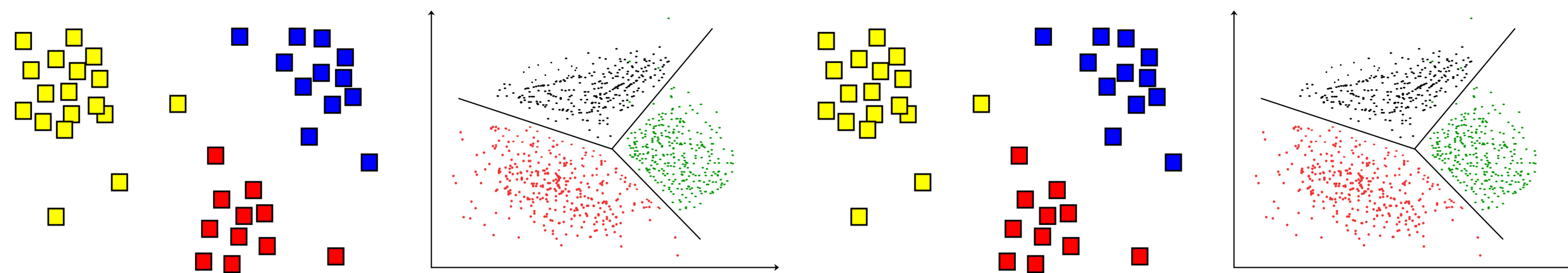


CPSC 425: Computer Vision



Lecture 29: Image Classification

Menu for Today (November 16, 2018)

Topics:

- Scene Classification
- Bag of Words Representation
- Decision Tree
- Boosting

Readings:

- **Today's** Lecture: Forsyth & Ponce (2nd ed.) 16.1.3, 16.1.4, 16.1.9
- **Next** Lecture: Forsyth & Ponce (2nd ed.) 17.1–17.2

Reminders:

- **Assignment 5:** Scene Recognition with Bag of Words due **last day of classes**

Image Classification

We next discuss **image classification**, where we pass a whole image into a classifier and obtain a class label as output.

What Makes Image Classification **Hard**?



Intra-class variation, viewpoint, illumination, clutter, and occlusion (among others!)

Image Classification

In addition to images containing single objects, the same techniques can be applied to classify natural scenes (e.g. beach, forest, harbour, library).

Why might classifying scenes be useful?

Image Classification

In addition to images containing single objects, the same techniques can be applied to classify natural scenes (e.g. beach, forest, harbour, library).

Why might classifying scenes be useful?

Visual perception is influenced by expectation. Our expectations are often conditioned on the **context**.



What is This **Object**?



What is This **Object**?



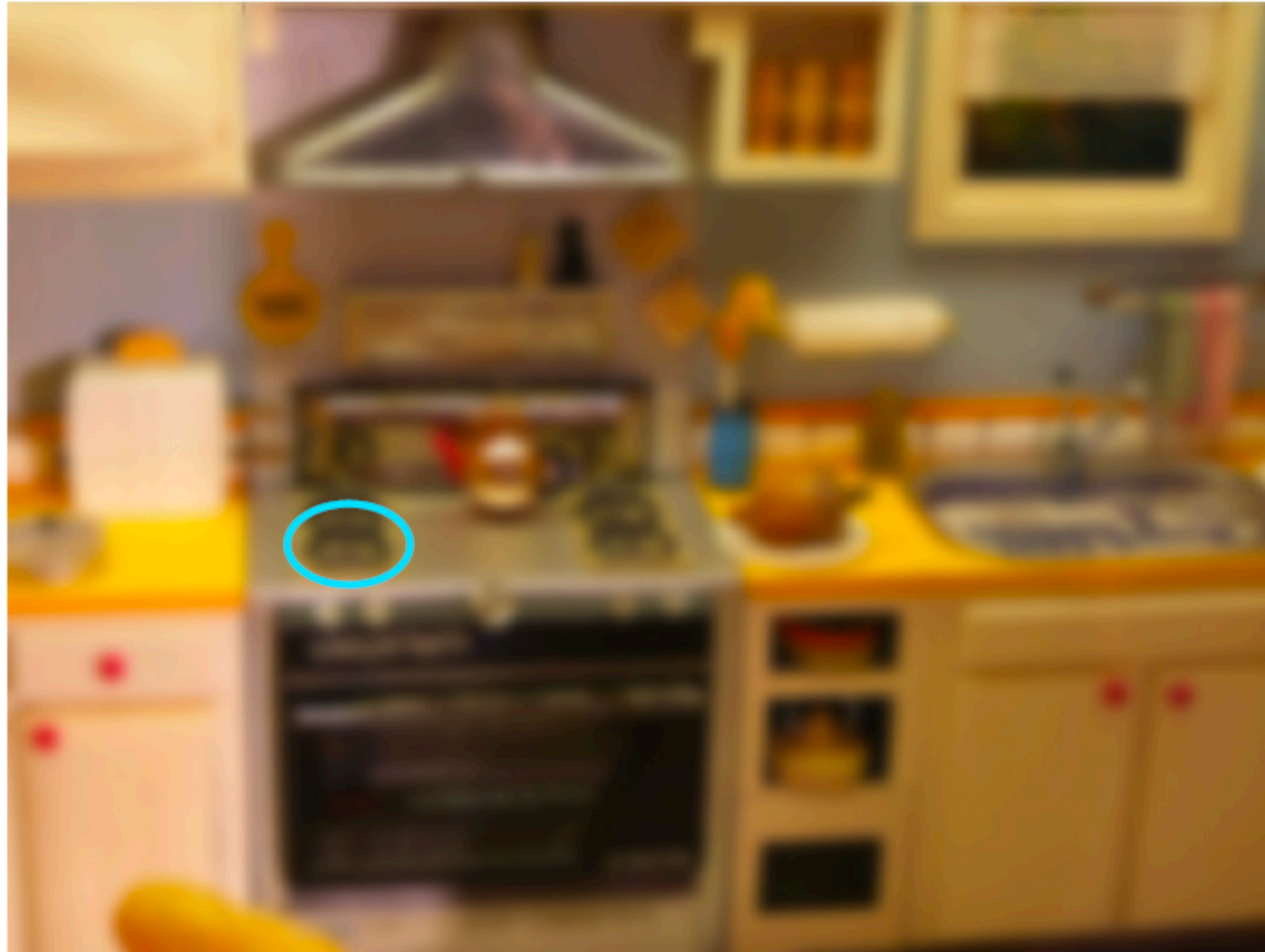
What is This **Object**?



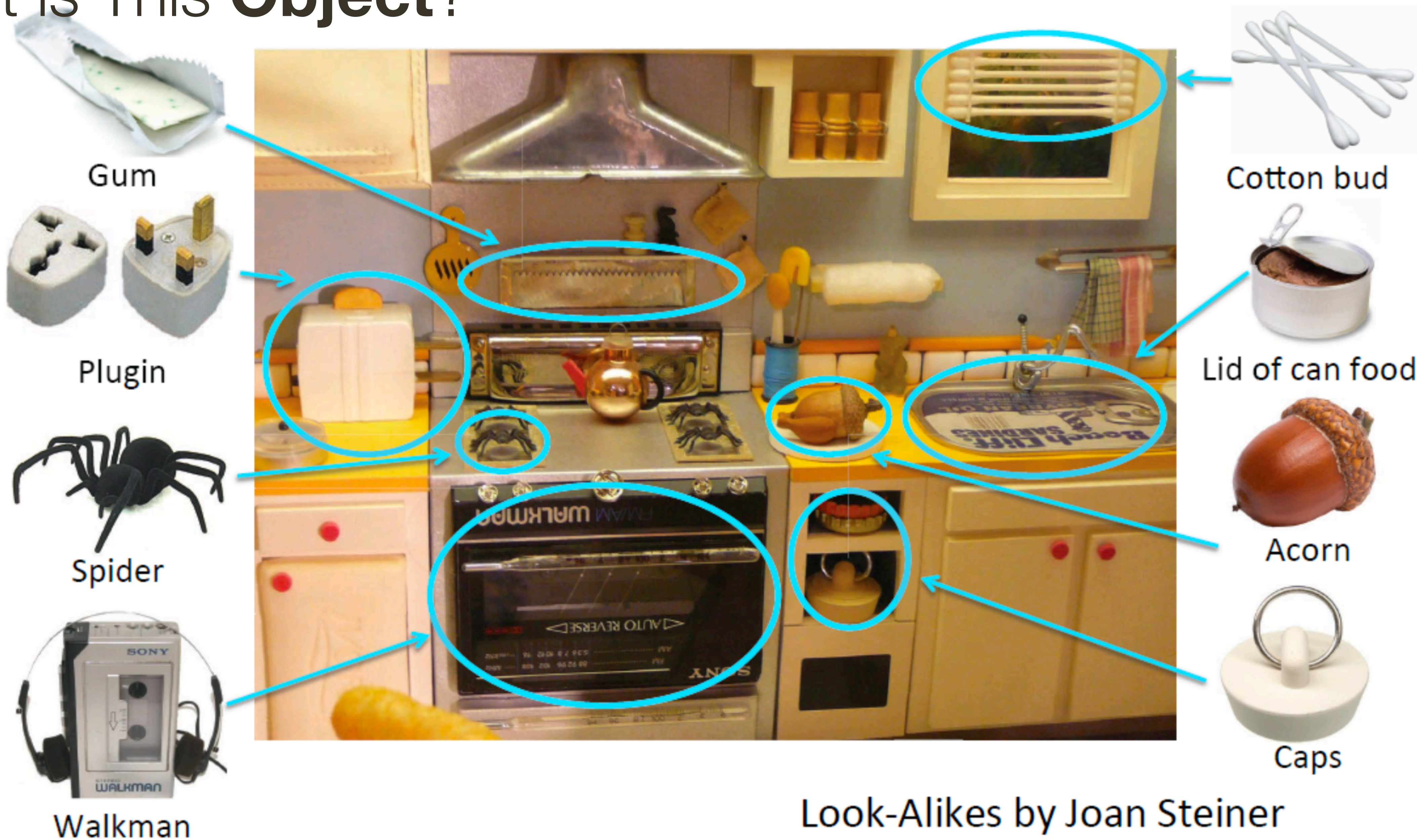
What is This **Object**?



What is This **Object**?



What is This **Object**?



Look-Alikes by Joan Steiner

Figure source: Jianxiong Xiao

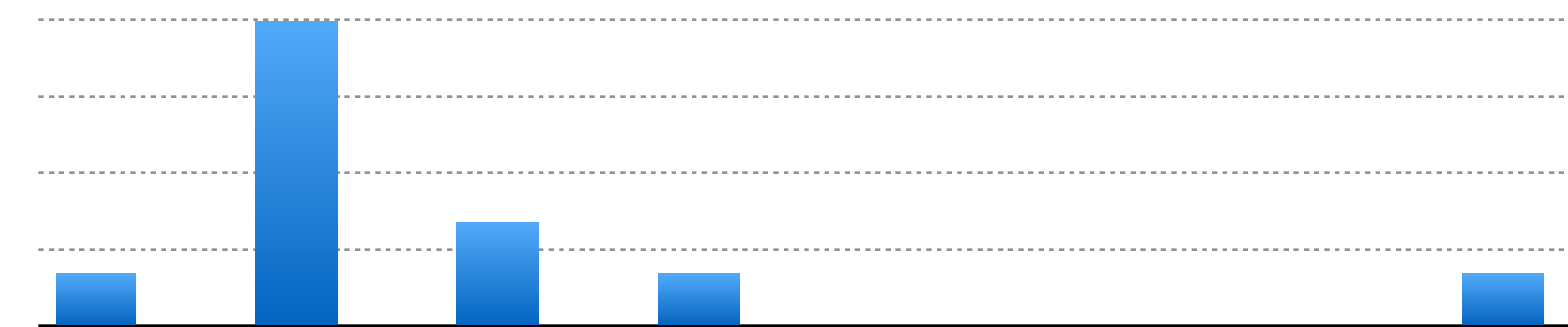
Visual **Words**

Many algorithms for image classification accumulate evidence on the basis of **visual words**.

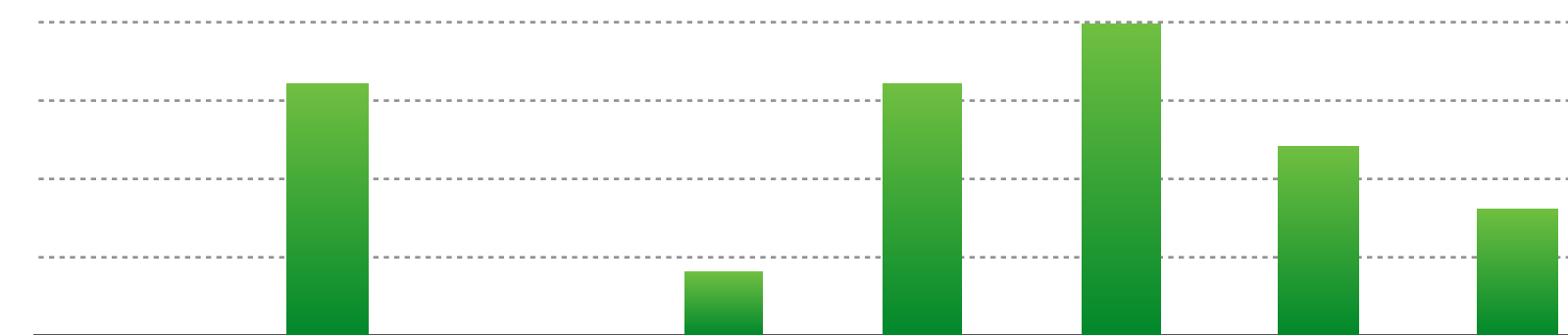
To classify a text document (e.g. as an article on sports, entertainment, business, politics) we might find patterns in the occurrences of certain words.

Vector Space Model

G. Salton. ‘Mathematics and Information Retrieval’ Journal of Documentation, 1979



1	6	2	1	0	0	0	1
Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor



0	4	0	1	4	5	3	2
Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor

<http://www.fodey.com/generators/newspaper/snippet.asp>

Vector Space Model

A document (datapoint) is a vector of counts over each word (feature)

$$\mathbf{v}_d = [n(w_{1,d}) \quad n(w_{2,d}) \quad \cdots \quad n(w_{T,d})]$$

$n(\cdot)$ counts the number of occurrences

just a histogram over words

What is the similarity between two documents?



Vector Space Model

A document (datapoint) is a vector of counts over each word (feature)

$$\mathbf{v}_d = [n(w_{1,d}) \quad n(w_{2,d}) \quad \cdots \quad n(w_{T,d})]$$

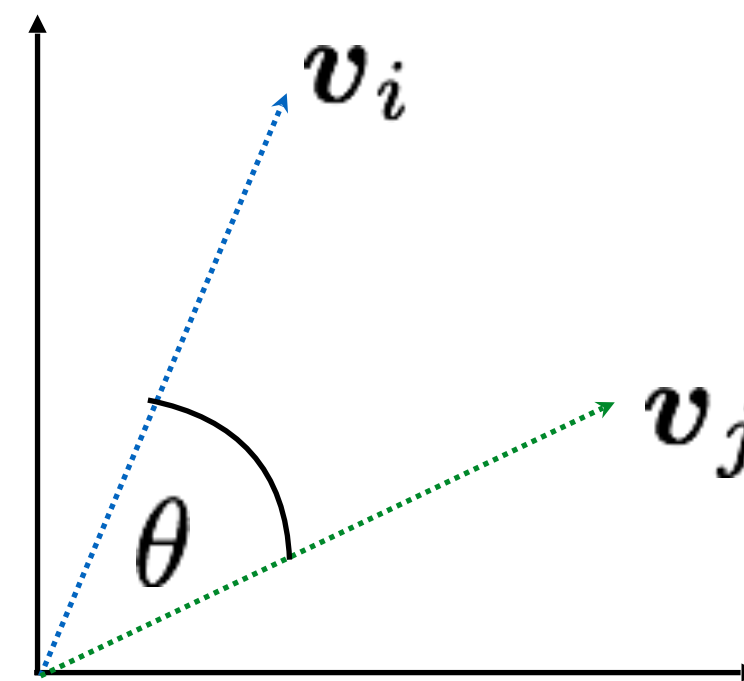
$n(\cdot)$ counts the number of occurrences

just a histogram over words

What is the similarity between two documents?

Use any distance you want but the cosine distance is fast and well designed for high-dimensional vector spaces:

$$\begin{aligned} d(\mathbf{v}_i, \mathbf{v}_j) &= \cos \theta \\ &= \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|} \end{aligned}$$

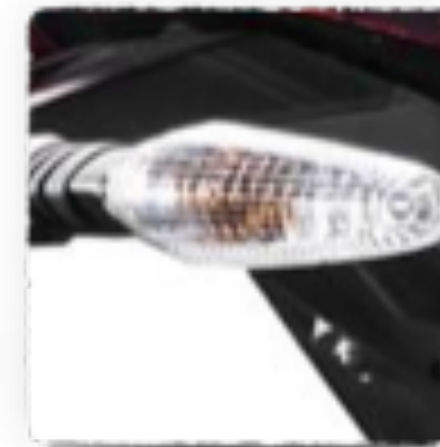
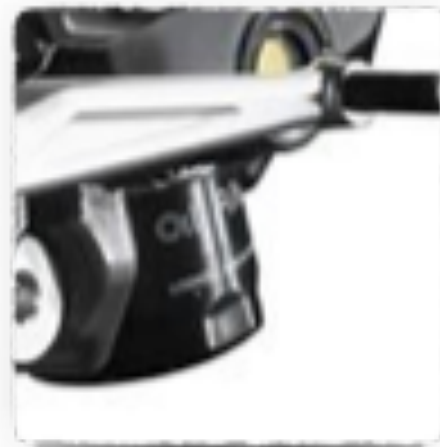


Visual **Words**

In images, the equivalent of a word is a local image patch. The local image patch is described using a descriptor such as SIFT.

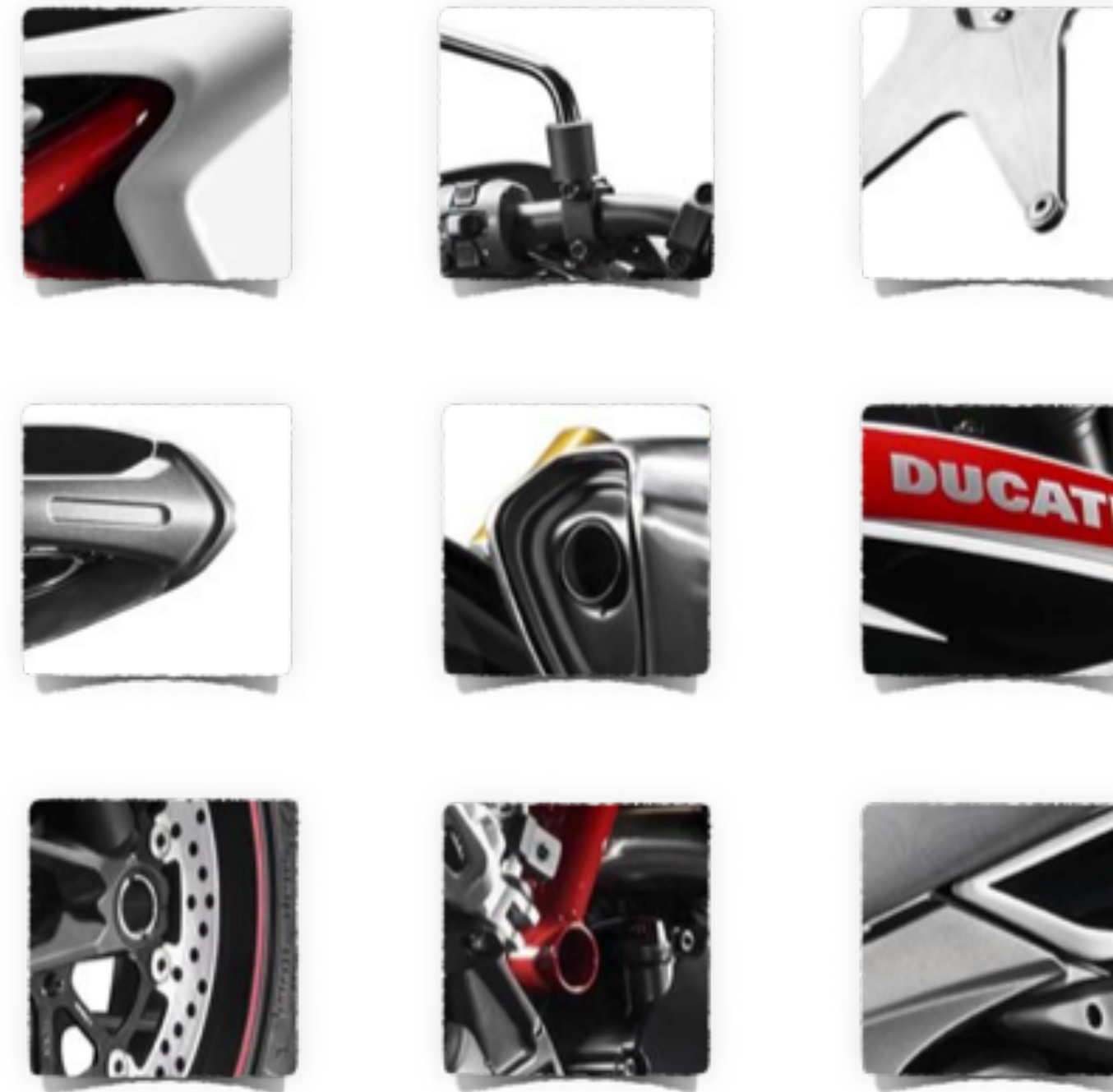
We construct a **vocabulary** or **codebook** of local descriptors, containing representative local descriptors.

What **Objects** do These Parts Belong To?



Some local feature are
very informative

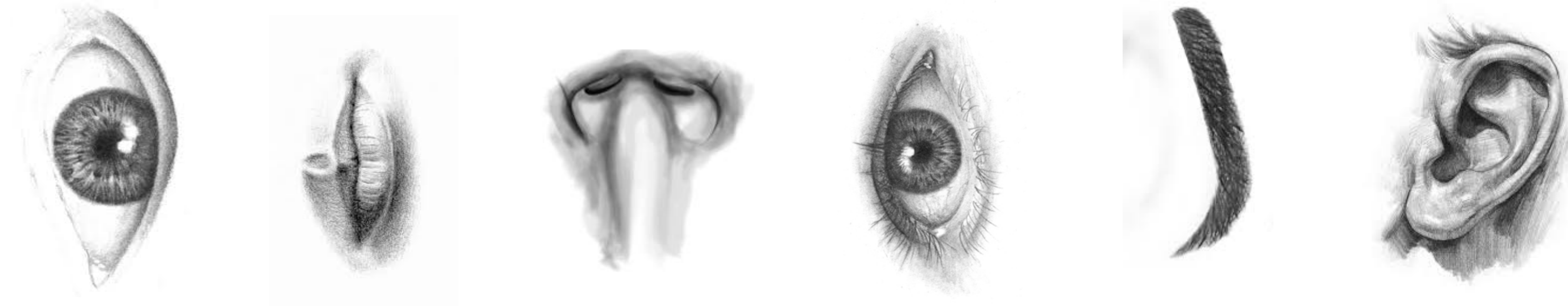
An object as



a collection of local features
(bag-of-features)

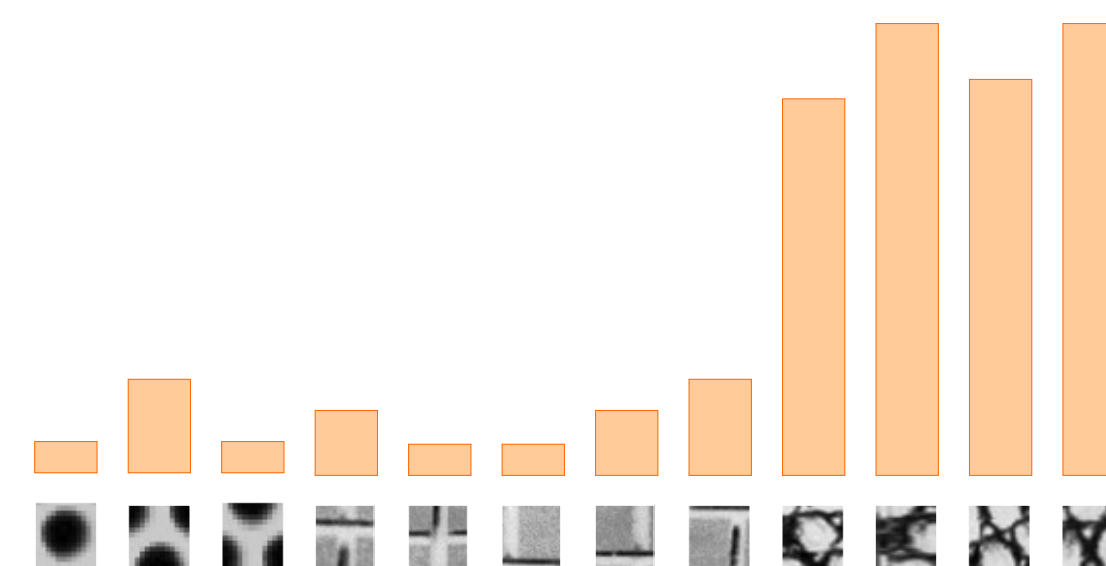
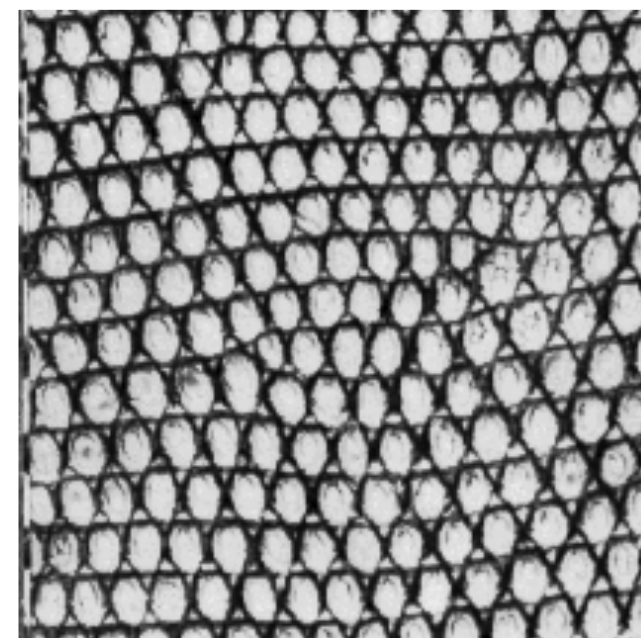
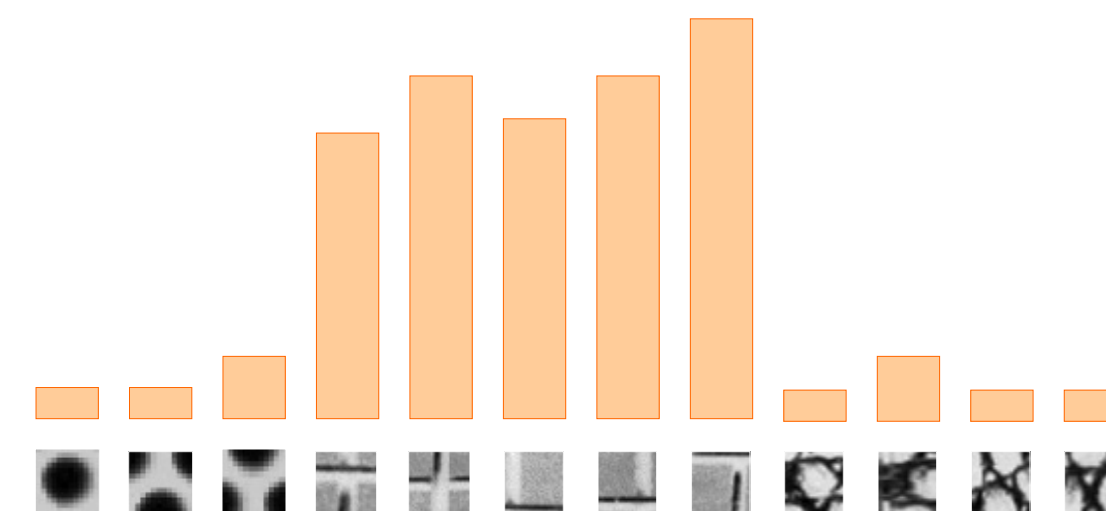
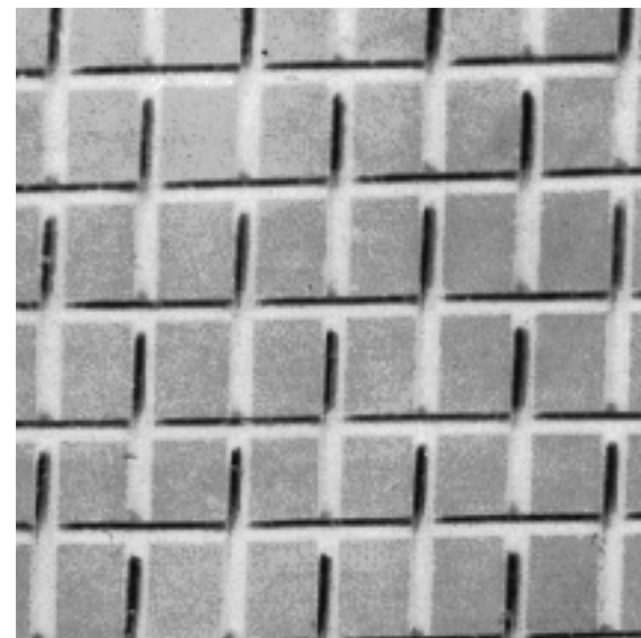
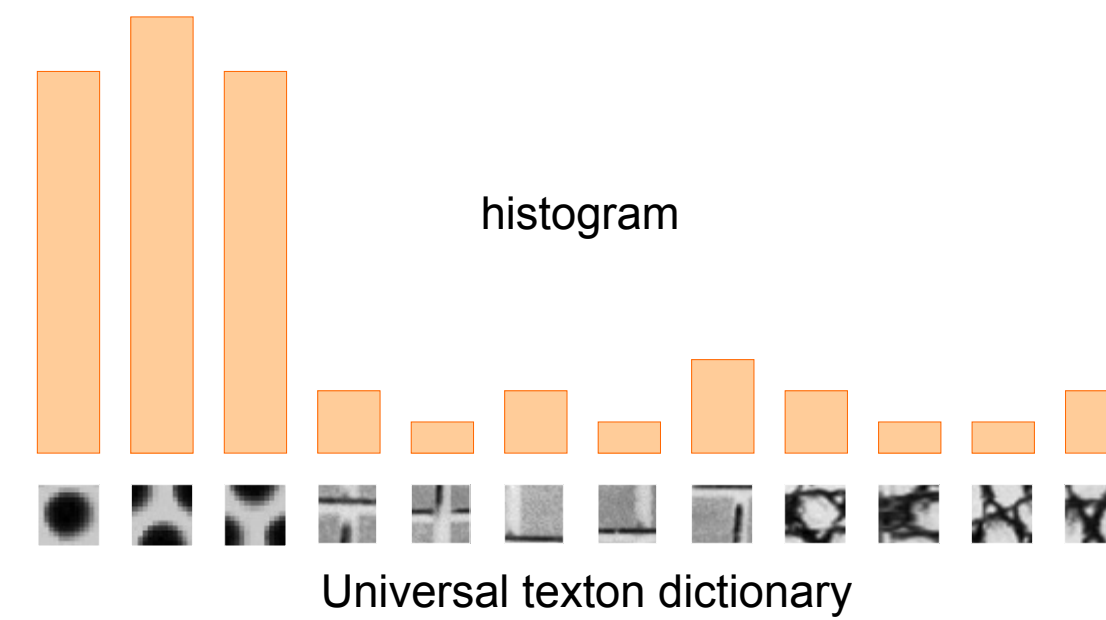
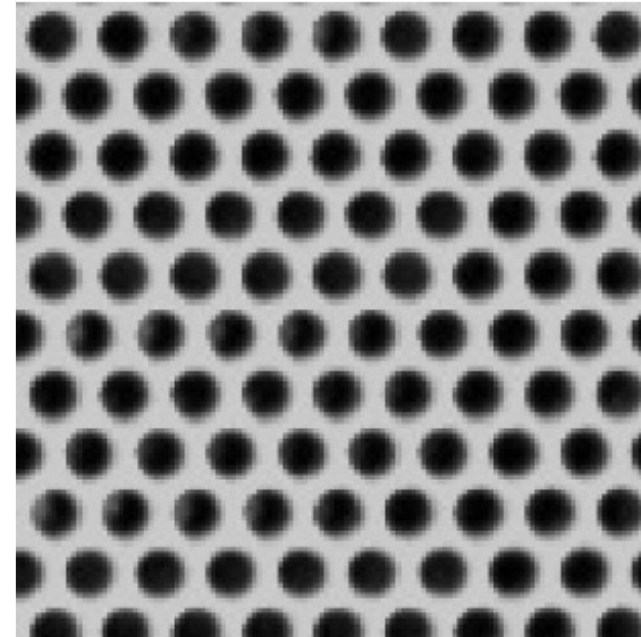
- deals well with occlusion
- scale invariant
- rotation invariant

(not so) Crazy Assumption



spatial information of local features
can be ignored for object recognition (i.e., verification)

Recall: Texture Representation



Visual **Words**

In images, the equivalent of a word is a local image patch. The local image patch is described using a descriptor such as SIFT.

We construct a **vocabulary** or **codebook** of local descriptors, containing representative local descriptors.

Question: How might we construct such a codebook? Given a large sample of SIFT descriptors, say 1 million, how can we choose a small number of ‘representative’ SIFT codewords, say 1000?

Standard **Bag-of-Words** Pipeline (for image classification)

Dictionary Learning:

Learn Visual Words using clustering

Encode:

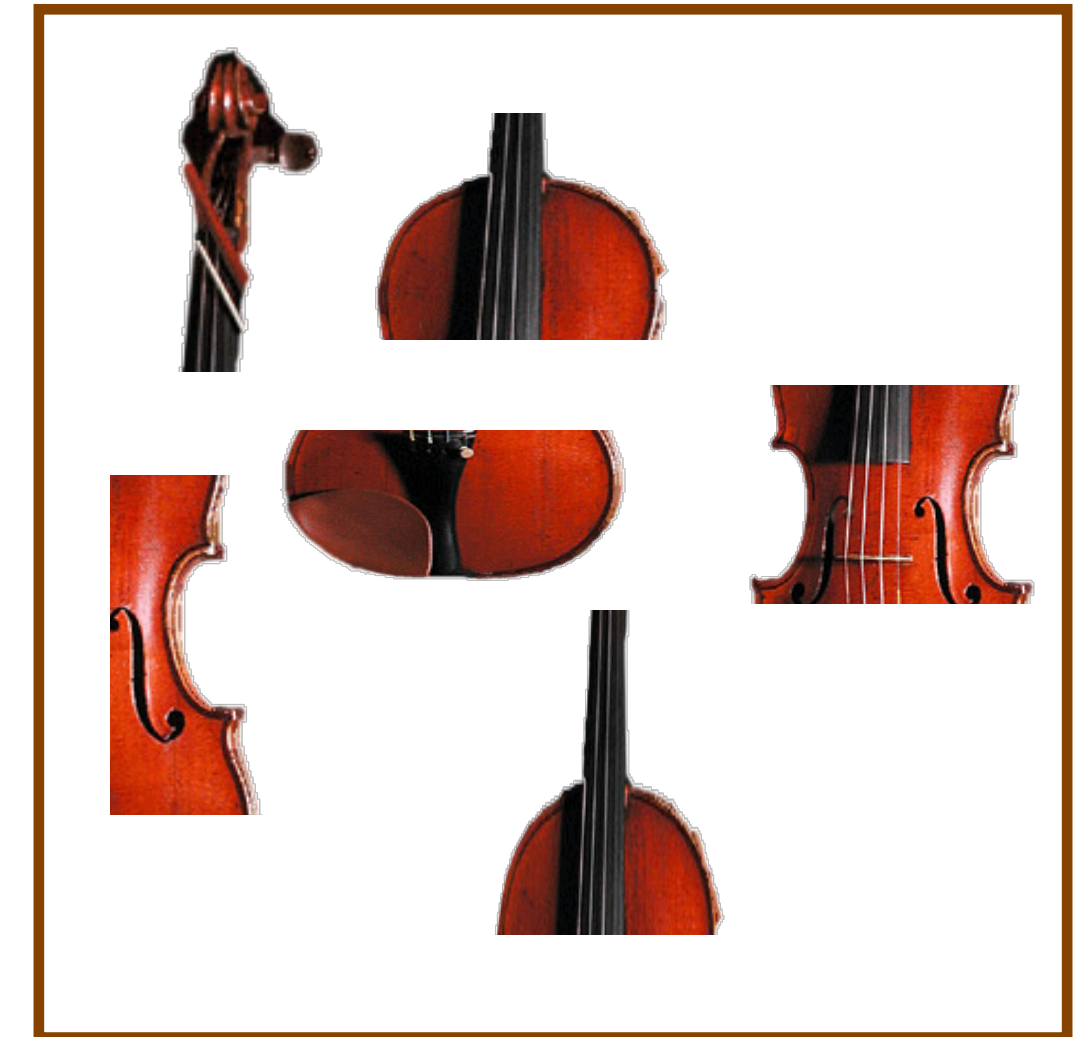
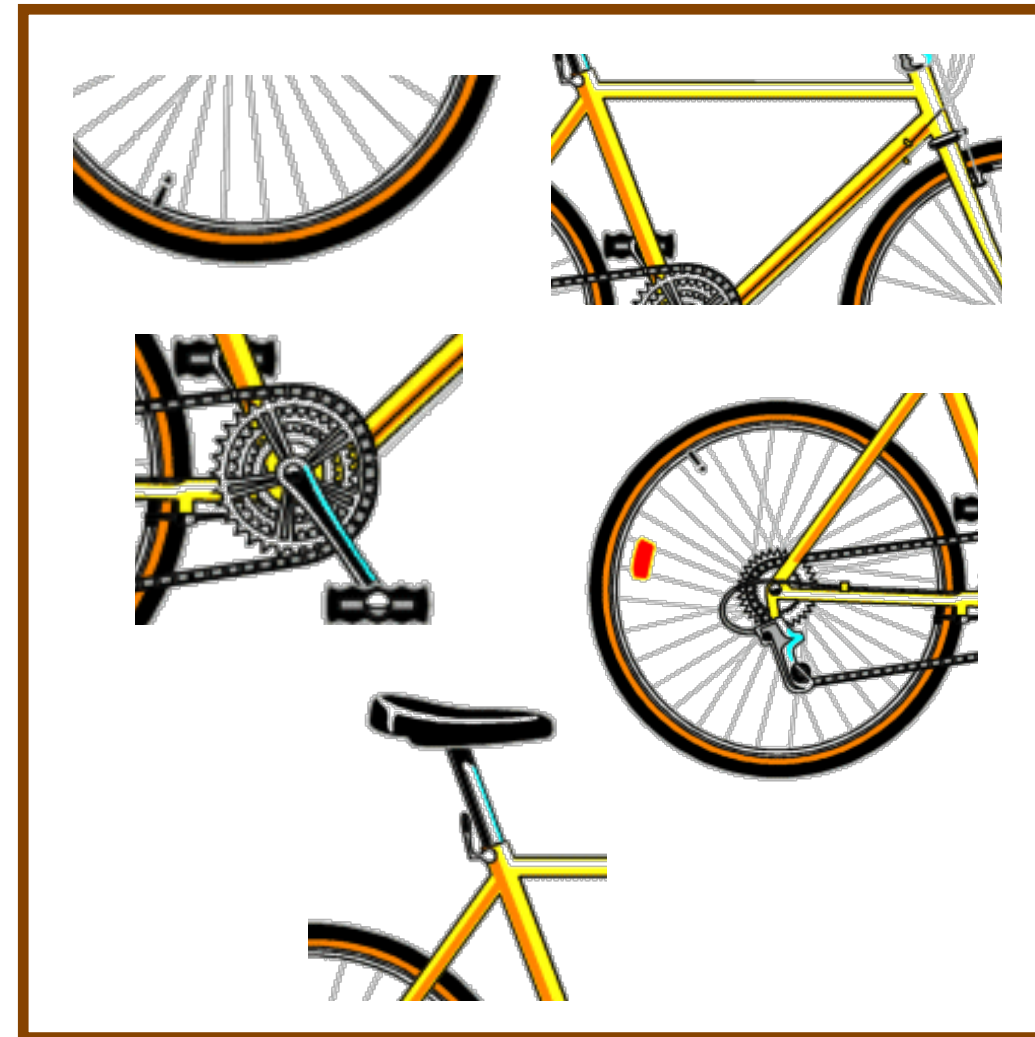
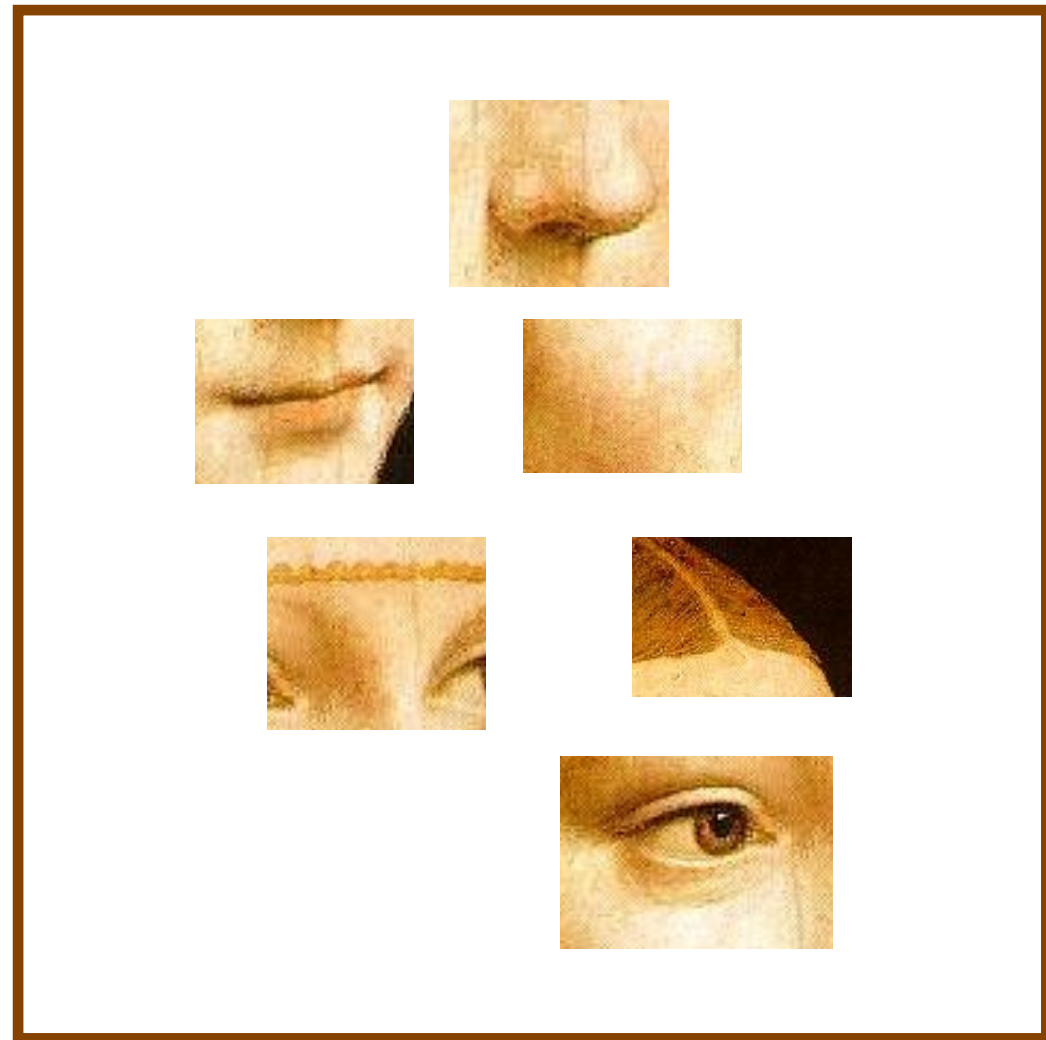
build Bags-of-Words (BOW) vectors
for each image

Classify:

Train and test data using BOWs

1. Dictionary Learning: Learn Visual Words using Clustering

1. **extract features** (e.g., SIFT) from images



1. Dictionary Learning: Learn Visual Words using Clustering

2. **Learn visual dictionary** (e.g., K-means clustering)



What **Features** Should We Extract?

- Regular grid

Vogel & Schiele, 2003

Fei-Fei & Perona, 2005

- Interest point detector

Csurka et al. 2004

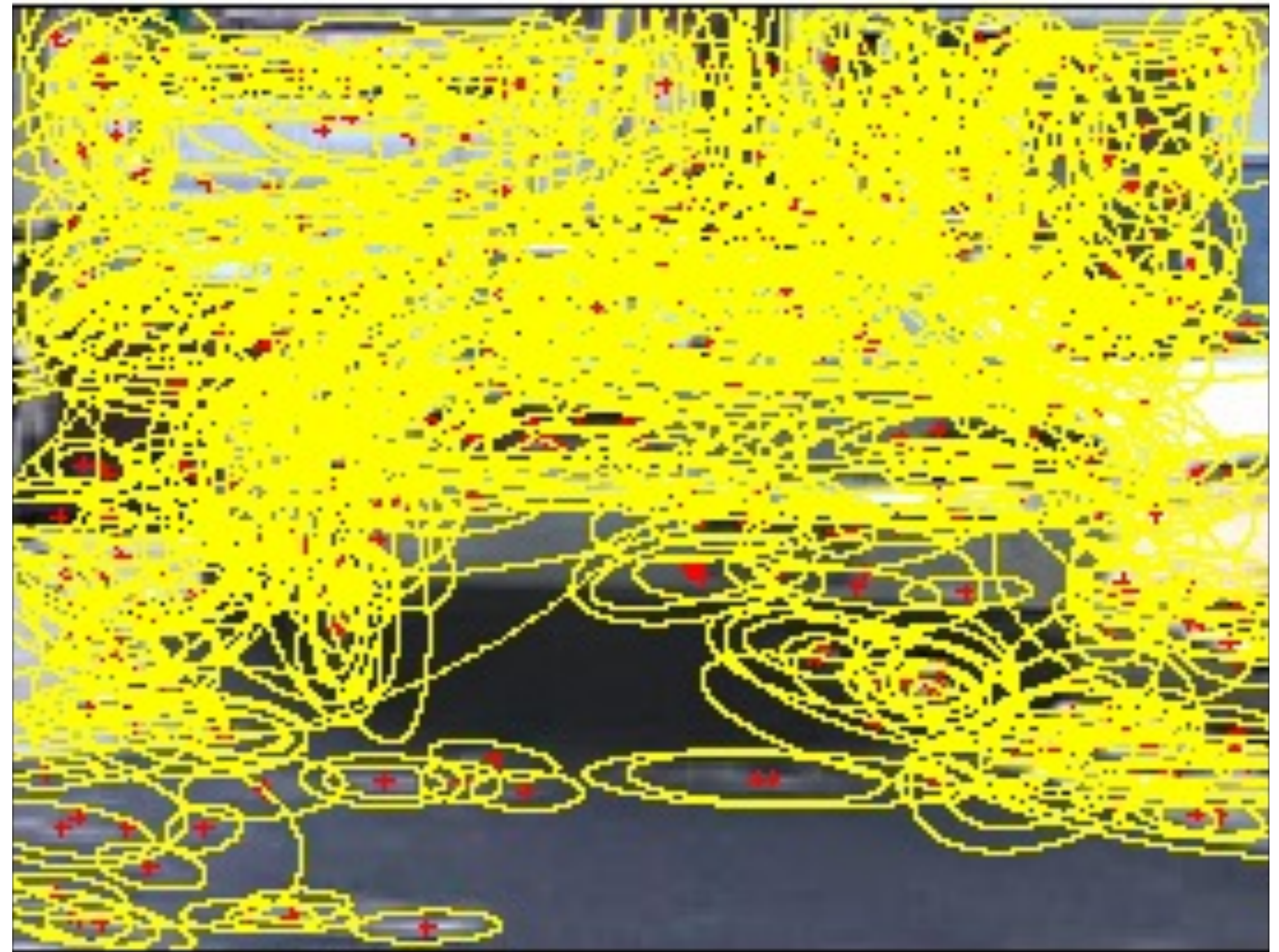
Fei-Fei & Perona, 2005

Sivic et al. 2005

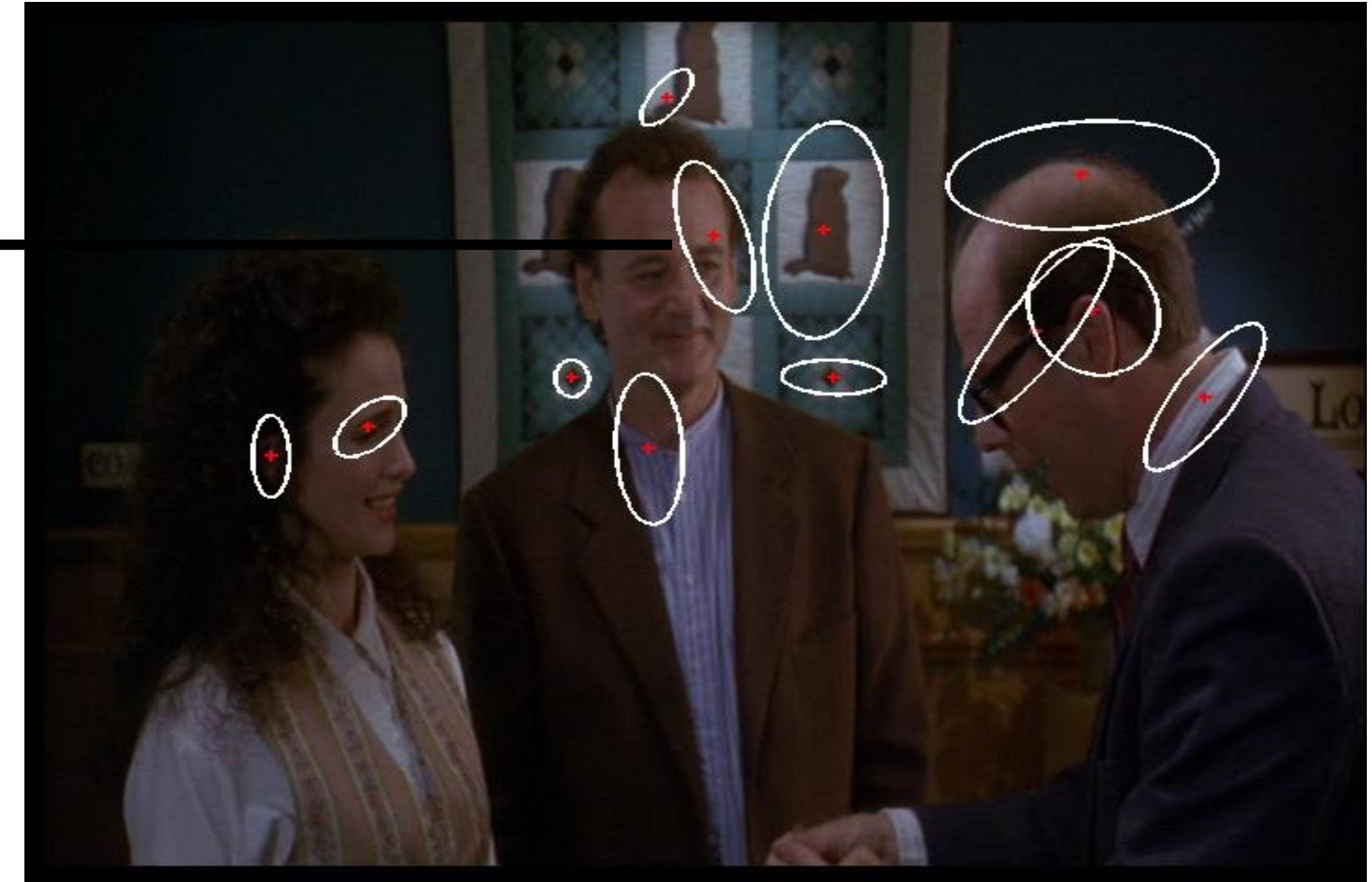
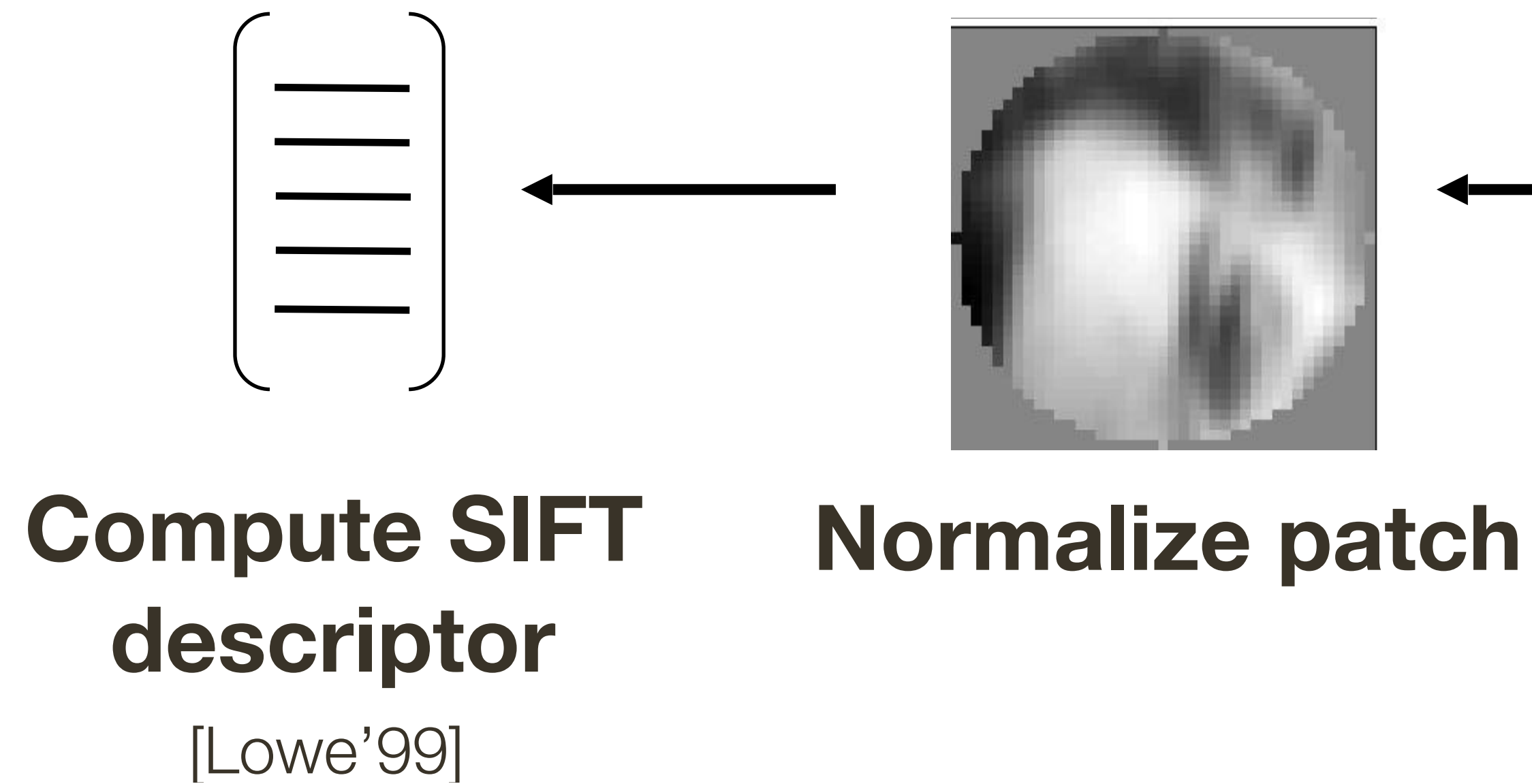
- Other methods

Random sampling (Vidal-Naquet & Ullman, 2002)

Segmentation-based patches (Barnard et al. 2003)



Extracting **SIFT** Patches



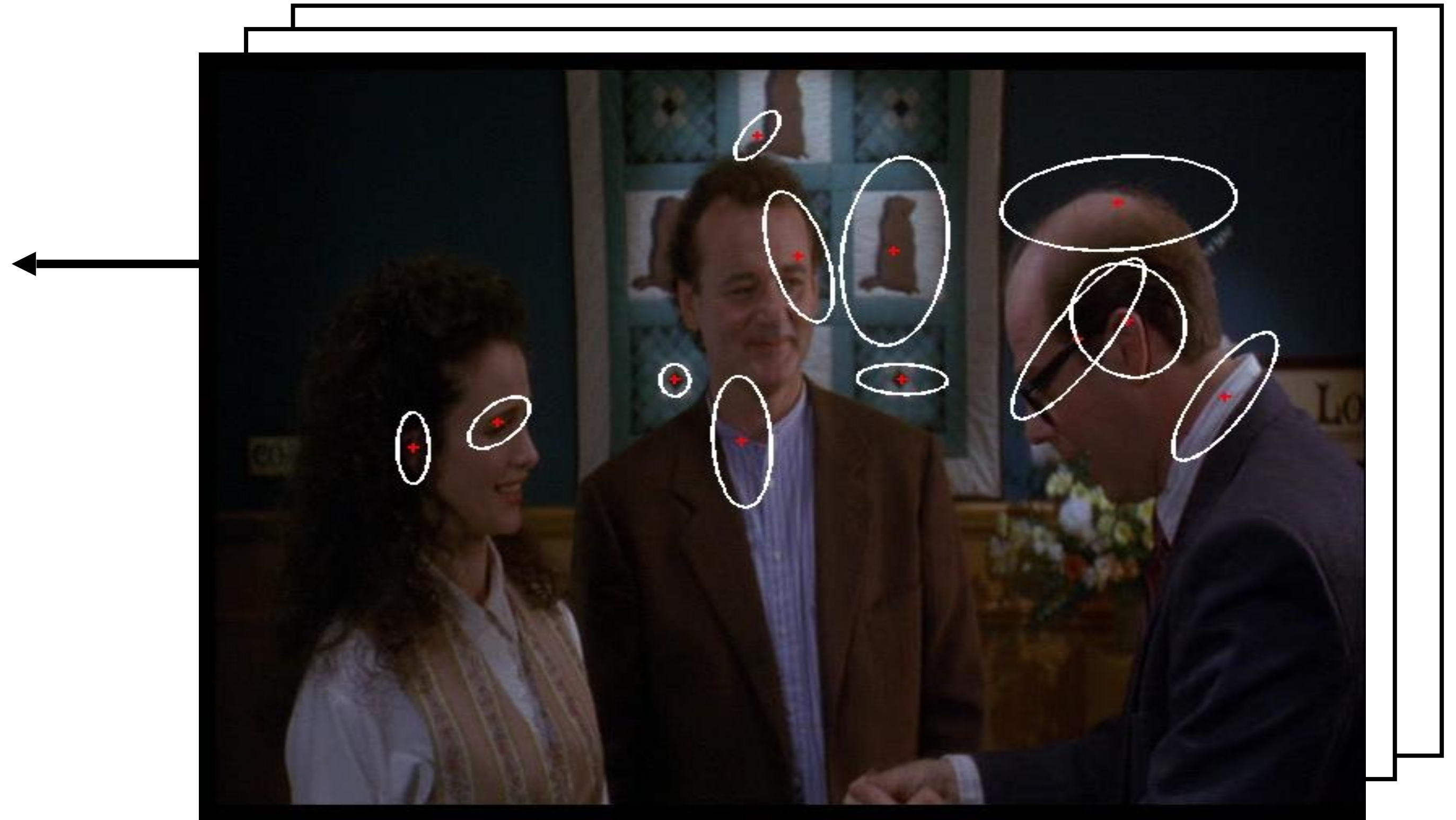
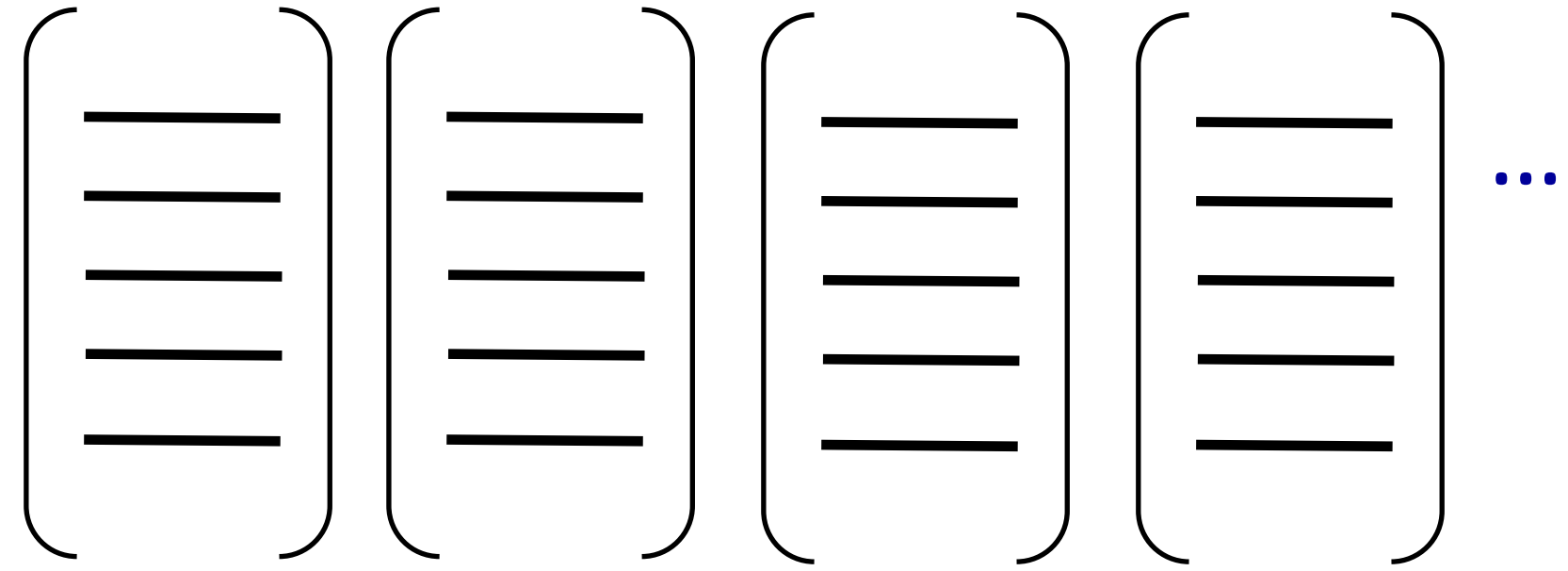
Detect patches

[Mikojaczuk and Schmid '02]

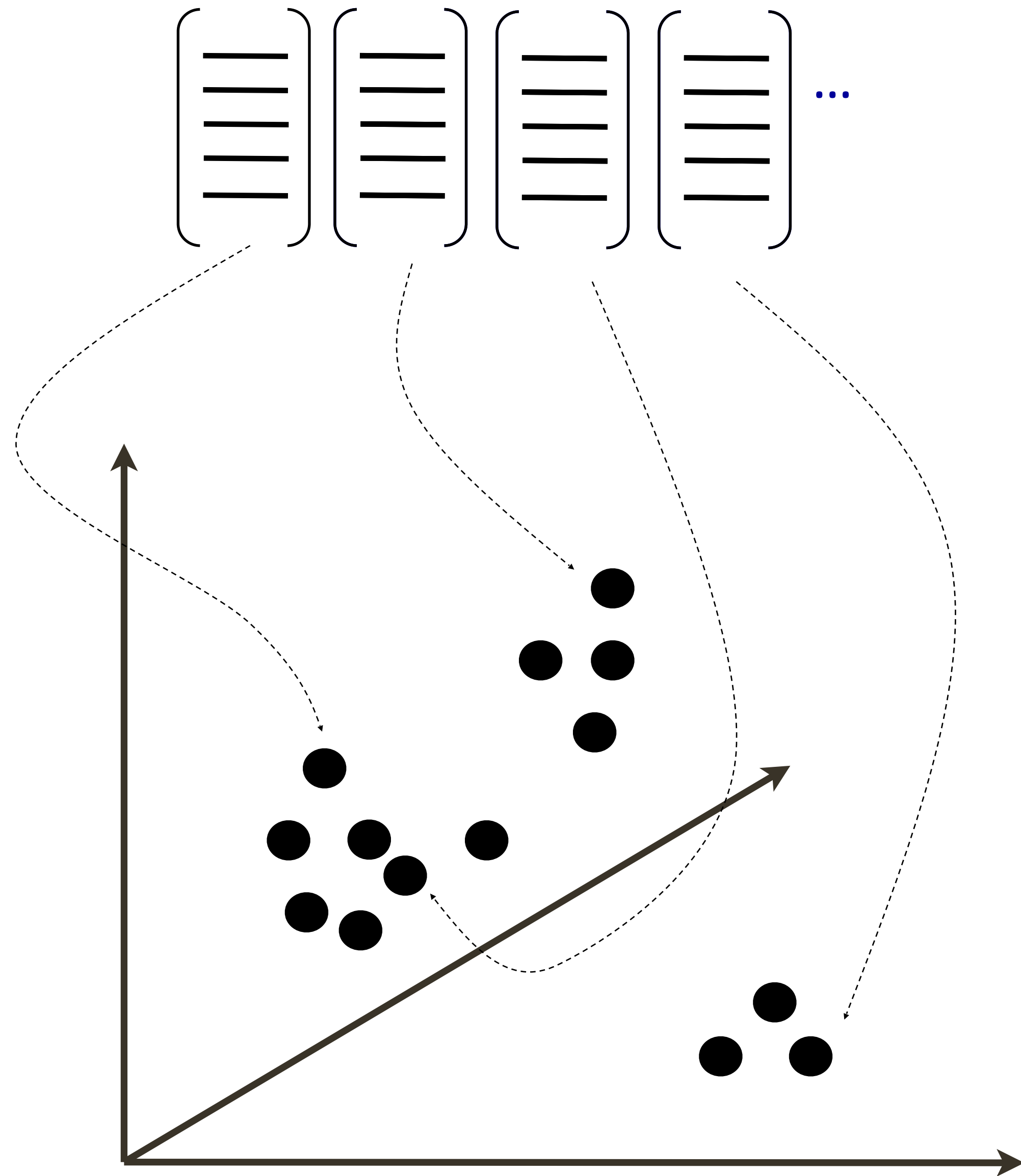
[Mata, Chum, Urban & Pajdla, '02]

[Sivic & Zisserman, '03]

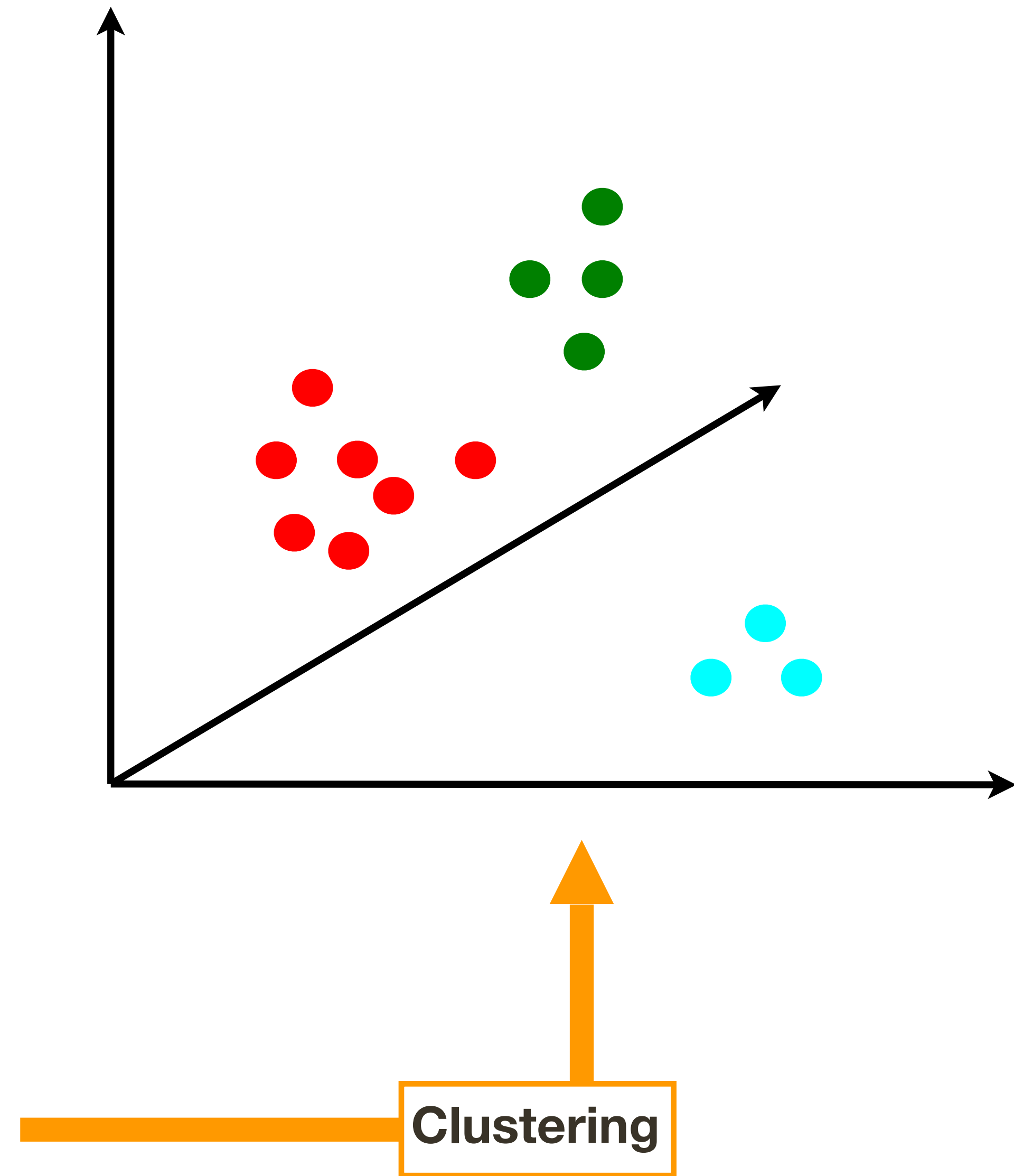
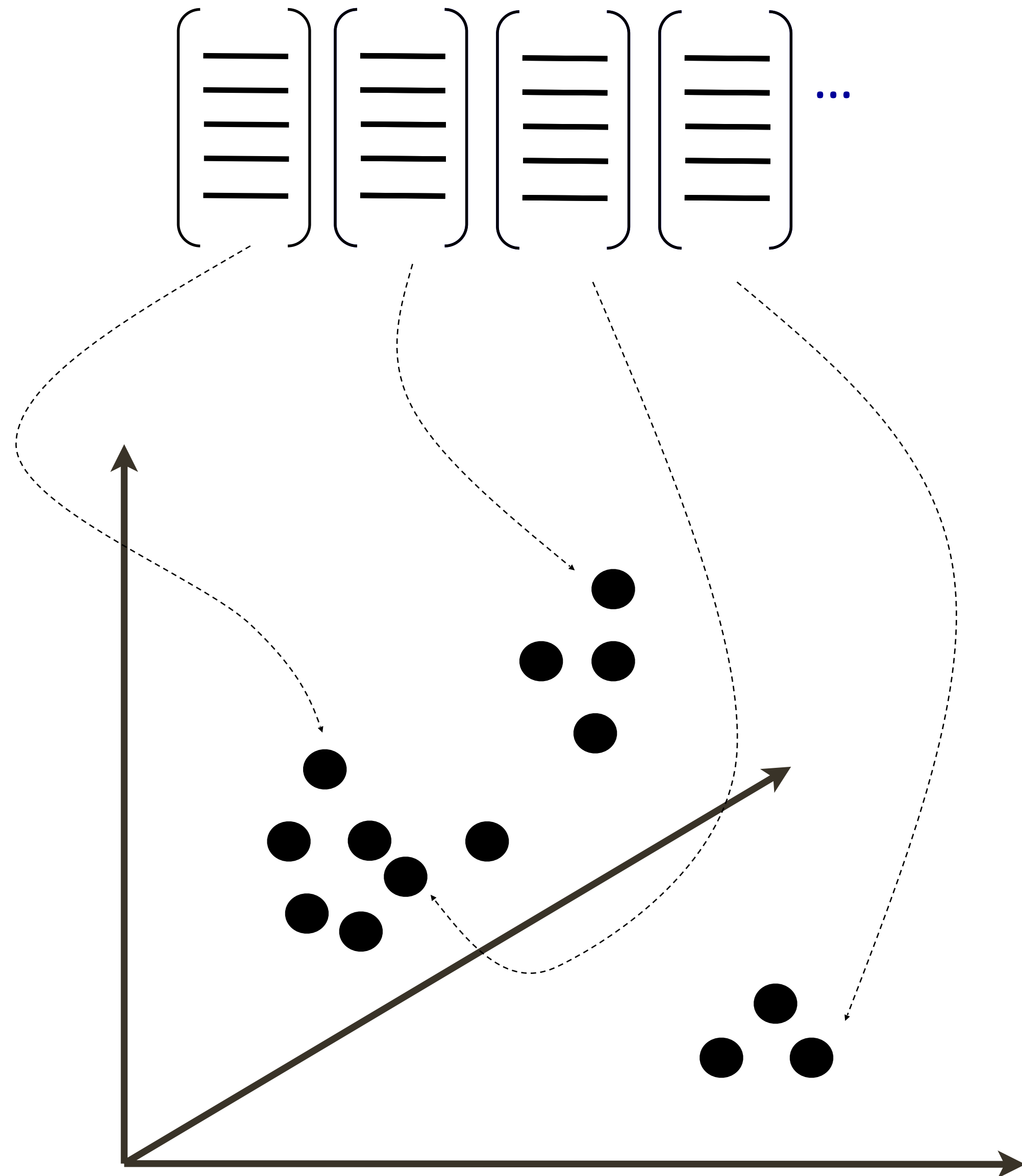
Extracting **SIFT** Patches



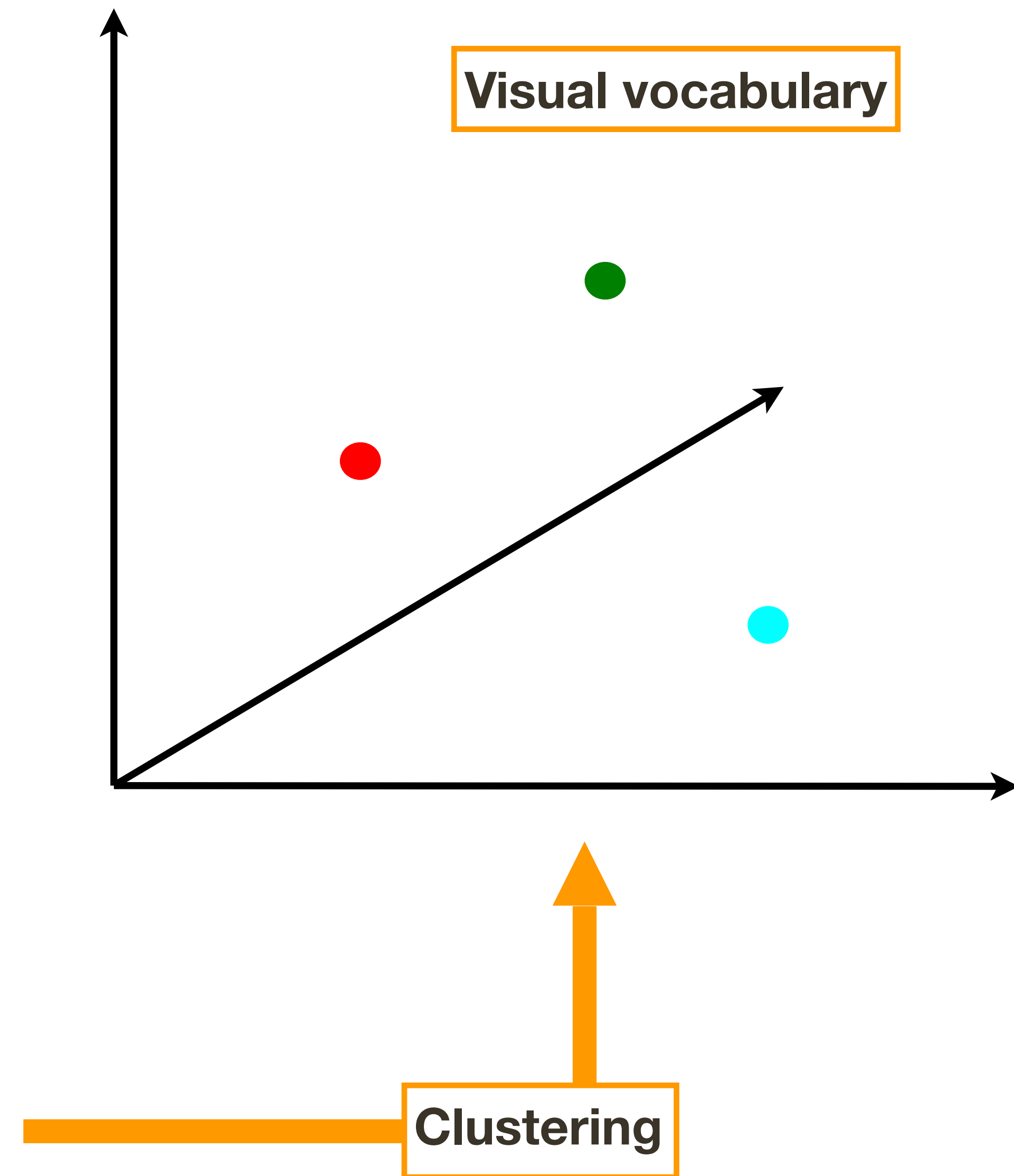
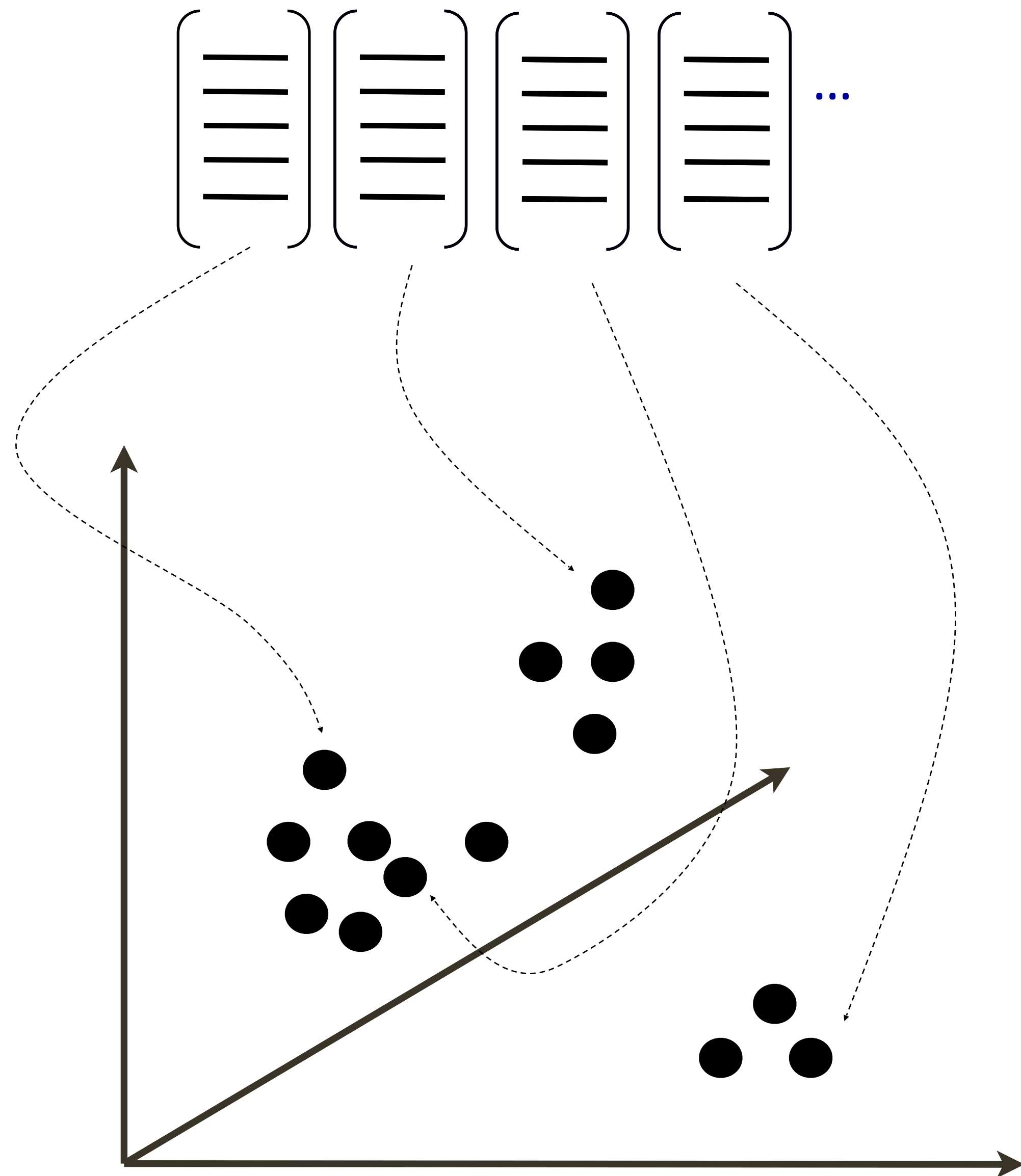
Creating **Dictionary**



Creating **Dictionary**



Creating Dictionary



K-means clustering

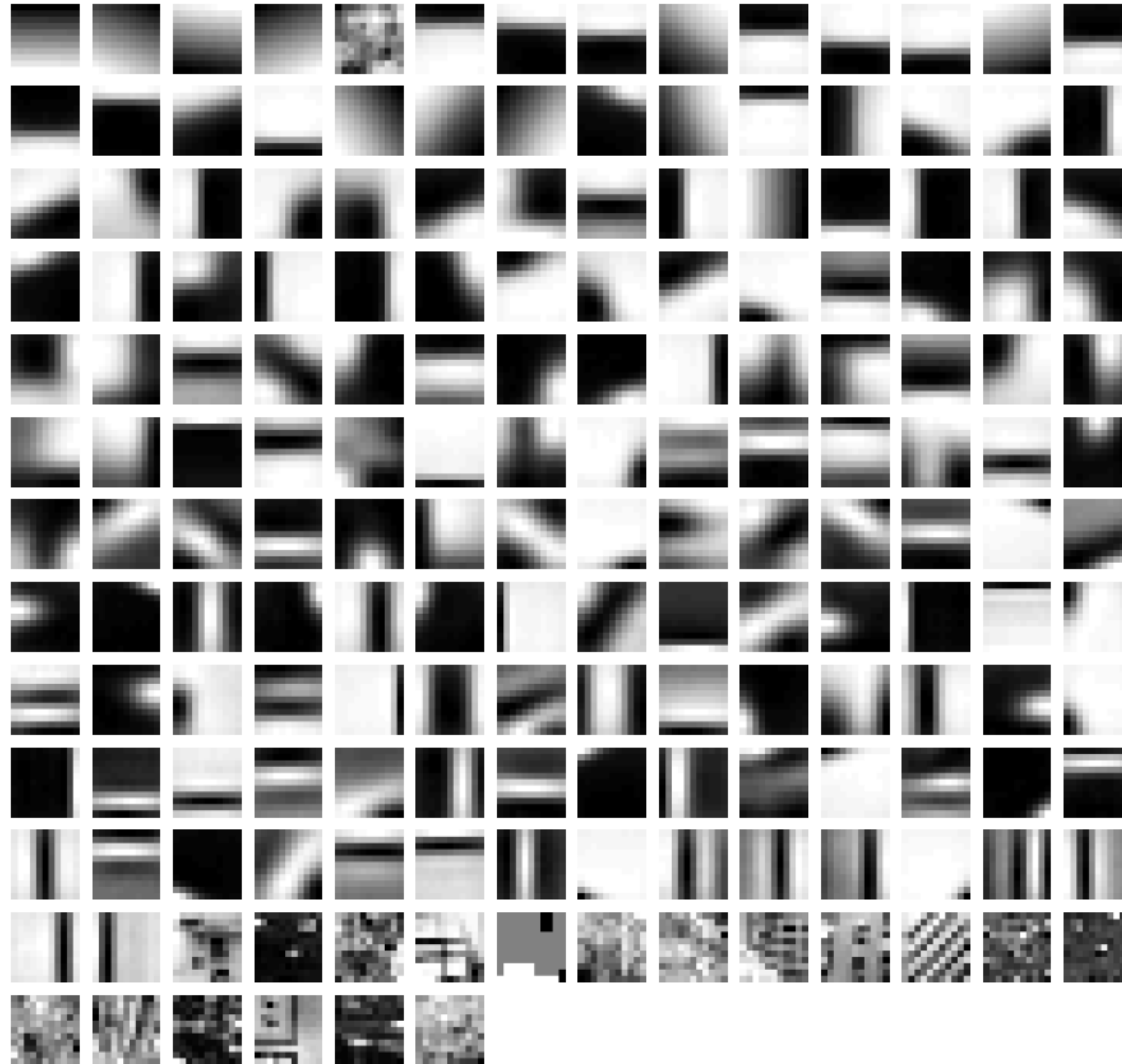
Lecture 26: Re-cap

K-means is a clustering technique that iterates between

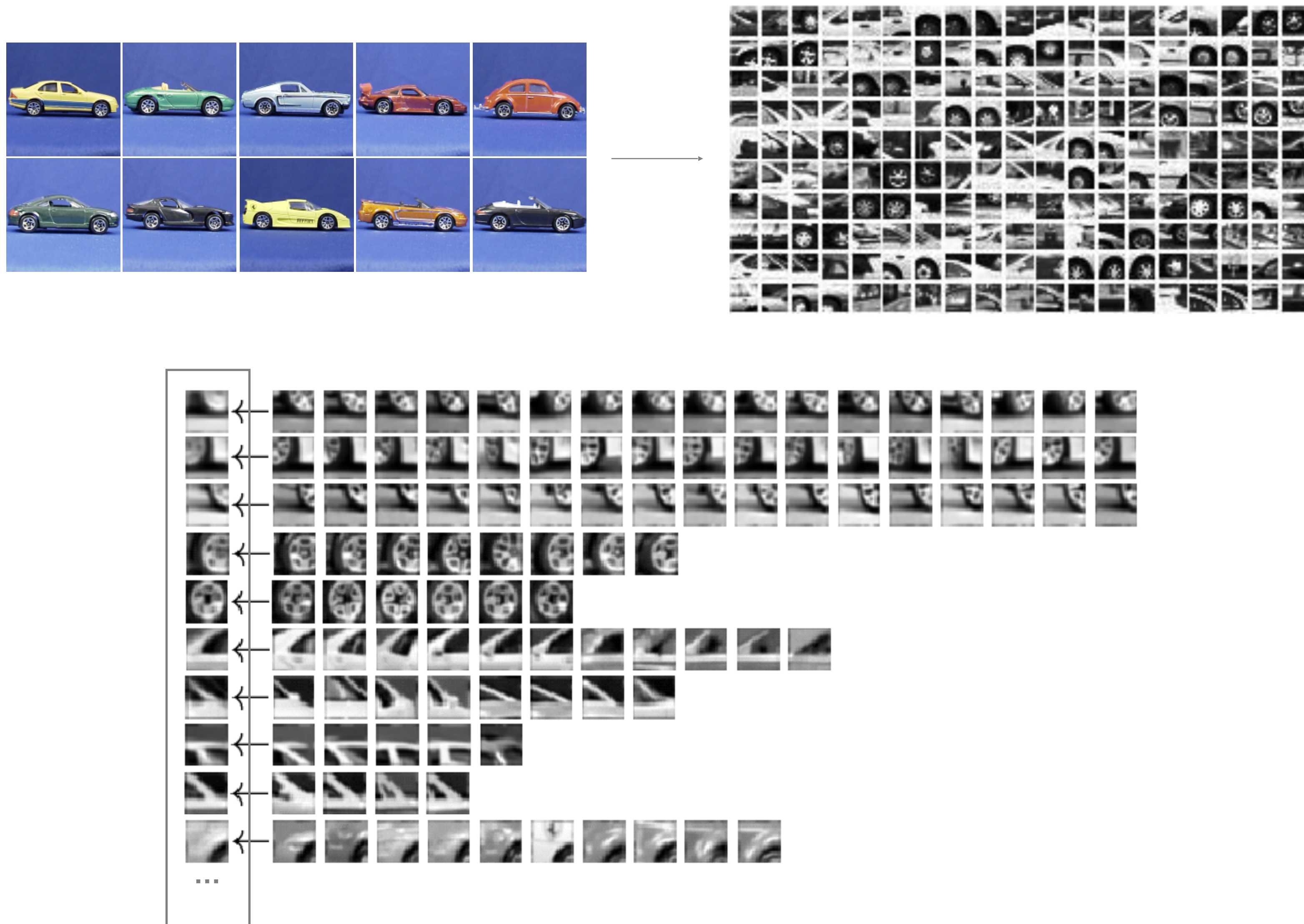
1. Assume the cluster centers are known. Assign each point to the closest cluster center.
2. Assume the assignment of points to clusters is known. Compute the best cluster center for each cluster (as the mean).

K-means clustering is initialization dependent and converges to a local minimum

Example **Visual Dictionary**

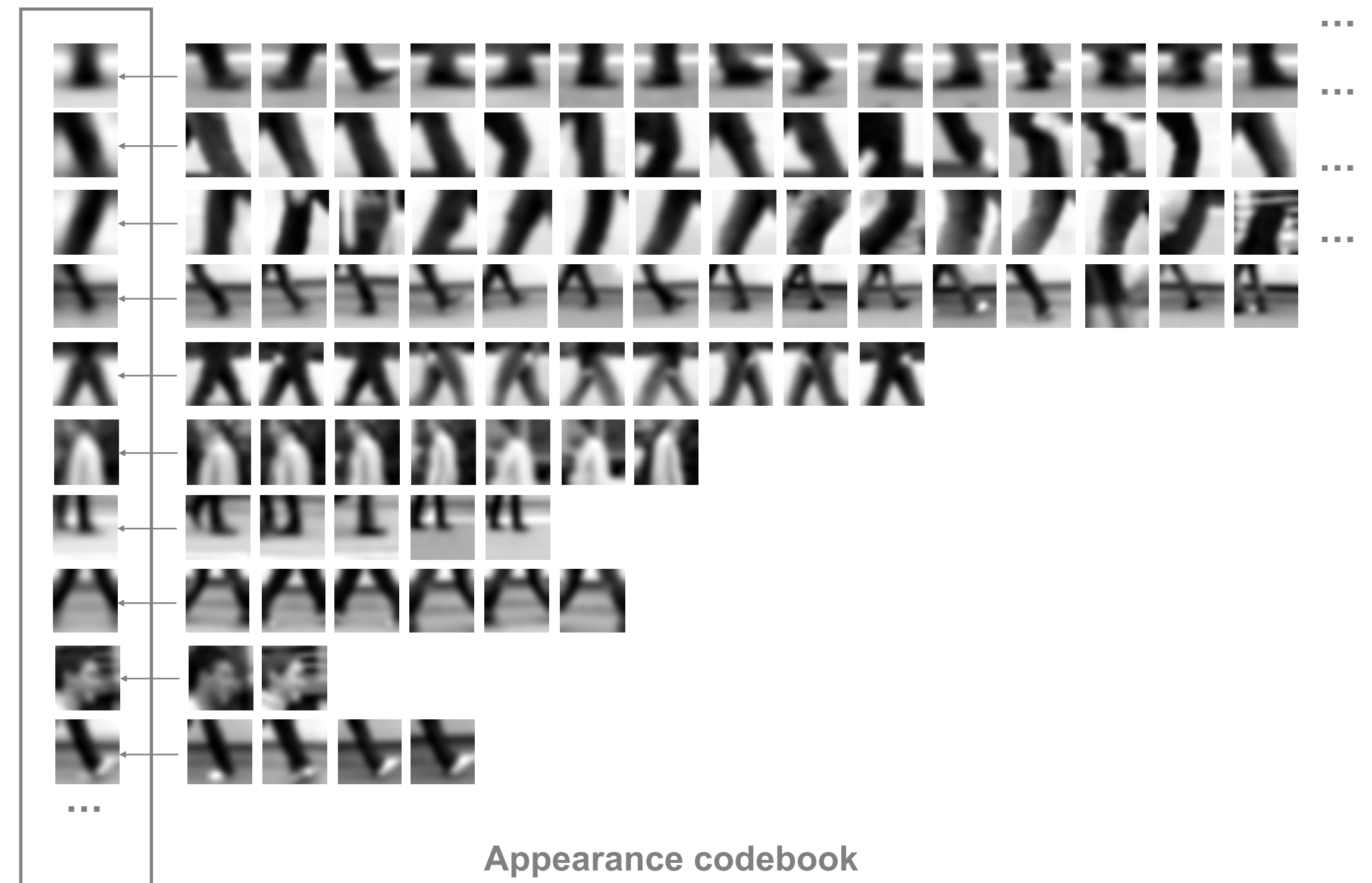
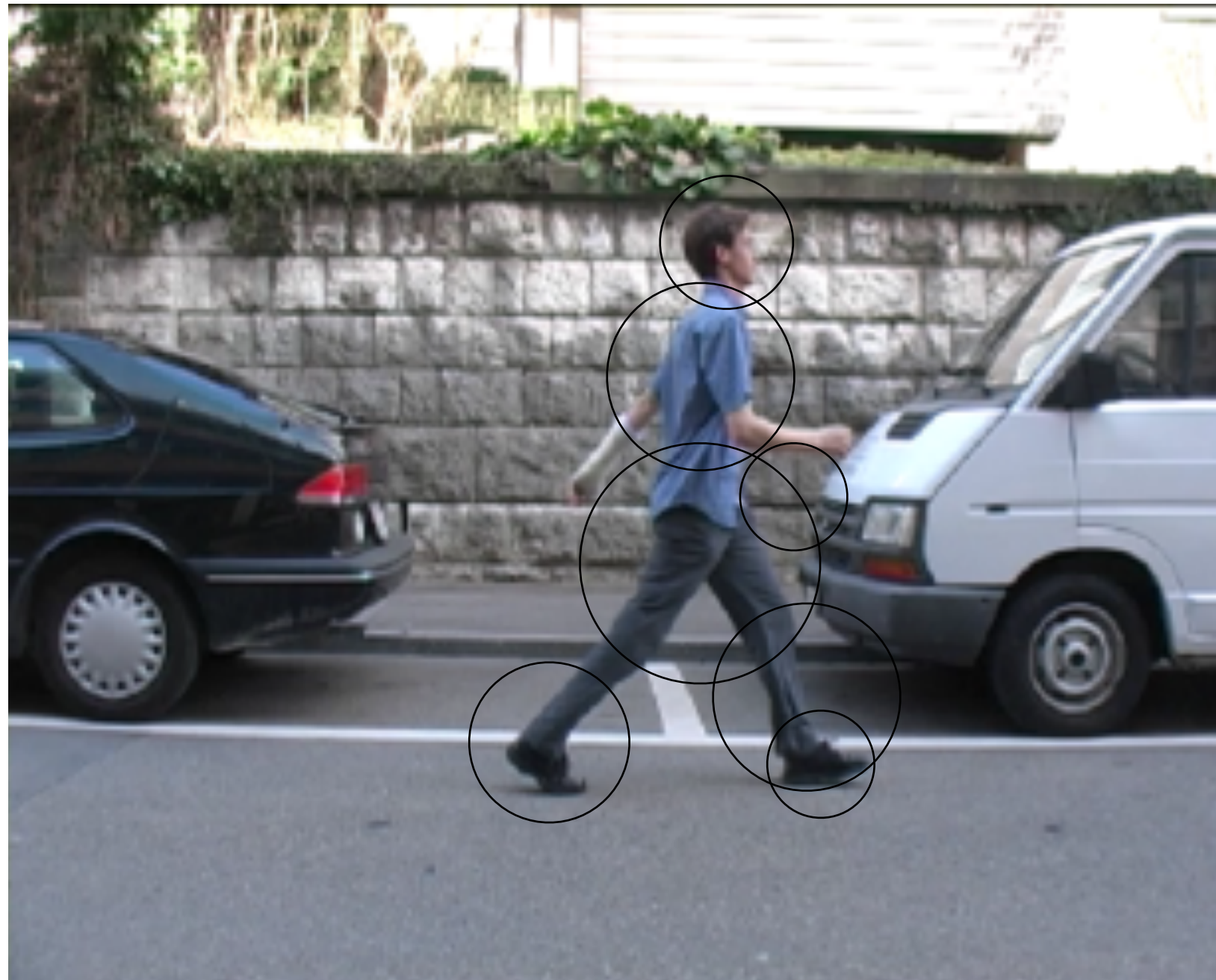


Example **Visual Dictionary**



Source: B. Leibe

Example **Visual Dictionary**



Source: B. Leibe

Standard **Bag-of-Words** Pipeline (for image classification)

Dictionary Learning:

Learn Visual Words using clustering

Encode:

build Bags-of-Words (BOW) vectors
for each image

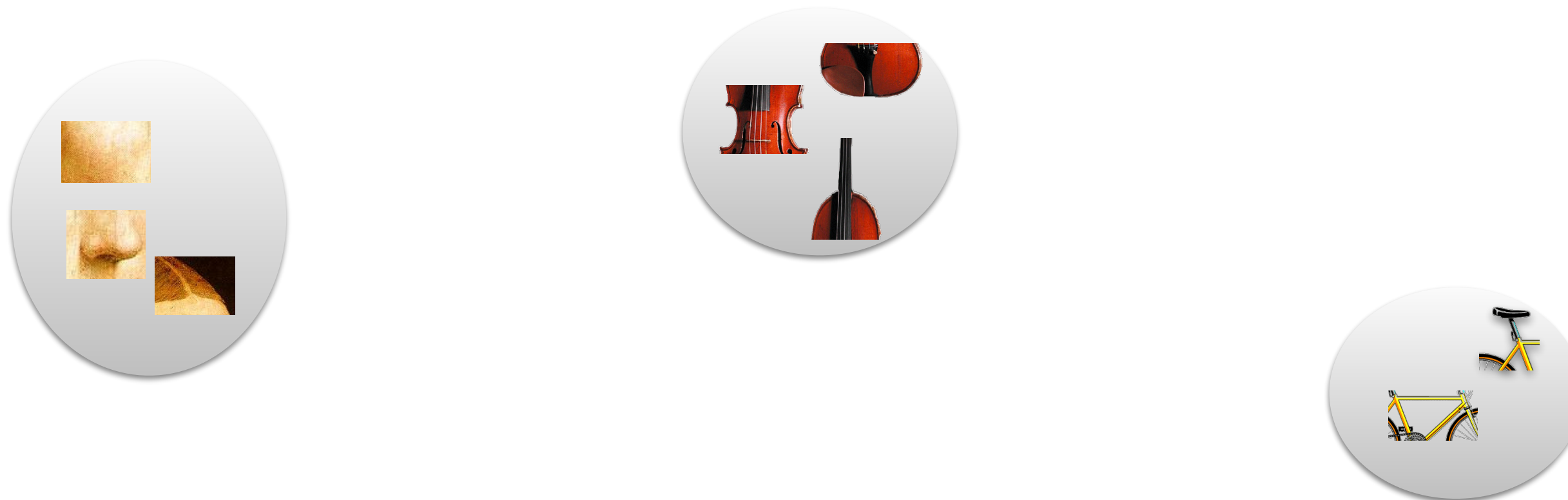
Classify:

Train and test data using BOWs

2. **Encode:** build Bag-of-Words (BOW) vectors for each image

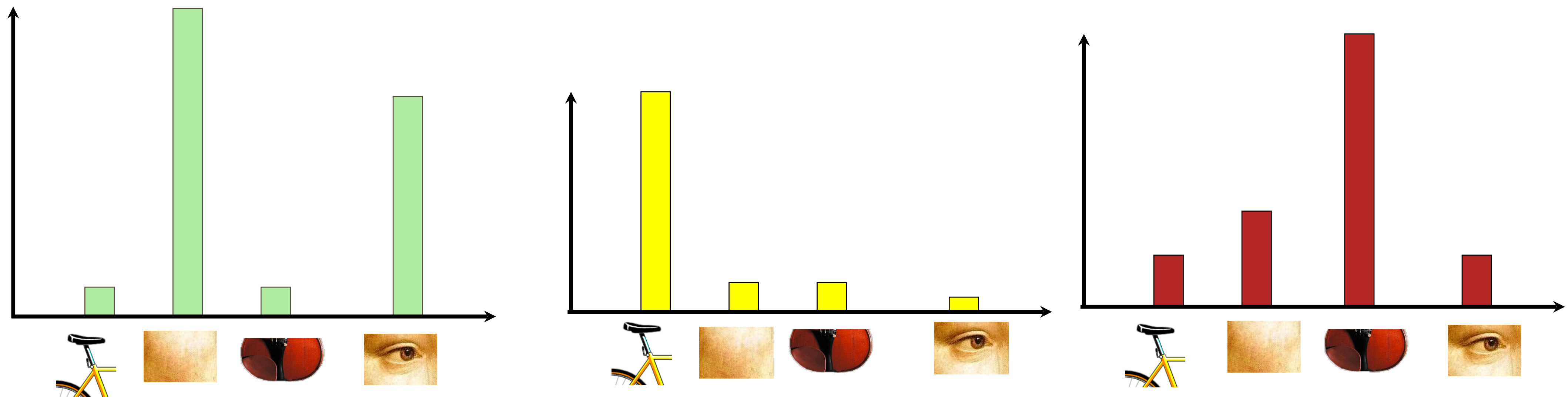


1. **Quantization:** image features gets associated to a visual word (nearest cluster center)

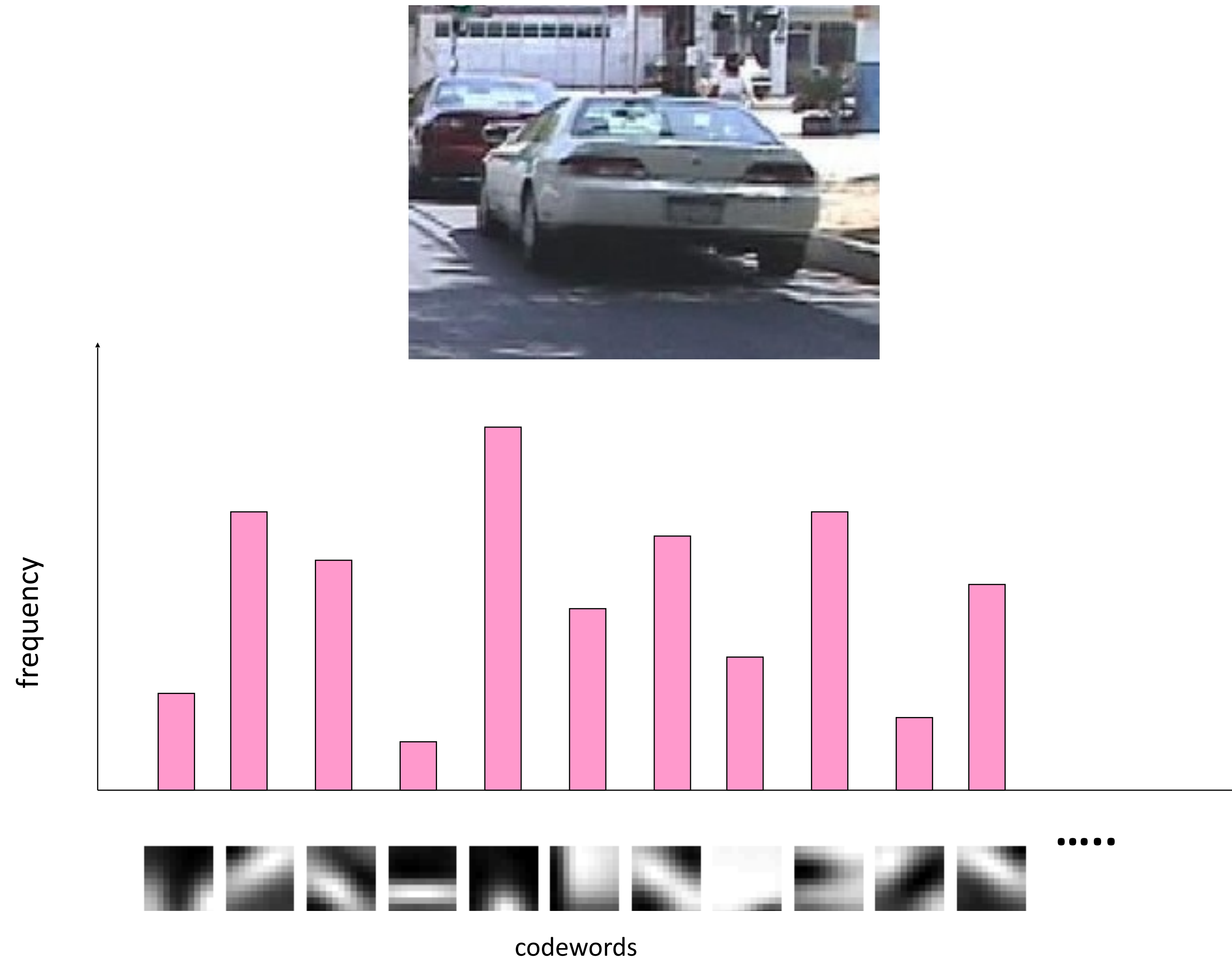


2. **Encode:** build Bag-of-Words (BOW) vectors for each image

2. **Histogram:** count the number of visual word occurrences



2. Encode: build Bag-of-Words (BOW) vectors for each image



Standard **Bag-of-Words** Pipeline (for image classification)

Dictionary Learning:

Learn Visual Words using clustering

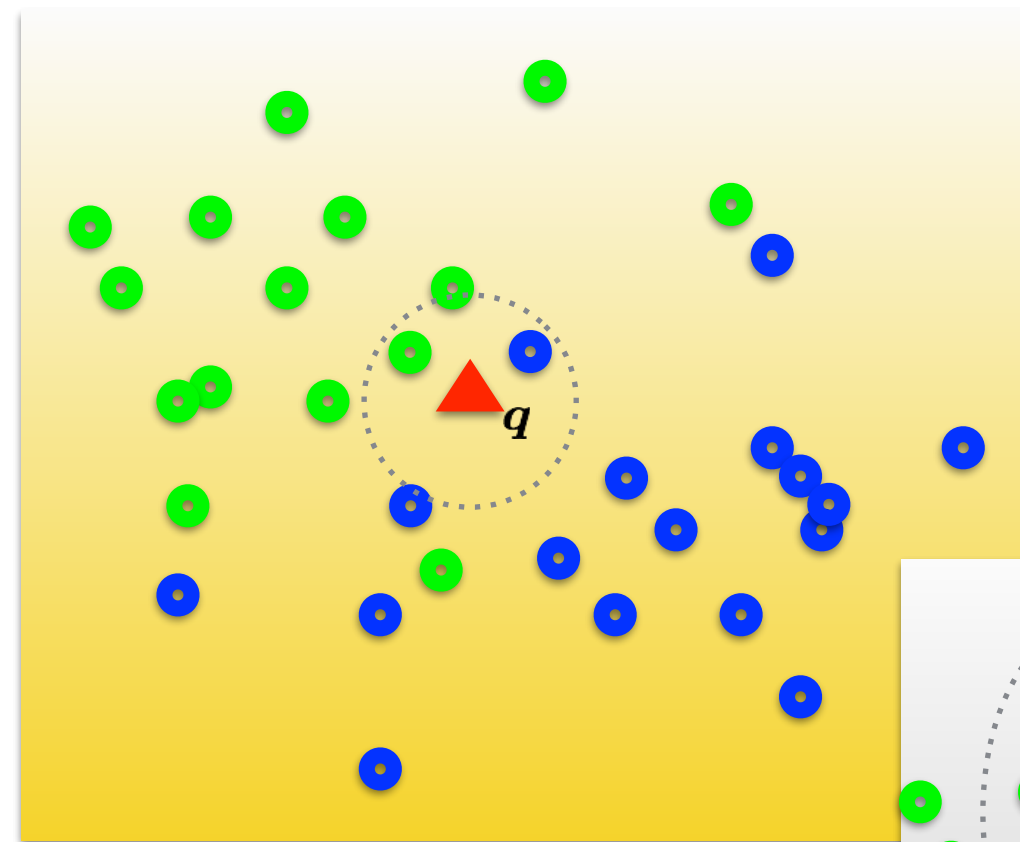
Encode:

build Bags-of-Words (BOW) vectors
for each image

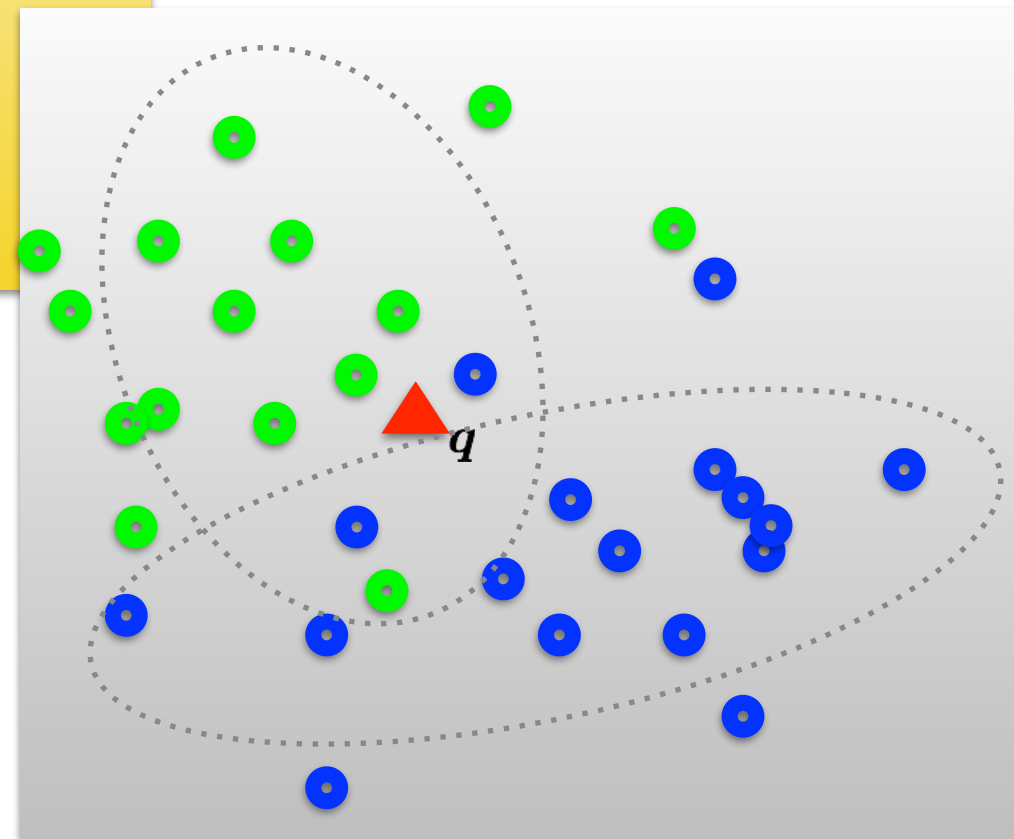
Classify:

Train and test data using BOWs

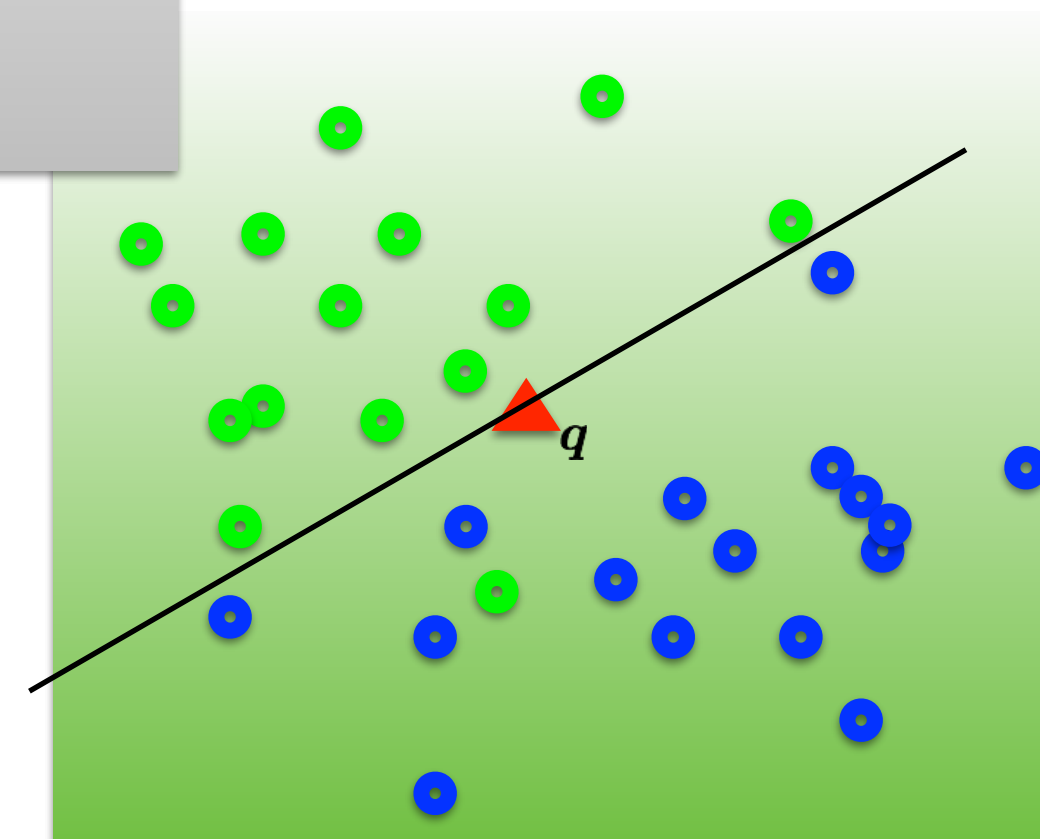
3. Classify: Train and text classifier using BOWs



K nearest neighbors



Naïve Bayes



Support Vector Machine