


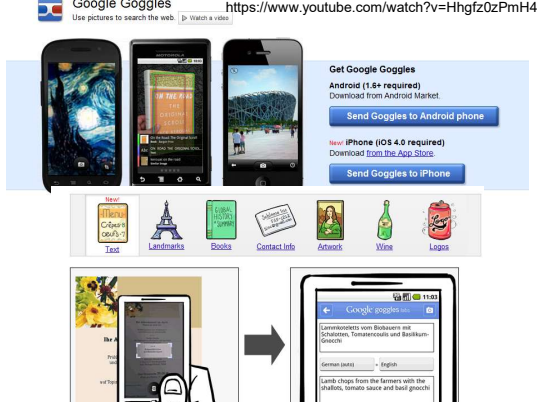
Today

- Instance recognition
 - Indexing local features efficiently
 - Spatial verification models



Google Goggles <https://www.youtube.com/watch?v=HhgFz0zPmH4>

Use pictures to search the web. | Watch a video





Recognizing or retrieving specific objects


Example I: Visual search in feature films

Visually defined query

"Groundhog Day" [Rammis, 1993]


"Find this clock" 

"Find this place" 



Recognizing or retrieving specific objects

Example II: Search photos on the web for particular places

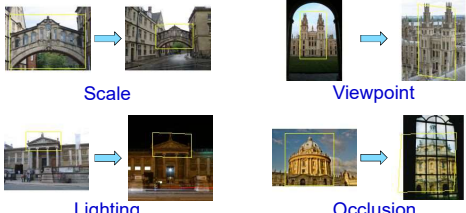


Find these landmarks ...in these images and 1M more

Slide credit: J. Sivic

Why is it difficult?

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



Slide credit: J. Sivic

Recall: matching local features

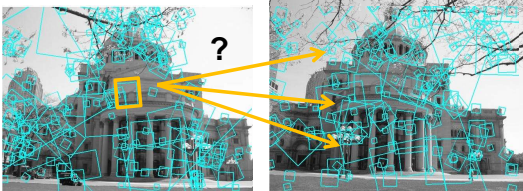


Image 1 Image 2

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)

Slide credit: Kristen Grauman

Multi-view matching

The diagram is divided into two parts. The left part, titled 'Matching two given views for depth', shows two camera frustums and two corresponding images of a building. Red lines connect corresponding points between the two images. The right part, titled 'Search for a matching view for recognition', shows a central image of a building with a large question mark. Red lines radiate from this central image to several other images of the same building from different angles, illustrating a search process. A vertical ellipsis indicates a large number of potential views.

Matching two given views for depth

Search for a matching view for recognition

Slide credit: Kristen Graumar

Indexing local features

The diagram shows two images of a building on the left. Red arrows point from various local features in these images to a central vertical column of small feature patches. A yellow arrow points from a feature patch in this column to a corresponding feature in a third image on the right. A vertical ellipsis at the bottom of the central column indicates many more features.

Slide credit: Kristen Graumar

Indexing local features

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

The diagram shows two images of a building on the left with several small patches highlighted. Arrows point from these patches to a 3D coordinate system on the right. Inside this coordinate system, several black dots represent the descriptors for the patches. The text 'Descriptor's feature space' is written below the coordinate system.

Descriptor's feature space

Slide credit: Kristen Graumar

Indexing local features

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

Slide credit: Kristen Graumar



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?
- Possible solutions:
 - Inverted file
 - Nearest neighbor data structures
 - Kd-trees
 - Hashing

Slide credit: Kristen Graumar



Indexing local features: inverted file index

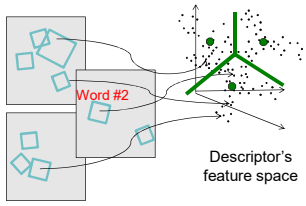
- For text documents, an efficient way to find all *pages* on which a *word* occurs is to use an index...
- We want to find all *images* in which a *feature* occurs.
- To use this idea, we'll need to map our features to "visual words".

Slide credit: Kristen Graumar



Visual words

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



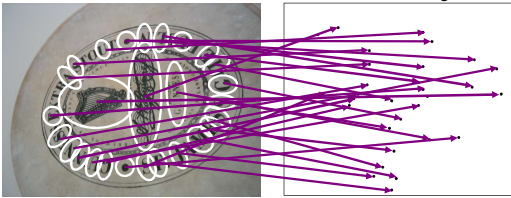
Descriptor's feature space

- 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.

Slide credit: Kristen Grauman

Visual words: main idea

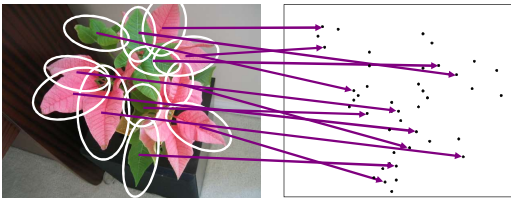
- Extract some local features from a number of images ...

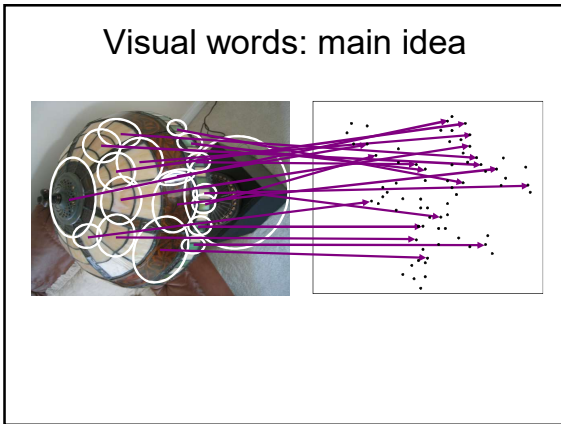


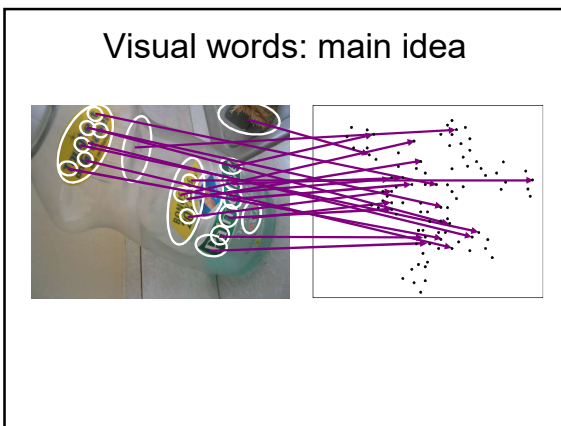
e.g., SIFT descriptor space: each point is 128-dimensional

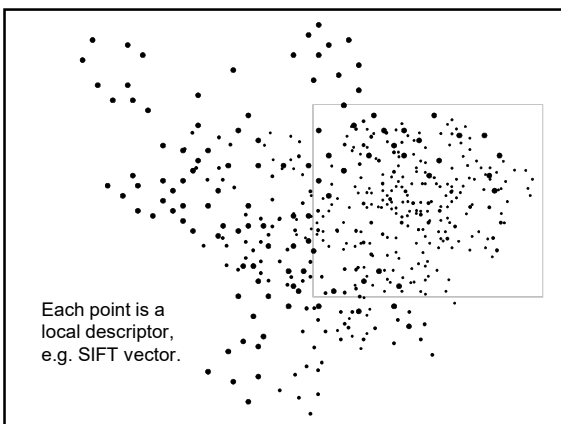
Slide credit: D. Nister, CVPR 2006

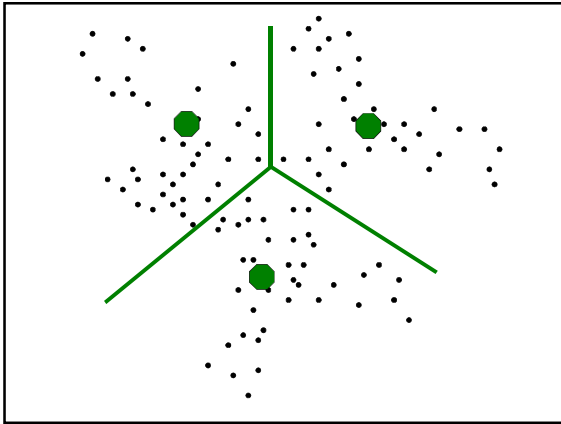
Visual words: main idea





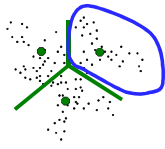






Visual words

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



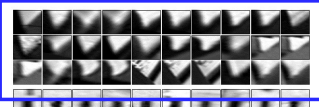
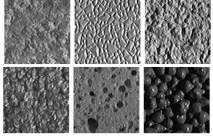
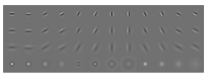


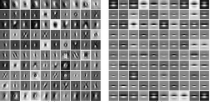
Figure from Sivic & Zisserman, ICCV 2003

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
Slide credit: Kristen Grauman

Recall: Texture representation example

	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

Slide credit: Kristen Grauman

Visual vocabulary formation

Issues:

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

Slide credit: 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

Slide credit: Kristen Grauman

Inverted file index

When will this give us a significant gain in efficiency?

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.

Slide credit: Kristen Graumar

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?

Slide credit: Kristen Graumar

Analogy to documents

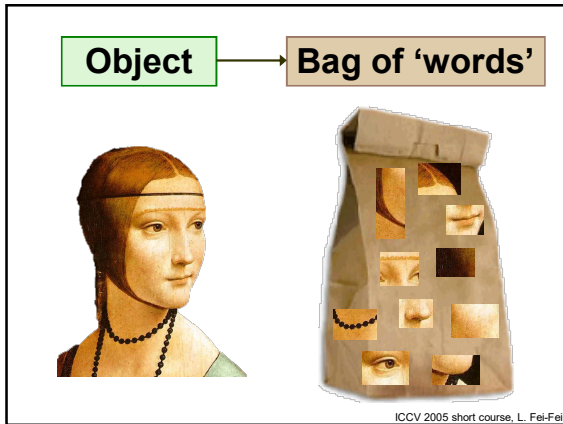
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the visual centers of the brain were considered as a movie screen. It was not until the discovery of the receptive field that we know the perception is more complex. 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 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.

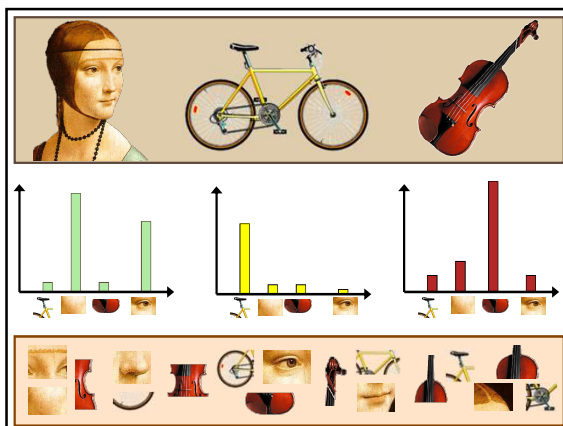
China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports, compared with \$660bn. The government also needs to curb demand for yuan against the dollar. China's government also needs to curb demand for yuan against the dollar. China's government also needs to curb demand for yuan against the dollar.

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 inner 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

Slide credit: Kristen Graumar

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)

Number of occurrences of word i in document d \rightarrow $t_i = \frac{n_{id}}{n_d} \log \frac{N}{n_i}$ \leftarrow Total number of documents in database

Number of words in document d \rightarrow n_d \leftarrow Number of documents word i occurs in, in whole database

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

Slide credit: Kristen Graumar

Bags of words for content-based image retrieval

Visually defined query


Find this clock



Find this place




"Groundhog Day" [Rammis, 1993]



Slide from Andrew Zisserman
Sivic & Zisserman, ICCV 2003

Example



retrieved shots

Start frame 52907	Key frame 53026	End frame 53028
Start frame 54342	Key frame 54376	End frame 54644
Start frame 51170	Key frame 52251	End frame 52248
Start frame 54079	Key frame 54201	End frame 54201
Start frame 39097	Key frame 39126	End frame 39306
Start frame 40760	Key frame 40826	End frame 41041
Start frame 39701	Key frame 39676	End frame 39730


Slide from Andrew Zisserman
Sivic & Zisserman, ICCV 2003

Video Google System

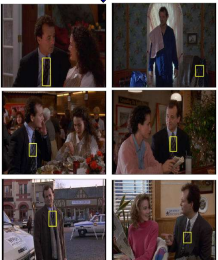
1. Collect all words within query region
2. Inverted file index to find relevant frames
3. Compare word counts
4. Spatial verification

Sivic & Zisserman, ICCV 2003

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



Query region



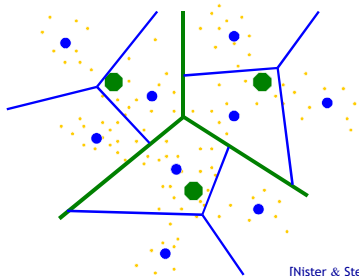
Retrieved frames

K. Grauman, B. Leibe

Visual Object Recognition Tutorial

Vocabulary Trees: hierarchical clustering for large vocabularies

- Tree construction:



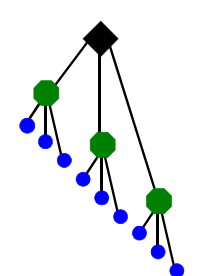
[Nister & Stewenius, CVPR'06]
K. Grauman, B. Leibe Slide credit: David Nister

The diagram illustrates hierarchical clustering. It shows a collection of points (yellow dots) in a 2D space. Some points are highlighted in green and blue. Lines connect these points to form a tree structure, with green lines representing the primary branches and blue lines representing secondary branches. The root of the tree is at the top, and it branches out downwards.

Visual Object Recognition Tutorial

Vocabulary Tree

- Training: Filling the tree



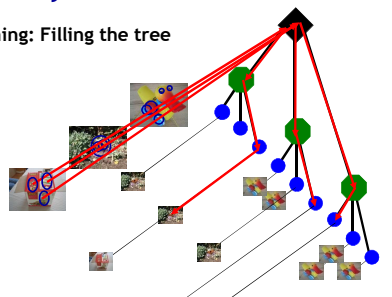
[Nister & Stewenius, CVPR'06]
K. Grauman, B. Leibe Slide credit: David Nister

The diagram shows a tree structure with a black diamond at the root. Three green circles are connected to the root. Each green circle has several blue circles connected to it, representing a hierarchical structure of nodes.

Visual Object Recognition Tutorial

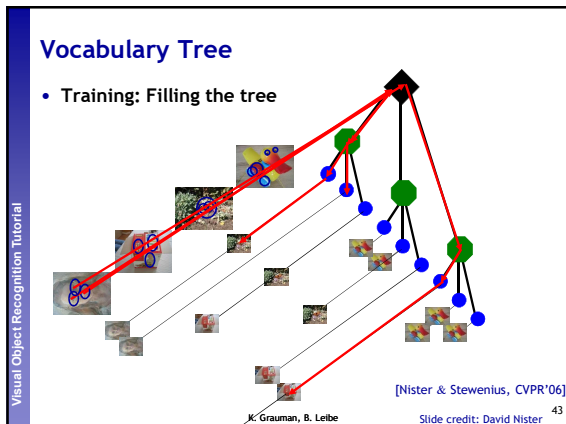
Vocabulary Tree

- Training: Filling the tree

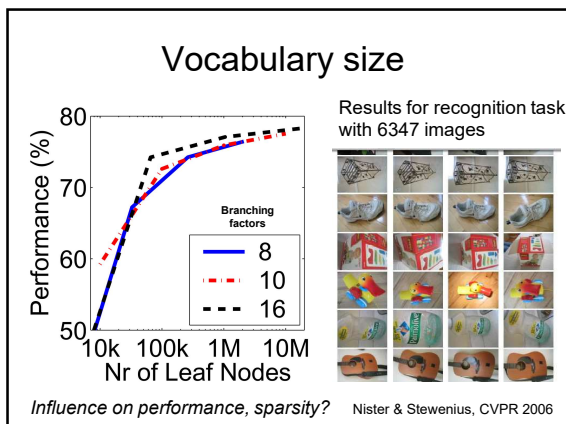


[Nister & Stewenius, CVPR'06]
K. Grauman, B. Leibe Slide credit: David Nister

The diagram shows a tree structure similar to the previous one, but with small images (e.g., a car, a person) attached to the nodes. Red lines connect the root to the green nodes, and black lines connect the green nodes to the blue nodes. The images are placed near the nodes they represent.



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



Bags of words: pros and cons

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

- basic model ignores geometry – must verify afterwards, or encode via features
- background and foreground mixed when bag covers whole image
- optimal vocabulary formation remains unclear

Slide credit: Kristen Grauman

Summary so far

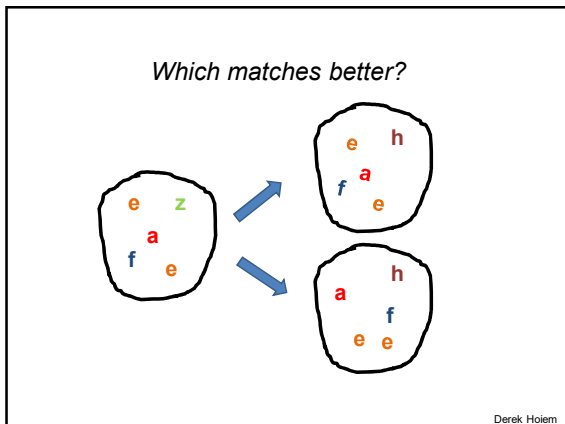
- **Matching local invariant features**
 - Useful not only to provide matches for multi-view geometry, but also to find objects and scenes.
- **Bag of words** representation: quantize feature space to make discrete set of visual words
 - Summarize image by distribution of words
 - Index individual words
- **Inverted index**: pre-compute index to enable faster search at query time
- **Recognition of instances via alignment**: matching local features

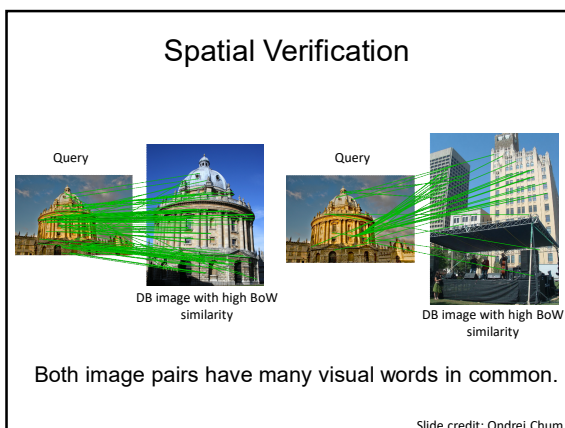
Kristen Grauman

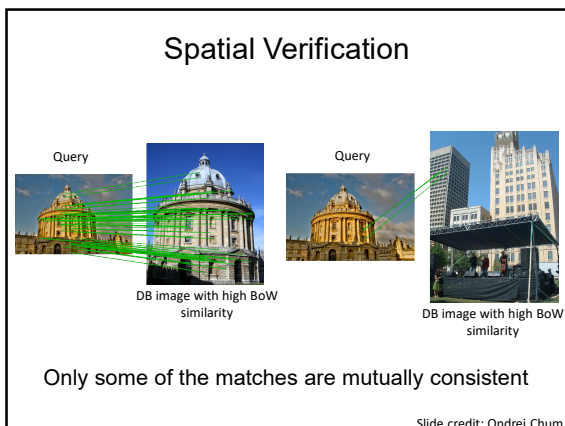
Instance recognition: remaining issues

- 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?

Slide credit: Kristen Grauman







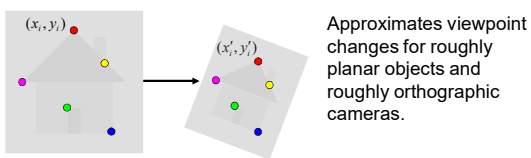
Spatial Verification: two basic strategies

- 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

RANSAC verification

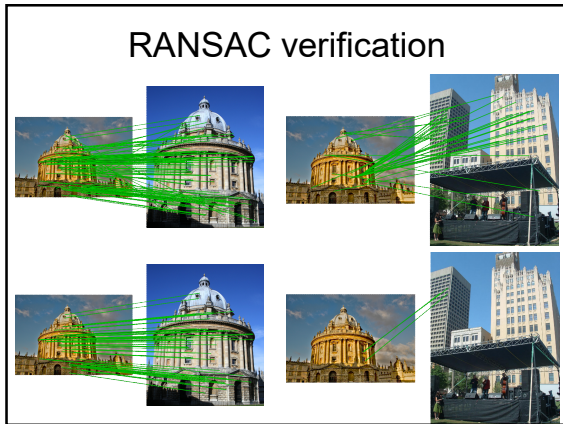


Recall: Fitting an affine transformation



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

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




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

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

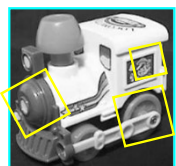


Novel image

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

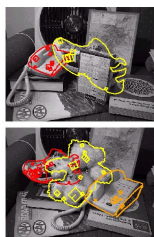
Example result



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

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

Slide credit: Kristen Graumar

Example Applications



- Mobile tourist guide**
- Self-localization
 - Object/building recognition
 - Photo/video augmentation



Visual Object Recognition Tutorial

B. Leibe

[Quack, Leibe, Van Gool, CVR'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

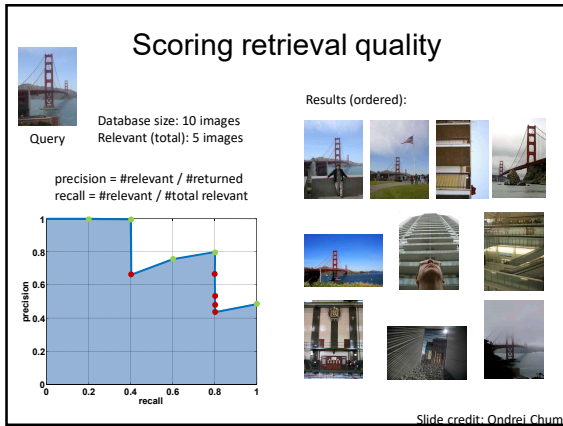
1. Take a picture with your mobile phone camera
2. Search it
 - in Switzerland to 8008 (Charge Customers 079-334 5100)
 - in Germany to 80000
 - available also to @kooba.ch
3. Search result is sent straight to your phone.

http://www.kooba.com/en/products_engine.html#

**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



- ### 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
- Kristen Grauman
