



Convolutional Neural Networks (Part 1)

COMP90086 Computer Vision

QiuHong Ke

Copyright: University of Melbourne

Before we start...

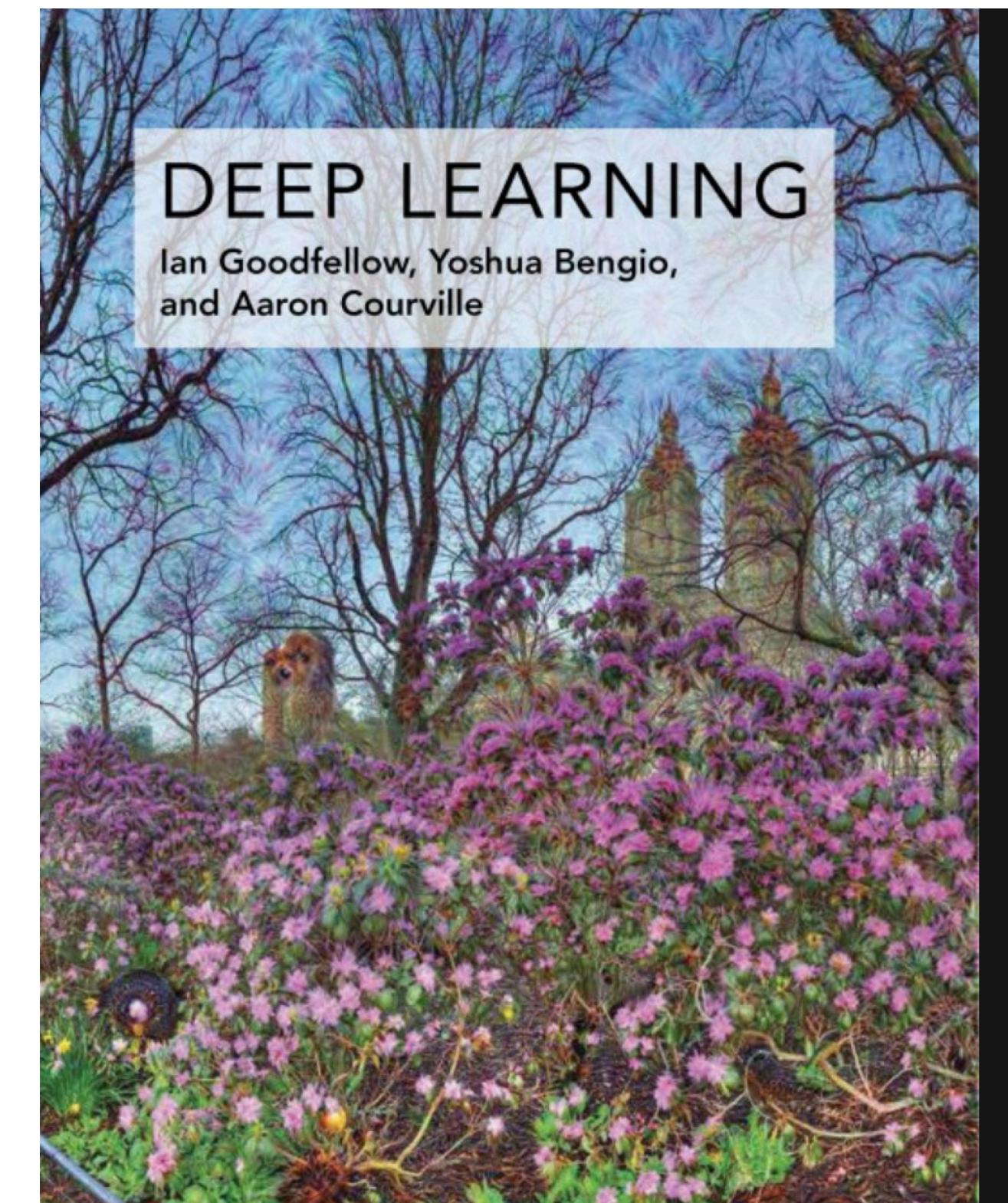
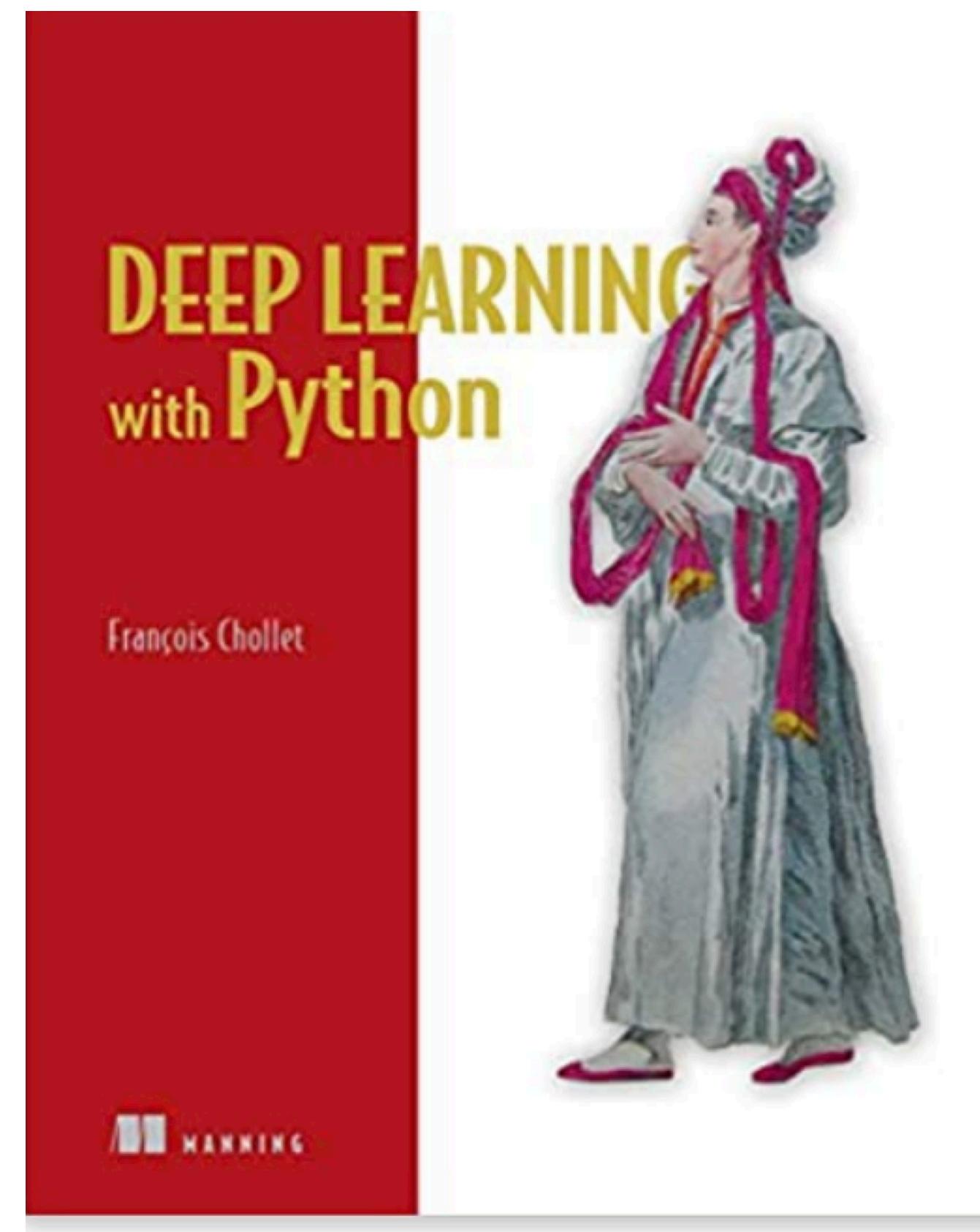
About me

- 2015.02-2018.04: PhD in UWA
- 2018.05-2019.12: Post-doc in MPII
- From 2020.01: Lecturer in UniMelb
- Research: Action recognition and prediction
- Contact: qiu.hong.ke@unimelb.edu.com;
qiu.hong.ke@unimelb.edu.au

Before we start

Books & resources

- Books:
 - Deep Learning with Python, by Francois Chollet, available in unimelb library
 - Deep Learning, by Ian Goodfellow and Yoshua Bengio and Aaron Courville, <https://www.deeplearningbook.org>
- Tool: Keras, <https://keras.io/>



Outline

- Introduction on Deep Learning
- Neural Networks: Recap of Multi-layer Perceptron
- Convolutional Neural Networks

Angry birds

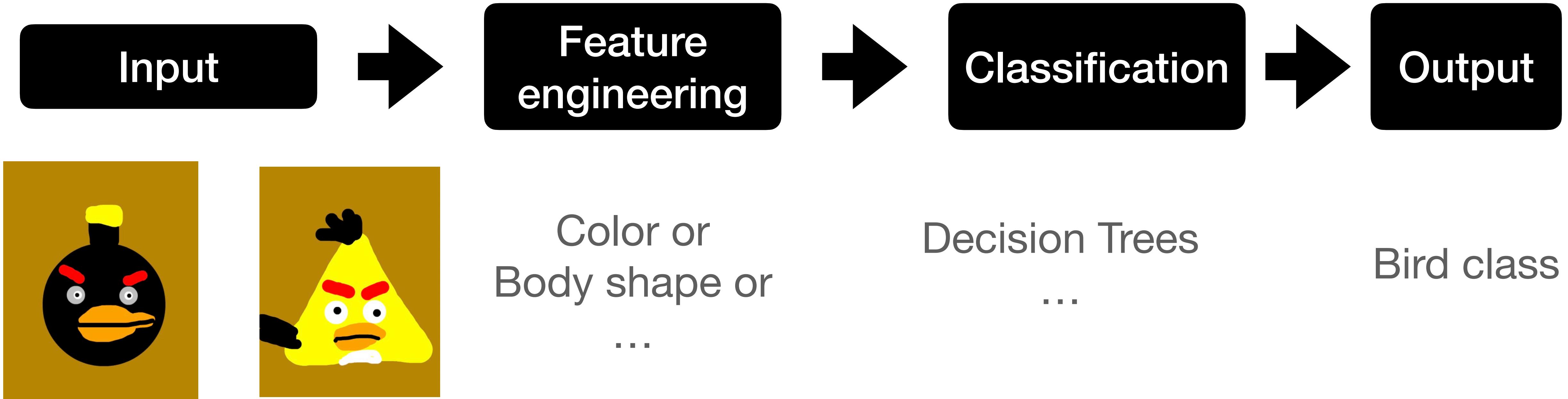


VS



Source: by Papachan (CC BY)
<https://www.sketchport.com/drawing/4524248689278976/bomb-and-chuck>

Bird classification



Introduction

Neural Networks

Convolutional Neural Networks

In real life...

Boat tailed Grackle



Bobolink



Bohemian Waxwing



Brandt Cormorant



Brewer Blackbird



Pomarine Jaeger



Prairie Warbler



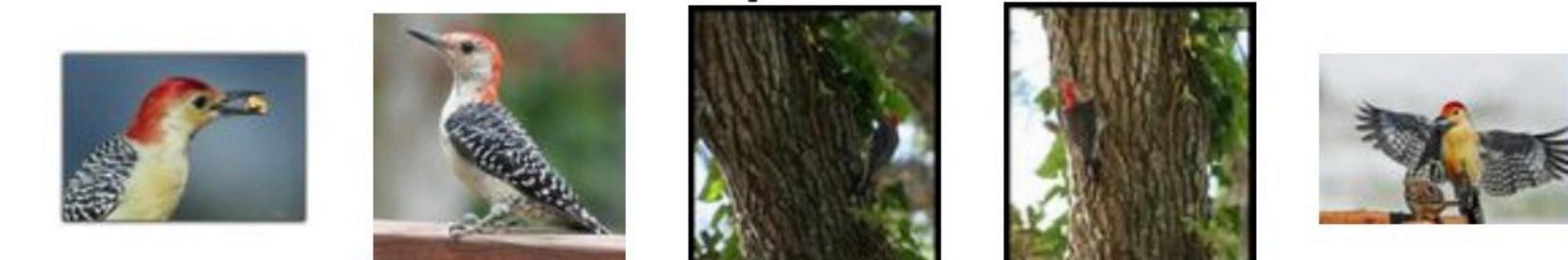
Prothonotary Warbler



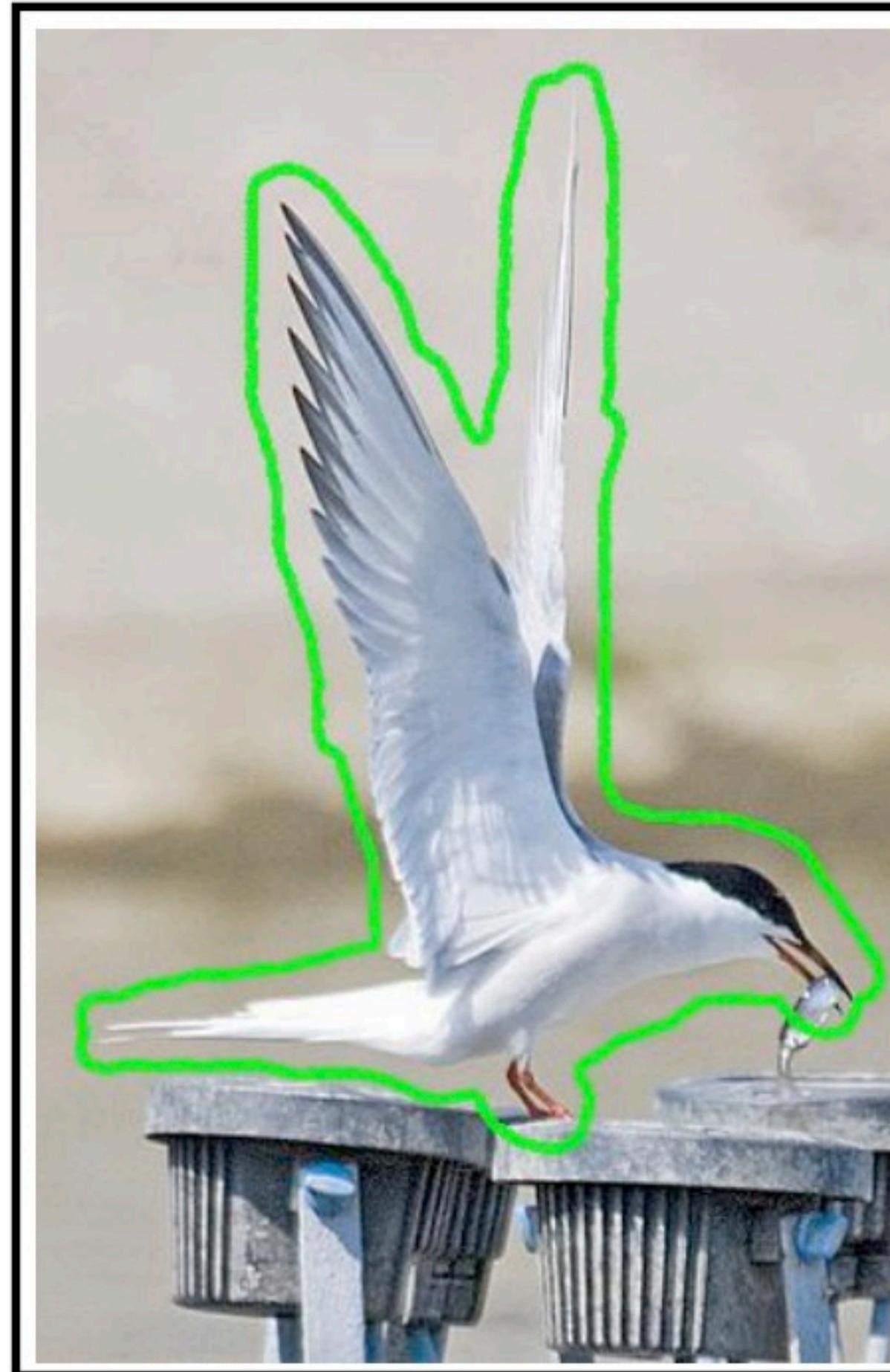
Purple Finch



Red bellied Woodpecker



More features for real birds...

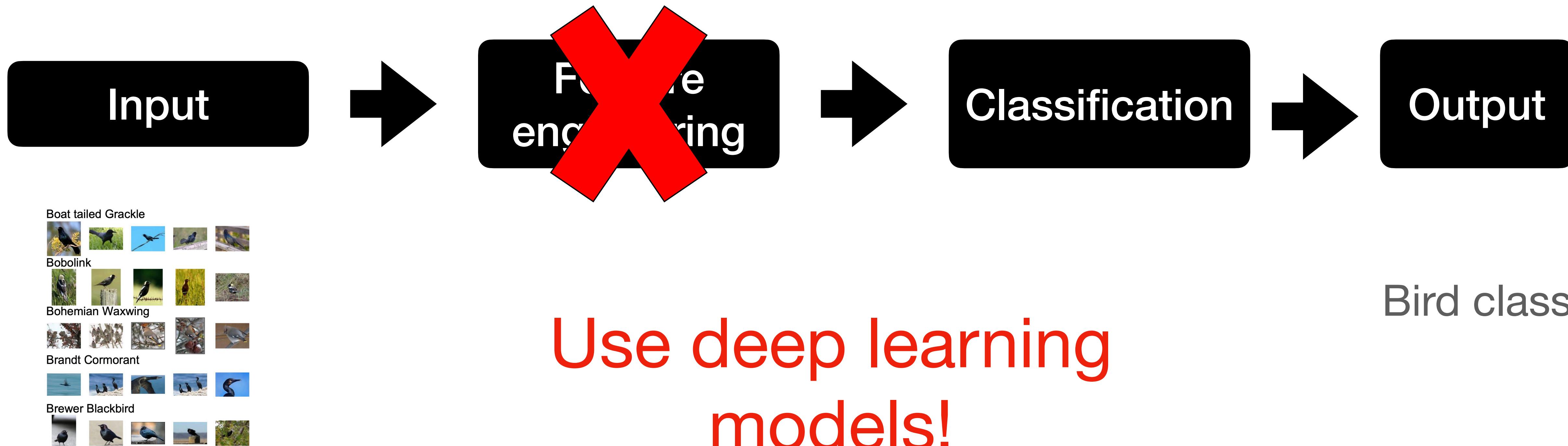


forehead_color	black	black	black
breast_pattern	solid	solid	solid
breast_color	white	white	white
head_pattern	plain	capped	plain
back_color	white	white	black
wing_color	grey/white	grey	white
leg_color	orange	orange	orange
size	medium	large	medium
bill_shape	needle	dagger	dagger
wing_shape	pointed	tapered	long
...
primary_color	white	white	white

Problems?

- time-consuming
- Require expert knowledge
- Require feature selection

'REAL Bird' classification



Deep learning is everywhere

Face recognition



By Mike MacKenzie

<https://www.flickr.com/photos/mikemacmarketing/30188201497>

Deep learning is everywhere

Self driving cars



Source: <https://usa.streetsblog.org/2020/12/15/was-2020-the-year-driverless-cars-became-inevitable/>

Deep learning is everywhere

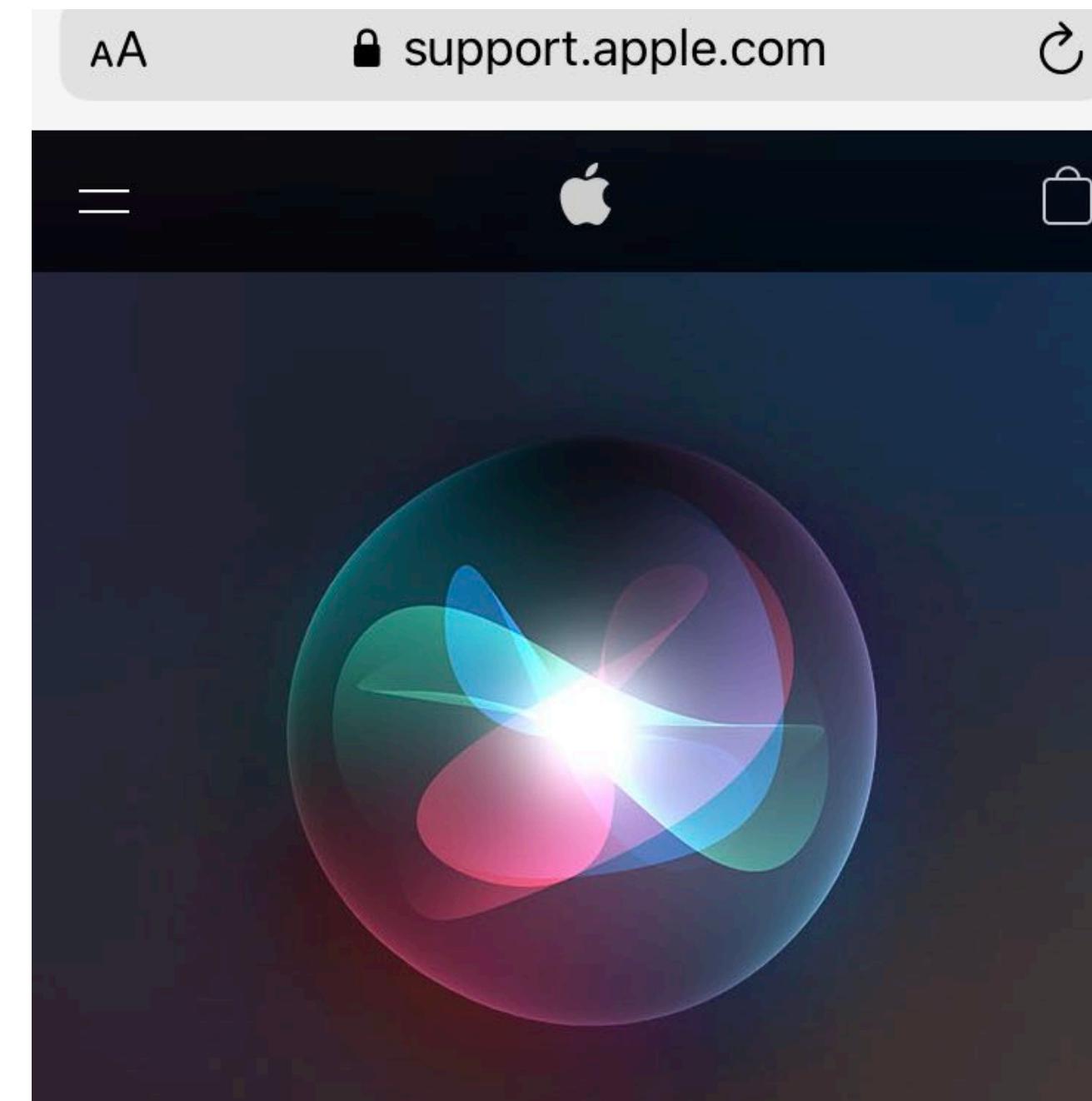
Virtual Assistants:

Amazon Alexa

Cortana

Google Assistant

Siri



What can I ask Siri?

Siri can help you do the things you do every day — and on any of your Apple devices. All you have to do is ask.

Deep learning is everywhere

Chatbots: AI application
to chat online (to solve
customer problems)

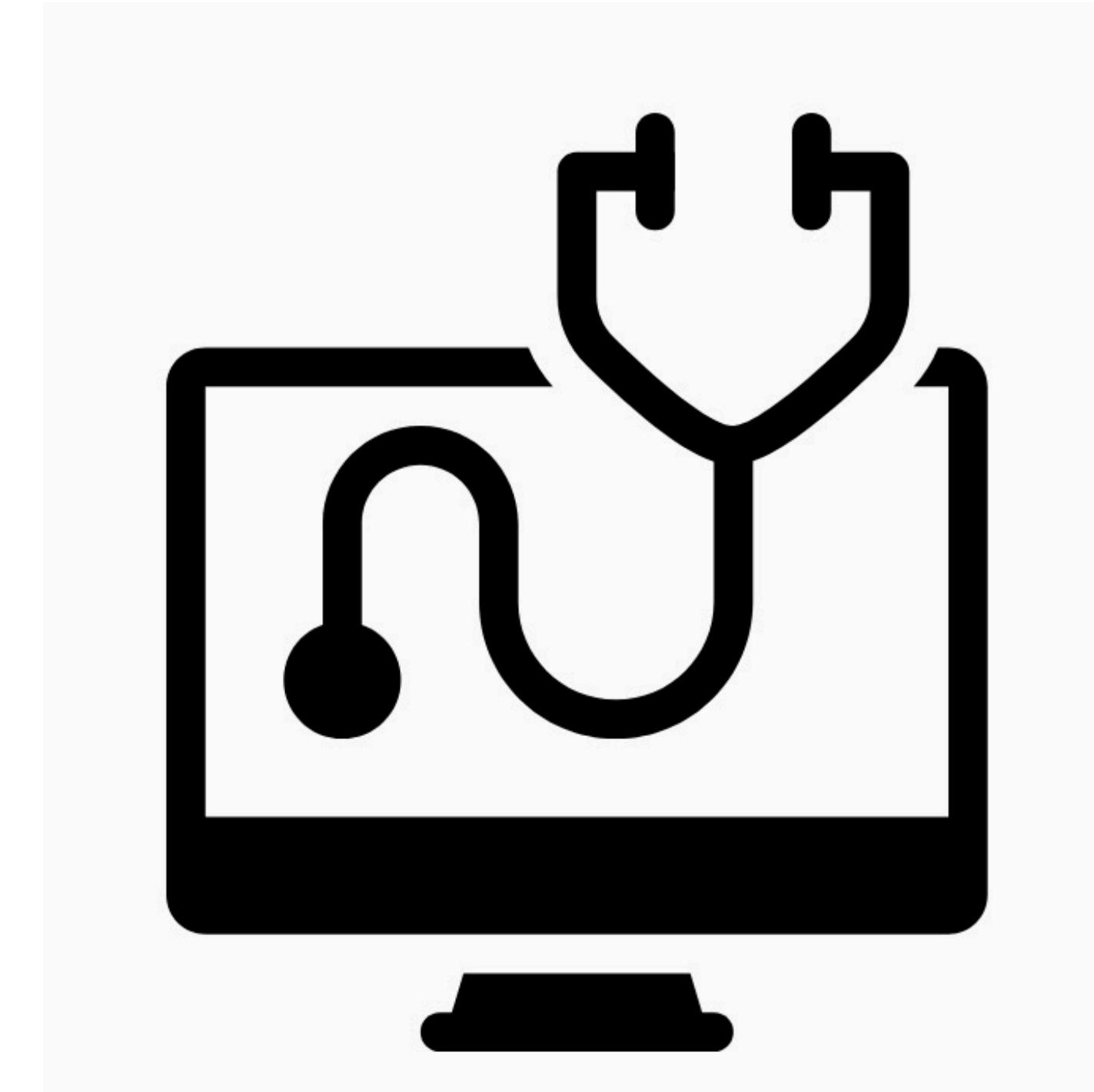


by mohamed hassan

source: <https://www.stockvault.net/photo/259374/chatbot-technology>

Deep learning is everywhere

Healthcare: medical research, medicine discovery, and disease diagnosis

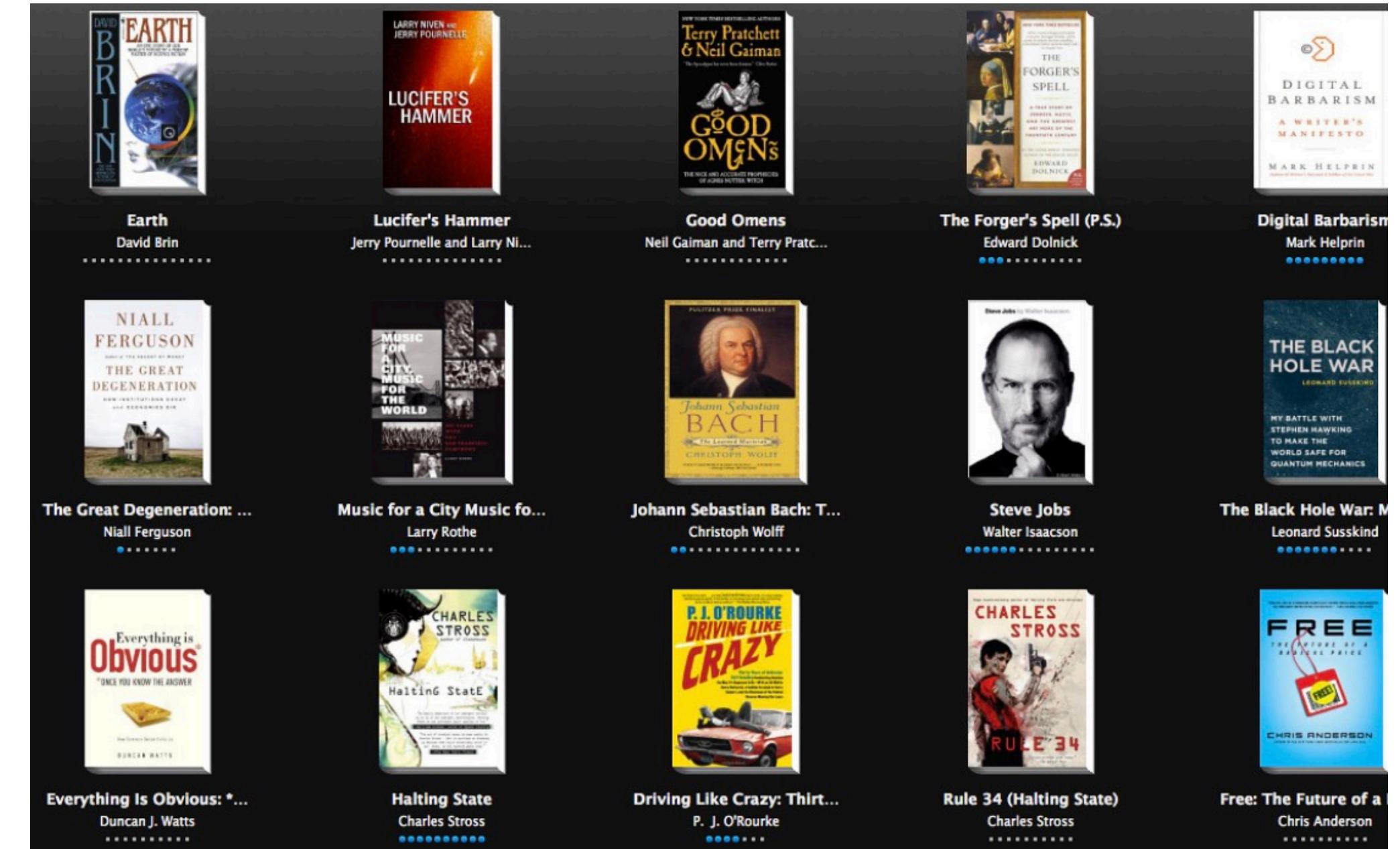


by Xinh Studio

source:<https://thenounproject.com/term/healthcare/1243846/>

Deep learning is everywhere

Recommendation:
books, videos...



by Jose Camões Silva

<https://www.flickr.com/photos/josecamoessilva/9115095115>

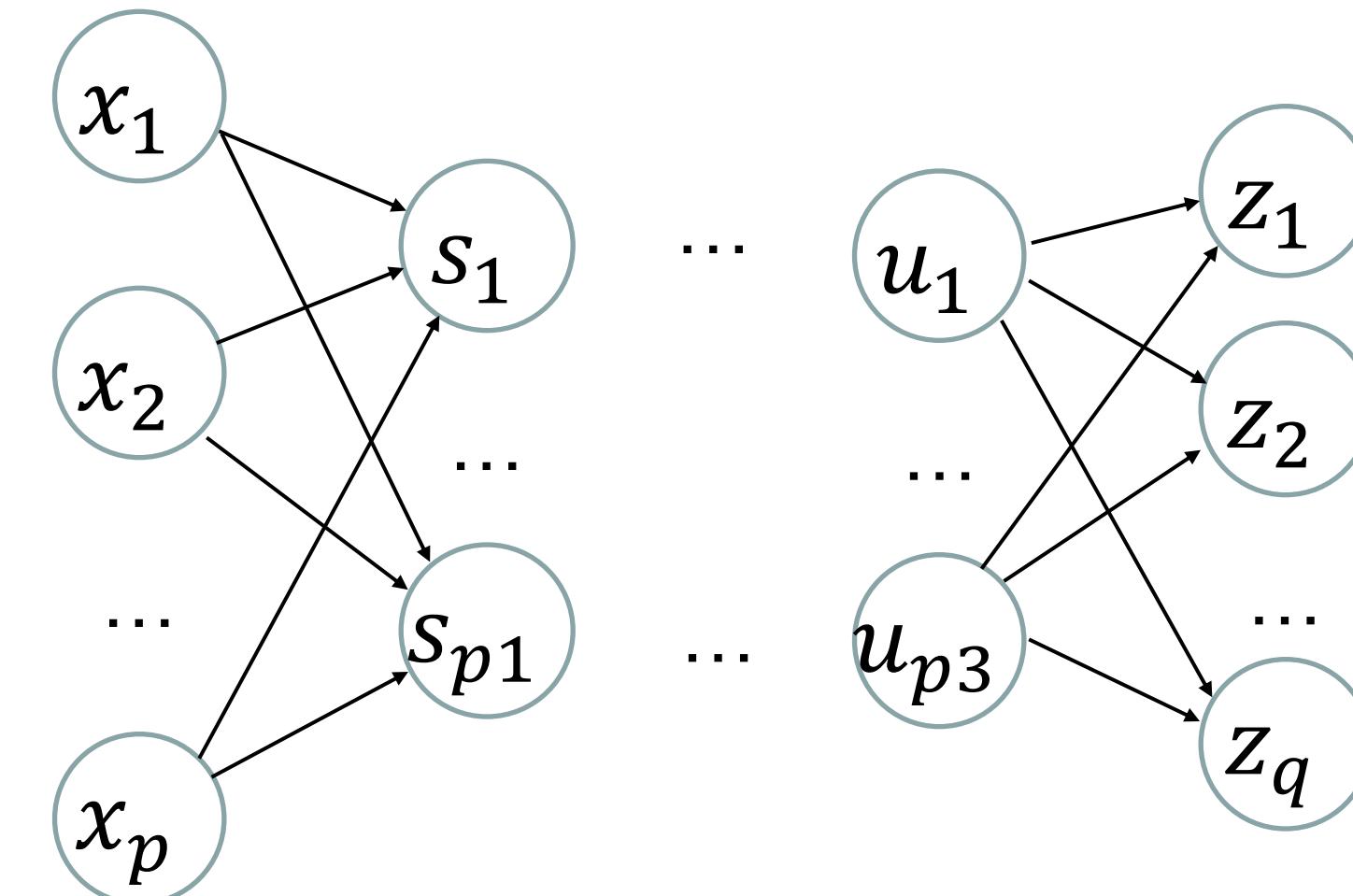
What is Deep learning

Input →

Feature learning + classification

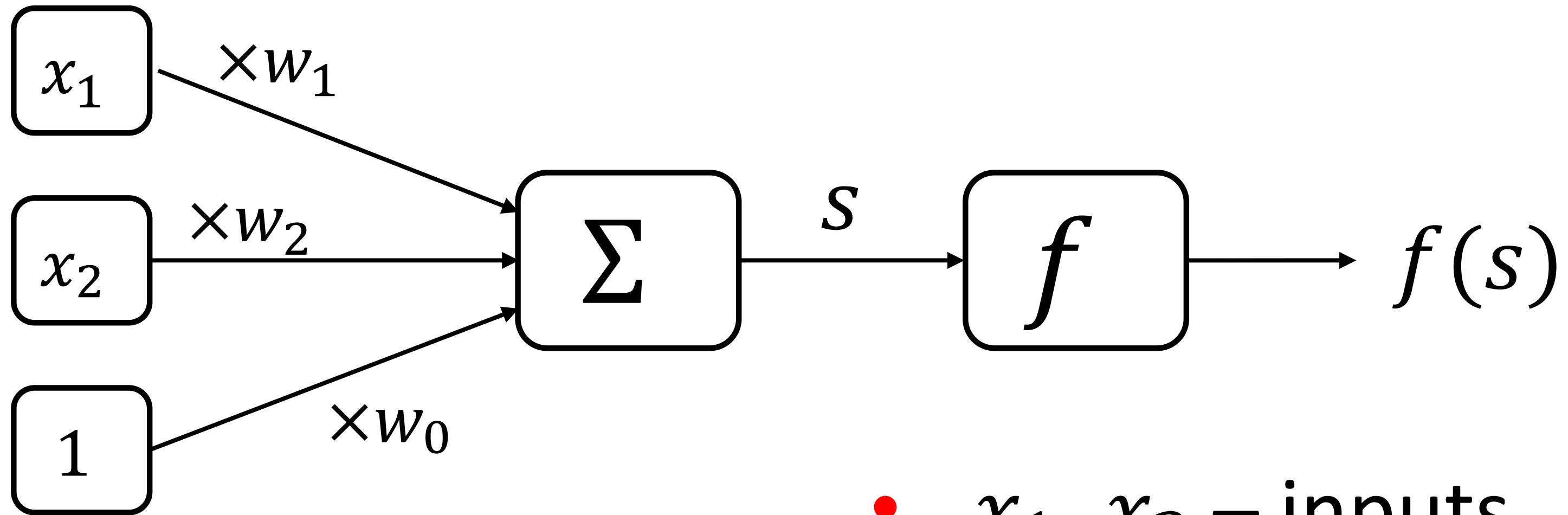
→ Output

Deep neural network



many hidden layers

Perceptron



$$s = \sum_{i=0}^m x_i w_i$$

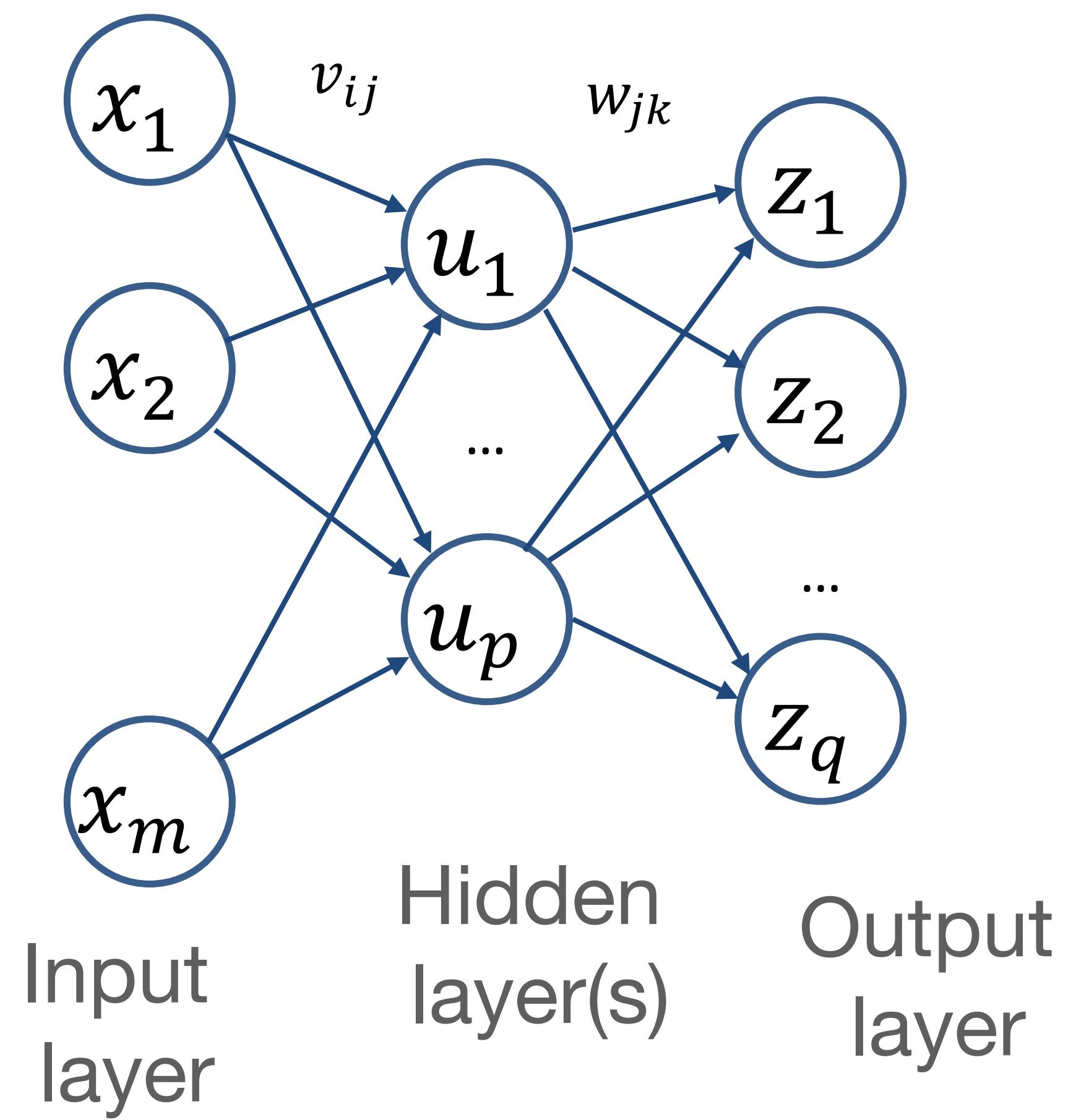
- x_1, x_2 – inputs
- w_1, w_2 – synaptic weights
- w_0 – bias weight
- f – activation function

Multi-layer perceptron: function composition

$$u_j = g(r_j)$$

$$r_j = \sum_{i=0}^m x_i v_{ij}$$

$g(x)$: activation function

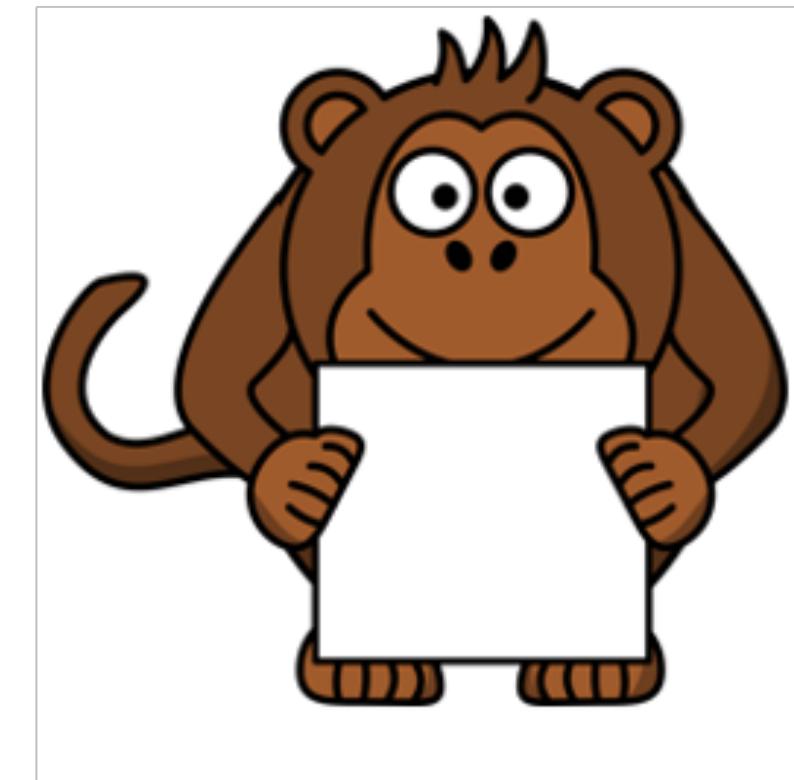
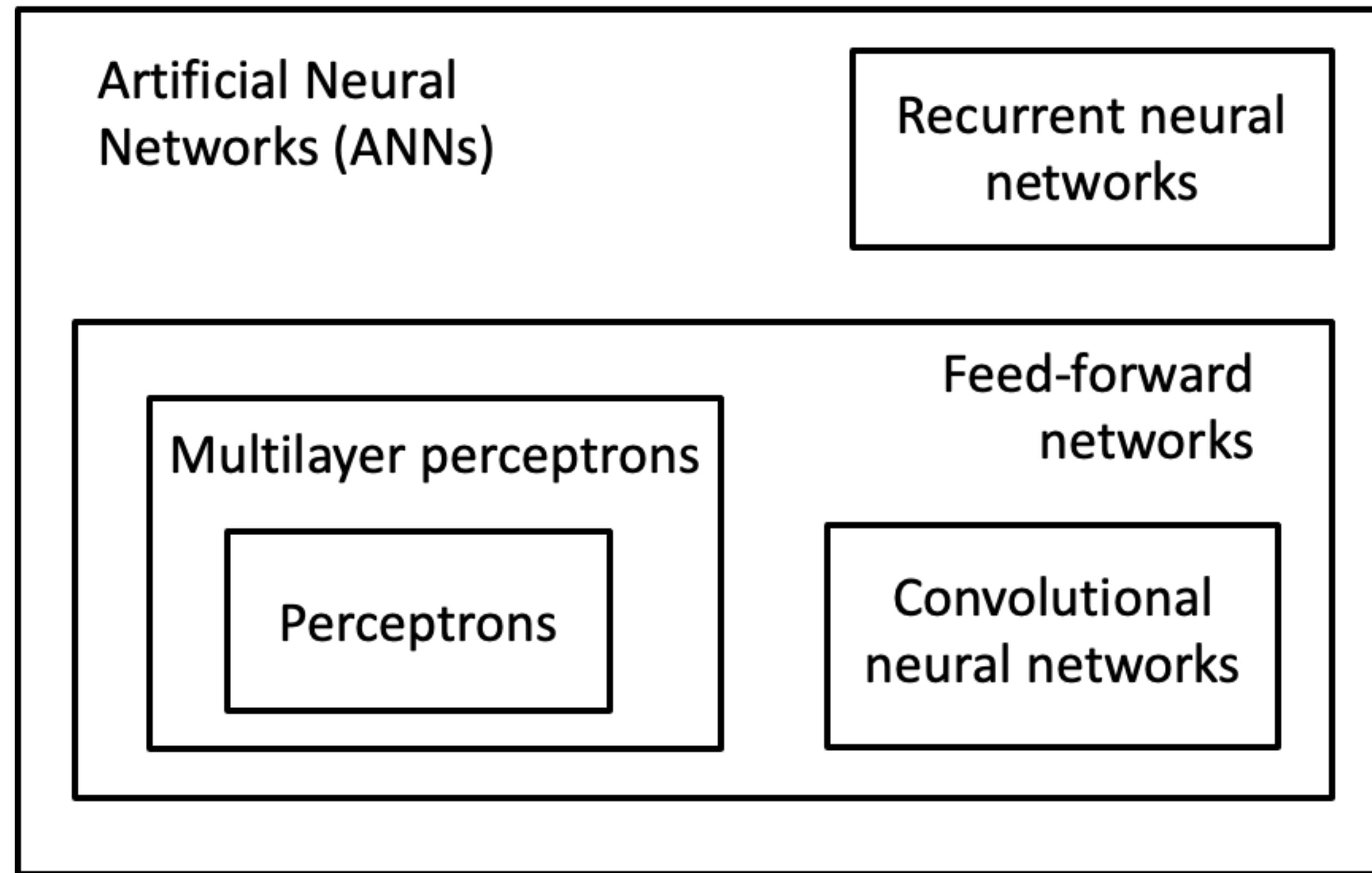


$$z_k = h(s_k)$$

$$s_k = \sum_{j=0}^p u_j w_{jk}$$

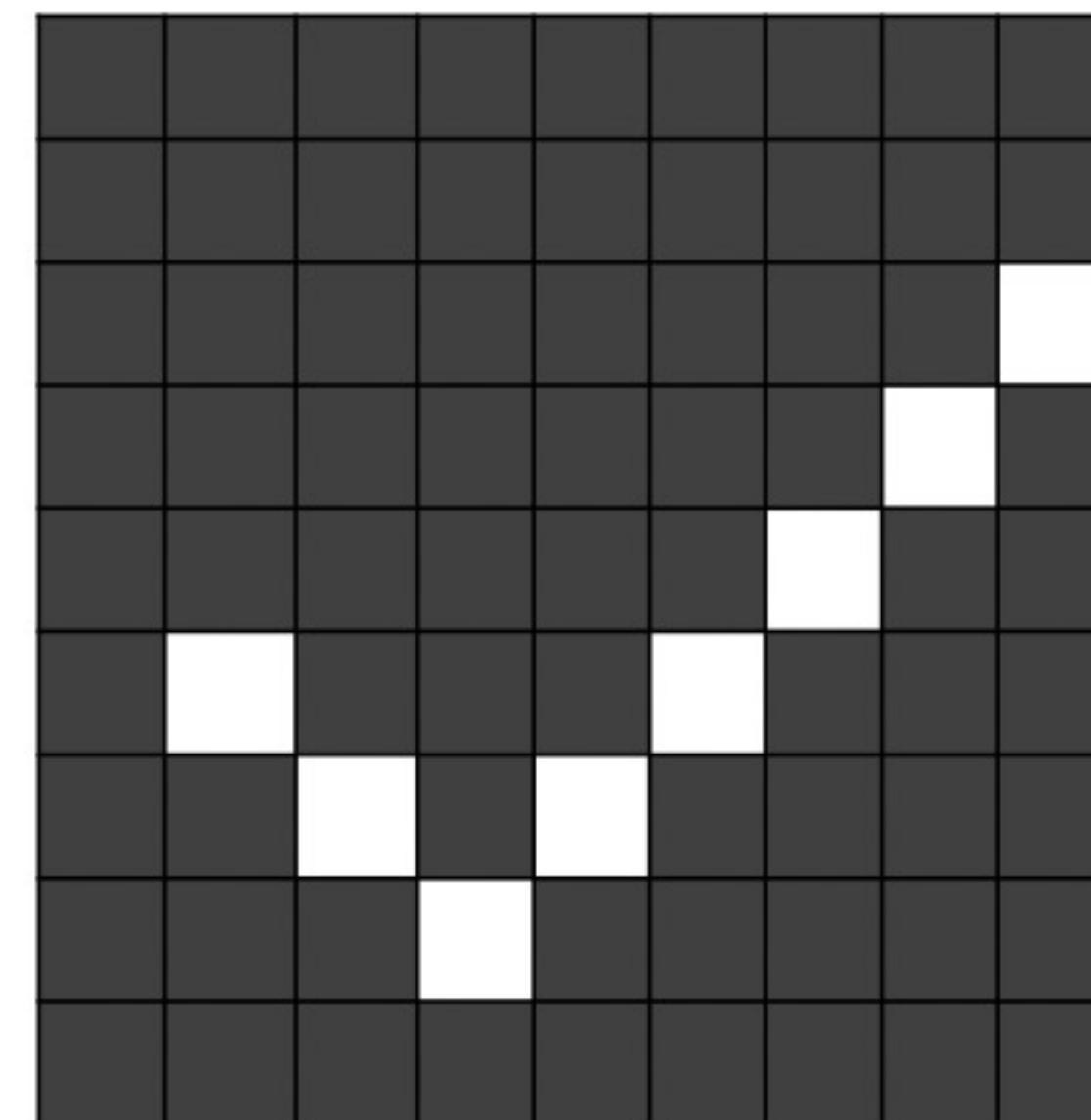
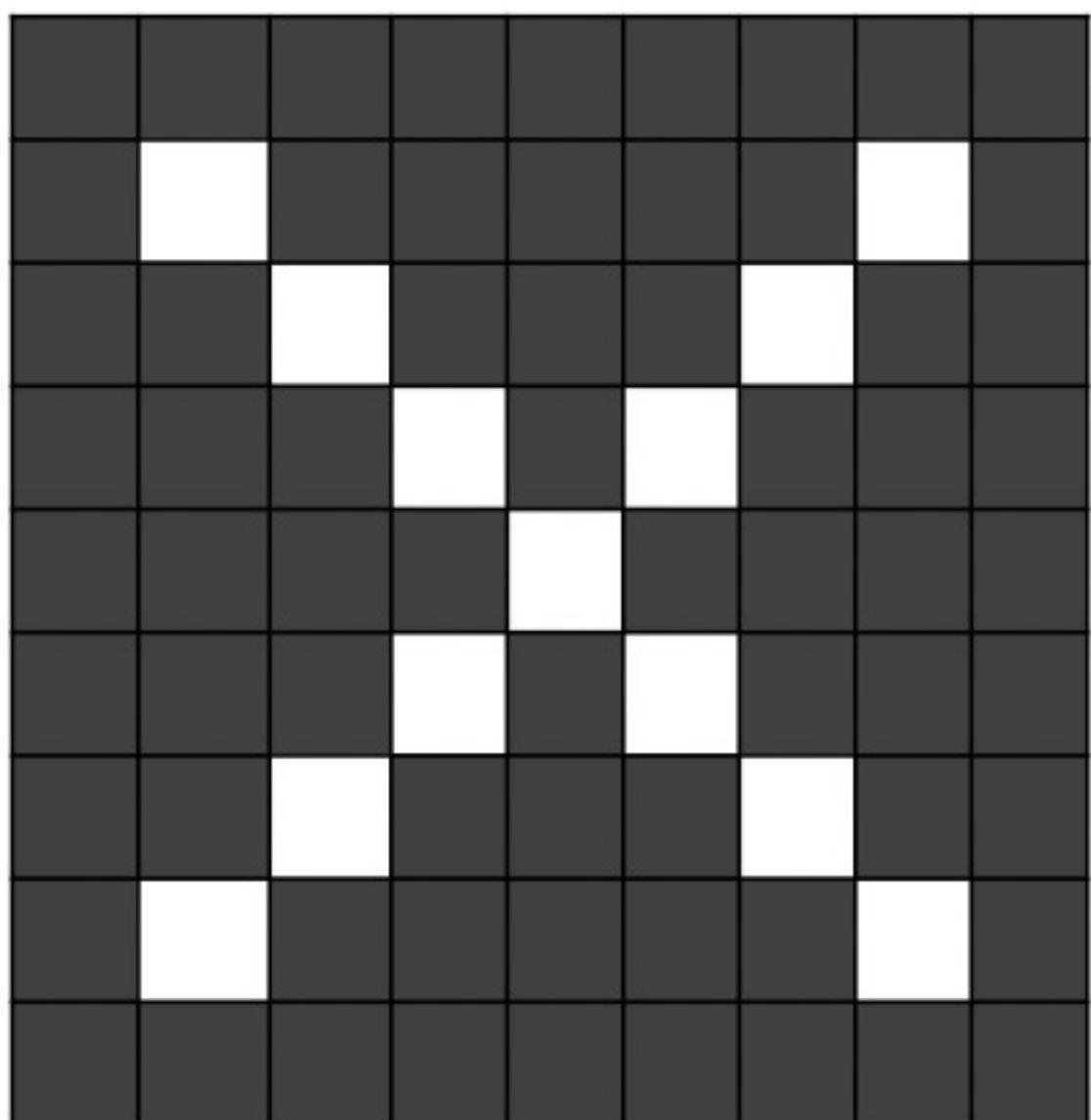
$h(x)$: activation function

Animals in the zoo

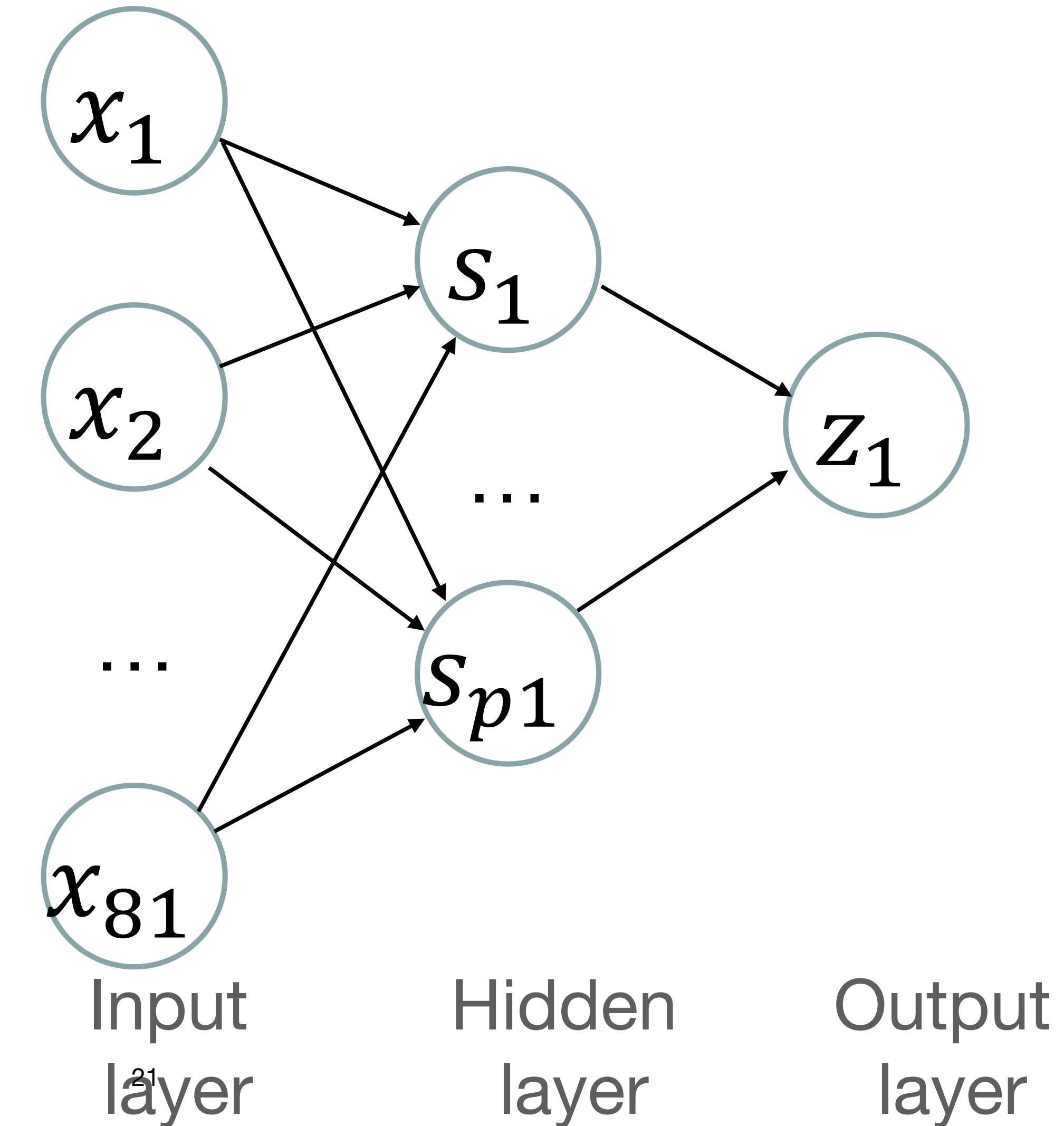
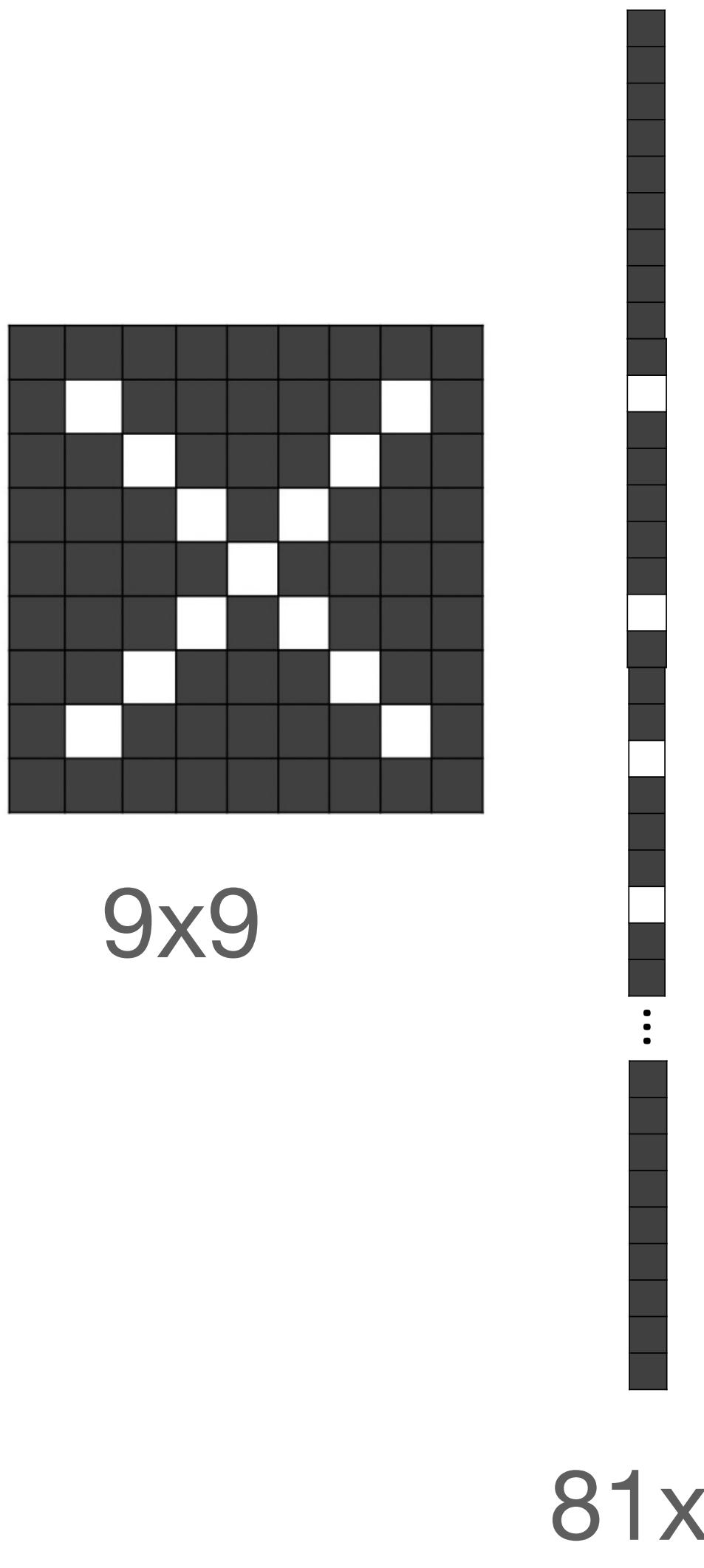


art: OpenClipartVectors
at pixabay.com (CC0)

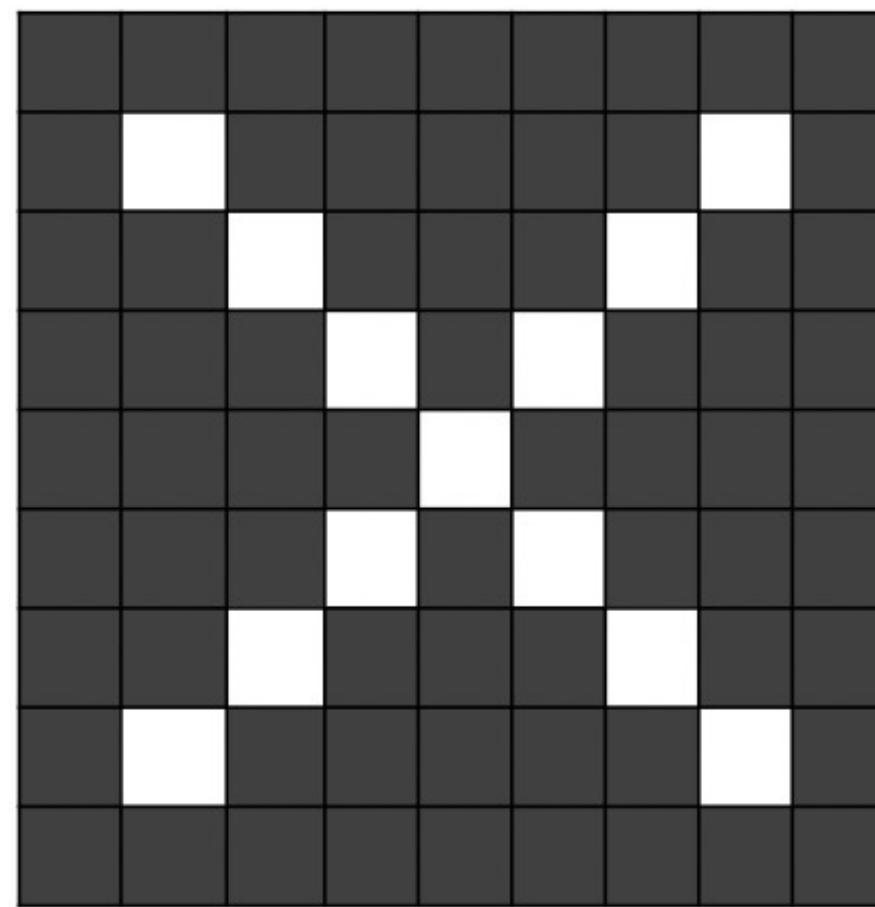
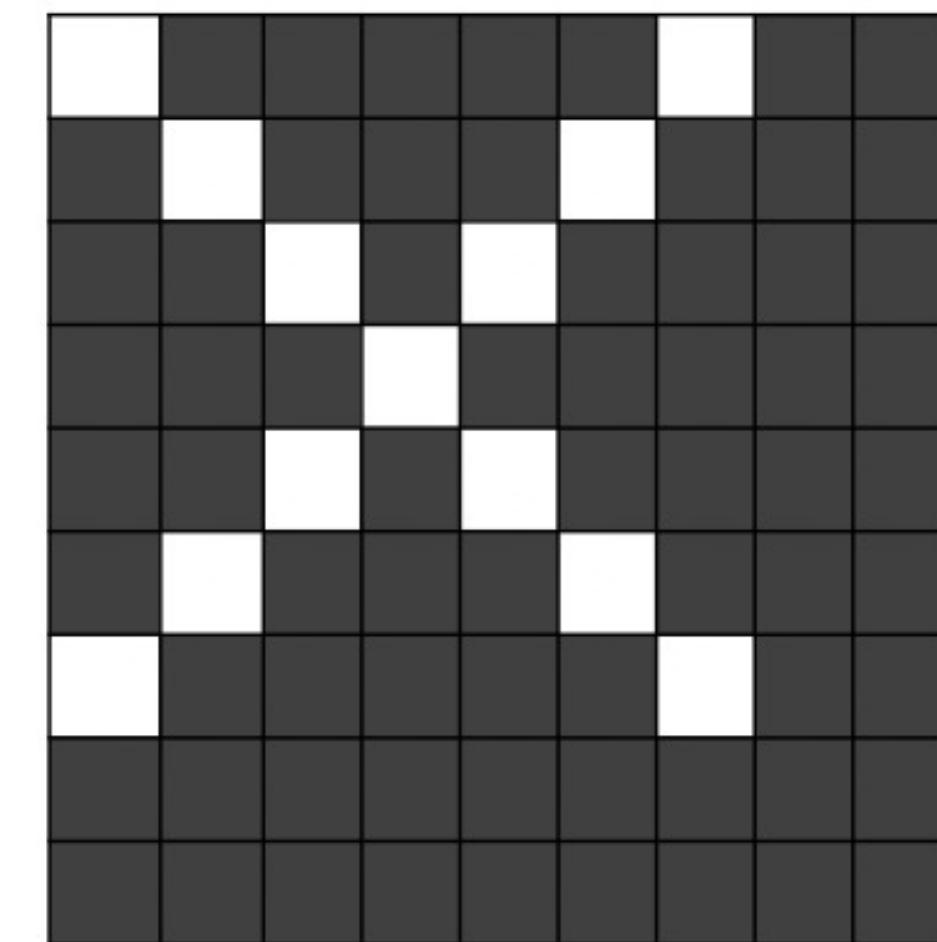
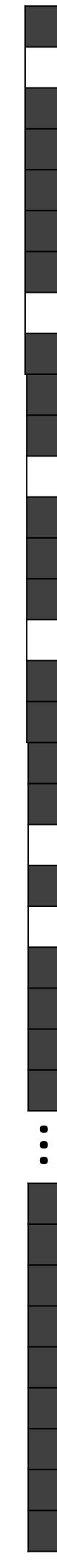
✗ **vs** ✓



Multi-layer perceptron: A fully connected network



Multi-layer perceptron: Not spatial invariant

 \neq 

Multi-layer perceptron: cannot build deep neural network

Boat tailed Grackle



Bobolink



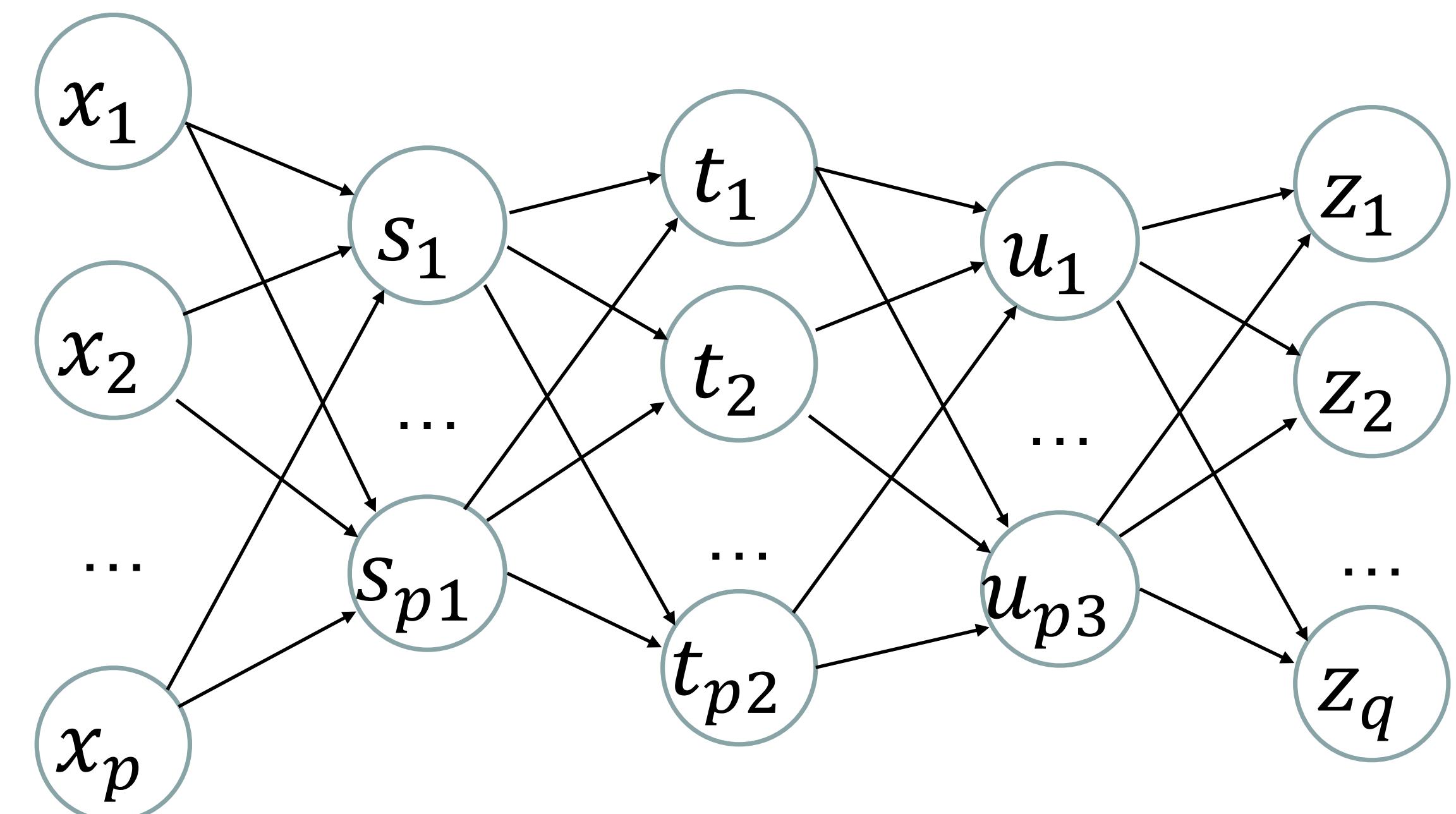
Bohemian Waxwing



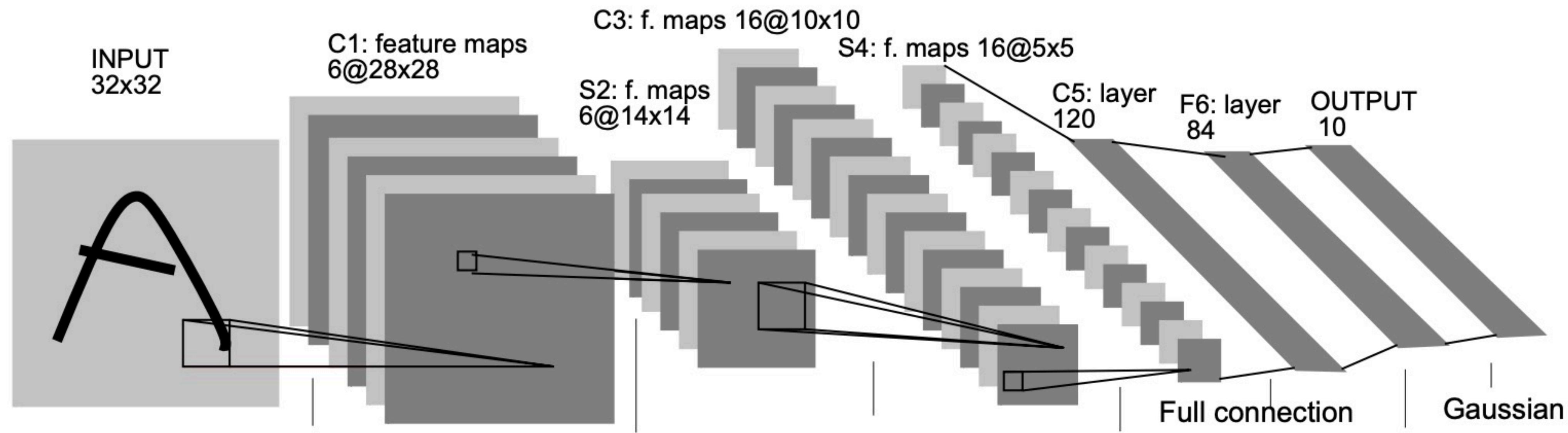
Brandt Cormorant



Brewer Blackbird



Convolutional, Max-Pooling, and Fully Connected (FC) layers

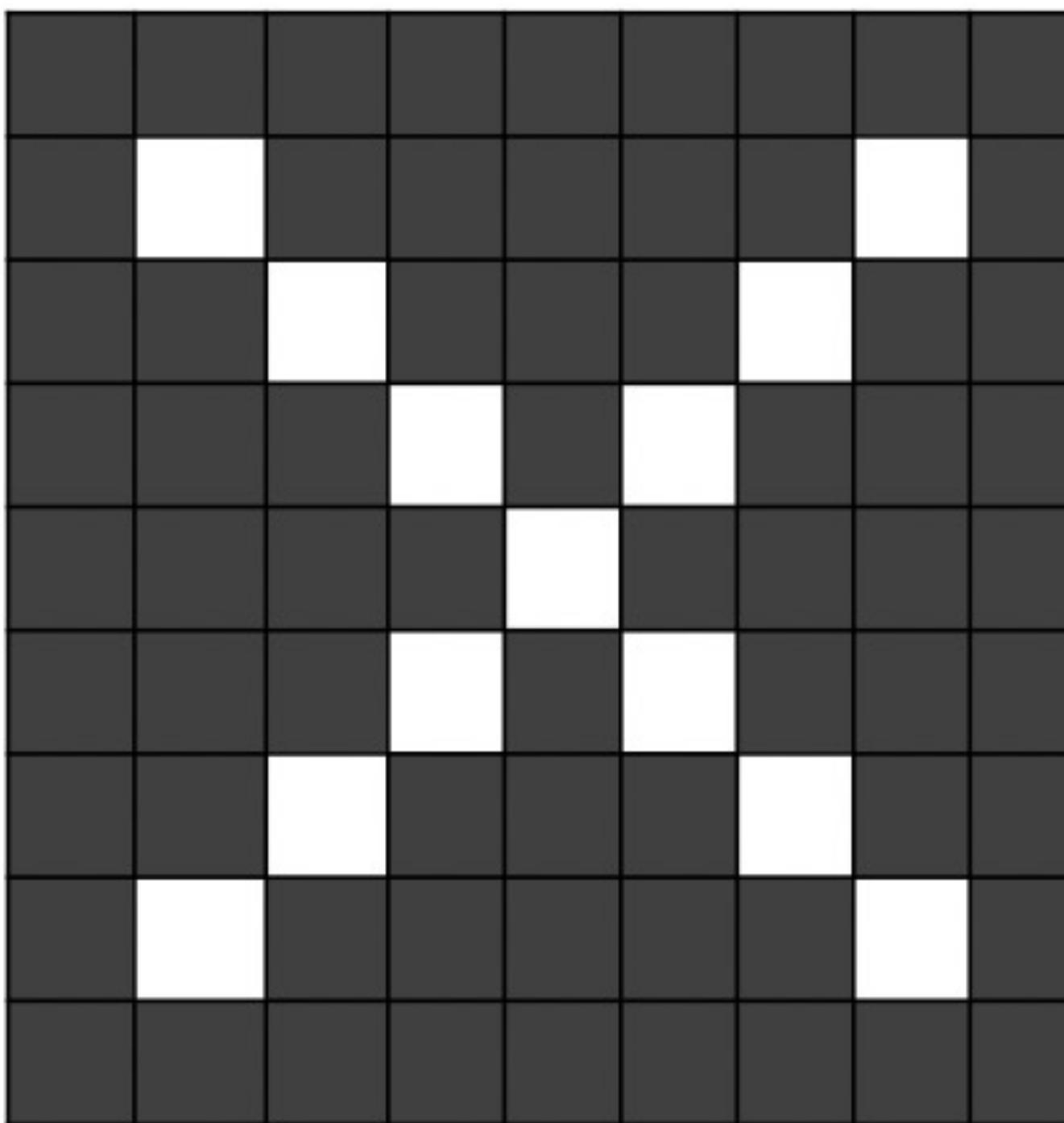


Motivation of convolution

- An image == several patches
- Combine the local features of the patches-> image feature



Identify different patterns at local patches



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

Identify different patterns

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

Element-wise
multiplication

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) \times \begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = 2$$

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right) \times \begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = 1$$

Filter (kernel)

Input and kernel have the same pattern: high response

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

Element-wise
multiplication

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) \times \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right) = 1$$

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right) \times \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{array} \right) = 2$$

Filter (kernel)

Identify different patterns

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

Element-wise
multiplication

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) \times \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = 2$$

$$\text{Sum} \left(\begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right) \times \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = 2$$

Filter (kernel)

Different kernels identify different patterns

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

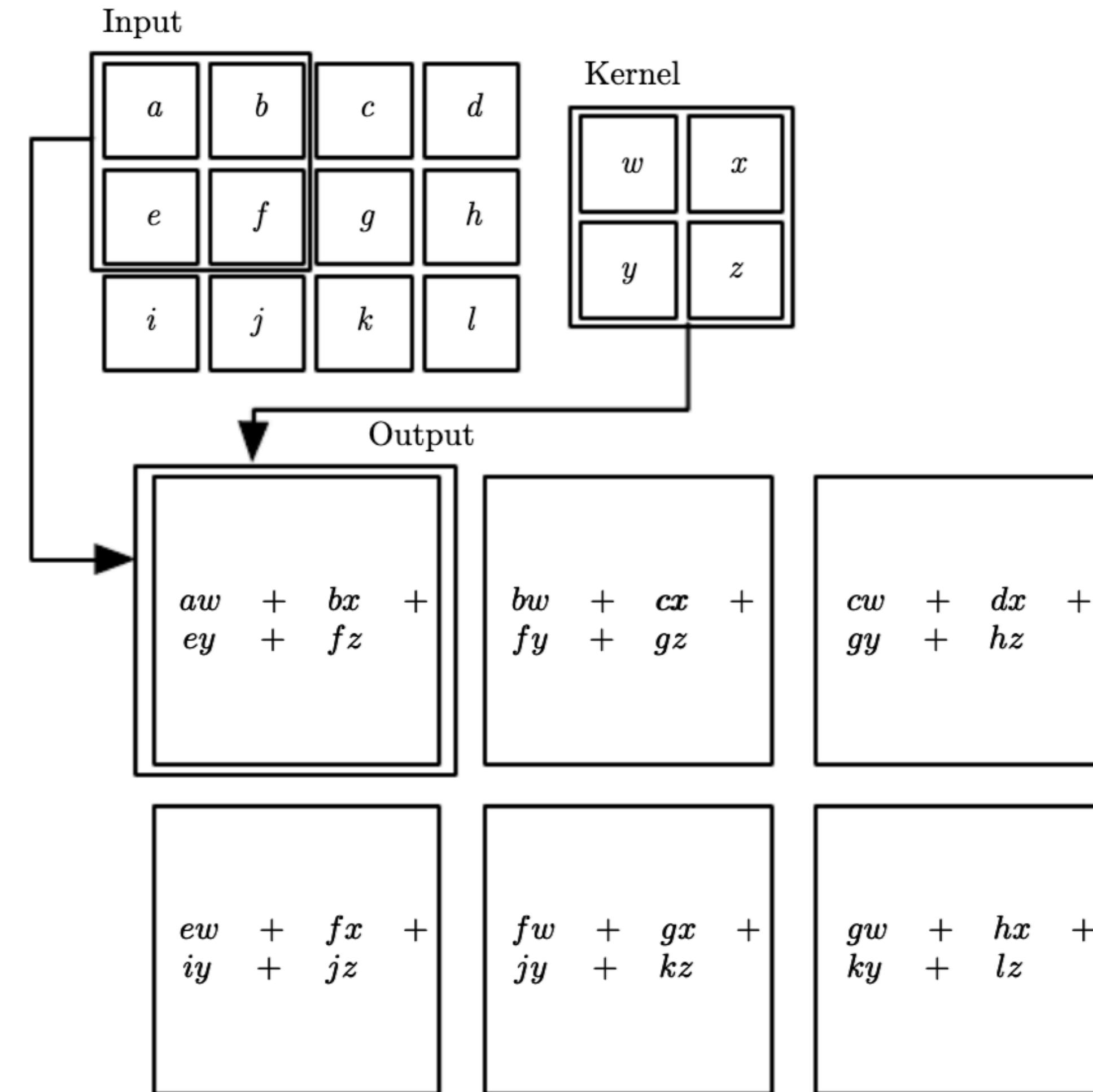
Element-wise
multiplication

$$\text{Sum} \left(\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) \times \left(\begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right) = 2$$
$$\text{Sum} \left(\begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right) \times \left(\begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right) = 5$$

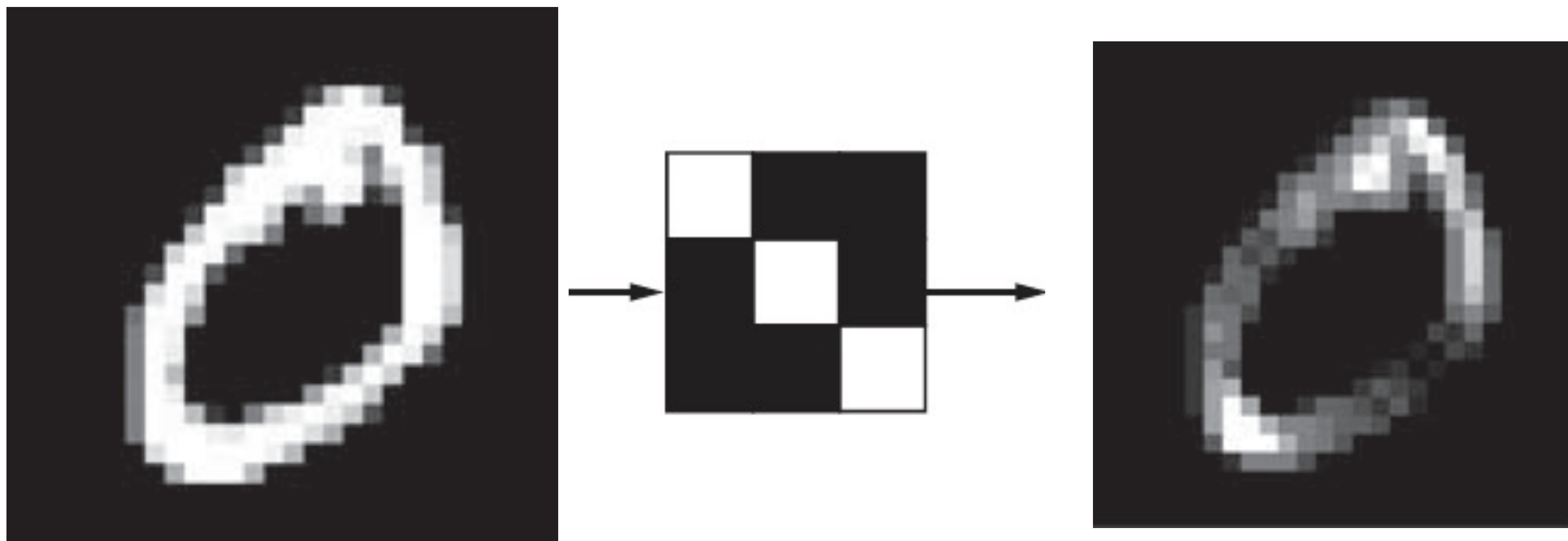
Filter (kernel)

Convolution on 2D

Use kernel to perform element-wise multiplication and sum for every local patch



Response map (Feature map)



Input

kernel

Feature map: 2D map of the presence of a pattern at different locations in an input

Different kernels identify different patterns: use multiple filters in each layer

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

1	0	0
0	1	0
0	0	1

0	0	1
0	1	0
1	0	0

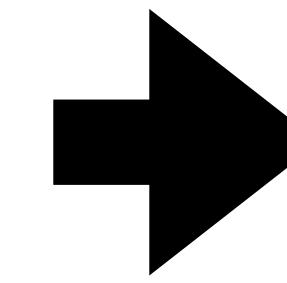
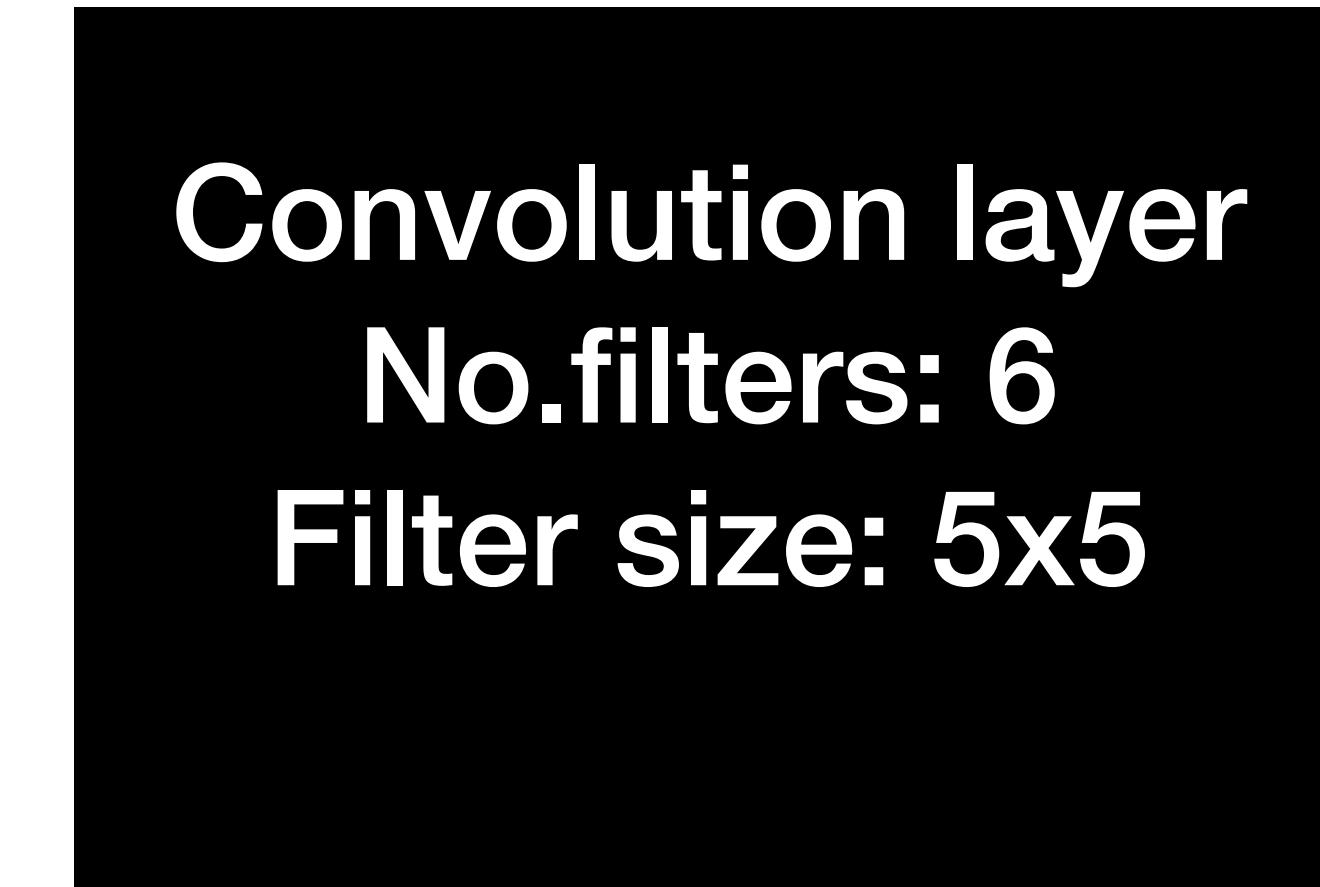
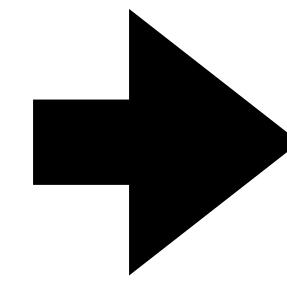
0	1	0
1	0	1
0	0	0

The number of filters decides the number of output feature maps

Two key hyperparameters in Convolution

Filter (kernel) size: Size of the patches extracted from the inputs

Number of filters: Depth (channel) of the output feature map



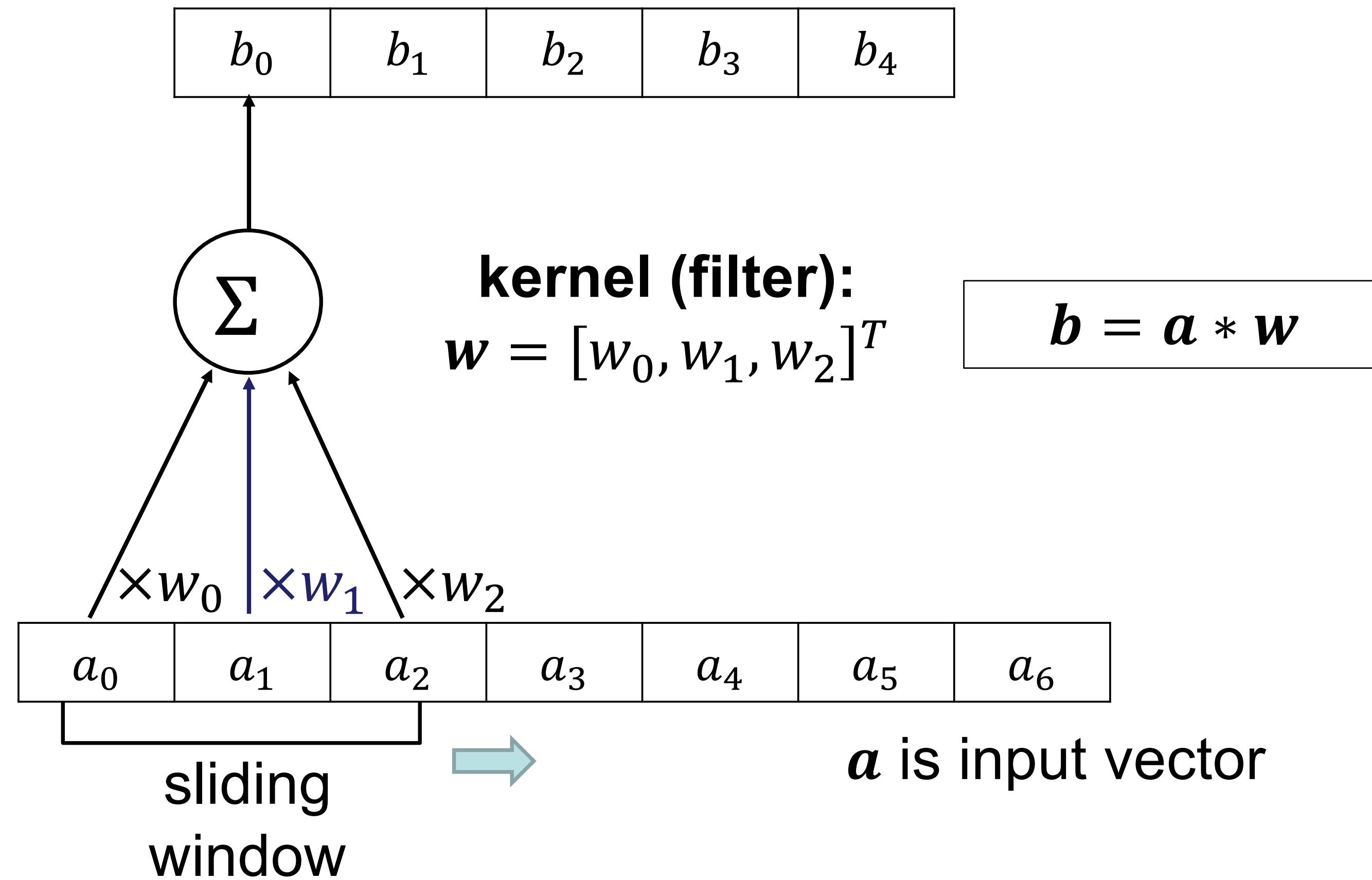
32x32x1

Input: 1 channel

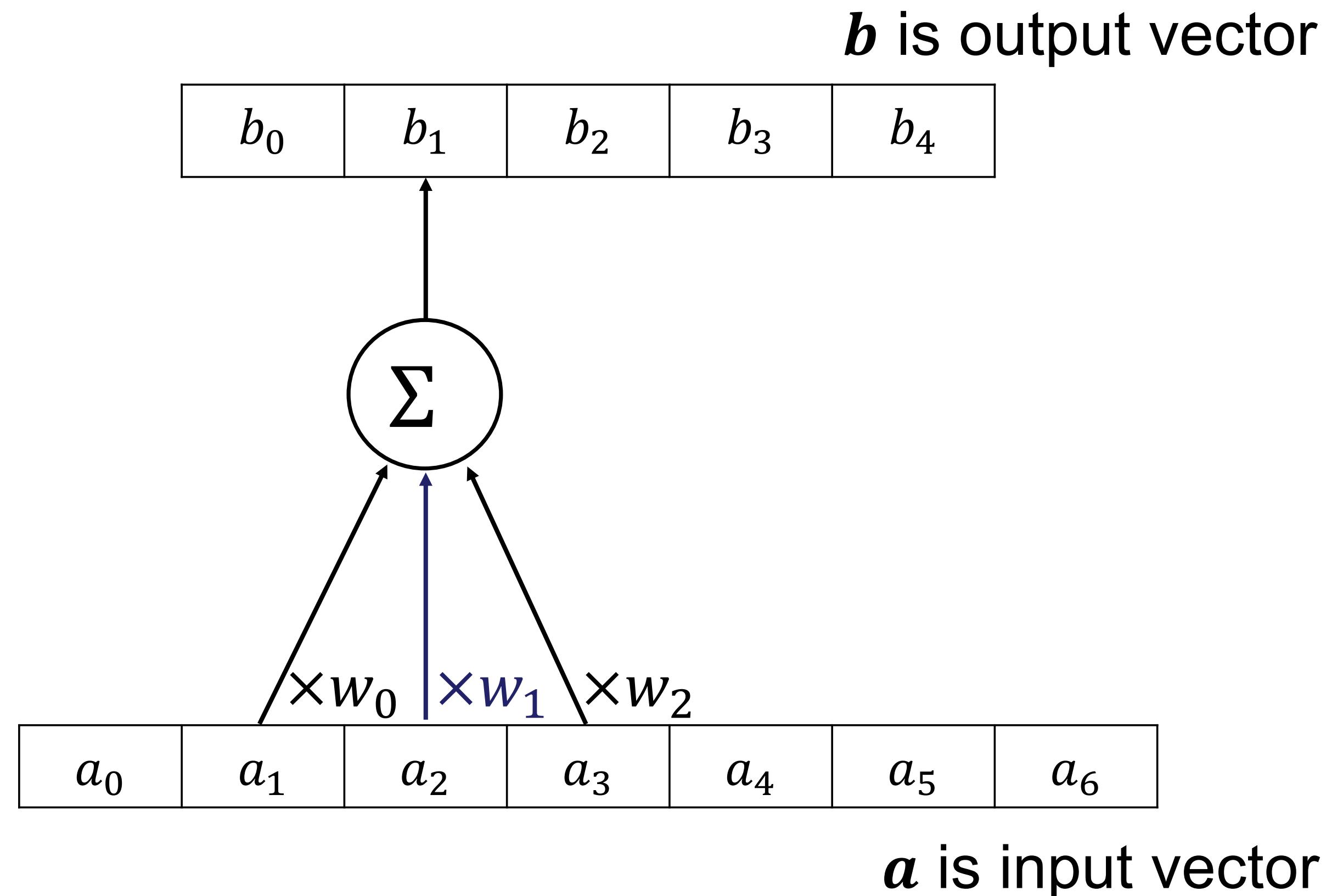
output: 6 channel

Convolution on 1D

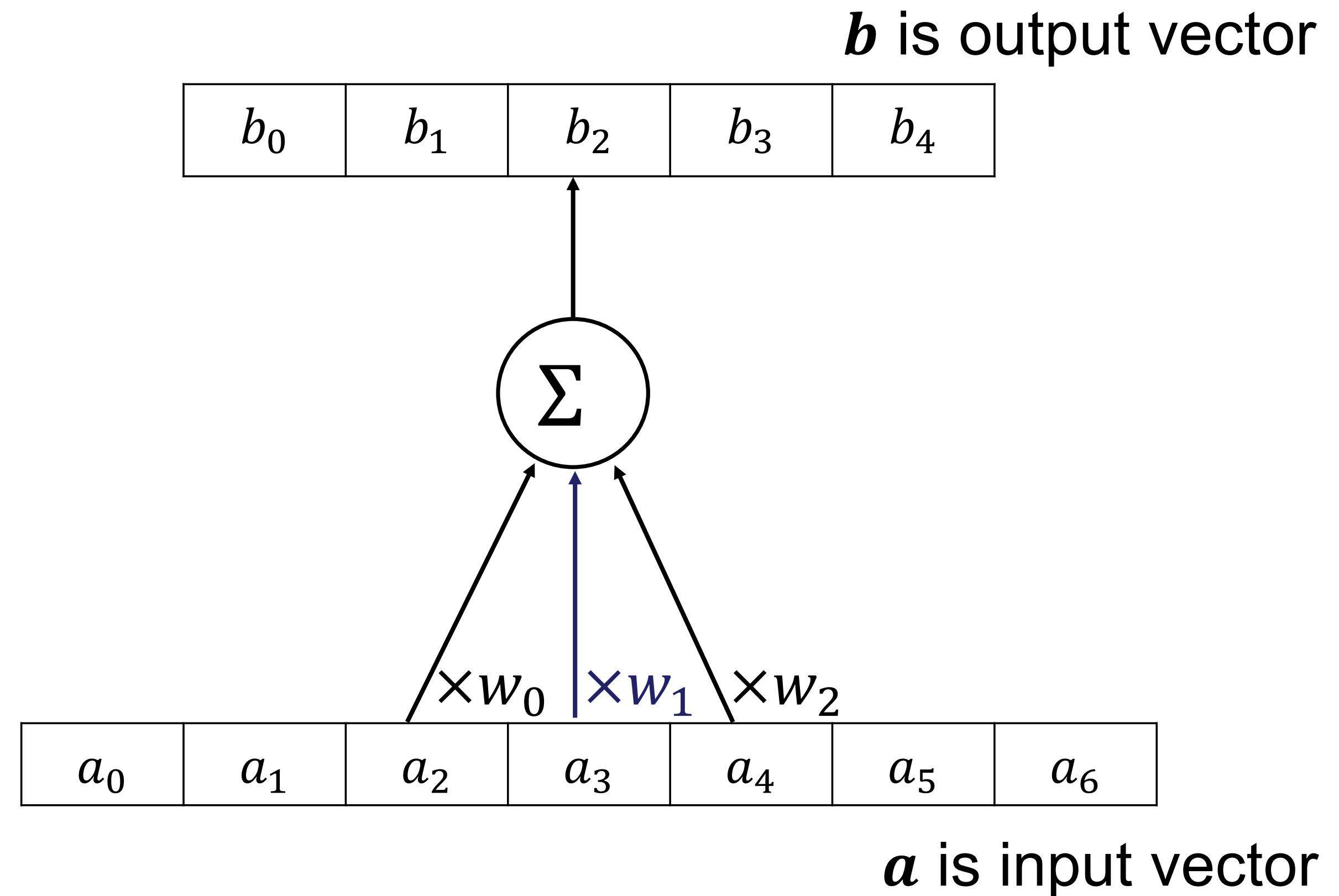
\mathbf{b} is output vector (feature map)



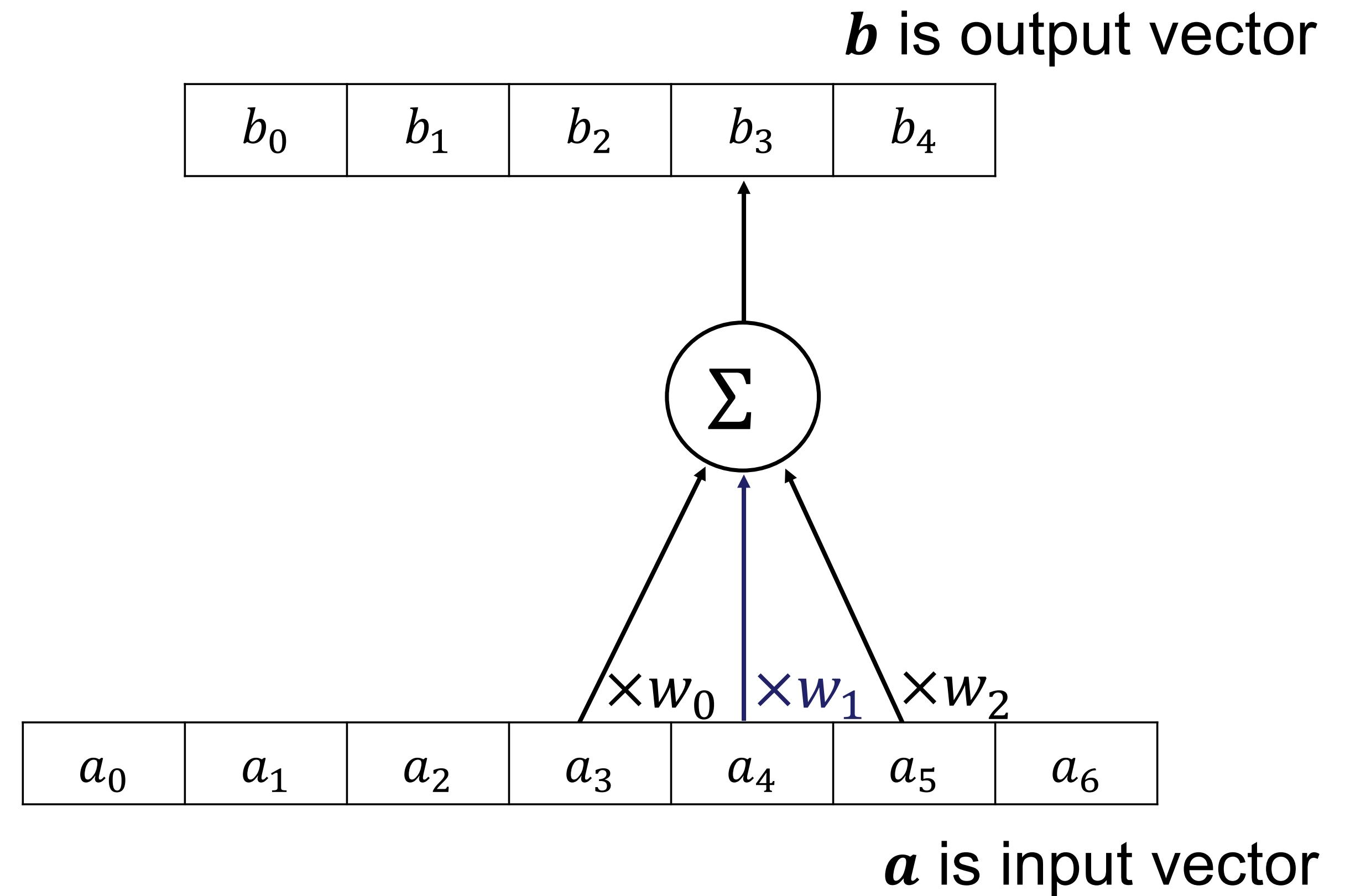
Convolution on 1D



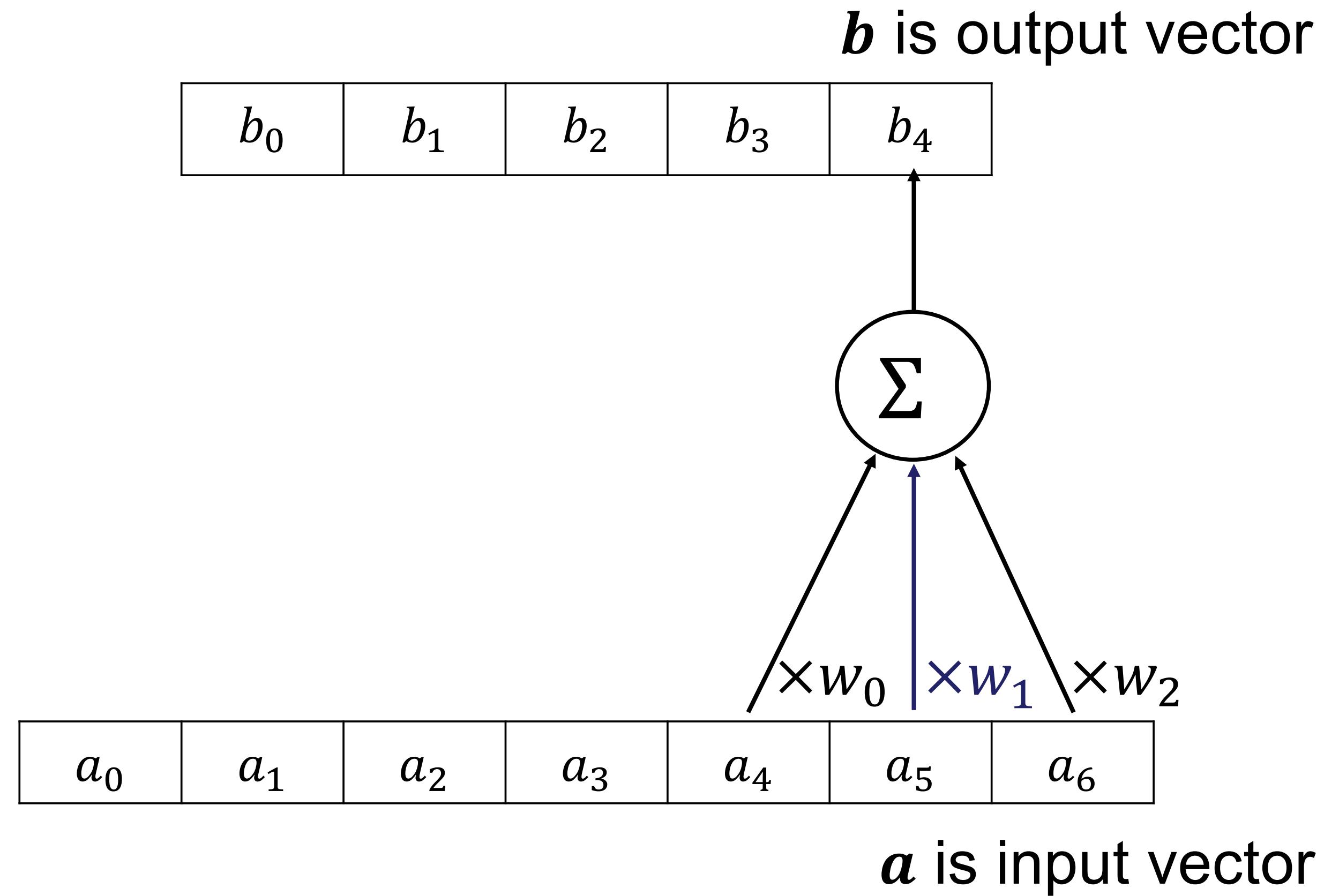
Convolution on 1D



Convolution on 1D



Convolution on 1D

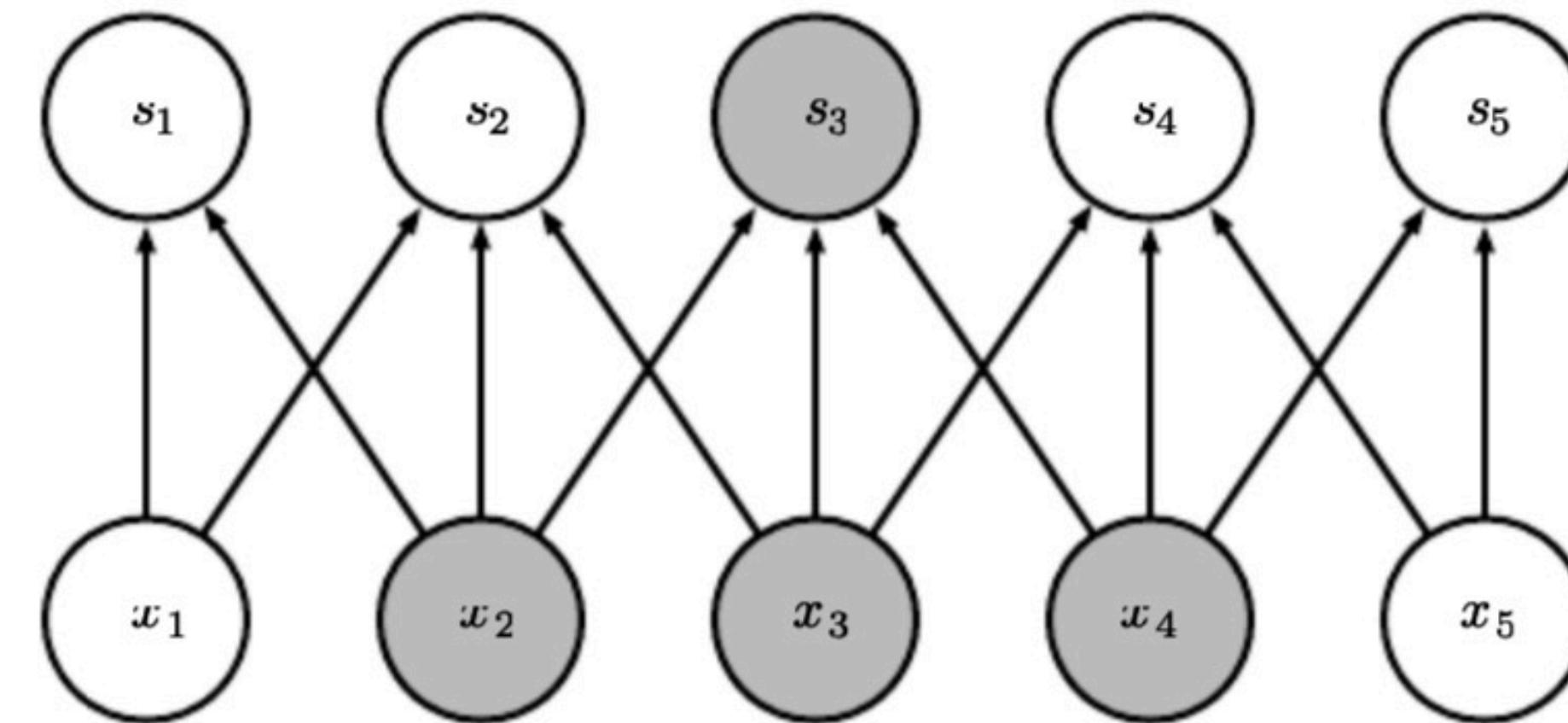


Advantage: learn translation-invariant pattern



Advantage: weight sharing and sparse connection

Convolutional layer:



Fully connected layer:
Each arrow is a
weight (no sharing)

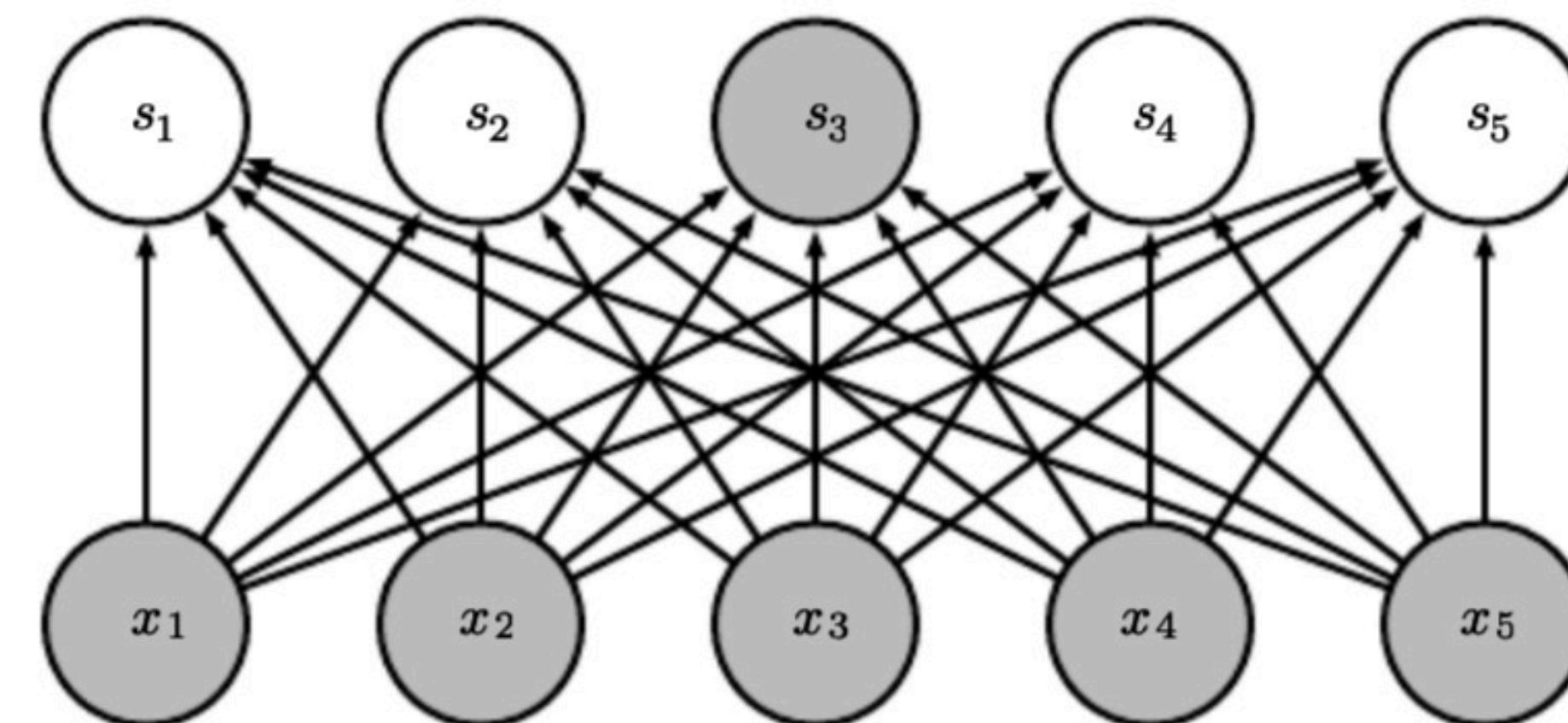


Figure 9.3 in Deep learning by Ian Goodfellow and Yoshua Bengio and Aaron Courville

Advantage: learn hierarchical pattern

More layers: larger size of receptive field
(larger window of the input is seen)

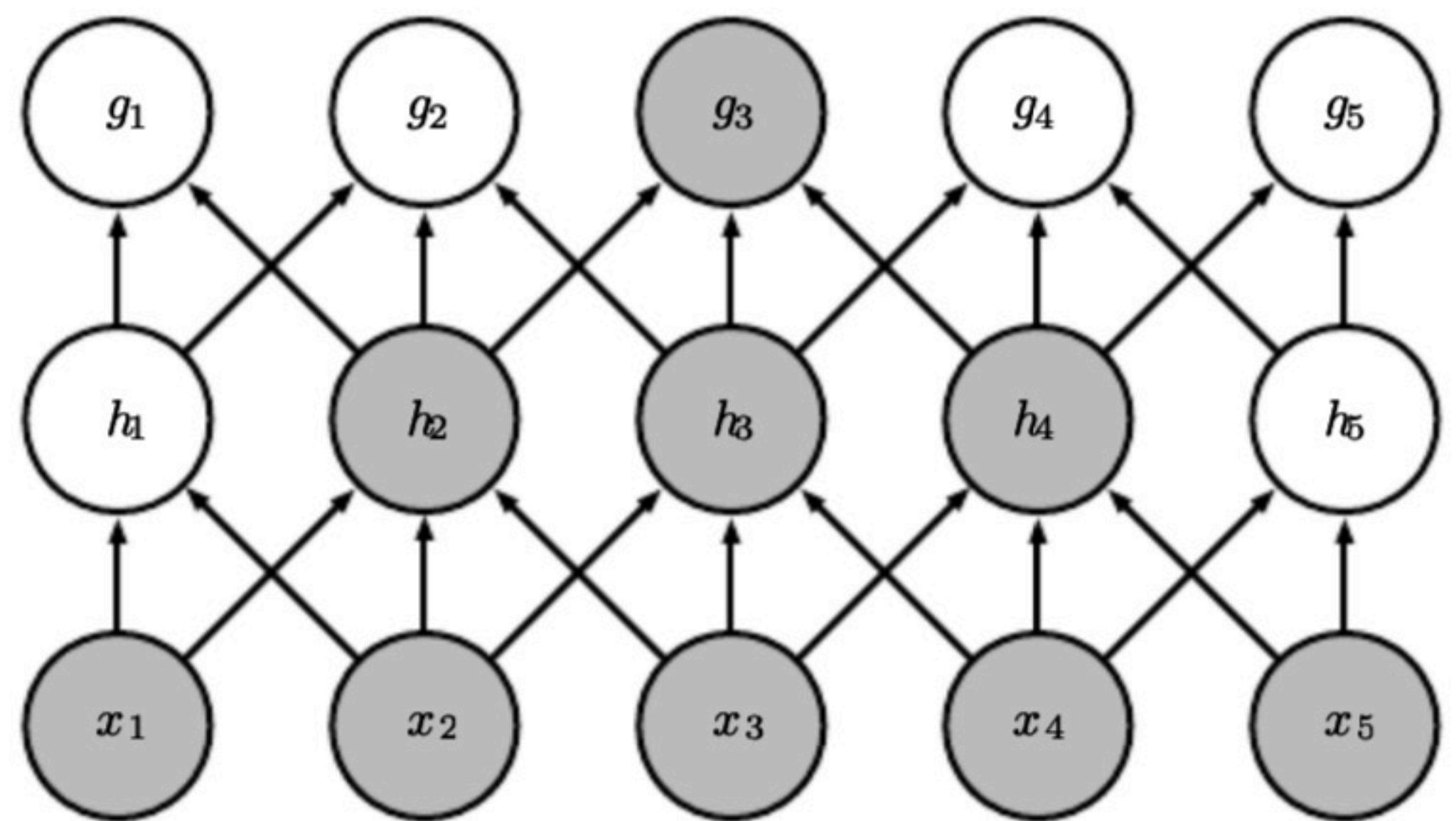


Figure 9.4 in Deep learning by Ian Goodfellow and Yoshua Bengio and Aaron Courville

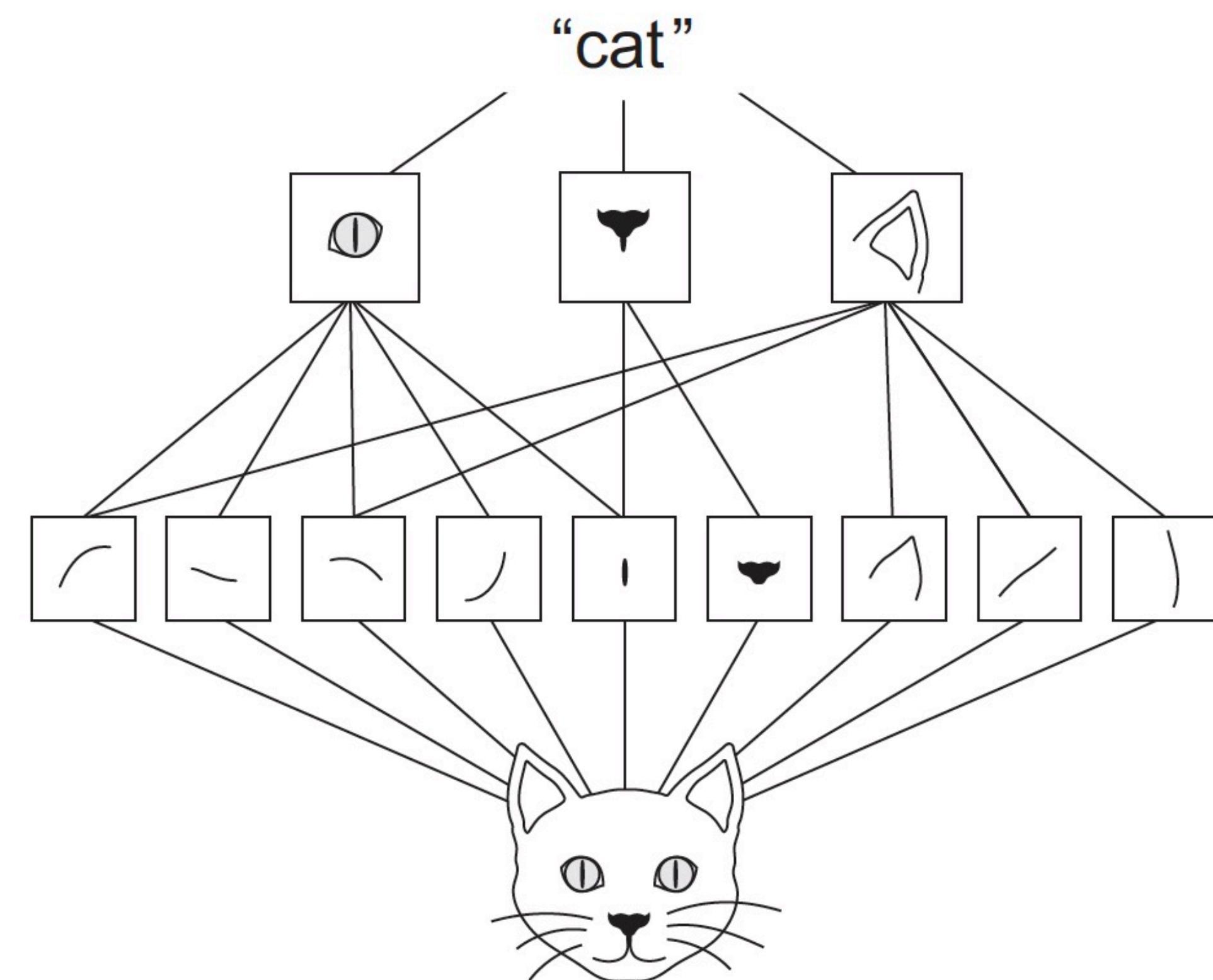
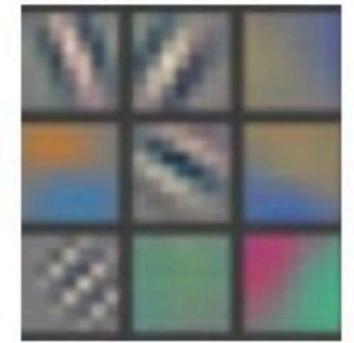


Figure 5.2 in Deep learning with python by Francois Chollet

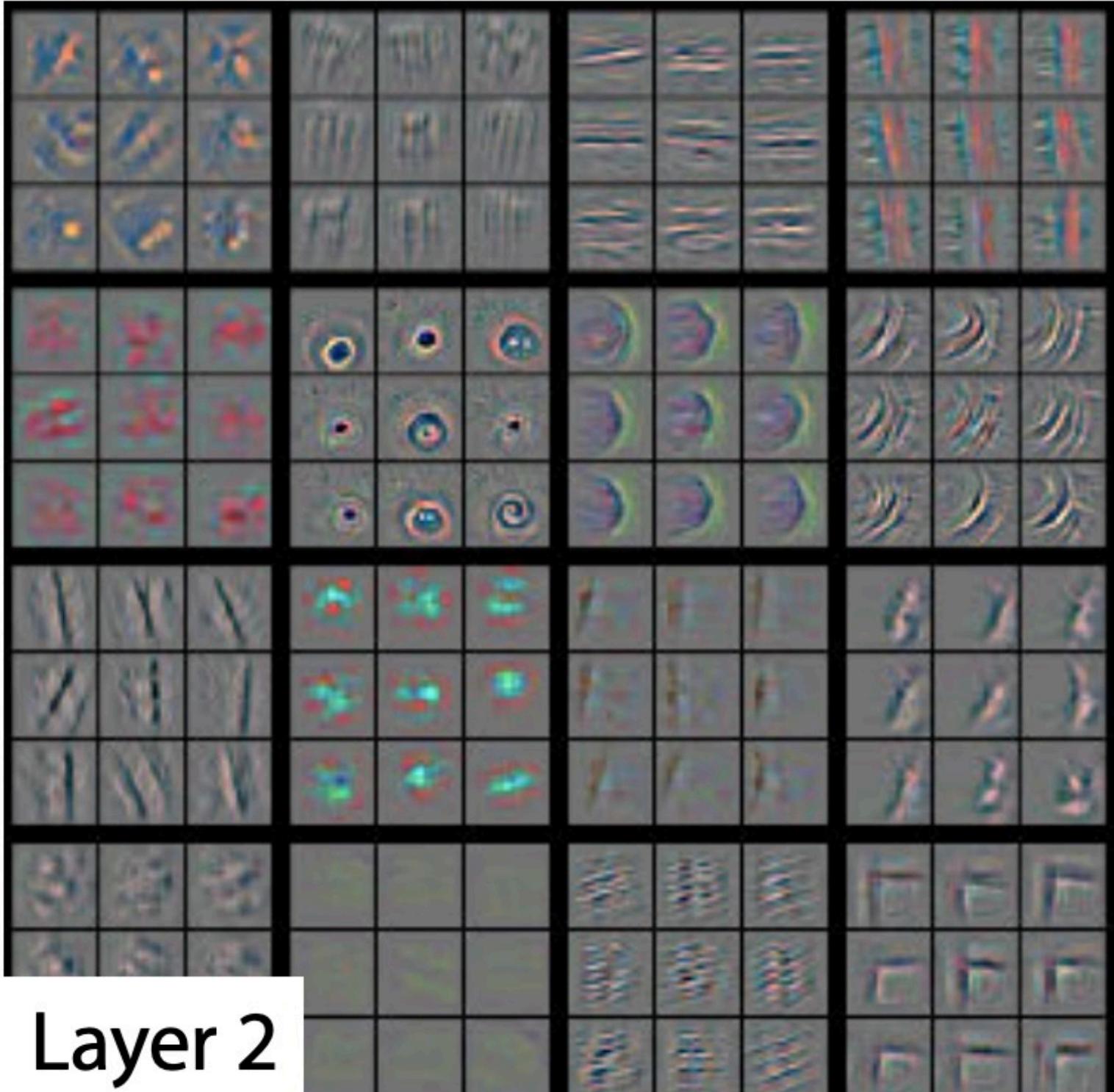
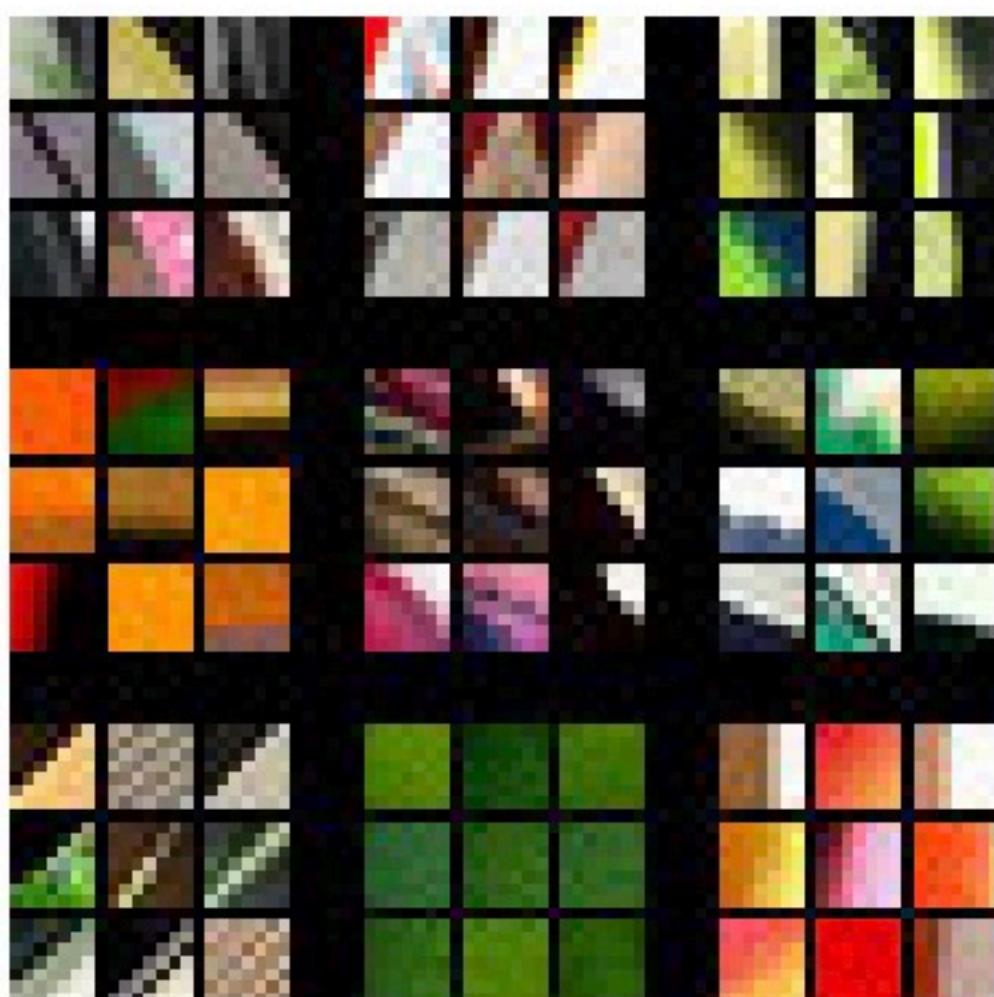
Introduction

Neural Networks

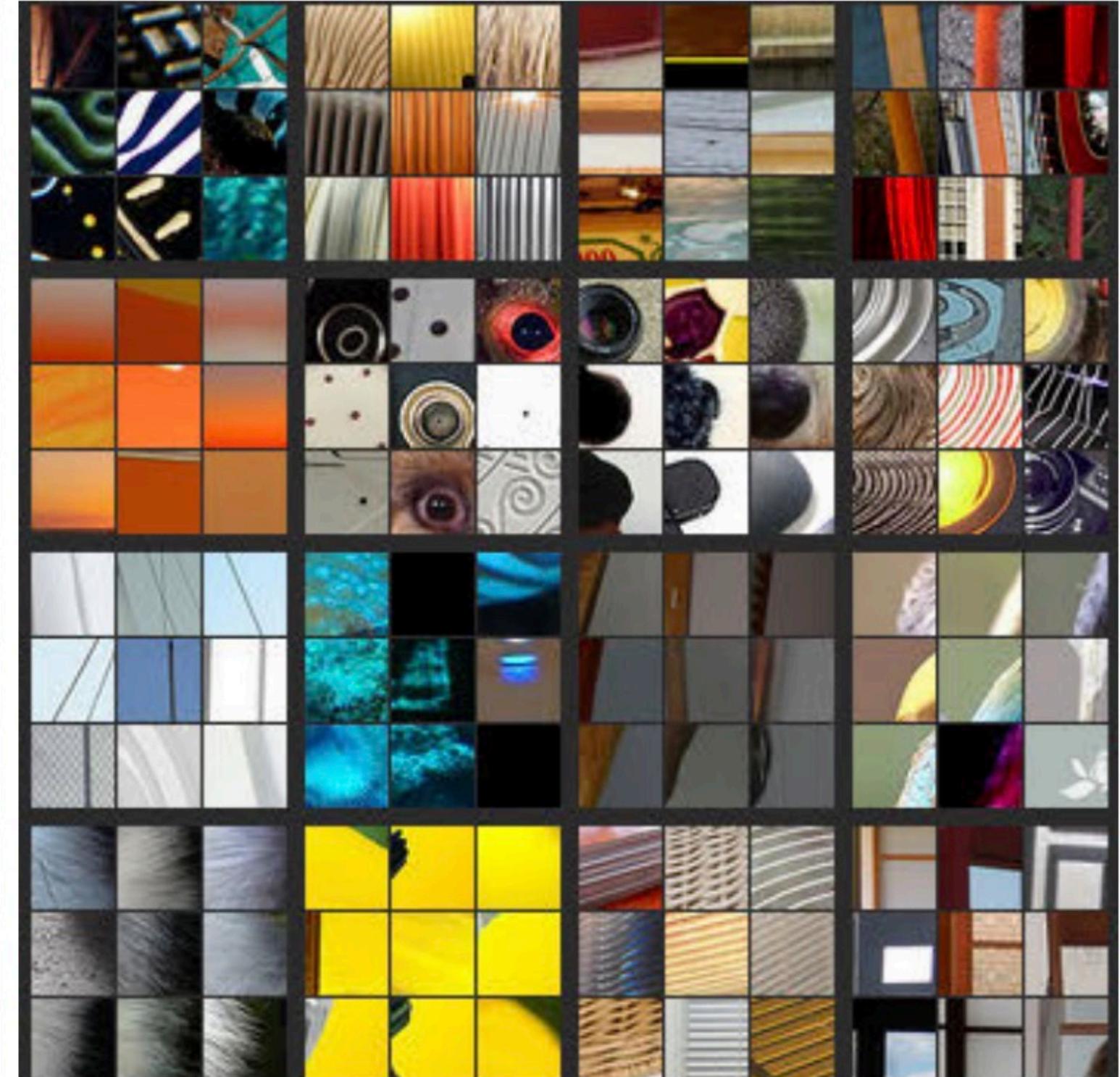
Convolutional Neural Networks



Layer 1



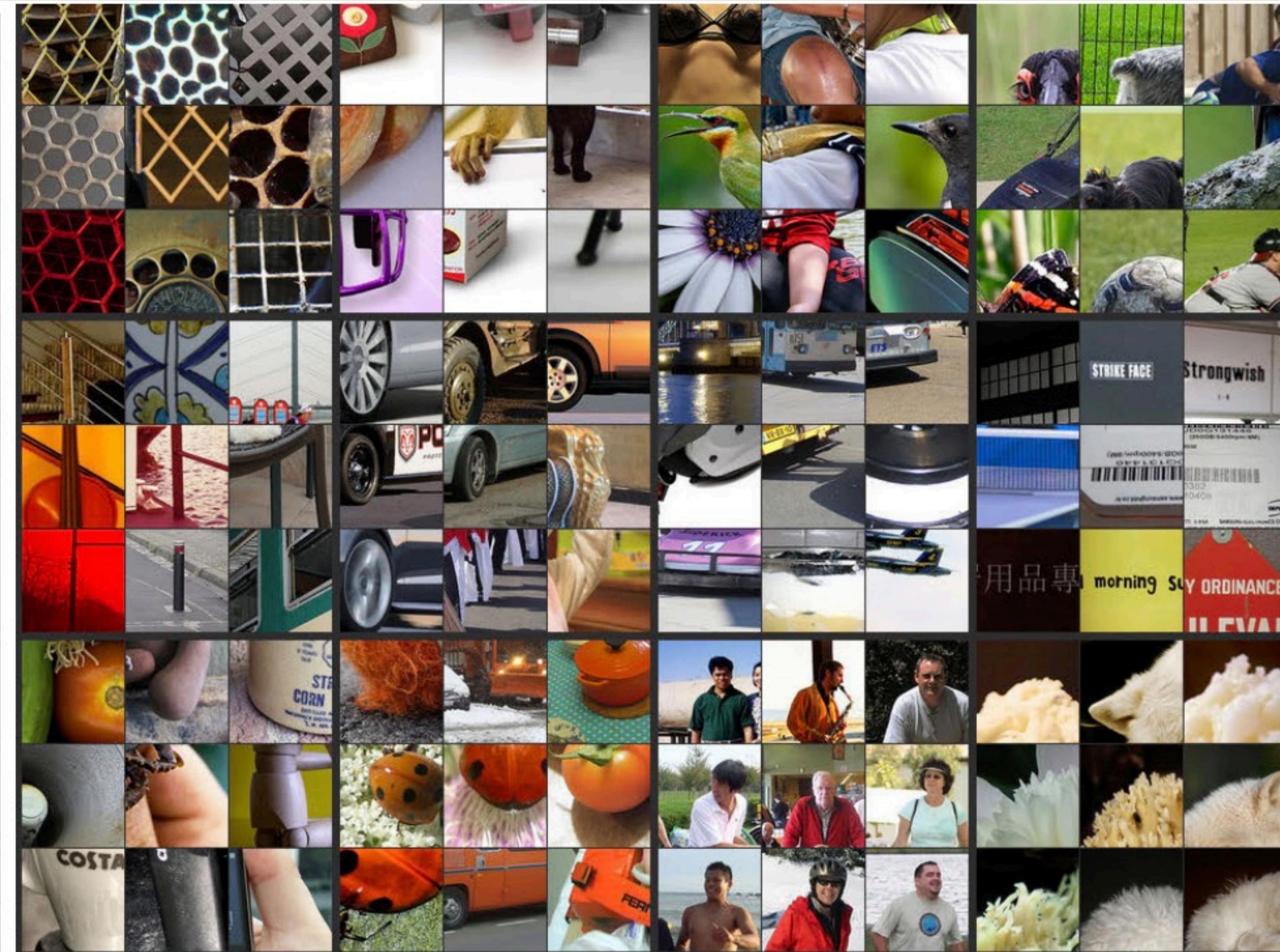
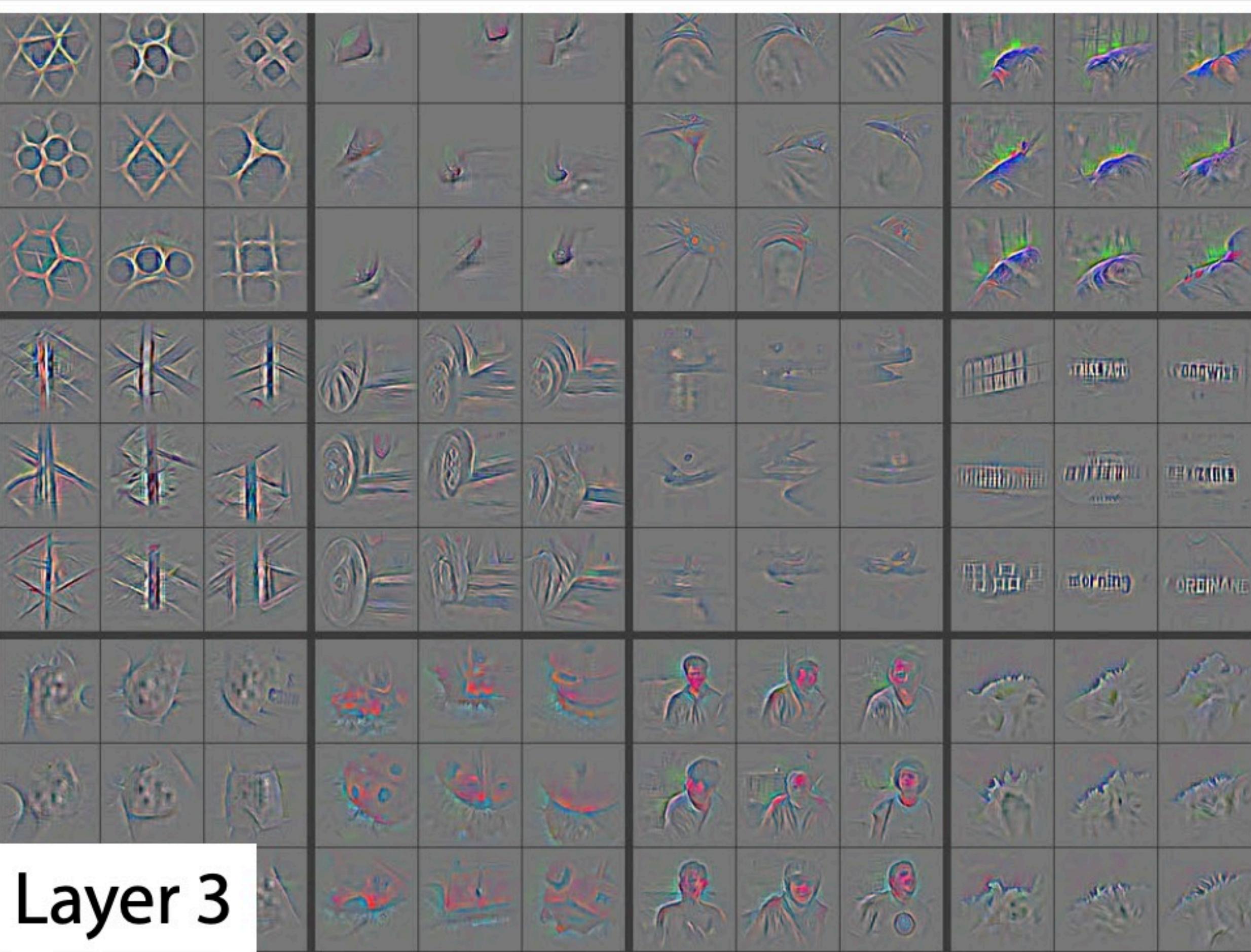
Layer 2



Introduction

Neural Networks

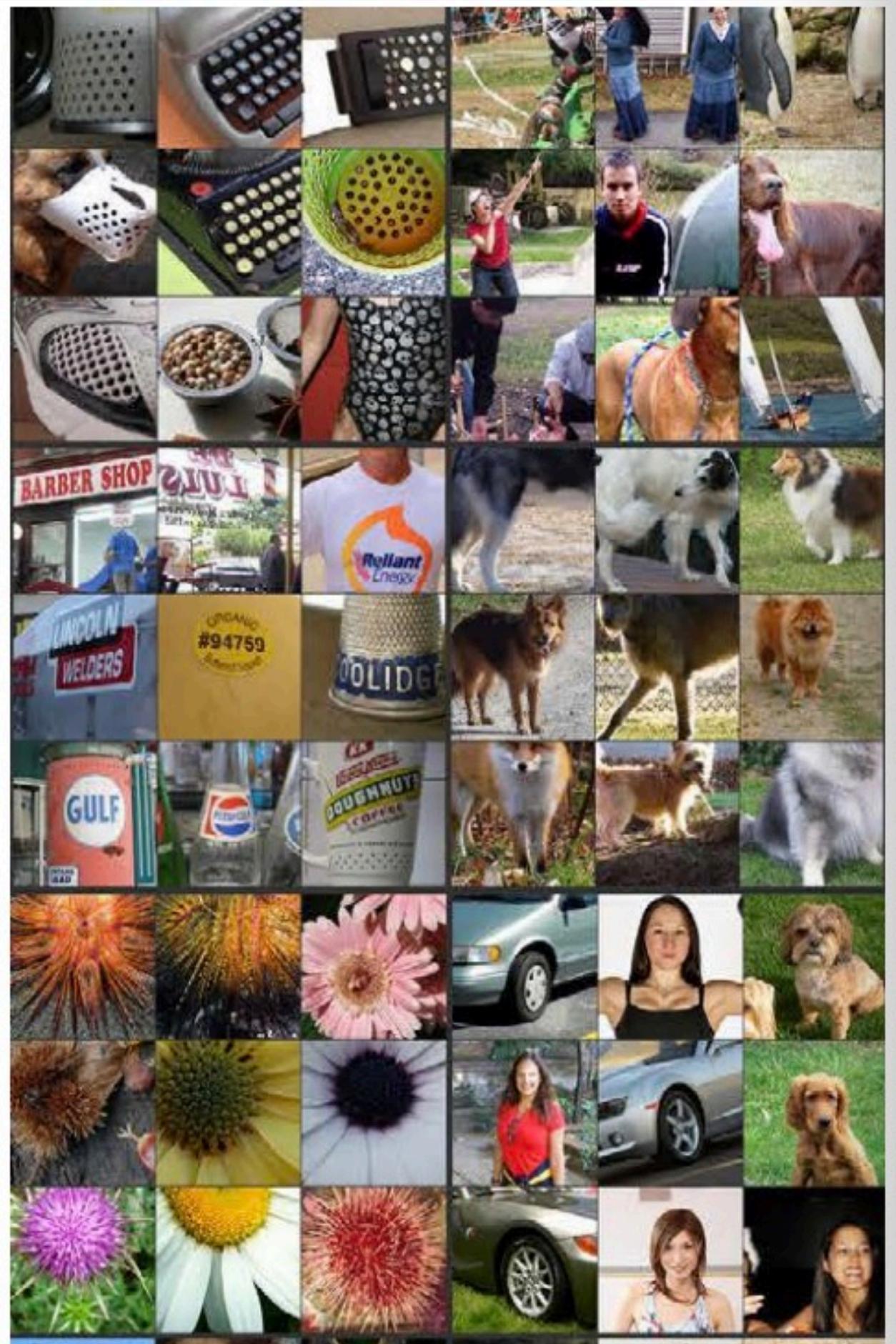
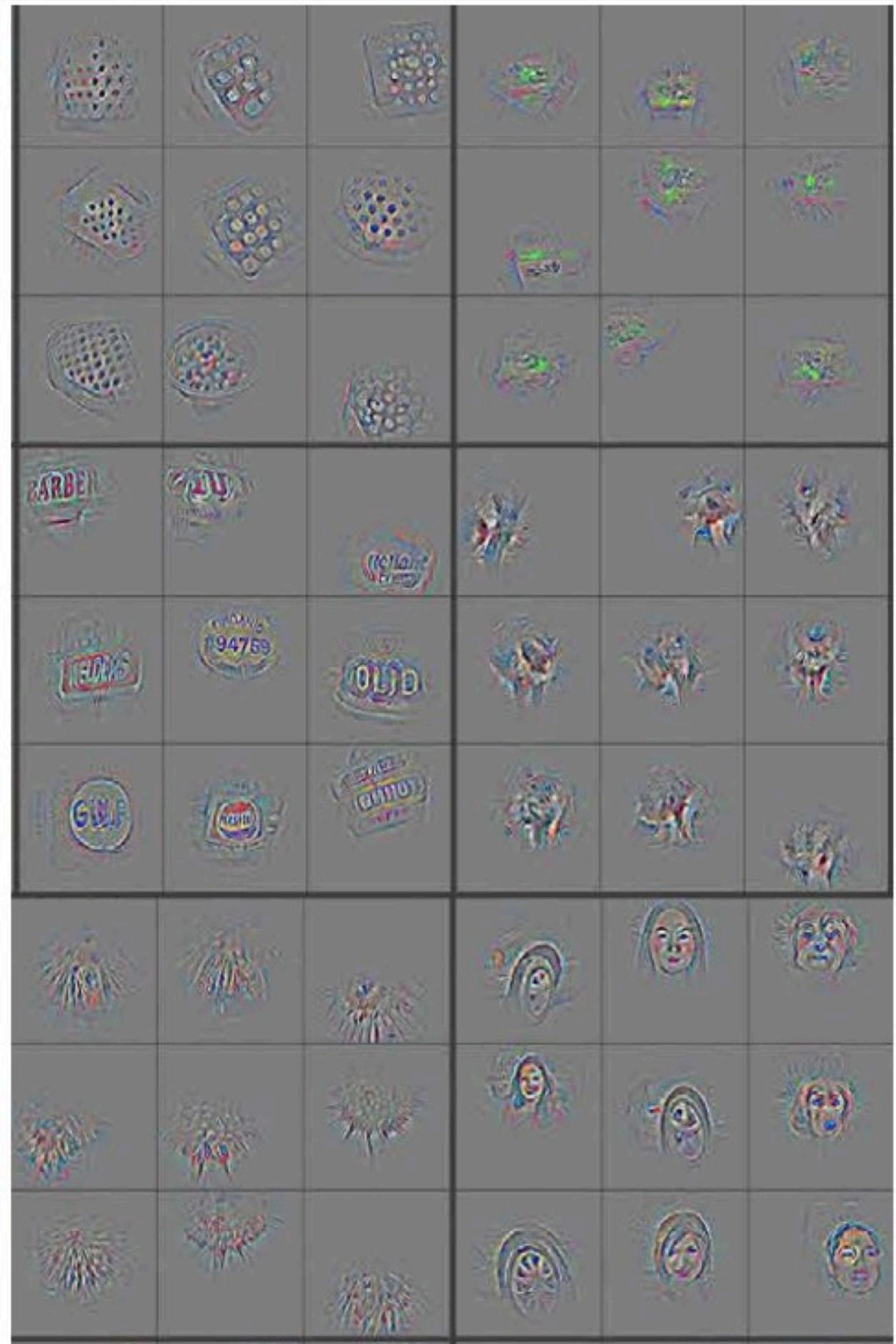
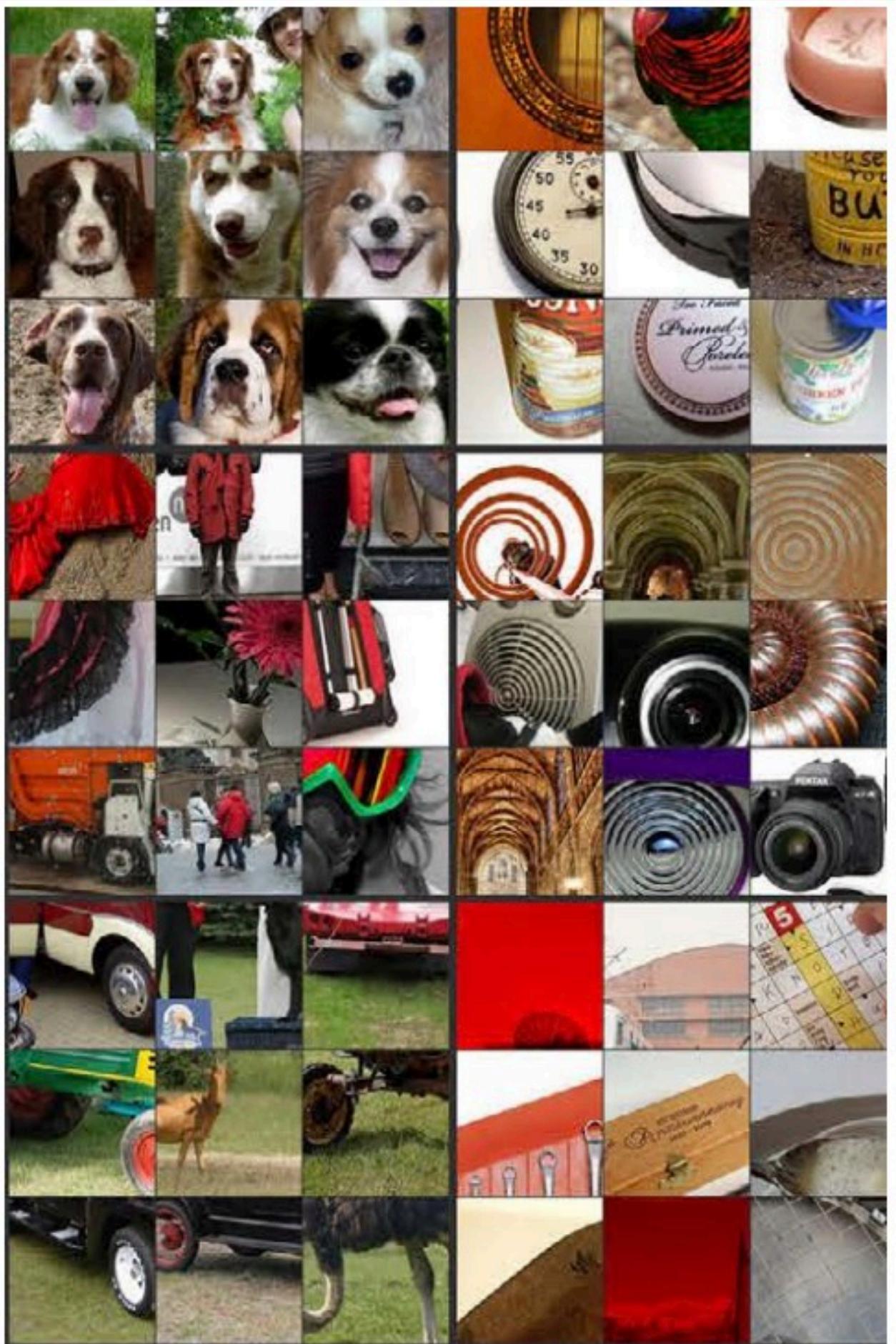
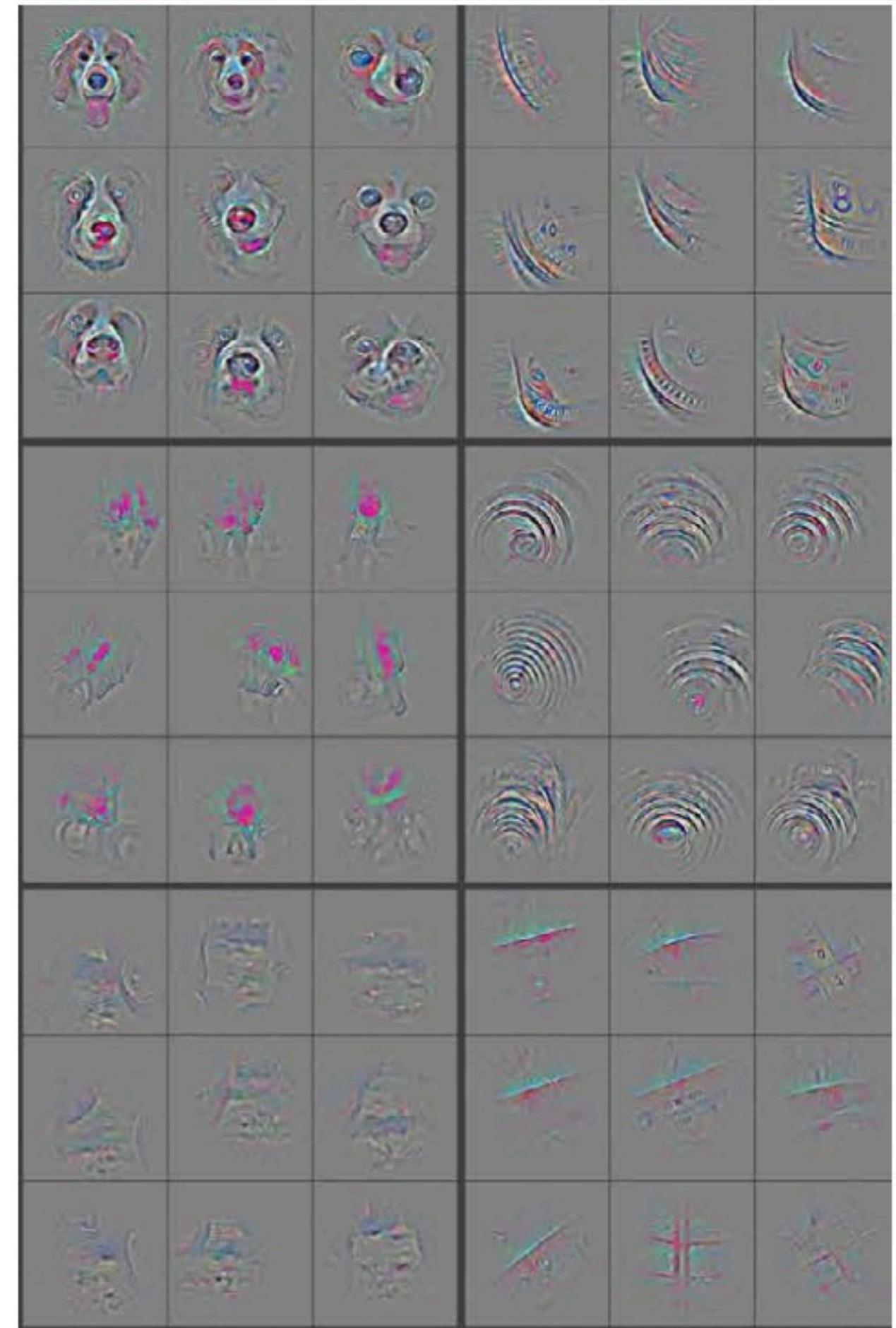
Convolutional Neural Networks



Introduction

Neural Networks

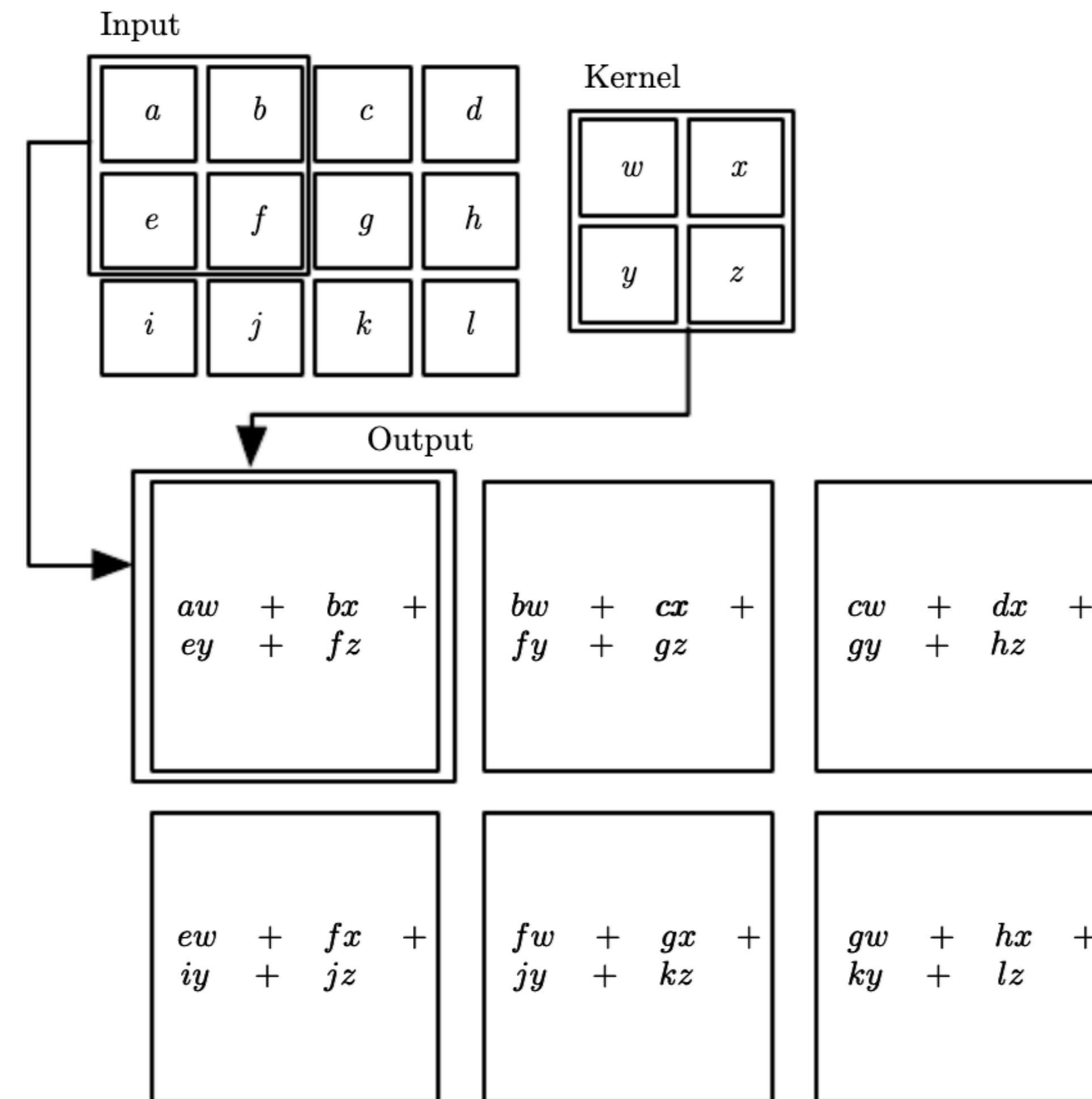
Convolutional Neural Networks



Layer 4

Layer 5

output size \neq input size



Padding

adding an appropriate number of rows and columns on each side of the input feature map

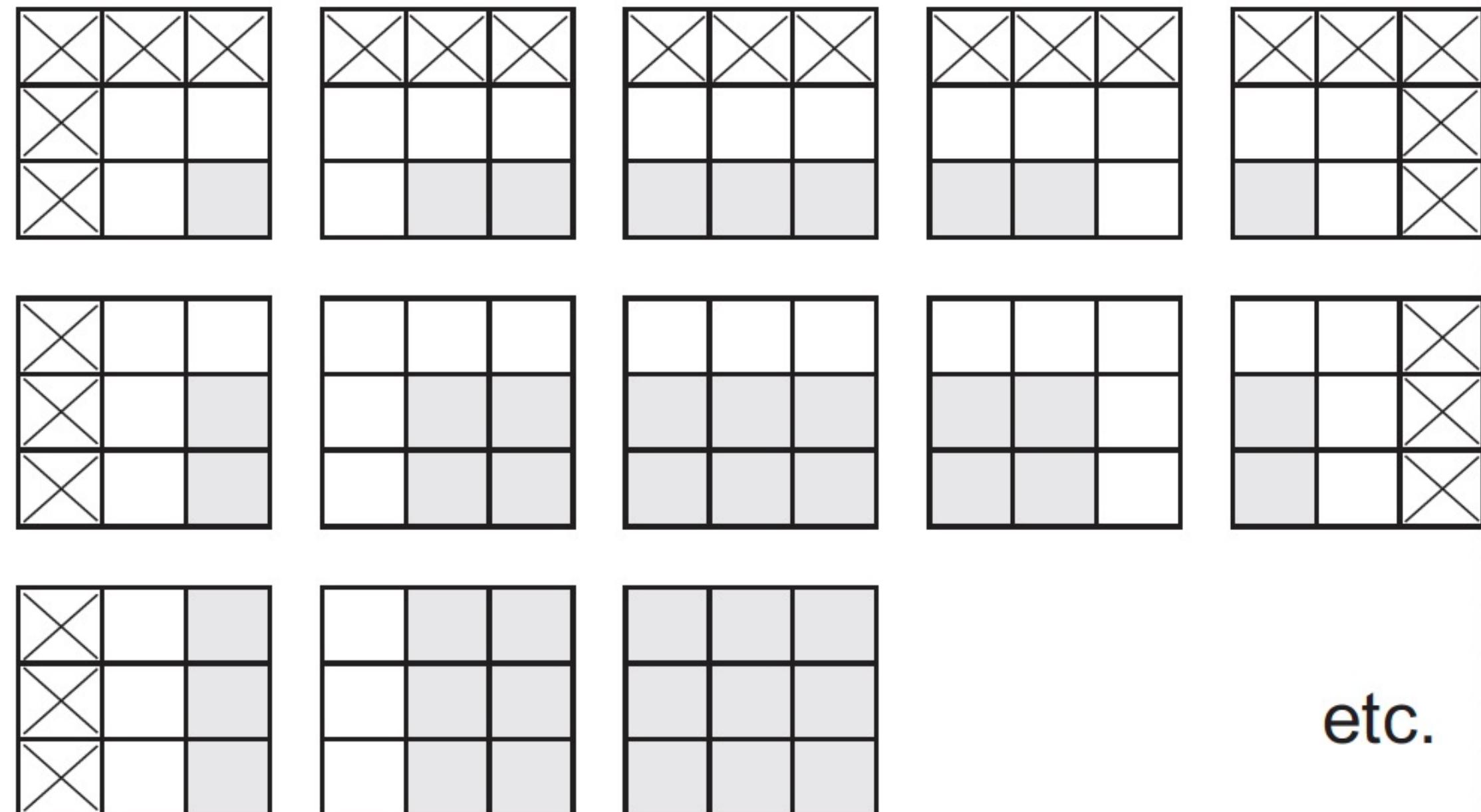
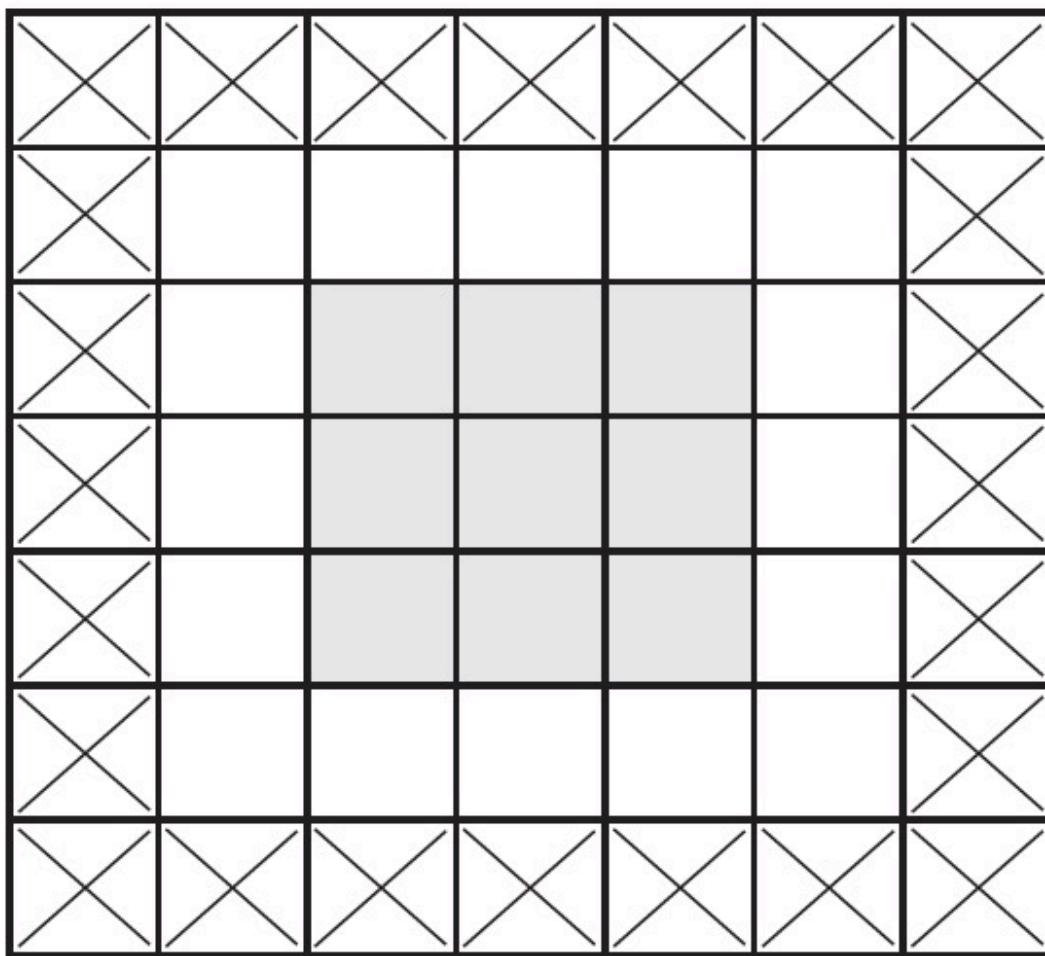


Figure 5.6 in Deep learning with python by Francois Chollet

Stride

the distance between two successive windows

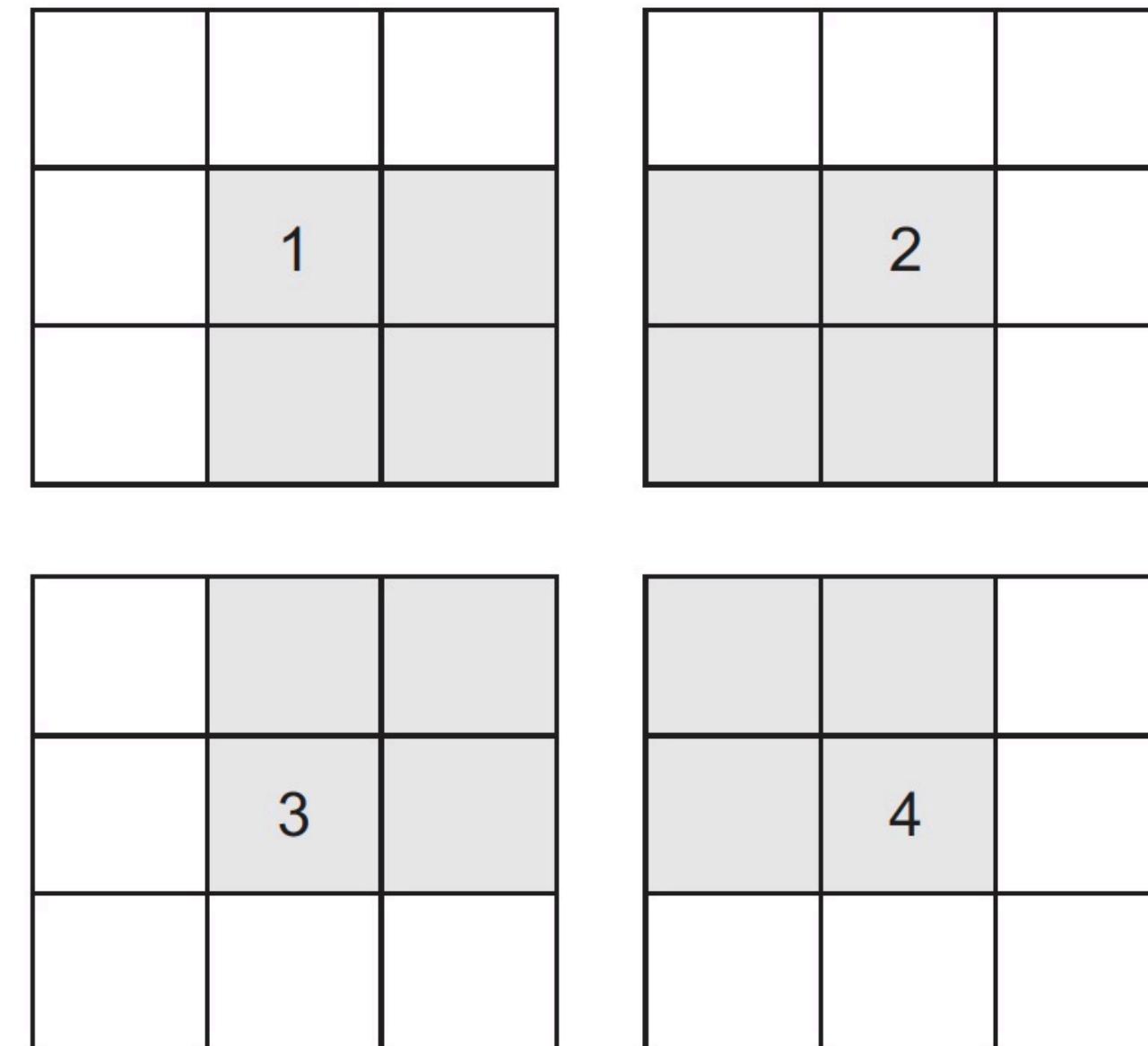
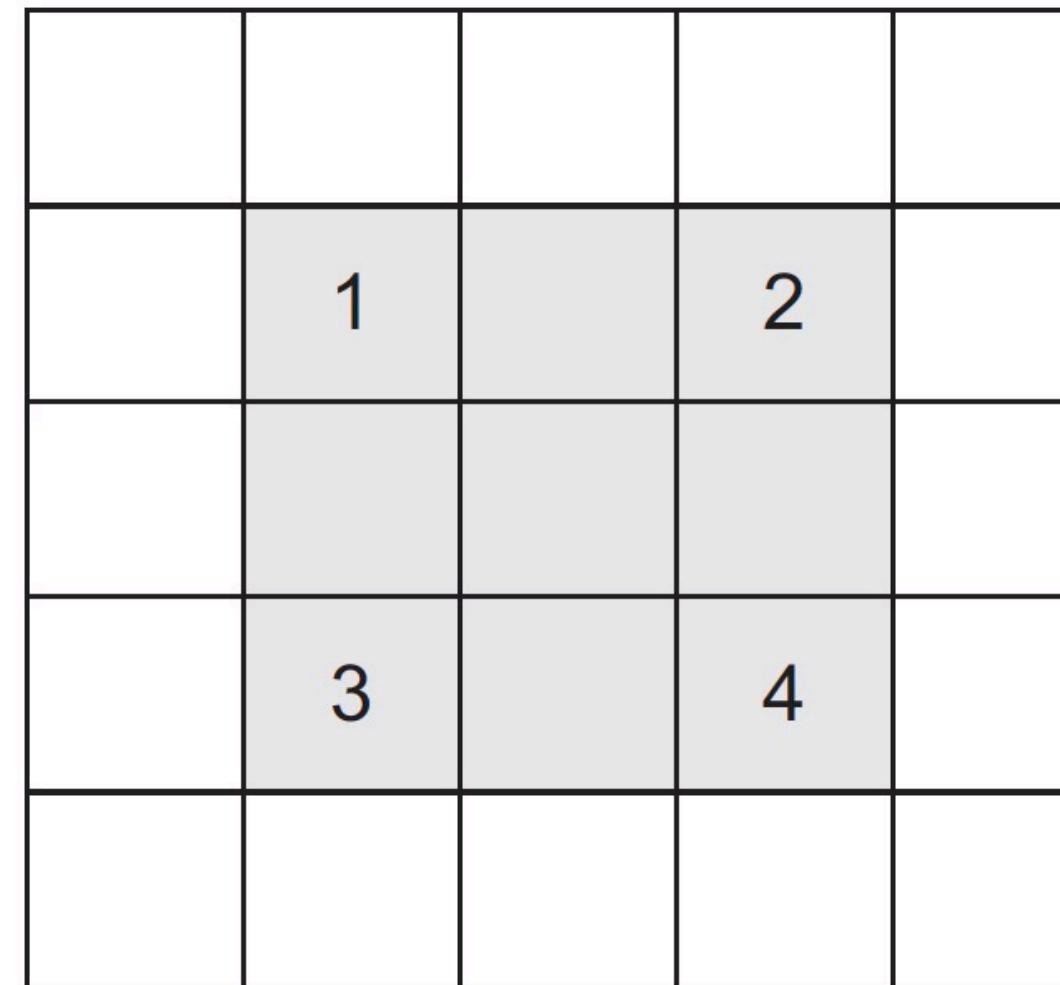


Figure 5.5 in Deep learning with python by Francois Chollet

No padding: `output_size=ceiling((input_size-kernel_size+1)/stride)`

padding: `output_size=ceiling(input_size /stride)`

Convolutional layer

filters: the number of filters in the convolution

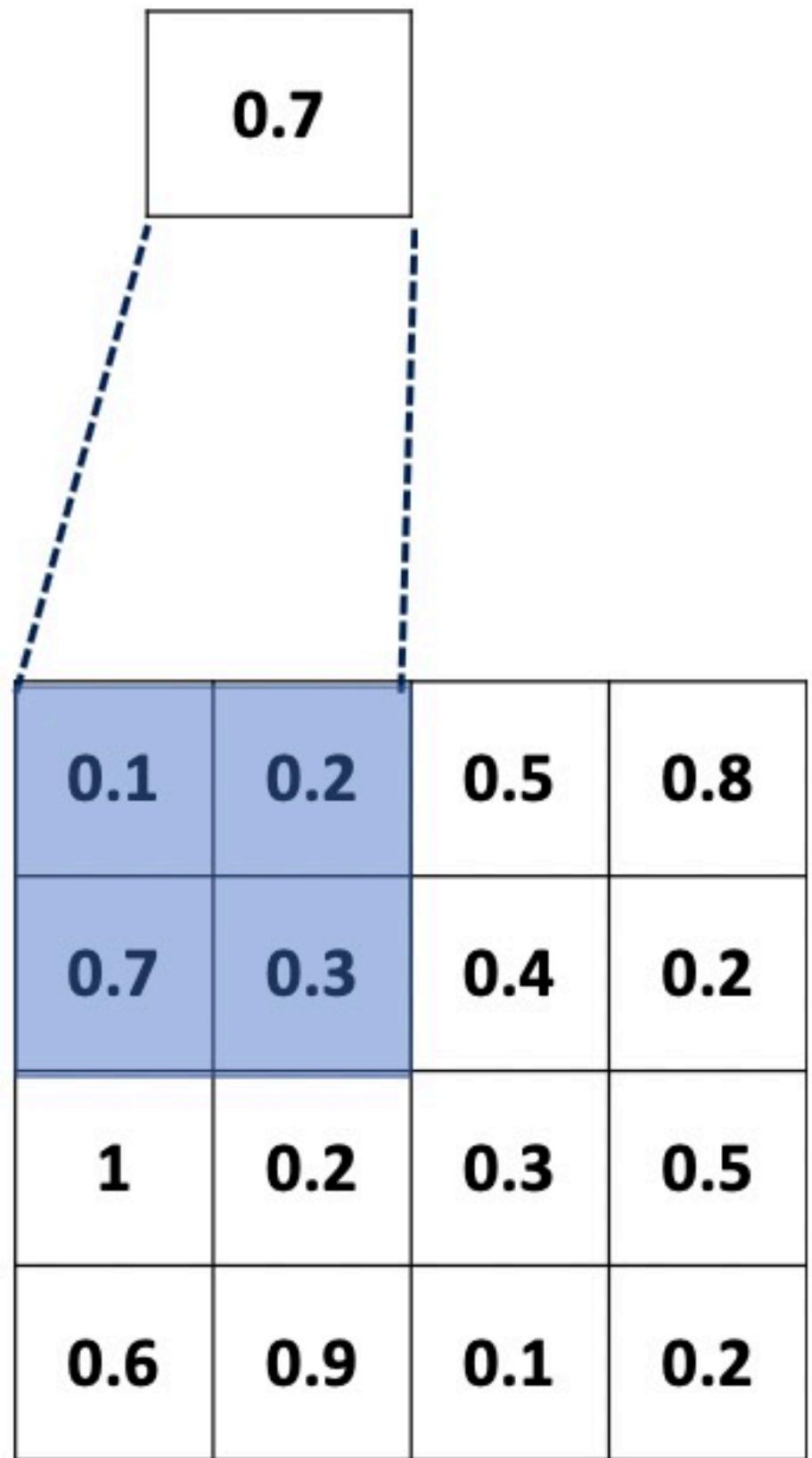
kernel_size: the height and width of the 2D convolution window

padding: one of "valid" or "same"

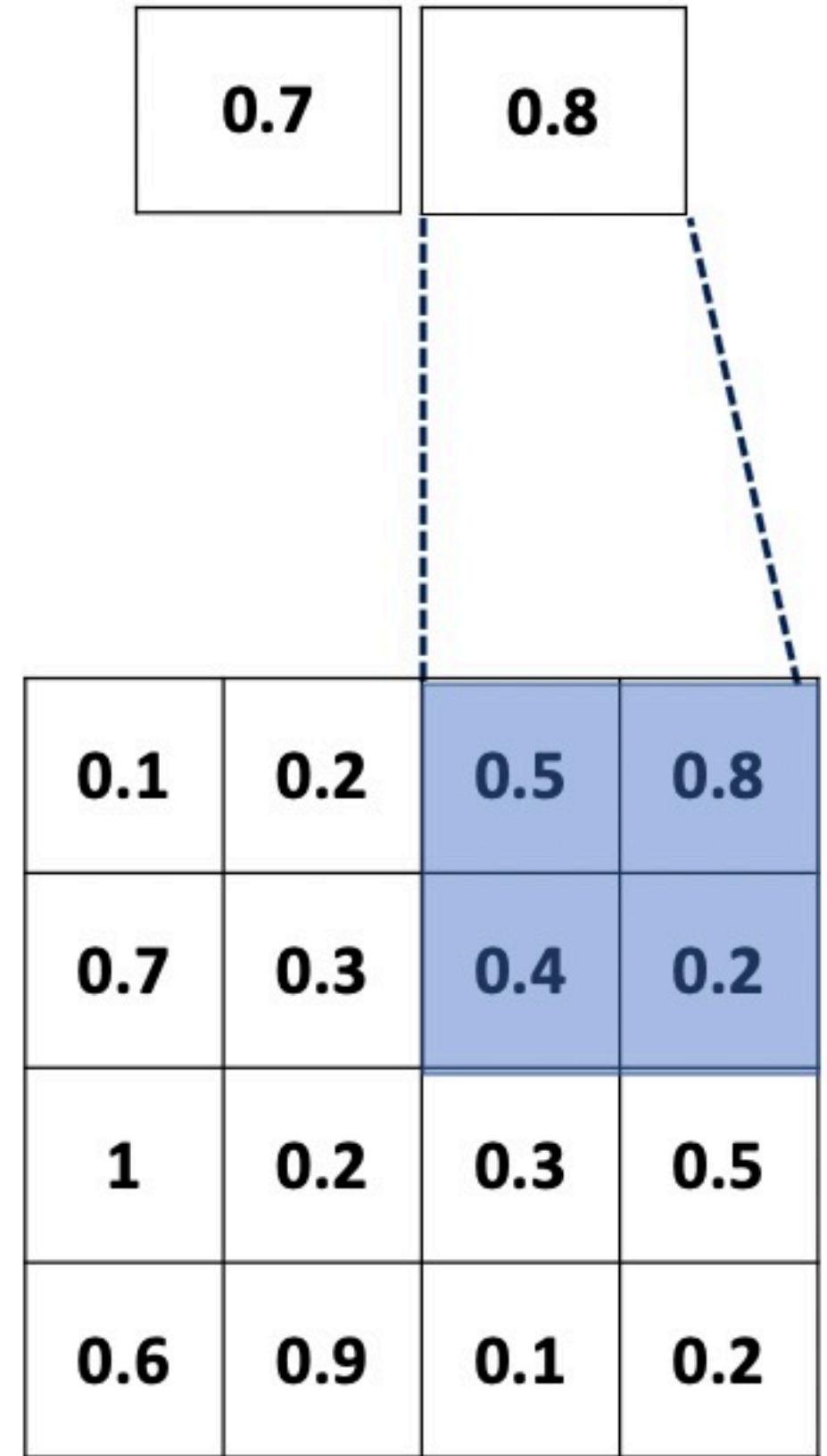
stride: the strides of the convolution along the height and width

```
layer=keras.layers.Conv2D(  
    32,  
    (3, 3),  
    padding='same',  
    strides=1,  
    activation='relu')
```

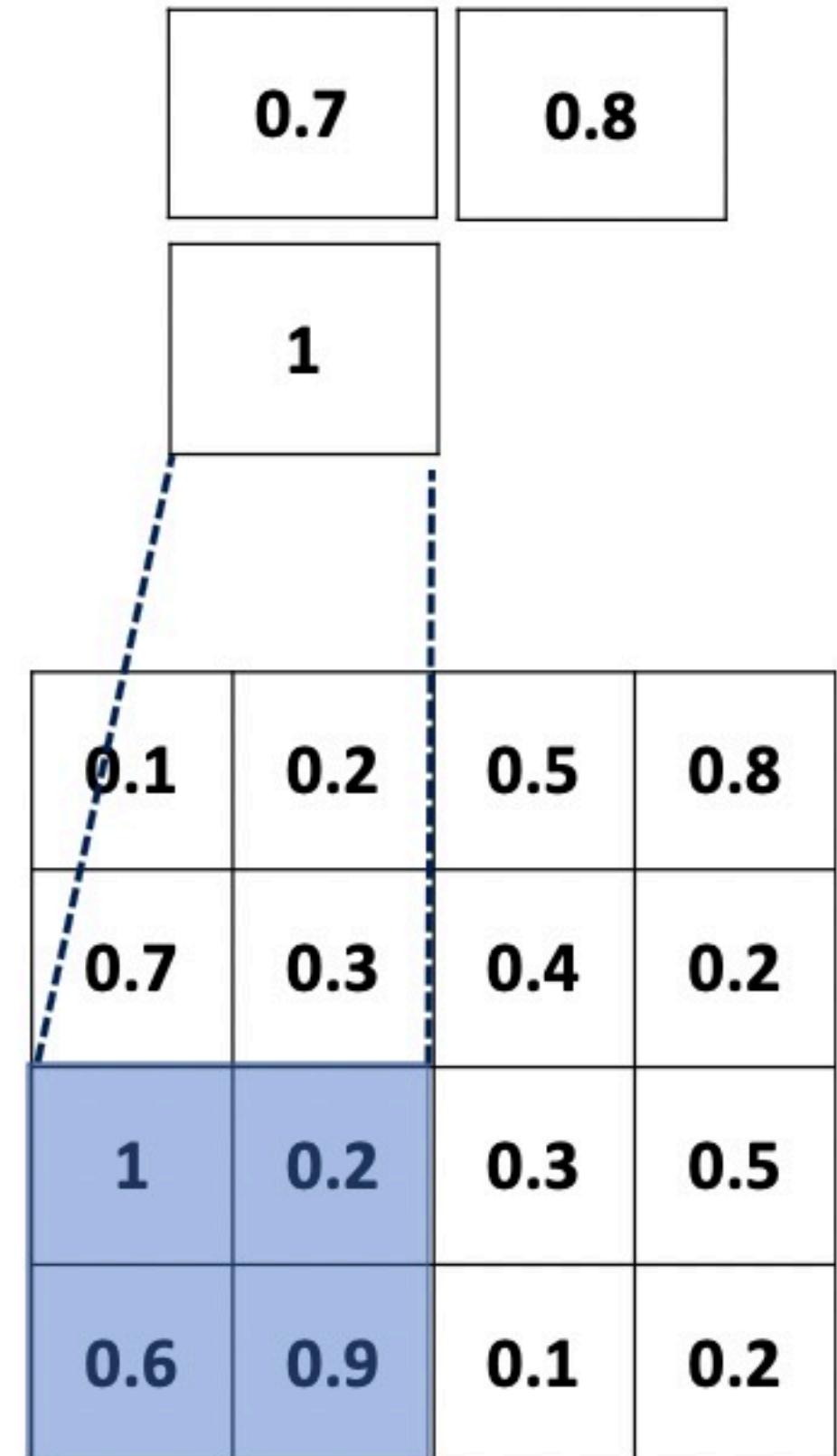
Max-pooling layer



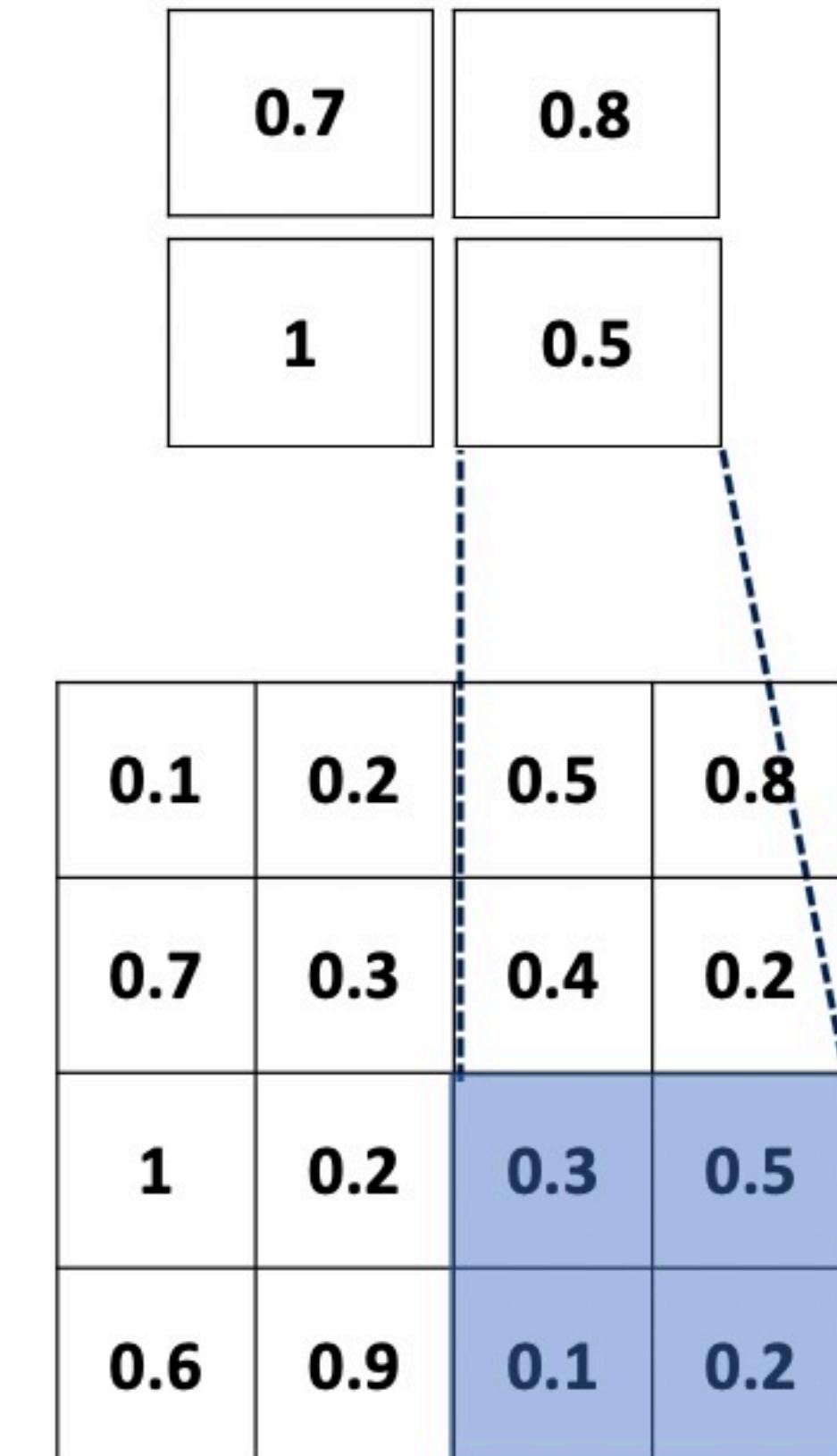
Max-pooling layer



Max-pooling layer

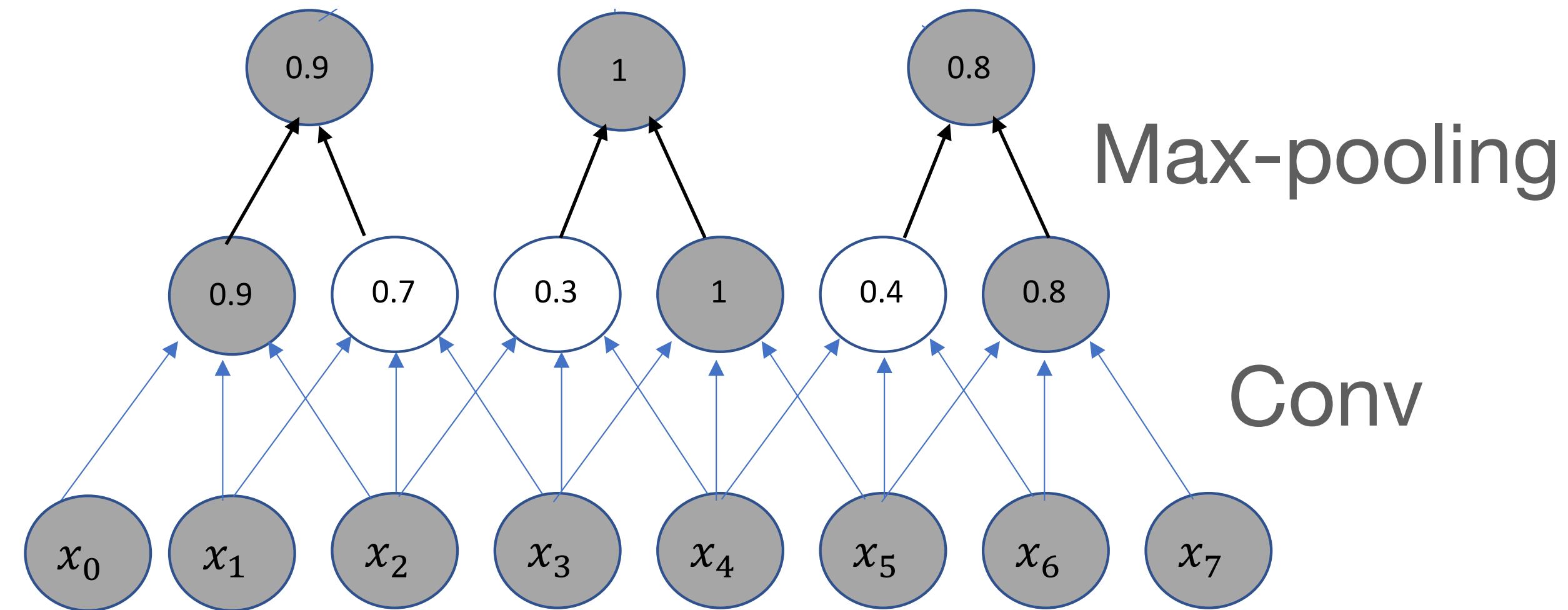


Max-pooling layer

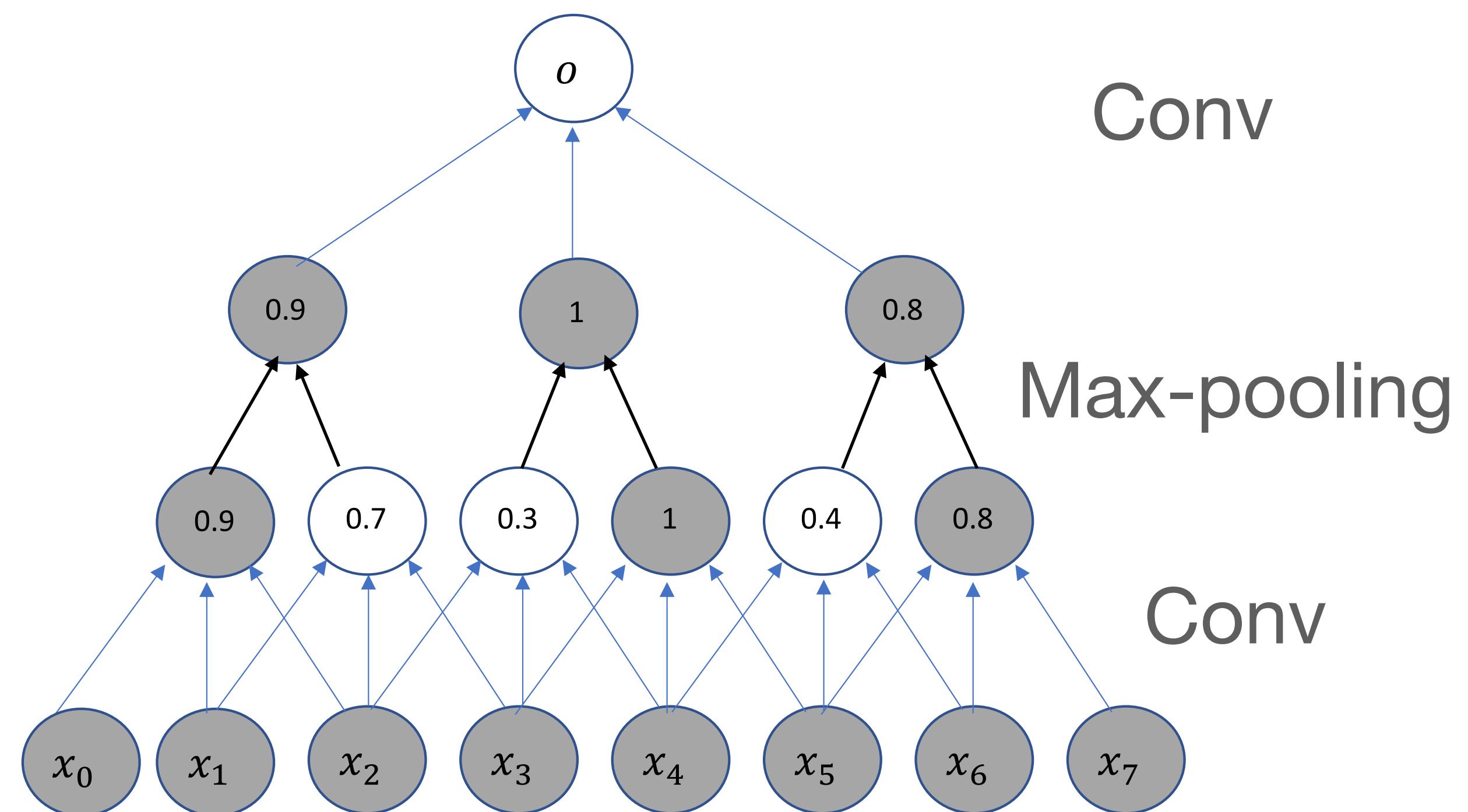
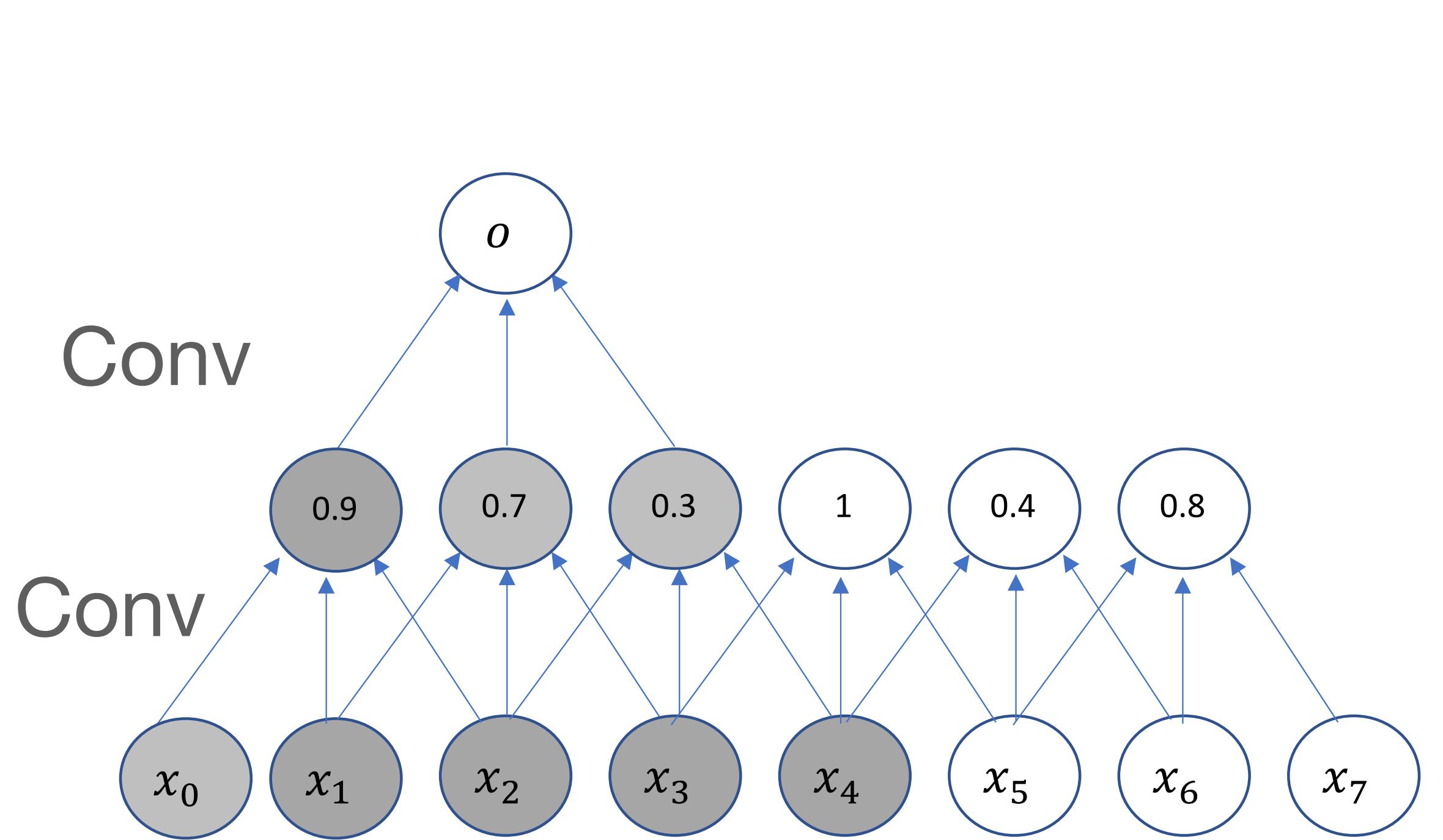


```
tf.keras.layers.MaxPooling2D(  
    pool_size=(2, 2), strides=None, padding="valid")
```

Advantage: downsample feature map, reduce computational burden

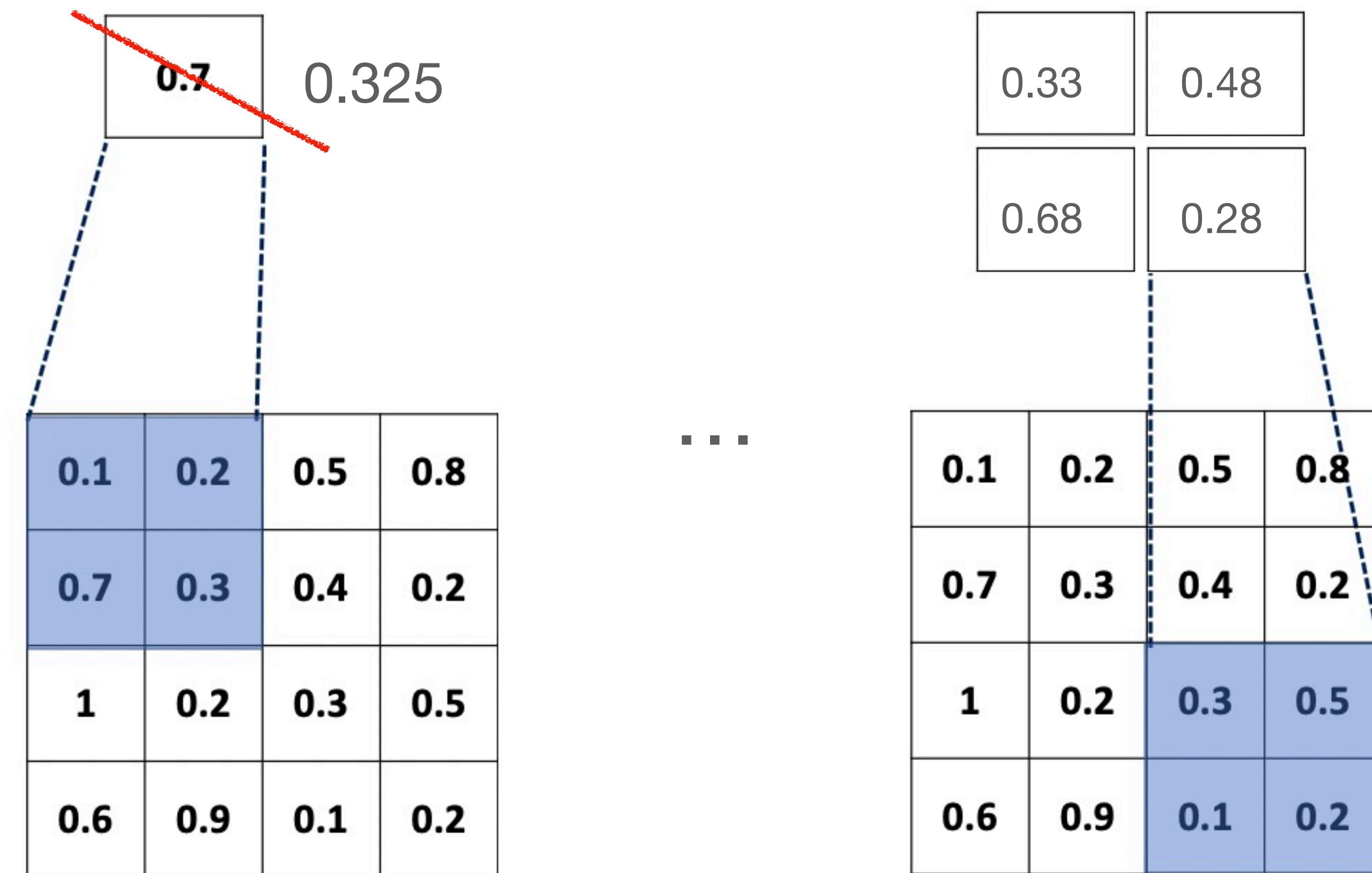


Advantage: increase size of receptive field (window of the input is seen)



Other pooling method: Average pooling

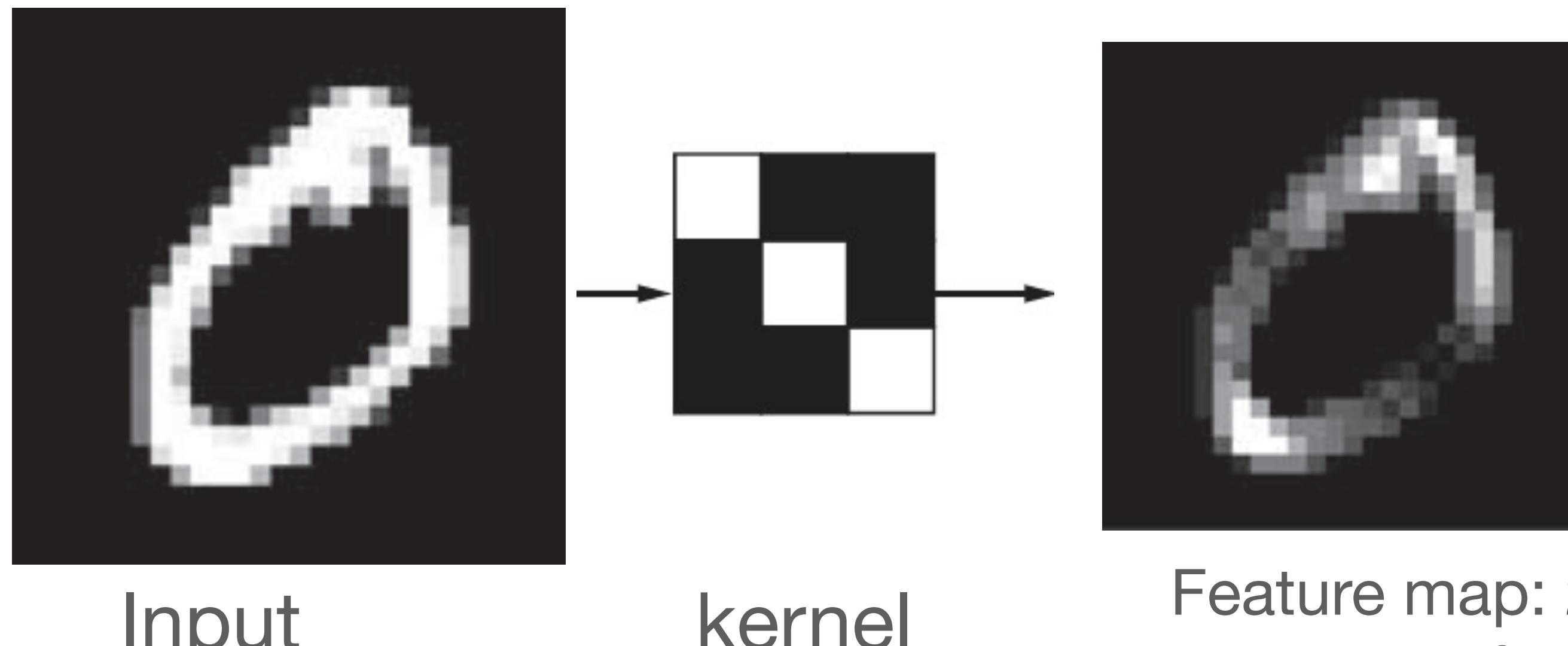
taking the average value over the patch



Why max-pooling to downsample feature map?

Average pooling: dilute feature-presence information

Convolution with stride>1: miss feature-presence information



Feature map: 2D map of the presence of a pattern at different locations in an input

Summary

- Why/what is deep learning?
- Disadvantages of multilayer perceptron?
- How to perform convolution?
- Advantages of convolutional layer
- How to perform max-pooling?
- Advantages of max-pooling?
- What is stride/padding?
- How to compute output size of feature maps?