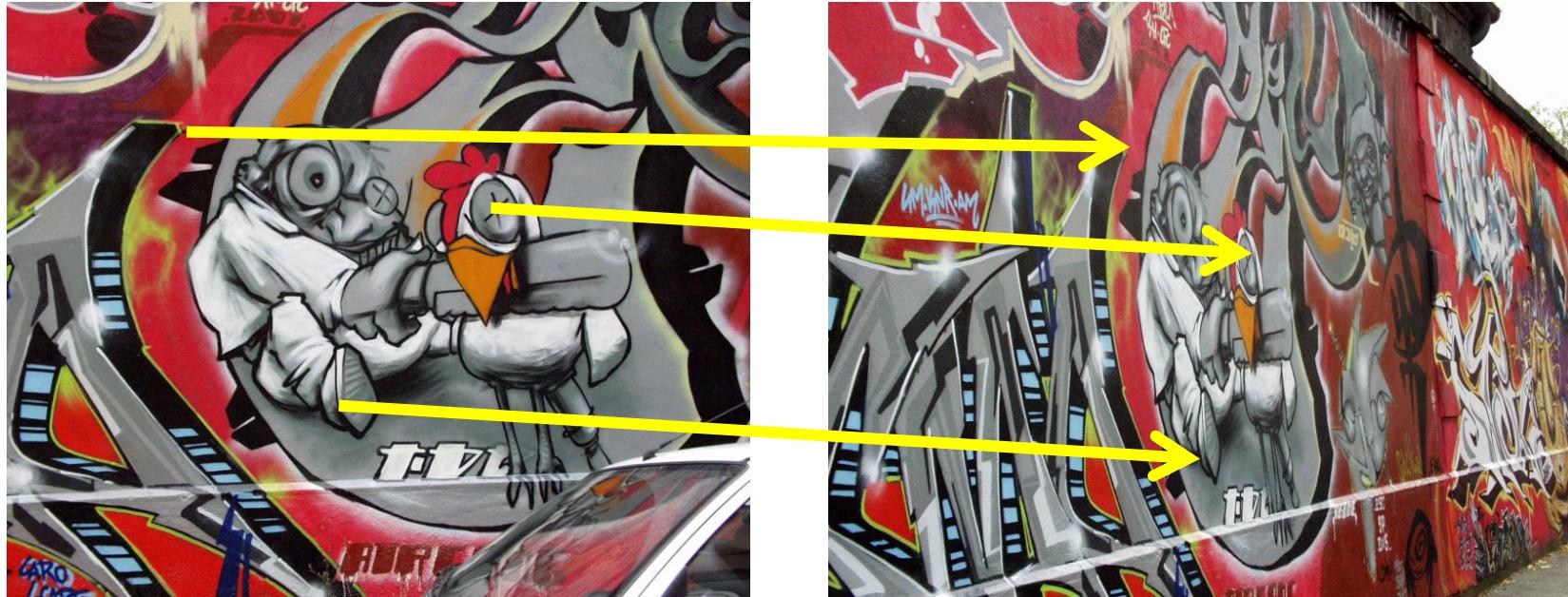


特征匹配与光流

陶煜波，鲍虎军，章国锋
浙江大学CAD&CG国家重点实验室

■ How can we find corresponding points?



稀疏对应——特征点匹配

根据匹配程度分为：

稠密对应——光流估计

■ What is a feature?

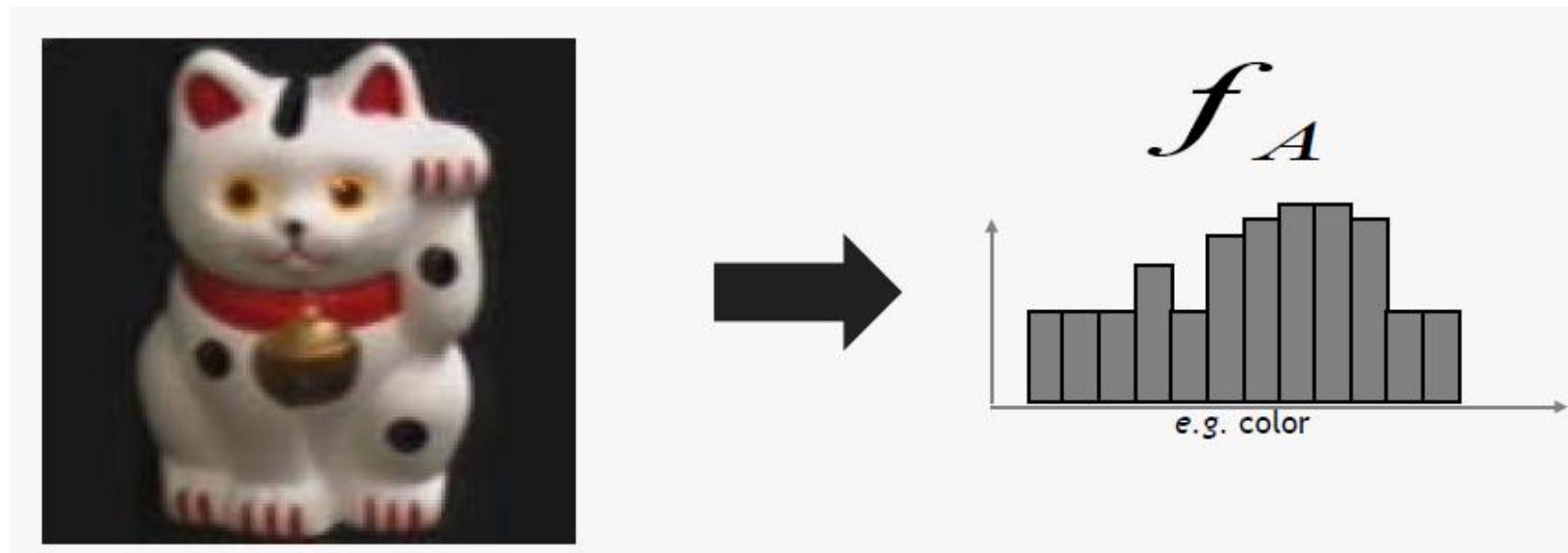
- Point
- Line: straight line, curve,...
- Edge: 2D, 3D
- Corner
- Shape: rectangle, circle, ellipse, sphere,...
- Texture
- Object
- Motion...

■ Invariance

- View point (scale, orientation, translation)
- Lighting condition
- Object deformations
- Partial occlusion

What is feature extraction?

- From image space to feature space:
 - Dimensionality reduction
 - Find a latent semantic space where the pattern can be recognized easily
 - Feature detection and representation



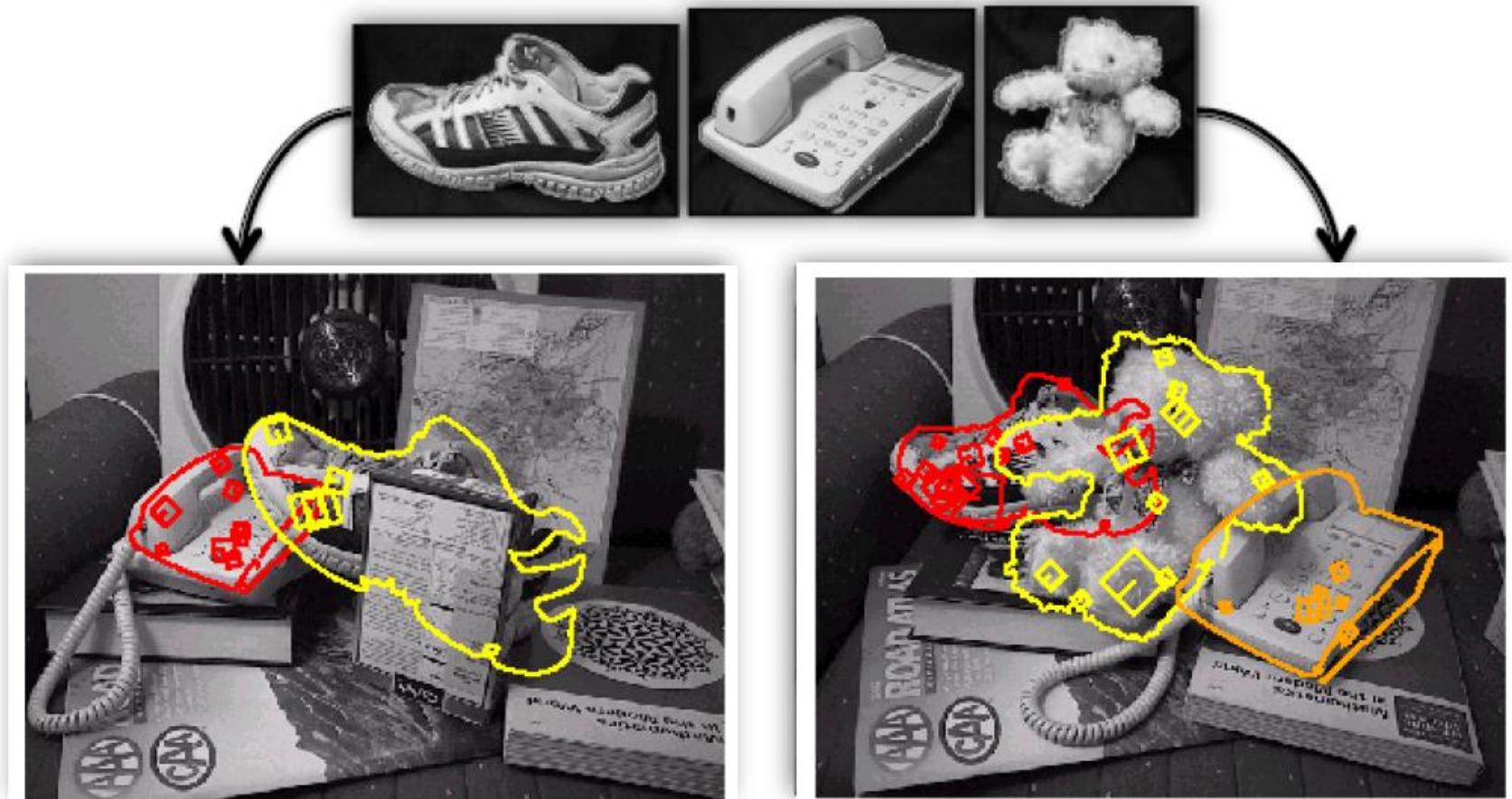
Motivation

- For object recognition



Motivation

- For object recognition



Motivation

■ Panorama creation



Motivation

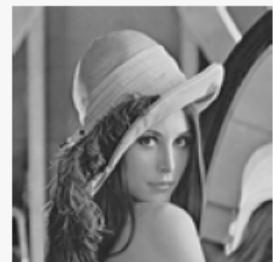
■ Stereo pair matching



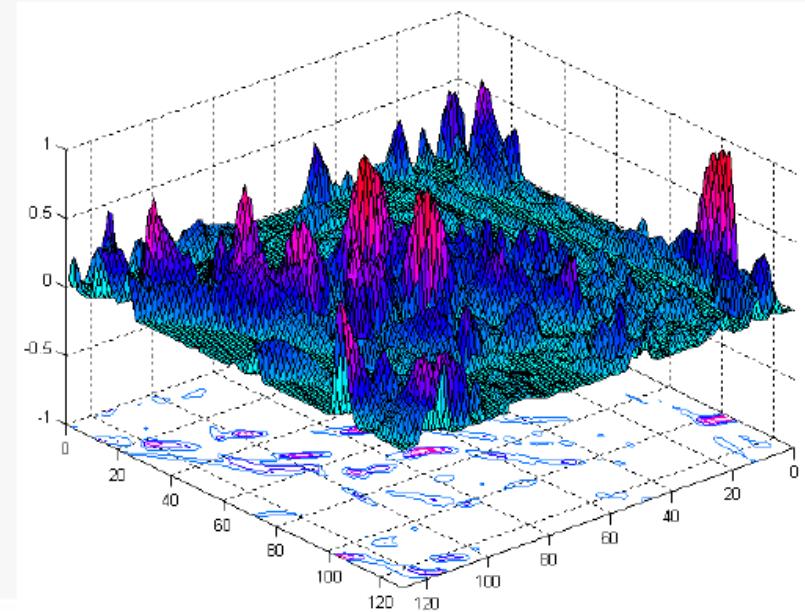
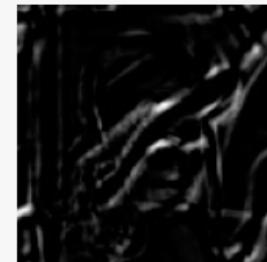
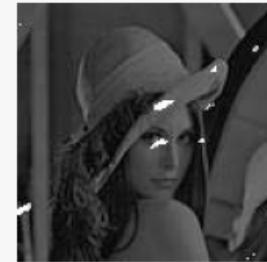
Motivation

■ Template matching

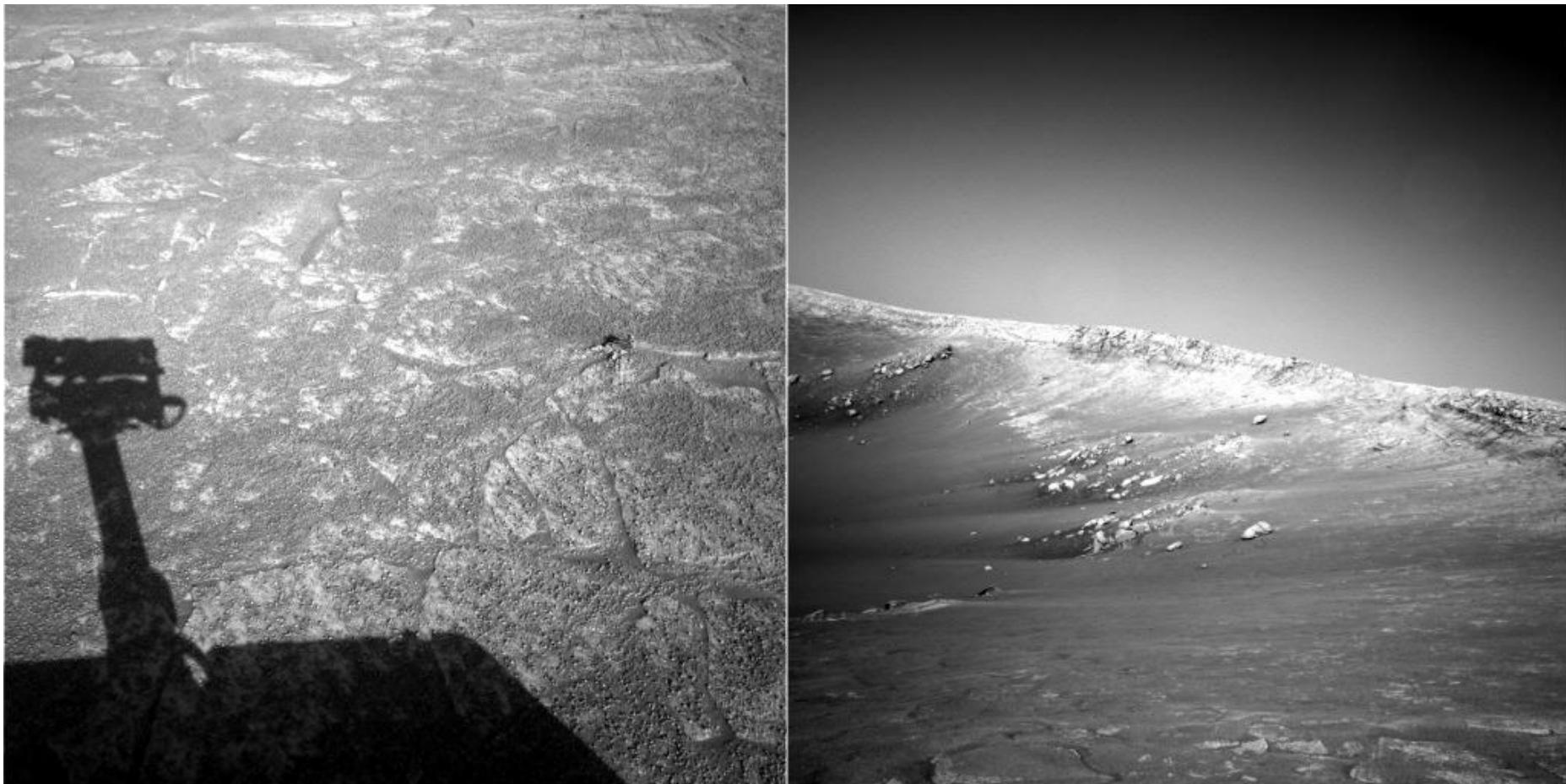
- Pick a template-rectangular/square region of an image
- Goal: find it in the same image/images of the same from different viewpoint



1

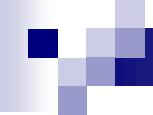


Find corresponding points-Not always easy

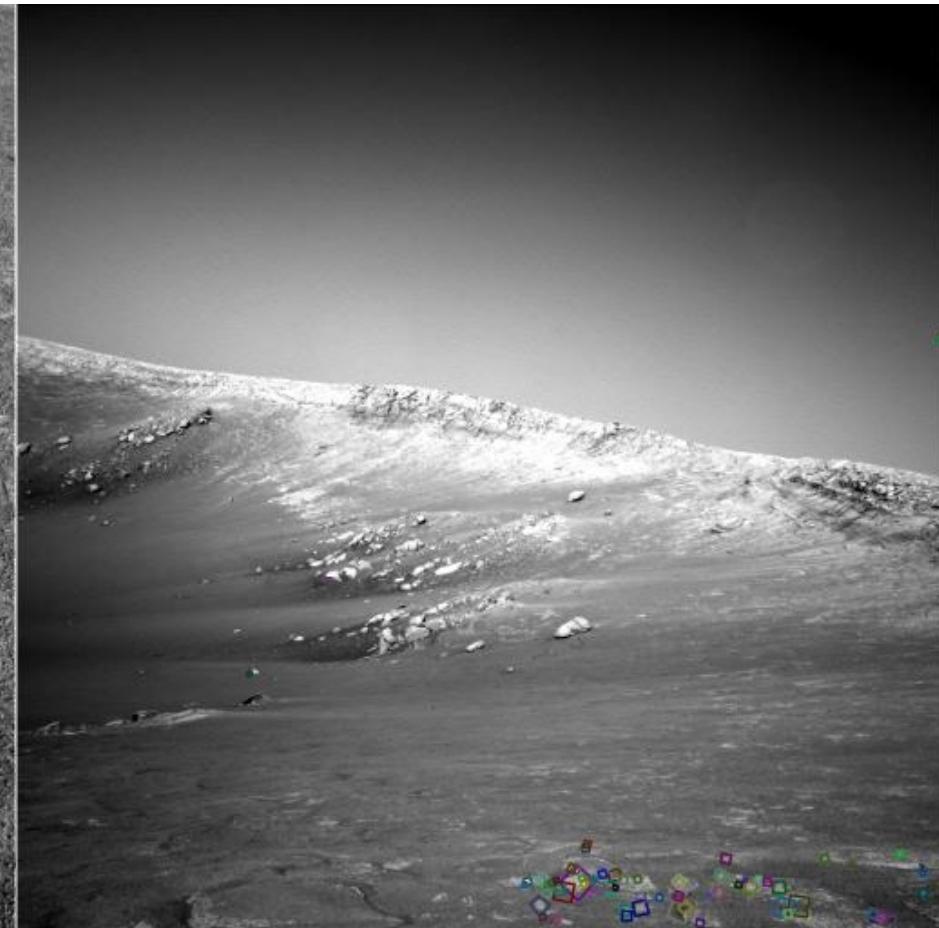
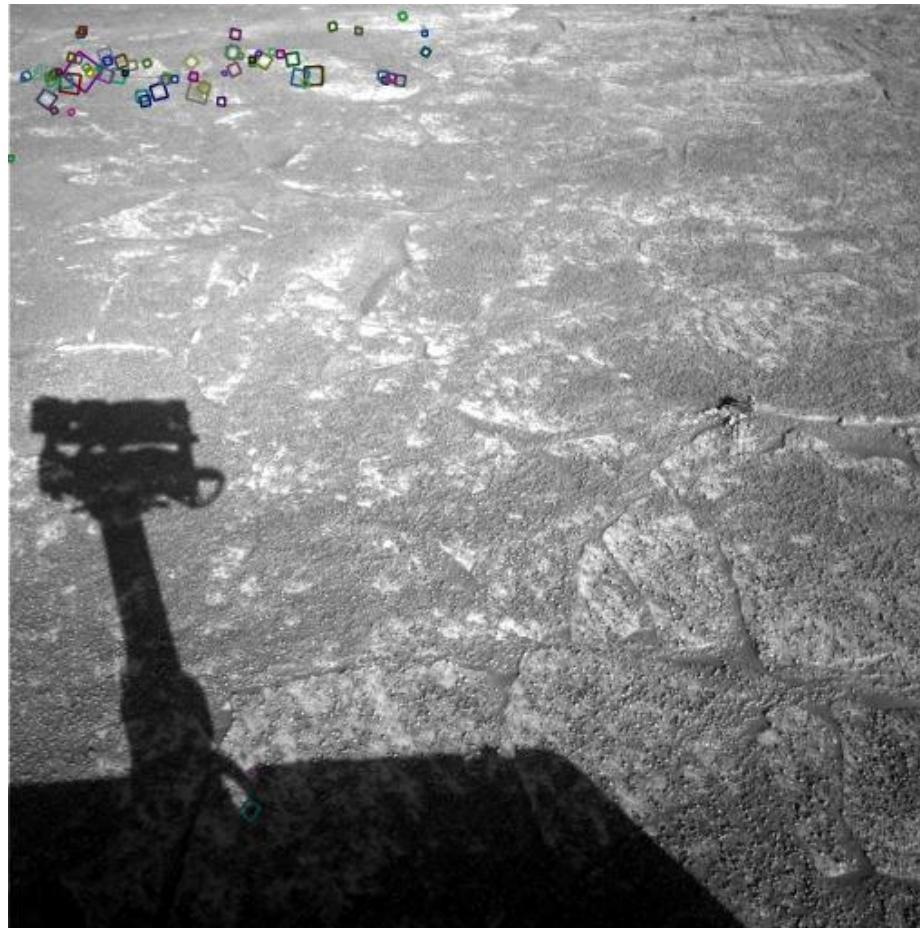


NASA Mars Rover images

from “<http://szeliski.org/Book/>”



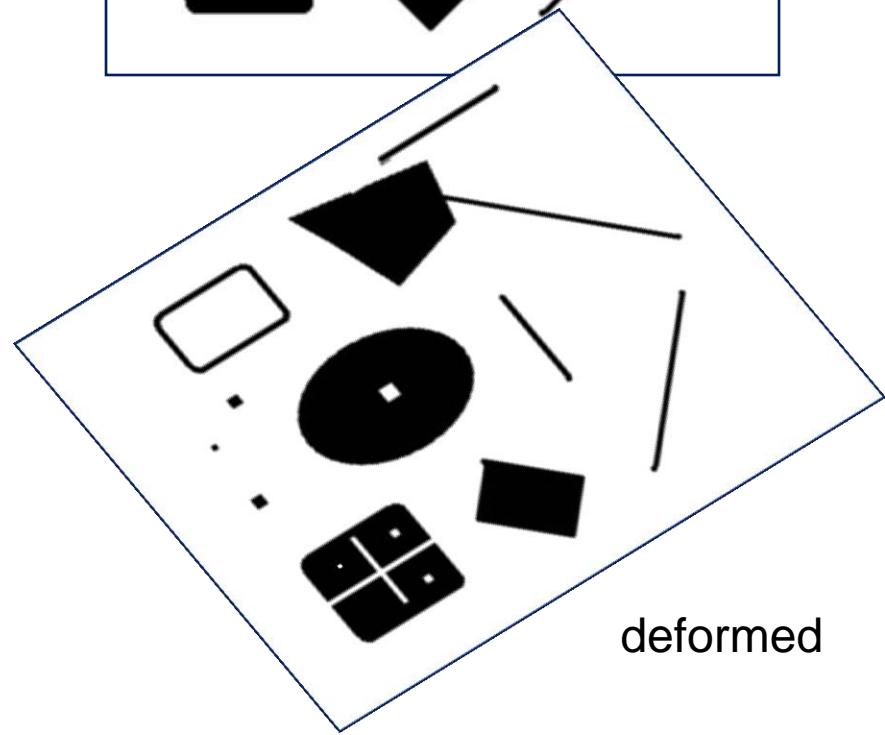
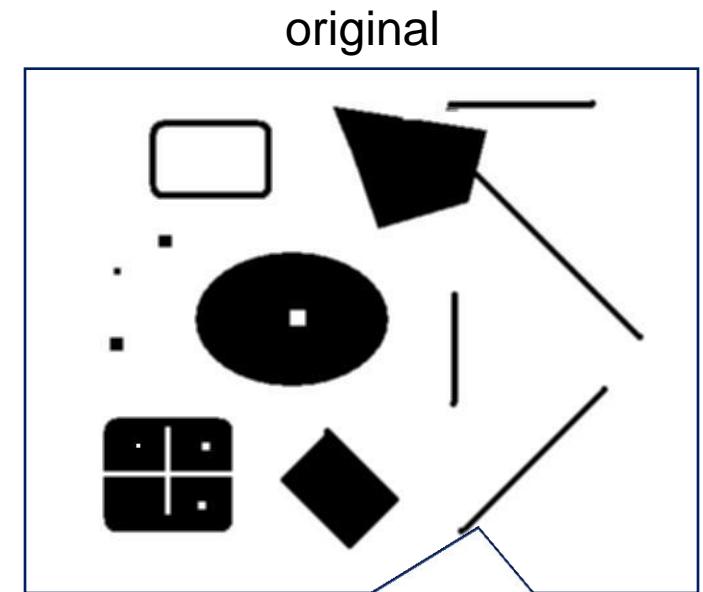
Answer below (look for tiny colored squares...)



NASA Mars Rover images
with SIFT feature matches
Figure by Noah Snavely

Interest points

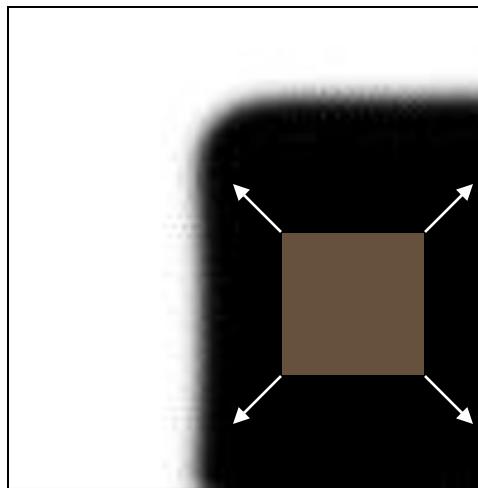
- Suppose you have to click on some point, go away and come back after I deform the image, and click on the same points again.
 - Which points would you choose?



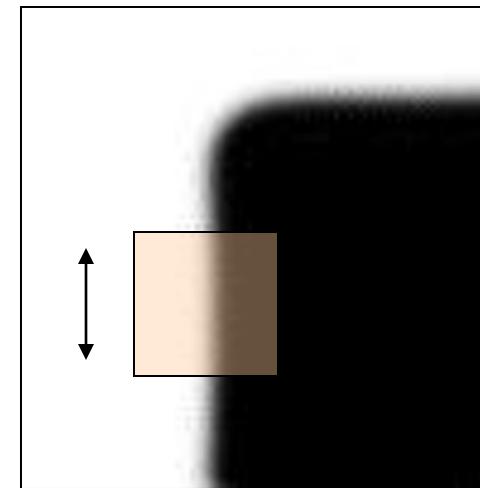
deformed

Corners

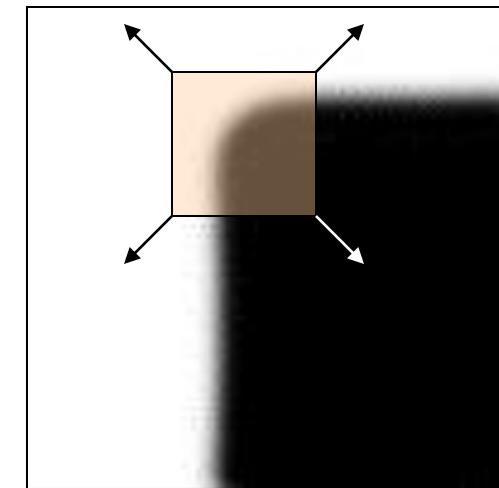
- We should easily recognize the point by looking through a small window
- Shifting a window in *any direction* should give a *large change* in intensity



“flat” region:
no change in
all directions

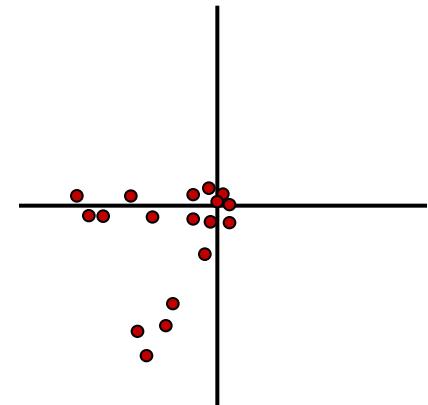
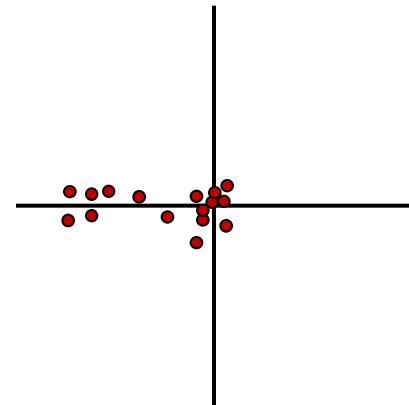
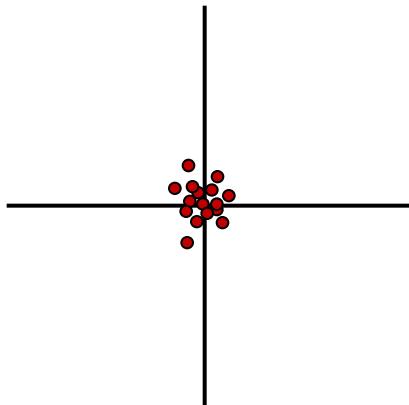
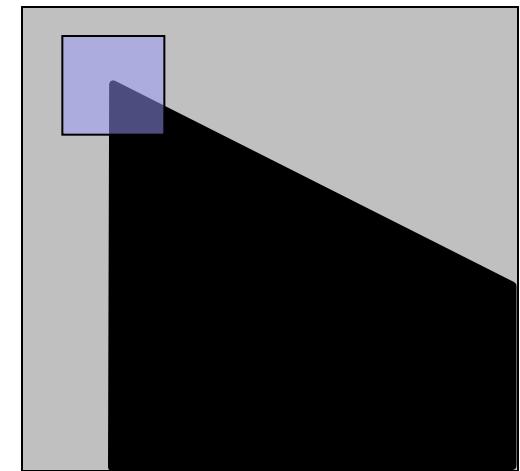
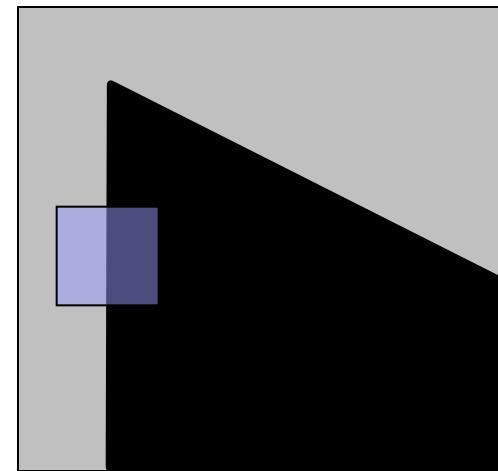
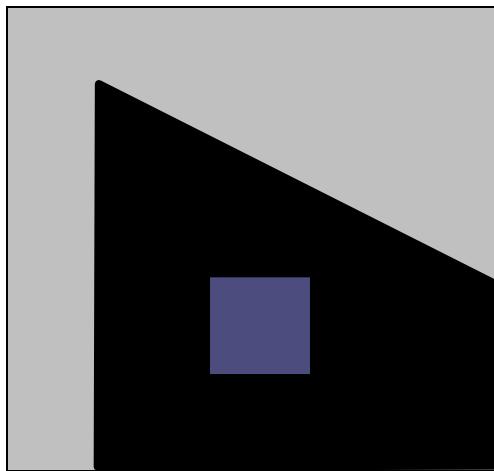


“edge”:
no change
along the edge
direction



“corner”:
significant
change in all
directions

Let's look at the gradient distributions



Principle Component Analysis

Principal component is the direction of highest variance.

Next, highest component is the direction with highest variance *orthogonal* to the previous components.

How to compute PCA components:

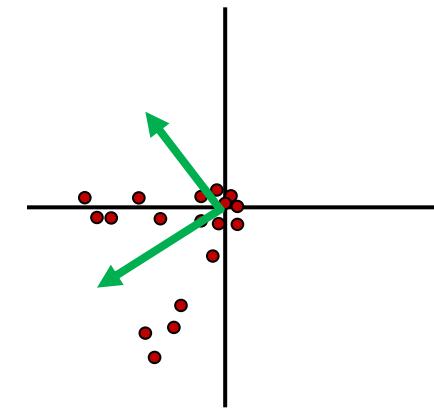
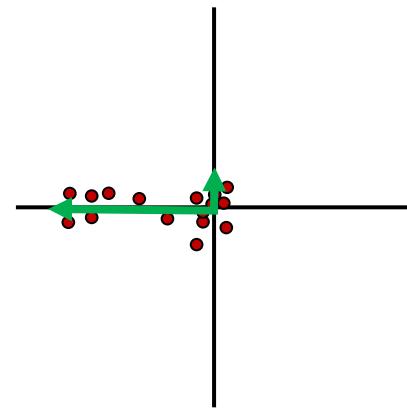
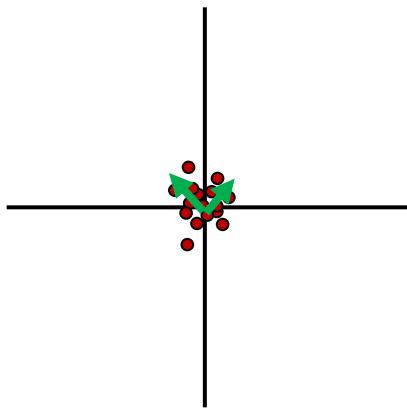
1. Subtract off the mean for each data point.

2. Compute the covariance matrix.

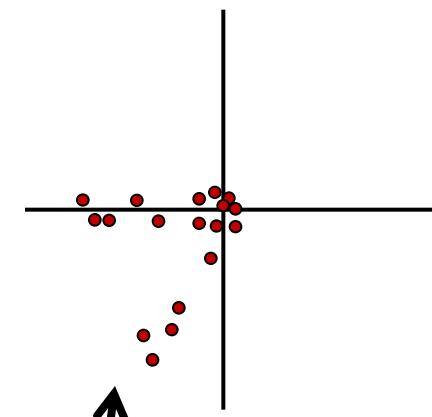
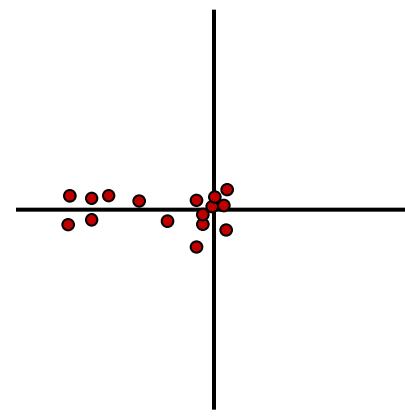
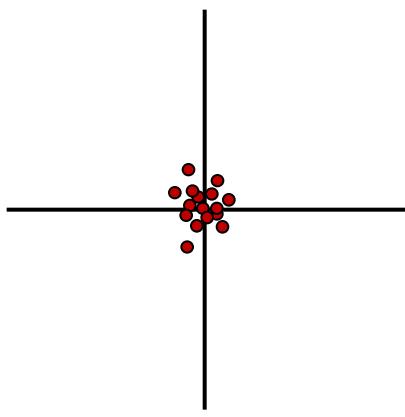
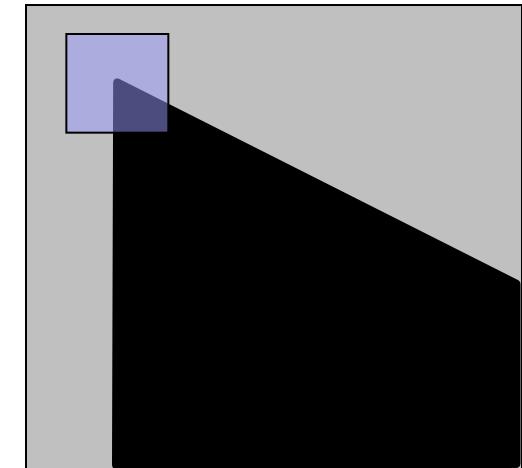
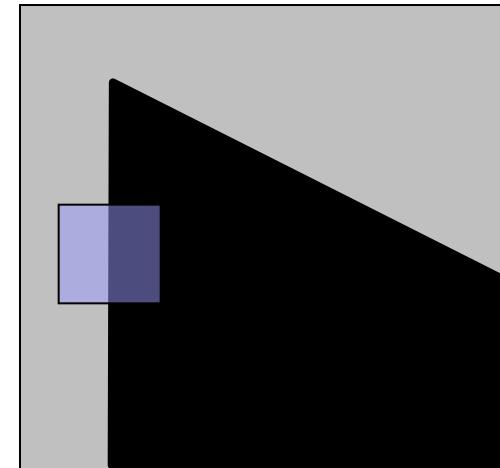
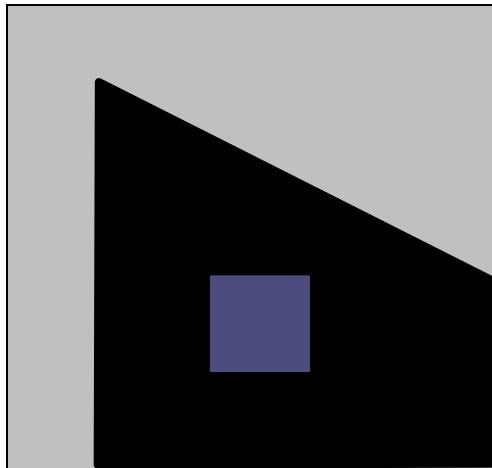
3. Compute eigenvectors and eigenvalues.

4. The components are the eigenvectors ranked by the eigenvalues.

$$Hx = \lambda x$$



Corners have ...



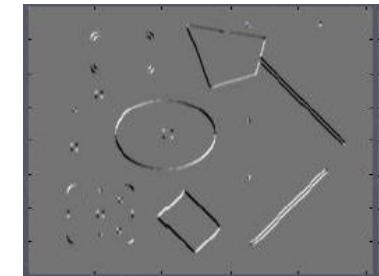
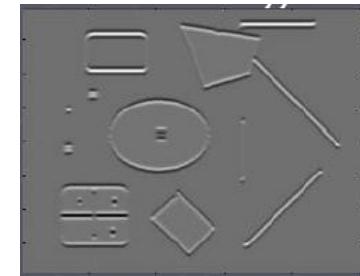
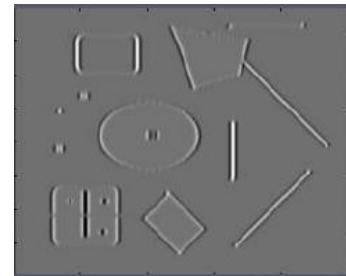
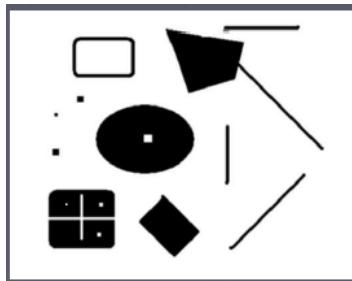
Both eigenvalues are large!



Second Moment Matrix

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

2 x 2 matrix of image derivatives (averaged in neighborhood of a point).



Notation:

$$I_x \Leftrightarrow \frac{\partial I}{\partial x}$$

$$I_y \Leftrightarrow \frac{\partial I}{\partial y}$$

$$I_x I_y \Leftrightarrow \frac{\partial I}{\partial x} \frac{\partial I}{\partial y}$$

The math

To compute the eigenvalues:

1. Compute the covariance matrix.

$$H = \sum_{(u,v)} w(u,v) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \quad I_x = \frac{\partial f}{\partial x}, I_y = \frac{\partial f}{\partial y}$$

↑
Typically Gaussian weights

2. Compute eigenvalues.

$$H = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad \lambda_{\pm} = \frac{1}{2} \left((a + d) \pm \sqrt{4bc + (a - d)^2} \right)$$

Corner Response Function

- Computing eigenvalues are expensive
- Harris corner detector uses the following alternative

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

Reminder:

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc \quad \text{trace} \begin{pmatrix} a & b \\ c & d \end{pmatrix} = a + d$$

Harris detector: Steps

1. Compute Gaussian derivatives at each pixel
2. Compute second moment matrix M in a Gaussian window around each pixel
3. Compute corner response function R
4. Threshold R
5. Find local maxima of response function (nonmaximum suppression)

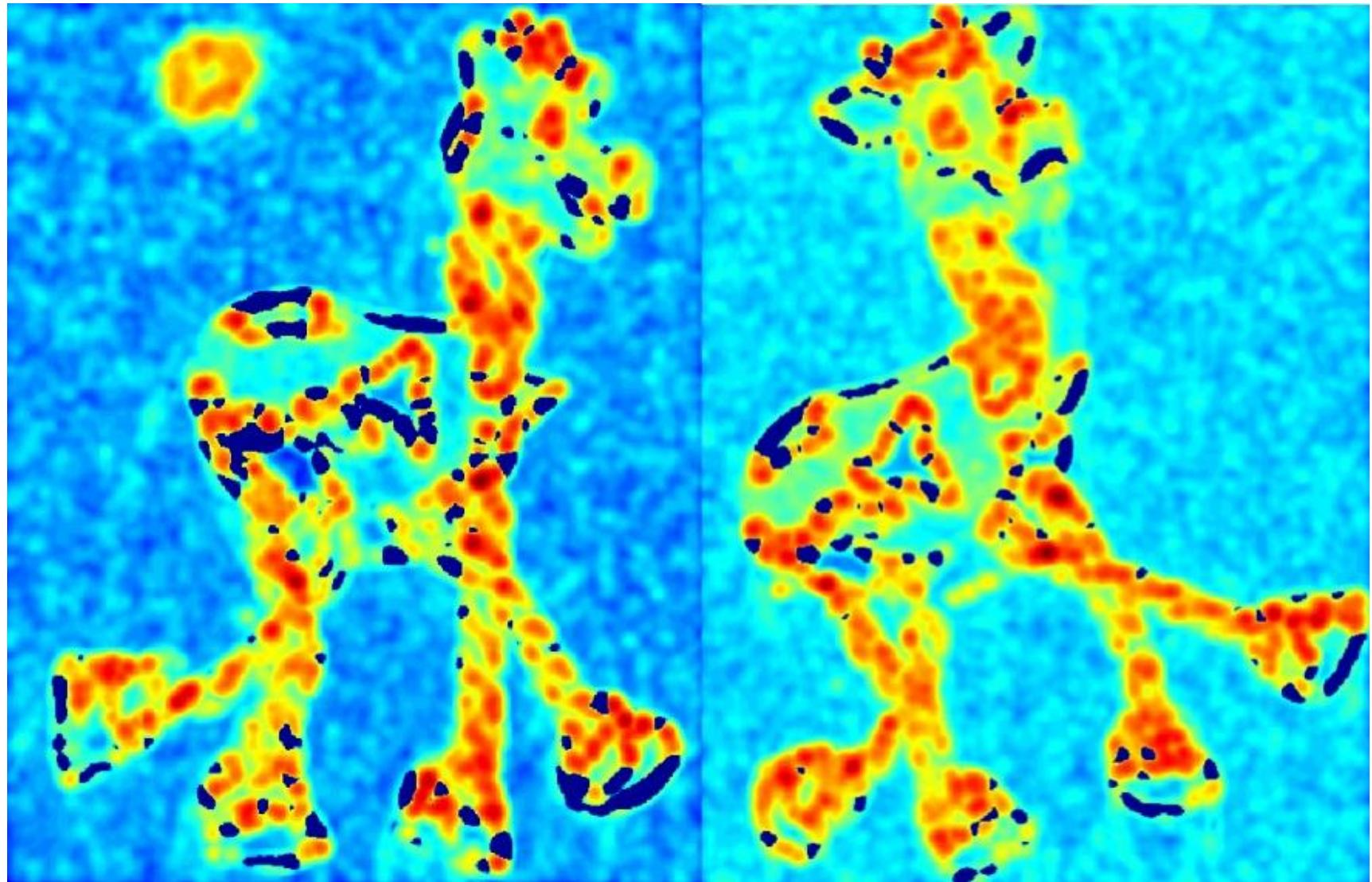
C.Harris and M.Stephens. "[A Combined Corner and Edge Detector.](#)"
Proceedings of the 4th Alvey Vision Conference: pages 147—151, 1988.

Harris Detector: Steps



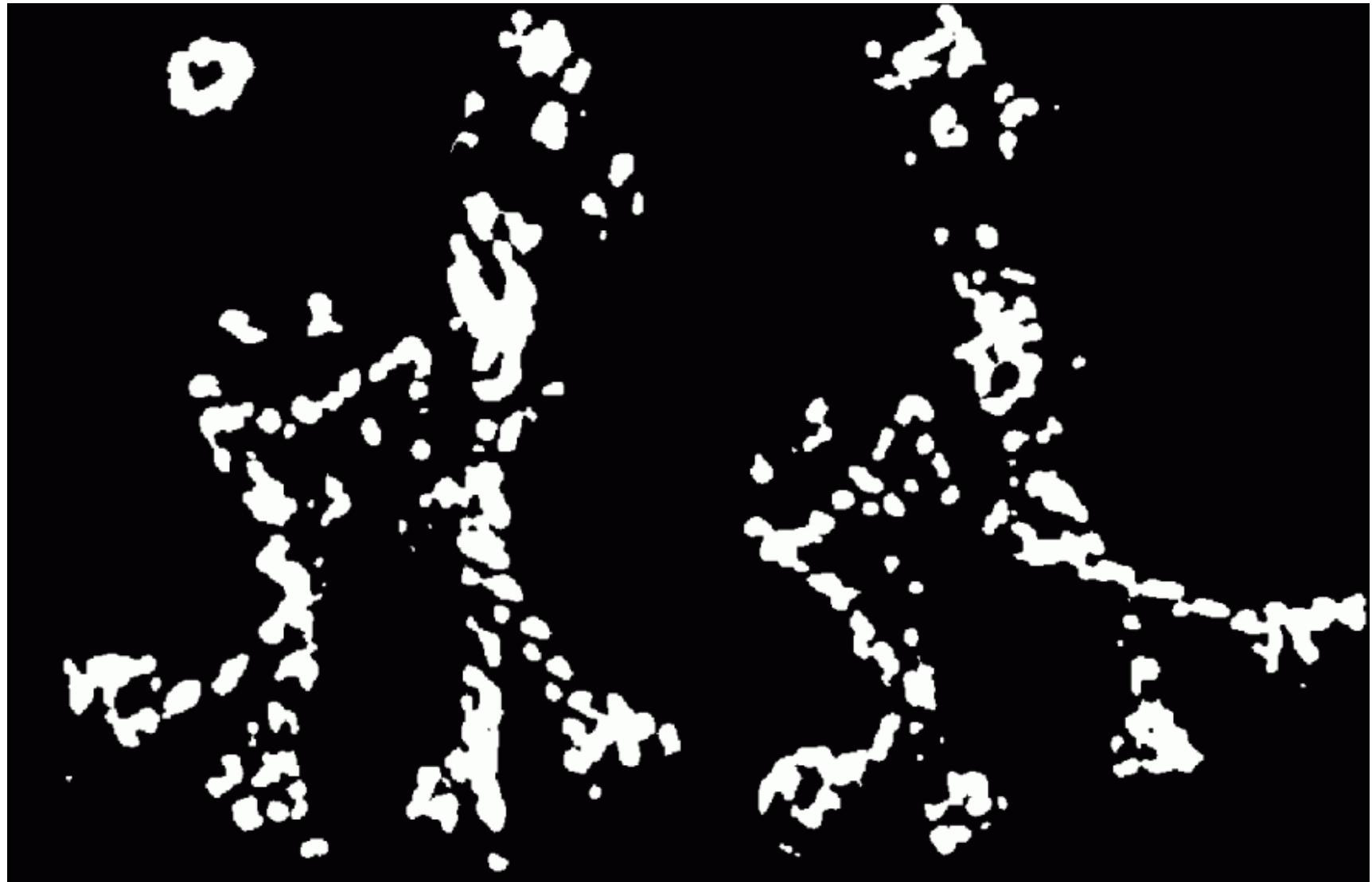
Harris Detector: Steps

Compute corner response R



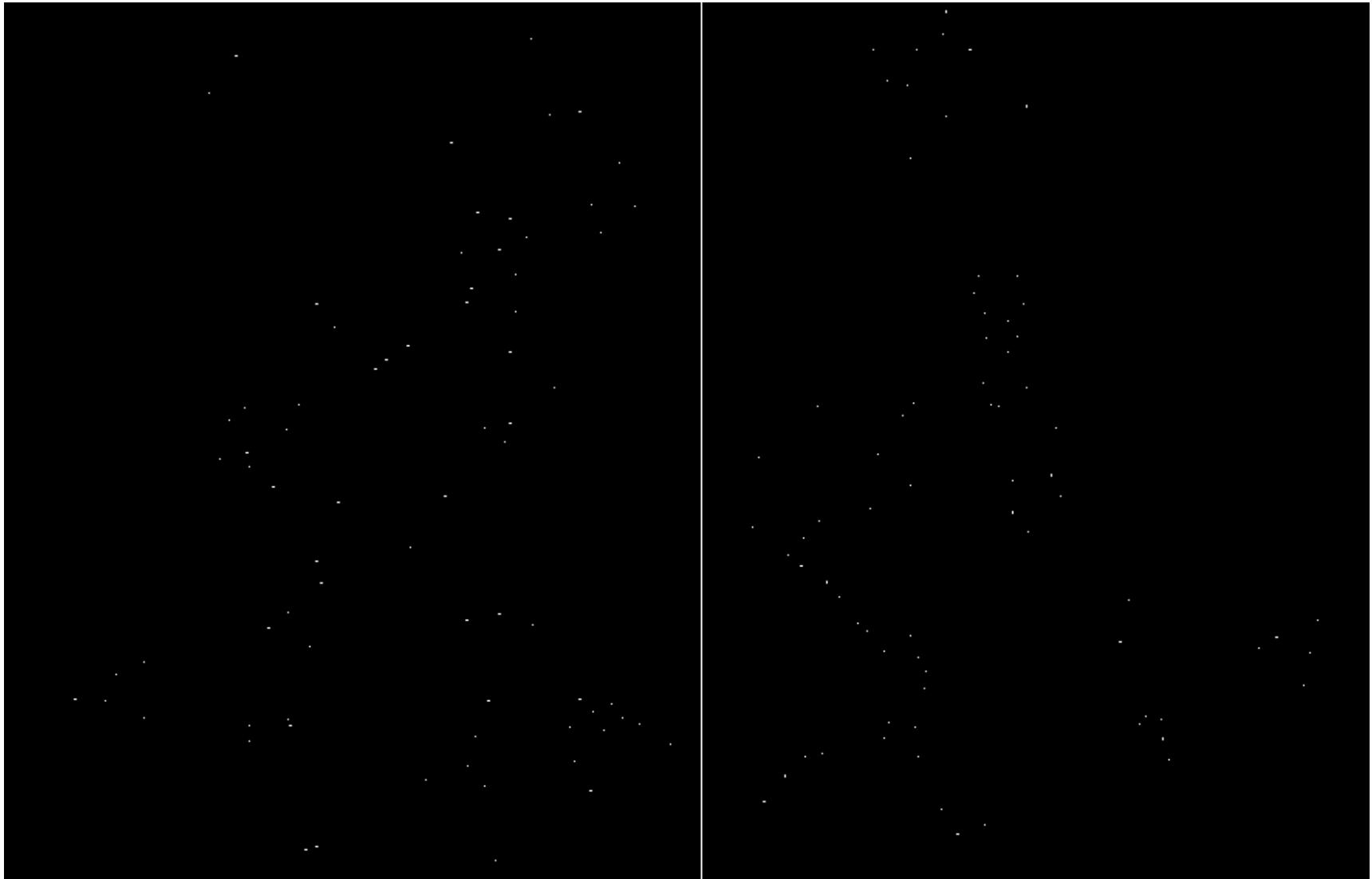
Harris Detector: Steps

Find points with large corner response: $R > \text{threshold}$



Harris Detector: Steps

Take only the points of local maxima of R

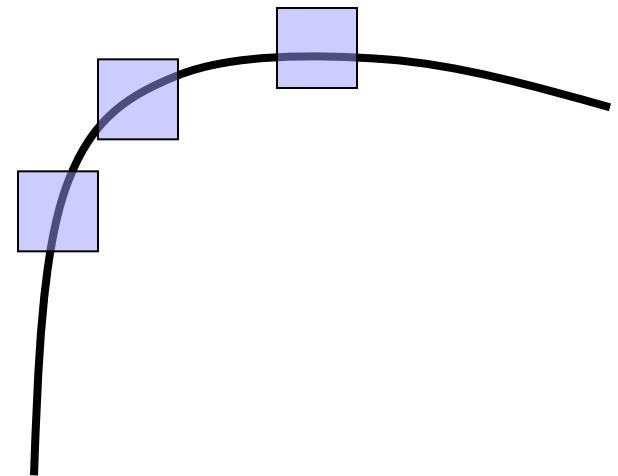
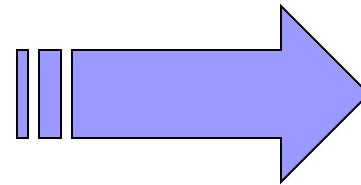


Harris Detector: Steps



Properties of the Harris corner detector

- Translation invariant? Yes
- Rotation invariant? Yes
- Scale invariant? No

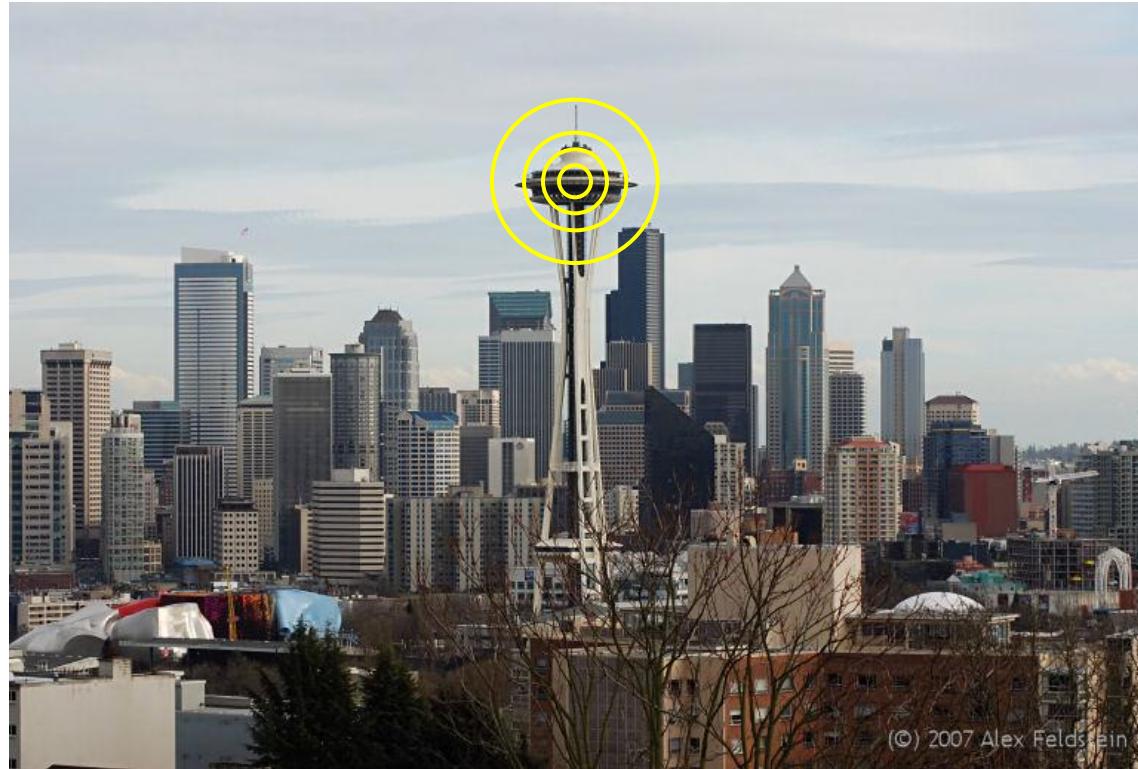


Corner !

All points will be
classified as edges

Scale

Let's look at scale first:



What is the “best” scale?

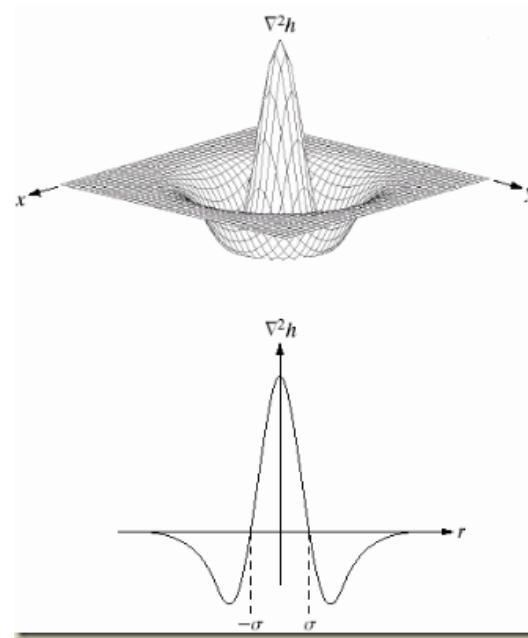
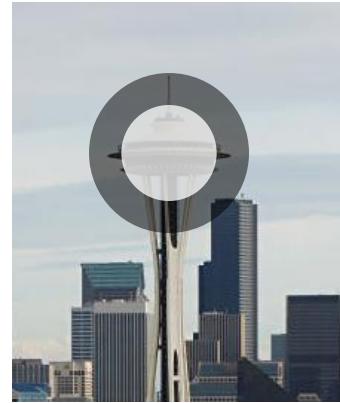
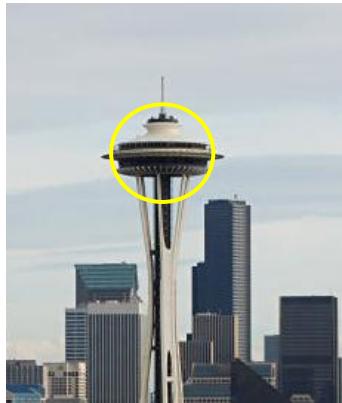
Scale Invariance



$$f(I_{i_1 \dots i_m}(x, \sigma)) = f(I_{i_1 \dots i_m}(x', \sigma'))$$

How can we independently select interest points in each image, such that the detections are repeatable across different scales?

Differences between Inside and Outside

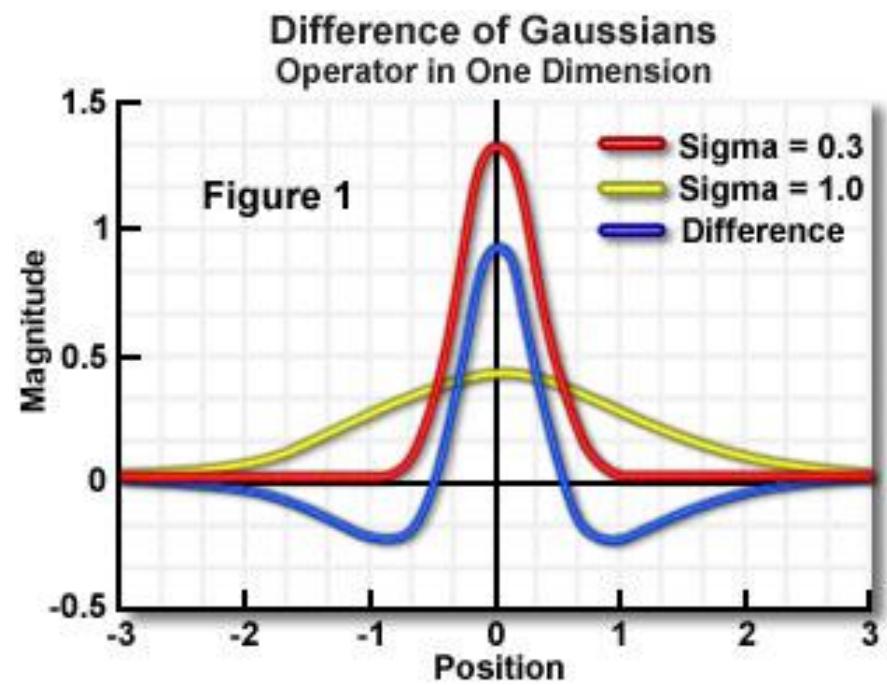


Scale

Why Gaussian?

It is invariant to scale change,
i.e., $f * \mathcal{G}_\sigma * \mathcal{G}_{\sigma'} = f * \mathcal{G}_{\sigma''}$
and has several other nice
properties. Lindeberg, 1994

In practice, the Laplacian
is approximated using a
Difference of Gaussian
(DoG).



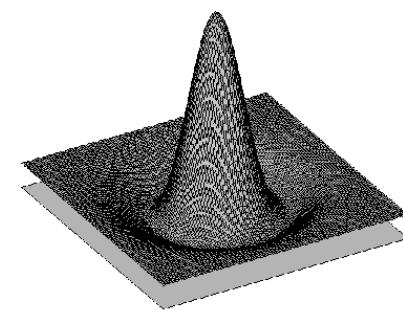
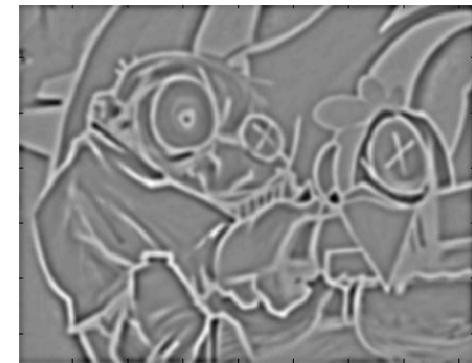
Difference-of-Gaussian (DoG)



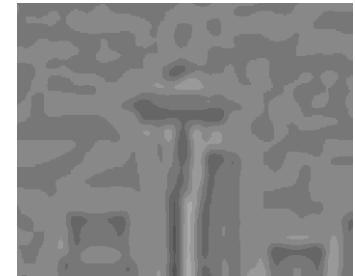
-



=



DoG example



$\sigma = 1$



$\sigma = 66$

Scale invariant interest points

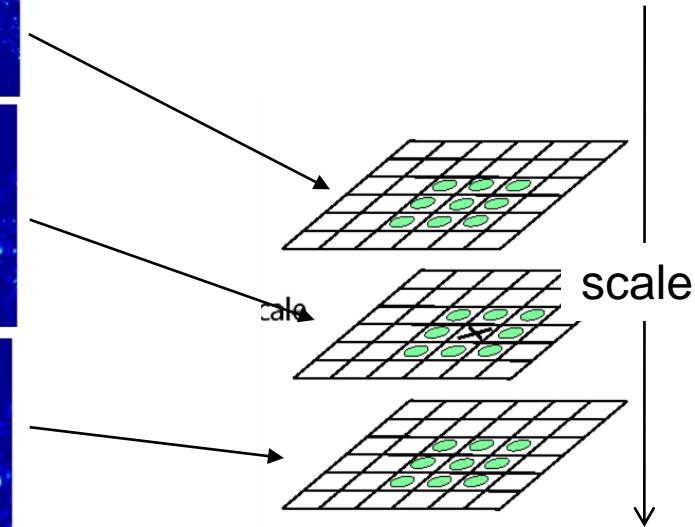
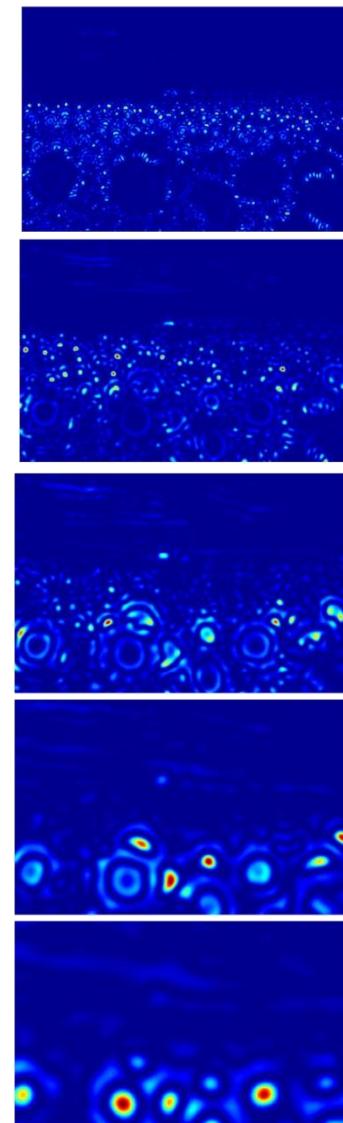
Interest points are local maxima in both position and scale.



$$L_{xx}(\sigma) + L_{yy}(\sigma)$$

$$\rightarrow \sigma_1 \quad \sigma_2 \quad \sigma_3 \quad \sigma_4 \quad \sigma_5$$

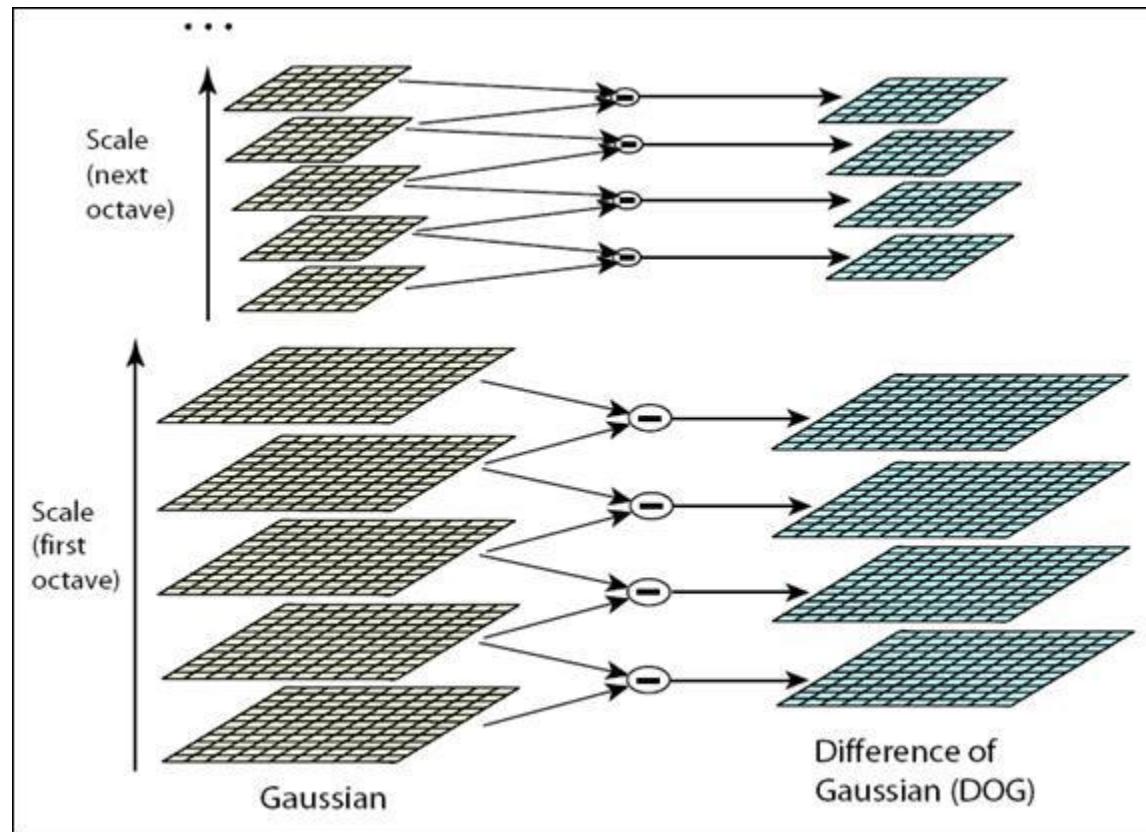
Squared filter response maps



\Rightarrow List of
 (x, y, σ)

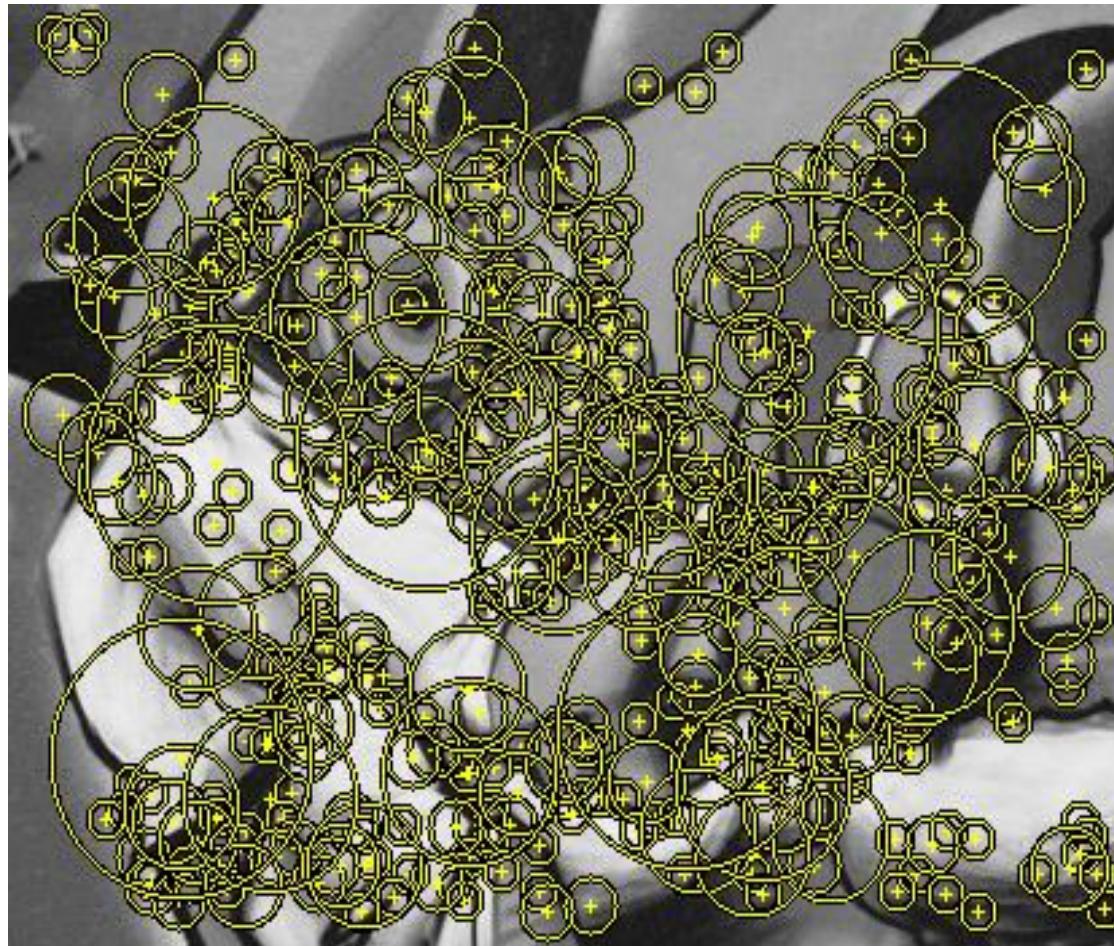
Scale

In practice the image is downsampled for larger sigmas.

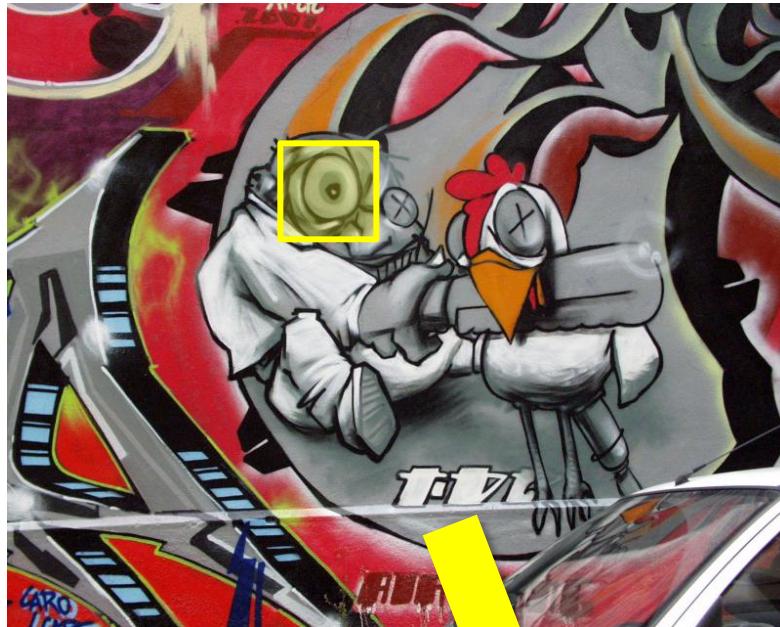


Lowe, 2004.

Results: Difference-of-Gaussian

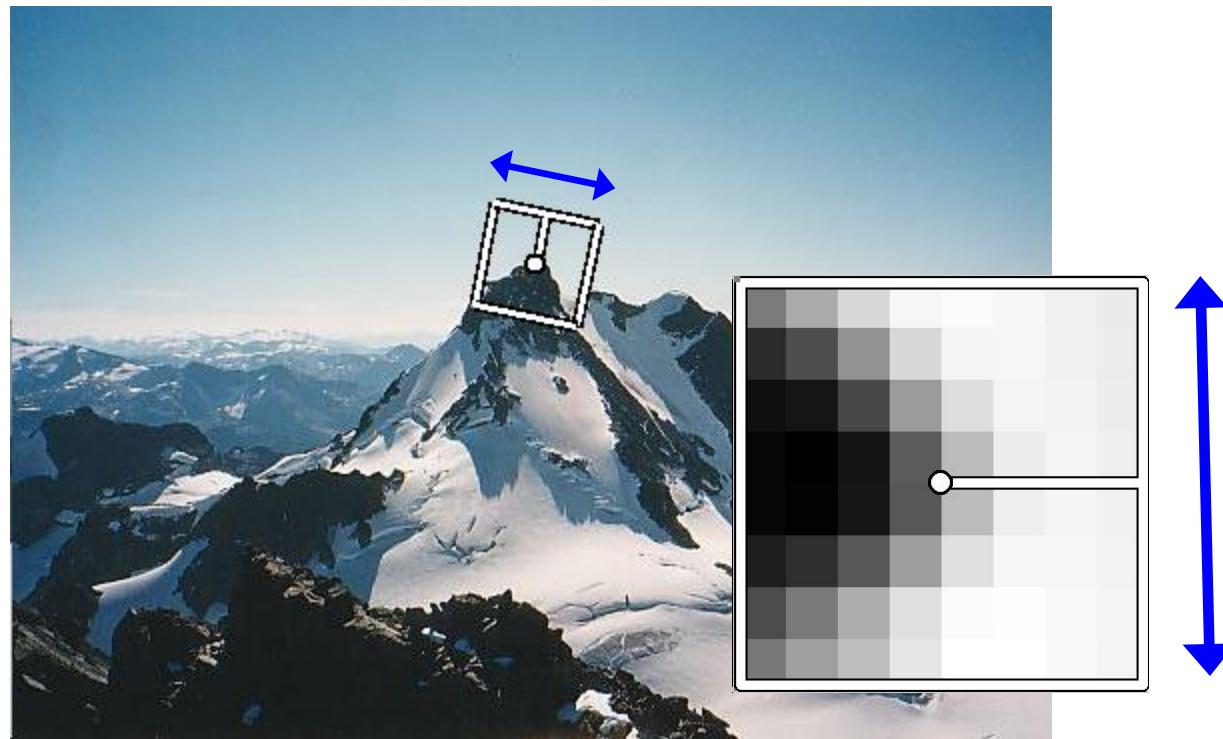


How can we find correspondences?



Similarity transform

Rotation invariance

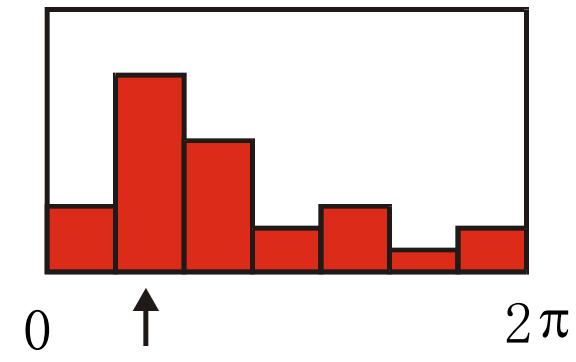
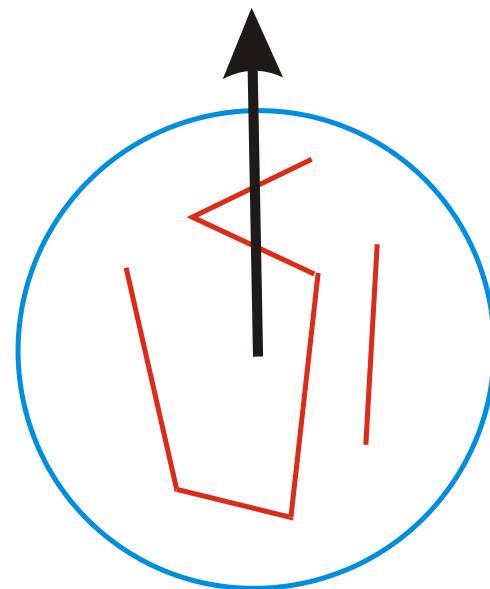


- Rotate patch according to its dominant gradient orientation
- This puts the patches into a canonical orientation.

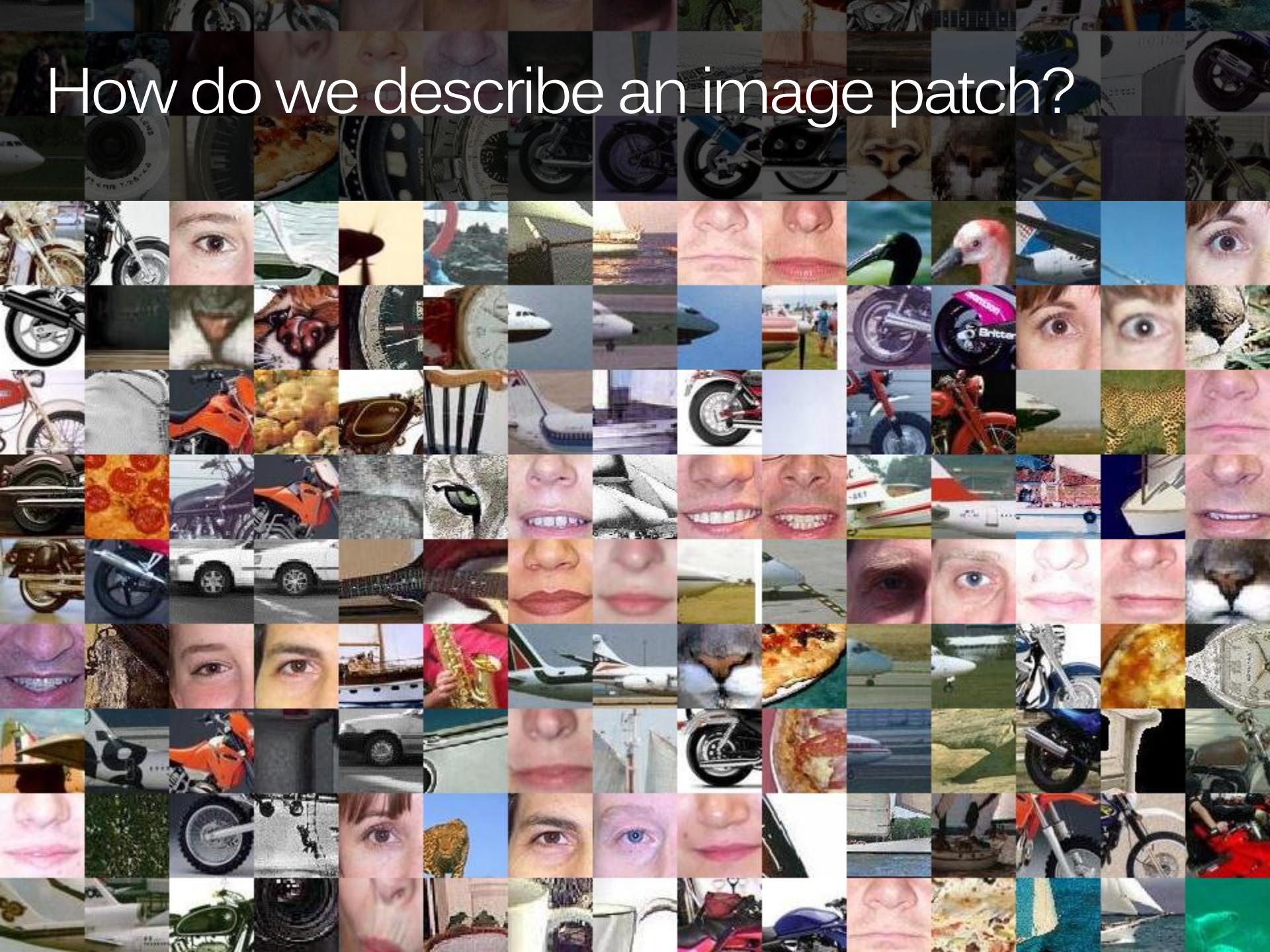
Orientation Normalization

[Lowe, SIFT, 1999]

- Compute orientation histogram
- Select dominant orientation
- Normalize: rotate to fixed orientation

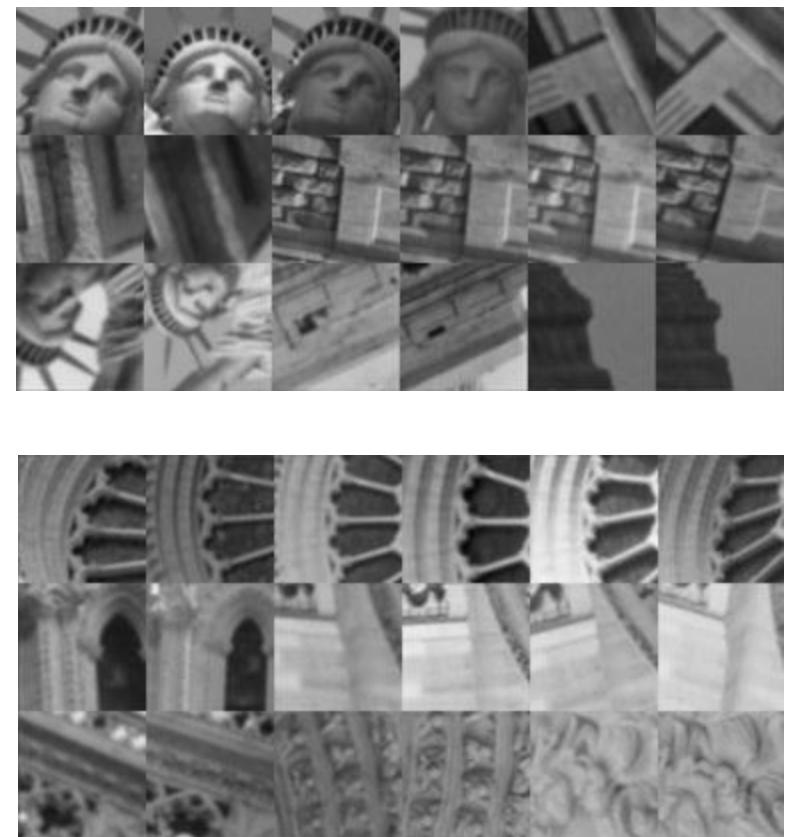


How do we describe an image patch?

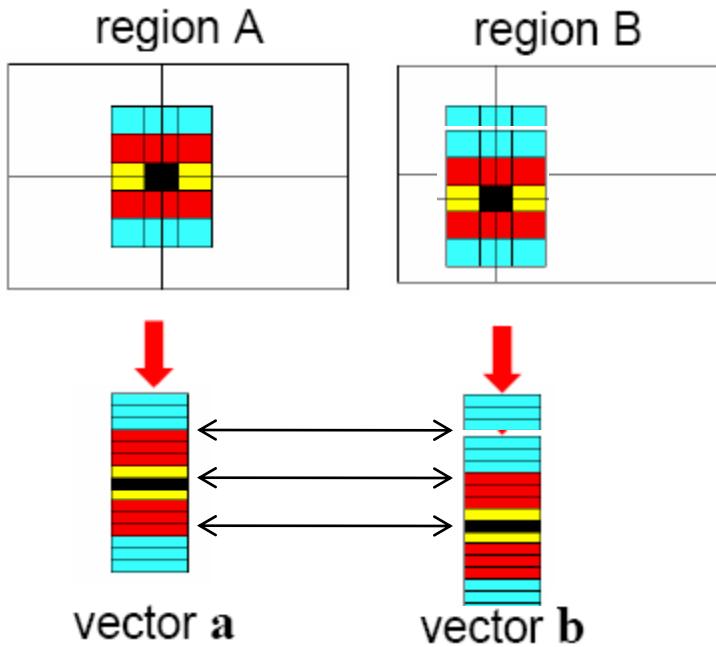


How do we describe an image patch?

Patches with similar content should have similar descriptors.



Raw patches as local descriptors



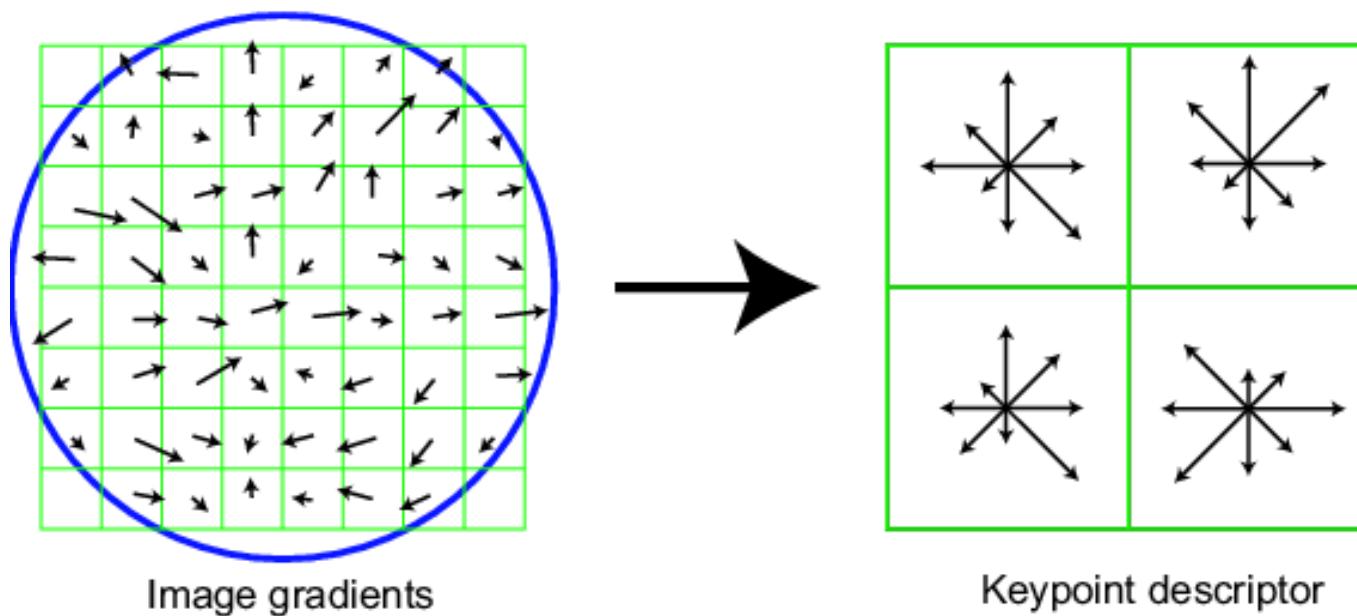
The simplest way to describe the neighborhood around an interest point is to write down the list of intensities to form a feature vector.

But this is very sensitive to even small shifts, rotations.

SIFT descriptor

Full version

- Divide the 16x16 window into a 4x4 grid of cells (2x2 case shown below)
- Compute an orientation histogram for each cell
- 16 cells * 8 orientations = 128 dimensional descriptor



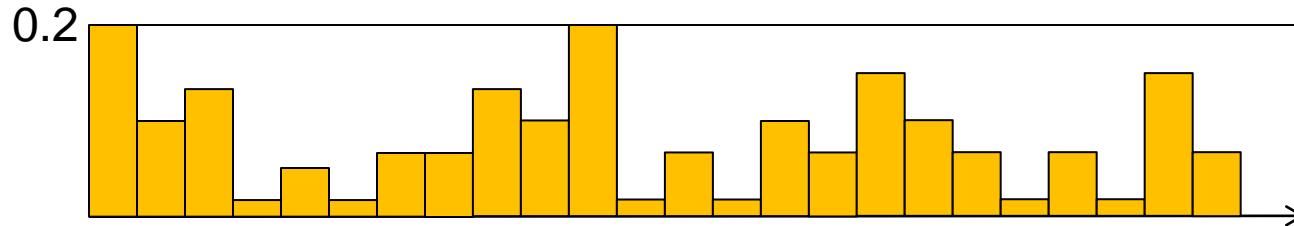
Adapted from slide by David Lowe

SIFT descriptor

Full version

- Divide the 16x16 window into a 4x4 grid of cells (2x2 case shown below)
- Compute an orientation histogram for each cell
- 16 cells * 8 orientations = 128 dimensional descriptor
- Threshold normalize the descriptor:

$$\sum_i d_i^2 = 1 \quad \text{such that: } d_i < 0.2$$

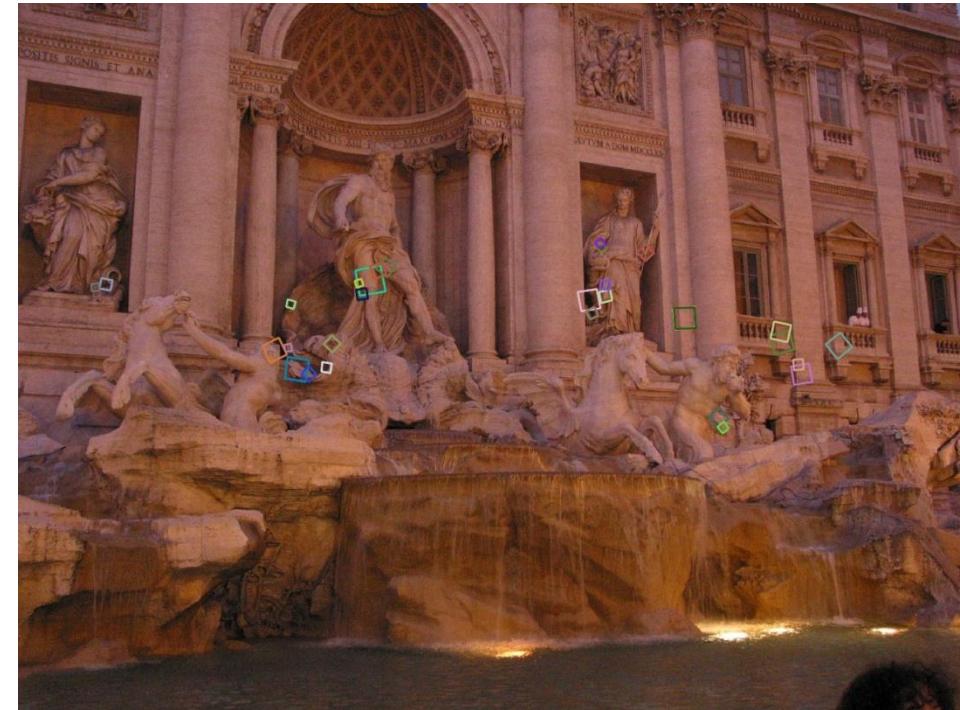


Adapted from slide by David Lowe

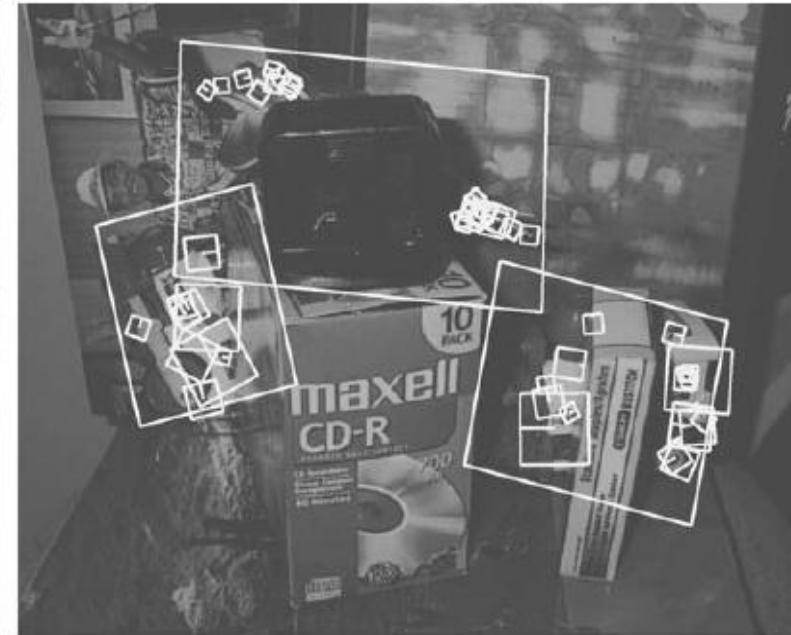
Properties of SIFT

Extraordinarily robust matching technique

- Can handle changes in viewpoint
 - Up to about 30 degree out of plane rotation
- Can handle significant changes in illumination
 - Sometimes even day vs. night (below)
- Fast and efficient—can run in real time
- Lots of code available
 - http://people.csail.mit.edu/albert/ladypack/wiki/index.php/Known_implementations_of_SIFT



Example: Object Recognition



SIFT is extremely powerful for object instance recognition, especially for well-textured objects

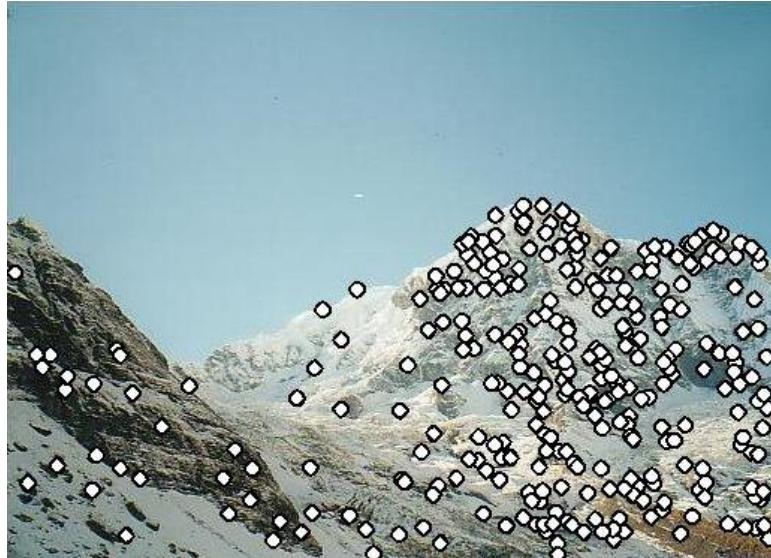
panorama?

- We need to match (align) images



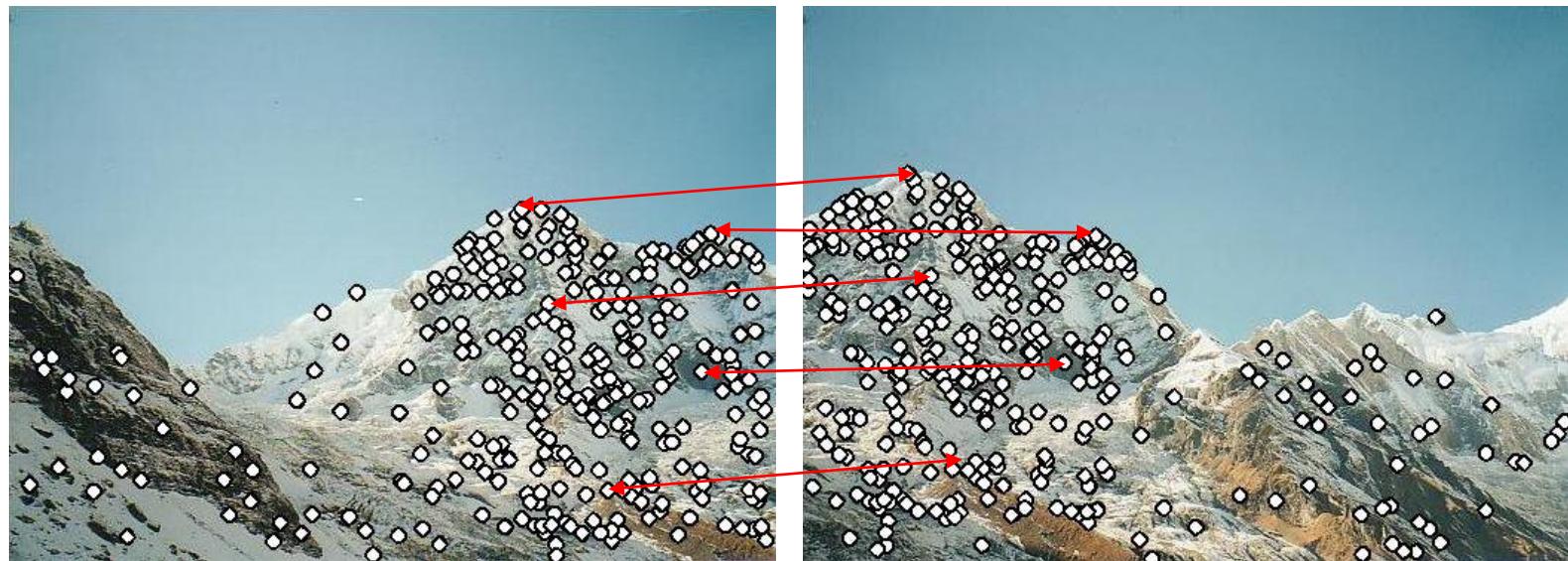
Matching with Features

- Detect feature points in both images



Matching with Features

- Detect feature points in both images
- Find corresponding pairs

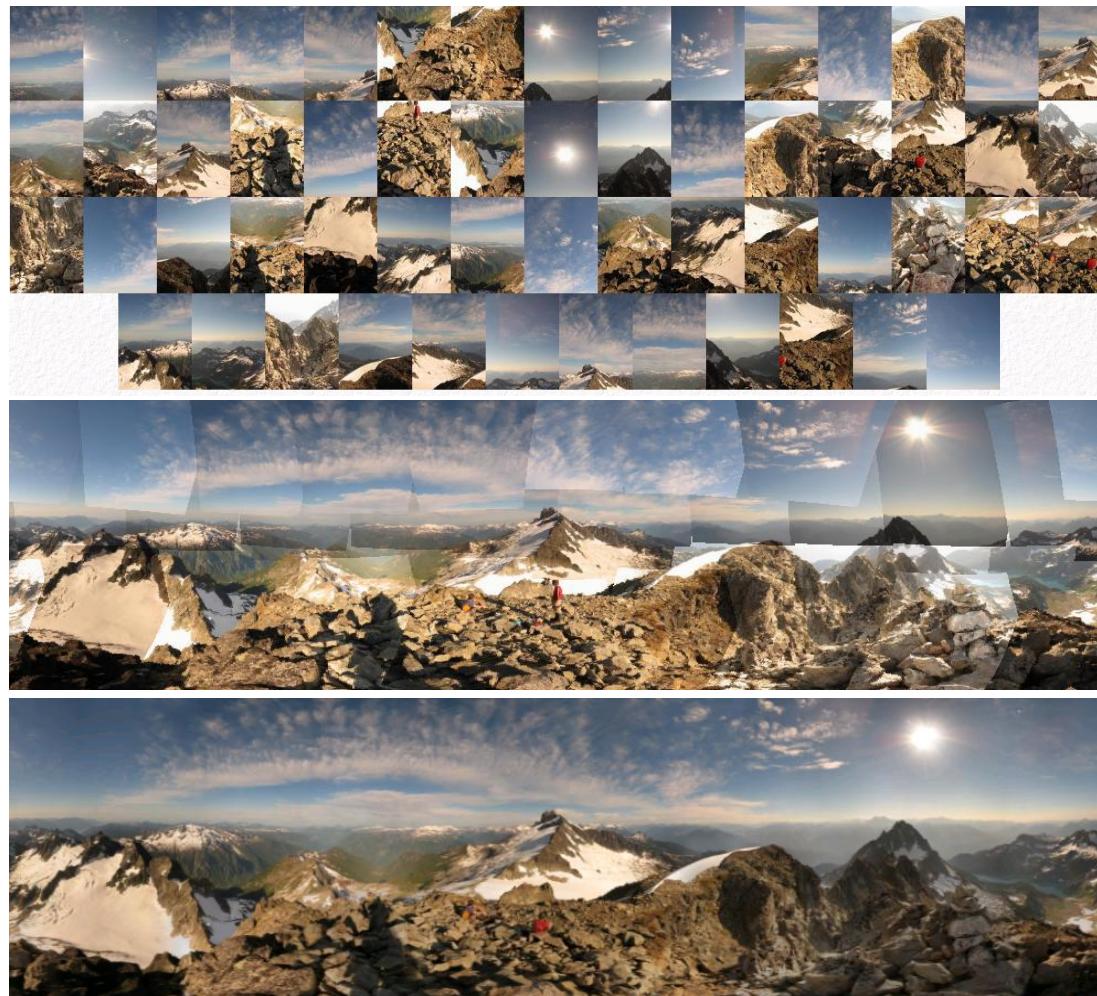


Matching with Features

- Detect feature points in both images
- Find corresponding pairs
- Use these matching pairs to align images - the required mapping is called a homography.



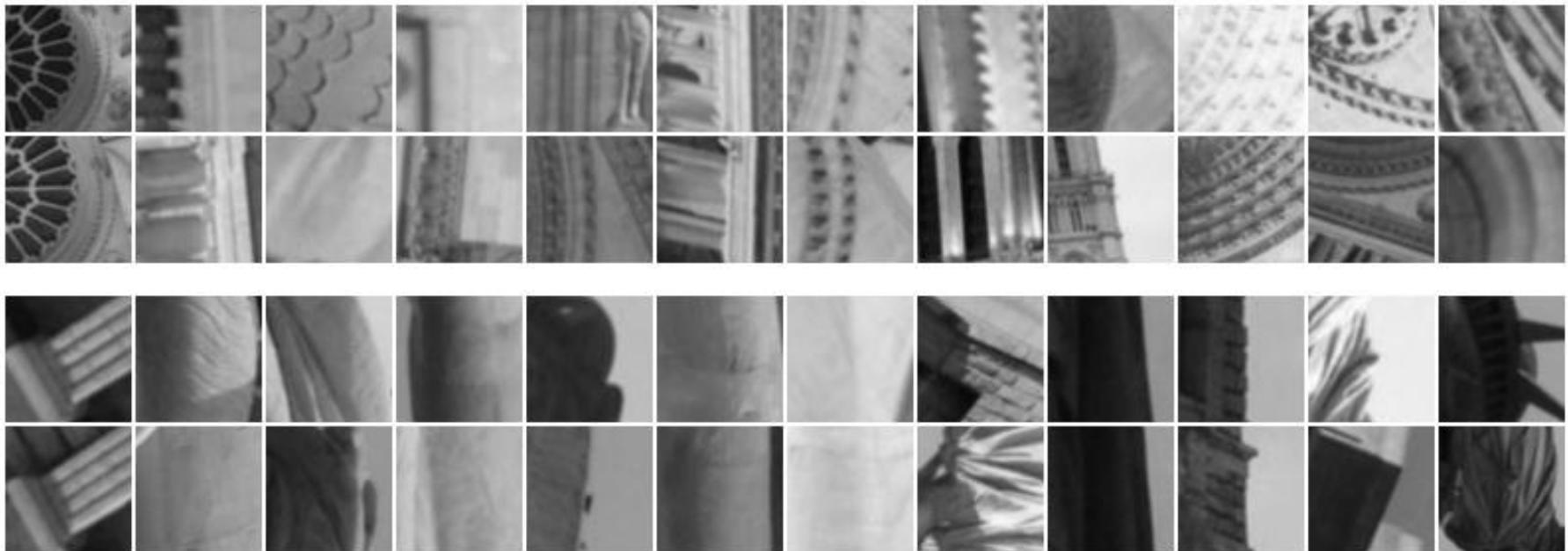
Automatic mosaicing



<http://www.cs.ubc.ca/~mbrown/autostitch/autostitch.html>

When does SIFT fail?

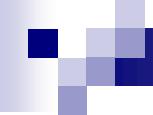
Patches SIFT thought were the same but aren't:



SURF: Speeded Up Robust Features

Background

- Local invariant Interest point detector-descriptor
 - For finding *correspondences* between two images of the same scene or object
 - Many applications, including 3D reconstruction, image retrieval and object recognition
 - SIFT is one of the best but slow
 - Image of size 1000 x 700 described in around **6** seconds (actual cost depends on the # features generated, 4000 in this case)
 - 128-D feature vectors

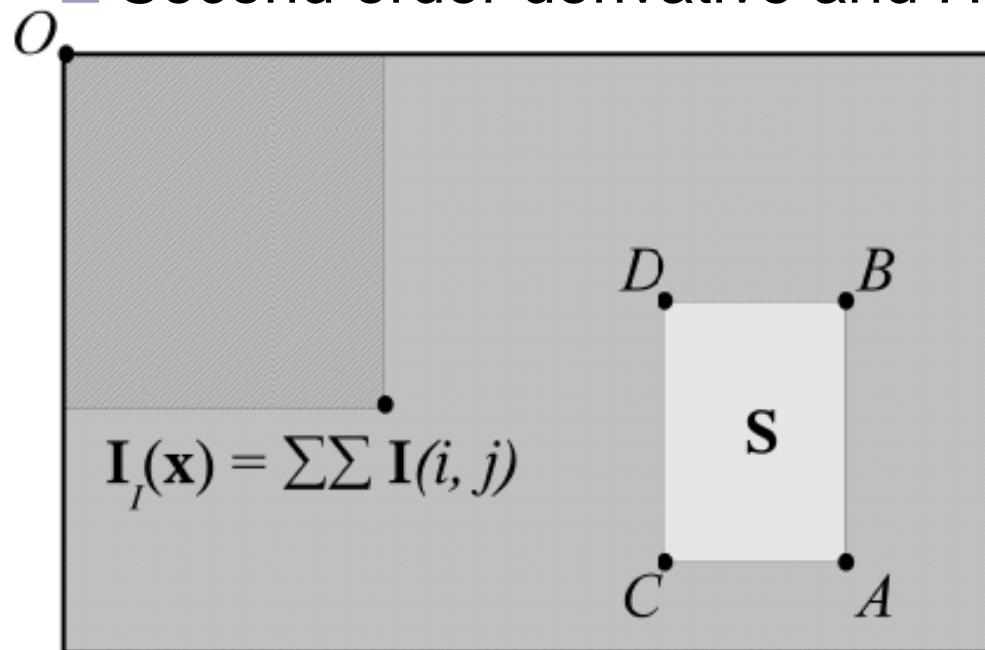


Motivation

- Fast interest point **detection**
- Distinctive interest point **description**
- Speeded-up descriptor **matching**
- Invariant to common image transformations:
 - ✓ Image rotation
 - ✓ Scale changes
 - ✓ Illumination change
 - ✓ *Small* change in Viewpoint

Methodology

- Using integral images for major speed up
 - Integral Image (summed area tables) is an intermediate representation for the image and contains the **sum of gray scale pixel values of image**
 - Second order derivative and Haar-wavelet response



$$\mathbf{S} = A - B - C + D$$

Cost *four* additions
operation only

Detection

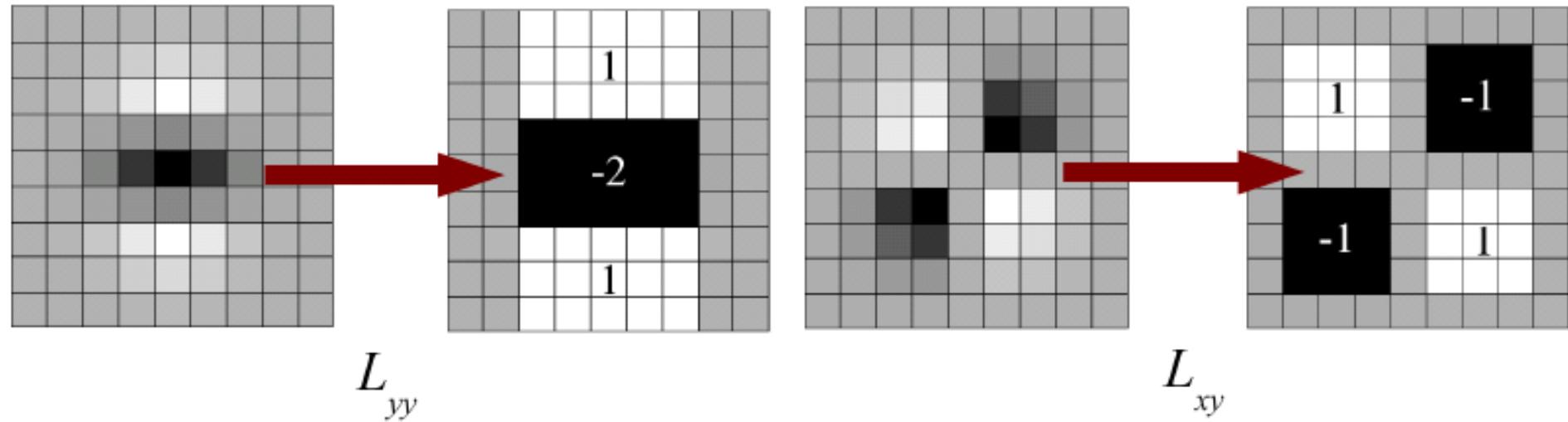
- Hessian-based interest point localization

$$H = \begin{bmatrix} L_{xx} & L_{xy} \\ L_{xy} & L_{yy} \end{bmatrix}$$

- $L_{xx}(x,y,\sigma)$ is the **Laplacian of Gaussian** of the image
- It is the convolution of the *Gaussian* second order derivative with the image
- Lindeberg showed Gaussian function is optimal for scale-space analysis
- This paper argues that Gaussian is overrated since the **property that no new structures can appear while going to lower resolution** is not proven in 2D case

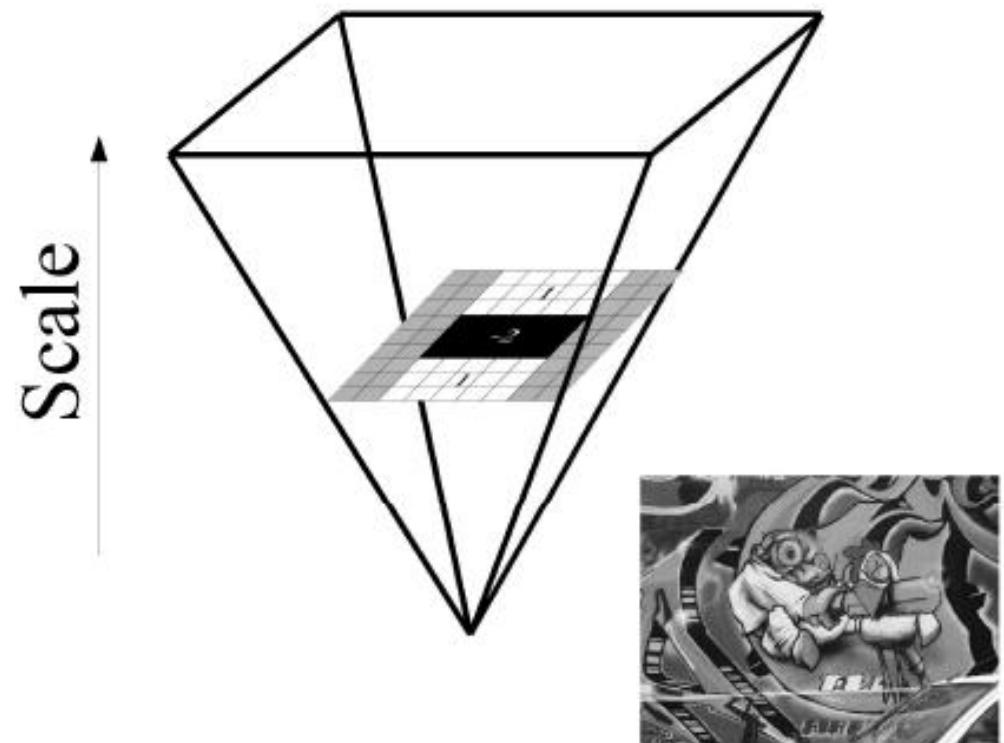
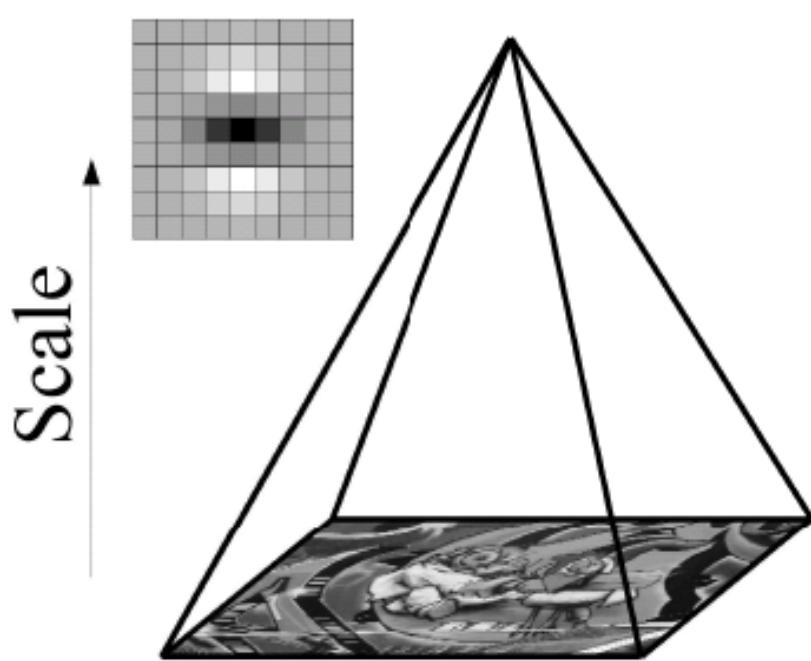
Detection

- Approximated second order derivatives with box filters (mean/average filter)



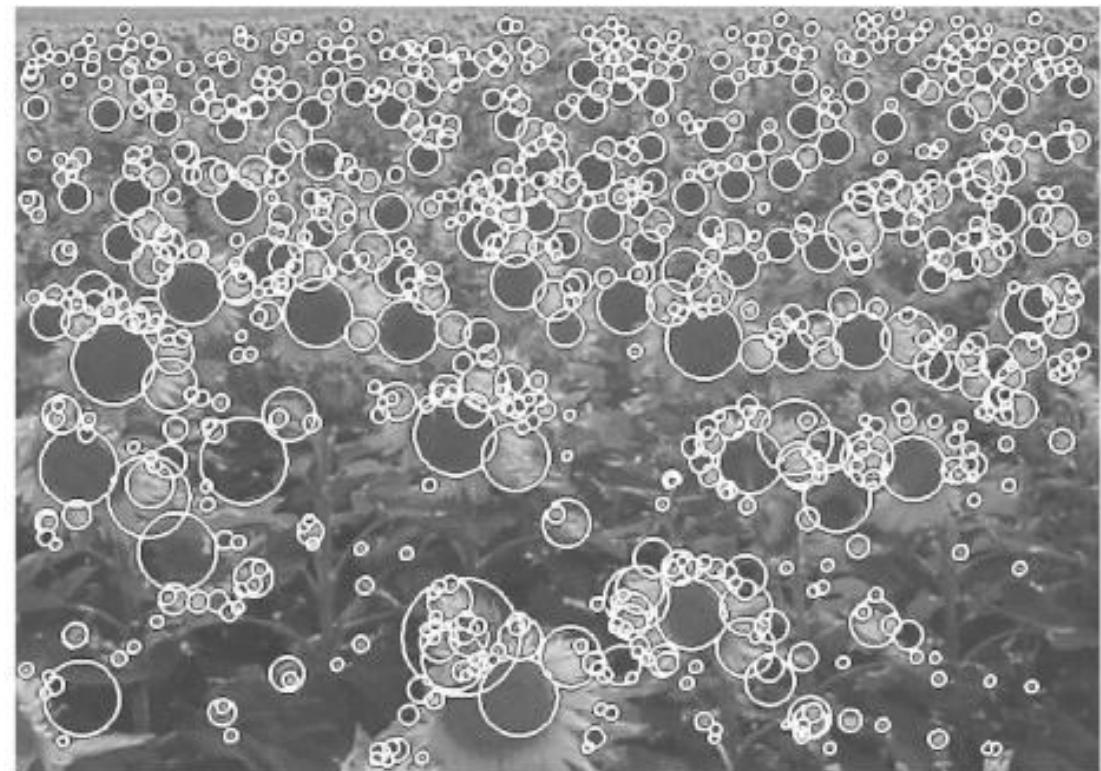
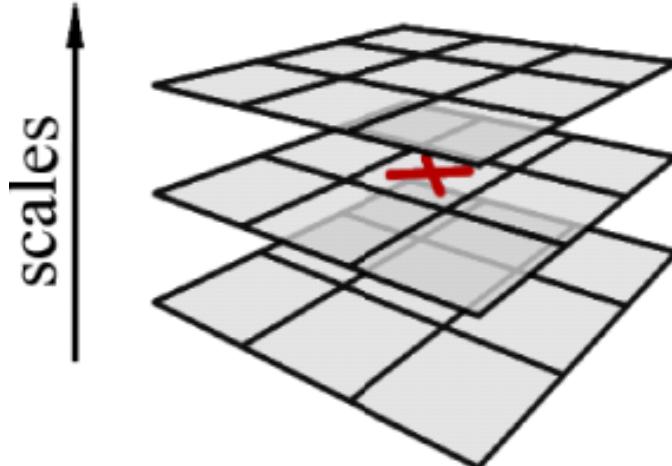
Detection

- ## ■ Scale analysis with constant image size



Detection

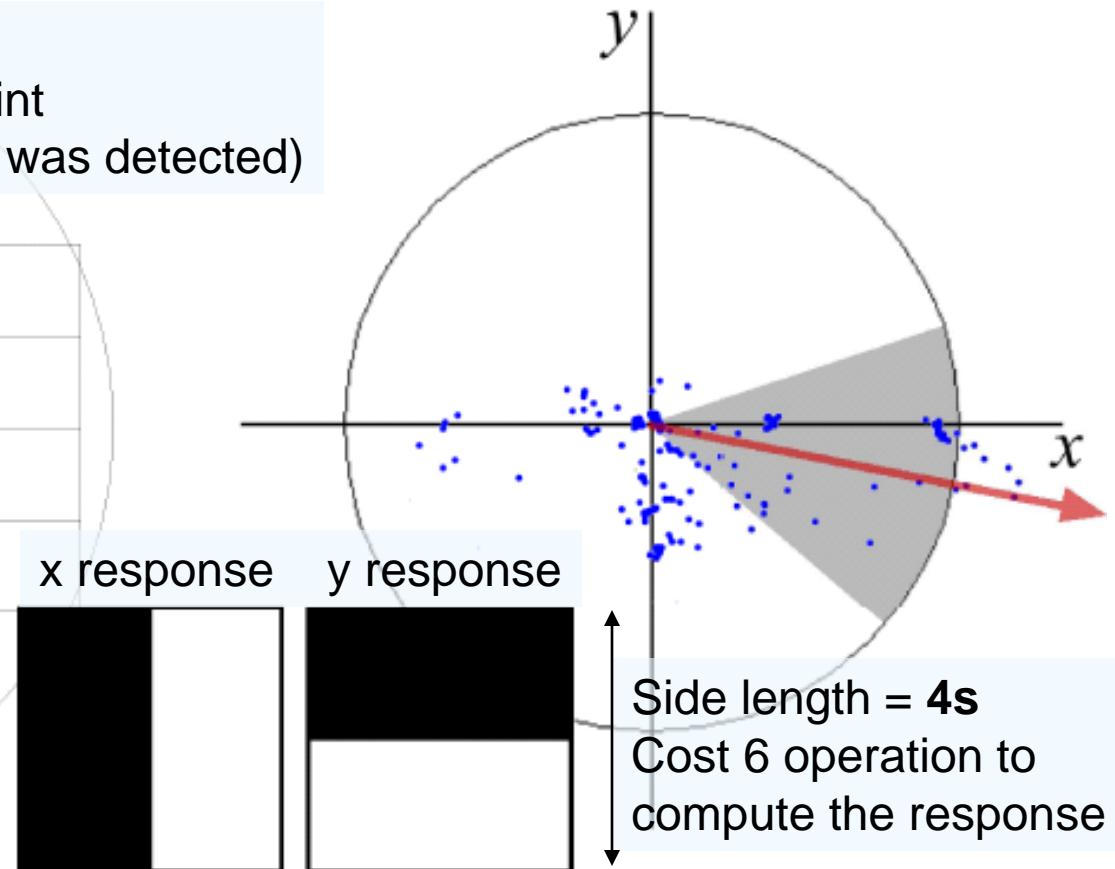
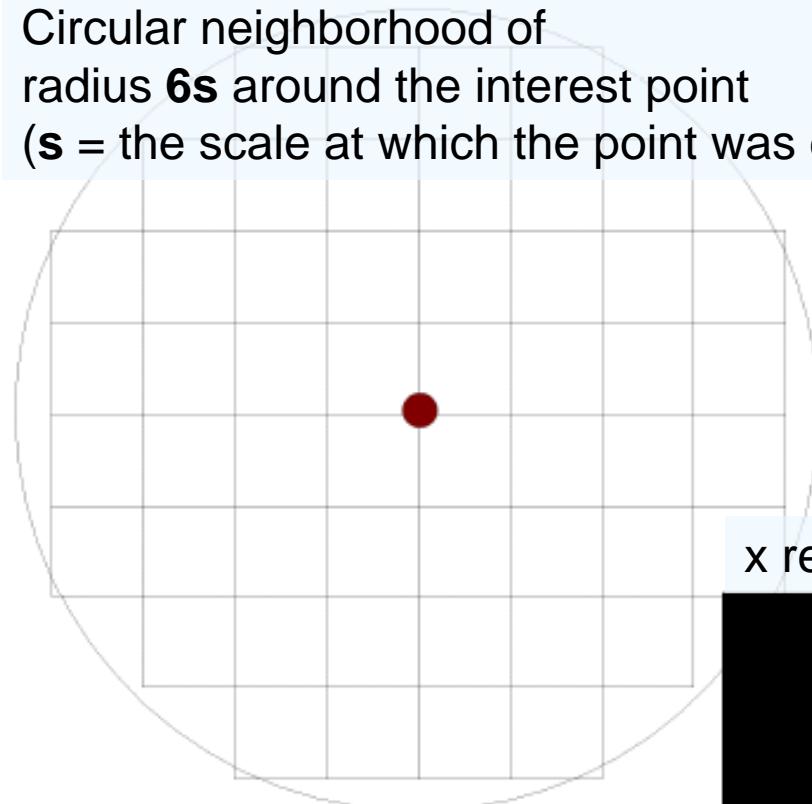
- Non-maximum suppression and interpolation
 - Blob-like feature detector



Description

■ Orientation Assignment

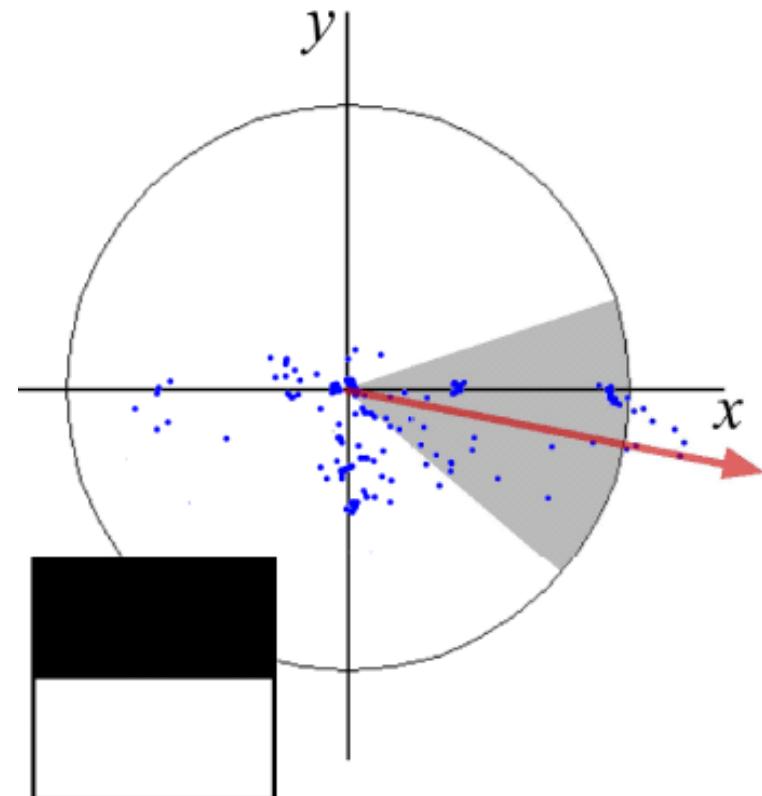
Circular neighborhood of radius **6s** around the interest point (**s** = the scale at which the point was detected)



Description

■ Dominant orientation

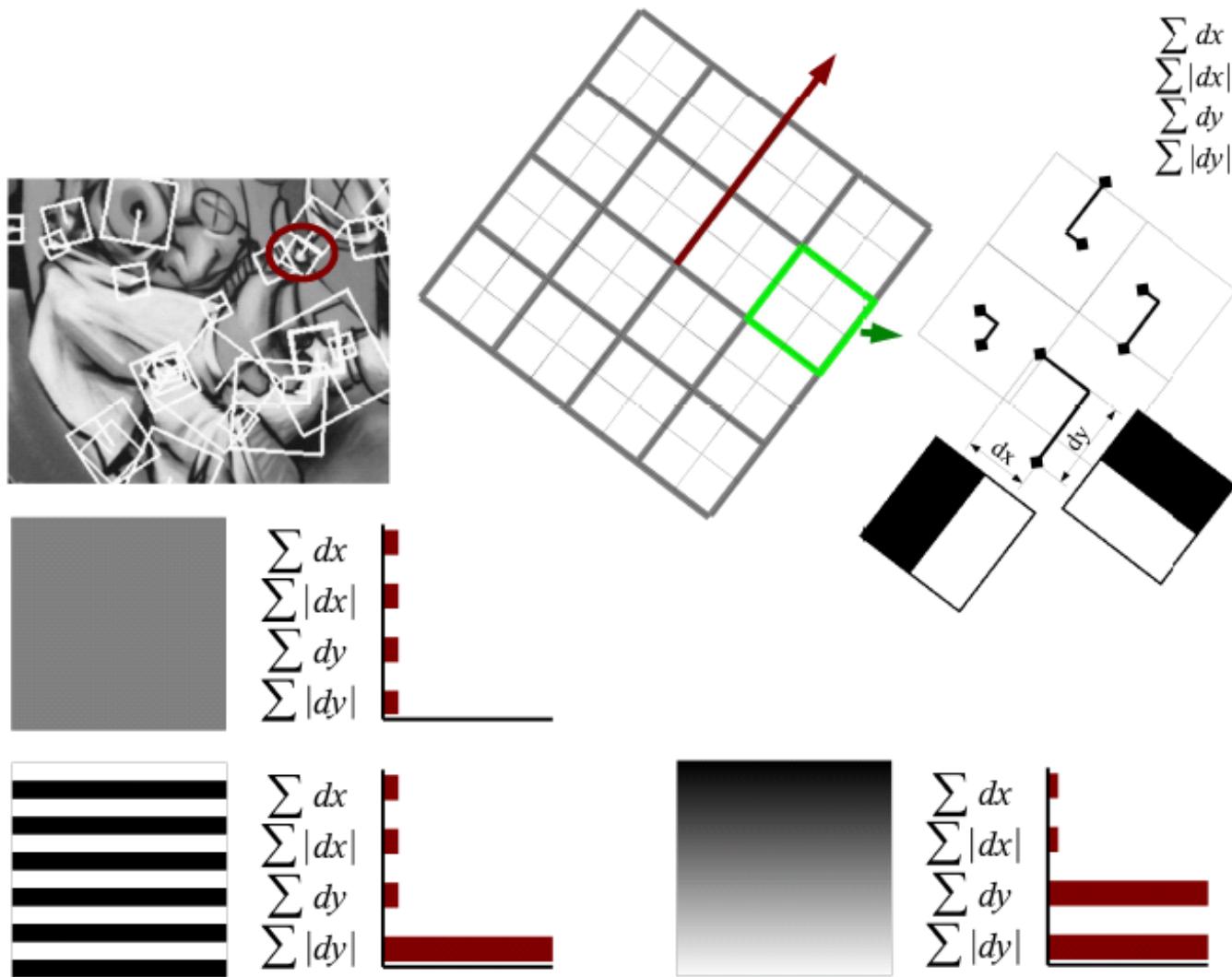
- The Haar wavelet responses are represented as vectors
- Sum all responses within a sliding orientation window covering an angle of 60 degree
- The two summed response yield a new vector
- **The longest vector** is the dominant orientation
- Second longest is ...
ignored



Description

- Split the interest region up into 4×4 square sub-regions with 5×5 regularly spaced sample points inside
- Calculate Haar wavelet response d_x and d_y
- Weight the response with a Gaussian kernel centered at the interest point
- Sum the response over each sub-region for d_x and d_y separately → **feature vector of length 32**
- In order to bring in information about the polarity of the intensity changes, extract the sum of absolute value of the responses → **feature vector of length 64**
- Normalize the vector into unit length

Description

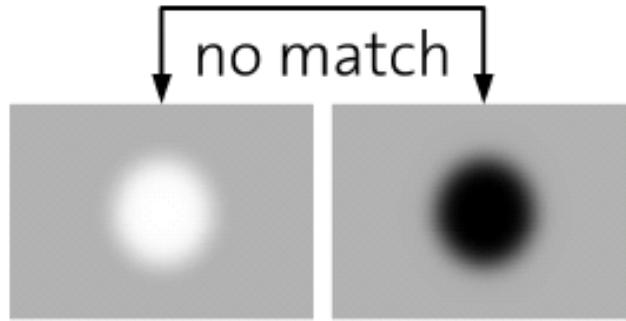


Description

■ SURF-128

- The sum of d_x and $|d_x|$ are computed separately for $d_y < 0$ and $d_y > 0$
- Similarly for the sum of d_y and $|d_y|$
- This doubles the length of a feature vector

Matching

- Fast indexing through the sign of the Laplacian for the underlying interest point
 - The sign of trace of the Hessian matrix
 - $\text{Trace} = L_{xx} + L_{yy}$
 - Either 0 or 1 (Hard thresholding, may have boundary effect ...)
 - In the matching stage, compare features if they have the same type of contrast (sign)
- 

Analysis

- SURF is good at
 - handling serious blurring
 - handling image rotation
- SURF is poor at
 - handling viewpoint change
 - handling illumination change
- SURF is always better than the SIFT implemented by VGG but not the original SIFT

img#	Bikes	Boat	graf	leuven	wall
2	o	++	--	-	---
3	o	o	-	--	---
4	+	+++	-	--	---
5	++	+++	o	--	o
6	+++	+++	o	---	o

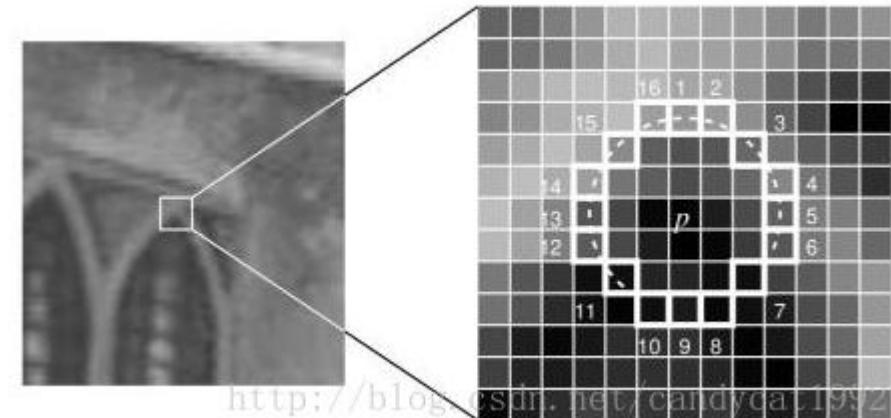
Conclusion

- SURF describes image faster than SIFT by 3 times
- SURF is not as well as SIFT on invariance to illumination change and viewpoint change

Other methods: Fast corner Detection

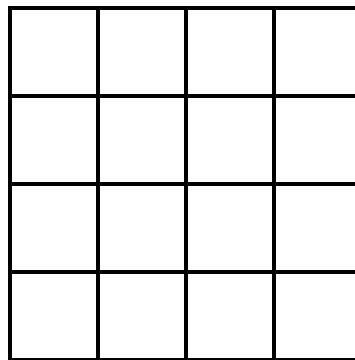
1. 从图片中选取一个像素点 p , 下面我们将判断它是否是一个特征点。我们首先把它的密度(即像素值)设为 l_p 。
2. 设定一个合适的阈值 t 。
3. 考虑该像素点周围的16个像素。
4. 如果在这个大小为16个像素的圆上有 n 个连续的像素点, 它们的像素值要么都比 l_p+t 大, 要么都比 l_p-t 小, 那么它就是一个角点。(如上图中白色虚线所示)。 n 这里被设定为12。
5. 检查在位置1、9、5和13四个位置的像素(首先检查1和9, 看它们是否亮于或暗于阈值。如果是, 再检查5和13)。如果是一个角点, 那么上述四个像素点中至少有3个应该必须都大于或者小于(因为若是一个角点, 超过四分之三个圆的部分应该满足判断条件, 半圆比包含上述某三个点)。如果不满足, 那么不可能是一个角点。完整的分段测试可以被用于接受的所有候选点, 通过检测圆上的所有点。

总结: **FAST**算法比其他已知的角点检测法要快很多倍。但是当图片的噪声较多时, 它的健壮性并不好。这依靠一个阈值。

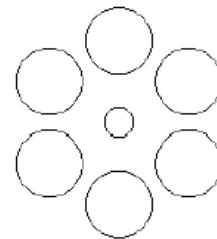


Other methods: Daisy

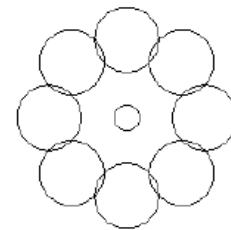
Circular gradient binning



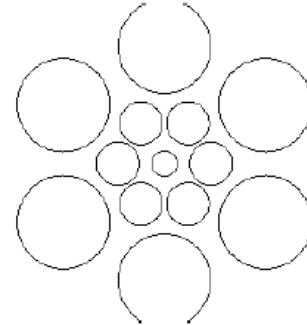
SIFT



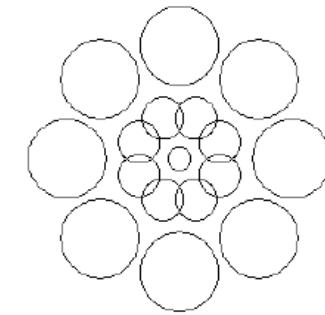
1 Ring 6 Segments



1 Ring 8 Segments



2 Rings 6 Segments



2 Rings 8 Segments

Daisy

Other methods: HOG

- **Hog:** Histogram of Oriented Gradient, 通过计算和统计图像局部区域的梯度方向直方图来构成特征。
- **本质:** 梯度的统计信息，而梯度主要存在于边缘的地方。
- **实现方法:** 首先将图像分成小的连通区域，我们把它叫细胞单元。然后采集细胞单元中各像素点的梯度的或边缘的方向直方图。最后把这些直方图组合起来就可以构成特征描述器。
- **优点:** 首先，由于**HOG**是在图像的局部方格单元上操作，所以它对图像几何的和光学的形变都能保持很好的不变性，这两种形变只会出现在更大的空间领域上。其次，在粗的空域抽样、精细的方向抽样以及较强的局部光学归一化等条件下，只要行人大体上能够保持直立的姿势，可以容许行人有一些细微的肢体动作，这些细微的动作可以被忽略而不影响检测效果。因此**HOG**特征是特别适合于做图像中的人体检测的。
- **应用:** **Hog**特征结合**SVM**分类器已经被广泛应用于图像识别中，尤其在行人检测中获得了极大的成功。

Other methods: BRIEF

Randomly sample pair of pixels a and b.
1 if $a > b$, else 0. Store binary vector.

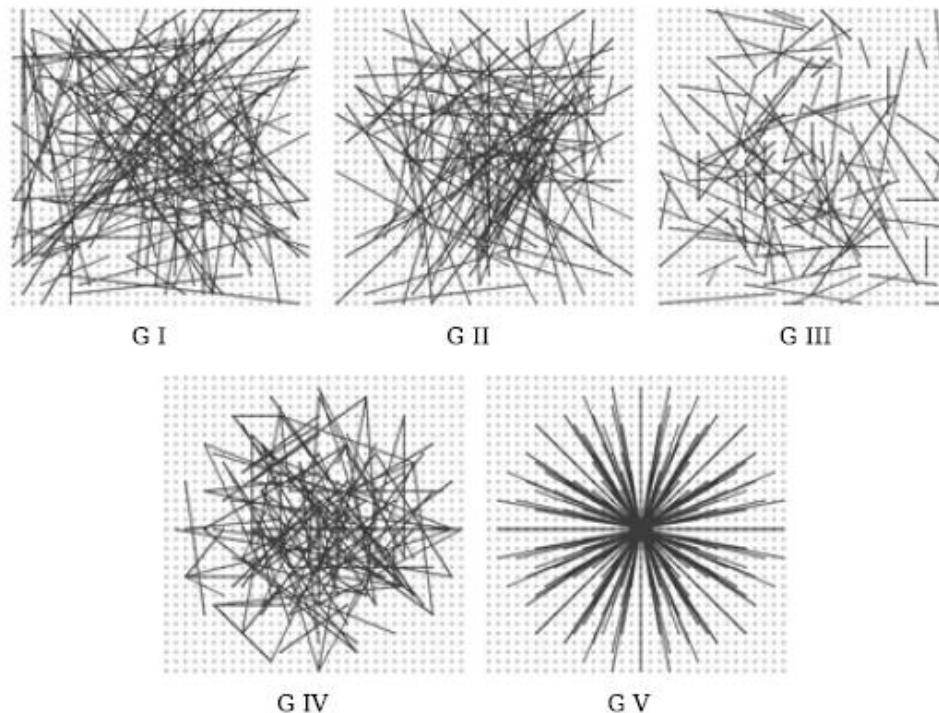


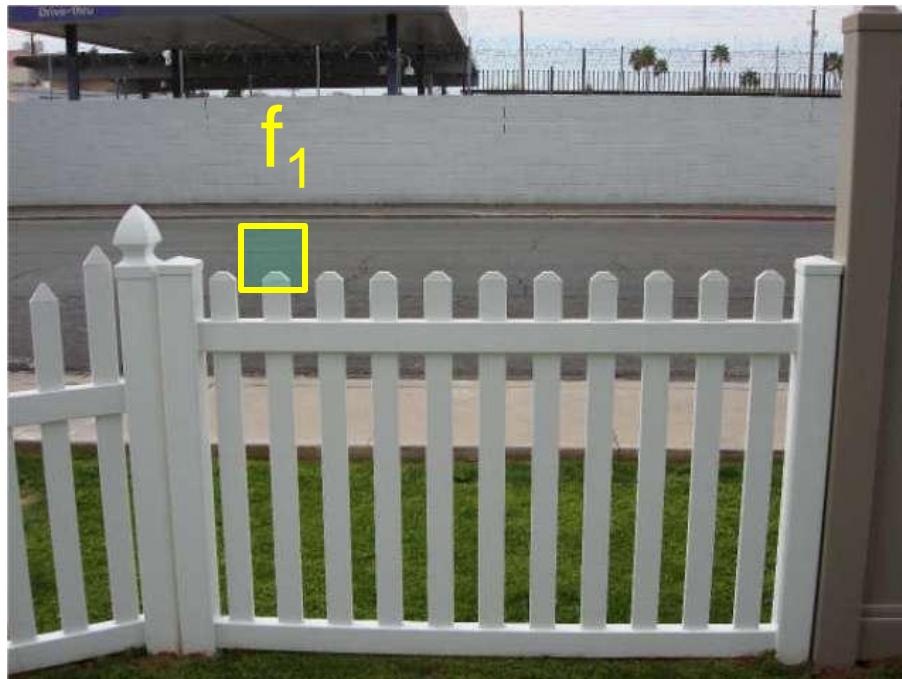
Fig. 2. Different approaches to choosing the test locations. All except the rightmost one are selected by random sampling. Showing 128 tests in every image.

BRIEF: binary robust independent elementary features,
Calonder, V Lepetit, C Strecha, ECCV 2010

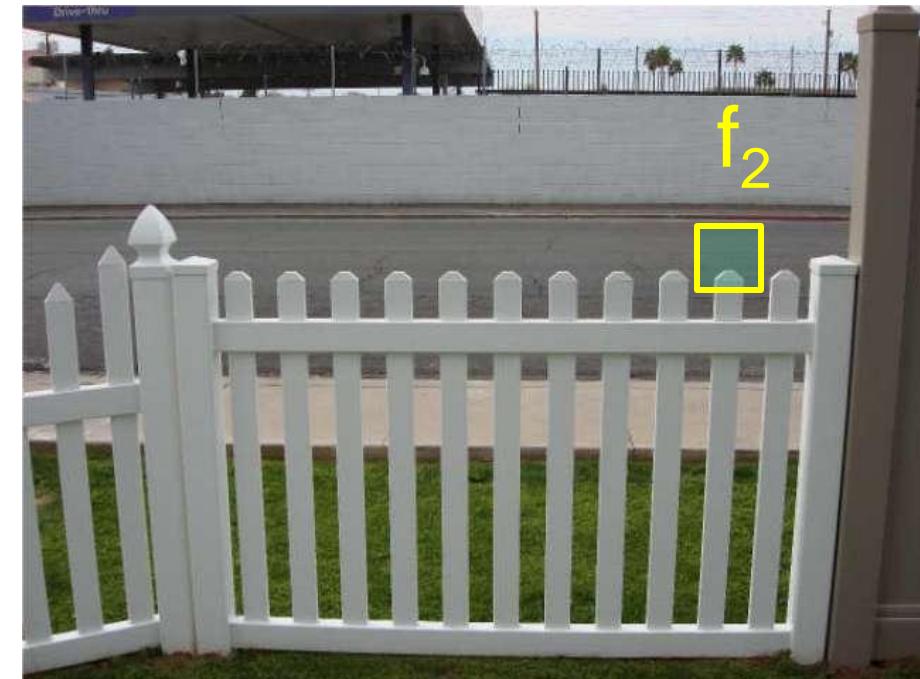
Feature distance

How to define the difference between two features f_1 , f_2 ?

- Simple approach is $\text{SSD}(f_1, f_2)$
 - sum of square differences between entries of the two descriptors
 - can give good scores to very ambiguous (bad) matches



$|_1$

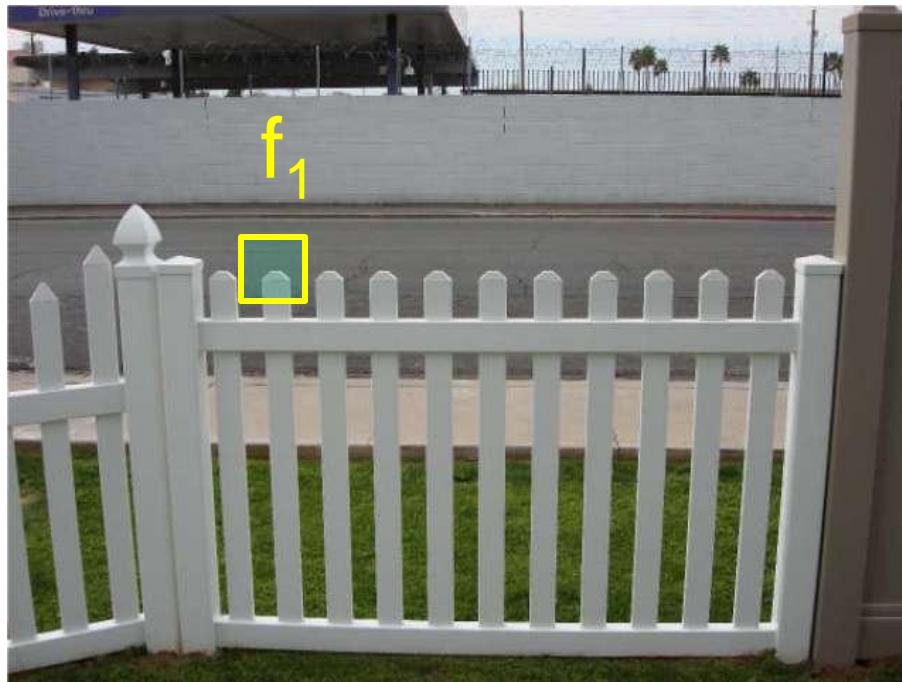


$|_2$

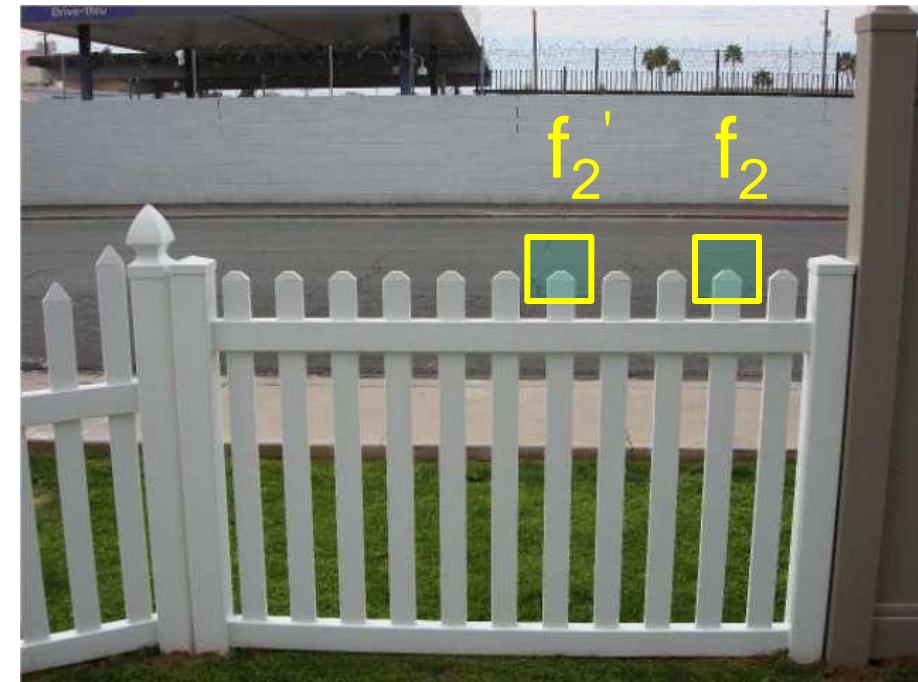
Feature distance

How to define the difference between two features f_1, f_2 ?

- Better approach: ratio distance = $\text{SSD}(f_1, f_2) / \text{SSD}(f_1, f_2')$
 - f_2 is best SSD match to f_1 in I_2
 - f_2' is 2nd best SSD match to f_1 in I_2
 - gives large values (~1) for ambiguous matches

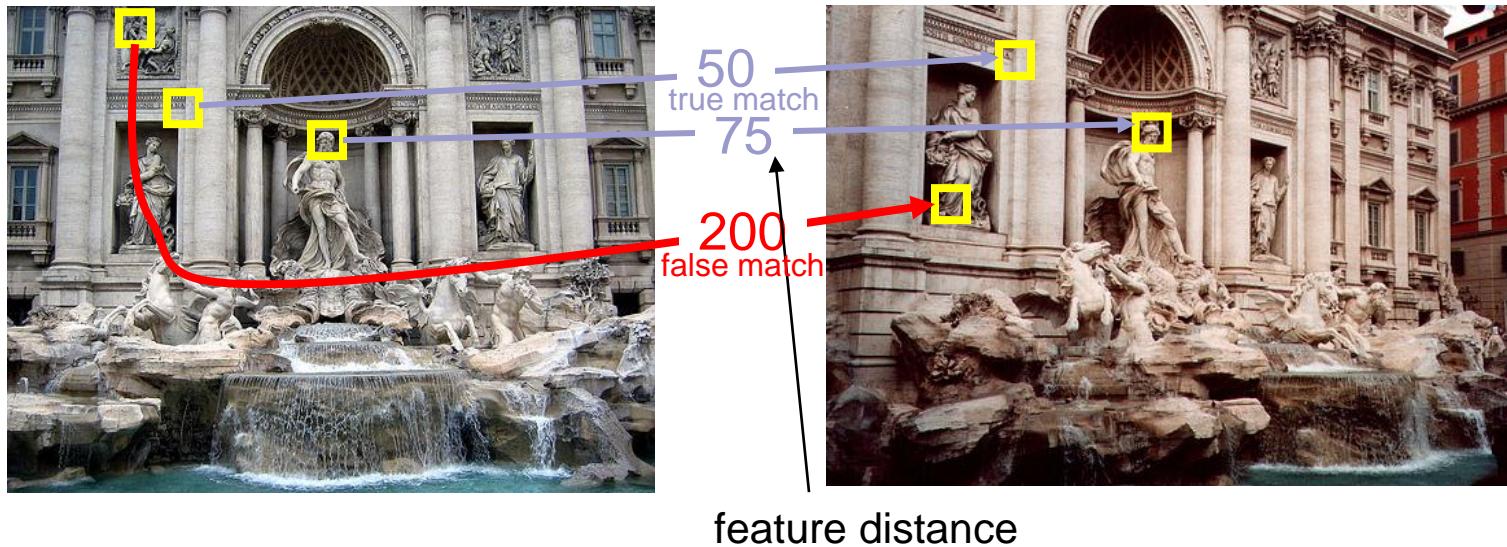


I_1



I_2

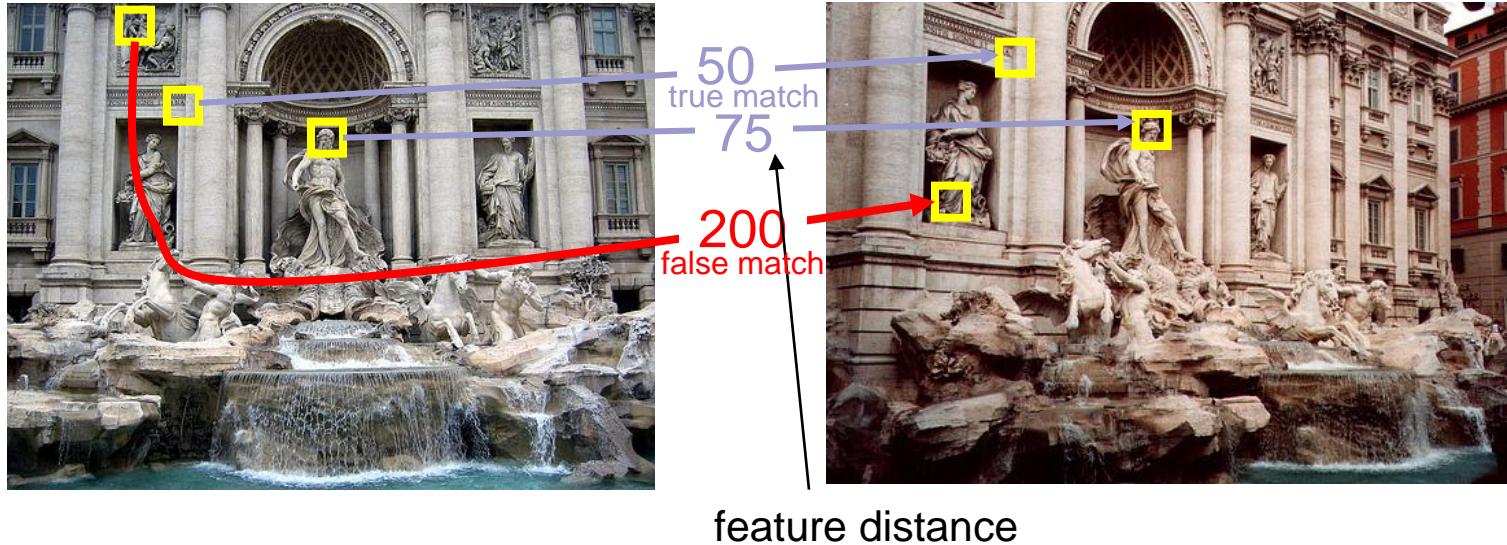
Eliminating bad matches



Throw out features with $\text{distance} > \text{threshold}$

- How to choose the threshold?

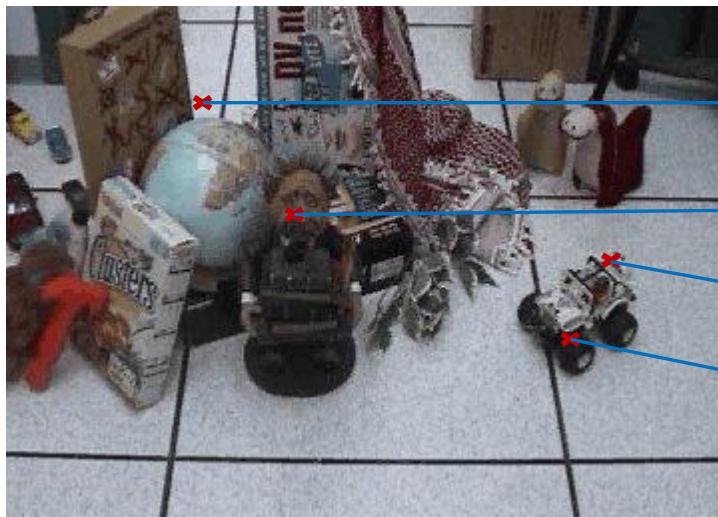
True/false positives



The distance threshold affects performance

- True positives = # of detected matches that are correct
 - Suppose we want to maximize these—how to choose threshold?
- False positives = # of detected matches that are incorrect
 - Suppose we want to minimize these—how to choose threshold?

What is Optical Flow ?



Where did **each pixel** in image 1 go to in image 2:

Goal:

Find for **each pixel** a velocity vector $\vec{u} = (u, v)$ which says:

- How quickly(**magnitude**) is the pixel moving across the image
- In which **direction** it is moving

Dense!

We live in a moving world

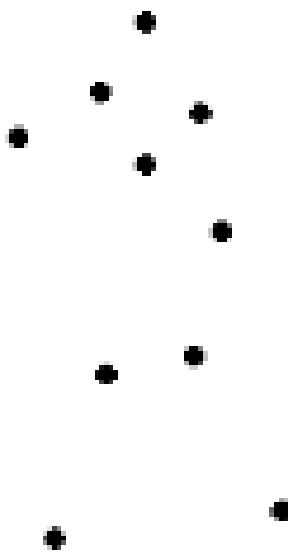
- Perceiving, understanding and predicting motion is an important part of our daily lives



from “<http://szeliski.org/Book/>”

Motion and perceptual organization

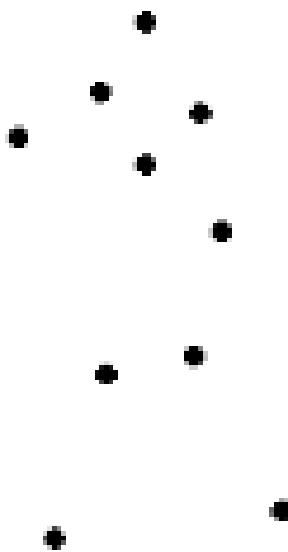
- Even “impoverished” motion data can evoke a strong percept



G. Johansson, “Visual Perception of Biological Motion and a Model For Its Analysis”, *Perception and Psychophysics* 14, 201-211, 1973.

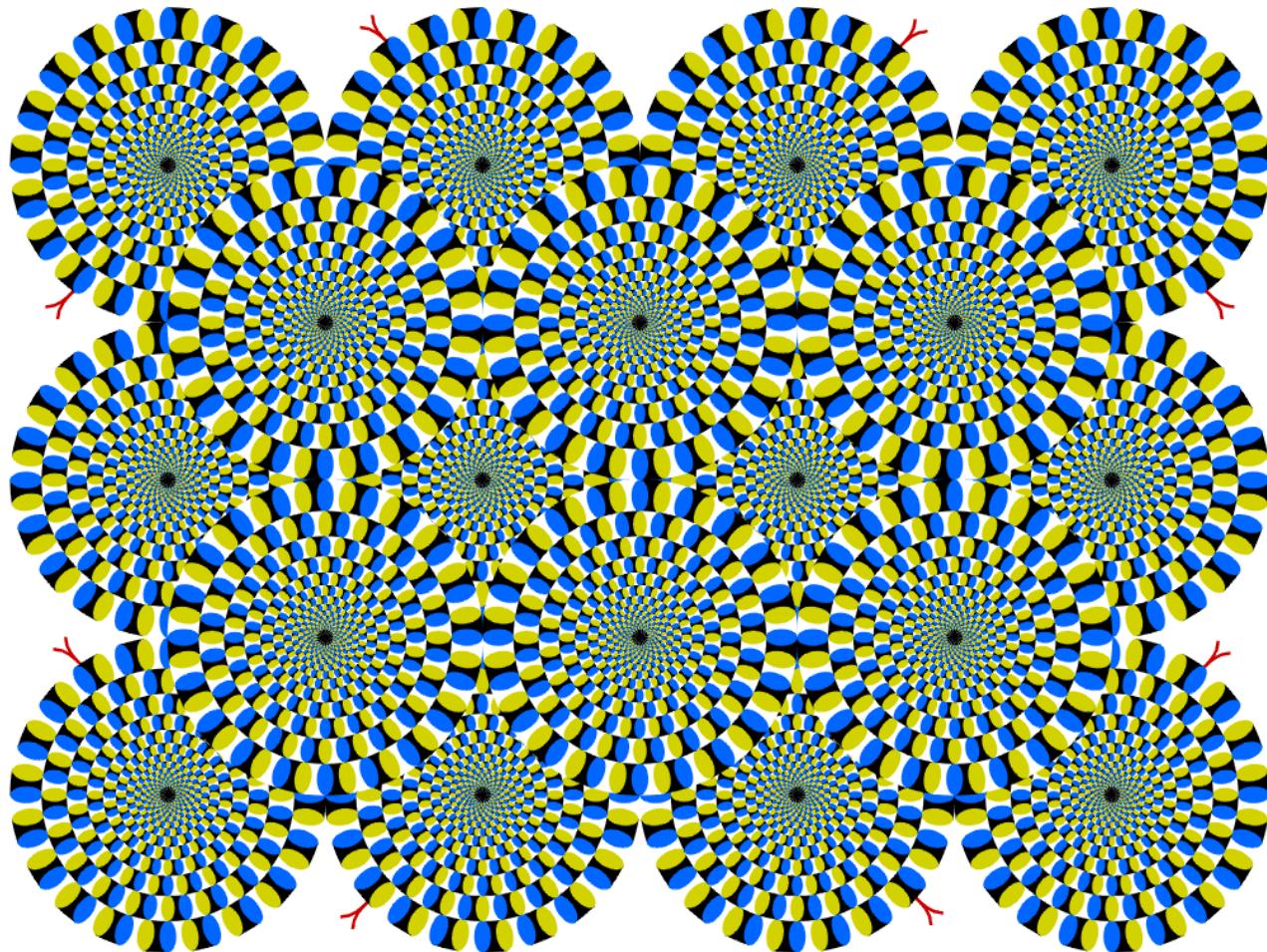
Motion and perceptual organization

- Even “impoverished” motion data can evoke a strong percept

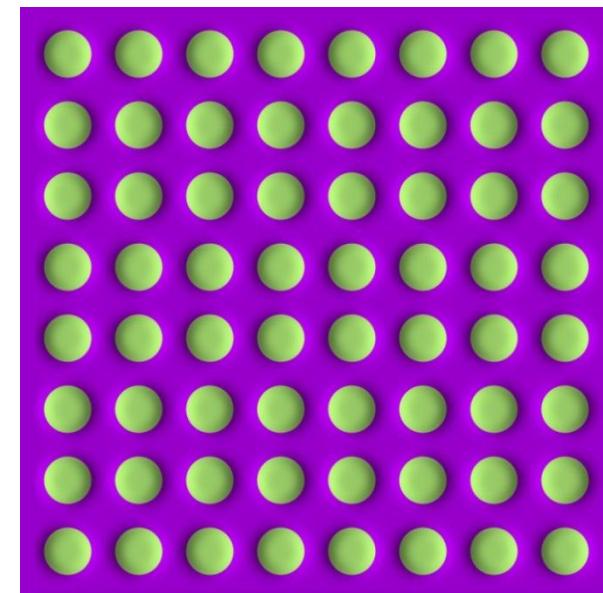
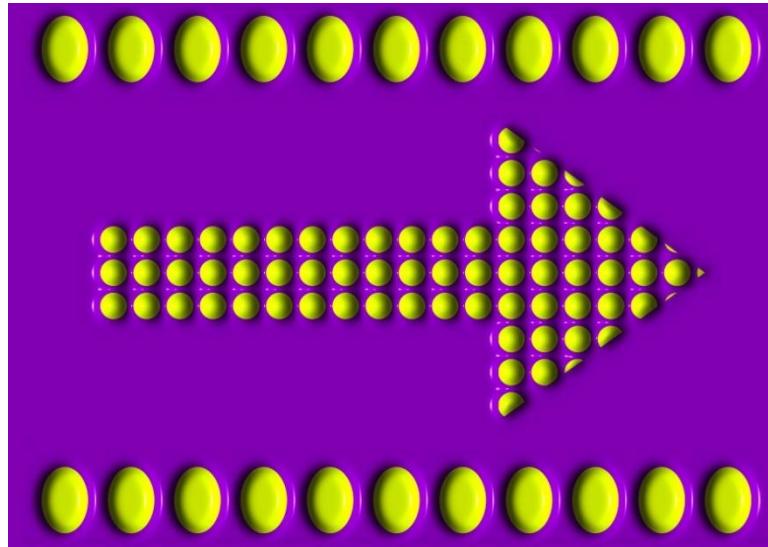
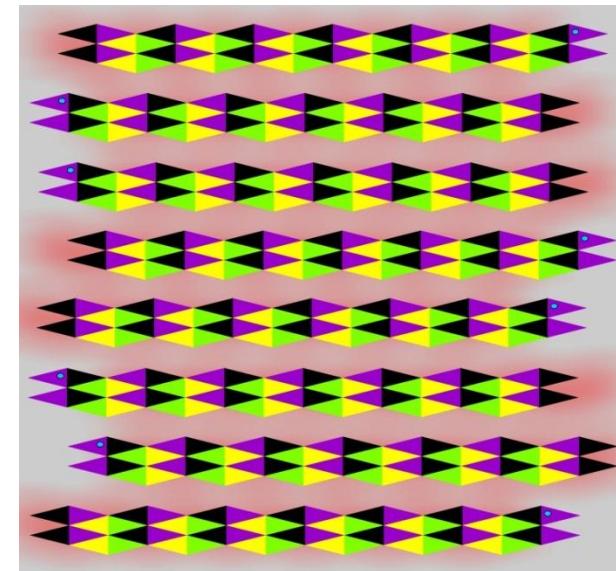
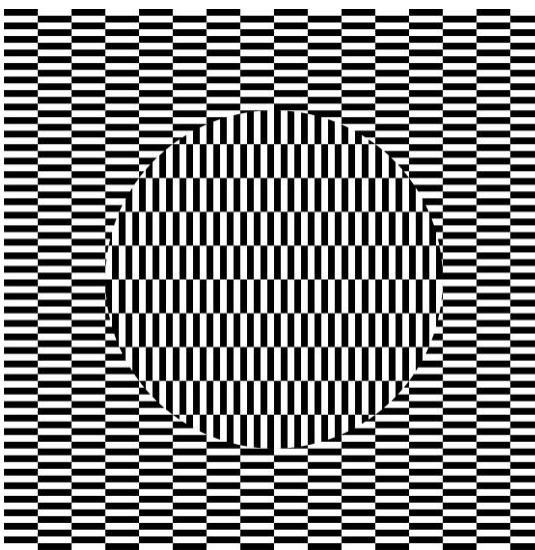


G. Johansson, “Visual Perception of Biological Motion and a Model For Its Analysis”, *Perception and Psychophysics* 14, 201-211, 1973.

Seeing motion from a static picture?

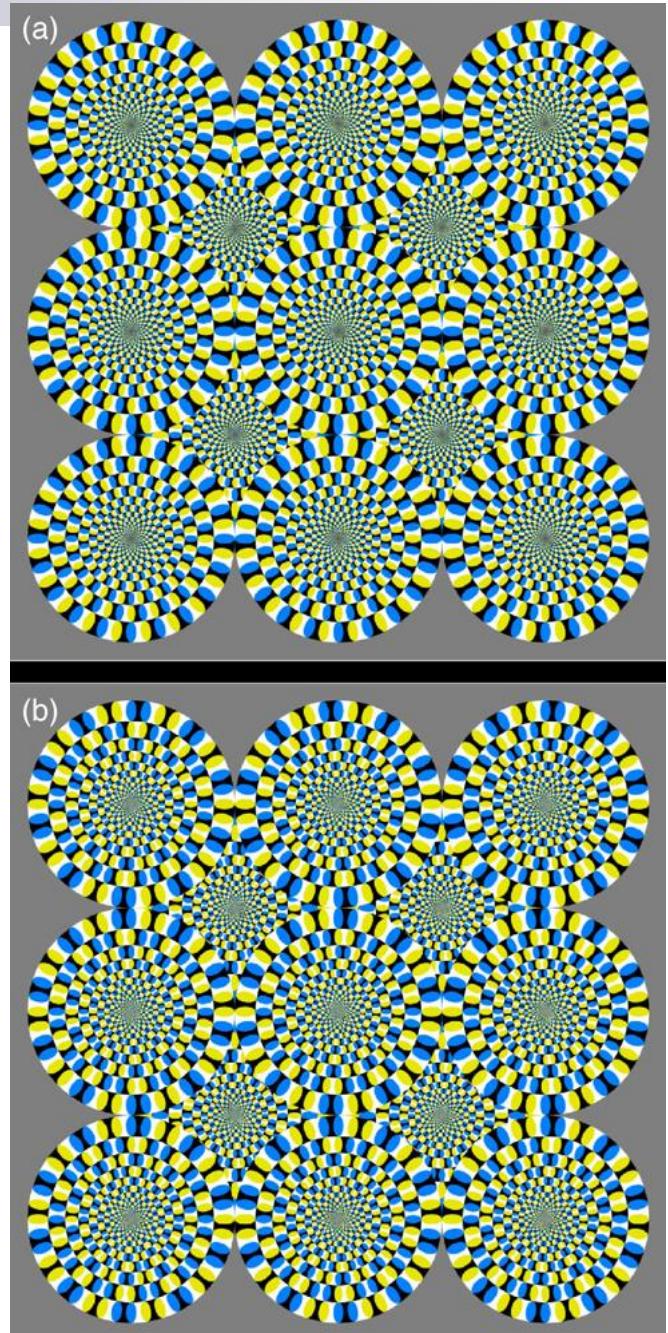


More examples



How is this possible?

- The true mechanism is to be revealed
- FMRI data suggest that illusion is related to some component of eye movements
- We don't expect computer vision to "see" motion from these stimuli, yet



The cause of motion

- Three factors in imaging process
 - Light
 - Object
 - Camera
- Varying either of them causes motion
 - Static camera, moving objects (surveillance)
 - Moving camera, static scene (3D capture)
 - Moving camera, moving scene (sports, movie)
 - Static camera, moving objects, moving light (time lapse)



Motion scenarios (priors)



Static camera, moving scene



Moving camera, static scene



Moving camera, moving scene

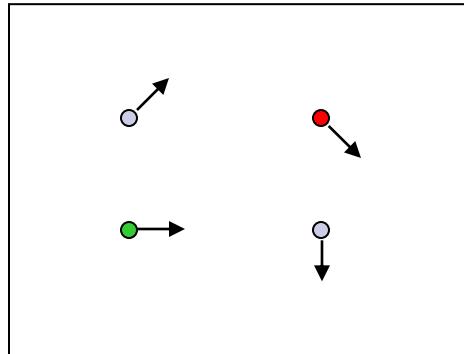


Static camera, moving scene, moving light

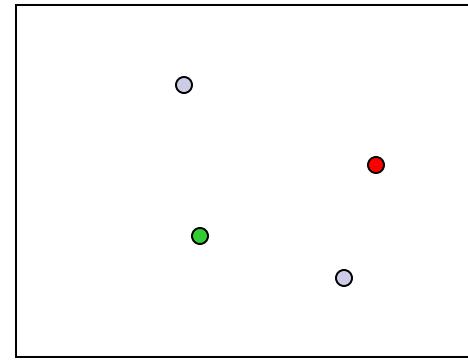
We still don't touch these areas



Recovering motion



$I(x,y,t)$

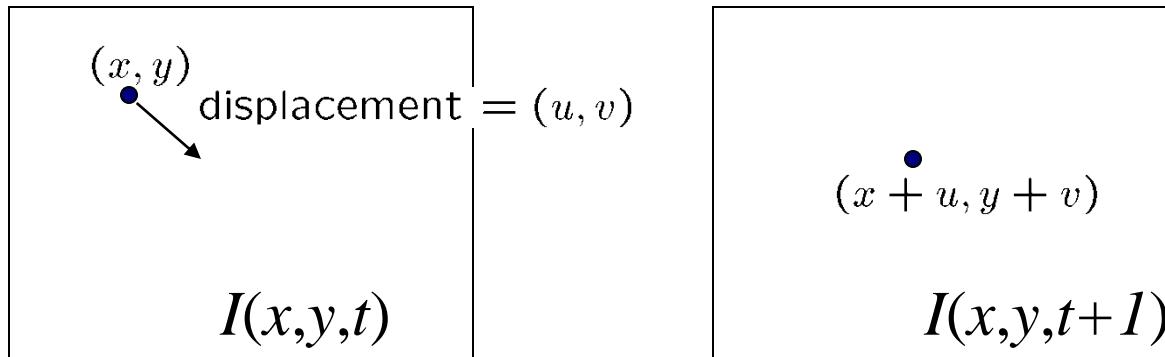


$I(x,y,t+1)$

- Given two subsequent frames, estimate the point translation
 - Key assumptions of Lucas-Kanade Tracker
 - **Brightness constancy:** projection of the same point looks the same in every frame
 - **Small motion:** points do not move very far
 - **Spatial coherence:** points move like their neighbors

B. Lucas and T. Kanade. [An iterative image registration technique with an application to stereo vision](#). In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

The brightness constancy constraint



■ Brightness Constancy Equation:

$$I(x, y, t) = I(x + u, y + v, t + 1)$$

Take Taylor expansion of $I(x+u, y+v, t+1)$ at (x, y, t) to linearize the right side:

$$I(x + u, y + v, t + 1) \approx I(x, y, t) + \boxed{I_x} \cdot u + I_y \cdot v + \boxed{I_t}$$

Image derivative along x Difference over frames

$$I(x + u, y + v, t + 1) - I(x, y, t) = I_x \times u + I_y \times v + I_t$$

So: $I_x \cdot u + I_y \cdot v + I_t \approx 0 \rightarrow \nabla I \cdot [u \ v]^T + I_t = 0$

The brightness constancy constraint

Can we use this equation to recover image motion (u, v) at each pixel?

$$\nabla I \cdot [u \ v]^T + I_t = 0$$

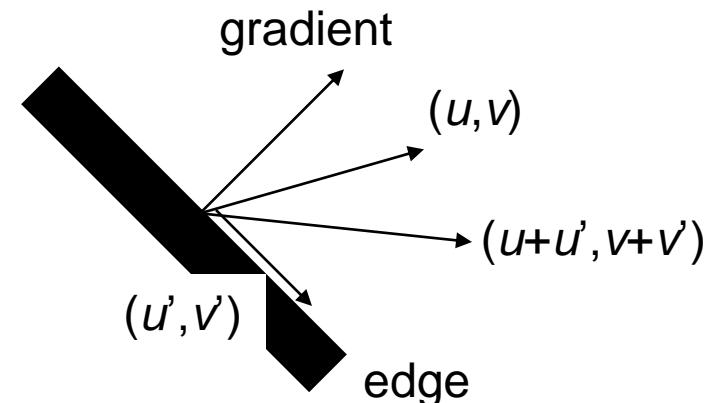
- How many equations and unknowns per pixel?

- One equation (this is a scalar equation!), two unknowns (u, v)

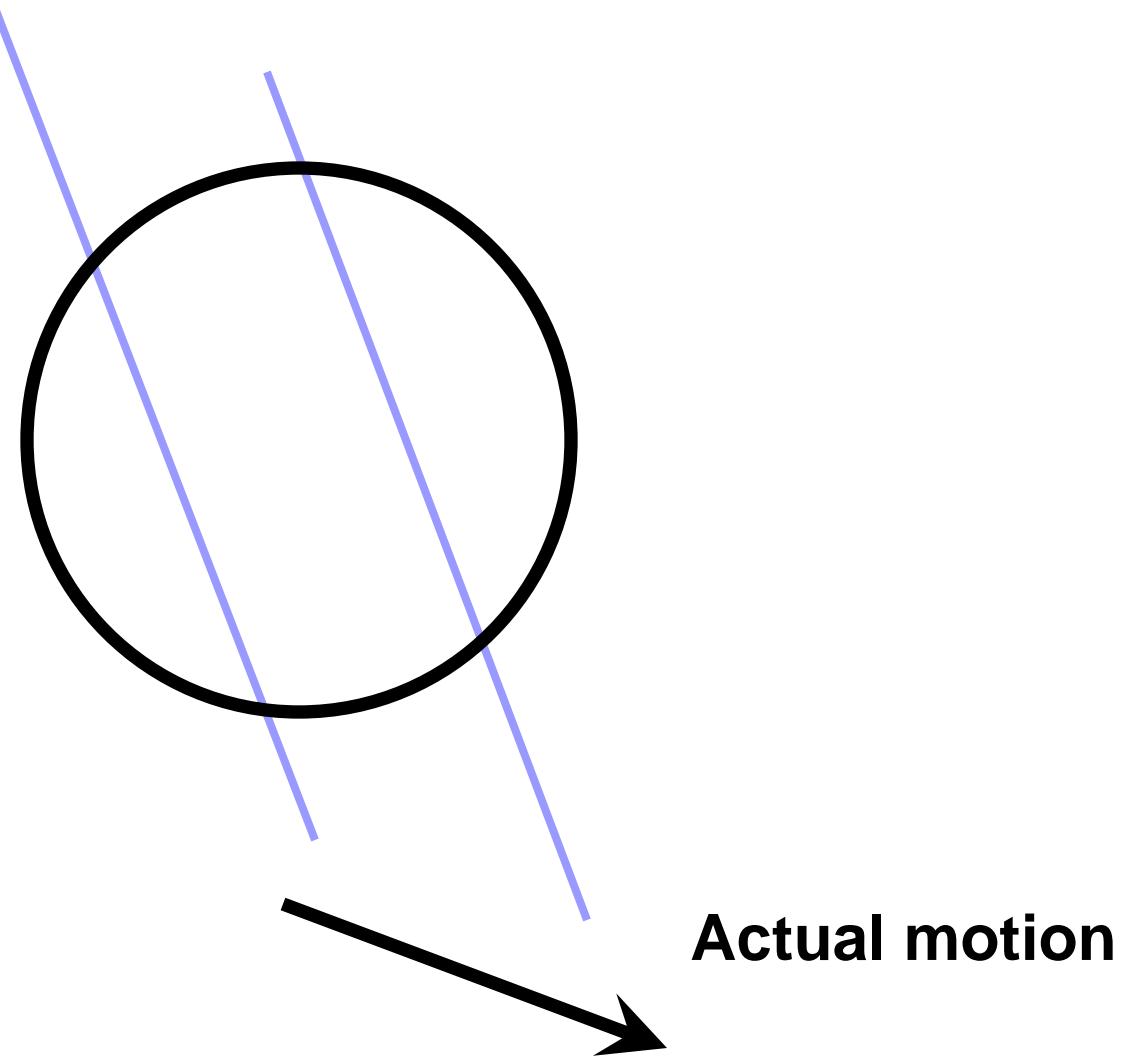
The component of the motion perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

If (u, v) satisfies the equation,
so does $(u+u', v+v')$ if

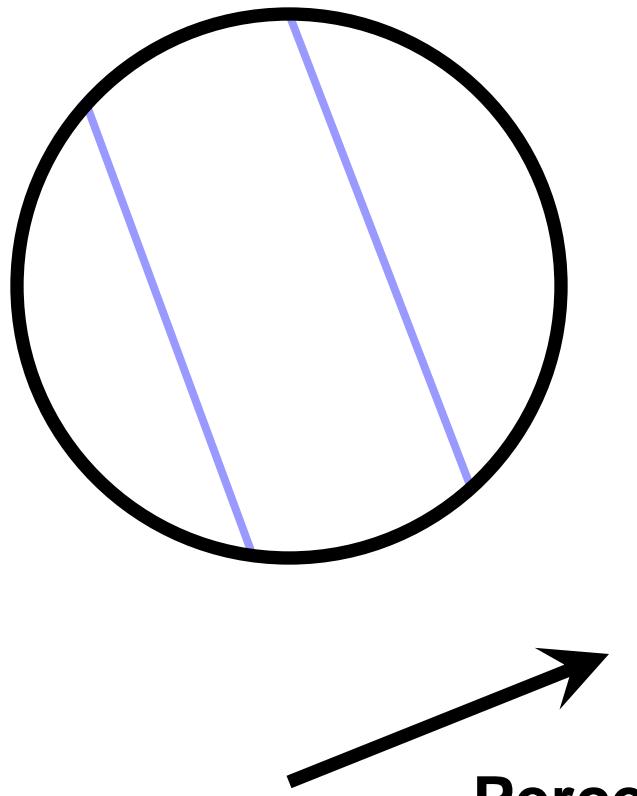
$$\nabla I \cdot [u' \ v']^T = 0$$



The aperture problem



The aperture problem



Solving the ambiguity...

B. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

- How to get more equations for a pixel?
- **Spatial coherence constraint**
- Assume the pixel's neighbors have the same (u,v)
 - If we use a 5x5 window, that gives us 25 equations per pixel

$$0 = I_t(\mathbf{p}_i) + \nabla I(\mathbf{p}_i) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix}$$

Solving the ambiguity...

- Least squares problem:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$A \quad d = b$
 $25 \times 2 \quad 2 \times 1 \quad 25 \times 1$



Local Smoothness

Matching patches across images

- Overconstrained linear system

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix} \quad \begin{matrix} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{matrix}$$

Least squares solution for d given by $(A^T A)^{-1} A^T b$

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$
$$A^T A \qquad \qquad \qquad A^T b$$

The summations are over all pixels in the $K \times K$ window

Conditions for solvability

Optimal (u, v) satisfies Lucas-Kanade

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$ $A^T b$

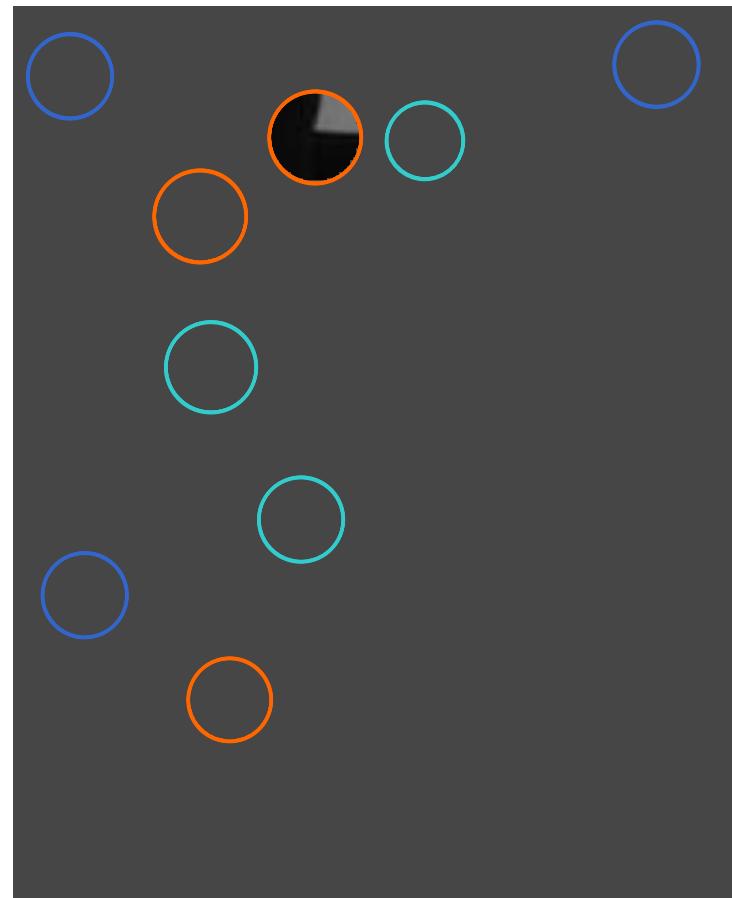
When is this solvable? I.e., what are good points to track?

- $\mathbf{A}^T \mathbf{A}$ should be invertible
- $\mathbf{A}^T \mathbf{A}$ should not be too small due to noise
 - eigenvalues λ_1 and λ_2 of $\mathbf{A}^T \mathbf{A}$ should not be too small
- $\mathbf{A}^T \mathbf{A}$ should be well-conditioned
 - λ_1 / λ_2 should not be too large (λ_1 = larger eigenvalue)

Does this remind you of anything?

Criteria for Harris corner detector

Aperture problem



Corners

Lines

Flat regions

Low-texture region



$$\sum \nabla I (\nabla I)^T$$

- gradients have small magnitude
- small λ_1 , small λ_2

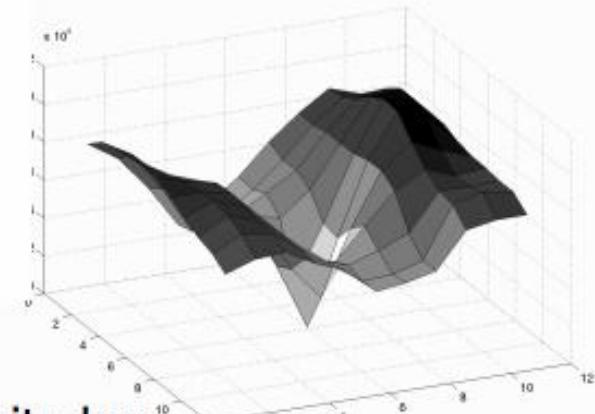
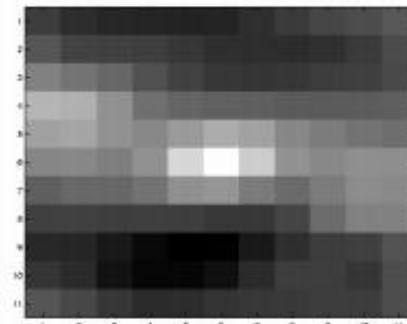
Edge



$$\sum \nabla I (\nabla I)^T$$

- gradients very large or very small
- large λ_1 , small λ_2

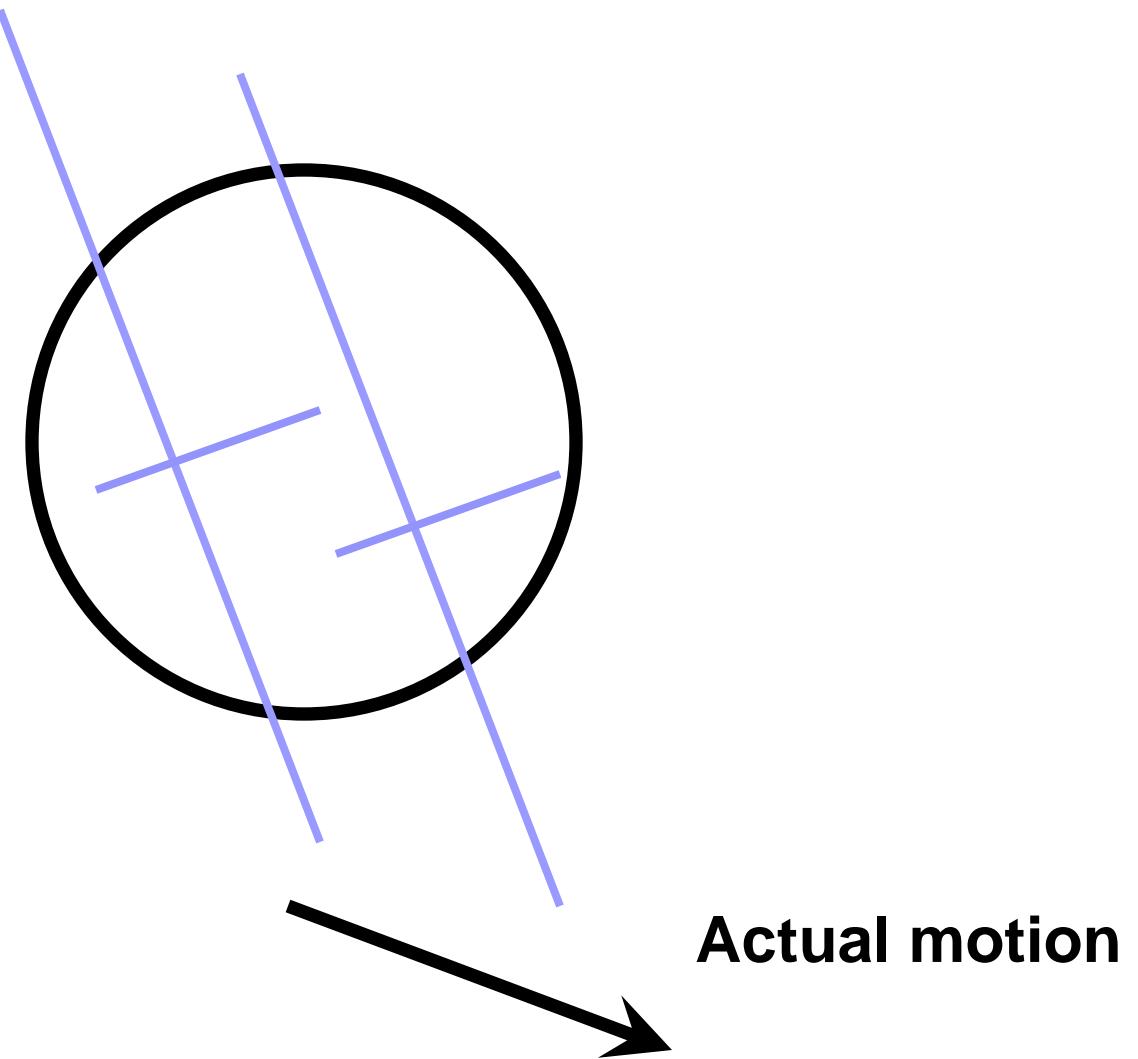
High Texture Region



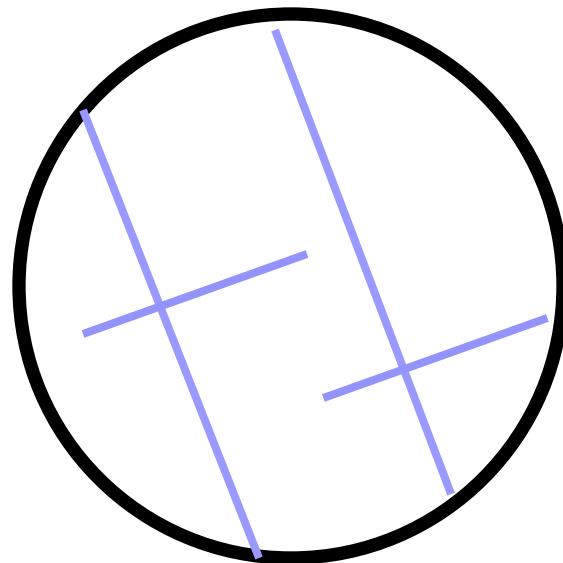
$$\sum \nabla I (\nabla I)^T$$

- gradients are different, large magnitudes
- large λ_1 , large λ_2

The aperture problem resolved



The aperture problem resolved



Perceived motion

Global Smoothness Method

Determining Optical Flow, Horn B K, Schunck B G. 1981

Lucas and Kanade assume the pixel's neighbors have the same motion——local smooth method, Horn and Schunck put forward a global smooth method, formed as:

$$E = E_{data} + \lambda E_{Smooth}$$

- Data term: $E_{Data} = \sum_x (I_x u + I_y v + I_t)^2$
- Smooth term: $E_{Smooth} = \sum_x (u_x^2 + u_y^2 + v_x^2 + v_y^2)$
- Solved with calculus of variations

Break-downs

- Brightness constancy is not perfectly satisfied.



Robust penalty function.

- A point does not move like its neighbors.

Motion discontinuity



anisotropic, bilateral weight

- The motion is not small (Taylor expansion doesn't hold).



Coarse-to-fine(multi-scale) estimation.

Robust penalty function

- 1. Quadratic: $\rho(x) = x^2$
- 2. Truncated quadratic: $\rho_{\alpha,\lambda}(x) = \begin{cases} \lambda x^2 & \text{if } x < \frac{\sqrt{\alpha}}{\sqrt{\lambda}} \\ \alpha & \text{otherwise} \end{cases}$
- 3. Lorentzian: $\rho_\sigma(x) = \log\left(1 + \frac{1}{2}\left(\frac{x}{\sigma}\right)^2\right)$ GNC
- 4. L1 Norm: $\rho(x) = \|x\|$ TV
- 5. Charbonnier: $\rho(x) = \sqrt{(x^2 + \varepsilon^2)}$
- 6. Generalized Charbonnier: $\rho(x) = (x^2 + \varepsilon^2)^\alpha$
- 7. M-estimator:
$$\rho(t;\sigma) = \frac{t^2}{t^2 + \sigma^2}$$

Calculus
of
variation

□ Huber: $\rho(t;\sigma) = \begin{cases} t^2 & |t| < \sigma \\ \sigma^2 & |t| \geq \sigma \end{cases}$ $\rho(t;\sigma) = \begin{cases} t^2, & |t| < \sigma \\ \sigma/|t|, & \sigma \leq |t| < 3\sigma \\ 0, & |t| \geq 3\sigma \end{cases}$

GNC: graduated non convex

TV: total variance

IRLS: iterated reweighted least square

IRLS

Motion discontinuity

- 1. Anisotropic Smoothness: define a structure adaptive map according to the brightness derivatives.

$$w_{Smooth} = \exp(-\|\nabla I\|^\kappa)$$

“structure and motion adaptive regularization for high accuracy optical flow”(ICCV2009)

- 2.Bilateral weight:

$$w_{Smooth}(x, y) = \exp \left\{ -\frac{\|x - y\|^2}{2\sigma_1^2} - \frac{\|I(x) - I(y)\|^2}{2\sigma_2^2} \right\}$$

“Efficient Nonlocal Regularization for optical flow”(ECCV2012)

Dealing with larger movements: Iterative refinement

1. Initialize $(x', y') = (x, y)$

2. Compute (u, v) by

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

2nd moment matrix for feature
patch in first image

Original (x,y) position

$$I_t = I(x', y', t+1) - I(x, y, t)$$

displacement

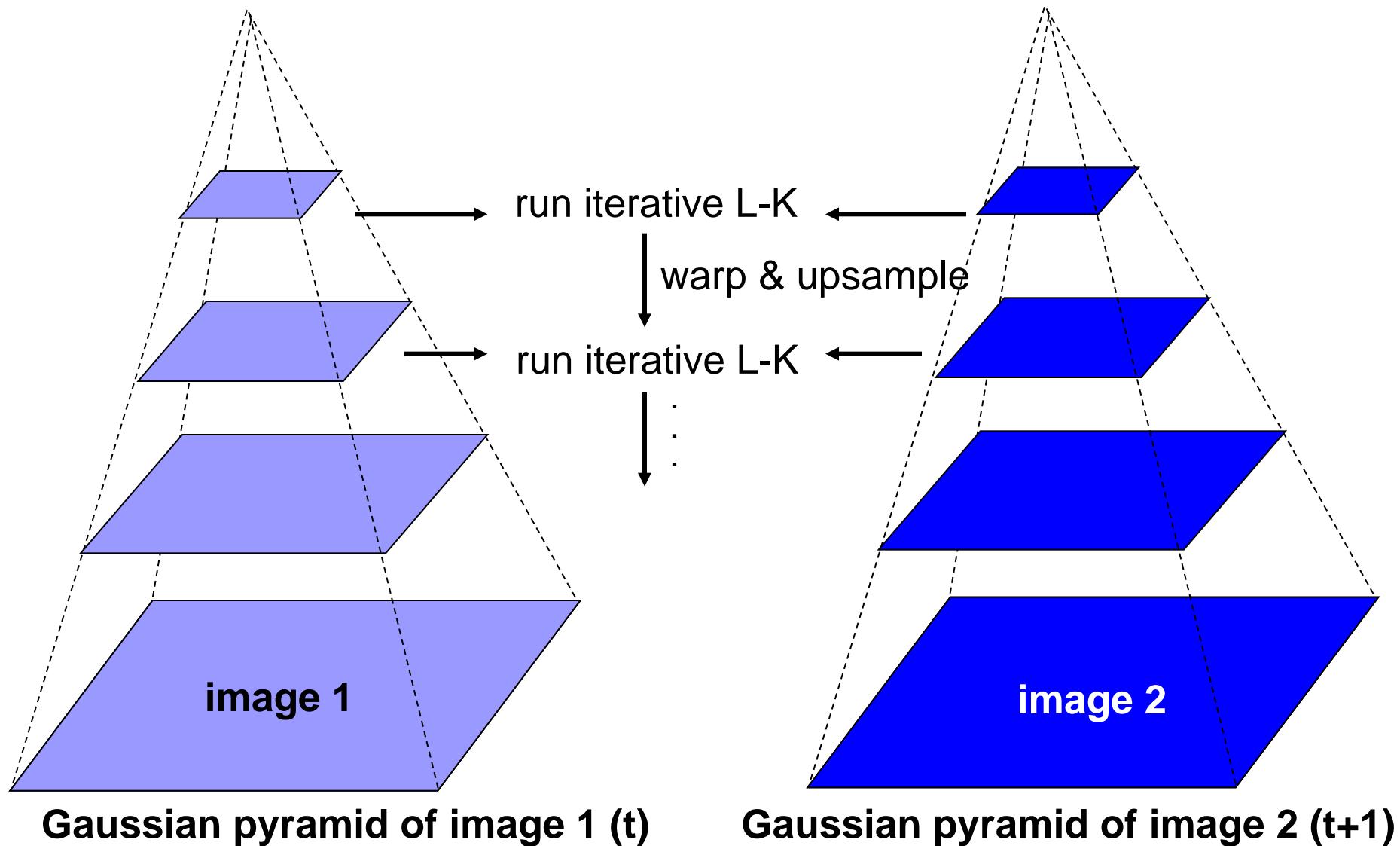
3. Shift window by (u, v) : $x' = x' + u$; $y' = y' + v$;

4. Recalculate I_t

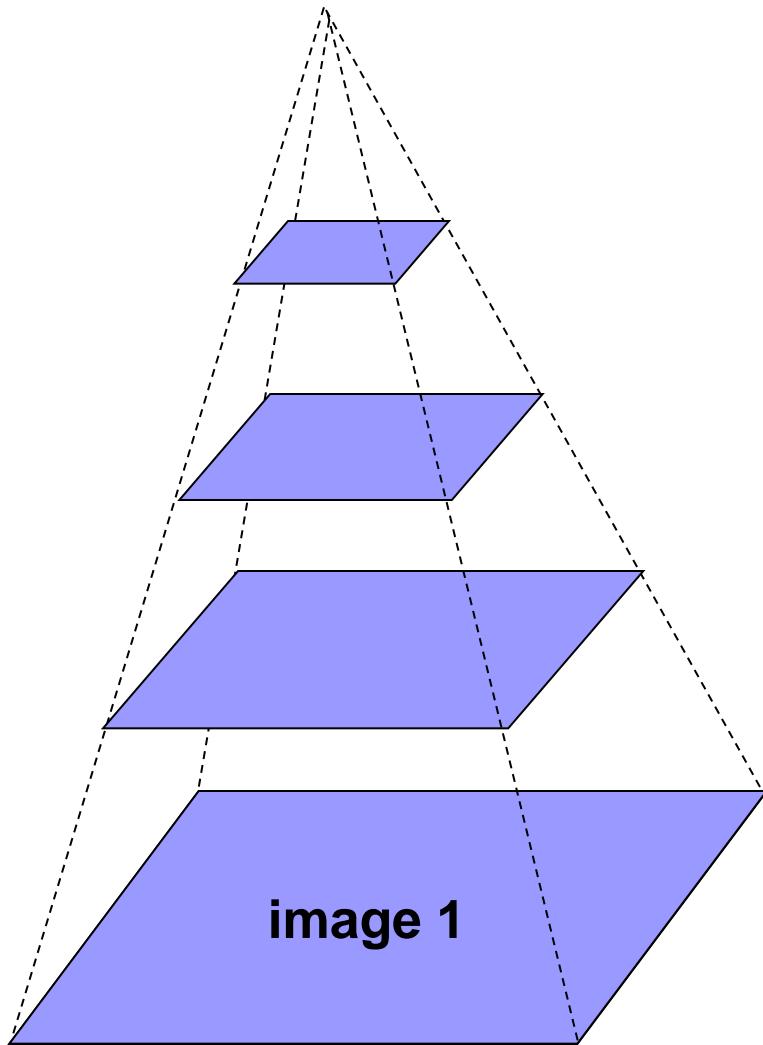
5. Repeat steps 2-4 until small change

- Use interpolation for subpixel values

Coarse-to-fine optical flow estimation



Coarse-to-fine optical flow estimation



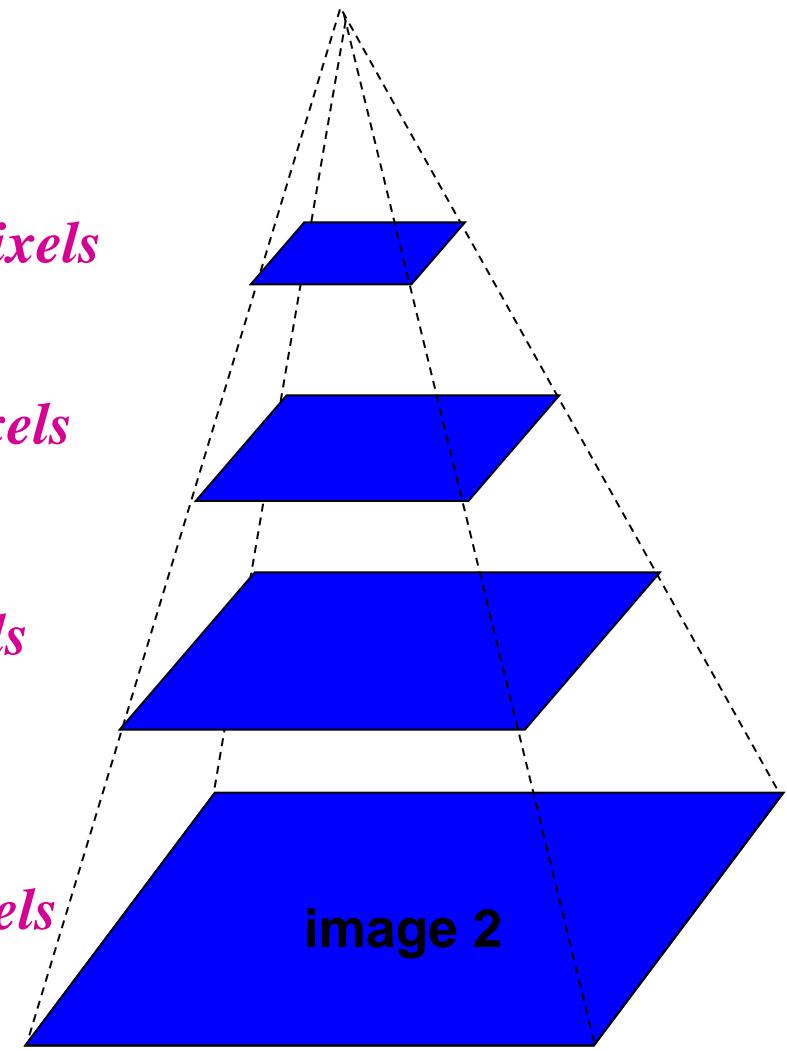
Gaussian pyramid of image 1

$u=1.25 \text{ pixels}$

$u=2.5 \text{ pixels}$

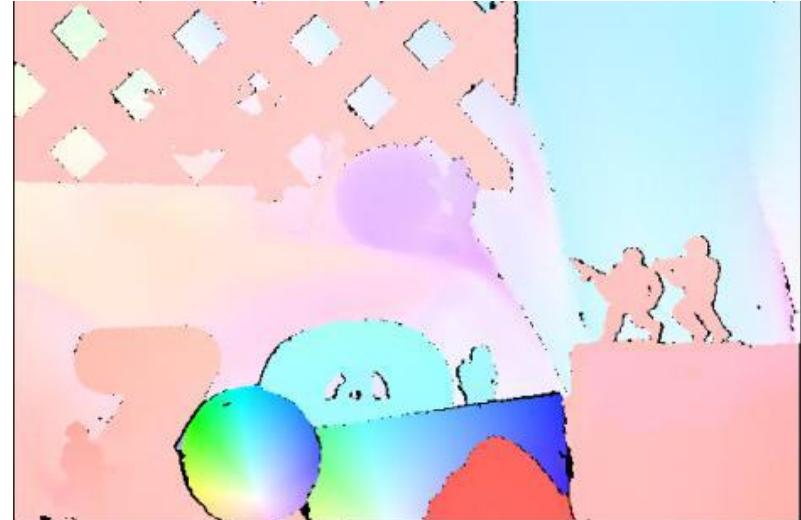
$u=5 \text{ pixels}$

$u=10 \text{ pixels}$



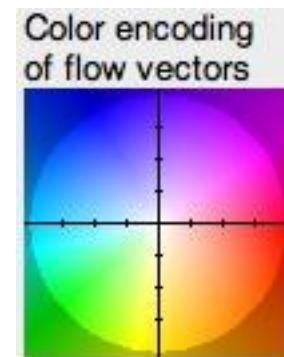
Gaussian pyramid of image 2

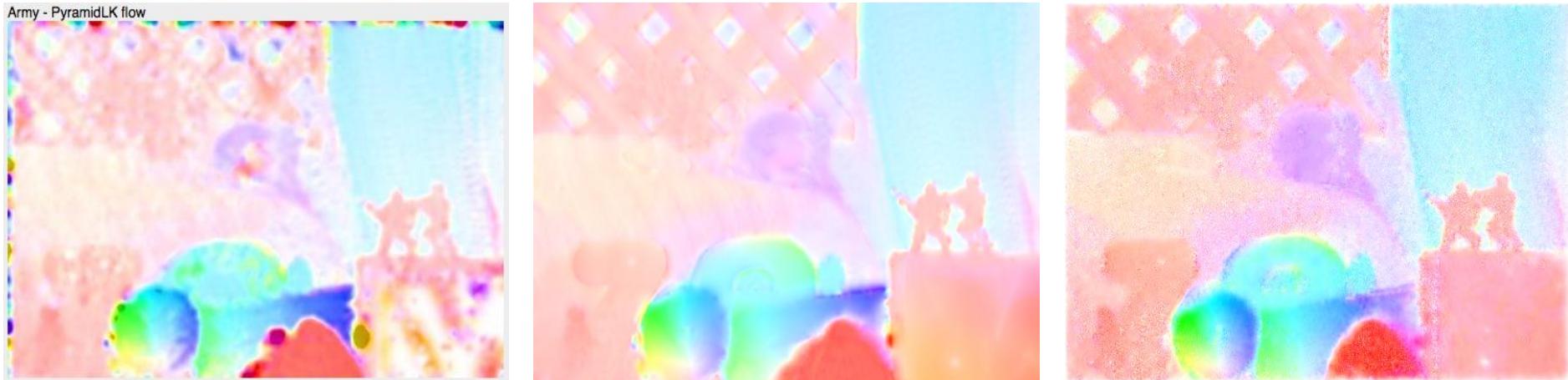
Flow quality evaluation



Ground Truth

- Middlebury flow page
 - <http://vision.middlebury.edu/flow/>





Lucas-Kanade

HS

GNC



Brox

TV-L1 Improved

MDP

Video stabilization



Video denoising

Original



Denoised



Video super resolution

Low-Res



Motion Detail Preserving Optical Flow Estimation

- Framework
 - Extended coarse-to-fine motion estimation for both large and small displacement optical flow
- Model
 - A new data term to selectively combine constraints
- Solver
 - Efficient numerical solver for discrete-continuous optimization

Extended Flow Initialization

- Sparse feature matching for each level



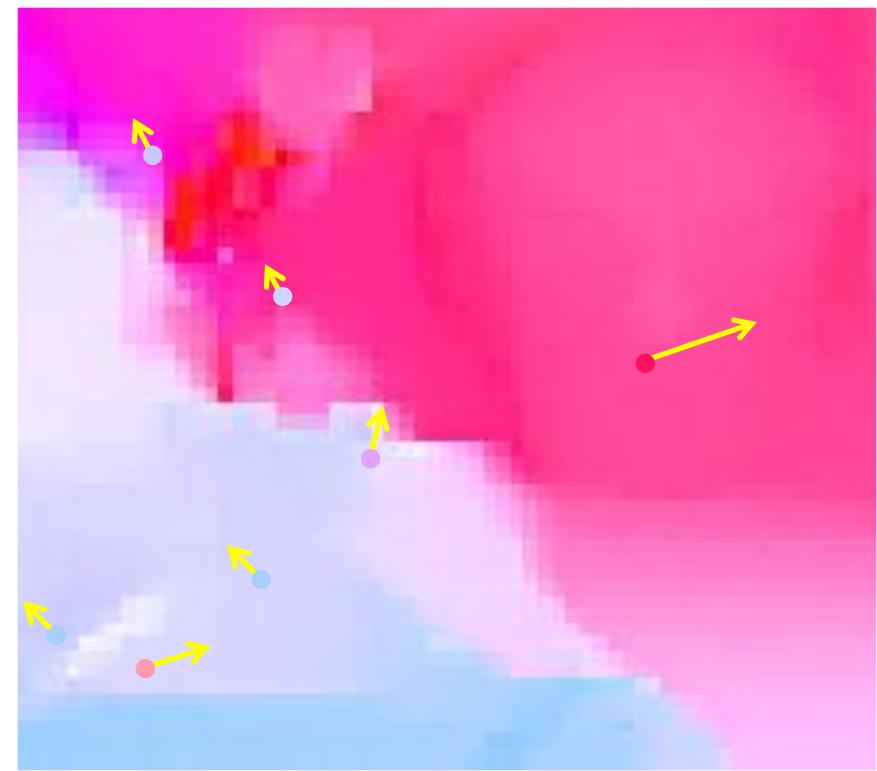
Extended Flow Initialization

- Identify missing motion vectors

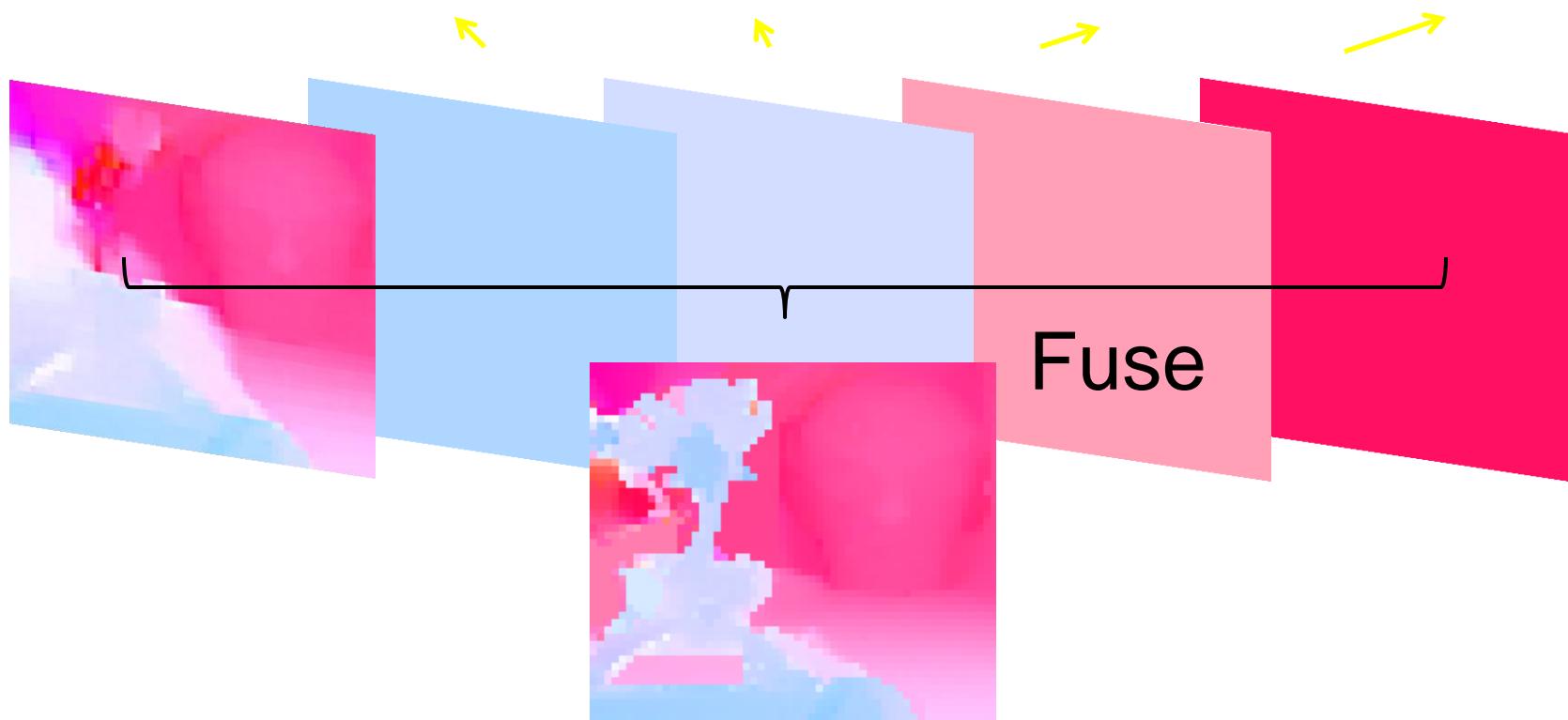


Extended Flow Initialization

- Identify missing motion vectors



Extended Flow Initialization



Results from Different Steps



Coarse-to-fine



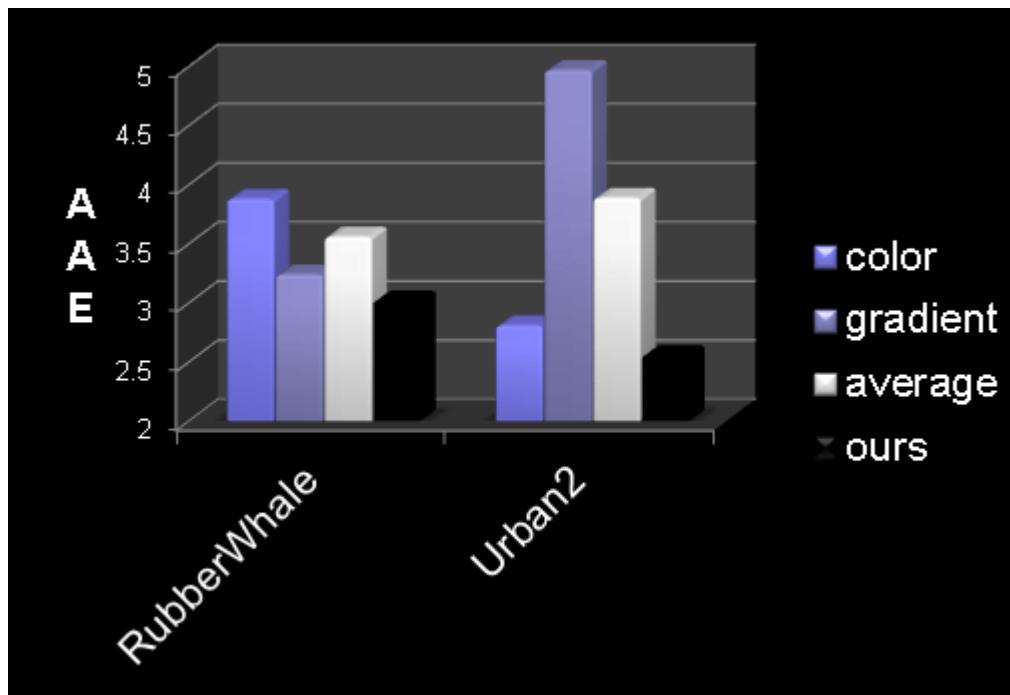
Extended coarse-to-fine

■ Selectively combine the constraints

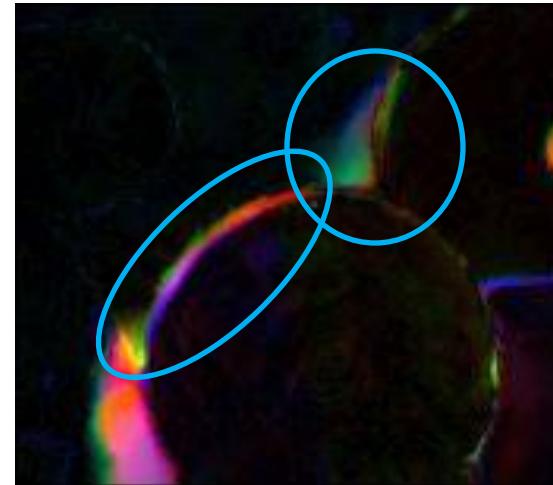
$$E_D(u, \alpha) = \sum_x \alpha(x) D_I(u, x) + (1 - \alpha(x)) D_{VI}(u, x)$$

where $\alpha(x) : Z^2 \mapsto \{0, 1\}$

Comparisons



Results



Difference



Averaging constraints



Ours

More Results



Overlaid Input



Coarse-to-fine



Our Results

Challenge Problems

LARGE DISPLACEMENT OCCLUSION

always appear
simultaneously

Average endpoint error	avg.	Army (Hidden texture)				Mequon (Hidden texture)				Schefflera (Hidden texture)				Wooden (Hidden texture)				Grove (Synthetic)				Urban (Synthetic)				Yosemite (Synthetic)				Teddy (Stereo)			
		GT		im0	im1	GT		im0	im1	GT		im0	im1	GT		im0	im1	GT		im0	im1	GT		im0	im1	GT		im0	im1	GT		im0	im1
		rank	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	all	disc	untext	
IVANN [87]	2.7	0.071	0.202	0.051	0.151	0.513	0.125	0.181	0.371	0.141	0.102	0.493	0.062	0.411	0.611	0.212	0.232	0.662	0.191	0.104	0.129	0.1712	0.341	0.804	0.232	0.111	0.129	0.2130	0.427	0.782	0.6312		
OFLAF [77]	6.9	0.087	0.213	0.065	0.165	0.534	0.125	0.192	0.371	0.141	0.147	0.7723	0.074	0.514	0.785	0.253	0.315	0.763	0.257	0.1111	0.129	0.2130	0.427	0.782	0.6312	0.111	0.129	0.2130	0.427	0.782	0.6312		
MDP-Flow2 [68]	7.9	0.087	0.213	0.0714	0.151	0.481	0.111	0.204	0.404	0.141	0.1518	0.8029	0.0810	0.6314	0.9314	0.4315	0.263	0.763	0.236	0.1111	0.129	0.1712	0.383	0.793	0.444	0.111	0.129	0.2130	0.427	0.782	0.6312		
NN-field [71]	8.7	0.087	0.2214	0.051	0.177	0.556	0.1310	0.192	0.393	0.156	0.091	0.482	0.051	0.411	0.611	0.201	0.5246	0.641	0.2610	0.1331	0.1326	0.2023	0.352	0.835	0.211	0.111	0.129	0.2130	0.427	0.782	0.6312		
ComponentFusion [96]	10.0	0.071	0.213	0.051	0.165	0.556	0.125	0.204	0.447	0.156	0.113	0.656	0.062	0.7127	1.0732	0.5329	0.327	0.6620	0.2813	0.1111	0.1326	0.157	0.416	0.889	0.545	0.111	0.129	0.2130	0.427	0.782	0.6312		
WLIF-Flow [93]	14.7	0.087	0.213	0.065	0.189	0.556	0.1519	0.2513	0.5614	0.1712	0.147	0.687	0.0810	0.6112	0.9113	0.4113	0.4323	0.9611	0.2918	0.1331	0.129	0.2130	0.5127	1.0328	0.7226	0.111	0.129	0.2130	0.427	0.782	0.6312		
TC/T-Flow [76]	14.9	0.071	0.213	0.051	0.1914	0.6826	0.125	0.2818	0.6623	0.141	0.147	0.8637	0.074	0.6724	0.9823	0.4924	0.221	0.825	0.191	0.1111	0.111	0.3067	0.5022	1.0225	0.6414	0.111	0.129	0.2130	0.427	0.782	0.6312		
Layers++ [37]	16.5	0.087	0.213	0.0714	0.1914	0.569	0.1726	0.204	0.404	0.1818	0.136	0.584	0.074	0.483	0.703	0.336	0.4734	0.0114	0.3335	0.1552	0.1447	0.2442	0.4612	0.889	0.7226	0.111	0.129	0.2130	0.427	0.782	0.6312		
LME [70]	17.0	0.087	0.2214	0.065	0.151	0.492	0.111	0.3026	0.6418	0.3169	0.1518	0.7826	0.0923	0.6620	0.9619	0.5329	0.338	0.1832	0.2813	0.1222	0.129	0.1816	0.448	0.9111	0.6110	0.111	0.129	0.2130	0.427	0.782	0.6312		
IROF++ [58]	17.4	0.087	0.2319	0.0714	0.2126	0.6826	0.1726	0.2818	0.6317	0.1930	0.1518	0.7319	0.0923	0.6010	0.8910	0.4214	0.4323	0.0823	0.3125	0.104	0.129	0.124	0.4714	0.9818	0.6821	0.111	0.129	0.2130	0.427	0.782	0.6312		
nLayers [57]	17.6	0.071	0.191	0.065	0.2232	0.5912	0.1945	0.2513	0.5411	0.2038	0.1518	0.8434	0.0810	0.535	0.785	0.348	0.4427	0.846	0.3022	0.1331	0.1326	0.2023	0.4714	0.9717	0.6719	0.111	0.129	0.2130	0.427	0.782	0.6312		
FC-2Layers-FF [74]	19.6	0.087	0.213	0.0714	0.2126	0.7030	0.1726	0.204	0.404	0.1818	0.1518	0.7622	0.0810	0.535	0.774	0.379	0.4940	0.0215	0.3335	0.1662	0.1326	0.2963	0.448	0.878	0.6414	0.111	0.129	0.2130	0.427	0.782	0.6312		



Recent Advances

- For Occlusion: explicitly detection and handling
 - **Cross-check:** "Fast cost volume filtering for visual correspondence and beyond"(CVPR2011)
 - **Layered model** (depth order):" Layered image motion with explicit occlusions, temporal consistency, and depth ordering"(NIPS2010)
 - **Exploit motion cues:** "Occlusion boundary detection and figure-ground assignment from optical flow"(CVPR2011)
 - **Occ layer:** "Joint Motion Estimation and Segmentation of Complex Scenes with Label Costs and Occlusion Modeling"(CVPR2012)
 - **Binary classification:** "Locally affine sparse-to-dense matching for motion and occlusion estimation"(ICCV2013)
 - **Occ relationship inference:** "Local Layering for Joint Motion Estimation and Occlusion Detection"(CVPR2014)

Recent Advances

■ For Large Displacement: add correspondence constraint.

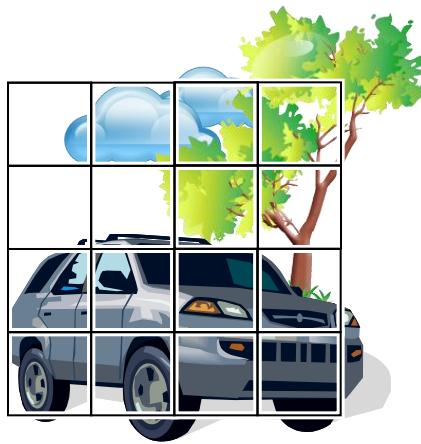
- “Large displacement Pattern Recognition from **nearest neighbor fields**” ”(CVPR2013) **CSH**
 - “Large displacement optical flow computation without warping”(ICCV2009)
 - “Large displacement optical flow(CVPR2009)
 - “Large displacement optical flow: descriptor matching in variational motion estimation”(PAMI2011)
 - “Motion detail preserving optical flow estimation”(CVPR2010,PAMI2012):**Fusion multi-candidate flow field**
- “Deepflow: Large displacement optical flow with **deep matching**”(ICCV2013)
 - “Fast Edge-preserving **PatchMatch** for Large Displacement Optical Flow”(CVPR2014): random sample, bilateral weight.

H
O
G

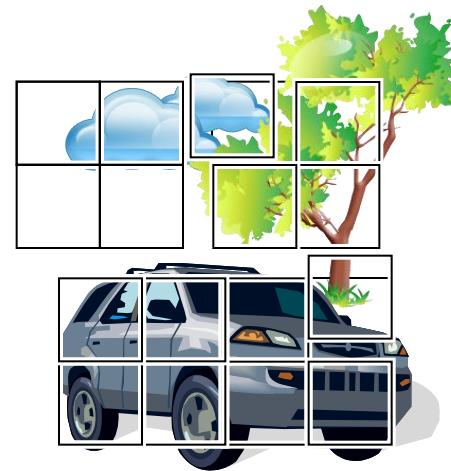
DeepFlow

- Classical optical flow [Horn and Schunck 1981]
 - ▶ energy $E(\mathbf{w}) = \iint E_{data} + \alpha E_{smooth} \mathbf{dx}$
- Integration of Deep Matching
 - ▶ energy $E(\mathbf{w}) = \iint E_{data} + \alpha E_{smooth} + \beta E_{match} \mathbf{dx}$
 - ▶ matches guide the flow
 - ▶ similar to [Brox and Malik 2011]
- Minimization using:
 - ▶ coarse-to-fine strategy
 - ▶ fixed point iterations
 - ▶ Successive Over Relaxation (SOR)

Deep Matching: main idea



First image

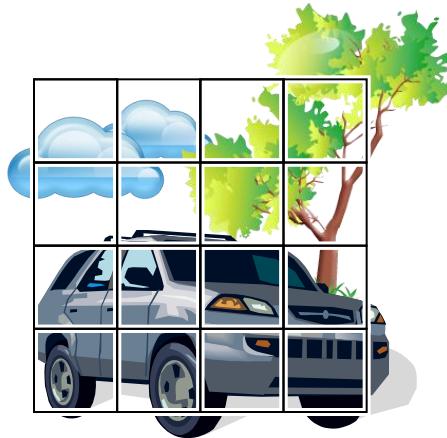


Second image

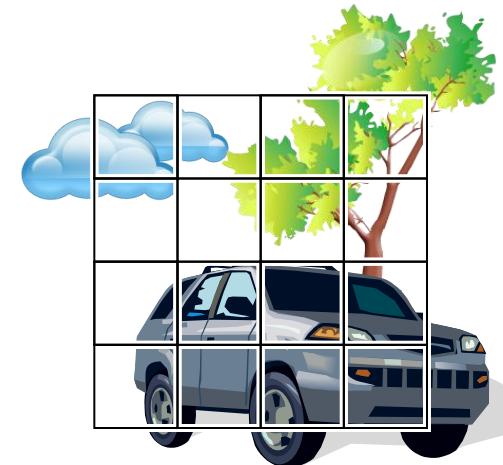
- We allow each subpatch to move:
 - ▶ independently
 - ▶ in a limited range depending on its size
- Our approach is fast to compute using convolution and max-pooling
- We apply this idea recursively

Related Work: matching

- HOG [Dalal and Triggs 2005] + Nearest Neighbor
 - ▶ rigid patches



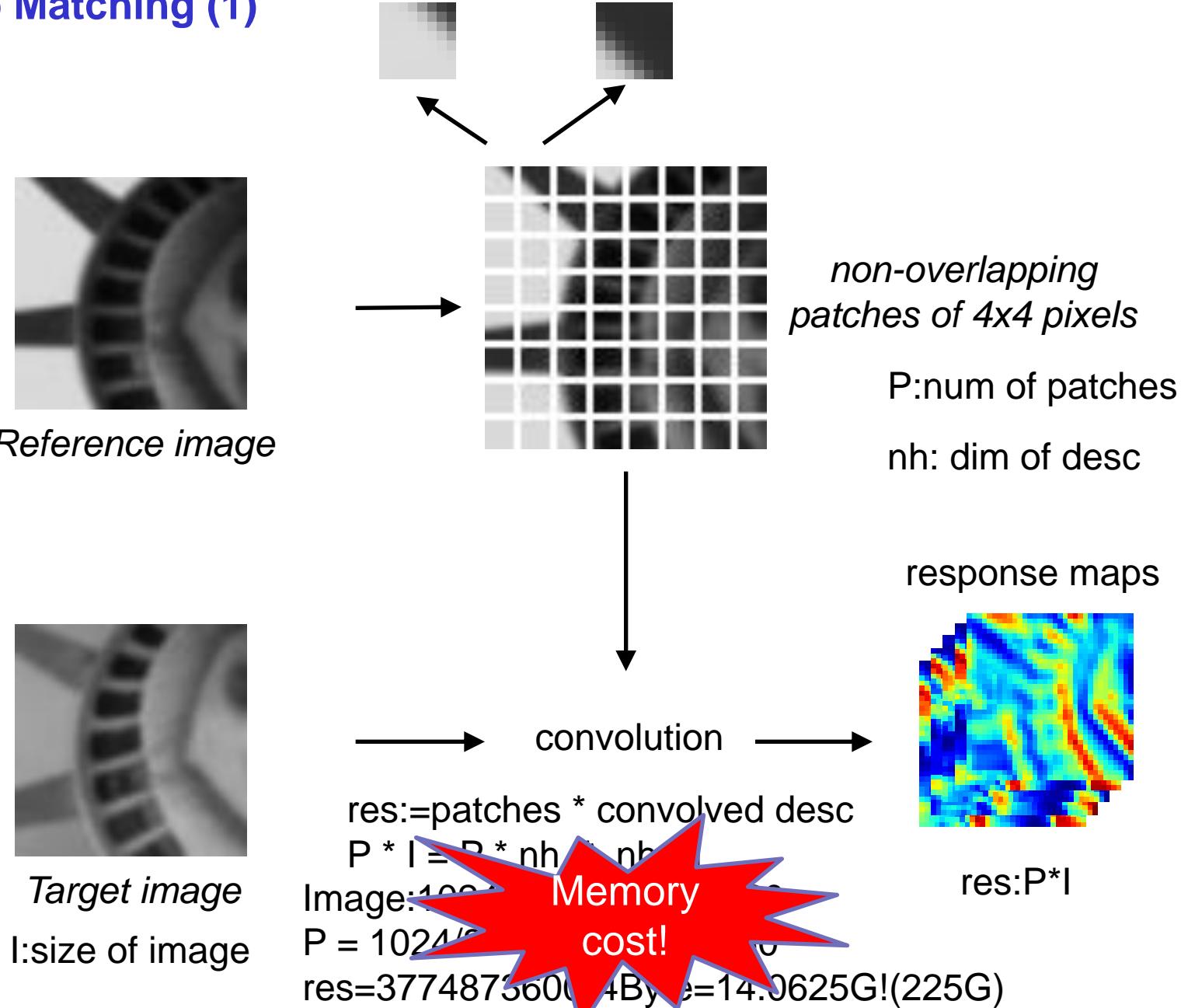
Input patch



Rigid matching

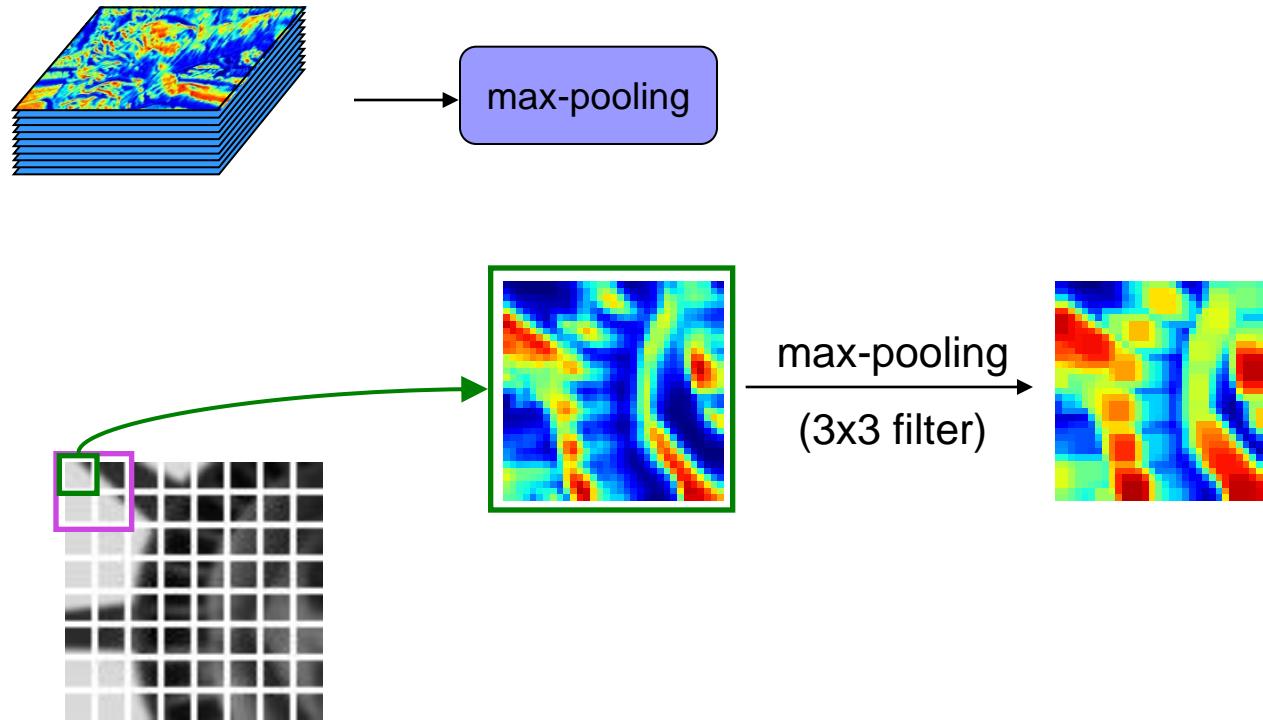
- Small patches and propagation: PatchMatch [Barnes *et al.* 2010]
 - ▶ patches not sufficiently discriminative

Deep Matching (1)



Deep Matching (2)

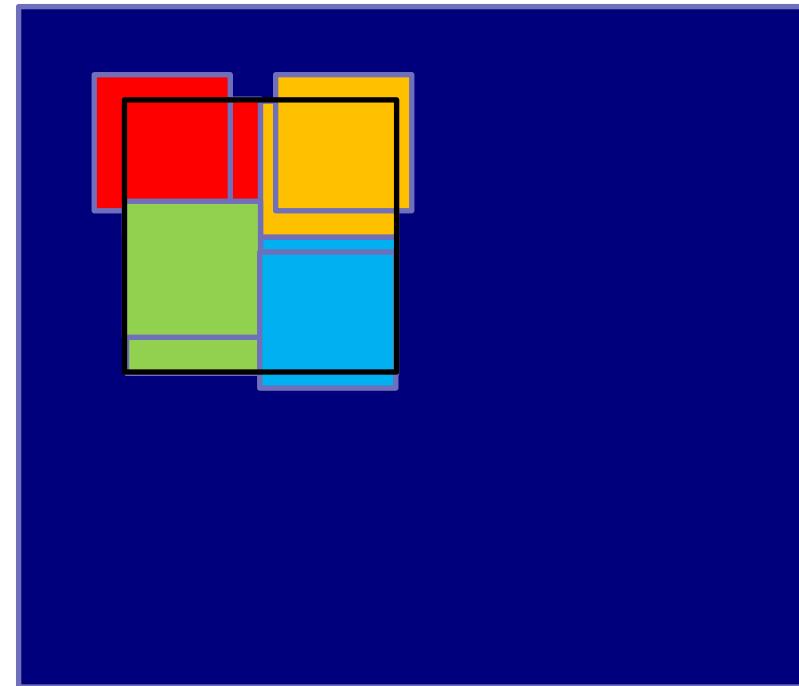
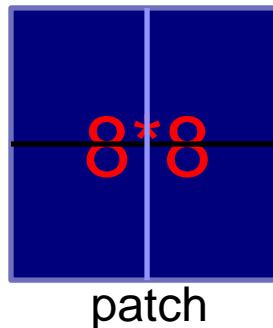
response maps for
each 4x4 patch



Max pooling

$$S(N) = \max_{t \in [-\frac{N}{8}, \frac{N}{8}]} S_{\text{left}}\left(\frac{N}{2}, t\right) + \max_{t \in [-\frac{N}{8}, \frac{N}{8}]} S_{\text{right}}\left(\frac{N}{2}, t\right). \quad (\text{1D case})$$

4*4



For each 4*4 patch, small offset: $[-1, 1] * [-1, 1]$, i.e:

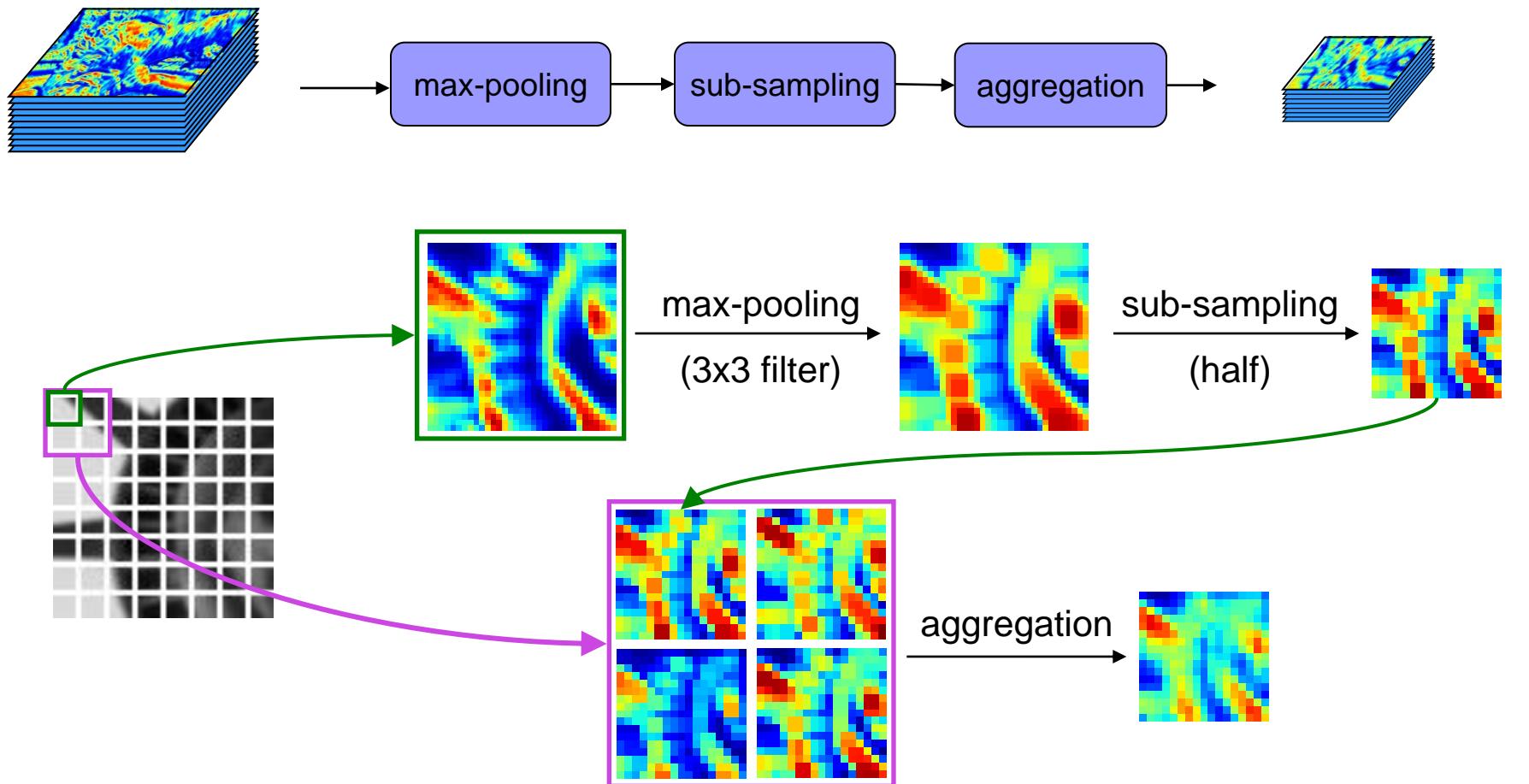
$$\begin{Bmatrix} (-1,-1) & (-1,0) & (-1,1) \\ (0,-1) & (0,0) & (0,1) \\ (1,-1) & (1,0) & (1,1) \end{Bmatrix}$$

$\rightarrow \begin{Bmatrix} -1 \\ 0 \\ 1 \end{Bmatrix} \times \{-1, 0, 1\}$ e.g :

$\begin{Bmatrix} 1 & 21 & 15 \\ 27 & 40 & 50 \\ 54 & 80 & 36 \end{Bmatrix} \rightarrow \begin{Bmatrix} 27 & 40 & 50 \\ 54 & 80 & 50 \\ 54 & 80 & 50 \end{Bmatrix} \rightarrow \begin{Bmatrix} 40 & 50 & 50 \\ 80 & 80 & 80 \\ 80 & 80 & 80 \end{Bmatrix}$

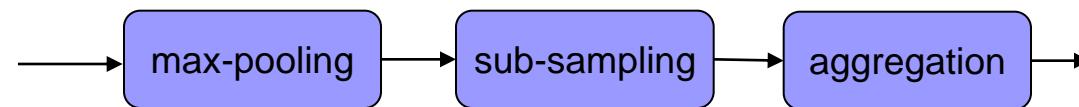
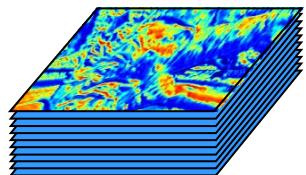
Deep Matching (2)

response maps for
each 4x4 patch

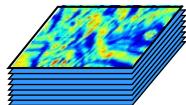


Deep Matching (2)

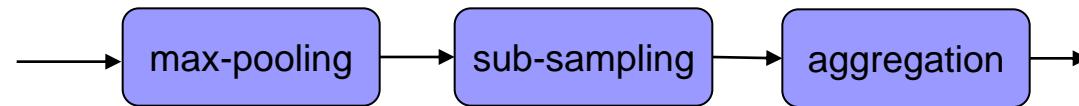
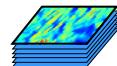
response maps for
each 4x4 patch



response maps
of 8x8 patches



response maps
of 16x16 patches

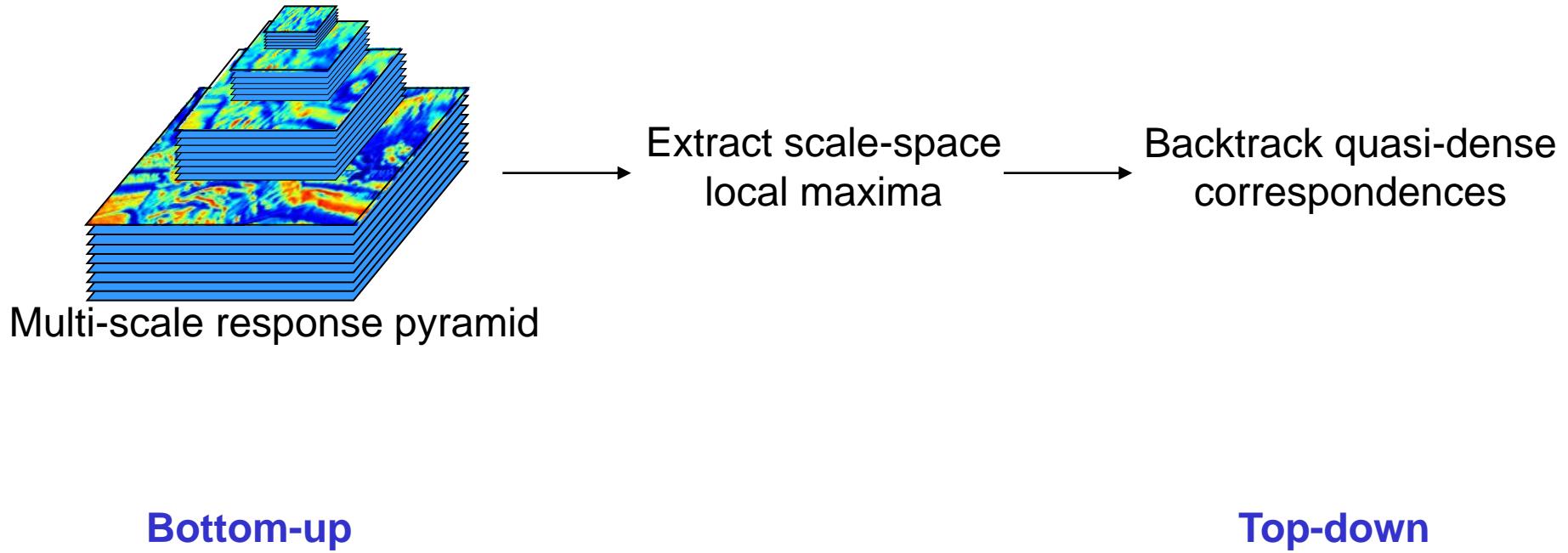


response maps
of 32x32 patches

...

Pipeline similar in spirit to **deep** convolutional nets [Lecun *et al.* 1998]

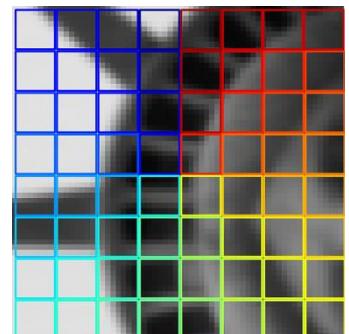
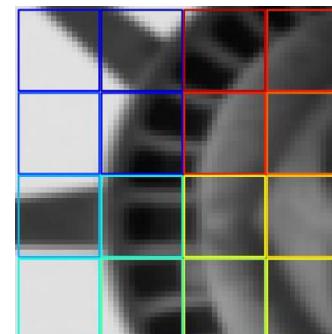
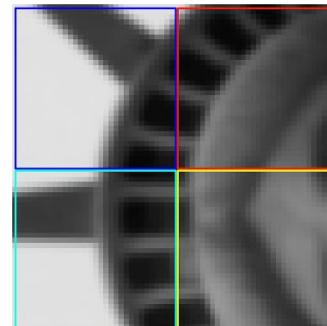
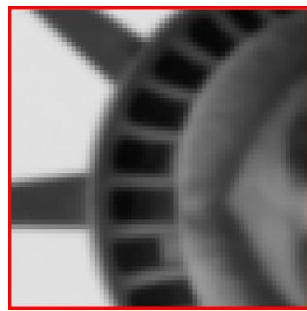
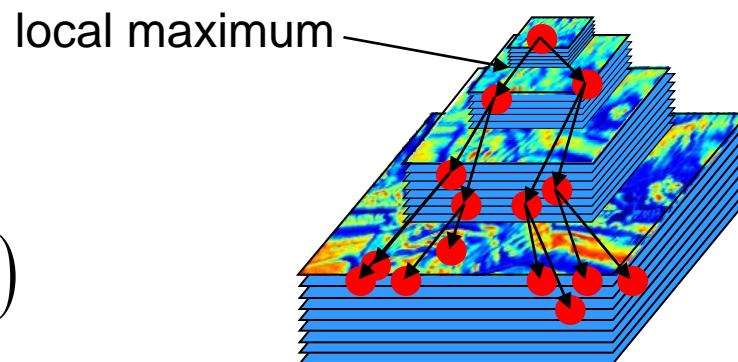
Deep Matching (3)



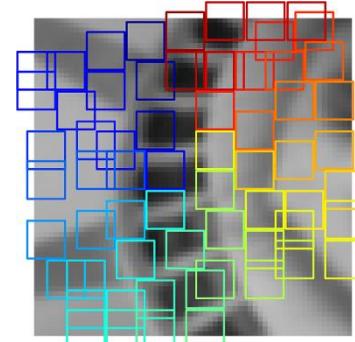
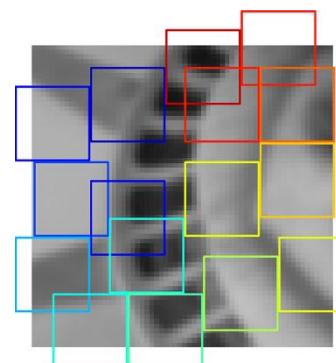
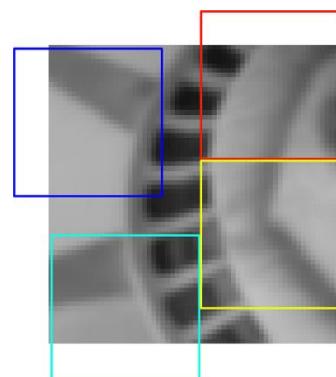
Deep Matching (3)

Merge the correspondences extracted from all local maxima

$$weight = sim(P_4, P'_4) \cdot l \cdot S(w^*)$$



First
image



Second
image

Deep Matching: example results

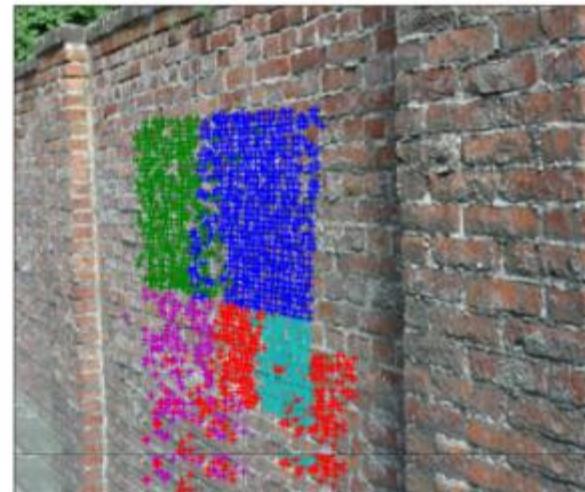
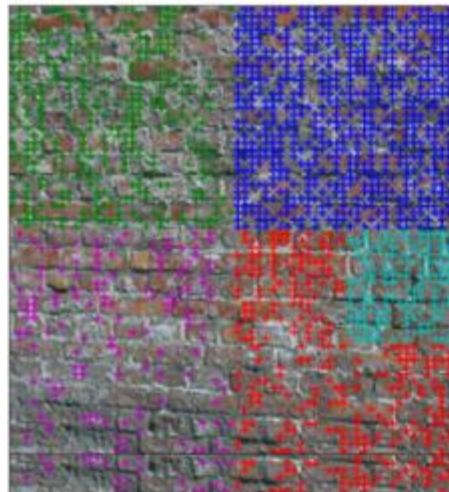
- Repetitive textures



First image



Second image



Deep Matching: example results

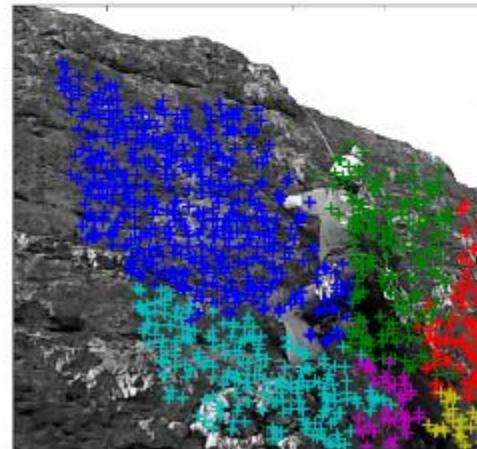
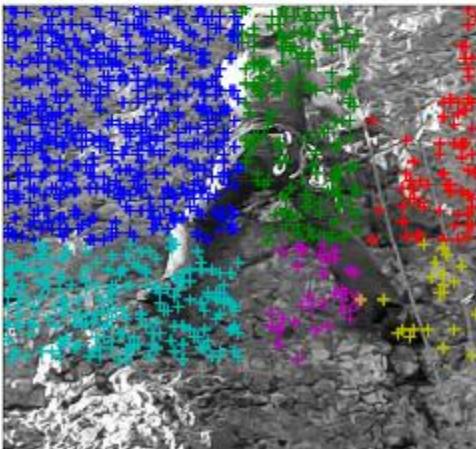
- Non-rigid deformation



First image



Second image



特征匹配与光流总结

■ 稀疏对应 - 特征点匹配

- Invariance: Translation, Orientation
(orientation normalization), Scale **(difference of Gaussian)**, Lighting condition **(gradient)**, Object deformations, Partial occlusion
- Interest point detection, description, descriptor matching
- Corner (spatial non-maximum suppression), SIFT (scale space), SURF, HOG,

特征匹配与光流总结

■ 稠密对应—光流估计

- Brightness constancy

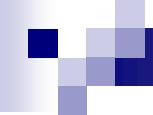
- Robust penalty function

- Small motion

- Coarse-to-fine (multi-scale) estimation (Taylor expansion doesn't hold)

- Spatial coherence

- Anisotropic, bilateral weight



Thank you!