

[https://www.sli.do/  
#073374](https://www.sli.do/#073374)



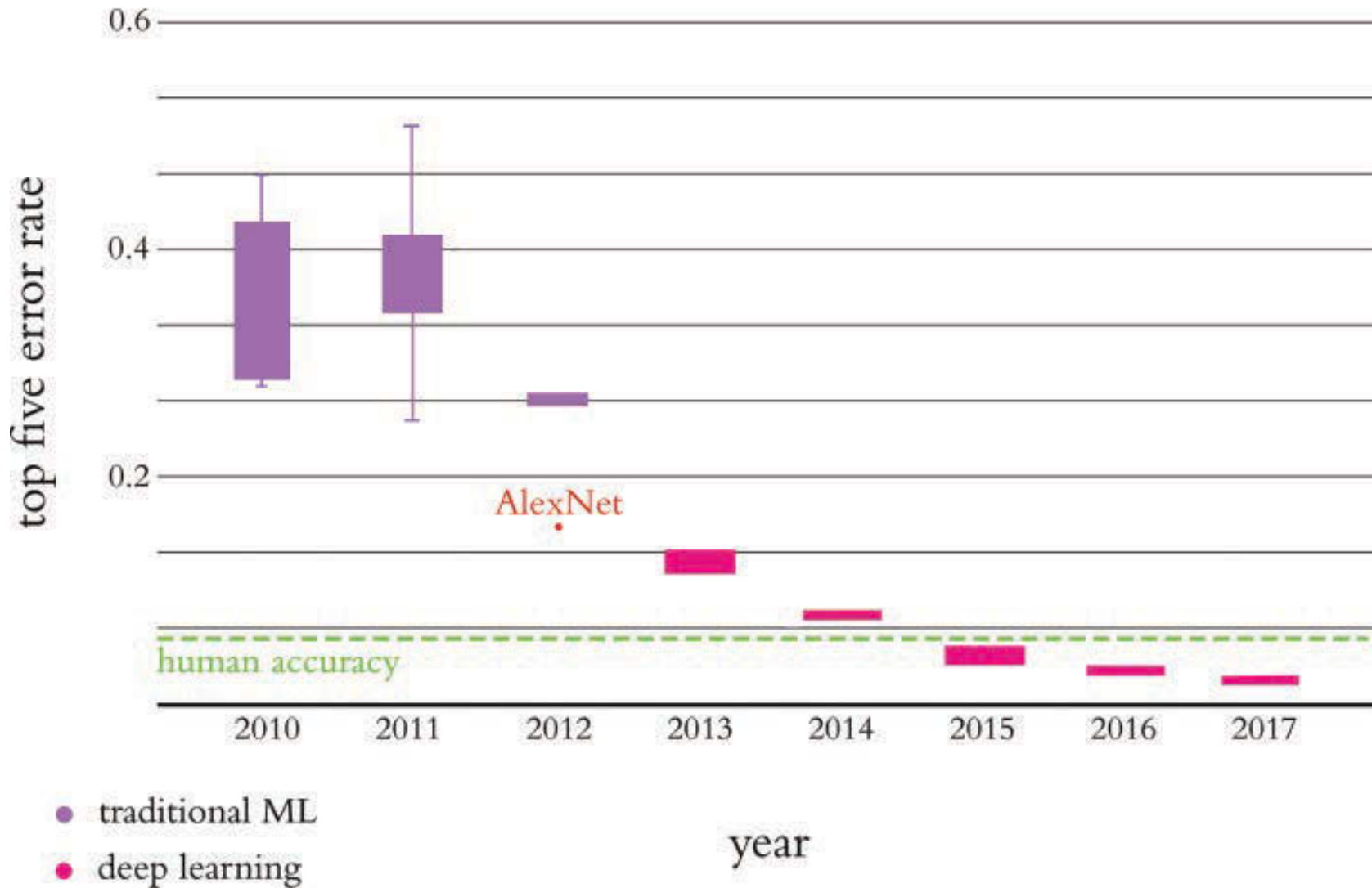
## Deep Learning for Computer Vision (I)

---

### Learning Objectives

- Learn the vast applications of deep learning for computer vision.
- Learn the basics of digital image representation
- Learn the workhorse: convolutional neural networks (CNNs)
- Learn the basic components and theories behind CNNs.

# ILSVRC (the ImageNet Large Scale Visual Recognition Challenge)



# 沈向洋：以 Deep Learning 為核心的 Computer Vision，**十年內**將全面取代人眼 (2019.10.31 【與 AI 大師沈向洋博士對話】)





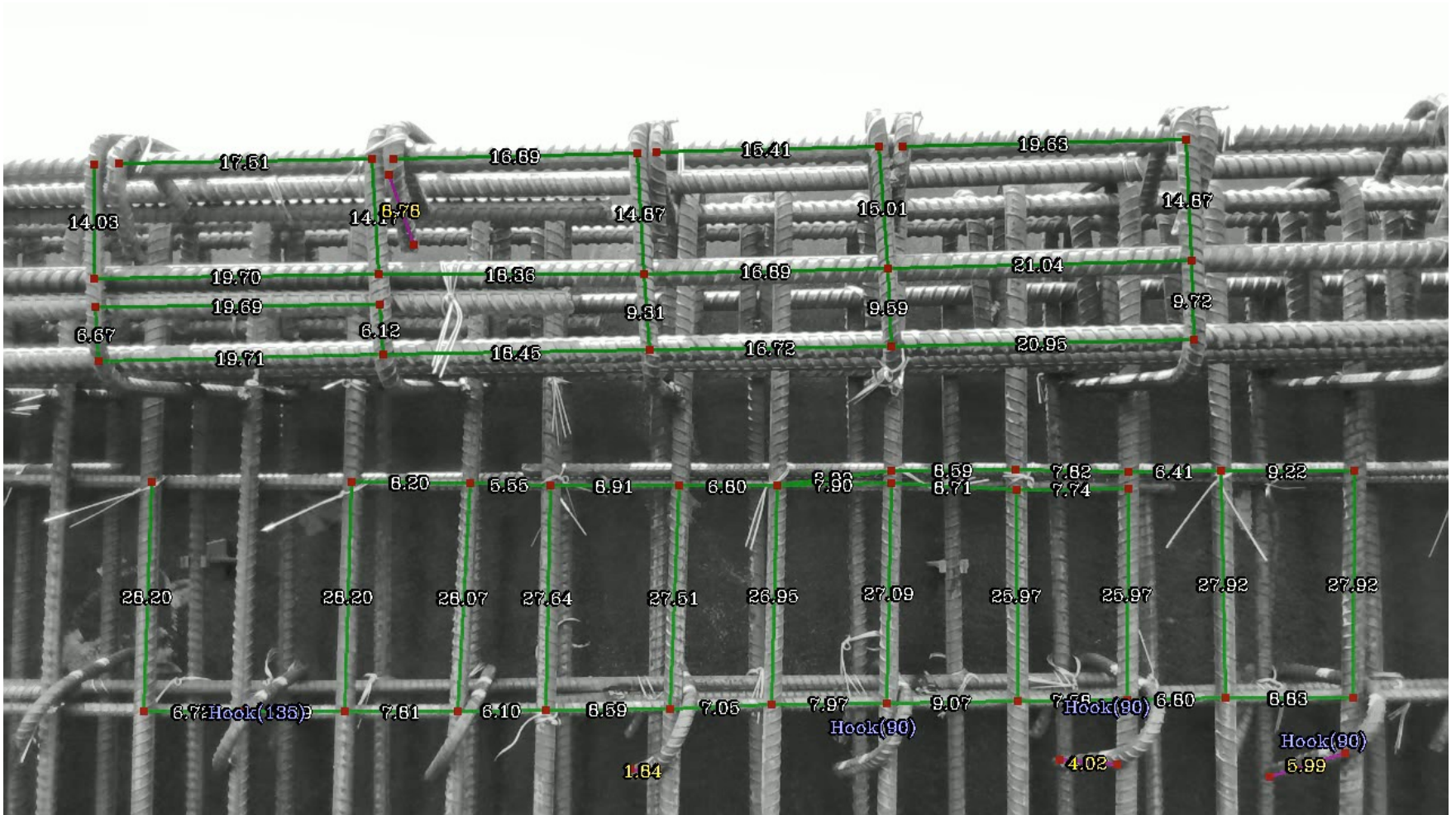
# Tesla released what Autopilot's neural net can **see** (2020.01.31)



Apply cutting-edge research to train **deep neural networks** on problems ranging from perception to control. Our per-camera networks analyze raw images to perform **semantic segmentation**, **object detection** and **monocular depth estimation**. Our birds-eye-view networks take video from all cameras to output the **road layout**, **static infrastructure** and **3D objects** directly in the top-down view.

Our networks learn from the most complicated and diverse scenarios in the world, iteratively sourced from our fleet of nearly 1M vehicles in real time. A full build of Autopilot neural networks involves 48 networks that take 70,000 GPU hours to train 🔥. Together, they **output 1,000 distinct tensors (predictions)** at each timestep.

# 工地鋼筋全檢測





# 監造影像智慧加值：運用AI輔助現場工安管理

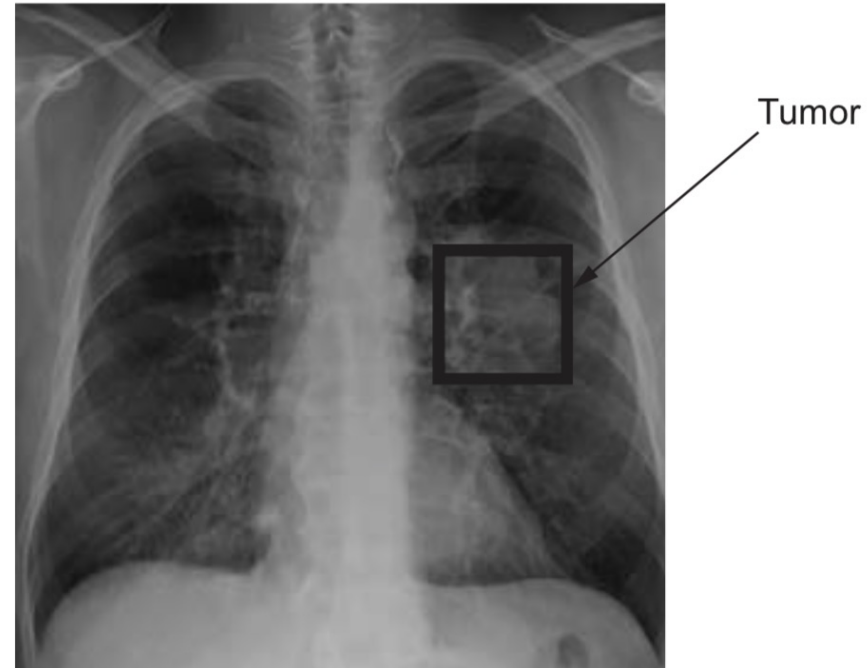


Source：趙志偉 研發工程師 中興工程顧問研發及資訊部

# Applications of Computer Vision



CT scan



X-ray

**Figure 1.5** Vision systems are now able to learn patterns in X-ray images to identify tumors in earlier stages of development.



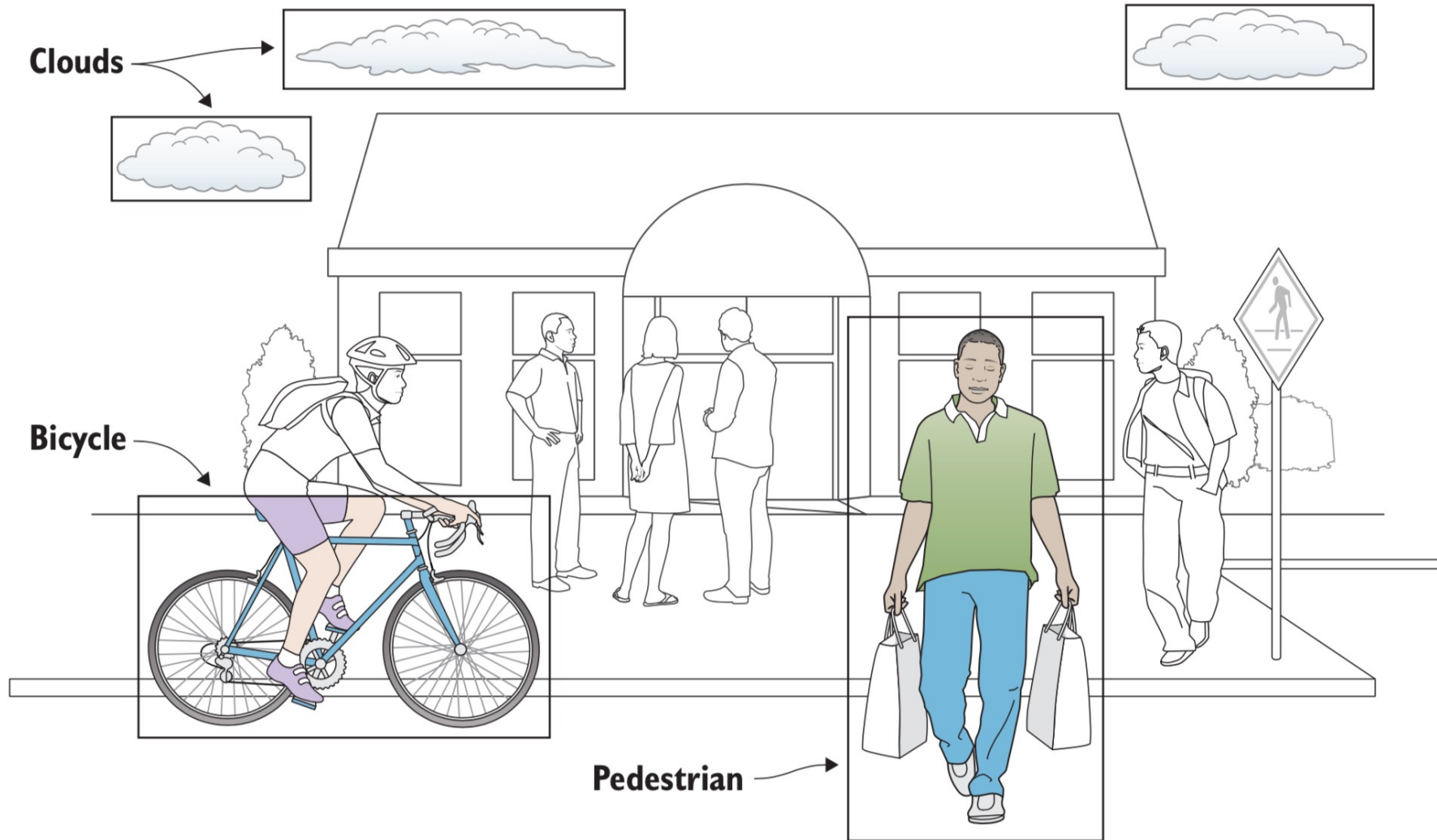
# Applications of Computer Vision



**Figure 1.6** Vision systems can detect traffic signs with very high performance.



# Applications of Computer Vision



**Figure 1.7** Deep learning systems can segment objects in an image.

# Applications of Computer Vision

Original image

Style

Generated art



**Figure 1.8** Style transfer from Van Gogh's *The Starry Night* onto the original image, producing a piece of art that feels as though it was created by the original artist

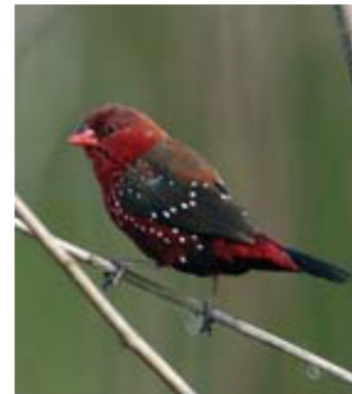


# Applications of Computer Vision

This small blue bird has a short, pointy beak and brown on its wings.



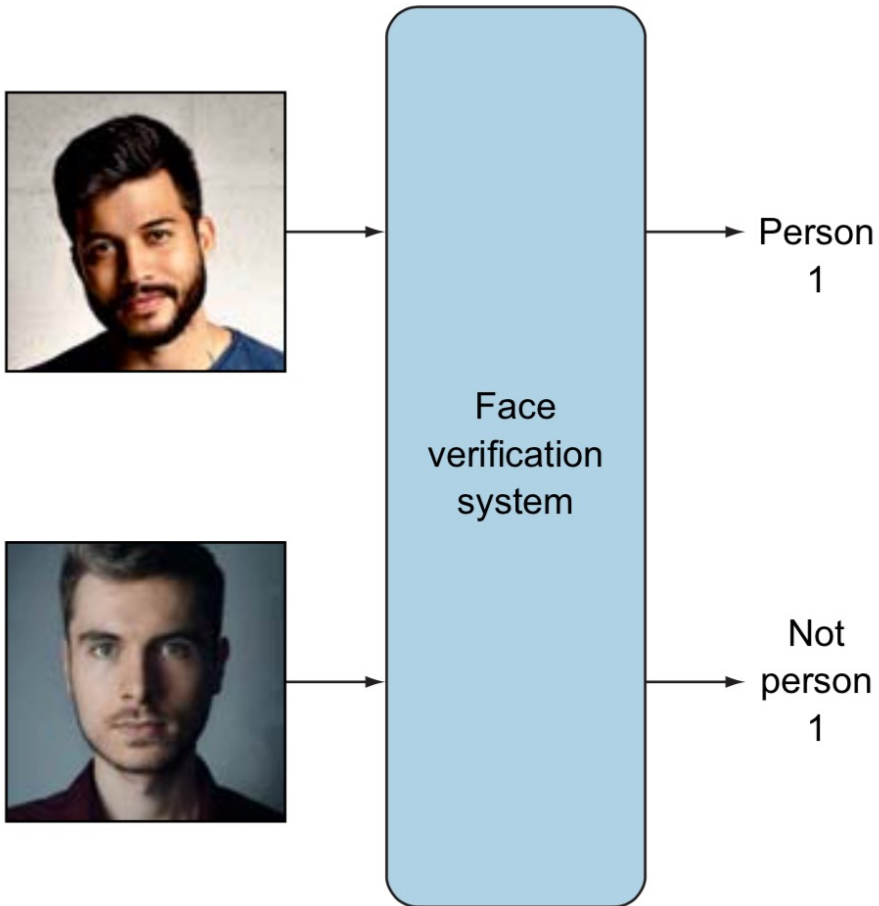
This bird is completely red with black wings and a pointy beak.



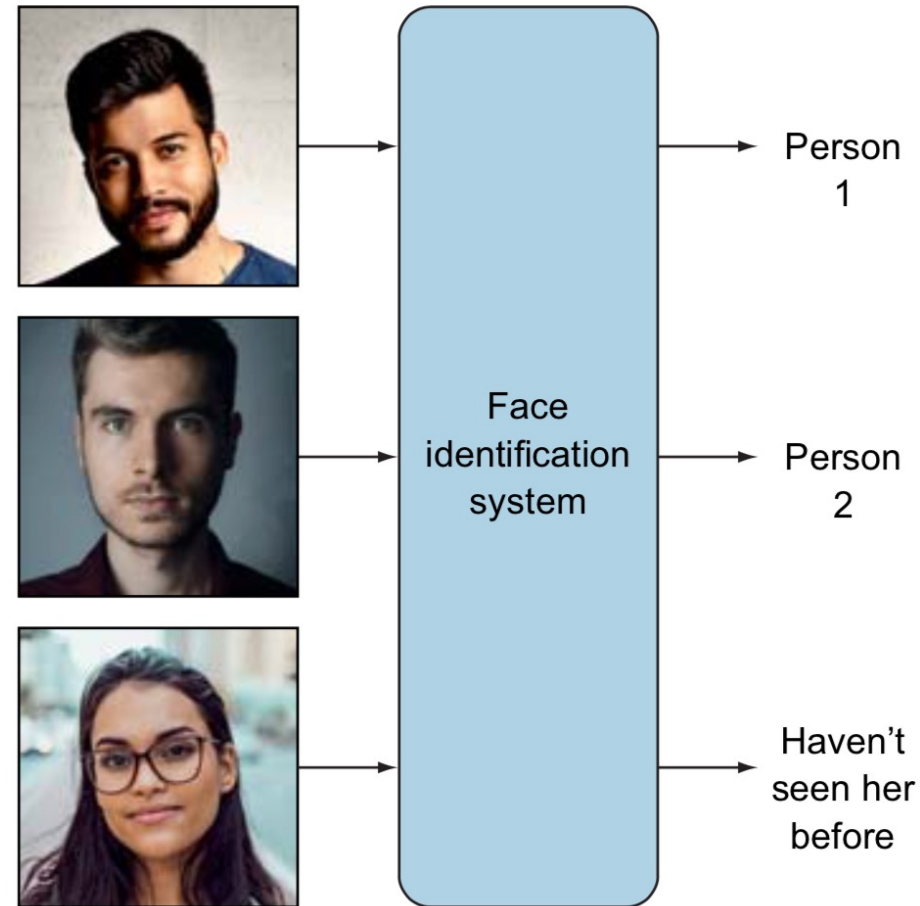
**Figure 1.9** Generative adversarial networks (GANS) can create new, “made-up” images from a set of existing images.

# Applications of Computer Vision

Face verification



Face identification





Query

Retrievals

## Applications of Computer Vision

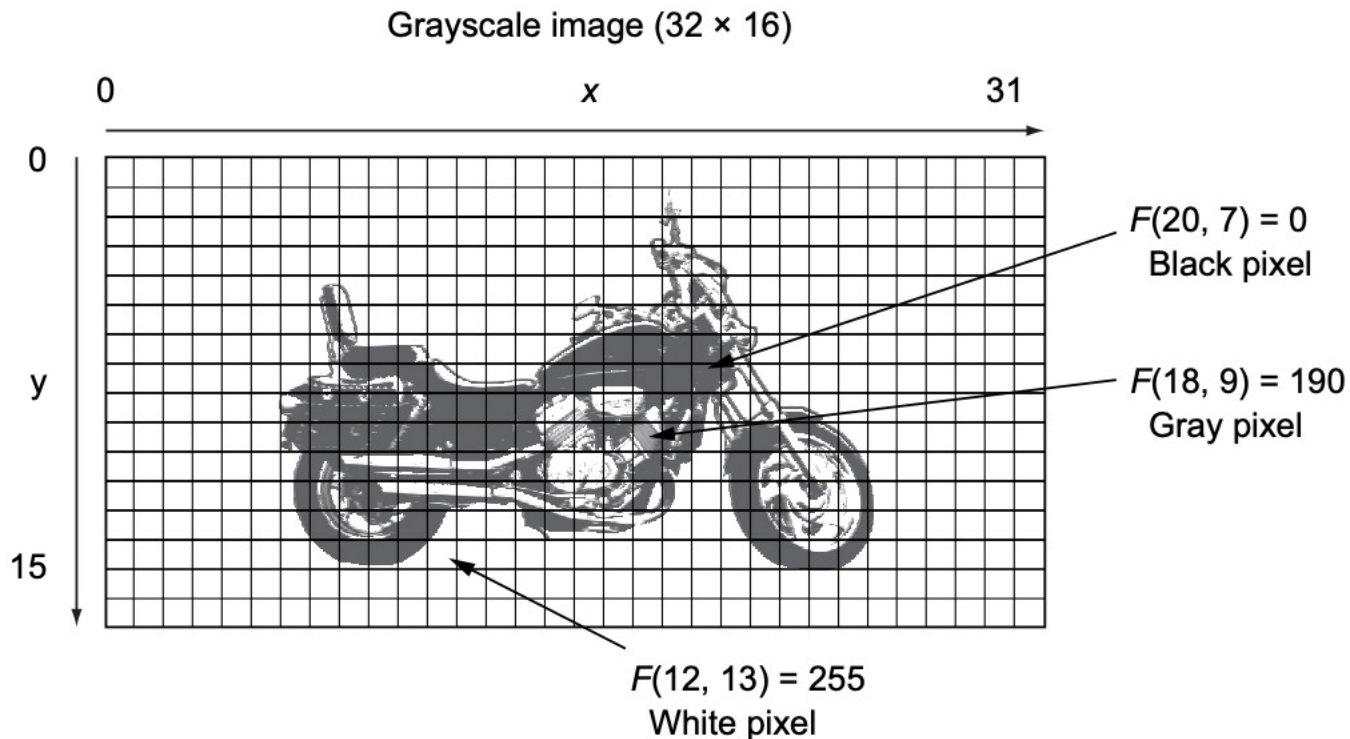


# Image Representation 101



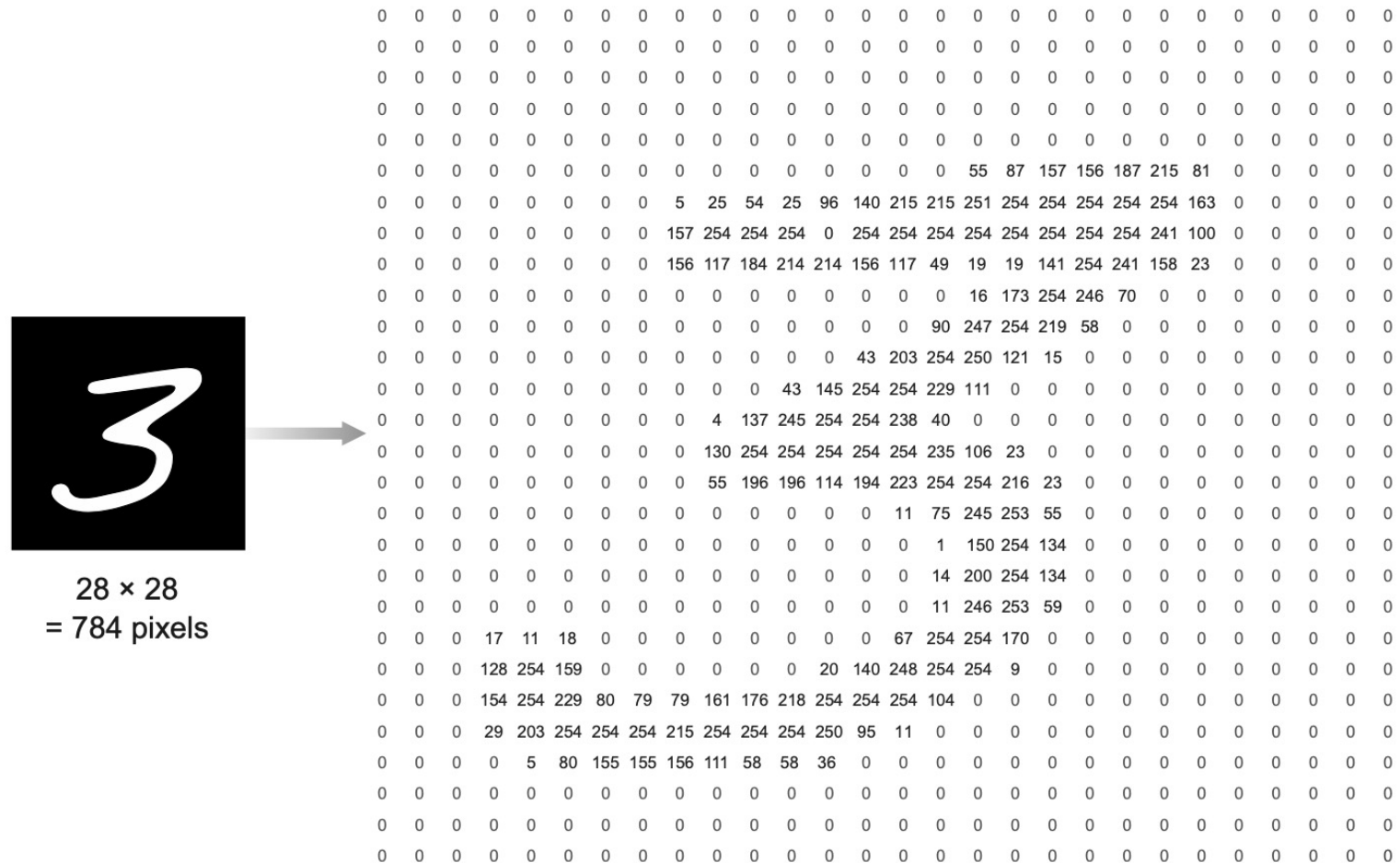
# Image as a Function

- An image can be represented as a function of two variables  $x$  and  $y$ , which define a two- dimensional area. A digital image is made of a grid of pixels.
- The ***pixel*** is the raw building block of an image. Every image consists of a set of pixels in which their values represent the ***intensity*** of light that appears in a given place in the image.



# How Computer See Image

- To a computer, the image looks like a 2D matrix of the pixels' values, which represent intensities. There is no context here, just a massive pile of data.



## How Computer See Color Image

- Color images have three channels (red, green, and blue) and are often represented by three matrices: one represents the intensity of red in the pixel, one represents green, and one represents blue

### Color image

## RGB channels

$$F(0, 0) = [11, 102, 35]$$


Channel 3  
Blue intensity  
values

Channel 2  
Green intensity  
values

Channel 1  
Red intensity  
values

Channel 3

Blue intensity values

2

Intensity

1

Intensity

35

102

11



# Image Flattening

- We often need to do image flattening when we use MLP for classification.

To help visualize the flattened input vector, let's look at a much smaller matrix (4, 4):

Blue	White	White	White
Blue	Blue	Blue	White
Blue	Blue	White	Blue
Blue	White	Blue	Blue

The input ( $x$ ) is a flattened vector with the dimensions (1, 16):

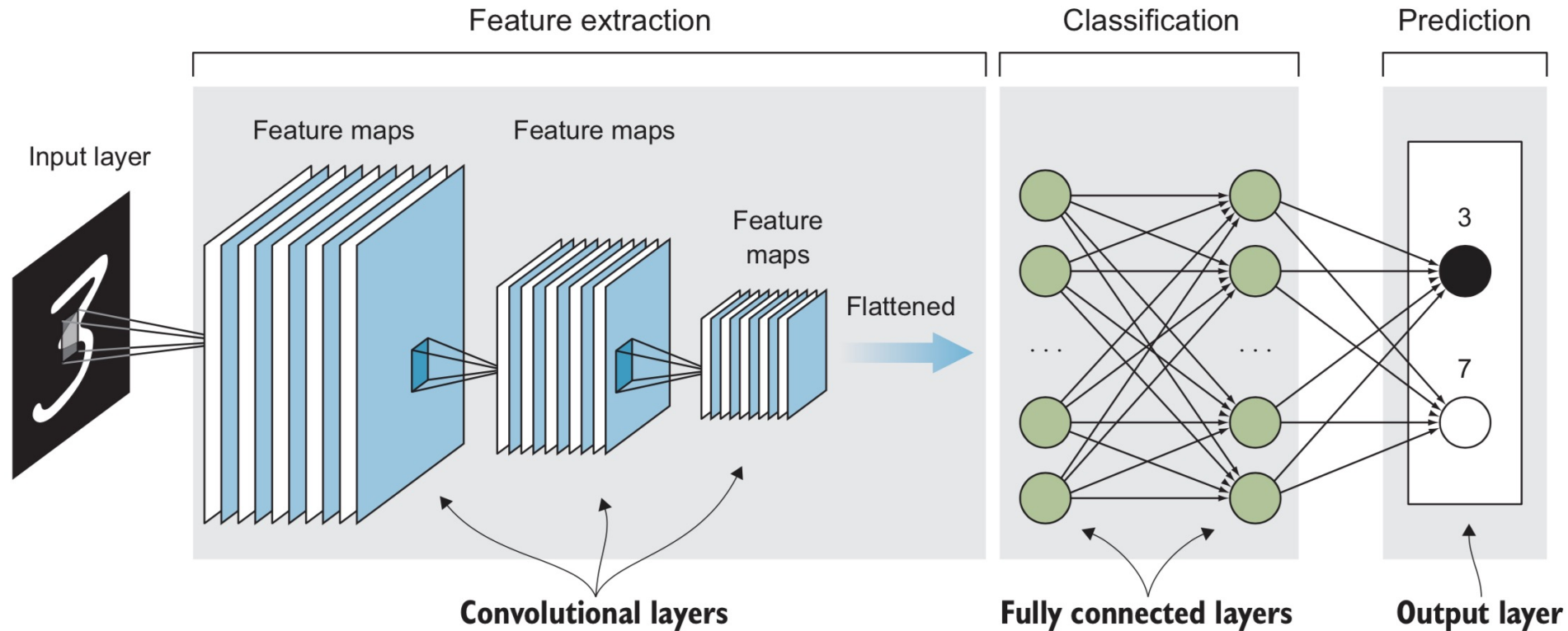
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$
Row 1				Row 2				Row 3				Row 4			

So, if we have pixel values of 0 for black and 255 for white, the input vector will be as follows:

Input = [0, 255, 255, 255, 0, 0, 0, 255, 0, 0, 255, 0, 0, 255, 0, 0]

# **Convolutional Neural Networks (CNNs) : An Overview**

# High-level architecture of CNNs: **input layer**, **convolutional layers**, **fully connected layers**, and **output prediction**





## DNNs vs. CNNs: See CNN\_MNIST.ipynb

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 784)	615440
dense_1 (Dense)	(None, 10)	7850
Total params: 623,290		
Trainable params: 623,290		
Non-trainable params: 0		

Layer (type)	Output Shape	Param #
conv2d_3 (Conv2D)	(None, 26, 26, 32)	320
max_pooling2d_2 (MaxPooling2D)	(None, 13, 13, 32)	0
conv2d_4 (Conv2D)	(None, 11, 11, 64)	18496
max_pooling2d_3 (MaxPooling2D)	(None, 5, 5, 64)	0
conv2d_5 (Conv2D)	(None, 3, 3, 64)	36928
flatten_1 (Flatten)	(None, 576)	0
dense_4 (Dense)	(None, 64)	36928
dense_5 (Dense)	(None, 10)	650
Total params: 93,322		
Trainable params: 93,322		
Non-trainable params: 0		

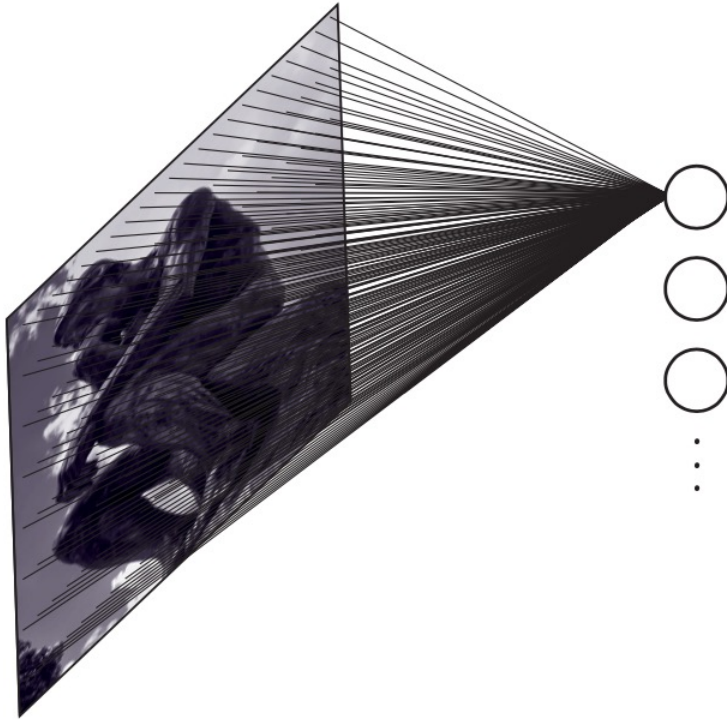
CNNs yield better results with less training parameters

DNNs and CNNs do not usually yield comparable results; MNIST dataset is an exception. In messy real-world image data, CNNs truly outshine DNNs.

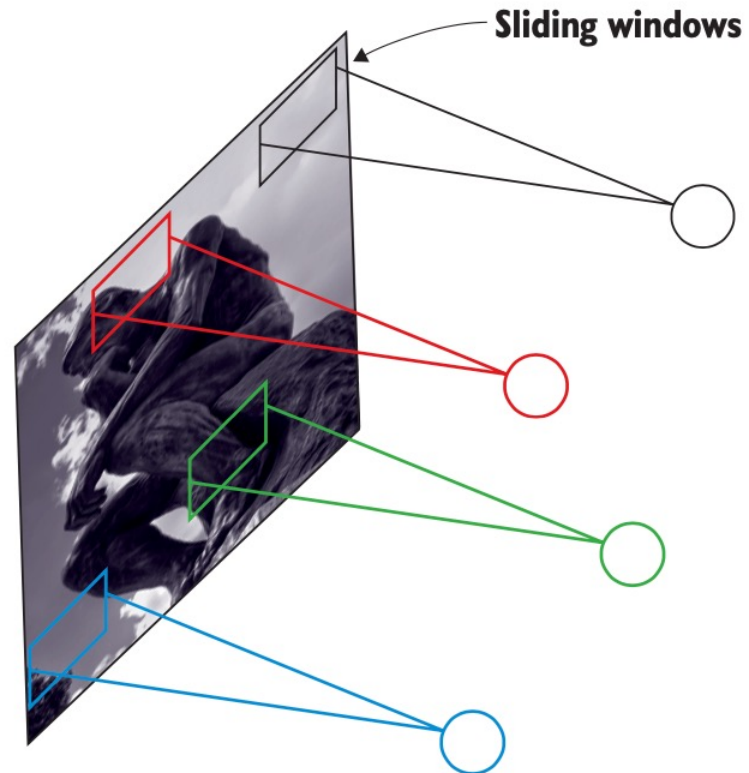
# Fundamental Difference between DNNs and CNNs

## Global vs. Local

Fully connected neural net



Locally connected neural net

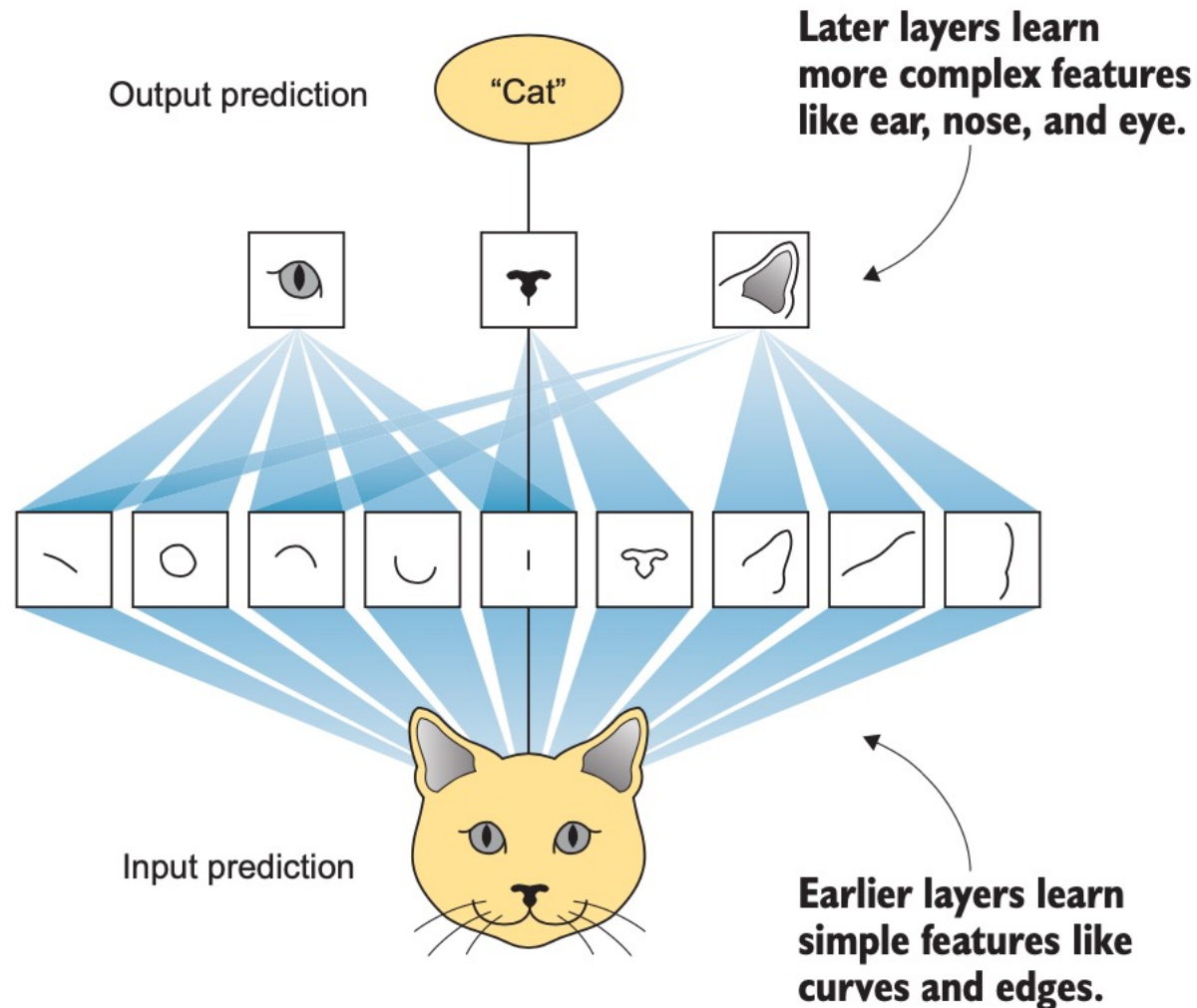


**The patterns CNNs learn are translation invariant.  
CNNs can learn spatial hierarchies of patterns.**

# CNNs learn the image features through its layers.

The patterns CNNs learn are translation invariant.

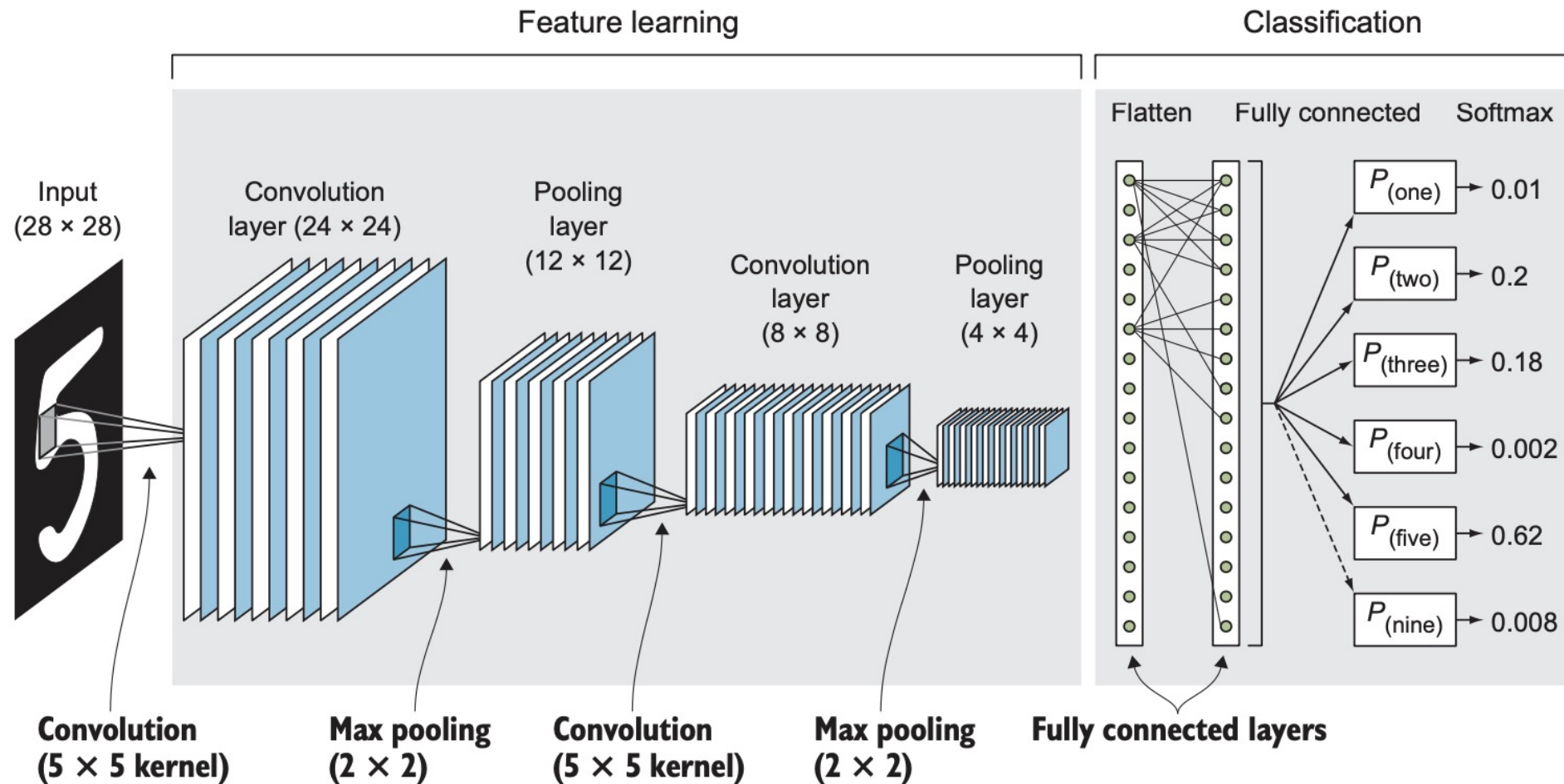
CNNs can learn spatial hierarchies of patterns.





# Basic Components of Convolutional Neural Networks (CNNs)

The basic components of convolutional networks are **convolutional layers** and **pooling layers** to perform **feature extraction**, and **fully connected layers** for classification



# Basic components of CNNs: theoretical minimum and example

- The phrase “theoretical minimum” is taken from a very successful book series written by Leonard Susskind, a great physicist at Stanford University.
- “Theoretical minimum” means just the minimum theories and equations you need to know in order to proceed to the next level.
- See CNN\_Basics.pdf