

Object Classification

<Vision System>

Department of Robot Engineering
Prof. Younggun Cho

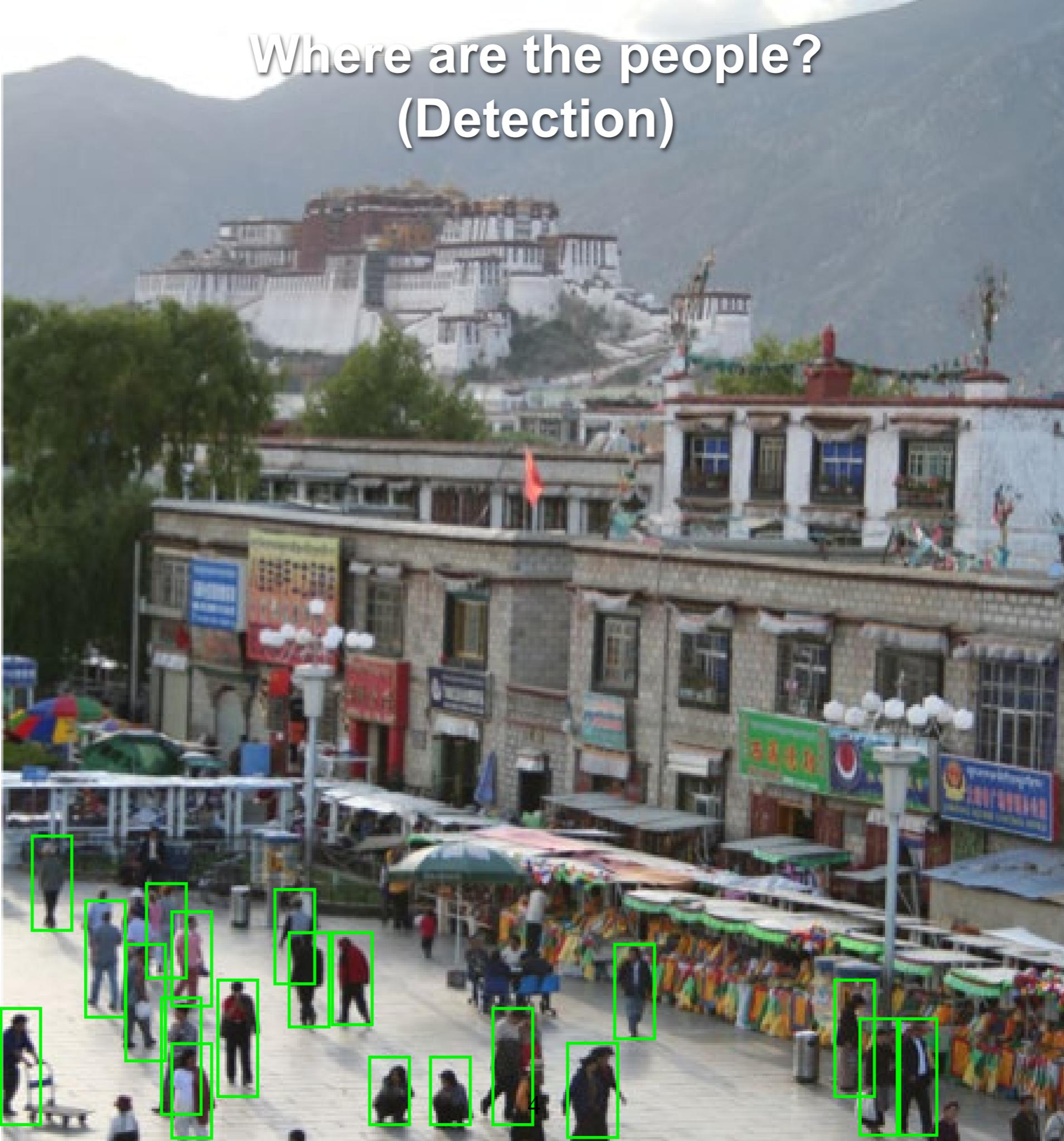


What do we mean by learning-based vision or ‘semantic vision’?

Is this a street light? (Recognition / classification)



Where are the people? (Detection)

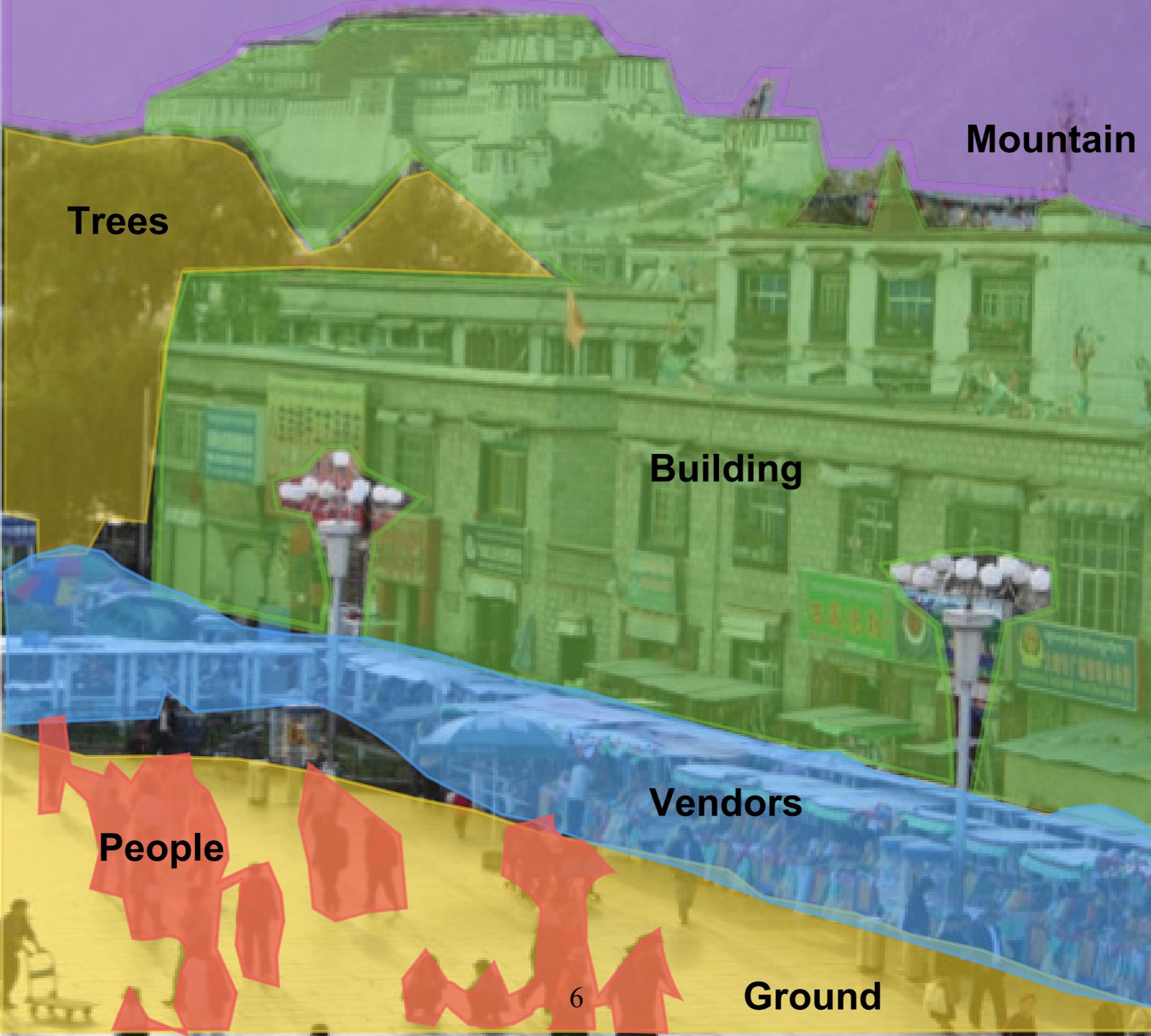


Is that Potala palace? (Identification)



Sky

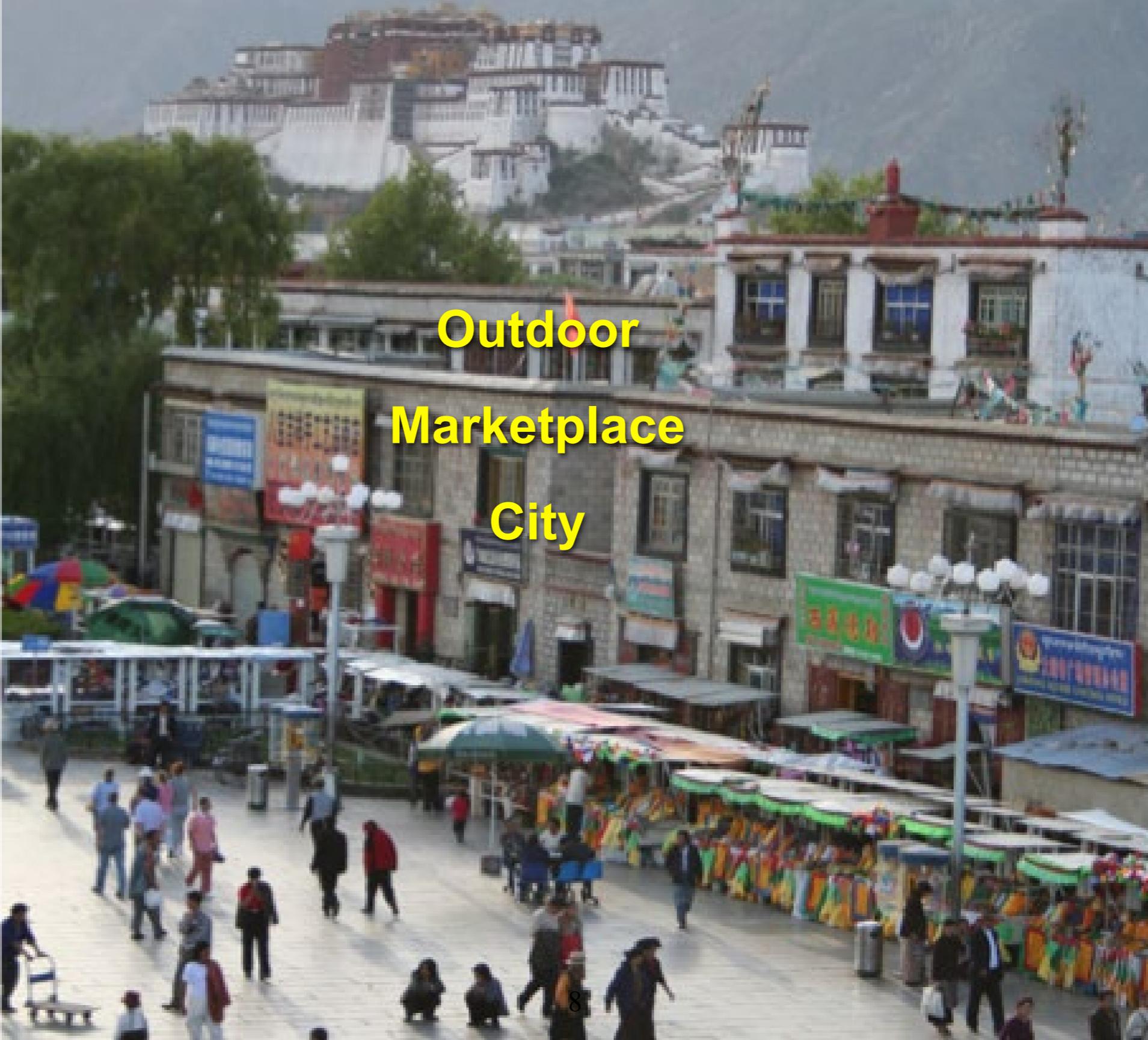
What's in the scene? (semantic segmentation)



Object categorization



What type of scene is it? (Scene categorization)



Activity / Event Recognition



Object recognition

Is it really so hard?

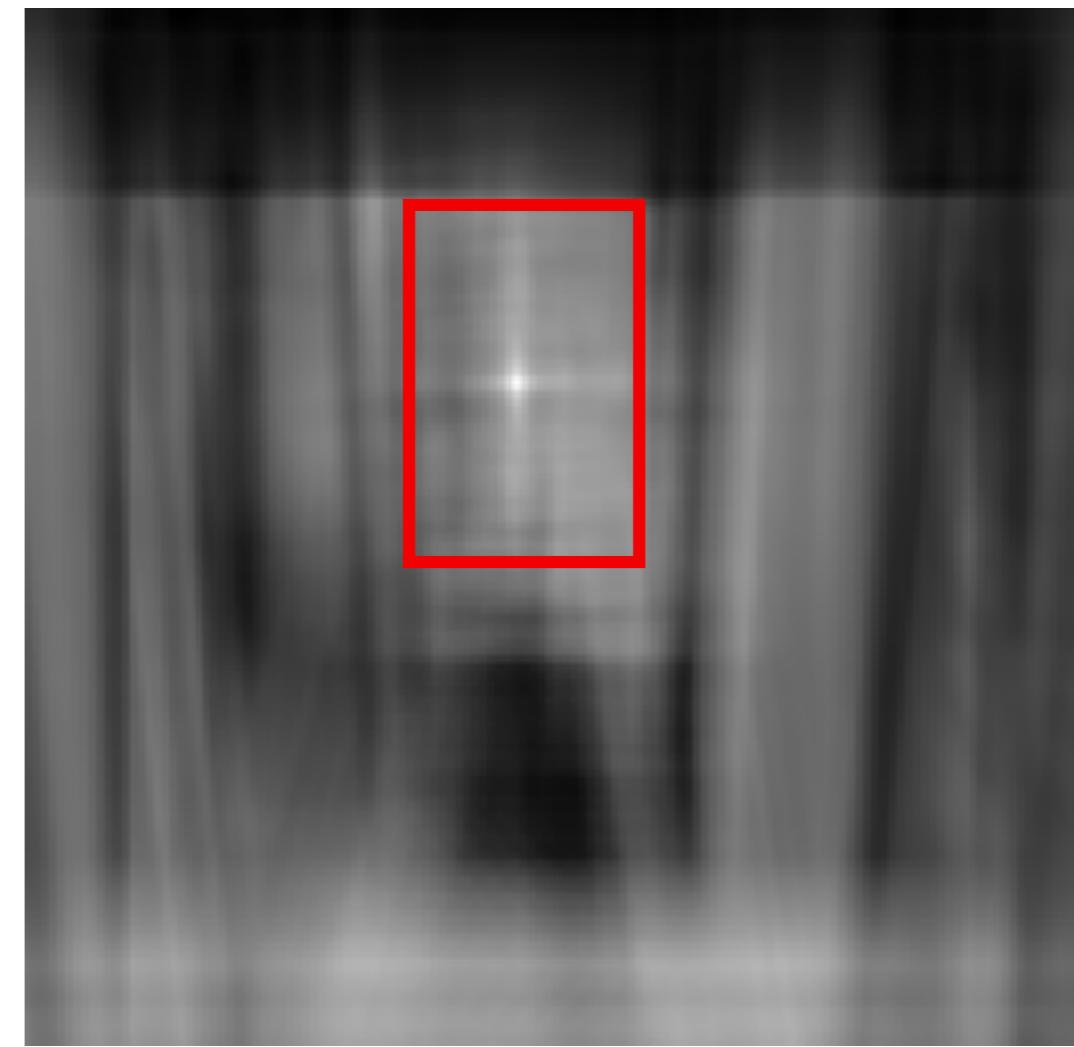
This is a chair



Find the chair in this image



Output of normalized correlation

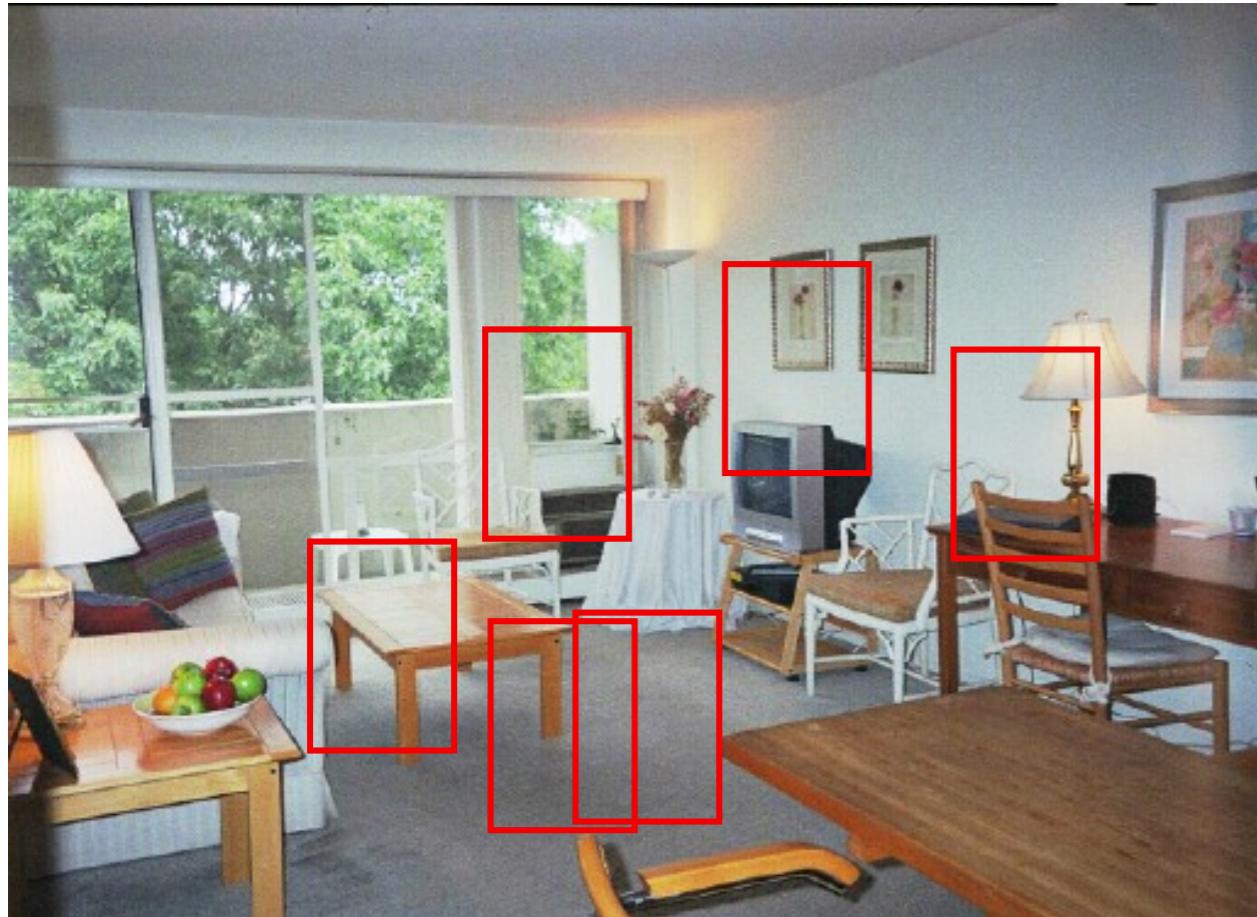




Object recognition

Is it really so hard?

Find the chair in this image



Pretty much garbage

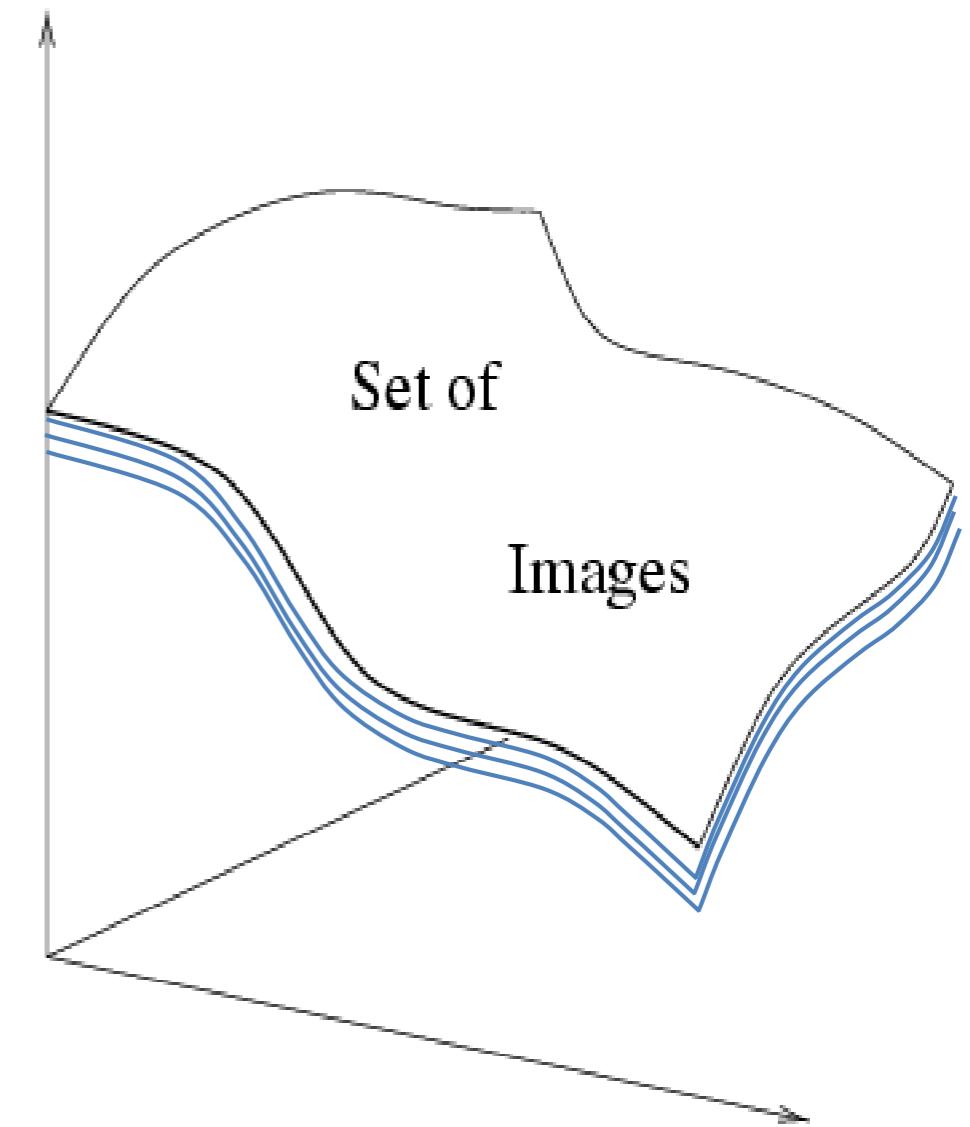
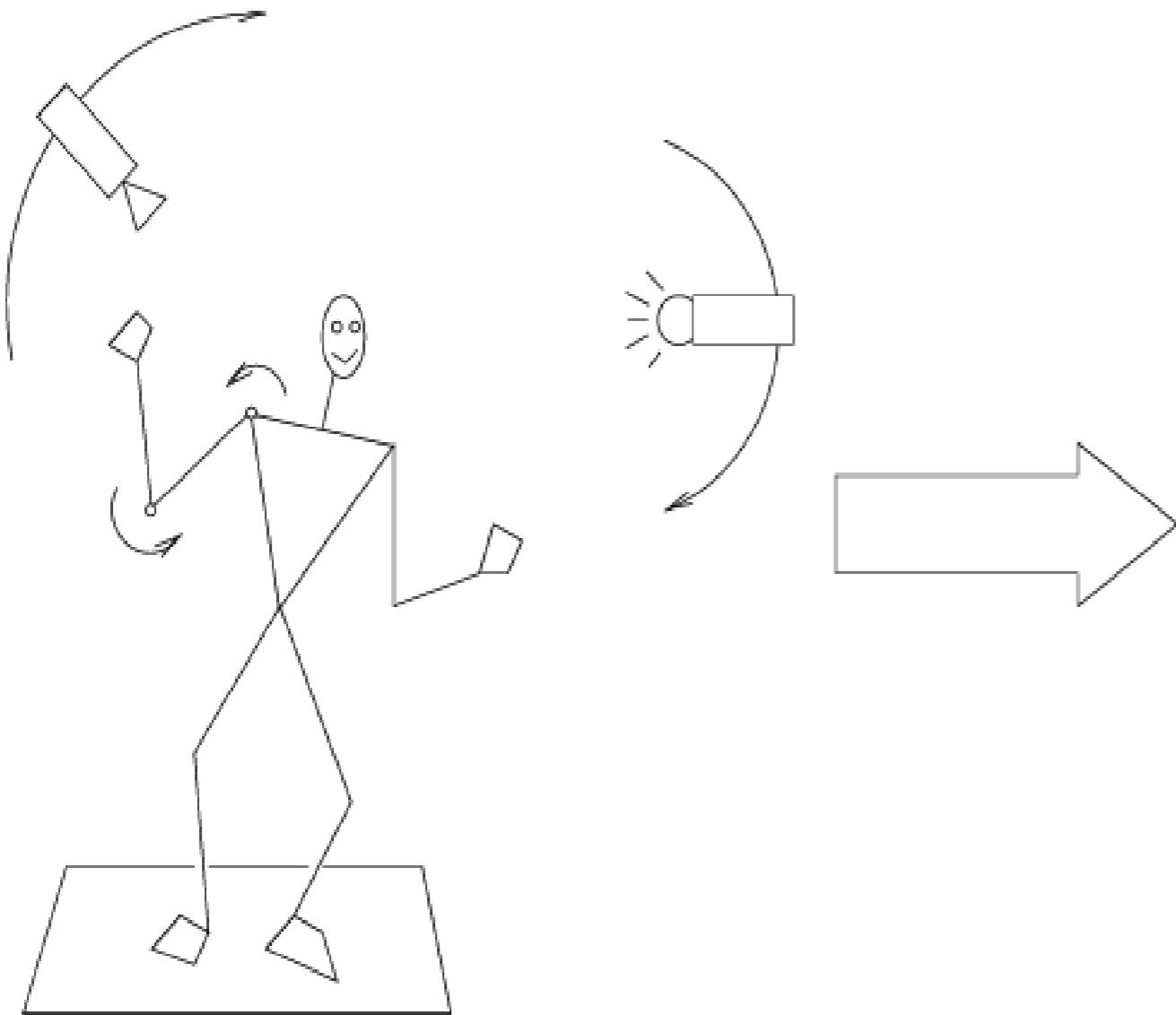
Simple template matching is not going to make it

A “popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts.” Nivatia & Binford, 1977.

And it can get a lot harder



Why is this hard?



Variability:

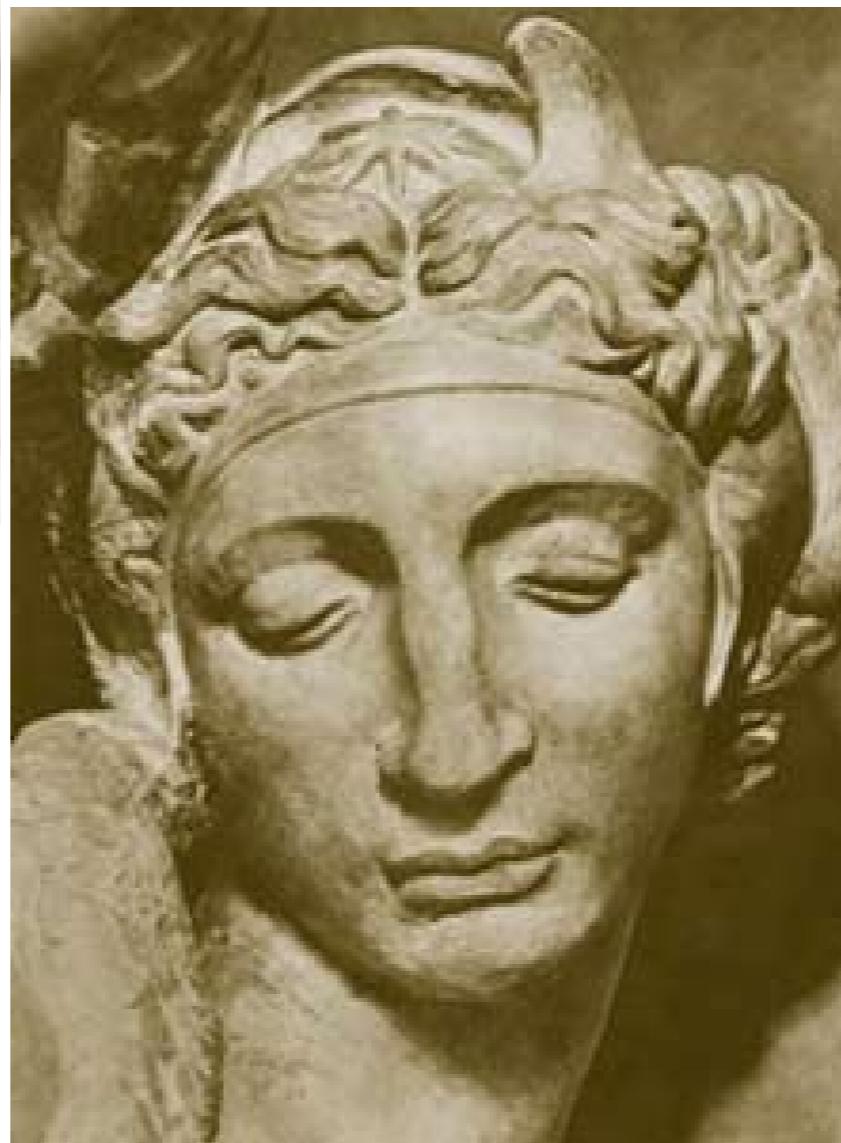
- Camera position
- Illumination
- Shape parameters

How many object categories are there?

~10,000 to 30,000

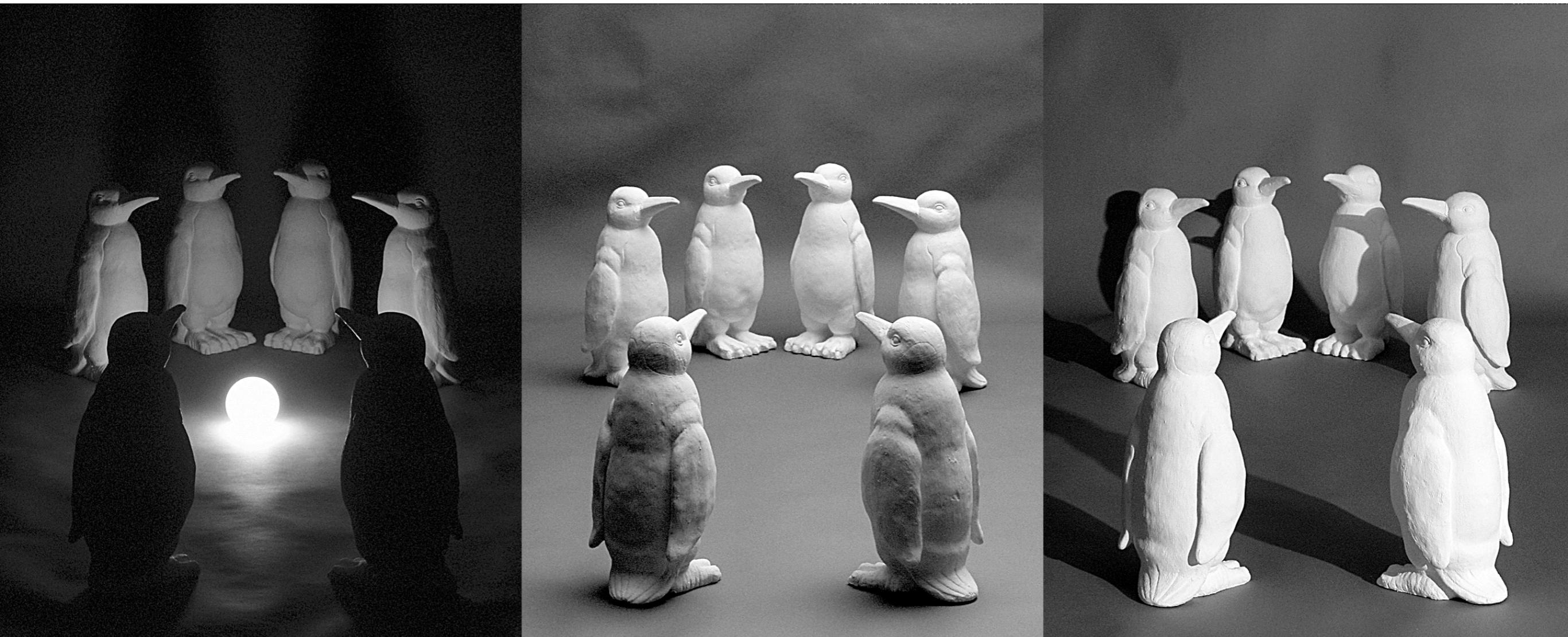


Challenge: variable viewpoint



Michelangelo 1475-1564

Challenge: variable illumination



and small things

from Apple.

(Actual size)



Challenge: scale

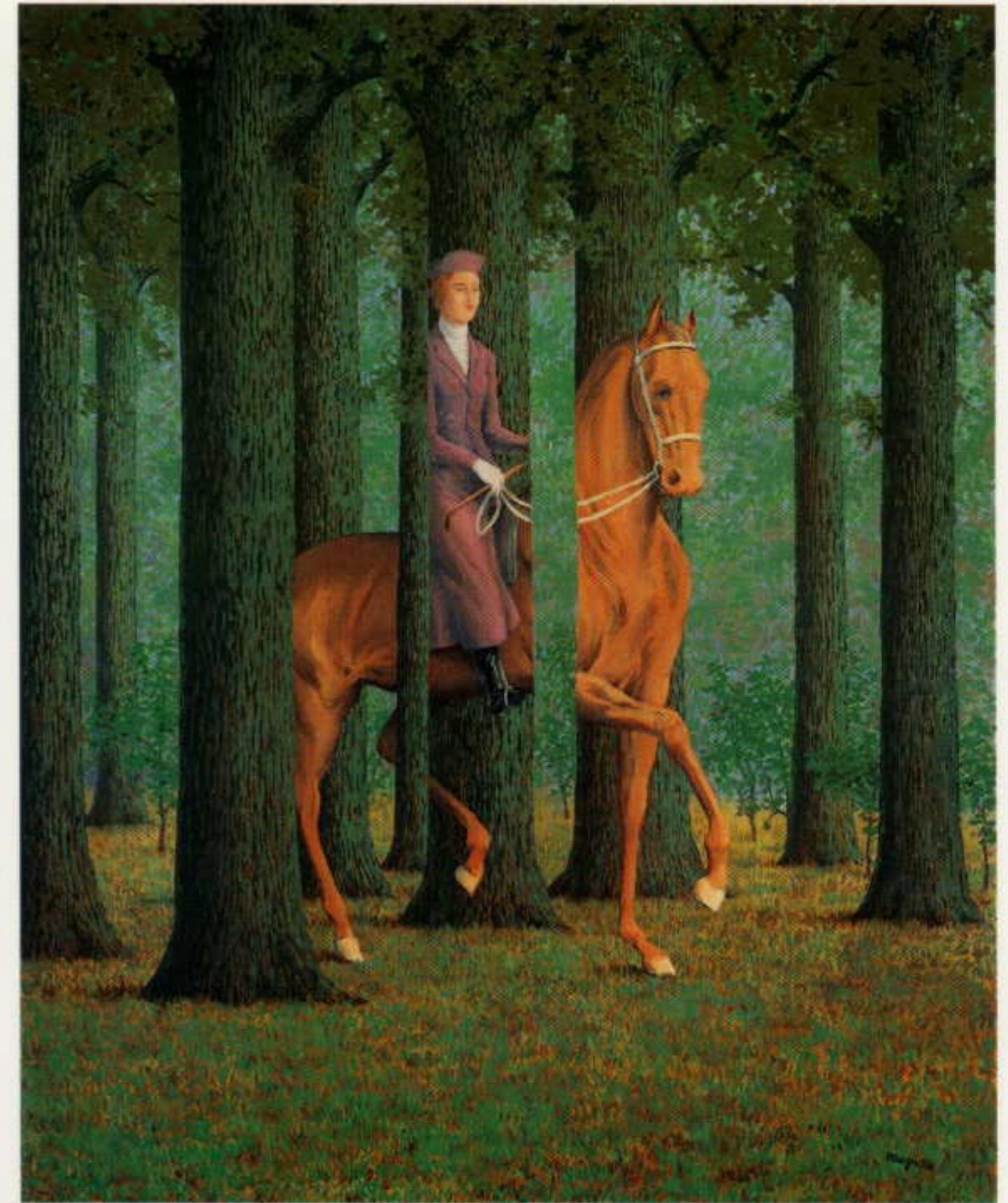
Challenge: deformation



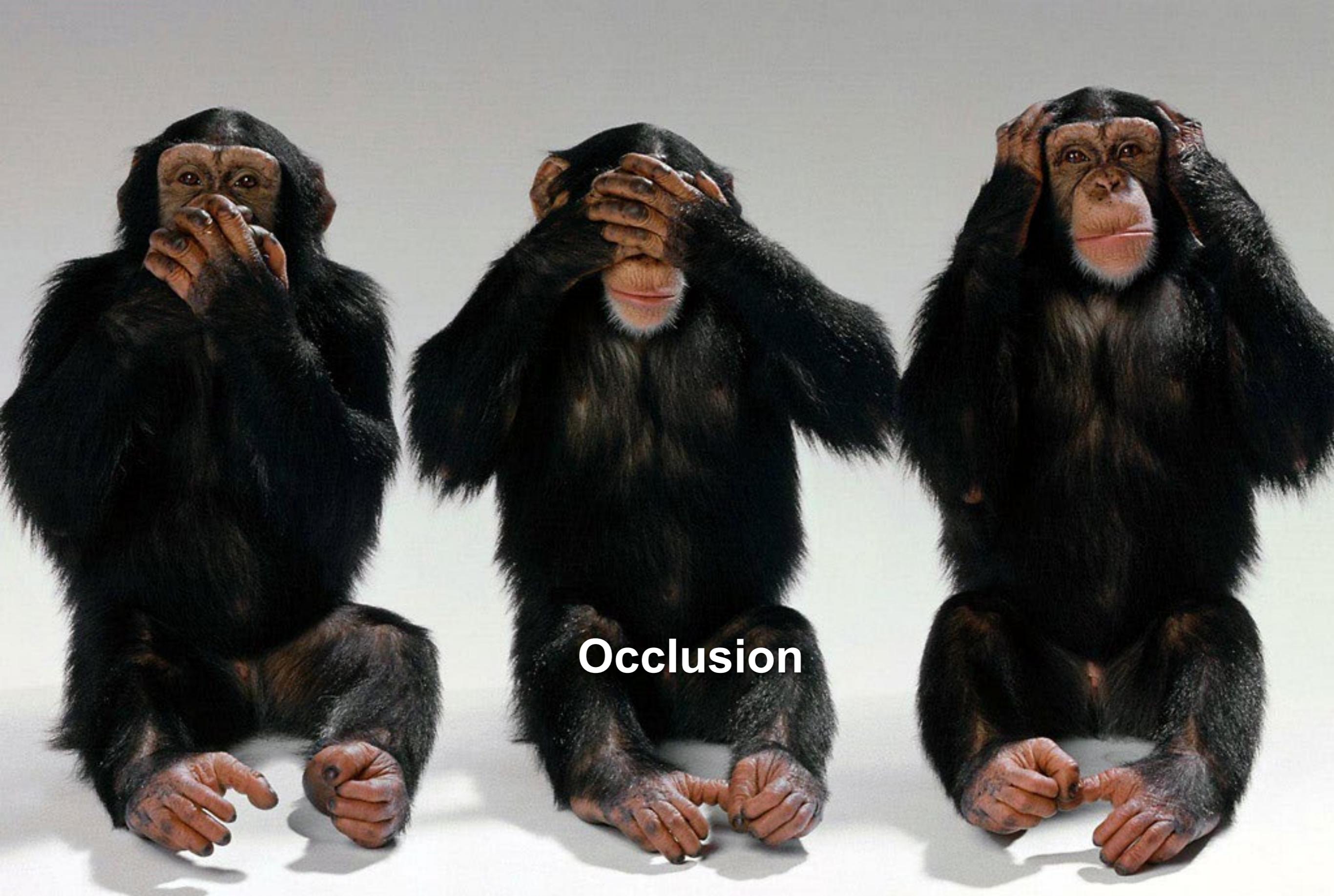


Deformation

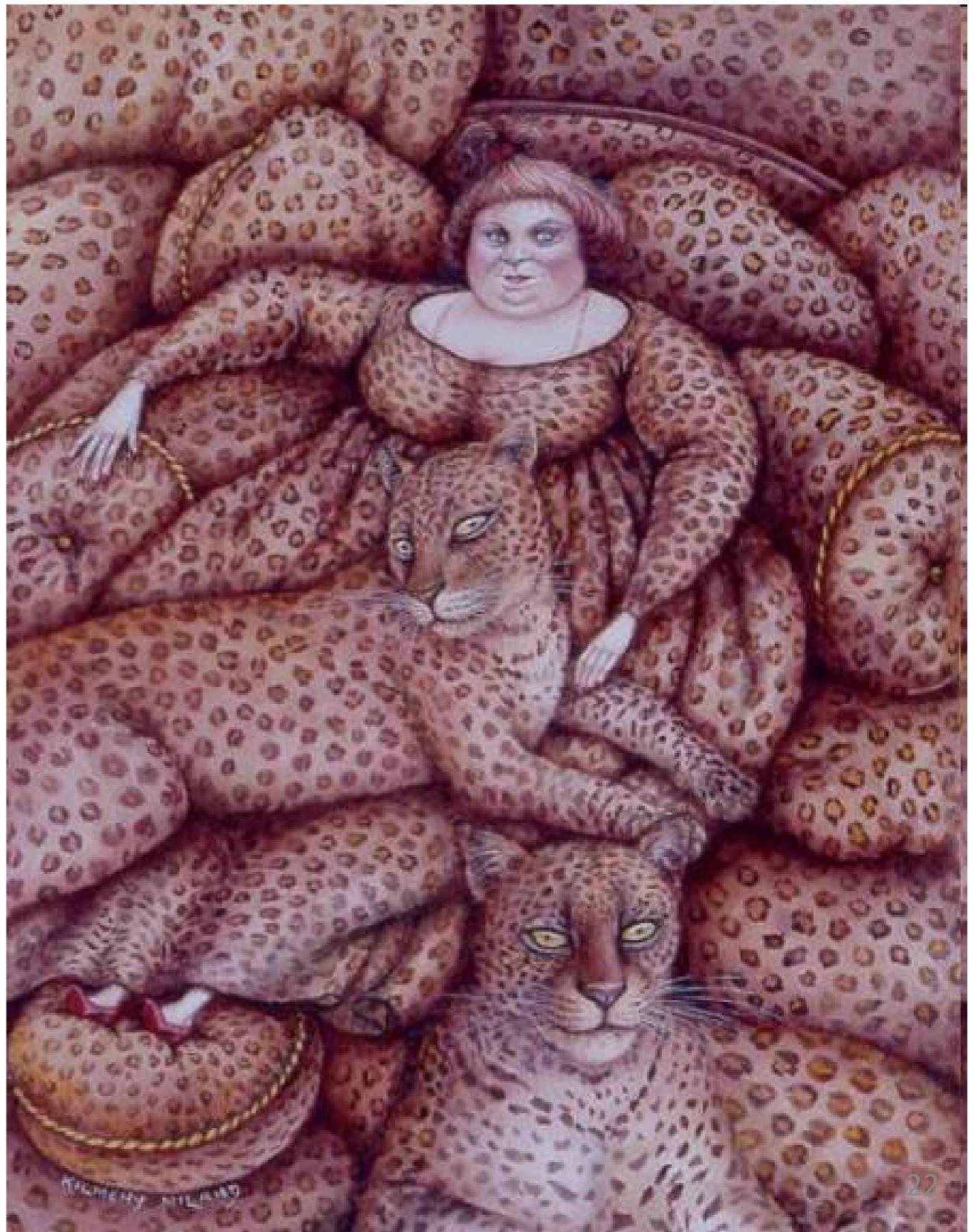
Challenge: Occlusion



Magritte, 1957



Challenge: background clutter



Kilmeny Niland. 1995



Challenge: Background clutter

Challenge: intra-class variations



Image Classification



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



cat

Image Classification: Problem



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	81	00
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
01	49	31	73	55	79	14	29	93	71	40	67	55	88	30	03	49	13	36	65
32	70	95	23	04	60	11	42	69	47	68	56	01	32	56	71	97	02	36	91
22	31	16	71	51	62	05	99	41	92	36	54	22	40	40	28	66	33	13	80
24	47	29	60	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	35	58	05	66	73	99	26	97	17	78	78	96	03	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	38	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	06	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	37	98	66
01	44	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	38	85	39	11	24	94	72	18	05	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	31	68	99	82	67	59	85	74	04	36	16	
20	73	35	29	78	31	90	01	74	31	49	71	48	66	31	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	17	42	48

What the computer sees

image classification

82% cat
15% dog
2% hat
1% mug

Data-driven approach

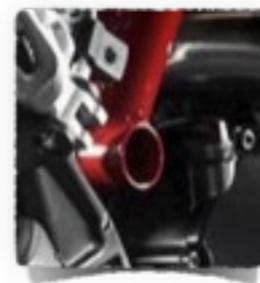
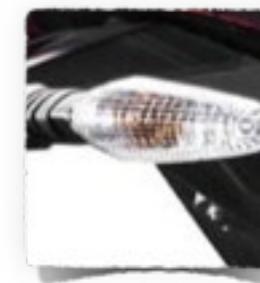
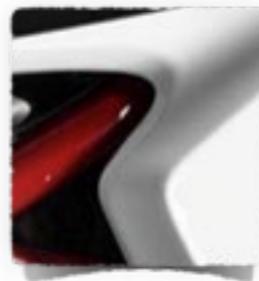
- Collect a database of images with labels
- Use ML to train an image classifier
- Evaluate the classifier on test images

Example training set

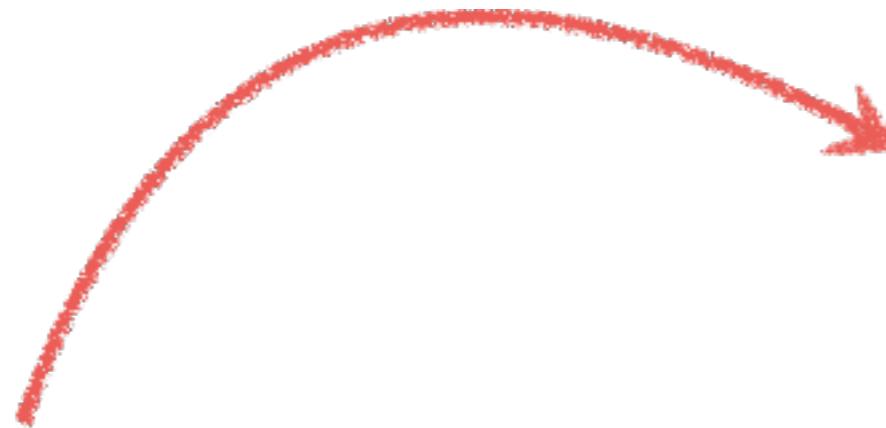


Bag of words

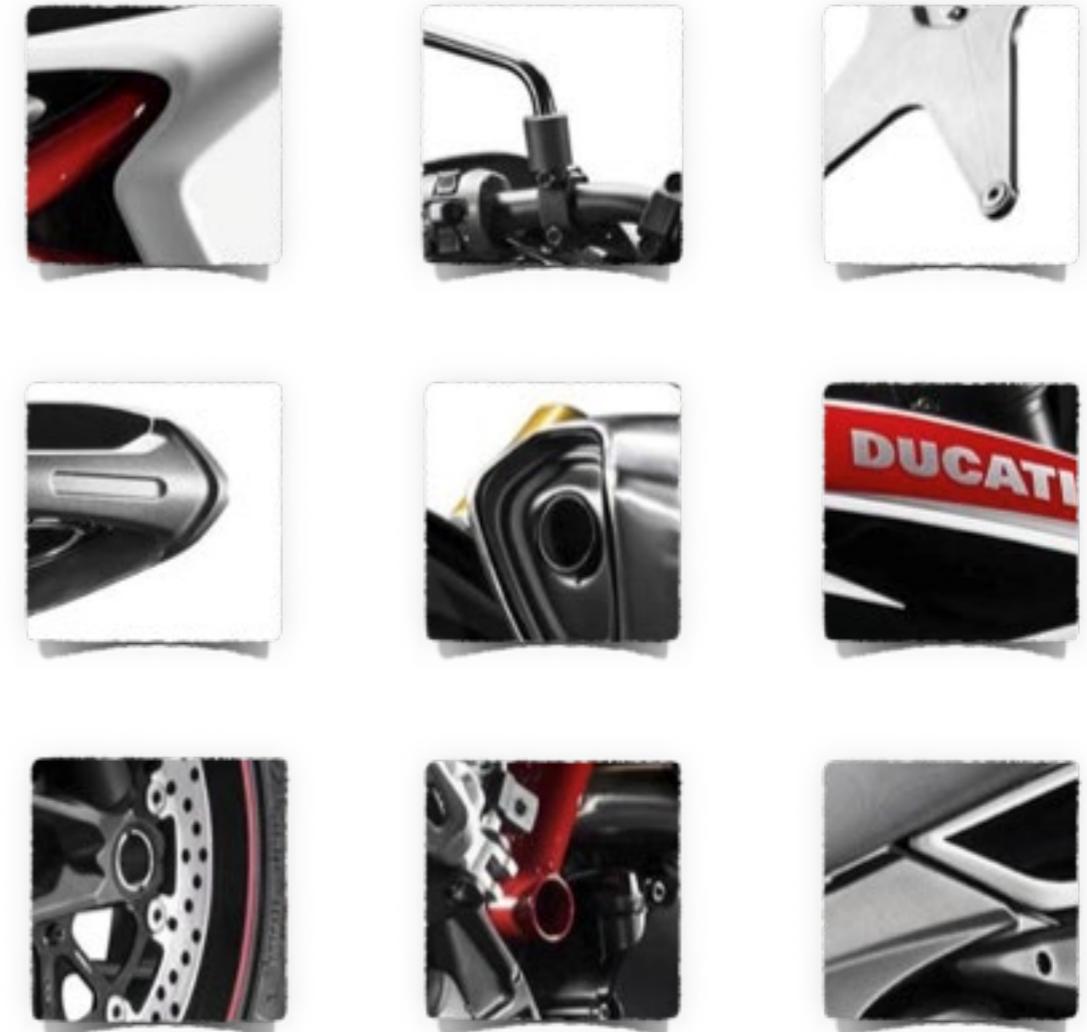
What object do these parts belong to?



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)

CalTech6 dataset



class	bag of features		Parts-and-shape model
	Zhang et al. (2005)	Willamowski et al. (2004)	Fergus et al. (2003)
airplanes	98.8	97.1	90.2
cars (rear)	98.3	98.6	90.3
cars (side)	95.0	87.3	88.5
faces	100	99.3	96.4
motorbikes	98.5	98.0	92.5
spotted cats	97.0	—	90.0

Works pretty well for image-level classification

Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)

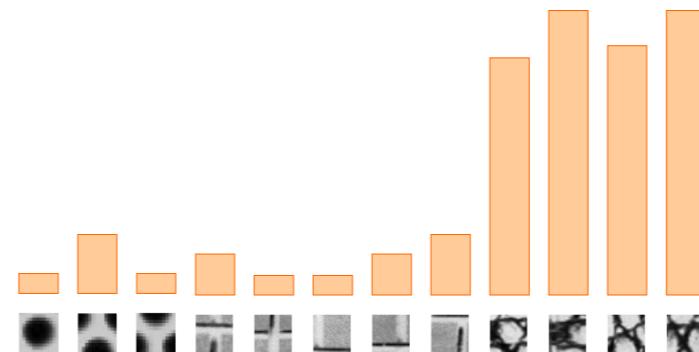
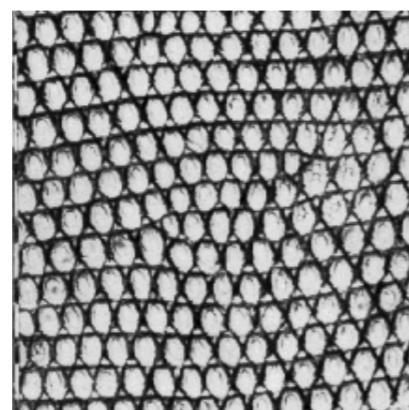
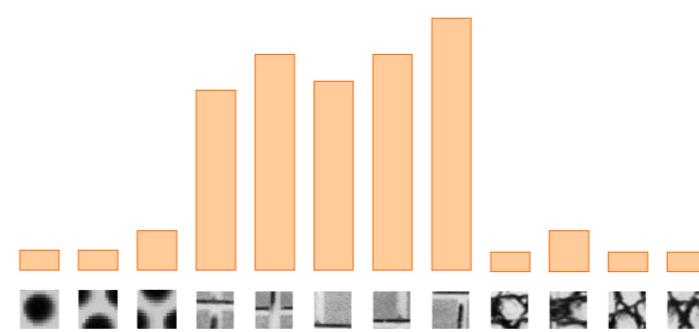
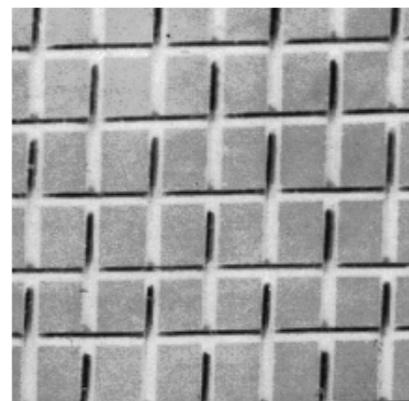
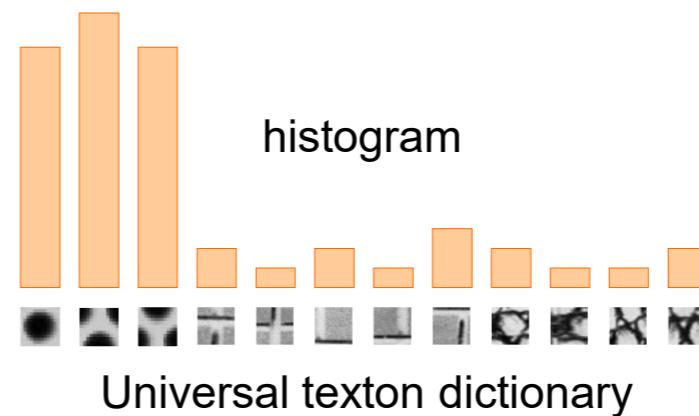
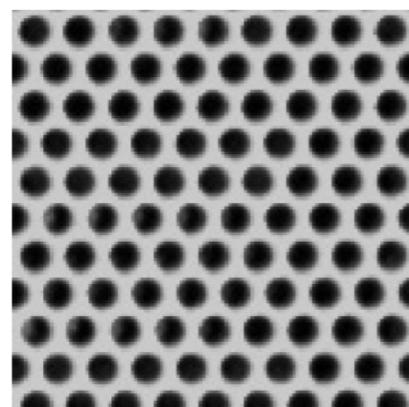
Bag-of-features

represent a data item (document, texture, image)
as a histogram over features

an old idea

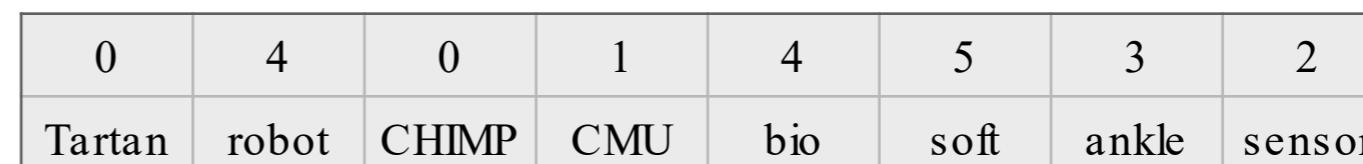
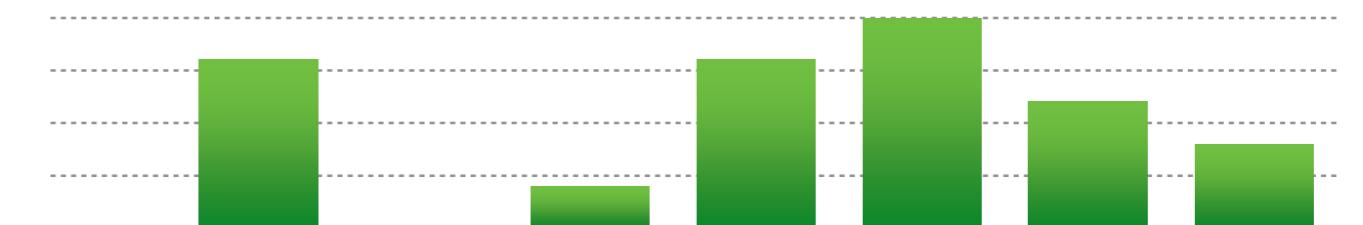
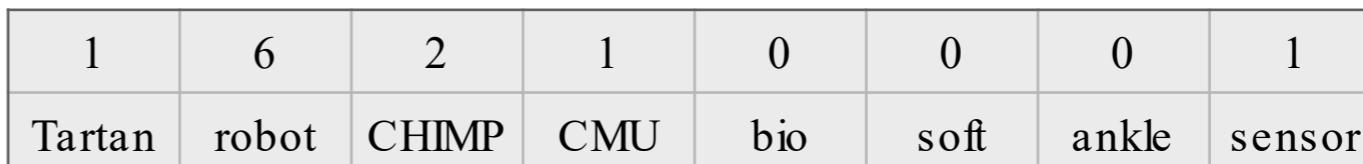
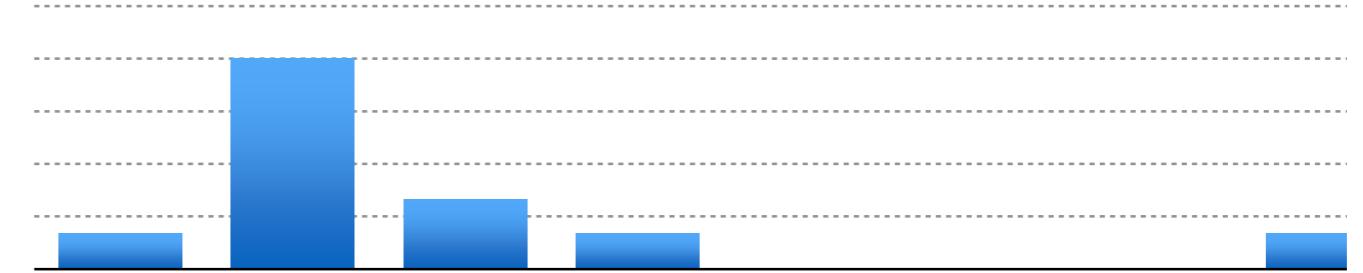
(e.g., texture recognition and information retrieval)

Texture recognition



Vector Space Model

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



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

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

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

just a histogram over words

What is the similarity between two documents?



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

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

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

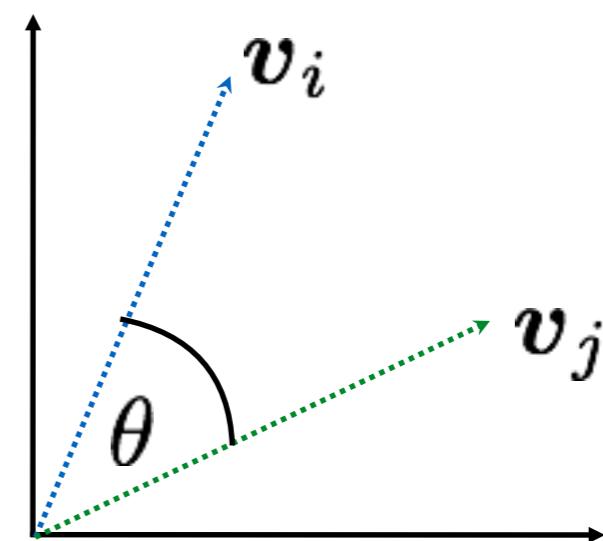
just a histogram over words

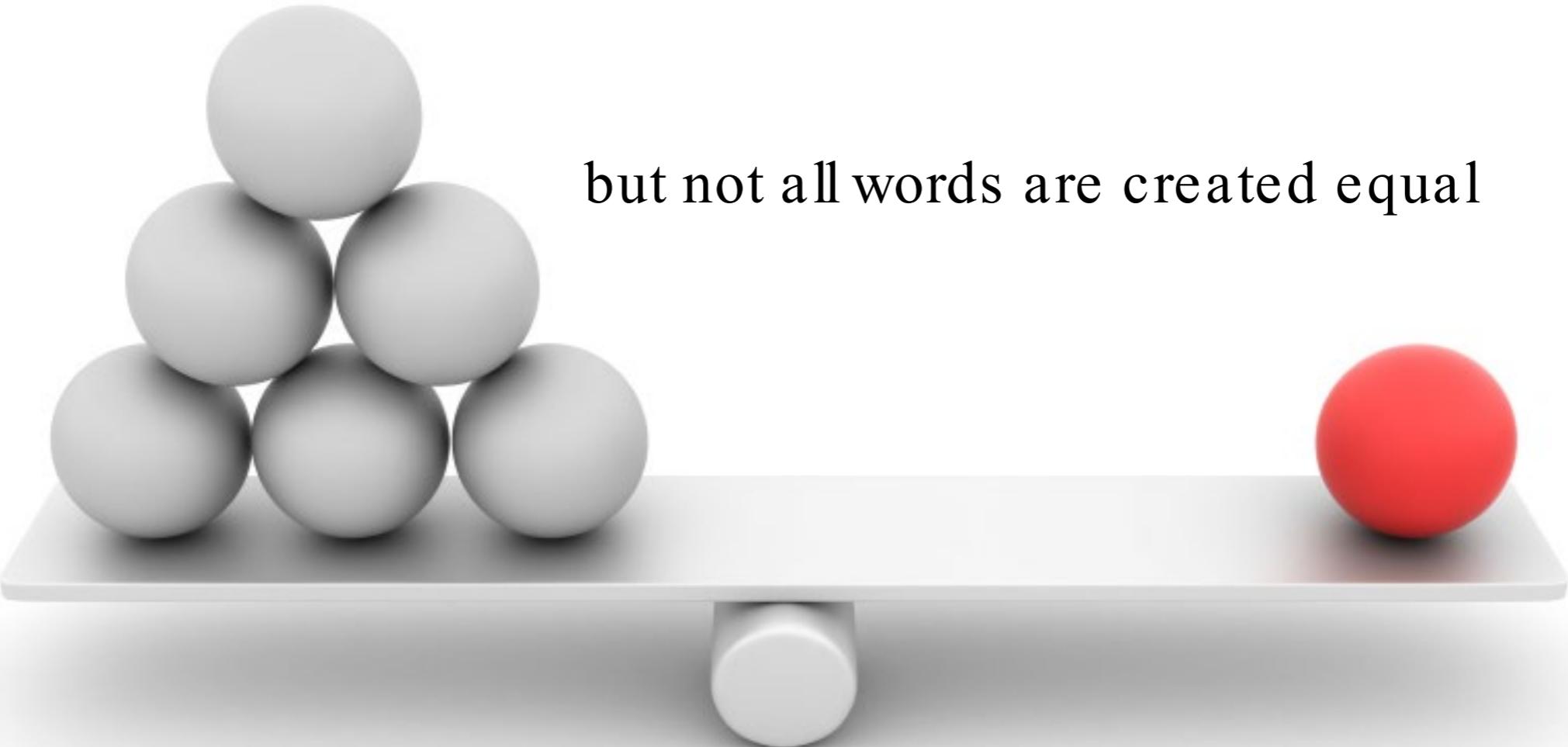
What is the similarity between two documents?



Use any distance you want but the cosine distance is fast.

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





but not all words are created equal

TF-IDF

Term Frequency Inverse Document Frequency

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

weigh each word by a heuristic

$$\mathbf{v}_d = [n(w_{1,d})\alpha_1 \ n(w_{2,d})\alpha_2 \ \cdots \ n(w_{T,d})\alpha_T]$$

$$n(w_{i,d})\alpha_i = n(w_{i,d}) \log \left\{ \frac{D}{\sum_{d'} \mathbf{1}[w_i \in d']} \right\}$$

term frequency inverse document frequency

(down-weights **common** terms)

Standard BOW pipeline

(for image classification)

Dictionary Learning:

Learn Visual Words using clustering

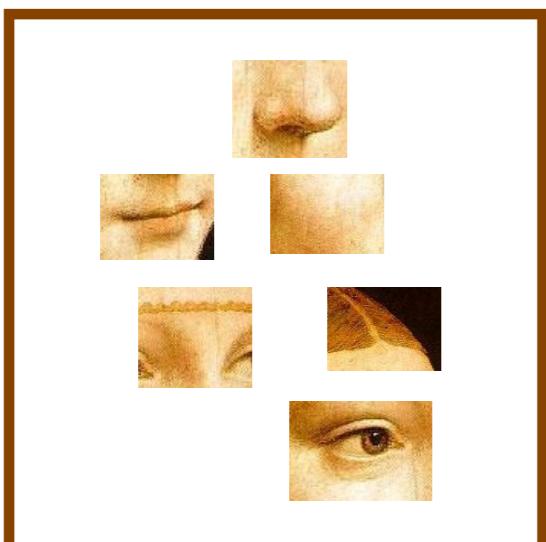
Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs

Dictionary Learning:

Learn Visual Words using clustering

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



Dictionary Learning:

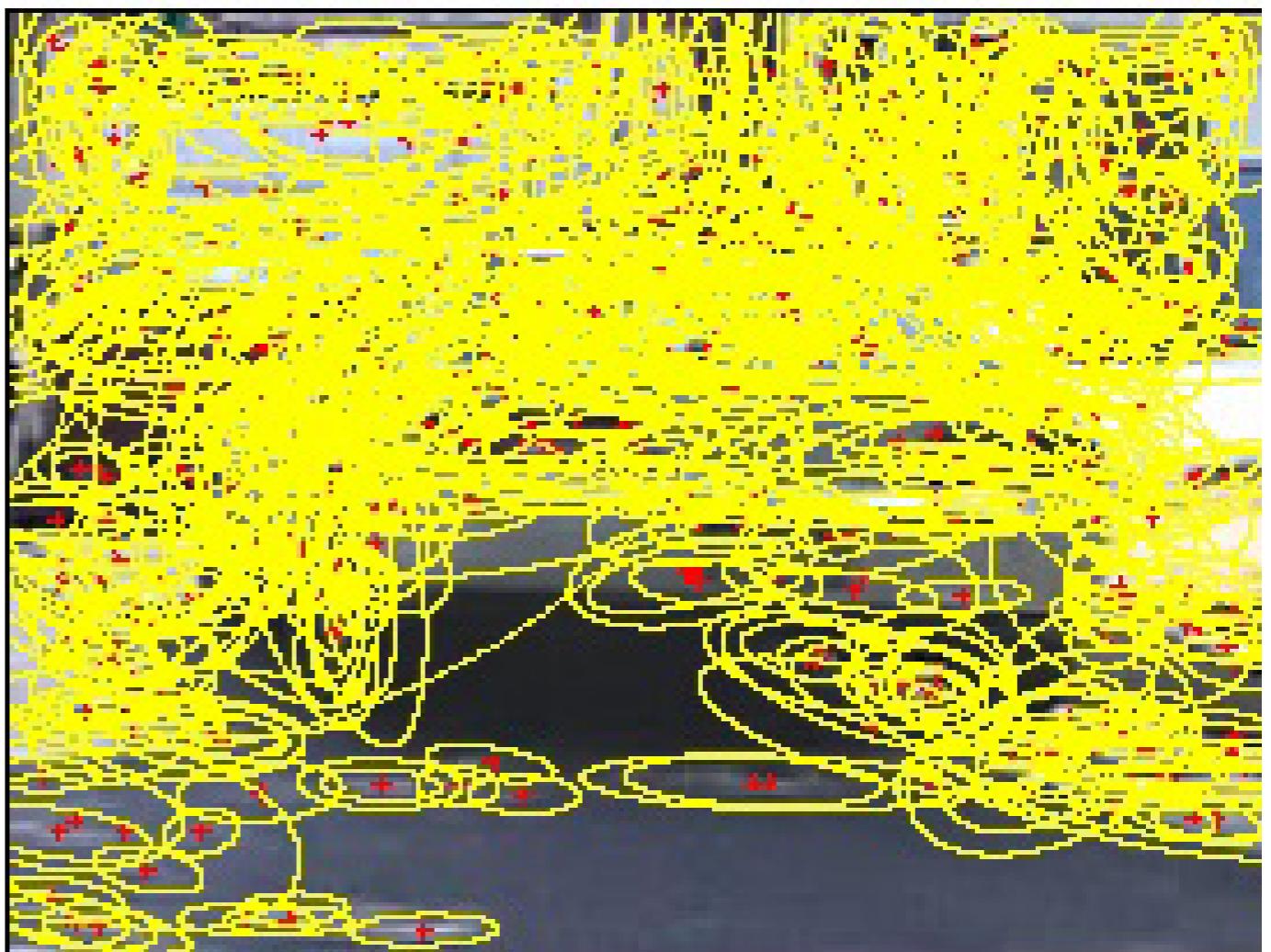
Learn Visual Words using clustering

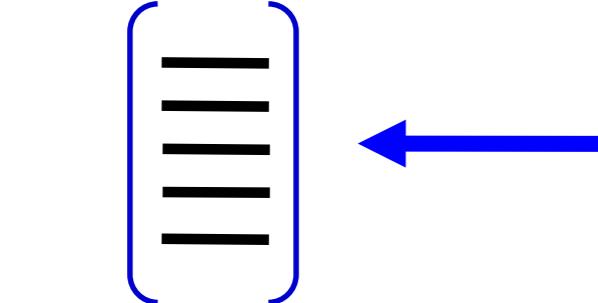
2. Learn visual dictionary (e.g., K-means clustering)



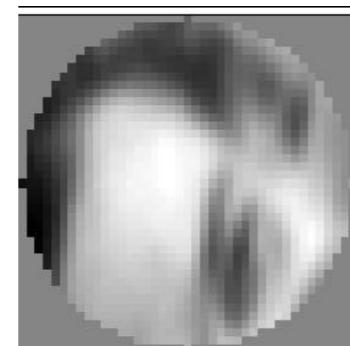
What kinds of features can we extract?

- 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)





**Compute SIFT
descriptor**
[Lowe'99]



Normalize patch

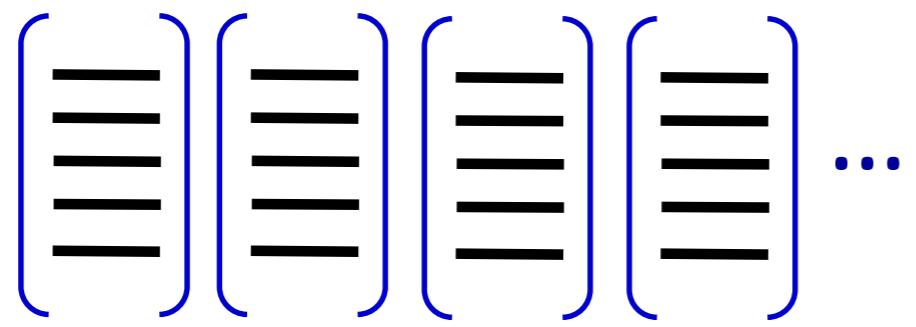


Detect patches

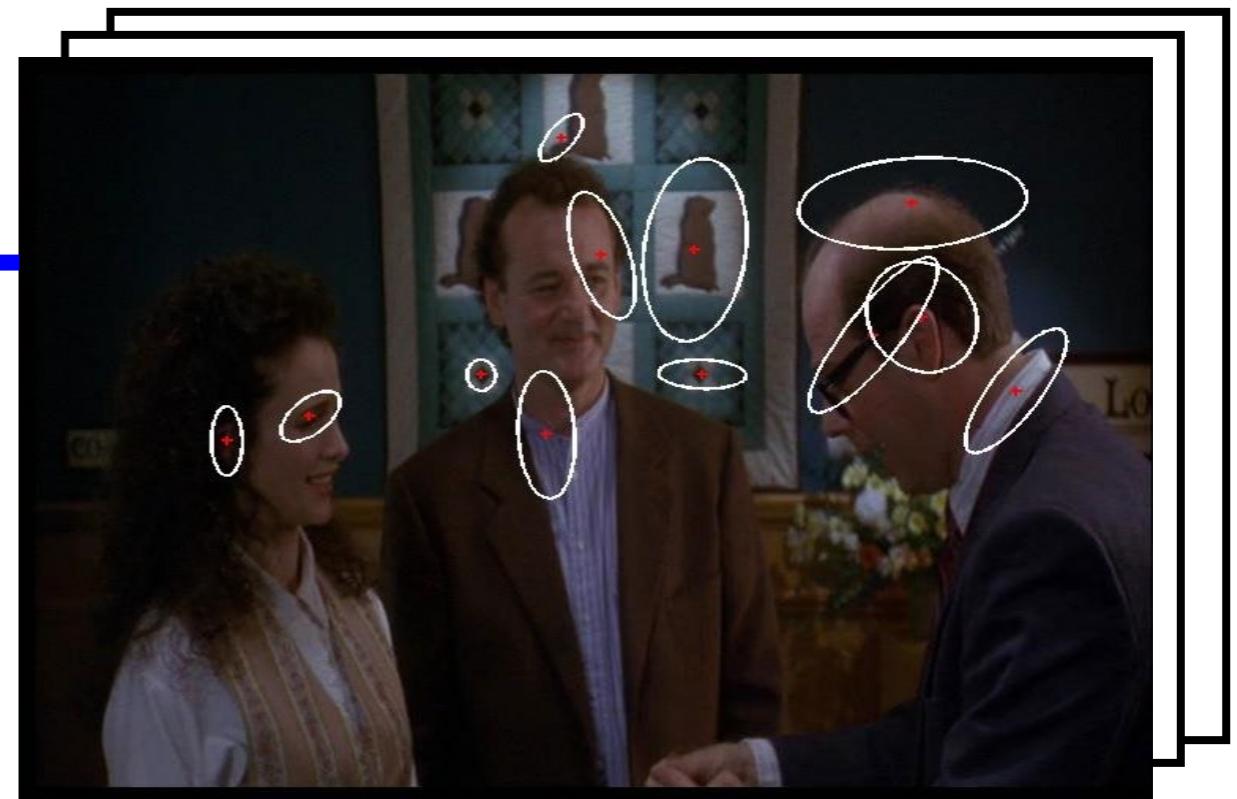
[Mikojaczyk and Schmid '02]

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

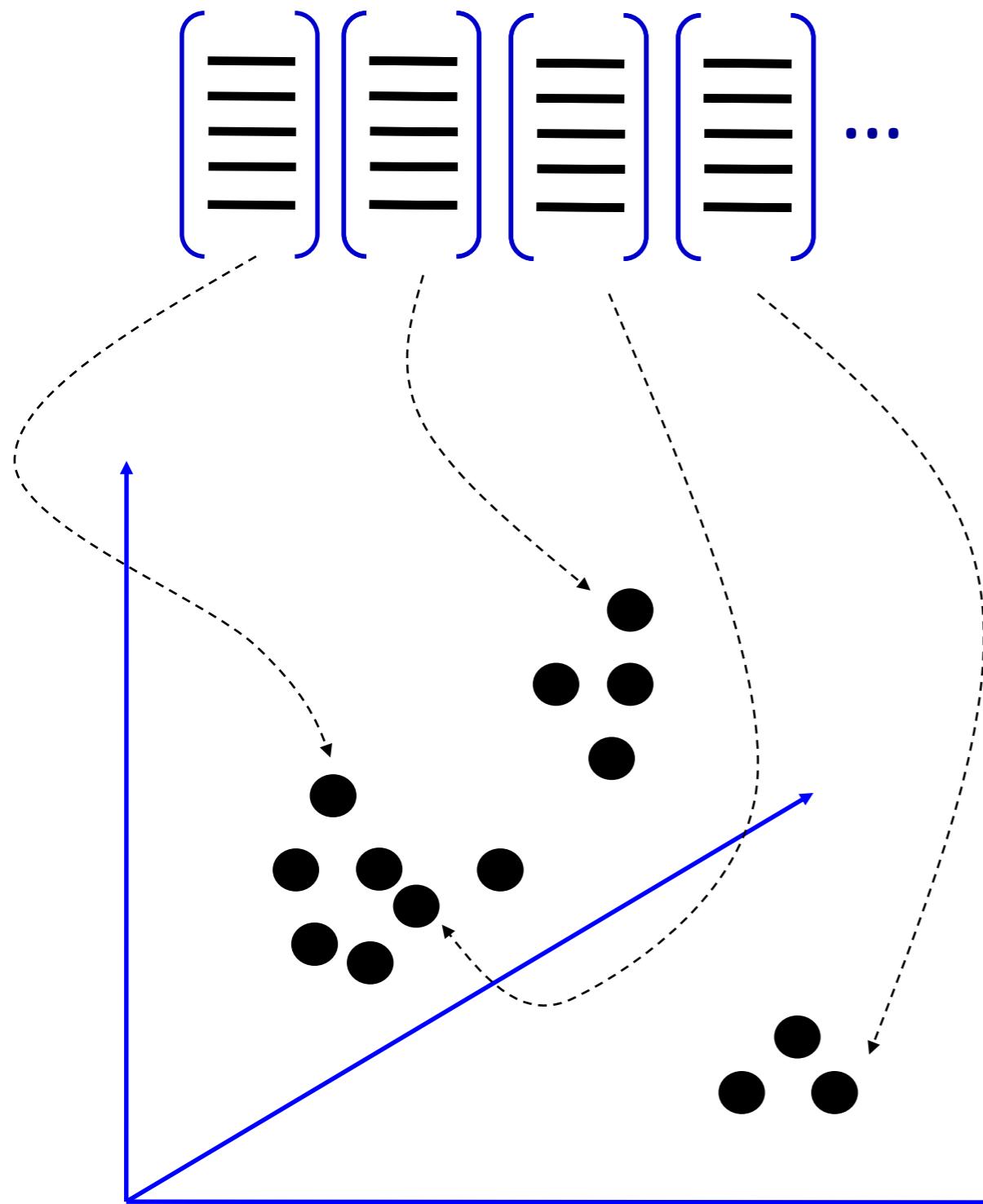
[Sivic & Zisserman, '03]

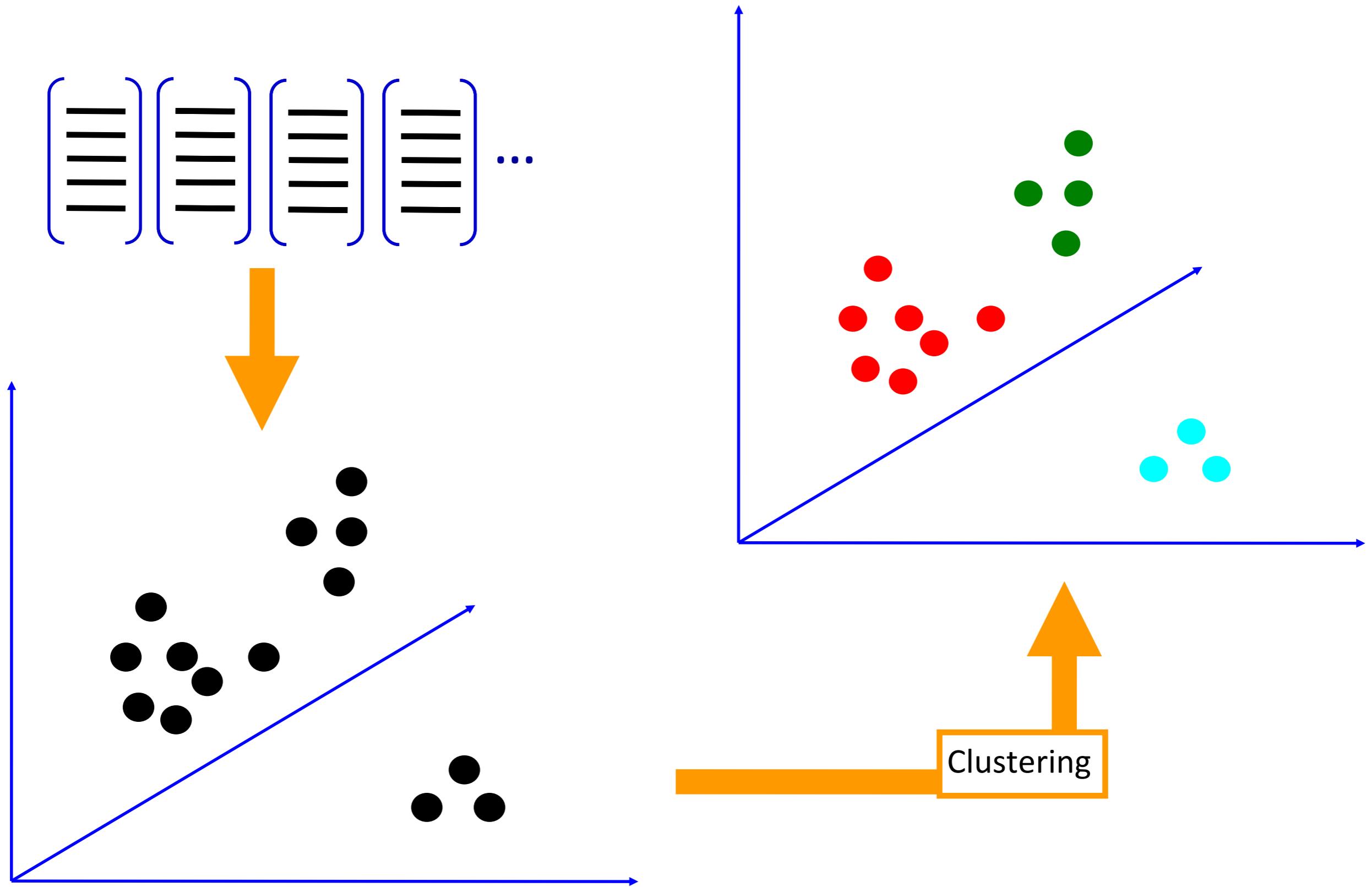


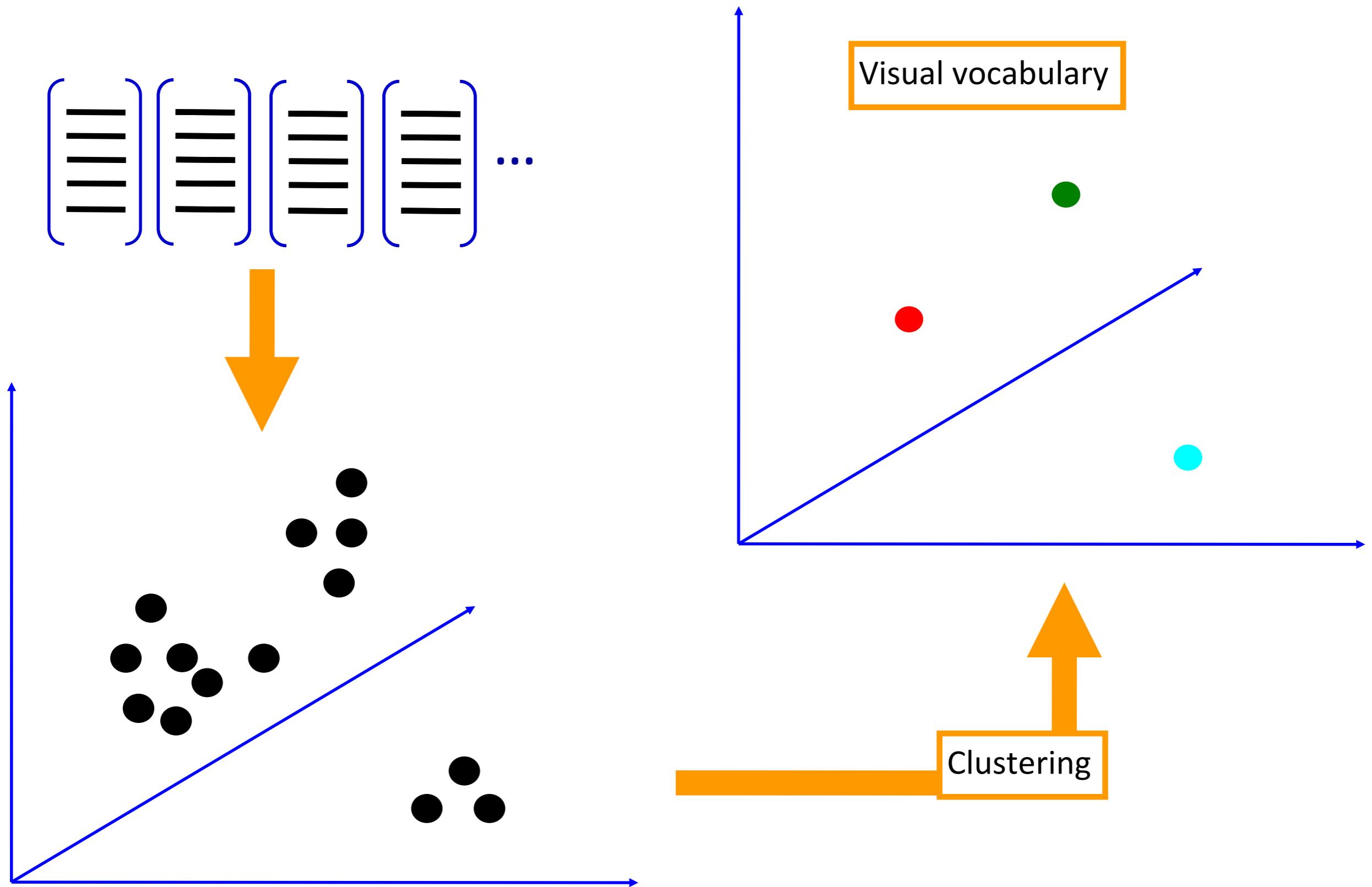
...



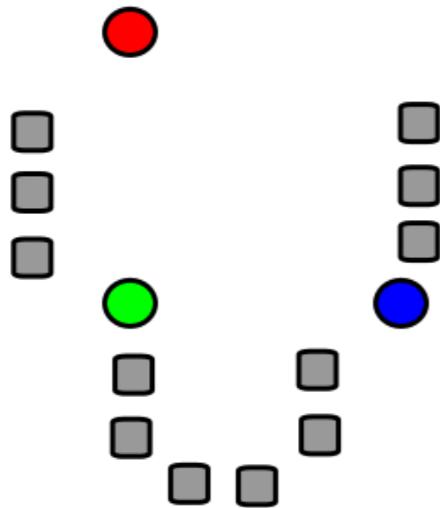
How do we learn the dictionary?



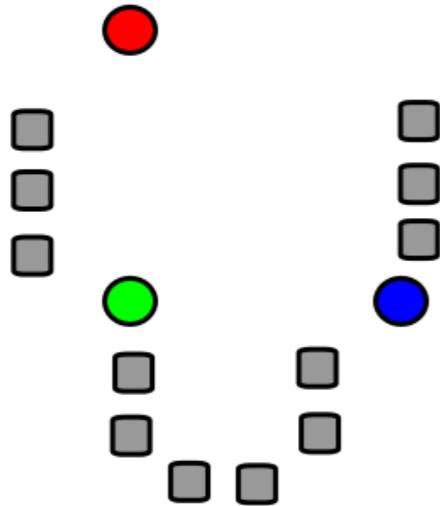




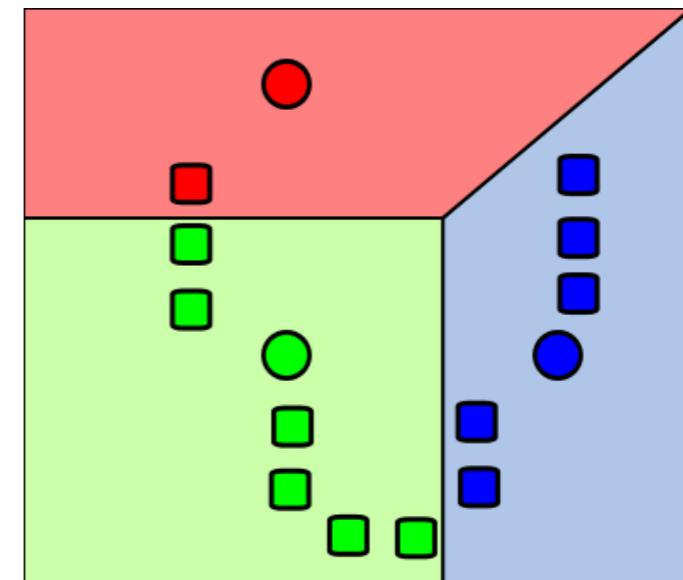
K-means clustering



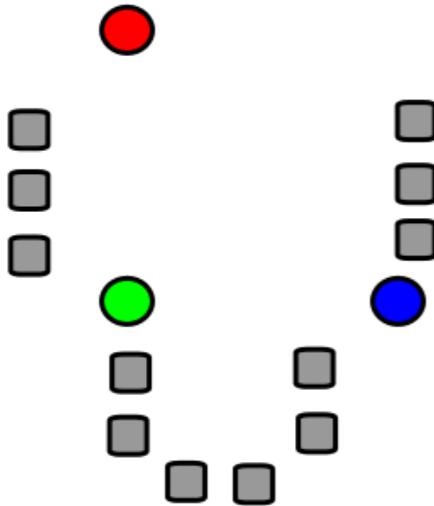
1. Select initial
centroids at random



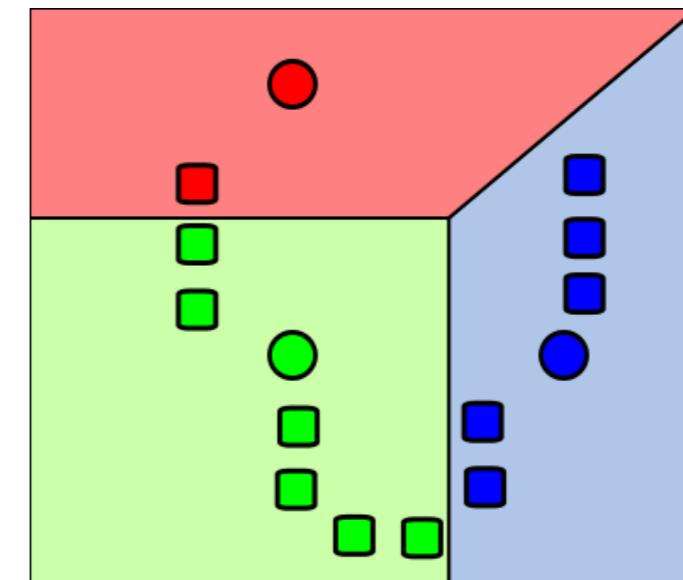
1. Select initial centroids at random



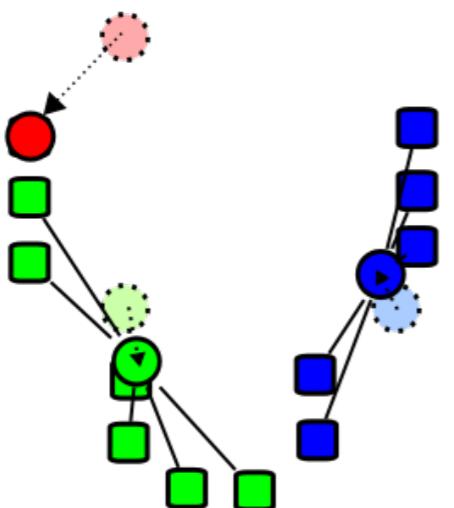
2. Assign each object to the cluster with the nearest centroid.



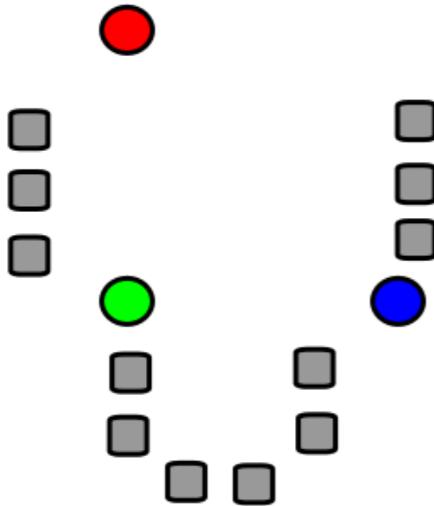
1. Select initial centroids at random



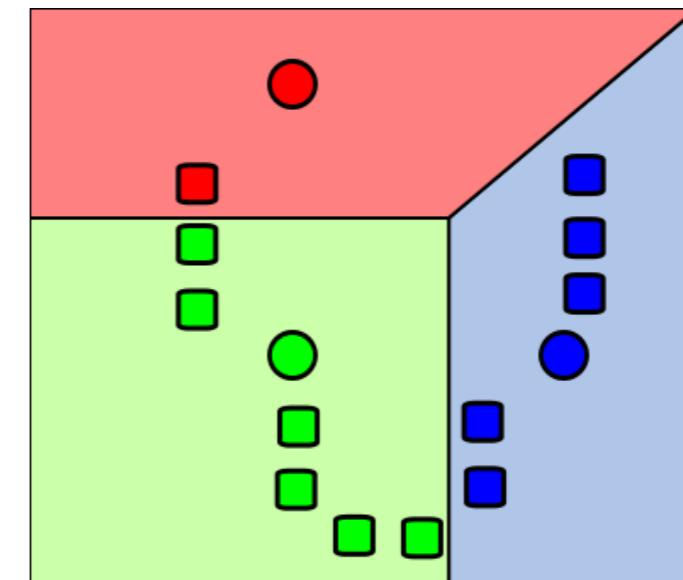
2. Assign each object to the cluster with the nearest centroid.



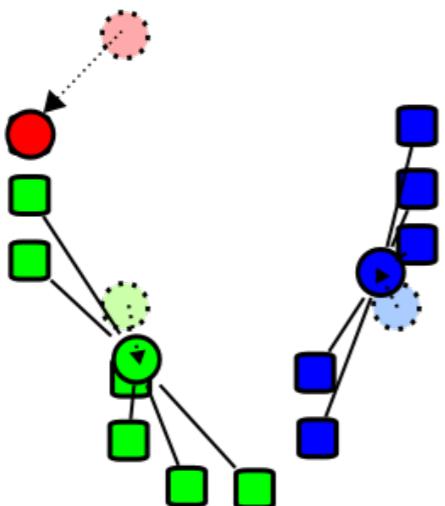
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



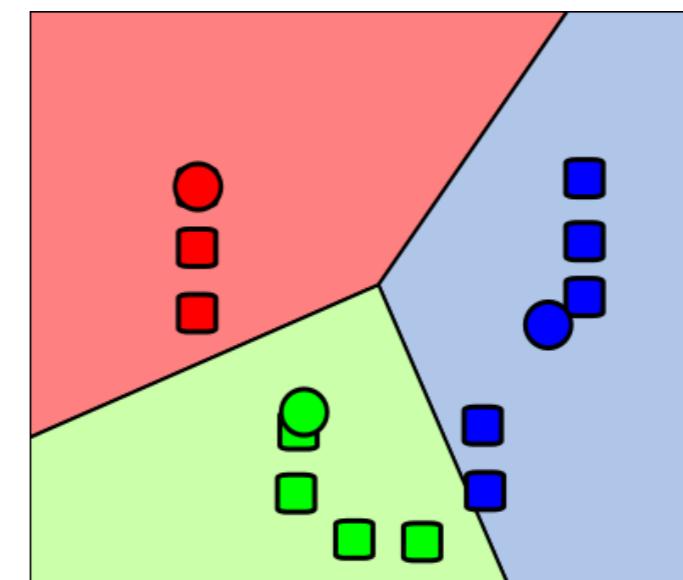
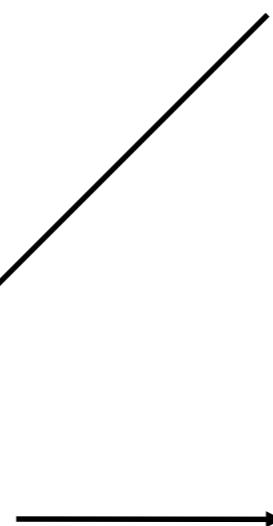
1. Select initial centroids at random



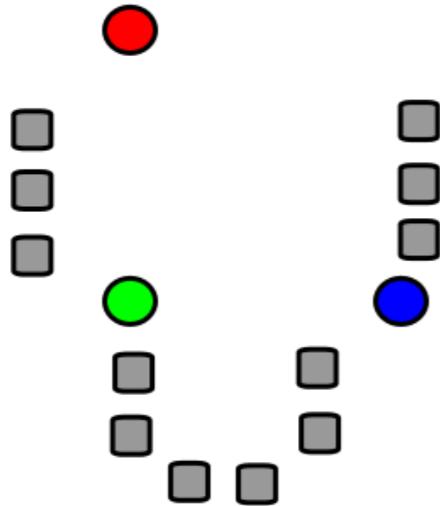
2. Assign each object to the cluster with the nearest centroid.



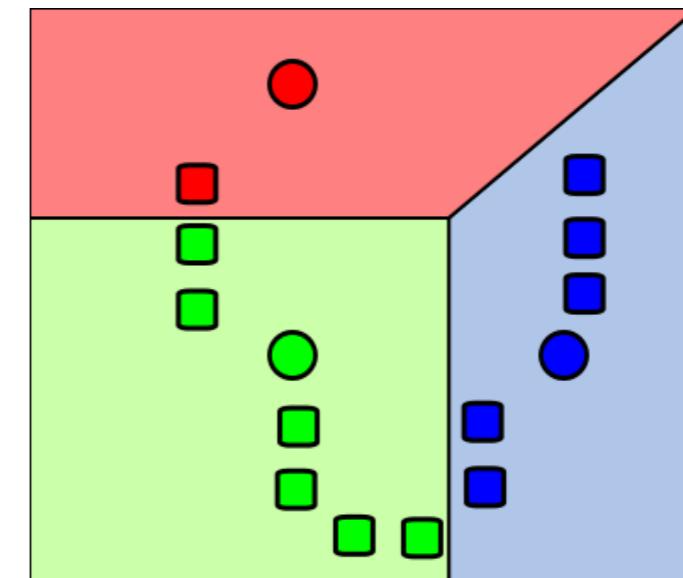
3. Compute each centroid as the mean of the objects assigned to it (go to 2)



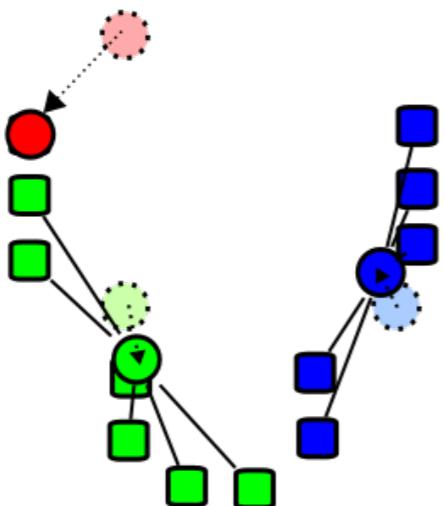
2. Assign each object to the cluster with the nearest centroid.



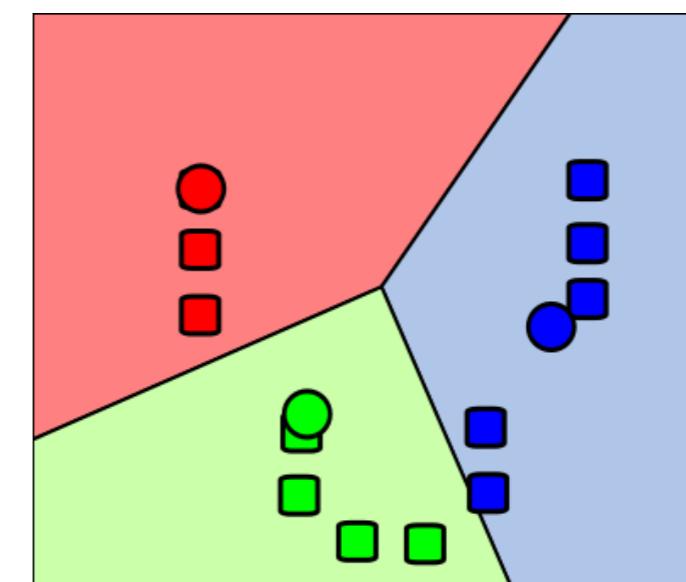
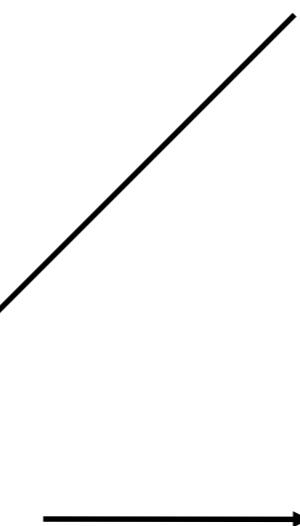
1. Select initial centroids at random



2. Assign each object to the cluster with the nearest centroid.



3. Compute each centroid as the mean of the objects assigned to it (go to 2)



2. Assign each object to the cluster with the nearest centroid.

K-means Clustering

Given k:

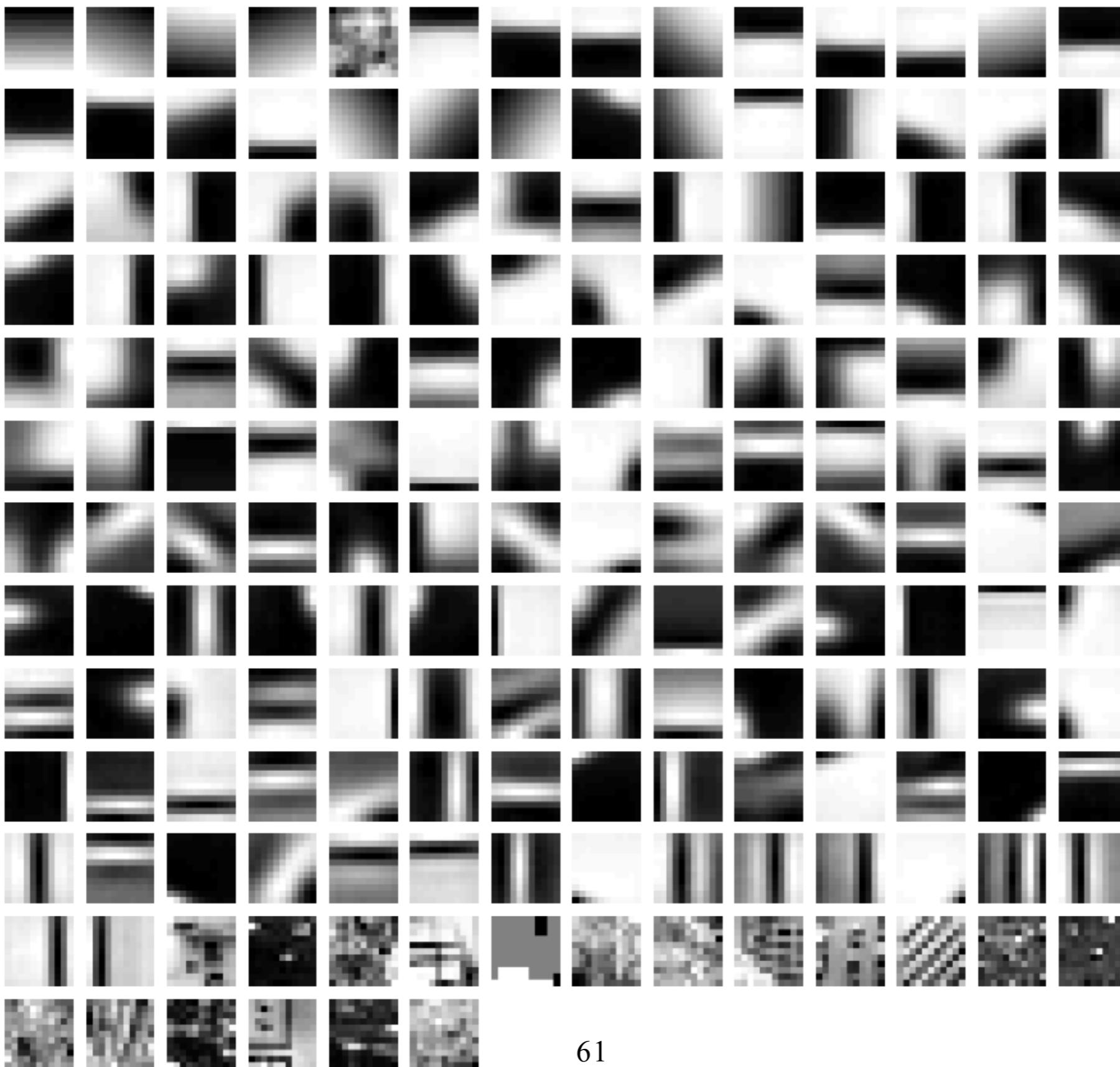
1. Select initial centroids at random.
2. Assign each object to the cluster with the nearest centroid.
3. Compute each centroid as the mean of the objects assigned to it.
4. Repeat previous 2 steps until no change.

*From what **data** should I learn the dictionary?*

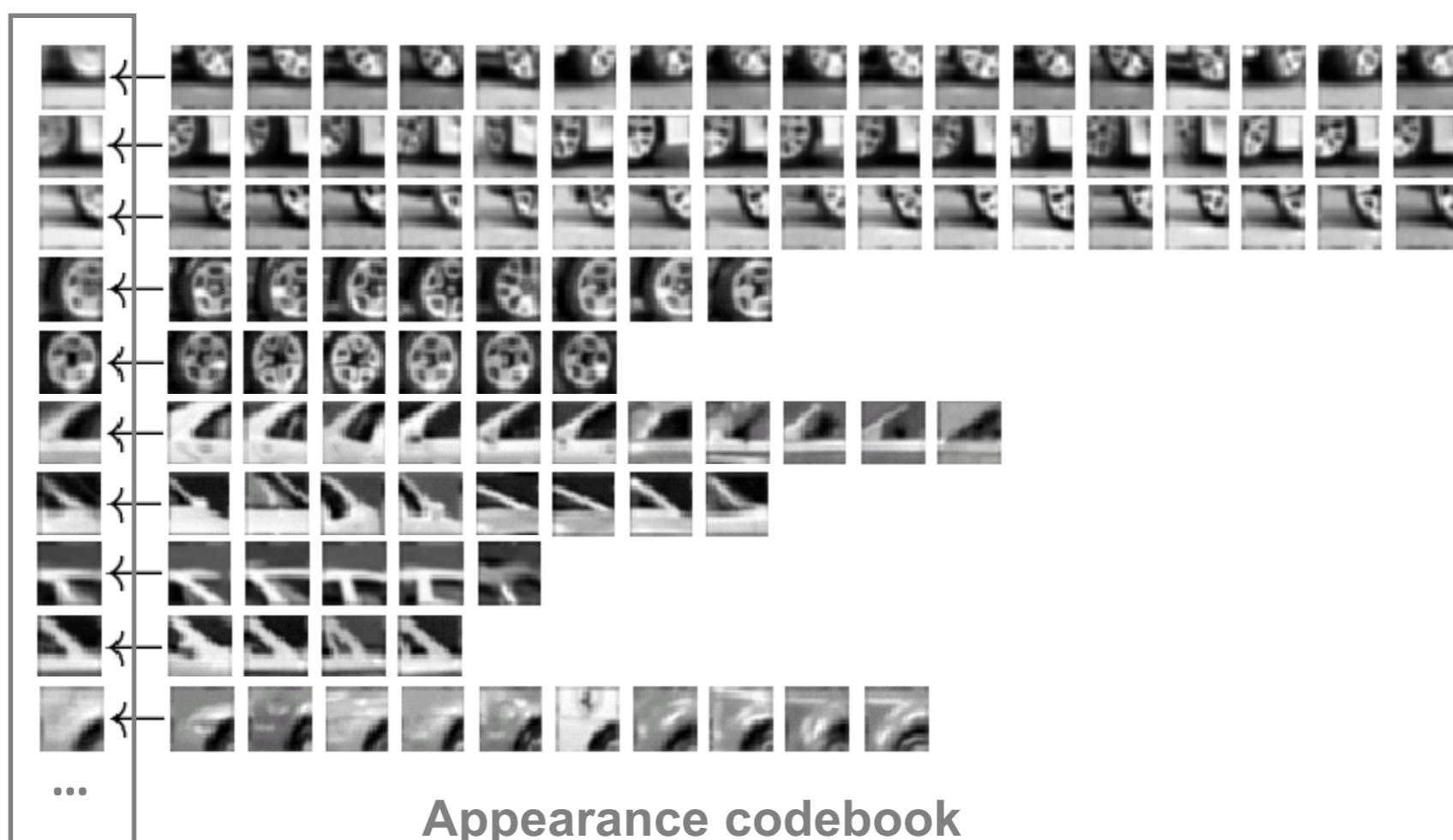
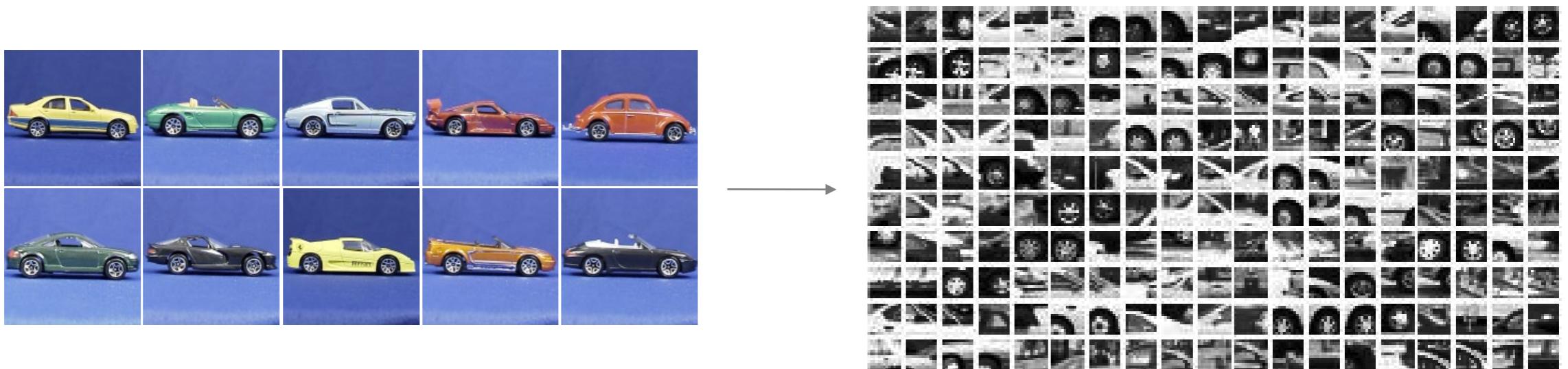
*From what **data** should I learn the dictionary?*

- Dictionary can be learned on separate training set
- Provided the training set is sufficiently representative, the dictionary will be “universal”

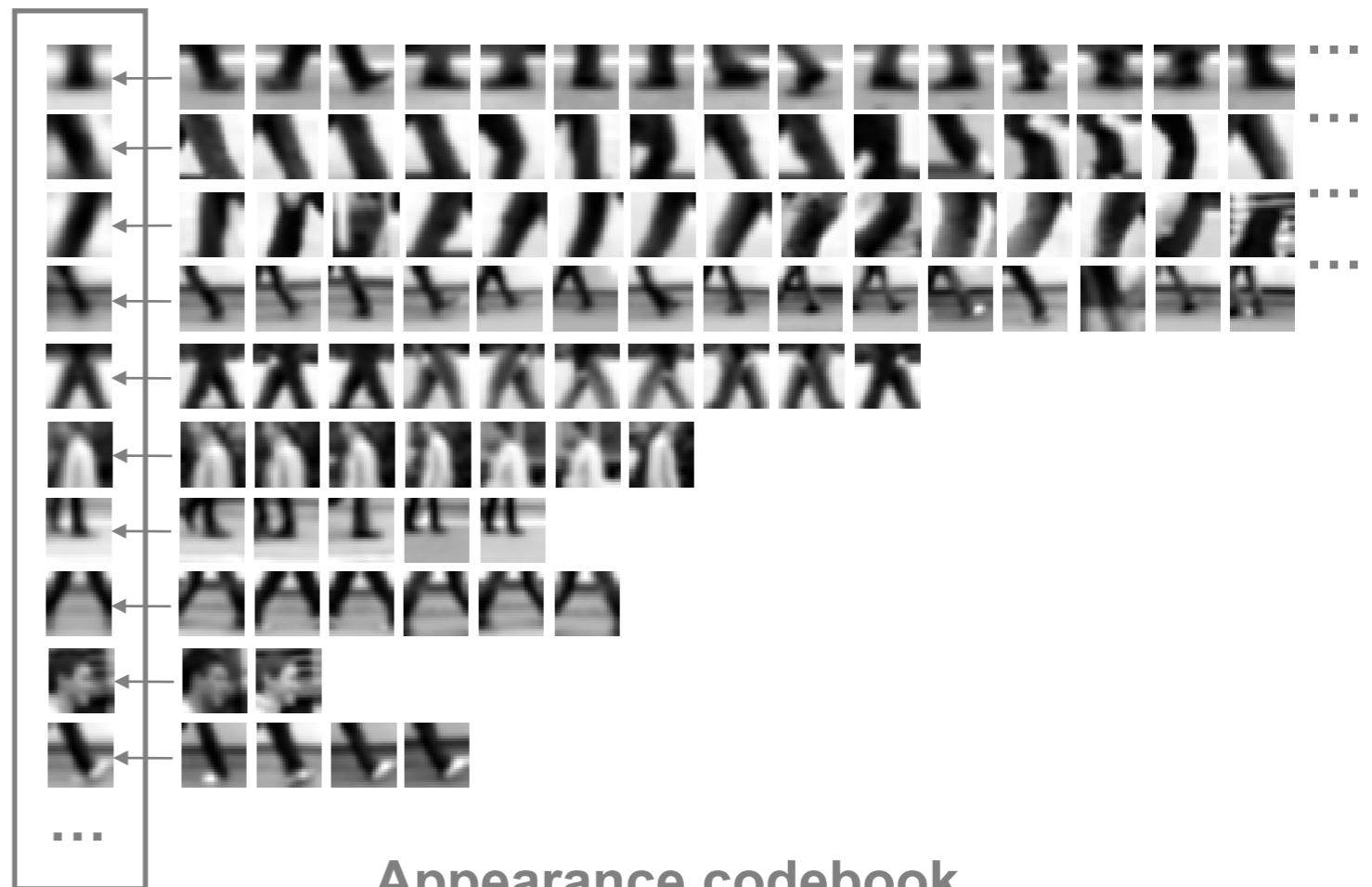
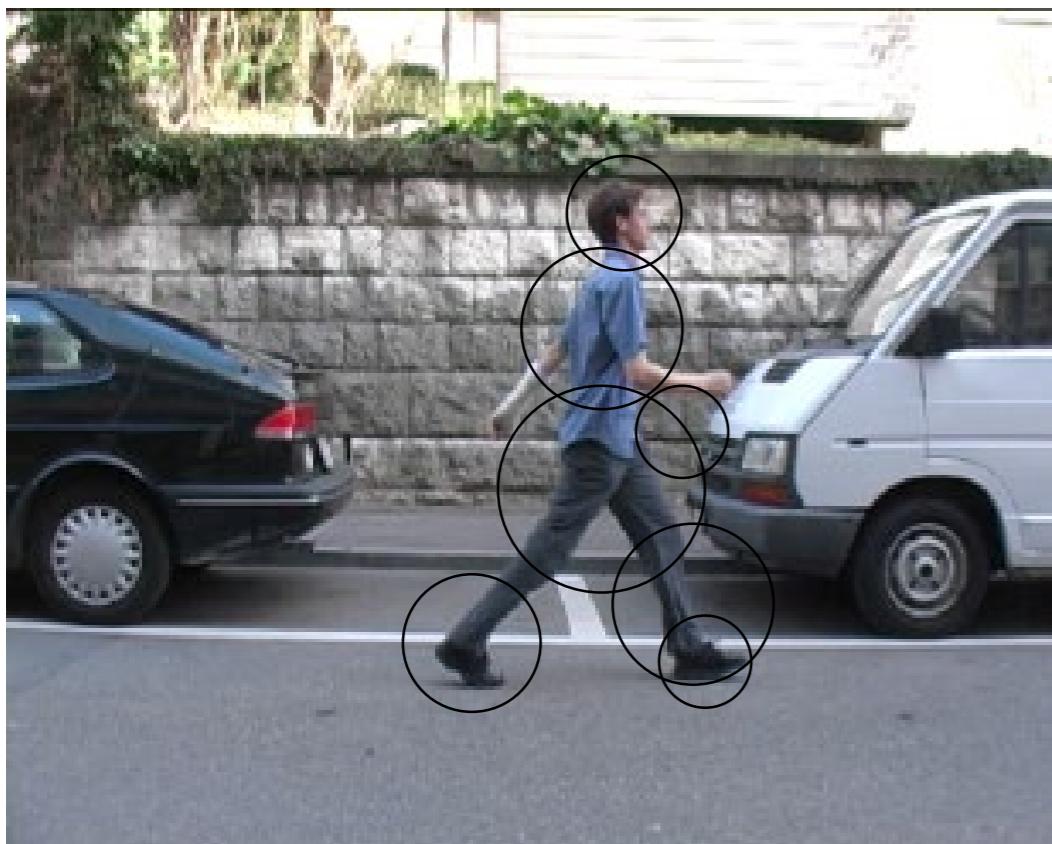
Example visual dictionary



Example dictionary



Another dictionary

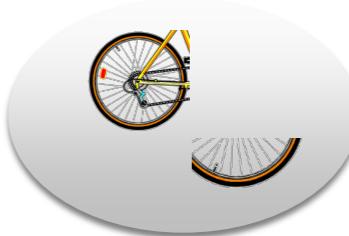


Appearance codebook

Dictionary Learning: Learn Visual Words using clustering

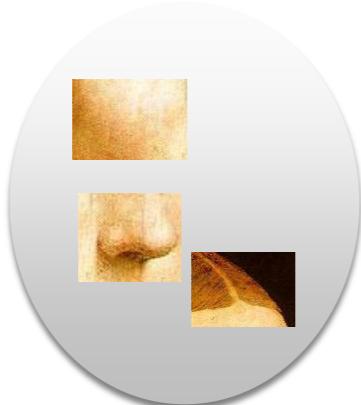
Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs



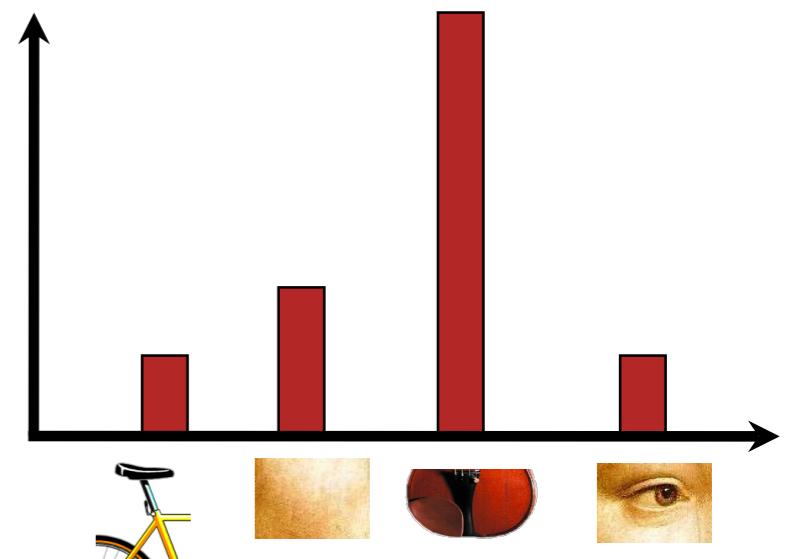
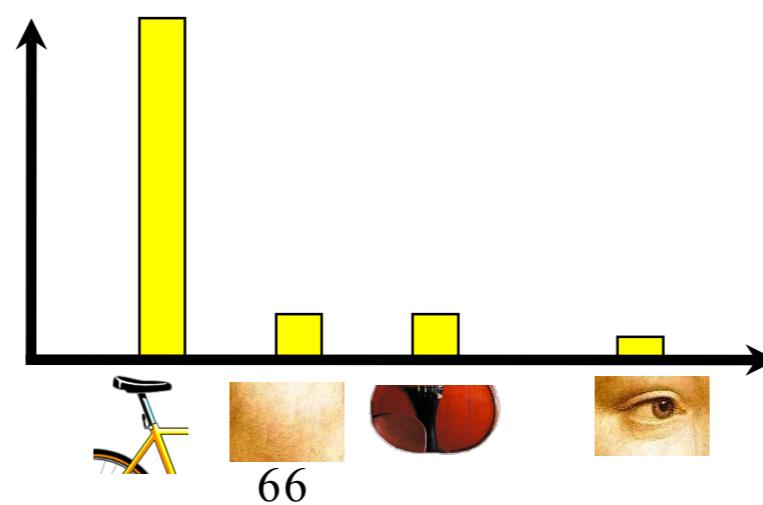
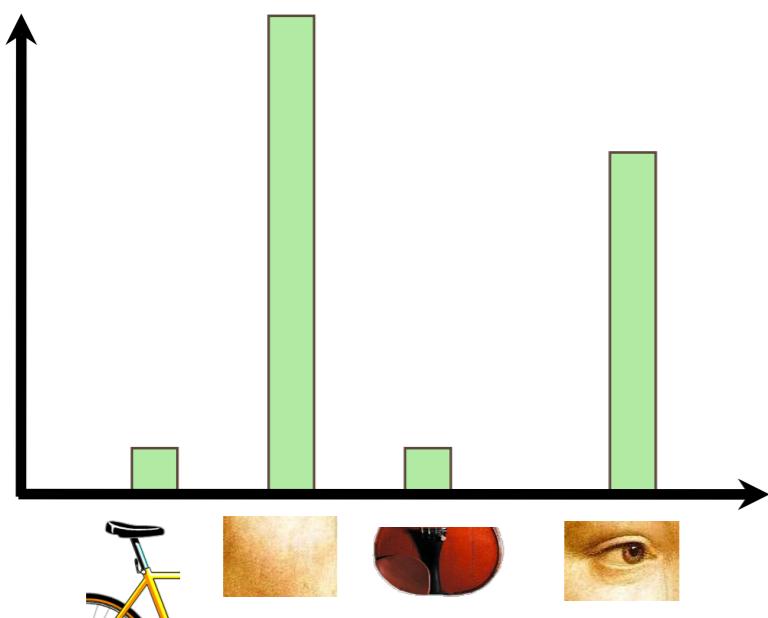
1. Quantization: image features gets associated to a visual word (nearest cluster center)

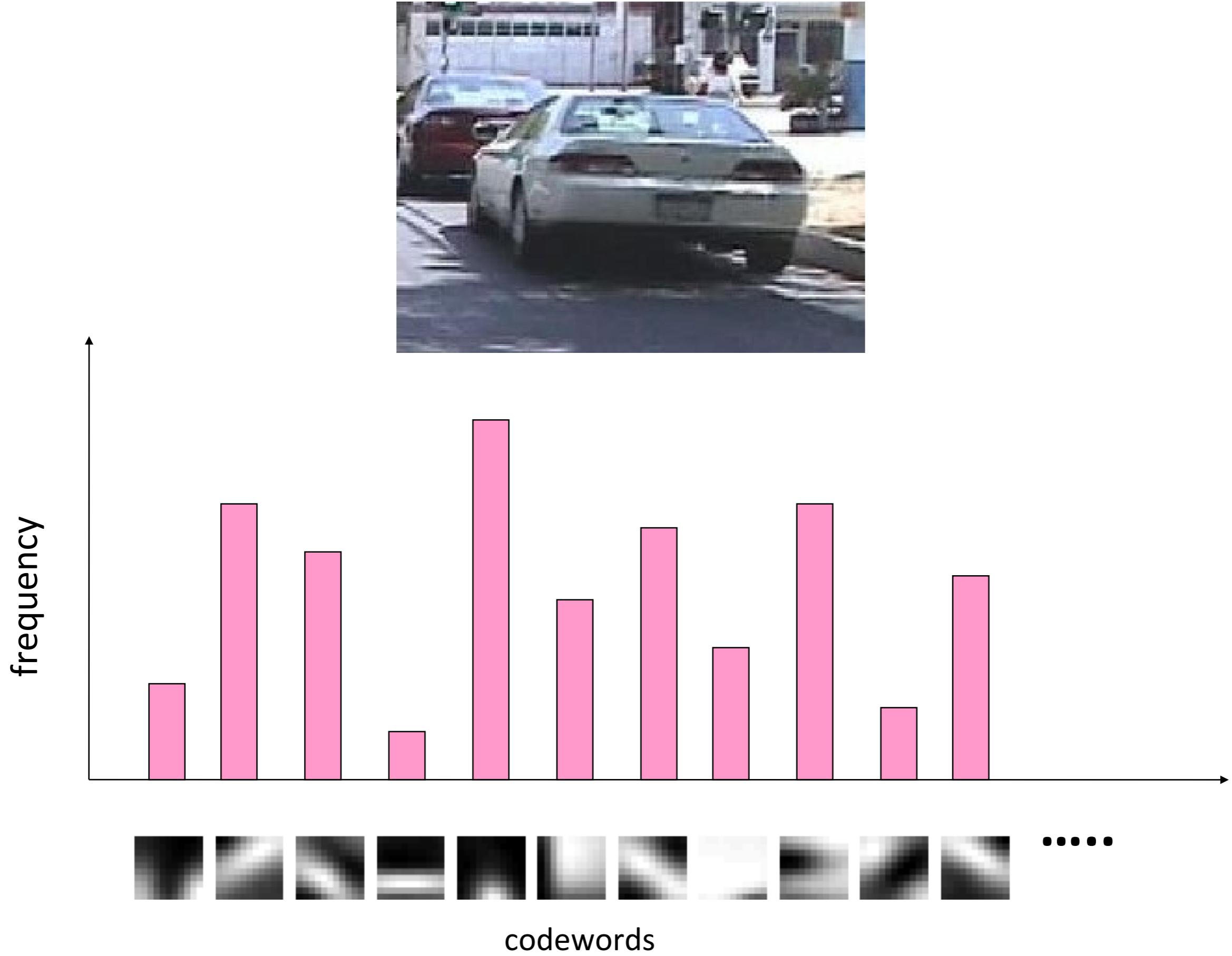
Encode:
build Bags-of-Words (BOW) vectors
for each image



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

2. Histogram: count the
number of visual word
occurrences

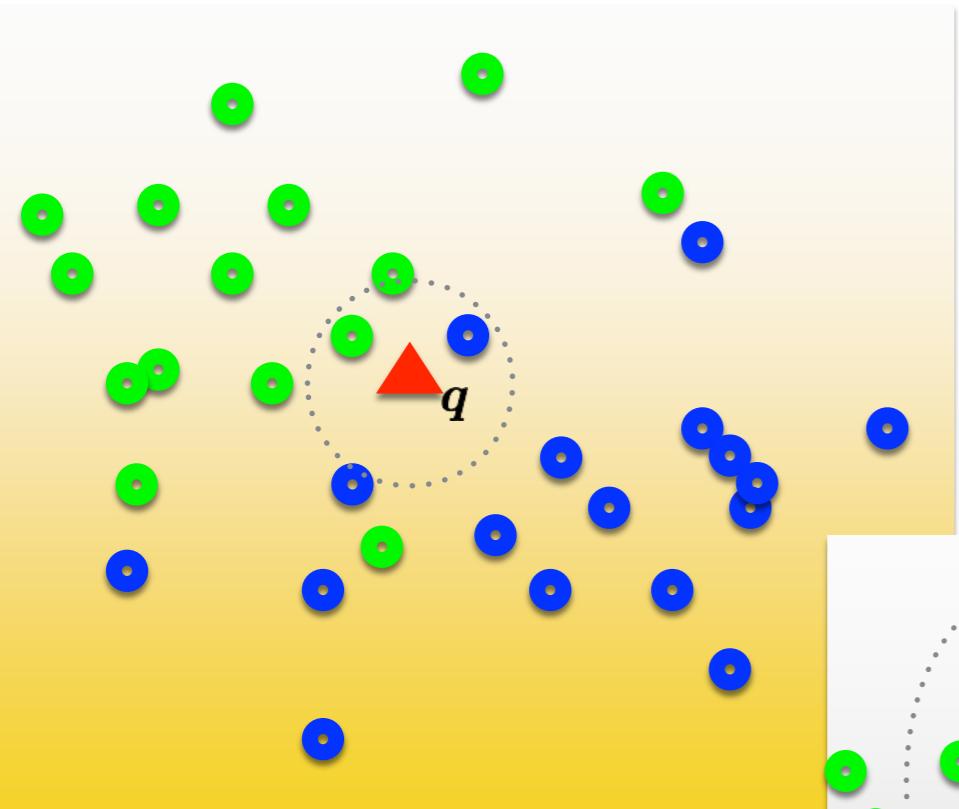




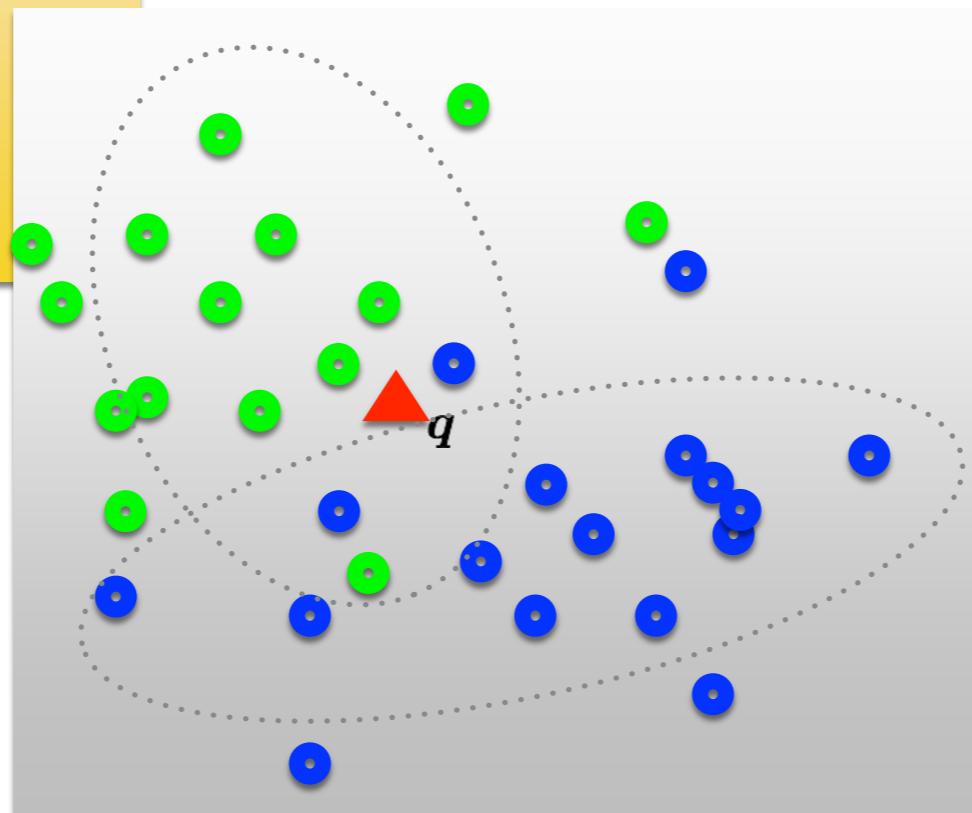
Dictionary Learning: Learn Visual Words using clustering

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

Classify:
Train and test data using BOWs

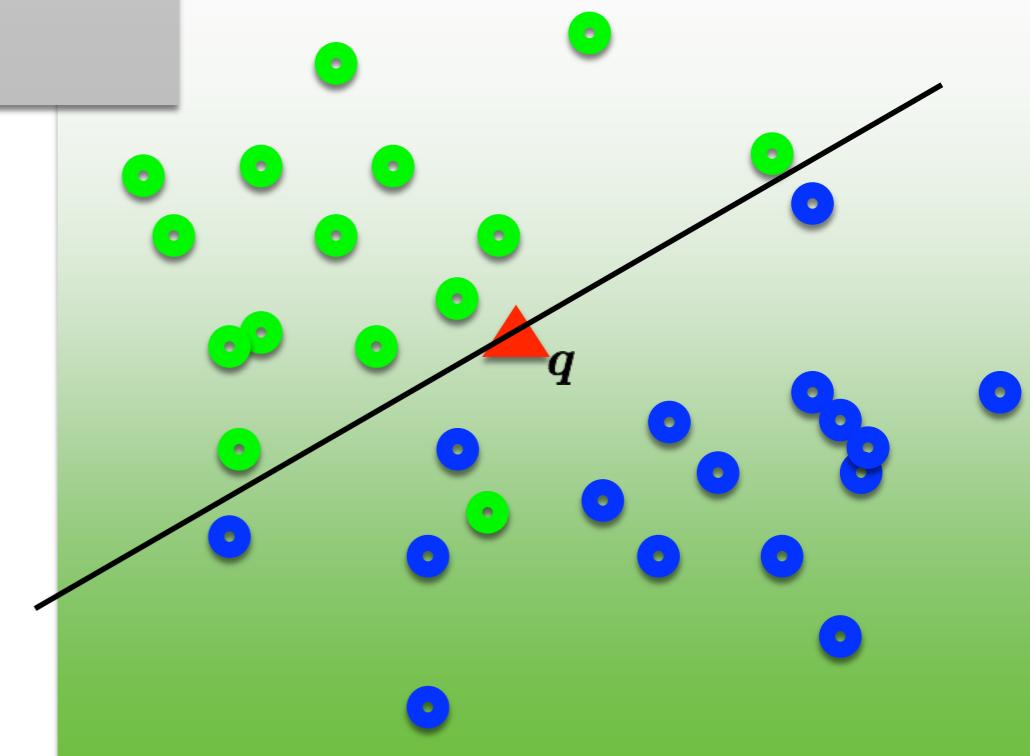


Knearest neighbors



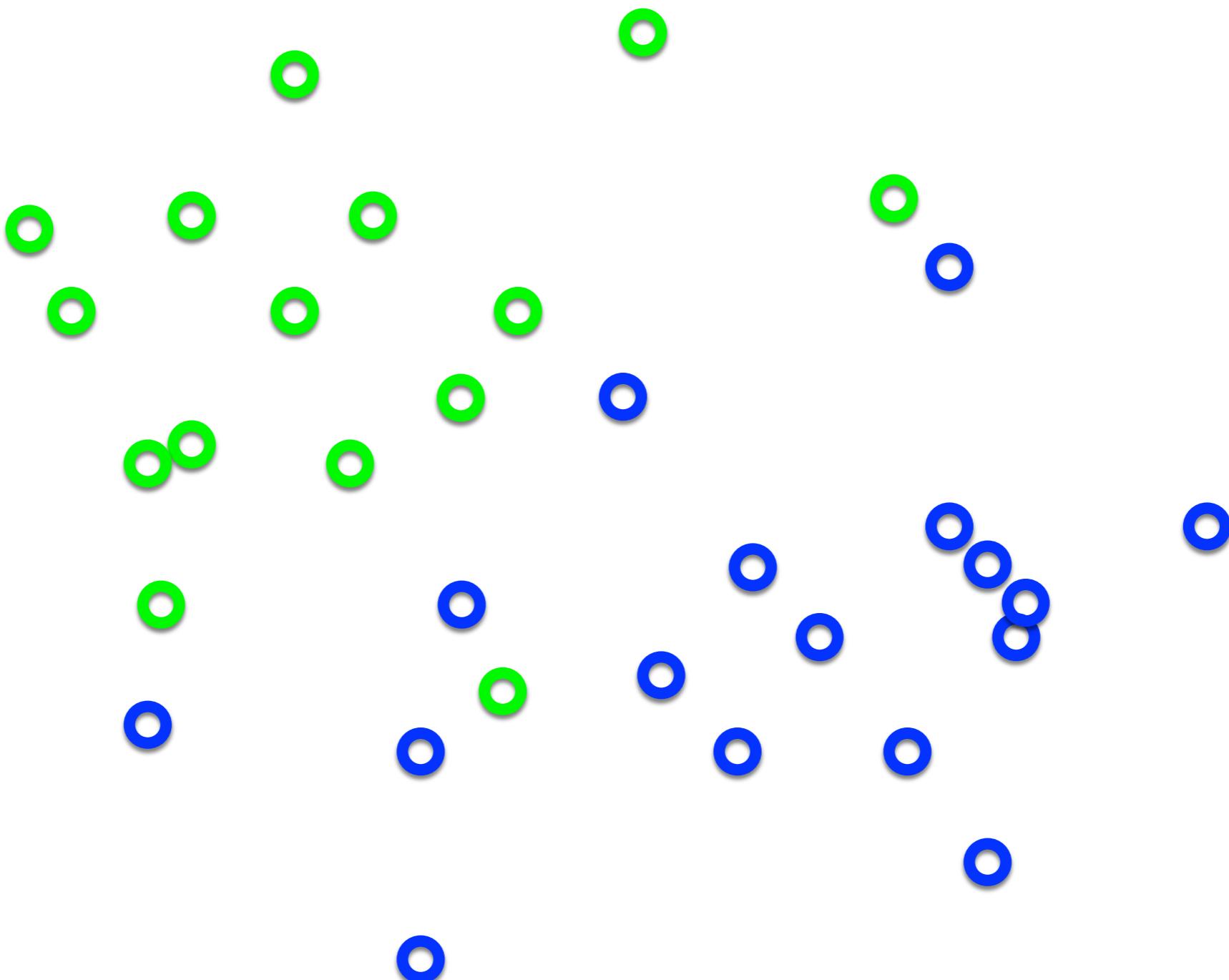
Naïve Bayes

Support Vector Machine

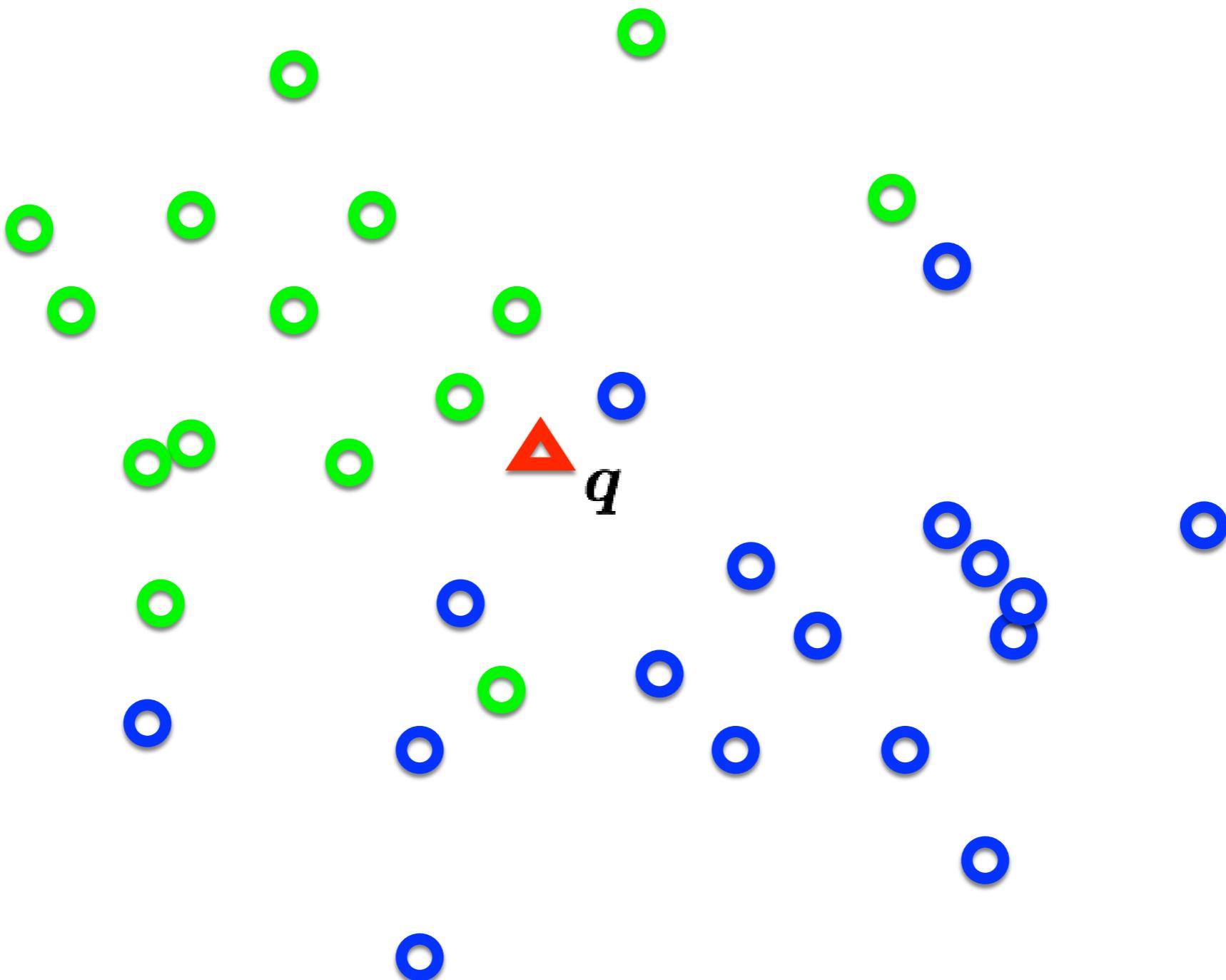


K nearest neighbors

Distribution of data from two classes

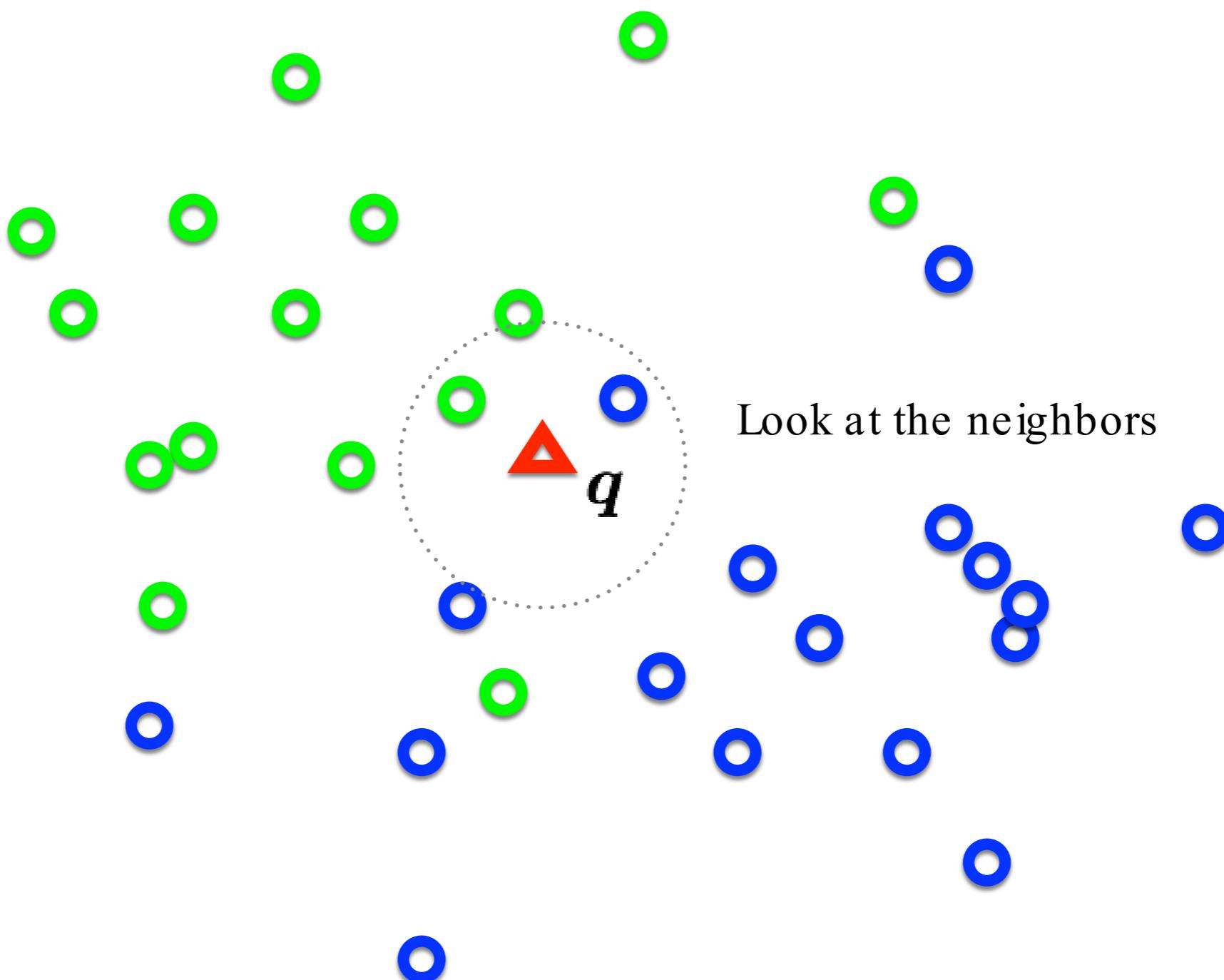


Distribution of data from two classes

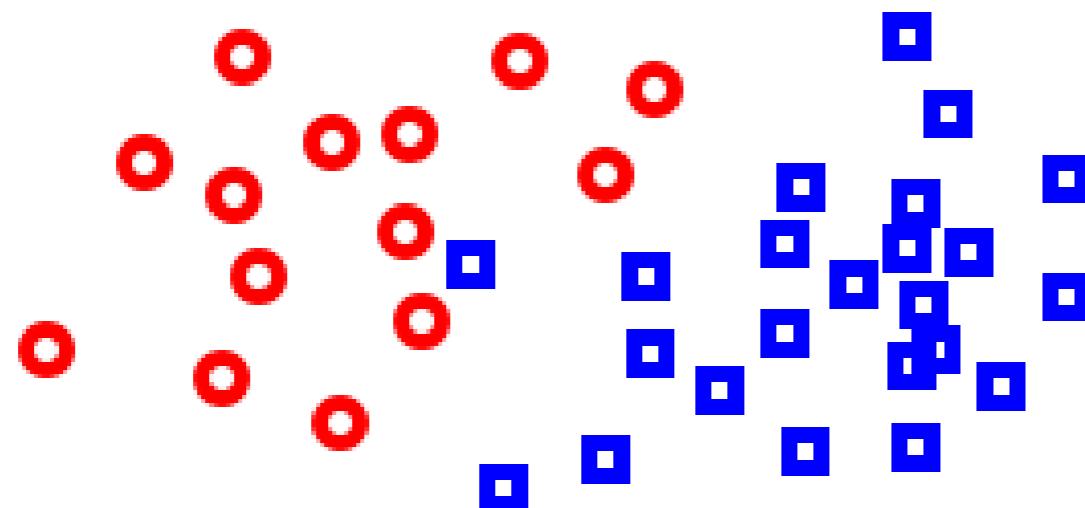


Which class does q belong to?

Distribution of data from two classes



K-Nearest Neighbor (KNN) Classifier



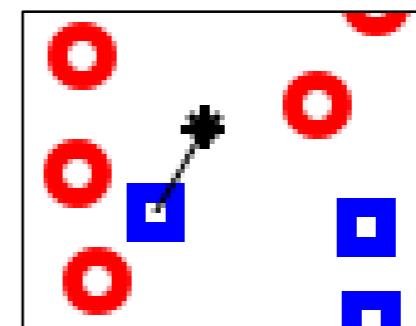
For a given query point q ,
assign the class of the nearest
neighbor

Compute the k nearest
neighbors and assign the
class by majority vote.

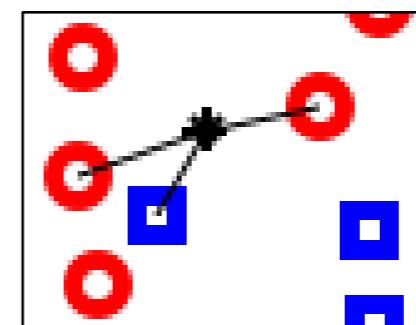
Non-parametric pattern classification
approach

Consider a two class problem where
each sample consists of two
measurements (x, y) .

$k = 1$



$k = 3$



Nearest Neighbor is competitive

40281508803272726475529284686500876171127400776386420940578274711366
 50711167679664143112241082634006330117113109975414895351982339901029
 8468682467933943144705960444612336459685608641865284554770782237018
 76953465018828357808571101378507110114527623028596972136418240510226
 93477149064842728100783331376131665747595849918501320348220251514889
 82049962335648092836457294912840709116759914592504108908989425198980
 35517216919955162286714604033223689853854520563283995794671313660901
 94368160413174951001162198403649071657525185470670258104571851900607
 8857389886823975629288168879180172075190209862393802111142972512199
 148534347507488153959769036398212868553941251514414435912233029009
 931909754920105149336152522026601203025579580895032590884884546549
 6928545799216340783934656219260061287982047750564674307507420899404
 12845278113035703193631773084826529739099642972116747590821445161325
 90666367728608302983253980019513960141712379749939282718091017796999
 21010452828351781129784050788477858498138031795516574935471208160734
 2830878408445856630937689349589128868137901147081745712113621280766
 41992780136134111560707232522949812161274000822922799275134941856283

MNIST Digit Recognition

- Handwritten digits
- 28x28 pixel images: $d = 784$
- 60,000 training samples
- 10,000 test samples

Yann LeCunn

	Test Error Rate (%)
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
K-NN, Tangent Distance, 16x16	1.1
K-NN, shape context matching	0.67
1000 RBF + linear classifier	3.6
SVM deg 4 polynomial	1.1
2-layer NN, 300 hidden units	4.7
2-layer NN, 300 HU, [deskewing]	1.6
LeNet-5, [distortions]	0.8
Boosted LeNet-4, [distortions]	0.7

What is the best distance metric between data points?

- Typically Euclidean distance
- Locality sensitive distance metrics
- Important to normalize.
Dimensions have different scales

How many K?

- Typically $k=1$ is good
- Cross-validation (try different k !)

Distance metrics

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + \cdots + (x_N - y_N)^2} \quad \text{Euclidean}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{x_1 y_1 + \cdots + x_N y_N}{\sqrt{\sum_n x_n^2} \sqrt{\sum_n y_n^2}} \quad \text{Cosine}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_n \frac{(x_n - y_n)^2}{(x_n + y_n)} \quad \text{Chi-squared}$$

Choice of distance metric

- Hyperparameter

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

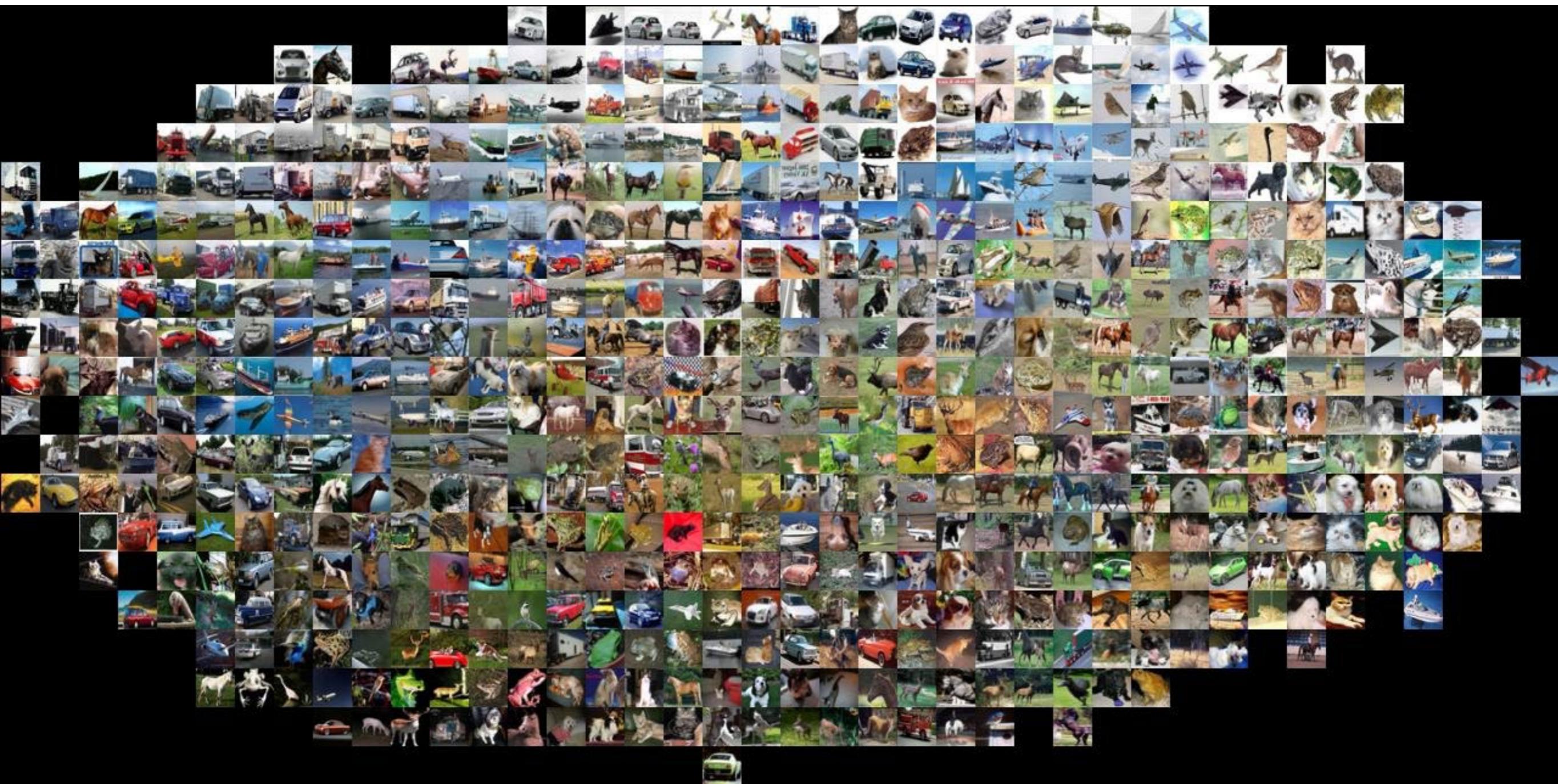
L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

- Two most commonly used special cases of p-norm

$$\|x\|_p = \left(|x_1|^p + \cdots + |x_n|^p \right)^{\frac{1}{p}} \quad p \geq 1, x \in \mathbb{R}^n$$

Visualization: L2 distance



CIFAR-10 and NN results

Example dataset: **CIFAR-10**

10 labels

50,000 training images

10,000 test images.

airplane



automobile



bird



cat



deer



dog



frog



horse



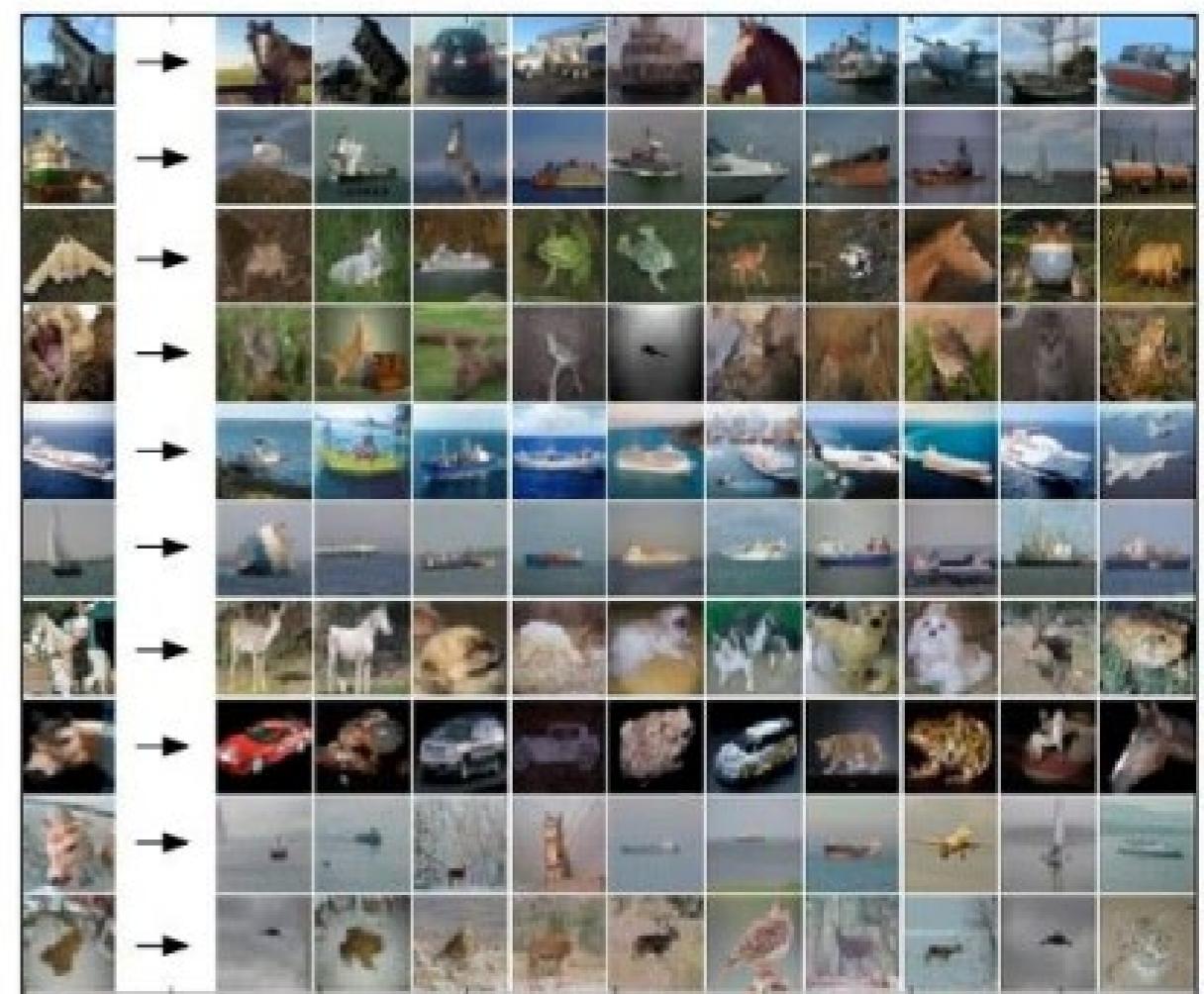
ship



truck

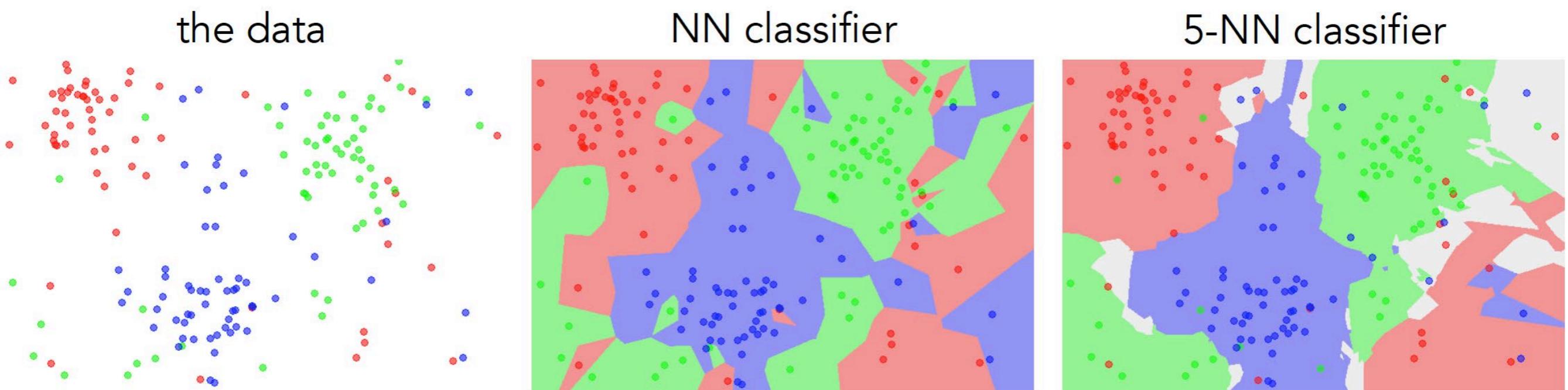


For every test image (first column),
examples of nearest neighbors in rows



k-nearest neighbor

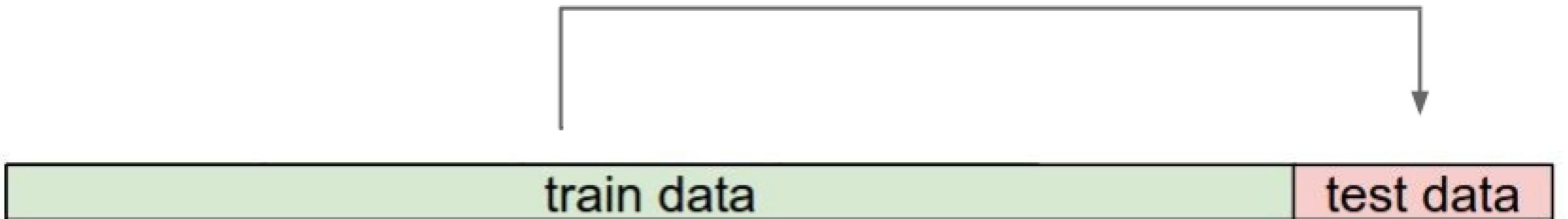
- Find the k closest points from training data
- Labels of the k points “vote” to classify



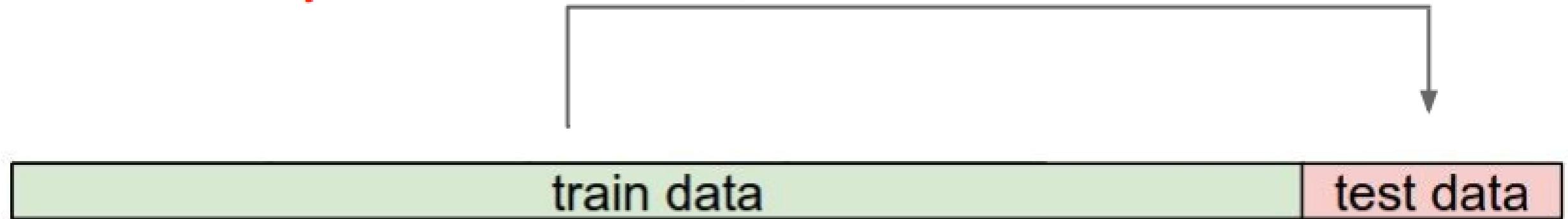
Hyperparameters

- What is the best distance to use?
- What is the best value of k to use?
- i.e., how do we set the hyperparameters?
- Very problem-dependent
- Must try them all and see what works best

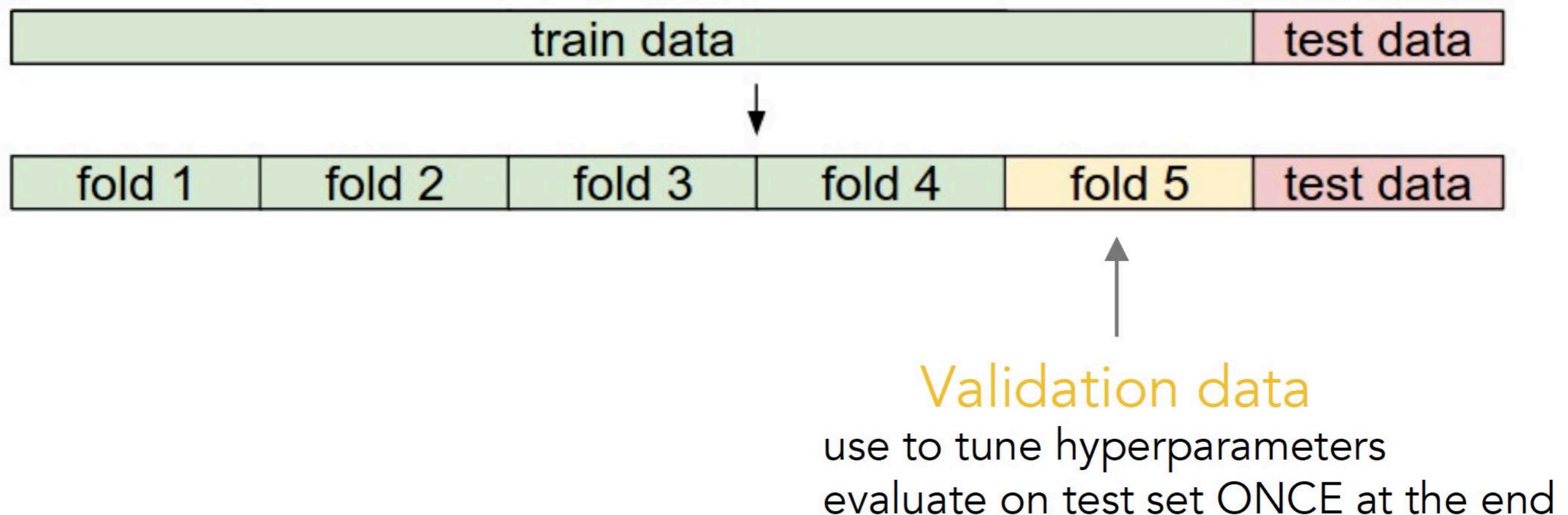
Try out what hyperparameters work best on test set.



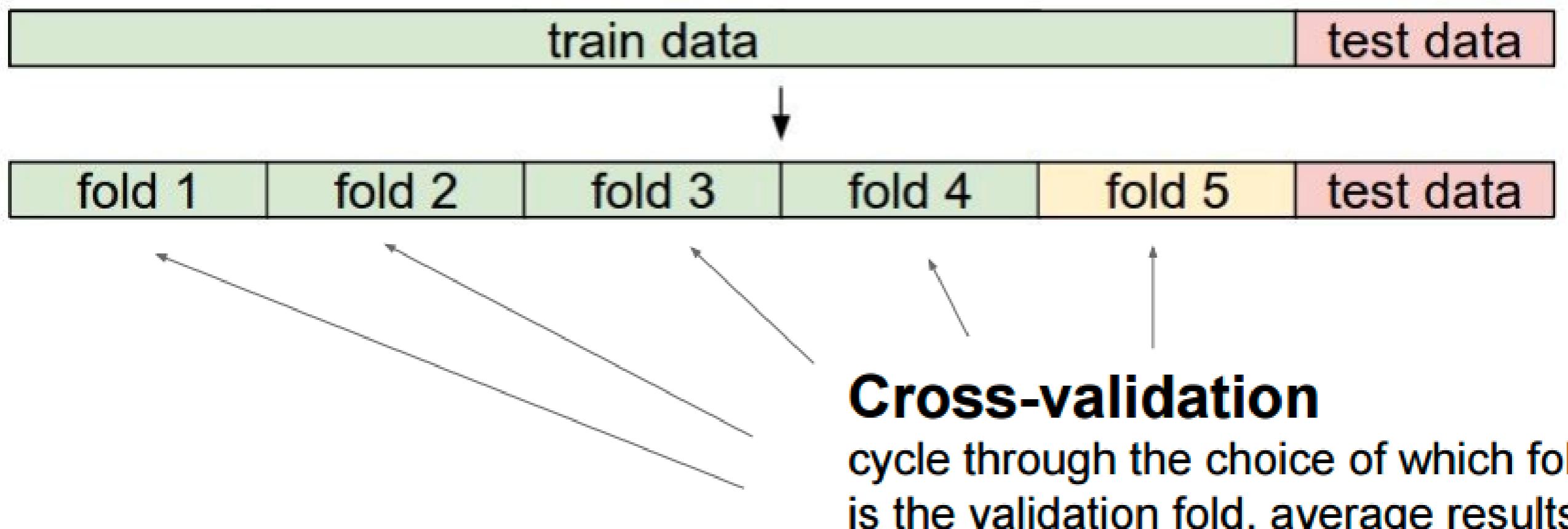
Trying out what hyperparameters work best on test set:
Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY**, at the end.

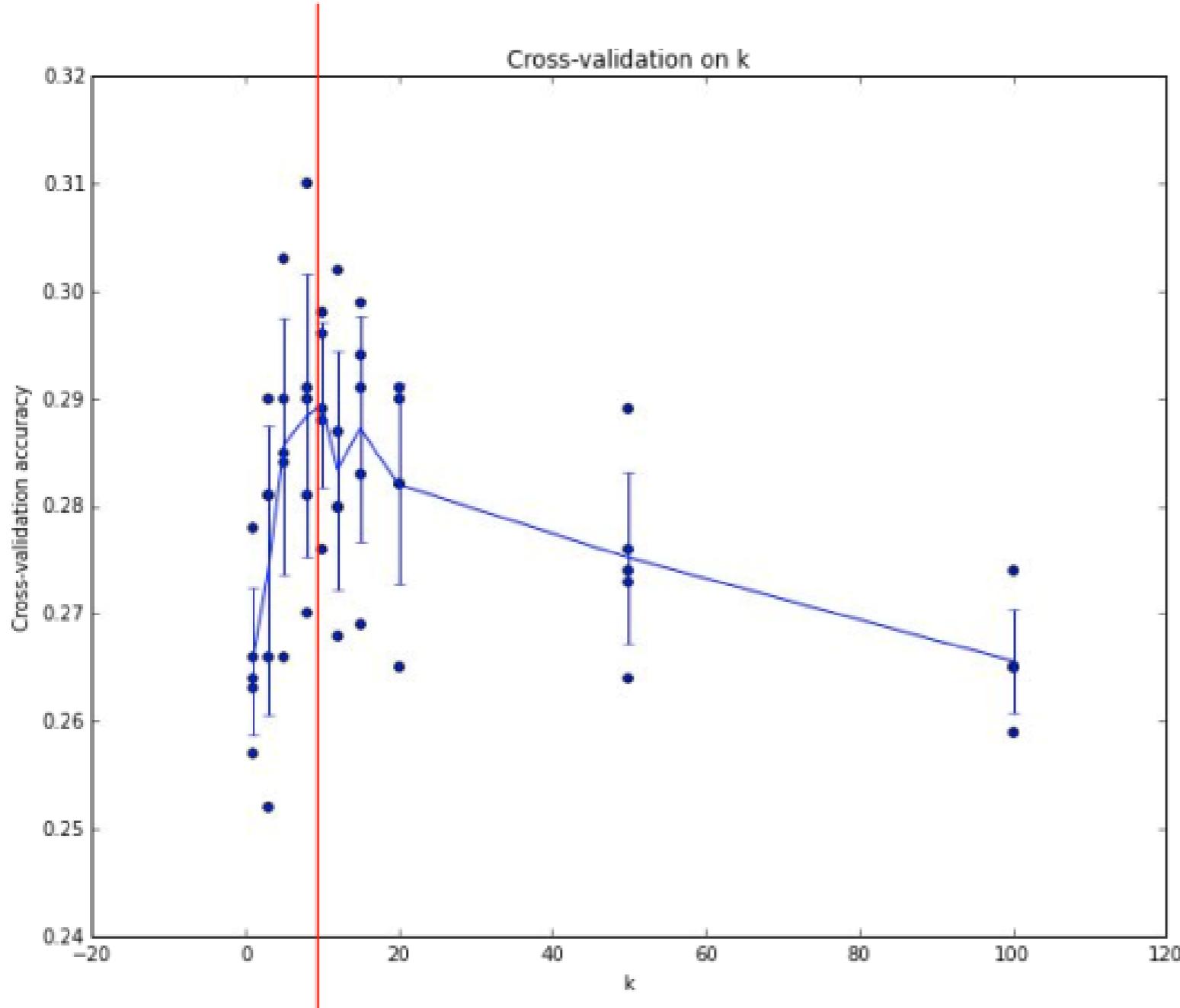


Validation



Cross-validation





Example of
5-fold cross-validation
for the value of k .

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)

How to pick hyperparameters?

- Methodology
 - Train and test
 - Train, validate, test
- Train for original model
- Validate to find hyperparameters
- Test to understand generalizability

Pros

- simple yet effective

Cons

- search is expensive (can be sped-up)
- storage requirements
- difficulties with high-dimensional data

Timeline of recognition

Before deep learning

- 1965-late 1980s: alignment, **geometric** primitives
 - Early 1990s: **invariants**, appearance-based methods
 - Mid-late 1990s: **sliding window** approaches
 - Late 1990s: **feature**-based methods
 - Early 2000s: **parts-and-shape** models
 - 2003 – 2010: **bags of features**
-

After deep learning

- Present trends: combination of local and global methods, modeling context, integrating recognition and segmentation, **deep learning**

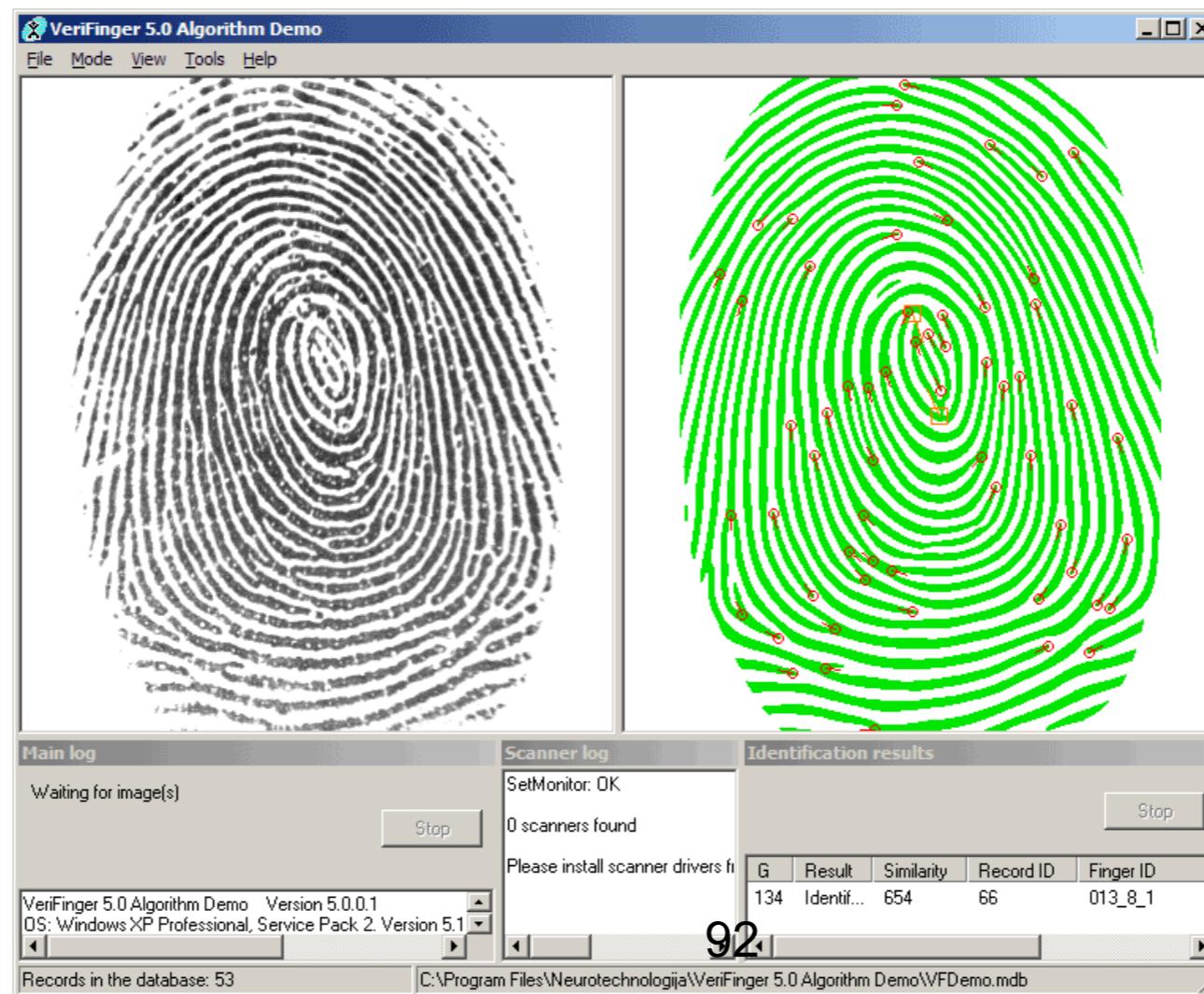
What “works” today

- Reading license plates, zip codes, checks

3 6 8 1 7 9 6 6 9 1
6 7 5 7 8 6 3 4 8 5
2 1 7 9 7 1 2 8 4 6
4 8 1 9 0 1 8 8 9 4
7 6 1 8 6 4 1 5 6 0
7 5 9 2 6 5 8 1 9 7
2 2 2 2 2 3 4 4 8 0
0 2 3 8 0 7 3 8 5 7
0 1 4 6 4 6 0 2 4 3
7 1 2 8 7 6 9 8 6 1

What “works” today

- Reading license plates, zip codes, checks
- Fingerprint recognition



What “works” today

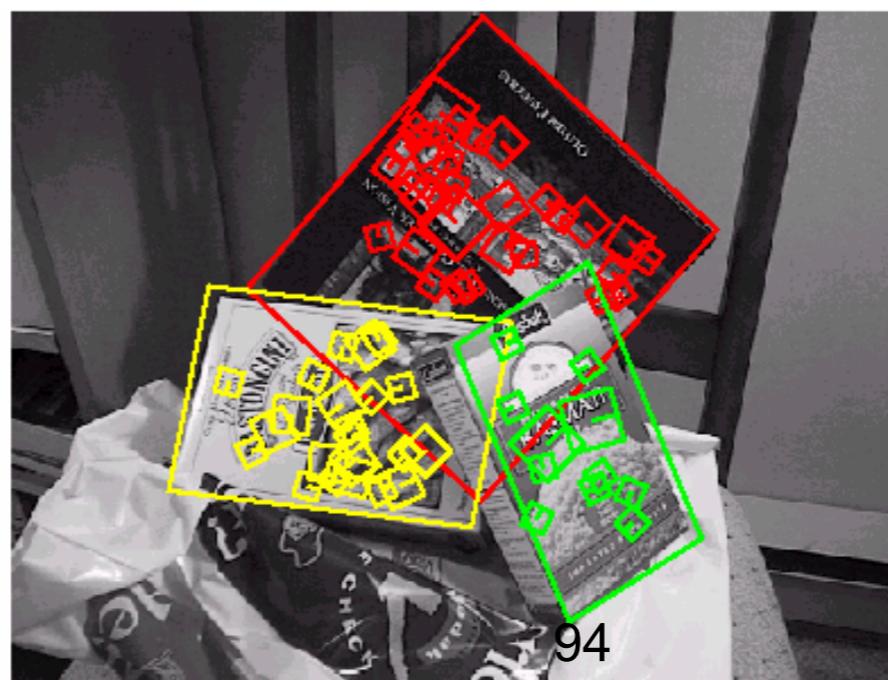
- Reading license plates, zip codes, checks
- Fingerprint recognition
- Face detection



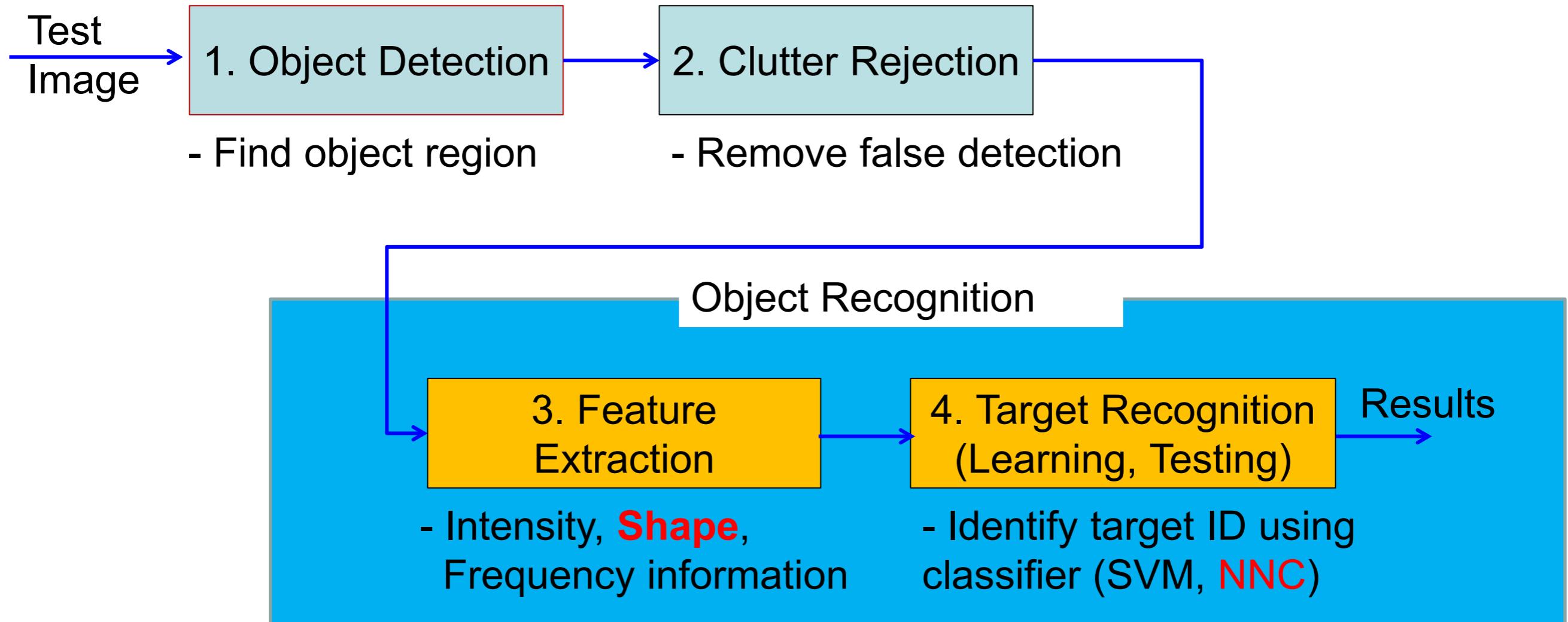
[Face priority AE] When a bright part of the face is too bright

What “works” today

- Reading license plates, zip codes, checks
- Fingerprint recognition
- Face detection
- Recognition of flat textured objects (CD covers, book covers, etc.)



Previous Object Detection/Recognition Flow



What Matters in Recognition?

- Learning Techniques
 - E.g. choice of classifier or inference method
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - **Deep learned features**
- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part