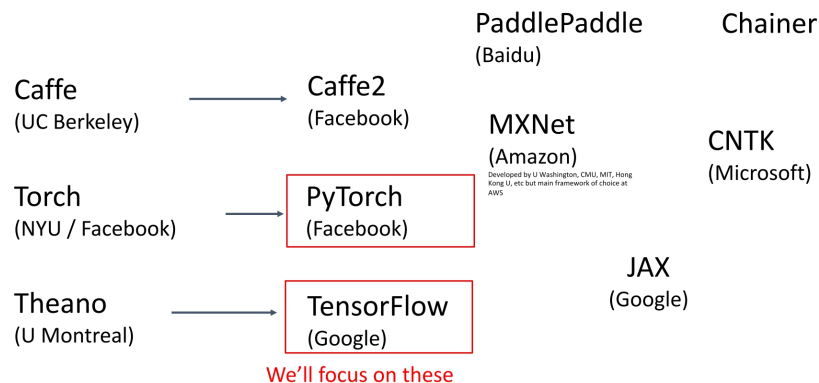


## 12. Lecture 12: Deep Learning Software

### 12.1 Deep Learning Frameworks: Evolution and Landscape

Deep learning software frameworks enable researchers and engineers to efficiently prototype, train, and deploy neural networks. This chapter explores key frameworks, their underlying computational structures, and comparisons between static and dynamic computation graphs. Each framework is providing different trade-offs between usability, performance, and scalability.

A zoo of frameworks!



Justin Johnson

Lecture 12 - 5

February 16, 2022

Figure 12.1: Overview of major deep learning frameworks and their affiliations.

Some notable frameworks include:

- **Caffe** (UC Berkeley) – One of the earliest frameworks, optimized for speed but limited in flexibility.
- **Theano** (U. Montreal) – A pioneer in automatic differentiation, but now discontinued.

- **TensorFlow** (Google) – Popular for production deployments; originally focused on static computation graphs.
- **PyTorch** (Facebook) – An imperative, Pythonic framework with dynamic computation graphs, widely used in research.
- **MXNet** (Amazon) – Developed by multiple institutions, designed for distributed deep learning.
- **JAX** (Google) – A newer framework optimized for high-performance computing and auto-differentiation.

While many frameworks exist, **PyTorch and TensorFlow** dominate deep learning research and deployment. The following sections explore these frameworks in detail, starting with computational graphs and automatic differentiation.

### 12.1.1 The Purpose of Deep Learning Frameworks

Deep learning frameworks provide essential tools that simplify the implementation, training, and deployment of neural networks. They abstract away low-level operations, enabling users to focus on model design and experimentation rather than manual gradient computations or hardware-specific optimizations. The three primary goals of deep learning frameworks are:

- **Rapid Prototyping:** Frameworks allow researchers to quickly experiment with new architectures, optimization techniques, and data pipelines. High-level APIs simplify model definition, while flexible debugging tools enable faster iteration.
- **Automatic Differentiation:** Modern frameworks automatically compute gradients via back-propagation, eliminating the need for manual derivative calculations. This accelerates research and reduces implementation errors.
- **Efficient Execution on Hardware:** Frameworks optimize computations for GPUs & TPUs, leveraging parallel processing and efficient memory management to accelerate training and inference.

### 12.1.2 Recall: Computational Graphs

Recall: Computational Graphs

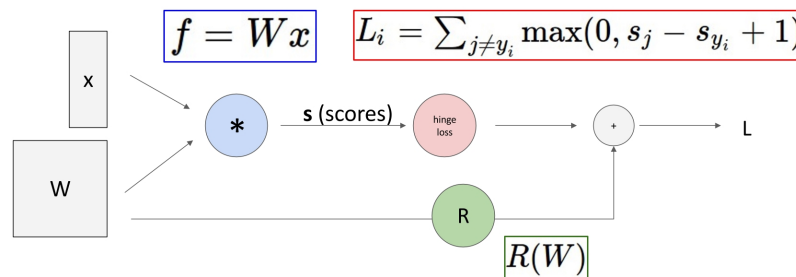


Figure 12.2: **Computational graphs in deep learning.** These graphs define the sequence of operations for training and inference, enabling automatic differentiation and optimization.

Neural networks are represented as **computational graphs**. The graphs define the sequence of operations required to compute outputs and gradients during training. A graph consists of:

- **Nodes:** Represent mathematical operations (e.g., Sigmoid).
- **Edges:** Represent data flow between operations, forming a directed acyclic graph (DAG).

During training, frameworks use computational graphs to:

1. **Forward Pass:** Compute the output by passing data through the graph.
2. **Backward Pass:** Compute gradients via backpropagation, traversing the graph in reverse.
3. **Optimization Step:** Update parameters using computed gradients.

Understanding computational graphs is crucial, as different frameworks implement them in distinct ways. The next sections explore how PyTorch and TensorFlow utilize these graphs, comparing **dynamic** vs. **static** computation strategies.

## 12.2 PyTorch: Fundamental Concepts

PyTorch is a deep learning framework that provides flexibility, dynamic computation graphs, and efficient execution on both CPUs and GPUs. It introduces key abstractions:

- **Tensors:** Multi-dimensional arrays similar to NumPy arrays but capable of running on GPUs.
- **Modules:** Objects representing layers of a neural network, potentially storing learnable parameters.
- **Autograd:** A system that automatically computes gradients by building computational graphs dynamically.

### 12.2.1 Tensors and Basic Computation

To illustrate PyTorch's fundamentals, consider a simple two-layer ReLU network trained using gradient descent on random data.

```

1 import torch
2 device = torch.device('cpu') # Change to 'cuda:0' to run on GPU
3 N, D_in, H, D_out = 64, 1000, 100, 10 # Batch size, input, hidden, output
   ↳ dimensions
4
5 #Create random tensors for data and weights
6 x = torch.randn(N, D_in, device=device)
7 y = torch.randn(N, D_out, device=device)
8 w1 = torch.randn(D_in, H, device=device)
9 w2 = torch.randn(H, D_out, device=device)
10 learning_rate = 1e-6
11
12 for t in range(500):
13     # Forward pass: compute predictions and loss
14     h = x.mm(w1) # Matrix multiply (fully connected layer)
15     h_relu = h.clamp(min=0) # Apply ReLU non-linearity
16     y_pred = h_relu.mm(w2) # Output prediction
17     loss = (y_pred - y).pow(2).sum() # Compute L2 loss
18

```

```

19 # Backward pass: manually compute gradients
20 grad_y_pred = 2.0 * (y_pred - y)
21 grad_w2 = h_relu.t().mm(grad_y_pred)
22 grad_h_relu = grad_y_pred.mm(w2.t())
23 grad_h = grad_h_relu.clone()
24 grad_h[h < 0] = 0 # Backpropagate ReLU
25 grad_w1 = x.t().mm(grad_h)
26
27 # Gradient descent step on weights
28 w1 -= learning_rate * grad_w1 # Gradient update
29 w2 -= learning_rate * grad_w2

```

PyTorch tensors operate efficiently on GPUs by simply setting:  
`device = torch.device('cuda:0')`.

### 12.2.2 Autograd: Automatic Differentiation

PyTorch's **autograd** system automatically builds computational graphs when performing operations on tensors with `requires_grad=True`. These graphs allow automatic computation of gradients via backpropagation.

```

1 x = torch.randn(N, D_in)
2 y = torch.randn(N, D_out)
3 w1 = torch.randn(D_in, H, requires_grad=True)
4 w2 = torch.randn(H, D_out, requires_grad=True)

```

The forward pass remains unchanged:

```

1 h = x.mm(w1)
2 h_relu = h.clamp(min=0)
3 y_pred = h_relu.mm(w2)
4 loss = (y_pred - y).pow(2).sum() # Compute loss

```

PyTorch automatically tracks operations and maintains intermediate values, eliminating the need for manual gradient computation. We backpropagate as follows:

```

1 loss.backward() # Computes gradients for w1 and w2

```

Gradients are accumulated in `w1.grad` and `w2.grad`, so we must clear them manually before the next update:

```

1 with torch.no_grad(): # Prevents unnecessary graph construction
2     w1 -= learning_rate * w1.grad
3     w2 -= learning_rate * w2.grad
4
5     w1.grad.zero_()
6     w2.grad.zero_()

```

Forgetting to reset gradients is a common PyTorch bug, as gradients accumulate by default.



### 12.2.3 Computational Graphs and Modular Computation

PyTorch dynamically constructs **computational graphs** during forward passes, enabling automatic differentiation and backpropagation. Each tensor operation that involves `requires_grad=True` contributes to the computational graph.

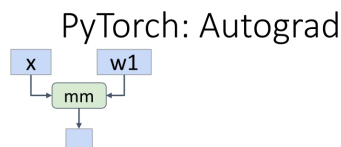
#### Building the Computational Graph

The computation graph begins when we perform operations on tensors with `requires_grad=True`. Consider the following forward pass:

```
1 h = x.mm(w1) # Matrix multiply (fully connected layer)
2 h_relu = h.clamp(min=0) # Apply ReLU non-linearity
3 y_pred = h_relu.mm(w2) # Output prediction
4 loss = (y_pred - y).pow(2).sum() # Compute L2 loss
```

This sequence of operations results in the following computational graph:

- `x.mm(w1)` creates a matrix multiplication node with inputs `x` and `w1`, producing an output tensor with `requires_grad=True`.



Every operation on a tensor with `requires_grad=True` will add to the computational graph, and the resulting tensors will also have `requires_grad=True`

```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

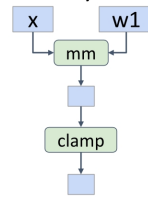
    loss.backward()

    with torch.no_grad():
        w1 -= learning_rate * w1.grad
        w2 -= learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

Figure 12.3: **First computational node in the graph.** The matrix multiplication `x.mm(w1)` creates the first node in the computational graph.

- `.clamp(min=0)` applies a ReLU activation, forming another node.

## PyTorch: Autograd



Every operation on a tensor with `requires_grad=True` will add to the computational graph, and the resulting tensors will also have `requires_grad=True`

```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()

    with torch.no_grad():
        w1 -= learning_rate * w1.grad
        w2 -= learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

Justin Johnson

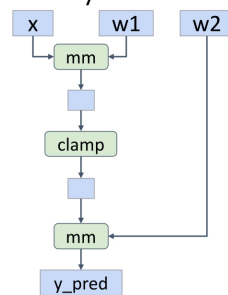
Lecture 12 - 23

February 16, 2022

Figure 12.4: **ReLU activation node.** The ReLU function introduces a non-linearity while maintaining the computational graph structure.

- `.mm(w2)` applies another matrix multiplication, producing the final prediction.

## PyTorch: Autograd



```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()

    with torch.no_grad():
        w1 -= learning_rate * w1.grad
        w2 -= learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

Justin Johnson

Lecture 12 - 24

February 16, 2022

Figure 12.5: **Final matrix multiplication node.** The output prediction `y_pred` is produced by matrix multiplication with `w2`.

### Loss Computation and Backpropagation

After computing the loss, we backpropagate through the graph to compute gradients:

```
1 loss.backward() # Computes gradients for w1 and w2
```

During this process:

- `(y_pred - y)` creates a subtraction node with inputs `y_pred` and `y`.
- `.pow(2)` squares the result, creating a new node.

- `.sum()` sums the squared differences, outputting a scalar loss.

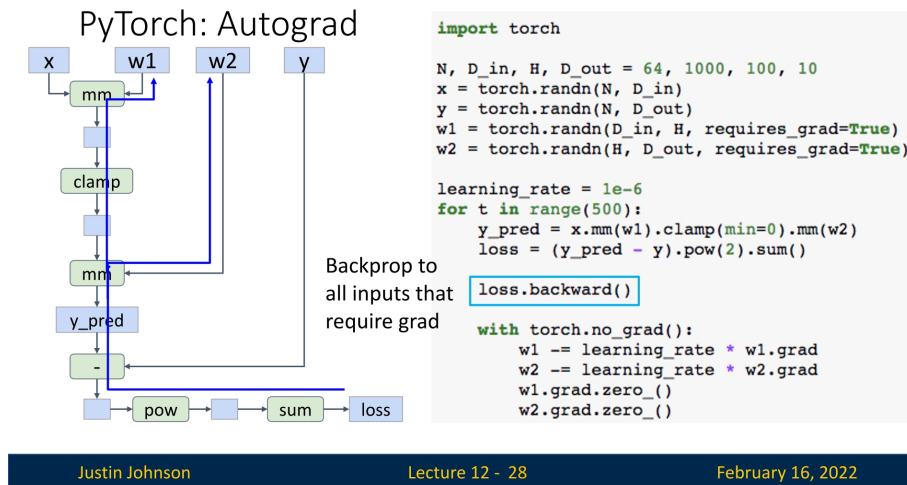


Figure 12.6: **Loss computation node.** The final loss is computed as a scalar output in the computational graph, allowing backpropagation to all inputs requiring gradients.

Once gradients are computed, they are stored in `w1.grad` and `w2.grad`. However, PyTorch accumulates gradients by default, so they must be cleared before the next update (`grad.zero_()`):

```

1 with torch.no_grad(): # Prevents unnecessary graph construction
2     w1 -= learning_rate * w1.grad
3     w2 -= learning_rate * w2.grad
4     w1.grad.zero_()
5     w2.grad.zero_()
  
```

Forgetting to reset gradients is a common mistake in PyTorch. Although probably a design flaw in PyTorch, as we usually don't want to accumulate gradients, we need to be aware of that when we create models.

### Extending Computational Graphs with Python Functions

PyTorch's autograd system allows users to construct computational graphs dynamically using Python functions. When a function is called inside a forward pass, PyTorch records all tensor operations occurring within it.

```

1 def custom_relu(x):
2     return x.clamp(min=0) # Element-wise ReLU
3     h_relu = custom_relu(h)
  
```

Although this function improves code readability, PyTorch still constructs the same computational graph as if we had used `.clamp(min=0)` directly.

### Custom Autograd Functions

PyTorch's automatic differentiation works by building a computational graph out of primitive operations (e.g., add, mul, exp) and then applying the chain rule. In most cases this is sufficient, but sometimes we want:

- To **treat a whole computation as a single semantic unit** in the graph (cleaner, fewer nodes, less bookkeeping).
- To **override the automatically derived backward** with a numerically more stable or more efficient formula.

For this, PyTorch lets us define custom operations by subclassing `torch.autograd.Function` and explicitly specifying forward and backward.

#### Motivating Example: Sigmoid

A naive Python implementation of the sigmoid is:

```
1 def sigmoid(x):
2     return 1.0 / (1.0 + torch.exp(-x))
```

This looks harmless, but it can introduce numerical issues in deep networks:

- For very **large negative** inputs  $x \ll 0$ , we compute `torch.exp(-x) = exp(large positive)`, which overflows to `inf` in `float32`. The forward result is still  $1/(1 + \infty) \approx 0$ , so we might not notice.
- However, during **backward**, autograd differentiates through these primitives and uses the same intermediate `inf` values. Expressions such as  $\frac{\infty}{(1+\infty)^2}$  or  $\infty \cdot 0$  can appear, which numerically become `nan`, even though the true derivative is 0.

Mathematically, the derivative of the sigmoid is

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)),$$

and this is perfectly stable: once we know  $y = \sigma(x) \in (0, 1)$ , the product  $y(1 - y)$  is always bounded in  $[0, 0.25]$  and never overflows. So a more stable strategy is:

1. Compute  $y = \sigma(x)$  in the forward pass.
2. Save  $y$ .
3. Compute the gradient in backward using  $y(1 - y)$  instead of recomputing exponentials.

This is exactly what a custom autograd function allows us to do.

```
1 class Sigmoid(torch.autograd.Function):
2     @staticmethod
3     def forward(ctx, x):
4         # Forward as usual (PyTorch's built-in sigmoid is already stable;
5         # here we reimplement it for illustration).
6         y = 1.0 / (1.0 + torch.exp(-x))
7         # Save only the stable output y for backward.
8         ctx.save_for_backward(y)
9         return y
```



```

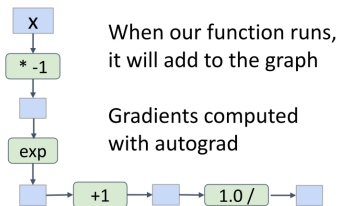
1 @staticmethod
2     def backward(ctx, grad_y):
3         # Retrieve saved output
4         (y,) = ctx.saved_tensors
5         # Use the stable formula seen earlier
6         grad_x = grad_y * y * (1.0 - y)
7         return grad_x
8
9 def sigmoid(x):
10     return Sigmoid.apply(x)

```

### PyTorch: New functions

Can define new operations  
using Python functions

```
def sigmoid(x):
    return 1.0 / (1.0 + (-x).exp())
```



Define new autograd operators  
by subclassing Function, define  
forward and backward

```

class Sigmoid(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x):
        y = 1.0 / (1.0 + (-x).exp())
        ctx.save_for_backward(y)
        return y

    @staticmethod
    def backward(ctx, grad_y):
        y, = ctx.saved_tensors
        grad_x = grad_y * y * (1.0 - y)
        return grad_x

def sigmoid(x):
    return Sigmoid.apply(x)

```

Now when our function runs,  
it adds one node to the graph!



Justin Johnson

Lecture 12 - 36

February 16, 2022

Figure 12.7: **Custom autograd function for sigmoid.** Left: the naive implementation expands into several primitive nodes (exp, add, div), each with its own backward. Right: the custom Sigmoid is a single node with a hand-crafted, numerically stable backward.

Once defined, we can use the new sigmoid as any other PyTorch operation:

```

1 x = torch.randn(10, requires_grad=True)
2 sigmoid_out = sigmoid(x)
3 sigmoid_out.sum().backward()

```

In practice, this level of control is rarely needed for basic operations: PyTorch's built-in functions (`torch.sigmoid`, `torch.softmax`, etc.) are already implemented internally using optimized and stable autograd functions.

Custom Functions become most useful when implementing new layers, composite operations, or specialized losses where we know a better backward formula than the one autograd would derive automatically.

**Summary: Backpropagation and Graph Optimization**

- **Any operation on a tensor with `requires_grad=True` extends the computational graph.**
- **PyTorch dynamically records these operations** and stores just enough context (saved tensors) to evaluate gradients efficiently via the chain rule.
- **Forgetting to reset gradients** (e.g., omitting `optimizer.zero_grad()`) causes gradients to accumulate across iterations, leading to incorrect updates.
- **Graph structure can be optimized** using custom autograd functions: they fuse multiple primitive ops into a single node, can implement numerically stable backward formulas, and provide more meaningful graph semantics than low-level primitives alone.

A solid understanding of PyTorch’s computational graphs—and how to customize them when necessary—is essential for debugging, improving numerical robustness, and optimizing the performance of deep learning models.

**12.2.4 High-Level Abstractions in PyTorch: `torch.nn` and Optimizers**

PyTorch provides a high-level wrapper, `torch.nn`, which simplifies neural network construction by offering an object-oriented API for defining models. This abstraction allows for more structured and maintainable code, making deep learning models easier to build and extend.

**Using `torch.nn.Sequential`**

The `torch.nn.Sequential` container allows defining models as a sequence of layers. Below, we define a simple two-layer network with ReLU activation:

```

1 import torch
2
3 N, D_in, H, D_out = 64, 1000, 100, 10
4 x = torch.randn(N, D_in)
5 y = torch.randn(N, D_out)
6
7 model = torch.nn.Sequential(
8     torch.nn.Linear(D_in, H),
9     torch.nn.ReLU(),
10    torch.nn.Linear(H, D_out)
11 )
12
13 learning_rate = 1e-2
14 for t in range(500):
15     y_pred = model(x)
16     loss = torch.nn.functional.mse_loss(y_pred, y)
17     loss.backward()
18
19     with torch.no_grad():
20         for param in model.parameters():
21             param -= learning_rate * param.grad
22
23 model.zero_grad()
```

- The `model` object is a container holding layers. Each layer manages its own parameters.
- Calling `model(x)` performs the forward pass.
- The loss is computed using `torch.nn.functional.mse_loss()`.

- Calling `loss.backward()` computes gradients for all model parameters.
- Parameter updates are performed manually in a loop over `model.parameters()`.
- Calling `model.zero_grad()` resets gradients for all parameters.

### Using Optimizers: Automating Gradient Descent

Instead of manually implementing gradient descent, PyTorch provides optimizer classes that handle parameter updates. Below, we use the Adam optimizer:

```
1 import torch
2
3 N, D_in, H, D_out = 64, 1000, 100, 10
4 x = torch.randn(N, D_in)
5 y = torch.randn(N, D_out)
6
7 model = torch.nn.Sequential(
8     torch.nn.Linear(D_in, H),
9     torch.nn.ReLU(),
10    torch.nn.Linear(H, D_out)
11 )
12
13 learning_rate = 1e-4
14 optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
15
16 for t in range(500):
17     y_pred = model(x)
18     loss = torch.nn.functional.mse_loss(y_pred, y)
19     loss.backward()
20
21     optimizer.step()
22     optimizer.zero_grad()
```

- The optimizer is instantiated with `torch.optim.Adam()` and receives model parameters.
- Calling `optimizer.step()` updates all parameters automatically.
- Calling `optimizer.zero_grad()` resets gradients before the next step.

This approach is both cleaner and less error-prone than manual updates.

### Defining Custom `nn.Module` Subclasses

For more complex architectures, we can define custom `nn.Module` subclasses:

```

1 import torch
2
3 class TwoLayerNet(torch.nn.Module):
4     def __init__(self, D_in, H, D_out):
5         super(TwoLayerNet, self).__init__()
6         self.linear1 = torch.nn.Linear(D_in, H)
7         self.linear2 = torch.nn.Linear(H, D_out)
8
9     def forward(self, x):
10         h_relu = self.linear1(x).clamp(min=0)
11         y_pred = self.linear2(h_relu)
12         return y_pred
13
14 N, D_in, H, D_out = 64, 1000, 100, 10
15 x = torch.randn(N, D_in)
16 y = torch.randn(N, D_out)
17
18 model = TwoLayerNet(D_in, H, D_out)
19 optimizer = torch.optim.SGD(model.parameters(), lr=1e-4)
20
21 for t in range(500):
22     y_pred = model(x)
23     loss = torch.nn.functional.mse_loss(y_pred, y)
24     loss.backward()
25
26     optimizer.step()
27     optimizer.zero_grad()

```

- **Model Initialization:** The `__init__` method defines layers as class attributes.
- **Forward Pass:** The `forward()` method specifies how inputs are transformed.
- **Autograd Integration:** PyTorch automatically tracks gradients for model parameters.
- **Training Loop:** The optimizer updates weights based on computed gradients.

### Key Takeaways

- `torch.nn.Sequential` simplifies defining networks as a stack of layers.
- **Optimizers automate gradient descent**, making training loops cleaner.
- **Custom `nn.Module` subclasses** provide flexibility for complex architectures.
- **Autograd handles differentiation automatically**, eliminating the need for manual backward computations.

Using `torch.nn` and optimizers streamlines model development, making PyTorch a powerful and expressive framework for deep learning.

#### 12.2.5 Combining Custom Modules with Sequential Models

A common practice in PyTorch is to combine custom `nn.Module` subclasses with `torch.nn.Sequential` containers. This enables modular and scalable architectures while maintaining the expressiveness of object-oriented model design.



**Example: Parallel Block**

The following example defines a `ParallelBlock` module that applies two linear transformations to the input independently and then multiplies the results element-wise:

```

1 import torch
2
3 class ParallelBlock(torch.nn.Module):
4     def __init__(self, D_in, D_out):
5         super(ParallelBlock, self).__init__()
6         self.linear1 = torch.nn.Linear(D_in, D_out)
7         self.linear2 = torch.nn.Linear(D_in, D_out)
8
9     def forward(self, x):
10        h1 = self.linear1(x)
11        h2 = self.linear2(x)
12        return (h1 * h2).clamp(min=0) # Element-wise multiplication followed
        ↪ by ReLU
13
14 N, D_in, H, D_out = 64, 1000, 100, 10
15 x = torch.randn(N, D_in)
16 y = torch.randn(N, D_out)
17
18 model = torch.nn.Sequential(
19     ParallelBlock(D_in, H),
20     ParallelBlock(H, H),
21     torch.nn.Linear(H, D_out)
22 )
23
24 optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
25
26 for t in range(500):
27     y_pred = model(x)
28     loss = torch.nn.functional.mse_loss(y_pred, y)
29     loss.backward()
30     optimizer.step()
31     optimizer.zero_grad()

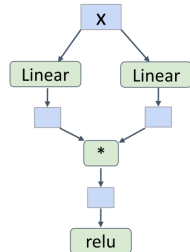
```

- The `ParallelBlock` applies two separate linear layers to the input.
- The outputs are multiplied element-wise before applying ReLU.
- The `Sequential` container stacks multiple `ParallelBlock` instances, followed by a final linear layer.
- Using this approach allows rapid experimentation with modular neural network components.

Although this example is not very smart and not thing we should in practice, it demonstrates well the ability to create building blocks using torch and thus create using this abstraction some complex neural networks with ease.

## PyTorch: nn Defining Modules

Define network component  
as a Module subclass



```
import torch

class ParallelBlock(torch.nn.Module):
    def __init__(self, D_in, D_out):
        super(ParallelBlock, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, D_out)
        self.linear2 = torch.nn.Linear(D_in, D_out)

    def forward(self, x):
        h1 = self.linear1(x)
        h2 = self.linear2(x)
        return (h1 * h2).clamp(min=0)

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

model = torch.nn.Sequential(
    ParallelBlock(D_in, H),
    ParallelBlock(H, H),
    torch.nn.Linear(H, D_out))

optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = torch.nn.functional.mse_loss(y_pred, y)
    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

Justin Johnson

Lecture 12 - 51

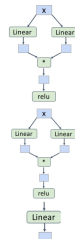
February 16, 2022

Figure 12.8: **ParallelBlock module design:** The implementation of the ParallelBlock and its corresponding computational graph visualization.

## PyTorch: nn Defining Modules

Stack multiple instances of the  
component in a sequential

Very easy to quickly  
build complex network  
architectures!



```
import torch

class ParallelBlock(torch.nn.Module):
    def __init__(self, D_in, D_out):
        super(ParallelBlock, self).__init__()
        self.linear1 = torch.nn.Linear(D_in, D_out)
        self.linear2 = torch.nn.Linear(D_in, D_out)

    def forward(self, x):
        h1 = self.linear1(x)
        h2 = self.linear2(x)
        return (h1 * h2).clamp(min=0)

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

model = torch.nn.Sequential(
    ParallelBlock(D_in, H),
    ParallelBlock(H, H),
    torch.nn.Linear(H, D_out))

optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
for t in range(500):
    y_pred = model(x)
    loss = torch.nn.functional.mse_loss(y_pred, y)
    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

Justin Johnson

Lecture 12 - 52

February 16, 2022

Figure 12.9: **Stacking multiple ParallelBlock instances in a Sequential model.** The left side of the figure shows the computational graph produced.

### 12.2.6 Efficient Data Loading with `torch.utils.data`

Training deep neural networks efficiently requires a robust data pipeline. PyTorch provides the `torch.utils.data` module, which abstracts away data loading, shuffling, batching, and parallelization—ensuring that model computation and data preparation can run concurrently. The two key components are:

- **Dataset:** Represents a collection of samples. You can use built-in classes like `TensorDataset` for in-memory tensors or implement a custom `Dataset` that reads from files or databases.
- **DataLoader:** Wraps a `Dataset` to provide mini-batching, shuffling, and multi-process data loading. It also supports pinned memory for faster GPU transfer.

**Example: Using DataLoader for Mini-batching**

The example below demonstrates how to use DataLoader with synthetic data for mini-batch training.

```

1      import torch
2      from torch.utils.data import TensorDataset, DataLoader
3
4      # 1. Create a simple in-memory dataset
5      N, D_in, H, D_out = 64, 1000, 100, 10
6      x = torch.randn(N, D_in)
7      y = torch.randn(N, D_out)
8
9      dataset = TensorDataset(x, y)
10
11     # 2. Create a DataLoader with batching and parallel loading
12     loader = DataLoader(
13         dataset,
14         batch_size=8,
15         shuffle=True,          # Shuffle each epoch for stable training
16         num_workers=2,        # Parallel CPU workers for background loading
17         pin_memory=True       # Speeds up host→GPU transfers
18     )
19
20     model = TwoLayerNet(D_in, H, D_out)
21     optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
22
23     # 3. Training loop using the DataLoader
24     for epoch in range(20):
25         for x_batch, y_batch in loader:
26             y_pred = model(x_batch)
27             loss = torch.nn.functional.mse_loss(y_pred, y_batch)
28
29             loss.backward()
30             optimizer.step()
31             optimizer.zero_grad()

```

This setup automatically handles mini-batch creation, shuffling, and memory prefetching. With `num_workers > 0`, the CPU preloads data while the GPU trains on the previous batch, preventing GPU idle time—a crucial optimization for large datasets.

**Best Practices**

- Use `shuffle=True` to avoid order bias and improve gradient diversity.
- Adjust `num_workers` to match your CPU cores (typical range: 2–8) for best throughput.
- Set `pin_memory=True` when training on GPU to accelerate host–device transfers.

**Handling Multiple Datasets**

In practice, data often comes from multiple sources—different domains, modalities, or tasks. PyTorch offers flexible tools to combine and balance these datasets efficiently.

### Concatenating Datasets

When datasets share the same structure (e.g., same feature dimensions), use `ConcatDataset` to merge them into a single unified dataset.

```
1 from torch.utils.data import ConcatDataset, DataLoader
2
3 dataset_a = TensorDataset(torch.randn(100, 20), torch.randn(100, 1))
4 dataset_b = TensorDataset(torch.randn(200, 20), torch.randn(200, 1))
5
6 combined = ConcatDataset([dataset_a, dataset_b])
7
8 loader = DataLoader(
9     combined,
10    batch_size=16,
11    shuffle=True,
12    num_workers=4
13 )
```

This approach interleaves samples from all datasets proportionally to their sizes. It is ideal for combining related sources (e.g., merging multiple corpora or image datasets).

### Weighted Sampling Across Datasets

If some datasets are much smaller or more important, you can balance sampling probabilities using `WeightedRandomSampler`. This ensures underrepresented data appears more frequently in training batches.

```
1 from torch.utils.data import WeightedRandomSampler
2
3 # Example: emphasize smaller dataset (dataset_a)
4 weights = [1.0 / len(dataset_a)] * len(dataset_a) + \
5 [1.0 / len(dataset_b)] * len(dataset_b)
6
7 sampler = WeightedRandomSampler(weights, num_samples=len(weights),
8     ↪ replacement=True)
9
10 balanced_loader = DataLoader(
11     combined,
12     batch_size=16,
13     sampler=sampler,
14     num_workers=4
15 )
```

Weighted sampling is especially useful for:

- **Imbalanced datasets.** For example, when rare classes need more representation during training.
- **Multi-source training.** Combining labeled and unlabeled data or datasets from distinct domains.
- **Curriculum learning.** Gradually increasing sample difficulty or diversity over time.



*Streaming or Multi-modal Data*

For more dynamic or heterogeneous sources (e.g., loading text and image pairs), subclass `IterableDataset` to yield samples from multiple streams in real time, or define a custom `Sampler` to coordinate multi-modal alignment.

```
1 from torch.utils.data import IterableDataset
2
3 class MultiSourceStream(IterableDataset):
4     def __iter__(self):
5         for x_img, x_txt in zip(image_stream(), text_stream()):
6             yield preprocess(x_img, x_txt)
```

This design is common in large-scale vision–language or multi-task training pipelines, where data arrives asynchronously or from external APIs.

*Summary*

`DataLoader` and its related utilities form the backbone of efficient training in PyTorch. They decouple data I/O from model computation, provide clean abstractions for multi-source or imbalanced data, and make large-scale experiments reproducible and scalable across CPUs and GPUs.

### 12.2.7 Using Pretrained Models with TorchVision

PyTorch provides access to many pretrained models through the `torchvision` package, making it easy to leverage existing architectures for various vision tasks.

Using pretrained models is as simple as:

```
1 import torchvision.models as models
2
3 alexnet = models.alexnet(pretrained=True)
4 vgg16 = models.vgg16(pretrained=True)
5 resnet101 = models.resnet101(pretrained=True)
```

- These models come with pretrained weights on ImageNet, making them suitable for transfer learning.
- Fine-tuning pretrained models often leads to faster convergence and better performance on new tasks.
- `torchvision.models` provides a wide variety of architectures beyond AlexNet, VGG, and ResNet.

**Key Takeaways**

- **Custom modules and** `torch.nn.Sequential` **can be combined** to quickly build complex models while maintaining modularity.
- **Data loading utilities** such as `torch.utils.data.DataLoader` facilitate efficient mini-batching and dataset management.
- **TorchVision provides pretrained models**, making it easy to leverage state-of-the-art architectures for various vision tasks.

### 12.3 Dynamic vs. Static Computational Graphs in PyTorch

A fundamental design choice in PyTorch is its use of **dynamic computational graphs**. Unlike static graphs, which are constructed once and reused, PyTorch builds a fresh computational graph for each forward pass. Once `loss.backward()` is called, the graph is discarded, and a new one is constructed in the next iteration.

While dynamically building graphs in every iteration may seem inefficient, this approach provides a crucial advantage: *the ability to use standard Python control flow during model execution*. This enables complex architectures that modify their behavior on-the-fly based on intermediate results.

#### Example: Dynamic Graph Construction

Consider a model where the choice of weight matrix for backpropagation depends on the previous loss value. This scenario, though impractical, demonstrates PyTorch's ability to create different computational graphs in each iteration.

#### PyTorch: Dynamic Computation Graphs

Dynamic graphs let you use regular Python control flow during the forward pass!

Decide which one to use at each layer based on loss at previous iteration

(this model doesn't make sense! Just a simple dynamic example)

```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2a = torch.randn(H, D_out, requires_grad=True)
w2b = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
prev_loss = 5.0
for t in range(500):
    w2 = w2a if prev_loss < 5.0 else w2b
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
    prev_loss = loss.item()
```

Justin Johnson

Lecture 12 - 68

February 16, 2022

Figure 12.10: **Example of a dynamically constructed graph:** The model structure changes at each iteration based on previous loss values.

In dynamic graphs, every forward pass constructs a unique computation graph, allowing for models with **varying execution paths** across different iterations.

#### 12.3.1 Static Graphs and Just-in-Time (JIT) Compilation

In contrast, **static computational graphs** follow a two-step process:

1. **Graph Construction:** Define the computational graph once, allowing the framework to optimize it before execution.
2. **Graph Execution:** The same pre-optimized graph is reused for all forward passes.

While PyTorch natively operates with dynamic graphs, it also supports static graphs through **TorchScript** using **Just-in-Time (JIT) compilation**. This allows PyTorch to analyze the model's source code, compile it into an optimized static graph, and reuse it for improved efficiency.

### 12.3.2 Using JIT to Create Static Graphs

To convert a function into a static computational graph, PyTorch provides `torch.jit.script()`:

```
1 import torch
2
3 def model(x):
4     return x * torch.sin(x)
5
6 scripted_model = torch.jit.script(model)  # Convert to static graph
```

Alternatively, PyTorch allows automatic graph compilation using the `@torch.jit.script` annotation:

```
1 import torch
2
3 @torch.jit.script
4 def model(x):
5     return x * torch.sin(x)
```

### PyTorch: Static Graphs with JIT

Even easier: add **annotation** to function, Python function compiled to a graph when it is defined

Calling function uses graph

```
import torch

@torch.jit.script
def model(x, y, w1, w2a, w2b, prev_loss):
    w2 = w2a if prev_loss < 5.0 else w2b
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    return loss

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2a = torch.randn(H, D_out, requires_grad=True)
w2b = torch.randn(H, D_out, requires_grad=True)

prev_loss = 5.0
learning_rate = 1e-6
for t in range(500):
    loss = model(x, y, w1, w2a, w2b, prev_loss)

    loss.backward()
    prev_loss = loss.item()
```

Justin Johnson

Lecture 12 - 74

February 16, 2022

Figure 12.11: **TorchScript**: Using JIT compilation to convert PyTorch models into static graphs for optimization.

### 12.3.3 Handling Conditionals in Static Graphs

Static graphs struggle with conditionals because they are typically **fixed at compile time**. However, PyTorch's JIT can represent conditionals as graph nodes, enabling runtime flexibility.

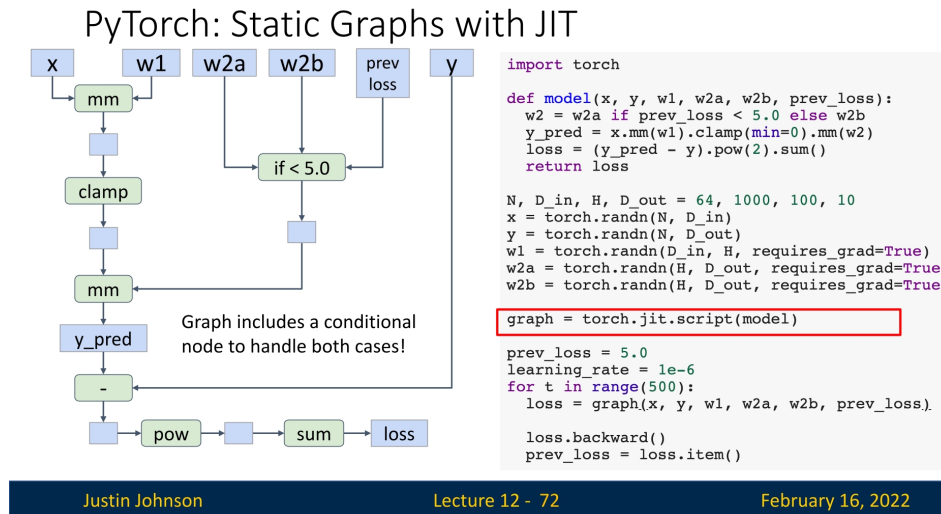


Figure 12.12: **Conditionals in static graphs:** JIT inserts a conditional node to handle different execution paths.

This allows some degree of flexibility while retaining the benefits of graph optimization.

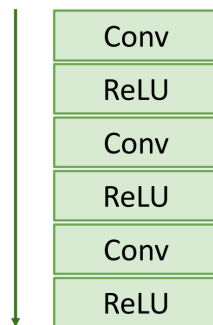
#### 12.3.4 Optimizing Computation Graphs with JIT

One advantage of static graphs is that they enable **graph-level optimizations**. PyTorch JIT can automatically **fuse operations** such as convolution and activation layers into a single efficient operation.

#### Static vs Dynamic Graphs: Optimization

With static graphs, framework can **optimize** the graph for you before it runs!

The graph you wrote



Equivalent graph with **fused operations**

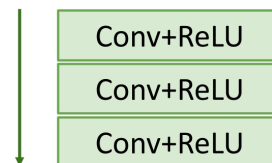


Figure 12.13: **Operation fusion in static graphs:** Layers such as Conv + ReLU are combined into a single operation to improve efficiency.

This optimization is performed once, eliminating the need to optimize in every iteration.



### 12.3.5 Benefits and Limitations of Static Graphs

#### Advantages of Static Graphs:

- **Graph Optimization:** The framework optimizes computation before execution, improving speed.
- **Operation Fusion:** Frequently used layers (e.g., Conv + ReLU) are merged into a single operation.
- **Serialization:** Models can be saved to disk and loaded in non-Python environments (e.g., C++).

#### Challenges of Static Graphs:

- **Difficult Debugging:** Debugging static graphs can be challenging due to indirection between graph construction and execution.
- **Less Flexibility:** Unlike dynamic graphs, static graphs struggle with models that modify their execution path.
- **Rebuilding Required:** Any model change requires reconstructing the entire graph.

### 12.3.6 When Are Dynamic Graphs Necessary?

Certain architectures *require* dynamic graphs due to their execution dependencies on input data:

- **Recurrent Neural Networks (RNNs):** The number of computation steps depends on input sequence length.
- **Recursive Networks:** Hierarchical models, such as parse trees in NLP, require dynamic execution paths.
- **Modular Networks:** Some architectures dynamically select which sub-network to execute.

A well-known example is the model in [270], where part of the network predicts which module should execute next.

## 12.4 TensorFlow: Dynamic and Static Computational Graphs

TensorFlow originally adopted **static computational graphs** by default (TensorFlow 1.0), requiring users to explicitly define a computation graph before running it. However, in **TensorFlow 2.0**, the framework transitioned to **dynamic graphs** by default, making the API more similar to PyTorch. This shift caused a significant divide in the TensorFlow ecosystem, as older static-graph code intertwined with newer dynamic-graph code, creating confusion and bugs.

### 12.4.1 Defining Computational Graphs in TensorFlow 2.0

In PyTorch, the computational graph is built *implicitly*: any operation performed on a tensor with `requires_grad=True` is automatically tracked. TensorFlow 2.0 (TF2), by contrast, introduced **eager execution** as the default mode—operations execute immediately like standard Python code, producing concrete values rather than symbolic graph nodes. This makes TF2 intuitive and debuggable but requires an explicit mechanism for recording operations when gradients are needed. That mechanism is the `tf.GradientTape`.

#### Understanding `tf.GradientTape`

The `GradientTape` is TensorFlow’s dynamic autodiff engine, analogous to PyTorch’s implicit `autograd`. It acts like a “*recorder*”: while active, it logs all operations on watched tensors (typically all `tf.Variable` objects) and can later “play back” those operations to compute gradients.

- Entering a `with tf.GradientTape() as tape:` block begins recording.

- Any operation involving watched variables is logged on the tape.
- Exiting the block stops recording.
- Calling `tape.gradient(loss, [vars])` replays the tape backward to compute exact gradients via the chain rule.

This explicit opt-in design prevents unnecessary gradient tracking (e.g., during inference) and gives developers fine-grained control over which computations are differentiable.

```

1 import tensorflow as tf
2
3 # Setup data and parameters
4 N, Din, H, Dout = 16, 1000, 100, 10
5 x = tf.random.normal((N, Din))
6 y = tf.random.normal((N, Dout))
7 w1 = tf.Variable(tf.random.normal((Din, H)))
8 w2 = tf.Variable(tf.random.normal((H, Dout)))
9
10 learning_rate = 1e-6
11
12 for t in range(1000):
13     # Begin recording operations on the tape
14     with tf.GradientTape() as tape:
15         h = tf.maximum(tf.matmul(x, w1), 0) # ReLU
16         y_pred = tf.matmul(h, w2)
17         diff = y_pred - y
18         loss = tf.reduce_mean(tf.reduce_sum(diff ** 2, axis=1))
19
20     # Compute gradients of loss w.r.t parameters
21     grad_w1, grad_w2 = tape.gradient(loss, [w1, w2])
22
23     # Parameter updates (in-place, safe for tf.Variables)
24     w1.assign_sub(learning_rate * grad_w1)
25     w2.assign_sub(learning_rate * grad_w2)

```

This process mirrors PyTorch’s autograd but with more explicit control: `GradientTape` defines the graph’s lifetime (inside the `with` block), rather than relying on implicit global tracking. The resulting computation graph is ephemeral—destroyed after gradient computation unless the tape is declared as `persistent=True` (allowing multiple gradient calls).

#### Key differences from PyTorch

- PyTorch automatically tracks gradients for all tensors with `requires_grad=True`. TensorFlow records only within the `GradientTape` context.
- TensorFlow’s graph is discarded after use unless marked persistent.
- `GradientTape` offers fine-grained control: you can record subsets of operations or specific variables only.

### 12.4.2 Static Graphs with `@tf.function`

While TF2 defaults to eager (imperative) execution for flexibility, static computation graphs are still essential for deployment and optimization. To combine both worlds, TensorFlow introduces the `@tf.function` decorator, which traces Python functions into optimized static graphs—comparable to `torch.jit.script()` in PyTorch.

#### *Motivation*

Eager execution simplifies experimentation but adds Python overhead per operation. Static graphs, on the other hand, allow TensorFlow to perform ahead-of-time optimizations: operation fusion (e.g., combining `matmul` + `bias_add`), kernel selection, memory reuse, and XLA compilation. Using `@tf.function`, developers write natural Python code while TensorFlow transparently traces and compiles it.

```
1 @tf.function # Compiles to a static graph on first call
2 def training_step(x, y, w1, w2, lr):
3     with tf.GradientTape() as tape:
4         h = tf.maximum(tf.matmul(x, w1), 0)
5         y_pred = tf.matmul(h, w2)
6         loss = tf.reduce_mean(tf.reduce_sum((y_pred - y) ** 2, axis=1))
7
8     grad_w1, grad_w2 = tape.gradient(loss, [w1, w2])
9     w1.assign_sub(lr * grad_w1)
10    w2.assign_sub(lr * grad_w2)
11    return loss
12
13 # Regular Python loop, but graph executes under the hood
14 for t in range(1000):
15     current_loss = training_step(x, y, w1, w2, learning_rate)
```

Here, `@tf.function` traces the computation during its first execution, then caches the resulting static graph for reuse—removing Python overhead and enabling runtime optimizations. This achieves up to 2–10× speedups for heavy workloads while preserving eager-like syntax.

#### *Summary of Modes*

- **Eager mode.** Operations run immediately, ideal for debugging and experimentation.
- **GradientTape.** Dynamically records operations for automatic differentiation, similar to PyTorch’s `autograd`.
- **@tf.function.** Converts eager code into a reusable static graph, fusing and optimizing operations for deployment.

Together, these tools give TensorFlow 2.0 both the interactivity of PyTorch and the performance advantages of static compilation—bridging the flexibility–efficiency trade-off that defined earlier deep learning frameworks.

## 12.5 Keras: High-Level API for TensorFlow

Keras provides a high-level API for building deep learning models, simplifying working with models.

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import InputLayer, Dense
4
5 N, Din, H, Dout = 16, 1000, 100, 10
6
7 model = Sequential([
8     InputLayer(input_shape=(Din,)),
9     Dense(units=H, activation='relu'),
10    Dense(units=Dout)
11 ])
12
13 loss_fn = tf.keras.losses.MeanSquaredError()
14 opt = tf.keras.optimizers.SGD(learning_rate=1e-6)
15
16 x = tf.random.normal((N, Din))
17 y = tf.random.normal((N, Dout))
18
19 for t in range(1000):
20     with tf.GradientTape() as tape:
21         y_pred = model(x)
22         loss = loss_fn(y_pred, y)
23         grads = tape.gradient(loss, model.trainable_variables)
24         opt.apply_gradients(zip(grads, model.trainable_variables))
```

Keras simplifies training by providing:

- **Predefined layers:** Easily stack layers with `Sequential()`.
- **Common loss functions and optimizers:** Use built-in losses and optimizers like Adam.
- **Automatic gradient handling:** `opt.apply_gradients()` simplifies parameter updates.

We can further simplify the training loop using `opt.minimize()` by defining a step function:

```
1 def step():
2     y_pred = model(x)
3     loss = loss_fn(y_pred, y)
4     return loss
5
6 for t in range(1000):
7     opt.minimize(step, model.trainable_variables)
```

## 12.6 TensorBoard: Visualizing Training Metrics

### TensorBoard

Add logging to code to record loss, stats, etc  
Run server and get pretty graphs!



Justin Johnson

Lecture 12 - 101

February 16, 2022

Figure 12.14: **TensorBoard visualization:** Loss curves and weight distributions during training.

**TensorBoard** is a visualization tool that helps monitor deep learning experiments. It allows users to track:

- Loss curves and accuracy during training.
- Weight distributions and parameter updates.
- Computational graphs of the model.

While originally designed for TensorFlow, TensorBoard now support **PyTorch** via the `torch.utils.tensorboard` API. However, modern alternatives such as **Weights and Biases (wandb)** and **MLFlow** provide additional functionality, making them popular choices for tracking experiments.

## 12.7 Comparison: PyTorch vs. TensorFlow

- **PyTorch:**
  - Imperative API that is easy to debug.
  - Dynamic computation graphs enable flexibility.
  - `torch.jit.script()` allows for static graph compilation.
  - Harder to optimize for TPUs.
  - Deployment on mobile is less streamlined.
- **TensorFlow 1.0:**
  - Static graphs by default.
  - Faster execution but difficult debugging.
  - API inconsistencies made it less user-friendly.
- **TensorFlow 2.0:**
  - Defaulted to dynamic graphs, similar to PyTorch.
  - Standardized Keras API for ease of use.
  - Still retains static graph capability with `tf.function`.

*Conclusion*

Both PyTorch and TensorFlow 2.0 now support both dynamic and static graphs, offering flexibility for different use cases. PyTorch remains the preferred choice for research due to its intuitive imperative style, while TensorFlow is still widely used in production, particularly in environments requiring static graph optimization.