# Indexing with local features, Bag of words models

Thursday, Oct 29

Kristen Grauman

UT-Austin

# Last time

- Interest point detection
  - Harris corner detector
  - Laplacian of Gaussian, automatic scale selection

# Local features: main components
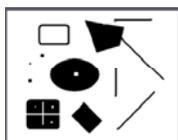
1) **Detection**: Identify the interest points



2) **Description**:Extract vector feature descriptor surrounding each interest point.

3) **Matching**: Determine correspondence between descriptors in two views

---

**Corners** as distinctive interest points

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

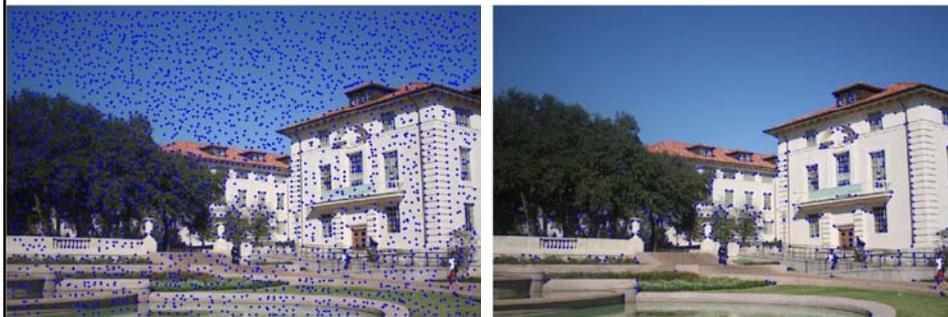2 x 2 matrix of image derivatives (averaged in neighborhood of a point).



Notation: $\quad I_x \Leftrightarrow \dfrac{\partial I}{\partial x} \qquad I_y \Leftrightarrow \dfrac{\partial I}{\partial y} \qquad I_x I_y \Leftrightarrow \dfrac{\partial I}{\partial x}\dfrac{\partial I}{\partial y}$
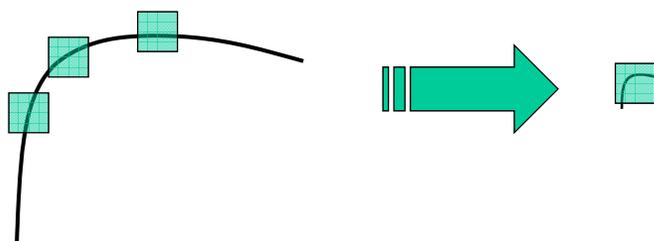
# Harris corners example



Any local max in 3 x 3 window
from the R map



Only local maxes exceeding
average R (thresholded)

---

## Properties of the Harris corner detector
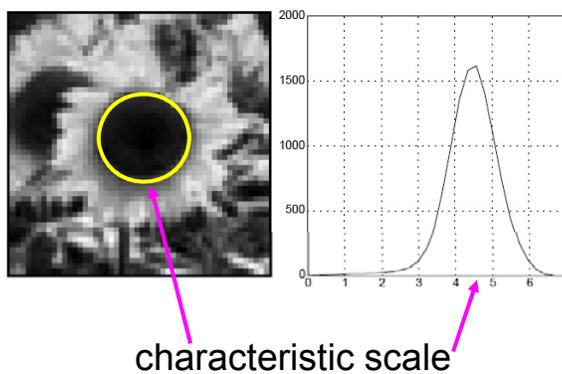
Rotation invariant?    Yes

Scale invariant?        No



All points will be
classified as edges

Corner !

# Automatic scale selection

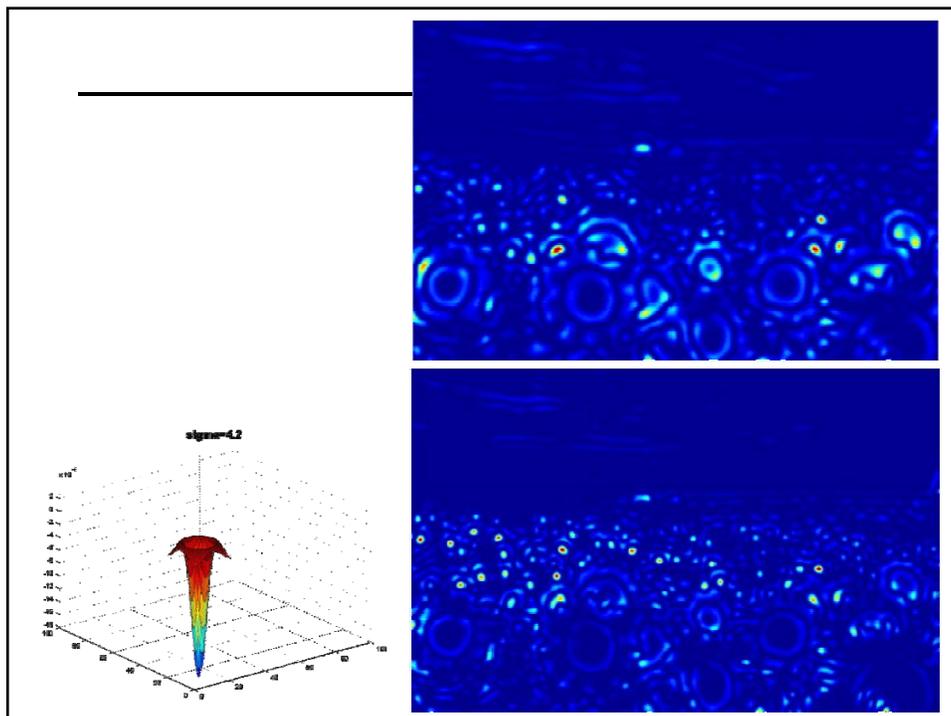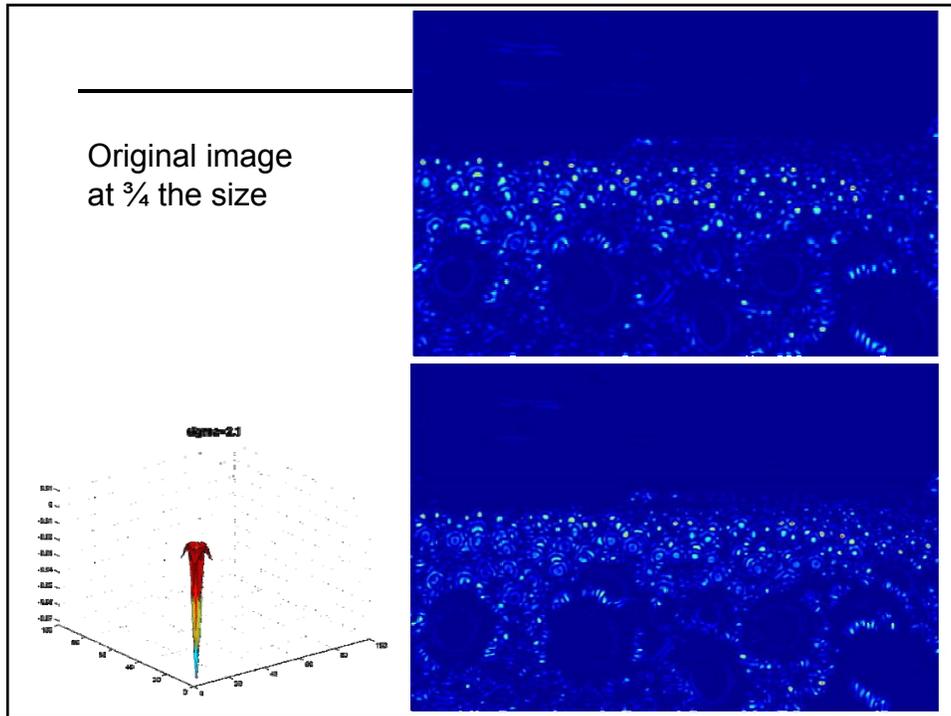We define the *characteristic scale* as the scale that produces peak of Laplacian response
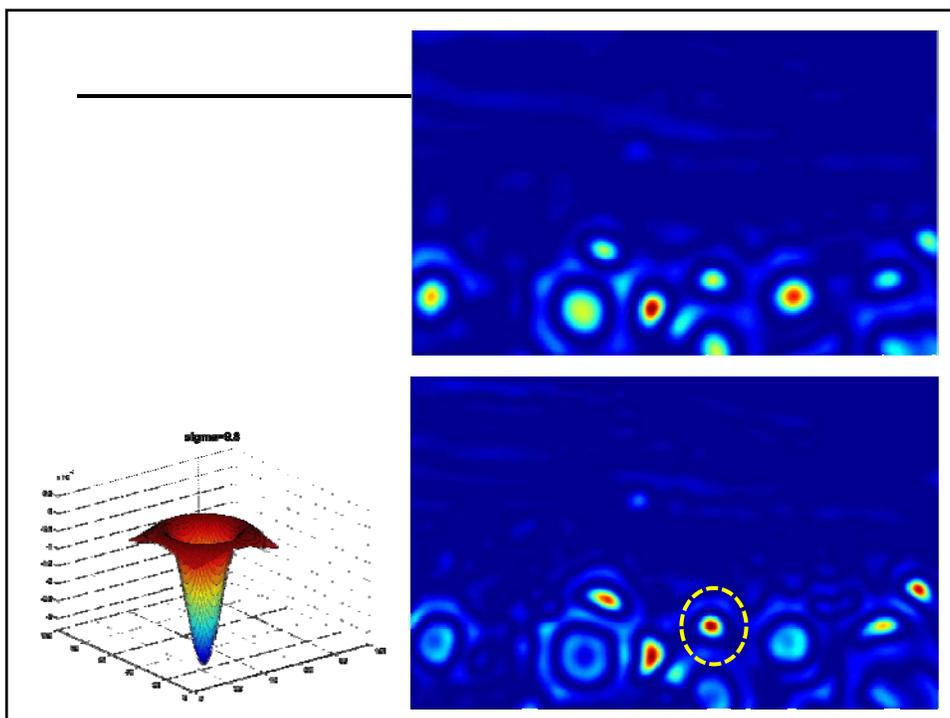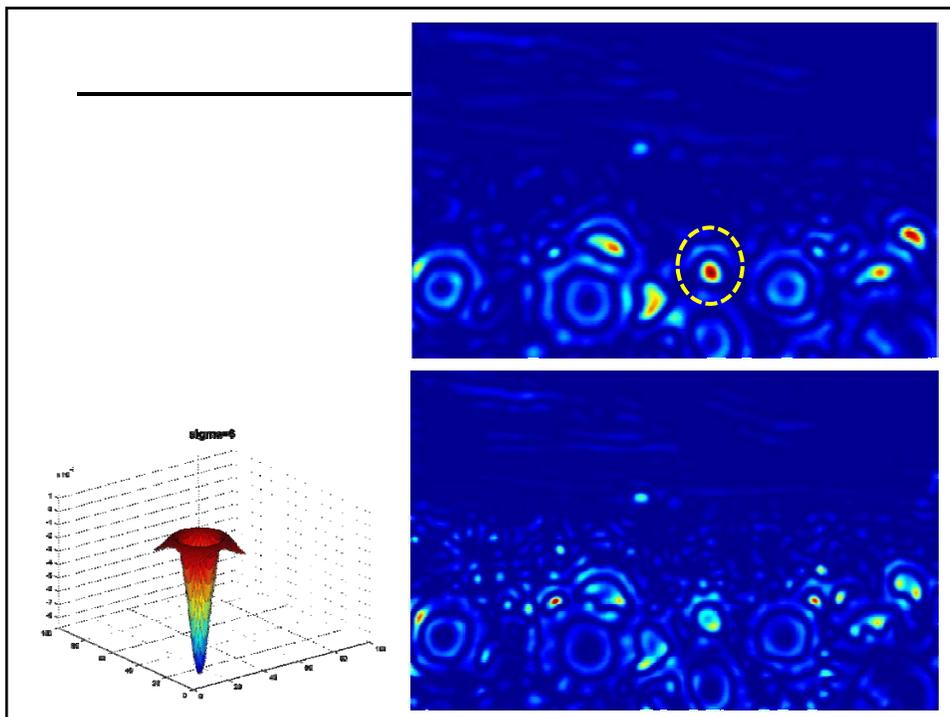


characteristic scale

# Example



Original image
at ¾ the size

Original image
at ¾ the size

sigma=2.1

sigma=4.2

## Scale invariant interest points

Interest points are local maxima in both position and scale.

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

σ5
σ4
σ3
σ2
σ1

scale

⇒ **List of**
**(x, y,** σ**)**

Squared filter
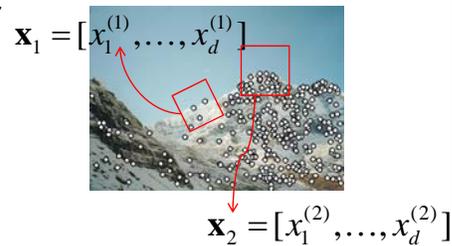response maps

# Today

- Matching local features
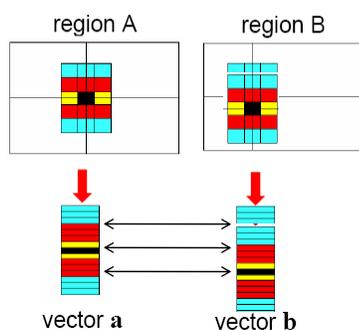- Indexing features
- Bag of words model

# Local features: main components

1) Detection: Identify the interest points

2) Description: Extract vector feature descriptor surrounding each interest point.

$$\mathbf{x}_1 = [x_1^{(1)}, \ldots, x_d^{(1)}]$$

$$\mathbf{x}_2 = [x_1^{(2)}, \ldots, x_d^{(2)}]$$

3) Matching: Determine correspondence between descriptors in two views

# Raw patches as local descriptors

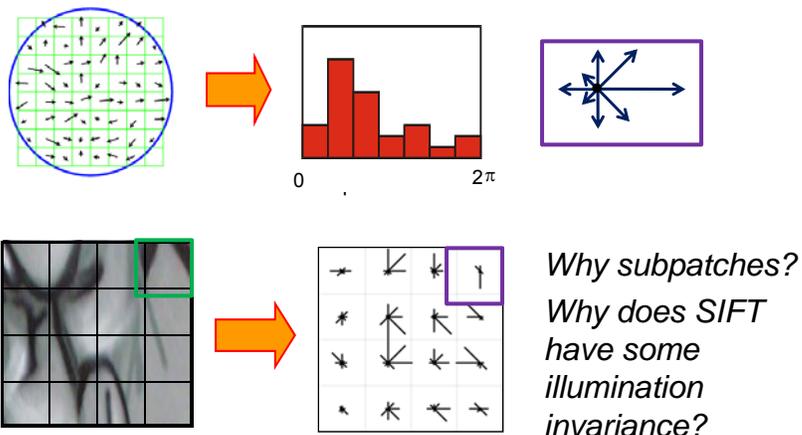region A      region B

vector **a**      vector **b**

The simplest way to describe the neighborhood around an interest point is to write down the list of intensities to form a feature vector.
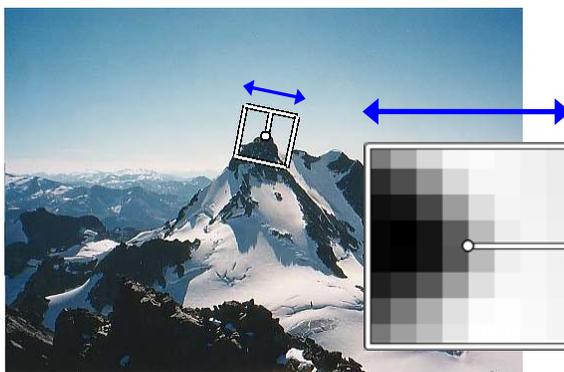
But this is very sensitive to even small shifts, rotations.

# SIFT descriptor
# [Lowe 2004]

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



0     2π

*Why subpatches?*
*Why does SIFT have some illumination invariance?*

# Making the 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
    - http://people.csail.mit.edu/albert/ladypack/wiki/index.php/Known_implementations_of_SIFT



Steve Seitz

# 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

# Matching local features



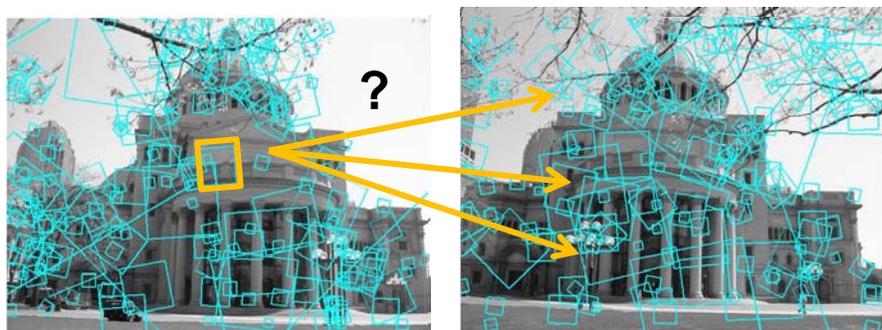# Matching local features

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)

# Matching local features



Image 1                    Image 2
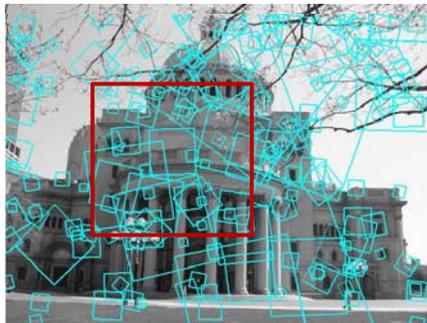
In stereo case, may constrain by proximity if we make assumptions on max disparities.

# Ambiguous matches



Image 1                    Image 2
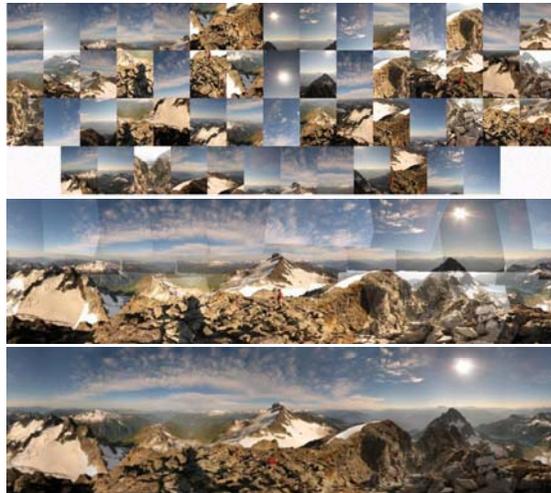
At what SSD value do we have a good match?

To add robustness to matching, can consider **ratio** : distance to best match / distance to second best match

If high, could be ambiguous match.

# 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]



# Recognition

Schmid and Mohr 1997
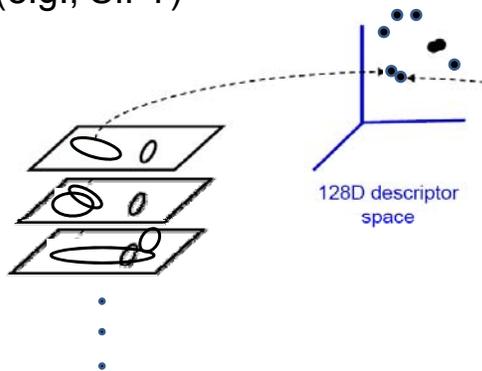
Sivic and Zisserman, 2003

Rothganger et al. 2003

Lowe 2002

# Today

- Matching local features
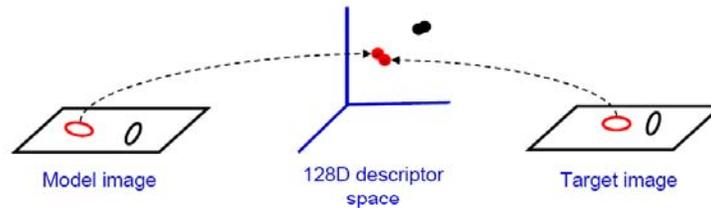- Indexing features
- Bag of words model

# Indexing local features

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



128D descriptor space
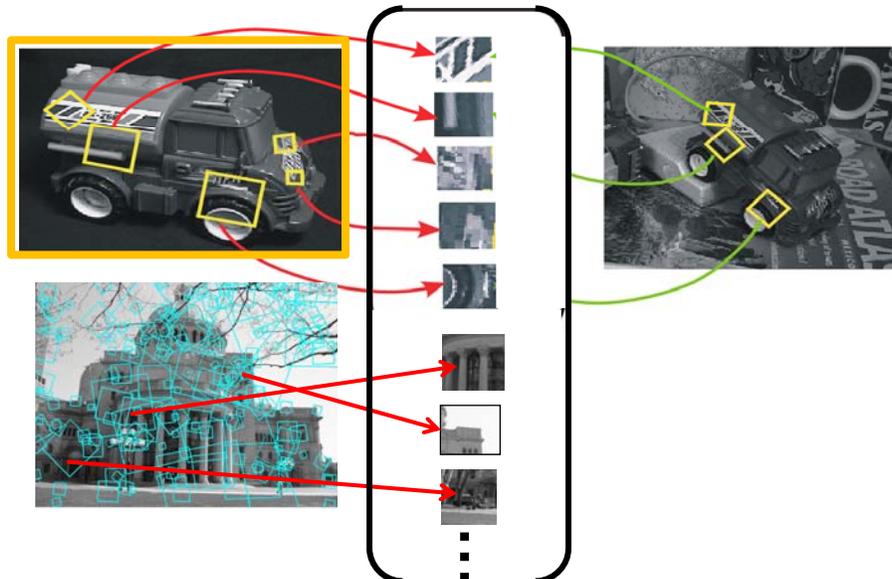
# Indexing local features

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



- This is of interest not only for 3d reconstruction, but also for retrieving images of similar objects.

Figure credit: A. Zisserman

# Indexing local features

# Indexing local features

- With potentially thousands of features per image, and hundreds to millions of images to search, how to efficiently find those that are relevant to a new image?

# 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".

# Text retrieval vs. image search

- What makes the problems similar, different?

# 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

# Visual words: main idea



# Visual words: main idea

# Visual words: main idea



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

# Visual words

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

• Quantize via clustering, let cluster centers be the prototype "words"

Descriptor space

# Visual words

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

• Determine which word to assign to each new image region by finding the closest cluster center.

Descriptor space

# Visual words

• Example: each group of patches belongs to the same visual word

Figure from Sivic & Zisserman, ICCV 2003

# Visual words and textons

- First explored for texture and material representations

- *Texton* = cluster center of filter responses over collection of images

- Describe textures and materials based on distribution of prototypical texture elements.

Leung & Malik 1999; Varma & Zisserman, 2002; Lazebnik, Schmid & Ponce, 2003;



---

# Recall: Texture representation example

Windows with primarily horizontal edges

Both

Dimension 2 (mean d/dy value)

Dimension 1 (mean d/dx value)

Windows with small gradient in both directions

Windows with primarily vertical edges

| | mean d/dx value | mean d/dy value |
|---|---|---|
| Win. #1 | 4 | 10 |
| Win.#2 | 18 | 7 |
| Win.#9 | 20 | 20 |
| | | |

**statistics to summarize patterns in small windows**

# Visual words

- More recently used for describing scenes and objects for the sake of indexing or classification.

Sivic & Zisserman 2003; Csurka, Bray, Dance, & Fan 2004; many others.



# Inverted file index



| 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

# Inverted file index



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

New query image

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

---

- If a local image region is a visual word, how can we summarize an image (the document)?

## 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 retinal image was visual centers in a movie image discovered know th perception more com following the to the various cortex, Hubel and Wiesel demonstrate that the *message about image falling on the retina undergoes wise analysis in a system of nerve cell stored in columns. In this system each has its specific function and is responsible a specific detail in the pattern of the retinal image.*

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

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% $750bn, compared $660bn. T annoy th China's deliber agrees yuan is governo also need demand so country. China yuan against the permitted it to trade within a narrow but the US wants the yuan to be allowed freely. However, Beijing has made it it will take its time and tread carefully be allowing the yuan to rise further in value.

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

ICCV 2005 short course, L. Fei-Fei

---

**Object** → **Bag of 'words'**



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]' ° [5 1 1 0]



$$sim(d_j, q) = \frac{\vec{d_j} \bullet \vec{q}}{|\vec{d_j}| \times |\vec{q}|}$$

$$= \frac{\sum_{i=1}^{t} w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^{t} w_{i,j}^2} \times \sqrt{\sum_{j=1}^{t} w_{i,q}^2}}$$

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

# *tf-idf* weighting

- **T**erm **f**requency – **i**nverse **d**ocument **f**requency
- Describe frame by frequency of each word within it, downweight words that appear often in the database
- (Standard weighting for text retrieval)

Number of occurrences of word i in document d

Total number of documents in database

$$t_i = \frac{n_{id}}{n_d} \log \frac{N}{n_i}$$

Number of words in document d

Number of documents word i occurs in, in whole database

---

# Bags of words for content-based image retrieval

What if query of interest is a portion of a frame?

Visually defined query

"Groundhog Day" [Rammis, 1993]

"Find this clock"

"Find this place"

Slide from Andrew Zisserman
Sivic & Zisserman, ICCV 2003

# 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

Query region

Retrieved frames

K. Grauman, B. Leibe

59

*Visual Object Recognition Tutorial*

---

- Collecting words within a query region

Query region:
pull out only the SIFT descriptors whose positions are within the polygon

60

raw nn 1sim=0.56697    raw nn 2sim=0.56163    raw nn 5sim=0.54917



raw nn 1sim=0.67818    raw nn 2sim=0.66144    raw nn 3sim=0.66023    raw nn 4sim=0.65774    raw nn 5sim=0.6546

# Bag of words and spatial info

- A bag of words is an orderless representation: throwing out spatial relationships between features

- Middle ground:
  - Visual "phrases" : frequently co-occurring words
  - Semi-local features : describe configuration, neighborhood
  - Let position be part of each feature
  - Count bags of words only within sub-grids of an image
  - After matching, verify spatial consistency (e.g., look at neighbors – are they the same too?)

# Visual vocabulary formation

Issues:
- Sampling strategy: where to extract features?

# Sampling strategies



Specific object                                    Category

# Sampling strategies



Sparse, at interest points



Dense, uniformly



Randomly



Multiple interest operators

- To find specific, textured objects, sparse sampling from interest points often more reliable.
- Multiple complementary interest operators offer more image coverage.
- For object categorization, dense sampling offers better coverage.

Image credits: F-F. Li, E. Nowak, J. Sivic

# Visual vocabulary formation

Issues:
- Sampling strategy: where to extract features?
- Unsupervised vs. supervised
- What corpus provides features (universal vocabulary?)
- Vocabulary size, number of words
- Clustering / quantization algorithm

## Vocabulary Trees: hierarchical clustering for large vocabularies

- Tree construction:



[Nister & Stewenius, CVPR'06]

Visual Object Recognition Tutorial

K. Grauman, B. Leibe

Slide credit: David Nister

## Vocabulary Tree

- Training: Filling the tree



[Nister & Stewenius, CVPR'06]

Visual Object Recognition Tutorial

K. Grauman, B. Leibe

Slide credit: David Nister

# Vocabulary Tree

- Training: Filling the tree



[Nister & Stewenius, CVPR'06]

K. Grauman, B. Leibe

Slide credit: David Nister

Visual Object Recognition Tutorial

# Vocabulary Tree

- Training: Filling the tree



[Nister & Stewenius, CVPR'06]

K. Grauman, B. Leibe

Slide credit: David Nister

Visual Object Recognition Tutorial

# Vocabulary Tree

- Training: Filling the tree



[Nister & Stewenius, CVPR'06]

K. Grauman, B. Leibe

Slide credit: David Nister

Visual Object Recognition Tutorial

# Vocabulary Tree

- Training: Filling the tree



[Nister & Stewenius, CVPR'06]

74

K. Grauman, B. Leibe

Slide credit: David Nister

Visual Object Recognition Tutorial

What is the computational advantage of the
hierarchical representation bag of words, vs.
a flat vocabulary?

## Vocabulary Tree

- Recognition

RANSAC
verification

[Nister & Stewenius, CVPR'06]

76

K. Grauman, B. Leibe          Slide credit: David Nister

# Bags of words: pros and cons

+ flexible to geometry / deformations / viewpoint
+ compact summary of image content
+ provides vector representation for sets
+ has yielded good recognition results in practice

- basic model ignores geometry – must verify afterwards, or encode via features
- background and foreground mixed when bag covers whole image
- interest points or sampling: no guarantee to capture object-level parts
- optimal vocabulary formation remains unclear

# Summary

- **Local invariant features**: distinctive matches possible in spite of significant view change, useful not only to provide matches for multi-view geometry, but also to find objects and scenes.

- To find **correspondences** among detected features, measure distance between descriptors, and look for most similar patches.

- **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