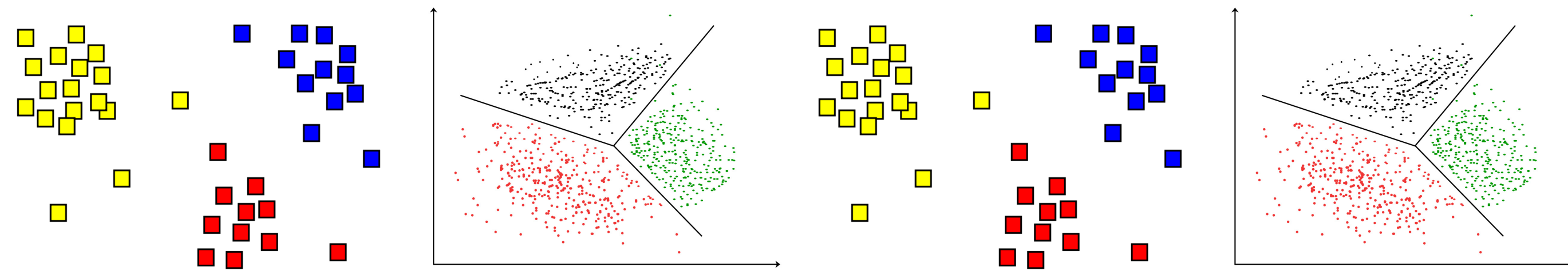


# CPSC 425: Computer Vision



## Lecture 19: Image Classification (cont.)

# Menu for Today (March 19, 2020)

## 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:

- No more **fun** examples (sorry) ... it is difficult to get them to work
- **Assignment 5** is out

# Lecture 18: Re-cap

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

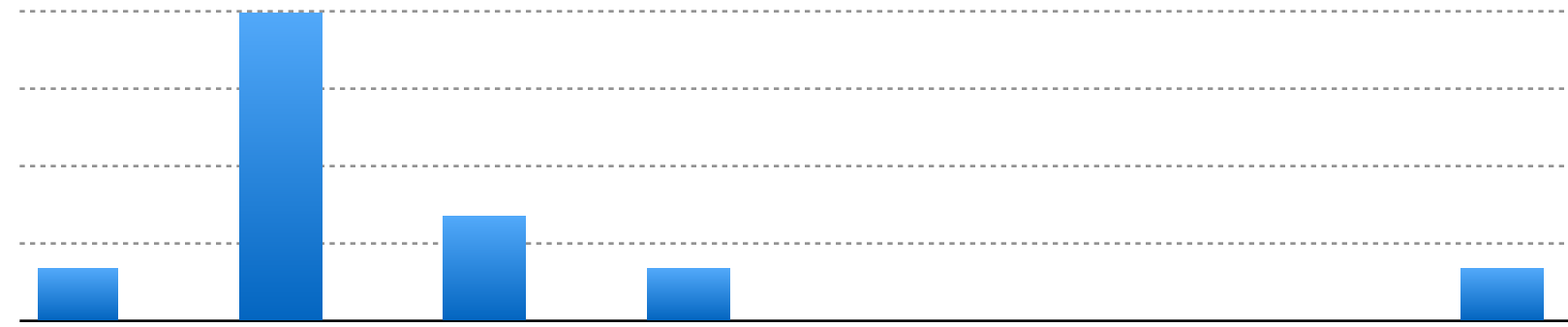
# Lecture 18: Re-cap (Vector Space Model)

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

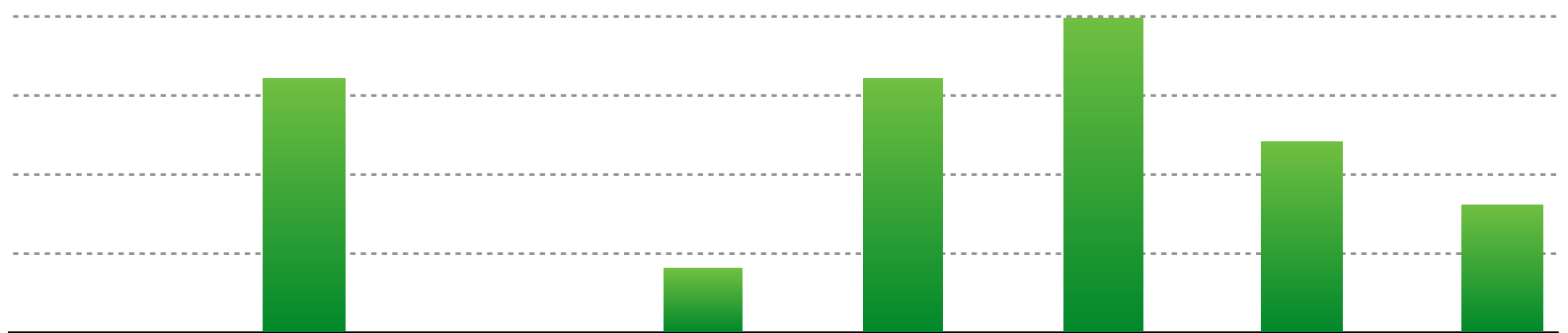
To classify a text document (e.g. as an article on sports, entertainment, business, politics) we might find patterns in the occurrences of certain words.

# Lecture 18: Re-cap (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>

# Lecture 18: Re-cap (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?



# Lecture 18: Re-cap (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

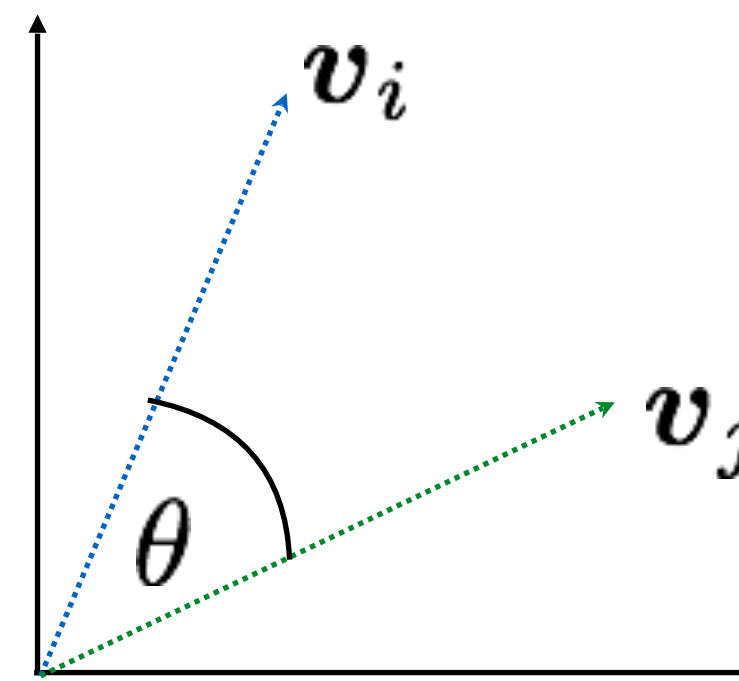
just a histogram over words



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:

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



Slide Credit: Ioannis (Yannis) Gkioulekas (CMU)

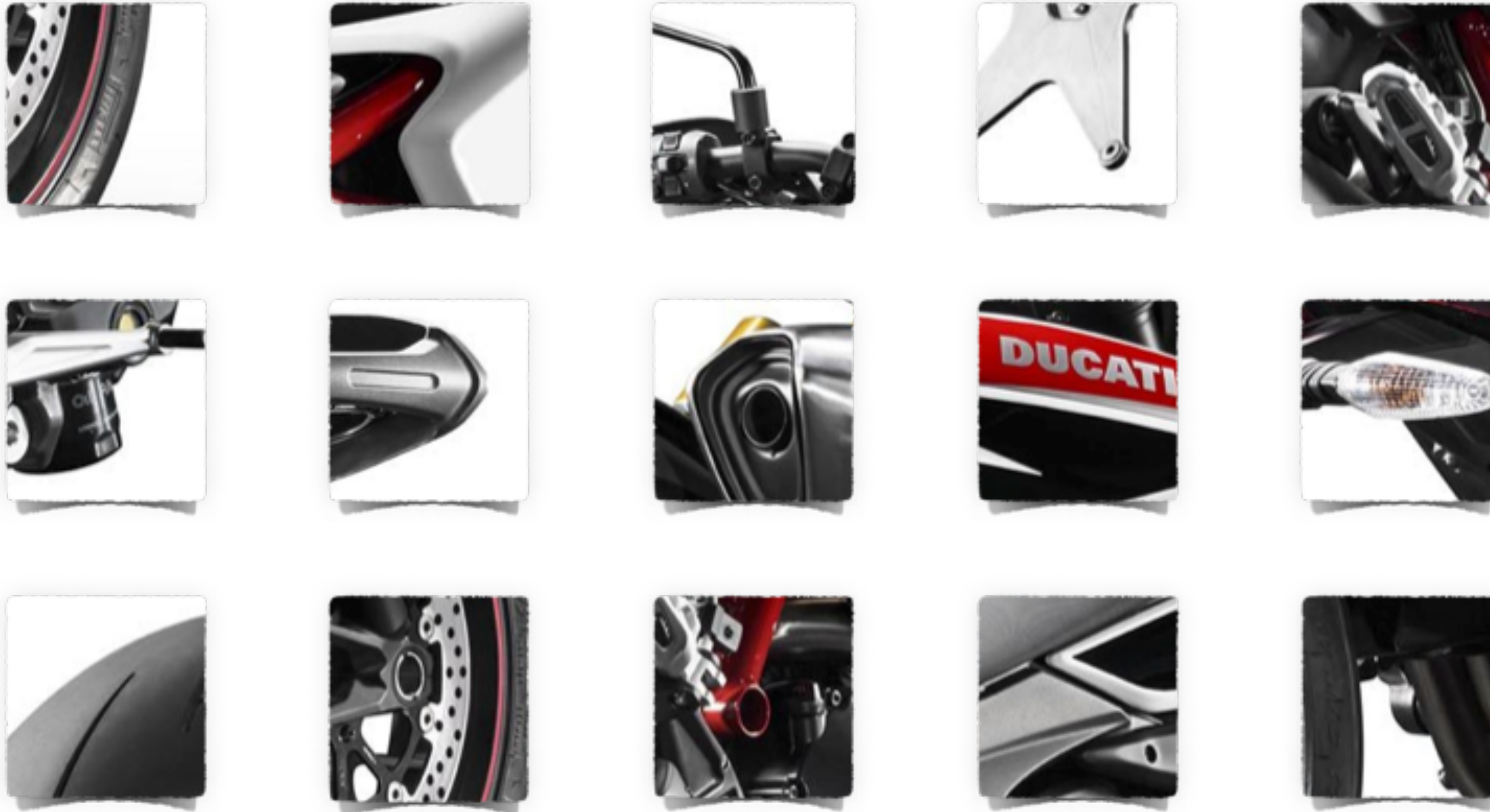
# Visual **Words**

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.

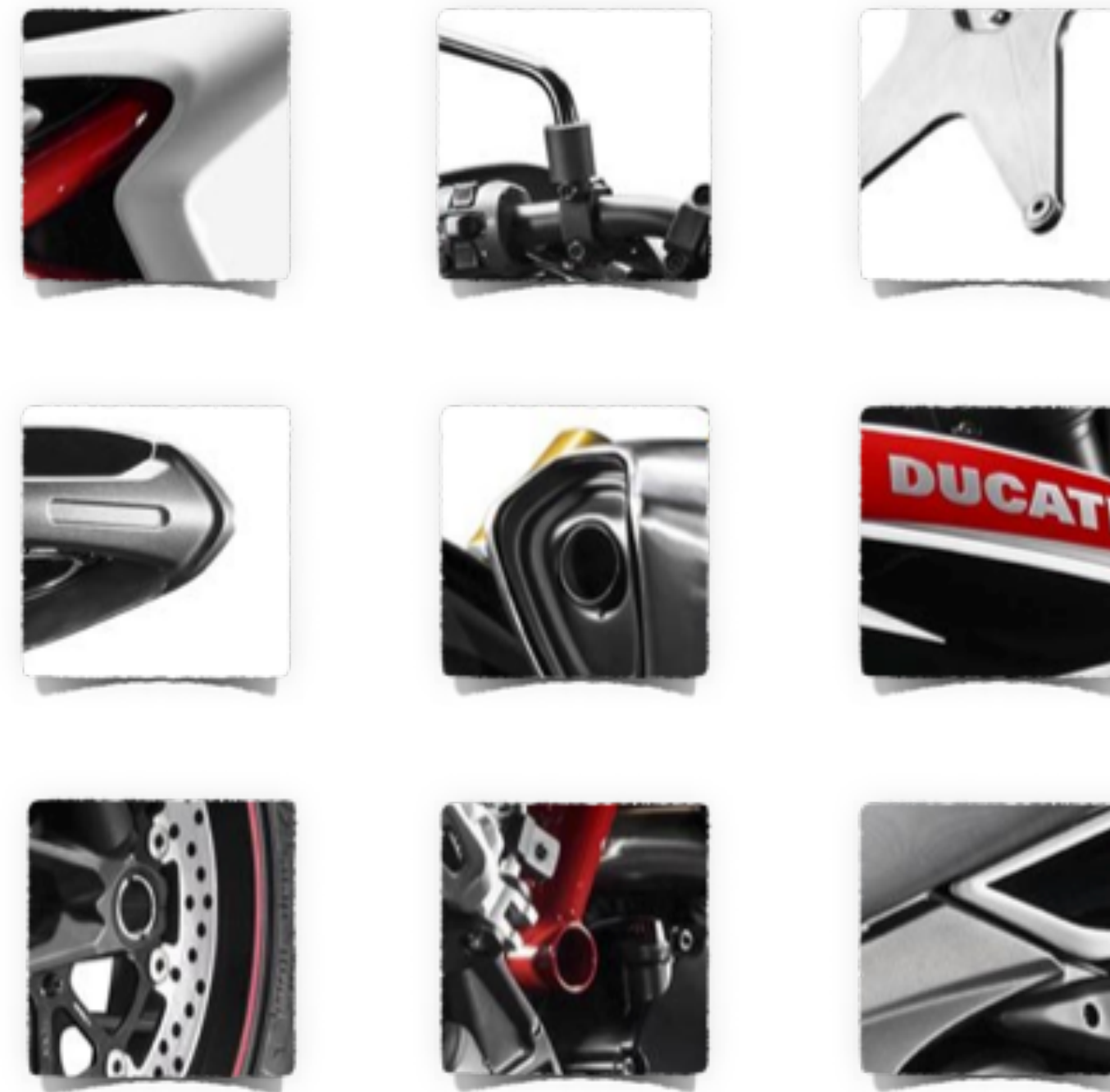


# What **Objects** do These Parts Belong To?



Some local feature are very informative

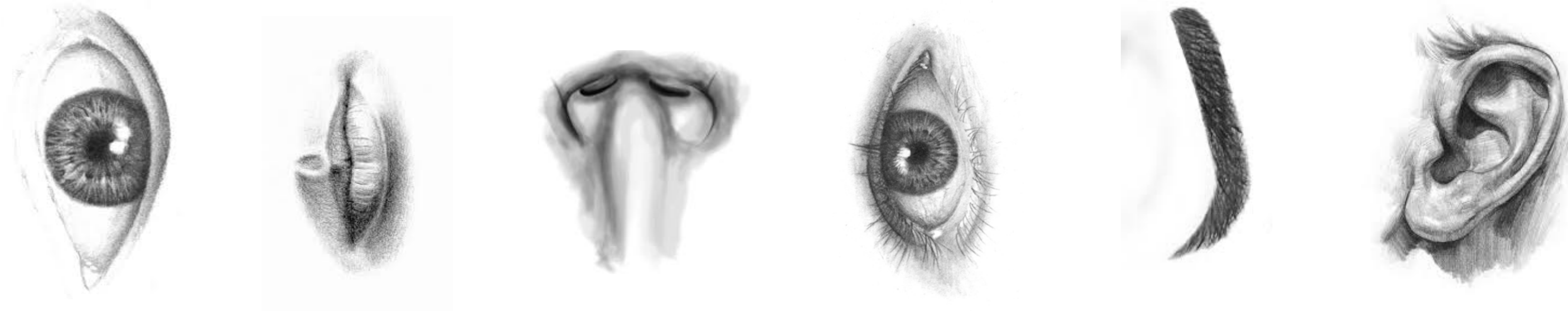
An object as



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

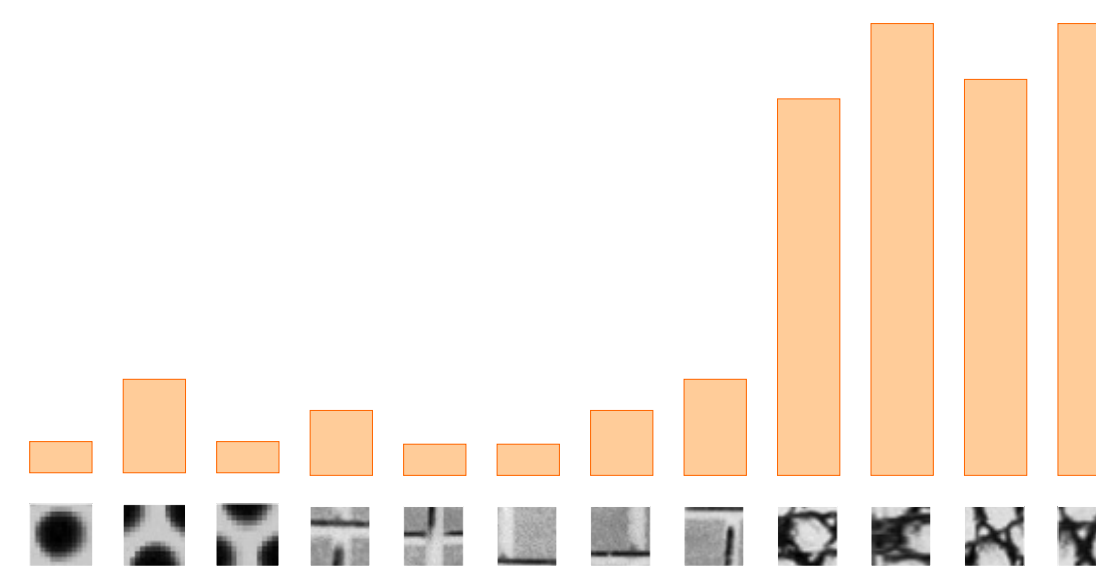
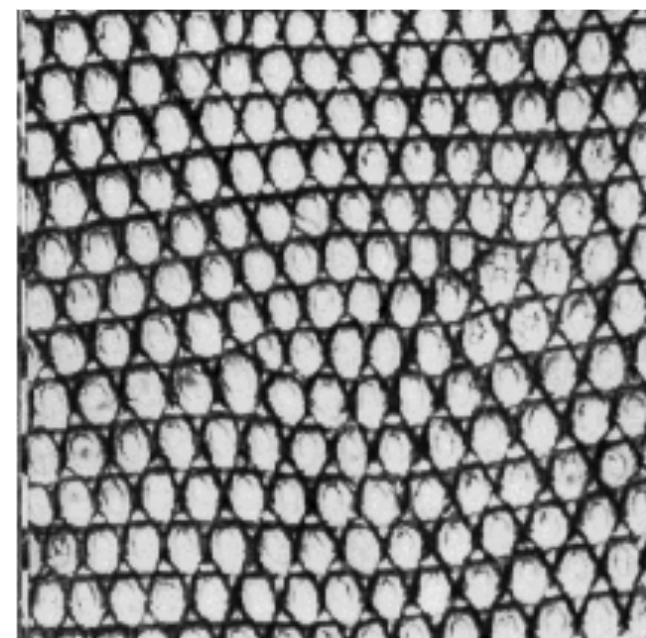
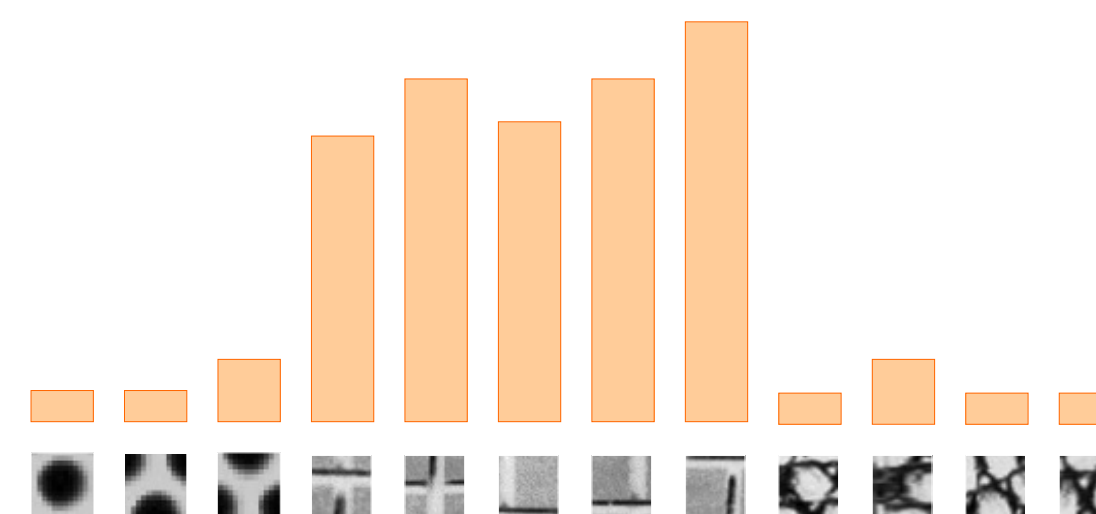
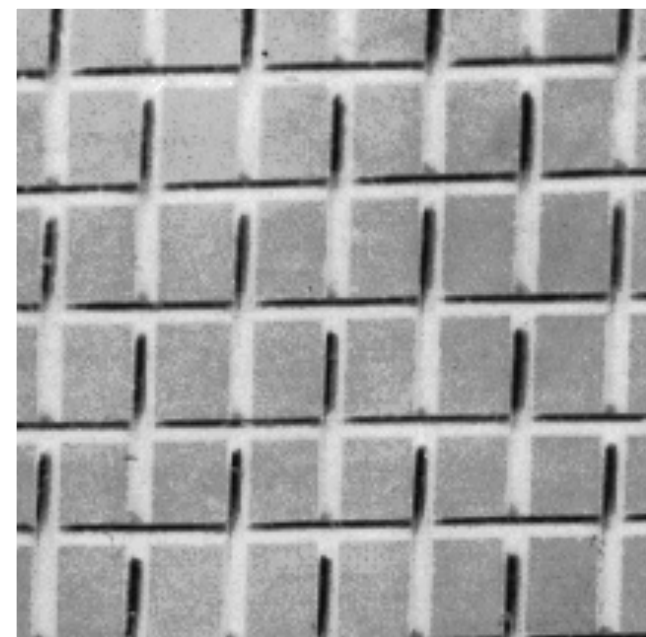
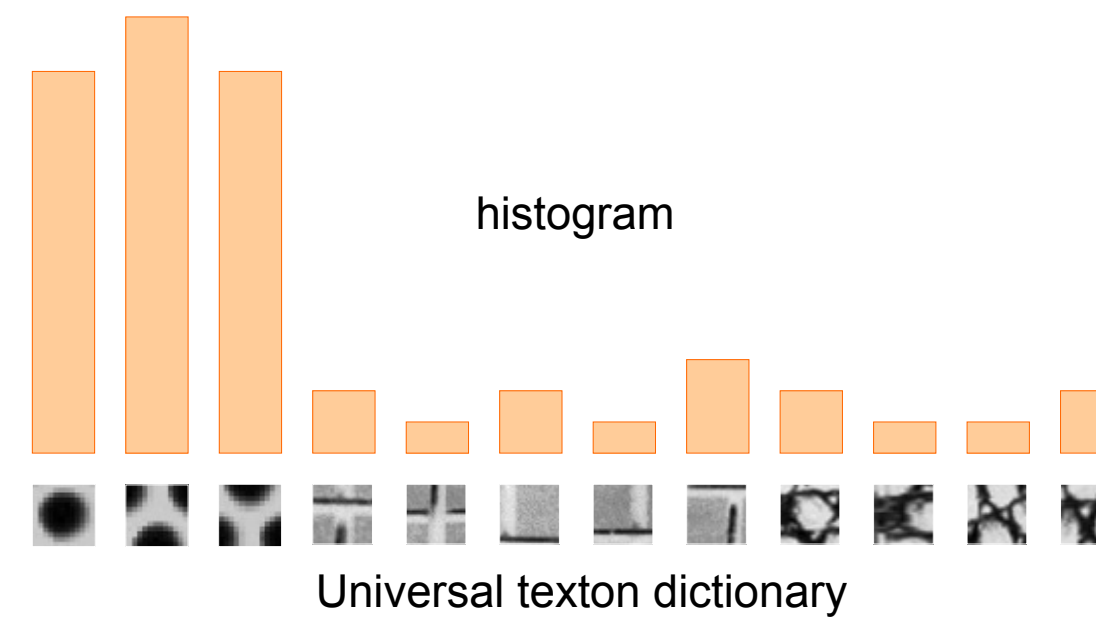
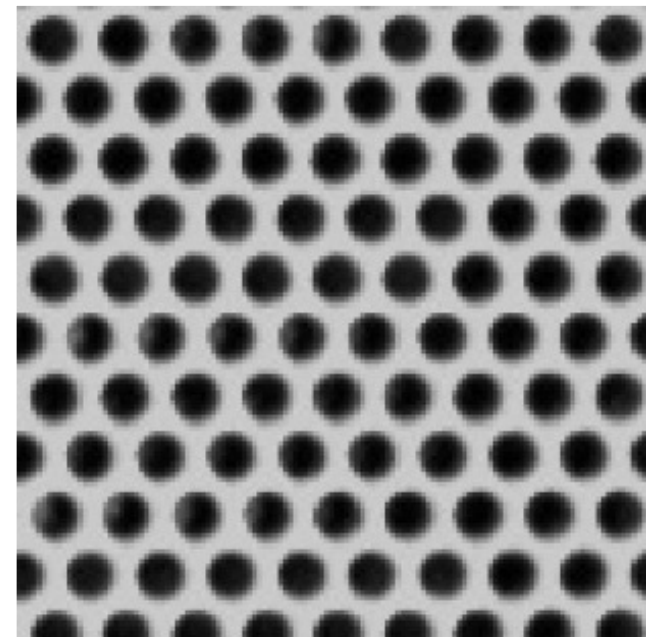
- deals well with occlusion
- scale invariant
- rotation invariant

# (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:**

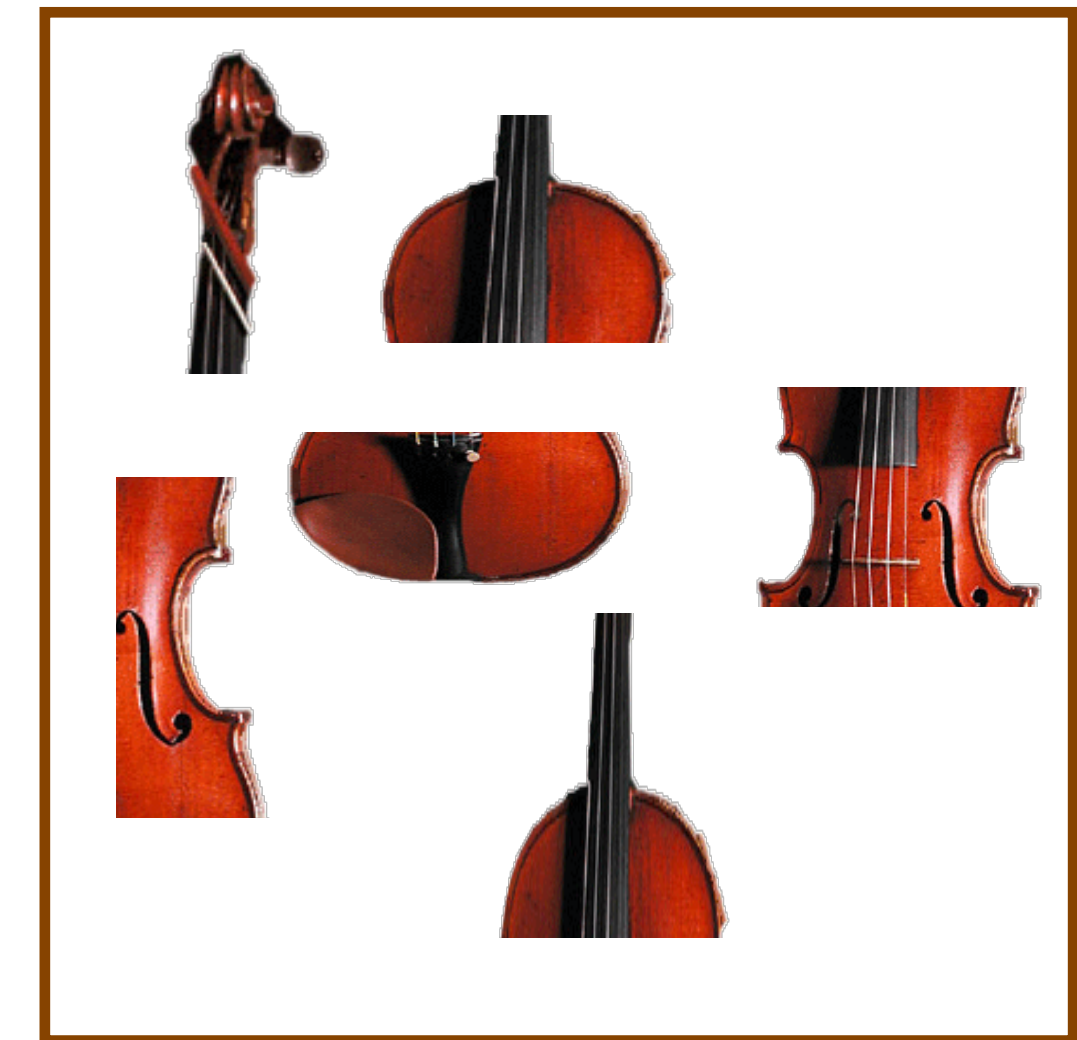
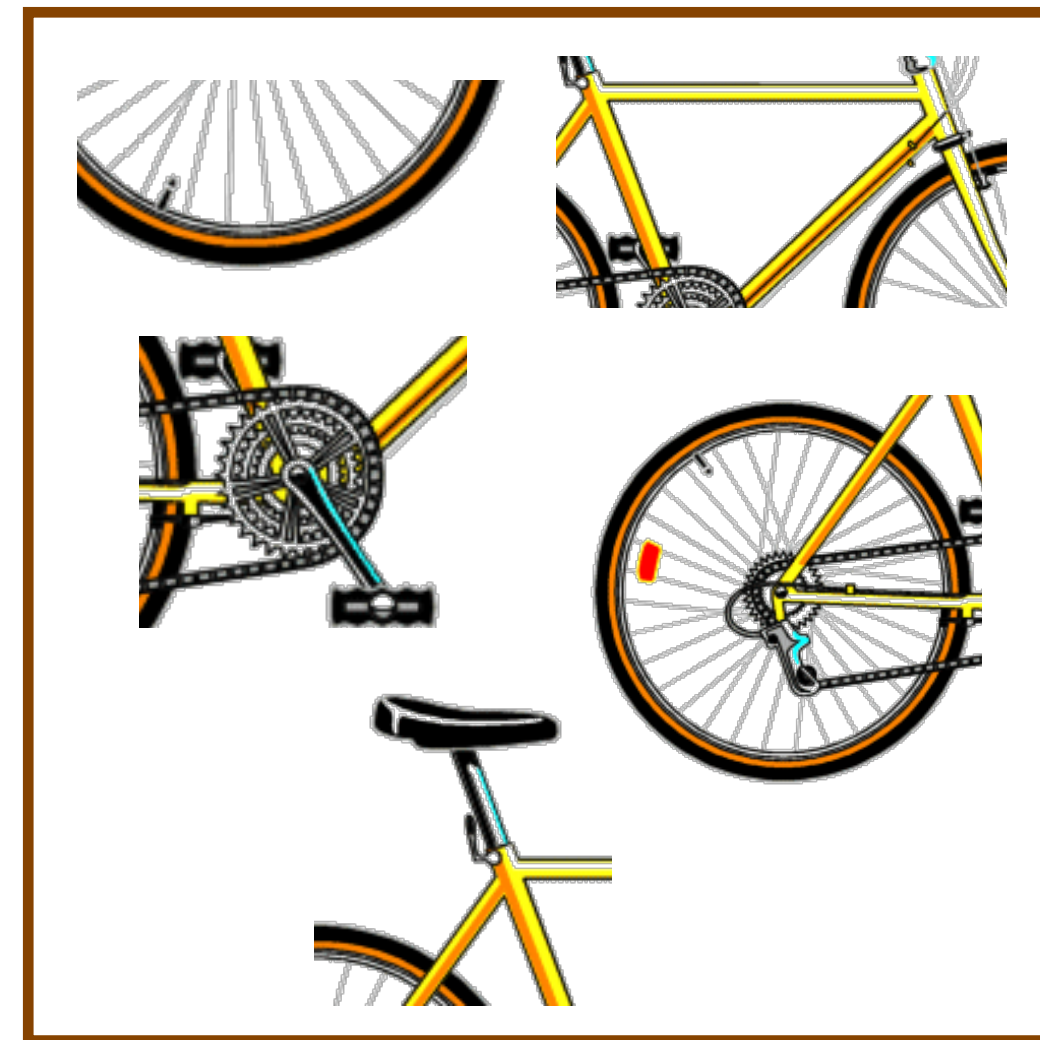
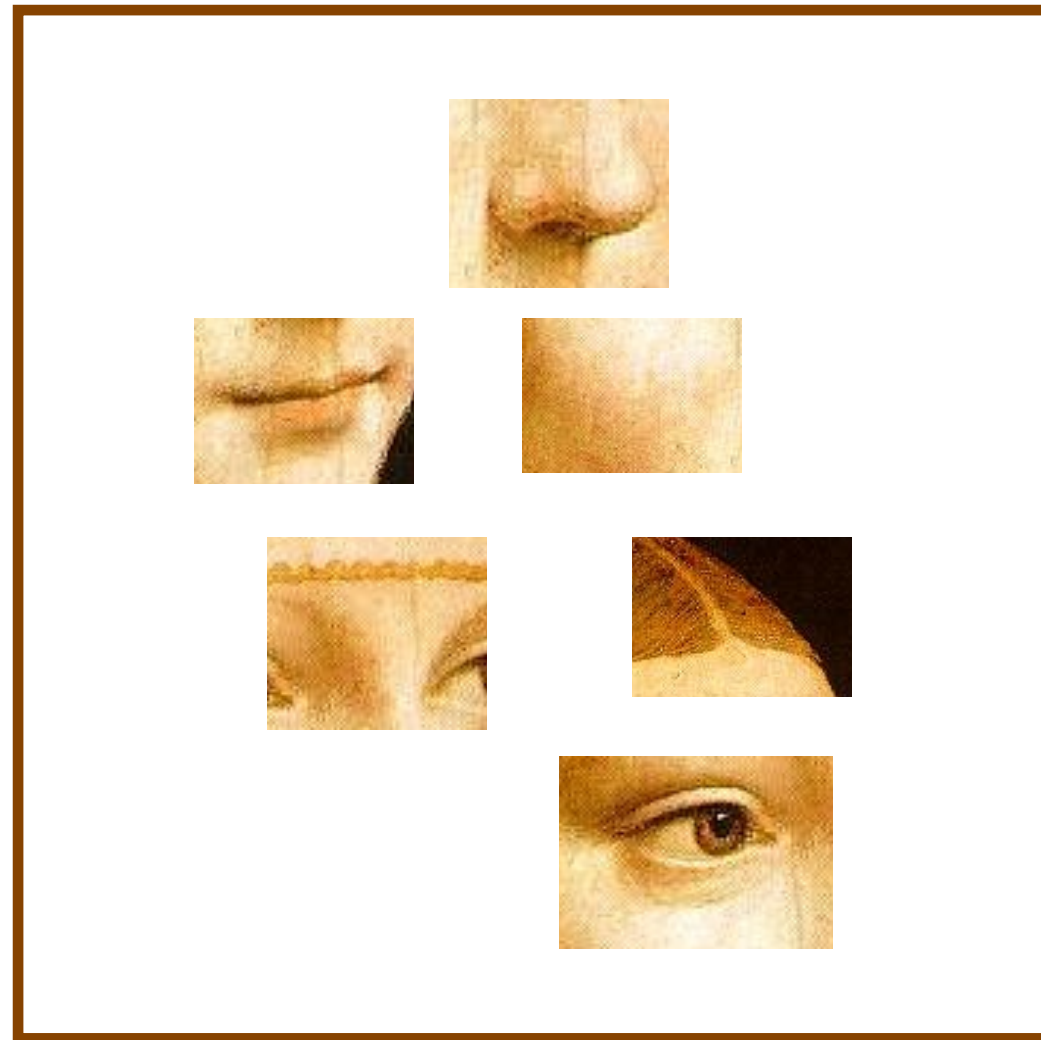
build Bags-of-Words (BOW) vectors  
for each image

## **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

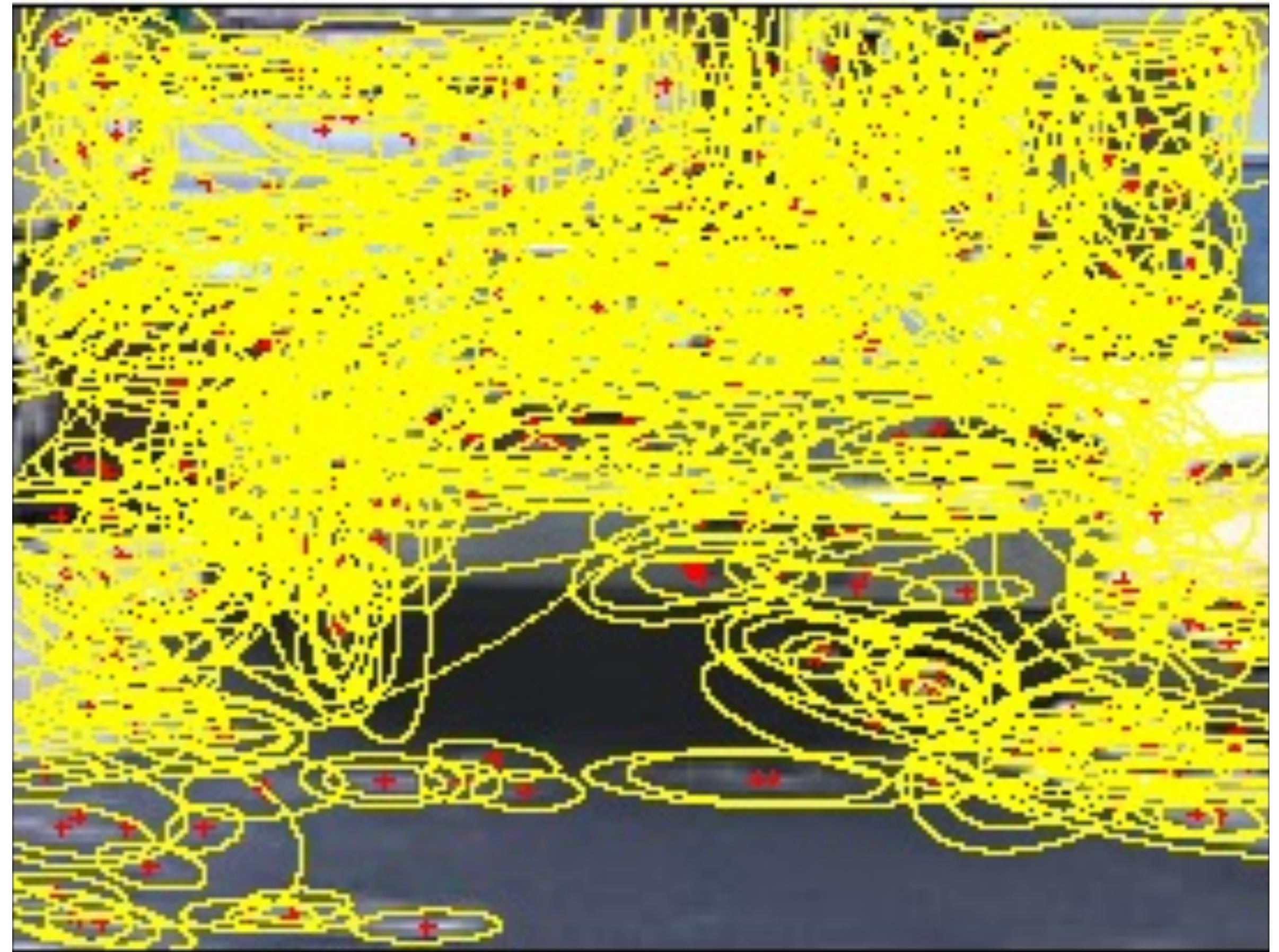
Fei-Fei & Perona, 2005

Sivic et al. 2005

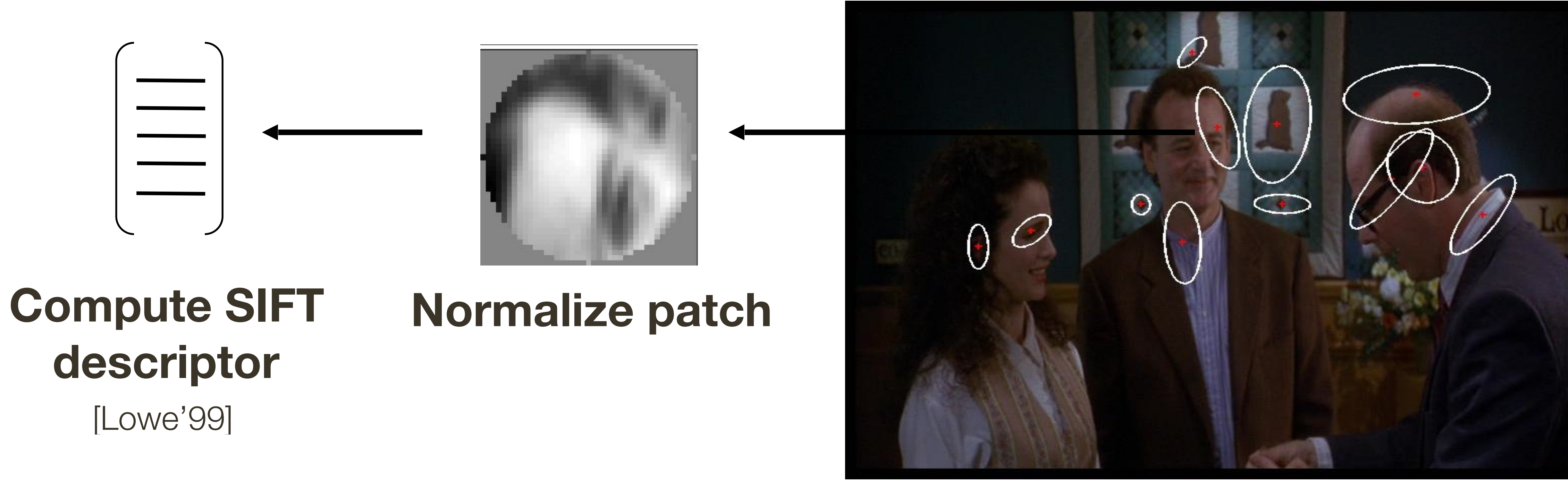
- Other methods

Random sampling (Vidal-Naquet & Ullman, 2002)

Segmentation-based patches (Barnard et al. 2003)



# Extracting **SIFT** Patches



**Compute SIFT descriptor**

[Lowe'99]

**Normalize patch**

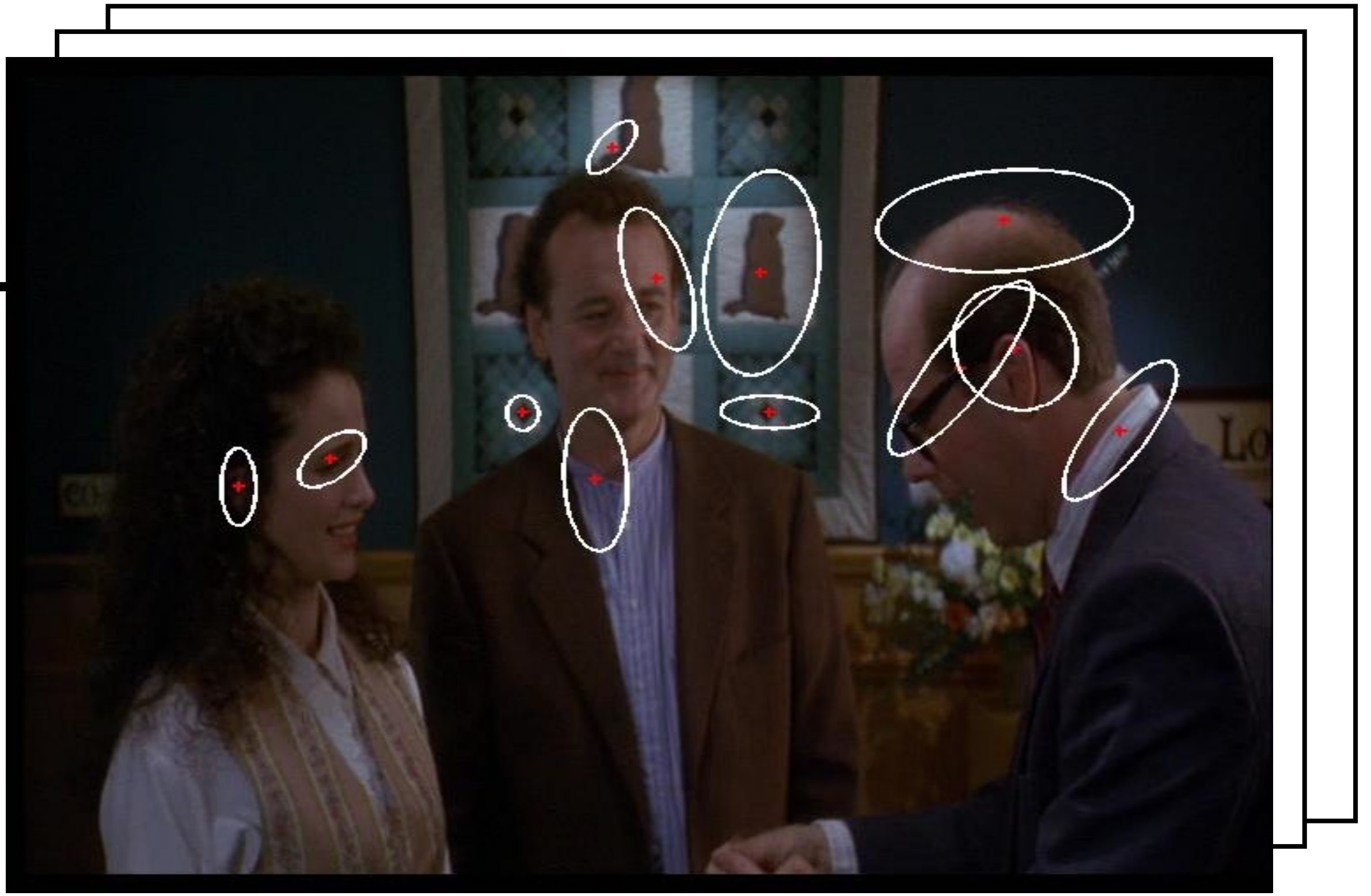
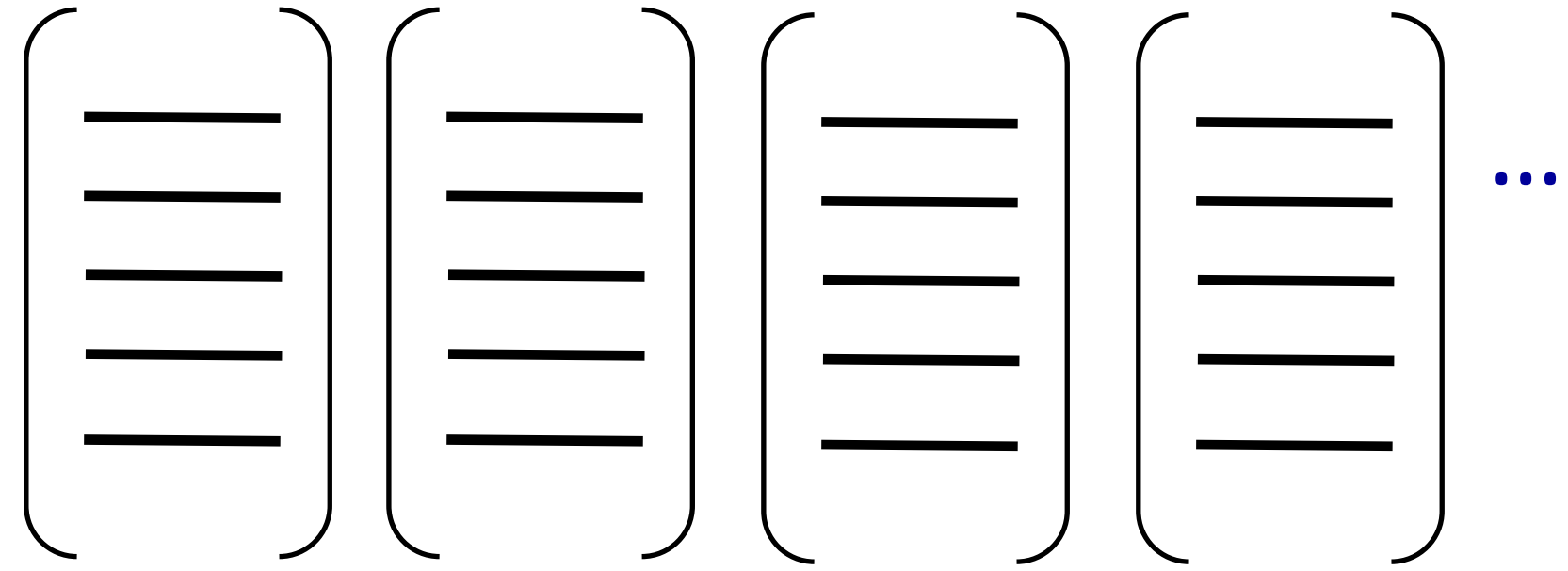
**Detect patches**

[Mikojaczyk and Schmid '02]

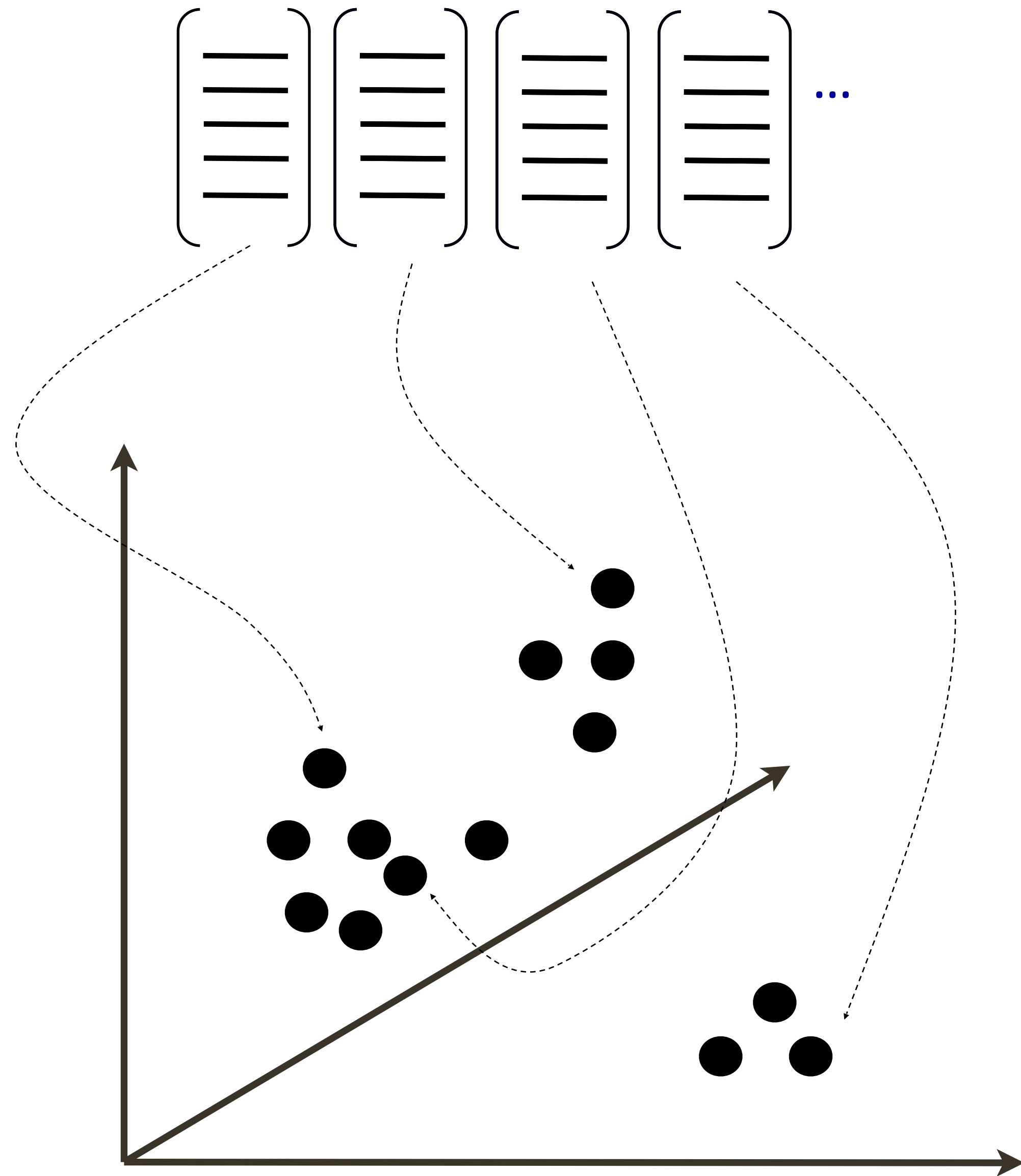
[Mata, Chum, Urban & Pajdla, '02]

[Sivic & Zisserman, '03]

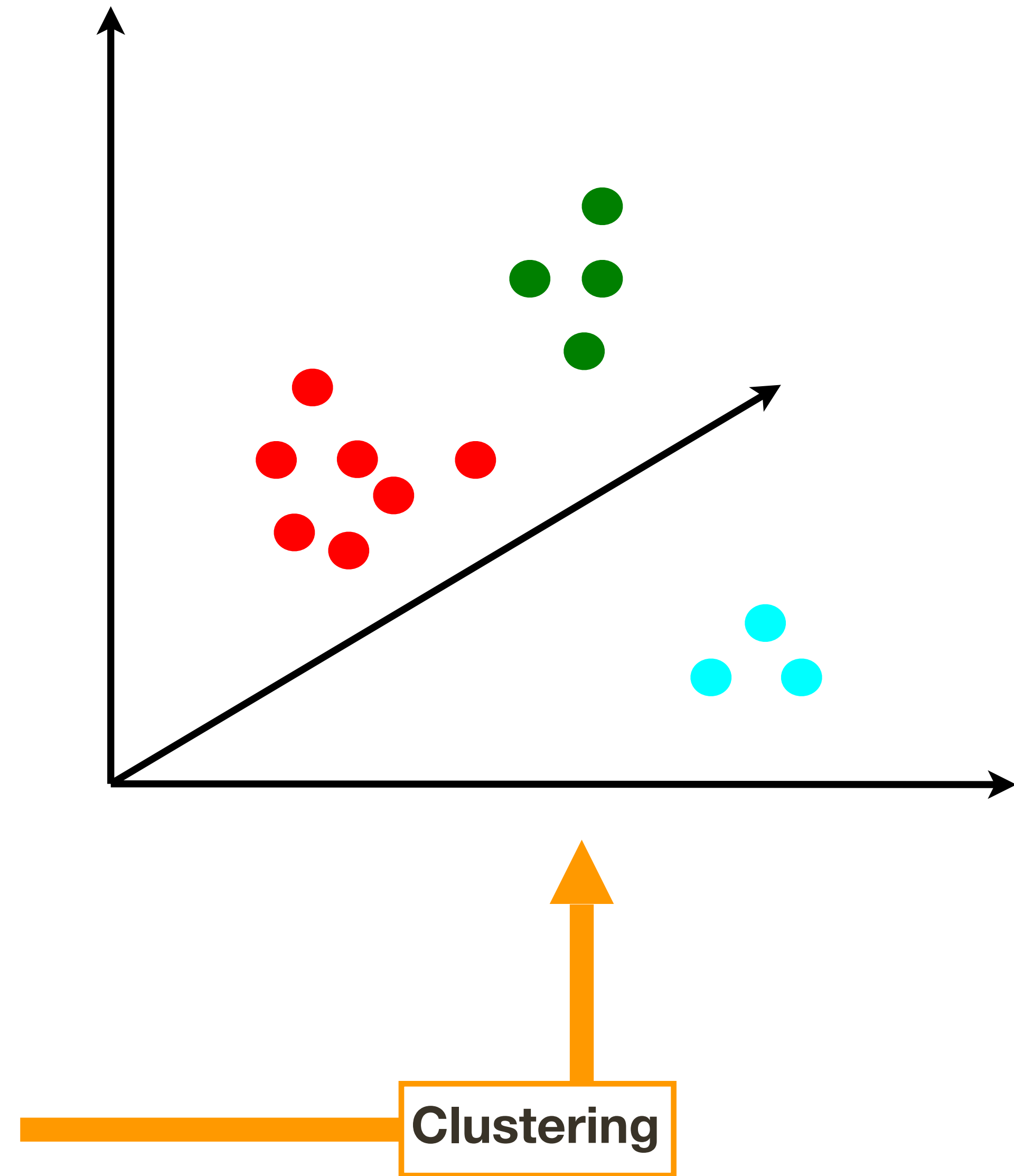
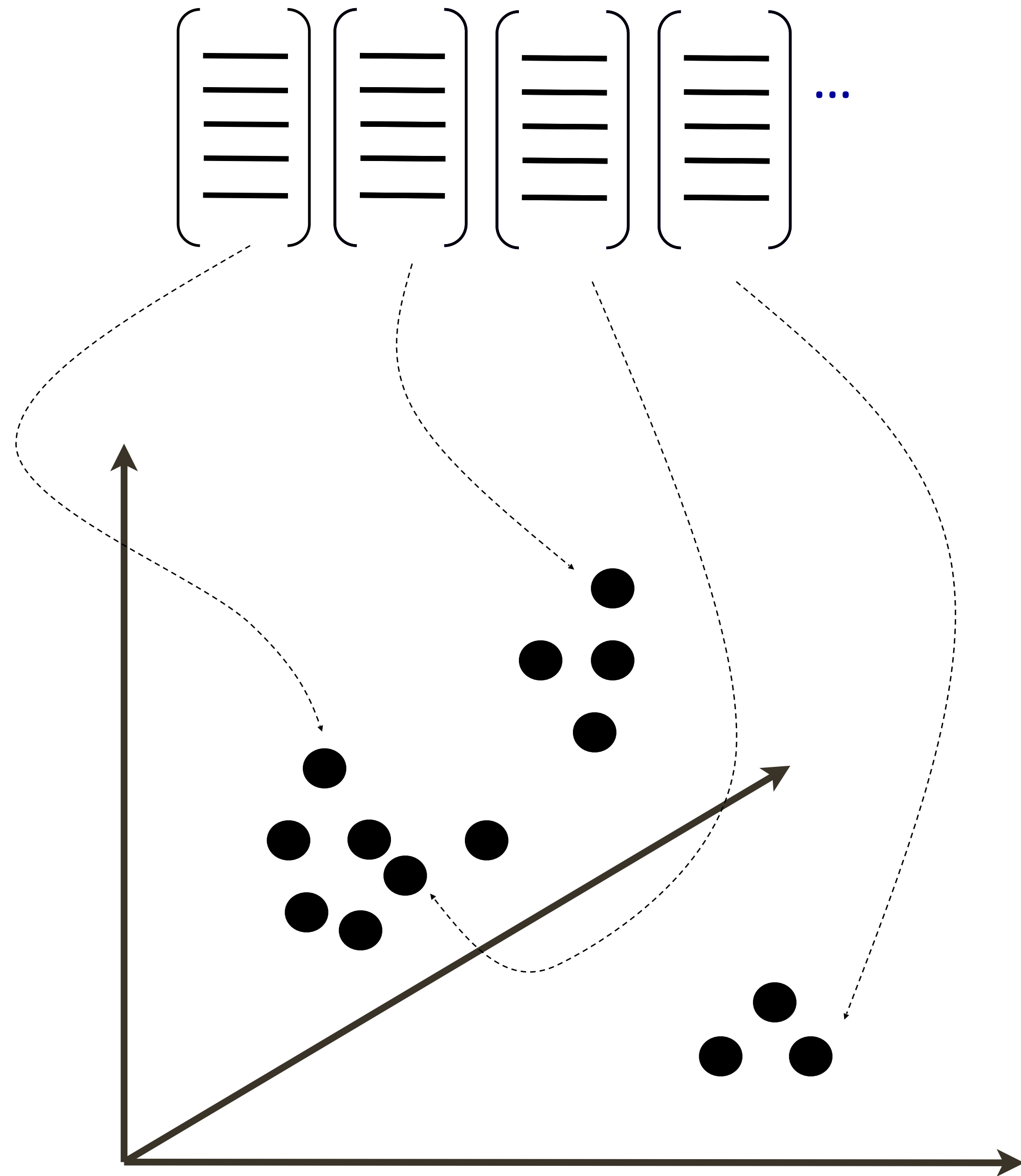
# Extracting **SIFT** Patches



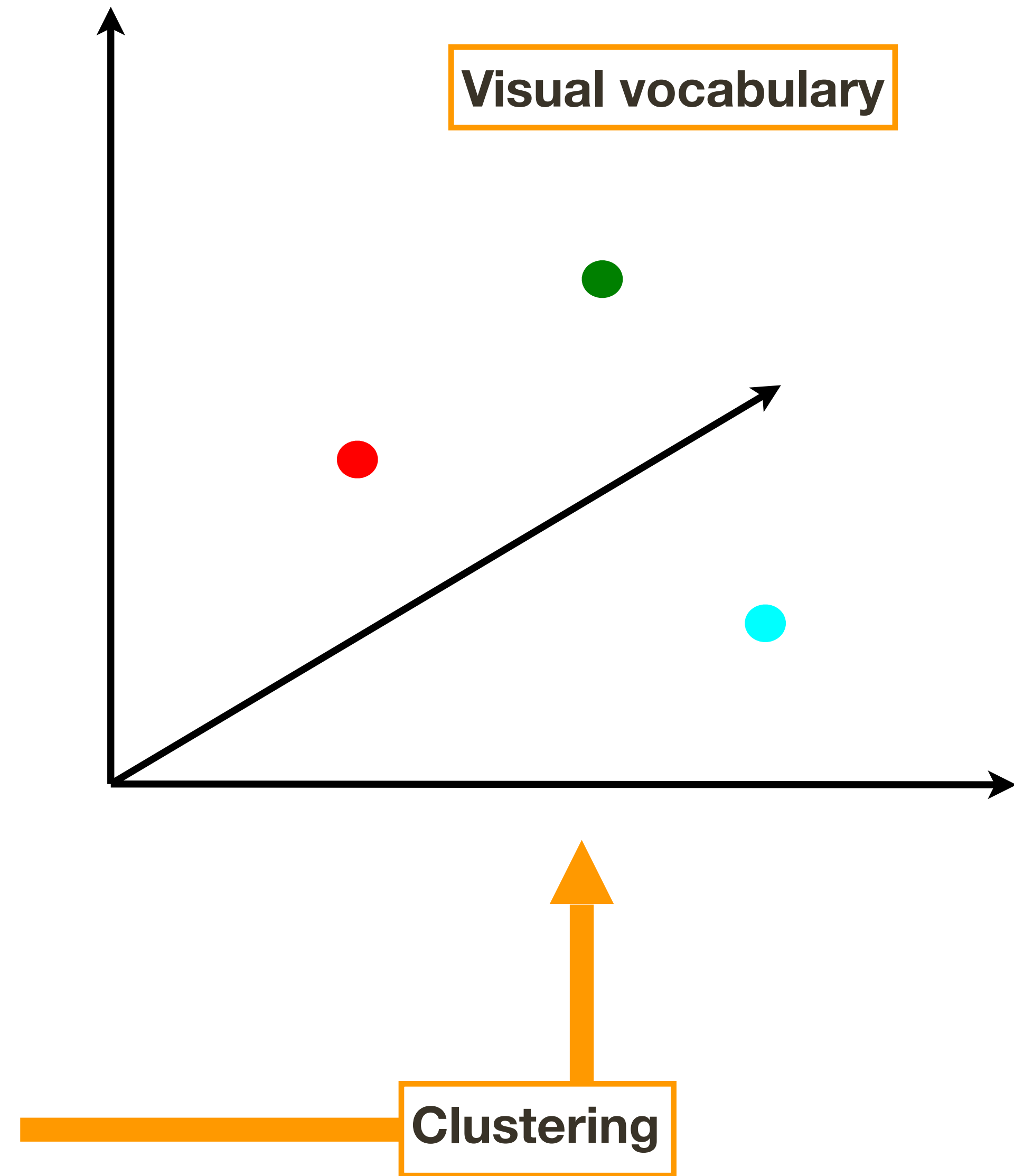
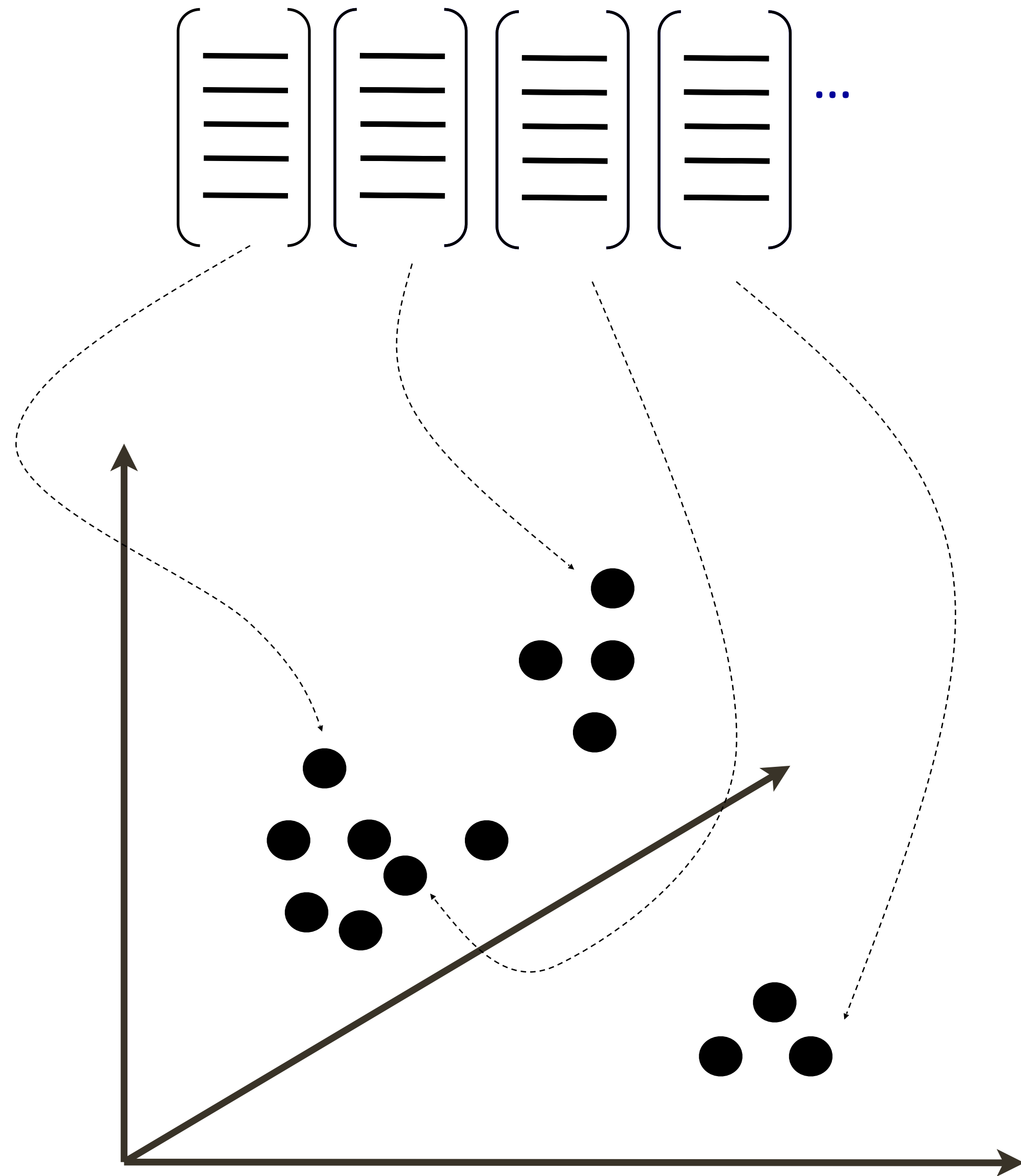
# Creating Dictionary



# Creating Dictionary



# Creating Dictionary



# **K-means** clustering

# K-means Clustering

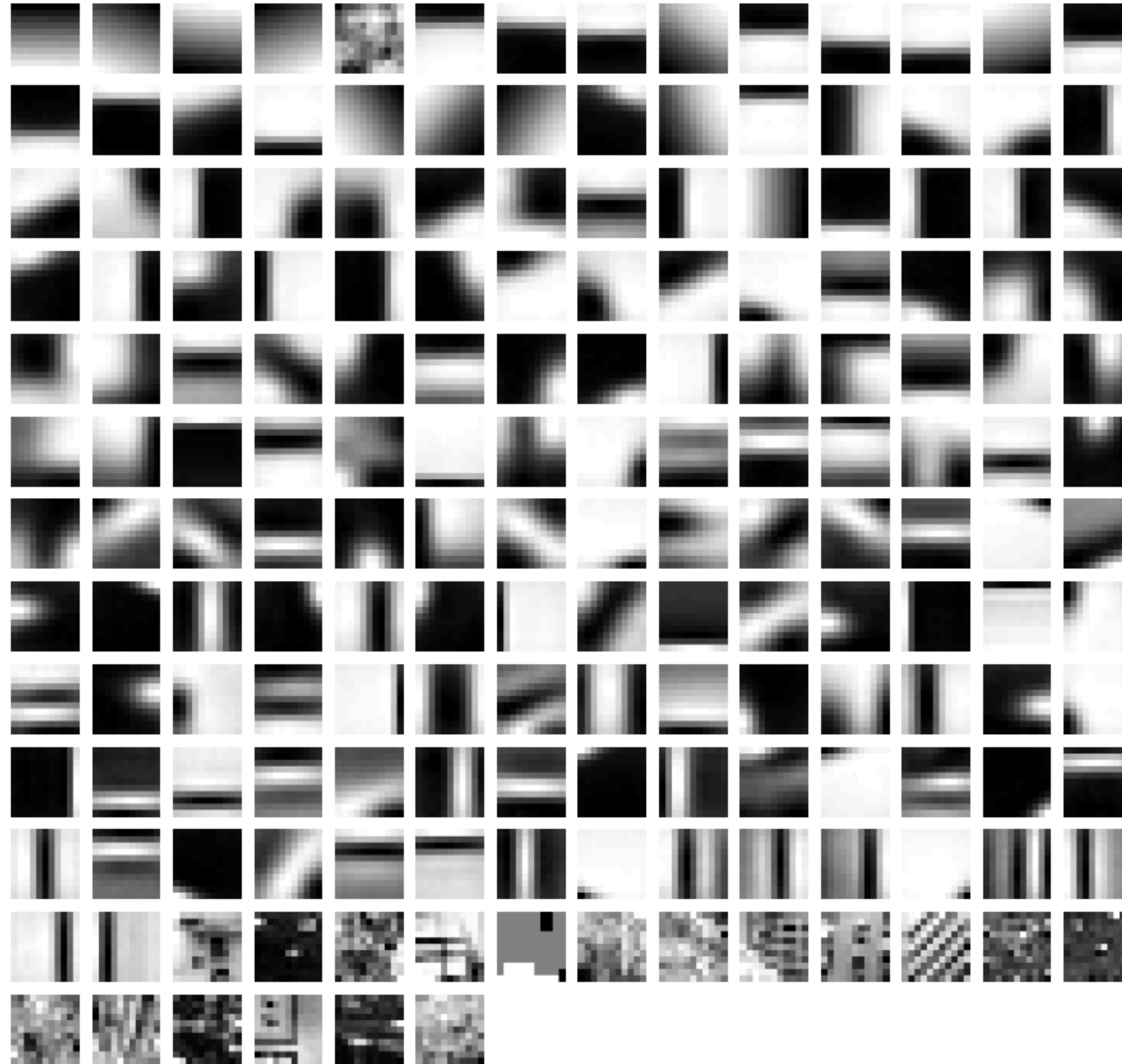
**K-means** is a clustering technique that iterates between

- 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).

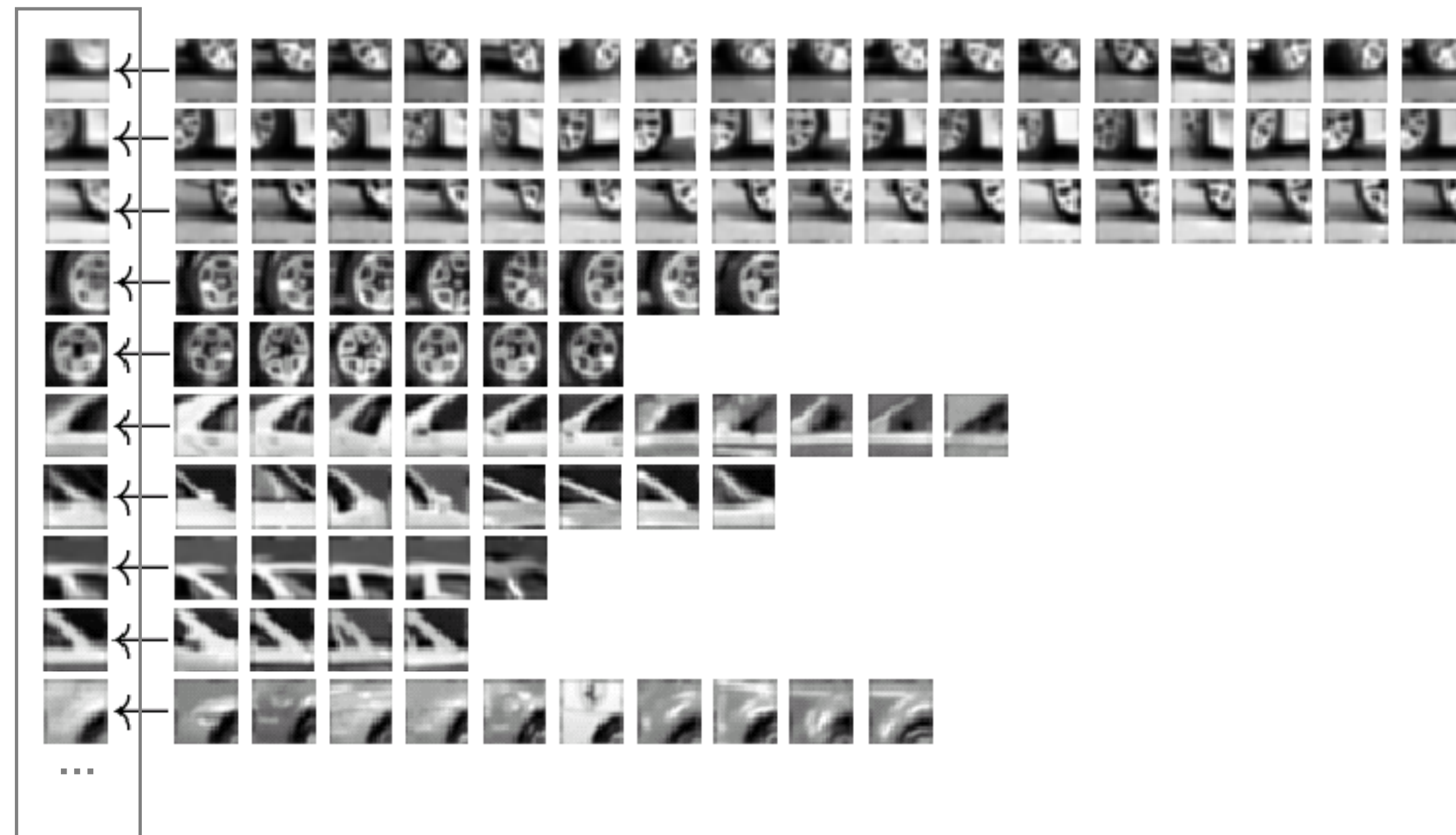
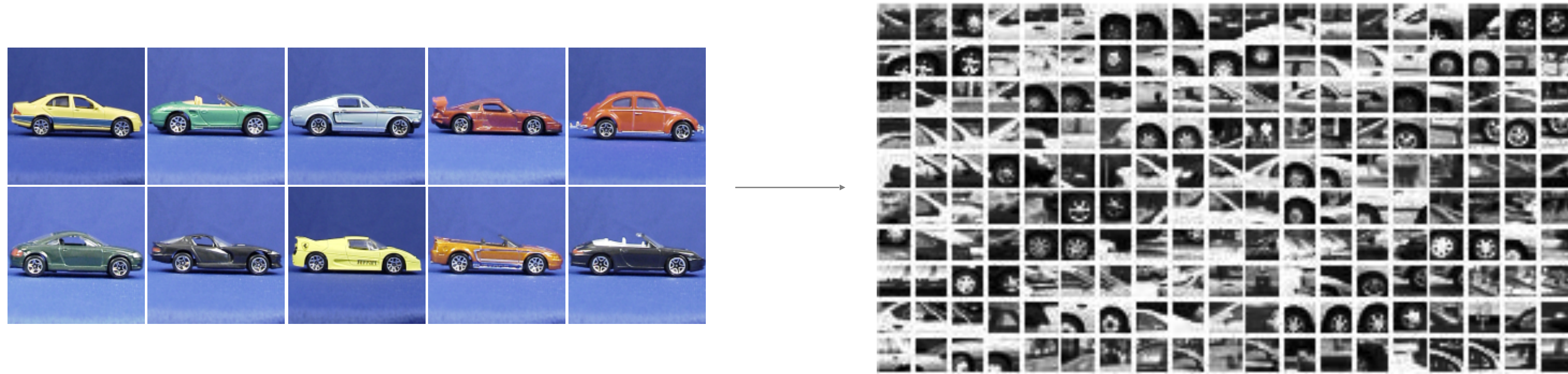
**K-means** clustering is initialization dependent and converges to a local minimum



# Example **Visual Dictionary**

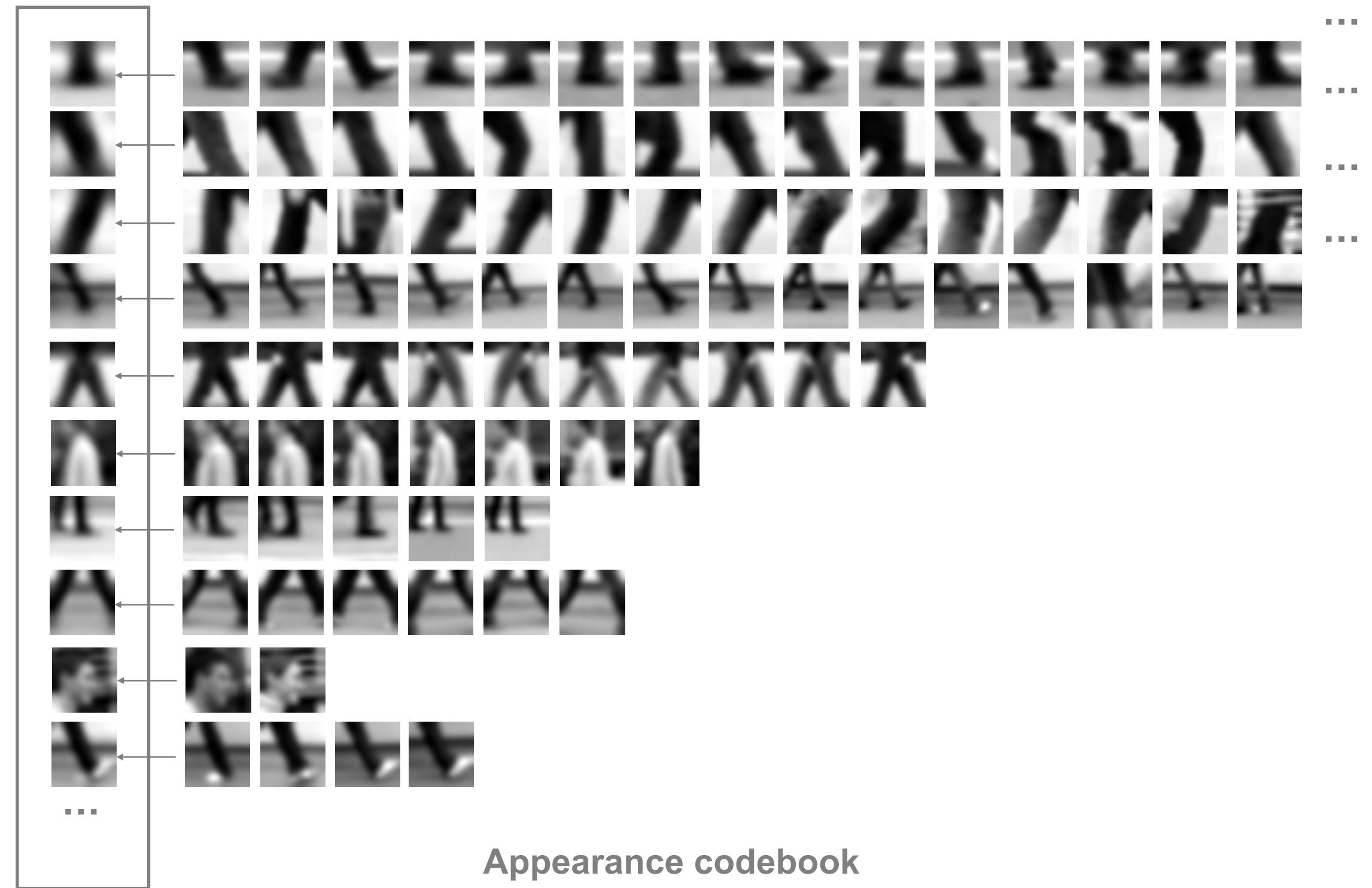


# Example **Visual Dictionary**



**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

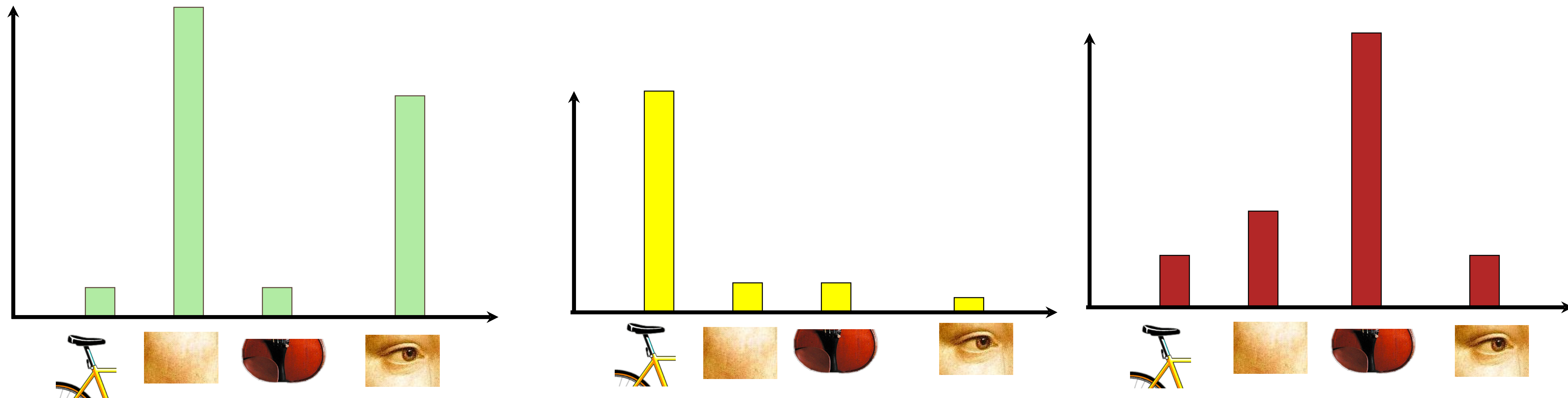


1. **Quantization:** image features gets associated to a visual word (nearest cluster center)

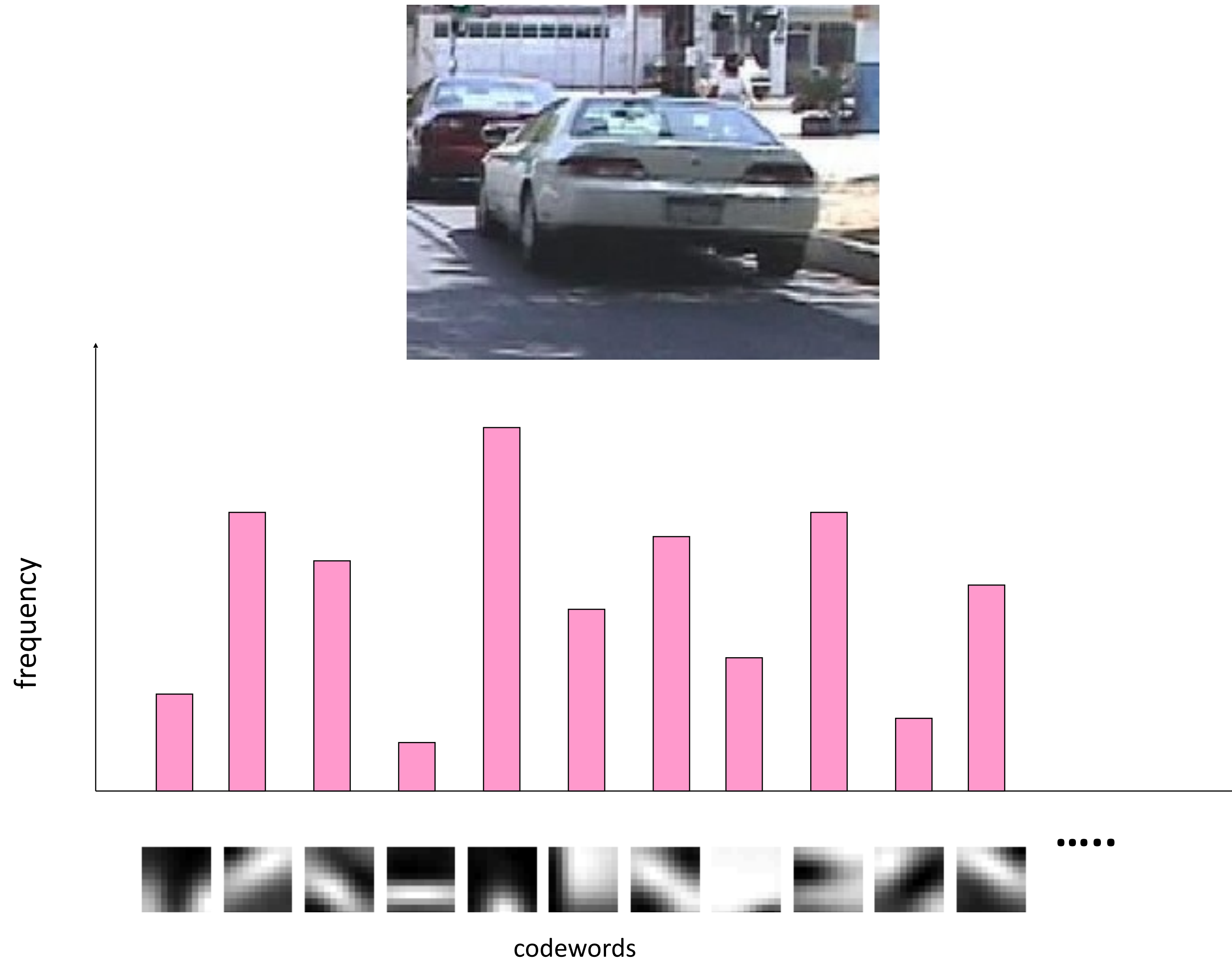


## 2. Encode: build Bag-of-Words (BOW) vectors for each image

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:**

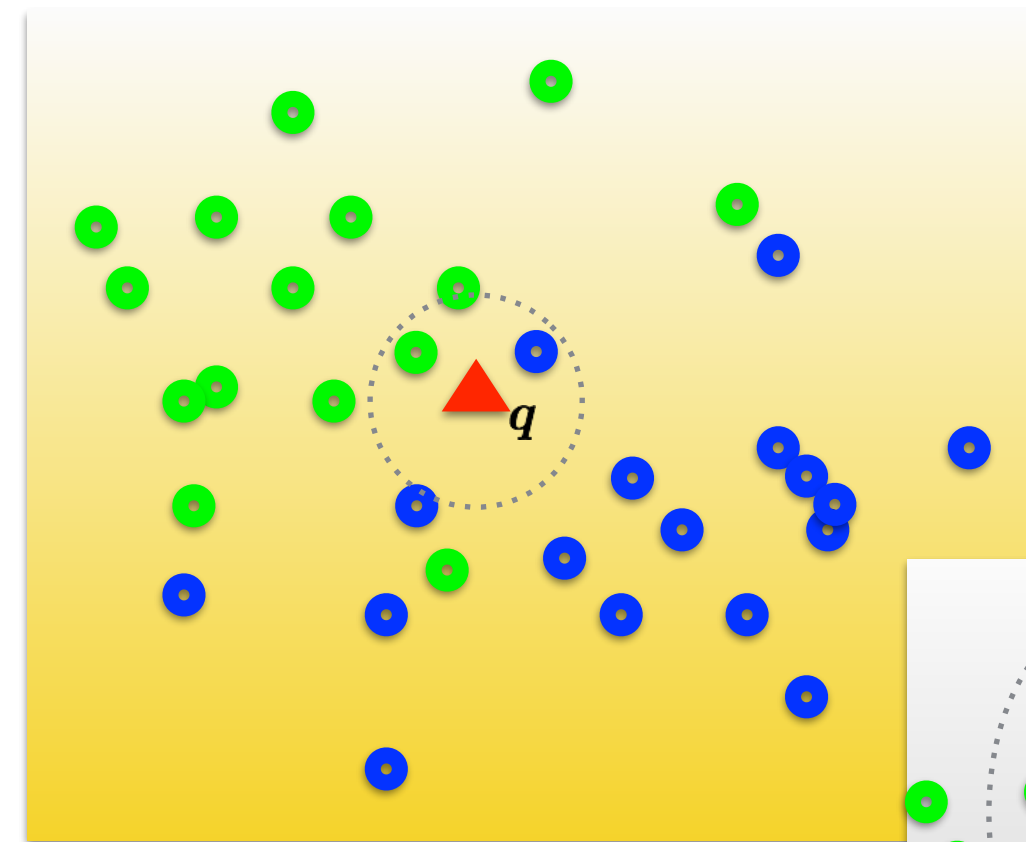
build Bags-of-Words (BOW) vectors  
for each image

## **Classify:**

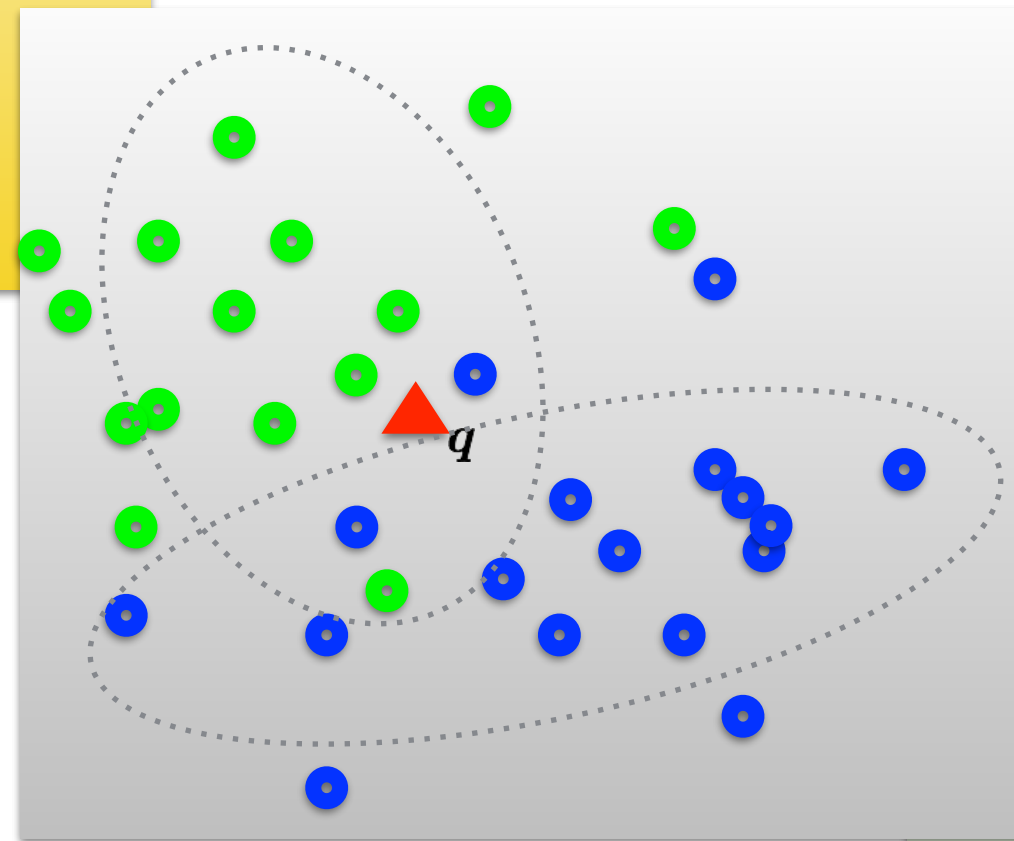
Train and test data using BOWs



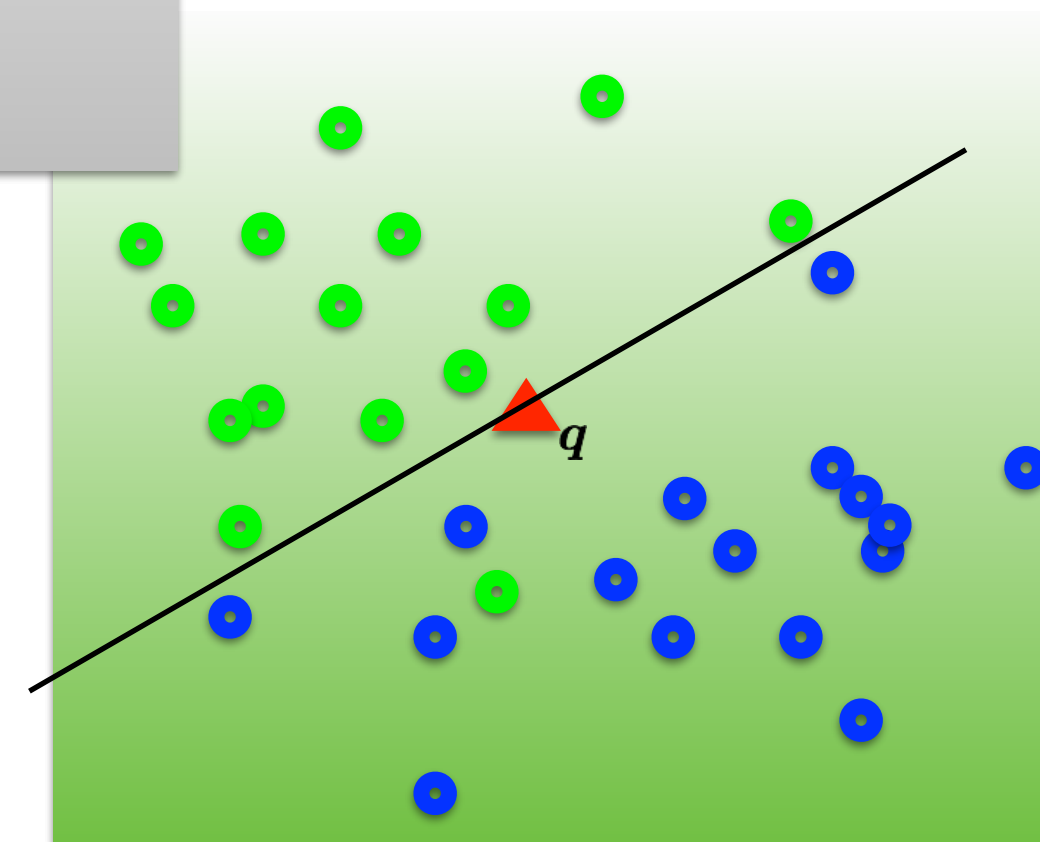
### 3. Classify: Train and text classifier using BOWs



K nearest neighbors



Naïve Bayes



Support Vector Machine

# Bag-of-Words Representation

## Algorithm:

Initialize an empty  $K$  -bin histogram, where  $K$  is the number of codewords

Extract local descriptors (e.g. SIFT) from the image

For each local descriptor  $\mathbf{x}$

    Map (Quantize)  $\mathbf{x}$  to its closest codeword  $\rightarrow \mathbf{c}(\mathbf{x})$

    Increment the histogram bin for  $\mathbf{c}(\mathbf{x})$

Return histogram

We can then classify the histogram using a trained classifier, e.g. a support vector machine or k-Nearest Neighbor classifier

# Spatial Pyramid

The bag of words representation does not preserve any spatial information

The **spatial pyramid** is one way to incorporate spatial information into the image descriptor.

A spatial pyramid partitions the image and counts codewords within each grid box; this is performed at multiple levels

# Spatial Pyramid

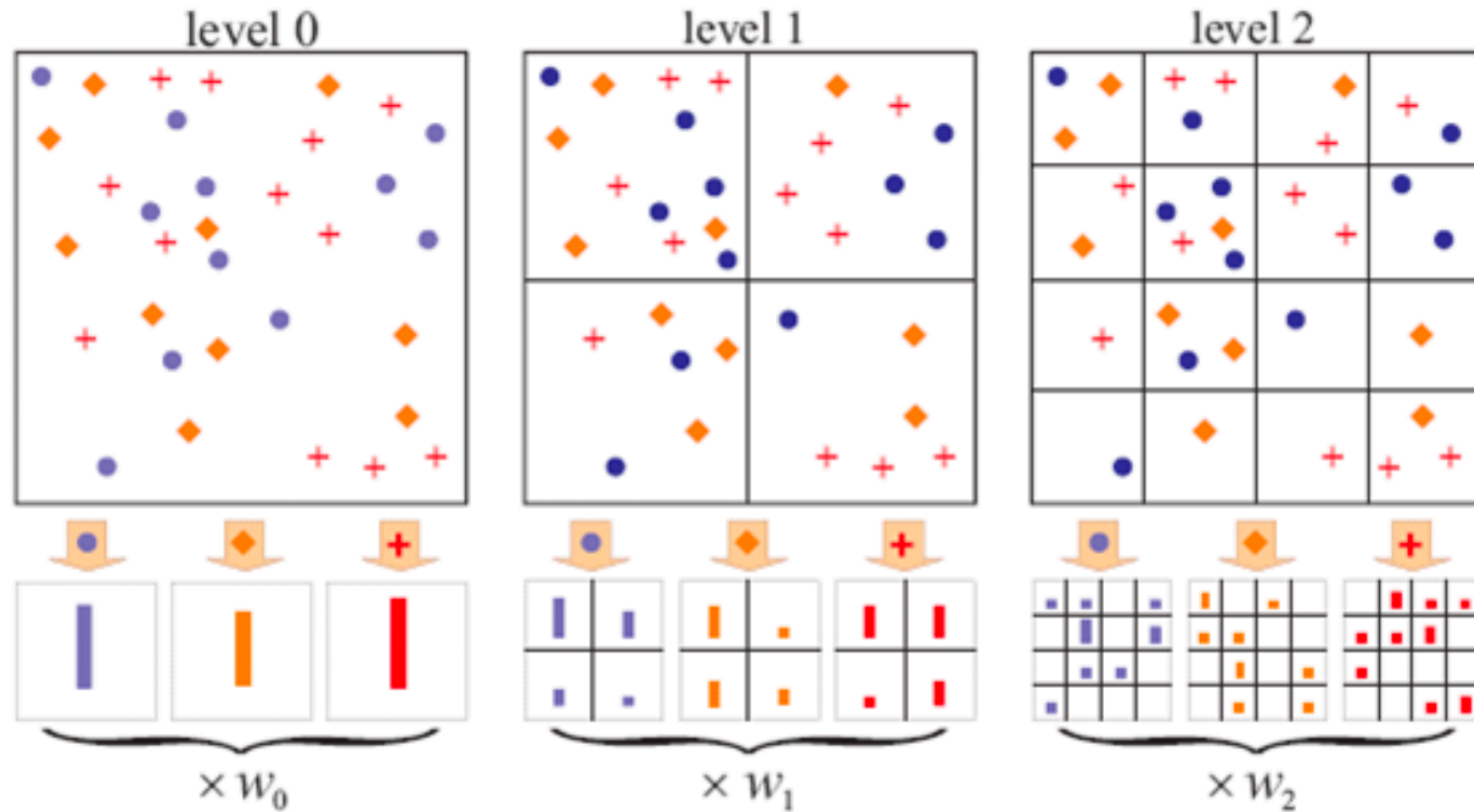


Fig. 16.8 in Forsyth & Ponce (2nd ed.).  
Original credit: Lazebnik et al., 2006

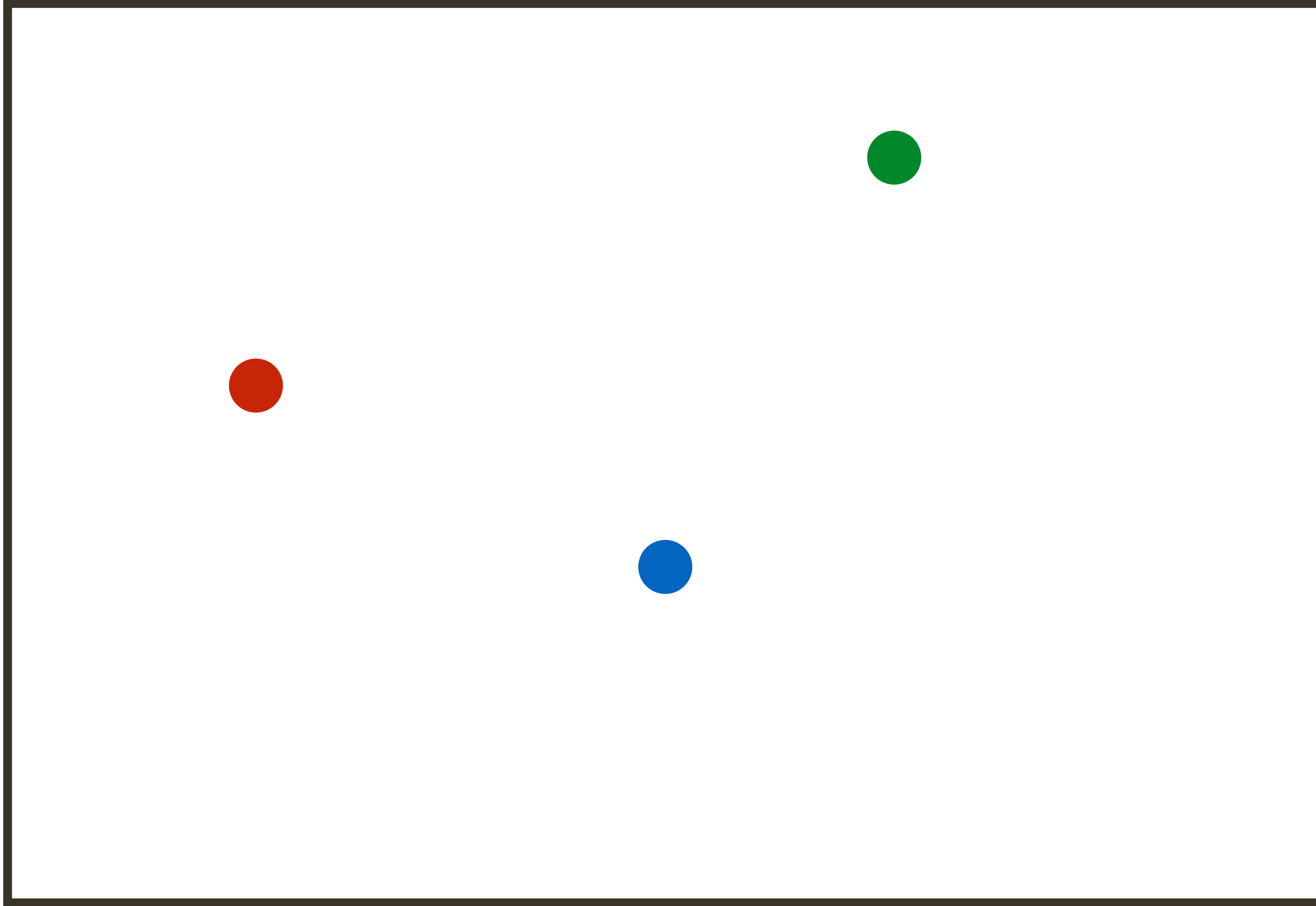
# **VLAD** (Vector of Locally Aggregated Descriptors)

There are more advanced ways to ‘count’ visual words than incrementing its histogram bin

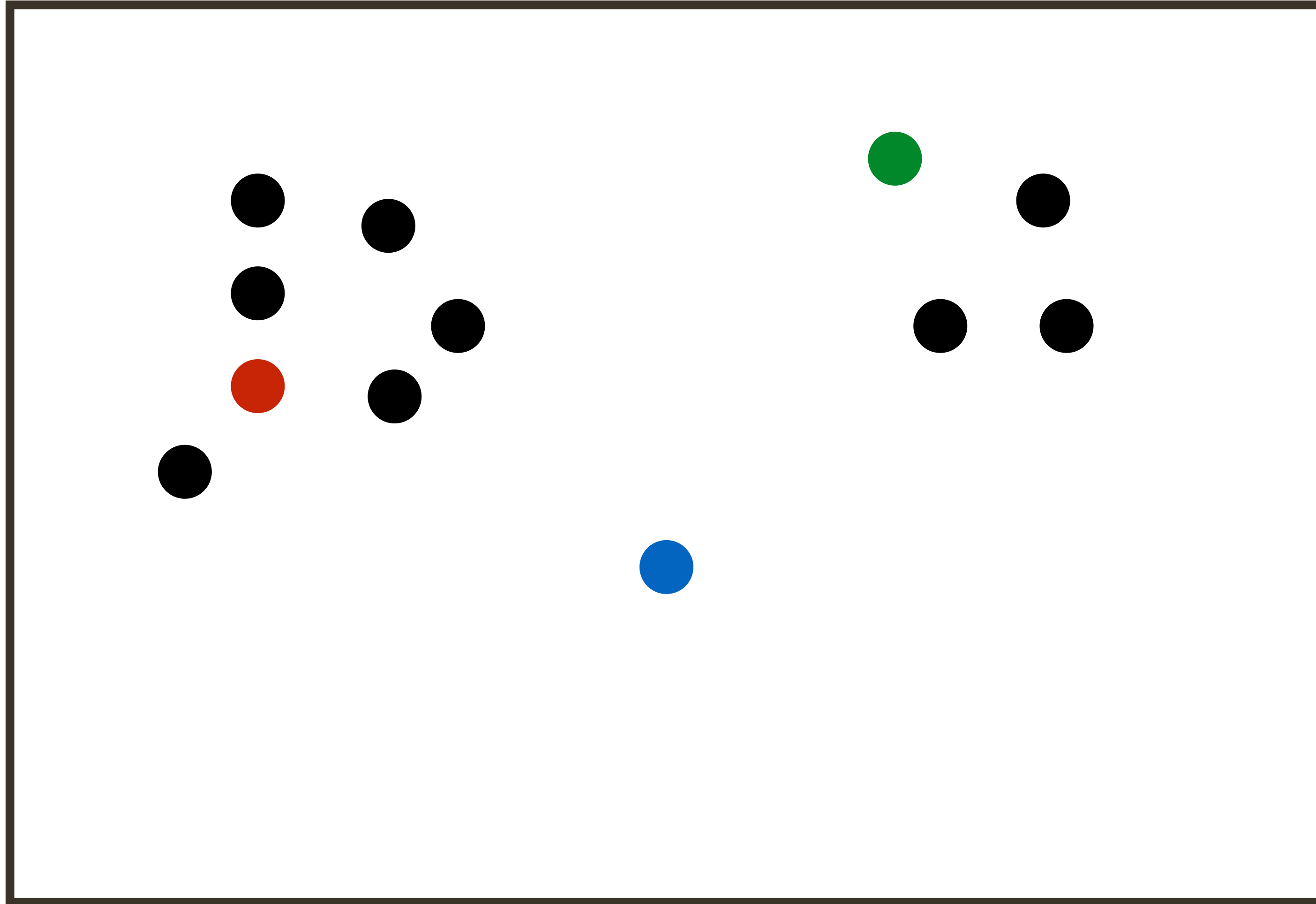
For example, it might be useful to describe how local descriptors are quantized to their visual words

In the VLAD representation, instead of incrementing the histogram bin by one, we increment it by the **residual** vector  $\mathbf{x} - \mathbf{c}(\mathbf{x})$

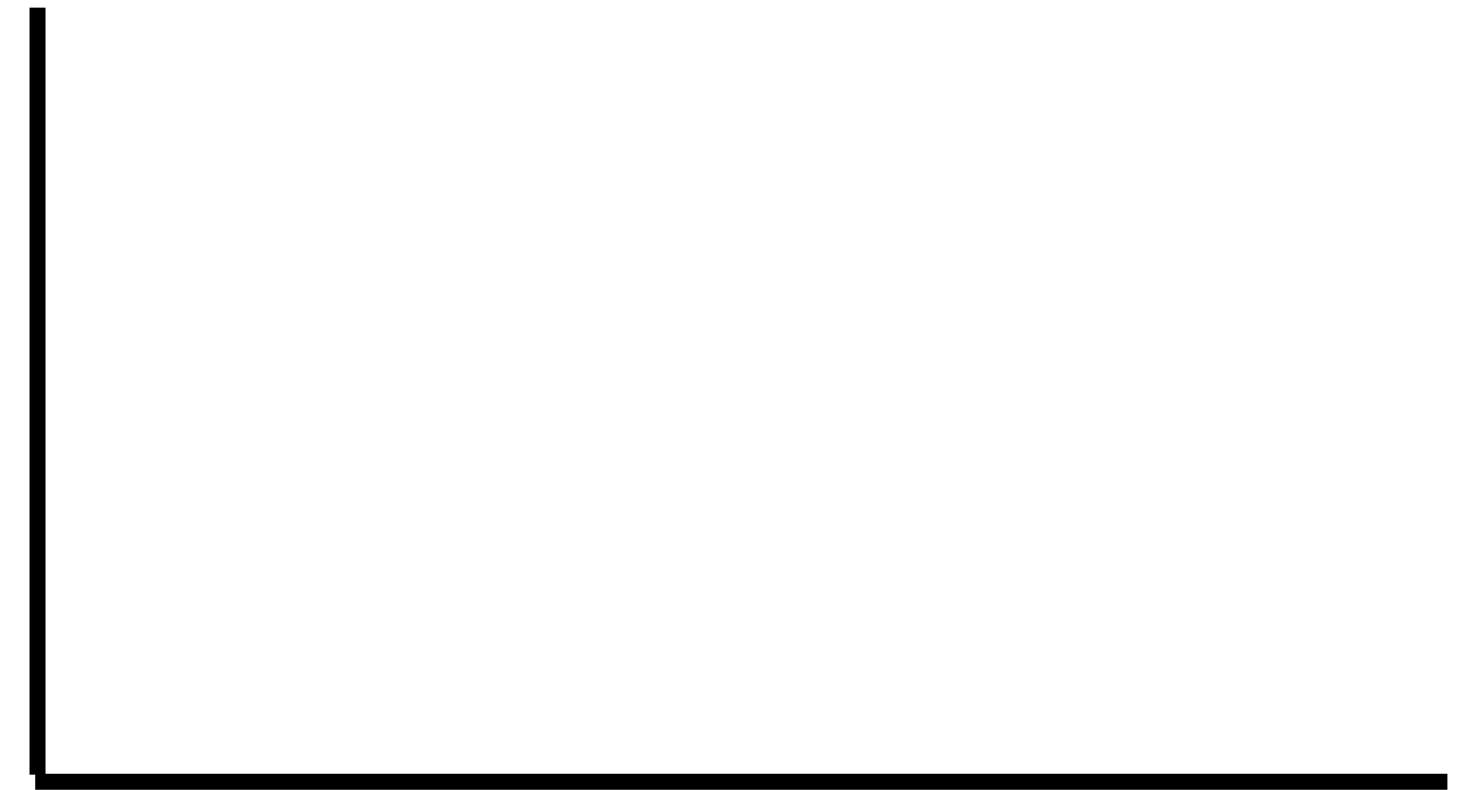
# Example: VLAD



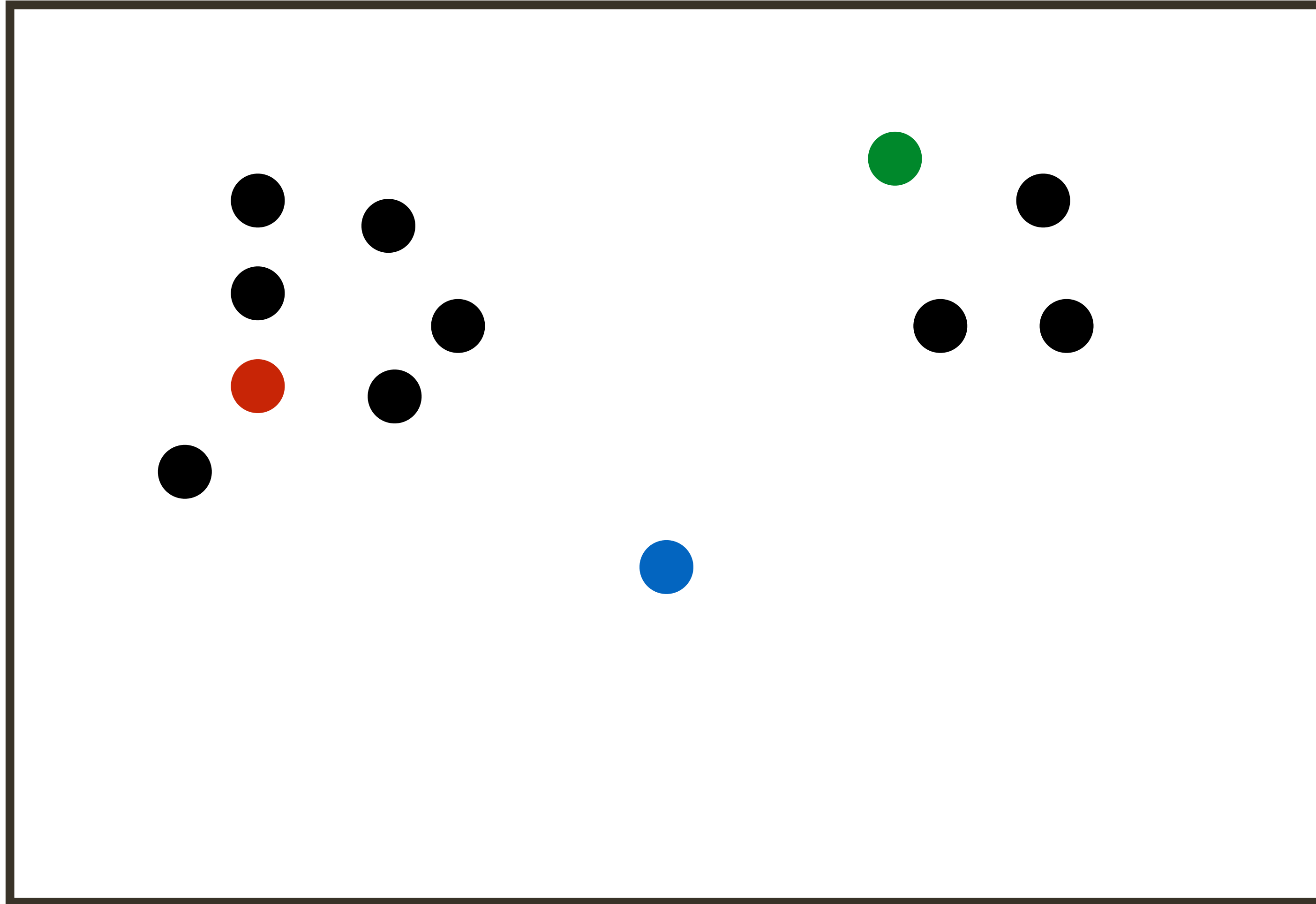
# Example: VLAD



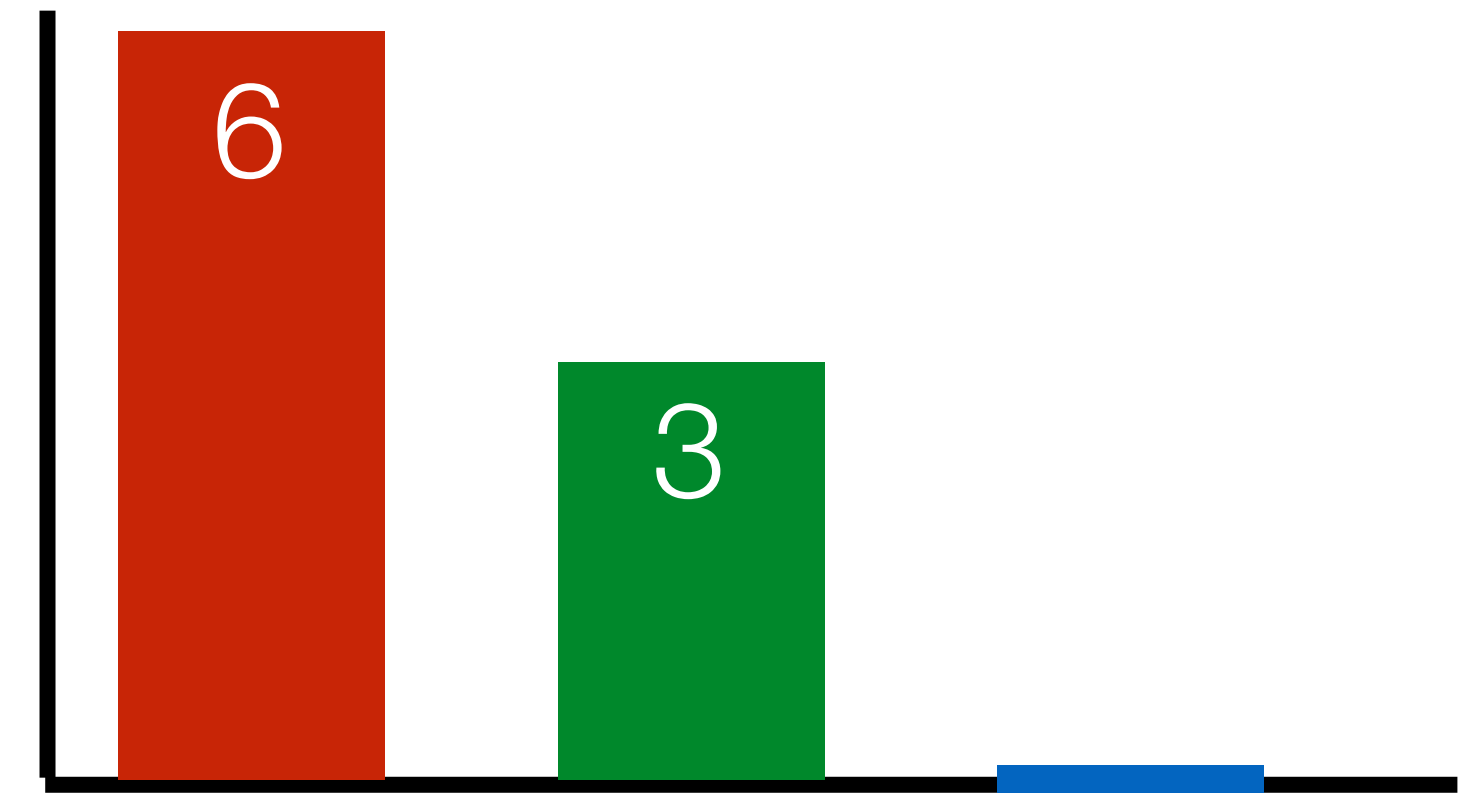
Bag of Word



# Example: VLAD

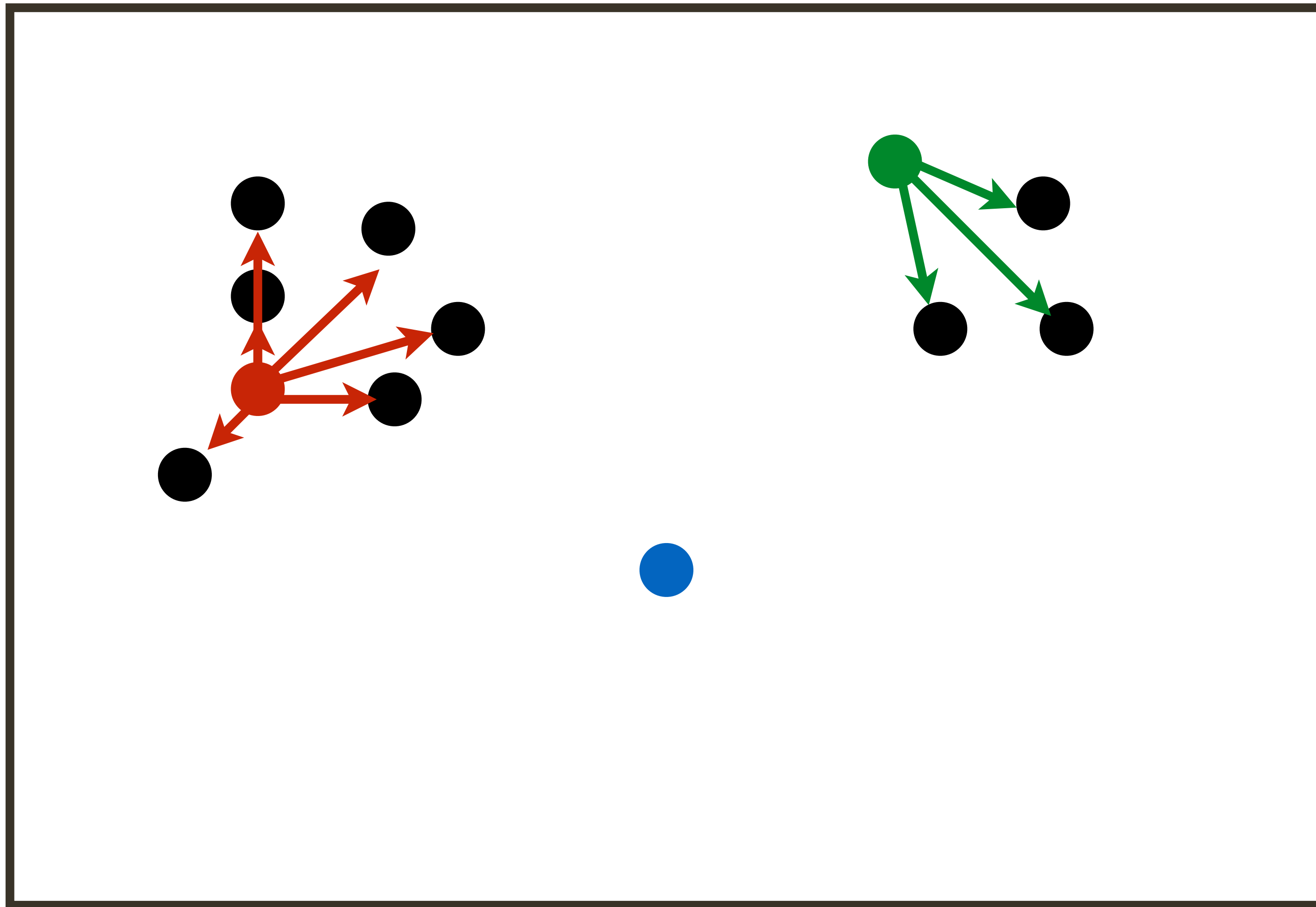


Bag of Word

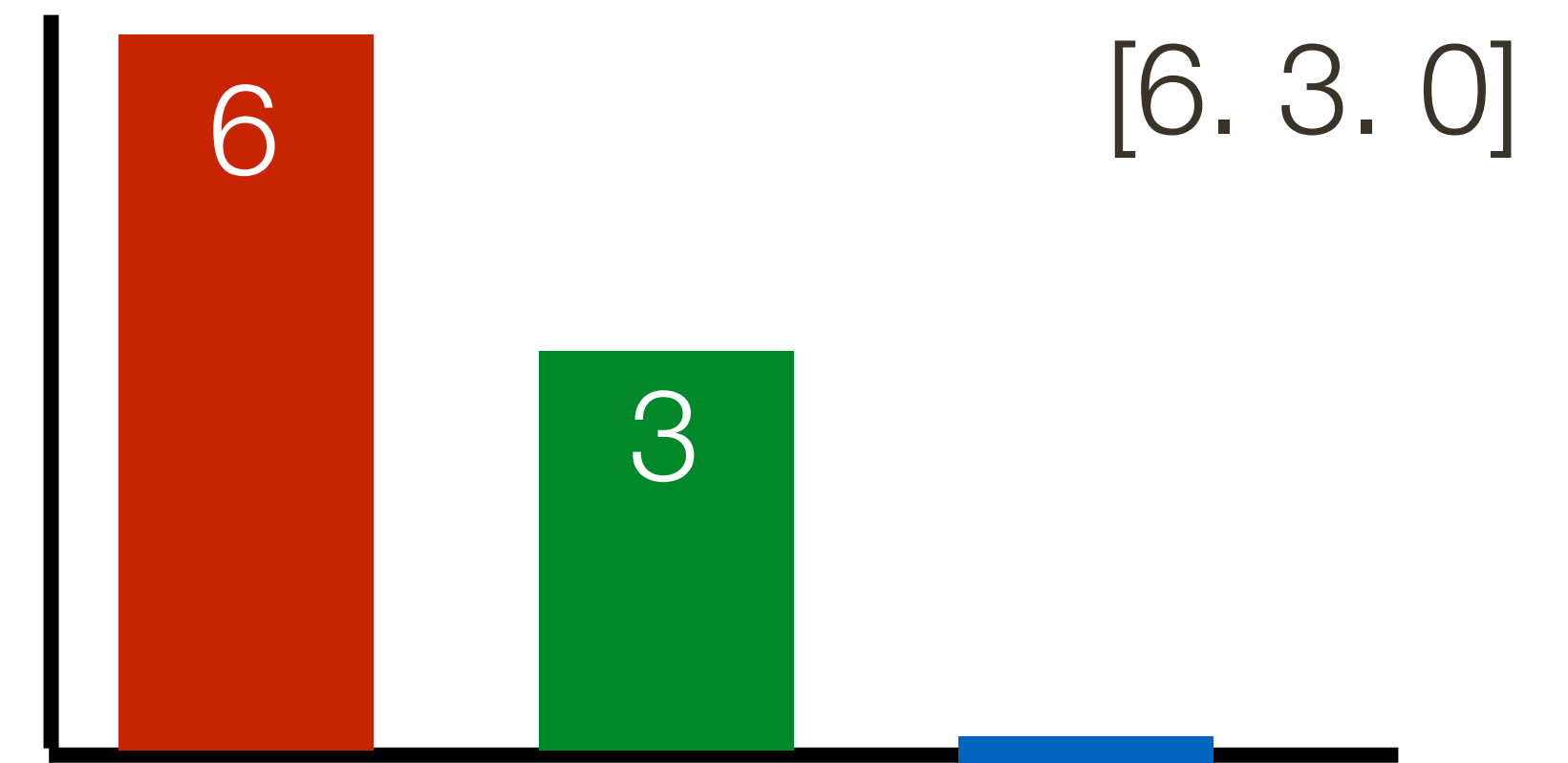




# Example: VLAD

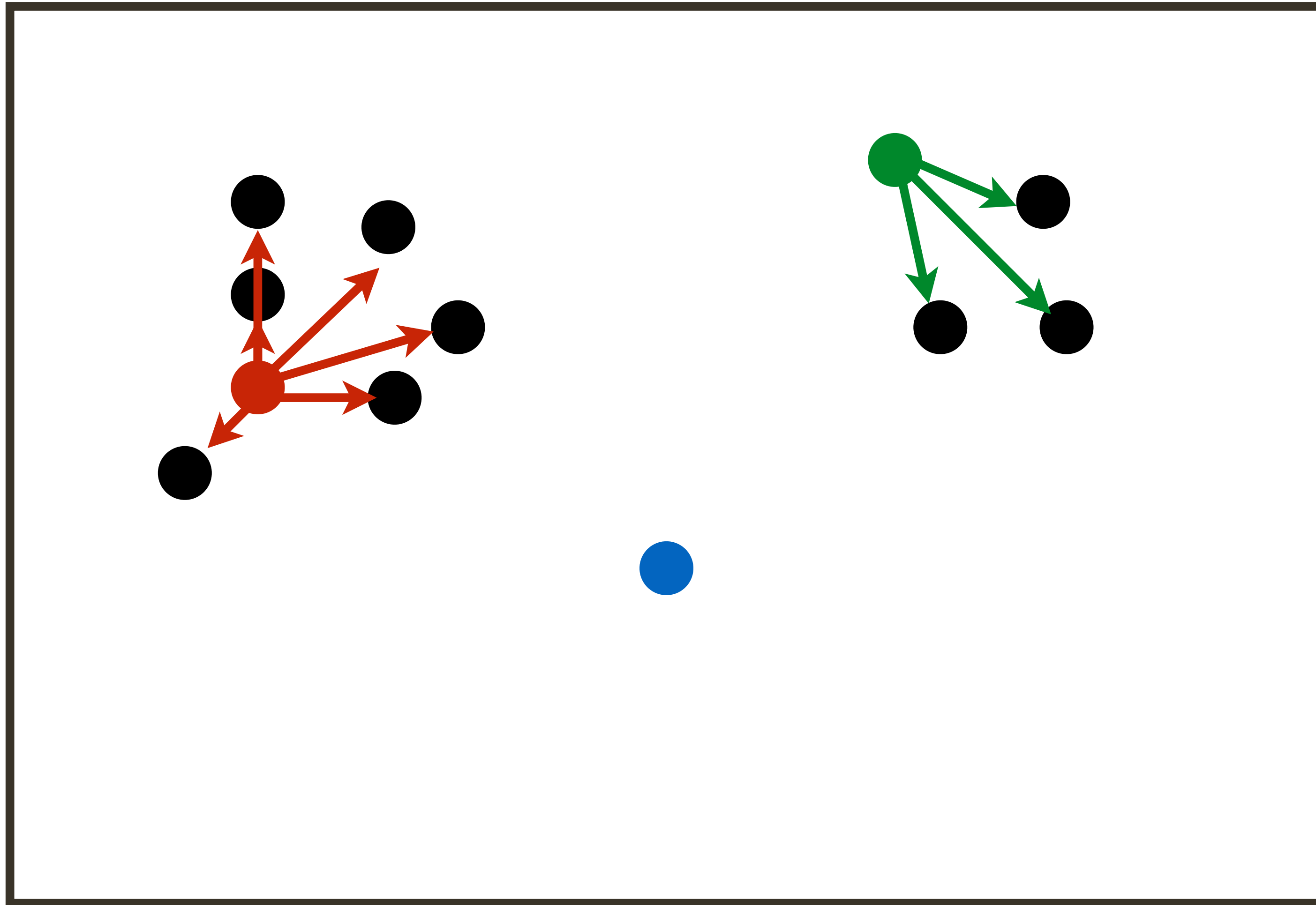


Bag of Word

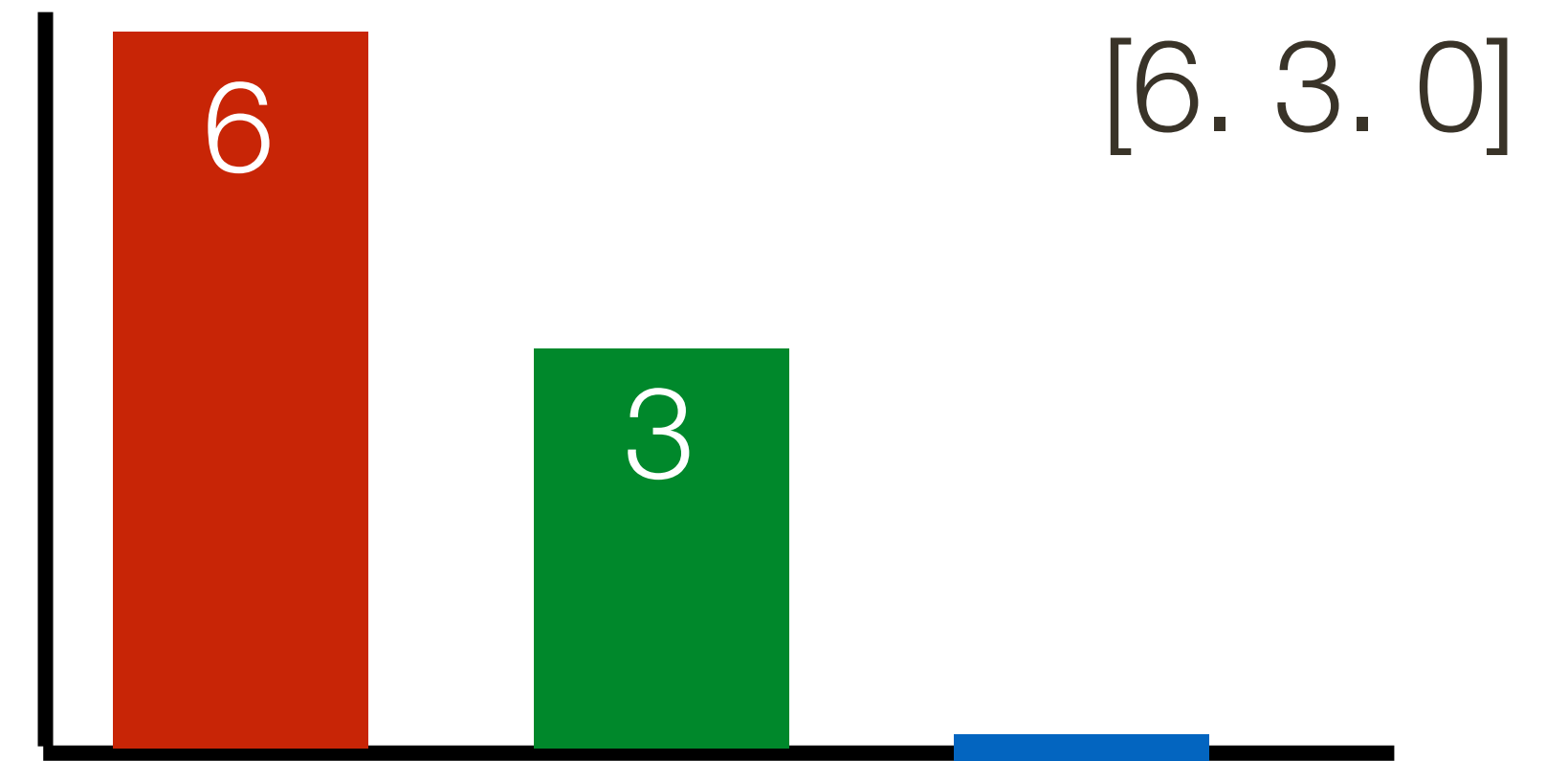


VLAD

# Example: VLAD



Bag of Word



VLAD



# **VLAD** (Vector of Locally Aggregated Descriptors)

The dimensionality of a **VLAD** descriptor is  $Kd$

- $K$  : number of codewords
- $d$  : dimensionality of the local descriptor

**VLAD** characterizes the distribution of local descriptors with respect to the codewords

# Summary

Factors that make image classification hard

— intra-class variation, viewpoint, illumination, clutter, occlusion...

A codebook of **visual words** contains representative local patch descriptors

— can be constructed by clustering local descriptors (e.g. SIFT) in training images

The **bag of words** model accumulates a histogram of occurrences of each visual word

The **spatial pyramid** partitions the image and counts visual words within each grid box; this is repeated at multiple levels

# Back to **Classification**

# Decision Tree

A **decision tree** is a simple non-linear parametric classifier

Consists of a tree in which each internal node is associated with a feature test

A data point starts at the root and recursively proceeds to the child node determined by the feature test, until it reaches a leaf node

The leaf node stores a class label or a probability distribution over class labels

# Decision Tree

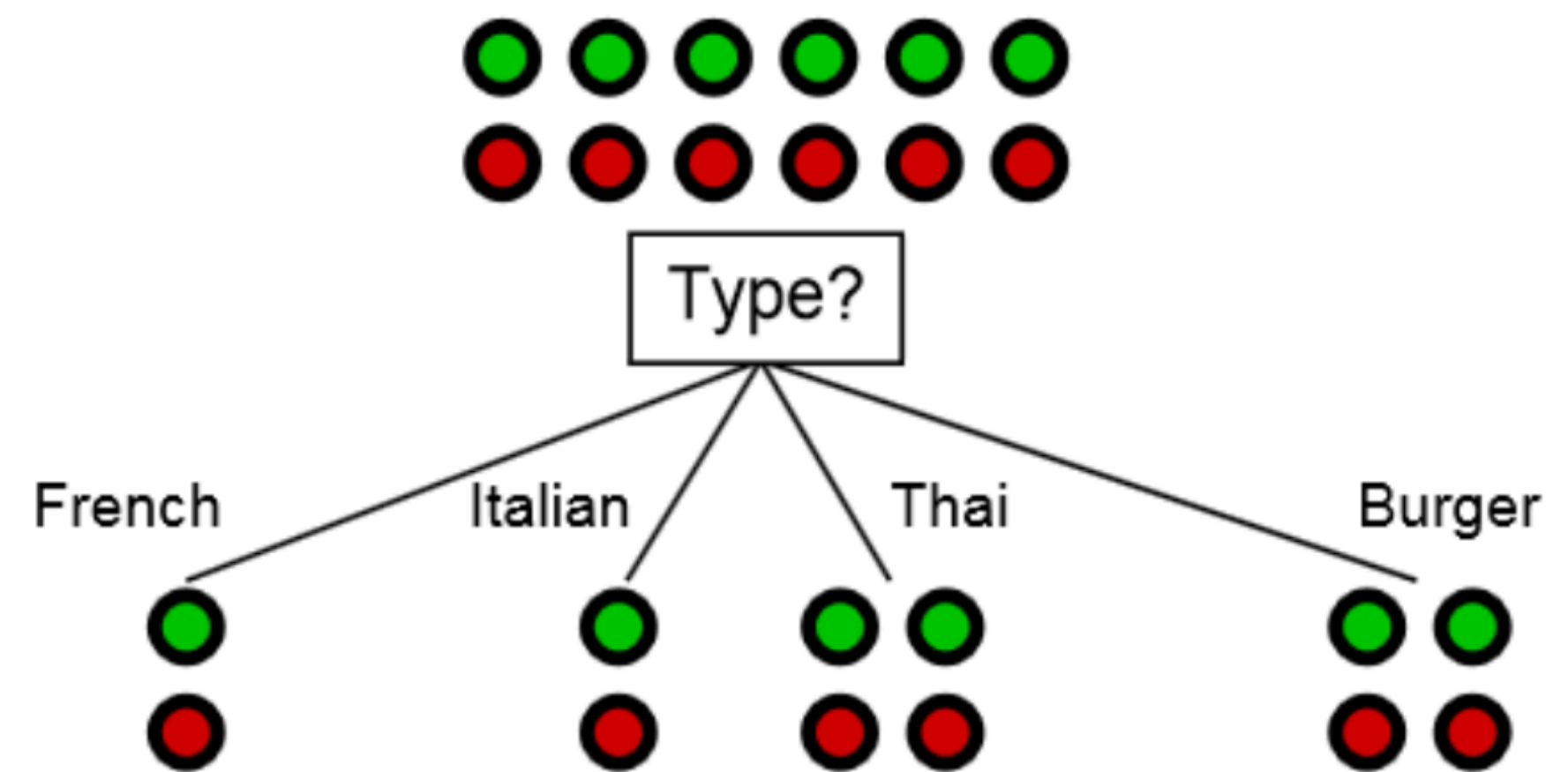
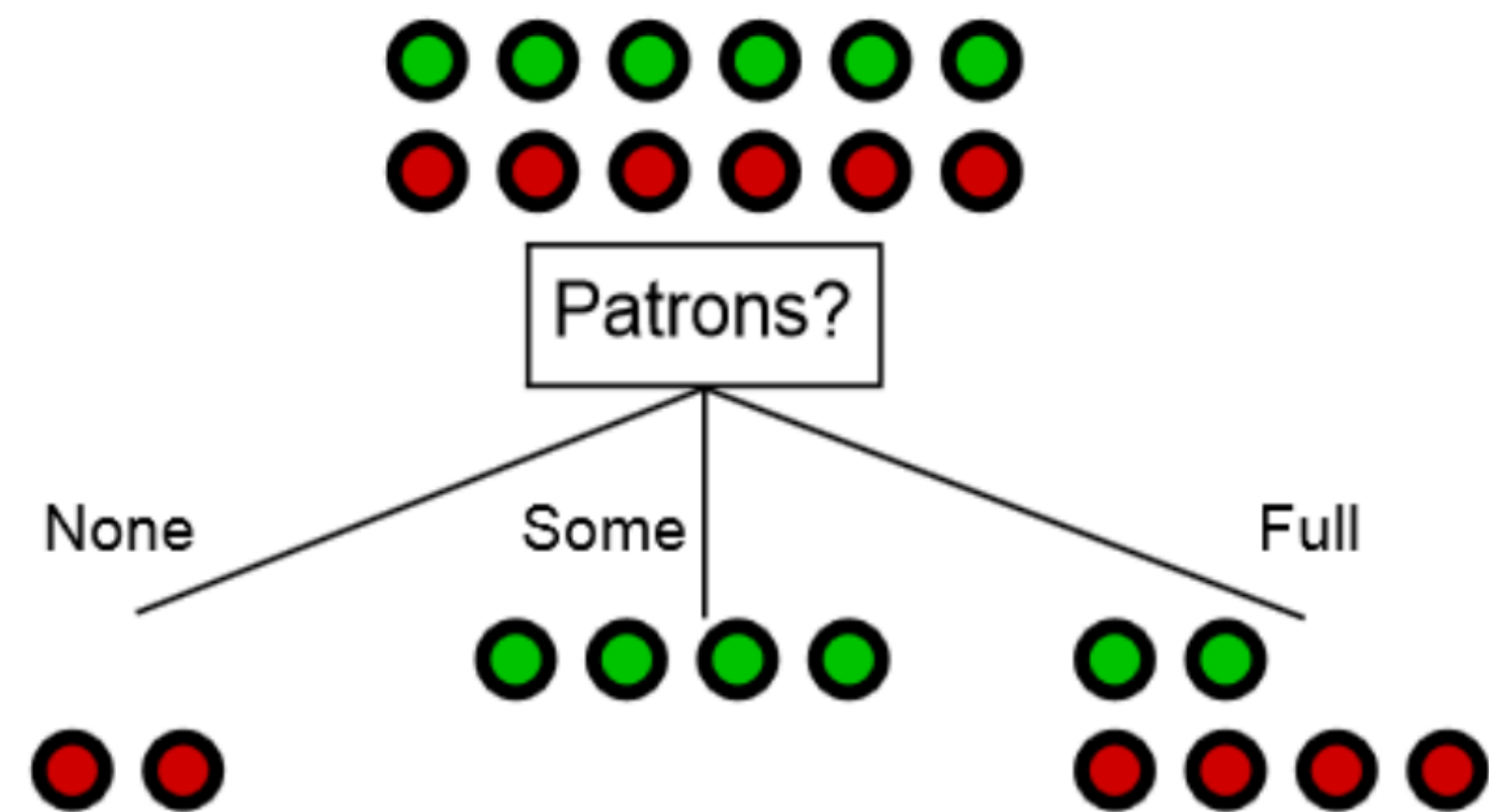
Learning a decision tree from a training set involves selecting an efficient sequence of feature tests

**Example:** Waiting for a restaurant table

Example	Attributes										Target
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$X_1$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0-10</i>	<i>T</i>
$X_2$	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30-60</i>	<i>F</i>
$X_3$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>T</i>
$X_4$	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10-30</i>	<i>T</i>
$X_5$	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>&gt;60</i>	<i>F</i>
$X_6$	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0-10</i>	<i>T</i>
$X_7$	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>F</i>
$X_8$	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0-10</i>	<i>T</i>
$X_9$	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>&gt;60</i>	<i>F</i>
$X_{10}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10-30</i>	<i>F</i>
$X_{11}$	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0-10</i>	<i>F</i>
$X_{12}$	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30-60</i>	<i>T</i>

# Decision Tree

Which test is more helpful?



**Figure credit:** Russell and Norvig (3rd ed.)



# Decision Tree

The **entropy** of a set  $S$  of data samples is defined as

$$H(S) = - \sum_{c \in C} p(c) \log(p(c))$$

where  $C$  is the set of classes represented in  $S$ , and  $p(c)$  is the empirical distribution of class  $c$  in  $S$

Entropy is highest when data samples are spread equally across all classes, and zero when all data samples are from the same class.

# Decision Tree

In general we try to select the feature test that maximizes the **information gain**:

$$I = H(S) - \sum_{i \in \{children\}} \frac{|S^i|}{|S|} H(S^i)$$

In the previous example, the information gains of the two candidate tests are:

$$I_{Patrons} = 0.541 \qquad I_{Type} = 0$$

So we choose the 'Patrons' test.

# Decision Tree

Following this construction procedure we obtain the final decision tree:

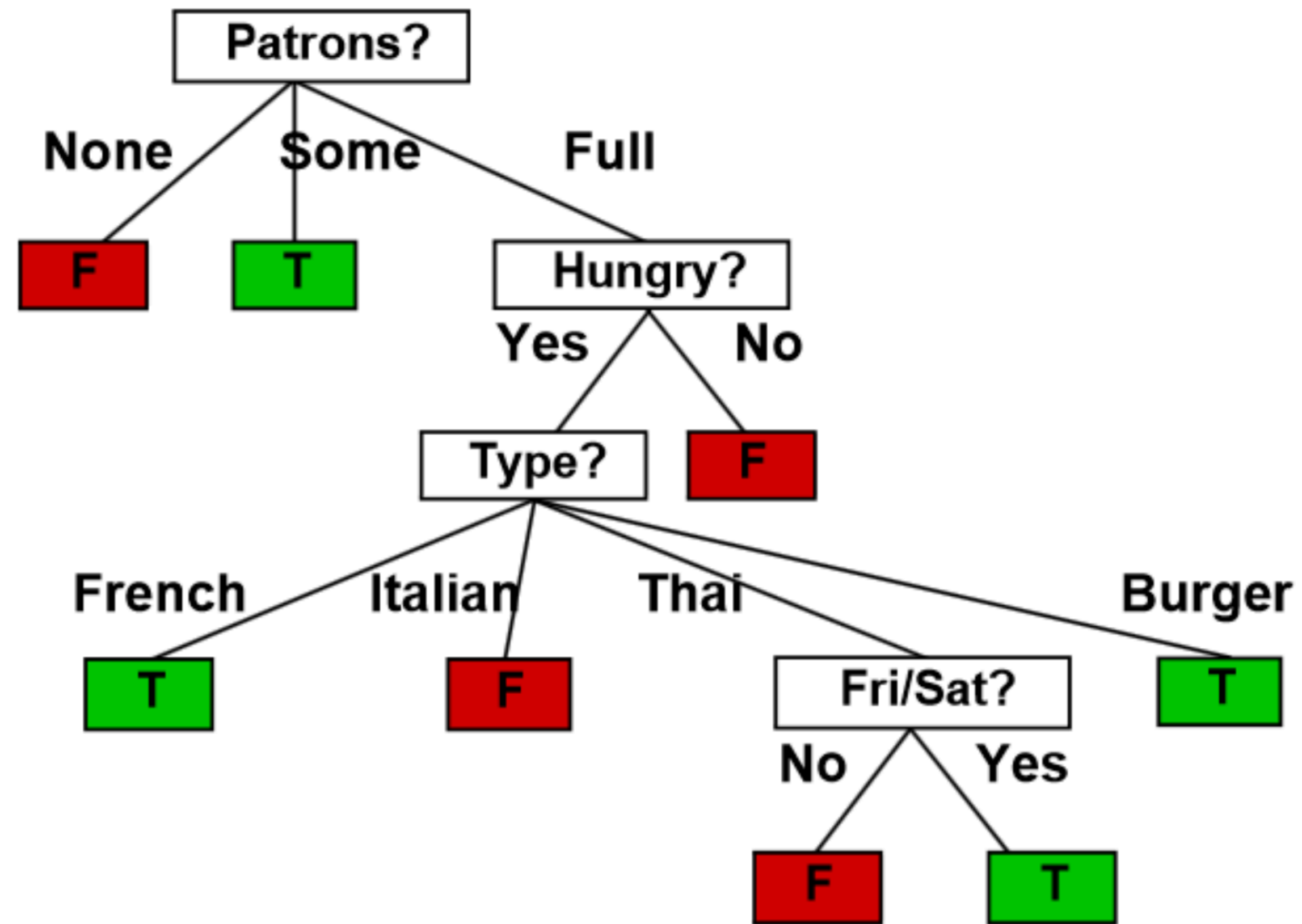


Figure credit: Russell and Norvig (3rd ed.)