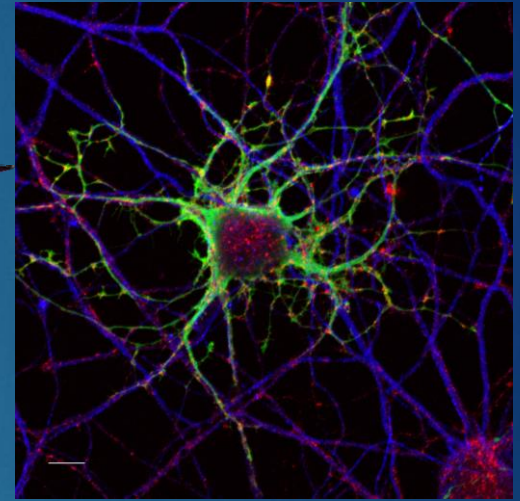


# Convolutional Neural Networks

INTRODUCTION TO IMAGE CLASSIFICATION  
(MODULE 6 – VISION MODELS) – By Christelle JULIAS



# Content

- ❑ Discovery and Overview
- ❑ [Computer Vision - Images ]
- ❑ Use Cases of CNNs IRL
  - ❑ Image Classification
  - ❑ Object Recognition
- ❑ What are Convolutional Neural Networks?
- ❑ How do they work?
  - ❑ Core Components: Convolution, Filters, and Pooling layers
- ❑ Let's explore!
- ❑ How do we manipulate Convolutional Neural Networks?
  - ❑ Creation
  - ❑ Training
  - ❑ Fine-tuning
  - ❑ A bit of Maths

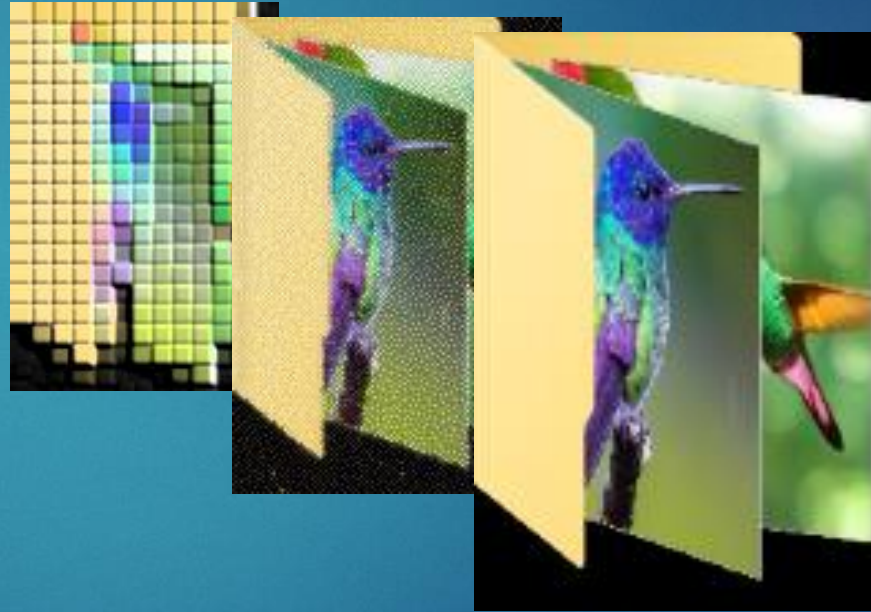
# Muffins or Chihuahuas?





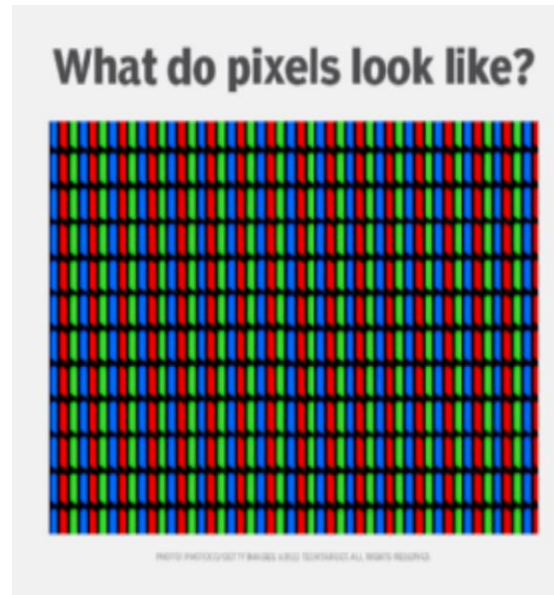
# Introduction

WHAT CAN WE DO WITH IMAGES?



# Images are « sets of pixels »

1. A **Pixel** (short for "picture element") : smallest unit of a digital image or display.
2. Represents a **single point** in an image. Measure of size for calibration. Larger pixels capture more light, smaller provide finer details.
3. Displays various colors based on its **RGB** (red, green, blue) **components**. In Black & White: binary, in Grayscale: B, W & shades of gray.



- **Composition:** made up of subpixels, typically three for RGB color representation.
- **Resolution:** total number of pixels in an image. Ex: a 1920 x 1080 display has over 2 million pixels.
- **Function:** More pixels generally lead to clearer and more detailed images=> impact the overall quality of visual displays (ex.FullHD)

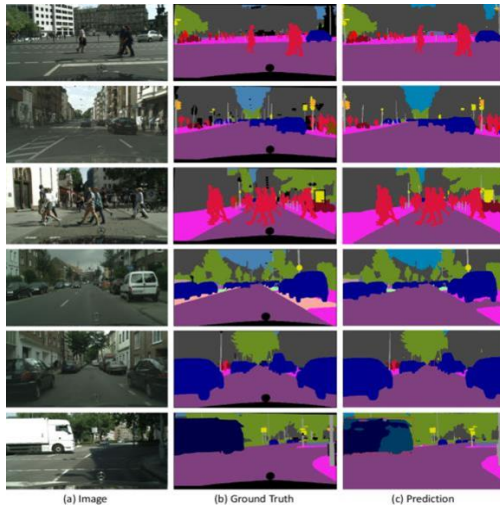
# Image Processing techniques include

- **Enhancement:** parameters such as brightness, contrast, and sharpness
- **Restoration:** revamp degraded original images
- **Segmentation:** divide into "semantic/meaningful" parts according to the pixel distribution.



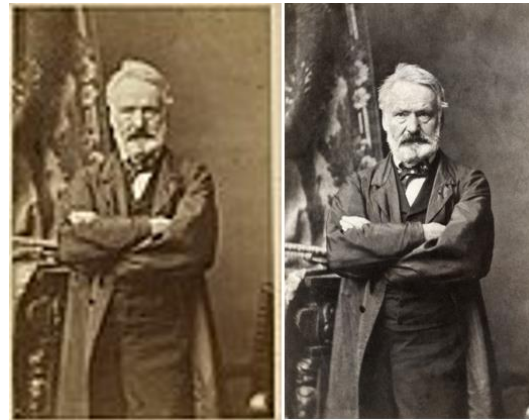
- **Compression:** to reduce the file size through lossy or lossless methods (MPEG, HEIF)
- **Generation:** GANs (Generative Adversarial Networks) to create new images or enhance
- **Morphological Processing:** shape image structures (twist, wrinkle etc.)

# Examples



## Semantic Segmentation

Ex: spatial data within contexts (streets, hours, references with self-driving cars), pixel level classification



## Restorative Filters

Ex: Antique Photos (here Victor Hugo, the French Novelist of the 19th Century, by Bertall 1867)



## Morphological Processing + Generation

Ex: Beauty filters (Instagram, Adobe, PicsArt ...), Background removal and replacement



# Types of Object Recognition Algorithms



## Image classification

- **Classifies** / labels the object contained in the picture (what is it?)
- Provides a **probability** of the guess~(it's 90% likely to be a bird)

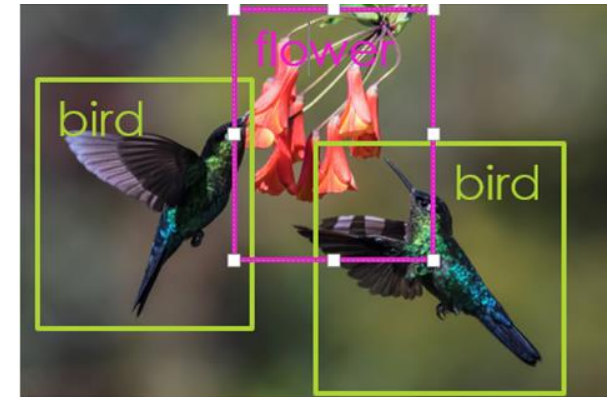
Traditional CNNs



## Both Classification & Localization

- **Detects** the objects and also **localizes** it in the picture simultaneously(what is the bird?)

Simplified Yolo, R- CNNs



## Object Detection

- **Detects** the content of an image, presence of **several objects** (many bounding boxes) (what are those?)

Yolo, R- CNNs



# QUIZZ



Identify real world  
application of Image  
Classification



# Image Classification / Recognition

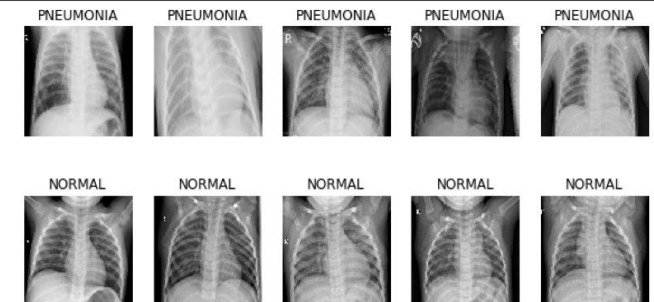
[Is it a [label, name]?]

- Classifying images into predefined **categories**.
- Recognizing **patterns** and **features**.
- Widely used in applications like:
  - **photo tagging** in social media
  - and **medical imaging** for diagnostics identifying pathologies.



```
[64]: def plot(image_batch, label_batch):  
    plt.figure(figsize=(10,5))  
    for i in range(10):  
        ax = plt.subplot(2,5,i+1)  
        img = cv2.imread(str(image_batch[i]))  
        img = cv2.resize(img, (224,224))  
        plt.imshow(img)  
        plt.title(label_batch[i])  
        plt.axis("off")
```

```
[65]: plot(image_batch, label_batch)
```



# Object Detection

(Presence: Is there an [object name] in the image / landscape / supermarket aisle?)

- ➔ - **Identifying** and **locating** objects within a picture simultaneously.
- Applications include **self-driving cars**, **facial recognition**, **security cameras**, **sorting/ recycling wastes** etc

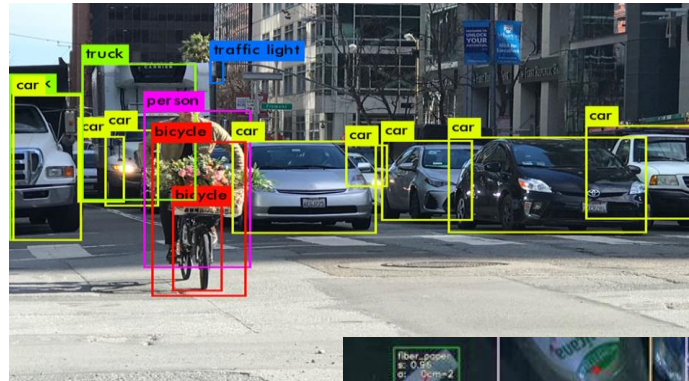
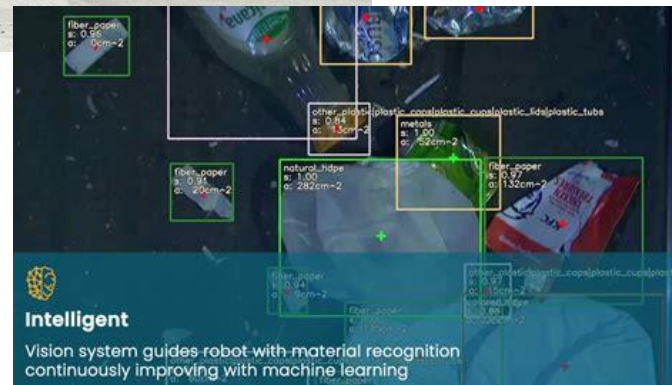


Image: ywseo

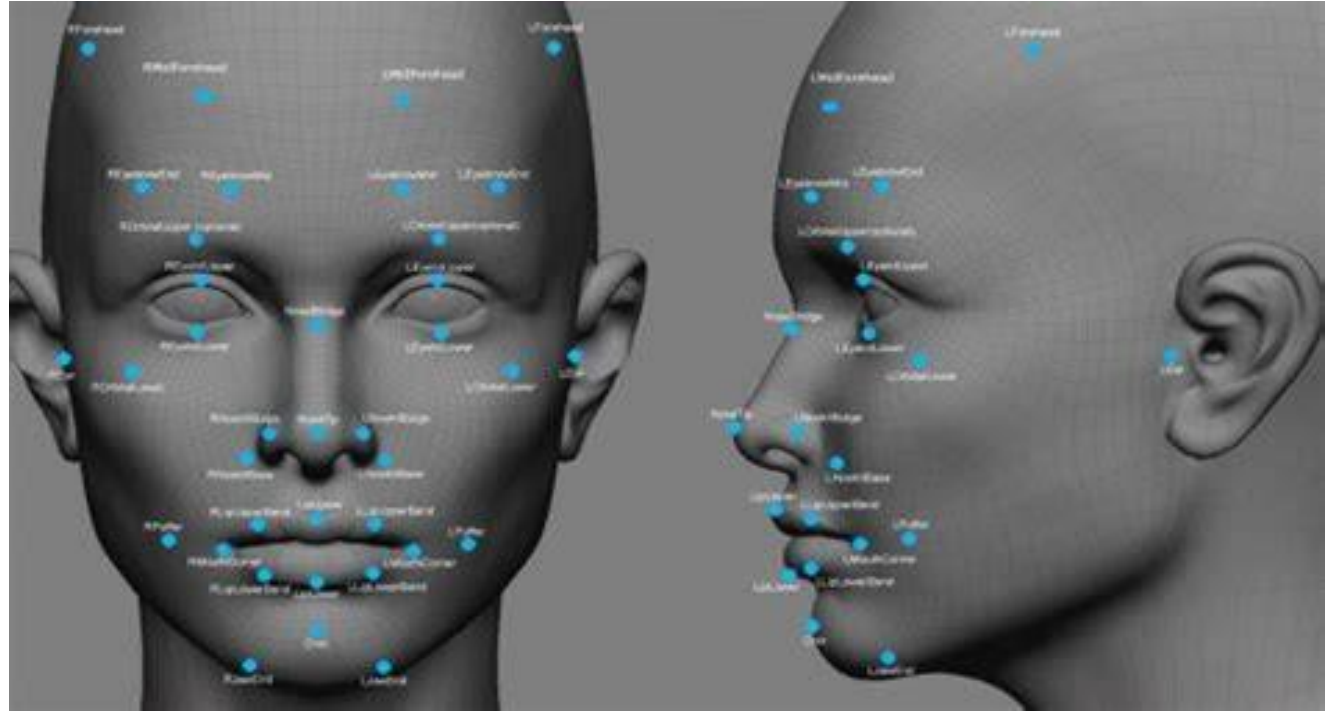




# Face Recognition

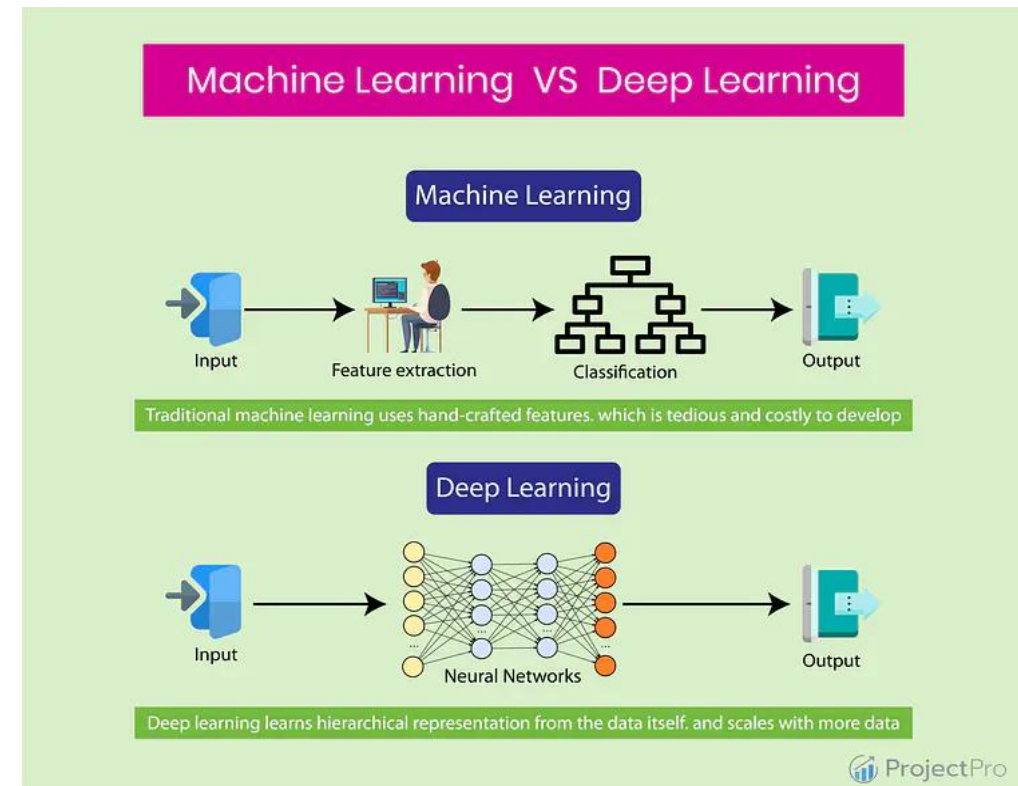
[is it {person/celebrity\_name}/?]

- Attaching clues or “landmarks”, to distinctive features
- **Biometric screening**
- **ID tracking** across video frames
- **Emotional intelligence** (microexpressions)



# Machine Learning vs Deep Learning

- ▶ Traditional machine learning uses hand-crafted features which is tedious and costly to develop
- ▶ Deep Learning learns hierarchical representation from the data itself and scales with more data
- ▶ Deep Learning is a subset of Machine Learning. They are complementary



Deep Learning approach #1

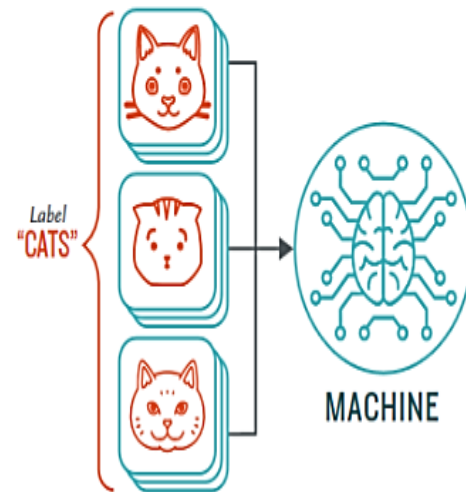
# Supervised Learning

Requires **labeled datasets** for training models (e.g., convolutional neural network for image classification)

## How Supervised Machine Learning Works

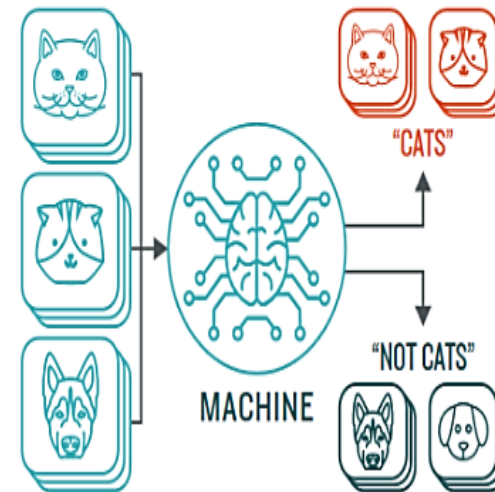
### STEP 1

Provide the machine learning algorithm categorized or "labeled" input and output data from to learn



### STEP 2

Feed the machine new, unlabeled information to see if it tags new data appropriately. If not, continue refining the algorithm

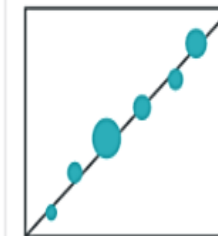


### TYPES OF PROBLEMS TO WHICH IT'S SUITED



#### CLASSIFICATION

Sorting items into categories



#### REGRESSION

Identifying real values (dollars, weight, etc.)

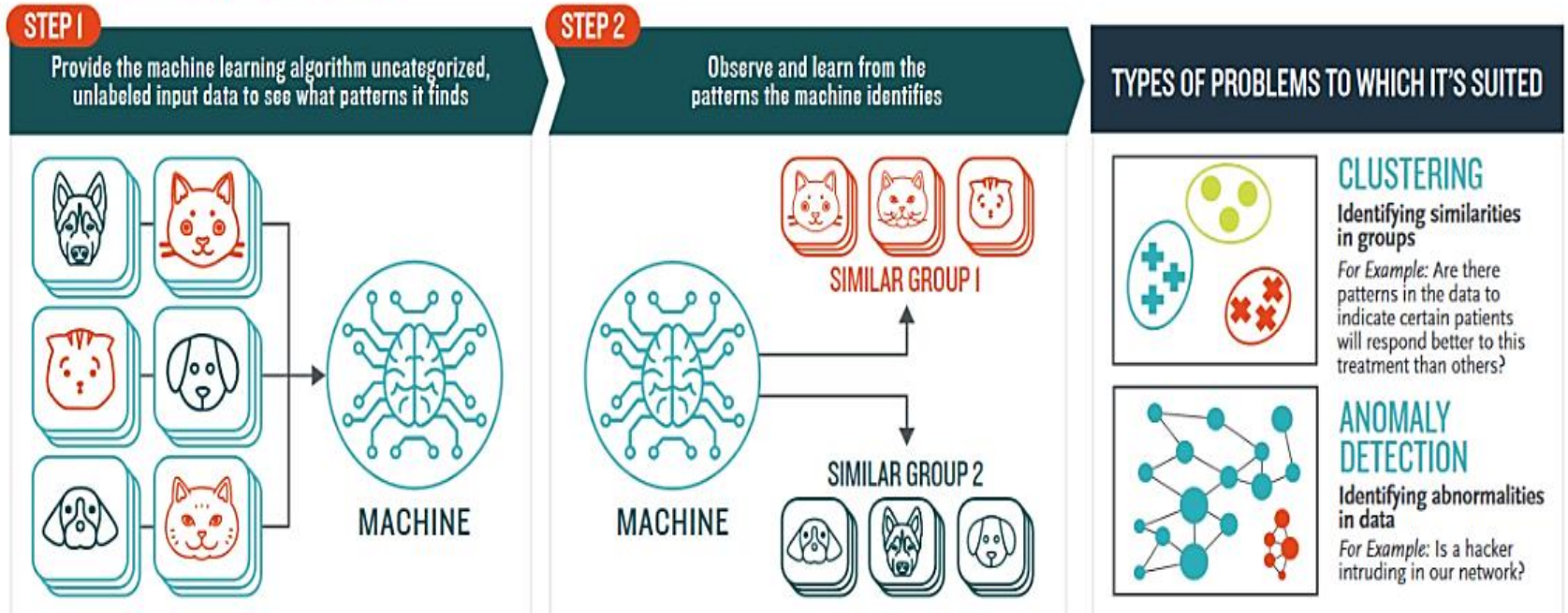


Deep Learning approach #2

# Unsupervised Learning

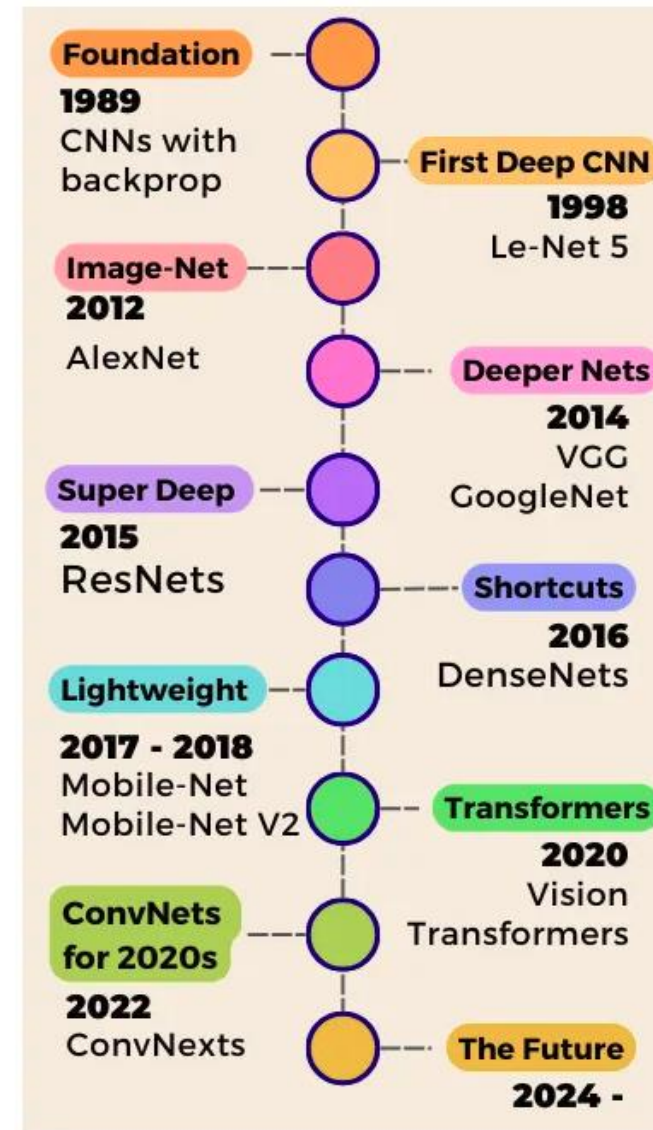
utilizes methods, such as autoencoders and generative adversarial networks (GANs) to learn from unlabeled data.

## How **Unsupervised** Machine Learning Works



# What are Convolutional Neural Networks?

CONVNETS OR CNNS

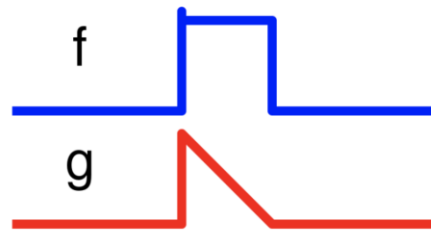


# In the context of mathematics

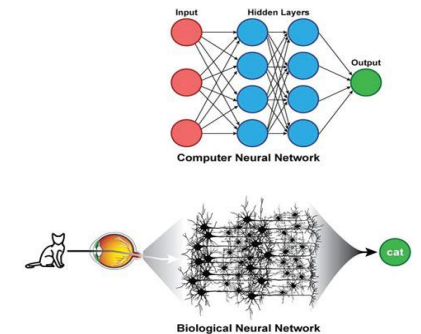
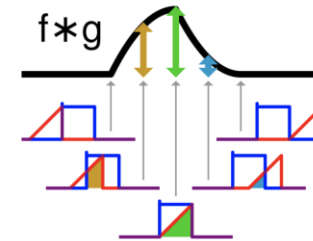
$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau$$

To convolve a kernel with an input signal:  
flip the signal, move to the desired time,  
and accumulate every interaction with the kernel

Example functions



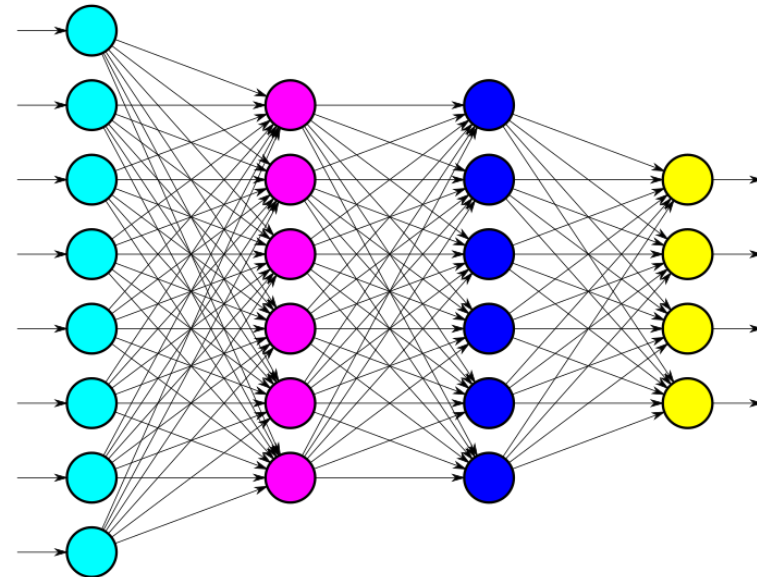
Convolution operation result





# In the context of AI, what is a Convolutional Neural Network?

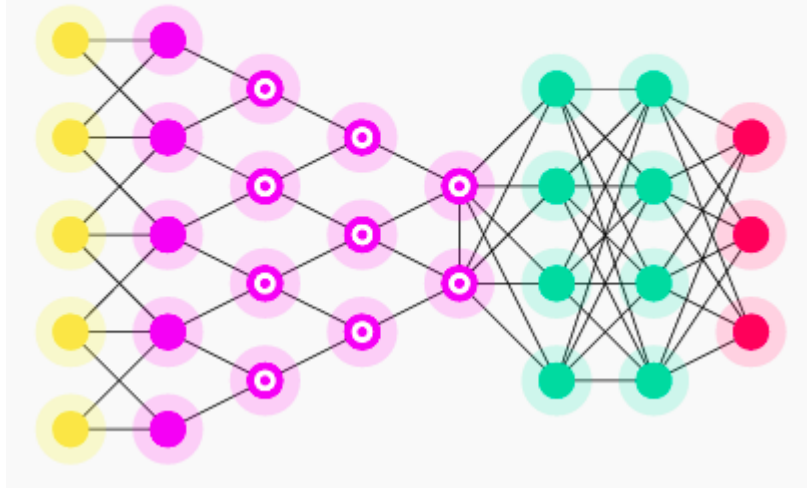
1. A type of **Deep Learning model** designed for **image processing**, also called ConvNet or CNN
2. **Mimics visual perception** in **humans**
3. A powerful tool for **visual recognition** and **style manipulation**
4. Consists of multiple **sequential layers** (convolutions, filters, and pooling) to **extract features** (edges, texture) from images ex: the cat has pointy ears, and a round face
5. Essential for modern applications in **Computer Vision** (multimodal multimedia: video, images)



# Architecture of a CNN

## → Deep CNN

Deep Convolutional  
Network **(DCN)**

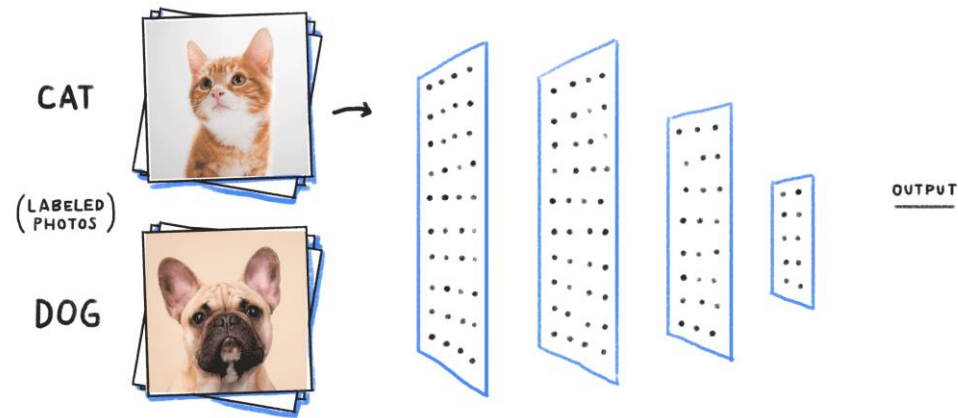


### Index

- Backfed Input Cell
- Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablisticc Hidden Cell
- Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Different Memory Cell
- Kernel
- Convolutional or Pool

# Image Classification with CNN

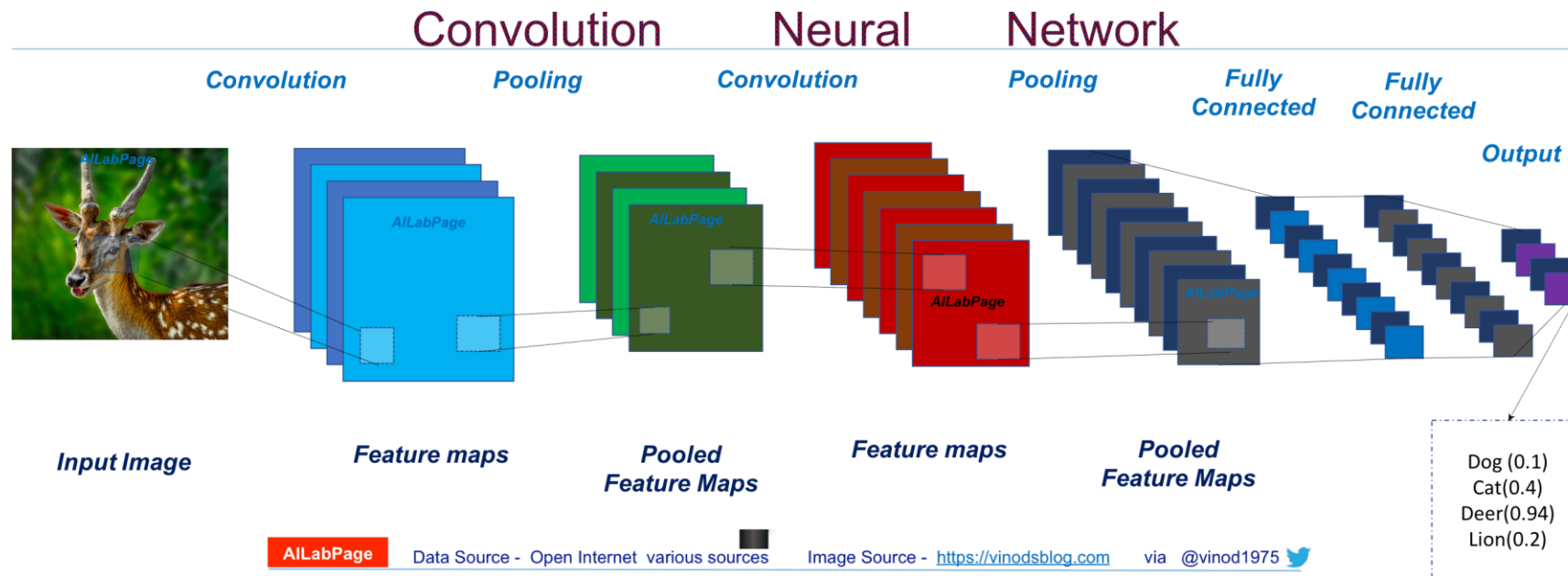
Is essentially a guess  
deducted from the  
extracted features



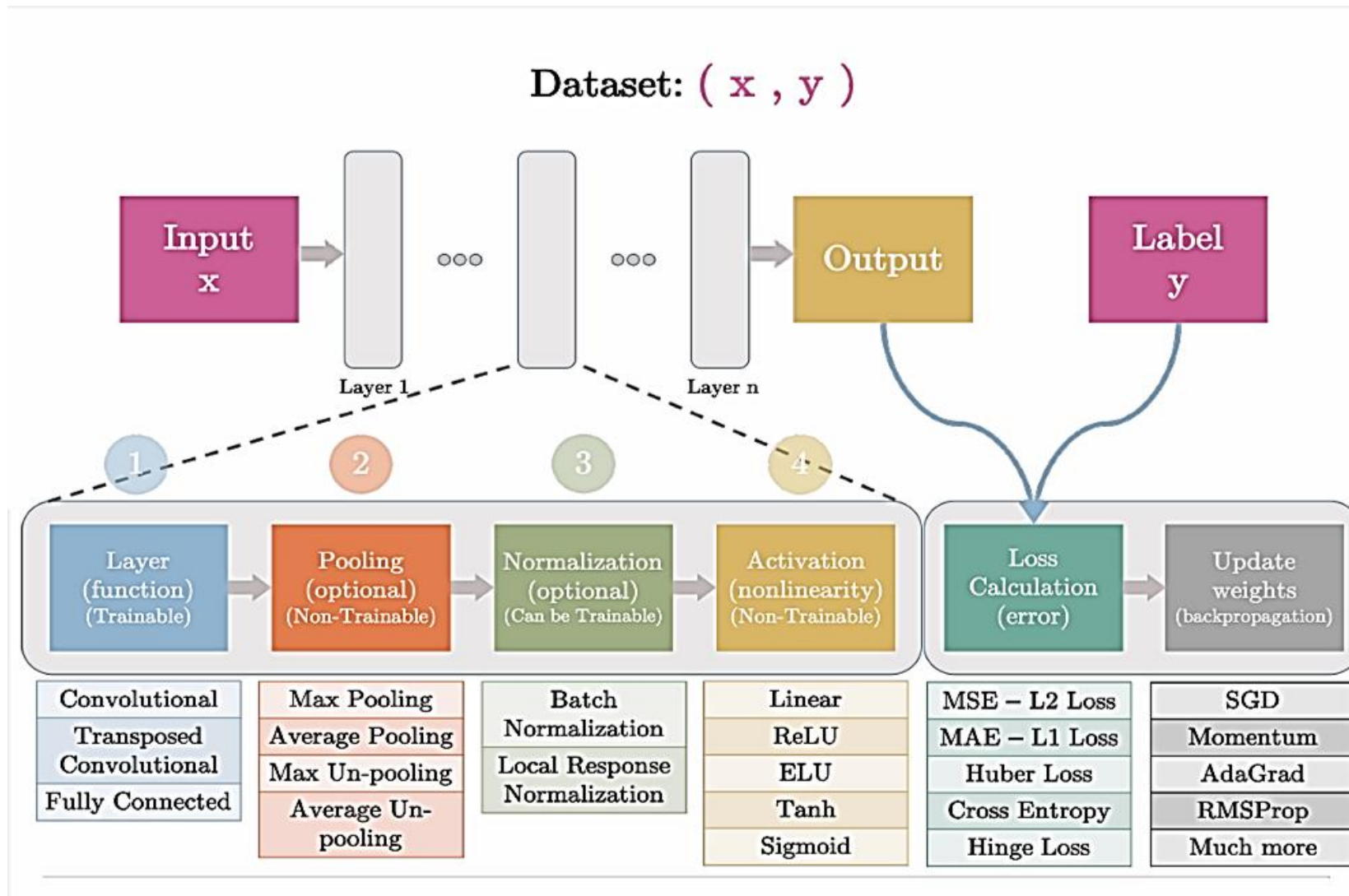




# CNN for image classification

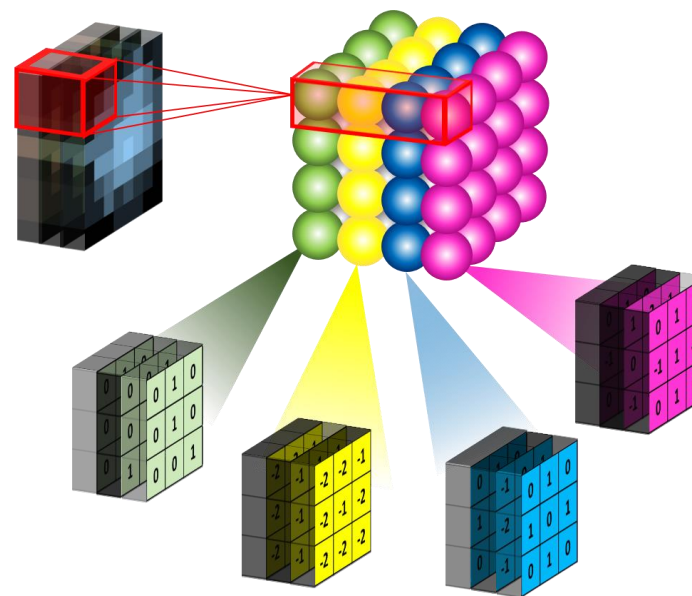
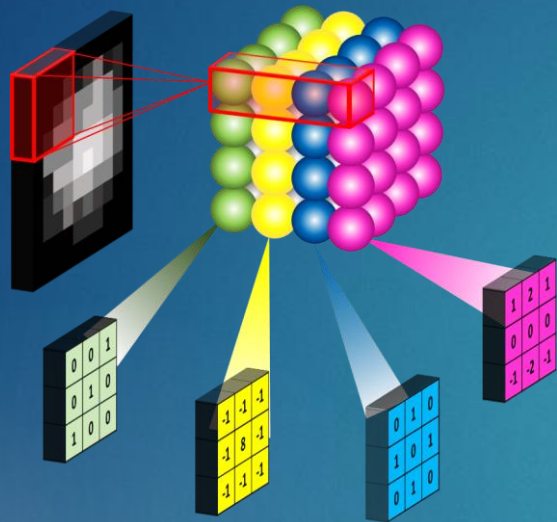


The Deer class gets 94% of accuracy



Trainable means  
can be fine-tuned:  
Hyperparameters  
weights and biases

# Convolutions

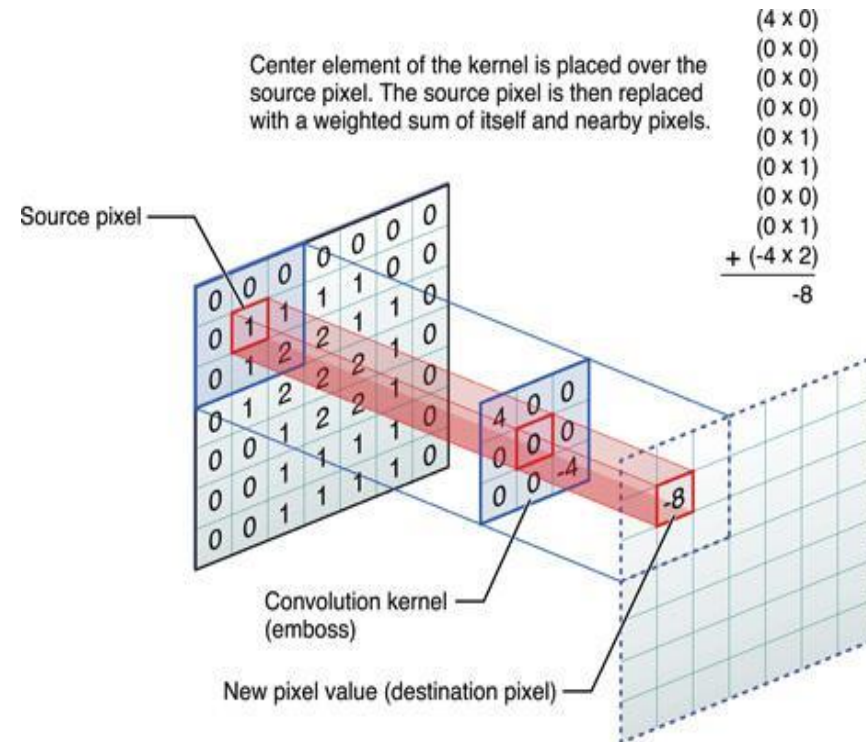


1D AND 3D



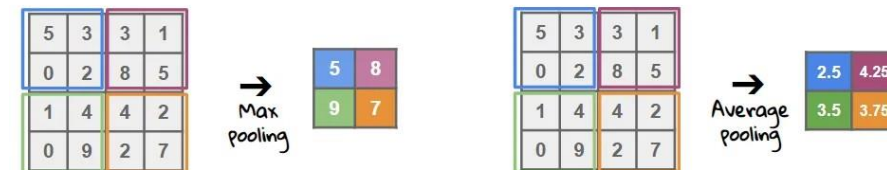
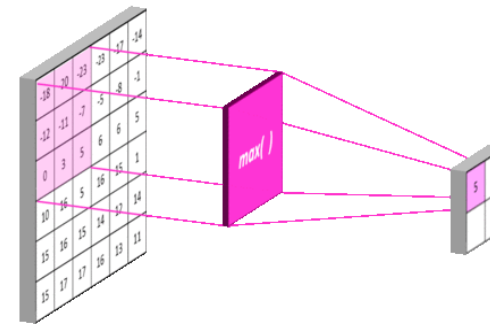
# What are Convolutions?

- ➔ - Definition: A mathematical operation that combines two functions.
- Purpose: To extract high-level features from images.
- Process: Apply a filter (kernel) over the image, computing dot products.
- ➔ -



# Pooling Layers?

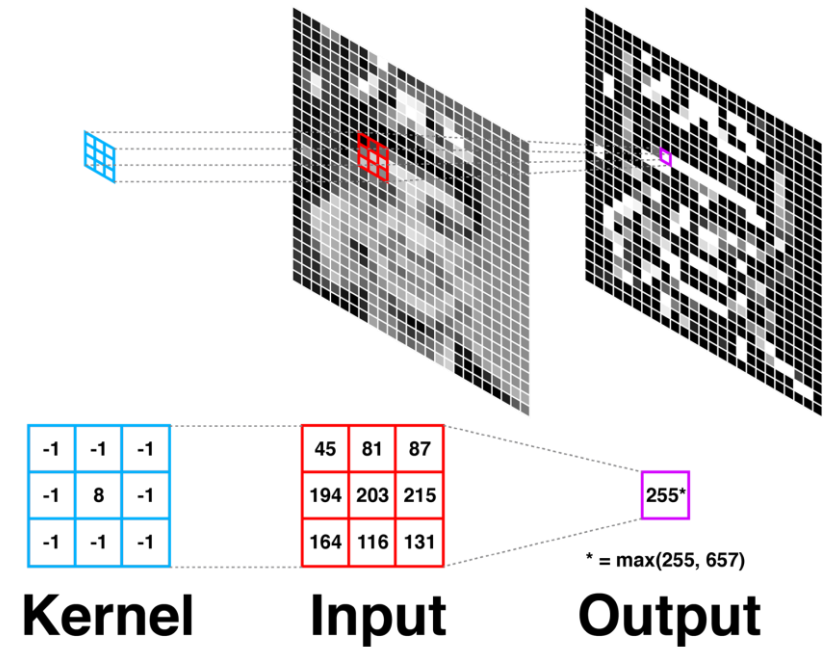
- ➔ **Function?** Reduces **spatial dimensions** of the feature maps.
- ➔ **Why?** Images can be voluminous (resolutions, size)
- ➔ **Solution:** Focus on the essential features
- ➔ **How?** Sliding a window and taking only one value
  - Types:
    - **Max Pooling:** Takes the maximum value from a set of values.
    - **Average Pooling:** Takes the average value.
    - **Benefits:** Decreases computational load and helps prevent overfitting.



## Filters / Kernels

Small matrices used in convolutions.

- Detect specific features (Edges, textures).
- Different filters highlight different aspects of the image.

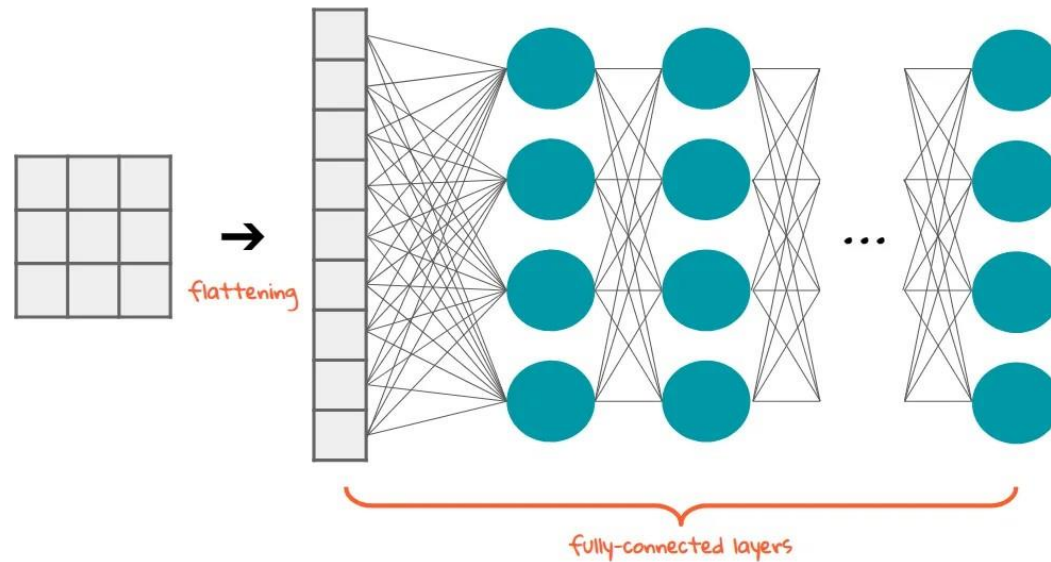


# Flattening Layer

Converting multi-dimensional feature maps from the convolutional and pooling layers into a one-dimensional array.

**Why?** To connect extracted features to the final classification task. Fully-connected layers require linear input format.

**How?** Takes the output from the last pooling layer (dimensions: Height, Width, Depth) and reshapes into a single vector with a size Height x Width x Depth





# Tools for Computer Vision using CNNs

## Frameworks

- **TensorFlow:** An open-source library for deep learning developed by Google, ideal for building CNNs.
- **PyTorch:** A flexible deep learning framework from Facebook, popular for research and production.
- **Keras:** A high-level API that simplifies building CNNs, often used with TensorFlow.
- **MATLAB:** Provides tools for designing and deploying CNNs through its Deep Learning Toolbox

## Object Detection models

**YOLO** (You Only Look Once)

**R-CNN** (Region with CNN)

can be implemented using those frameworks



## Techniques

- **ReLU** activation function introduces non-linearity
- **Dropout** prevents overfitting
- **Data Augmentation** provides more data (perspectives, side views etc)
- **Transfer Learning** reduces time of training and resources

# Notion of Padding and Stride

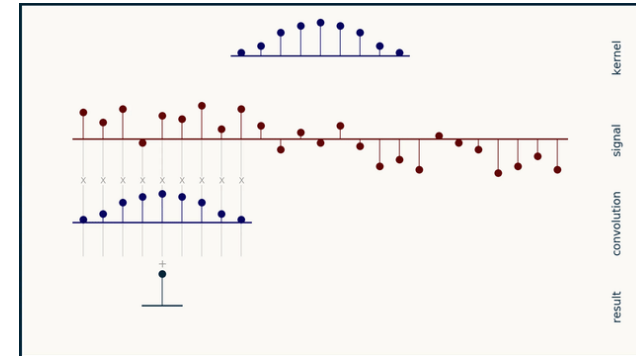
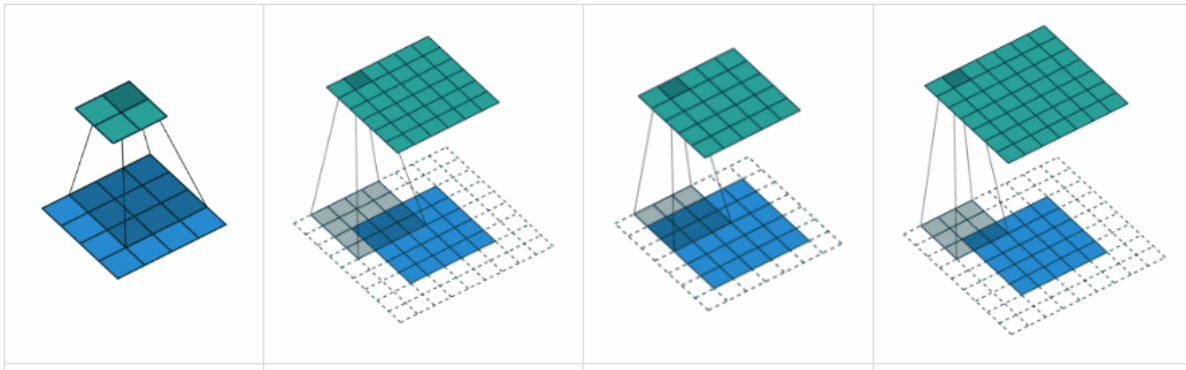


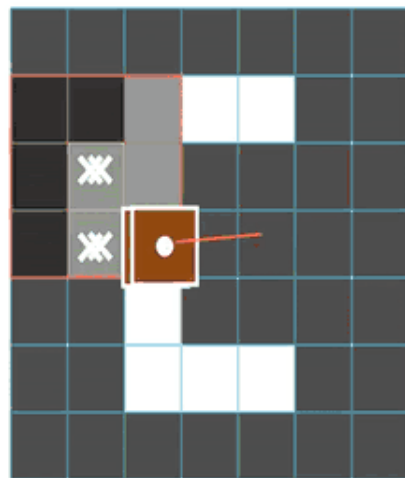
Image: Jiwon Jeong, e2eml.schools

**Padding:** adding extra pixels around input data before convolution ('valid' when none, 'same' when output match input)

**Stride:** number of pixels a convolution filter moves/ steps at a time

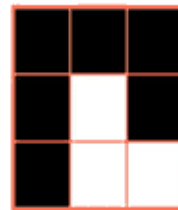
Crucial for effective feature extraction (relevant features) and model performance

# How to get a feature map?



7x7 Image (padded)

\*

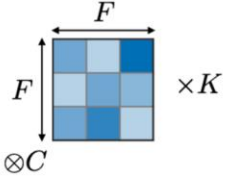
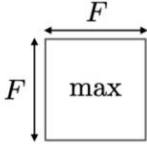
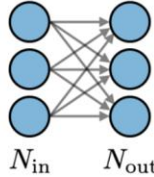


3x3 kernel

1	2	1	1	0

5x5 Feature Map

# Notion of complexity in training

	CONV	POOL	FC
Illustration			
Input size	$I \times I \times C$	$I \times I \times C$	$N_{\text{in}}$
Output size	$O \times O \times K$	$O \times O \times C$	$N_{\text{out}}$
Number of parameters	$(F \times F \times C + 1) \cdot K$	0	$(N_{\text{in}} + 1) \times N_{\text{out}}$
Remarks	<ul style="list-style-type: none"><li>• One bias parameter per filter</li><li>• In most cases, <math>S &lt; F</math></li><li>• A common choice for <math>K</math> is <math>2C</math></li></ul>	<ul style="list-style-type: none"><li>• Pooling operation done channel-wise</li><li>• In most cases, <math>S = F</math></li></ul>	<ul style="list-style-type: none"><li>• Input is flattened</li><li>• One bias parameter per neuron</li><li>• The number of FC neurons is free of structural constraints</li></ul>

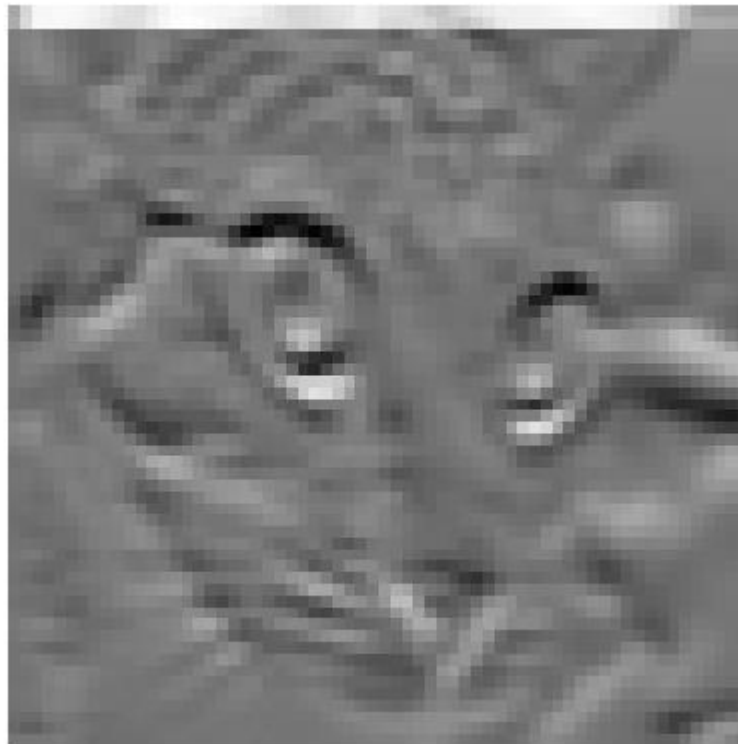


# Filters

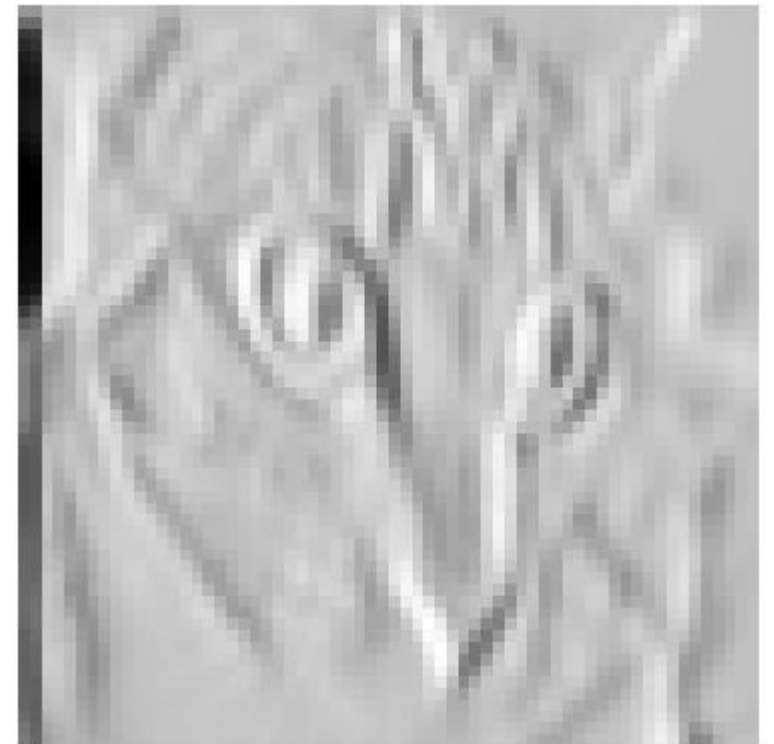
Original  
input image



➤ Horizontal edge filter

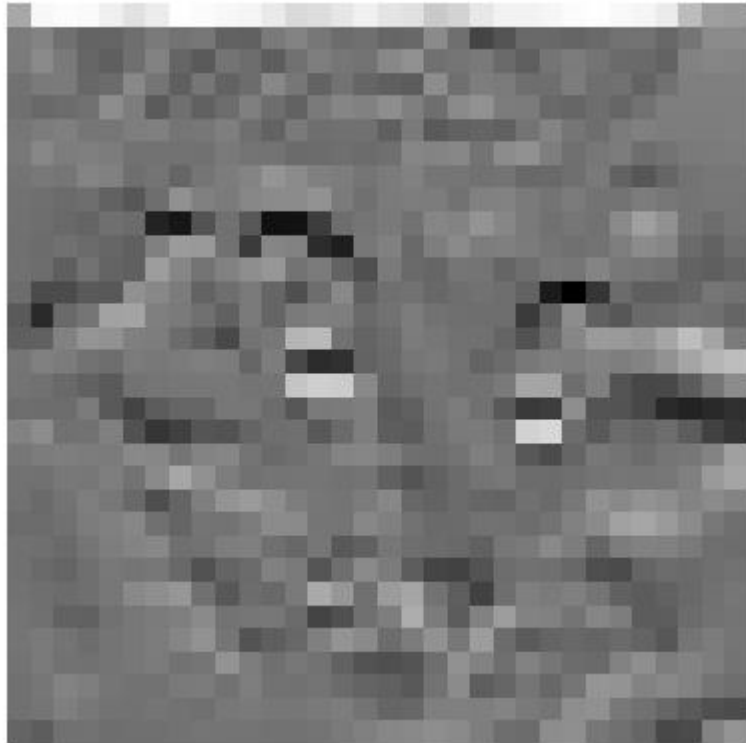


➤ Vertical edge filter



# STRIDES

Horizontal edge filter with stride2



Vertical edge filter with stride2

