

Schedule

Block 4 (3:30 - 5:00 pm)

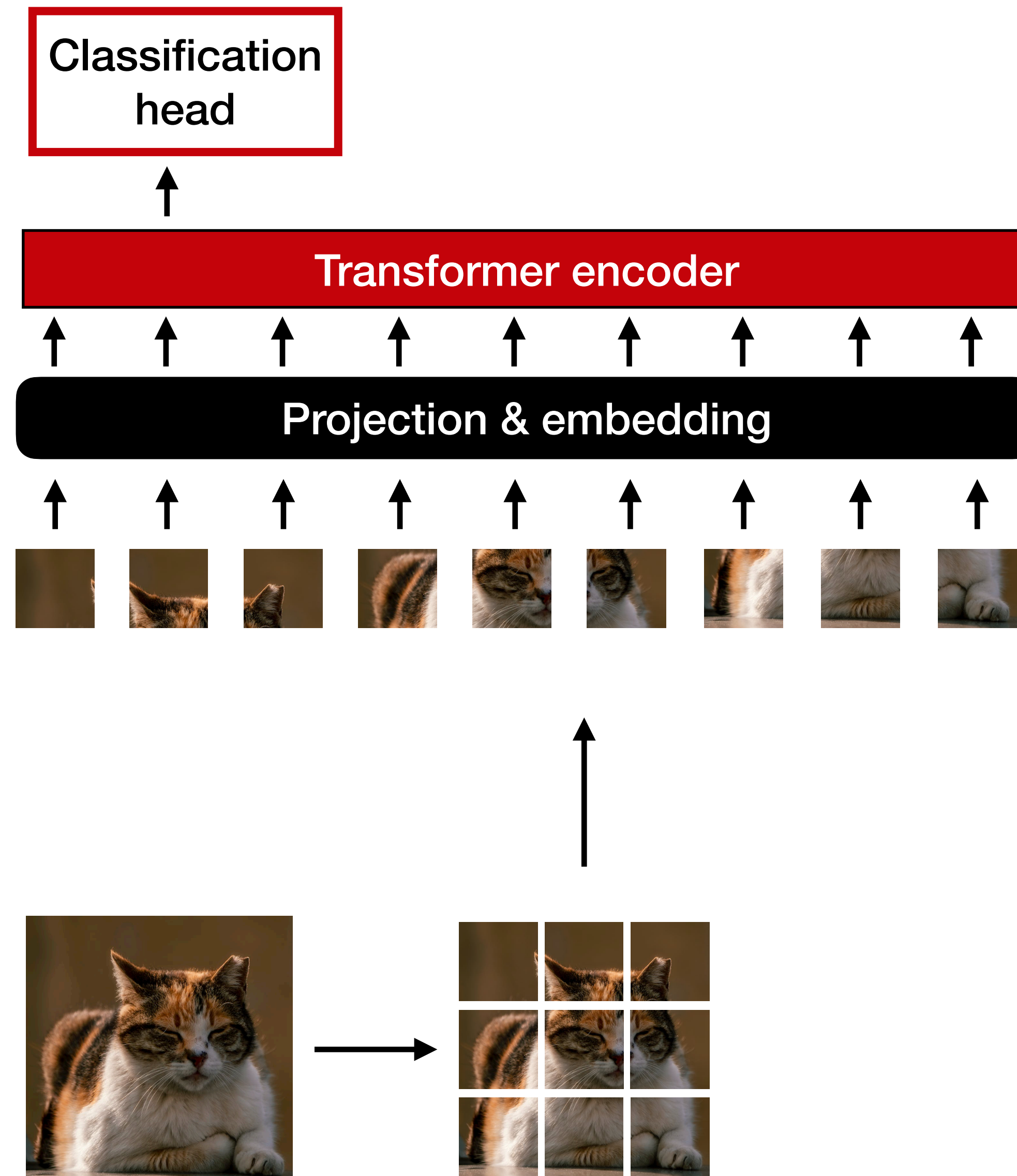
(8) Accelerating PyTorch model training

(9) Finetuning large language models

(10) Conclusion

Vision Transformer

More details at
<https://magazine.sebastianraschka.com/p/ahead-of-ai-10-state-of-computer>



Training a vision transformer **from scratch**

```
Epoch: 0001/0010 | Batch 0000/2812 | Loss: 2.5971
Epoch: 0001/0010 | Batch 0300/2812 | Loss: 1.9497
Epoch: 0001/0010 | Batch 0600/2812 | Loss: 1.6241
...
Epoch: 0001/0010 | Batch 2700/2812 | Loss: 1.5693
Epoch: 0001/0010 | Train acc.: 34.81% | Val acc.: 47.22%
Epoch: 0002/0010 | Batch 0000/2812 | Loss: 1.0045
...
Epoch: 0004/0010 | Batch 2700/2812 | Loss: 1.1540
Epoch: 0004/0010 | Train acc.: 59.33% | Val acc.: 59.14%
...
Epoch: 0005/0010 | Train acc.: 62.04% | Val acc.: 60.06%
...
Epoch: 0006/0010 | Train acc.: 64.22% | Val acc.: 61.80%
...
Epoch: 0007/0010 | Train acc.: 66.12% | Val acc.: 61.14%
Epoch: 0008/0010 | Batch 0000/2812 | Loss: 0.5717
...

Epoch: 0010/0010 | Batch 2700/2812 | Loss: 0.6933
Epoch: 0010/0010 | Train acc.: 71.16% | Val acc.: 62.80%
Time elapsed 61.48 min
Test accuracy 62.85%
```

00_pytorch-vit-random-init.py

```
from torchvision.models import vit_b_16

model = vit_b_16(weights=None)
model.heads.head = torch.nn.Linear(in_features=768, out_features=10)
```

**It's 2023, training from scratch often
doesn't make sense anymore**

Training a vision transformer from scratch

00_pytorch-vit-random-init.py

```
Epoch: 0001/0010 | Batch 0000/2812 | Loss: 2.5971
Epoch: 0001/0010 | Batch 0300/2812 | Loss: 1.9497
Epoch: 0001/0010 | Batch 0600/2812 | Loss: 1.6241
...
Epoch: 0001/0010 | Batch 2700/2812 | Loss: 1.5693
Epoch: 0001/0010 | Train acc.: 34.81% | Val acc.: 47.22%
Epoch: 0002/0010 | Batch 0000/2812 | Loss: 1.0045
...
Epoch: 0004/0010 | Batch 2700/2812 | Loss: 1.1540
Epoch: 0004/0010 | Train acc.: 59.33% | Val acc.: 59.14%
...
Epoch: 0005/0010 | Train acc.: 62.04% | Val acc.: 60.06%
...
Epoch: 0006/0010 | Train acc.: 64.22% | Val acc.: 61.80%
...
Epoch: 0007/0010 | Train acc.: 66.12% | Val acc.: 61.14%
Epoch: 0008/0010 | Batch 0000/2812 | Loss: 0.5717
...

Epoch: 0010/0010 | Batch 2700/2812 | Loss: 0.6933
Epoch: 0010/0010 | Train acc.: 71.16% | Val acc.: 62.80%
Time elapsed 61.48 min
Test accuracy 62.85%
```

```
from torchvision.models import vit_b_16

model = vit_b_16(weights=None)
model.heads.head = torch.nn.Linear(in_features=768, out_features=10)
```

Finetuning a pretrained vision transformer

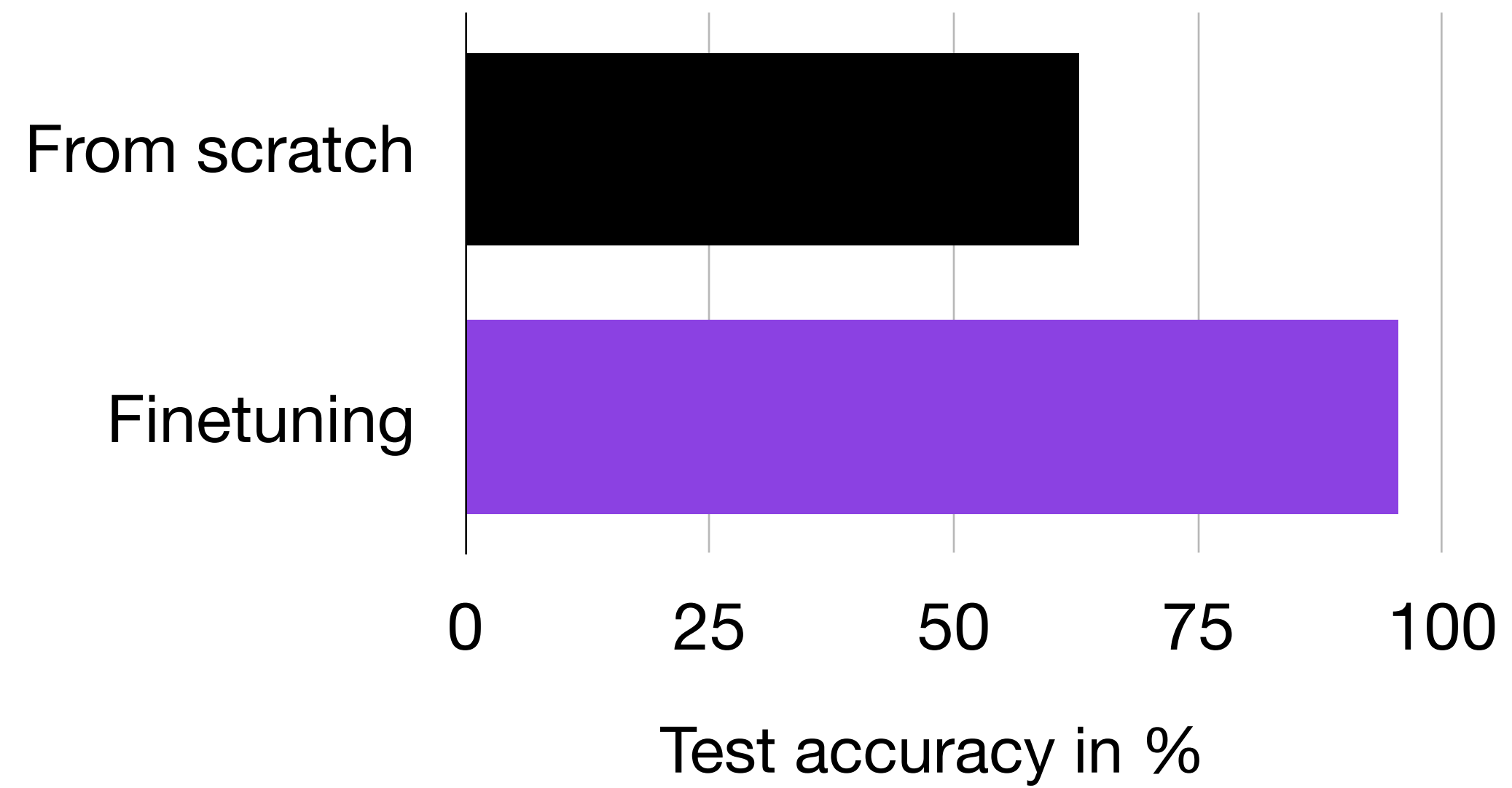
01_pytorch-vit.py

```
Epoch: 0001/0003 | Batch 0000/2812 | Loss: 2.4934
Epoch: 0001/0003 | Batch 0300/2812 | Loss: 0.0954
Epoch: 0001/0003 | Batch 0600/2812 | Loss: 0.0981
Epoch: 0001/0003 | Batch 0900/2812 | Loss: 0.2078
Epoch: 0001/0003 | Batch 1200/2812 | Loss: 0.3588
Epoch: 0001/0003 | Batch 1500/2812 | Loss: 0.0104
Epoch: 0001/0003 | Batch 1800/2812 | Loss: 0.1560
Epoch: 0001/0003 | Batch 2100/2812 | Loss: 0.0474
Epoch: 0001/0003 | Batch 2400/2812 | Loss: 0.4250
Epoch: 0001/0003 | Batch 2700/2812 | Loss: 0.4414
Epoch: 0001/0003 | Train acc.: 92.40% | Val acc.: 94.12%
Epoch: 0002/0003 | Batch 0000/2812 | Loss: 0.0912
Epoch: 0002/0003 | Batch 0300/2812 | Loss: 0.0337
Epoch: 0002/0003 | Batch 0600/2812 | Loss: 0.1545
Epoch: 0002/0003 | Batch 0900/2812 | Loss: 0.0478
...
Epoch: 0003/0003 | Batch 2100/2812 | Loss: 0.0060
Epoch: 0003/0003 | Batch 2400/2812 | Loss: 0.1395
Epoch: 0003/0003 | Batch 2700/2812 | Loss: 0.1128
Epoch: 0003/0003 | Train acc.: 97.24% | Val acc.: 95.74%
Time elapsed 18.70 min
Test accuracy 95.37%
```

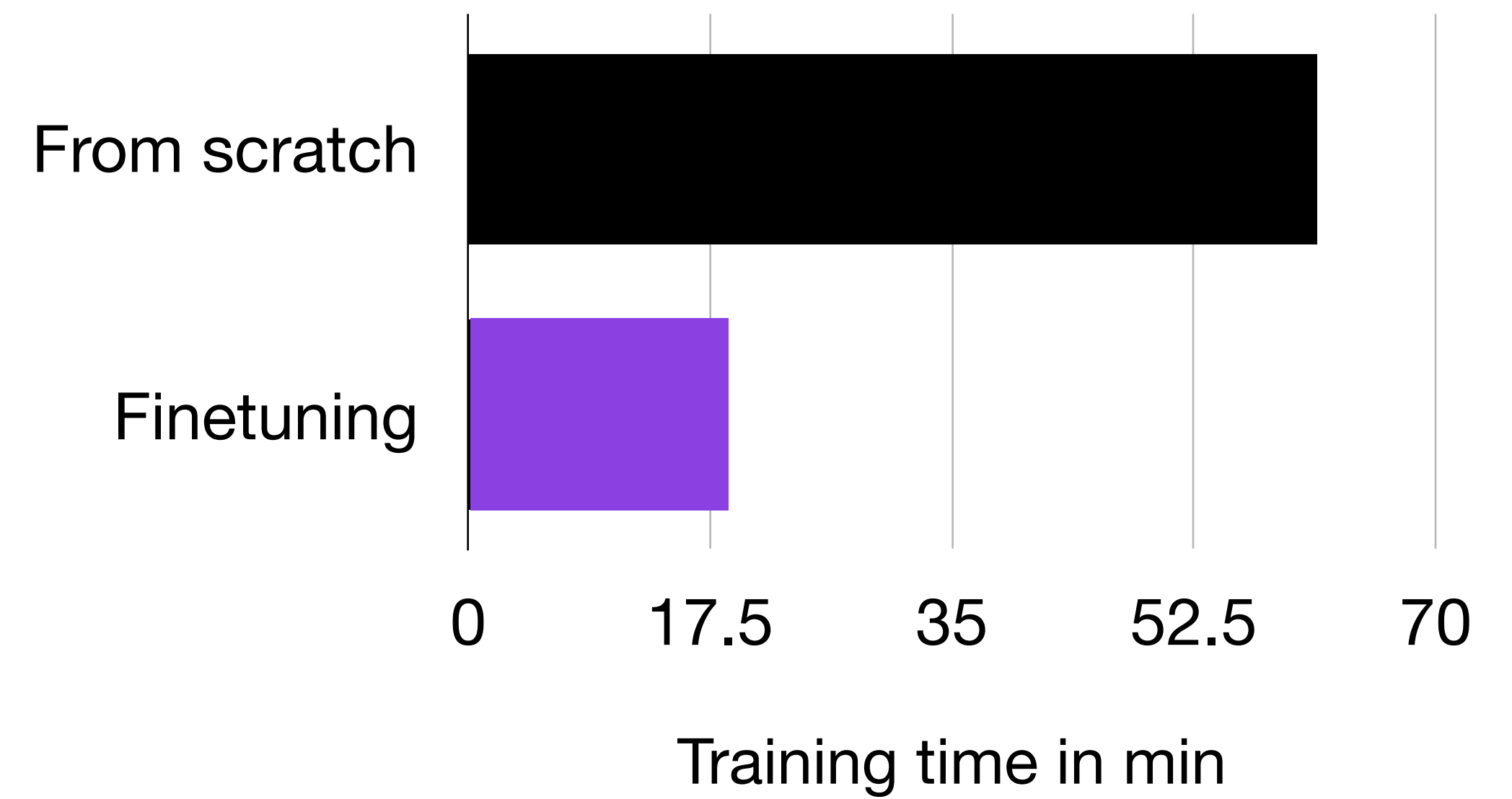
```
from torchvision.models import vit_b_16
from torchvision.models import ViT_B_16_Weights

model = vit_b_16(weights=ViT_B_16_Weights.IMAGENET1K_V1)
model.heads.head = torch.nn.Linear(in_features=768, out_features=10)
```

Accuracy (larger is better)



Time to convergence (lower is better)



How can we accelerate the training even further?

Lightning Fabric is a fast and lightweight way to scale **PyTorch** models without boilerplate code

100% free and open source

```
pip install lightning
```

```
import lightning.fabric as Fabric
```


Suppose we
have the following
plain **PyTorch** code

```
1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5
6 class PyTorchModel(nn.Module):
7     # ...
8
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13
14 num_epochs = 10
15 loss_fn = torch.nn.functional.cross_entropy()
16
17 device = "cuda" if torch.cuda.is_available() else "cpu"
18 model = PyTorchModel(...)
19 optimizer = torch.optim.SGD(model.parameters())
20
21 dataloader = DataLoader(PyTorchDataset(...), ...)
22
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27         input, target = input.to(device), target.to(device)
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31         loss.backward()
32         optimizer.step()
33
34
```

PyTorch

PyTorch + Fabric

```
pytorch.py ↔ pytorch-fabric.py

Users > sebastian > Desktop > pytorch-fabric.py > ...

1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5+ from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12+
13+ fabric = Fabric(accelerator="cuda")
14+ fabric.launch()
15+
16 loss_fn = torch.nn.functional.cross_entropy()
17
18 (variable) dataloader: DataLoader
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20+
21+ model, optimizer = fabric.setup(model, optimizer)
22+ dataloader = fabric.setup_data_loaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27+         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31+         fabric.backward(loss)
32         optimizer.step()
33+
```

Convert PyTorch to Fabric Step by Step

```
1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5 from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13 fabric = Fabric(accelerator="cuda")
14 fabric.launch()
15
16 loss_fn = torch.nn.functional.cross_entropy()
17 model = PyTorchModel(...)
18 optimizer = torch.optim.SGD(model.parameters())
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20
21 model, optimizer = fabric.setup(model, optimizer)
22 dataloader = fabric.setup_data_loaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31         fabric.backward(loss)
32         optimizer.step()
33
```

```
1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5 from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13 fabric = Fabric(accelerator="cuda")
14 fabric.launch()
15
16 loss_fn = torch.nn.functional.cross_entropy()
17 model = PyTorchModel(...)
18 optimizer = torch.optim.SGD(model.parameters())
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20
21 model, optimizer = fabric.setup(model, optimizer)
22 dataloader = fabric.setup_dataloaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31         fabric.backward(loss)
32         optimizer.step()
33
```

Import, instantiate, and launch Fabric

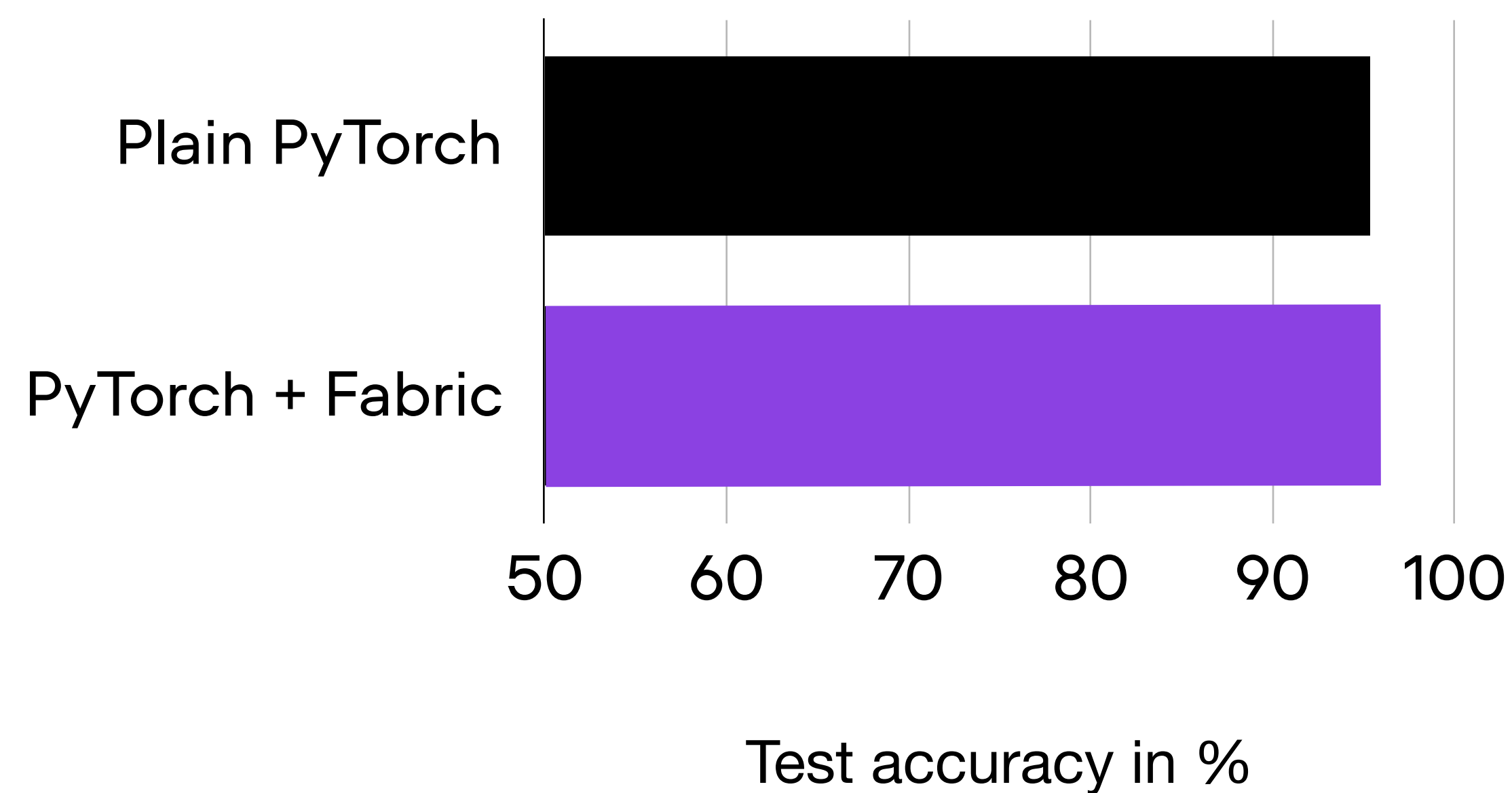
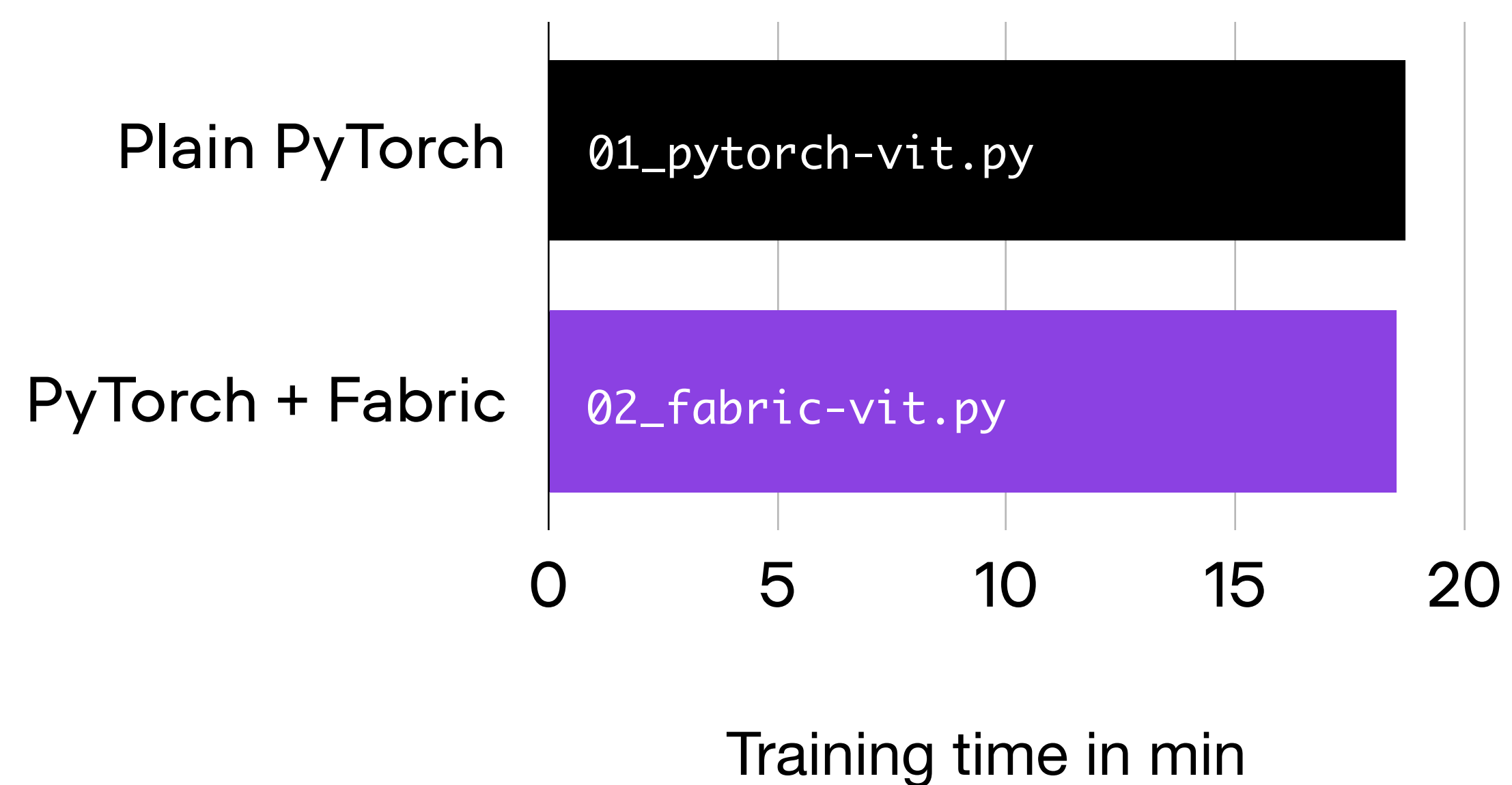
```
1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5 from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13 fabric = Fabric(accelerator="cuda")
14 fabric.launch()
15
16 loss_fn = torch.nn.functional.cross_entropy()
17 model = PyTorchModel(...)
18 optimizer = torch.optim.SGD(model.parameters())
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20
21 model, optimizer = fabric.setup(model, optimizer)
22 dataloader = fabric.setup_data_loaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31         fabric.backward(loss)
32         optimizer.step()
33
```

Set up model, optimizer, and data loader

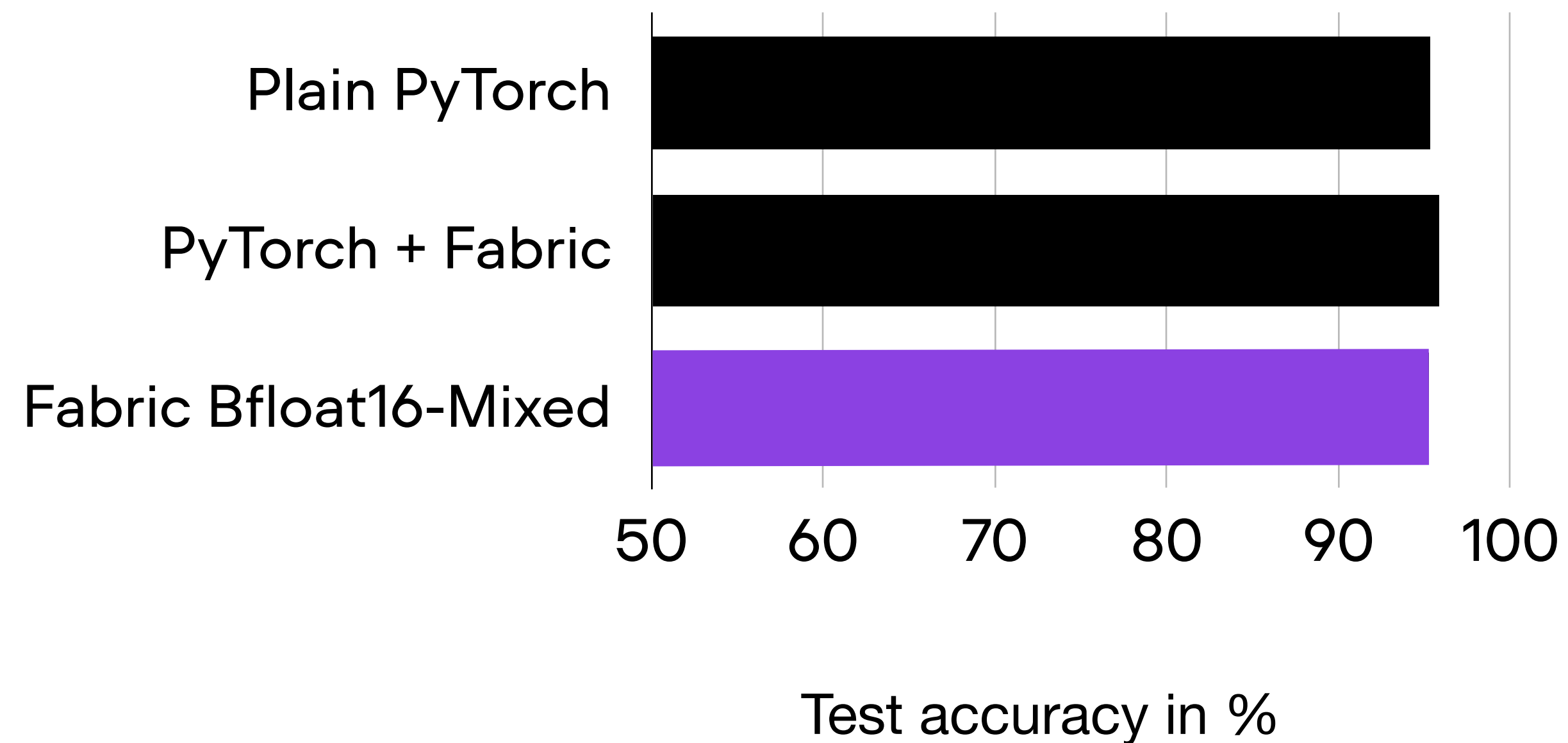
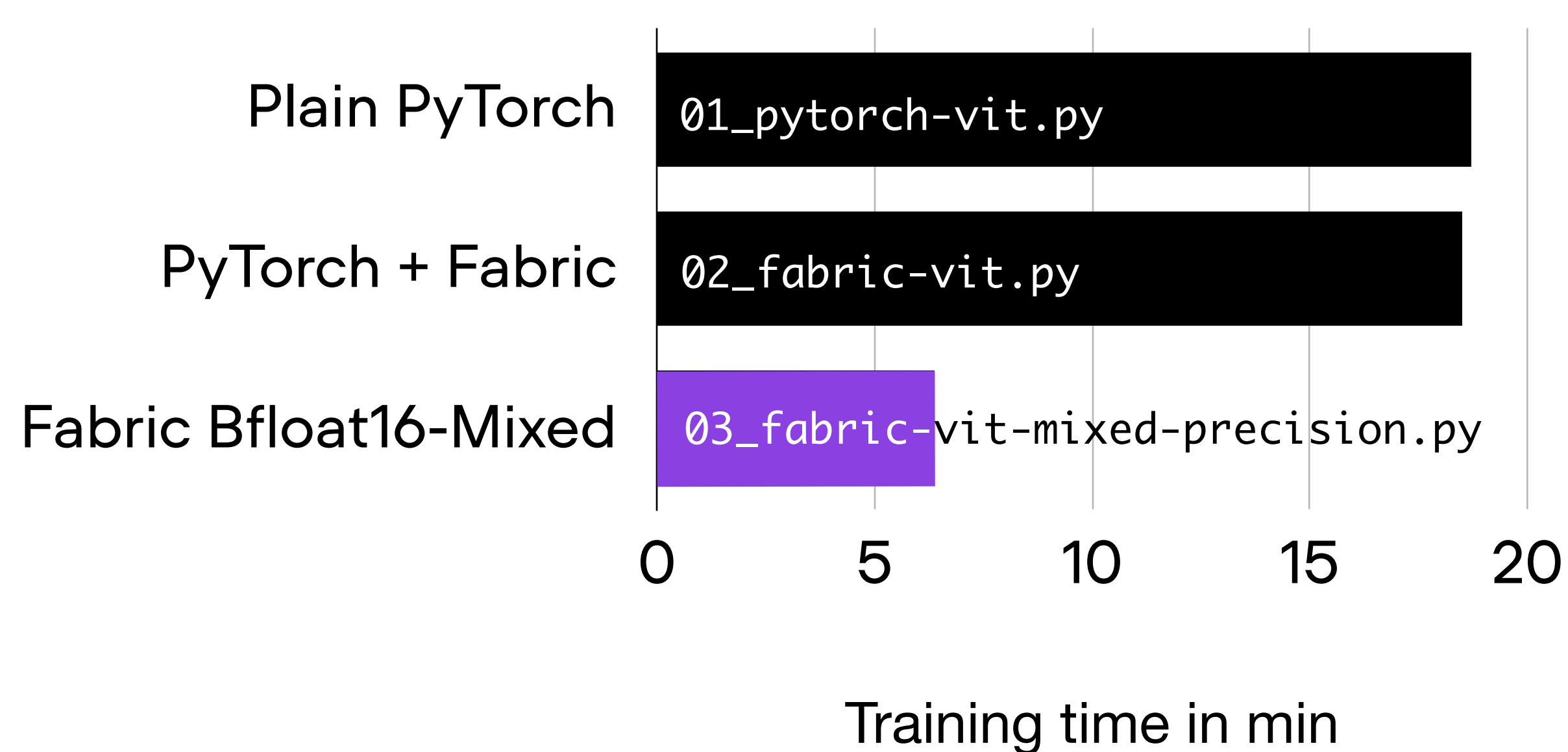
```
1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5 from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13 fabric = Fabric(accelerator="cuda")
14 fabric.launch()
15
16 loss_fn = torch.nn.functional.cross_entropy()
17 model = PyTorchModel(...)
18 optimizer = torch.optim.SGD(model.parameters())
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20
21 model, optimizer = fabric.setup(model, optimizer)
22 dataloader = fabric.setup_data_loaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31         fabric.backward(loss)
32         optimizer.step()
33
```

Use `fabric.backward(loss)`

A Performance Baseline



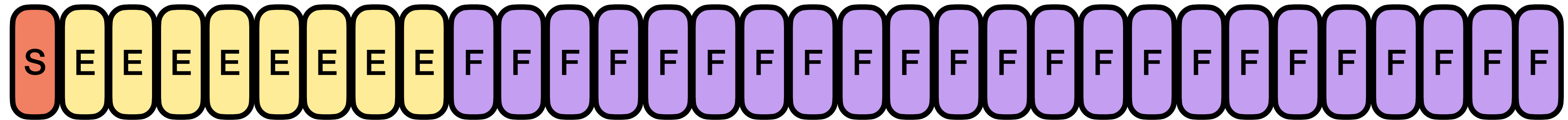
Using **Mixed-Precision** Training



Required Change: 1/2 Lines of Code

```
fabric = Fabric(accelerator="cuda", devices=1, precision="bf16-mixed")
```

Float 32

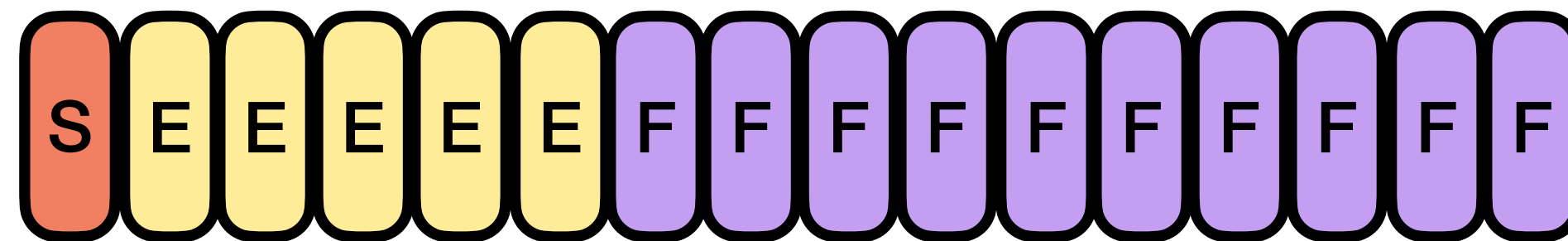


Sign
(1 bit)

Exponent
(8 bits)

Fraction
(23 bits)

Float 16 ("half" precision)

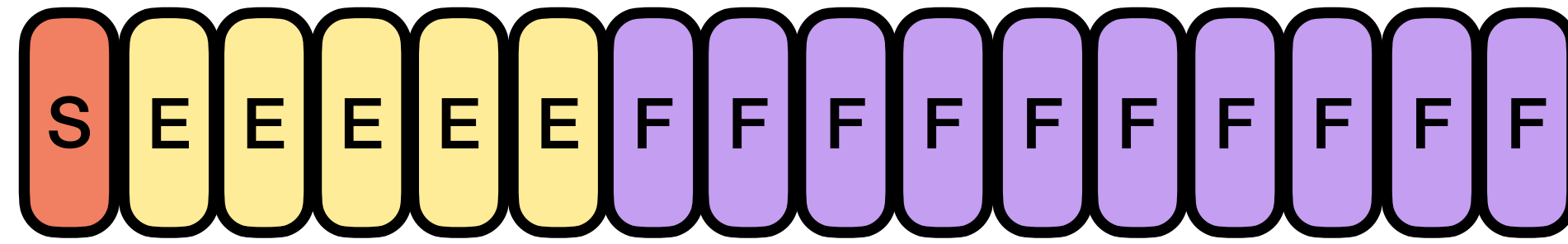


Sign
(1 bit)

Exponent
(5 bits)

Fraction
(10 bits)

Float 16 ("half" precision)

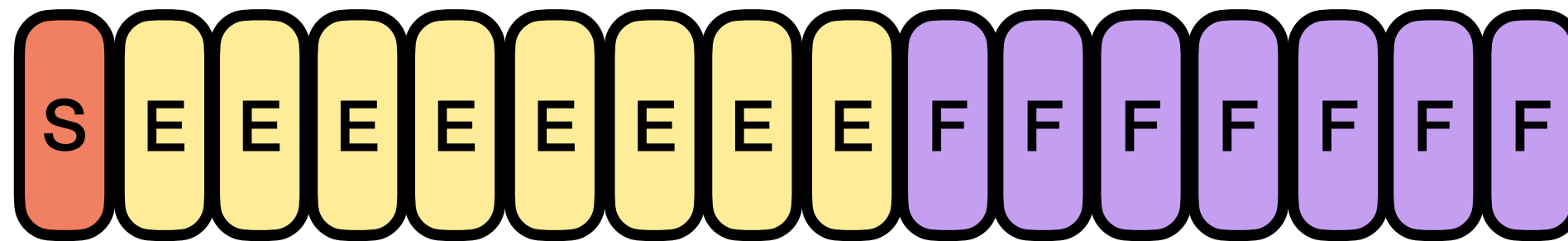


Sign
(1 bit)

Exponent
(5 bits)

Fraction
(10 bits)

Bfloat 16 ("brain" floating point, more "dynamic range" like float 32)

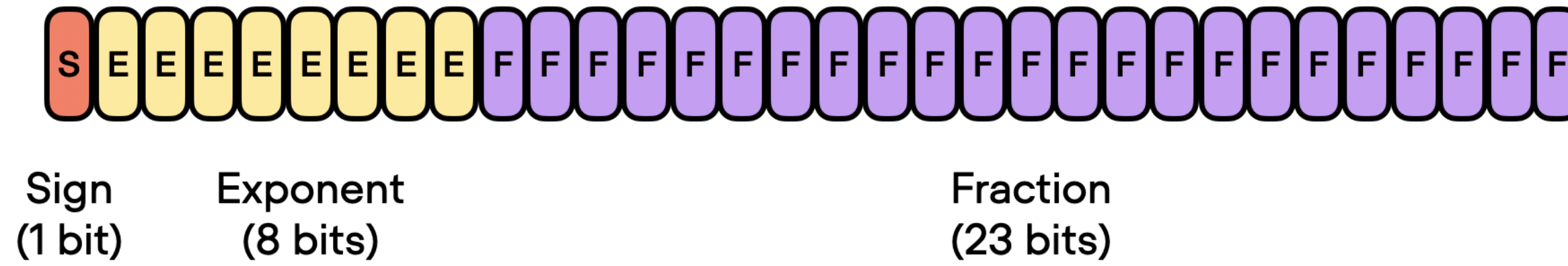


Sign
(1 bit)

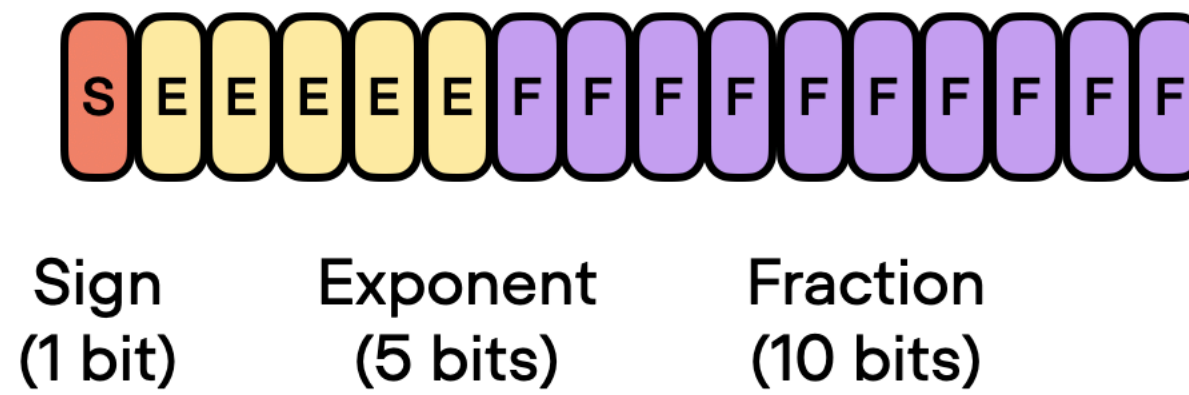
Exponent
(8 bits)

Fraction
(7 bits)

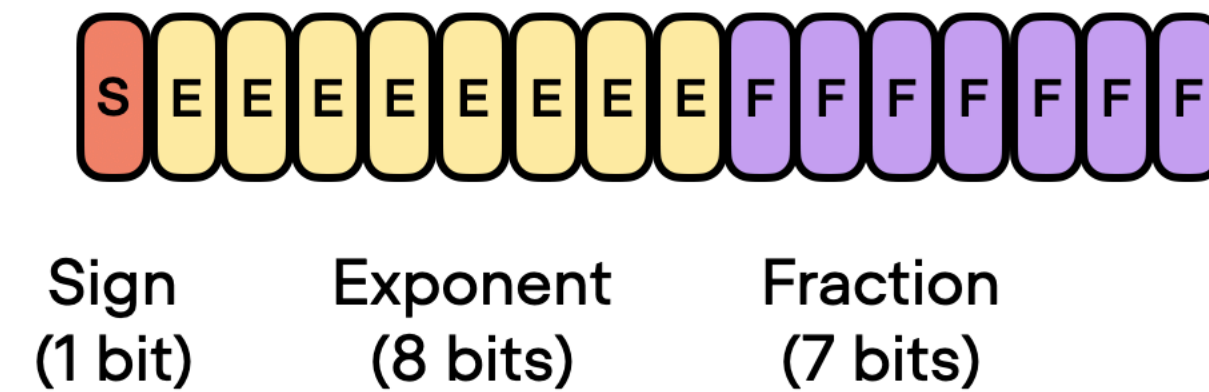
Float 32



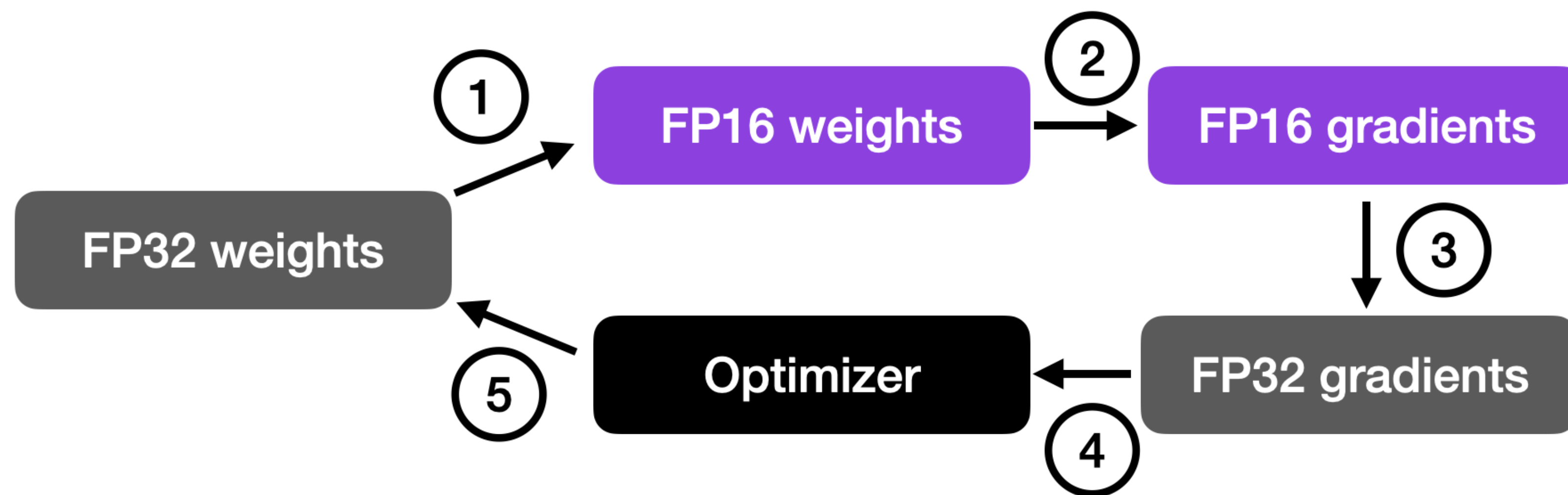
Float 16 ("half" precision)



Bfloat 16 ("brain" floating point, more "dynamic range" like float 32)



Using **Mixed-Precision Training**



Multi-GPU Training with Fully Sharded Data Parallelism

Broad Categories of Parallelism for Multi-GPU Training

Model parallelism

Data parallelism

Pipeline parallelism

Tensor parallelism

Sequence parallelism

Broad Categories of Parallelism for Multi-GPU Training

Model parallelism

Data parallelism

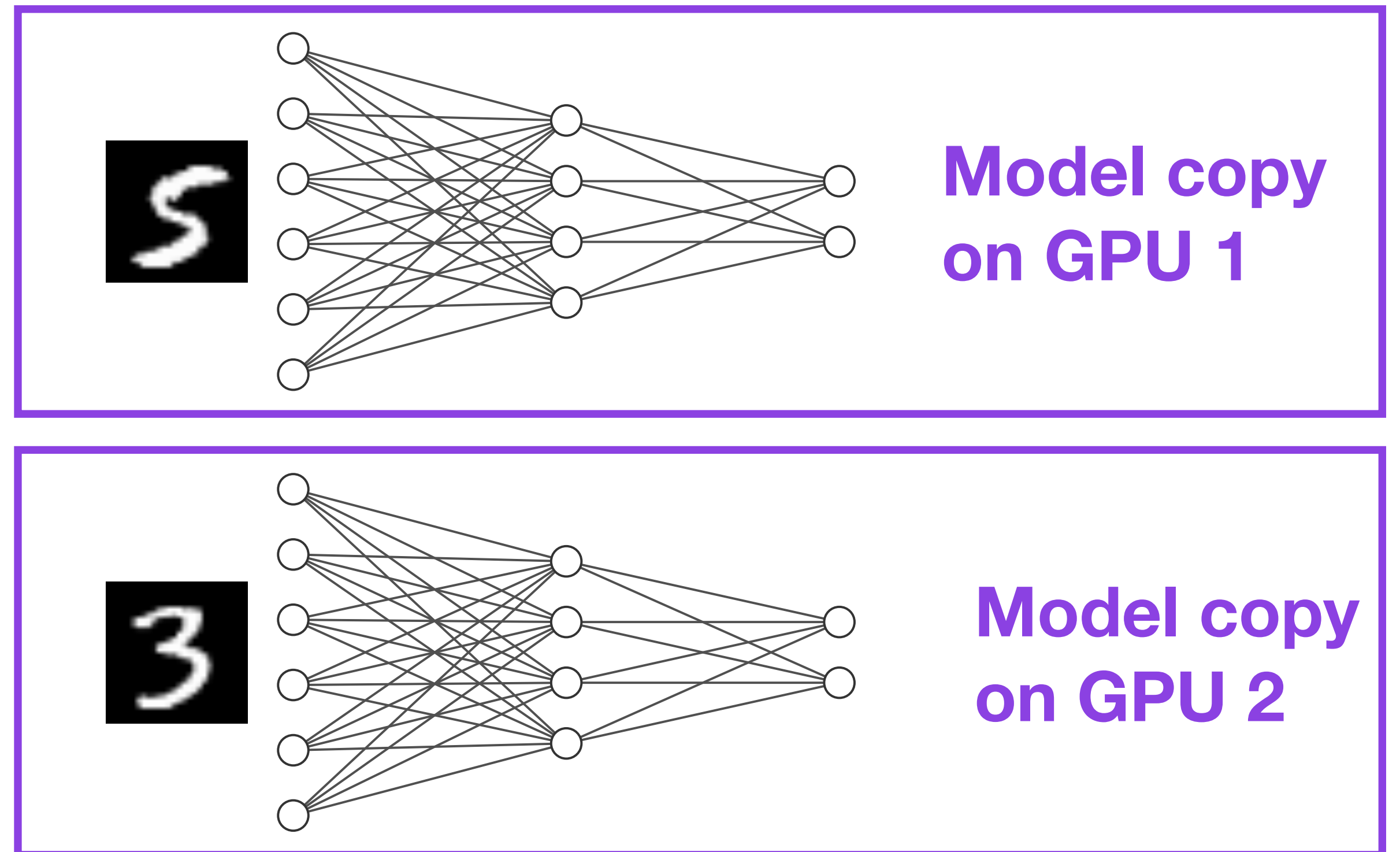
Pipeline parallelism

Tensor parallelism

Sequence parallelism

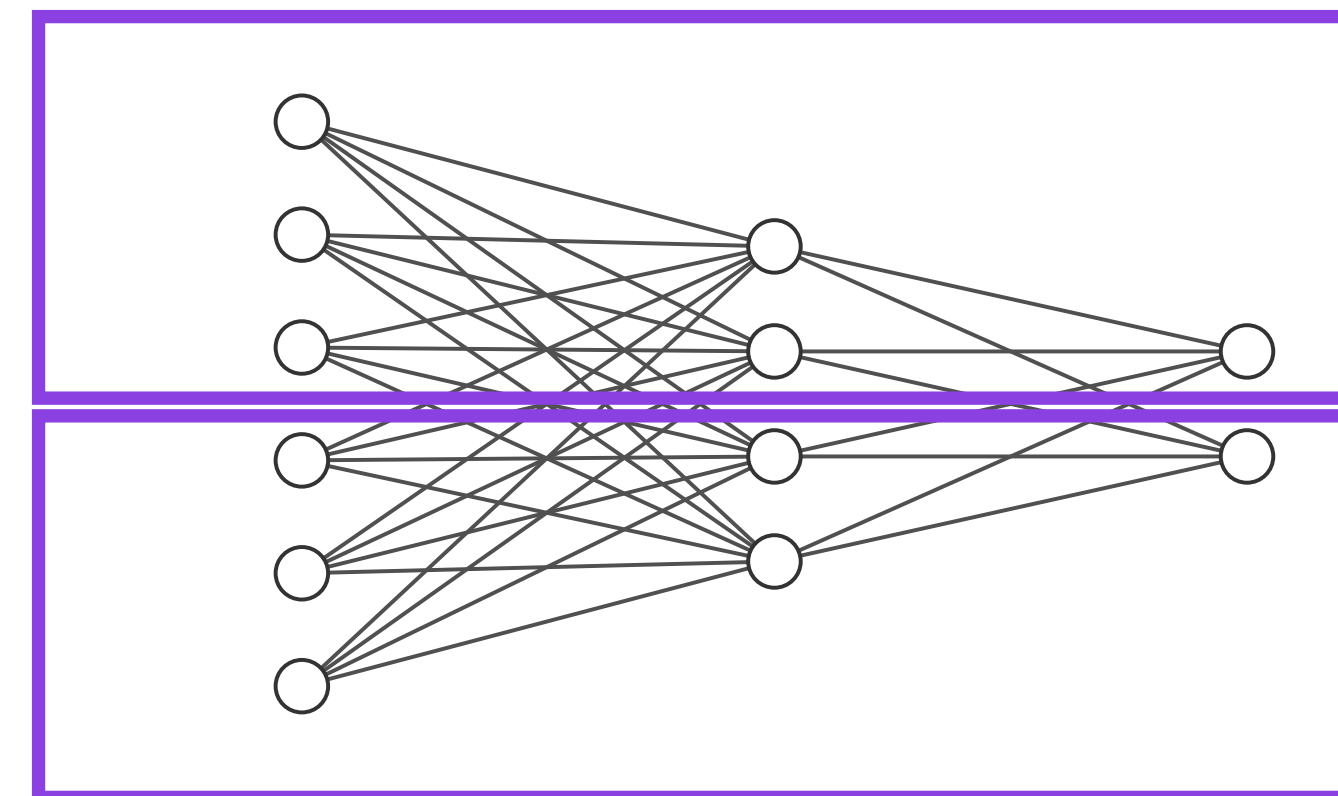
Data Parallelism

**Split batch to train
model(s) on more data
in parallel**



Tensor Parallelism

Related to model parallelism, but split horizontally instead of vertically



GPU 1

GPU 2

(Depending on the implementation, you may split weight matrices, optimizer states, gradients)

Tensor Parallelism

For example, split the matrix multiplication **by column**

The diagram shows the following matrix multiplication:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix}$$

The second column of the second matrix (values 2, 4, 6) and the second column of the resulting matrix (values 28, 64) are highlighted with dashed borders, indicating the column-wise split.

Tensor Parallelism

For example, split the matrix multiplication **by column**

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix}$$

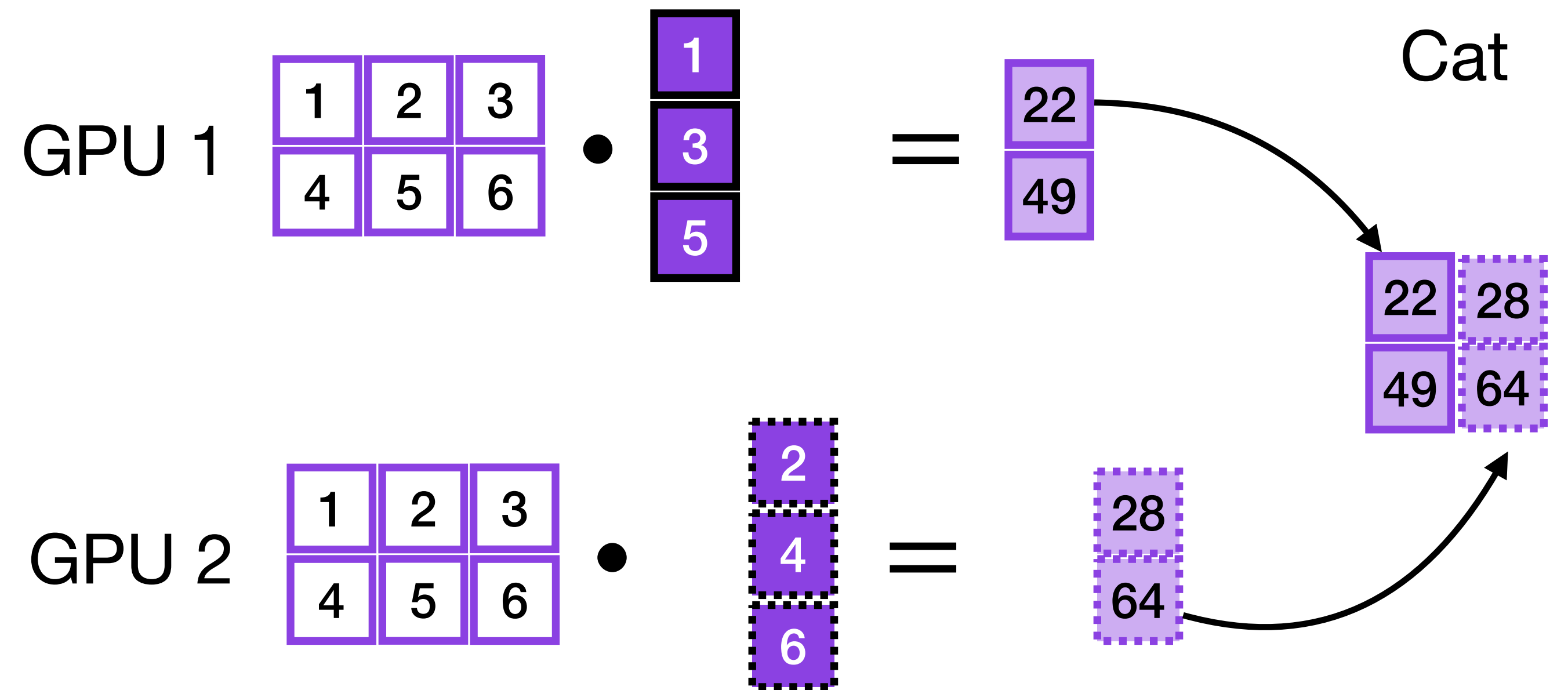
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} = \begin{bmatrix} 22 \\ 49 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix} = \begin{bmatrix} 28 \\ 64 \end{bmatrix}$$

Tensor Parallelism

For example, split the matrix multiplication **by column**

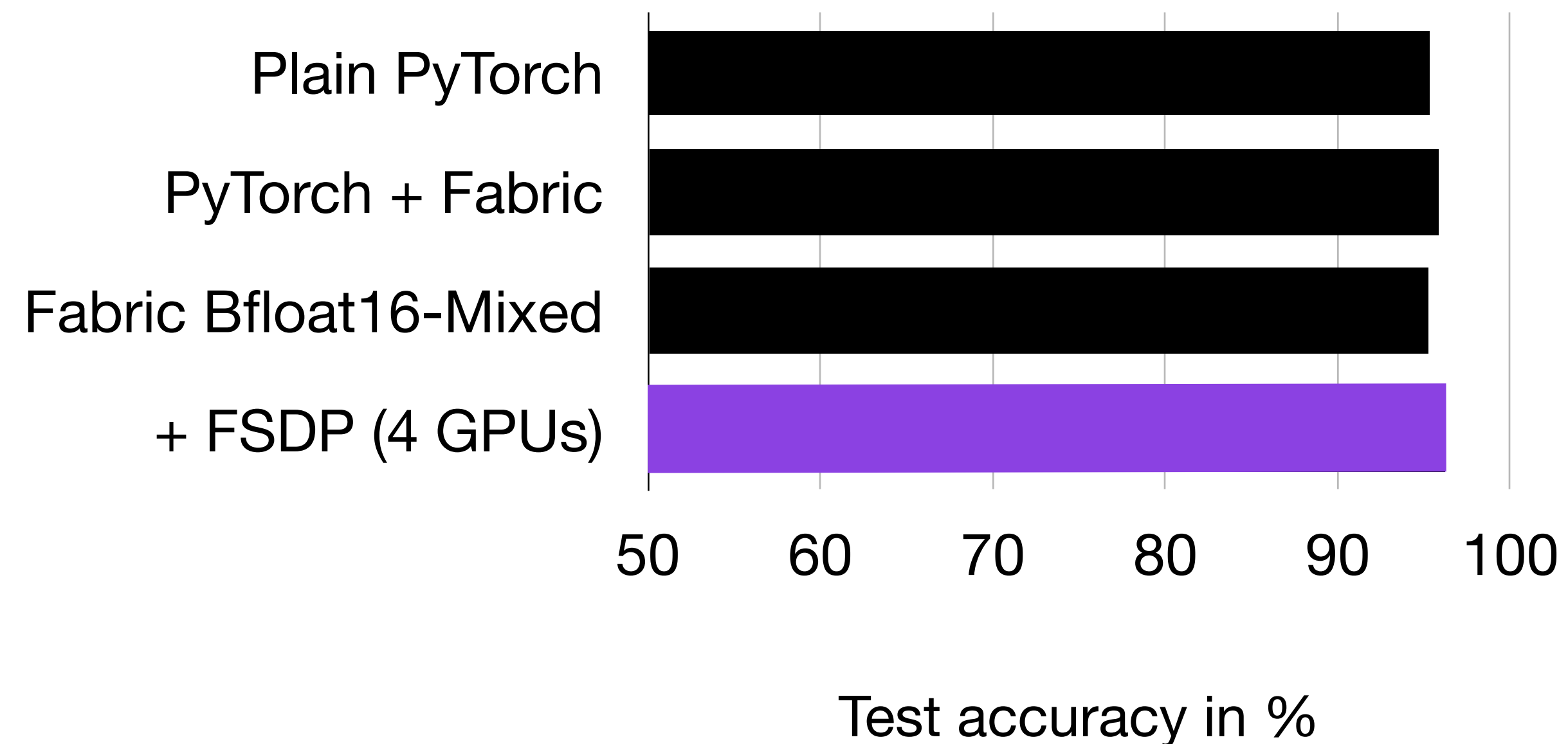
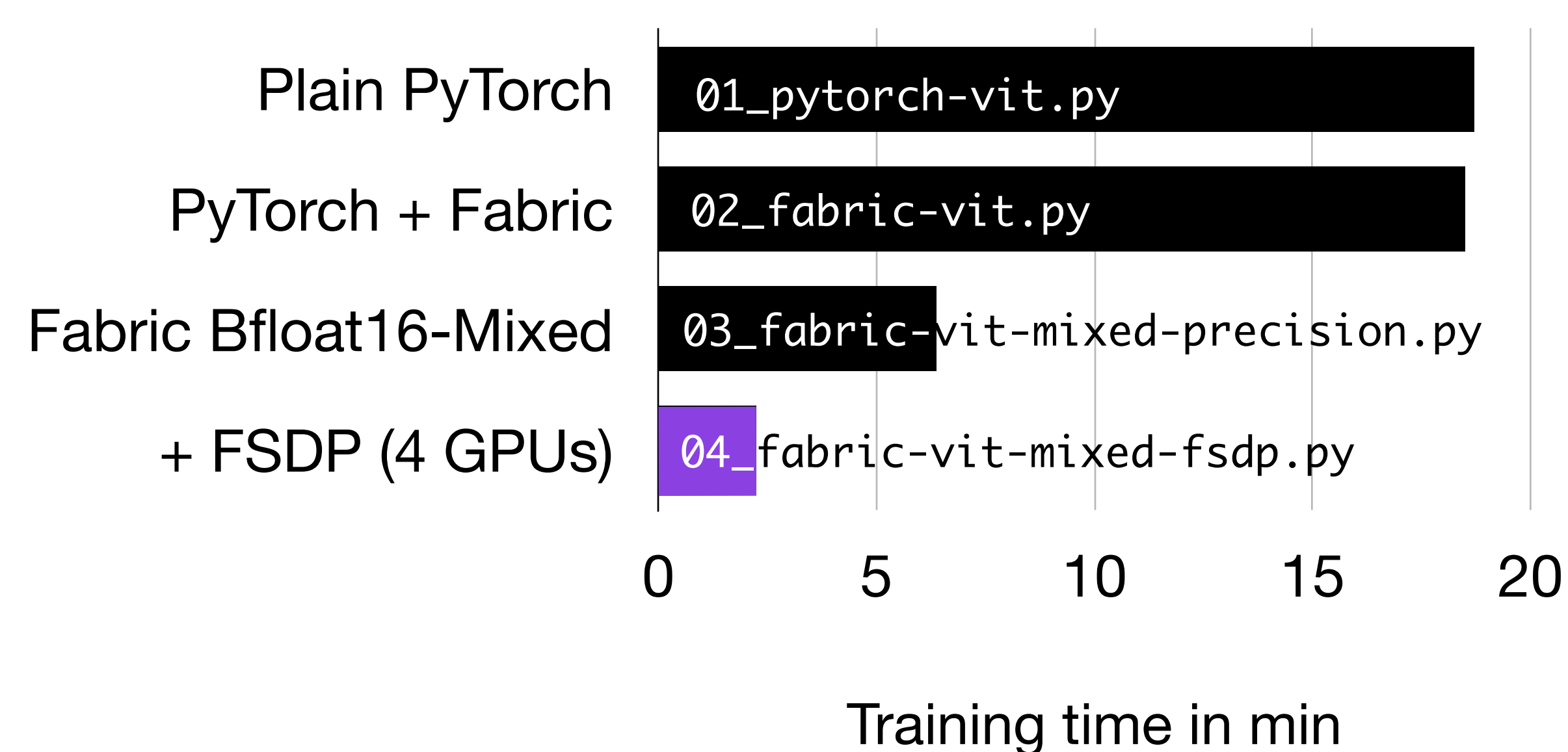
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix}$$



Required Change: 1/2 Lines of Code

```
fabric = Fabric(  
    accelerator="cuda", precision="bf16-mixed",  
    devices=4, strategy="FSDP"  
)
```

Multi-GPU Training with Fully Sharded Data Parallelism



Exercise:

Convert PyTorch script (01_pytorch-vit.py) to Fabric

PyTorch

PyTorch + Fabric

```
pytorch.py ↔ pytorch-fabric.py

Users > sebastian > Desktop > pytorch-fabric.py > ...

1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

1 import torch
2 import torch.nn as nn
3 from torch.utils.data import DataLoader, Dataset
4
5+ from lightning.fabric import Fabric
6
7 class PyTorchModel(nn.Module):
8     # ...
9
10 class PyTorchDataset(Dataset):
11     # ...
12+
13+ fabric = Fabric(accelerator="cuda")
14+ fabric.launch()
15+
16 loss_fn = torch.nn.functional.cross_entropy()
17
18 (variable) dataloader: DataLoader
19 dataloader = DataLoader(PyTorchDataset(...), ...)
20+
21+ model, optimizer = fabric.setup(model, optimizer)
22+ dataloader = fabric.setup_data_loaders(dataloader)
23 model.train()
24
25 for epoch in range(num_epochs):
26     for batch in dataloader:
27+         input, target = batch
28         optimizer.zero_grad()
29         output = model(input)
30         loss = loss_fn(output, target)
31+         fabric.backward(loss)
32         optimizer.step()
33+
```