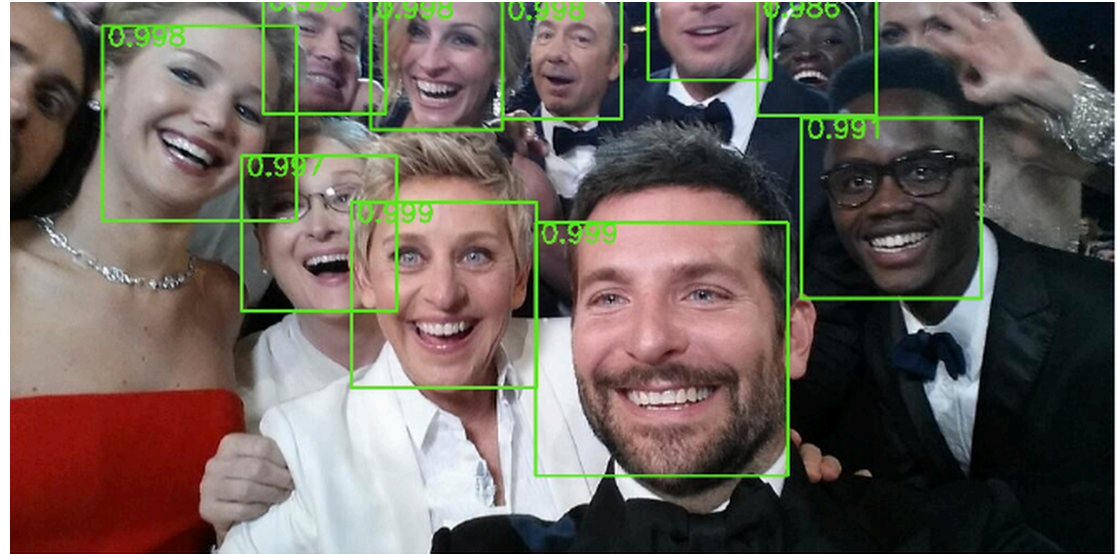


Convolutional Neural Networks

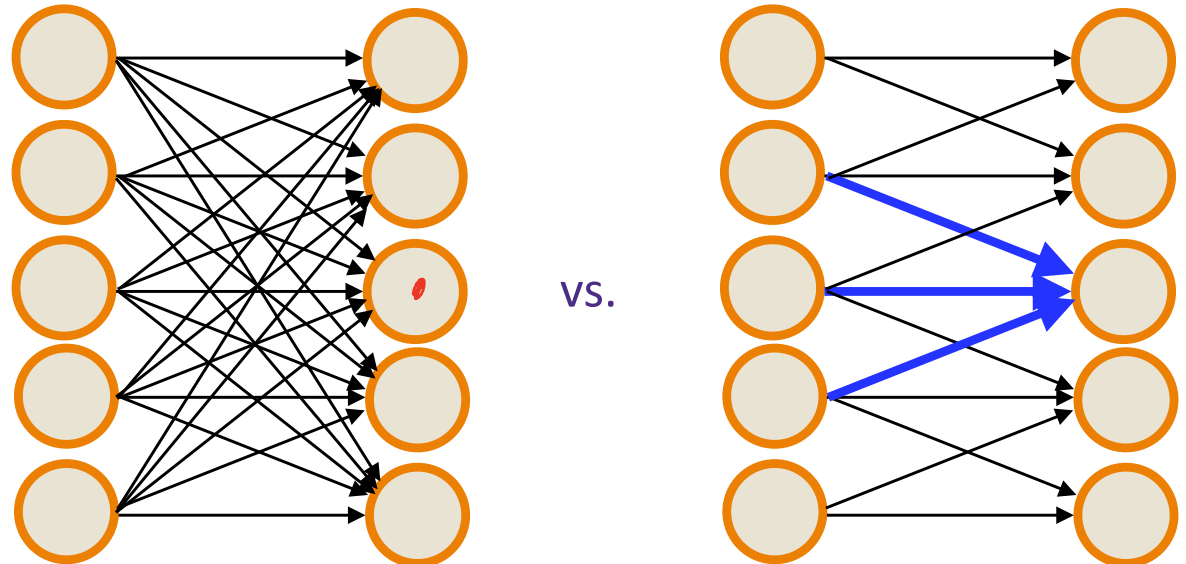


Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.



Similarly, to identify edges or other local structure, it makes sense to only look at **local information**

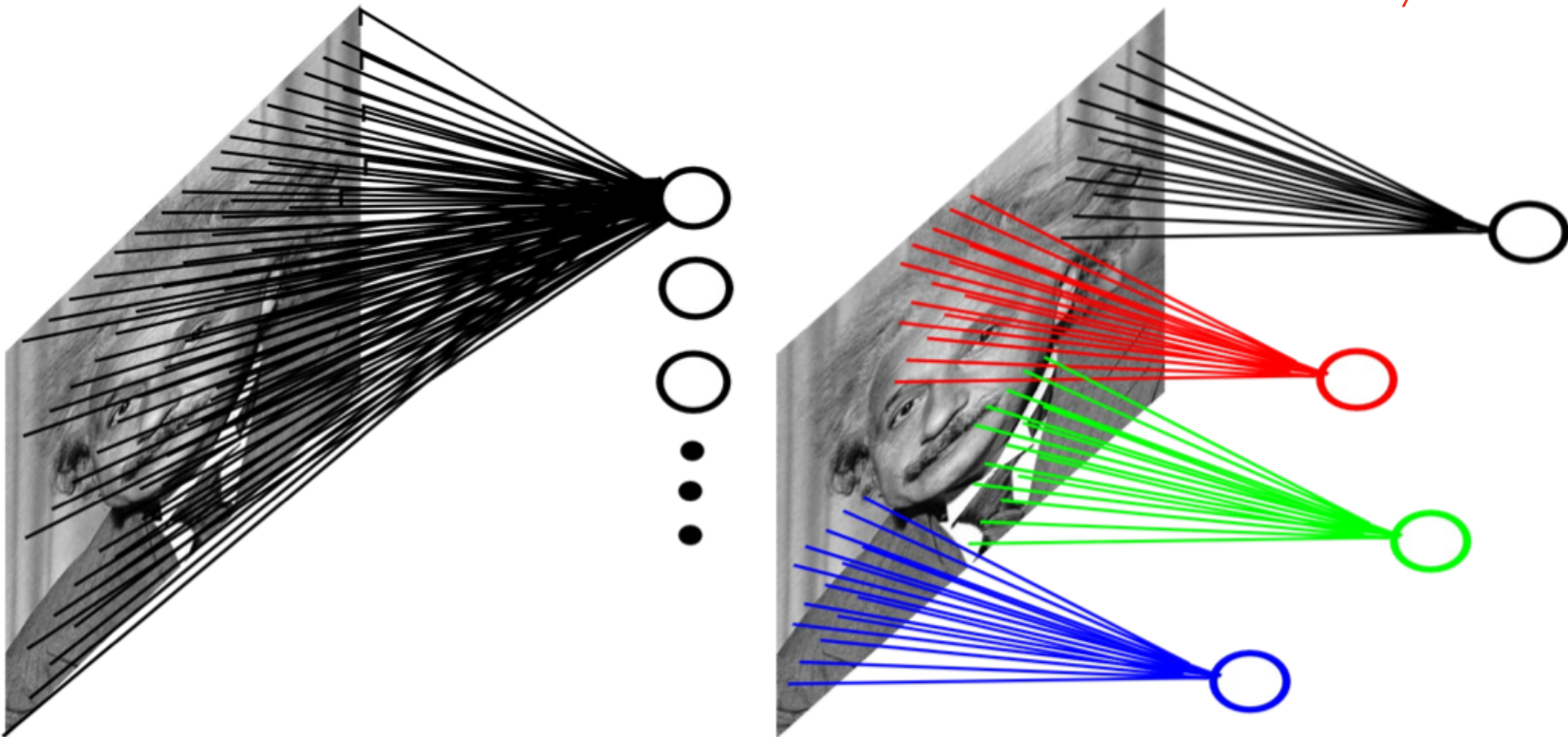


2d Convolution Layer

■ Example: 200x200 image

- ▶ Fully-connected, 400,000 hidden units = 16 billion parameters
- ▶ Locally-connected, 400,000 hidden units 10x10 fields = 40 million params
- ▶ Local connections capture local dependencies

weight-sharing



Convolution of images (2d convolution)

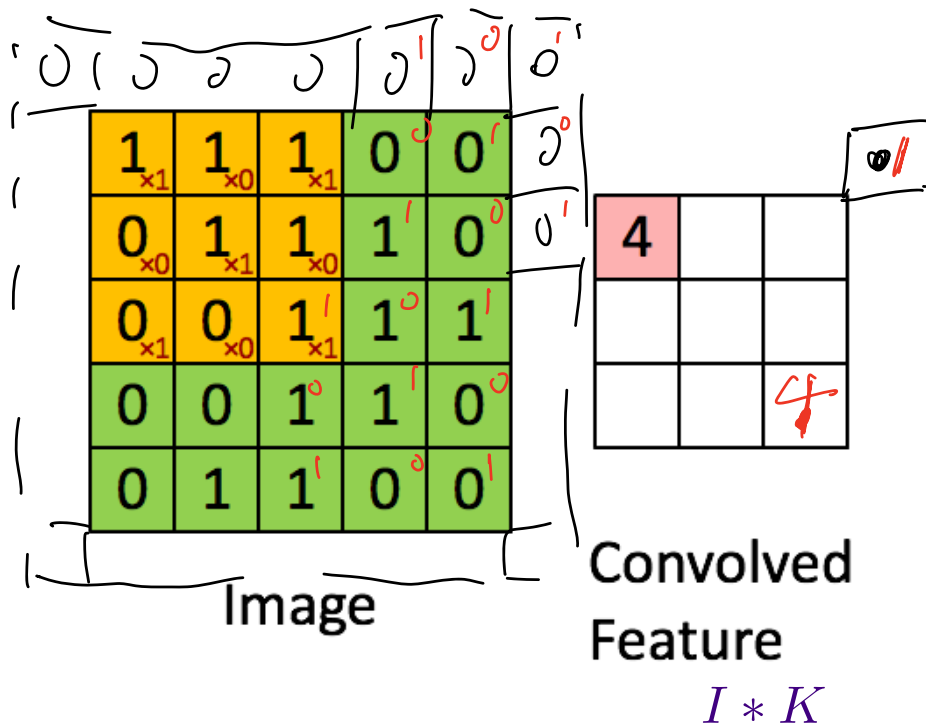
$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image I

1	0	1
0	1	0
1	0	1

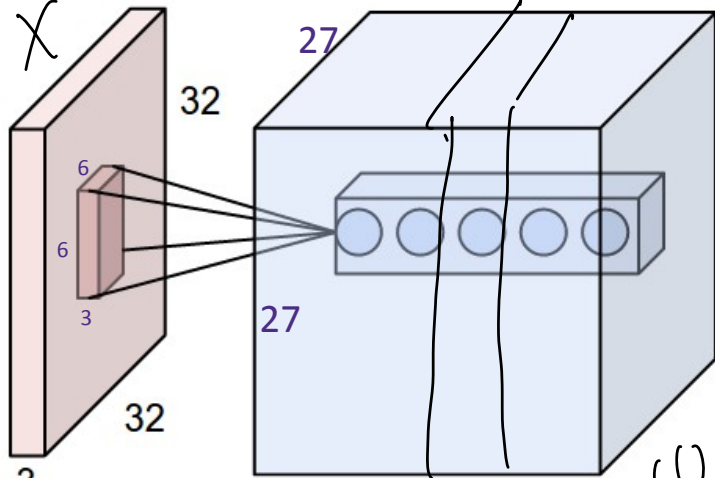
Filter K



Stacking convolved images

g : activation function

$\bar{\gamma} = 1, \dots, d$
 $a_{\bar{\gamma}}^{(1)}$

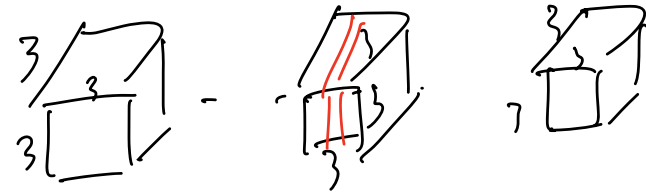


R G B
 (channels)

d filters

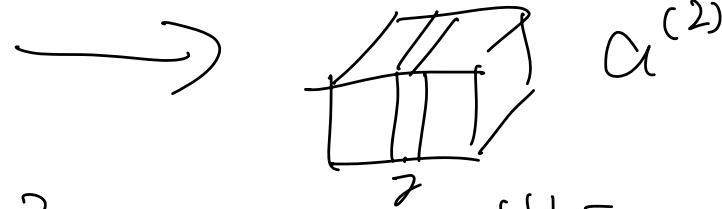
$a^{(1)}$

$$a_{\bar{\gamma}}^{(1)} = g \left(\sum_{\alpha=1}^3 X[:, :, \alpha] * K_{\bar{\gamma}}^{(1)}[:, :, \alpha] \right)$$



Kernel i

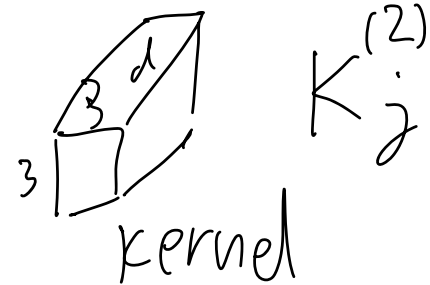
Repeat with d filters!



$a^{(2)}$

$$a_{\bar{j}}^{(2)} = g \left(\sum_{\bar{\gamma}=1}^d a^{(1)}[:, :, \bar{\gamma}] * K_{\bar{j}}^{(2)}[:, :, \bar{\gamma}] \right)$$

Second layer filters



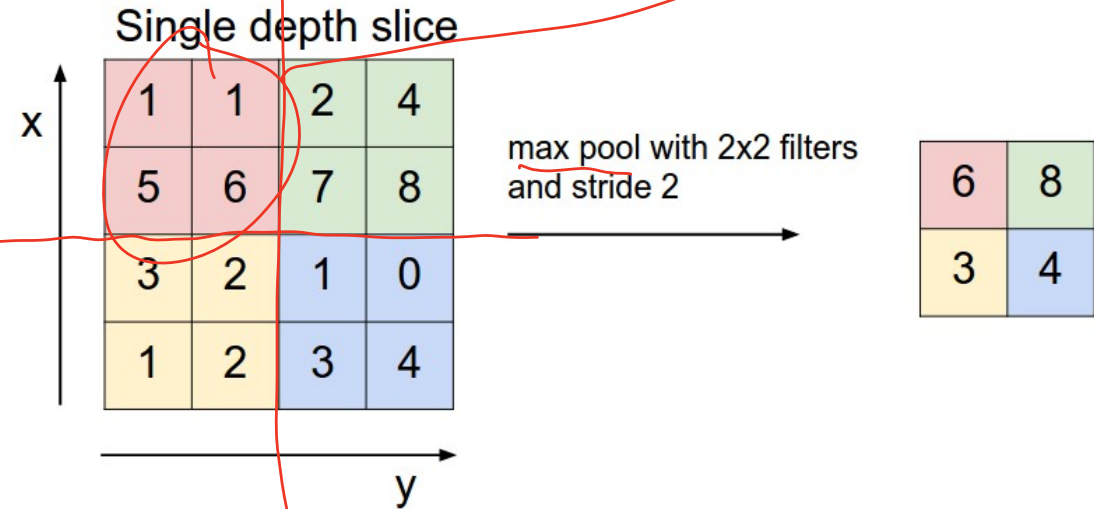
kernel

Pooling

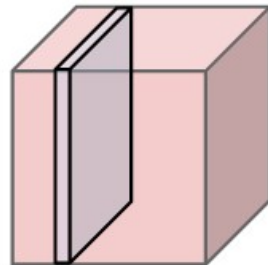
$$\frac{13}{4}$$

average pooling

Pooling reduces the dimension and can be interpreted as "This filter had a high response in this general region"

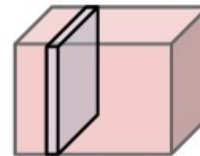


27x27x64

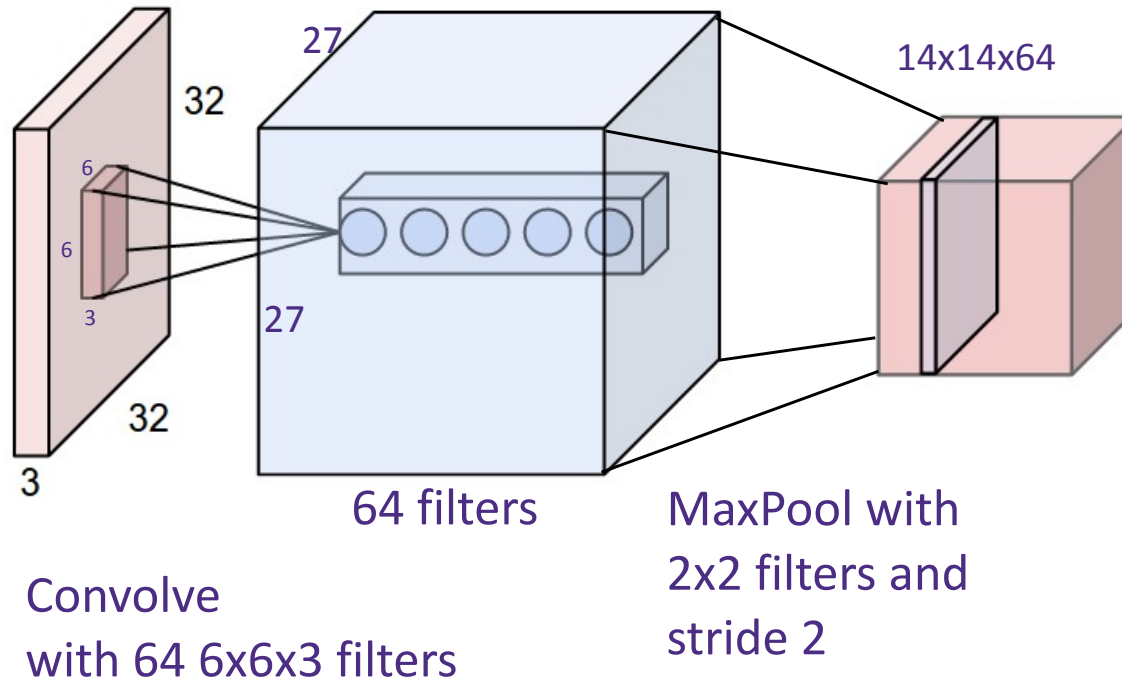


pool

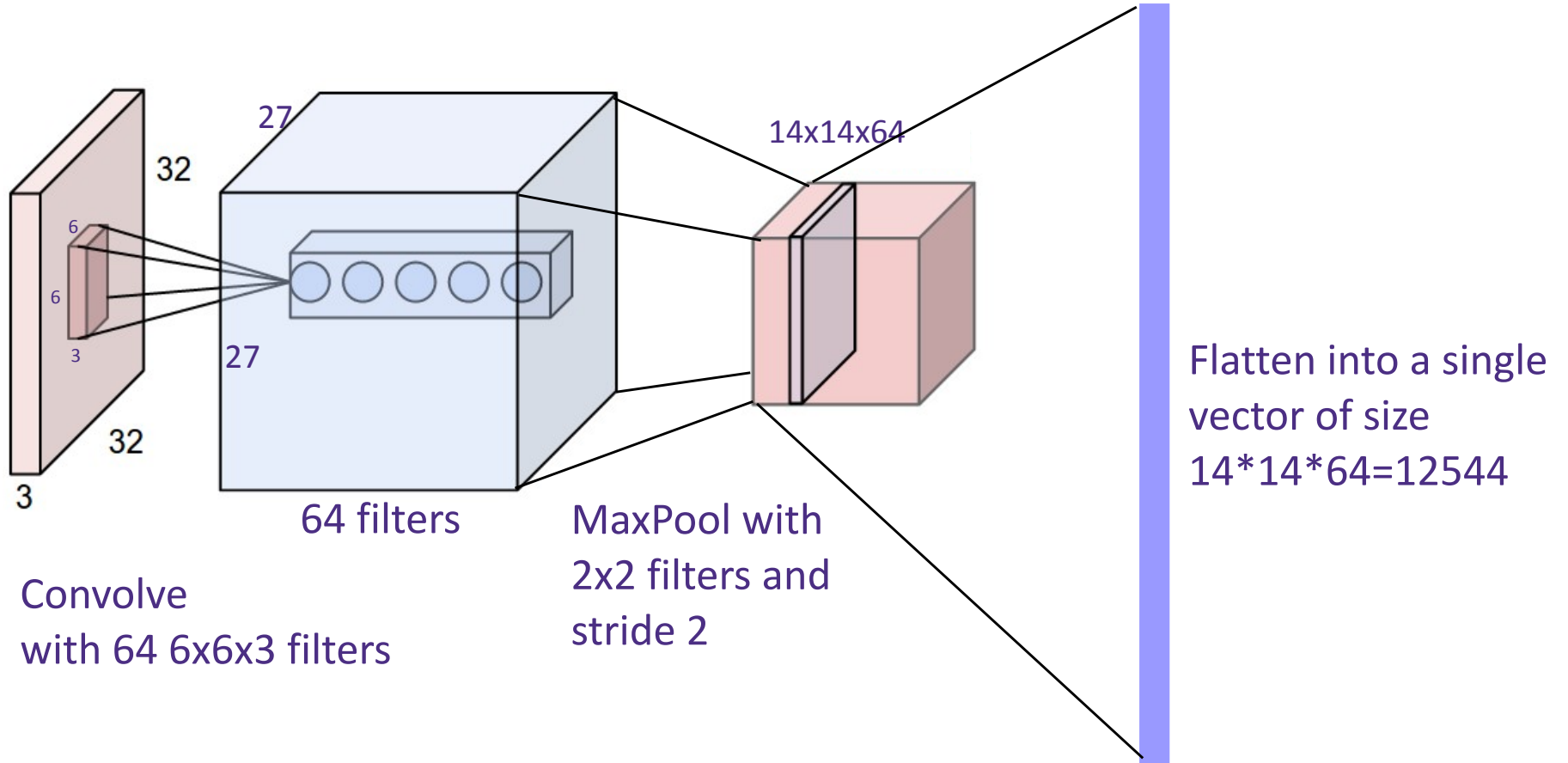
14x14x64



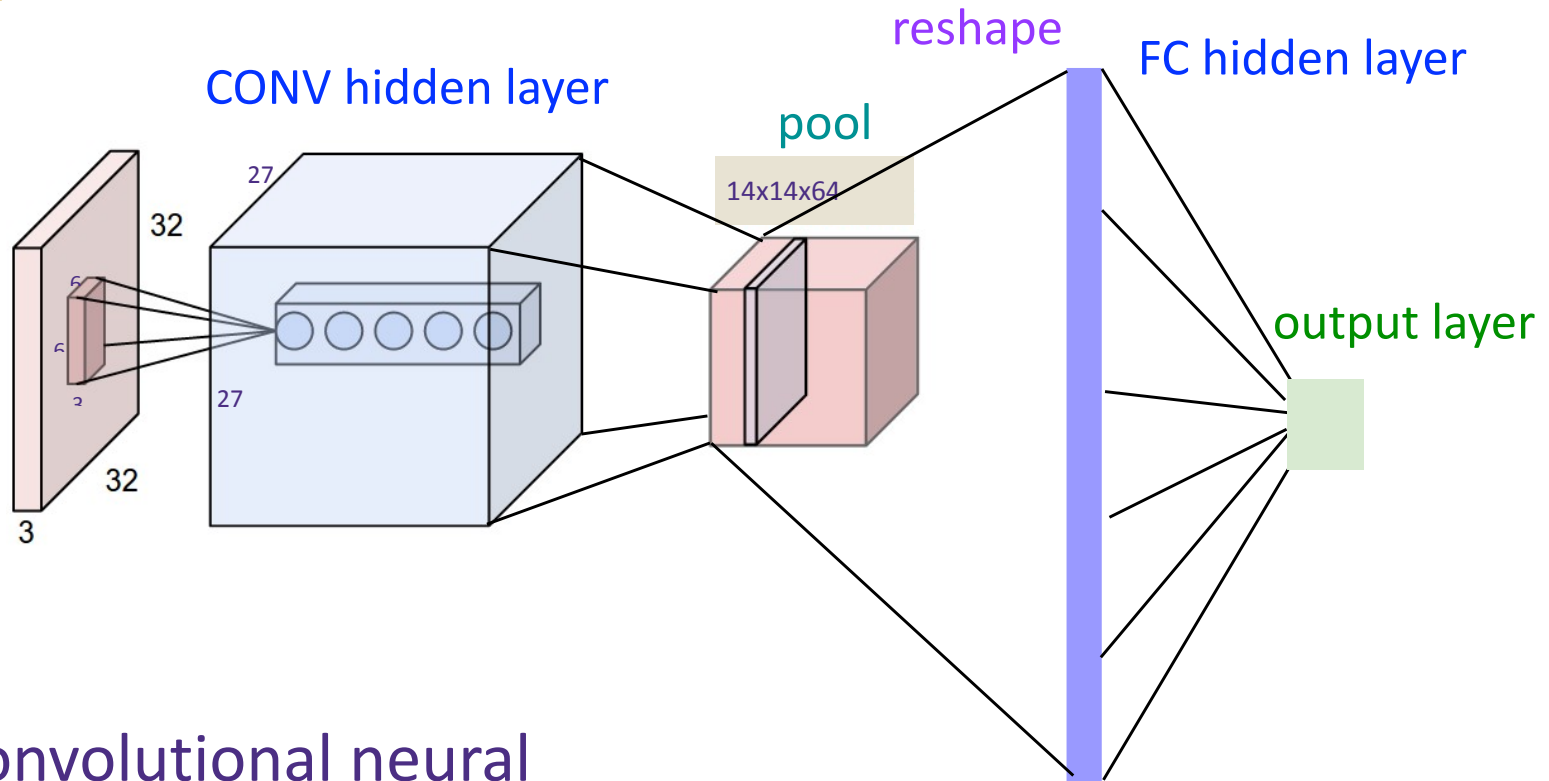
Pooling Convolution layer



Flattening

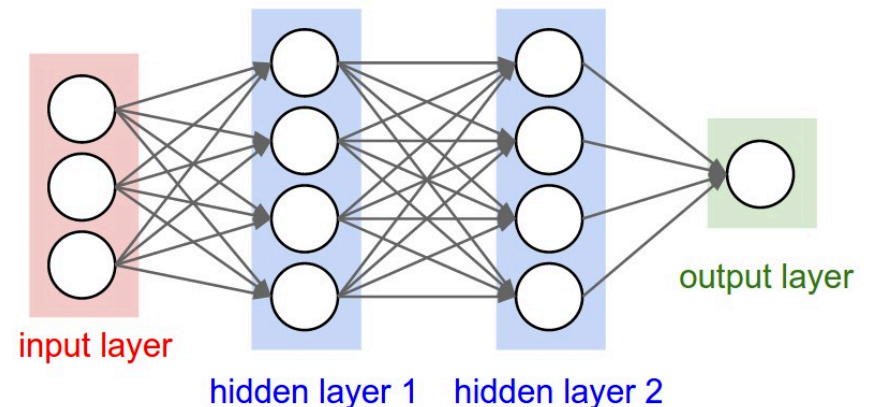


Training Convolutional Networks

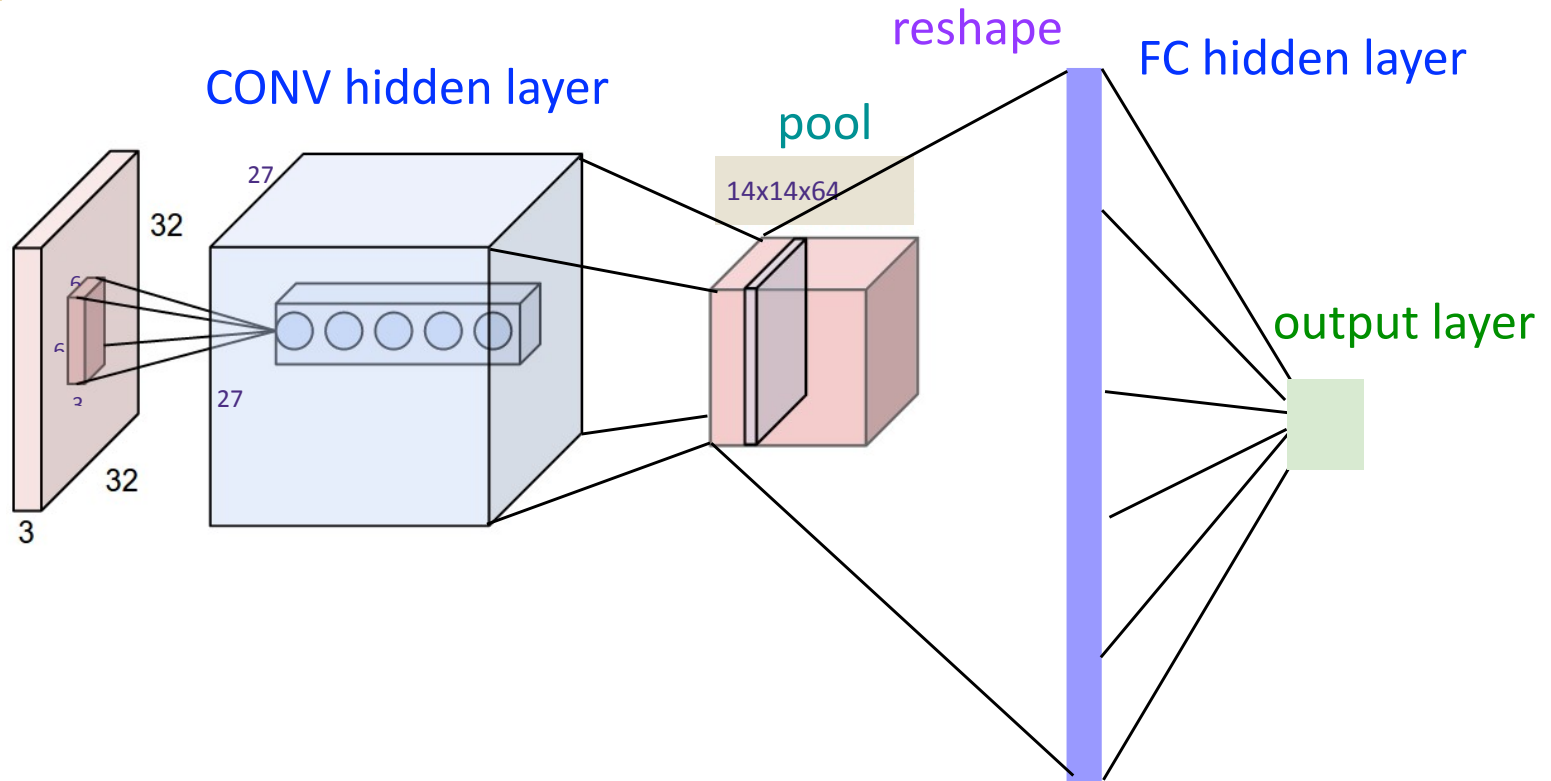


Recall: Convolutional neural networks (CNN) are just regular fully connected (FC) neural networks with some connections removed.

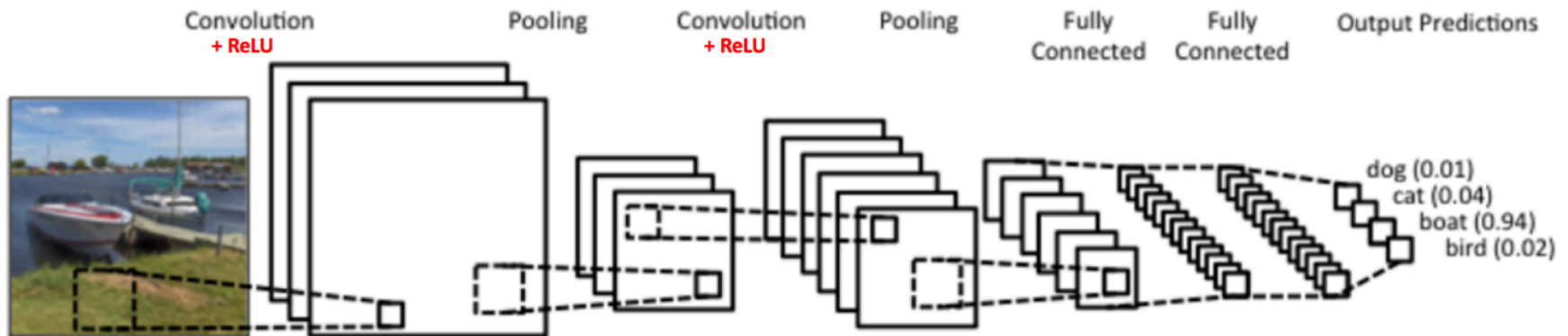
Train with SGD!

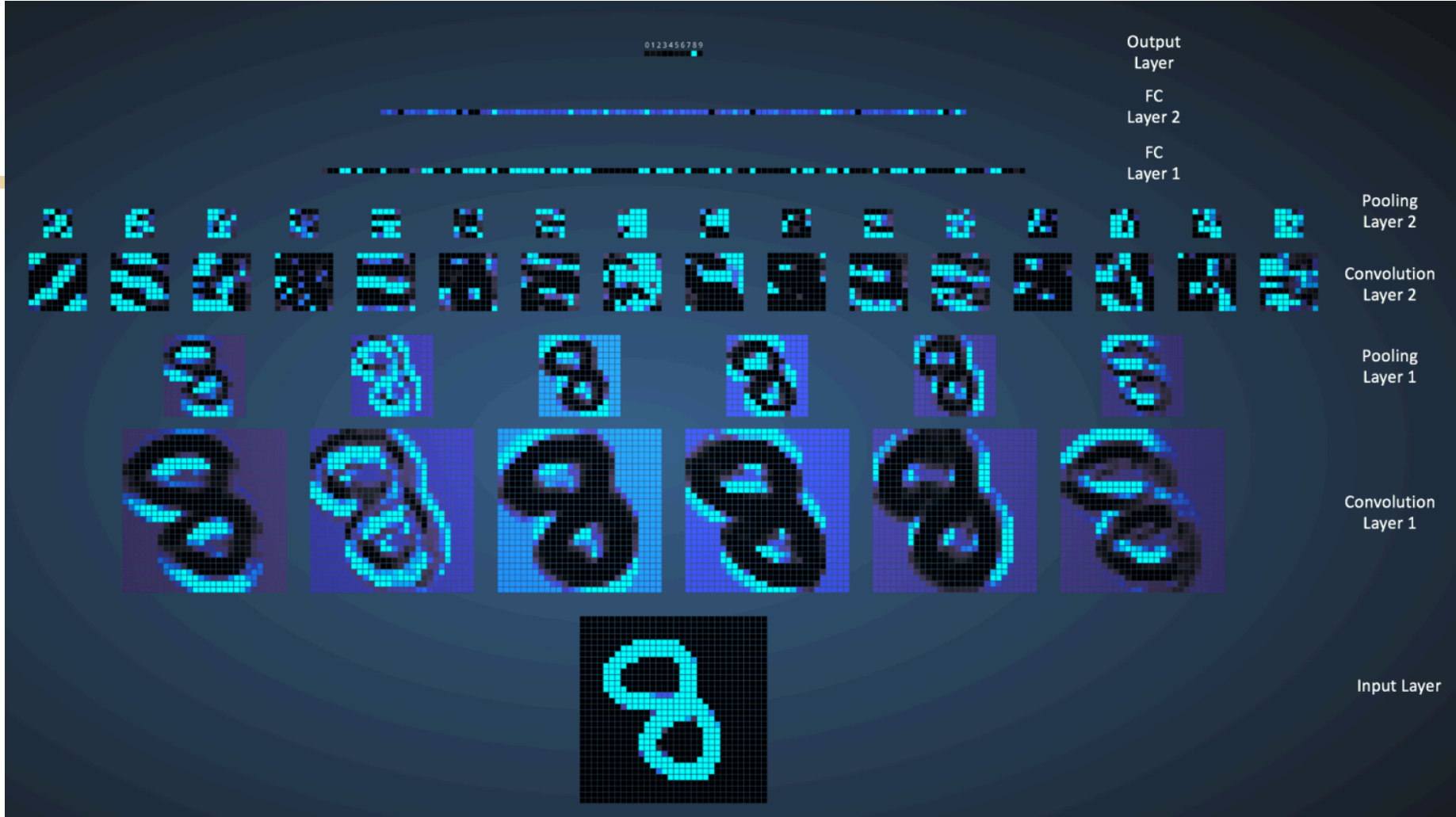


Training Convolutional Networks

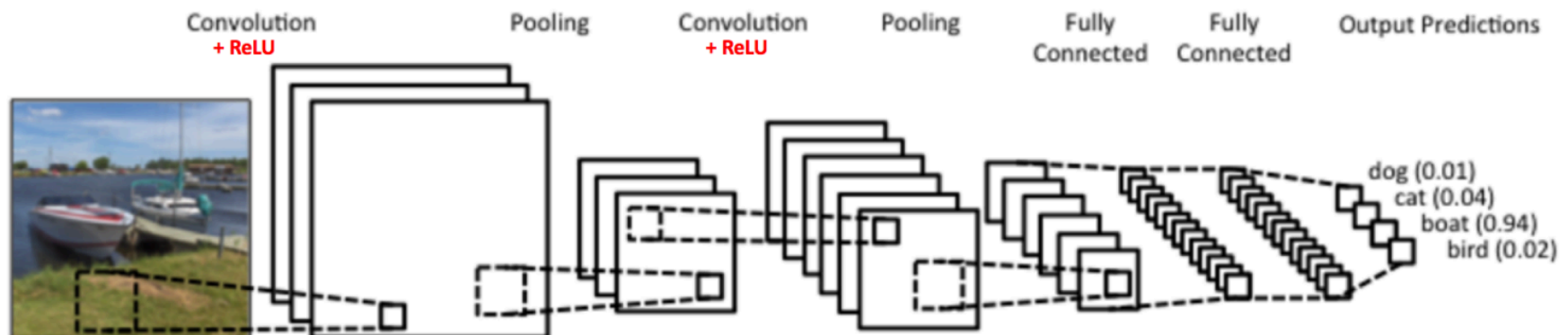


Real example network: LeNet





Real example network: LeNet



Famous CNNs

W

ImageNet Dataset

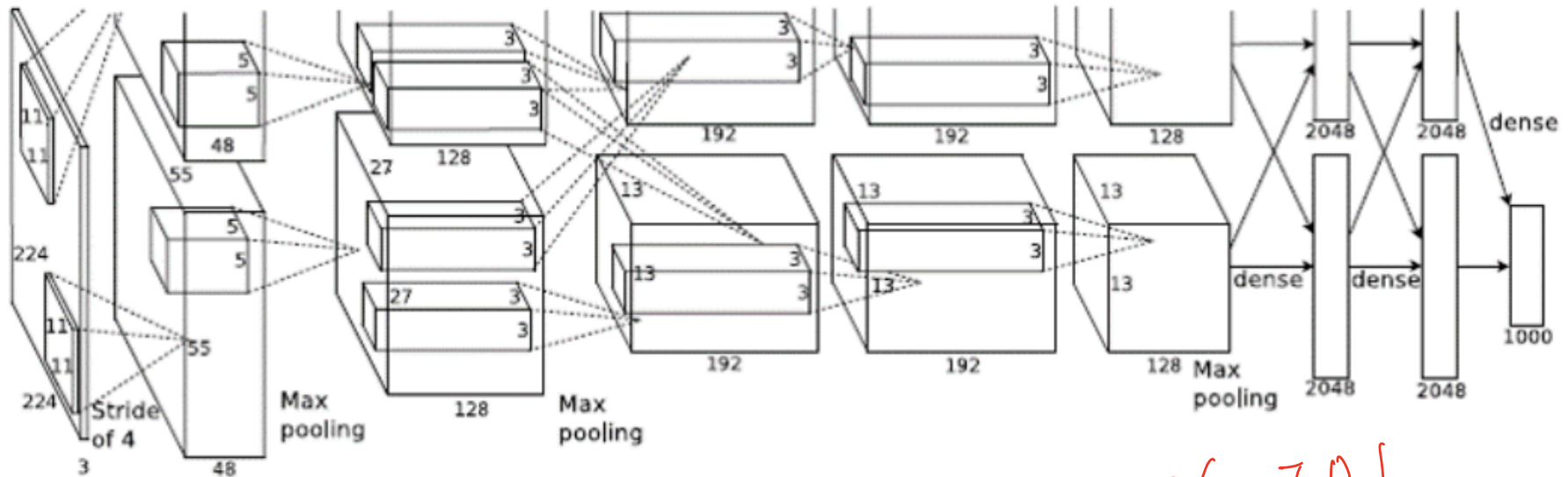
~14 million images, 20k classes



Deng et al. "Imagenet: a large scale hierarchical image database" '09

AlexNet

Breakthrough on ImageNet: ~the beginning of deep learning era



2012 New I/P)
Test of time

Krizhevsky, Sutskever, Hinton “ImageNet Classification with Deep Convolutional Neural Networks”, NIPS 2012.

AlexNet

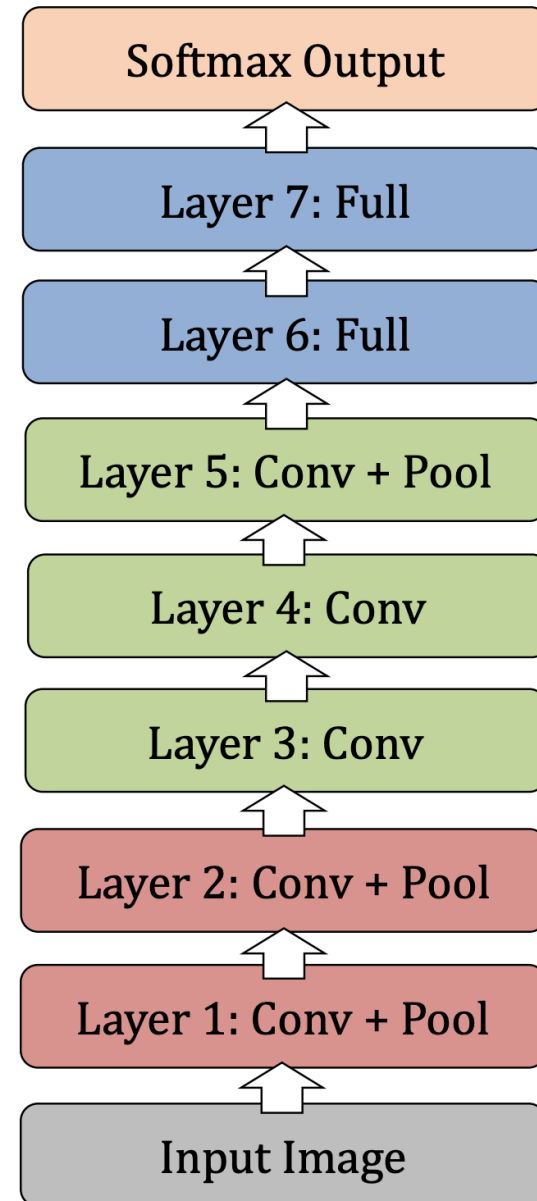
8 layers, ~60M parameters

Top5 error: 18.2%

Top1

Techniques used:

ReLU activation, overlapping pooling, dropout, ensemble (create 10 $\hat{y}_1, \dots, \hat{y}_{10}$ patches by cropping and average the predictions), data-augmentation (intensity of RGB channels)



[From Rob Fergus' CIFAR 2016 tutorial]

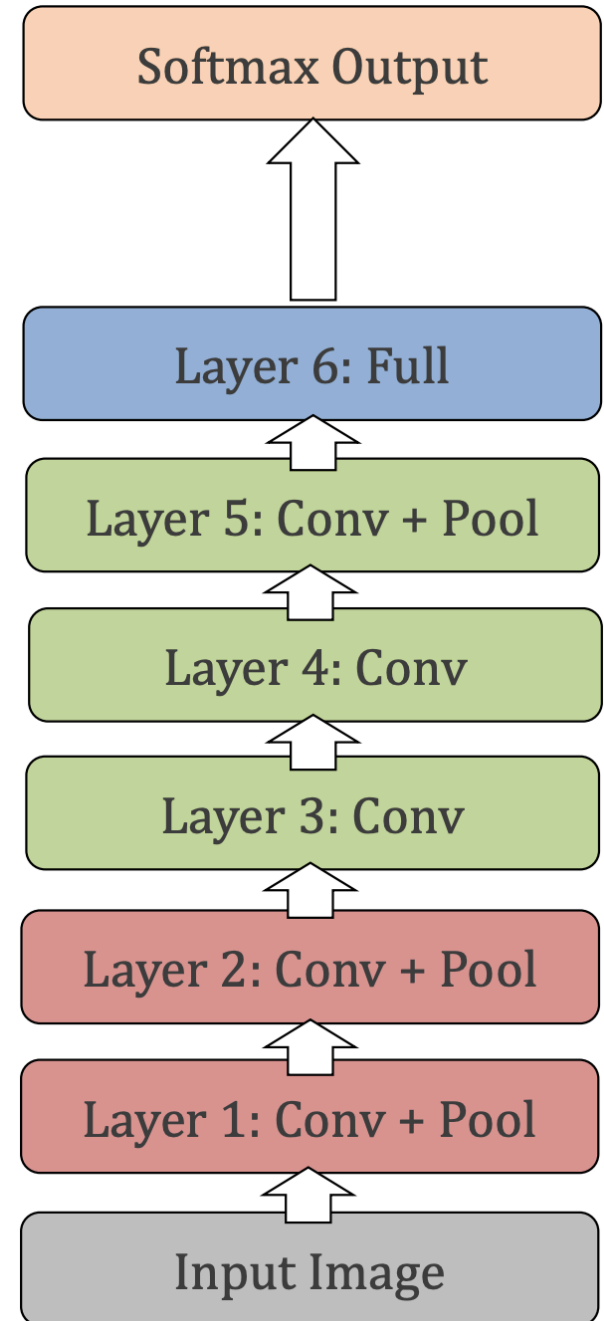
AlexNet

Remove top fully-connected layer 7

Drop ~16 million parameters

1.1% drop in performance

[From Rob Fergus' CIFAR 2016 tutorial]

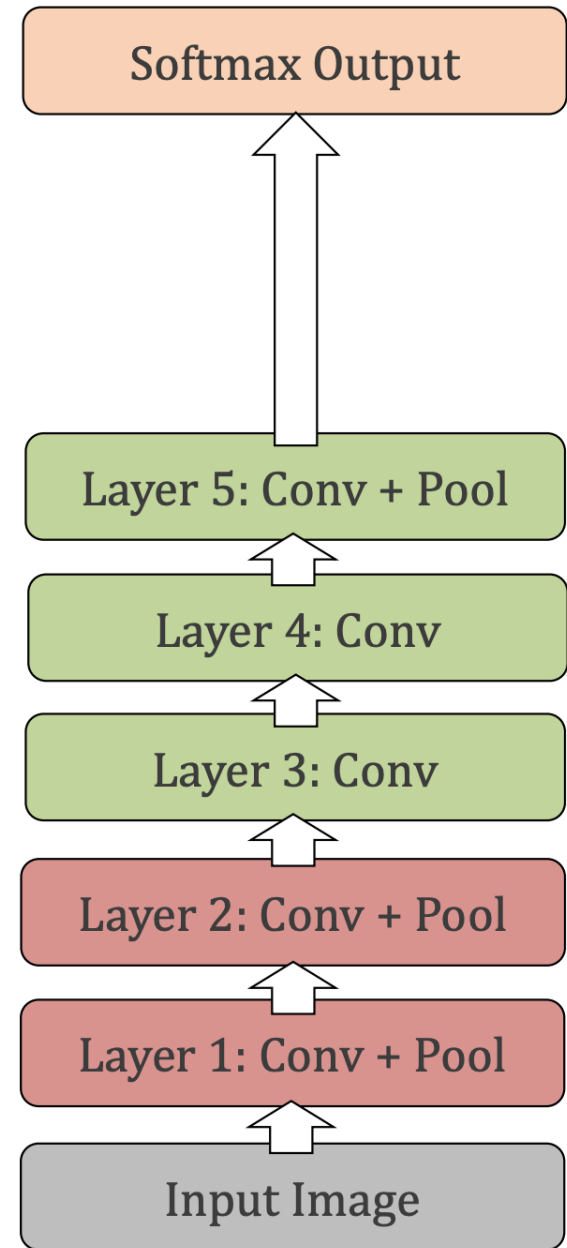


AlexNet

Remove both fully connected layers 6 and 7

Drop ~50 million parameters

5.7% drop in performance



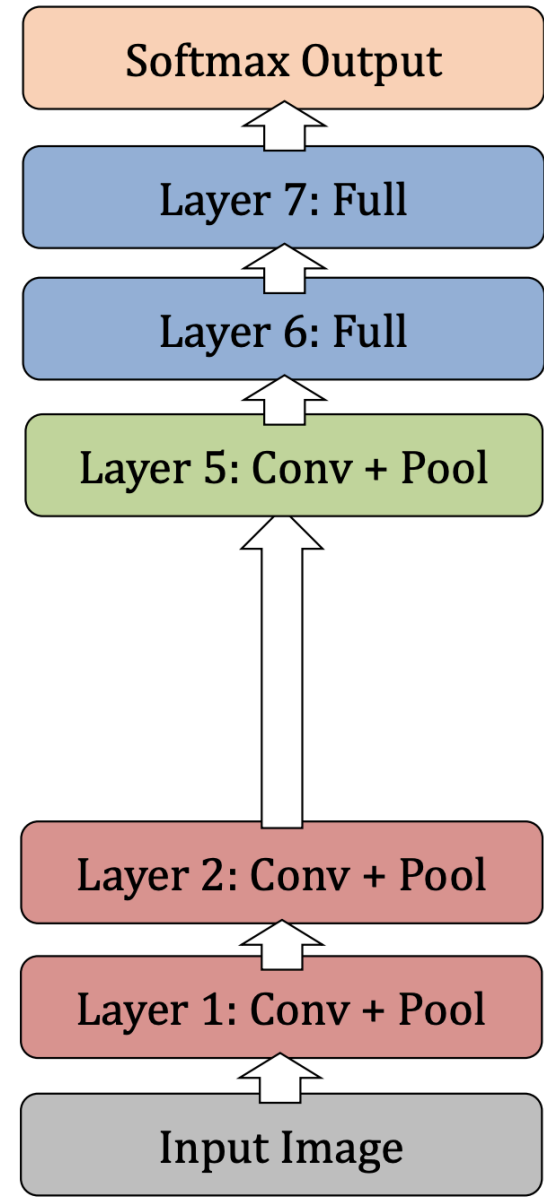
[From Rob Fergus' CIFAR 2016 tutorial]

AlexNet

Remove upper convolutio / feature extractor layers (layer 3 and 4)

Drop ~1 million parameters

3% drop in performance



[From Rob Fergus' CIFAR 2016 tutorial]

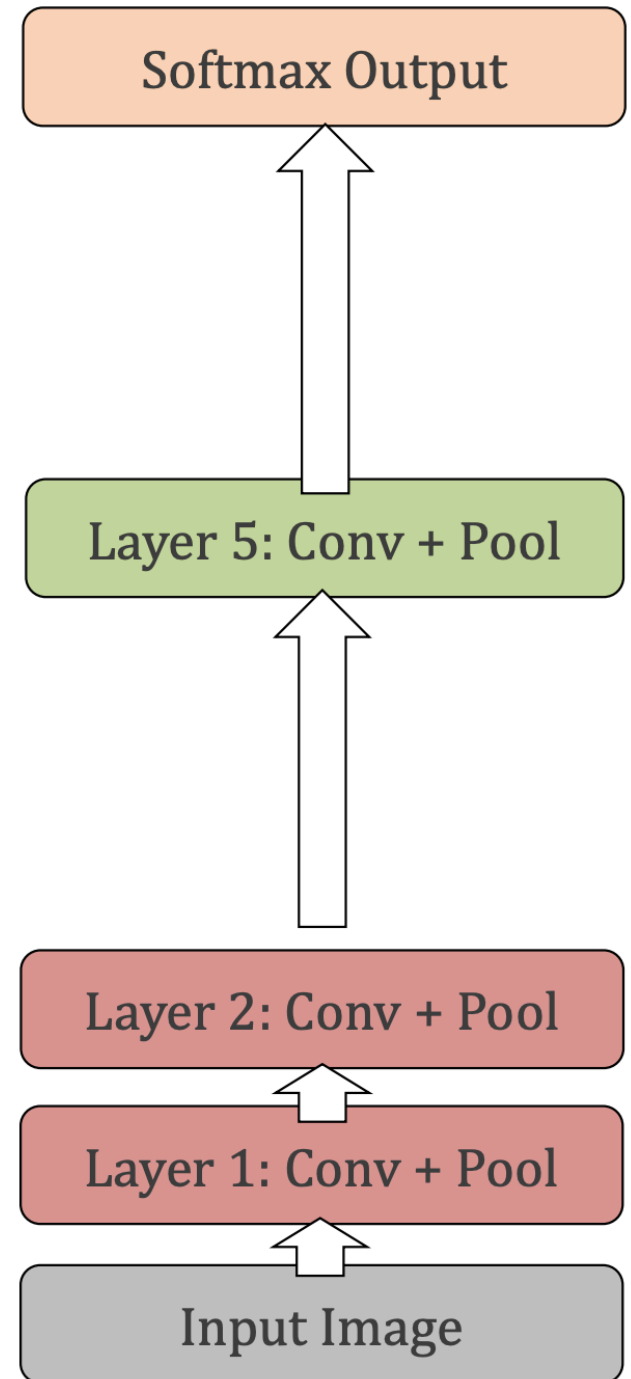
AlexNet

Remove top fully connected layer 6,7 and upper convolution layers 3,4.

33.5% drop in performance.

Depth of the network is the key.

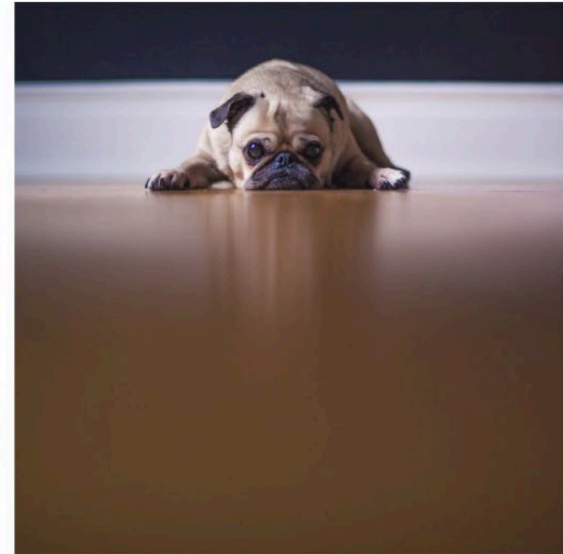
[From Rob Fergus' CIFAR 2016 tutorial]



GoogLeNet

7×7 , 5×5 , 3×3 , 1×1

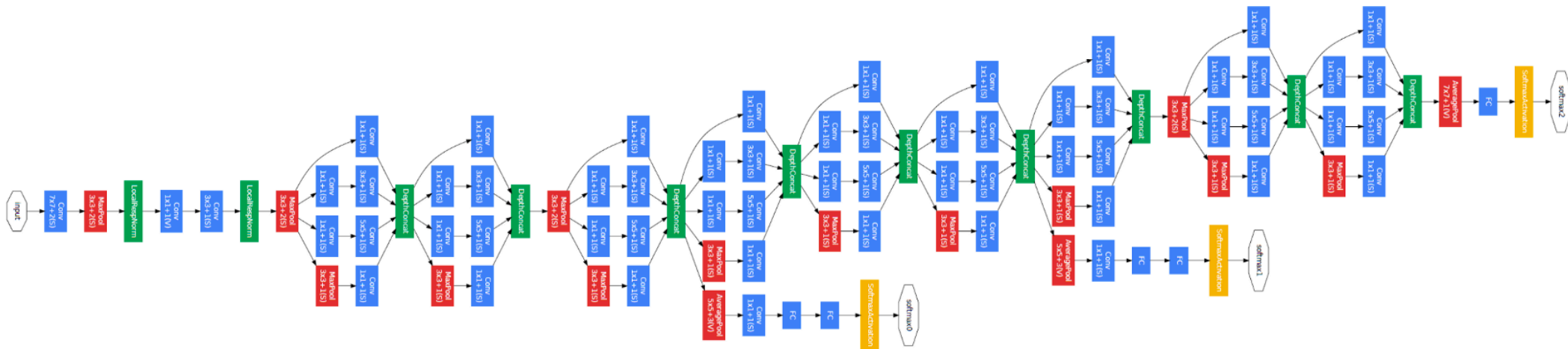
Motivation: multiscale nature of images



Large kernel for global features, and **smaller kernel** for local features.

Idea: have multiple different-size kernels at any layer.

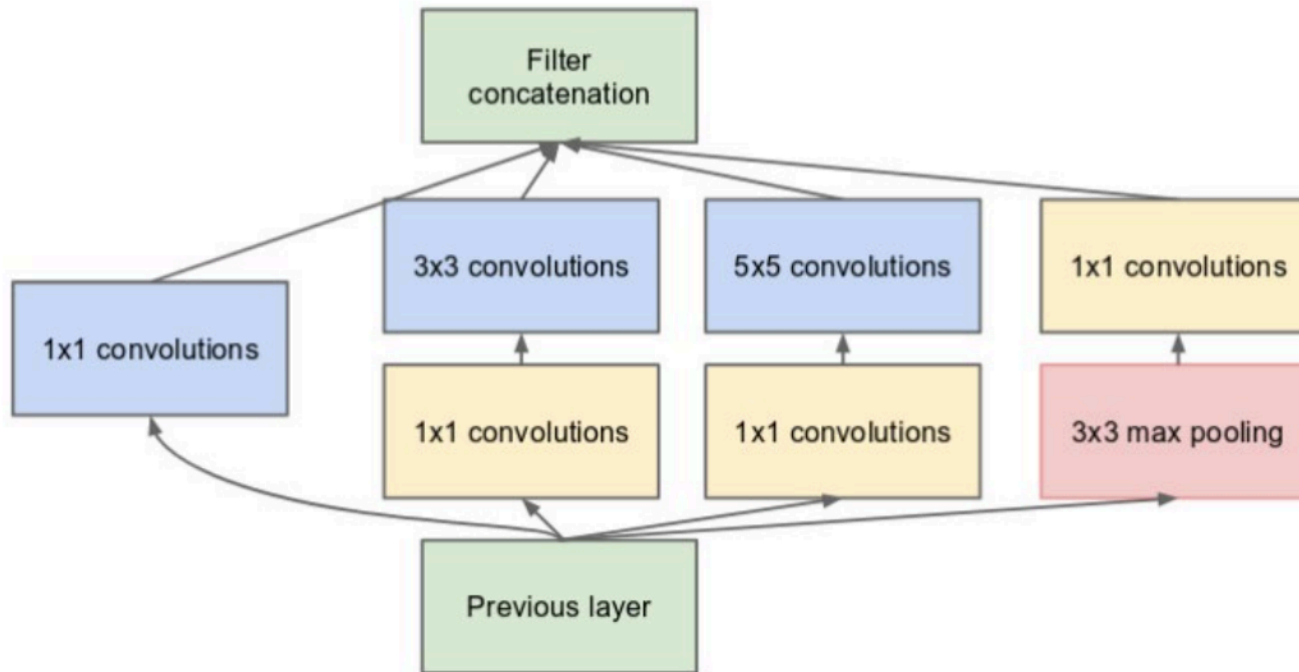
GoogLeNet



Large kernel for global features, and **smaller kernel** for local features.

Idea: have multiple different-size kernels at any layer.

Inception Module



Multiple filter scales at each layer

Dimensionality reduction to keep computational requirements down

Residual Networks

Motivation: extremely deep nets are hard to train (gradient explosion/vanishing)

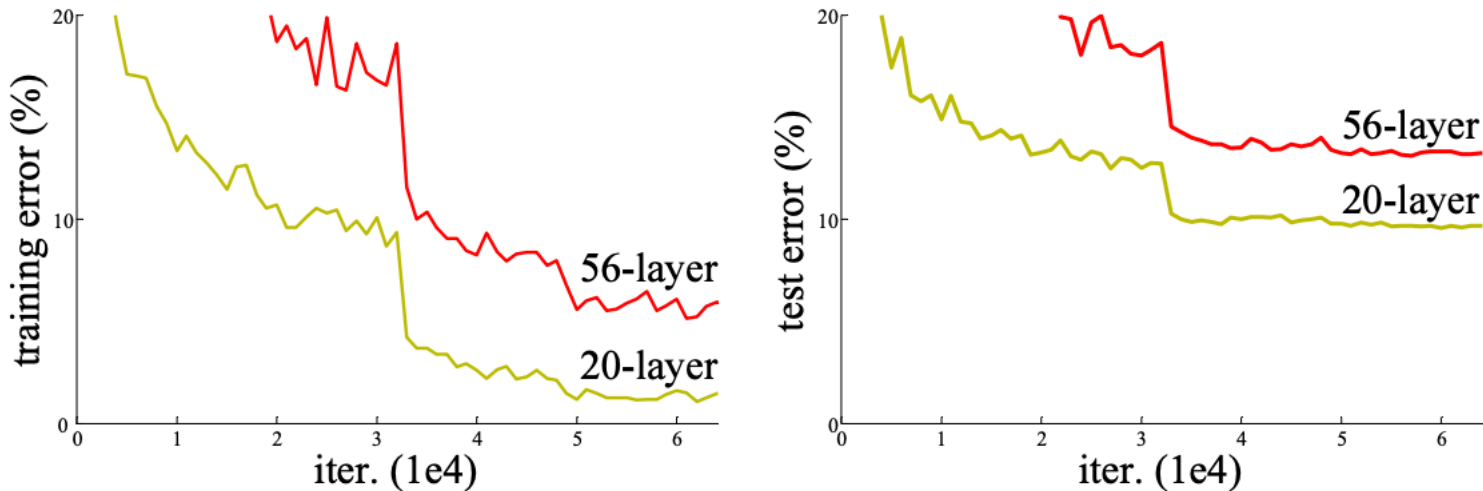
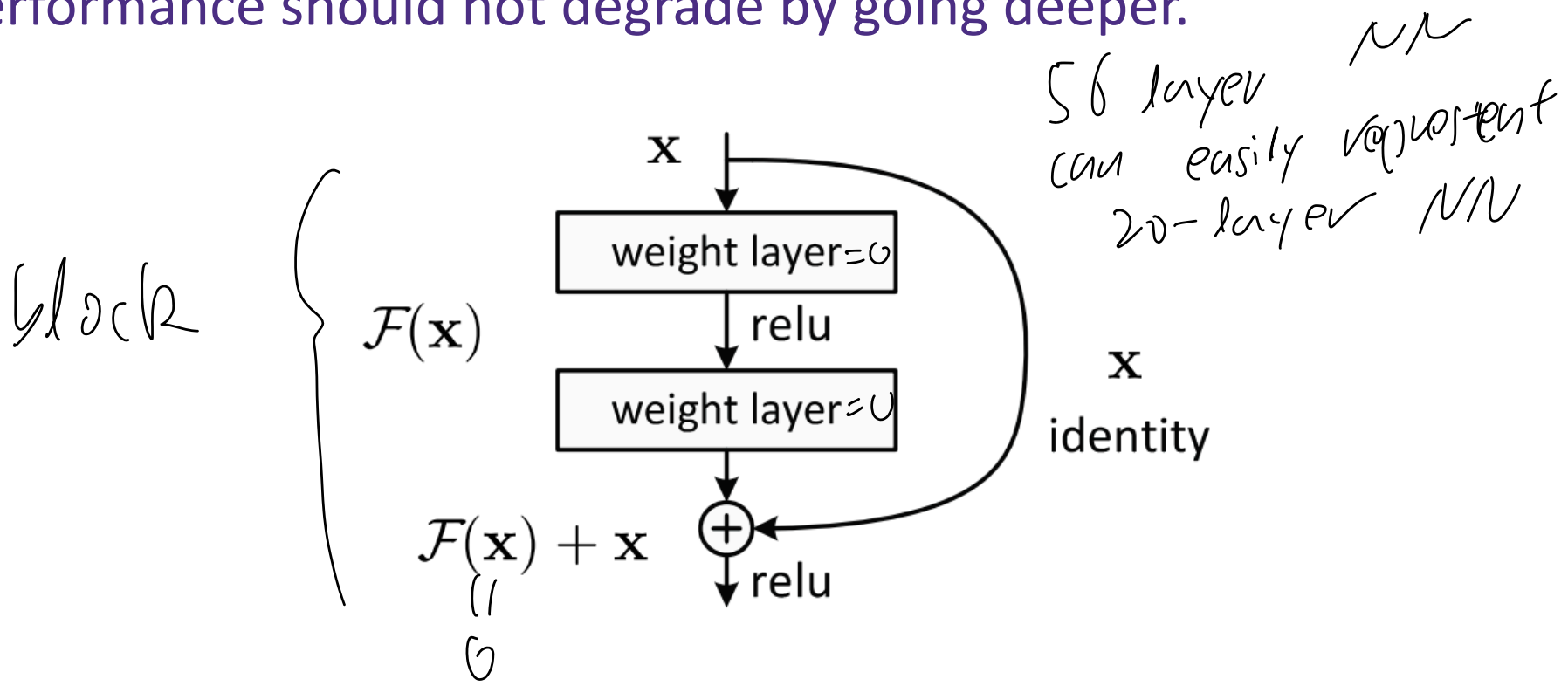


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

Residual Networks

Idea: identity shortcut, skip one or more layers.

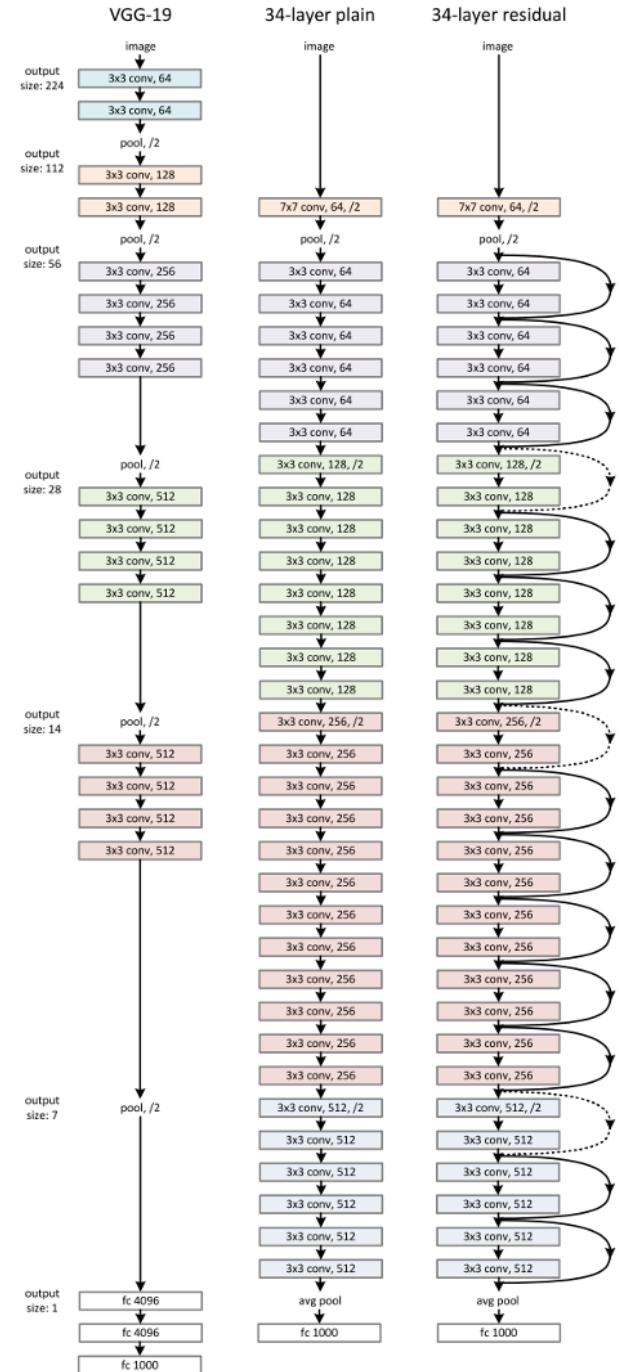
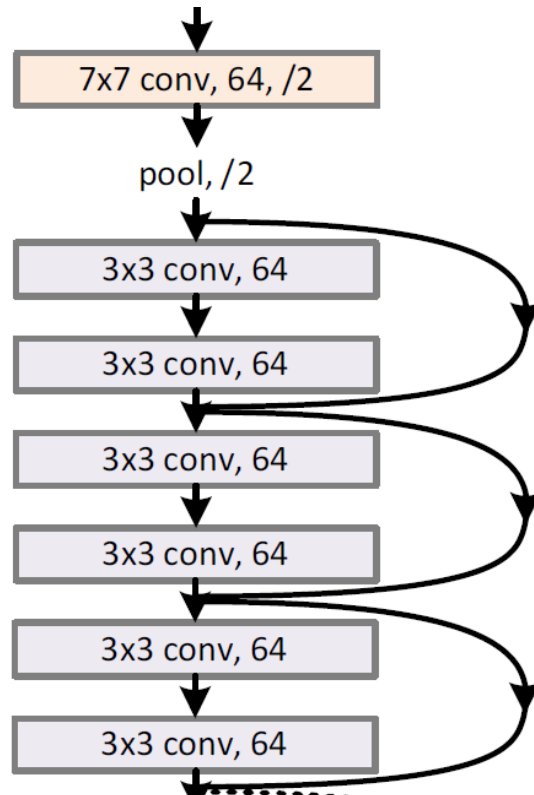
Justification: network can easily simulate shallow network ($F \approx 0$), so performance should not degrade by going deeper.



Residual Networks

vs. 18%

- 3.57% top-5 error on ImageNet
- First deep network with > 100 layers.
- Widely used in many domains (AlphaGo)

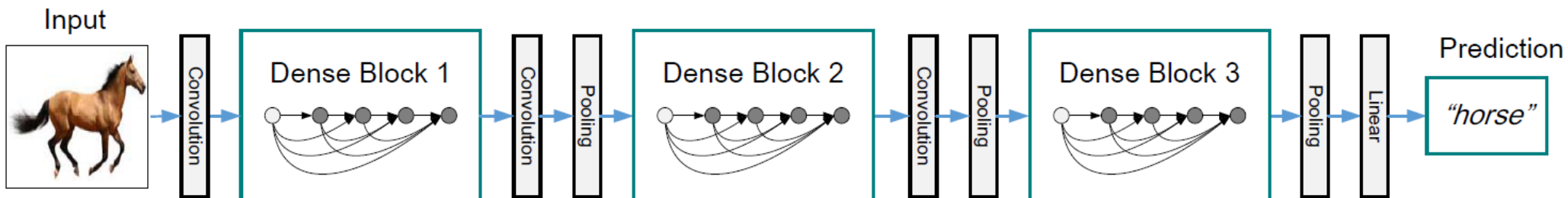
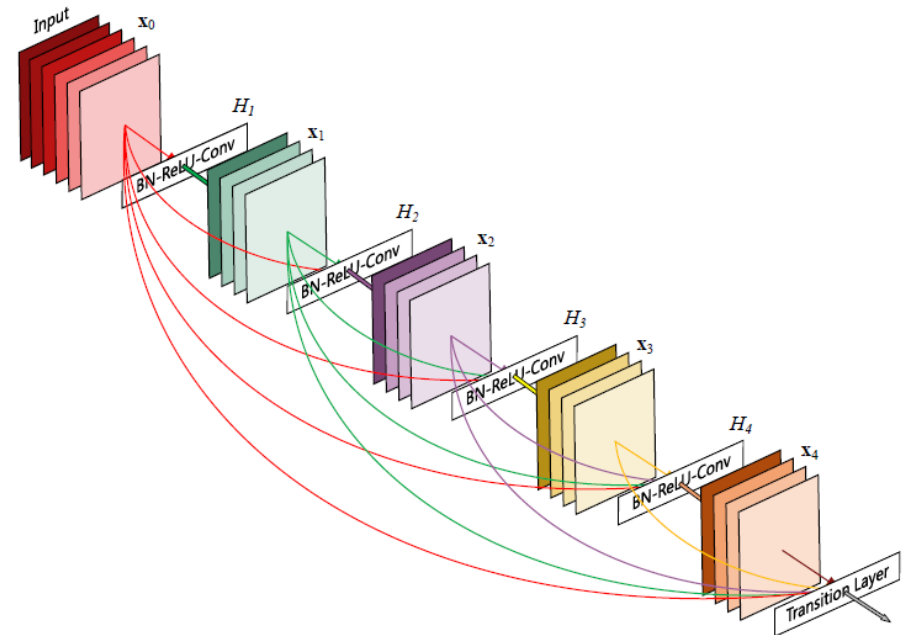


Densely Connected Network

Idea: explicit forward output of layer to all future layers (by concatenation)

Intuition: helps vanishing gradients, encourage reuse features (reduce parameter count)

Issues: network maybe too wide, need to be careful about memory consumption



Neural Architecture / Hyper-Parameter Search

Many design choices:

- Number of layers, width, kernel size, pooling, connections, etc.
- Normalization, learning rate, batch size, etc.

Strategies:

- Grid search
- Random search [Bergstra & Bengio '12]
- Bandit-based [Li et al. '16]
- Gradient-based (DARTS) [Liu et al. '19]
- Neural tangent kernel [Xu et al. '21]
- ...