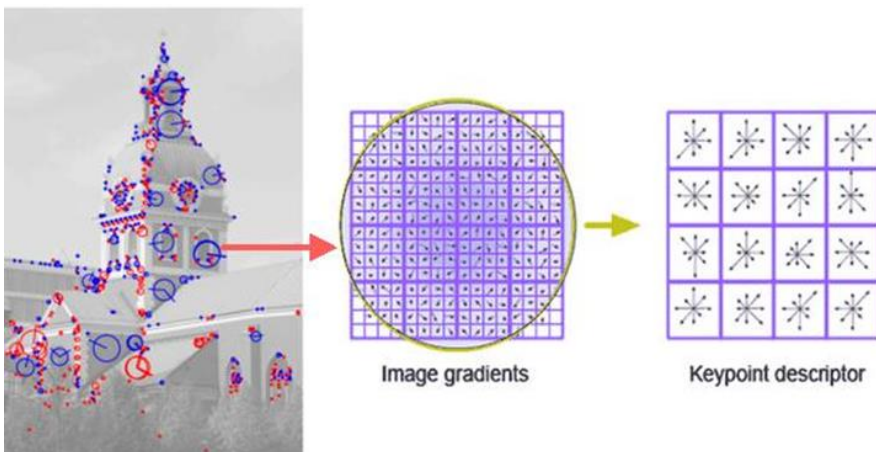# The SIFT feature

Stefano Ghidoni

- SIFT feature in detail
  - Keypoint detection
  - Descriptor calculation
- SIFT performance measurement

- Very reliable keypoint detector and descriptor
- Widely used

Although it's not always the case that a paper cited more contributes more to the field, a highly cited paper usually indicates that something interesting have been discovered. The following are the papers to my knowledge being cited the most in Computer Vision. (updated on 11/24/2013) If you want your "friend's" paper listed here, just comment below.

Cited by 21528 + 6830 (Object recognition from local scale-invariant features)

Distinctive image features from scale-invariant keypoints
DG Lowe - International journal of computer vision, 2004

Cited by 17671

A theory for multiresolution signal decomposition: The wavelet representation
SG Mallat – Pattern Analysis and Machine Intelligence, IEEE ..., 1989

Cited by 17611

A computational approach to edge detection
J Canny – Pattern Analysis and Machine Intelligence, IEEE ..., 1986

Cited by 15422

Snakes: Active contour models
M Kass, A Witkin, Demetri Terzopoulos - International journal of computer ..., 1988
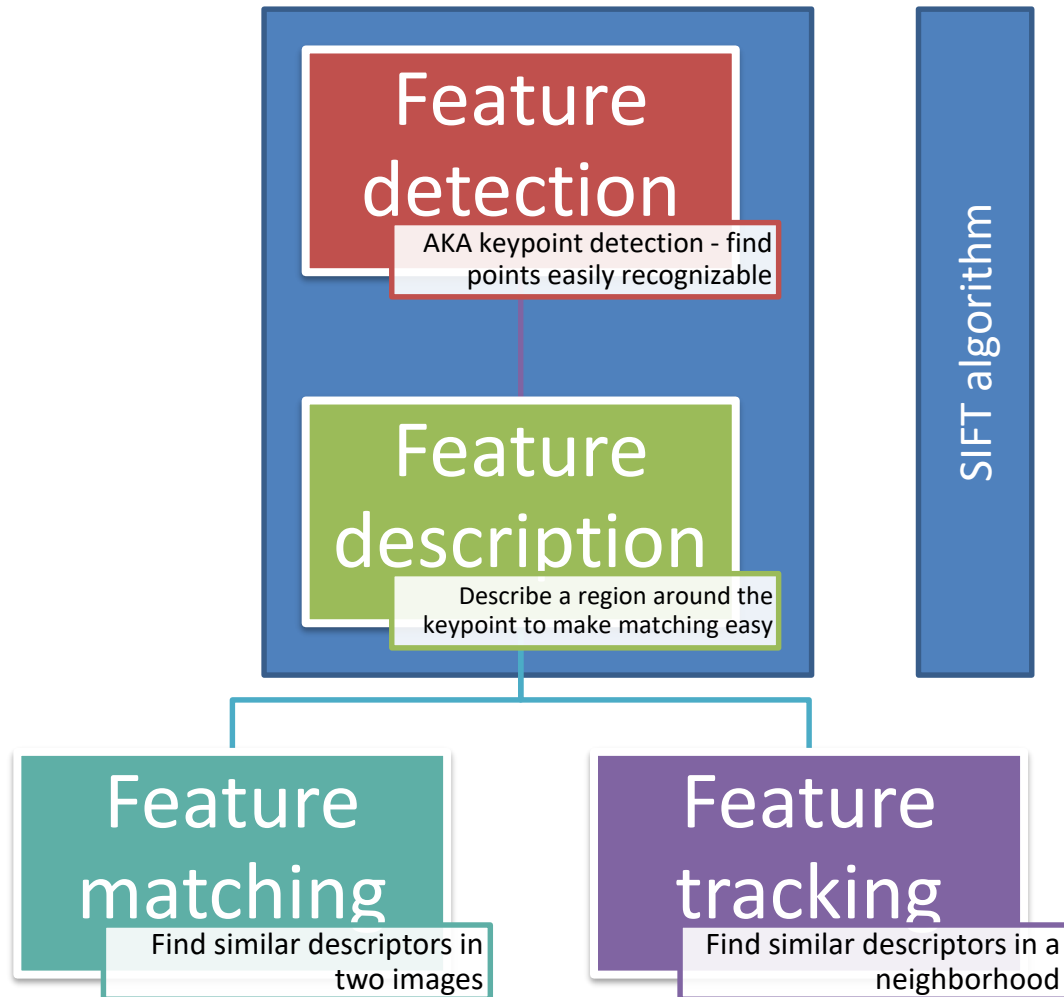
Cited by 15188

Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images
Geman and Geman - Pattern Analysis and Machine ..., 1984

Cited by 11630+ 4138 (Face Recognition using Eigenfaces)

Eigenfaces for Recognition
Turk and Pentland, Journal of cognitive neuroscience Vol. 3, No. 1, Pages 71-86, 1991 (9358 citations)

Cited by 8788

Determining optical flow
B.K.P. Horn and B.G. Schunck, Artificial Intelligence, vol 17, pp 185-203, 1981



Image gradients

Keypoint descriptor

**Feature detection**

AKA keypoint detection - find points easily recognizable

**Feature description**

Describe a region around the keypoint to make matching easy

SIFT algorithm

**Feature matching**

Find similar descriptors in two images

**Feature tracking**

Find similar descriptors in a neighborhood

- Why such a strong focus on scale?

- Different scales
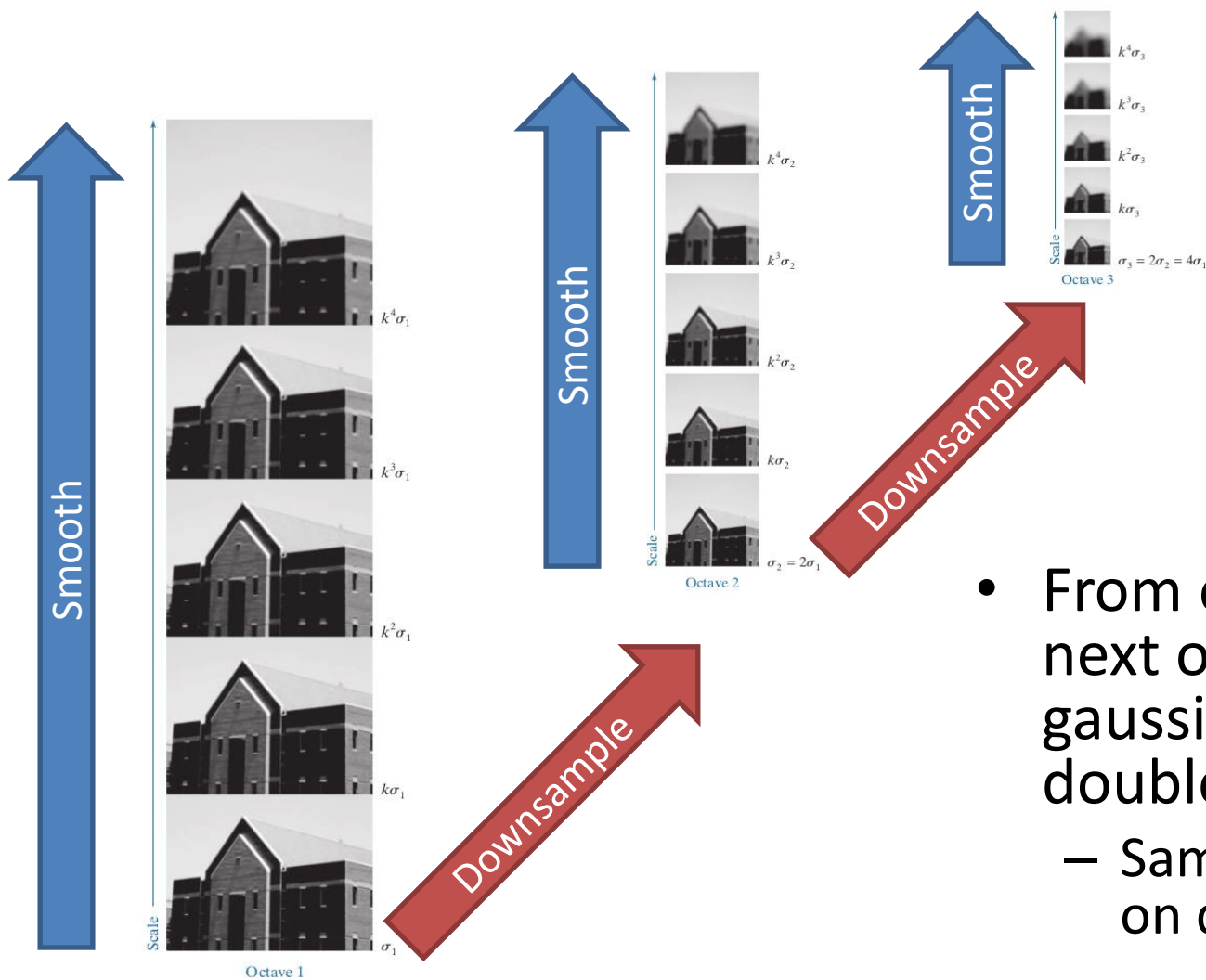- What are the main differences?

- SIFT features:
    - Local – robust to occlusions
    - Distinctive – distinguish objects in large databases
    - Dense – many features can be found even on small objects
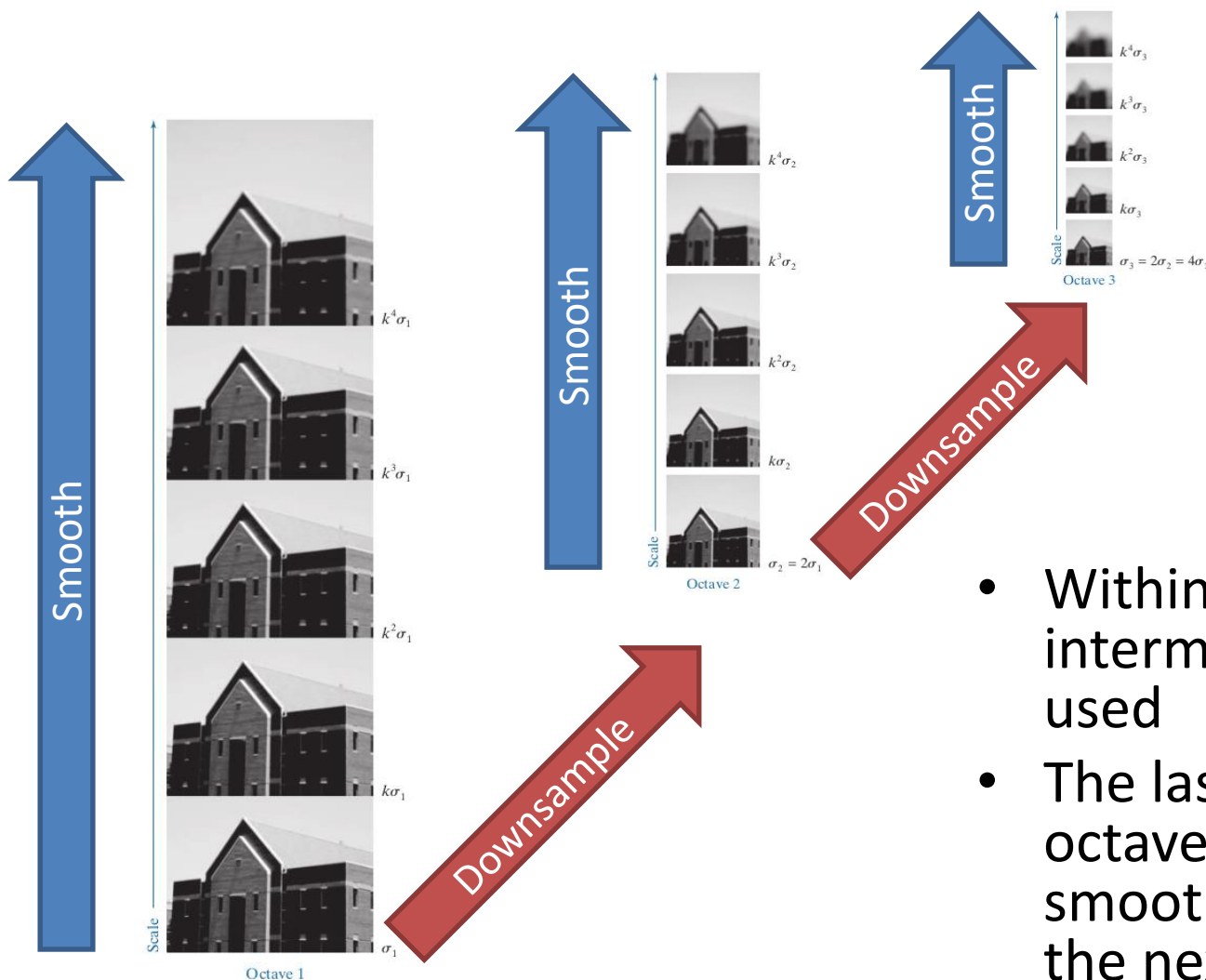    - Efficiency – (rather) fast computation (what's the meaning of fast?)

- Scale-space extrema detection

- Keypoint localization

- Orientation measurement

- Descriptor calculation

D.G. Lowe, "Distinctive image features from scale-invariant keypoints." Int. Journal of Computer Vision, 2004

Smooth

Smooth

Smooth

Downsample

Downsample

$k^4\sigma_1$

$k^3\sigma_1$

$k^2\sigma_1$

$k\sigma_1$

$\sigma_1$

Scale

Octave 1

$k^4\sigma_2$

$k^3\sigma_2$

$k^2\sigma_2$

$k\sigma_2$

$\sigma_2 = 2\sigma_1$

Scale

Octave 2

$k^4\sigma_3$

$k^3\sigma_3$

$k^2\sigma_3$

$k\sigma_3$

$\sigma_3 = 2\sigma_2 = 4\sigma_1$

Scale

Octave 3

- SIFT makes use of a scale space

- The scale space is organized in octaves

- From one octave to the next one the $\sigma$ of the gaussian smoothing is doubled
  - Same smoothing filter on downsampled image

Smooth

Smooth

Smooth

Downsample

Downsample

$k^4\sigma_1$

$k^3\sigma_1$

$k^2\sigma_1$

$k\sigma_1$

$\sigma_1$

Scale

Octave 1

$k^4\sigma_2$

$k^3\sigma_2$

$k^2\sigma_2$

$k\sigma_2$

$\sigma_2 = 2\sigma_1$

Scale

Octave 2

$k^4\sigma_3$

$k^3\sigma_3$

$k^2\sigma_3$

$k\sigma_3$

$\sigma_3 = 2\sigma_2 = 4\sigma_1$

Scale

Octave 3

- Within one octave intermediate scales are used
- The last image of one octave has the same smoothing as the first in the next octave

Within one octave:
- s intervals
  - s+1 images
- Standard deviations
  - $\sigma, k\sigma, \mathrm{k}^2\sigma, \dots, k^s\sigma$

Merging of two octaves:
$$k^s \sigma = 2\sigma \; \rightarrow k = 2^{1/s}$$

- E.g. for $s$=2 (3 images)
$$\sigma, \sqrt{2}\sigma, 2\sigma$$

- Layers in the scale space are combined by subtracting consecutive layers

L – gaussian filtered image
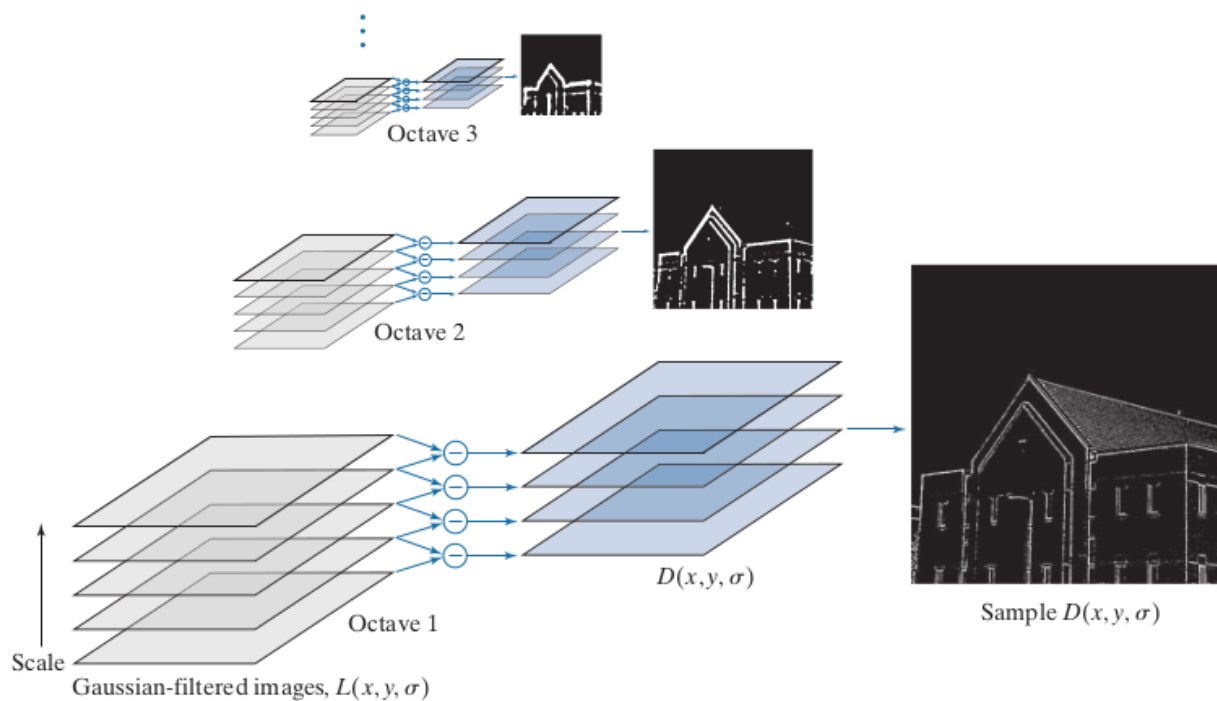D – difference of gaussian filtered images



Octave 3

Octave 2

$D(x, y, \sigma)$

Octave 1

Scale

Gaussian-filtered images, $L(x, y, \sigma)$

- This requires two additional images in the scale space to process the first and last image



Octaves for $s = 2$

| Octave | Scale | | | | |
|--------|-------|-------|-------|-------|--------|
|        | 1     | 2     | 3     | 4     | 5      |
| 1      | 0.707 | 1.000 | 1.414 | 2.000 | 2.828  |
| 2      | 1.414 | 2.000 | 2.828 | 4.000 | 5.657  |
| 3      | 2.828 | 4.000 | 5.657 | 8.000 | 11.314 |

- Output images



Octave 3

Octave 2

Octave 1

$D(x, y, \sigma)$

Sample $D(x, y, \sigma)$

Scale

Gaussian-filtered images, $L(x, y, \sigma)$

- Subtracting consecutive layers means: evaluate the difference between two smoothed images
  - Same smoothing, different smoothing intensity
- Such filtering is called Difference of Gaussians (DoG) and is represented by the function
$$D(x, y, \sigma)$$

- Subtracting consecutive layers means: evaluate the difference between two smoothed images
  - Same smoothing, different smoothing intensity
- What is the meaning of this filter?
  - Let's go one step back: derivative filters

- Recall – we observed many times the pattern:
  - Smoothing
  - Derivative filter (edge detection filter)
- They can be combined into one filter
  - Derivative of the smoothing filter
- This is the concept exploited in the Laplacian of Gaussian (LoG) filter

- Laplacian of a Gaussian (LoG) is obtained using the gaussian filter:

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

And considering its laplacian $\nabla^2 G(x, y)$ to filter the image:

- Filtering with LoG filter corresponds to:
  - Filter the input using $G(x, y)$
  - Compute the Laplacian of the resulting image

- Smoothing
- Noise removal at scales smaller than $\sigma$
- Zero-crossing
- Isotropic
- Typ filter size: $n \times n$ s.t. $n > 6\sigma$

$\nabla^2 G$

x    y

$\nabla^2 G$

Zero crossing    Zero crossing

$2\sqrt{2}\sigma$

| 0 | 0 | −1 | 0 | 0 |
|---|---|---|---|---|
| 0 | −1 | −2 | −1 | 0 |
| −1 | −2 | 16 | −2 | −1 |
| 0 | −1 | −2 | −1 | 0 |
| 0 | 0 | −1 | 0 | 0 |

-LoG

- An approximation: DoG (Difference of Gaussians) $-\sigma_1 > \sigma_2$

$$D(x,y) = \frac{1}{2\pi\sigma_1^2} e^{-\frac{x^2+y^2}{2\sigma_1^2}} - \frac{1}{2\pi\sigma_2^2} e^{-\frac{x^2+y^2}{2\sigma_2^2}}$$

a  b

**FIGURE 10.23**
(a) Negatives of the LoG (solid) and DoG (dotted) profiles using a standard deviation ratio of 1.75:1.
(b) Profiles obtained using a ratio of 1.6:1.

- Back to SIFT: scale space + subtraction of consecutive layers is a DoG filter
  - Output images:

- Scale-space extrema detection

- Keypoint localization

- Orientation measurement

- Descriptor calculation

- Scale-space extrema detection
- Keypoint localization
- Orientation measurement
- Descriptor calculation

- Search for maxima and minima of the DoG

- Comparison with the 8-neighbors in the current, previous and next scale level

- This is done at **all scales**
  - This means: **scale independence**



Scale

Corresponding sections of three contiguous $D(x, y, \sigma)$ images

- Sub-pixel accuracy by means of interpolation
  - Taylor series expansion up to the quadratic term
- Interpolation along $x, y, \sigma$

Scale

Corresponding sections of three
contiguous $D(x, y, \sigma)$ images

- Recall: SIFT is based on the laplacian – 2nd order derivative

- The Hessian matrix provides all the 2nd order derivatives

$$H = \begin{bmatrix} \partial^2 D/\partial x^2 & \partial^2 D/\partial x \partial y \\ \partial^2 D/\partial y \partial x & \partial^2 D/\partial y^2 \end{bmatrix} = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{yx} & D_{yy} \end{bmatrix}$$

- Further processing is based on $H$

Similar to the auto-correlation matrix A

- The Hessian matrix is used to:
  - Discard points with a low gradient
  $$|D(\hat{x})| < 0.03$$
  - Discard edge points – require that both eigenvalues of $H$ are large
    - Means: require that both curvatures are high
    - Similar to the discussion about Harris corners
    - Similar considerations about eigenvalues ratios

- ## More scale levels considered
  - – More points
  - – More unstable
  - – Best value: s=3

- The output is a set of keypoints with associated scales

- Scale-space extrema detection
- Keypoint localization
- Orientation measurement
- Descriptor calculation

- Scale-space extrema detection
- Keypoint localization
- Orientation measurement
- Descriptor calculation

- Each keypoint comes with its scale
- The gaussian smoothed image closest to that scale is selected (L)
  - This process is **independent from the scale**
- Compute gradient magnitude and orientation in the keypoint neighborhood using L

- Gradient magnitude

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

- Gradient orientation angle

$$\theta(x,y) = \tan^{-1}\left[\frac{L(x,y+1) - L(x,y-1)}{L(x+1,y) - L(x-1,y)}\right]$$
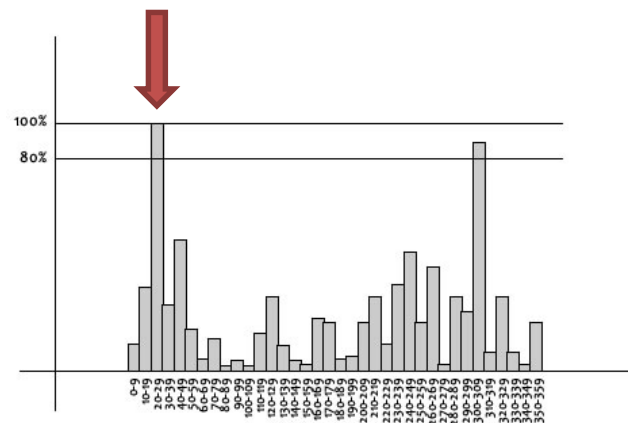
- To achieve **rotation invariance**, the dominant local direction shall be measured

- Build the histogram of gradient orientations
  - 36 bins of 10° each - region around the keypoint

- To achieve **rotation invariance**, the dominant local direction shall be measured
- Build the histogram of gradient orientations
  - 36 bins of 10° each - region around the keypoint
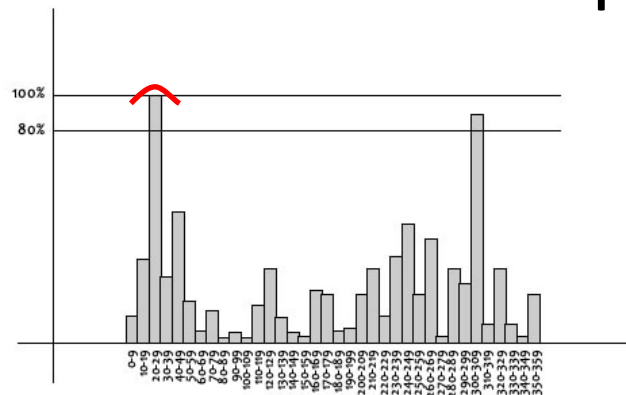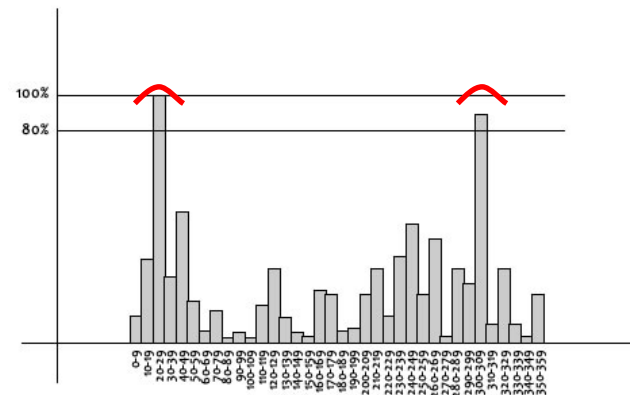- Look for the peak

- To achieve **rotation invariance**, the dominant local direction shall be measured

- Build the histogram of gradient orientations
    - 36 bins of 10° each - region around the keypoint

- Look for the peak

- Fit a parabola to the 3 bins close to the peak
    - Refine the peak location

- Secondary peaks are also considered
  - Peak value > 80% highest peak

- Secondary peaks duplicate the keypoint
  - A new keypoint is created, with the orientation set to the secondary peak

- # Keypoint detection example
  - ## Arrows represent keypoint orientation



*(a) 233x189 image*
*(b) 832 DOG extrema*
*(c) 729 left after peak value threshold*
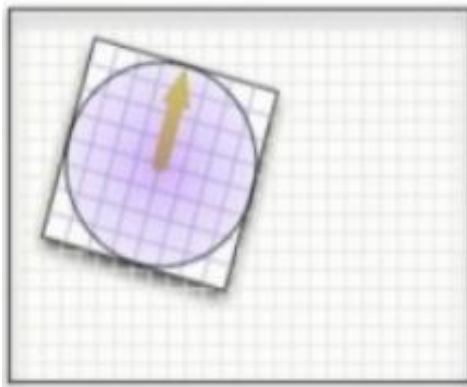*(d) 536 left after testing ratio of principle curvatures*

- Scale-space extrema detection
- Keypoint localization
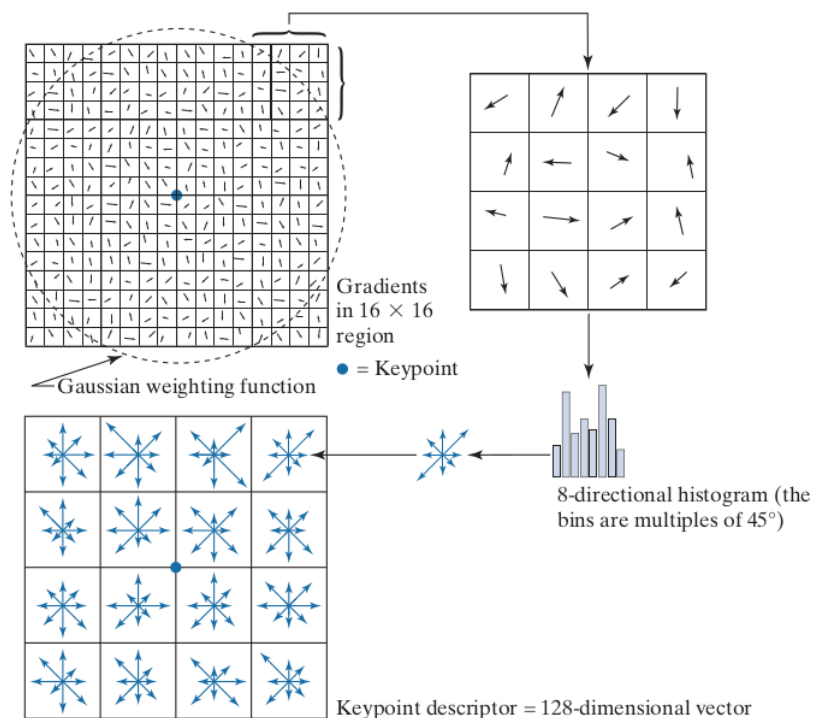- Orientation measurement
- Descriptor calculation

- Scale-space extrema detection

- Keypoint localization
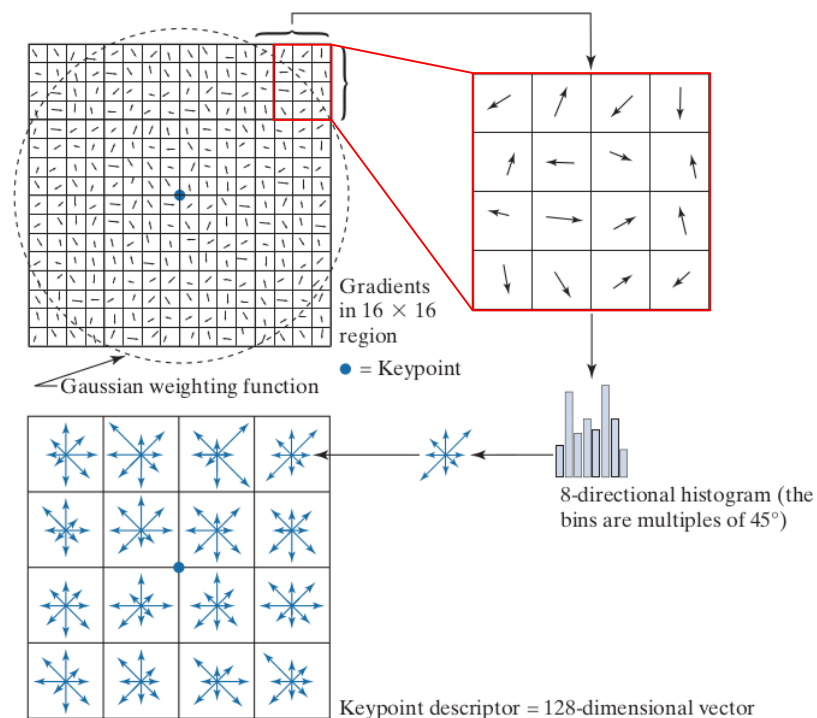
- Orientation measurement

- Descriptor calculation

- SIFT evaluates a descriptor for each keypoint
  - In the image $L$ corresponding to the keypoint scale
  - Considering coordinates that are rotated based on the measured keypoint orientation
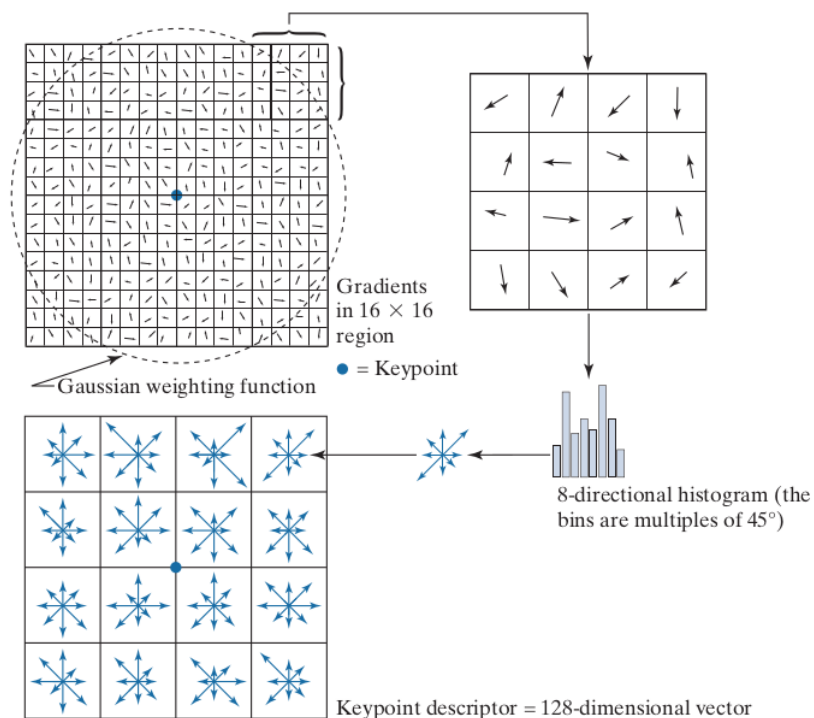
- Compute gradient magnitude and direction
  - neighborhood: 16×16
- Values weighted with a gaussian function centered on the keypoint



Gradients in 16 × 16 region

● = Keypoint

Gaussian weighting function

8-directional histogram (the bins are multiples of 45°)

Keypoint descriptor = 128-dimensional vector

- 16x16 neighborhood divided into regions of 4×4 pixels

- Evaluate histogram of each region
  - 8 bins (45° each)

- Save the values



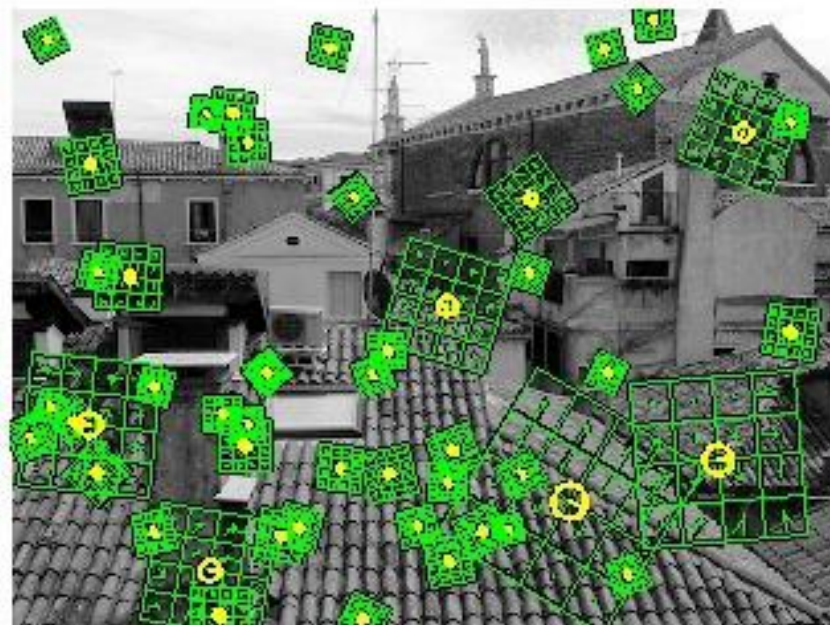Gradients in 16 × 16 region

Gaussian weighting function

● = Keypoint

8-directional histogram (the bins are multiples of 45°)

Keypoint descriptor = 128-dimensional vector

- Descriptor size:
  - 16 regions
  - 8 histogram values per region
- 128 elements
  - 1 byte per element



Gradients in 16 × 16 region

Gaussian weighting function

• = Keypoint

8-directional histogram (the bins are multiples of 45°)

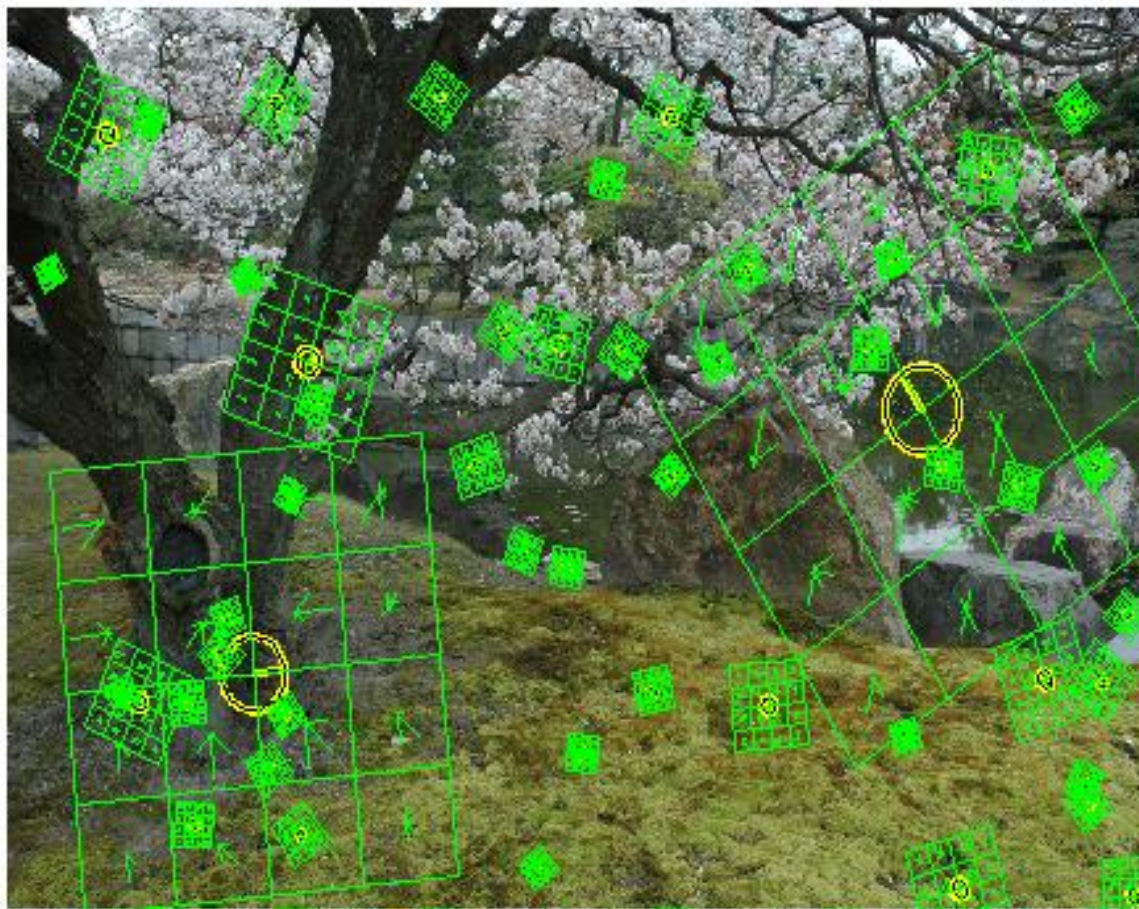Keypoint descriptor = 128-dimensional vector

- Additional details are considered by the SIFT algorithm

- Enhanced **invariance to illumination**: feature vector normalized to unit length

- Enhanced **robustness to large gradient magnitudes**: threshold to 0.2 on all the components
  - Nonlinear illumination effect – e.g., camera saturation
  - Further normalization to 1 after thresholding

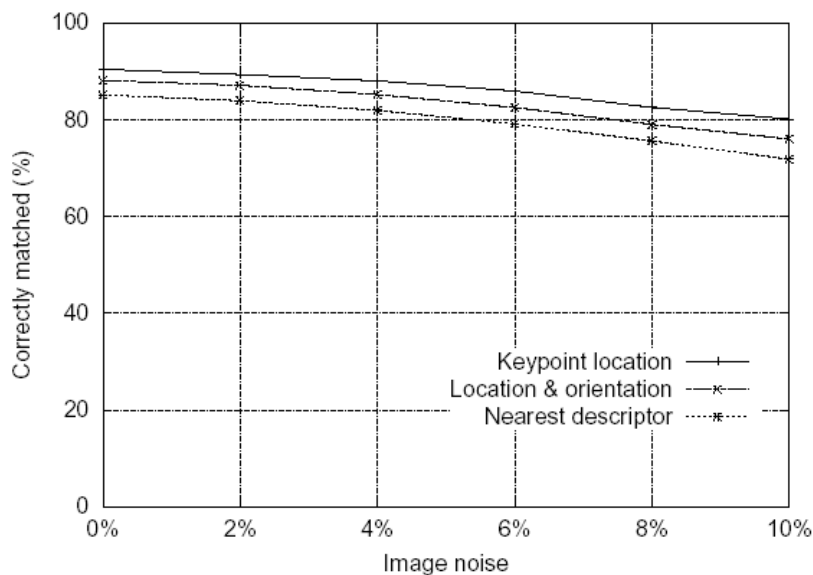- ## SIFT (VLFeat library)

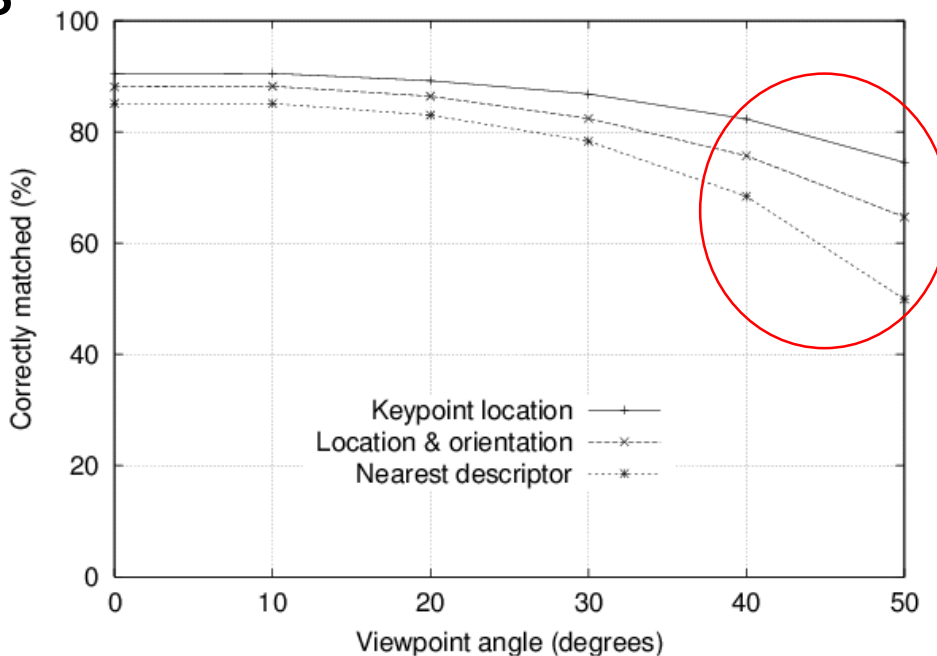- SIFT (VLFeat library)

- SIFT (OpenCV)

- The robustness of the SIFT features to noise factor was accurately measured

- Random rotation and scale change with different noise levels
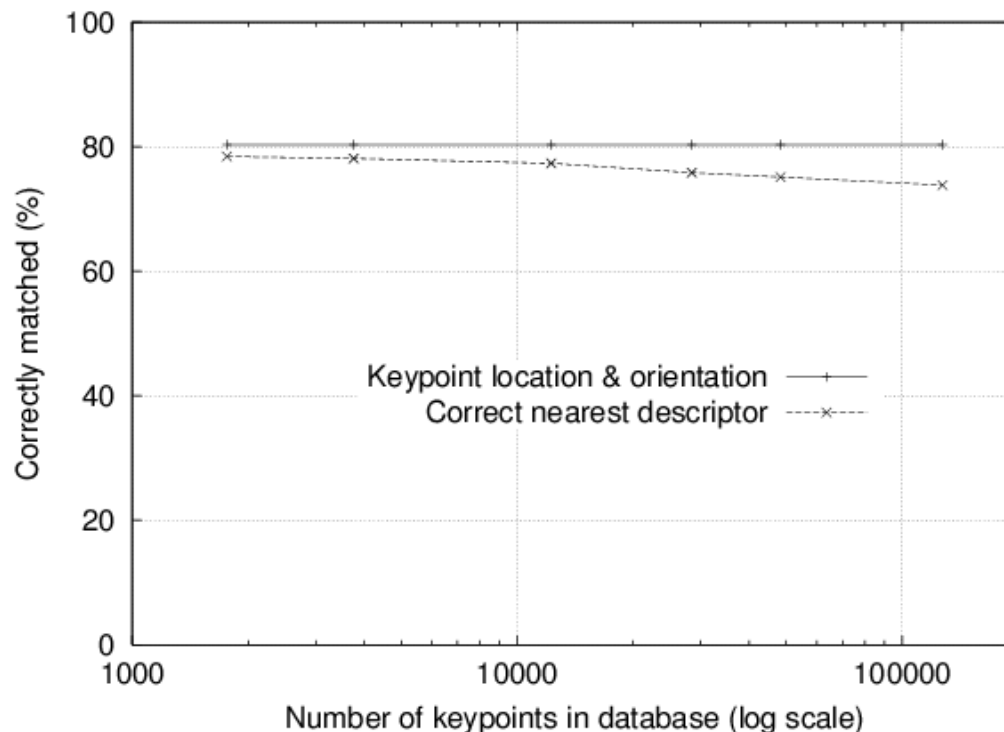  - Nearest neighbor on 30,000 feature database

- Random scale change + rotation with 2% noise and affine distortion

  – Nearest neighbor on 30,000 feature database

- Strong for small angles

- ## 2% Noise and 30° viewpoint change

- Why is SIFT invariant to scale?

- Why is SIFT invariant to illumination changes?

- Why is SIFT invariant to orientation?

# The SIFT feature

Stefano Ghidoni

DIPARTIMENTO
DI INGEGNERIA
DELL'INFORMAZIONE