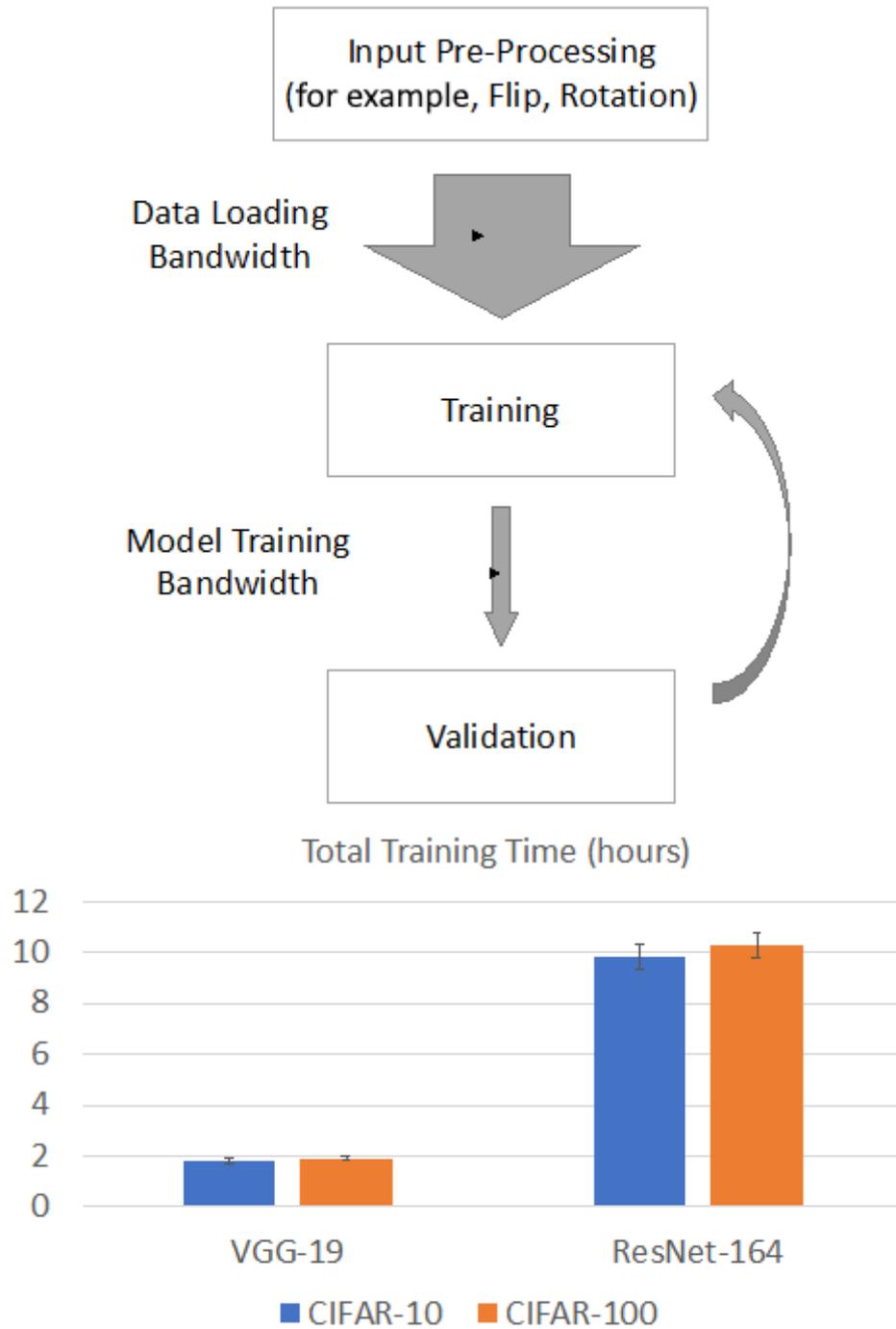
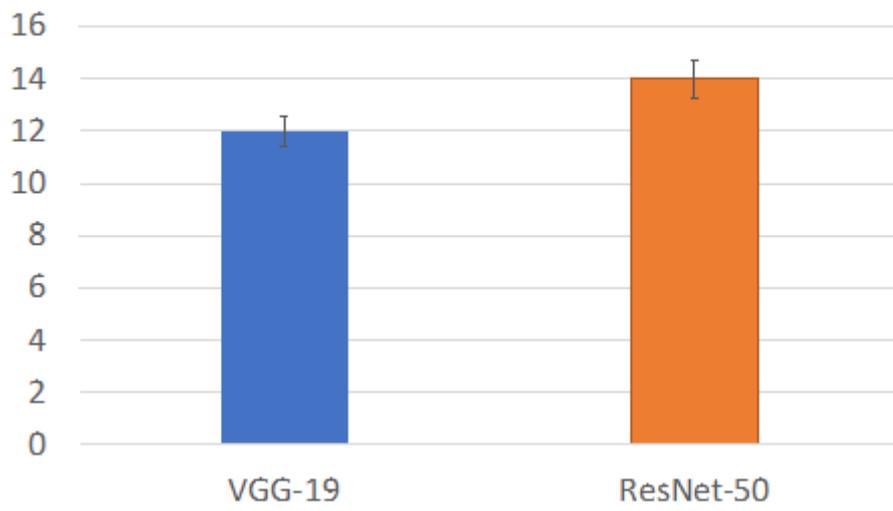


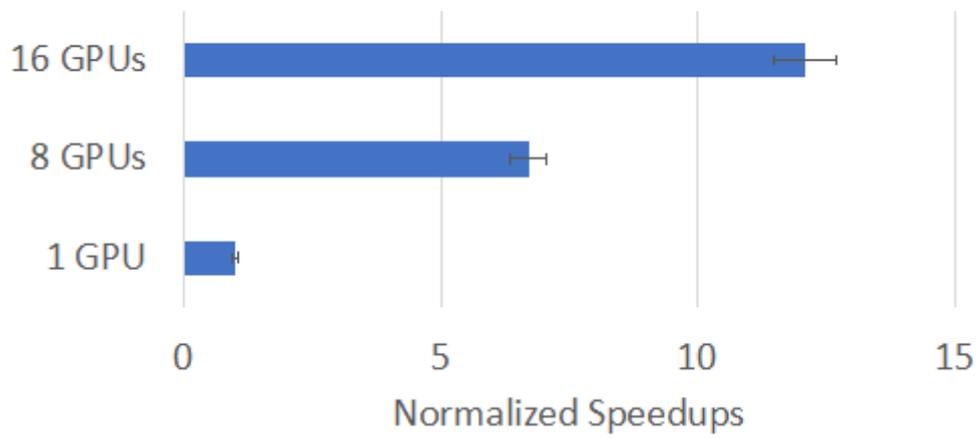
Chapter 1: Splitting Input Data

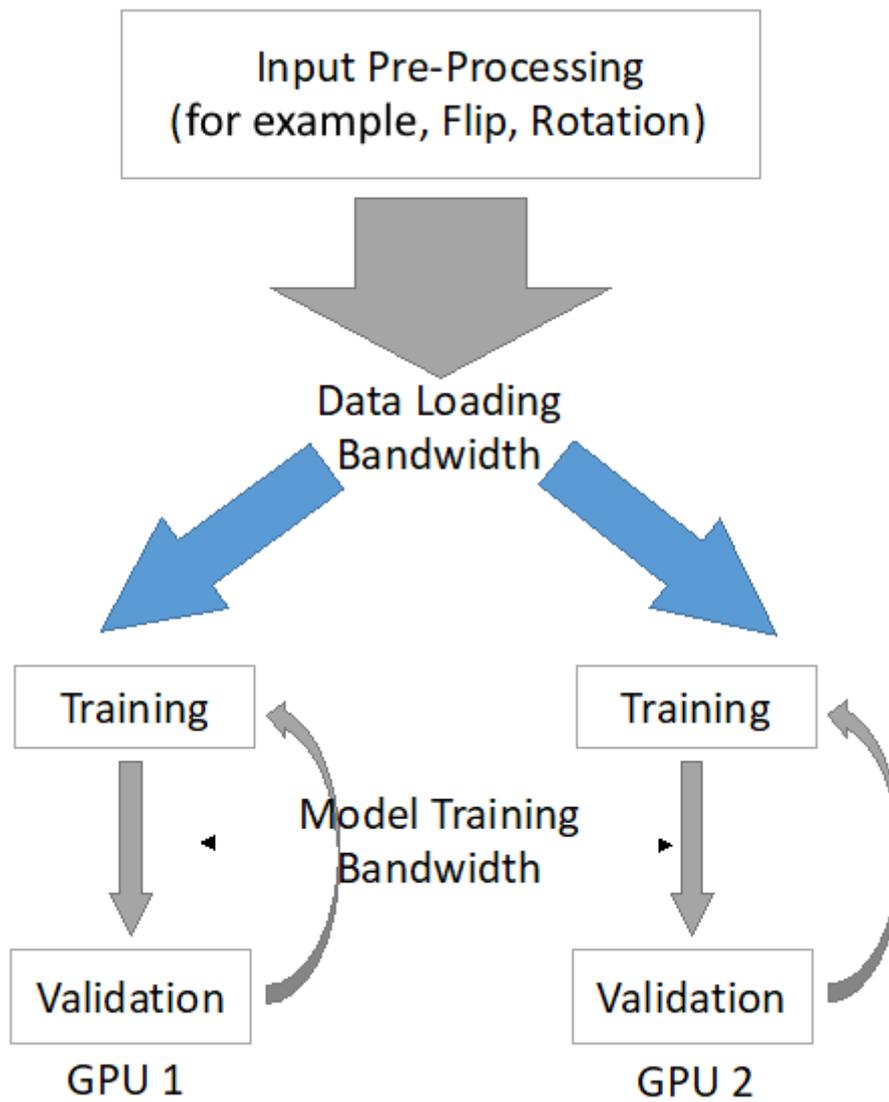


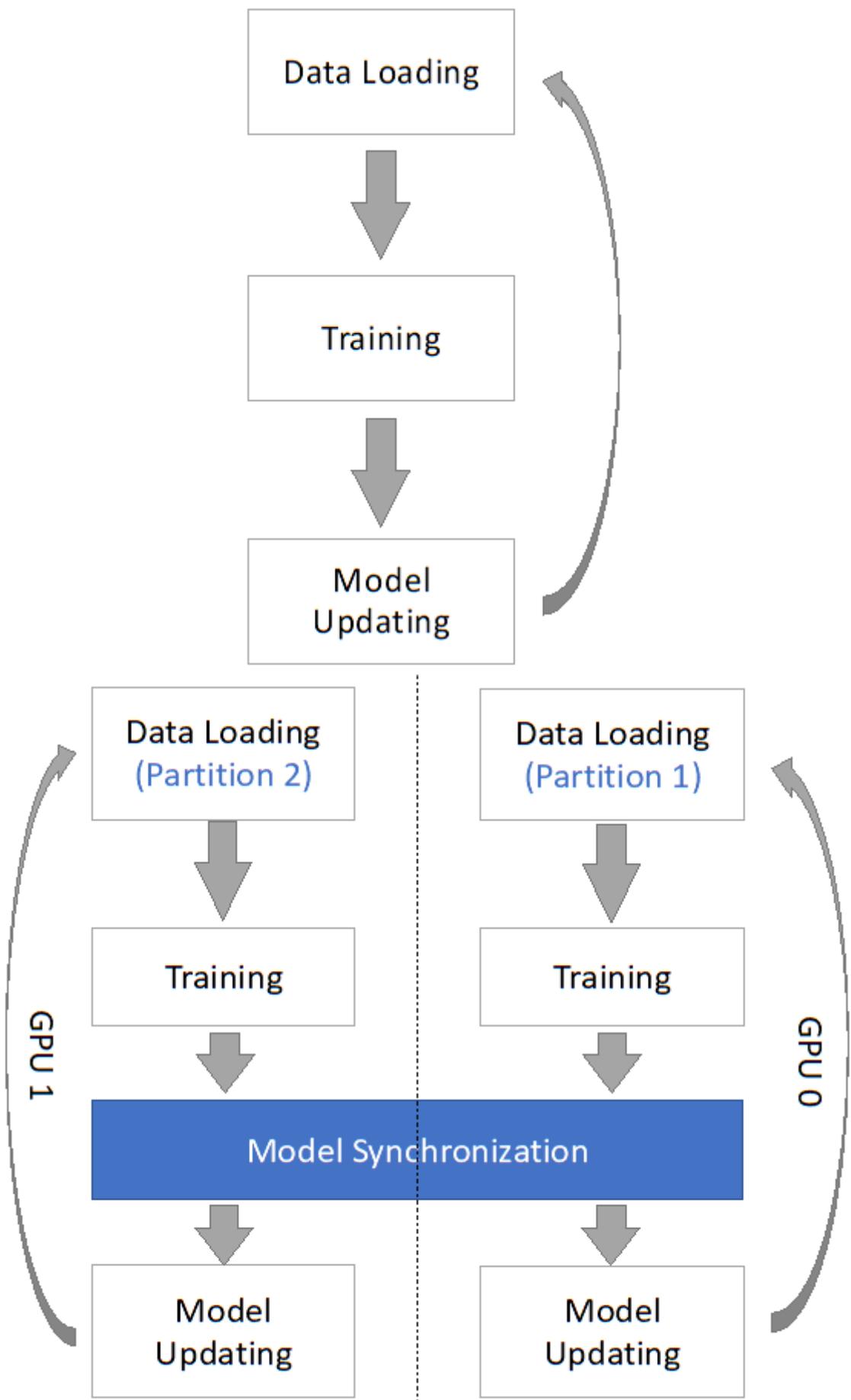
ImageNet-1K Training Time (Days)

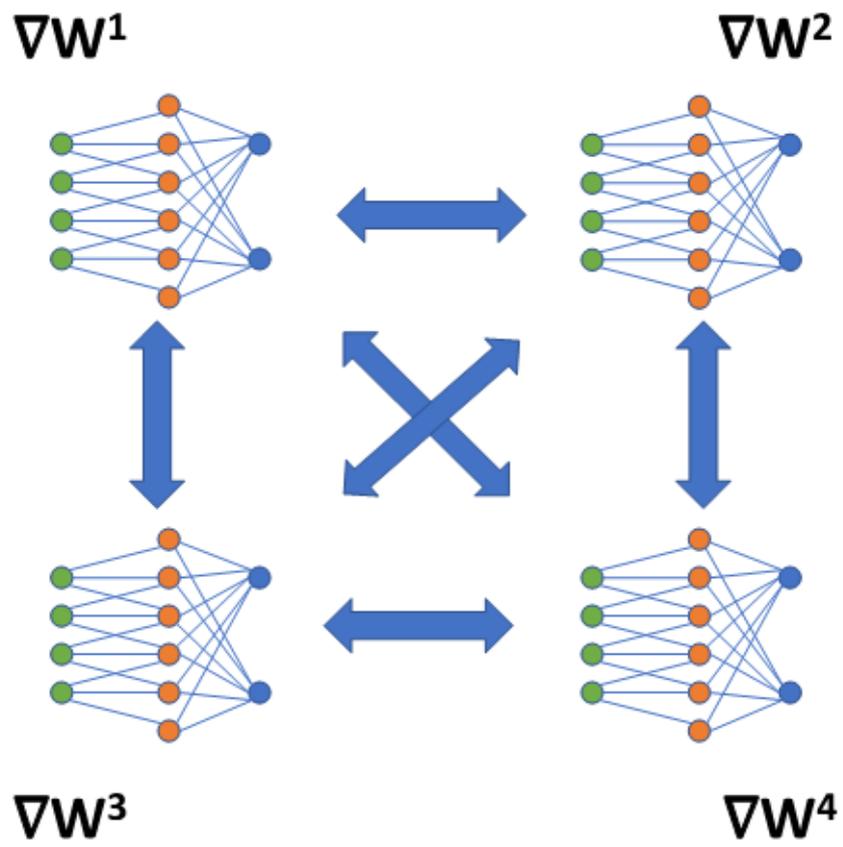


ResNet-50 Training on ImageNet-1K Dataset





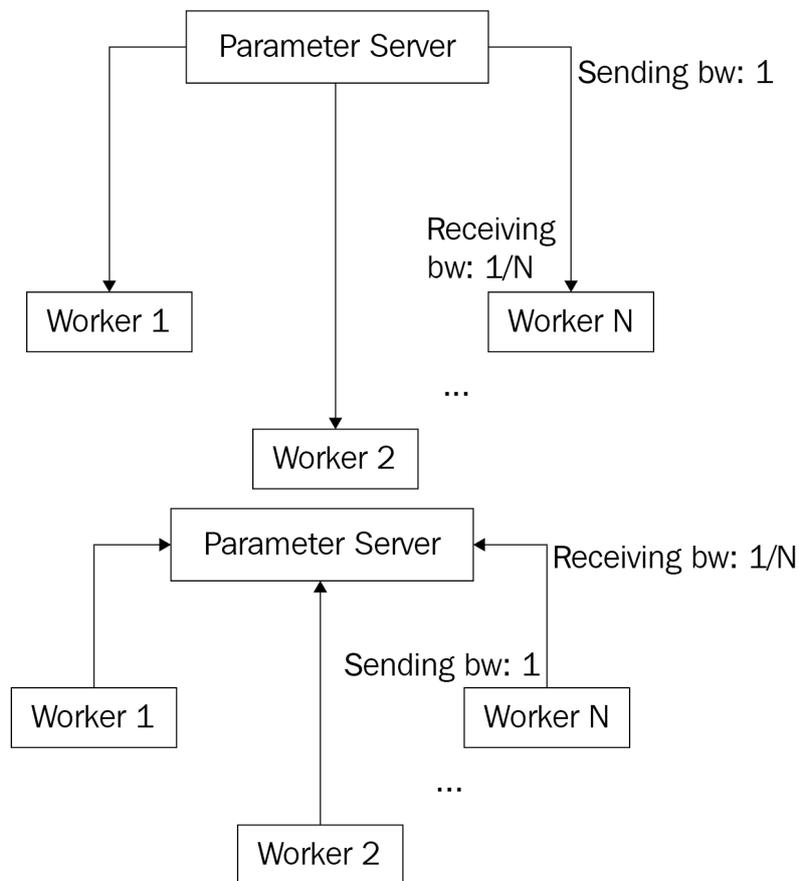
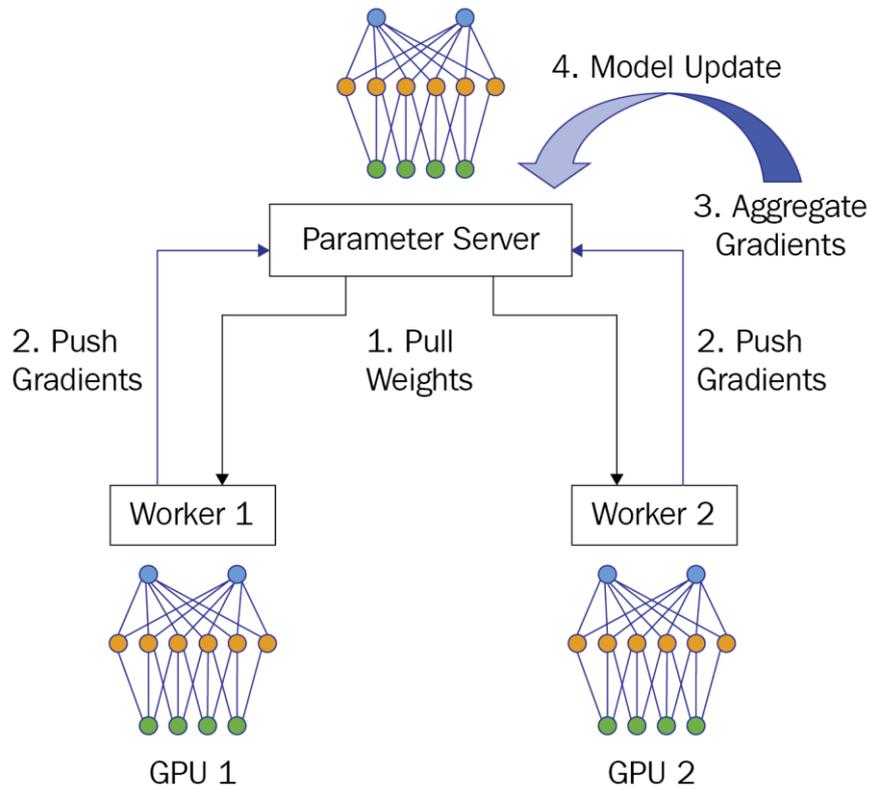


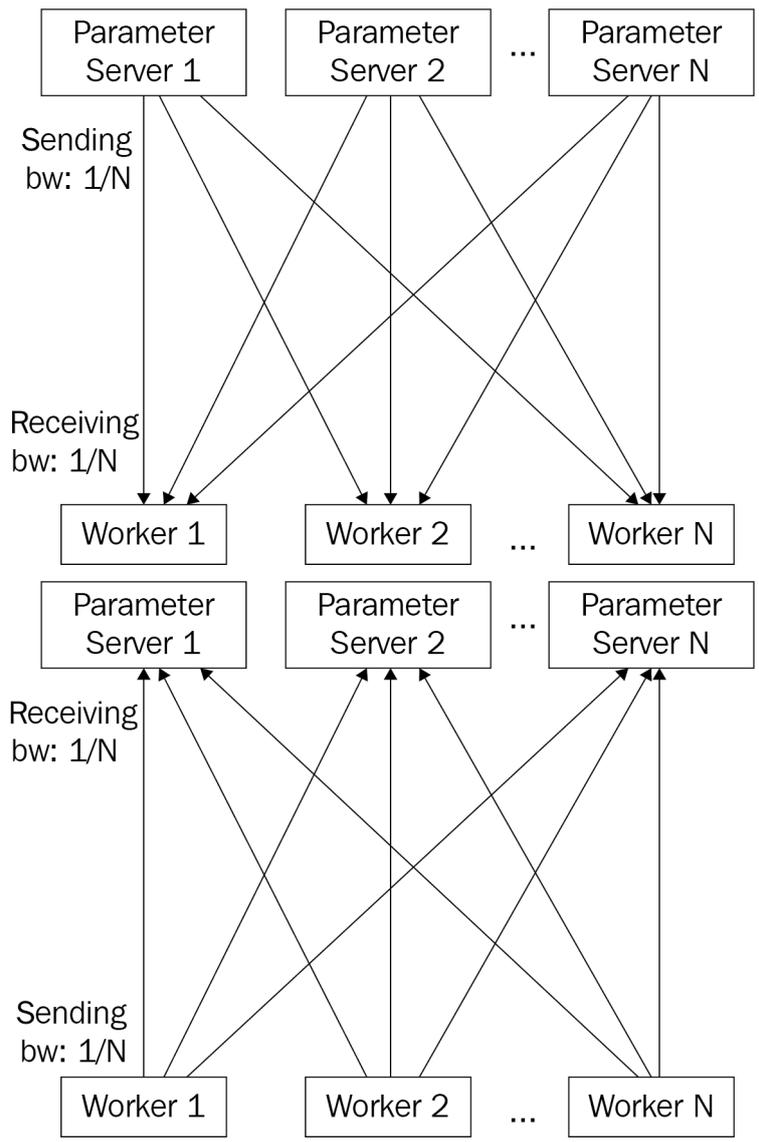


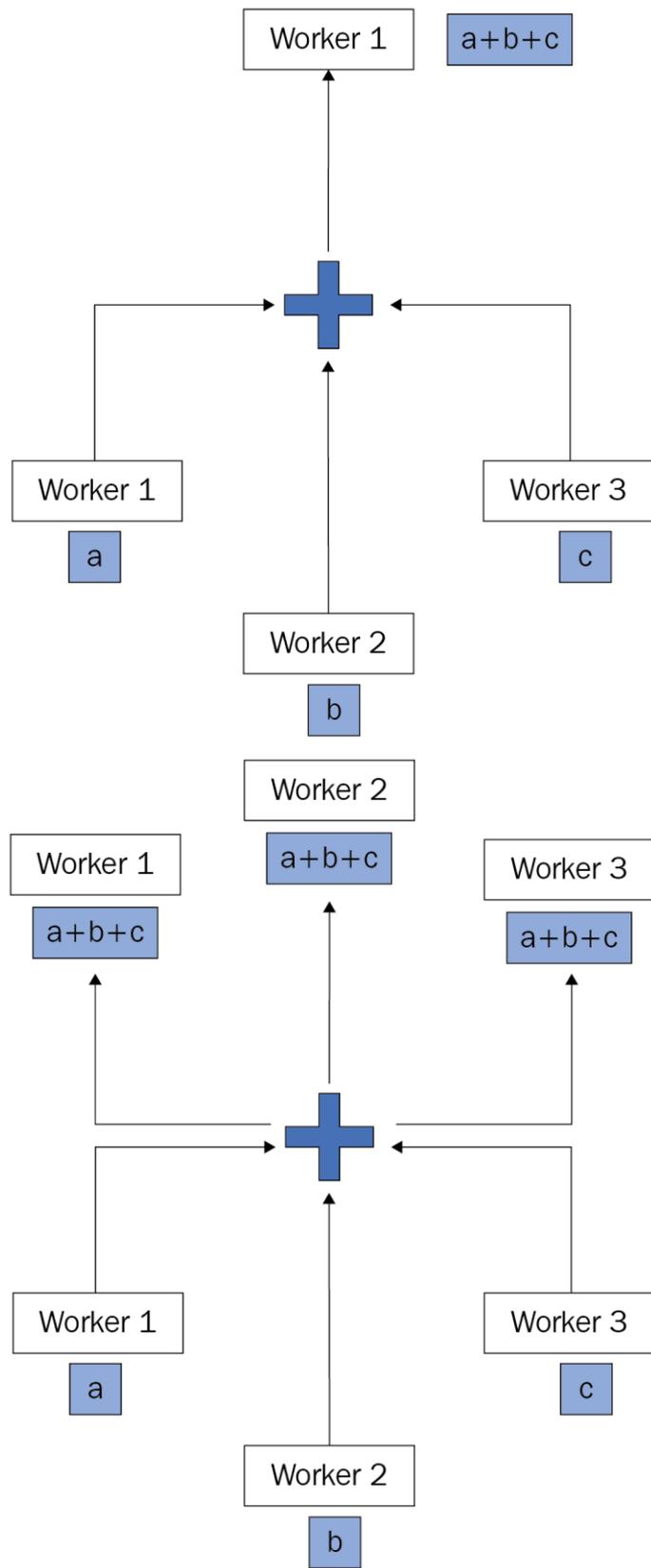
Synchronization

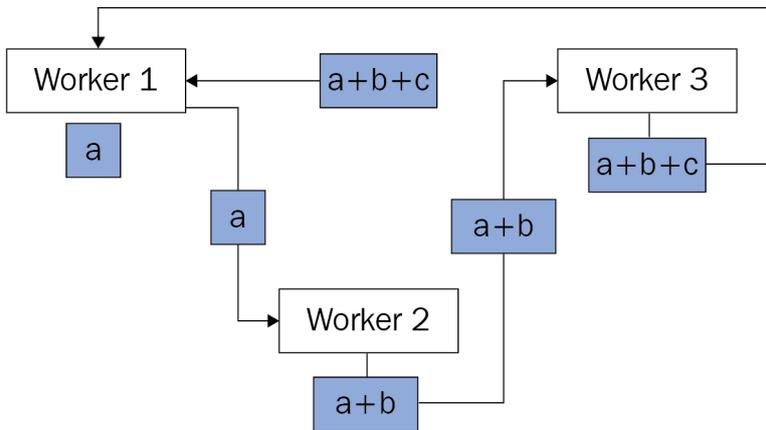
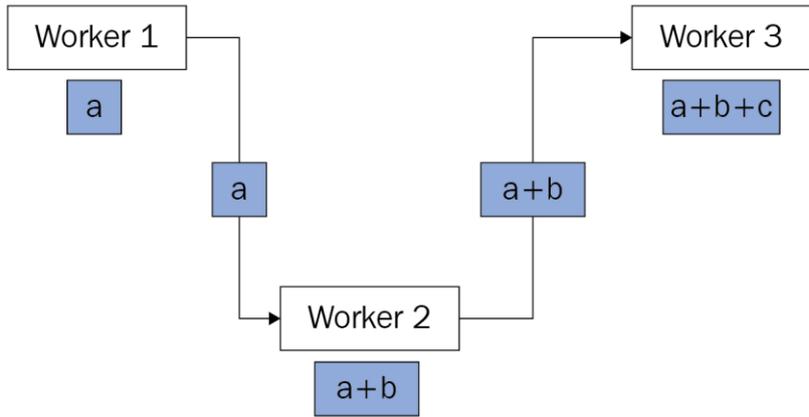
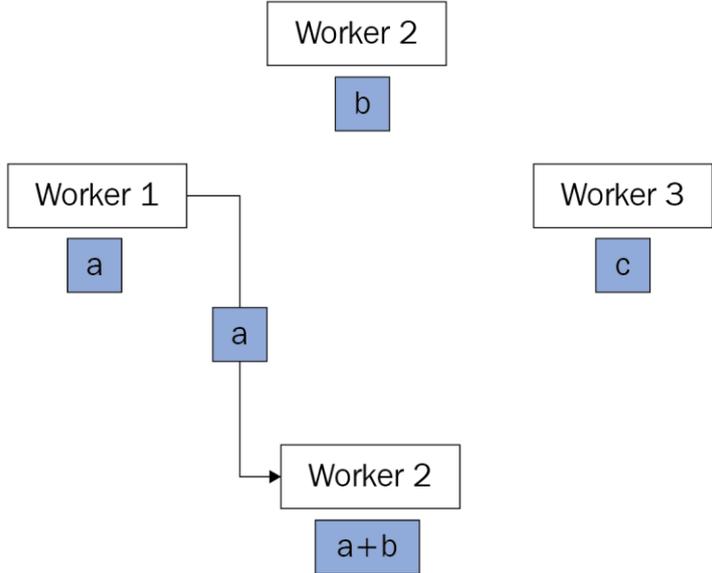
$$\nabla W = \nabla W^1 + \nabla W^2 + \dots + \nabla W^N$$

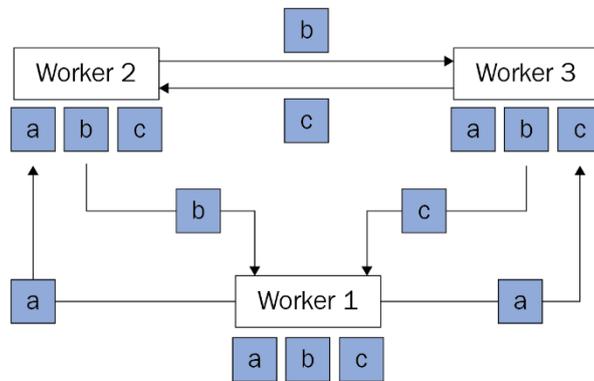
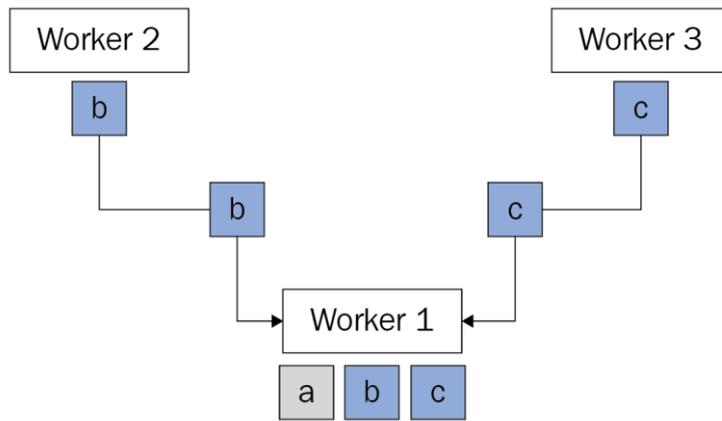
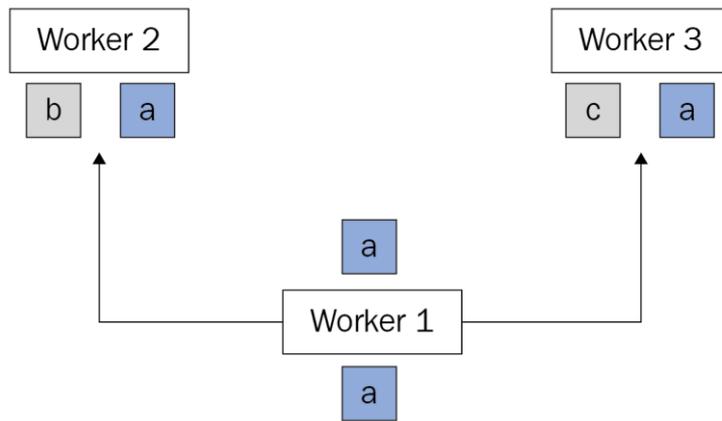
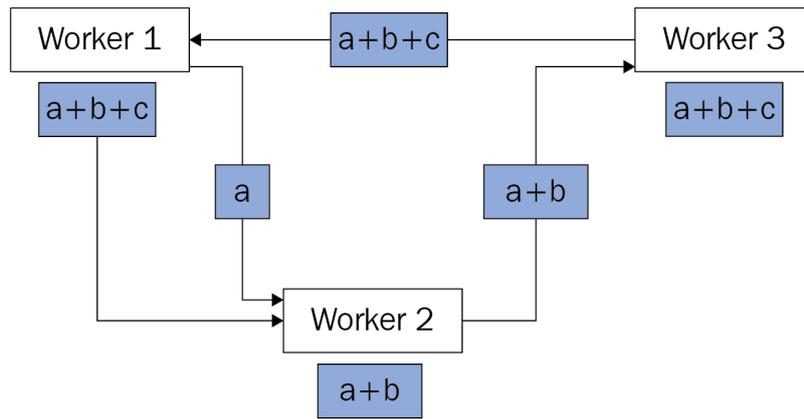
Chapter 2: Parameter Server and All-Reduce

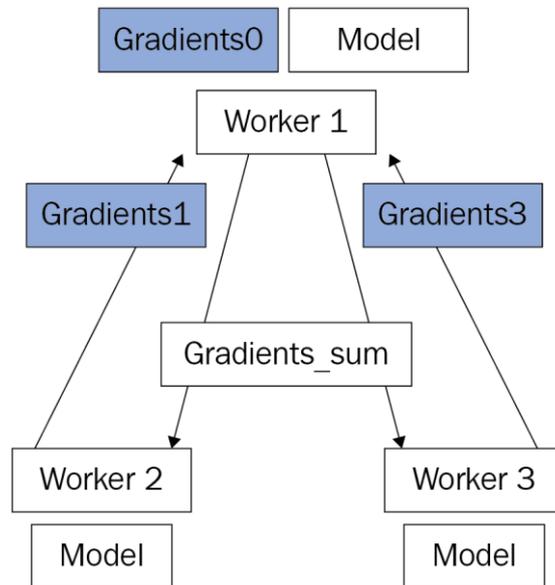
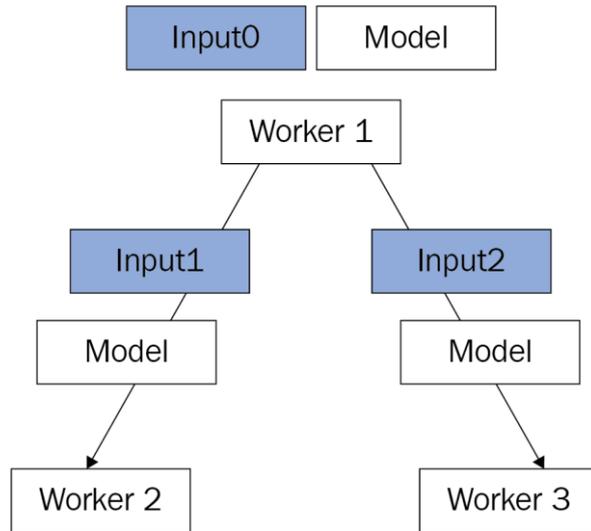








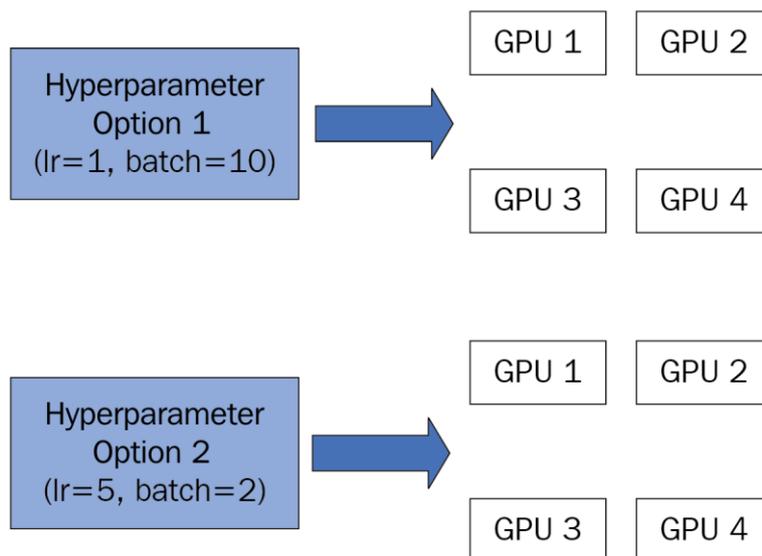




Processes:							
GPU	GI ID	CI ID	PID	Type	Process name	GPU Memory Usage	
0	N/A	N/A	7906	C	python	1393MiB	
1	N/A	N/A	7906	C	python	1393MiB	
2	N/A	N/A	7906	C	python	1393MiB	
3	N/A	N/A	7906	C	python	1035MiB	

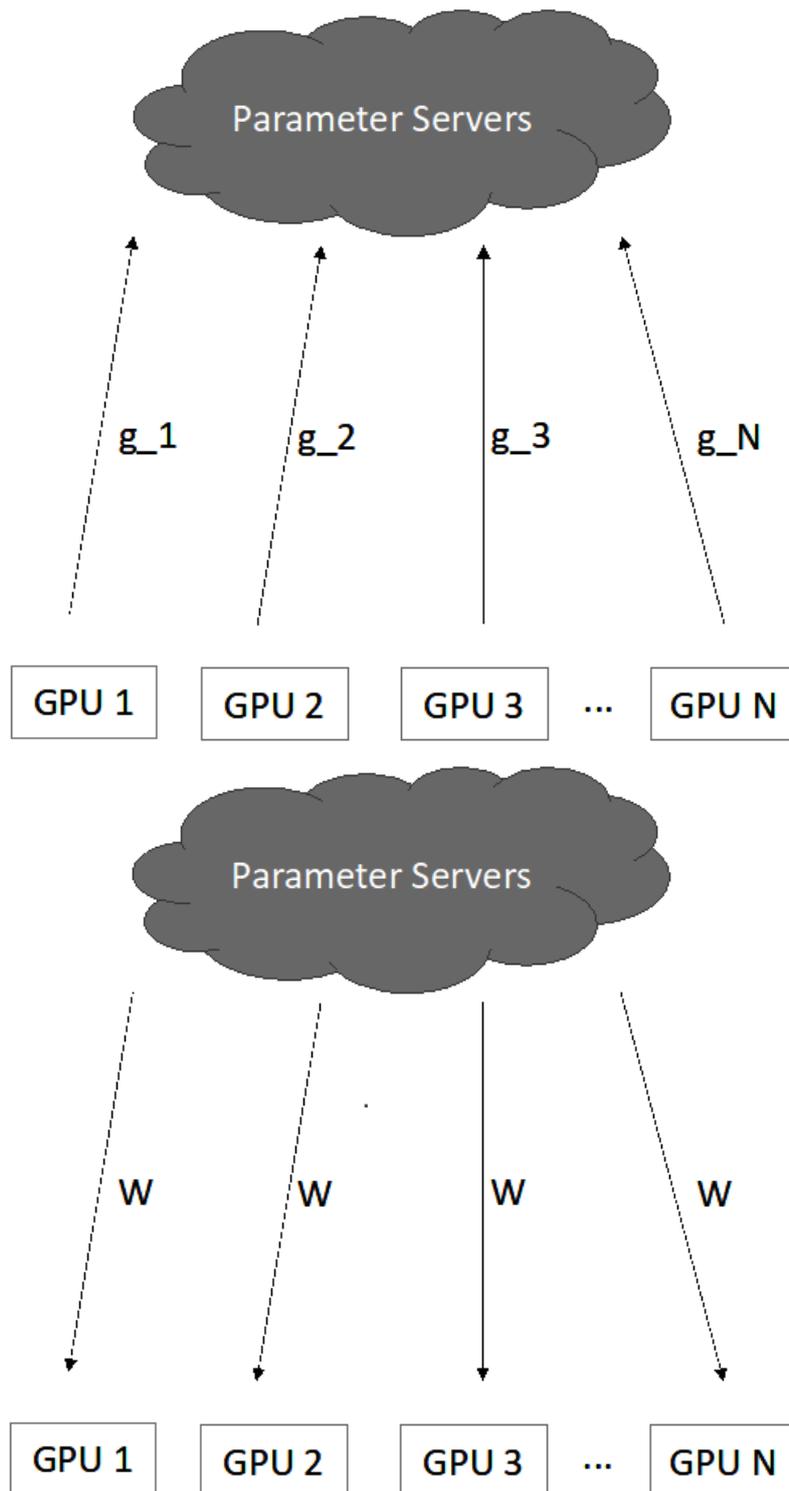
NVIDIA-SMI 450.142.00 Driver Version: 450.142.00 CUDA Version: 11.0							
GPU	Name	Persistence-MI	Bus-Id	Disp.A	Volatile	Uncorr. ECC	
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	MIG M.
0	Tesla V100-SXM2...	On	00000000:00:1B.0	Off		0	
N/A	45C	P0	71W / 300W	1768MiB / 16160MiB	25%	Default	N/A
1	Tesla V100-SXM2...	On	00000000:00:1C.0	Off		0	
N/A	46C	P0	67W / 300W	1792MiB / 16160MiB	24%	Default	N/A
2	Tesla V100-SXM2...	On	00000000:00:1D.0	Off		0	
N/A	47C	P0	82W / 300W	1792MiB / 16160MiB	26%	Default	N/A
3	Tesla V100-SXM2...	On	00000000:00:1E.0	Off		0	
N/A	45C	P0	73W / 300W	1768MiB / 16160MiB	29%	Default	N/A

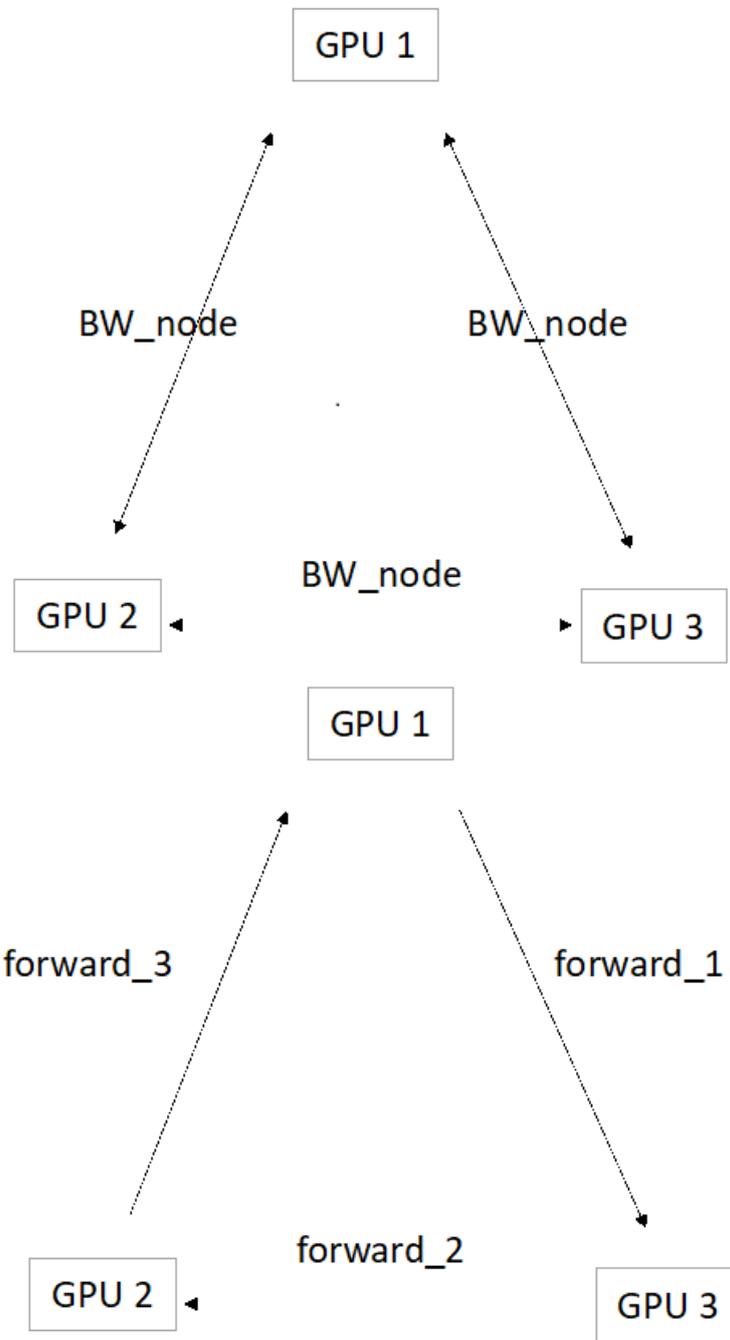
Processes:							
GPU	GI ID	CI ID	PID	Type	Process name	GPU Memory Usage	
0	N/A	N/A	27914	C	...rch_latest_p37/bin/python	1765MiB	
1	N/A	N/A	27915	C	...rch_latest_p37/bin/python	1789MiB	
2	N/A	N/A	27916	C	...rch_latest_p37/bin/python	1789MiB	
3	N/A	N/A	27917	C	...rch_latest_p37/bin/python	1765MiB	

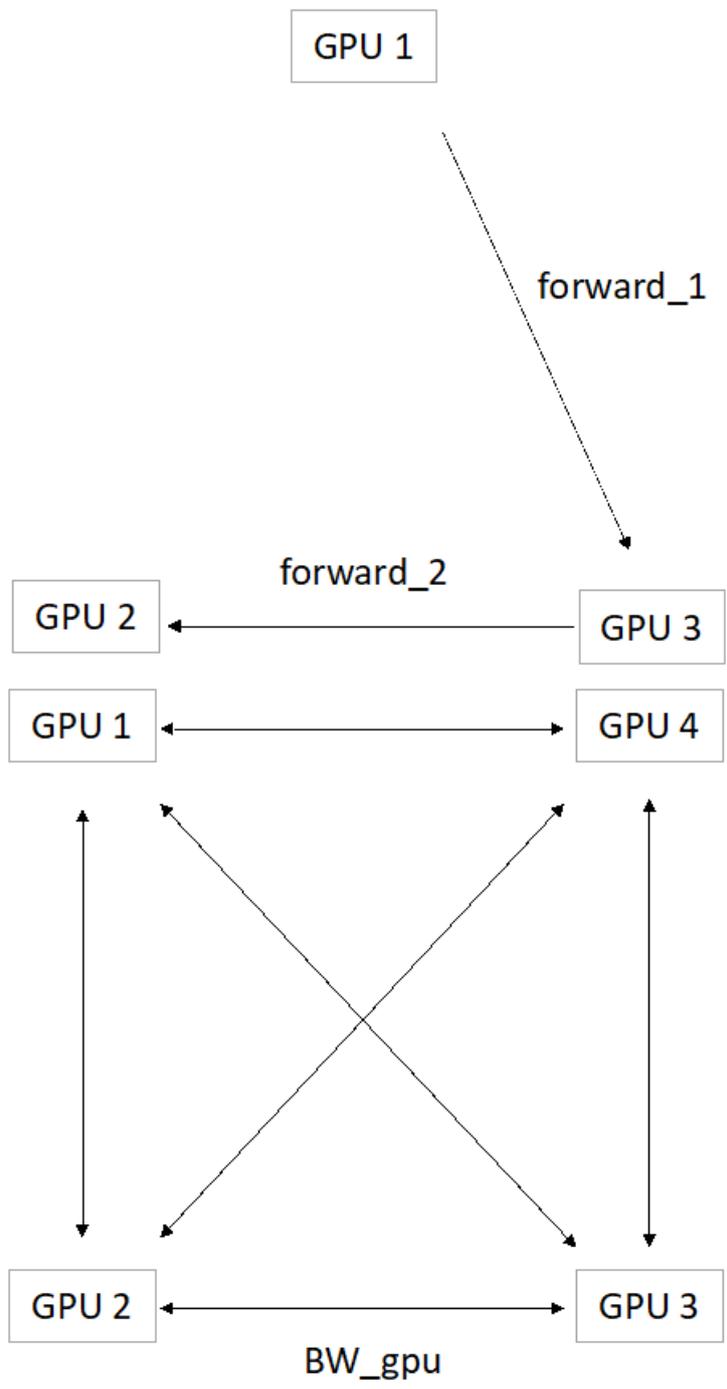


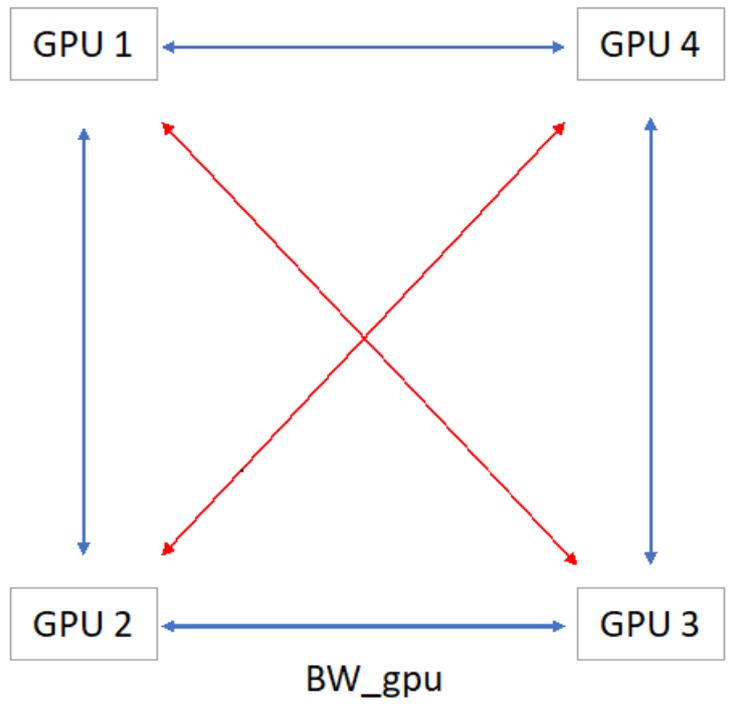
Load checkpoint 3	Load checkpoint 6
Load checkpoint 2	Load checkpoint 5
Load checkpoint 1	Load checkpoint 7
Load checkpoint 0	Load checkpoint 4
Checkpoint loading done!	Checkpoint loading done!
Checkpoint loading done!	Checkpoint loading done!
GPU 0, Test Accuracy 0.1198	GPU 5, Test Accuracy 0.119
Test Done!	Test Done!
GPU 1, Test Accuracy 0.1189	GPU 4, Test Accuracy 0.1202
Test Done!	Test Done!
Checkpoint loading done!	Checkpoint loading done!
Checkpoint loading done!	Checkpoint loading done!
GPU 2, Test Accuracy 0.1192	GPU 7, Test Accuracy 0.1195
Test Done!	Test Done!
GPU 3, Test Accuracy 0.1181	GPU 6, Test Accuracy 0.1193
Test Done!	Test Done!

Chapter 4: Bottlenecks and Solutions



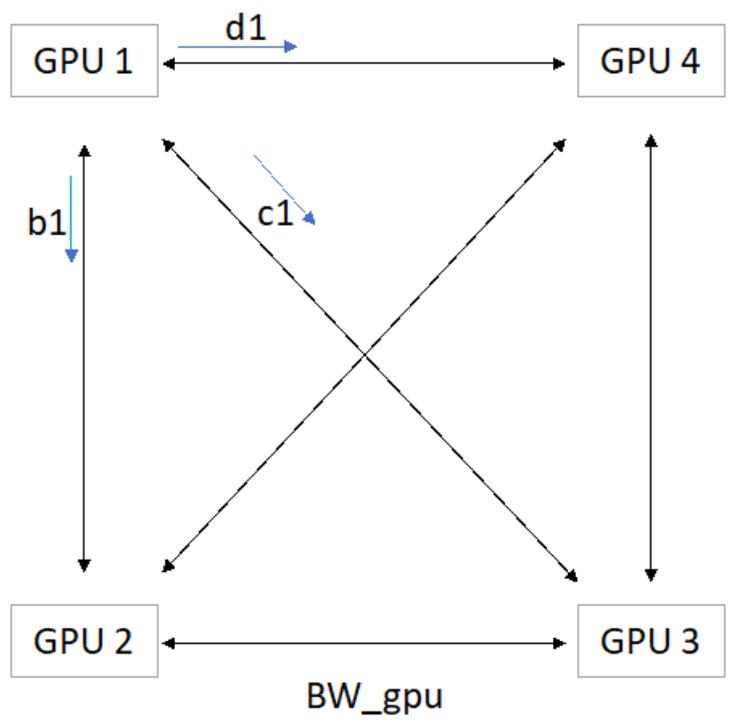






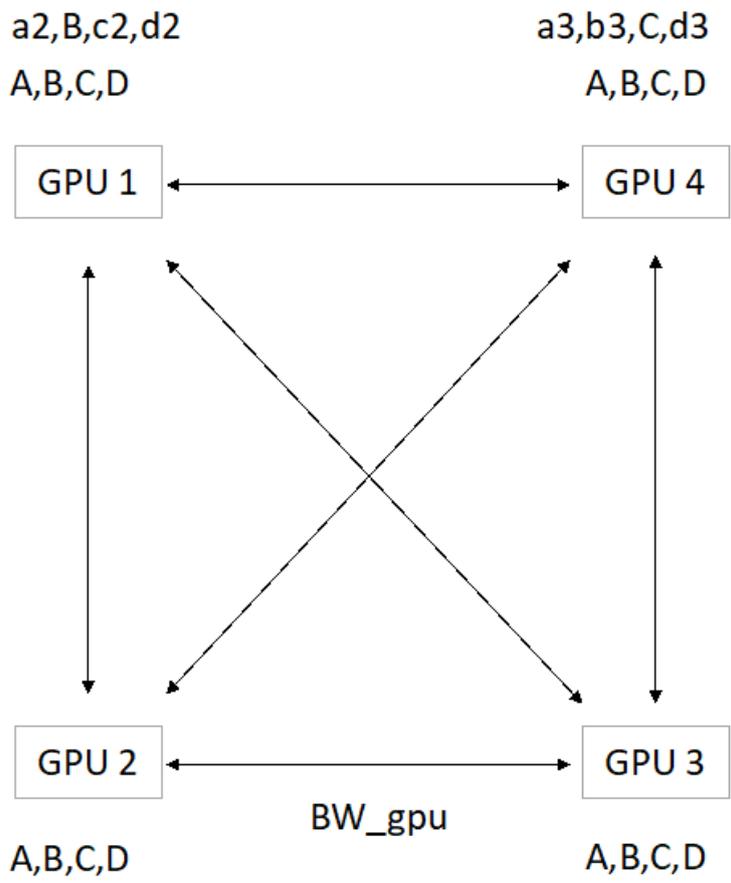
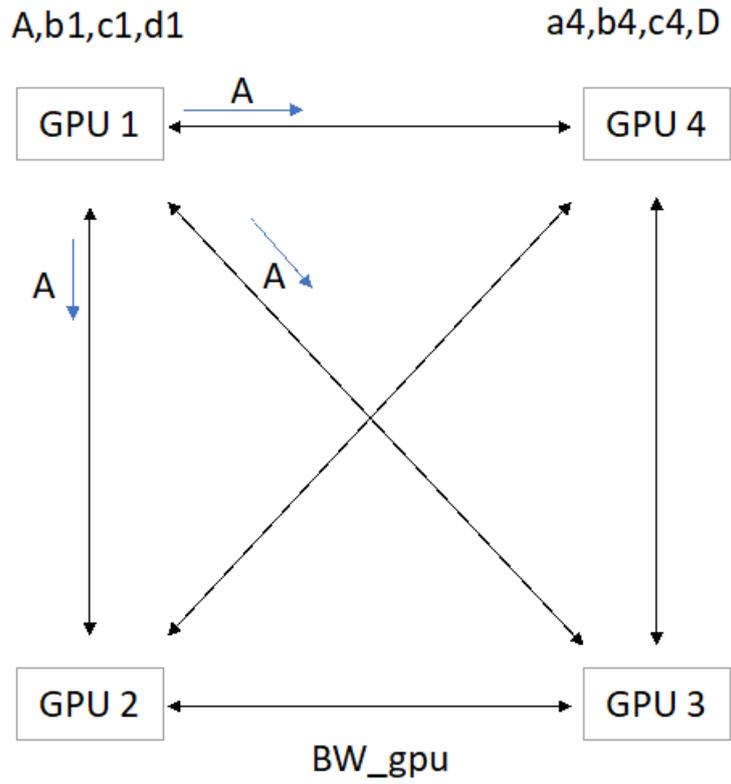
a1,b1,c1,d1

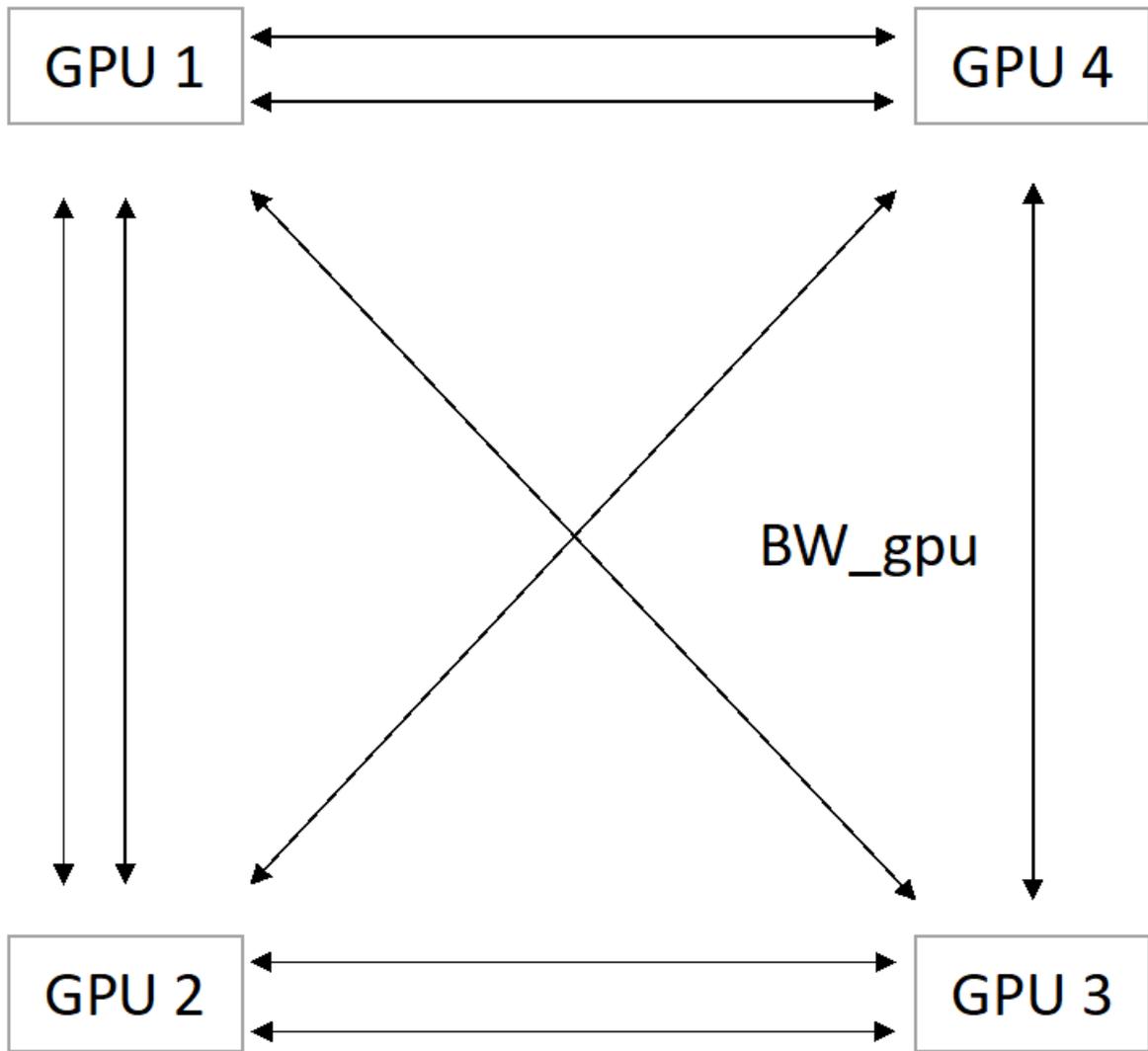
a4,b4,c4,d4

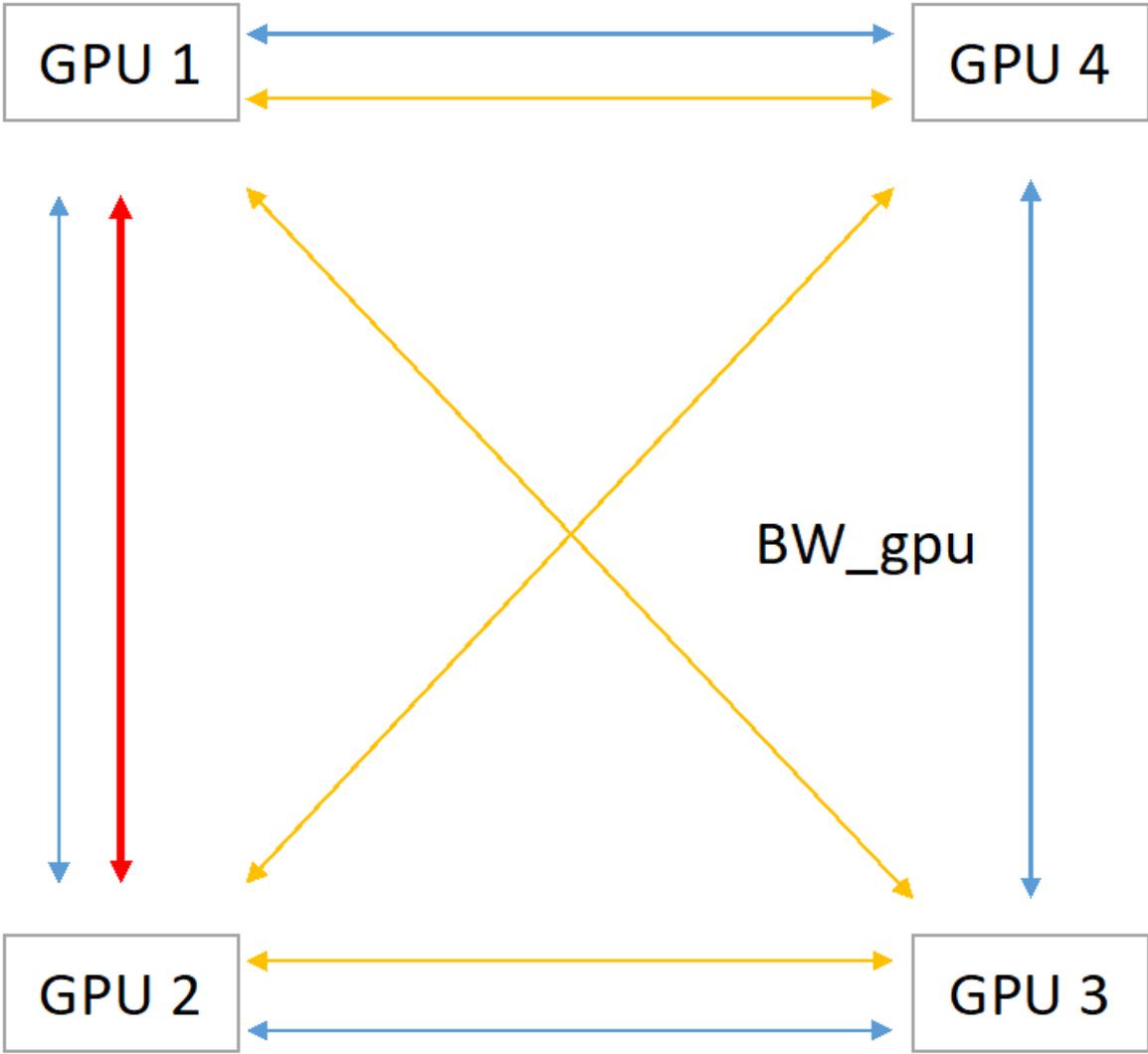


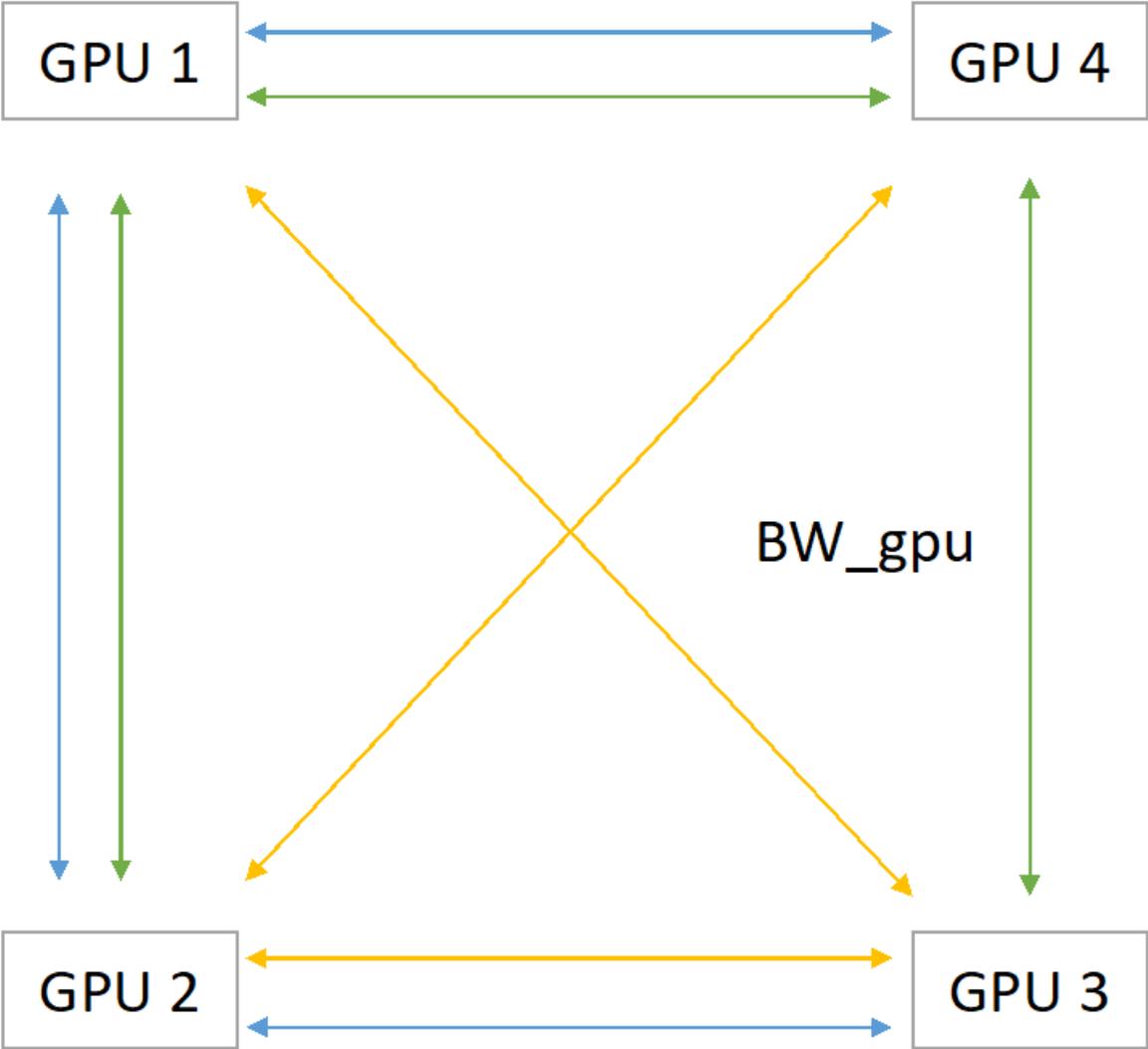
a2,b2,c2,d2

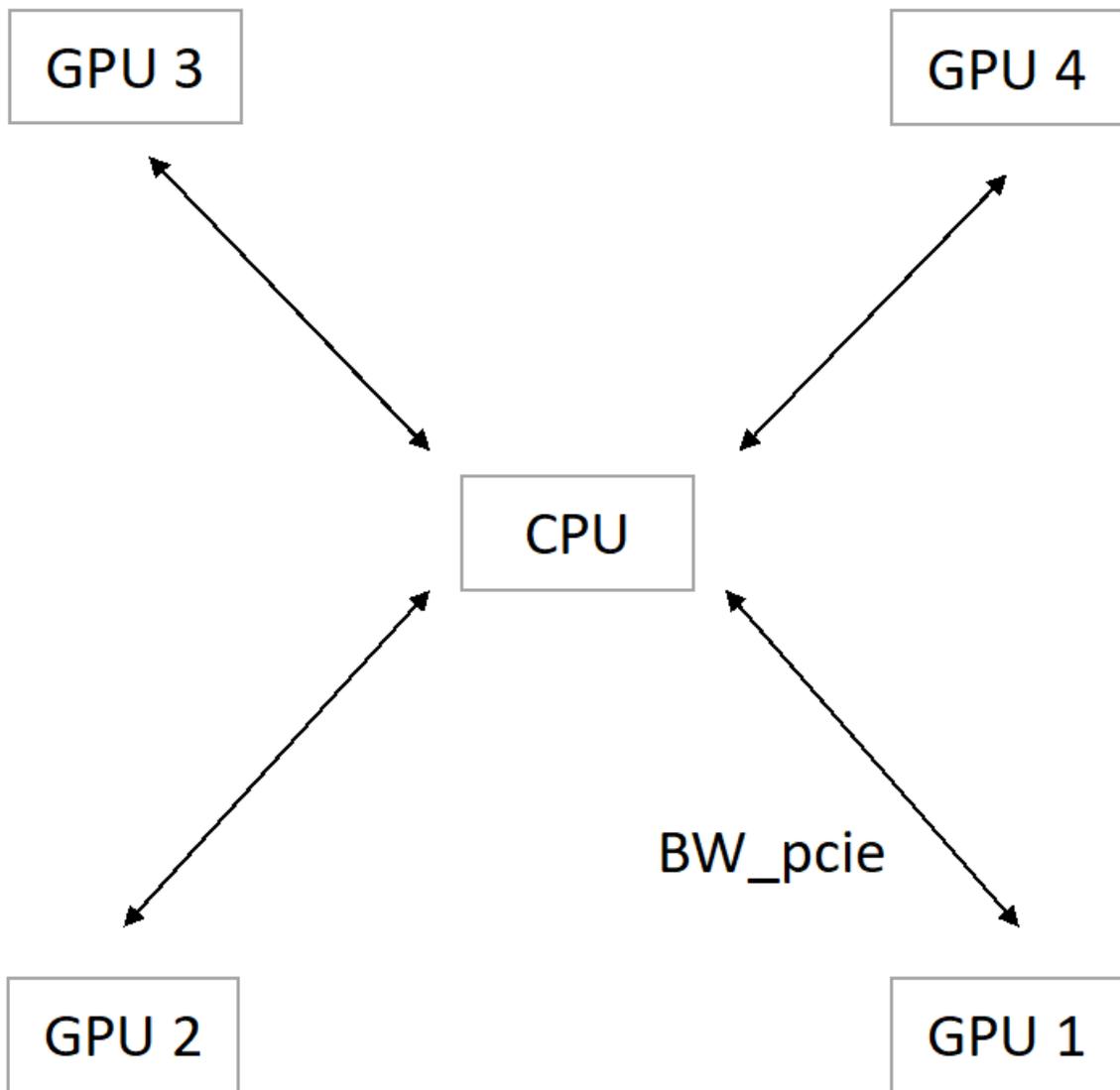
a3,b3,c3,d3





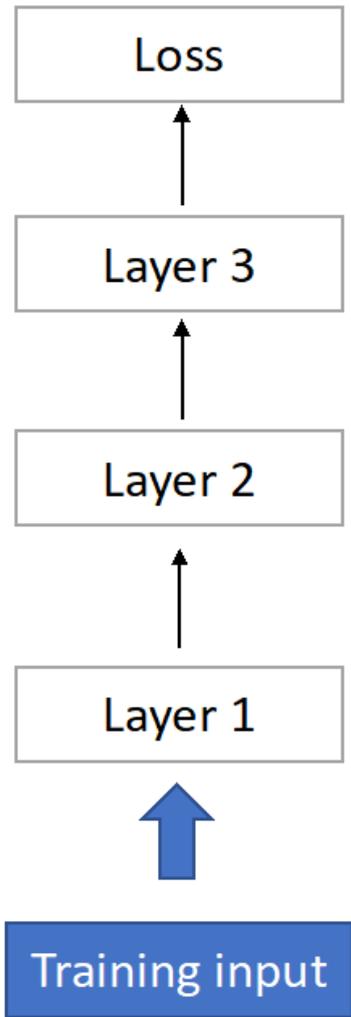


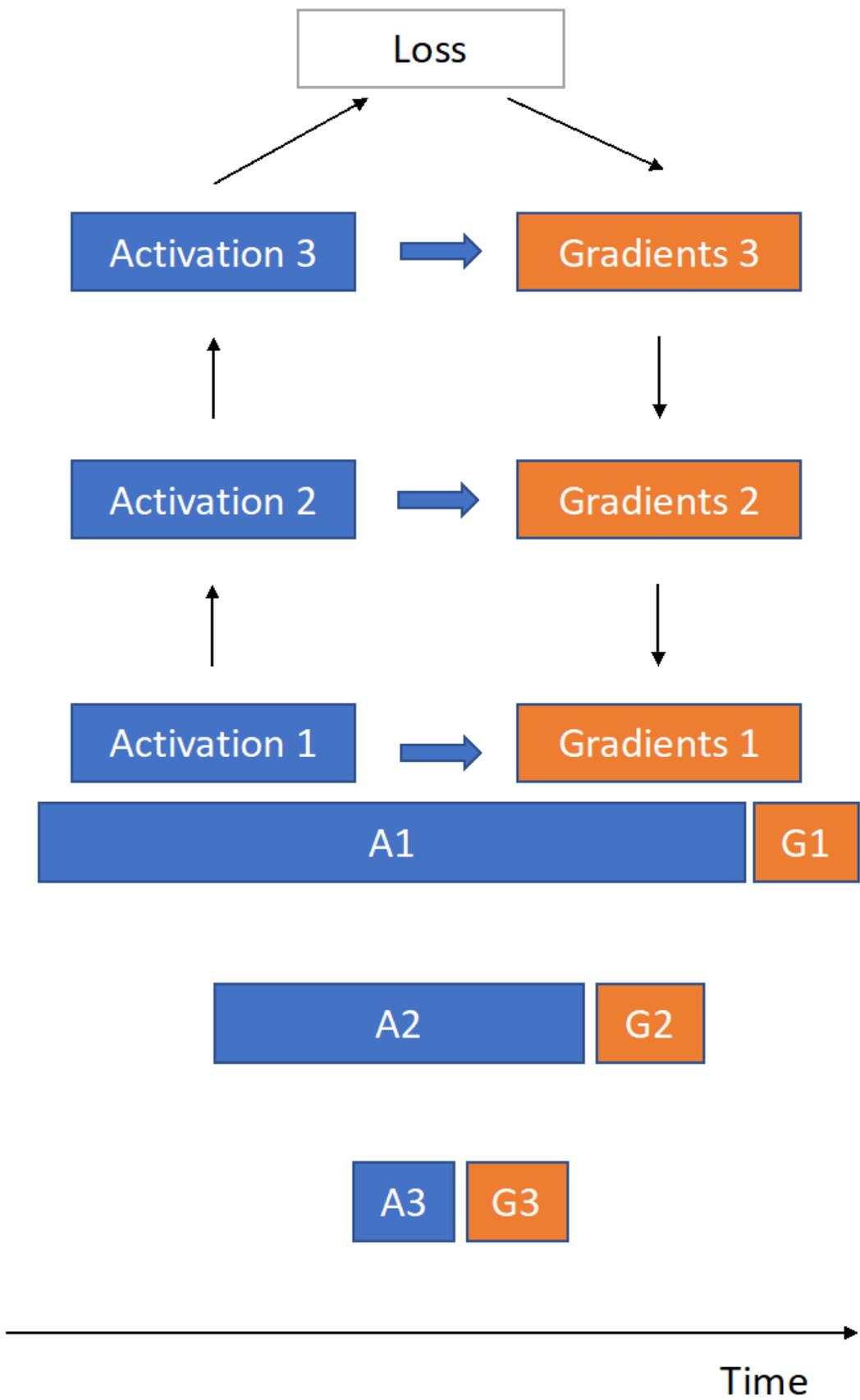


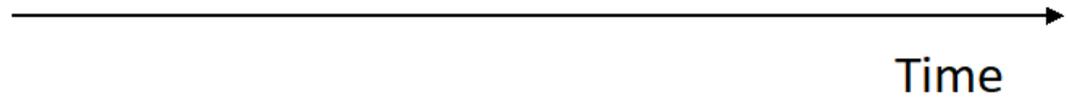
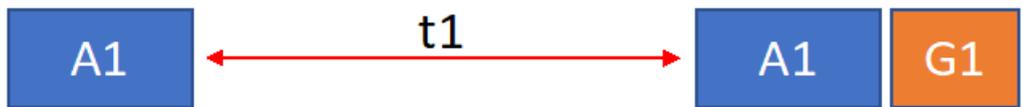
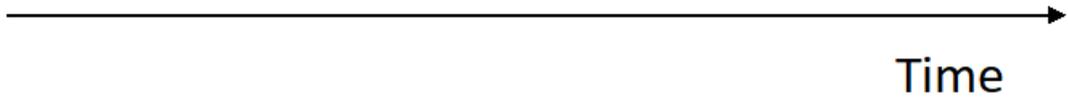


GPU type	On-device memory size
NVIDIA 1080	8 GB
NVIDIA RTX 2080	8 GB
NVIDIA K80	12 GB
NVIDIA V100	16 GB
NVIDIA A100	40 GB

	Recomputation	Quantization
Lossy/lossless	Lossless	Lossy
Reducing memory footprint	Yes	Yes
Computation overhead	Medium	Low







0



3

1 bit: 0

1 bit: 1

0

1

2

3

2 bits:
00

2 bits:
01

2 bits:
10

2 bits:
11

Chapter 5: Splitting the Model

```
Python 3.7.10 | packaged by conda-forge | (default, Feb 19 2021, 16:07:37)
[GCC 9.3.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import transformers
>>> print(transformers.__version__)
4.10.3
>>>
```

```
Training epoch 1
0% | 0/86136
[00:00<?, ?it/s]Traceback (most recent call last):
  File "bert.py", line 198, in <module>
    end_positions=end_token_idx, return_dict=False)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/modules/module.py",
line 889, in _call_impl
    result = self.forward(*input, **kwargs)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 1825, in forward
    return_dict=return_dict,
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/modules/module.py",
line 889, in _call_impl
    result = self.forward(*input, **kwargs)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 1000, in forward
    return_dict=return_dict,
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/modules/module.py",
line 889, in _call_impl
    result = self.forward(*input, **kwargs)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 589, in forward
    output_attentions,
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/modules/module.py",
line 889, in _call_impl
    result = self.forward(*input, **kwargs)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 511, in forward
    self.feed_forward_chunk, self.chunk_size_feed_forward, self.seq_len_dim, attention_output
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/modeling_utils.p
y", line 2196, in apply_chunking_to_forward
    return forward_fn(*input_tensors)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 522, in feed_forward_chunk
    intermediate_output = self.intermediate(attention_output)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/modules/module.py",
line 889, in _call_impl
    result = self.forward(*input, **kwargs)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/transformers/models/bert/mod
eling_bert.py", line 426, in forward
    hidden_states = self.intermediate_act_fn(hidden_states)
  File "/home/ubuntu/anaconda3/envs/pytorch_latest_p37/lib/python3.7/site-packages/torch/nn/functional.py", line
1459, in gelu
    return torch._C._nn.gelu(input)
RuntimeError: CUDA out of memory. Tried to allocate 144.00 MiB (GPU 0; 15.78 GiB total capacity; 14.19 GiB alrea
dy allocated; 102.75 MiB free; 14.25 GiB reserved in total by PyTorch)
0% | 0/86136
[00:04<?, ?it/s]
```



```
+-----+
| NVIDIA-SMI 450.119.03   Driver Version: 450.119.03   CUDA Version: 11.0   |
+-----+-----+
| GPU  Name          Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|  Memory-Usage | GPU-Util  Compute M. |
|                                           MIG M.         |
+-----+-----+
|   0   Tesla V100-SXM2...  On    | 00000000:00:1E.0 Off  |           0         |
| N/A   57C    P0     118W / 300W |  3727MiB / 16160MiB |    65%      Default  |
|                                           N/A              |
+-----+-----+
| Processes:                                                       GPU Memory |
|  GPU   GI    CI          PID  Type  Process name      Usage   |
|-----+-----+
|   0   N/A   N/A         14498   C   python             3725MiB |
+-----+-----+
```

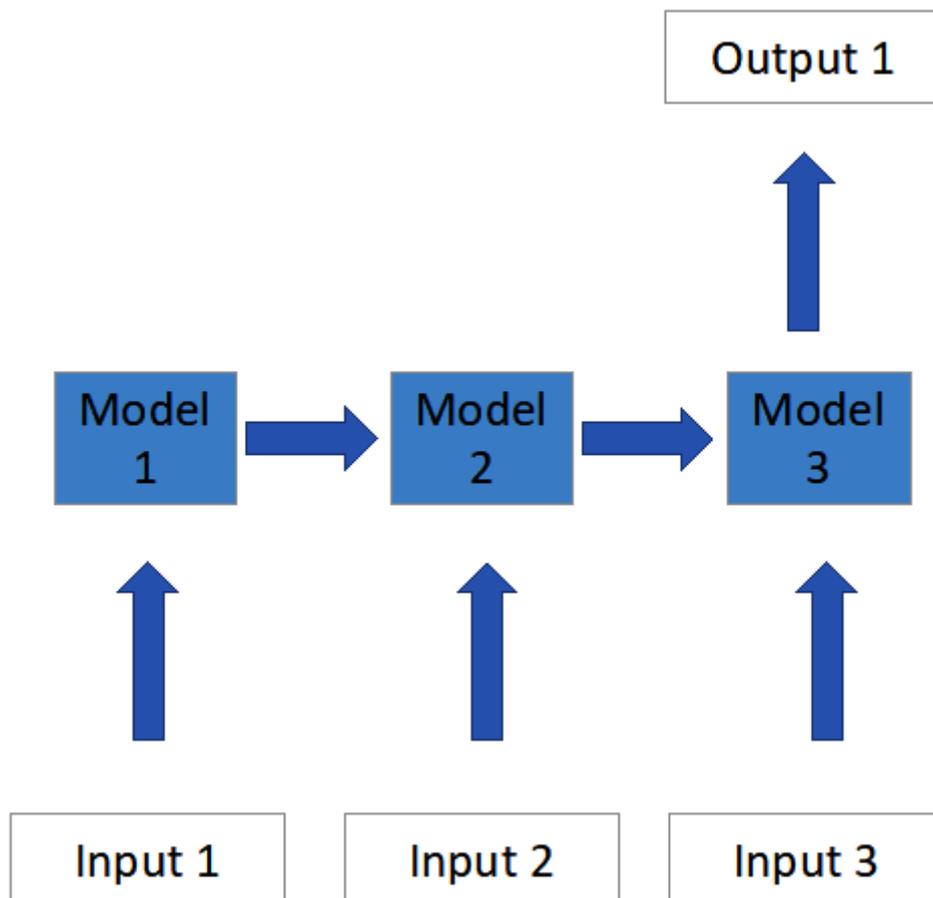
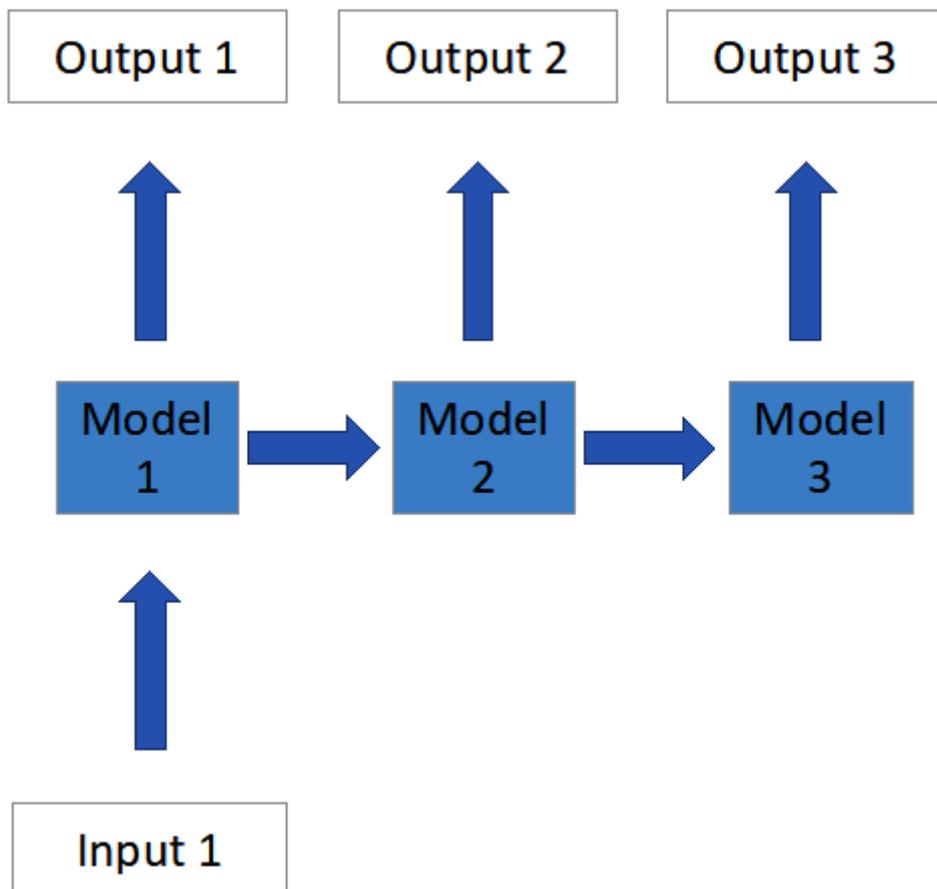
Output 1

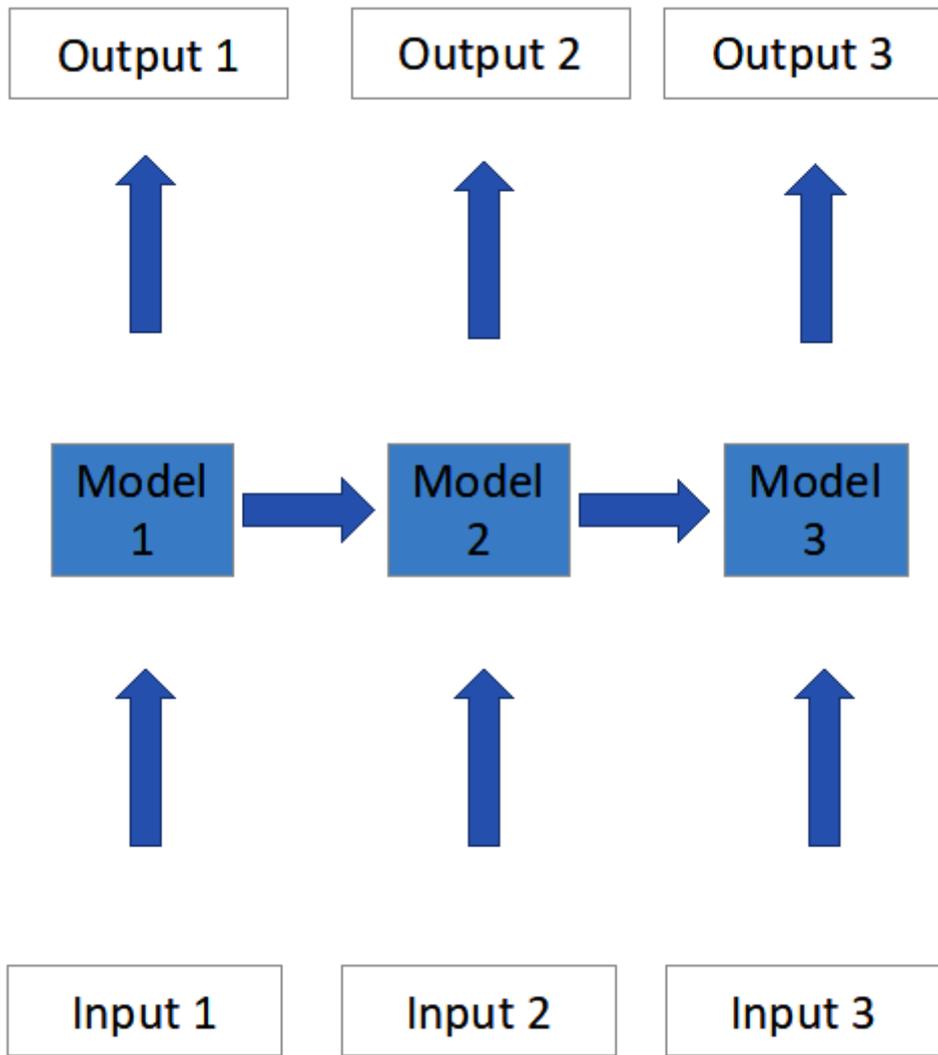


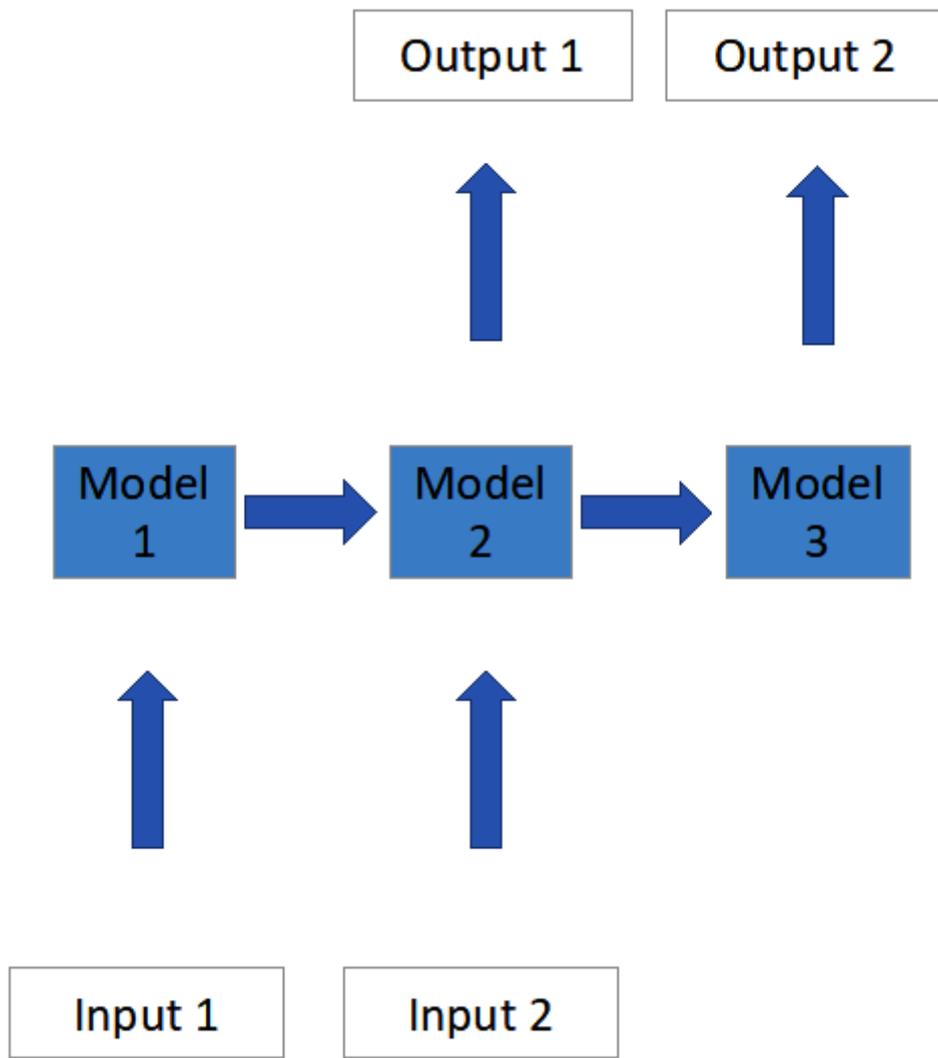
Model



Input 1







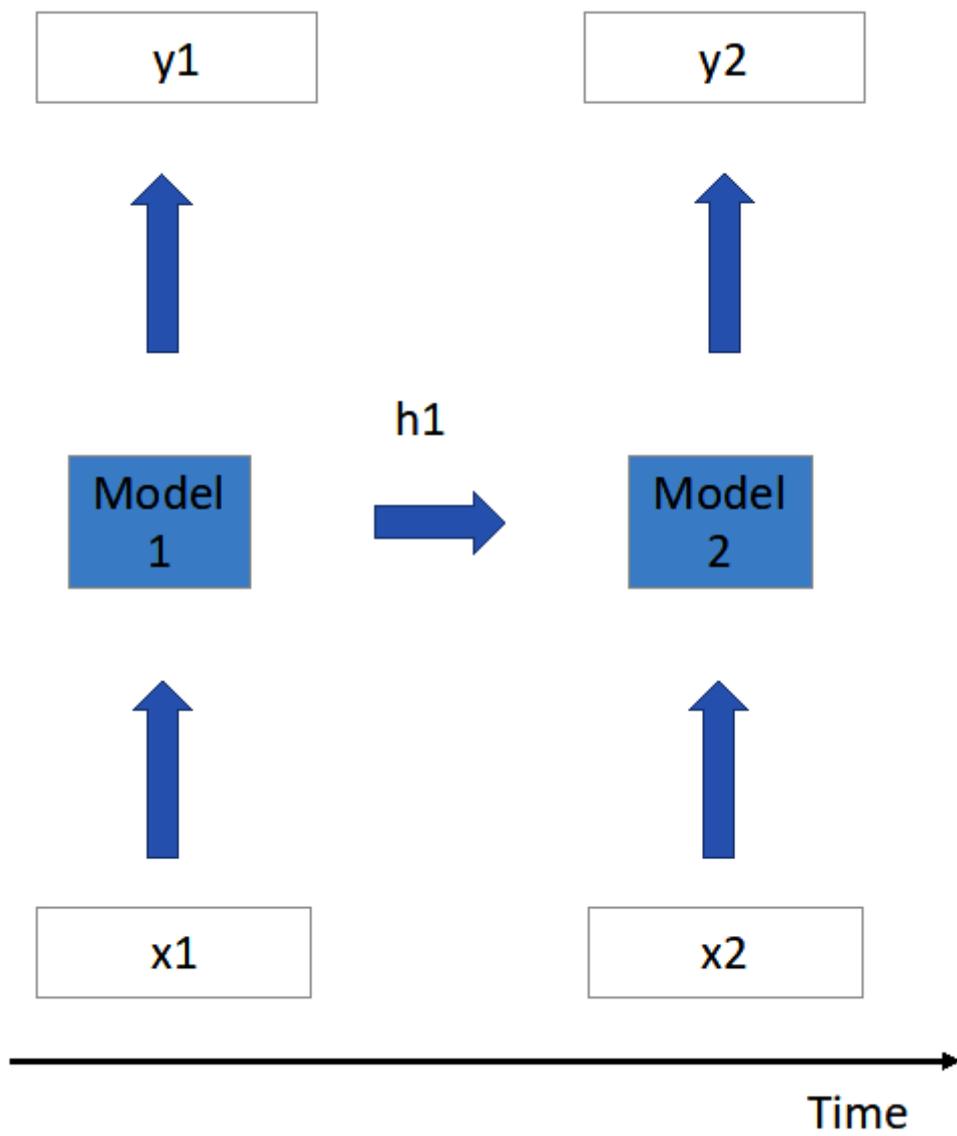
Output j

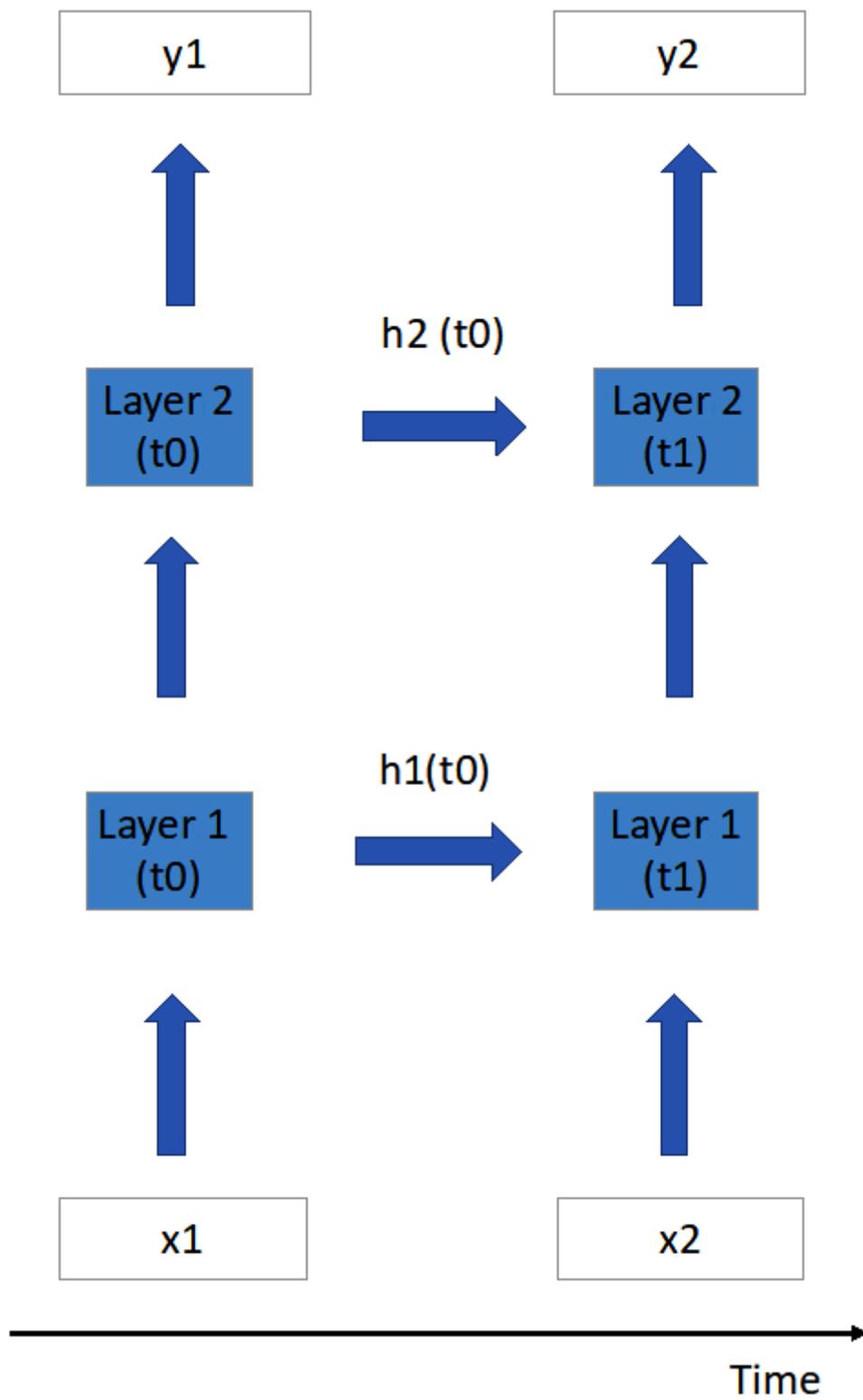


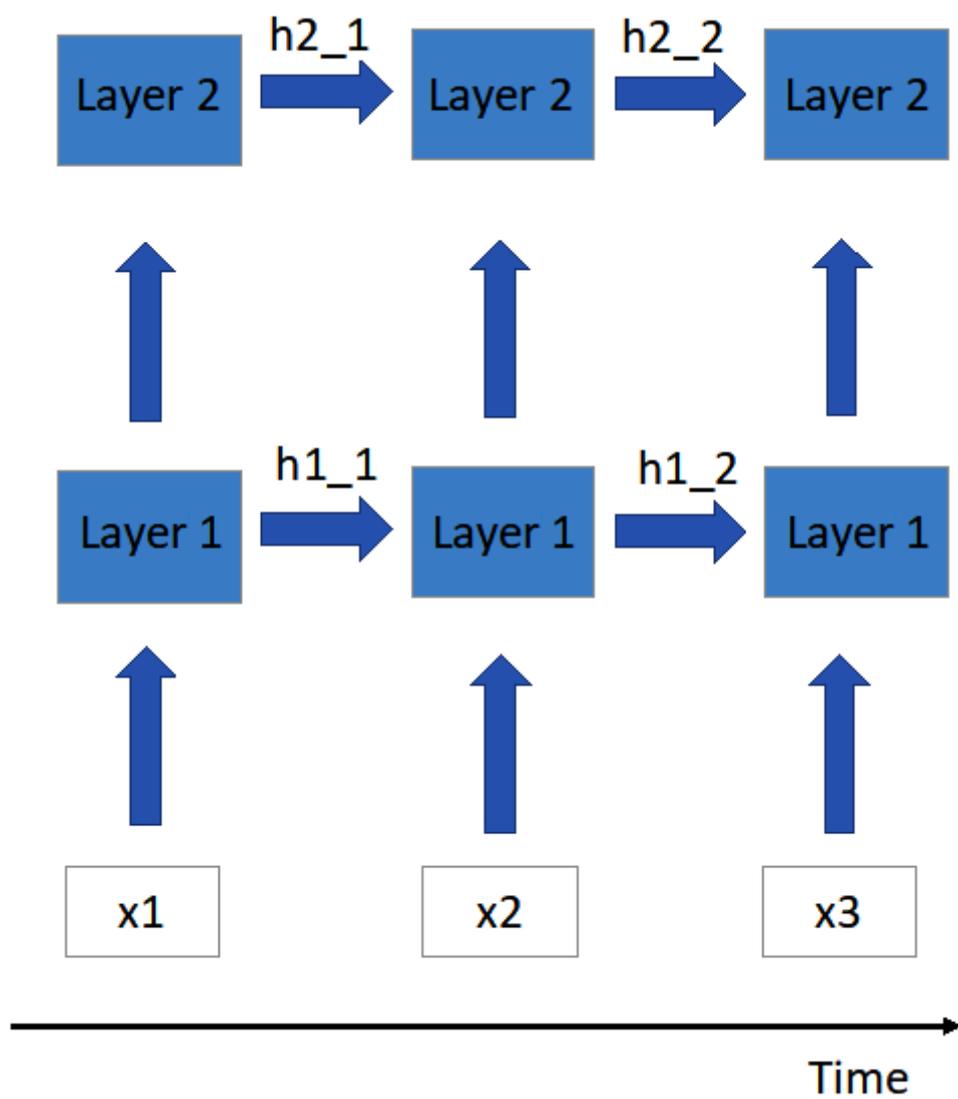
Model

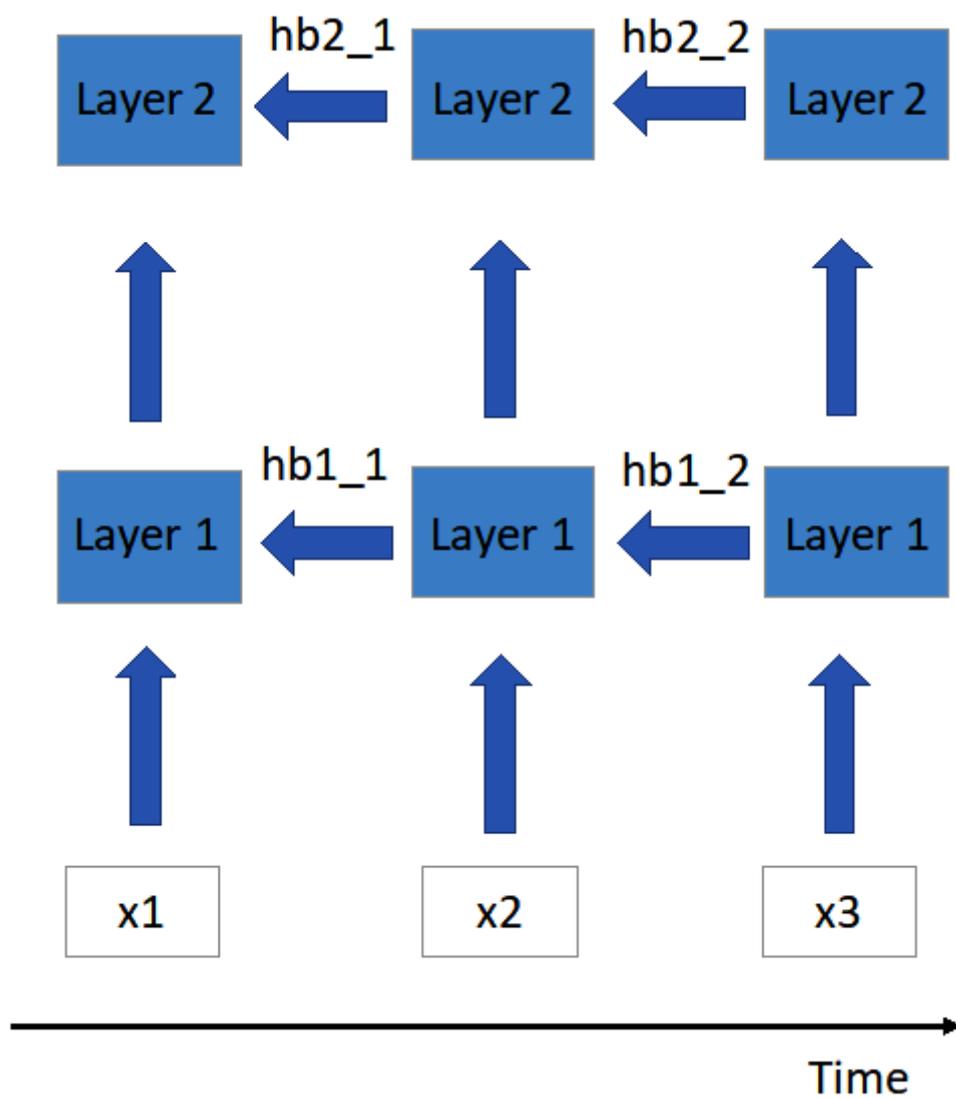


Input i

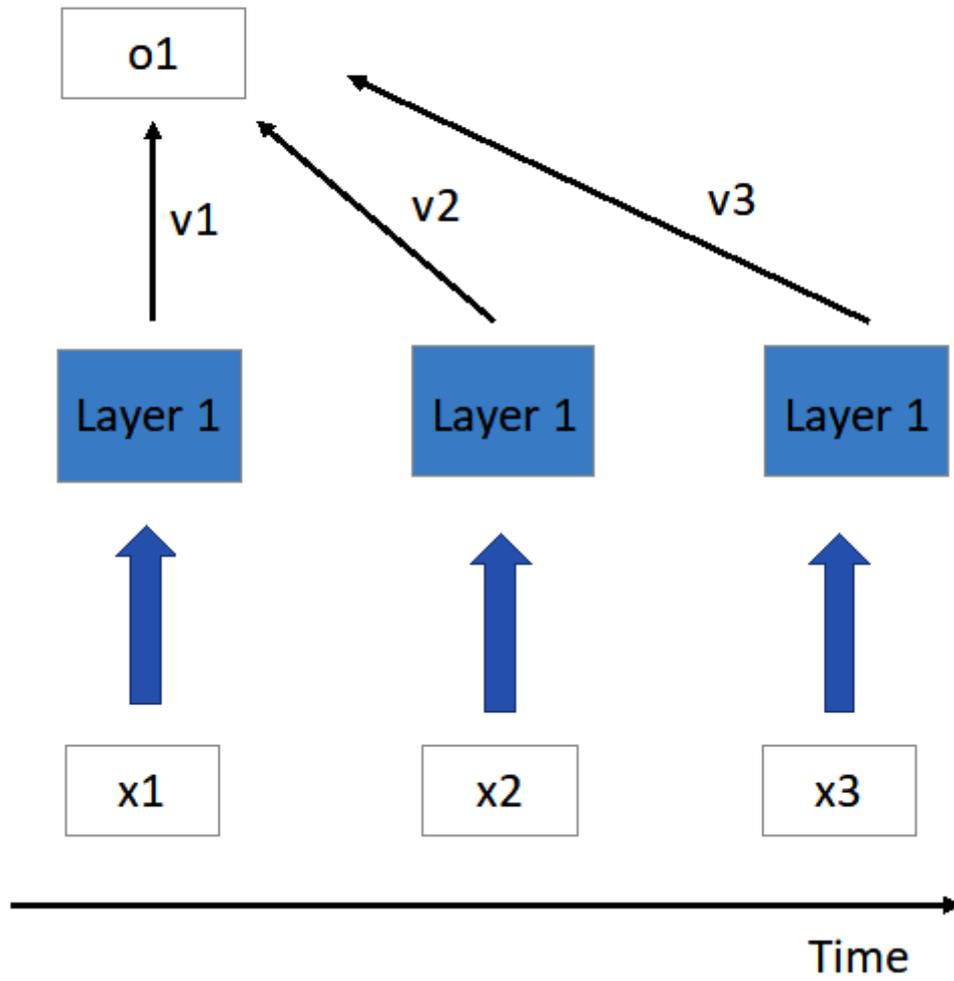


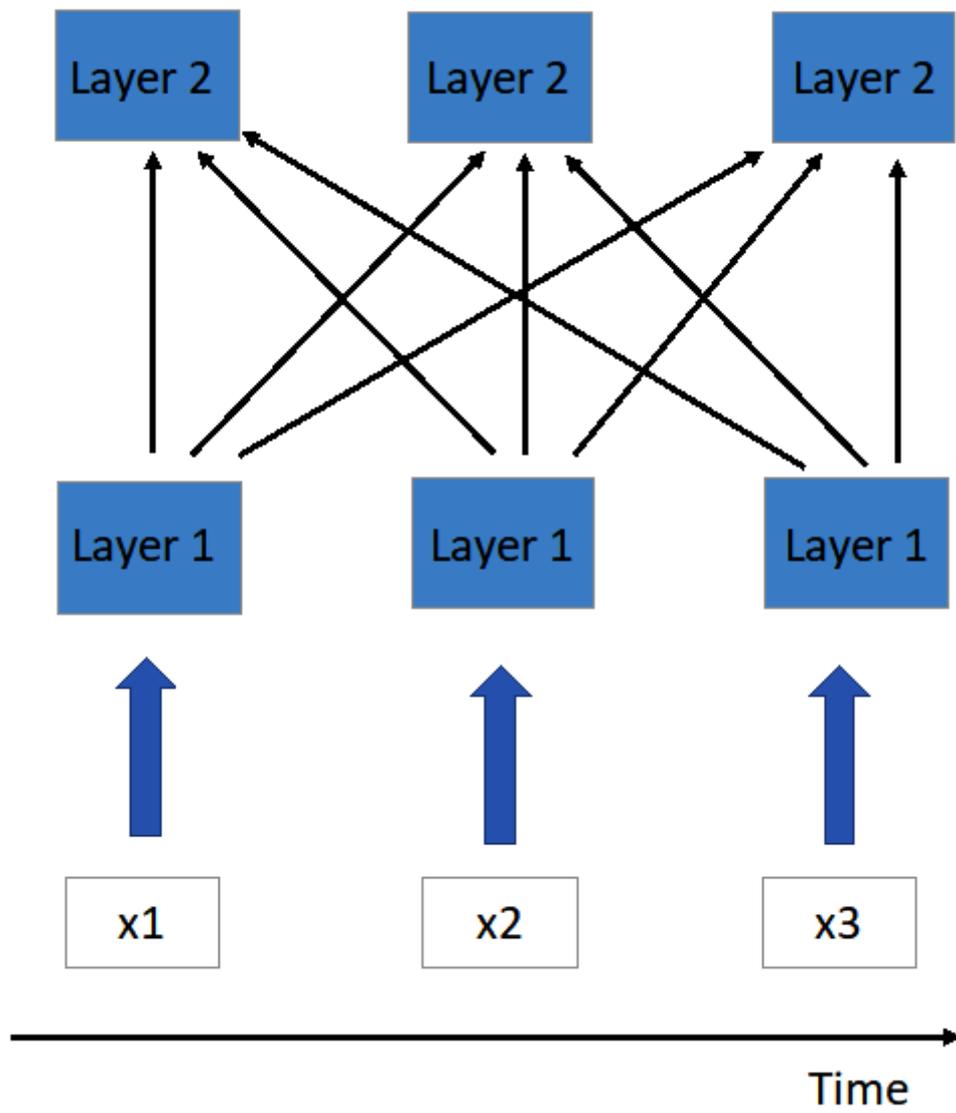


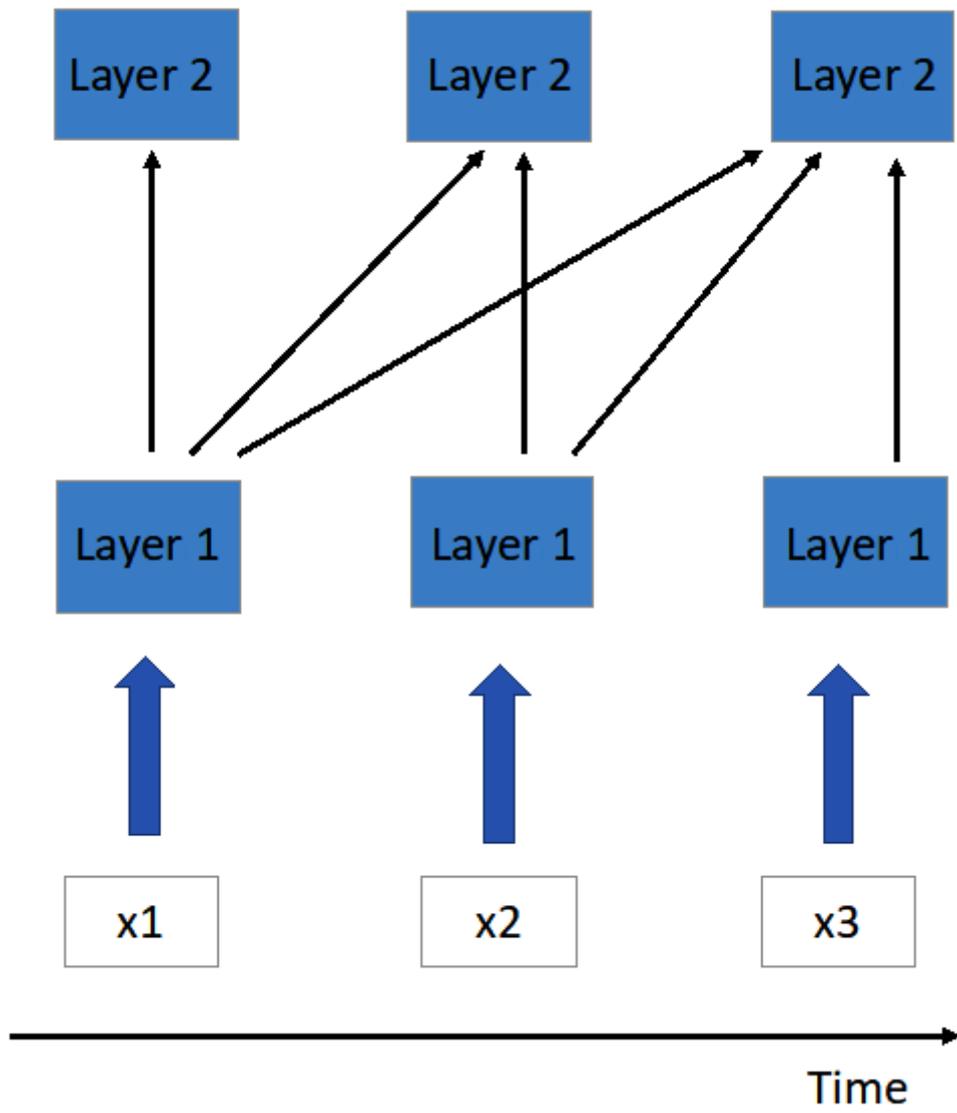


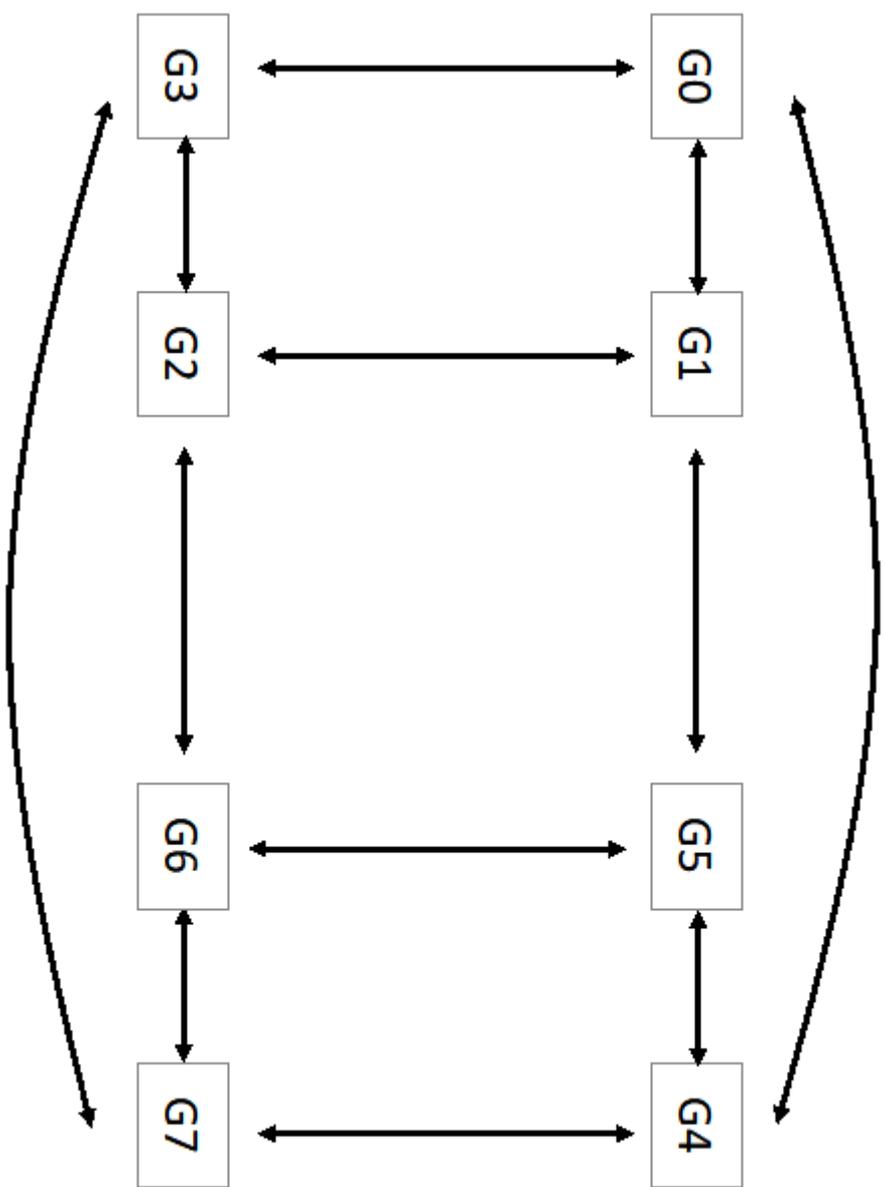


$$o1 = w1*v1 + w2*v2 + w3*v3$$

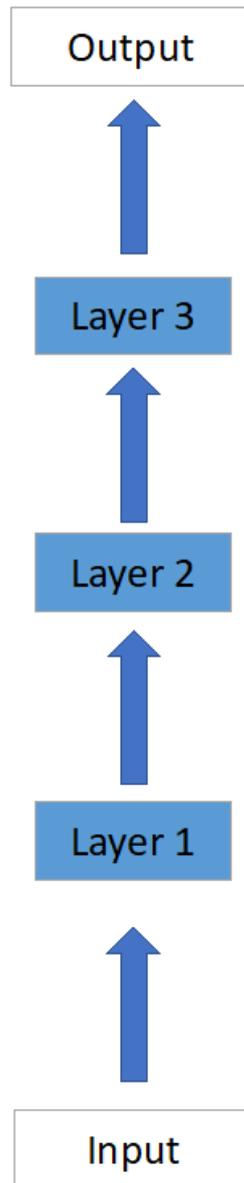


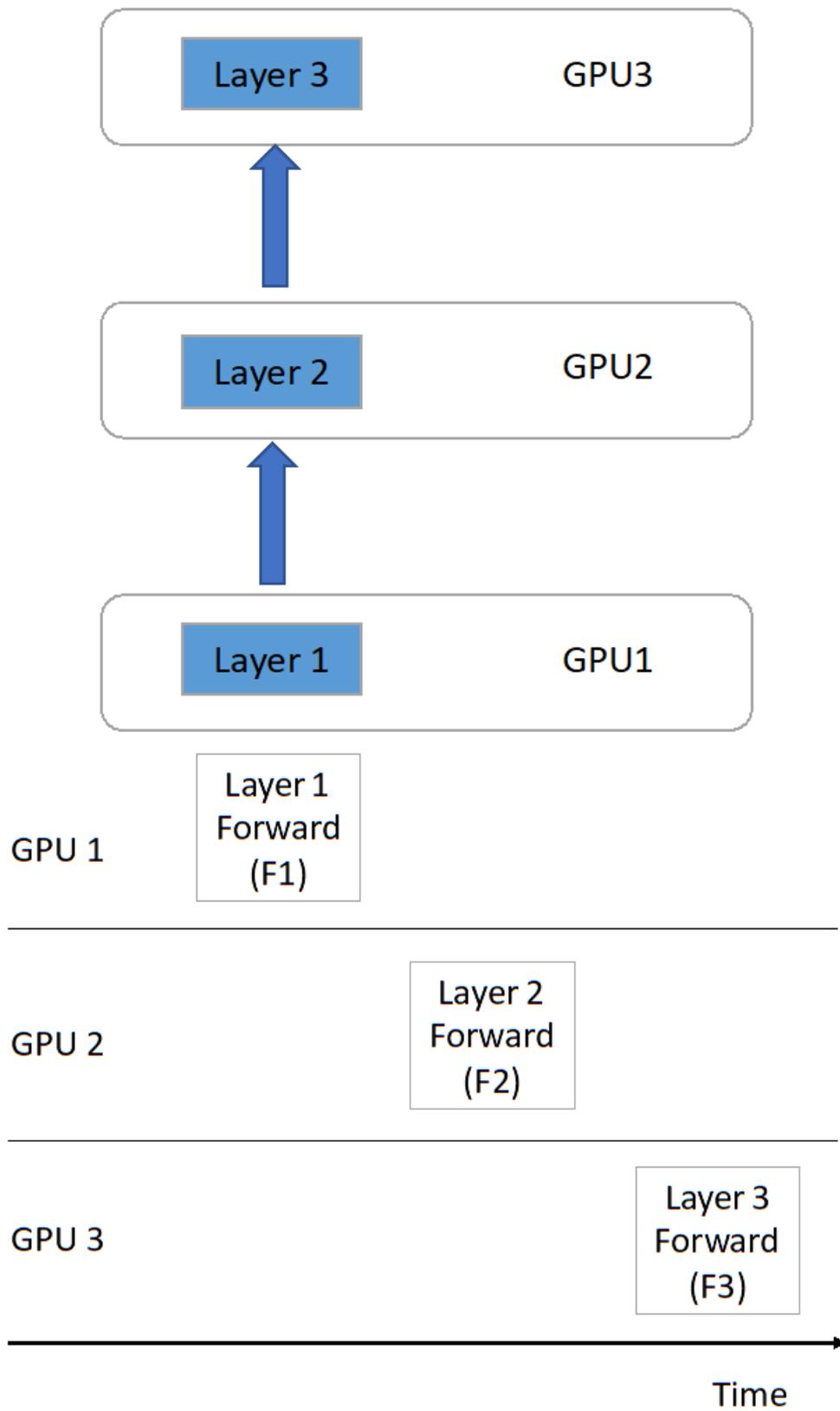


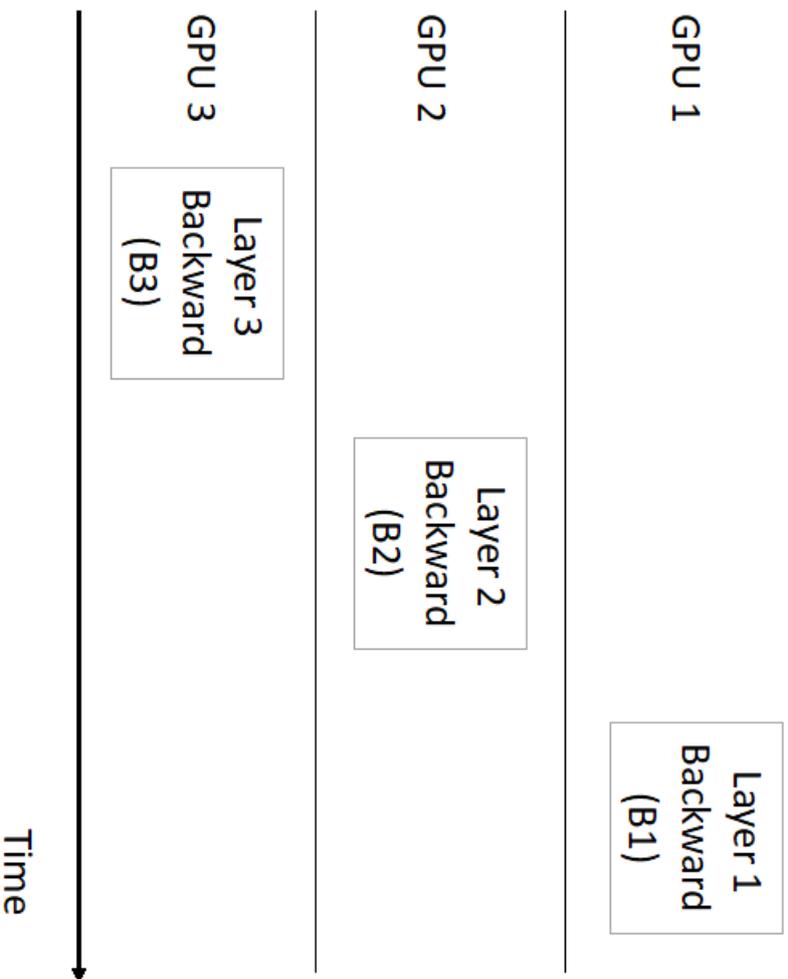


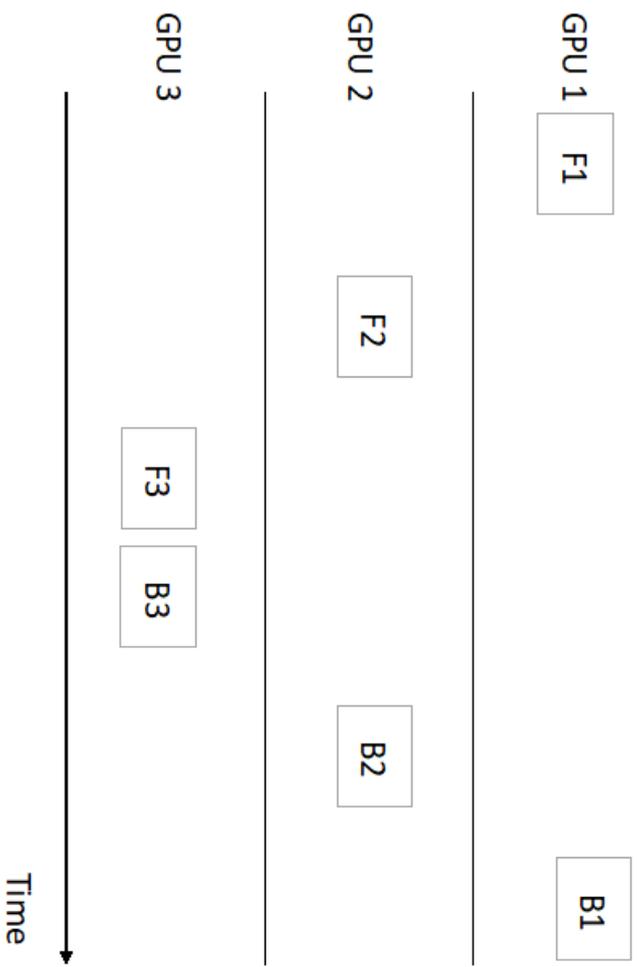


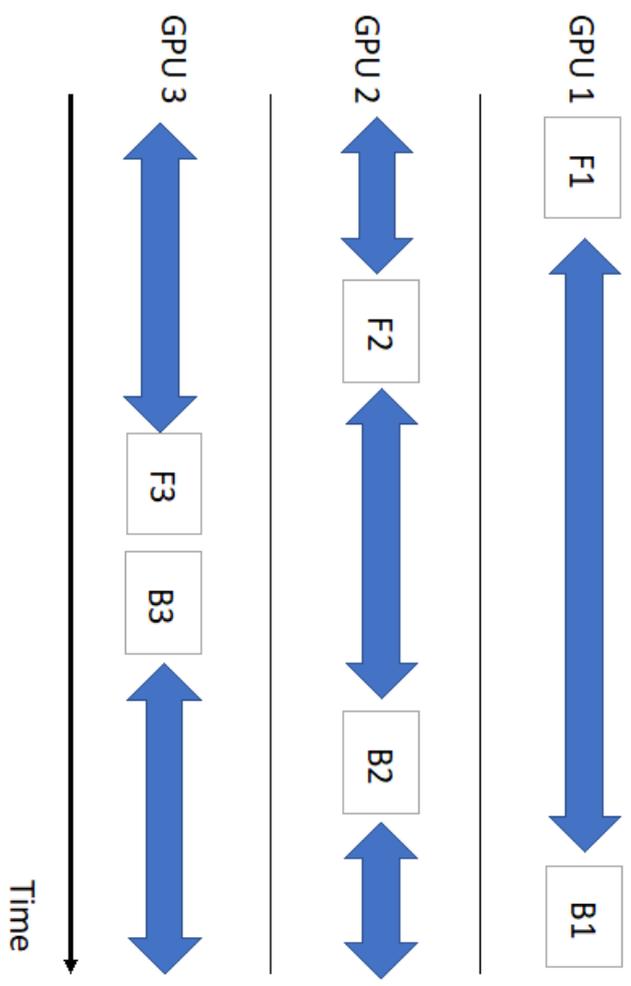
Chapter 6: Pipeline Input and Layer Split

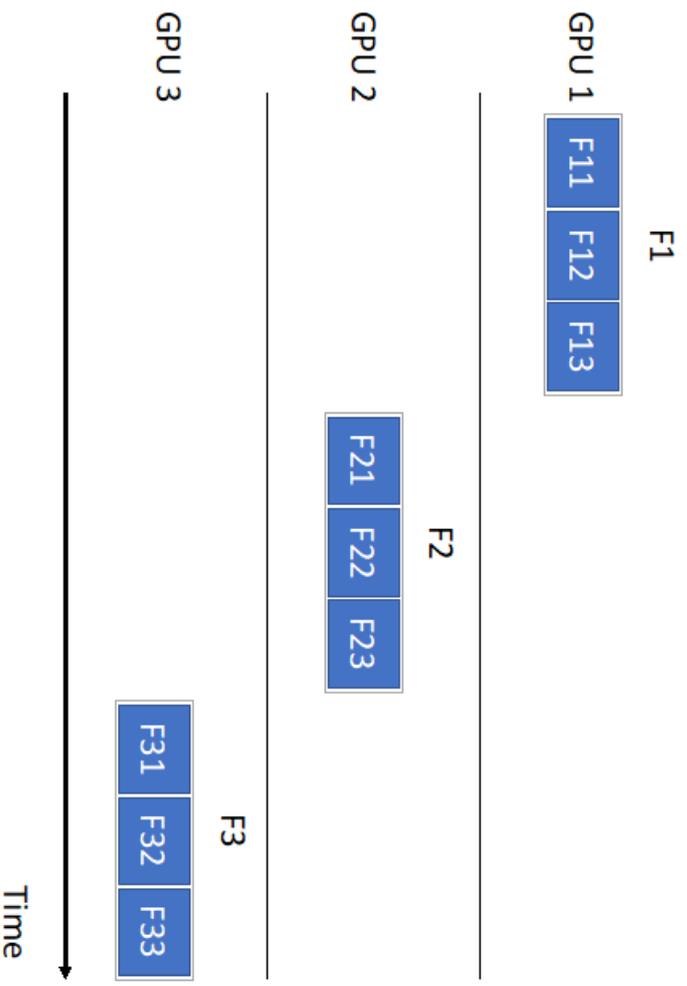


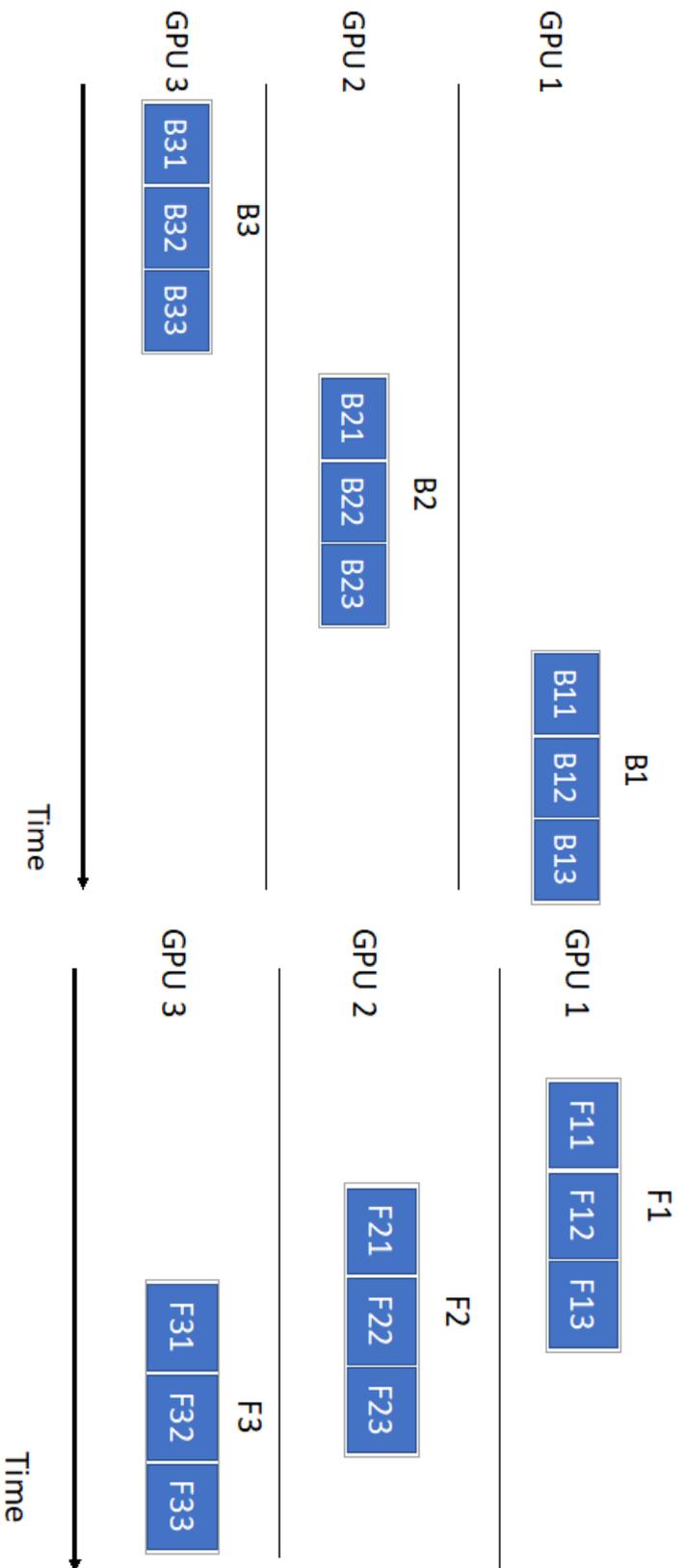


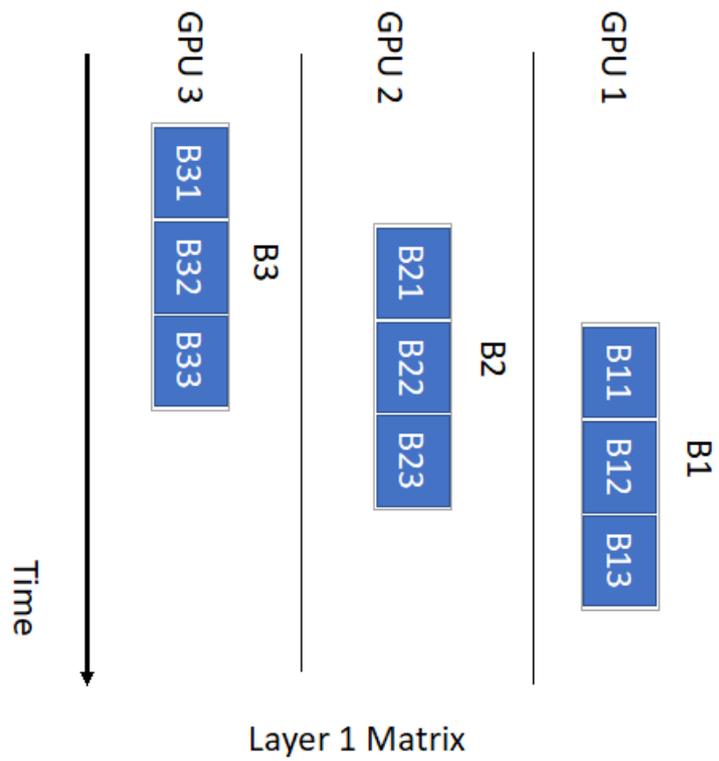












Layer 1 Matrix

$w(0,0)$	$w(1,0)$	$w(2,0)$	$w(3,0)$
$w(0,1)$	$w(1,1)$	$w(2,1)$	$w(3,1)$
$w(0,2)$	$w(1,2)$	$w(2,2)$	$w(3,2)$
$w(0,3)$	$w(1,3)$	$w(2,3)$	$w(3,3)$

Input Matrix (batch size = 4)

$x(0,0)$	$x(0,1)$	$w(0,2)$	$w(0,3)$
$x(1,0)$	$x(1,1)$	$x(1,2)$	$x(1,3)$
$x(2,0)$	$x(2,1)$	$x(2,2)$	$x(2,3)$
$x(3,0)$	$x(3,1)$	$x(3,2)$	$x(3,3)$

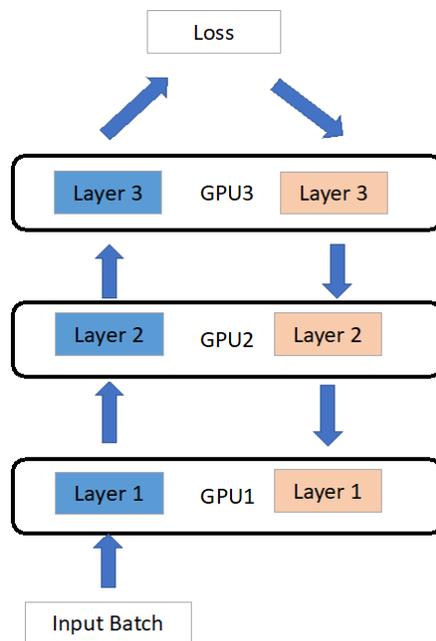
Layer 1 Matrix Splits
(Column-wise)

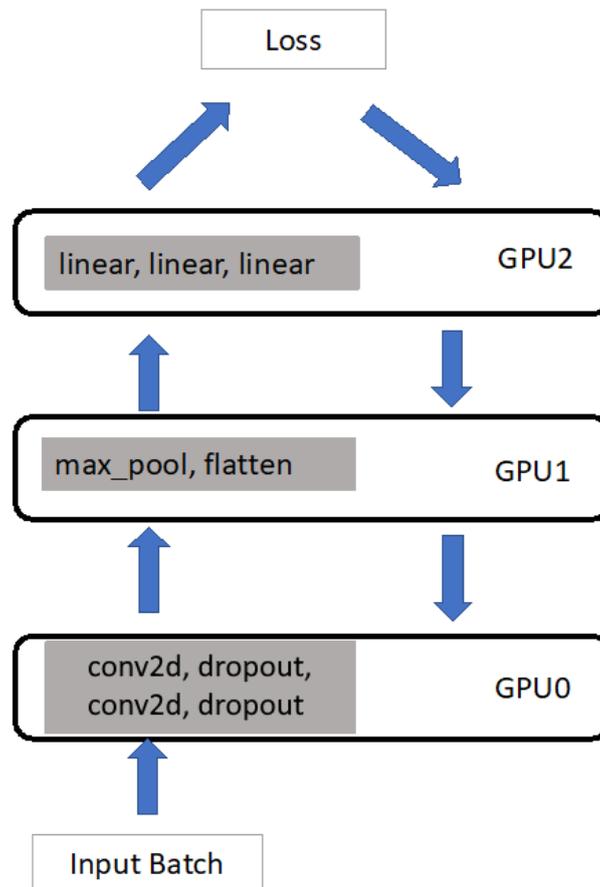
$w(0,0)$	$w(1,0)$	$w(2,0)$	$w(3,0)$
$w(0,1)$	$w(1,1)$	$w(2,1)$	$w(3,1)$
$w(0,2)$	$w(1,2)$	$w(2,2)$	$w(3,2)$
$w(0,3)$	$w(1,3)$	$w(2,3)$	$w(3,3)$

A_01

A_23

Chapter 7: Implementing Model Parallel Training and Serving Workflows





```

import torch
import torch.nn as nn
import torch.nn.functional as F

class MyNet(nn.Module):
    def __init__(self):
        super(MyNet, self).__init__()
        self.seq1 = nn.Sequential(
            nn.Conv2d(1, 32, 3, 1),
            nn.Dropout2d(0.5),
            nn.Conv2d(32, 64, 3, 1),
            nn.Dropout2d(0.75)).to('cuda:0')

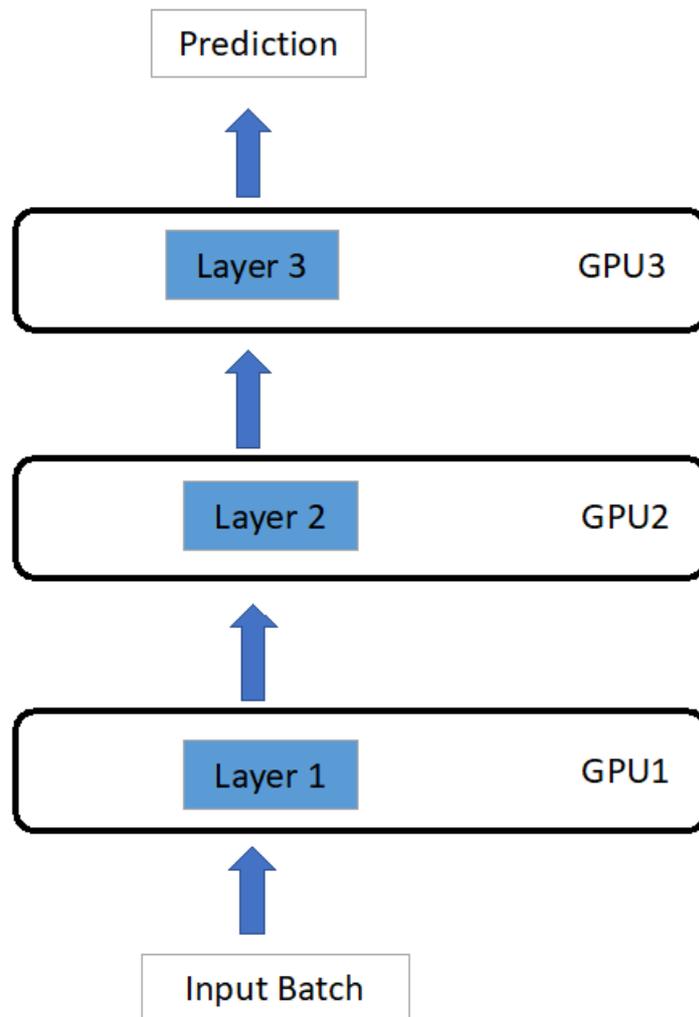
        self.seq2 = nn.Sequential(
            nn.Linear(9216, 128),
            nn.Linear(128, 20),
            nn.Linear(20, 10)).to('cuda:2')

    def forward(self, x):
        x = self.seq1(x.to('cuda:0'))
        x = F.max_pool2d(x, 2).to('cuda:1')
        x = torch.flatten(x, 1).to('cuda:1')
        x = self.seq2(x.to('cuda:2'))
        output = F.log_softmax(x, dim = 1)
        return output
  
```

Epoch 0
batch 0 training :: loss 2.367696523666382
batch 1 training :: loss 2.358067274093628
batch 2 training :: loss 2.3166911602020264
batch 3 training :: loss 2.3472657203674316
batch 4 training :: loss 2.3291213512420654
batch 5 training :: loss 2.341862201690674
batch 6 training :: loss 2.3476767539978027
batch 7 training :: loss 2.3589253425598145
batch 8 training :: loss 2.3385939598083496
batch 9 training :: loss 2.314199209213257
batch 10 training :: loss 2.357100486755371
batch 11 training :: loss 2.341332197189331
batch 12 training :: loss 2.3510727882385254
batch 13 training :: loss 2.305490732192993
batch 14 training :: loss 2.2896692752838135
batch 15 training :: loss 2.2965853214263916
batch 16 training :: loss 2.289027452468872
batch 17 training :: loss 2.318589687347412
batch 18 training :: loss 2.314786911010742
batch 19 training :: loss 2.292377471923828
batch 20 training :: loss 2.311783790588379
batch 21 training :: loss 2.3006303310394287
batch 22 training :: loss 2.2897908687591553
batch 23 training :: loss 2.309767246246338
batch 24 training :: loss 2.326434373855591
batch 25 training :: loss 2.3054540157318115
batch 26 training :: loss 2.3287947177886963
batch 27 training :: loss 2.309558391571045
batch 28 training :: loss 2.289318084716797
batch 29 training :: loss 2.3383259773254395
batch 30 training :: loss 2.2959561347961426
batch 31 training :: loss 2.2574143409729004
batch 32 training :: loss 2.293168783187866
batch 33 training :: loss 2.2815194129943848
batch 34 training :: loss 2.2899670600891113
batch 35 training :: loss 2.2440366744995117
batch 36 training :: loss 2.2733407020568848
batch 37 training :: loss 2.2578611373901367
batch 38 training :: loss 2.2523033618927
batch 39 training :: loss 2.2961771488189697
batch 40 training :: loss 2.269951820373535

batch 450 training :: loss 0.4500477612018585
batch 451 training :: loss 0.41694992780685425
batch 452 training :: loss 0.5335432291030884
batch 453 training :: loss 0.3785797357559204
batch 454 training :: loss 0.5250097513198853
batch 455 training :: loss 0.48590853810310364
batch 456 training :: loss 0.4359087646007538
batch 457 training :: loss 0.5516181588172913
batch 458 training :: loss 0.4193853735923767
batch 459 training :: loss 0.24893827736377716
batch 460 training :: loss 0.4412848949432373
batch 461 training :: loss 0.6855418086051941
batch 462 training :: loss 0.60863196849823
batch 463 training :: loss 0.6327939629554749
batch 464 training :: loss 0.4109138548374176
batch 465 training :: loss 0.3921489715576172
batch 466 training :: loss 0.3610058128833771
batch 467 training :: loss 0.4983370304107666
batch 468 training :: loss 0.497358113527298
Training Done!

NVIDIA-SMI 450.142.00 Driver Version: 450.142.00 CUDA Version: 11.0							
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC	
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute M.	MIG M.
0	Tesla M60	On	00000000:00:1B.0	Off		0	
N/A	32C P0	82W / 150W	1008MiB / 7618MiB		69%	Default	N/A
1	Tesla M60	On	00000000:00:1C.0	Off		0	
N/A	25C P0	38W / 150W	708MiB / 7618MiB		18%	Default	N/A
2	Tesla M60	On	00000000:00:1D.0	Off		0	
N/A	29C P0	39W / 150W	772MiB / 7618MiB		21%	Default	N/A
3	Tesla M60	On	00000000:00:1E.0	Off		0	
N/A	23C P8	14W / 150W	3MiB / 7618MiB		0%	Default	N/A
Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory	
ID	ID	ID				Usage	
0	N/A	N/A	722	C	python	1005MiB	
1	N/A	N/A	722	C	python	705MiB	
2	N/A	N/A	722	C	python	769MiB	



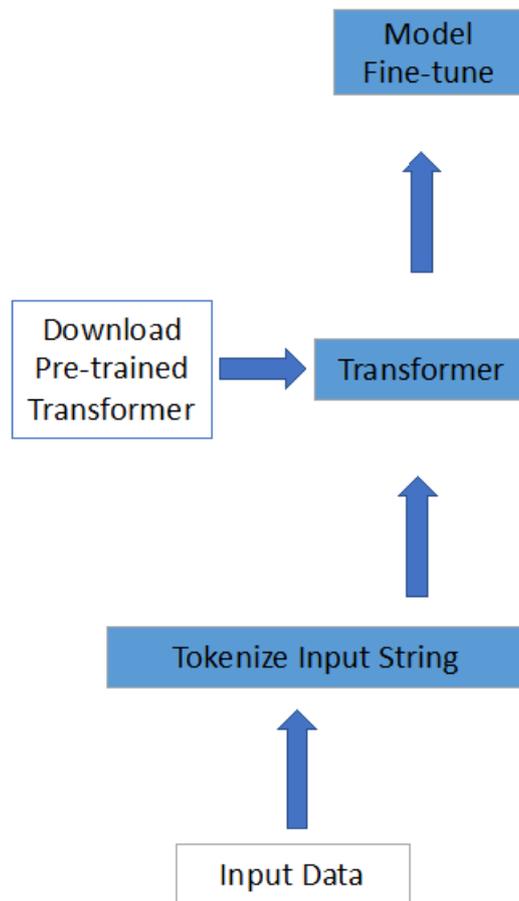
```
Test Accuracy 0.0019666666666666665
Test Accuracy 0.0037333333333333333
Test Accuracy 0.0056833333333333334
Test Accuracy 0.0076
Test Accuracy 0.0094833333333333333
Test Accuracy 0.0114
Test Accuracy 0.0132166666666666666
Test Accuracy 0.0150666666666666667
Test Accuracy 0.0168
Test Accuracy 0.01875
Test Accuracy 0.0205666666666666667
Test Accuracy 0.0225
Test Accuracy 0.0243
Test Accuracy 0.0263333333333333334
```

```

Test Accuracy 0.8768166666666667
Test Accuracy 0.8786333333333334
Test Accuracy 0.88055
Test Accuracy 0.8826333333333334
Test Accuracy 0.8844666666666666
Test Accuracy 0.8865666666666666
Test Accuracy 0.88795
Test Done!

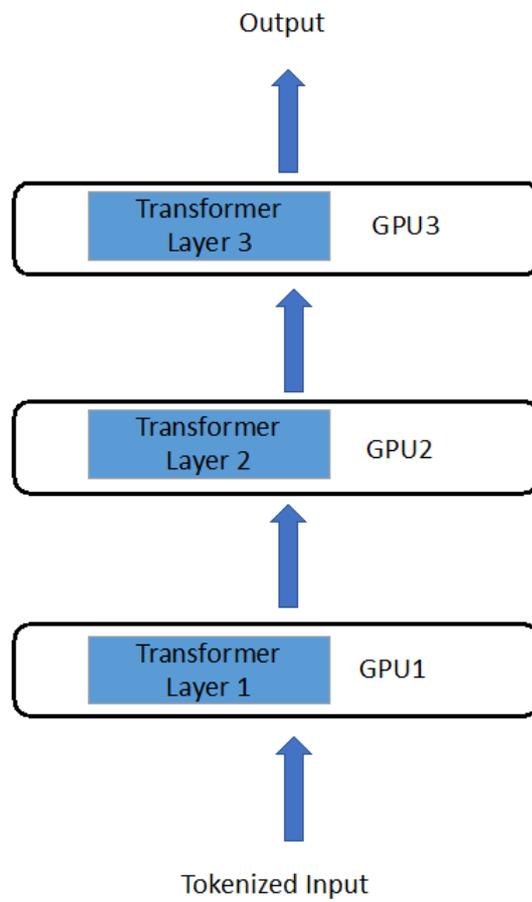
```

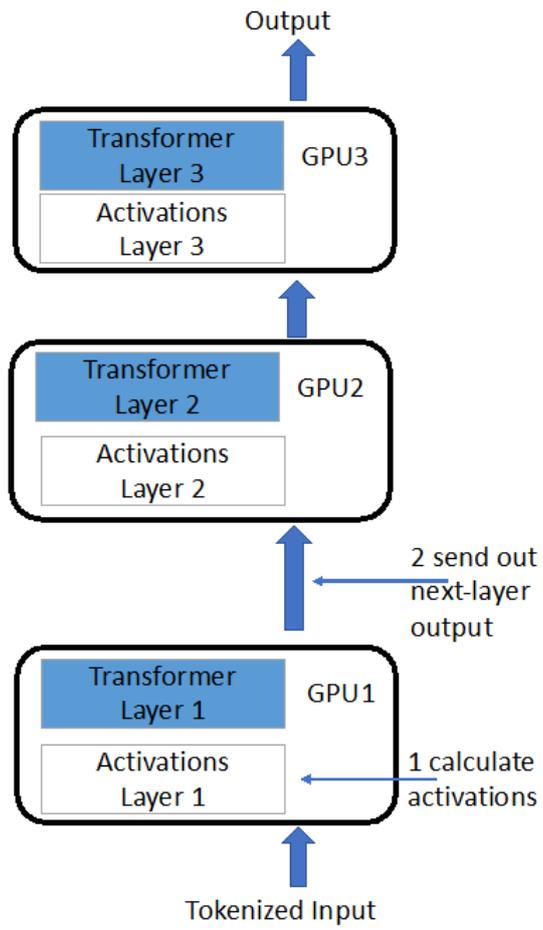
NVIDIA-SMI 450.142.00 Driver Version: 450.142.00 CUDA Version: 11.0							
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr.	ECC
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute	M. MIG M.
0	Tesla M60	On	00000000:00:1B.0	Off	42%	Default	0
N/A	38C	P0	78W / 150W	1008MiB / 7618MiB		N/A	
1	Tesla M60	On	00000000:00:1C.0	Off	20%	Default	0
N/A	27C	P0	38W / 150W	708MiB / 7618MiB		N/A	
2	Tesla M60	On	00000000:00:1D.0	Off	21%	Default	0
N/A	30C	P0	39W / 150W	772MiB / 7618MiB		N/A	
3	Tesla M60	On	00000000:00:1E.0	Off	0%	Default	0
N/A	23C	P8	14W / 150W	3MiB / 7618MiB		N/A	
Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory	Usage
	ID	ID					
0	N/A	N/A	10485	C	python	1005MiB	
1	N/A	N/A	10485	C	python	705MiB	
2	N/A	N/A	10485	C	python	769MiB	

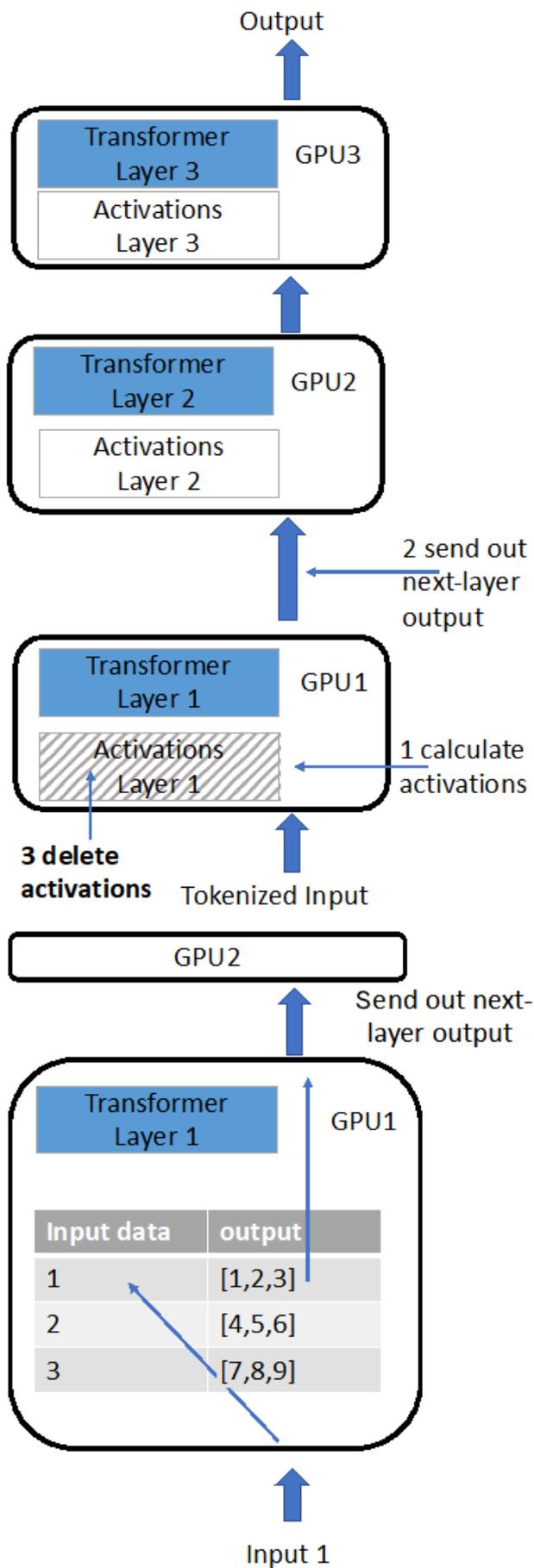


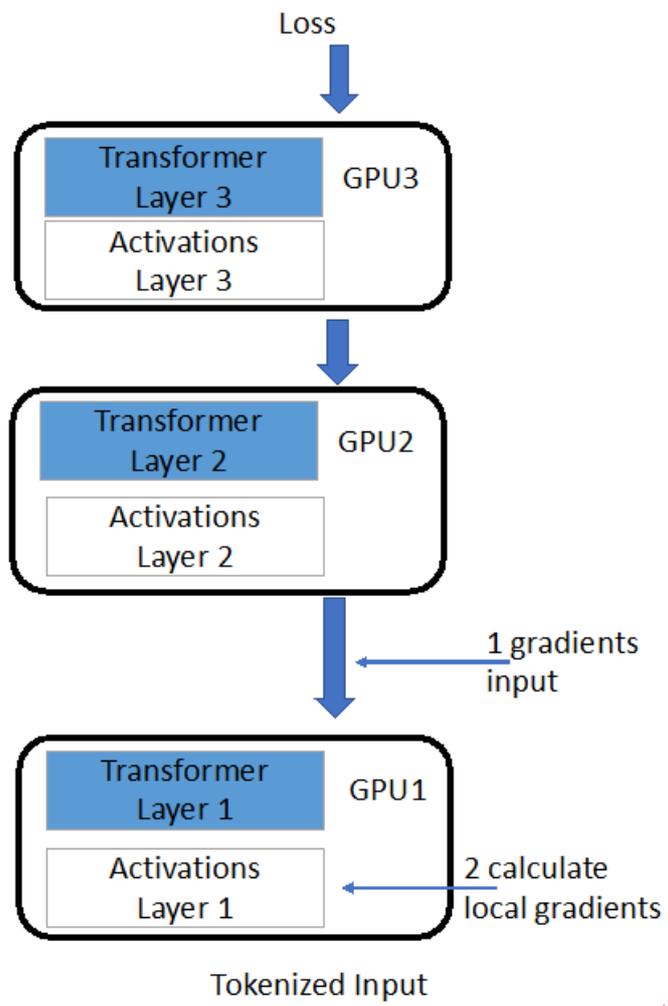
GPU Fan	Name Temp	Perf	Persistence-M Pwr:Usage/Cap	Bus-Id	Disp.A Memory-Usage	Volatile GPU-Util	Uncorr. Compute	ECC M. MIG M.
0 N/A	Tesla 38C	M60 P0	On 78W / 150W	00000000:00:1B.0	Off 1008MiB / 7618MiB	42%	Default	0 N/A
1 N/A	Tesla 27C	M60 P0	On 38W / 150W	00000000:00:1C.0	Off 708MiB / 7618MiB	20%	Default	0 N/A

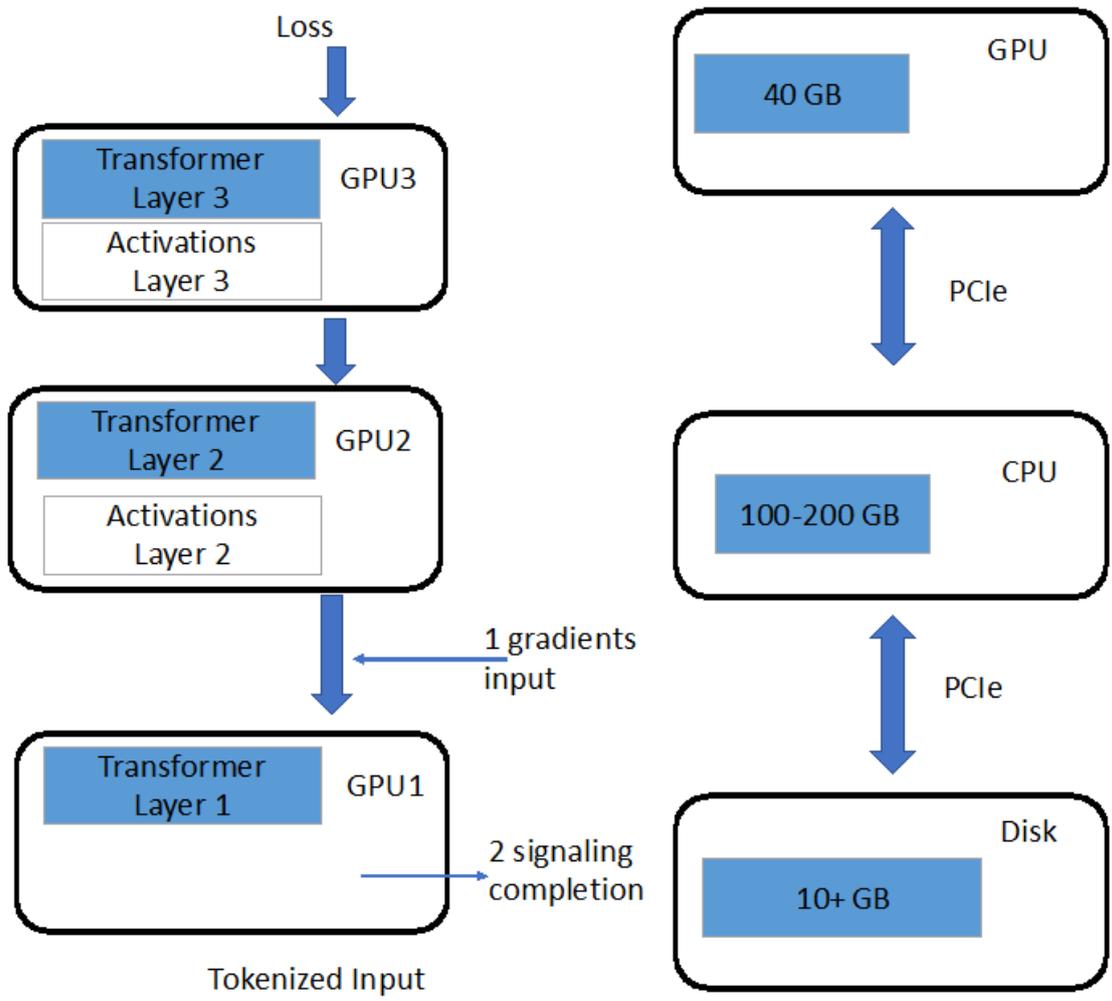
Chapter 8: Achieving Higher Throughput and Lower Latency

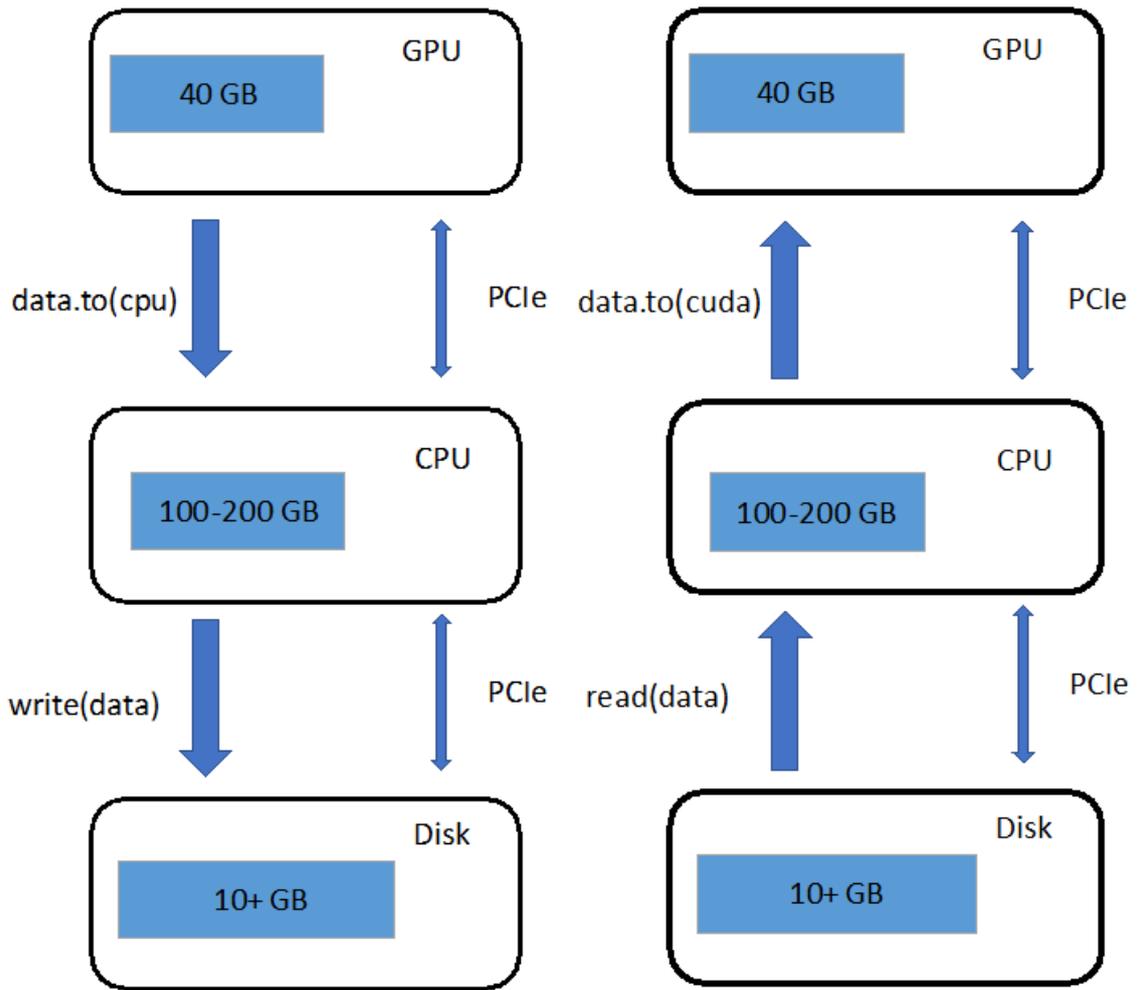


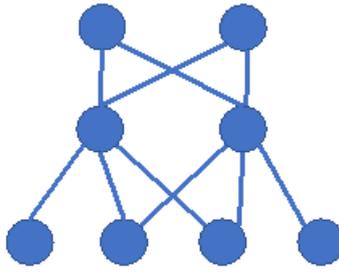




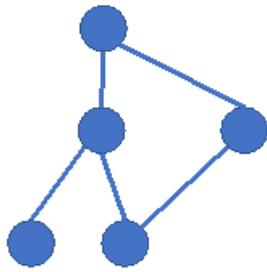




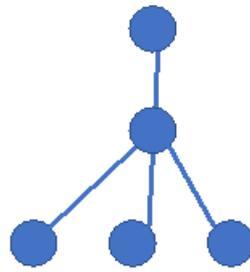




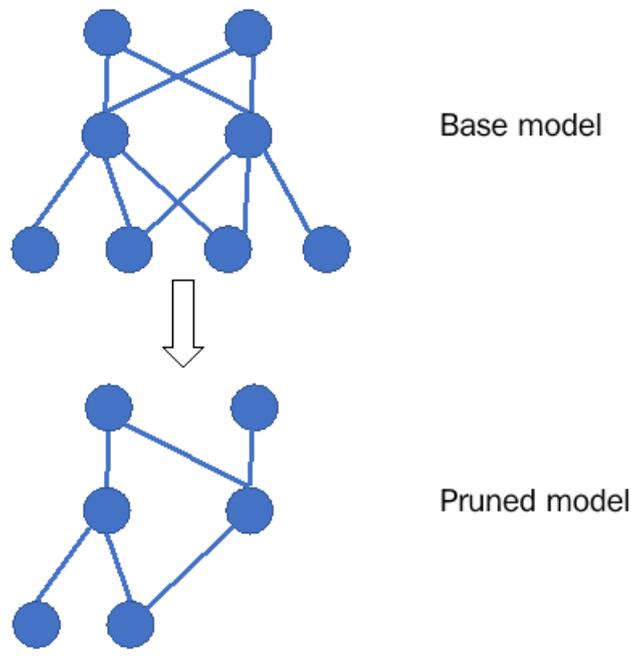
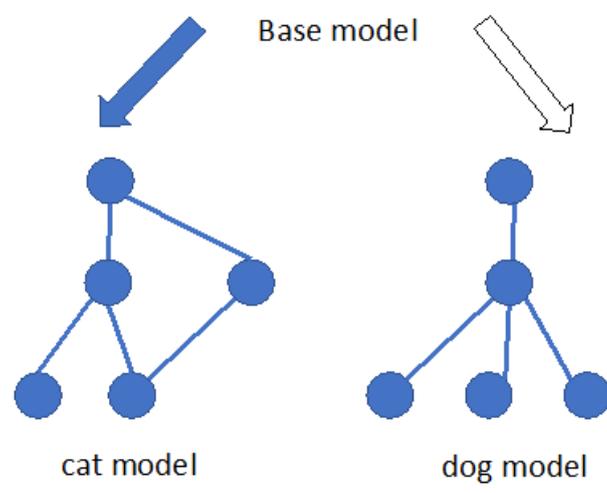
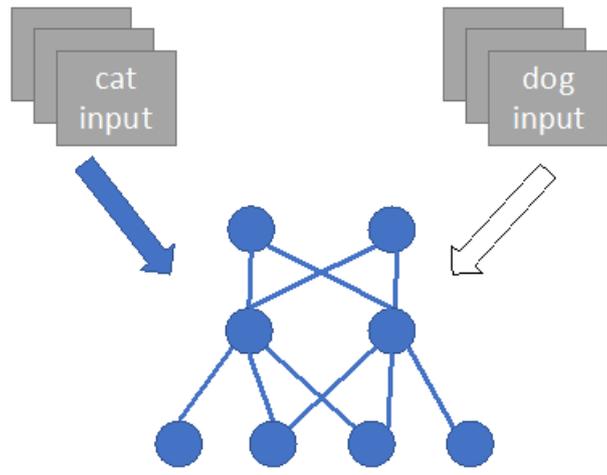
Base model



cat model



dog model



Representation of model weight

Normal

Reduced

Further reduced



Chapter 9: A Hybrid of Data and Model Parallelism

A(0,0)	A(0,1)	A(0,2)	A(0,3)
A(1,0)	A(1,1)	A(1,2)	A(1,3)
A(2,0)	A(2,1)	A(2,2)	A(2,3)
A(3,0)	A(3,1)	A(3,2)	A(3,3)

Matrix A

(weights)

$X * A$



Model Layer

Weight Matrix A



Input Matrix X

$X(0,0)$	$X(0,1)$	$X(0,2)$	$X(0,3)$
$X(1,0)$	$X(1,1)$	$X(1,2)$	$X(1,3)$
$X(2,0)$	$X(2,1)$	$X(2,2)$	$X(2,3)$
$X(3,0)$	$X(3,1)$	$X(3,2)$	$X(3,3)$

Matrix X

(Input)

$A(0,0)$	$A(0,1)$	$A(0,2)$	$A(0,3)$
$A(1,0)$	$A(1,1)$	$A(1,2)$	$A(1,3)$

A[0]

$A(2,0)$	$A(2,1)$	$A(2,2)$	$A(2,3)$
$A(3,0)$	$A(3,1)$	$A(3,2)$	$A(3,3)$

A[1]

Matrix A
(weights)

$X(0,0)$	$X(0,1)$
$X(1,0)$	$X(1,1)$
$X(2,0)$	$X(2,1)$
$X(3,0)$	$X(3,1)$

$X[0]$

$X(0,2)$	$X(0,3)$
$X(1,2)$	$X(1,3)$
$X(2,2)$	$X(2,3)$
$X(3,2)$	$X(3,3)$

$X[1]$

Matrix X

(Input)

$A(0,0)$	$A(0,1)$	$A(0,2)$	$A(0,3)$
$A(1,0)$	$A(1,1)$	$A(1,2)$	$A(1,3)$

$A[0]$

$X(0,0)$	$X(0,1)$
$X(1,0)$	$X(1,1)$
$X(2,0)$	$X(2,1)$
$X(3,0)$	$X(3,1)$

$X[0]$

$X[0]*A[0]$

A(2,0)	A(2,1)	A(2,2)	A(2,3)
A(3,0)	A(3,1)	A(3,2)	A(3,3)

A[1]

X(0,2)	X(0,3)
X(1,2)	X(1,3)
X(2,2)	X(2,3)
X(3,2)	X(3,3)

X[1]*A[1]

X[1]

A(0,0)	A(0,1)
A(1,0)	A(1,1)
A(2,0)	A(2,1)
A(3,0)	A(3,1)

A'[0]

A(0,2)	A(0,3)
A(1,2)	A(1,3)
A(2,2)	A(2,3)
A(3,2)	A(3,3)

A'[1]

Matrix A

(weights)

$X(0,0)$	$X(0,1)$	$X(0,2)$	$X(0,3)$	
$X(1,0)$	$X(1,1)$	$X(1,2)$	$X(1,3)$	X
$X(2,0)$	$X(2,1)$	$X(2,2)$	$X(2,3)$	
$X(3,0)$	$X(3,1)$	$X(3,2)$	$X(3,3)$	

$A(0,0)$	$A(0,1)$
$A(1,0)$	$A(1,1)$
$A(2,0)$	$A(2,1)$
$A(3,0)$	$A(3,1)$

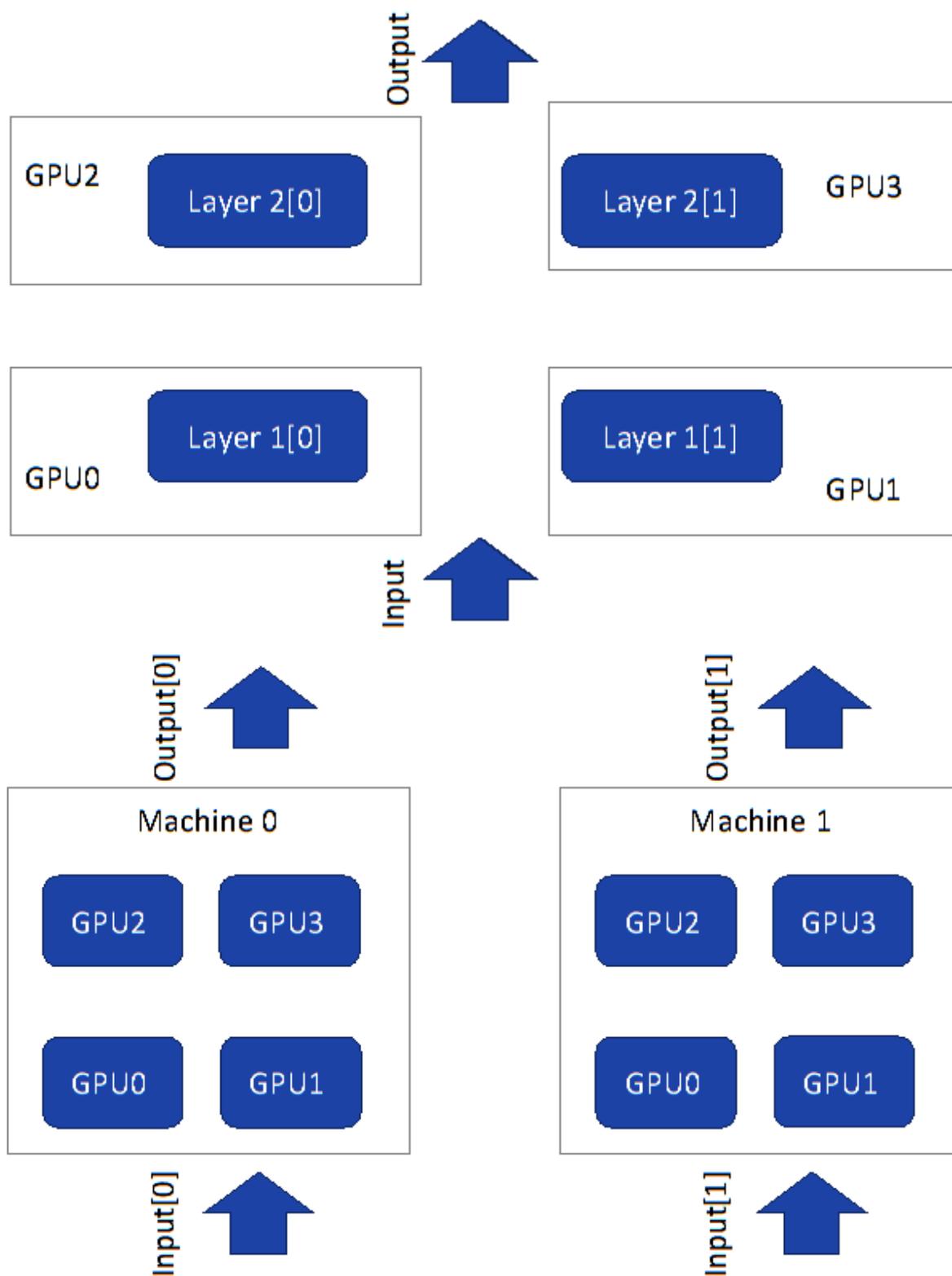
$A'[0]$

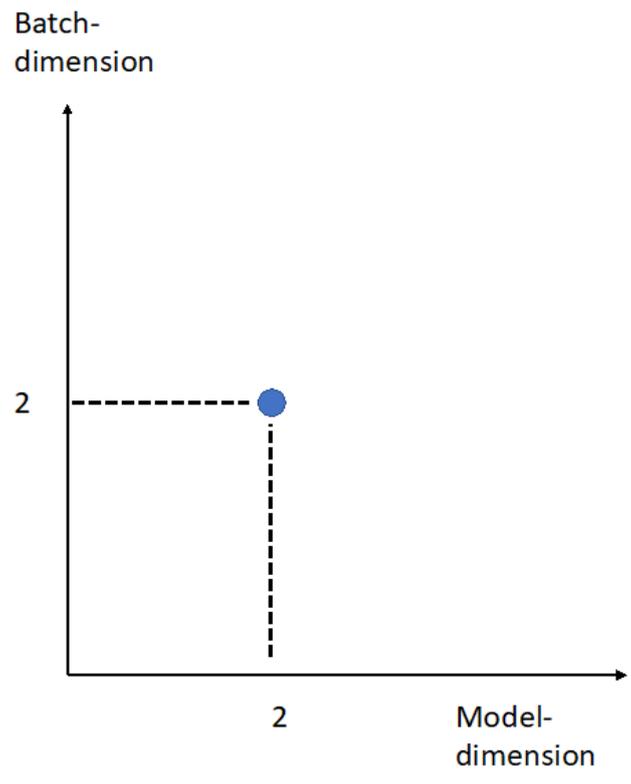
$X(0,0)$	$X(0,1)$	$X(0,2)$	$X(0,3)$
$X(1,0)$	$X(1,1)$	$X(1,2)$	$X(1,3)$
$X(2,0)$	$X(2,1)$	$X(2,2)$	$X(2,3)$
$X(3,0)$	$X(3,1)$	$X(3,2)$	$X(3,3)$

X

$A(0,2)$	$A(0,3)$
$A(1,2)$	$A(1,3)$
$A(2,2)$	$A(2,3)$
$A(3,2)$	$A(3,3)$

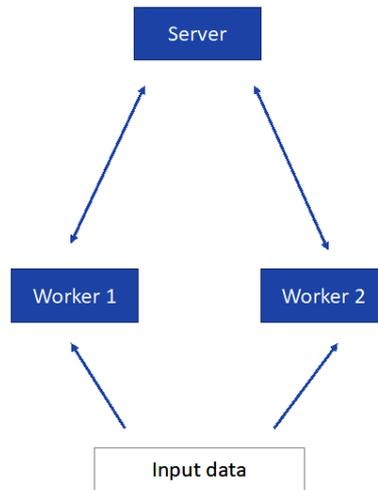
$A'[1]$



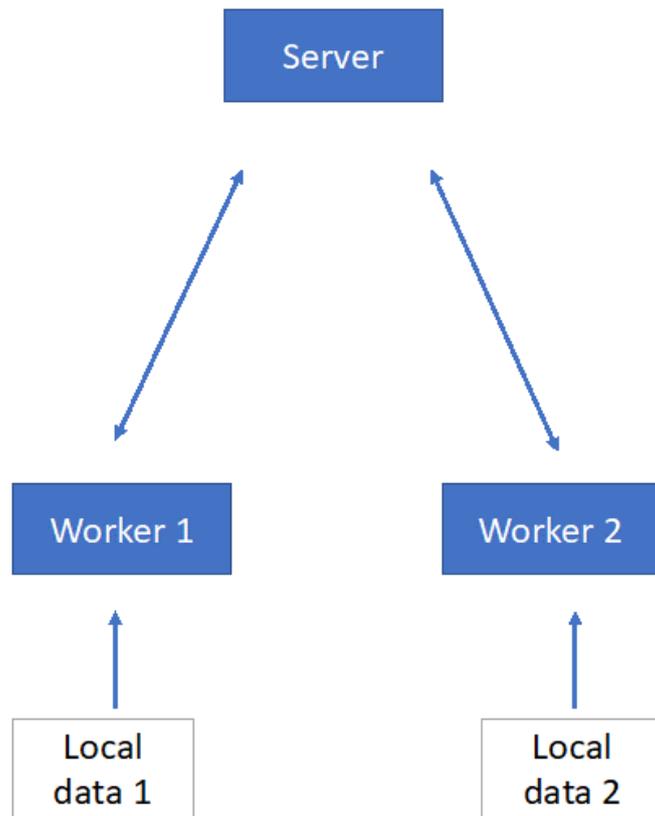


Mesh-2D

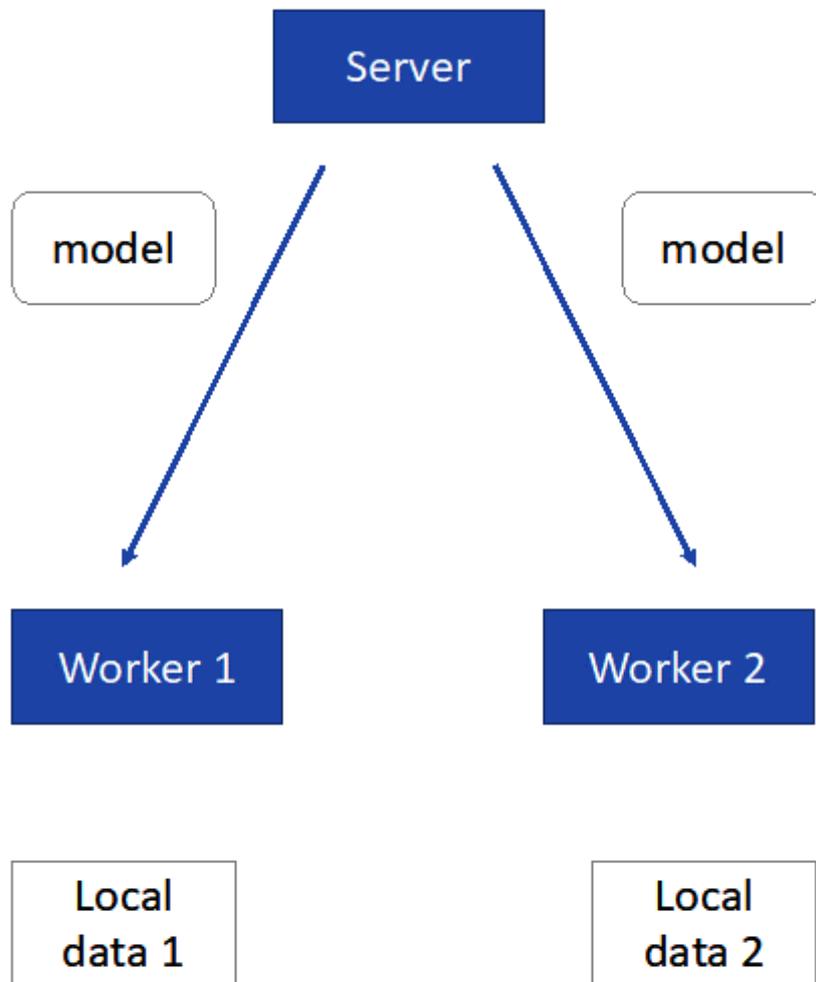
Chapter 10: Federated Learning and Edge Devices



Data parallel training
(2 workers and 1 server)

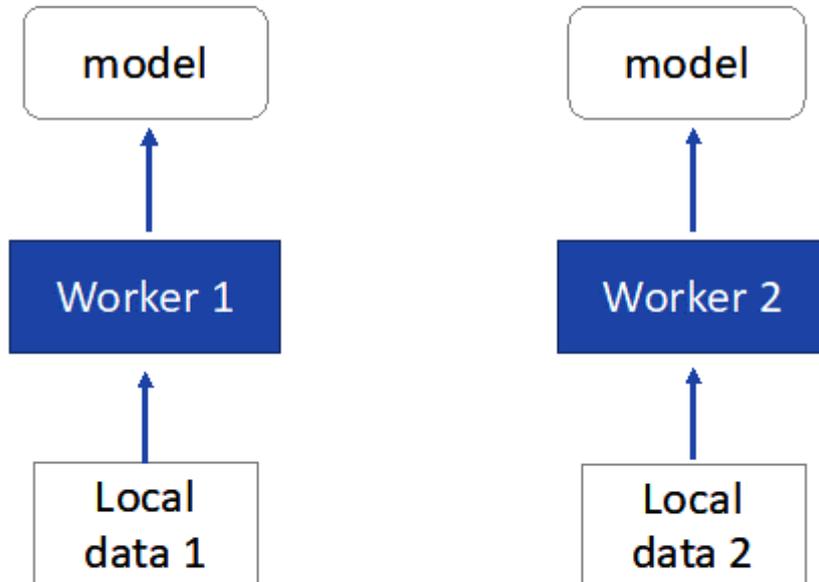


Federated Learning
(2 workers and 1 server)

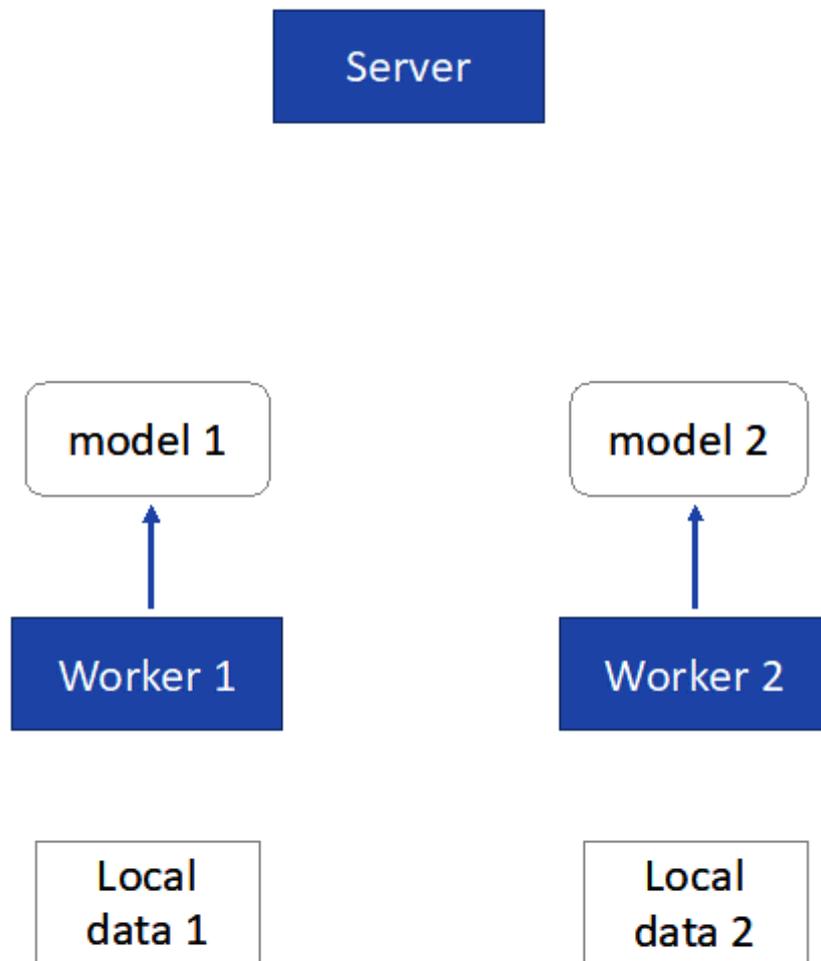


Federated Learning
(2 workers and 1 server)

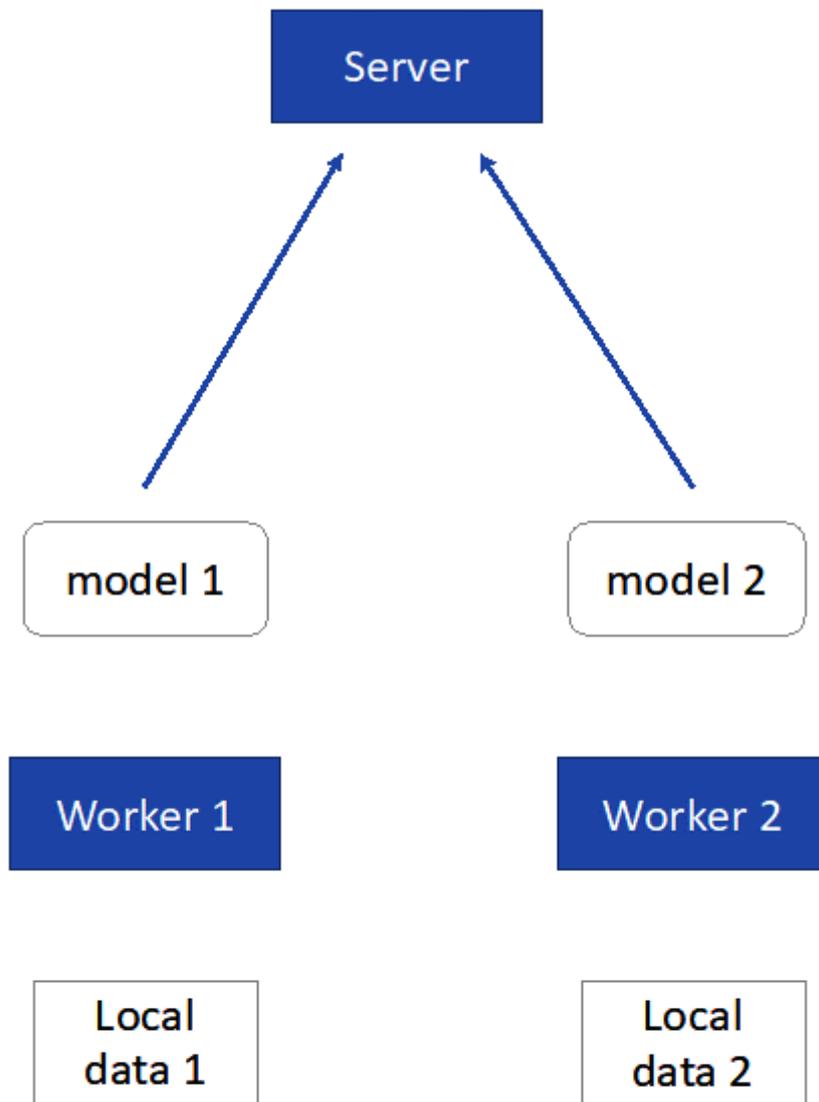
Server



Federated Learning
(2 workers and 1 server)



Federated Learning
(2 workers and 1 server)



Federated Learning
(2 workers and 1 server)

Federated Learning (FL)
APIs



Federated Core (FC)
APIs

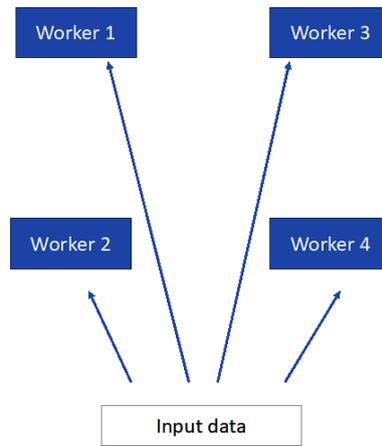


Distributed TensorFlow

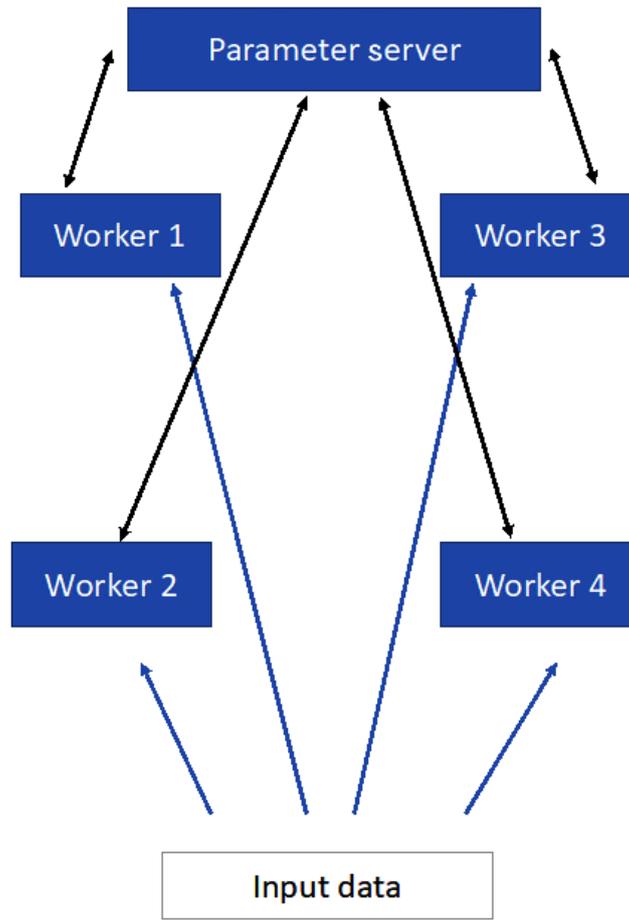
Collective Communications

TensorFlow Federated
(two top layers)

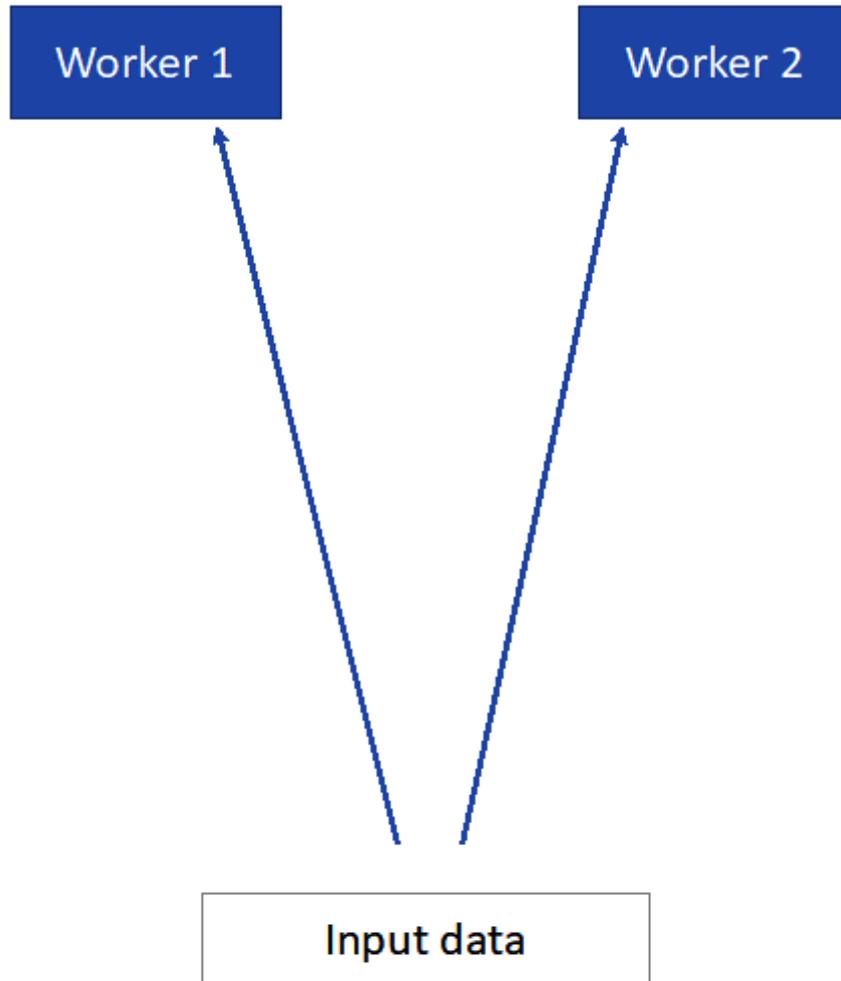
Chapter 11: Elastic Model Training and Serving



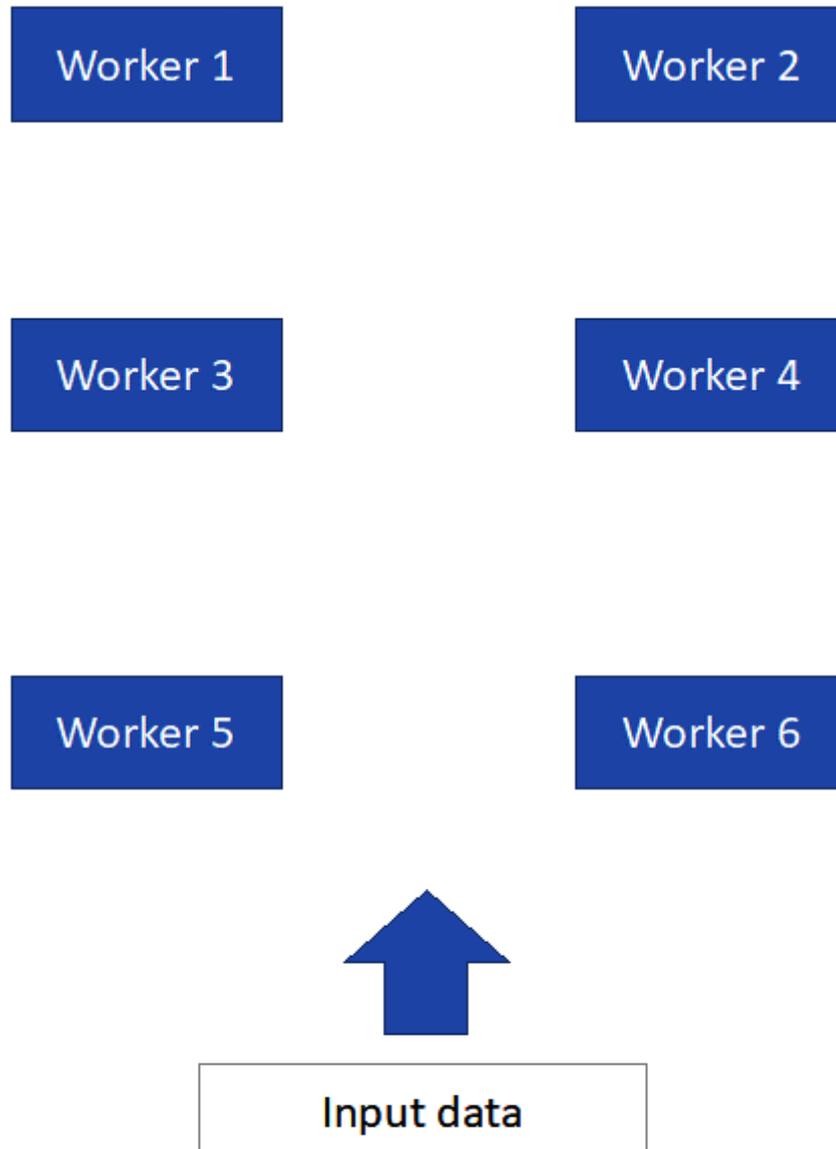
Data-parallel training
(AllReduce-based)



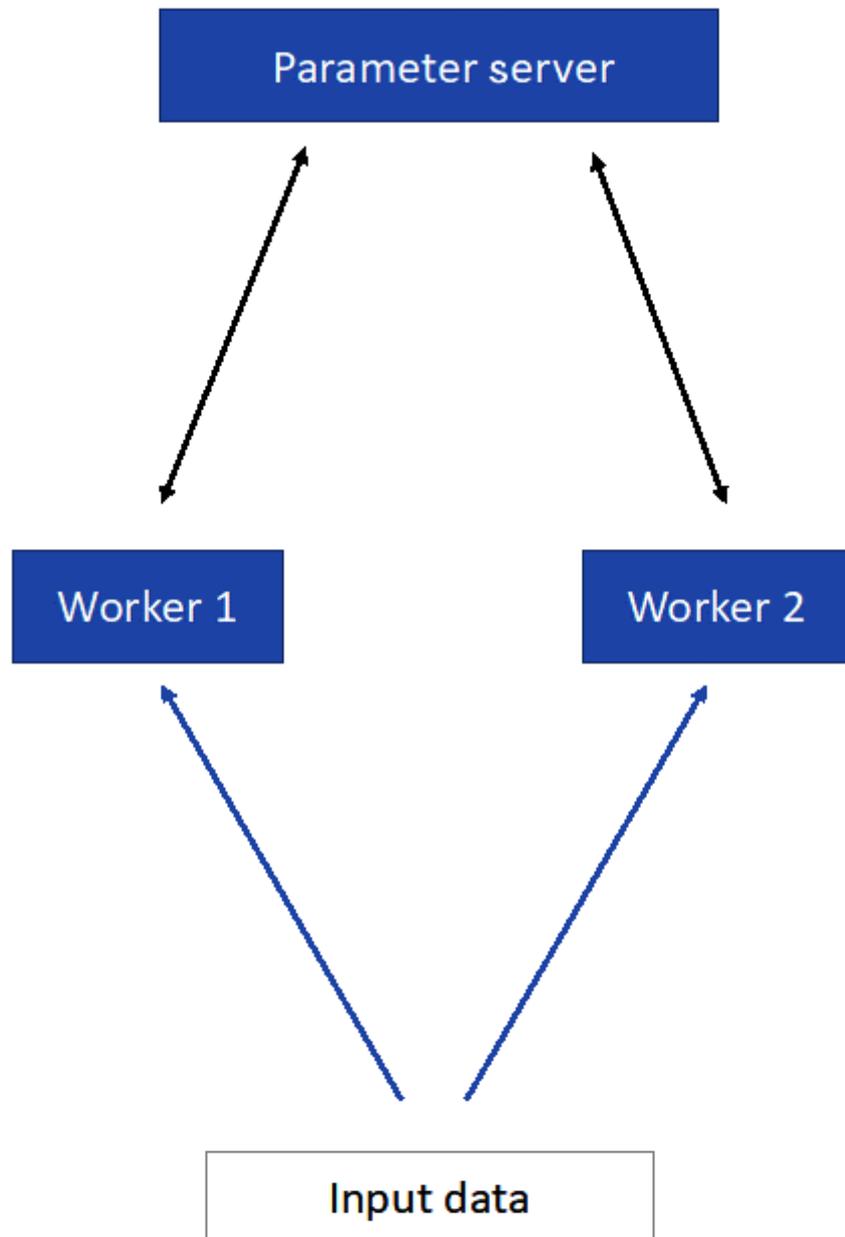
Data-parallel training
(Parameter server-based)



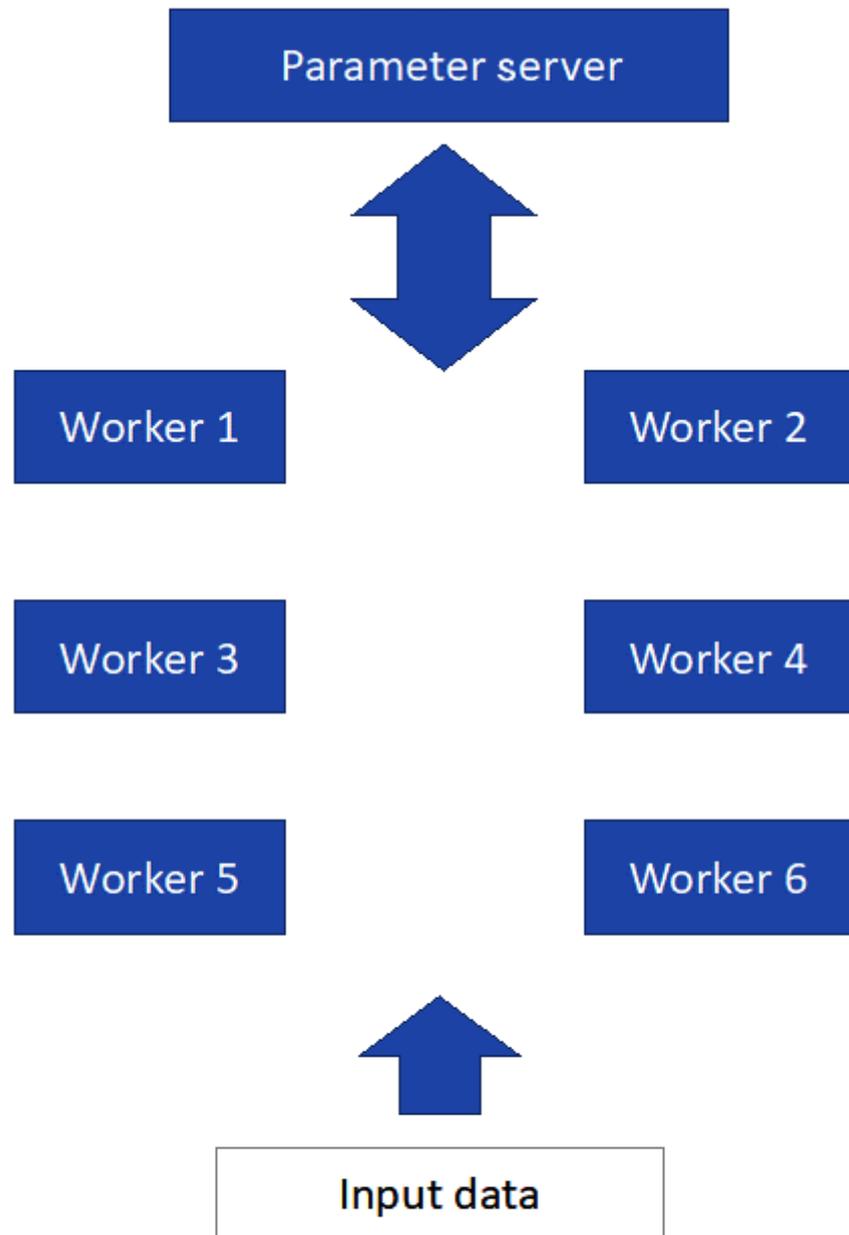
Adaptive data-parallel training
(AllReduce-based, early stage)



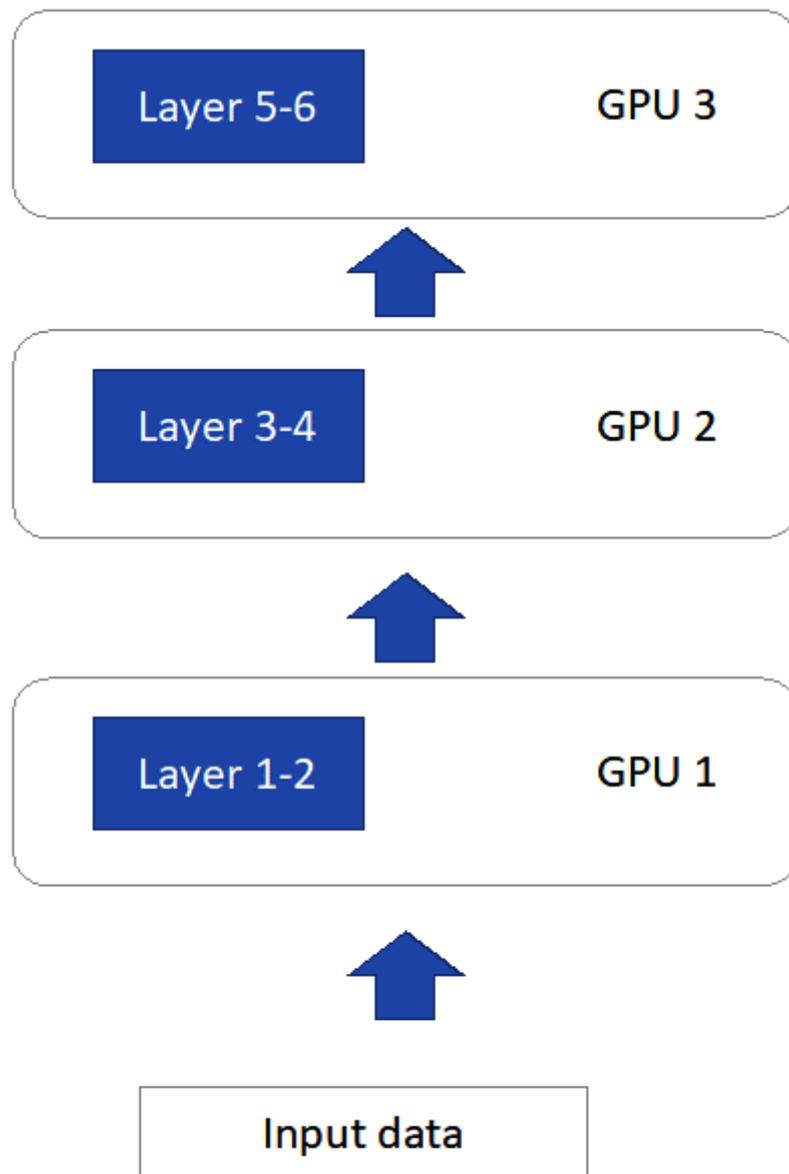
**Adaptive data-parallel training
(AllReduce-based, late stage)**



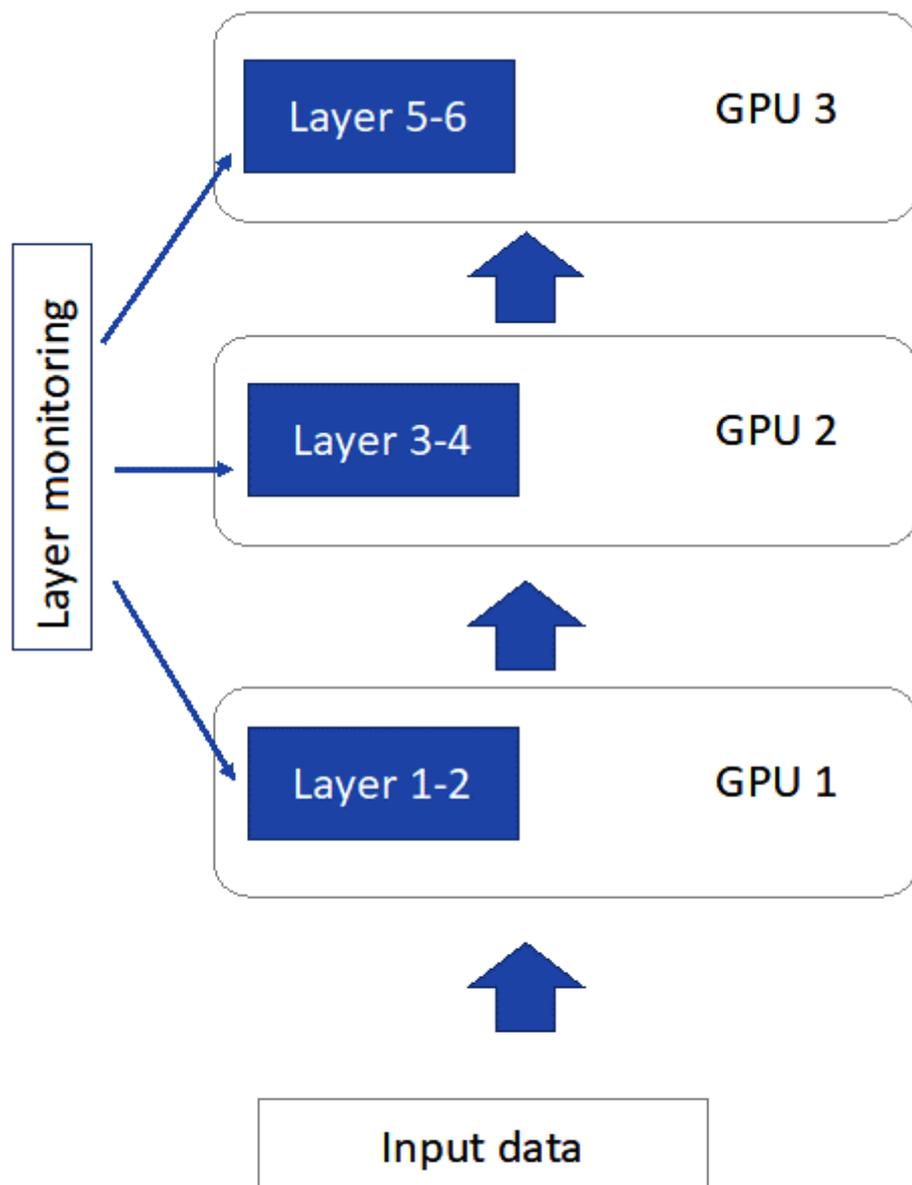
Data-parallel training
(Parameter server-based, early stage)



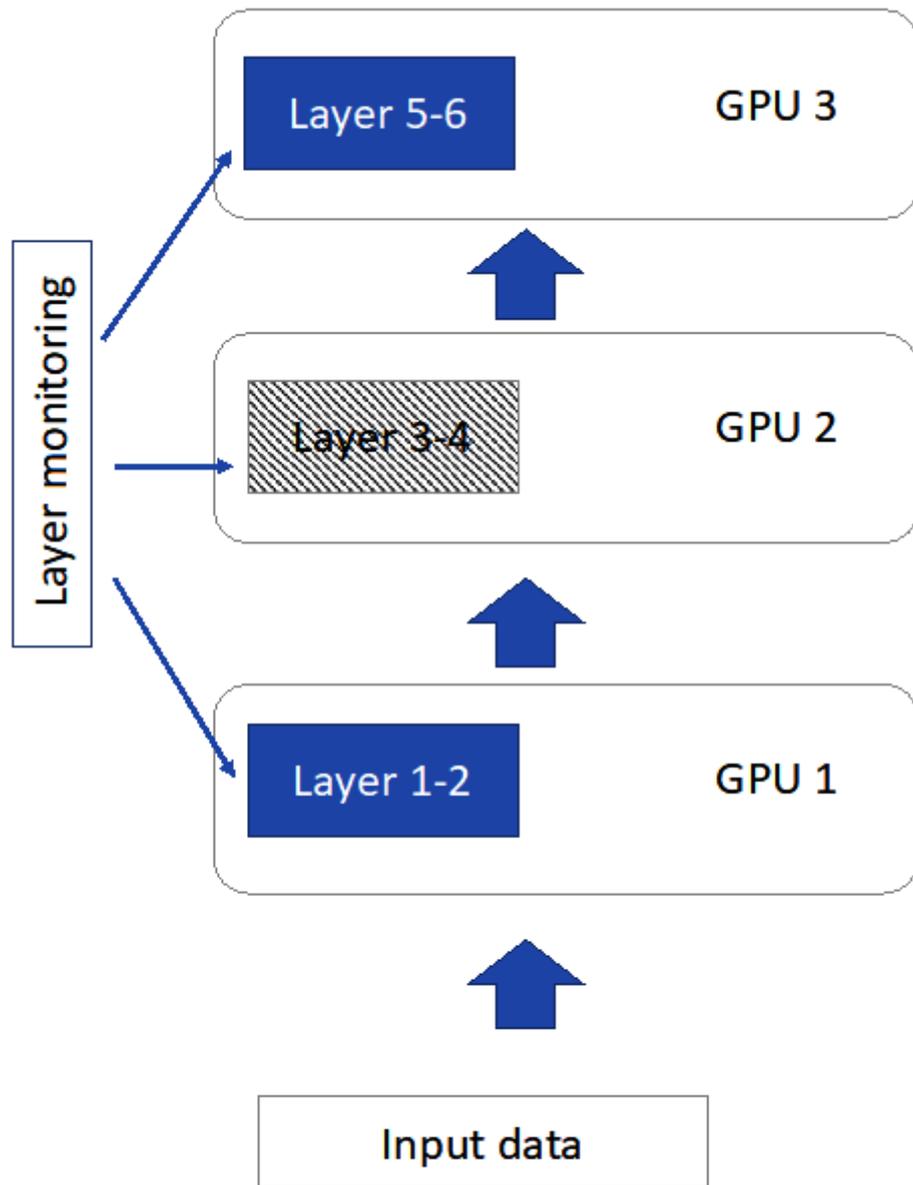
Data-parallel training
(Parameter server-based, late stage)



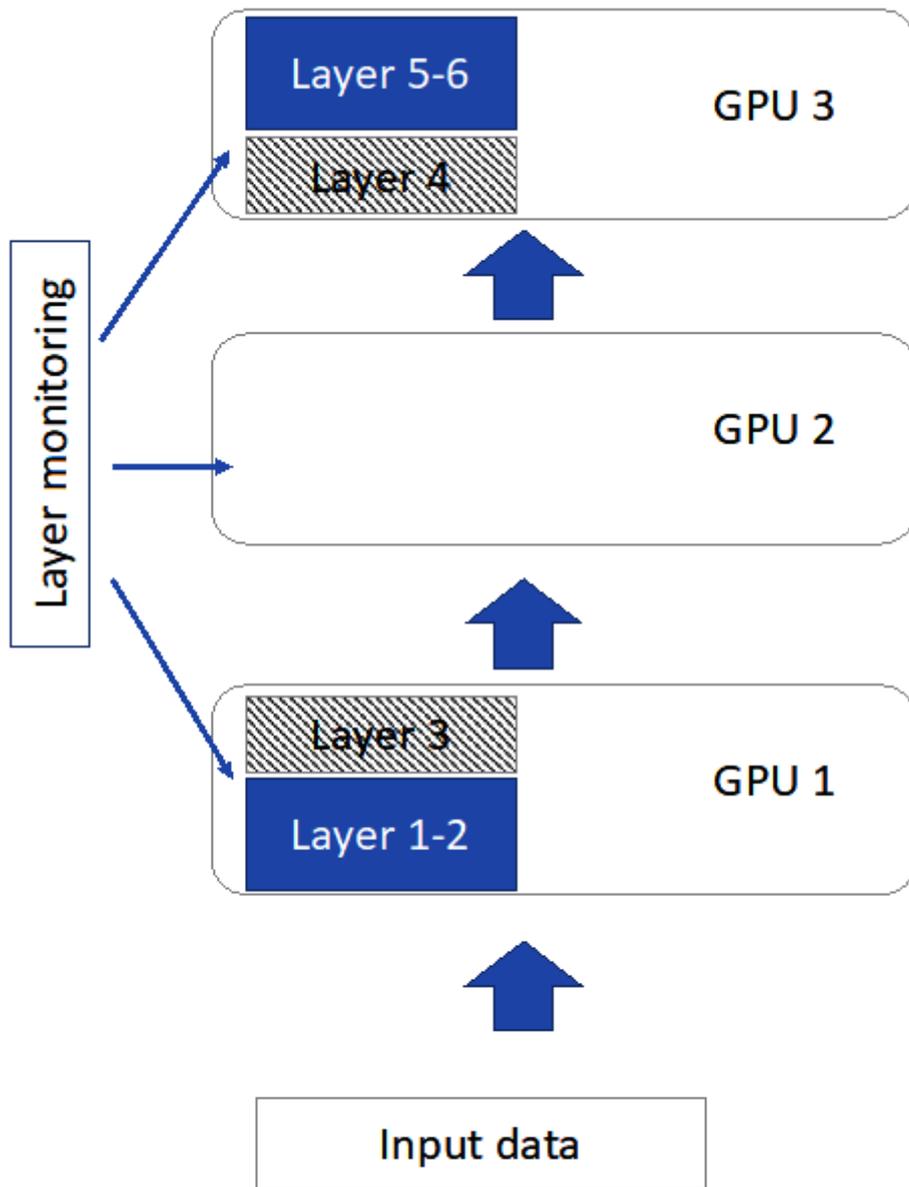
Model-parallel Training



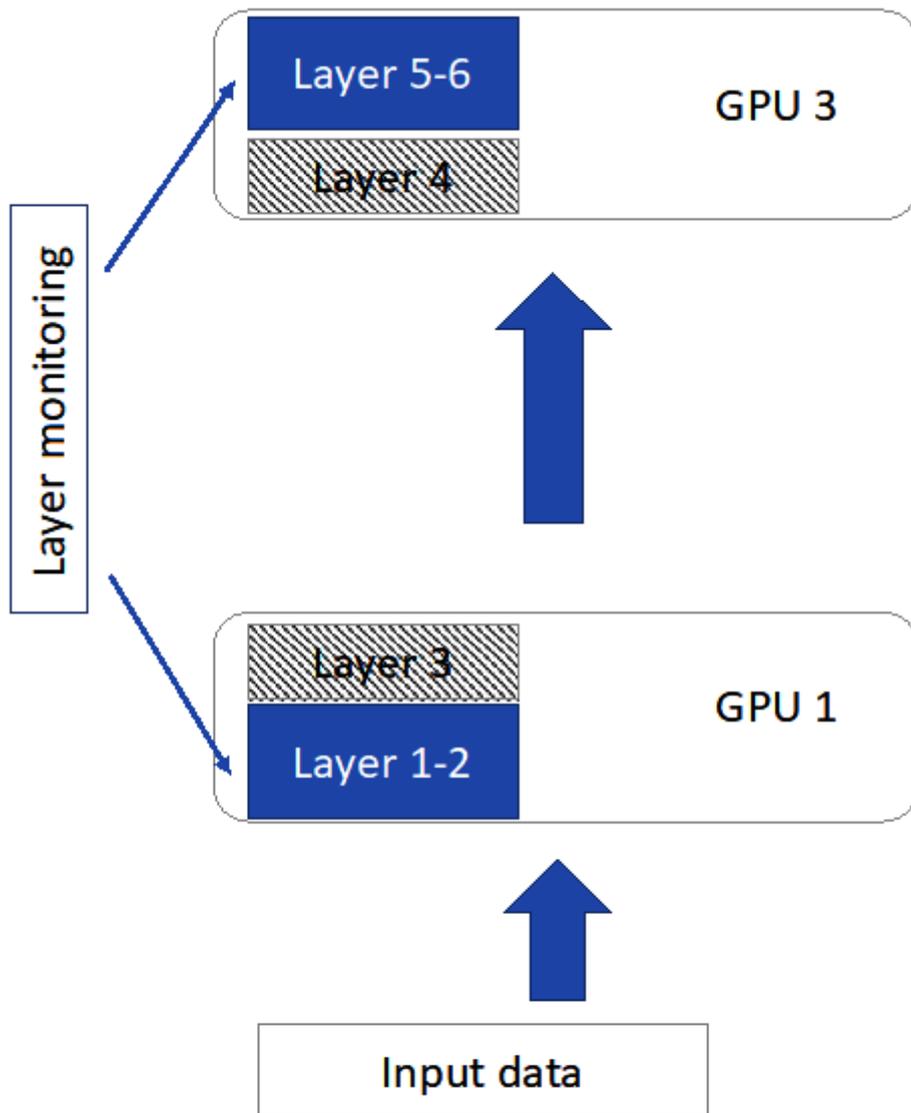
Adaptive model-parallel training
(early stage)



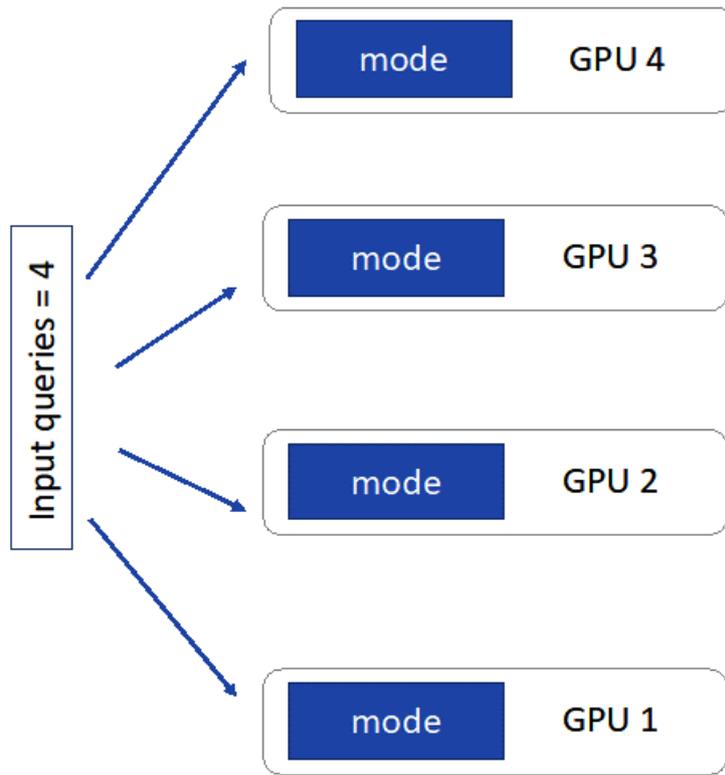
Adaptive model-parallel training
(after early stage)



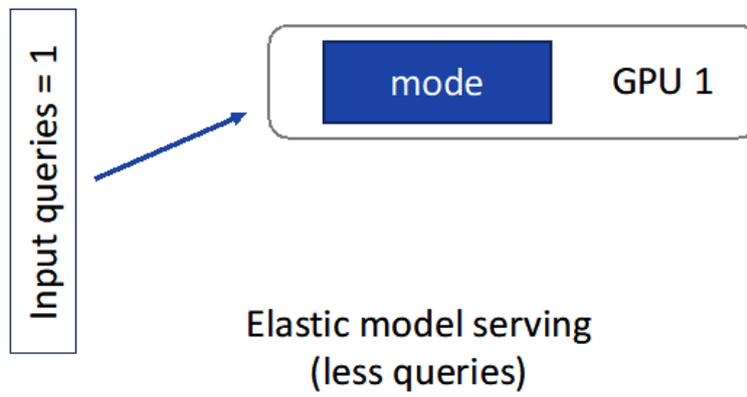
Adaptive model-parallel training
(before late stage)



Adaptive model-parallel training
(late stage)



Elastic model serving
(more queries)



Chapter 12: Advanced Techniques for Further Speed-Ups

