سركيب

**THEASSEMBLY**

DATA SCIENCE

# About The Assembly

- A smart lab based out of **in5 Tech** since Dec 2014

- Almost 300 free workshops done

- ASSEMBLY: HACK - Embedded systems, IoT and hardware

- ASSEMBLY: CODE - Software projects - APIs, frameworks, apps

- Age range: 16-60 - students, professionals, entrepreneurs

- Focus on smart technology and practical applications

- Forum: **members.theassembly.ae**

THEASSEMBLY
MAKE | SMART | THINGS

# Tag Us On Social Media

FACEBOOK    The Assembly (@MakeSmartThings)

TWITTER    @MakeSmartThings

INSTAGRAM    @MakeSmartThings

YOUTUBE    The Assembly

# Introducing The Assembly: Data Science

- Covers
  - **Data Science Platforms & Techniques**
  - **Machine Learning & Artificial Intelligence Paradigms**
- More sophisticated content
- Deep dive into the topics – lots of technical knowledge
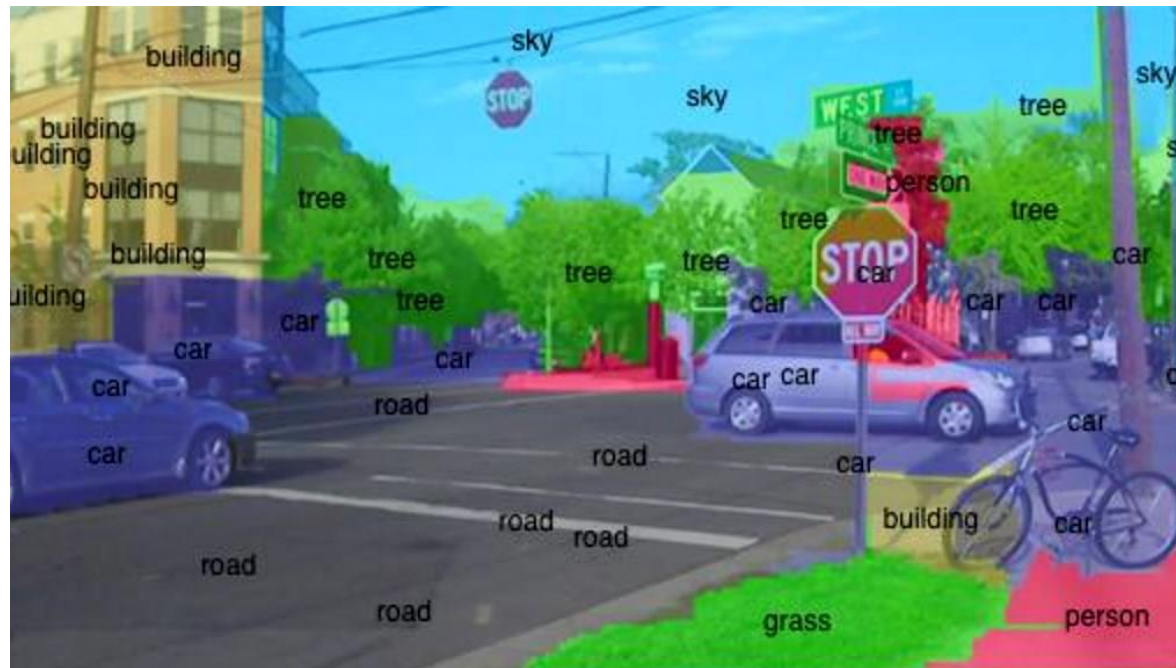- Hands on guidance from experts

# GitHub Link

https://github.com/anudeepsekhar/The-Assembly-Computer-Vision-Worshop

# If We Want Machines to Think, We Need to Teach Them to See

By Anudeepsekhar Bolimera

# What is Computer Vision?

An interdisciplinary scientific field that deals with how computers can be made to gain high-level understanding from digital images or videos. From the perspective of engineering, it seeks to automate tasks that the human visual system can do.

# How does it work?

The process involves

- Acquiring digital images

- Processing images to enhance features

- Analyzing & understanding the images

- Extracting high-dimensional data from the real world in order to produce numerical or symbolic information
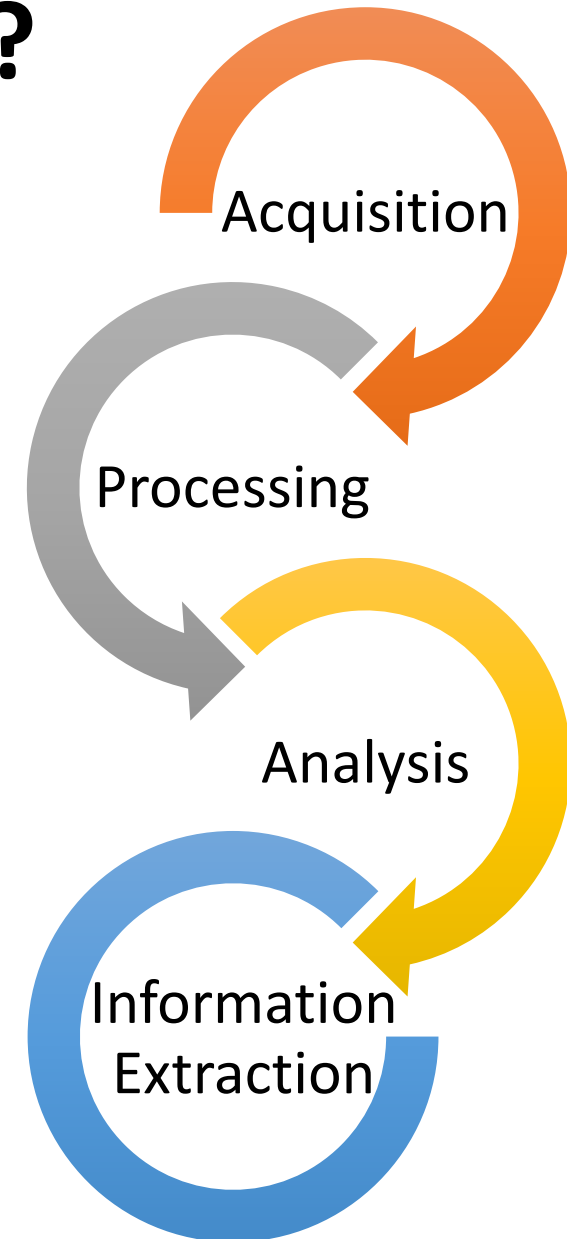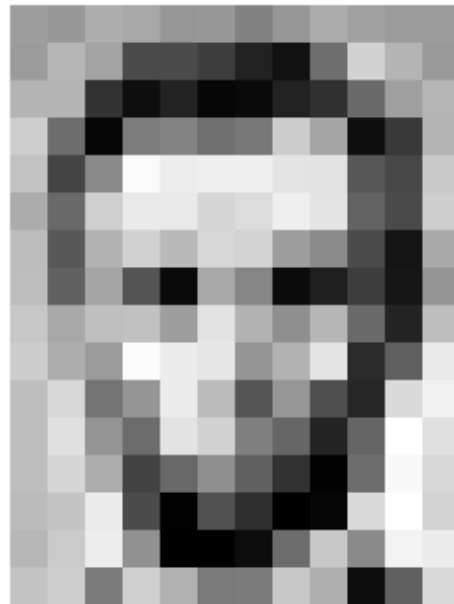
Acquisition

Processing

Analysis

Information Extraction

تركيب
THEASSEMBLY
MAKE | SMART | THINGS

# Image Acquisition

*The primary sensor machines use to see the world around them as images is **the camera**.* Machines interpret images (acquired from camera) very simply: **as a series of pixels**, each with their own set of color values.

# Image Processing

Operations are performed on images to get intended manipulations – which will enhance useful information available in the images.

# Image Analysis & Data Extraction

Once the image has been primed, we can analyze the contents and extract meaningful data from it

THEASSEMBLY
MAKE | SMART | THINGS

# Core Tasks In Computer Vision

- **Object classification**
- **Object detection**
- **Image segmentation**
- **Scene reconstruction**
- **Image restoration**
- **Motion analysis**
- **Visual QnA**
- **Image generation**

Any other application that involves understanding pixels through software can safely be labeled as computer vision.



CAT

# Core Tasks In Computer Vision

- **Object classification**
- **Object detection**
- **Image segmentation**
- **Scene reconstruction**
- **Image restoration**
- **Motion analysis**
- **Visual QnA**
- **Image generation**

Any other application that involves understanding pixels through software can safely be labeled as computer vision.



CAT, DOG, DUCK

# Core Tasks In Computer Vision

- **Object classification**
- **Object detection**
- **Image segmentation**
- **Scene reconstruction**
- **Image restoration**
- **Motion analysis**
- **Visual QnA**
- **Image generation**

Any other application that involves can safely be labeled as computer vision.

# Core Tasks In Computer Vision

- **Object classification**

- **Object detection**

- **Image segmentation**

- **Scene reconstruction**

- **Image restoration**

- **Motion analysis**

- **Visual QnA**

- **Image generation**

Any other application that involve
can safely be labeled as computer vision.

# Core Tasks In Computer Vision

- **Object classification**

- **Object detection**

- **Image segmentation**

- **Scene reconstruction**

- **Image restoration**

- **Motion analysis**

- **Visual QnA**

- **Image generation**

Any other application that involv
can safely be labeled as computer

# Core Tasks In Computer Vision

- **Object classification**
- **Object detection**
- **Image segmentation**
- **Scene reconstruction**
- **Image restoration**
- **Motion analysis**
- **Visual QnA**
- **Image generation**

Any other application that involv
can safely be labeled as computel



By Anudeepsekhar Bolimera

# Core Tasks In Computer Vision

- **Object classification**

- **Object detection**

- **Image segmentation**

- **Scene reconstruction**

- **Image restoration**

- **Motion analysis**

- **Visual QnA**

- **Image generation**

Any other application that invo
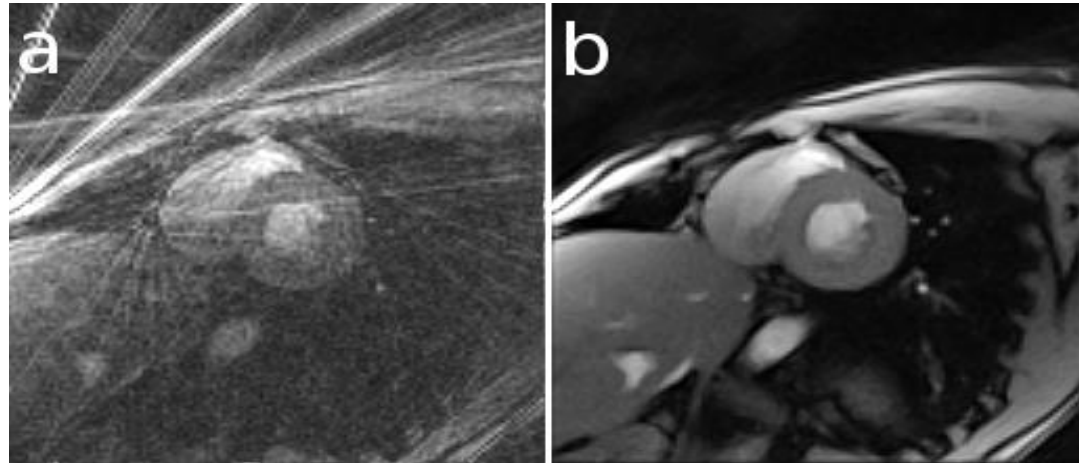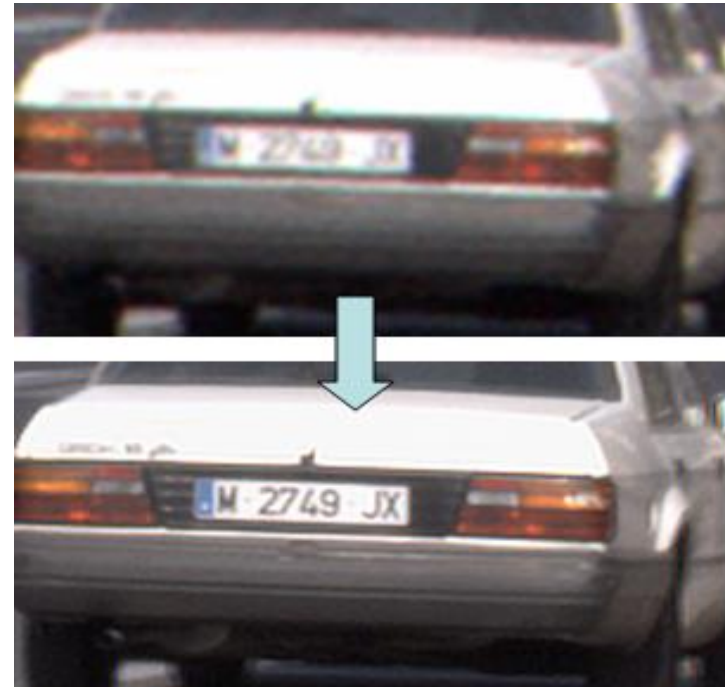can safely be labeled as compute
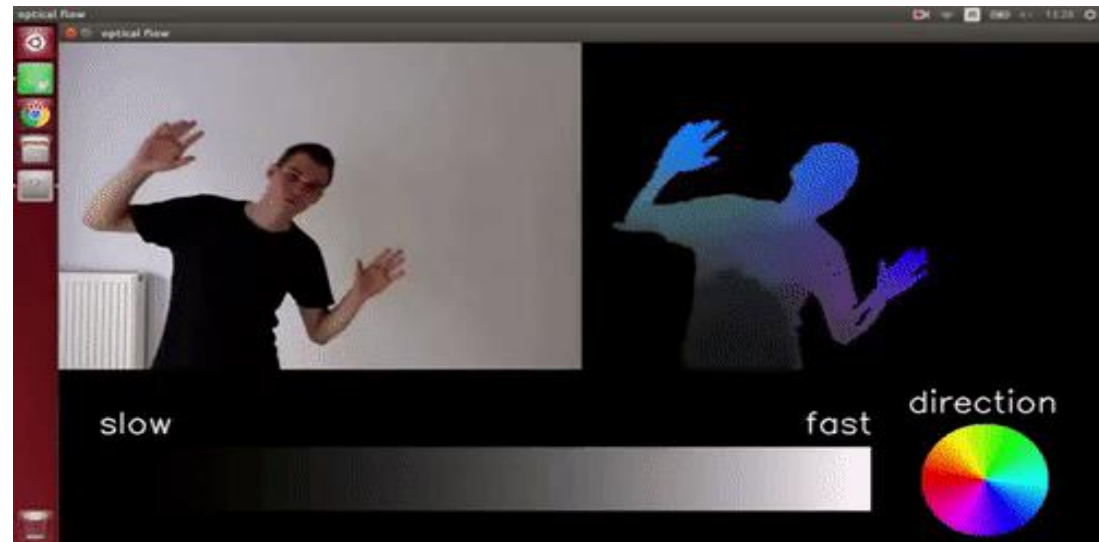
# Core Tasks In Computer Vision

- **Object classification**
- **Object detection**
- **Image segmentation**
- **Scene reconstruction**
- **Image restoration**
- **Motion analysis**
- **Visual QnA**
- **Image generation**

Any other application that invol
can safely be labeled as compute

By Anudeepsekhar Bolimera

@MakeSmartThings     THEASSEMBLY.AE

THE**ASSEMBLY**
MAKE | SMART | THINGS

# Introduction to OpenCV

- OpenCV stands for *__Open Source Computer Vision__* library and it was invented by Intel in 1999.

- Has over **2500 modules** to help you with image-processing and other computer vision tasks

- Free and open-source.

- **Multiplatform** - Windows, Linux, MacOS, Android, IOS.

- **Multilingual** - Python, Java, JavaScript, etc.

# Basics Of OpenCV

- **Getting started with images and OpenCV**: First step – import an image into OpenCV

- **Getting started with video**: Then we modify our code to work with video (streams of images)

- **Changing Color spaces**: We will explore different color spaces and learn how we can change between them using OpenCV

- **Working with Drawing functions**: We will use OpenCV to add text and draw on our frames

- **Geometric Transformations**: We learn how to apply perform basic geometric transformations to our images

# Image Analysis

How does the human brain recognize an object in an image?

Elongated oval head with orange beak

Oval-shaped white blob

Long white rectangular neck

COMMUNITY INNOVATION WORKSPACE

تركيب
THEASSEMBLY
MAKE | SMART | THINGS

# Image Analysis – Naïve Detector

Features

Elongated oval head with orange beak

Oval-shaped white blob

Long white rectangular neck

Naïve Detector

Is a Swan

Is not a Swan

THE ASSEMBLY
MAKE | SMART | THINGS

# Let's Test Our Classifier

# Better Feature Definitions?

Researchers built multiple computer vision techniques to deal with issues

- SIFT
- FAST
- SURF
- BRIEF - etc

The problem?

- The detectors were too general OR
- Too over engineered to generalize

# Can These Features Be Learned?

# Introducing Deep Learning

## Artificial Intelligence

### Machine Learning

#### Deep Learning

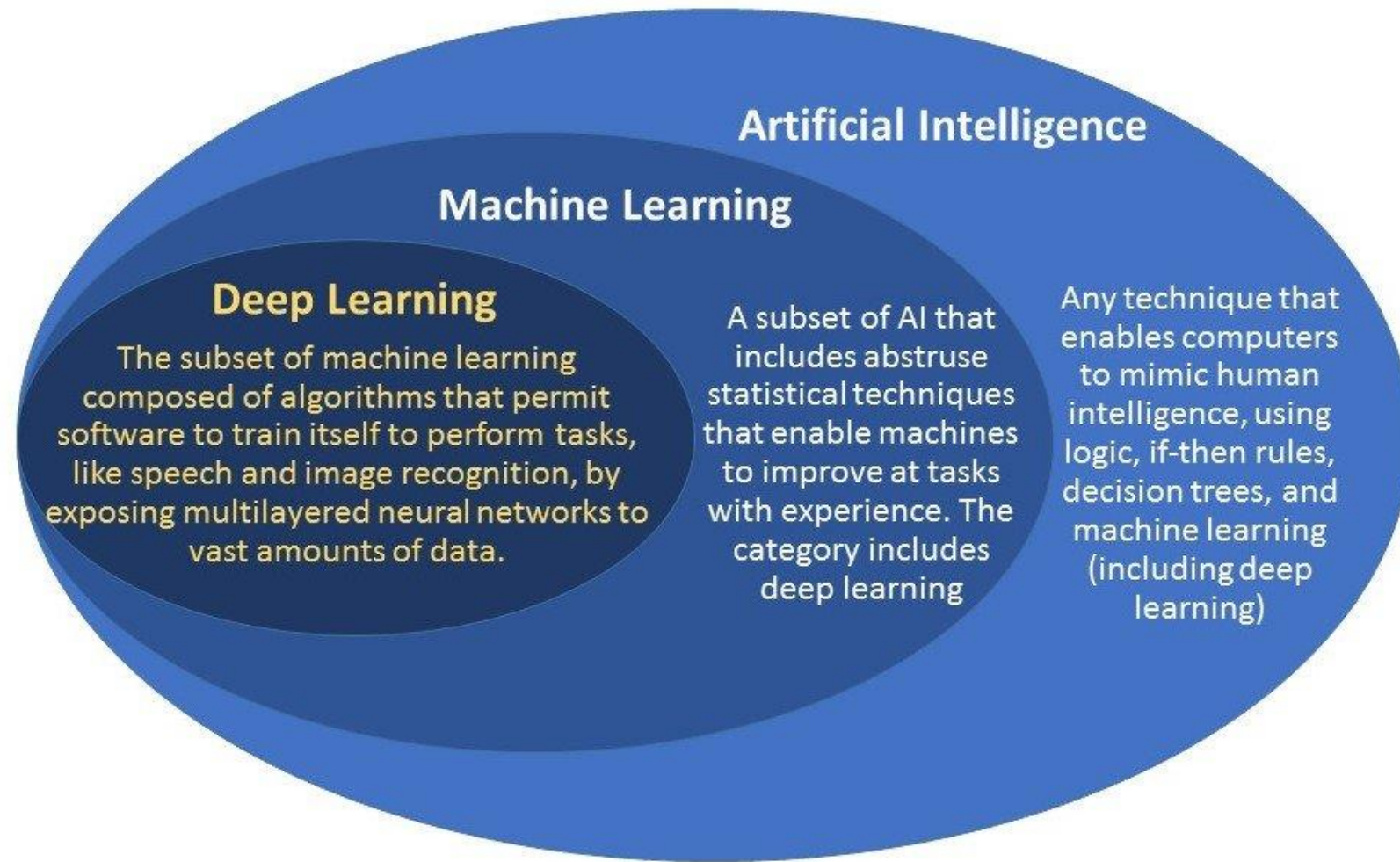The subset of machine learning composed of algorithms that permit software to train itself to perform tasks, like speech and image recognition, by exposing multilayered neural networks to vast amounts of data.

A subset of AI that includes abstruse statistical techniques that enable machines to improve at tasks with experience. The category includes deep learning

Any technique that enables computers to mimic human intelligence, using logic, if-then rules, decision trees, and machine learning (including deep learning)

THE ASSEMBLY
MAKE | SMART | THINGS
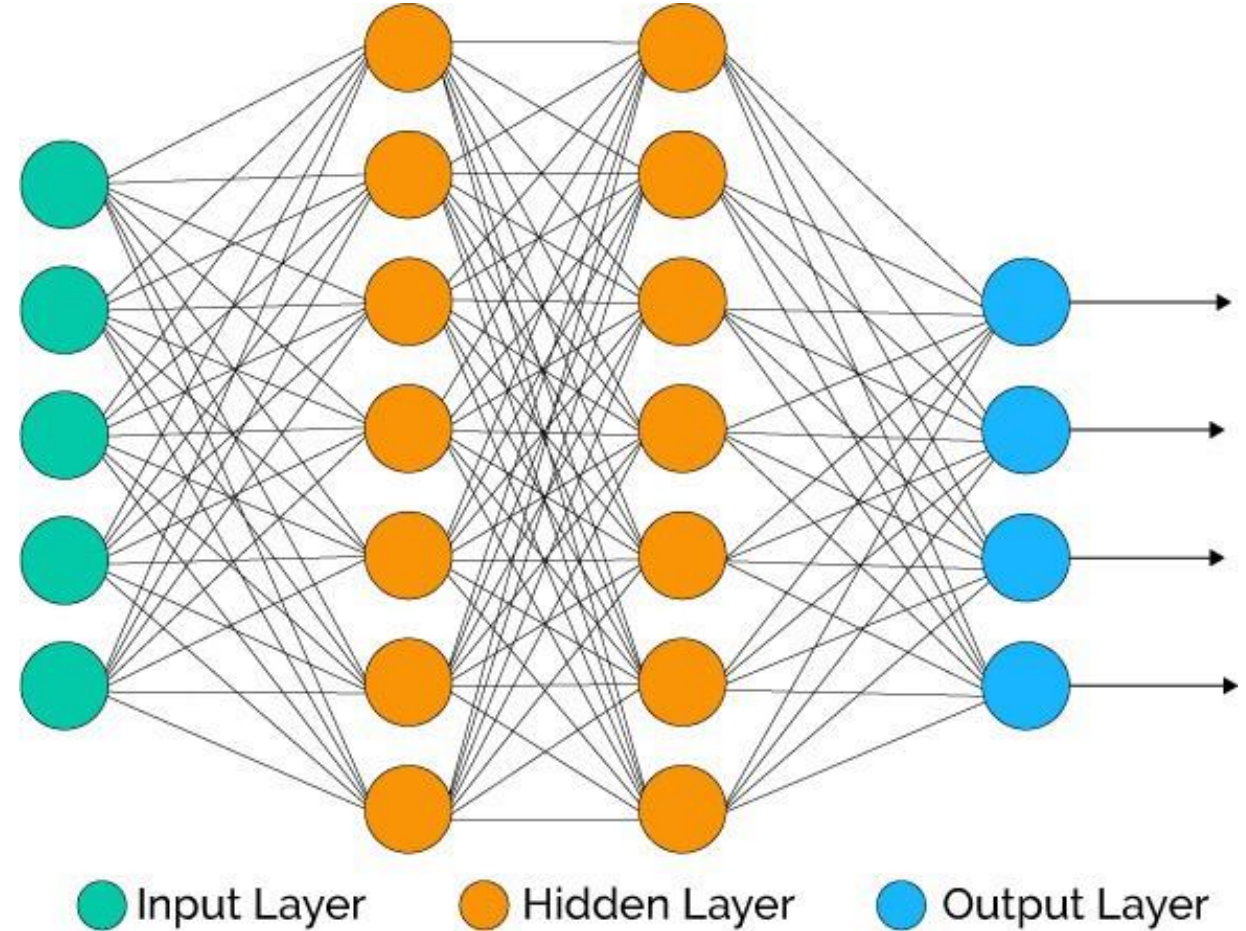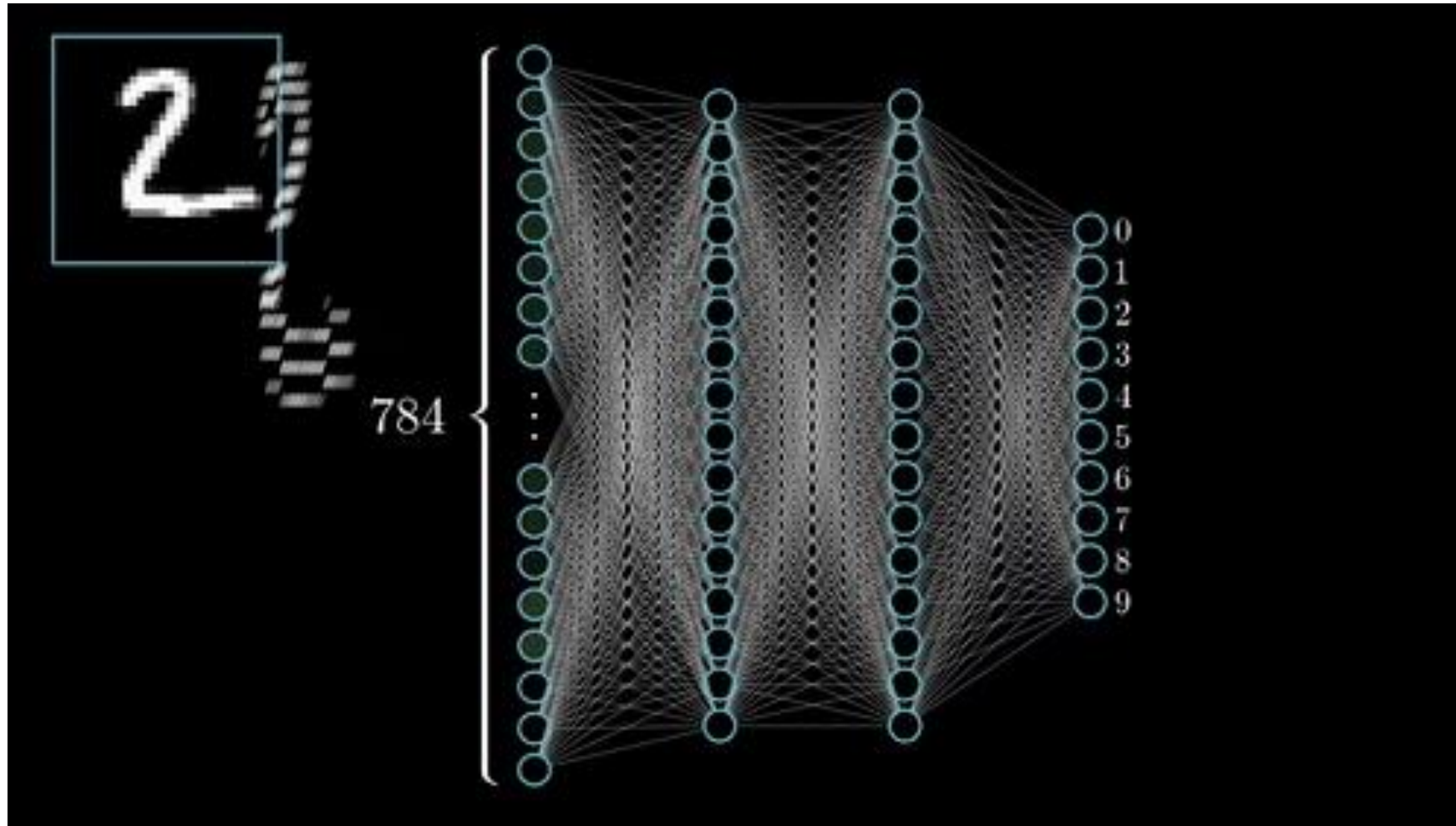
# What Are Neural Networks?

- Neural networks are a set of algorithms, modeled loosely after the human brain, that are designed to recognize patterns
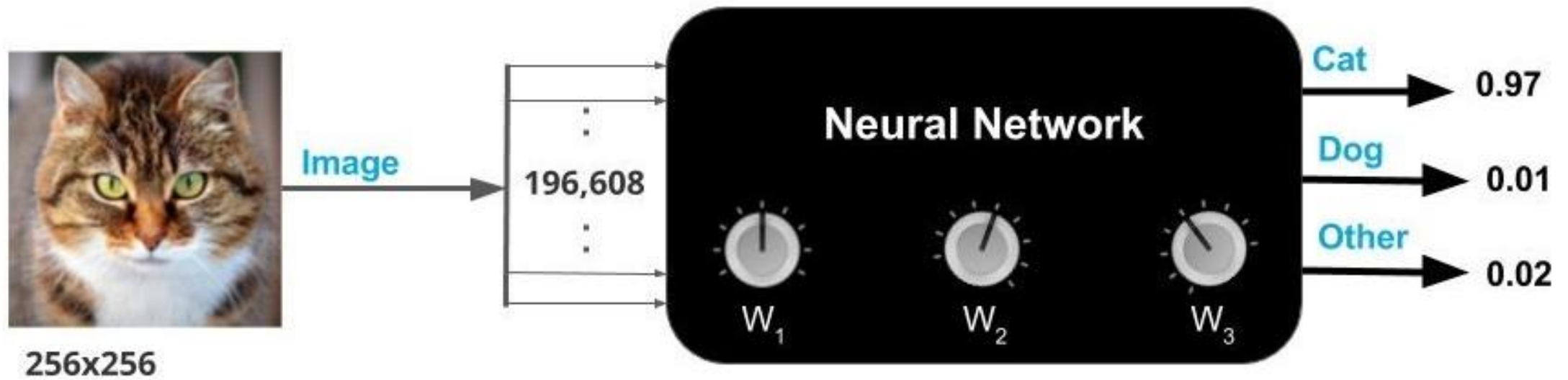


● Input Layer    ● Hidden Layer    ● Output Layer

# Introducing Artificial Neural Networks

# The Knobs Analogy

256x256

Image

196,608

Neural Network

$W_1$   $W_2$   $W_3$

Cat   0.97

Dog   0.01

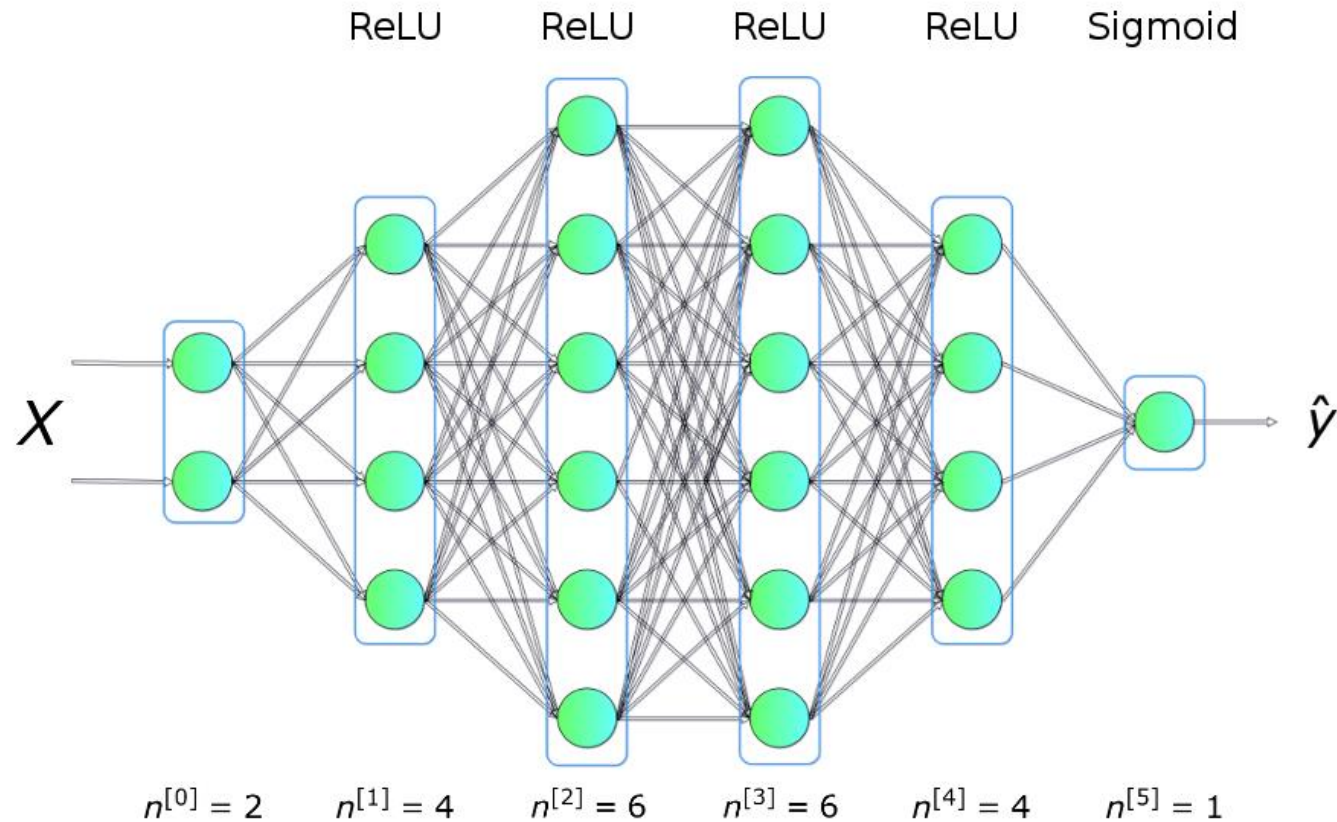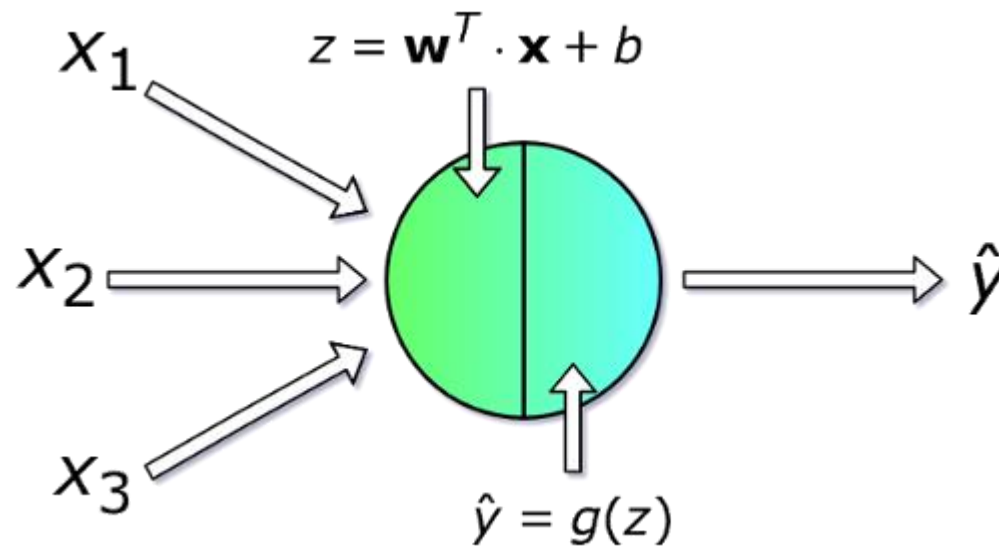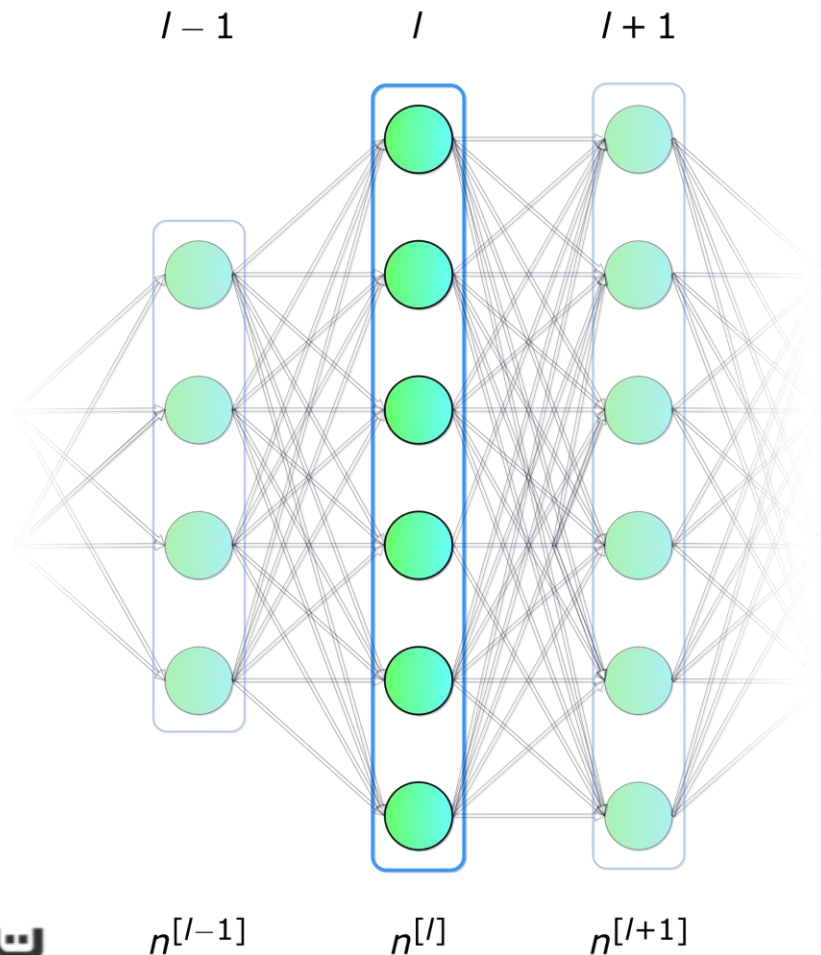Other   0.02

# The Neural Network Architecture

# The Artificial Neuron

$$z = w_1x_1 + w_2x_2 + w_3x_3 + \ldots + w_nx_n = \mathbf{w}^T \cdot \mathbf{x}$$

# Single layer

$l-1$  $l$  $l+1$



$n^{[l-1]}$  $n^{[l]}$  $n^{[l+1]}$

$$z_i^{[l]} = \mathbf{w}_i^T \cdot \mathbf{a}^{[l-1]} + b_i \qquad a_i^{[l]} = g^{[l]}(z_i^{[l]})$$

$$z_1^{[2]} = \mathbf{w}_1^T \cdot \mathbf{a}^{[1]} + b_1 \qquad a_1^{[2]} = g^{[2]}(z_1^{[2]})$$
$$z_2^{[2]} = \mathbf{w}_2^T \cdot \mathbf{a}^{[1]} + b_2 \qquad a_2^{[2]} = g^{[2]}(z_2^{[2]})$$
$$z_3^{[2]} = \mathbf{w}_3^T \cdot \mathbf{a}^{[1]} + b_3 \qquad a_3^{[2]} = g^{[2]}(z_3^{[2]})$$
$$z_4^{[2]} = \mathbf{w}_4^T \cdot \mathbf{a}^{[1]} + b_4 \qquad a_4^{[2]} = g^{[2]}(z_4^{[2]})$$
$$z_5^{[2]} = \mathbf{w}_5^T \cdot \mathbf{a}^{[1]} + b_5 \qquad a_5^{[2]} = g^{[2]}(z_5^{[2]})$$
$$z_6^{[2]} = \mathbf{w}_6^T \cdot \mathbf{a}^{[1]} + b_6 \qquad a_6^{[2]} = g^{[2]}(z_6^{[2]})$$
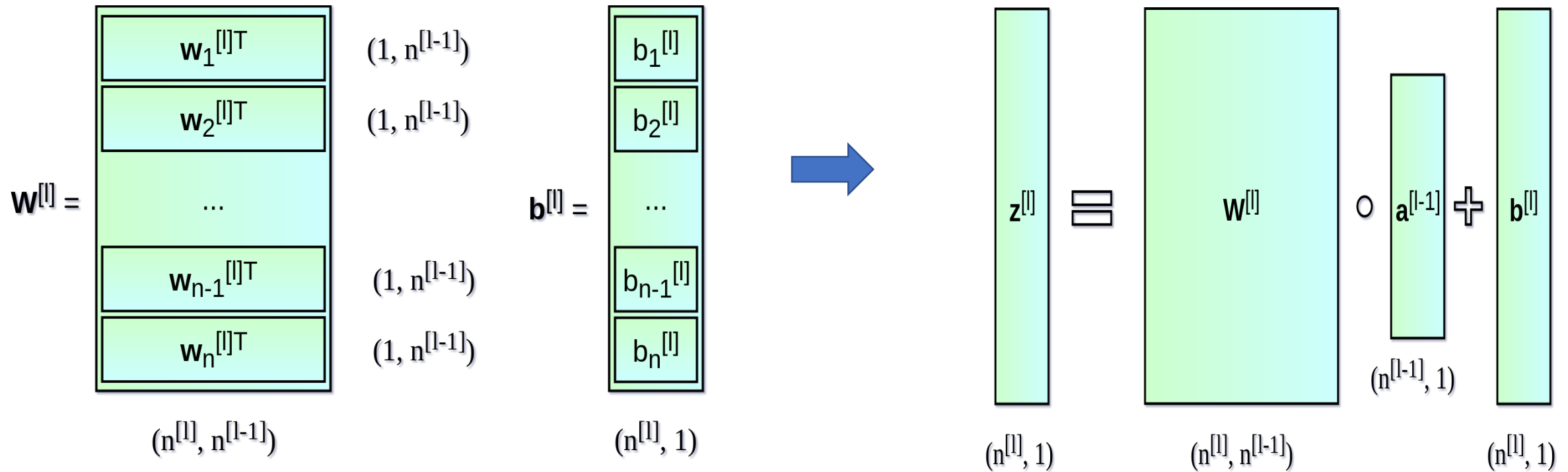
# Vectorizing the equations

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \cdot \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]} \qquad \mathbf{a}^{[l]} = g^{[l]}(\mathbf{z}^{[l]})$$



$$\mathbf{W}^{[l]} = \begin{array}{c} \mathbf{w}_1^{[l]T} \\ \mathbf{w}_2^{[l]T} \\ \ldots \\ \mathbf{w}_{n-1}^{[l]T} \\ \mathbf{w}_n^{[l]T} \end{array} \qquad \begin{array}{c} (1, n^{[l-1]}) \\ (1, n^{[l-1]}) \\ \\ (1, n^{[l-1]}) \\ (1, n^{[l-1]}) \end{array}$$

$$(n^{[l]}, n^{[l-1]})$$

$$\mathbf{b}^{[l]} = \begin{array}{c} b_1^{[l]} \\ b_2^{[l]} \\ \ldots \\ b_{n-1}^{[l]} \\ b_n^{[l]} \end{array}$$

$$(n^{[l]}, 1)$$

$$\mathbf{z}^{[l]} \; \boxed{=} \; \mathbf{W}^{[l]} \; \circ \; \mathbf{a}^{[l-1]} \; + \; \mathbf{b}^{[l]}$$

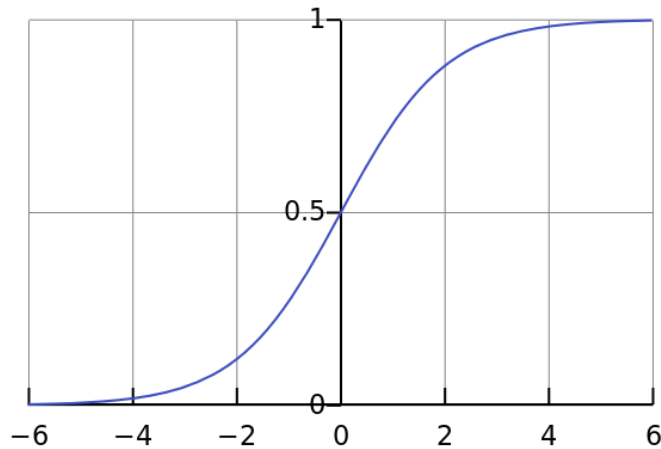$$(n^{[l]}, 1) \qquad\qquad (n^{[l]}, n^{[l-1]}) \qquad (n^{[l-1]}, 1) \qquad (n^{[l]}, 1)$$

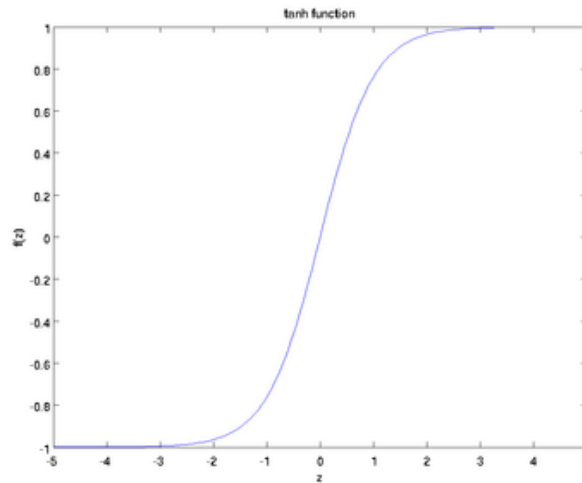# Activation Function And Why Do We Need It?

- **Without them, our neural network would become a combination of linear functions, so it would be just a linear function itself**

- The activation function also has a significant impact on the speed of learning, which is one of the main criteria for their selection
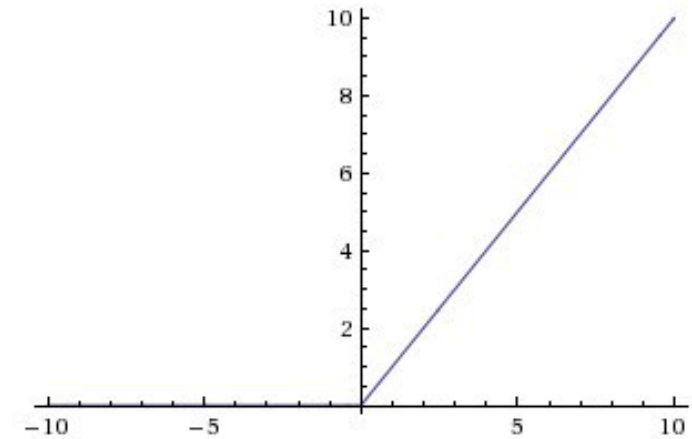
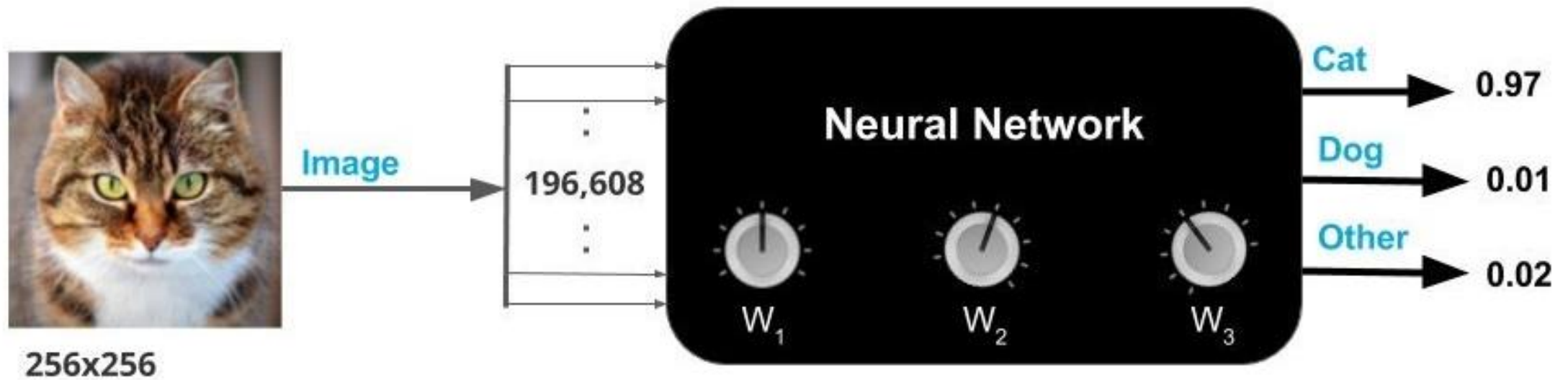# Activation Function And Why Do We Need It?

$$A = \frac{1}{1+e^{-x}}$$



tanh function

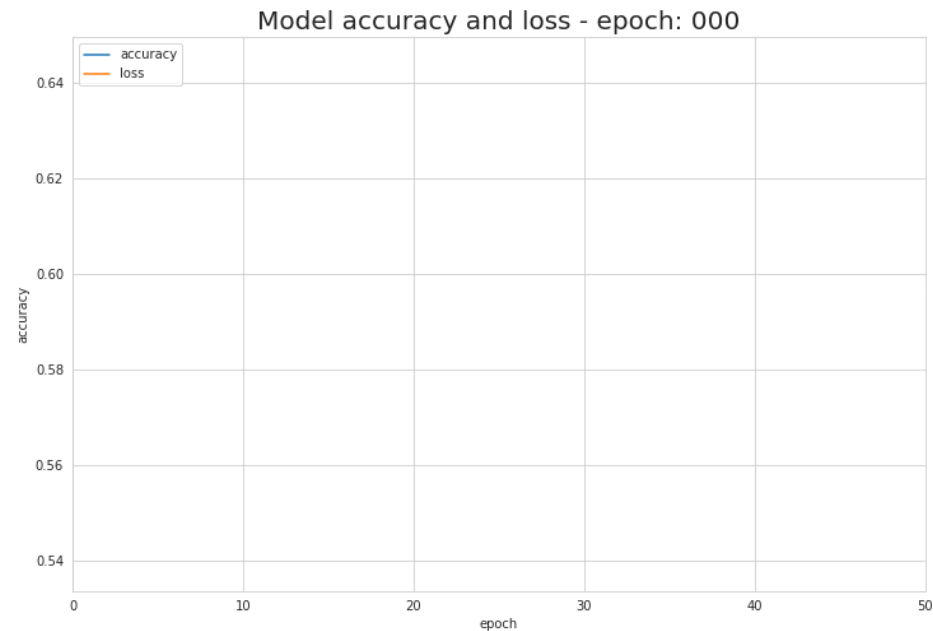$$f(x) = tanh(x) = \frac{2}{1+e^{-2x}} - 1$$
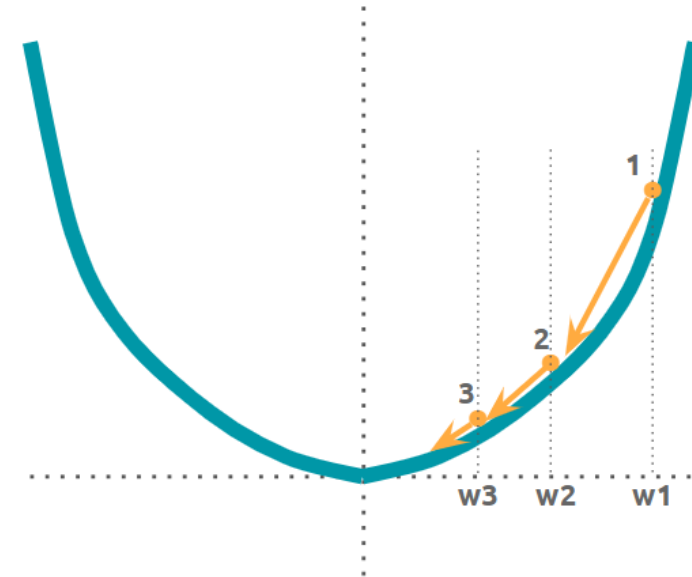


$$A(x) = max(0,x)$$

# The Knobs Analogy

# The Loss Function

- The basic source of information on the progress of the learning process is the value of the loss function.



Model accuracy and loss - epoch: 000

# How do neural networks learn?

- The learning process is about changing the values of the **W** and **b** parameters so that the loss function is minimized. In order to achieve this goal, we will turn for help to calculus and **use gradient descent method to find a function minimum**.



$$W = W - \alpha \times dW$$

W2 = W1 - α*dW1
W3 = W2 - α*dW2

# Backpropagation

- Backpropagation is an algorithm that allows us to calculate a very complicated gradient, like the one we need.

- The parameters of the neural network are adjusted according to the following formulae.

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha \mathbf{dW}^{[l]}$$
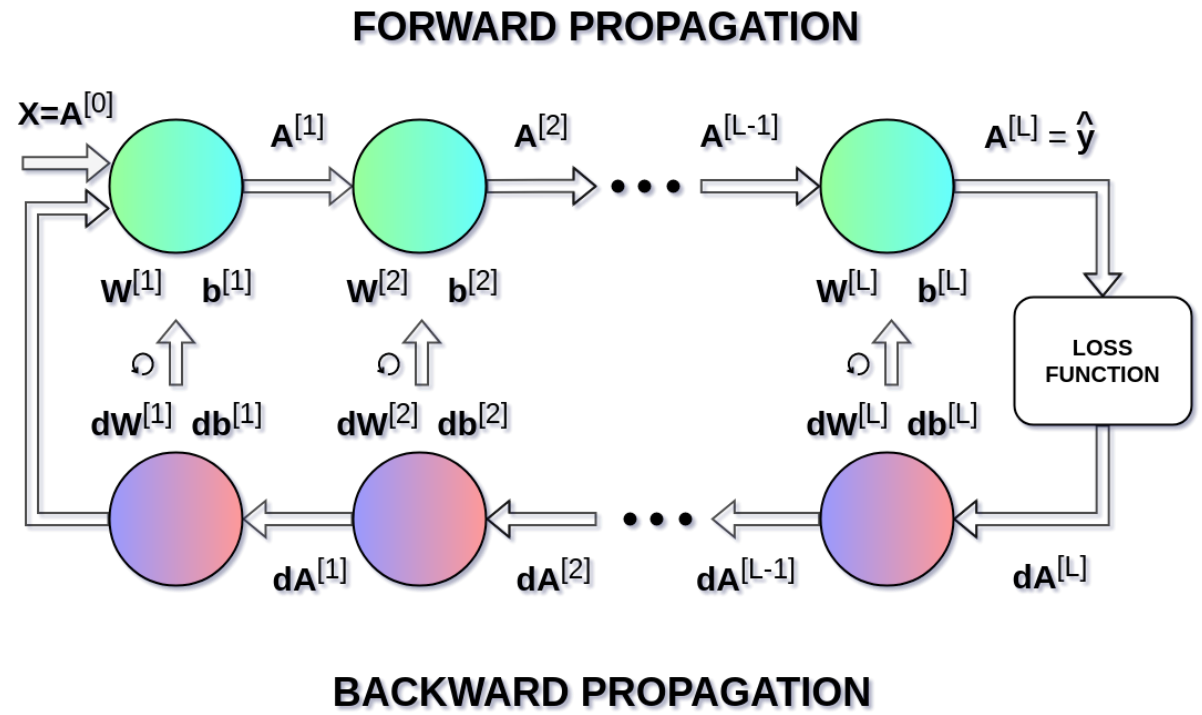$$\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \mathbf{db}^{[l]}$$

- In the equations above, α represents learning rate - a hyperparameter which allows you to control the value of performed adjustment

# The Training Process

- **dW** and **db** are calculated using the chain rule, partial derivatives of loss function with respect to **W** and **b**.

$$dW^{[l]} = \frac{\partial L}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]T}$$

$$db^{[l]} = \frac{\partial L}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^{m} dZ^{[l](i)}$$

$$dA^{[l-1]} = \frac{\partial L}{\partial A^{[l-1]}} = W^{[l]T} dZ^{[l]}$$

$$dZ^{[l]} = dA^{[l]} * g'(Z^{[l]})$$



**FORWARD PROPAGATION**

**BACKWARD PROPAGATION**

# Let's Build Our First Neural Network

COMMUNITY
INNOVATION
WORKSPACE

THEASSEMBLY

MAKE | SMART | THINGS
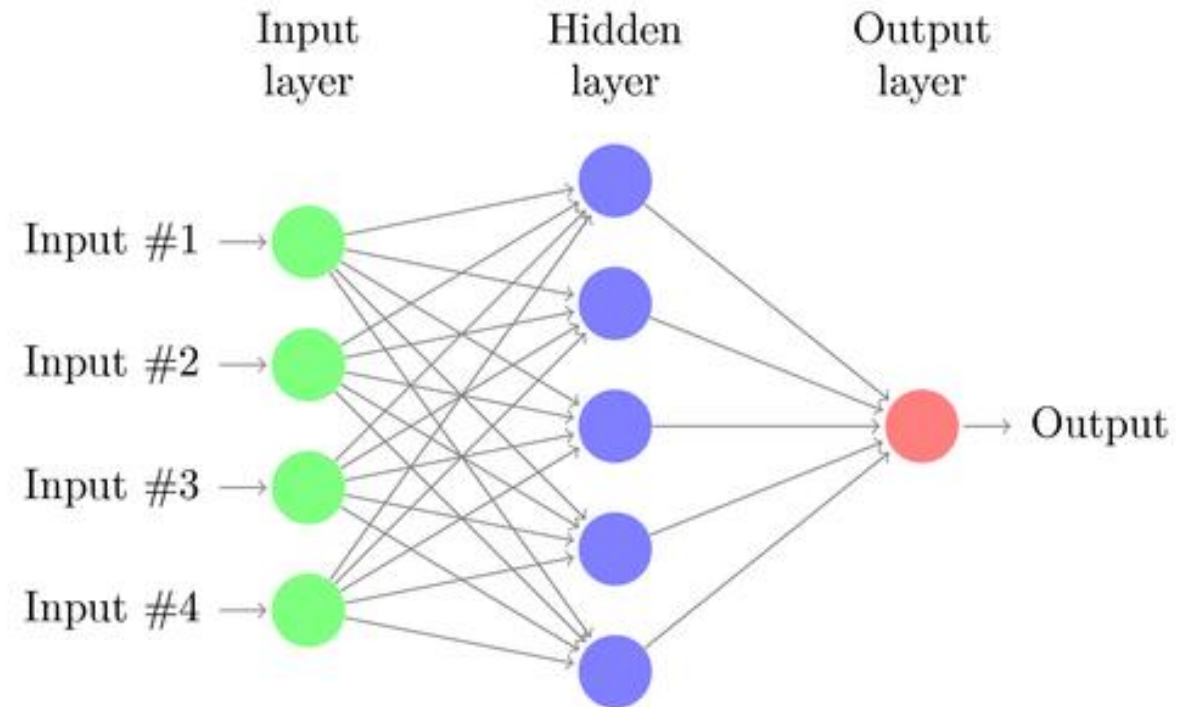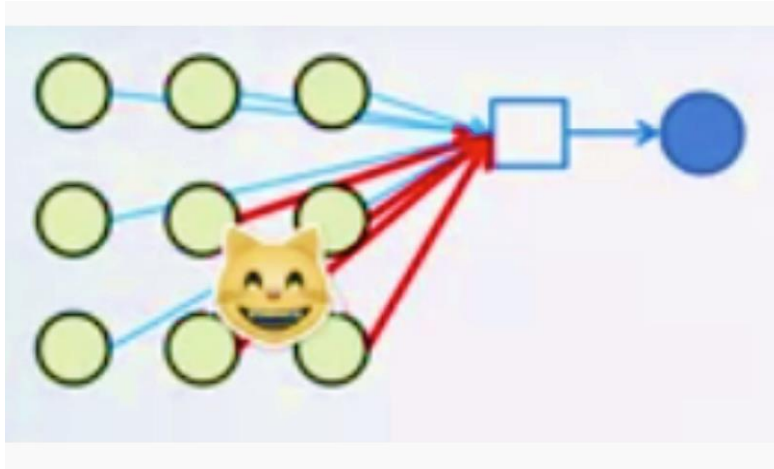
# The Problem With Traditional Neural Networks

- MLPs use one perceptron for each input (e.g. pixel in an image, multiplied by 3 in RGB case)

- The amount of weights rapidly becomes unmanageable for large images

- For a 224 x 224 pixel image with 3 color channels there are around 150,000 weights that must be trained!
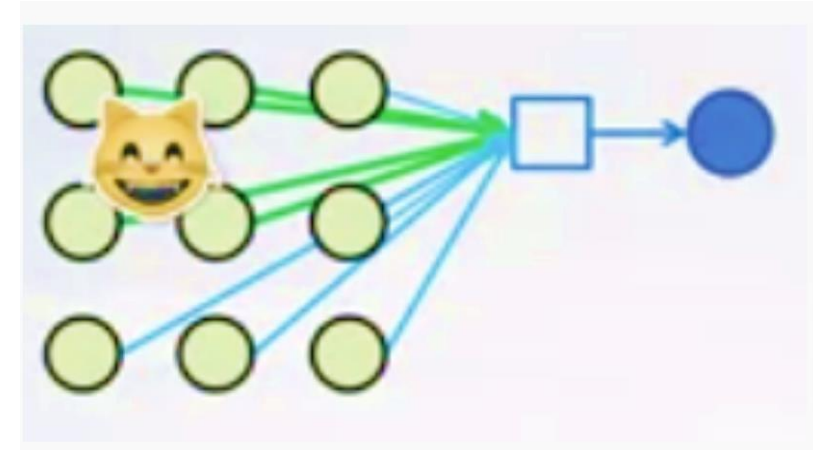
# The Problem With Traditional Neural Networks

- MLPs react differently to an input (images) and its shifted version — **they are NOT translation invariant**.

- MLPs do not consider spatial information as each pixel is considered as an individual input



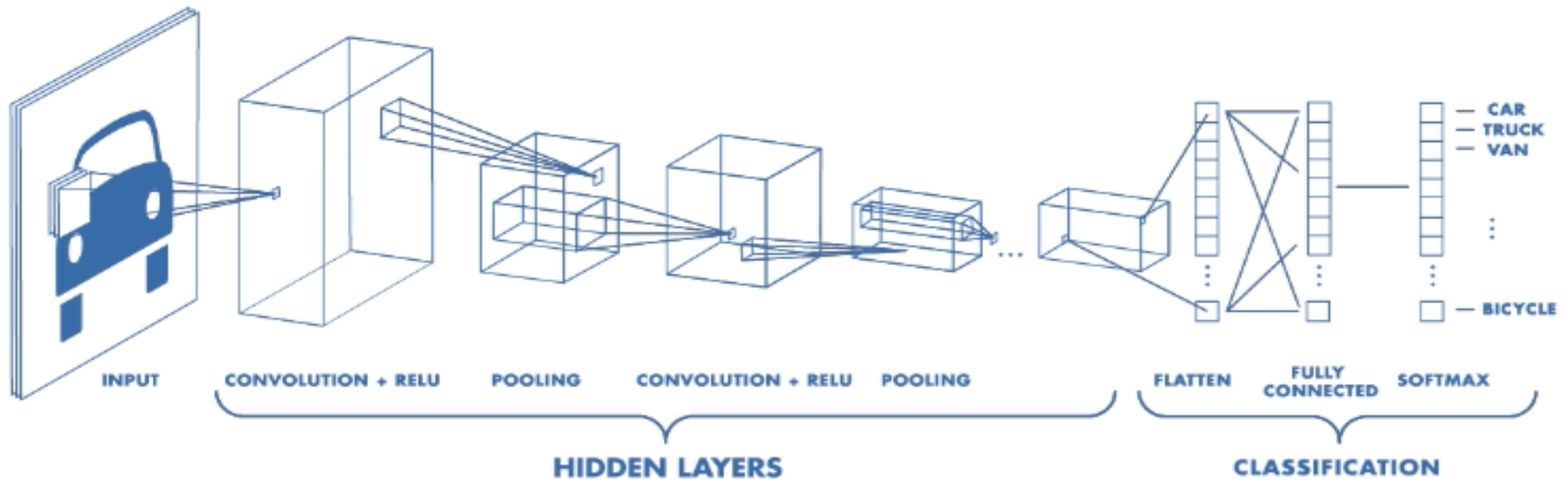In this case red weights will get modified to better recognize cats



In this case green weights will get modified to better recognize cats

# Enter the Convolutional Neural Network

# Core Idea Of Convolutional Neural Networks

- Our brain captures the patterns in figures to classify an object

- What if our machine can be trained to classify images by detecting some patterns like us?

- CNNs leverage the fact that nearby pixels are more strongly related than distant ones and analyze the influence of nearby pixels by using something called a **filter**

- In simpler terms, CNNs use filters to find patterns in an image just like our brain does

# What Are Filters?

- A filter is nothing more than a small 2-dimensional array of values which represents a pattern that we are looking for`.

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Pixel representation of a filter

Visualization of a filter

The sum of the multiplication value that is generated is = 4*(50*30)+(20*30) = 6600 (large number)



Visualization of the receptive field

Pixel representation of the receptive field

Pixel representation of filter

Large Value

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

The sum of the multiplication value that is generated is = 0 (small number)

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 40 | 0 | 0 | 0 | 0 | 0 |
| 40 | 0 | 40 | 0 | 0 | 0 | 0 |
| 40 | 20 | 0 | 0 | 0 | 0 | 0 |
| 0 | 50 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 50 | 0 | 0 | 0 | 0 |
| 25 | 25 | 0 | 50 | 0 | 0 | 0 |

\*

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Small Value

Visualization of the filter on the image

Pixel representation of receptive field

Pixel representation of filter

THEASSEMBLY
MAKE | SMART | THINGS

# What is Convolution?

- By definition, **convolution** is a mathematical operation on two objects to produce an outcome that expresses how the shape of one is modified by the other.

- With this computation, we detect a particular feature from the input image and get the result having information about that feature

- This is called a '***feature map***.' Outcome with a real world example below.

Input image

Convolution Kernel

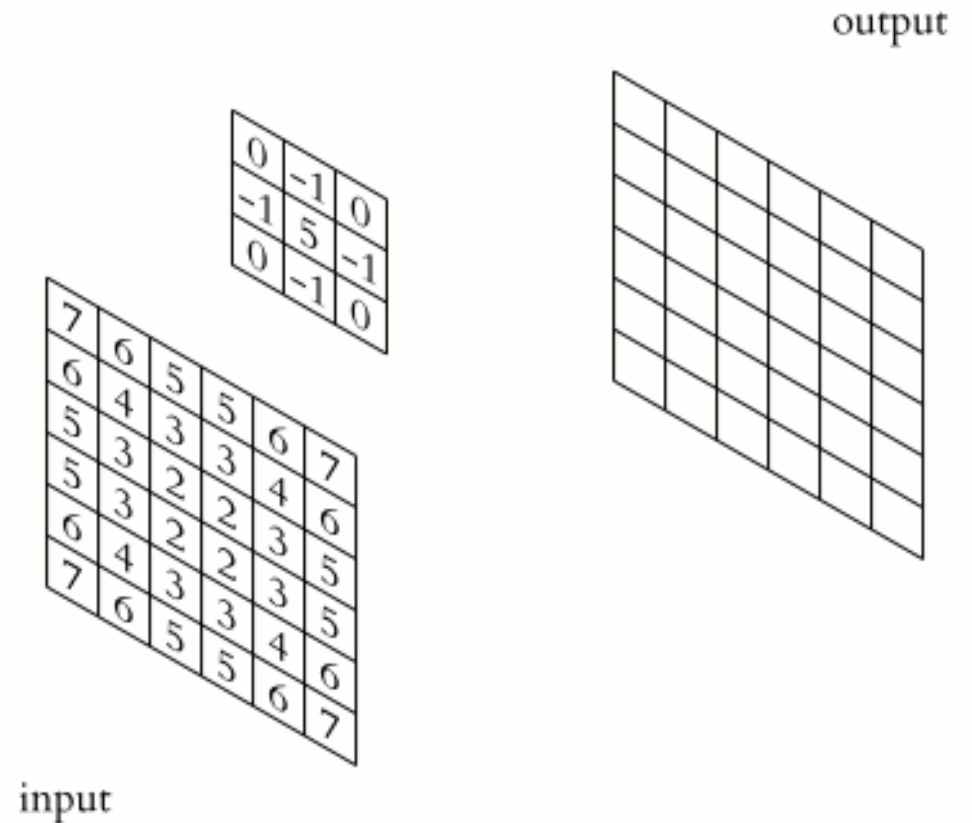$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

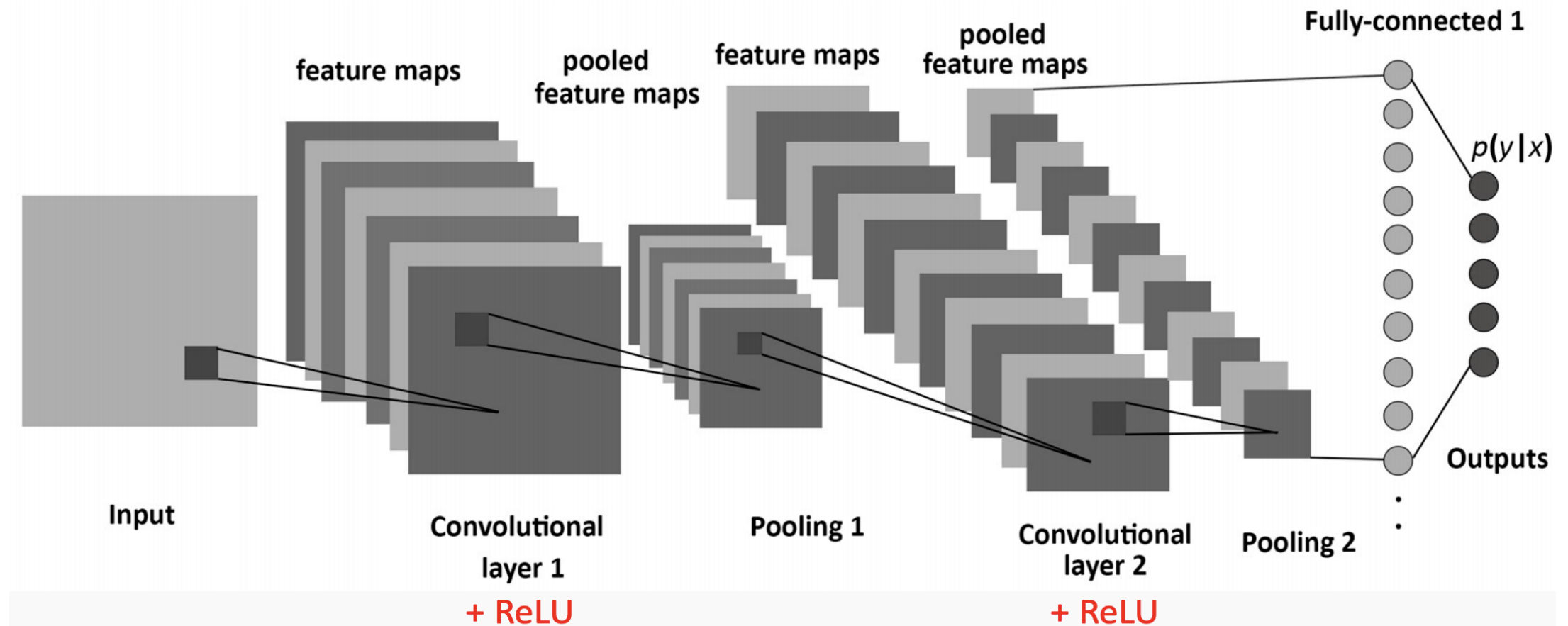Feature map

# How Does Convolution Work?

The output image, C, is computed from the input image, A, and the convolution kernel, H, according to the following expression.

$$C[m][n] = \sum_u \sum_v (A[m+u][n+v] \cdot H[u][v])$$



output

input

# The Vanilla CNN

# Layers Of A CNN

- We use three main types of layers to build CNN architectures: **Convolutional Layer**, **Pooling Layer**, and **Fully-Connected Layer**

- We can stack these layers to form a full CNN **architecture**

- **INPUT** will hold the raw pixel values of the image

- **CONV** layer will compute the output of neurons that are connected to local regions in the input

- **RELU** layer will apply an elementwise activation function, such as the max(0,x) thresholding at zero

- **POOL** layer will perform a downsampling operation along the spatial dimensions (width, height)

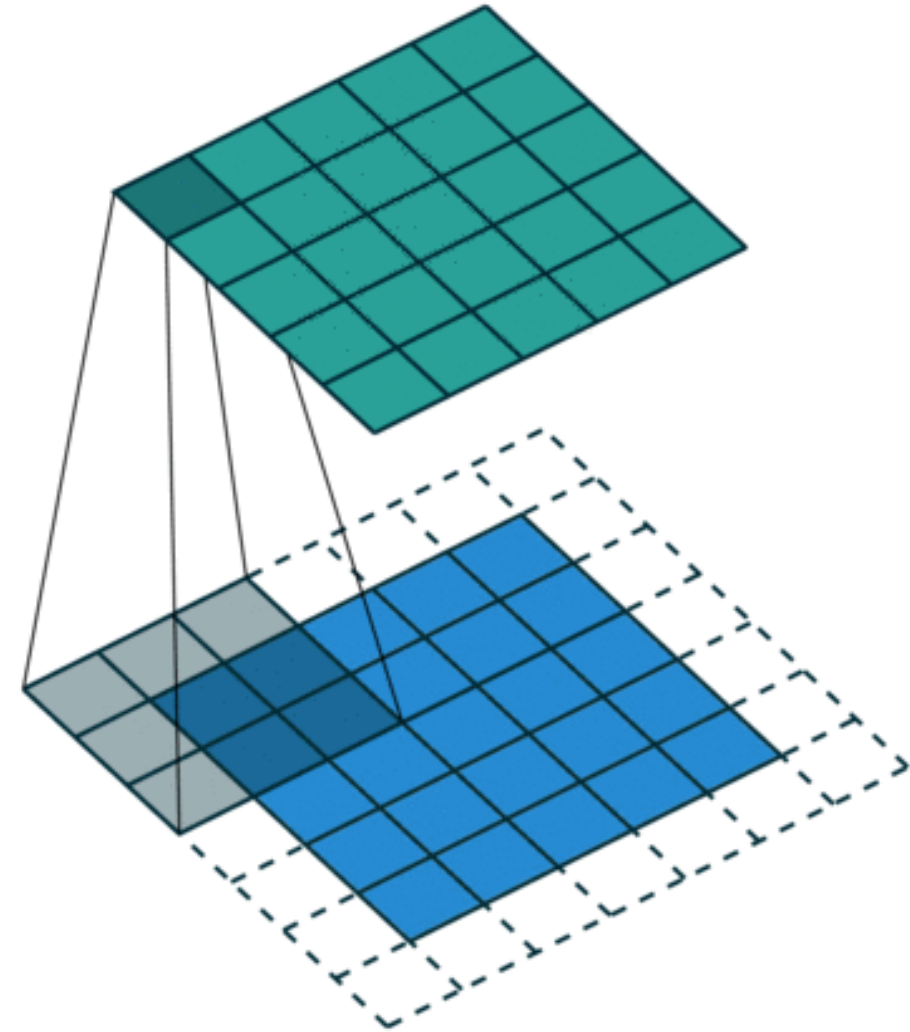- **FC** (i.e. fully-connected) layer will compute the class scores

By Anudeepsekhar Bolimera

# Convolution Layer

- Accepts a volume of size W1×H1×D1

- Requires four hyperparameters:
  - Number of filters K
  - their spatial extent F,
  - the stride S,
  - the amount of zero padding P.

- Produces a volume of size W2×H2×D2
  where:
  - W2=(W1−F+2P)/S+1
  - H2=(H1−F+2P)/S+1
  (i.e. width and height are computed equally by symmetry)
  - D2=K
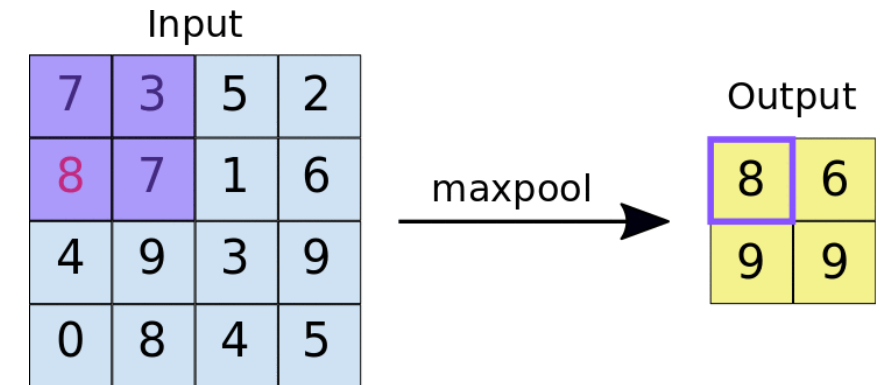
# Padding And Convolution

- If we apply convolutions on a normal image, the result will be down-sampled by an amount depending on the size of the filter

- Padding essentially makes the feature maps produced by the filter kernels the same size as the original image

By Anudeepsekhar Bolimera

# Pooling Layer

- Accepts a volume of size $W1 \times H1 \times D1$
- Requires two hyperparameters:
  - their spatial extent F,
  - the stride S,
- Produces a volume of size $W2 \times H2 \times D2$
  where:
  - $W2 = (W1 - F)/S + 1$
  - $H2 = (H1 - F)/S + 1$
  - $D2 = D1$
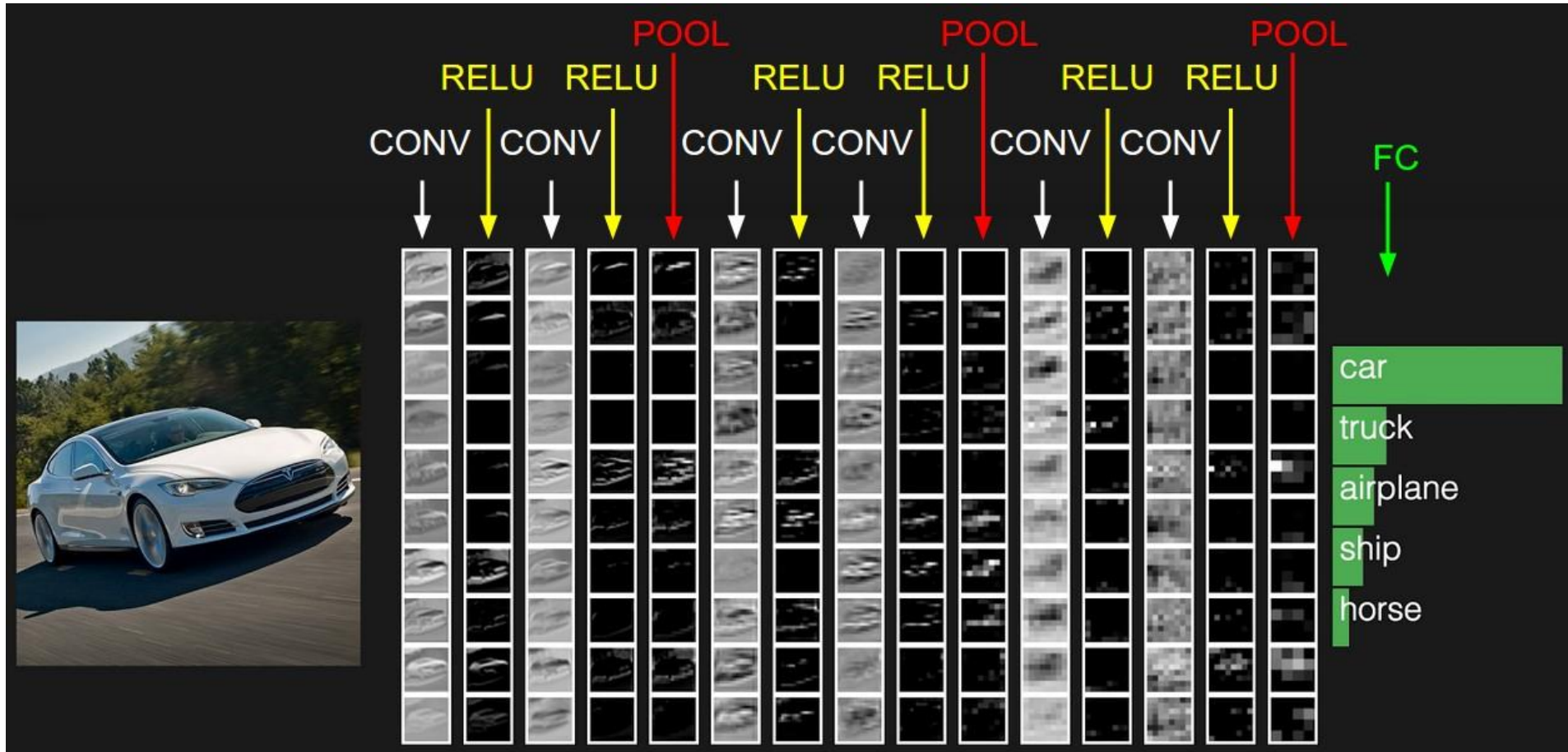- Introduces zero parameters since it computes a fixed function of the input

Input

| 7 | 3 | 5 | 2 |
|---|---|---|---|
| 8 | 7 | 1 | 6 |
| 4 | 9 | 3 | 9 |
| 0 | 8 | 4 | 5 |

maxpool →

Output

| 8 | 6 |
|---|---|
| 9 | 9 |

# Activation Layer And Relu

# Inside The Layers Of A CNN

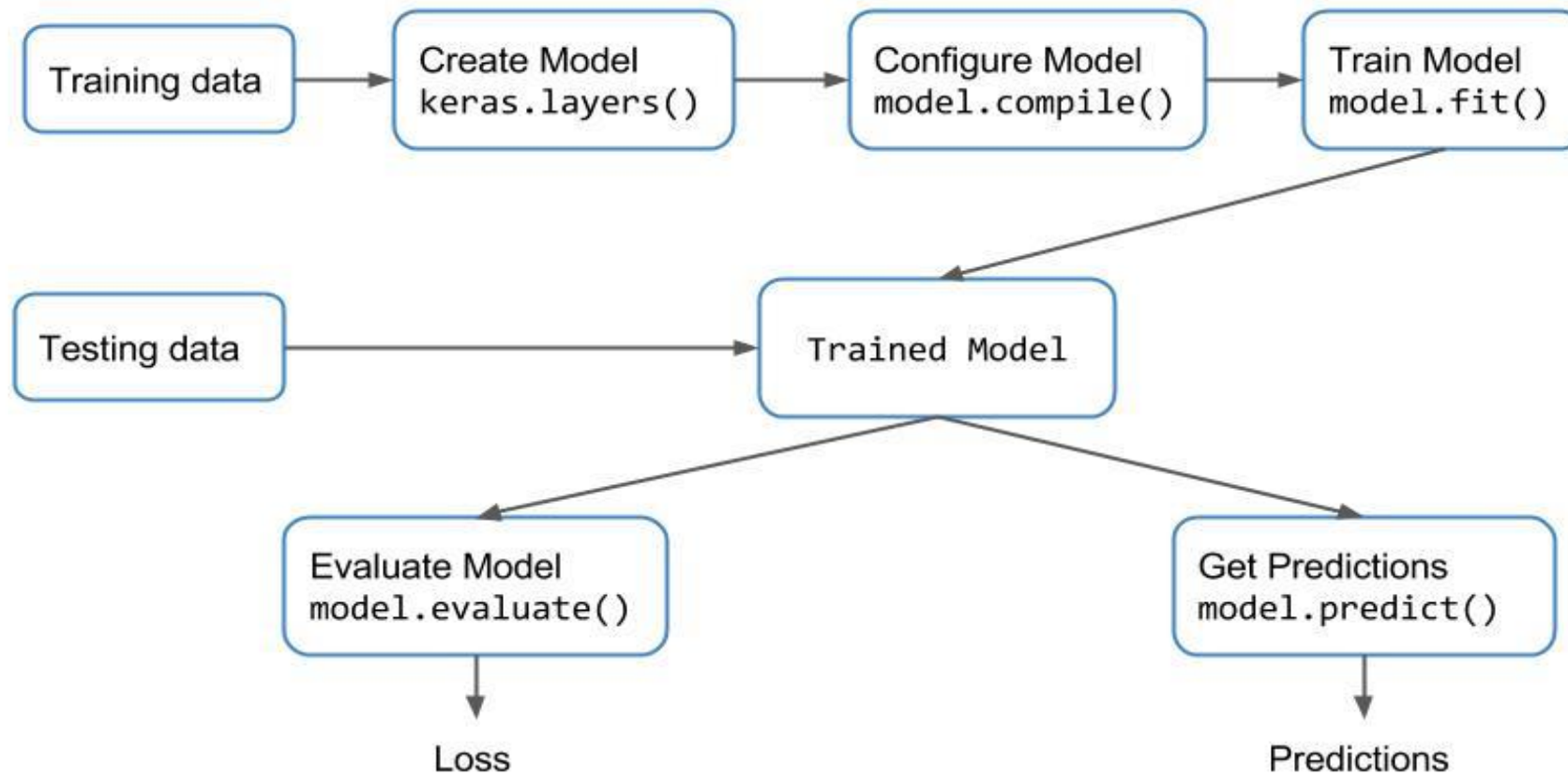# 3D CNN Demo
http://www.scs.Ryerson.ca/~aharley/vis/conv/

# Keras Workflow for training the network

# Let's Build Your Own CNN with Keras

# Thank You