# CS 1674/2074: Local features: detection, description and matching

**PhD. Nils Murrugarra-Llerena**
nem177@pitt.edu

University of
Pittsburgh

# [Motivation] Local Features

**The "Where's Waldo?" [Game]**

The goal is to find Waldo in a crowded, complex scene.
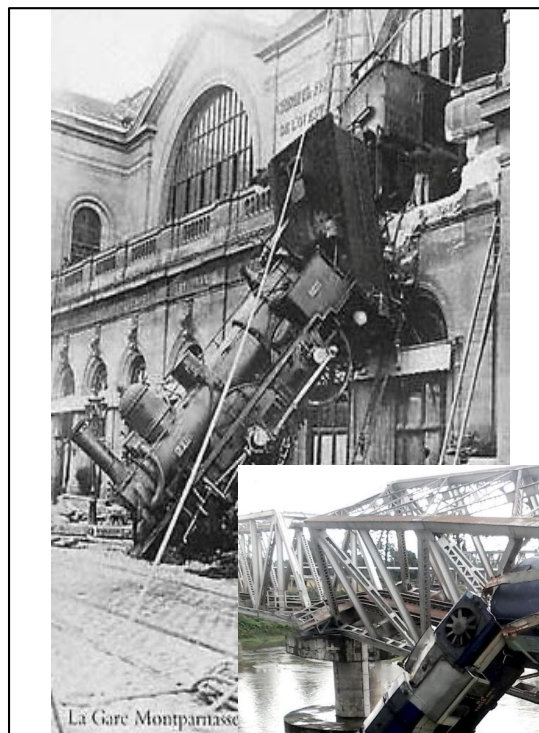
How do you search for Waldo?



What are Waldo's distinctive features?

# Plan for this lecture

- Feature detection / keypoint extraction
  - Corner detection
- Feature description (of detected features)
- Matching features across images

# An image is a set of pixels



La Gare Montparnasse, 1895

| 6 | 9 | 8 |
|---|---|---|
| 5 | 6 | 7 |
| 4 | 7 | 6 |
| 3 | 4 | 5 |
| 2 | 5 | 4 |
| 1 | 2 | 3 |

Adapted from S. Narasimhan

# Problems with pixel representation

- Not invariant to small changes
    - Translation
    - Illumination
    - etc.
- Some parts of an image are more important than others
- What do we want to represent?

Adriana Kovashka
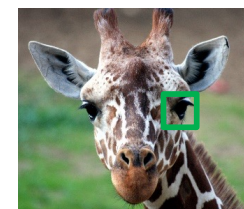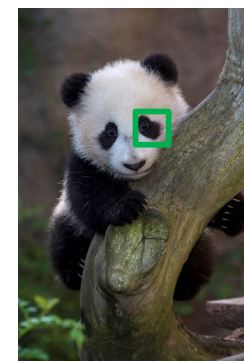
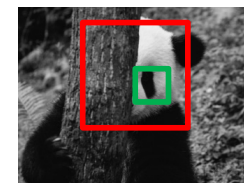# Human eye movements



Yarbus eye tracking

D. Hoiem

# Local features

- *Local* means that they only cover a small part of the image

- There will be many local features detected in an image; later we'll use those to compute a representation of the whole image

- Local features usually exploit image gradients, ignore color

- Feature ~= *vector* of gradient statistics for a window with *particular location and size*

Adriana Kovashka

# Local features: desired properties

- **Locality**
  - A feature occupies a relatively small area of the image; robust to clutter and occlusion
- **Repeatability and flexibility**
  - Robustness to expected variations: the same feature can be found in several images despite geometric/photometric transformations
  - Maximize correct matches (panda to panda)
- **Distinctiveness**
  - Each feature has a distinctive description
  - Minimize wrong matches (panda to giraffe)
- **Compactness and efficiency**
  - Many fewer features than image pixels

Adapted from K. Grauman and D. Hoiem

# Interest(ing) points

- Note: "interest points" = "keypoints", also sometimes called "features"

- Many applications
    - Recognition: which patches are likely to tell us something about the object category?

    - Image search: which points would allow us to match images between query and database?

    - 3D reconstruction: how to find correspondences across different views?

    - Tracking: which points are good to track?

Adapted from D. Hoiem

# Interest points

original

- 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

D. Hoiem

# Choosing interest points

Where would you tell
your friend to meet you?

→ Corner detection



D. Hoiem

# Choosing interest points

Where would you tell
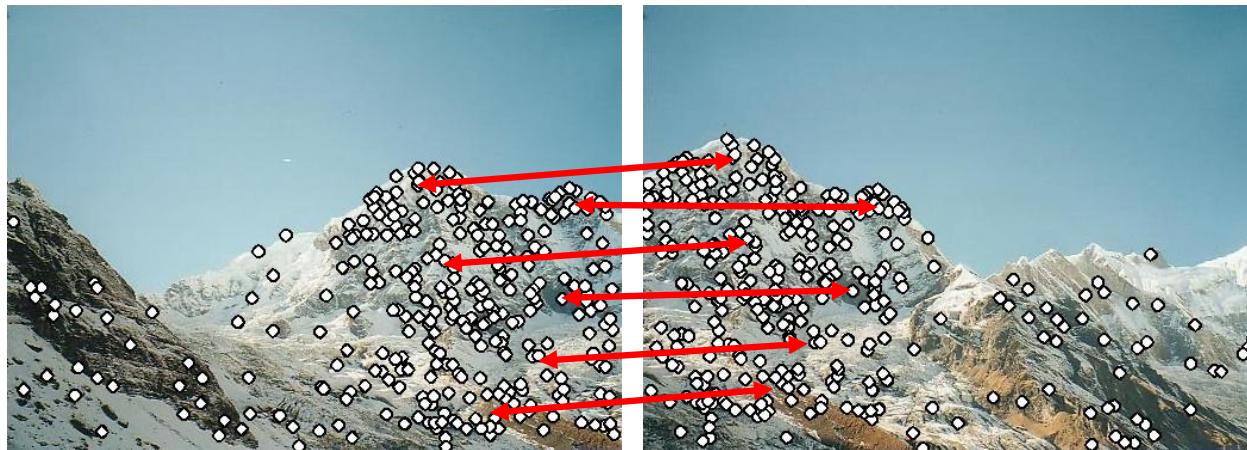your friend to meet you?

→ Blob detection



D. Hoiem

# Application: Panorama stitching

- We have two images – how do we combine them?

L. Lazebnik

# Application: Panorama stitching

- We have two images – how do we combine them?
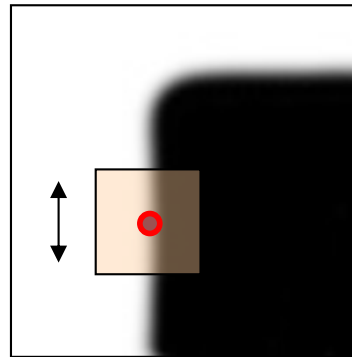


Step 1: extract features
Step 2: match features

L. Lazebnik

# Application: Panorama stitching

- We have two images – how do we combine them?



Step 1: extract features
Step 2: match features
Step 3: align images

L. Lazebnik

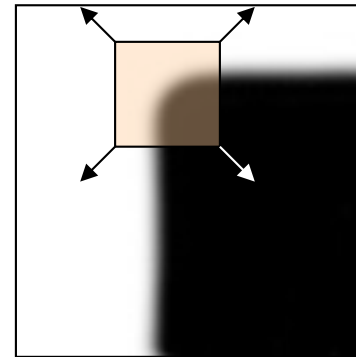# **Corners** are distinctive interest points

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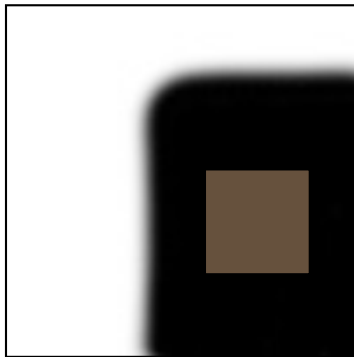- Shifting a window in *any direction* should give *a large change* in intensity

🔴 Candidate keypoint

"flat" region:
no change in
all directions

"edge":
no change along
the edge direction

"corner":
significant change
in all directions

Adapted from A. Efros, D. Frolova, D. Simakov

# **Corners** are distinctive interest points

- We should easily recognize the keypoint 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
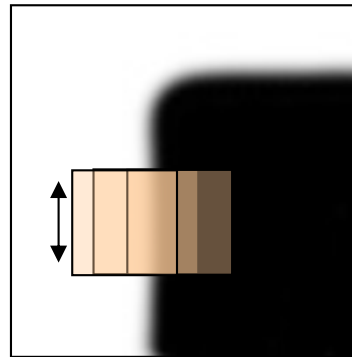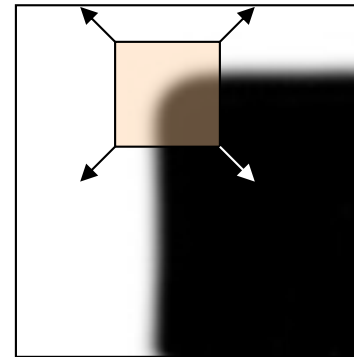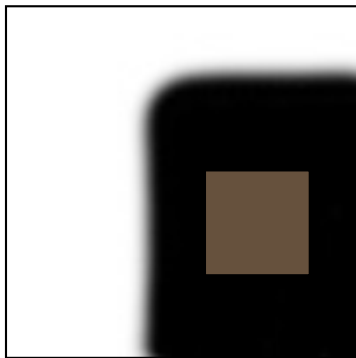
"edge":
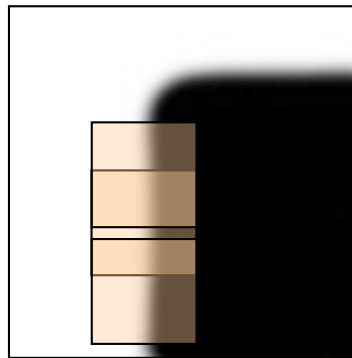no change along
the edge direction

"corner":
significant change
in all directions

Adapted from A. Efros, D. Frolova, D. Simakov
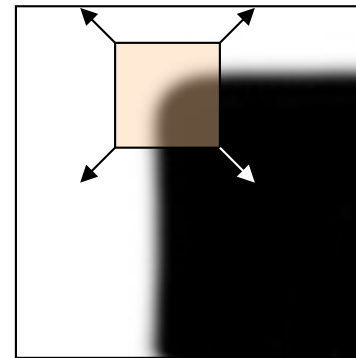
# **Corners** are distinctive interest points

- We should easily recognize the keypoint by looking through a small window
- Shifting a window in *any direction* should give *a large change* in intensity
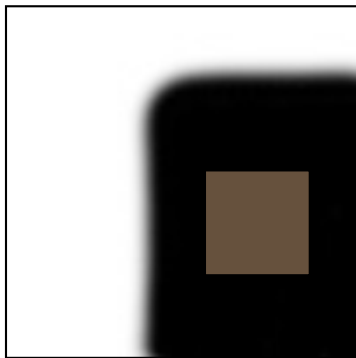
"flat" region:
no change in
all directions

"edge":
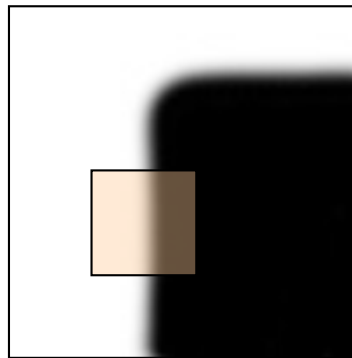no change along
the edge direction

"corner":
significant change
in all directions

Adapted from A. Efros, D. Frolova, D. Simakov
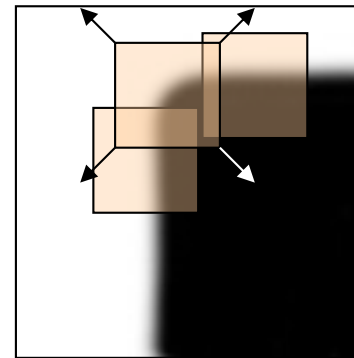
# **Corners** are distinctive interest points

- We should easily recognize the keypoint 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

Adapted from A. Efros, D. Frolova, D. Simakov

# **Corners** are distinctive interest points

- We should easily recognize the keypoint 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

Adapted from A. Efros, D. Frolova, D. Simakov

# Harris Detector: Mathematics

Window-averaged squared change of intensity induced by shifting the patch for a fixed candidate keypoint by [u,v]:
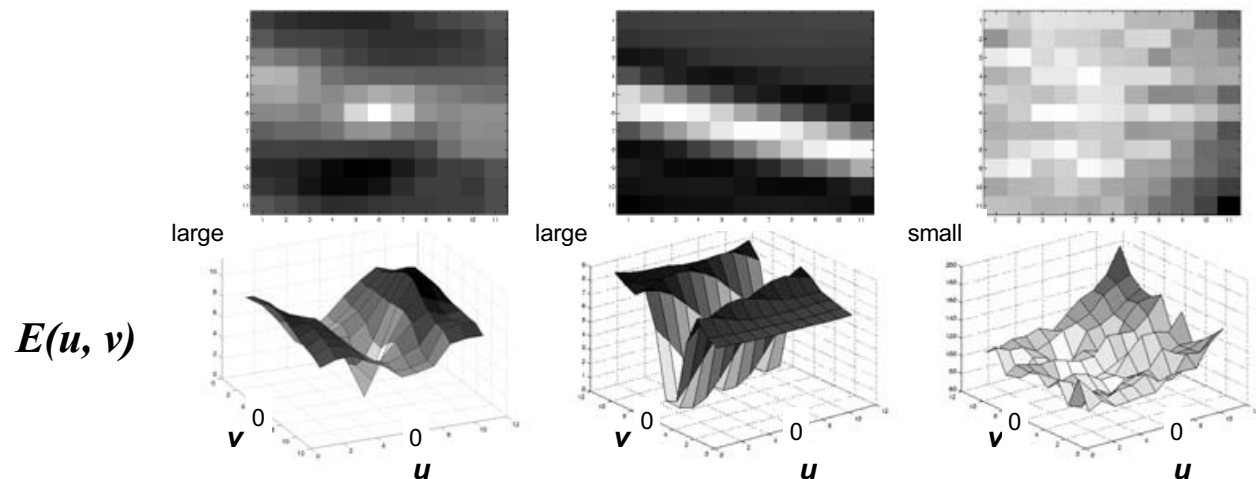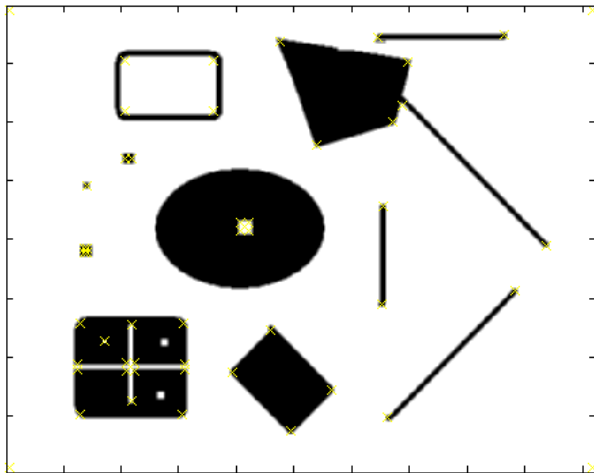
$$E(u,v) = \sum_{x,y} \left[ I(x+u, y+v) - I(x,y) \right]^2$$

Shifted intensity

Intensity

Adapted from D. Frolova, D. Simakov

# Harris Detector: Mathematics

Window-averaged squared change of intensity induced by shifting
the patch for a fixed candidate keypoint by [u,v]:

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

large          large          small

$E(u, v)$

$v$ 0        0
$u$

$v$ 0        0
$u$

$v$ 0        0
$u$

Adapted from D. Frolova, D. Simakov

# Example of Harris Application



K. Grauman

# More Harris Responses



***Effect:* A very precise corner detector.**

D. Hoiem

# Plan for this lecture

- Feature detection / keypoint extraction
    - Corner detection
- Feature description (of detected features)
- Matching features across images
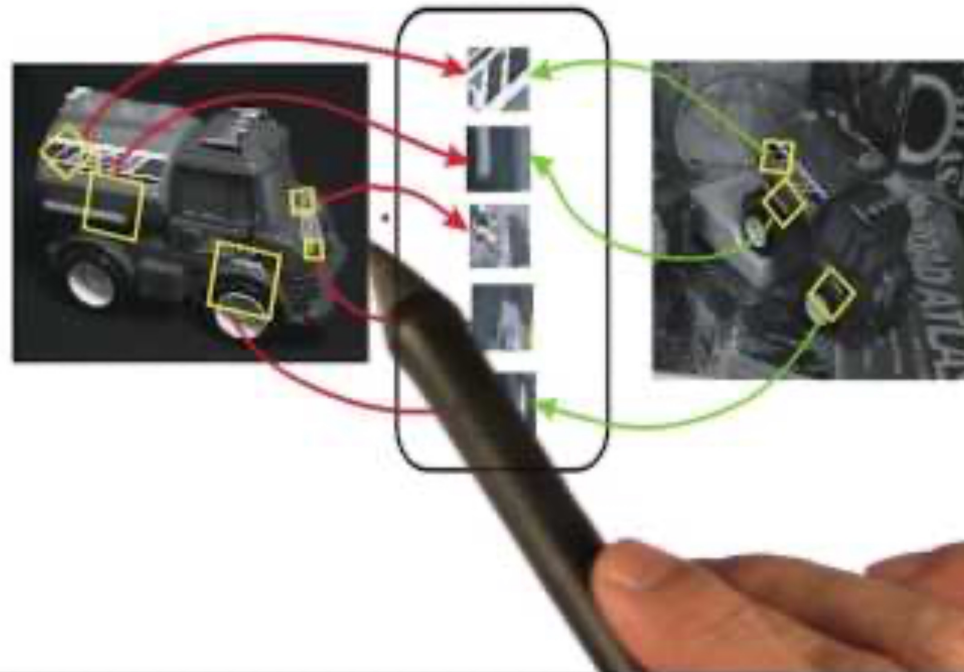
# Geometric transformations



e.g. scale, translation , rotation

K. Grauman

# Short Video SIFT Descriptor

# Scale-Invariant Feature Transform (SIFT) descriptor

Journal + conference versions: 87,527 citations (AlexNet paper has 93,821)



Image gradients → Keypoint descriptor

**Histogram of oriented gradients**

• **Captures important texture information**

• **Robust to small translations / affine deformations**
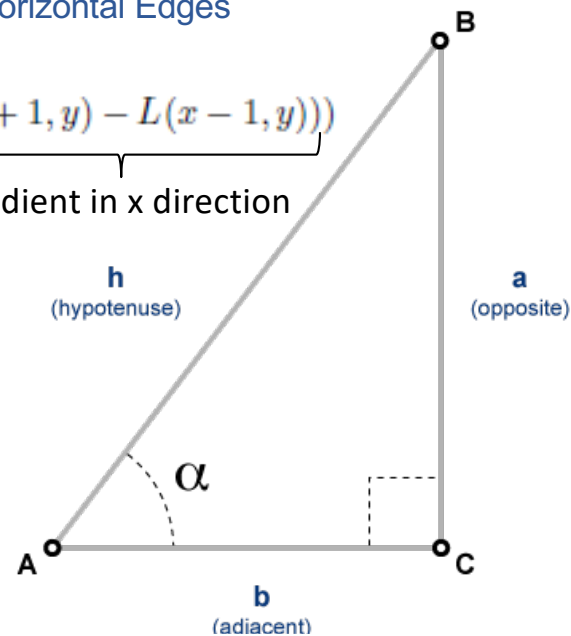
[Lowe, ICCV 1999]

K. Grauman, B. Leibe

# Computing gradients

L = the image intensity

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

gradient in x direction      gradient in y direction

Vertical Edges      Horizontal Edges

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

gradient in y direction      gradient in x direction

- $\tan(\alpha) = \dfrac{opposite\ side}{adjacent\ side}$

**h**
(hypotenuse)

**a**
(opposite)

$\alpha$

A

**b**
(adjacent)

B

C

Adriana Kovashka

# Gradients



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

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

m(x, y) = sqrt(1 + 0) = 1
Θ(x, y) = atan(0/-1) = 0

0: Black
1: White

y

(0, 0)                                                    x

Adriana Kovashka

# Gradients

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

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

m(x, y) = sqrt(0 + 1) = 1
Θ(x, y) = atan(1/0) = 90

0: Black
1: White

y

Adriana Kovashka    (0, 0)

x

# Gradients

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

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

(note length / magnitude)

m(x, y) = sqrt(1 + 1) = 1.41
Θ(x, y) = atan(-1/-1) = 45

0: Black
1: White

Adriana Kovashka   (0, 0)

# Scale Invariant Feature Transform

Basic idea:

- Take 16x16 square window around detected feature
- Compute gradient orientation for each pixel
- Create histogram over edge orientations weighted by magnitude
- That's your feature descriptor!

Image gradients

0    2π

angle histogram

Adapted from L. Zitnick, D. Lowe

# Scale Invariant Feature Transform

## Full version

- Divide the 16x16 window into a 4x4 grid of cells (2x2 case shown below)
- Quantize the gradient orientations i.e. snap each gradient to one of 8 angles
- Each gradient contributes not just 1, but magnitude(gradient) to the histogram, i.e. stronger gradients contribute more
- 16 cells * 8 orientations = 128 dimensional descriptor for each detected feature



Image gradients                    Keypoint descriptor

Adapted from L. Zitnick, D. Lowe

# Scale Invariant Feature Transform

Uniform weight (ignore magnitude)



Count

2     3     2     2

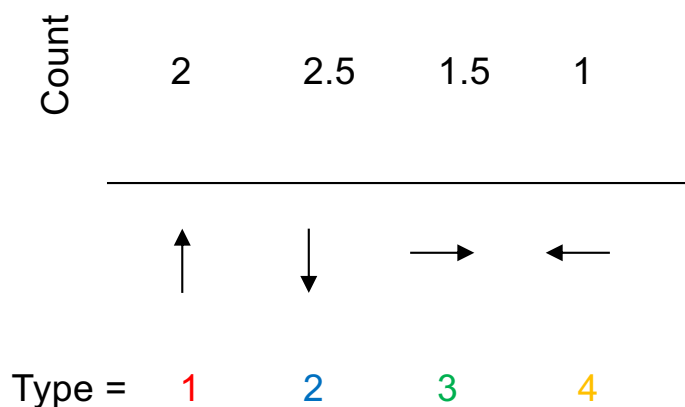Type =   1   2   3   4

Gradients           Histogram of gradients

Adriana Kovashka

# Scale Invariant Feature Transform



**Weight contribution by magnitude
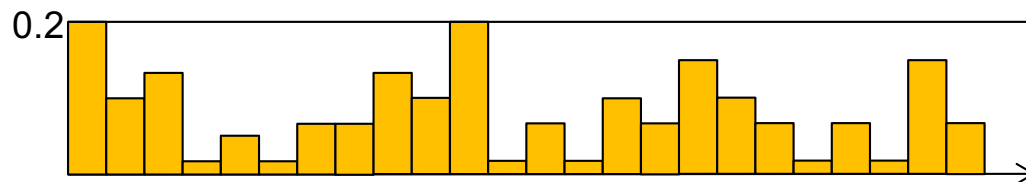(e.g. long = 1, short = 0.5)**

Gradients

Histogram of gradients

Adriana Kovashka

# Scale Invariant Feature Transform

## Full version

- Divide the 16x16 window into a 4x4 grid of cells (2x2 case shown below)
- Quantize the gradient orientations i.e. snap each gradient to one of 8 angles
- Each gradient contributes not just 1, but magnitude(gradient) to the histogram, i.e. stronger gradients contribute more
- 16 cells * 8 orientations = 128 dimensional descriptor for each detected feature
- Normalize + clip (threshold normalize to 0.2) + normalize the descriptor
- We want:

$$\sum_i d_i = 1 \quad \text{such that:} \quad d_i < 0.2$$
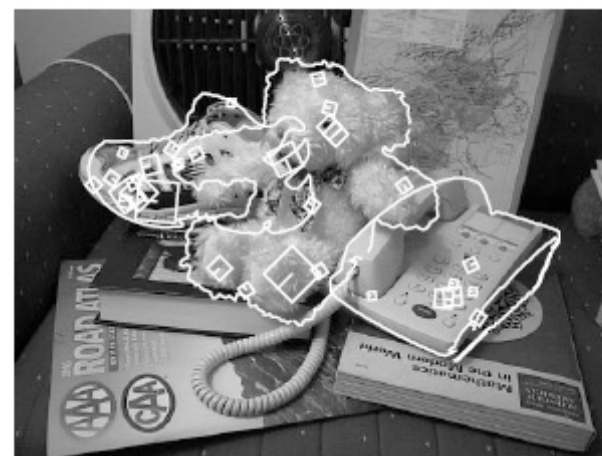
Adapted from L. Zitnick, D. Lowe

# Making descriptor rotation invariant



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

Adapted from K. Grauman, image from Matthew Brown

# SIFT is robust

- 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
- Fast and efficient—can run in real time

- Can be made to work without feature detection, resulting in "dense SIFT" (more points means robustness to occlusion)
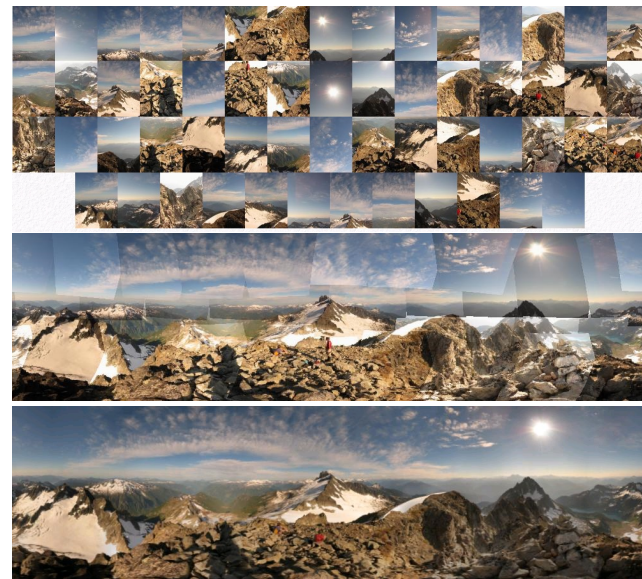
- One commonly used implementation
  - http://www.vlfeat.org/overview/sift.html

Adapted from S. Seitz

# Examples of using SIFT

# Applications of local invariant features

- Object recognition
- Indexing and retrieval
- Robot navigation
- 3D reconstruction
- Motion tracking
- Image alignment
- Panoramas and mosaics
- …



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

Adapted from K. Grauman and L. Lazebnik

# Lab 3: SIFT

Duration: 10 min

To join, go to: **ahaslides.com/OGYZC**  AhaSlides

# Please, select two images and draw SIFT matches among these images. Then, upload your result.

Get Feedback

☰  Ⓚ  🎉  |  ⧉ Group  ●

✋0  👤0/100  ☁

# Plan for this lecture

- Feature detection / keypoint extraction
  - Corner detection
  - Blob detection
- Feature description (of detected features)
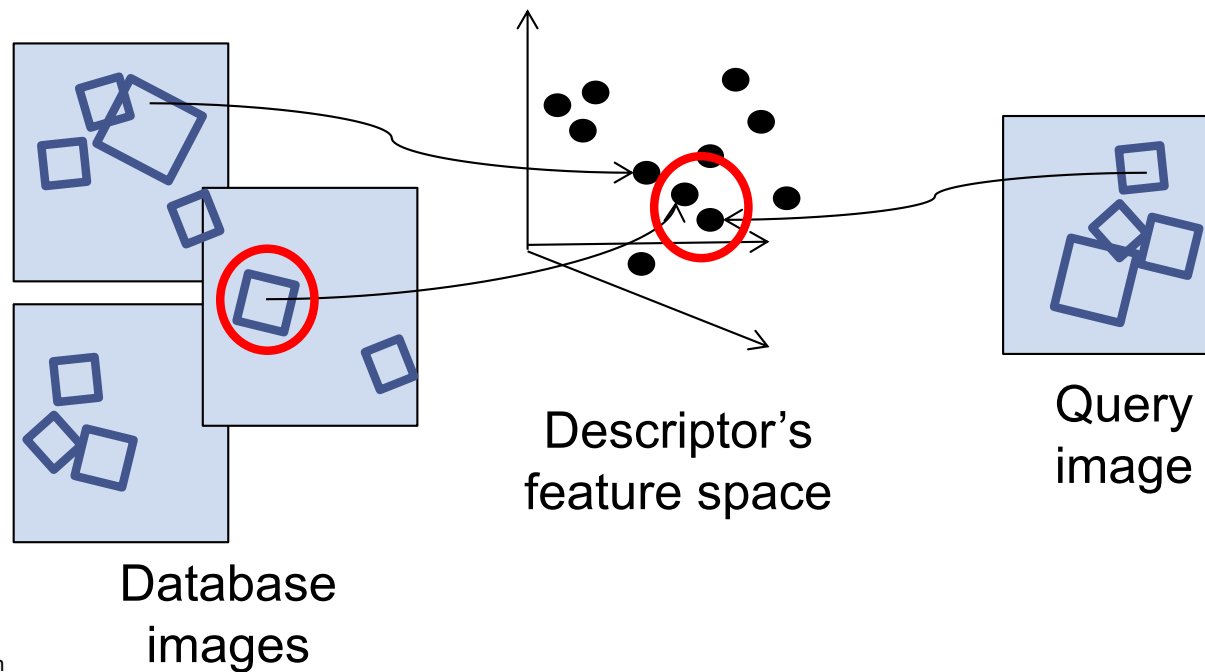- Matching features across images

# Matching Local Features Setup

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



Descriptor's
feature space

Database
images

K. Grauman

# Matching Local Features Setup

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



Database images

Descriptor's feature space

Query image

K. Grauman

# Indexing local features



- 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".

K. Grauman

# Visual Words: main idea

- Extract some local features from a number of images …



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

D. Nister, CVPR 2006

# Visual Words: main idea



D. Nister, CVPR 2006

# Visual Words: main idea



D. Nister, CVPR 2006

# Visual Words: main idea



D. Nister, CVPR 2006

Each point is a local
descriptor, e.g. SIFT
feature vector.

D. Nister, CVPR 2006

"Quantize" the space by grouping (*clustering*) the features.
Note: For now, we'll treat *clustering* as a black box.

D. Nister, CVPR 2006

# Visual Words

- Patches on the right = regions used to compute SIFT

- Each group of patches belongs to the same "visual word"



Figure from Sivic & Zisserman, ICCV 2003

Adapted from K. Grauman

# Visual Words for Indexing

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



Word #3

Query

Descriptor's feature space

Adapted from K. Grauman

- Each cluster has a center

- To determine which word to assign to new image region (e.q. query), find closest cluster center

- *To compare features:* Only compare query to others in same cluster, or just compare word IDs

- *To compare images:* see next few slides

# How to describe documents with words

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 r̶e̶a̶c̶h̶ ̶t̶h̶e̶ ̶b̶r̶a̶i̶n̶ from our eyes. For a l̶o̶n̶g̶ ̶t̶i̶m̶e̶ ̶i̶t̶ ̶w̶a̶s̶ ̶t̶h̶o̶u̶g̶h̶t̶ that the retinal im̶a̶g̶e̶ ̶w̶a̶s̶ ̶t̶r̶a̶n̶s̶m̶i̶t̶t̶e̶d̶ ̶p̶o̶int to visual̶ ̶c̶e̶n̶t̶e̶r̶s̶ ̶i̶n̶ ̶t̶h̶e̶ ̶c̶e̶r̶e̶b̶r̶a̶l̶ cortex ̶t̶h̶e̶ ̶i̶m̶a̶g̶e̶ ̶w̶o̶u̶l̶d̶ ̶b̶e̶ ̶p̶r̶o̶j̶e̶c̶t̶e̶d̶ upon ̶t̶h̶e̶ ̶s̶c̶r̶e̶e̶n̶.̶ ̶

**sensory, brain, visual, perception, retinal, cerebral cortex, eye, cell, optical nerve, image Hubel, Wiesel**

Wiesel have been able to demonstr̶a̶t̶e̶ ̶t̶h̶at the *message about the image falling̶ ̶o̶n̶ ̶the retina undergoes a step-wise analysi̶s̶ ̶i̶n̶ ̶a̶ system of nerve cells stored in columns̶.̶ ̶I̶n̶ this system each cell has its specific func̶t̶i̶o̶n̶ and is responsible for a specific detail in th̶e̶ 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% i̶n̶c̶r̶e̶a̶s̶e̶ ̶i̶n̶ exports to $750bn, comp̶a̶r̶e̶d̶ ̶w̶i̶t̶h̶ ̶a̶ ̶2̶0̶%̶ ̶r̶i̶s̶e̶ ̶i̶n̶ imports to ̶$̶6̶6̶0̶b̶n̶.̶ ̶T̶h̶e̶ ̶f̶i̶g̶u̶r̶e̶s̶ ̶a̶r̶e̶ ̶l̶i̶k̶e̶l̶y̶ to further a̶n̶n̶o̶y̶ ̶t̶h̶e̶ ̶U̶S̶,̶ ̶w̶h̶i̶c̶h̶ ̶h̶a̶s̶ ̶a̶r̶g̶ued that Ch̶i̶n̶a̶'̶s̶ ̶e̶x̶p̶o̶r̶t̶s̶ ̶a̶r̶e̶ ̶u̶n̶f̶a̶i̶r̶l̶y̶ ̶h̶e̶l̶p̶e̶d̶ ̶b̶y̶ ̶a̶ deliber̶a̶t̶e̶l̶y̶ ̶u̶n̶d̶e̶r̶v̶a̶l̶u̶e̶d̶ ̶

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**

yuan against the ̶d̶o̶l̶l̶a̶r̶ ̶l̶a̶s̶t̶ ̶J̶u̶l̶y̶ and permitted it to trade within a narro̶w̶ ̶b̶a̶n̶d̶, but the US wants the yuan to be all̶o̶w̶e̶d̶ ̶to trade freely. However, Beijing has ma̶d̶e̶ clear that it will take its time and tread̶ ̶c̶a̶r̶e̶f̶u̶l̶l̶y̶ before allowing the yuan to ris̶e̶ further in value.

# Describing images with visual words

- Summarize entire image based on its distribution (histogram) of word occurrences

- Analogous to bag of words representation commonly used for documents

Cluster 1

Cluster 2

Cluster 3

Cluster 4

Feature patches:

Adapted from K. Grauman

# Describing images with visual words

- Summarize entire image based on its distribution (histogram) of word occurrences

- Analogous to bag of words representation commonly used for documents

Feature patches:

K. Grauman

times appearing

times appearing

times appearing

Visual words

# Comparing bags of words

- Similarity of images measured as normalized scalar product between their word occurrence counts

- Can be used to rank results (nearest neighbors of query)

[1  8  1   4]          [5  1  1   0]

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

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

for vocabulary of *V* words

$\vec{d}_j$          $\vec{q}$

Adapted from K. Grauman

# Bags of words: pros and cons

+ Flexible to geometry / deformations / viewpoint

+ Compact summary of image content

- Basic model ignores geometry – verify afterwards

- What is the optimal vocabulary size?

- Background and foreground mixed when bag covers whole image

Adapted from K. Grauman

# Summary: Inverted file index and bags of words similarity

Offline:

- Extract features in database images, cluster them to find words = cluster centers, make index

Online (during search):

1. Extract words in query (extract features and map each to closest cluster center)
2. Use inverted file index to find database images relevant to query
3. Rank database images by comparing word counts of query and database image

Adapted from K. Grauman

# Summary

- Keypoint detection: repeatable and distinctive
  - Corners, blobs, stable regions
  - Laplacian of Gaussian, automatic scale selection



- Descriptors: robust and selective
  - Histograms for robustness to small shifts and translations (SIFT descriptor)

- Matching: cluster and index
  - Compare images through their feature distribution



Image gradients        Keypoint descriptor

Adapted from D. Hoiem, K. Grauman