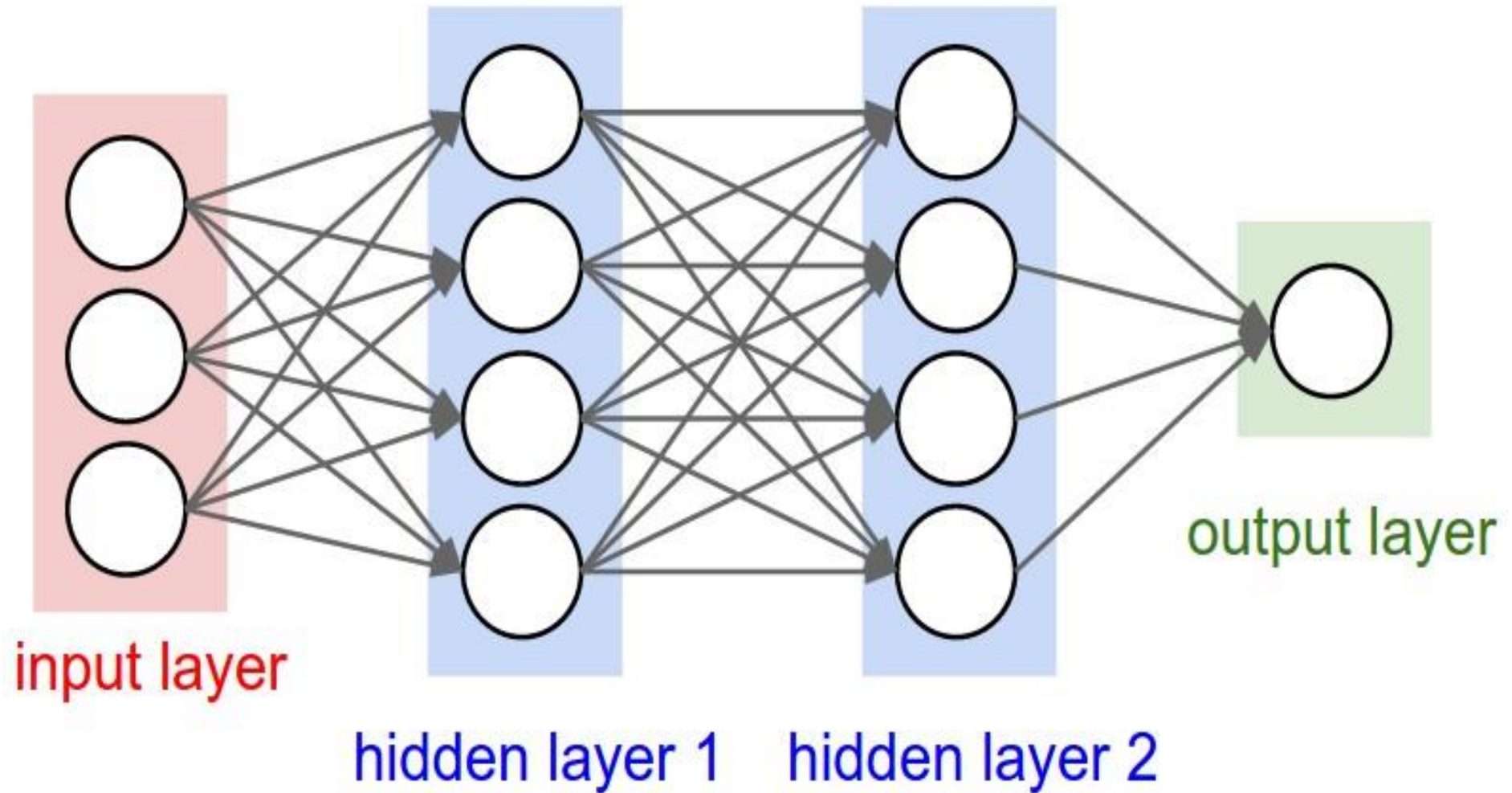


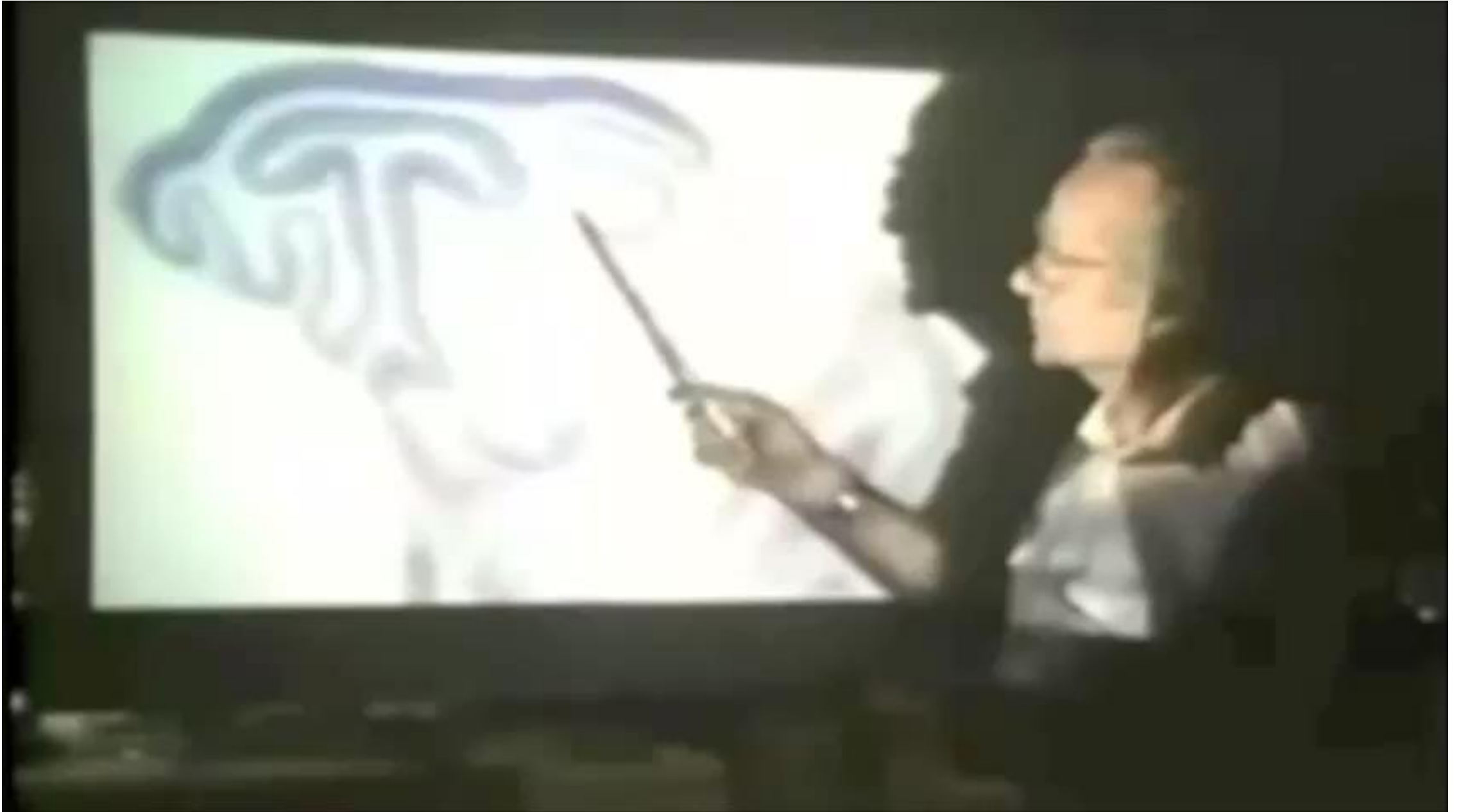
# Deep Learning

# Multi-layer perceptron



David Rumelhart, Geoffrey Hinton and Ronald Williams, 1986

# Convolution



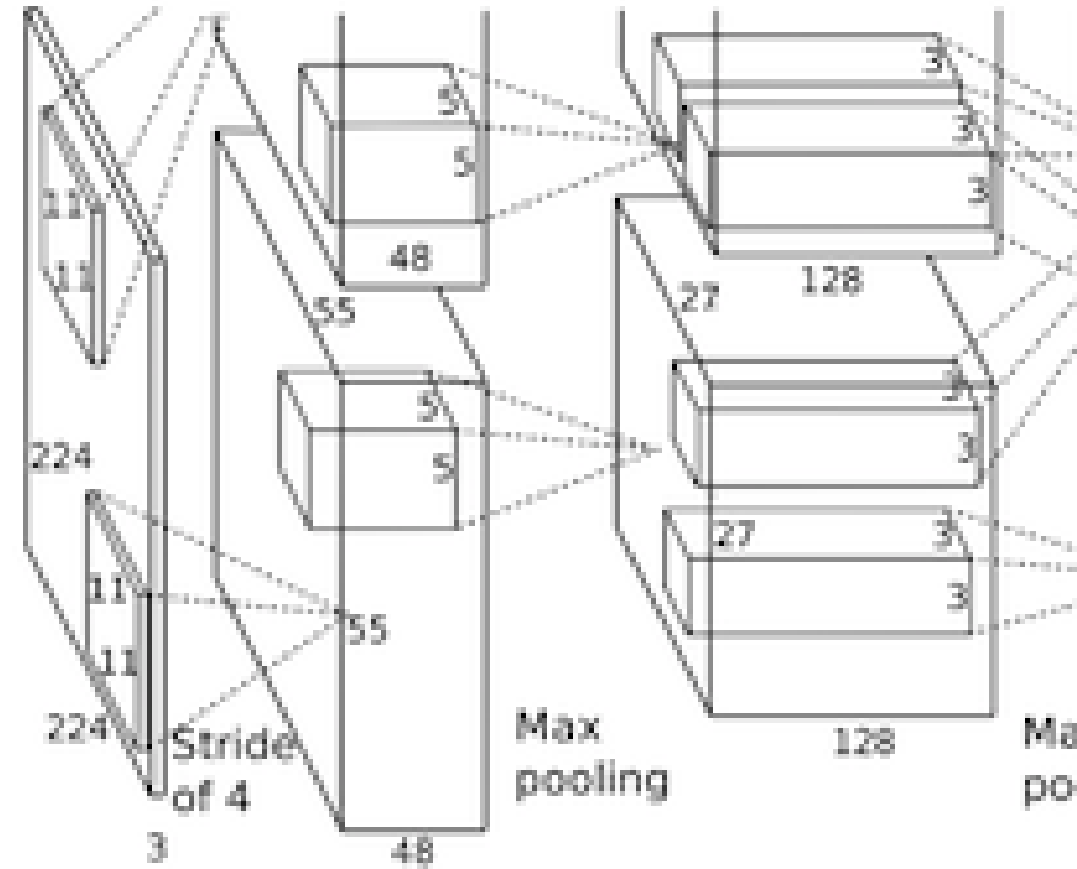
# Convolution

1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature



# Convolution

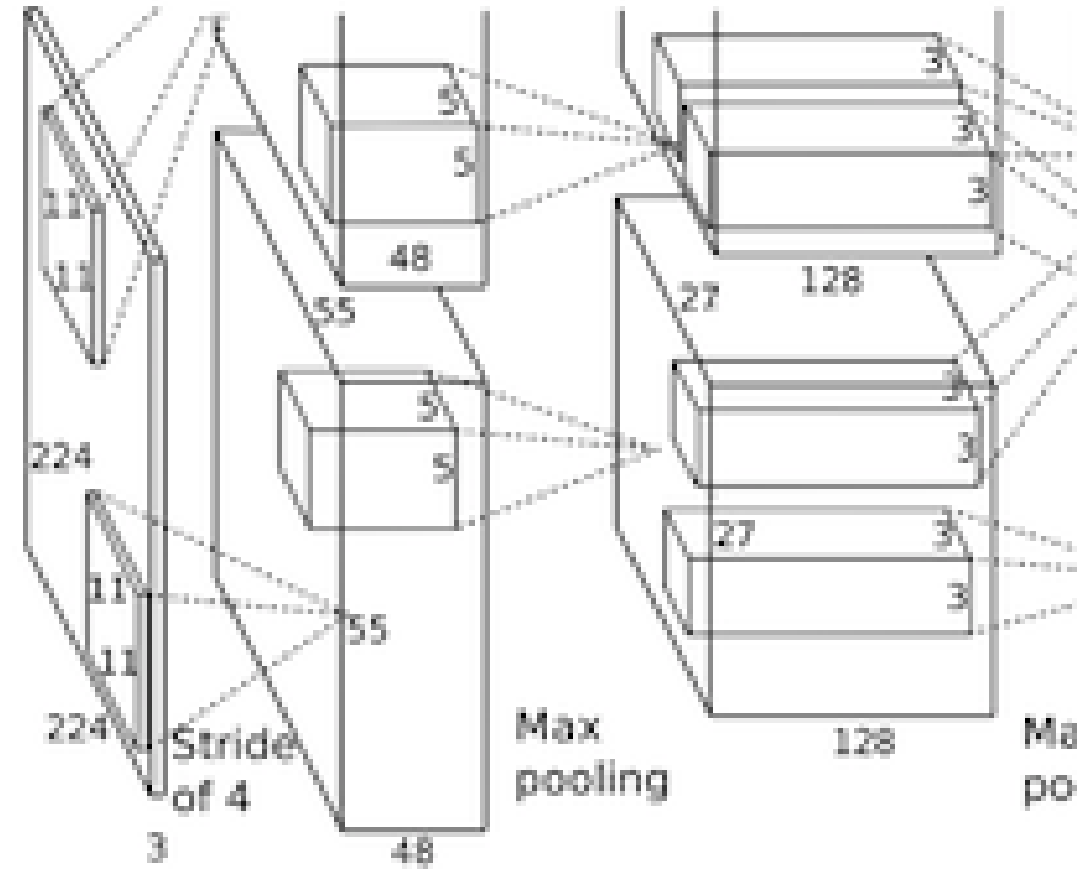
1 <sub>x1</sub>	1 <sub>x0</sub>	1 <sub>x1</sub>	0	0
0 <sub>x0</sub>	1 <sub>x1</sub>	1 <sub>x0</sub>	1	0
0 <sub>x1</sub>	0 <sub>x0</sub>	1 <sub>x1</sub>	1	1
0	0	1	1	0
0	1	1	0	0

Image

Kernel 3x3, Stride 1

4		

Convolved  
Feature

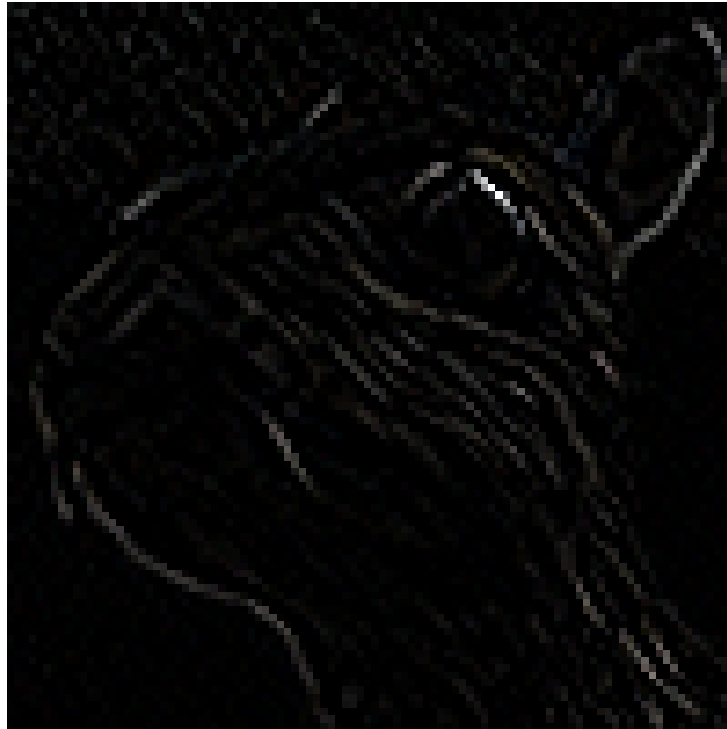


# Convolution examples

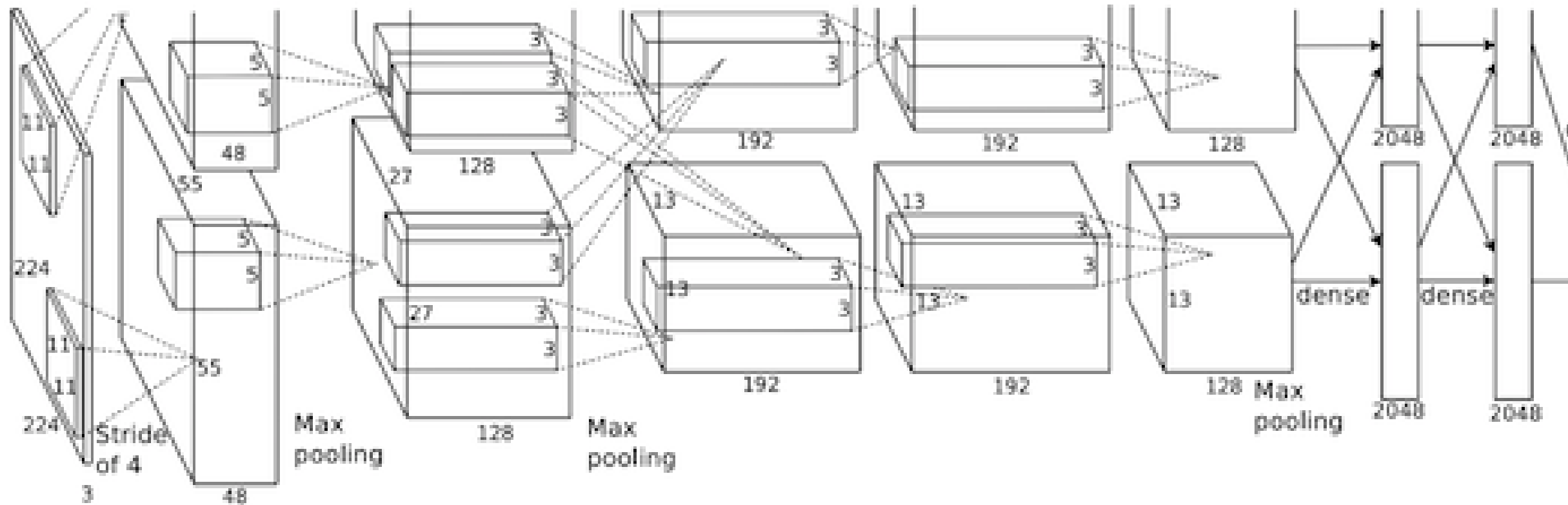
$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$

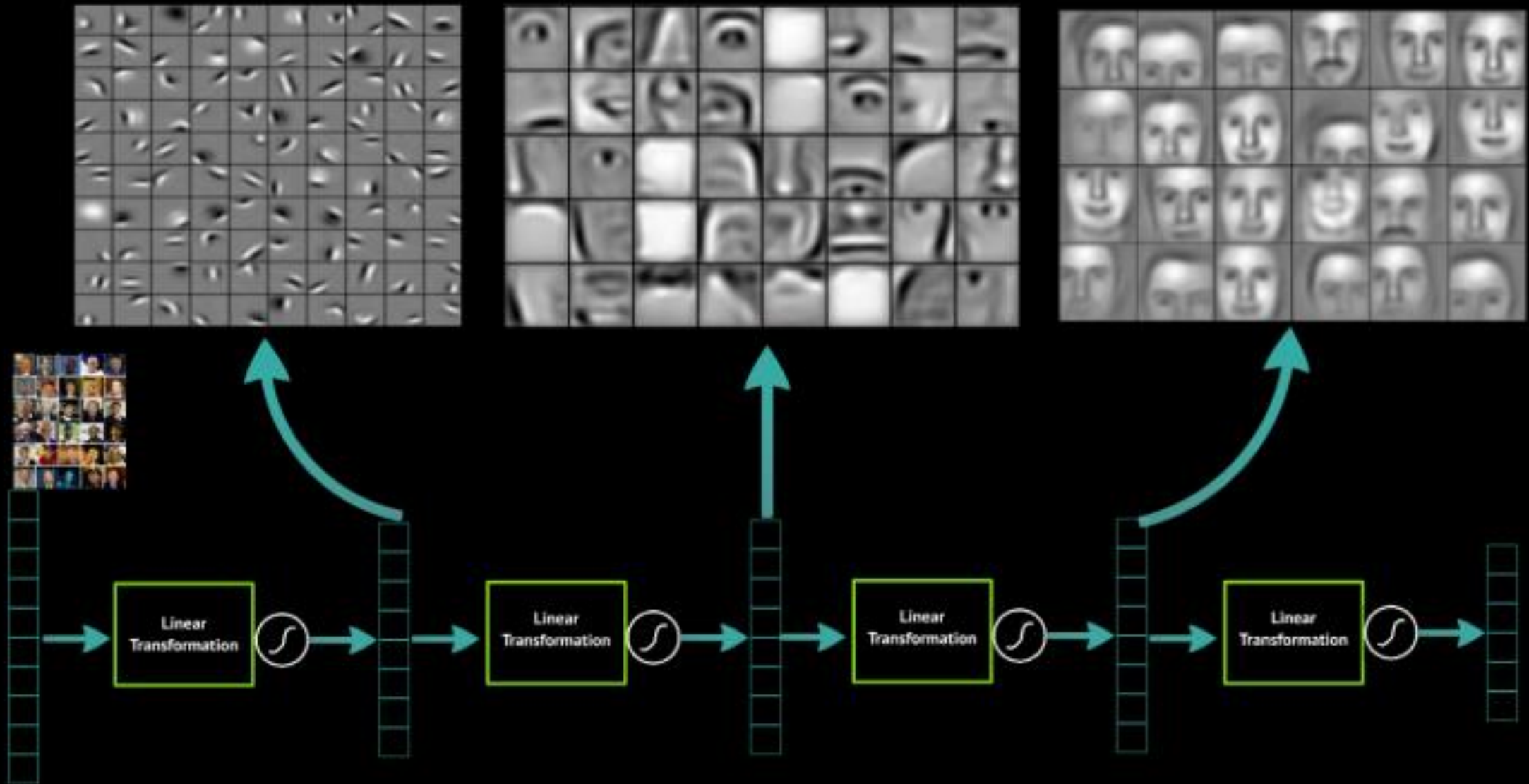
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



# Next convolutions

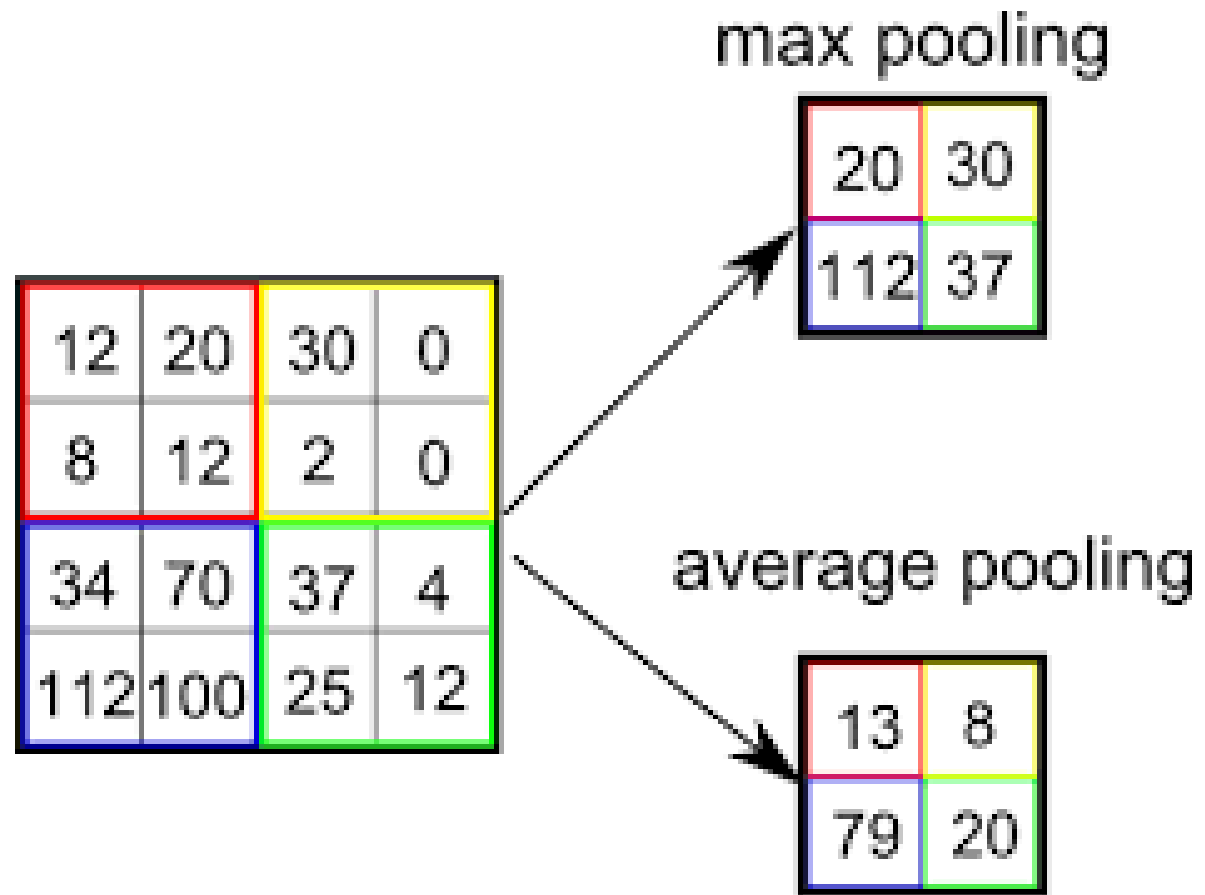


# Convolution layers





# Pooling



# Pooling

Kernel 2x2

Stride 2

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

max pooling

20	30
112	37

average pooling

13	8
79	20

# Padding

0 <sub>2</sub>	0 <sub>0</sub>	0 <sub>1</sub>	0	0	0	0
0 <sub>1</sub>	2 <sub>0</sub>	2 <sub>0</sub>	3	3	3	0
0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>1</sub>	3	0	3	0
0	2	3	0	1	3	0
0	3	3	2	1	2	0
0	3	3	0	2	3	0
0	0	0	0	0	0	0

1	6	5
7	10	9
7	10	8

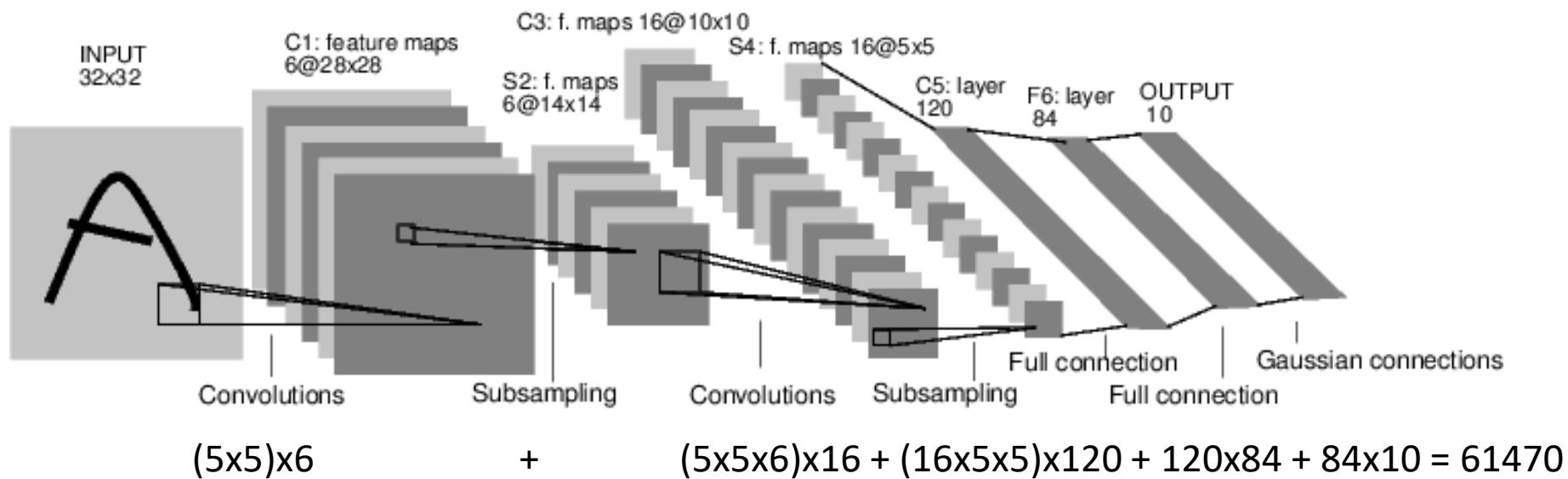
32

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	28 x 28						0	0
0	0							0	0
0	0							0	0
0	0							0	0
0	0							0	0
0	0							0	0
0	0							0	0
0	0							0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

32

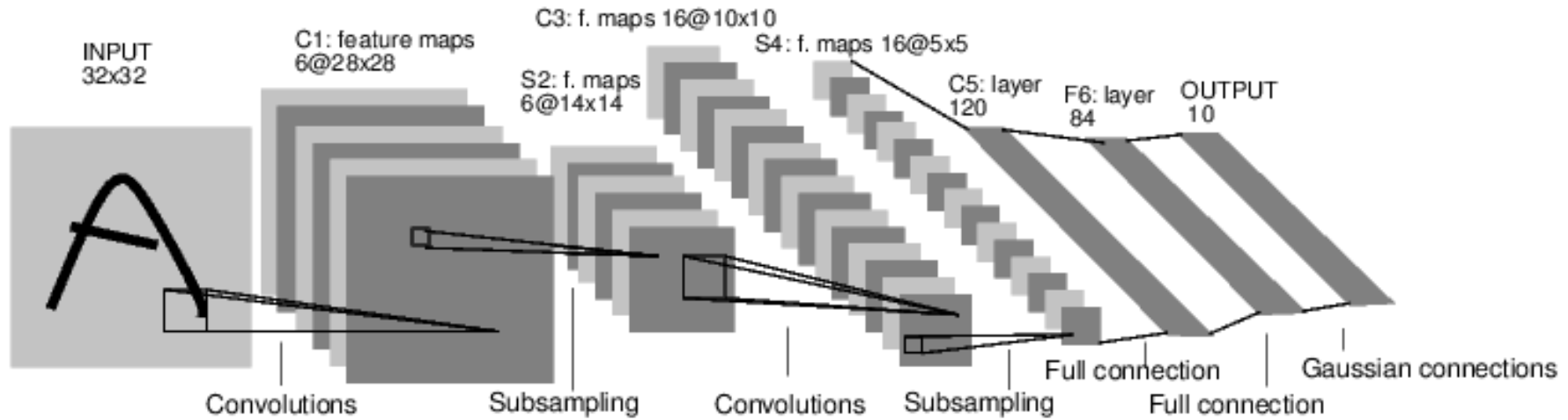
# Deep learning

LeNet-5 (1998)



# Deep learning

LeNet-5 (1998)

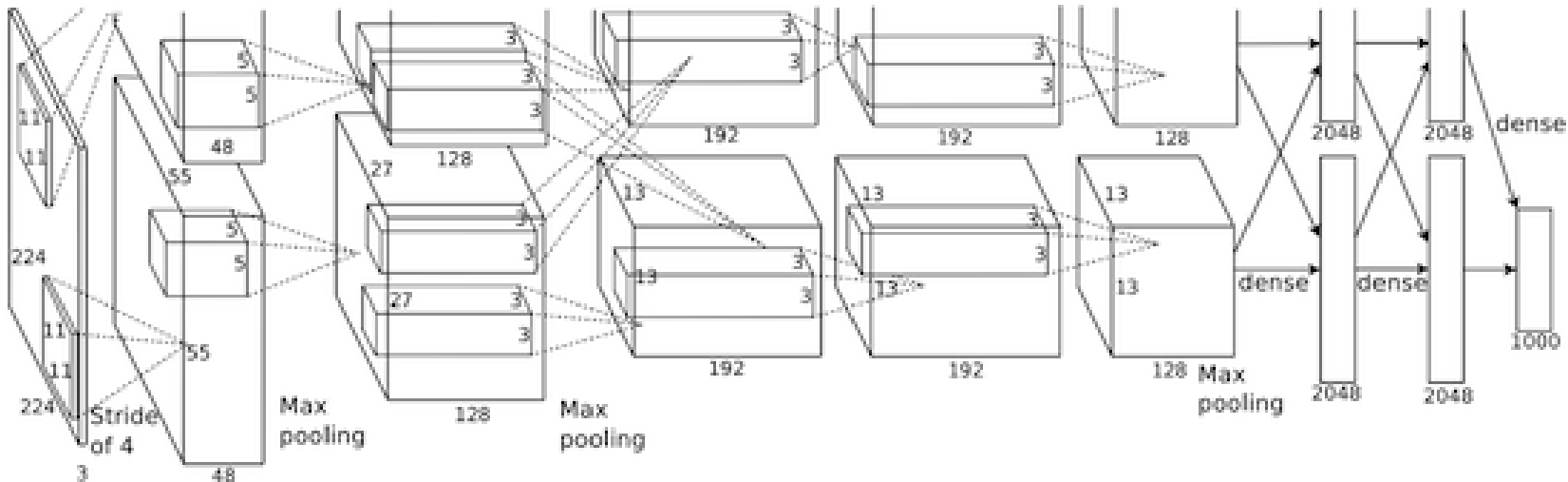


$$(5 \times 5) \times 6$$

+

$$(5 \times 5 \times 6) \times 16 + (16 \times 5 \times 5) \times 120 + 120 \times 84 + 84 \times 10 = 61470$$

AlexNet (2012)

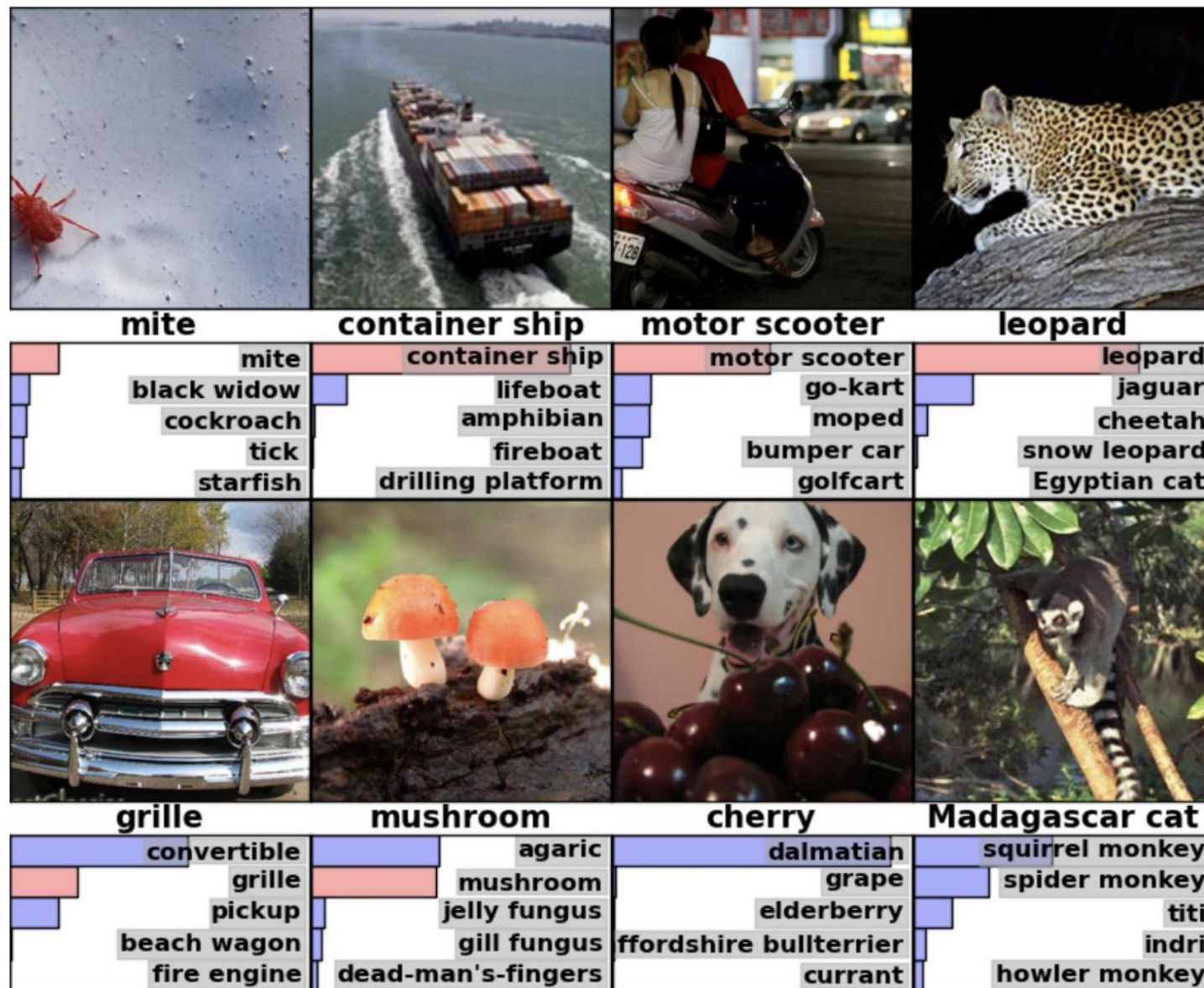


$$(11 \times 11 \times 3) \times 48 + (5 \times 5 \times 48) \times 128 + (3 \times 3 \times 128) \times 192 \times 2 + (3 \times 3 \times 192) \times 192 + (3 \times 3 \times 192) \times 128 + (13 \times 13 \times 128) \times 2048 \times 2 + 2048 \times 2048 \times 2 + 2048 \times 1000 = 100\,207\,632$$

# ImageNet Challenge

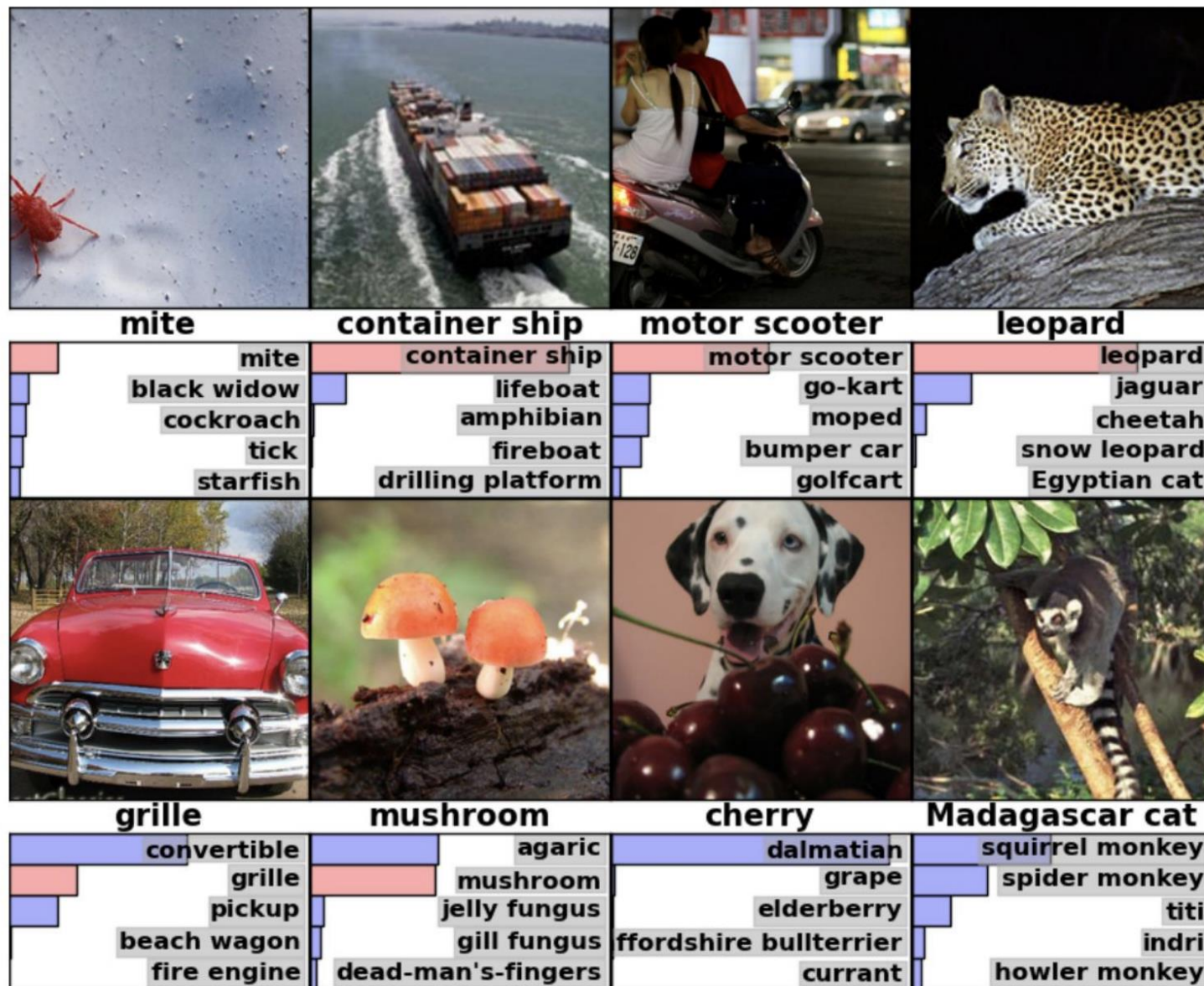
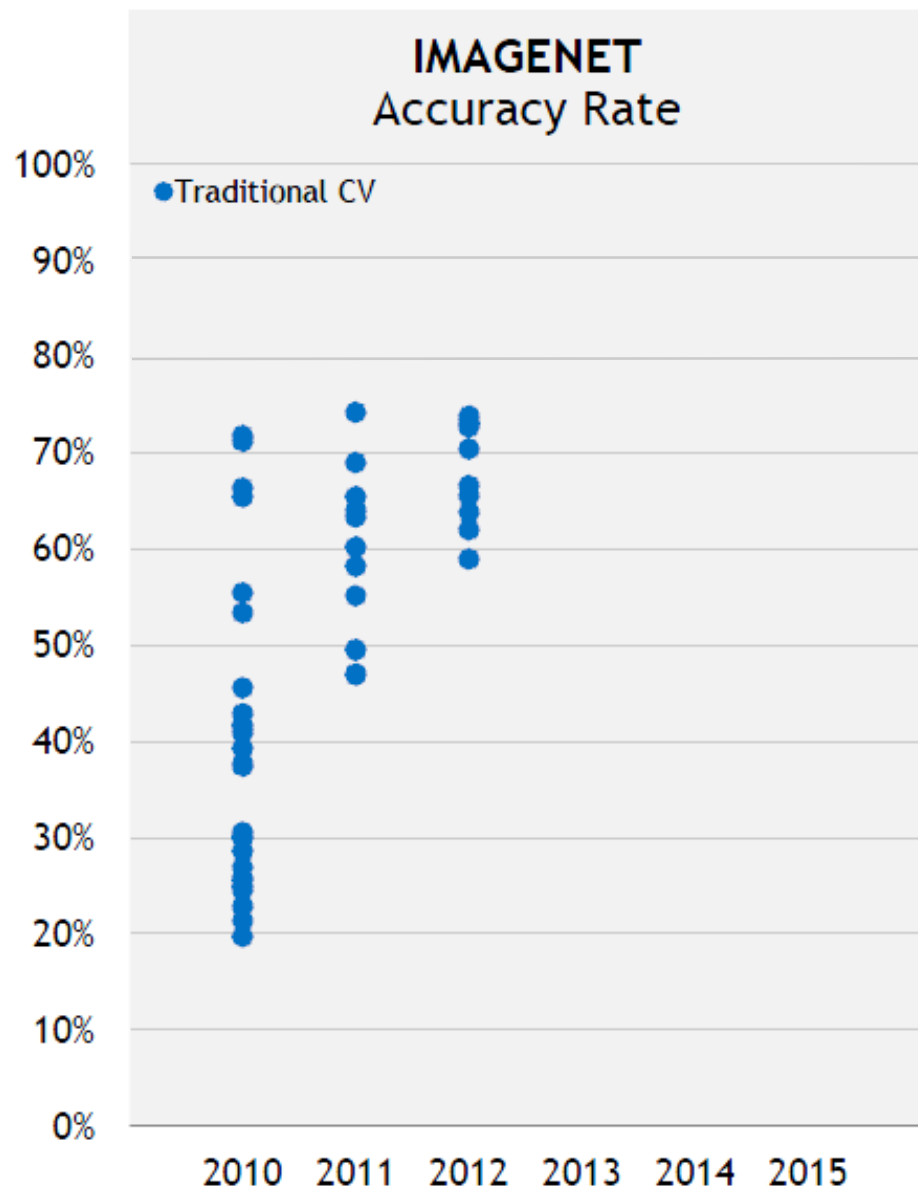


- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

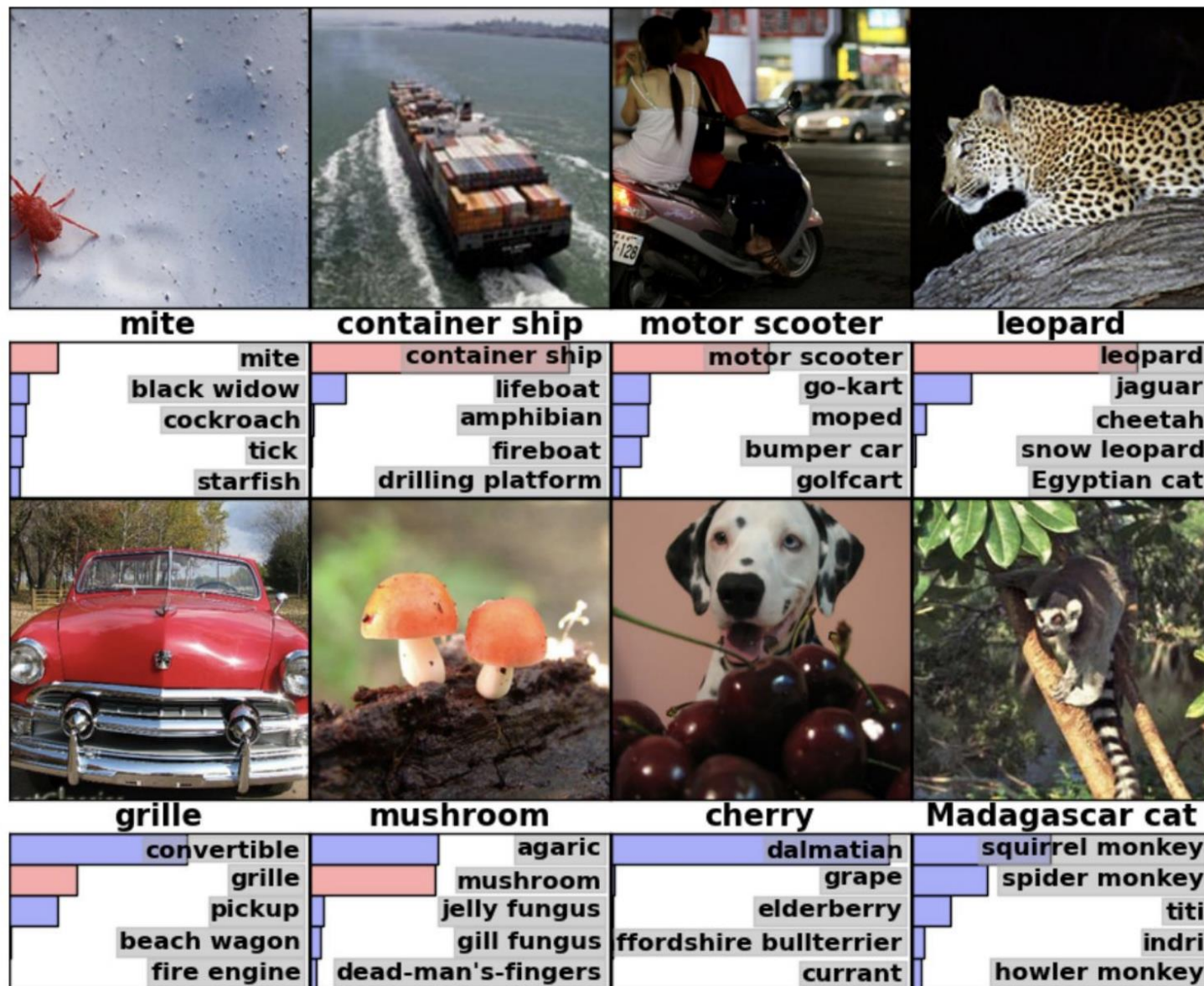
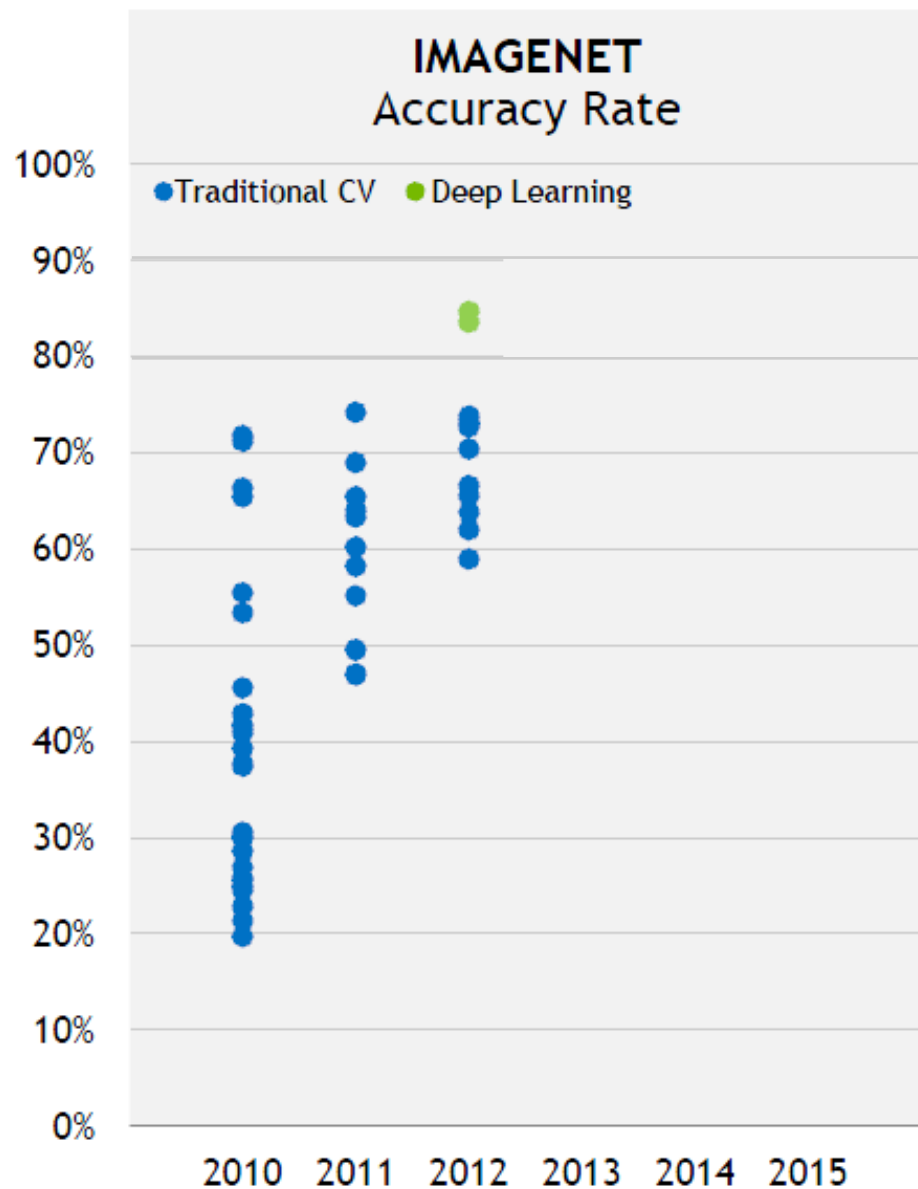




# ImageNet Challenge

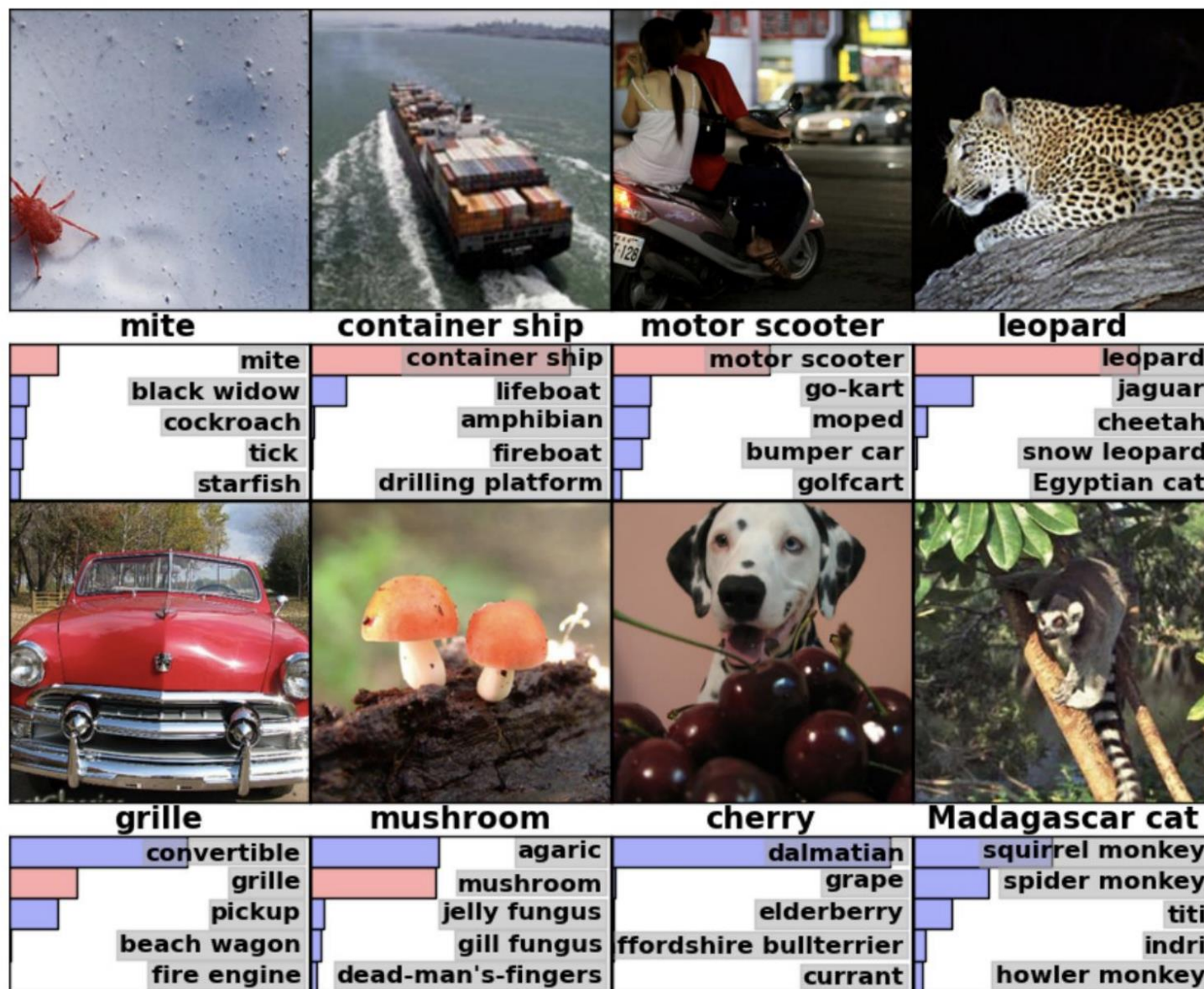
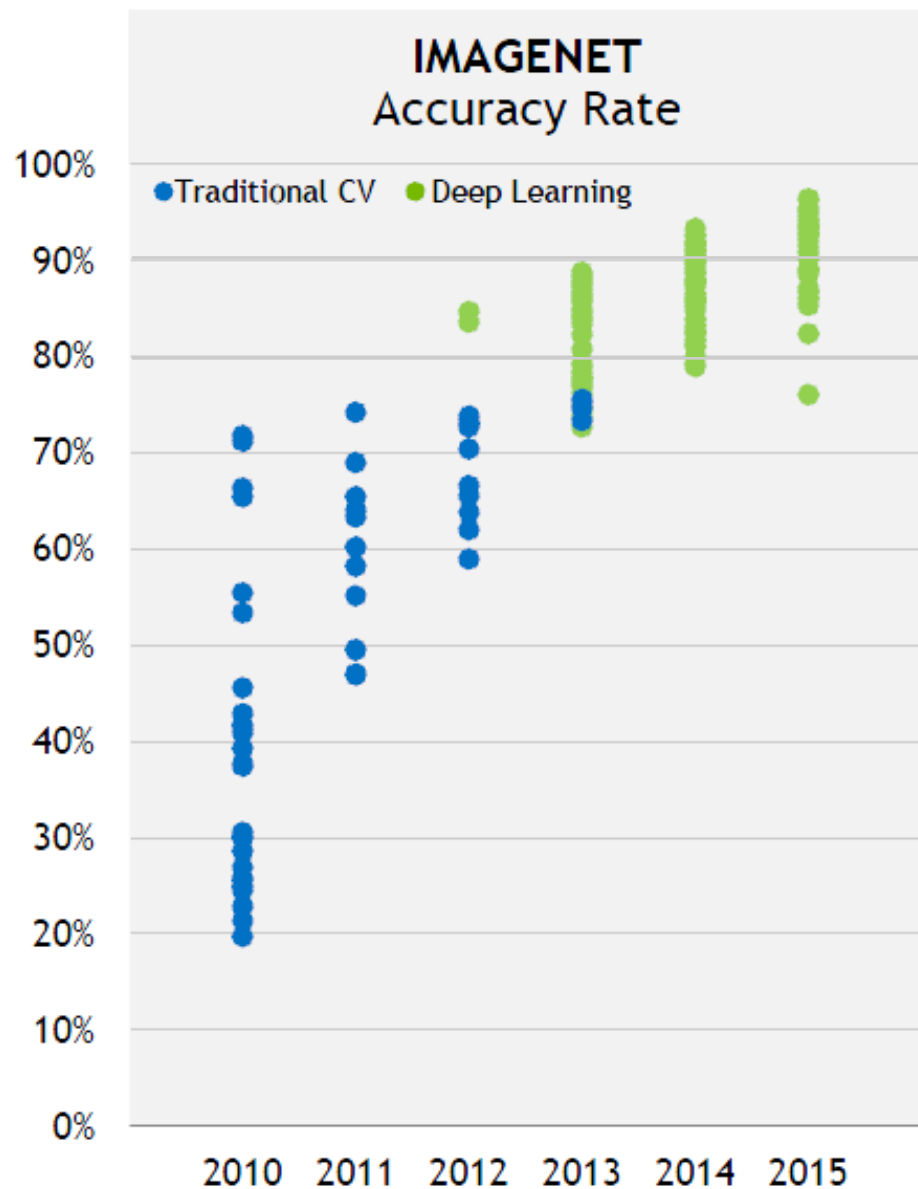


# ImageNet Challenge





# ImageNet Challenge



# AlexNet

- Scale all images to 256x256, then take random 224x224 patches and mirror them.
- Subtract average pixel value from each pixel.
- ReLU ( $f(x) = \max(0, x)$ )
- Dropout 0.5
- Batch size 128, SGD with momentum (0.9), L2 weight decay ( $\lambda = 0.0005$ )

# AlexNet, first convolution



“The kernels on GPU 1 are largely color-agnostic, while the kernels on GPU 2 are largely color-specific. This kind of specialization occurs during every run and is independent of any particular random weight initialization (modulo a renumbering of the GPUs)”

# Visualization of convolutions

## Deep Visualization Toolbox

[yosinski.com/deepvis](http://yosinski.com/deepvis)

#deepvis



Jason Yosinski



Jeff Clune



Anh Nguyen



Thomas Fuchs



Hod Lipson





# Very Deep Convolutional Networks (VGGNet), 2014

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 ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	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 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

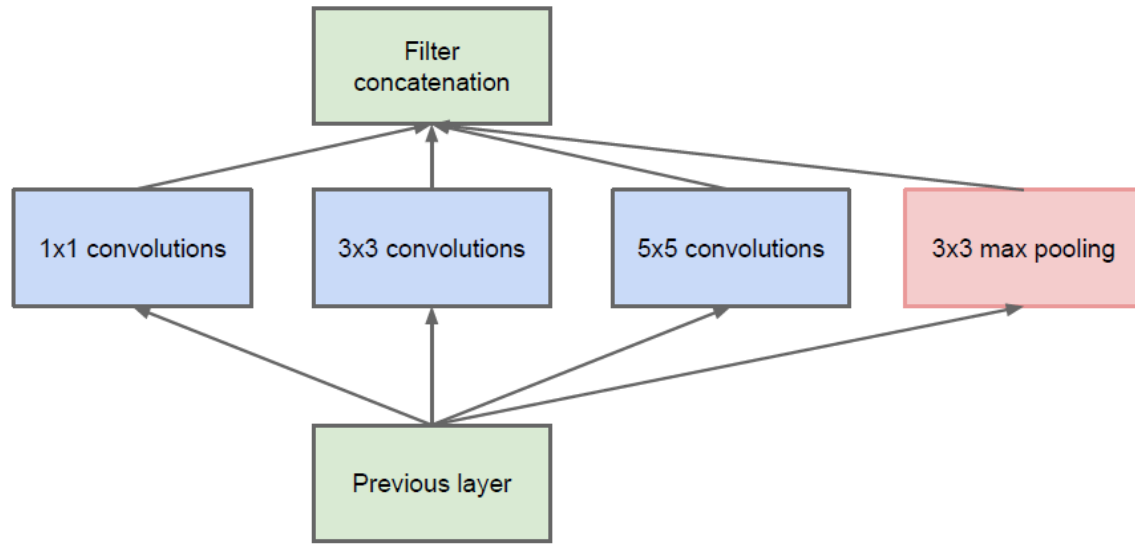
# Very Deep Convolutional Networks (VGGNet), 2014

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 ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	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 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

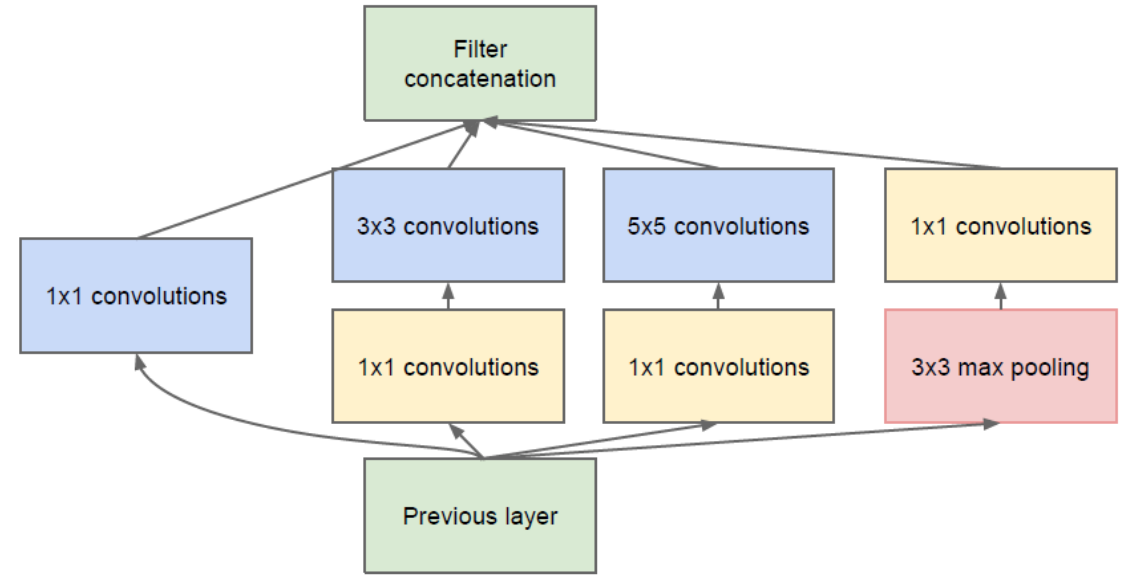
ImageNet Challeng result:

VGGNet (D+E) – 7.3% (2<sup>nd</sup>)

# Inception (Google 2014)



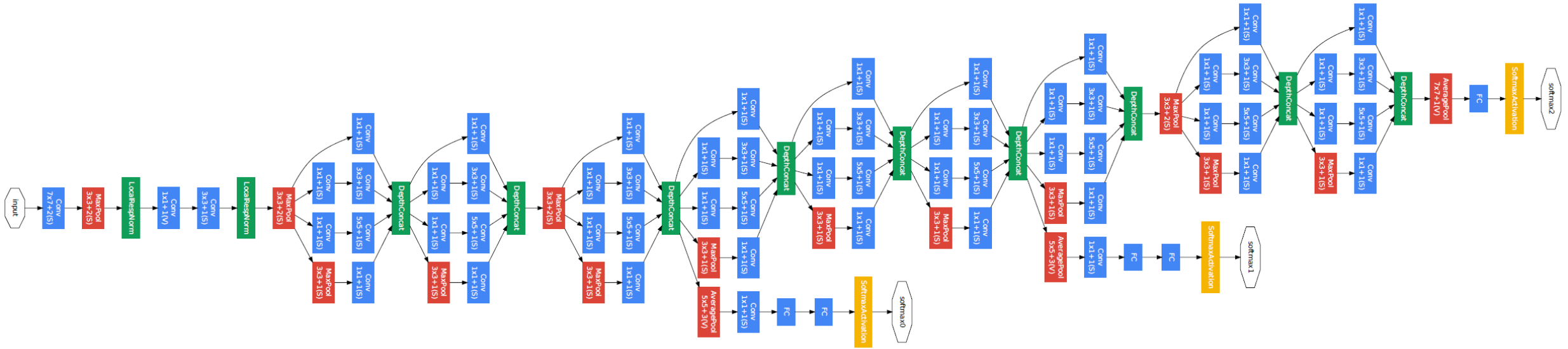
(a) Inception module, naïve version



(b) Inception module with dimension reductions

Using 1x1 convolutions to decrease a number of maps in layers.

# GoogLeNet



Convolutions

MaxPooling

Concatenation

FC

AveragePooling

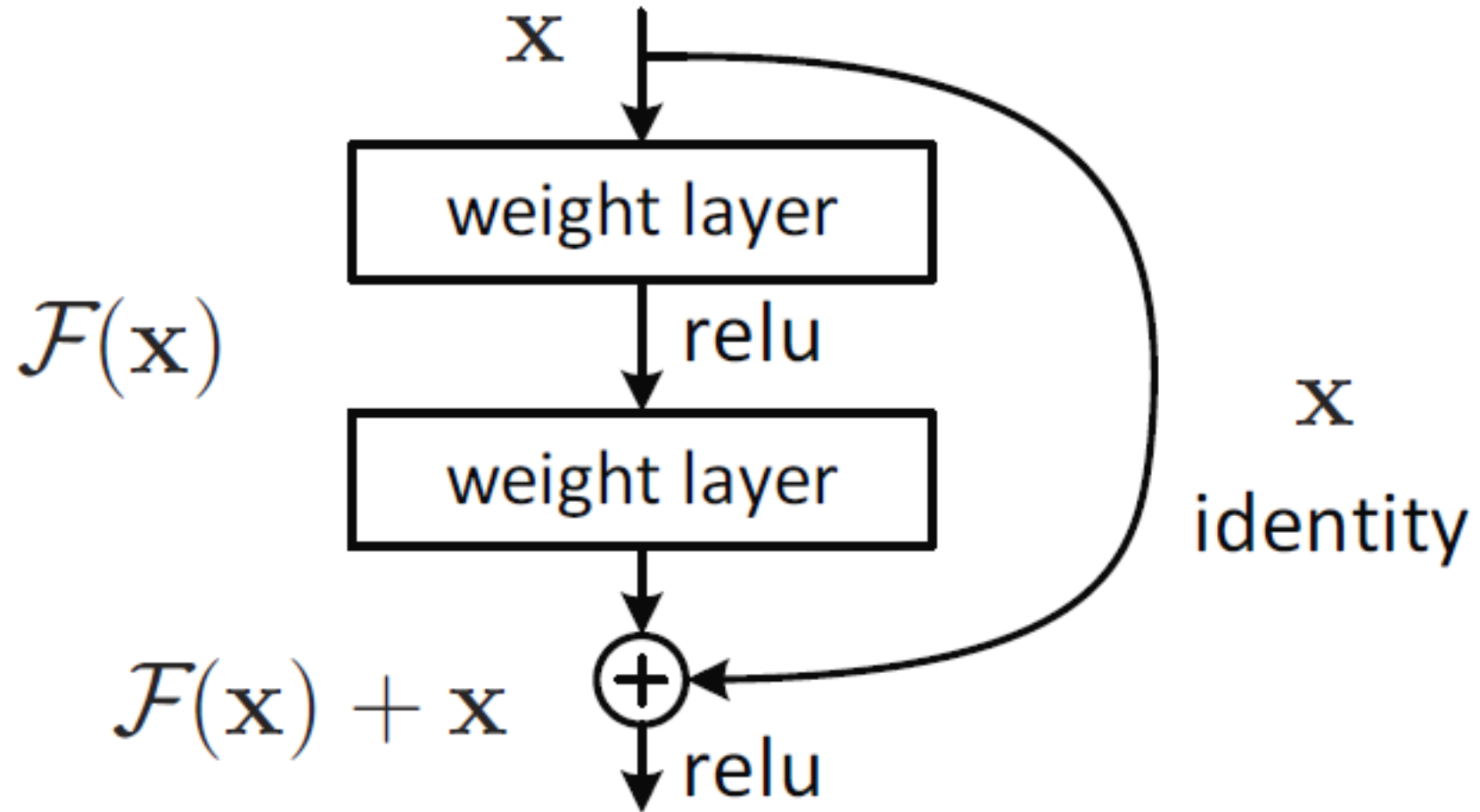
Softmax



# GoogLeNet (2014)

Team	Year	Place	Error (top-5)
SuperVision	2012	1 <sup>st</sup>	16.4%
Clarifai	2013	1 <sup>st</sup>	11.7%
MSRA	2014	3 <sup>rd</sup>	7.35%
VGG	2014	2 <sup>nd</sup>	7.32%
GoogLeNet	2014	1 <sup>st</sup>	6.67%

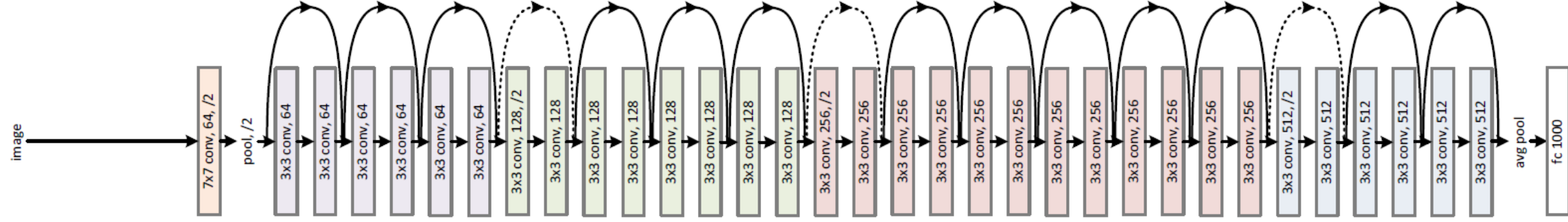
# Residual Learning



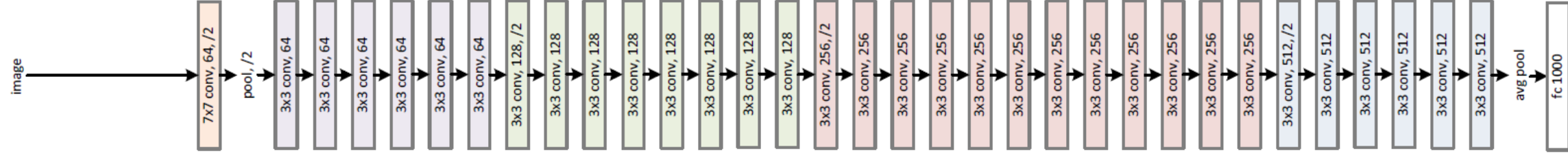
Now called "skip-connections".

# ResNet (2015)

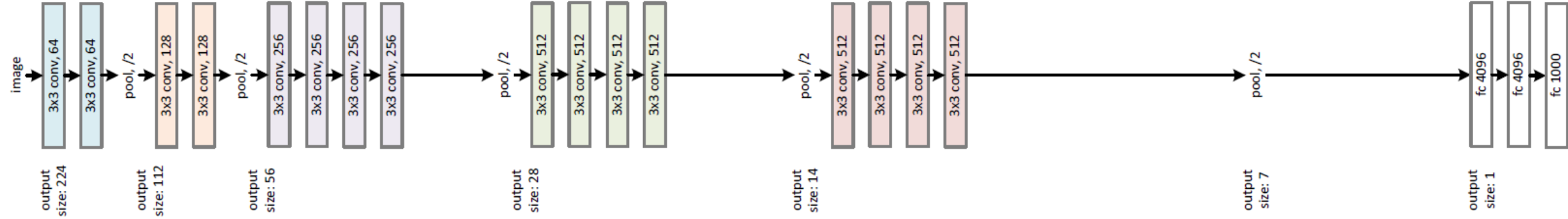
34-layer residual



34-layer plain



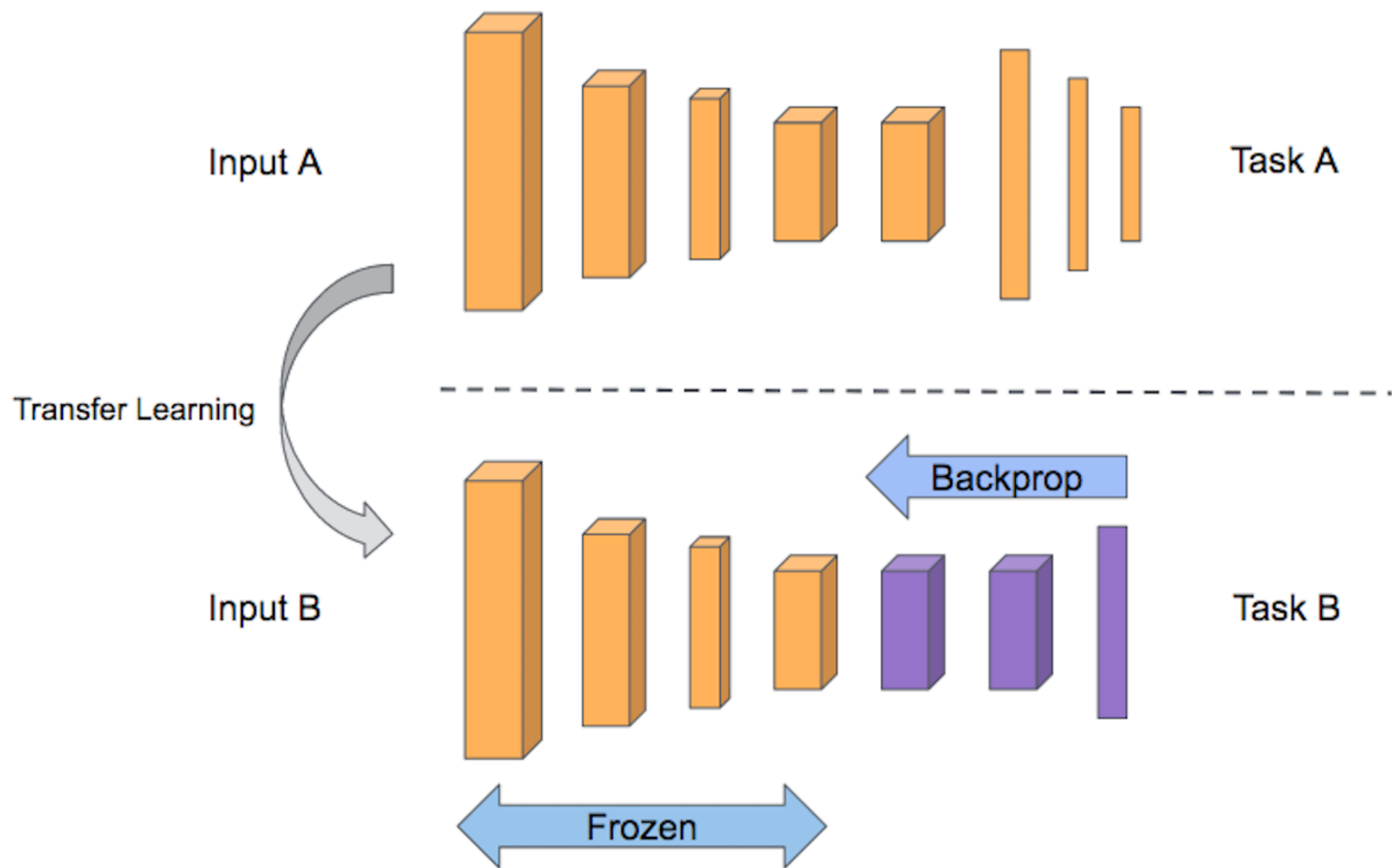
VGG-19

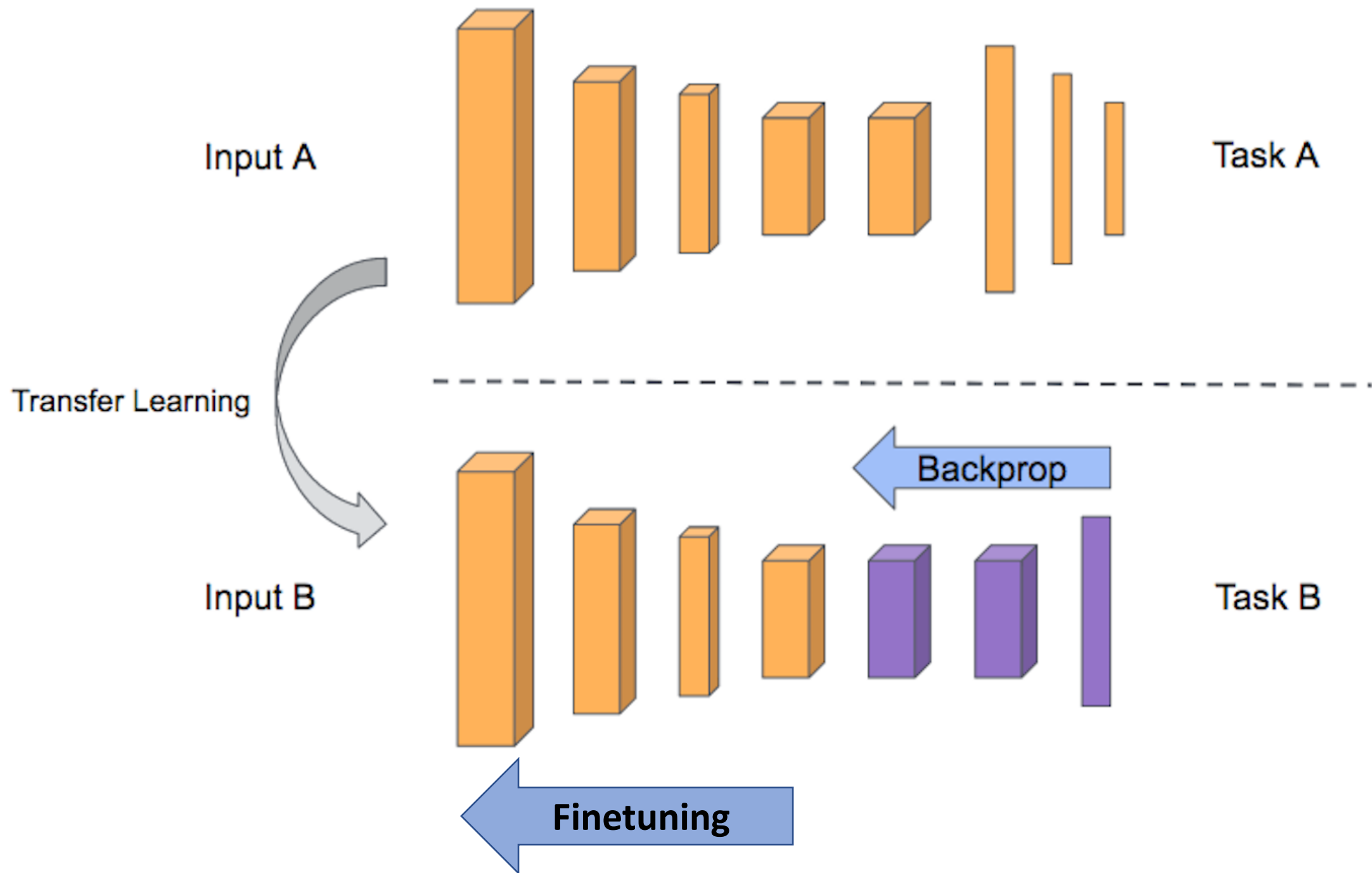


# ResNet

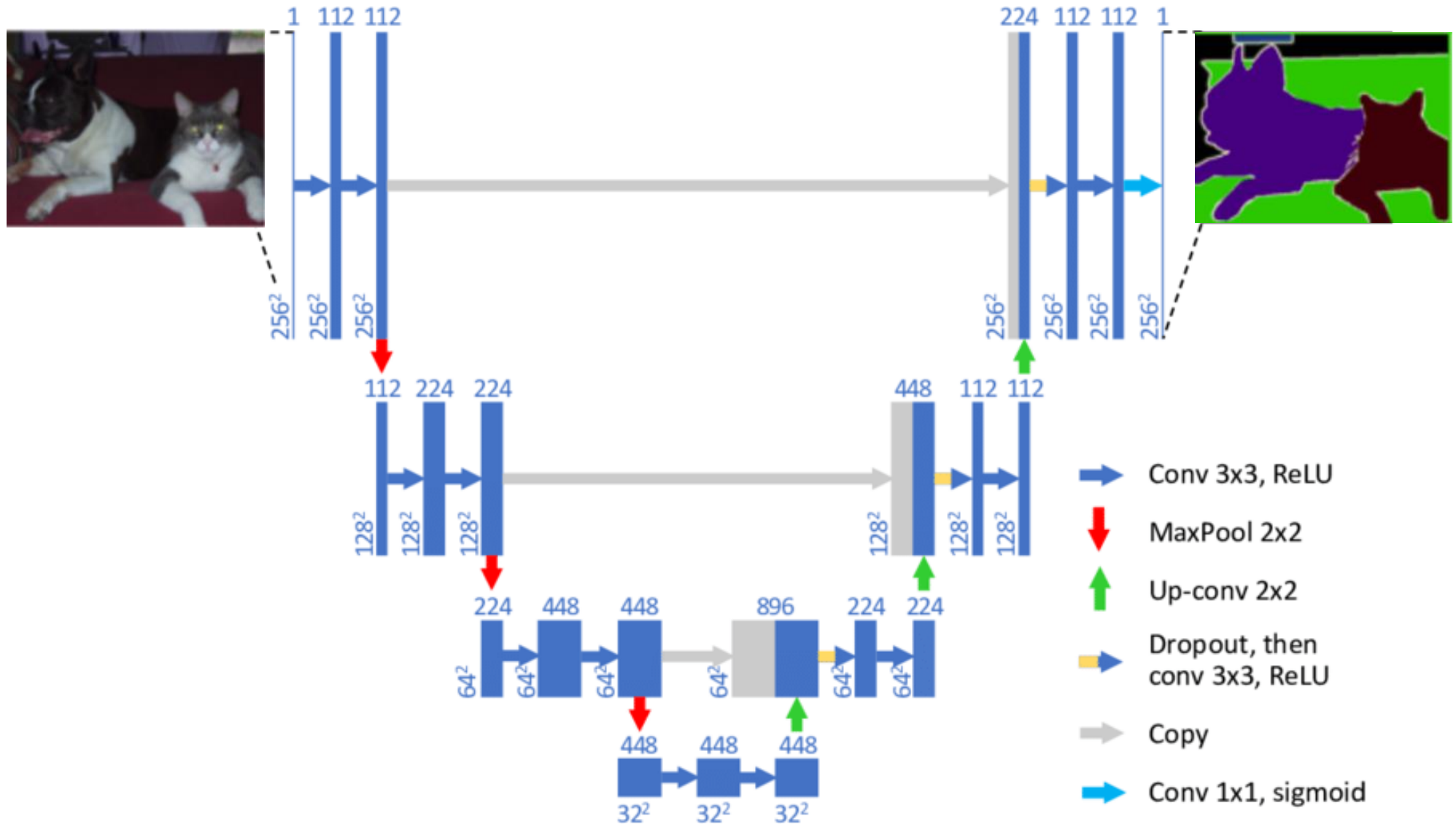
method	top-5 err. ( <b>test</b> )
VGG [41] (ILSVRC'14)	7.32
GoogLeNet [44] (ILSVRC'14)	6.66
VGG [41] (v5)	6.8
PReLU-net [13]	4.94
BN-inception [16]	4.82
<b>ResNet (ILSVRC'15)</b>	<b>3.57</b>

Transfer learning and finetuning



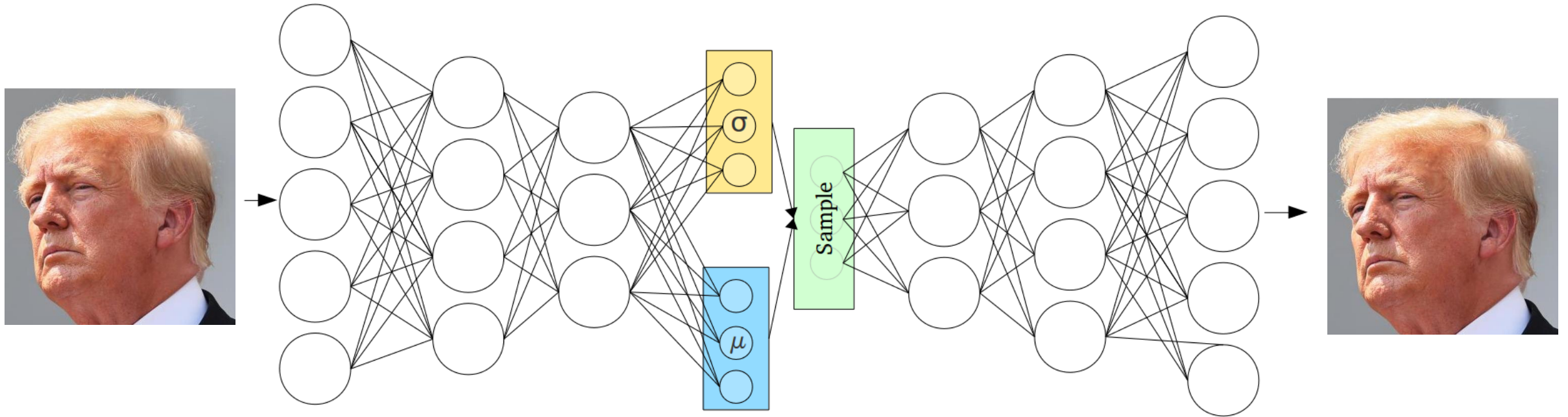


# Image segmentation and U-Net





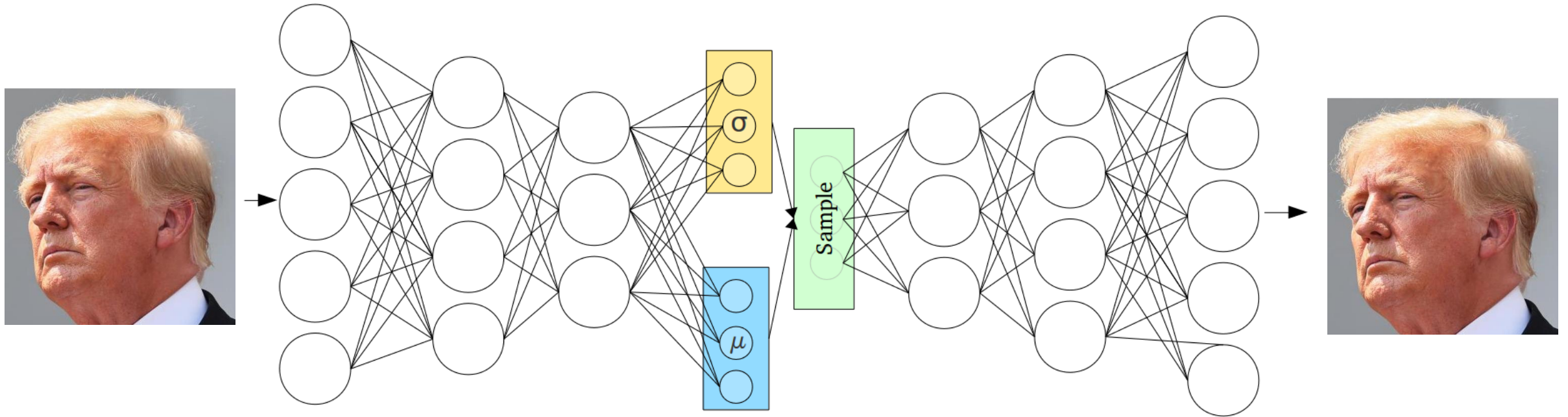
# Variational Autoencoder (VAE)



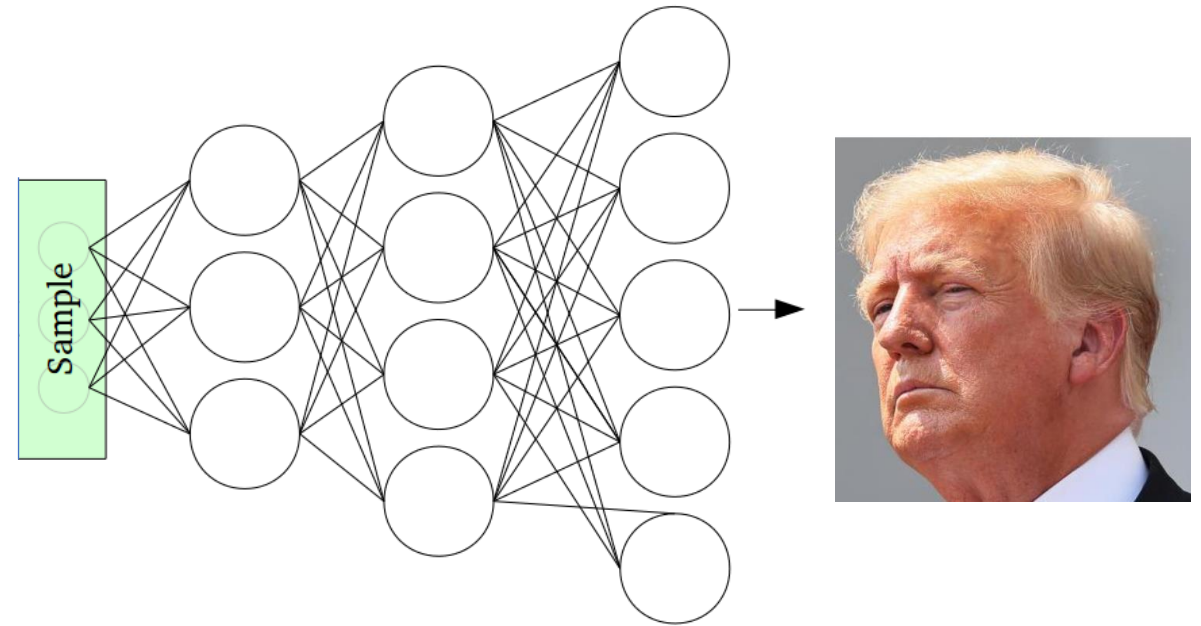
# Variational Autoencoder (VAE)



# Variational Autoencoder (VAE)



# Generative Adversarial Network (GAN)



# Generative Adversarial Network (GAN)

