

Introduction to practical aspects of Deep Learning

Yuping Luo



Introduction to ~~practical aspects~~ ~~of Deep Learning~~ PyTorch

Yuping Luo

Why using a framework?

- Performance;
- Fewer bugs;
- Code reuse (backpropagation, convolution, etc.);
- Community;
- ...

Basic Pipeline

1. (Installation and Import)
2. Data Loading
3. Network Architecture
4. Optimization(train)
5. Evaluation

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torch.optim as optim
5 from torch.utils.data import DataLoader
6 from torchvision import datasets, transforms
7
8 train_loader = torch.utils.data.DataLoader(
9     datasets.MNIST('../data', train=True, download=True,
10         transform=transforms.Compose([transforms.ToTensor(),
11             transforms.Normalize((0.1307,), (0.3081,))])),
12     batch_size=32, shuffle=True)
13 test_loader = torch.utils.data.DataLoader(
14     datasets.MNIST('../data', train=False,
15         transform=transforms.Compose([transforms.ToTensor(),
16             transforms.Normalize((0.1307,), (0.3081,))])),
17     batch_size=32)
18
19 device = torch.device("cpu") # "cuda"
20
21 model = nn.Sequential(
22     nn.Linear(784, 200),
23     nn.ReLU(),
24     nn.Linear(200, 10),
25     nn.LogSoftmax(dim=-1),
26 ).to(device)
27
28 optimizer = optim.SGD(model.parameters(), lr=0.01)
29
30 def train():
31     model.train()
32     for data, target in train_loader:
33         data, target = data.to(device).view(-1, 28 * 28), target.to(device)
34         optimizer.zero_grad()
35         output = model(data)
36         loss = F.nll_loss(output, target)
37         loss.backward()
38         optimizer.step()
39
40 def test():
41     model.eval()
42     correct = 0
43     with torch.no_grad():
44         for data, target in test_loader:
45             data, target = data.to(device).view(-1, 28 * 28), target.to(device)
46             output = model(data)
47             pred = output.max(1, keepdim=True)[1]
48             correct += (pred == target.view_as(pred)).sum().item()
49     print('accuracy', correct / len(test_loader.dataset))
50
51 for epoch in range(100):
52     train()
53     test()
```

Installation and Import

```
$ pip install torch torchvision
```

```
1 import torch      # Basic Tensor operations
2 import torch.nn as nn    # Modules and layers (class style)
3 import torch.nn.functional as F    # Modules and layers (function style)
4 import torch.optim as optim    # Optimizers
5 from torch.utils.data import DataLoader
6 from torchvision import datasets, transforms
```

Data Loading

Data preprocessing

```
8 train_loader = torch.utils.data.DataLoader(  
9     datasets.MNIST('../data', train=True, download=True,  
10         transform=transforms.Compose([transforms.ToTensor(),  
11         transforms.Normalize((0.1307,), (0.3081,))])),  
12     batch_size=32, shuffle=True)  
13 test_loader = torch.utils.data.DataLoader(  
14     datasets.MNIST('../data', train=False,  
15         transform=transforms.Compose([transforms.ToTensor(),  
16         transforms.Normalize((0.1307,), (0.3081,))])),  
17     batch_size=32)
```

Network Architecture and Optimizer

```
17 device = torch.device("cpu") # "cuda"
18
19 model = nn.Sequential(
20     nn.Linear(784, 200),
21     nn.ReLU(),
22     nn.Linear(200, 10),
23     nn.LogSoftmax(dim=-1),
24 ).to(device) ← Move parameters to device (cpu or cuda)
25
26 optimizer = optim.SGD(model.parameters(), lr=0.01)
```

Training

```
28 def train():
29     model.train() ← Use train mode: Dropout/BatchNorm/etc.
30     for data, target in train_loader:
31         data, target = data.to(device).view(-1, 28 * 28), target.to(device)
32         optimizer.zero_grad() ← compute gradient
33         output = model(data)
34         loss = F.nll_loss(output, target)
35         loss.backward() ← compute gradient
36         optimizer.step() ← apply gradient
```

Evaluation

+ Average over multiple random seeds!

```
40 def test():
41     model.eval() ← Use train mode: Dropout/BatchNorm/etc.
42     correct = 0
43     with torch.no_grad(): don't need gradient
44         for data, target in test_loader:
45             data, target = data.to(device).view(-1, 28 * 28), target.to(device)
46             output = model(data)
47             pred = output.max(1, keepdim=True)[1]
48             correct += (pred == target.view_as(pred)).sum().item()
49     print('accuracy', correct / len(test_loader.dataset))
```

Main Loop

```
51 for epoch in range(100):  
52     train()  
53     test()  
54
```

Run!

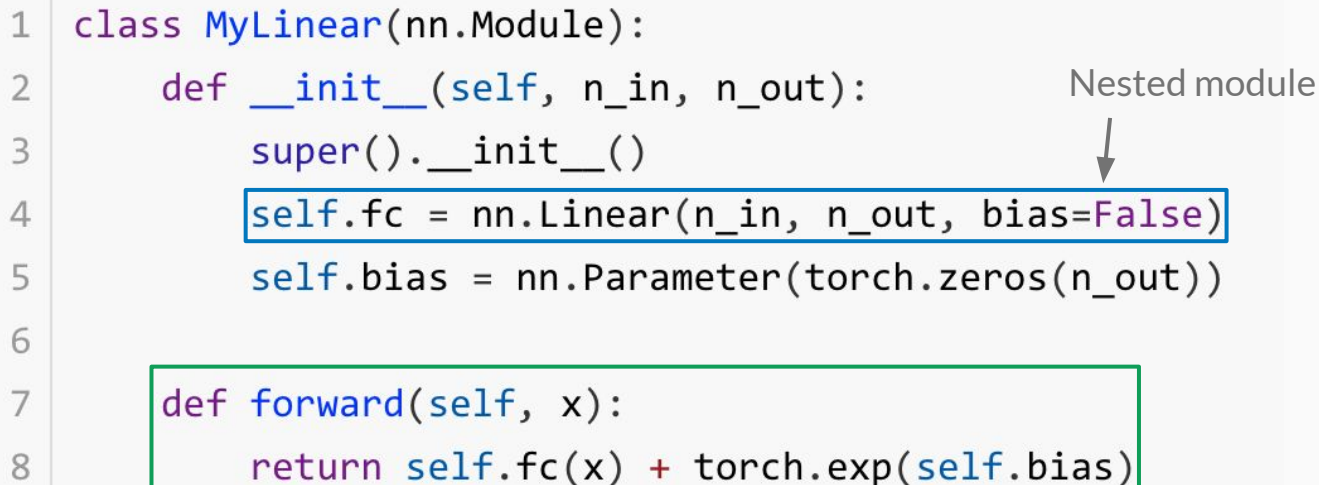


nn.Module: Flexible Network Architecture

A container of parameters.

```
1 class MyLinear(nn.Module):
2     def __init__(self, n_in, n_out):
3         super().__init__()
4         self.fc = nn.Linear(n_in, n_out, bias=False)
5         self.bias = nn.Parameter(torch.zeros(n_out))
6
7     def forward(self, x):
8         return self.fc(x) + torch.exp(self.bias)
```

Nested module
↓



nn.Module: Flexible Network Architecture

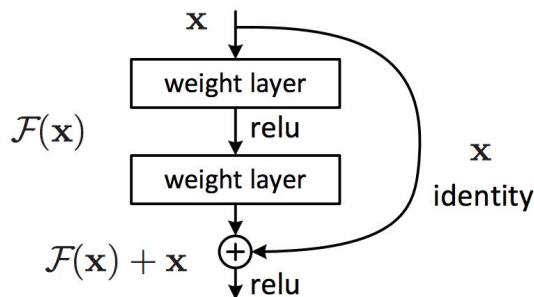


Figure 2. Residual learning: a building block.

He, Kaiming, et al.

```
1 class ResBlock(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.fc1 = nn.Linear(100, 100)
5         self.fc2 = nn.Linear(100, 100)
6
7     def forward(self, x):
8         y = self.fc1(x)
9         y = F.relu(y)
10        y = self.fc2(y)
11        return F.relu(y + x)
12
13 net = nn.Sequential(
14     # ...
15     ResBlock(),
16     ResBlock(),
17     # ...
18 )
```

Inplace Operations

Do NOT use inplace operations if you require grads.

Live examples!

GPU

1. GPU/CPU interaction is slow.
2. Large batch size.
 - a. Data parallel (mostly used)
 - b. Async gradient update (ASGD, etc.)
3. Floating point: 32-bit (float) vs 64-bit (double) vs 16-bit (half)
4. `nvidia-smi`
5. Async preprocessing (by CPU).

Reproducibility

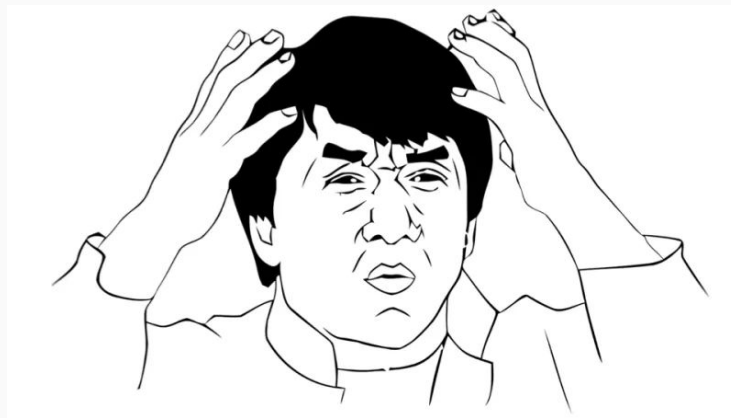
1. Easier to debug.
2. Fix random seed!

```
1 np.random.seed(seed)
2 tf.set_random_seed(np.random.randint(2**30))
3 torch.manual_seed(np.random.randint(2**30))
4 random.seed(np.random.randint(2**30))
5 torch.cuda.manual_seed_all(np.random.randint(2**30))
```

3. Make a copy of source code / command lines.

...TensorFlow

1. Fewer calls to `sess.run` due to large overhead in TensorFlow.
2. Debug?
 - a. `tf.Print`
 - b. `sess.run('Add:0')`
 - c. ...or eager mode



Overfitting

1. Regularization (L2, etc.);
2. Dropout;
3. Data augmentation;
4. Smaller network;
5. Early stop;
6. ...

Hyperparameter Tuning

1. Coordinate Descent;
2. Grid Search;
3. Random Search;
4. ...



Advanced Techniques

Hessian-vector product

- + Why not storing the whole matrix?

- + Quadratic Form $f(x) \approx f(x_0) + (x - x_0)^\top \nabla f(x)|_{x_0} + \frac{1}{2}(x - x_0)^\top \nabla^2 f(x)|_{x_0} (x - x_0)$

- + Minimizer $x = H^{-1} \nabla f(x)|_{x_0}$

- + Conjugate Gradient only requires to compute Hv

Hessian-vector product

$$\mathbf{H}v = \frac{\partial}{\partial x} \left(v^T \frac{\partial f}{\partial x} \right)$$

```
1 def hessian_vec_prod(f, params, v: torch.Tensor) -> torch.Tensor:
2     grads = torch.autograd.grad(f, params, create_graph=True)
3
4     dot = nn.utils.parameters_to_vector(grads).mul(v).sum()
5     grads = [g.contiguous() for g in torch.autograd.grad(dot, params,
6                   retain_graph=True)]
7     return nn.utils.parameters_to_vector(grads)
```

symbolic

Hessian-vector product

$$\mathbf{H}v = \frac{\partial}{\partial x} \left(v^\top \frac{\partial f}{\partial x} \right)$$

```
1 def hessian_vec_prod(f, params, v: torch.Tensor) -> torch.Tensor:
2     grads = torch.autograd.grad(f, params, create_graph=True)
3
4     dot = nn.utils.parameters_to_vector(grads).mul(v).sum()
5     grads = [g.contiguous() for g in torch.autograd.grad(dot, params,
6     retain_graph=True)]
7     return nn.utils.parameters_to_vector(grads)
```

Hessian-vector product

$$\mathbf{H}v = \frac{\partial}{\partial x} \left(v^T \frac{\partial f}{\partial x} \right)$$

```
1 def hessian_vec_prod(f, params, v: torch.Tensor) -> torch.Tensor:
2     grads = torch.autograd.grad(f, params, create_graph=True)
3
4     dot = nn.utils.parameters_to_vector(grads).mul(v).sum()
5     grads = [g.contiguous() for g in torch.autograd.grad(dot, params,
6     retain_graph=True)]
7     return nn.utils.parameters_to_vector(grads)
```

Gradient Checkpointing

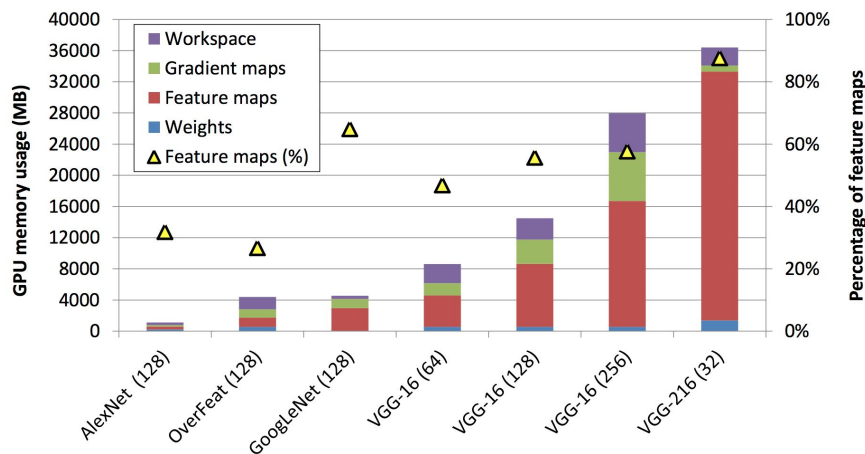
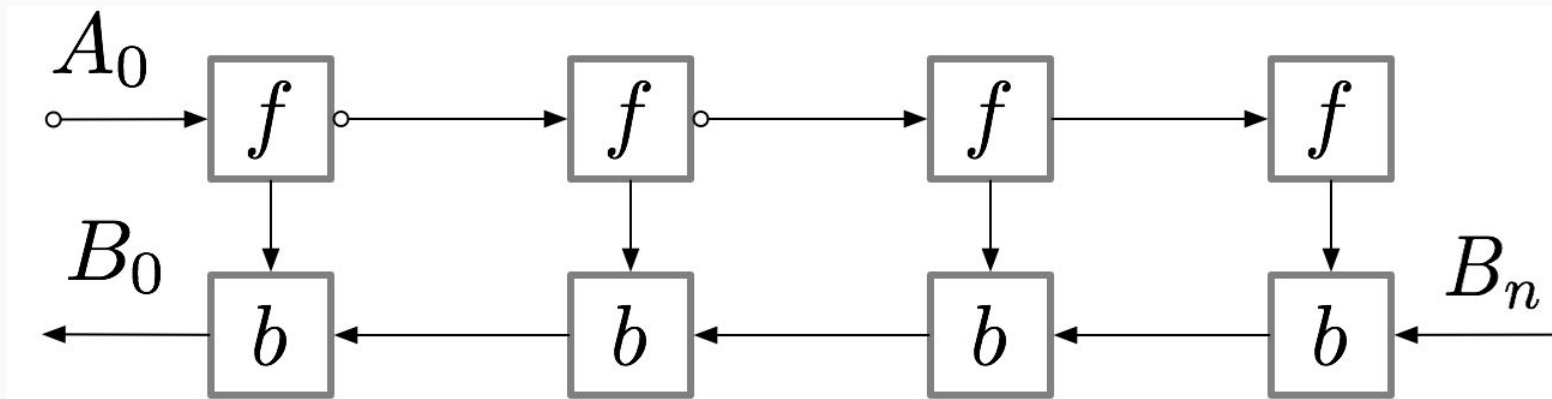


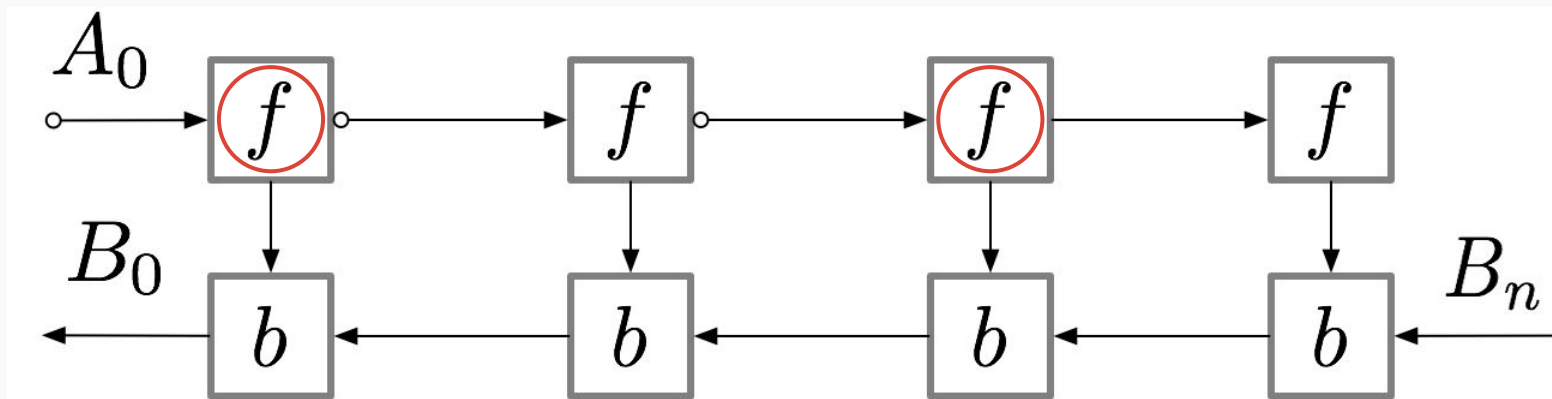
Fig. 4: Breakdown of GPU memory usage based on its functionality (left axis). The right axis shows the fraction of allocated memory consumed by feature maps.

Gradient Checkpointing

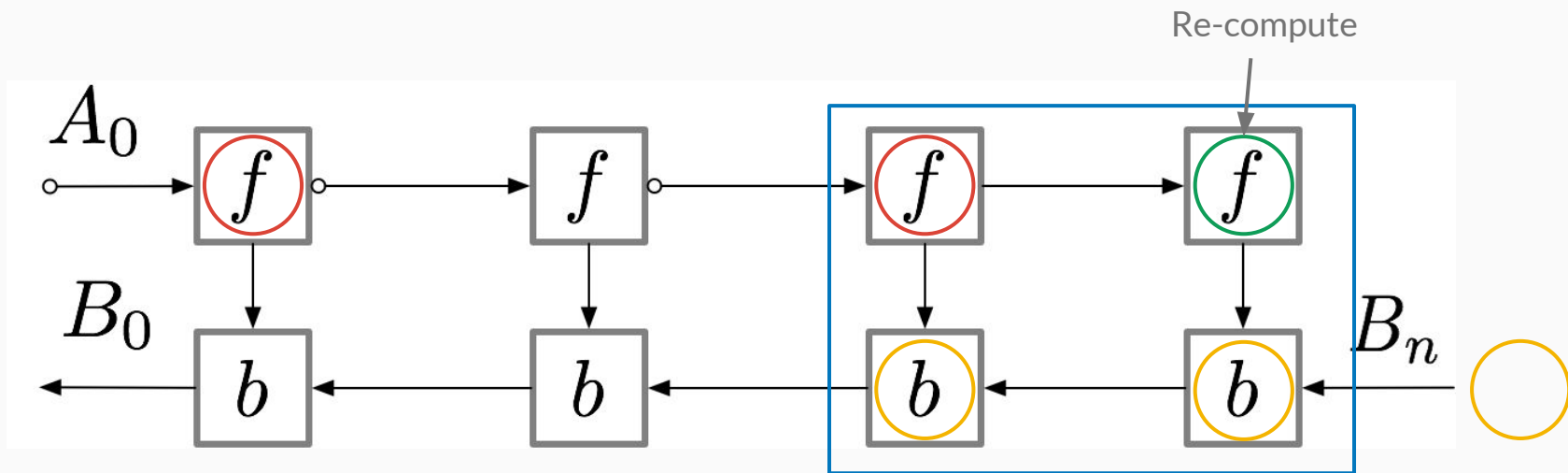


<https://github.com/openai/gradient-checkpointing>

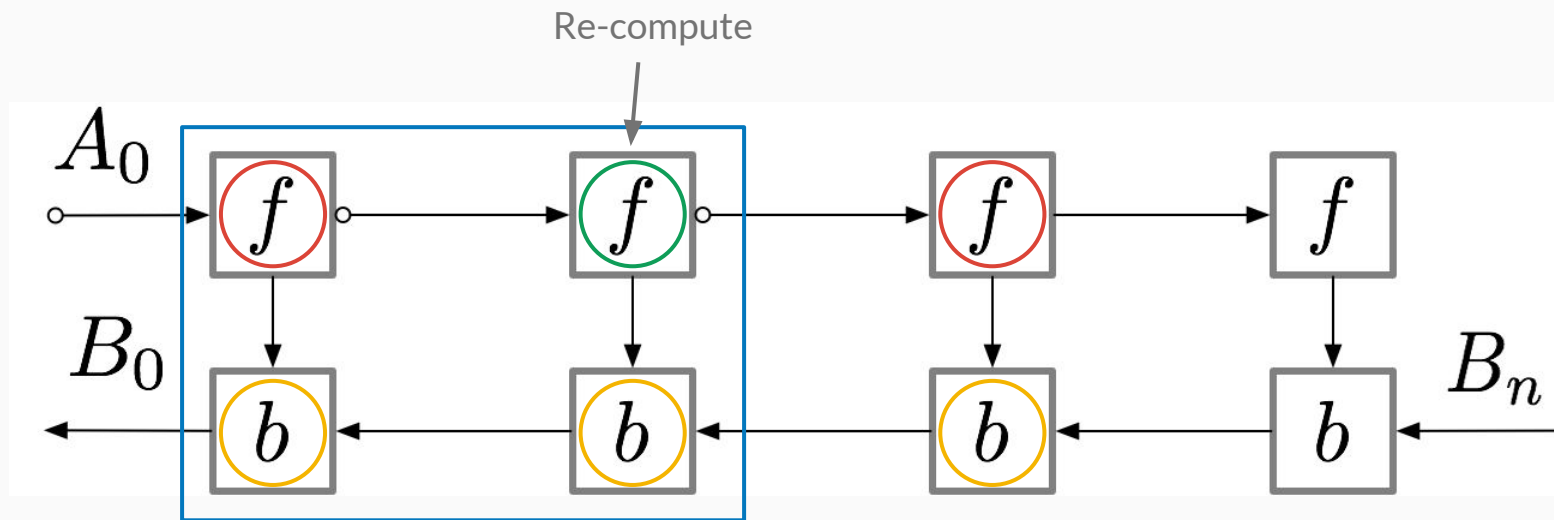
Gradient Checkpointing



Gradient Checkpointing



Gradient Checkpointing



Gradient Checkpointing

