

DEEP LEARNING FOR COMPUTER VISION

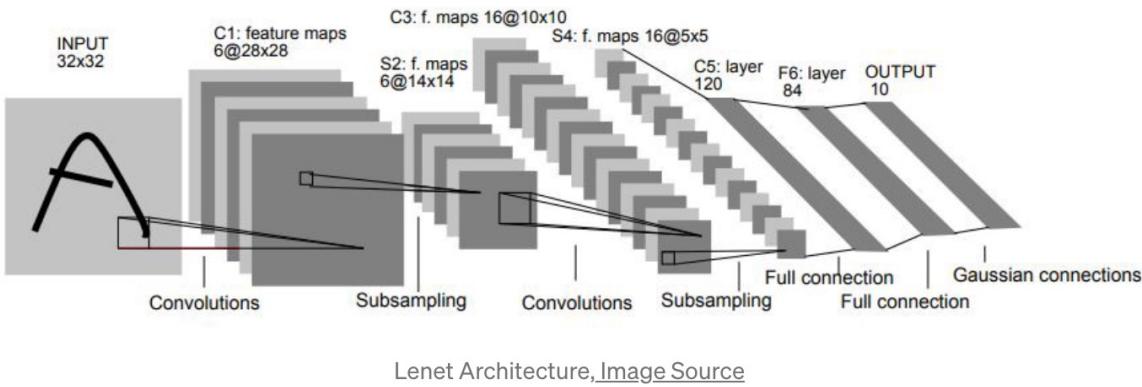
Week10



Dr. Tuchsanai. PloySuwan

Lenet:

Lenet 5 is considered as the first architecture for Convolutional Neural Networks, which are used to identify handwritten digits in the zip codes in the US. It was introduced in the paper, “[Gradient-Based Learning Applied To Document Recognition.](#)”



The LeNet-5 architecture consists of two sets of convolutional, activation, and pooling layers, followed by two fully-connected layers, tanh activation, and finally a softmax classifier which classifies the MNIST digits[1]

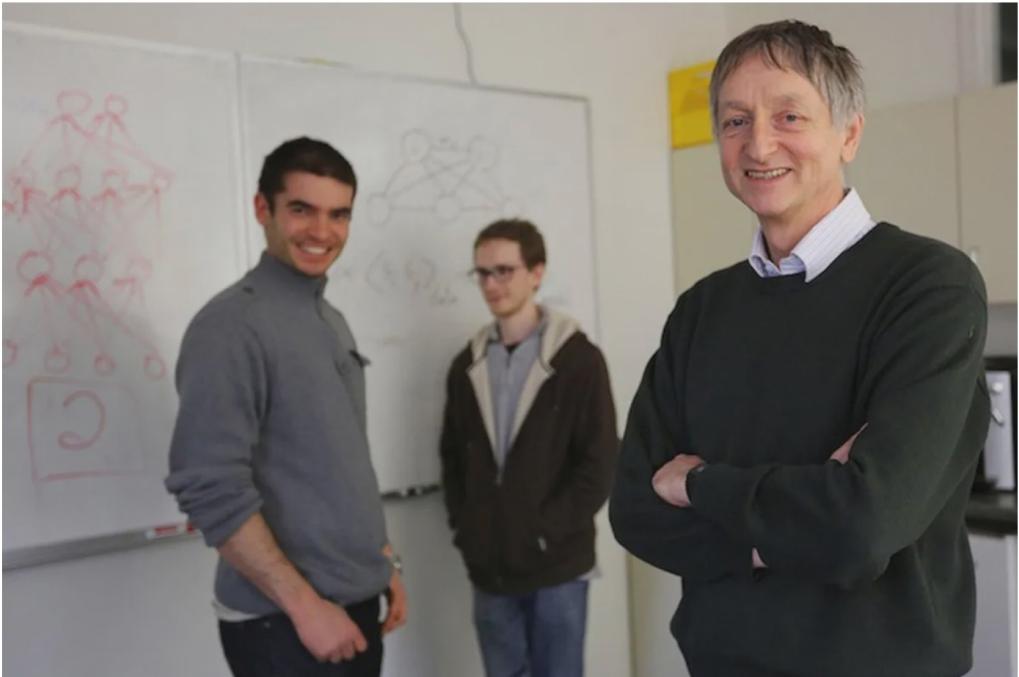
```
class Lenet(nn.Module):
    def __init__(self):
        super(Lenet, self).__init__()
        self.tanh = nn.Tanh()
        self.pool = nn.AvgPool2d(kernel_size=(2,2), stride=(2,2))
        self.conv1 = nn.Conv2d(in_channels=1, out_channels=6, kernel_size=(5,5), stride=(1,1))
        self.conv2 = nn.Conv2d(in_channels=6, out_channels=16, kernel_size=(5,5), stride=(1,1))
        self.conv3 = nn.Conv2d(in_channels=16, out_channels=120, kernel_size=(5,5), stride=(1,1))
        self.linear1 = nn.Linear(120, 84)
        self.linear2 = nn.Linear(84, 10)

    def forward(self, x):
        x = self.tanh(self.conv1(x))
        x = self.pool(x)
        x = self.tanh(self.conv2(x))
        x = self.pool(x)
        x = self.tanh(self.conv3(x))
        x = x.reshape(x.shape[0], -1)
        x = self.tanh(self.linear1(x))
        x = self.linear2(x)
        return x

model = Lenet()
```

Classical Deep Learning Research Paper — ALEXNET

One of the Most Influential Papers of Deep Convolutional Neural Network



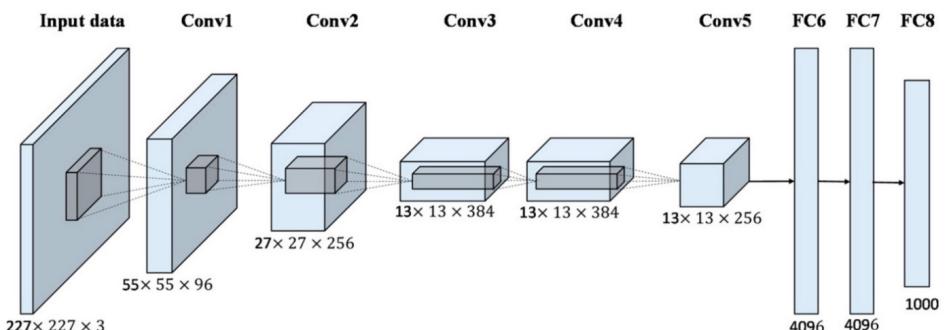
Geoffrey Hinton, Ilya Sutskever, Alex Krizhevsky Credit: Wired

“ImageNet Classification with Deep Convolutional Neural Networks” — was published in the year 2012 by **Alex Krizhevsky**, **Ilya Sutskever**, and the *Father of Backpropagation, Geoffrey Hinton*. You can download the paper from [here](#).

AlexNet, the winner of ImageNet 2012 and the model that apparently kick started the focus on deep learning had only 8 convolutional layers,

AlexNet Summary

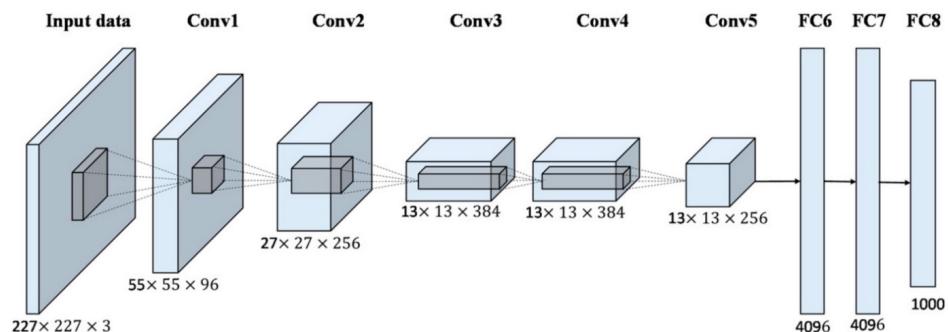
Layer	Feature Map	Size	Kernel Size	Stride	Activation
Input	Image	1	227x227x3	-	-
1	Convolution	96	55 x 55 x 96	11x11	4
	Max Pooling	96	27 x 27 x 96	3x3	2
2	Convolution	256	27 x 27 x 256	5x5	1
	Max Pooling	256	13 x 13 x 256	3x3	2
3	Convolution	384	13 x 13 x 384	3x3	1
4	Convolution	384	13 x 13 x 384	3x3	1
5	Convolution	256	13 x 13 x 256	3x3	1
	Max Pooling	256	6 x 6 x 256	3x3	2
6	FC	-	9216	-	relu
7	FC	-	4096	-	relu
8	FC	-	4096	-	relu
Output	FC	-	1000	-	Softmax



This article describes about highlights of ALEXNET architecture.

AlexNet Summary

Layer		Feature Map	Size	Kernel Size	Stride	Activation
Input	Image	1	227x227x3	-	-	-
1	Convolution	96	55 x 55 x 96	11x11	4	relu
	Max Pooling	96	27 x 27 x 96	3x3	2	relu
2	Convolution	256	27 x 27 x 256	5x5	1	relu
	Max Pooling	256	13 x 13 x 256	3x3	2	relu
3	Convolution	384	13 x 13 x 384	3x3	1	relu
4	Convolution	384	13 x 13 x 384	3x3	1	relu
5	Convolution	256	13 x 13 x 256	3x3	1	relu
	Max Pooling	256	6 x 6 x 256	3x3	2	relu
6	FC	-	9216	-	-	relu
7	FC	-	4096	-	-	relu
8	FC	-	4096	-	-	relu
Output	FC	-	1000	-	-	Softmax



This article describes about highlights of ALEXNET architecture.

```

import torch
import torch.nn as nn

class AlexNet(nn.Module):
    def __init__(self, in_channels=3, classes=1000):
        super().__init__()
        self.c1 = nn.Conv2d(in_channels=3, out_channels=96, kernel_size=11, stride=4, padding=0)
        self.c2 = nn.Conv2d(in_channels=96, out_channels=256, kernel_size=5, stride=1, padding=2)
        self.c3 = nn.Conv2d(in_channels=256, out_channels=384, kernel_size=3, stride=1, padding=1)
        self.c4 = nn.Conv2d(in_channels=384, out_channels=384, kernel_size=3, stride=1, padding=1)
        self.c5 = nn.Conv2d(in_channels=384, out_channels=256, kernel_size=3, stride=1, padding=1)

        self.fc1 = nn.Linear(6*6*256, 4096)
        self.fc2 = nn.Linear(4096, 4096)
        self.fc3 = nn.Linear(4096, classes)

        self.localnorm = nn.LocalResponseNorm(size=5)
        self.maxpool = nn.MaxPool2d(kernel_size=3, stride=2)
        self.relu = nn.ReLU(inplace=True)
        self.dropout = nn.Dropout(0.5)

    def forward(self, x):
        # x shape: [batch, 3, 227, 227]
        x = self.relu(self.c1(x))
        # x shape: [batch, 96, 55, 55]
        x = self.maxpool(self.localnorm(x))
        # x shape: [batch, 96, 27, 27]
        x = self.relu(self.c2(x))
        # x shape: [batch, 256, 27, 27]
        x = self.maxpool(self.localnorm(x))
        # x shape: [batch, 256, 13, 13]
        x = self.relu(self.c3(x))
        # x shape: [batch, 384, 13, 13]
        x = self.relu(self.c4(x))
        # x shape: [batch, 384, 13, 13]
        x = self.maxpool(self.relu(self.c5(x)))
        # x shape: [batch, 256, 6, 6]
        x = torch.flatten(x,1)
        # x shape: [batch, 256*6*6]
        x = self.relu(self.dropout(self.fc1(x)))
        # x shape: [batch, 4096]
        x = self.relu(self.dropout(self.fc2(x)))
        # x shape: [batch, 4096]
        x = self.fc3(x)
        # x shape: [batch, classes]
        return x

model = AlexNet()

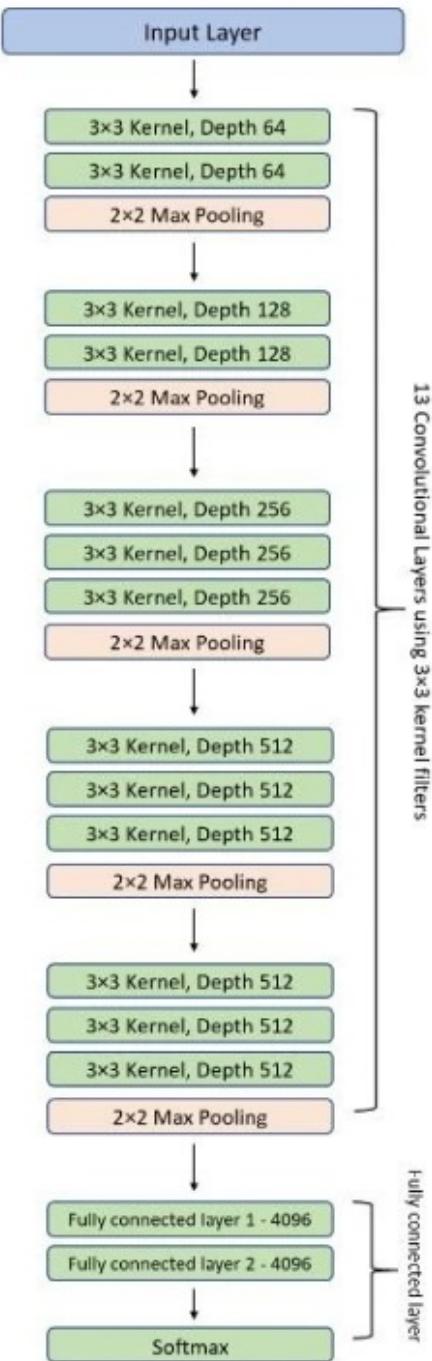
```

VGGNet Architecture

Comparison					
Network	Year	Salient Feature	top5 accuracy	Parameters	FLOP
AlexNet	2012	Deeper	84.70%	62M	1.5B
VGGNet	2014	Fixed-size kernels	92.30%	138M	19.6B
Inception	2014	Wider - Parallel kernels	93.30%	6.4M	2B
ResNet-152	2015	Shortcut connections	95.51%	60.3M	11B

VGGNet Architecture

The standard
VGG-16
network

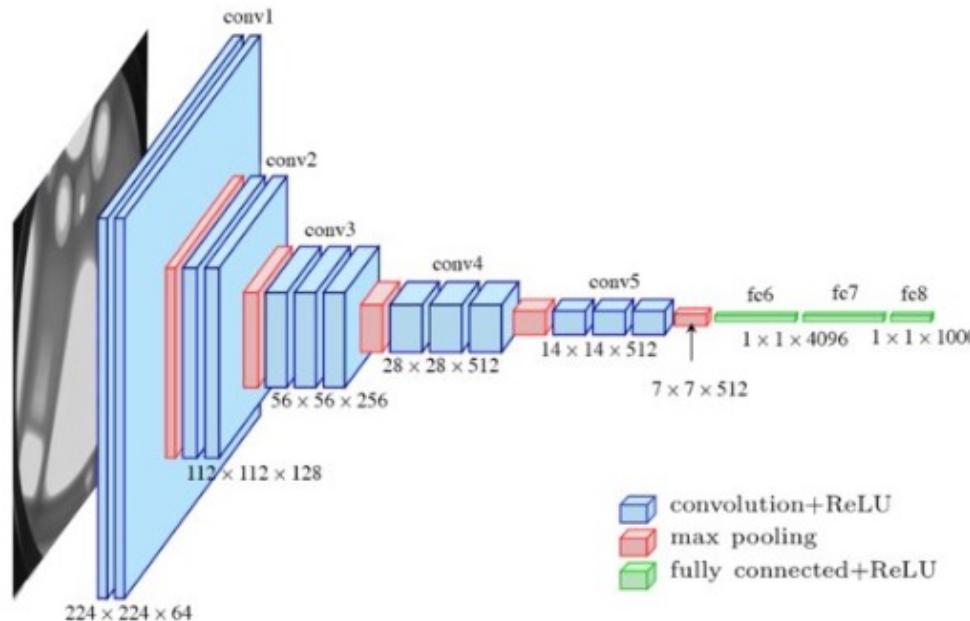


ConvNet Configuration									
A	A-LRN	B	C	D	E				
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers				
input (224 × 224 RGB image)									
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64	conv3-64	conv3-64 conv3-64				
maxpool									
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128 conv3-128				
conv3-128									
maxpool									
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256 conv3-256				
conv3-256	conv3-256		conv3-256		conv3-256 conv3-256				
conv1-256									
maxpool									
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512 conv3-512				
conv3-512	conv3-512		conv3-512		conv3-512 conv3-512				
conv1-512									
maxpool									
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512 conv3-512				
conv3-512	conv3-512		conv3-512		conv3-512 conv3-512				
conv1-512									
maxpool									
FC-4096	FC-4096		FC-4096		FC-1000				
soft-max									

VGGNet Architecture

VGGNet is a Convolutional Neural Network architecture proposed by Karen Simonyan and Andrew Zisserman from the University of Oxford in 2014. This paper mainly focuses on the effect of the convolutional neural network depth on its accuracy. You can find the original paper of VGGNet which is titled Very Deep Convolutional Networks for Large Scale Image Recognition.

<https://arxiv.org/pdf/1409.1556.pdf>



The standard VGG-16 network

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

```

configs = {
    "VGG11" : [1,1,2,2,2], # 1x64, 1x128, 2x256, 2x512, 2x512
    "VGG13" : [2,2,2,2,2], # 2x64, 2x128, 2x256, 2x512, 2x512
    "VGG16" : [2,2,3,3,3], # 2x64, 2x128, 3x256, 3x512, 3x512
    "VGG19" : [2,2,4,4,4] # 2x64, 2x128, 4x256, 4x512, 4x512
}

```

```

import torch
import torch.nn as nn

class Block(nn.Module):
    def __init__(self, in_channels, out_channels, no_layers):
        super().__init__()
        self.layers = [nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=3, stride=1, padding=1), nn.ReLU()]
        for i in range(no_layers-1):
            self.layers.append(nn.Conv2d(in_channels=out_channels, out_channels=out_channels, kernel_size=3, stride=1, padding=1))
            self.layers.append(nn.ReLU())

        self.layers.append(nn.MaxPool2d(kernel_size=2, stride=2))
        self.seq = nn.Sequential(*self.layers)

    def forward(self, x):
        return self.seq(x)

```

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

```

class VGG(nn.Module):
    def __init__(self,
                 model_size,
                 in_channels : int = 3,
                 classes : int = 1000):
        super().__init__()

        config = configs[model_size]
        channels = [in_channels, 64, 128, 256, 512, 512]
        self.blocks = nn.ModuleList([])
        for i in range(len(config)):
            self.blocks.append(Block(channels[i], channels[i+1], config[i]))

        self.fc1 = nn.Linear(7*7*512, 4096)
        self.fc2 = nn.Linear(4096, 4096)
        self.fc3 = nn.Linear(4096, classes)
        self.relu = nn.ReLU()
        self.dropout = nn.Dropout(0.5)

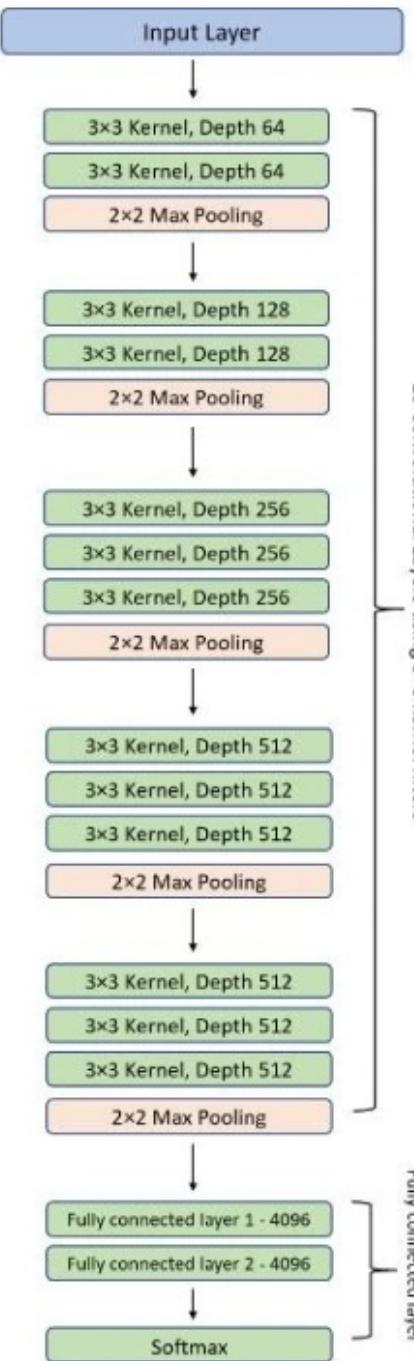
    def forward(self, x):
        for block in self.blocks:
            x = block(x)

        x = torch.flatten(x,1)
        x = self.relu(self.dropout(self.fc1(x)))
        x = self.relu(self.dropout(self.fc2(x)))
        x = self.fc3(x)

    return x

```

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256	conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

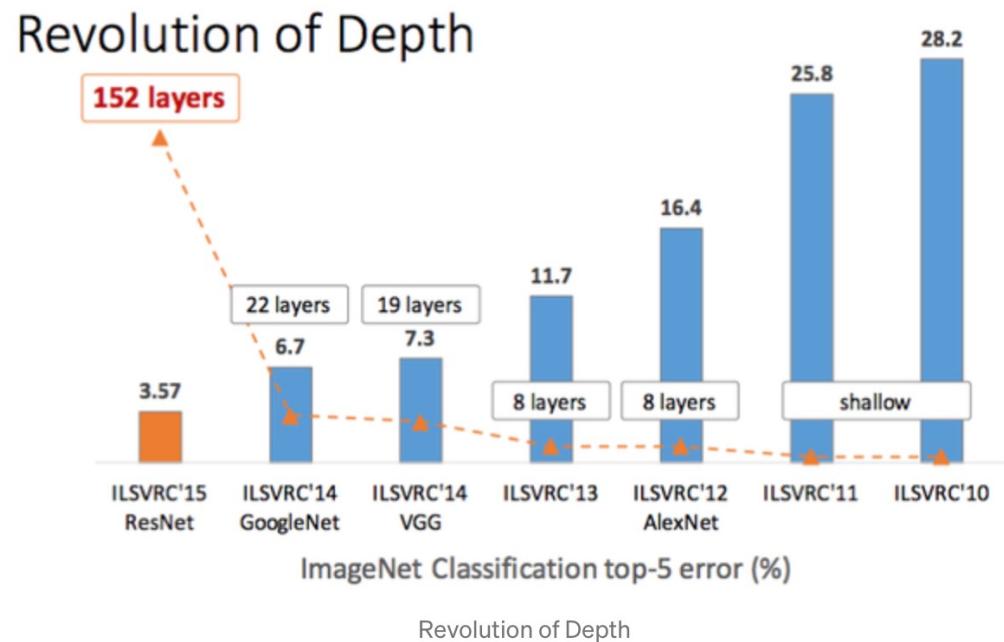


```

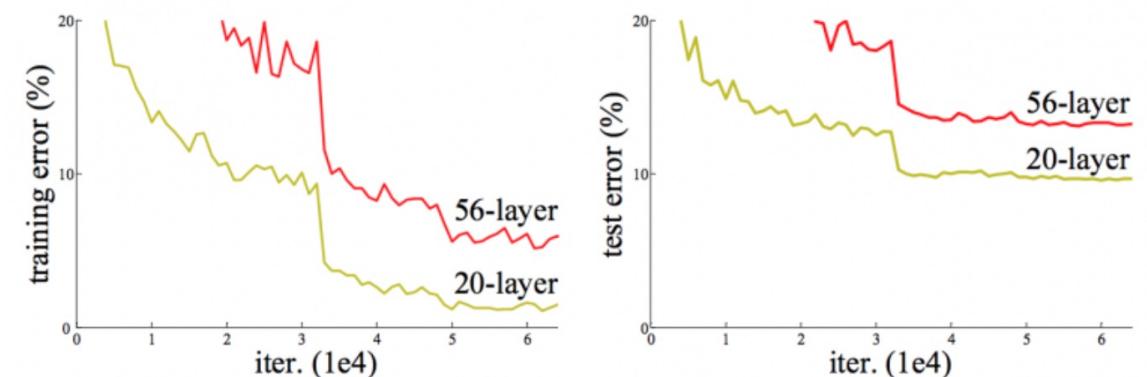
vgg = VGG("VGG16")
vgg
✓ 0.4s

Output exceeds the size limit. Open the full output data in a text editor
VGG(
  (blocks): ModuleList(
    (0): Block(
      (seq): Sequential(
        (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        (1): ReLU()
        (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        (3): ReLU()
        (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
      )
    )
    (1): Block(
      (seq): Sequential(
        (0): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        (1): ReLU()
        (2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        (3): ReLU()
        (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
      )
    )
    (2): Block(
      (seq): Sequential(
        (0): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
        (1): ReLU()
        (2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
      )
    )
    ...
    (fc2): Linear(in_features=4096, out_features=4096, bias=True)
    (fc3): Linear(in_features=4096, out_features=1000, bias=True)
    (relu): ReLU()
    (dropout): Dropout(p=0.5, inplace=False)
  )
)
  
```

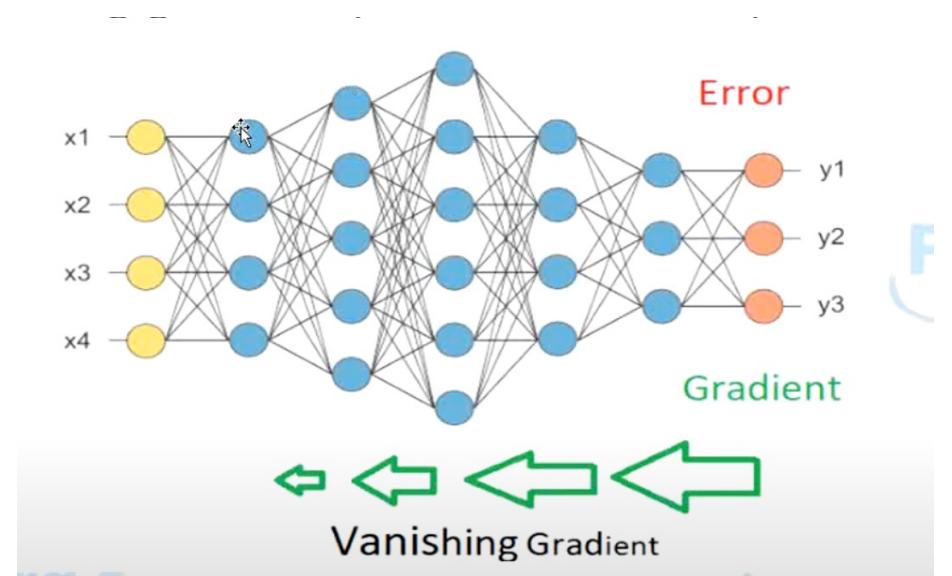
ResNet:

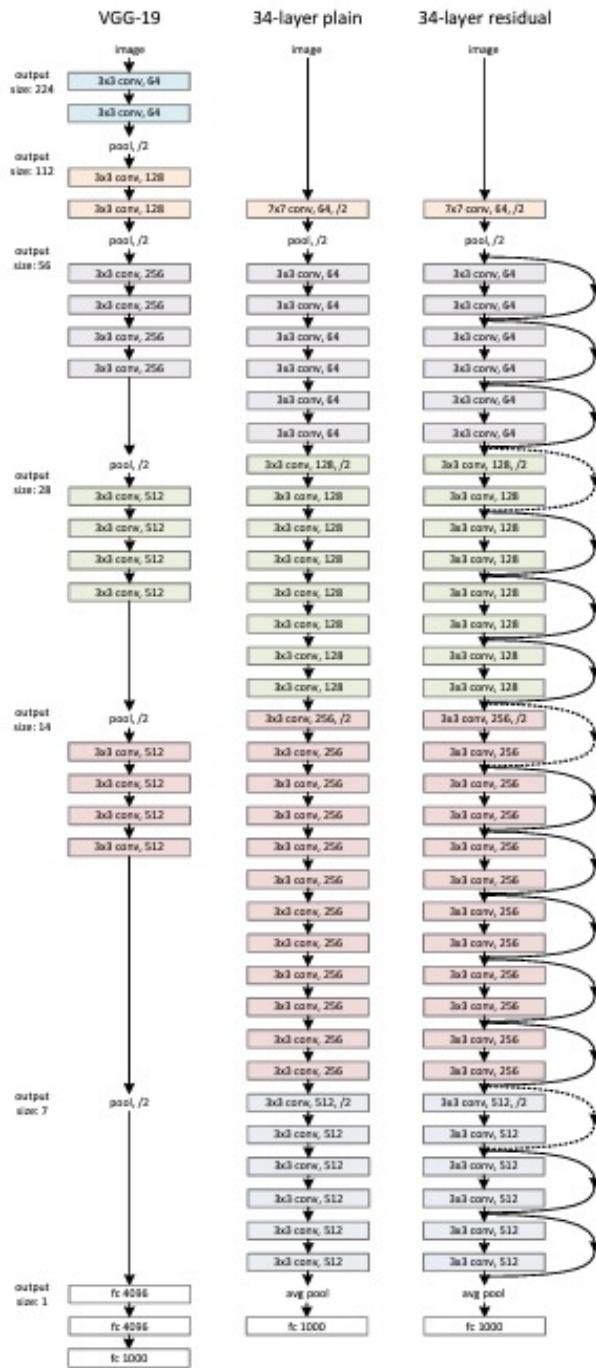


However, increasing network depth does not work by simply stacking layers together. Deep networks are hard to train because of the notorious vanishing gradient problem — as the gradient is back-propagated to earlier layers, repeated multiplication may make the gradient extremely small. As a result, as the network goes deeper, its performance gets saturated or even starts degrading rapidly.



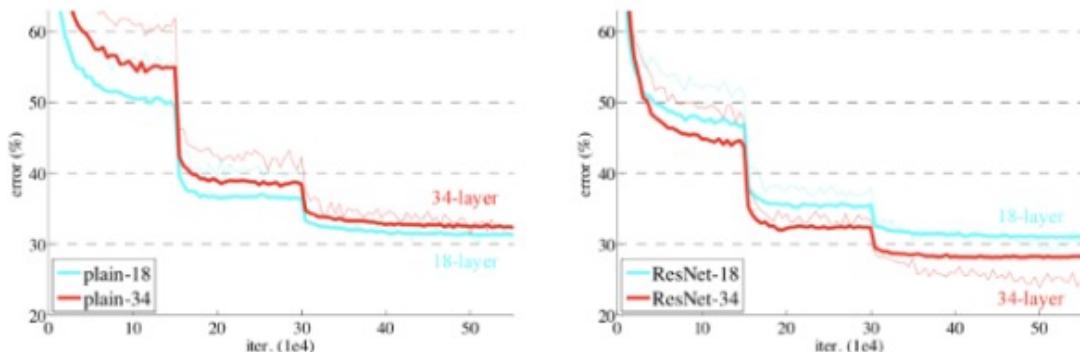
It is observed that the network having a higher count (56-layer) of layers are resulting in higher training error in contrast to the network having a much lower count (20-layer) of layers thus resulting in higher test errors! Image Credits to the authors of original ResNet paper([Source](#))



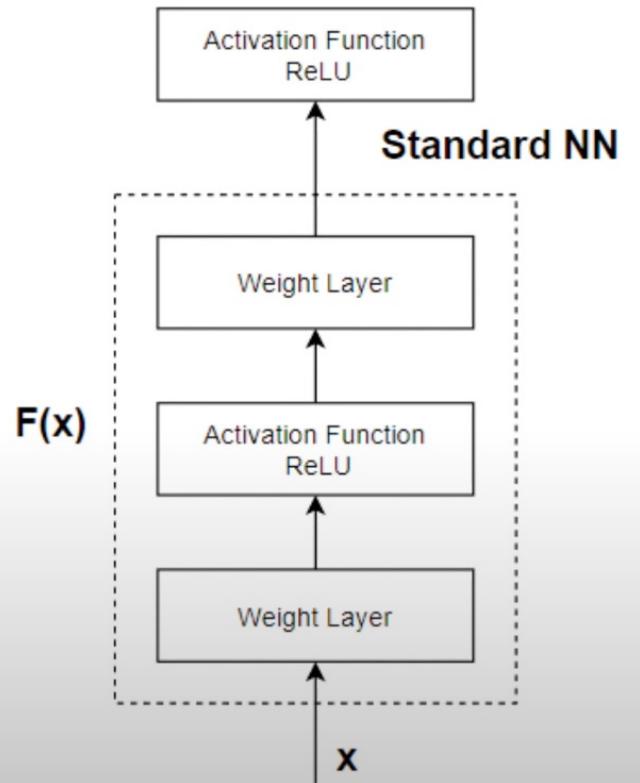
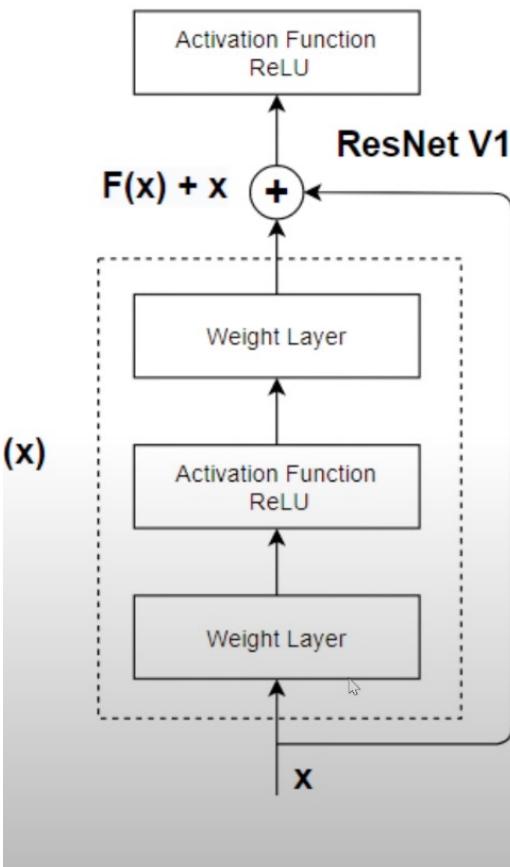
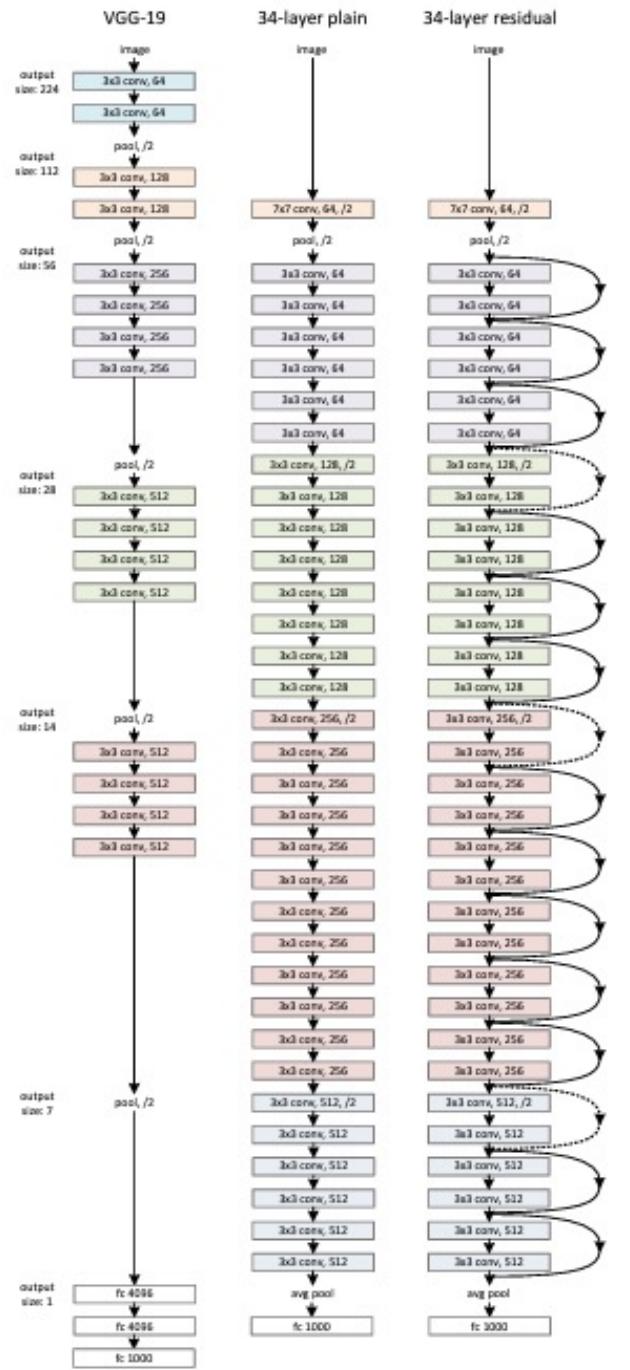


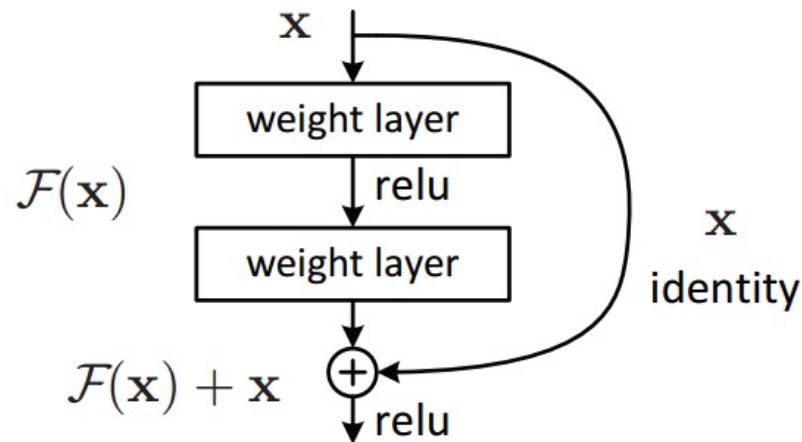
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\left[\begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$	$\left[\begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$
conv3_x	28×28	$\left[\begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$	$\left[\begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 4$	$\left[\begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[\begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 8$	
conv4_x	14×14	$\left[\begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$	$\left[\begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$	$\left[\begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$	$\left[\begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 23$	$\left[\begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 36$
conv5_x	7×7	$\left[\begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[\begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

	plain	ResNet
18 layers	27.94	27.88
34 layers	28.54	25.03



As seen ResNet performs better than plain neural network models, Image Credits to the authors of original ResNet paper([Source](#))





Motivation

Since neural networks are good function approximators, they should be able to easily solve the identity function, where the output of a function becomes the input itself.

$$f(x) = x$$

Following the same logic, if we bypass the input to the first layer of the model to be the output of the last layer of the model, the network should be able to predict whatever function it was learning before with the input added to it.

Skip Connect or Shortcut Connection, Image Credits to the authors of original ResNet paper([Source](#))

$$f(x) + x = h(x)$$

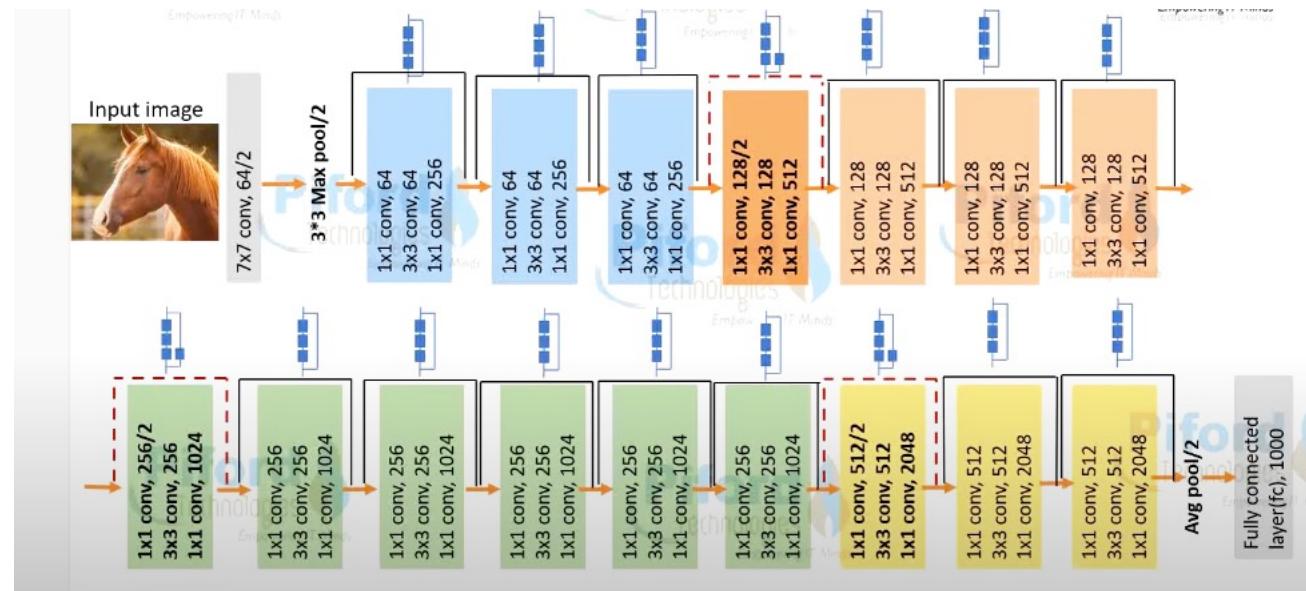
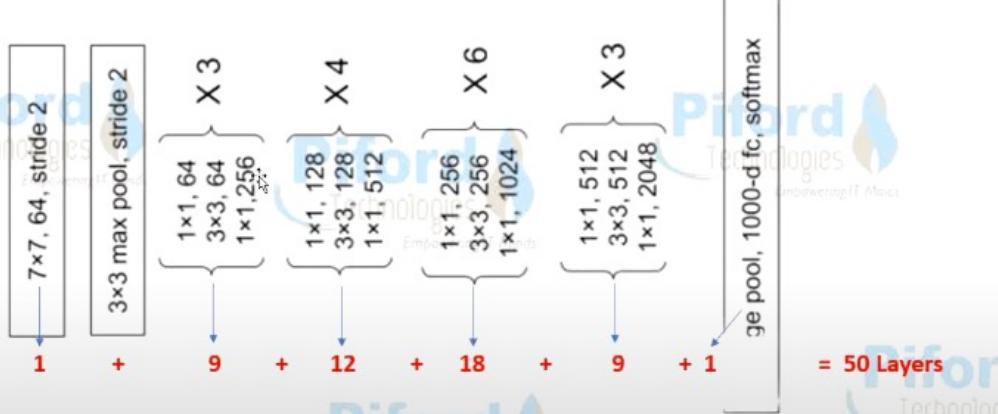
The intuition is that learning $f(x) = 0$ has to be easy for the network.

<https://medium.com/swlh/resnet-a-simple-understanding-of-the-residual-networks-bfd8a1b4a447>

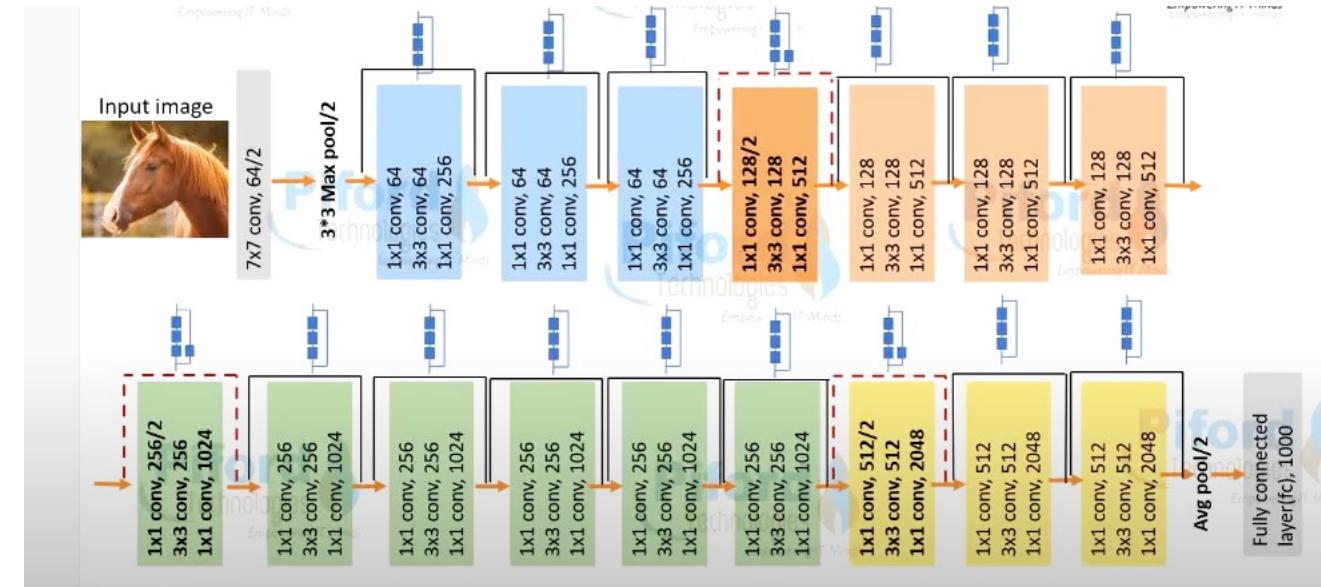
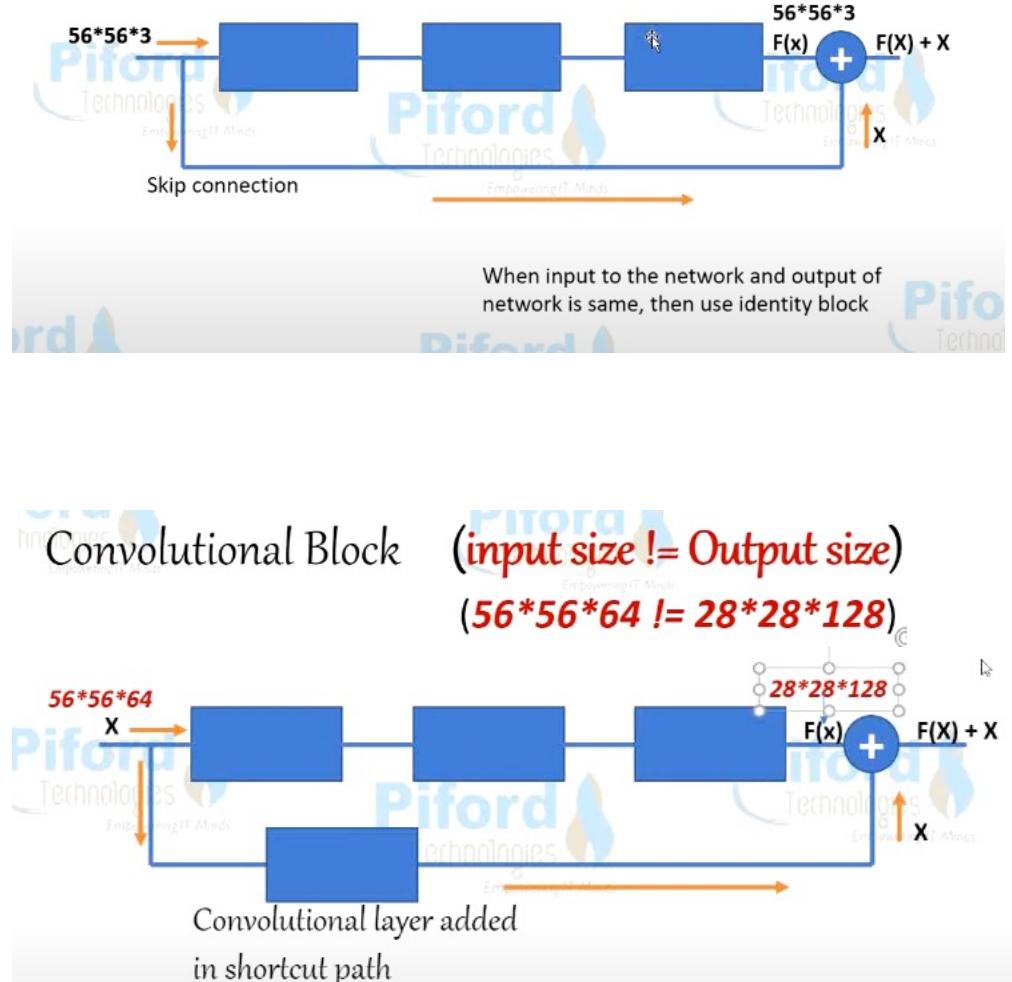
Example of ResNet-50

Pad = 3 →

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
				3×3 max pool, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9



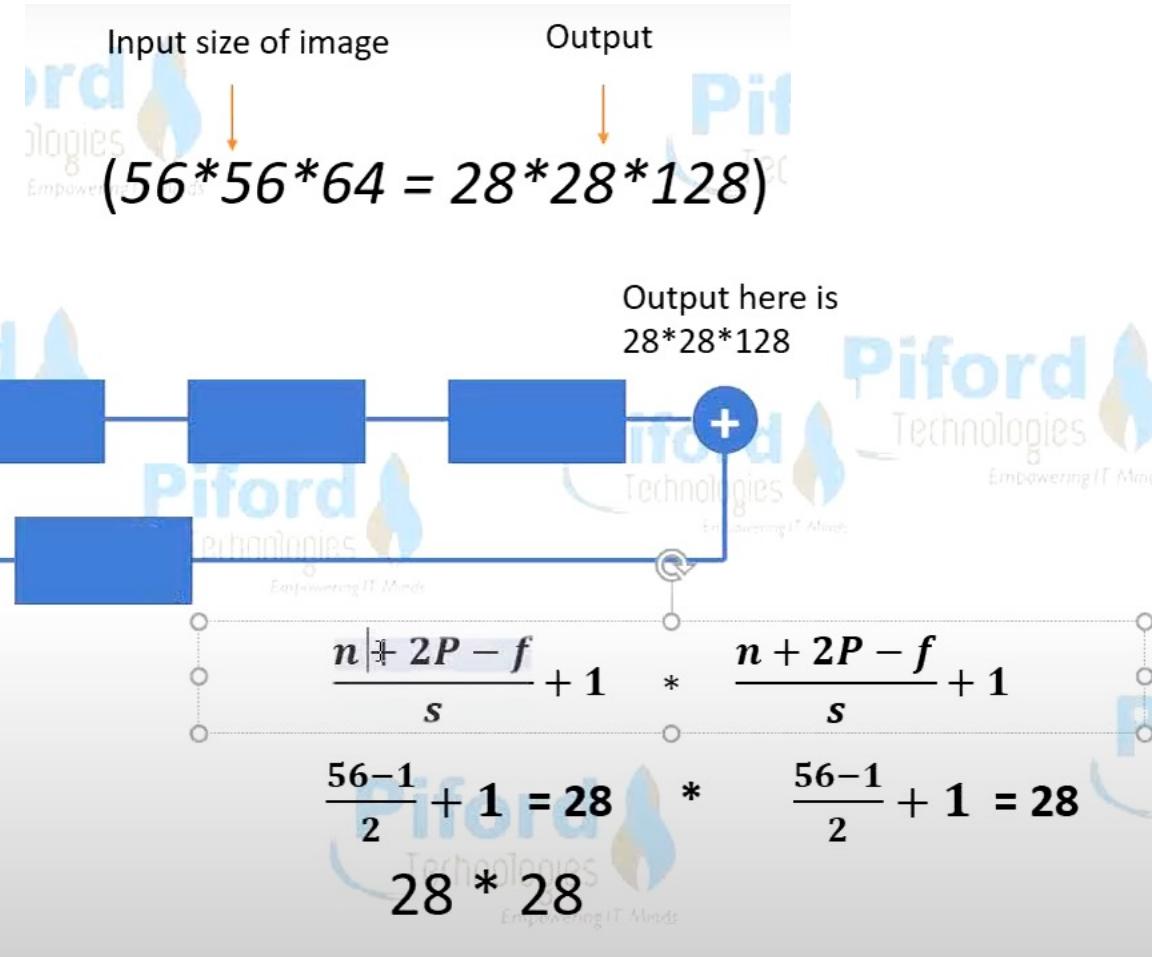
Identity Block (input size == output size)



layer name		output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1		112×112			$7 \times 7, 64, \text{stride } 2$		
			$3 \times 3 \text{ max pool, stride } 2$				
conv2_x	56×56		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28		$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14		$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7		$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
		1×1			average pool, 1000-d fc, softmax		
FLOPs			1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

Convolutional Block (**input size != Output size**)

(**$56*56*64 \neq 28*28*128$**)



The **1x1 convolution approach** is shown

300*300*3



$$\frac{n+2P-f}{s} + 1 \quad * \quad \frac{n+2P-f}{s} + 1$$

$$\frac{300 + 2 * 3 - 7}{2} + 1$$

150*150*64 is the output size

1

7×7, 64, stride 2

3×3 max pool, stride 2

$$\left\{ \begin{array}{l} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{array} \right\} \times 3$$

$$\left\{ \begin{array}{l} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{array} \right\} \times 4$$

$$\left\{ \begin{array}{l} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{array} \right\} \times 6$$

$$\left\{ \begin{array}{l} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{array} \right\} \times 3$$

average pool, 1000-d fc, softmax

Input Image : 150*150*64

$$\frac{150 + 2 * 1 - 3}{2} + 1 * \frac{150 + 2 * 1 - 3}{2} + 1 = 75 * 75$$

75*75*64 is the output of pooling layer

7x7, 64, stride 2

2 3x3 max pool, stride 2

$\begin{Bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 128 \end{Bmatrix} \times 3$

$\begin{Bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{Bmatrix} \times 4$

$\begin{Bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{Bmatrix} \times 6$

$\begin{Bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{Bmatrix} \times 3$

average pool, 1000-d fc, softmax

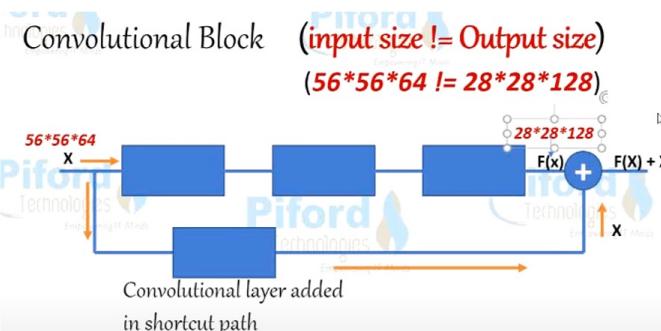
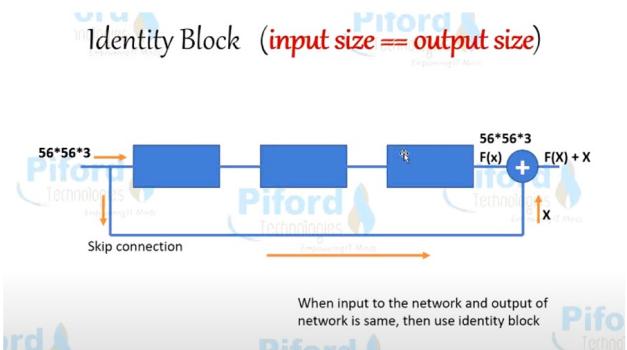
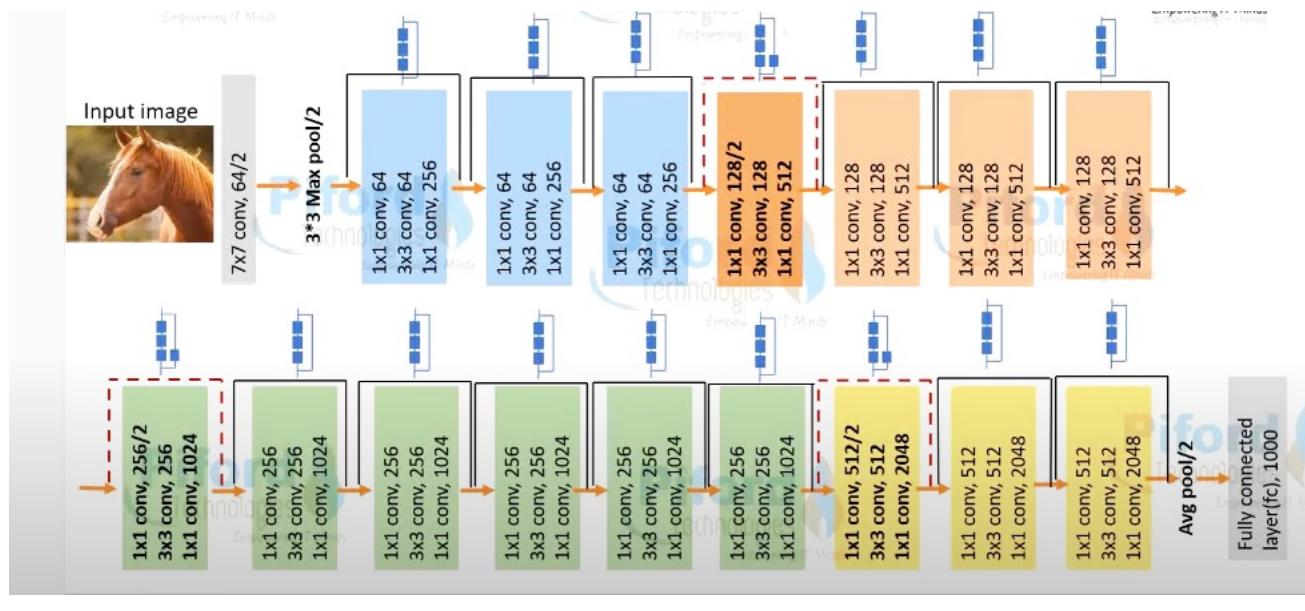
```

import torch
import torch.nn as nn

# ConvBlock
class ConvBlock(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size, stride, padding):
        super().__init__()
        self.c = nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=kernel_size, stride=stride, padding=padding)
        self.bn = nn.BatchNorm2d(out_channels)

    def forward(self, x):
        return self.bn(self.c(x))

```



```

# Bottleneck ResidualBlock
class ResidualBlock(nn.Module):
    def __init__(self, in_channels, out_channels, first=False):
        super().__init__()
        res_channels = in_channels // 4
        stride = 1

        self.projection = in_channels!=out_channels
        if self.projection:
            self.p = ConvBlock(in_channels, out_channels, 1, 2, 0)
            stride = 2
            res_channels = in_channels // 2

        if first:
            self.p = ConvBlock(in_channels, out_channels, 1, 1, 0)
            stride = 1
            res_channels = in_channels

        self.c1 = ConvBlock(in_channels, res_channels, 1, 1, 0)
        self.c2 = ConvBlock(res_channels, res_channels, 3, stride, 1)
        self.c3 = ConvBlock(res_channels, out_channels, 1, 1, 0)
        self.relu = nn.ReLU()

    def forward(self, x):
        f = self.relu(self.c1(x))
        f = self.relu(self.c2(f))
        f = self.c3(f)

        if self.projection:
            x = self.p(x)

        h = self.relu(torch.add(f, x))
        return h

```

```

# ResNetx
class ResNet(nn.Module):
    def __init__(self,
                 config_name : int,
                 in_channels=3,
                 classes=1000
                 ):
        super().__init__()

        configurations = {
            50 : [3, 4, 6, 3],
            101 : [3, 4, 23, 3],
            152 : [3, 8, 36, 3]
        }

        no_blocks = configurations[config_name]

        out_features = [256, 512, 1024, 2048]
        self.blocks = nn.ModuleList([ResidualBlock(64, 256, True)])

        for i in range(len(out_features)):
            if i > 0:
                self.blocks.append(ResidualBlock(out_features[i-1], out_features[i]))
            for _ in range(no_blocks[i]-1):
                self.blocks.append(ResidualBlock(out_features[i], out_features[i]))

        self.conv1 = ConvBlock(in_channels, 64, 7, 2, 3)
        self.maxpool = nn.MaxPool2d(kernel_size=3, stride=2, padding=1)
        self.avgpool = nn.AdaptiveAvgPool2d((1,1))
        self.fc = nn.Linear(2048, classes)

        self.relu = nn.ReLU()

        self.init_weight()

    def forward(self, x):
        x = self.relu(self.conv1(x))
        x = self.maxpool(x)
        for block in self.blocks:
            x = block(x)

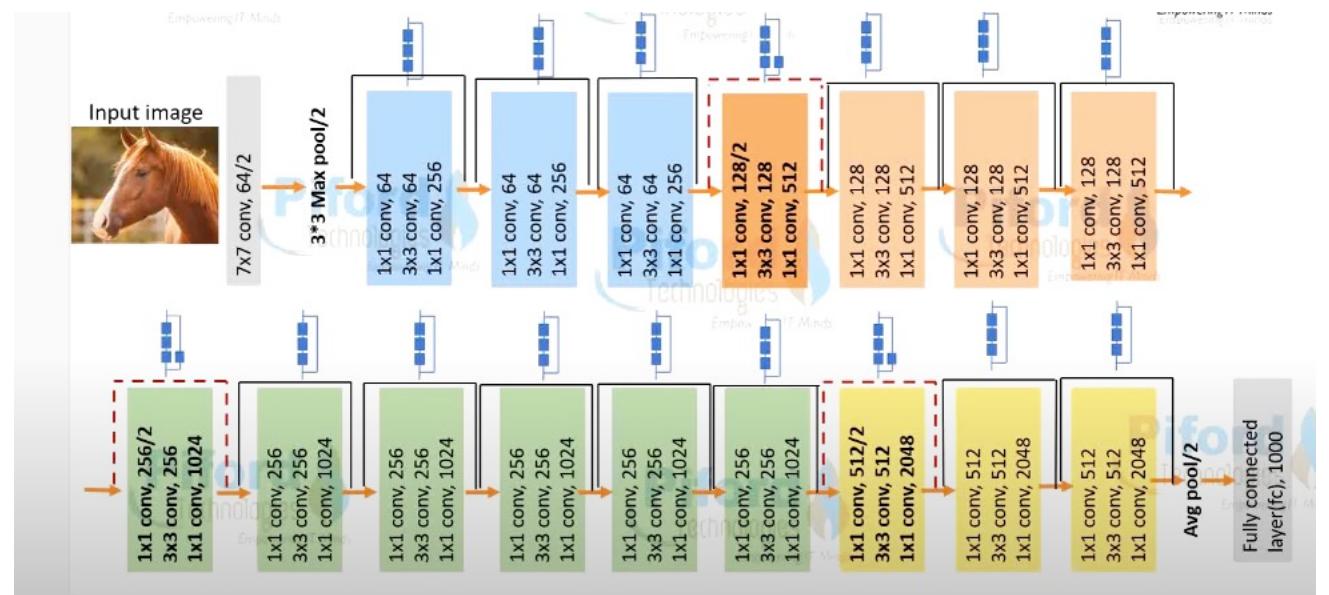
        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.fc(x)
        return x

```

```

config_name = 50 # 50-layer
resnet50 = ResNet(50)
image = torch.rand(1, 3, 224, 224)
print(resnet50(image).shape)

```



Pad=3 →

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

} Stride 1 } Stride 2 } -11- } -11-

The screenshot shows the PyTorch Hub for Researchers interface. At the top, there's a navigation bar with links for PyTorch, Get Started, Ecosystem, Mobile, Blog, Tutorials, Docs, Resources, GitHub, and a search icon. Below the header, the title "PYTORCH HUB FOR RESEARCHERS" is displayed in large red and black text. A subtitle "Explore and extend models from the latest cutting edge research." follows. Underneath, there are category filters: All, Audio, Generative, Nlp, Scriptable, and Vision. To the right of these filters are sorting options (Sort) and a search icon. The main content area displays the code for the VGG model structure.

```
VGG(  
    (features): Sequential(  
        (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (1): ReLU(inplace=True)  
        (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (3): ReLU(inplace=True)  
        (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (5): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (6): ReLU(inplace=True)  
        (7): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (8): ReLU(inplace=True)  
        (9): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (10): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (11): ReLU(inplace=True)  
        (12): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (13): ReLU(inplace=True)  
        (14): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (15): ReLU(inplace=True)  
        (16): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (17): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (18): ReLU(inplace=True)  
        (19): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (20): ReLU(inplace=True)  
        (21): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (22): ReLU(inplace=True)  
    ...  
        (4): ReLU(inplace=True)  
        (5): Dropout(p=0.5, inplace=False)  
        (6): Linear(in_features=4096, out_features=1000, bias=True)  
    )
```

```
import torch  
#model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg11', pretrained=True)  
# or any of these variants  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg11_bn', pretrained=True)  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg13', pretrained=True)  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg13_bn', pretrained=True)  
model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg16', pretrained=True)  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg16_bn', pretrained=True)  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg19', pretrained=True)  
# model = torch.hub.load('pytorch/vision:v0.10.0', 'vgg19_bn', pretrained=True)  
model.eval()
```

2. Pre-trained Models from torch vision <https://pytorch.org/vision/main/models.html>

```
import torchvision.models as models  
  
model = models.vgg16(pretrained=True)
```

```
# Look at the model's layers  
model  
  
Output exceeds the size limit. Open the full output data in a text editor  
VGG(  
    features: Sequential(  
        (0): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (1): ReLU(inplace=True)  
        (2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (3): ReLU(inplace=True)  
        (4): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (5): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (6): ReLU(inplace=True)  
        (7): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (8): ReLU(inplace=True)  
        (9): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (10): Conv2d(128, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (11): ReLU(inplace=True)  
        (12): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (13): ReLU(inplace=True)  
        (14): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (15): ReLU(inplace=True)  
        (16): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (17): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (18): ReLU(inplace=True)  
        (19): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (20): ReLU(inplace=True)  
        (21): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (22): ReLU(inplace=True)  
    ...  
        (4): ReLU(inplace=True)  
        (5): Dropout(p=0.5, inplace=False)  
        (6): Linear(in_features=4096, out_features=1000, bias=True)  
    )  
)
```

```
import torchvision.models as models  
  
model = models.mobilenet_v2(pretrained=True)
```

- AlexNet
- ConvNeXt
- DenseNet
- EfficientNet
- EfficientNetV2
- GoogLeNet
- Inception V3
- MaxVit
- MNASNet
- MobileNet V2
- MobileNet V3
- RegNet
- ResNet
- ResNeXt
- ShuffleNet V2
- SqueezeNet
- SwinTransformer
- VGG
- VisionTransformer
- Wide ResNet