



Content-Based Image Retrieval

(以图搜图)

Jun Wu (邬俊)

School of Computer & Information Technology

Beijing Jiaotong University, Beijing, China

URL: <http://faculty.bjtu.edu.cn/8620/>



A Multifaceted Definition of CV

☐ Automatic *understanding* of image/video (pixels)

■ Interpretation

- ☐ *Recognizing* objects, people, scenes and activities from visual data

■ Measurement

- ☐ Computing *3D properties* from visual data

■ Search & management

- ☐ Finding out useful information from *massive* visual data

Retrieving Specific Objects

Search photos on the web for particular places



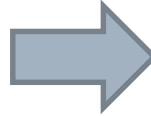
Find these landmarks

...in these images and 1M more

Retrieving Specific Objects

Visual search in feature films

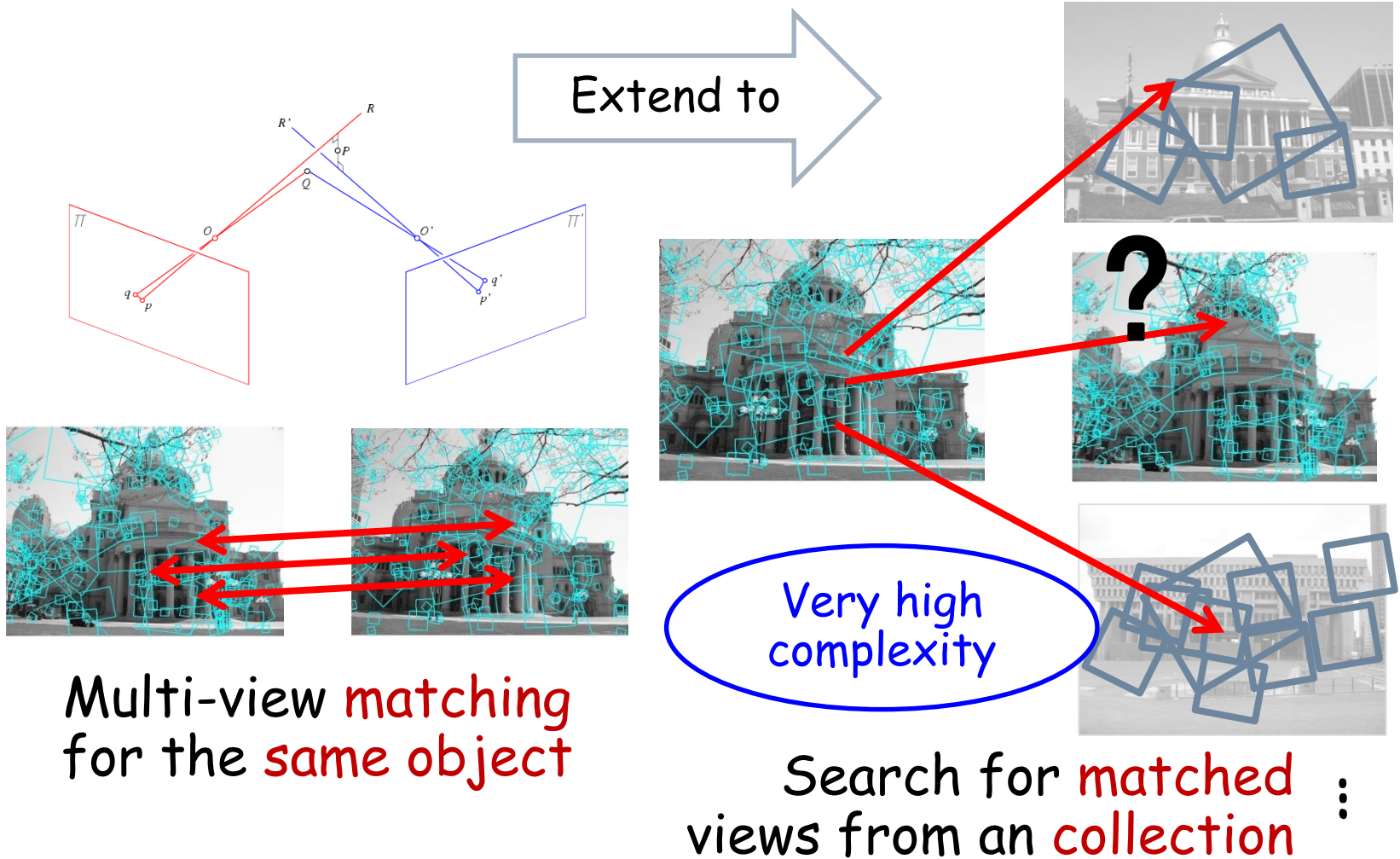
Visually defined query



Retrieved frames

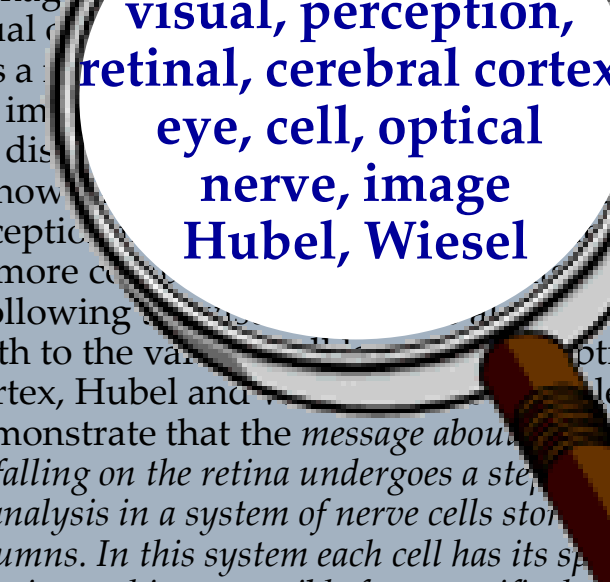
"Groundhog Day" [Rammis, 1993]

Matching Local Features




Bag of Words Models

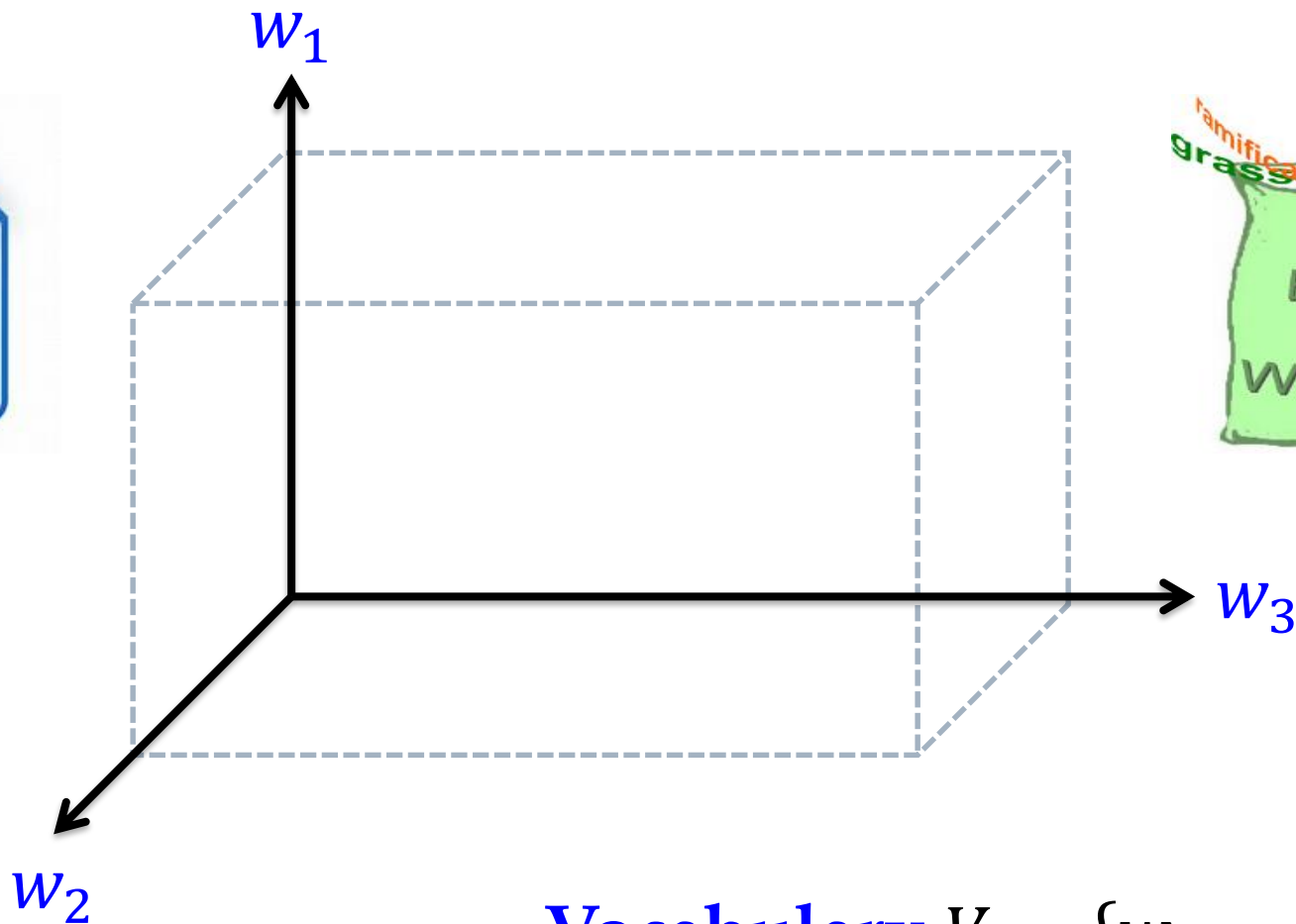




**sensory, brain,
visual, perception,
retinal, cerebral cortex,
eye, cell, optical
nerve, image
Hubel, Wiesel**

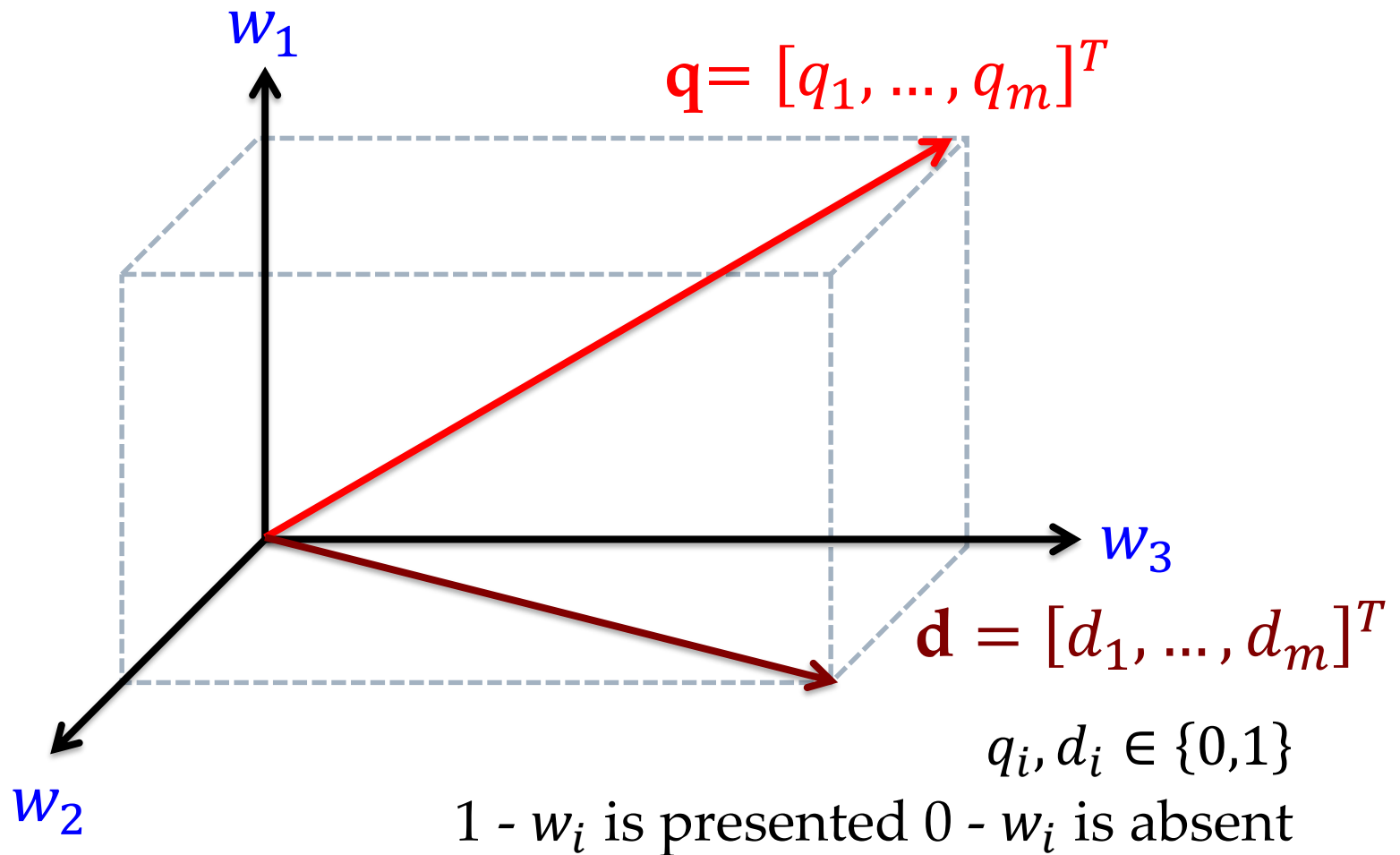


Dimension Instantiation: BoW

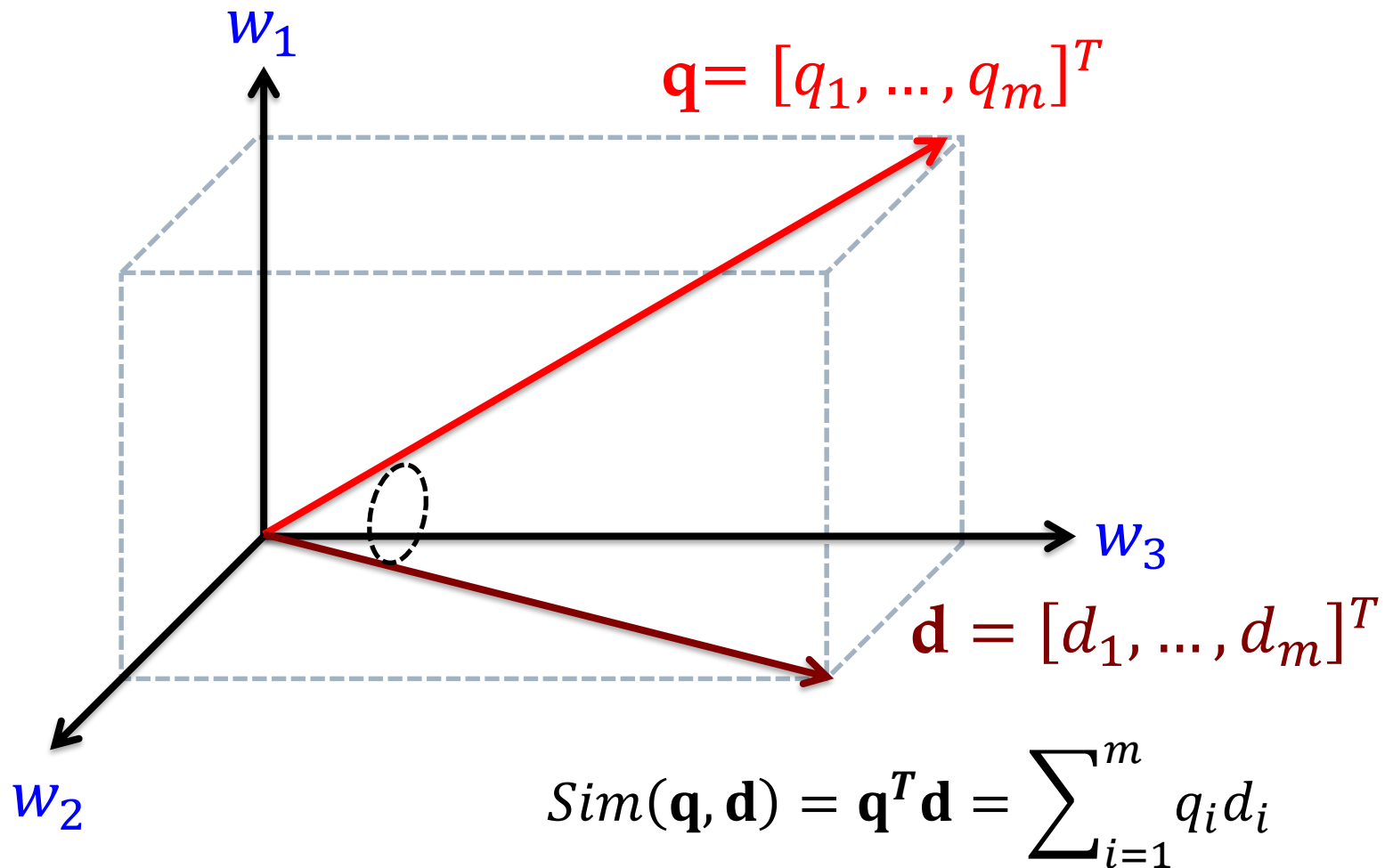


Vocabulary $V = \{w_1, \dots, w_m\}$

Vector Placement: Bit Vector



Vector Placement: Bit Vector



Simplest VSM

BoW + Bit-Vector + Dot-Product

$$\mathbf{q} = [q_1, \dots, q_m]^T \quad q_i, d_i \in \{0,1\}$$

$$\mathbf{d} = [d_1, \dots, d_m]^T \quad 1 - w_i \text{ is present; } 0 - w_i \text{ is absent}$$

$$\text{Sim}(\mathbf{q}, \mathbf{d}) = \mathbf{q}^T \mathbf{d} = \sum_{i=1}^m q_i d_i$$

What does this ranking function intuitively capture?
Is this a good ranking function?

How Would You Rank These Documents?

Query = “news about presidential campaign”

Ideal Rank

d1	... news about ...	d4	+
d2	... news about organic food campaign...	d3	+
d3	... news of presidential campaign ...	d1	0
d4	... news of presidential campaign presidential candidate ...	d2	-

Ranking Using the Simplest VSM

Query = “news about presidential campaign”

d3 ... news of presidential campaign ...

$V = \{\text{news, about, presidential, campaign, food ...}\}$

q=	(1,	1,	1,	0, ...)
d3=	(1,	0,	1,	1, 0, ...)

$$f(q, d3) = 1*1 + 1*0 + 1*1 + 1*1 + 0*0 + \dots = 3$$

of distinct query words matched in a document

Is the Simplest VSM Effective?

Query = “news about presidential campaign”

Score

d1	... news about ...	$f(q, d1)=2$
d2	... news about organic food campaign...	$f(q, d2)=3$
d3	... news of presidential campaign ...	$f(q, d3)=3$
d4	... news of presidential campaign presidential candidate ...	$f(q, d4)=3$

Speculation

□ How to achieve $f(q,d4) > f(q,d3) > f(q,d2)$

d2	... news about organic food campaign...	$f(q,d2)=3$
d3	... news of presidential campaign ...	$f(q,d3)=3$
d4	... news of presidential campaign presidential candidate ...	$f(q,d4)=3$

1. Matching “**presidential**” more times deserves more credit
(Comparing d3 & d4)
2. Matching “**presidential**” is more important than matching
“**about**” (Comparing d2 & d3)

Improved VSM with TF Weighting

$\mathbf{q} = [q_1, \dots, q_m]^T$ \square q_i : count of w_i in query

$\mathbf{d} = [d_1, \dots, d_m]^T$ \square d_i : count of w_i in document

$$\text{Sim}(\mathbf{q}, \mathbf{d}) = \mathbf{q}^T \mathbf{d} = \sum_{i=1}^m q_i d_i$$

- \square What does this ranking function intuitively capture?
- \square Does it fix the problems of the simplest VSM?

Query = “news about presidential campaign”

d2

... news about organic food campaign...

$f(q, d2)=3$

q=	(1,	1,	1,	0,	...)
d2=	(1,	1,	1,	1,	...)

d3

... news of presidential campaign ...

$f(q, d3)=3$

q=	(1,	1,	1,	0,	...)
d3=	(1,	0,	1,	1,	...)

d4

... news of presidential campaign ...
... presidential candidate ...

$f(q, d4)=4$

q=	(1,	1,	1,	0,	...)
d4=	(1,	0,	2,	1,	...)

How to Fix Problem 2 (“presidential” vs. “about”)

d2 ... news about organic food campaign...


d3 ... news of presidential campaign ...

$V = \{\text{news, about, presidential, campaign, food ...}\}$

q =	(1,	1,	1,	1,	0,	...)	$f(q, d2) < 3$
d2 =	(1,	1,	0,	1,	1,	...)	



q =	(1,	1,	1,	1,	0,	...)	$f(q, d3) > 3$
d3 =	(1,	0,	1,	1,	0,	...)	



Further Improved **VSM** with **TF-IDF**

$$\mathbf{q} = [q_1, \dots, q_m]^T$$

□ q_i : count of w_i in query

$$\mathbf{d} = [d_1, \dots, d_m]^T$$

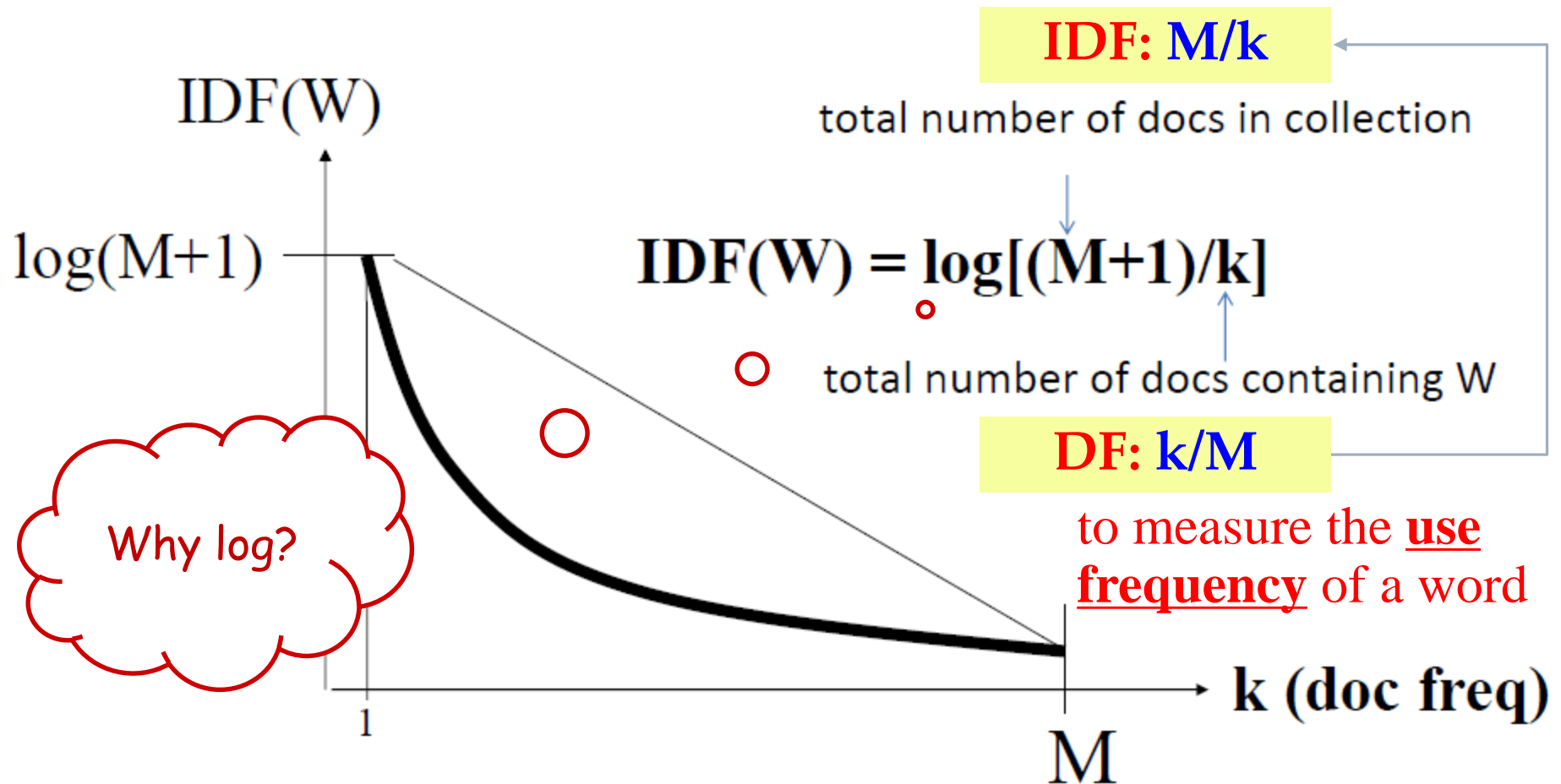
□ d_i : $c(w_i, \mathbf{d}) * \mathbf{IDF}(w_i)$

$$Sim(\mathbf{q}, \mathbf{d}) = \mathbf{q}^T \mathbf{d} = \sum_{i=1}^m q_i d_i$$

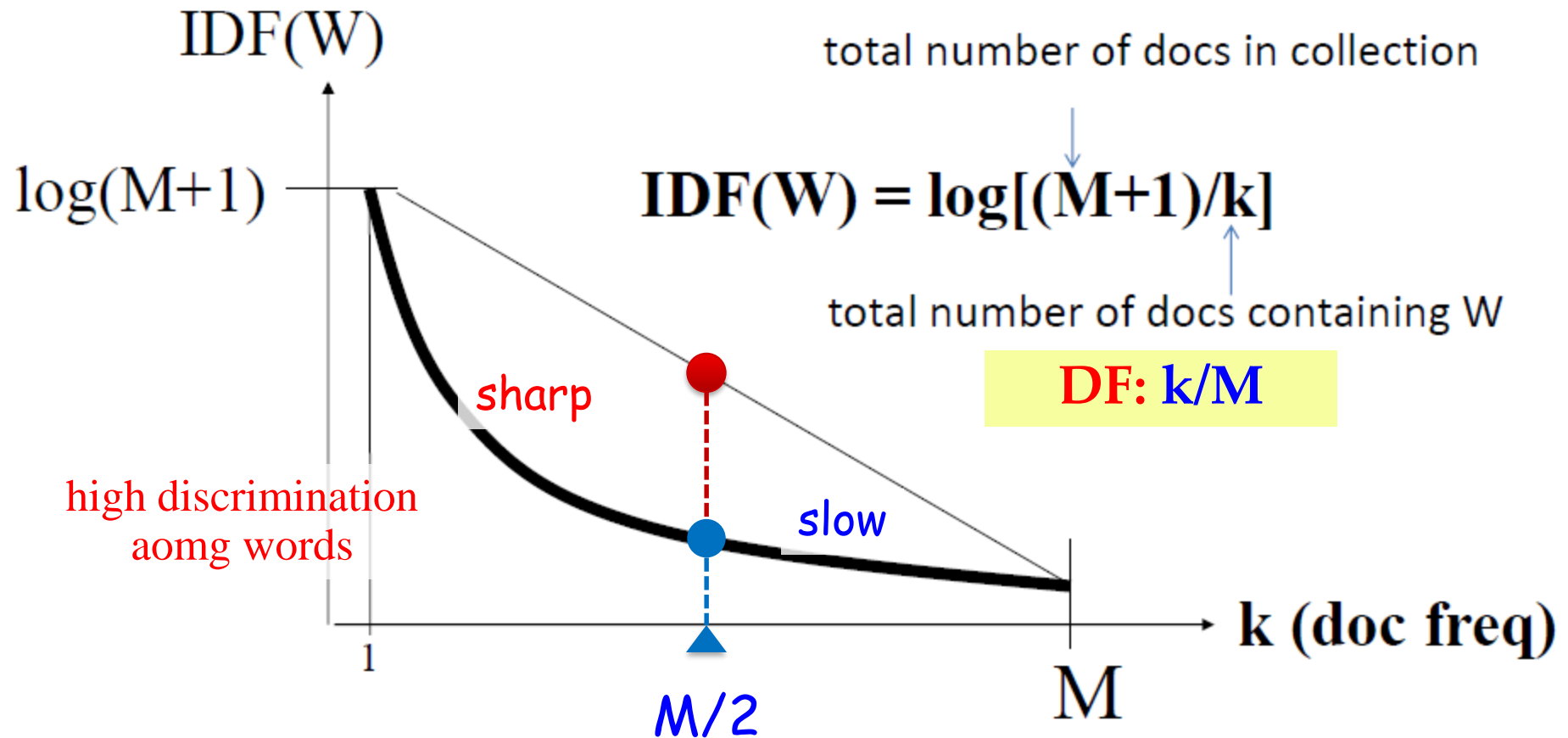
□ **Adding Inverse Document Frequency (IDF)**

□ Does it fix the problems of the simplest VSM?

IDF Weighting: Penalizing Popular Terms



IDF Weighting: Penalizing Popular Terms



Solving Problem 2 (“presidential” vs. “about”)

d2 ... news about organic food campaign...

d3 ... news of presidential campaign ...

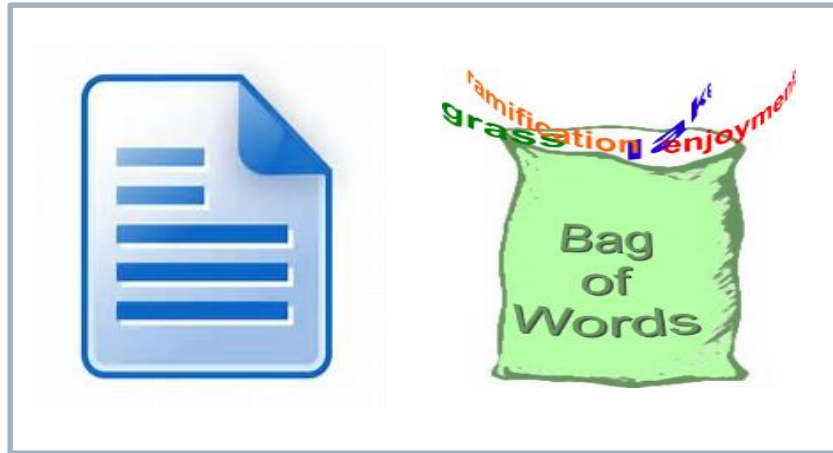
V= {news, about, presidential, campaign, food }				
IDF(W)= 1.5	1.0	2.5	3.1	1.8
q= (1,	1,	1,	1,	0, ...)
d2= (1*1.5,	1*1.0	0,	1*3.1,	0, ...)
q= (1,	1,	1,	1,	0, ...)
d3= (1*1.5,	0,	1*2.5	1*3.1,	0, ...)

$$f(q,d2) = 5.6 < f(q,d3)=7.1$$

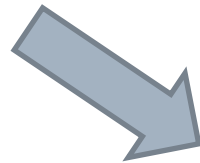
How Effective is VSM with TF-IDF?

Query = “news about presidential campaign”		Score
d1	... news about ...	$f(q,d1)=2.5$
d2	... news about organic food campaign...	$f(q,d2)=5.6$
d3	... news of presidential campaign ...	$f(q,d3)=7.1$
d4	... news of presidential campaign presidential candidate ...	$f(q,d4)=9.6$

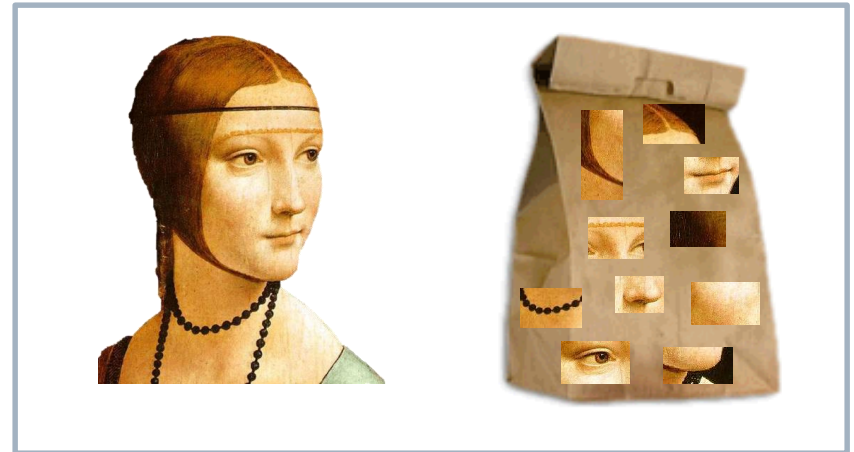
Bag-of-VisualWord (Features)

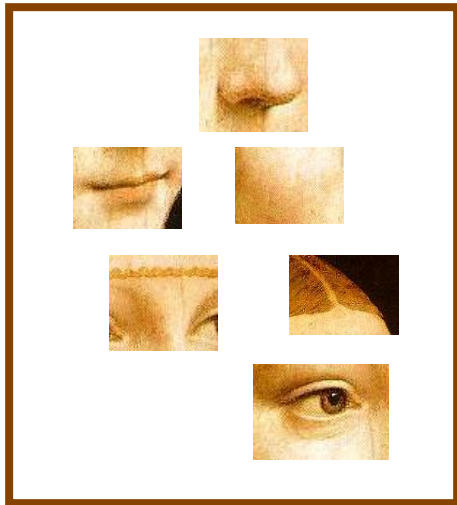


Bag-of-Word



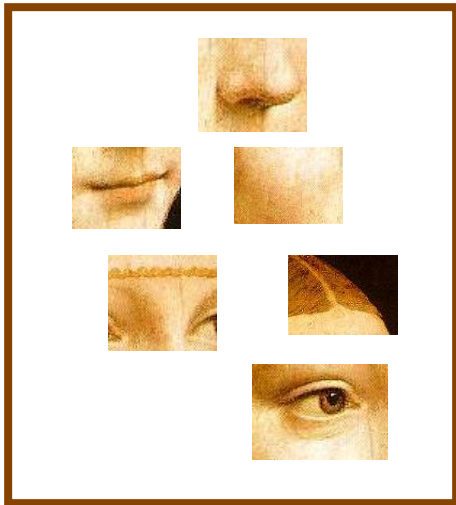
Bag-of-Visual-Word





Step-1: Extract local features from each image

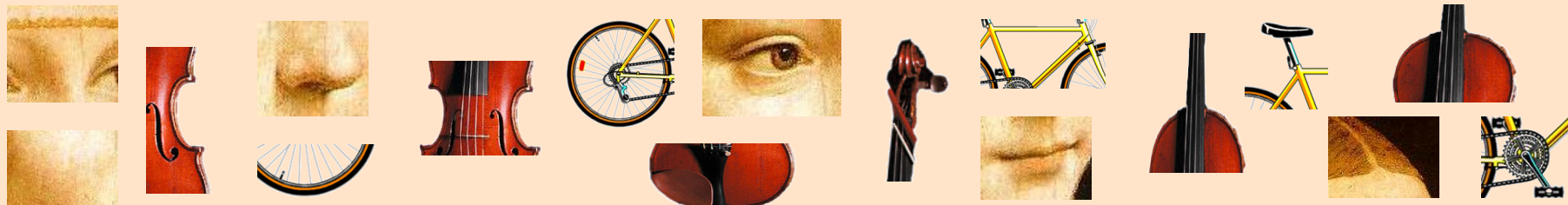
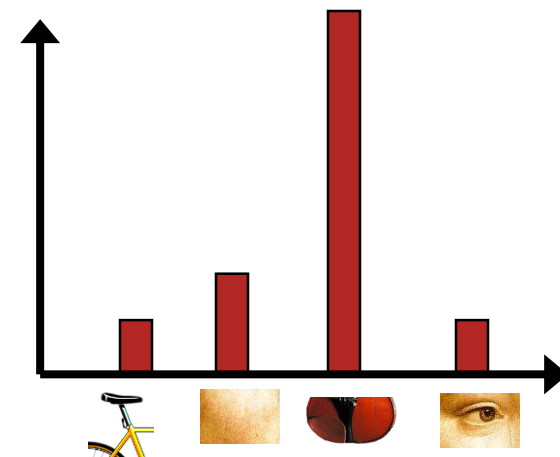
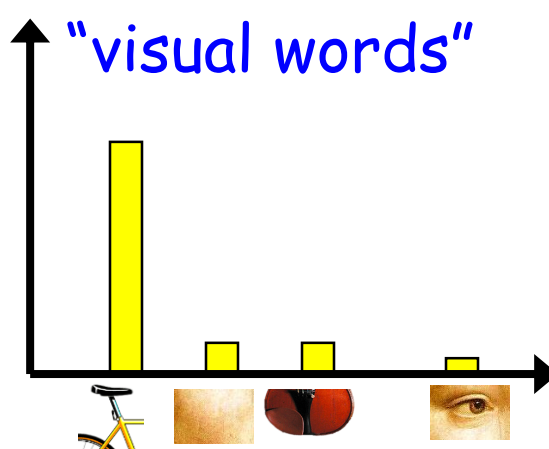
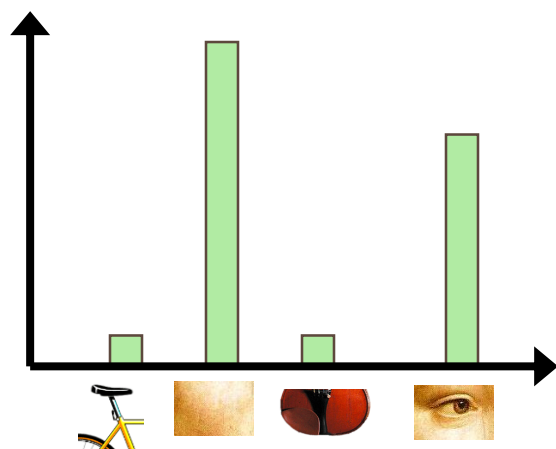
Step-2: Learn "visual vocabulary" based on extracted local features





Step-3: Represent images by frequencies of

"visual words"





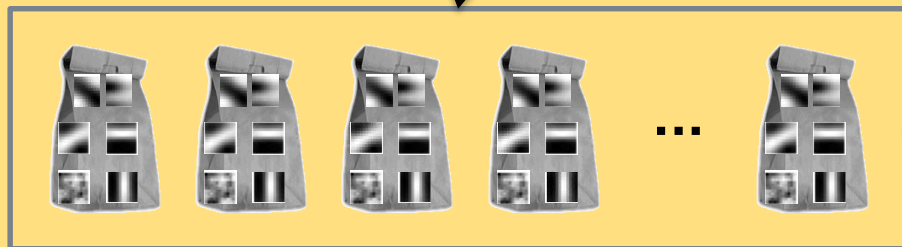
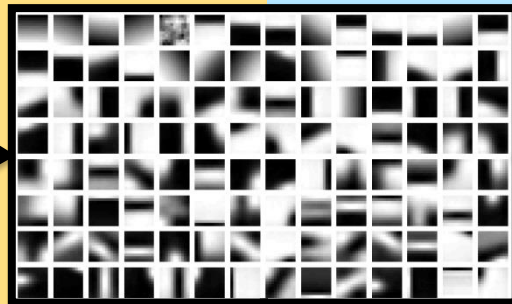
Learning

Querying



1. feature detection
& description

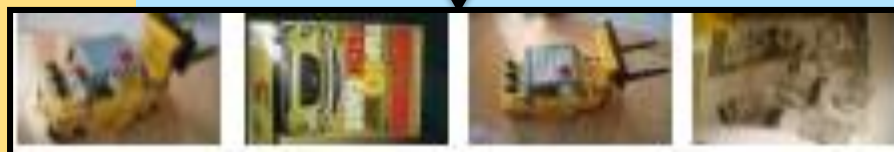
2. visual-word dictionary



3. image representation

Similarity
Computation

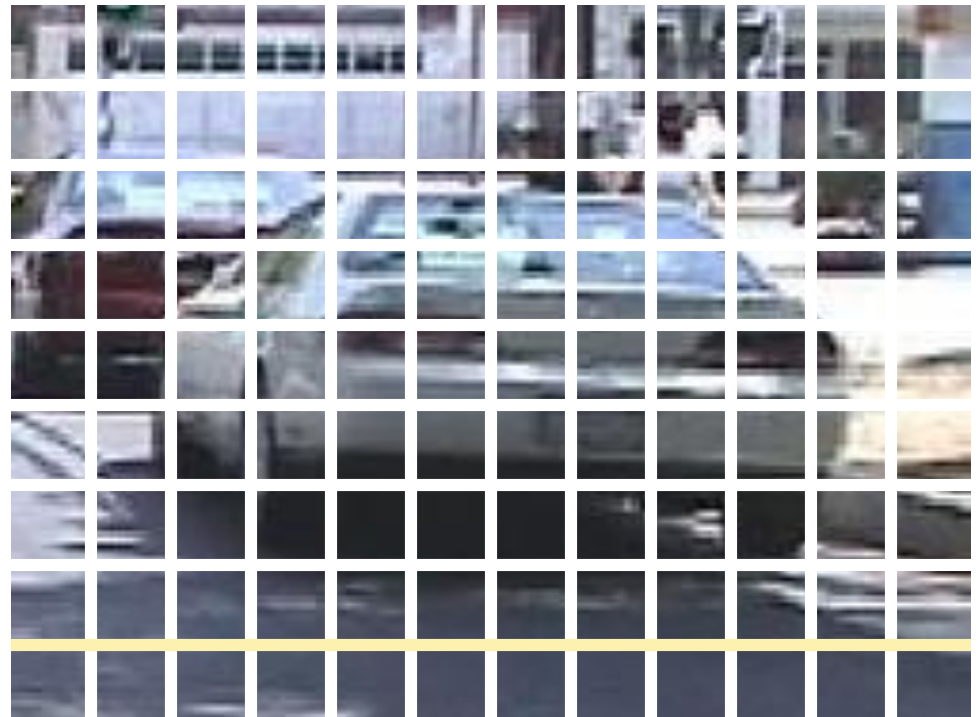
Ranked image list



Step-1: Feature extraction

□ Regular grid

- Vogel & Schiele, 2003
- Fei-Fei & Perona, 2005



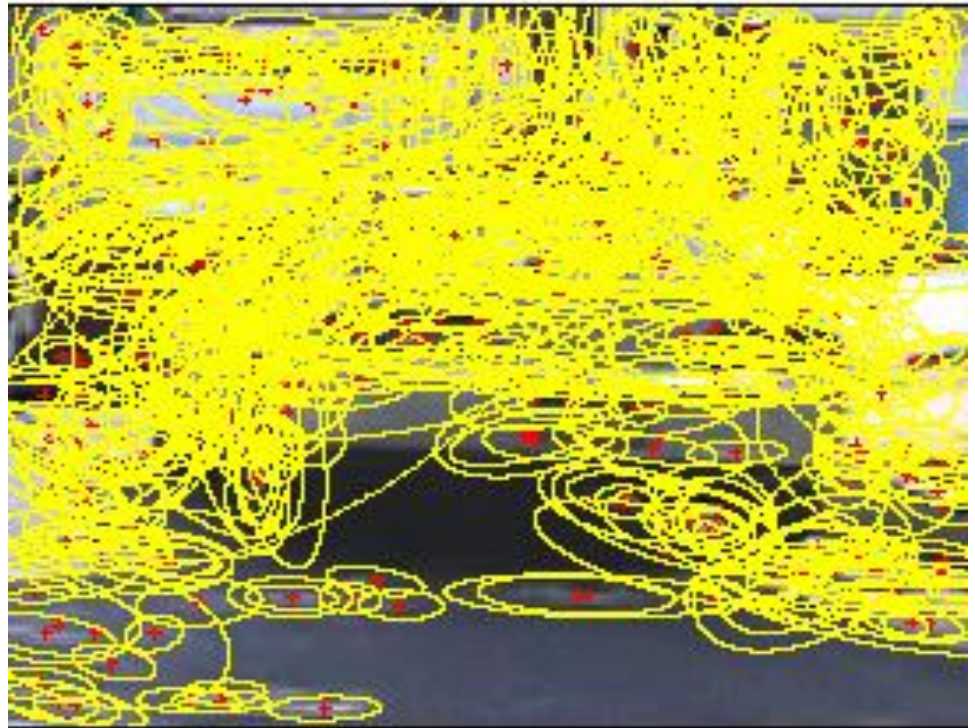
Step-1: Feature extraction

□ Regular grid

- Vogel & Schiele, 2003
- Fei-Fei & Perona, 2005

□ Interest point detector

- Csurka et al. 2004
- Fei-Fei & Perona, 2005
- Sivic et al. 2005



Step-1: Feature extraction

□ Regular grid

- Vogel & Schiele, 2003
- Fei-Fei & Perona, 2005

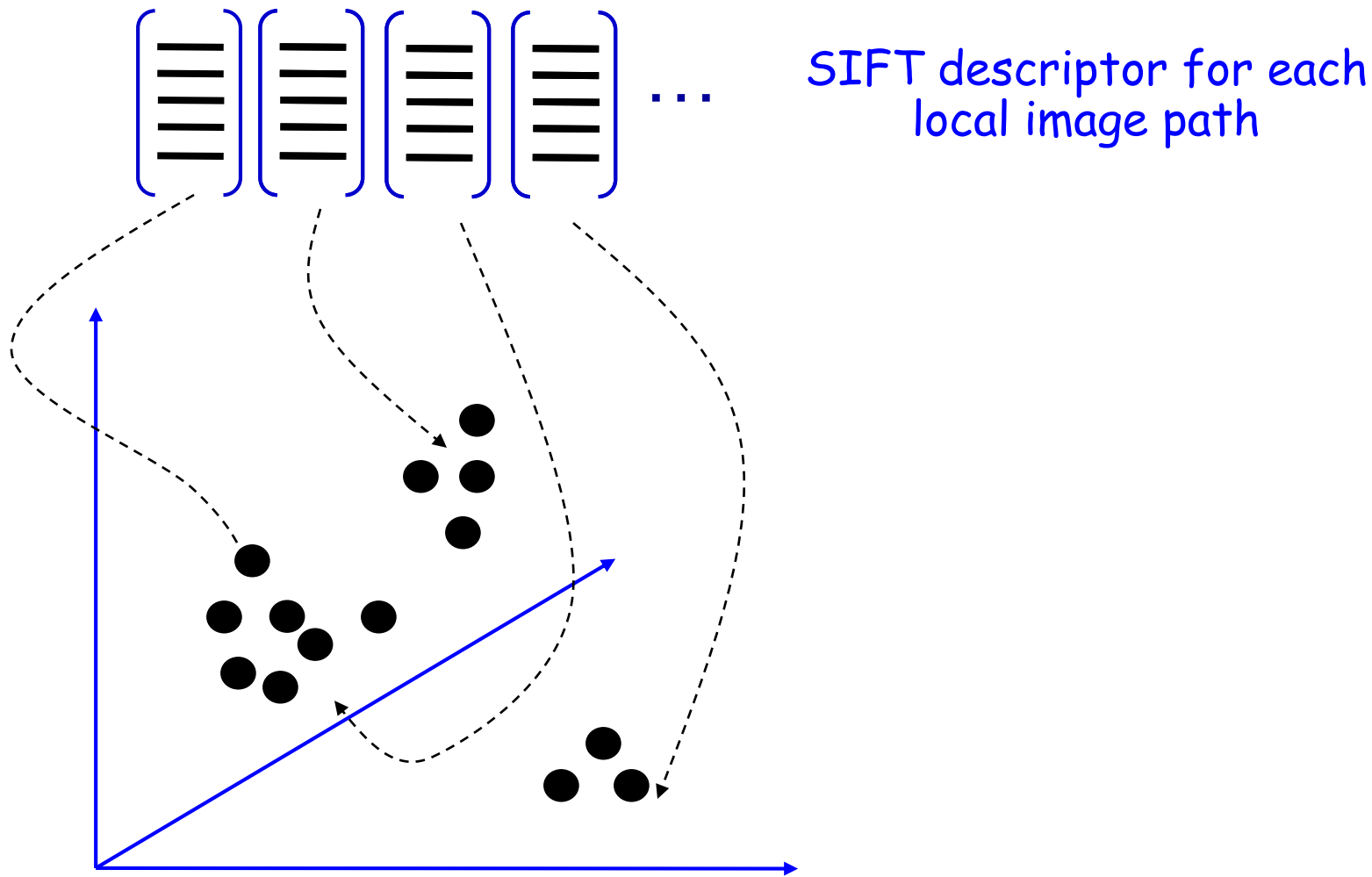
□ Interest point detector

- Csurka et al. 2004
- Fei-Fei & Perona, 2005
- Sivic et al. 2005

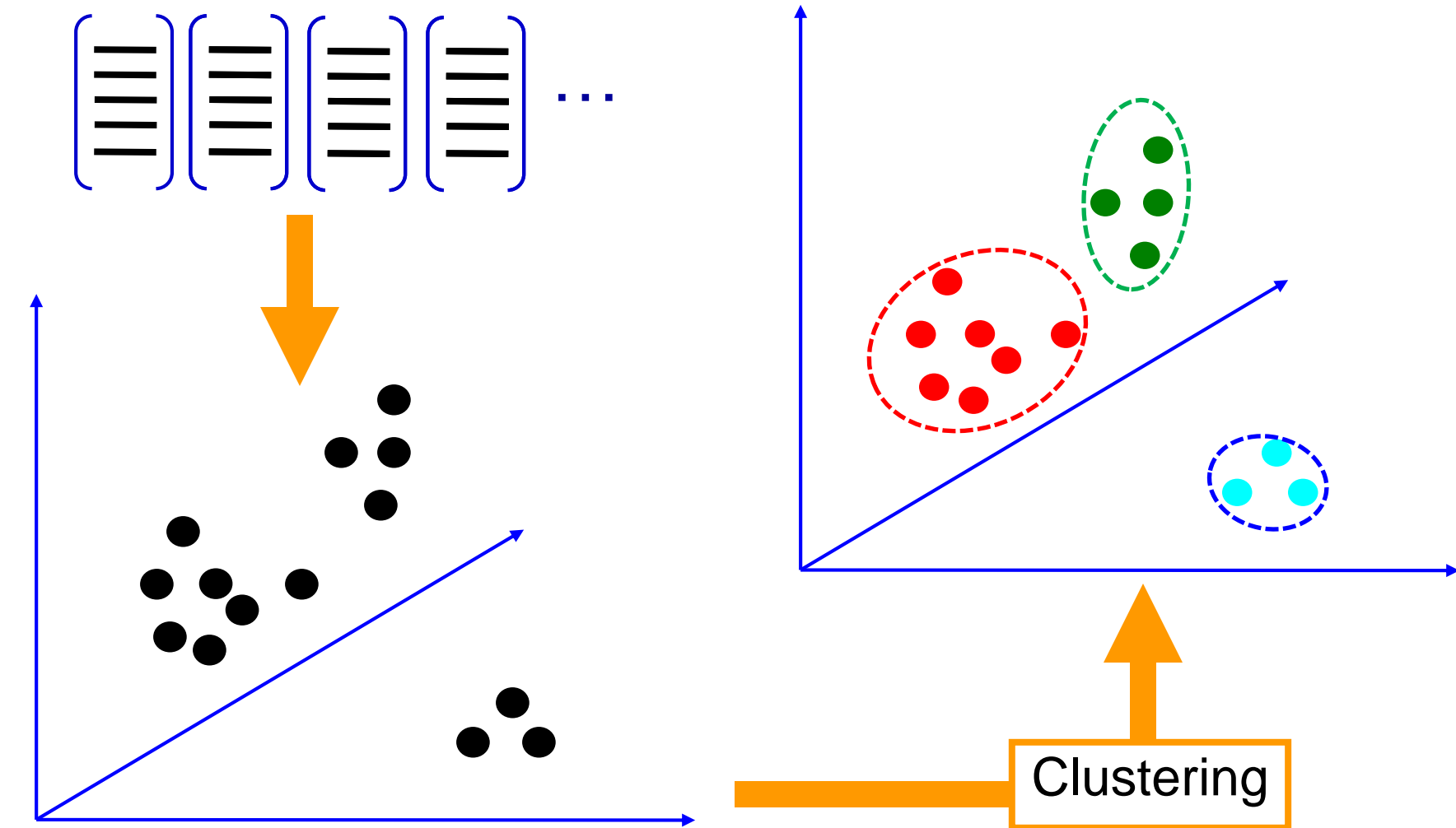
□ Other methods

- Random sampling (Vidal-Naquet & Ullman, 2002)
- Segmentation-based patches (Barnard et al. 2003)

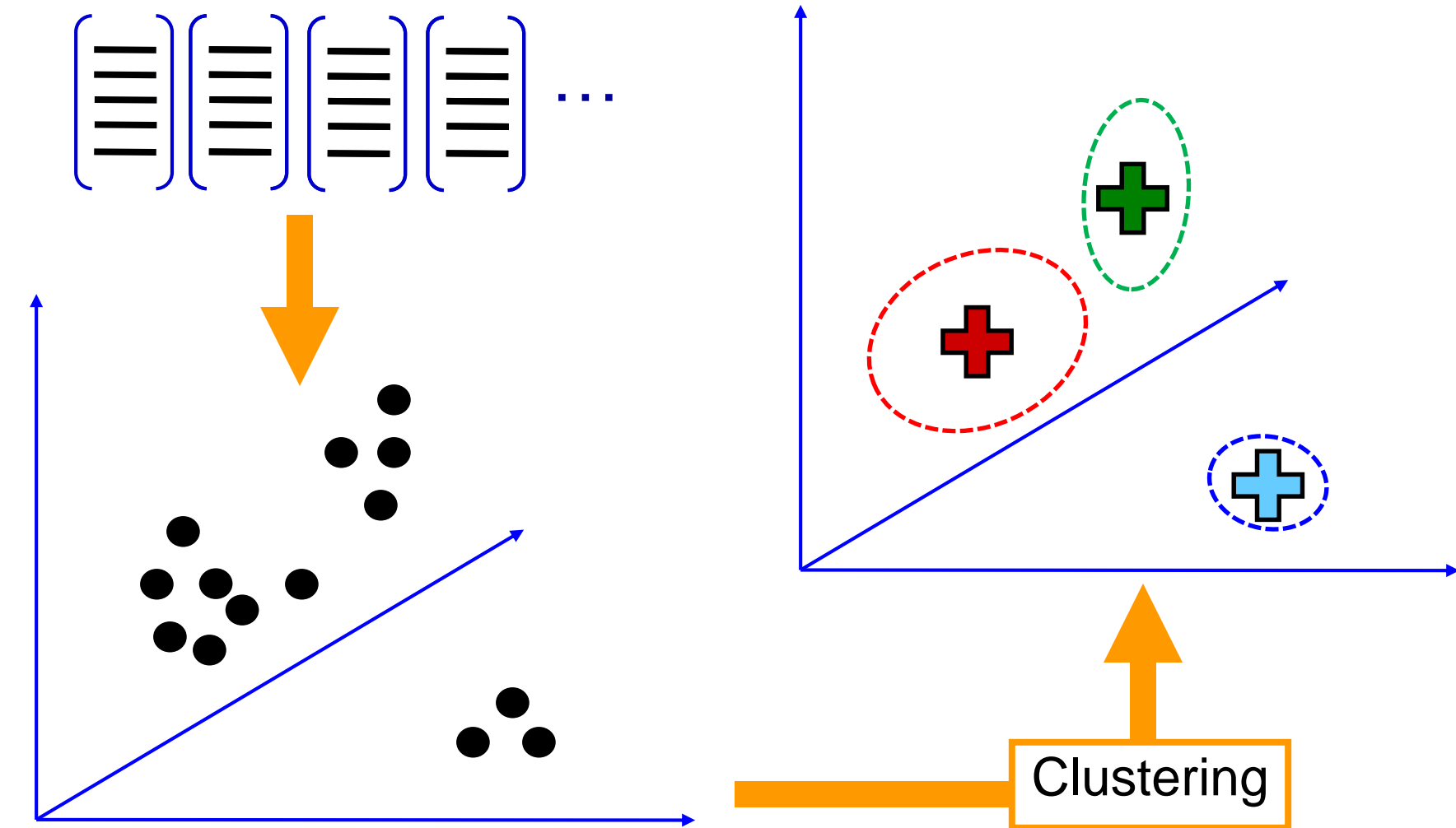
Step-2: Learning the visual vocabulary



Step-2: Learning the visual vocabulary



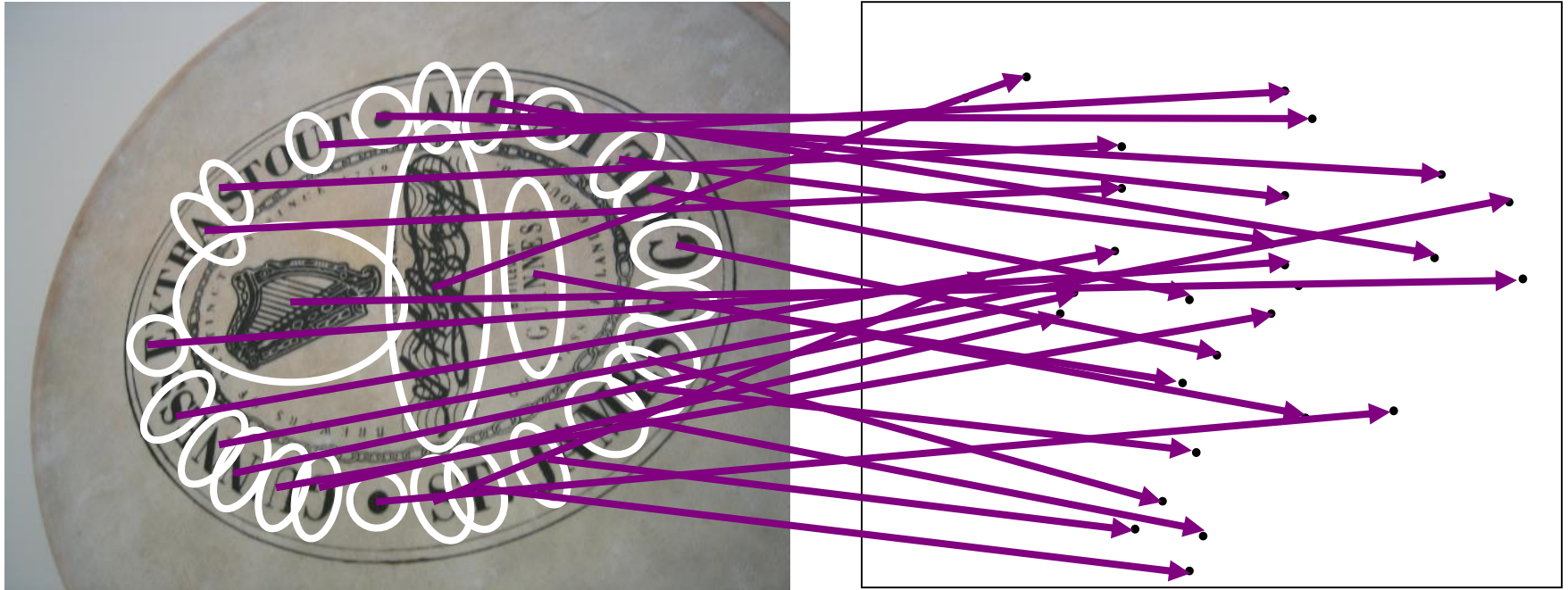
Step-2: Learning the visual vocabulary



From Clustering to Vector Quantization

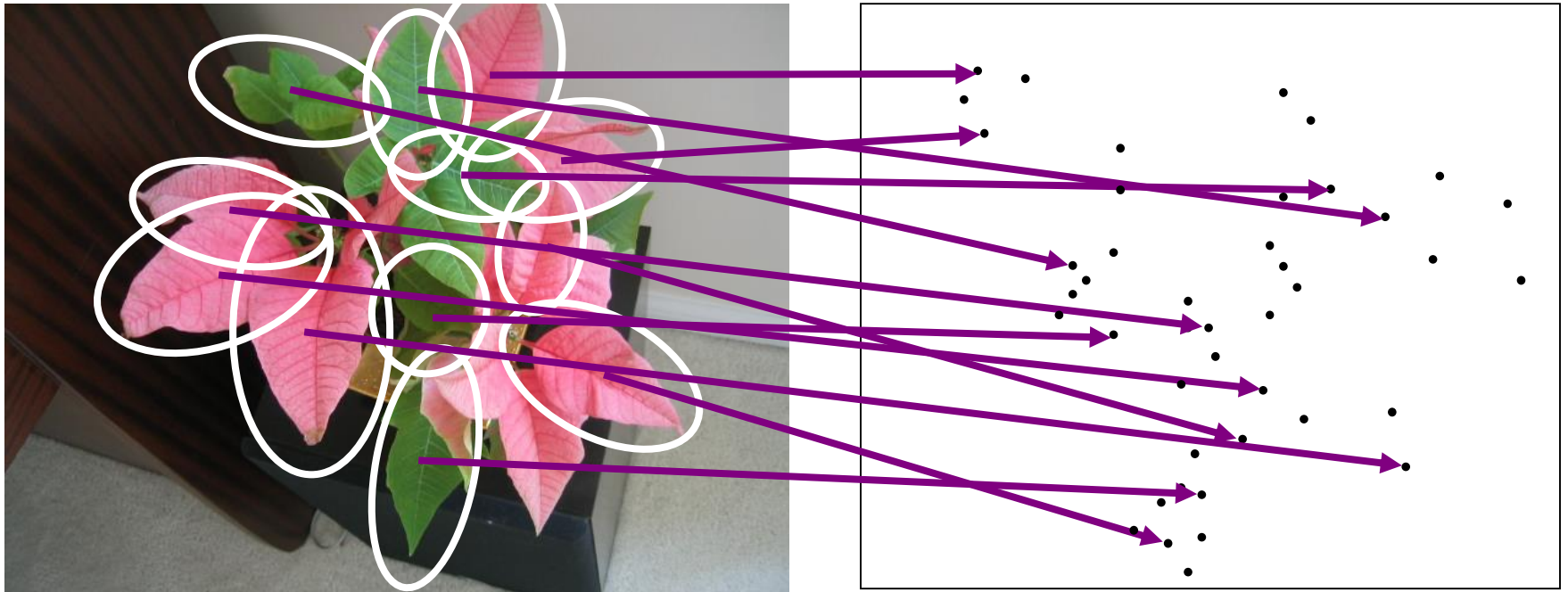
- **Clustering** is a common method for learning a visual **vocabulary or codebook**
 - Each **cluster center** produced by k-means becomes a **codevector (Quantization)**
 - Codebook can be **learned** on **separate** training set
 - Provided the training set is **sufficiently** representative, the codebook will be “**universal**”
 - **Codebook = visual vocabulary**
 - **Codevector = visual word**

Extract local features from images

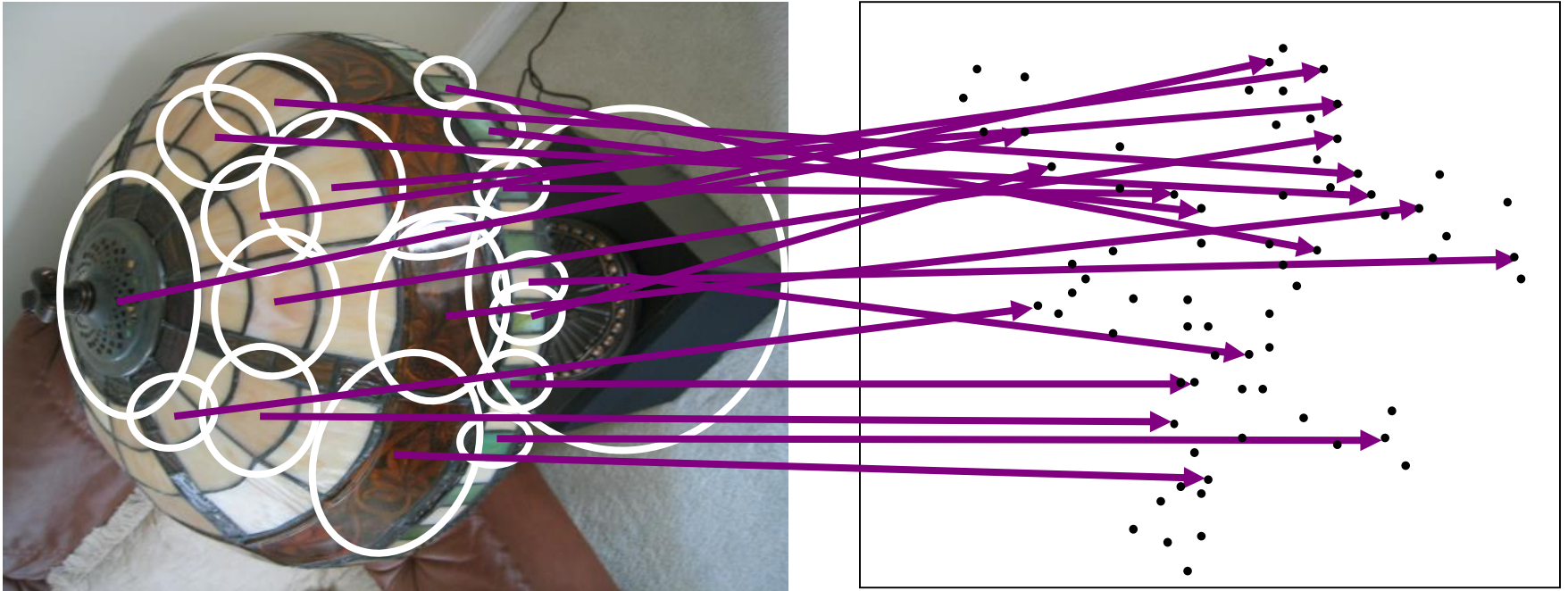


e.g., SIFT descriptor space:
each point is 128-dimensional

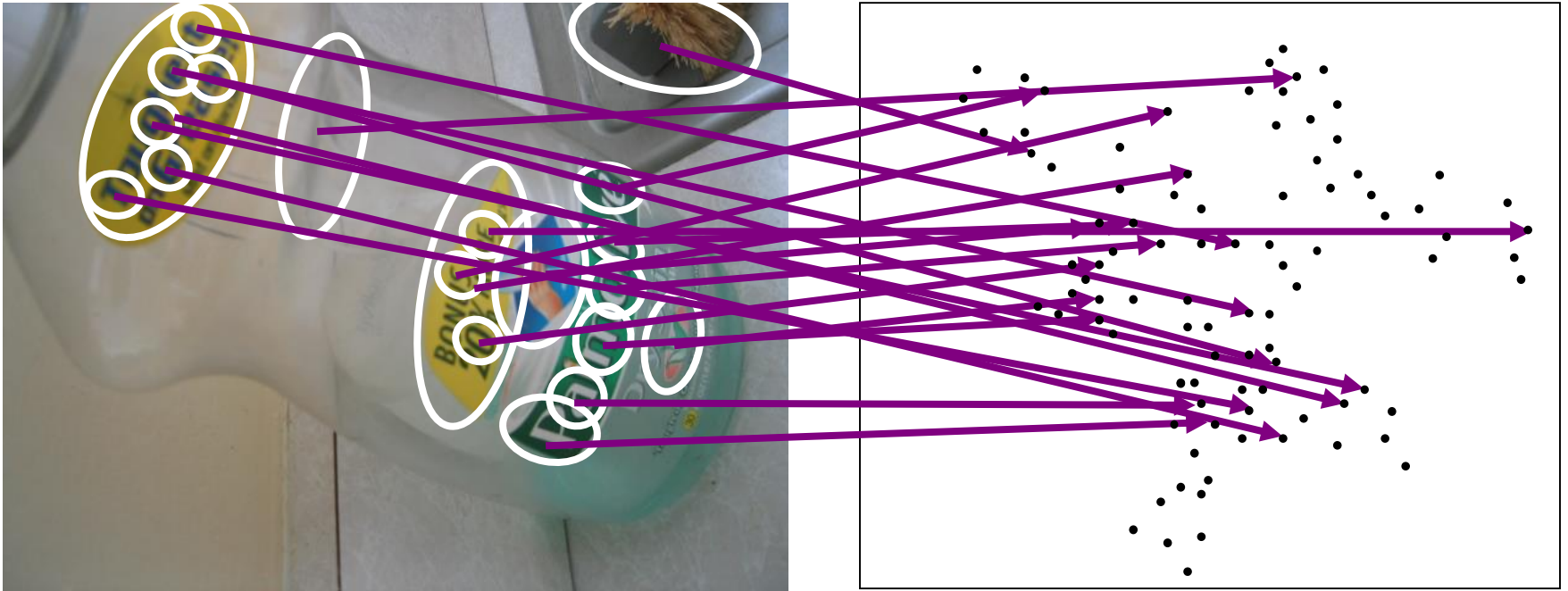
Extract local features from images



Extract local features from images

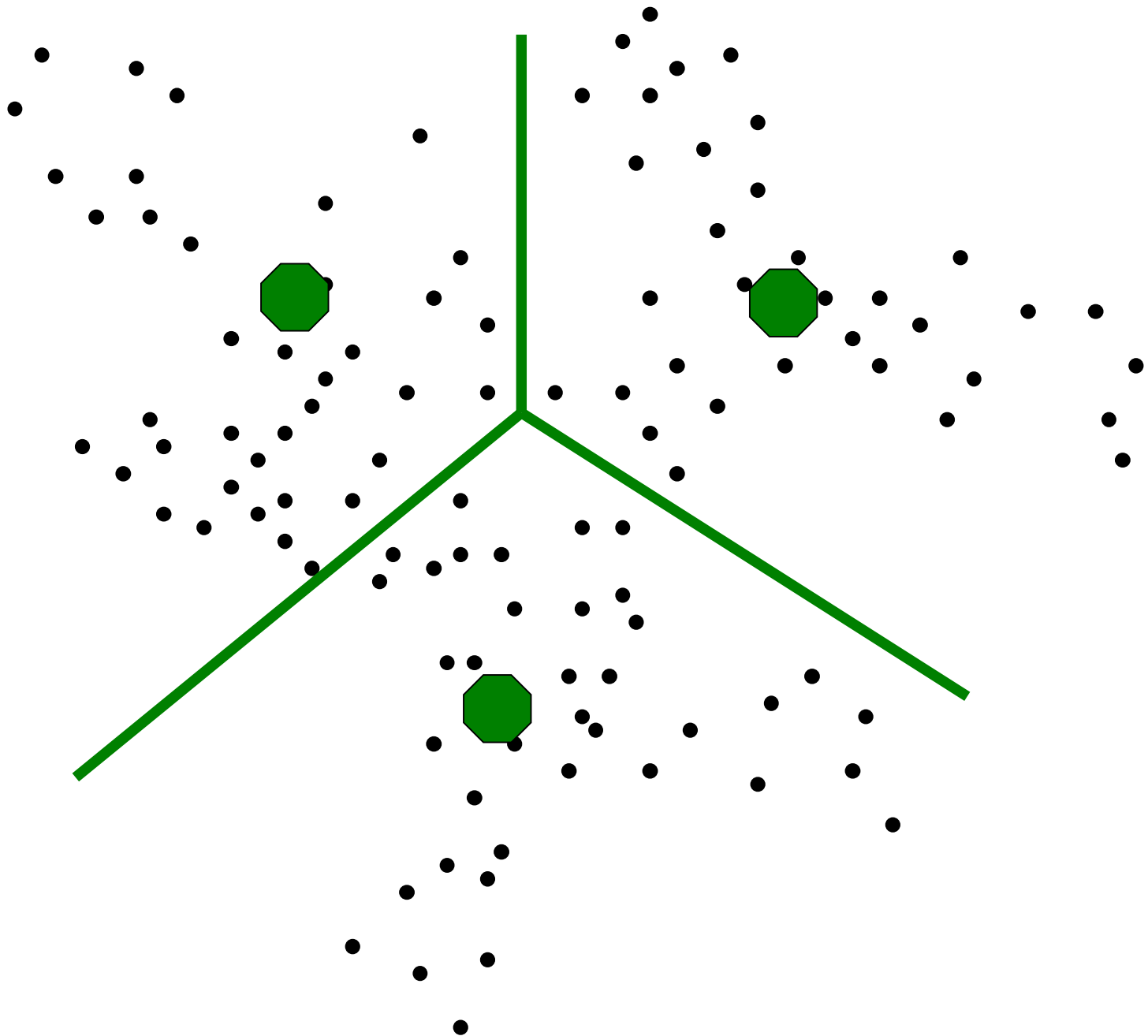


Extract local features from images



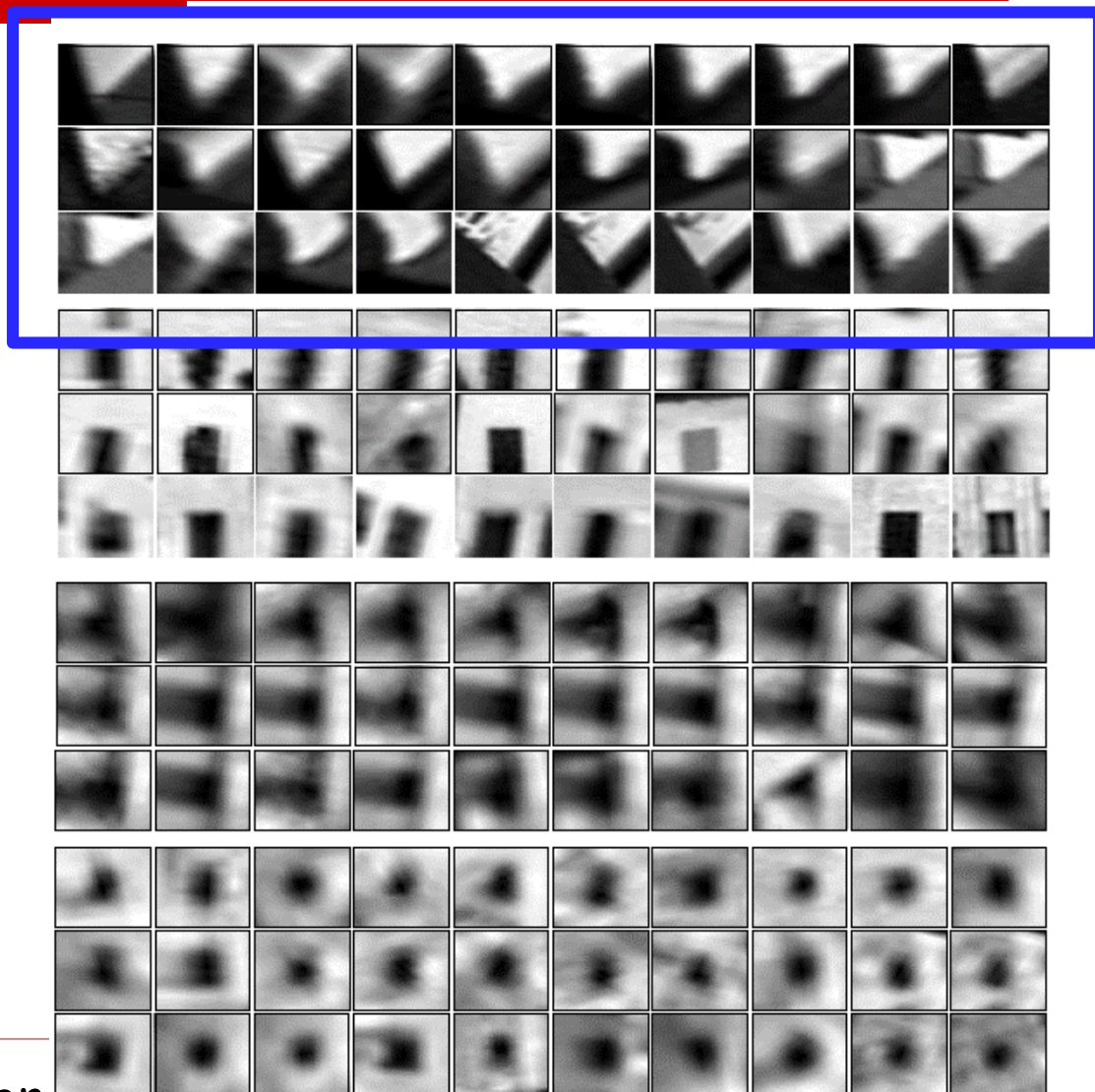
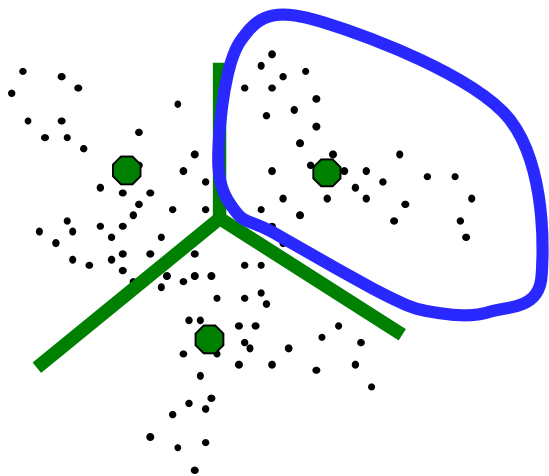


Each point is a
local descriptor,
e.g. SIFT vector.



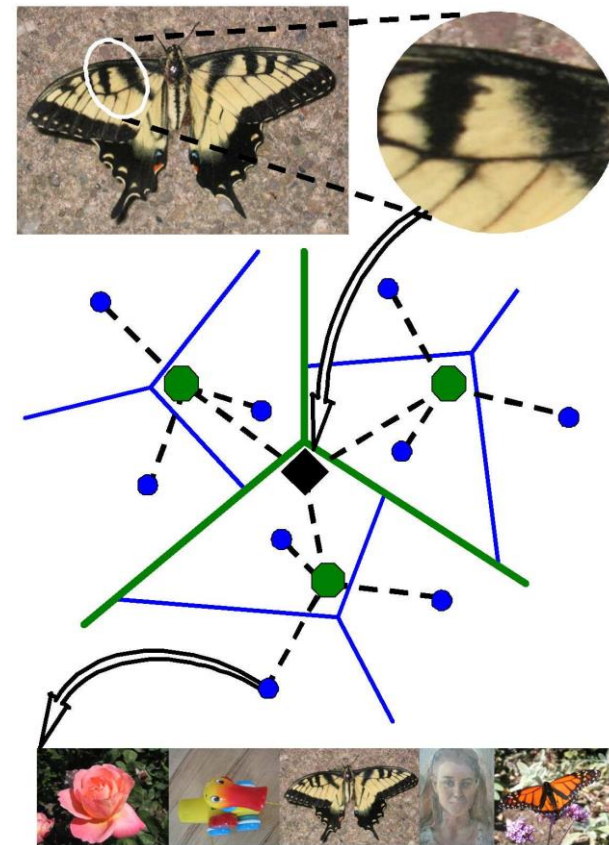
Visual words

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

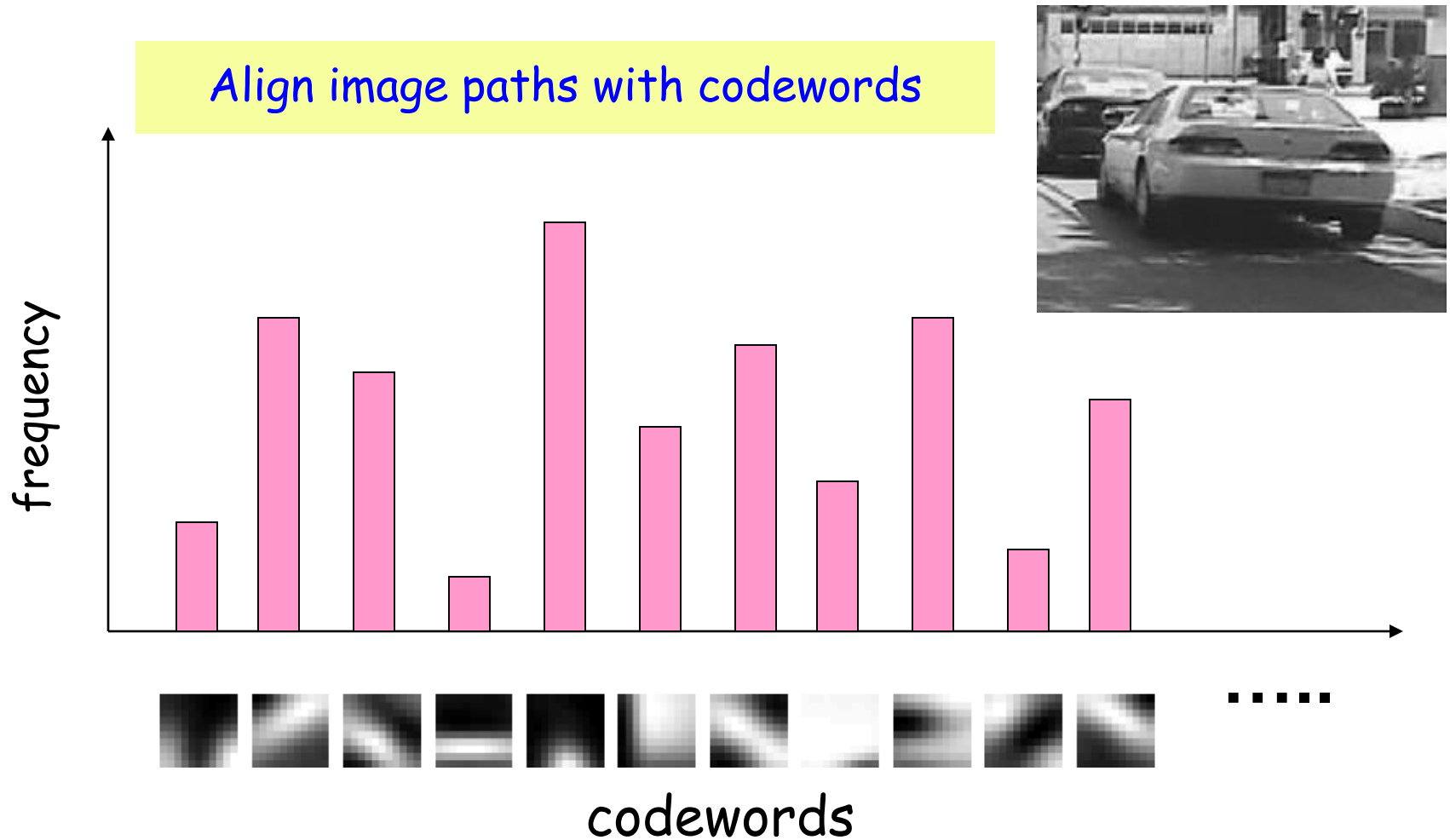


Visual vocabularies: Issues

- How to choose vocabulary size?
 - Too small: visual words not representative of all patches
 - Too large: quantization artifacts, overfitting
- Computational efficiency
 - Vocabulary trees
(Nister & Stewenius, 2006)

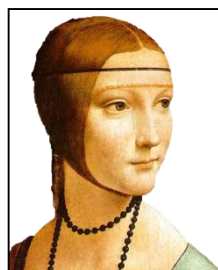
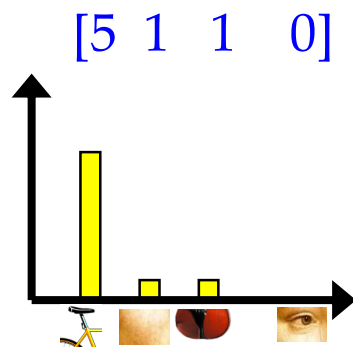
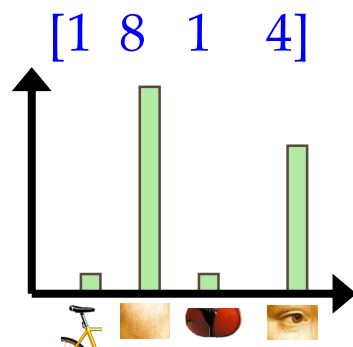


Step-3: Image representation

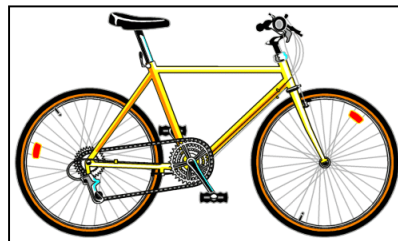


Step-4: Comparing bags of words

- Rank images by normalized inner product between their (possibly weighted) occurrence counts---*nearest neighbor* search



\vec{d}_j



\vec{q}

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

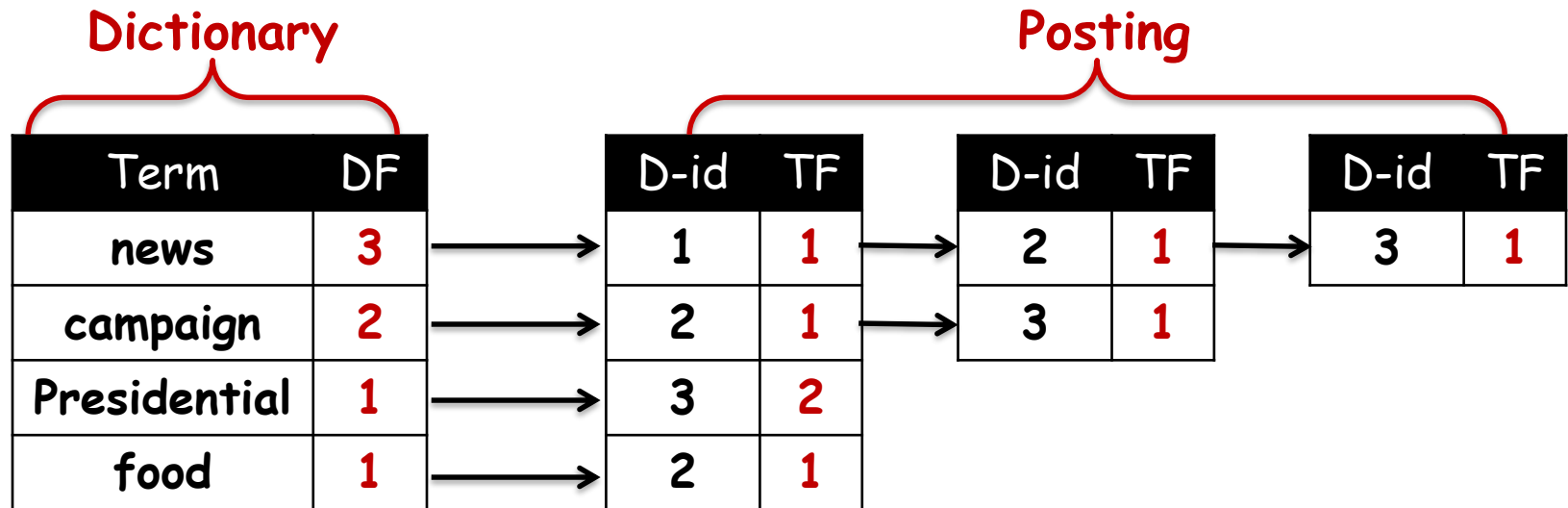
$$= \frac{\sum_{i=1}^V d_j(i) * q(i)}{\sqrt{\sum_{i=1}^V d_j(i)^2} * \sqrt{\sum_{i=1}^V q(i)^2}}$$

for vocabulary of V words

Inverted indexing for faster BoW similarity computation



An Example of Inverted Index

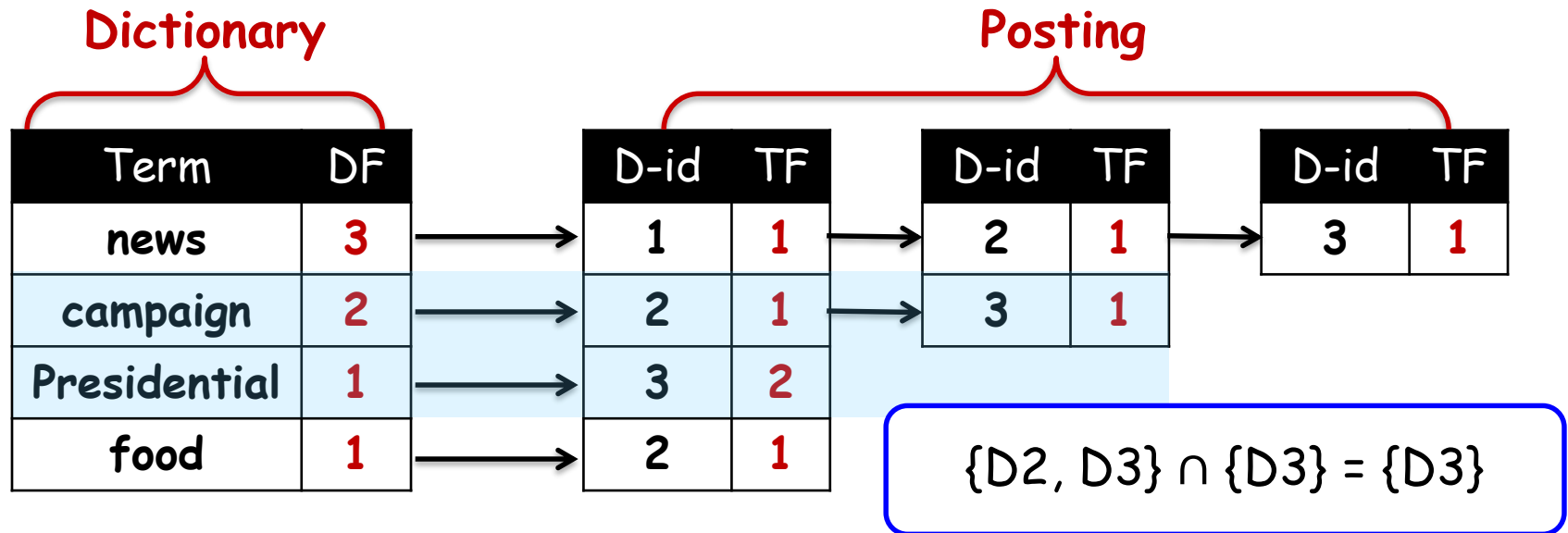


D1 ... news about

D2 ... news about organic food campaign...

D3 ... news of presidential campaign ... presidential candidate ...

An Example of Inverted Index



- ❑ Transform **Multi-term query** into **Boolean** operation
 - Match term "A" **AND/OR** term "B"
 - E.g. query="presidential campaign"

Inverted indexing for BoVW



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. Rank images with word counts

How to score the retrieval results?



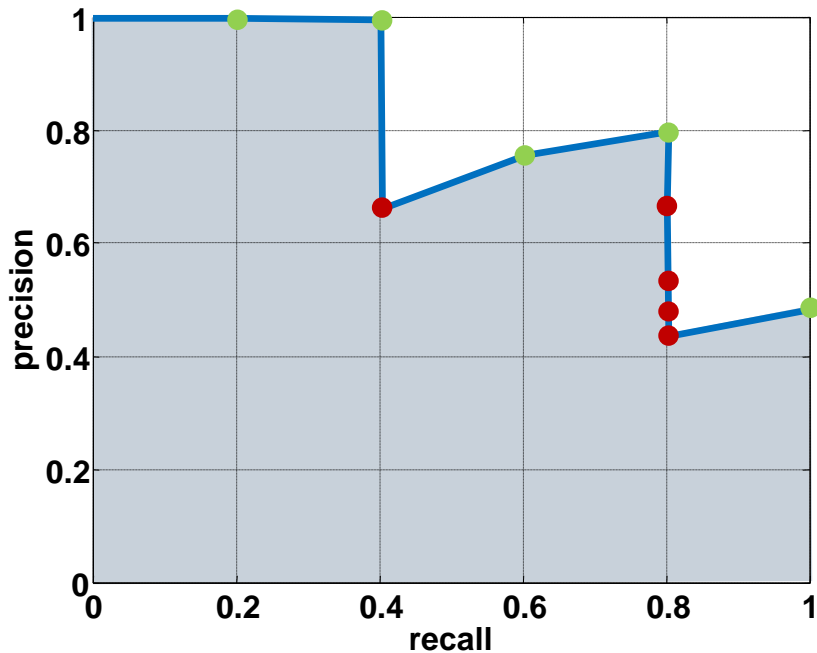


Query

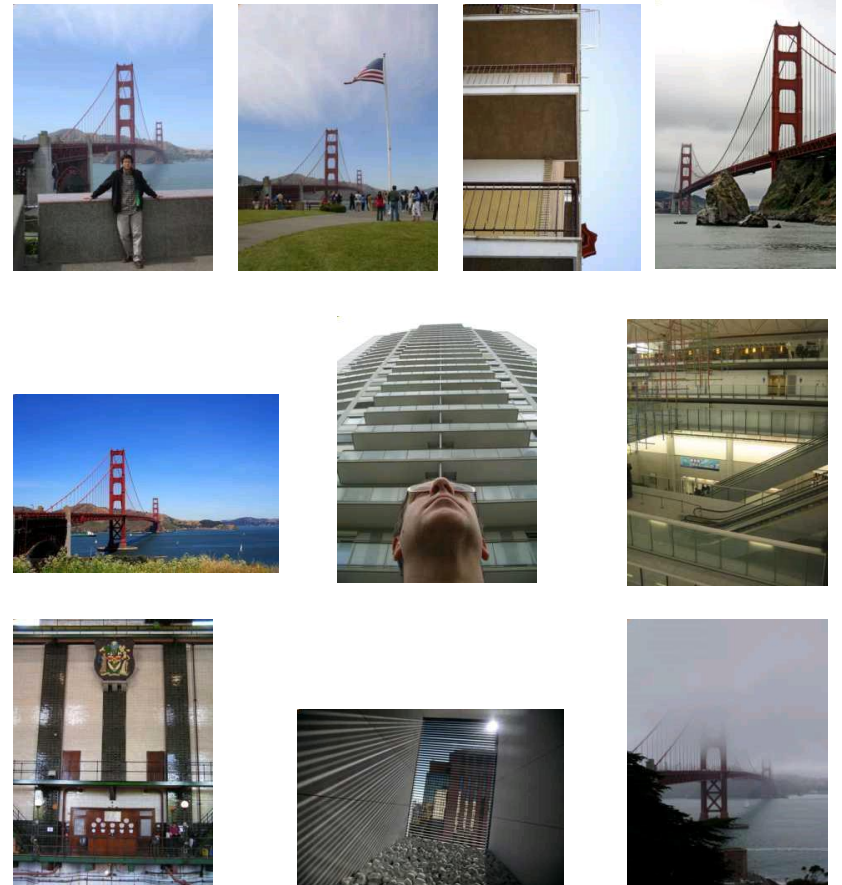
Scoring retrieval quality

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

$\text{precision} = \# \text{relevant} / \# \text{returned}$
 $\text{recall} = \# \text{relevant} / \# \text{total relevant}$



Results (ordered):



Summary on BoW Model

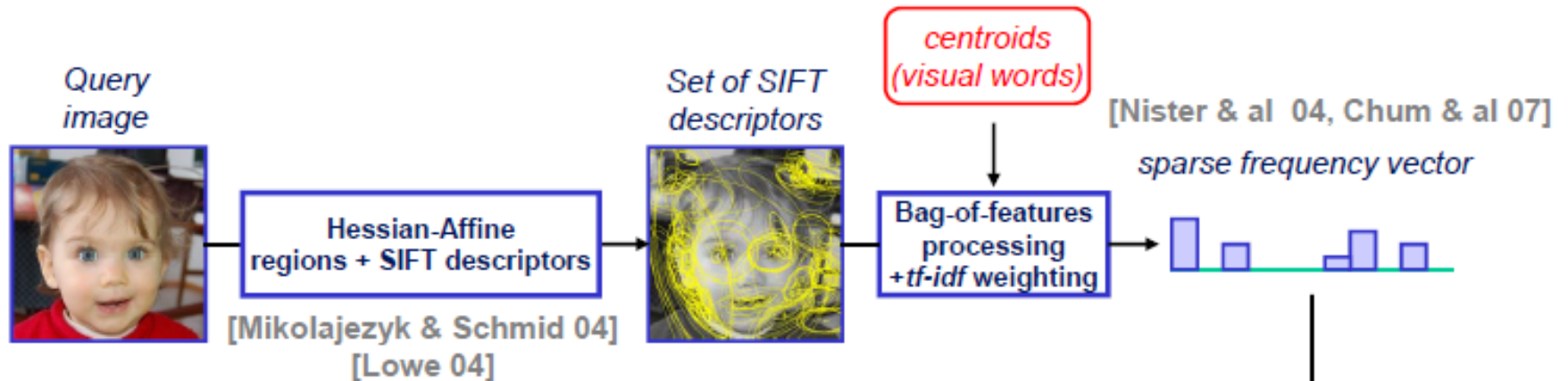
Positives

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

Negatives

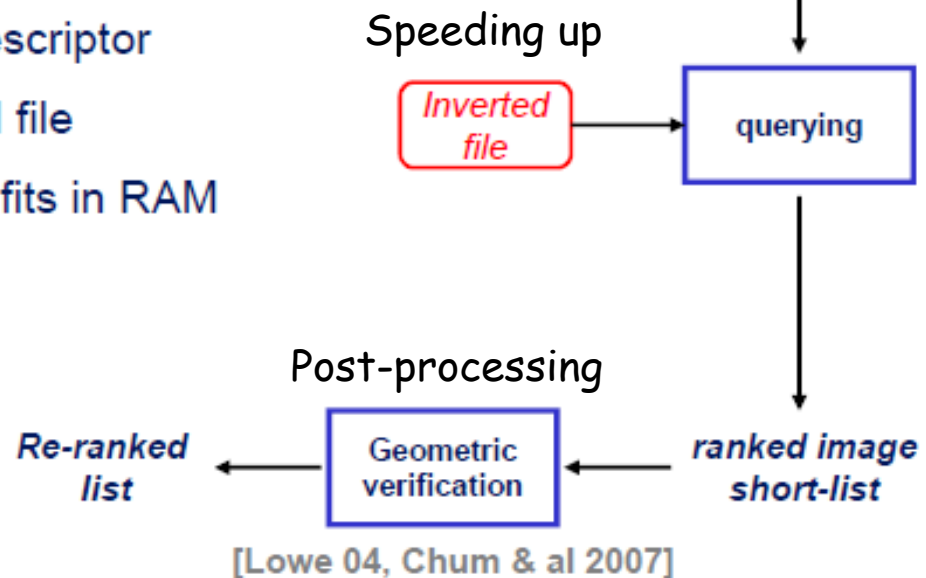
- background and foreground mixed when bag covers whole image
- optimal vocabulary formation remains unclear

State-of-the-Art BoVW Flowchart for Search



Visual Words

- 1 word (index) per local descriptor
 - only images ids in inverted file
- ⇒ 8 GB for a million images, fits in RAM



Search on 10.2 M Image Database

Demo of Similar Image Search