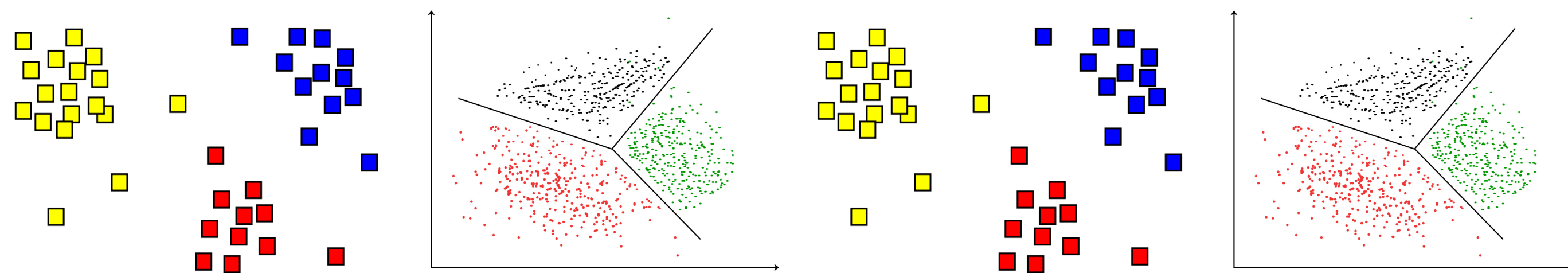


CPSC 425: Computer Vision



Lecture 22: Image Classification (cont.)

Menu for Today (March 26, 2019)

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 **April 4**

Today's “**fun**” Example:

Audio-Visual Scene Analysis with Self-Supervised Multisensory Features

Andrew Owens Alexei A. Efros
UC Berkeley



Lecture 21: Re-cap

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

Dictionary Learning:

Learn Visual Words using clustering

Encode:

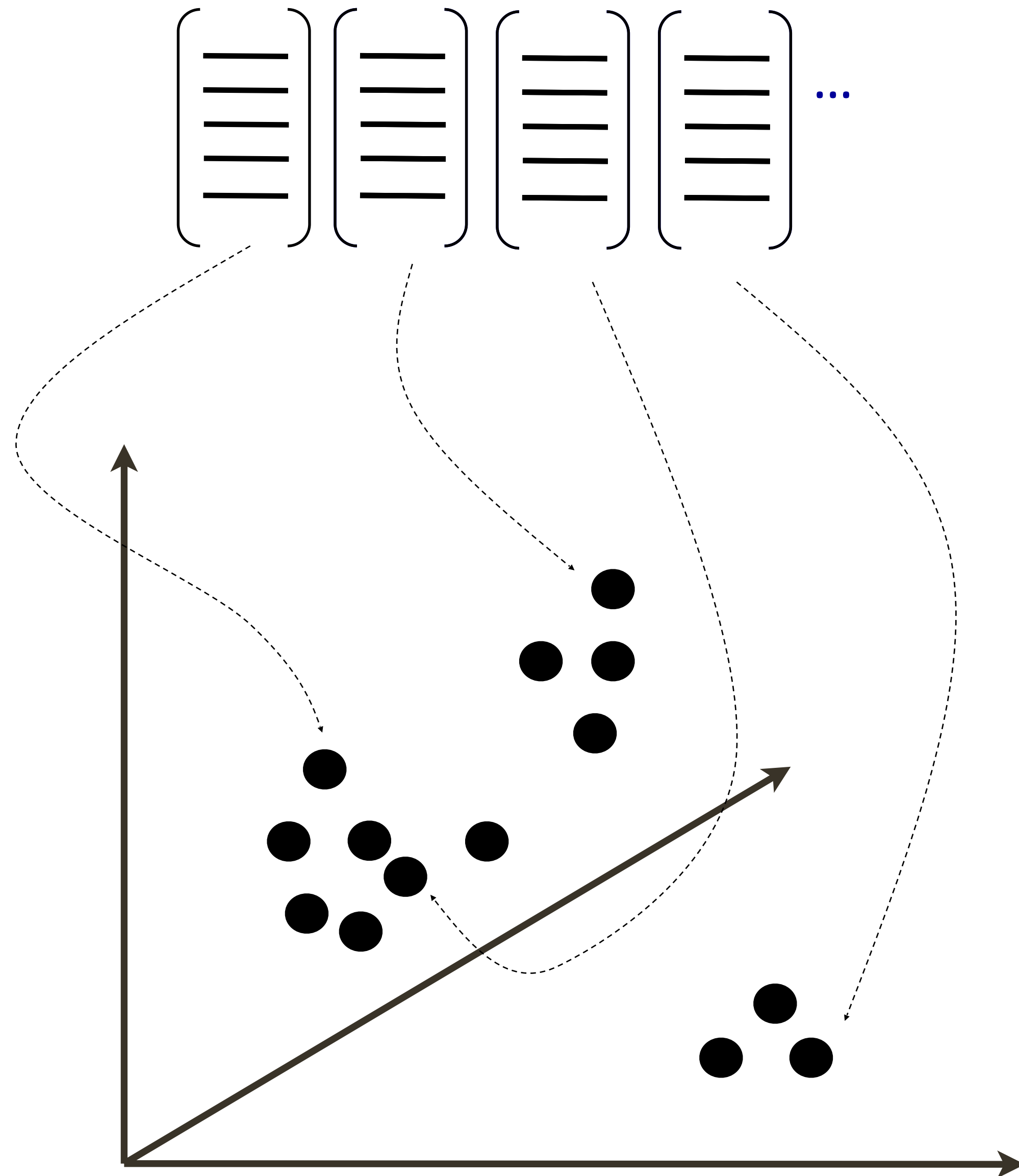
build Bags-of-Words (BOW) vectors
for each image

Classify:

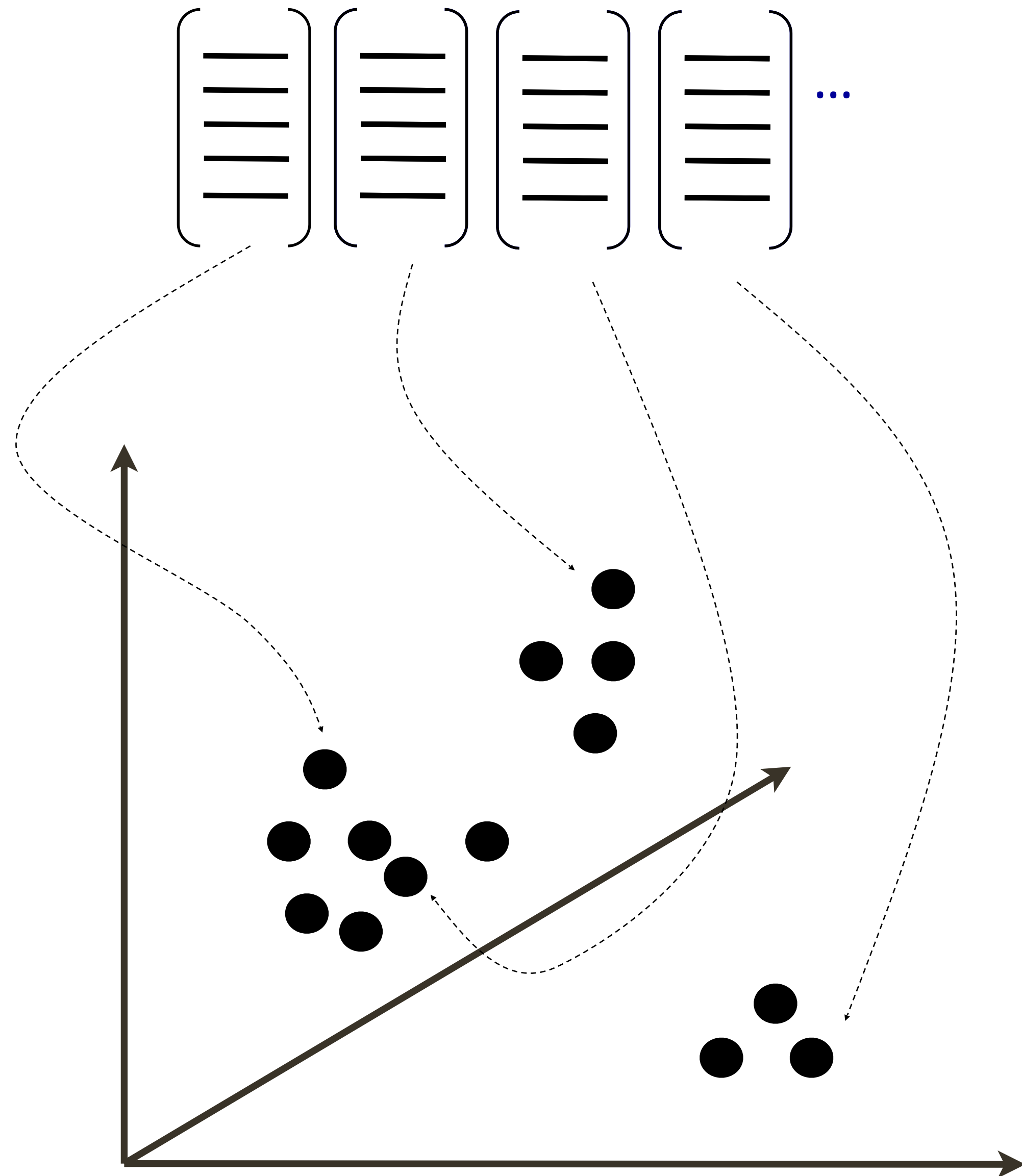
Train and test data using BOWs

Lecture 21: Re-cap

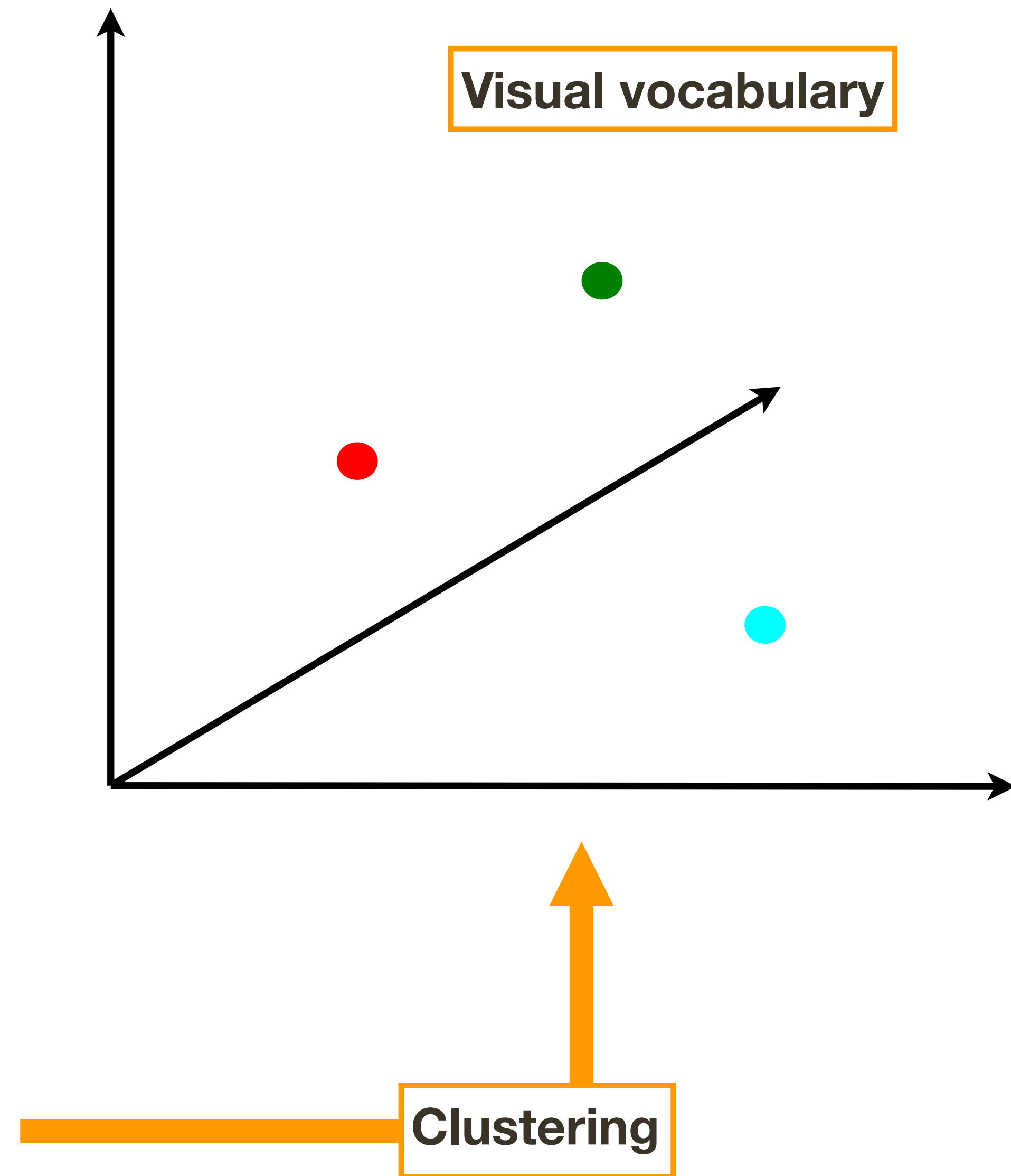
Bag-of-Words Representation



Lecture 21: Re-cap



Bag-of-Words Representation

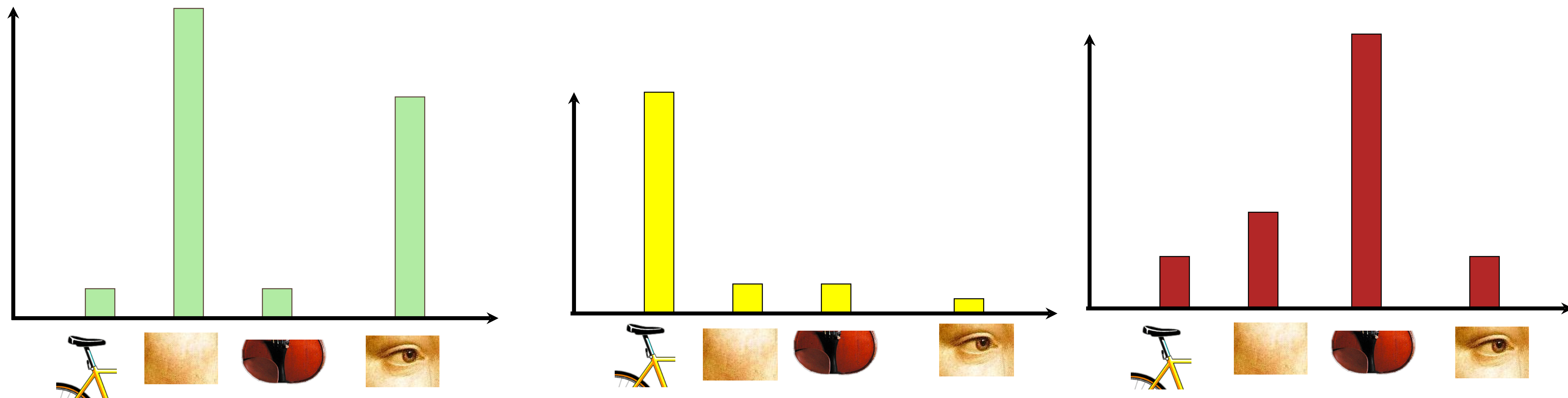


Lecture 21: Re-cap

Bag-of-Words Representation

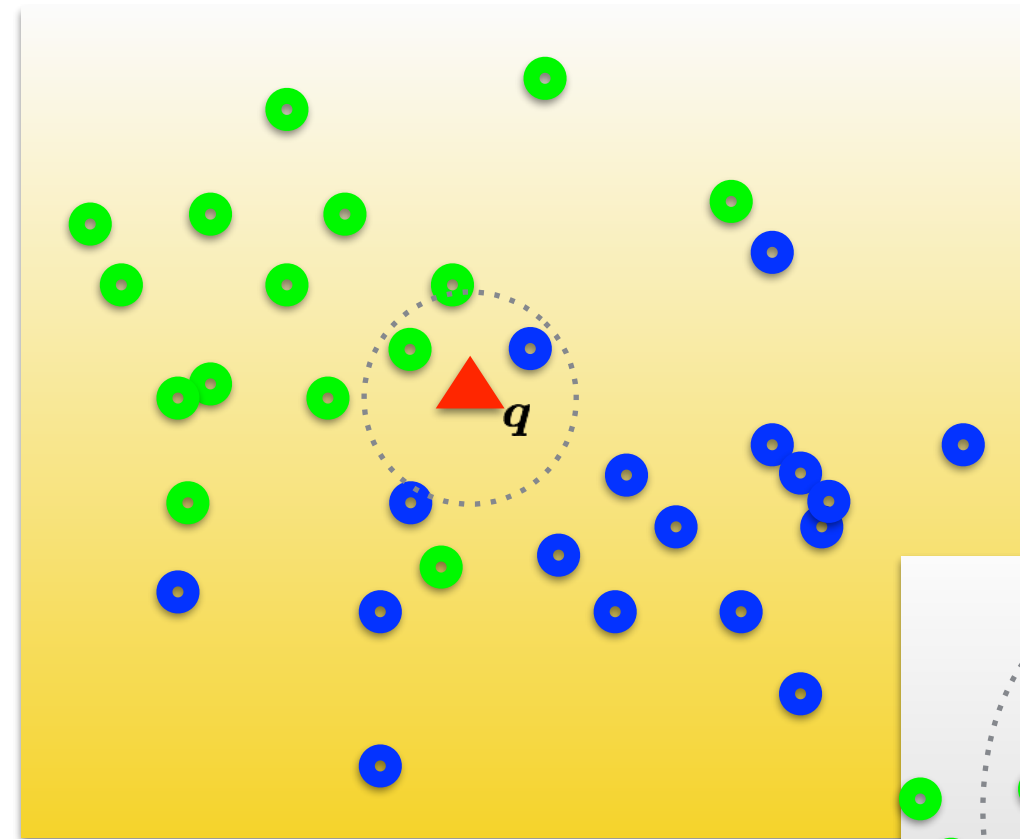
Quantization: image features gets associated to a visual word (nearest cluster center)

Histogram: count the number of visual word occurrences

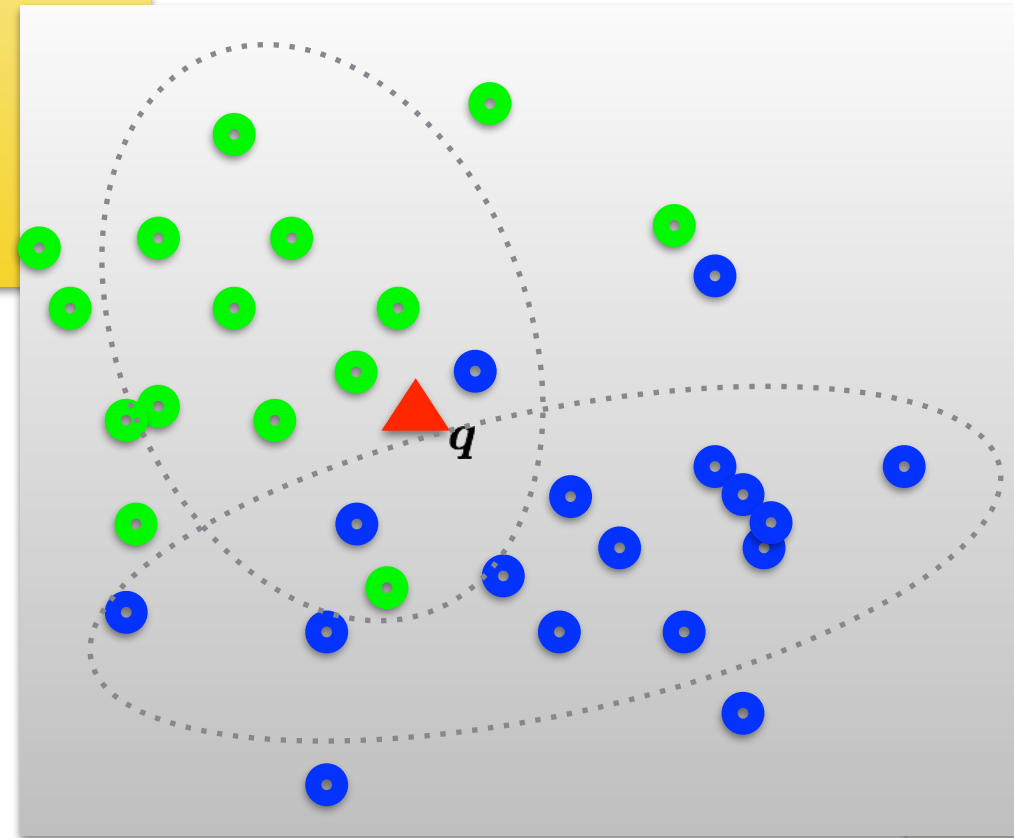


Lecture 21: Re-cap

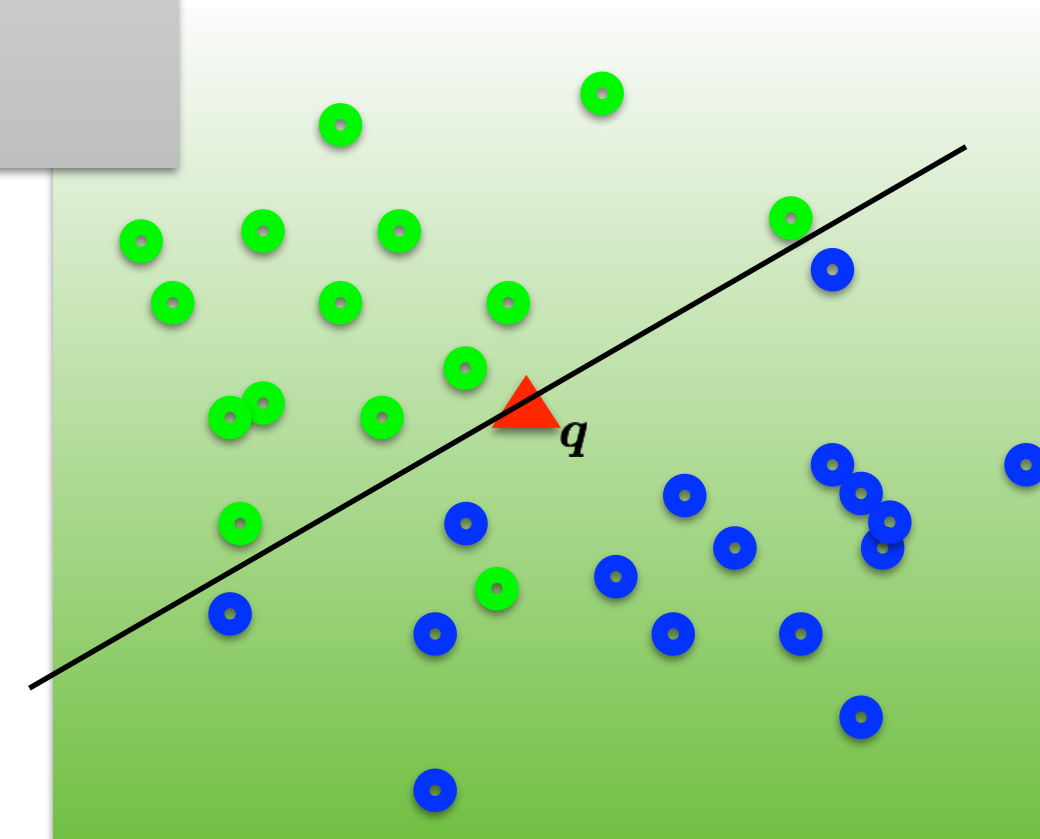
Bag-of-Words Representation



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

Please get your **iClickers** — Quiz

Quiz 5, Question 1

Suppose we have a codebook of 1,000 SIFT visual words (recall that SIFT is a 128-dimensional feature). Now we are given a new image and we extract 2,000 SIFT descriptors from it. What is the dimensionality of a bag of words descriptor?

- A) 128
- B) 1,000
- C) 2,000
- D) 128,000
- E) It is not possible to construct a bag of words because there are more SIFT descriptors in the image than visual words

Quiz 5, Question 1

Suppose we have a codebook of 1,000 SIFT visual words (recall that SIFT is a 128-dimensional feature). Now we are given a new image and we extract 2,000 SIFT descriptors from it. What is the dimensionality of a bag of words descriptor?

A) 128

B) 1,000

C) 2,000

D) 128,000

E) It is not possible to construct a bag of words because there are more SIFT descriptors in the image than visual words

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

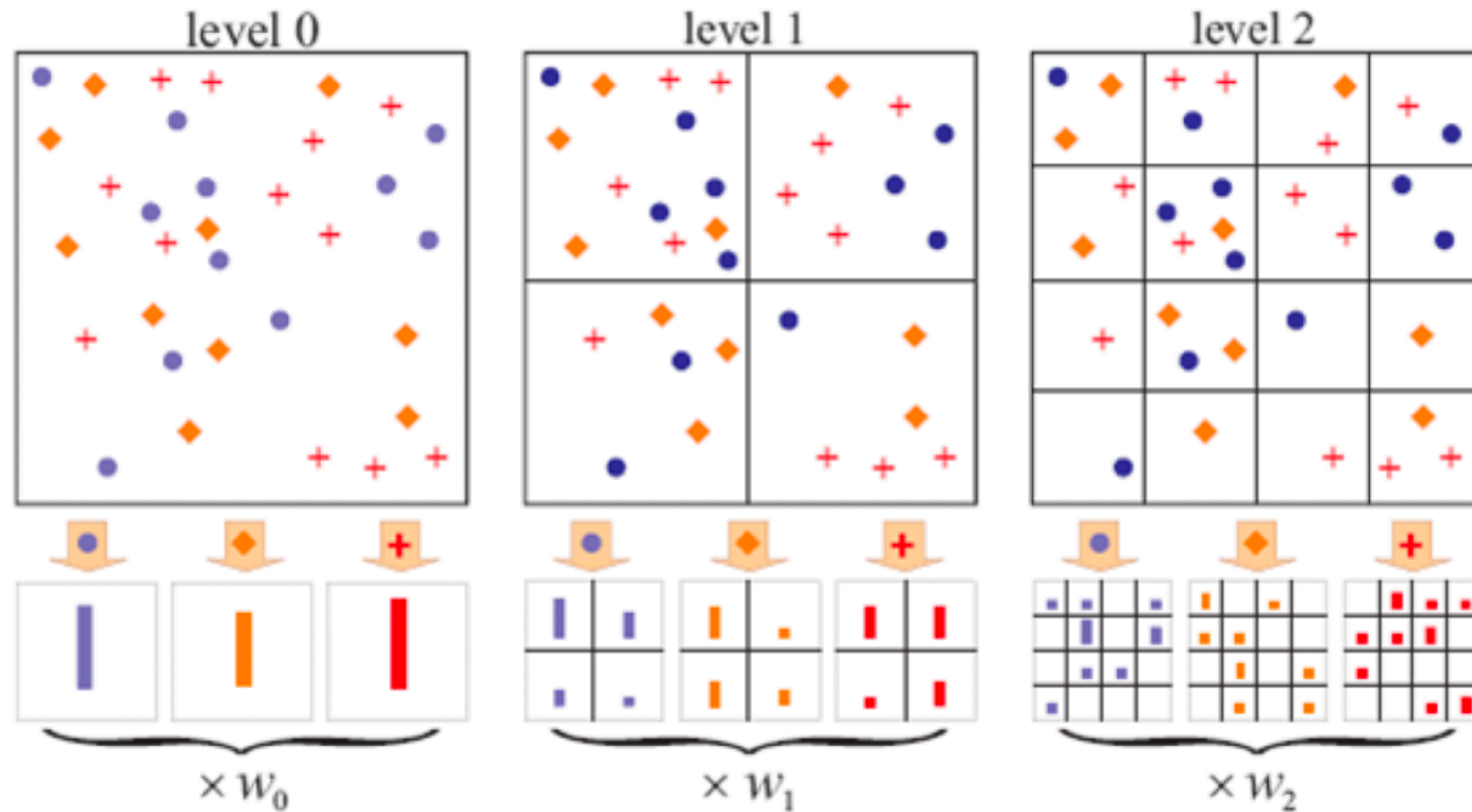


Fig. 16.8 in Forsyth & Ponce (2nd ed.).
Original credit: Lazebnik et al., 2006

Please get your **iClickers** — Quiz

Quiz 5, Question 2

We have a codebook of 1,000 SIFT visual words (recall that SIFT is a 128-dimensional feature). We are given a new image and we extract 2,000 SIFT descriptors from it. What is the dimensionality of a spatial pyramid descriptor with 1x1 and 2x2 grids?

- A) 1,000
- B) 2,000
- C) 5,000**
- D) 10,000
- E) Not enough information to tell - we need to know the spatial locations from which the SIFT descriptors are extracted

Quiz 5, Question 2

We have a codebook of 1,000 SIFT visual words (recall that SIFT is a 128-dimensional feature). We are given a new image and we extract 2,000 SIFT descriptors from it. What is the dimensionality of a spatial pyramid descriptor with 1x1 and 2x2 grids?

A) 1,000

B) 2,000

C) 5,000

D) 10,000

E) Not enough information to tell - we need to know the spatial locations from which the SIFT descriptors are extracted

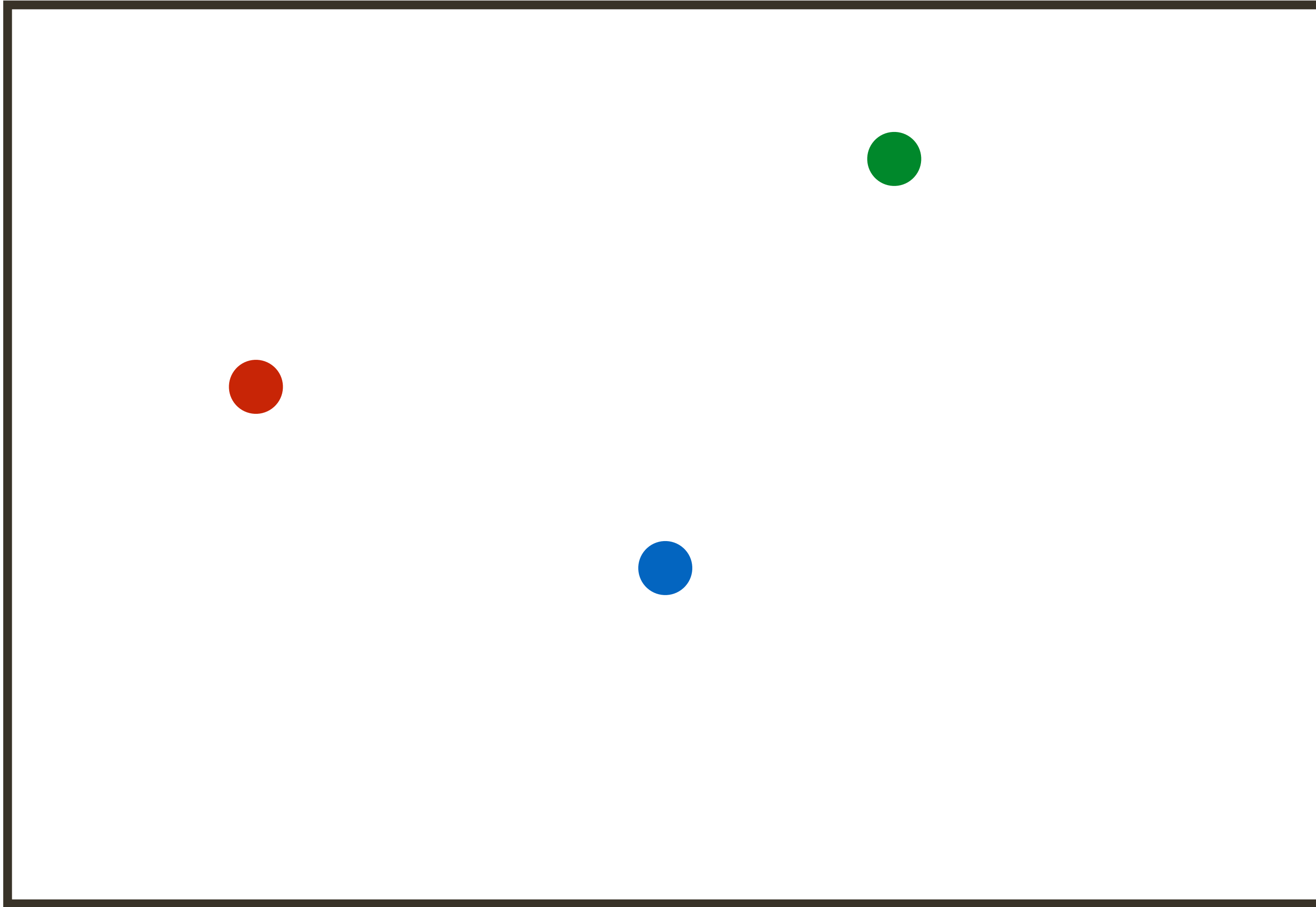
VLAD (Vector of Locally Aggregated Descriptors)

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

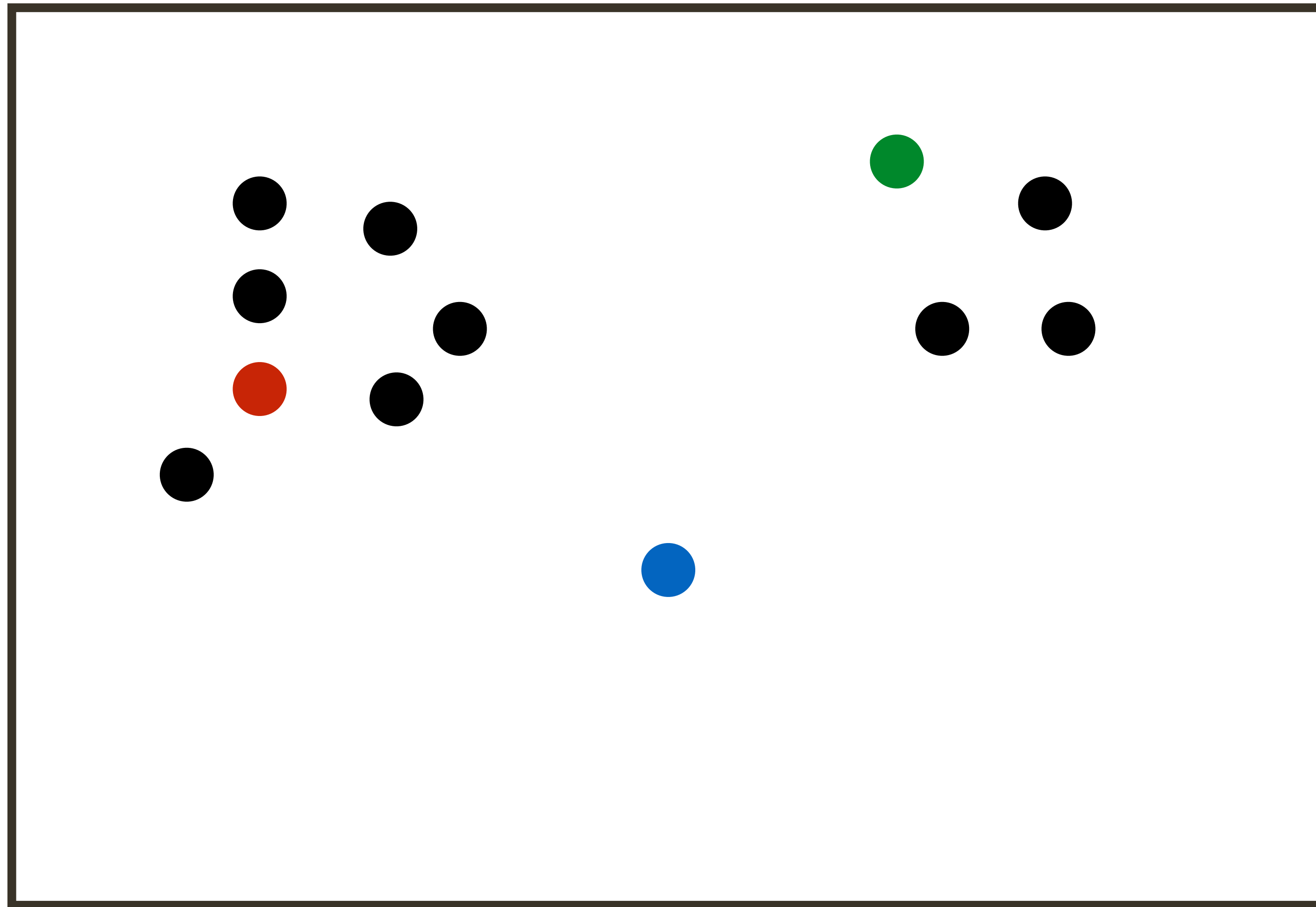
For example, it might be useful to describe how local descriptors are quantized to their visual words

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

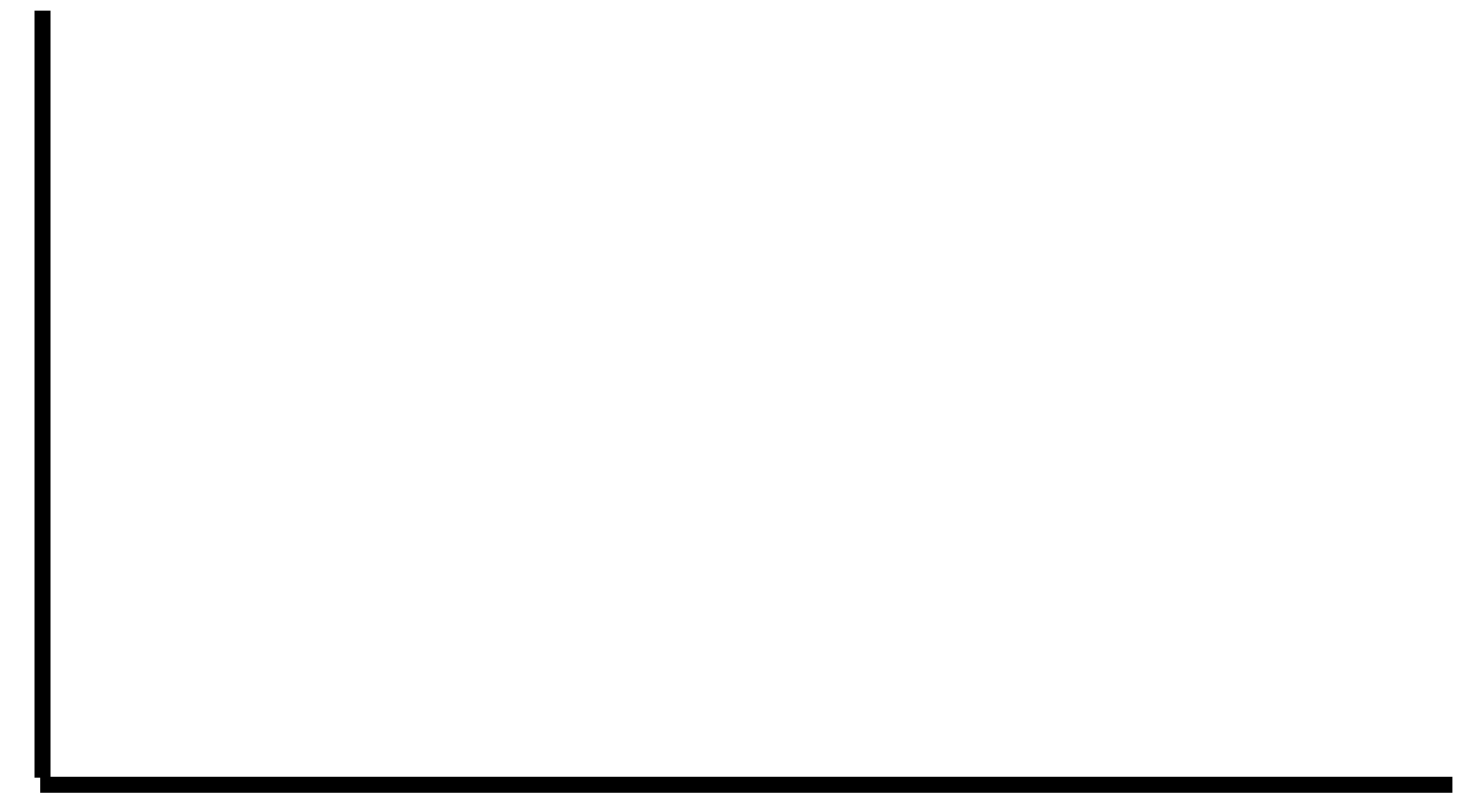
Example: VLAD



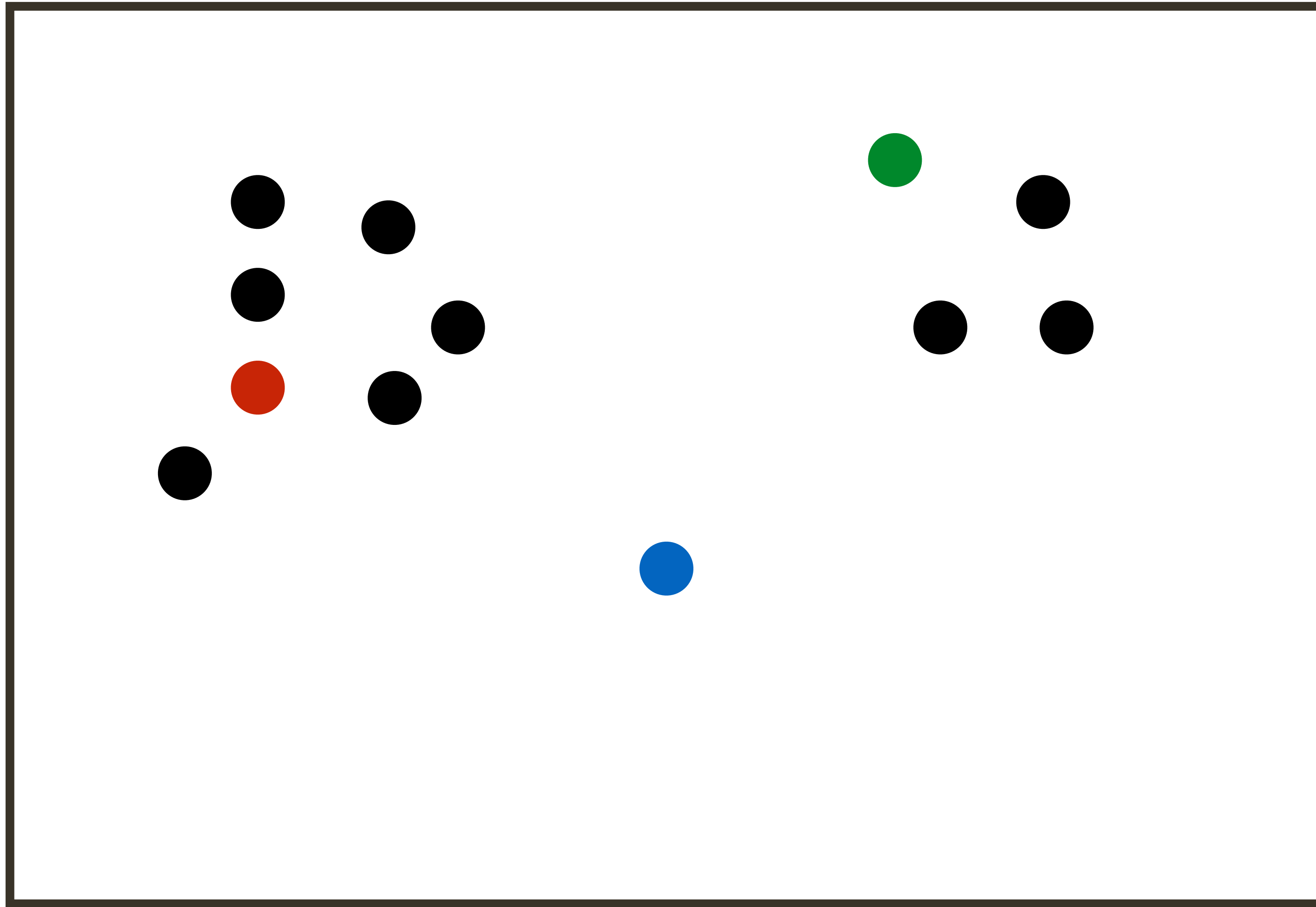
Example: VLAD



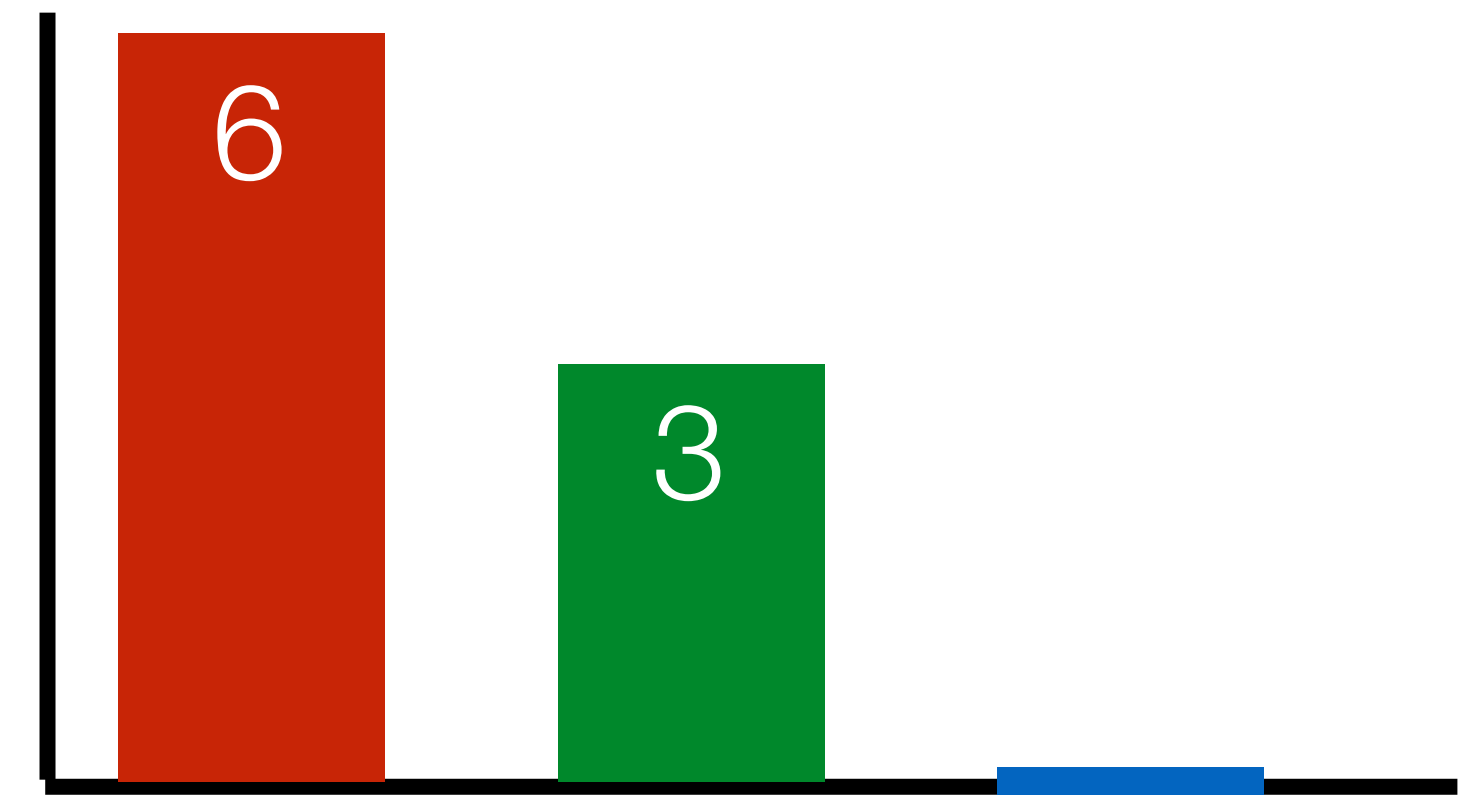
Bag of Word



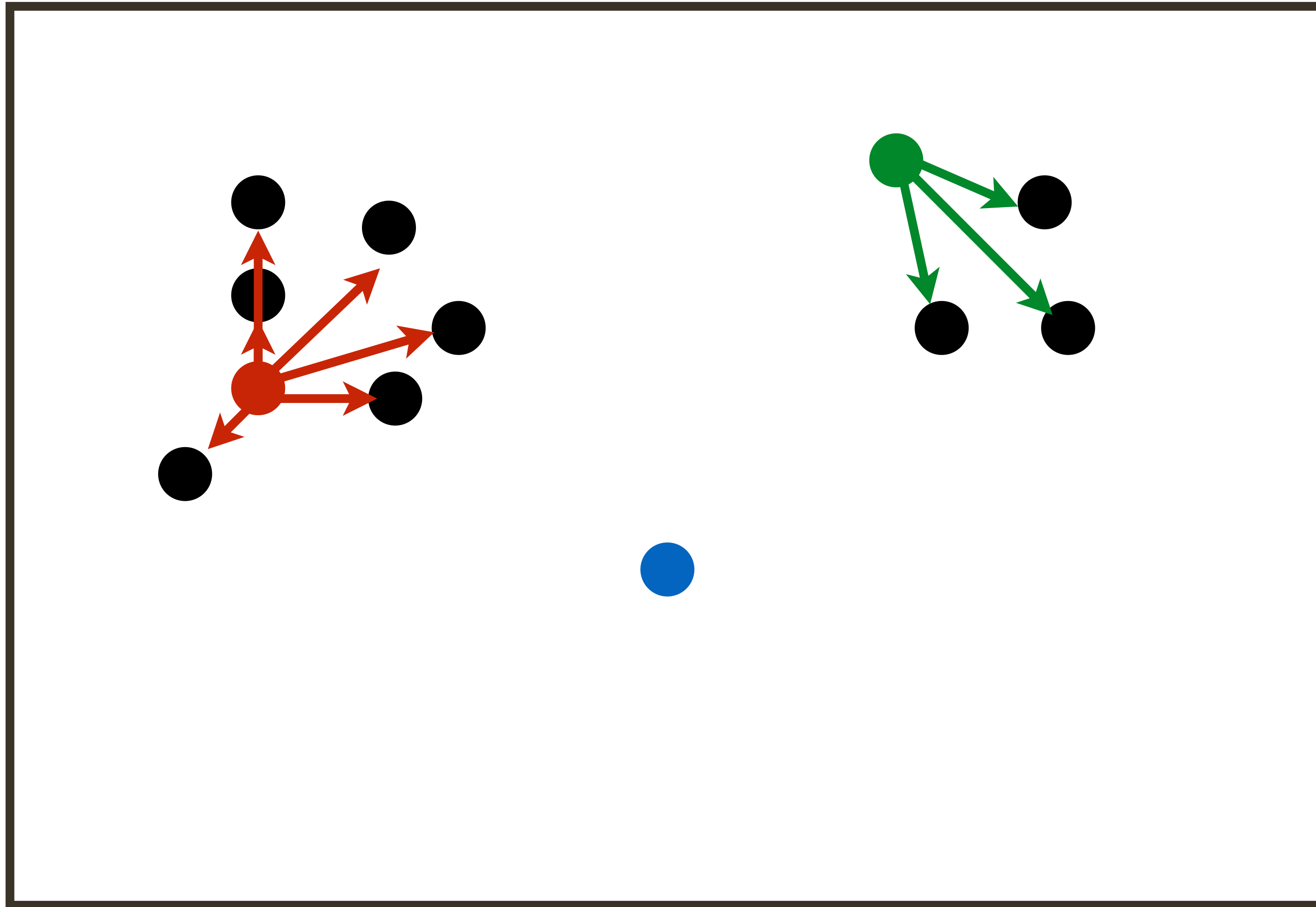
Example: VLAD



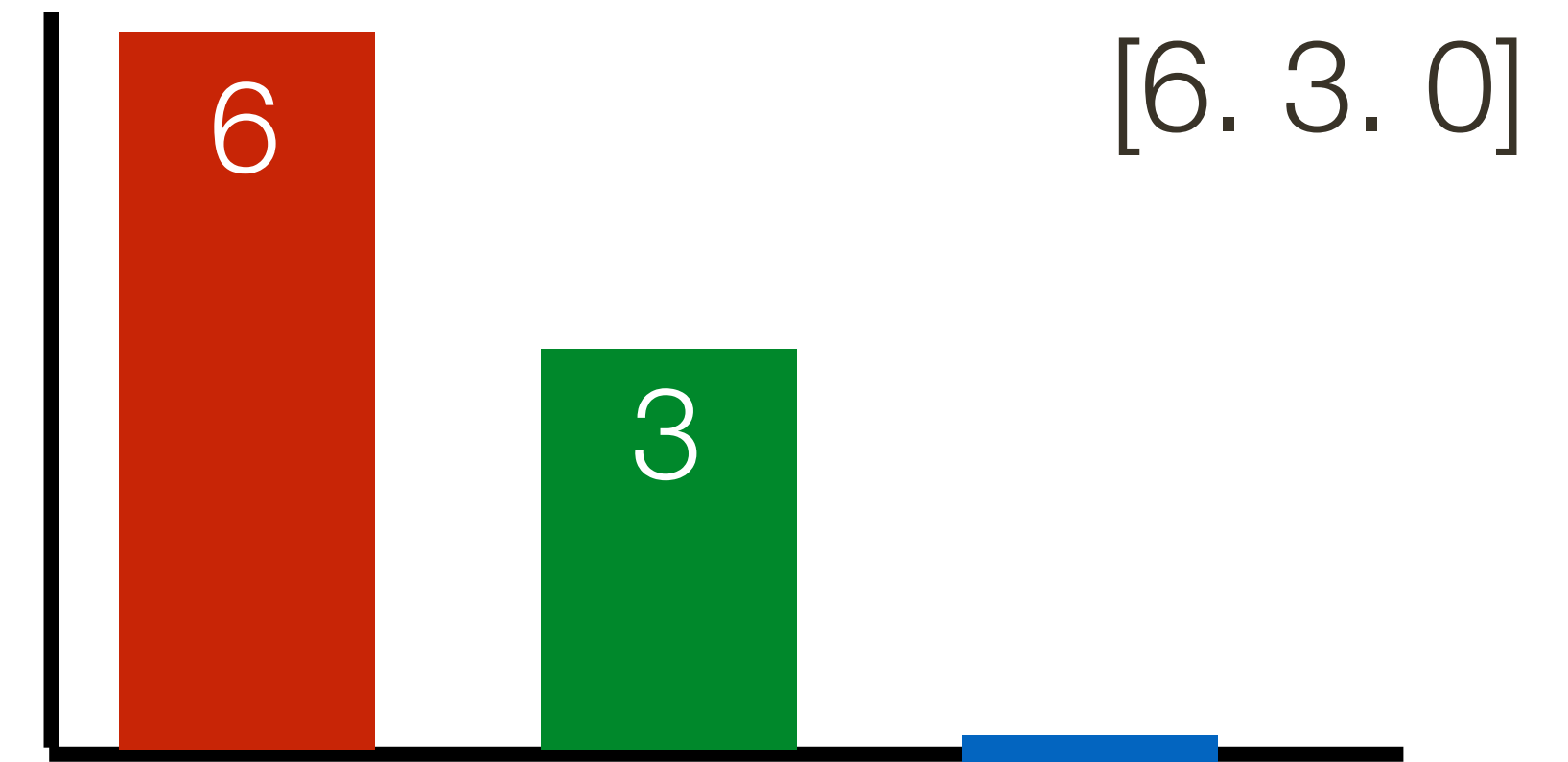
Bag of Word



Example: VLAD

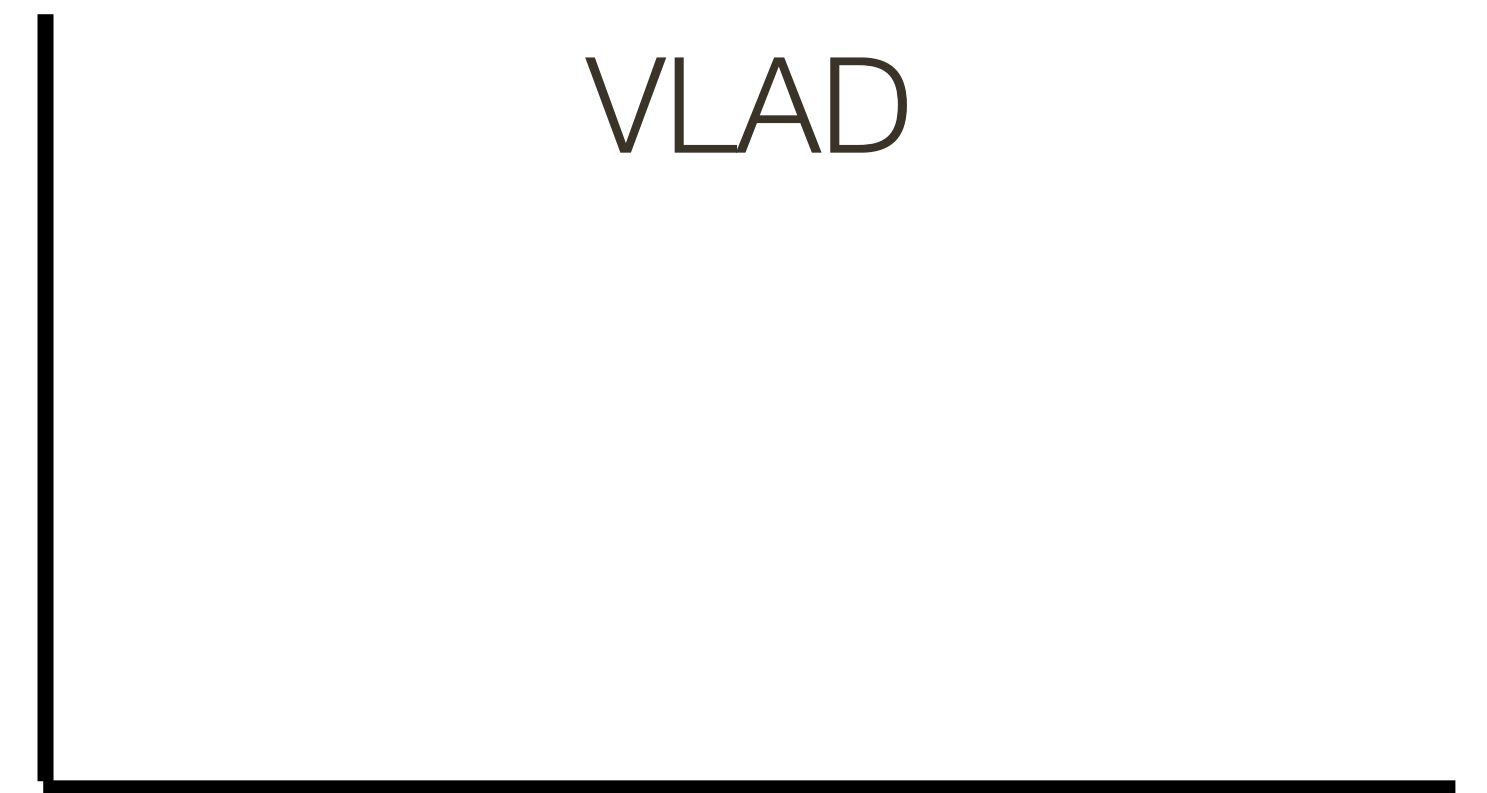


Bag of Word

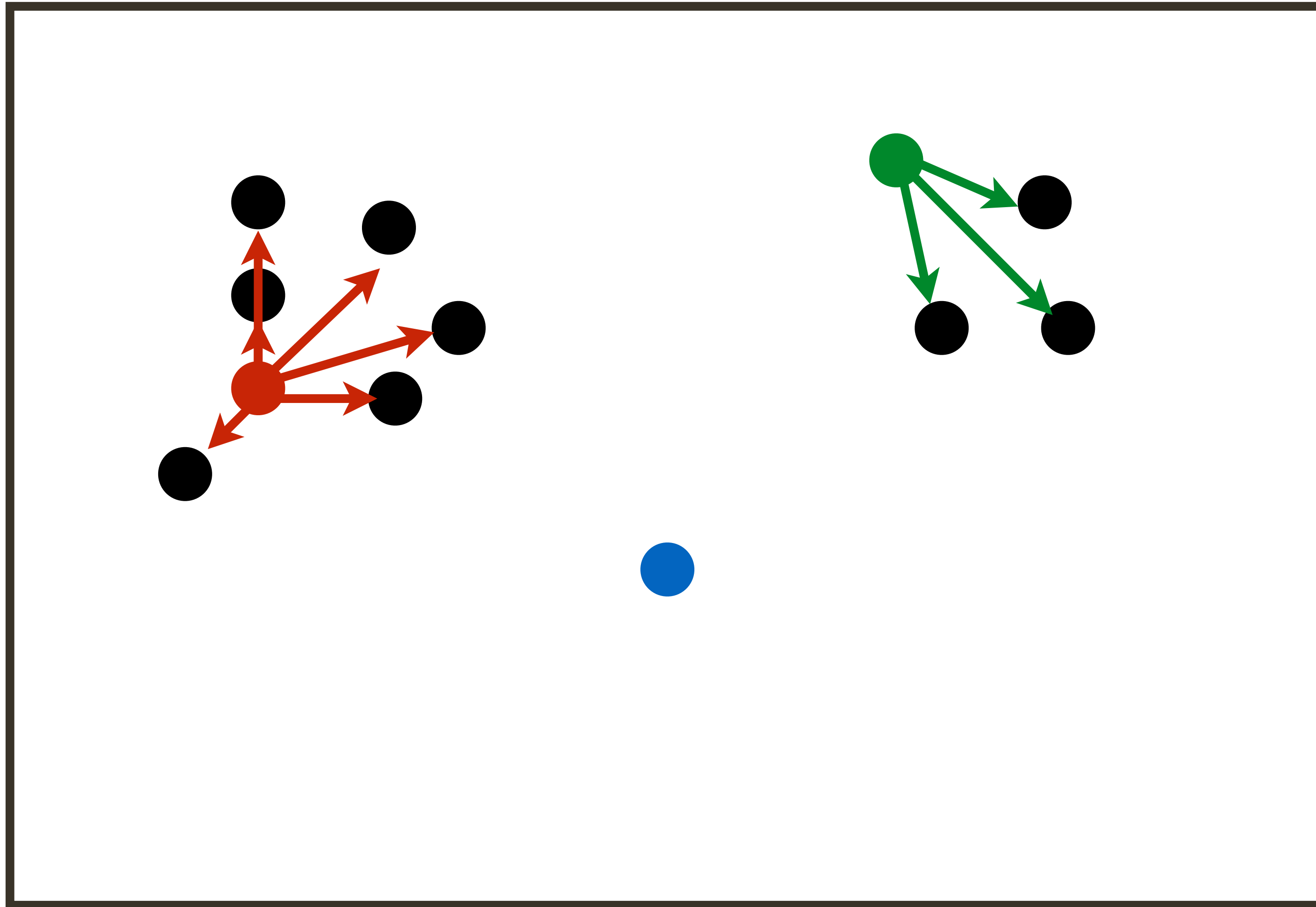


[6. 3. 0]

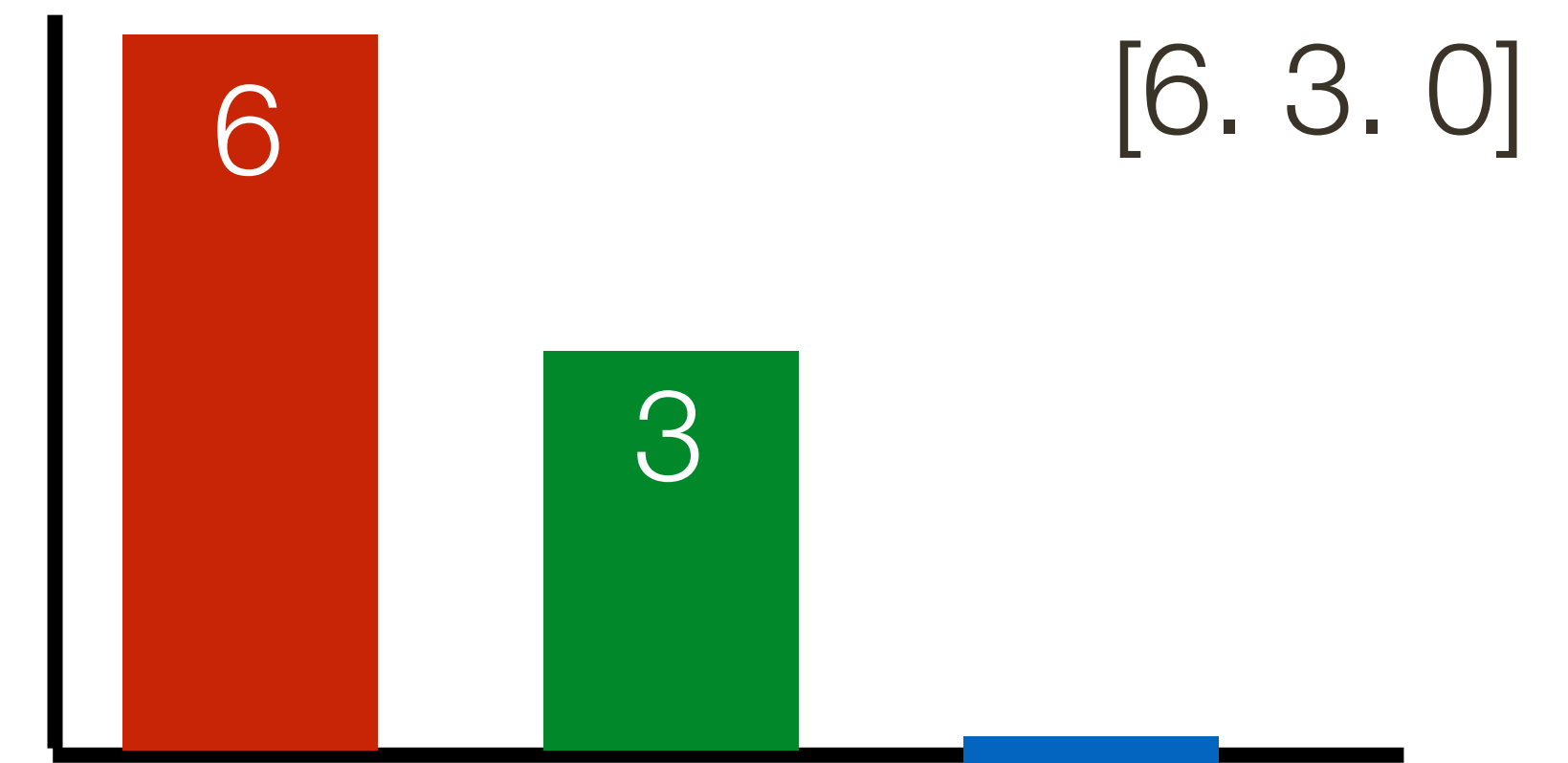
VLAD



Example: VLAD



Bag of Word



[6. 3. 0]

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>>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>>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?

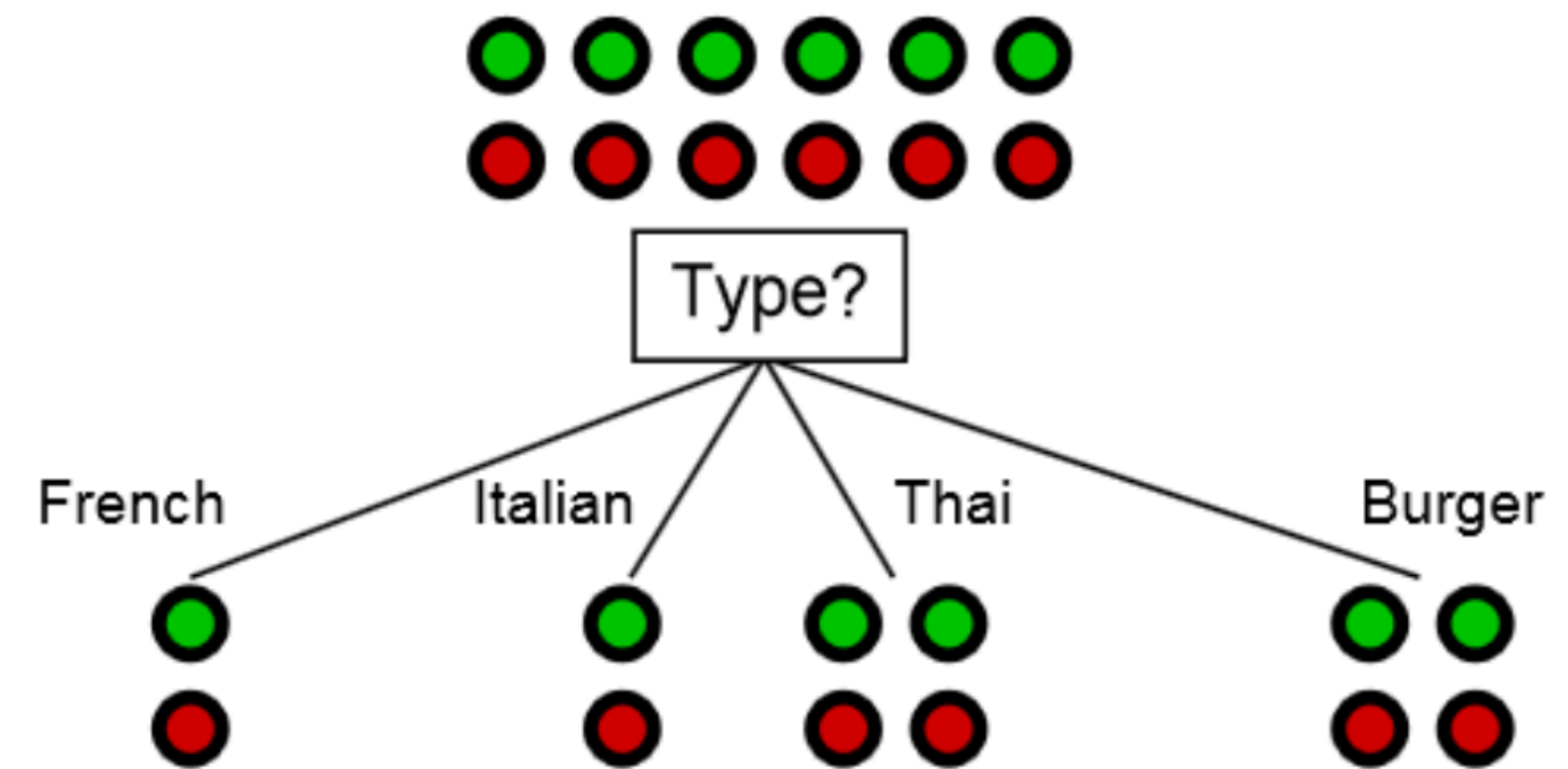
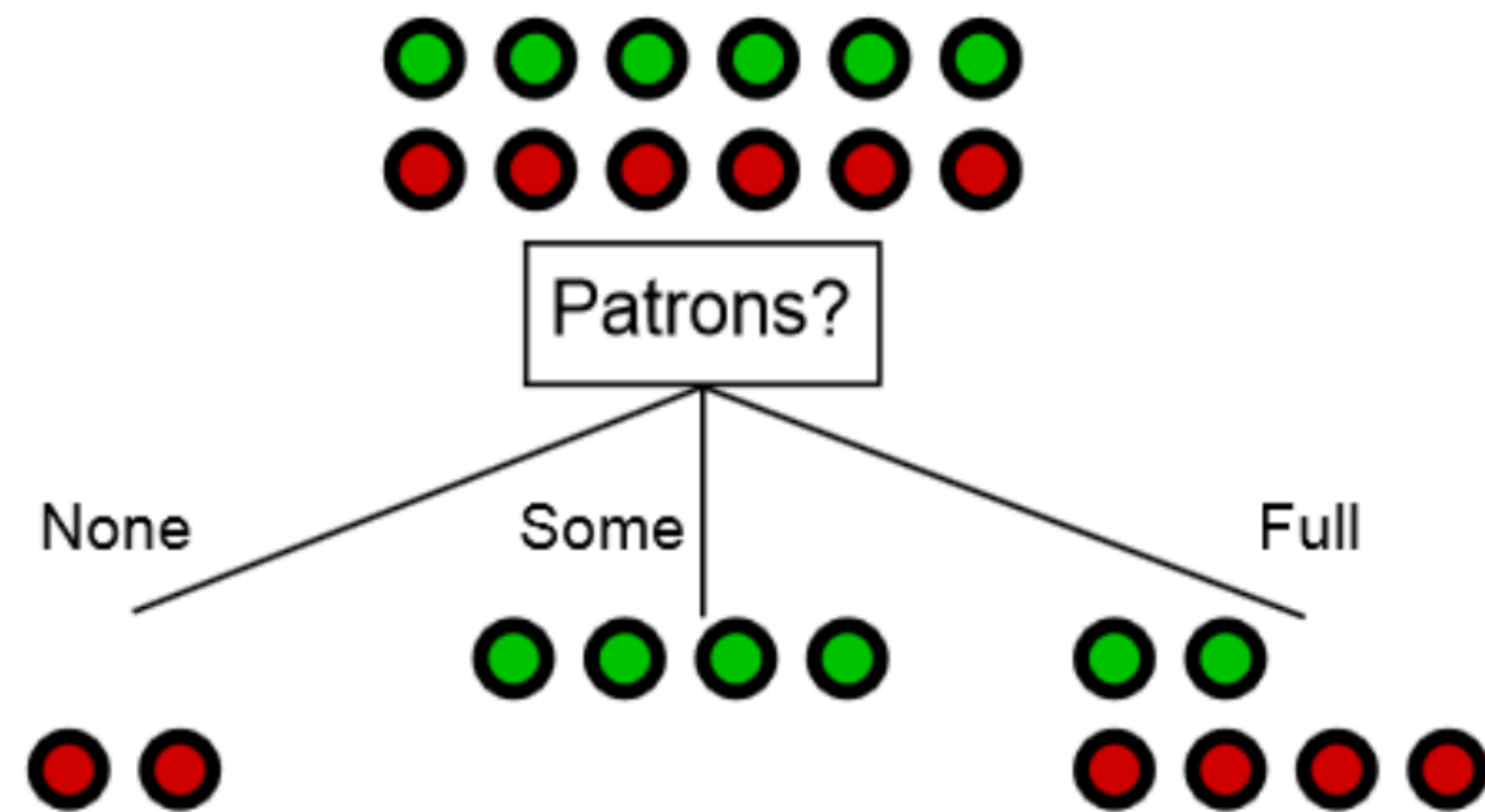


Figure credit: Russell and Norvig (3rd ed.)

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:

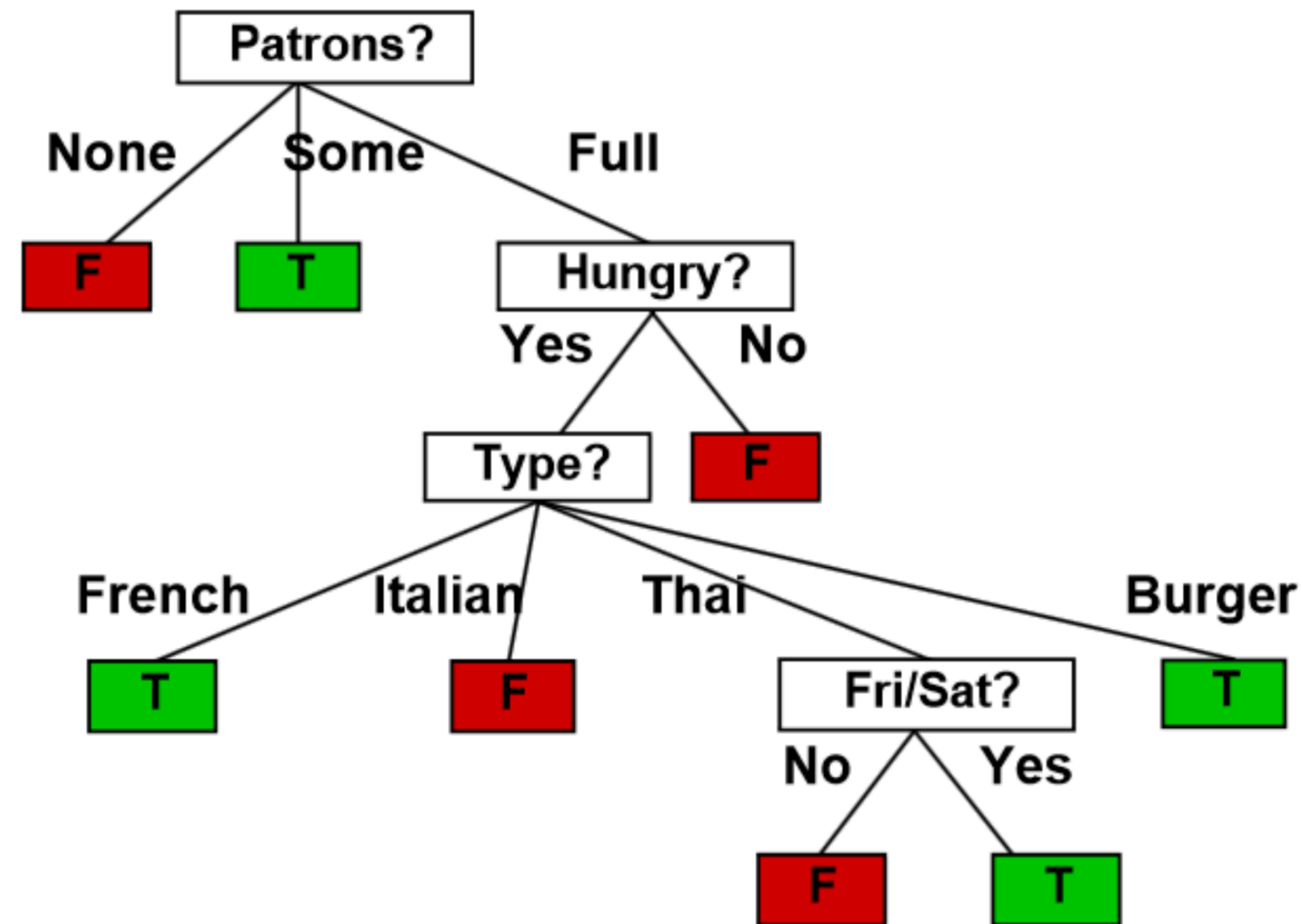


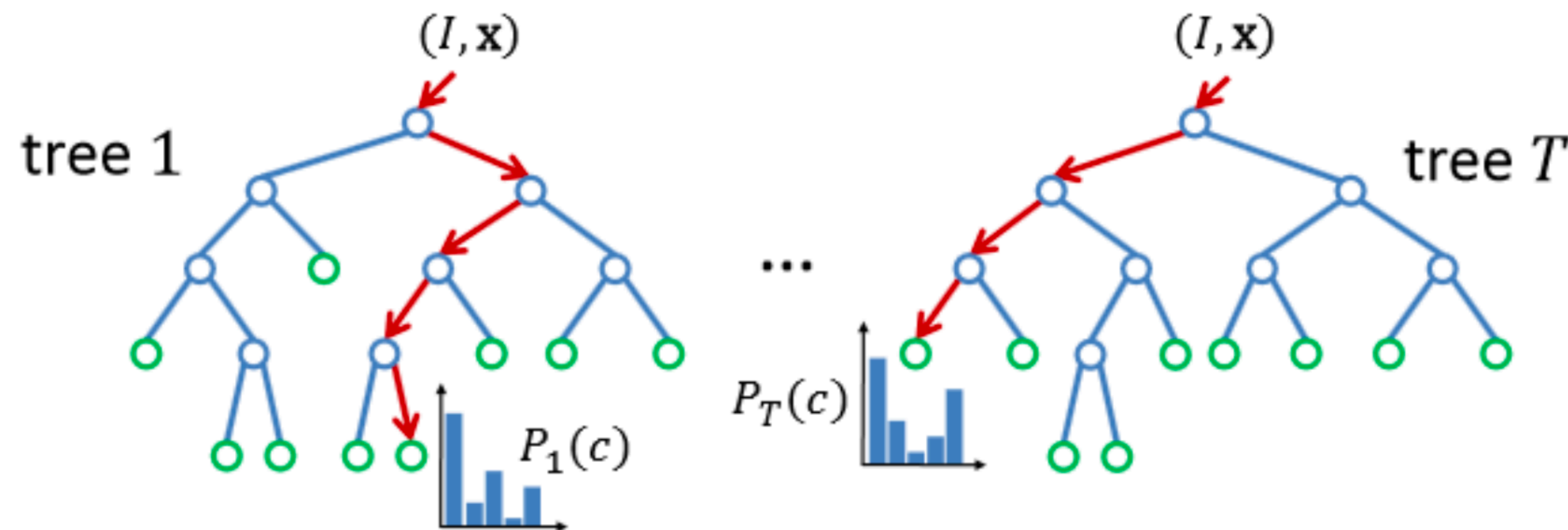
Figure credit: Russell and Norvig (3rd ed.)

Decision Tree

A **random forest** is an ensemble of decision trees.

Randomness is incorporated via training set sampling and/or generation of the candidate binary tests

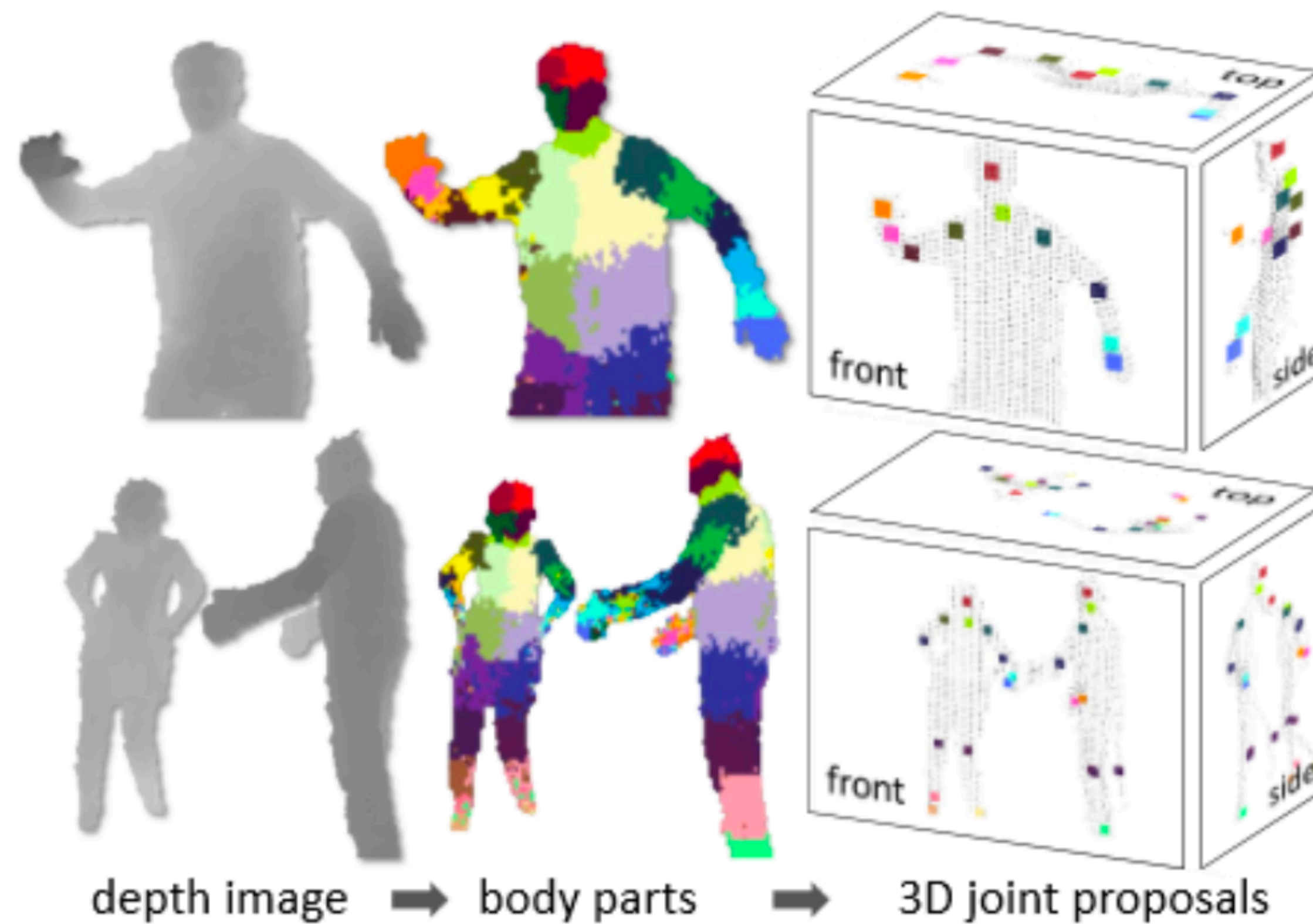
The prediction of the random forest is obtained by averaging over all decision trees.



Forsyth & Ponce (2nd ed.) Figure 14.19. Original credit: J. Shotton et al., 2011

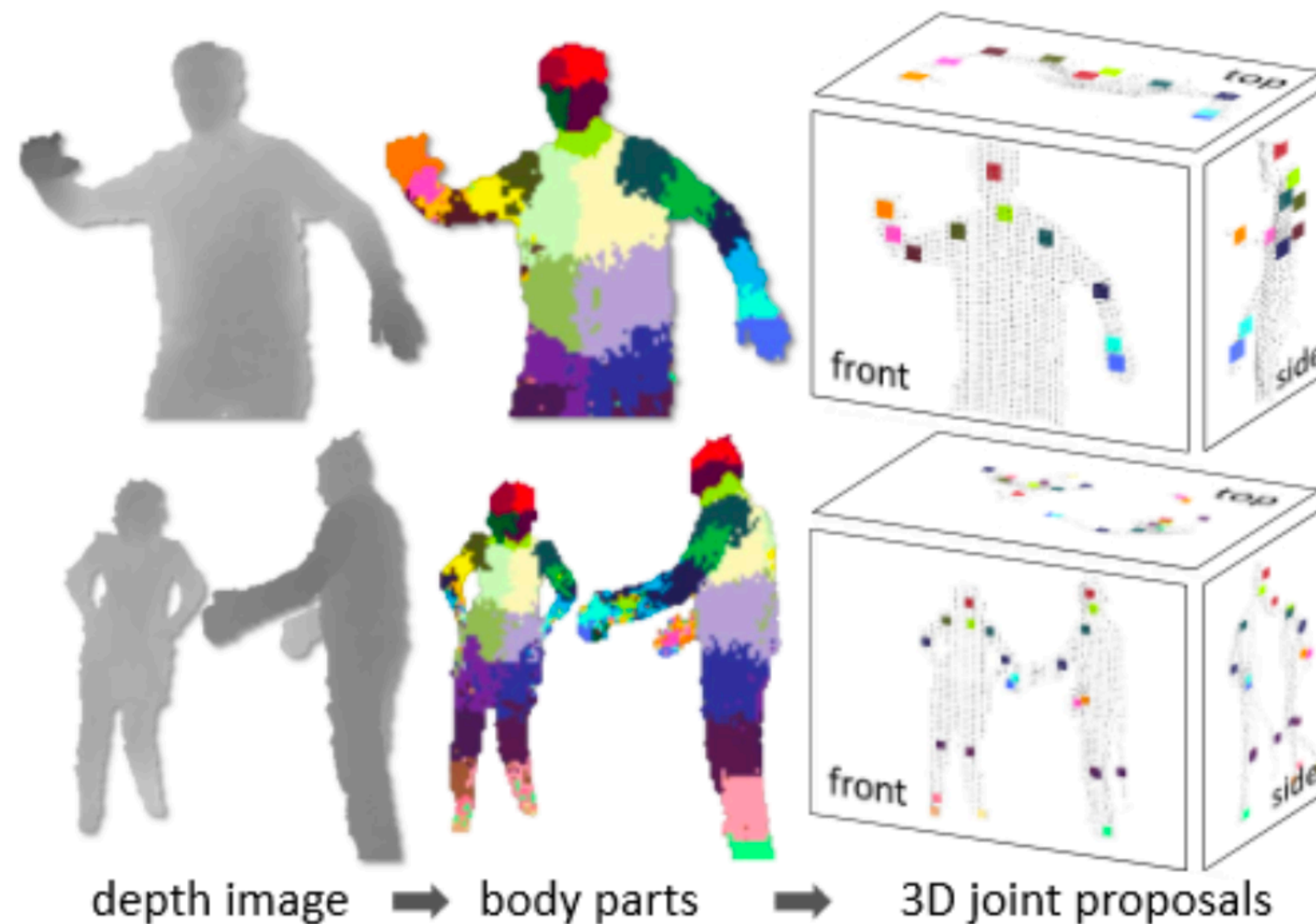
Example 1: Kinect

Kinect allows users of Microsoft's Xbox 360 console to interact with games using natural body motions instead of a traditional handheld controller. The pose (joint positions) of the user is predicted using a random forest trained on depth features.



Example 1: Kinect

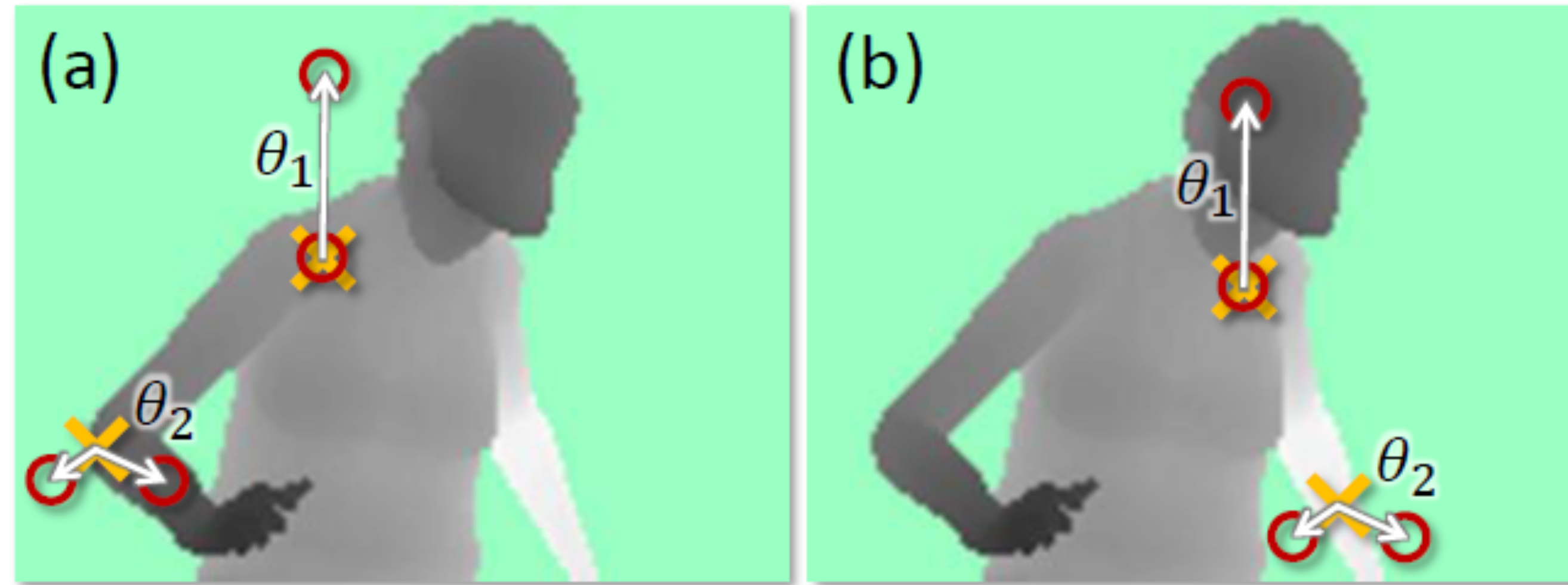
Kinect allows users of Microsoft's Xbox 360 console to interact with games using natural body motions instead of a traditional handheld controller. The pose (joint positions) of the user is predicted using a random forest trained on depth features.



Jamie Shotton

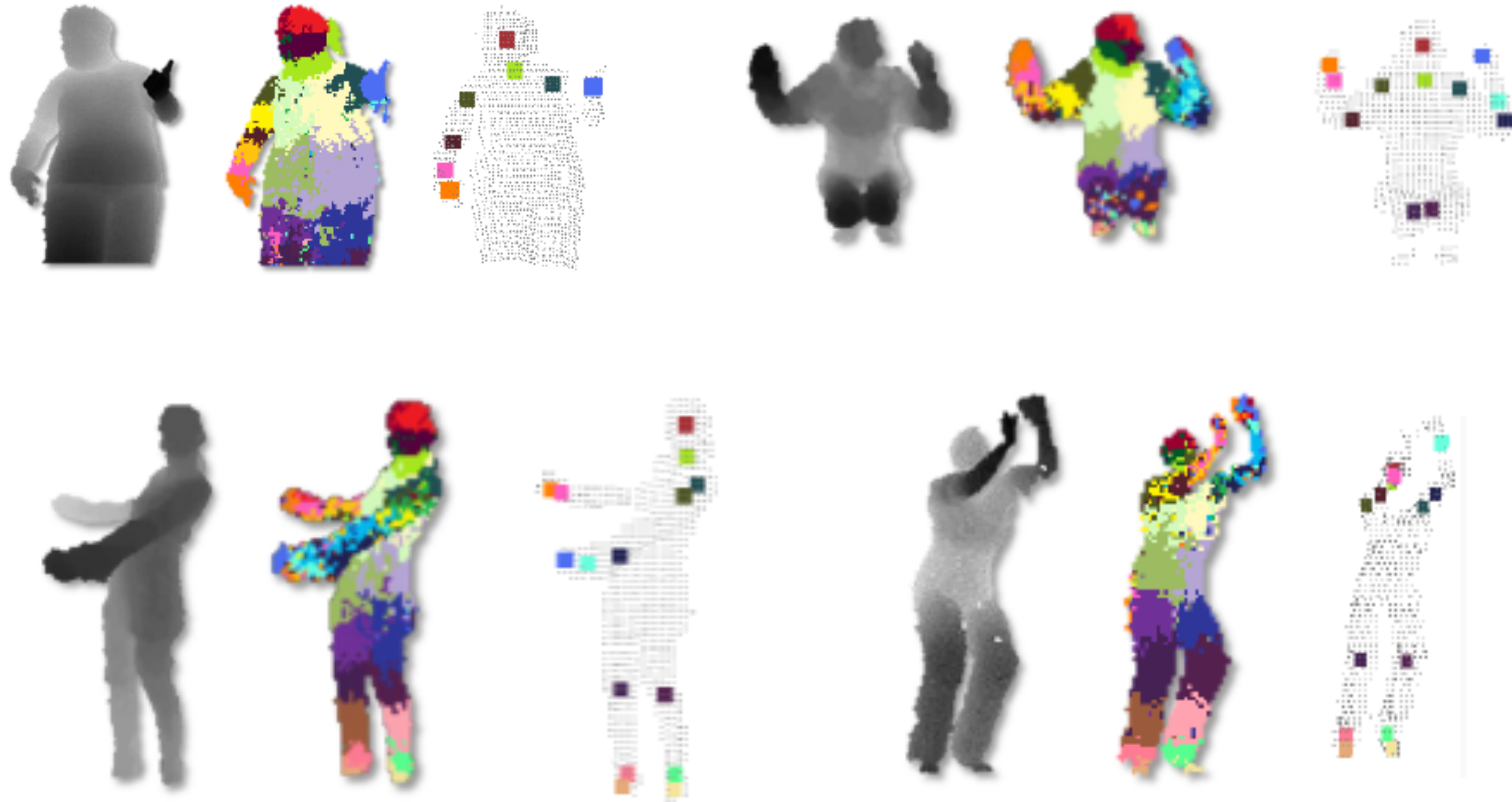
Figure credit: J. Shotton et al., 2011

Example 1: Kinect



$$f_{\theta}(I, \mathbf{x}) = d_I \left(\mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left(\mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$

Example 1: Kinect



Combining **Classifiers**

One common strategy to obtain a better classifier is to combine multiple classifiers.

A simple approach is to train an ensemble of independent classifiers, and average their predictions.

Boosting is another approach.

- Train an ensemble of classifiers sequentially.
- Bias subsequent classifiers to correctly predict training examples that previous classifiers got wrong.
- The final boosted classifier is a weighted combination of the individual classifiers.

Combining Classifiers: **Boosting**

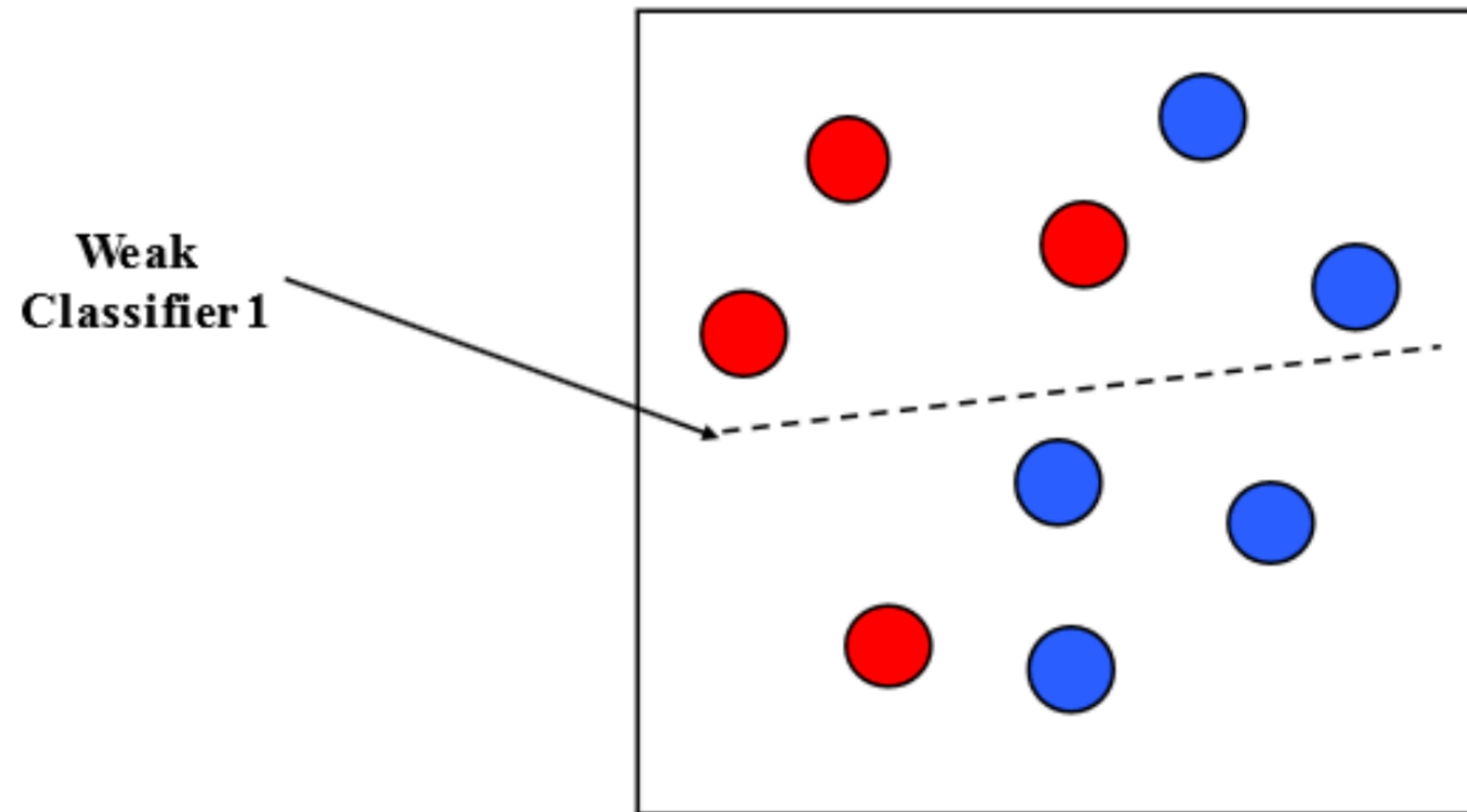


Figure credit: Paul Viola

Combining Classifiers: **Boosting**

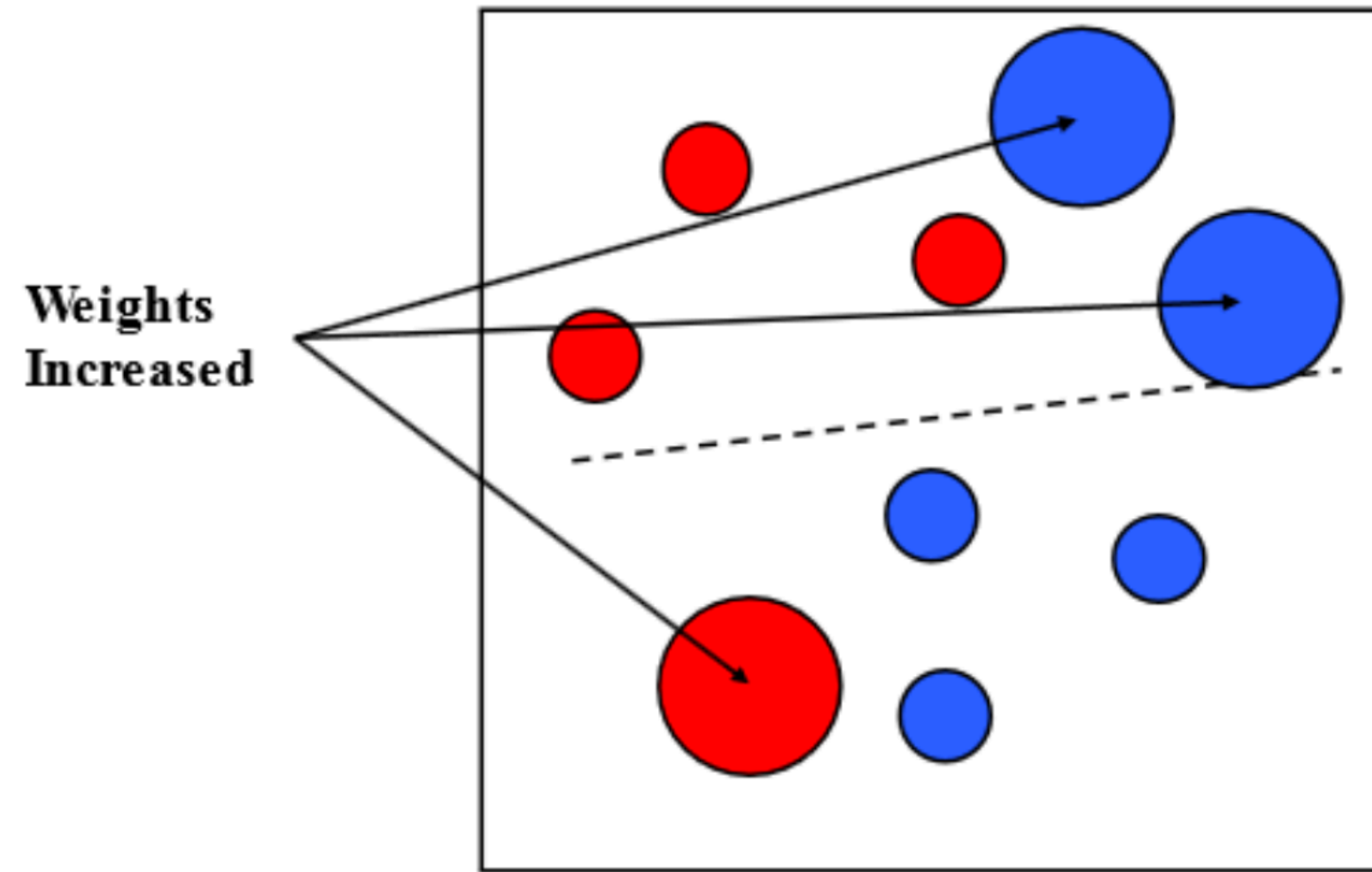


Figure credit: Paul Viola

Combining Classifiers: **Boosting**

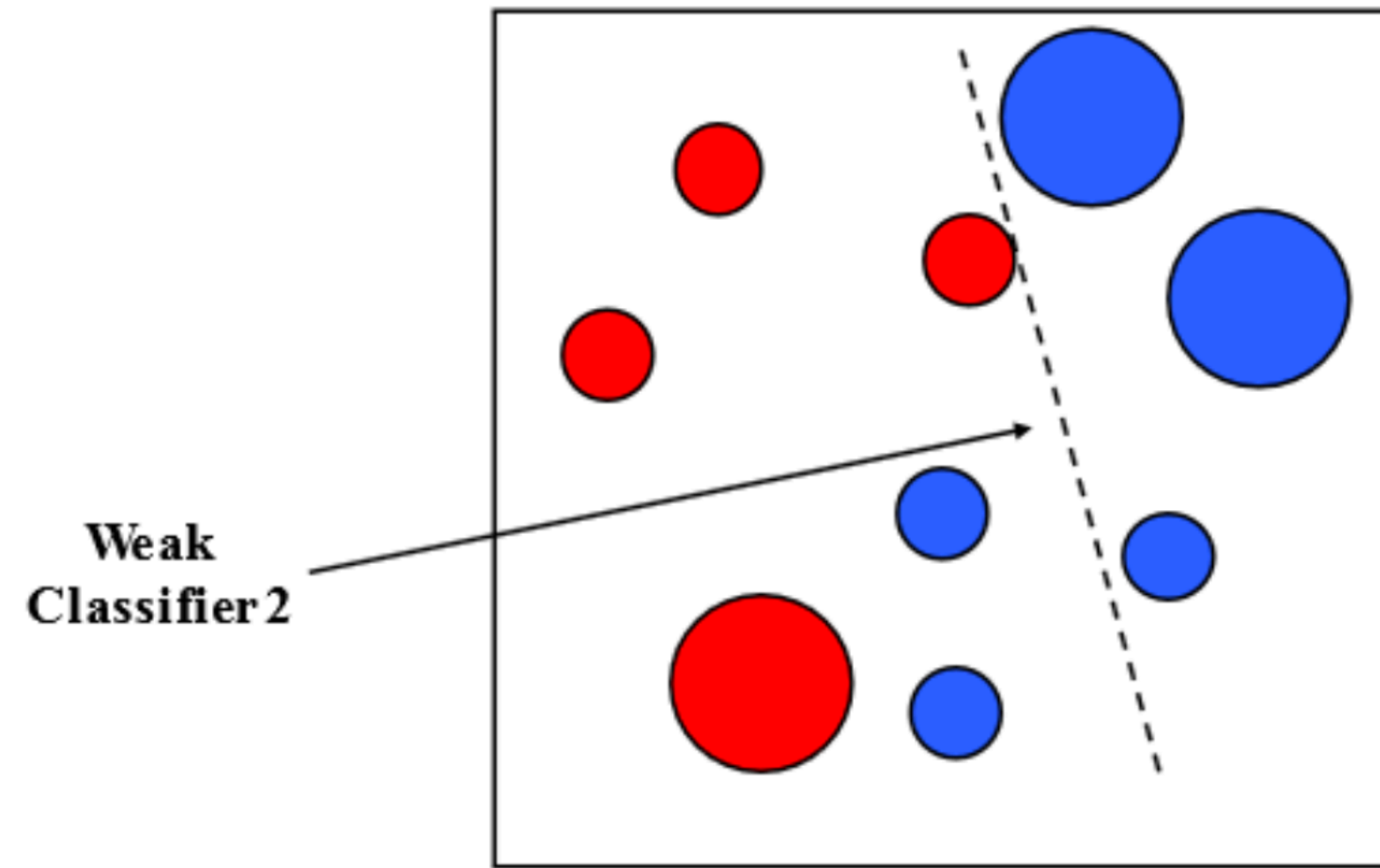


Figure credit: Paul Viola

Combining Classifiers: **Boosting**

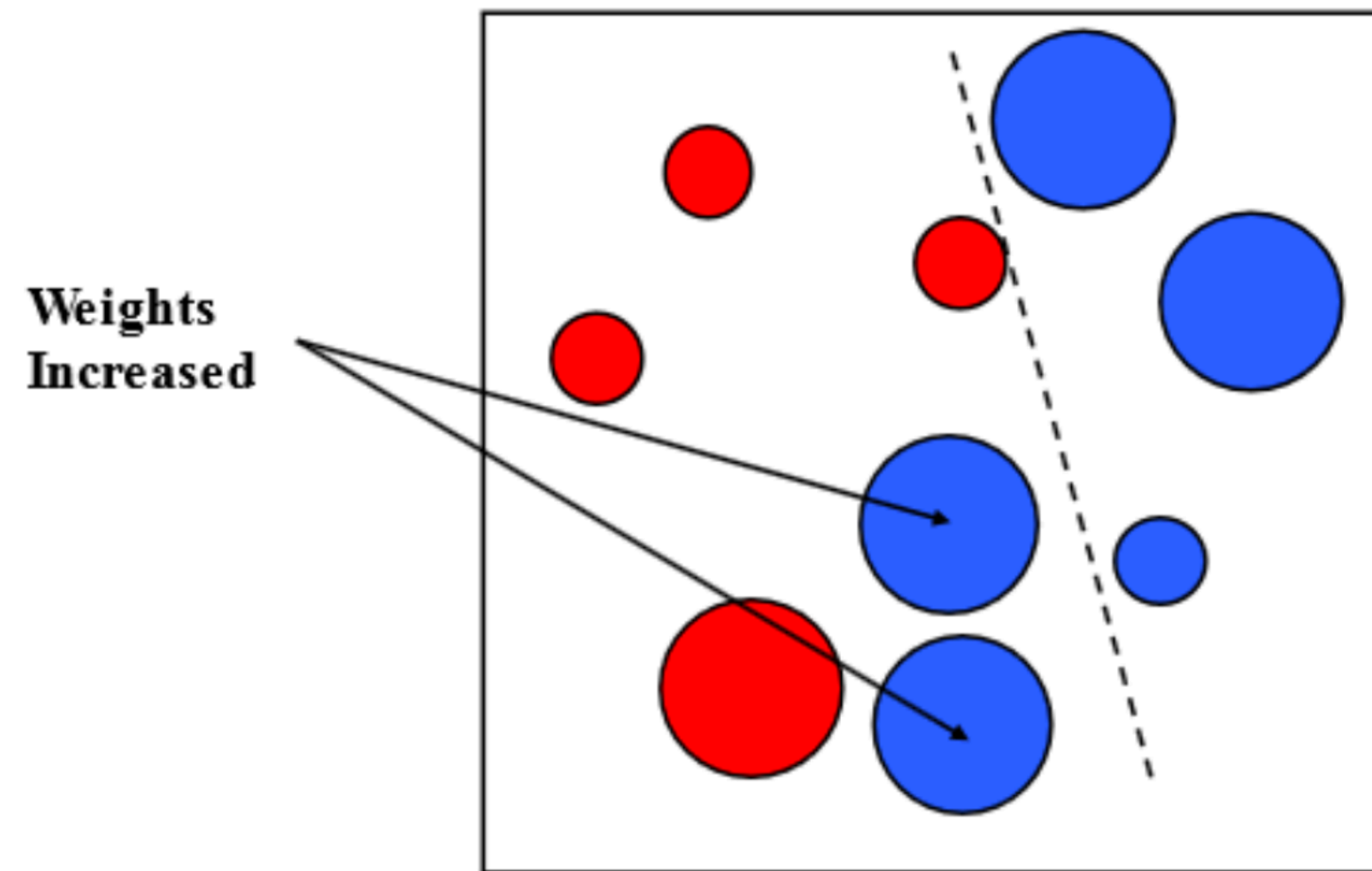


Figure credit: Paul Viola

Combining Classifiers: **Boosting**

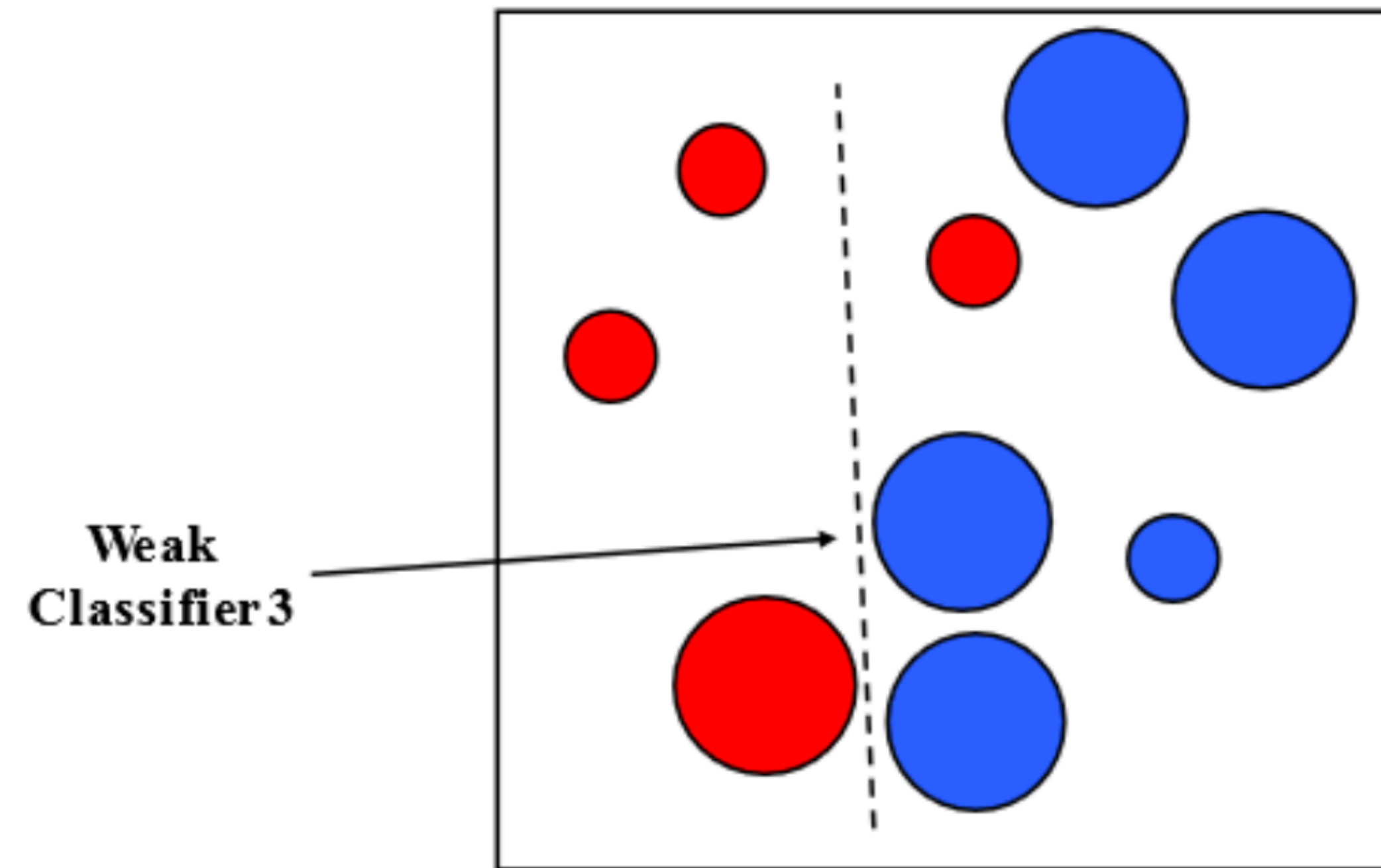


Figure credit: Paul Viola

Combining Classifiers: **Boosting**

**Final classifier is
a combination of weak
classifiers**

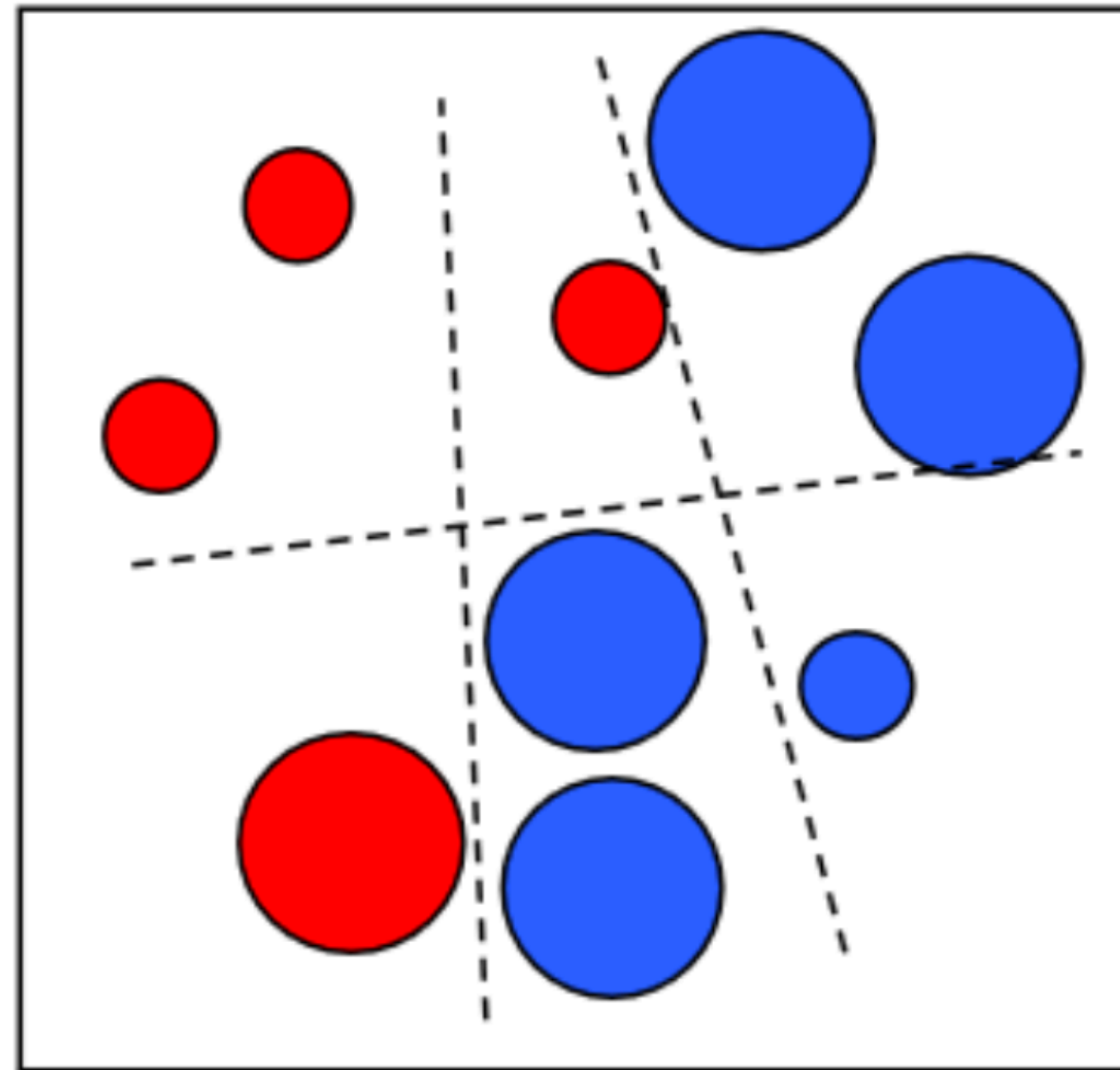


Figure credit: Paul Viola

Summary

A **decision tree** passes a data point through a sequence of feature tests. A random forest is an ensemble of decision trees.

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