

Recognizing object instances

Monday, April 4
Prof. Kristen Grauman
UT-Austin

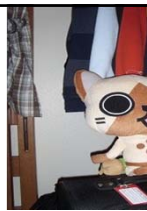
Some pset 3 results!



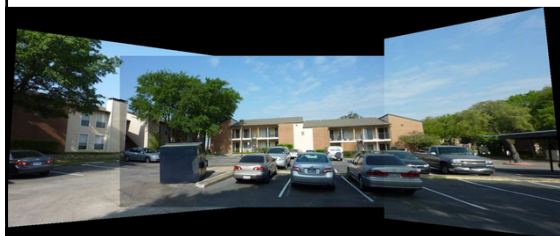
Brian Bates



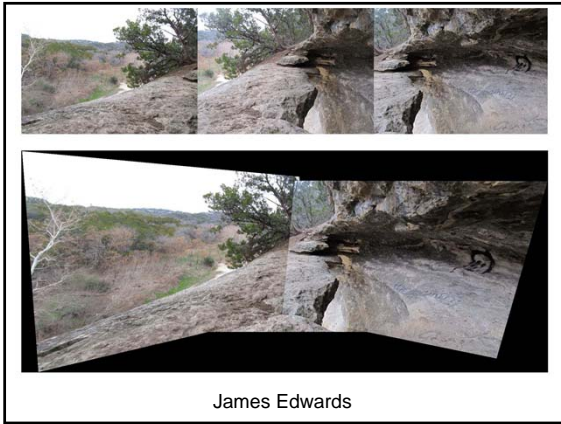
Christopher Tosh



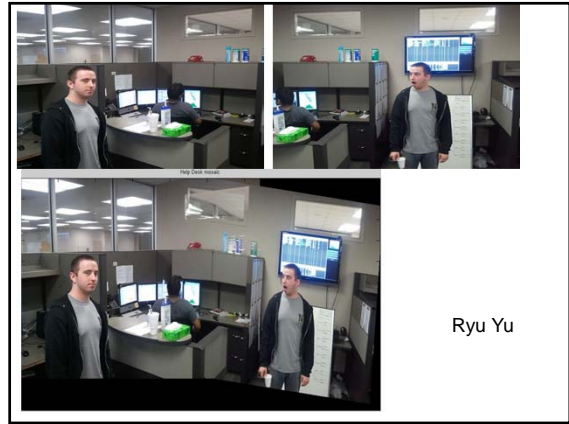
Brian Nguyen



Che-Chun Su



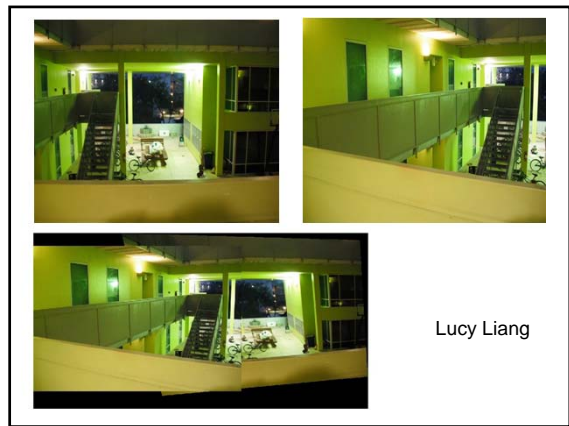
James Edwards



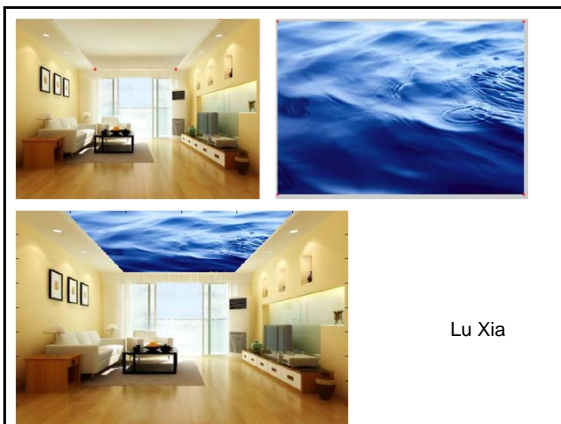
Ryu Yu



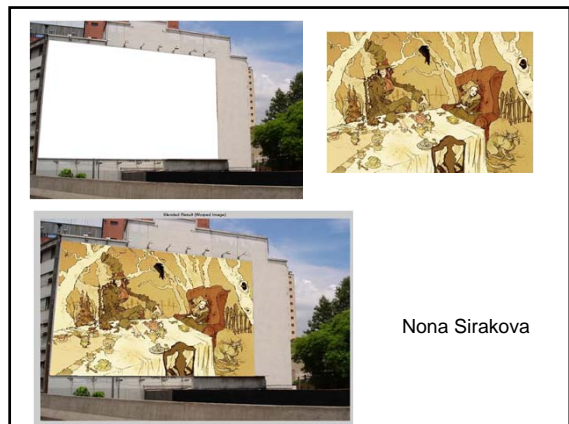
Kevin Harkness



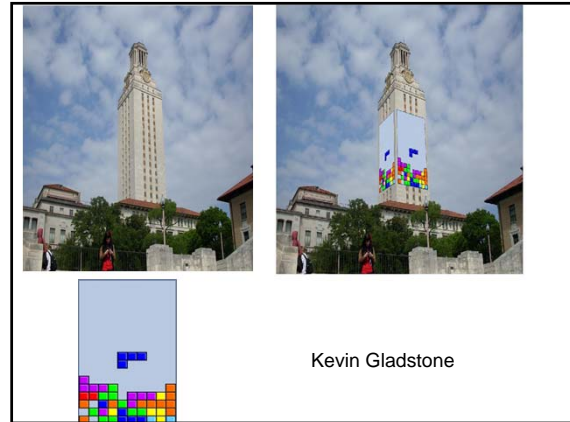
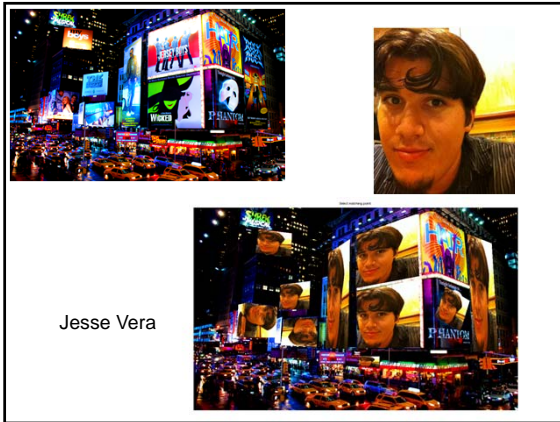
Lucy Liang



Lu Xia



Nona Sirakova



Today: instance recognition

- Motivation – visual search
- **Visual words**
 - quantization, index, bags of words
- **Spatial verification**
 - affine; RANSAC, Hough
- **Other text retrieval tools**
 - tf-idf, query expansion
- **Example applications**

Multi-view matching

Matching two given views for depth

vs

Search for a matching view for recognition

Kristen Grauman



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

Descriptor's feature space

Query image

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

Database images

Descriptor's feature space

Word #2

- 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.

Kristen Grauman

Visual words

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

Figure from Sivic & Zisserman, ICCV 2003

Kristen Grauman

Visual vocabulary formation

Issues:

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

Kristen Grauman


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

Kristen Grauman

Inverted file index



| Word # | Image # |
|--------|---------|
| 1 | 3 |
| 2 | |
| 7 | 1, 2 |
| 8 | 5 |
| 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?
- How to score the retrieval results?

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 retinal image was considered as a movie of the world. It was discovered that the visual system knows the more complex, following the various stages of the visual pathway. Hubel and Wiesel demonstrated that the message about the image falling on the retina undergoes a coarse analysis in a system of nerve cells, stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

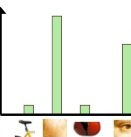
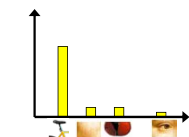
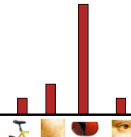
China is forecasting a trade surplus of \$90bn (€51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports of \$750bn, compared with \$560bn. The government also needs to demand so much more from the country. China's yuan against the dollar has been 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.


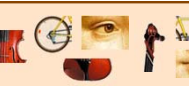

sensory, brain, visual, perception, retinal, cerebral cortex, eye, cell, optical nerve, image Hubel, Wiesel

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

ICCV 2005 short course, L. Fei-Fei

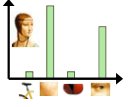

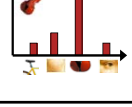









Bags of visual words

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

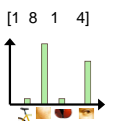







Comparing bags of words

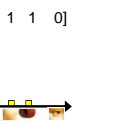

- Rank frames by normalized scalar product between their (possibly weighted) occurrence counts---nearest neighbor search for similar images.

[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

Kristen Grauman

Inverted file index and bags of words similarity

| 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?
- How to score the retrieval results?

Kristen Grauman

Vocabulary size

Results for recognition task with 6347 images

Influence on performance, sparsity

Nister & Stewenius, CVPR 2006
Kristen Grauman

Vocabulary Trees: hierarchical clustering for large vocabularies

- Tree construction:

[Nister & Stewenius, CVPR'06]
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

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]
K. Grauman, B. Leibe Slide credit: David Nister 39

Visual Object Recognition Tutorial

Vocabulary Tree

- Recognition

[Nister & Stewenius, CVPR'06]
K. Grauman, B. Leibe Slide credit: David Nister 40

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?
- How to score the retrieval results?

Kristen Grauman

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

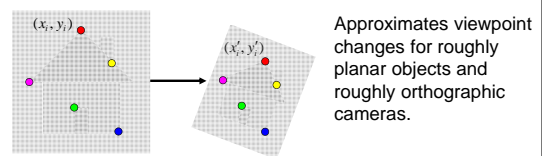
- RANSAC
 - Typically sort by BoW similarity as initial filter
 - Verify by checking support (inliers) for possible transformations
 - e.g., “success” if find a transformation with > N inlier correspondences
- Generalized Hough Transform
 - Let each matched feature cast a vote on location, scale, orientation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

Kristen Grauman

RANSAC verification



Recall: Fitting an affine transformation



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

$$\begin{bmatrix} x_i & y_i & 0 & 0 & 1 & 0 \\ 0 & 0 & x_i & y_i & 0 & 1 \\ \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ t_1 \\ t_2 \end{bmatrix} = \begin{bmatrix} \dots \\ x'_i \\ y'_i \\ \dots \end{bmatrix}$$

RANSAC verification

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

Kristen Grauman

Example Applications

Mobile tourist guide

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

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

Application: Large-Scale Retrieval

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

[Philbin CVPR'07]

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#

Google Goggles

Use pictures to search the web | Watch a video

Send Goggles to Android phone

Send Goggles to iPhone

Text

Landmarks

Books


Contact info

Artwork

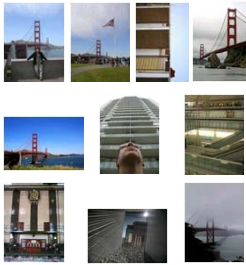
Wine

Logos

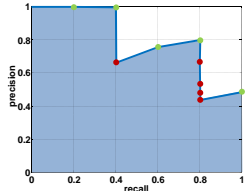
Scoring retrieval quality

Query 

Database size: 10 images
Relevant (total): 5 images

Results (ordered): 

precision = #relevant / #returned
recall = #relevant / #total relevant



Slide credit: Ondrej Chum


Spatial Verification: two basic strategies

- RANSAC
 - Typically sort by BoW similarity as initial filter
 - Verify by checking support (inliers) for possible transformations
 - e.g., "success" if find a 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: 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).

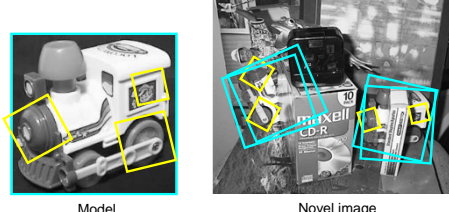


Model Novel image

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]

Recall: difficulties of voting

- Noise/clutter can lead to as many votes as true target
- Bin size for the accumulator array must be chosen carefully
- In practice, good idea to make broad bins and spread votes to nearby bins, since verification stage can prune bad vote peaks.

Gen Hough vs RANSAC

GH

- 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

What else can we borrow from text retrieval?

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 to \$750bn, compared with \$580bn in 2004. The yuan is also expected to rise to 6.85 per dollar, from 6.70 in 2004. China's government also needs to raise interest rates to curb inflation. China's government also needs to raise interest rates to curb inflation. China's government also needs to raise interest rates to curb inflation.

tf-idf weighting

- Term frequency – inverse document frequency
- Describe frame by frequency of each word within it, downweight words that appear often in the database
- (Standard weighting for text retrieval)

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

Number of occurrences of word i in document d

Number of words in document d

Total number of documents in database

Number of documents word i occurs in, in whole database

Kristen Grauman

Query expansion

Query: **golf green**

Results:

- How can the grass on the **greens** at a **golf** course be so perfect?
- For example, a skilled **golfer** expects to reach the **green** on a par-four hole in ...
- Manufactures and sells synthetic **golf** putting **greens** and mats.

Irrelevant result can cause a 'topic drift':

- Volkswagen **Golf**, 1999, **Green**, 2000cc, petrol, manual, hatchback, 94000miles, 2.0 GTI, 2 Registered Keepers, HPI Checked, Air-Conditioning, Front and Rear Parking Sensors, ABS, Alarm, Alloy

Slide credit: Ondrej Chum

Query Expansion



Slide credit: Ondrej Chum

Recognition via alignment

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

- **Matching local invariant features**
 - Useful not only to provide matches for multi-view geometry, but also to find 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

CCFP
CENTER FOR COSMOLOGY
AND PARTICLE PHYSICS

Astrometry.net

Making the Sky Searchable: Fast Geometric Hashing for Automated Astrometry

Sam Roweis, Dustin Lang & Keir Mierle
University of Toronto

David Hogg & Michael Blanton
New York University

<http://astrometry.net>
roweis@ca.toronto.edu

Example

A shot of the Great Nebula, by Jerry Lodriguss (c.2006), from astropix.com
<http://astrometry.net/gallery.html>

Example

An amateur shot of M100, by Filippo Ciferri (c.2007) from lickr.com
<http://astrometry.net/gallery.html>

Example

A beautiful image of Bode's nebula (c.2007) by Peter Bresseler, from slightfriend.de
<http://astrometry.net/gallery.html>