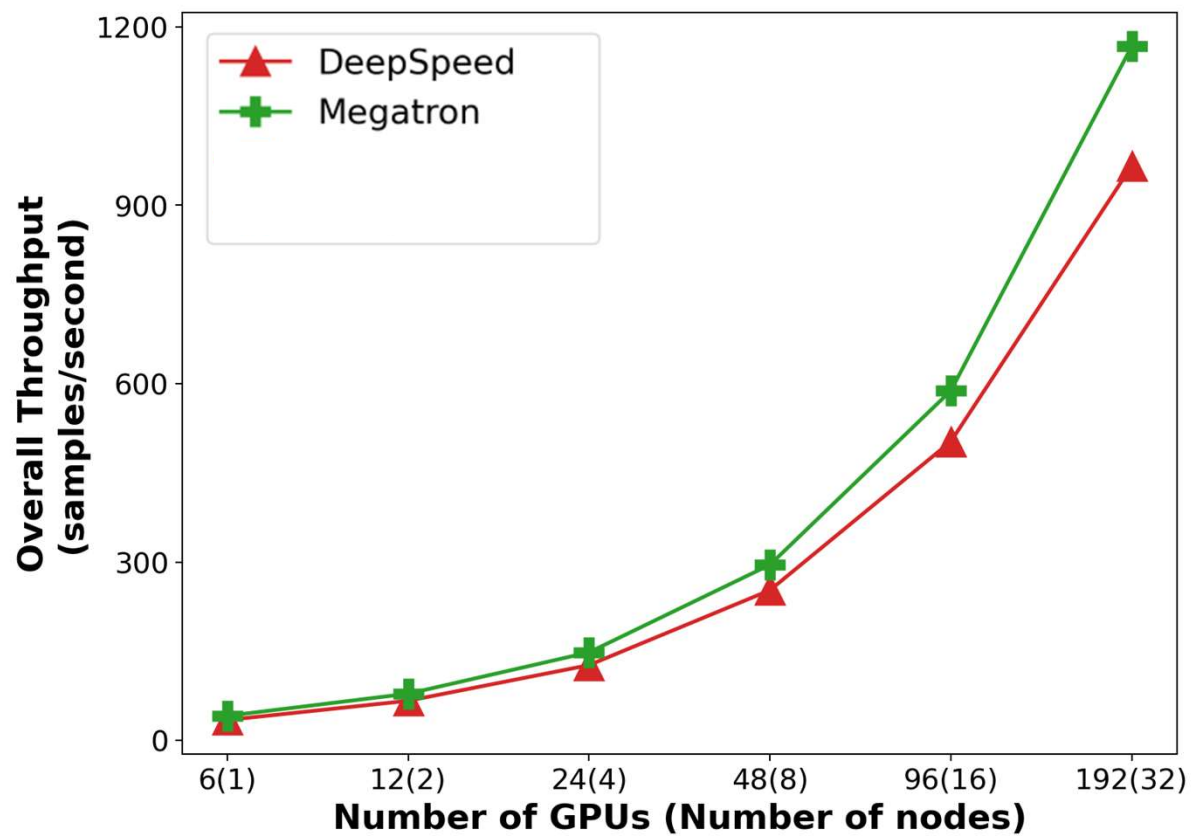
## Baselines



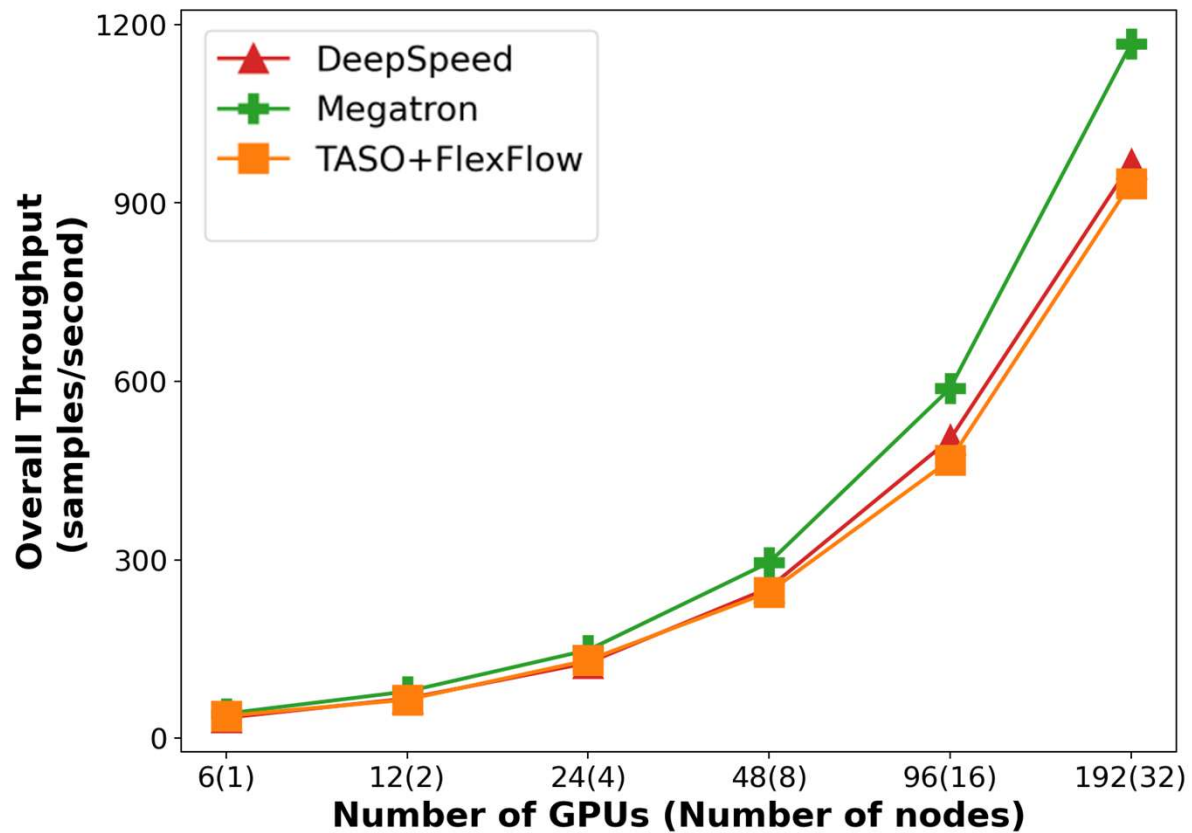
# BERT-Large



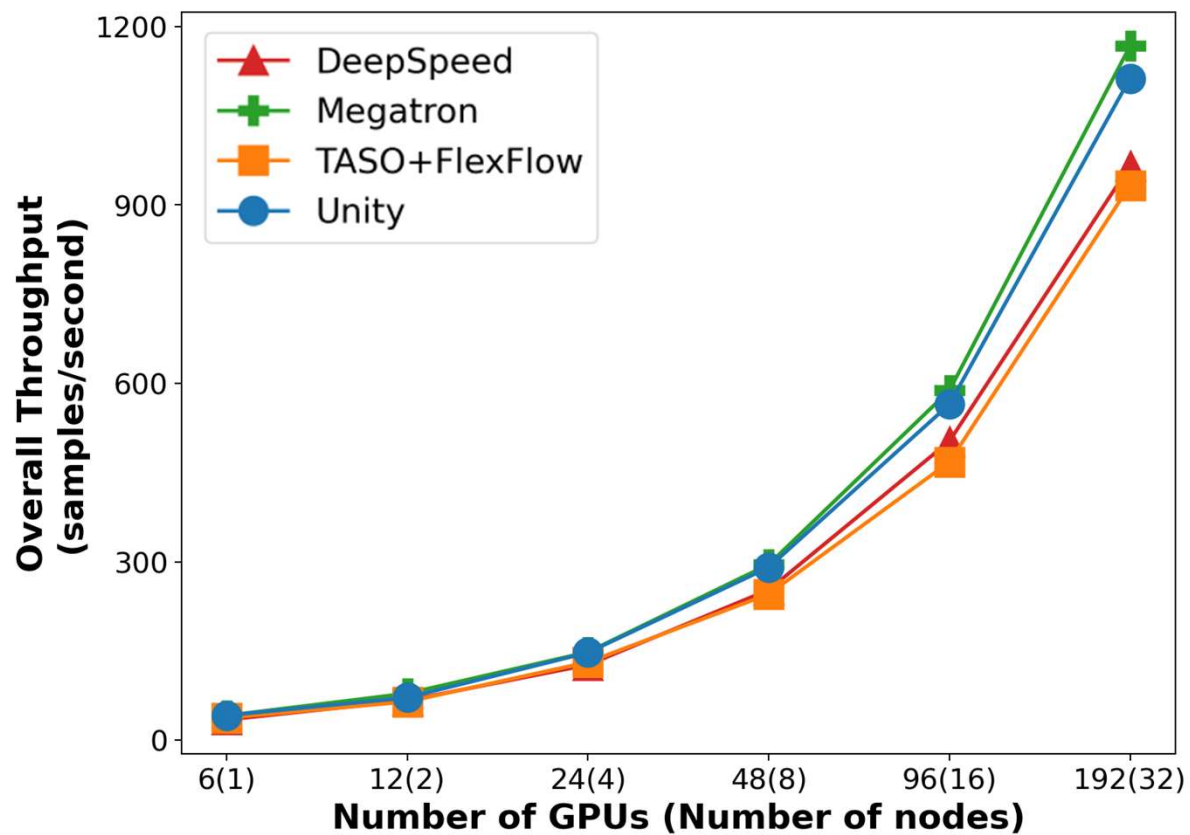
# BERT-Large



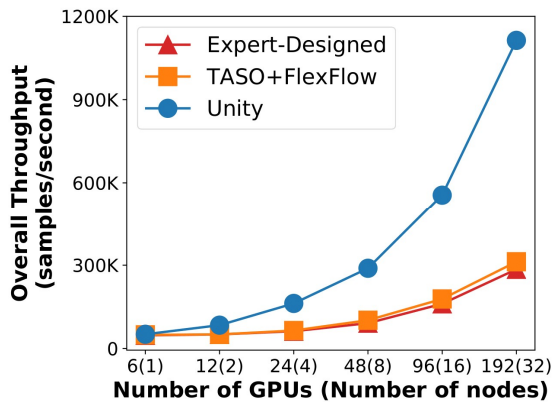
# BERT-Large



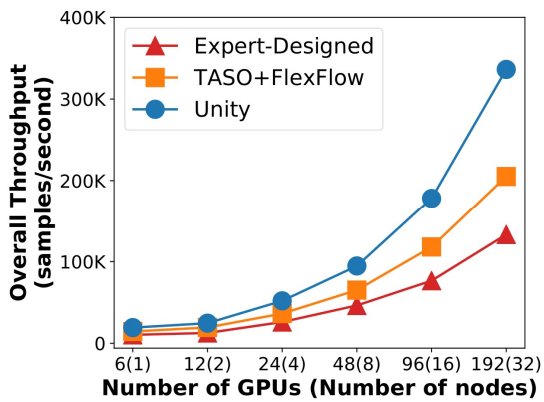
# BERT-Large



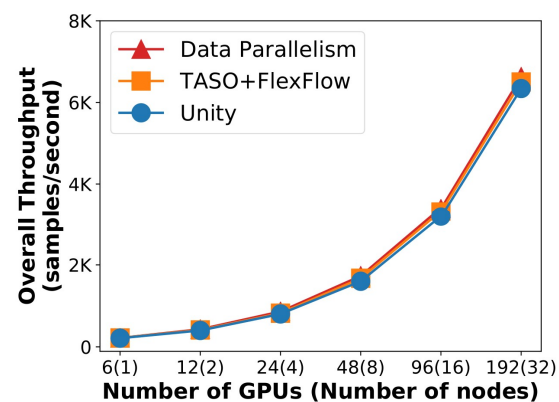
### DLRM



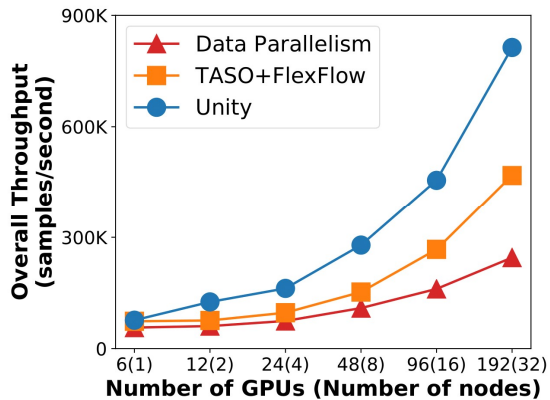
### CANDLE-Uono



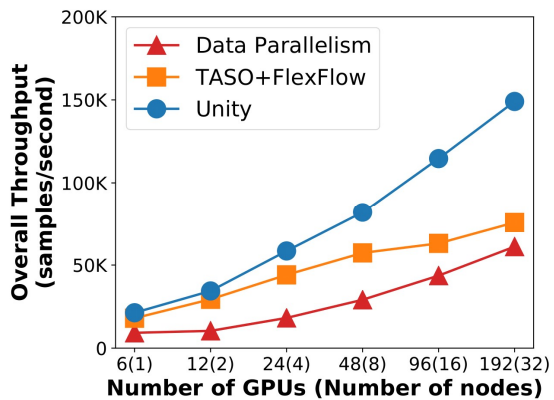
### ResNeXt-50



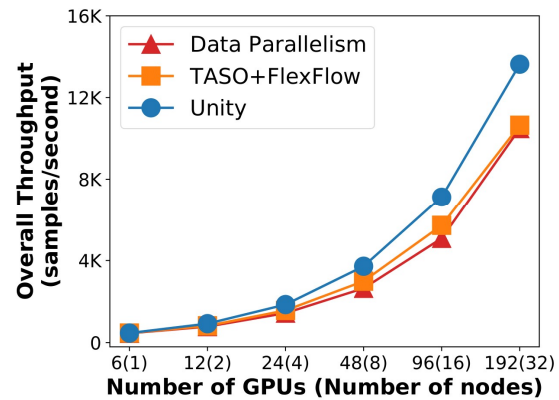
### XDL



### MLP

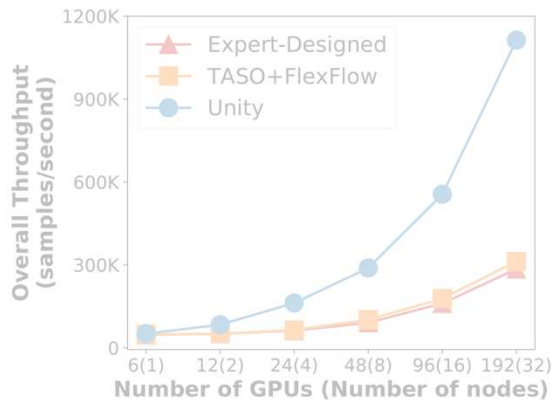


### Inception-v3

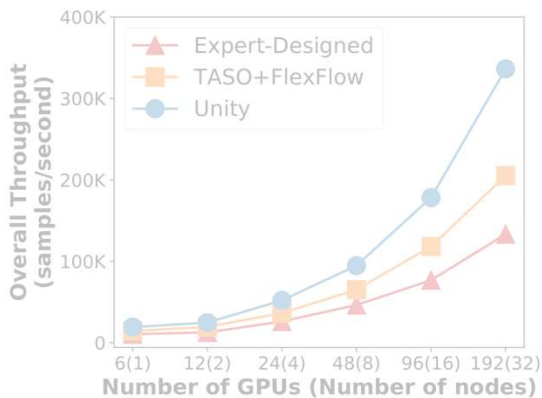




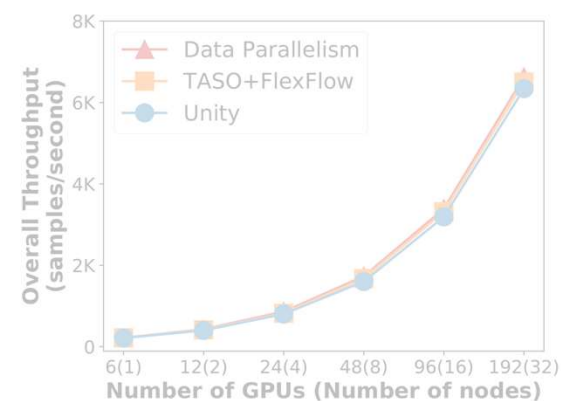
### DLRM



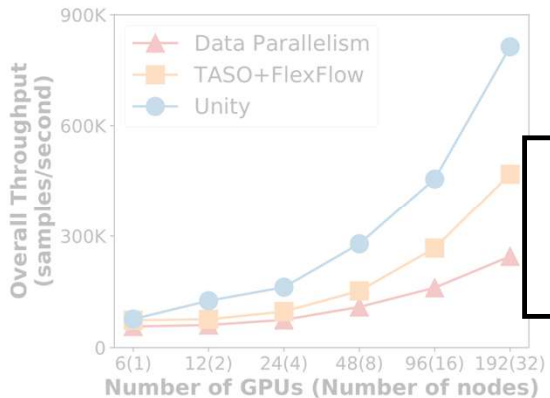
### CANDLE-Uono



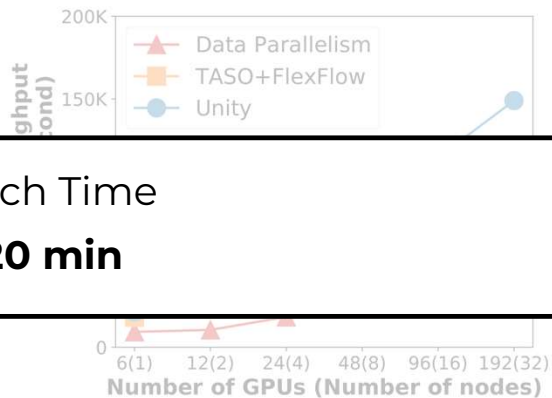
### ResNeXt-50



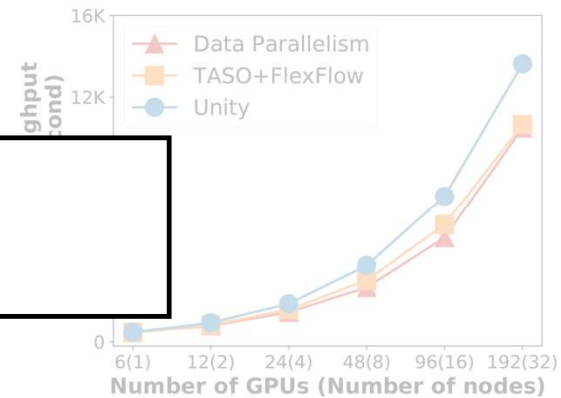
### XDL



### MLP

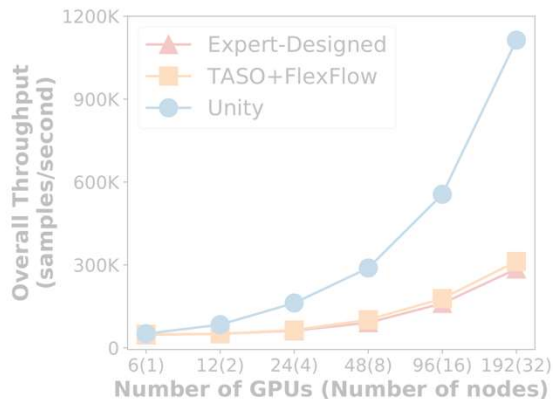


### Inception-v3

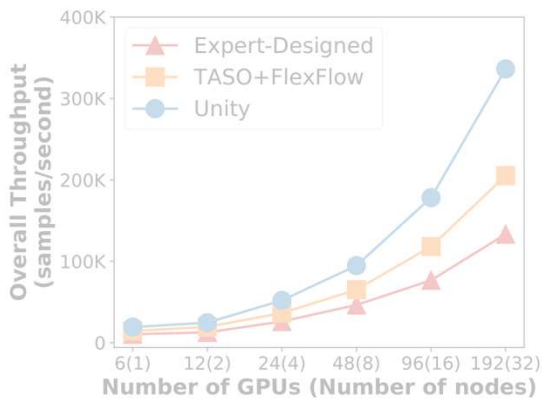


Search Time  
**< 20 min**

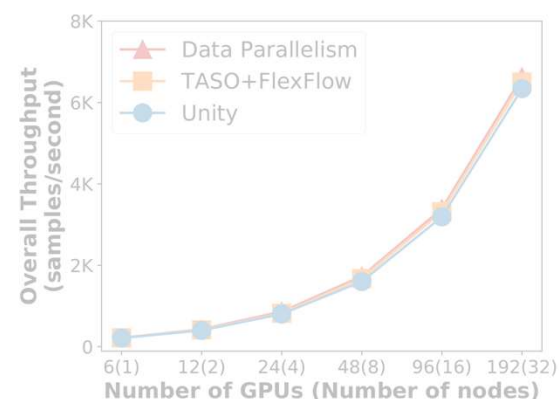
### DLRM



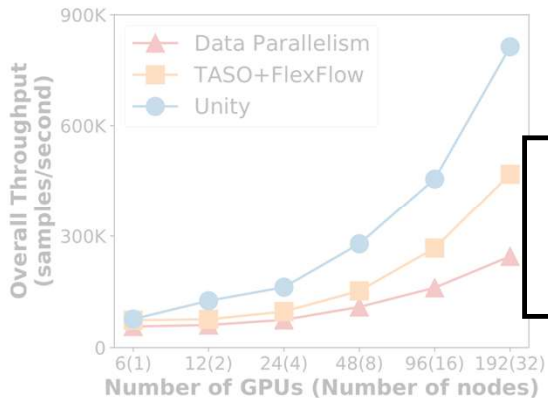
### CANDLE-Uno



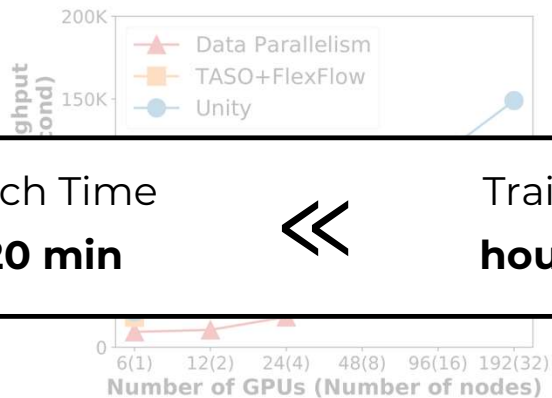
### ResNeXt-50



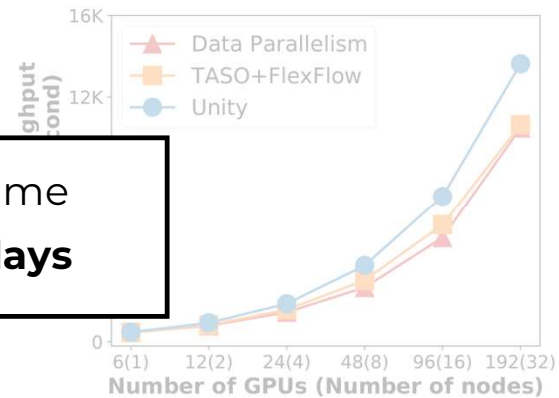
### XDL



### MLP



### Inception-v3



Search Time < 20 min << Training Time hours or days



# 1

## Joint Optimization

1

Joint  
Optimization

2

Unified Graph  
Representation  
(PCG)

1

Joint  
Optimization

2

Unified Graph  
Representation  
(PCG)

3

Hierarchical  
Search  
Algorithm

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

A screenshot of the FlexFlow GitHub repository page. The top section shows a commit by kadinzhang: "Added tests for Linear operator in align/linear (#264)" from 19 days ago with 1,339 commits. Below this is a list of folders: circleci, align, bootcamp\_demo, cmake, conda, config, deps, and docker, each with a brief description and a timestamp. To the right, the repository description reads: "A distributed deep learning framework that supports flexible parallelization strategies." It also lists metadata: Apache-2.0 license, 458 stars, 21 watchers, and 100 forks. A release section shows "Release 21.09 (September 30th ...)" with a "Latest" button. At the bottom, another commit by kadinzhang is visible: "Merging the public branch to the master branch (#74)" from 16 months ago.

 Keras  PyTorch  ONNX



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

# Questions?



Keras

PyTorch



ONNX



CMU

