



UNIVERSITÀ DEGLI STUDI DI PADOVA

**High-level vision: object recognition,
template matching, bag of words**

Stefano Ghidoni





- Finding object in an image
- Template matching
- Histogram of oriented gradients
- Bag of words



- Consider the high-level task of *getting some information from an image about objects*
- **Object recognition** is a general term to describe a collection of related computer vision tasks that involve identifying objects in images or videos



Classification



CAT

Output:
A label

Classification



CAT

Classification + Localization



CAT

Single object

Output:
A label + a bounding box

Classification



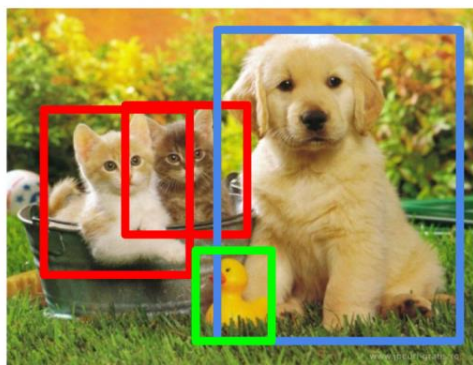
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Single object

Output:
Multiple bounding boxes with label

Classification



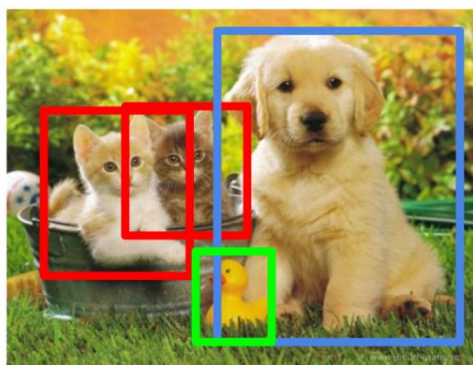
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



CAT, DOG, DUCK

Single object

Multiple objects

Output:
Multiple areas with label



- The tasks discussed so far shall cope with
 - Different camera positions
 - Perspective deformations
 - Illumination changes
 - Intra-class variations



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Camera position

IAS-LAB





UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Perspective deformation

IAS-LAB





UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Illumination changes

IAS-LAB

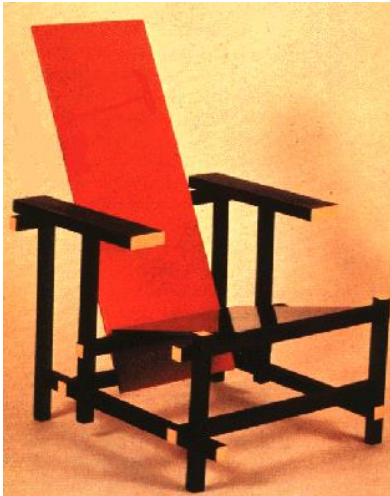




UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Intra-class variations

IAS-LAB





- We already covered a method for object detection
 - Which one?



- Anti spoiler 😊



- We already covered a method for object detection
 - Which one?
- Boosting for face detection (Viola and Jones)
 - Can be applied to other targets

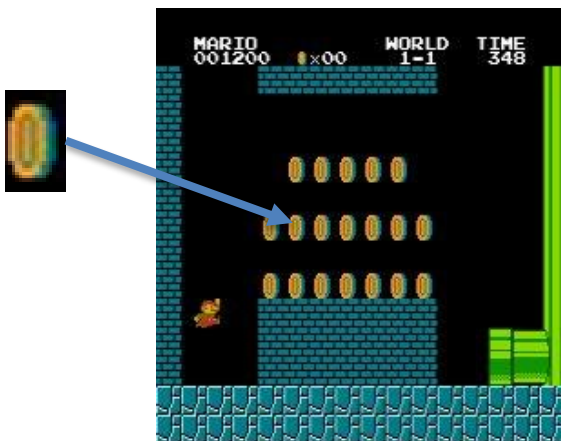


- We already covered a method for object detection
 - Which one?
- Other approaches are available
 - Template matching
 - Histogram of Oriented Gradients (HOG)
 - Bag of Words
 - Machine learning / deep learning

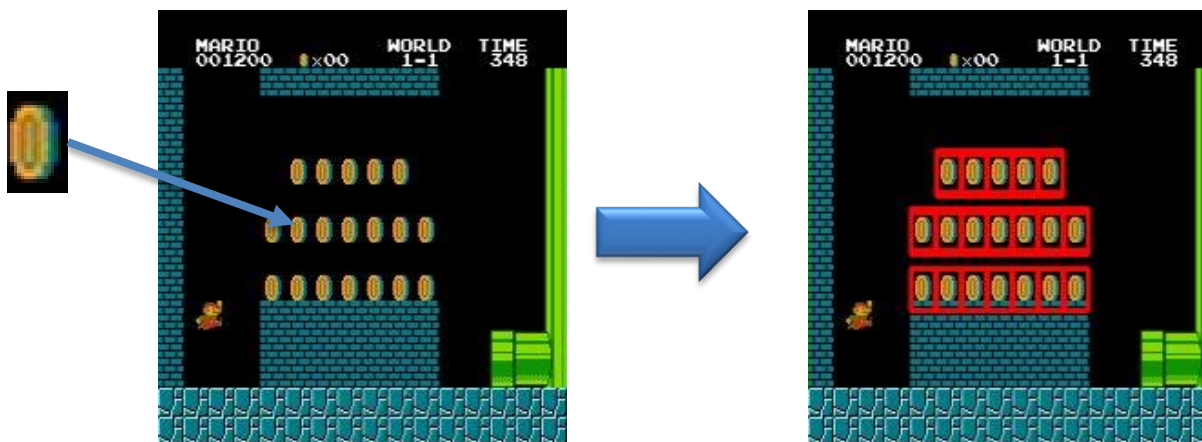


Template matching

- A template is
 - Something fashioned, shaped or designed to serve as a model
 - Something formed after a model
 - A representative instance (an example)

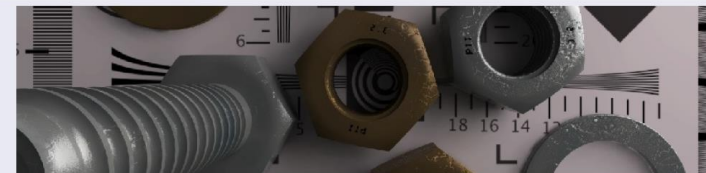


- "Where is a given object?"
- Find instances of the templates in the image
- A similarity measure should be chosen

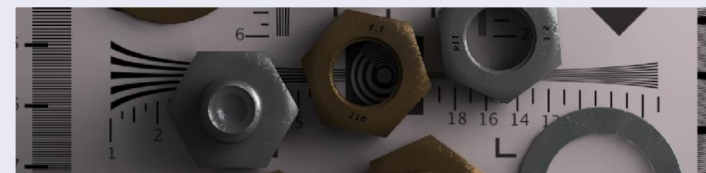


- Template variability, deformable objects
- Imaging device properties
- Viewpoint changes
- Affine transforms – scale, rotations, translations, ...
- Noise
- Illumination

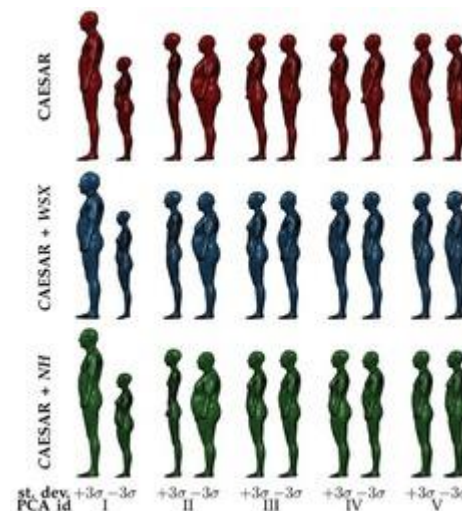
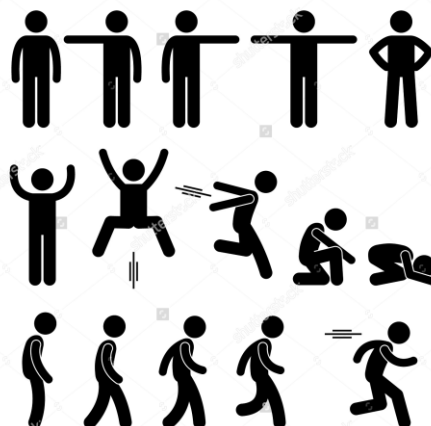
Perspective camera



Telecentric camera



a a a a





- Given
 - An image
 - A template
- How can we evaluate the match?

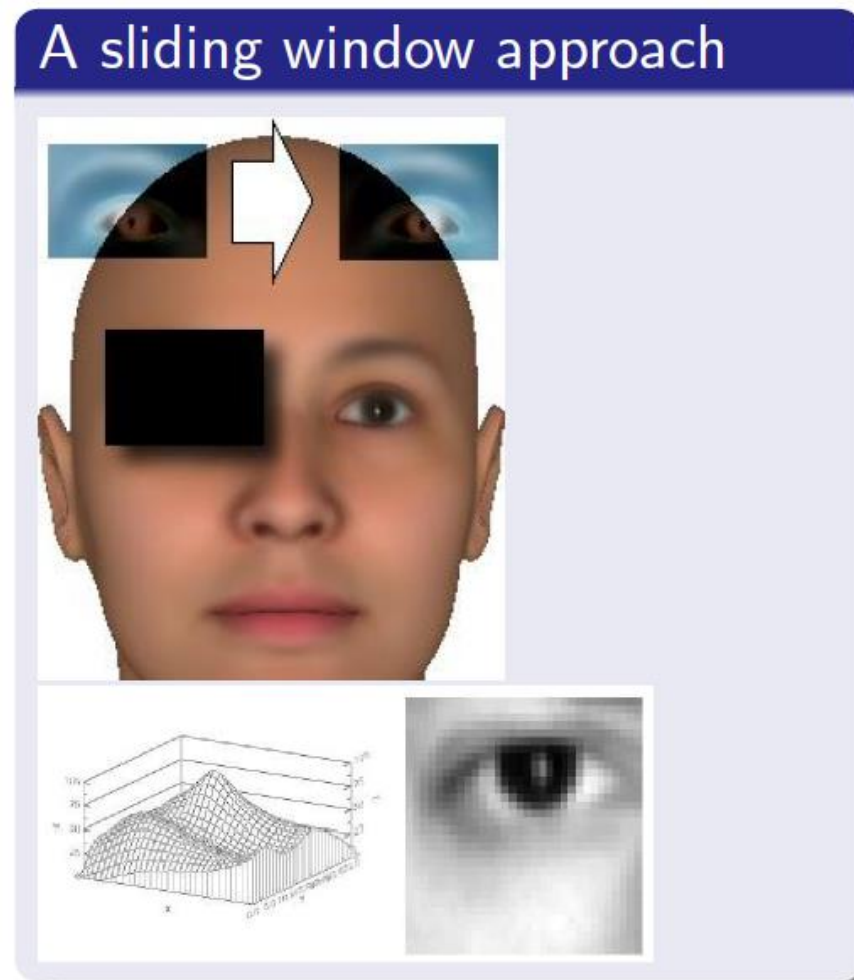


- Anti spoiler



- Common option: correlation-based approach
- Template T: rigid object, often a small image
- Sliding window
- Comparison of
 - Pixel values
 - Features
 - Edges or gradient orientation
- Similarity metrics: SSD, SAD, ZNCC

- Template-based approaches are often based on a sliding window
- The template is placed in every possible position across the image
 - Basic approach w/out rotation and scaling



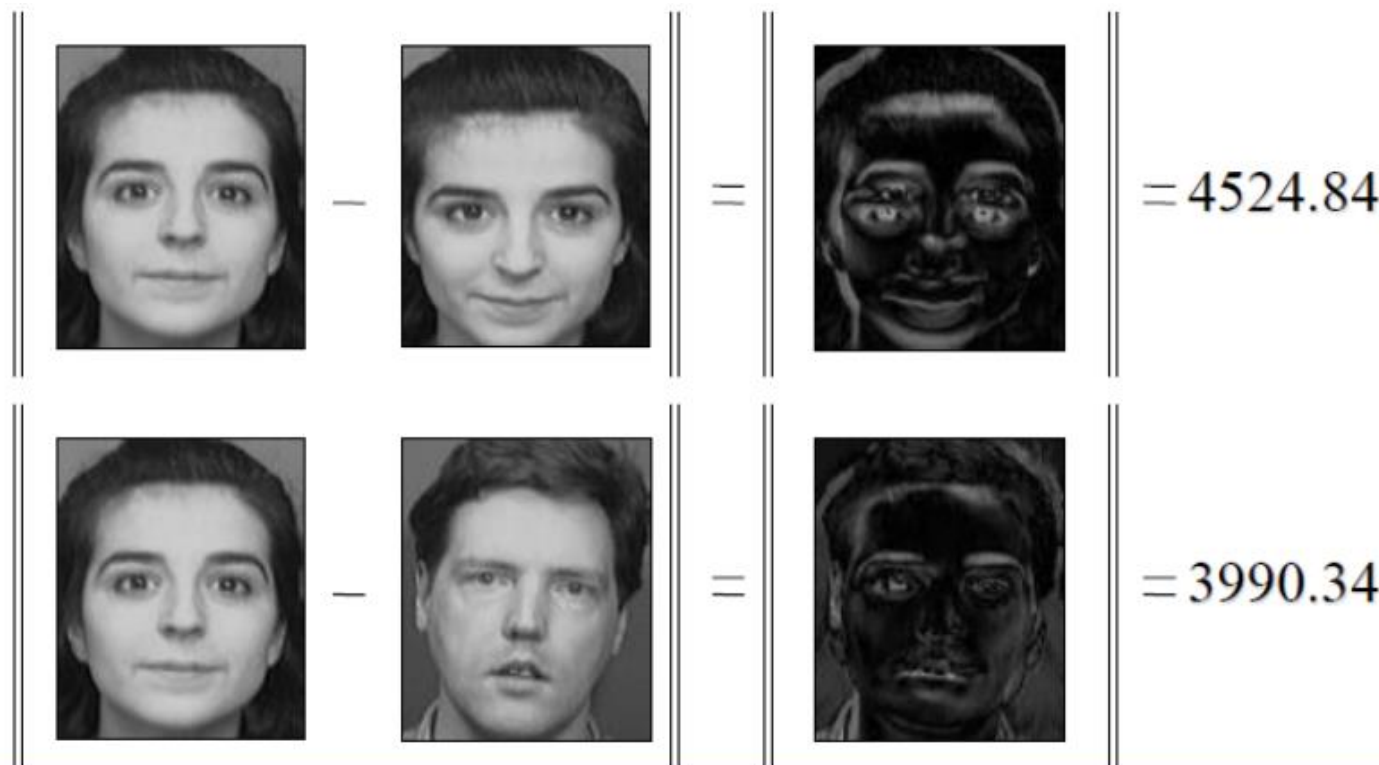


- Common option: correlation-based approach
- Template T: rigid object, often a small image
- Sliding window
- Comparison of
 - Pixel values
 - Features
 - Edges or gradient orientation
- Similarity metrics: SSD, SAD, ZNCC



- Common option: correlation-based approach
- Template T: rigid object, often a small image
- Sliding window
- Comparison of
 - Pixel values
 - Features
 - Edges or gradient orientation
- Similarity metrics: SSD, SAD, ZNCC

- Simple differencing does not always provide reliable results – **why in this case?**





- Common option: correlation-based approach
- Template T: rigid object, often a small image
- Sliding window
- Comparison of
 - Pixel values
 - Features
 - Edges or gradient orientation
- Similarity metrics: SSD, SAD, ZNCC



- Common option: correlation-based approach
- Template T: rigid object, often a small image
- Sliding window
- Comparison of
 - Pixel values
 - Features
 - Edges or gradient orientation
- Similarity metrics: SSD, SAD, ZNCC

- Sum of Squared Differences (SSD)

$$\phi(x, y) = \sum_{u, v \in T} (I(x + u, y + v) - T(u, v))^2$$

- Sum of Absolute Differences (SAD)

$$\phi(x, y) = \sum_{u, v \in T} |I(x + u, y + v) - T(u, v)|$$

- Zero-mean Normalized Cross-Correlation (ZNCC)

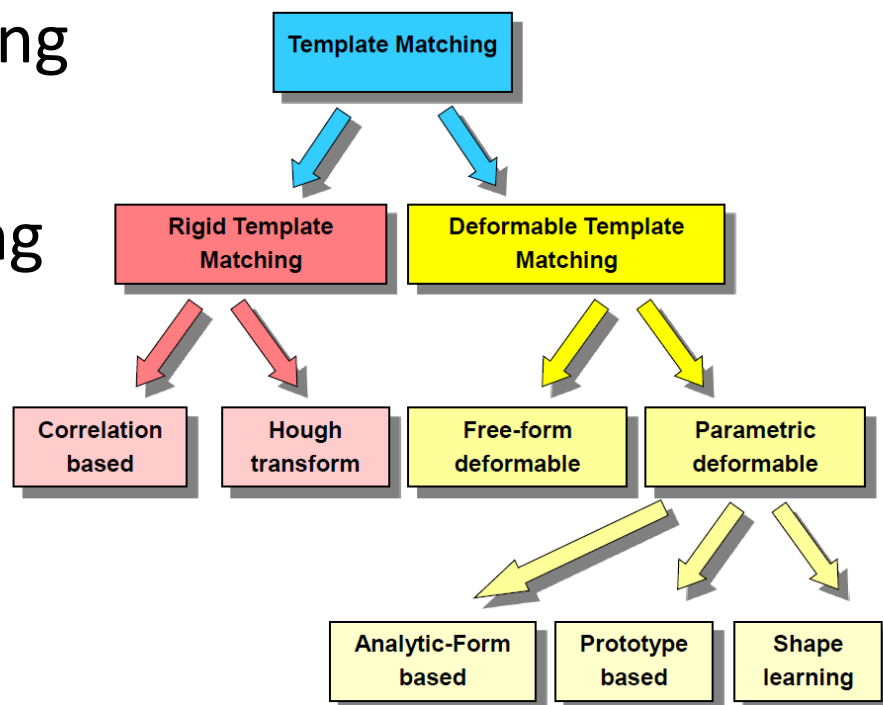
$$\phi(x, y) = \frac{\sum_{u, v \in T} (I(x + u, y + v) - \bar{I}(x, y))(T(u, v) - \bar{T})}{\sigma_I(x, y)\sigma_T}$$

- $\bar{I}(x, y)$: average on window, \bar{T} : template average, $\sigma_I(x, y)$, σ_T : standard deviation on image window and on template

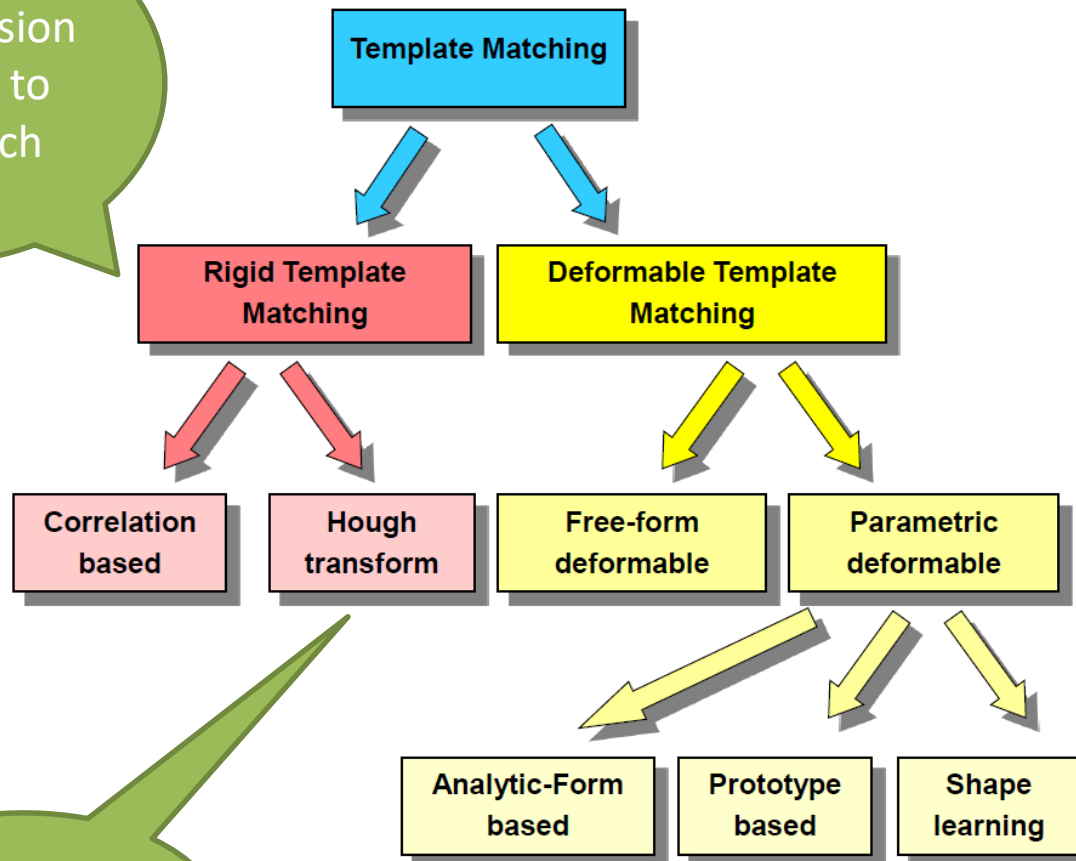


- Dealing with illumination changes
 - Use edge maps instead of images
 - Use ZNCC: subtracts the uniform illumination component
- Dealing with scale changes
 - Matching with several scaled versions of the template
 - Work with multiple rescaled copies of the image
- Dealing with rotation
 - Matching with several rotated versions of the template

- The basic idea of TM generated a family of approaches
 - Different ways of defining the template
 - Different ways of dealing with the template



Our discussion
is related to
this branch



Check this!



- The generalized Hough transform can be seen as a form of template matching
- The Hough transform works for more complex shapes
- General equation:

$$g(\boldsymbol{v}, \boldsymbol{c}) = 0$$

- Where \boldsymbol{v} is a vector of coordinates and \boldsymbol{c} a vector of coefficients

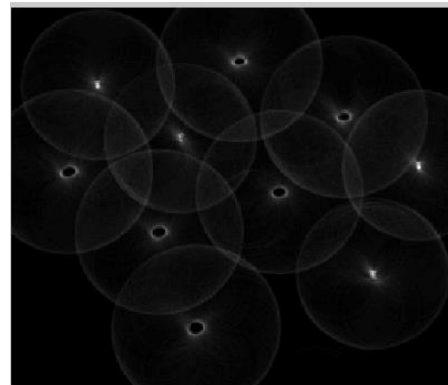
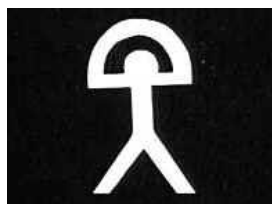


- E.g. (circle):

$$(x - c_1)^2 + (y - c_2)^2 = c_3^2$$

- The parameter space might have high dimensionality

Generalized Hough transform – ex.

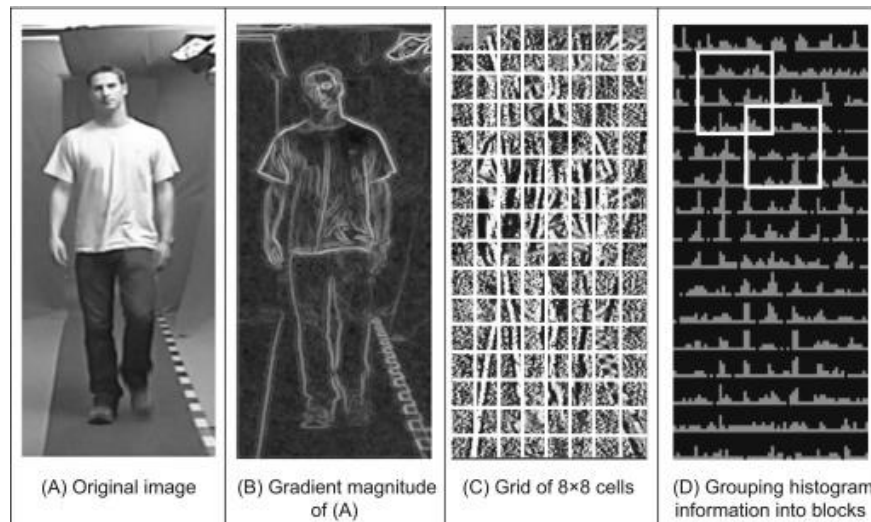


Histogram of Oriented Gradients (HOG)

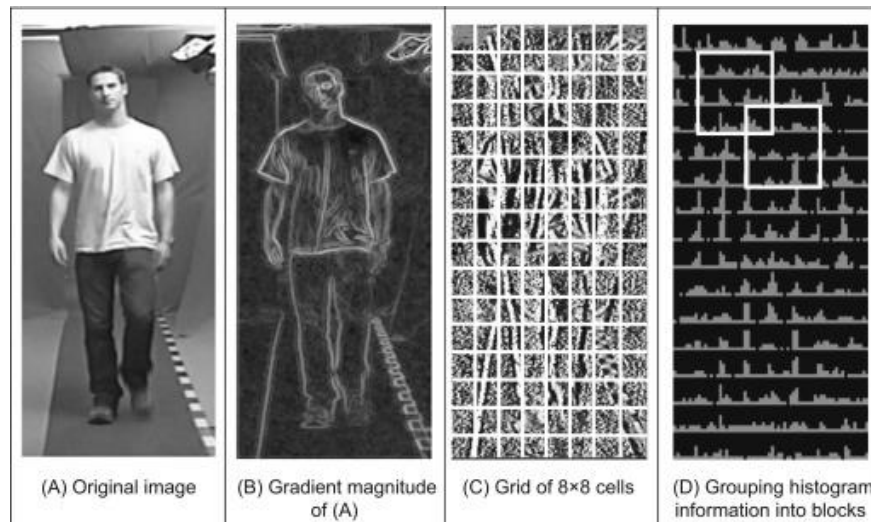


- HOG-based detectors work by
 - Sliding a window (similarly to TM)
 - Characterizing the window by evaluating the edge magnitude and phase – this produces a descriptor (similar to feature descriptor)
 - Classifying the descriptor

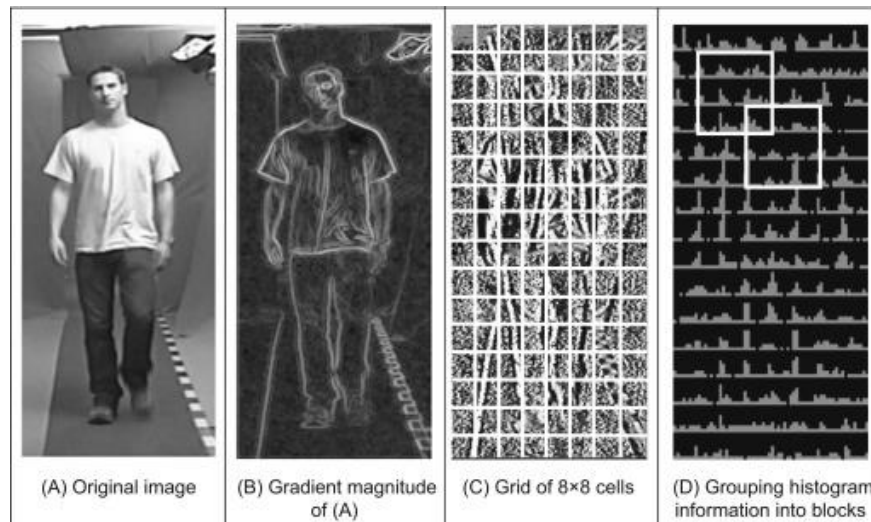
- HOG descriptor evaluation
 1. Intensity normalization/histogram equalization + smoothing
 2. Calculate edge map (magnitude + phase)



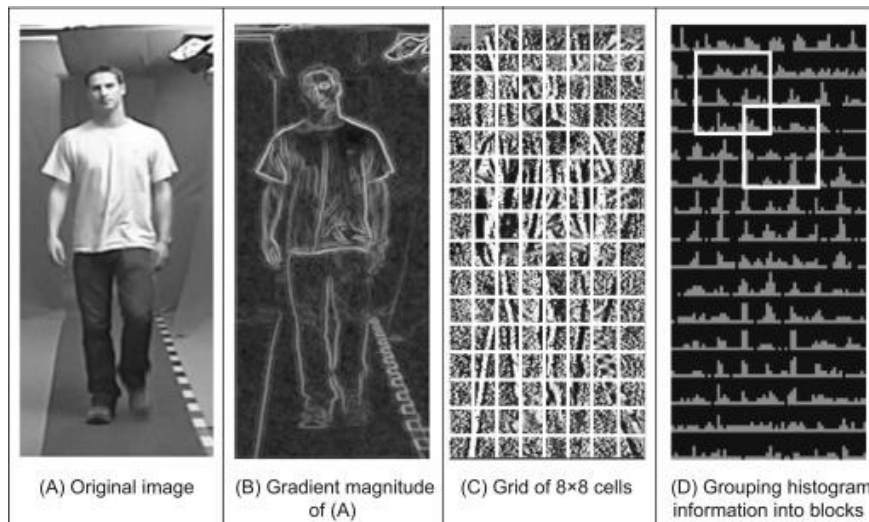
- HOG descriptor evaluation
 1. Detect edges in the image
 2. Compute the gradient magnitude and direction
 3. Evaluate edge histogram on 8x8 non-overlapping cells – this creates voting vectors (e.g., 9 bins of 20° for covering 180°)



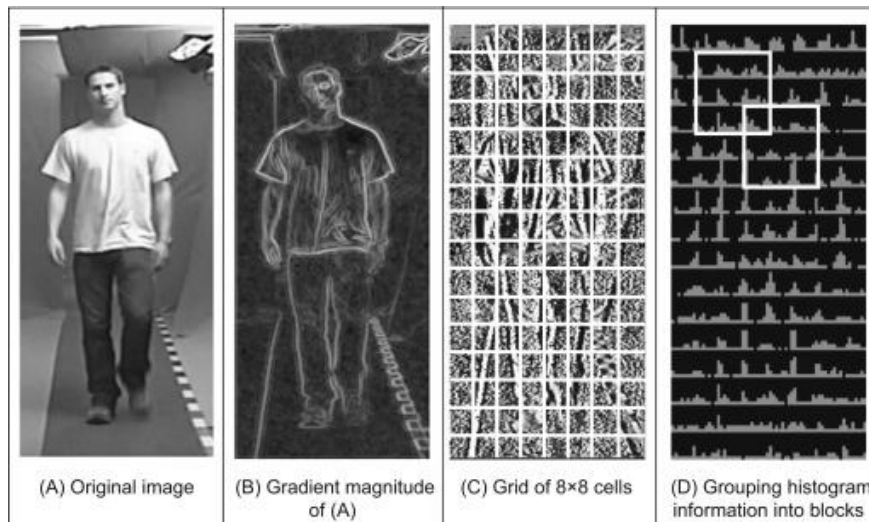
- HOG descriptor evaluation
 4. Create overlapping blocks of 2x2 cells
 5. Normalize voting vectors over each block and create block vectors (36 elements)



- HOG descriptor evaluation
6. Serialize block vectors



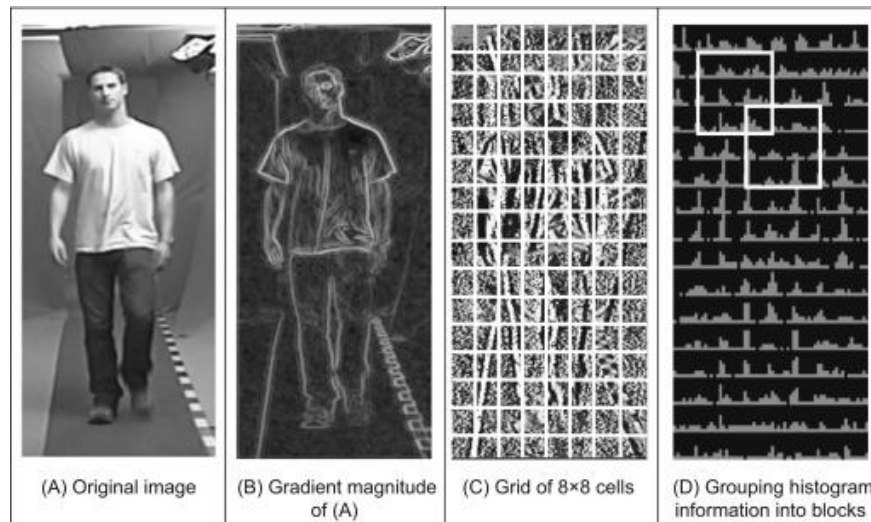
- Consider a 64x128 window
 - How many bins?
 - How many blocks?
 - What is the feature size?





- Anti spoiler 😊

- Consider a 64×128 window
 - How many bins? – 8×16 bins
 - How many blocks? – 7×15 blocks
 - What is the feature size? – 3780 elements





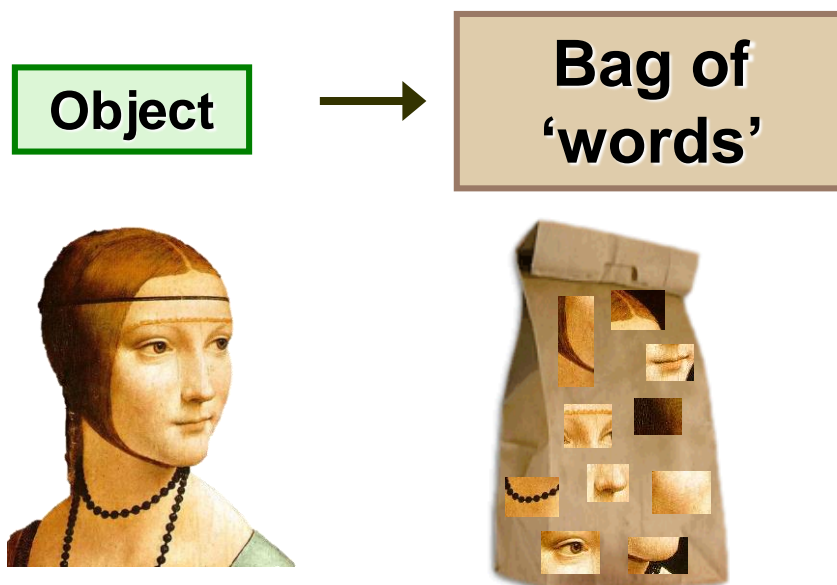
- The HOG descriptor characterizes the content of a bounding box
- The HOG approach needs BBoxes normalized to a standard size
- Multiple scales can be managed by resizing BBoxes of different dimensions to the standard size



- The HOG descriptor is commonly used to train a classifier
- **Note** – the number (cell and block size) described above are one possible implementation for HOG
 - This has an influence on the descriptor size and meaning

Bag of words

- Approach taken from document analysis

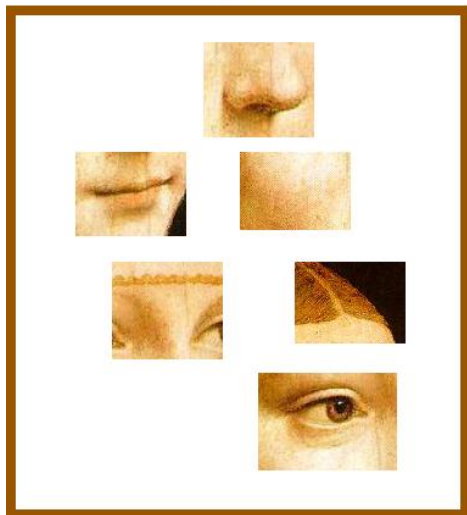




- Image and object classification
- Designed to be invariant to several factors
 - Mainly viewpoint and deformations
- Decomposes complex patterns into (semi) independent features

- Decomposition into visual words

face



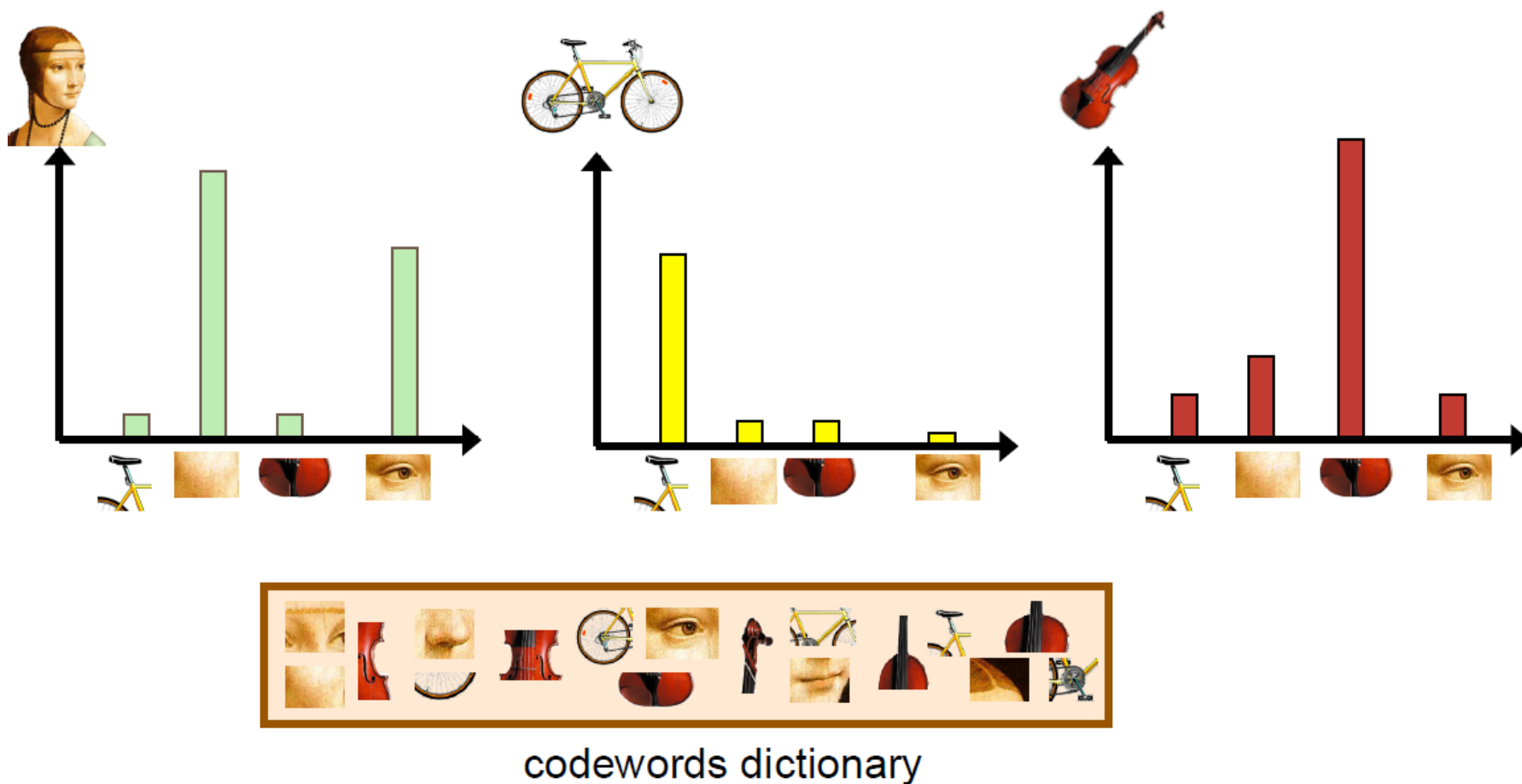
bike



violin



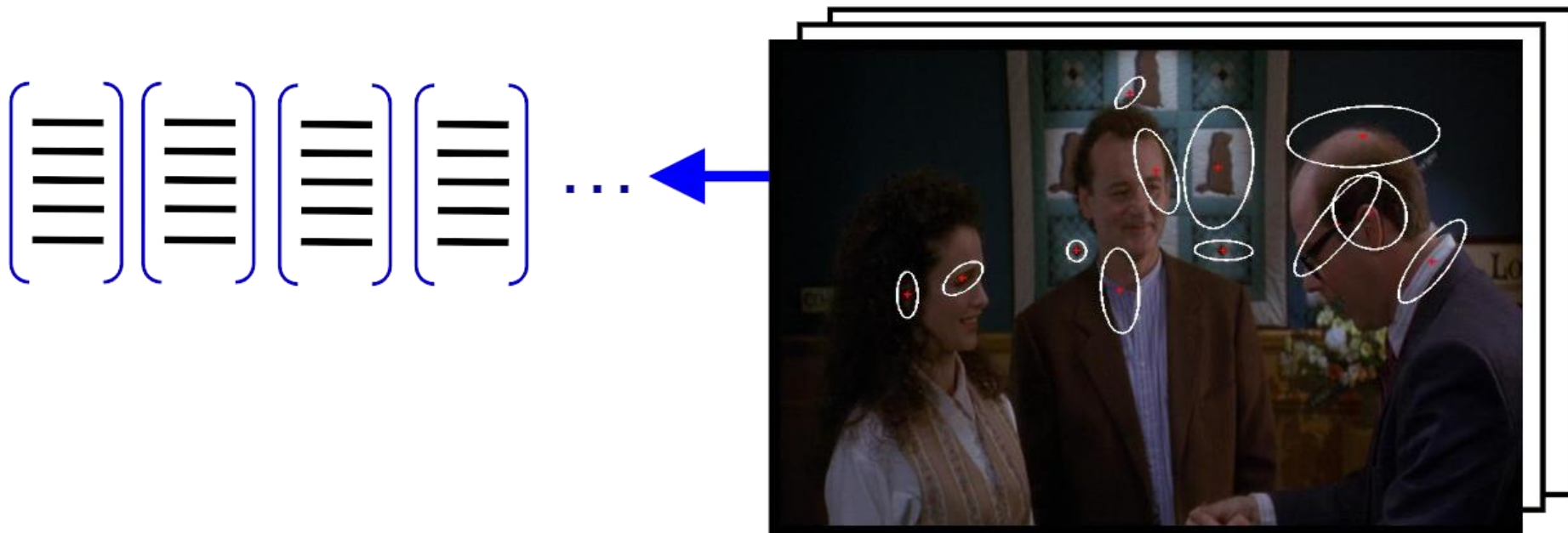
- Histogram representation



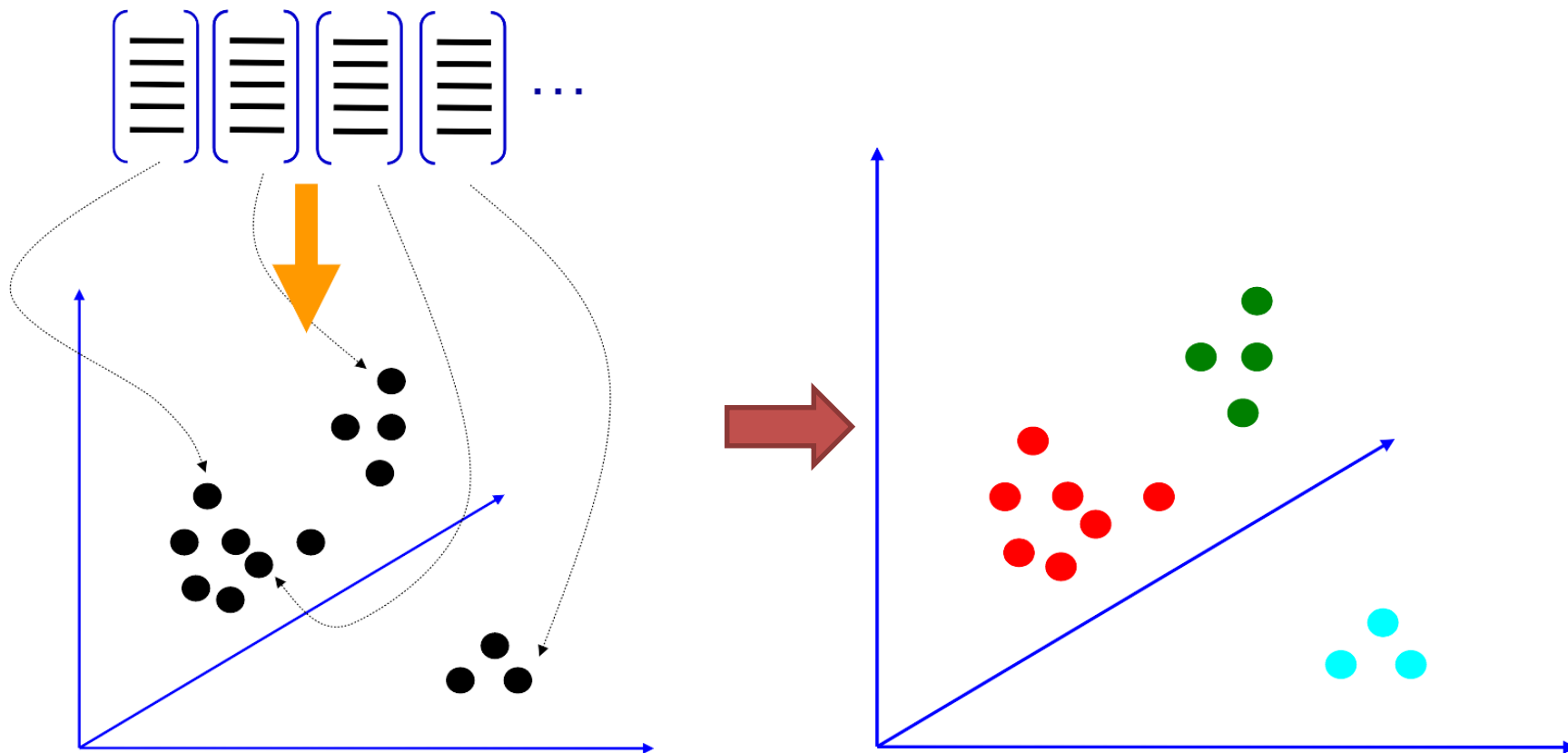


- Words can be represented using features
 - Exploit discriminative properties
 - Exploit invariance properties
 - Re-use an efficient description

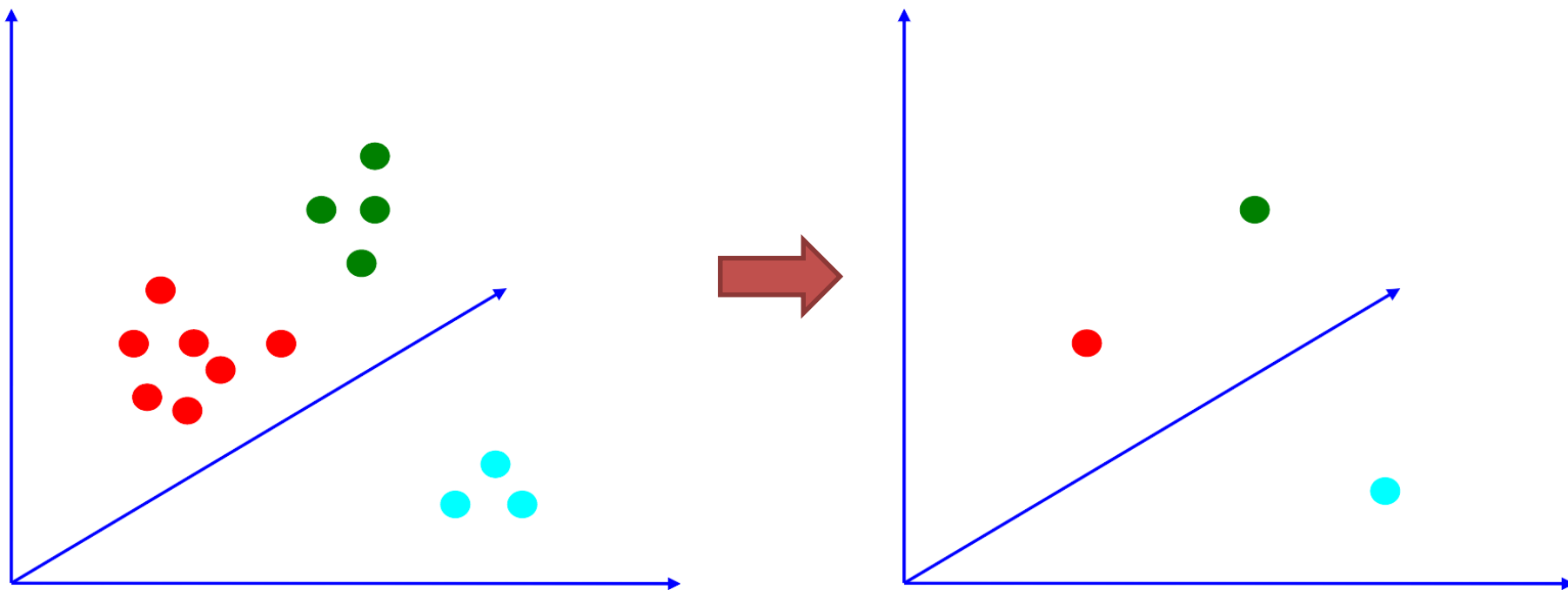
1. Extract features – keypoints and descriptors



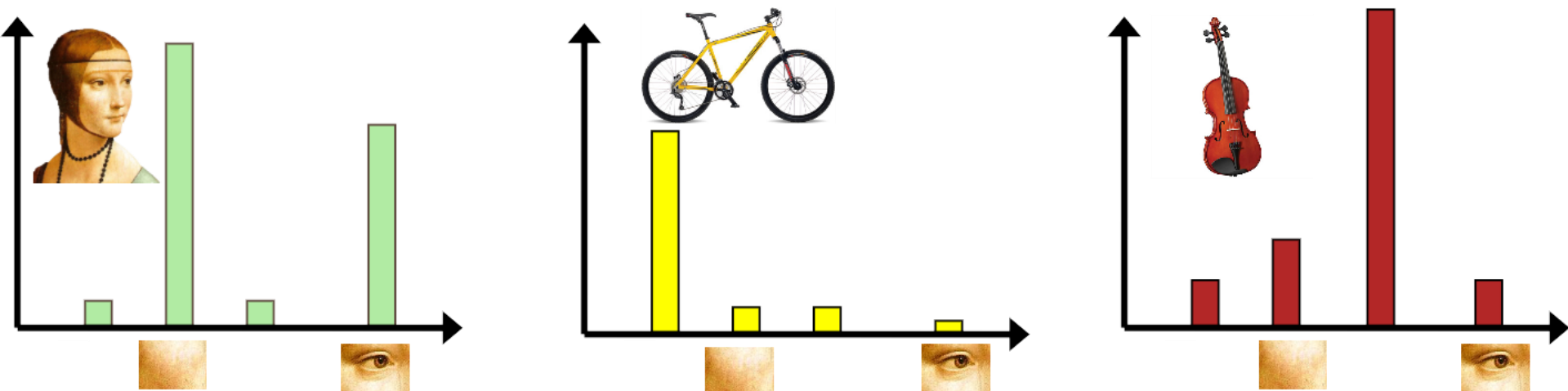
2. Clustering in the feature space (e.g., K-means)



3. Codebook generation: each cluster generates a representative sample (e.g., centroid)



- Image classification:
 - Evaluate the occurrence of each word in the codeword
 - Classify based on histogram





UNIVERSITÀ DEGLI STUDI DI PADOVA

**High-level vision: object recognition,
template matching, bag of words**

Stefano Ghidoni

