

# Computer Vision II - Recognition: *Image Categorization*

Michael Yang

# Roadmap (3 lectures)

- Object Detection
- Scene Understanding
- Image Categorization

# Roadmap (last lecture)

- Part Based Detector (cont. last last lecture)
  - Deformable Part Model
  - Poselets
- Scene Understanding Problem
- Context
- Spatial Layout
- 3D Scene Understanding

# Class-based recognition: Level of Detail

- **Image Categorization**

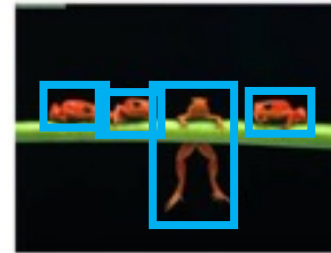
- One or more categories per image



Frog, branch

- **Object Class Detection**

- Also find bounding box



2D bounding box for each frog

- **Part-based Object Detection**

- Find parts of the object  
(and in this way the full object)



- **Semantic Segmentation (see last lecture)**  
(segmentation implies pixel-wise accuracy)

- Object-class segmentation

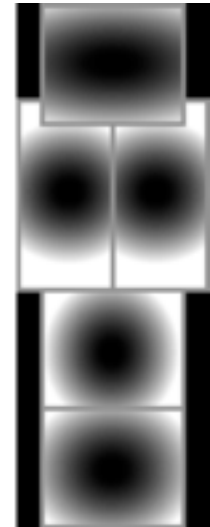
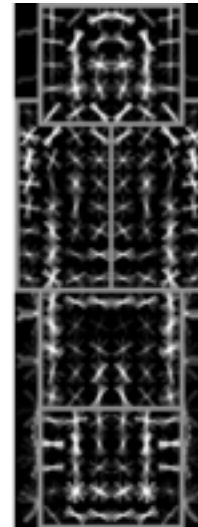
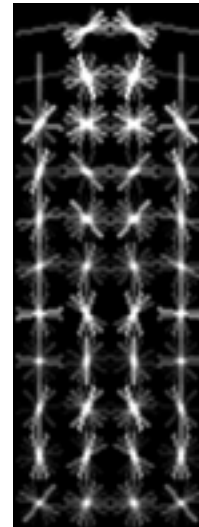
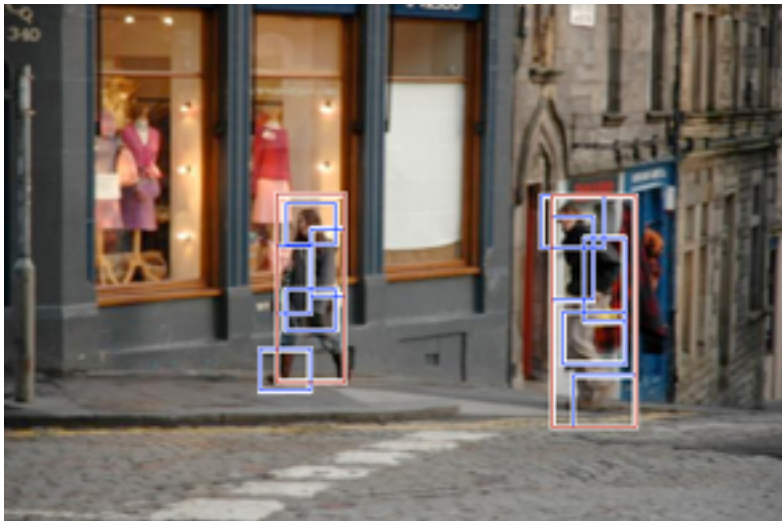




# Task: Generic object detection



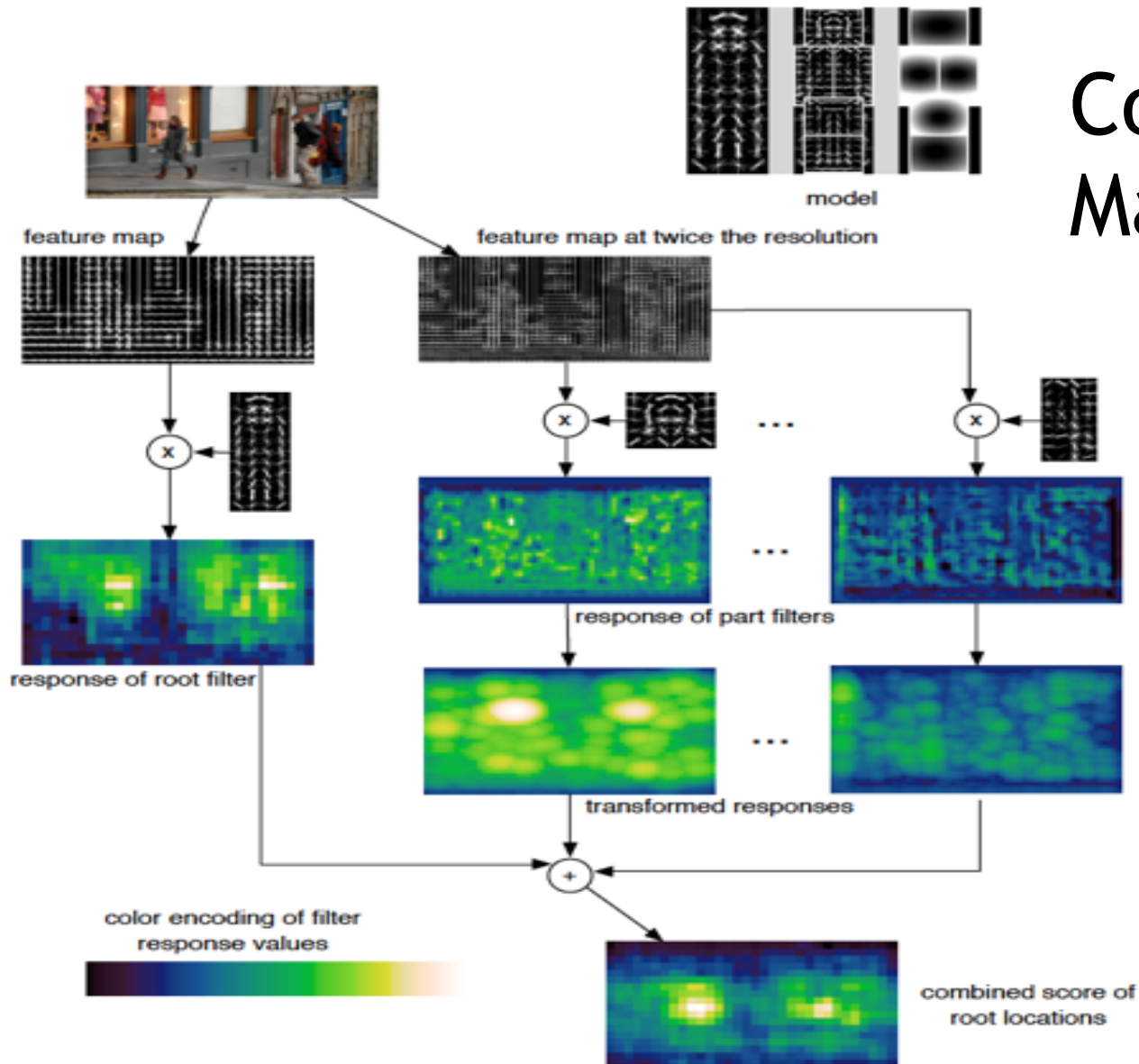
# DPM : Object Detection with Discriminatively Trained Part Based Models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part Based Models](#), PAMI 32(9), 2010

- Each category detector has mixture of deformable part models (components)
- Each component has global template + deformable parts
- Fully trained from bounding boxes alone (Latent SVM)

## Combine Many Parts





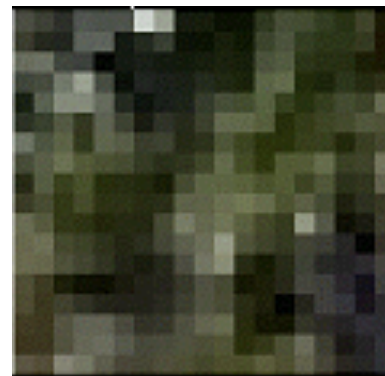
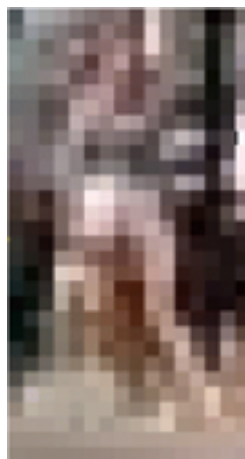
# Poselet



One poselet one classifier  
not a model for whole human body

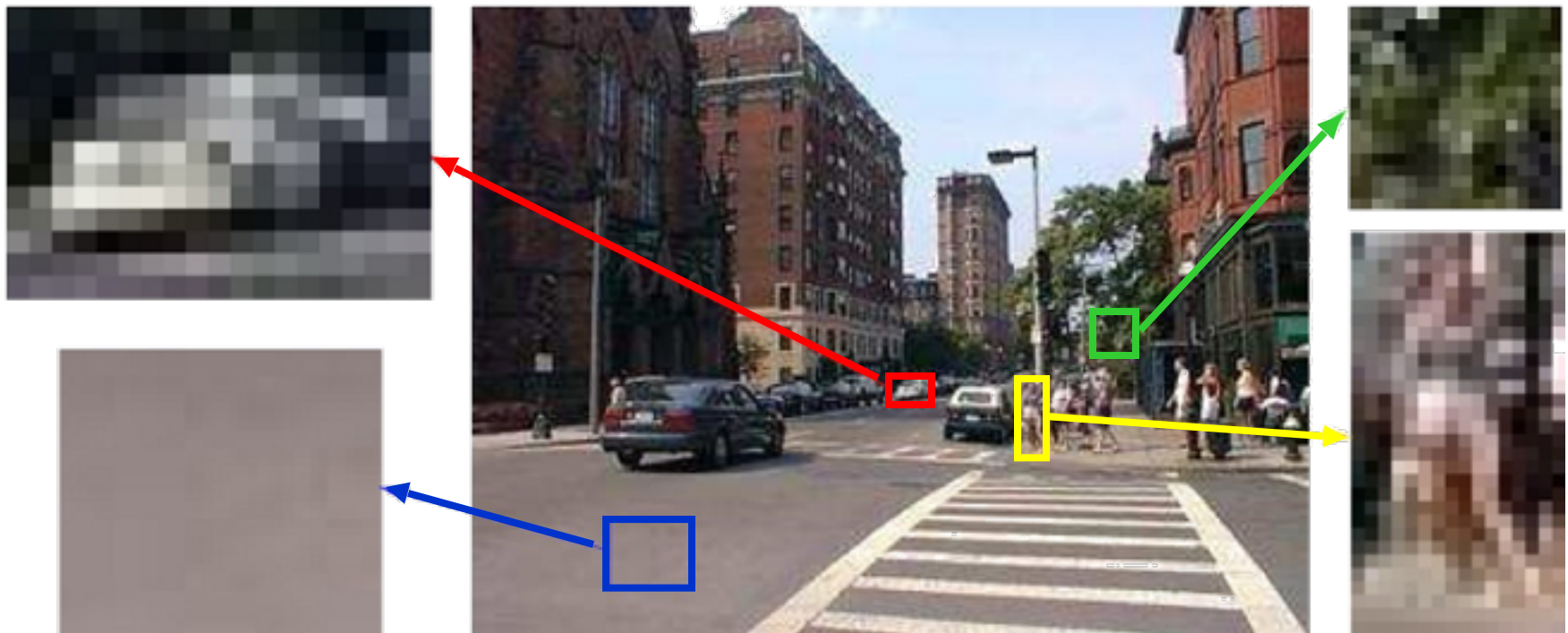
# Spatial layout is especially important

## 1. Context for recognition



# Spatial layout is especially important

## 1. Context for recognition





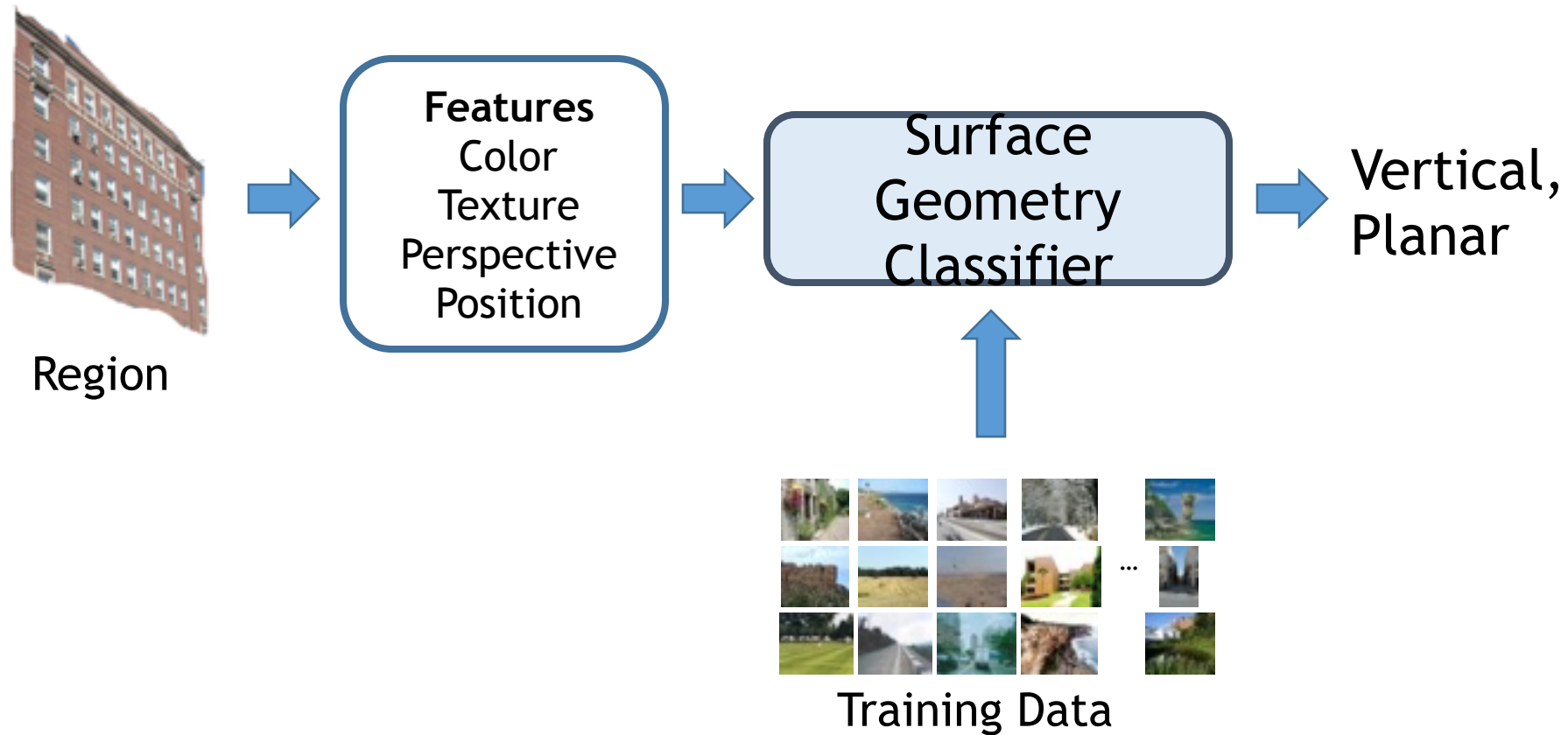
# Spatial layout is especially important

1. Context for recognition
2. Scene understanding





# Geometry estimation as recognition



# Surface Layout Algorithm

Input Image

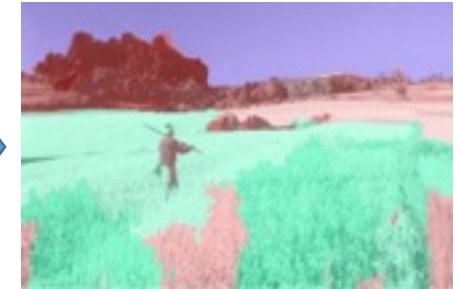


Segmentation

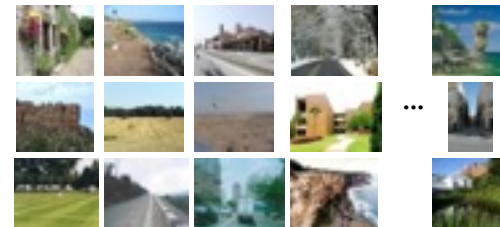


Features  
Perspective  
Color  
Texture  
Position

Surface Labels



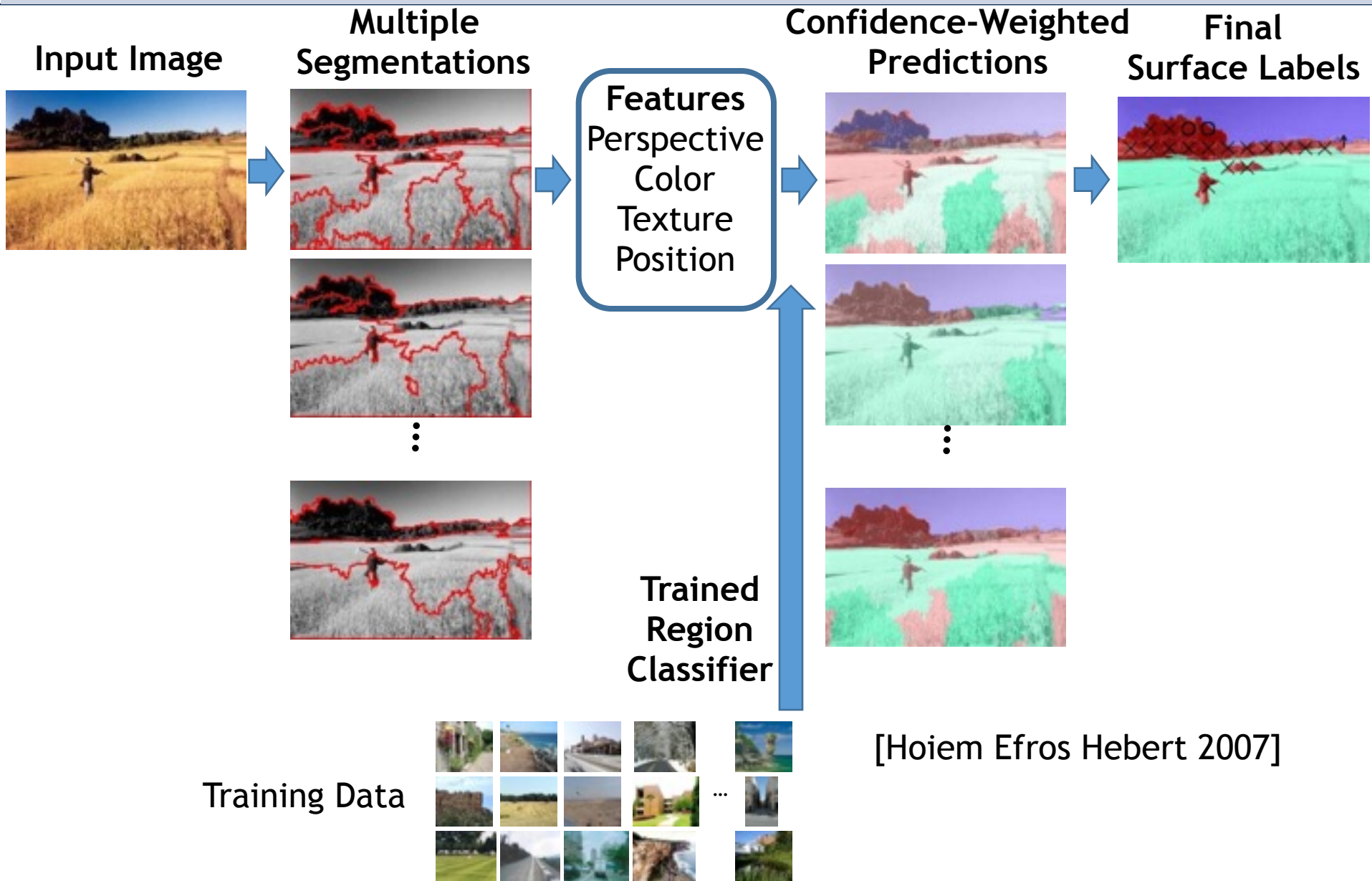
Trained  
Region  
Classifier



Training Data

[Hoiem Efros Hebert 2007]

# Surface Layout Algorithm



# Roadmap (this lecture)

- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching

# Class-based recognition: Level of Detail

- **Image Categorization**

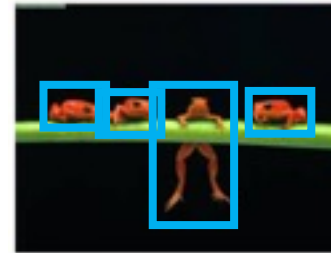
- One or more categories per image



Frog (branch)

- **Object Class Detection**

- Also find bounding box



2D bounding box for each frog

- **Part-based Object Detection**

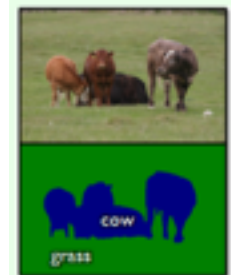
- Find parts of the object  
(and in this way the full object)



- **Semantic Segmentation**

(segmentation implies pixel-wise accuracy)

- Object-class segmentation

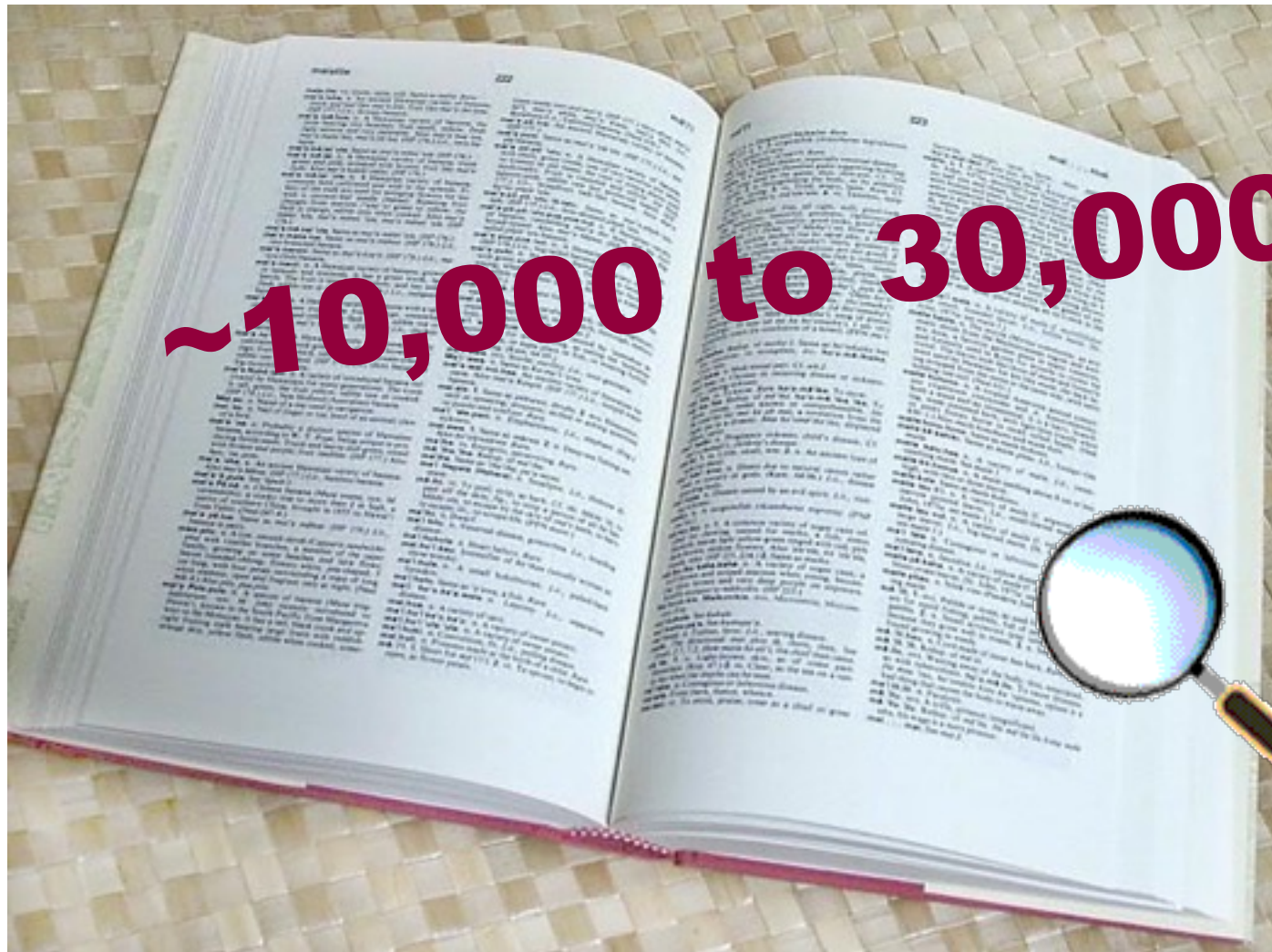


Why?

Application



# How many visual object categories are there?



Biederman 1987



~10,000 to 30,000





# OBJECTS

ANIMALS

PLANTS

INANIMATE

.....

VERTEBRATE

NATURAL

MAN-MADE

MAMMALS

BIRDS

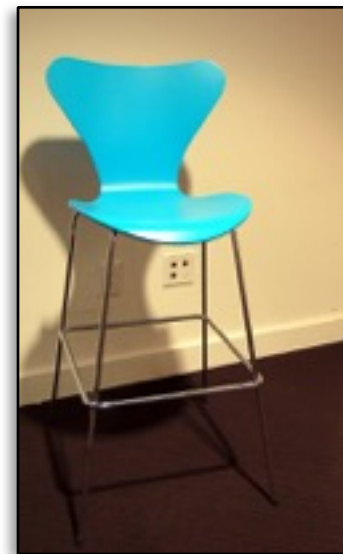
TAPIR

BOAR

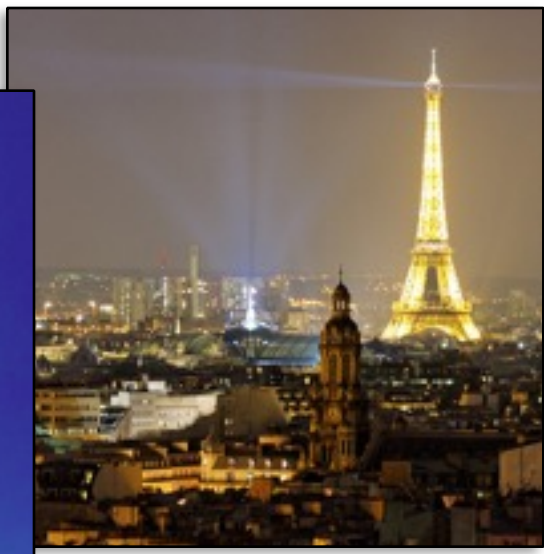
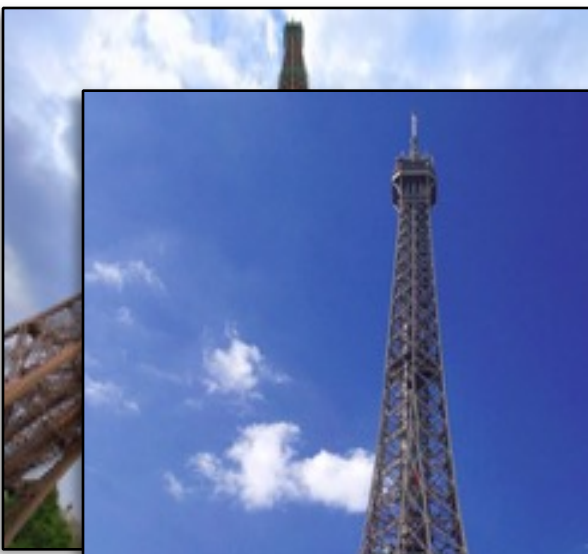
GROUSE

CAMERA





Variation within an object class



Viewpoint/Scales/Illumination Variability

Images from Flickr



# Recognition: A machine learning approach



Slides adapted from Fei-Fei Li, Rob Fergus, Antonio Torralba, Kristen Grauman, and Derek Hoiem

# The machine learning framework

- Apply a prediction function to a feature representation of the image to get the desired output:

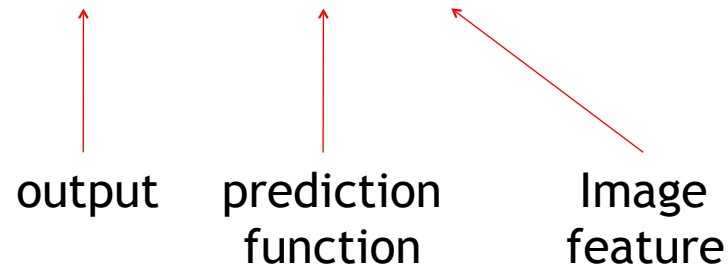
$f(\text{apple image}) = \text{“apple”}$

$f(\text{tomato image}) = \text{“tomato”}$

$f(\text{cow image}) = \text{“cow”}$

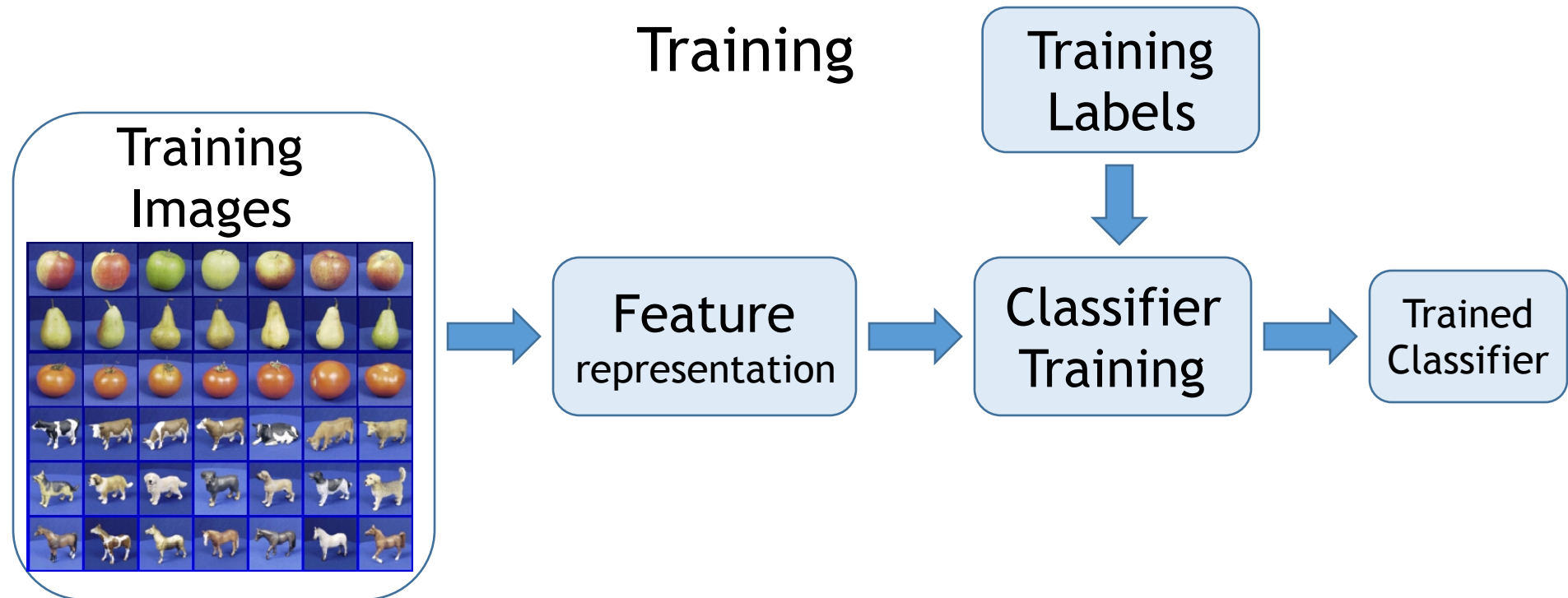
# The machine learning framework

$$y = f(x)$$

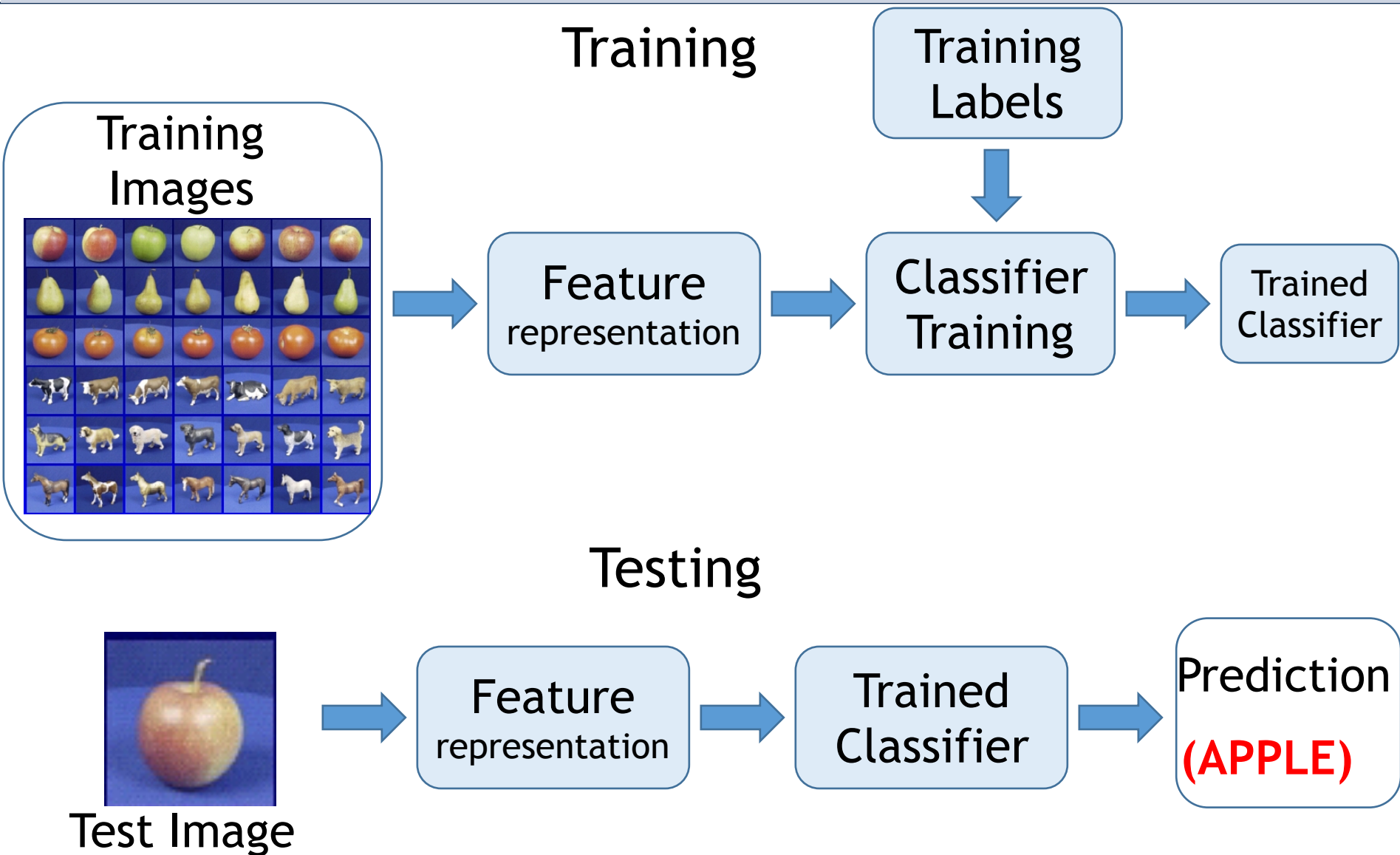


- **Training:** given a *training set* of labeled examples  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , estimate the prediction function  $f$  by minimizing the prediction error on the training set
- **Testing:** apply  $f$  to a never before seen *test example*  $\mathbf{x}$  and output the predicted value  $y = f(\mathbf{x})$

# Image Categorization-Steps



# Image Categorization-Steps

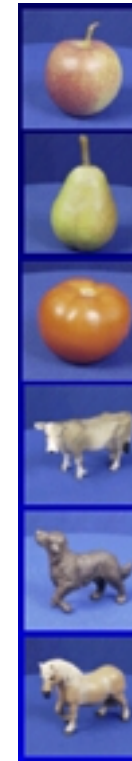




# Generalization



Training set (labels known)



Test set (labels unknown)

- How well does a learned model *generalize* from the data it was trained on to a new test set?

## Generalization depends on:

- Invariance properties of the feature representation
  - There is a tradeoff between invariance and discriminability
- Training data
  - Some intra-class variations must be adequately represented in the training data (hard to model analytically)
- Statistical model
  - Some models are more powerful than others and able to generalize better.

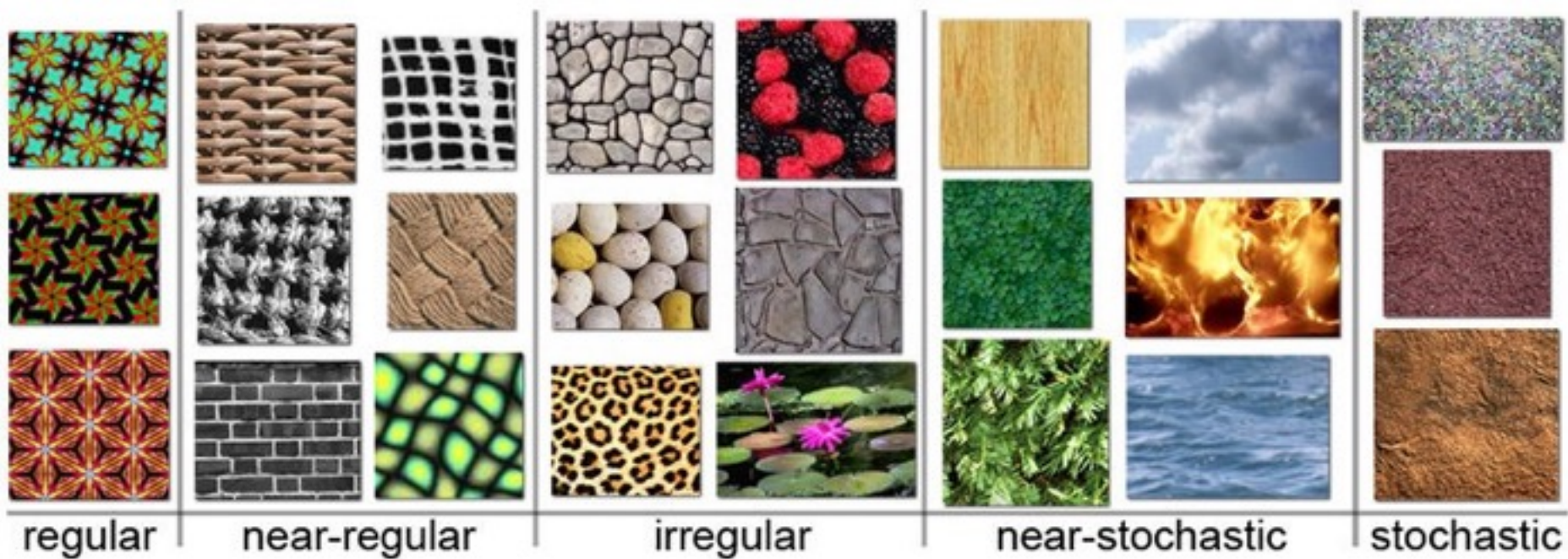
# Roadmap (this lecture)

- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching

# Image Categorization - Bag of Words Approach



# Origin 1: Texture recognition

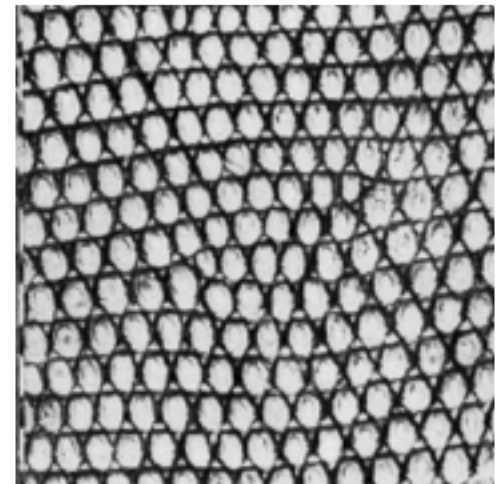
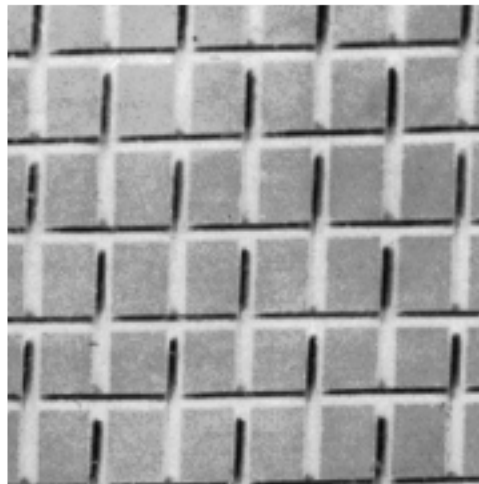
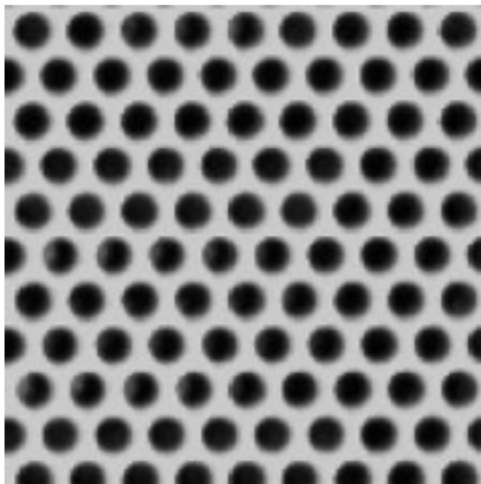


Example textures (from Wikipedia)



# Origin 1: Texture recognition

- Texture is characterized by the repetition of basic elements or *textons*
- For stochastic textures, it is the identity of the textons, not their spatial arrangement, that matters



Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

2007-01-23: State of the Union Address

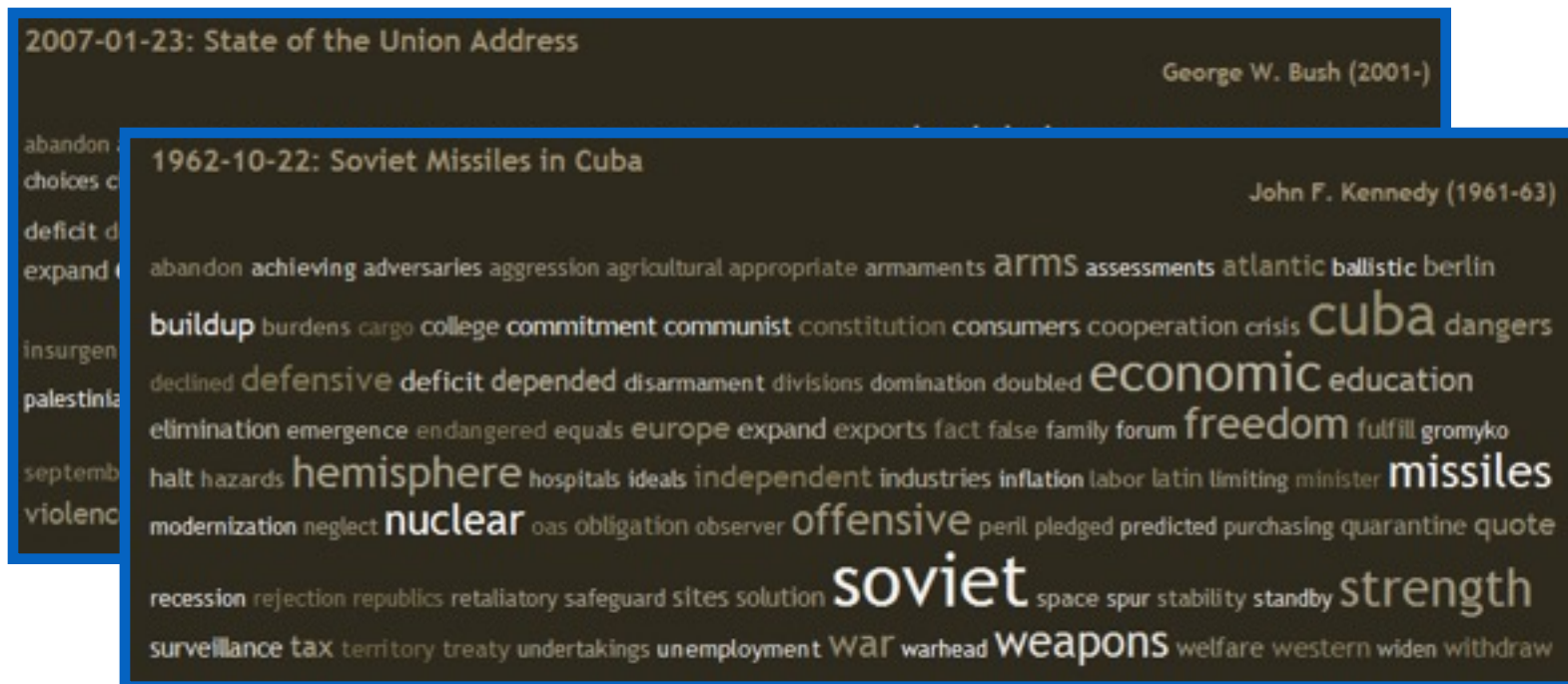
George W. Bush (2001-)

abandon accountable affordable afghanistan africa aided ally anbar armed army **baghdad** bless **challenges** chamber chaos  
choices civilians coalition commanders **commitment** confident confront congressman constitution corps debates deduction  
deficit deliver **democratic** deploy dikembe diplomacy disruptions earmarks **economy** einstein elections eliminates  
expand **extremists** failing faithful families **freedom** fuel funding god haven ideology immigration impose  
insurgents iran **iraq** islam julie lebanon love madam marine math **medicare** moderation neighborhoods **nuclear** offensive  
palestinian payroll province pursuing **qaeda** radical **regimes** resolve retreat rleman sacrifices science sectarian senate  
september **shia** stays strength students succeed **sunni** tax territories **terrorists** threats uphold victory  
violence violent WAF washington weapons wesley



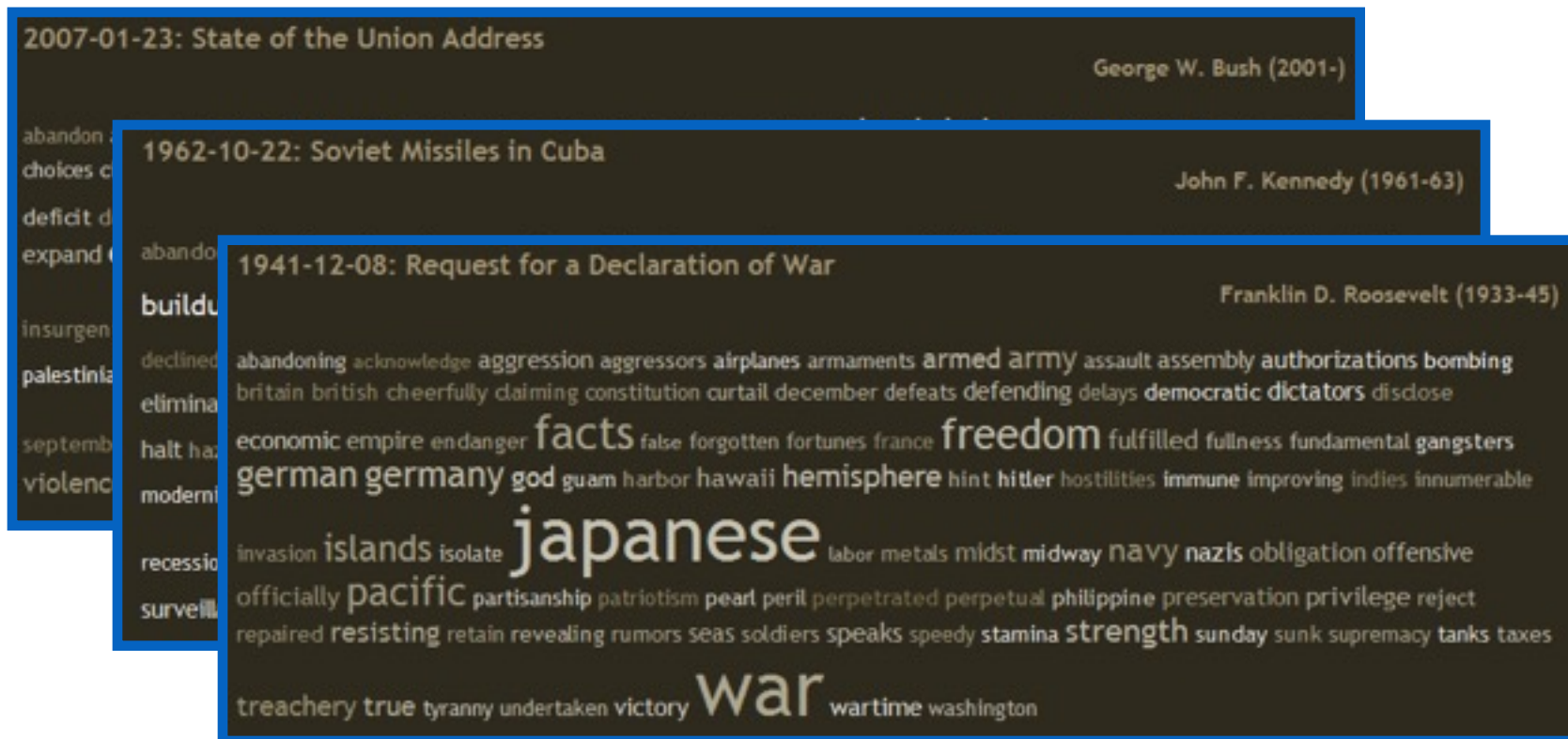
# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

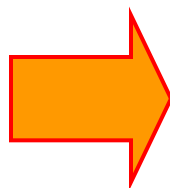


# Origin 2: Bag-of-words models

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



# Bags of words for object recognition



face, flowers, building

- Works pretty well for image-level classification and for recognizing object *instances*

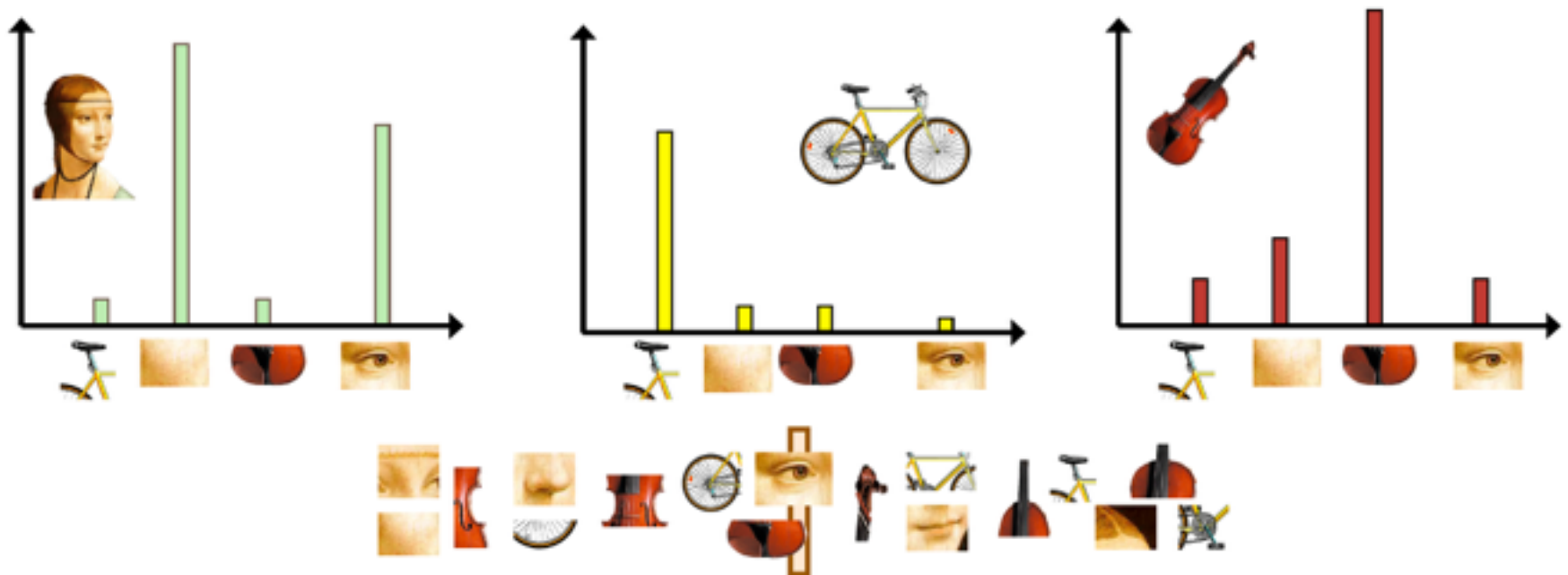
# Bags of words for object recognition



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

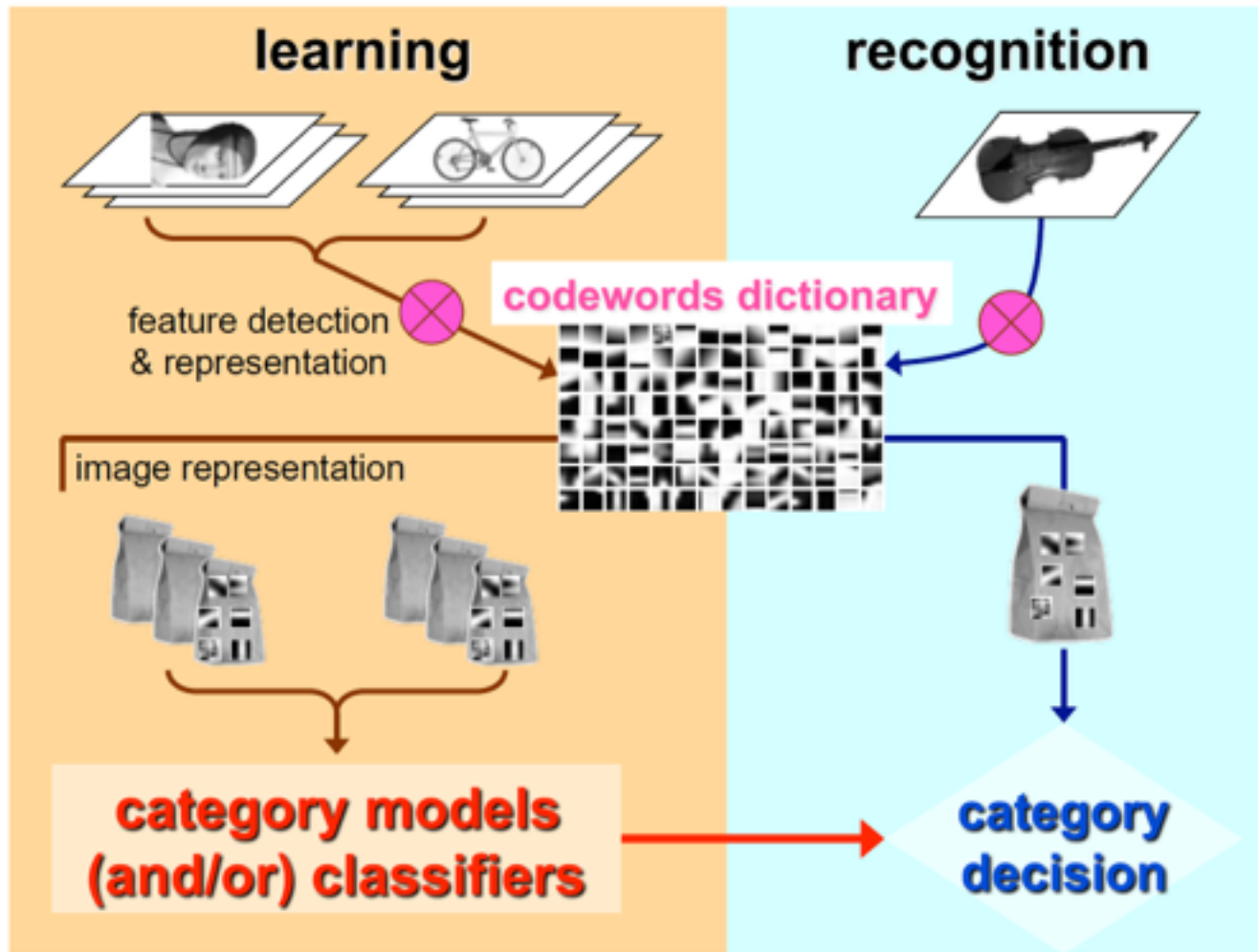
# Bag of Words

- Independent features
- Histogram representation

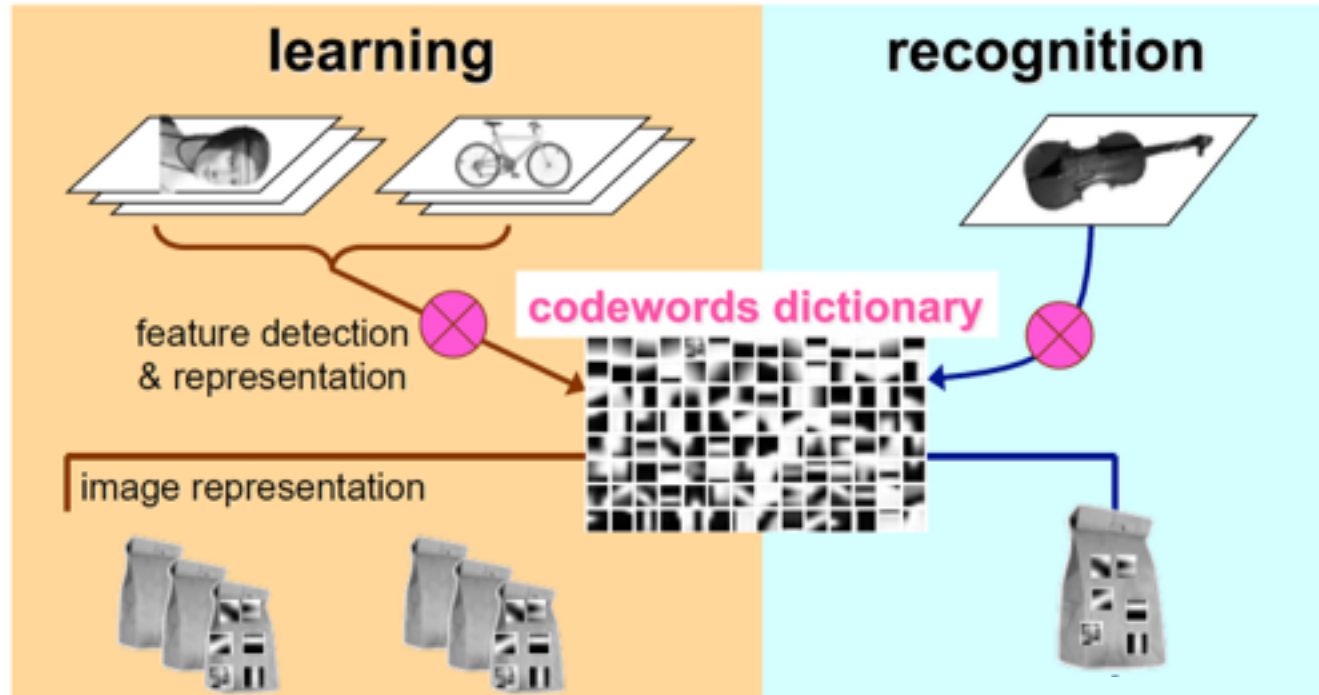




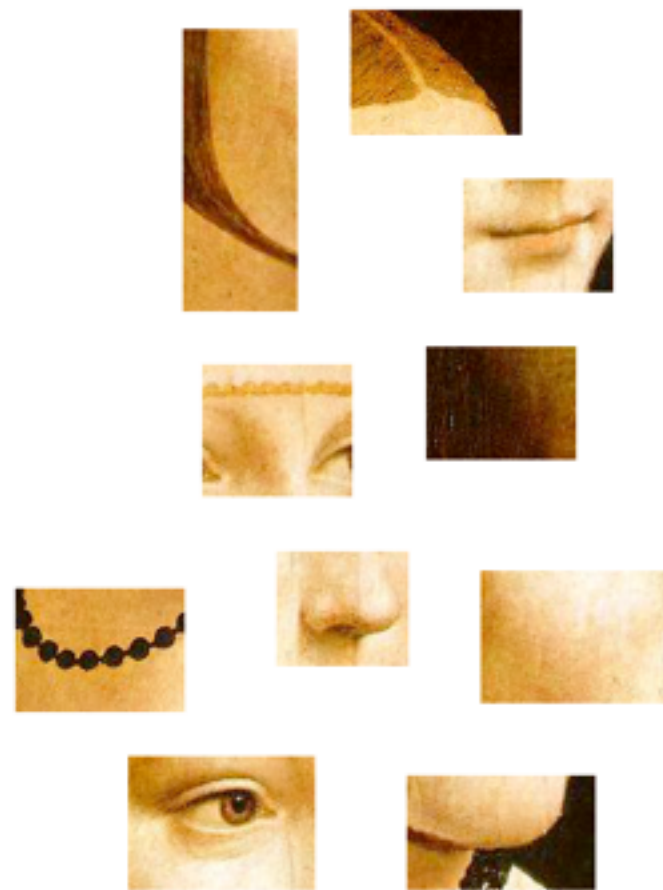
# Bag of Words - Overview



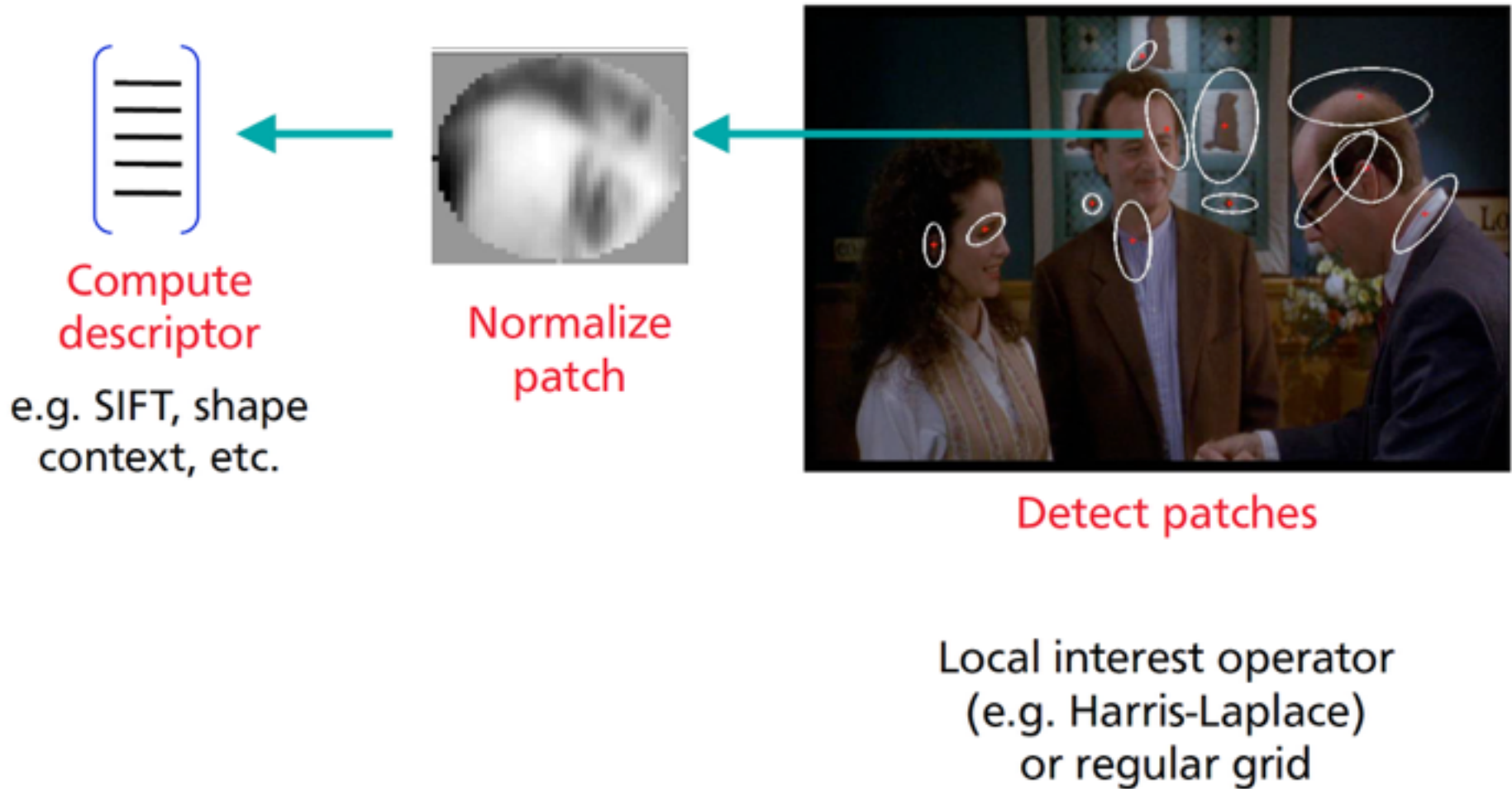
# Object Representation



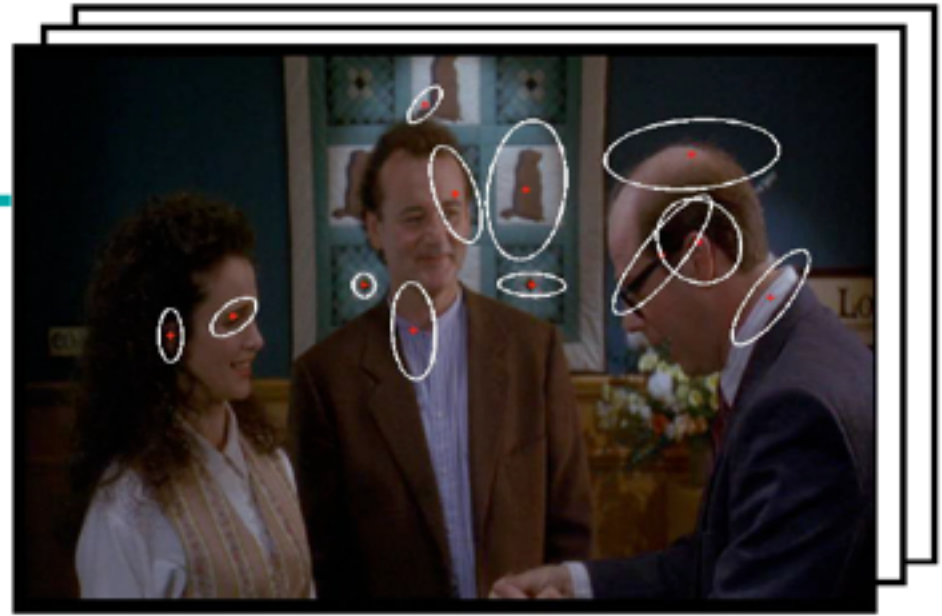
# Feature Detection and Representation



# Feature Detection and Representation



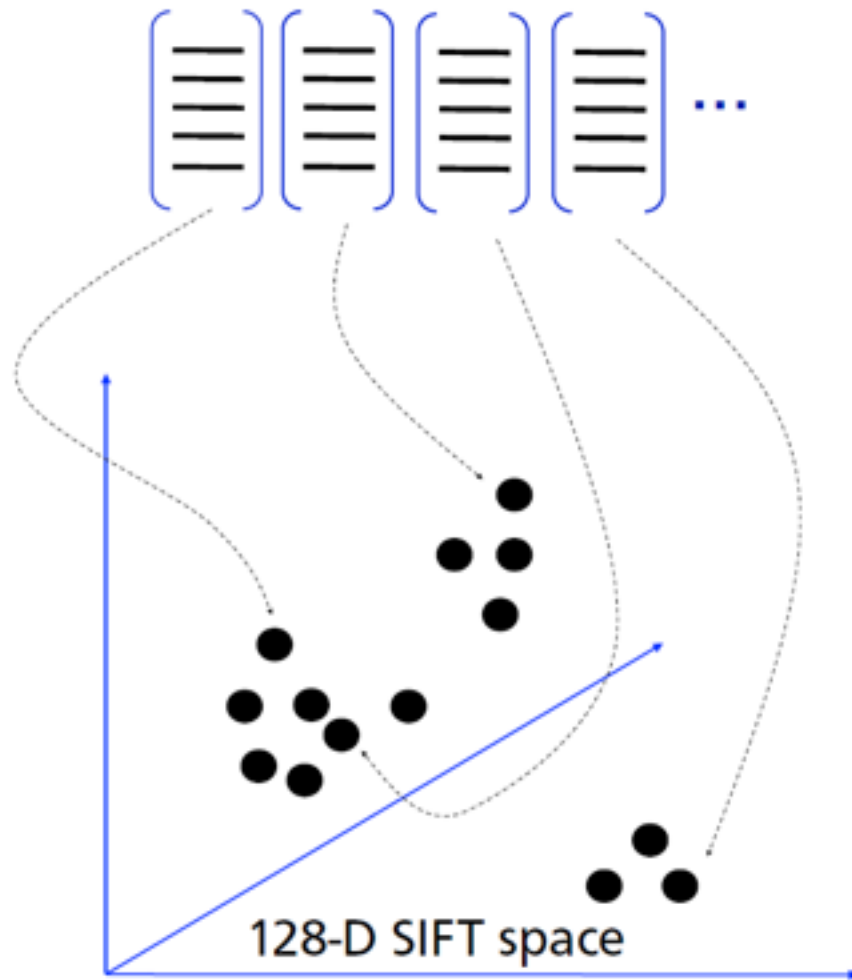
# Feature Detection and Representation



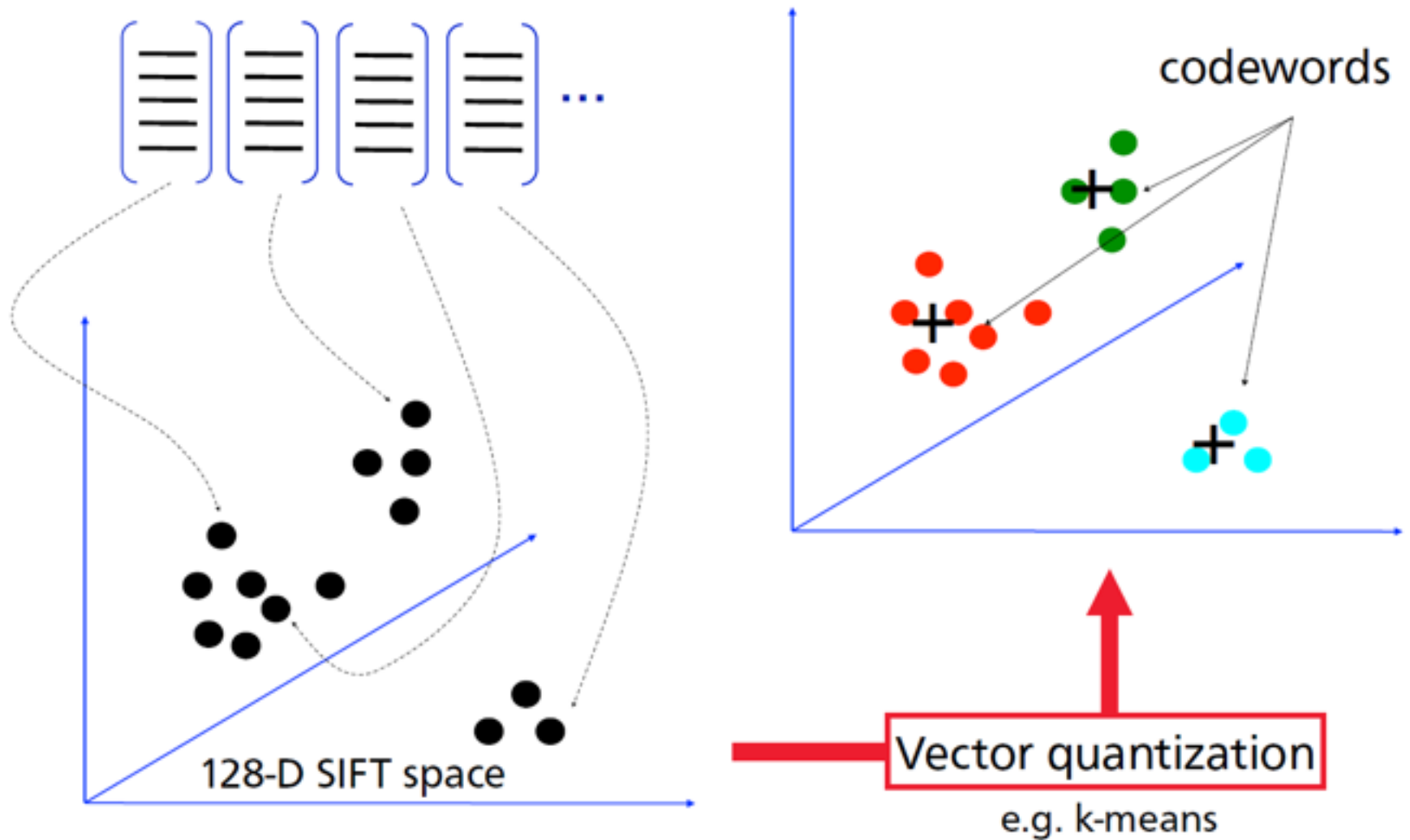
Take all training images



# Feature Detection and Representation

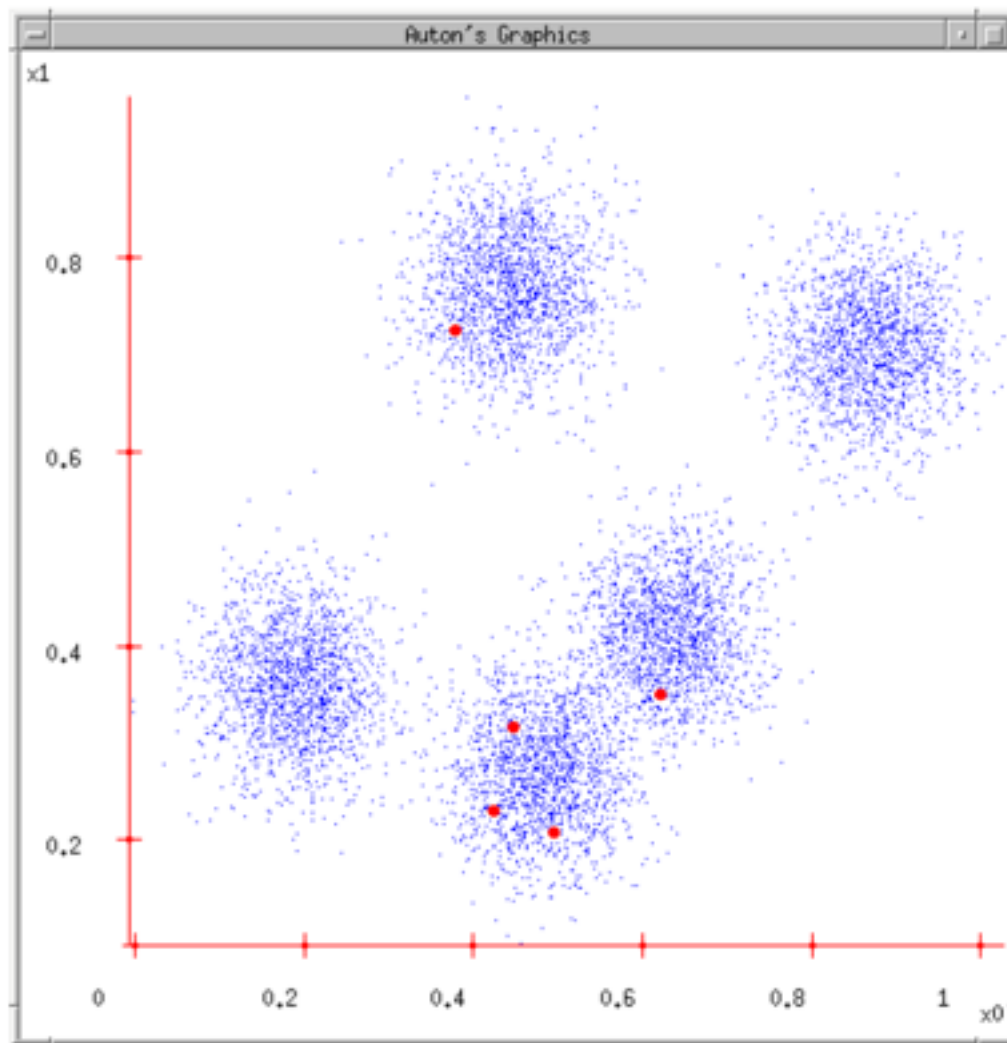


# Codeword dictionary formation



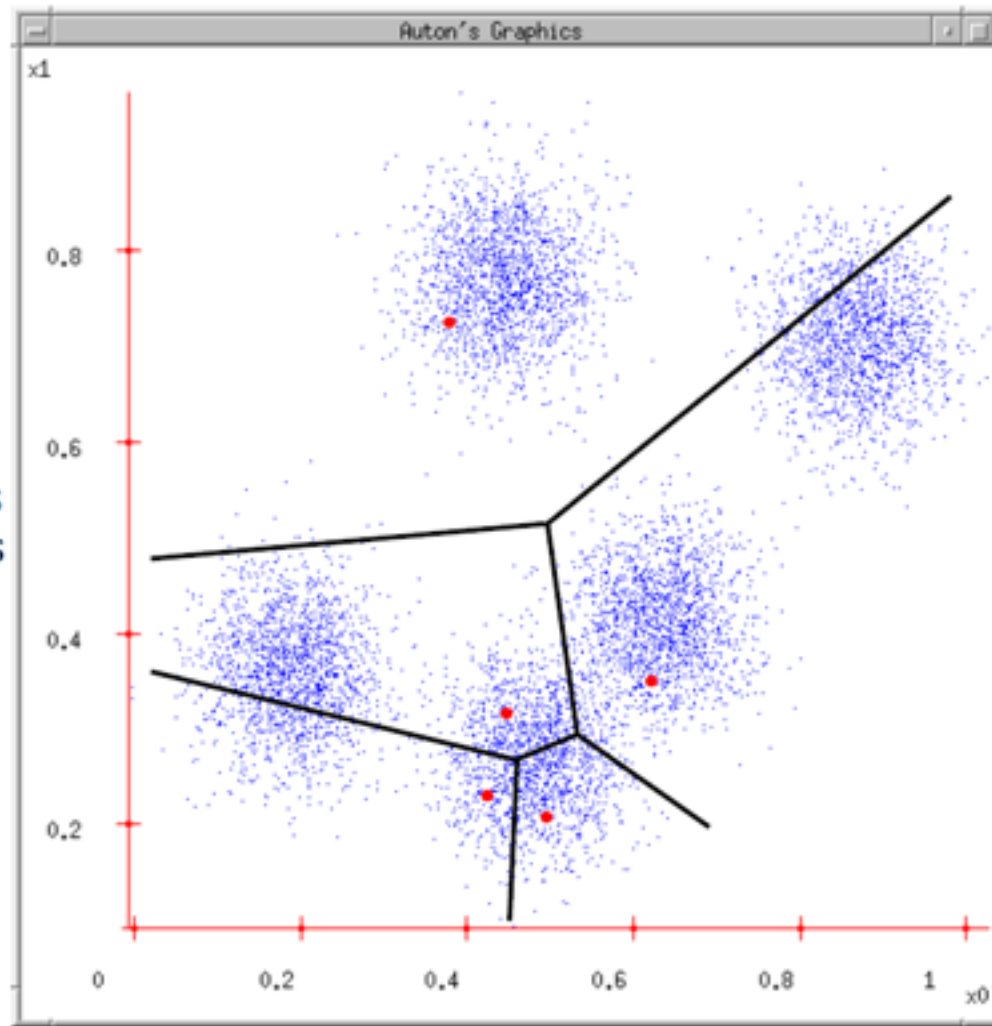
## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



## K-means

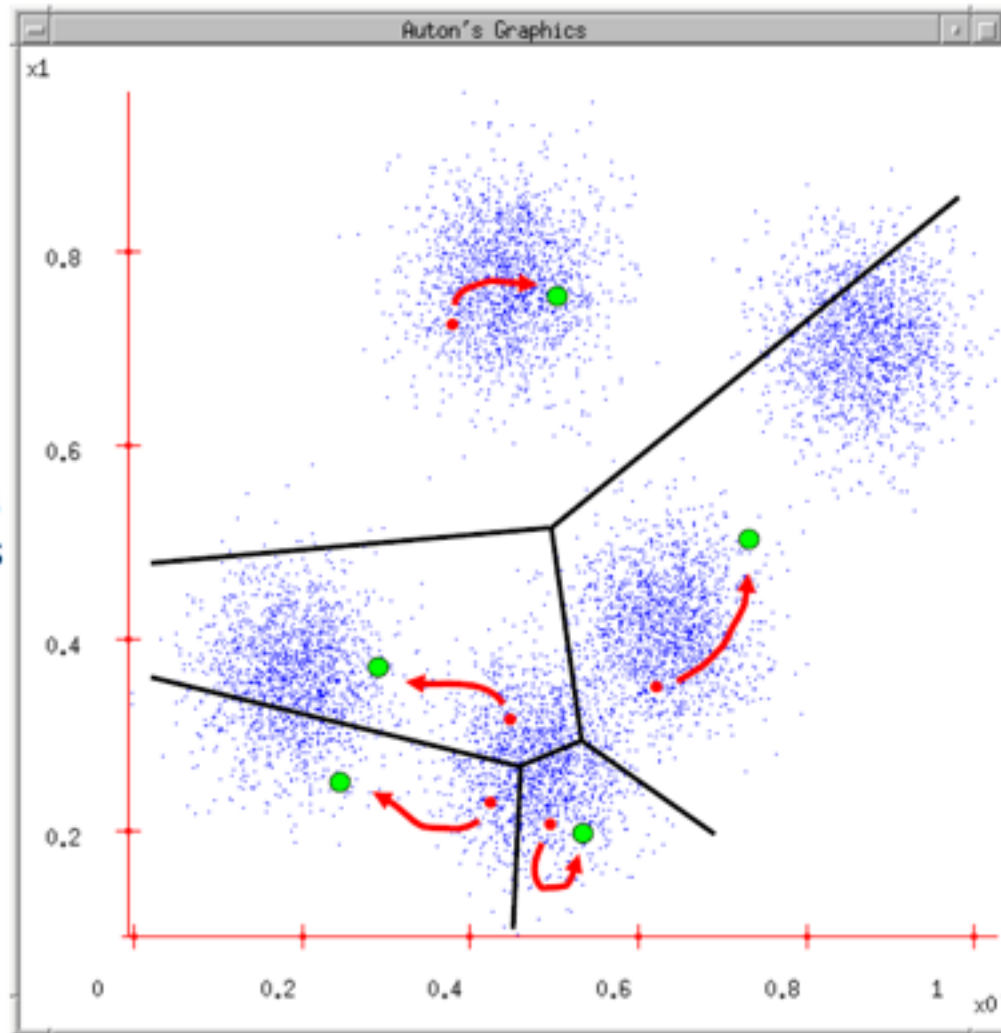
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



# Reminder: K-means

## K-means

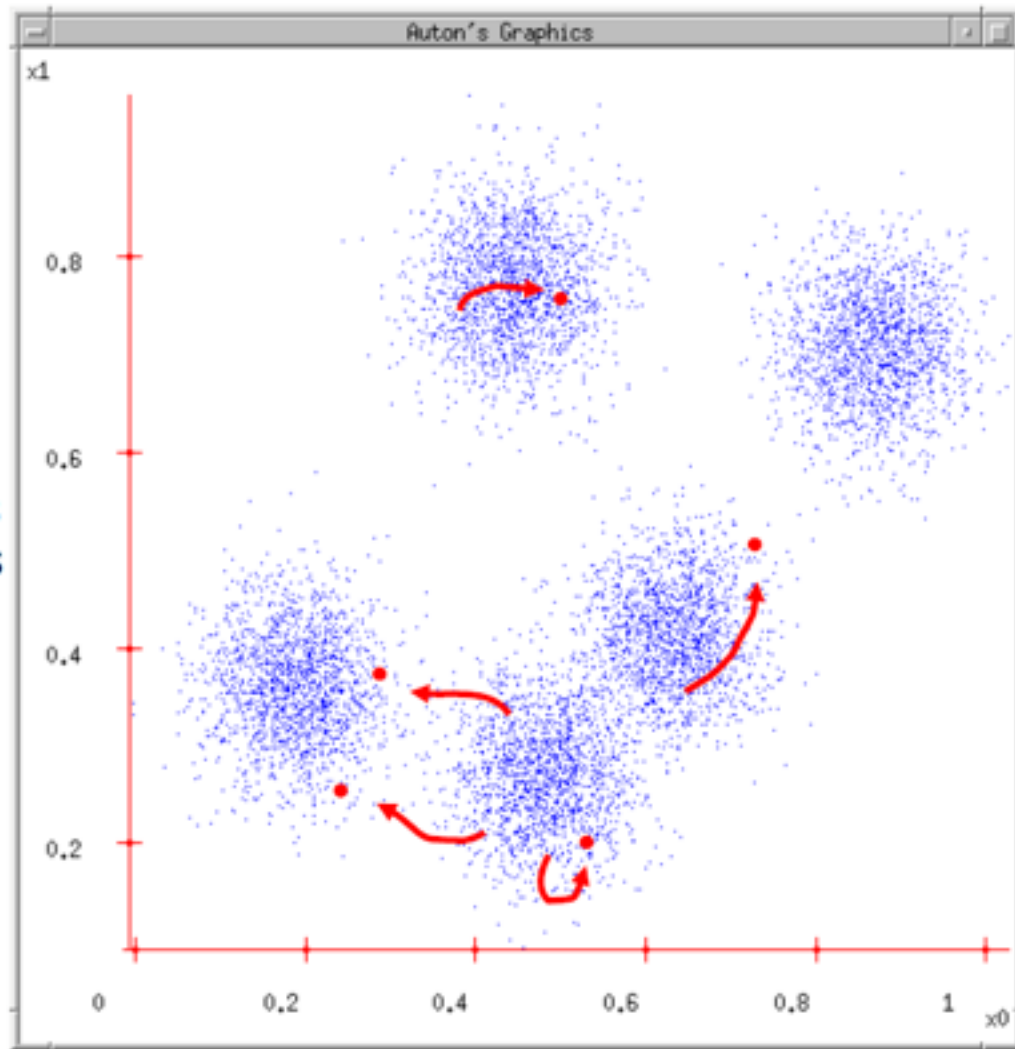
1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



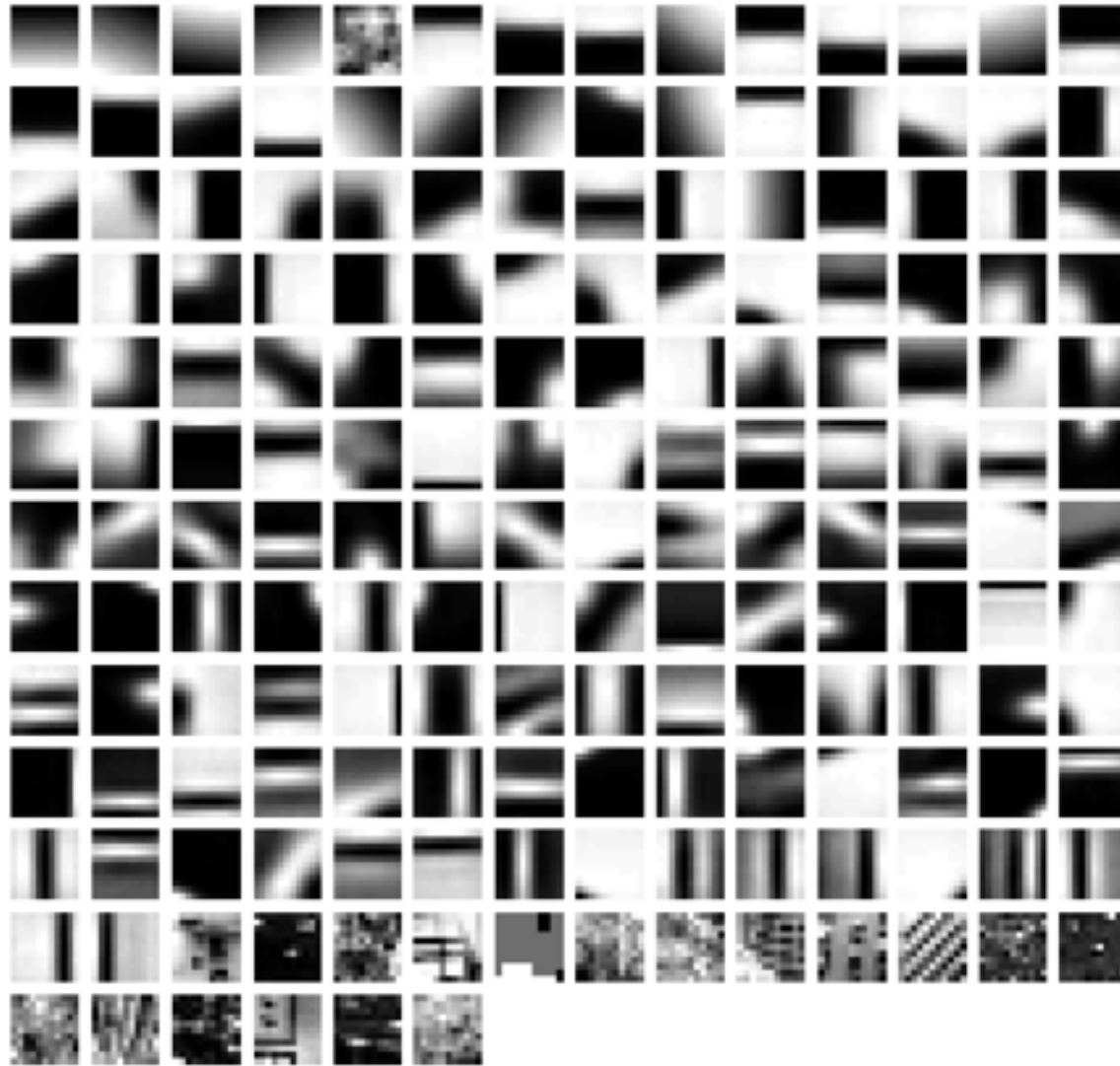


## K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!  
(Repeat means go to step 3)

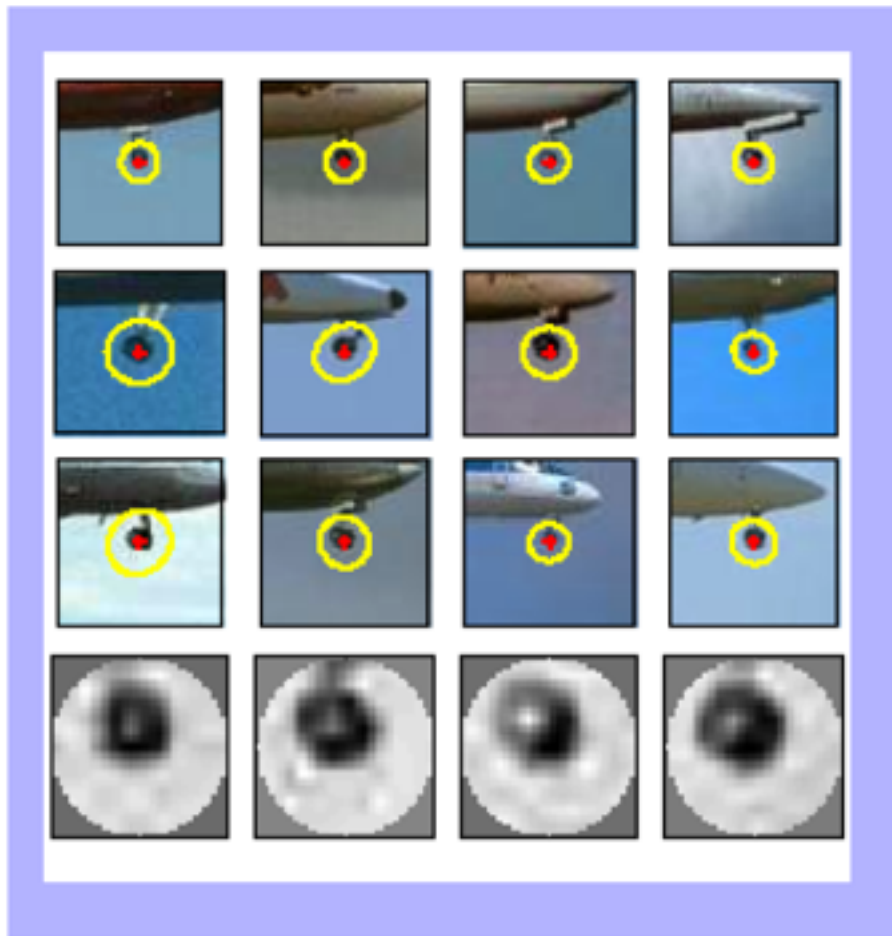


# Codeword dictionary visualization

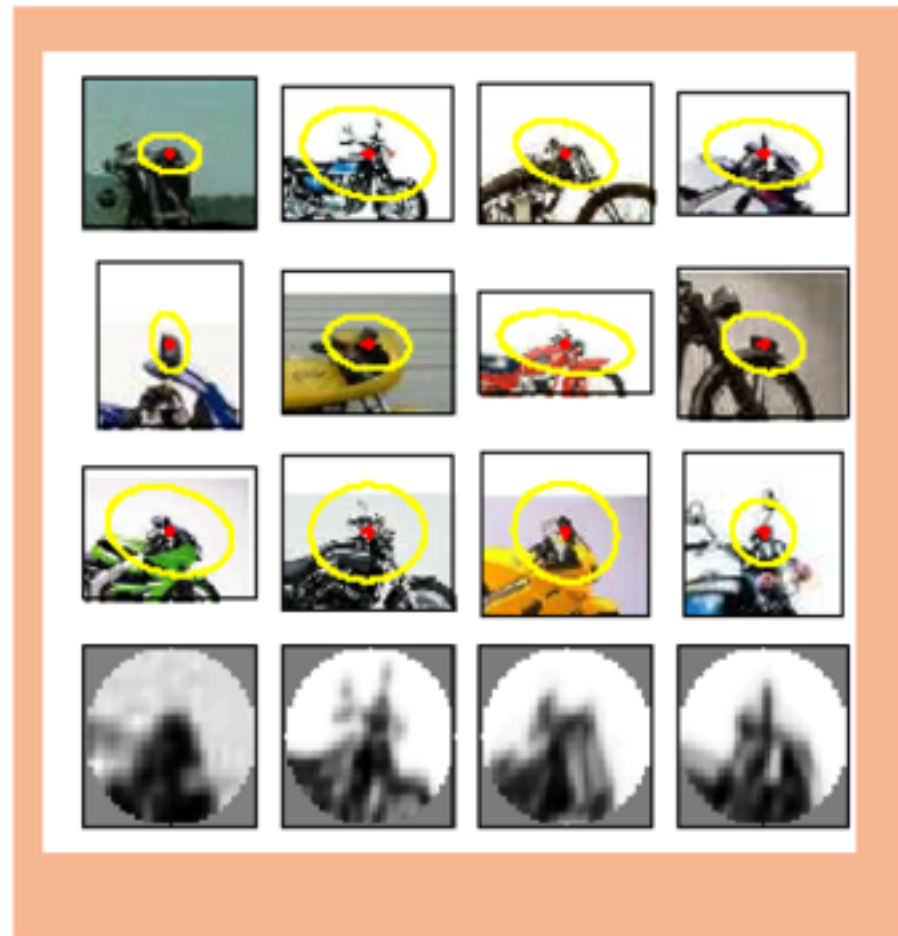


$K = 174$  (averaged patches for each cluster) [from Fei Fei Li]

# Image Patch examples of Codewords



Examples which are assigned to same codeword

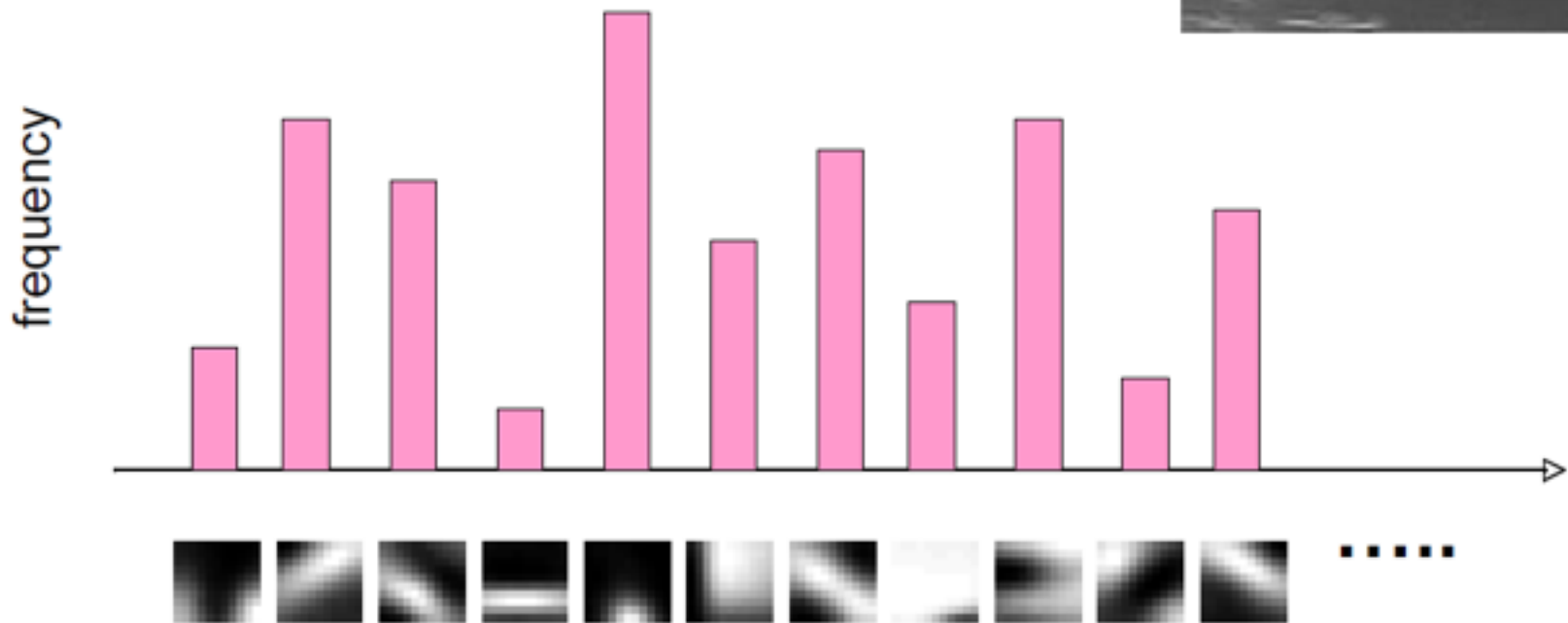


Examples which are assigned to same codeword

[from Josef Sivic]

# Bag of Words - Image Representation

- Histogram of features assigned to each cluster



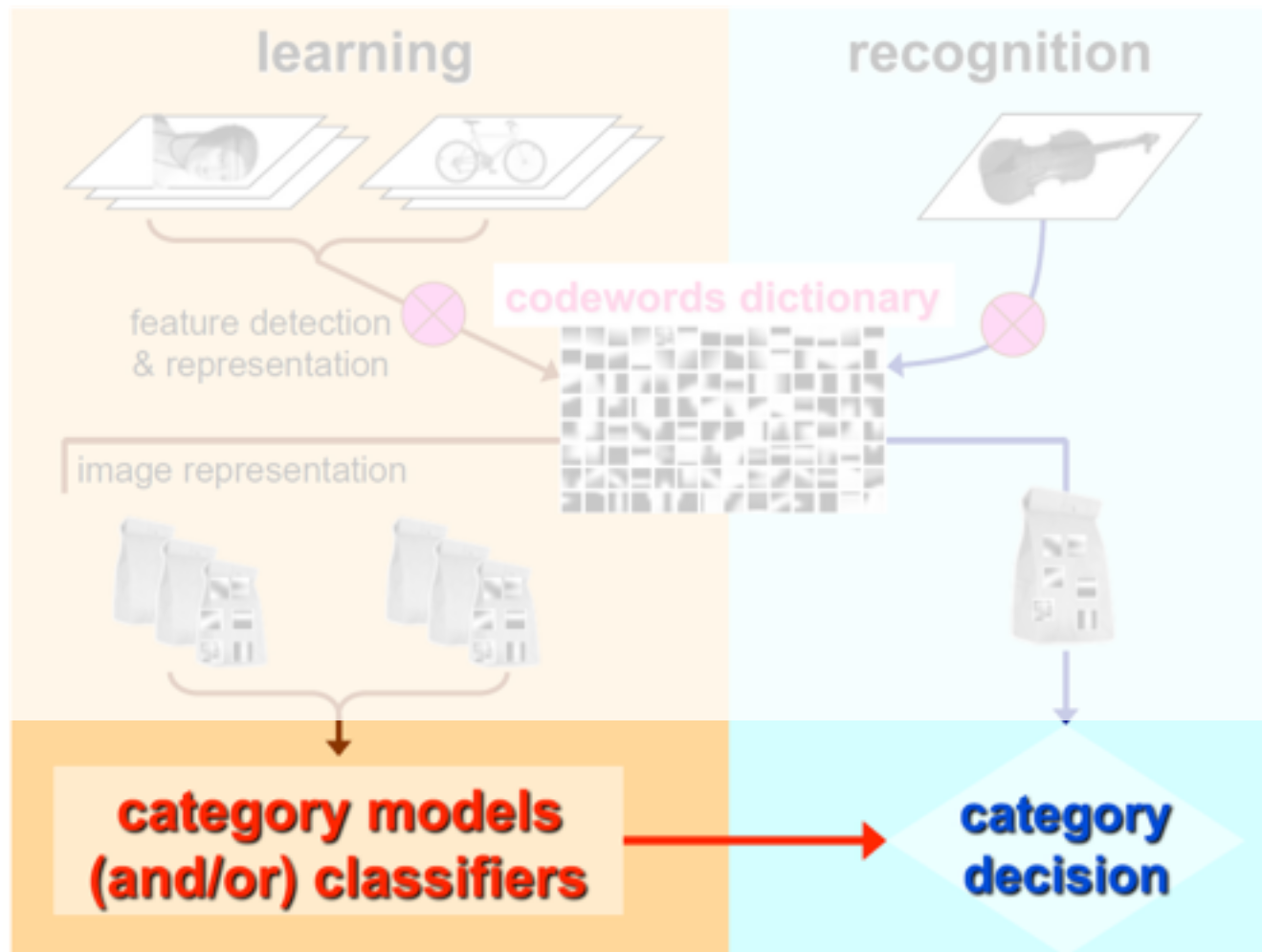
$K = 174$

# Roadmap (this lecture)

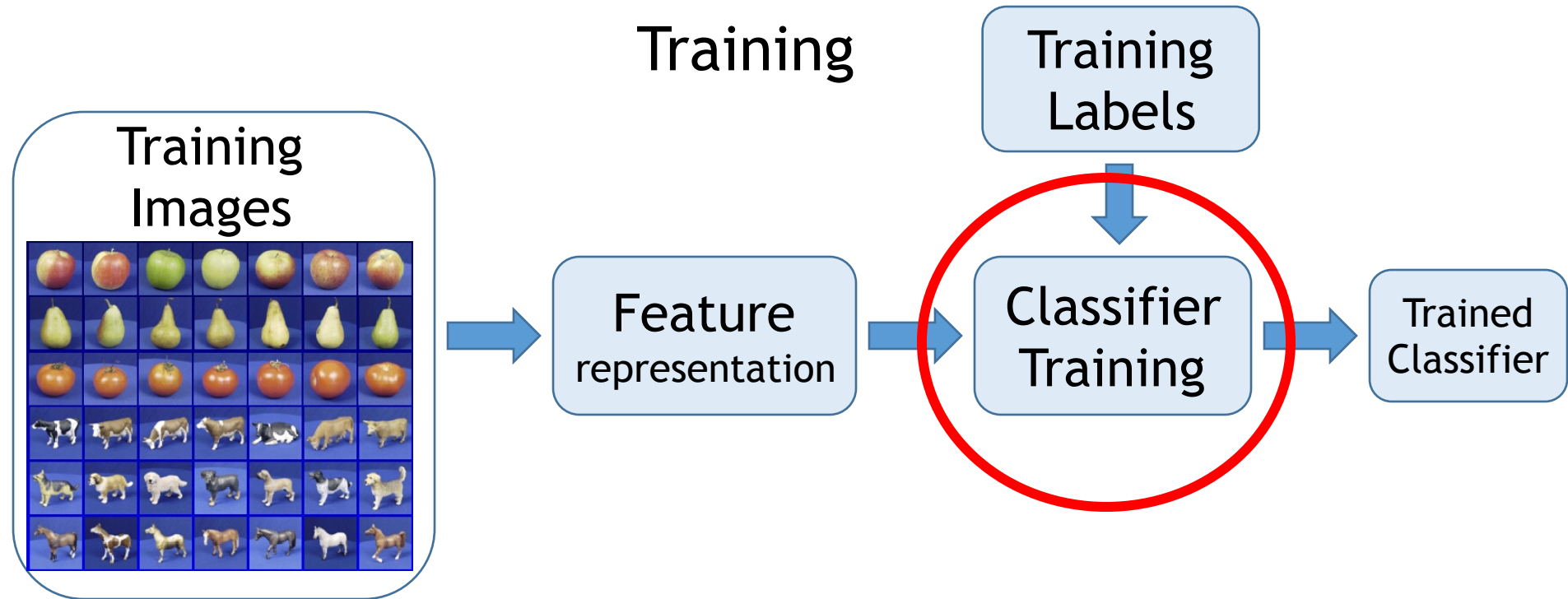
- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching



# Bag of Words - Overview

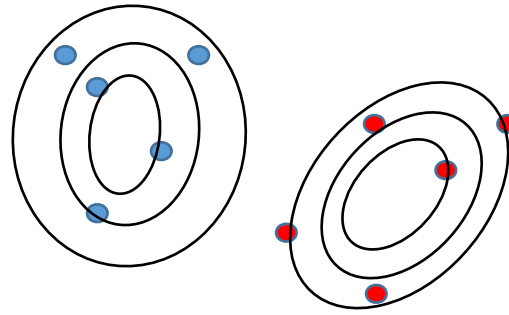


# Classifiers

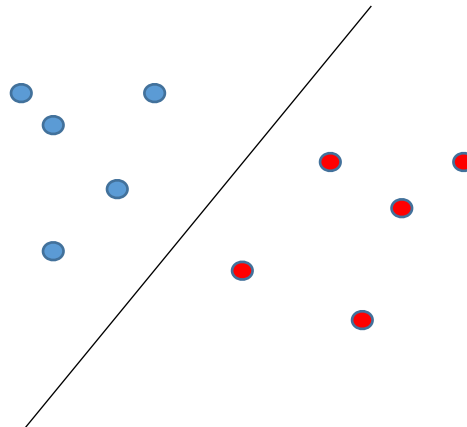


# Two approaches

Generative approach:  
*models distributions*



Discriminative function:  
*models decision function*



## Generative

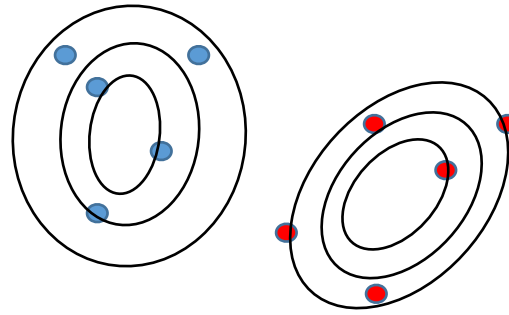
- Training
  - Maximize joint likelihood of data and labels
  - Assume (or learn) probability distribution and dependency structure
  - Can impose priors
- Testing
  - $P(y=1, x) / P(y=0, x) > t?$
- Examples
  - Foreground/background GMM
  - Naïve Bayes classifier
  - Bayesian network

## Discriminative

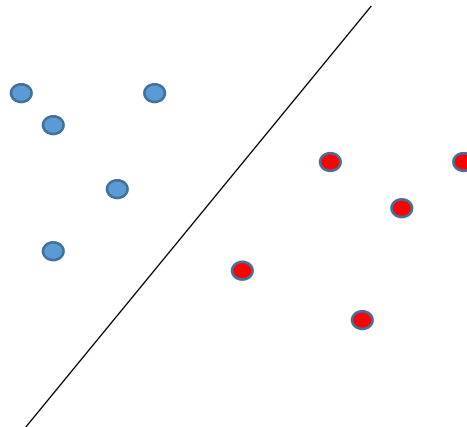
- Training
  - Learn to directly predict the labels from the data
  - Assume form of boundary
  - Margin maximization or parameter regularization
- Testing
  - $f(x) > t$  ; e.g.,  $w^T x > t$
- Examples
  - Logistic regression
  - SVM
  - Boosted decision trees

# Two approaches

Generative approach:  
*models distributions*



Discriminative function:  
*models decision function*





# Discriminative functions

## ■ Linear discriminant function:

- Linear hyperplane:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

- trained on samples of both classes  $C_1$  (■) and  $C_2$  (●)

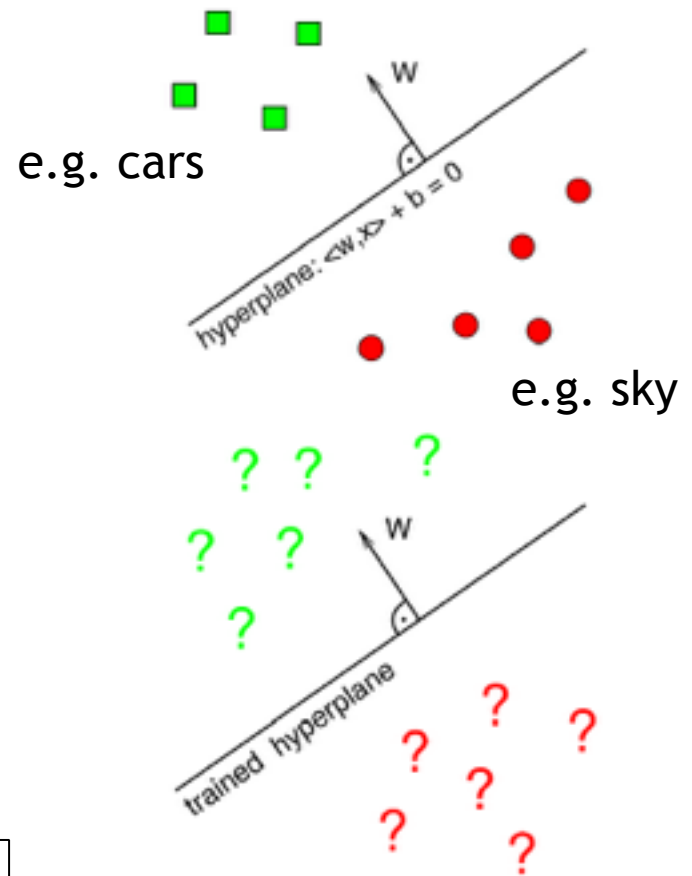
## ■ Classification:

- decide class  $C_1$  (?) when  $y(\mathbf{x}) > 0$
- decide class  $C_2$  (?) when  $y(\mathbf{x}) < 0$

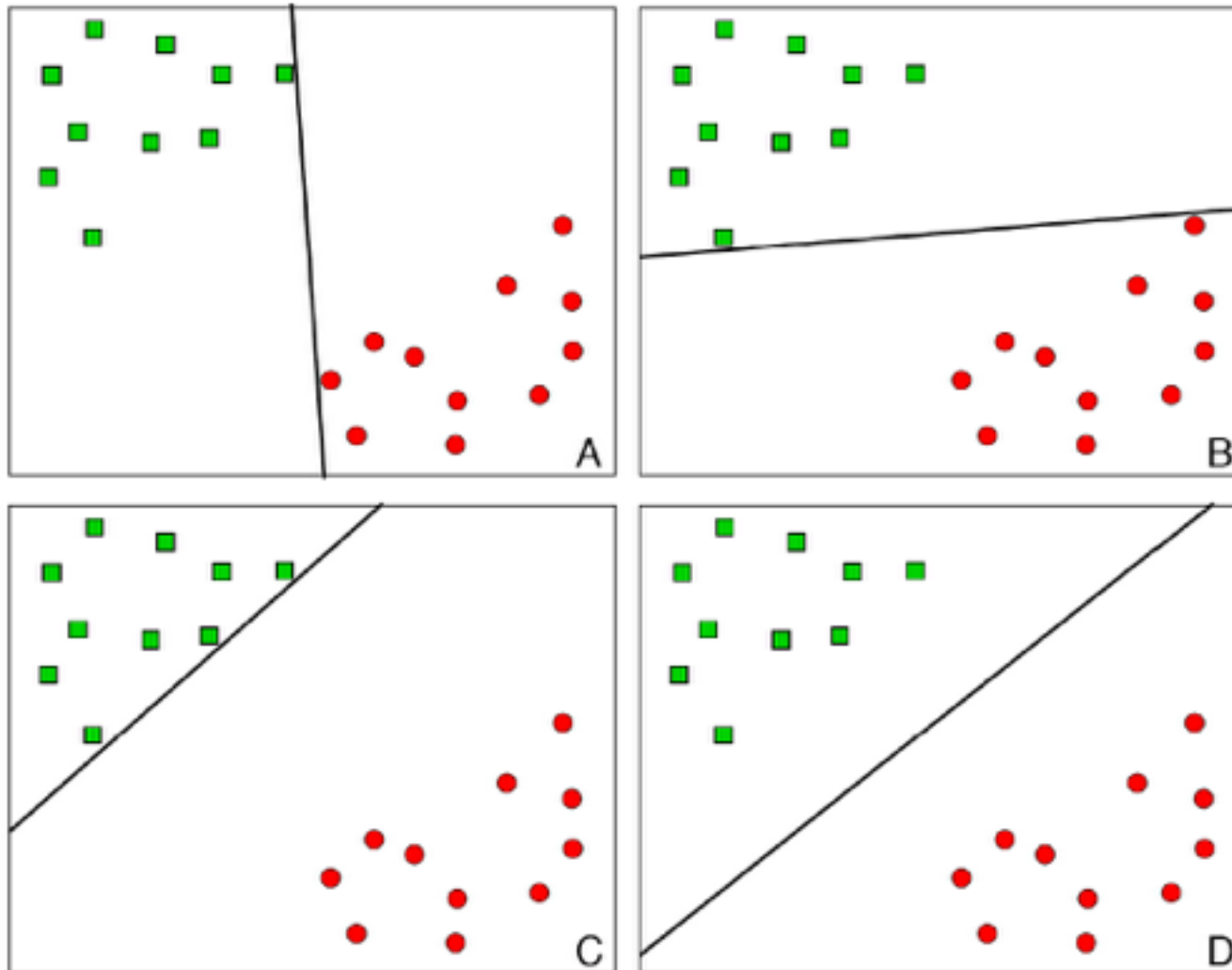
Support Vector Machine is the optimal classifier

-> see Machine Learning 1

“2D space (two codewords)”



# Which Hyperplane is best and why?



SVM classifier: Max-Margin behavior  
- best generalization

# Support Vector Machines

- For now: **linearly separable data**

$$x_i \in \mathbb{R}^d$$

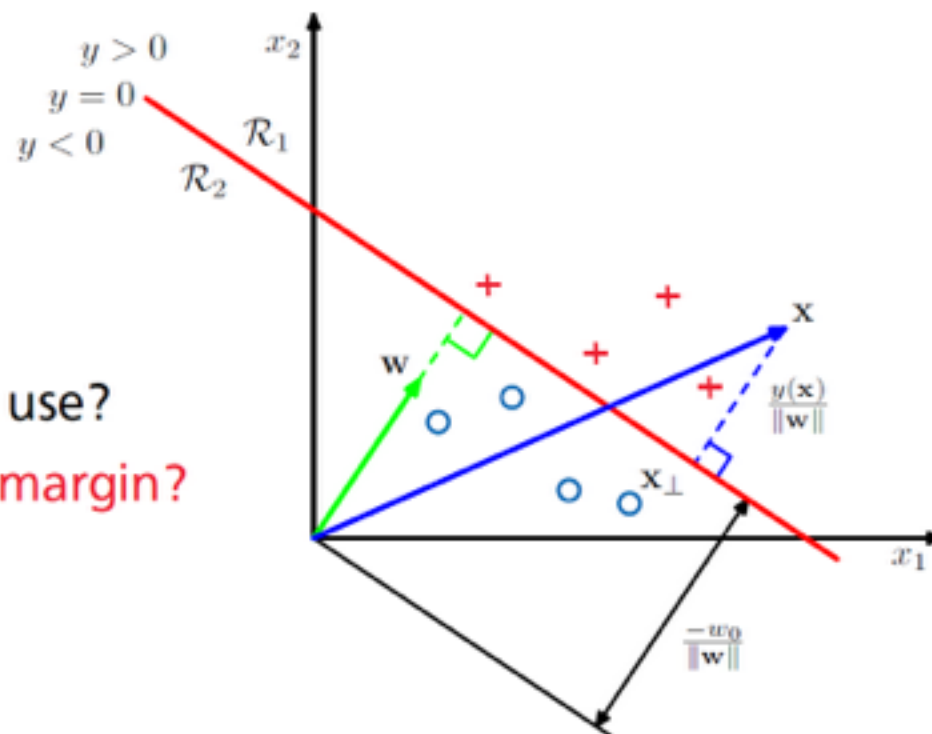
- N training data points:  $\{x_i, y_i\}_{i=1}^N$

$$y_i \in \{-1, 1\}$$

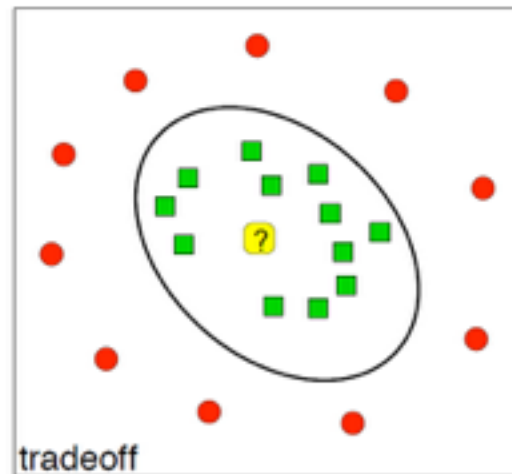
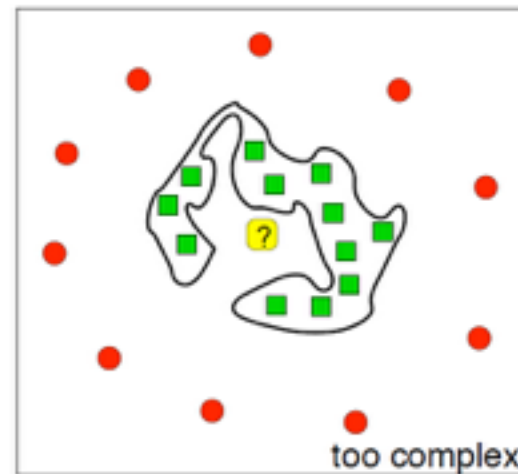
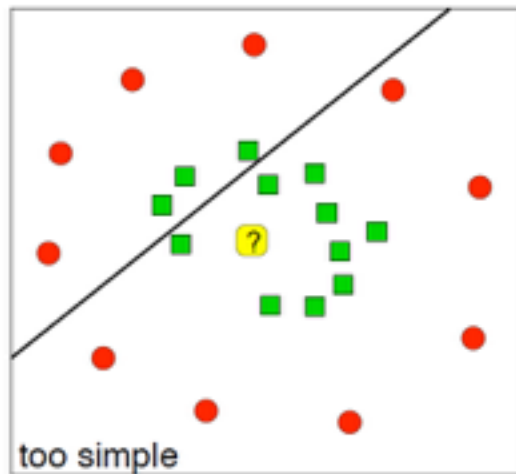
- Hyperplane that separates the data:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- Which hyperplane shall we use?
- How can we maximize the margin?**



# Simpler decision functions are better



- negative example
- positive example
- ? new patient

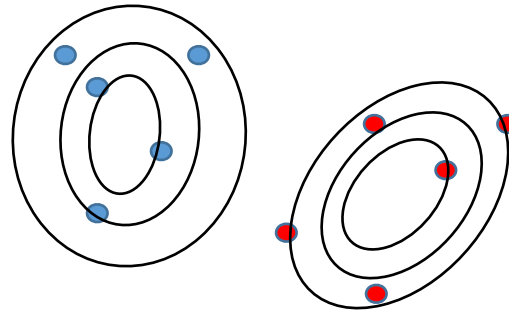
Best generalization

[Florian Markowetz]

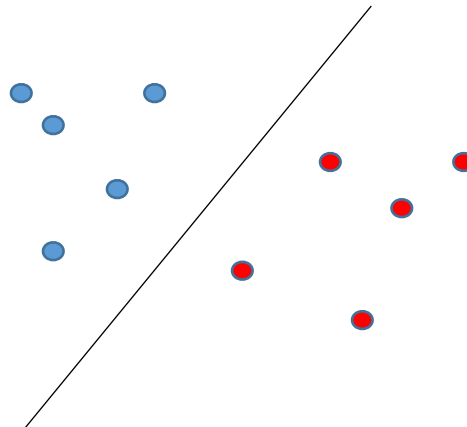
3 Minutes break

# Two approaches

Generative approach:  
*models distributions*



Discriminative function:  
*models decision function*



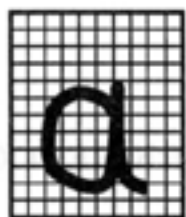


# Bayesian Decision Theory

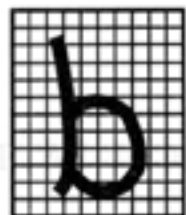
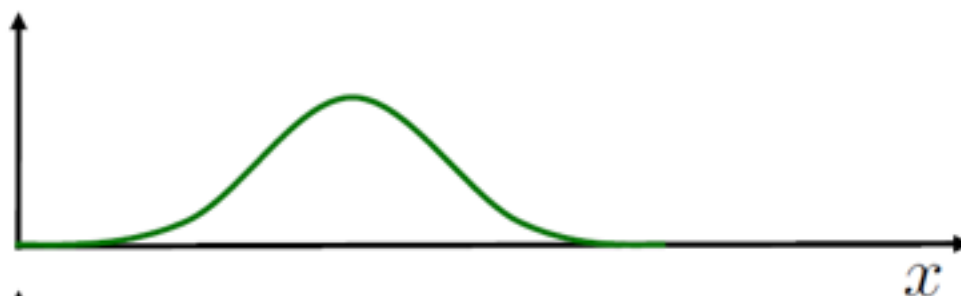
- 1st concept: **Class conditional probabilities**
  - Probability of making an observation  $x$  knowing that it comes from some class  $C_k$ .
  - Here  $x$  is a feature (vector).
  - $x$  measures / describes properties of the data.

$$p(x|C_k)$$

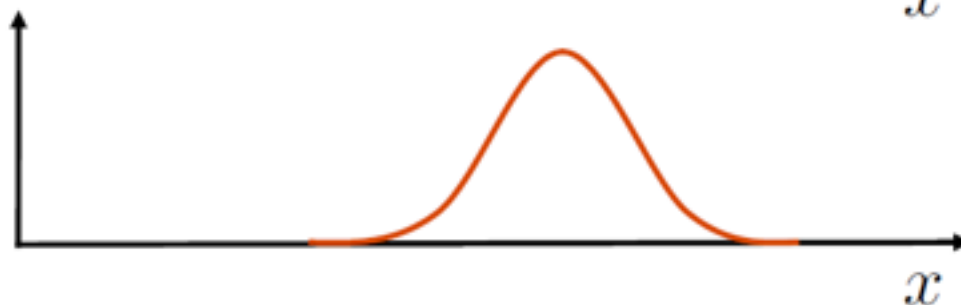
(Likelihood)



$$p(x|a)$$



$$p(x|b)$$



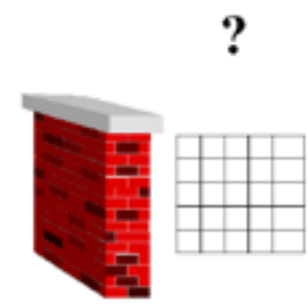
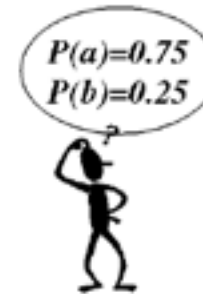
# Bayesian Decision Theory

- 2nd concept: **Class priors**

$$p(C_k)$$

(a priori probability of a data point belonging to a particular class)

- Example:



$$C_1 = a$$

$$p(C_1) = 0.75$$

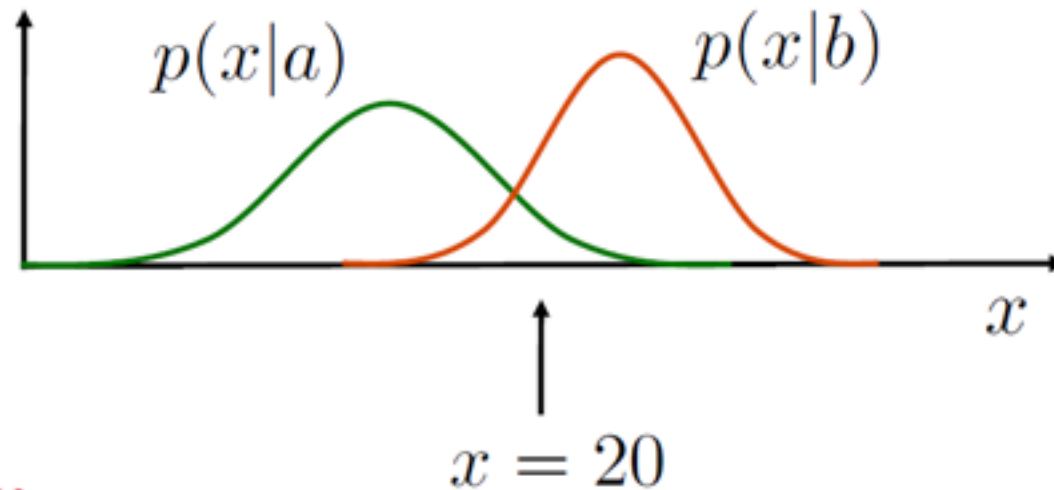
$$C_2 = b$$

$$p(C_2) = 0.25$$

- Generally:

$$\sum_k p(C_k) = 1$$

## ■ Example:



## ■ Question:

- How do we decide which class the data point belongs to?
- Remember that  $p(a) = 0.75$  and  $p(b) = 0.25$
- This means we **may** decide class  $a$ .

- Formalize this using Bayes' theorem:
  - We want to find the **a-posteriori probability** (posterior) of the class  $C_k$  given the observation (feature)  $x$

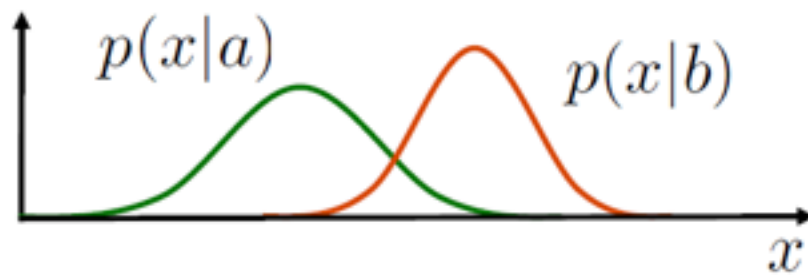
$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$$

Diagram illustrating Bayes' theorem with labels:

- class posterior (points to  $p(C_k|x)$ )
- class-conditional probability (likelihood) (points to  $p(x|C_k)$ )
- class prior (points to  $p(C_k)$ )
- normalization term (points to  $p(x)$ )

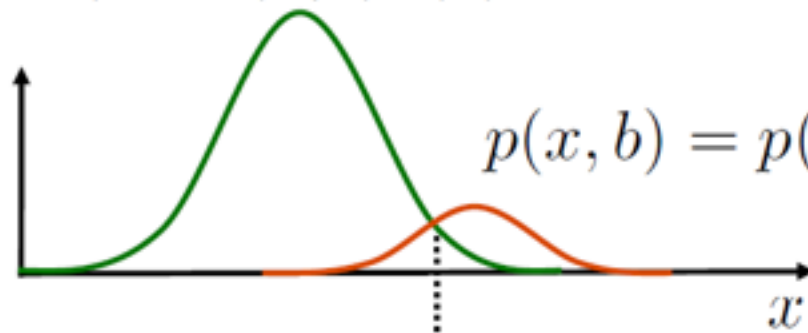
$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)} = \frac{p(x|C_k)p(C_k)}{\sum_j p(x|C_j)p(C_j)}$$

# Bayesian Decision Theory



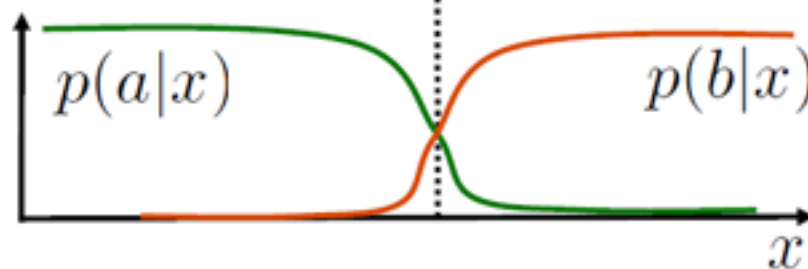
*Likelihood*

$$p(x, a) = p(x|a)p(a)$$



*Likelihood  $\times$  Prior*

← **decision boundary**



$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Normalization factor}}$$

## ■ Decision rule:

- Decide  $C_1$  if  $p(C_1|x) > p(C_2|x)$

We do not need  
the normalization!

- This is equivalent to

$$p(x|C_1)p(C_1) > p(x|C_2)p(C_2)$$

## ■ MAP classifier:

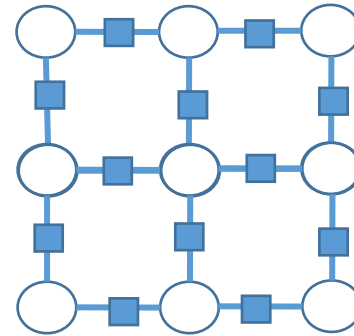
- A classifier obeying this rule is called a MAP classifier (sometimes called Bayes optimal classifier)



# Relation to previous lectures



- Image gets a label (class):  
 $K$  labelings

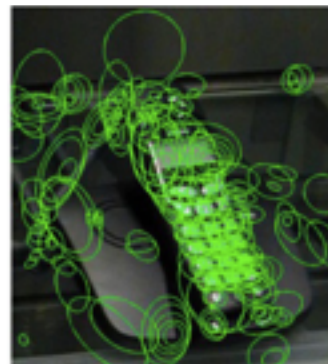


- Each pixel gets a label (class):  
 $K^n$  labelings
- Pixels are structured

# Naive Bayes Classifier

A naive Bayes classifier is a simple probabilistic [classifier](#) based on applying [Bayes' theorem](#) with strong (naive) [independence](#) assumptions

- Encode each image as a feature vector  $\mathbf{x} = (x^1, \dots, x^n)$  where  $n$  is the number of interest points.
- $x^j \in \{w_1, \dots, w_m\}$ . Here  $m$  visual words.



~200 interest points



Interest points for codewords (visual words)

- Naive Bayes Classifier assumes that visual words are **conditionally independent** given object class:  $P(\mathbf{x}|c) = \prod_j P(x^j|c)$  (which is rarely true in practice)
- **Naive Bayes Classifier:**  
$$c^* = \operatorname{argmax}_c P(c|\mathbf{x}) = \operatorname{argmax}_c P(c) P(\mathbf{x}|c) = \operatorname{argmax}_c P(c) \prod_j P(x^j|c)$$

# Image Classification with Naive Bayes

- Image dataset: 7 object categories, arbitrary views, partial occlusions

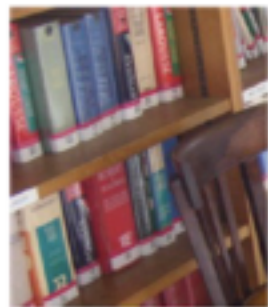


# Image Classification with Naive Bayes

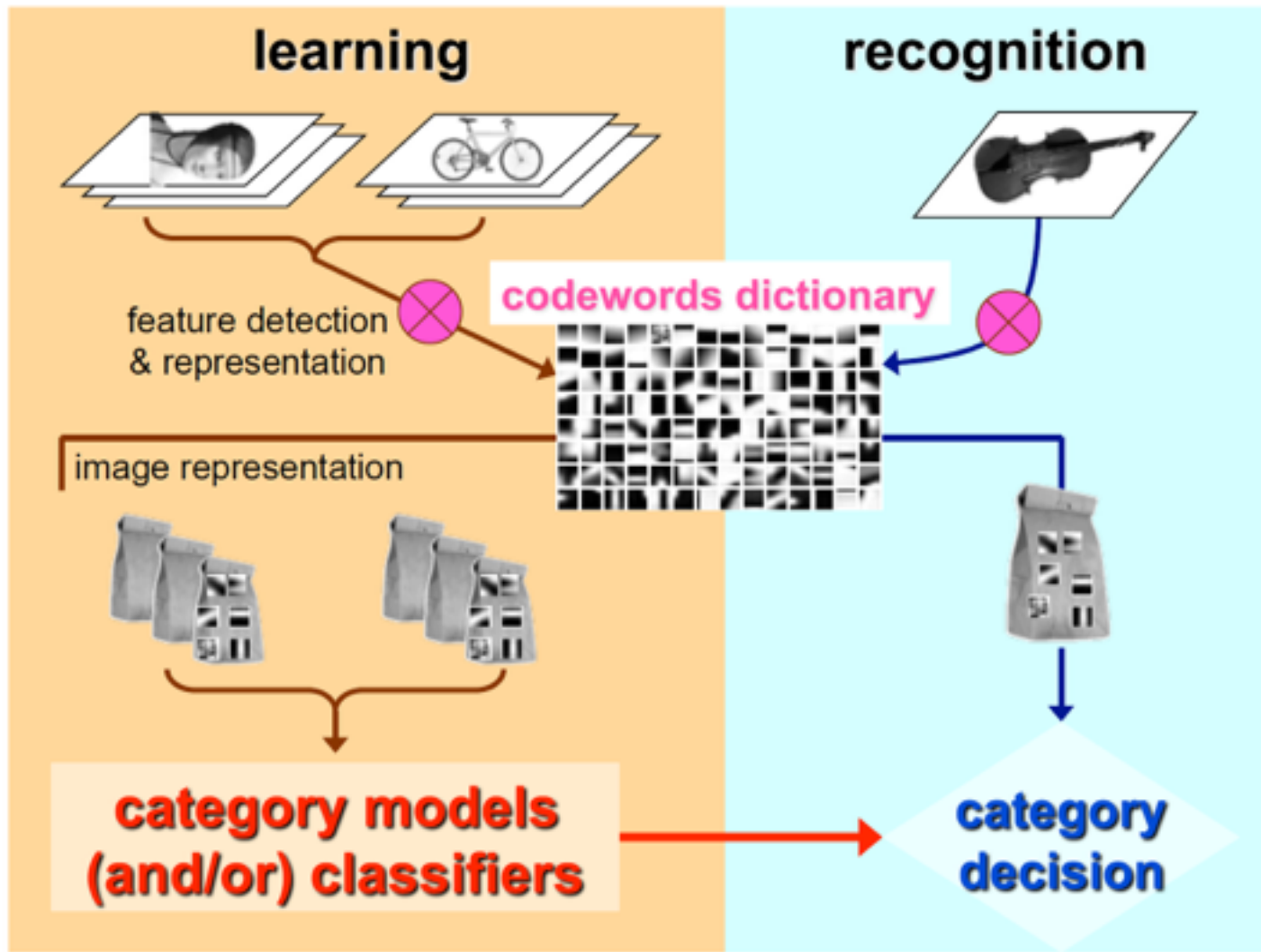
**Table 1.** Confusion matrix and the mean rank for the best vocabulary ( $k=1000$ ).

True classes →	<i>faces</i>	<i>buildings</i>	<i>trees</i>	<i>cars</i>	<i>phones</i>	<i>bikes</i>	<i>books</i>
<i>faces</i>	76	4	2	3	4	4	13
<i>buildings</i>	2	44	5	0	5	1	3
<i>trees</i>	3	2	80	0	0	5	0
<i>cars</i>	4	1	0	75	3	1	4
<i>phones</i>	9	15	1	16	70	14	11
<i>bikes</i>	2	15	12	0	8	73	0
<i>books</i>	4	19	0	6	7	2	69
<i>Mean ranks</i>	1.49	1.88	1.33	1.33	1.63	1.57	1.57

Examples of correctly classified images:



# Bag of words - Done!





# Summary and Discussion

- **Bag of words representation:**
  - Sparse representation of object categories
  - Many Machine learning techniques can be applied (here naïve Bayes and SVM)
  - Robust to occlusion
  - Allows sharing of representation between multiple classes (via codeword dictionary)
- **Problems:**
  - Spatial distribution of visual works is not modelled.





# Roadmap (this lecture)

- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching

# Spatial Pyramid Matching

- Add spatial information to the bag-of-features
- Perform matching in 2D image space



[Lazebnik, Schmid & Ponce, CVPR 2006]

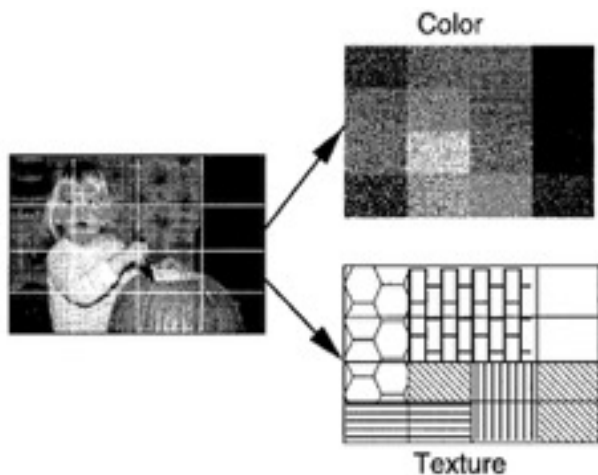
# Spatial Pyramid Matching

Similar approaches:

Subblock description [Szummer & Picard, 1997]

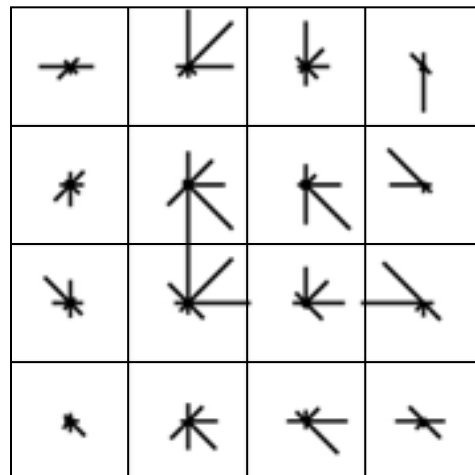
SIFT [Lowe, 1999]

GIST [Torralba et al., 2003]



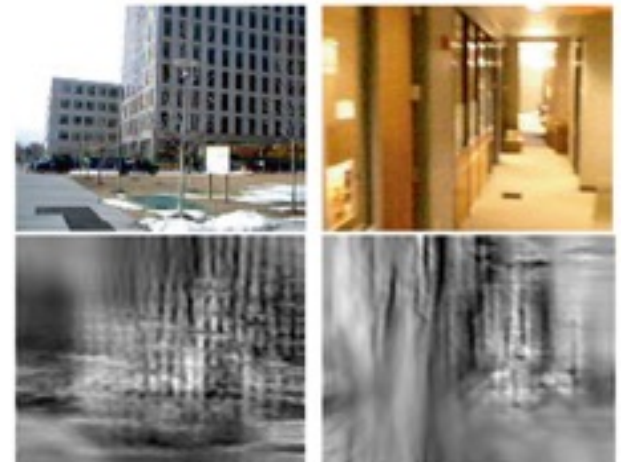
Szummer & Picard (1997)

SIFT



Lowe (1999, 2004)

Gist

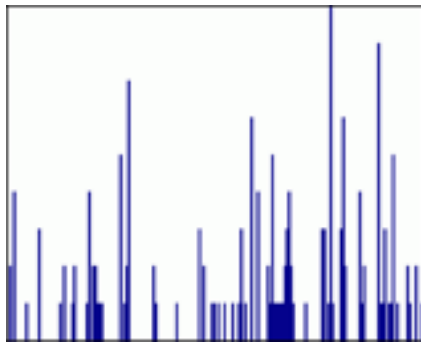


Torralba et al. (2003)

# Spatial pyramid representation



Locally orderless  
representation at  
several levels of  
spatial resolution

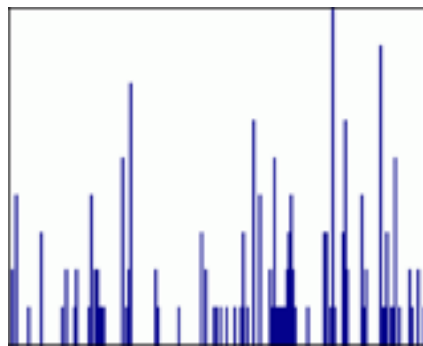


level 0

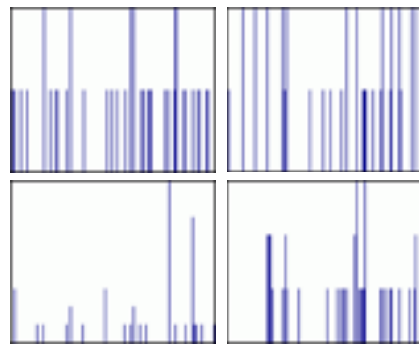
# Spatial pyramid representation



Locally orderless  
representation at  
several levels of  
spatial resolution

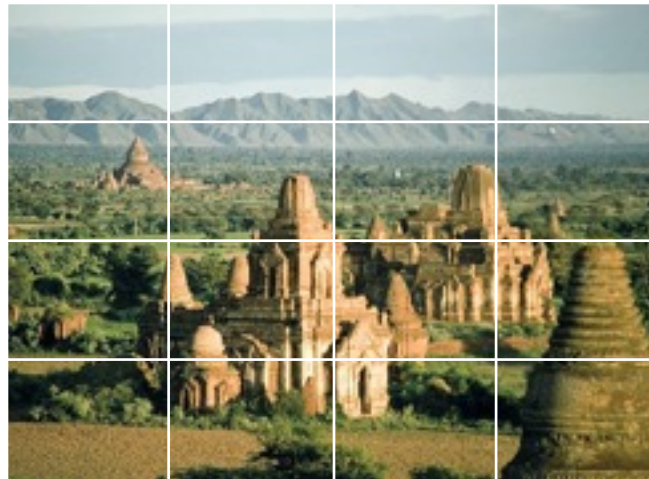


level 0

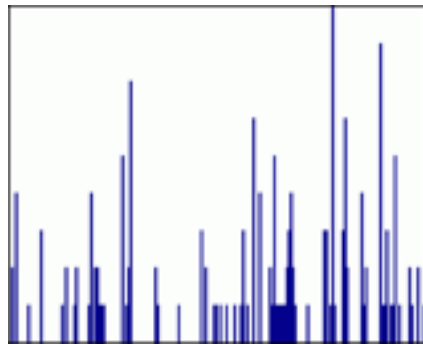


level 1

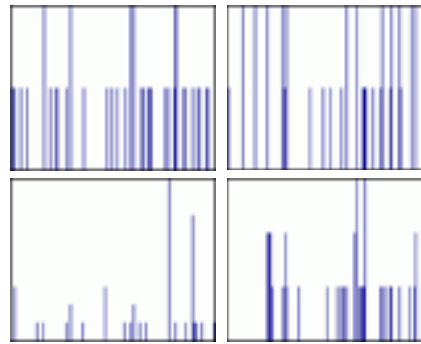
# Spatial pyramid representation



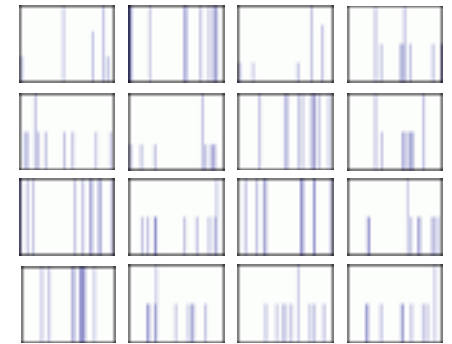
Locally orderless  
representation at  
several levels of  
spatial resolution



level 0



level 1

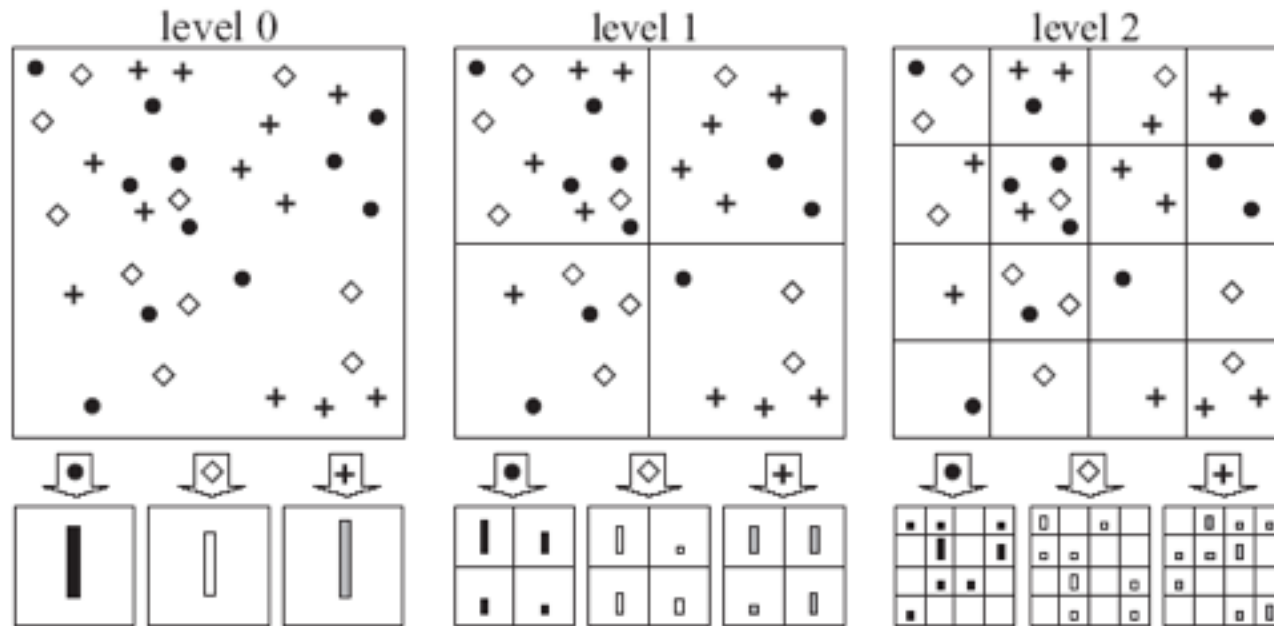


level 2

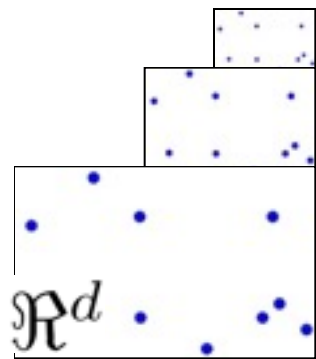


# Spatial Pyramid Matching

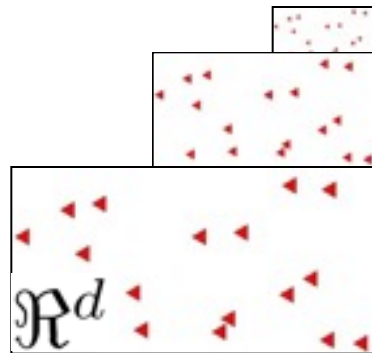
- Combination of spatial levels with pyramid match kernel  
[Grauman & Darrell'05]



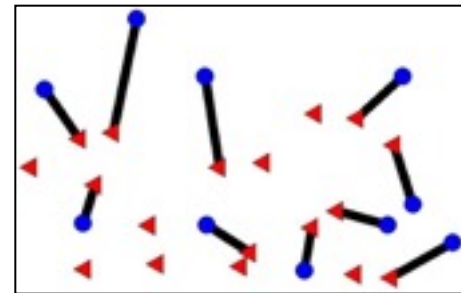
# Pyramid Matching Kernel



$\cup$



$\approx$

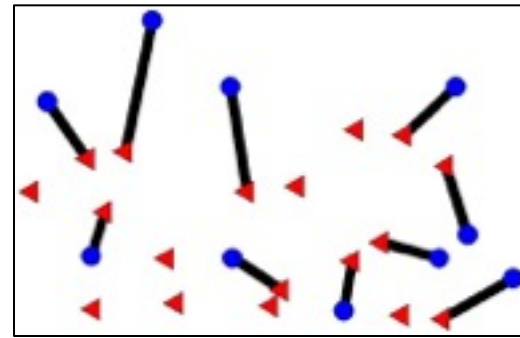
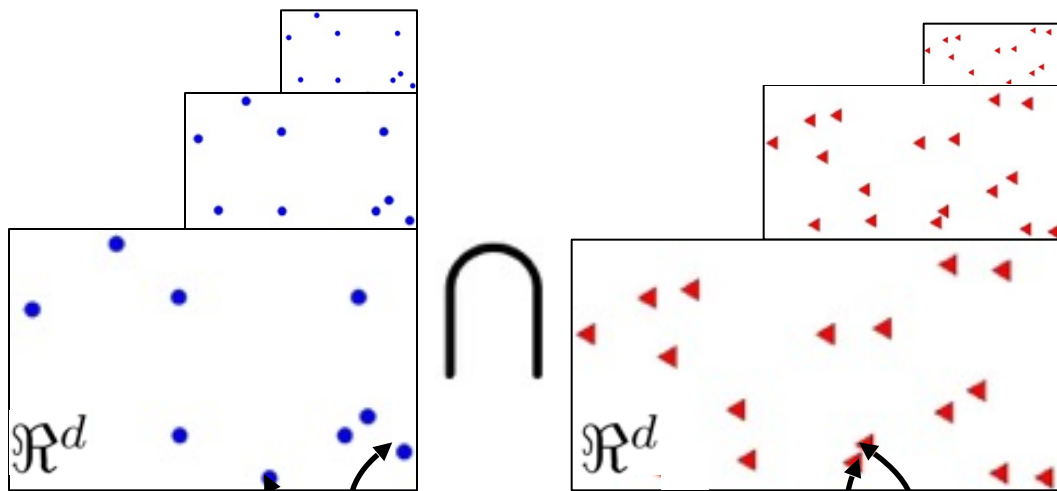


optimal partial  
matching between  
sets of features



Slides Credit: Kristen Grauman

# Pyramid Matching Kernel



optimal partial matching

$$\max_{\pi: X \rightarrow Y} \sum_{x_i \in X} \mathcal{S}(x_i, \pi(x_i))$$

$$X = \{\vec{x}_1, \dots, \vec{x}_m\}$$
$$\vec{x}_i \in \mathbb{R}^d$$

$$Y = \{\vec{y}_1, \dots, \vec{y}_n\}$$
$$\vec{y}_i \in \mathbb{R}^d$$

# Pyramid match overview

---

Pyramid match kernel measures similarity of a partial matching between two sets:

- Place multi-dimensional, multi-resolution grid over point sets
- Consider points matched at finest resolution where they fall into same grid cell
- Approximate similarity between matched points with worst case similarity at given level

# Pyramid match kernel

Number of newly  
matched pairs at level  $i$

Approximate  
partial match  
similarity

$$K_{\Delta} = \sum_{i=0}^L w_i N_i$$

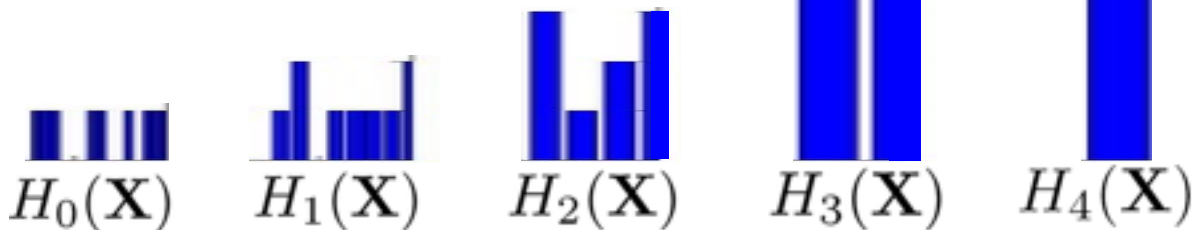
Measure of difficulty  
of a match at level  $i$

# Feature extraction

$$\mathbf{X} = \{\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_m\}, \quad \vec{\mathbf{x}}_i \in \mathbb{R}^d$$



Histogram pyramid:  
level  $i$  has bins of size  $2^i$



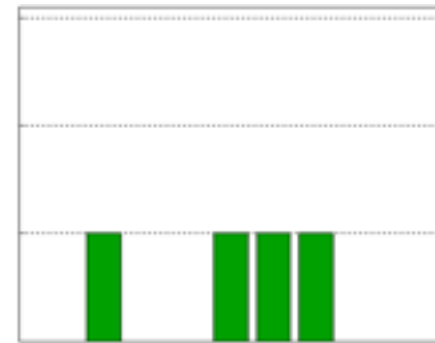
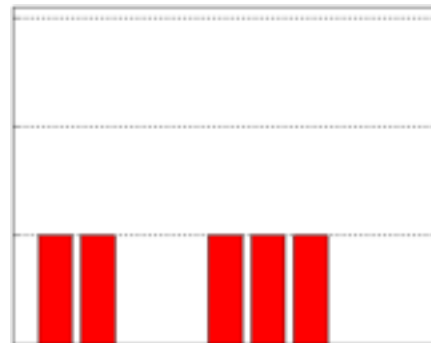
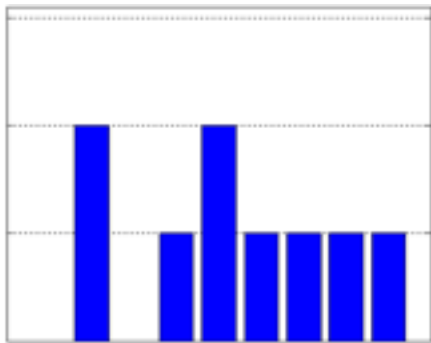
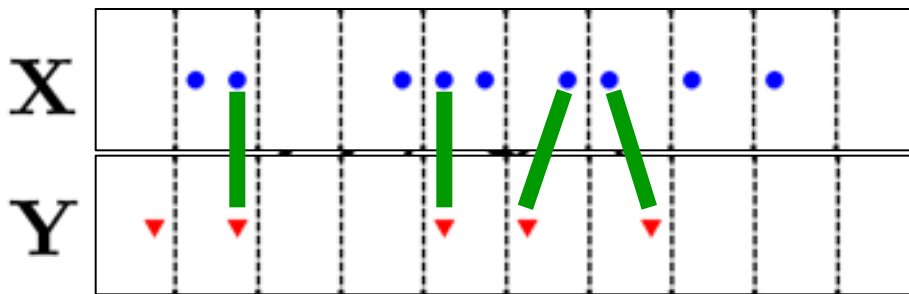
$$\Psi(\mathbf{X}) = [H_0(\mathbf{X}), \dots, H_L(\mathbf{X})]$$



# Counting matches

Histogram  
intersection

$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = \sum_{j=1}^r \min(H(\mathbf{X})_j, H(\mathbf{Y})_j)$$



$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = 4$$

# Counting new matches

Histogram  
intersection

$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = \sum_{j=1}^r \min(H(\mathbf{X})_j, H(\mathbf{Y})_j)$$

matches at this level

matches at previous level

$$N_i = \mathcal{I}(H_i(\mathbf{X}), H_i(\mathbf{Y})) - \mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y}))$$

Difference in histogram intersections across levels counts *number of new pairs* matched

# Pyramid match kernel

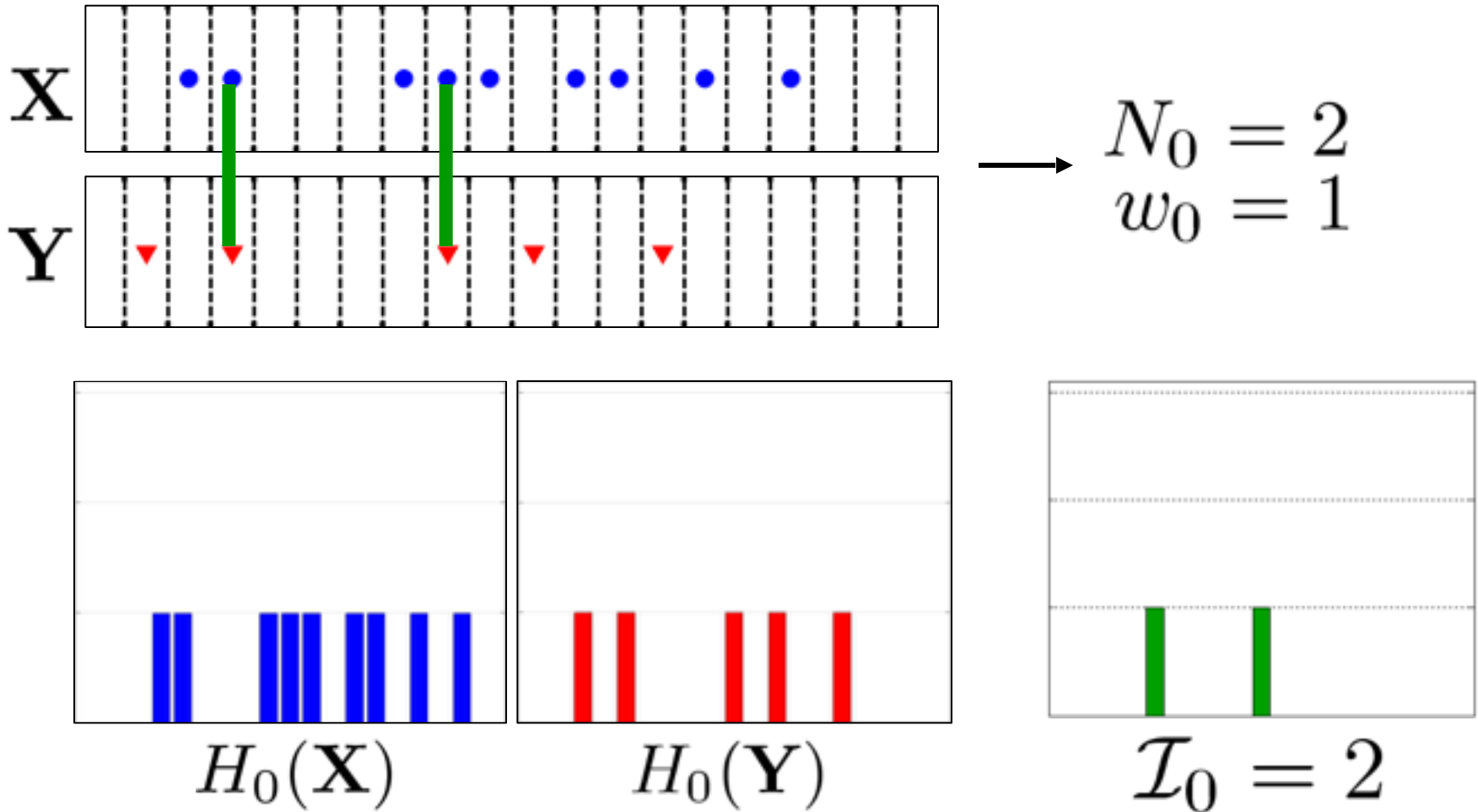
$$K_{\Delta} (\overbrace{\Psi(\mathbf{X}), \Psi(\mathbf{Y})}^{\text{histogram pyramids}}) = \sum_{i=0}^L \frac{1}{2^i} \left( \underbrace{\mathcal{I}(H_i(\mathbf{X}), H_i(\mathbf{Y})) - \mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y}))}_{\text{number of newly matched pairs at level } i} \right)$$

↑  
measure of difficulty of a  
match at level  $i$

- Weights inversely proportional to bin size
- Normalize kernel values to avoid favoring large sets

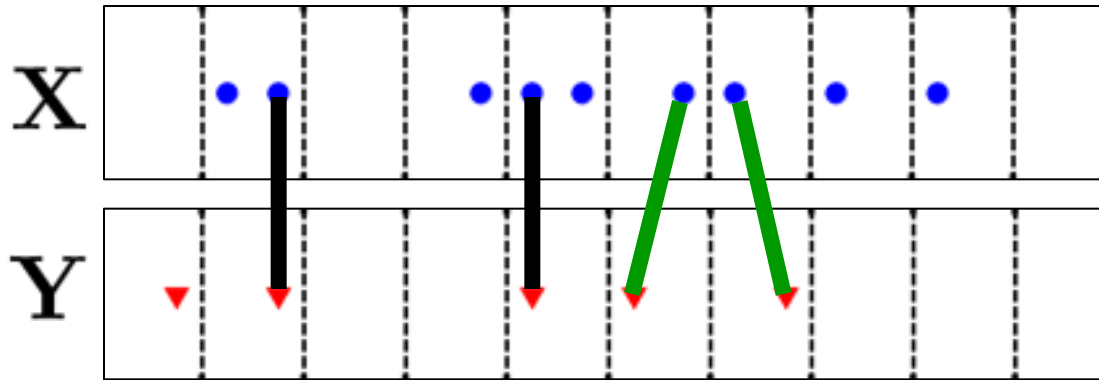
# Example pyramid match

Level 0

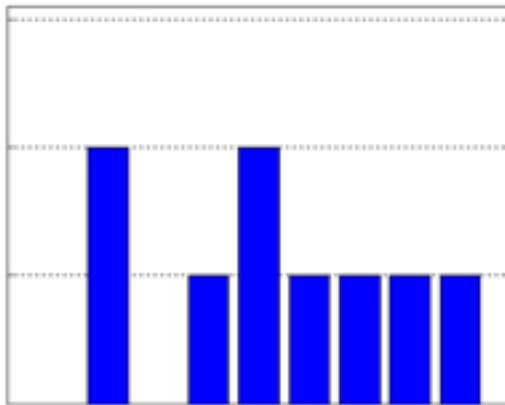


# Example pyramid match

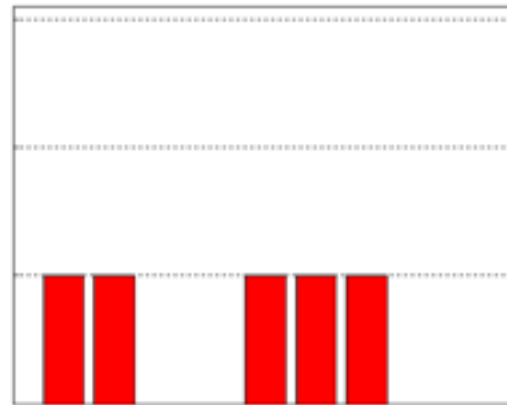
Level 1



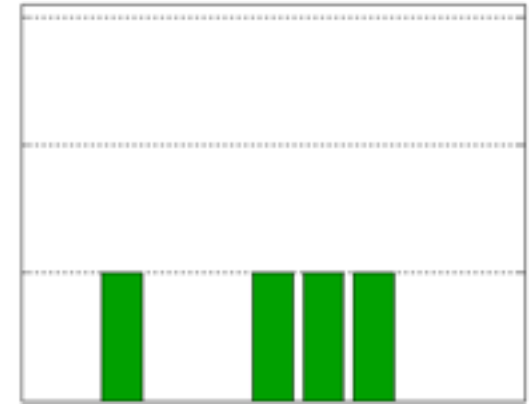
$$\begin{aligned} N_1 &= 4 - 2 = 2 \\ w_1 &= \frac{1}{2} \end{aligned}$$



$H_1(\mathbf{X})$



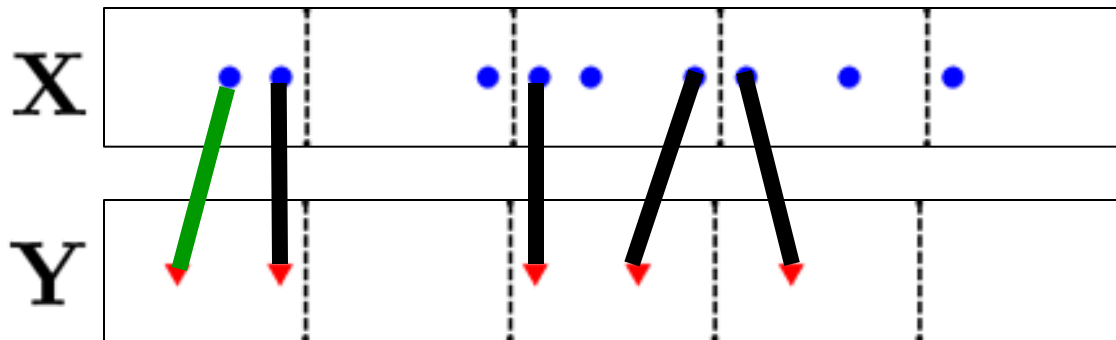
$H_1(\mathbf{Y})$



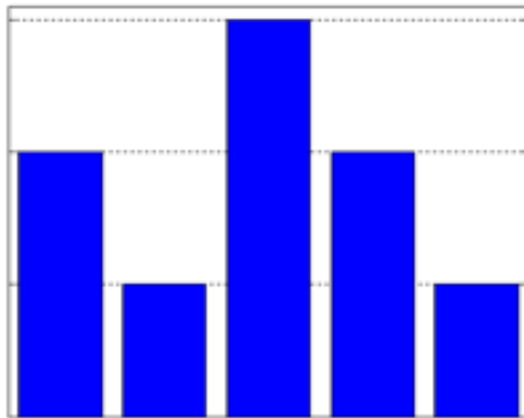
$\mathcal{I}_1 = 4$

# Example pyramid match

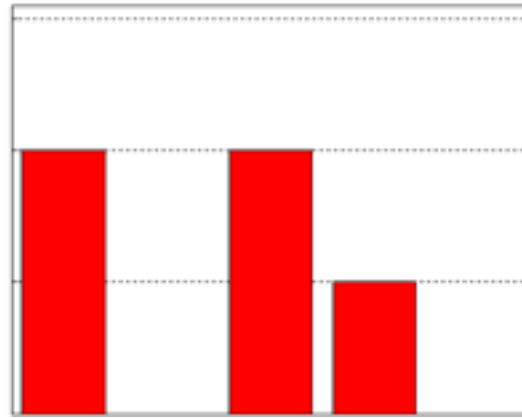
Level 2



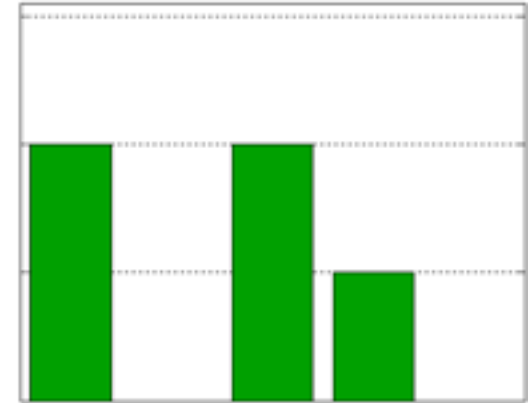
$$\begin{aligned} N_2 &= 5 - 4 = 1 \\ w_2 &= \frac{1}{4} \end{aligned}$$



$H_2(\mathbf{X})$



$H_2(\mathbf{Y})$

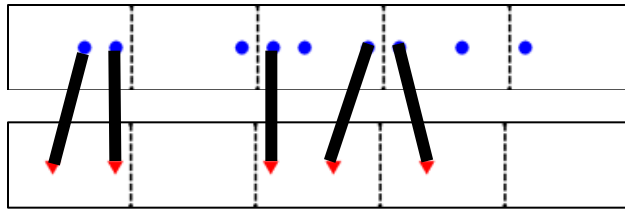


$\mathcal{I}_2 = 5$



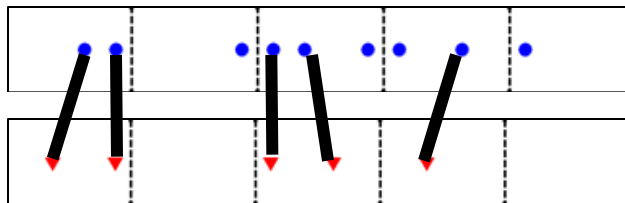
# Example pyramid match

pyramid match



$$\begin{aligned} K_{\Delta} &= \sum_{i=0}^L w_i N_i \\ &= 1(2) + \frac{1}{2}(2) + \frac{1}{4}(1) = 3.25 \end{aligned}$$

optimal match



$$\begin{aligned} K &= \max_{\pi: \mathbf{X} \rightarrow \mathbf{Y}} \sum_{\mathbf{x}_i \in \mathbf{X}} \mathcal{S}(\mathbf{x}_i, \pi(\mathbf{x}_i)) \\ &= 1(2) + \frac{1}{2}(3) = 3.5 \end{aligned}$$

# Scene Classification



L	Single-level	Pyramid
0(1x1)	72.2±0.6	
1(2x2)	77.9±0.6	79.0 ±0.5
2(4x4)	79.4±0.3	81.1 ±0.3
3(8x8)	77.2±0.4	80.7 ±0.3

# Retrieval Examples



(a) kitchen



living room



living room



living room



office



living room



living room



living room



living room



(b) kitchen



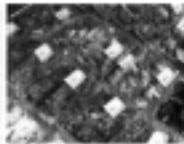
office



inside city



(c) store



mountain



forest



(d) tall bldg



inside city

inside city



(e) tall bldg



inside city

mountain

mountain

mountain

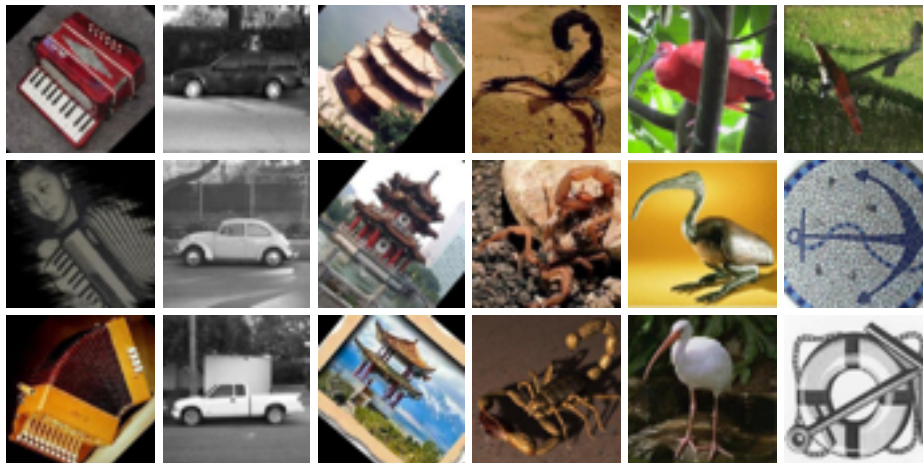


(f) inside city



tall bldg

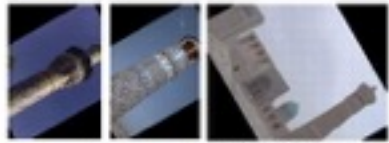
# Category classification - CalTech101



L	Single-level	Pyramid
0(1x1)	41.2±1.2	
1(2x2)	55.9±0.9	57.0 ±0.8
2(4x4)	63.6±0.9	64.6 ±0.8
3(8x8)	60.3±0.9	64.6 ±0.7

Bag-of-words approach by Zhang et al.'07: 54 %

## Easiest and hardest classes



minaret (97.6%)



windsor chair (94.6%)



joshua tree (87.9%)



okapi (87.8%)



cougar body (27.6%)



beaver (27.5%)



crocodile (25.0%)



ant (25.0%)

- Sources of difficulty:
  - Lack of texture
  - Camouflage
  - Thin, articulated limbs
  - Highly deformable shape

- **Summary**

- Spatial pyramid representation: appearance of local image patches + coarse global position information
- Substantial improvement over bag of features
- Depends on the similarity of image layout

- **Extensions**

- Integrating different types of features, learning weights, use of different grids
- Flexible, object-centered grid



# Roadmap (this lecture)

- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching

# Roadmap (this lecture)

- Image Categorization
- Bag-of-Words (BOW)
- Generative vs. Discriminative Approach
- Spatial Pyramid Matching
- Application: Remote Sensing Image Classification