
CS688: Web-Scale Image Retrieval

Bag-of-Words (BoW) Models for Local Descriptors

Sung-Eui Yoon
(윤성익)

Course URL:
<http://sglab.kaist.ac.kr/~sungeui/IR>

KAIST



Class Objectives

- **Bag-of-visual-Word (BoW) model**
 - Pooling operation
- **Understand approximate nearest neighbor search**
 - Inverted index
 - Inverted multi-index

Object

Bag of 'words'

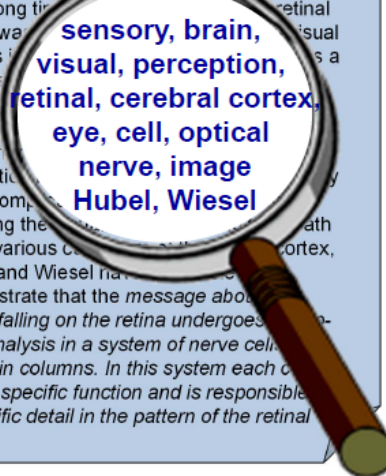


Represent an image
with a histogram of
words

Inspired by text search

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 image was considered as a movie of retinal centers. It was only after the discovery of the eye, cell, optical nerve, image Hubel, Wiesel

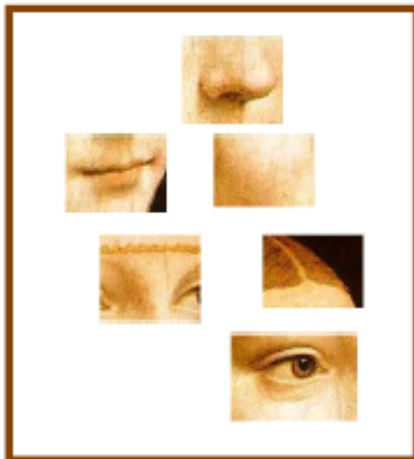
China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports, compared with \$750bn, \$660bn. The yuan is also needed to meet the demand so the country. China has permitted it to trade within a narrow band but the US wants the yuan to be allowed to trade 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.



definition of “BoW”

– Independent features

face



bike

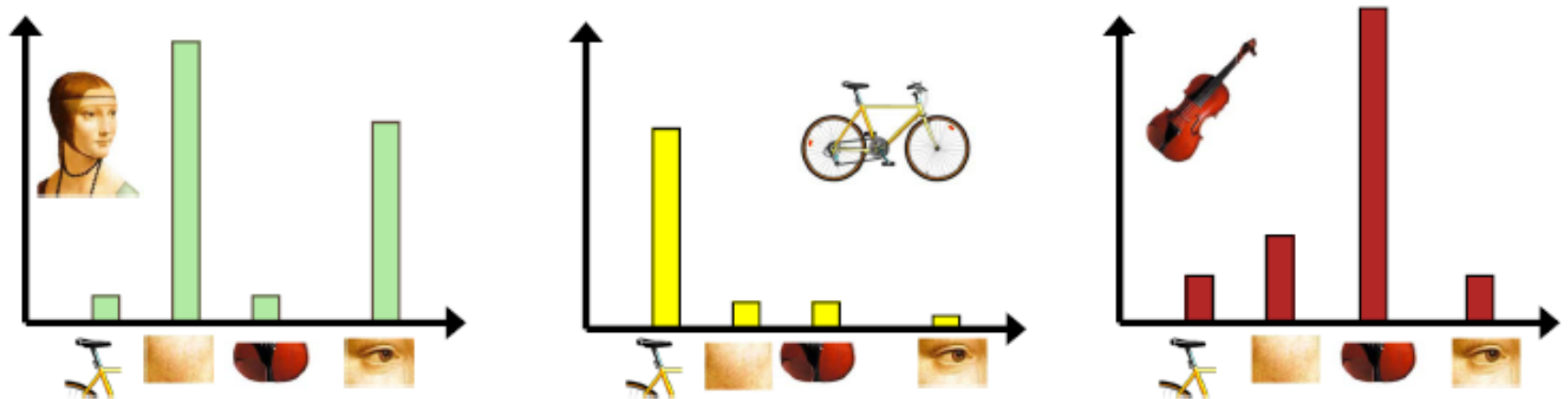


violin



definition of “BoW”

- Independent features
- histogram representation



codewords dictionary

Representation



feature detection
& representation

codewords dictionary



image representation



learning

category models
(and/or) classifiers

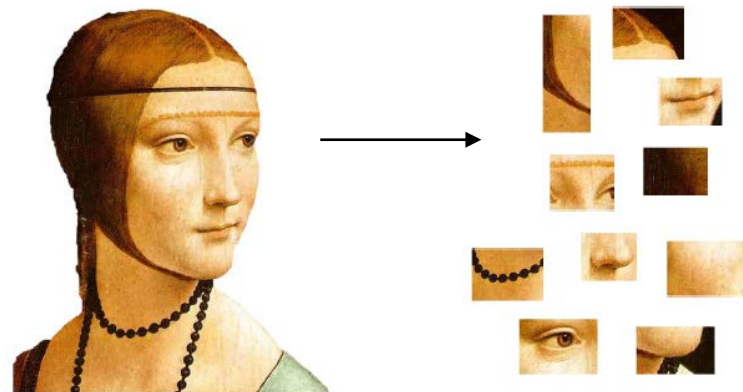
recognition



category
decision

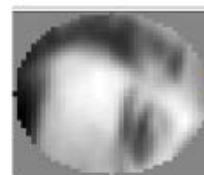
1. Feature Detection and Representations

- Assume many local features as an aggregation model
 - Global feature is not used
- Densely sampled or sampled only at key points
 - Detect patches extract features from them

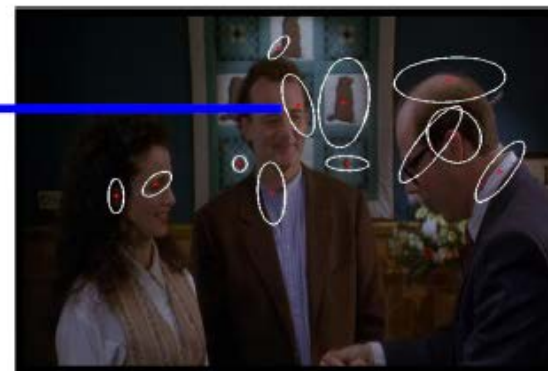


Ack.: Josef Sivic and Li Fei-Fei

Compute
SIFT
descriptor
[Lowe'99]



Normalize
patch

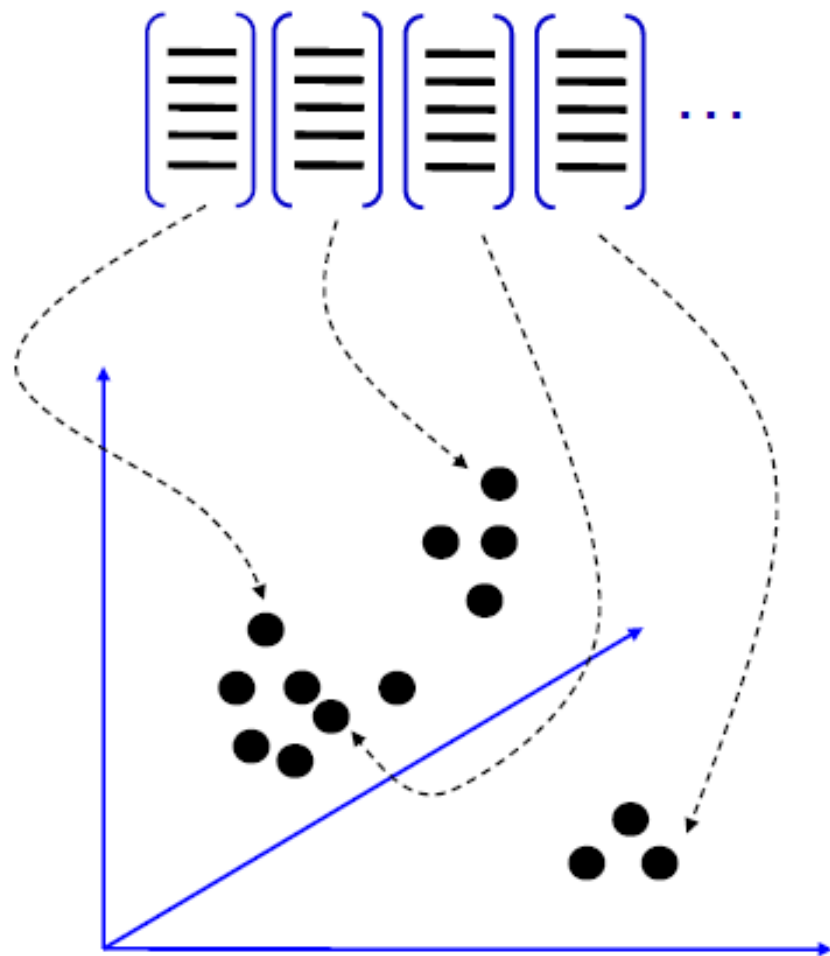


Detect patches

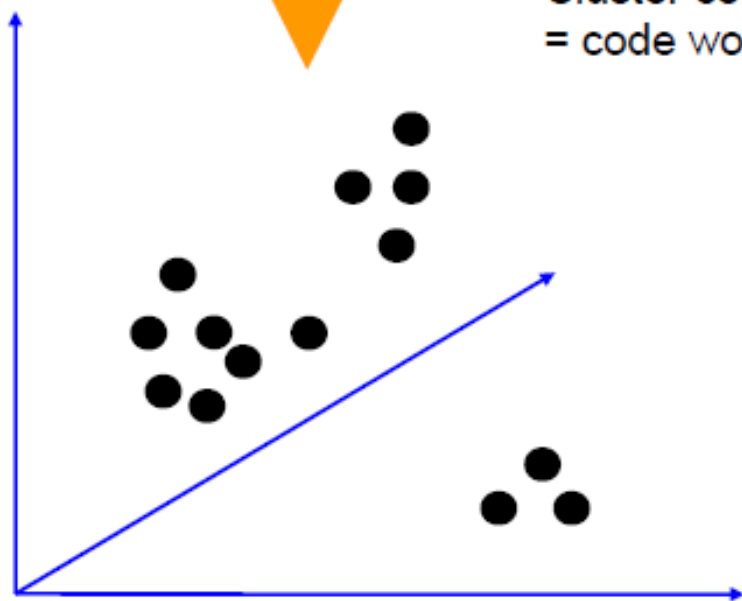
[Mikojaczyk and Schmid '02]

[Mata, Chum, Urban & Pajdla, '02]

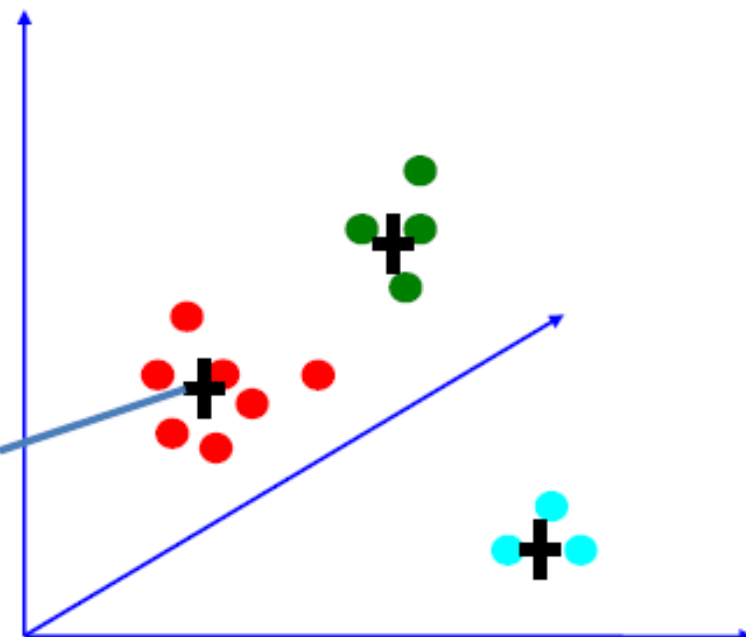
2. Codewords dictionary formation



2. Codewords dictionary formation



Cluster center
= code word



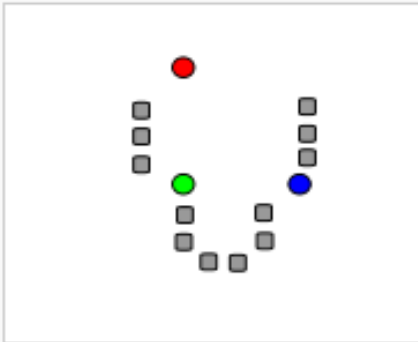
Clustering/
vector quantization

K-Means Clustering

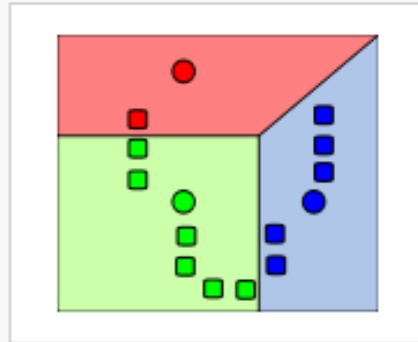
- An unsupervised learning
- Minimize the within-cluster sum of squares

$$\operatorname{argmin}_{\mathcal{S}} \sum_{i=1}^k \sum_{\mathbf{x}_j \in \mathcal{S}_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

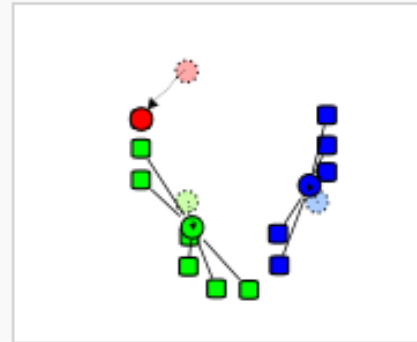
Demonstration of the standard algorithm



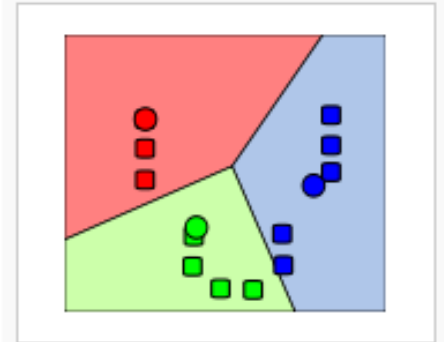
1) k initial "means" (in this case $k=3$) are randomly selected from the data set (shown in color).



2) k clusters are created by associating every observation with the nearest mean. The partitions here represent the [Voronoi diagram](#) generated by the means.

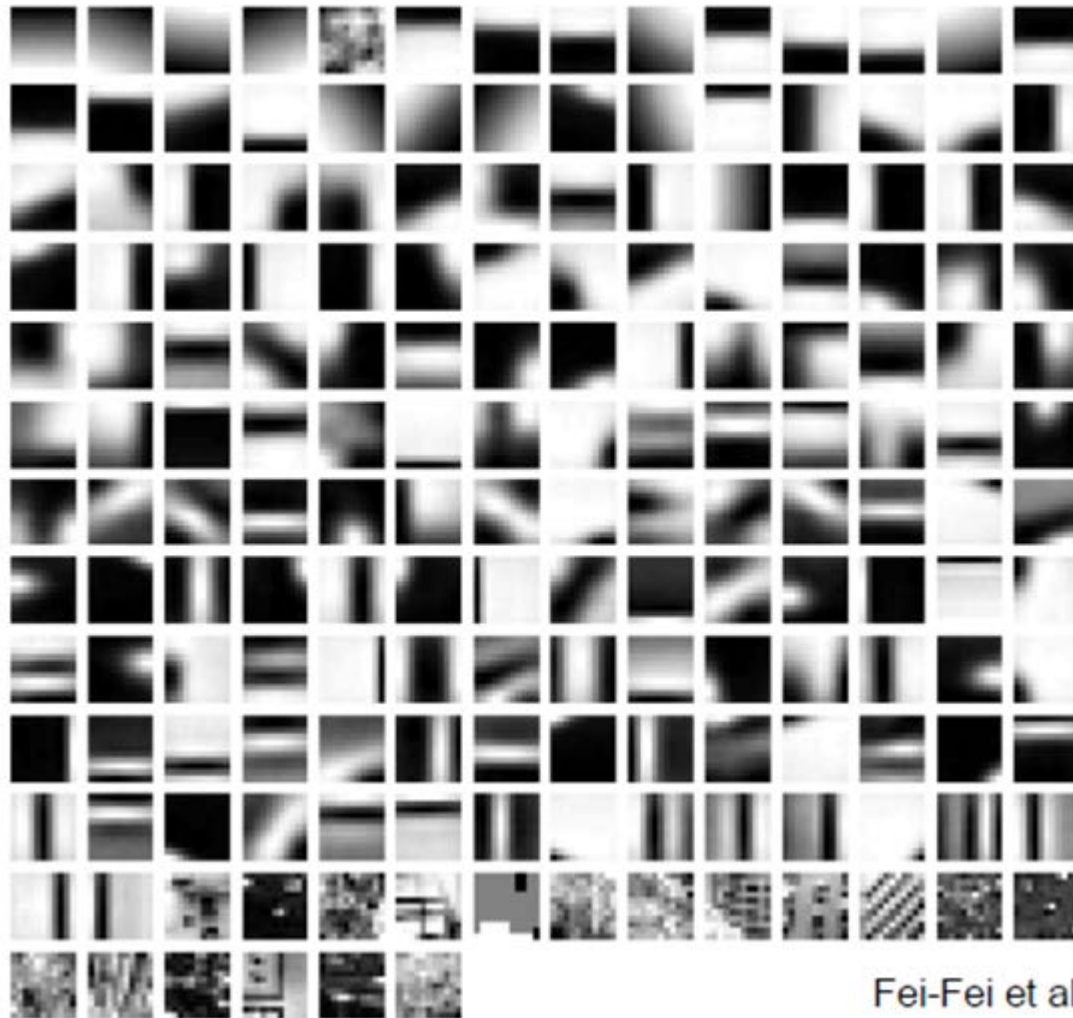


3) The [centroid](#) of each of the k clusters becomes the new means.



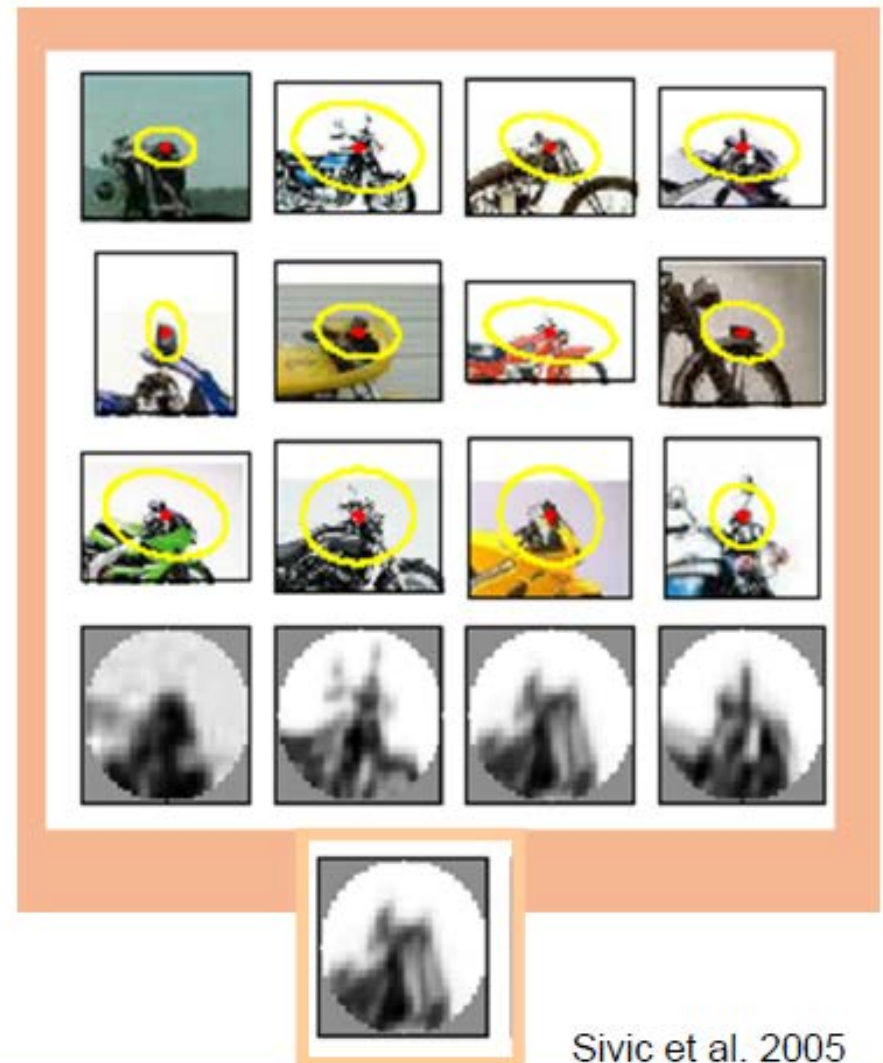
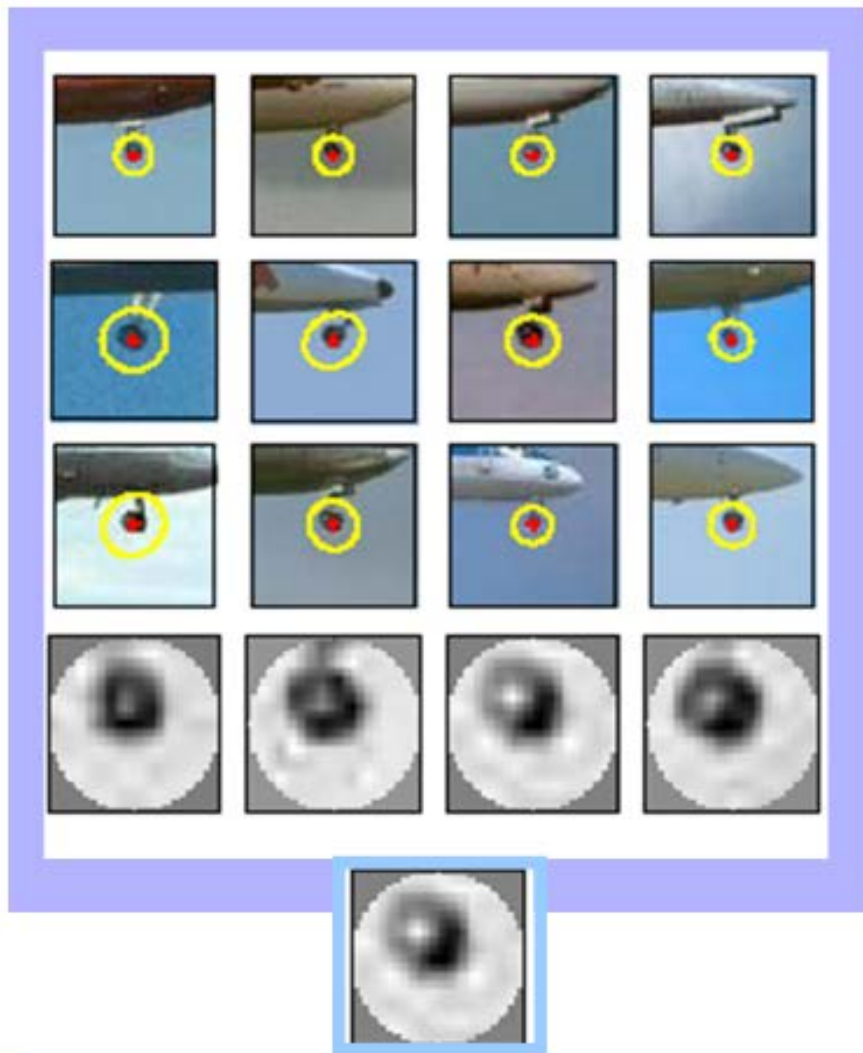
4) Steps 2 and 3 are repeated until convergence has been reached.

Codewords Dictionary Formation



Fei-Fei et al. 2005

Image patch examples of codewords

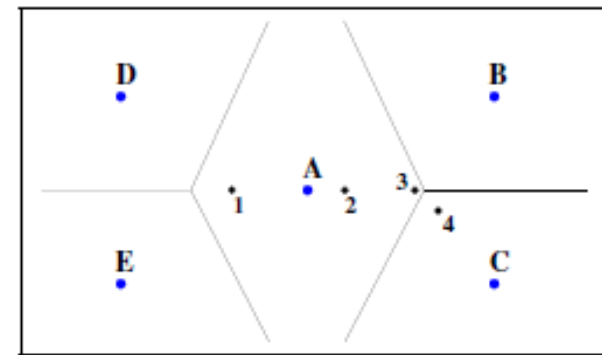
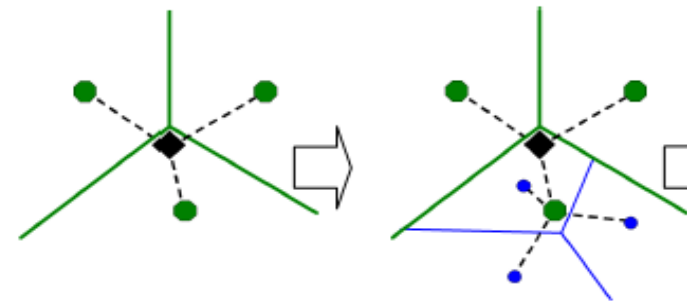


Issues of Visual Vocabulary

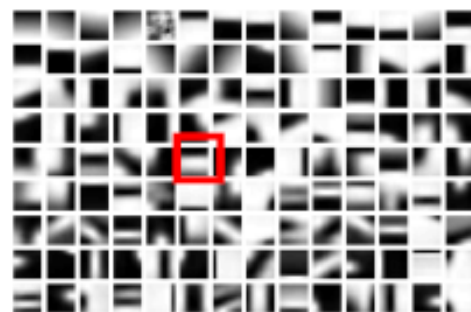
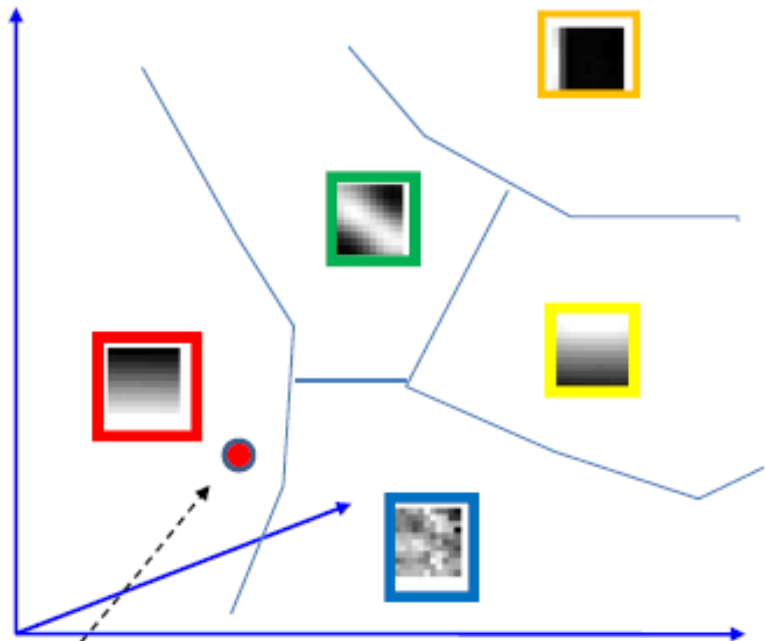
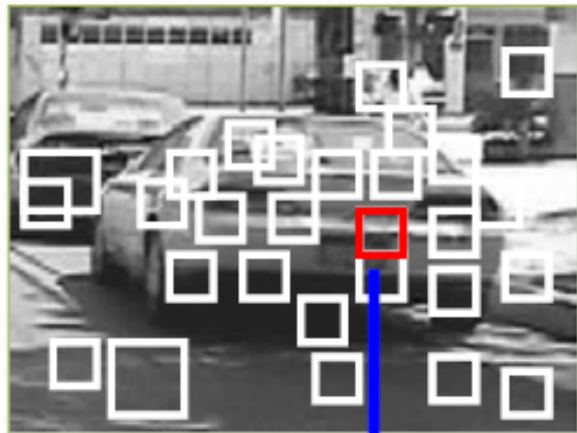
- Related to quantization
 - Too many words: quantization artifacts
 - Too small words: not representative
- K-means also takes long computation times

- Alternatives

- Faster performance: vocabulary tree, Nister et al.
- Low quantization artifacts: soft quantization, Philbin et al.



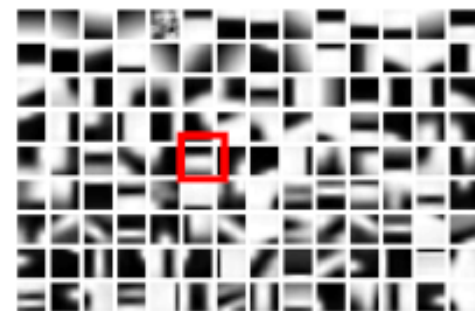
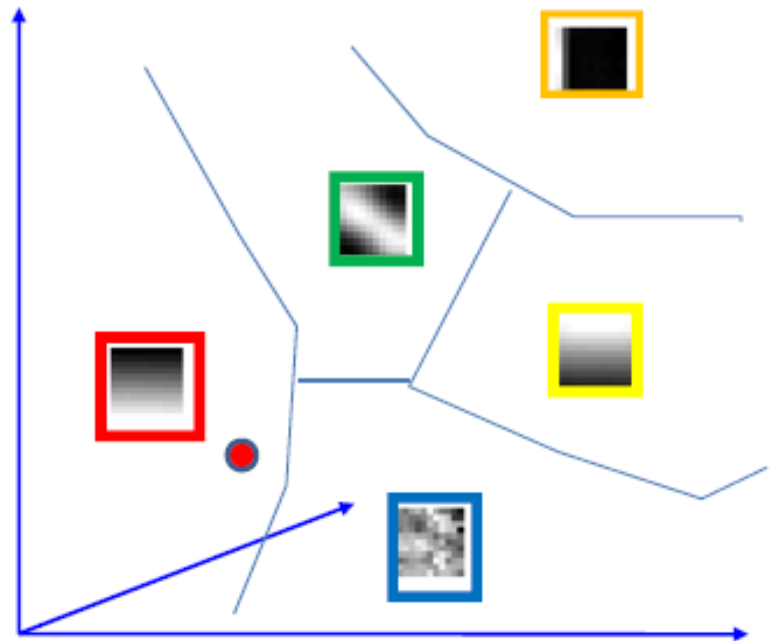
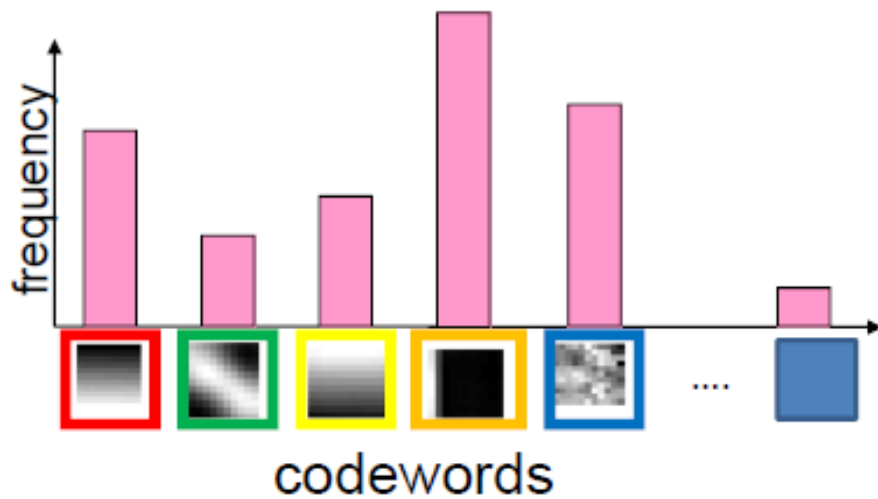
3. Bag of word representation



Codewords dictionary

- Nearest neighbors assignment
- K-D tree search strategy

3. Bag of word representation



Codewords dictionary

Representation



1. feature detection & representation



2. codewords dictionary

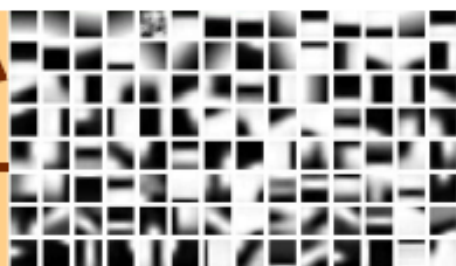
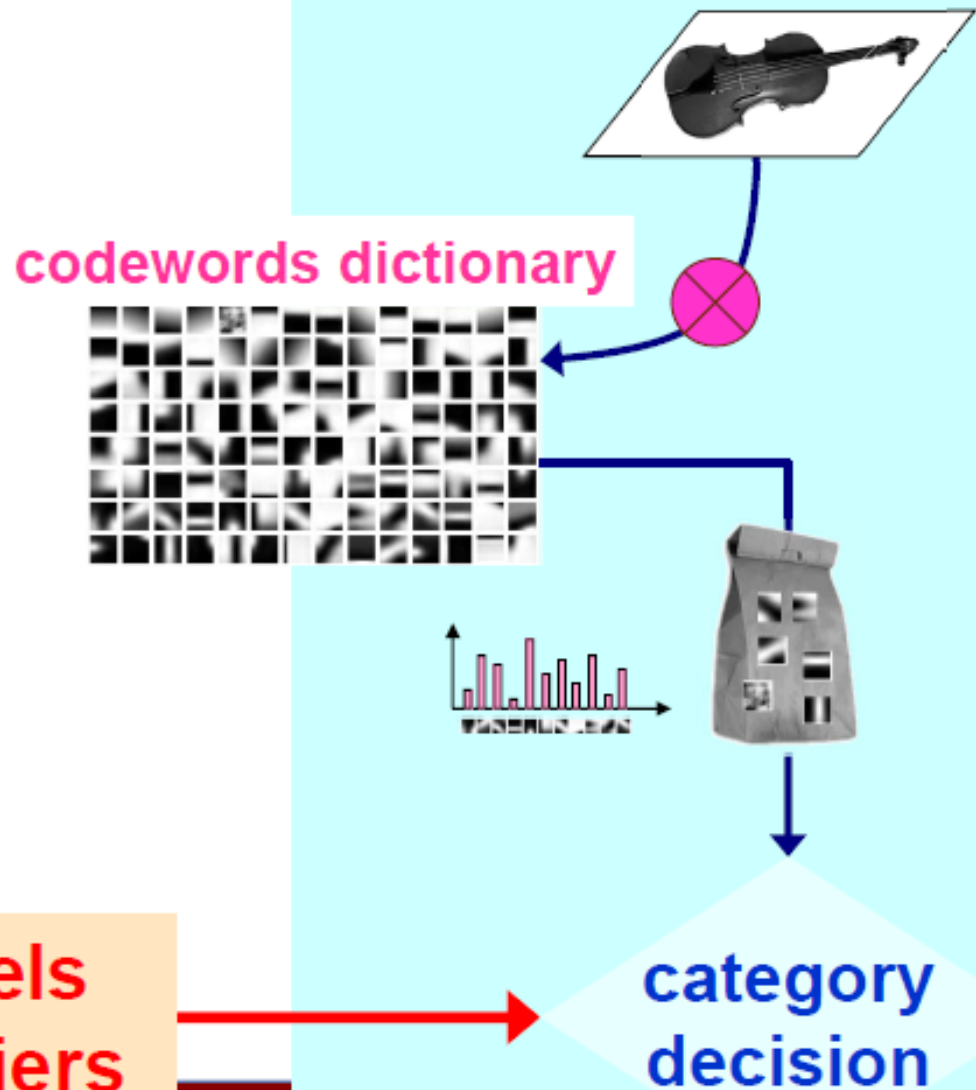


image representation

3.



Learning and Recognition



**category models
(and/or) classifiers**

TF-IDF

- Adopted from text search
 - A kind of weighting and normalization process
- Assume a document to be represented by $(t_1, \dots, t_i, \dots, t_k)^T$
- Weighted by TF (Term frequency) * log (IDF (Inverse Document Frequency))

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

- n_{id} : # of occurrences of word i in document d
- n_d : total # of words in the document d
- n_i : # of occurrences of term i in the whole database
- N : # of documents in the whole database

Similarity and Distance Functions

- Dot product measuring the angle between two vectors
- **L1 or Euclidean distance**

$$L1 (h_1, h_2) = \sum_i |h_1^i - h_2^i|$$

- χ^2 distance

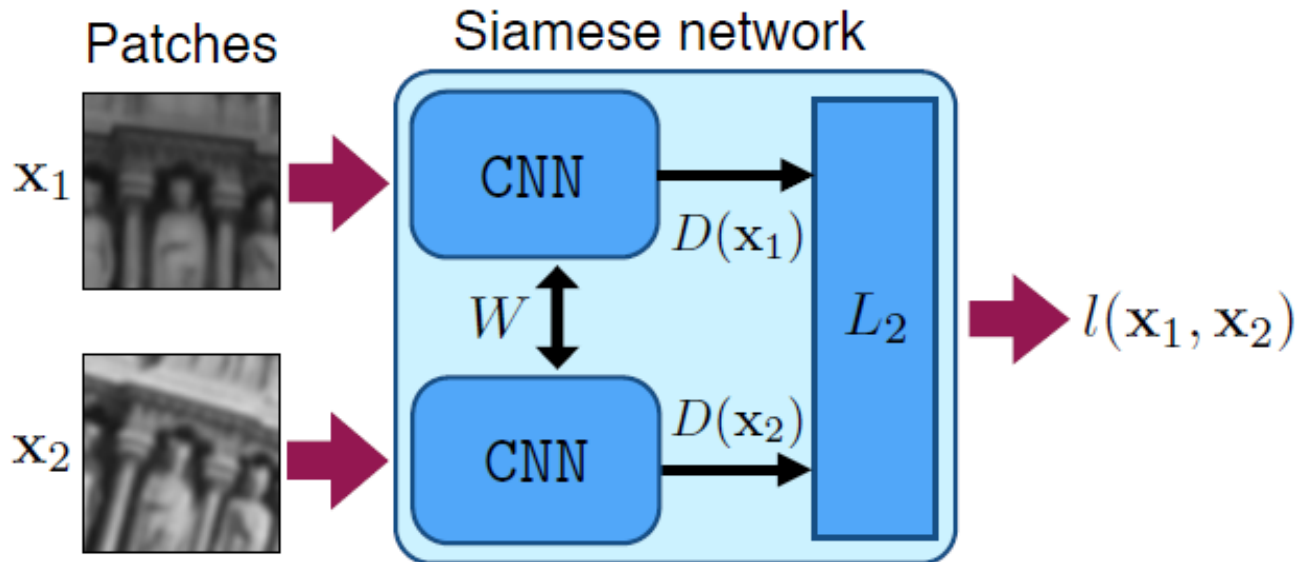
$$D(h_1, h_2) = \sum_{i=1}^N \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$$

- Quadratic distance (*cross-bin*)

$$D(h_1, h_2) = \sum_{i,j} A_{ij} (h_1(i) - h_2(j))^2$$

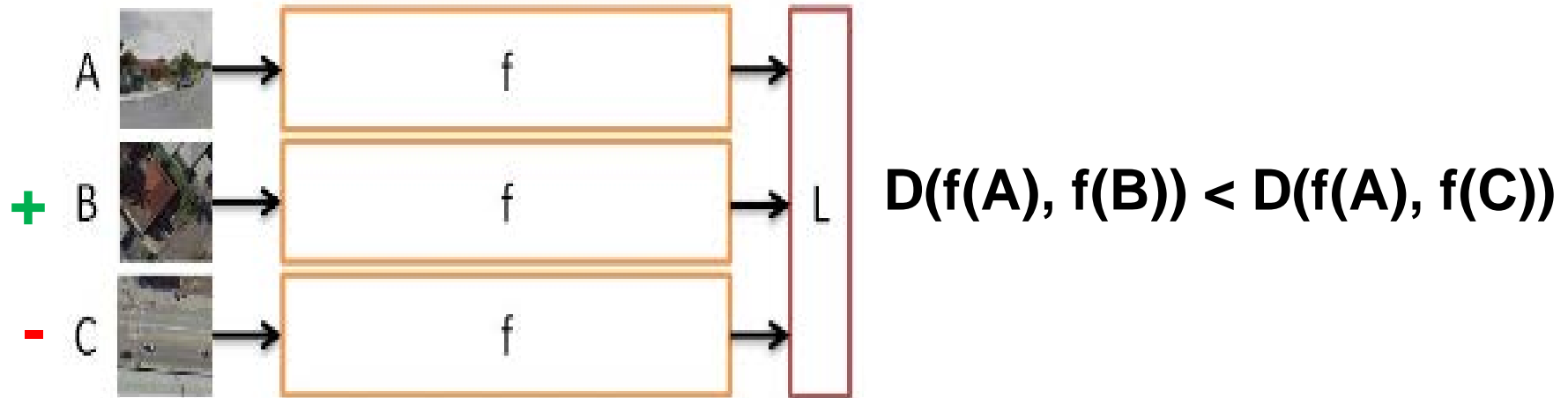
Similarity Learning: Siamese CNN

- Learn a feature representation mapping the sample patches with the L2 distance



Simo-Serra, E., Trulls, E., Ferraz, L., Kokkinos, I., Fua, P. and Moreno-Noguer, F., 2015. Discriminative learning of deep convolutional feature point descriptors. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 118-126).

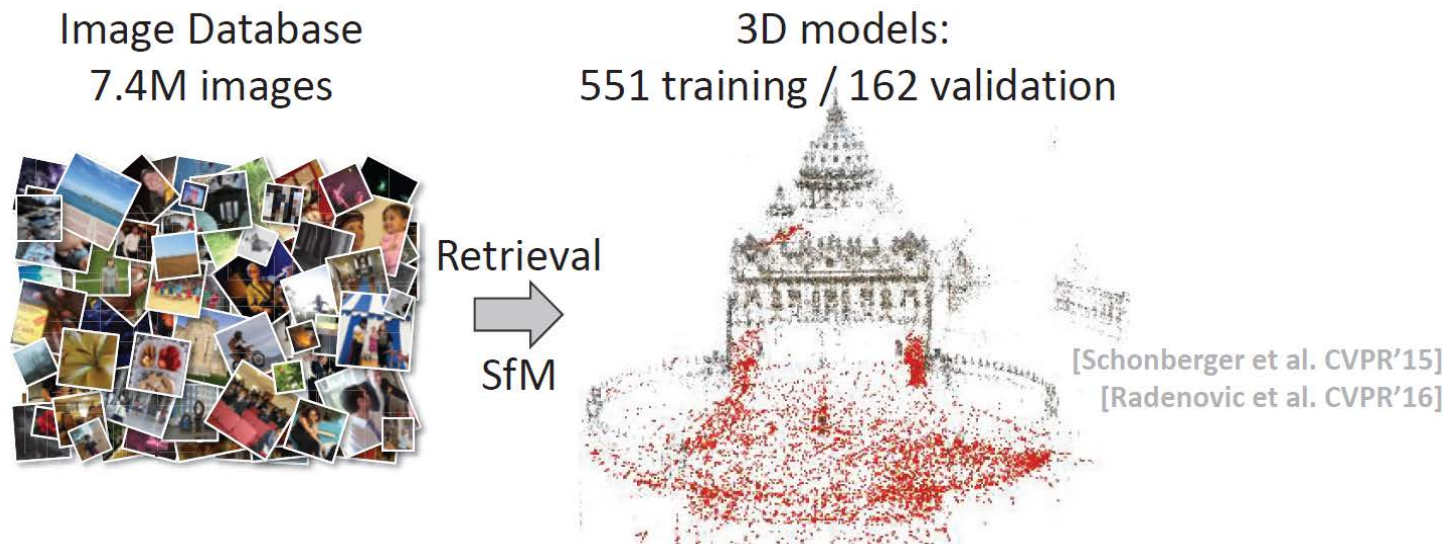
Siamese CNN Variants: Triplet Network or Loss



- Allows us to learn ranking between samples
 - Known as a ranking loss

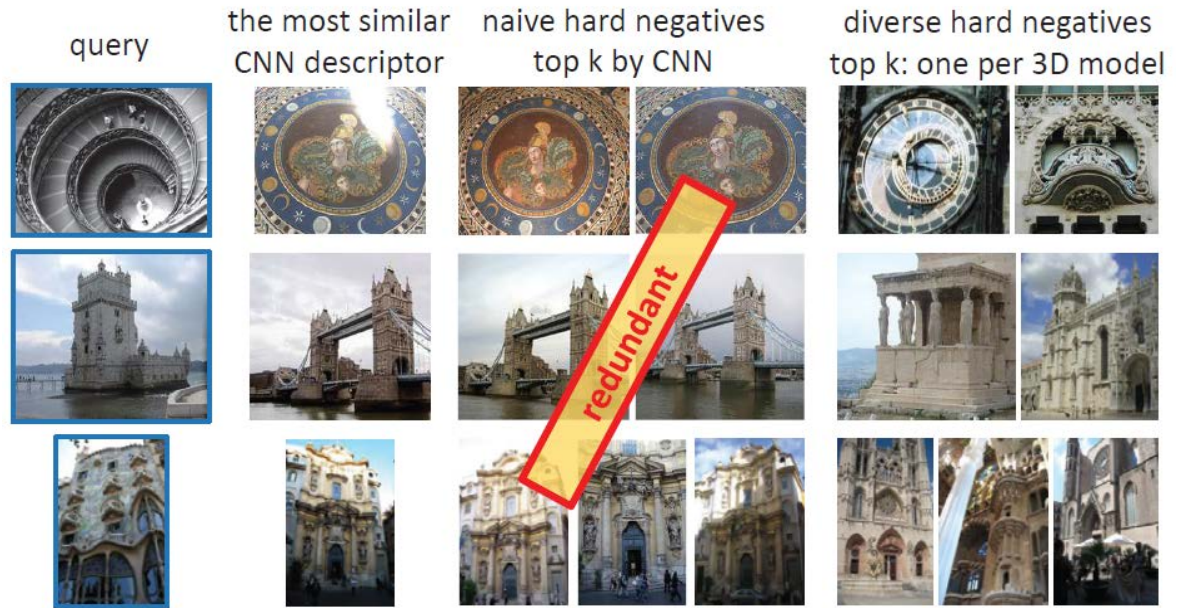
Utilize BoW for CNN Image Retrieval

- Construct 3D models from BoW based image retrieval
 - Unsupervised fine tuning with hard examples



- Given a query, identify its positive (same cluster or city) and its negative image given a query

Negative images

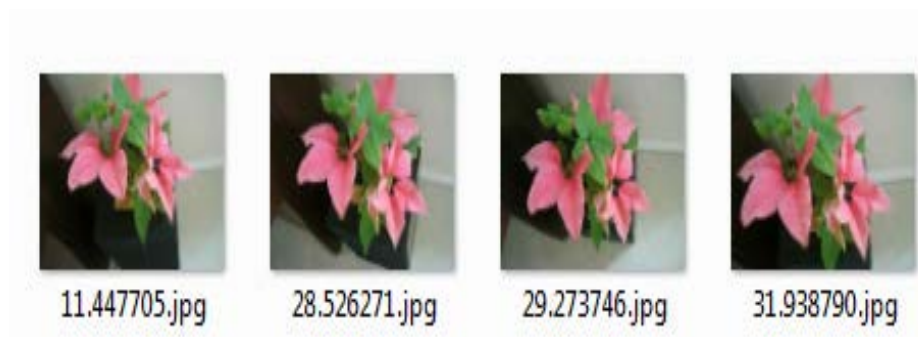
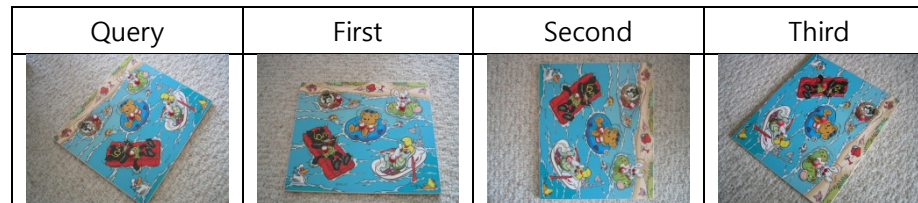


Positive images



PA2

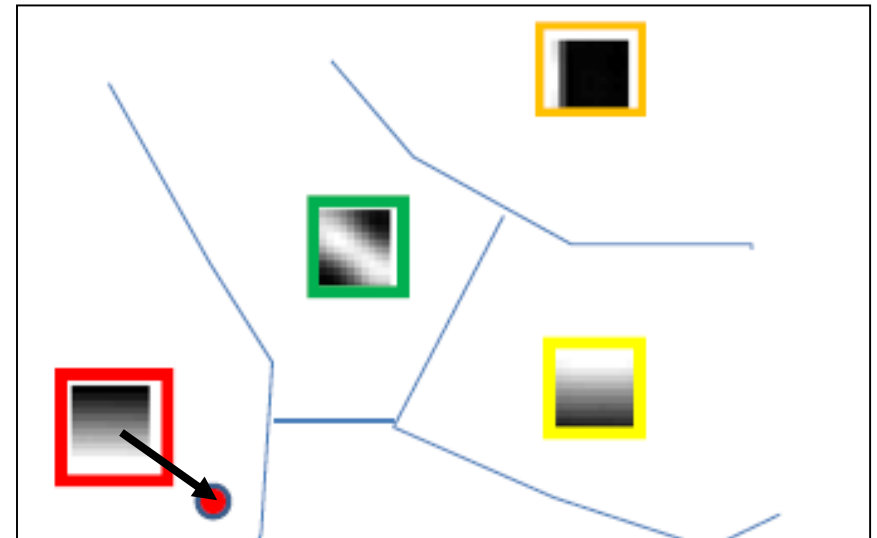
- Understand and implement a basic image retrieval system
- Use the original UKBenchmark
- Measure its accuracy



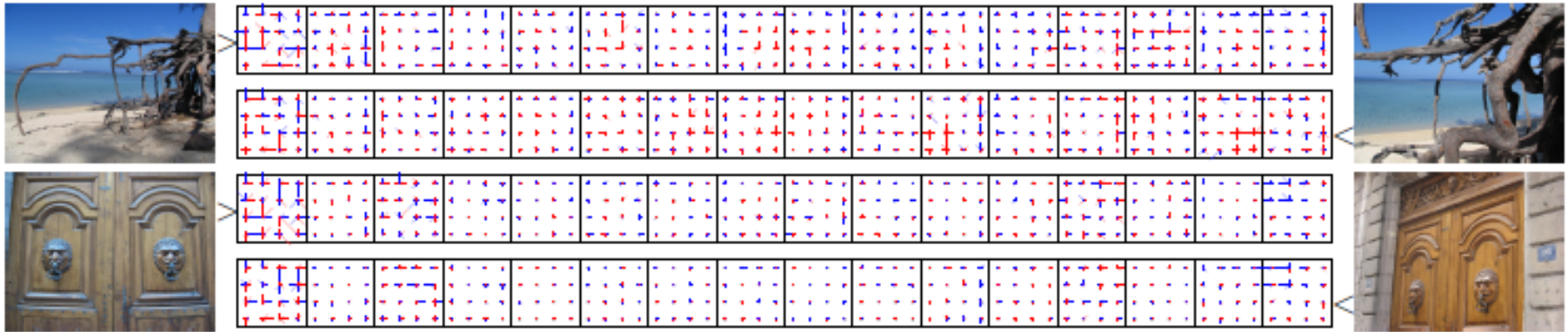
VLAD (Vector of Locally Aggregated Descriptors)

- BoW
 - Count the number of SIFTs assigned to each cluster
- VLAD
 - Compute the difference between a SIFT and its cluster center

$$v_{i,j} = \sum_{x \text{ such that } \text{NN}(x)=c_i} x_j - c_{i,j}$$



VLAD



- VLAD descriptors w/ 16 clusters
- Show better accuracy than BoW

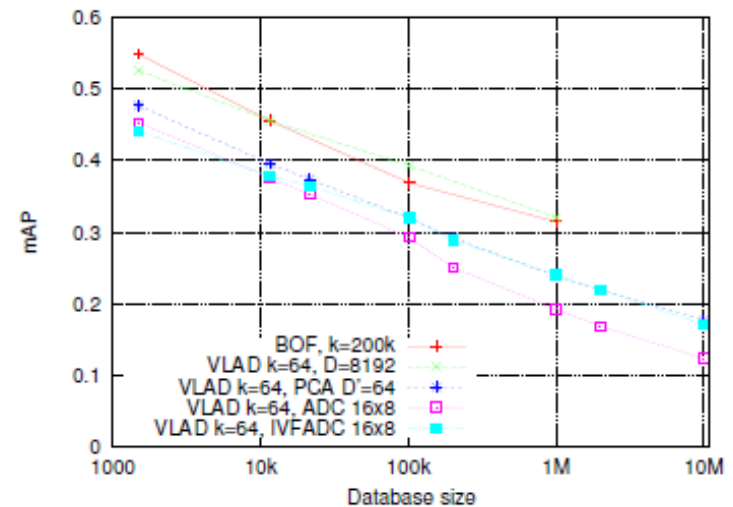
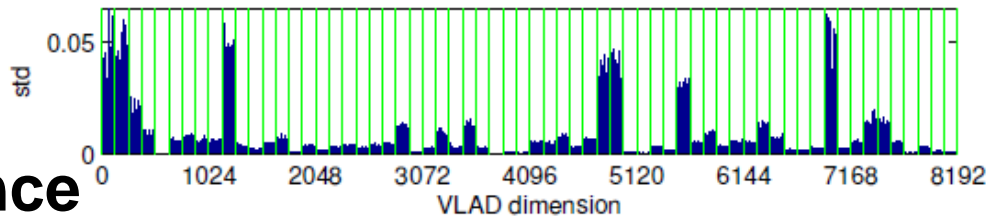


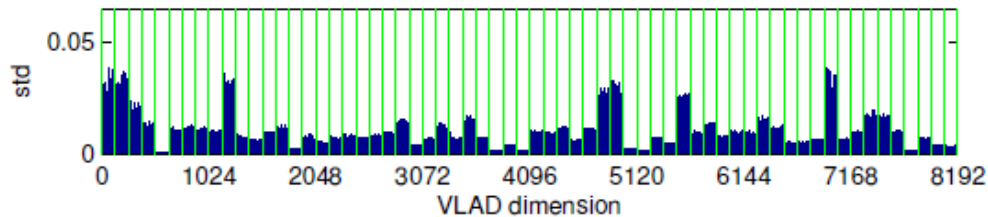
Figure 5. Search accuracy as a function of the database size.

Normalization for VLAD

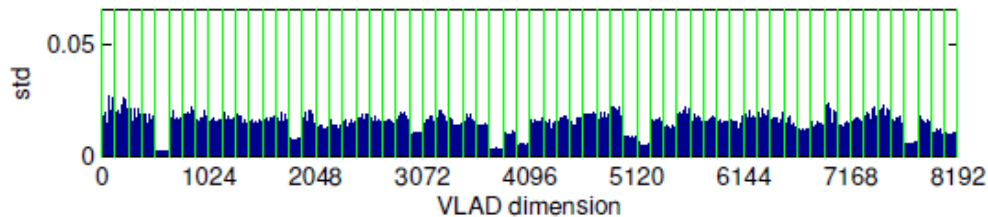
- Results in better accuracy



(a) Original VLAD normalization (L2)



(b) Signed square rooting (SSR) followed by L2



(c) Intra-normalization (innorm) followed by L2

L2 normalization,
i.e., $\frac{v}{|v|}$

Square rooting
for burstiness



L2 normalization
within each VLAD
block

NetVLAD: CNN architecture for weakly supervised place recognition

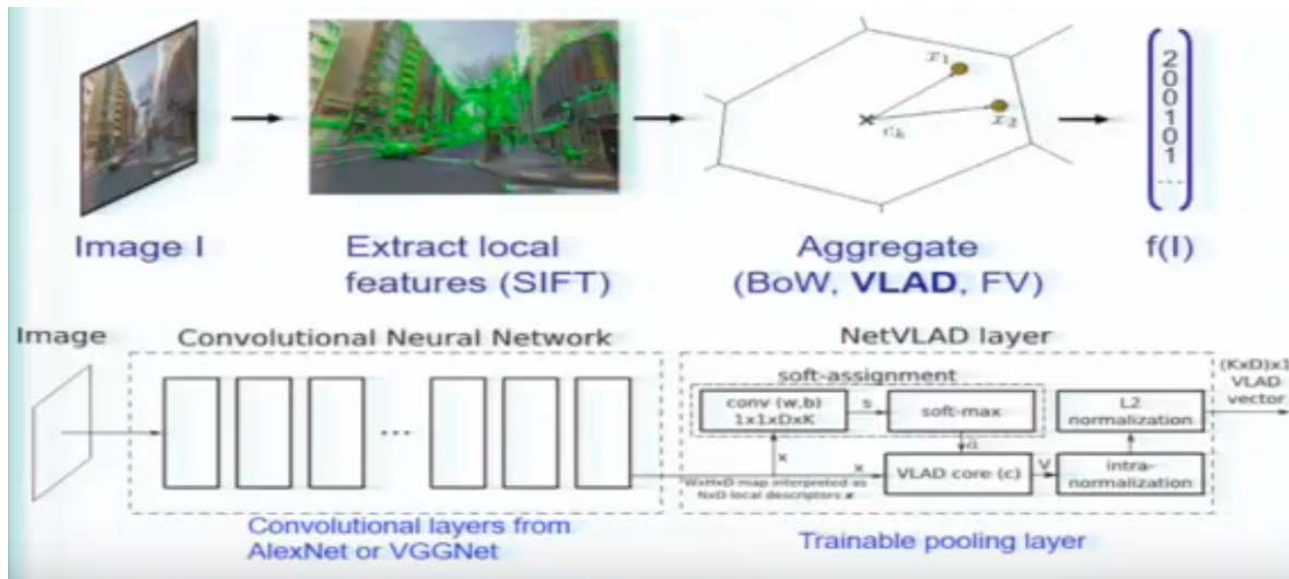
- Identify its location given an query image
- Application of place recognition



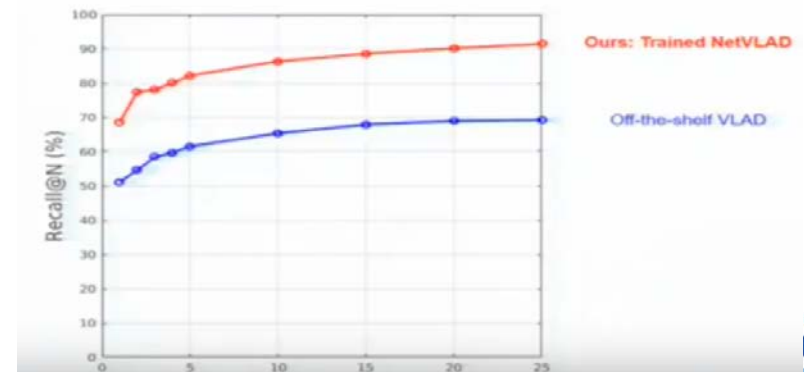
From the author talk

Mimic the classical approach

- Make it end-to-end trainable



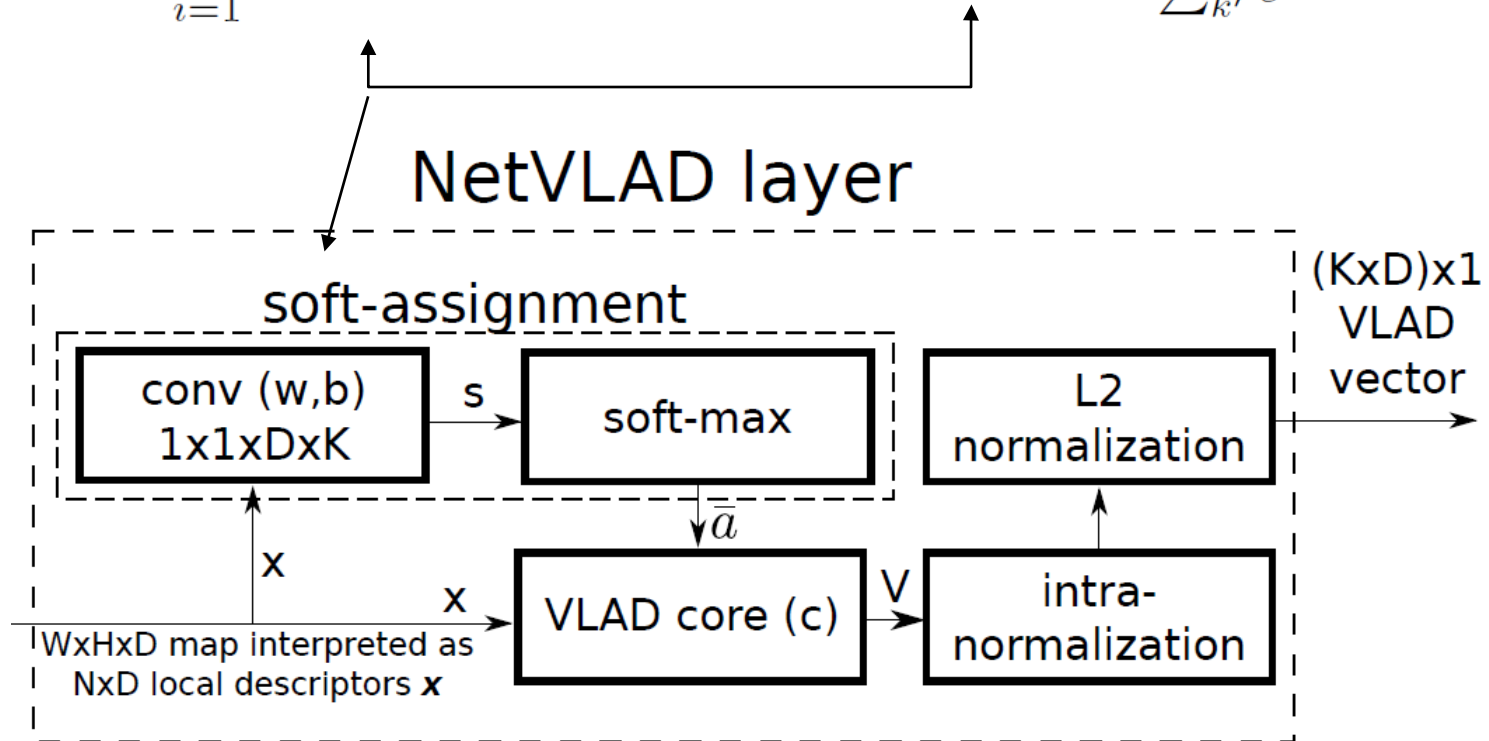
- Training is important



Trainable VLAD

- Hard assignment to soft assignment using the soft-max

$$V(j, k) = \sum_{i=1}^N a_k(\mathbf{x}_i) (x_i(j) - c_k(j)), \quad \bar{a}_k(\mathbf{x}_i) = \frac{e^{-\alpha \|\mathbf{x}_i - \mathbf{c}_k\|^2}}{\sum_{k'} e^{-\alpha \|\mathbf{x}_i - \mathbf{c}_{k'}\|^2}},$$



Problems of BoW Model

- No spatial relationship between words
- How can we perform segmentation and localization?



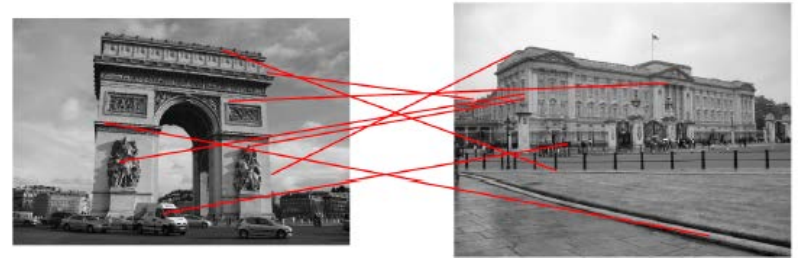
Ack.: Fei-Fei Li

Post-Processing or Reranking



Post-Processing

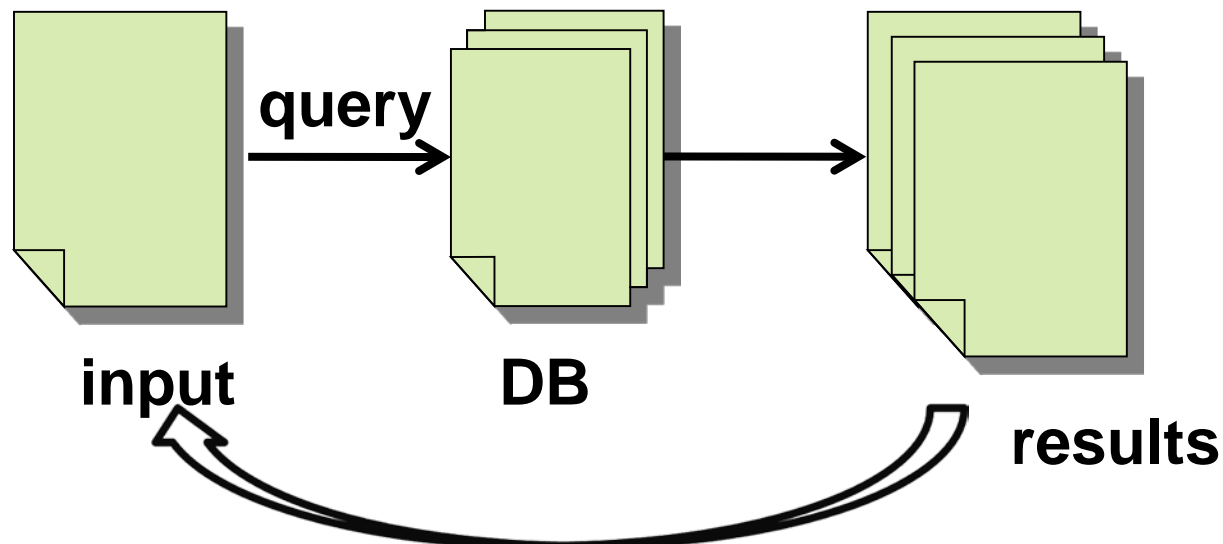
- Geometric verification
 - RANSAC



Matching w/o spatial matching

(Ack: Edward Johns et al.)

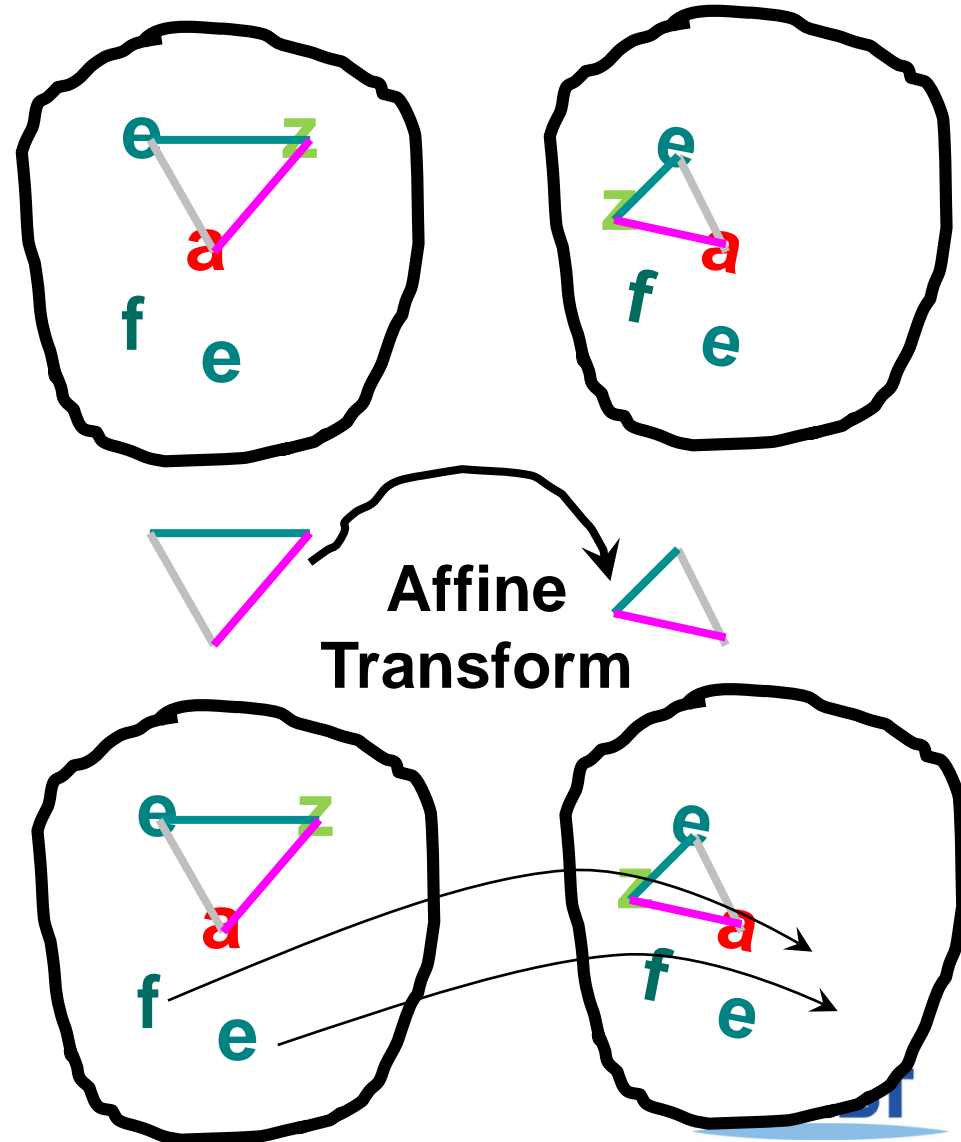
- Query expansion



Geometric Verification using RANSAC for Affine Transform

Repeat N times:

- Randomly choose 3 matching pairs
- Estimate transformation
- Predict remaining points and count “inliers”

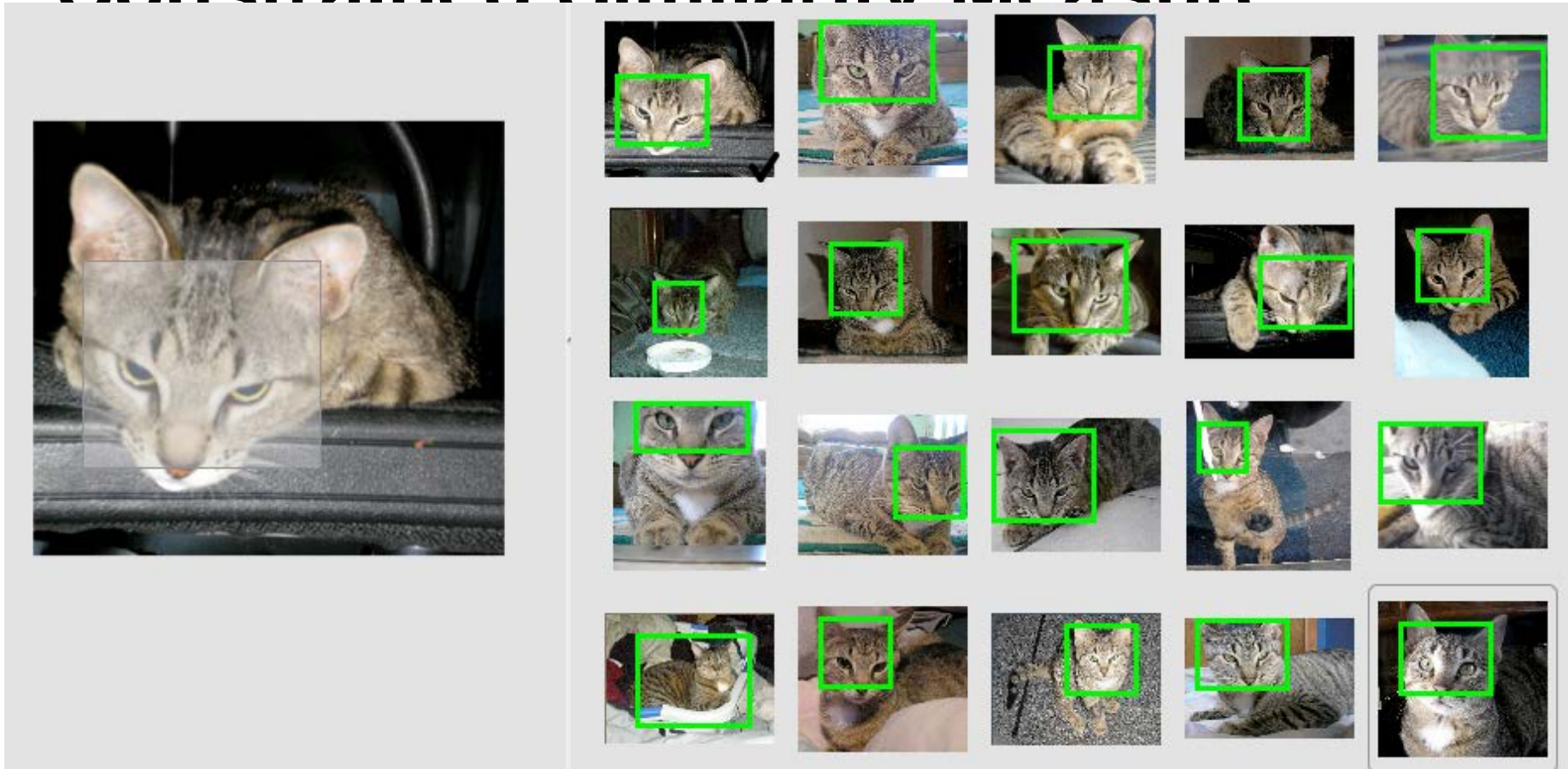


Pattern matching

- Drones surveying city
 - Identify a particular car



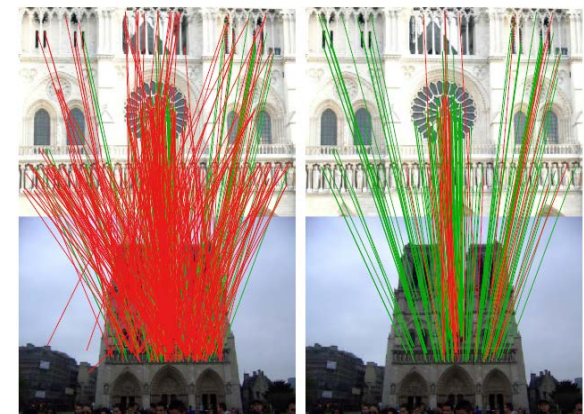
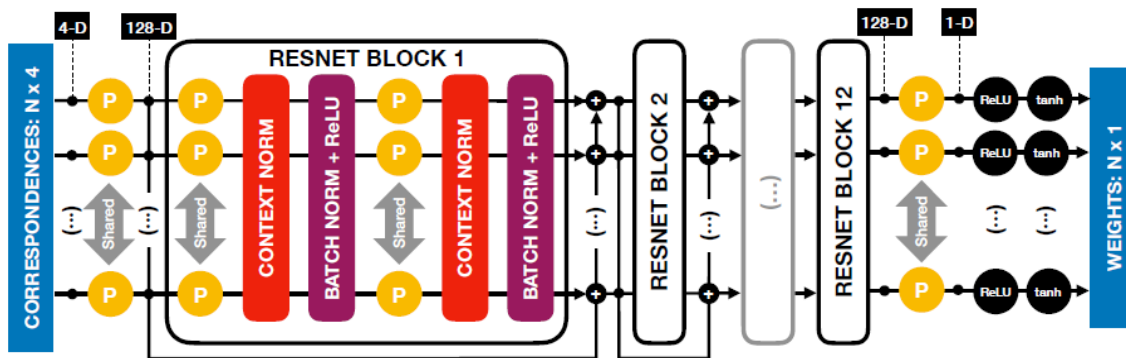
Image Retrieval with Spatially Constrained Similarity Measure



[Xiaohui Shen, Zhe Lin, Jon Brandt, Shai Avidan and Ying Wu, CVPR 2012]

Learning to Find Good Correspondences, CVPR 18

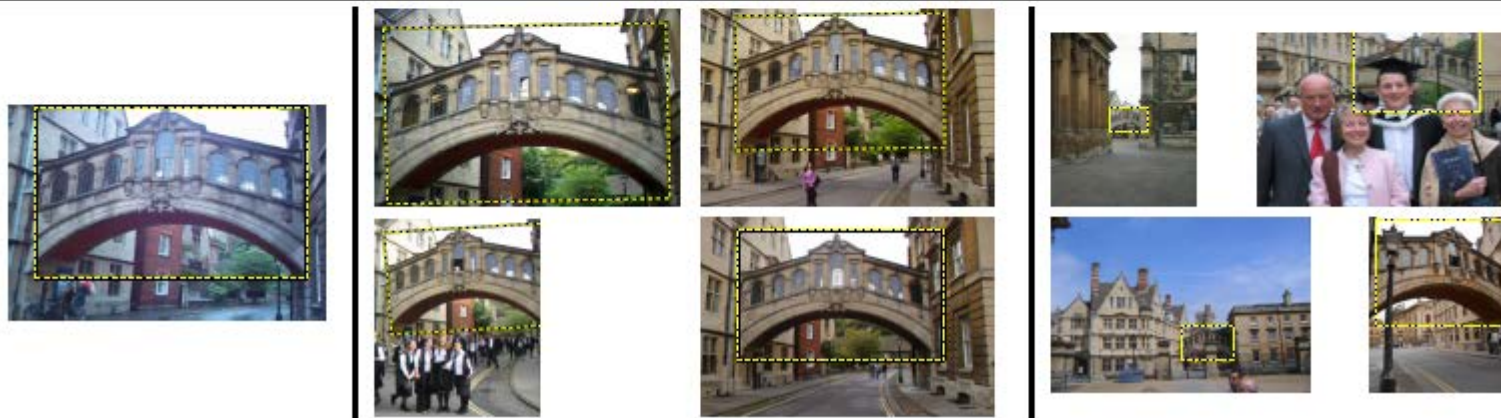
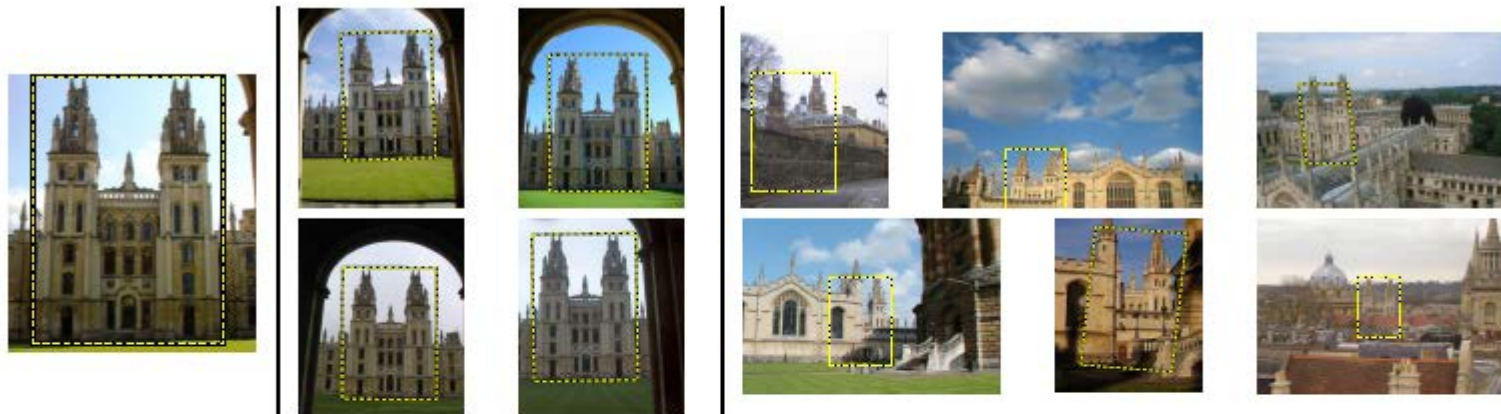
- Given two input features (e.g., SIFTs), return a prob. of being inliers for each feature
 - Adopt the classification approach
 - Additionally perform the regression for pose estimation



(a) RANSAC

(b) Our approach

Query Expansion [Chum et al. 07]



Original query

Top 4 images

Expanded results that were not identified by the original query

Different Outputs of Image Search

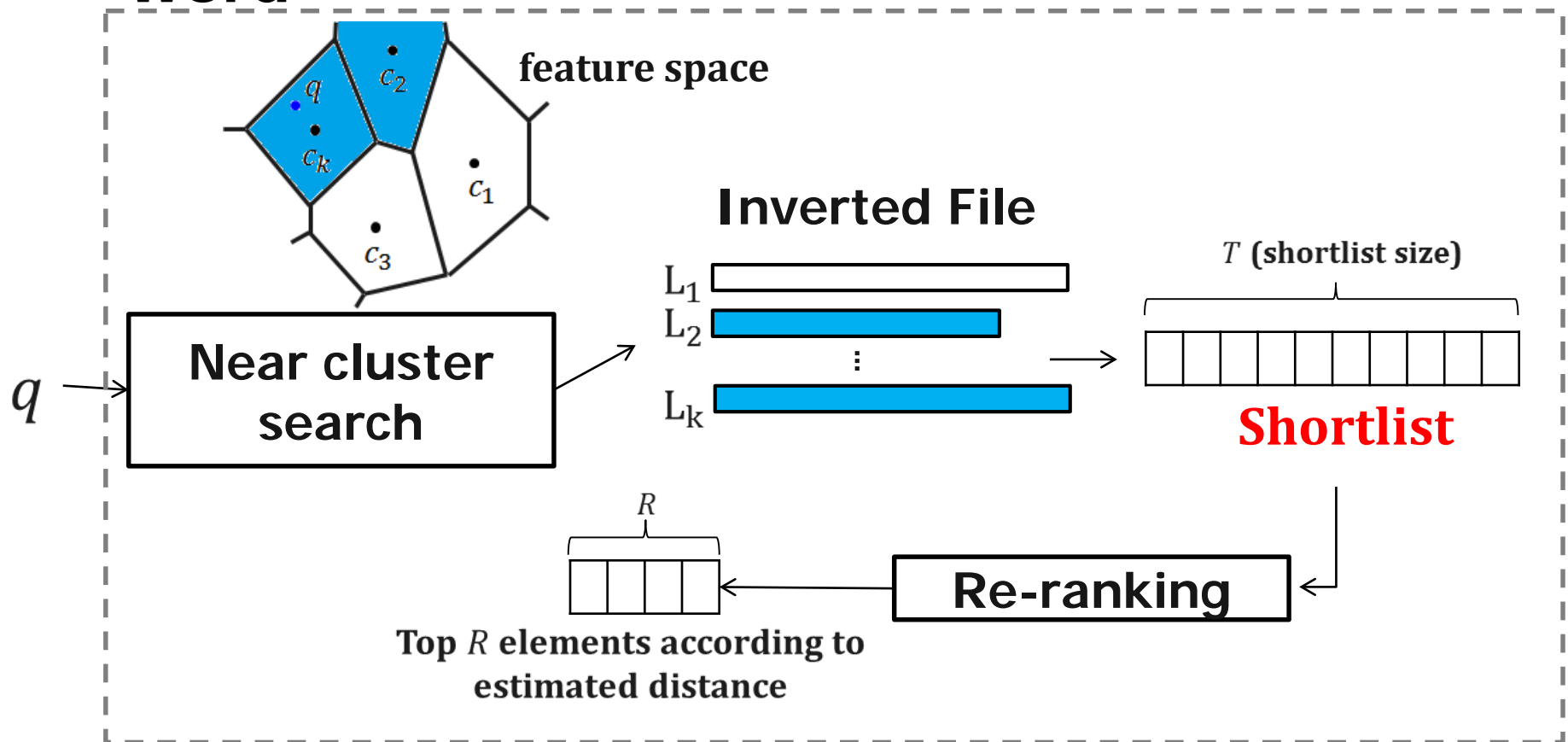
- Generate image sequences from paragraphs



Kim et al., Ranking and Retrieval of Image Sequences from Multiple Paragraph Queries

Inverted File or Index for Efficient Search

- For each word, list images containing the word



Inverted Index

Construction time:

- Generate a codebook by quantization
 - e.g. k-means clustering
- Build an inverted index
 - Quantize each descriptor into the closest word
 - Organize desc. IDs in terms of words

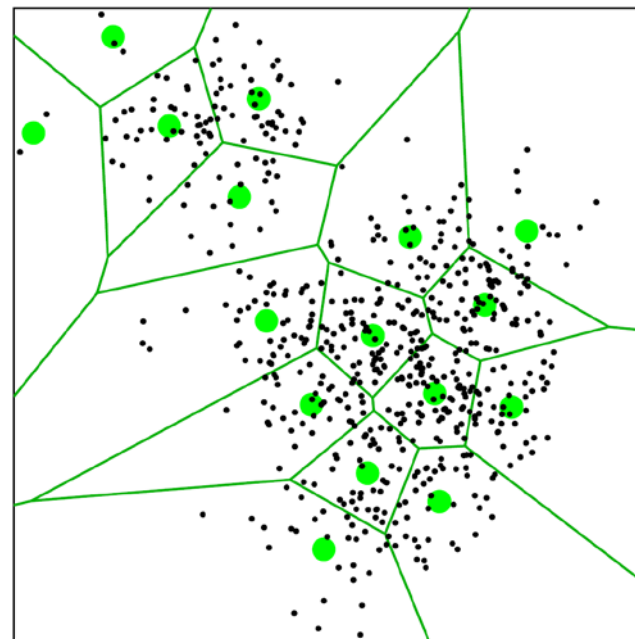
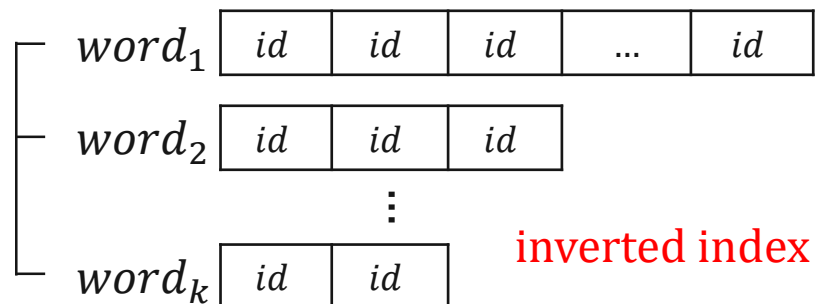


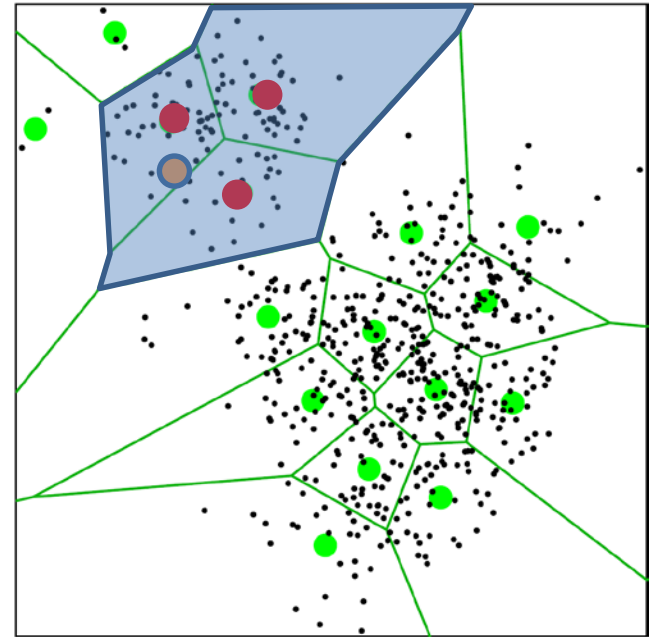
Figure from Lempitsky's slides



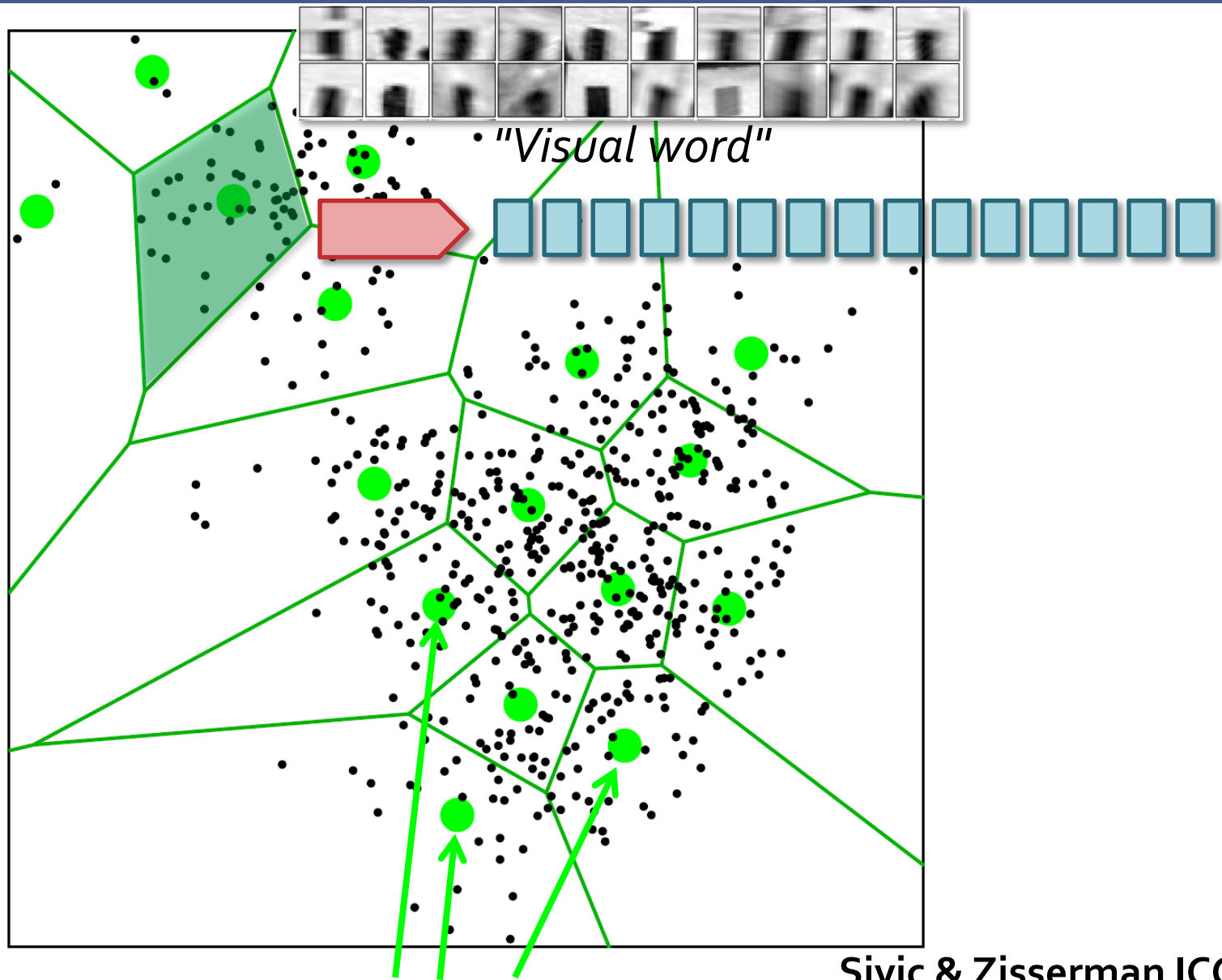
Inverted Index

Query time:

- Given a query,
 - Find its K closest **words**
 - Retrieve all the data in the K lists corresponding to the words
- Large K
 - Low quantization distortion
 - Expensive to find kNN words



The inverted index



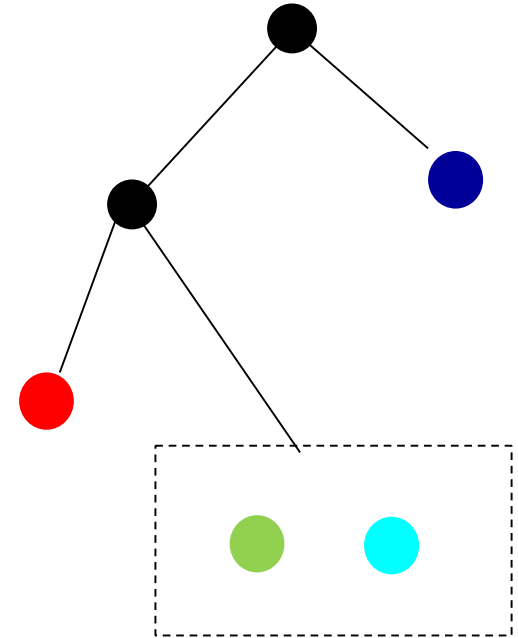
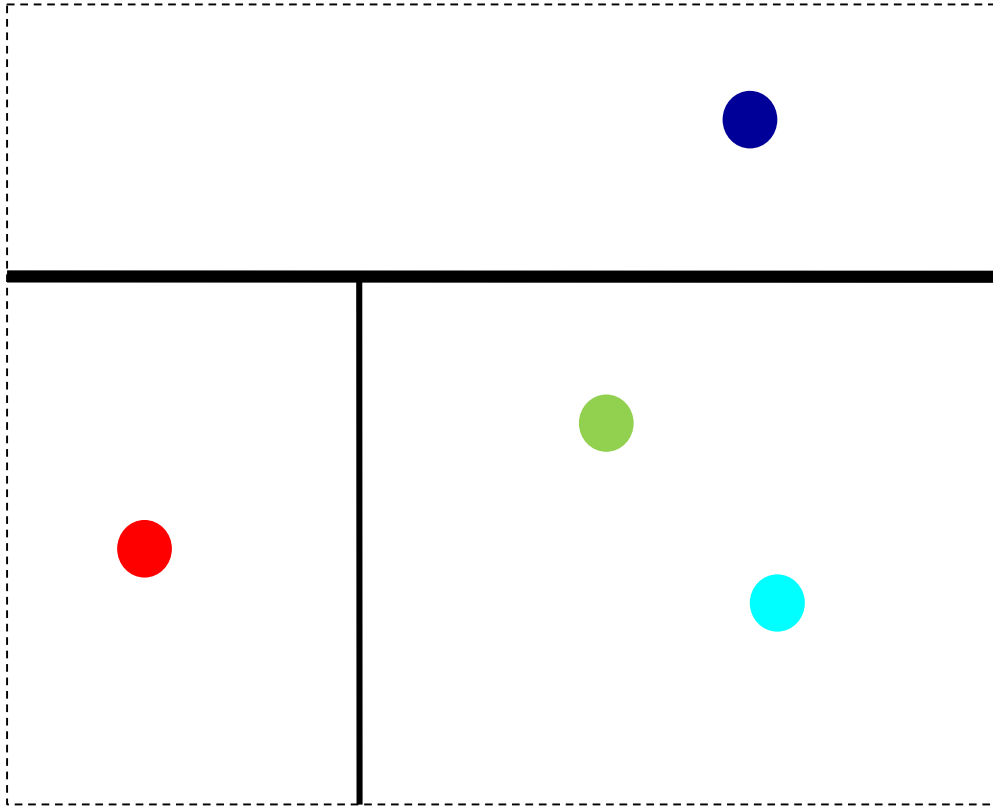
Sivic & Zisserman ICCV 2003

Visual codebook

Approximate Nearest Neighbor (ANN) Search

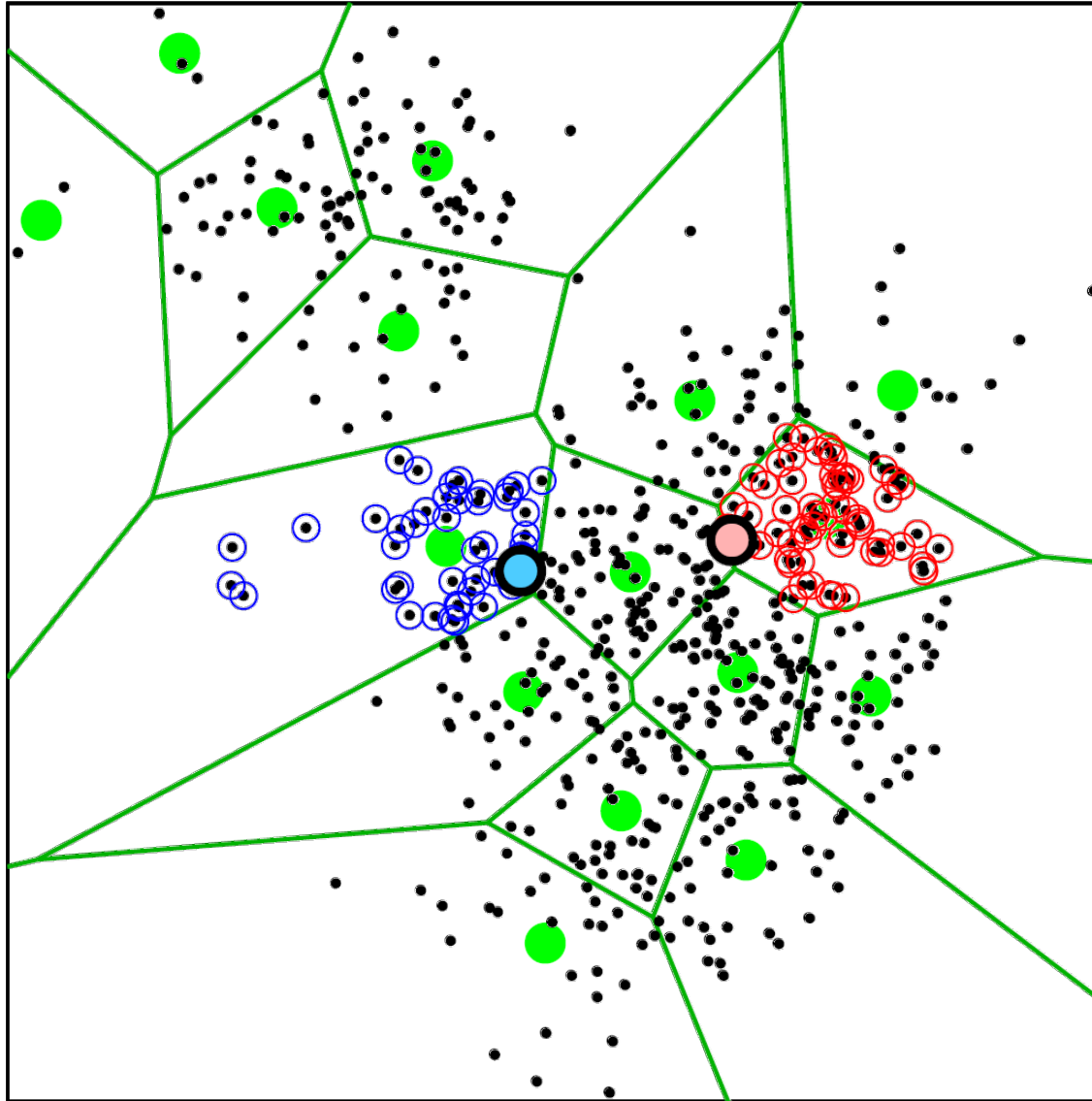
- For large K
 - Takes time to find clusters given the query
 - Use those ANN techniques for efficiently finding near clusters
- ANN search techniques
 - kd-trees: hierarchical approaches for low-dimensional problems
 - Hashing for high dimensional problems; will be discussed later with binary code embedding
 - Quantization (k-means cluster and product quantization)

kd-tree Example

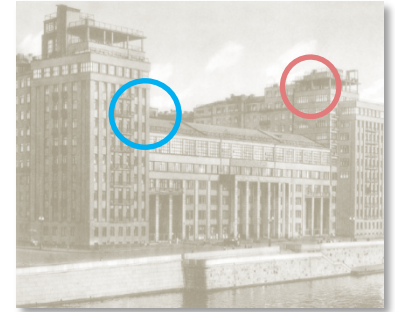


- Many good implementations (e.g., vl-feat)

Querying the inverted index



Query:



- Have to consider several words for best accuracy
- Want to use as big codebook as possible

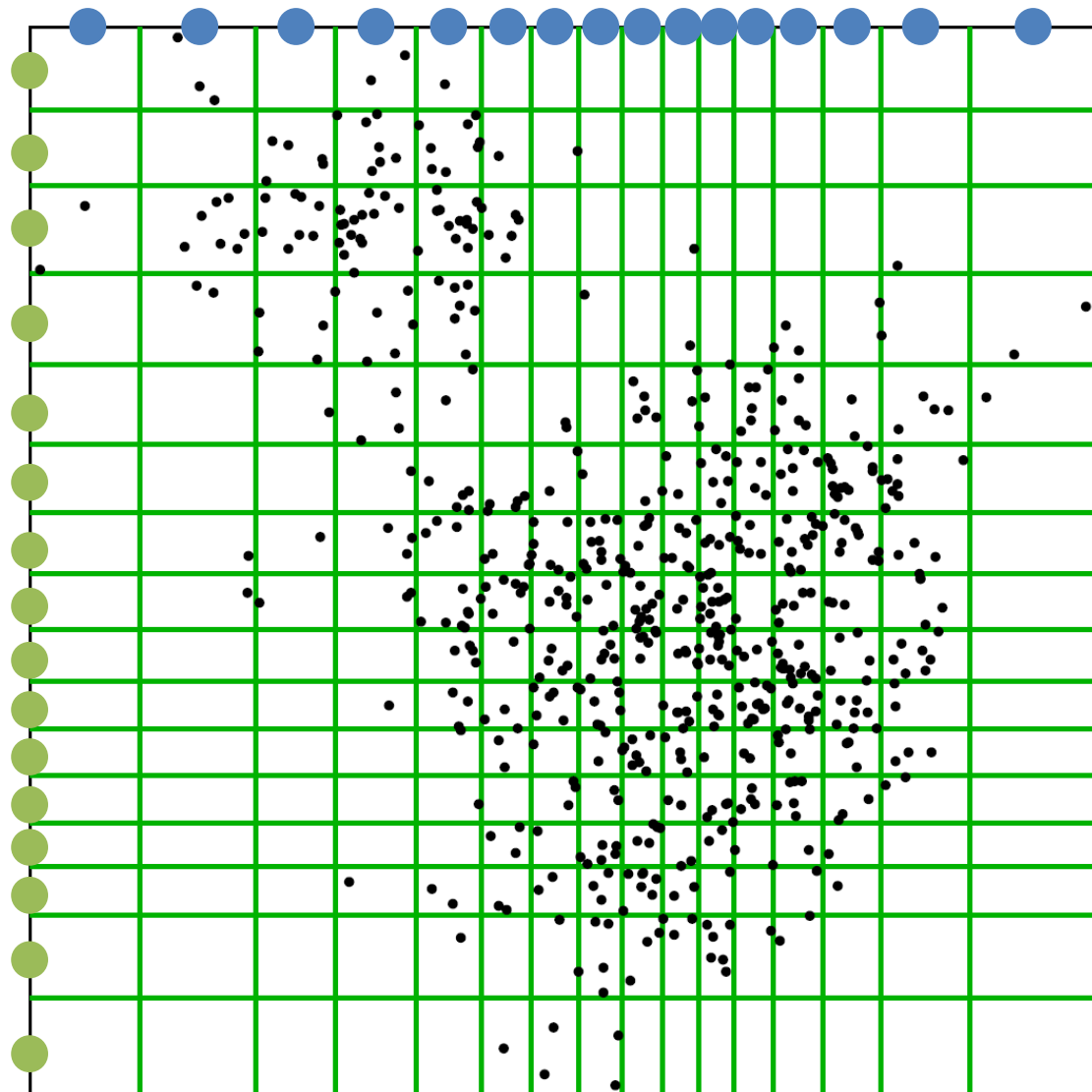


- Want to spend as little time as possible for matching to codebooks

Ack.: Lempitsky

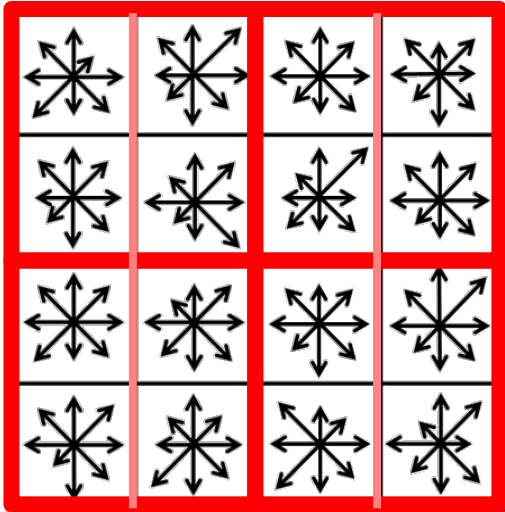
Inverted Multi-Index [Babenko and Lempitsky, CVPR 2012]

- **Product quantization for indexing**
- **Main advantage:**
 - For the same K , much finer subdivision
 - Very efficient in finding k NN codewords



Ack.: Lempitsky

Product quantization

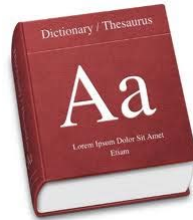


1. Split vector into correlated subvectors
2. use separate small codebook for each chunk

Quantization vs. Product quantization:

For a budget of 4 bytes per descriptor:

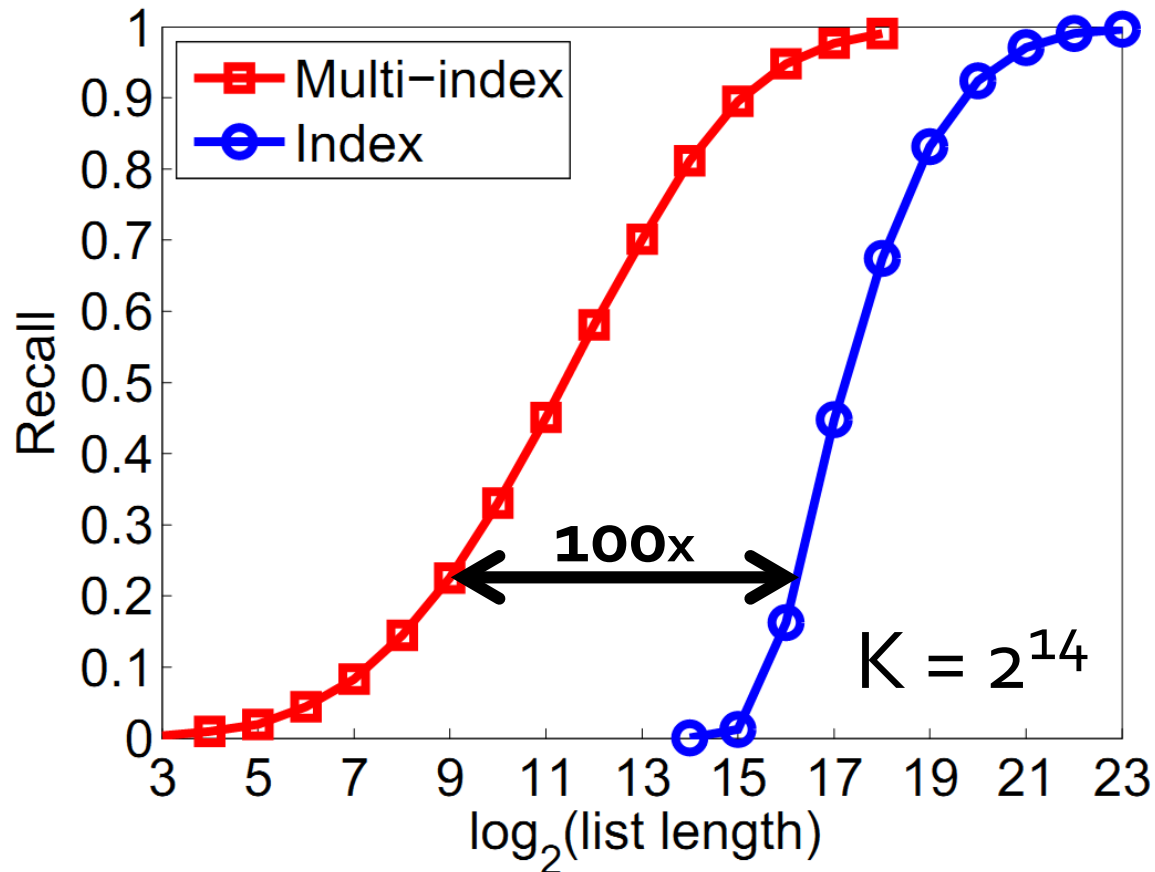
1. Use a single codebook with 1 billion codewords or
2. Use 4 different codebooks with 256 codewords each



many minutes **128GB**

< 1 millisecond **32KB**

Performance comparison on 1 B SIFT descriptors

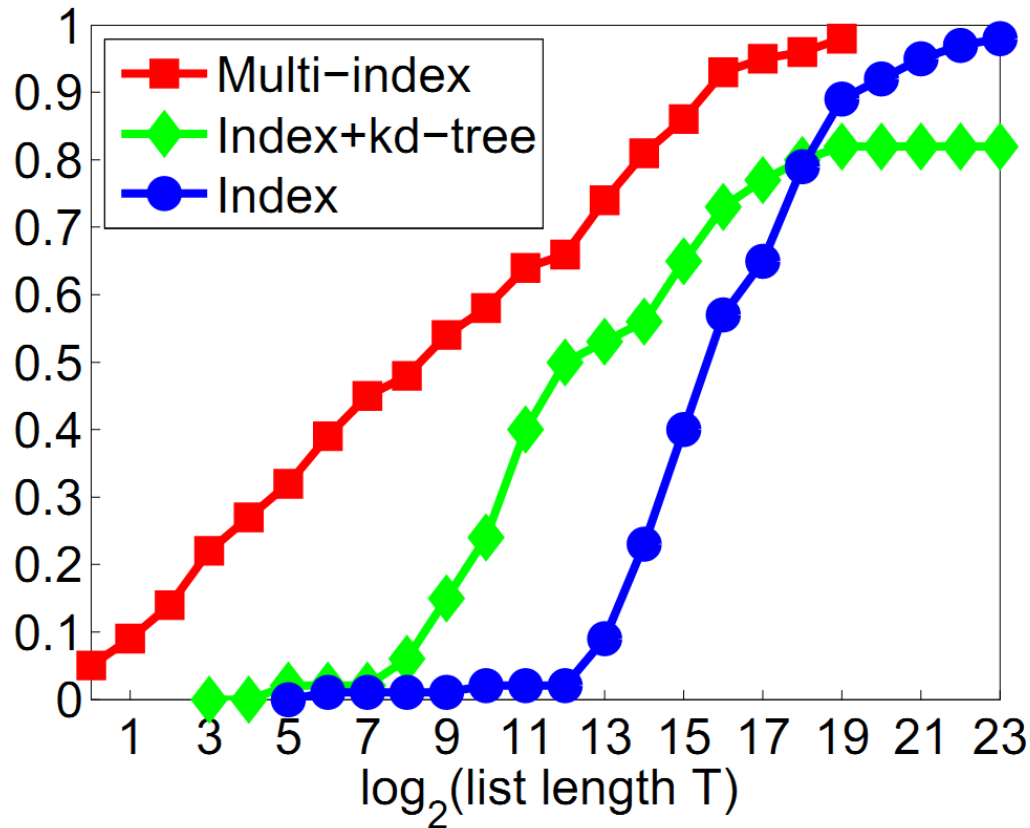


Time increase: 1.4 msec -> 2.2 msec on a single core
(with BLAS instructions)

Ack.: Lempitsky

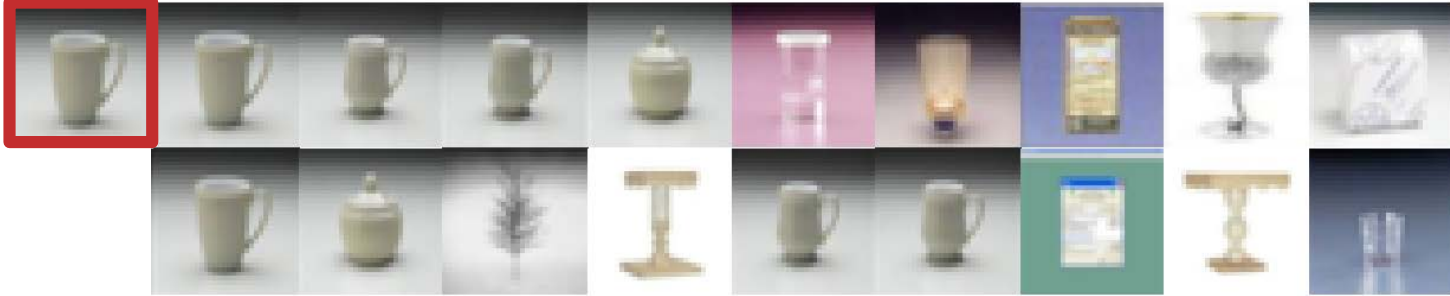
Performance on 80 million GISTs

Index vs Multi-index:



Tests on 80 million GISTs (384 dimensions) of Tiny Images [Torralba et al. PAMI'o8]

Retrieval examples



Exact NN
Uncompressed GIST

Multi-D-ADC
16 bytes



Exact NN
Uncompressed GIST

Multi-D-ADC
16 bytes



Exact NN
Uncompressed GIST

Multi-D-ADC
16 bytes



Exact NN
Uncompressed GIST

Multi-D-ADC
16 bytes

Ack.: Lempitsky

Scalability

- **Issues with billions of images?**
 - Searching speed → inverted index
 - Accuracy → larger codebooks, spatial verification, expansion, features
 - Memory → compact representations
 - Easy to use?
 - Applications?
 - A new aspect?

Class Objectives were:

- Bag-of-visual-Word (BoW) model
- Understand approximate nearest neighbor search
 - Inverted index
 - Inverted multi-index

Next Time...

- Learning techniques

Homework for Every Class

- Go over the next lecture slides
- Come up with one question on what we have discussed today
 - 1 for typical questions (that were answered in the class)
 - 2 for questions with thoughts or that surprised me
- Write questions at least 4 times