

---

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**



# Tentative Schedule

---

● **10/25: mid-term exam**

**10/27: Student pre. I**

**11/1**

**11/3**

**11/8**

**11/10: mid-pre**

**11/15: mid-pre**

**11/17: no class**

**11/22: Student pre. II**

**11/24**

**11/29**

**12/1**

**12/6**

**12/8: reserved**

**12/13: final pre.**

**12/15: final pre.**

**D: 10/4 (student schedule)**

# Deadlines

---

- **Declare project team members**
  - **By 10/3 at Noah**
- **Confirm schedules of paper talks and project talks at 10/4**
- **Declare two papers for student presentations**
  - **by 10/17 at Noah**
  - **Discuss them at the class of 10/18**

# Class Objectives

---

- **Bag-of-visual-Word (BoW) model**
- **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

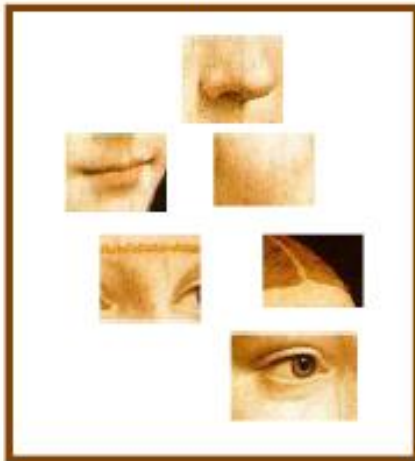
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the visual image was considered as a picture on the retina. The visual centers in the brain are a movie screen. The visual image is a photograph on the screen. The discovery of the visual cortex, 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 to \$750bn, compared with \$570bn in 2004. The surplus is \$660bn. The surplus will be used to buy US goods. It will annoy the US. China's surplus is a deliberate policy. China's government agrees to buy US goods. The yuan is undervalued. The government also needs to buy US goods. The demand so far is \$100bn. The country. China's surplus is the yuan against the dollar. The surplus is permitted it to trade within a narrow range 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

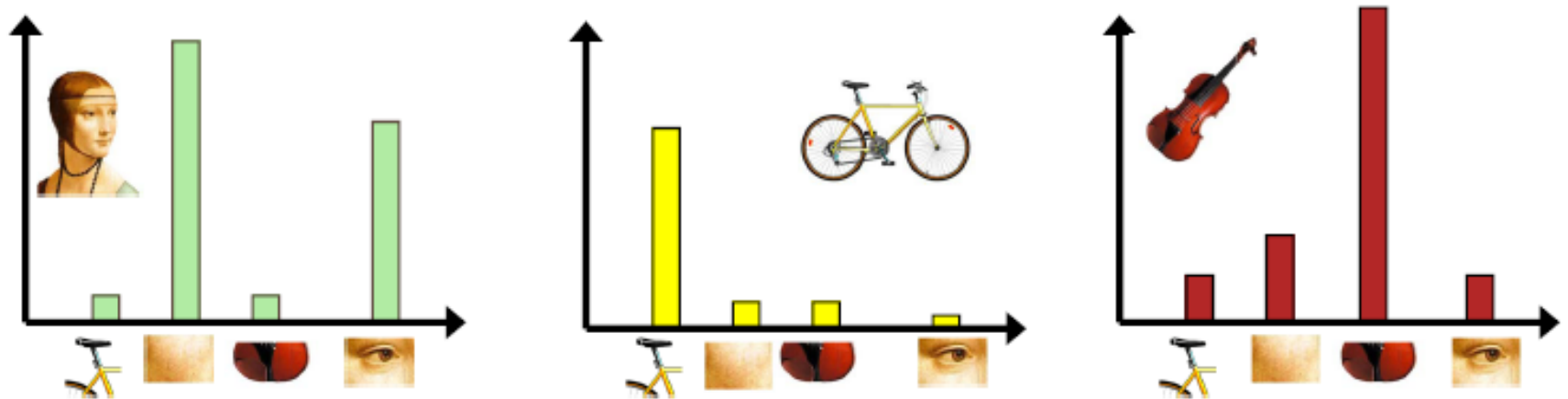


violin



# definition of “BoW”

- Independent features
- histogram representation



codewords dictionary

# Representation



feature detection & representation

codewords dictionary



image representation



learning

category models (and/or) classifiers

# recognition

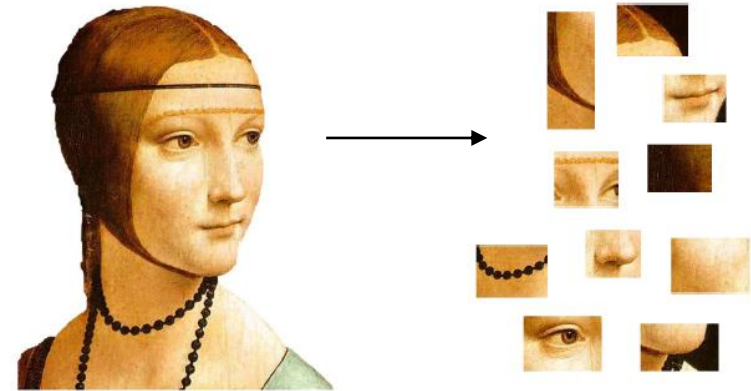


category decision



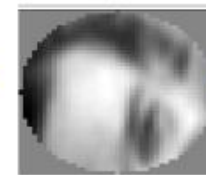
# 1. Feature Detection and Representations

- 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

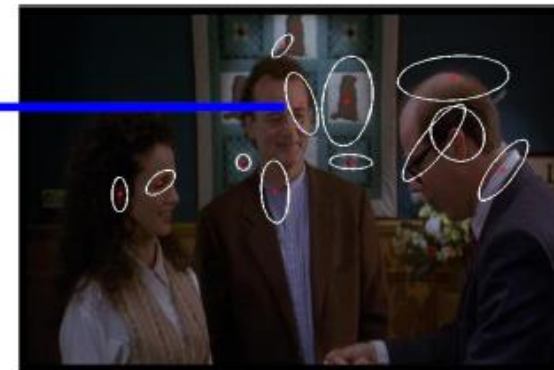


Ack.: Josef Sivic and Li Fei-Fei

  
Compute  
SIFT  
descriptor  
[Lowe'99]



Normalize  
patch

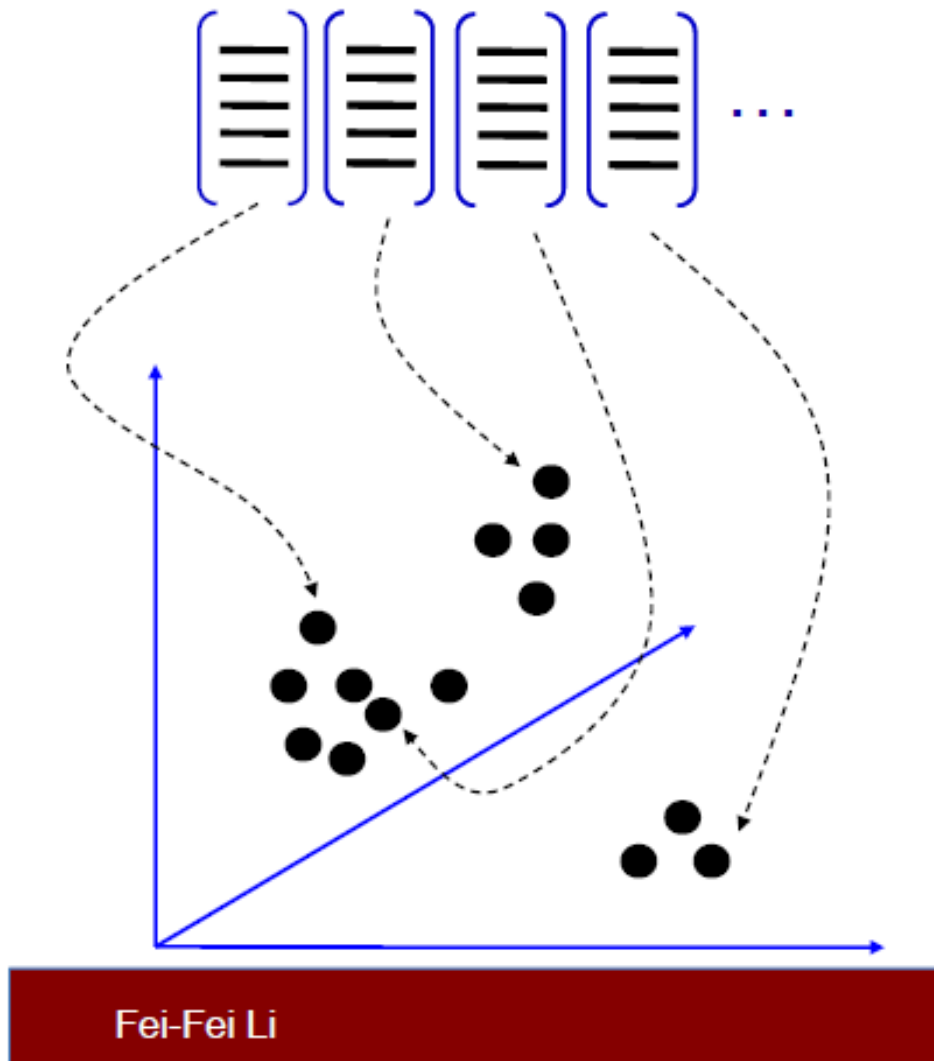


Detect patches

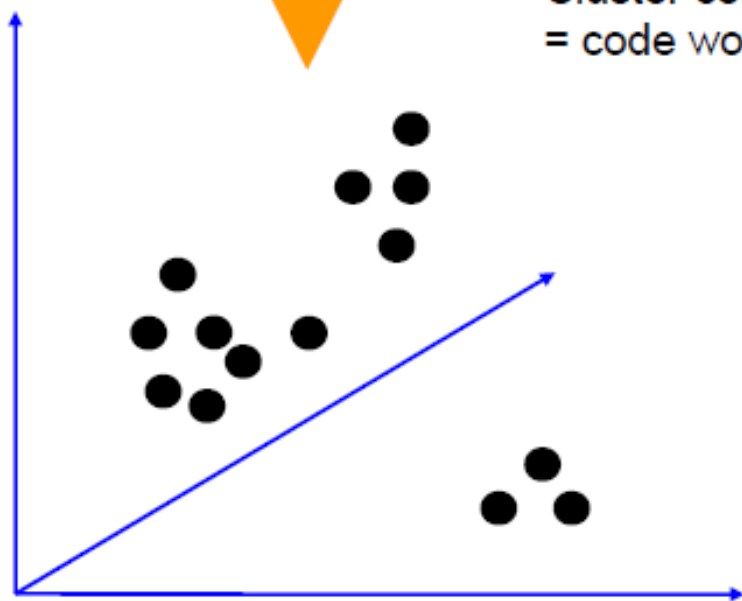
[Mikojczyk and Schmid '02]

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

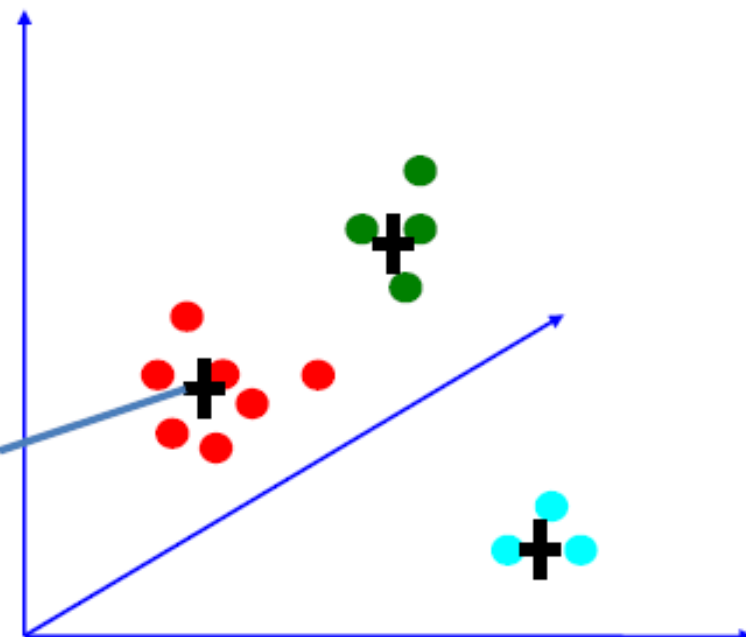
## 2. Codewords dictionary formation



## 2. Codewords dictionary formation



Cluster center  
= code word



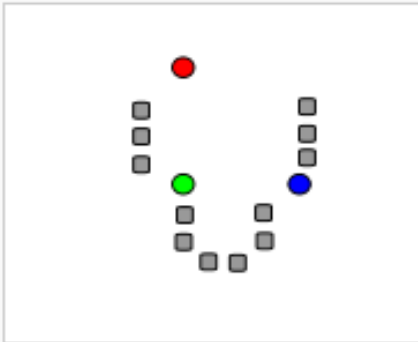
Clustering/  
vector quantization

# K-Means Clustering

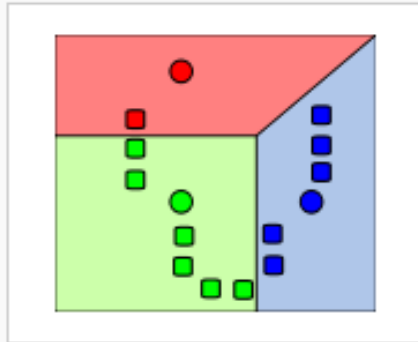
- An unsupervised learning
- Minimize the within-cluster sum of squares

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

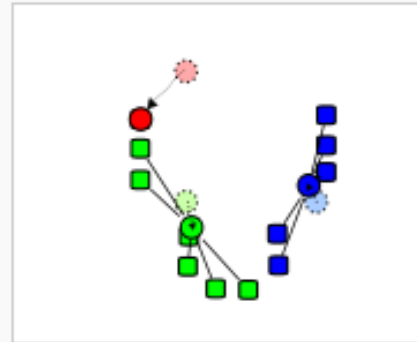
Demonstration of the standard algorithm



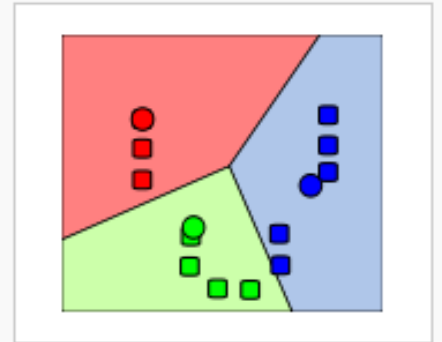
1)  $k$  initial "means" (in this case  $k=3$ ) are randomly selected from the data set (shown in color).



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



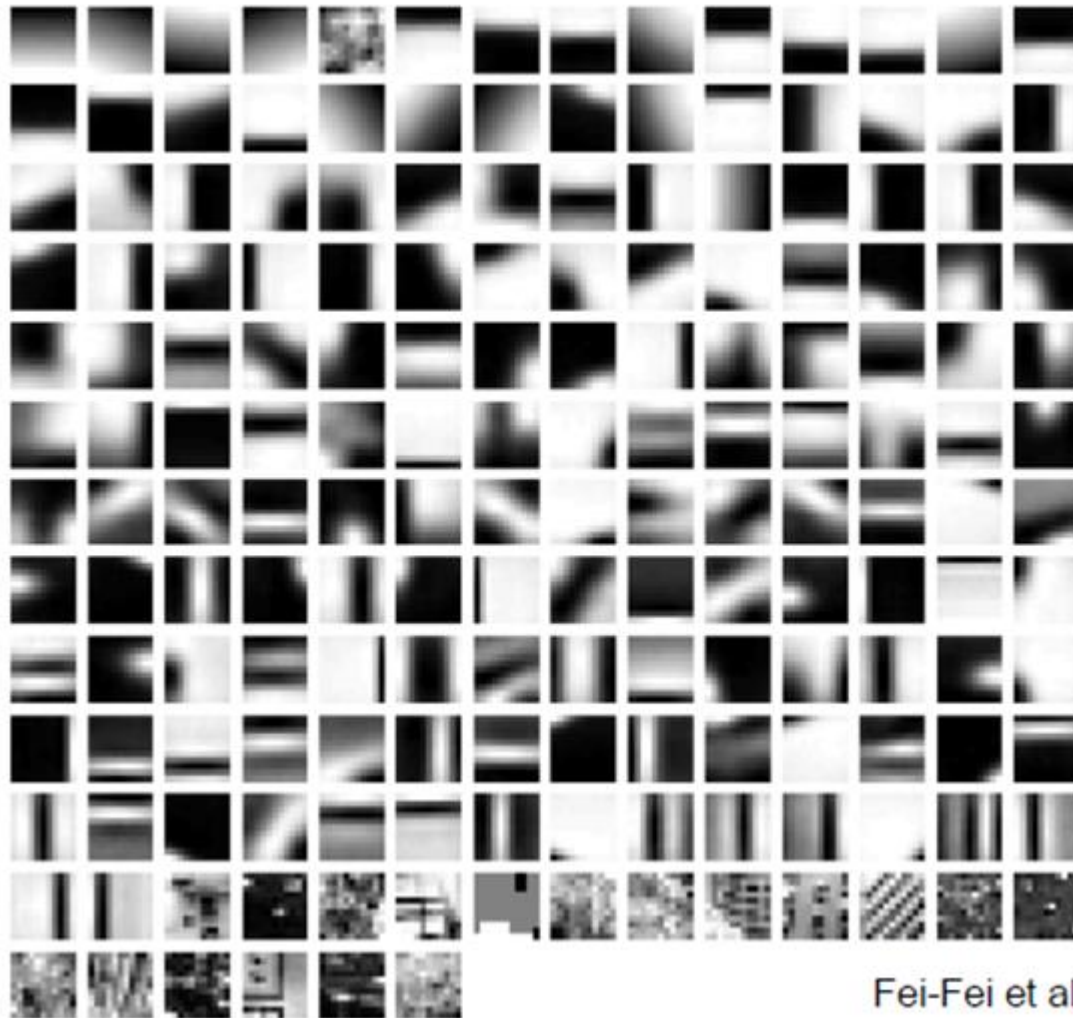
3) The [centroid](#) of each of the  $k$  clusters becomes the new means.



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

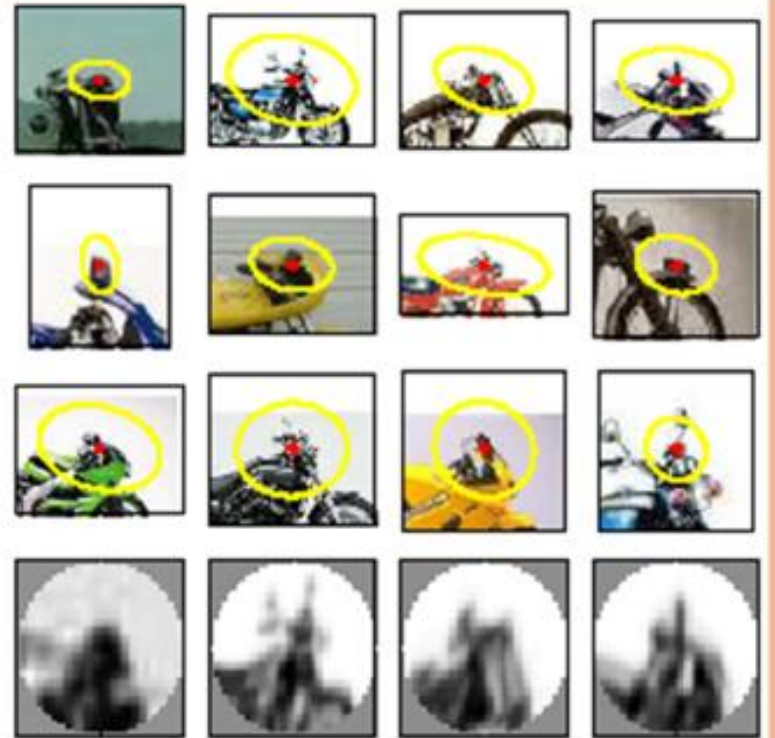
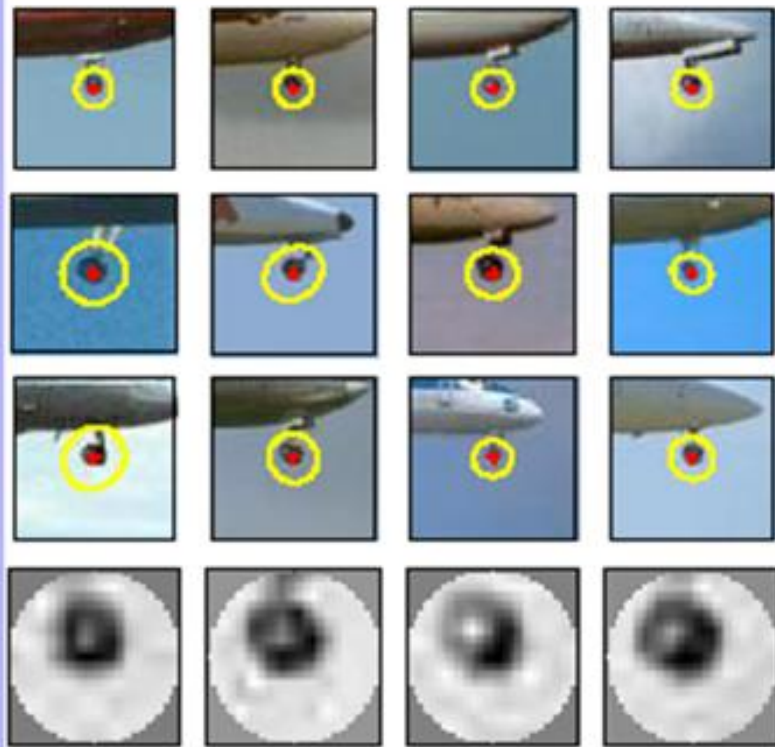
# Codewords Dictionary Formation

---



Fei-Fei et al. 2005

# Image patch examples of codewords

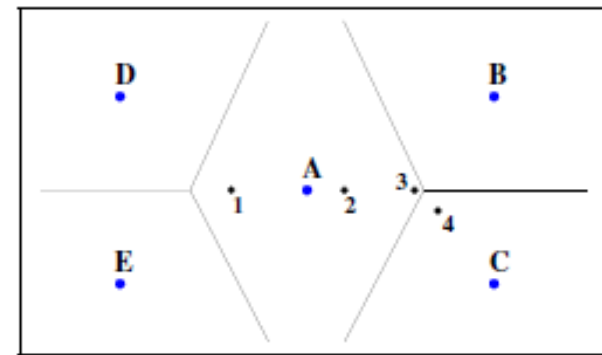
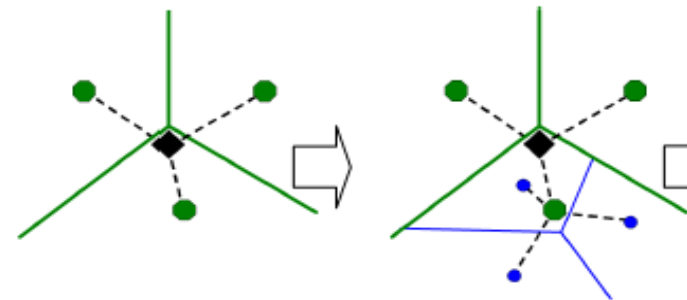


# Issues of Visual Vocabulary

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

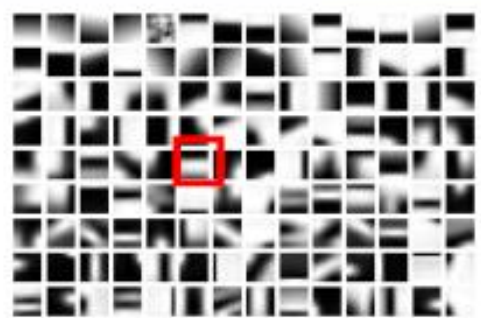
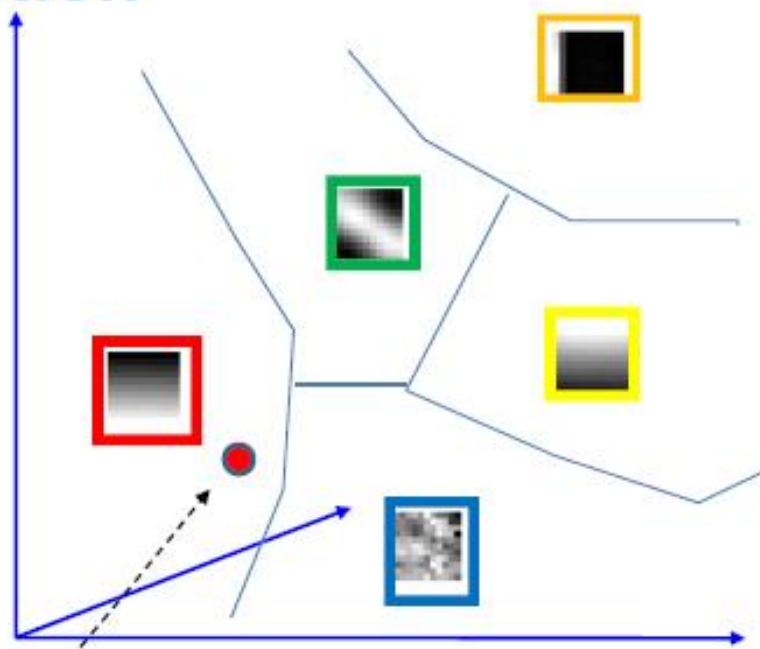
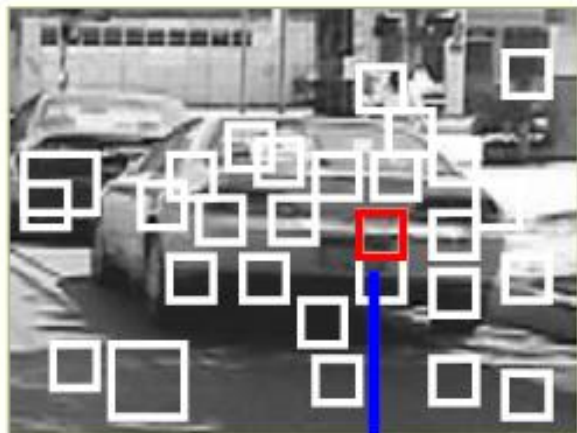
- **Alternatives**

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





### 3. Bag of word representation

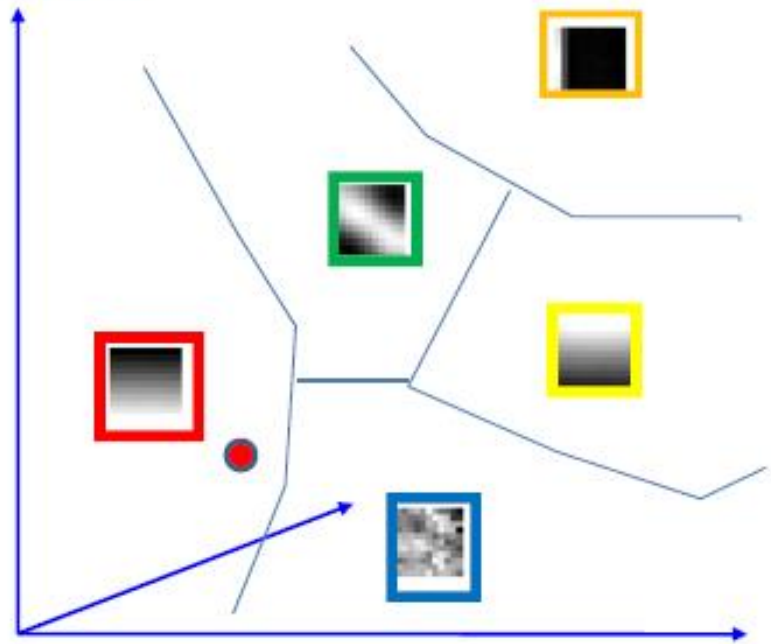


Codewords dictionary

- Nearest neighbors assignment
- K-D tree search strategy



### 3. Bag of word representation



Codewords dictionary

# Representation



1. feature detection & representation



2. codewords dictionary

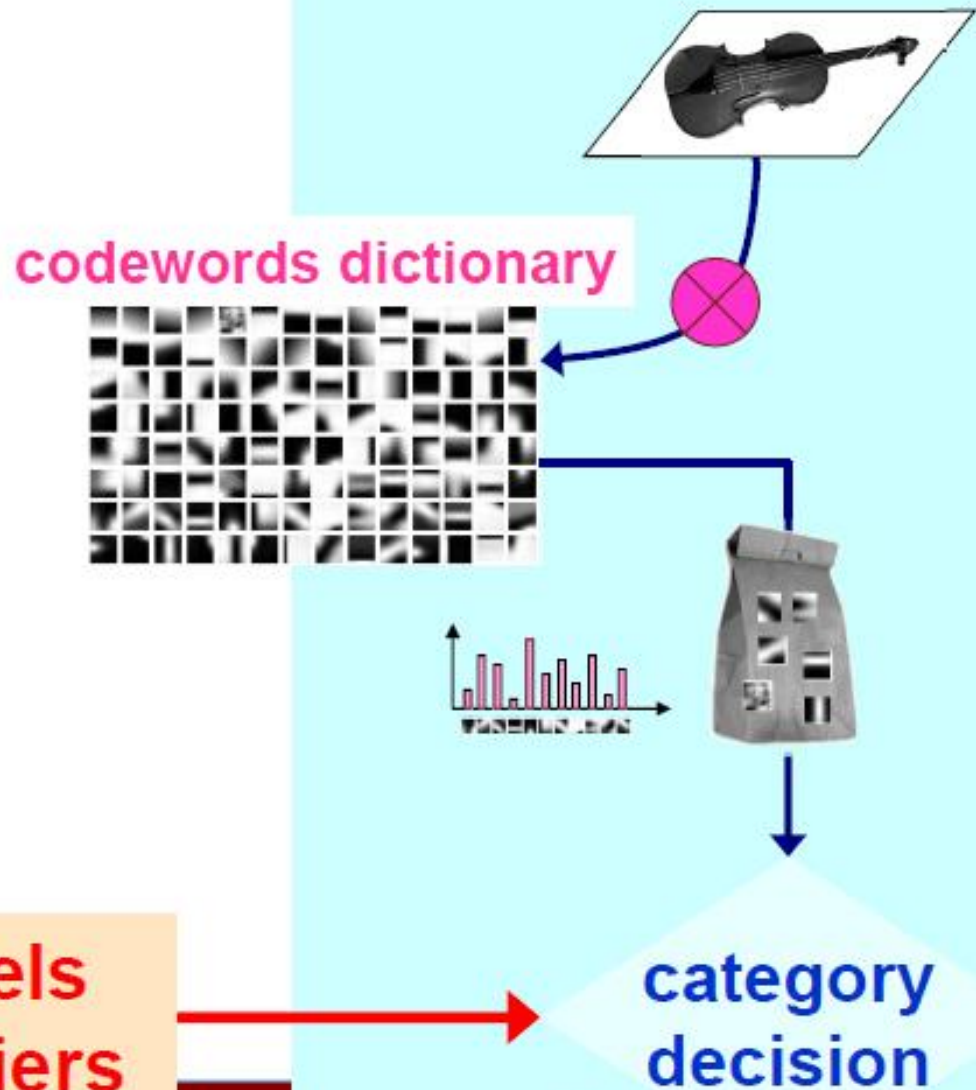


image representation

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**

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

- $\chi^2$  distance

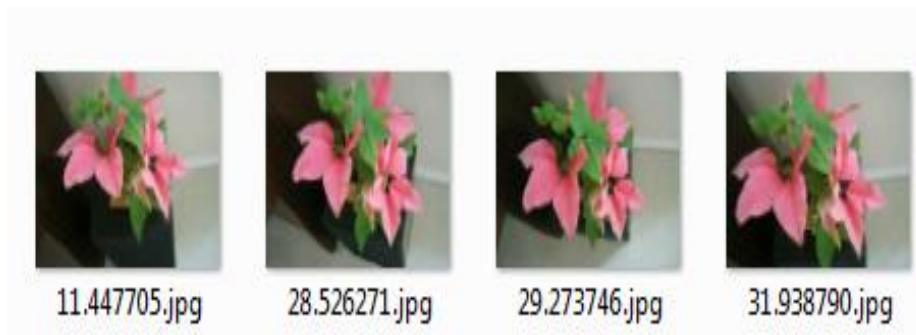
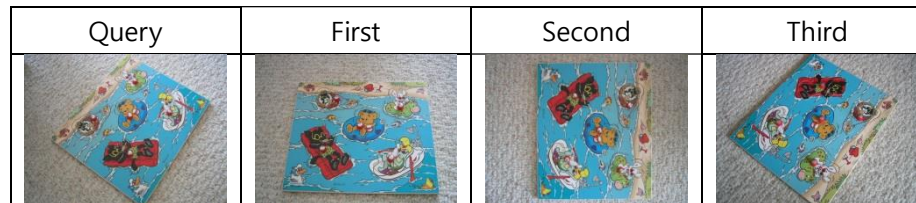
$$D(h_1, h_2) = \sum_{i=1}^N \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$$

- Quadratic distance (*cross-bin*)

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

# PA2

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



# VLAD (Vector of Locally Aggregated Descriptors)

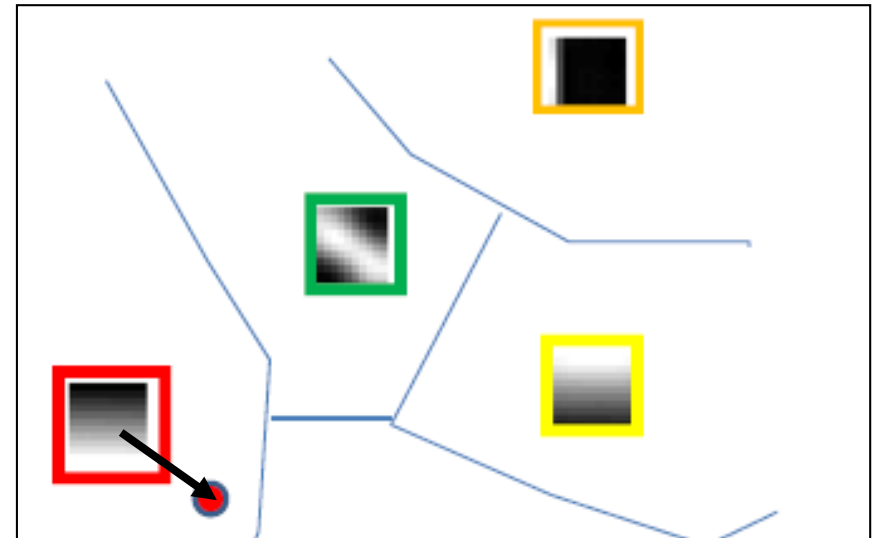
- **BoW**

- Count the number of SIFTs assigned to each cluster

- **VLAD**

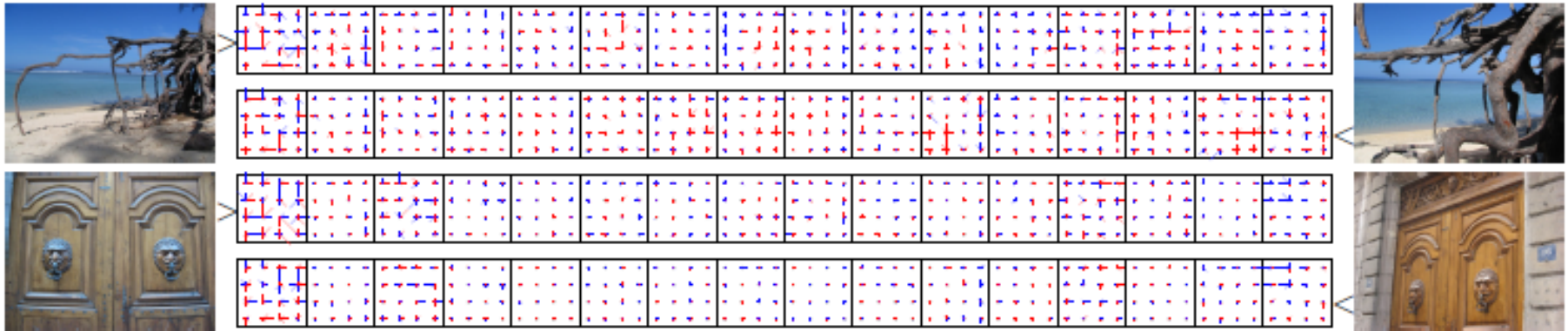
- Compute the difference between a SIFT and its cluster center

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





# VLAD



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

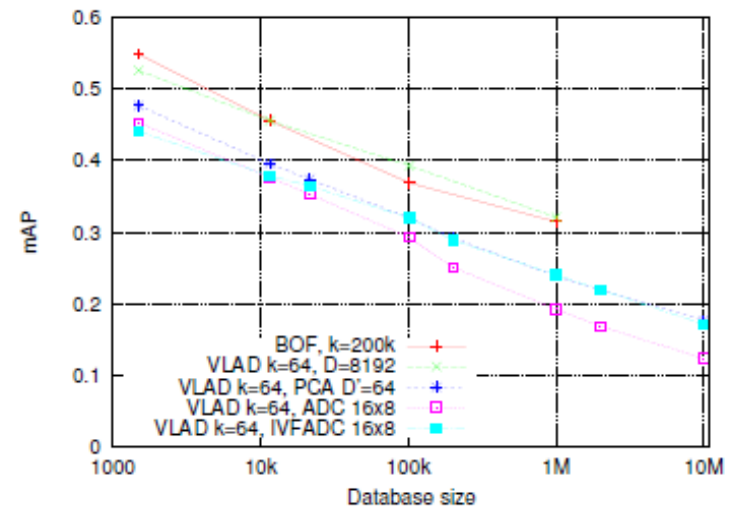
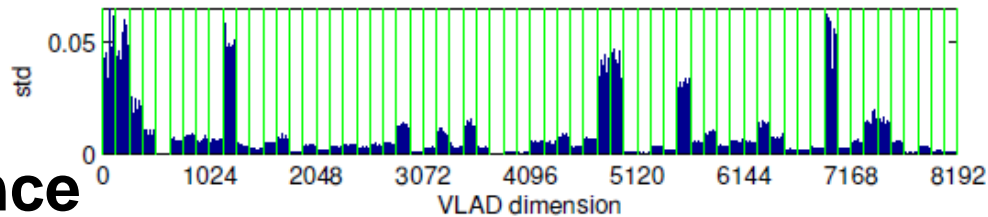


Figure 5. Search accuracy as a function of the database size.

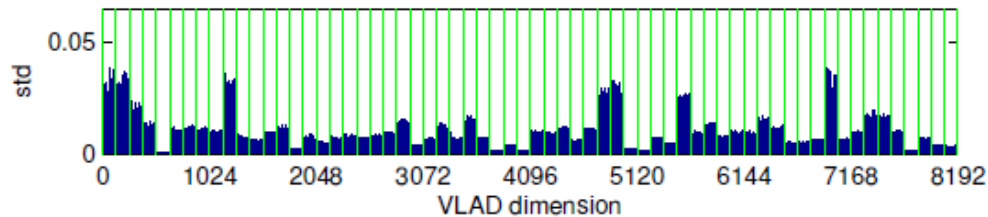


# Normalization for VLAD

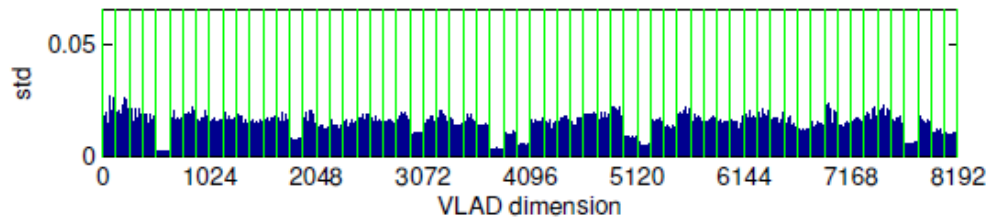
- Results in better accuracy



(a) Original VLAD normalization (L2)



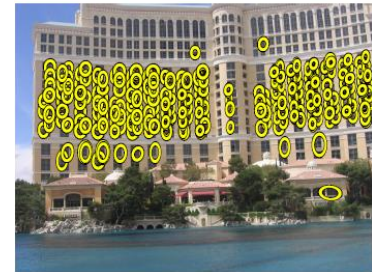
(b) Signed square rooting (SSR) followed by L2



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

L2 normalization,  
i.e.,  $\frac{v}{|v|^2}$

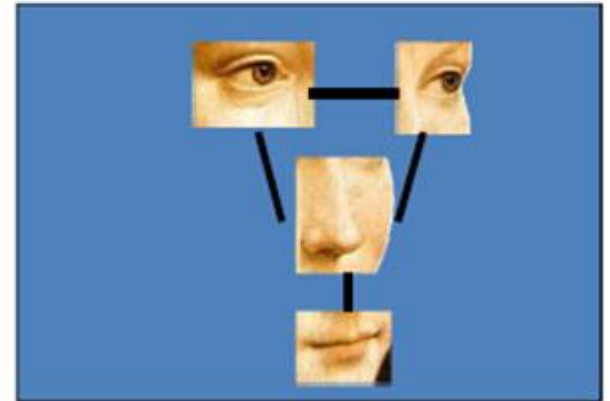
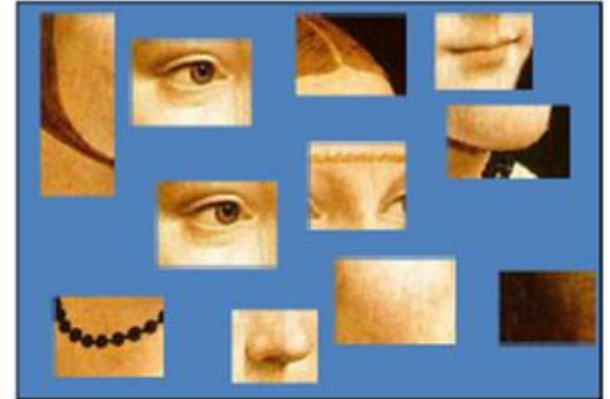
Square rooting  
for burstiness



L2 normalization  
within each VLAD  
block

# Problems of BoW Model

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



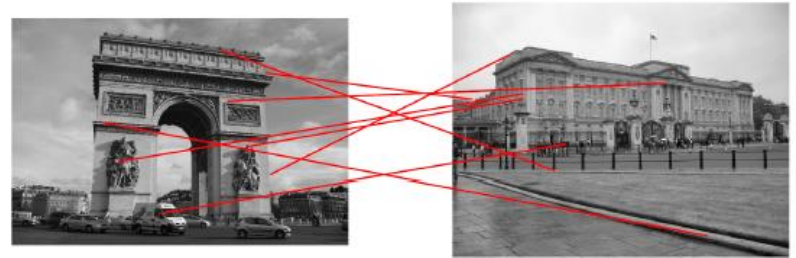
Ack.: Fei-Fei Li

# Post-Processing or Reranking



# Post-Processing

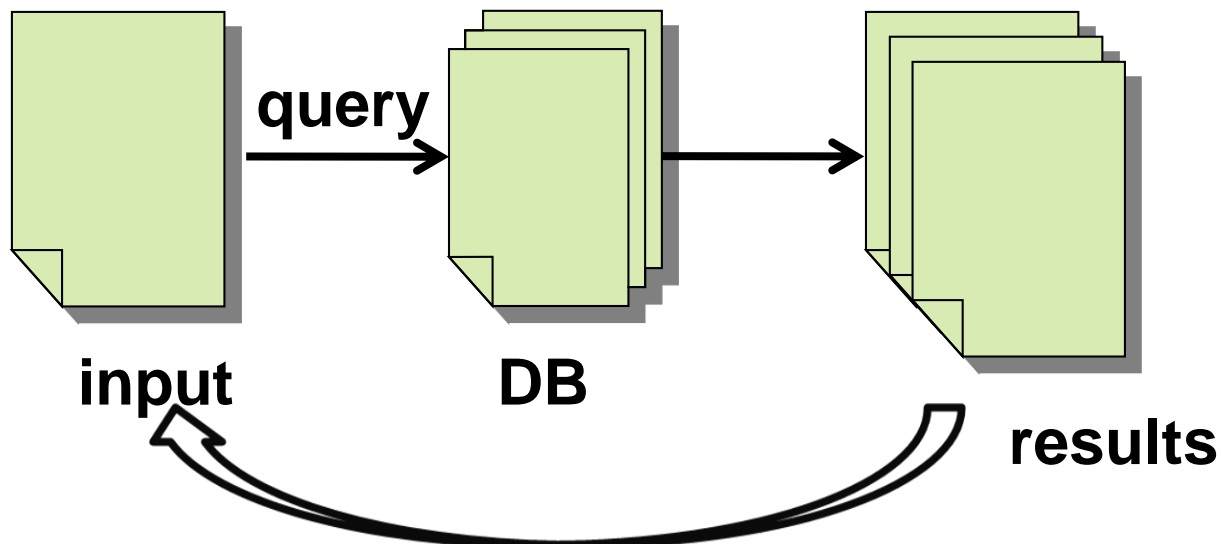
- Geometric verification
  - RANSAC



Matching w/o spatial matching

(Ack: Edward Johns et al.)

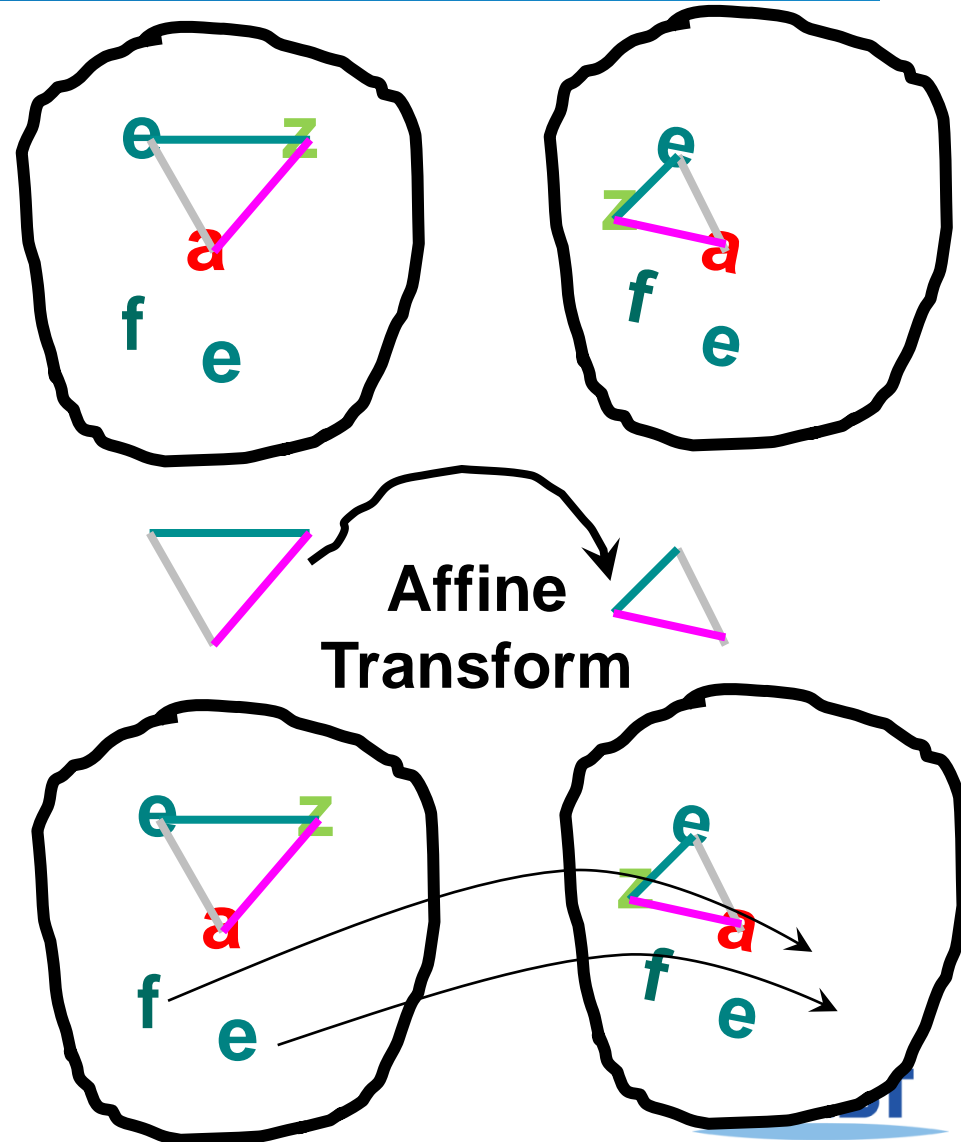
- Query expansion



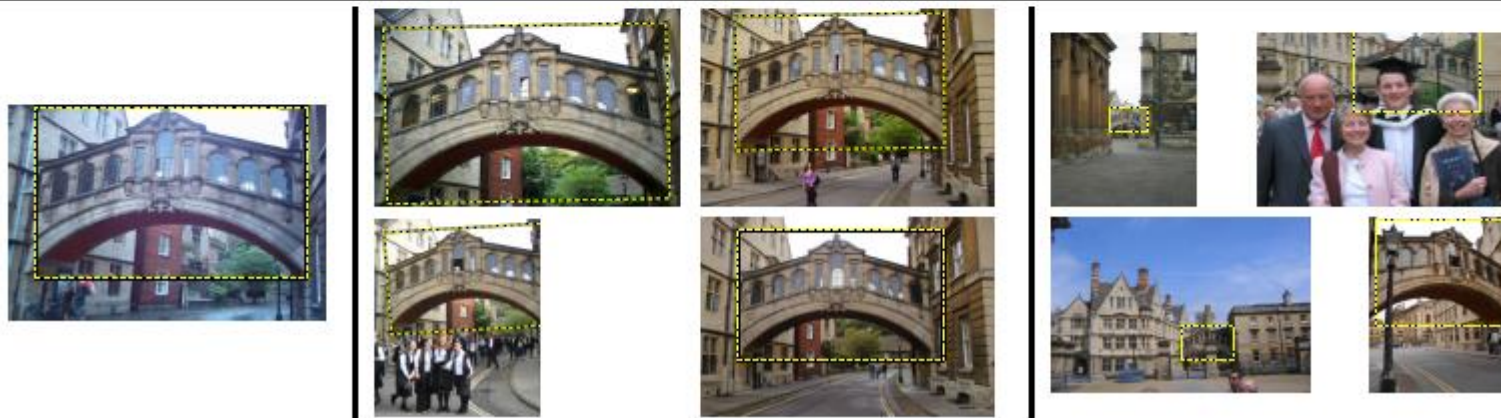
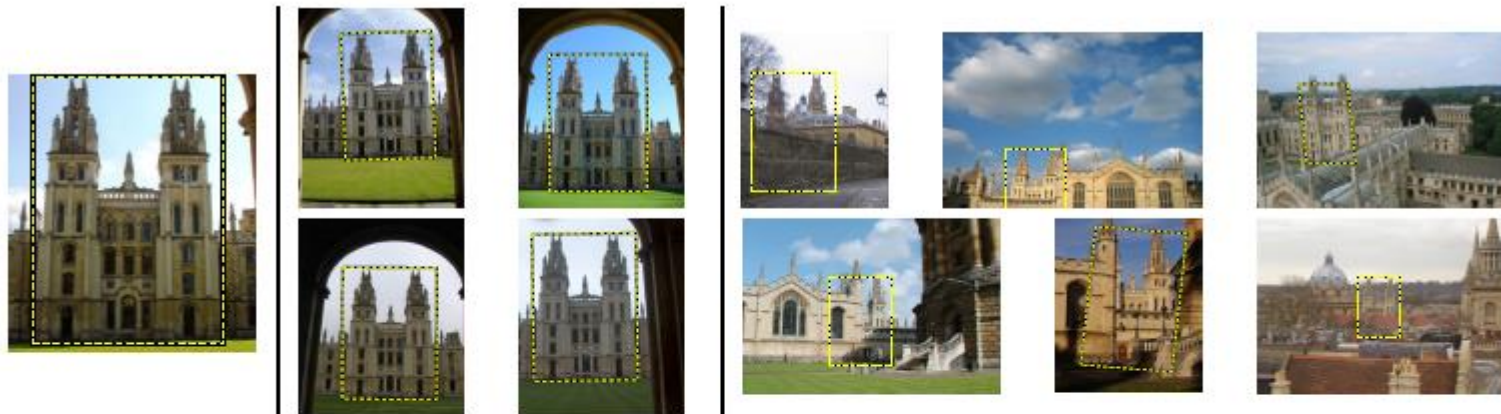
# Geometric Verification using RANSAC for Affine Transform

Repeat N times:

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



# Query Expansion [Chum et al. 07]



Original query

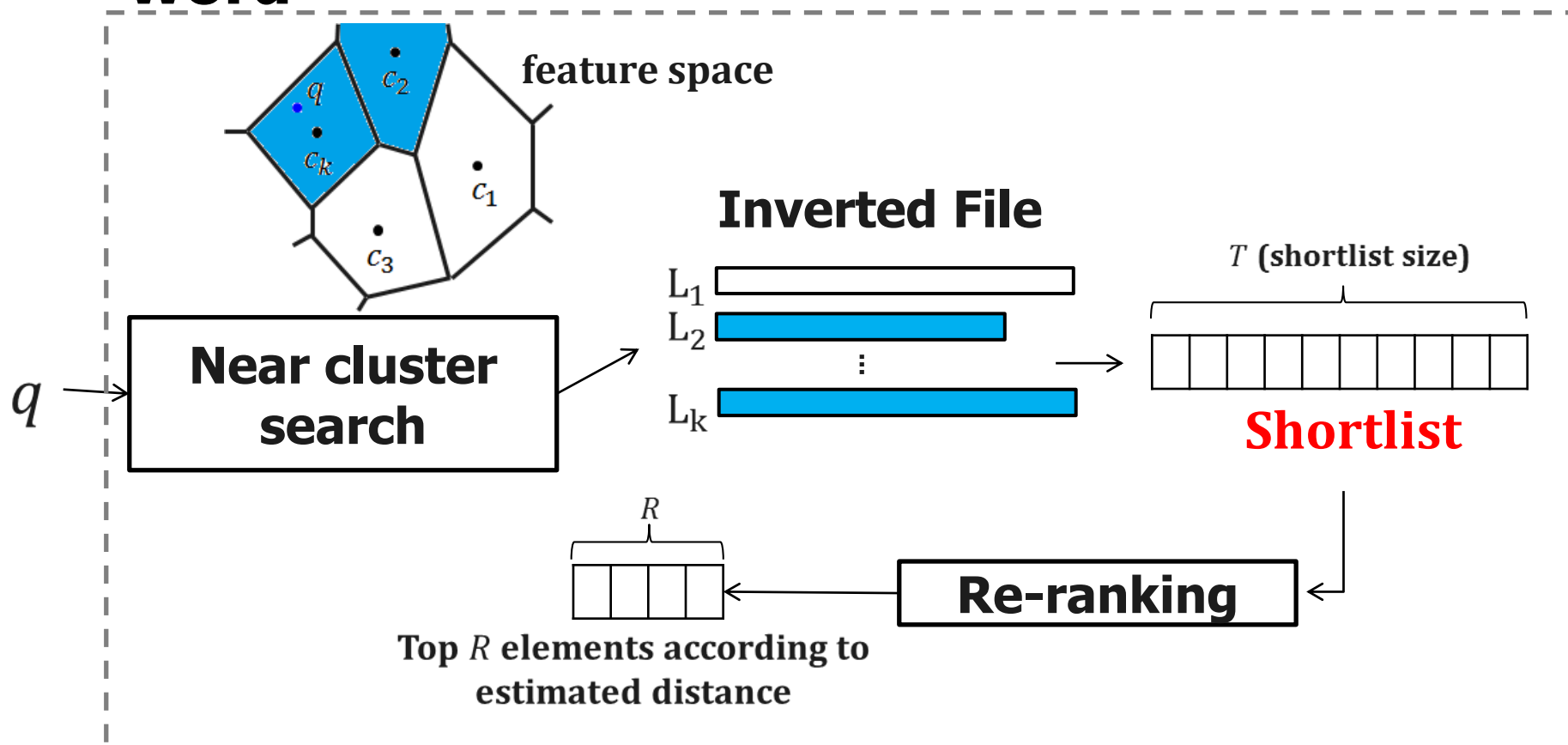
Top 4 images

Expanded results that were not identified by the original query



# 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

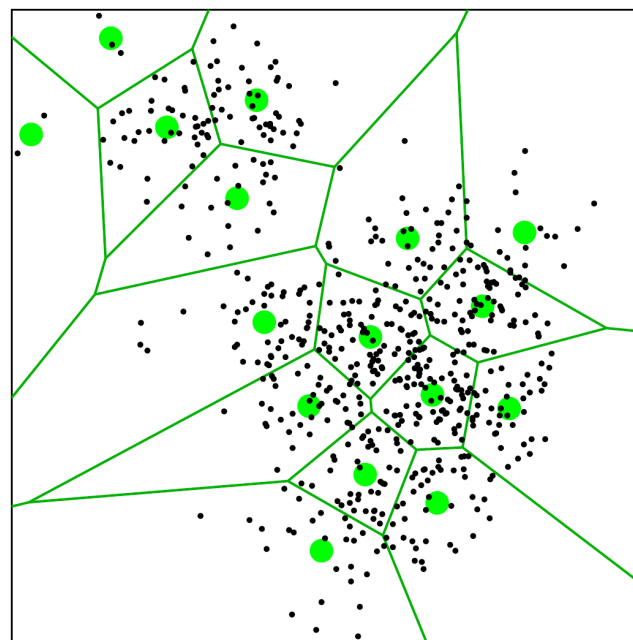
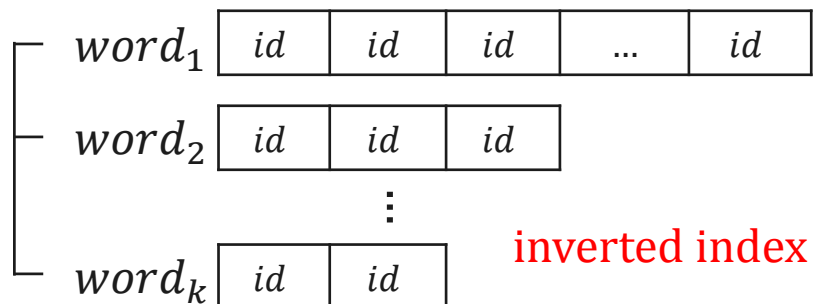


Figure from Lempitsky's slides

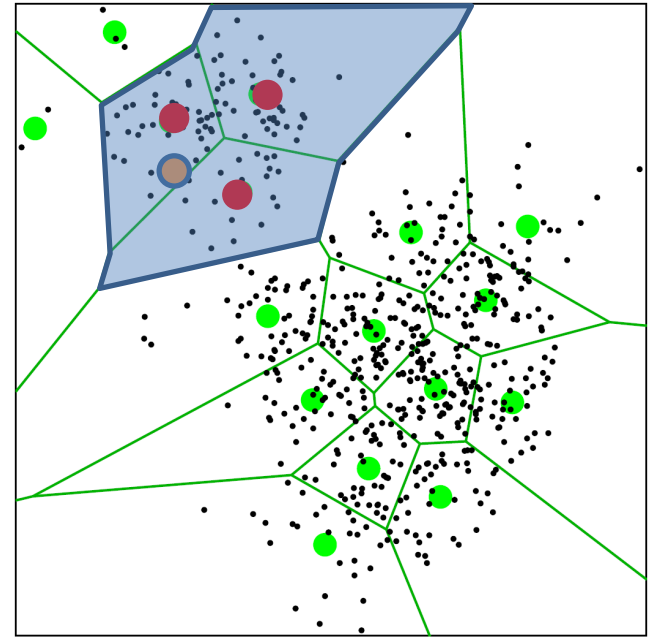




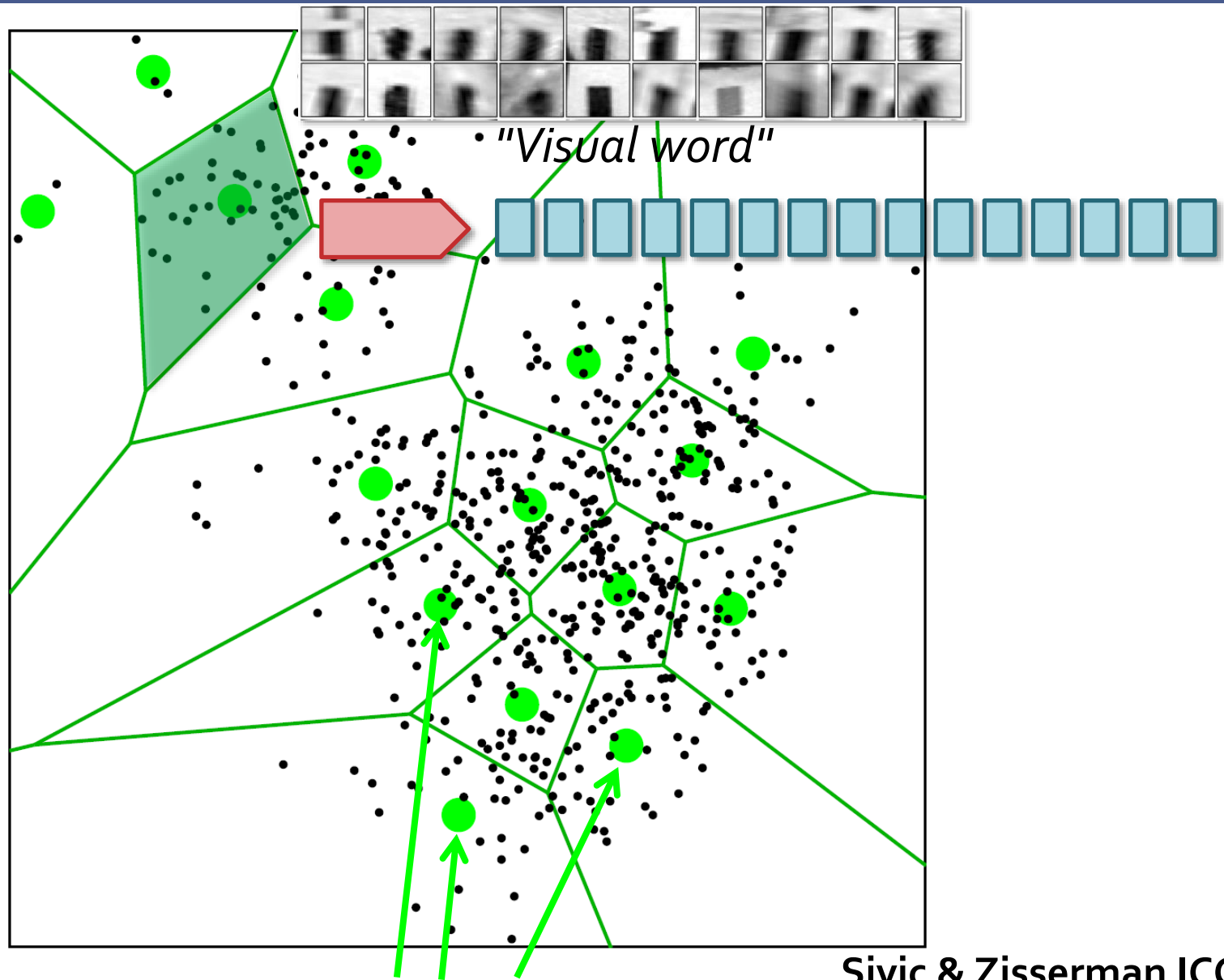
# Inverted Index

## 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



Visual codebook

Sivic & Zisserman ICCV 2003

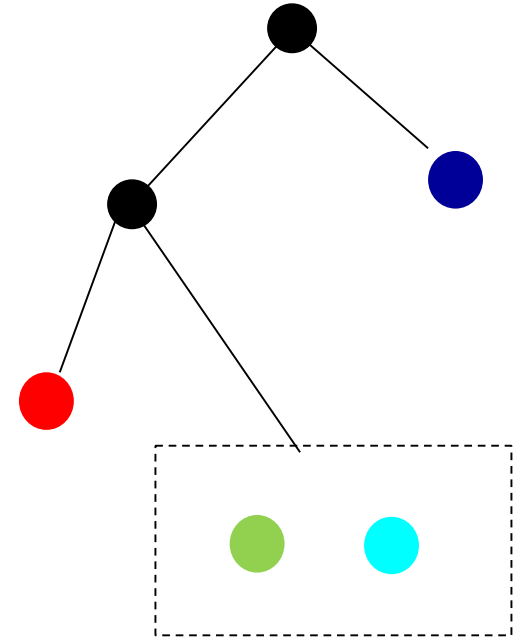
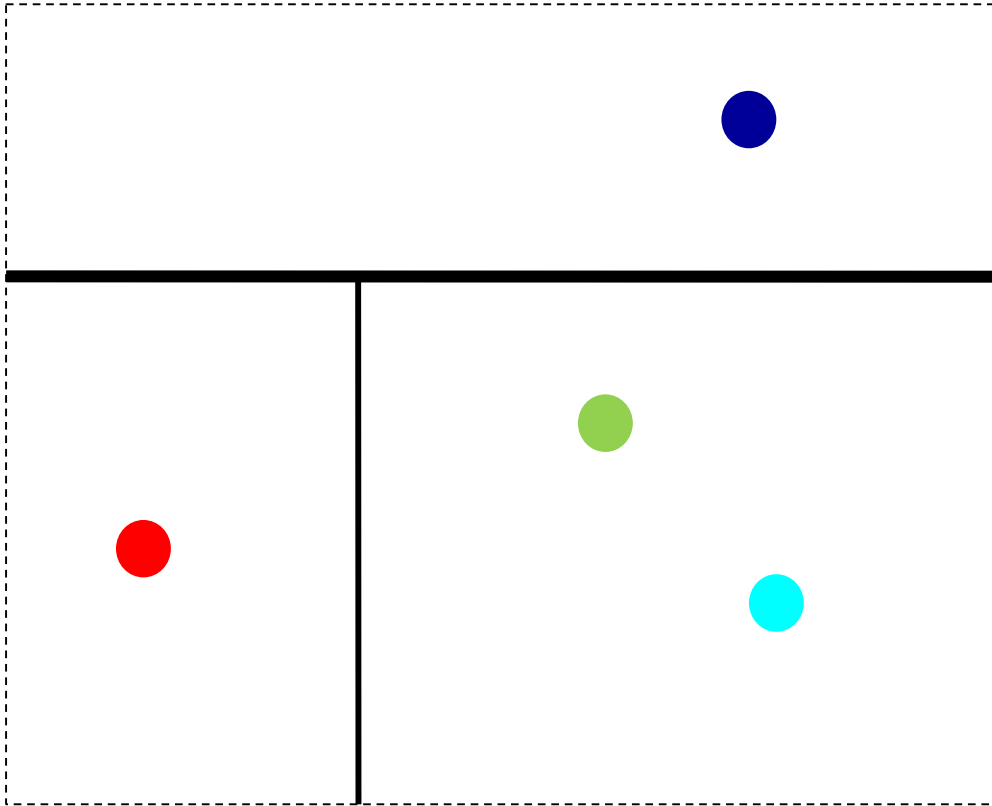
# Approximate Nearest Neighbor (ANN) Search

---

- **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)**

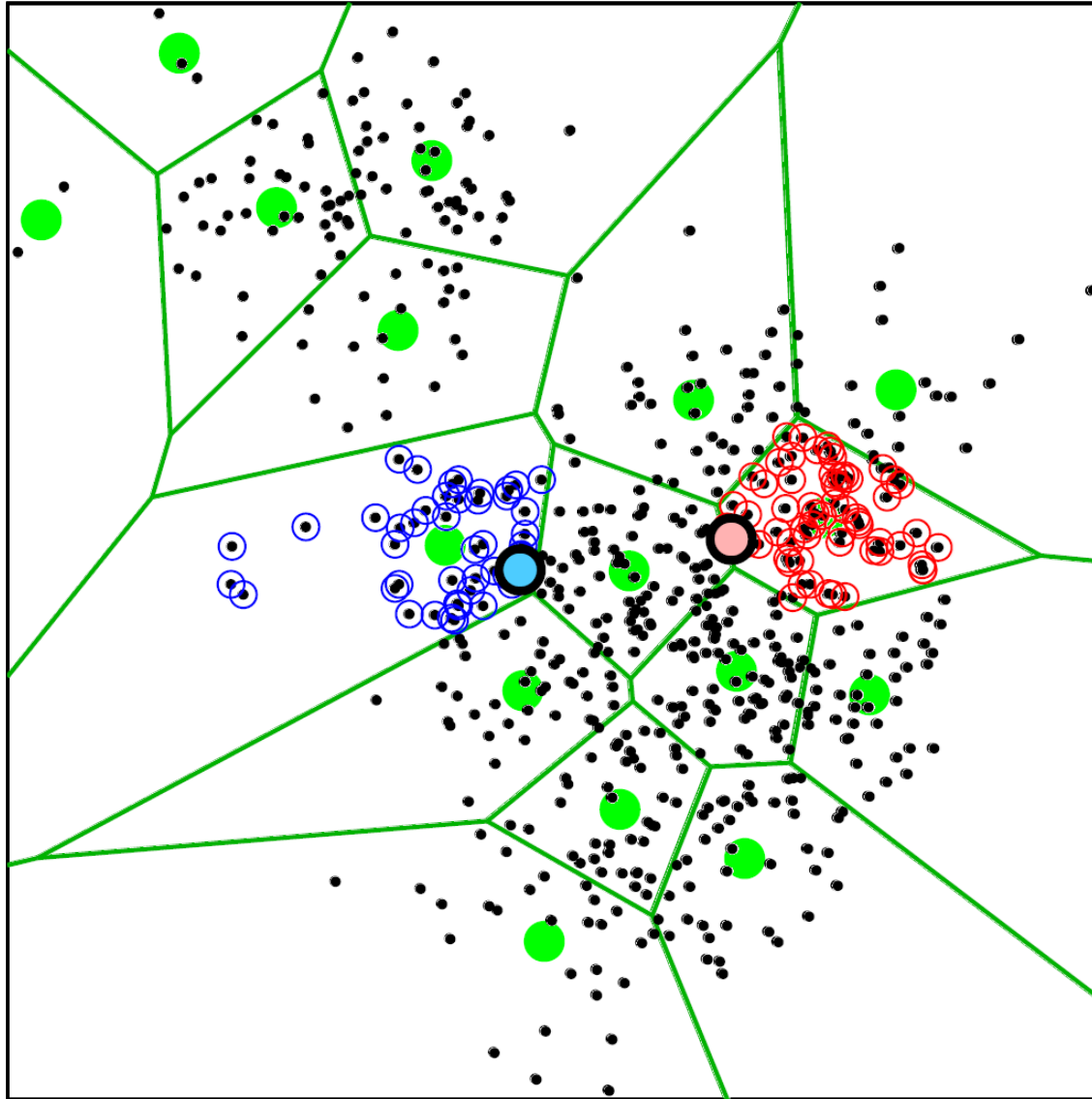
# kd-tree Example

---

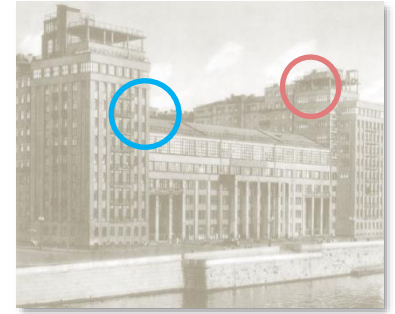


- Many good implementations (e.g., vl-feat)

# Querying the inverted index



Query:



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

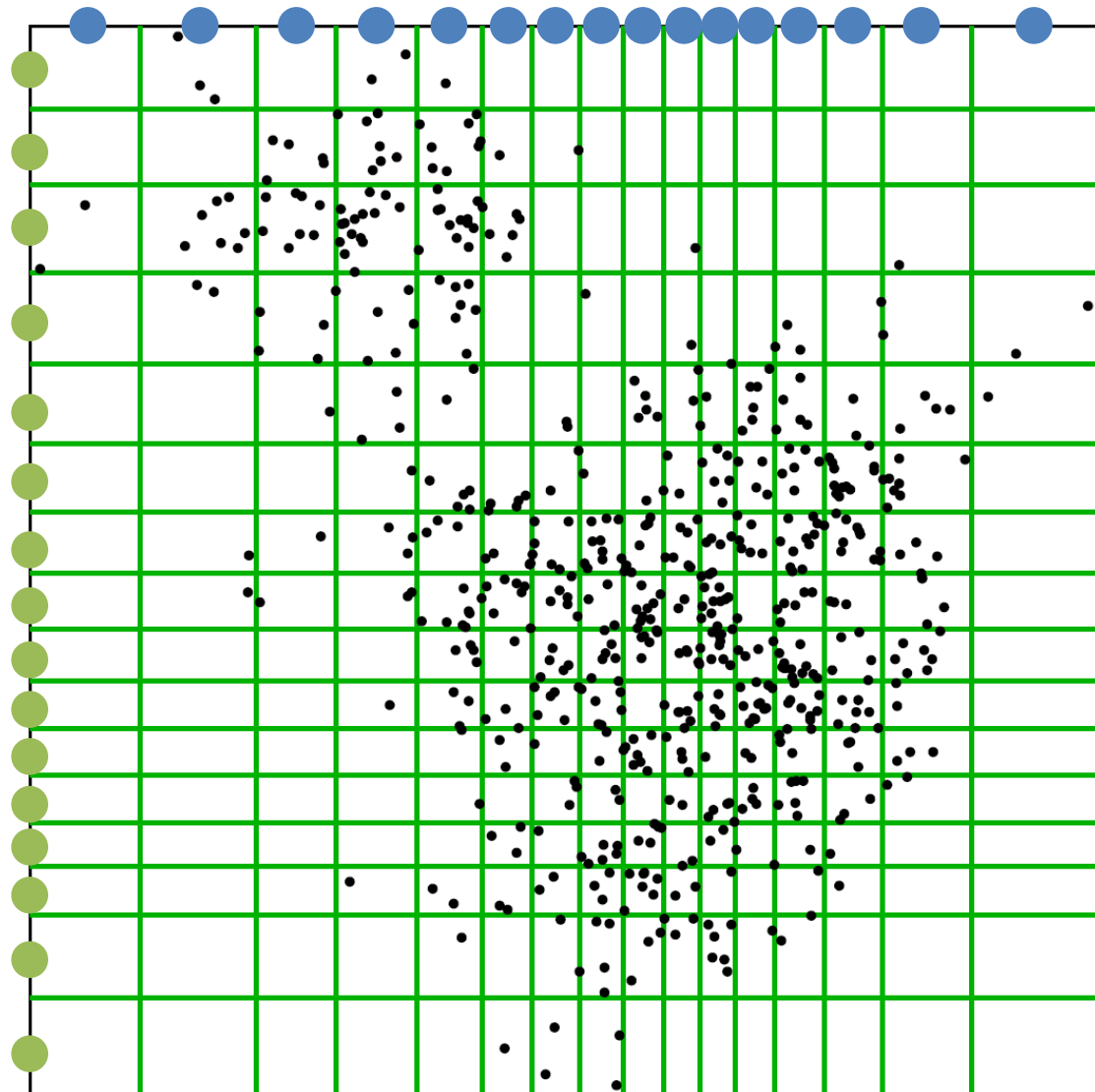


- Want to spend as little time as possible for matching to codebooks

Ack.: Lempitsky

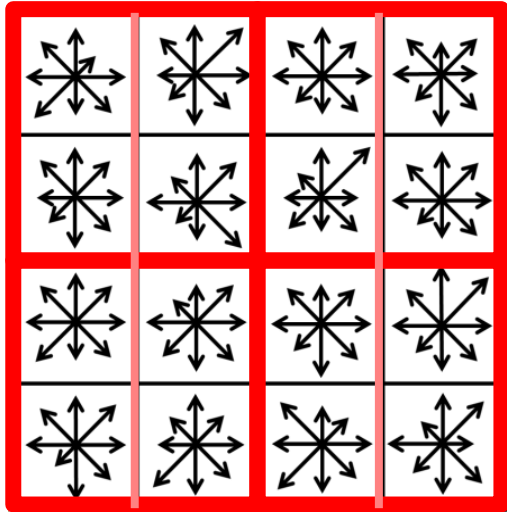
# Inverted Multi-Index [Babenko and Lempitsky, CVPR 2012]

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



Ack.: Lempitsky

# Product quantization



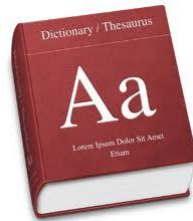
[Jegou, Douze, Schmid // TPAMI 2011]:

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

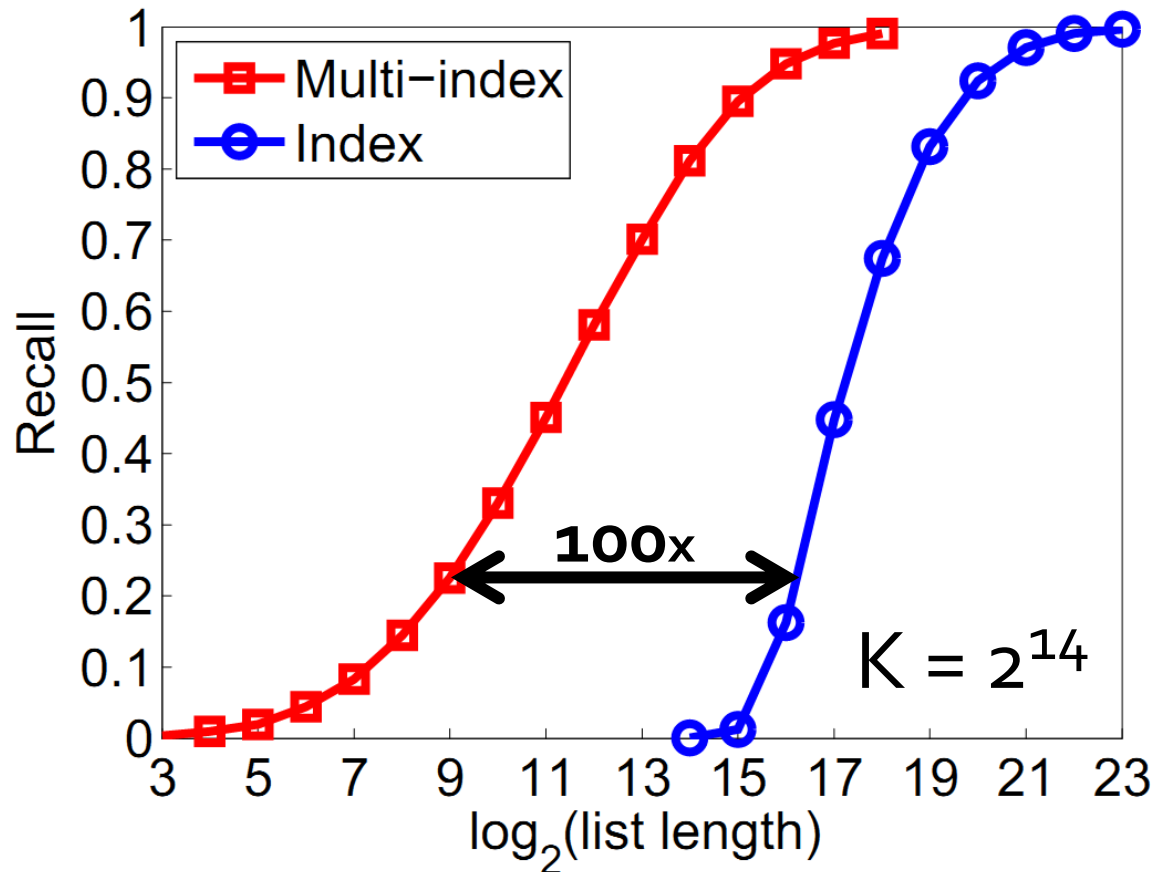


many minutes 128GB

< 1 millisecond 32KB



# Performance comparison on 1 B SIFT descriptors

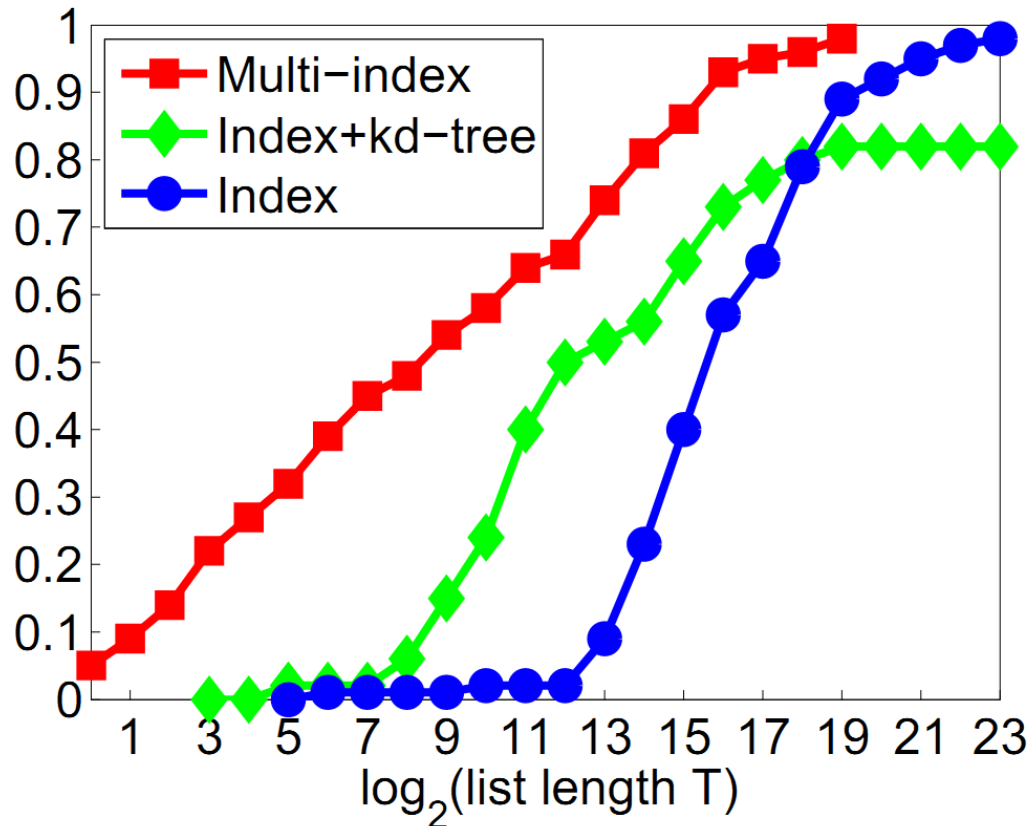


Time increase: 1.4 msec -> 2.2 msec on a single core  
(with BLAS instructions)

Ack.: Lempitsky

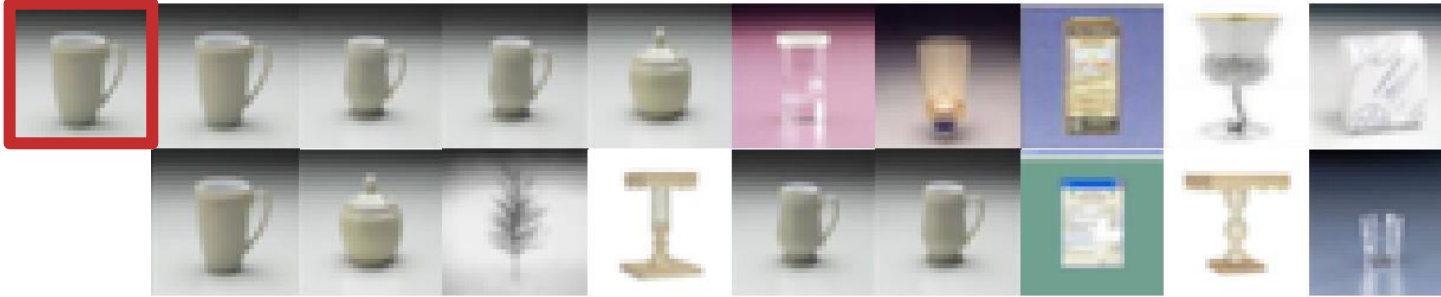
# Performance on 80 million GISTs

## 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**