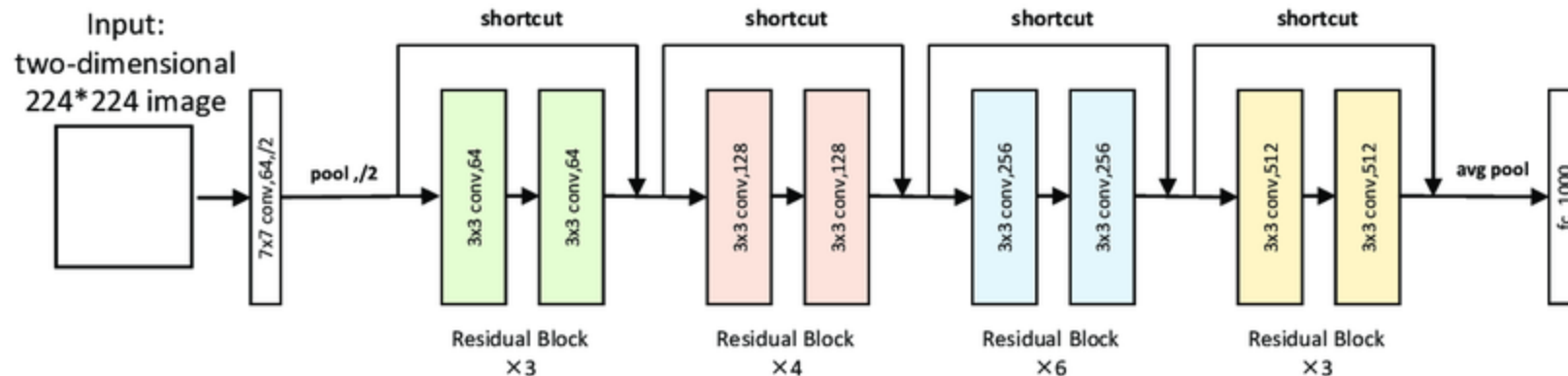


# Implementing your own Neural Network

Explanation of building blocks



## nn.Module

- The base class for all neural network components
- Allows you to organize layers and their parameters modularly
- Automatically tracks all parameters for optimization
- Must implement forward() method defining computation

## nn.Parameter

- Special tensor that gets automatically tracked
- Used for learnable weights and biases
- Will be updated during backpropagation
- Example: 

```
self.weight = nn.Parameter(torch.randn(3, 4))
```

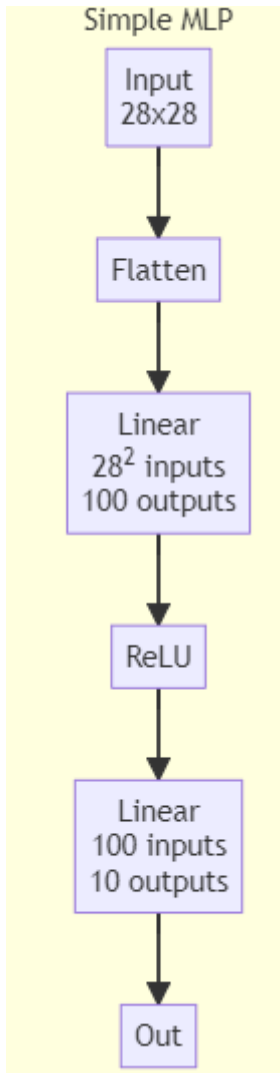
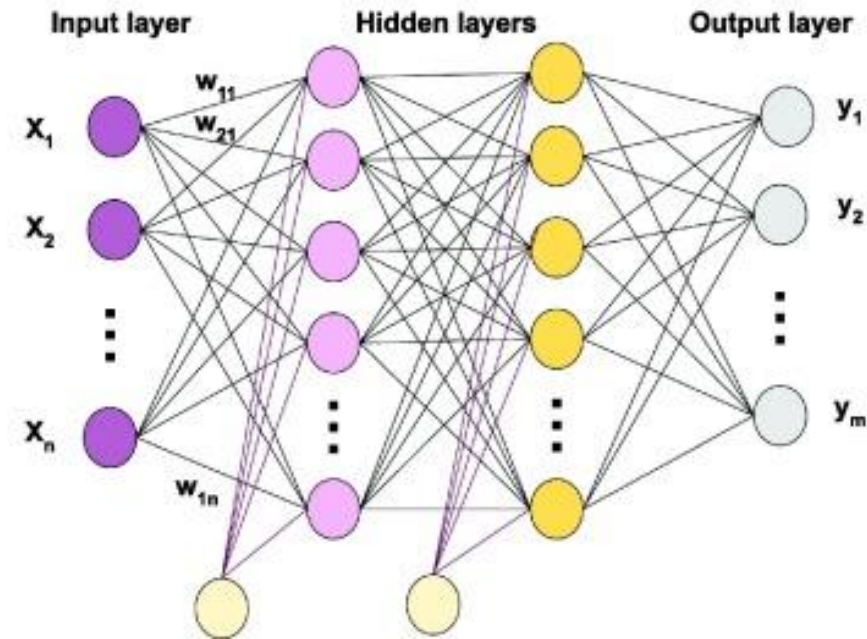
## nn.Sequential

- Container to stack modules in sequence
- Automatically chains forward() calls
- Useful for linear architectures
- Example:

```
nn.Sequential(  
    nn.Linear(10, 20),  
    nn.ReLU(),  
    nn.Linear(20, 30)  
)
```

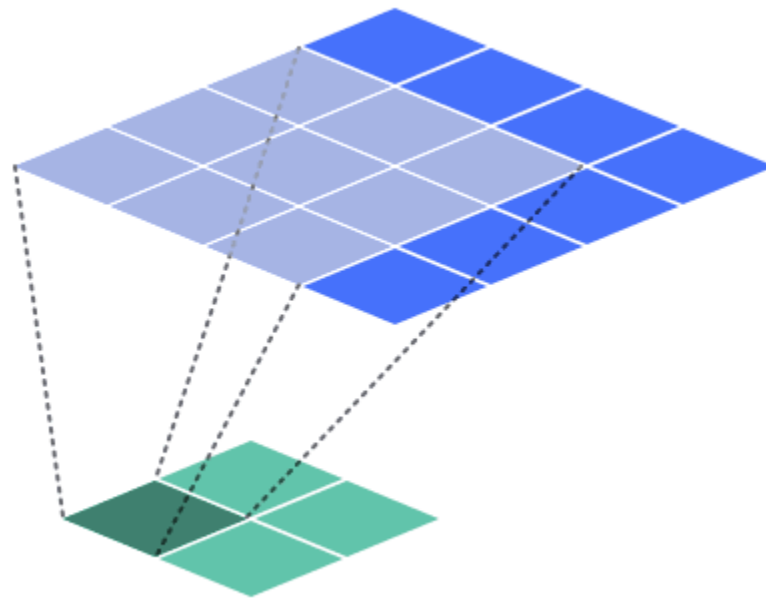
# Multilayer Perceptron

- The classic simplest neural network
- A stack of fully connected layers with activation functions
- No sense of space in the neural network structure
- Works with any kind of data
- Also, still part of Transformer architecture



# Convolution

- You can make use of the “local” structure of the image
- An “edge detector” (for example) acts the same on parts of the image



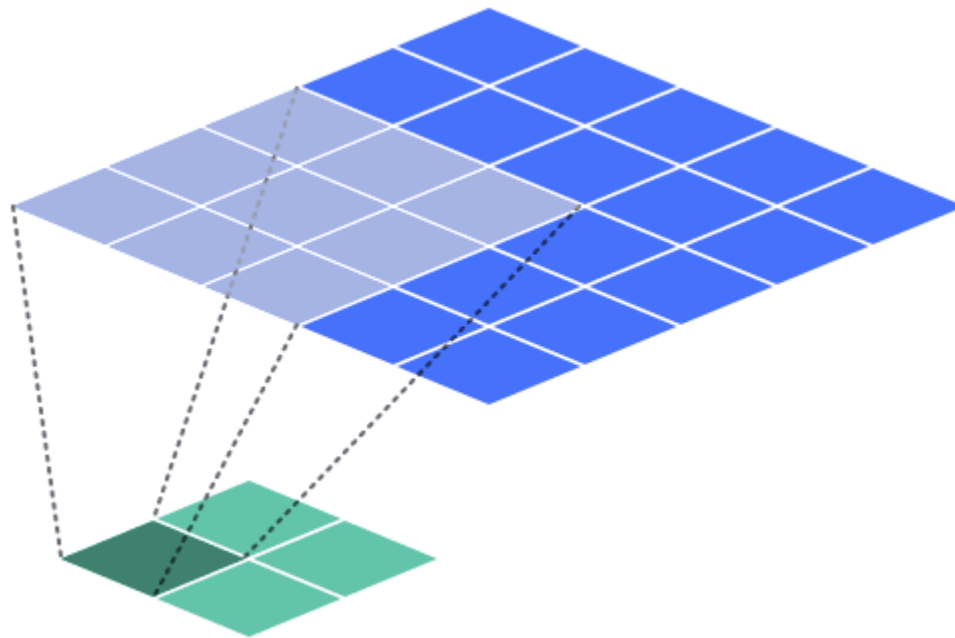
## Parameters

```
kernal: 3  
stride: 1  
padding: 0  
output_padding: 0  
dilation: 1
```

■ Kernel  
■ Output  
■ Input

# Convolution - stride

- You can give the image, as it passes through the network, lower resolution, by skipping steps
- Here, no pixel gets completely lost



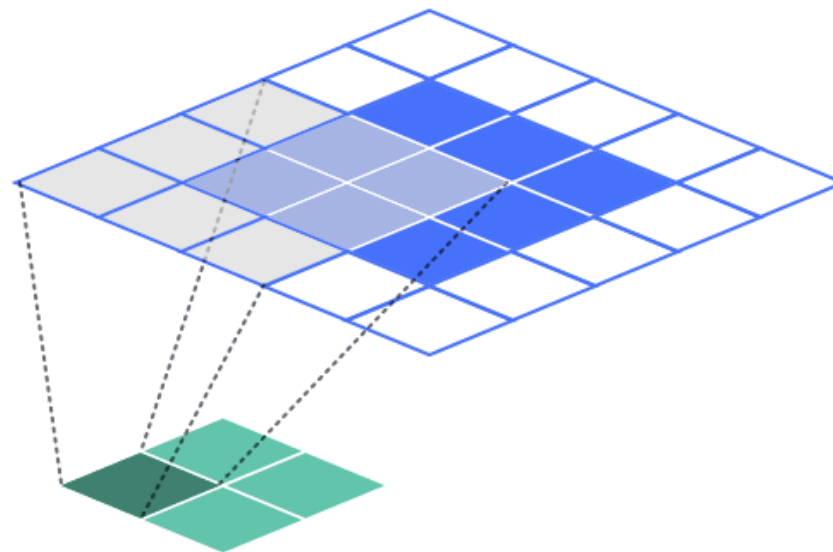
## Parameters

```
kernal: 3  
stride: 2  
padding: 0  
output_padding: 0  
dilation: 1
```

■ Kernel  
■ Output  
■ Input

# Convolution - padding

- Pixels on the border of the image get viewed only once, image gets smaller
- Image gets smaller without stride
- Add “empty” pixels outside the image to cancel this effect



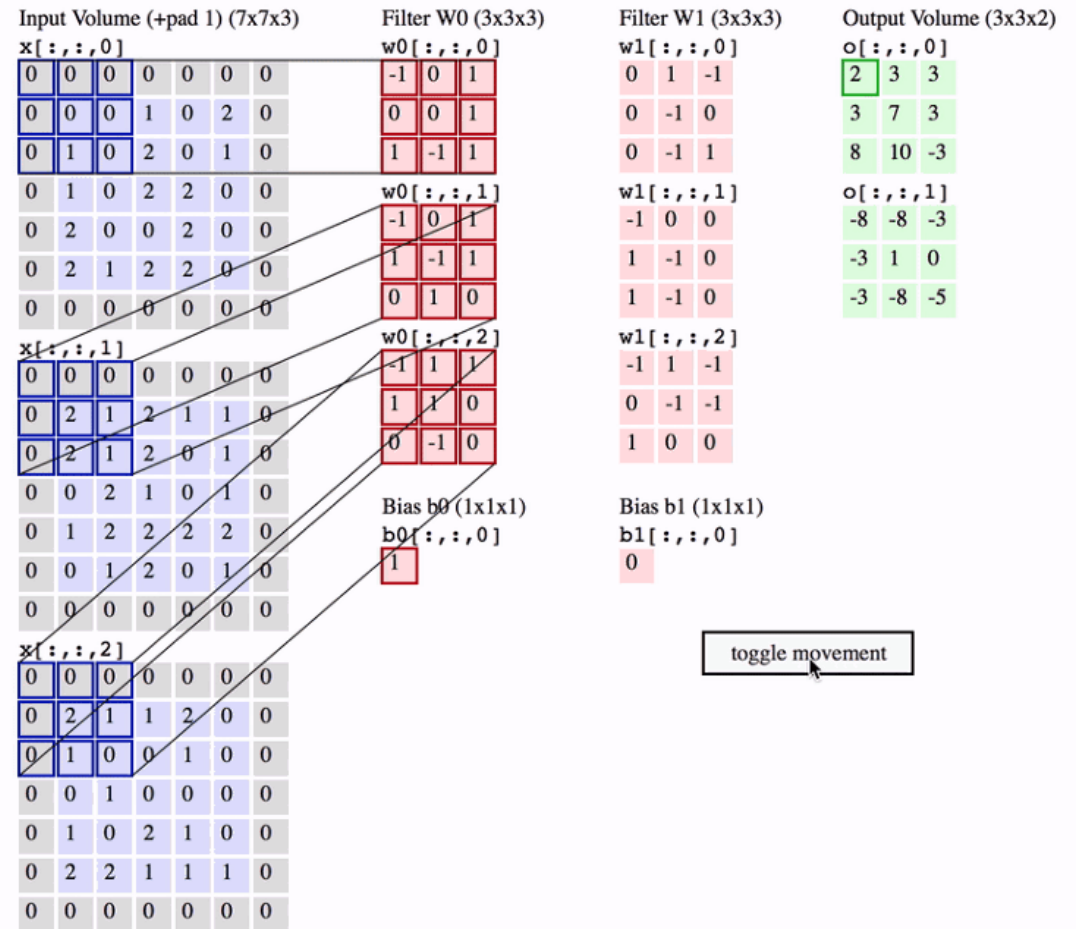
## Parameters

```
kernal: 3  
stride: 2  
padding: 1  
output_padding: 0  
dilation: 1
```

- Padding
- Kernel
- Output
- Input

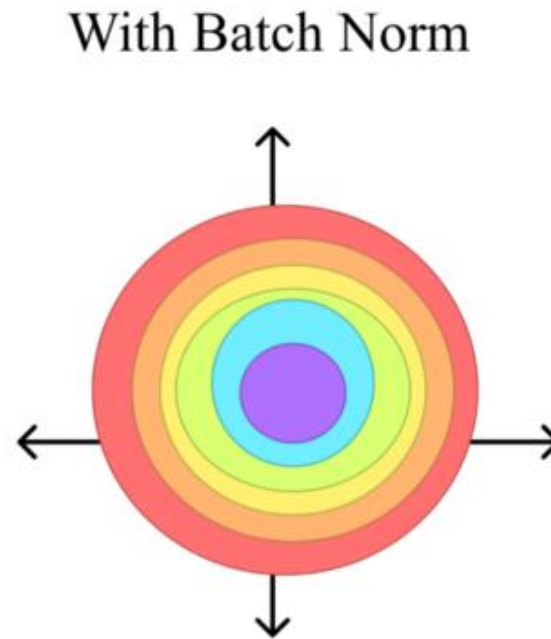
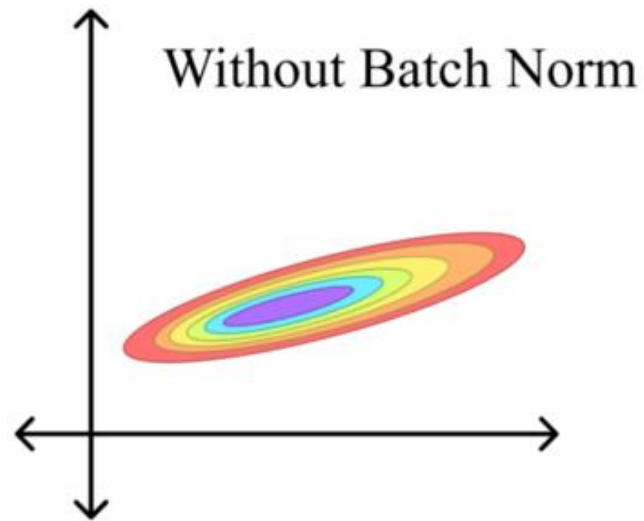
# Convolution - channels

- Typically, you have more than one convolution kernel active at once
- This leads to an additional output dimension per picture: channel
- Example: Colour or different features like different edges/objects



# Batch Norm

- The value of the activation is less important than its relative size
- Example: brightness in an image / across images
- Normalize distribution





# Batch Norm– how this gets calculated

Learnable parameters

The diagram shows the Batch Normalization formula: 
$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$
 Annotations include: Two blue arrows pointing to  $\gamma$  and  $\beta$  from the text "Learnable parameters". Two orange arrows pointing to  $E[x]$  and  $\text{Var}[x]$  from the text "In Training: Updated as a rolling average" and "In Deployment: Frozen". Two green arrows point from the bottom left towards the denominator of the formula.

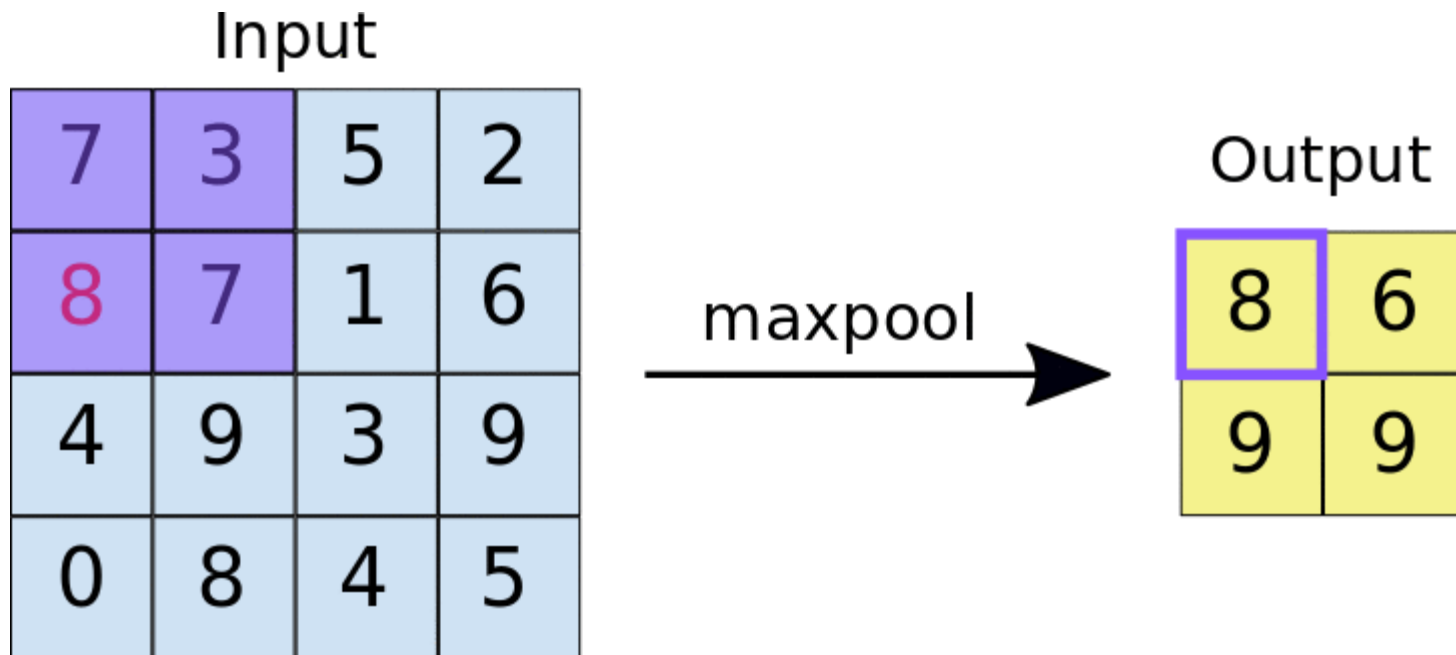
$$y = \frac{x - E[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

In Training: Updated as a rolling average  
In Deployment: Frozen

Activations: For taking activation statistics, all activations from the same channel are taken as draws from the same distribution (independent of location in picture)

# MaxPooling (2d)

- Reduce the size of the resulting Image
- Taking the maximum of each kernel window
- Example: Object detection in a part of the image



# AveragePool

- Reduce resulting image to a single number
- Example: classifying an image with a single label
- Average over space dimensions

# ResNet 34

