

Image Classification

Giacomo Boracchi,
DEIB, Politecnico di Milano
October, 21th, 2020

giacomo.boracchi@polimi.it
home.deib.polimi.it/boracchi/

Who I am

Giacomo Boracchi (giacomo.boracchi@polimi.it)

Mathematician (Università Statale degli Studi di Milano 2004),
PhD in Information Technology (DEIB, Politecnico di Milano 2008)
Associate Professor since 2019 at DEIB, Polimi (Computer Science)



My research interests are mathematical and statistical methods for:

- Image analysis and processing
 - Machine Learning and in particular unsupervised learning, change and anomaly detection
- ... and the two combined

Materials

<https://boracchi.faculty.polimi.it/teaching/AN2DL.htm>

What is image classification
(and computer vision) about

Computer Vision

An interdisciplinary scientific field that deals with how
computers can be made to gain high-level
understanding from digital images or videos

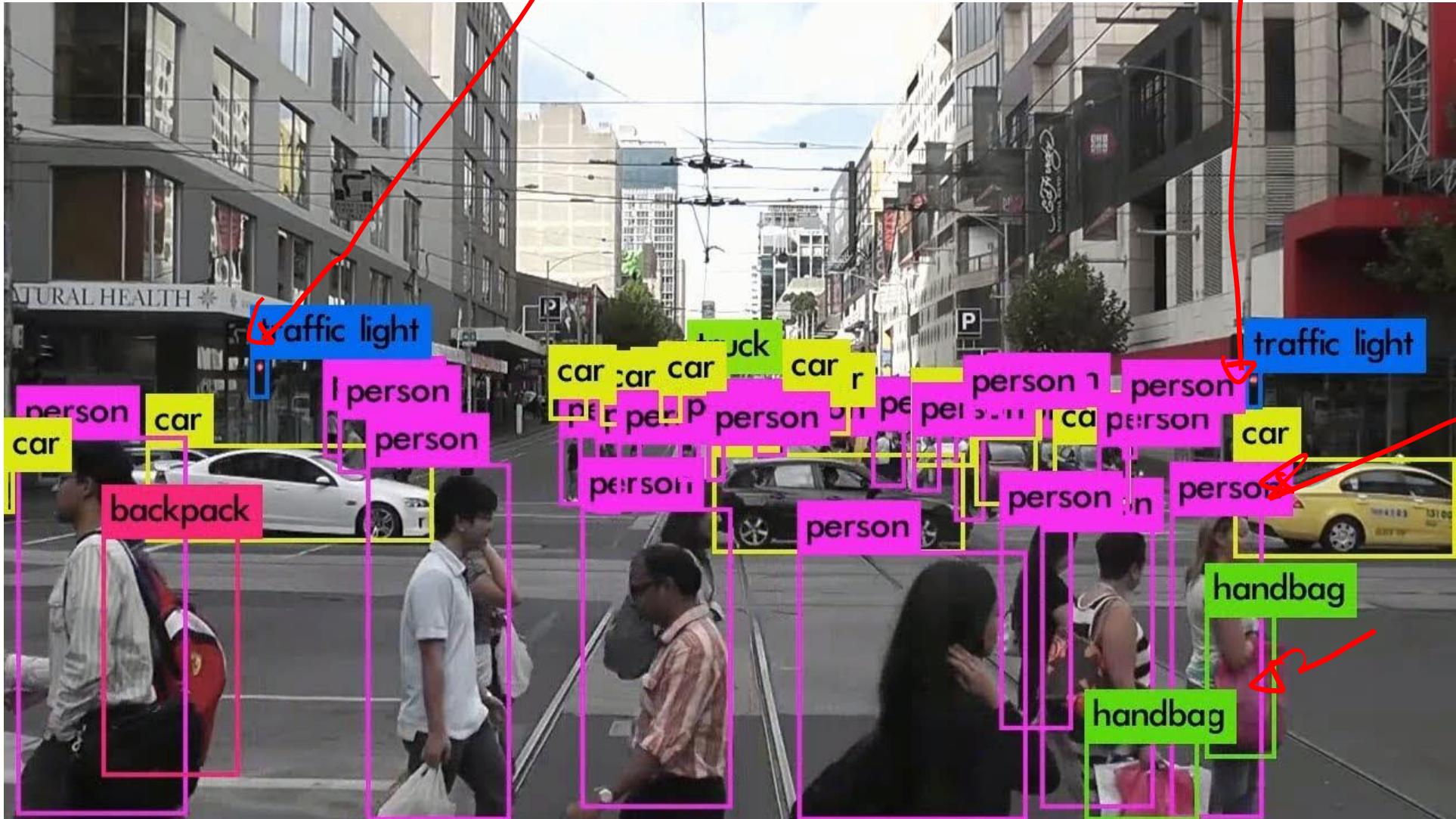
... which has grown incredibly fast in the last years

Computer Vision

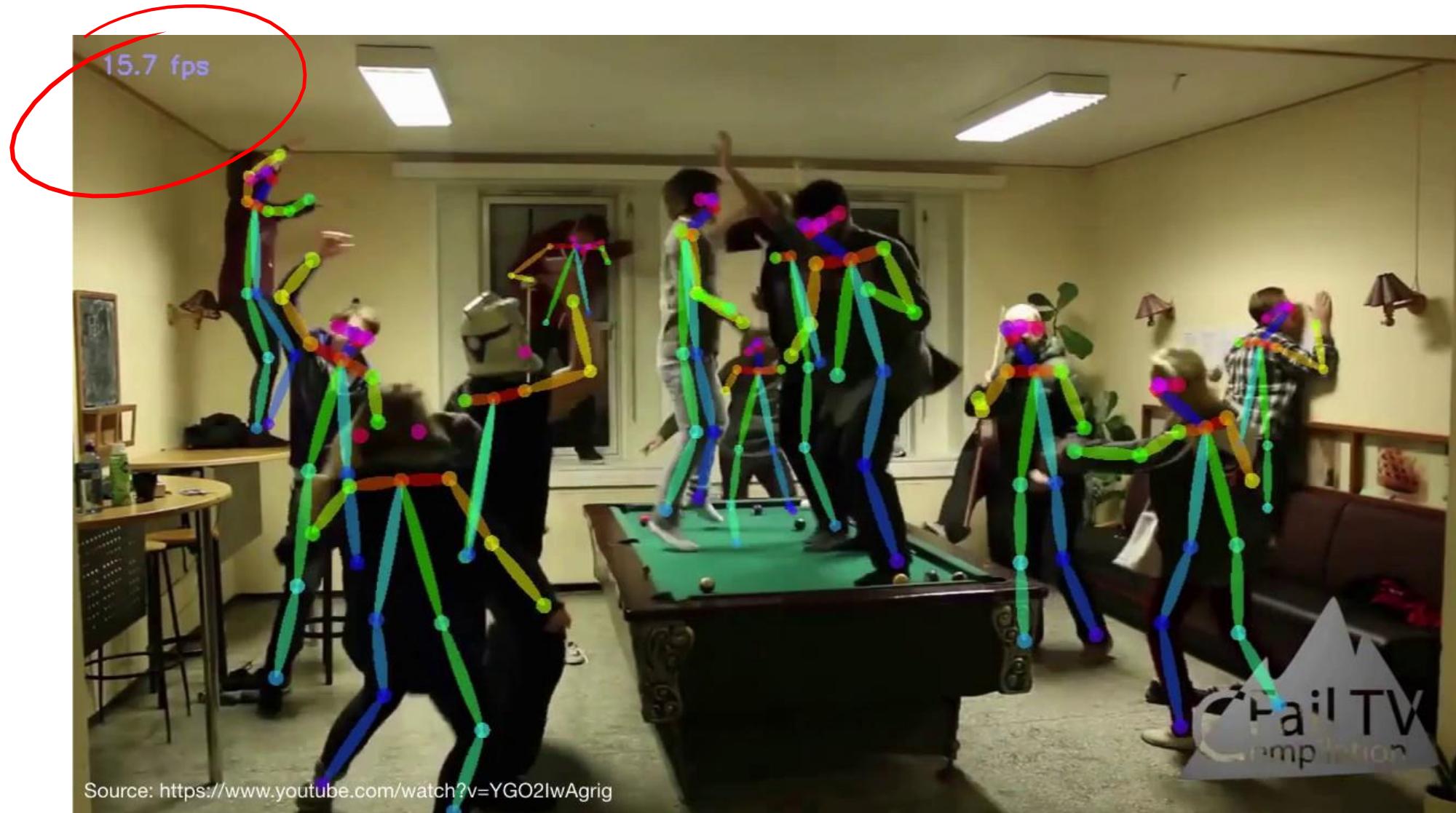
An interdisciplinary scientific field that deals with **how computers** can be made to gain **high-level understanding from digital images or videos**

... which has grown incredibly fast in the last years

Object Detection



Pose Estimation



Cao, Z., Simon, T., Wei, S. E., & Sheikh, Y. (2017). Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 7291-7299).

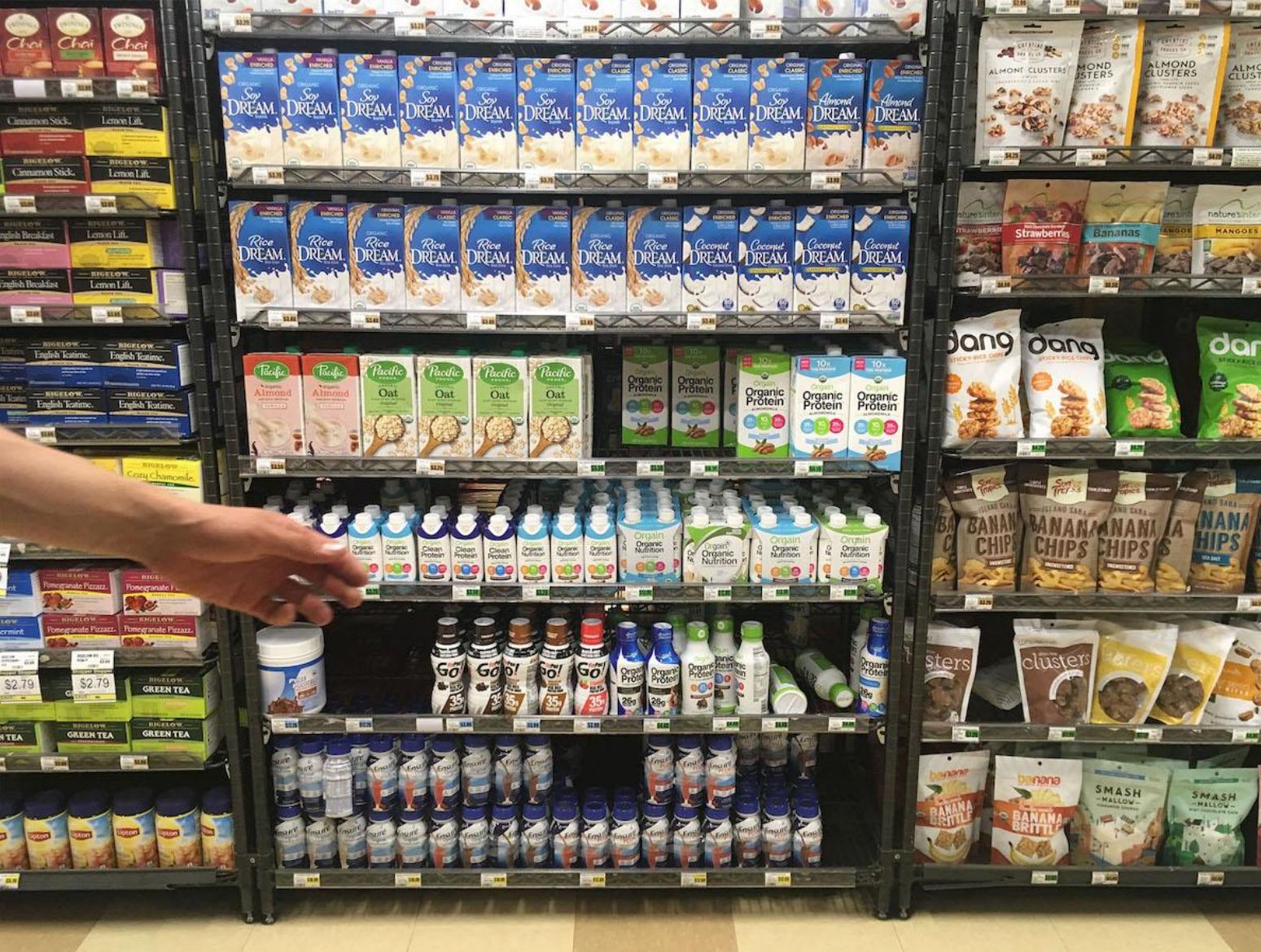
Autonomous Driving

By Dllu - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=64517567>



Ian Maddox [CC BY-SA (<https://creativecommons.org/licenses/by-sa/4.0/>)]

Automatic Shelf Analysis

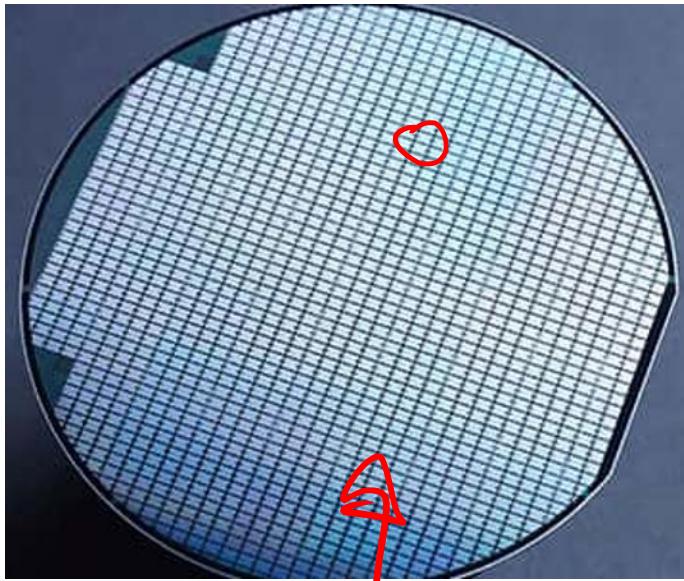


Automatic Shelf Analysis



Quality Inspection

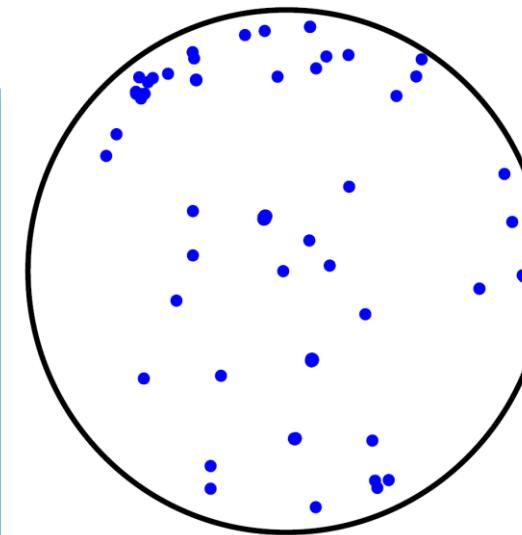
Wafer Defect Map (WDM)



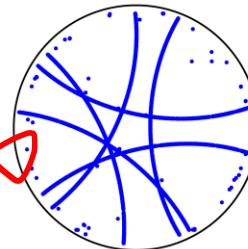
Silicon Wafer

Inspection
Tool

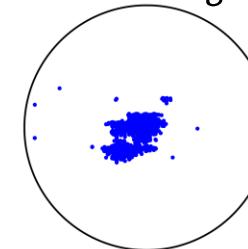
| X | Y |
|------|-------|
| 1863 | 709 |
| 1346 | 3067 |
| 2858 | 17095 |
| 3392 | 3508 |
| ... | ... |
| 282 | 6532 |
| 892 | 18888 |
| 4427 | 9873 |



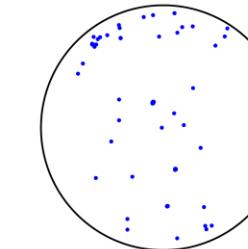
BasketBall



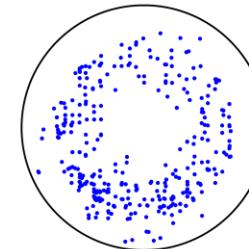
ClusterBig



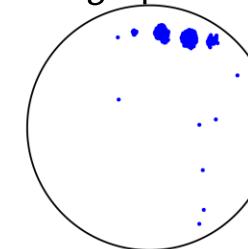
ClusterSmall



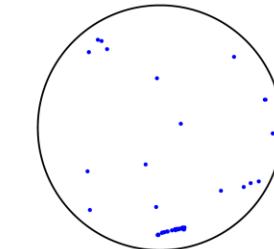
Donut



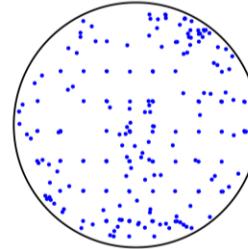
Fingerprints



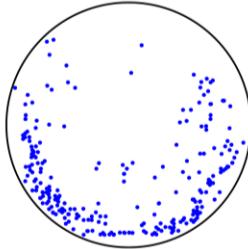
GeometricScratch



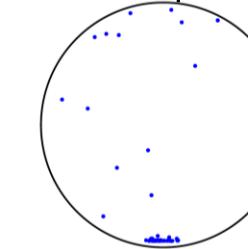
Grid



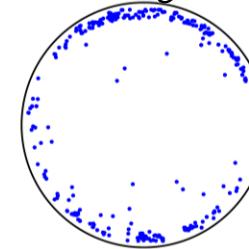
HalfMoon



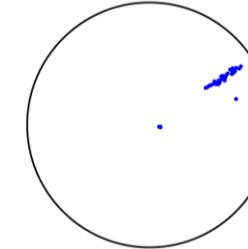
Incomplete



Ring



Slice



ZigZag

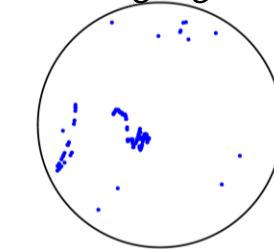


Image Captioning

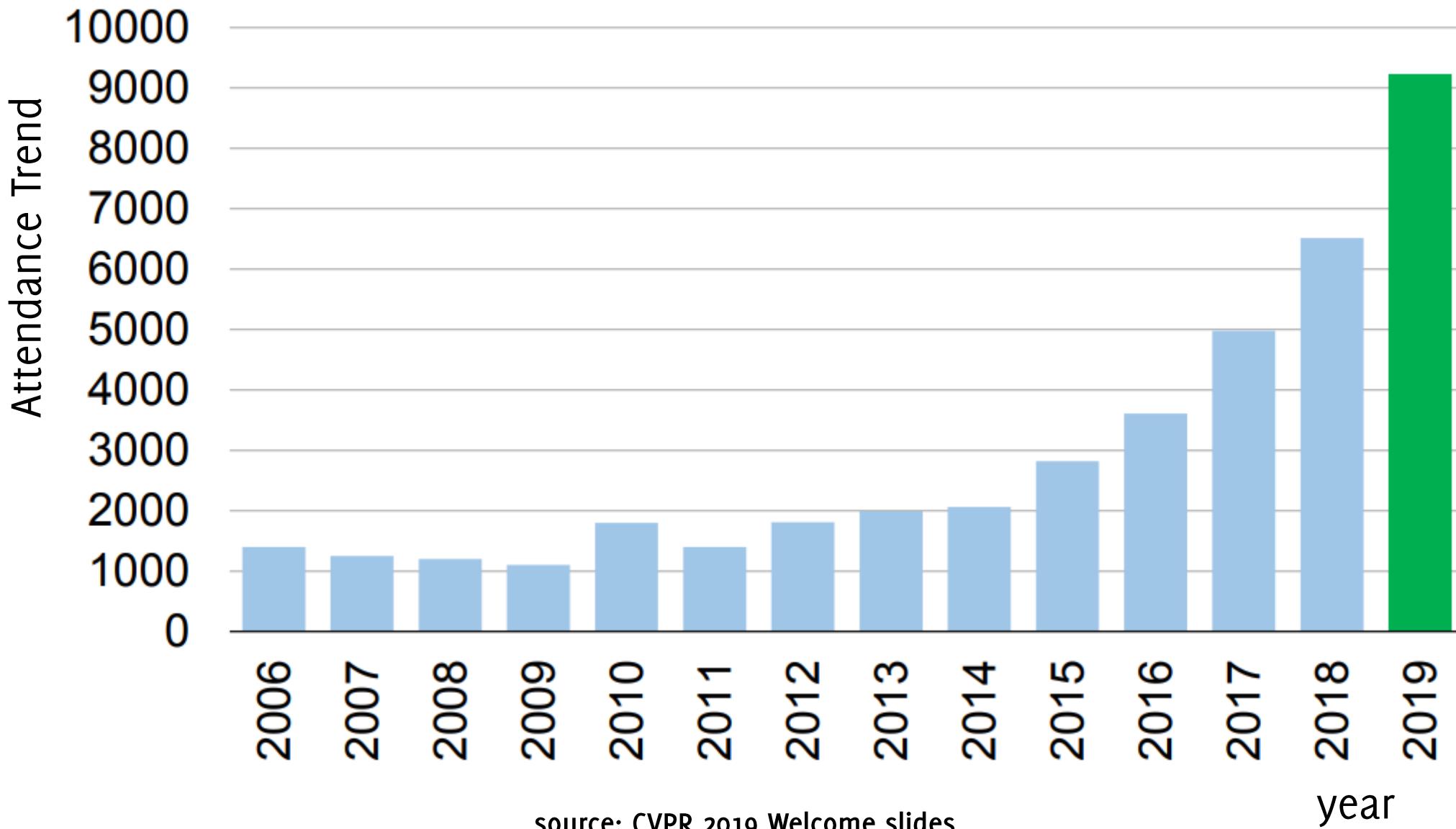


↗ "little girl is eating piece of cake."



"black cat is sitting on top of suitcase."

CVPR

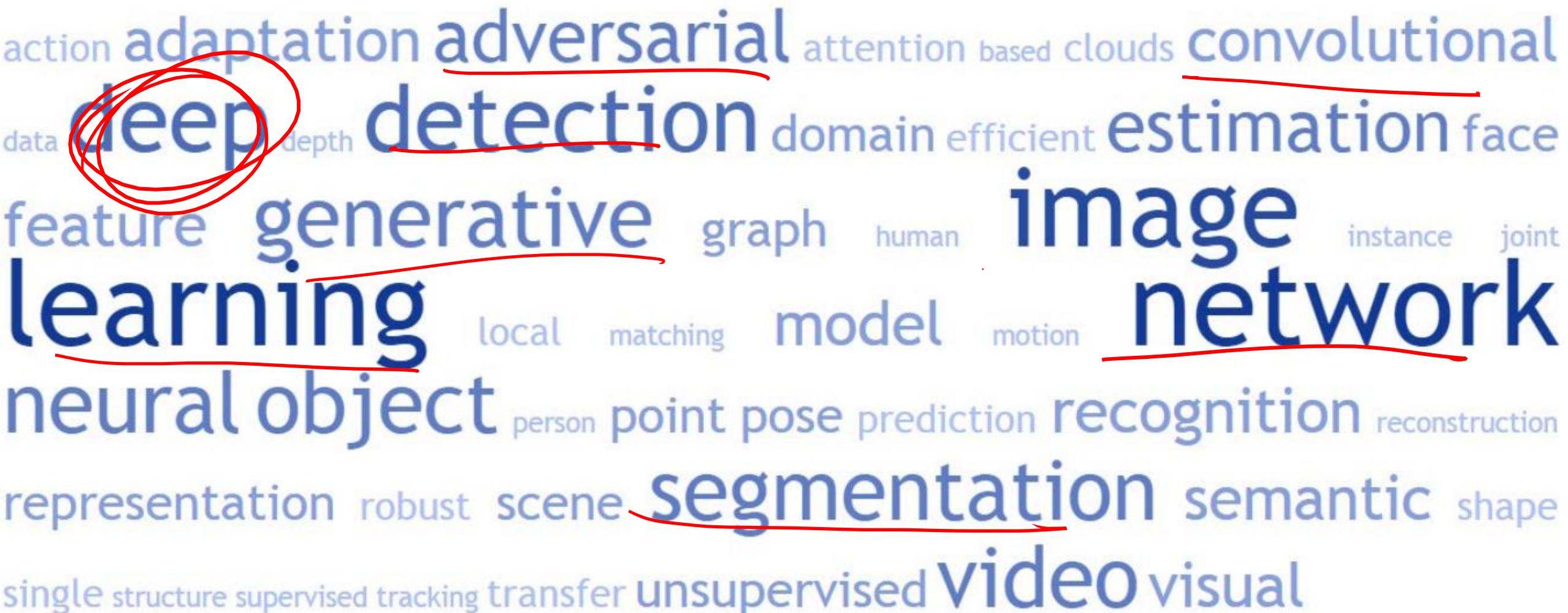


CVPR



1,294 papers in CVPR'19 (25.2% acceptance rate)

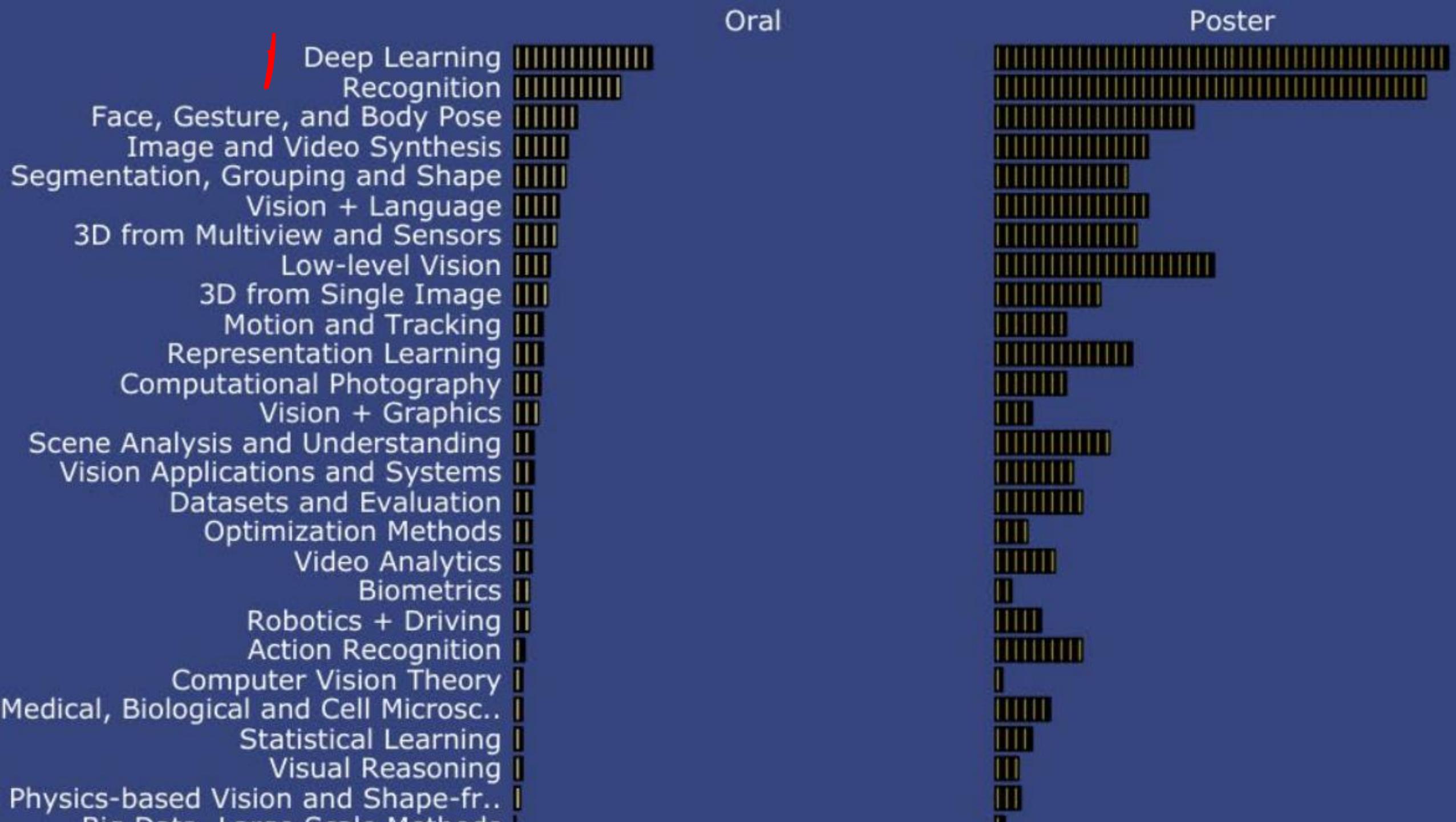
Acceptance rate is roughly even across topics



Lately, connection with ML

There has seen a dramatic change in CV:

- Once, most of techniques and algorithms **build upon a mathematical/statistical description of images**
- Nowadays, **machine-learning** methods are much more popular

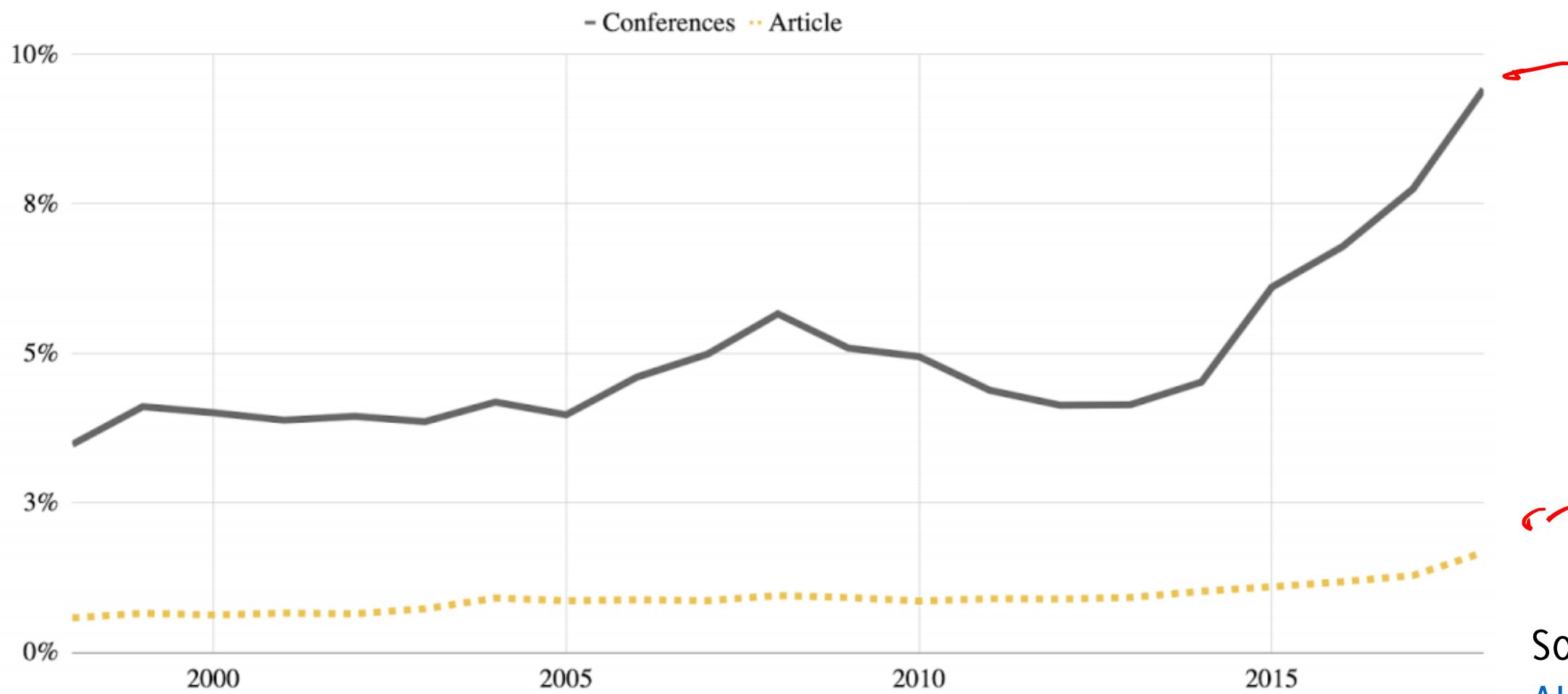


All in all...

If you plan pursuing a research-oriented career:

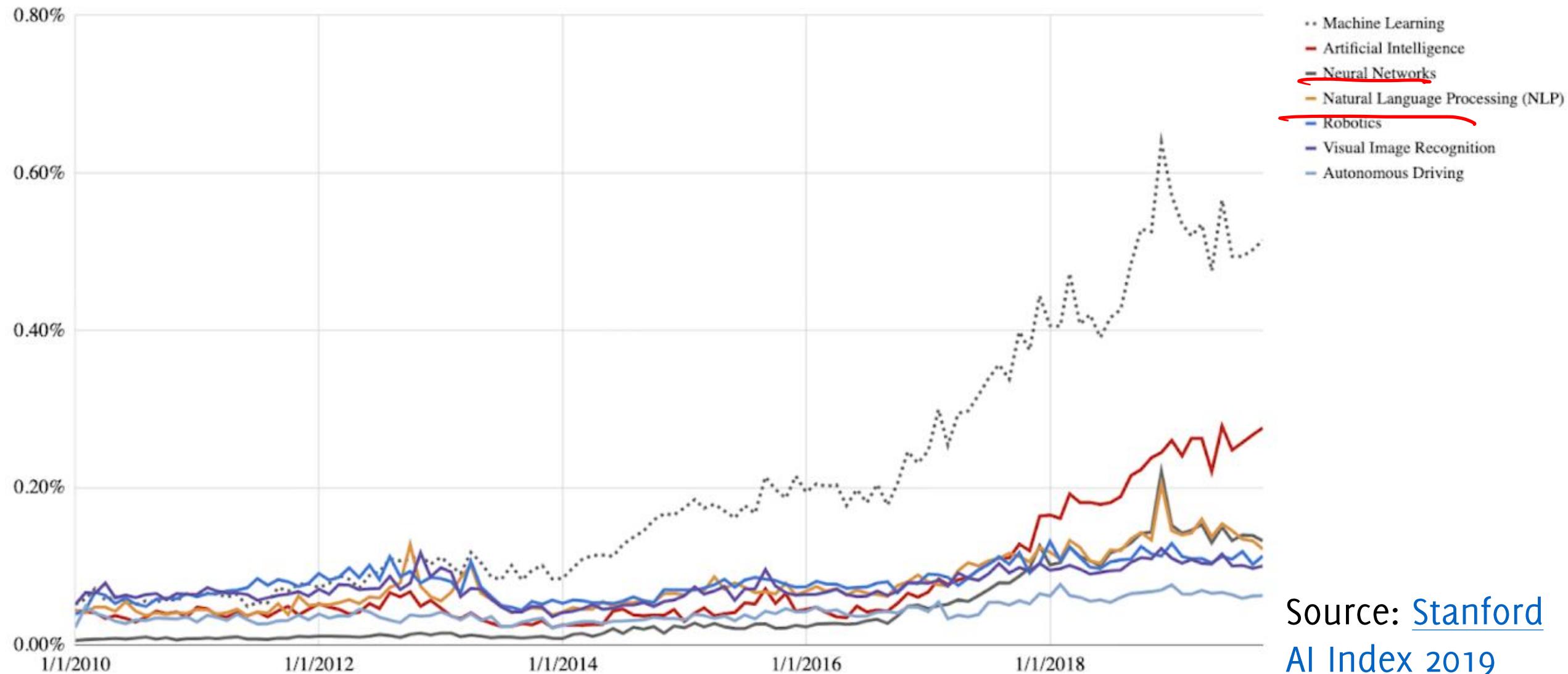
- If you go towards ML and/or CV, consider these fields are increasingly contaminated
- In any case a strong background in Computer Vision and Pattern Recognition is definitively a plus considering the widespread use of imaging data and the need of automation
- This course is a good way to practice more fundamental theory aspects related to neural networks
- Companies are increasingly interested in Computer Vision and Pattern Recognition application

AI publications over all the publications



Source: [Stanford AI Index 2019](#).

Percentages of AI posted job positions



Images



Photo Credits: Andrea Sanfilippo

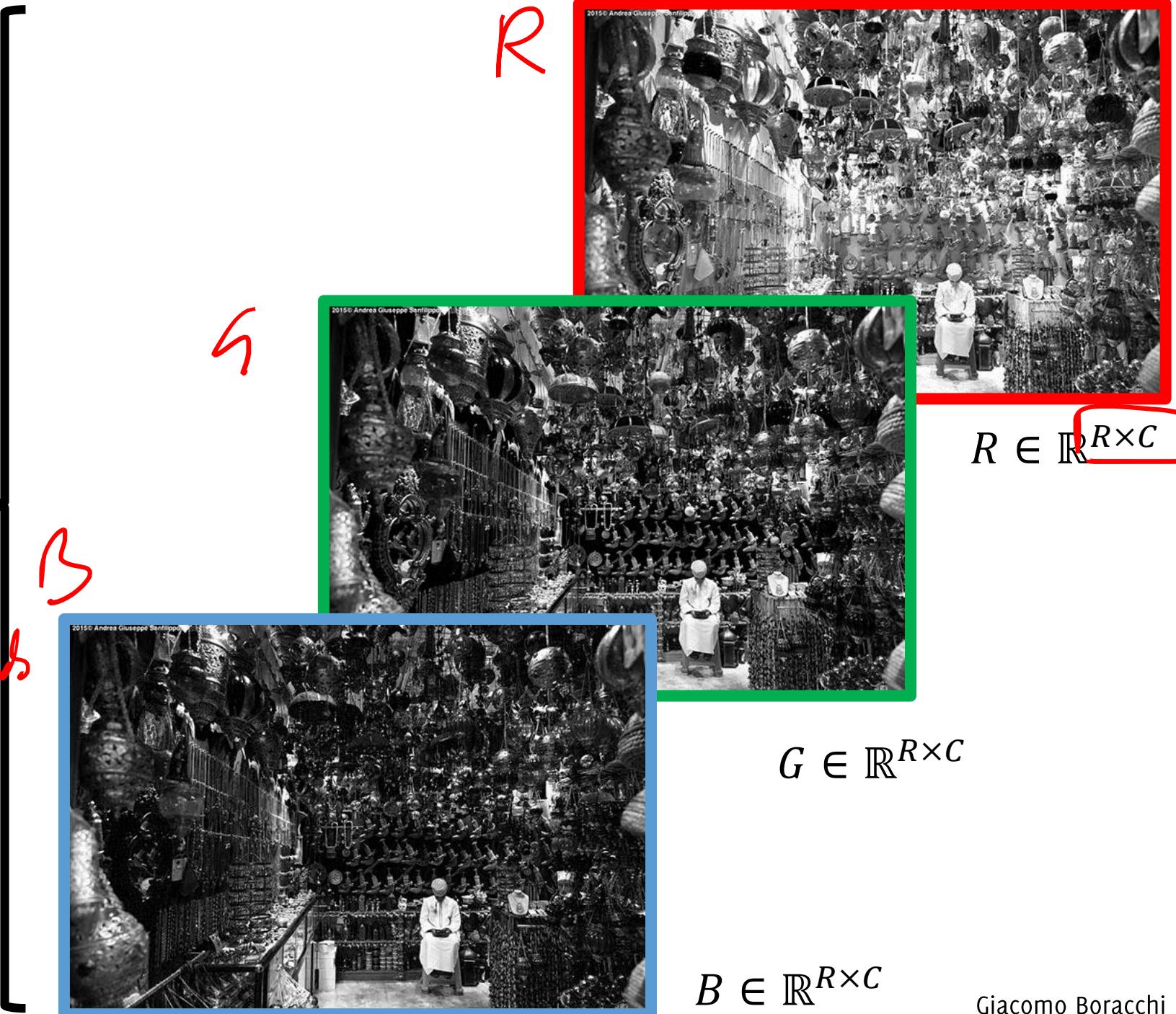
RGB Images



$I \in \mathbb{R}^{R \times C \times 3}$

Diagram illustrating the dimensions of an RGB image:

- The image is represented as $I \in \mathbb{R}^{R \times C \times 3}$.
- The width is labeled as $\# \text{cols}$.
- The height is labeled as $\# \text{rows}$.
- The depth is labeled as $\# \text{cols/rows}$.

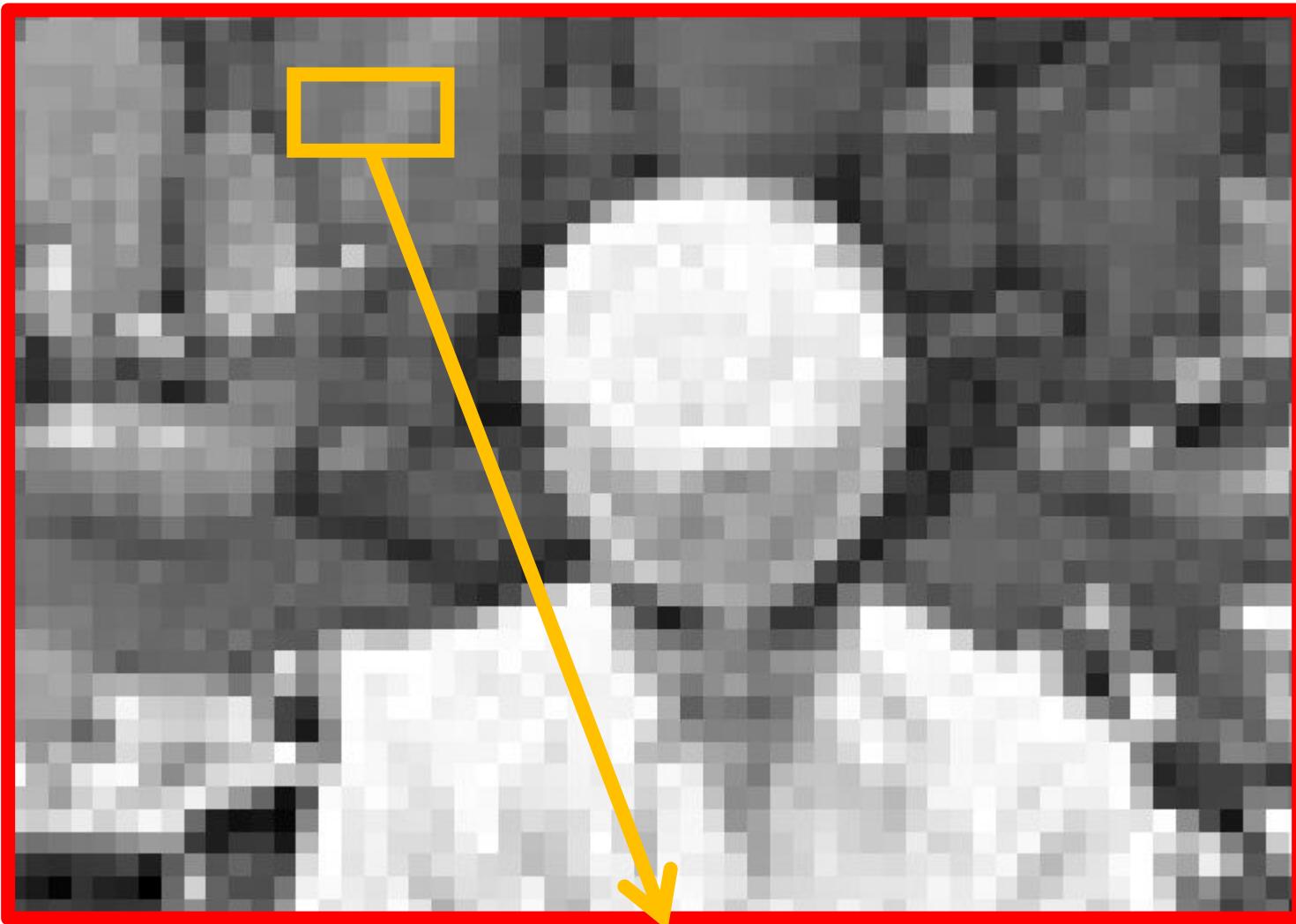


RGB Images

Images are saved by encoding each color information in 8 bits. So images are rescaled and casted in [0,255]



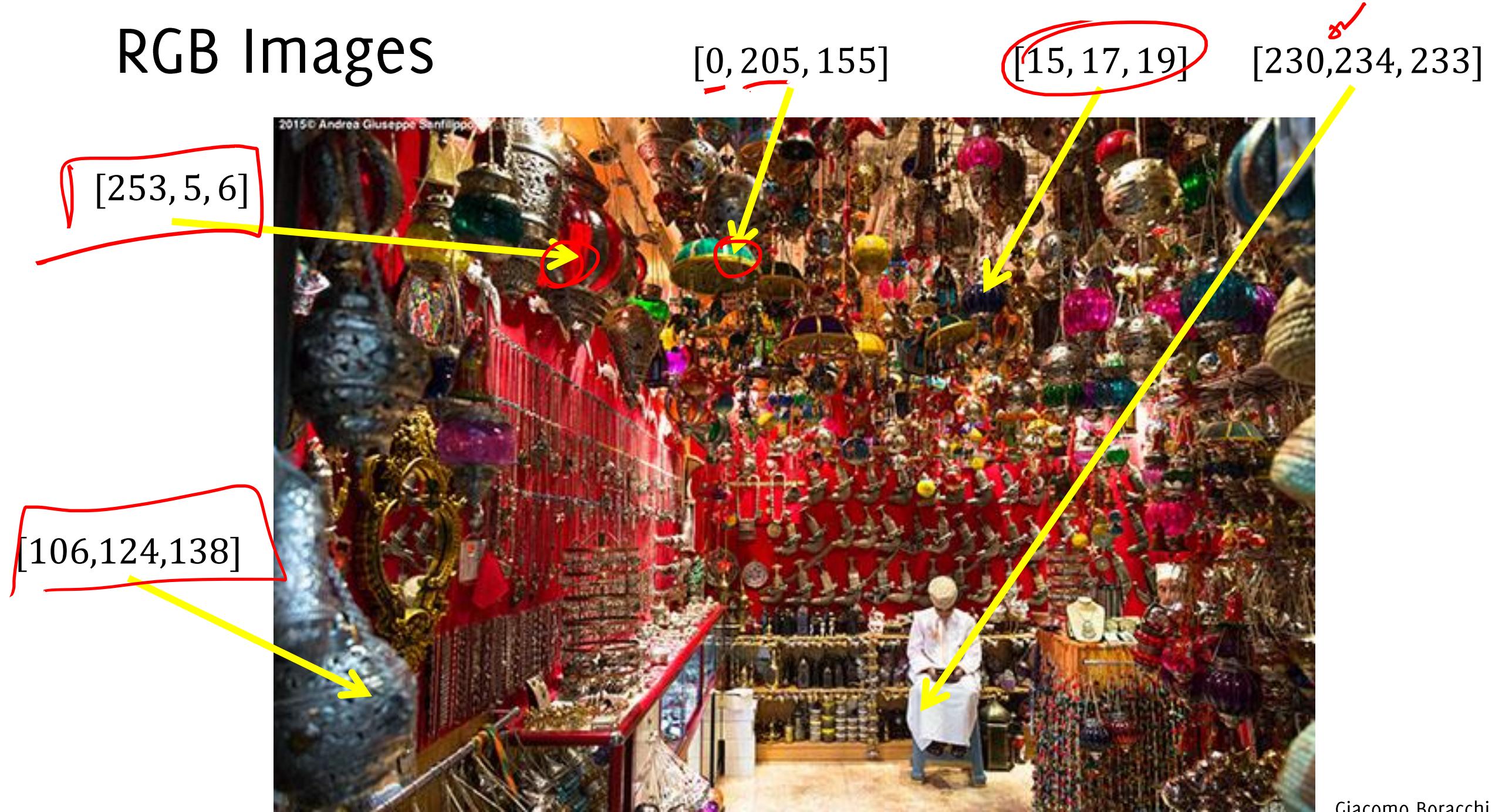
$$R \in \mathbb{R}^{R \times C}$$



| | | | | |
|-----|-----|-----|-----|-----|
| 123 | 122 | 134 | 121 | 132 |
| 122 | 121 | 125 | 132 | 124 |
| 119 | 127 | 137 | 119 | 139 |

1 Byte

RGB Images



The input of our classifier!

Three dimensional arrays $I \in \mathbb{R}^{R \times C \times 3}$

```
from skimage.io import imread  
  
# Read the image  
I = imread('bazar.jpg')  
  
# Extract the color channels of the image  
R = I[:, :, 0]  
G = I[:, :, 1]  
B = I[:, :, 2]
```

Videos

Higher dimensional images

Videos are sequences of images (frames)

If a frame is

$$I \in \mathbb{R}^{R \times C \times 3}$$

a video of T frames is

$$144 \times 180 \times 3 \times 30$$

$$V \in \mathbb{R}^{R \times C \times 3 \times T}$$

`print(v.shape)`
`(144, 180, 3, 30)`



In this example: $R = 144, C = 180$, thus these 5 color frames contains: 388.800 values in [0,255], thus in principle, 388 KB

Dimension Increases very quickly

Without compression: 1Byte per color per pixel

1 frame in full HD: $R = 1080, C = 1920 \approx 6MB$

1 sec in full HD (24fps) $\approx 150MB$

Fortunately, visual data are very redundant, thus compressible

This has to be taken into account when you design a Machine learning algorithm to be used on images

Higher dimensional images

Multispectral imaging and remote sensing

Higher dimensional Images

These images are stacking multiple layers as color-planes

Multi-spectral or Hyper-spectral images:

each band covers a certain wavelength that is meaningful
for interpretation soil in areal images

- 0.45-0.52 µm Blue-Green
- 0.52-0.60 µm Green
- 0.63-0.69 µm Red
- 0.76-0.90 µm Near IR
- 1.55-1.75 µm Mid-IR
- 10.40-12.50 µm Thermal IR
- 2.08-2.35 µm Mid-IR

**Classification of
multispectral images:**
infer the soil coverage
from the values in the
different spectral bands

In hyperspectral images the number of bands becomes
larger and you get a whole energy spectrum in each pixel

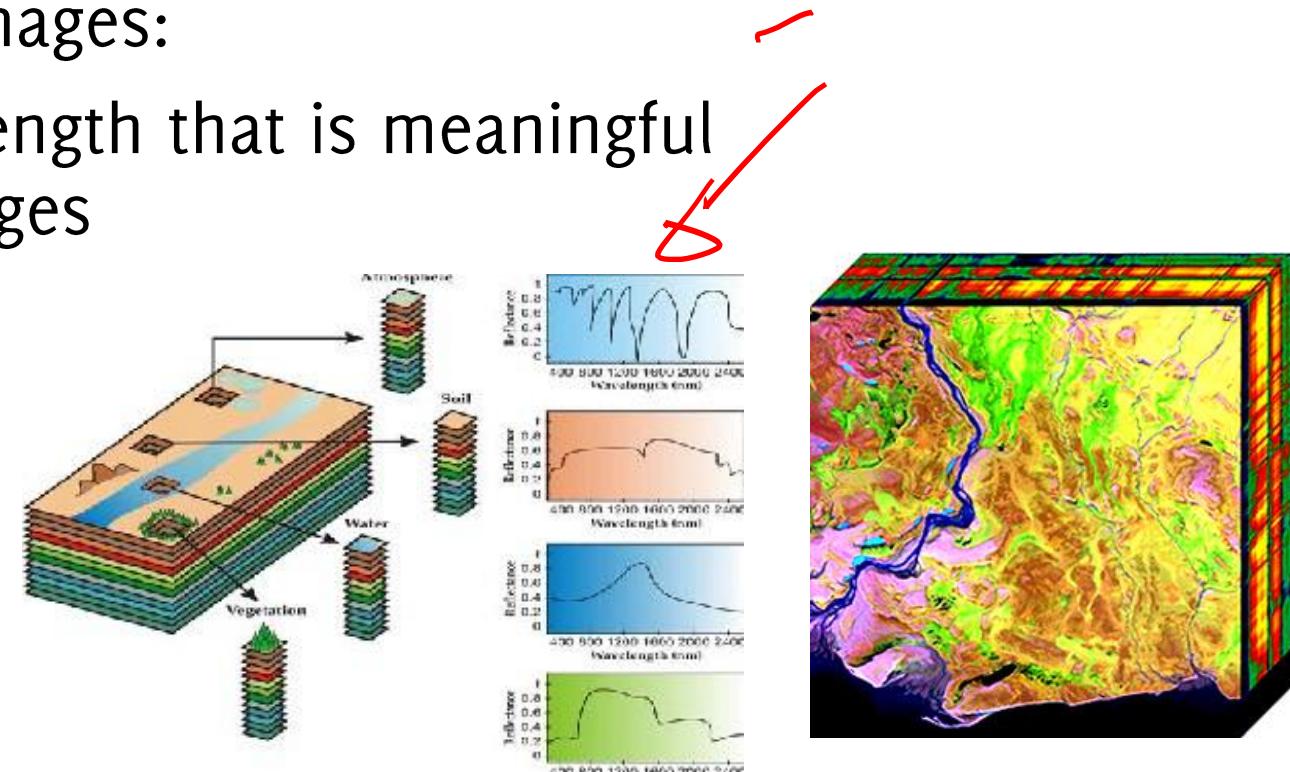
Higher dimensional Images

These images are stacking multiple layers as color-planes

Multi-spectral or Hyper-spectral images:

each band covers a certain wavelength that is meaningful for interpretation soil in areal images

- 0.45-0.52 µm Blue-Green
- 0.52-0.60 µm Green
- 0.63-0.69 µm Red
- 0.76-0.90 µm Near IR
- 1.55-1.75 µm Mid-IR
- 10.40-12.50 µm Thermal IR
- 2.08-2.35 µm Mid-IR



In hyperspectral images the number of bands becomes larger and you get a whole energy spectrum in each pixel

Band 1



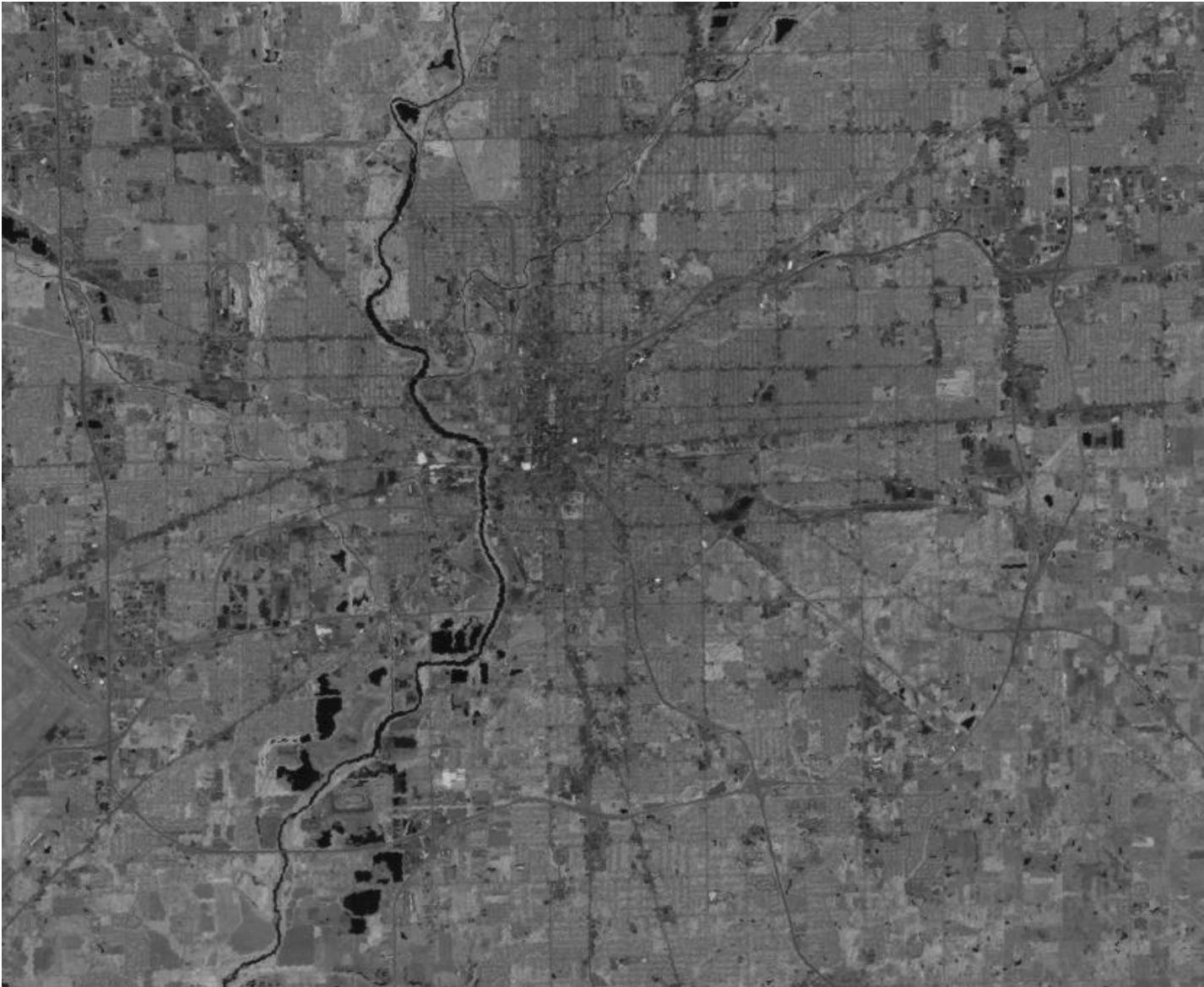
Band 2



Band 3



Band 4



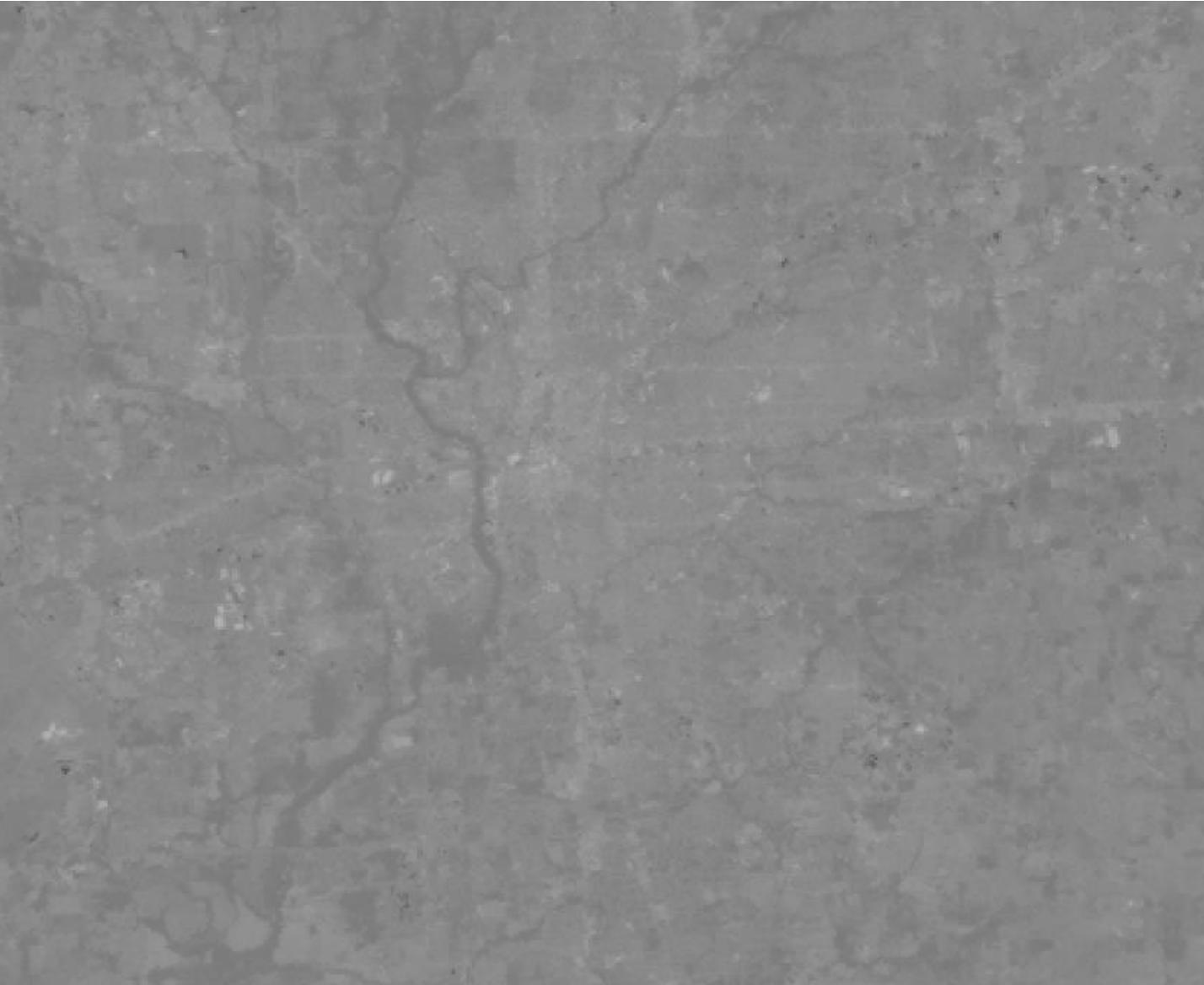
Band 5



Band 6



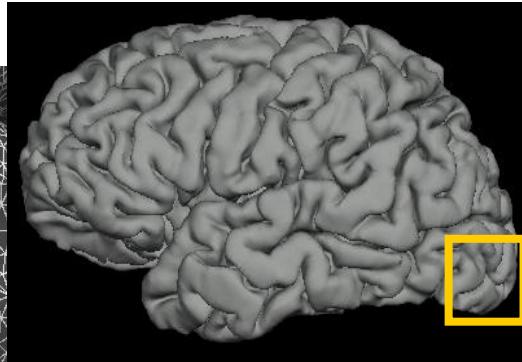
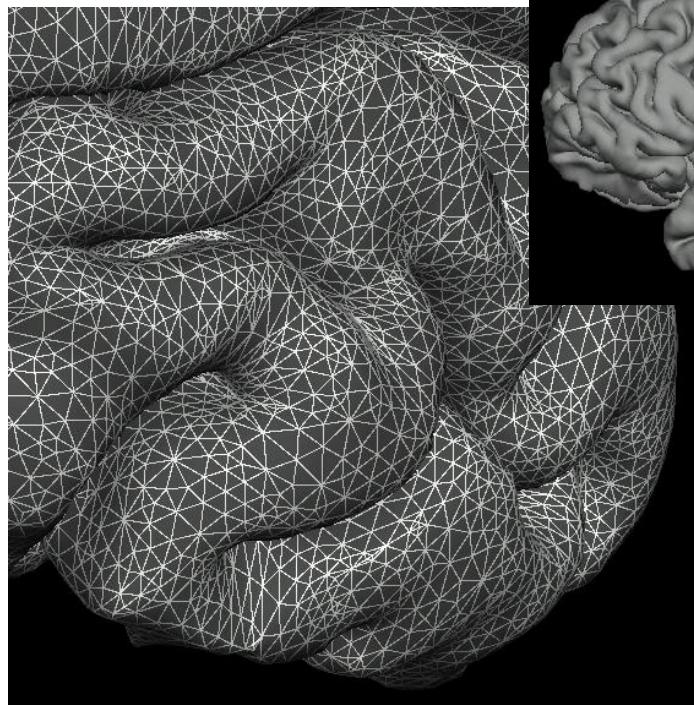
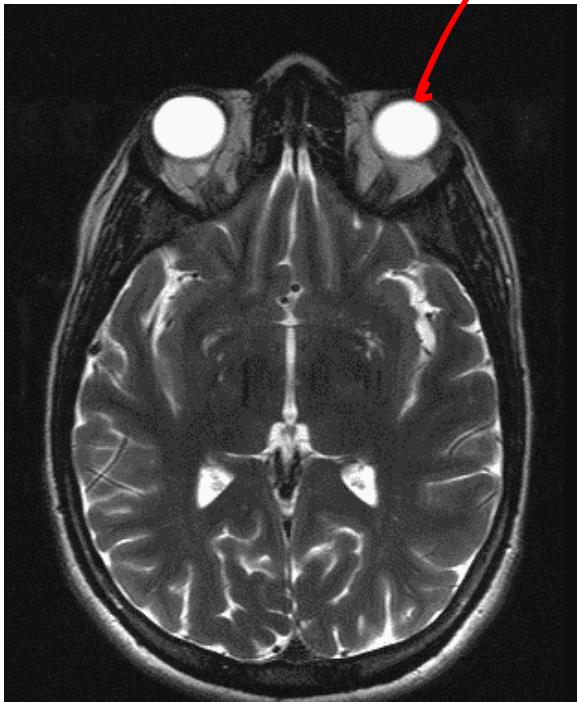
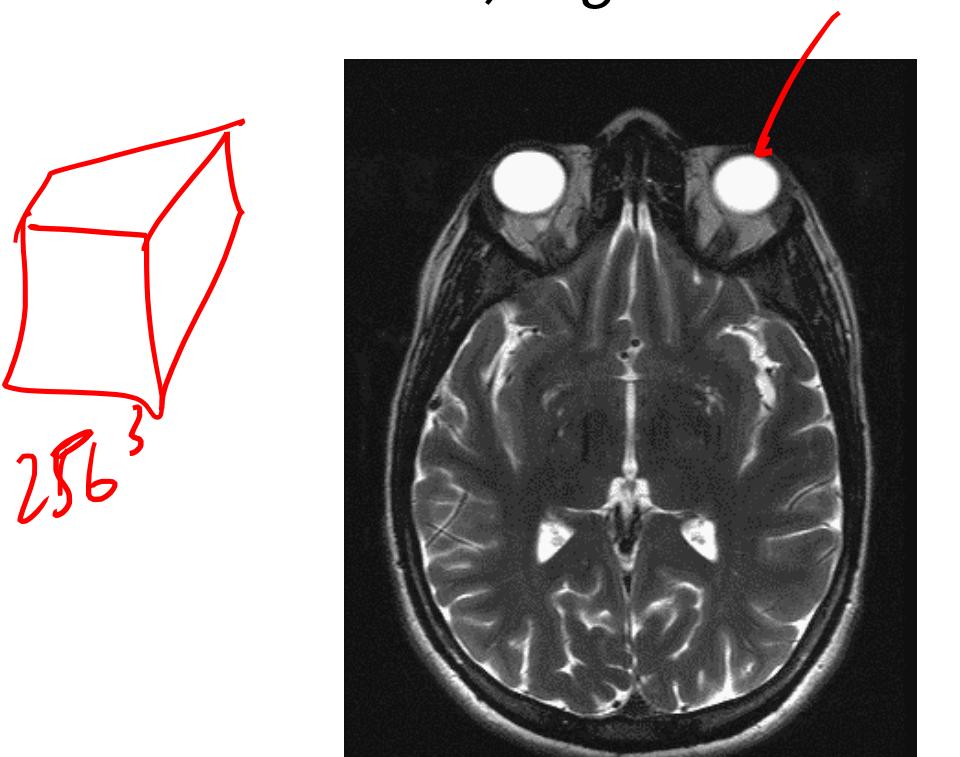
Band 7



MRI and TAC

(Structural) MRI provide a 3D stack of grayscale images that are analyzing the body at different depth levels

These can be used to perform 3D reconstruction of bones, organs and membranes, e.g. cortical surfaces



Local (Spatial) Transformations: Correlation (and Convolution)

Most important image processing operations for
classification problems

Local (Spatial) Transformation

In general, these can be written as

$$G(r, c) = T_U[I(r, c)]$$

Where

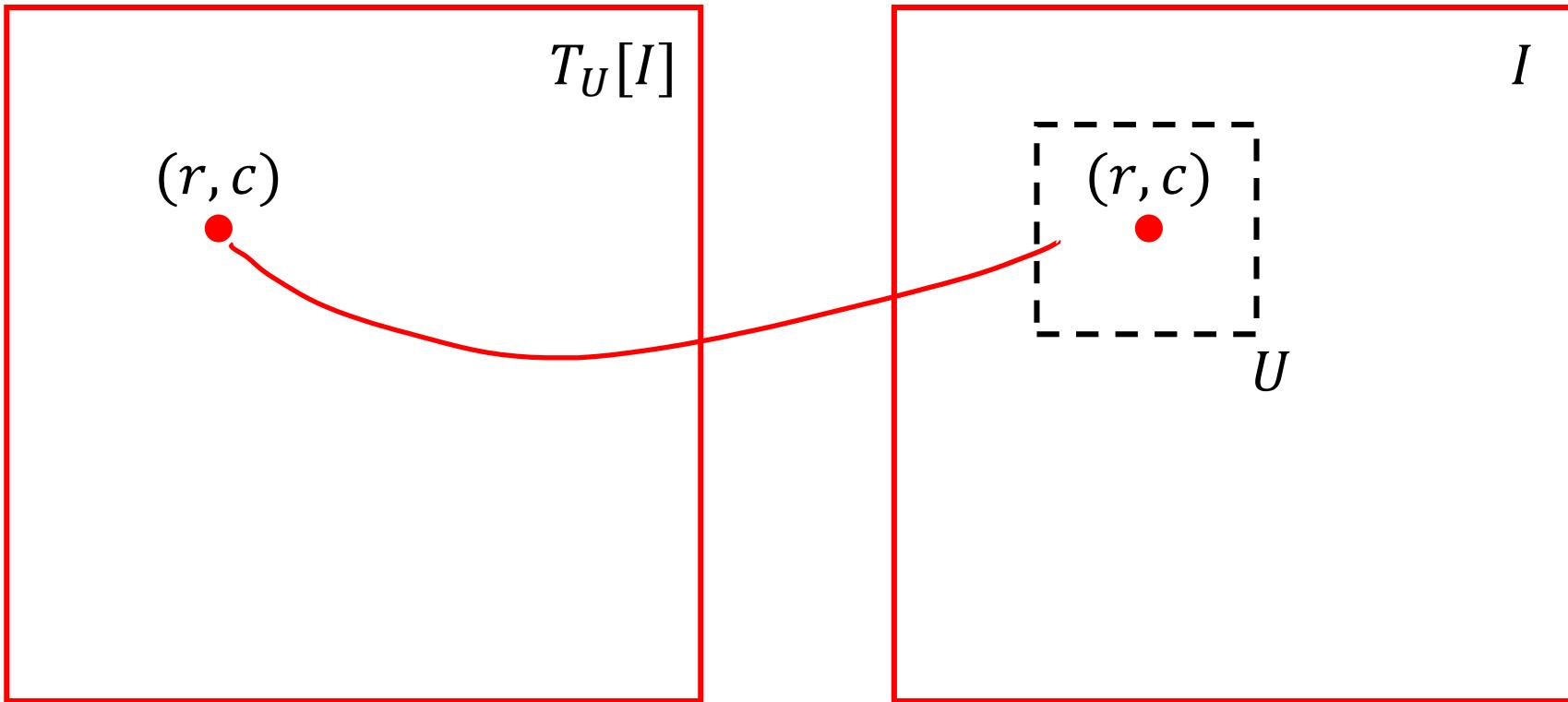
- I is the input image to be transformed
- G is the output
- $T_U: \mathbb{R}^3 \rightarrow \mathbb{R}^3$ or $T_U: \mathbb{R}^3 \rightarrow \mathbb{R}$ is a function
- U is a neighbourhood, identifies a region of the image that will concur in the output definition

T operates on I “around” U

The output at pixel (r, c) i.e., $T_U[I(r, c)]$ is defined by all pixels $\{I(x, y), (x - r, y - c) \in U\}$

Local (Spatial) Filters

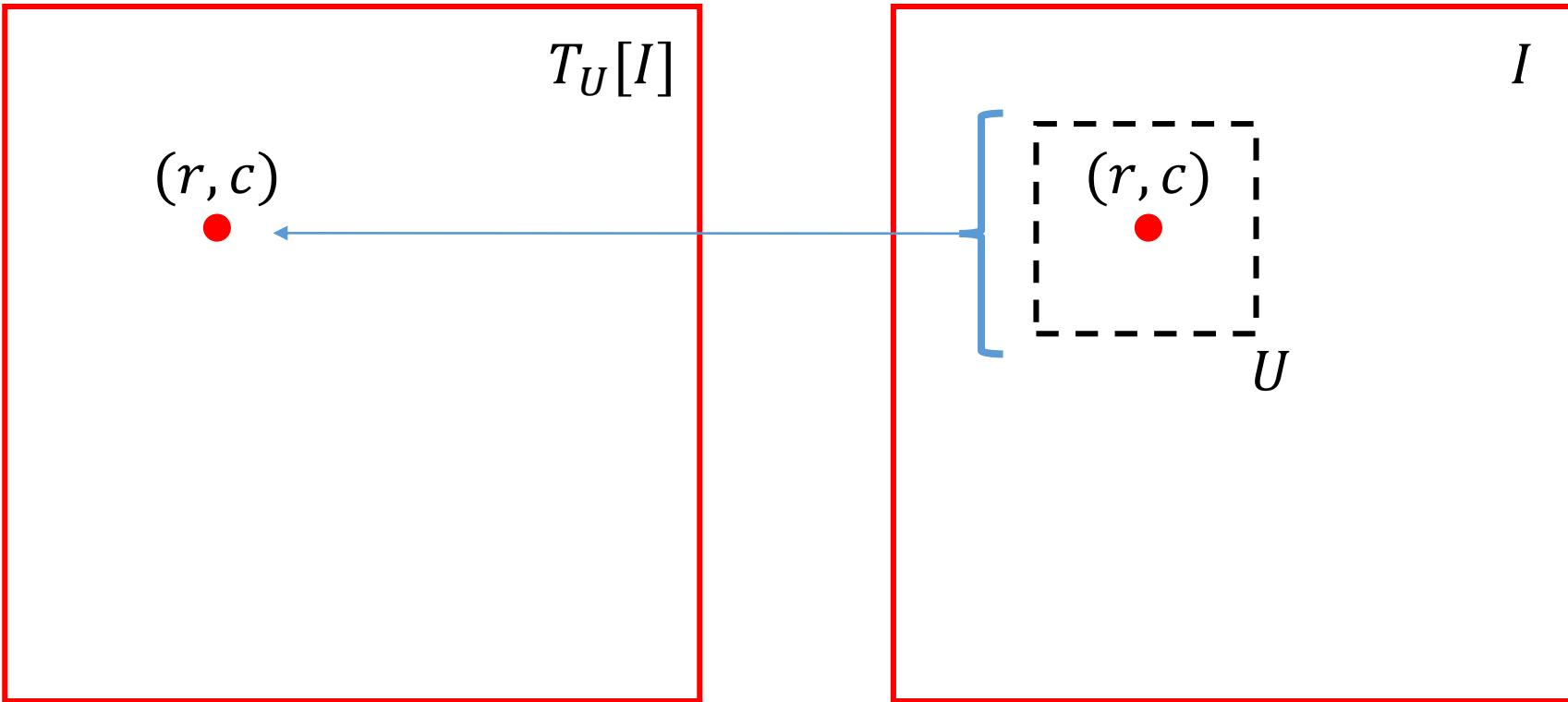
The dashed square represents $\{I(x, y), (x - r, y - c) \in U\}$



- The location of the output does not change
- The operation is repeated for each pixel
- T can be either linear or nonlinear

Local (Spatial) Filters

The dashed square represents $\{I(x, y), (x - r, y - c) \in U\}$



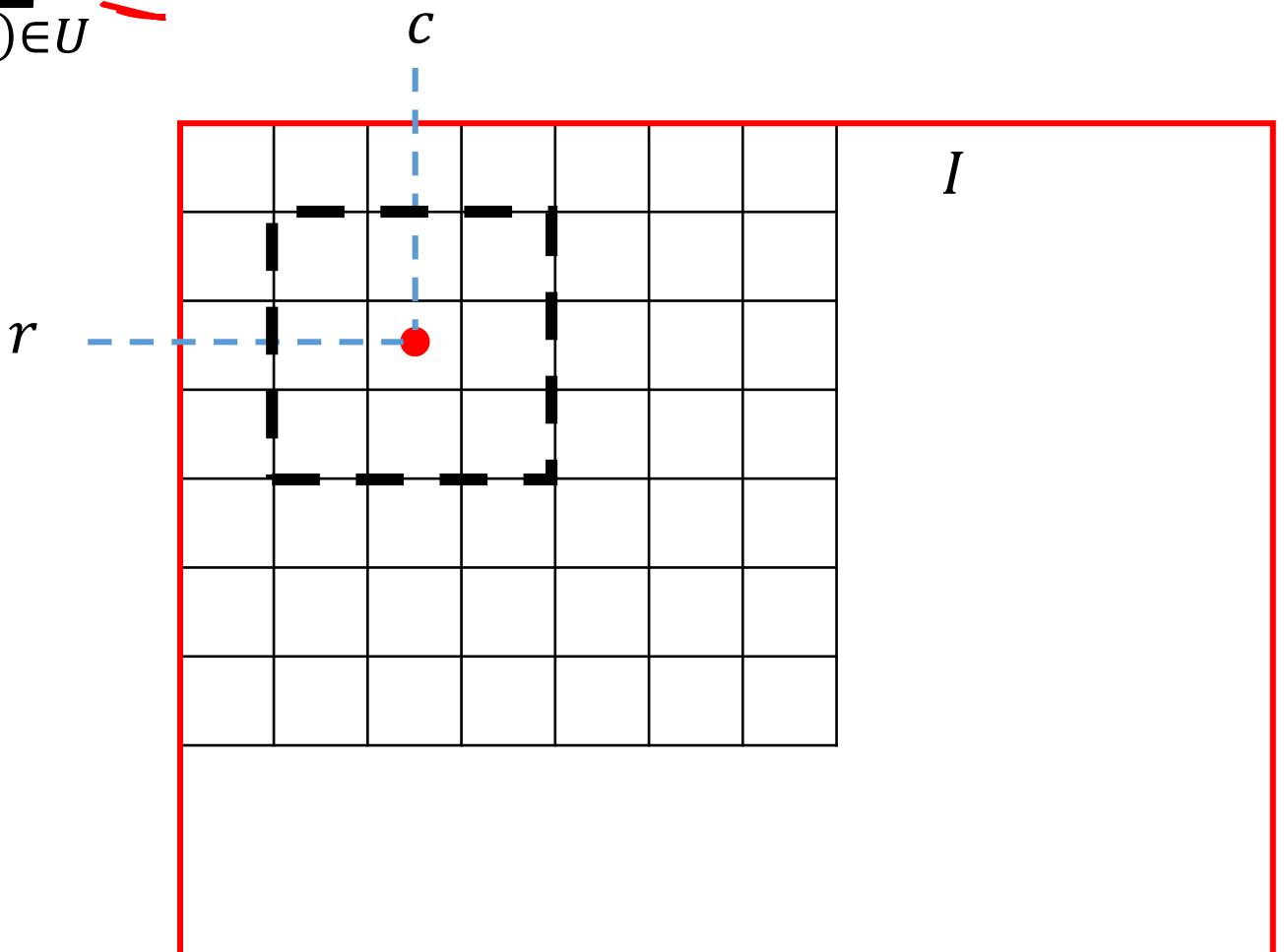
- The location of the output does not change
- The operation is repeated for each pixel
- T can be either linear or nonlinear

Local Linear Filters

Linear Transformation: Linearity implies that

$$T(I(r, c)) = \sum_{(x,y) \in U} w_i * I(r + x, c + y)$$

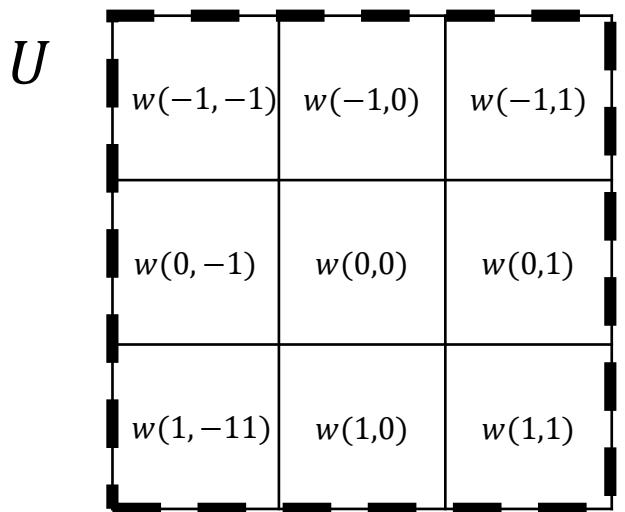
Considering some weights $\{w_i\}$



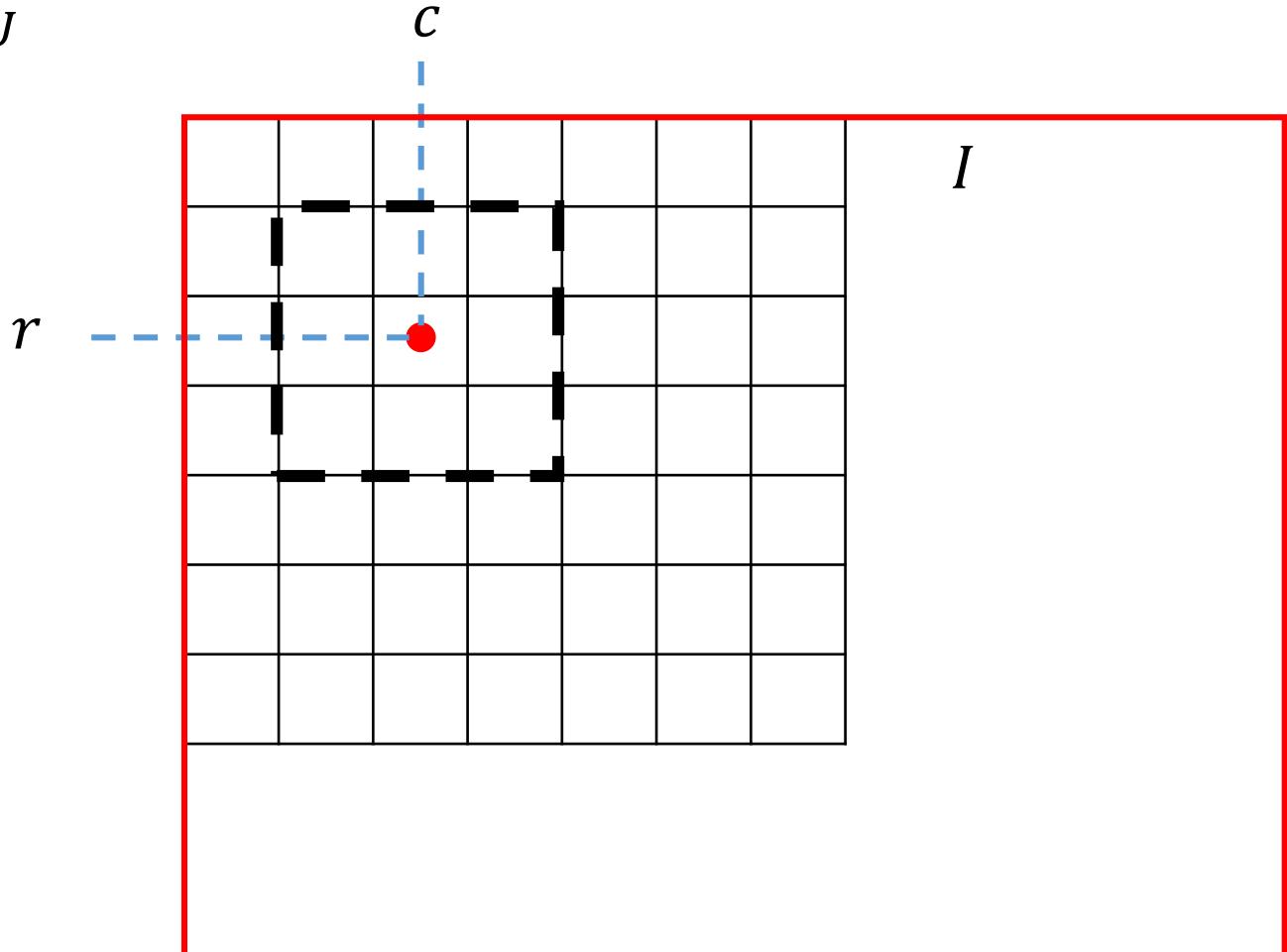
Local Linear Filters

Linear Transformation: Linearity implies that

$$T(I(r, c)) = \sum_{(x,y) \in U} w(x, y) * I(r + x, c + y)$$



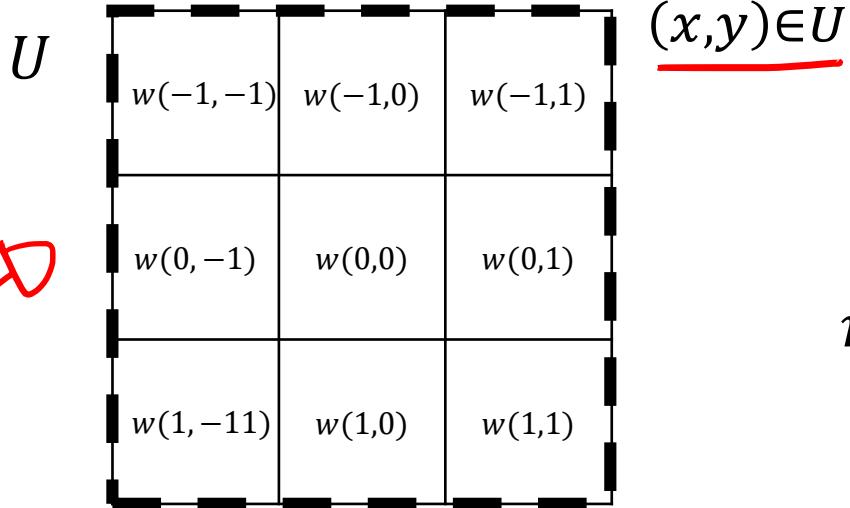
We can consider weights
as an image, or a filter h
The filter h entirely defines
this operation



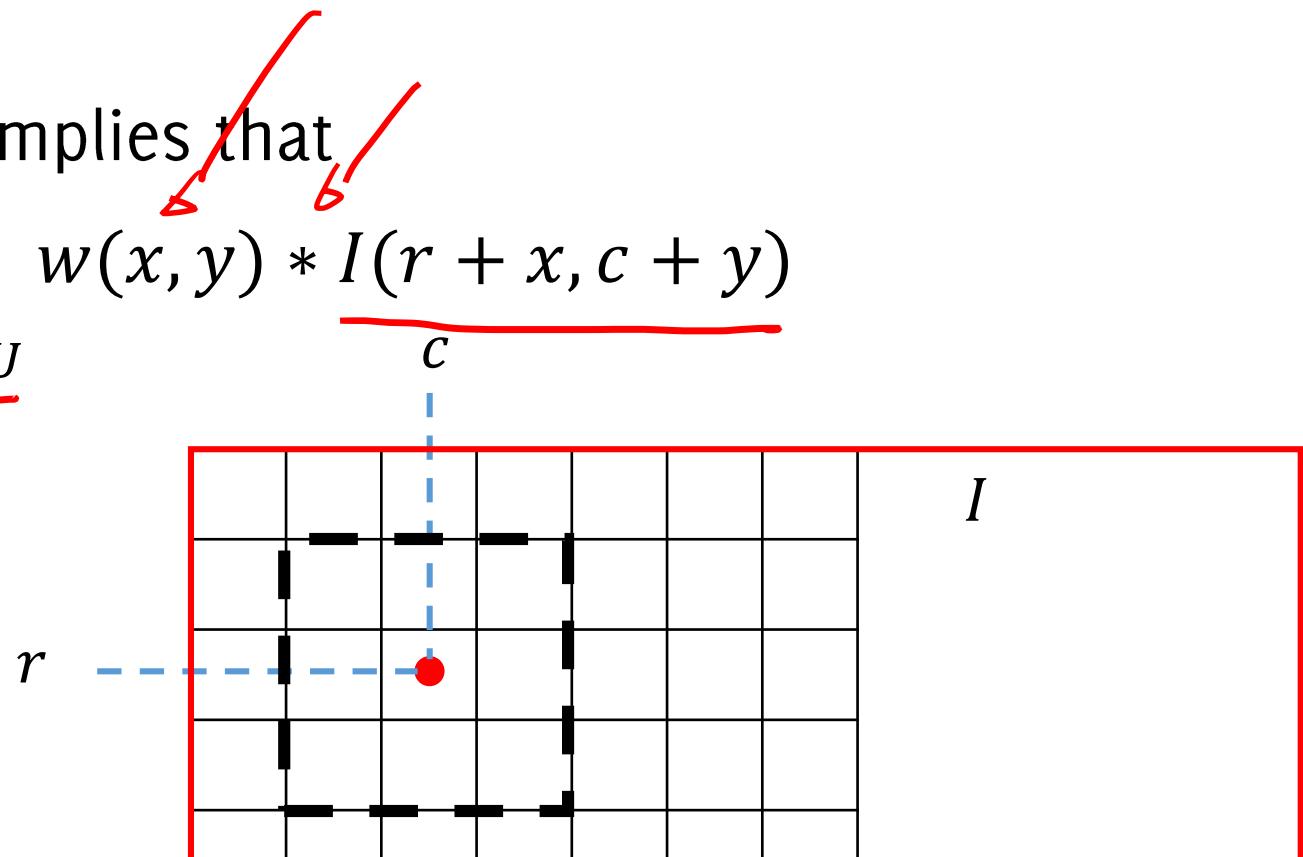
Local Linear Filters

Linear Transformation: Linearity implies that

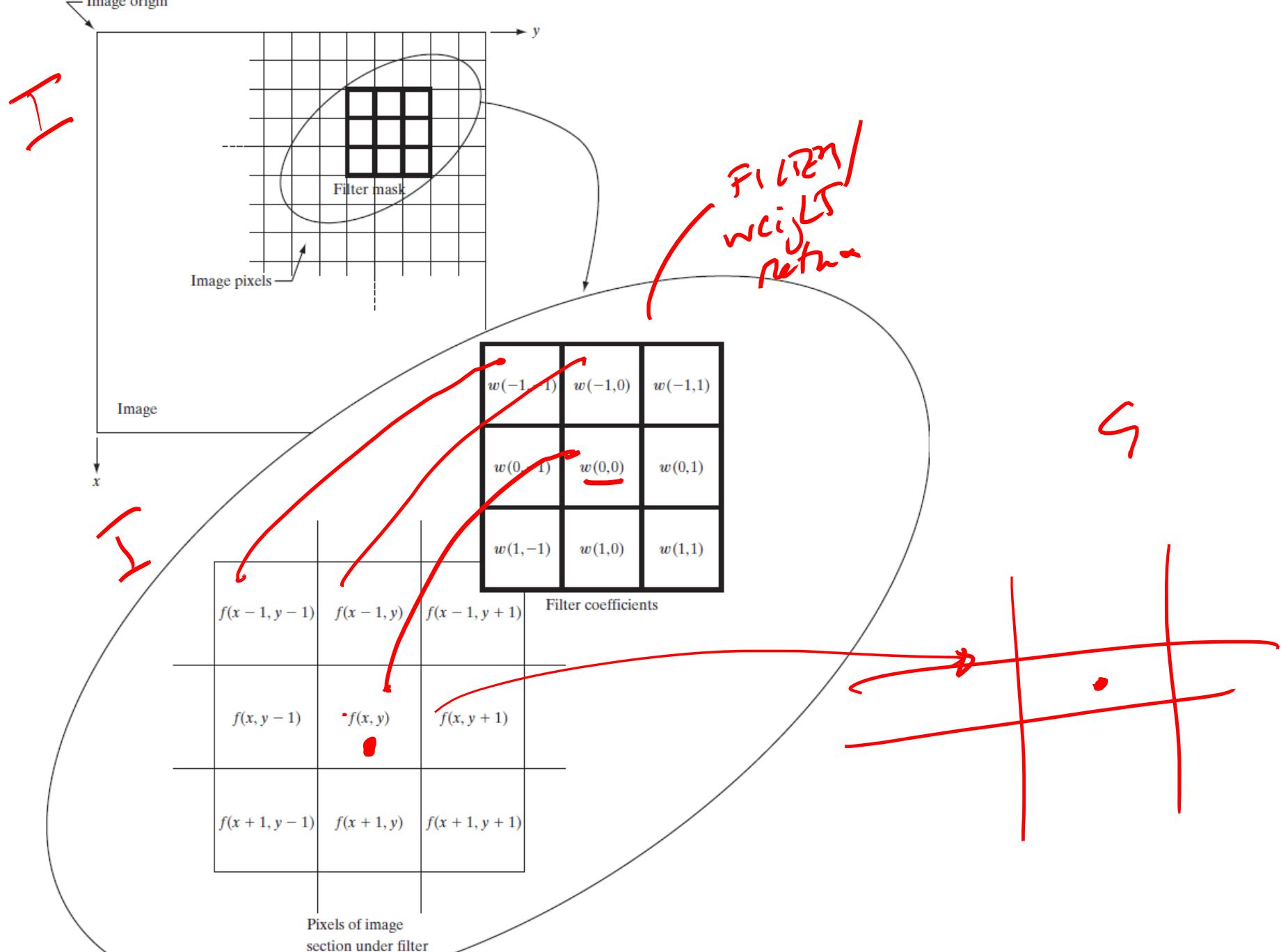
$$T(I(r, c)) = \sum_{(x,y) \in U} w(x, y) * I(r + x, c + y)$$



We can consider weights as an image, or a filter h
The filter h entirely defines this operation



This operation is repeated for each and every pixel in the input image



Correlation

The **correlation** among a filter h and an image is defined as

$$(I \otimes h) = \sum_{u=-L}^L \sum_{v=-L}^L w(u, v) * I(r + u, c + v)$$

where the filter w is of size $(2L + 1) \times (2L + 1)$

Correlation

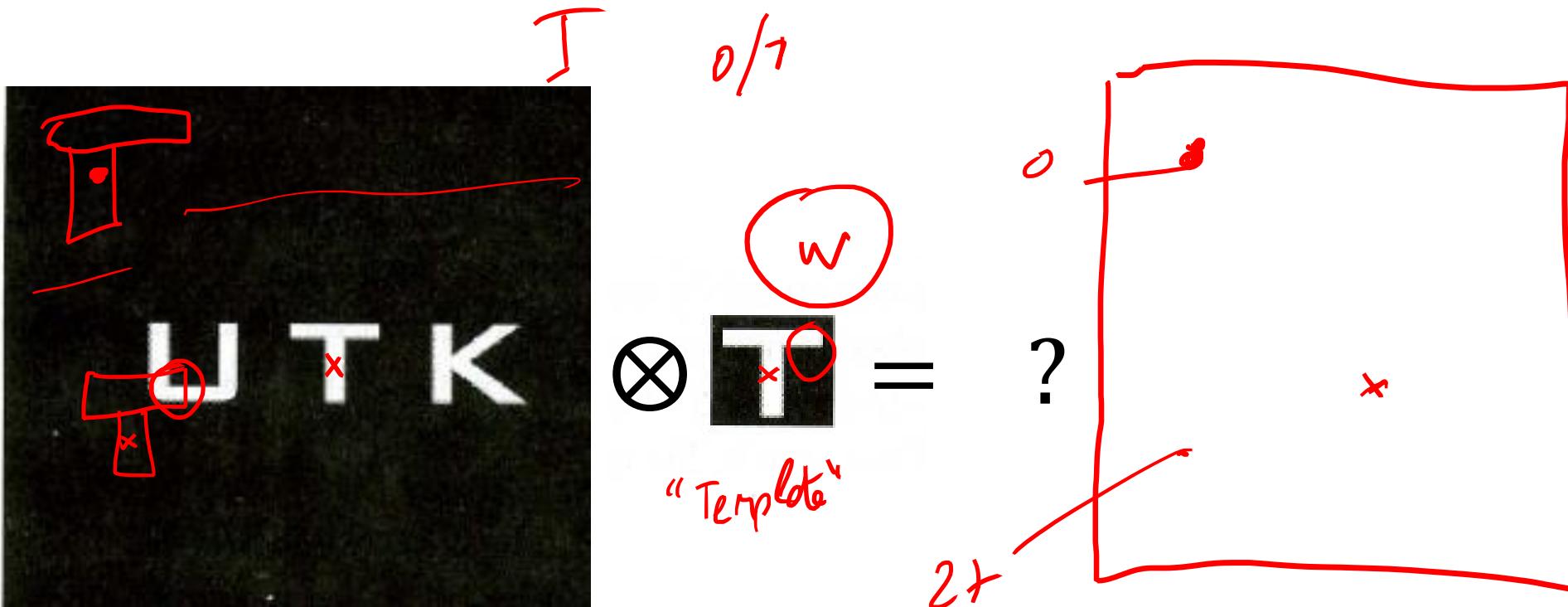
The **correlation** among a filter h and an image is defined as

$$(I \otimes h) = \sum_{u=-L}^L \sum_{v=-L}^L w(u, v) * I(r + u, c + v)$$

where the filter w is of size $(2L + 1) \times (2L + 1)$

```
from scipy import signal  
  
# correlation  
T_corr = signal.correlate(I, h)
```

Correlation for BINARY target matching



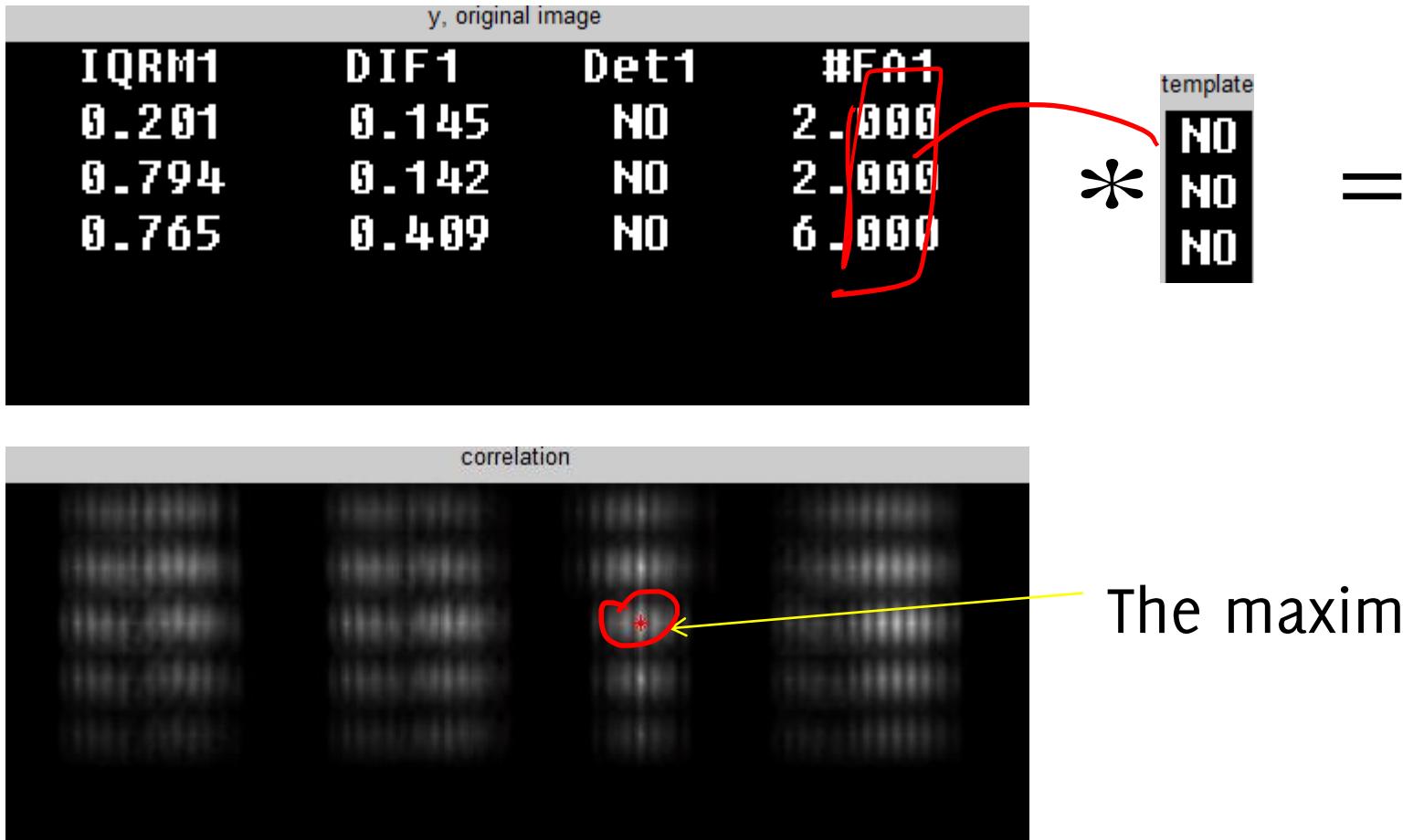
Easy to understand with binary images

Target used as a filter

Another example

$$\begin{matrix} & & & \text{y, original image} \\ \text{IQRM1} & \text{DIF1} & \text{Det1} & \#FA1 \\ 0.201 & 0.145 & \text{NO} & 2.000 \\ 0.794 & 0.142 & \text{NO} & 2.000 \\ 0.765 & 0.409 & \text{NO} & 6.000 \end{matrix} \quad * \quad \begin{matrix} \text{template} \\ \text{NO} \\ \text{NO} \\ \text{NO} \end{matrix} =$$

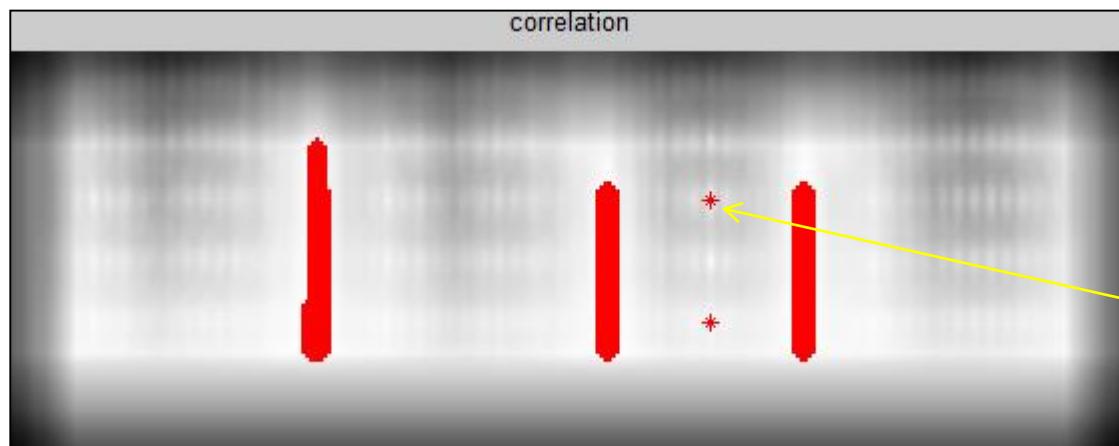
Another example



However...

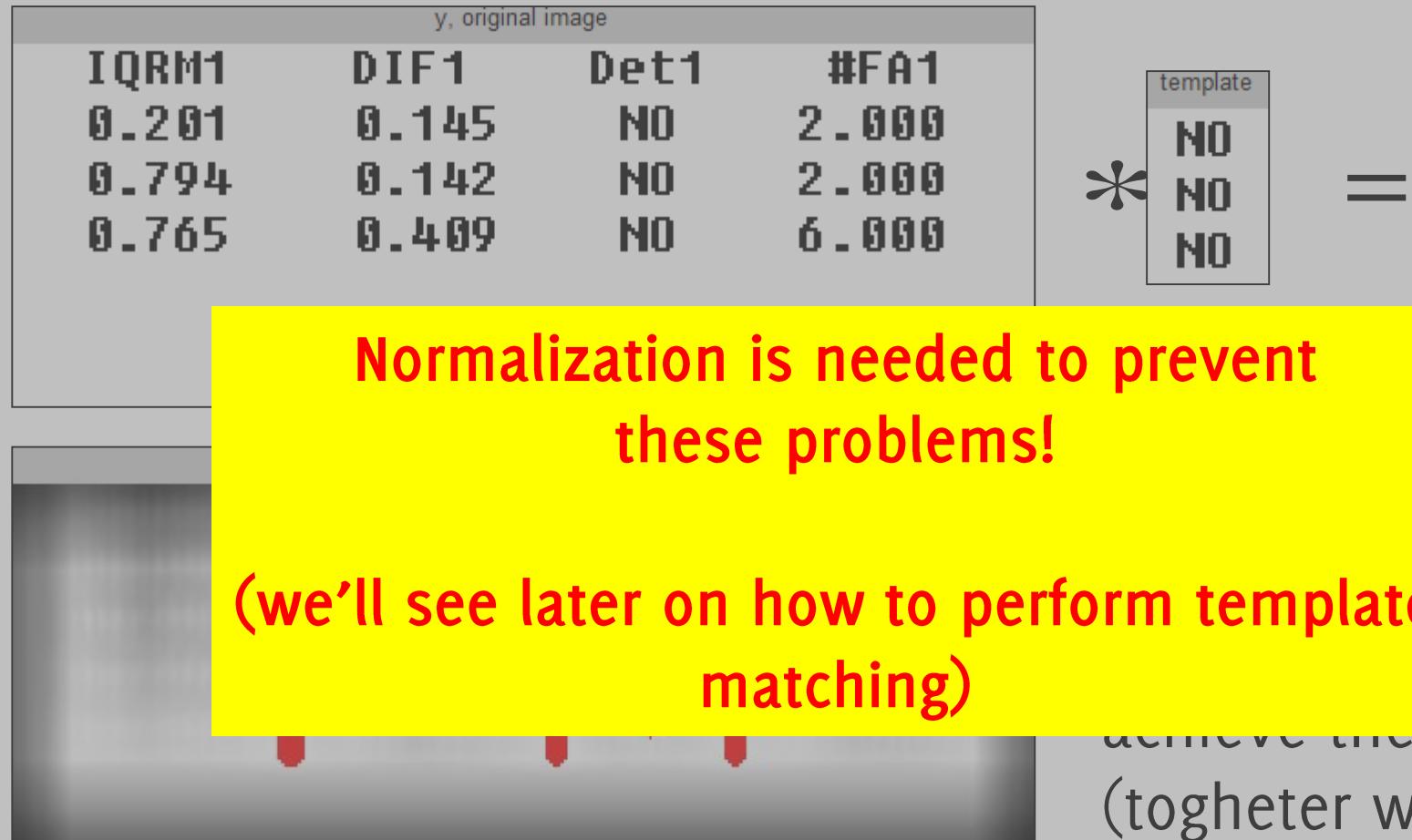
| y, original image | | | |
|-------------------|-------|------|-------|
| IQRM1 | DIF1 | Det1 | #FA1 |
| 0.201 | 0.145 | NO | 2.000 |
| 0.794 | 0.142 | NO | 2.000 |
| 0.765 | 0.409 | NO | 6.000 |

$$\begin{matrix} * \\ \text{template} \\ \text{NO} \\ \text{NO} \\ \text{NO} \end{matrix} =$$



Each point in a white area is as big as the template achieve the maximum value (together with the perfect match)

However...



a white area is template maximum value
(together with the perfect match)

Setting up the stage..

Image Classification



I

$$\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$$

“wheel”



I

“castle”

Image Classification


$$\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$$

→ “wheel” 65%, “tyre” 30%..



→ “castle” 55%, “tower” 43%..

Image Classification, the problem

Assign to an input image $I \in \mathbb{R}^{R \times C \times 3}$:

- a label l from a fixed set of categories
 $\Lambda = \{"wheel", "cars", ..., "castle", "baboon"\}$

$$I \rightarrow l \in \Lambda$$

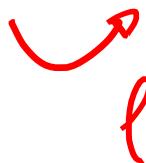


Image Classification Example

Sat, Apr 6



Thu, Apr 4



Sat, Apr 9, 2016



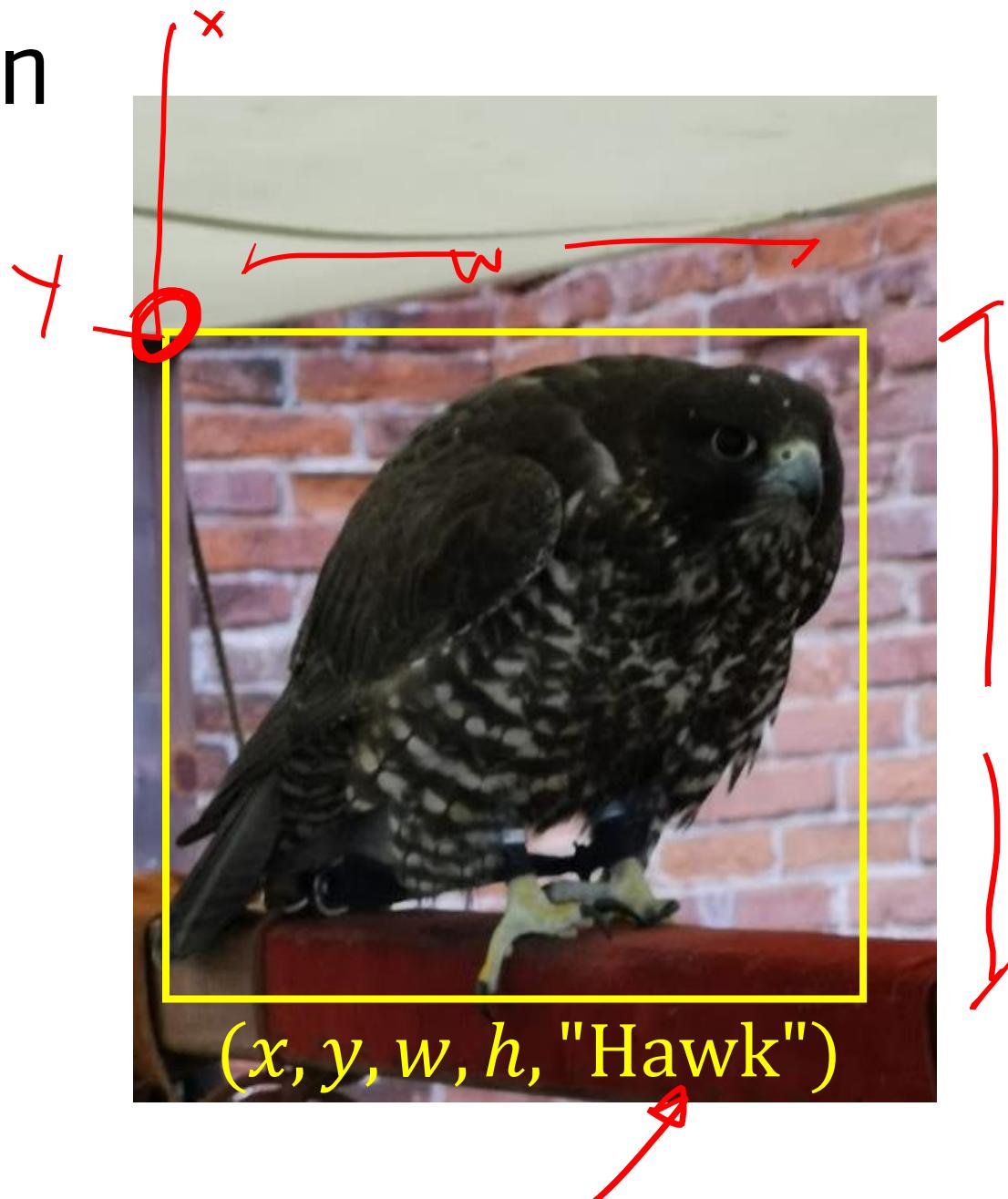
Inbox (39) - giacomo79@gmail.c rabbit - Google Photos

photos.google.com/search/rabbit

← +

Search bar: rabbit

Localization



Localization, the problem

Assign to an input image $I \in \mathbb{R}^{R \times C \times 3}$:

- a label l from a fixed set of categories
 $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$
- the coordinates (x, y, h, w) of the bounding box enclosing that object

$$I \rightarrow (x, y, h, w, l)$$

Object Detection



Object Detection, the problem

Assign to an input image $I \in \mathbb{R}^{R \times C \times 3}$:

- **multiple labels $\{l_i\}$ from a fixed set of categories $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$, each corresponding to an instance of that object**
- the coordinates $\{(x, y, h, w)_i\}$ of the bounding box enclosing each object

$$I \rightarrow \{(x, y, h, w, l)_1, \dots, (x, y, h, w, l)_N\}$$

Segmentation

Objects appearing in the image:

Boat

Dining table

Person



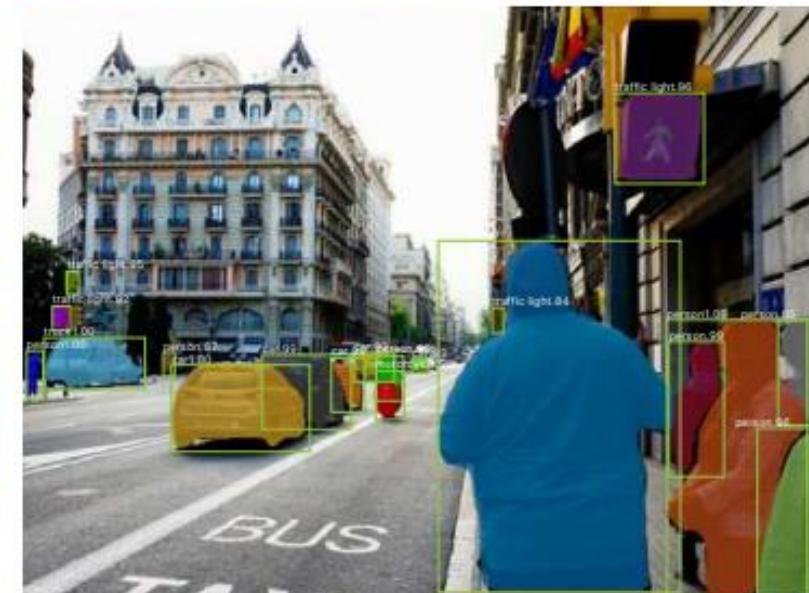
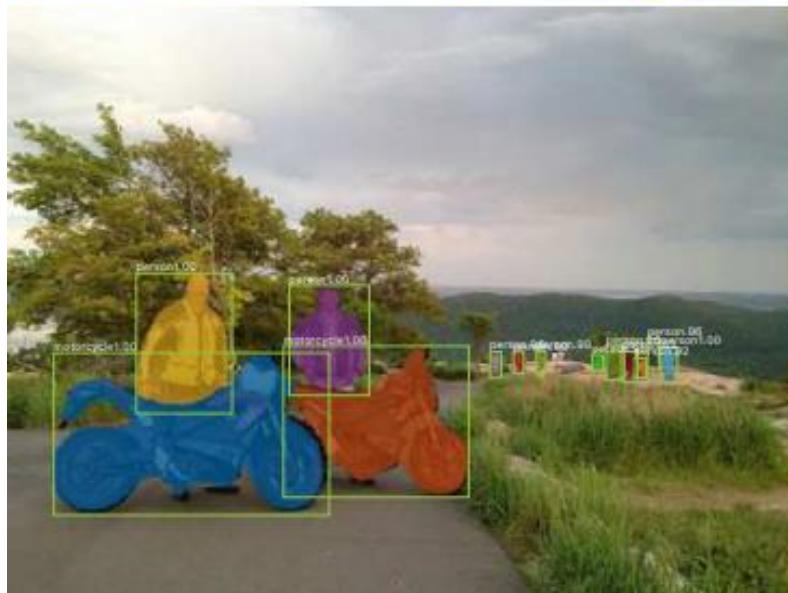
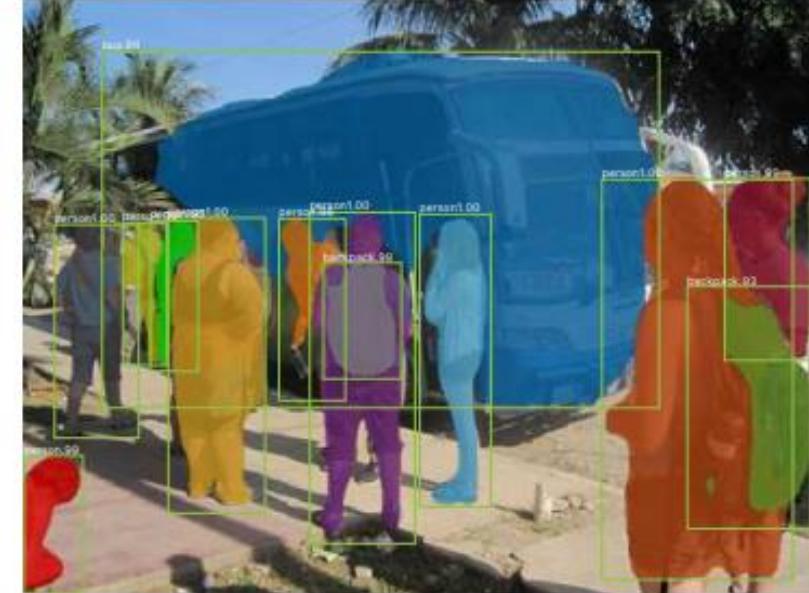
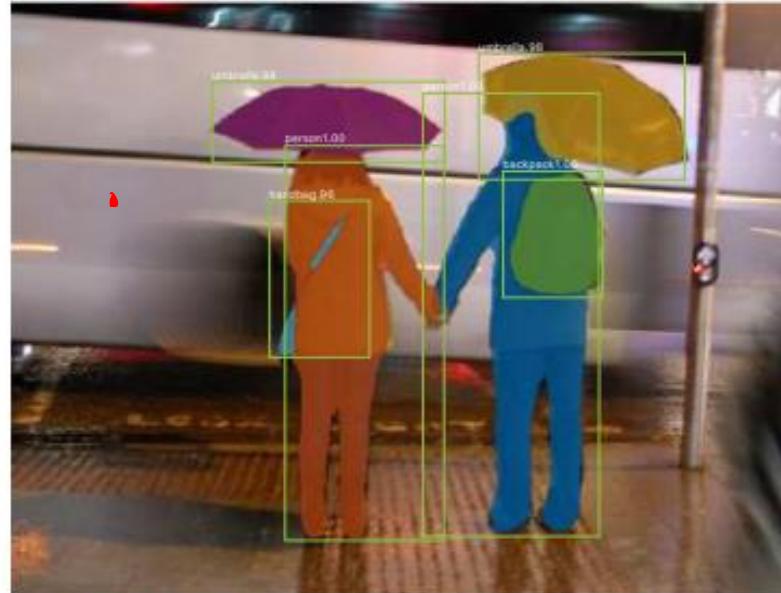
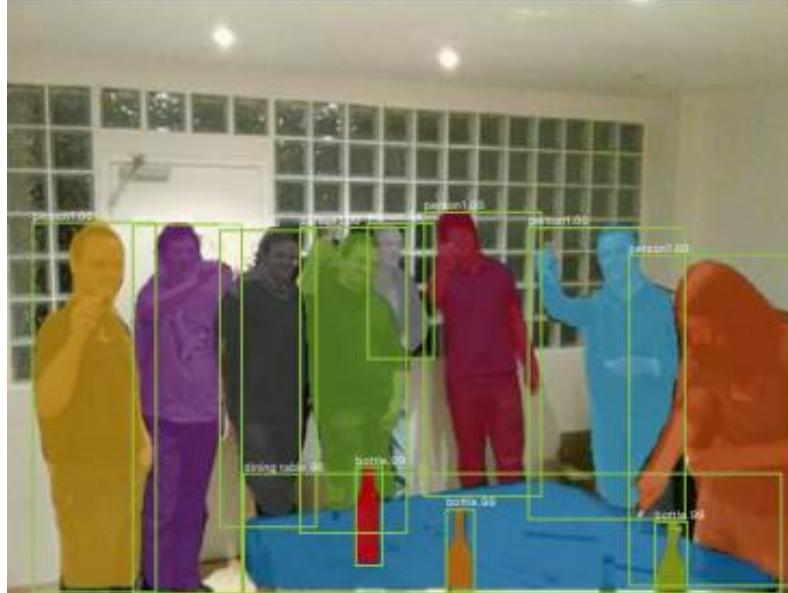
Image Segmentation, the problem

Assign to each pixel of an image $I \in \mathbb{R}^{R \times C \times 3}$:

- a label $\{l_i\}$ from a fixed set of categories
 $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\},$
 $I \rightarrow \underline{S} \in \Lambda^{R \times C}$

where $S(x, y) \in \Lambda$ denotes the class associated to the pixel (x, y)

Instance Segmentation: Mask R-CNN



Instance Segmentation, the problem

Assign to an input image I :

- **multiple labels** $\{l_i\}$ from a fixed set of categories $\Lambda = \{"wheel", "cars", \dots, "castle", "baboon"\}$, each corresponding to **an instance of that object**
- the coordinates $\{(x, y, h, w)_i\}$ of the **bounding box** enclosing each **object**
- the **set of pixels** S in each bounding box corresponding to that **label**

$$I \rightarrow \{(x, y, h, w, l, S)_1, \dots, (x, y, h, w, l, S)_N\}$$

Instance Segmentation



Is Image Classification a Challenging Problem?

Yes, it is...

First challenge: dimensionality

Images are very high-dimensional image data

CIFAR-10 dataset

The CIFAR-10 dataset contains 60000 images:

Each image is 32x32 RGB
Images are in 10 classes
6000 images per class

Extremely small images, but high-dimensional:

$$d = 32 \times 32 \times 3 = 3072$$

airplane



automobile



bird



cat



deer



dog



frog



horse



ship

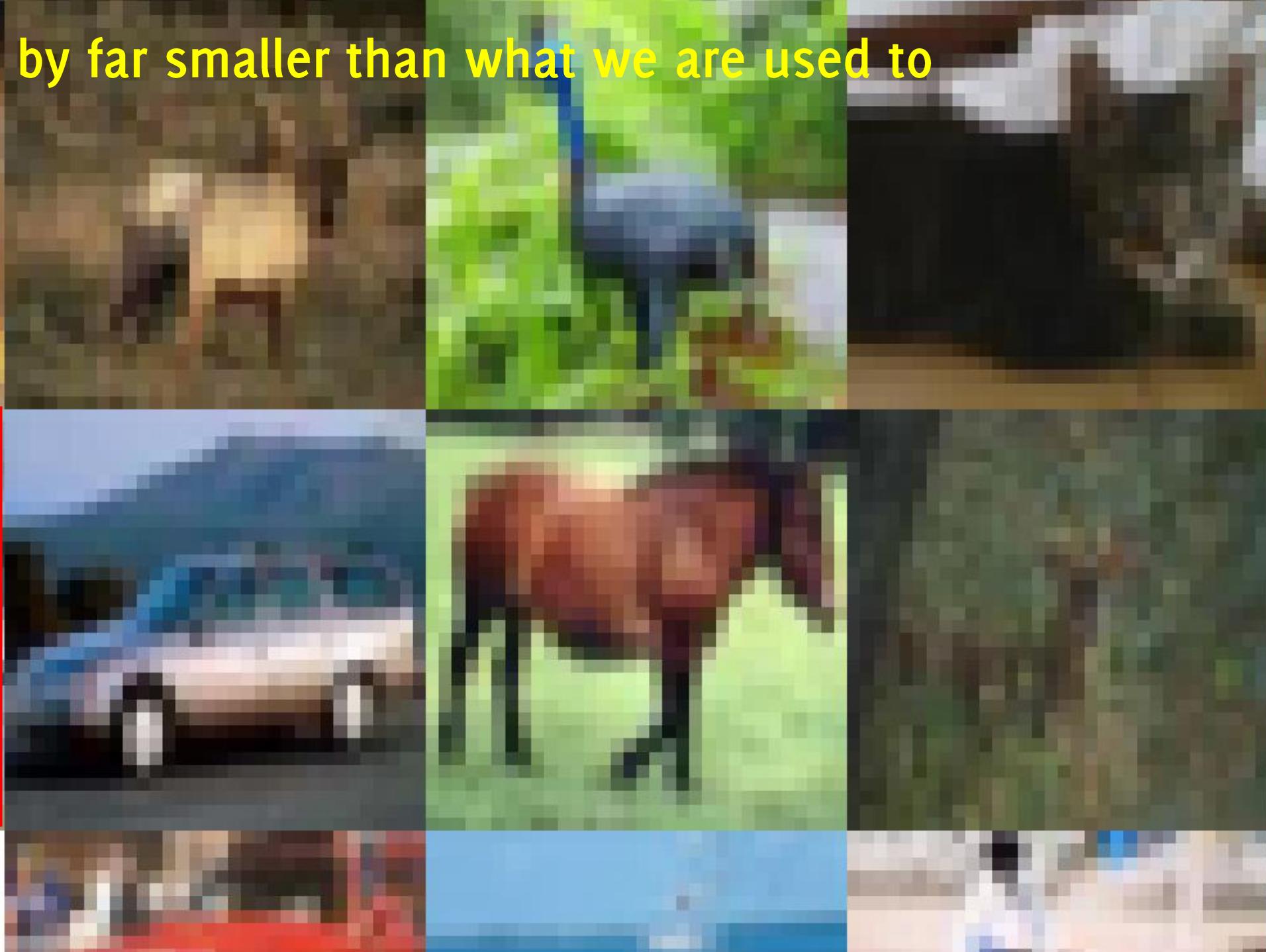


truck

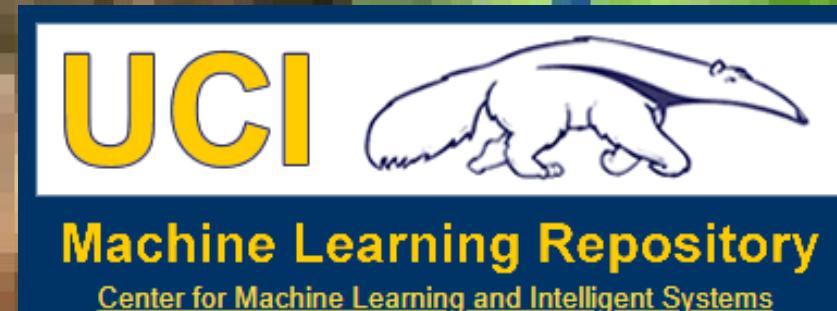


This resolution is by far smaller than what we are used to

$$d = 3072$$



Former standard repository for ML research



$d = 3072$

Attributes

Less than 10 (116)
10 to 100 (218)
Greater than 100 (86)

- 88% < 500 attributes
- 92% < 3.2K attributes

$d = 3072$

Former standard repository for ML research

Bear in mind how large an image is (in term of Bytes) when you'll be implementing your CNN... the whole batch and the corresponding activations have to be stored in memory!

2K attributes

Second challenge: label ambiguity

A label might not uniquely identify the image

Second problem: label ambiguity

Man?
Beer?
Dinner?
Restaurant?
Sausages?

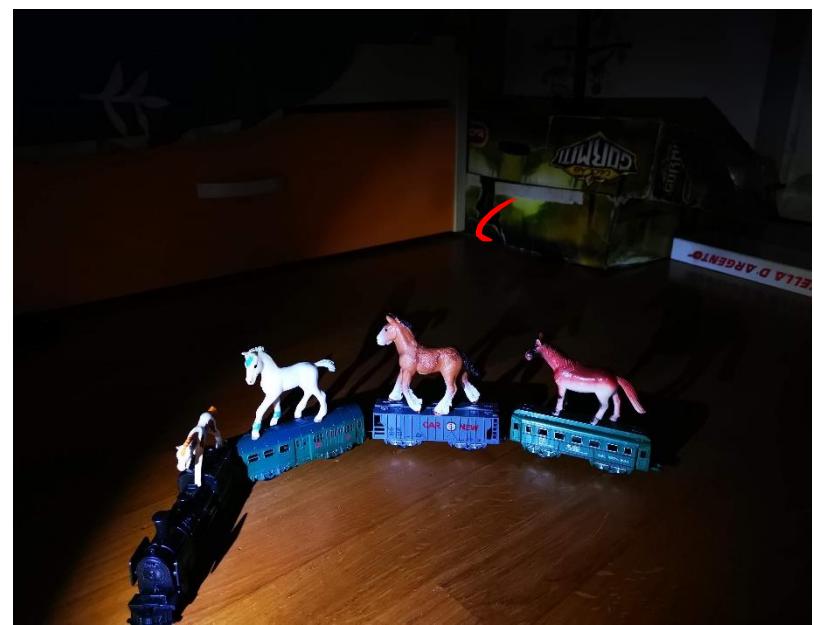
....



Third challenge: transformations

There are many transformations that change the image dramatically, while not its label

Illumination conditions changes



Giacomo Boracchi

Deformations



Copyright Christine Matthews



© Copyright Patrick Roper

Giacomo Boracchi

View Point Change



... and many others

Occlusion



Background clutter



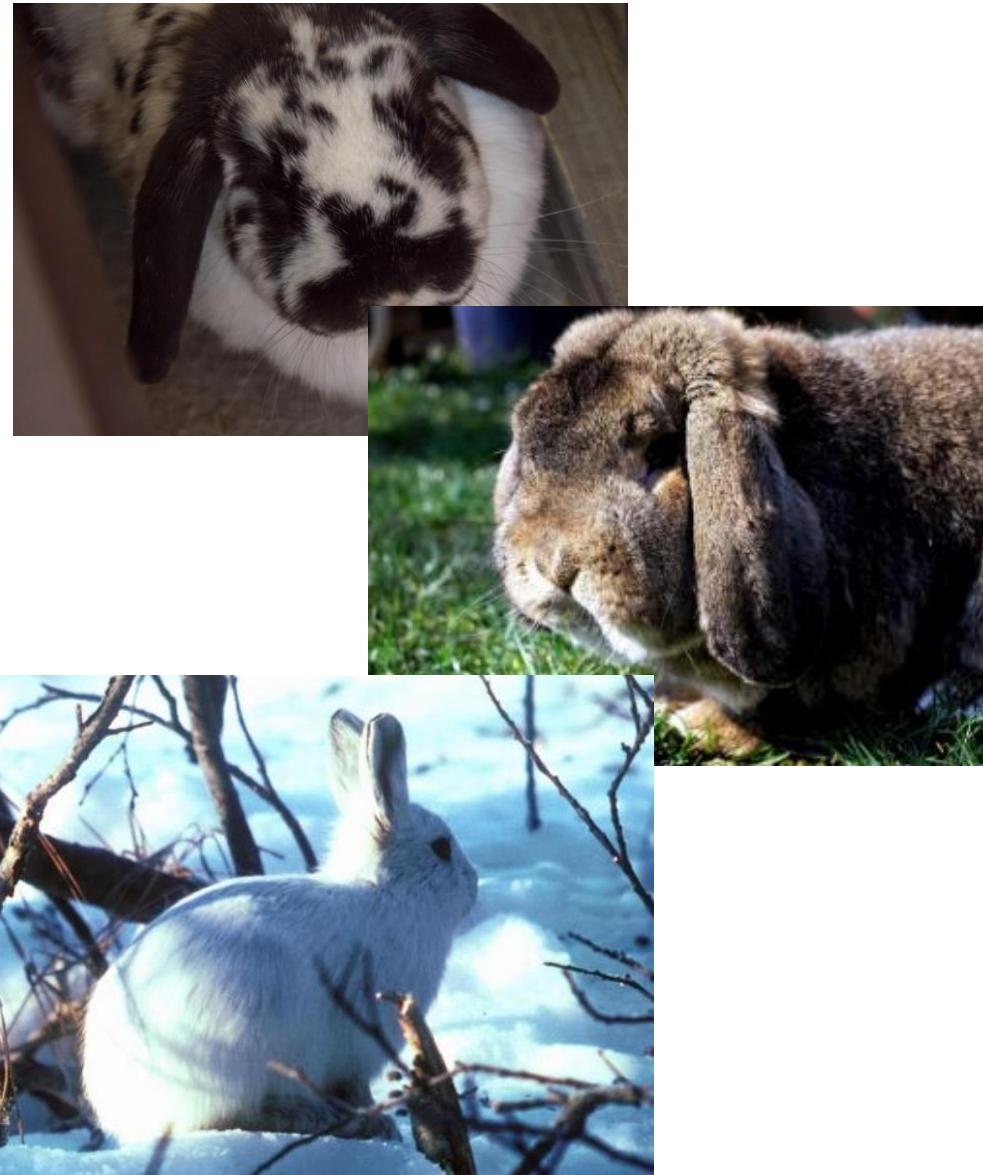
Scale variation



Fourth challenge: inter-class variability

Images in the same class might be
dramatically different

Inter-class variability



Fifth problem: perceptual similarity

Perceptual similarity in images is not
related to pixel-similarity

Nearest Neighborhood Classifiers for Images

Assign to a each test image, the label of the closest image in the training set

$$\hat{y}_j = y_{j^*}, \quad \text{being } j^* = \underset{i=1 \dots N}{\operatorname{argmin}} d(x_j, x_i)$$

Distances are typically measured as

$$d(x_j, x_i) = \|x_j - x_i\|_2 = \sqrt{\sum_k ([x_j]_k - [x_i]_k)^2}$$

Pixel-wise distance among images

↓

| | | | |
|------------|----|-----|-----|
| test image | | | |
| 56 | 32 | 10 | 18 |
| 90 | 23 | 128 | 133 |
| 24 | 26 | 178 | 200 |
| 2 | 0 | 255 | 220 |

$$-$$

| | | | |
|----------------|----|-----|-----|
| training image | | | |
| 10 | 20 | 24 | 17 |
| 8 | 10 | 89 | 100 |
| 12 | 16 | 178 | 170 |
| 4 | 32 | 233 | 112 |

$$=$$

| | | | |
|---------------------------------------|----|----|-----|
| pixel-wise absolute value differences | | | |
| 46 | 12 | 14 | 1 |
| 82 | 13 | 39 | 33 |
| 12 | 10 | 0 | 30 |
| 2 | 32 | 22 | 108 |

→ 456
Σ

K-Nearest Neighborhood Classifiers for Images

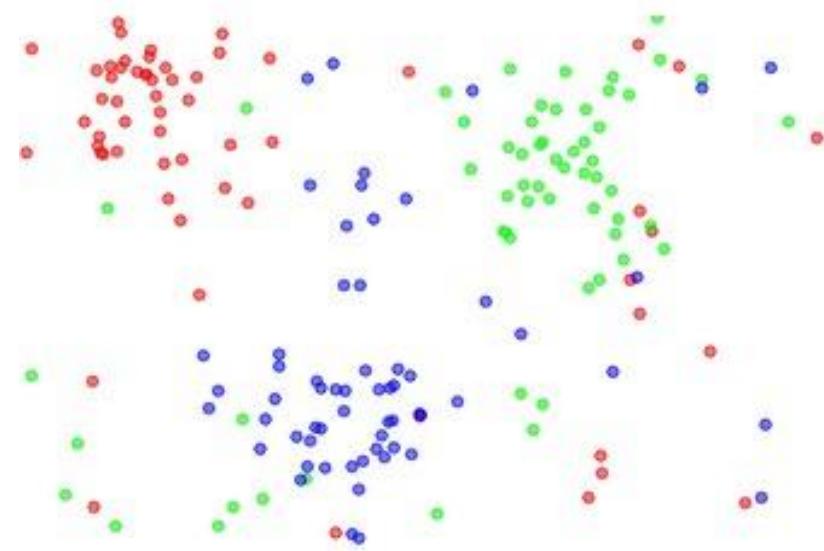
Assign to a each test image, the most frequent label among the K –closest images in the training set

$$\hat{y}_j = y_{j^*}, \quad \text{being } j^* \text{ the mode of } \mathcal{U}_K(x_j)$$

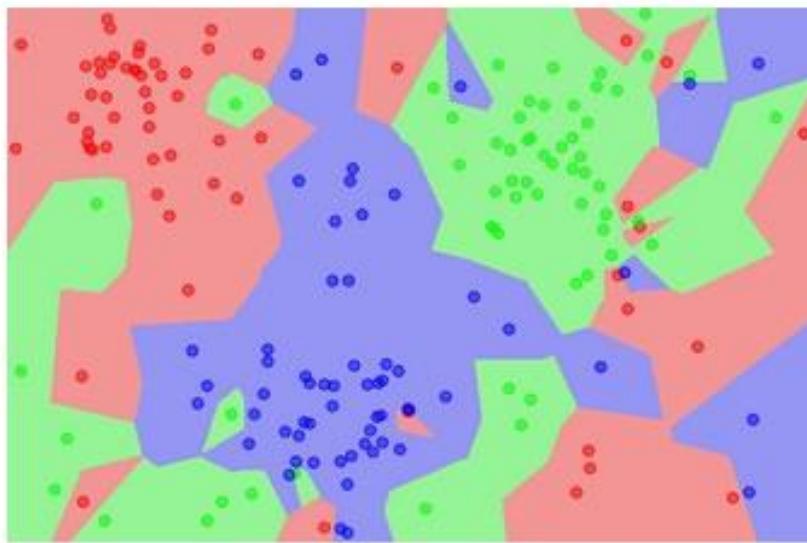
where $\mathcal{U}_K(x_j)$ contains the K closest training images to x_j

Choosing the right parameter K and the distance measure is often an issue

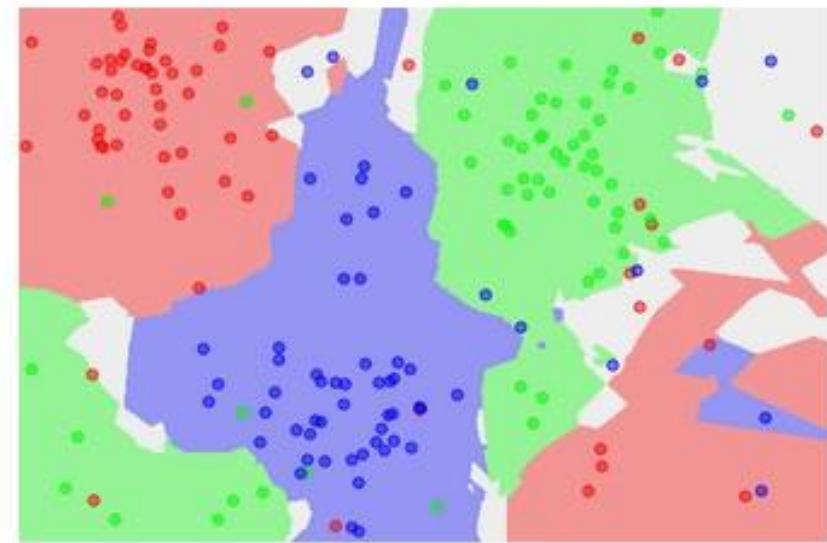
the data



NN classifier

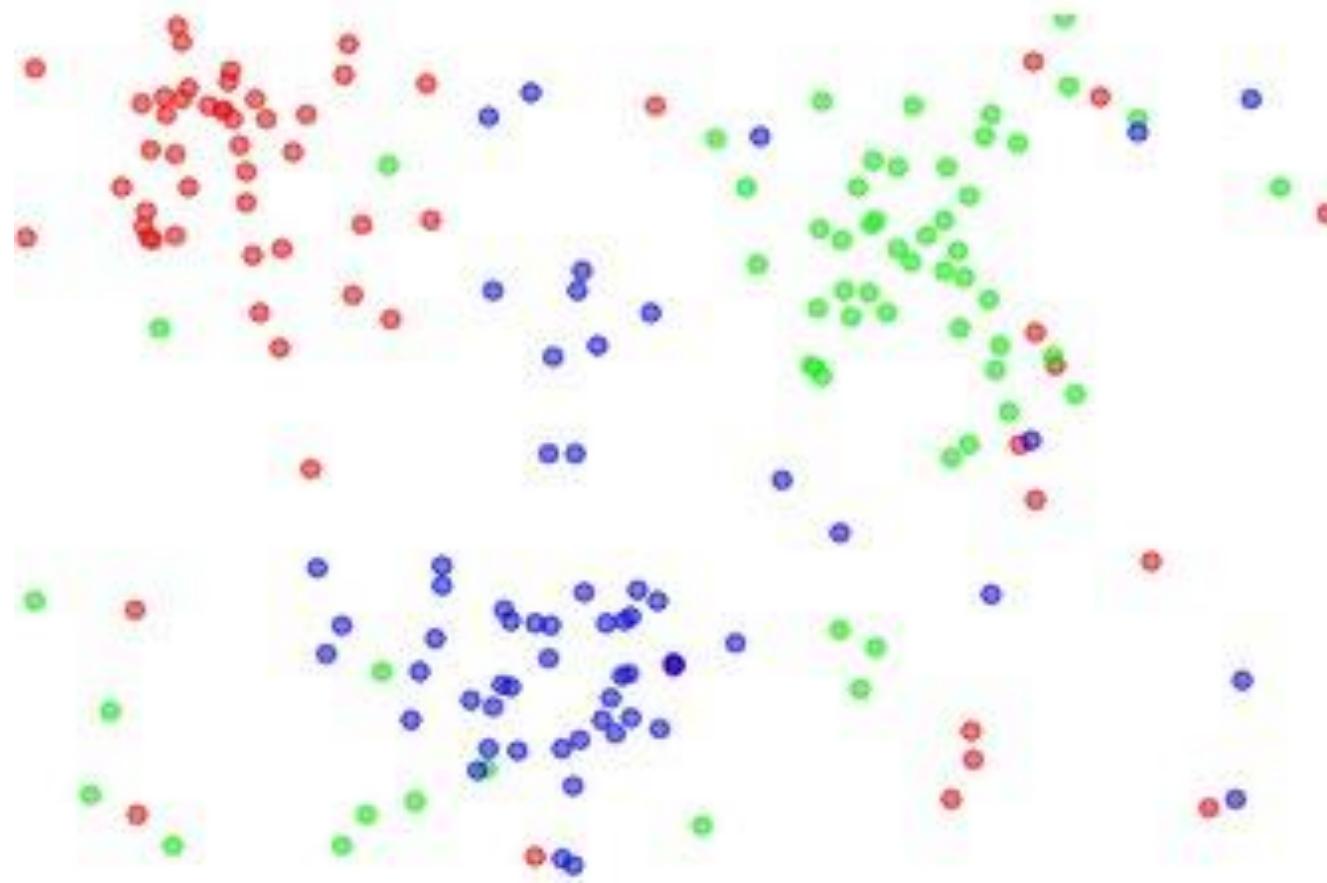


5-NN classifier



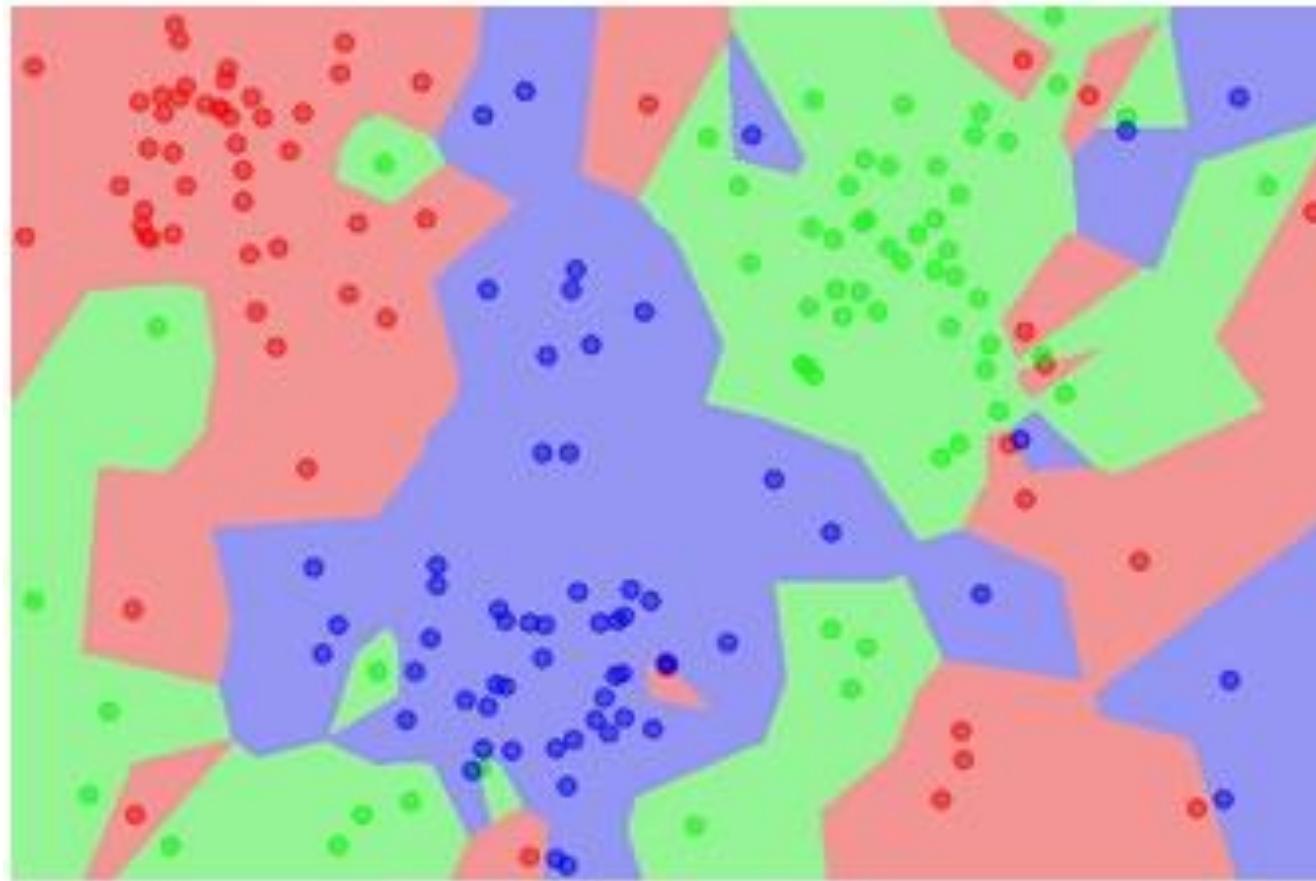
Nearest Neighborhood Classifiers for Images

the data



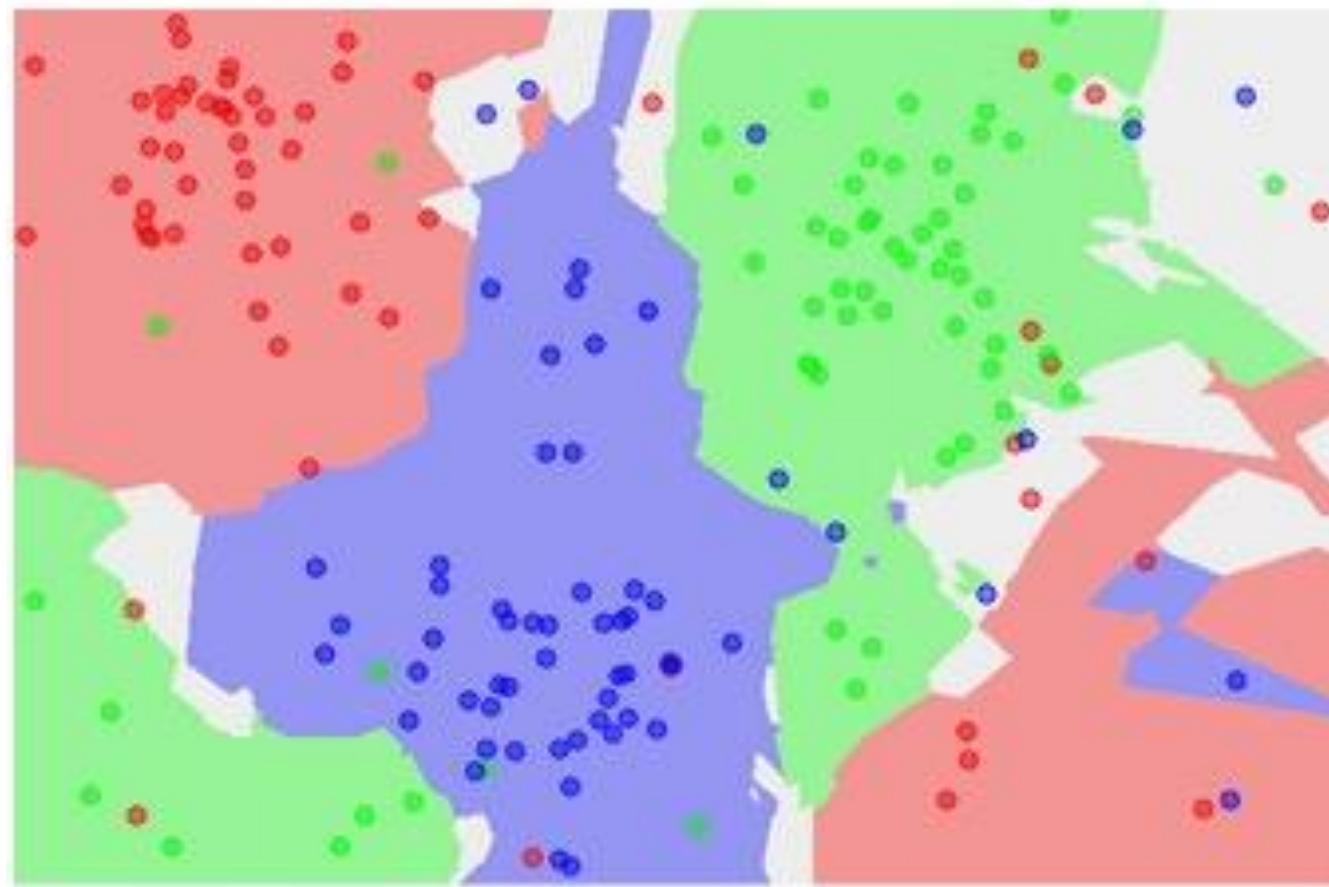
Nearest Neighborhood Classifiers for Images

NN classifier



Nearest Neighborhood Classifiers for Images

5-NN classifier



Nearest Neighborhood Classifiers for Images

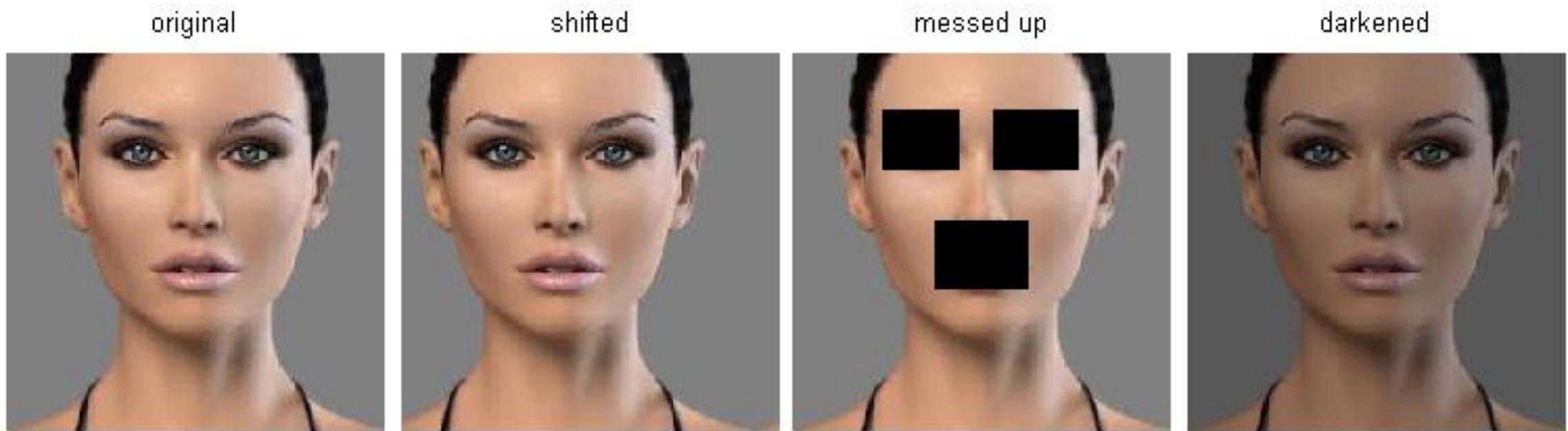
Pros:

- Easy to understand and implement
- It takes no training time

Cons:

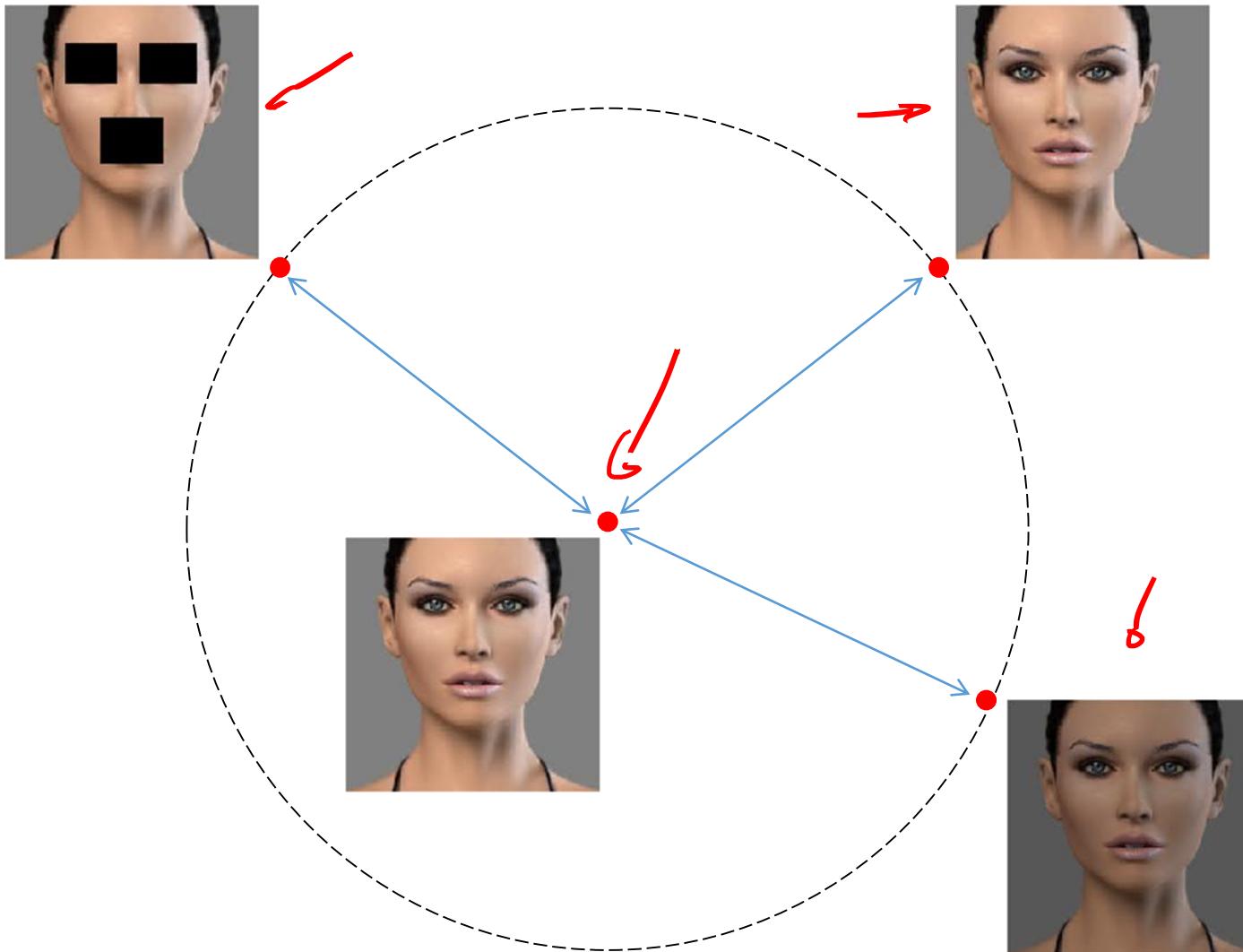
- Computationally demanding at test time (in particular when TR is large and d is also large)
- Large training sets have to be stored in memory
- Rarely practical on images: distances on high-dimensional objects are difficult to interpret

Perceptual Similarity vs Pixel Similarity



The three images have the same pixel-wise distance from the original one...
...but perceptually they are very different

Perceptual Similarity vs Pixel Similarity



Let's see what happens on the whole CIFAR10



On CIFAR10 we see exactly this problem



On CIFAR10 we see exactly this problem



Linear Classifier

the basic building block for deep architectures

Linear Classifiers for Images

A classifier can be seen as a function that maps an image \mathbf{x} to a confidence scores for each of the L class:

$$\mathcal{K}: \mathbb{R}^d \rightarrow \mathbb{R}^L$$

where $\mathcal{K}(\mathbf{x})$ is a L – dimensional vector and the i – th component

$$s_i = [\mathcal{K}(\mathbf{x})]_i$$

contains a score of how likely \mathbf{x} belongs to class i .

Intuitively, a good classifier associates to the correct class a score that is larger than scores associated to incorrect classes.

Linear Classifiers for Images

In linear classification \mathcal{K} is a linear function:

$$\mathcal{K}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$$

where $\mathbf{W} \in \mathbb{R}^{L \times d}$, $\mathbf{b} \in \mathbb{R}^L$ are the parameters of the classifier \mathcal{K} .

\mathbf{W} are referred to as the weights, \mathbf{b} the bias.

| | | | | | | | | |
|------|-----|-----|------|-----|------|------|-----|------|
| -8.1 | ... | 2.7 | 9.5 | ... | -9.0 | -5.4 | ... | 4.8 |
| 9.0 | ... | 5.4 | 4.8 | ... | 1.2 | 9.5 | ... | -8.0 |
| 1.2 | ... | 9.5 | -8.0 | ... | 8.1 | -2.7 | ... | 9.5 |

\mathbf{W}

$$\begin{array}{c} * \\ \mathbf{b} \end{array} \quad \begin{array}{c} 23 \\ \dots \\ 21 \\ 34 \\ \dots \\ 12 \\ 34 \\ \dots \\ 23 \end{array} \quad + \quad \begin{array}{c} -2 \\ 32 \\ -1 \end{array} \quad = \quad \begin{array}{c} -4 \\ 22 \\ 33 \end{array} \quad \begin{array}{l} s_1 \text{ dog score} \\ s_2 \text{ cat score} \\ s_3 \text{ rabbit score} \end{array}$$

$\mathcal{K}(\mathbf{x}_i; \mathbf{W}, \mathbf{b})$

Colors recall the color plane
where images are from

Linear Classifiers for Images



Unroll the image column-wise

| | | | | | | | | |
|------|-----|-----|------|-----|------|------|-----|------|
| -8.1 | ... | 2.7 | 9.5 | ... | -9.0 | -5.4 | ... | 4.8 |
| 9.0 | ... | 5.4 | 4.8 | ... | 1.2 | 9.5 | ... | -8.0 |
| 1.2 | ... | 9.5 | -8.0 | ... | 8.1 | -2.7 | ... | 9.5 |

W

*

| | | | | | | | | |
|----|-----|----|----|-----|----|----|-----|----|
| 23 | ... | 21 | 34 | ... | 12 | 34 | ... | 23 |
|----|-----|----|----|-----|----|----|-----|----|

| | | |
|----|----|----|
| -2 | 32 | -1 |
|----|----|----|

| | | |
|----|----|----|
| -4 | 22 | 33 |
|----|----|----|

s_1 dog score
 s_2 cat score
 s_3 rabbit score

b

$\mathcal{K}(x_i; W, b)$

x_i

Linear Classifiers for Images

The classifier assign to an input image the class corresponding to the largest score

$$\hat{y}_j = \operatorname{argmax}_{i=1,\dots,L} [s_j]_i$$

being $[s_j]_i$ the i -th component of the vector

$$\mathcal{K}(x_j) = Wx_j + b$$

$$\begin{array}{|c|c|c|c|c|c|c|c|c|c|} \hline -8.1 & \dots & 2.7 & 9.5 & \dots & -9.0 & -5.4 & \dots & 4.8 \\ \hline 9.0 & \dots & 5.4 & 4.8 & \dots & 1.2 & 9.5 & \dots & -8.0 \\ \hline 1.2 & \dots & 9.5 & -8.0 & \dots & 8.1 & -2.7 & \dots & 9.5 \\ \hline \end{array} * \begin{array}{|c|} \hline 23 \\ \hline \dots \\ \hline 21 \\ \hline 34 \\ \hline \dots \\ \hline 12 \\ \hline 34 \\ \hline \end{array} + \begin{array}{|c|} \hline -2 \\ \hline 32 \\ \hline -1 \\ \hline \end{array} = \begin{array}{|c|} \hline -4 \\ \hline 22 \\ \hline 33 \\ \hline \end{array}$$

W \mathbf{x}_i b $\mathcal{K}(x_i; W, b)$

s_1 dog score
 s_2 cat score
 s_3 rabbit score

Linear Classifiers for Images

The score of a class is the weighted sum of all the image pixels.
Weights are actually the classifier parameters.

Weights indicate which are the most important pixels / colors

| | | | | | | | | |
|------|-----|-----|------|-----|------|------|-----|------|
| -8.1 | ... | 2.7 | 9.5 | ... | -9.0 | -5.4 | ... | 4.8 |
| 9.0 | ... | 5.4 | 4.8 | ... | 1.2 | 9.5 | ... | -8.0 |
| 1.2 | ... | 9.5 | -8.0 | ... | 8.1 | -2.7 | ... | 9.5 |

W

| |
|----|
| -2 |
| 32 |
| -1 |

b

Training a Classifier

Given a training set TR and a loss function, define the parameters that minimize the loss function over the whole TR

In case of linear classifier

$$[W, b] = \operatorname{argmin}_{W \in \mathbb{R}^{L \times d}, b \in \mathbb{R}^L} \sum_{(x_i, y_i) \in TR} \mathcal{L}(x, y_i)$$

Solving this minimization problem provides the weights of our classifier

Loss Function

Loss function: a function \mathcal{L} that measures our unhappiness with the score assigned to training images

The loss \mathcal{L} will be high on a training image that is not correctly classifier, low otherwise.

Loss Function Minimization

Loss function can be minimized by specific algorithms (see Matteo' classes)

The loss function has to be typically regularized to achieve a unique solution satisfying some desired property

$$[W, b] = \operatorname{argmin}_{W \in \mathbb{R}^{L \times d}, b \in \mathbb{R}^L} \sum_{(x_i, y_i) \in TR} \mathcal{L}(x, y_i) + \lambda \mathcal{R}(W, b)$$

being $\lambda > 0$ a parameter balancing the two terms

Once Trained

The classifier is expected to provide to the correct class a score that is larger than that assigned to the incorrect classes.

The training data is used to learn the parameters W, b

Once the training is completed, it is possible to discard the whole training set and only keep the learned parameters.

Geometric Interpretation of a Linear Classifier

In Matlab notation:

$W(i, :)$ is a d –dimensional vector containing the weights of the score function for the i –th class

Computing the score function for the i –th class corresponds to computing the inner product

$$W(i, :) * \mathbf{x}$$

Thus, linear classifiers identify an hyperplane in a d –dimensional

(in Python $*$ is the element-wise product, here I mean the inner product of vectors):

Geometric Interpretation

Interpret each image as a point in \mathbb{R}^d .

Each classifier is a weighted sum of pixels, which corresponds to a linear function in \mathbb{R}^d

In \mathbb{R}^2 these would be

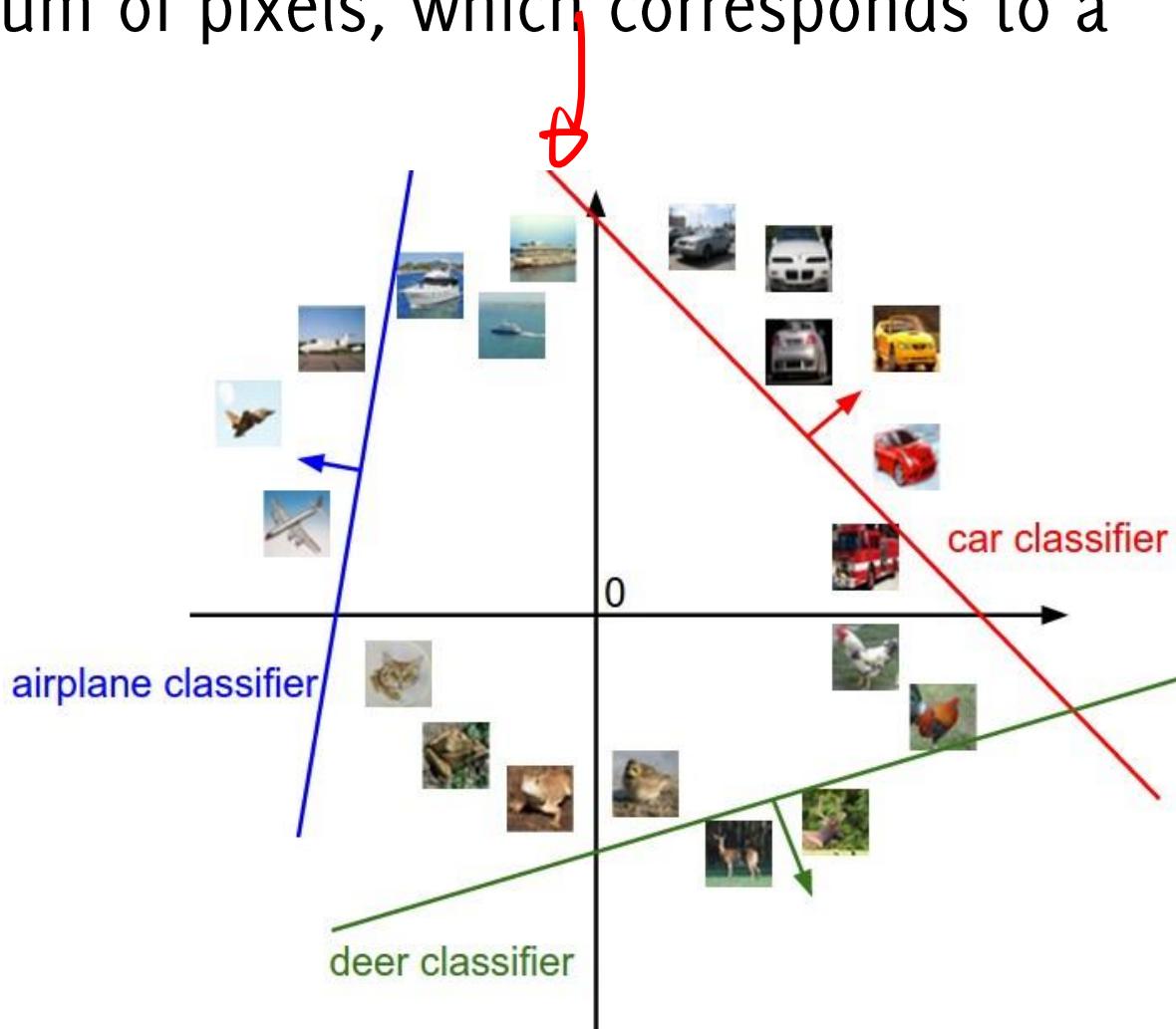
$$f([x_1, x_2]) = w_1 x_1 + w_2 x_2 + b$$

Then, points $[x_1, x_2]$ yielding

$$f([x_1, x_2]) = 0$$

would be lines.

Thus, in \mathbb{R}^2 the region that separates positive from negative scores for each class is a line. This region becomes an hyperplane in \mathbb{R}^d .



Template Matching Interpretation

Using Matlab notation:

- $W(i, :)$ is a d –dimensional vector containing the weights of the score function for the i –th class
- Computing the score function for the i –th class corresponds to computing the inner product

$$W(i, :) * \mathbf{x}$$

Then, $W(i, :)$ can be seen as a template used in matching (the output of correlation in the central pixel)

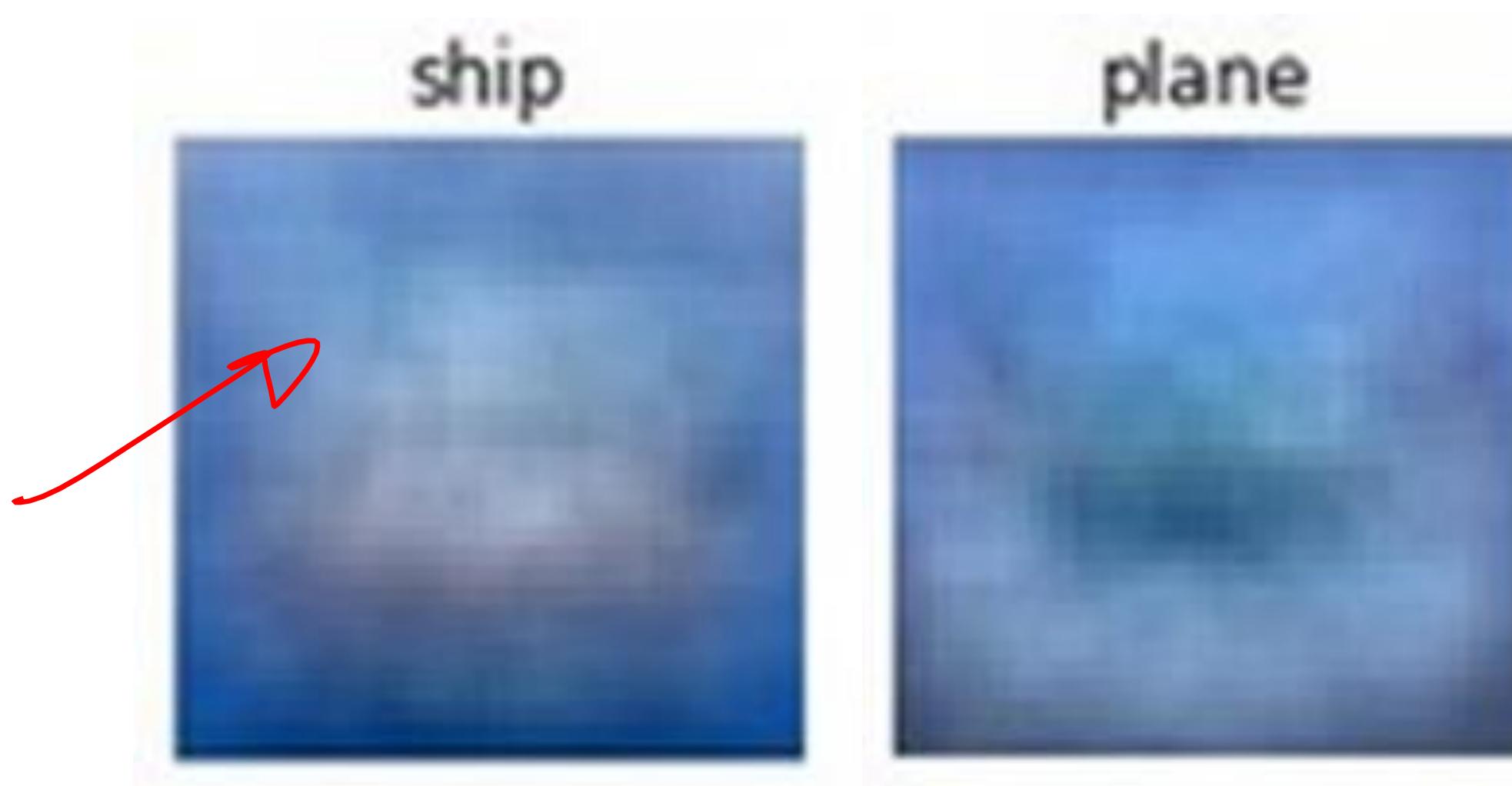
The template $W(i, :)$ is learned to match at best images belonging to the i –th class

Let's have a look at these templates

Templates Learned on the CIFAR-10 dataset



Templates Learned on the CIFAR-10 dataset



Templates Learned on the CIFAR-10 dataset

car

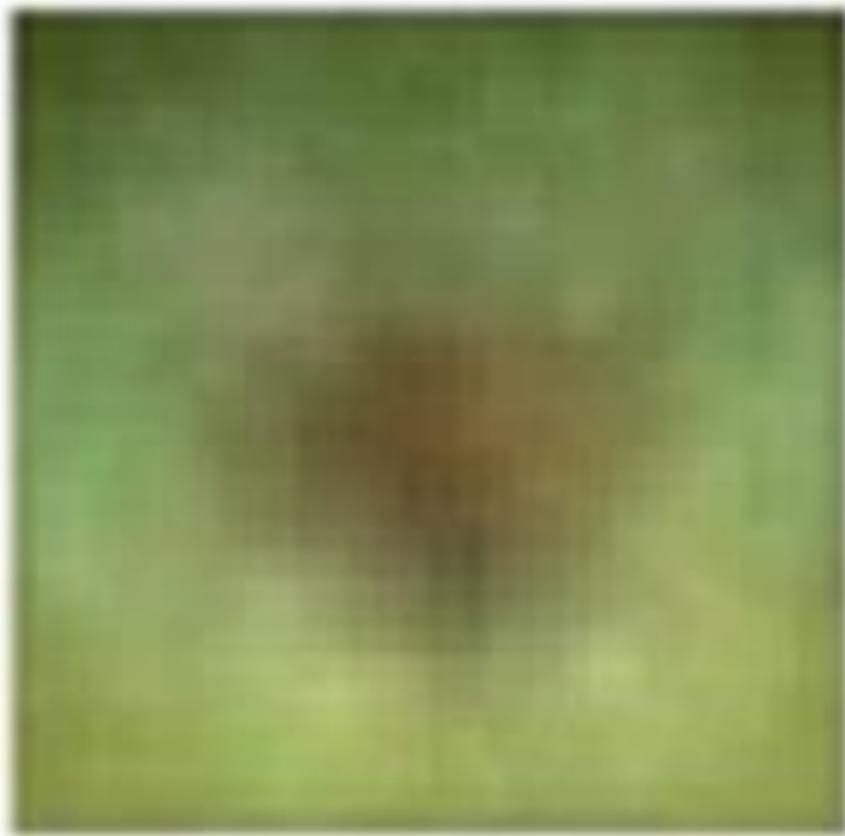


truck



Templates Learned on the CIFAR-10 dataset

deer



bird



Templates Learned on the CIFAR-10 dataset



Linear Classifier as a Template Matching

What has the classifier learned?

- That the background of bird and frog is green, (plane and boat is blue)
- Cars are typically red
- Horses have two heads! ☺

The model was definitely too simple / data were not enough for achieving higher performance and better templates

However:

- Linear Classifiers are among the most important layer of NN
- Such a simple model can be interpreted (with more sophisticated models you typically can't)

Linear Classifier as a Template Matching

What has the classifier learned?

- That the background of bird and frog is blue)
- Cars are typically red
- Horses have two heads! ☺

The model was definitely good for achieving higher performance / data were not enough for better templates

However:

- Linear Classifiers are among the most important layer of NN
- Such a simple model can be interpreted (with more sophisticated models you typically can't)

There should be a better way
for handling images