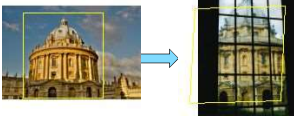


**Recognizing object instances**

Kristen Grauman  
UT-Austin



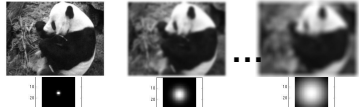
### Announcements

- Assignment 1 is out, due Fri Sept 22
- Presentation assignments - up this week
  
- Reminder – no laptops, phones, etc. in class please

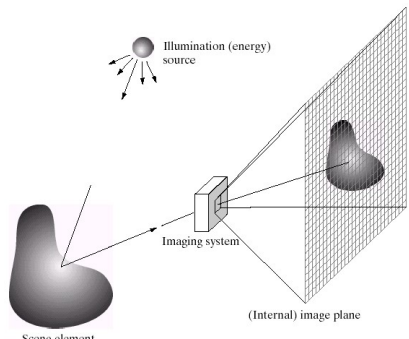
### Plan for today

- 1. Basics in feature extraction: filtering
- 2. Invariant local features
- 3. Recognizing object instances

### Basics in feature extraction

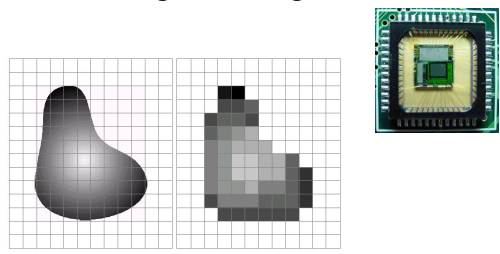


### Image Formation



Slide credit: Derek Hoiem

### Digital images



**FIGURE 2.17** (a) Continuous image projected onto a sensor array. (b) Result of image sampling and quantization.

Slide credit: Derek Hoiem

### Digital images

- **Sample** the 2D space on a regular grid
- **Quantize** each sample (round to nearest integer)
- Image thus represented as a matrix of integer values.

	$j$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32		
$i$	0	82	79	23	119	220	105	4	0																											
1	10	10	9	62	12	79	34	0																												
2	10	59	197	46	46	0	0	49																												
3	176	135	5	189	191	68	0	49																												
4	27	11	11	29	26	37	0	77																												
5	0	89	144	147	187	102	62	208																												
6	255	252	0	166	123	62	0	31																												
7	166	63	127	17	1	0	99	26																												

Adapted from S. Seitz

### Digital color images

Bayer filter

© 2000 How Stuff Works

### Digital color images

Color images, RGB color space

R G B

Kristen Grauman

### Main idea: image filtering

- Compute a function of the local neighborhood at each pixel in the image
  - Function specified by a “filter” or mask saying how to combine values from neighbors.
- Uses of filtering:
  - Enhance an image (denoise, resize, etc)
  - Extract information (texture, edges, etc)
  - Detect patterns (template matching)

Adapted from Derek Hoiem

### Motivation: noise reduction

- Even multiple images of the **same static scene** will not be identical.

Kristen Grauman

### Motivation: noise reduction

- Even multiple images of the same static scene will not be identical.
- How could we reduce the noise, i.e., give an estimate of the true intensities?
- **What if there's only one image?**

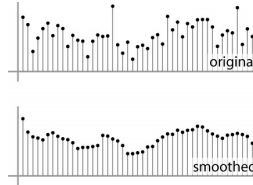
Kristen Grauman

First attempt at a solution

- Let's replace each pixel with an average of all the values in its neighborhood
- Assumptions:
  - Expect pixels to be like their neighbors
  - Expect noise processes to be independent from pixel to pixel

First attempt at a solution

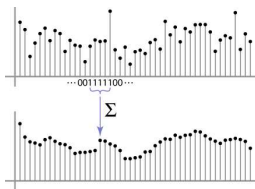
- Let's replace each pixel with an average of all the values in its neighborhood
- Moving average in 1D:



Source: S. Marschner

Weighted Moving Average

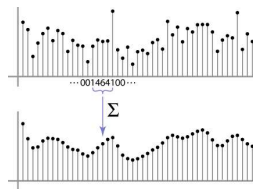
Can add weights to our moving average  
Weights [1, 1, 1, 1, 1] / 5



Source: S. Marschner

Weighted Moving Average

Non-uniform weights [1, 4, 6, 4, 1] / 16

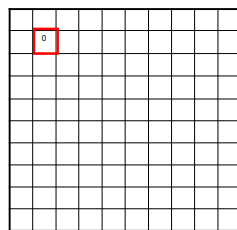
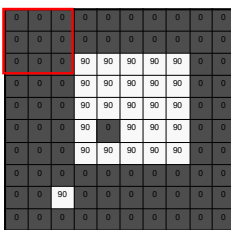


Source: S. Marschner

Moving Average In 2D

$F[x, y]$

$G[x, y]$

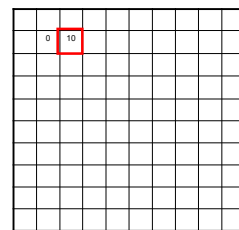
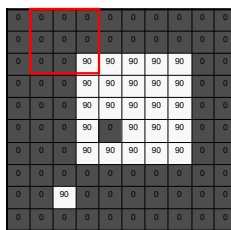


Source: S. Seitz

Moving Average In 2D

$F[x, y]$

$G[x, y]$



Source: S. Seitz

### Moving Average In 2D

$F[x, y]$

$G[x, y]$

Source: S. Seitz

### Moving Average In 2D

$F[x, y]$

$G[x, y]$

Source: S. Seitz

### Moving Average In 2D

$F[x, y]$

$G[x, y]$

Source: S. Seitz

### Moving Average In 2D

$F[x, y]$

$G[x, y]$

Source: S. Seitz

### Correlation filtering

Say the averaging window size is  $2k+1 \times 2k+1$ :

$$G[i, j] = \frac{1}{(2k+1)^2} \sum_{u=-k}^k \sum_{v=-k}^k F[i+u, j+v]$$

Attribute uniform weight to each pixel
Loop over all pixels in neighborhood around image pixel  $F[i, j]$

Now generalize to allow **different weights** depending on neighboring pixel's relative position:

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k \underbrace{H[u, v]}_{\text{Non-uniform weights}} F[i+u, j+v]$$

### Correlation filtering

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i+u, j+v]$$

This is called **cross-correlation**, denoted  $G = H \otimes F$

Filtering an image: replace each pixel with a linear combination of its neighbors.

The filter "**kernel**" or "**mask**"  $H[u, v]$  is the prescription for the weights in the linear combination.



### Averaging filter

- What values belong in the kernel  $H$  for the moving average example?

$F[x, y]$

$\otimes$

$H[u, v]$   

“box filter”

$G[x, y]$

$G = H \otimes F$

### Smoothing by averaging

depicts box filter:  
white = high value, black = low value

original filtered

What if the filter size was 5 x 5 instead of 3 x 3?

### Gaussian filter

- What if we want nearest neighboring pixels to have the most influence on the output?

$F[x, y]$

$\frac{1}{16}$

$H[u, v]$

This kernel is an approximation of a 2d Gaussian function:

$$h(u, v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{2\sigma^2}}$$

Source: S. Seltz

### Smoothing with a Gaussian

### Gaussian filters

- What parameters matter here?
- Variance** of Gaussian: determines extent of smoothing

$\sigma = 2$  with  
30 x 30  
kernel

$\sigma = 5$  with  
30 x 30  
kernel

Kristen Grauman

### Smoothing with a Gaussian

Parameter  $\sigma$  is the “scale” / “width” / “spread” of the Gaussian kernel, and controls the amount of smoothing.

...

```

for sigma=1:3:10
    h = fspecial('gaussian', fsize, sigma);
    out = imfilter(im, h);
    imshow(out);
    pause;
end
    
```


Kristen Grauman

### Properties of smoothing filters

- Smoothing
  - Values positive
  - Sum to 1 → \_\_\_\_\_
  - Amount of smoothing proportional to mask size
  - Remove "high-frequency" components; "low-pass" filter


Kristen Grauman

### Predict the outputs using correlation filtering


 $\times$ 


0	0	0
0	1	0
0	0	0

 $= ?$


 $\times$ 

0	0	0
0	0	1
0	0	0

 $= ?$


 $\times$ 


0	0	0
0	2	0
0	0	0

 $= \frac{1}{9}$ 

1	1	1
1	1	1
1	1	1

 $= ?$

### Practice with linear filters



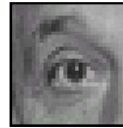
0	0	0
0	1	0
0	0	0

?


Original

Source: D. Lowe

### Practice with linear filters




0	0	0
0	1	0
0	0	0



Original Filtered (no change)

Source: D. Lowe

### Practice with linear filters




0	0	0
0	0	1
0	0	0

?


Original

Source: D. Lowe

### Practice with linear filters




0	0	0
0	0	1
0	0	0



Original Shifted left by 1 pixel with correlation

Source: D. Lowe

Practice with linear filters

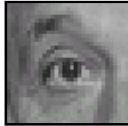


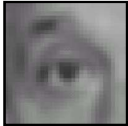
$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad ?$$

Original

Source: D. Lowe

Practice with linear filters




$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$


Original

Blur (with a box filter)

Source: D. Lowe

Practice with linear filters





$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad ?$$

Original

Source: D. Lowe


Practice with linear filters



$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$


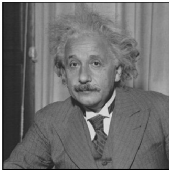
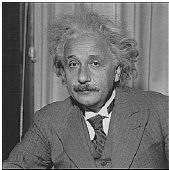
Original

Sharpening filter: accentuates differences with local average



Source: D. Lowe

Filtering examples: sharpening

before


after

Filtering application: Hybrid Images

What you see...

From Far Away

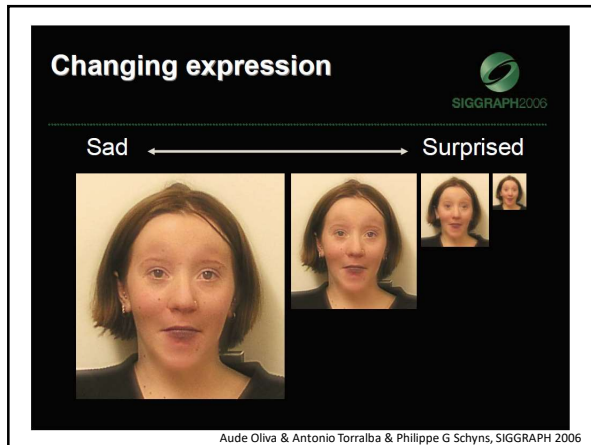
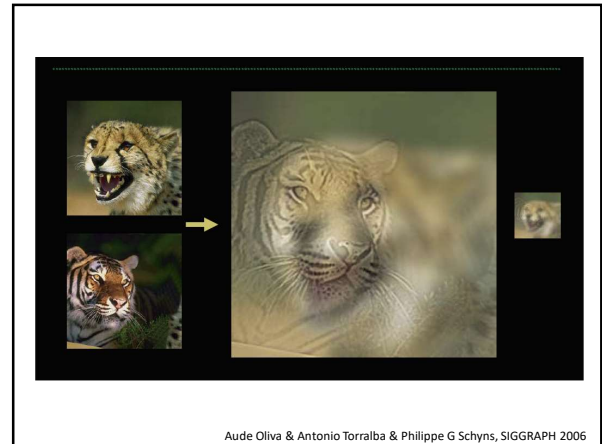
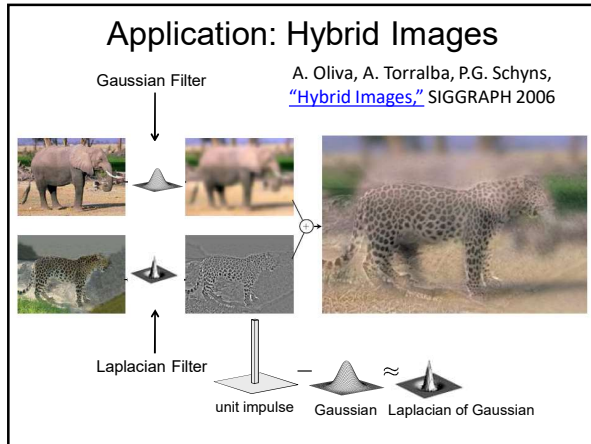
Up Close



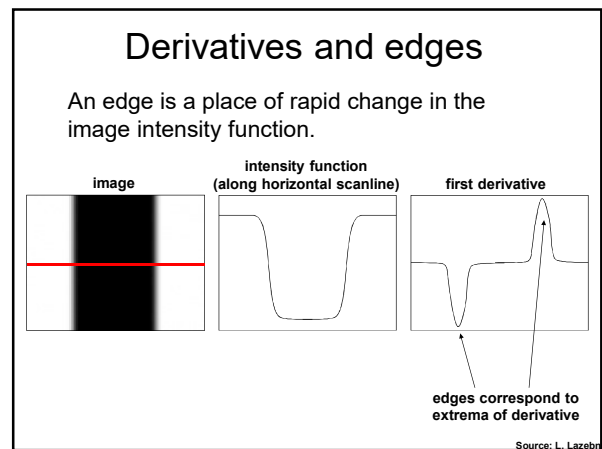
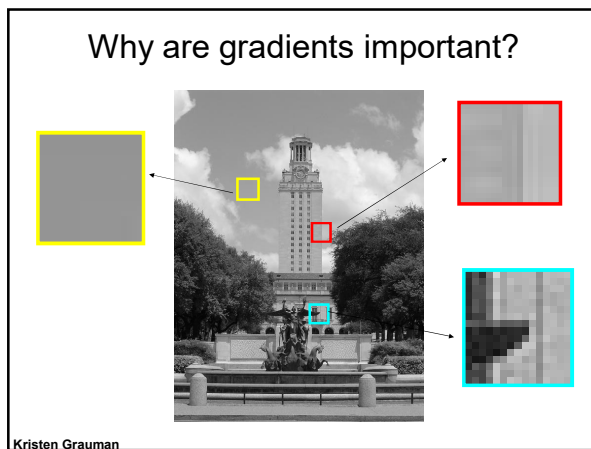
I see an angry guy

It's a woman!

Aude Oliva & Antonio Torralba & Philippe G Schyns, SIGGRAPH 2006



- ### Main idea: image filtering
- Compute a function of the local neighborhood at each pixel in the image
    - Function specified by a "filter" or mask saying how to combine values from neighbors.
  - Uses of filtering:
    - Enhance an image (denoise, resize, etc)
    - Extract information (texture, edges, etc)
    - Detect patterns (template matching)



### Derivatives with convolution

For 2D function,  $f(x,y)$ , the partial derivative is:

$$\frac{\partial f(x,y)}{\partial x} = \lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon, y) - f(x,y)}{\epsilon}$$


For discrete data, we can approximate using finite differences:

$$\frac{\partial f(x,y)}{\partial x} \approx \frac{f(x+1,y) - f(x,y)}{1}$$

To implement above as convolution, what would be the associated filter?

Kristen Grauman

### Partial derivatives of an image



$\frac{\partial f(x,y)}{\partial x}$ 
 $\frac{\partial f(x,y)}{\partial y}$

-1 1

-1 ? 1  
1 or -1

Which shows changes with respect to x?

(showing filters for correlation)

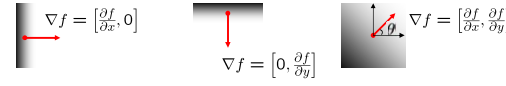
Kristen Grauman

### Image gradient

The gradient of an image:

$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

The gradient points in the direction of most rapid change in intensity



The **gradient direction** (orientation of edge normal) is given by:

$$\theta = \tan^{-1} \left( \frac{\partial f / \partial y}{\partial f / \partial x} \right)$$

Slide credit: Steve Seitz

### Mask properties

- **Smoothing**
  - Values positive
  - Sum to 1 → constant regions same as input
  - Amount of smoothing proportional to mask size
  - Remove "high-frequency" components; "low-pass" filter
- **Derivatives**
  - \_\_\_\_\_ signs used to get high response in regions of high contrast
  - Sum to \_\_\_\_ → no response in constant regions
  - High absolute value at points of high contrast

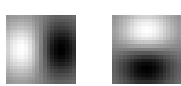
Kristen Grauman

### Main idea: image filtering

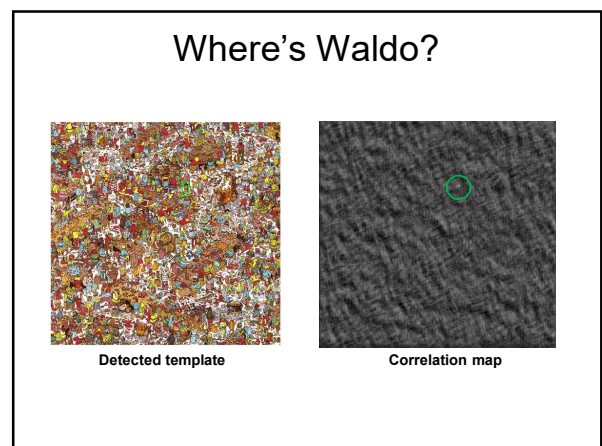
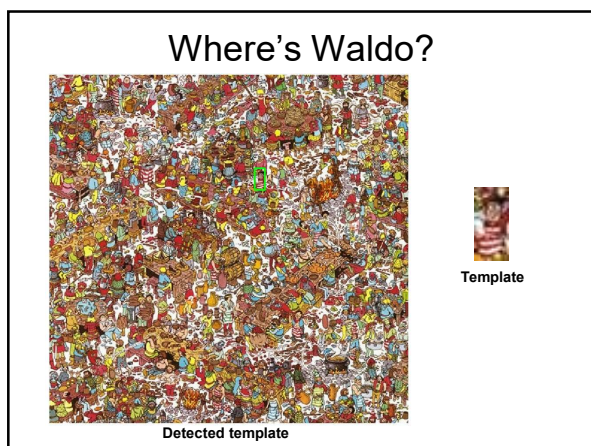
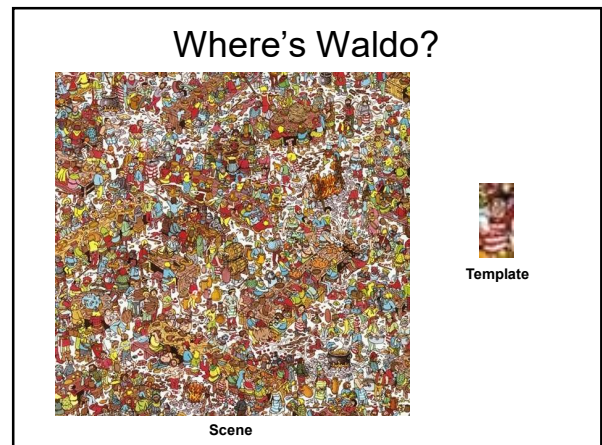
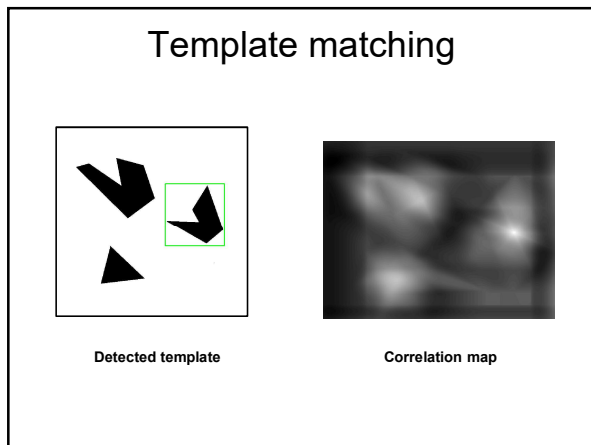
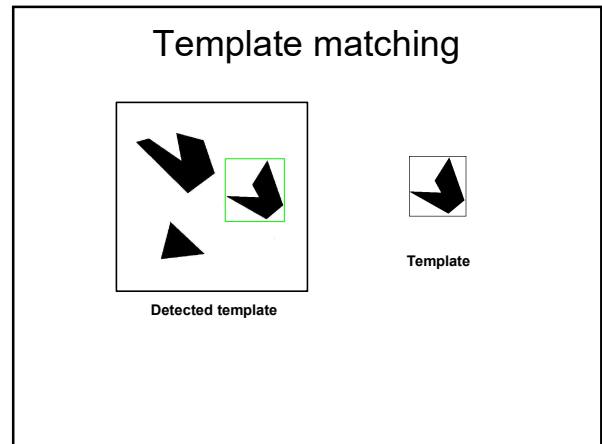
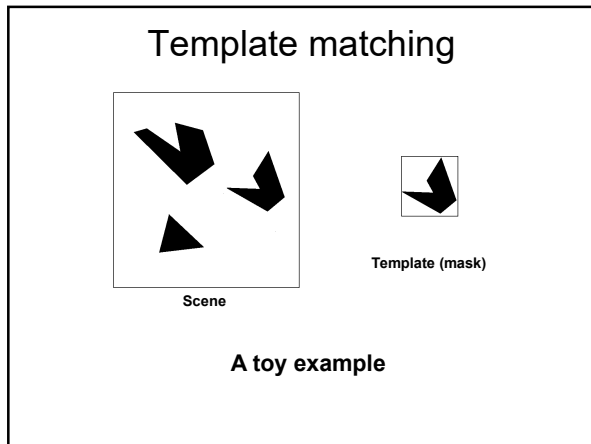
- Compute a function of the local neighborhood at each pixel in the image
  - Function specified by a "filter" or mask saying how to combine values from neighbors.
- Uses of filtering:
  - Enhance an image (denoise, resize, etc)
  - Extract information (texture, edges, etc)
  - Detect patterns (template matching)

### Template matching

- Filters as **templates**:
  - Note that filters look like the effects they are intended to find --- "matched filters"




- Use normalized cross-correlation score to find a given pattern (template) in the image.
- Normalization needed to control for relative brightnesses.






### Template matching



Scene                      Template

What if the template is not identical to some subimage in the scene?

### Template matching



Detected template                      Template

Match can be meaningful, if scale, orientation, and general appearance is right.  
...but we can do better!...


### Summary so far

- Compute a function of the local neighborhood at each pixel in the image
  - Function specified by a “filter” or mask saying how to combine values from neighbors.
- Uses of filtering:
  - Enhance an image (denoise, resize, etc)
  - Extract information (texture, edges, etc)
  - Detect patterns (template matching)

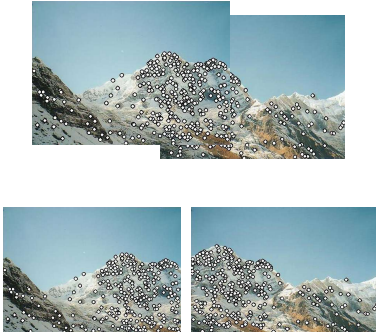
### Plan for today

- 1. Basics in feature extraction: filtering
- **2. Invariant local features**
- 3. Specific object recognition methods

### Local features: detection and description

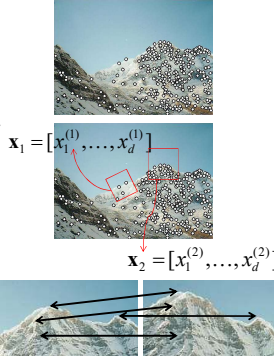


### Basic goal



### Local features: main components

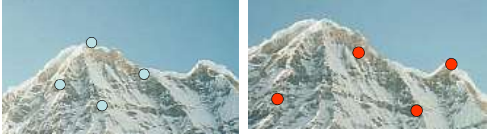
- 1) Detection: Identify the interest points
- 2) Description: Extract vector feature descriptor surrounding each interest point.  $\mathbf{x}_1 = [x_1^{(1)}, \dots, x_d^{(1)}]$
- 3) Matching: Determine correspondence between descriptors in two views  $\mathbf{x}_2 = [x_1^{(2)}, \dots, x_d^{(2)}]$



Kristen Grauman

### Goal: interest operator repeatability

- We want to detect (at least some of) the same points in both images.

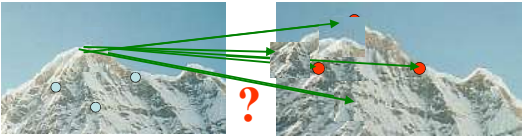


No chance to find true matches!

- Yet we have to be able to run the detection procedure *independently* per image.

### Goal: descriptor distinctiveness



- We want to be able to reliably determine which point goes with which.



- Must provide some invariance to geometric and photometric differences between the two views.


### Local features: main components

- 1) Detection: Identify the interest points
- 2) Description: Extract vector feature descriptor surrounding each interest point.
- 3) Matching: Determine correspondence between descriptors in two views

- What points would you choose?

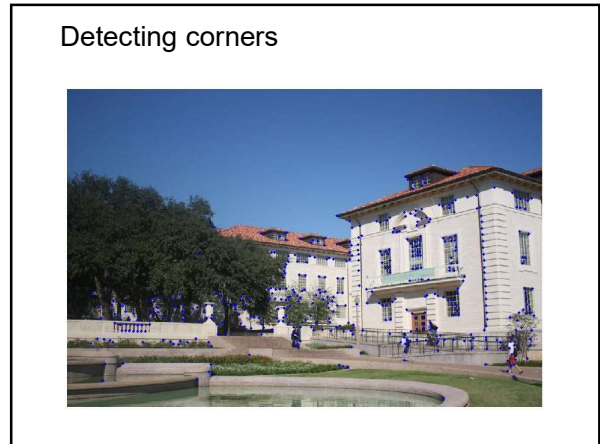
### Detecting corners





### Detecting corners

Compute "cornerness" response at every pixel.



### Detecting local invariant features

- Detection of interest points
  - Harris corner detection
  - Scale invariant blob detection: LoG

### Corners as distinctive interest points

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

“corner”: significant change in all directions

Slide credit: Alysha Efros, Darva Frolova, Denis Simakov

### Corners as distinctive interest points

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

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

Notation:  $I_x \leftrightarrow \frac{\partial I}{\partial x}$   $I_y \leftrightarrow \frac{\partial I}{\partial y}$   $I_x I_y \leftrightarrow \frac{\partial I}{\partial x} \frac{\partial I}{\partial y}$

### What does this matrix reveal?

First, consider an axis-aligned corner:

What does this matrix reveal?

First, consider an axis-aligned corner:

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

This means dominant gradient directions align with x or y axis

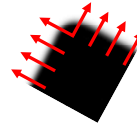
Look for locations where **both**  $\lambda$ 's are large.

If either  $\lambda$  is close to 0, then this is **not** corner-like.

What if we have a corner that is not aligned with the image axes?

What does this matrix reveal?

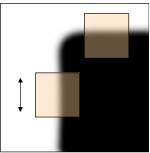
Since  $M$  is symmetric, we have  $M = X \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} X^T$



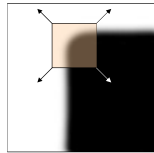
$$Mx_i = \lambda_i x_i$$

The *eigenvalues* of  $M$  reveal the amount of intensity change in the two principal orthogonal gradient directions in the window.

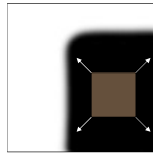
Corner response function



“edge”:  
 $\lambda_1 \gg \lambda_2$   
 $\lambda_2 \gg \lambda_1$



“corner”:  
 $\lambda_1$  and  $\lambda_2$  are large,  
 $\lambda_1 \sim \lambda_2$ ;



“flat” region  
 $\lambda_1$  and  $\lambda_2$  are small;

Cornerness score  
 (other variants possible)

$$f = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2}$$

Harris corner detector

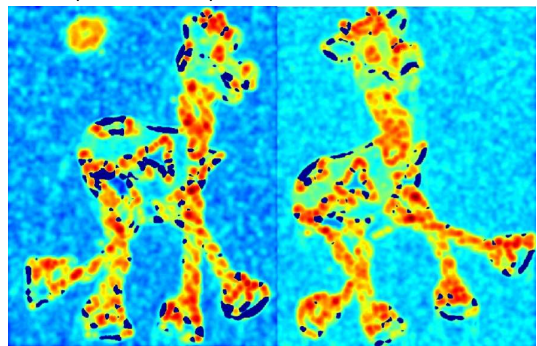
- 1) Compute  $M$  matrix for each image window to get their *cornerness* scores.
- 2) Find points whose surrounding window gave large corner response ( $f >$  threshold)
- 3) Take the points of local maxima, i.e., perform non-maximum suppression

Harris Detector: Steps

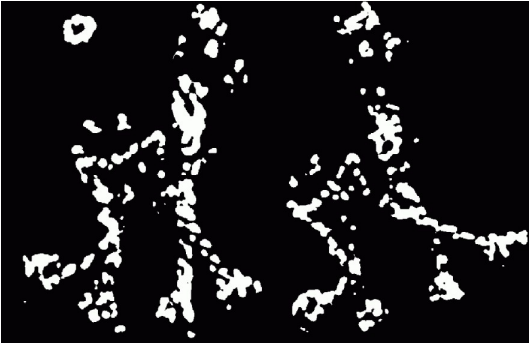


Harris Detector: Steps


Compute corner response  $f$



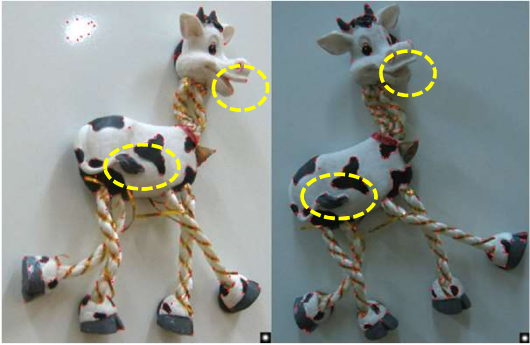
**Harris Detector: Steps**  
 Find points with large corner response:  $f > \text{threshold}$



**Harris Detector: Steps**  
 Take only the points of local maxima of  $f$



**Harris Detector: Steps**



**Properties of the Harris corner detector**

Rotation invariant? Yes

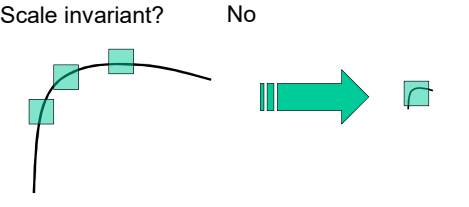
$$M = X \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} X^T$$

Scale invariant?

**Properties of the Harris corner detector**

Rotation invariant? Yes

Scale invariant? No




All points will be classified as edges


Corner !

**Scale invariant interest points**

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



### Automatic Scale Selection



How to find corresponding patch sizes, with only one image in hand?

K. Grauman, B. Leibe

### Automatic scale selection

**Intuition:**

- Find scale that gives local maxima of some function  $f$  in both position and scale.

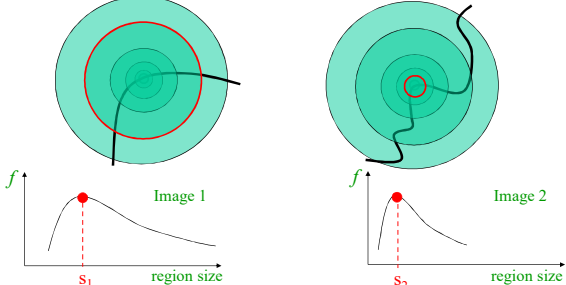


Image 1

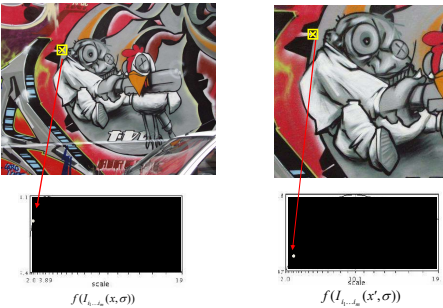
Image 2

$S_1$  region size

$S_2$  region size

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)



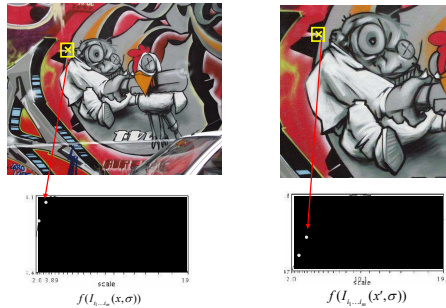
$f(U_{x \rightarrow x'}(x, \sigma))$

$f(U_{x' \rightarrow x}(x', \sigma))$

K. Grauman, B. Leibe

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)



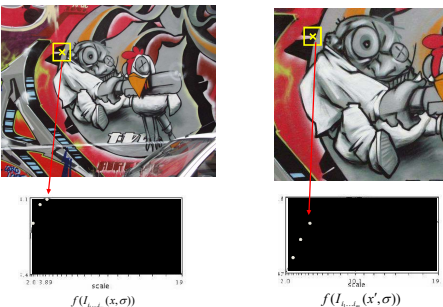
$f(U_{x \rightarrow x'}(x, \sigma))$

$f(U_{x' \rightarrow x}(x', \sigma))$

K. Grauman, B. Leibe

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)



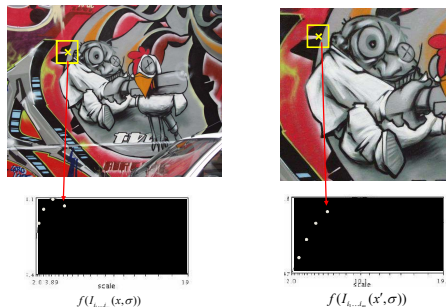
$f(U_{x \rightarrow x'}(x, \sigma))$

$f(U_{x' \rightarrow x}(x', \sigma))$

K. Grauman, B. Leibe

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)



$f(U_{x \rightarrow x'}(x, \sigma))$

$f(U_{x' \rightarrow x}(x', \sigma))$

K. Grauman, B. Leibe

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)

$f(U_{i,j}(x, \sigma))$   $f(U_{i,j}(x', \sigma'))$

K. Grauman, B. Leibe

### Automatic Scale Selection

- Function responses for increasing scale (scale signature)

$f(U_{i,j}(x, \sigma))$   $f(U_{i,j}(x', \sigma'))$

K. Grauman, B. Leibe

What can be the "signature" function?

### Blob detection in 2D

Laplacian of Gaussian: Circularly symmetric operator for blob detection in 2D

$$\nabla^2 g = \frac{\partial^2 g}{\partial x^2} + \frac{\partial^2 g}{\partial y^2}$$

### Blob detection in 2D: scale selection

Laplacian-of-Gaussian = "blob" detector  $\nabla^2 g = \frac{\partial^2 g}{\partial x^2} + \frac{\partial^2 g}{\partial y^2}$

filter scales

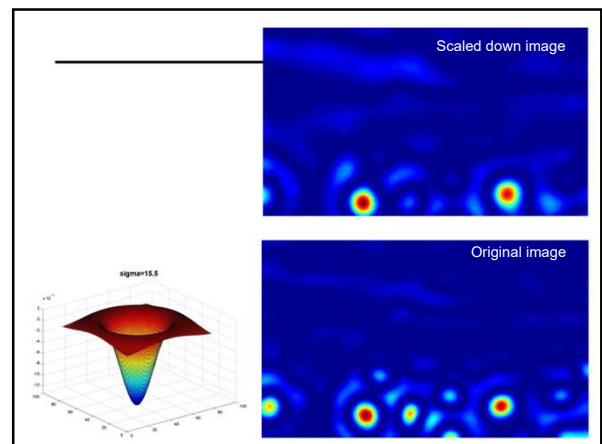
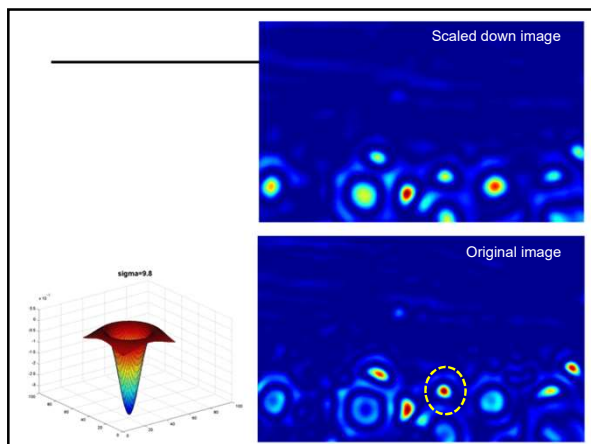
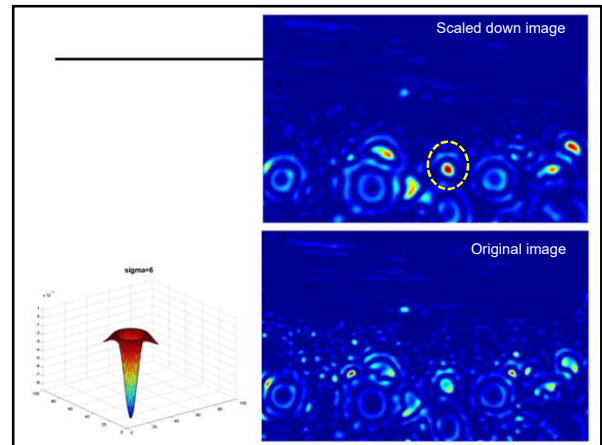
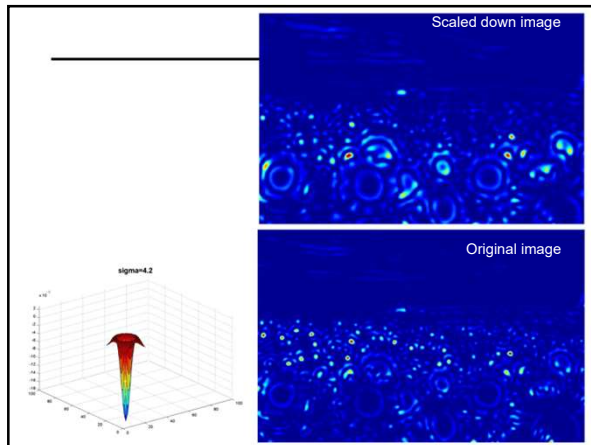
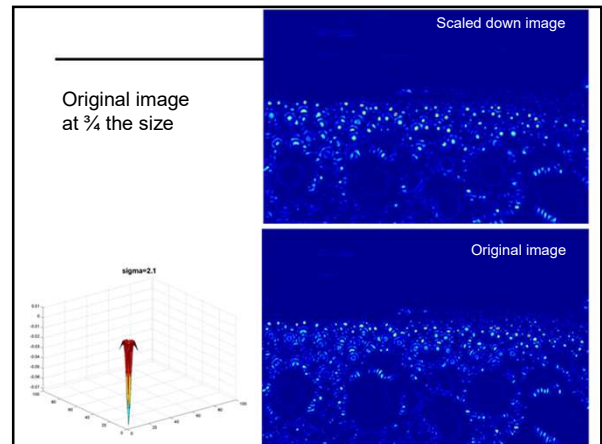
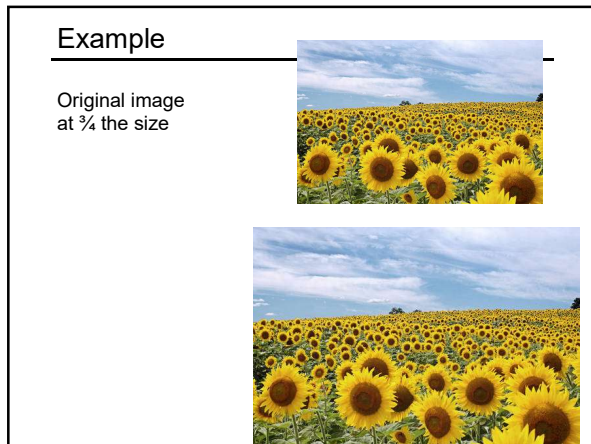
img1 img2 img3

### Blob detection in 2D

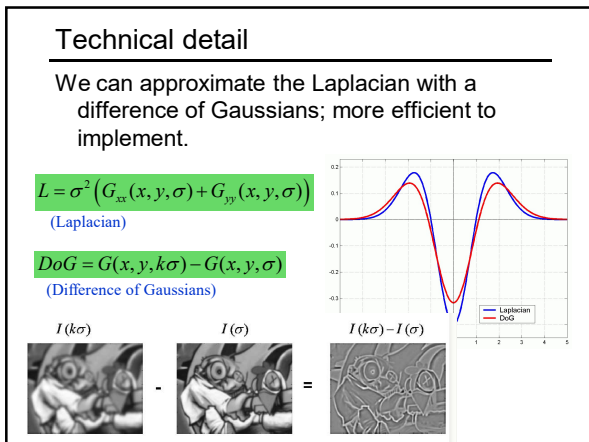
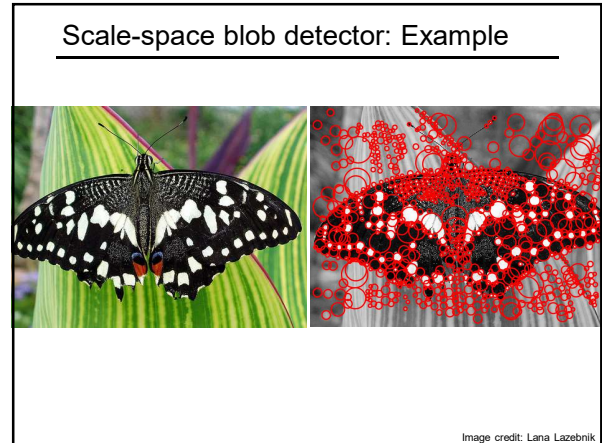
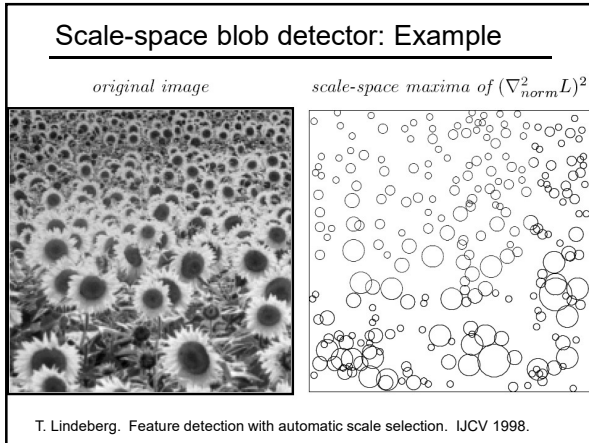
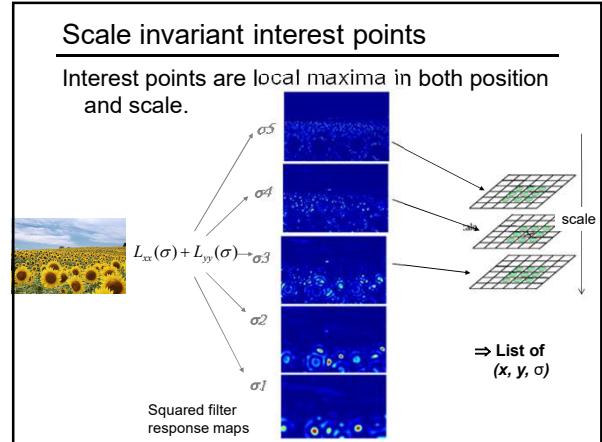
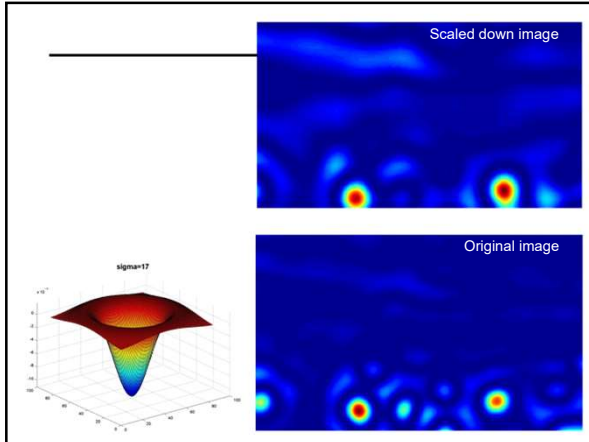
We define the *characteristic scale* as the scale that produces peak of Laplacian response

characteristic scale

Slide credit: Lana Lazebnik

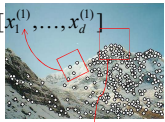






- ### Recap so far: interest points
- Interest point detection
    - Harris corner detector
    - Laplacian of Gaussian, automatic scale selection

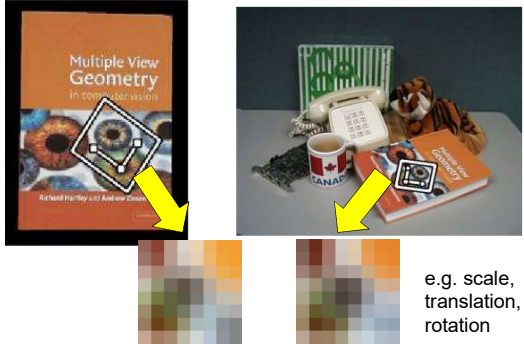
### Local features: main components

- 1) Detection: Identify the interest points
- 2) Description: Extract vector feature descriptor surrounding each interest point.
 
$$\mathbf{x}_1 = [x_1^{(1)}, \dots, x_d^{(1)}]$$


$$\mathbf{x}_2 = [x_1^{(2)}, \dots, x_d^{(2)}]$$
- 3) Matching: Determine correspondence between descriptors in two views

Kristen Grauman

### Geometric transformations



e.g. scale, translation, rotation

### Photometric transformations


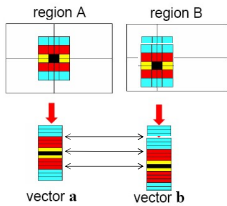


Figure from T. Tuytelaars ECCV 2006 tutorial

### Raw patches as local descriptors

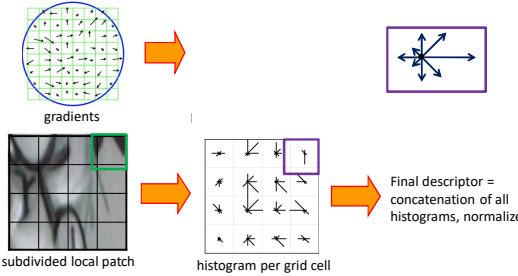


The simplest way to describe the neighborhood around an interest point is to write down the list of intensities to form a feature vector.

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

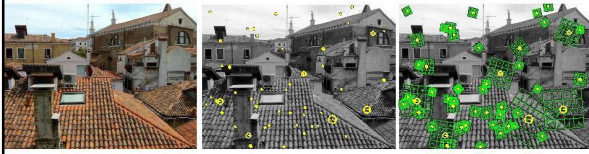
### Scale Invariant Feature Transform (SIFT) descriptor [Lowe 2004]

- Use histograms to bin pixels within sub-patches according to their orientation.



Final descriptor = concatenation of all histograms, normalized

### Scale Invariant Feature Transform (SIFT) descriptor [Lowe 2004]



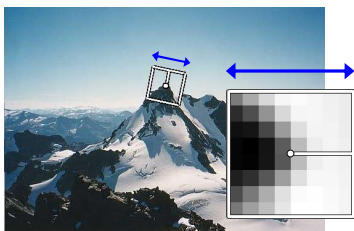
Interest points and their scales and orientations (random subset of 50)

SIFT descriptors

<http://www.vifecat.org/overview/sift.html>



### Making descriptor rotation invariant

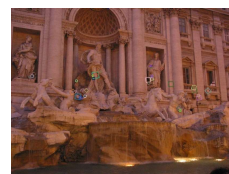


- Rotate patch according to its dominant gradient orientation
- This puts the patches into a canonical orientation.

Image from Matthew Brown

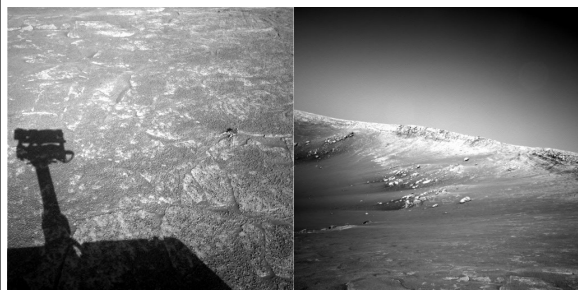
### SIFT descriptor [Lowe 2004]

- Extraordinarily robust matching technique
  - Can handle changes in viewpoint
    - Up to about 60 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, e.g. <http://www.vlfeat.org/overview/sift.html>



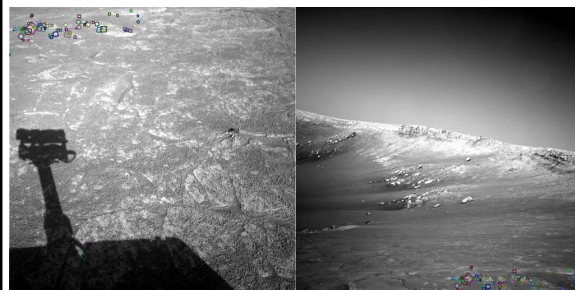
Steve Seitz

### Example



NASA Mars Rover images

### Example



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

### SIFT properties

- Invariant to
  - Scale
  - Rotation
- Partially invariant to
  - Illumination changes
  - Camera viewpoint
  - Occlusion, clutter


### Local features: main components

- 1) Detection: Identify the interest points
- 2) Description: Extract vector feature descriptor surrounding each interest point.
- 3) Matching: Determine correspondence between descriptors in two views



Kristen Grauman

### Matching local features



### Matching local features

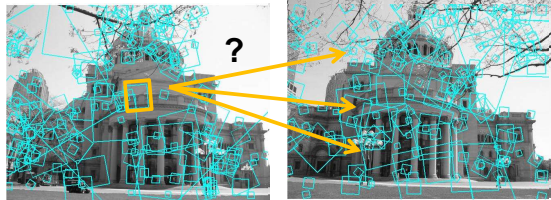


Image 1                      Image 2

To generate **candidate matches**, find patches that have the most similar appearance (e.g., lowest SSD)  
 Simplest approach: compare them all, take the closest (or closest k, or within a thresholded distance)

### Ambiguous matches

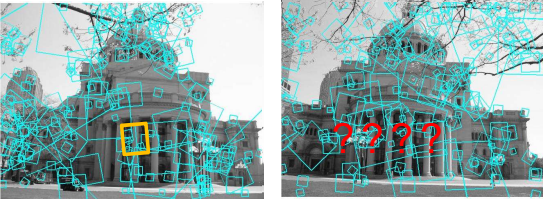
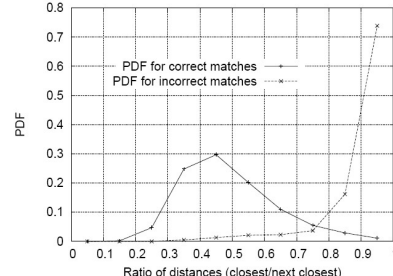


Image 1                      Image 2

At what SSD value do we have a good match?  
 To add robustness to matching, can consider **ratio** :  
 distance to best match / distance to second best match  
 If low, first match looks good.  
 If high, could be ambiguous match.

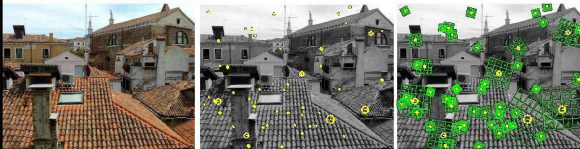
### Matching SIFT Descriptors

- Nearest neighbor (Euclidean distance)
- Threshold ratio of nearest to 2<sup>nd</sup> nearest descriptor



Lowe IJCV 2004

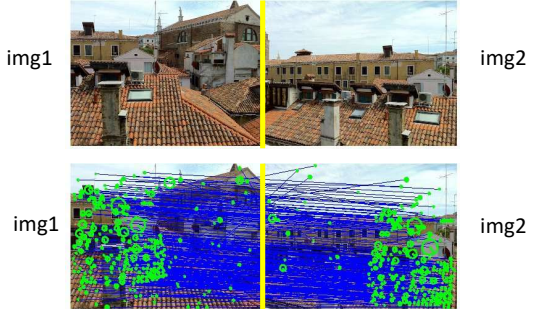
### Scale Invariant Feature Transform (SIFT) descriptor [Lowe 2004]



Interest points and their scales and orientations (random subset of 50)                      SIFT descriptors

<http://www.vifeat.org/overview/sift.html>

### SIFT (preliminary) matches



<http://www.vifeat.org/overview/sift.html>

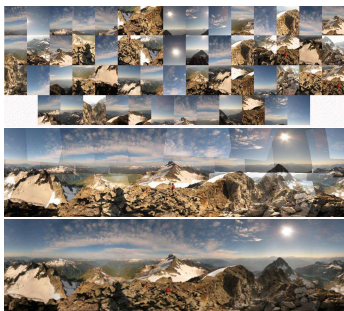
### Value of local (invariant) features

- Complexity reduction via selection of distinctive points
- Describe images, objects, parts without requiring segmentation
- Local character means robustness to clutter, occlusion
- Robustness: similar descriptors in spite of noise, blur, etc.

### Applications of local invariant features

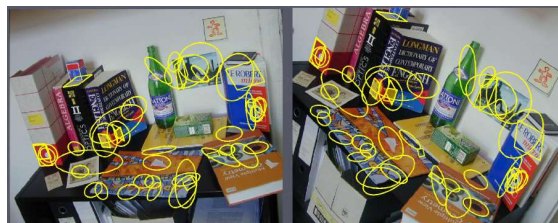
- Wide baseline stereo
- Motion tracking
- Panoramas
- Mobile robot navigation
- 3D reconstruction
- Recognition
- ...

### Automatic mosaicing



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

### Wide baseline stereo



[Image from T. Tuytelaars ECCV 2006 tutorial]

### Photo tourism [Snavely et al.]



### Recognition of specific objects, scenes



Many current applications - 2017

### Summary so far

- Interest point detection
  - Harris corner detector
  - Laplacian of Gaussian, automatic scale selection
- Invariant descriptors
  - Rotation according to dominant gradient direction
  - Histograms for robustness to small shifts and translations (SIFT descriptor)

### Plan for today

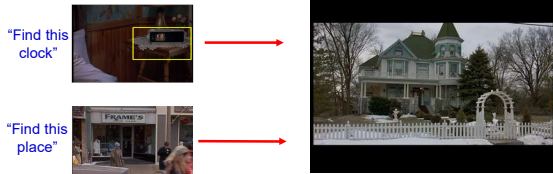
- 1. Basics in feature extraction: filtering
- 2. Invariant local features
- **3. Recognizing object instances**

### Recognizing or retrieving specific objects

Example I: Visual search in feature films

Visually defined query

"Groundhog Day" [Ramms, 1993]



Slide credit: J. Sivic

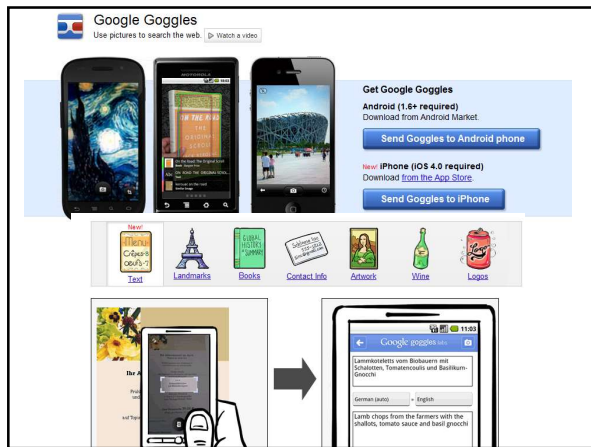
### Recognizing or retrieving specific objects

Example II: Search photos on the web for particular places



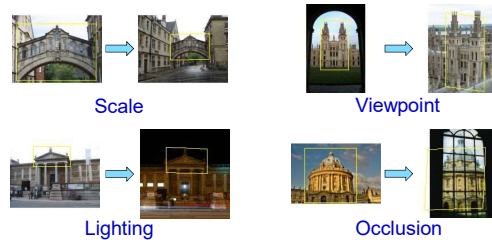
...in these images and 1M more

Slide credit: J. Sivic



### Why is it difficult?

Want to find the object despite possibly large changes in scale, viewpoint, lighting and partial occlusion



We can't expect to match such varied instances with a single global template...

Slide credit: J. Sivic



### Instance recognition: key new ideas

- **Visual words**
  - quantization, index, bags of words
- **Spatial verification**
  - affine; RANSAC, Hough

### Indexing local features

- Each patch / region has a descriptor, which is a point in some high-dimensional feature space (e.g., SIFT)

Descriptor's feature space

Kristen Grauman

### Indexing local features

- When we see close points in feature space, we have similar descriptors, which indicates similar local content.

Database images

Query image

Descriptor's feature space

Easily can have millions of features to search!

Kristen Grauman

### Indexing local features: inverted file index

- For text documents, an efficient way to find all pages on which a word occurs is to use an index...
- We want to find all images in which a feature occurs.
- To use this idea, we'll need to map our features to "visual words".

Kristen Grauman

### Visual words

- Map high-dimensional descriptors to tokens/words by quantizing the feature space

- Quantize via clustering, let cluster centers be the prototype "words"
- Determine which word to assign to each new image region by finding the closest cluster center.

Descriptor's feature space

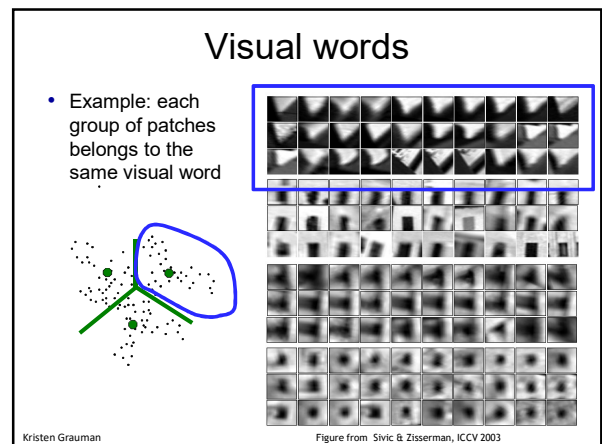
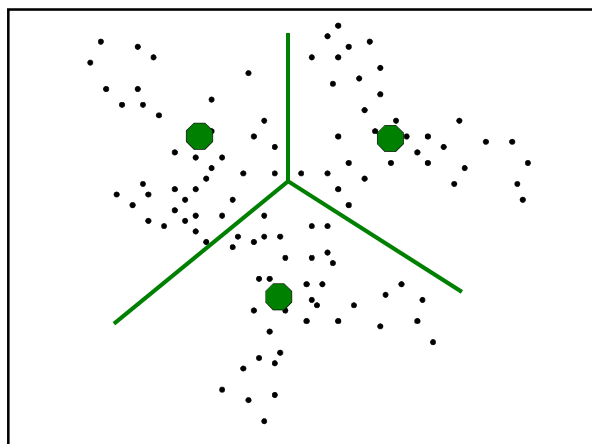
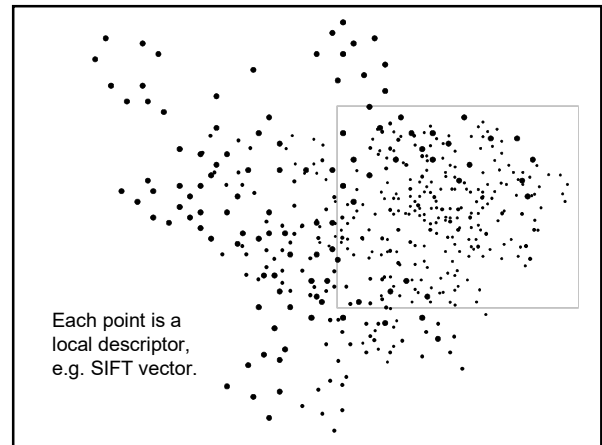
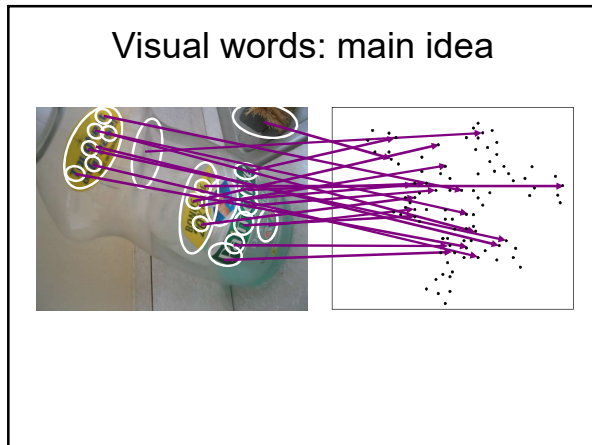
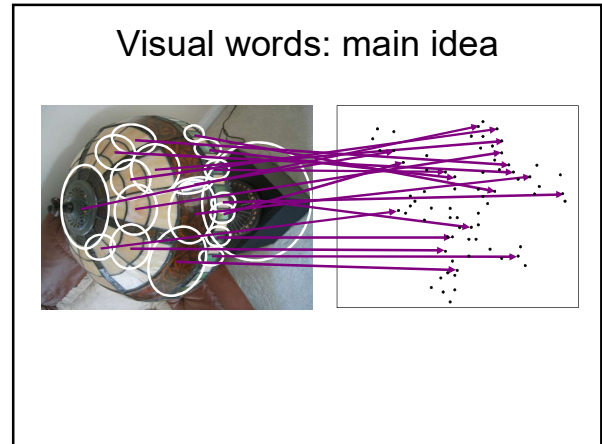
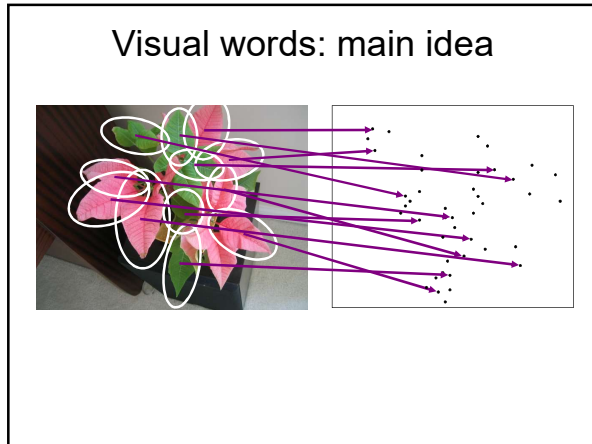
Kristen Grauman

### Visual words: main idea

- Extract some local features from a number of images ...

e.g., SIFT descriptor space: each point is 128-dimensional

Slide credit: D. Nister, CVPR 2006



### Inverted file index

Database images

Word #	Image #
1	3
2	
...	
7	1, 2
8	3
9	
10	
...	
91	2

- Database images are loaded into the index mapping words to image numbers

Kristen Grauman

### Inverted file index

New query image

Word #	Image #
1	3
2	
...	
7	1, 2
8	3
9	
10	
...	
91	2

- New query image is mapped to indices of database images that share a word.

Kristen Grauman

### Instance recognition: remaining issues

- How to summarize the content of an entire image? And gauge overall similarity?
- How large should the vocabulary be? How to perform quantization efficiently?
- Is having the same set of visual words enough to identify the object/scene? How to verify spatial agreement?

Kristen Grauman

### Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the visual centers of the brain were thought of as a movie camera that merely records the image falling on the retina. We now know that perception is a more complex process, involving the following the messages with the various parts of the cortex. Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a wise analysis in a system of nerve cells. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports, compared with \$660bn in 2004. China's deliberate yuan is also needed to meet the demand so far with the country. China's yuan against the dollar is permitted to trade within a narrow band, but the US wants the yuan to be allowed to rise freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

ICCV 2005 short course, L. Fei-Fei

### Bags of visual words

- Summarize entire image based on its distribution (histogram) of word occurrences.
- Analogous to bag of words representation commonly used for documents.

### Comparing bags of words

- Rank frames by normalized scalar product between their (possibly weighted) occurrence counts---nearest neighbor search for similar images.

[1 8 1 4]

$\vec{d}_j$

[5 1 1 0]

$\vec{q}$

$$sim(d_j, q) = \frac{\langle d_j, q \rangle}{\|d_j\| \|q\|}$$

$$= \frac{\sum_{i=1}^V d_j(i) * q(i)}{\sqrt{\sum_{i=1}^V d_j(i)^2} * \sqrt{\sum_{i=1}^V q(i)^2}}$$

for vocabulary of V words

### Inverted file index and bags of words similarity

New query image

Word #	Image #
1	3
2	
7	1, 2
8	3
9	
10	
91	2

1. Extract words in query
2. Inverted file index to find relevant frames
3. Compare word counts

Kristen Grauman

### Instance recognition: remaining issues

- How to summarize the content of an entire image? And gauge overall similarity?
- How large should the vocabulary be? How to perform quantization efficiently?
- Is having the same set of visual words enough to identify the object/scene? How to verify spatial agreement?

Kristen Grauman

### Vocabulary size

Results for recognition task with 6347 images

*Influence on performance, sparsity?* Nister & Stewenius, CVPR'06

### Vocabulary Trees: hierarchical clustering for large vocabularies

- Tree construction:

[Nister & Stewenius, CVPR'06]  
K. Grauman, B. Leibe  
Slide credit: David Nister

### Vocabulary Tree

[Nister & Stewenius, CVPR'06]  
K. Grauman, B. Leibe  
Slide credit: David Nister



### Vocabulary trees: complexity

Number of words given tree parameters:  
branching factor and number of levels

Word assignment cost vs. flat vocabulary

### Visual words/bags of words

- + flexible to geometry / deformations / viewpoint
- + compact summary of image content
- + provides vector representation for sets
- + very good results in practice
- background and foreground mixed when bag covers whole image
- optimal vocabulary formation remains unclear
- basic model ignores geometry – must verify afterwards, or encode via features

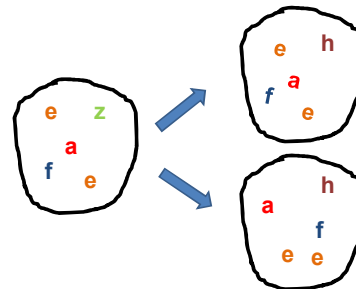
Kristen Grauman

### Instance recognition: remaining issues

- How to summarize the content of an entire image? And gauge overall similarity?
- How large should the vocabulary be? How to perform quantization efficiently?
- Is having the same set of visual words enough to identify the object/scene? How to verify spatial agreement?

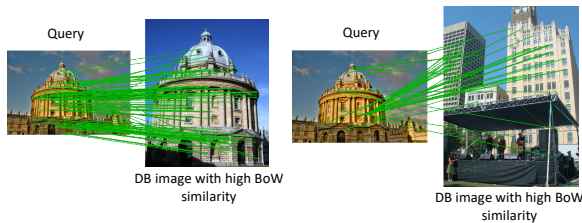
Kristen Grauman

*Which matches better?*



Derek Hoiem

### Spatial Verification



Both image pairs have many visual words in common.

Slide credit: Ondrej Chum

### Spatial Verification



Only some of the matches are mutually consistent

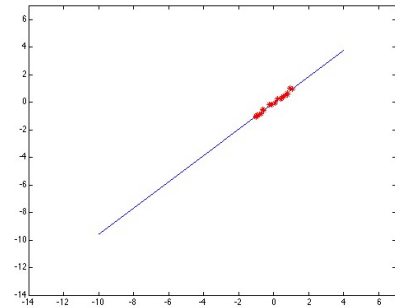
Slide credit: Ondrej Chum

### Spatial Verification: two basic strategies

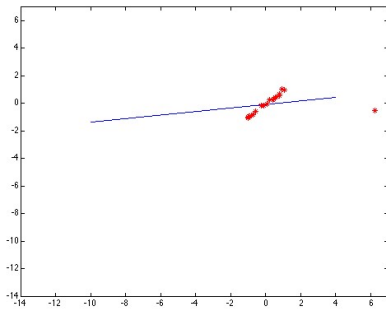
- RANSAC
- Generalized Hough Transform

Kristen Grauman

### Outliers affect least squares fit



### Outliers affect least squares fit



### RANSAC

- RANdom Sample Consensus
- **Approach:** we want to avoid the impact of outliers, so let's look for "inliers", and use those only.
- **Intuition:** if an outlier is chosen to compute the current fit, then the resulting line won't have much support from rest of the points.

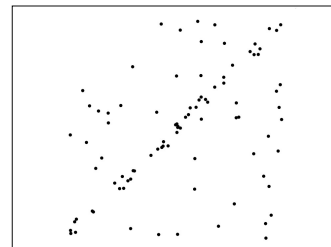
### RANSAC for line fitting

Repeat  $N$  times:

- Draw  $s$  points uniformly at random
- Fit line to these  $s$  points
- Find inliers to this line among the remaining points (i.e., points whose distance from the line is less than  $t$ )
- If there are  $d$  or more inliers, accept the line and refit using all inliers

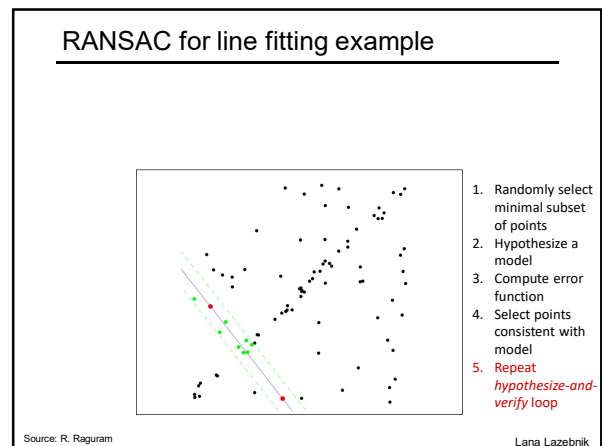
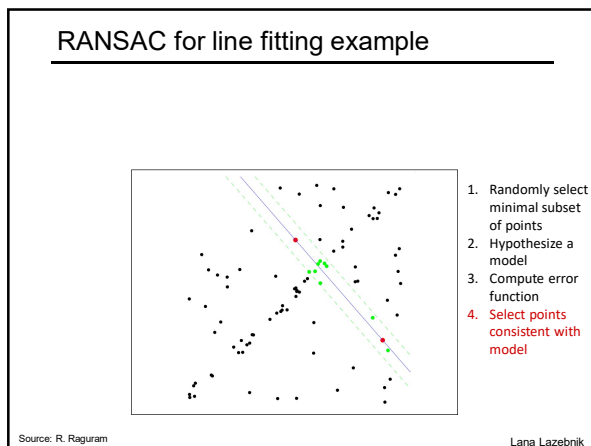
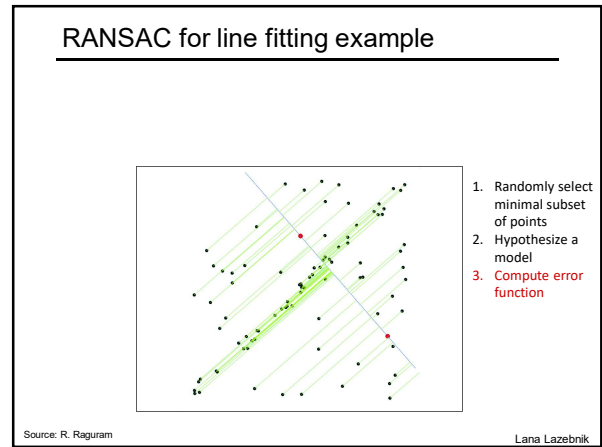
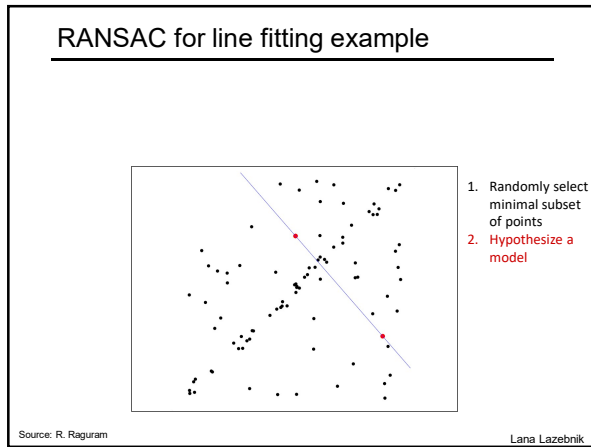
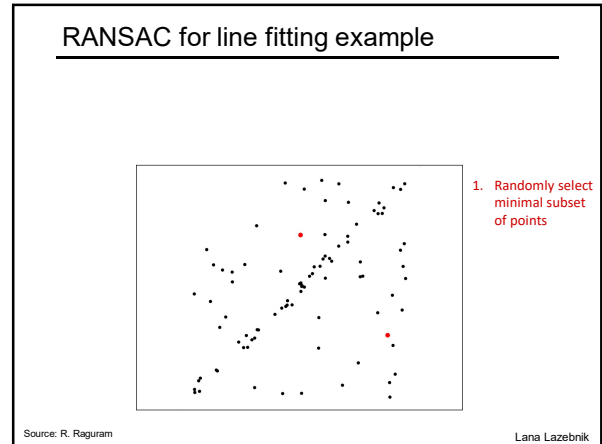
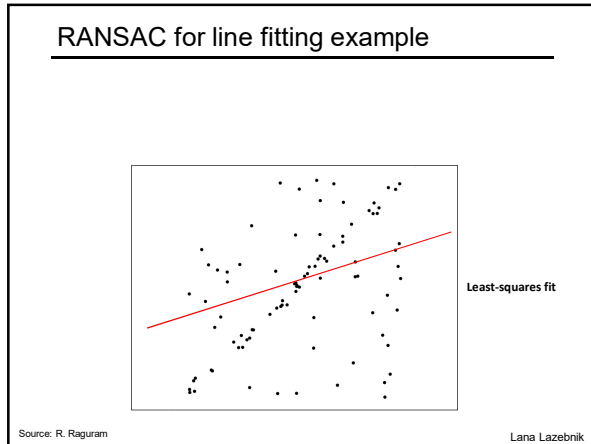
Lana Lazebnik

### RANSAC for line fitting example

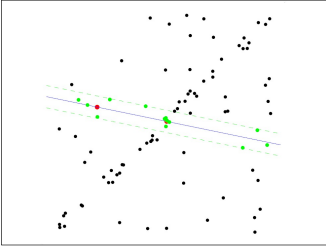


Source: R. Raguram

Lana Lazebnik



### RANSAC for line fitting example



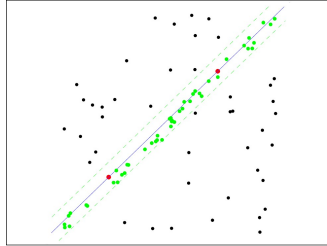
1. Randomly select minimal subset of points
2. Hypothesize a model
3. Compute error function
4. Select points consistent with model
5. Repeat hypothesize-and-verify loop

211

Source: R. Raguram Lana Lazebnik

### RANSAC for line fitting example

**Uncontaminated sample**

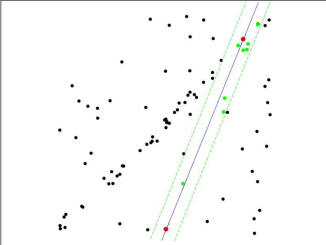


1. Randomly select minimal subset of points
2. Hypothesize a model
3. Compute error function
4. Select points consistent with model
5. Repeat hypothesize-and-verify loop

212

Source: R. Raguram Lana Lazebnik

### RANSAC for line fitting example



1. Randomly select minimal subset of points
2. Hypothesize a model
3. Compute error function
4. Select points consistent with model
5. Repeat hypothesize-and-verify loop

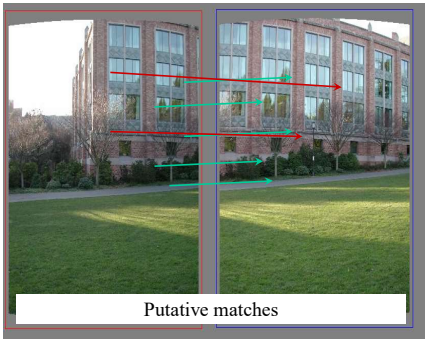
211

Source: R. Raguram Lana Lazebnik

That is an example fitting a model (line)...

What about fitting a transformation (translation)?

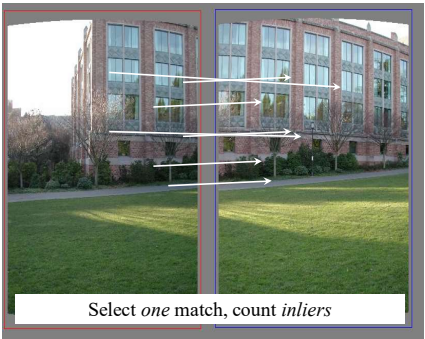
### RANSAC example: Translation



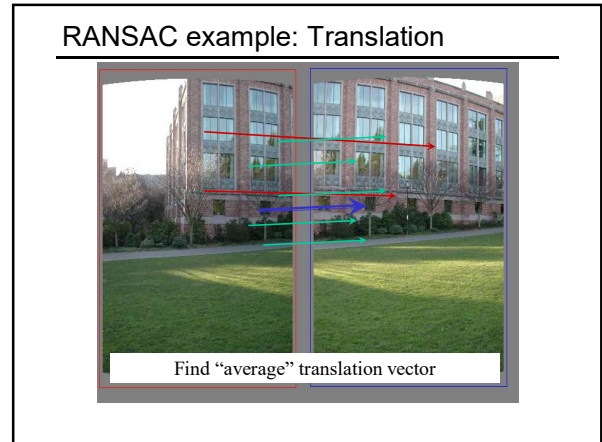
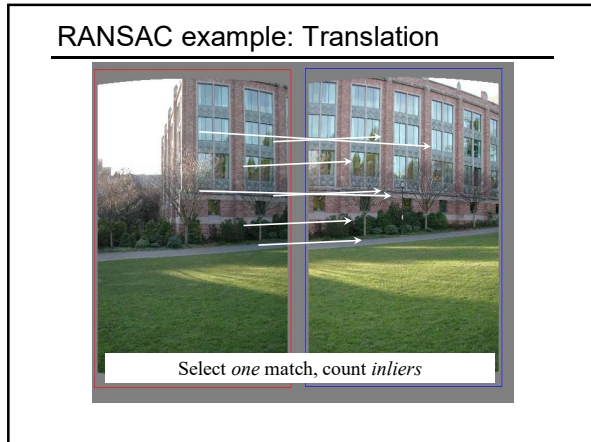
Putative matches

Source: Rick Szeliski

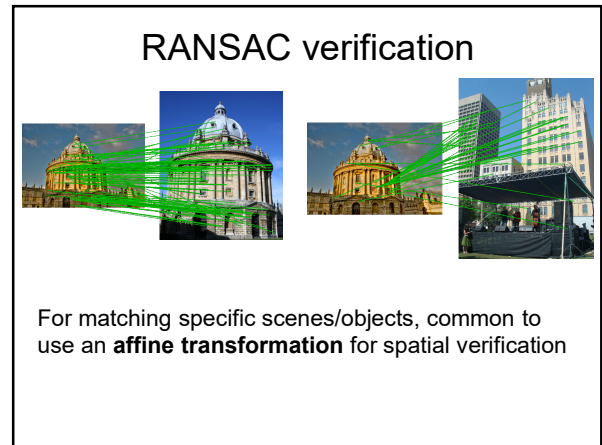
### RANSAC example: Translation



Select one match, count inliers

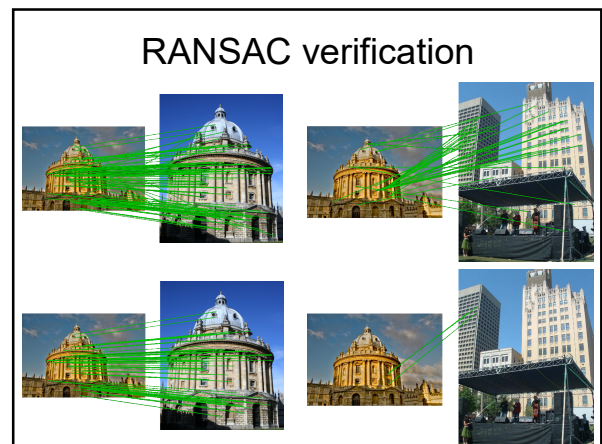


- RANSAC: General form**
- **RANSAC loop:**
    1. Randomly select a *seed group* of points on which to base transformation estimate
    2. Compute model from seed group
    3. Find *inliers* to this transformation
    4. If the number of inliers is sufficiently large, re-compute estimate of model on all of the inliers
  - Keep the model with the largest number of inliers



**Fitting an affine transformation**

Approximates viewpoint changes for roughly planar objects and roughly orthographic cameras.

$$\begin{bmatrix} x'_i \\ y'_i \end{bmatrix} = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x_i \\ y_i \end{bmatrix} + \begin{bmatrix} t_1 \\ t_2 \end{bmatrix}$$


### Spatial Verification: two basic strategies

- RANSAC
  - Typically sort by BoW similarity as initial filter
  - Verify by checking support (inliers) for possible affine transformations
    - e.g., "success" if find an affine transformation with > N inlier correspondences
- Generalized Hough Transform
  - Let each matched feature cast a vote on location, scale, orientation of the model object
  - Verify parameters with enough votes

Kristen Grauman

### Spatial Verification: two basic strategies

- RANSAC
  - Typically sort by BoW similarity as initial filter
  - Verify by checking support (inliers) for possible affine transformations
    - e.g., "success" if find an affine transformation with > N inlier correspondences
- Generalized Hough Transform
  - Let each matched feature cast a vote on location, scale, orientation of the model object
  - Verify parameters with enough votes

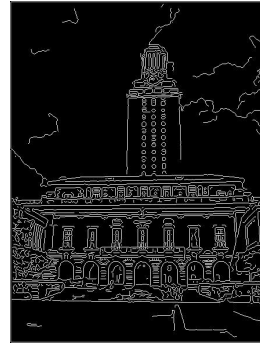
Kristen Grauman

### Voting

- It's not feasible to check all combinations of features by fitting a model to each possible subset.
- **Voting** is a general technique where we let the features vote for all models that are compatible with it.
  - Cycle through features, cast votes for model parameters.
  - Look for model parameters that receive a lot of votes.
- Noise & clutter features will cast votes too, but typically their votes should be inconsistent with the majority of "good" features.

Kristen Grauman

### Difficulty of line fitting



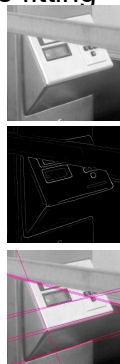
Kristen Grauman

### Hough Transform for line fitting

- Given points that belong to a line, what is the line?
- How many lines are there?
- Which points belong to which lines?
- **Hough Transform** is a voting technique that can be used to answer all of these questions.
 

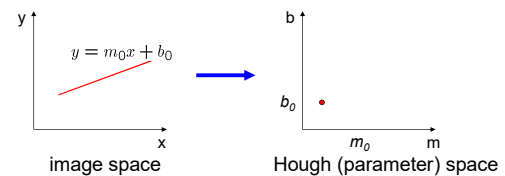
Main idea:

  1. Record vote for each possible line on which each edge point lies.
  2. Look for lines that get many votes.



Kristen Grauman

### Finding lines in an image: Hough space



#### Connection between image (x,y) and Hough (m,b) spaces

- A line in the image corresponds to a point in Hough space
- To go from image space to Hough space:
  - given a set of points (x,y), find all (m,b) such that  $y = mx + b$

Slide credit: Steve Seitz

### Finding lines in an image: Hough space

image space  $(x_0, y_0)$  → Hough (parameter) space  $b = -x_0m + y_0$

Connection between image  $(x,y)$  and Hough  $(m,b)$  spaces

- A line in the image corresponds to a point in Hough space
- To go from image space to Hough space:
  - given a set of points  $(x,y)$ , find all  $(m,b)$  such that  $y = mx + b$
- What does a point  $(x_0, y_0)$  in the image space map to?
  - Answer: the solutions of  $b = -x_0m + y_0$
  - this is a line in Hough space

Slide credit: Steve Seitz

### Finding lines in an image: Hough space

image space  $(x_0, y_0), (x_1, y_1)$  → Hough (parameter) space  $b = -x_0m + y_0, b = -x_1m + y_1$

What are the line parameters for the line that contains both  $(x_0, y_0)$  and  $(x_1, y_1)$ ?

- It is the intersection of the lines  $b = -x_0m + y_0$  and  $b = -x_1m + y_1$

### Finding lines in an image: Hough algorithm

image space → Hough (parameter) space

How can we use this to find the most likely parameters  $(m,b)$  for the most prominent line in the image space?

- Let each edge point in image space **vote** for a set of possible parameters in Hough space
- Accumulate votes in discrete set of bins; parameters with the most votes indicate line in image space.

### Voting: Generalized Hough Transform

- If we use scale, rotation, and translation invariant local features, then each feature match gives an alignment hypothesis (for scale, translation, and orientation of model in image).

Model Novel image

Adapted from Lana Lazebnik

### Voting: Generalized Hough Transform

- A hypothesis generated by a single match may be unreliable,
- So let each match **vote** for a hypothesis in Hough space

Model Novel image

### Gen Hough Transform details (Lowe's system)

- **Training phase:** For each model feature, record 2D location, scale, and orientation of model (relative to normalized feature frame)
- **Test phase:** Let each match btwn a test SIFT feature and a model feature vote in a 4D Hough space
  - Use broad bin sizes of 30 degrees for orientation, a factor of 2 for scale, and 0.25 times image size for location
  - Vote for two closest bins in each dimension
- Find all bins with at least three votes and perform geometric verification
  - Estimate least squares *affine* transformation
  - Search for additional features that agree with the alignment

David G. Lowe. "[Distinctive image features from scale-invariant keypoints.](#)" *IJCV* 60 (2), pp. 91-110, 2004.

Slide credit: Lana Lazebnik



### Example result

Background subtract for model boundaries

Objects recognized,

Recognition in spite of occlusion

[Lowe]

### Gen Hough vs RANSAC

**GHT**

- Single correspondence -> vote for all consistent parameters
- Represents uncertainty in the model parameter space
- Linear complexity in number of correspondences and number of voting cells; beyond 4D vote space impractical
- Can handle high outlier ratio

**RANSAC**

- Minimal subset of correspondences to estimate model -> count inliers
- Represents uncertainty in image space
- Must search all data points to check for inliers each iteration
- Scales better to high-d parameter spaces

Kristen Grauman

### Video Google System

1. Collect all words within query region
2. Inverted file index to find relevant frames
3. Compare word counts
4. Spatial verification

Sivic & Zisserman, ICCV 2003

- Demo online at : <http://www.robots.ox.ac.uk/~vgg/research/vgoogle/index.html>

Visual Object Recognition Tutorial

### Object retrieval with large vocabularies and fast spatial matching, Philbin et al., CVPR 2007

Query Results from 5k Flickr images (demo available for 100k set)

[Philbin CVPR'07]

### Instance recognition applications

- Snap, pick, pay

- <https://www.usatoday.com/videos/tech/2014/10/31/18261641/>

Slide credit: Kristen Grauman

### World-scale mining of objects and events from community photo collections, Quack et al., CIVR 2008

Auto-annotate by connecting to content on Wikipedia!



### Example Applications

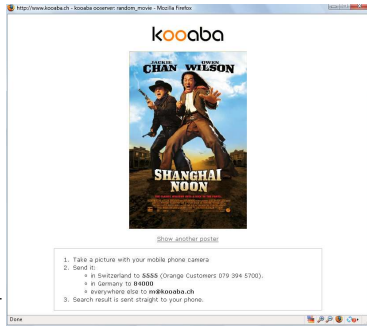


**Mobile tourist guide**

- Self-localization
- Object/building recognition
- Photo/video augmentation

B. Leibe [Quack, Leibe, Van Gool, CIVR'08]

### Web Demo: Movie Poster Recognition




50'000 movie posters indexed

Query-by-image from mobile phone available in Switzerland

[http://www.kooaba.com/en/products\\_engine.html#](http://www.kooaba.com/en/products_engine.html#)

### Google Goggles

Use pictures to search the web



Get Google Goggles

**Android (1.6+ required)**  
Download from the Android Market

**Send Goggles to Android phone**

**iPhone (iOS 4.0 required)**  
Download from the App Store

**Send Goggles to iPhone**

Text Landmarks Books Contact info Artwork Wine Logos

### Recognition via feature matching+spatial verification

**Pros:**

- Effective when we are able to find reliable features within clutter
- Great results for matching specific instances

**Cons:**

- Scaling with number of models
- Spatial verification as post-processing – not seamless, expensive for large-scale problems
- Not suited for category recognition.

Kristen Grauman

### Summary (Part 3)

- **Matching local invariant features**
  - To find specific objects and scenes.
- **Bag of words** representation: quantize feature space to make discrete set of visual words
  - Summarize image by distribution of words
  - Index individual words
- **Inverted index:** pre-compute index to enable faster search at query time
- **Recognition of instances via alignment:** matching local features followed by spatial verification
  - Robust fitting : RANSAC, GHT

Kristen Grauman

### Coming up

- Today - sign sheet if not registered / on wait list
- Read assigned papers, review 2
  - Don't be afraid of the ImageNet IJCV paper!
- Assignment 1 out now, due Sept 22