

Statistical Methods in AI (CSE 471)

Lecture 16: CNN

Vineet Gandhi
Centre for Visual Information Technology (CVIT)



Difficult problem

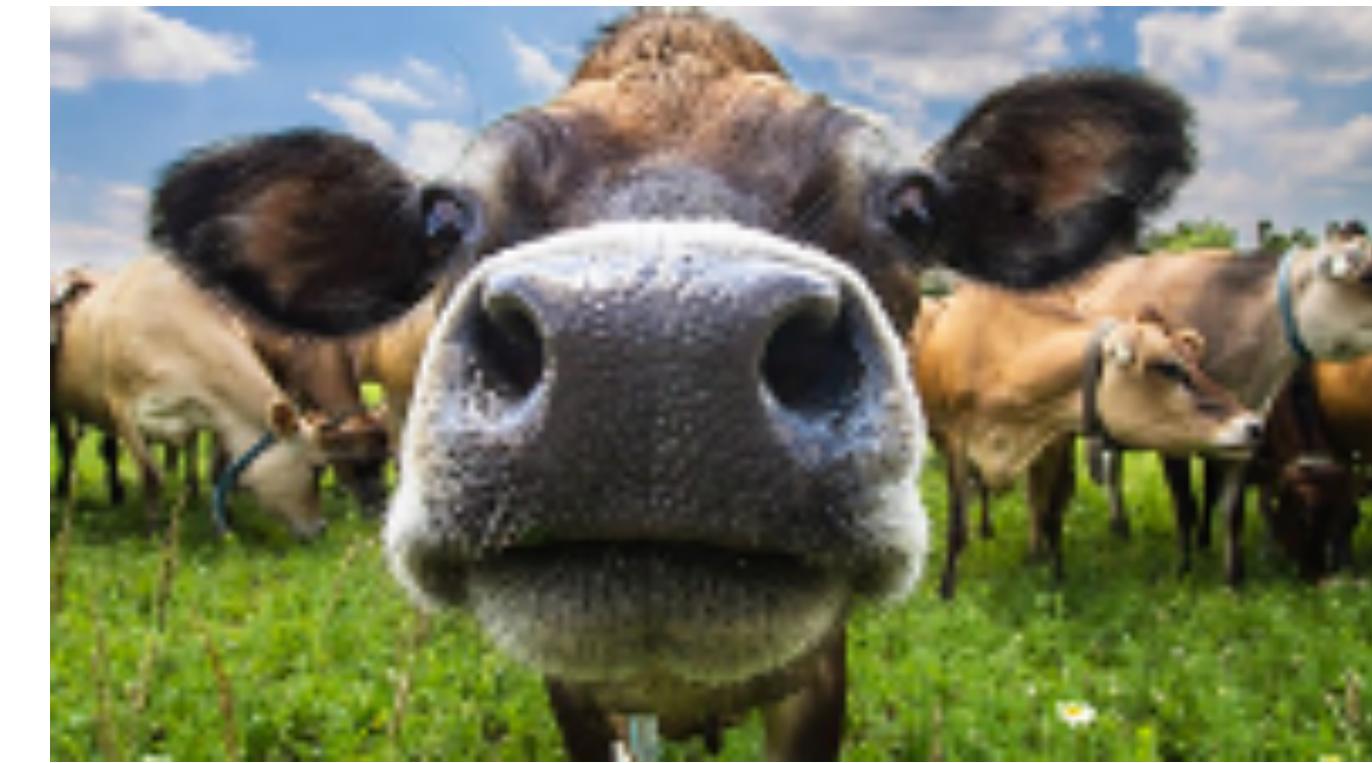
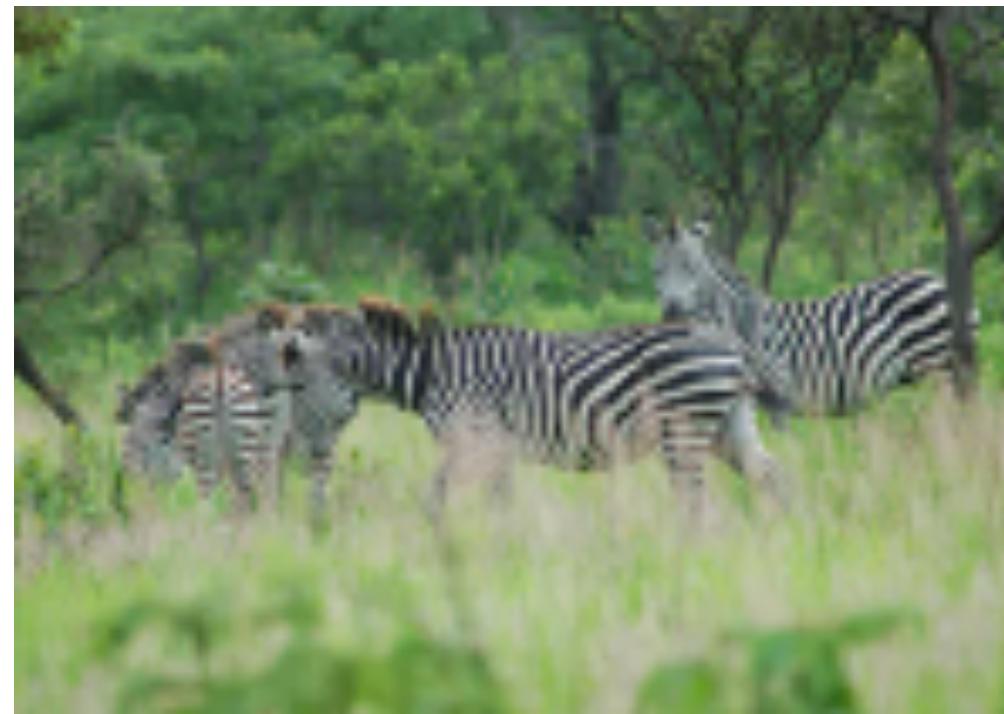
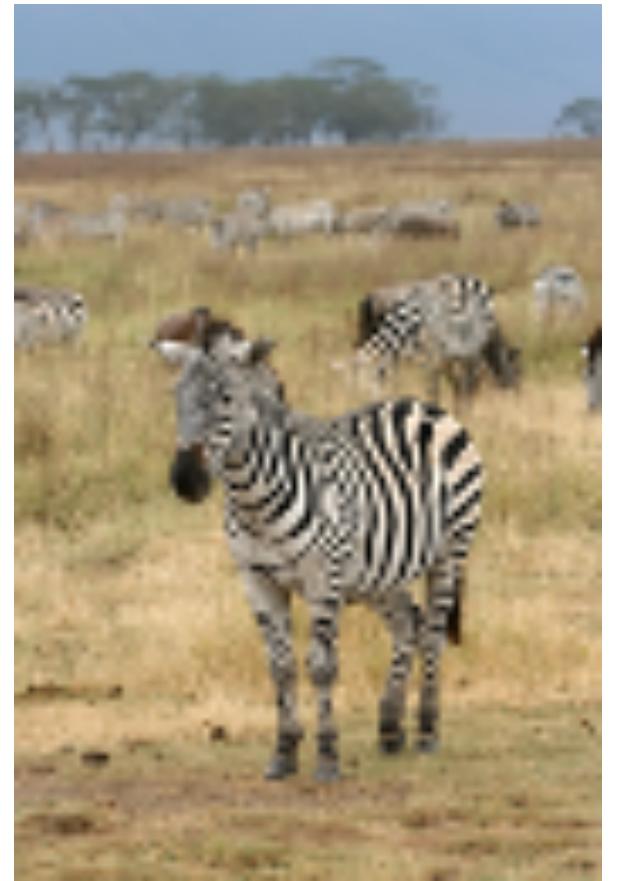
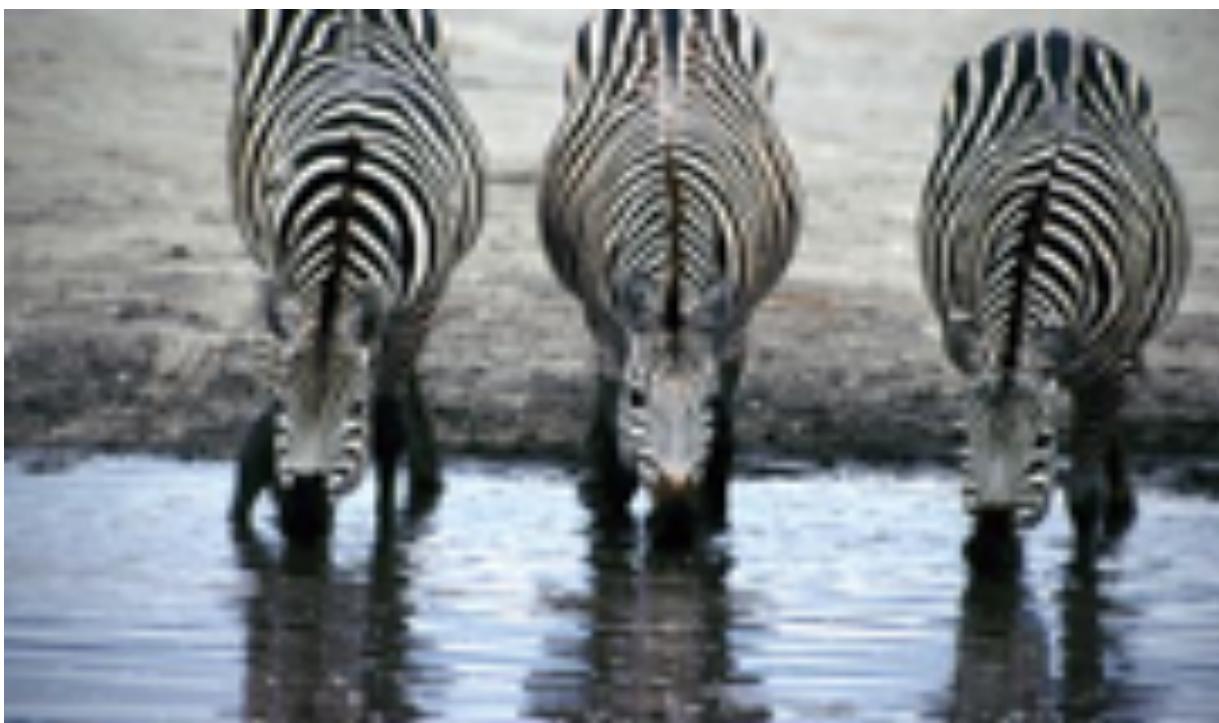


194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	67	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	60	99	63	45	60	63	50	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	66	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	78	76	73	59	75	69	50

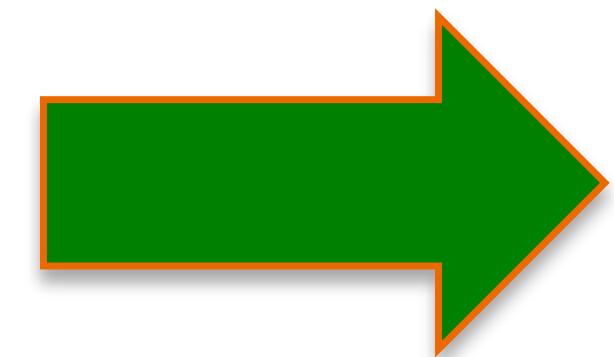
Cow vs Zebra



Cow vs Zebra



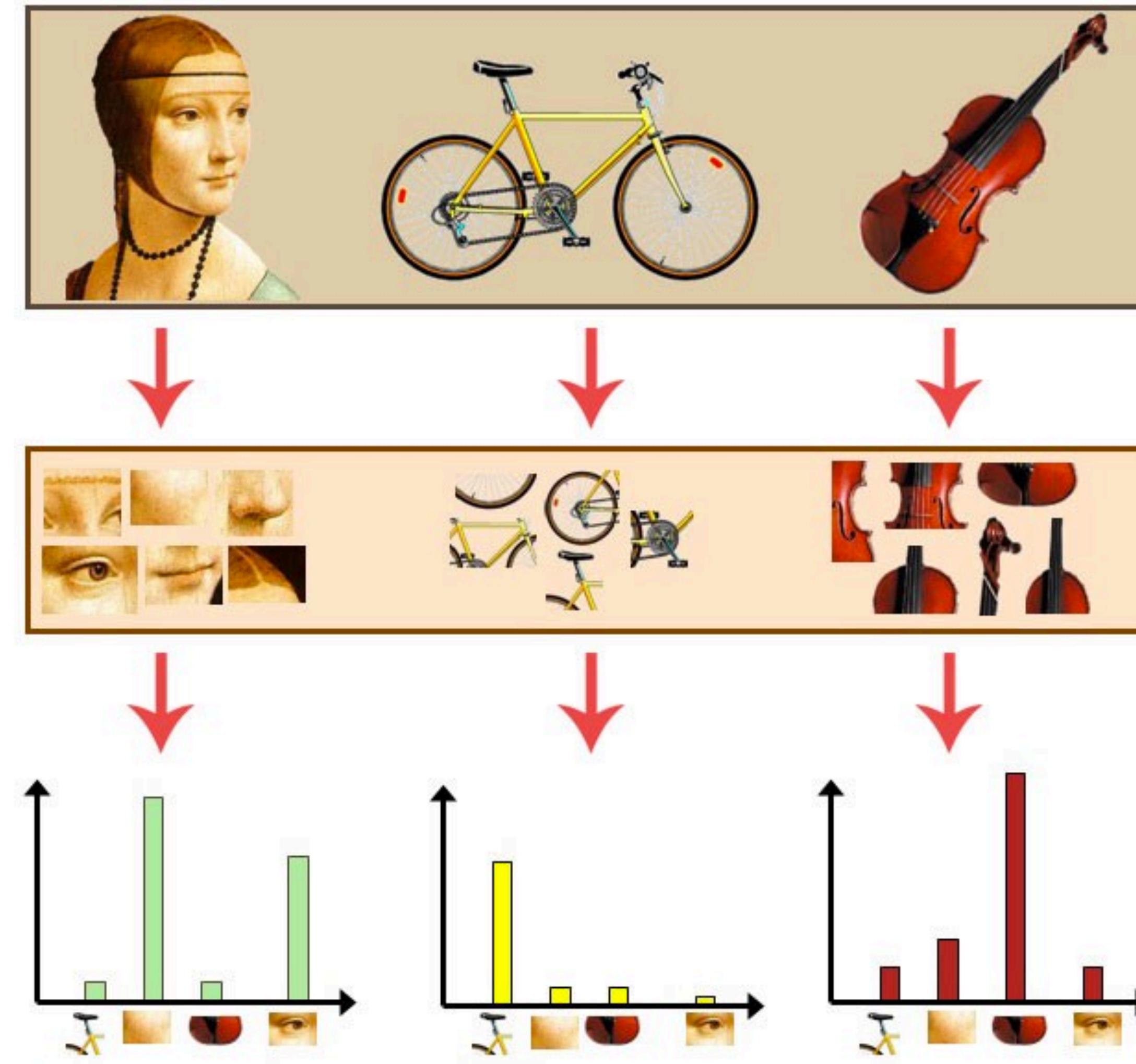
Raw image?



Other features

MIN RED
MAX RED
MEAN RED
MIN GREEN
MAX GREEN
MEAN GREEN
MIN BLUE
MAX BLUE
MEAN BLUE

Higher level representations: bag of words



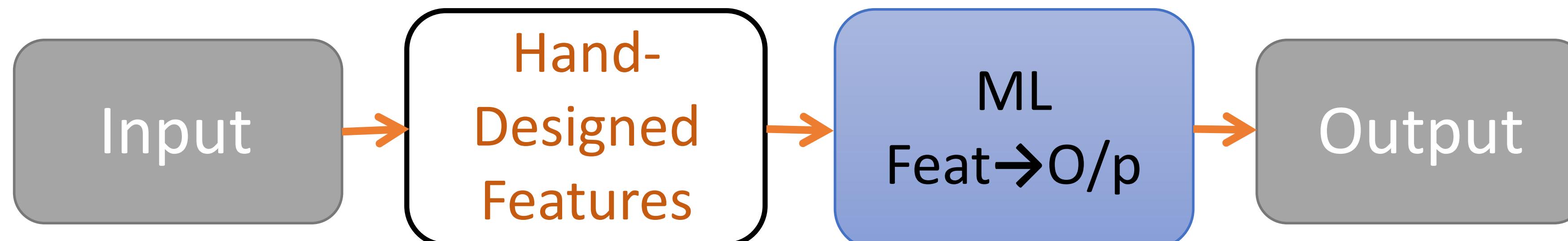
CNN

Expert Systems

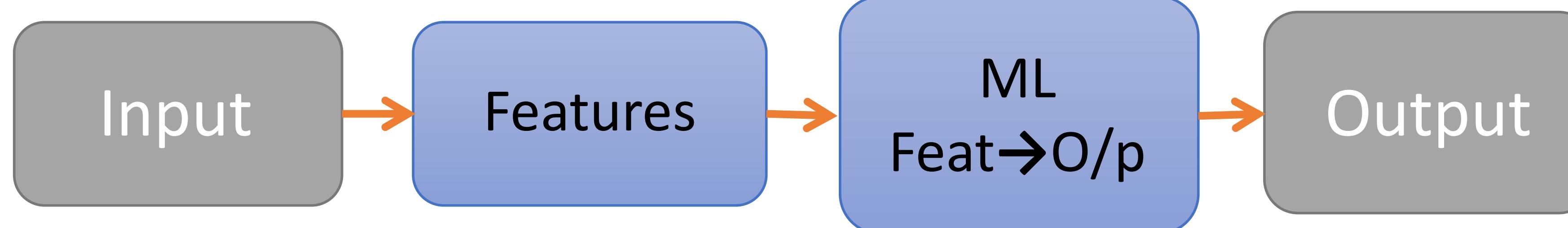


Y. Bengio et al,
“Deep
Learning”,
MIT Press, 2015

Classic
ML



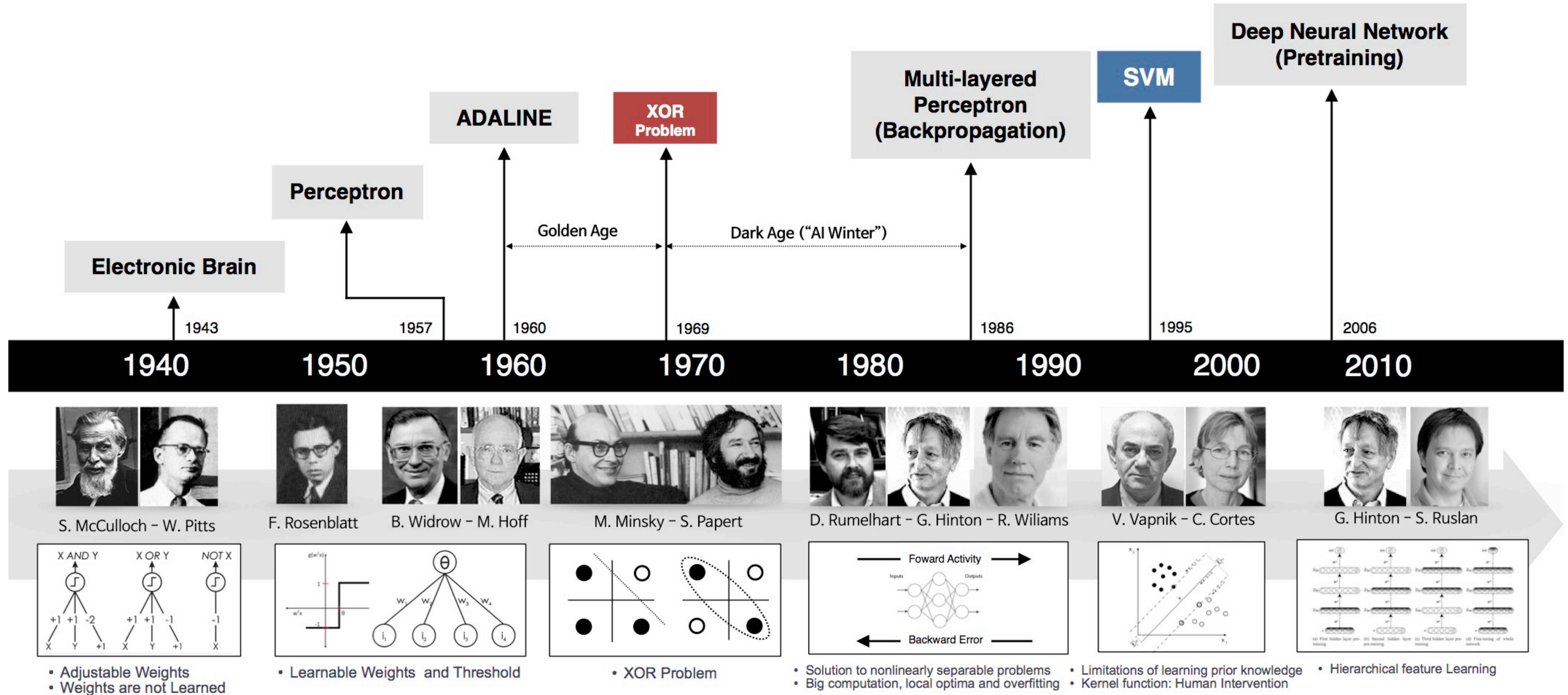
Repr'n
Learning



Deep
Learning

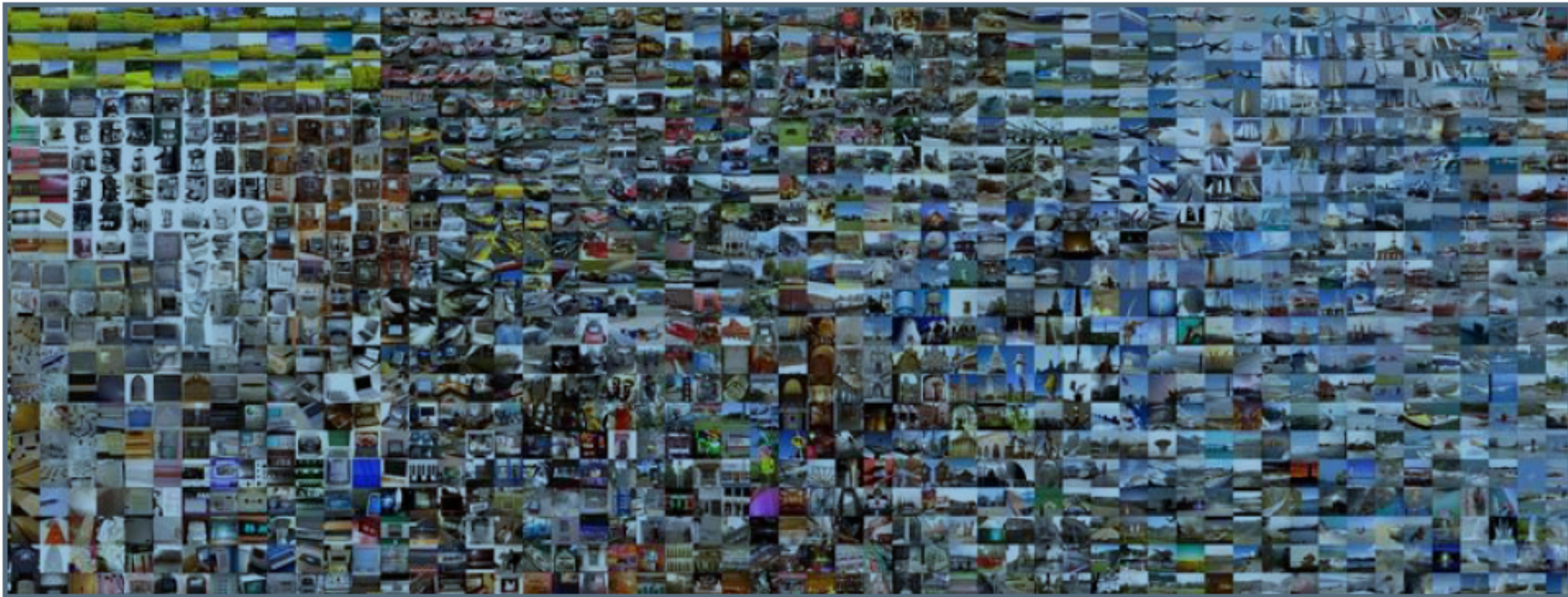


History

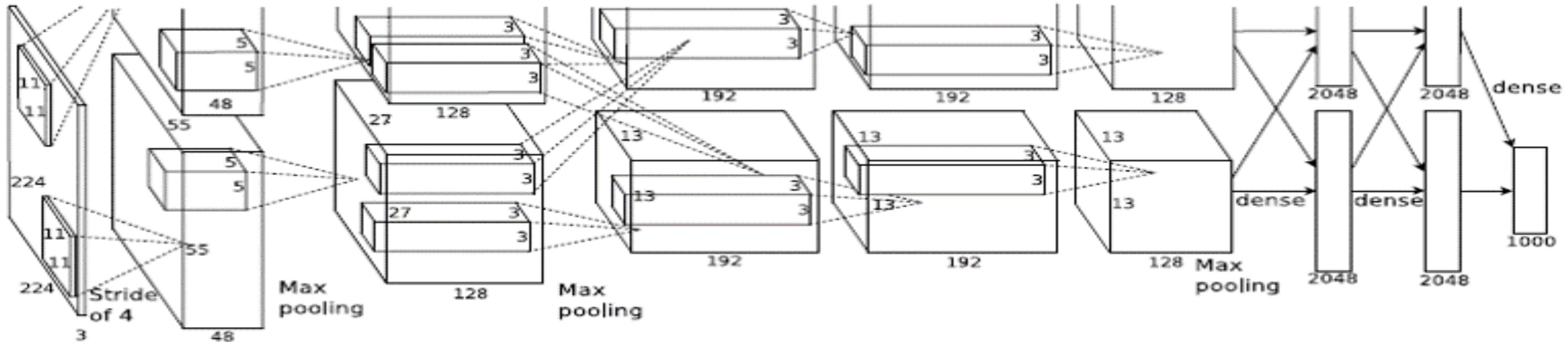


ImageNet

IMAGENET



CNN



ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

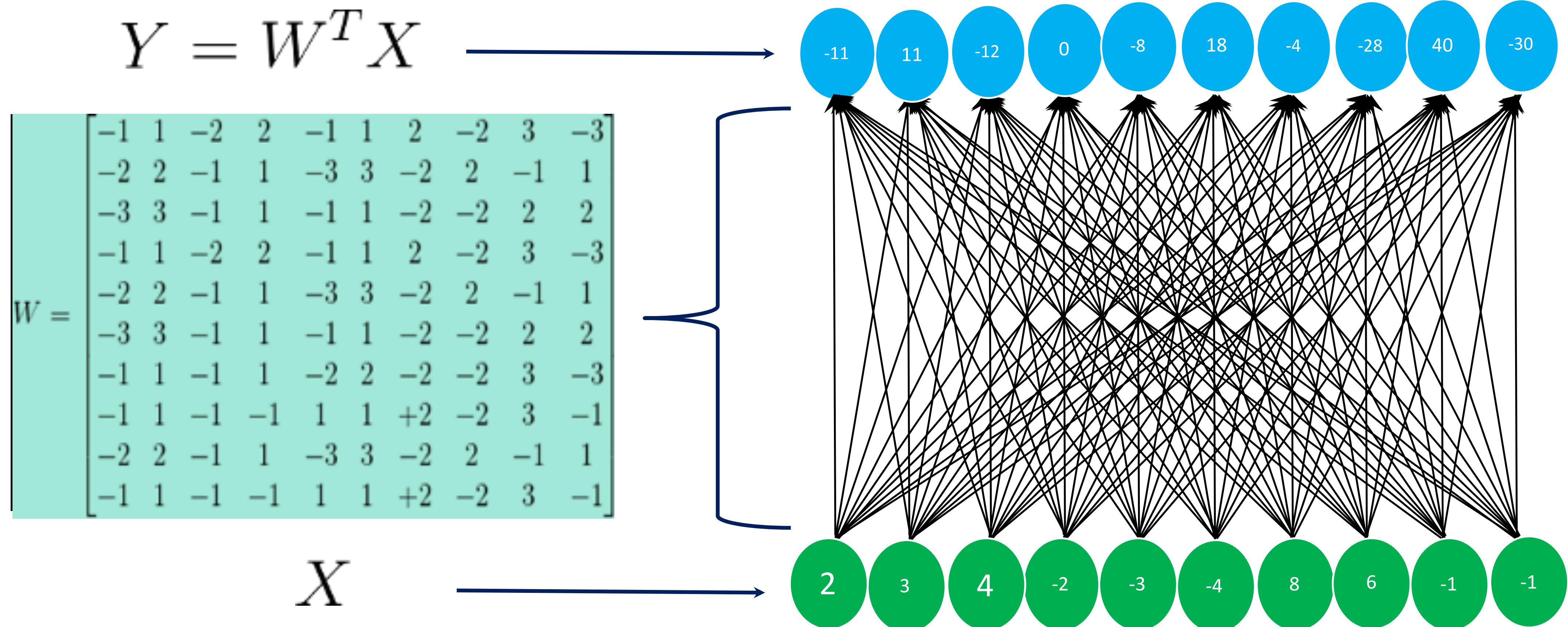
Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

ImageNet Classification Task:

Previous Best : ~25% (CVPR-2011)
AlexNet : ~15 % (NIPS-2012)

Dense connections



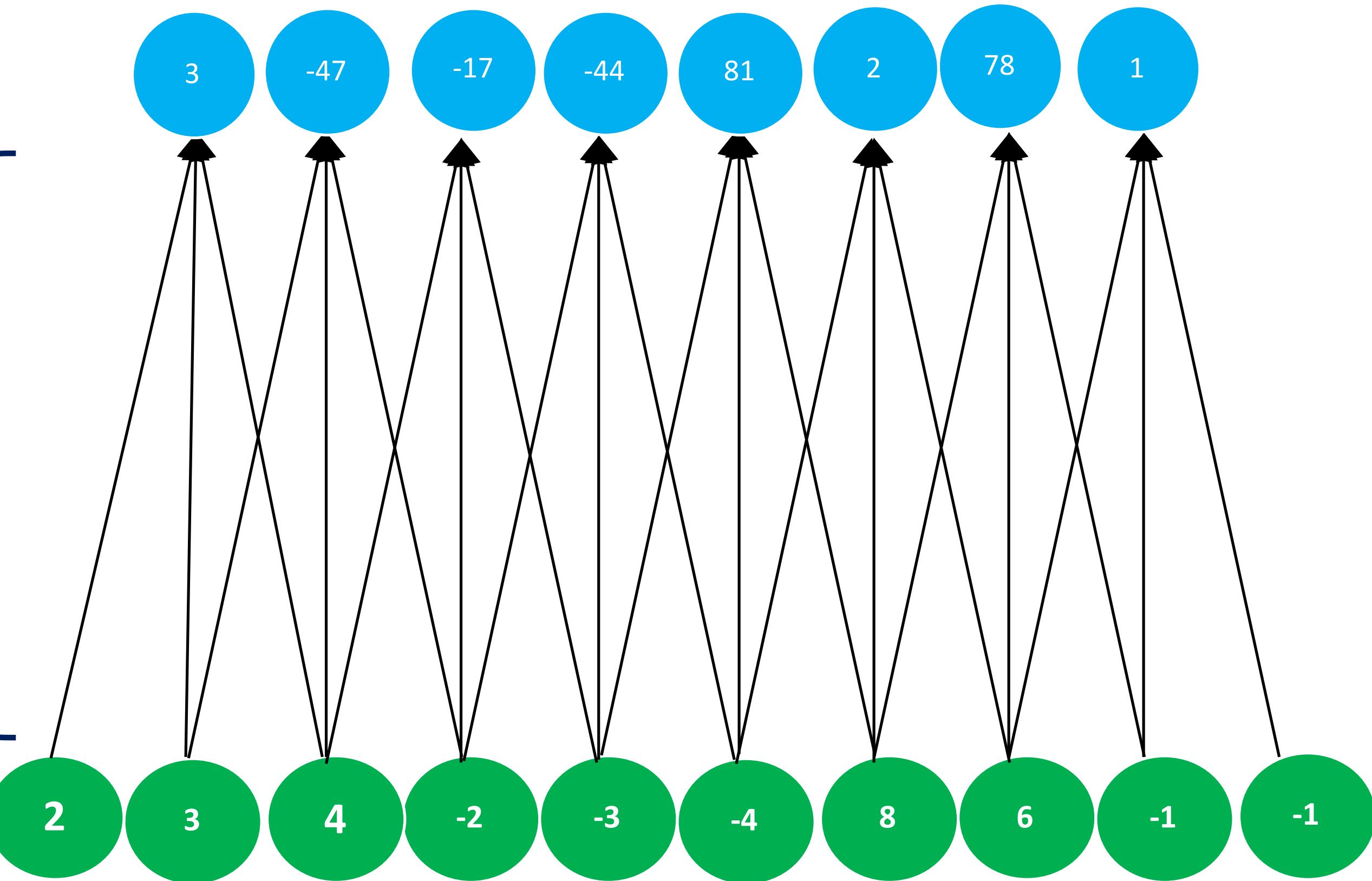
Dense connections

$$Y = W^T X$$

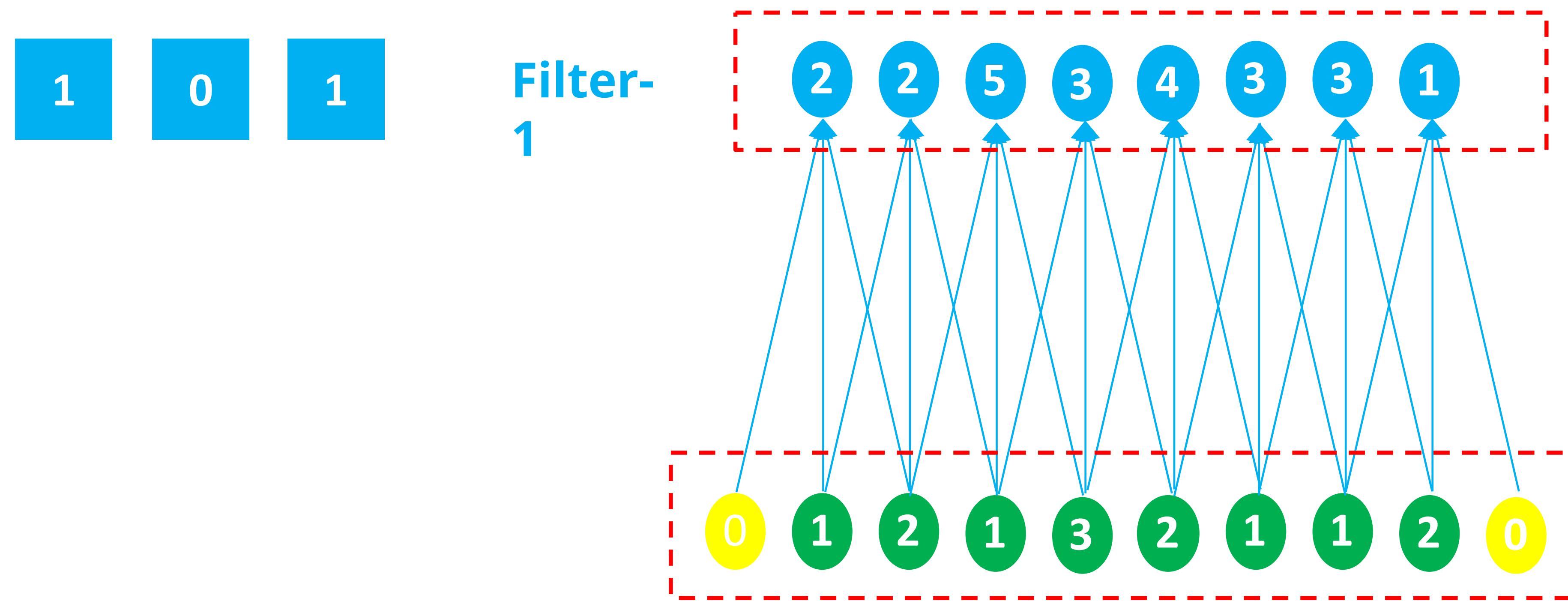
$$W =$$

$$\begin{bmatrix} 6 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & -7 & 0 & 0 & 0 & 0 & 0 & 0 \\ -3 & -4 & -9 & 0 & 0 & 0 & 0 & 0 \\ 0 & 5 & 4 & 4 & 0 & 0 & 0 & 0 \\ 0 & 0 & -9 & 0 & 5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 9 & -8 & 5 & 0 & 0 \\ 0 & 0 & 0 & 0 & 8 & -4 & 9 & 0 \\ 0 & 0 & 0 & 0 & 0 & 9 & 0 & -2 \\ 0 & 0 & 0 & 0 & 0 & 0 & -6 & -7 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -6 \end{bmatrix}$$

$$X$$

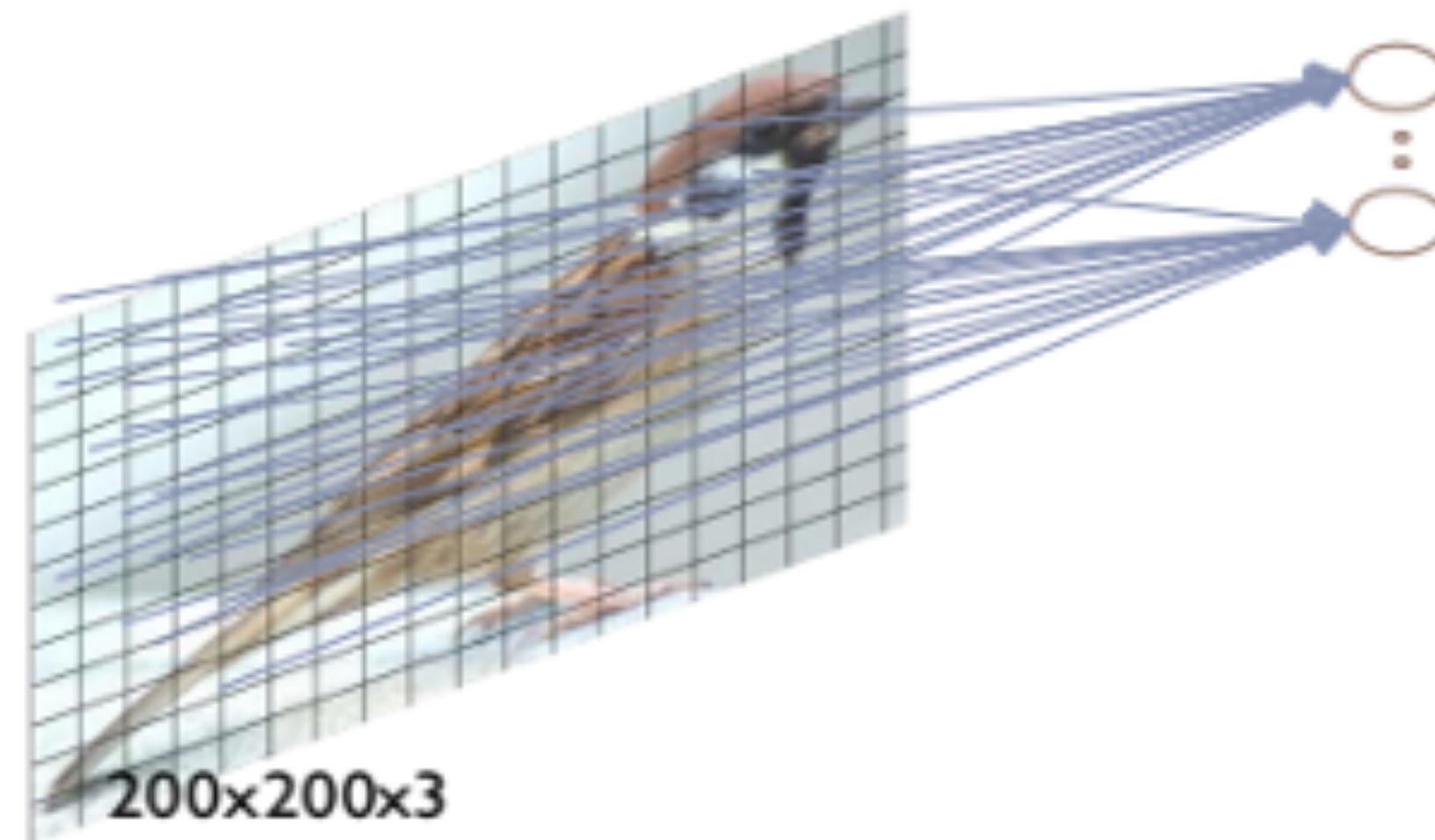


Convolutional connection

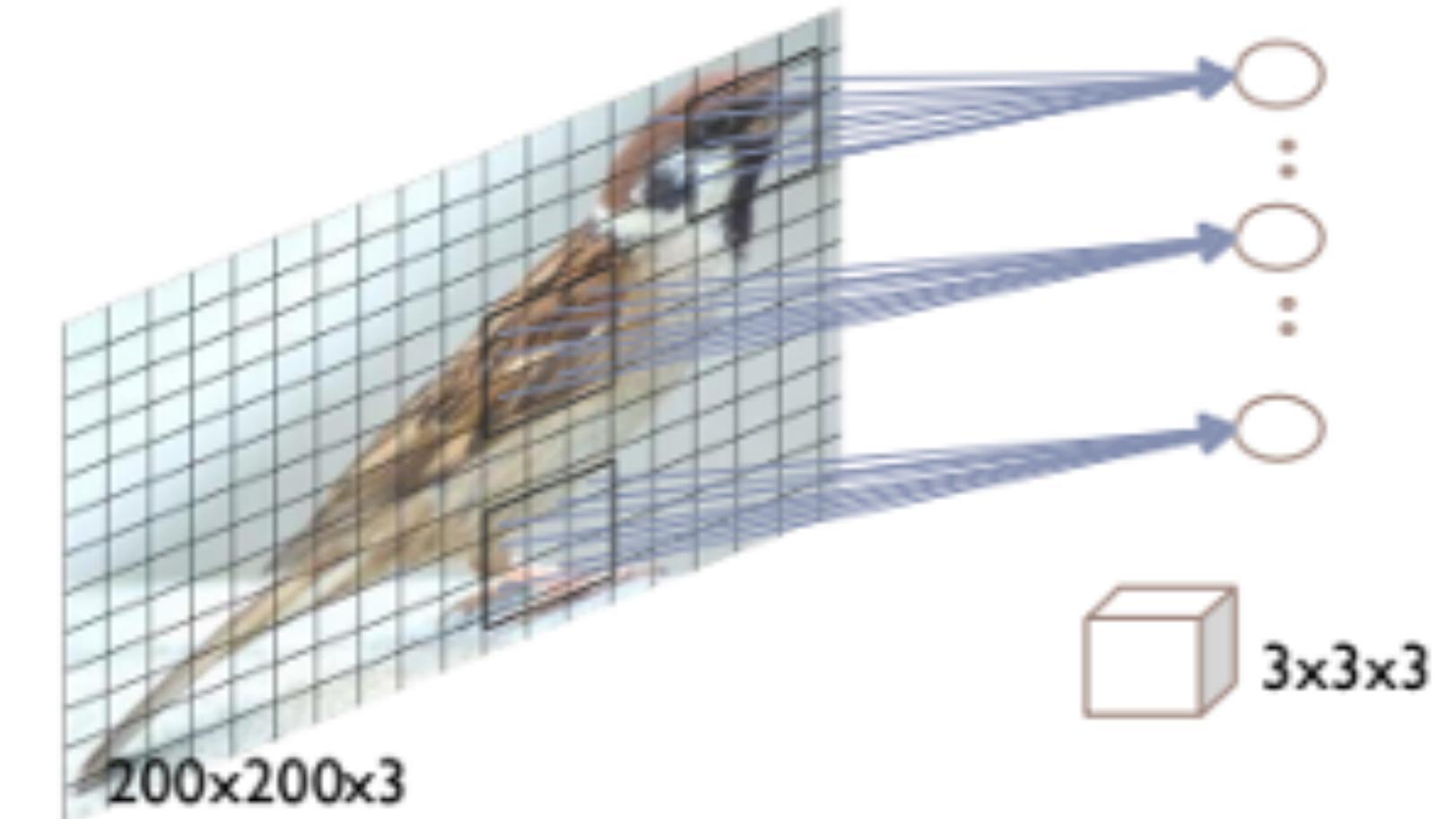


Convolution Layer

- Fully connected layer



- Locally connected layer

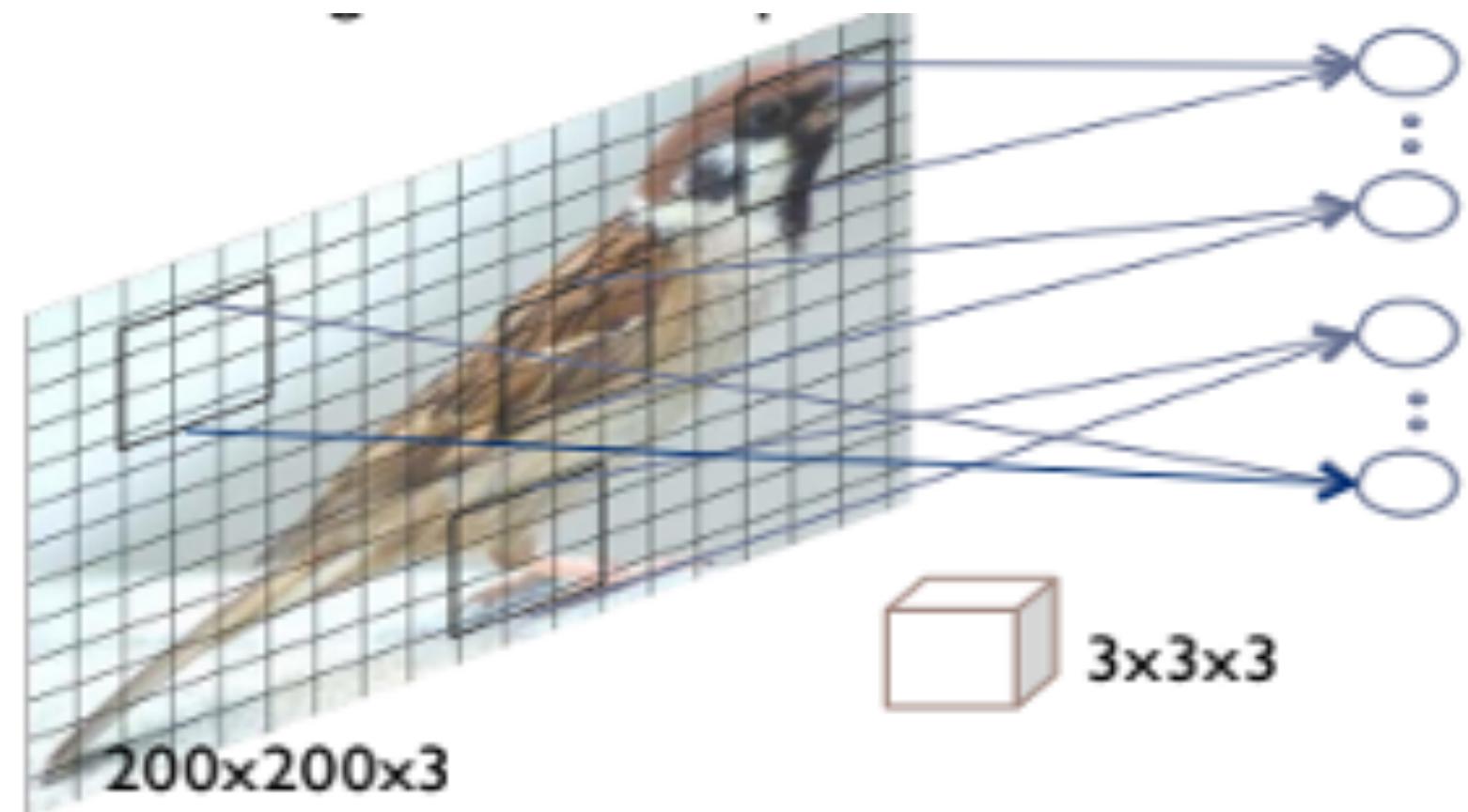


- Image of size 200×200 and 3 colours (RGB)
- #Hidden Units: 120,000 ($= 200 \times 200 \times 3$)
- #Params: 14.4 billion ($= 120K \times 120K$)
- Need huge training data to prevent overfitting!

- #Hidden Units: 120,000
- #Params: 3.2 Million ($= 120K \times 27$)
- Useful when the image is highly registered

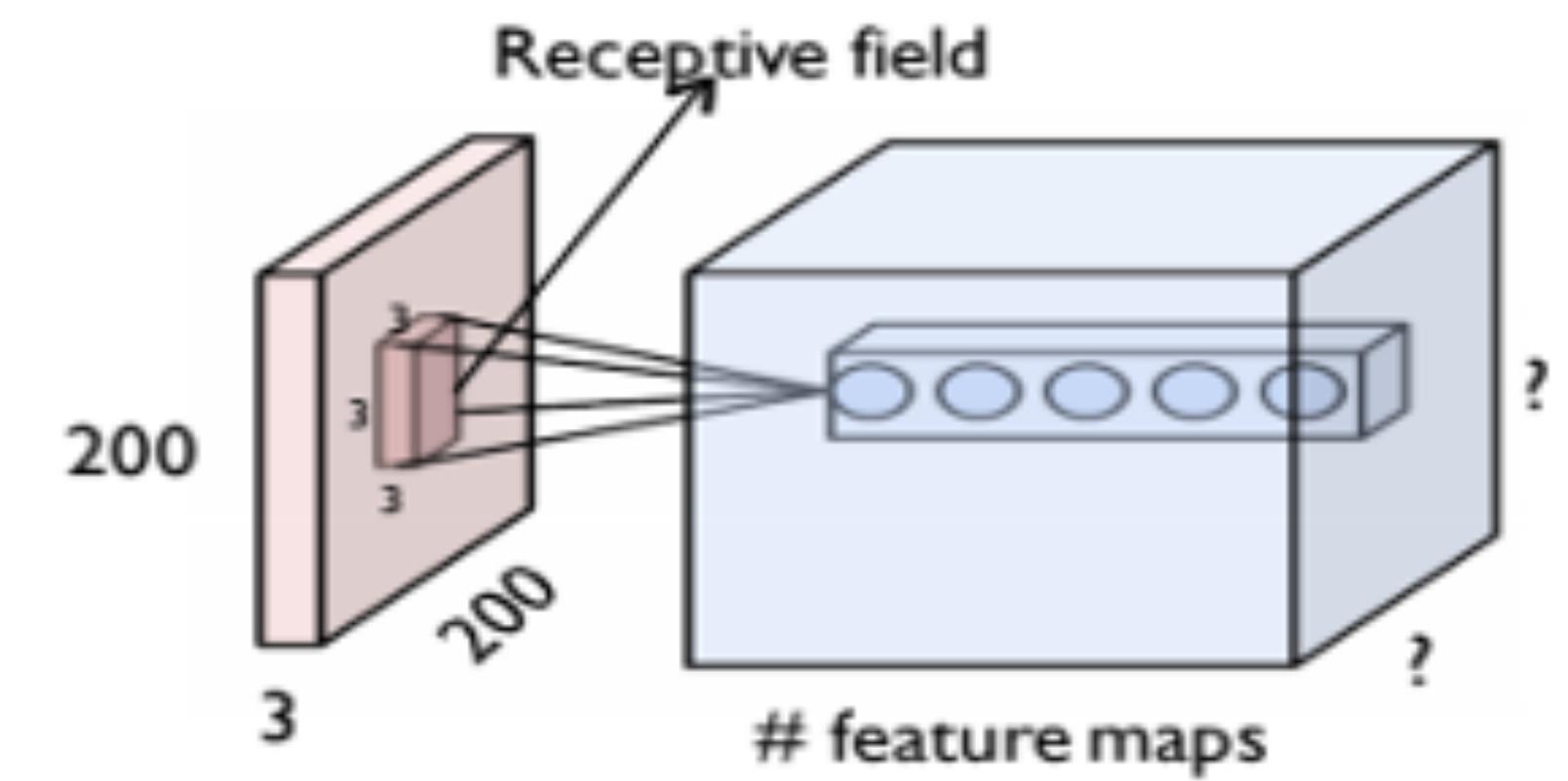
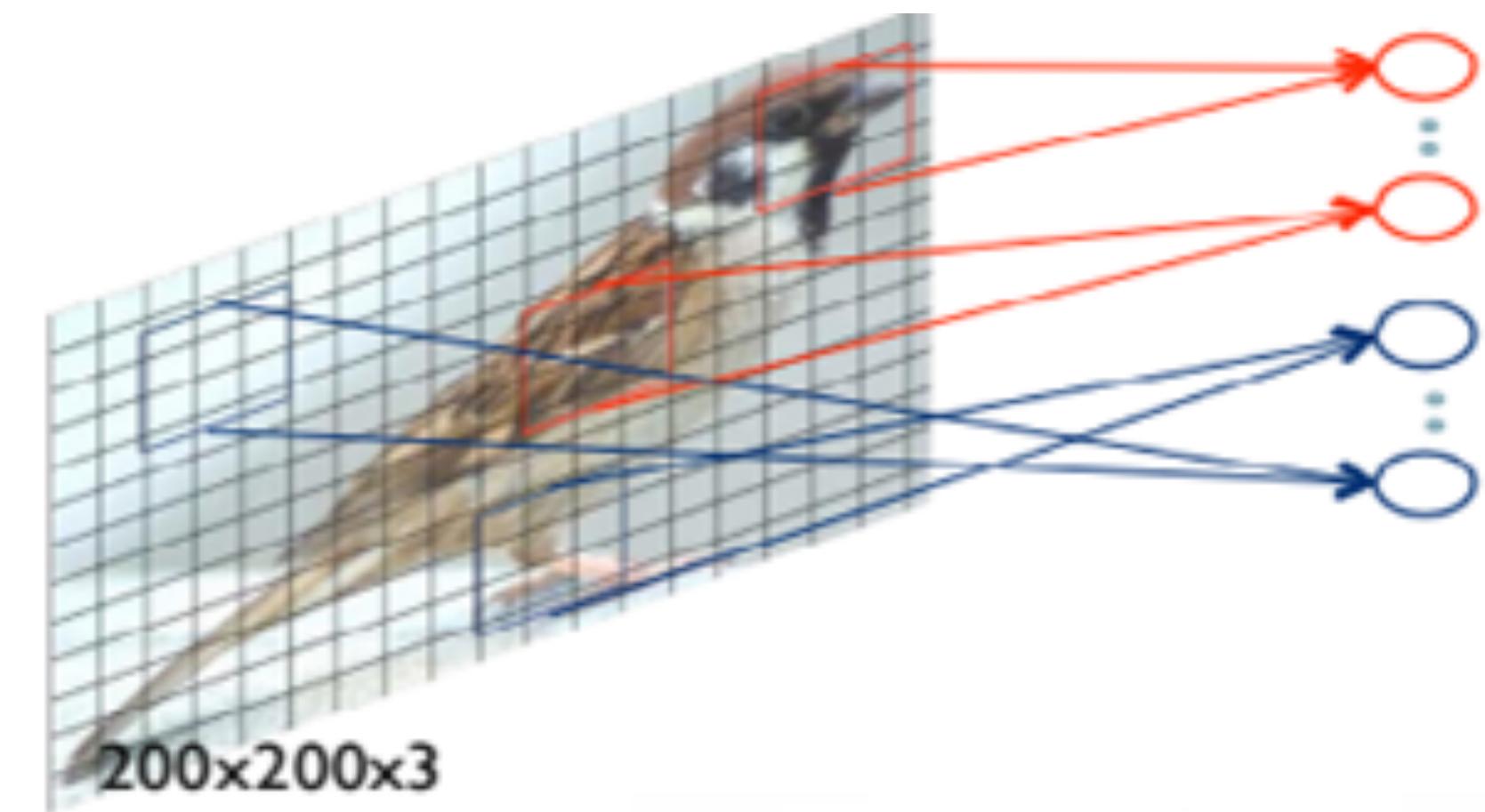
Convolution Layer

- Convolutional layer with single feature map.



- #Hidden Units: 120,000
- #Params: $27 \times \# \text{Feature Maps}$
- Sharing parameters
- Exploiting the stationarity property and preserves locality of pixel dependencies

- Convolutional layer with multiple feature maps



Stride and padding

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

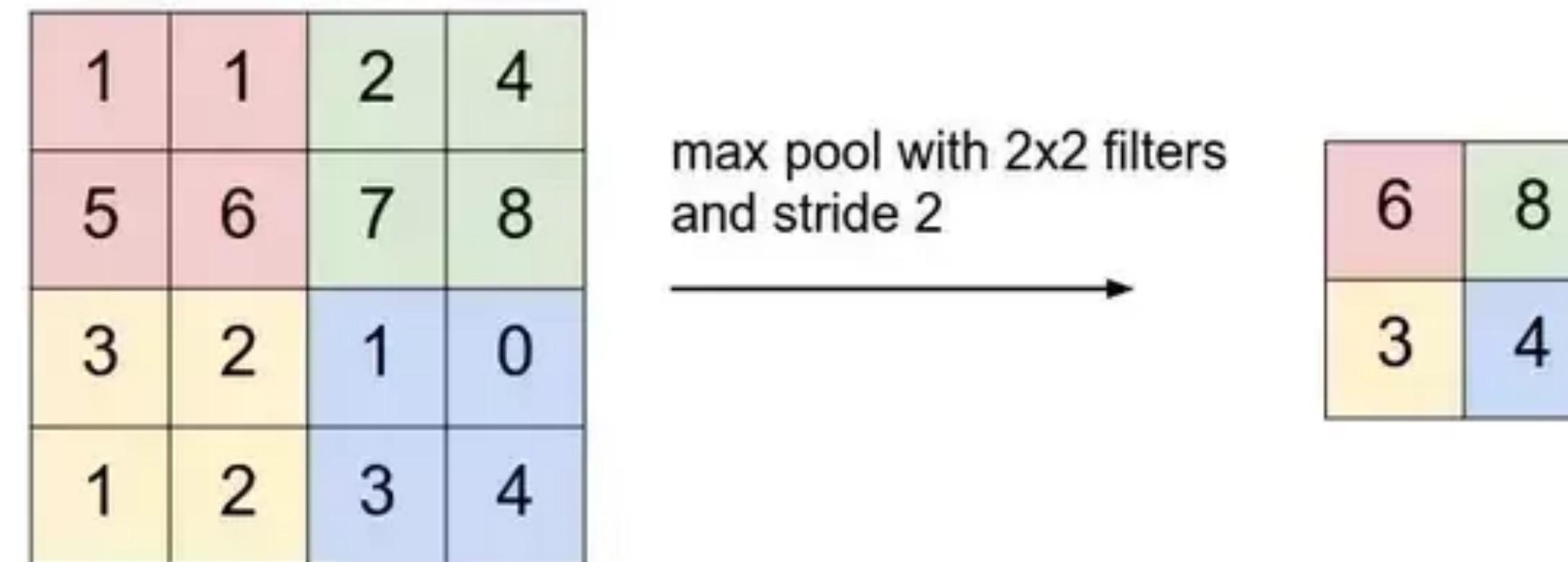
12	12	17
10	17	19
9	6	14

0 ₂	0 ₀	0 ₁	0	0	0	0
0 ₁	2 ₀	2 ₀	3	3	3	0
0 ₀	0 ₁	1 ₁	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

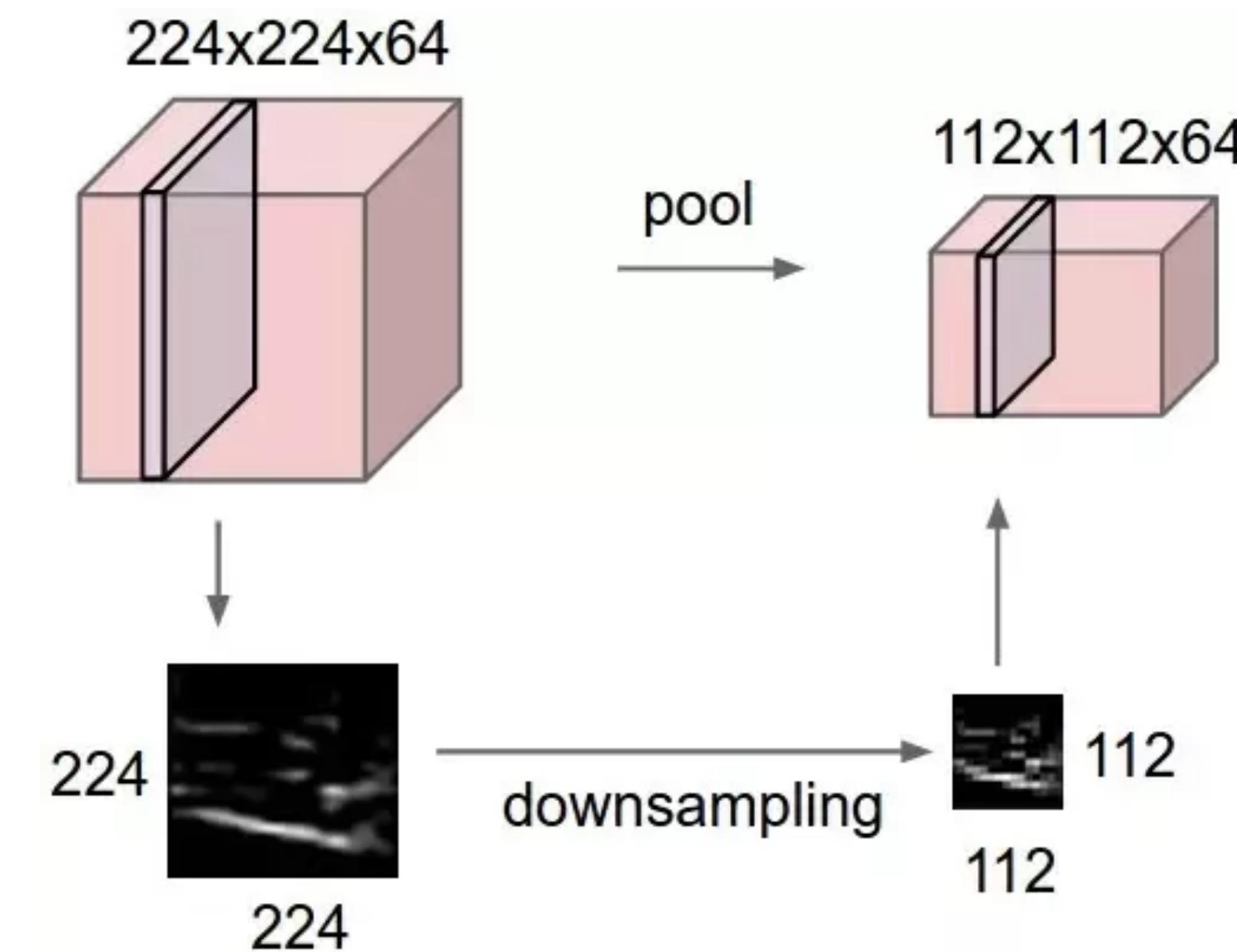
1	6	5
7	10	9
7	10	8

For instance, here is a 3×3 kernel applied to a 5×5 input padded with a 1×1 border of zeros using 2×2 strides:

Pooling

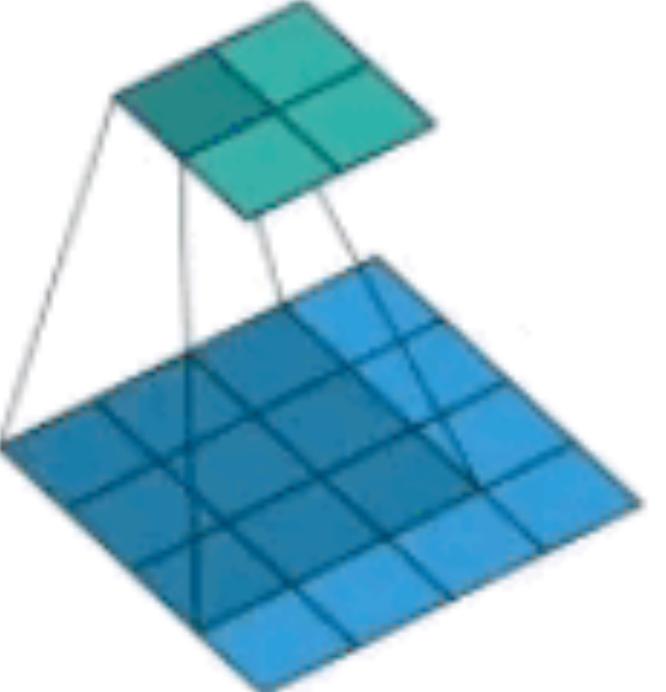
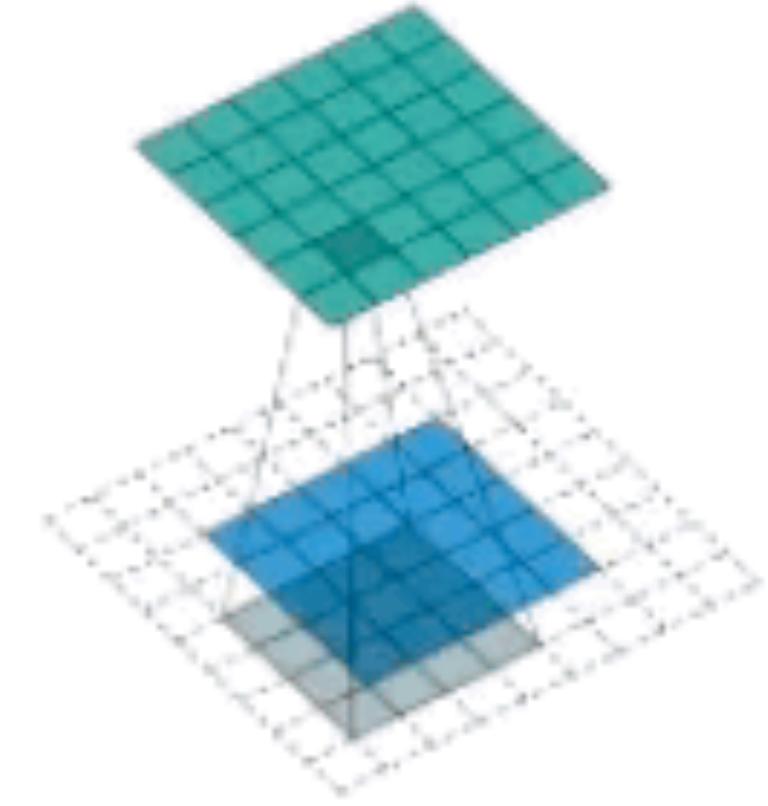
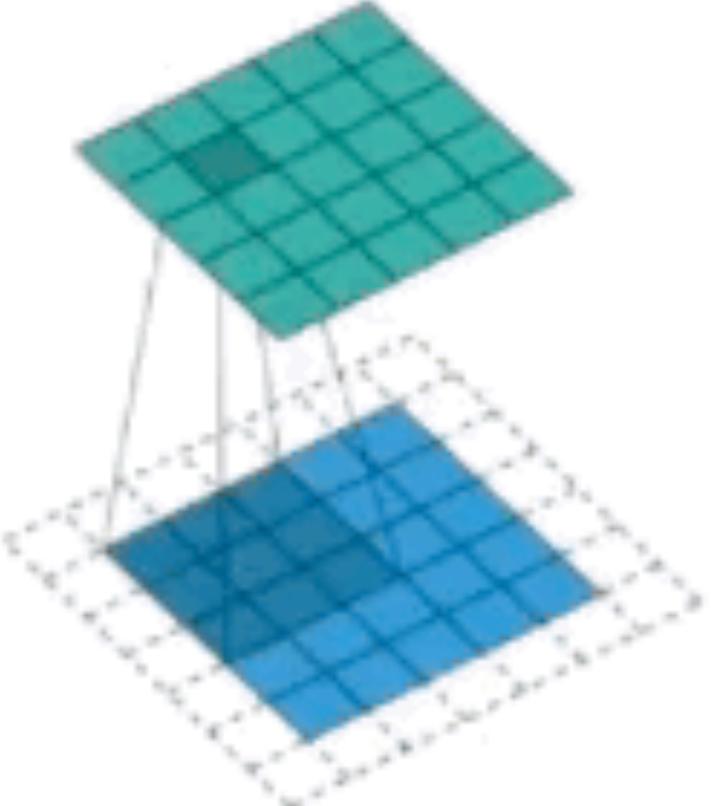
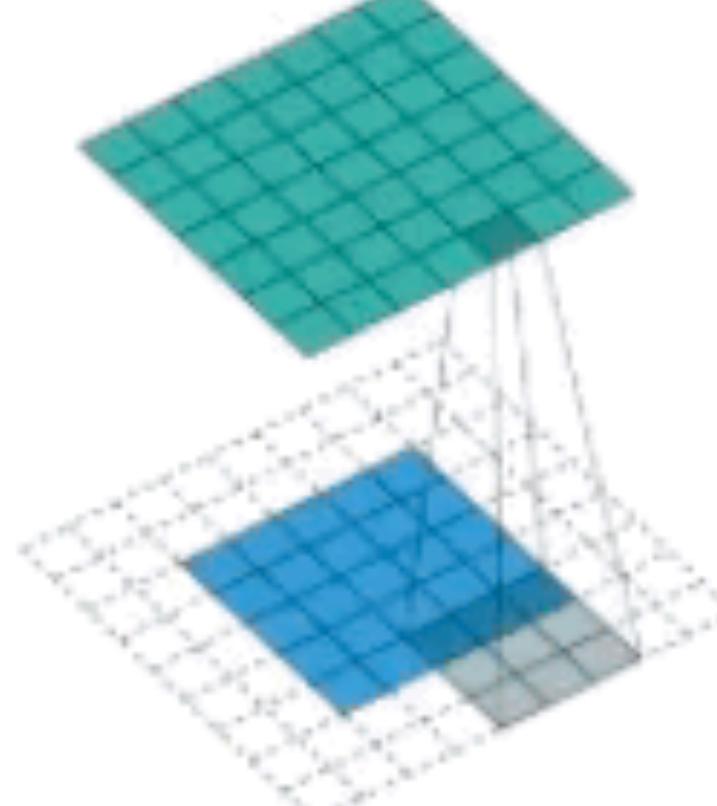
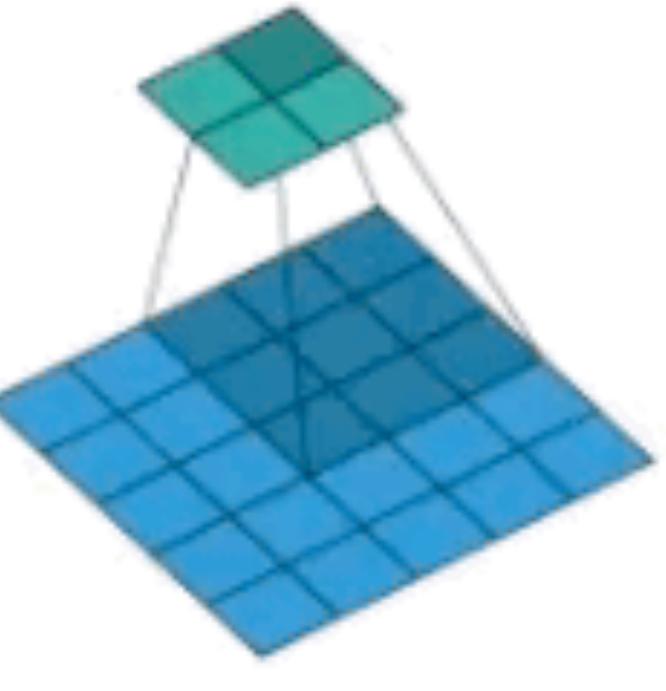
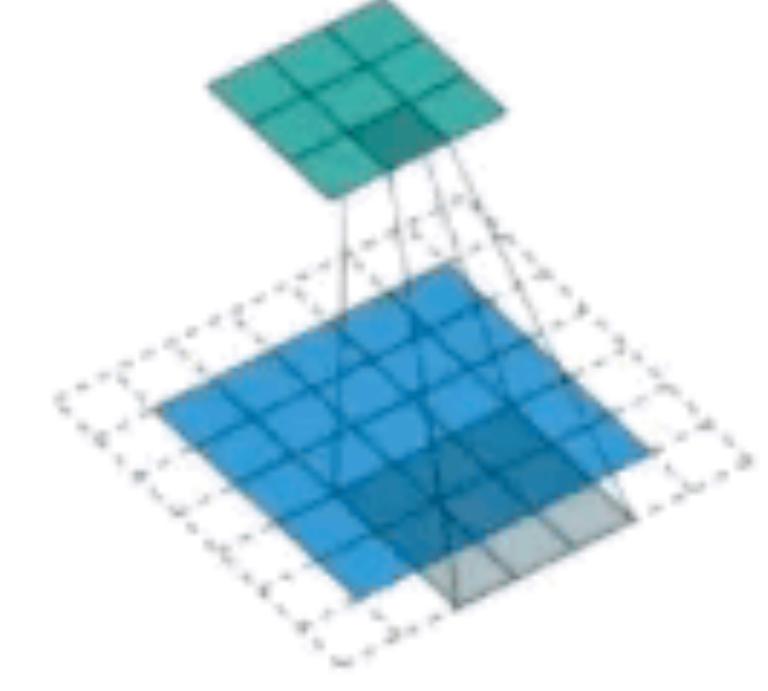
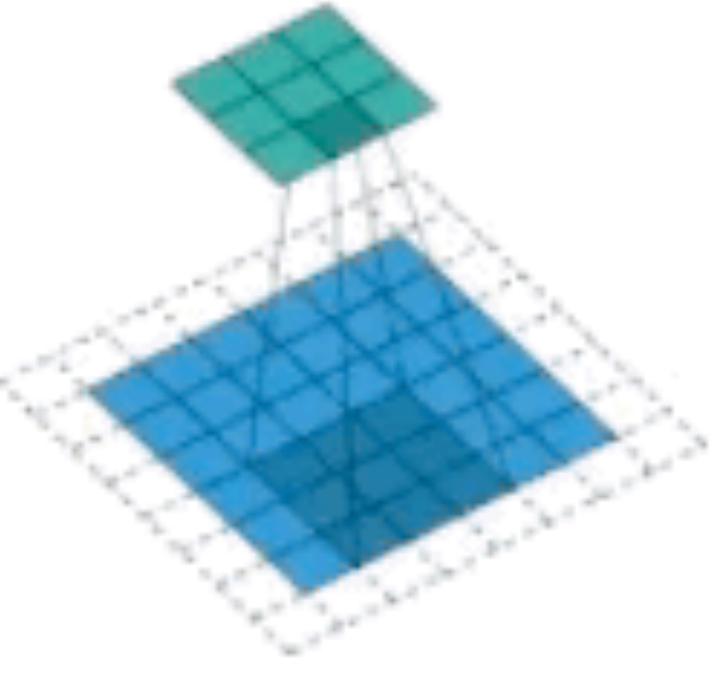


- **Motivation: We care about presence of features, not their exact location !**
- Dimensionality Reduction
- Prevents overfitting



Convolution animations

N.B.: Blue maps are inputs, and cyan maps are outputs.

			
No padding, no strides	Arbitrary padding, no strides	Half padding, no strides	Full padding, no strides
			
No padding, strides	Padding, strides	Padding, strides (odd)	

Input Volume (+pad 1) (7x7x3)

 $x[:, :, 0]$

0	0	0	0	0	0	0
0	0	2	1	2	0	0
0	2	1	2	2	0	0

 $x[:, :, 1]$

0	0	0	1	0	1	0
0	1	2	1	1	2	0
0	0	1	0	1	1	0
0	0	0	0	0	0	0

 $x[:, :, 2]$

0	0	0	0	0	0	0
0	0	2	2	1	0	0
0	1	0	0	0	0	0
0	2	2	2	1	2	0
0	2	0	0	0	2	0
0	0	2	2	0	0	0
0	0	0	0	0	0	0

 $x[:, :, 3]$

0	0	0	0	0	0	0
0	2	2	1	2	2	0
0	1	1	2	2	0	0
0	0	1	2	2	1	0
0	2	1	1	2	0	0
0	0	1	2	1	1	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

 $w0[:, :, 0]$

-1	1	0
1	1	0
1	0	1

 $w0[:, :, 1]$

0	-1	-1
1	1	1
0	0	-1

 $w0[:, :, 2]$

0	-1	0
1	0	0
-1	-1	0

Bias $b0 (1 \times 1 \times 1)$
 $b0[:, :, 0]$

1

Filter W1 (3x3x3)

 $w1[:, :, 0]$

1	-1	-1
0	1	-1
1	1	0

 $w1[:, :, 1]$

1	-1	0
0	0	0
-1	1	-1

 $w1[:, :, 2]$

1	1	1
-1	1	0
1	1	1

Bias $b1 (1 \times 1 \times 1)$ $b1[:, :, 0]$

0

Output Volume (3x3x2)

 $o[:, :, 0]$

3	11	6
5	8	4
0	5	3

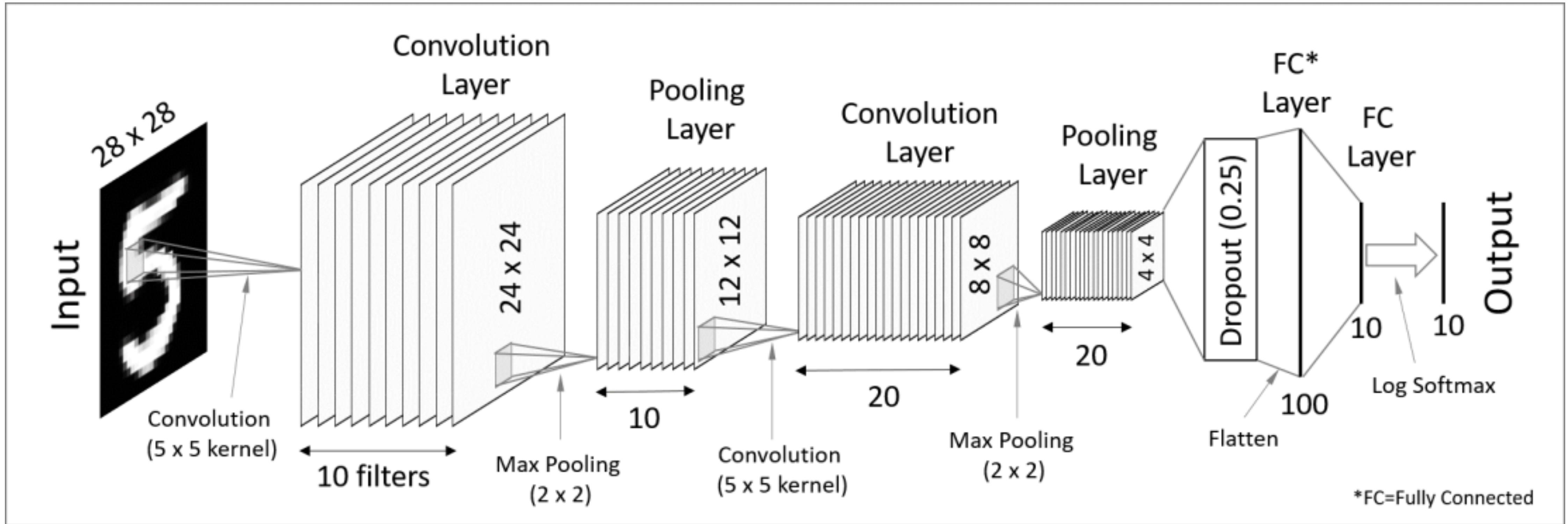
 $o[:, :, 1]$

5	6	4
4	11	11
-3	4	0

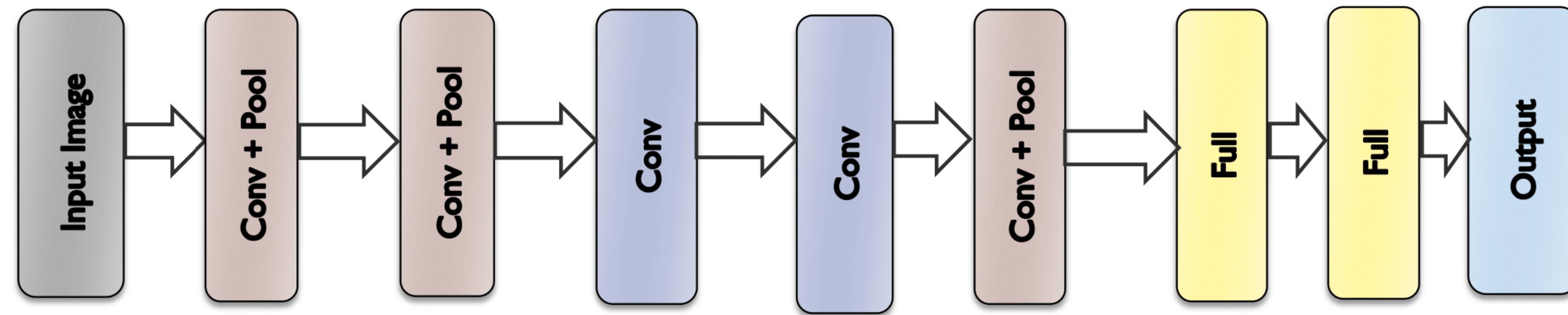
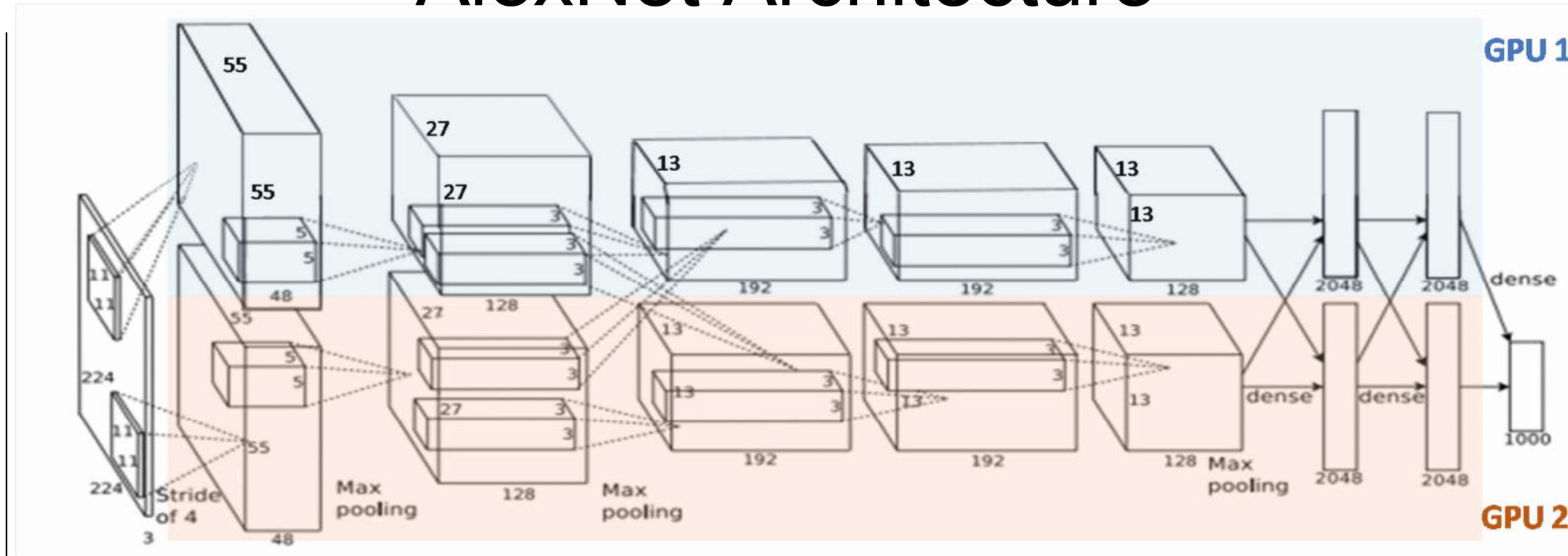
toggle movement

$$O = \frac{(W - K + 2P)}{S} + 1$$

where O is the output height/length, W is the input height/length, K is the filter size, P is the padding, and S is the stride.

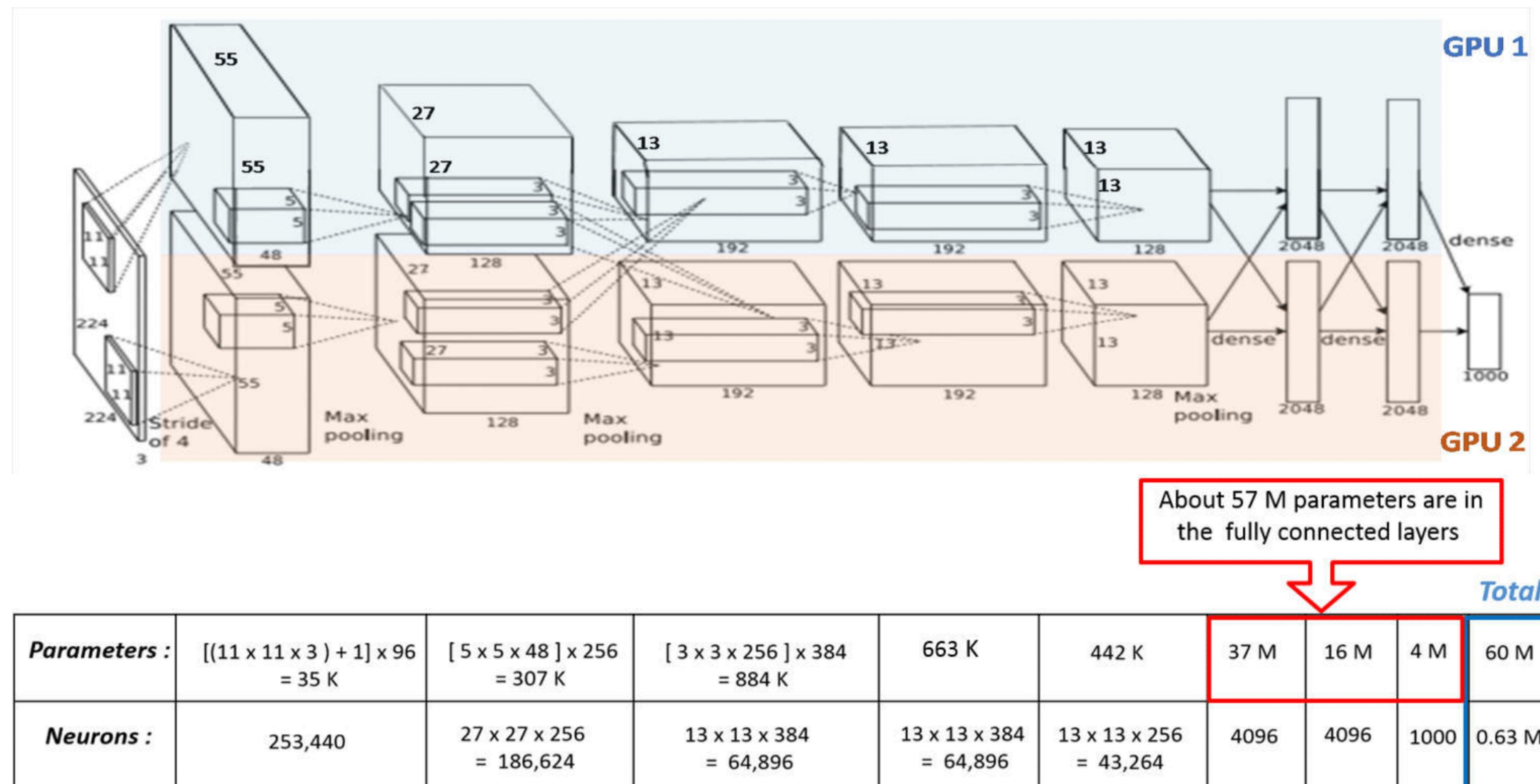


AlexNet Architecture





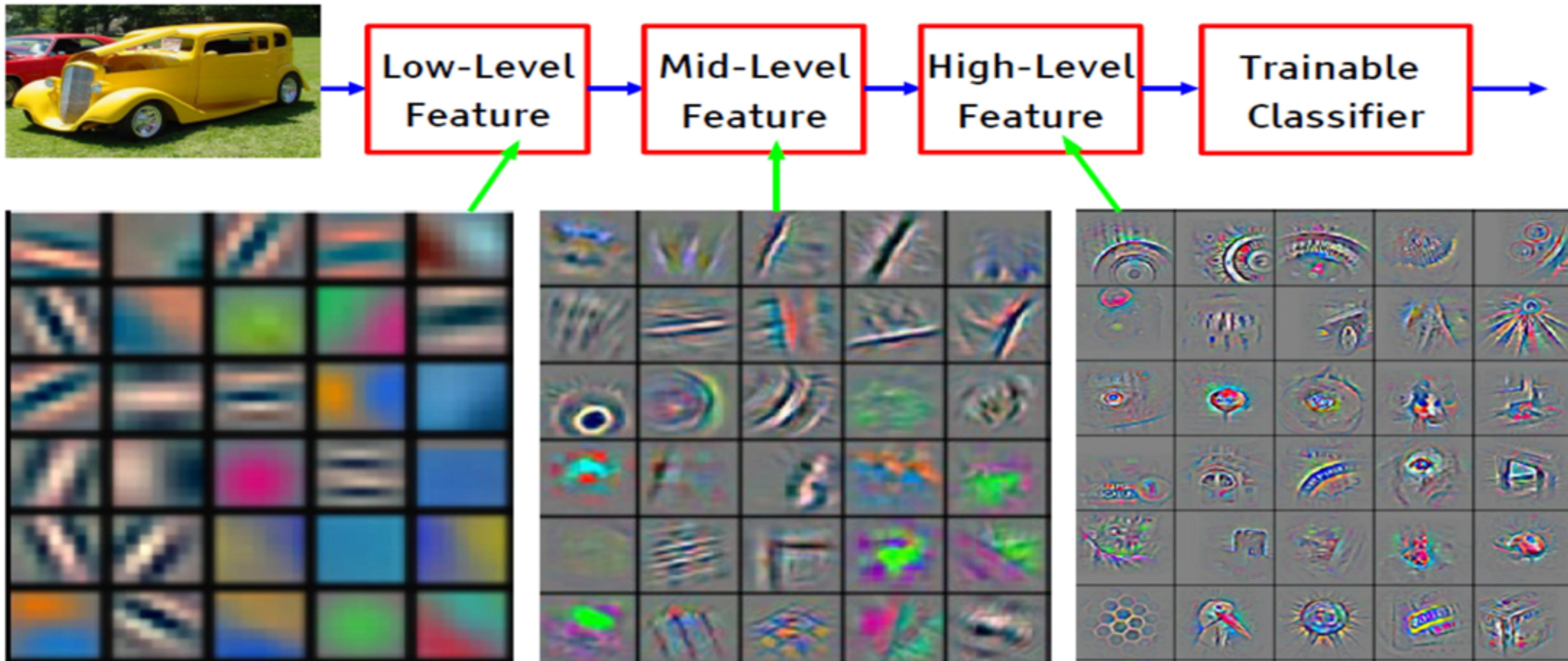
AlexNet Architecture



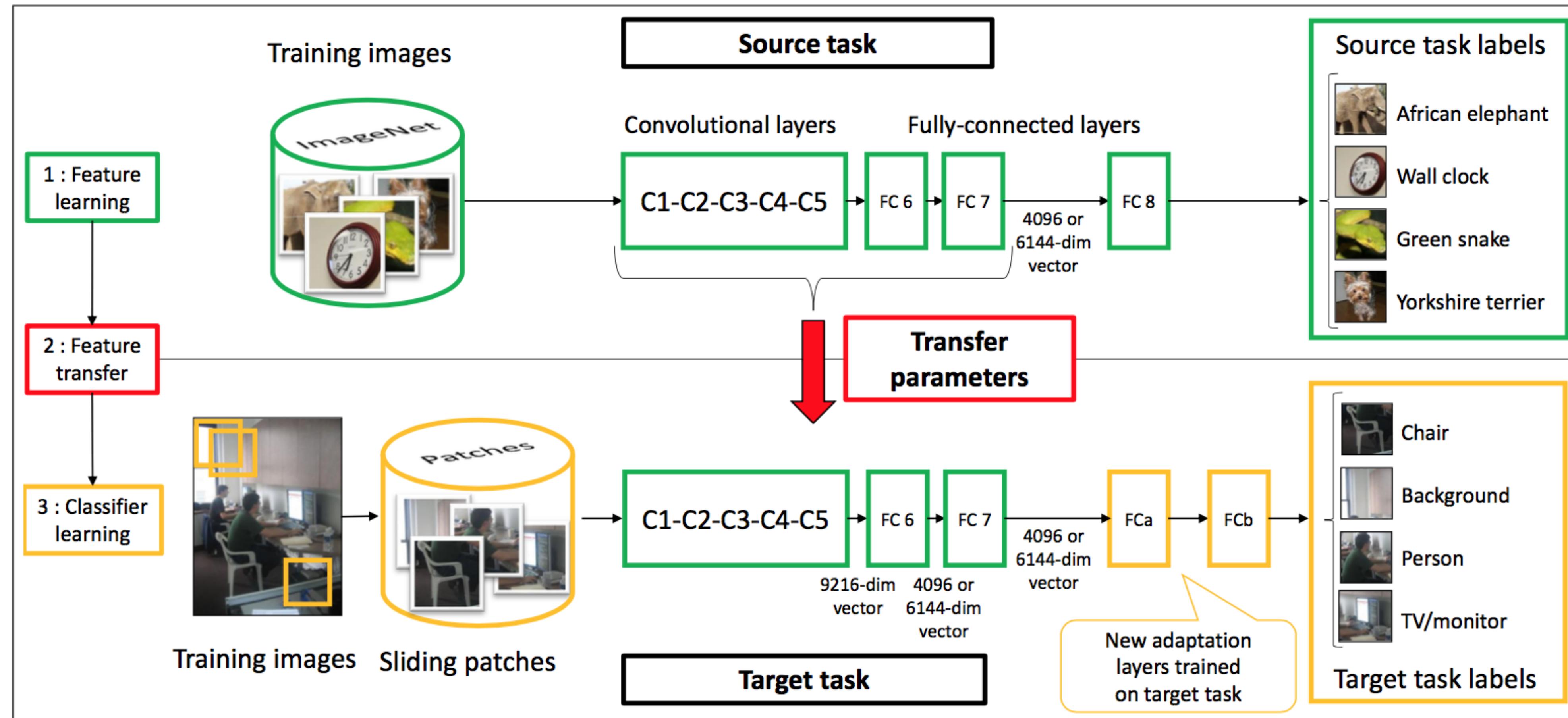
- Convolutional layers cumulatively contain about 90-95% of computation, only about 5% of the parameters
- Fully-connected layers contain about 95% of parameters.

Deep Learnt Features

- It's **deep** if it has **more than one stage** of non-linear feature transformation



Transfer Learning





Input



Segmentation [9]

