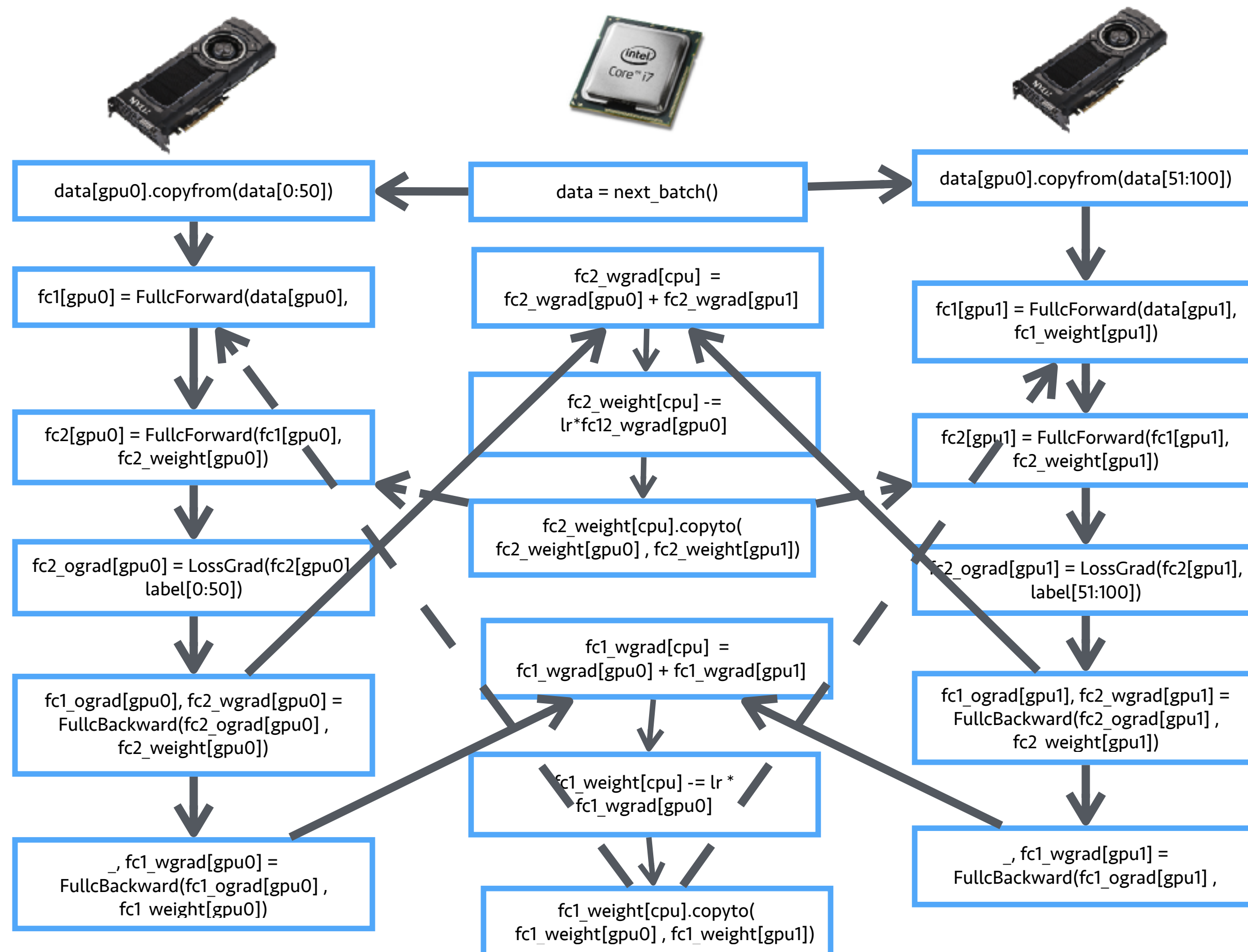# Writing Parallel Programs is Painful

2-layer neural networks with 2 GPUs
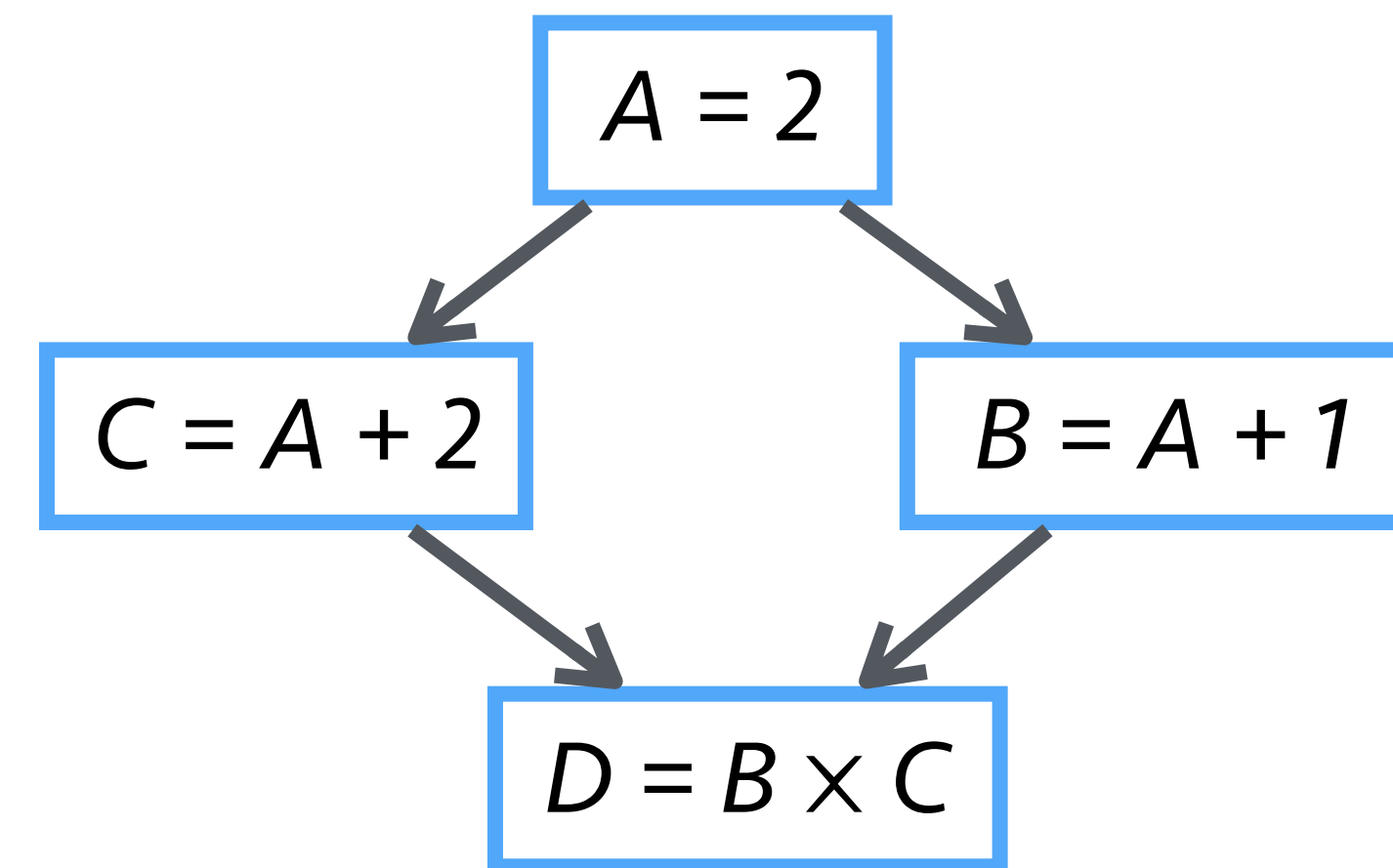


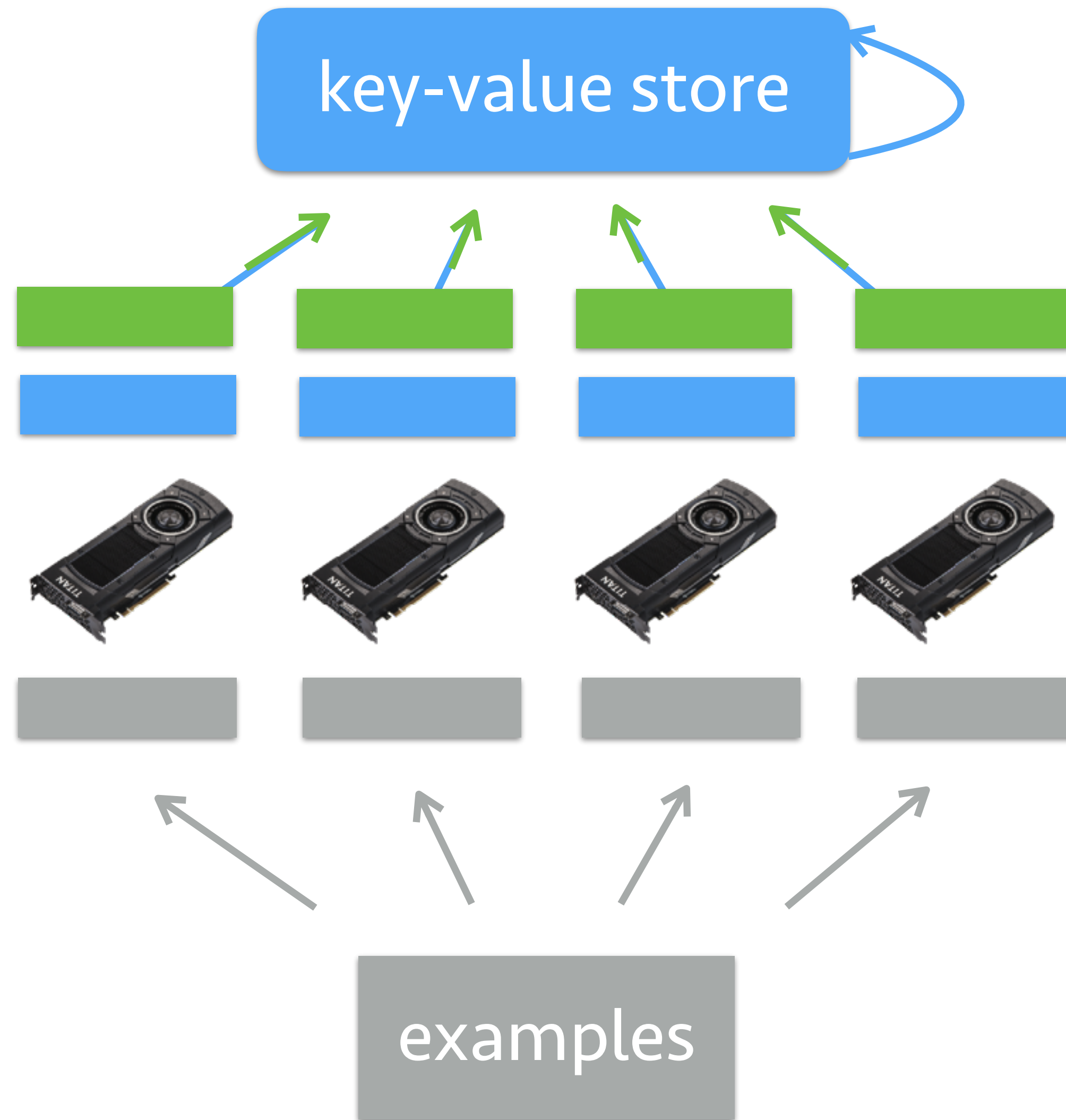A network may have
hundreds of layers

# Auto Parallelization

Write **serial** programs

```
>>> import mxnet as mx
>>> A = mx.nd.ones((2,2)) *2
>>> C = A + 2
>>> B = A + 1
>>> D = B * C
```
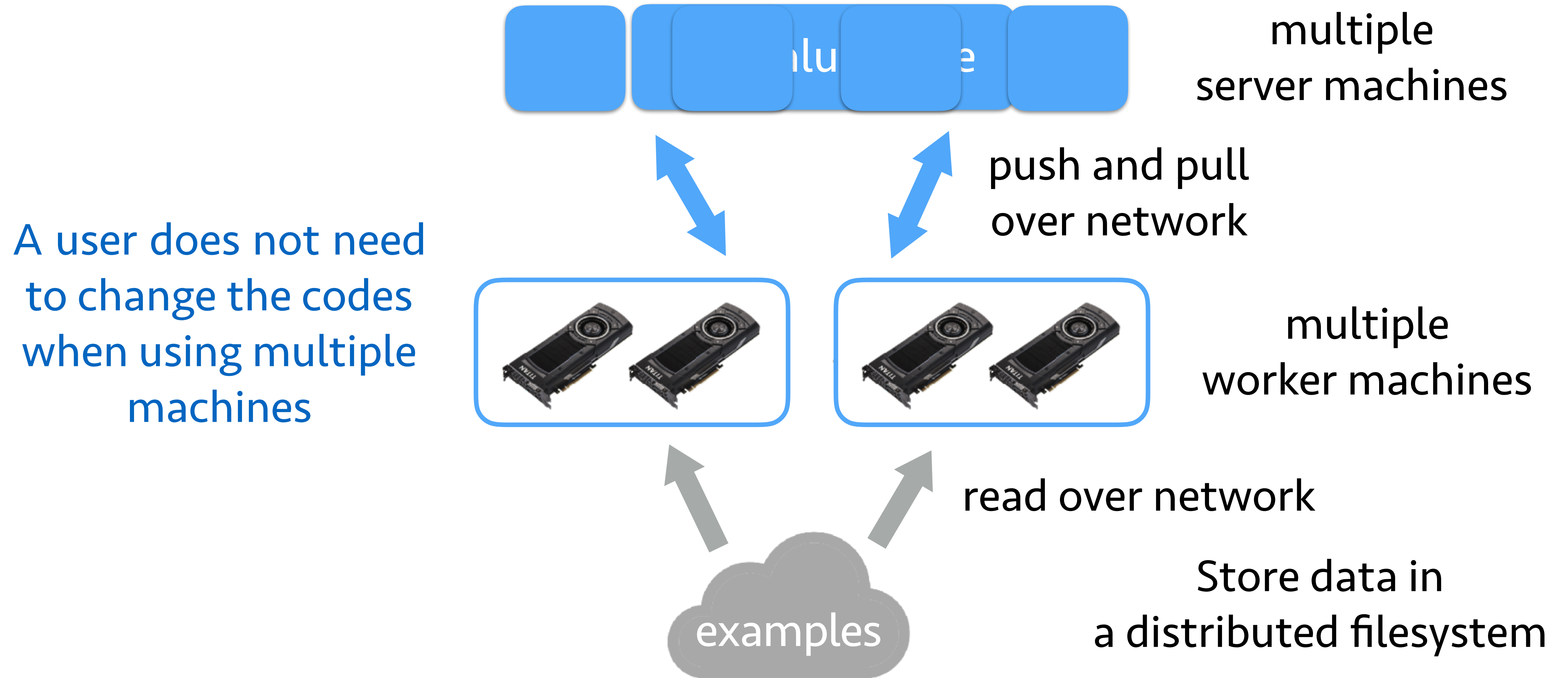
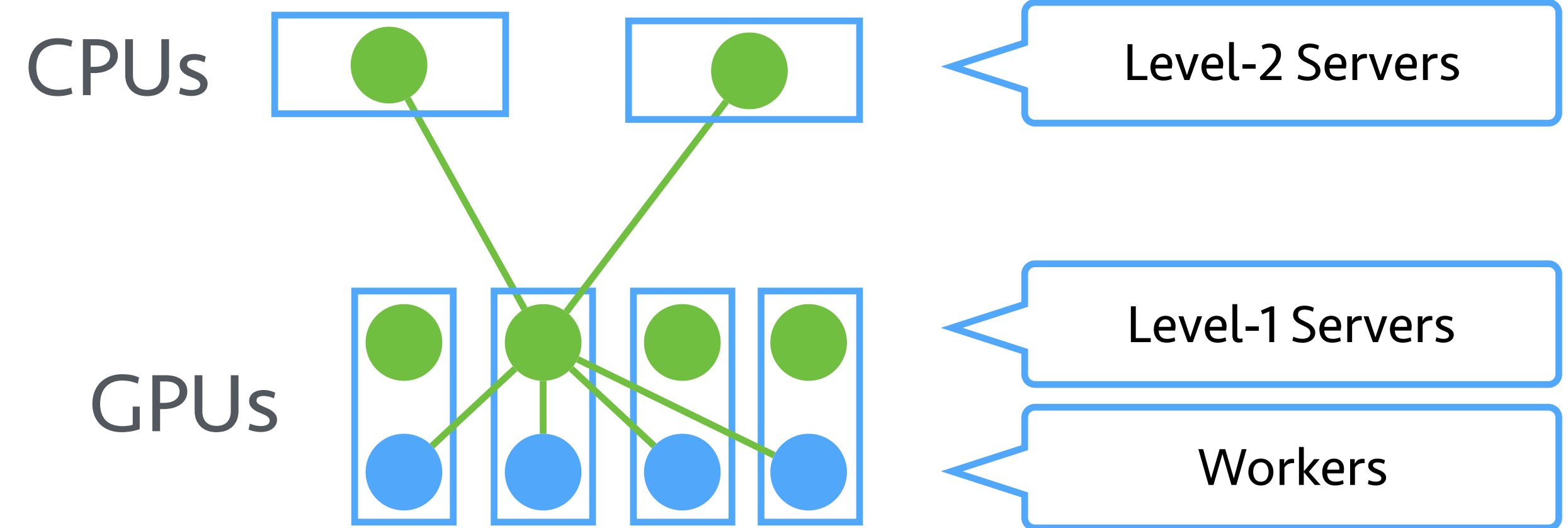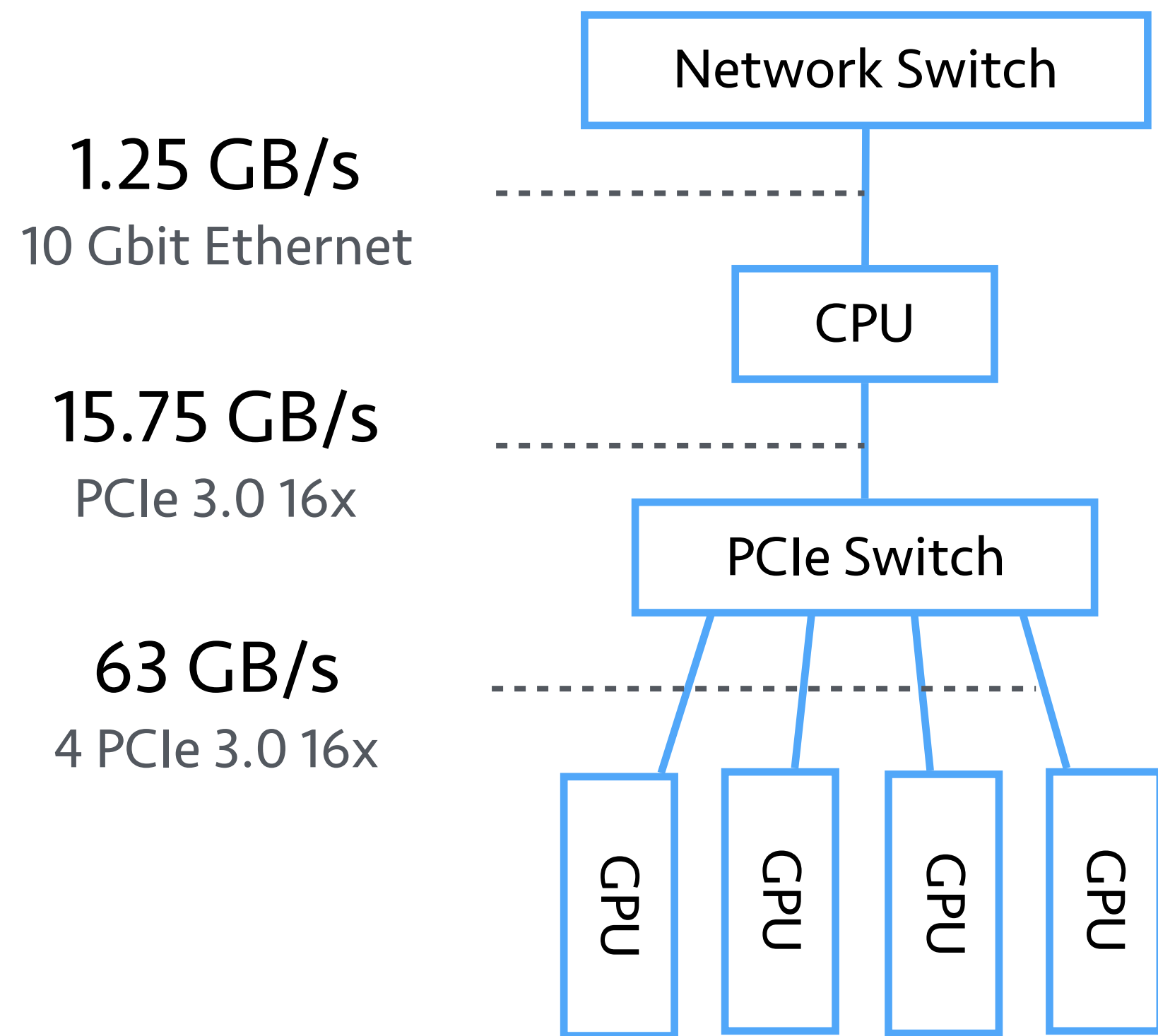Run in **parallel**

# Data Parallelism



1. Read a data partition
2. Pull the parameters
3. Compute the gradient
4. Push the gradient
5. Update the parameters

key-value store

examples

# Distributed Computing

multiple
server machines

push and pull
over network

A user does not need
to change the codes
when using multiple
machines

multiple
worker machines

read over network

examples

Store data in
a distributed filesystem

# Scale to Multiple GPU Machines

Hierarchical parameter server

Network Switch

1.25 GB/s
10 Gbit Ethernet

CPU

15.75 GB/s
PCIe 3.0 16x

PCIe Switch

63 GB/s
4 PCIe 3.0 16x

GPU  GPU  GPU  GPU

CPUs — Level-2 Servers

GPUs — Level-1 Servers
— Workers

# Experiment Setup

✧ IM▲GENET
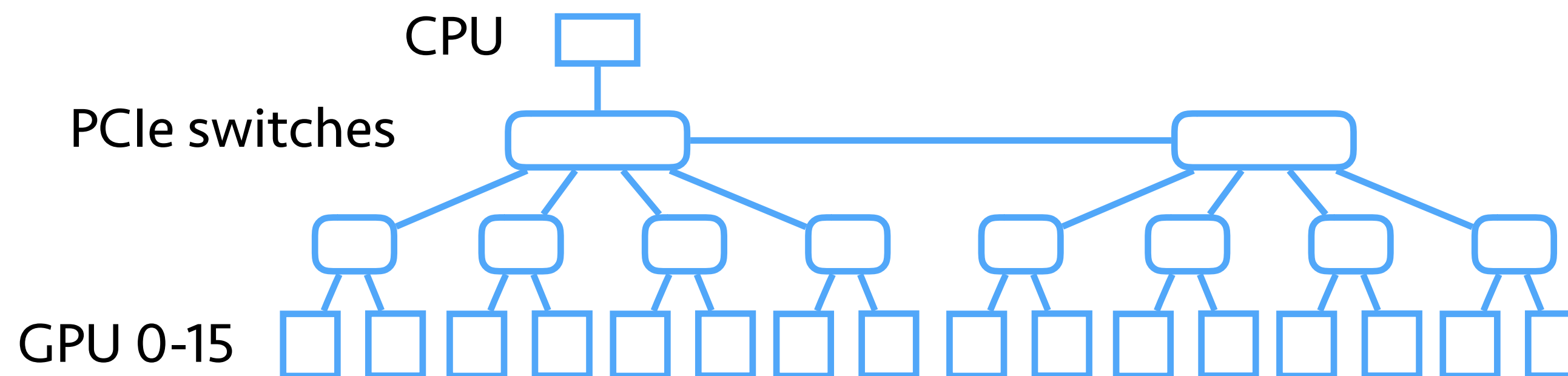
   ✓ 1.2 million images with 1000 classes

✧ Resnet 152-layer model

✧ EC2 P2.16xlarge

✧ Minibatch SGD

✧ Synchronized Updating
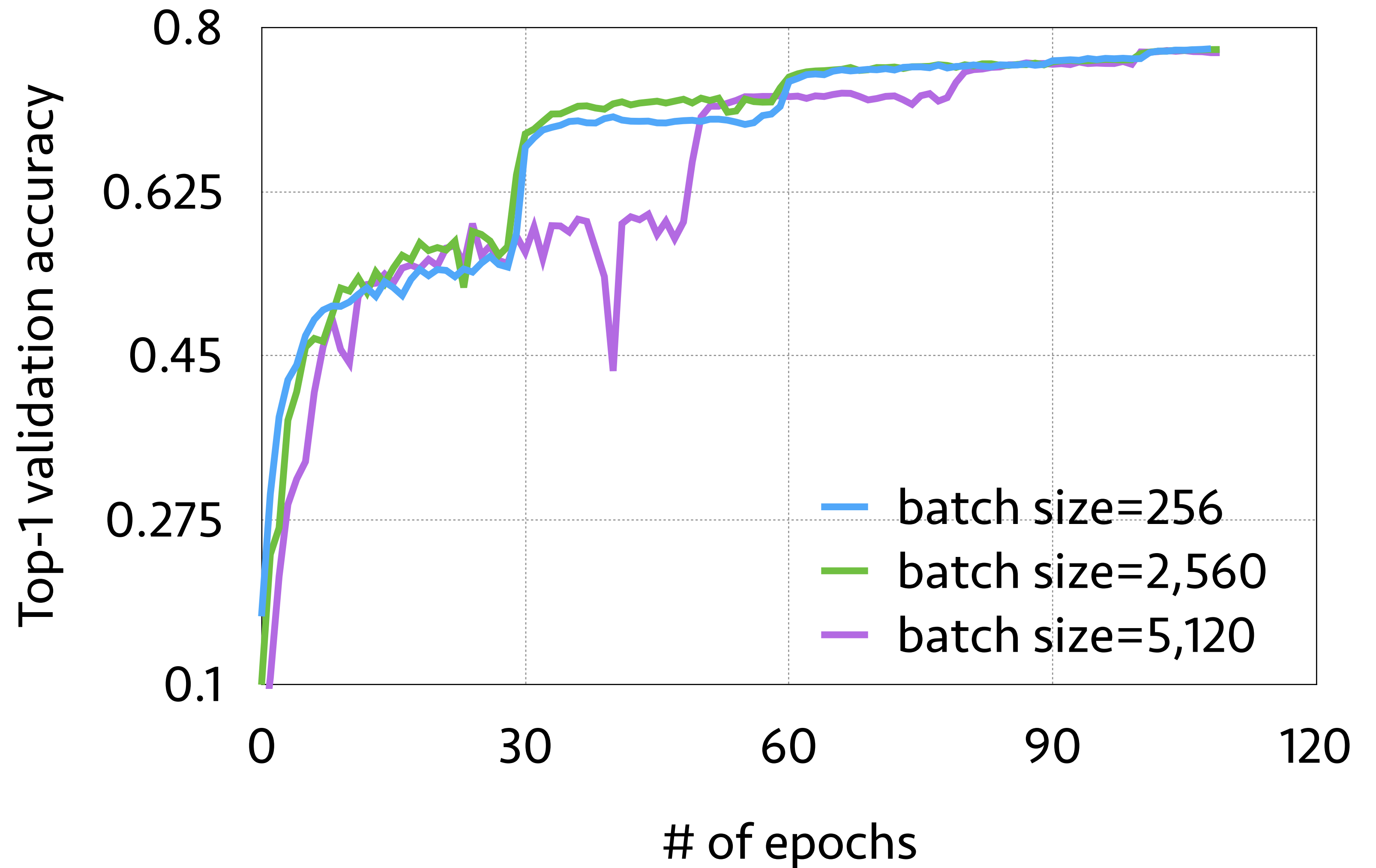
CPU

PCIe switches

GPU 0-15

# Scalability over Multiple Machines

# Convergence

✧ Increase learning rate by 5x

✧ Increase learning rate by 10x, decrease it at epoch 50, 80



24

# Time to achieve 22.5% top-1 accuracy

| | hour |
|---|---|
| 8 GPUs | ~ 1 week |
| 80 GPUs 9.6x | ~ 1 day |
| 160 GPUs 18.8x | ~ half day |

25