



Deep Learning

Alberto Ezpondaburu



Analytical Index

Convolutional Neural Networks (CNN)

- a. Introduction to Computer Vision
- b. Introduction to CNN



2

Convolutional Neural Networks (CNN)



2.1

Introduction to Computer Vision



Computer Vision I



- Massive growth in visual data:
 - 500 hours of video are uploaded to YouTube every minute.
 - TikTok has over 1 billion monthly active users that spend 1.5 hours on average per day.
 - 100 million post and 1 billion stories are uploaded to Instagram every day.
 - There are cameras everywhere.
- Techniques and ideas can be used in other types of ordered data.
- People create and share visual data continuously.



Computer Vision II

“Building artificial systems that process, perceive and reason about visual data”

Visual Perception: 500 millions of years of evolution, involves both the evaluation of stimulus features and attention (differential processing of selected information that is relevant to behavior)

Importance of computer vision:

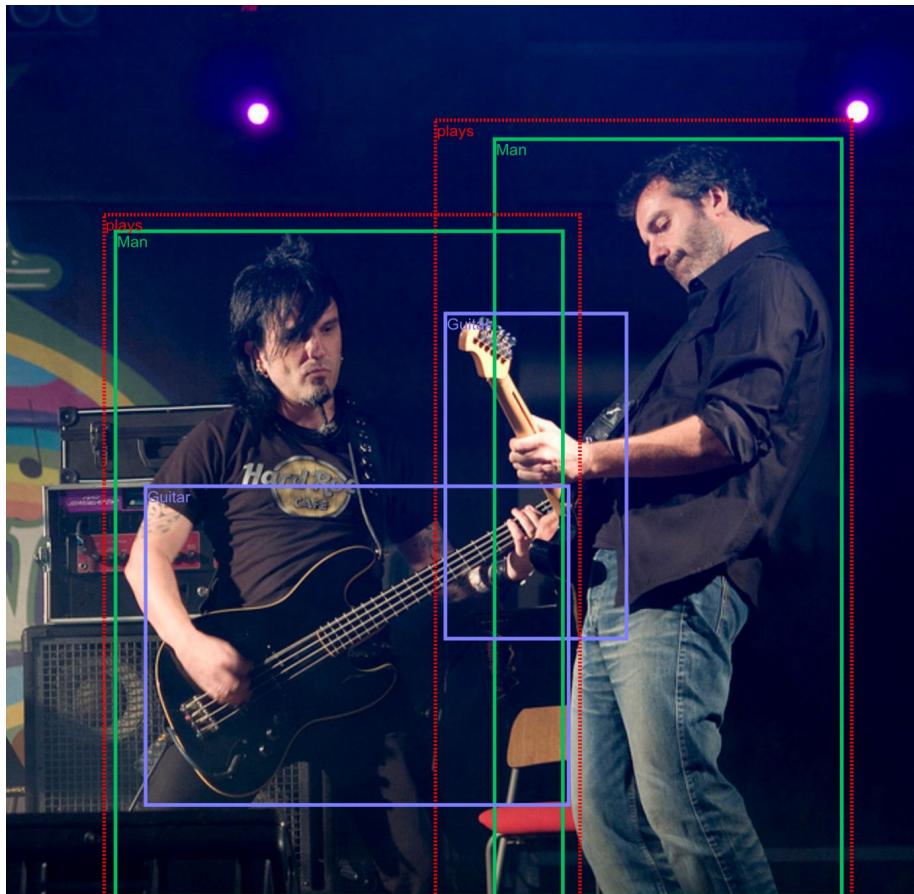
1. Automated analysis and understanding of visual data to enhance human capabilities.
2. Enables applications in diverse industries, including healthcare, automotive, retail, security, digital marketing...



Computer Vision: Image Classification & Location



Computer Vision: Relationship Detection



Computer Vision: Pose Estimation

Determines the position and orientation of objects or people, useful in AR, gaming, and healthcare



Computer Vision: Object Tracking

Identifies and locates objects within images or video streams, used in retail, security, and many other sectors



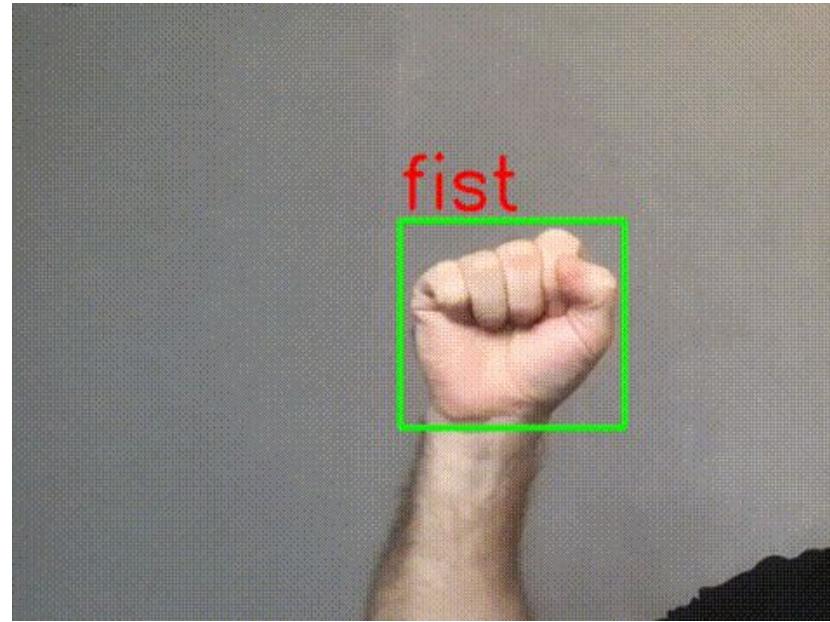
Computer Vision: Self-Driving Cars

Identifies and locates objects within images or video streams, used in retail, security, and many other sectors

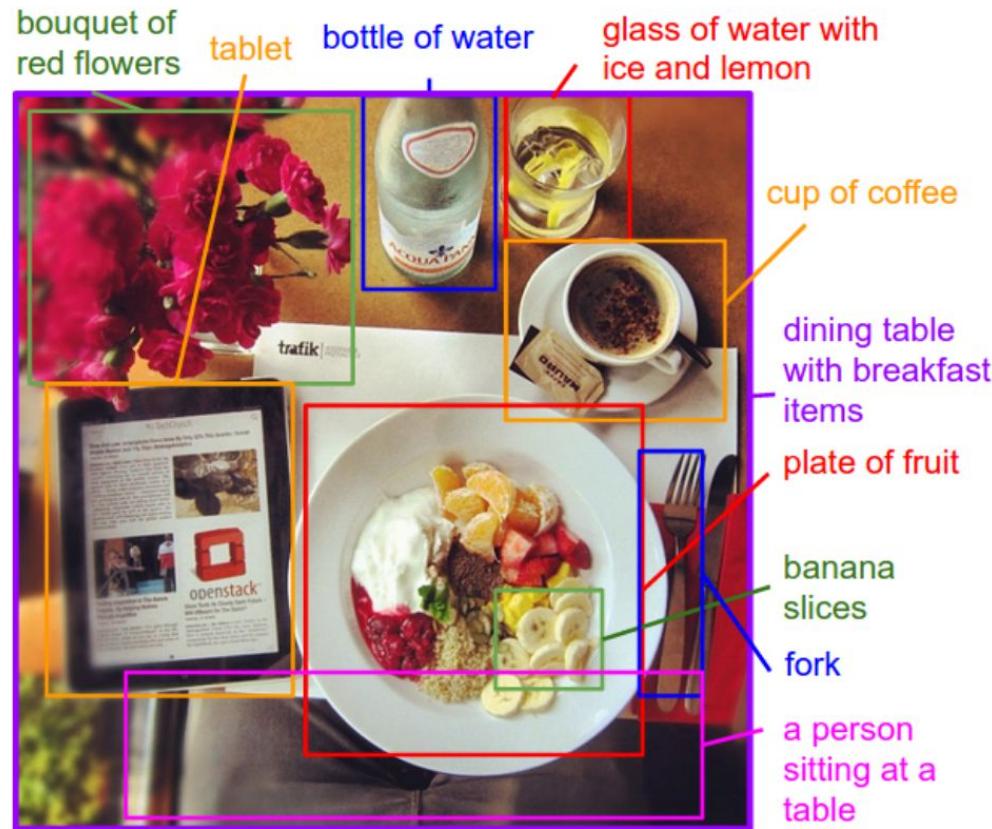


Computer Vision: Gesture Recognition

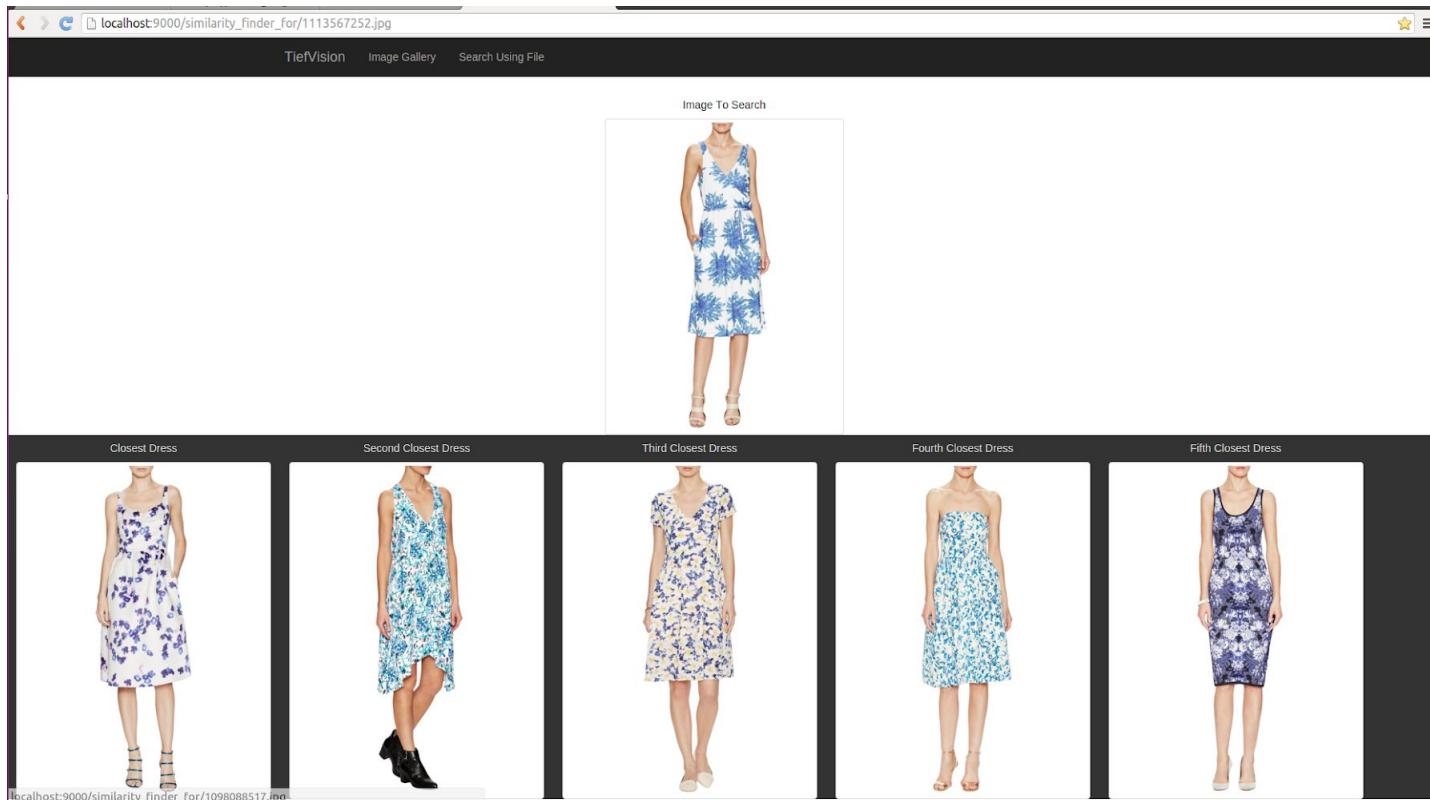
Interpreting human gestures to control devices or applications, used in gaming, healthcare, and smart home systems.



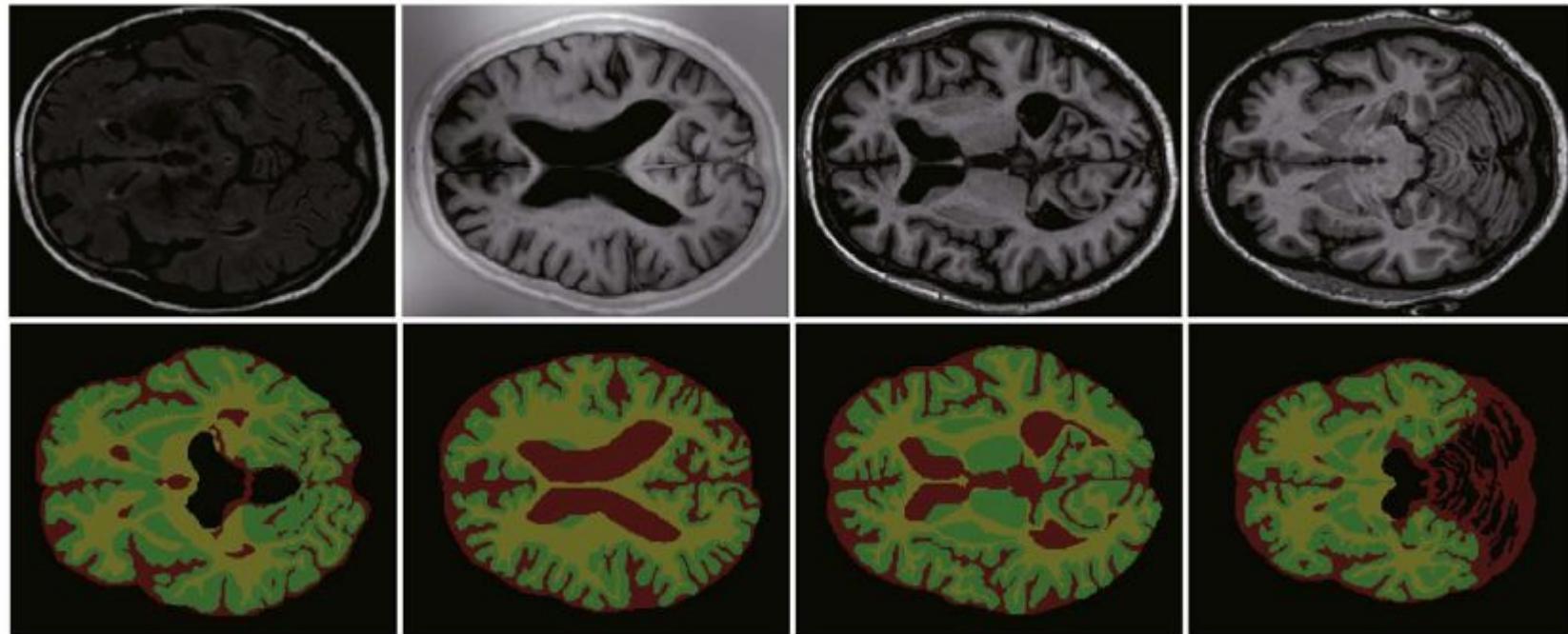
Computer Vision: Image annotation



Computer Vision: Image Similarity Search Engine



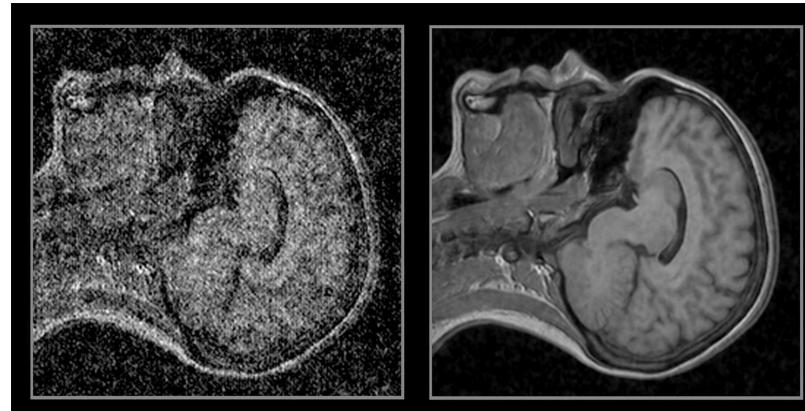
Computer Vision: Image segmentation



Computer Vision: Image Reconstruction



Computer Vision: Image Denoising



Computer Vision: Image search

```
image_search('Find me an image of a famous monument in India',  
Input query: Find me an image of a famous monument in India
```

0.3326866328716278

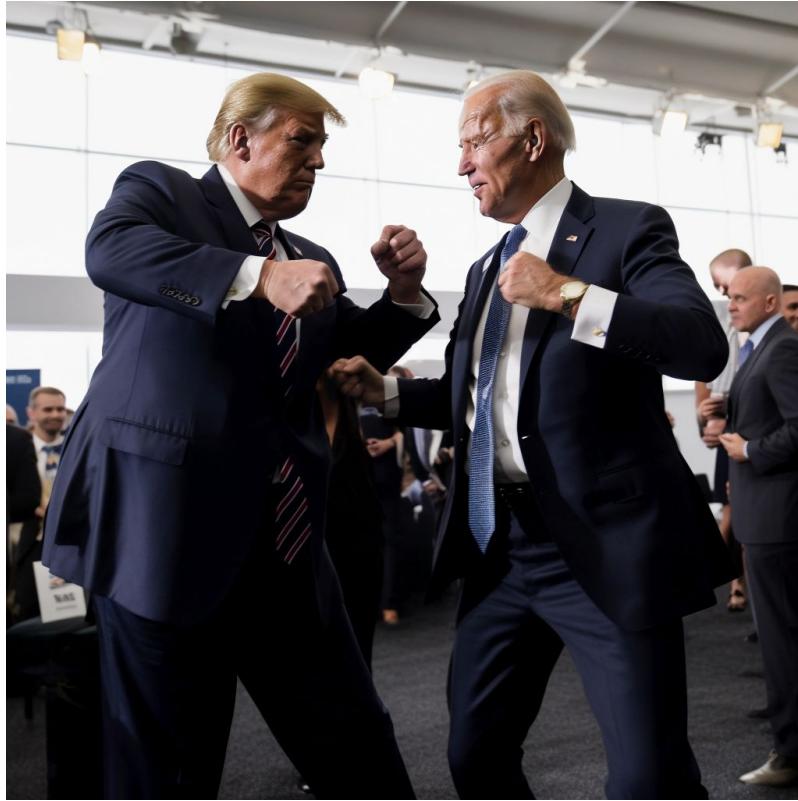


```
image_search_from_path('Two dogs playing in the snow', model, img_embeddings, img_  
Input query: Two dogs playing in the snow
```

0.31620916724205017

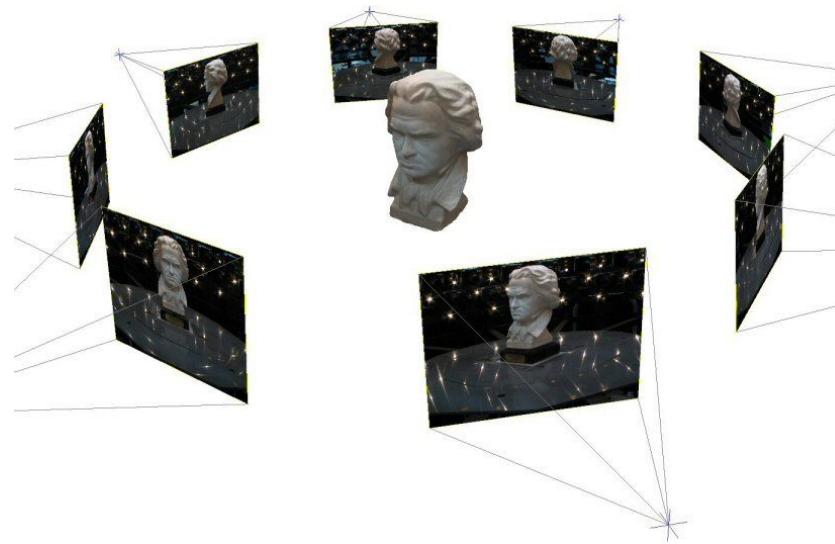


Computer Vision: Image Generation



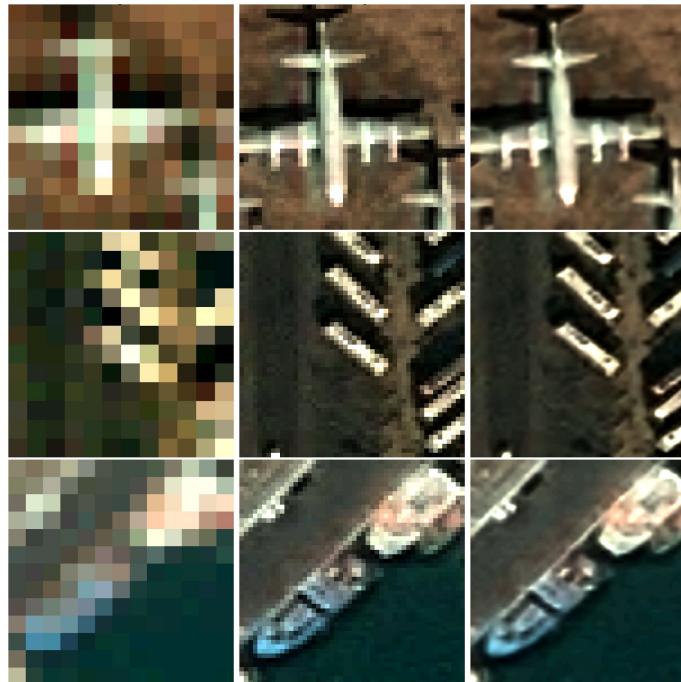
Computer Vision: 3D Reconstruction

Constructs 3D models from 2D images, used in virtual reality, architecture, and entertainment



Computer Vision: Super Resolution

Enhances the resolution of images, making them clearer, which is widely used in satellite imaging, medical imaging, and surveillance



2.2

Introduction CNN:



Image Data



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49	
13	217	243	255	155	35	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62	
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	21	252	255	248	144	6	0	
0	13	112	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0	
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4	
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	35	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3	
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	249	255	240	255	129	0	5	0	
0	0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	
0	0	6	1	0	52	153	233	255	252	147	37	0	4	1	0	
0	0	5	5	0	0	0	0	0	14	1	0	6	6	0	0	

- Coloured images are usually represented as mixes of three colours: Red, Green and Blue (RGB)

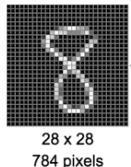


- Images are composed by pixels. 1 Megapixel can be a 1000x1000 matrix.
- Grayscale images can be seen as matrices of integers (0 -255) (black - white).

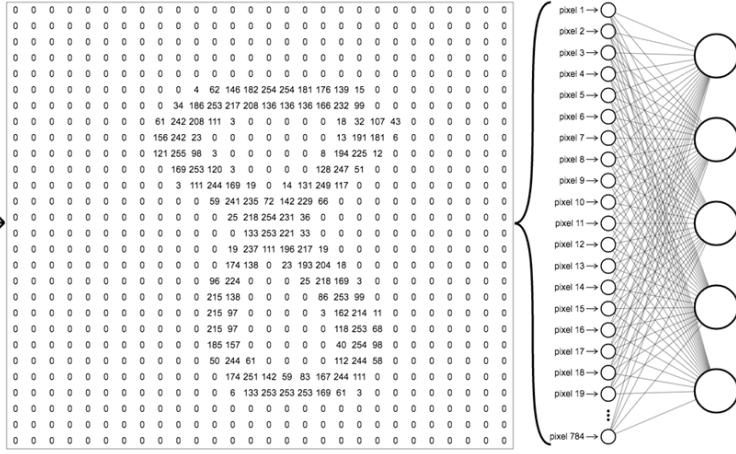
210	214	216
208	210	211
204	2	167
186	188	188
183	186	194
235	238	239
230	206	232



Problem with Dense Layers



28 x 28
784 pixels

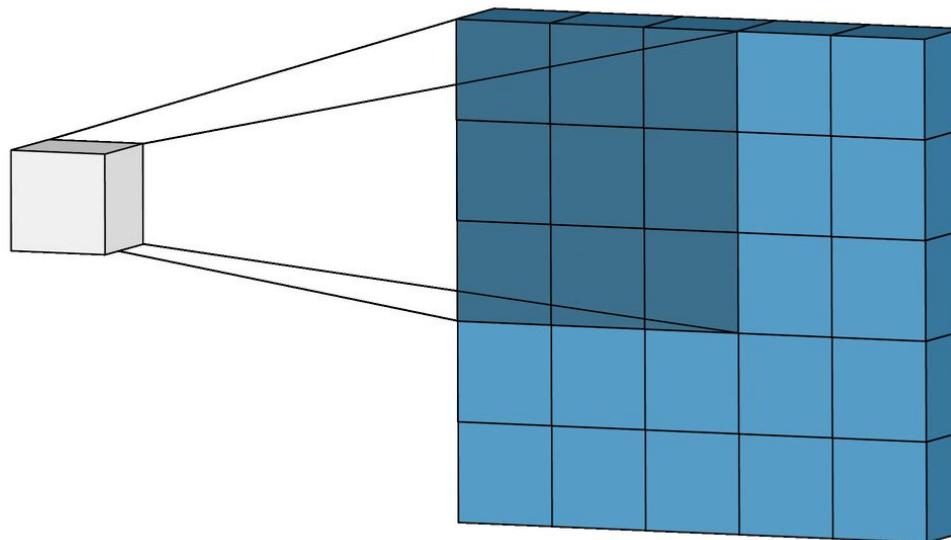


- Spotify: 44.1kHz sample rate.
- Same problems.

- 1 Megapixel = $1000 \times 1000 = 10^6$ inputs.
- 1000 units hidden layer = 10^9 parameters!!
- No spatial information. Collapse 2D => 1D
- In most tasks, you can share the same operations in different parts of the image, In a fully-connected layer you learn pixel-by-pixel operations.



Convolution Operation: Intuition



- Connect input patch to a single neuron to capture spatial information
- Use sliding window to order connections.
- 3x3 filter, stride = 1 (shift by 1 pixel the sliding window)

Convolution Operation: Filtering

$$\mathbf{X} \quad \begin{array}{|c|c|c|c|c|c|c|}\hline 3_1 & 0_0 & 1_{-1} & 2 & 7 & 4 \\ \hline 1_1 & 5_0 & 8_{-1} & 9 & 3 & 1 \\ \hline 2_1 & 7_0 & 2_{-1} & 5 & 1 & 3 \\ \hline 0 & 1 & 3 & 1 & 7 & 8 \\ \hline 4 & 2 & 1 & 6 & 2 & 8 \\ \hline 2 & 4 & 5 & 2 & 3 & 9 \\ \hline \end{array} \quad 6 \times 6$$

$*$ \mathbf{W} $=$ $\begin{array}{|c|c|c|c|}\hline -5 & & & \\ \hline \end{array} \quad 4 \times 4$

$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$

$\mathbf{Y} = \mathbf{X} \star \mathbf{W}$

$$y_{p,q} = \sum_{i=0}^2 \sum_{j=0}^2 w_{ij} x_{i+p,j+q}$$

- In mathematics it is called Cross-correlation, In Deep learning convolution.
- Operator: \star
- Dimensions: $(n \times n) \star (f \times f) \equiv (n - f + 1) \times (n - f + 1)$



Convolution Operation: Example

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

$$\frac{1}{9} *$$

1	1	1
1	1	1
1	1	1

“box filter”

0	10	20	30	30	30	20	10	
0	20	40	60	60	60	40	20	
0	30	60	90	90	90	60	30	
0	30	50	80	80	90	60	30	
0	30	50	80	80	90	60	30	
0	20	30	50	50	60	40	20	
10	20	30	30	30	30	20	10	
10	10	10	0	0	0	0	0	



Convolution Operation: Example

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	10	0	0
10	10	10	10	0	0
10	10	10	10	0	0

6 x 6

*

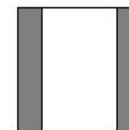
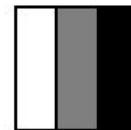
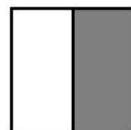
1	0	-1
1	0	-1
1	0	-1

3 x 3

=

-0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

4 x 4



Convolution Operation: Example

255	255	255	255	255	255
255	255	255	255	255	255
255	255	255	255	255	255
255	255	255	255	255	255
10	10	10	10	10	10
10	10	10	10	10	10
10	10	10	10	10	10
10	10	10	10	10	10

8 X 6

*

6 X 4

=

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

3 X 3

0	0	0	0
0	0	0	0
735	735	735	735
735	735	735	735
0	0	0	0
0	0	0	0



Convolution Operation: Example

$$\begin{array}{|c|c|c|c|c|} \hline 35 & 40 & 41 & 45 & 50 \\ \hline 40 & 40 & 42 & 46 & 52 \\ \hline 42 & 46 & 50 & 55 & 55 \\ \hline 48 & 52 & 56 & 58 & 60 \\ \hline 56 & 60 & 65 & 70 & 75 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline 0 & 42 & 0 \\ \hline \end{array}$$

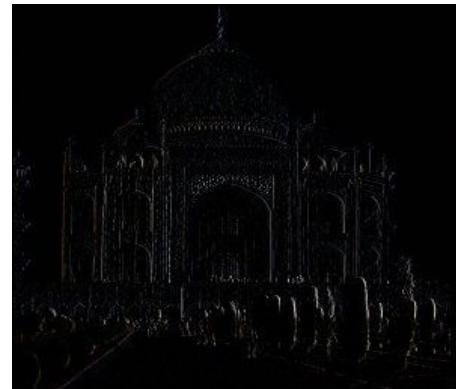
Edge Detection

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline 1 & -4 & 1 \\ \hline 0 & 1 & 0 \\ \hline \end{array}$$



Edge Detection
(right)

$$\begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline -1 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array}$$



Convolution Operation: Example



Blur

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

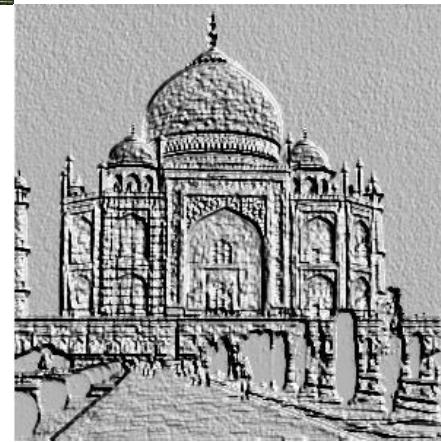


Sharpen

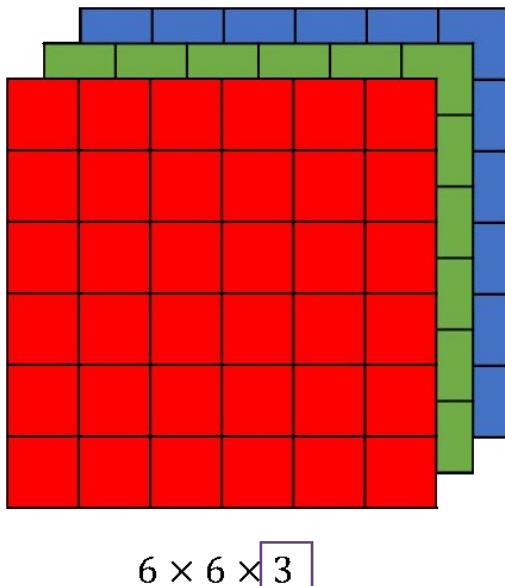
0	0	0	0	0	0
0	0	-1	0	0	0
0	-1	5	-1	0	0
0	0	-1	0	0	0
0	0	0	0	0	0

Emboss

-2	-1	0	
-1	1	1	
0	1	2	

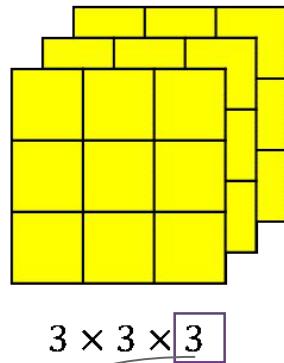


Convolution Operation: RGB

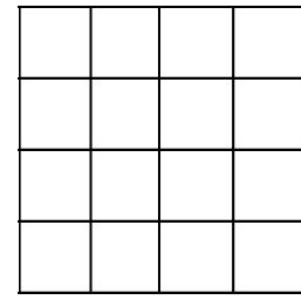


*

3 separate convolutions and then sum all 3

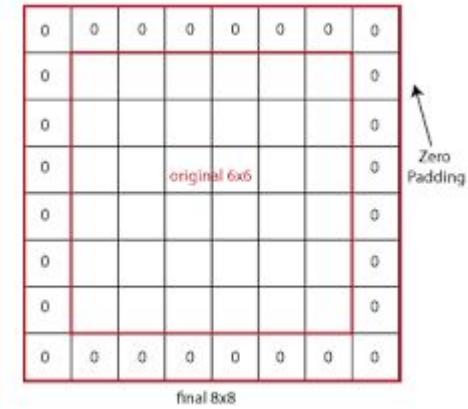
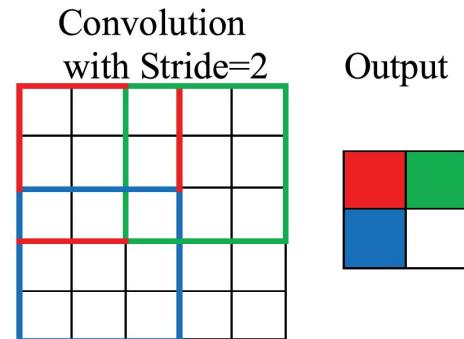
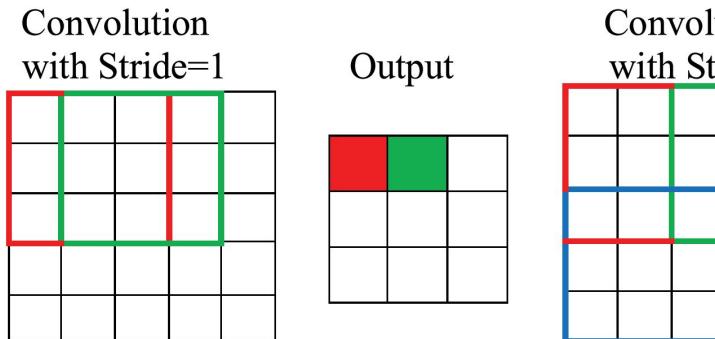


=



Convolution: Padding & Stride

- The image shrinks every time a convolution operation is performed
- Padding (p): amount of zero pixels added to an image: $(n - f + 1) \times (n - f + 1)$
- Types of convolution:
 - "valid": No padding ($p = 0$)
 - "same": Output size = Input size
$$p = \frac{f - 1}{2}$$



- Stride: (s) number of rows and columns traversed per slide

$$\left\lfloor \frac{n - f + 2p}{s} + 1 \right\rfloor$$

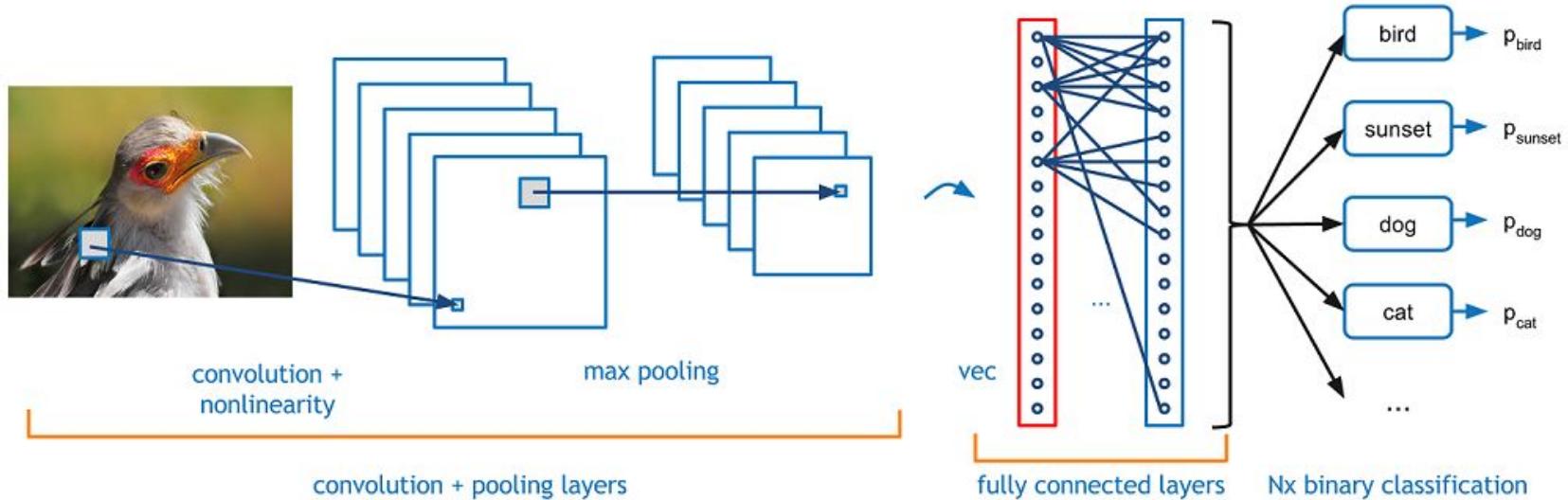


CNNs Advantages

- **Sparse Connections:** we need to store fewer parameters, which both reduces the memory requirements of the model and improves its statistical efficiency.
- **Shareable Weights:** use the same parameter for more than one function in a model. Rather than learning a separate set of parameters for every location, we learn only one set. **Equivariance to translation**, if the input changes, the output changes in the same way.



CNN Layers



- Convolution + non-linearity (ReLU)
- Pooling: Downsampling
- Fully-connected

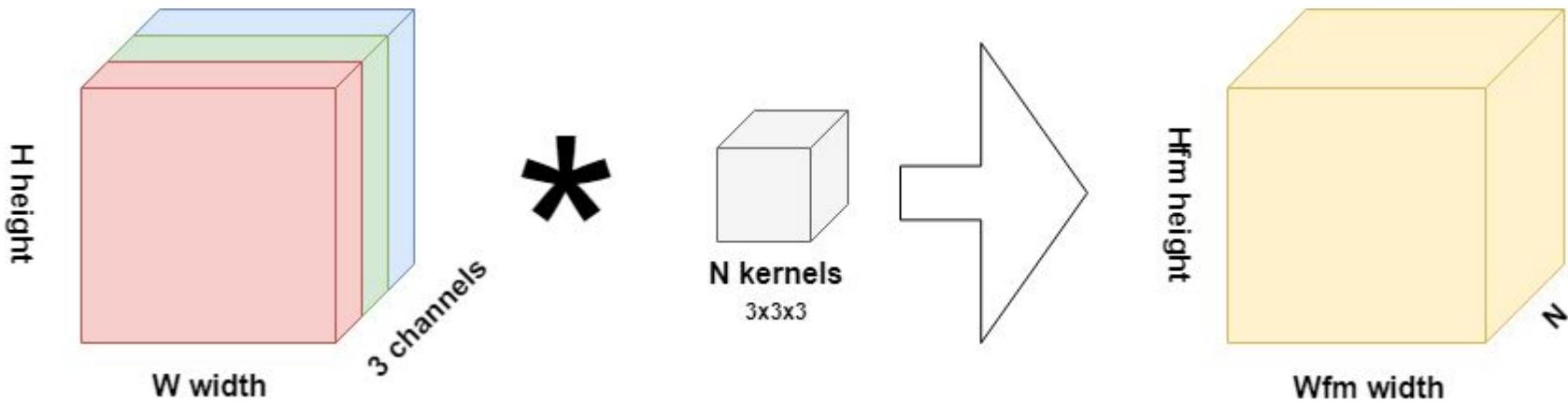
```
tf.keras.layers.Conv2D
```

```
tf.keras.layers.MaxPool2D
```

```
tf.keras.layers.Dense
```



CNN Layers: Convolution



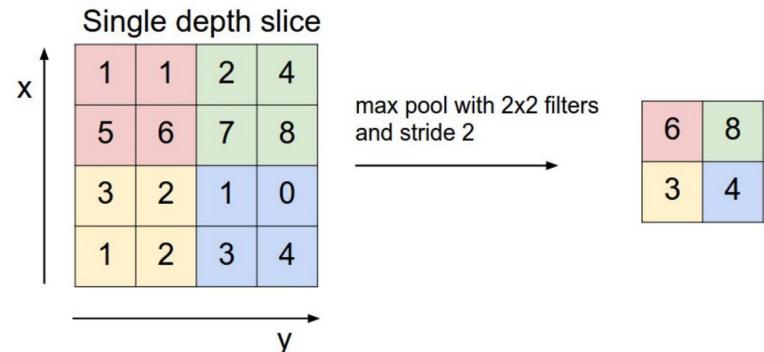
Layer dimensions: $(h \times w \times N)$ $\text{Relu}(\text{convolution} + b)$ for every filter (N different filters)

```
tf.keras.layers.Conv2D(  
    filters=N,  
    kernel_size=(h,w),  
    strides=(1, 1),  
    padding="valid"  
)
```



CNN Layers: Pooling

- **Dimensionality reduction:** representations turn smaller and more manageable.
- Reduces the number of parameters contributing to a smaller complexity.
- **Spatial invariance.**
- MaxPooling and average pooling are the most common.
- Backprop is not affected. **No parameters to tune.**

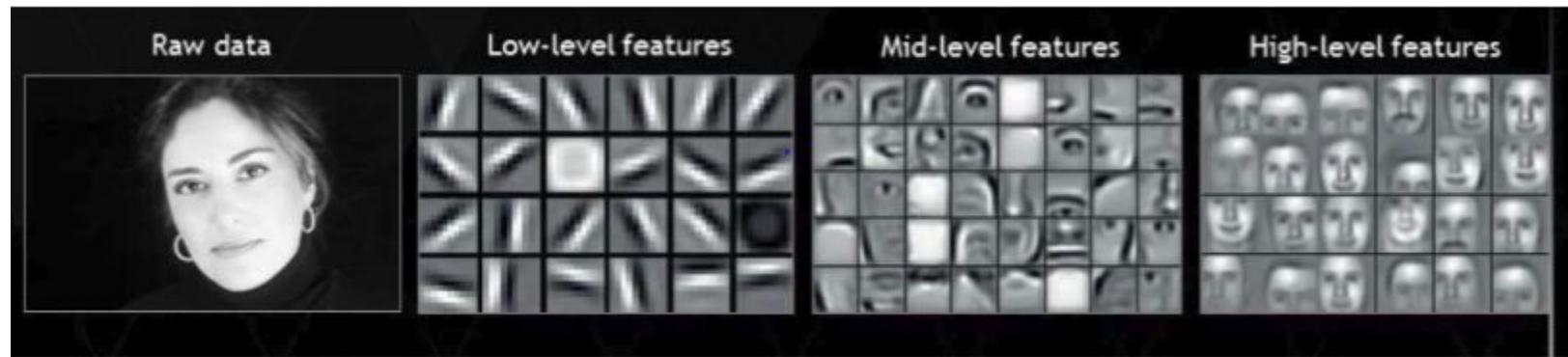


```
tf.keras.layers.MaxPooling2D(  
    pool_size=(2, 2), strides=(1, 1)  
)
```

```
tf.keras.layers.AveragePooling2D(  
    pool_size=(2, 2), strides=(1, 1)  
)
```



CNN Features



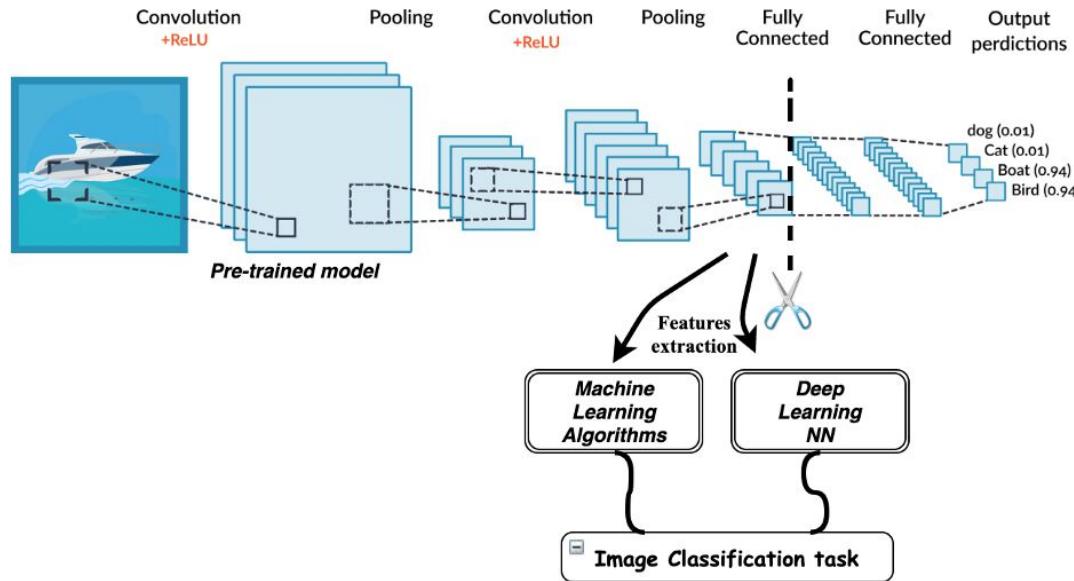
CNN: Data Augmentation



- Flip
- Rotation
- Scale
- Translation
- Crop
- Gaussian Noise
- ...



CNN: Transfer learning

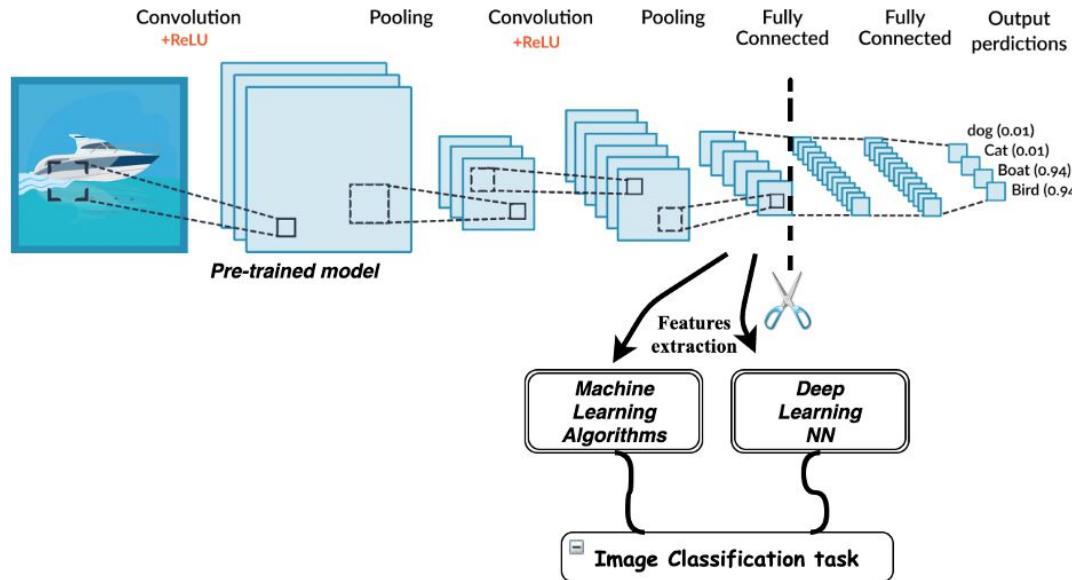


Instead of training a Deep Neural Net from scratch for your task:

- Take a network trained on a different domain for a different task
- Adapt it for your domain and your task



CNN: Fine-tuning



We not only update the CNN architecture but also re-train it to learn new object classes.





UNIVERSIDAD
COMPLUTENSE
DE MADRID

