# Large-Scale Image Retrieval

**Jianping Fan**
**Department of Computer Science**
**UNC-Charlotte**
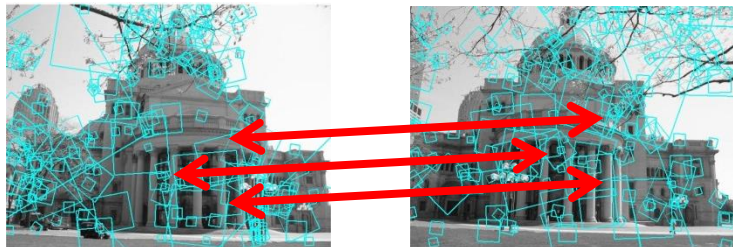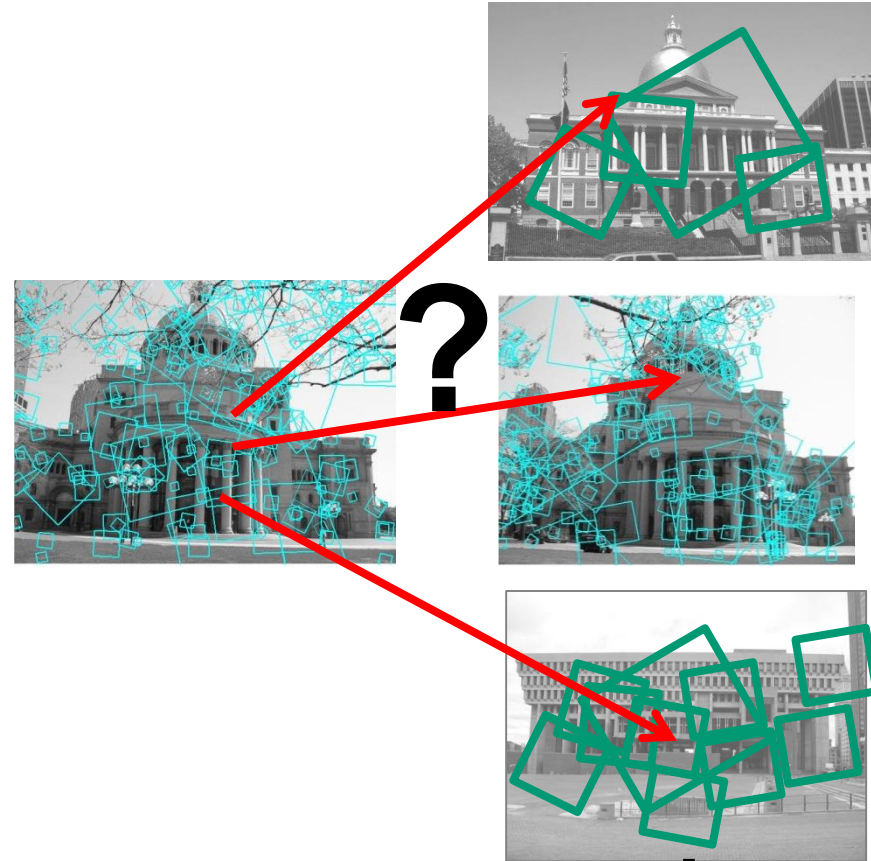
**Course Website:**
**http://webpages.uncc.edu/jfan/itcs5152.html**

# Multi-view matching



**vs**

**?**

Matching two given
views for depth

Search for a matching
view for recognition

# 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**
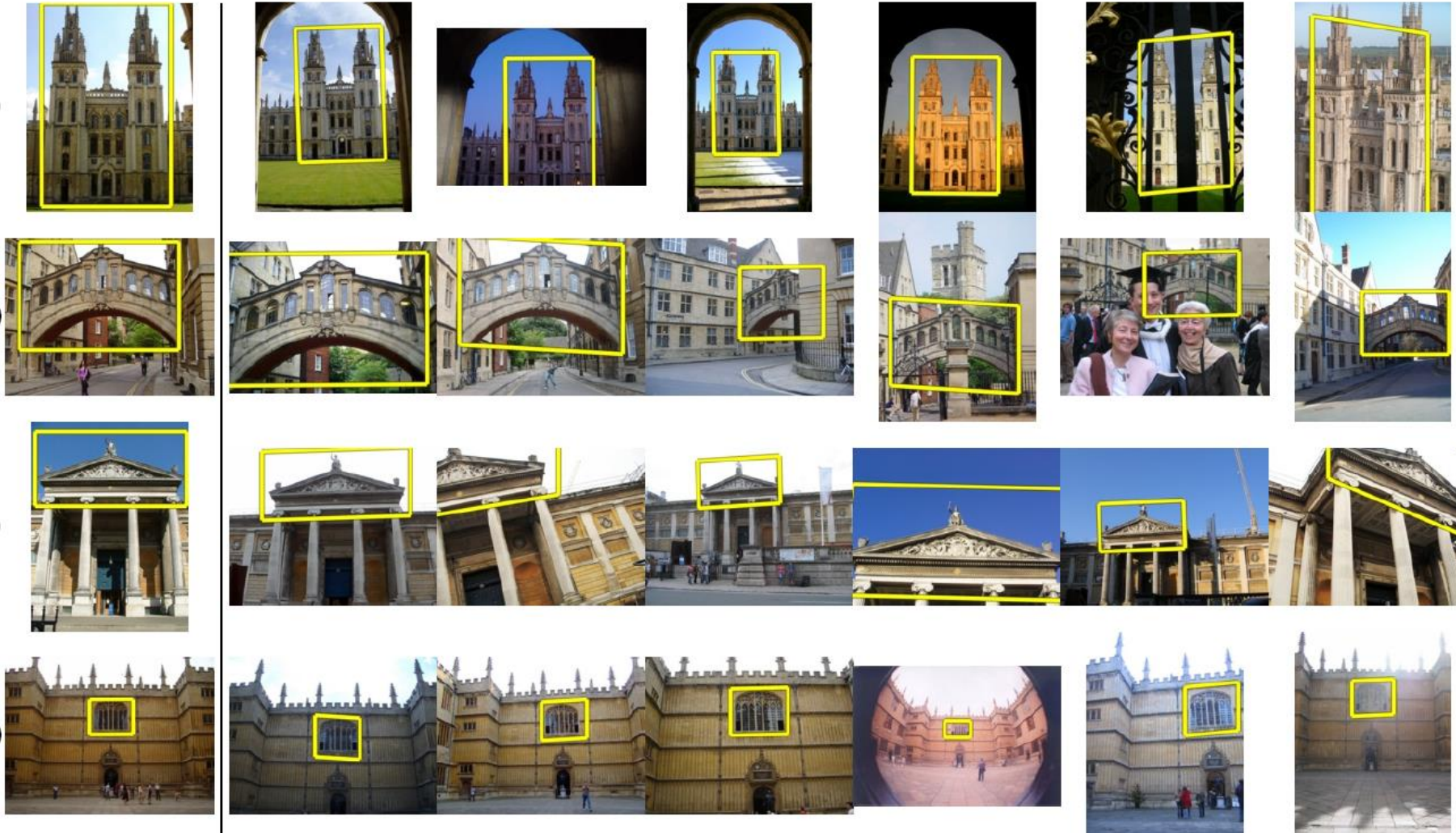
Query region

Retrieved frames

# Application: Large-Scale Retrieval

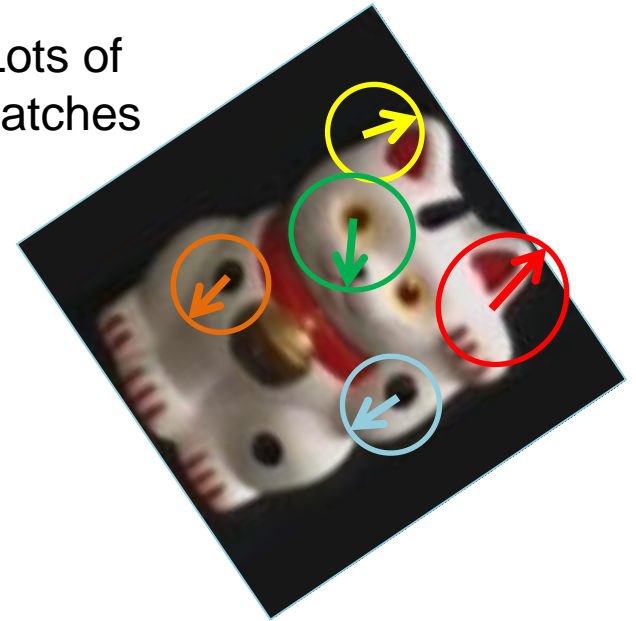Query          Results from 5k Flickr images (demo available for 100k set)

[Philbin CVPR'07]

# Simple idea

See how many keypoints are close to keypoints in each other image
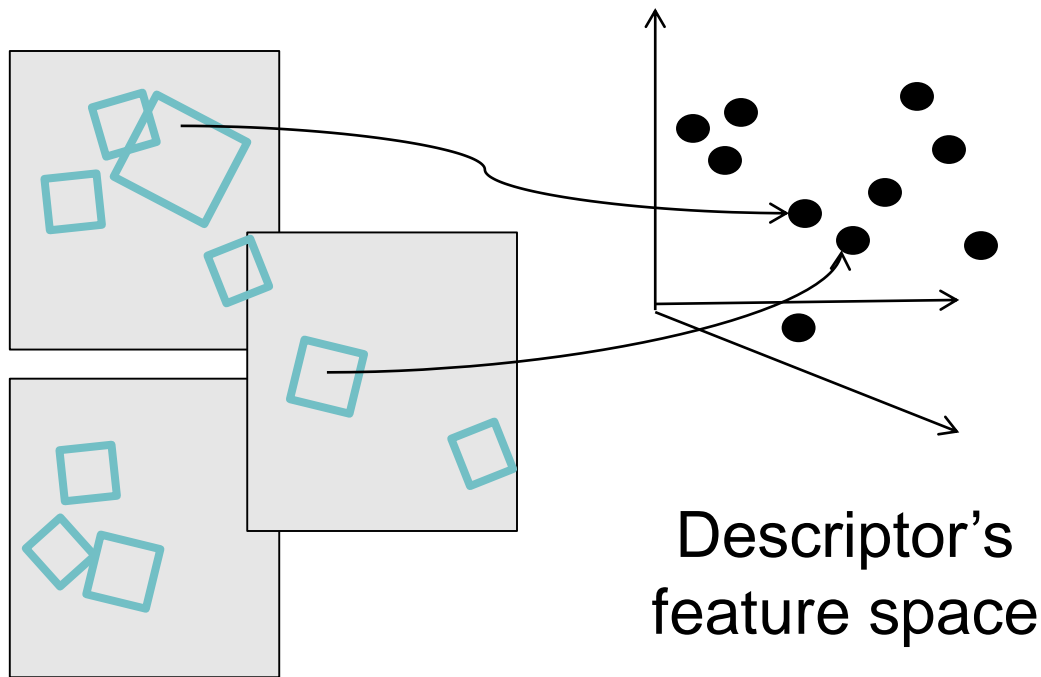
Lots of Matches



Few or No Matches
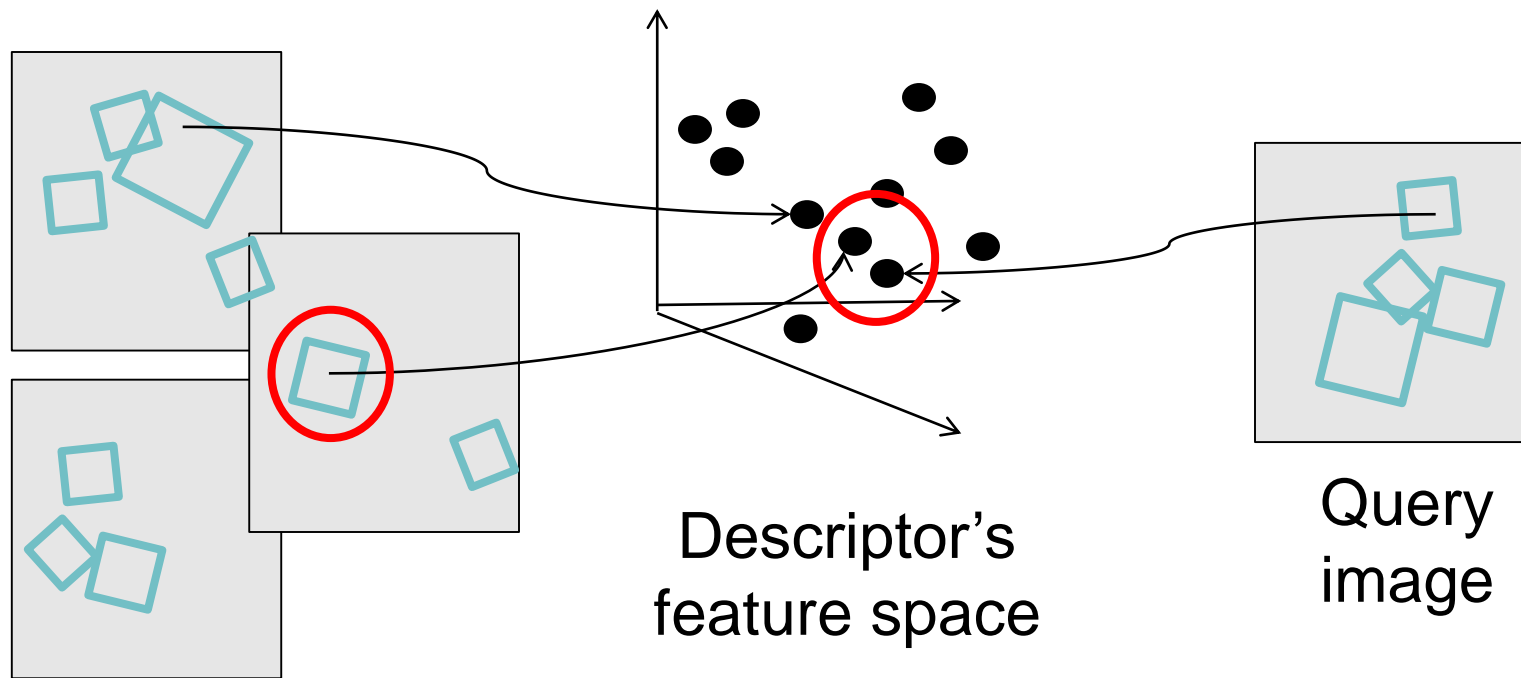


But this will be really, really slow!

# 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

# 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!*
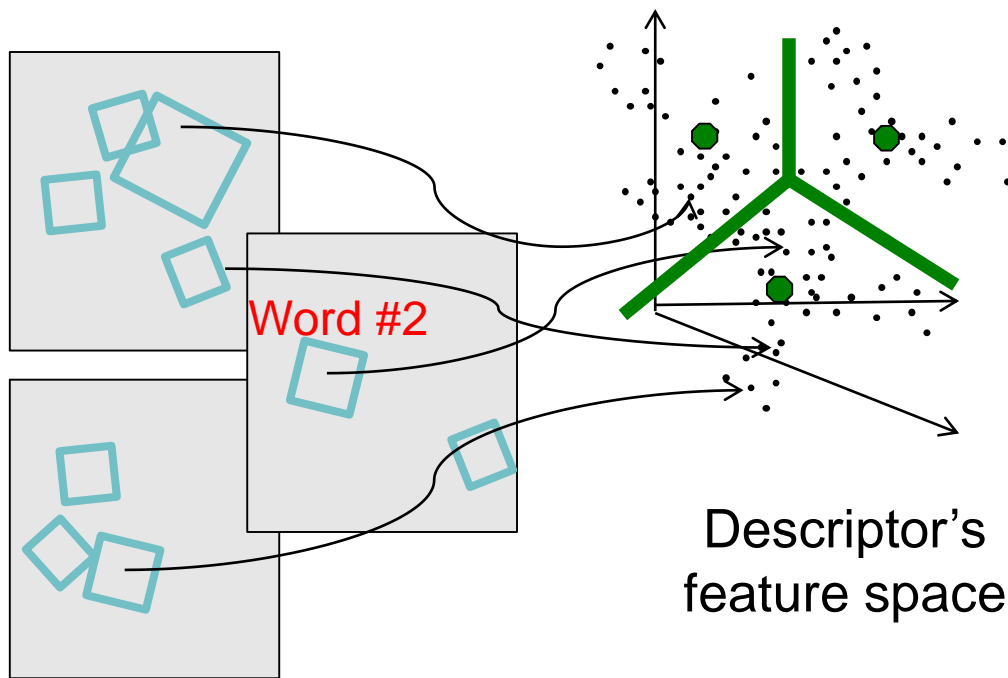
Kristen Grauman

# 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

- 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.

Word #2

Descriptor's feature space

Kristen Grauman

# Visual words

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

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

# 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.

[See Nowak, Jurie & Triggs, ECCV 2006]

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

K. Grauman, B. Leibe

# Inverted file index



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

Image #1

Image #2

Image #3

Database images

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

# 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

# Inverted file index

- Key requirement for inverted file index to be efficient: sparsity

- If most pages/images contain most words then you're no better off than exhaustive search.

  - Exhaustive search would mean comparing the word distribution of a query versus every page.
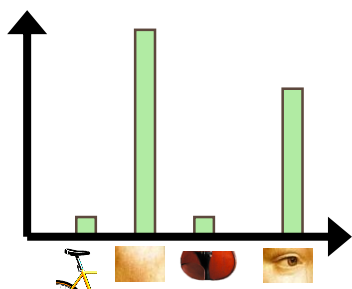
# Instance recognition: remaining issues

- 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?
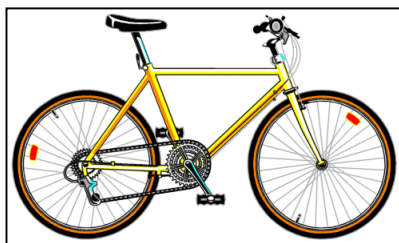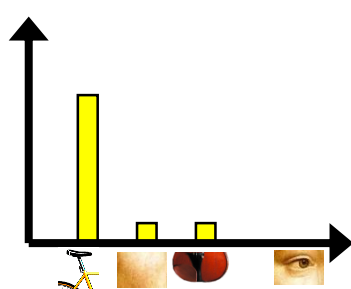
Kristen Grauman

# 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]

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

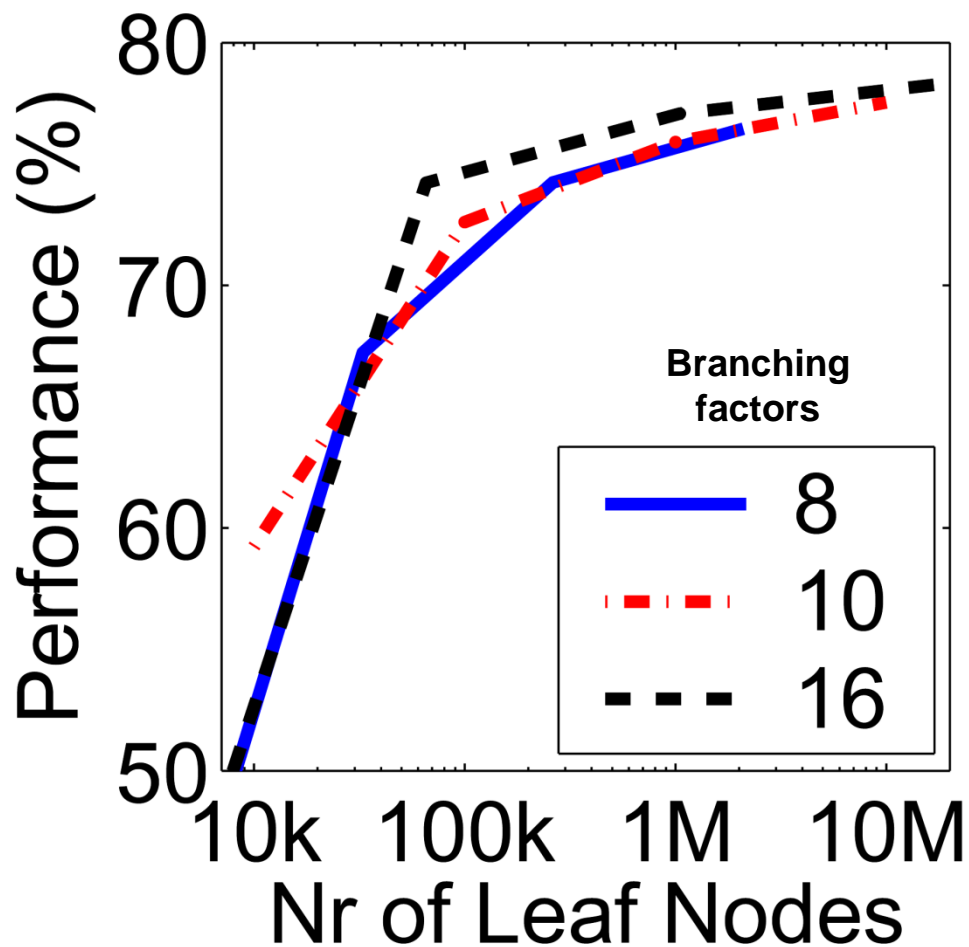$$= \frac{\sum_{i=1}^{V} d_j(i) * q(i)}{\sqrt{\sum_{i=1}^{V} d_j(i)^2} * \sqrt{\sum_{i=1}^{V} q(i}}$$

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

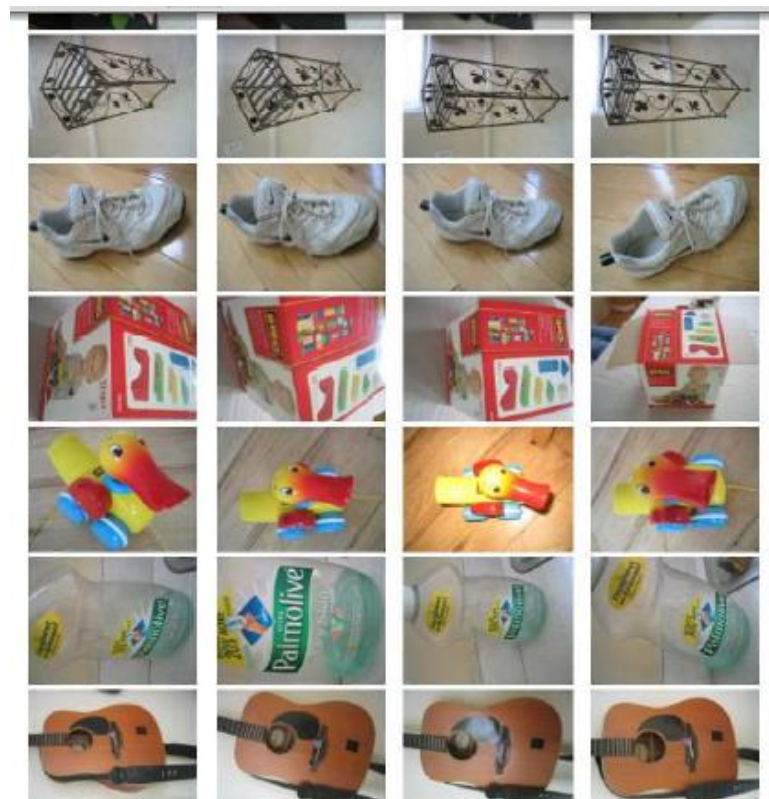for vocabulary of *V* words

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?
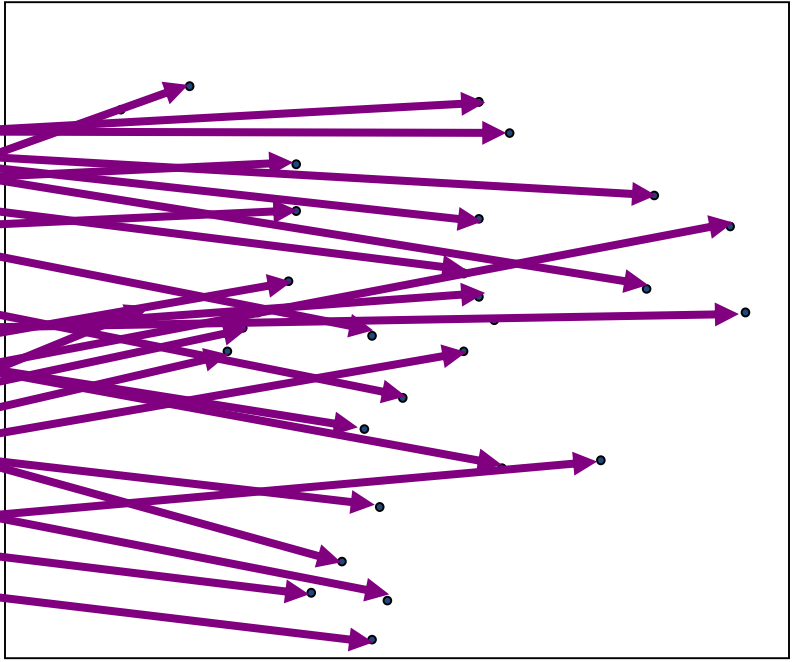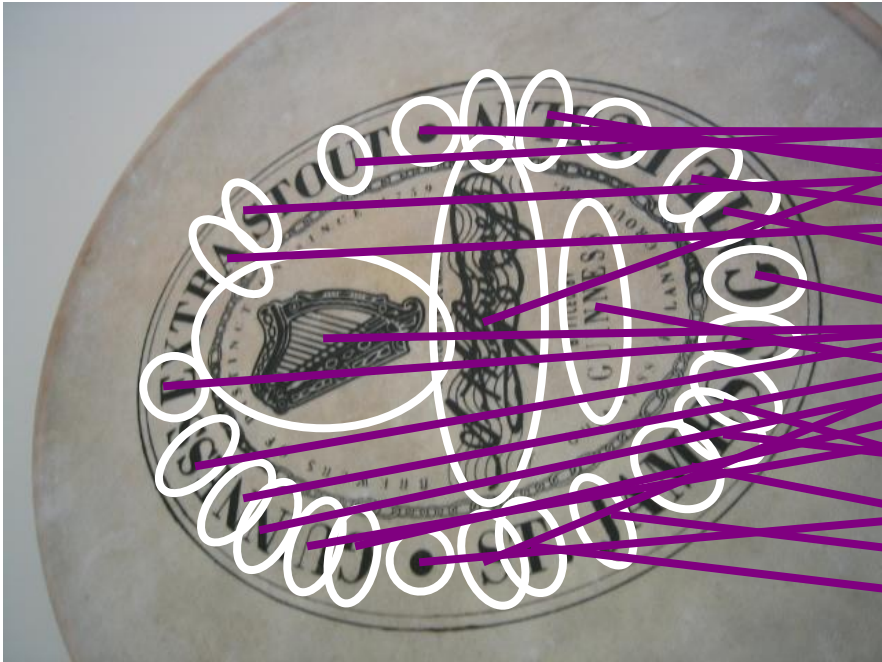
Kristen Grauman

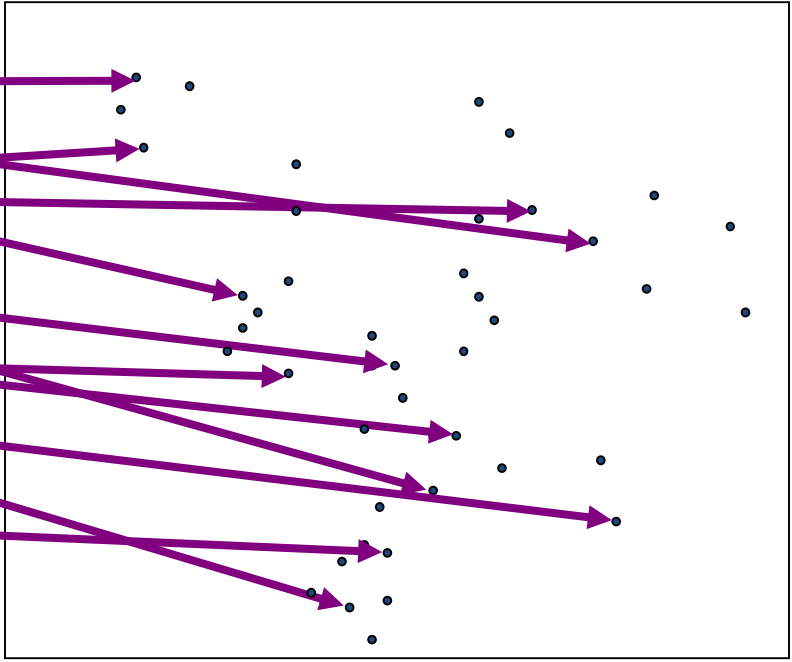# Vocabulary size

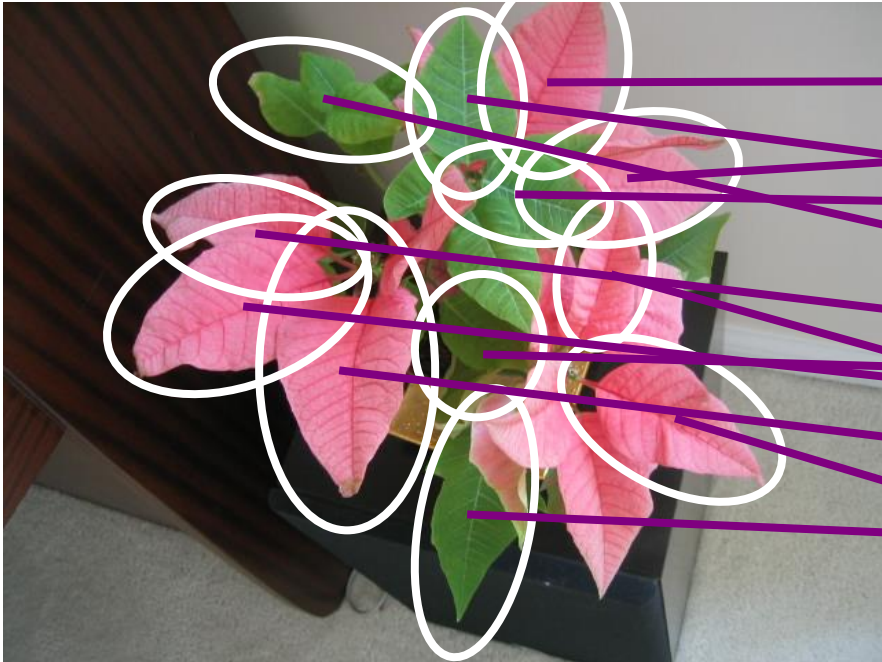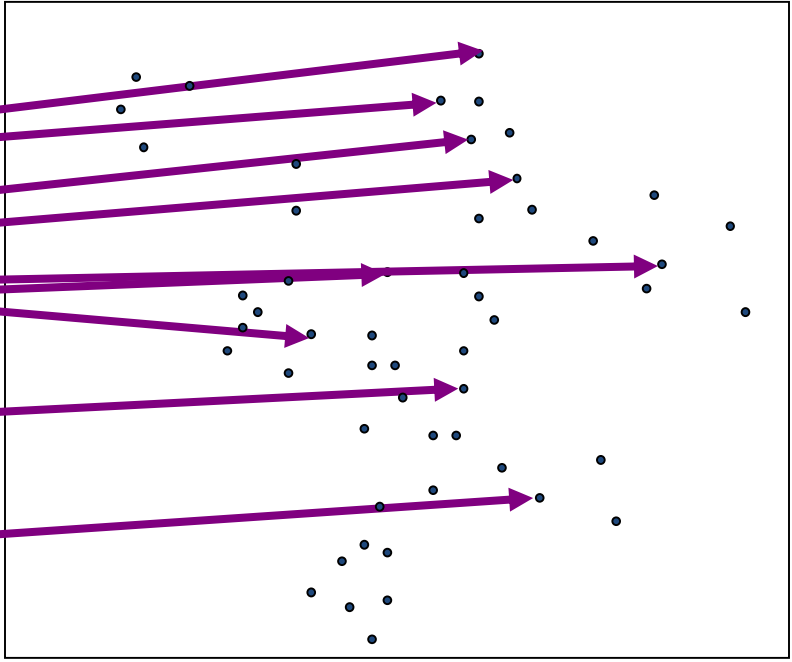Results for recognition task with 6347 images



*Influence on performance, sparsity*

Nister & Stewenius, CVPR 2006
Kristen Grauman

# Recognition with K-tree

# Vocabulary trees: complexity
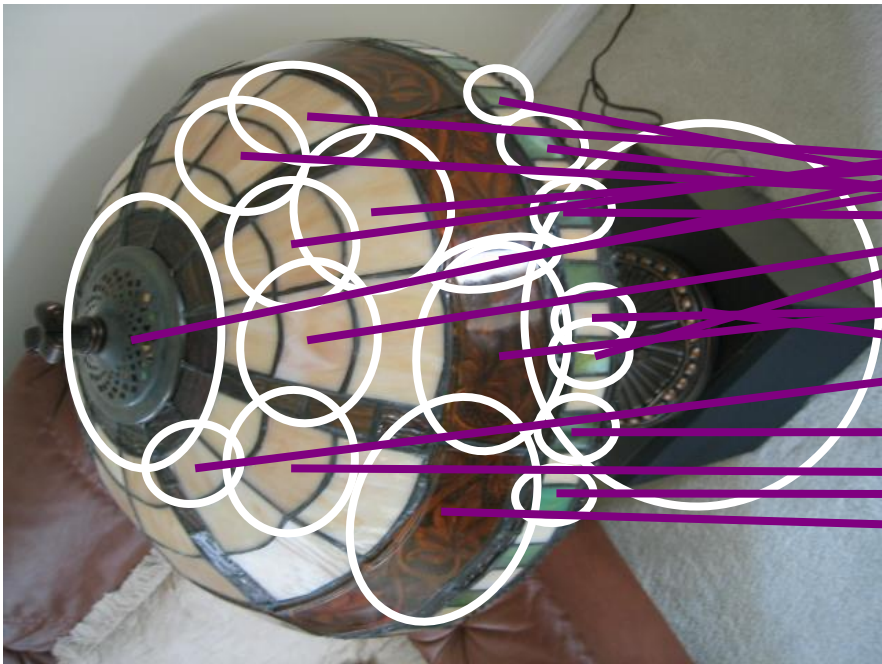
Number of words given tree parameters: branching factor and number of levels

$$\text{branching\_factor} \char`\^ \text{number\_of\_levels}$$

Word assignment cost vs. flat vocabulary

$O(k)$ for flat

$O(\log_{\text{branching\_factor}}(k) * \text{branching\_factor})$

Is this like a kd-tree?

Yes, but with better partitioning and defeatist search.

This hierarchical data structure is lossy – you might not find your true nearest cluster.

# 110,000,000 Images in 5.8 Seconds

# Higher branch factor works better
# (but slower)

# Visual words/bags of words

+ flexible to geometry / deformations / viewpoint
+ compact summary of image content
+ provides fixed dimensional 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

# 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

# Can we be more accurate?

So far, we treat each image as containing a "bag of words", with no spatial information

# Can we be more accurate?

So far, we treat each image as containing a "bag of words", with no spatial information



Real objects have consistent geometry

# Spatial Verification



Query

DB image with high BoW similarity

Query

DB image with high BoW similarity

Both image pairs have many visual words in common.

# Spatial Verification

Query



DB image with high BoW similarity

Query



DB image with high BoW similarity

# Only some of the matches are mutually consistent

# 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

$(x_i, y_i)$

$(x'_i, y'_i)$

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

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

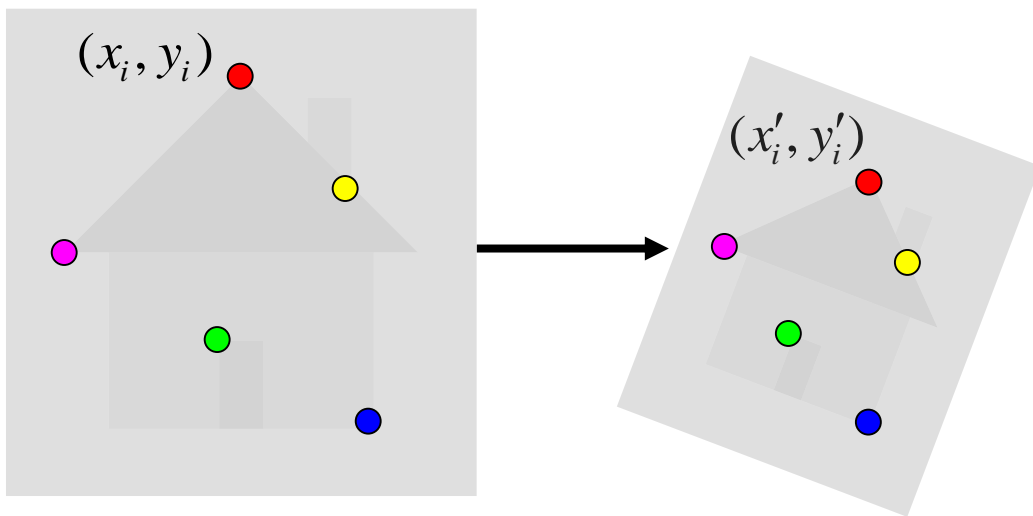$$\begin{bmatrix} & & \cdots & & & \\ x_i & y_i & 0 & 0 & 1 & 0 \\ 0 & 0 & x_i & y_i & 0 & 1 \\ & & \cdots & & & \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ t_1 \\ t_2 \end{bmatrix} = \begin{bmatrix} \cdots \\ x'_i \\ y'_i \\ \cdots \end{bmatrix}$$

# RANSAC verification

# 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?
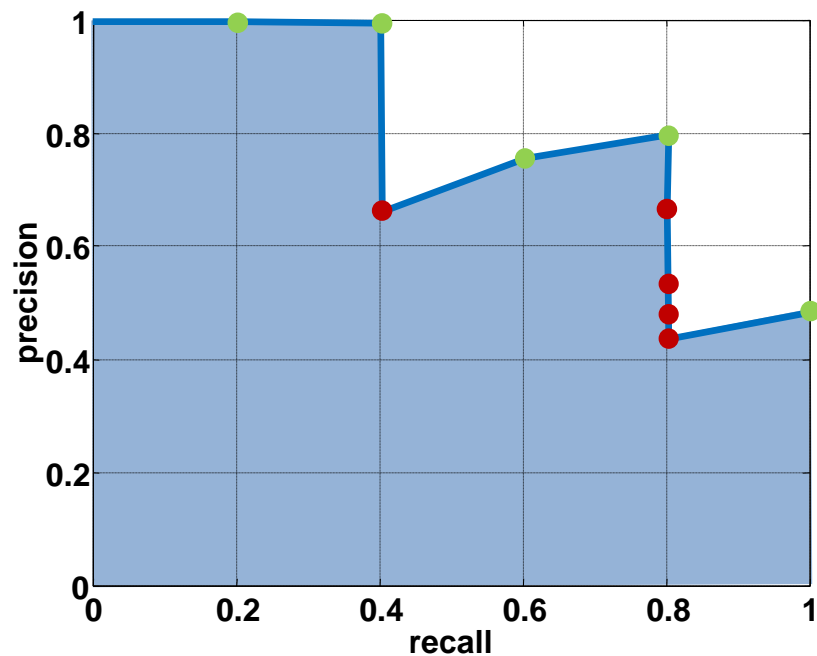
- How to score the retrieval results?

Kristen Grauman
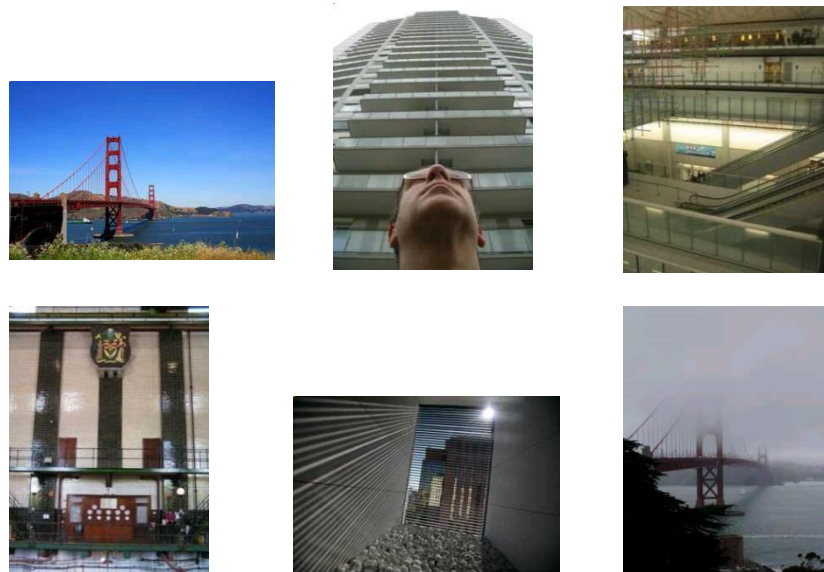
# Scoring retrieval quality


Query

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

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



Results (ordered):

# What else can we borrow from text retrieval?



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

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 with ... $660bn. ... annoy th... China's ... delibe... agrees... yuan is ... governo... also need... demand so ... country. China ... yuan against the doll... and permitted it to trade within a narrow ... but the US wants the yuan to be allowed ... freely. However, Beijing has made it c... it will take its time and tread carefully be... allowing the yuan to rise further in value.

# *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

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 *golf*er expects to reach the *green* on a par-four hole in **...**
- Manufactures and sells synthetic *golf* putting *green*s 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

# Query Expansion

Results



Spatial verification

Query image

New results

New query

Chum, Philbin, Sivic, Isard, Zisserman: Total Recall..., ICCV 2007

# 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.

# 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

# Lessons from a Decade Later

- For *Category* recognition (project 4)

  – Bag of Feature models remained the state of the art until Deep Learning.

  – Spatial layout either isn't that important or its too difficult to encode.

  – Quantization error is, in fact, the bigger problem. Advanced feature encoding methods address this.

  – Bag of feature models are nearly obsolete. At best they seem to be inspiring tweaks to deep models e.g. NetVLAD.

James Hays

# Lessons from a Decade Later

- For *instance* retrieval (this lecture)

  – deep learning is taking over.

  – learn better local features (replace SIFT) e.g. MatchNet

  – or learn better image embeddings (replace the histograms of visual features) e.g. Vo and Hays 2016.

  – or learn to do spatial verification e.g. DeTone, Malisiewicz, and Rabinovich 2016.

  – or learn a monolithic deep network to recognition all locations e.g. Google's PlaNet 2016.