# Keypoint detection and description

## Javier González Jiménez

Reference Books:

- *Computer Vision: Algorithms and Applications*. Richard Szeliski. Springer. 2010.

http://szeliski.org/Book

# Content

1. Introduction
2. Harris detector
   - Idea
   - Formulation
   - Implementation
3. KLT operator
4. Keypoint matching through correlation
5. SIFT operator
   - Scale Space
   - Detector
   - Descriptor

# 1. Introduction

What are *keypoints* (also, *interest* or *feature points*)?

Distinctive pixels in the image that can likely be projections on 3D entities.



Synthetic image
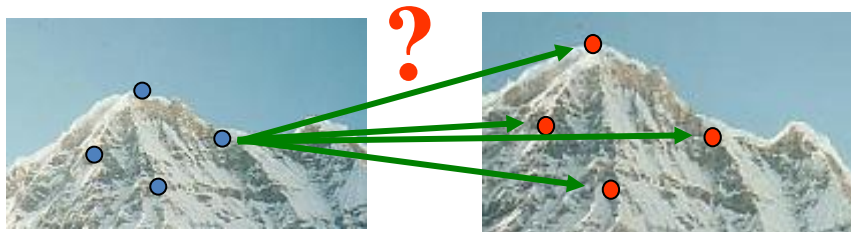


Real image

# 1. Introduction

Two problems involved: Detection and description

1. Detection



The operator must provide a reliable and repeatable response

2. Description: To match to its correspondence in other images



Invariant and discriminative description needed

**Harris** is just a detector (not a descriptor). The keypoint is described with a surrounding image patch
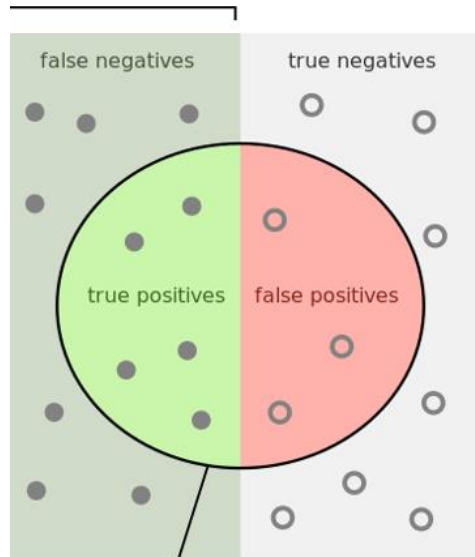**SIFT** is a detector and also provides its own descriptor

# 1. Introduction

## Desired properties:

**Detector:**
- Accurate (subpixel accuracy, if possible)
- Detect all the keypoints in the image (few FALSE NEGATIVES)
- NOT detect irrelevant points (few FALSE POSITIVES)

Keypoints in the image
(Ground truth)

false negatives

true negatives

true positives

false positives

Detected keypoints
(Prediction)

few FALSE NEGATIVES ➡ Wanted: High RECALL (to improve the SENSITIVITY of the detector)

$$Recall = \frac{}{} = \frac{TP}{TP + FN}$$

actually positive

Among all matches that are actually positive, what percentage is predicted positive

few FALSE POSITIVES ➡ Wanted: High PRECISION (to reduce the false alert rate)

$$Precision = \frac{}{} = \frac{TP}{TP + FP}$$

predicted positive

Among all matches that were predicted positive, what percentage is predicted positive

# Desired properties:

**Descriptor** should have *invariance* to:

&#10003; Illumination (changes in brightness and contrast)



&#10003; View point (scale, rotación, projective distortion)



**Detector and descriptor:** Computationally efficient and robust to image noise

# 1. Introduction

Applications

- Panoramic images (image sticking)
- 3D Reconstruction
- Object tracking
- Object recognition
- Image retrieval and indexing in database
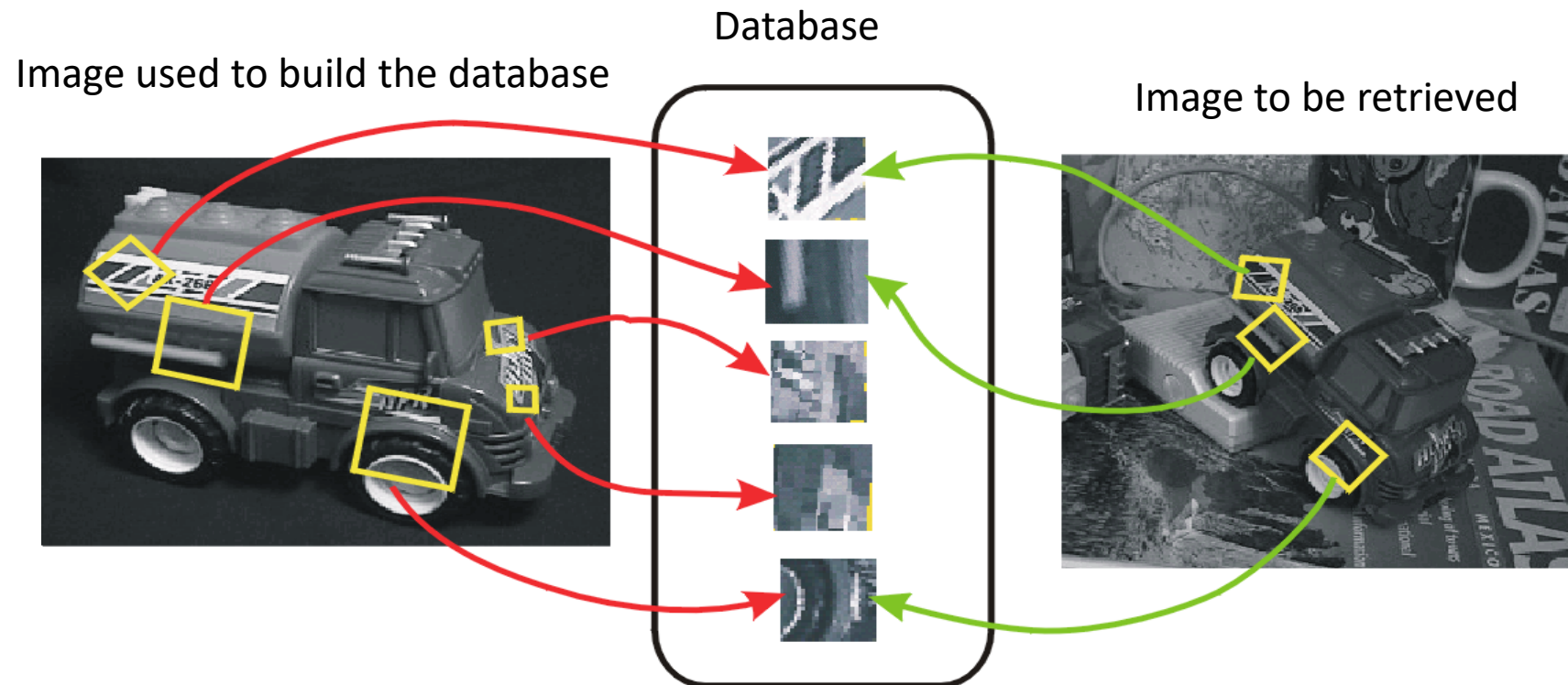- Robot navigation: mapping, localization, obstacle detection
- .... and many others

# 1. Introduction

**Applications:** Panoramic images (image stitching)



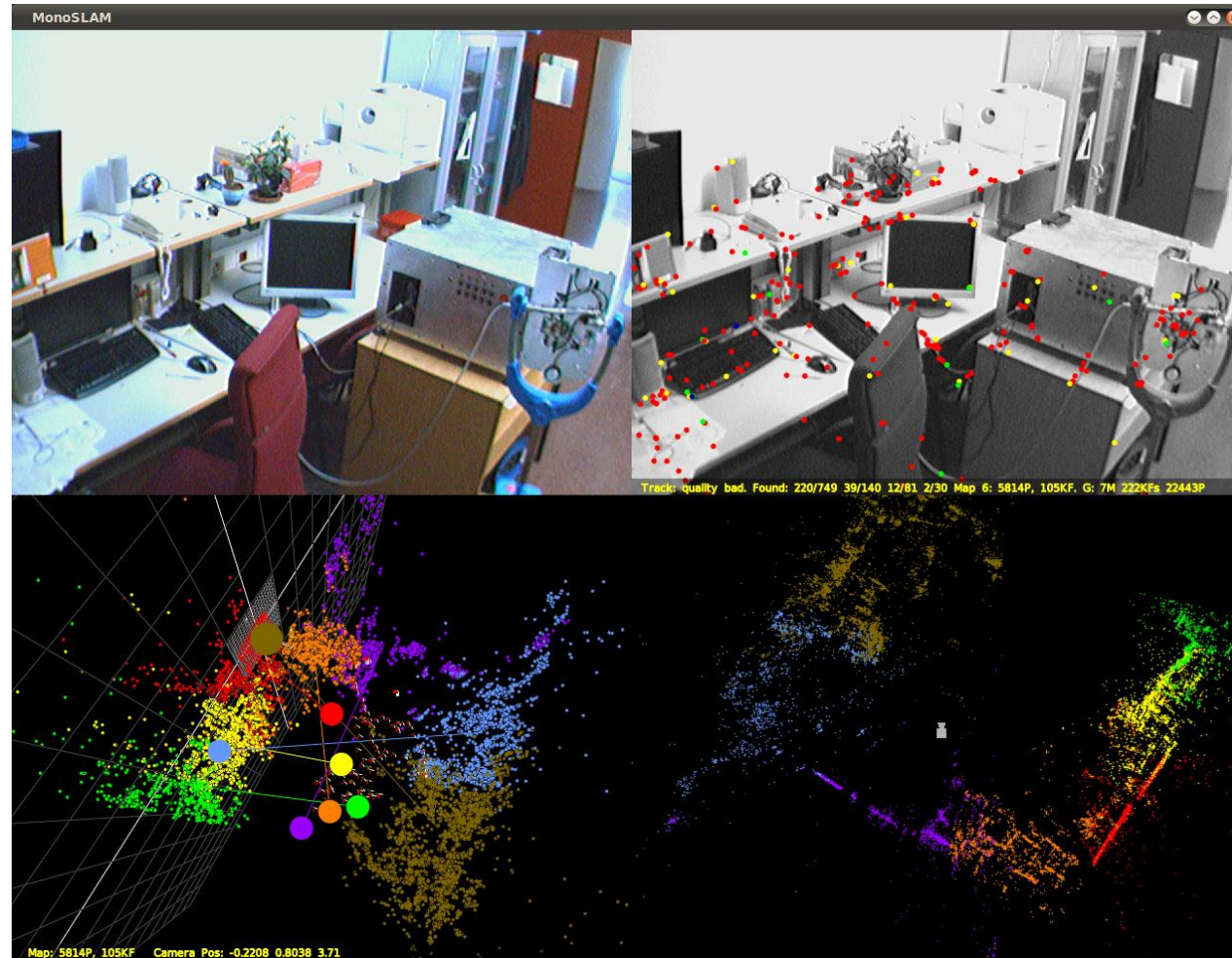M. Brown and D. G. Lowe. Recognising Panoramas. ICCV 2003

# 1. Introduction

Applications: Image retrieval and indexing in database

# 1. Introduction

**Applications:** Localization and 3D Reconstruction
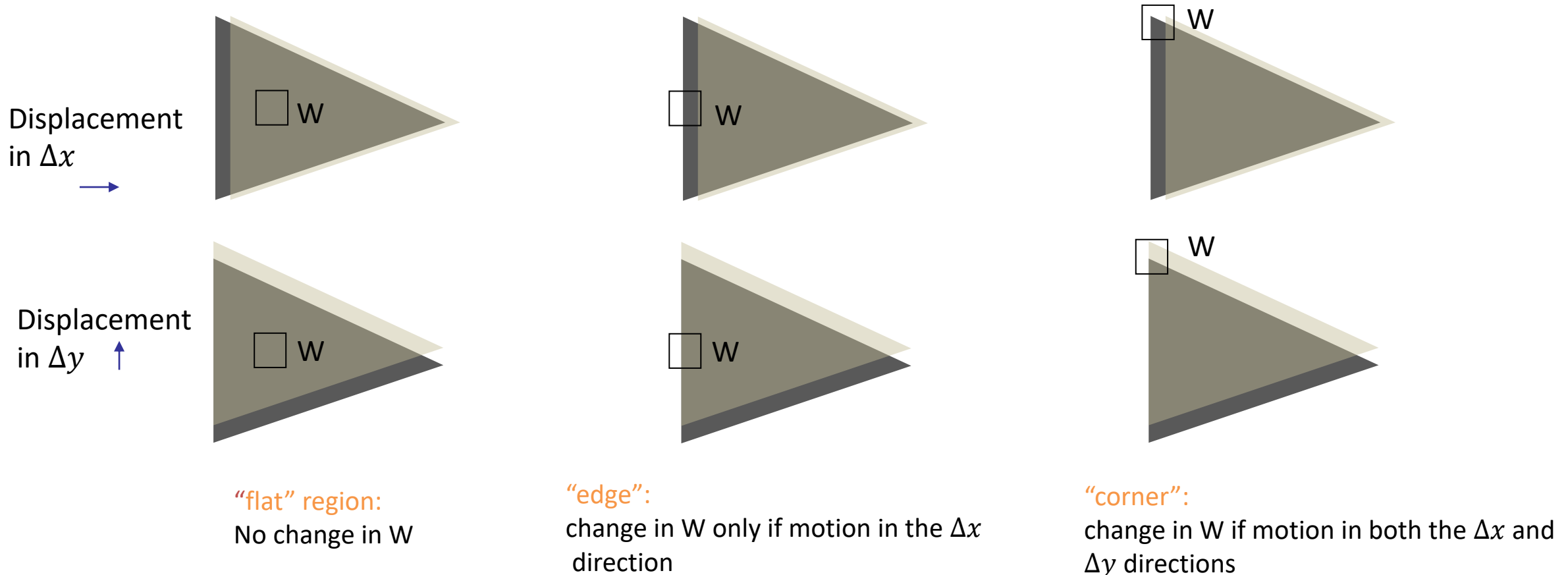
# 2. Harris detector

- Detect corners

- Simple and efficient implementation

- Robustness to noise (apply smoothing)

- Invariance to
  - Rotation: uses eigenvectors
  - Brightness (partially to contrast): uses derivatives  ☺

- **Not invariance to scale**  ☹

*C.Harris, M.Stephens. "A Combined Corner and Edge Detector".*
Proceedings of The Fourth Alvey Vision Conference, *1988*

# 2. Harris detector

- A keypoint is a corner ("esquina")
- A corner is a point with high variation of intensity in 2 spatial directions

**Basic idea**: look in a small window W around a pixel if the displacements of the image in two directions provoke changes

Displacement in $\Delta x$ →

Displacement in $\Delta y$ ↑

W

W

W

W

W

W

W

W

"flat" region:
No change in W

"edge":
change in W only if motion in the $\Delta x$ direction

"corner":
change in W if motion in both the $\Delta x$ and $\Delta y$ directions

# Detecting the local change in intensity due to a shift $(\Delta x, \Delta y)$:

Sum-of-square weighted difference at a pixel $[x_0 y_0]$ when a window image $I$ is shifted $(\Delta x, \Delta y)$:

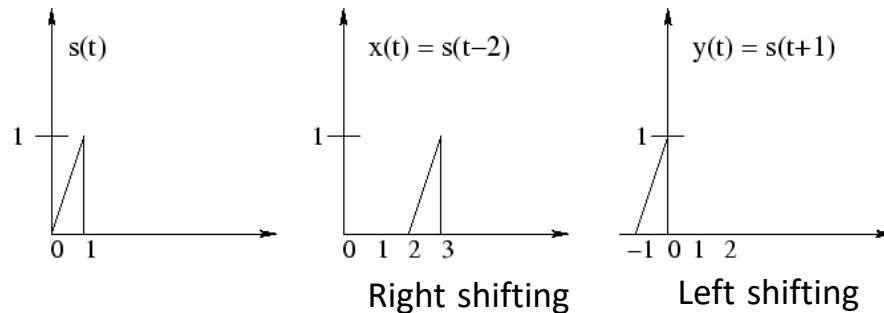$$E_{x_o y_0}(\Delta x, \Delta y) = \sum_{x,y} w(x,y)[I(x + \Delta x, y + \Delta y) - I(x,y)]^2$$

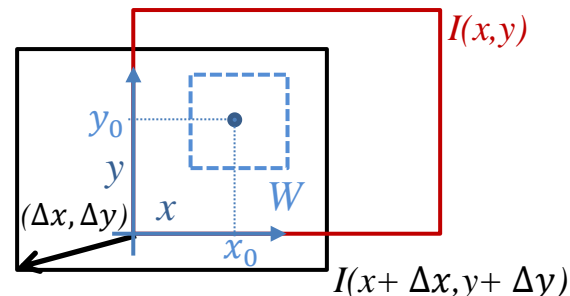Weighting window centered at $[x_0 y_0]$

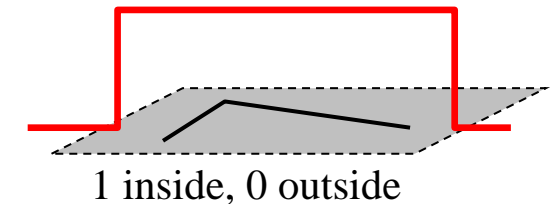Image shifted $(\Delta x, \Delta y)$

Image

**RECALL:**

1D signal shifting



s(t)

1

0 1

x(t) = s(t−2)

1

0 1 2 3

Right shifting

y(t) = s(t+1)

1

−1 0 1 2

Left shifting

2D Left-down image shifting



$I(x,y)$

$y_0$

$y$

$(\Delta x, \Delta y)$    $x$    $W$

$x_0$

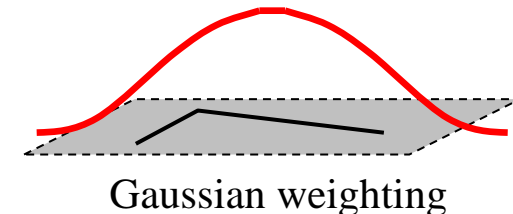$I(x+ \Delta x, y+ \Delta y)$

Weighting function *w(x,y)* may be:



Boolean

1 inside, 0 outside

Gaussian

Gaussian weighting

# Let's make more practical the computation of $E_{x_0 y_0}(\Delta x, \Delta y)$

$$E_{x_0 y_0}(\Delta x, \Delta y) = \sum_{x,y} w(x,y)[I(x + \Delta x, y + \Delta y) - I(x,y)]^2 = \sum_{(x_i,y_i)\in W} [I(x_i + \Delta x, y_i + \Delta y) - I(x_i, y_i)]^2$$

Sum only over a boolean window $W$

Image derivative at $(x_i, y_i)$ along the $y$ axis

First order Taylor approximation

$$E(\Delta x, \Delta y) \approx \sum_{(x_i,y_i)\in W} \left[ I(x_i, y_i) + [I_x(x_i, y_i) \quad I_y(x_i, y_i)] \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} - I(x_i, y_i) \right]^2$$

$$= [\Delta x \ \Delta y] \sum_{(x_i,y_i)\in W} \begin{bmatrix} (I_x(x_i, y_i))^2 & I_x(x_i, y_i)I_y(x_i, y_i) \\ I_x(x_i, y_i)I_y(x_i, y_i) & (I_y(x_i, y_i))^2 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

$(A^T B)^2 = (A^T B)^T (A^T B)$
$= B^T A A^T B = B^T M B$

$M$

$M$ is computed from the **first derivatives at each image point** $(x_0 y_0)$

$$E_{x_0 y_0}(\Delta x, \Delta y) \approx [\Delta x \ \Delta y] \begin{bmatrix} \sum_W (I_x(x_i, y_i))^2 & \sum_W I_x(x_i, y_i)I_y(x_i, y_i) \\ \sum_W I_x(x_i, y_i)I_y(x_i, y_i) & \sum_W (I_y(x_i, y_i))^2 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} = [\Delta x \ \Delta y] M \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$
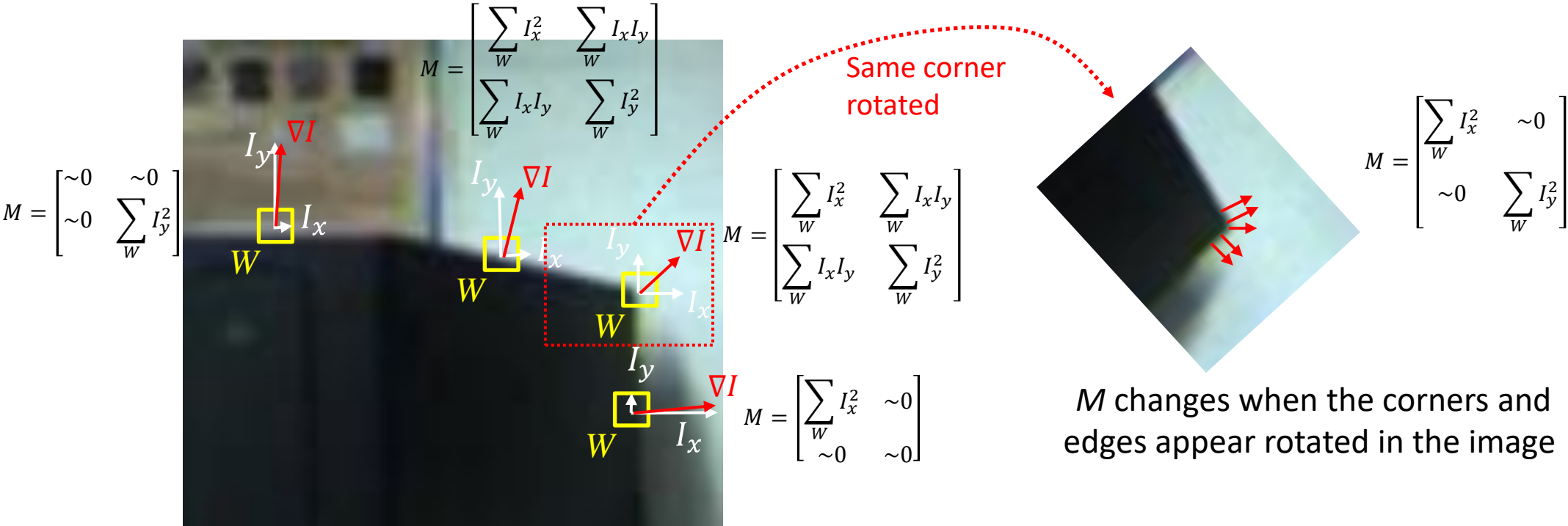
$E_{x_0 y_0}(\Delta x, \Delta y)$ is a **quadratic polynomial** which coefficients are the entries of $M$

# Understanding $M$

Summation of all the square derivatives **along the x-axis** of the points inside $W$

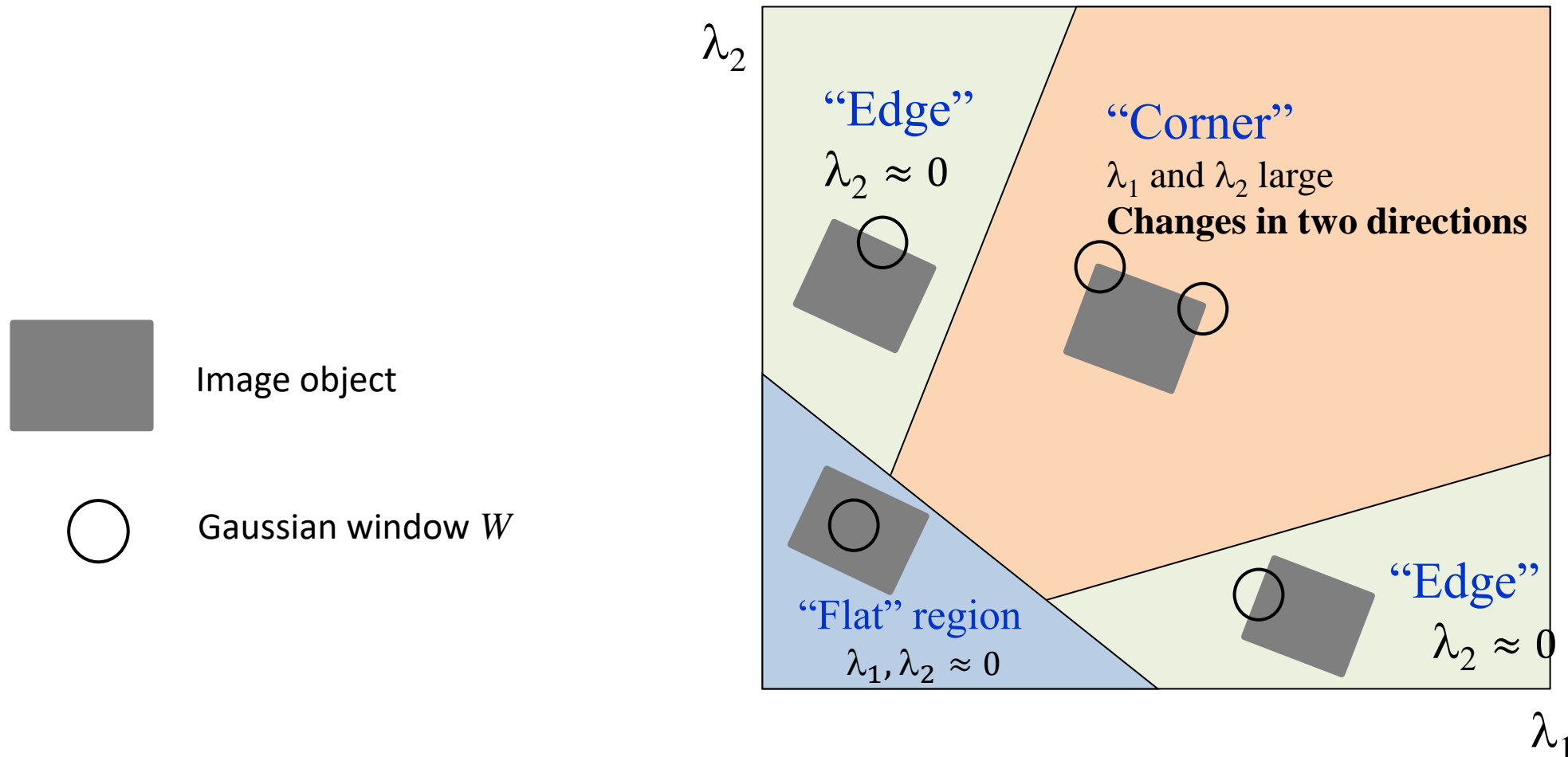$$M = \begin{bmatrix} \sum_W I_x^2 & \sum_W I_x I_y \\ \sum_W I_x I_y & \sum_W I_y^2 \end{bmatrix}$$

$M$ is computed from the **first derivatives at each image point $(x_o y_0)$**



$$M = \begin{bmatrix} \sum_W I_x^2 & \sum_W I_x I_y \\ \sum_W I_x I_y & \sum_W I_y^2 \end{bmatrix}$$

$$M = \begin{bmatrix} \sim0 & \sim0 \\ \sim0 & \sum_W I_y^2 \end{bmatrix}$$

$$M = \begin{bmatrix} \sum_W I_x^2 & \sum_W I_x I_y \\ \sum_W I_x I_y & \sum_W I_y^2 \end{bmatrix}$$

$$M = \begin{bmatrix} \sum_W I_x^2 & \sim0 \\ \sim0 & \sim0 \end{bmatrix}$$

Same corner rotated

$$M = \begin{bmatrix} \sum_W I_x^2 & \sim0 \\ \sim0 & \sum_W I_y^2 \end{bmatrix}$$

$M$ changes when the corners and edges appear rotated in the image

Derivative vector at different pixels (always perpendicular to the border)

# 2. Harris detector

The image points are now classified according to the eigenvalues



Image object

Gaussian window $W$

$\lambda_2$

"Edge"
$\lambda_2 \approx 0$

"Corner"
$\lambda_1$ and $\lambda_2$ large
**Changes in two directions**

"Flat" region
$\lambda_1, \lambda_2 \approx 0$

"Edge"
$\lambda_2 \approx 0$

$\lambda_1$

But computing the eigenvalues of a 2x2 matrix (M) at each pixel is costly!

# 2. Harris detector

Let's define a scalar variable $R$ that indexes the same domain:

$$R = \underbrace{\lambda_1 \lambda_2}_{\text{Determinant}} - k \underbrace{\left( \lambda_1 + \lambda_2 \right)^2}_{\text{Trace}}$$

($k = 0.04\text{-}0.06$ is an empiric constant)

- $R$ is large and positive at corners
- $R$ is negative at edges
- $|R|$ is small at flat regions

Trace and determinant of $M$ and $D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$ are the same:

$$\det M = \lambda_1 \lambda_2$$
$$\operatorname{trace} M = \lambda_1 + \lambda_2$$

$$\boxed{R = \det M - k \left( \operatorname{trace} M \right)^2}$$

No need to compute the eigenvalues!!



$\lambda_2$

"Edge" $R < 0$

"Corner"

$R > 0$

"Flat" $|R|$ small

"Edge" $R < 0$

$\lambda_1$

# 2. Harris detector

## Summary:

Original image



$I_y < 0$

$I_y > 0$

$I_y = 0$
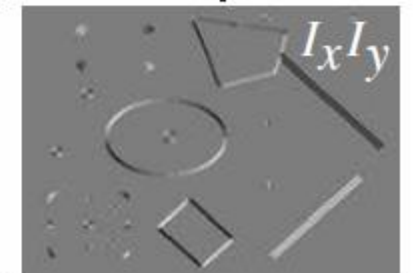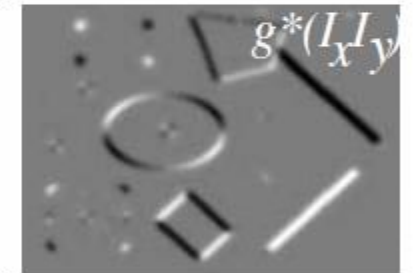
$I_x$

$I_y$

$I_x^2$

$I_y^2$

$I_x I_y$

# 2. Harris detector

Summary:

Original image

# 2. Harris detector

Summary:

Original image

$$\det(M) - \lambda \; \text{trace}(M) =$$

$$M = \begin{bmatrix} g*I_x^2 & g*(I_xI_y) \\ g*(I_xI_y) & g*I_y^2 \end{bmatrix}$$

$I_x$

$I_y$

$g*I_x^2$

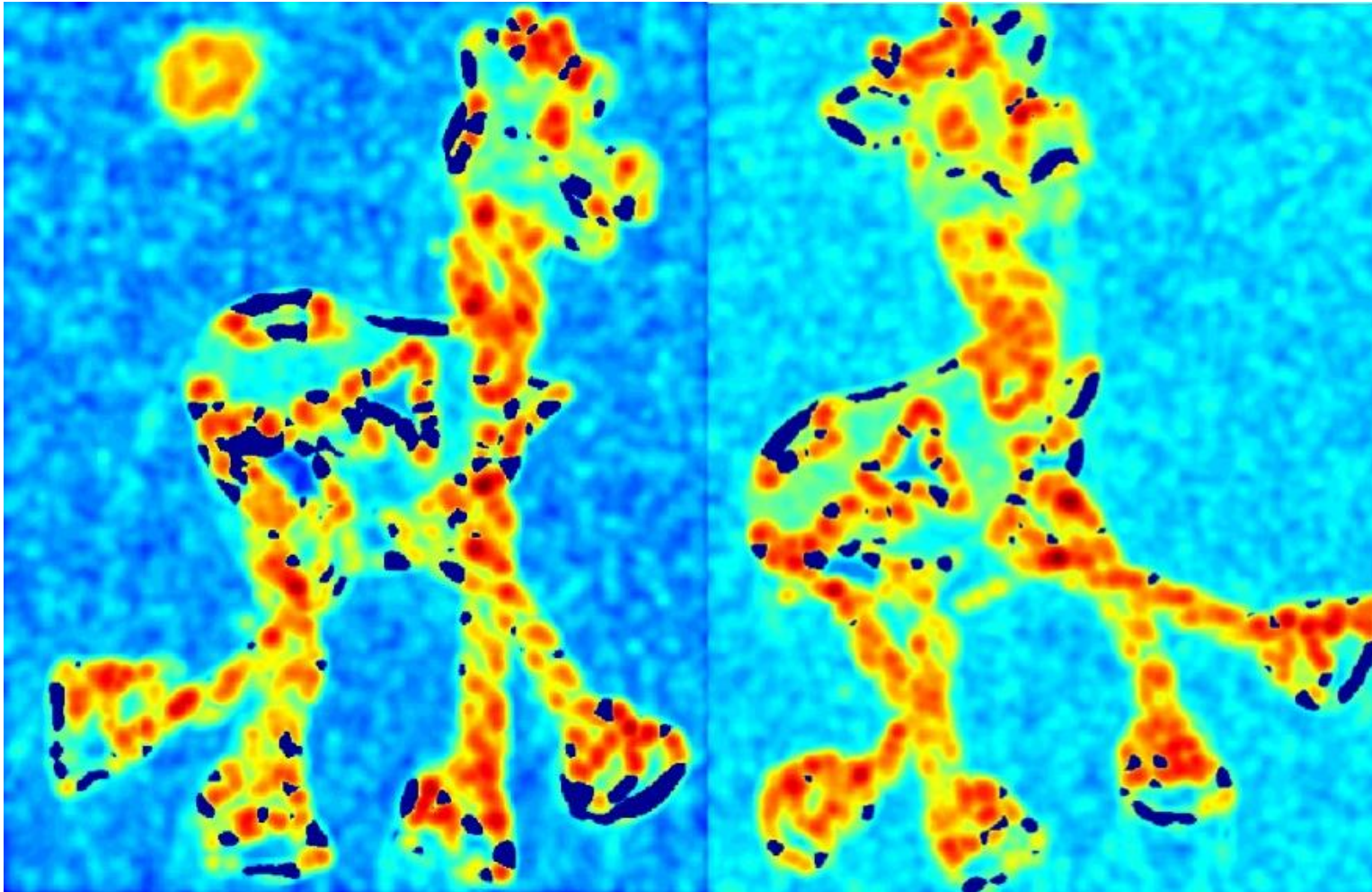$g*I_y^2$

$g*(I_xI_y)$

$I_x^2$

$I_y^2$

$I_xI_y$

# 2. Harris detector

Algorithm:

- Compute the image derivatives $I_x$, $I_y$ (i.e. Sobel)
- Gaussian smoothing of the 3 images: $(I_x)^2, (I_y)^2, I_x I_y$
- Compute the image *R* from the formula (trace and determinant)
- Find pixels where R is high (*R* > threshold)
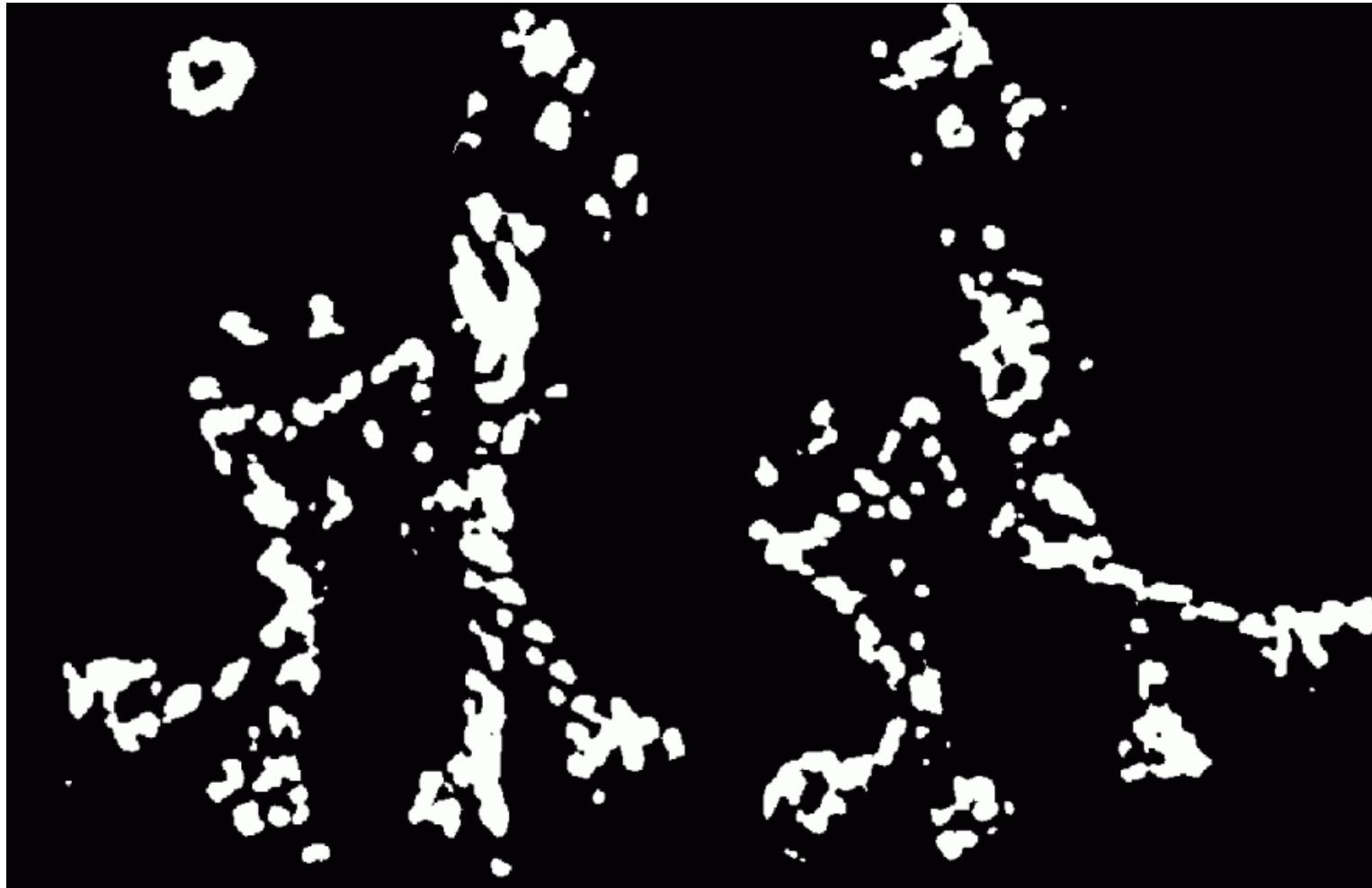- Select local maxima

# 2. Harris detector

# 2. Harris detector

$R$ image

# 2. Harris detector

$R >$ threshold

# 2. Harris detector

Local maxima of $R$

# 2. Harris detector

# 3. KLT operator (Kanade-Lucas-Tomasi)

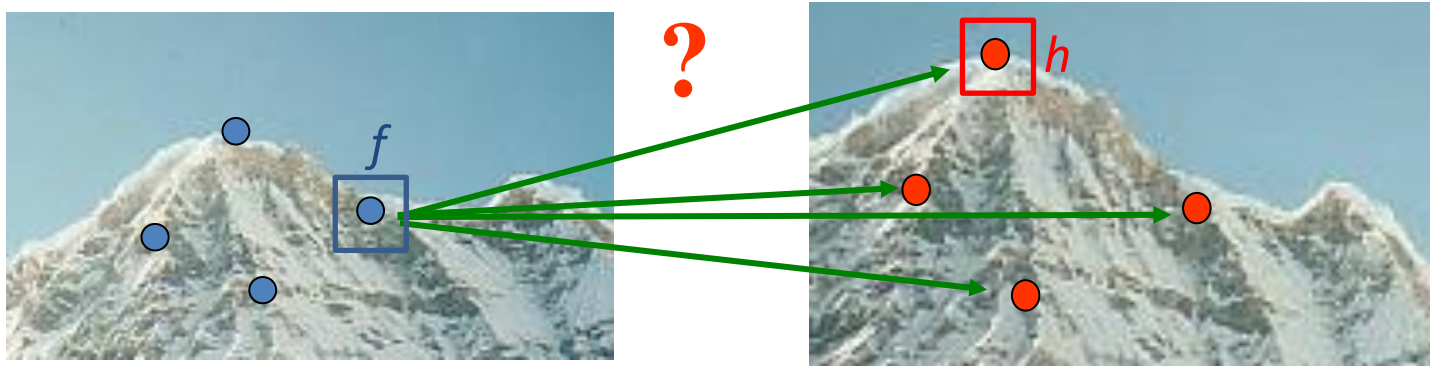**Objective**: Detect distintive points, suitable to be tracked in a image sequence

- Similar principle to Harris: "A good keypoint is that with a high intensity derivative in two directions" → **Min ($\lambda_1$, $\lambda_2$) > threshold**

- Also based on the first derivative matrix:

$$M = \begin{bmatrix} \sum_W I_x^2 & \sum_W I_x I_y \\ \sum_W I_x I_y & \sum_W I_y^2 \end{bmatrix} \implies D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

- But now, the eigenvalues are computed (no approximation with $R$) → better behavior under affine image deformation

# 4. Keypoint matching through correlation

## Which is the correspondence of each point in other image?



*Sum of squared differences (SSD):* $\quad SSD(f,h) = \sum \left[ f(i,j) - h(i,j) \right]^2$

SSD is approximated by the SAD (more efficient computationally):

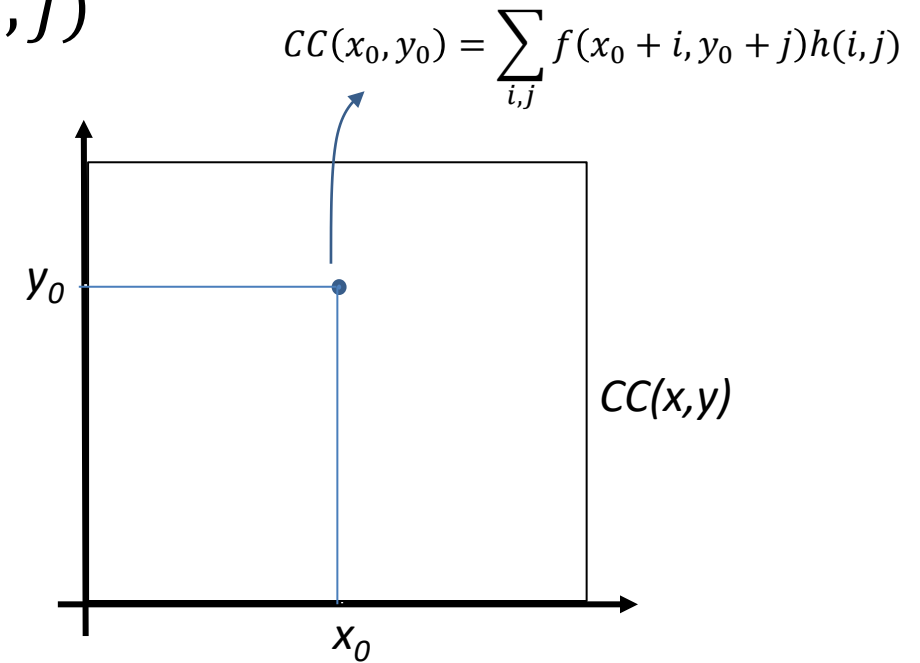*Sum of absolute differences (SAD):* $\quad \sum \left| f(i,j) - h(i,j) \right|$
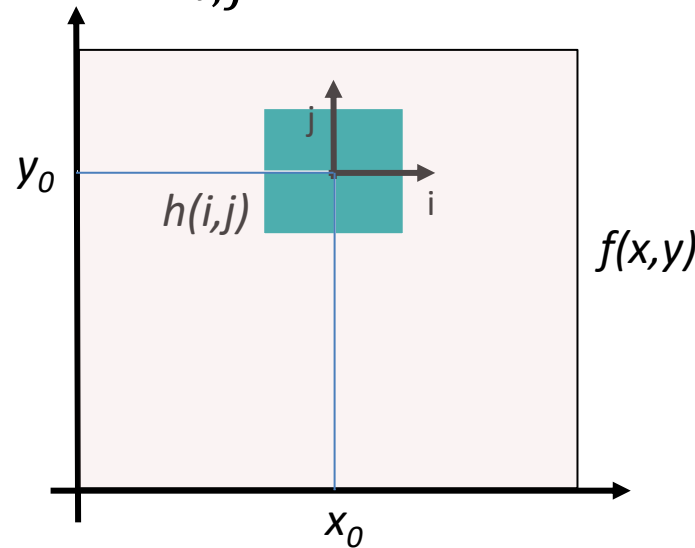
<u>Problem</u>: SSD and SAD are not invariant to brightness or contrast changes

Still, SAD is employed in keypoint tracking along an image sequence (where image brightness and contrast do not change very much.

# 4. Keypoint matching through correlation
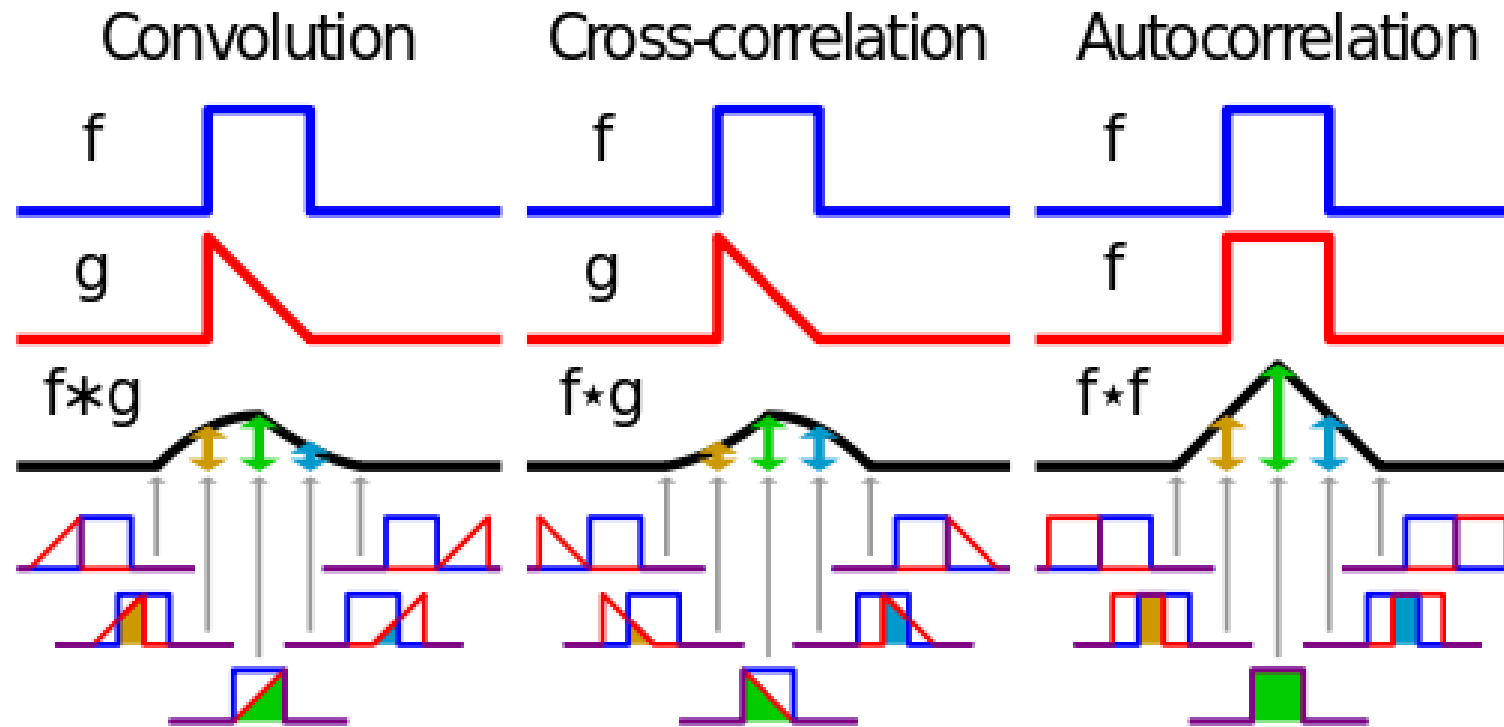
**Cross-Correlation:**

$$CC(x,y) = \sum_{i,j} f(x+i, y+j)h(i,j)$$

$$CC(x_0, y_0) = \sum_{i,j} f(x_0+i, y_0+j)h(i,j)$$



$y_0$

$h(i,j)$

$f(x,y)$

$x_0$

$y_0$

$CC(x,y)$

$x_0$

CC is similar to the convolution but without flipping the kernel

$$f \otimes h = \sum_{x,y} f(x-i, y-j)h(i,j)$$

equivalent to the cross-correlation of $f(-i, -j)$ and $h(i,j)$

Correlation vs. Convolution:



The convolution f * g  is equivalent to the cross-correlation of f(t) and g(−t)

**Normalized Cross-Correlation (NCC):**

Cross correlation is not invariant to change in brightness and contrast of f  y h
→ Normalization required

$$NCC(x, y) = \sum_{i,j} \hat{f}(x+i, y+j)\hat{h}(i, j)$$

**Normalization**: $\hat{f}$ and $\hat{h}$ have <span style="color:red">zero mean</span> and <span style="color:red">contrast one</span>

$$\hat{f}(x+i, y+j) = \frac{f(x+i, y+j) - \bar{f}}{\left\|f - \bar{f}\right\|_{W_m(x,y)}}$$  brightness

$$\hat{h}(x, y) = \frac{h(x, y) - \bar{h}}{\left\|h - \bar{h}\right\|_{W_m(x,y)}}$$  brightness

Mean brightness (intensity) of f and h in the window $W_m$

$$\bar{f} = \frac{1}{\left|W_m(x, y)\right|} \sum_{(i,j)\in W_m(x,y)} f(i, j)$$

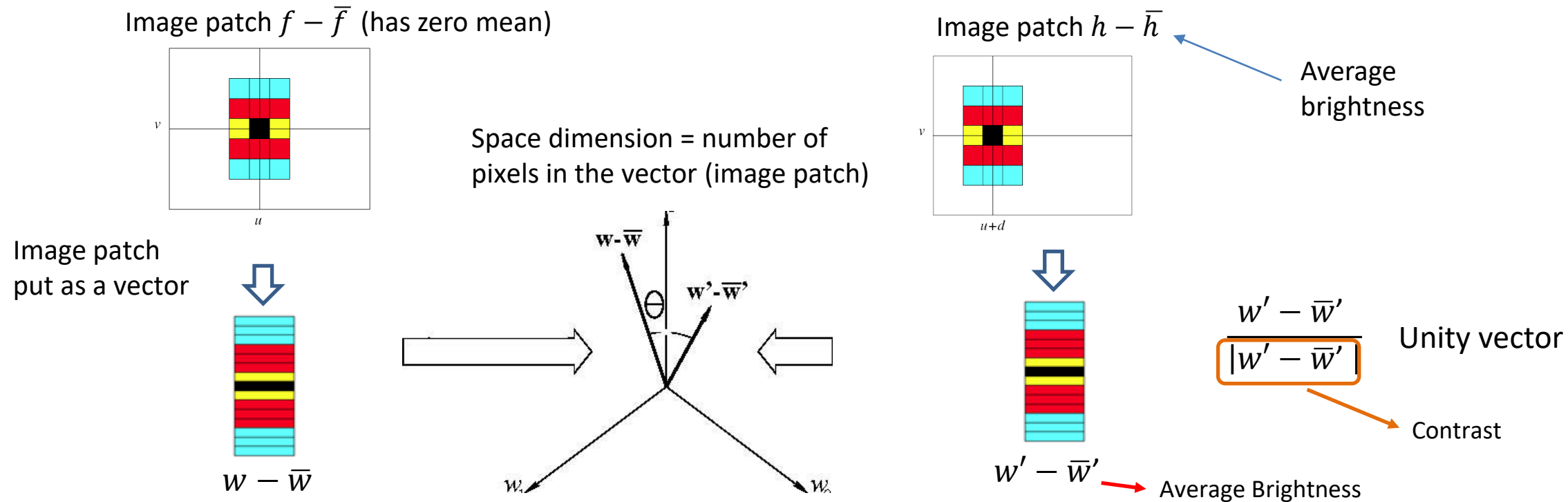$$\bar{h} = \frac{1}{\left|W_m(x, y)\right|} \sum_{(i,j)\in W_m(x,y)} h(i, j)$$

Contrast (norm) of brigntness in the window $W_m$

$$\left\|f - \bar{f}\right\|_{W_m(x,y)} = \sqrt{\sum_{(i,j)\in W_m(x,y)} \left[f(i, j) - \bar{f}(i, j)\right]^2}$$

$$\left\|h - \bar{h}\right\|_{W_m(x,y)} = \sqrt{\sum_{(i,j)\in W_m(x,y)} \left[h(i, j) - \bar{h}(i, j)\right]^2}$$

# Why does NCC measure similarity between two image patches?

## Let's consider an image patch as a vector

Image patch $f - \bar{f}$ (has zero mean)

Image patch $h - \bar{h}$

Average brightness

Space dimension = number of pixels in the vector (image patch)

Image patch put as a vector

$$w - \bar{w}$$

**w-$\bar{\text{w}}$**

$\theta$

**w'-$\bar{\text{w}}$'**

$w_1$

$w_2$

$$w' - \bar{w}'$$

$$\frac{w' - \bar{w}'}{|w' - \bar{w}'|}$$ Unity vector

Contrast

$w' - \bar{w}'$ → Average Brightness

$$C(d) = \frac{1}{|\boldsymbol{w} - \bar{\boldsymbol{w}}|} \frac{1}{|\boldsymbol{w}' - \bar{\boldsymbol{w}}'|} (\boldsymbol{w} - \bar{\boldsymbol{w}}) \cdot (\boldsymbol{w}' - \bar{\boldsymbol{w}}') = \cos\theta$$

The similarity between two unity vectors is given by the angle (or *cos*) between them
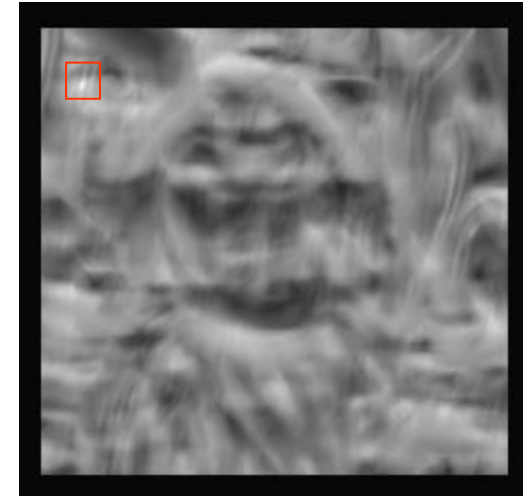
# 4. Keypoint matching through correlation



DEMO IN MATLAB

```
%Read and show the image
flowers = imread('flowers.tif'); figure, imshow(flowers)
% select a template from the image with the mouse
[sub_flowers,rect_flowers] = imcrop(flowers);
% Show the selected template
figure, imshow(sub_flowers)
% Do NCC with the blue channel and display the result
c = normxcorr2(sub_flowers(:,:,1),flowers(:,:,1));
figure, surf(c), shading flat
```
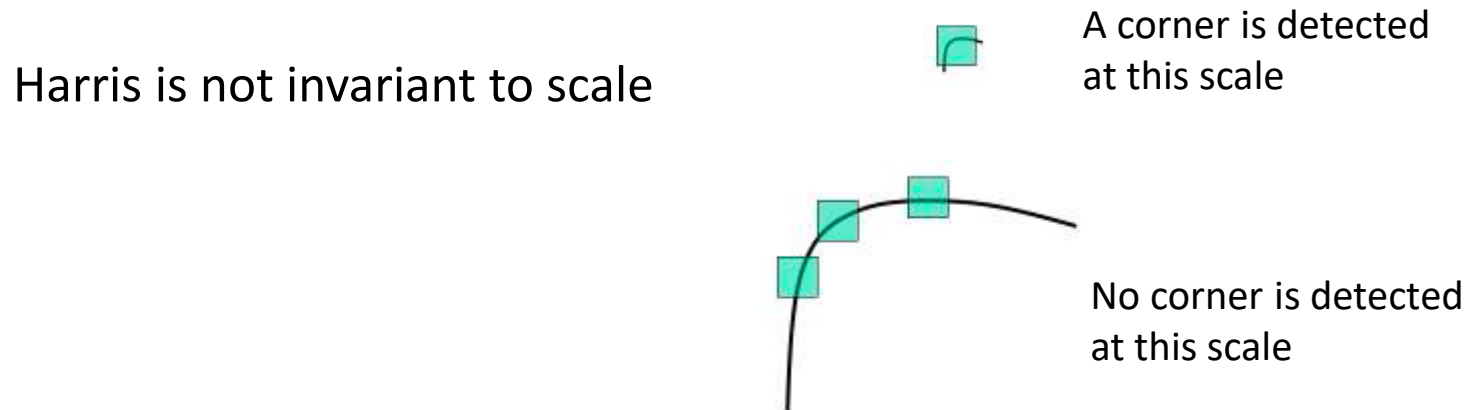
Correlation output

**Notice**: Output not invariant to the rotation of the patch

# Content

1. Introduction
2. Harris detector
   - Idea
   - Formulation
   - Implementation
3. KLT operator
4. Keypoint matching through correlation
5. SIFT operator
   - Scale Space
   - Detector
   - Descriptor

# 5. The SIFT (**S**cale **I**nvariant **F**eature **T**ransform) operator

- **Objective**: Find in image projections of **distintive** 3D points (not necessary corners!) that are <span style="color:red">invariant to scale</span>

Harris is not invariant to scale

A corner is detected at this scale

No corner is detected at this scale

- **Provides** both: Detector y descriptor of the detected keypoints

Proposed (and patented) by **David Lowe**: **"Distinctive image features from scale-invariant keypoints,"** *International Journal of Computer Vision,* 60, 2 (2004).

# 5. The SIFT operator

- Both, detector and descriptor have invariance to:
  - Scale        Important difference against Harris
  - Rotation   Important difference against NCC
    - Ilumination [constrast+brightness]
    - Affine transformation [parcially]
- Principle
  - Search for extrema in the scale space [**Detector**]
  - Normalized histogram of orientation [**Descriptor**]
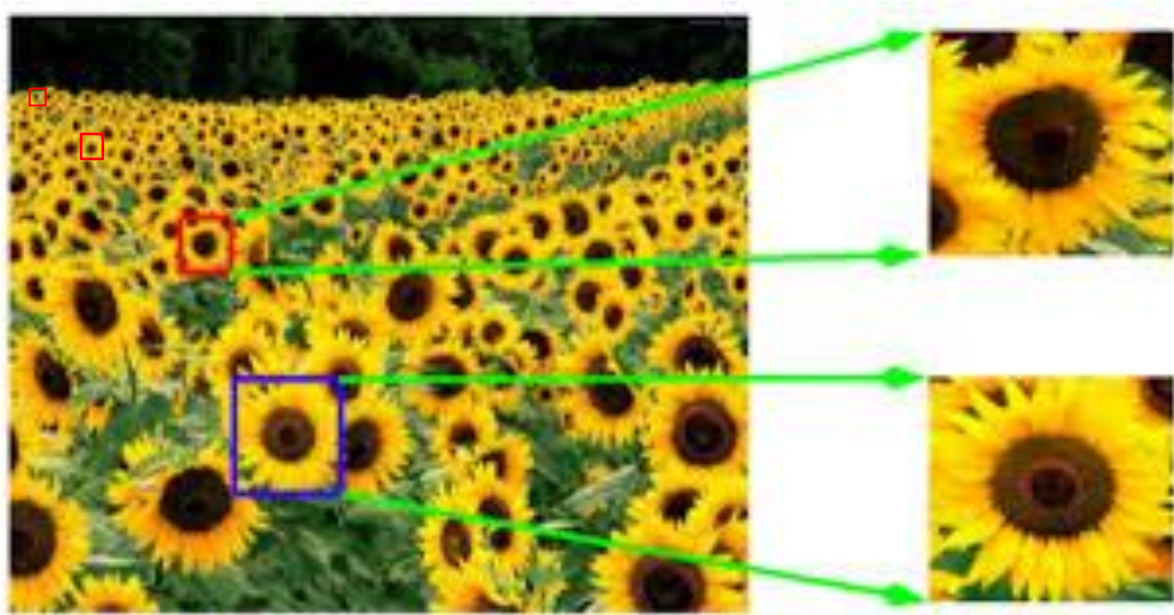- Descriptor is a vector up to 128 dimensions

# 5. The SIFT operator

## Overview of the method

# 5. The SIFT operator

Changing the scale (size) shows up different features

## Scale Space

In images, features emerge at different scales

**Low resolution**

Scale

**High resolution**

# 5. The SIFT operator

## Scale Space

Features show up at a given scale

Example in one-dimension



zoom out          zoom out

# 5. The SIFT operator

## Scale Space

We can change the scale by <span style="color:red">smoothing the signal with a Gaussian</span>

# 5. The SIFT operator

## Scale Space

Changing the scale of an image by <span style="color:red">smoothing with a Gaussian</span>



512 256 128 64 32 16 8

# Scale Space

Changing the scale of an image by <span style="color:red">smoothing with a Gaussian</span>

Gaussian operator:

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

Smoothed image:

$$L(x, y, t) = I(x, y) * G(x, y, t)$$

+           *Detail*           -

The starting image is always the original image (*t*=0)



$L(x,y,0) = I(x,y)$         $L(x,y,1)$         $L(x,y,4)$

starting image *t*=0         $\sigma$=1         $\sigma$=2

So, the scale space stores samples of the function $L(x,y,\sigma^2)$

# 5. The SIFT operator

**Objetive**: "continuous" scale space

Progressive convolution of the  input image with a Gaussian controlled with a constant factor K



$\sigma^2$

$k\sigma^2$

Incresing smoothness

$k^2\sigma^2$

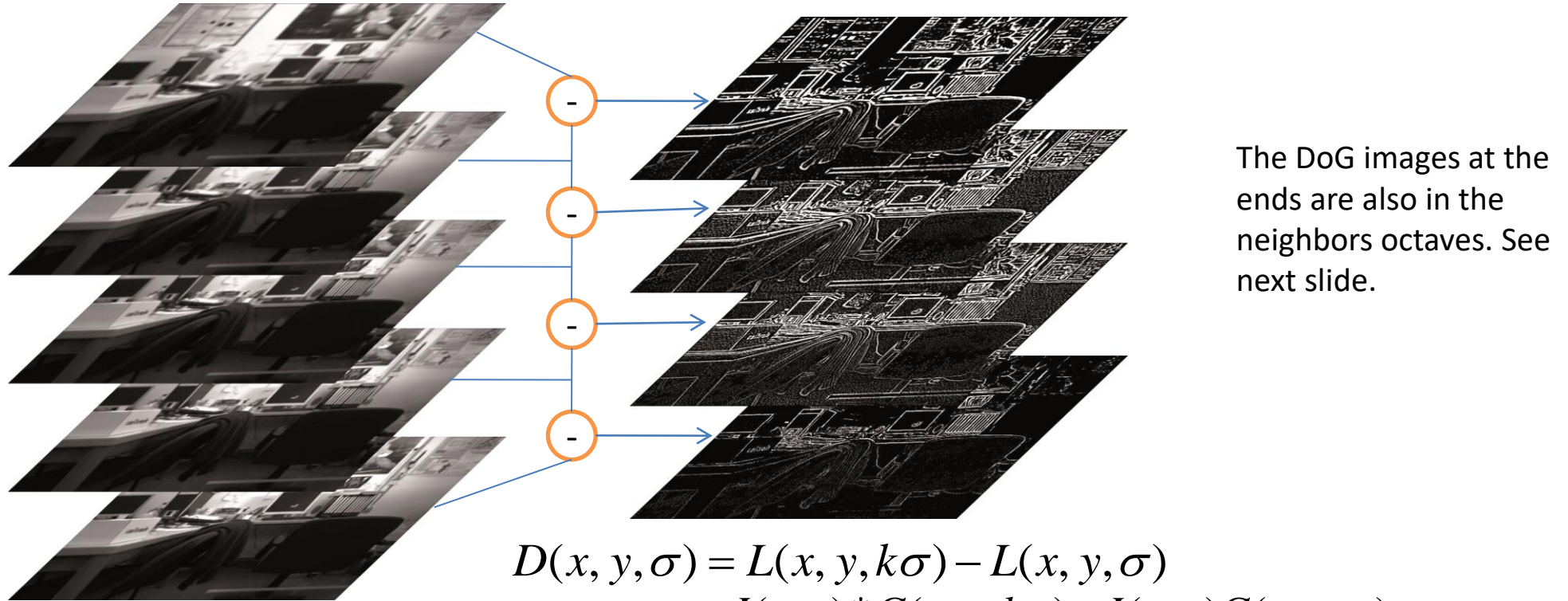$k^3\sigma^2$

$k^4\sigma^2$

# 5. The SIFT operator

- The scale space has the structure of a ***pyramid***: a collection of digital images sampled at progressively coarser spatial resolution and hence of progressively smaller size.

- The pyramid is built upon a number of ***Octaves.***

- Each ***Octave (o)*** consists of s+2 images of the same size (resolution) with increasing smoothness.

- In the **following *Octave*** the image has half the resolution (size) since it does not make sense to keep the resolution when small details have been removed.

# 5. The SIFT operator

From smoothed images to **Difference of Gaussians** (DoG)



The DoG images at the ends are also in the neighbors octaves. See next slide.

$$D(x, y, \sigma) = L(x, y, k\sigma) - L(x, y, \sigma)$$
$$= I(x, y) * G(x, y, k\sigma) - I(x, y)G(x, y, \sigma)$$
$$= I(x, y) * \left[G(x, y, k\sigma) - G(x, y, \sigma)\right]$$
$$= I(x, y) * DoG(x, y, \sigma)$$

Diference of smoothed images = Image convolved with a DoG

# 5. The SIFT operator

## Construction of the scale space through "octaves"

These two images are the same but with half resolution

Same DOG images but half resolution

Scale (next octave)

Scale (first octave)

Gaussian

Difference of Gaussian (DOG)

*DoG are used here to detect BLOBS (not edged)!!*

# 5. The SIFT operator

## Search for extreme points

Each pixel value is compared to its 26 neighbors along the full scale:
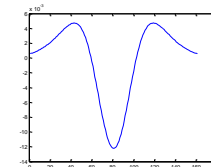- **8** in ithe same scale,
- **9** in the upper scale and
- **9** in the lower scale

A extreme point in the scale gives us a distintive point in (x,y) and in the pyramid

A extreme point in the scale gives us a distintive point in (x,y) and in the pyramid (k)
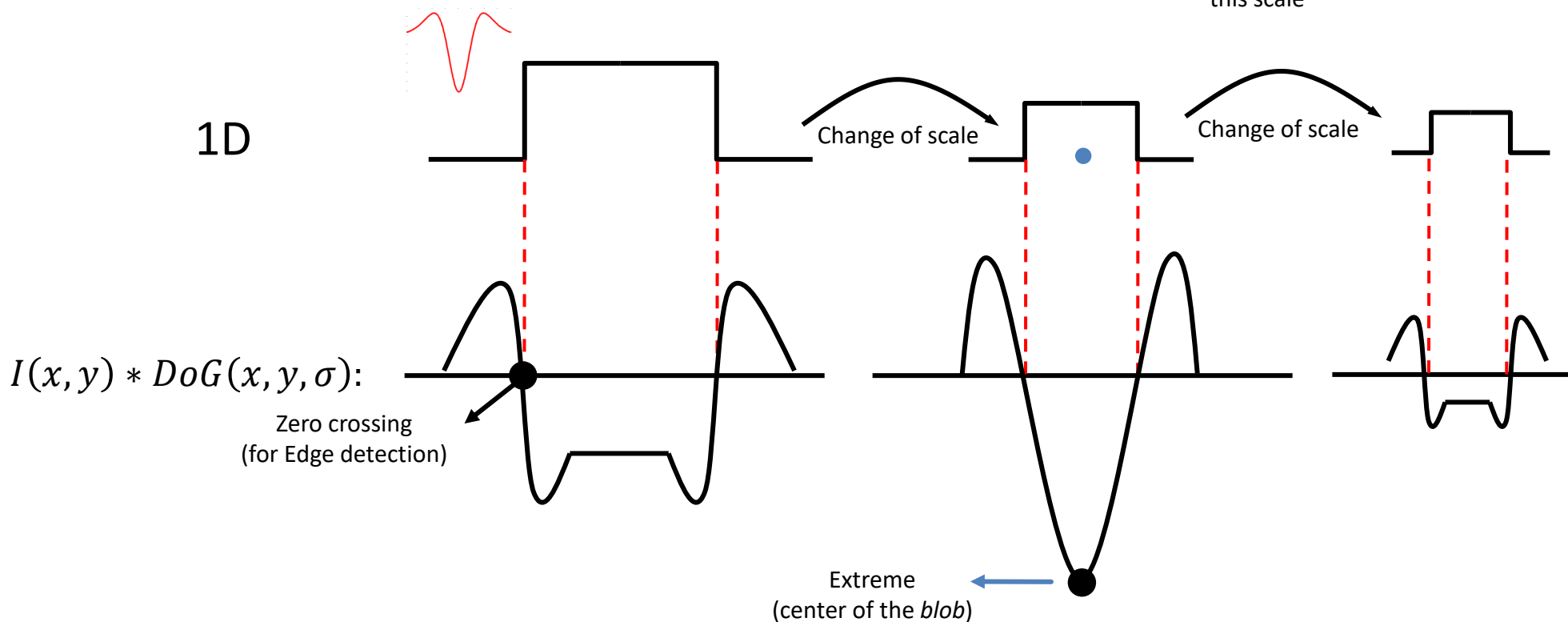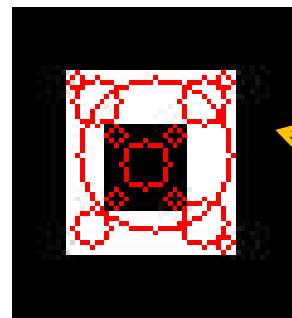
WHY?

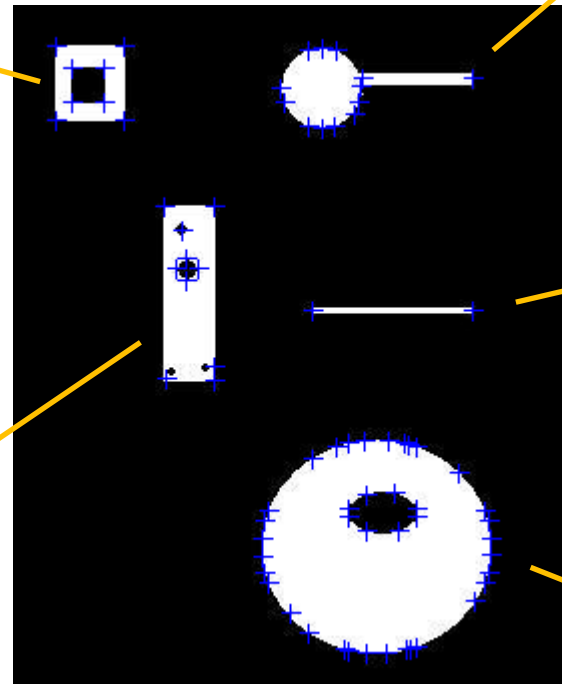Recall:



*DoG* operator (1D)

2D

Detected Blob at this scale

1D

Change of scale

Change of scale

$I(x,y) * DoG(x,y,\sigma)$:

Zero crossing (for Edge detection)

Extreme (center of the *blob*)

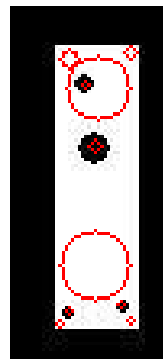# 5. The SIFT operator

Example: SIFT vs. Harris

SIFT

SIFT

SIFT

SIFT

Harris
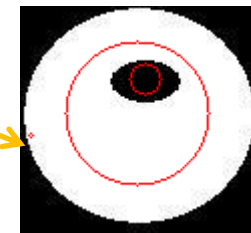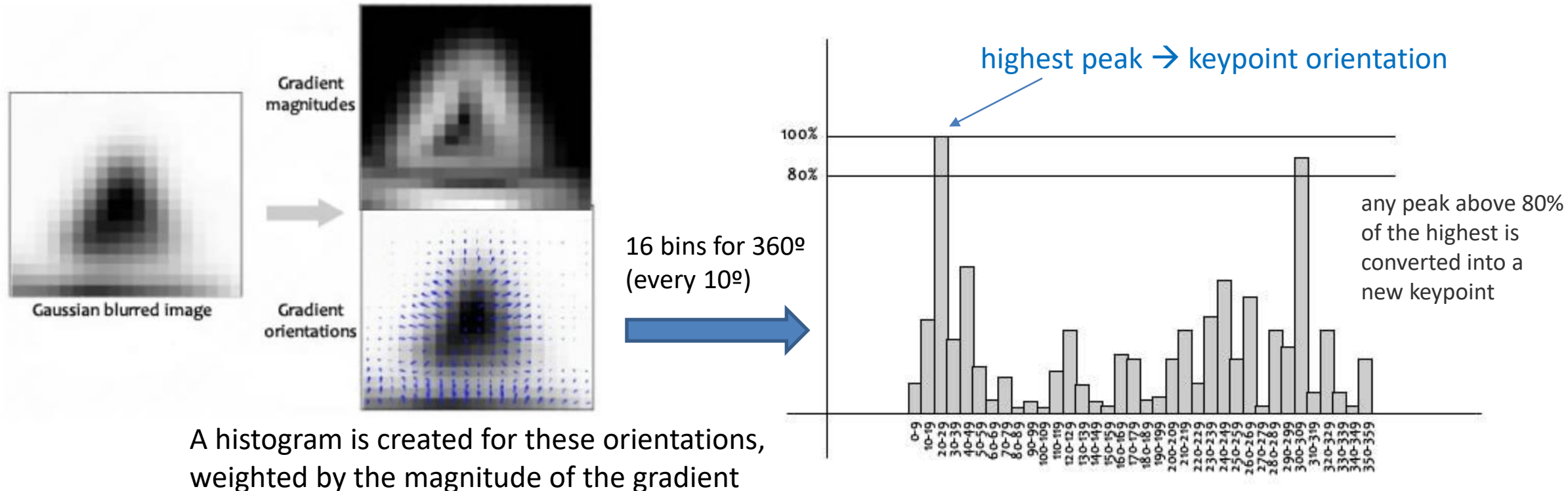
SIFT

SIFT

# SIFT Descriptor: Histogram of orientations around the extreme point

## Obtaining the descriptor orientation:

The magnitude ($m$) and orientation ($\theta$) of the gradient is calculated for all pixels around the keypoint.

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

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



Gradient magnitudes

Gaussian blurred image

Gradient orientations

16 bins for 360º (every 10º)

highest peak → keypoint orientation

100%

80%

any peak above 80% of the highest is converted into a new keypoint

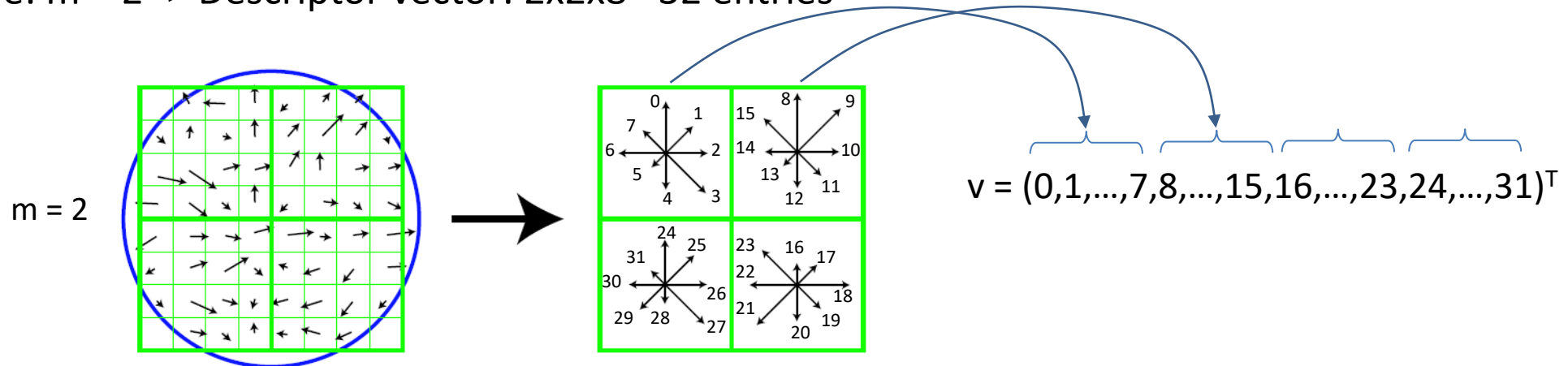A histogram is created for these orientations, weighted by the magnitude of the gradient

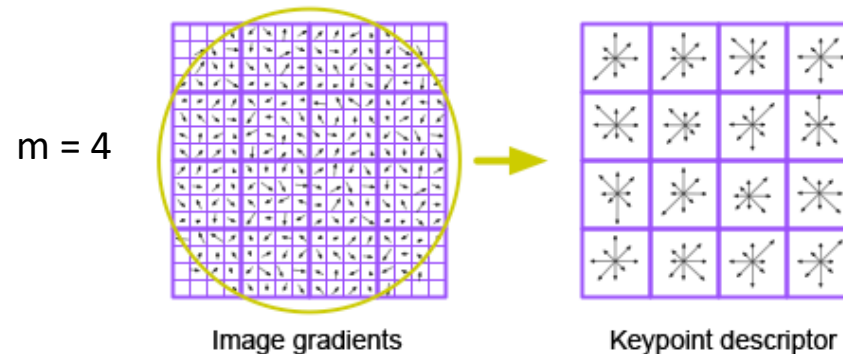# SIFT Descriptor: Histogram of orientations around the extreme point (EP)

## Obtaining the descriptor vector:

- The neighborhood of the EP is divided in **m** x **m cells** (a cell is 4x4 pixels)
- Histogram of 8 orientations for each cell → Size of descriptor vector: **mxmx8**
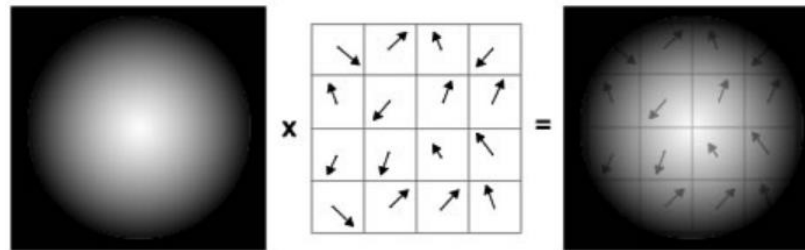
Example: m = 2 -> Descriptor vector: 2x2x8= 32 entries



$v = (0,1,...,7,8,...,15,16,...,23,24,...,31)^\mathsf{T}$

m = 2

In the Lowe's original paper (D. Lowe): **m = 4 ->** Descriptor: 4x4x8 = **128D**

m = 4

Image gradients          Keypoint descriptor

**Obtaining the descriptor vector:**

The histogram of orientations is weighted by

- Magnitude of the gradient: more importance to strong gradients
- Gaussian centered at the extreme point: more importance to close pixels



Descriptor: 4x4x8 = **128D**

## SIFT Invariances

- **Scale**: Window size based on the scale at which the extreme was found

- **Orientation:** Histogram rotated along the keypoint orientation: the keypoint orientation is subtracted from each orientation of the vector
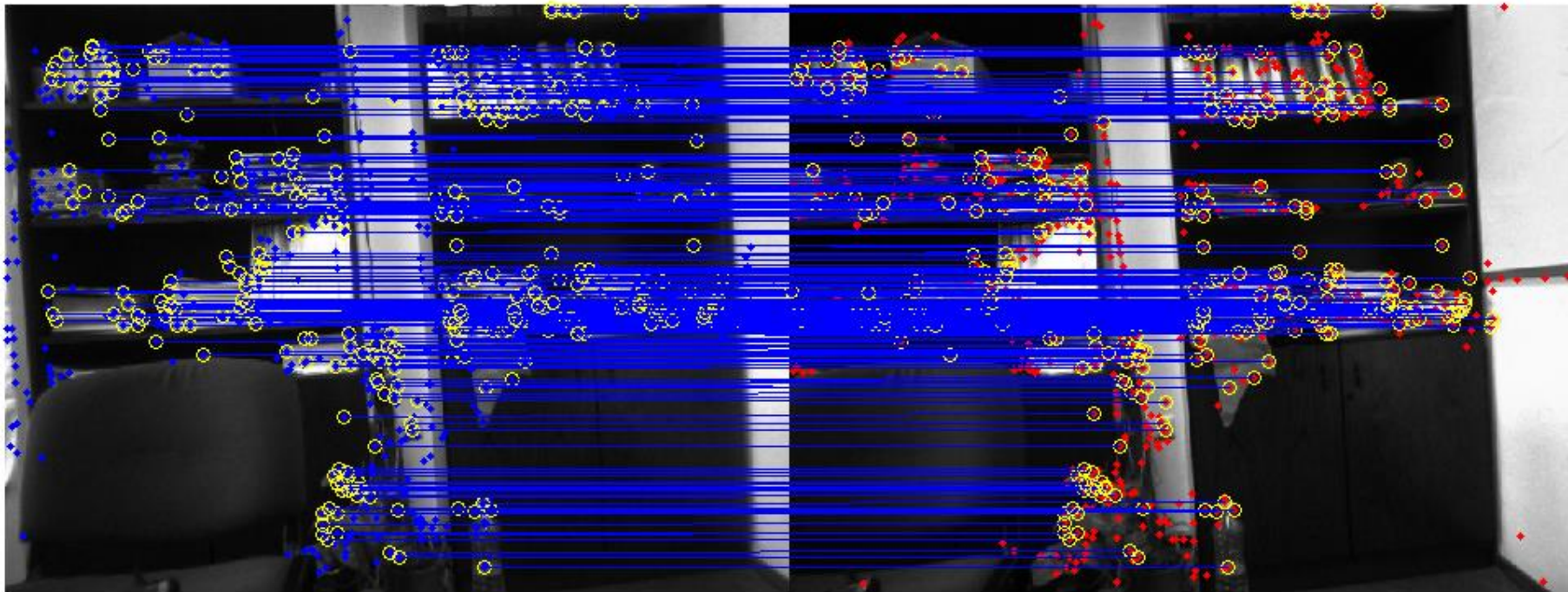
# 5. The SIFT operator

Example: Stereo matches from Euclidean distance

A keypoint of the lest image is matched to the one in the left with the closest descriptor

Points: 965/994          Matches: 359          37.20%
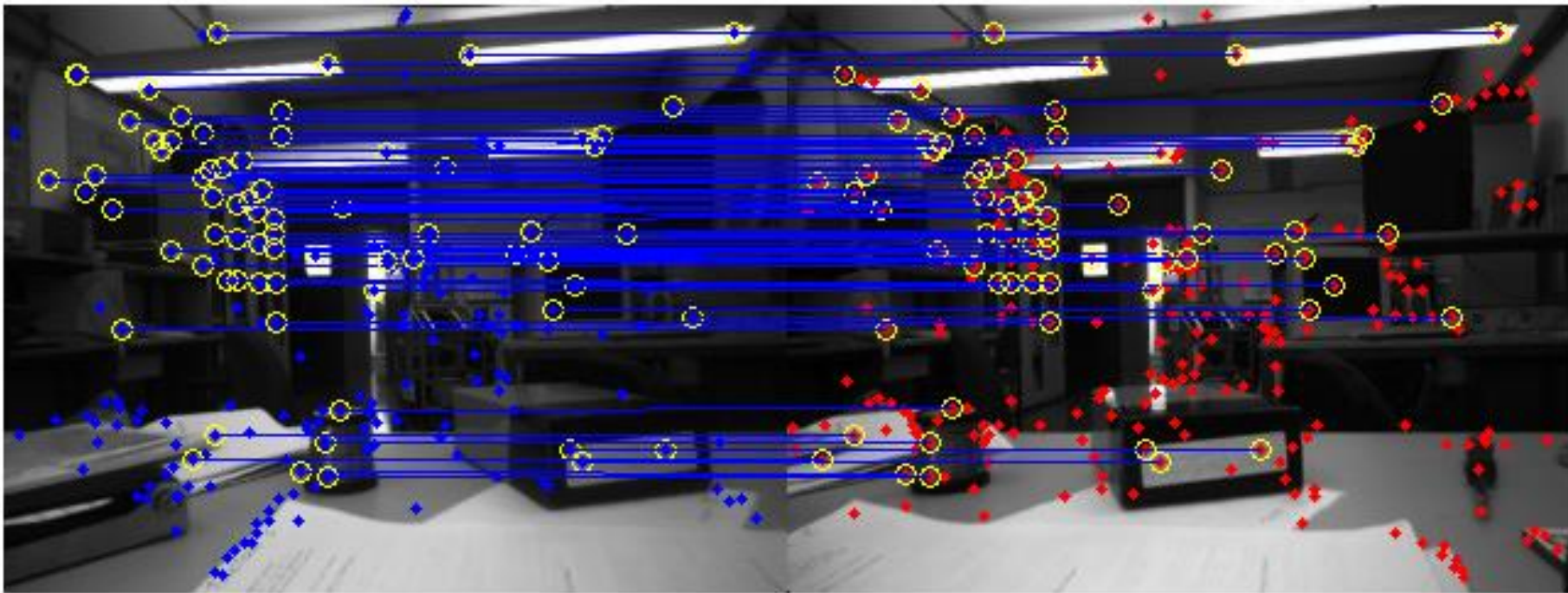
# 5. The SIFT operator

Example: Stereo matches

Points: 245/326          Matches: 89          36.33%
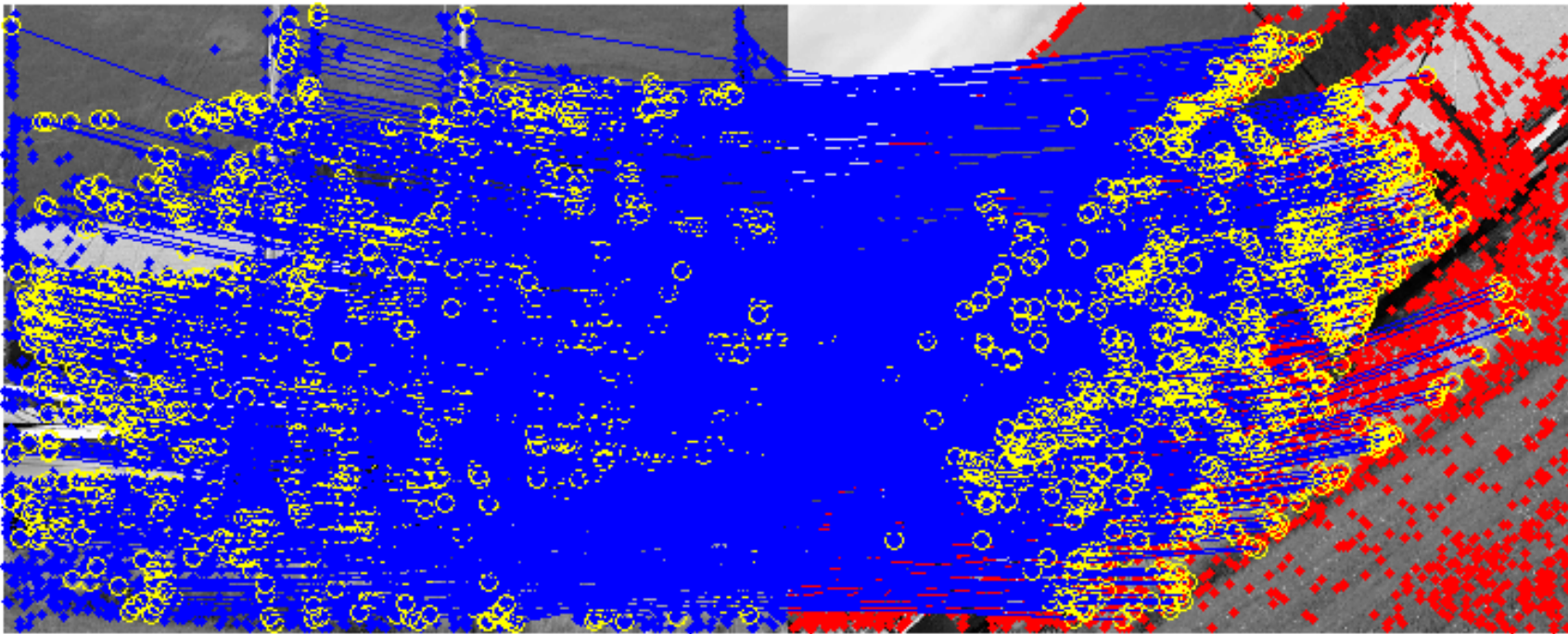
# 5. The SIFT operator

Example: Stereo matches (with rotation)
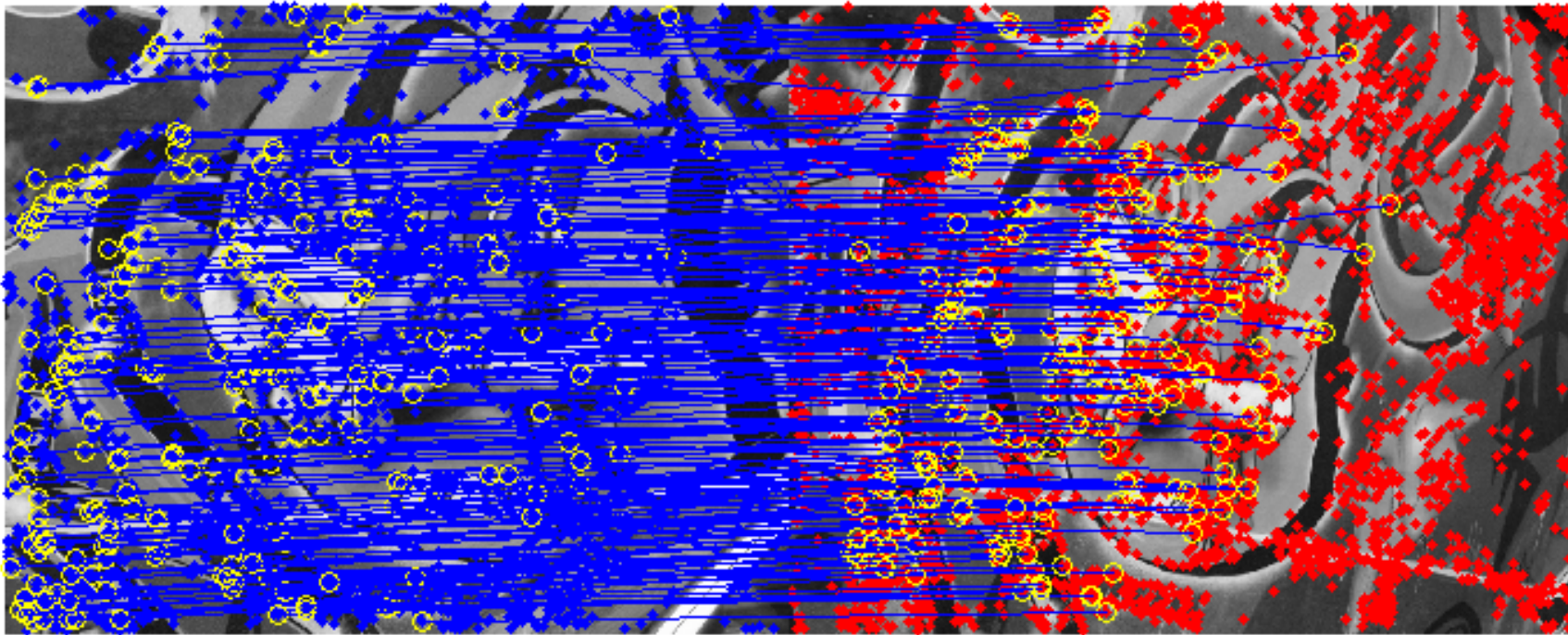
Points: 9687/7113          Matches: 1621          22.79%

# 5. The SIFT operator

Example: Stereo matches (with rotation)

Points: 3105/3990          Matches: 242                    7,79%

# Summary

- Harris operator is …
  - a corner detector which is combined with NCC for matching in other images
  - Based on first-order image derivatives
  - invariant to rotation (because derivatives along the eigenvectors)
  - NOT invariant to scale
  - Invariant to brightness (pixel intensities are not directly considered but derivatives)
  - Robust to noise (because of gaussian smoothing)
- KLT operator
  - Same idea as Harris but the two eigenvalues are used
- SIFT operator
  - Detect Blobs at different scales based on the DoG operator
  - Provides a descriptor with the information of the gradient around the detected keypoint at its scale