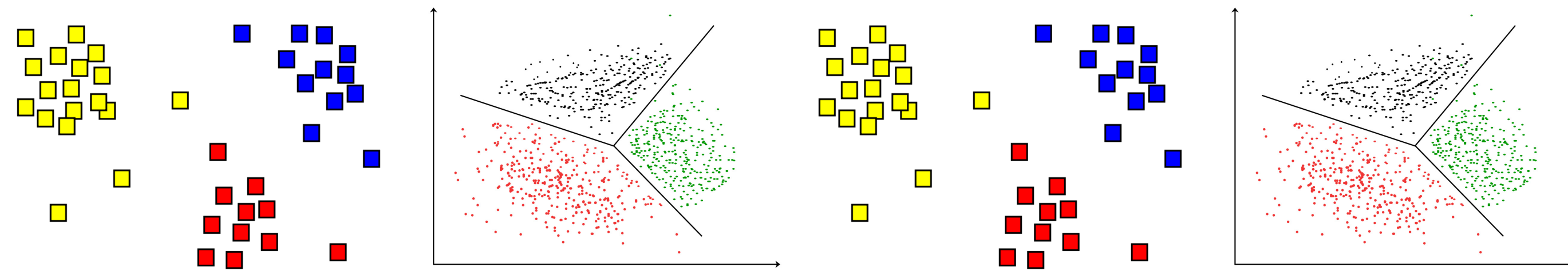




# CPSC 425: Computer Vision



## Lecture 26: Classification

# Menu for Today (November 13, 2020)

## Topics:

- Classification (cont)
- kNN, SVMs
- Bag of Words Representation
- Scene Classification

## Readings:

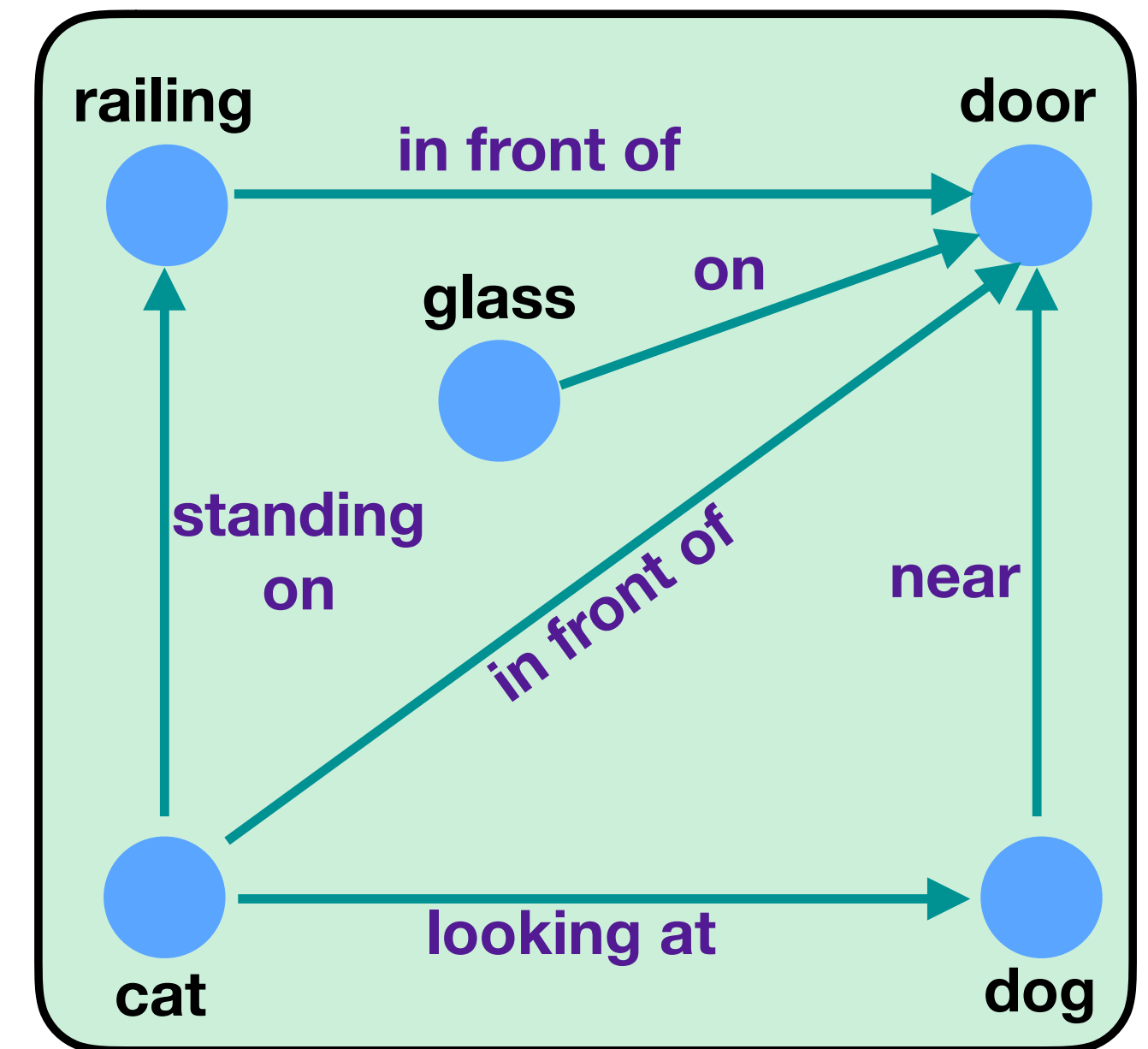
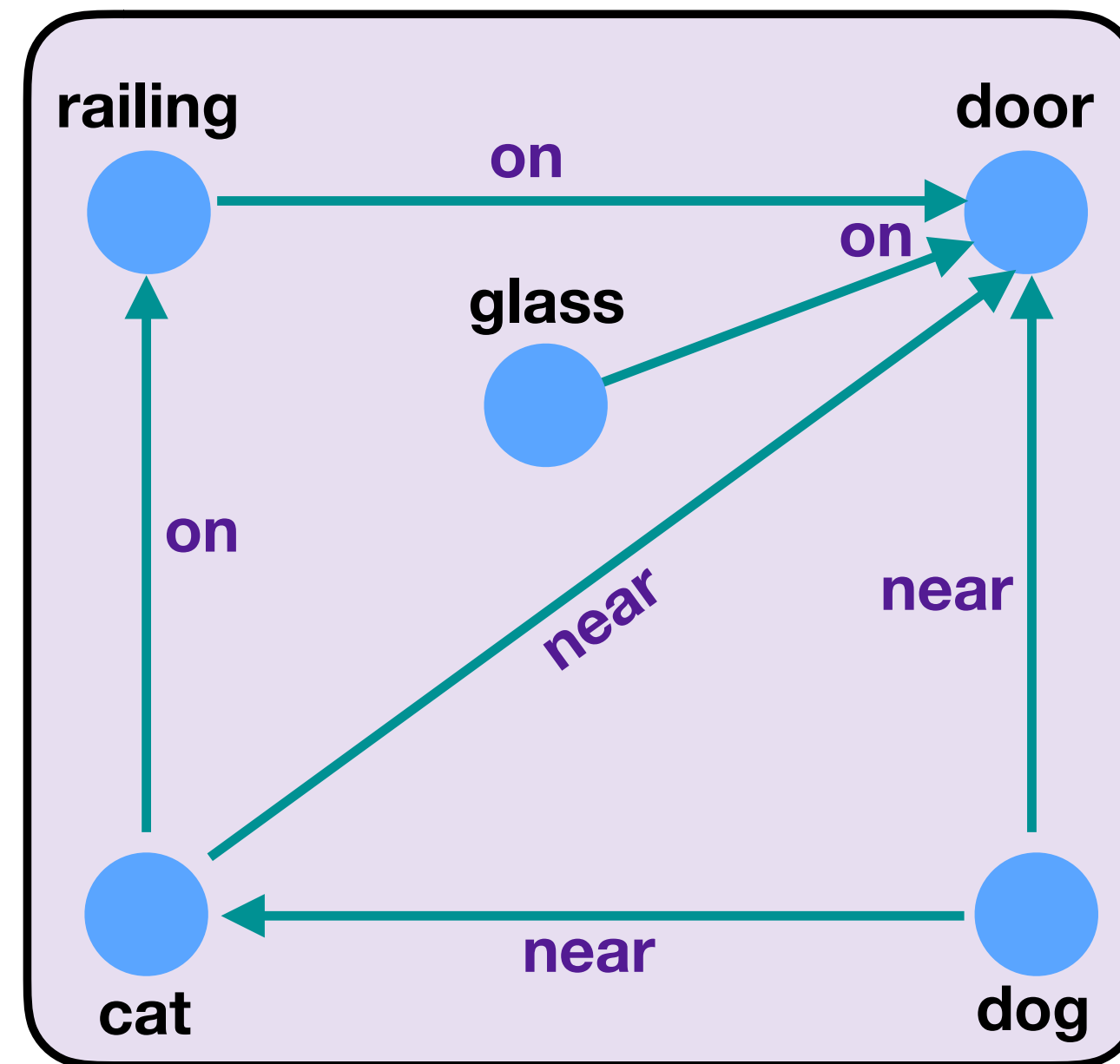
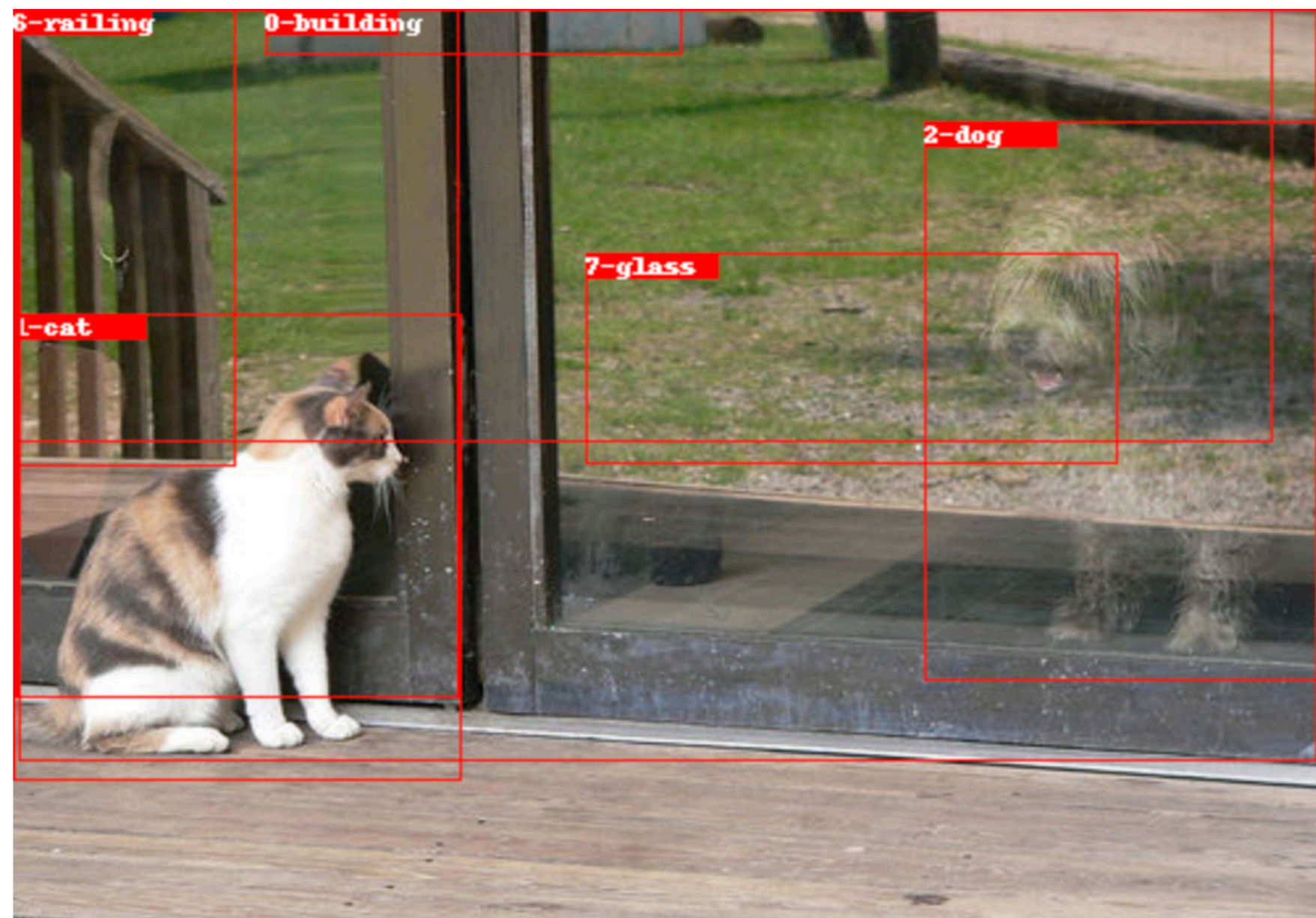
- **Today's** Lecture: Forsyth & Ponce (2nd ed.) 15
- **Next** Lecture: Forsyth & Ponce (2nd ed.) 16.1.3, 16.1.4, 16.1.9

## Reminders:

- **Assignment 5:** Scene Recognition with Bag of Words is **out**

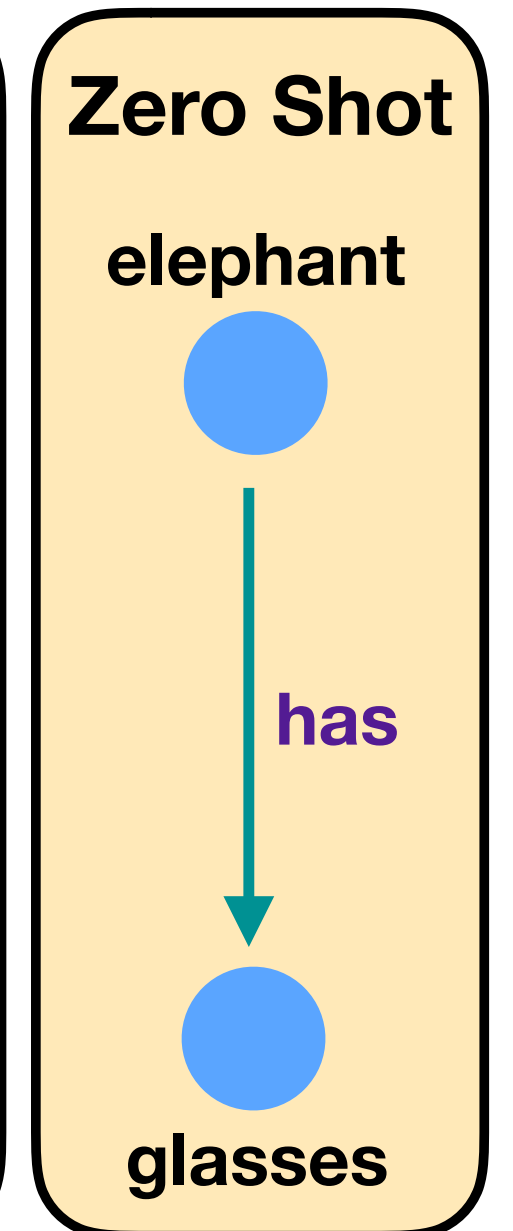
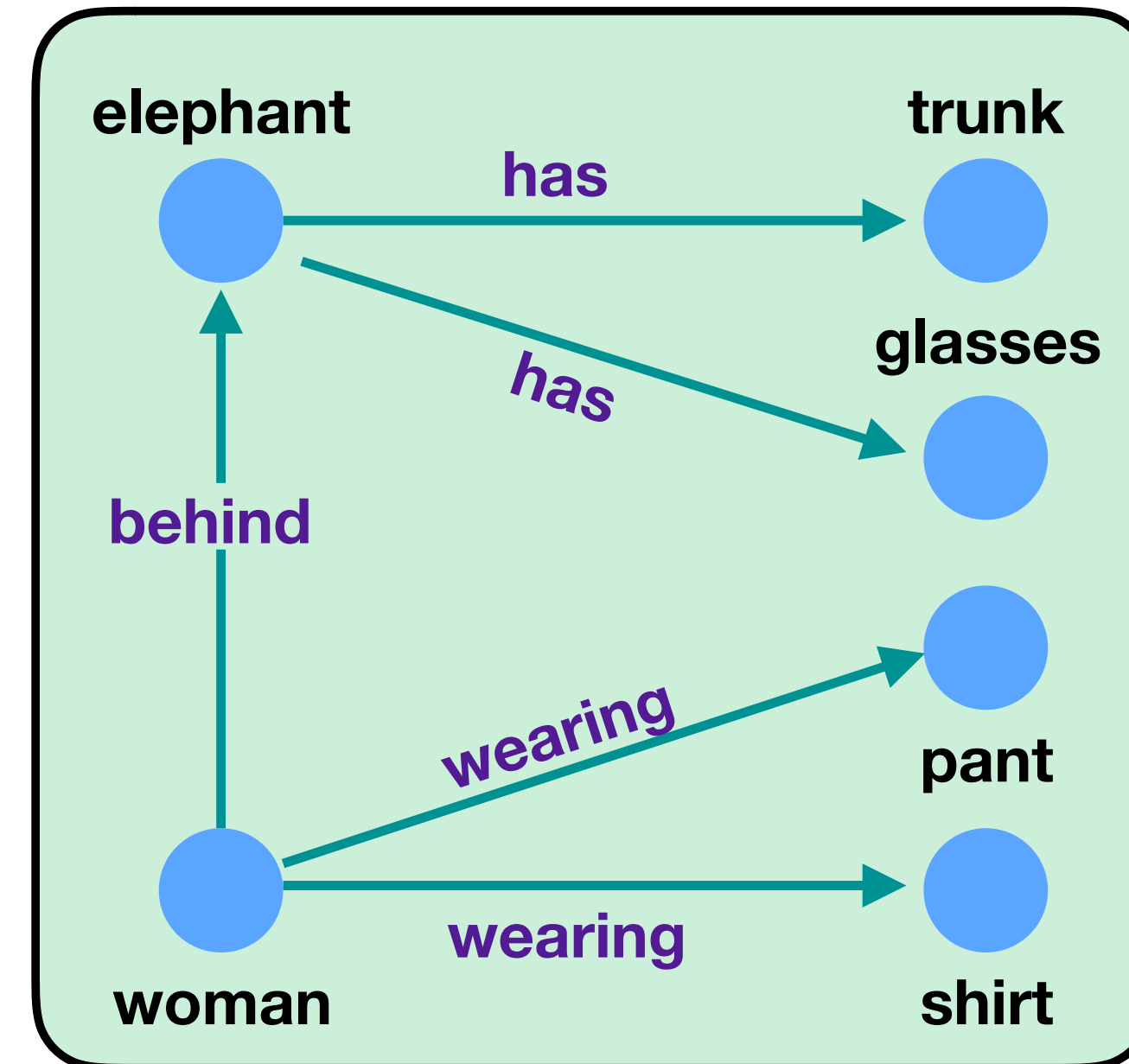
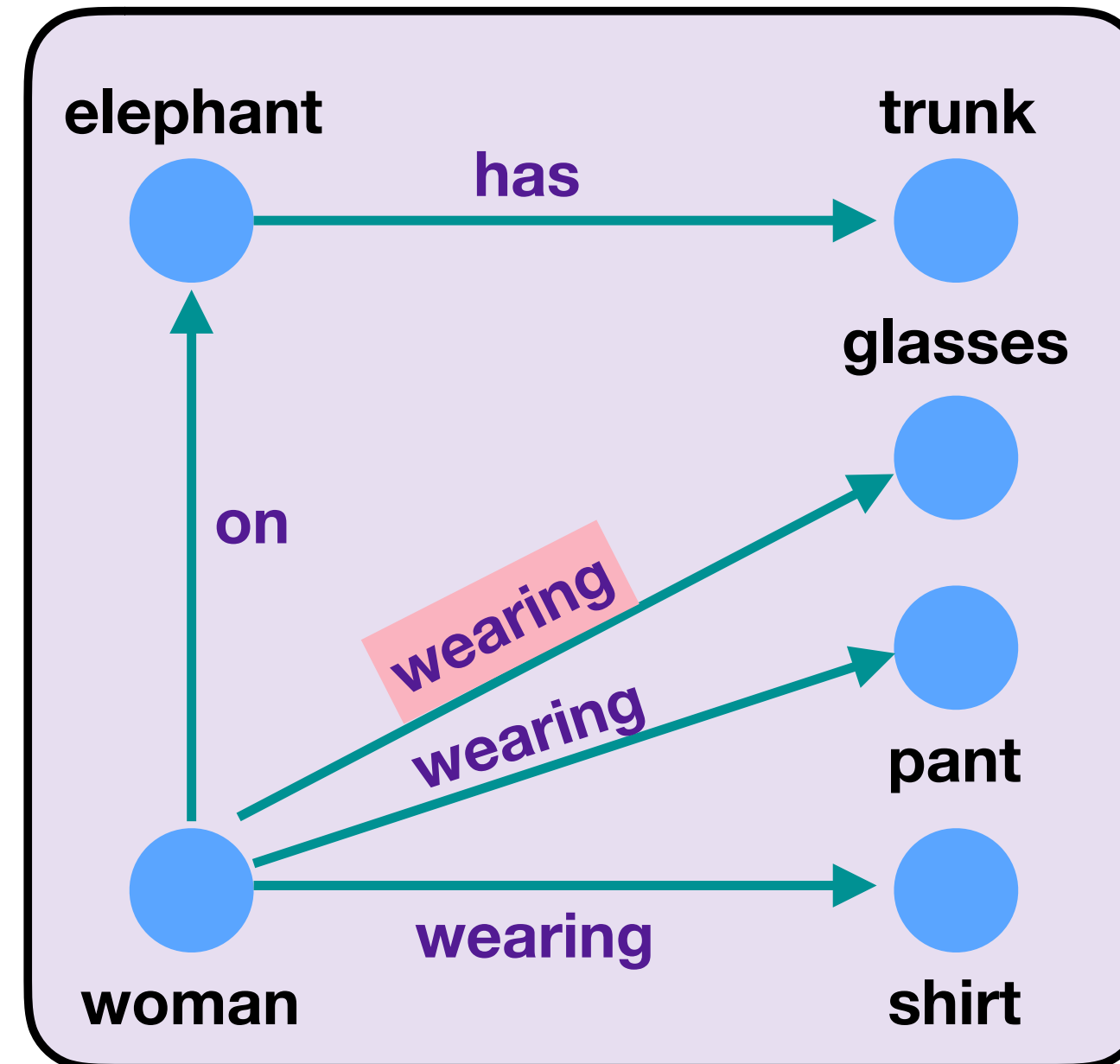
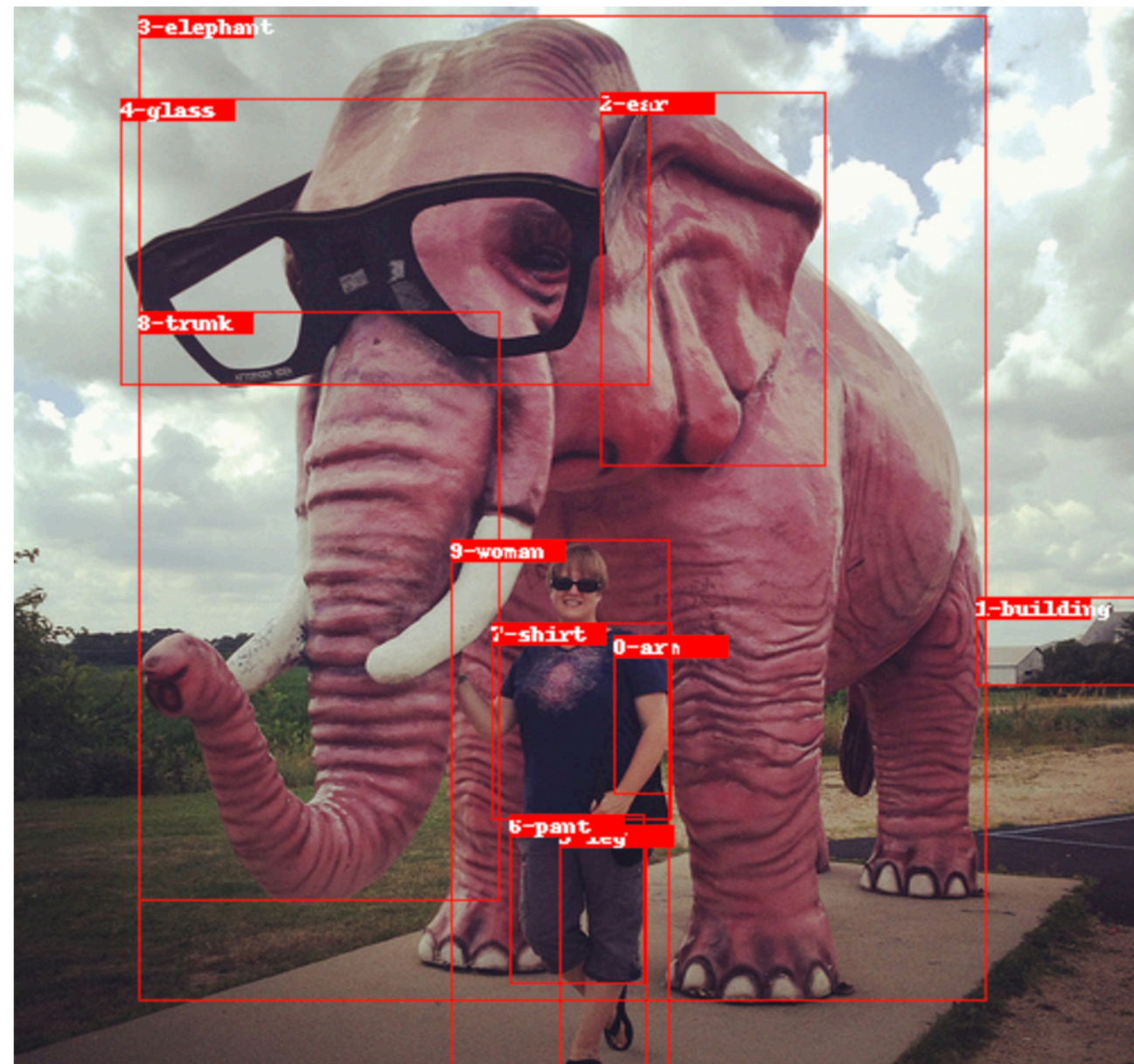


# Today's "fun" Example: Scene Graph Prediction



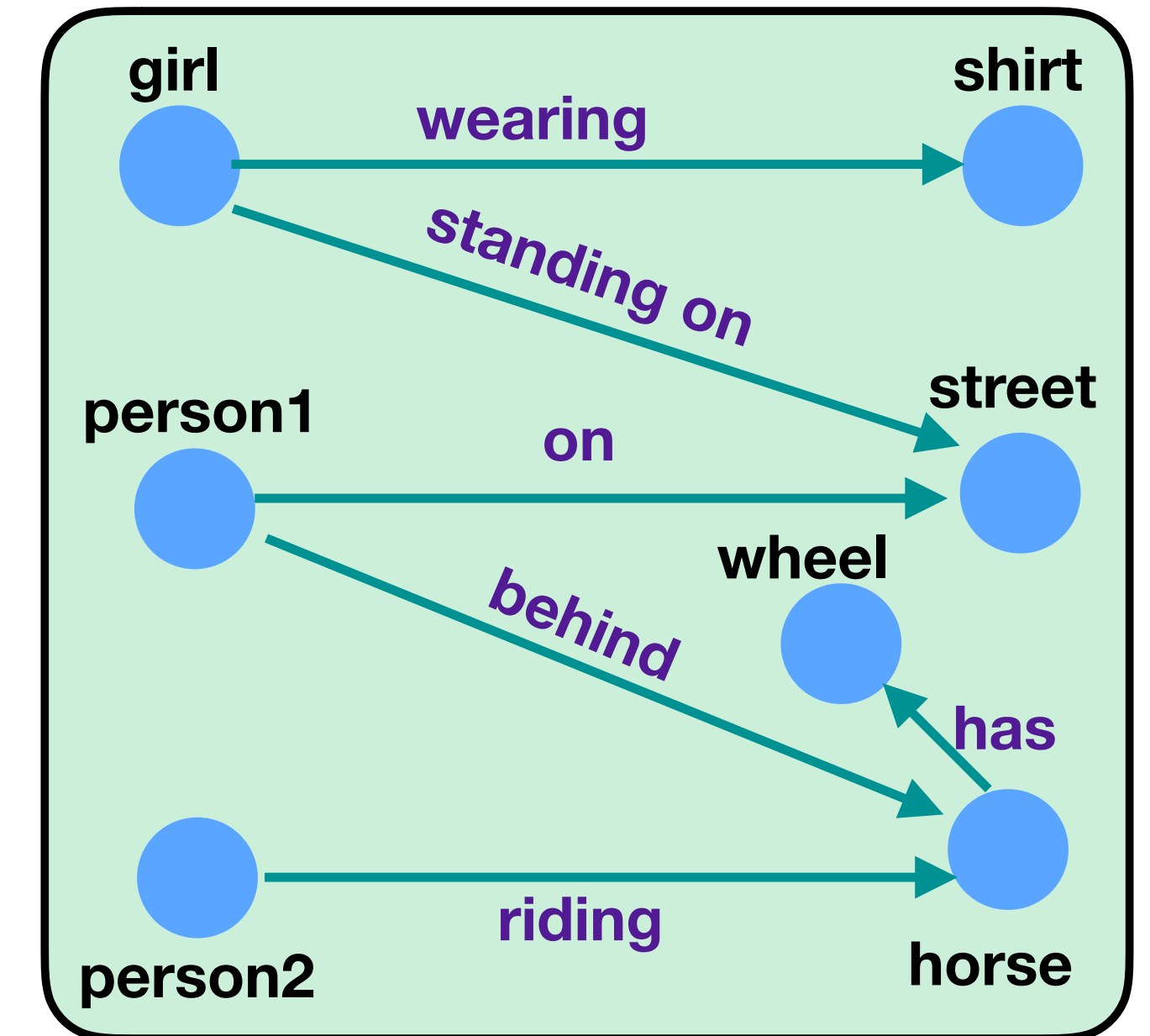
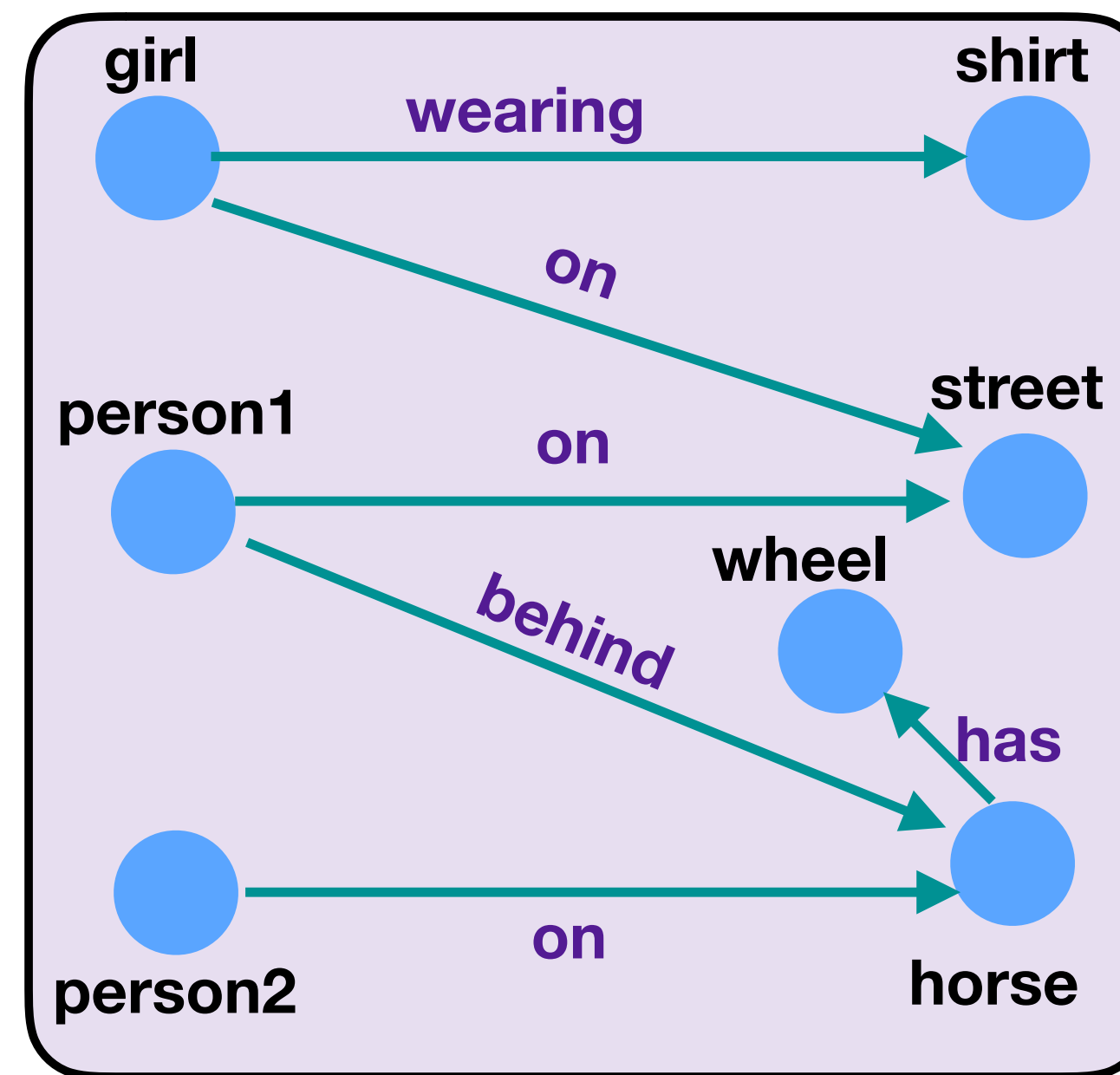
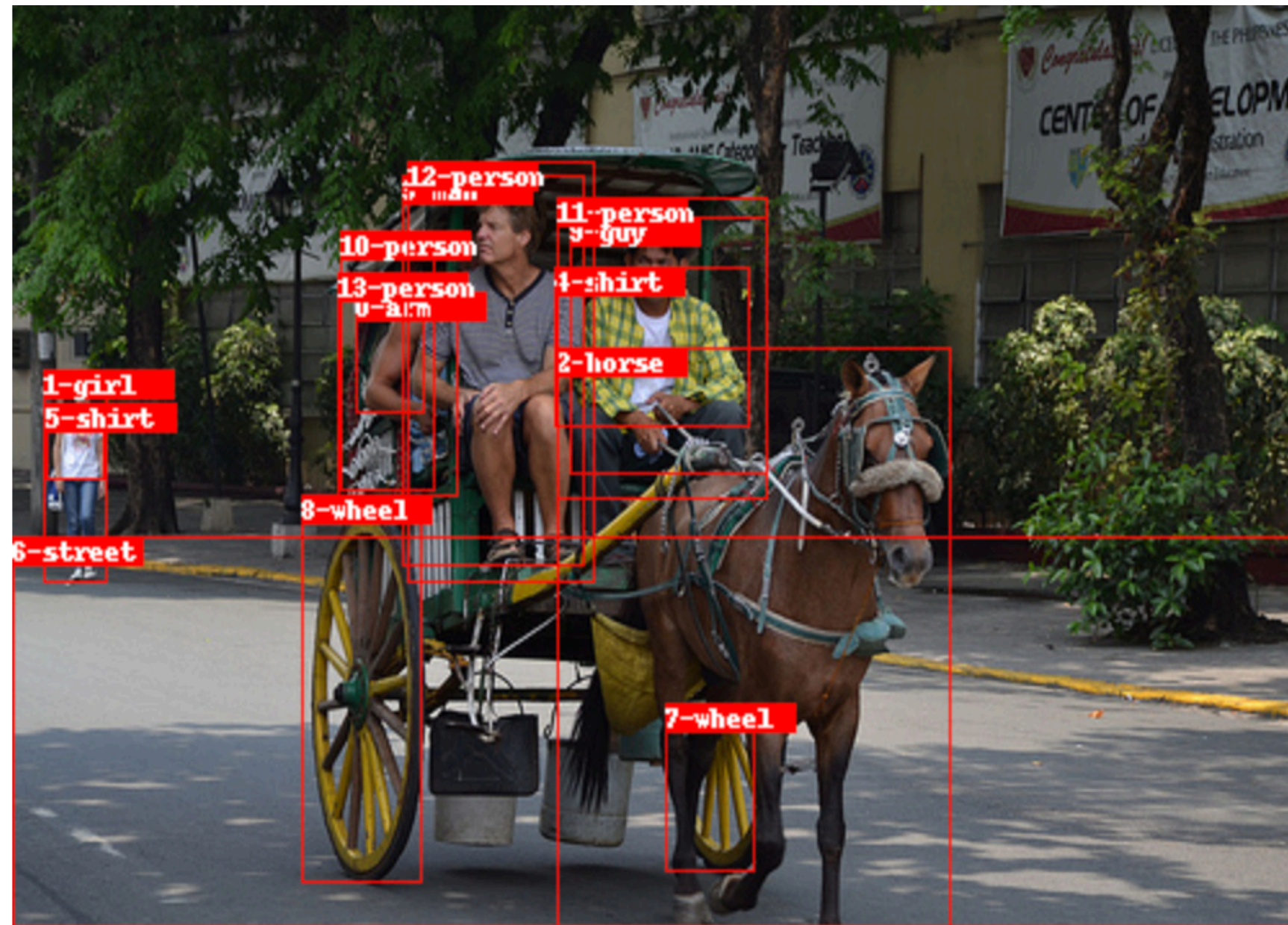


# Today's "fun" Example: Scene Graph Prediction





# Today's "fun" Example: Scene Graph Prediction





# Lecture 25: Classification



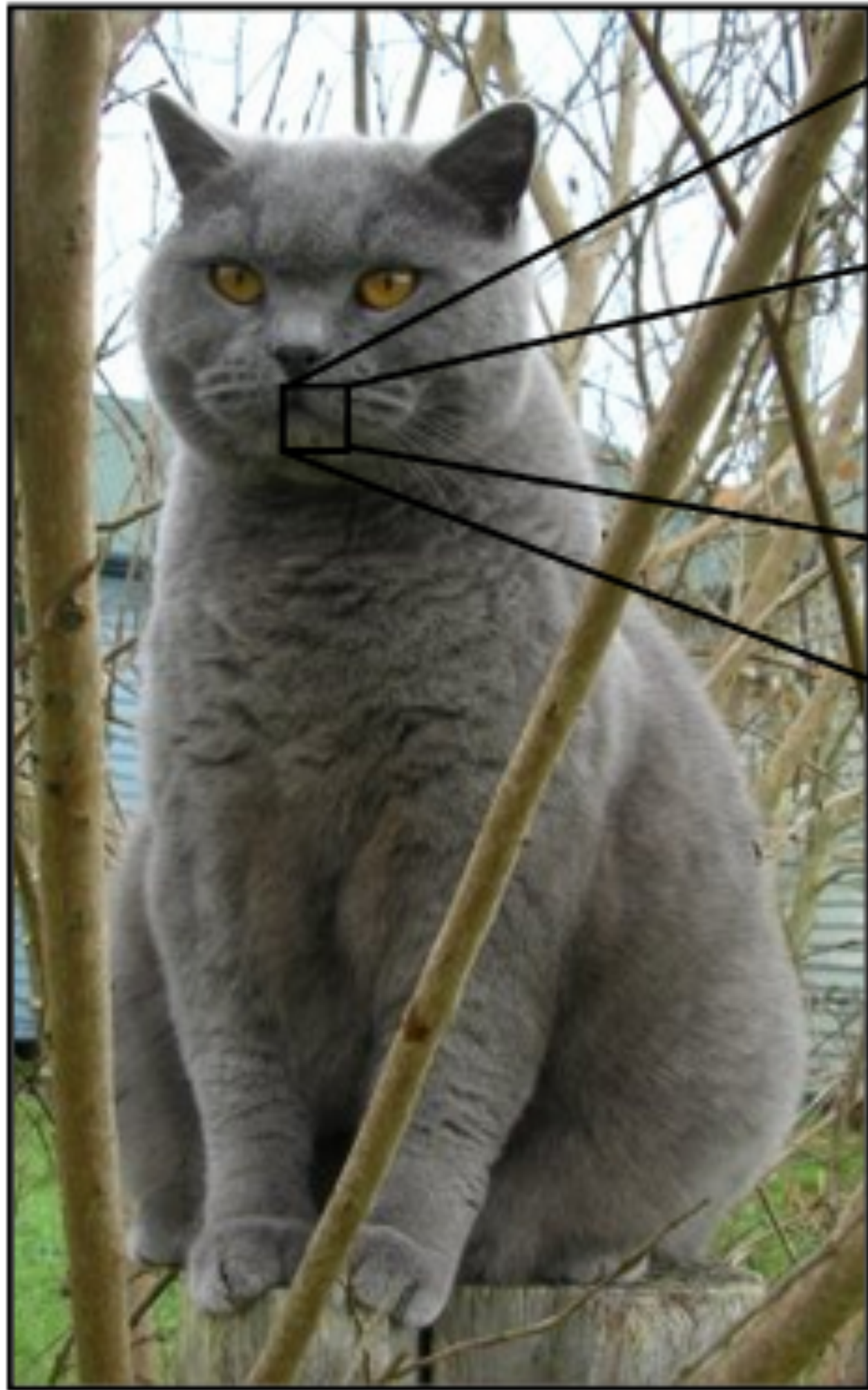
(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

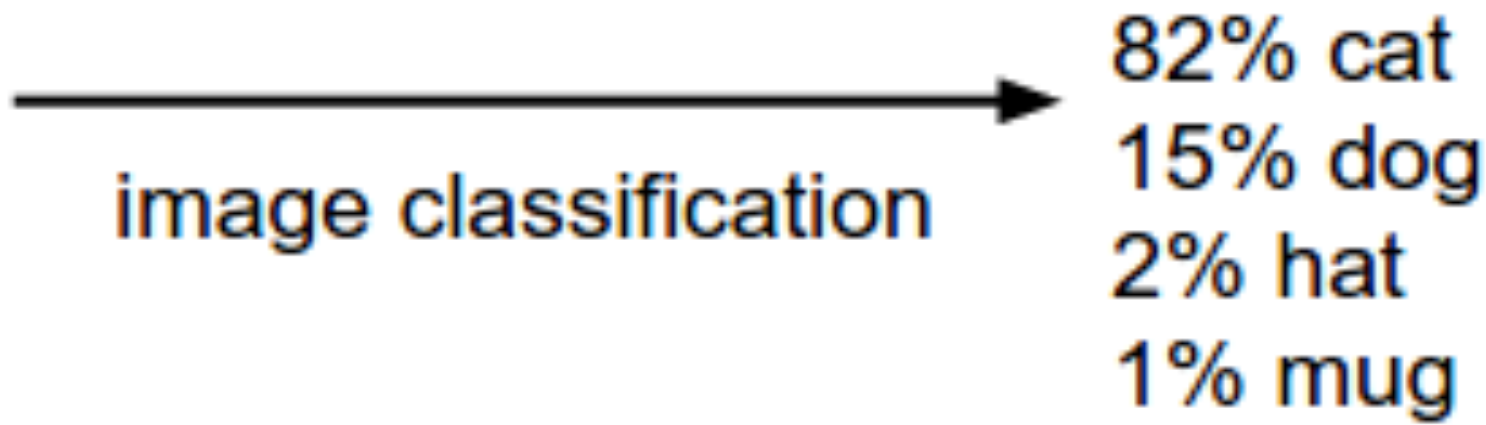


# Lecture 25: Classification



05	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91	88
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	52	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	69	21	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	03	59	41	92	36	54	22	40	40	28	66	33	13	80
24	47	39	80	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
59	36	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	35	35	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	63	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	86	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	13	67	48

What the computer sees



# Lecture 25: Bayes Classifier

Let  $c$  be the **class label** and let  $x$  be the **measurement** (i.e., evidence)

The diagram illustrates Bayes' theorem with color-coded components and labels:

- class-conditional probability (a.k.a. likelihood)**:  $P(x|c)$  (blue box)
- prior probability**:  $p(c)$  (green box)
- posterior probability**:  $P(c|x)$  (purple box)
- unconditional probability (a.k.a. marginal likelihood)**:  $P(x)$  (cyan box)

$$P(c|x) = \frac{P(x|c)p(c)}{P(x)}$$



# Lecture 25: Forms of Classifiers

Classification strategies fall under two broad types: **parametric** and **non-parametric**.

# Lecture 25: Forms of Classifiers

Classification strategies fall under two broad types: **parametric** and **non-parametric**.

Parametric classifiers are **model driven**. The parameters of the model are learned from training examples. New data points are classified by the learned model.

- fast, compact
- flexibility and accuracy depend on model assumptions



# Lecture 25: Forms of Classifiers

Classification strategies fall under two broad types: **parametric** and **non-parametric**.

Parametric classifiers are **model driven**. The parameters of the model are learned from training examples. New data points are classified by the learned model.

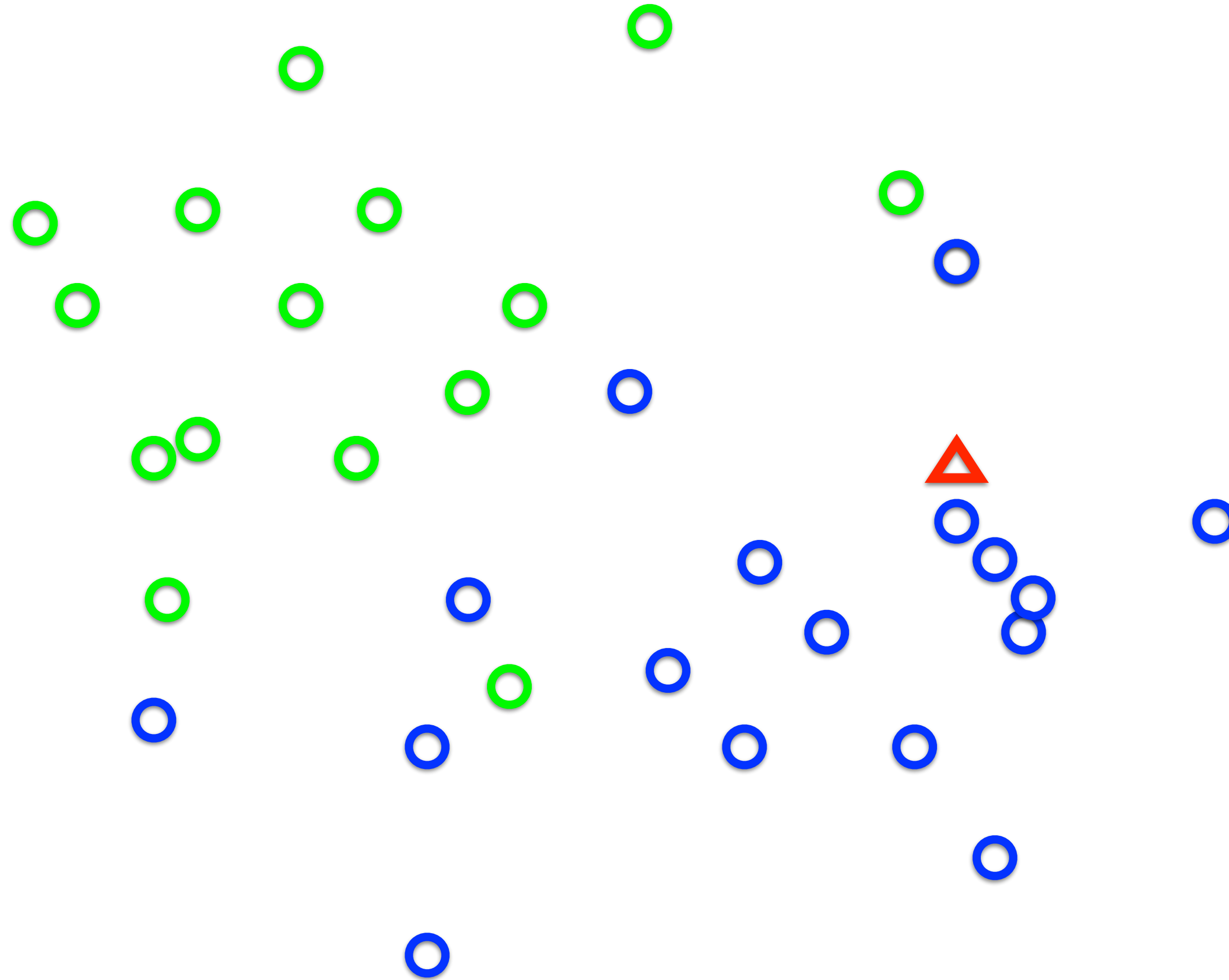
- fast, compact
- flexibility and accuracy depend on model assumptions

Non-parametric classifiers are **data driven**. New data points are classified by comparing to the training examples directly. "The data is the model".

- slow
- highly flexible decision boundaries

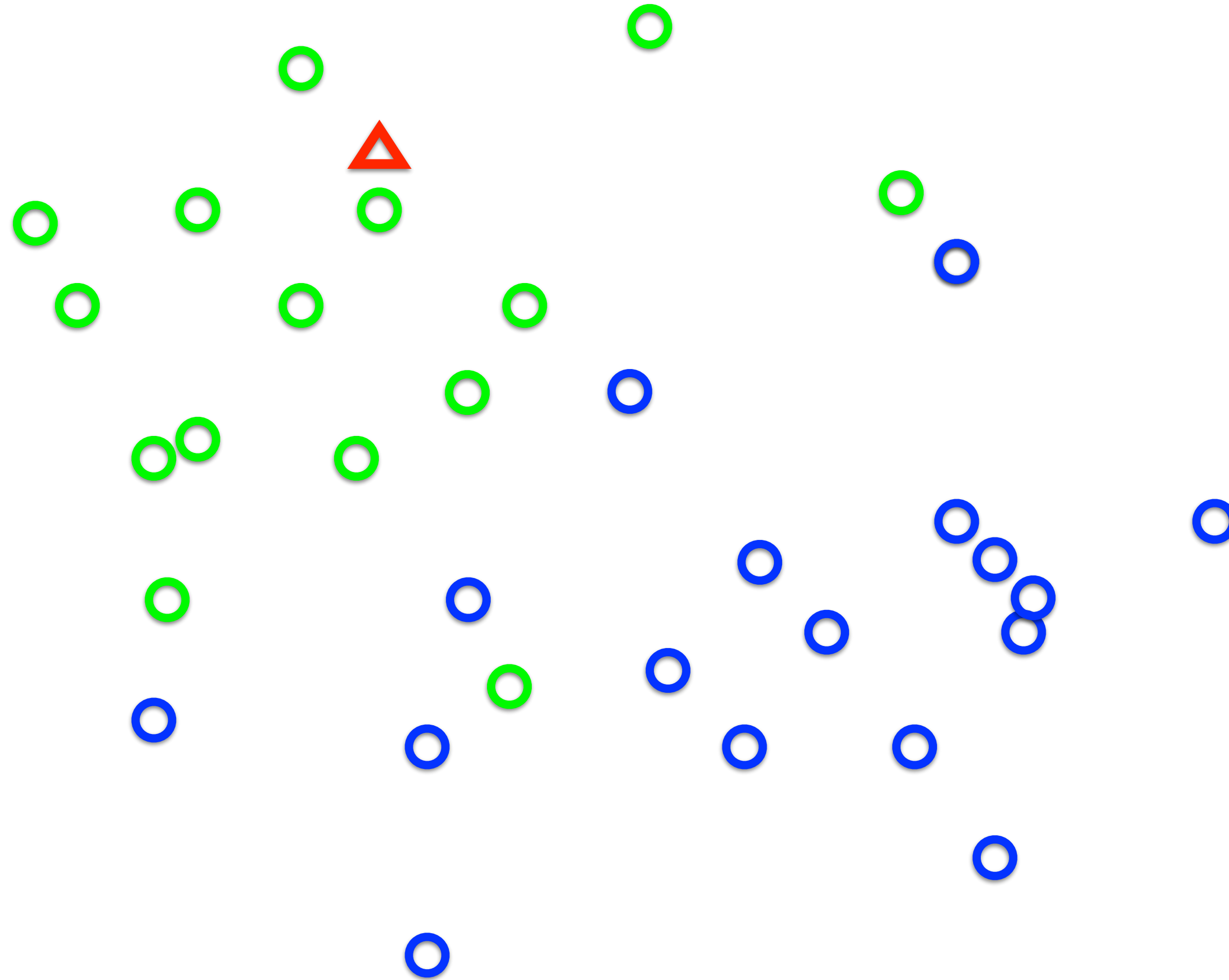
# Nearest Neighbor Classifier

Given a new data point, assign the label of nearest training example in feature space.



# Nearest Neighbor Classifier

Given a new data point, assign the label of nearest training example in feature space.



# k-Nearest Neighbor (kNN) Classifier

We can gain some robustness to noise by voting over **multiple** neighbours.

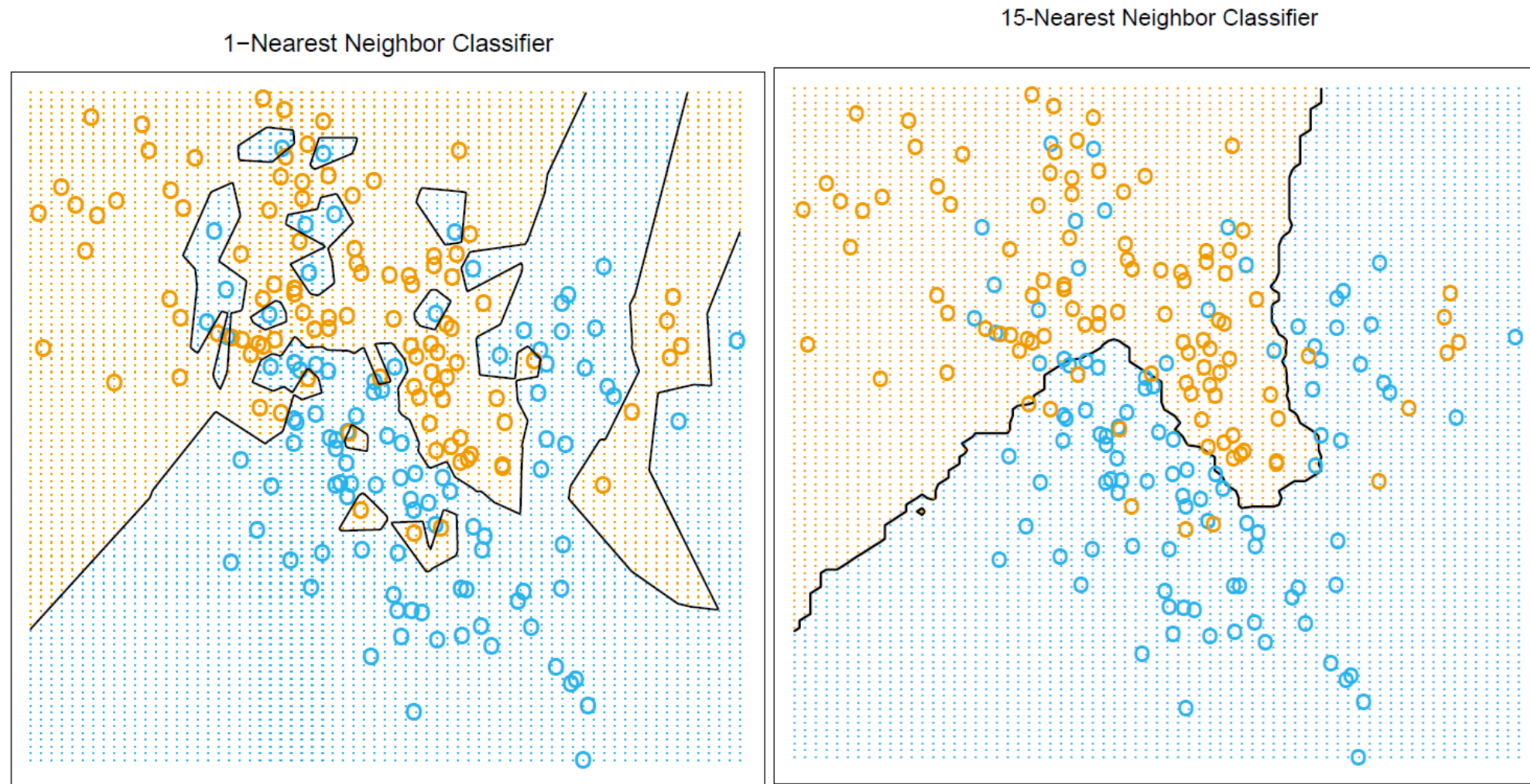
Given a **new** data point, find the k nearest training examples. Assign the label by **majority vote**.

Simple method that works well if the **distance measure** correctly weights the various dimensions

For **large data sets**, as k increases kNN approaches optimality in terms of minimizing probability of error



# k-Nearest Neighbor (kNN) Classifier



kNN decision boundaries respond to local clusters where one class dominates

**Figure credit:** Hastie, Tibshirani & Friedman (2nd ed.)



# Classifier Strategies

Classification strategies fall under two broad types: **parametric** and **non-parametric**.

Parametric classifiers are **model driven**. The parameters of the model are learned from training examples. New data points are classified by the learned model.

- fast, compact
- flexibility and accuracy depend on model assumptions

Non-parametric classifiers are **data driven**. New data points are classified by comparing to the training examples directly. "The data is the model".

- slow
- highly flexible decision boundaries

# Support Vector Machines (SVM)

**Idea:** Try to obtain the decision boundary directly

The decision boundary is parameterized as a **separating hyperplane** in feature space.

— e.g. a separating line in 2D

We choose the hyperplane that is as far as possible from each class - that maximizes the distance to the closest point from either class.

# Linear Classifier

Defines a score function:

$$f(\mathbf{x}_i, \mathbf{W}, \mathbf{b}) = \mathbf{W}\mathbf{x}_i + \mathbf{b}$$

image features

weights

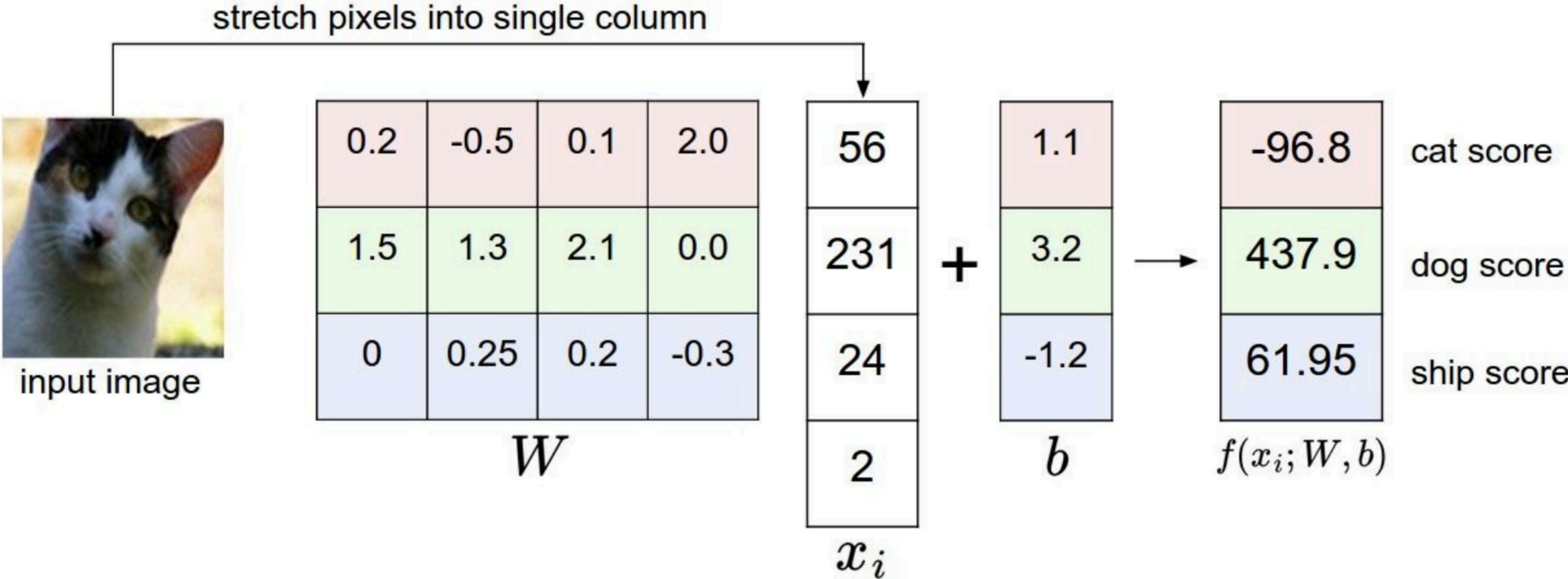
(parameters)

bias vector

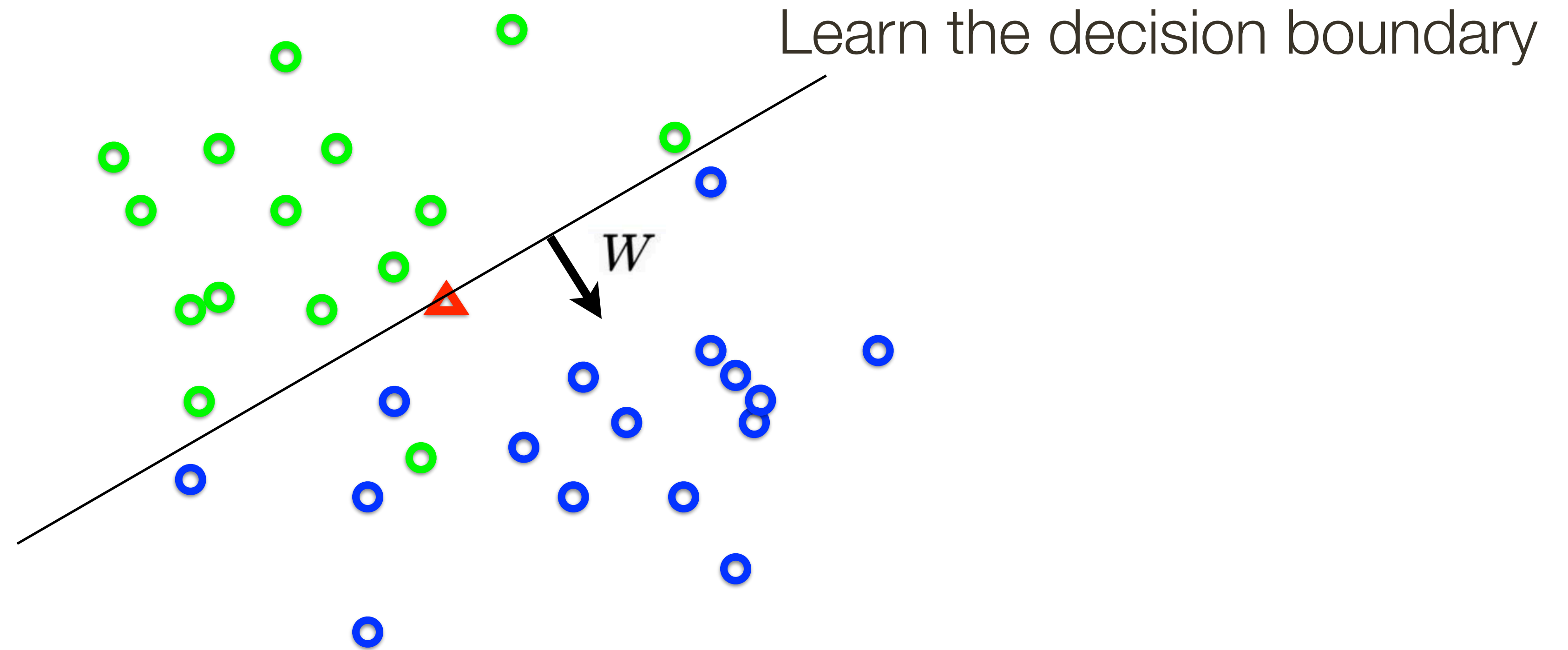


# Linear Classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

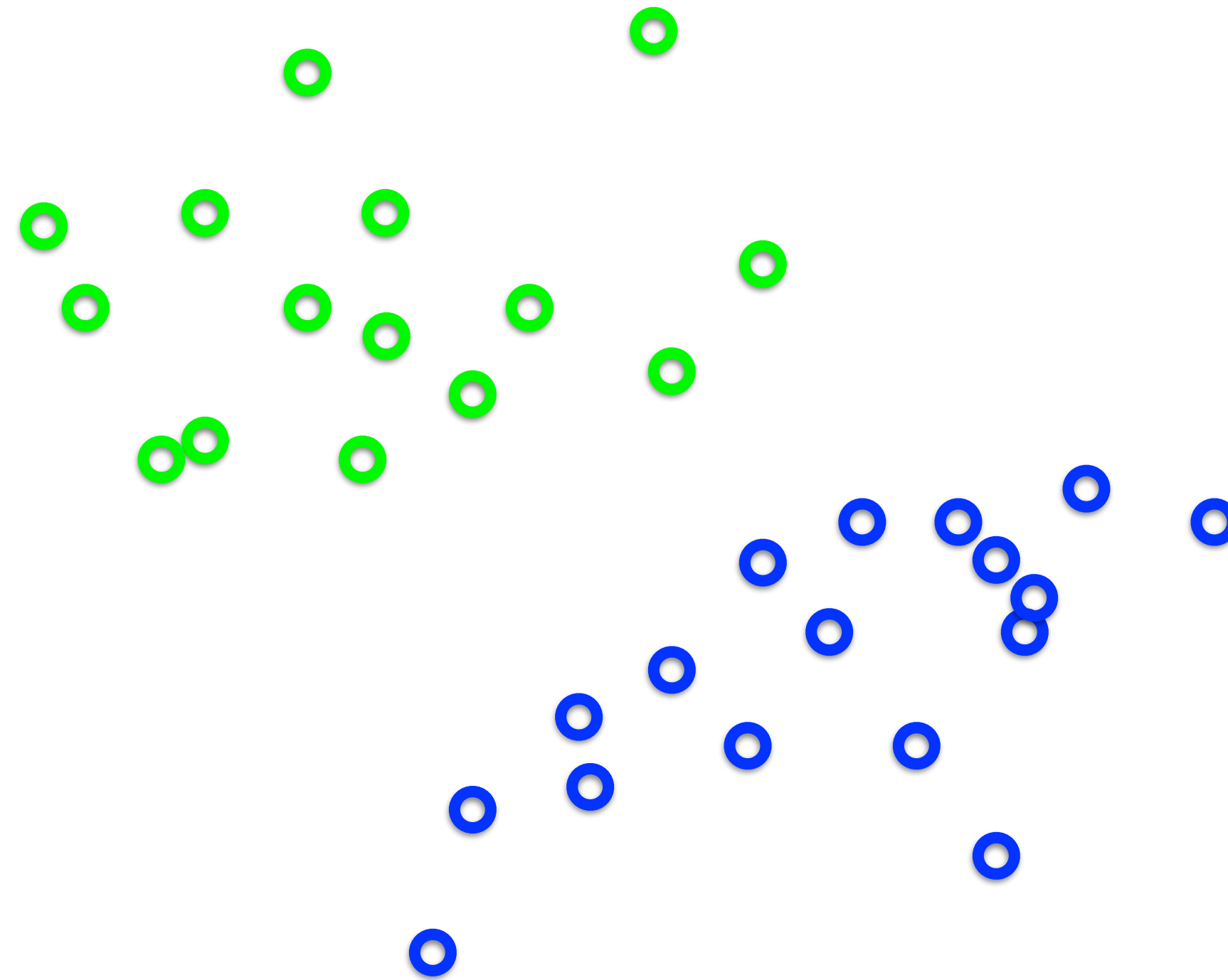


# Support Vector Machines (SVM)



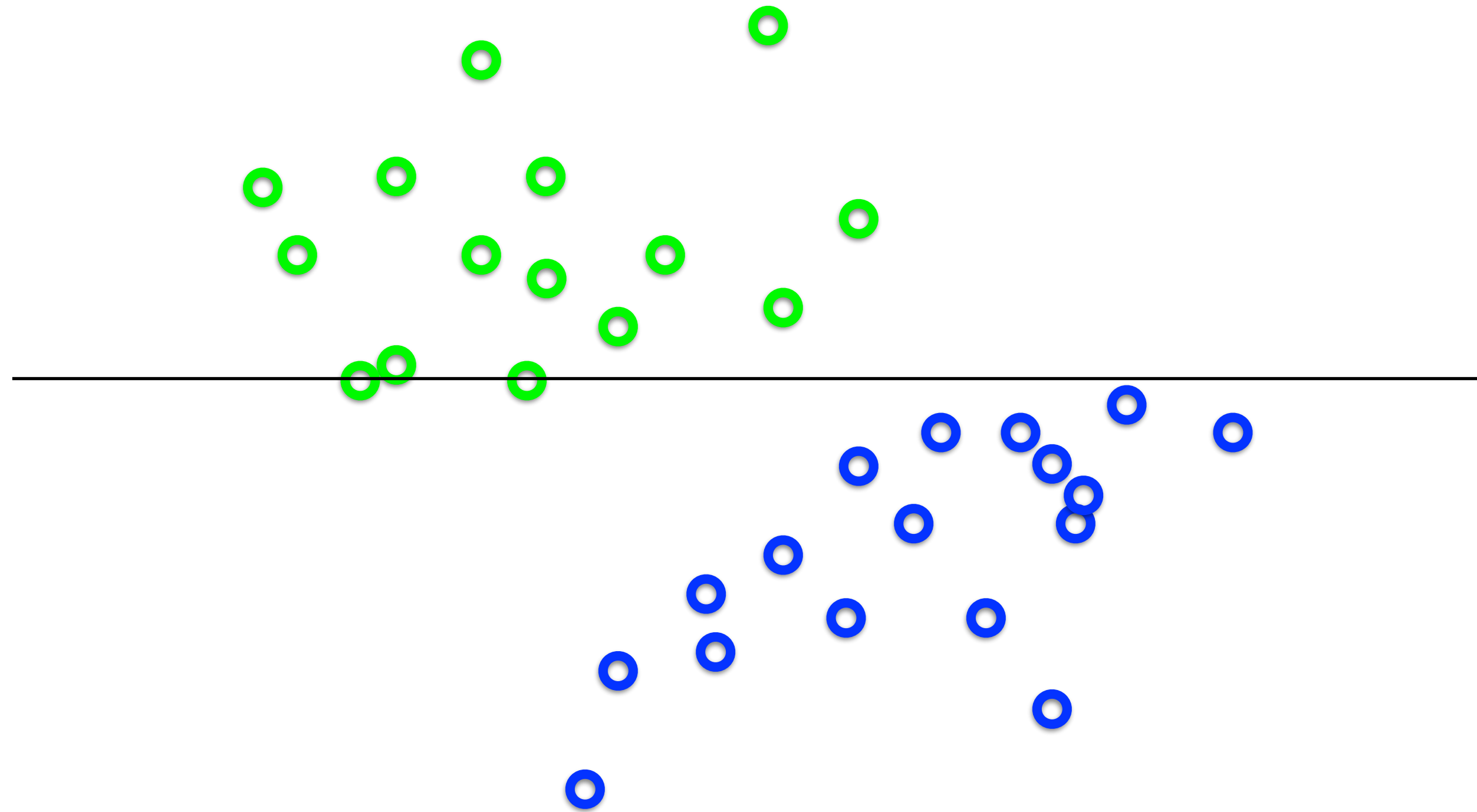
# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?



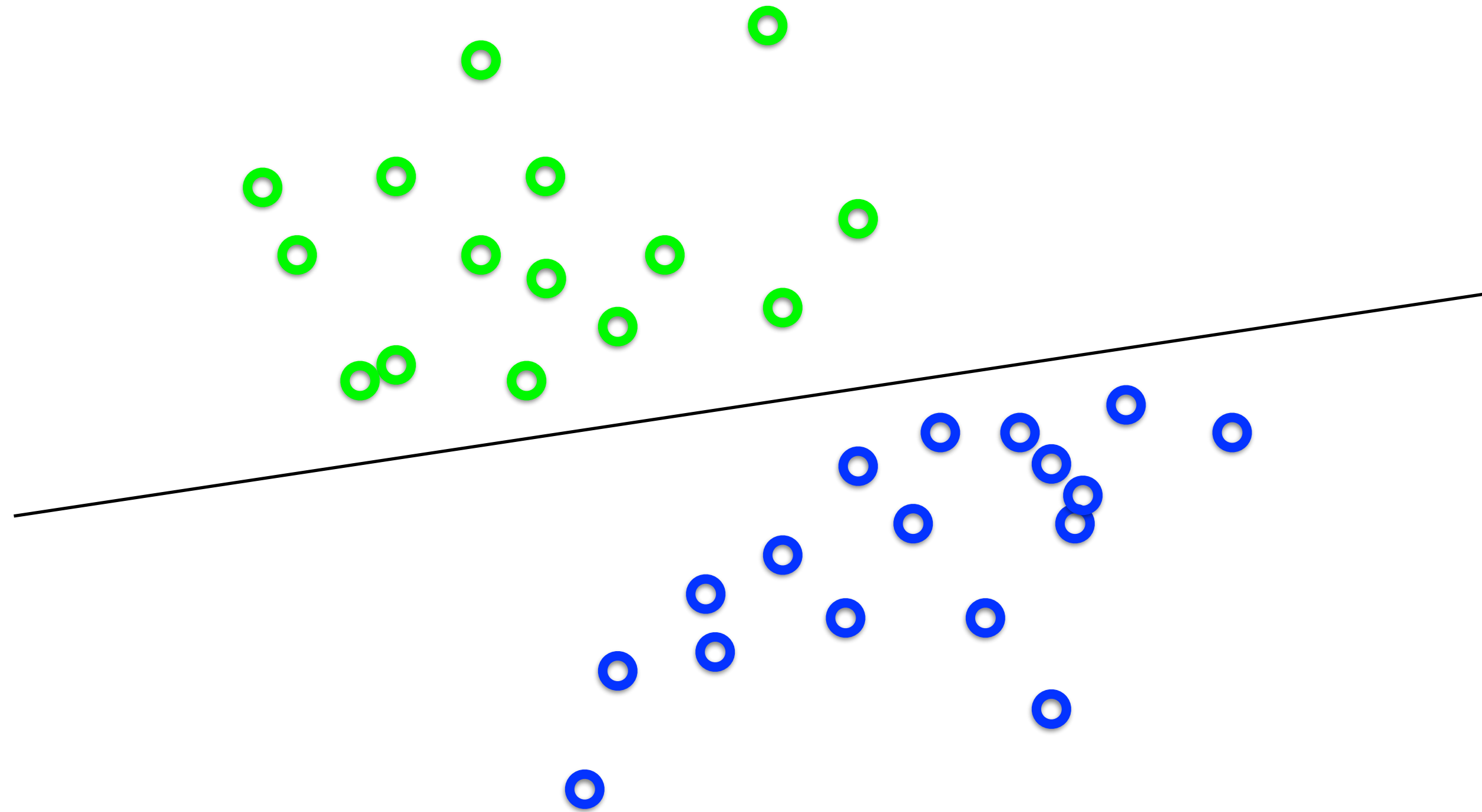
# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?



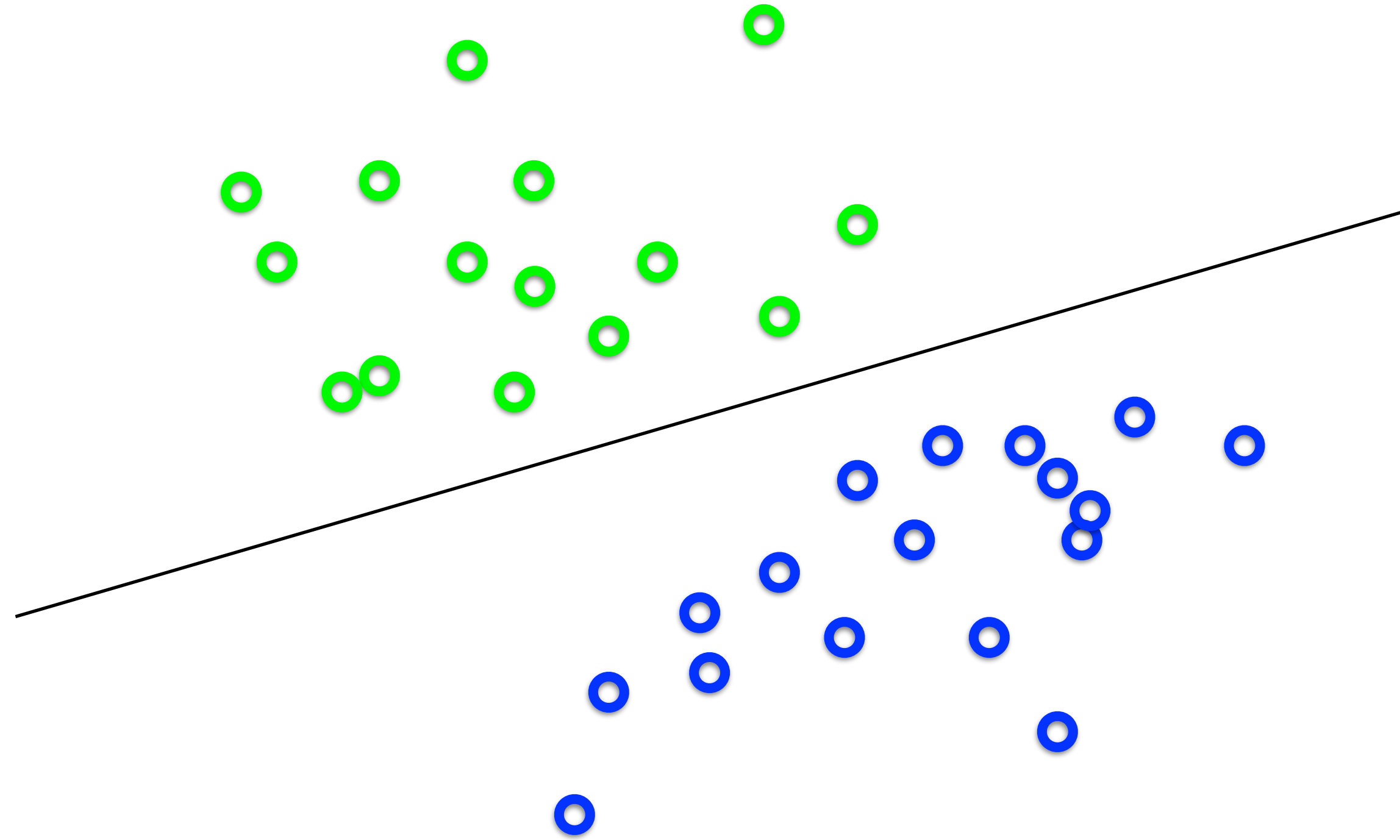
# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?



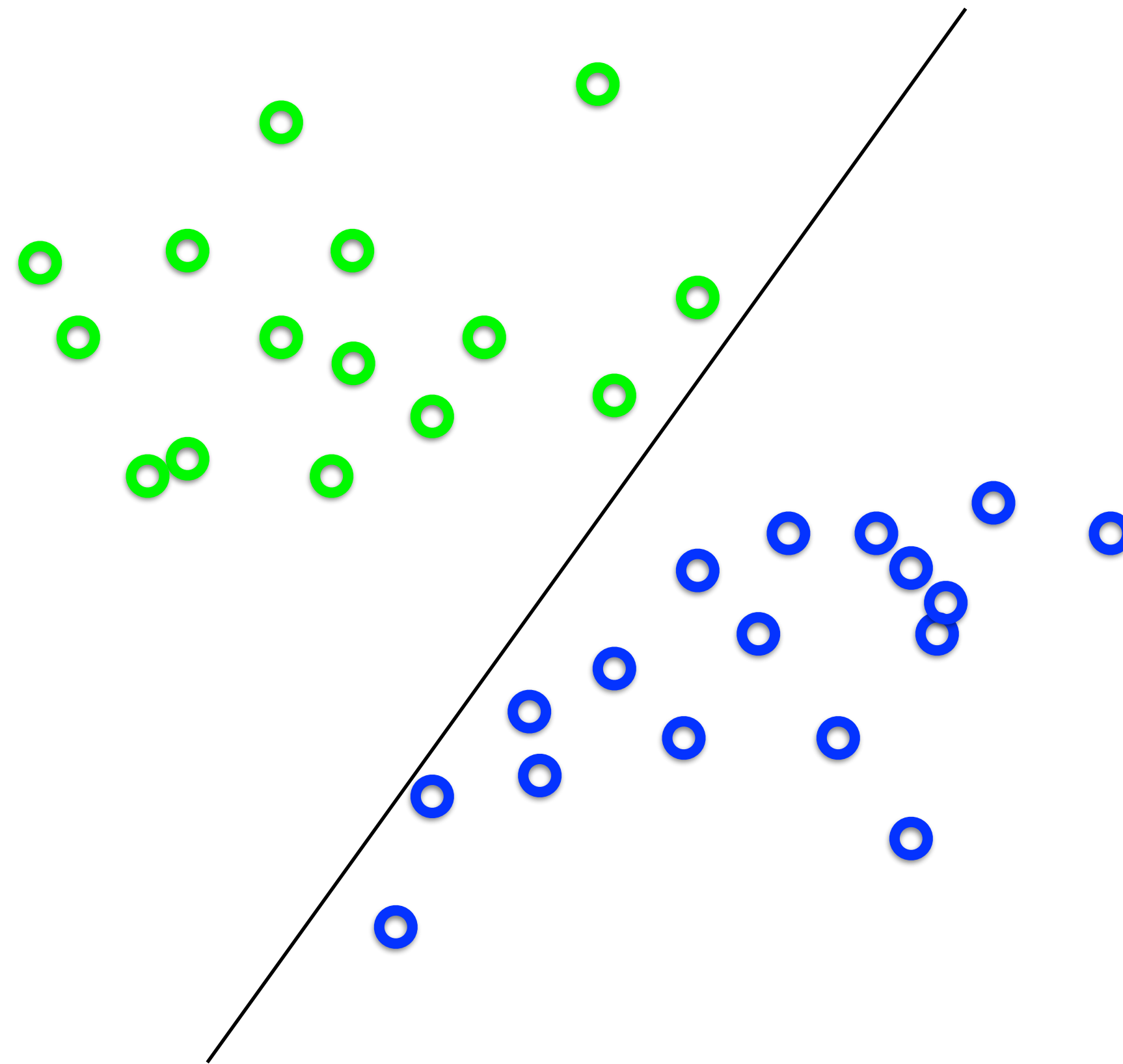
# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?



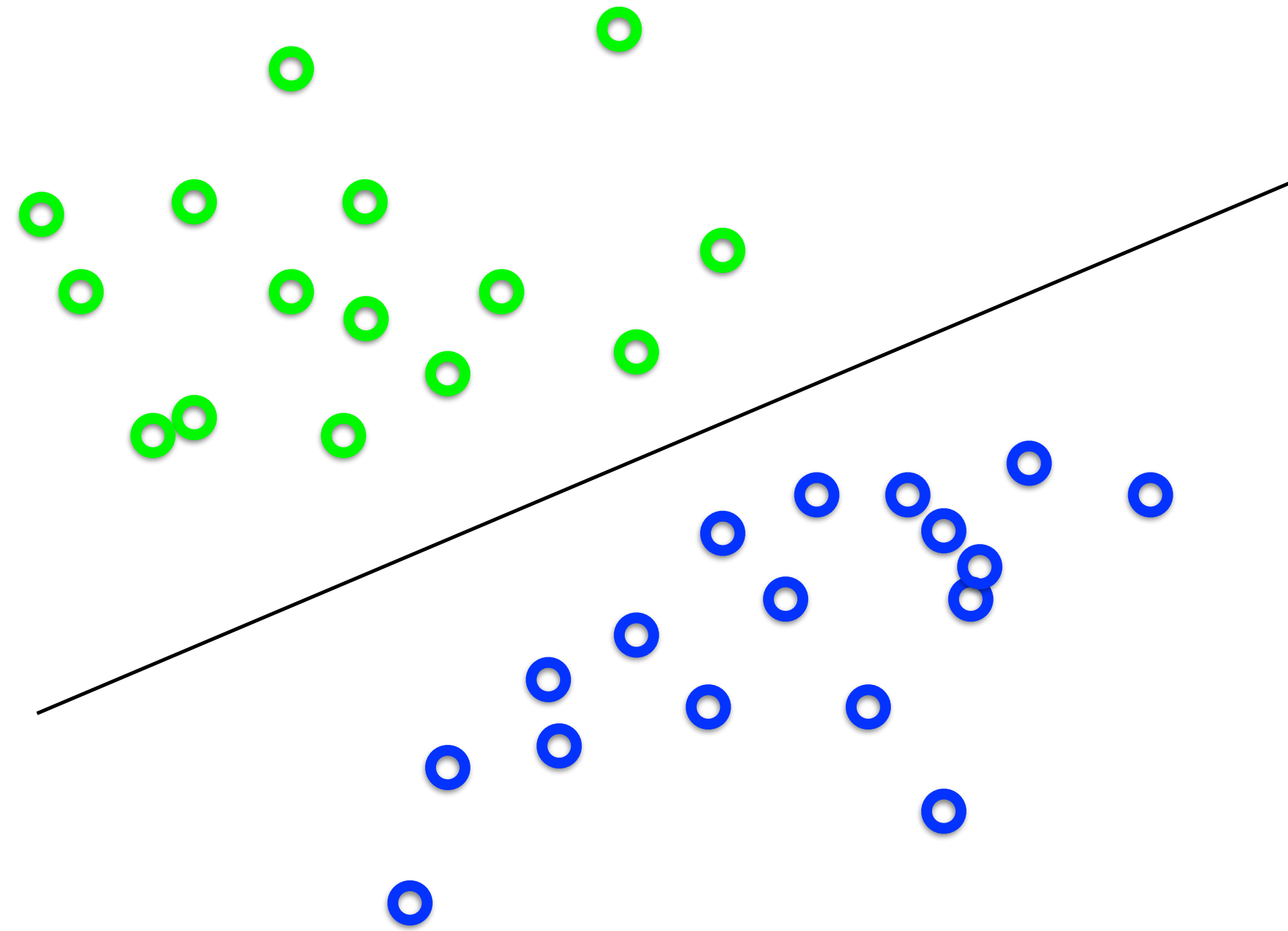
# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?



# Support Vector Machines (SVM)

What's the best  $\mathbf{w}$  ?

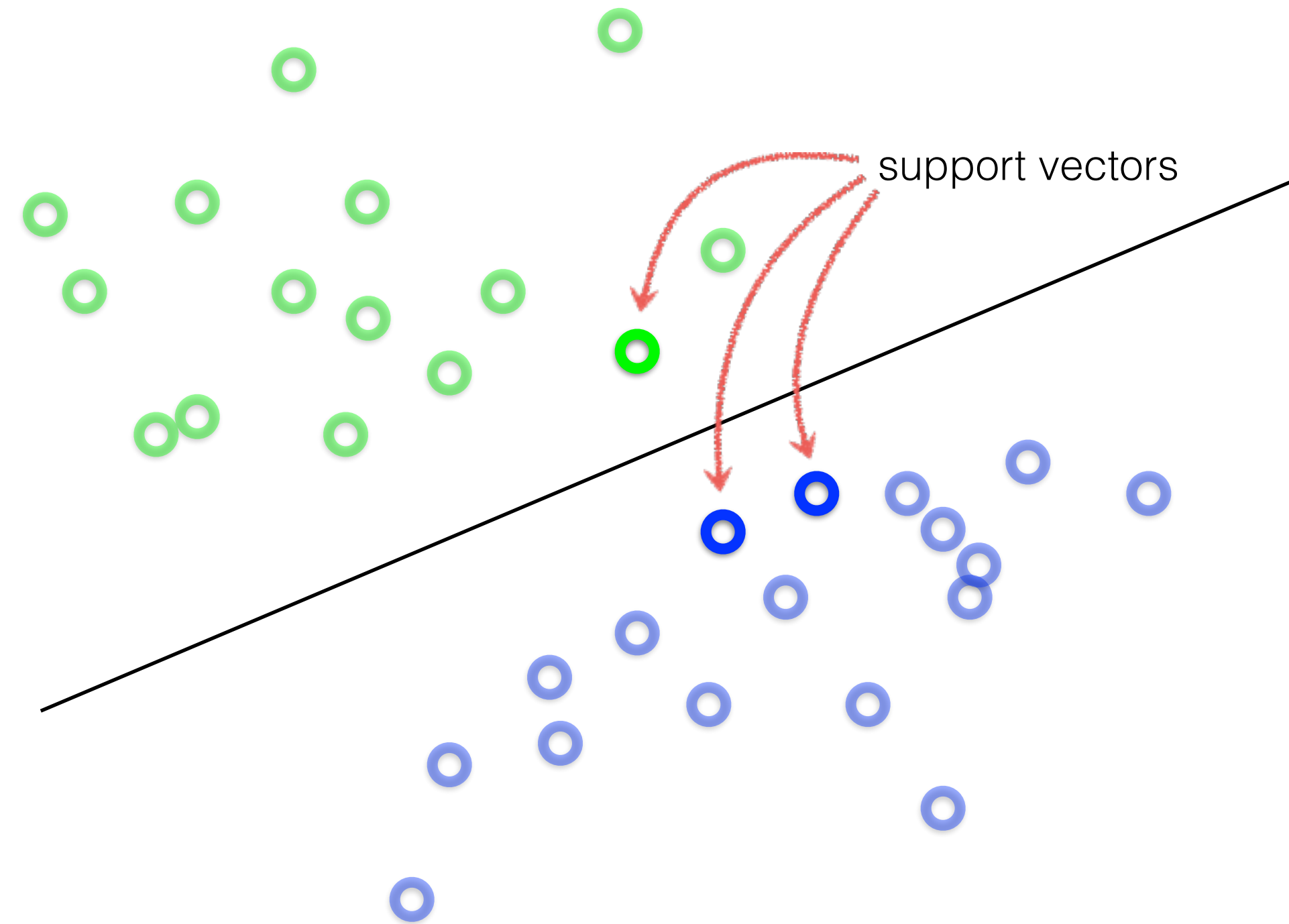


**Intuitively**, the line that is the farthest from all interior points



# Support Vector Machines (SVM)

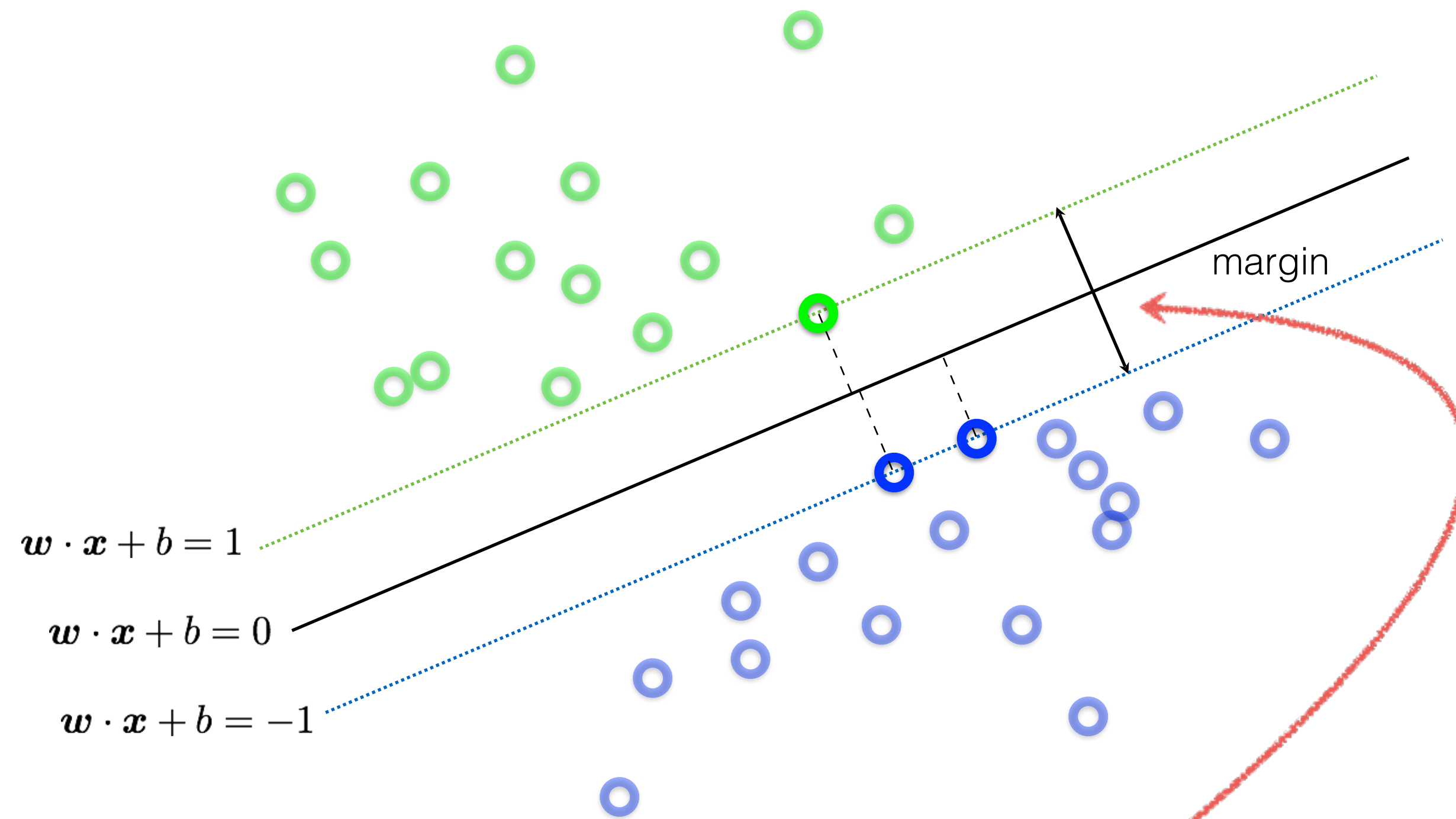
What's the best  $\mathbf{w}$  ?



Want a hyperplane that is far away from 'inner points'

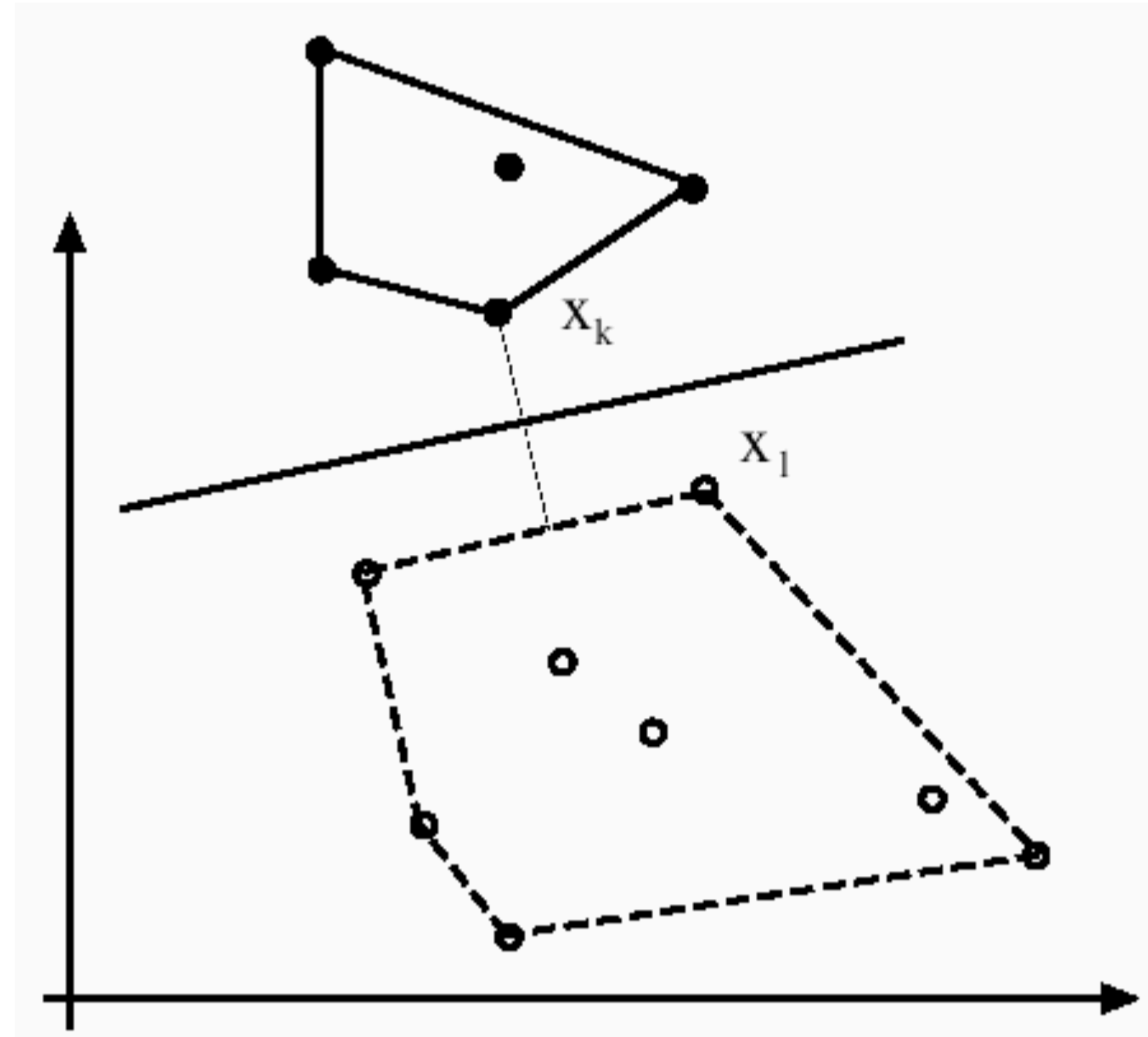
# Support Vector Machines (SVM)

Find hyperplane  $\mathbf{w}$  such that ...



the gap between parallel hyperplanes  $\frac{2}{\|\mathbf{w}\|}$  is maximized

# Support Vector Machines (SVM)



Forsyth & Ponce (2nd ed.) Figure 15.6



# Example: Pedestrian Detection with SVM

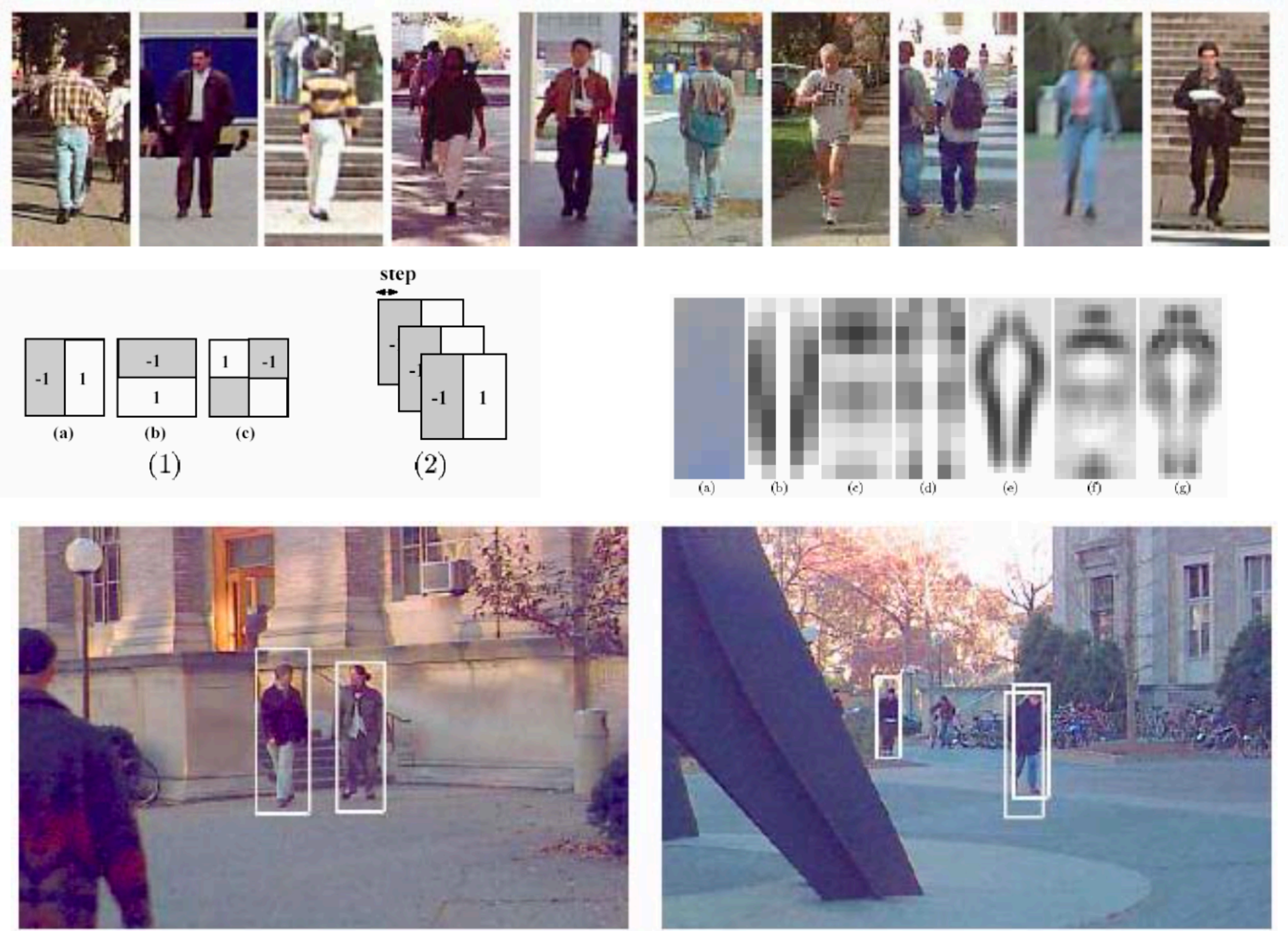


Figure credit: Papageorgiou, Oren, and Poggio, 1998



# Summary

A **classifier** accepts as input a set of features and outputs (predicts) a class label

Classifiers need to take into account “**loss**” associated with each kind of classification error

A Receiver Operating Characteristic (**ROC**) curve plots the trade-off between false negatives and false positives

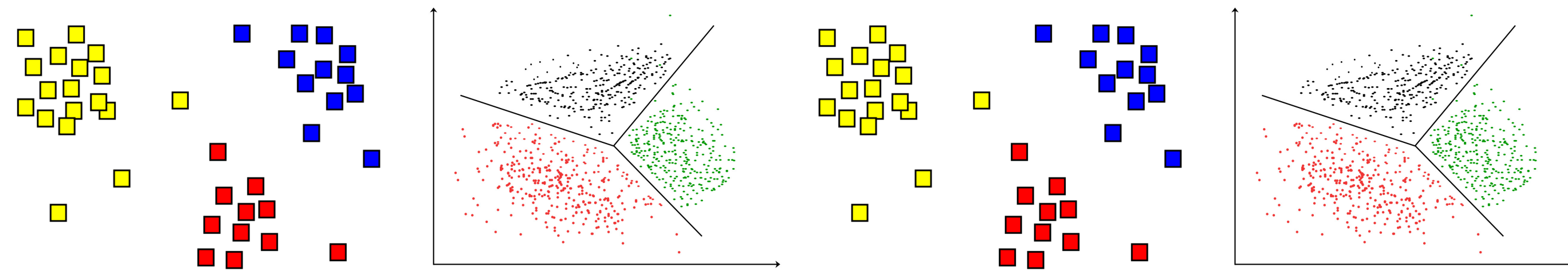
**Parametric** classifiers are model driven. The parameters of the model are learned from training examples

— e.g. support vector machine, decision tree

**Non-parametric** classifiers are data driven. New data points are classified by comparing to the training examples directly

— e.g. k-nearest neighbour

# CPSC 425: Computer Vision



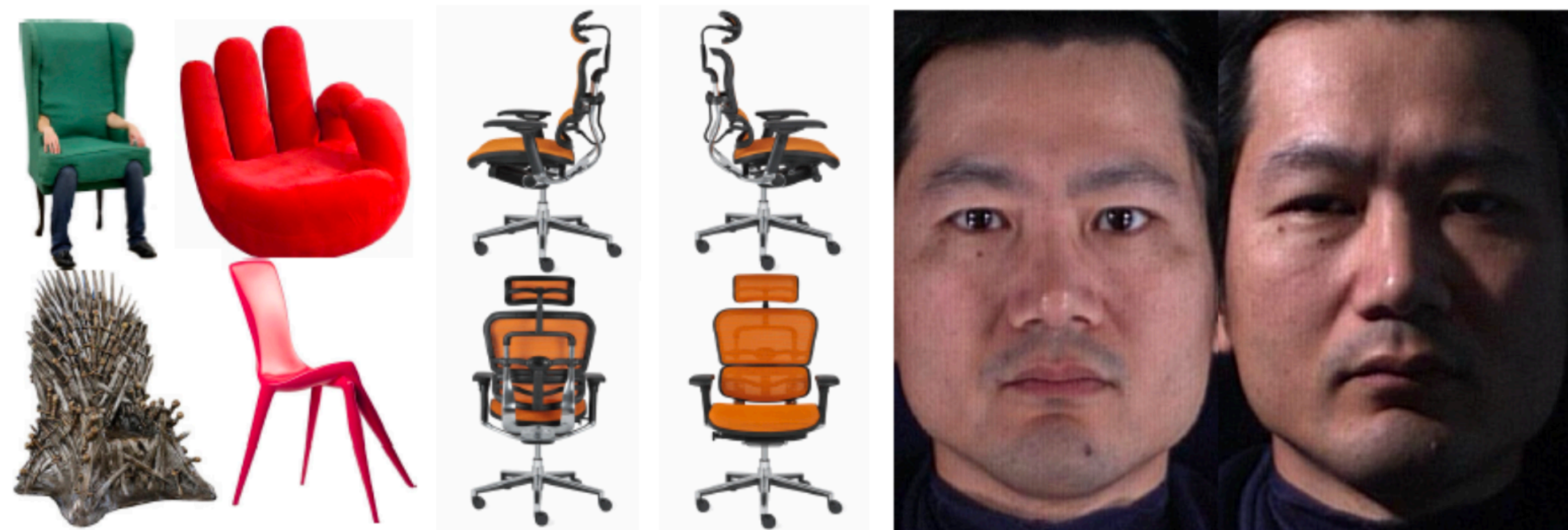
## Lecture 25: Image Classification

# 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?

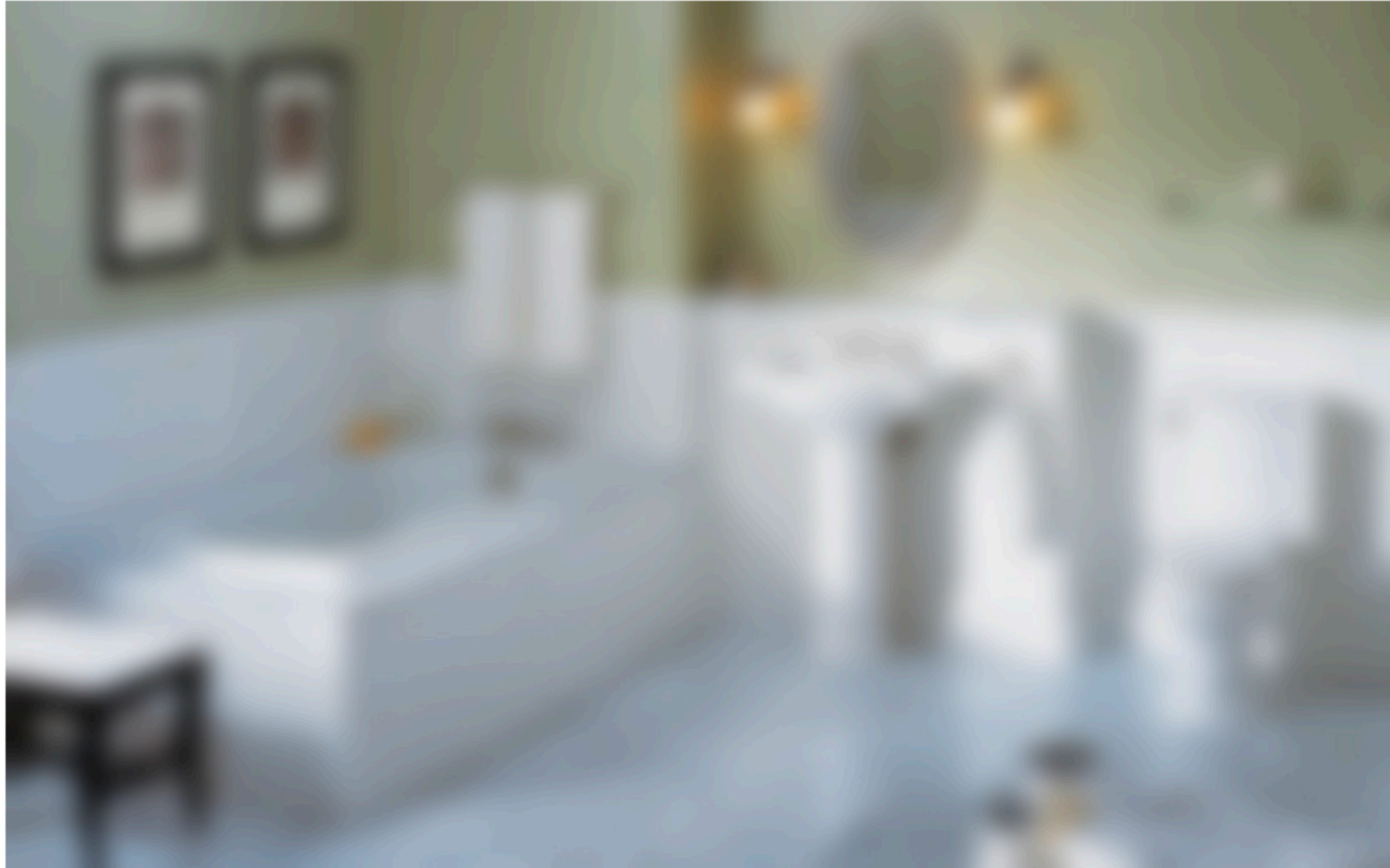
Visual perception is influenced by expectation. Our expectations are often conditioned on the **context**.

# What is This **Object**?





# What is This **Object**?





# What is This **Object**?

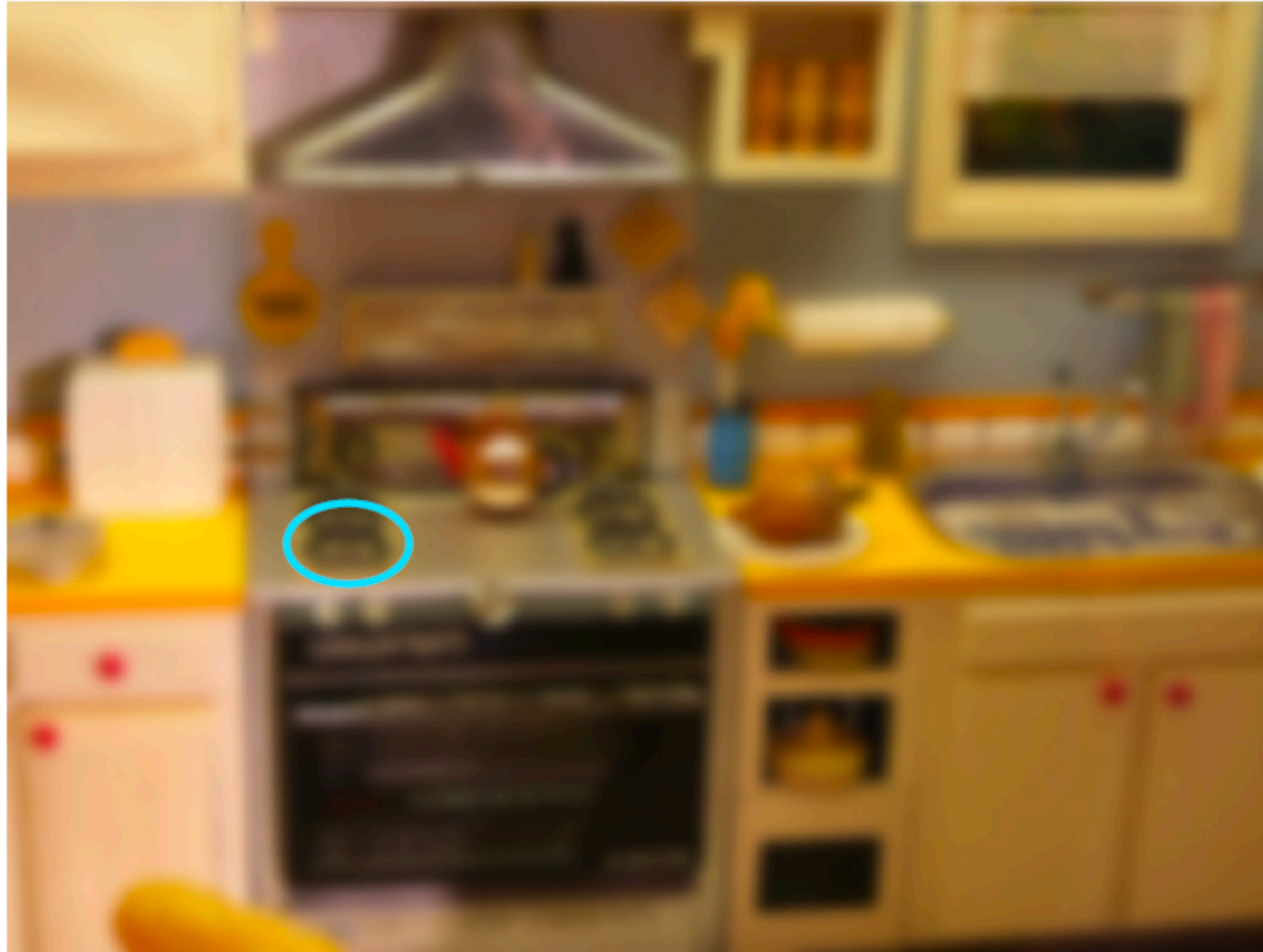


# What is This **Object**?



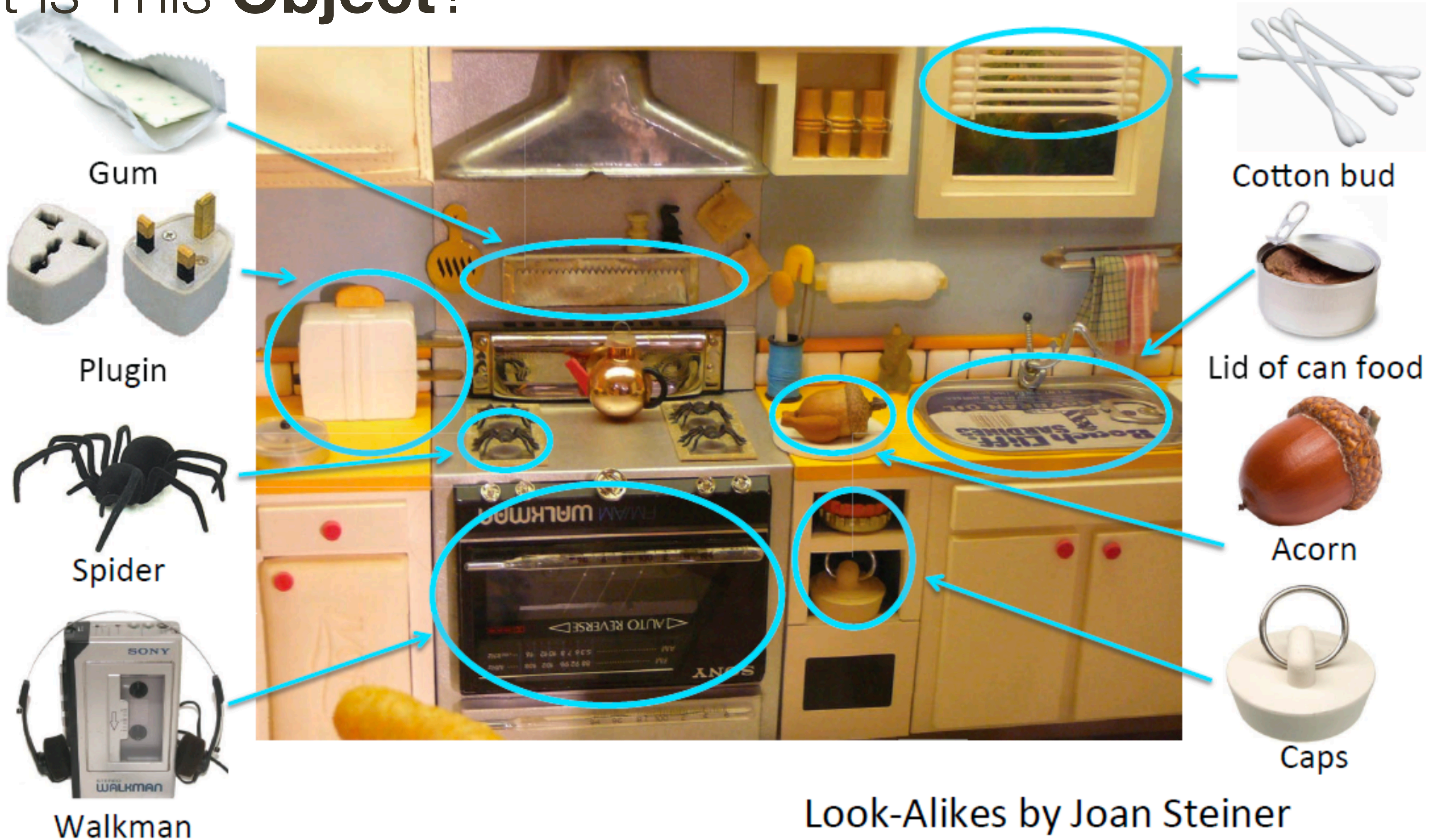


# What is This **Object**?





# What is This **Object**?



Look-Alikes by Joan Steiner

Figure source: Jianxiong Xiao



# Visual **Words**

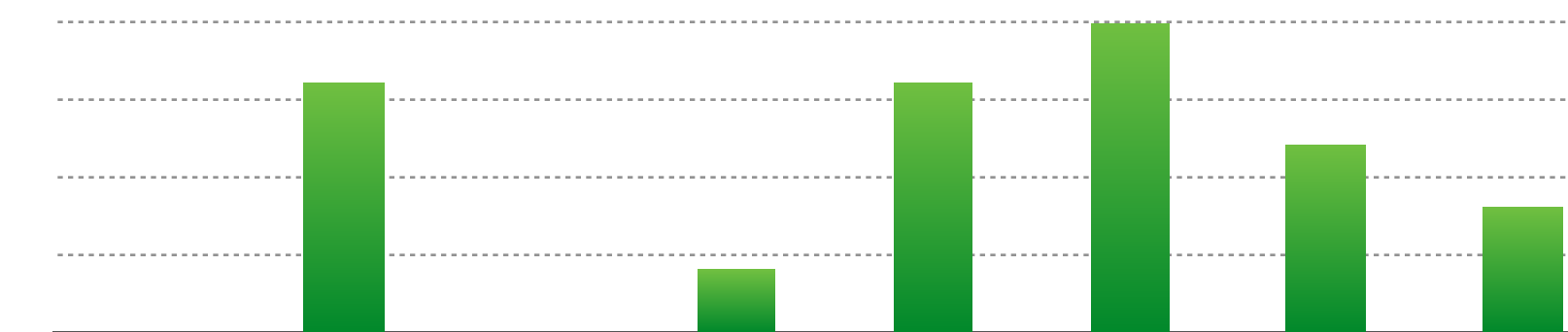
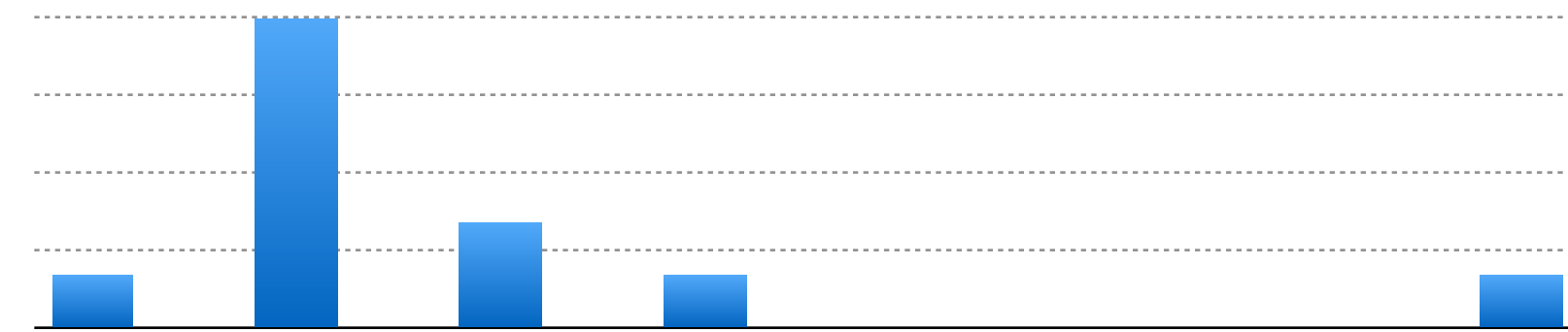
Many algorithms for image classification accumulate evidence on the basis of **visual words**.

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



<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

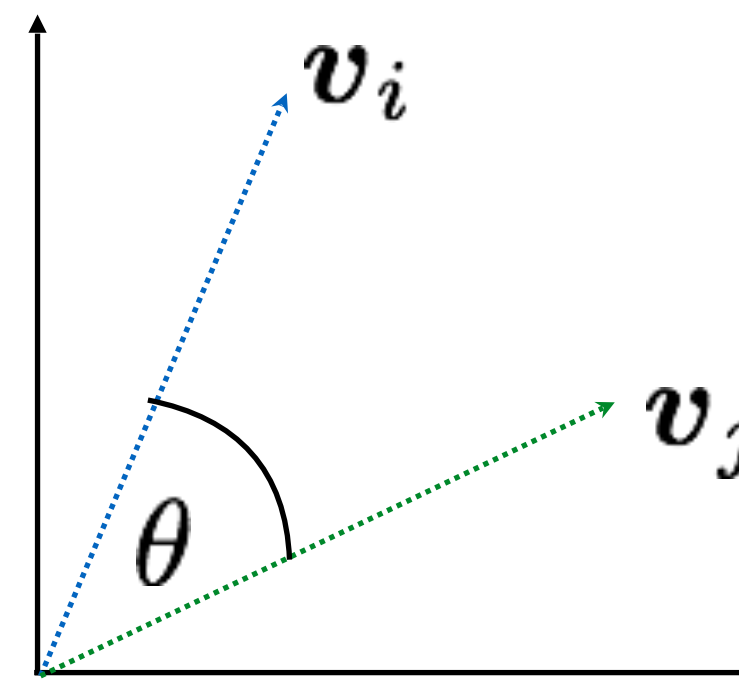
just a histogram over words



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:

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



Slide Credit: Ioannis (Yannis) Gkioulekas (CMU)

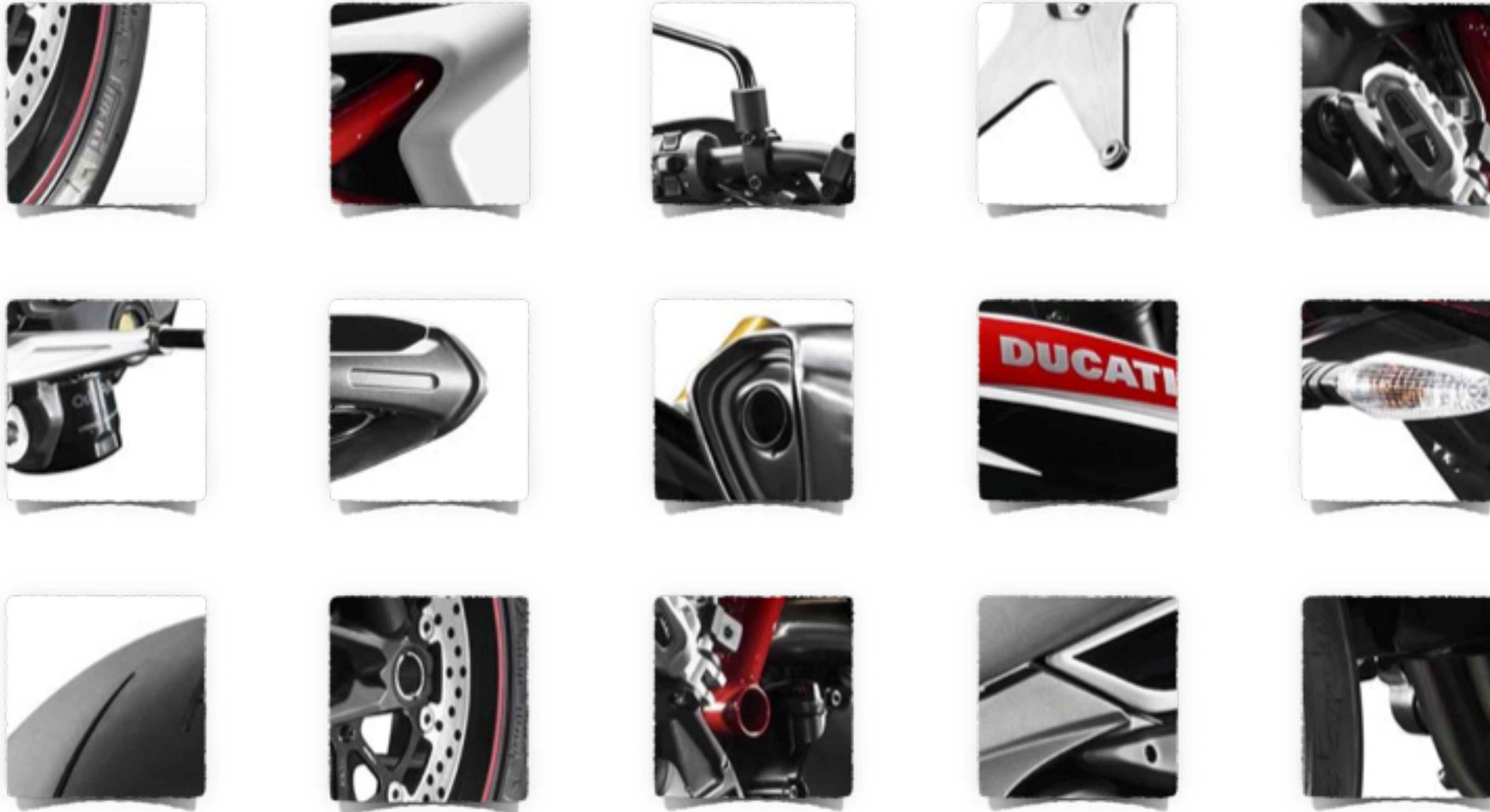


# Visual **Words**

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.

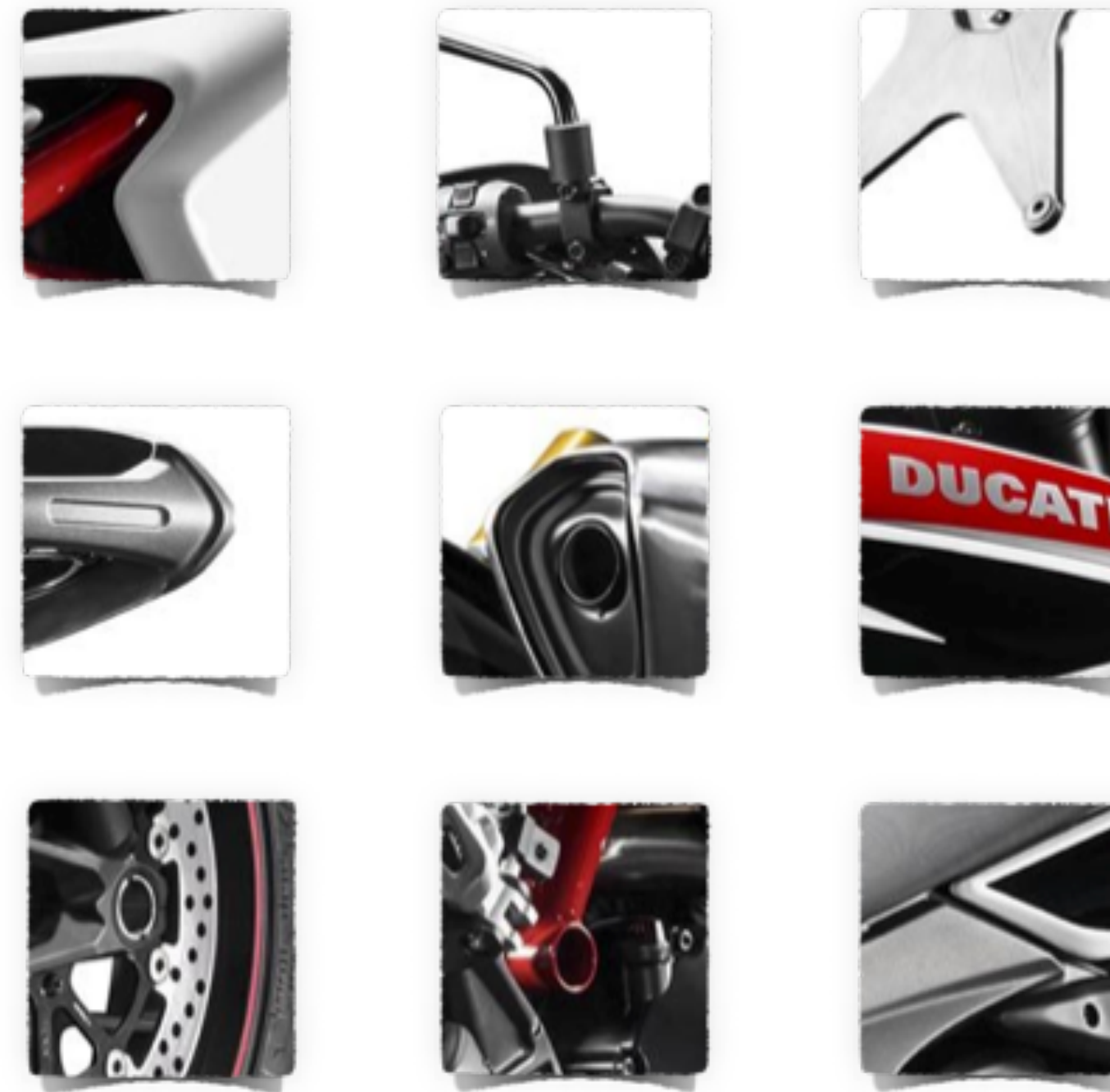
# What **Objects** do These Parts Belong To?





Some local feature are very informative

An object as

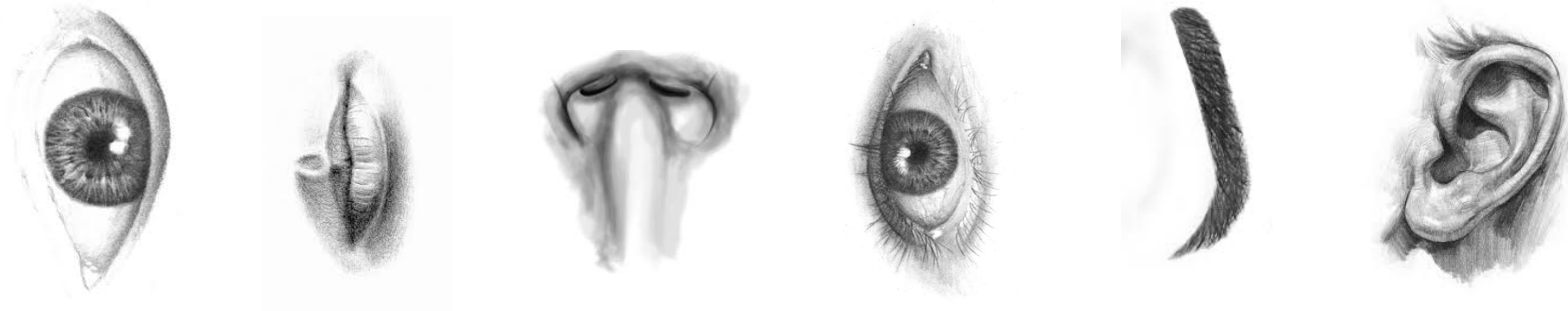


a collection of local features  
(bag-of-features)

- deals well with occlusion
- scale invariant
- rotation invariant

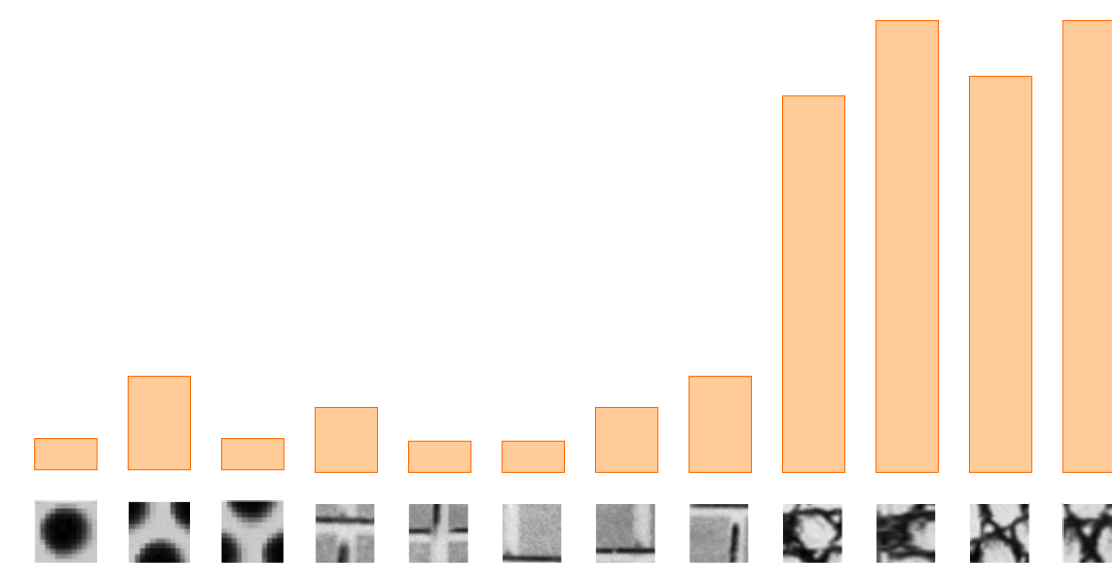
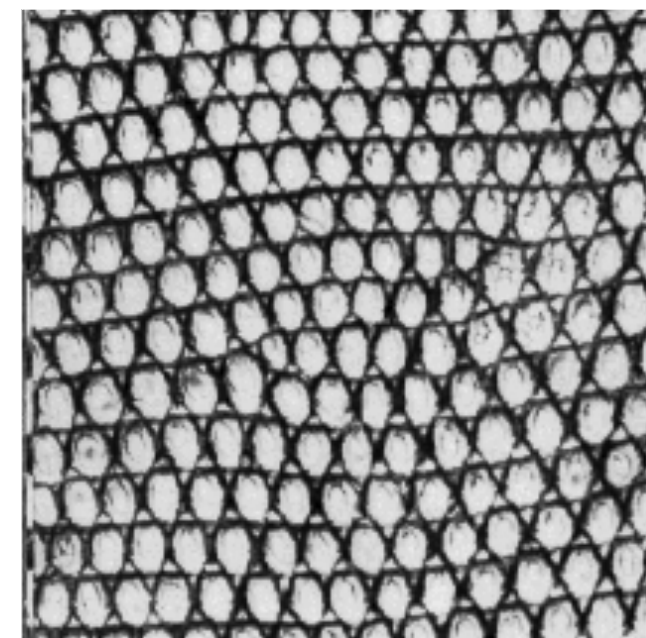
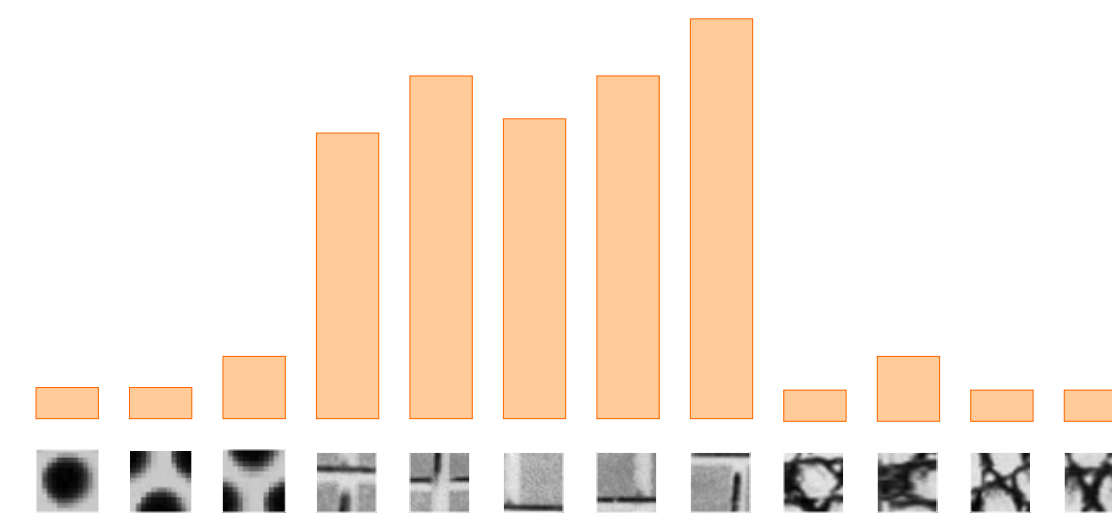
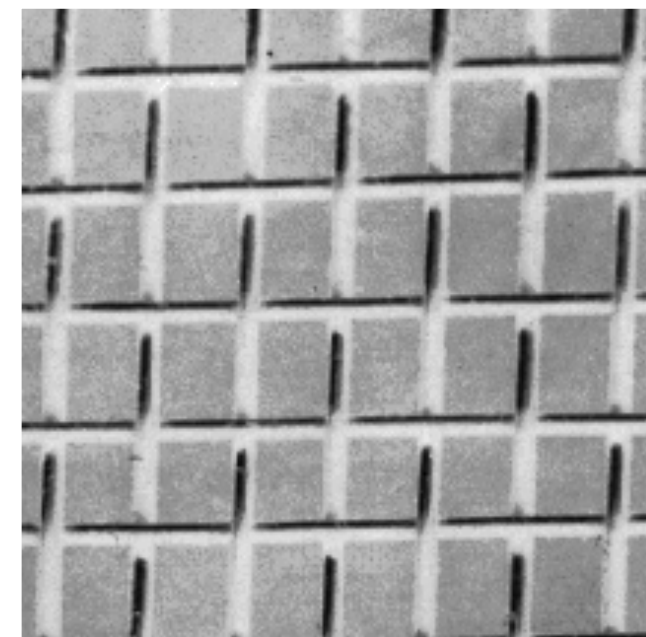
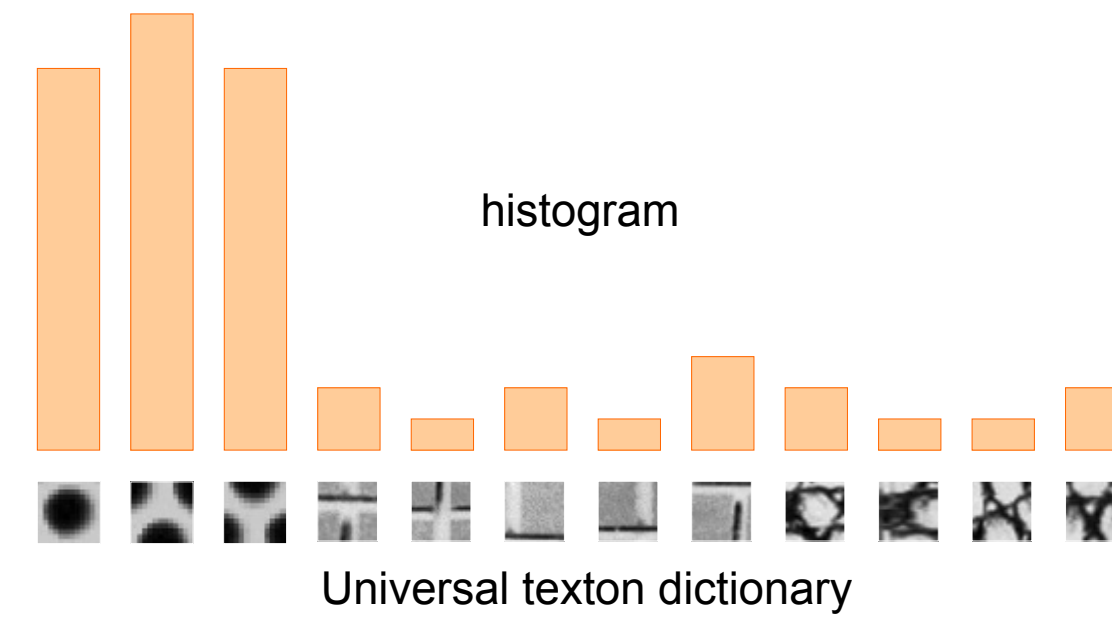
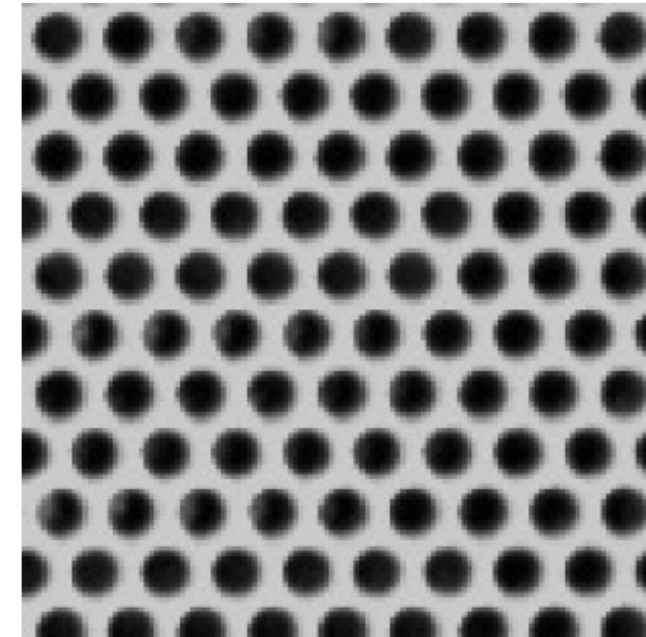


# (not so) Crazy Assumption



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:**

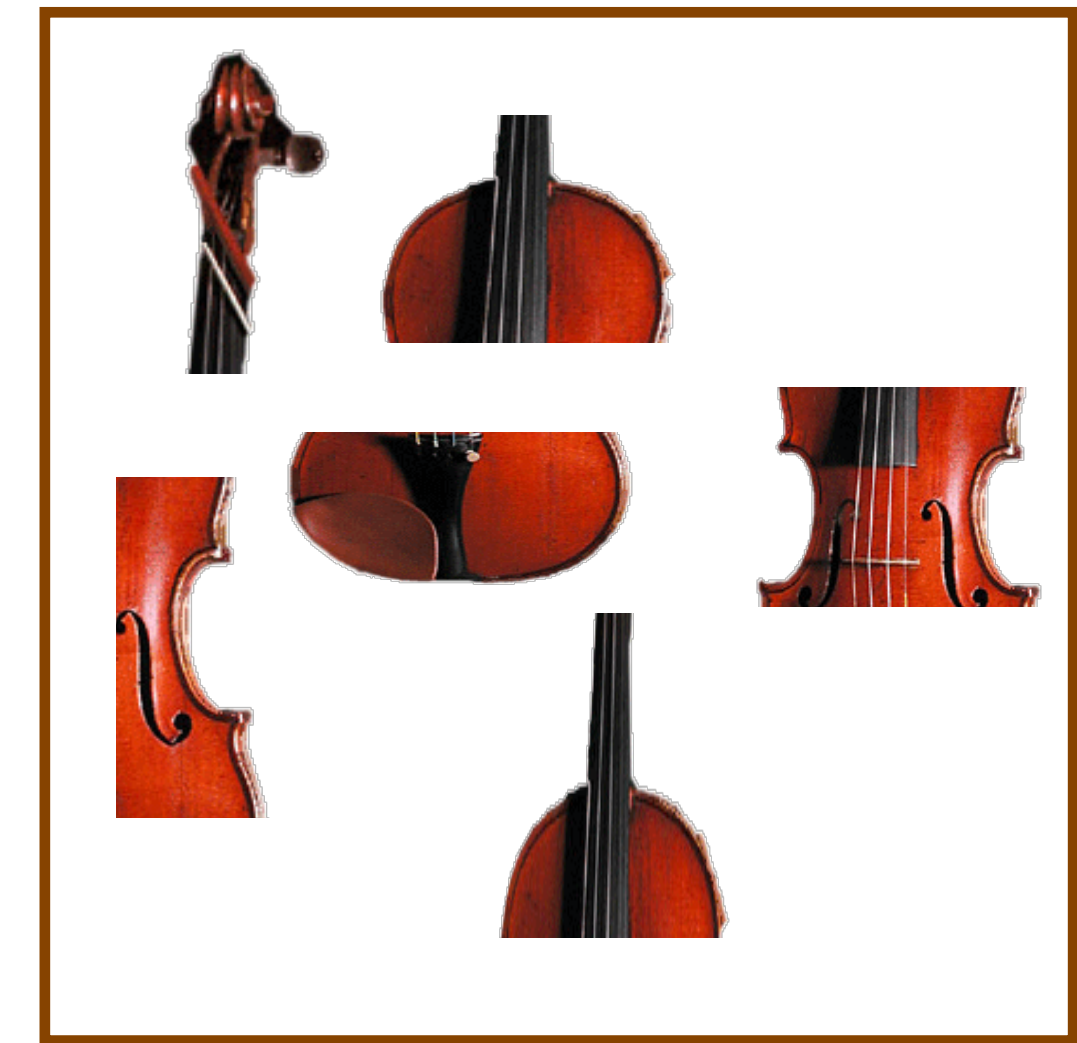
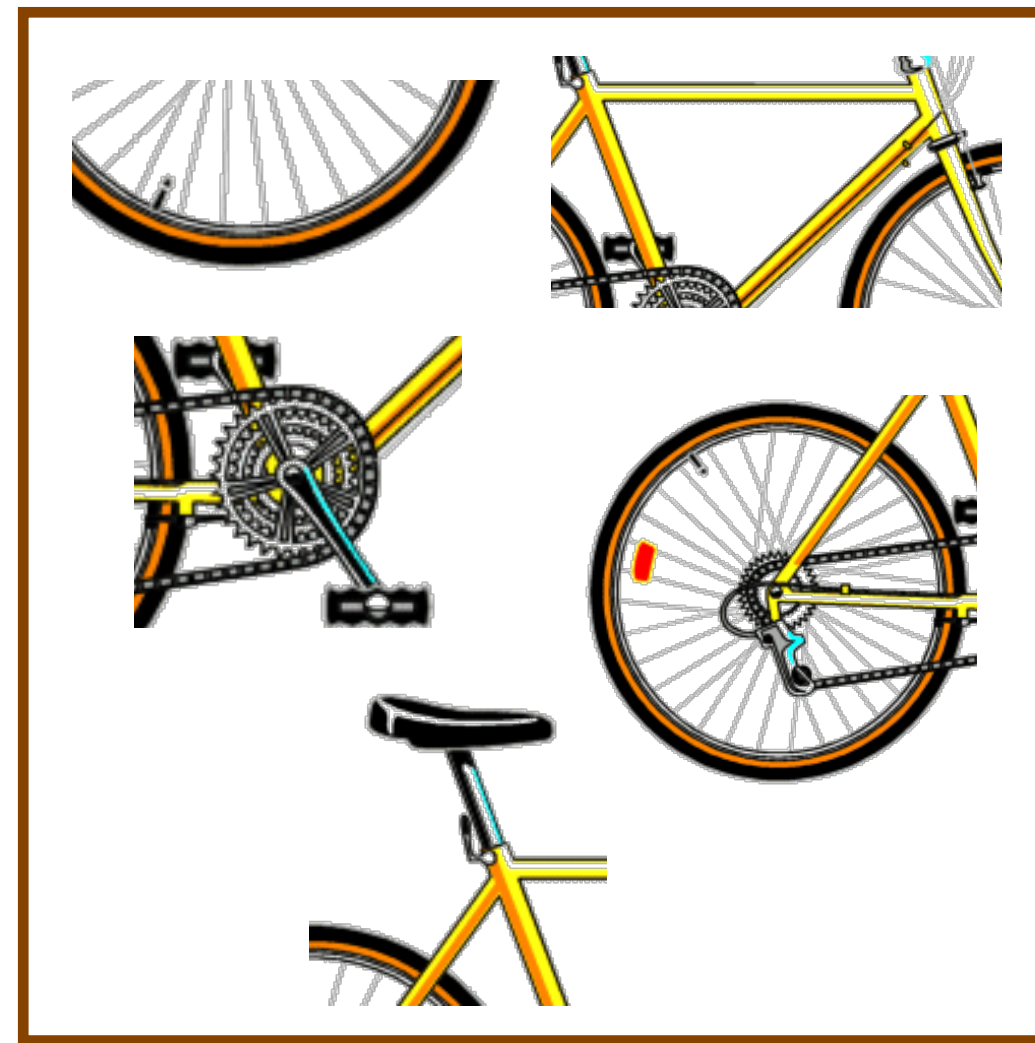
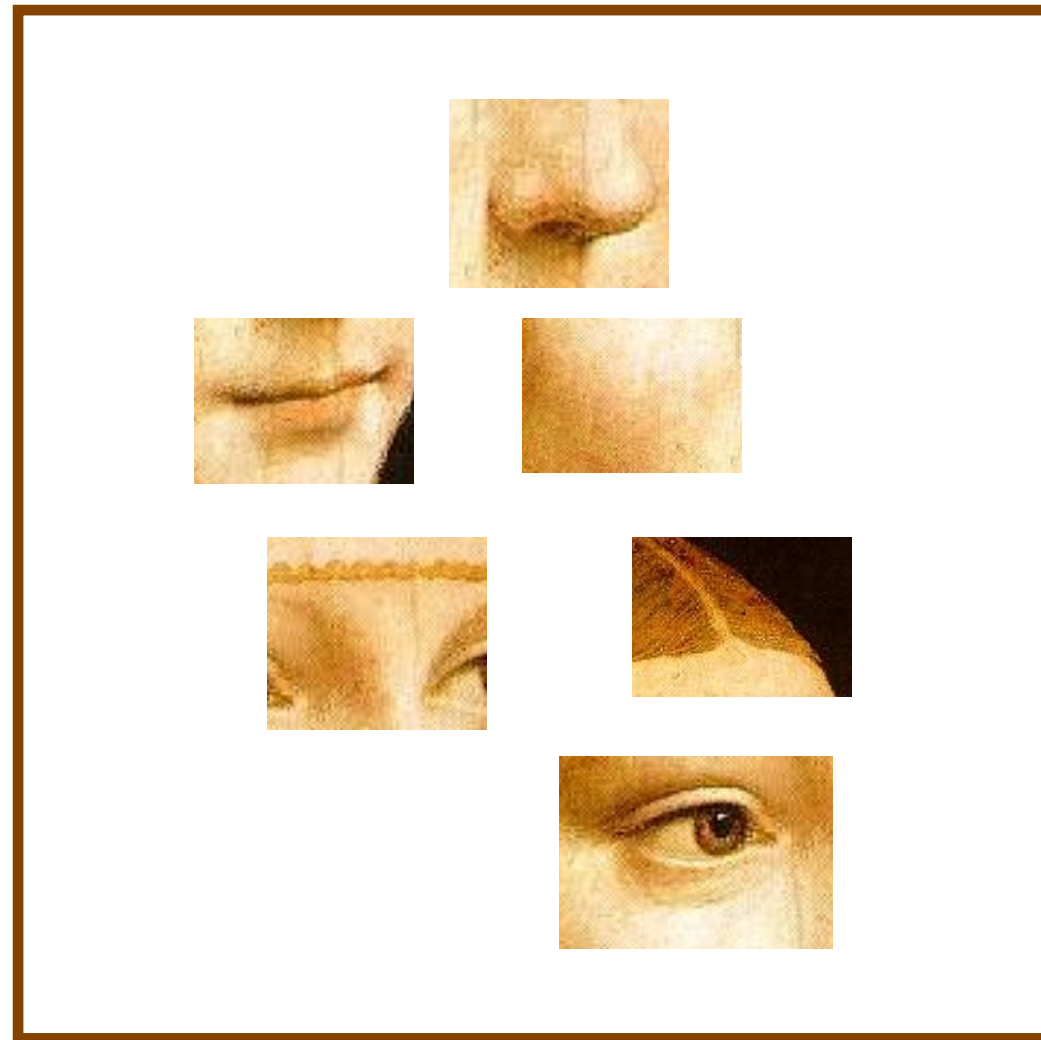
build Bags-of-Words (BOW) vectors  
for each image

## **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

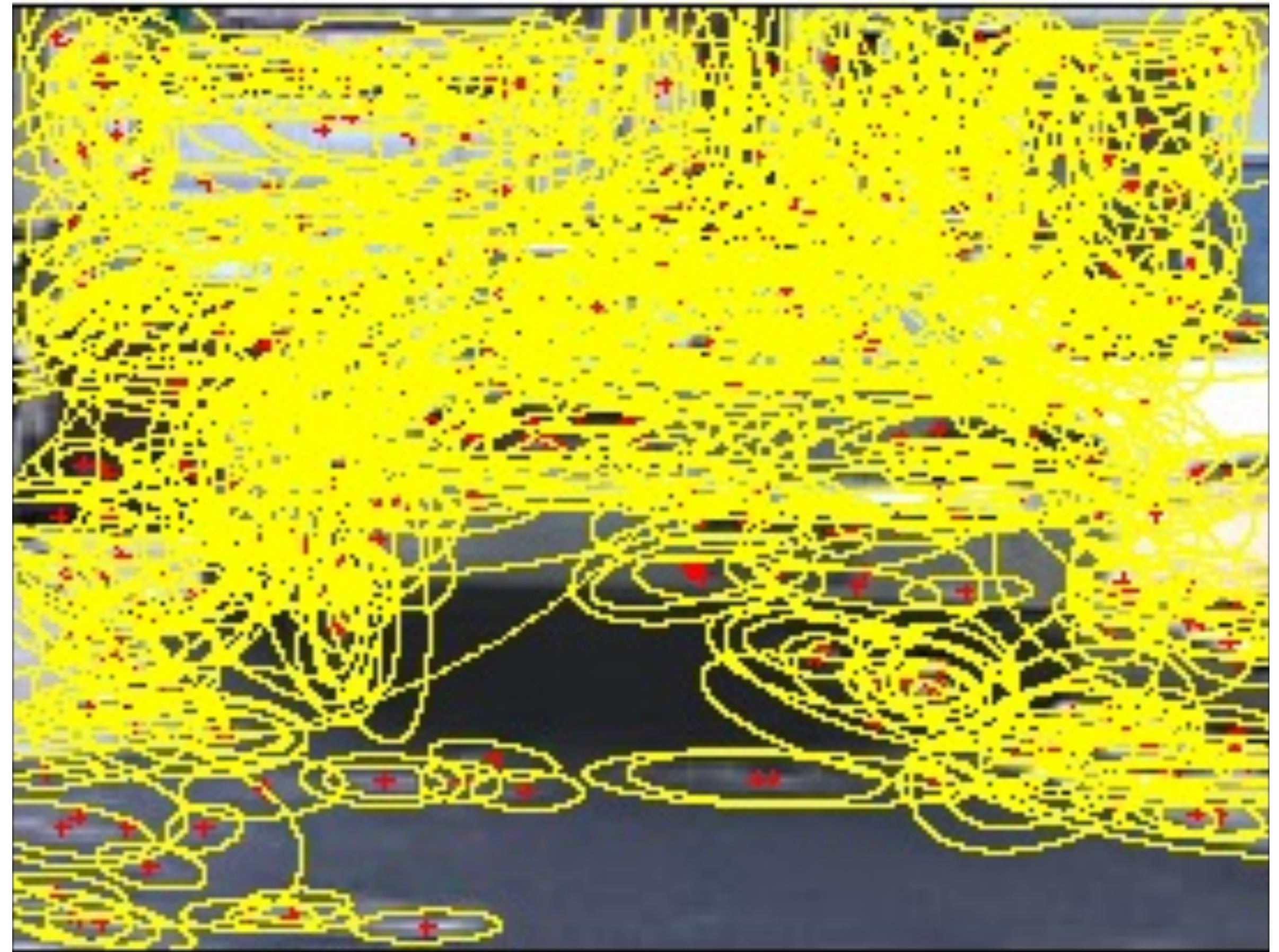
Fei-Fei & Perona, 2005

Sivic et al. 2005

- Other methods

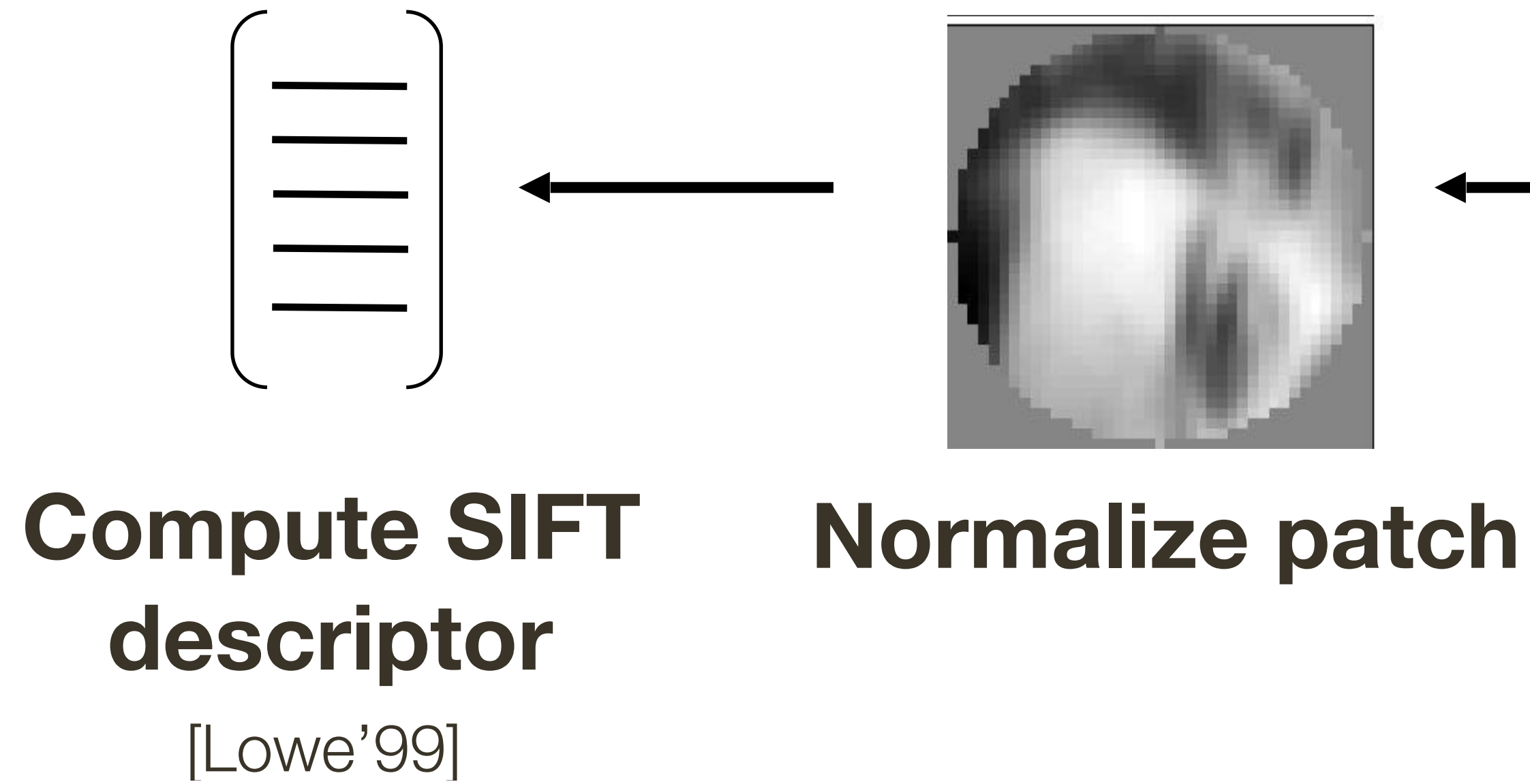
Random sampling (Vidal-Naquet & Ullman, 2002)

Segmentation-based patches (Barnard et al. 2003)

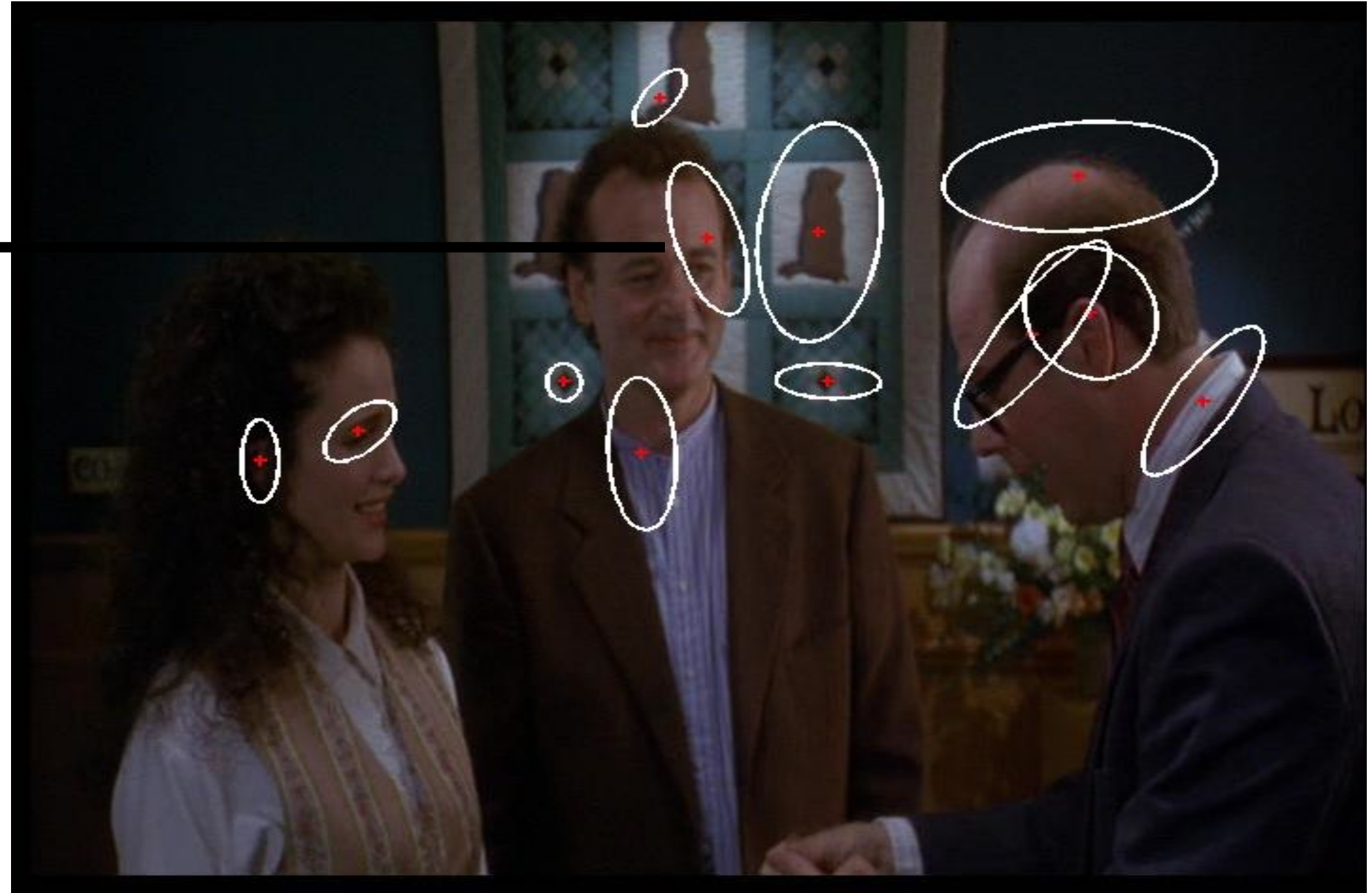




# Extracting **SIFT** Patches



**Normalize patch**



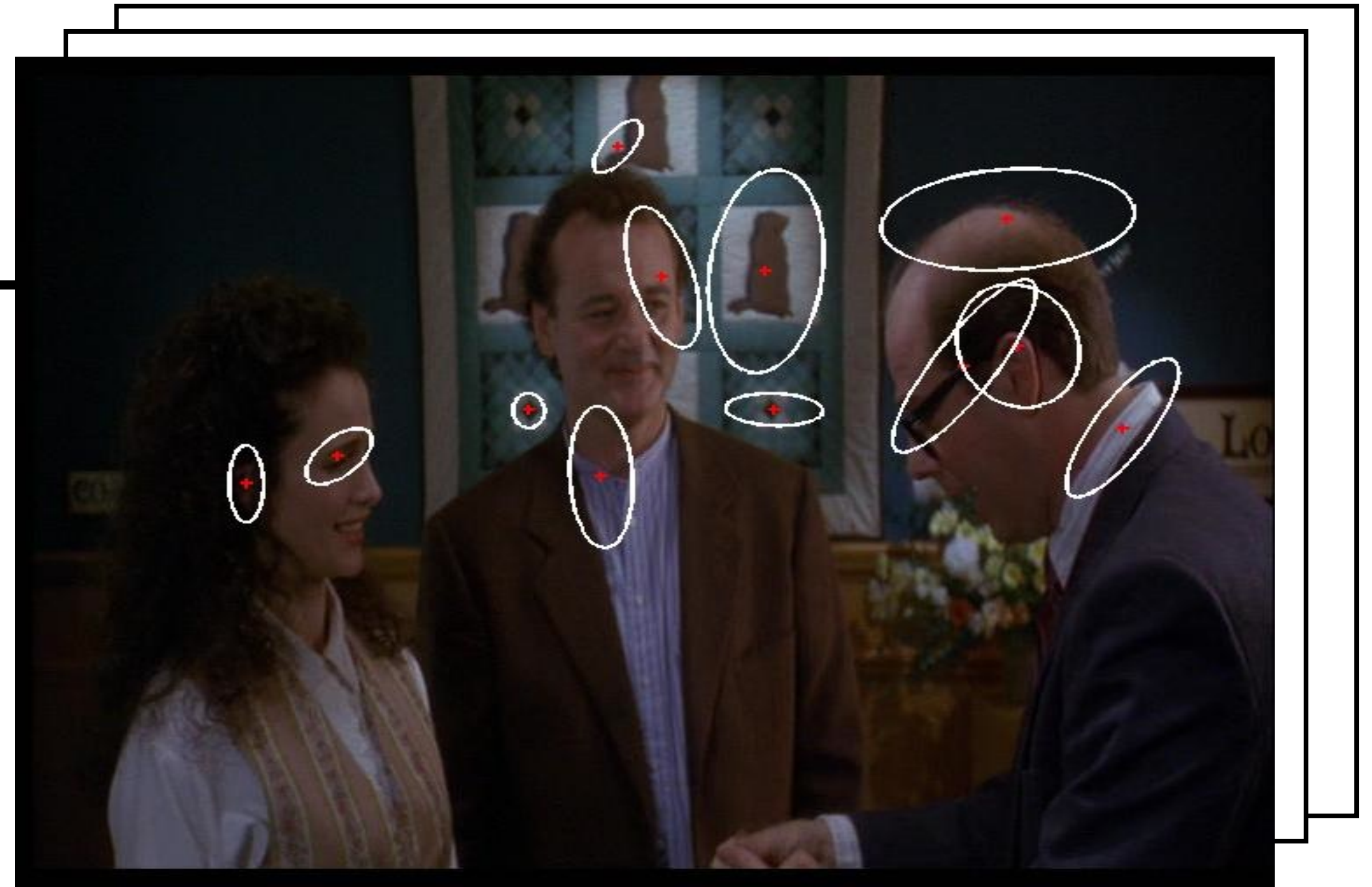
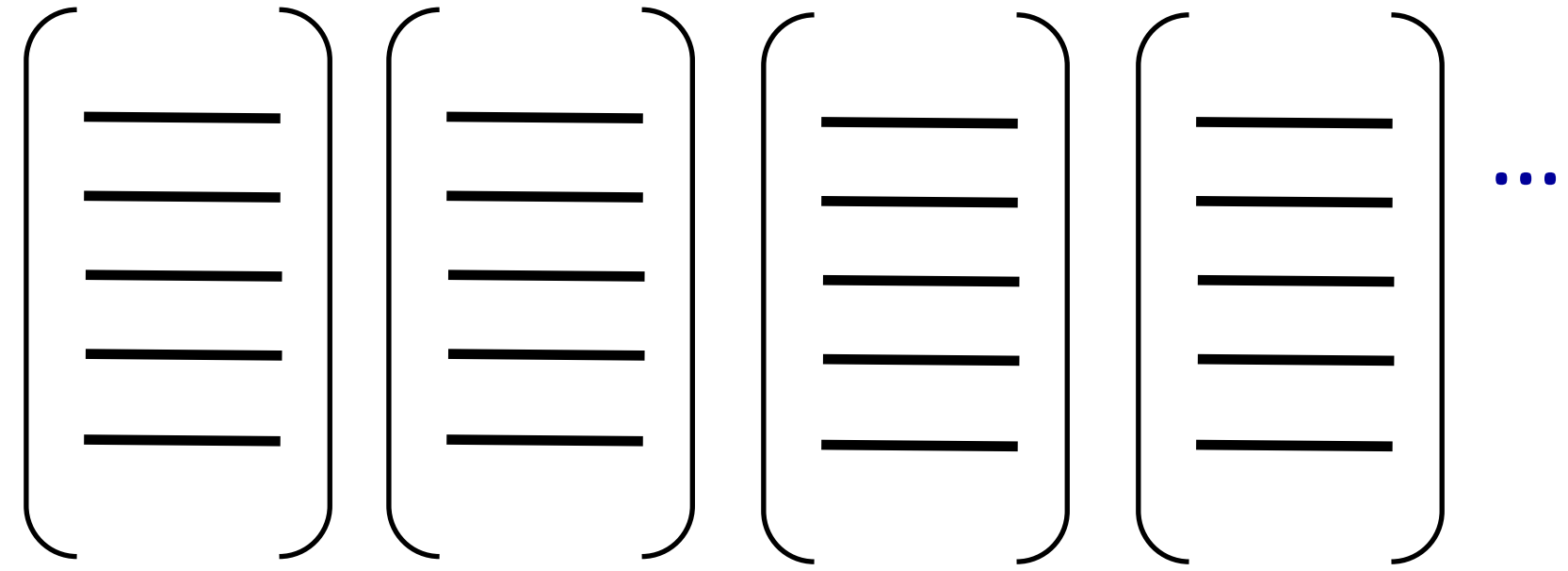
**Detect patches**

[Mikojczyk and Schmid '02]

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

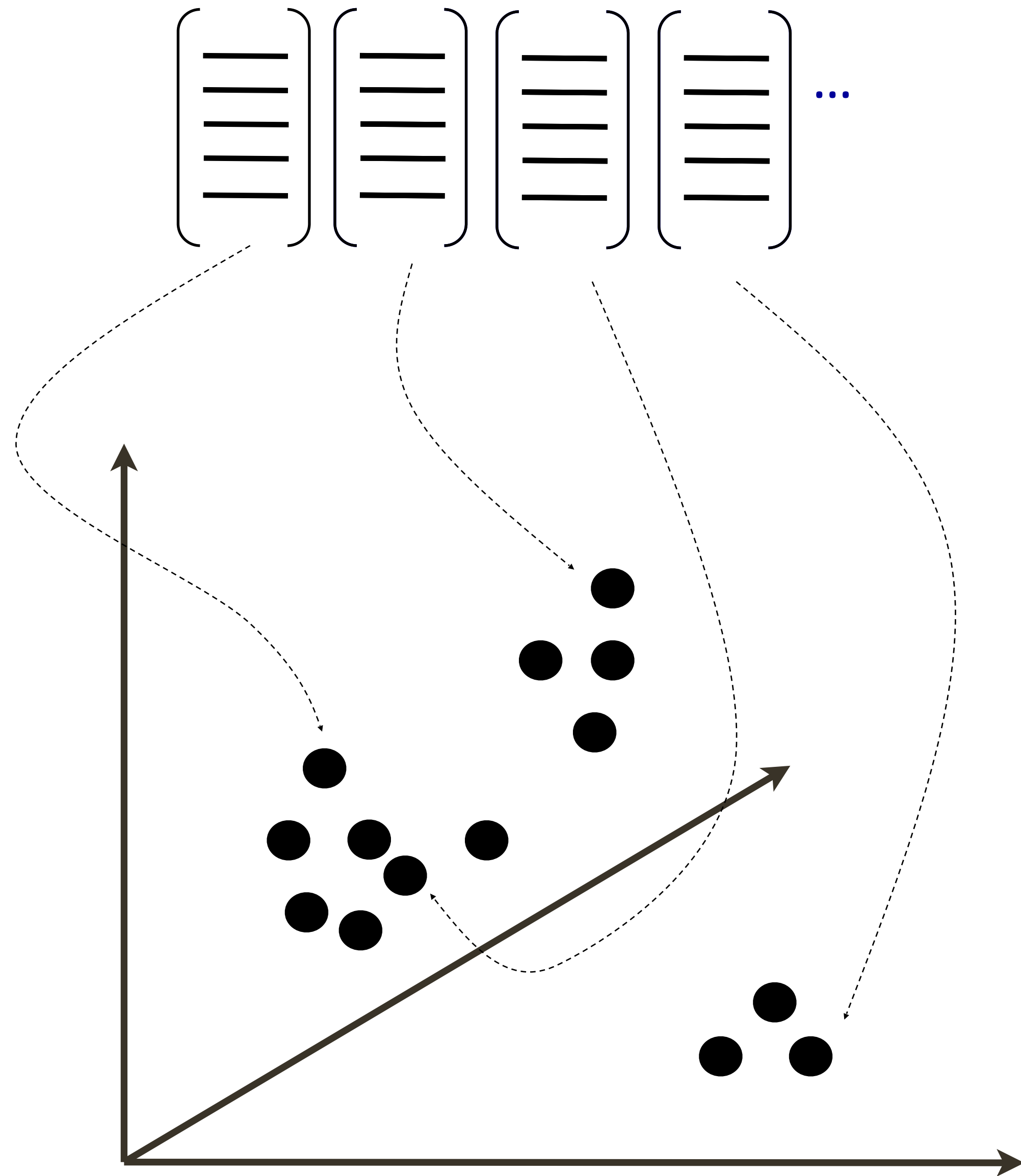
[Sivic & Zisserman, '03]

# Extracting **SIFT** Patches

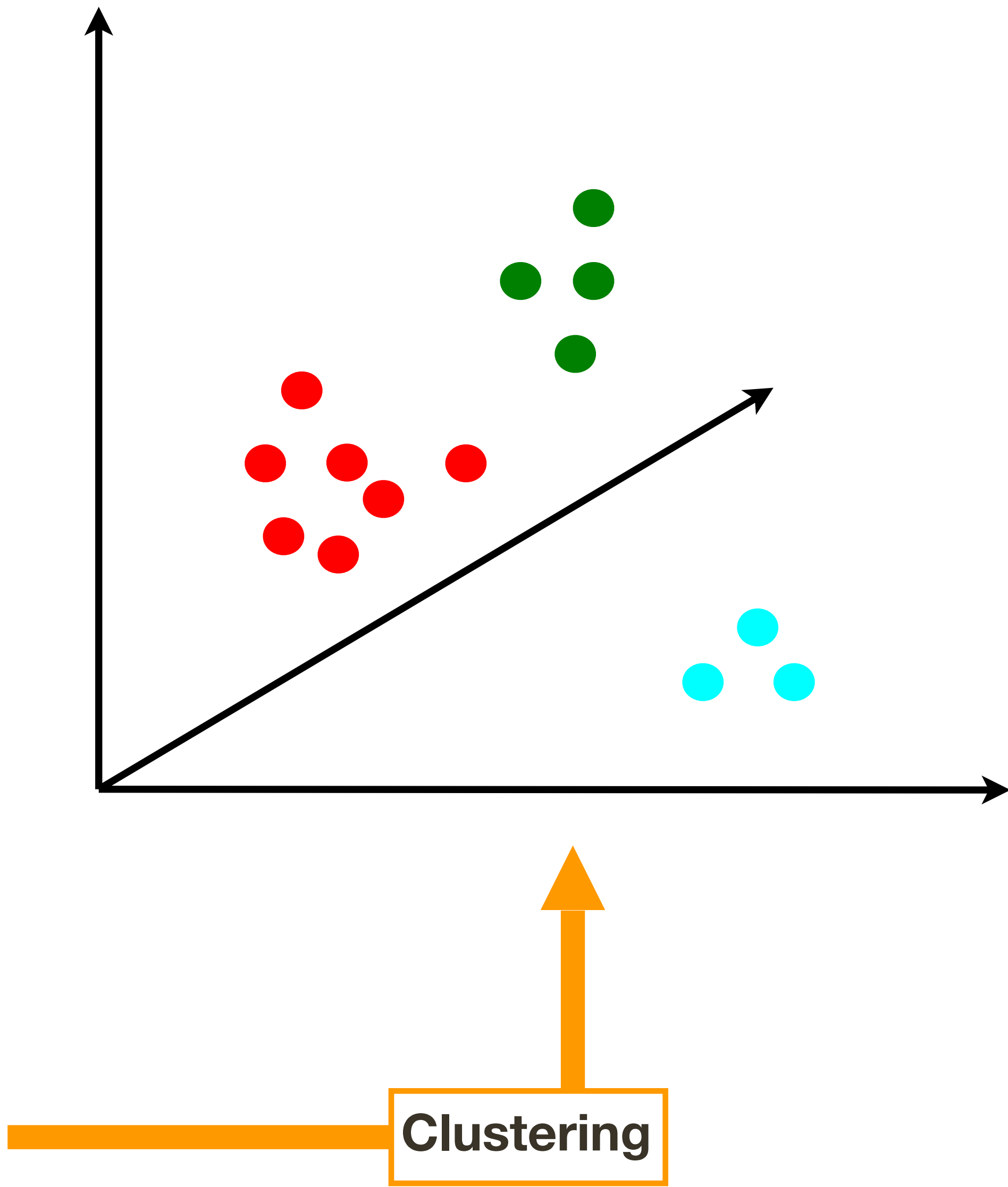
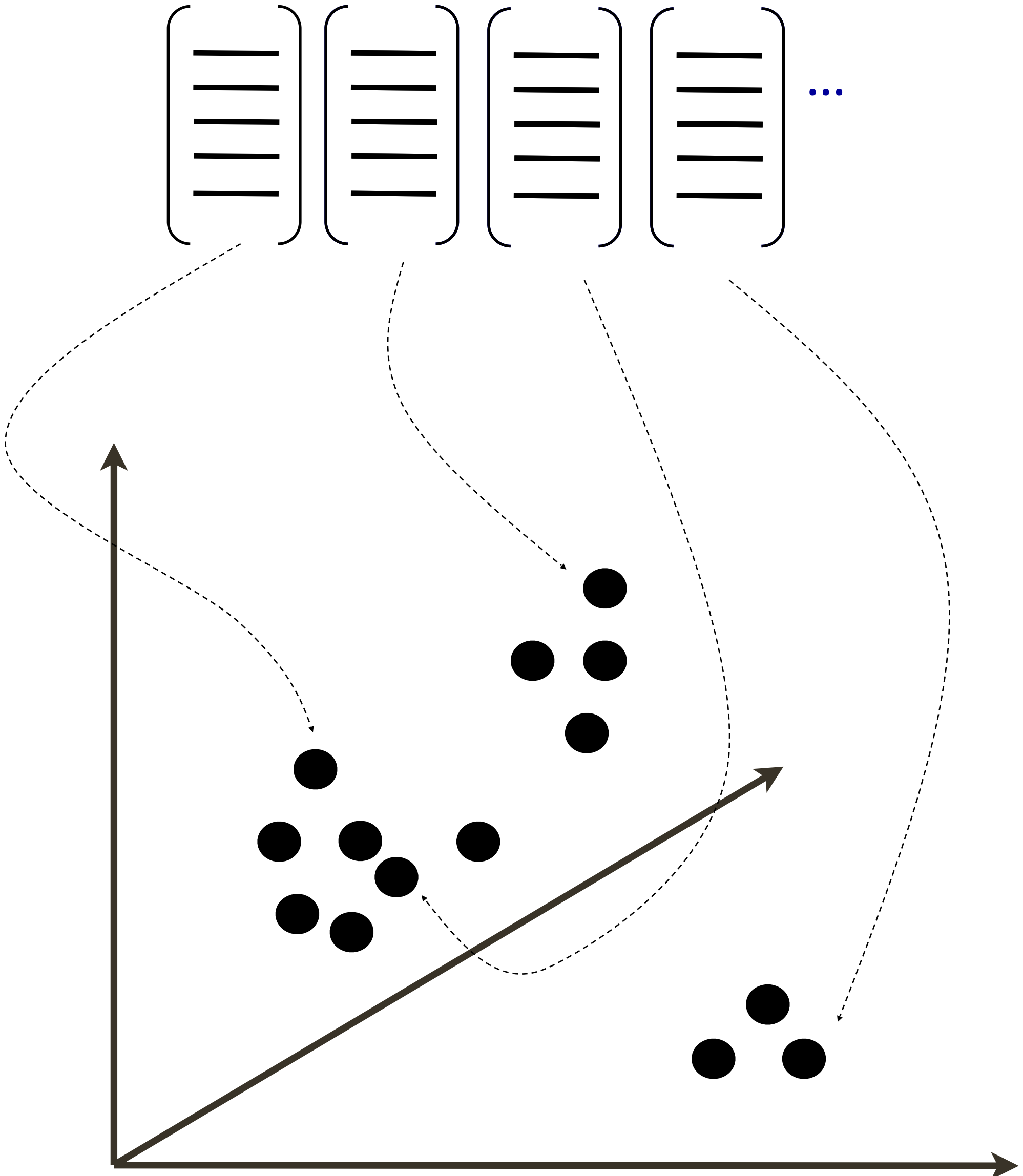




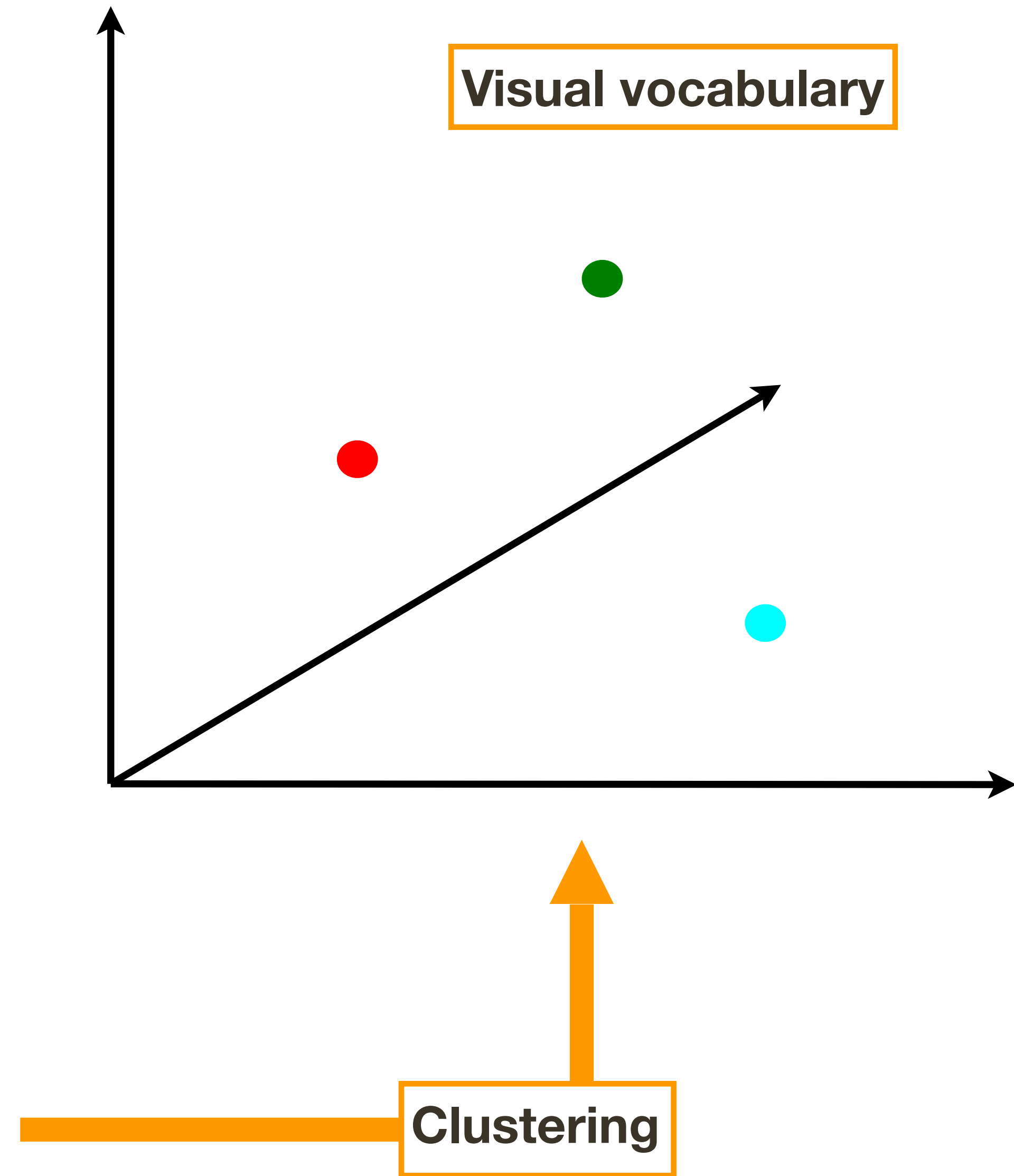
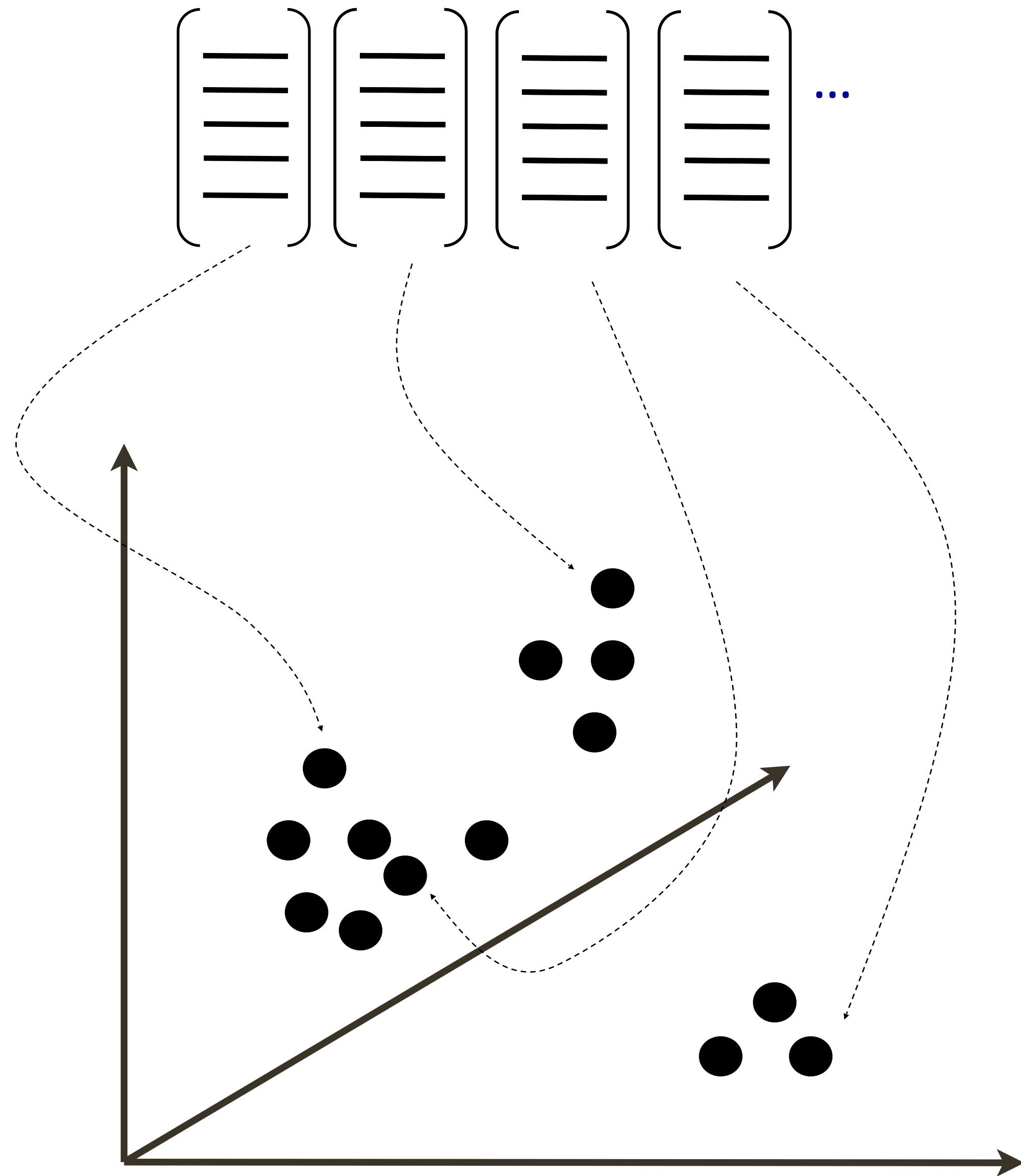
# Creating Dictionary



# Creating Dictionary



# Creating Dictionary





# **K-means** clustering

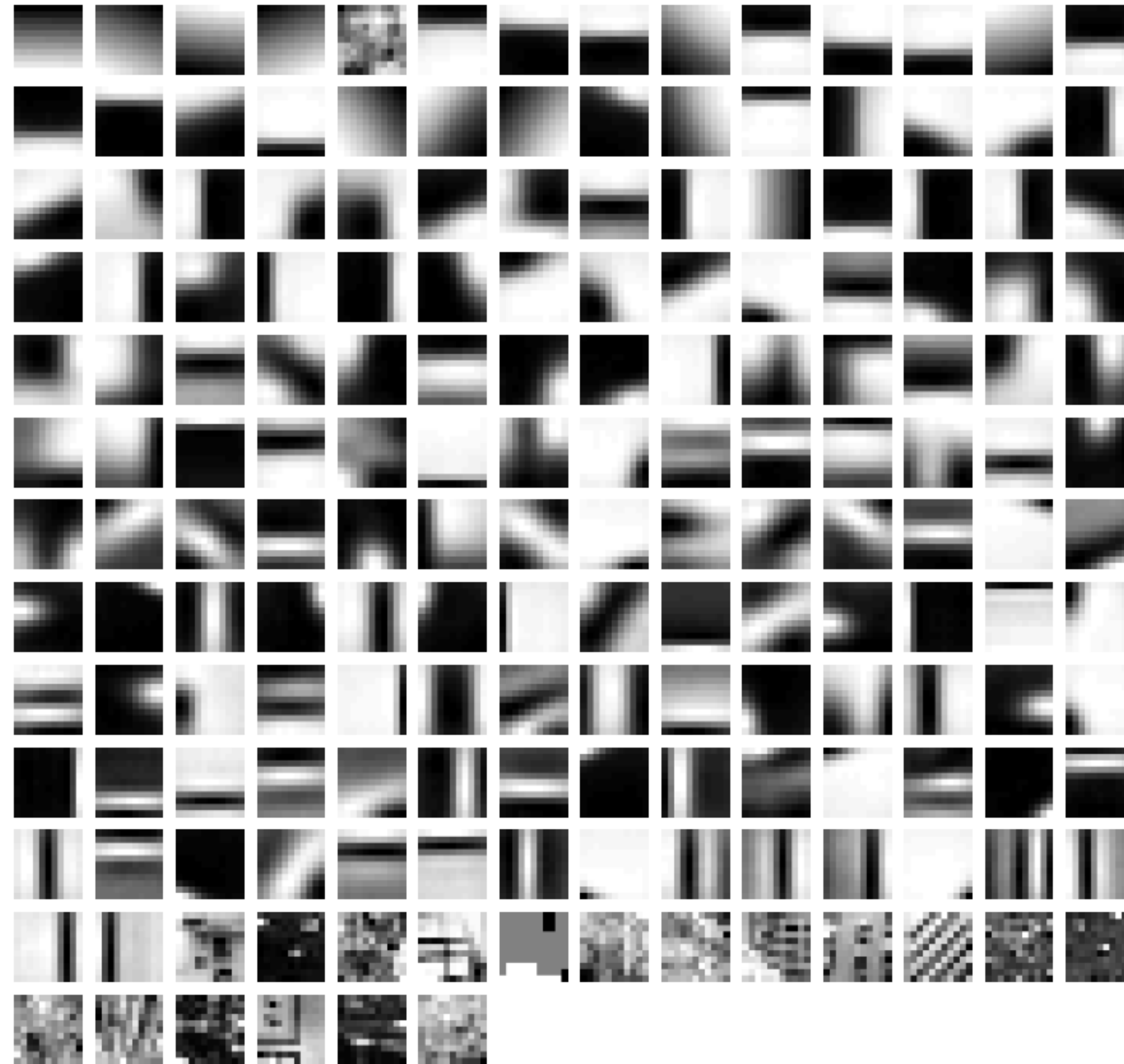
# K-means Clustering

**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).

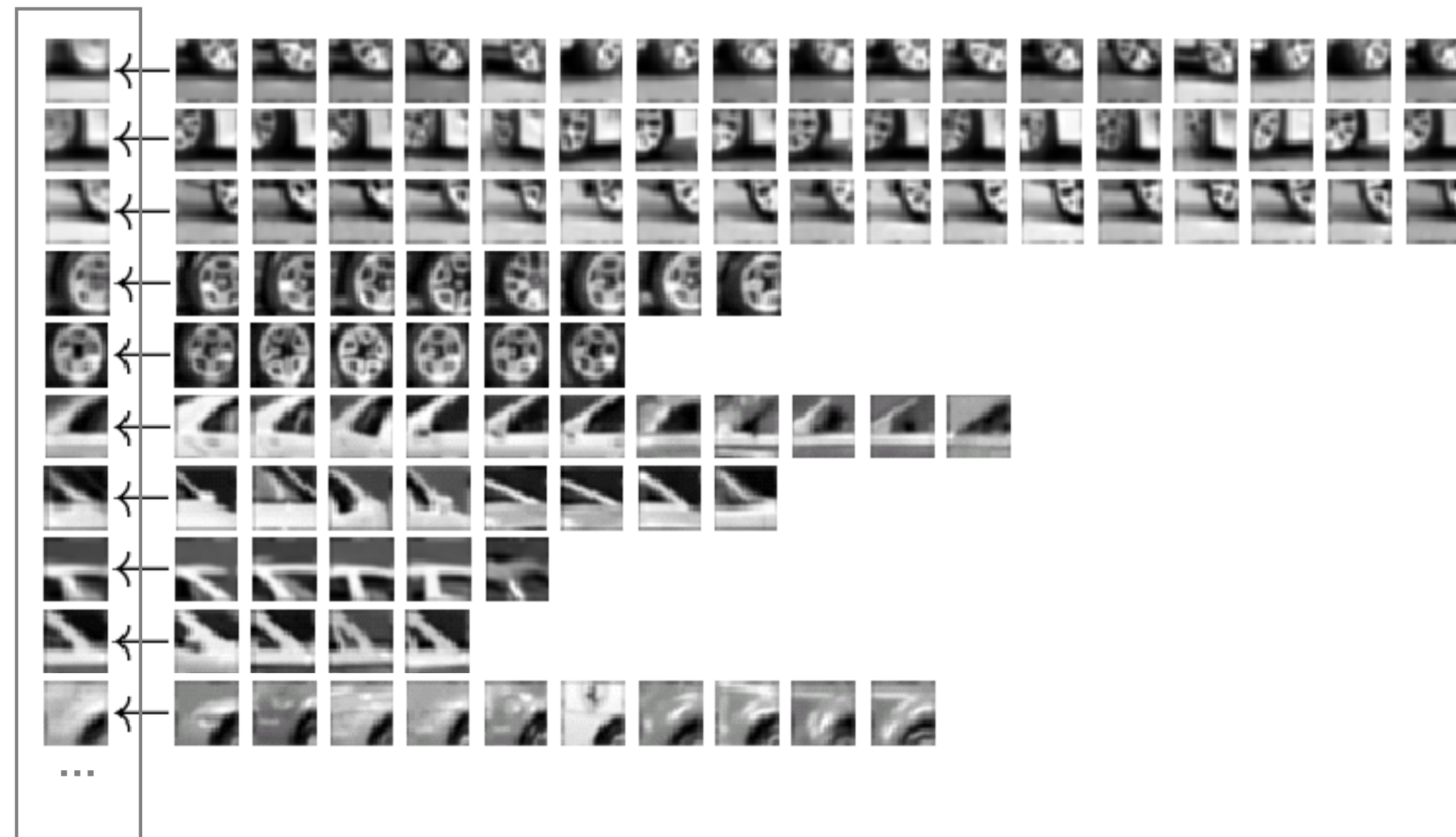
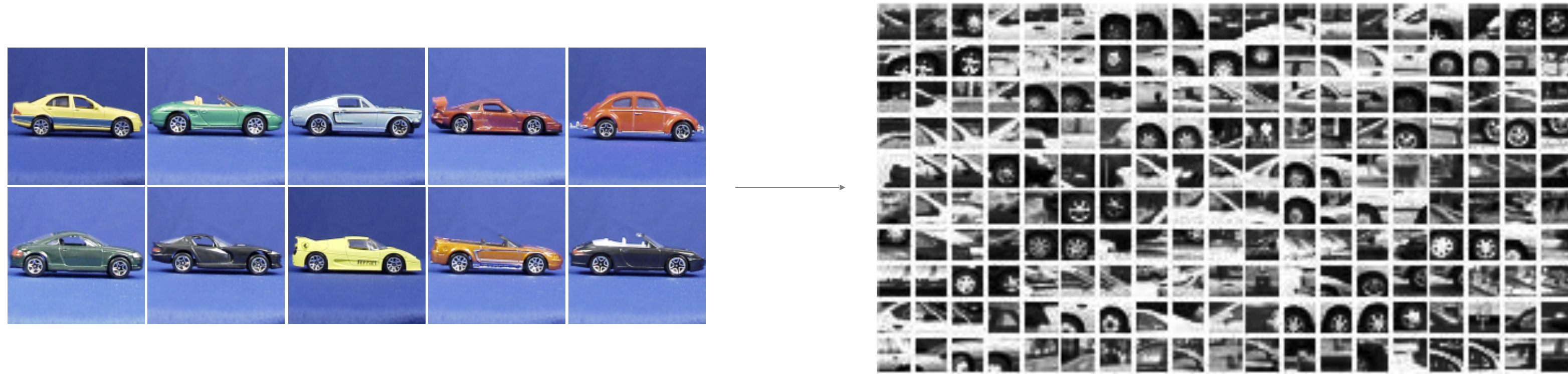
**K-means** clustering is initialization dependent and converges to a local minimum

# Example **Visual Dictionary**





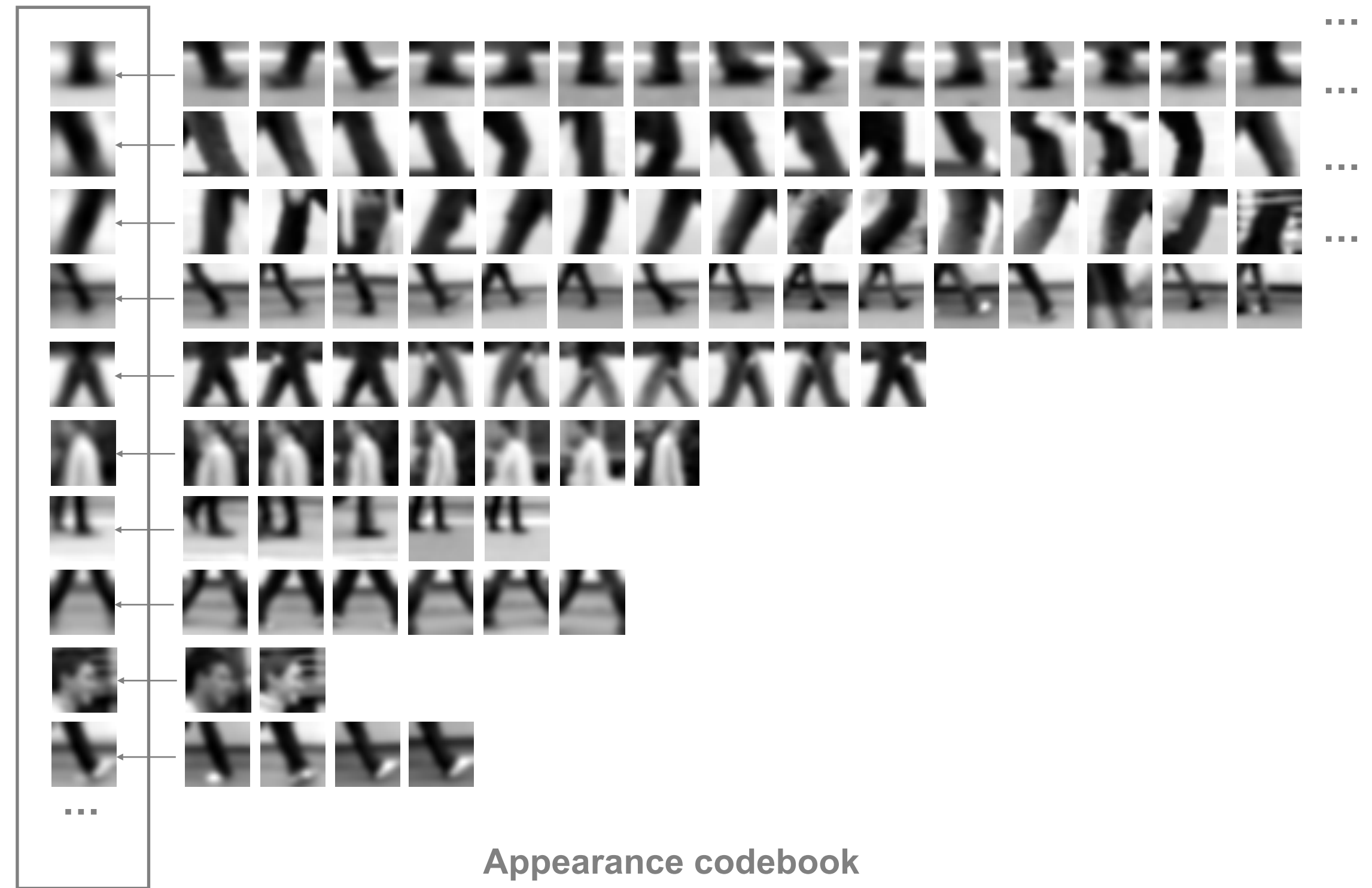
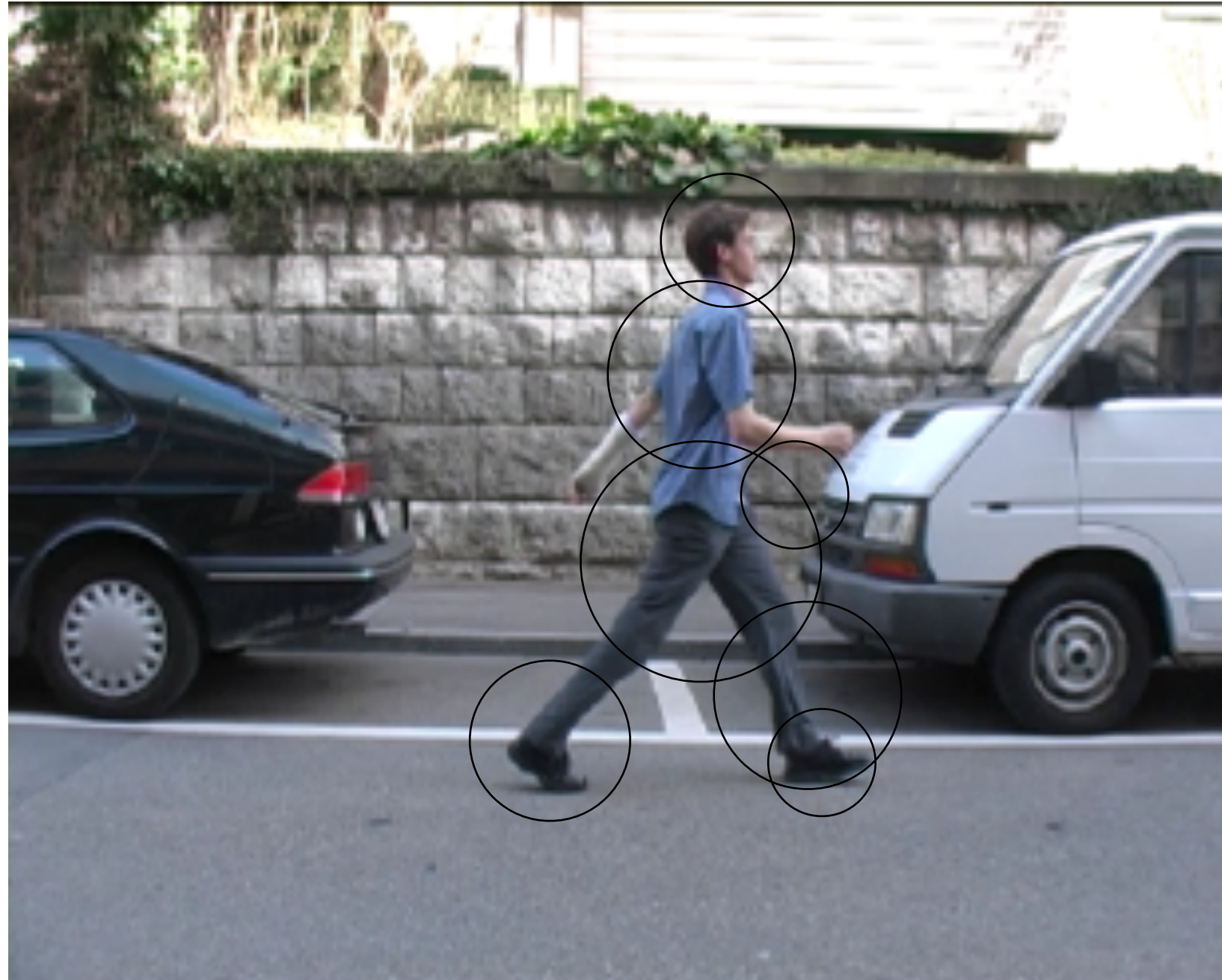
# Example **Visual Dictionary**



**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

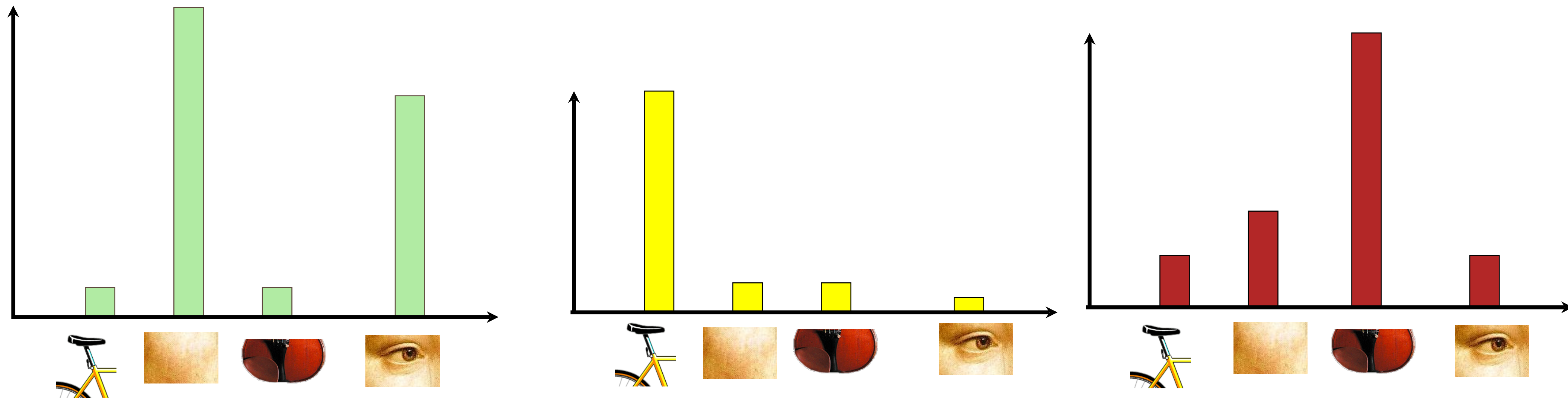


1. **Quantization:** image features gets associated to a visual word (nearest cluster center)

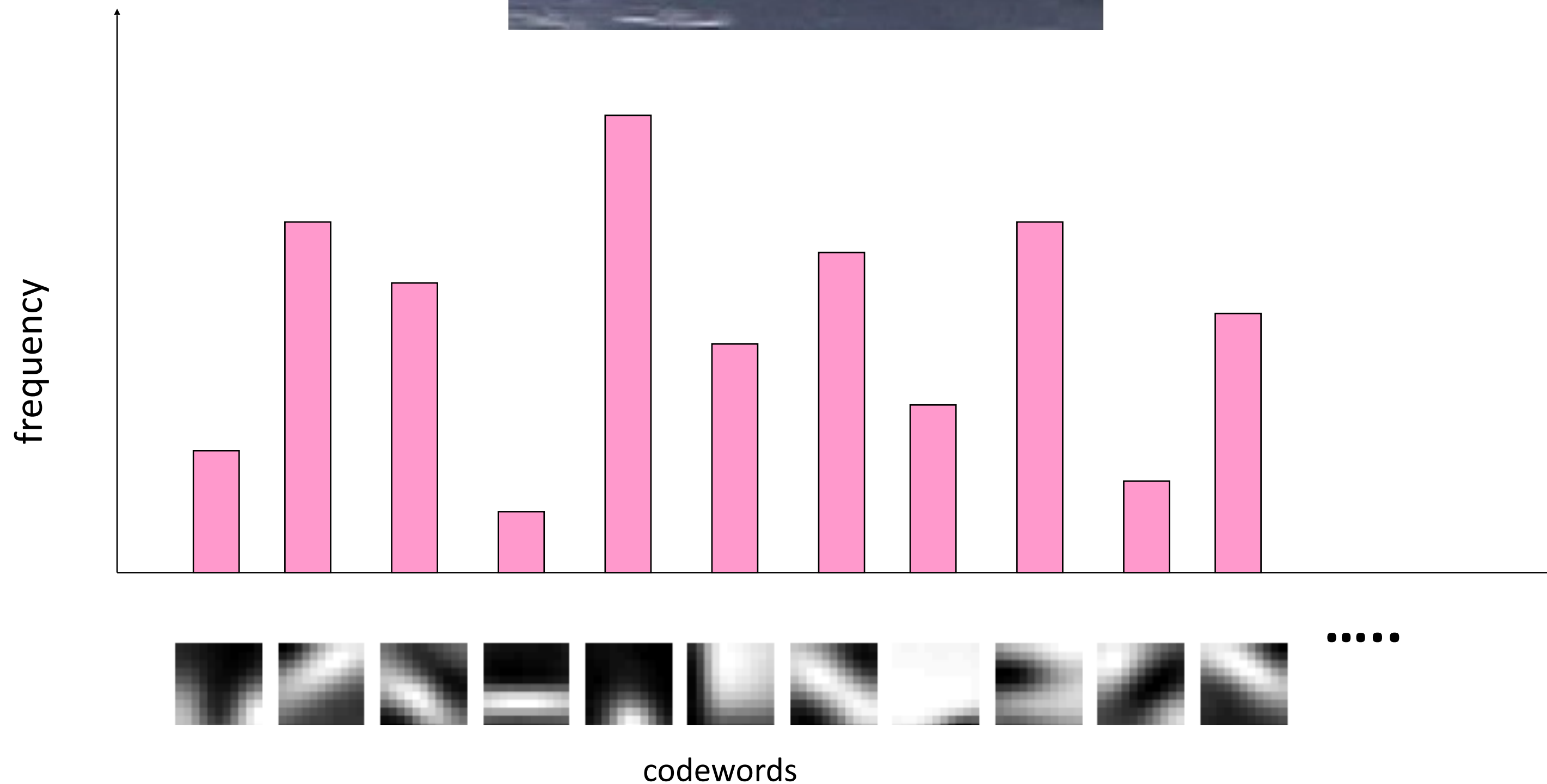


## 2. Encode: build Bag-of-Words (BOW) vectors for each image

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

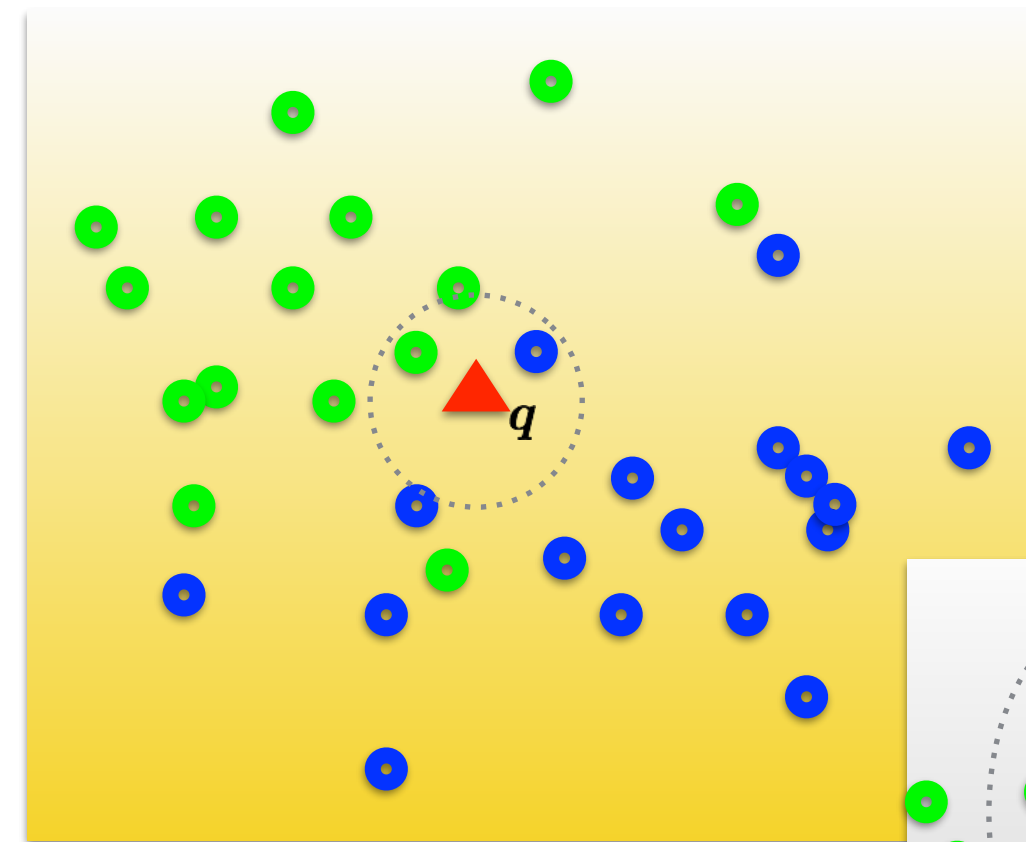
## **Encode:**

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

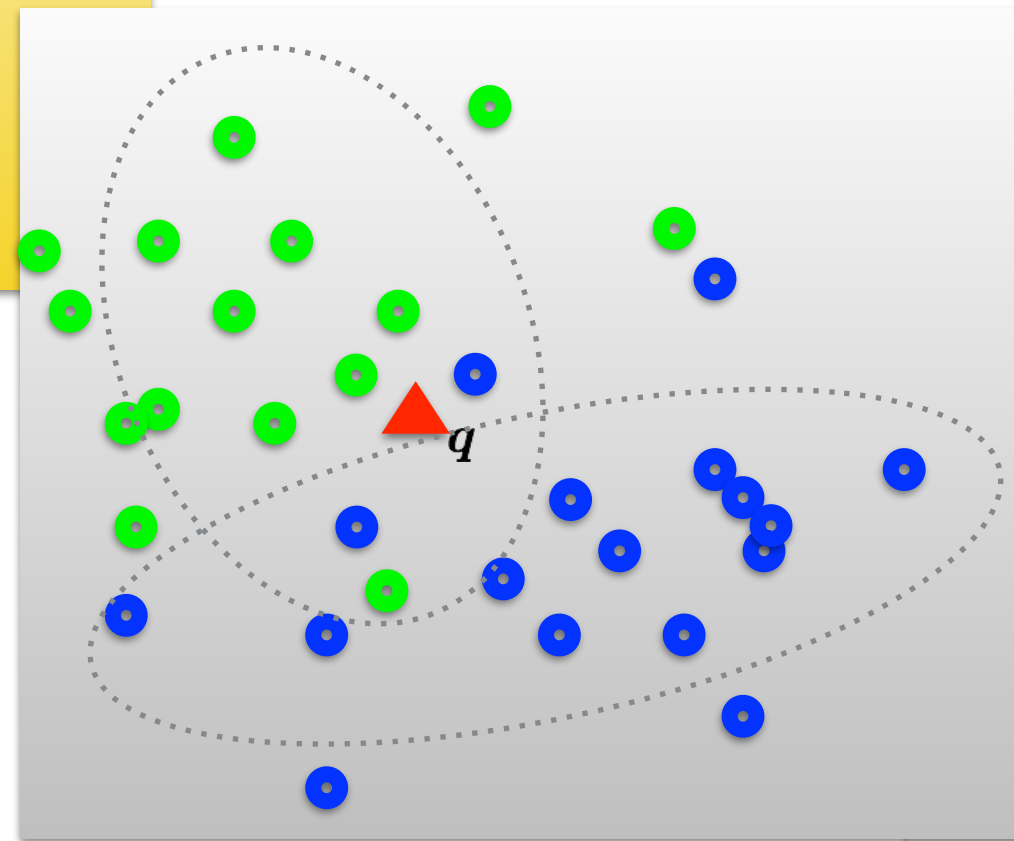
## **Classify:**

Train and test data using BOWs

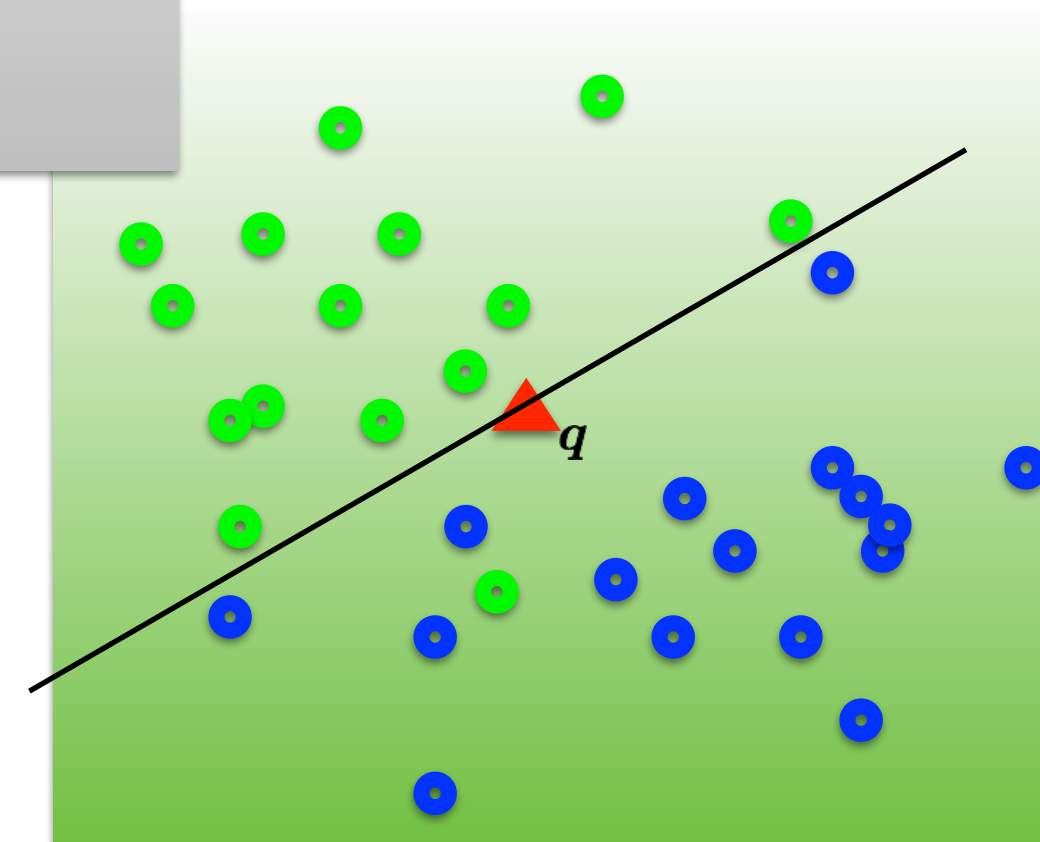
### 3. Classify: Train and text classifier using BOWs



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

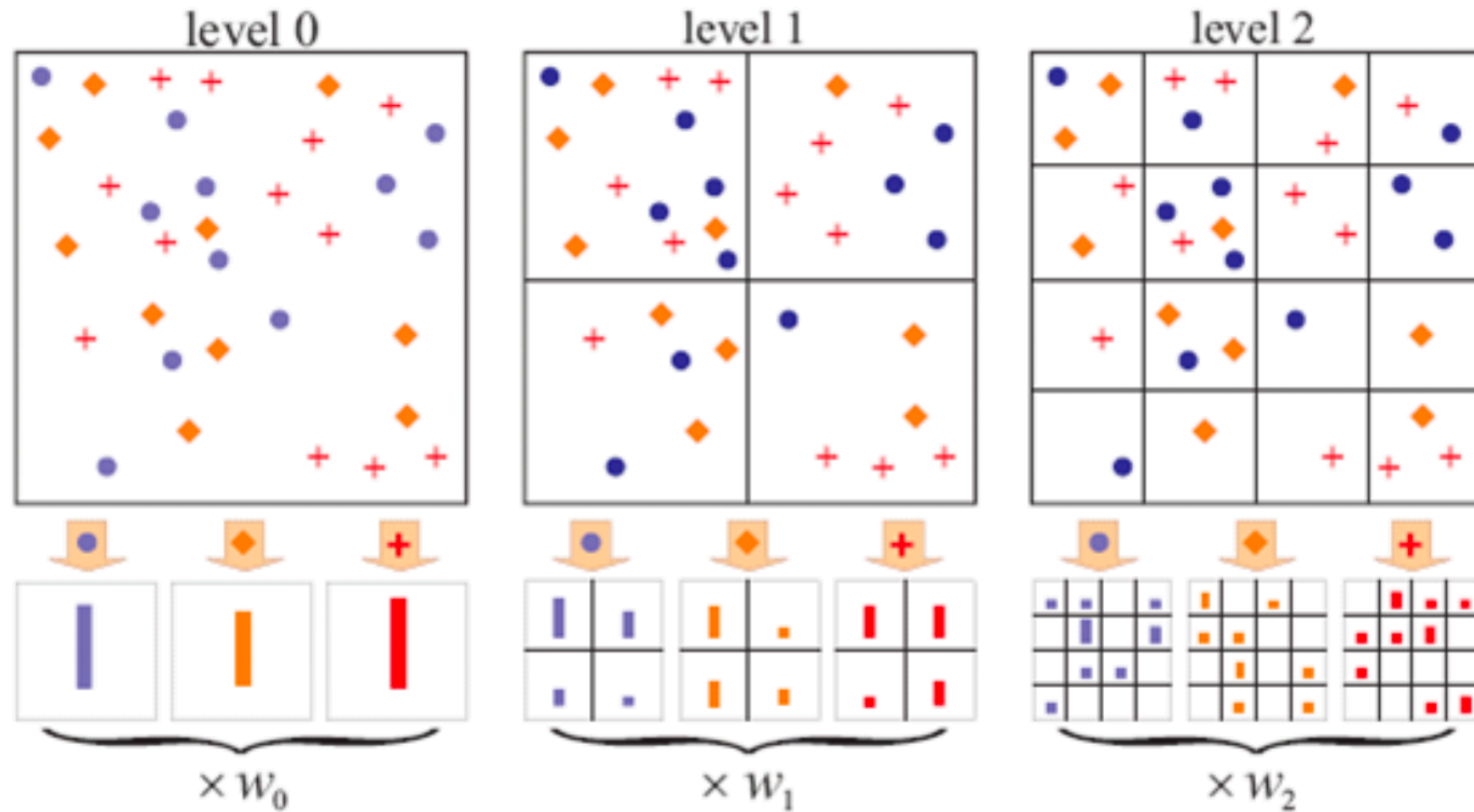


Fig. 16.8 in Forsyth & Ponce (2nd ed.).  
Original credit: Lazebnik et al., 2006

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

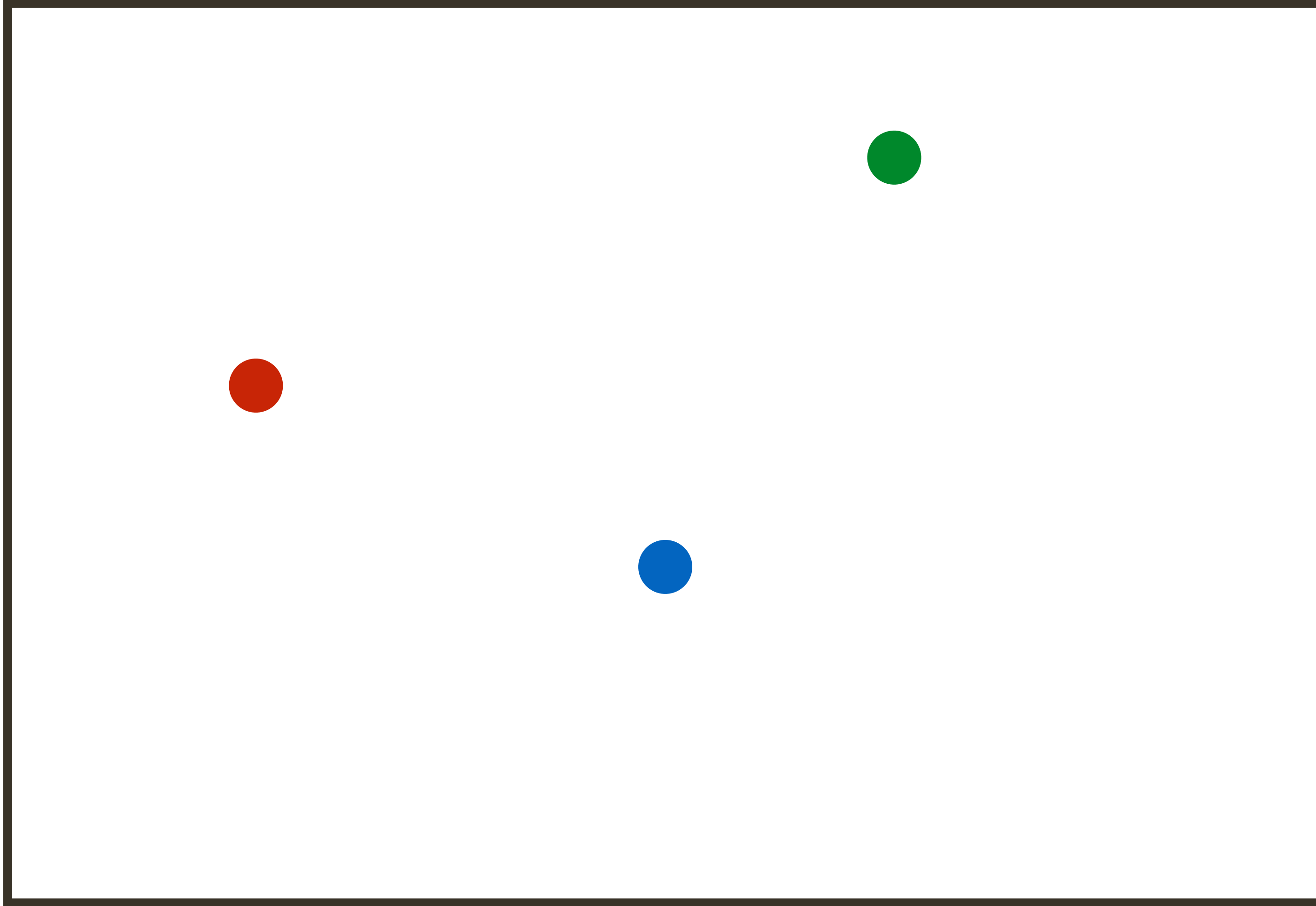
There are more advanced ways to ‘count’ visual words than incrementing its histogram bin

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

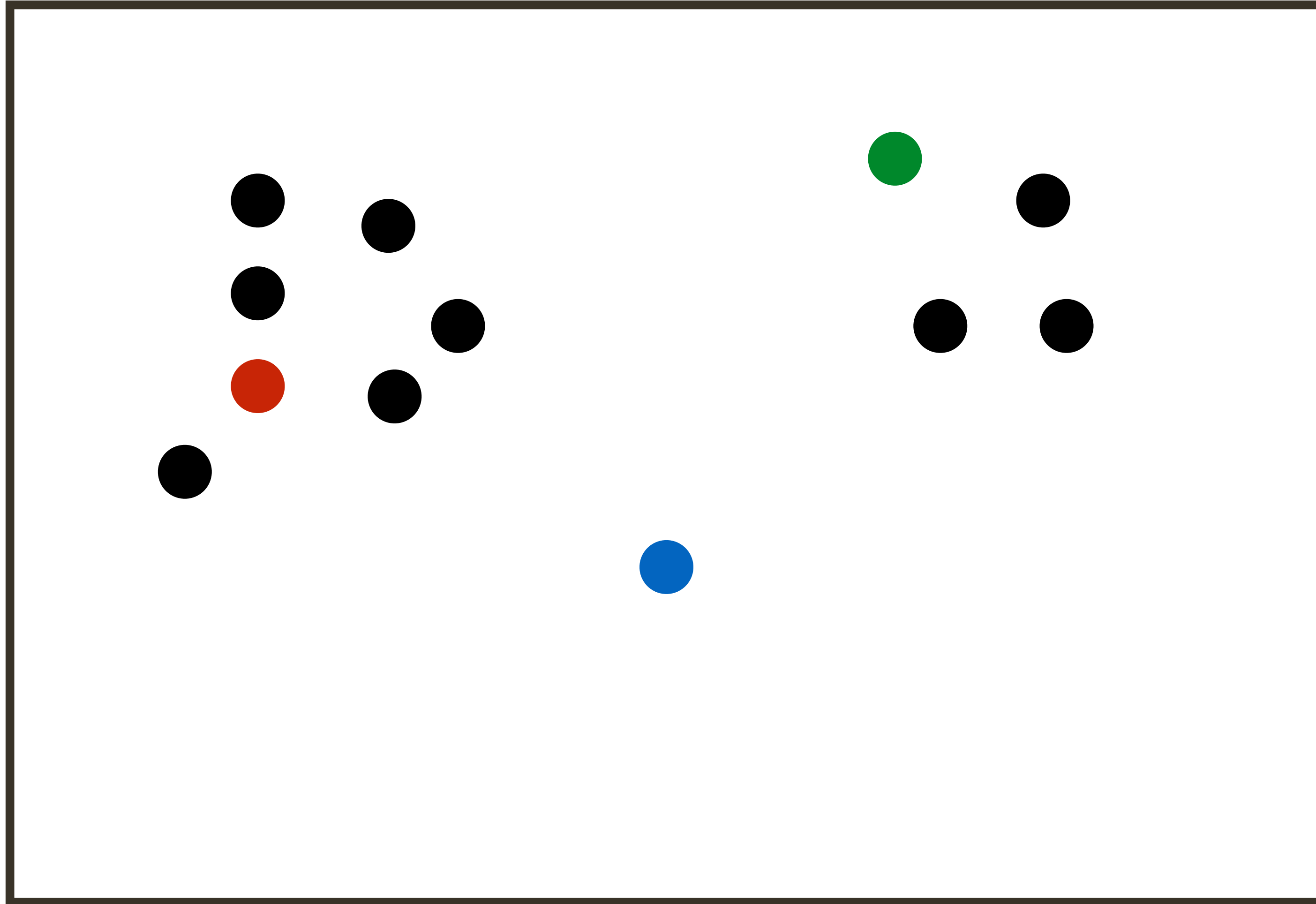
In the VLAD representation, instead of incrementing the histogram bin by one, we increment it by the **residual** vector  $\mathbf{x} - \mathbf{c}(\mathbf{x})$



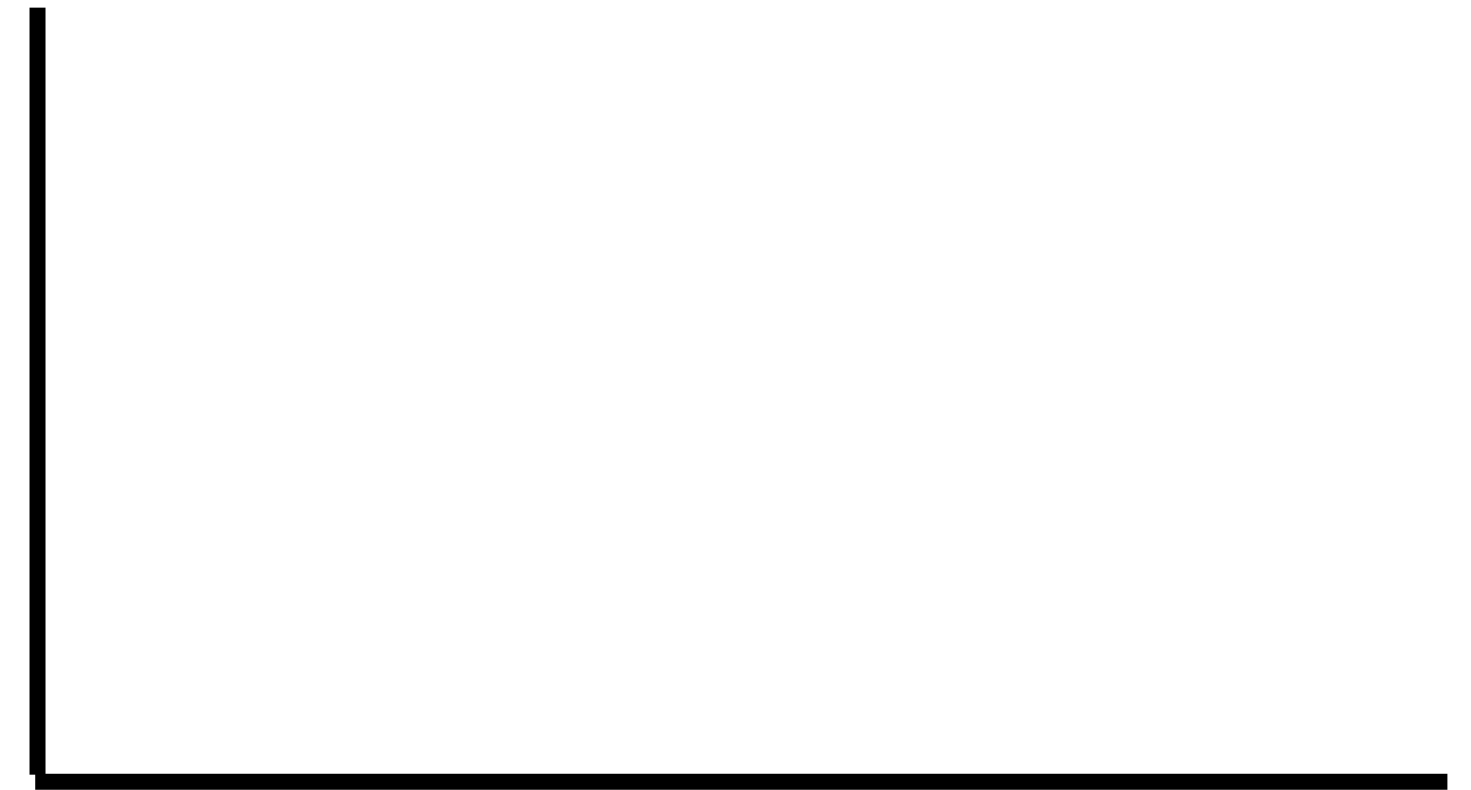
# Example: VLAD



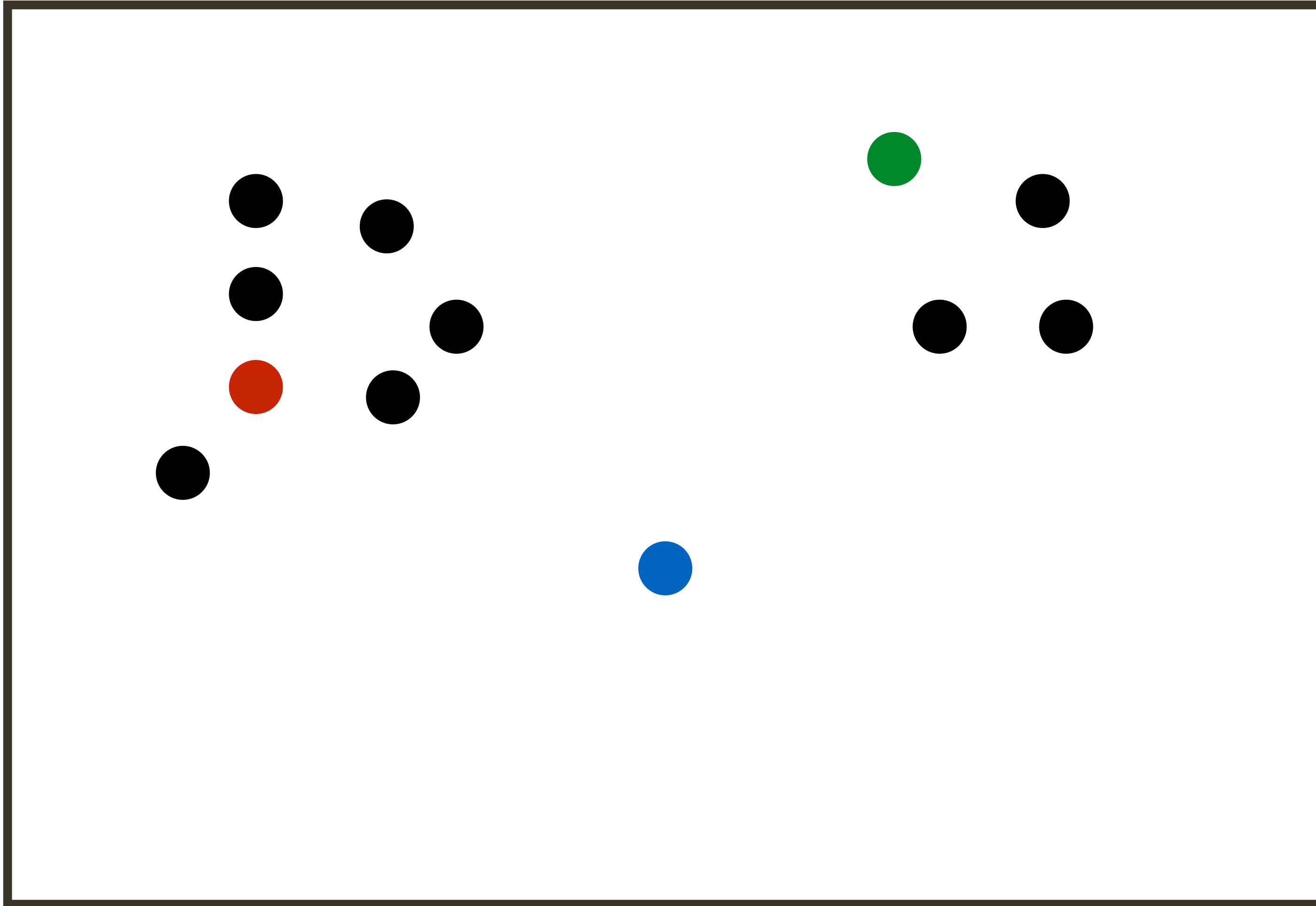
# Example: VLAD



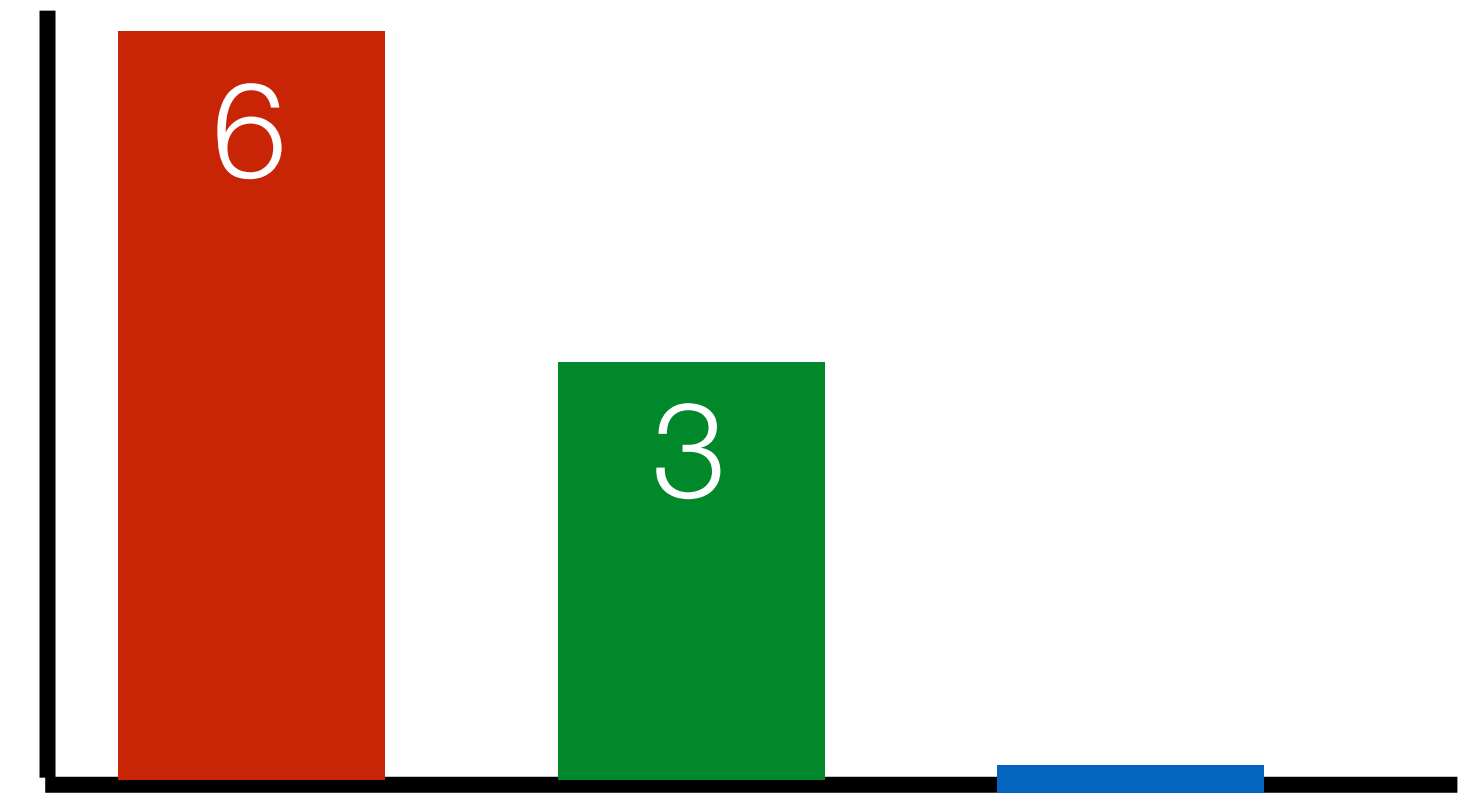
Bag of Word



# Example: VLAD

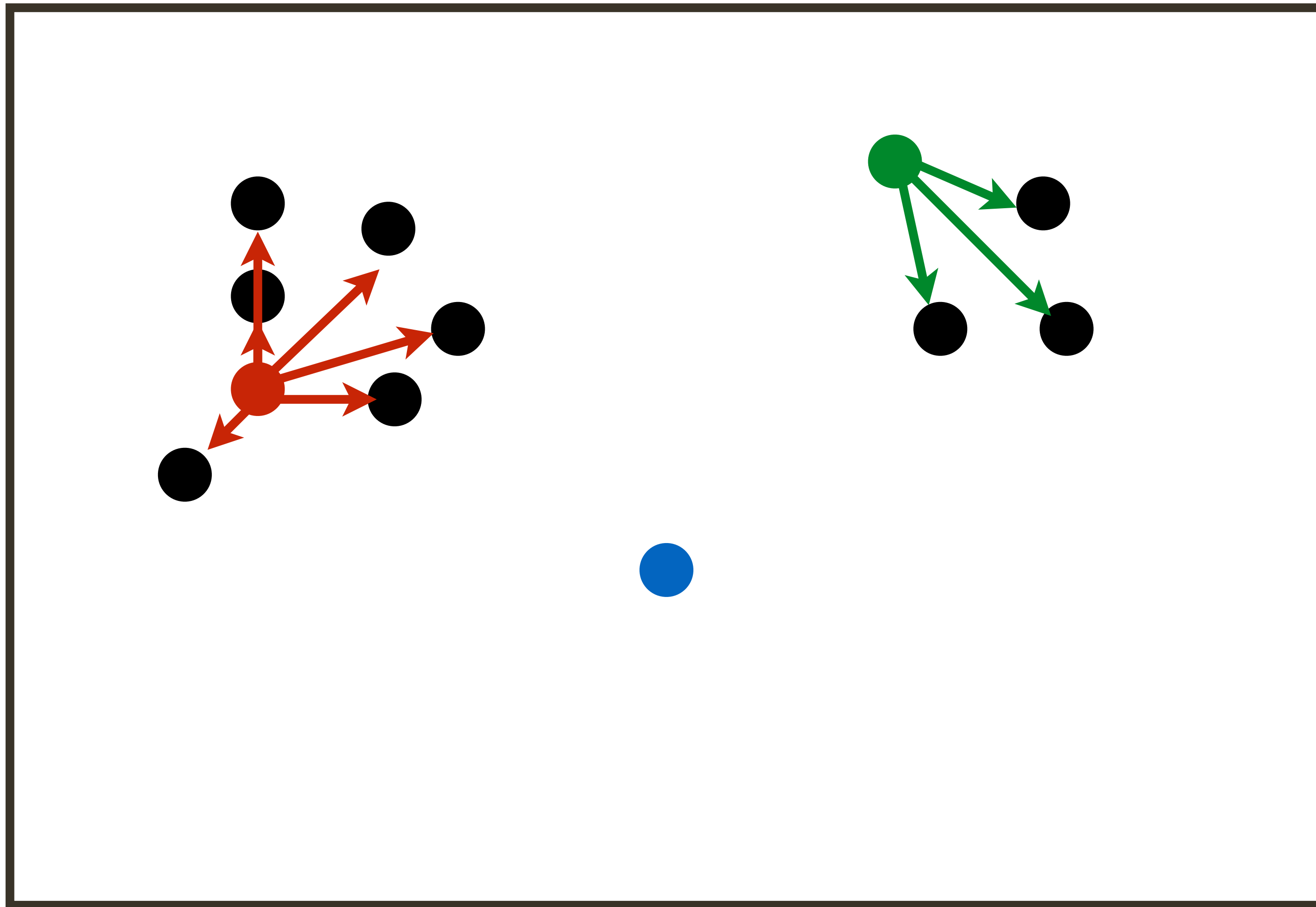


Bag of Word

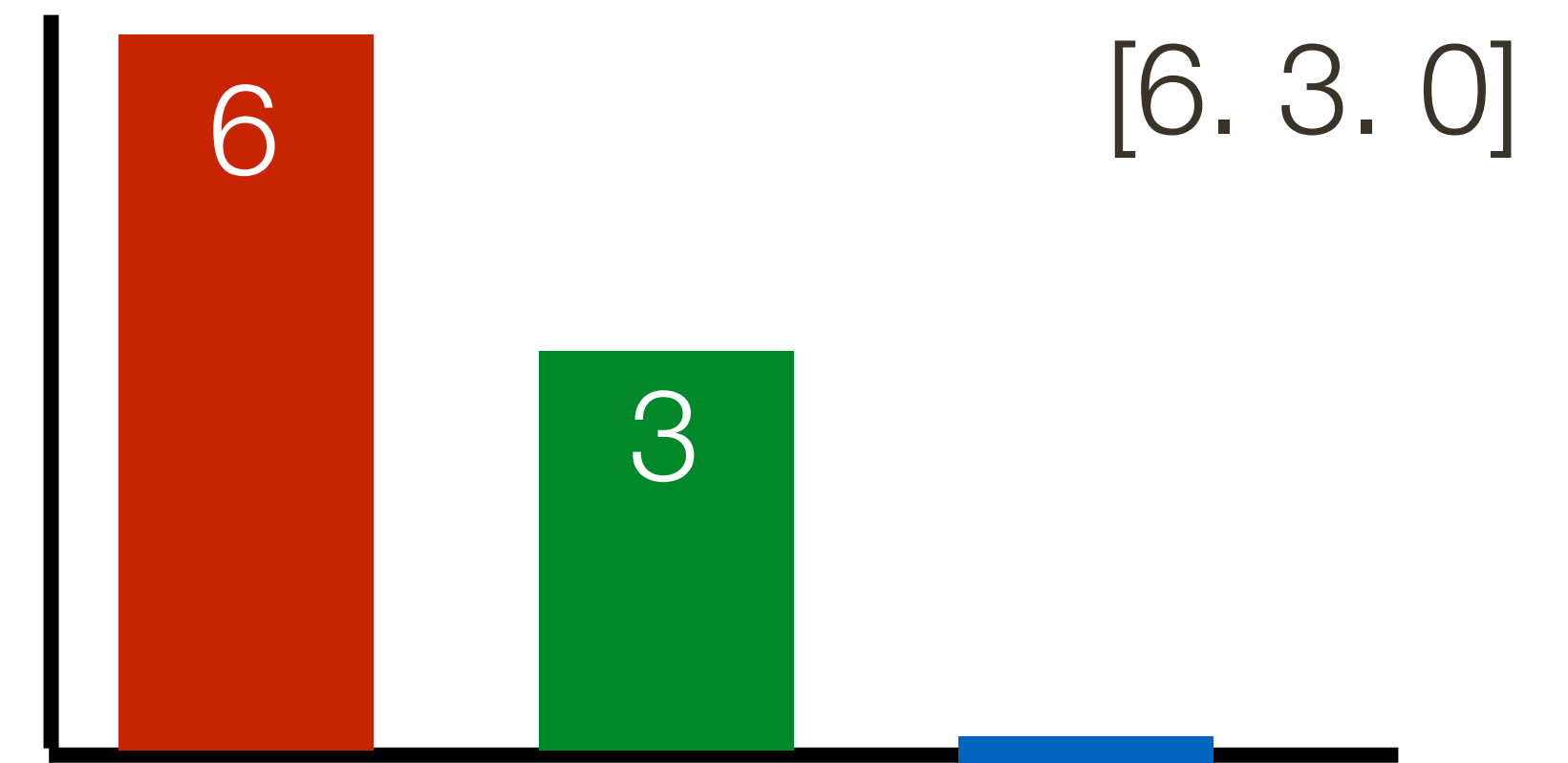




# Example: VLAD

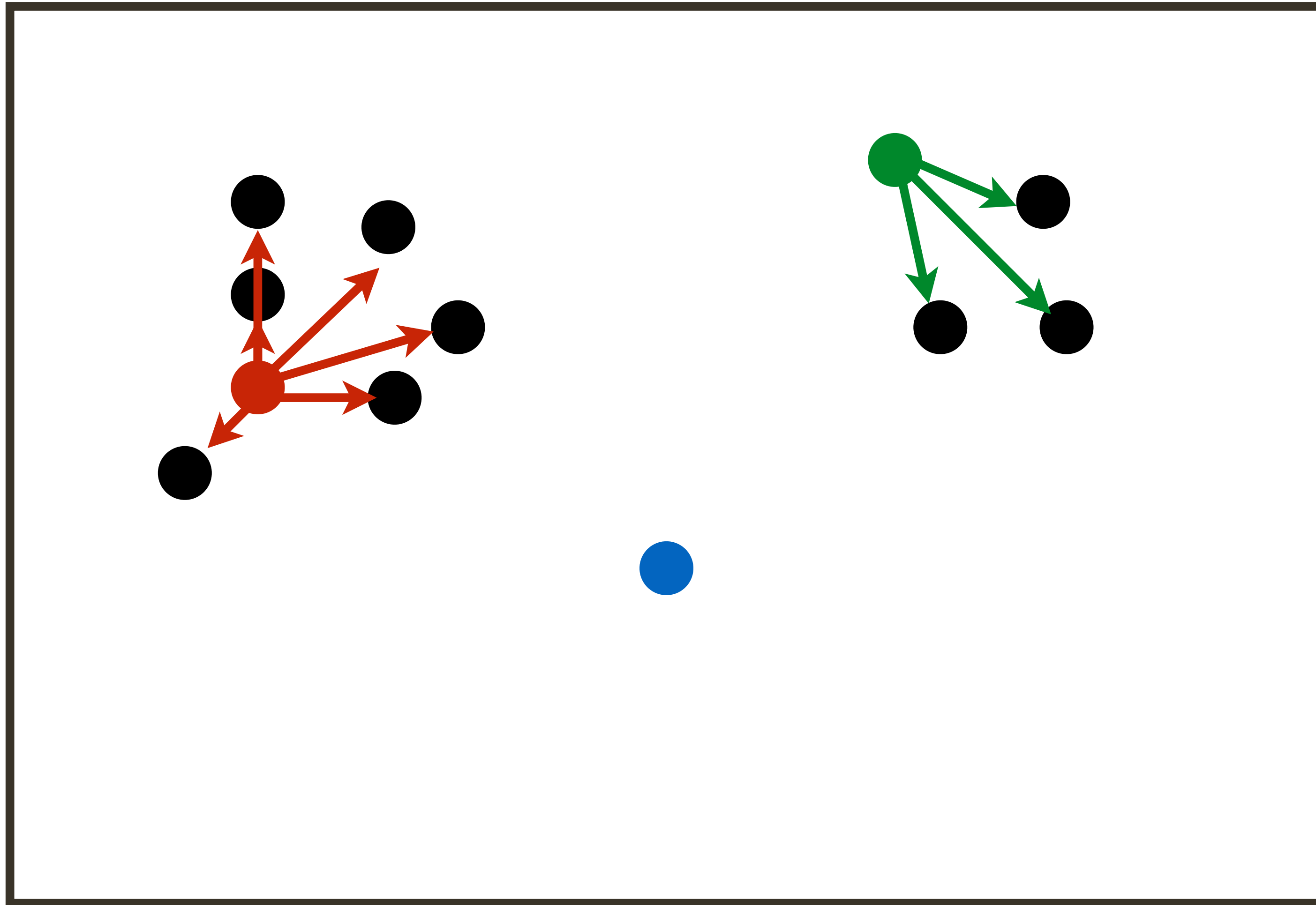


Bag of Word

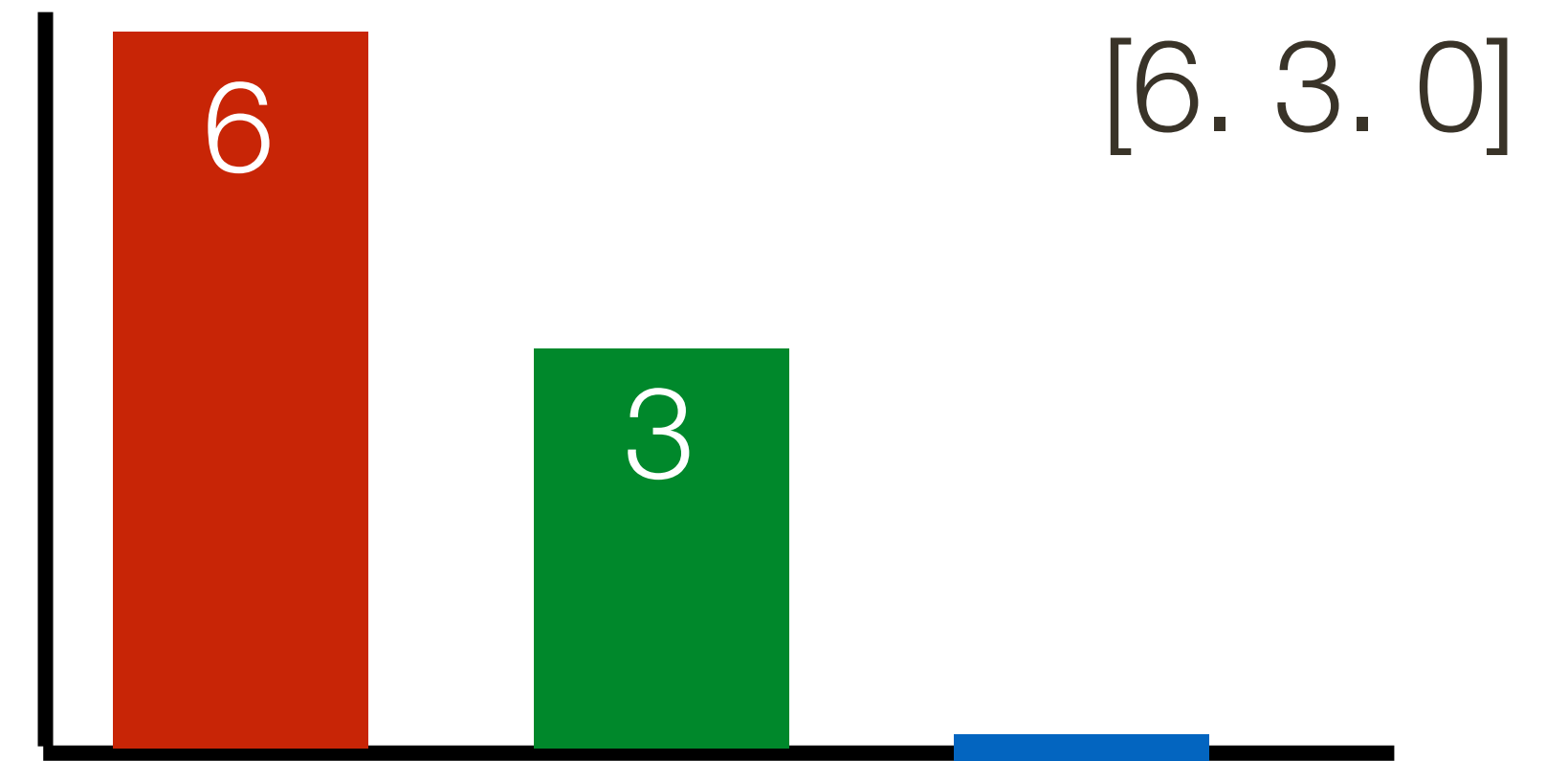


VLAD

# 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