

Machine Learning for Modelling: *Supervised Learning*

Simone Bianco

1

SIFT – Scale Invariant Feature Transform

2

Object Recognition

Classes of objects

Class 1: cars



~~Classification~~ Class 3: animals



3

Object Recognition

Object instances

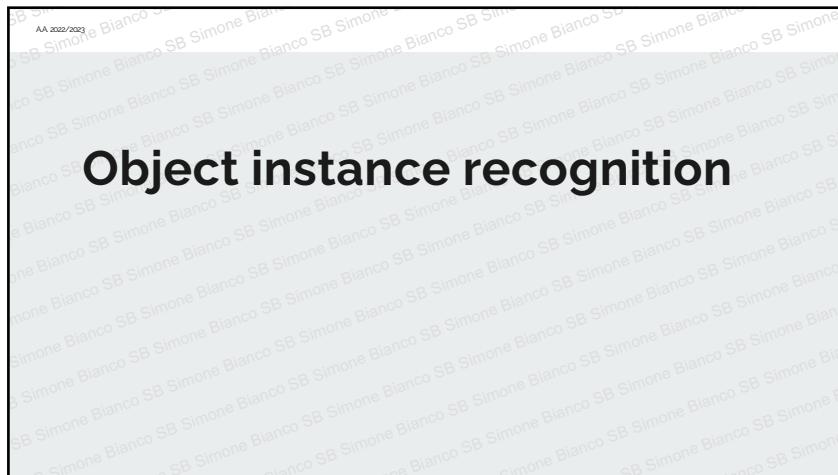
Instance 1: Tesla Model S



Instance 2: Book «The Martian» by Andy Weir



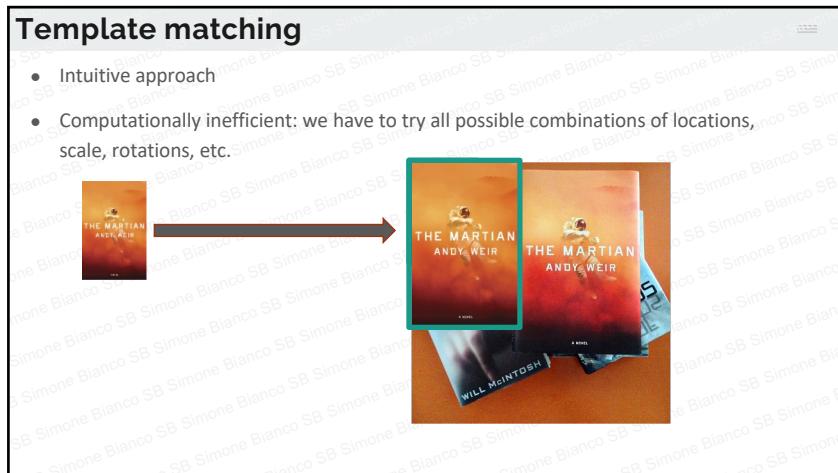
4



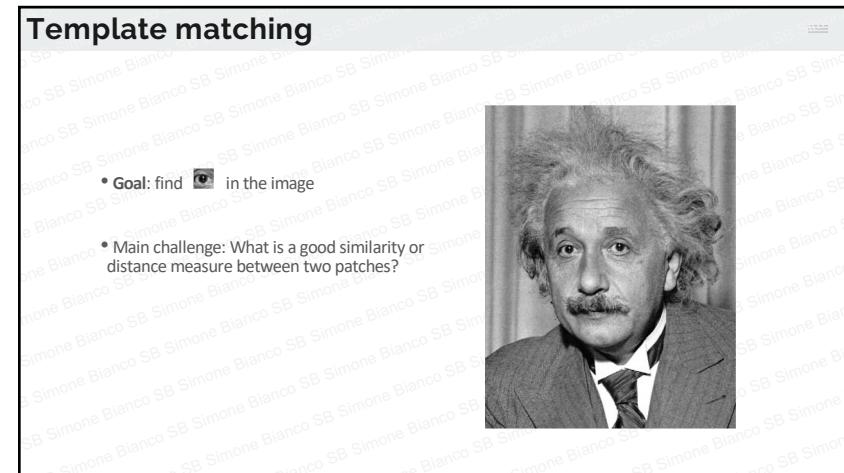
5



6



7



8

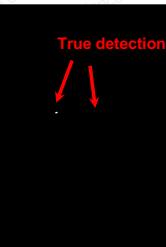
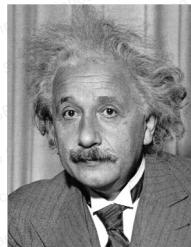
Euclidean distance

Matching with filtersGoal: find  in image

Method: SSD (Sum of Squared Differences)

 $g[k,l]$ template $f[m+k, n+l]$ local regional being analysed

$$h[m,n] = \sum_{k,l} (g[k,l] - f[m+k, n+l])^2$$



9

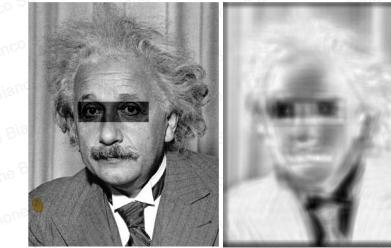
Input image now has a different contrast on the eyes

Now result of computation does not have maximum value on the eyes -> different technique is needed

Matching with filtersGoal: find  in image

Method: SSD (Sum of Squared Differences)

$$h[m,n] = \sum_{k,l} (g[k,l] - f[m+k, n+l])^2$$



1- sqrt(SSD)

10

Normalized cross-correlation uses the mean template and mean image patch.
The denominator ensures that this is a cross-correlation

Matching with filtersGoal: find  in image

Method : Normalized cross-correlation

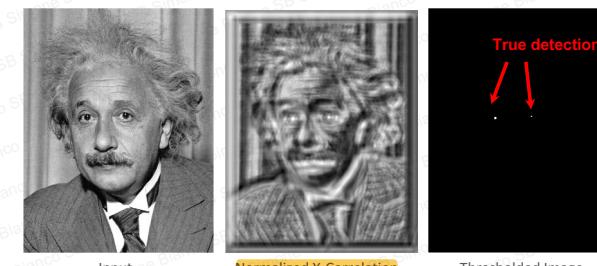
$$h[m,n] = \frac{\sum_{k,l} (g[k,l] - \bar{g})(f[m-k, n-l] - \bar{f}_{m,n})}{\left(\sum_{k,l} (g[k,l] - \bar{g})^2 \sum_{k,l} (f[m-k, n-l] - \bar{f}_{m,n})^2 \right)^{0.5}}$$

11

Normalized X-correlation for normal input image (no different regions of contrast)

Matching with filtersGoal: find  in image

Method : Normalized cross-correlation



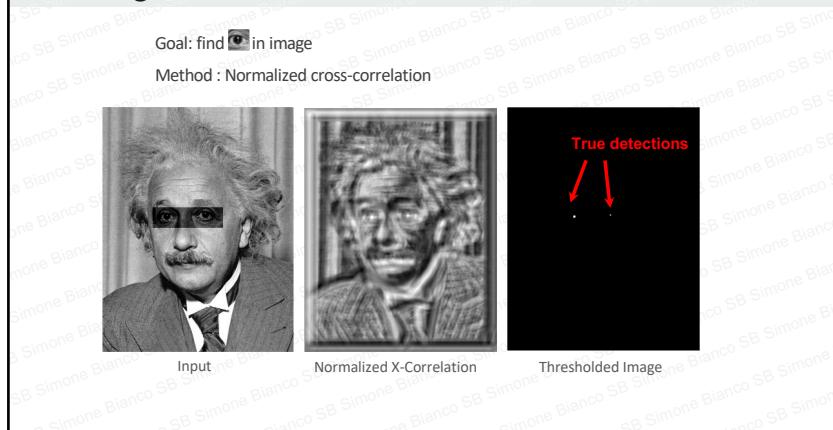
Normalized X-Correlation

12

Norm X-correlation for input with mask over eyes shows that the result is the same:

- invariant to local average intensity
- invariant to local contrast

Matching with filters



13

Matching with filters

Q: What is the best method to use?

SSD

- Faster
- Sensitive to overall intensity

Normalized cross-correlation

- Slower
- Invariant to local average intensity
- Invariant to local contrast

14

Template matching – Sliding window

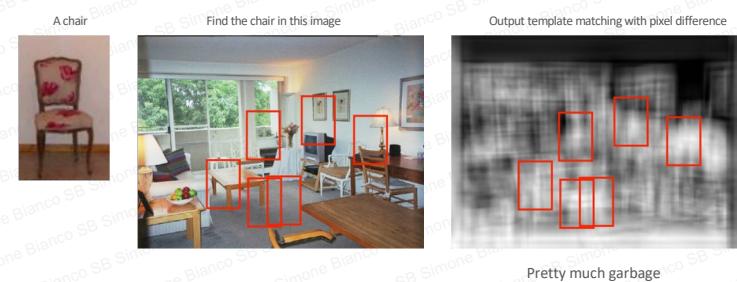
Example



15

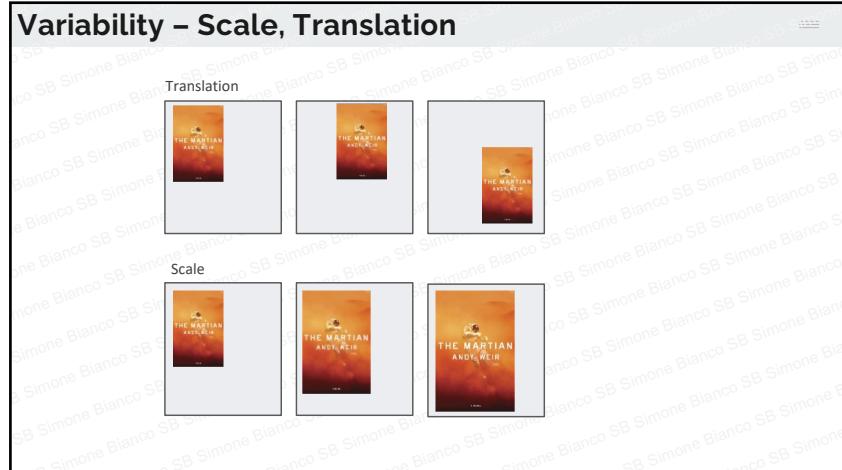
Template matching – Sliding window

Example



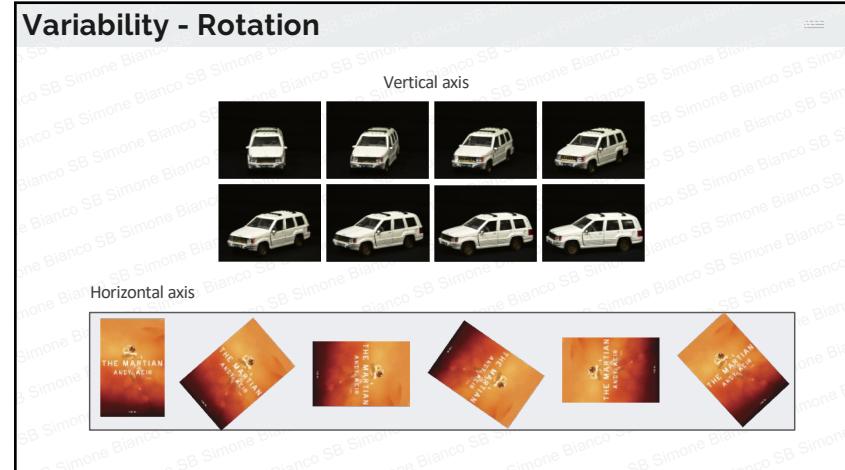
16

Variability – Scale, Translation



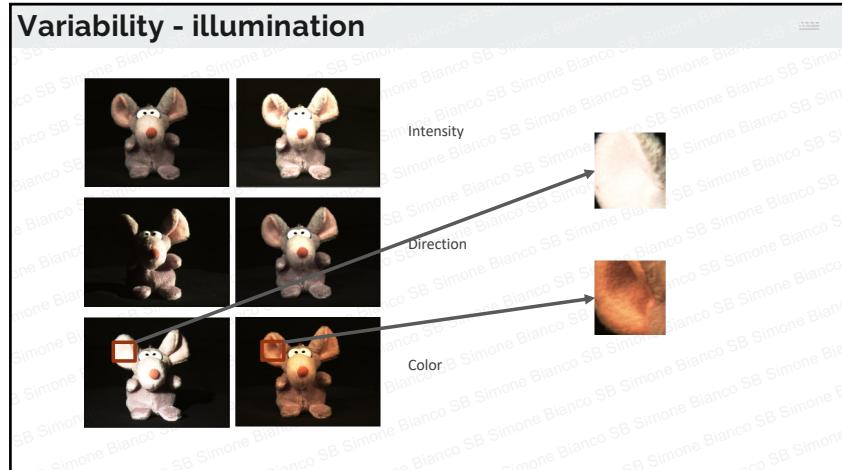
17

Variability - Rotation



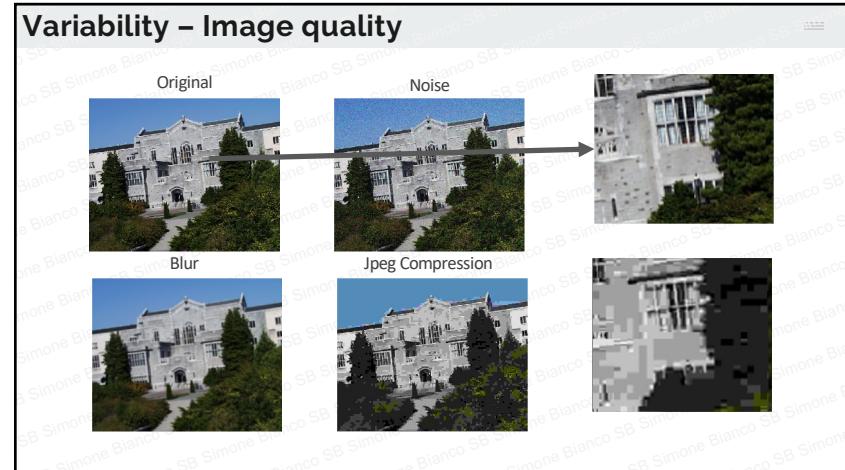
18

Variability - illumination



19

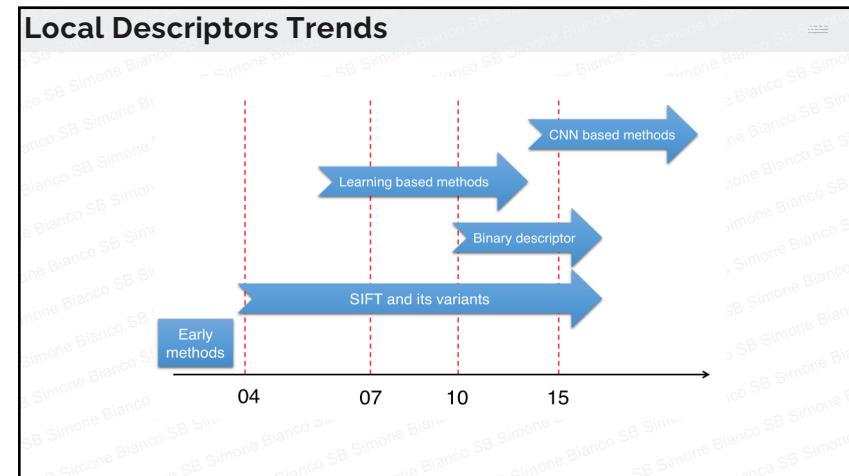
Variability – Image quality



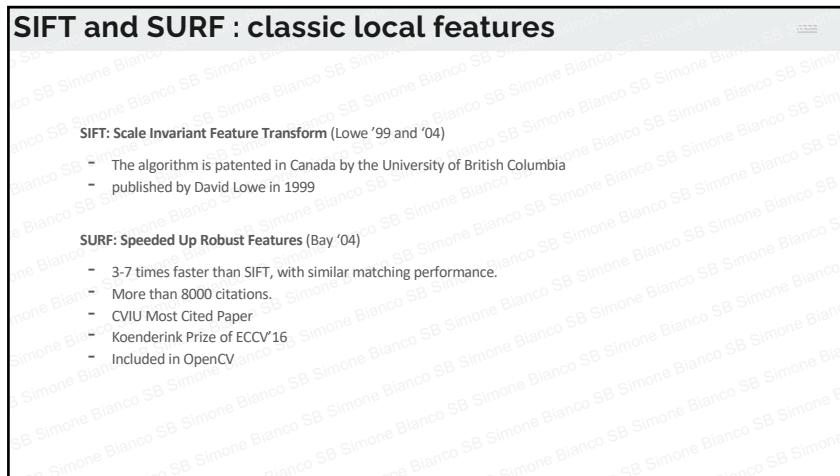
20



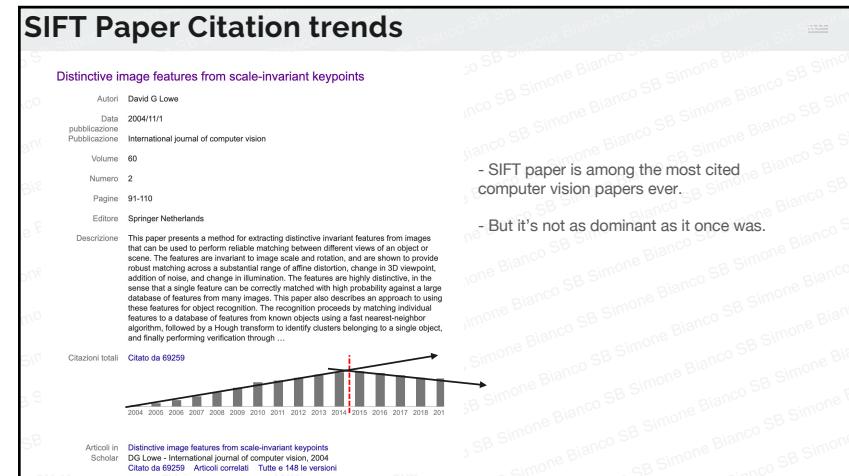
21



22



23



24

Why Keypoints-based methods?

Keypoints Remain Relevant

- When accurate geometric recovery matters, they remain unequalled.
- They are efficient for real-time applications.
- They provide an effective way
 - to compress the information present in large images,
 - to recognize specific locations.
- The algorithms do not need to be retrained for each new application.

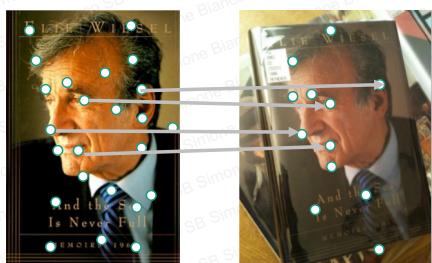
Future algorithms will combine Deep Learning and keypoint matching.

25

Keypoints-Based Approaches

26

Keypoints-based approaches



Steps

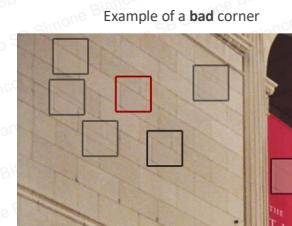
- Keypoint detection
- Keypoint description
- Matching of similar keypoints
- Similarity score based on matching points

- Use of Interest points (keypoints)
- Do not need to try all the combinations of transformations.

27

Harris Corner Detector

- Corners are better than edges!
- How to define a corner?
- Harris et al. [1] → Patches that generate a large variation when moved around



Example of a **bad** corner



Example of a **good** corner

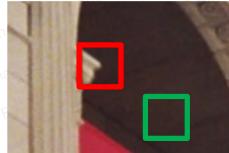
[1] Harris, Chris, and Mike Stephens. "A combined corner and edge detector." Alvey vision conference. Vol. 15. 1988.

28

Harris Corner Detector

Patches that generate a large variation when moved around

$$E(u, v) = \sum_{x,y} w(x, y)[I(x+u, y+v) - I(x, y)]^2$$



- E is the difference between the original and the moved window.
- u is the window's displacement in the x direction
- v is the window's displacement in the y direction
- w(x, y) is the window at position (x, y). This acts like a mask. Ensuring that only the desired window is used.
- I(x, y) is the intensity of the image at a position (x, y)
- I(x+u, y+v) is the intensity of the moved window
- I(x, y) is the intensity of the original window

[1] Harris, Chris, and Mike Stephens. "A combined corner and edge detector." Alvey vision conference. Vol. 15. 1988.

29

Harris Corner Detector

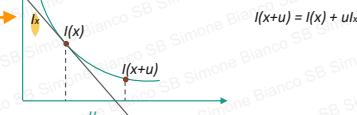
$$E(u, v) = \sum_{x,y} w(x, y)[I(x+u, y+v) - I(x, y)]^2$$

$$\sum_{x,y} [I(x+u, y+v) - I(x, y)]^2$$

$$E(u, v) \approx \sum_{x,y} [I(x, y) + uI_x + vI_y - I(x, y)]^2$$

Derivative

Math hints (derivative):



30

Harris Corner Detector

$$E(u, v) = \sum_{x,y} w(x, y)[I(x+u, y+v) - I(x, y)]^2$$

$$\sum_{x,y} [I(x+u, y+v) - I(x, y)]^2$$

$$E(u, v) \approx \sum_{x,y} [I(x, y) + uI_x + vI_y - I(x, y)]^2$$

$$E(u, v) \approx \sum_{x,y} [u^2 I_x^2 + 2uv I_x I_y + v^2 I_y^2]$$

Square expansion

$$E(u, v) \approx [u \ v] \left(\begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \right) \begin{bmatrix} u \\ v \end{bmatrix}$$

Conversion into a matrix

$$M = \sum_{x,y} w(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

$$E(u, v) \approx [u \ v] M \begin{bmatrix} u \\ v \end{bmatrix}$$

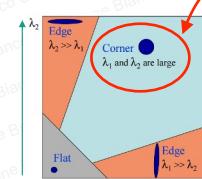
Another way to write it

Harris Corner Detector

$$M = \sum_{x,y} w(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix}$$

M matrix (see slides before)

Eigenvalues of the matrix can help determine the suitability of a window:



$$R = \det M - k(\text{trace } M)^2$$

$$\det M = \lambda_1 \lambda_2$$

$$\text{trace } M = \lambda_1 + \lambda_2$$

All windows that have a score R greater than a certain value are corners.

λ_1, λ_2 are the eigenvalues

[1] Images credits: <http://alishack.in/tutorials/harris-corner-detector/> very good website with lot of didactic material

31

32

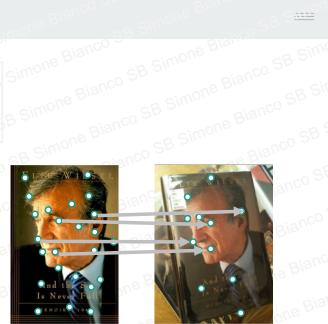
SIFT

- SIFT [1] (Scale Invariant Features Transform)
- Author: David Lowe. University of British Columbia –Canada

Main difference with Harris detector:

- Not only a Keypoint Detector but also a **Keypoint Descriptor**

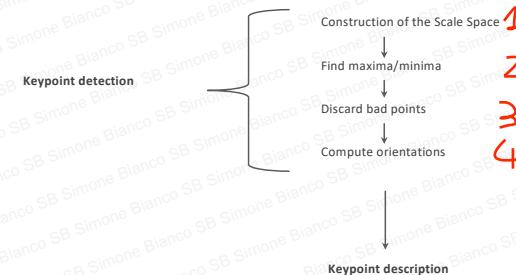
- Robust to:
- Scale
 - Rotation
 - Illumination
 - Viewpoint



[1] David G. Lowe, "Distinctive image features from scale-invariant keypoints," International Journal of Computer Vision, 60, 2 (2004), pp. 91-110.

33

Keypoint detection – Processing steps



34

Scale Space

Increasing blur simulates increasing the scale!



- Needed to obtain **Scale Invariance**

- Higher scales = higher blur
- Higher blur = less details

Increasing blur
Increasing scale

Blurred image Gaussian blur operator Original Image

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y)$$

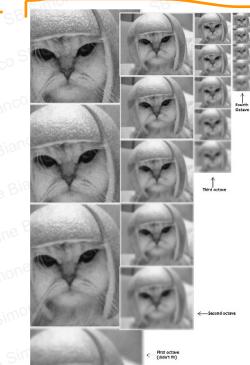
scale parameter: higher value higher blur

$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{-(x^2+y^2)/2\sigma^2}$$

Scale Space - Octaves

4 Octaves

5 blur levels



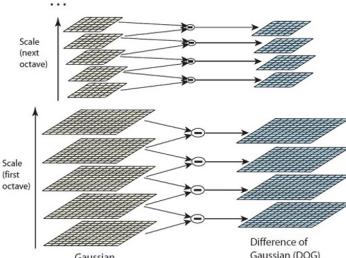
- Scales and Octaves
- Usually 4 octaves and 5 blur levels
- Reduces computations

35

36

Difference Of Gaussians

- Used to detect edges and corners
- Laplacian of Gaussian (LoG) approximation
- Very fast and efficient (only subtraction)



37

Locate maxima/minima in DoG images

- X is the current pixel
- Green circles are the considered neighbours
- $26 (=9+8)$ checks.
- X is a **keypoint** if it is the greatest of all its neighbours
- Usually there is no need to check all the neighbours.
- Discard happens after few checks.

39

Refinements

- To be removed
 - Low Contrast Points
 - Check for the value of pixel in the DoG image
 - If value is under threshold the point is rejected
 - Edges
 - Compute gradients in the two directions.
 - Gradient in one direction is just a difference of pixel values

Flat region

40	40	40
41	45	42
40	41	40

Edge

20	20	20
40	45	40
40	40	40

Corners:
have large difference
both vertically and
horizontally

Corner

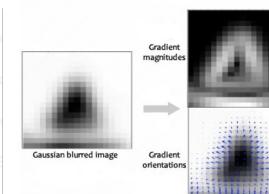
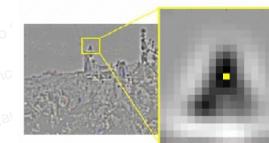
20	20	20
20	45	20
20	40	40

Searching for this!

40

Keypoint Orientation

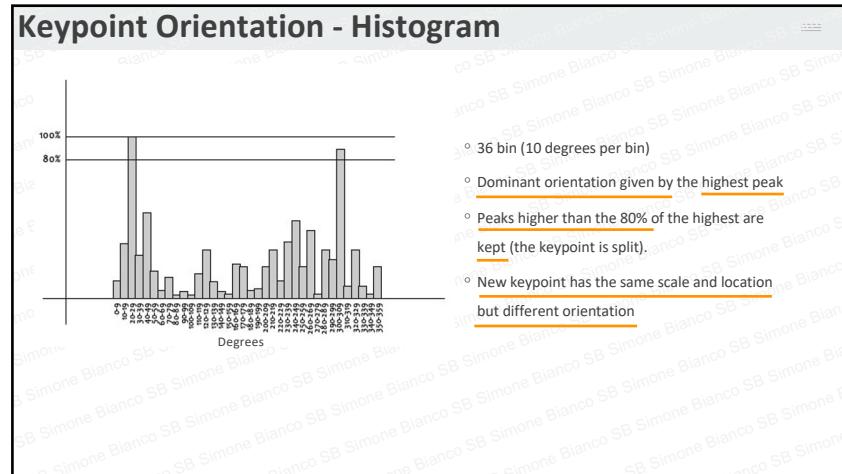
- Goal: achieve rotation invariance
- Histogram of Orientations
 - Weighted by Gradient magnitudes
- Rotation invariance is obtained by assigning an orientation value to the region



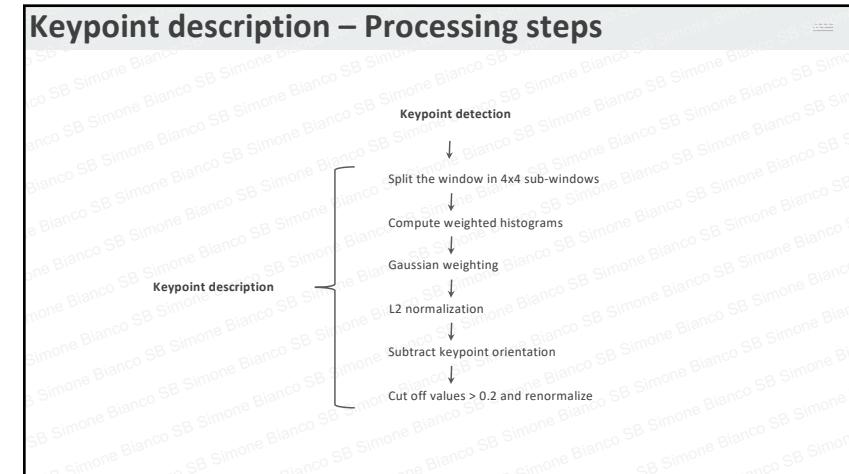
$$m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2}$$

$$\theta(x, y) = \tan^{-1}((L(x, y+1) - L(x, y-1)) / (L(x+1, y) - L(x-1, y)))$$

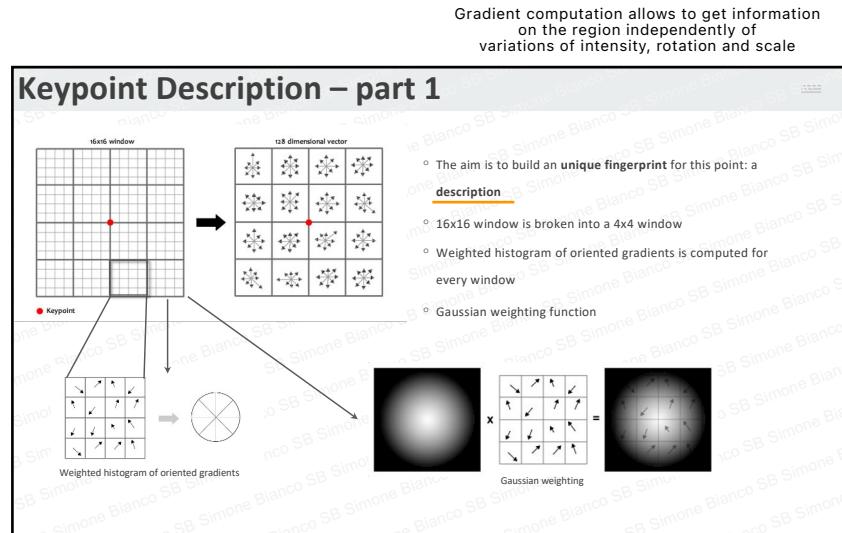
41



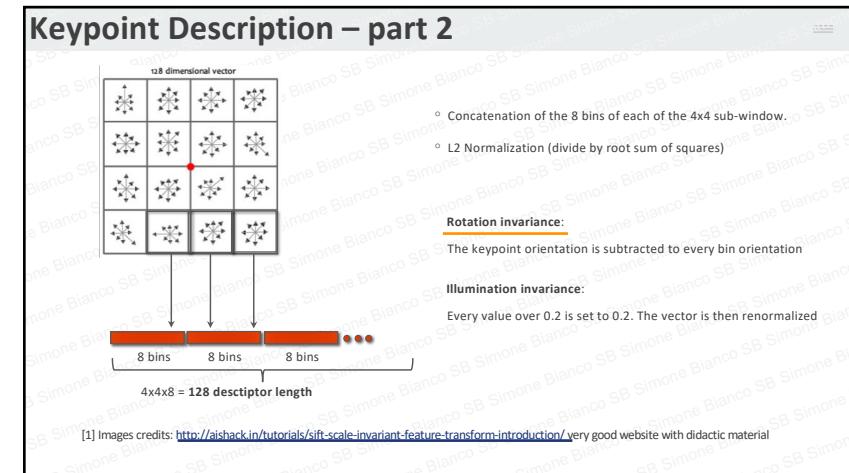
42



43

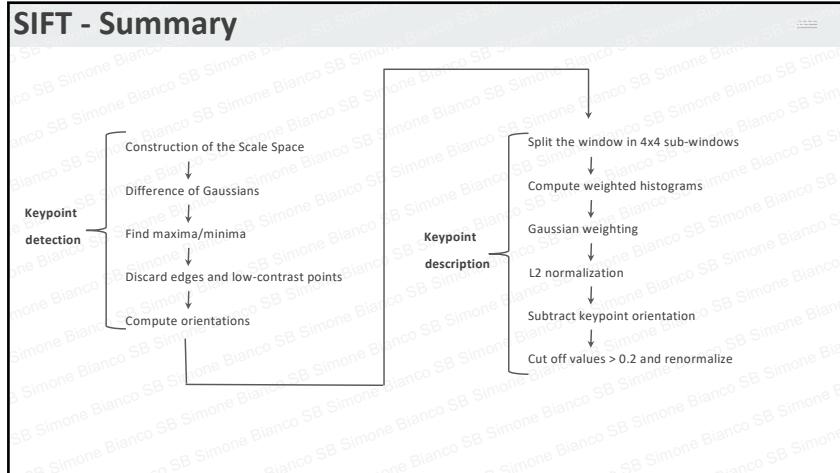


44



45

SIFT - Summary



46

List of keypoint detectors/descriptors

SIFT implementations

Detector	Descriptor	Dimensionality
SIFT Lowe (SIFT 1999)	SIFT Lowe	128
SIFT OpenCV	SIFT OpenCV	128
SIFT VLFeat (Vedaldi 2010)	SIFT VLFeat	128

Affine Invariant Detectors

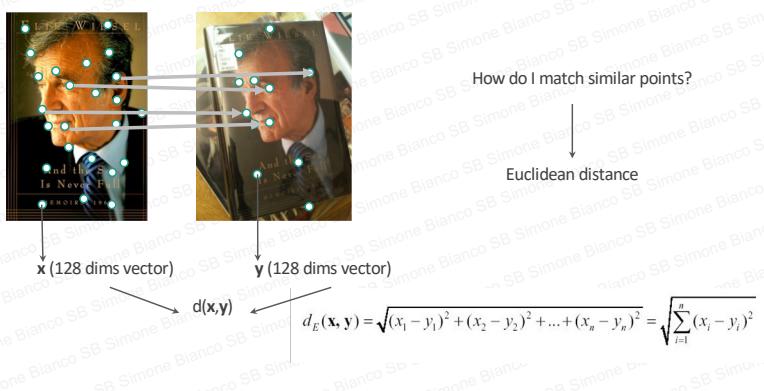
Detector	Descriptor	Dimensionality
DoG	SIFT	128
Multiscale-Harris	SIFT	128
Harris-Laplace	SIFT	128
Hessian	SIFT	128
Multiscale-Hessian	SIFT	128
Hessian-Laplace	SIFT	128

Others

Detector	Descriptor	Dimensionality
SURF	SURF (Bay 2006)	64
SURF	FREAK (Alahi 2012)	64
Kaze	Kaze (Alcantarilla 2012)	64

47

Keypoint matching



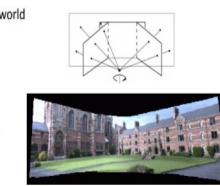
48

Homography

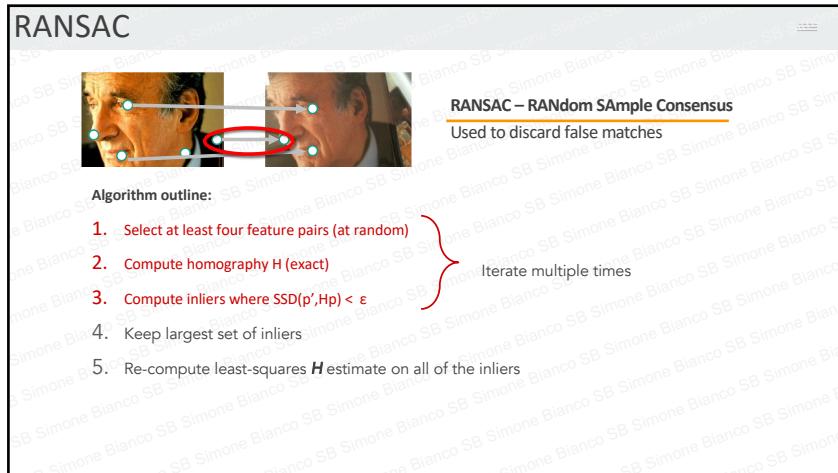
Briefly, the planar homography relates the transformation between two planes

$$s \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = H \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

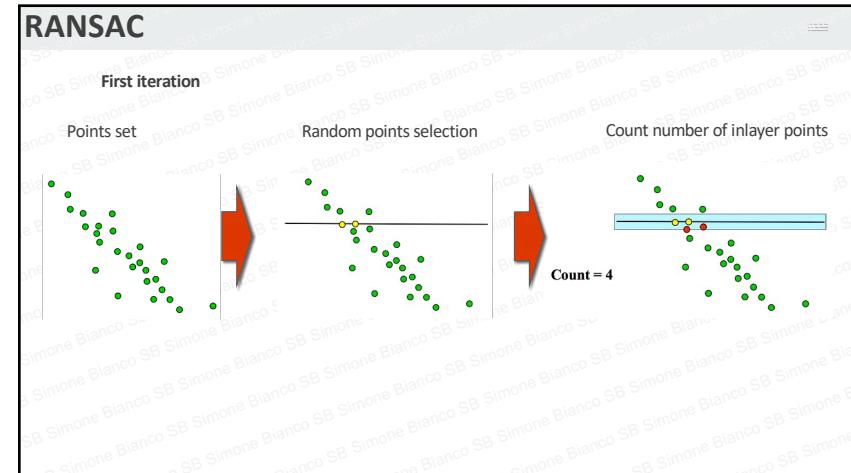
Rotating camera, arbitrary world



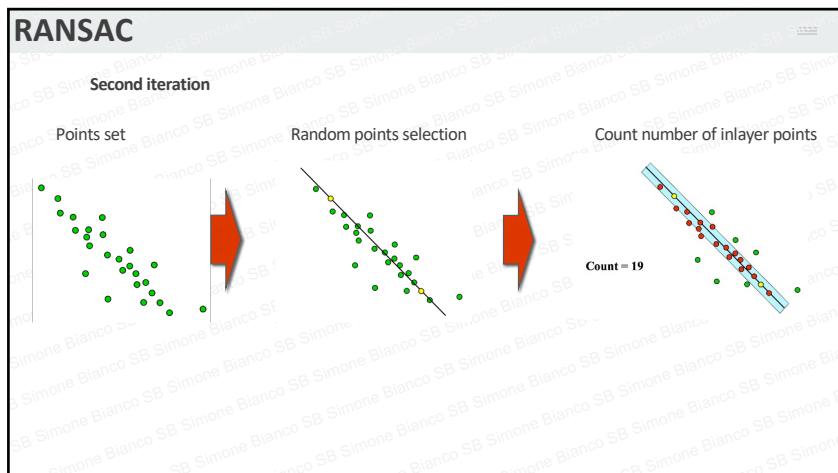
49



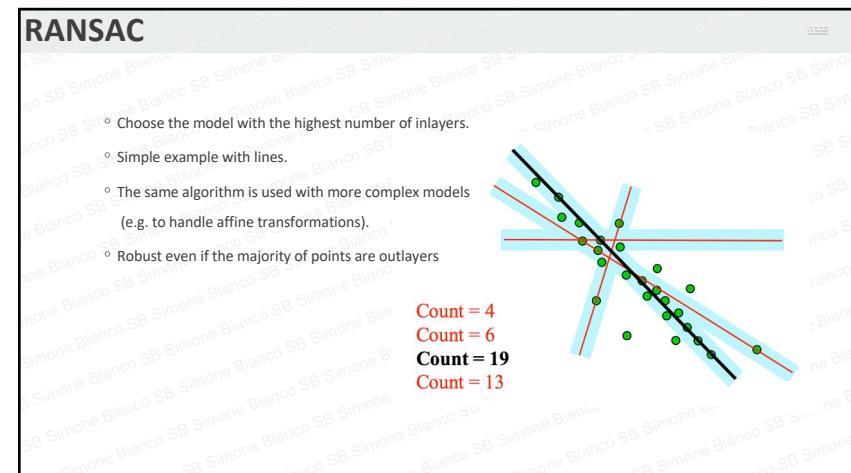
50



51

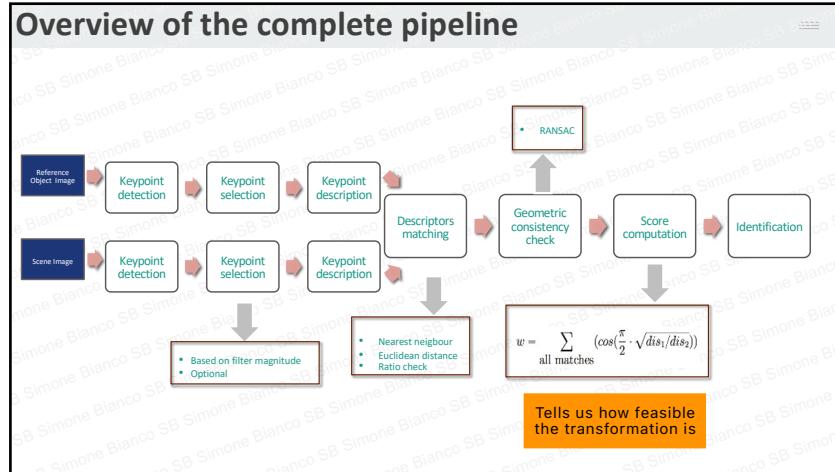


52



53

Overview of the complete pipeline



54

Identification

Example formula used to check the correspondence:

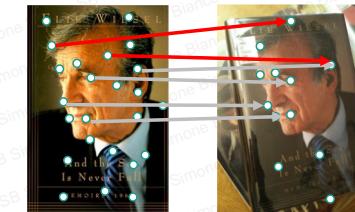
$$w = \sum_{\text{all matches}} (\cos(\frac{\pi}{2} - \sqrt{dis_1/dis_2}))$$

All correct matches
Distance from the nearest and second nearest.

To verify how strong the matching is.

Final score

- To be thresholded
- If w is over the threshold the two images contain the same object

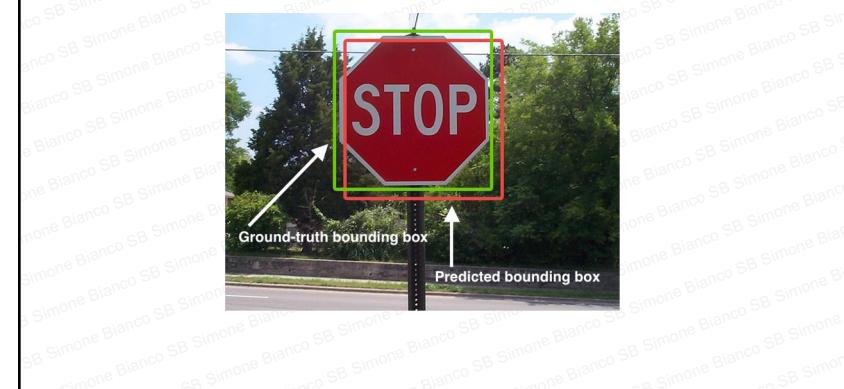


55

Object detection: evaluation metrics

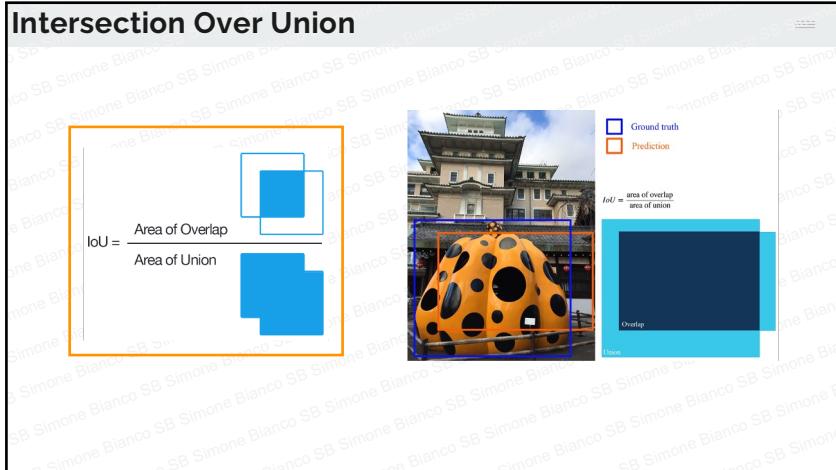
56

Intersection Over Union



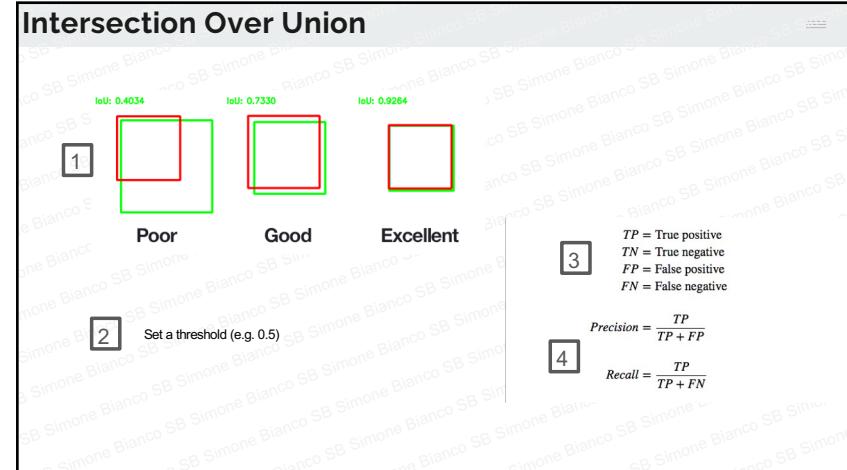
57

Intersection Over Union



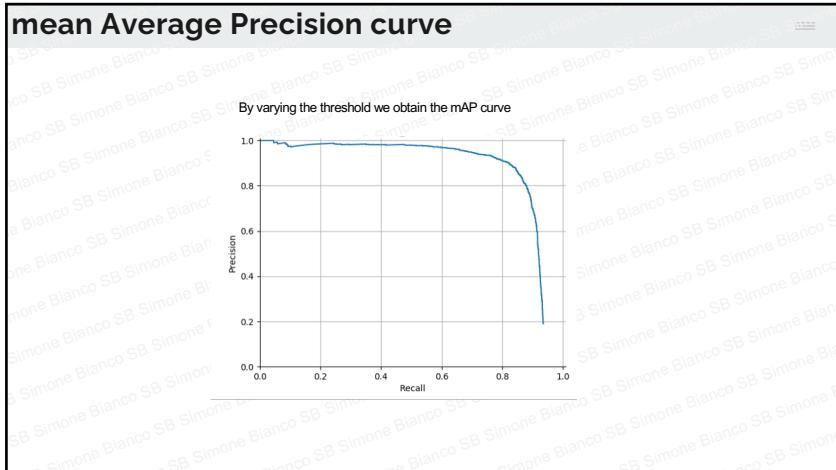
58

Intersection Over Union



59

mean Average Precision curve

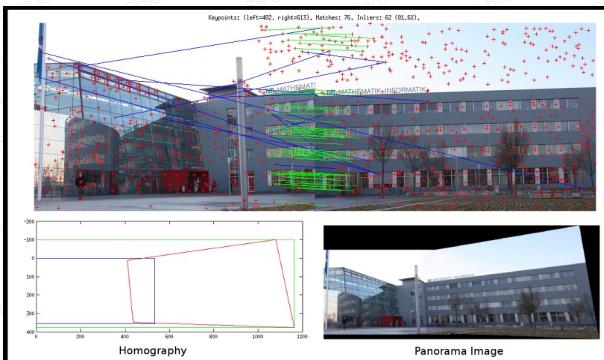


60

Applications

61

Applications – Panorama Stitching



62

Applications – Panorama Stitching



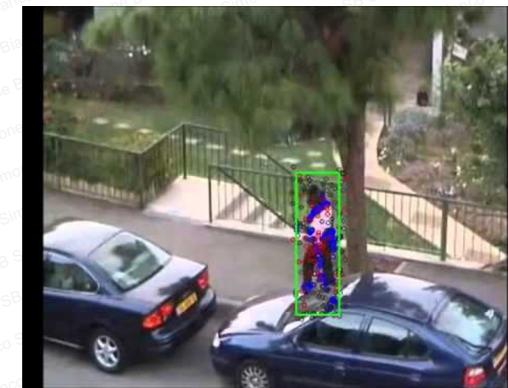
63

Applications – Video Stabilization



64

Applications – Tracking



65

Applications – Augmented reality - IKEA

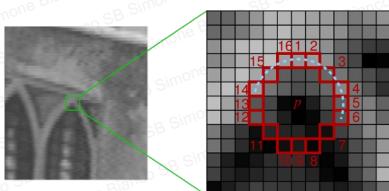


66

Variation on the theme

67

FAST – another keypoint detector



- FAST – Features from accelerated segment test
- For every pixel check the value of the neighbors
- If the value of at least N neighbors is less than the value of the center pixel, this is a corner

Tips for faster computation:

- First compare values of pixels 1,5,9,13
- At least 3 of these pixel values must be under the value of the central pixel
- If not the keypoint candidate is discarded
- Else everyone of the 16 neighbors is checked

68

Dense SIFT

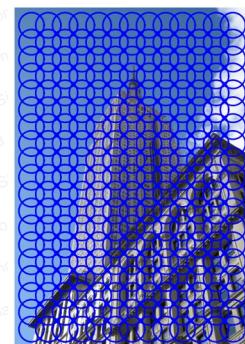
- Instead of computing keypoints use a fixed grid
- Every point in the fixed grid is described using the SIFT descriptor

PRO:

- Faster to compute (no keypoint detection)
- Can work with smooth surfaces (few corners)

CONS:

- Corners are usually more distinctive (reliable)
- Slower keypoint comparison (usually a lot more points to be compared)



69

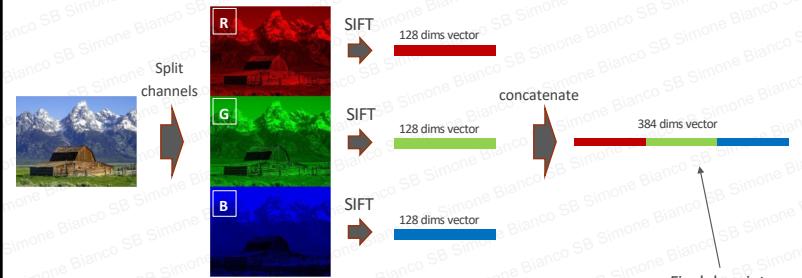
Color descriptors

- SIFT does not use color information
- It just works on grayscale images
- Some objects categories need color information to be distinguished
- Two ways of mixing the color and shape information:
 - Late Fusion approaches
 - Early Fusion approaches



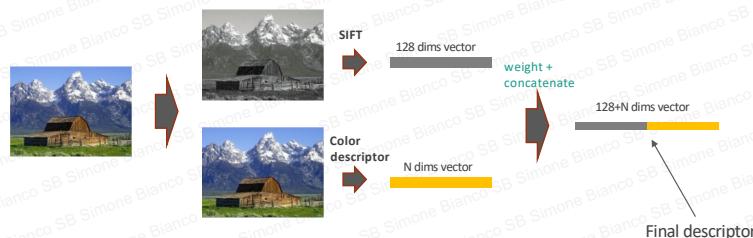
70

Early Fusion approaches



71

Late Fusion approaches



- N depends on the type of descriptor
- Example: simple histogram on color-space
- Before the concatenation of color and shape descriptor usually the two components are weighted
- Possibility to give more importance (higher weight) to color or shape information

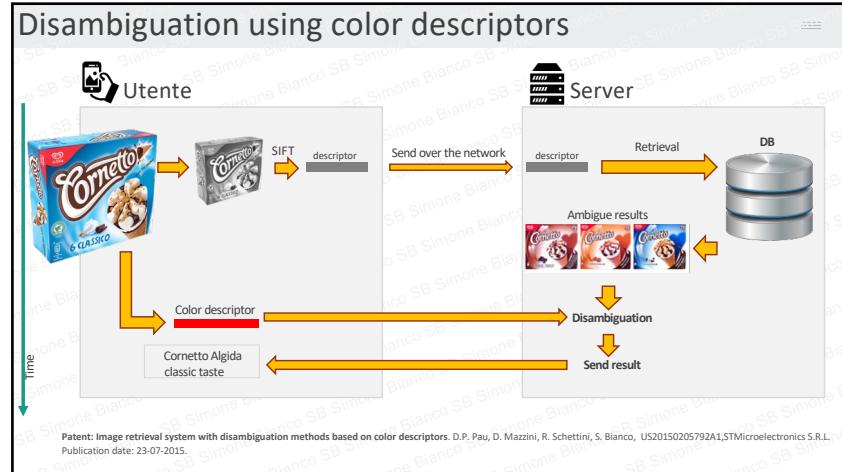
72

Color Descriptors

Detector	Dimension	Fusion
RGB SIFT [1]	384	Early
Opponent SIFT [1]	384	Early
Transformed Color SIFT [1]	384	Early
HSV SIFT [1]	384	Early
C-SIFT [2]	384	Early
rSIFT [1]	256	Early
cRGB-SIFT [1]	384	Early
Hue SIFT [1]	164	Late
Color Name [3]	139	Late
Fuzzy Sets Color Names [3]	139	Late
Discriminative Color [3]	139, 153, 178	Late

[1] K.E. Van De Sande, T. Gevers, C.G. Snoek, Evaluating color descriptors for object and scene recognition, IEEE Trans. Pattern Anal. Mach. Intell. 32 (9) (2010) 1582–1596.
 [2] A.E. Abdel-Hakim, A.A. Farag, Csift: a sift descriptor with color invariant characteristics, in: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, vol. 2, IEEE, 2006, pp. 1978–1983.
 [3] Van De Weijer, C. Schmid, Applying color names to image description, in: IEEE International Conference on Image Processing, 2007, ICIP 2007, vol. 3, IEEE, 2007, p. III-493.

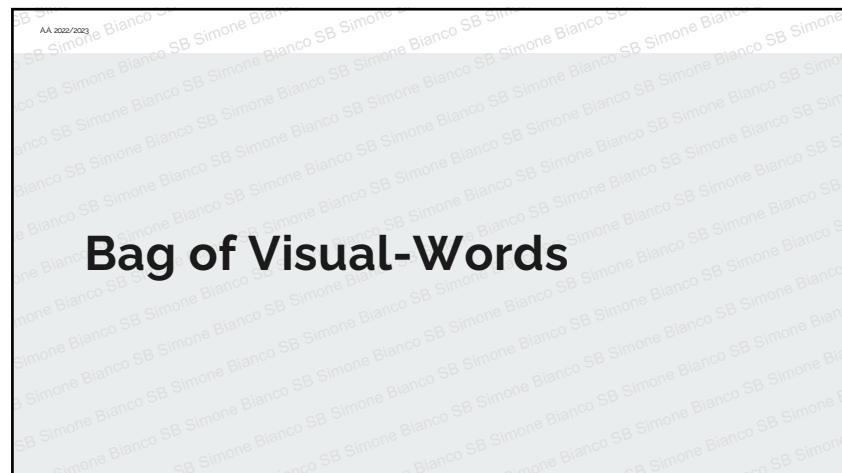
73



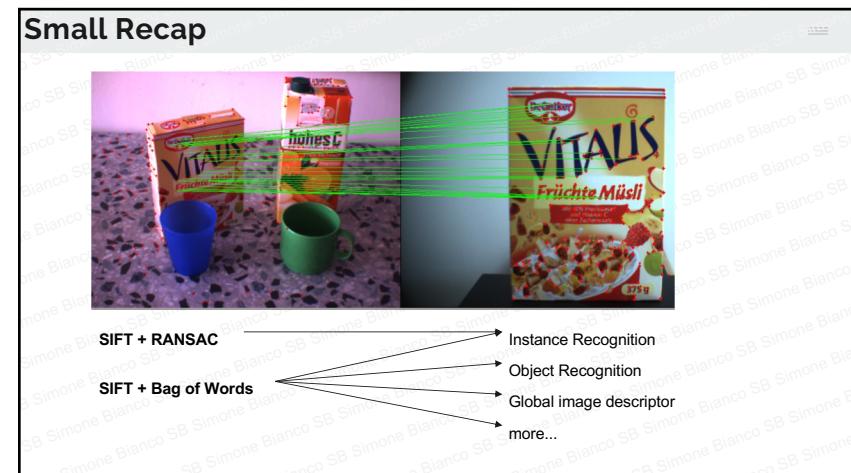
74



75



76



77

Bag-of-words: motivation

President George W. Bush Speech In 2001

ADDRESS TO THE JOINT SESSION OF THE 107TH CONGRESS
UNITED STATES CAPITOL
WASHINGTON, D.C. SEPTEMBER 20, 2001

Mr. Speaker, Mr. President Pro Tempore, members of Congress, and fellow Americans:

It is my privilege to speak to you tonight as the speaker to the joint session of the Congress of the United States. Tonight no such report is needed. It has already been delivered by the American people. We have seen it in the courage of passengers, who rushed terrorists to save others on the ground — passengers like the exceptional man named Todd Beamer. And would you please hear me say again, his wife, Lori Beamer, here tonight. We have seen it in the actions of our military, in the actions of our law enforcement agencies. We have seen the unfurling of the flags, the lighting of candles, the giving of blood, the saying of prayers — in English, Hebrew and Arabic. We have seen the decency of a loving and giving people who have made the grief of strangers their own. For the last two days, we have seen the love of our country for the safety of our Union — and it is clear that we are a country awakened to danger and called to defend freedom. Our grief has turned to anger, and anger to resolution. Whether we bring our enemies to justice, or we seek our enemies' justice, will be done. I thank you.

For those of us here, the words of the president of the United States were touching on the evening of the tragedy to see Republicans and Democrats joined together on the steps of this Capitol, singing "God Bless America." And you did more than sing; you acted, by delivering \$40 billion to rebuild our country, and meet the needs of our military. Speaker Hastert, Minority Leader Gephardt, Majority Leader Daschle and Senator Lott, I thank you for your friendship, for your leadership and for your service to our country... (you are continued)

- We want to represent the topic of the document in a compact way
 - **Orderless** document representation: frequencies of words from a dictionary
Salton & McGill (1983)

78

Bag-of-words: motivation



US Presidential Speeches Tag Cloud
<http://chir.ag/projects/preztags/>

79

Bag-of-words: motivation



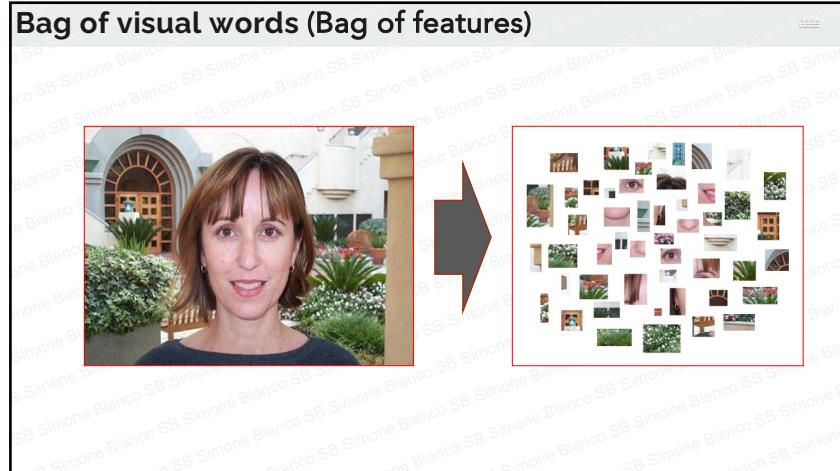
80

Bag-of-words: motivation



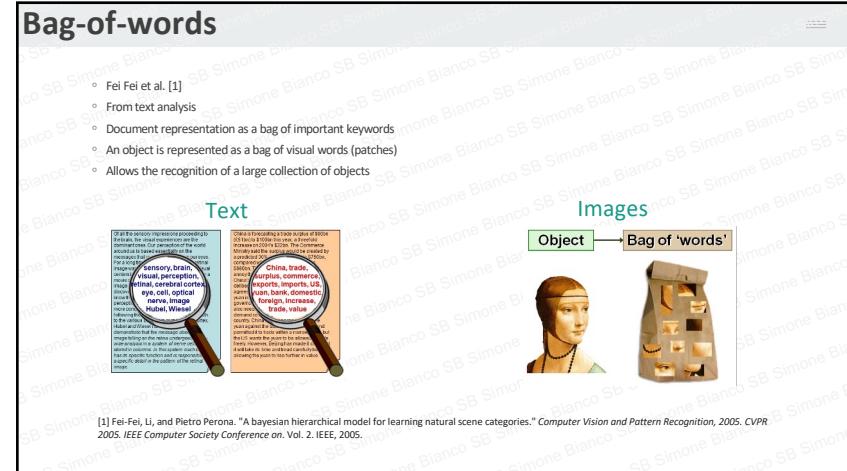
81

Bag of visual words (Bag of features)



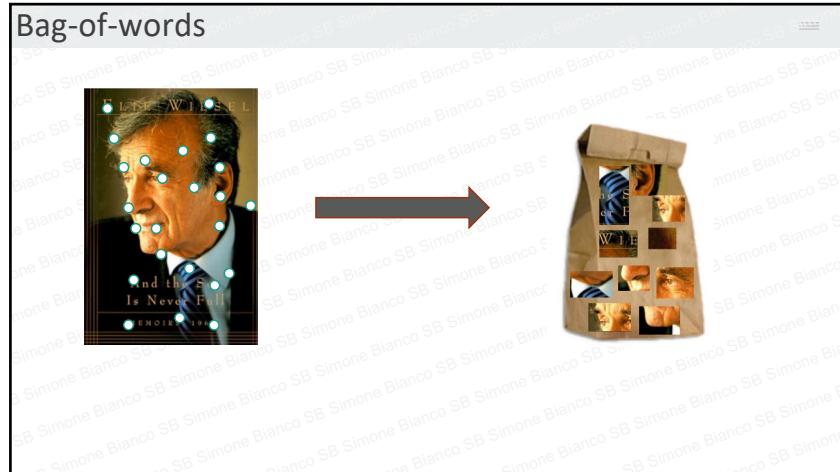
82

Bag-of-words



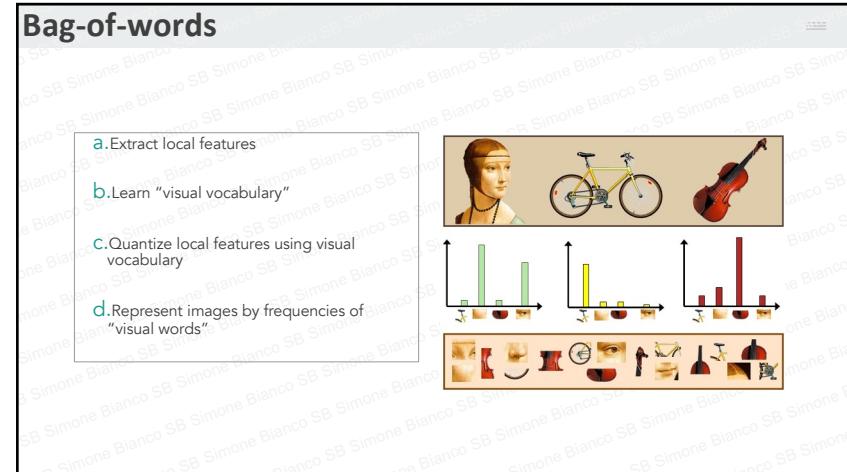
83

Bag-of-words



84

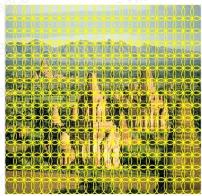
Bag-of-words



85

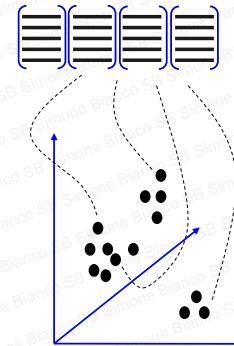
Local features extraction

Sample patches and extract descriptors



86

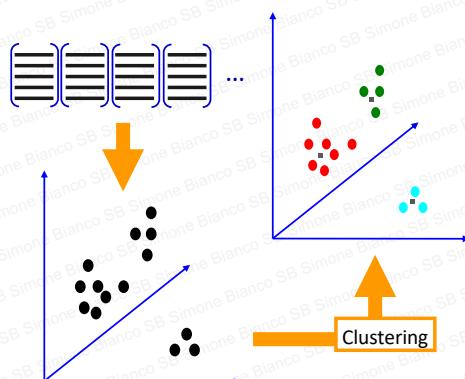
Learning the visual vocabulary (1/2)



Extracted descriptors from
the training set

87

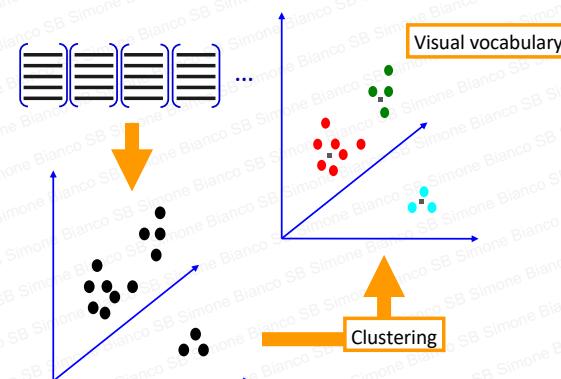
Learning the visual vocabulary (2/2)



Clustering

88

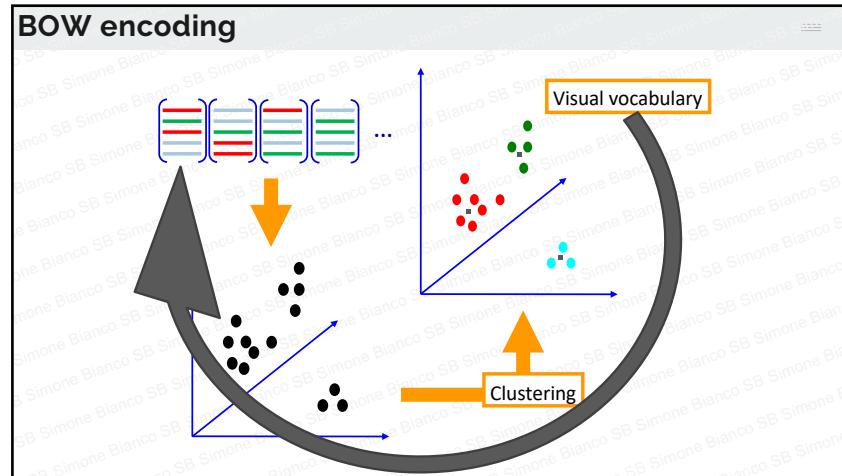
Learning the visual vocabulary (2/2)



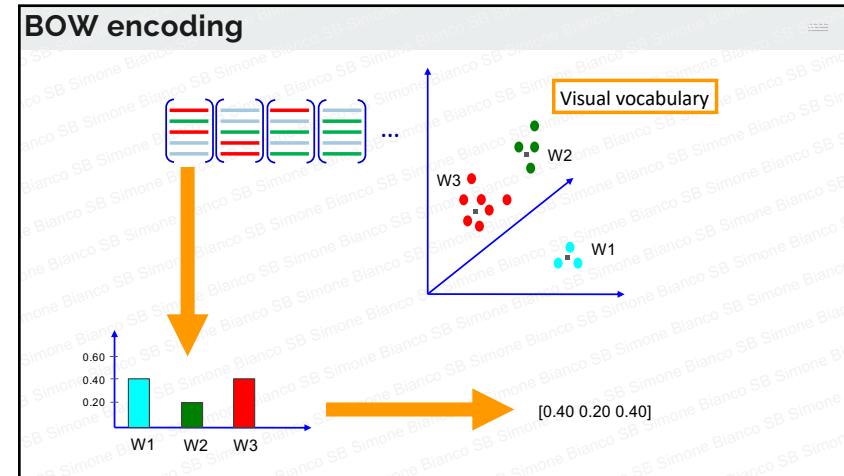
Clustering

Visual vocabulary

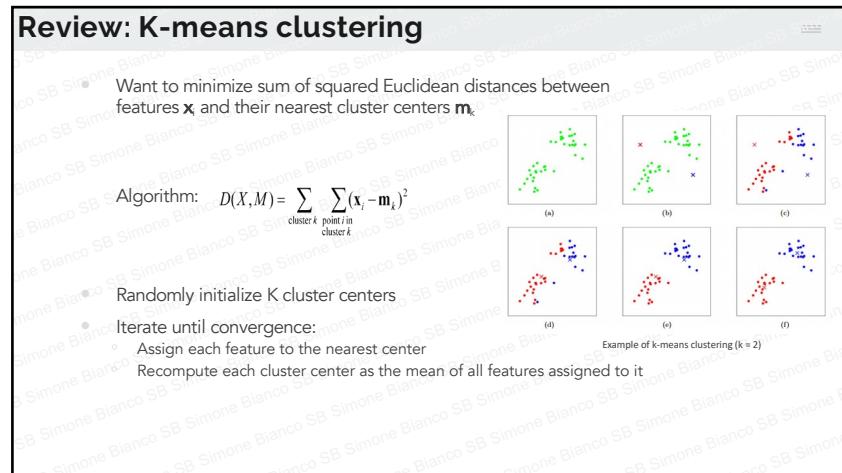
89



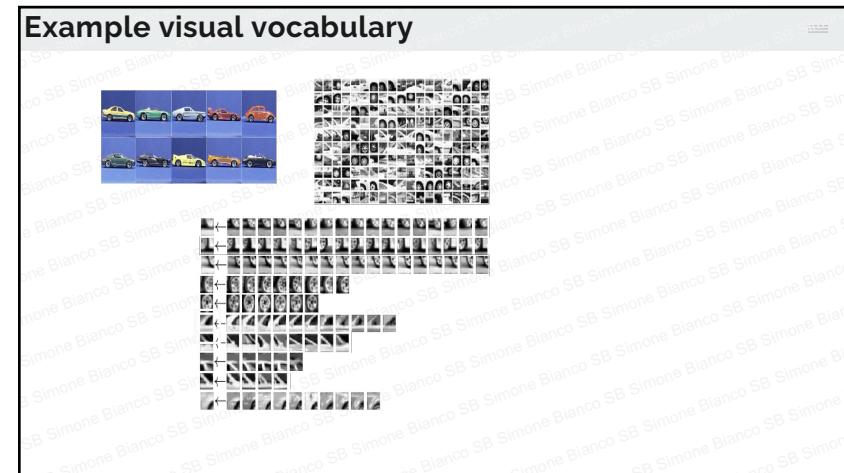
90



91

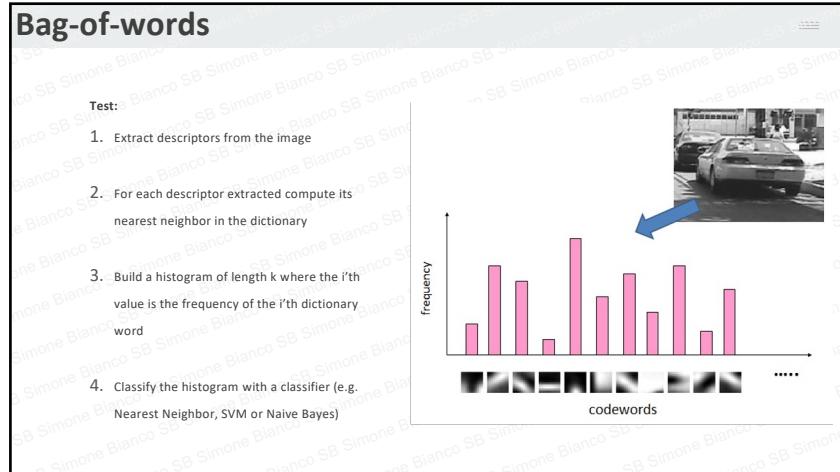


92



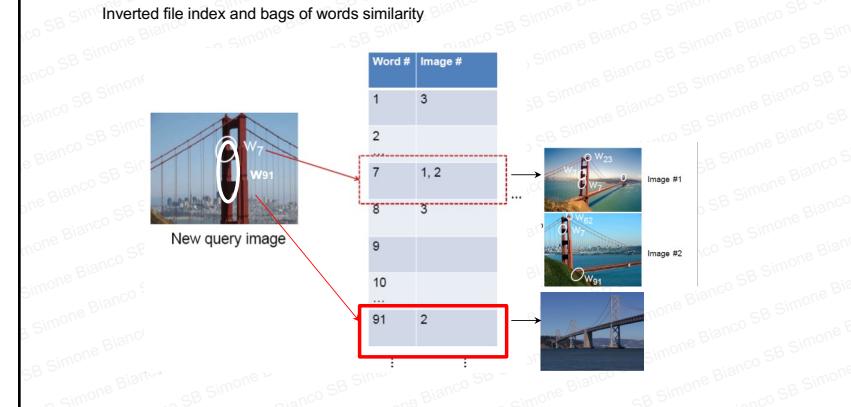
93

Bag-of-words



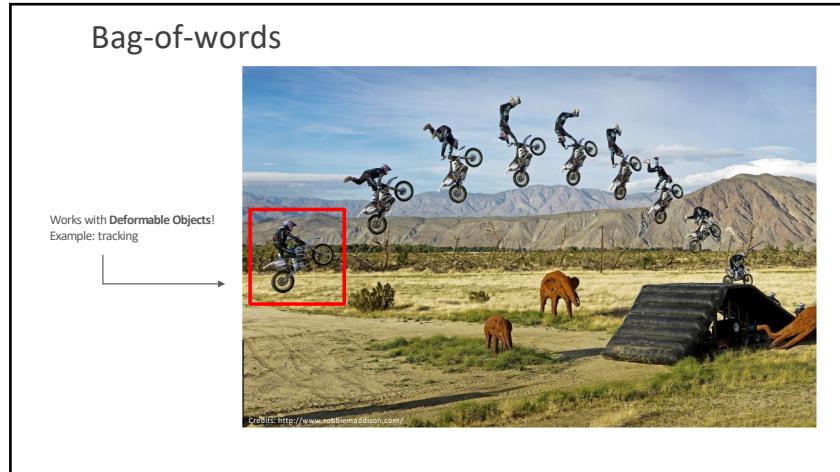
94

Applications of Bag of words



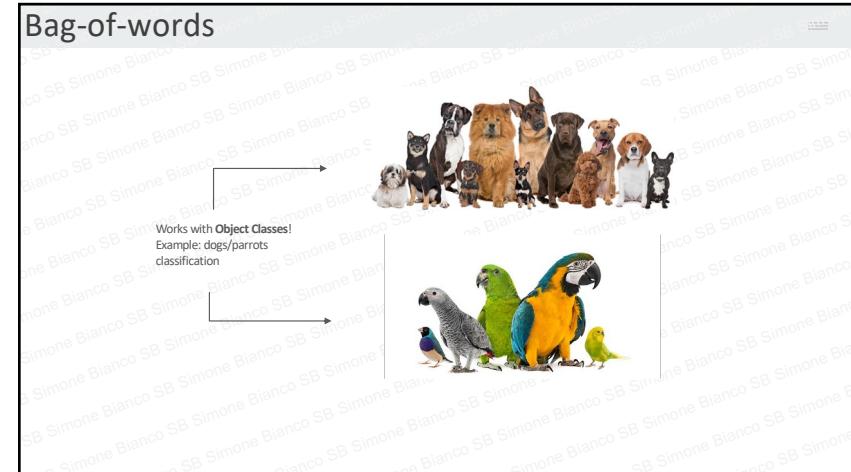
95

Bag-of-words



96

Bag-of-words



97