

上海交通大学

# 计算机视觉

教师: 赵旭

班级: AI4701

2024 春

# 16. 分类与识别

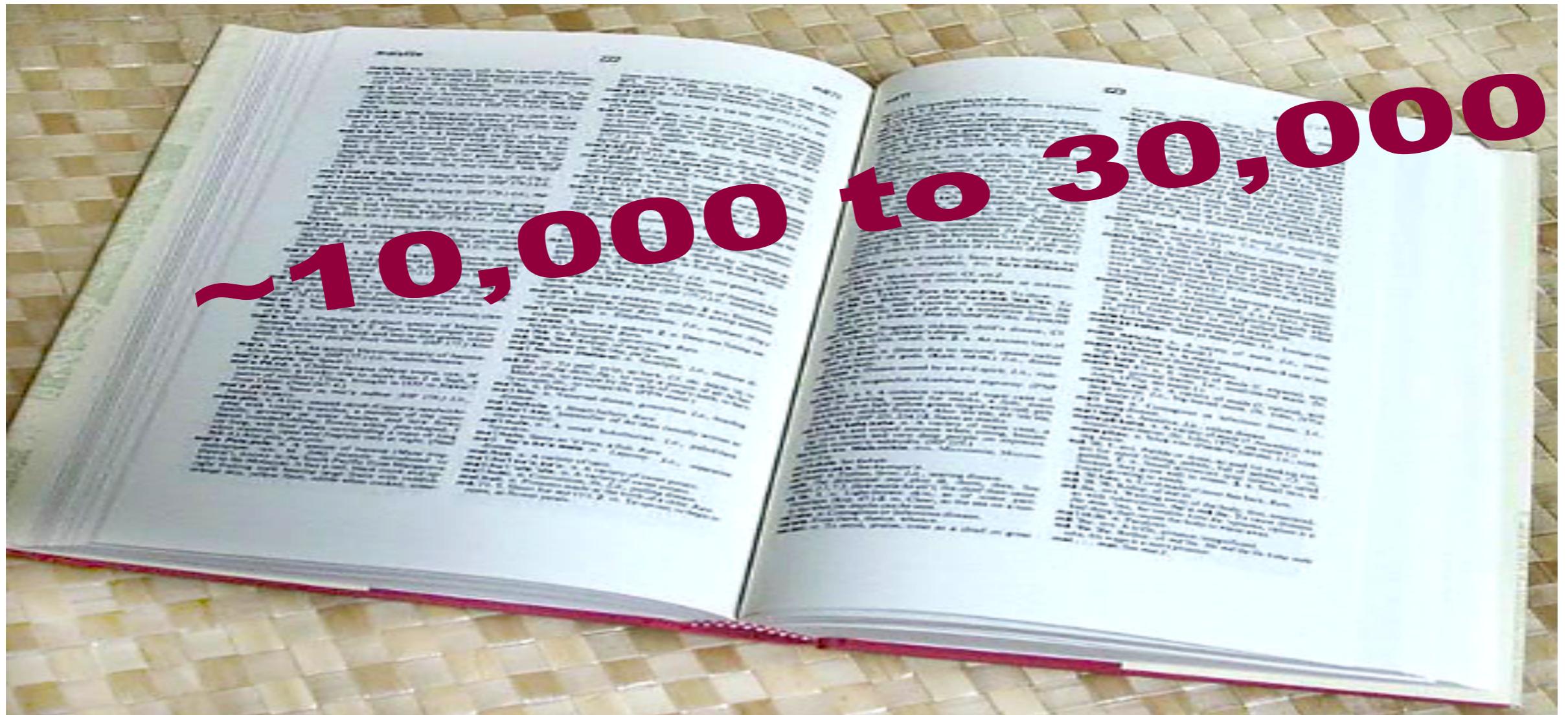
---

# Contents

---

- ❖ **Recognition: Overview and History**
- ❖ **Machine Learning for Recognition**

# How many visual object categories are there?



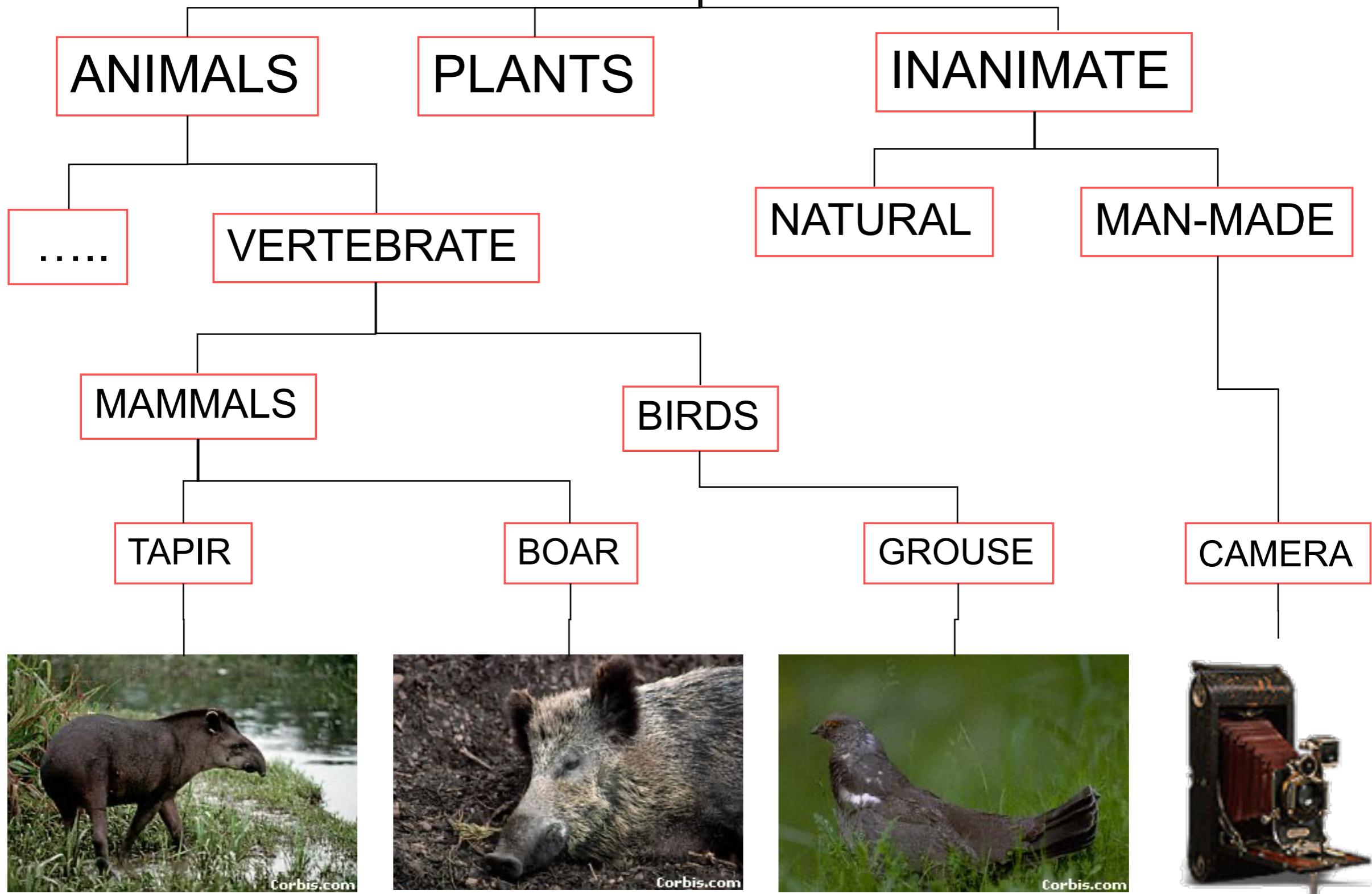
1500-3000 basic-level nouns, ~10 types per basic-level category

Alternative explanation (Perona): ~1000 names per domain (broad scene category), 20-30 domains

**~10,000 to 30,000**



# OBJECTS



# Recognition tasks



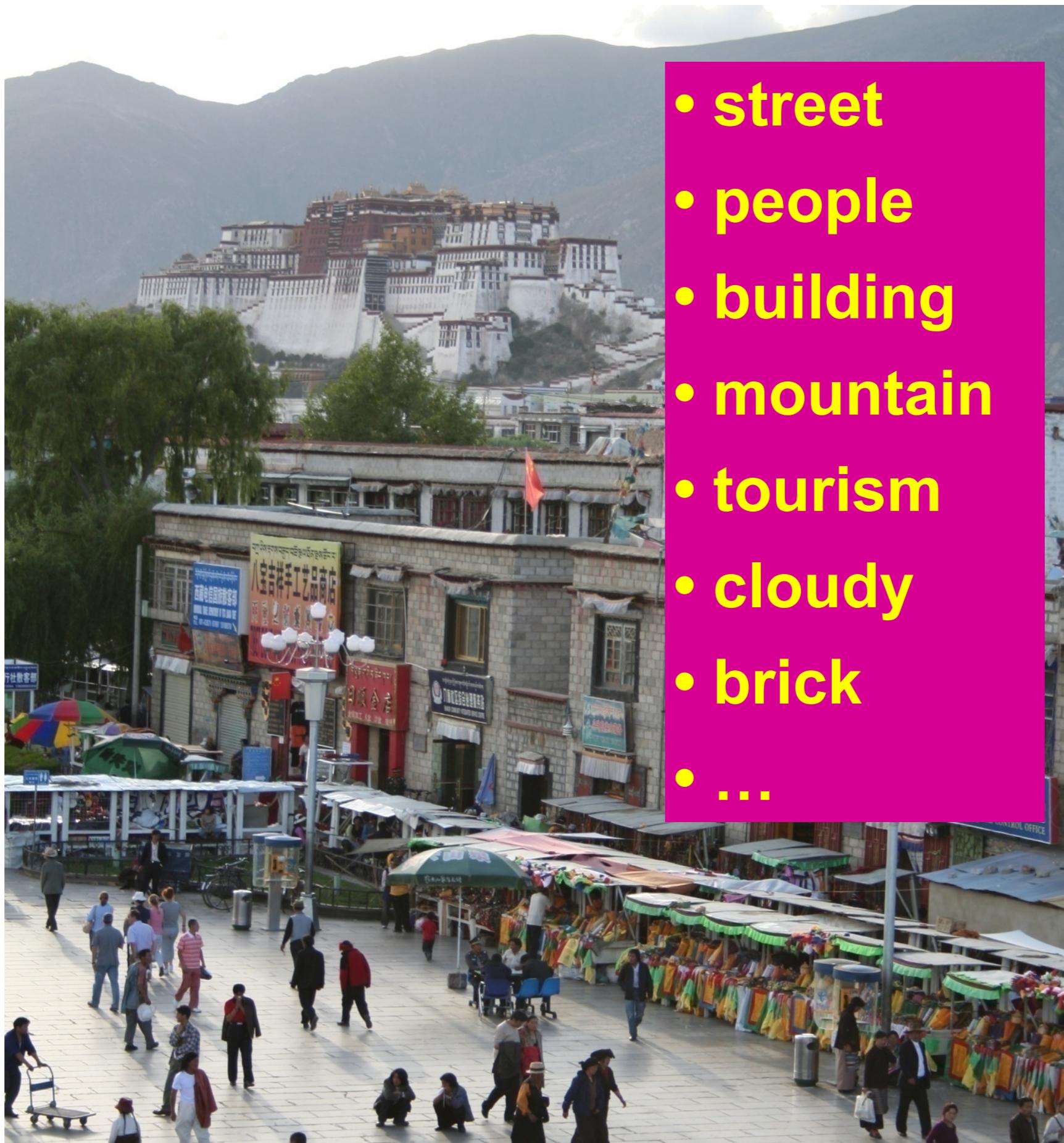
Svetlana Lazebnik

# Scene categorization or classification

- outdoor/indoor
- city/forest/factory/etc.



# Image annotation / tagging / attributes



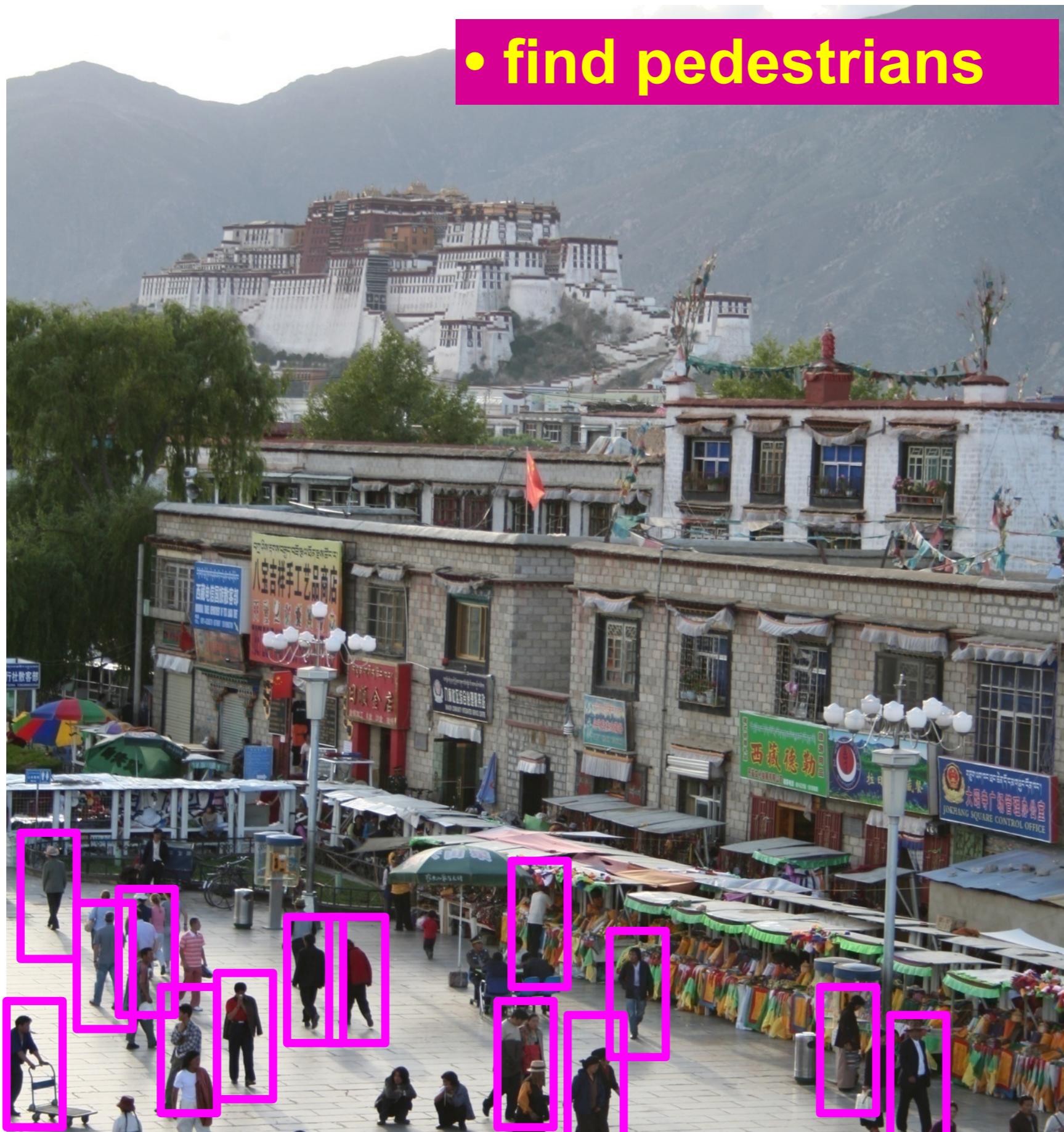
- street
- people
- building
- mountain
- tourism
- cloudy
- brick
- ...

# Image parsing / semantic segmentation



# Object detection

- find pedestrians



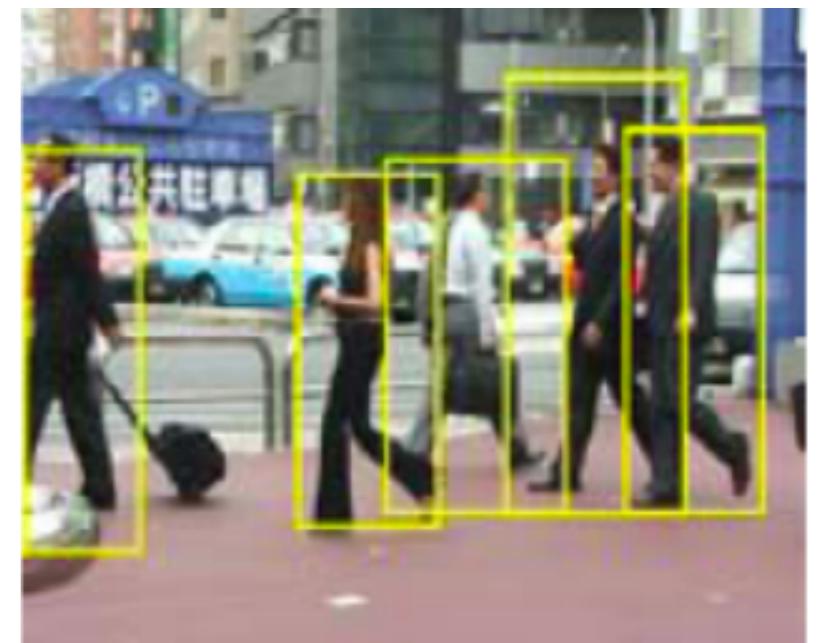
# Scene understanding?



# Category vs. instance recognition

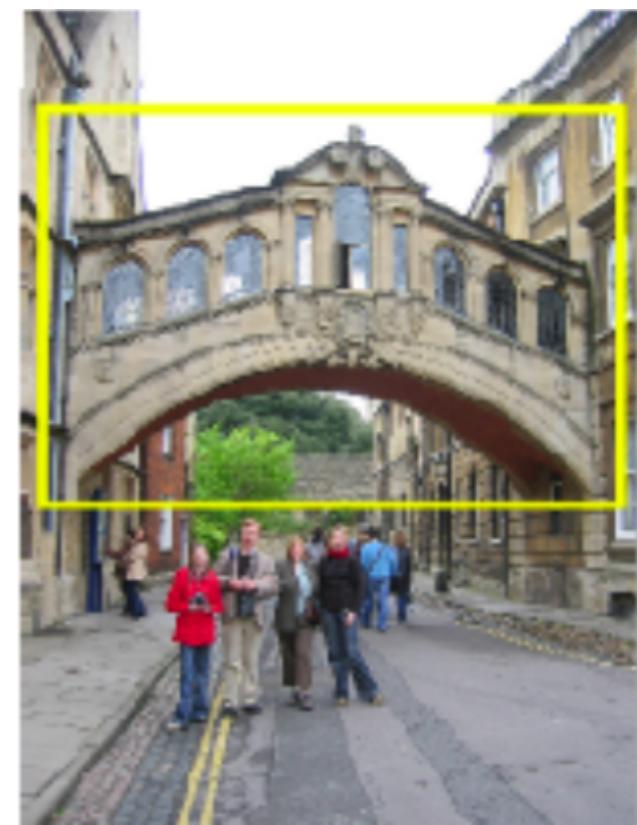
Category:

- ❖ Find all the people
- ❖ Find all the buildings
- ❖ Often within a single image
- ❖ Often ‘sliding window’

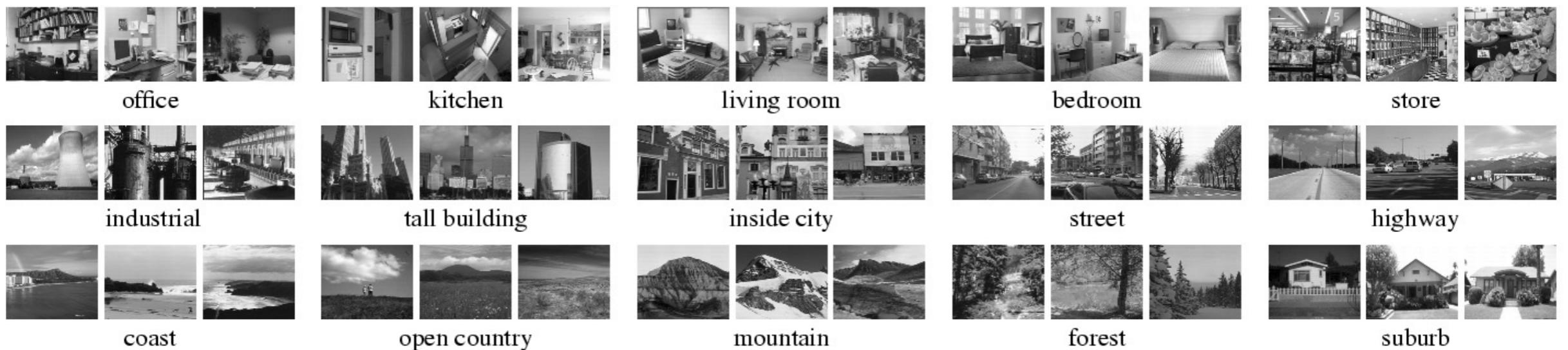


Instance:

- ❖ Is this face Trump?
- ❖ Find this specific famous building
- ❖ Often within a database of images



# Scene recognition dataset



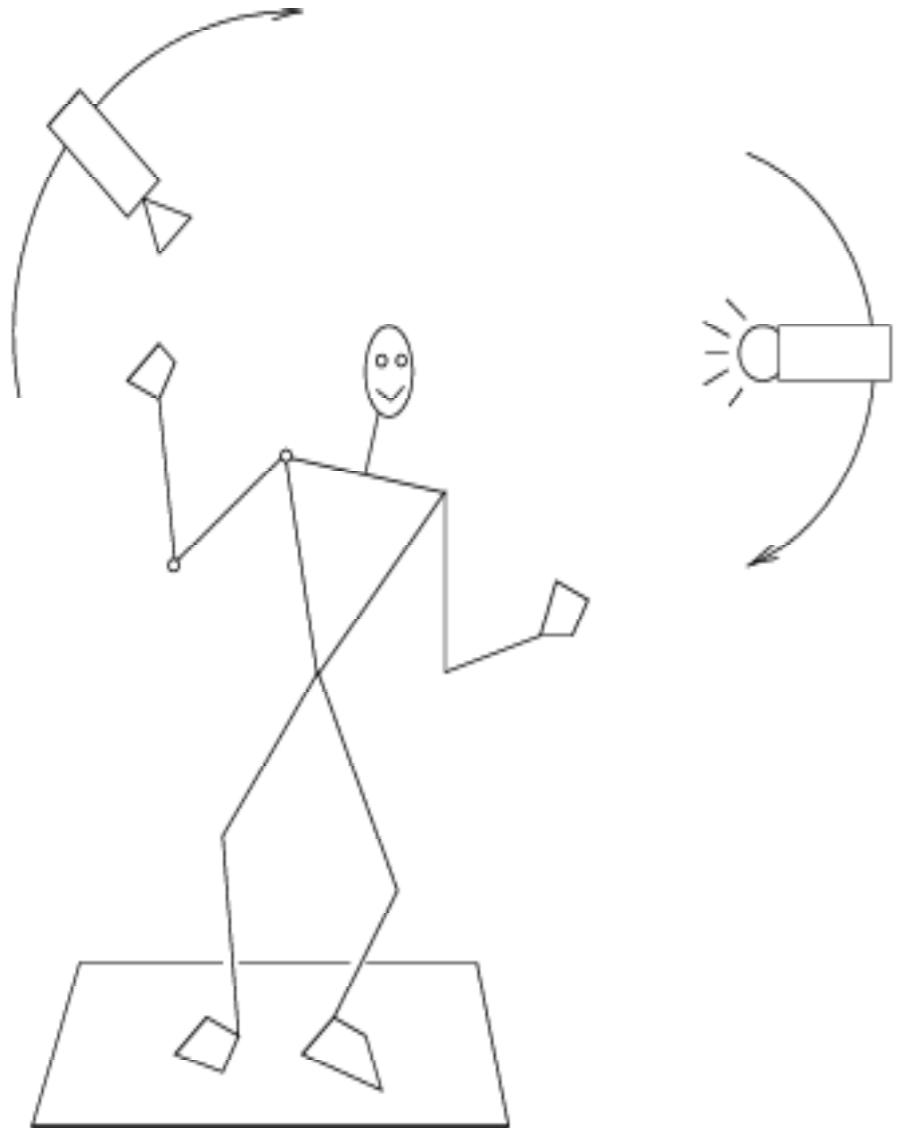
Instance or category?

# Recognition is all about modeling variability



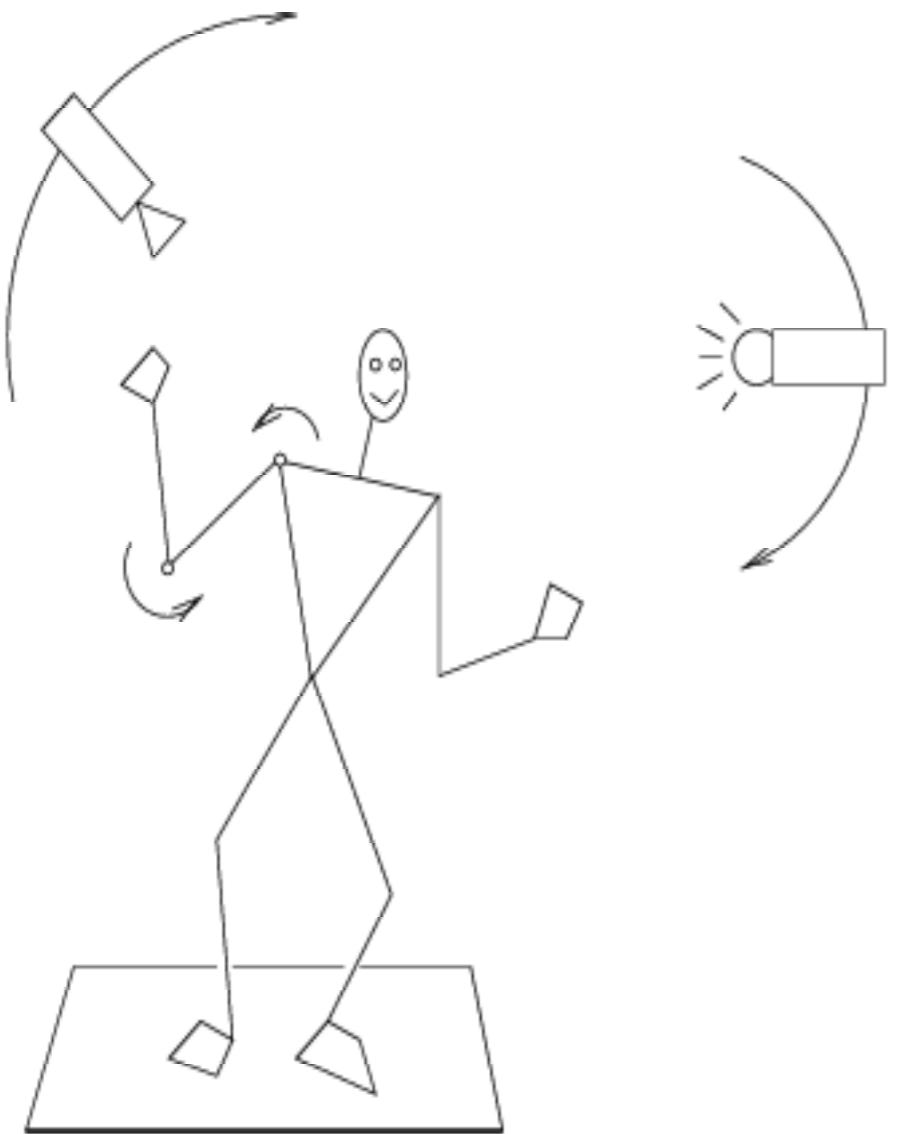
Variability: Camera position

# Recognition is all about modeling variability



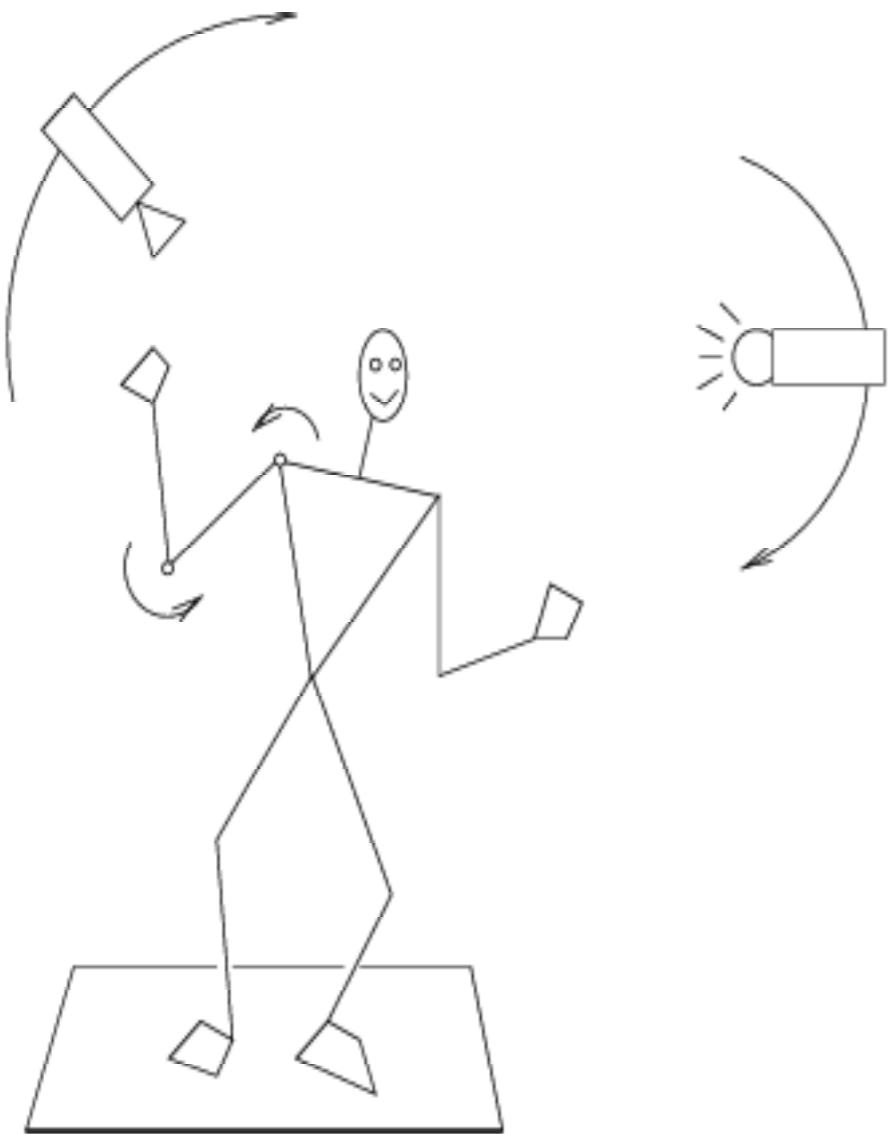
Variability:

- Camera position
- Illumination



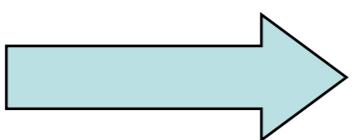
Variability:

- Camera position
- Illumination
- Pose/shape parameters

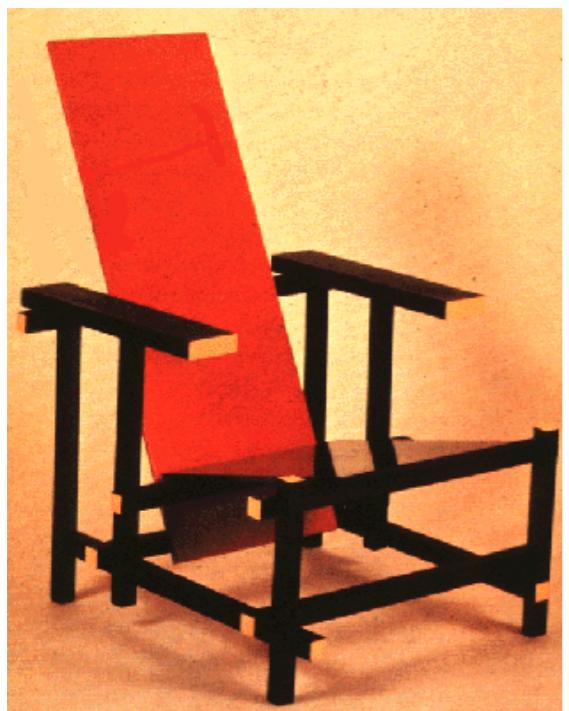


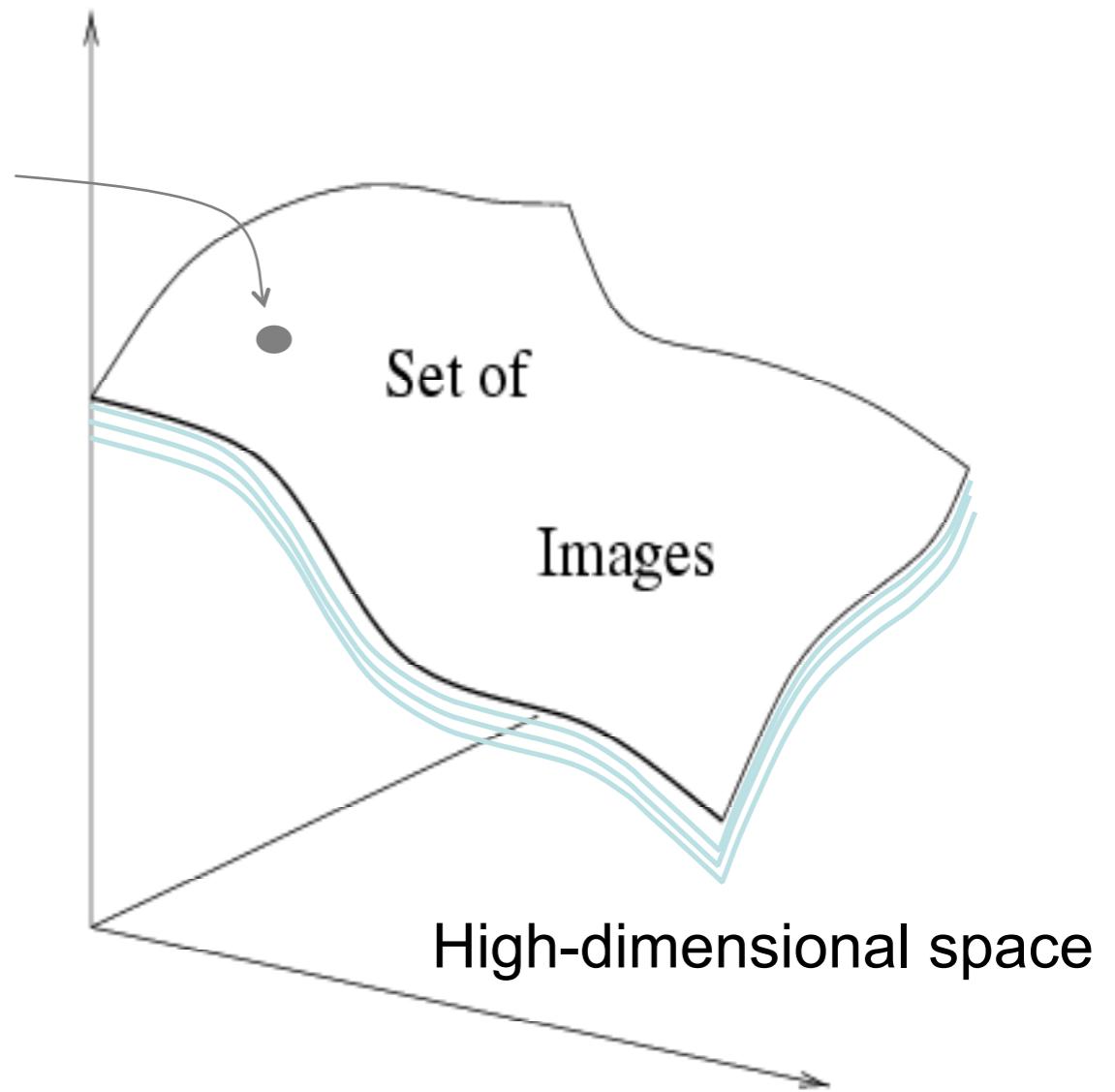
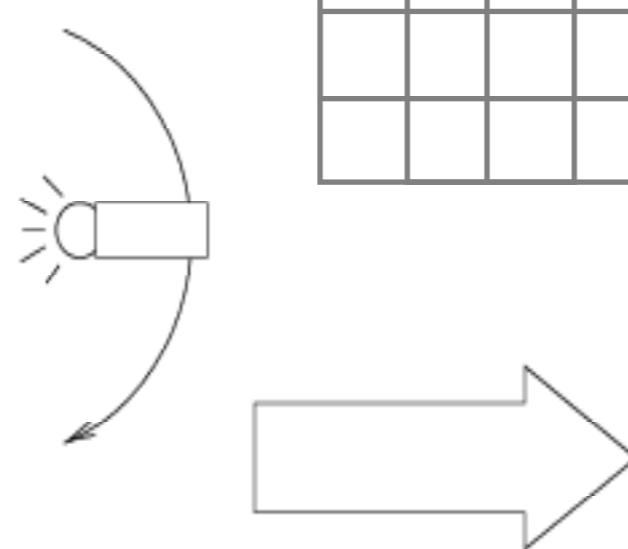
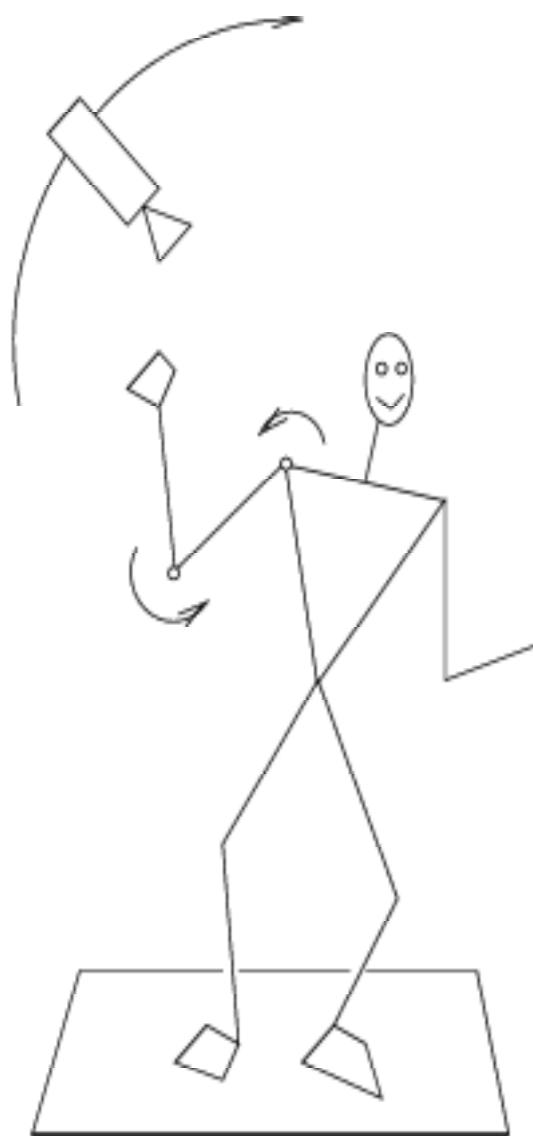
Variability:

Camera position  
Illumination  
Pose/shape parameters  
Within-class variations?



# Within-class variations





Variability:

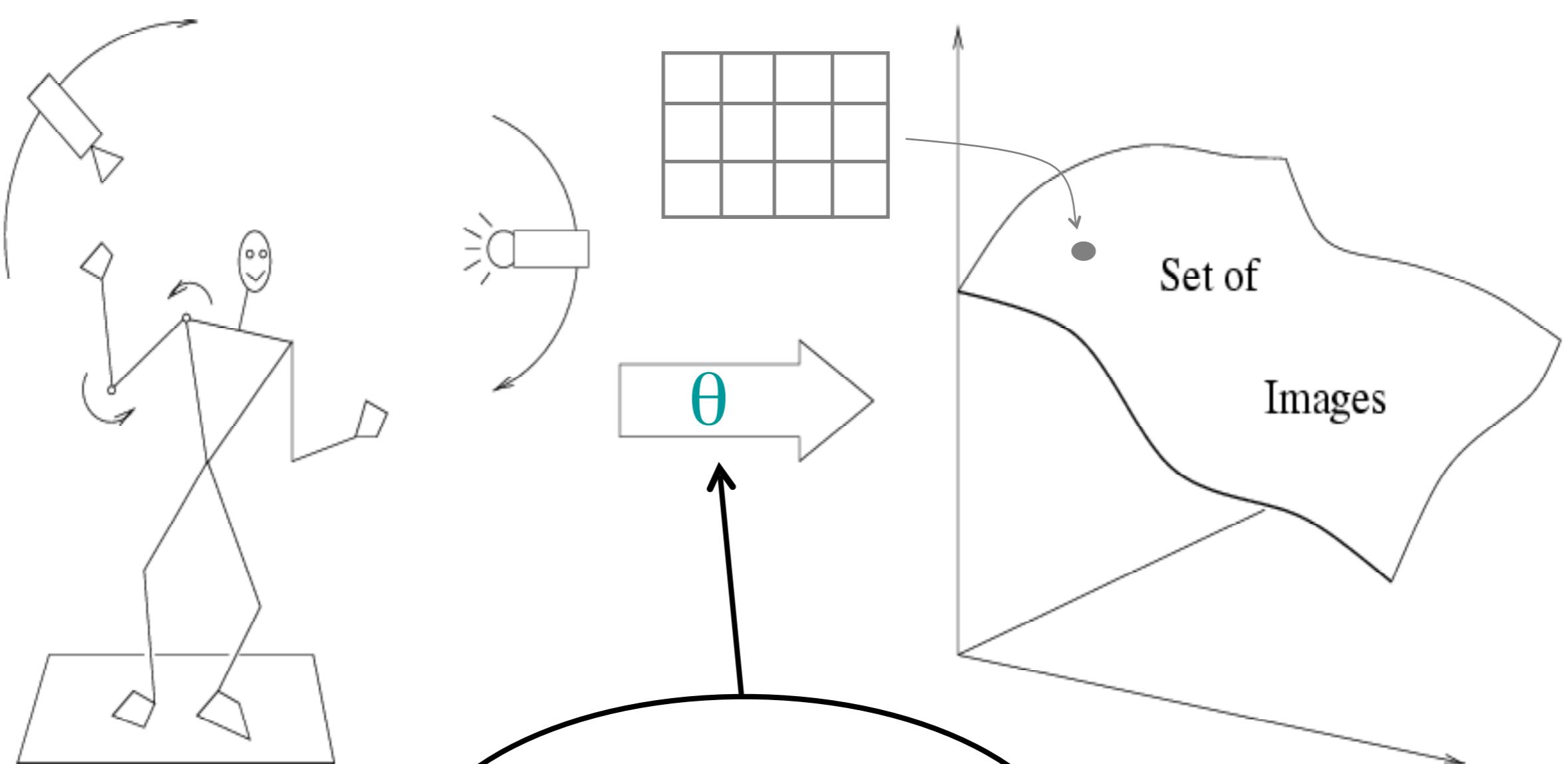
- Camera position
- Illumination
- Pose/shape parameters
- Within-class variation

# History of ideas in recognition

---

- ❖ 1960s – early 1990s: the geometric era

No digital cameras!  
Slow compute!

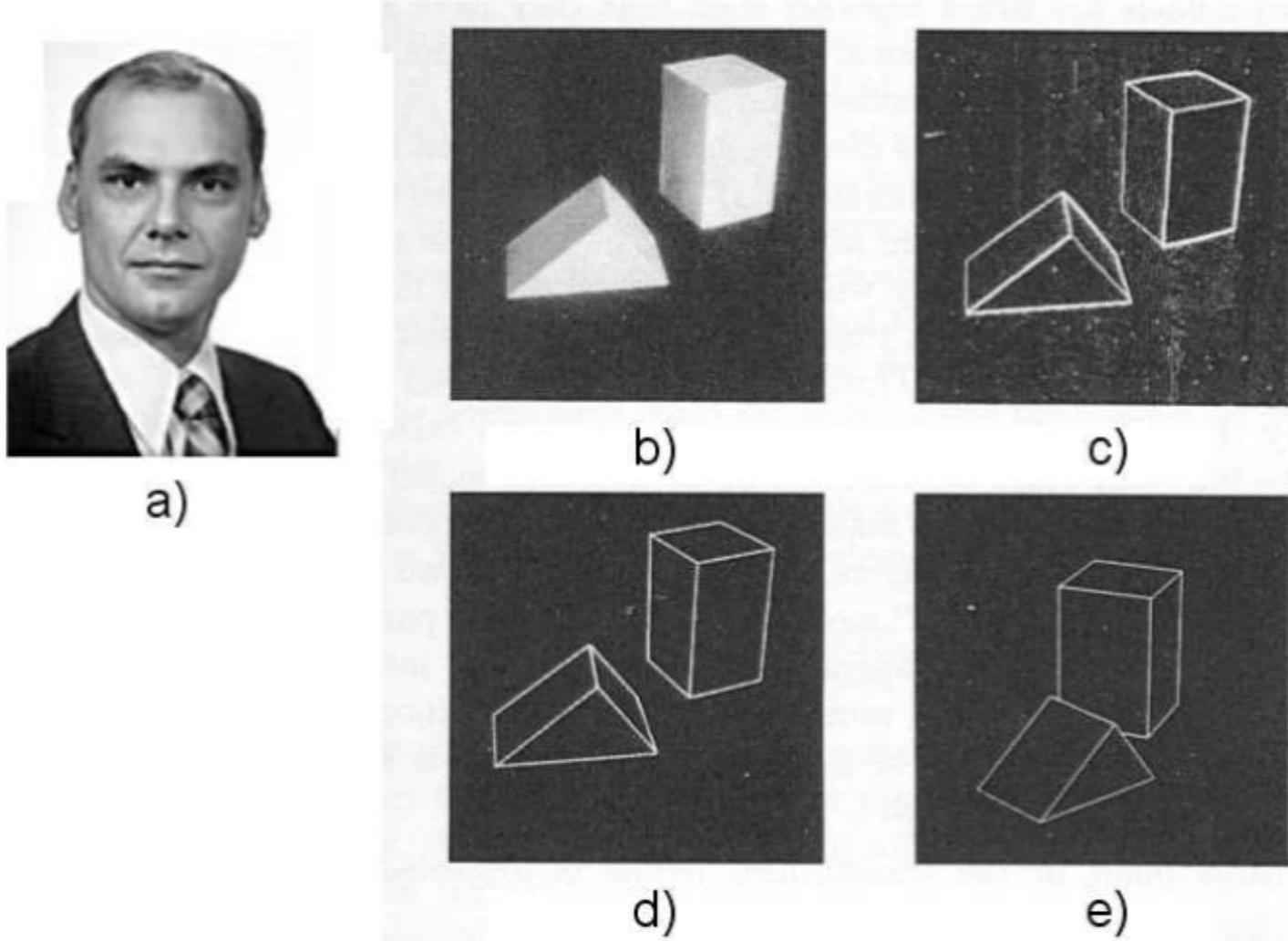


Variability:

Camera position  
Illumination

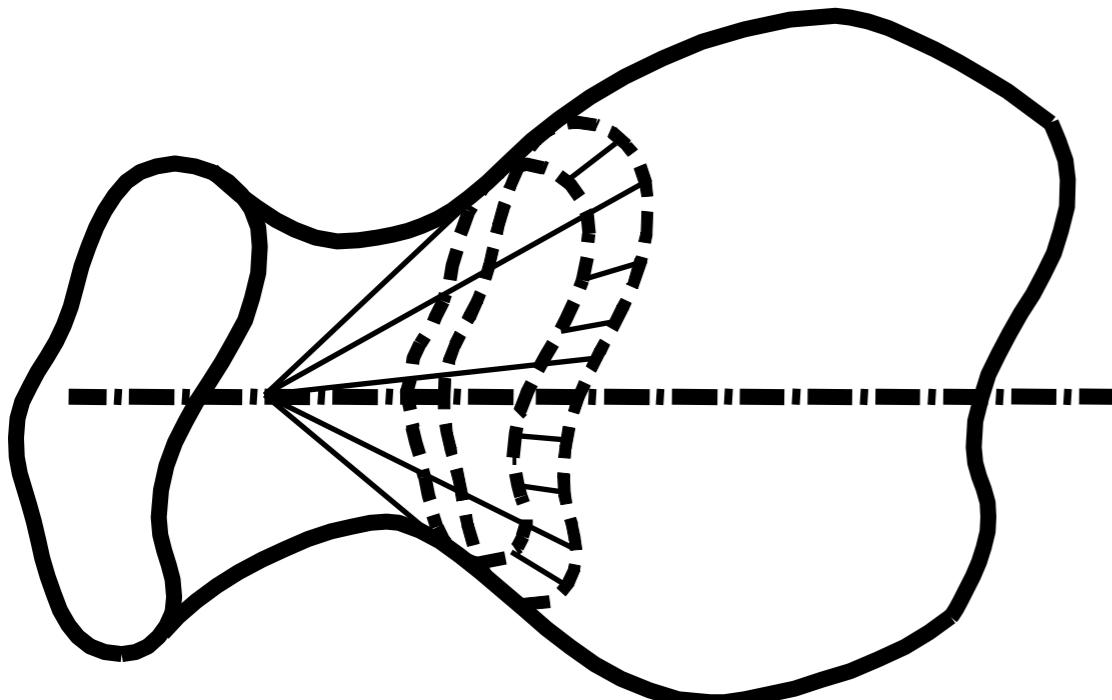
Shape is known

# Recognition as an alignment problem: Block world



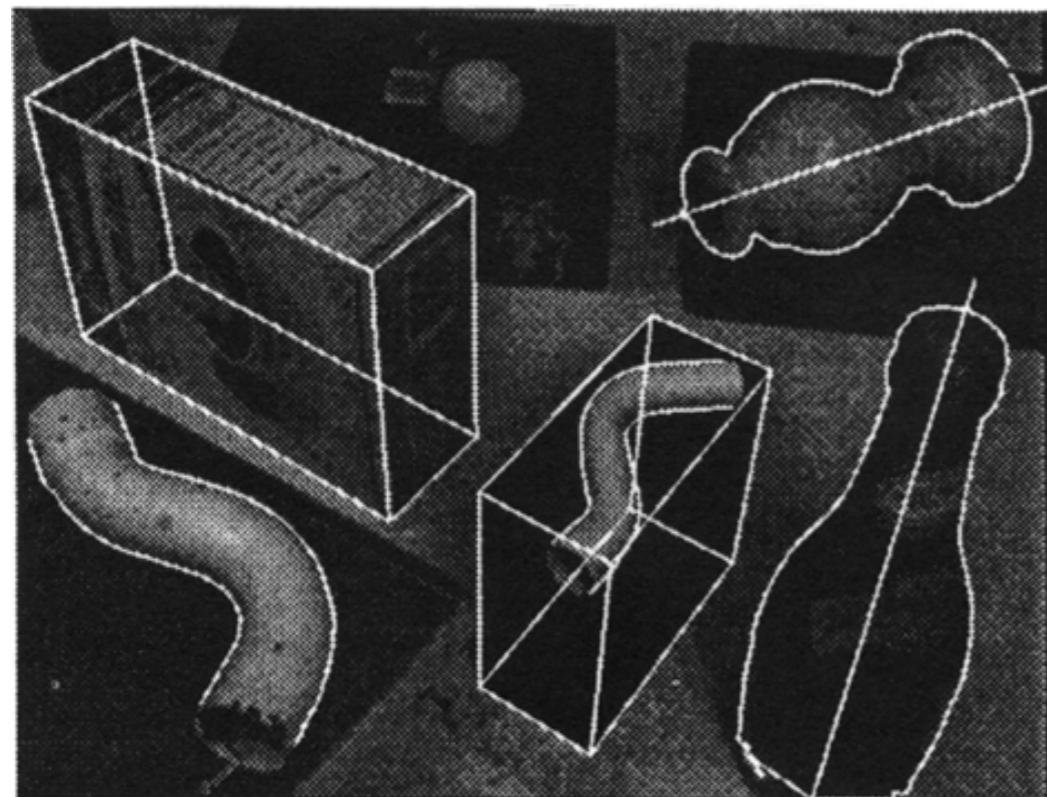
**Fig. 1.** A system for recognizing 3-d polyhedral scenes. a) L.G. Roberts. b) A blocks world scene. c) Detected edges using a 2x2 gradient operator. d) A 3-d polyhedral description of the scene, formed automatically from the single image. e) The 3-d scene displayed with a viewpoint different from the original image to demonstrate its accuracy and completeness. (b) - e) are taken from [64] with permission MIT Press.)

L. G. Roberts  
*Machine Perception of  
Three Dimensional Solids*,  
Ph.D. thesis, MIT  
Department of Electrical  
Engineering, 1963.

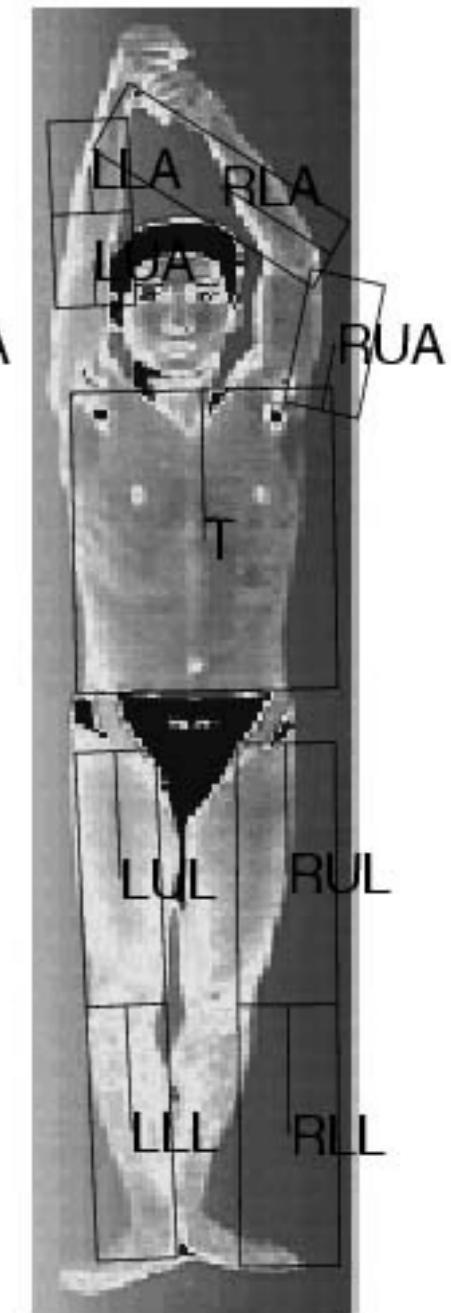
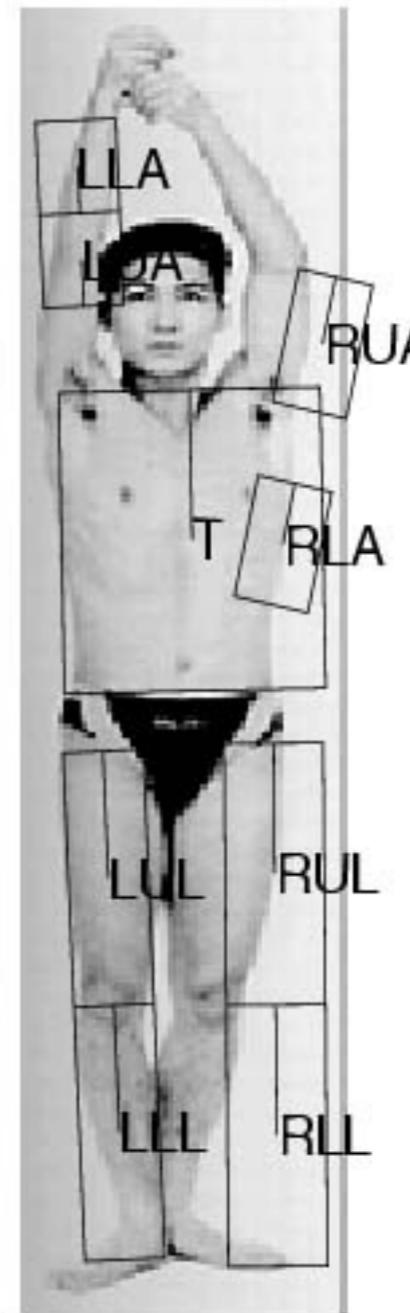


Generalized cylinders  
Ponce et al. (1989)

General shape primitives?



Zisserman et al. (1995)

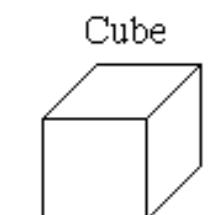


Forsyth (2000)

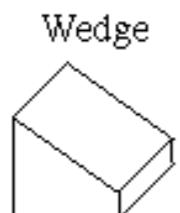
# Recognition by components

## Biederman (1987)

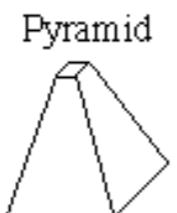
### Primitives (geons)



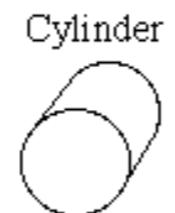
Cube



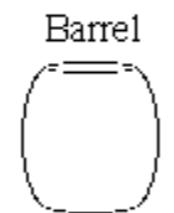
Wedge



Pyramid



Cylinder



Barrel1

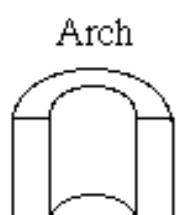
Straight Edge  
Straight Axis  
Constant

Straight Edge  
Straight Axis  
Expanded

Straight Edge  
Straight Axis  
Expanded

Curved Edge  
Straight Axis  
Constant

Curved Edge  
Straight Axis  
Exp & Cont



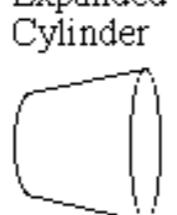
Arch

Straight Edge  
Curved Axis  
Constant



Cone

Curved Edge  
Straight Axis  
Expanded



Expanded  
Cylinder

Curved Edge  
Straight Axis  
Expanded



Handle

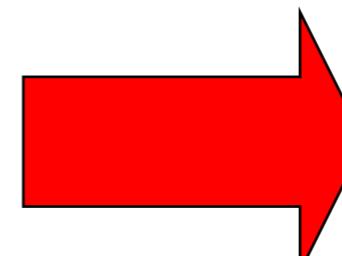
Curved Edge  
Curved Axis  
Constant



Expanded  
Handle

Curved Edge  
Curved Axis  
Expanded

### Objects



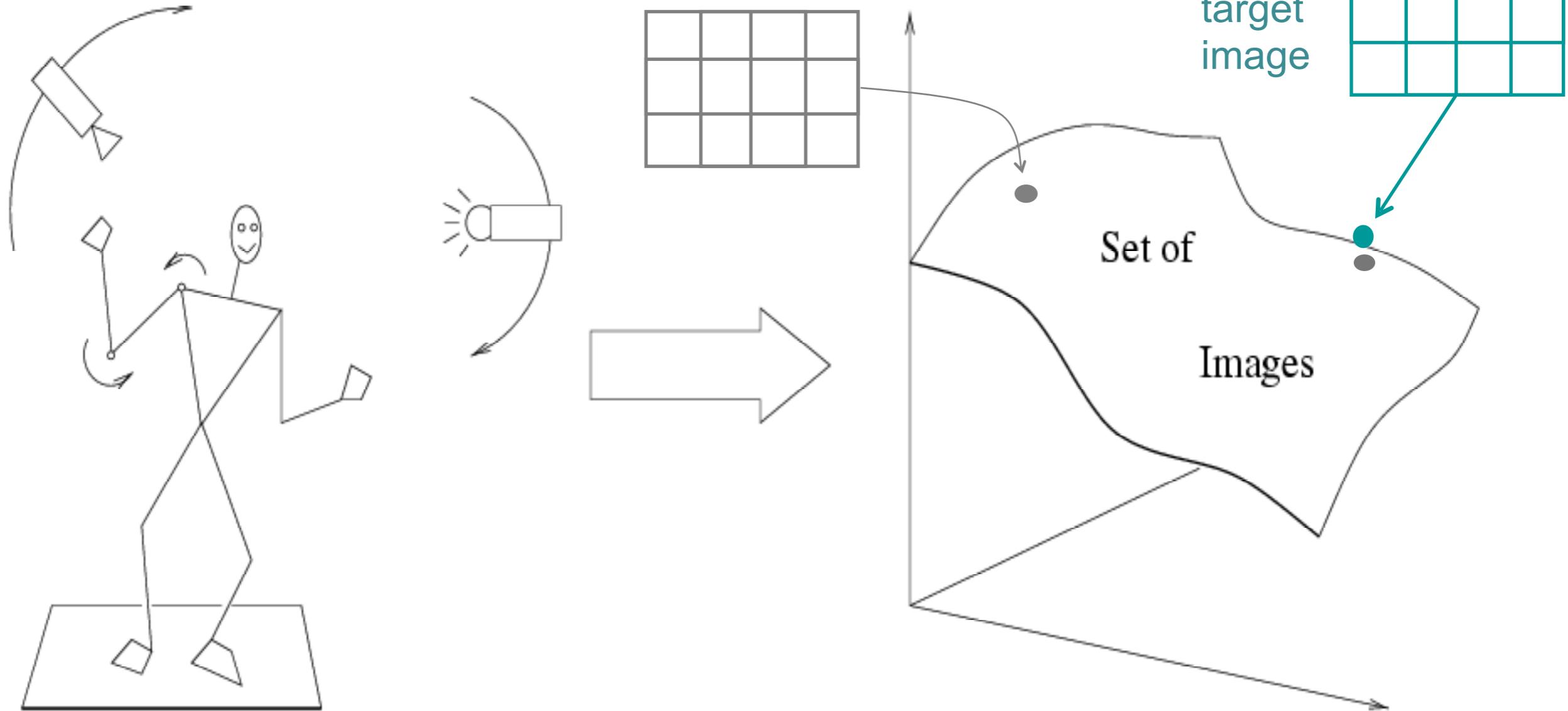
[http://en.wikipedia.org/wiki/Recognition\\_by\\_Components\\_Theory](http://en.wikipedia.org/wiki/Recognition_by_Components_Theory)

# History of ideas in recognition

- ❖ 1960s – early 1990s: the geometric era
- ❖ 1990s: appearance-based models

No digital cameras!  
Slow compute!

Slow compute!

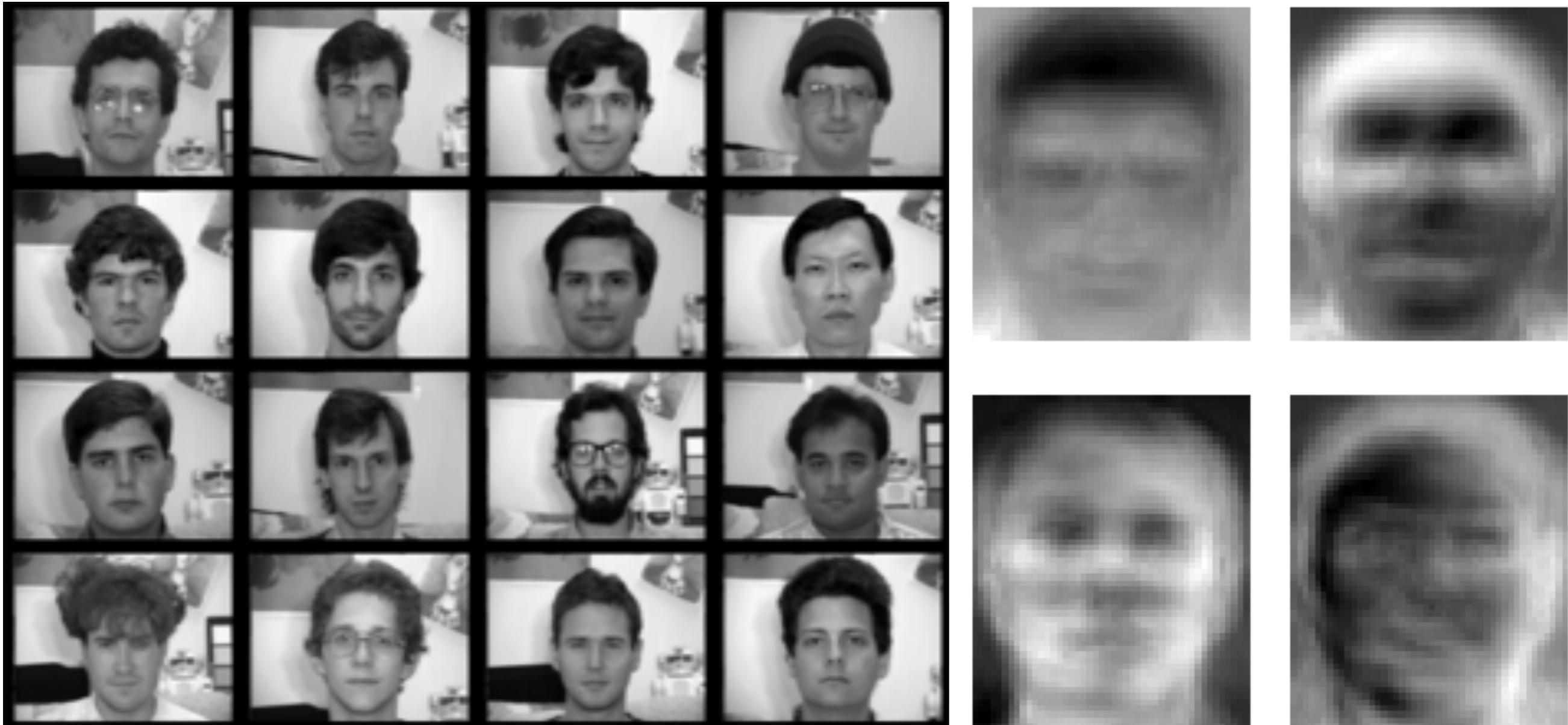


Empirical models of image variability

## Appearance-based techniques

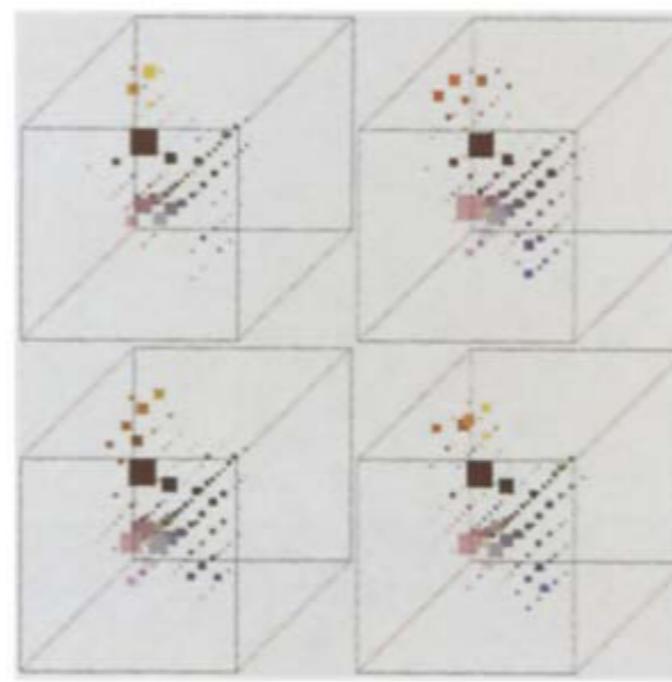
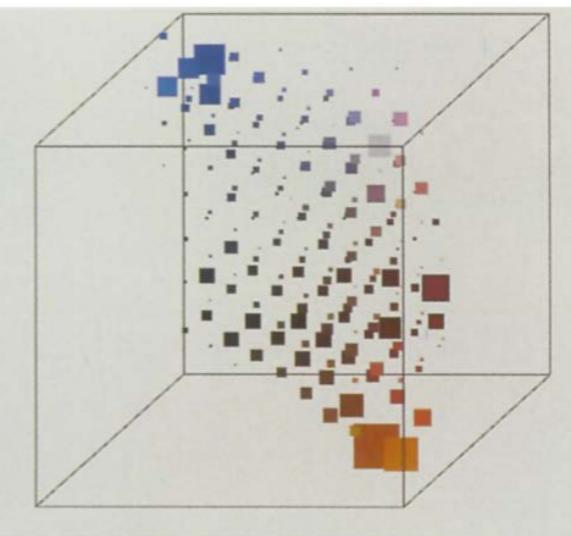
Turk & Pentland (1991); Murase & Nayar (1995); etc.

# Eigenfaces (Turk & Pentland, 1991)



Experimental Condition	Correct/Unknown Recognition Percentage		
Condition	Lighting	Orientation	Scale
Forced classification	96/0	85/0	64/0
Forced 100% accuracy	100/19	100/39	100/60
Forced 20% unknown rate	100/20	94/20	74/20

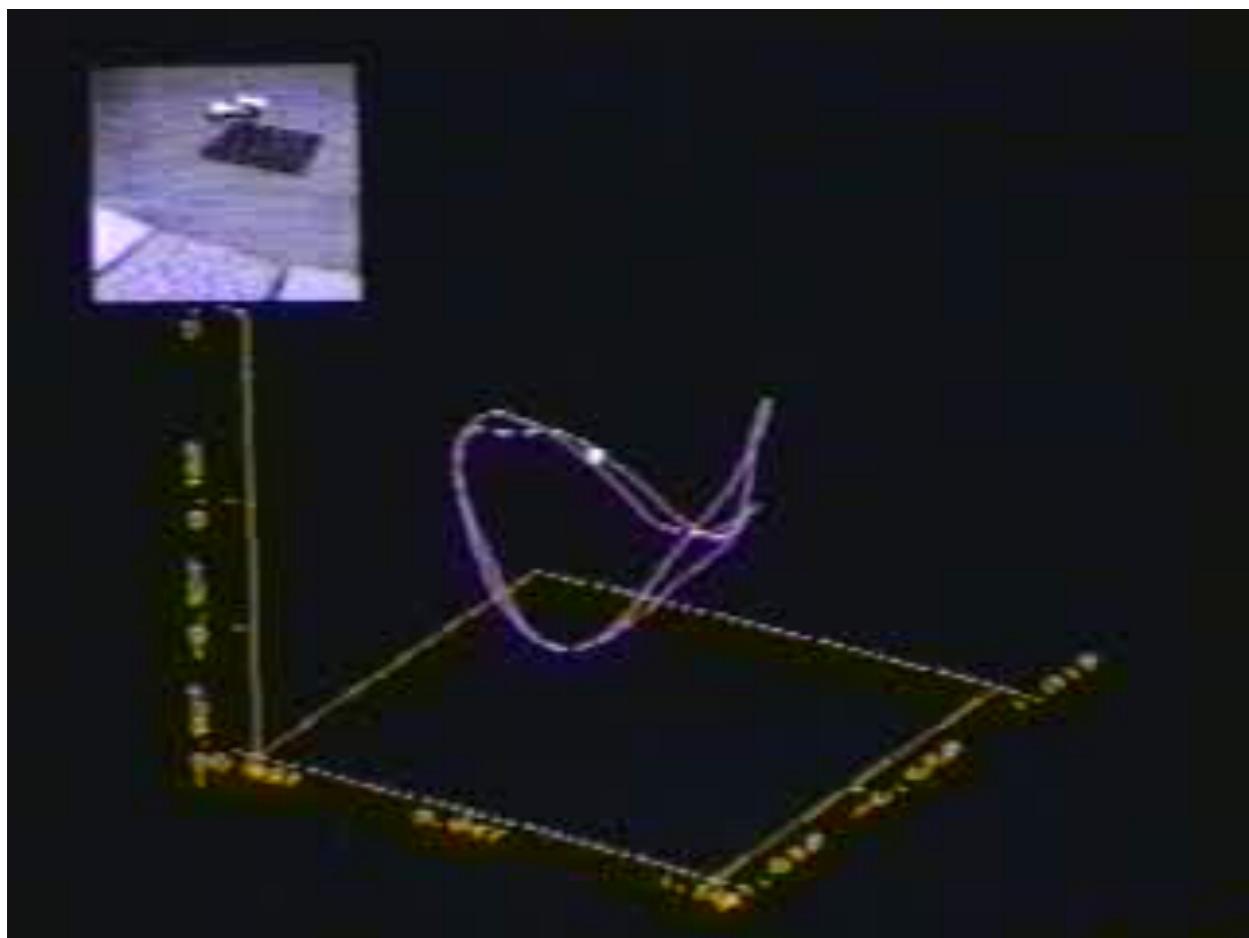
# Color Histograms



Swain and Ballard, [Color Indexing](#), IJCV 1991.

Svetlana Lazebnik

# Appearance manifolds



H. Murase and S. Nayar, Visual learning and recognition of 3-d objects from appearance,  
IJCV 1995

# Limitations of global appearance models

- ❖ Requires global registration of patterns
- ❖ Not robust to clutter, occlusion, geometric transformations



# History of ideas in recognition

- ❖ 1960s – early 1990s: the geometric era No digital cameras!  
Slow compute!
- ❖ 1990s: appearance-based models Slow compute!
- ❖ 1990s – present: sliding window approaches

# Sliding window approaches



# Sliding window approaches



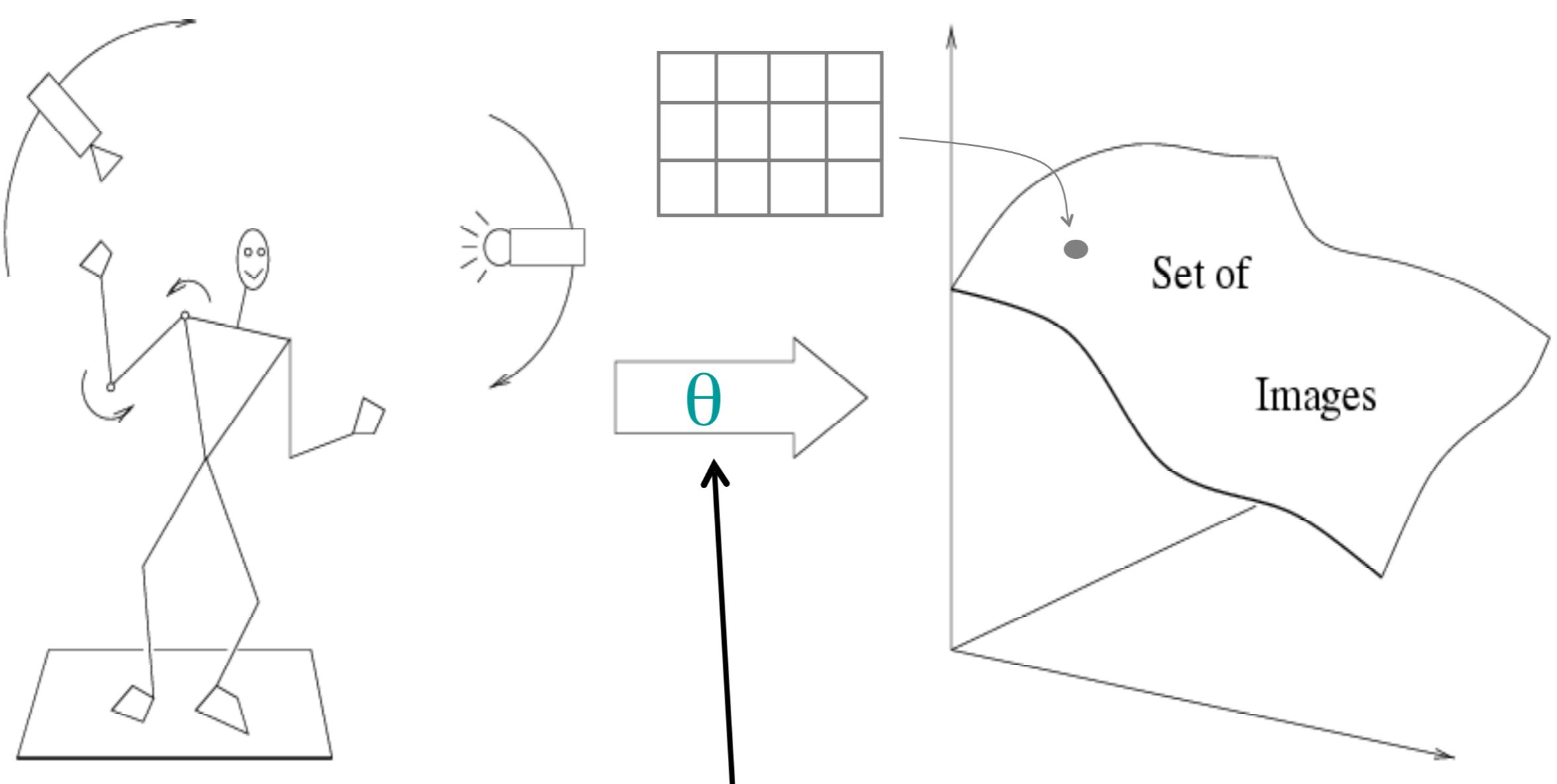
- ❖ Turk and Pentland, 1991
- ❖ Belhumeur, Hespanha, & Kriegman, 1997
- ❖ Schneiderman & Kanade 2004
- ❖ Viola and Jones, 2000



- ❖ Schneiderman & Kanade, 2004
- ❖ Argawal and Roth, 2002
- ❖ Poggio et al. 1993

# History of ideas in recognition

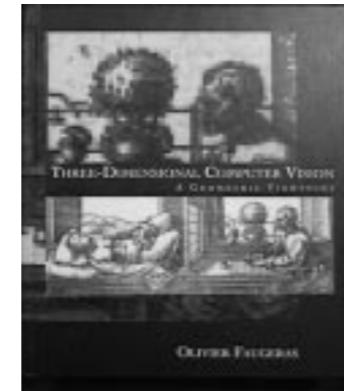
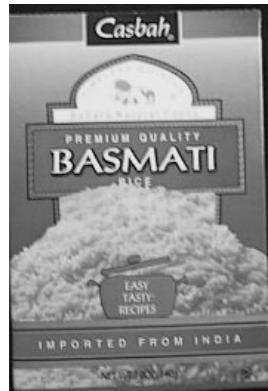
- ❖ 1960s – early 1990s: the geometric era
  - No digital cameras!  
Slow compute!
- ❖ 1990s: appearance-based models
  - Slow compute!
- ❖ Mid-1990s: sliding window approaches
- ❖ Late 1990s: local features



Variability:

Camera position  
Illumination  
Shape is partially known

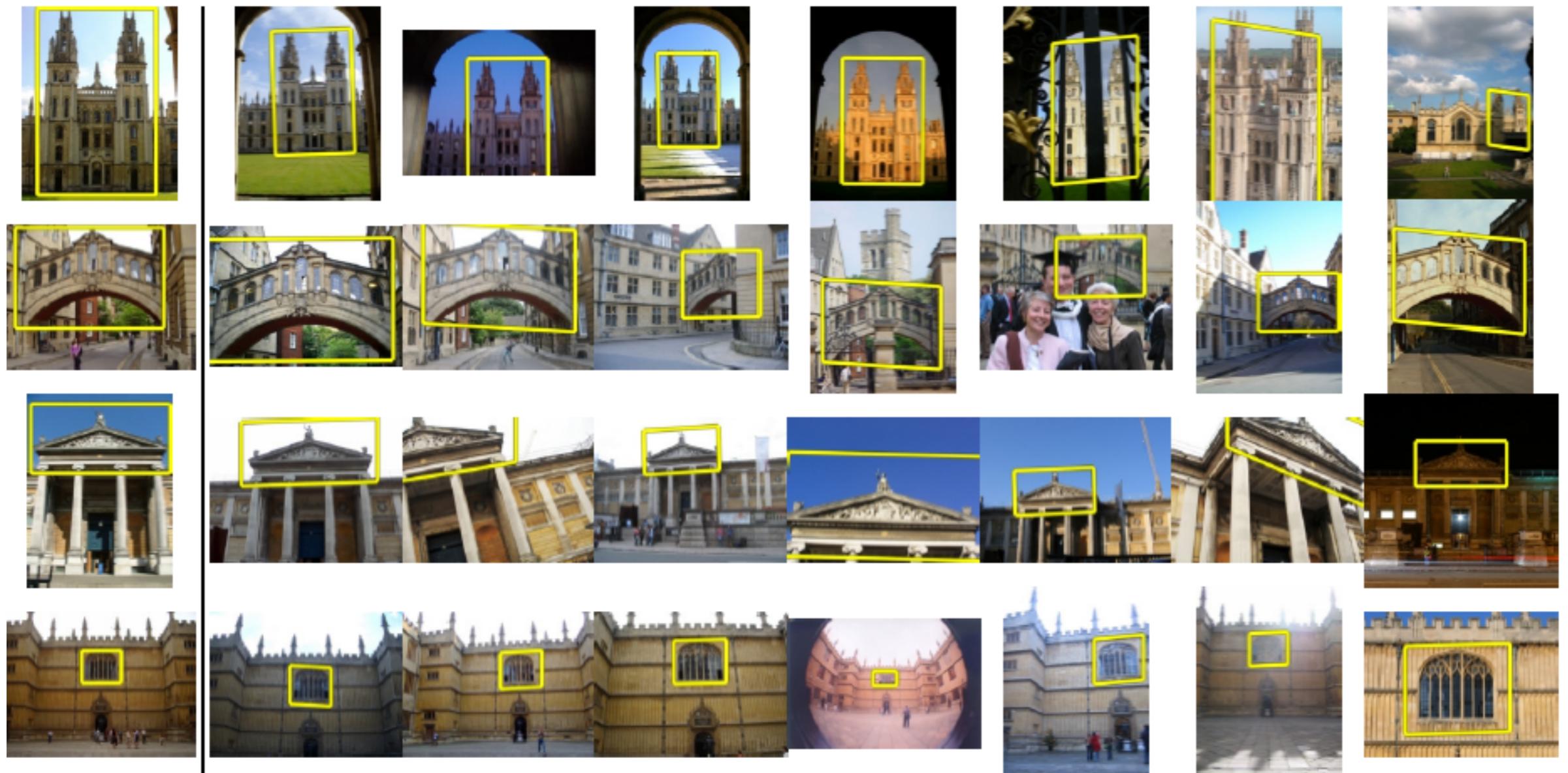
# Local features for object instance recognition



D. Lowe (1999, 2004)

# Large-scale image search

## Combining local features, indexing, and spatial constraints



Philbin et al. ‘07

# Large-scale image search

## Combining local features, indexing, and spatial constraints

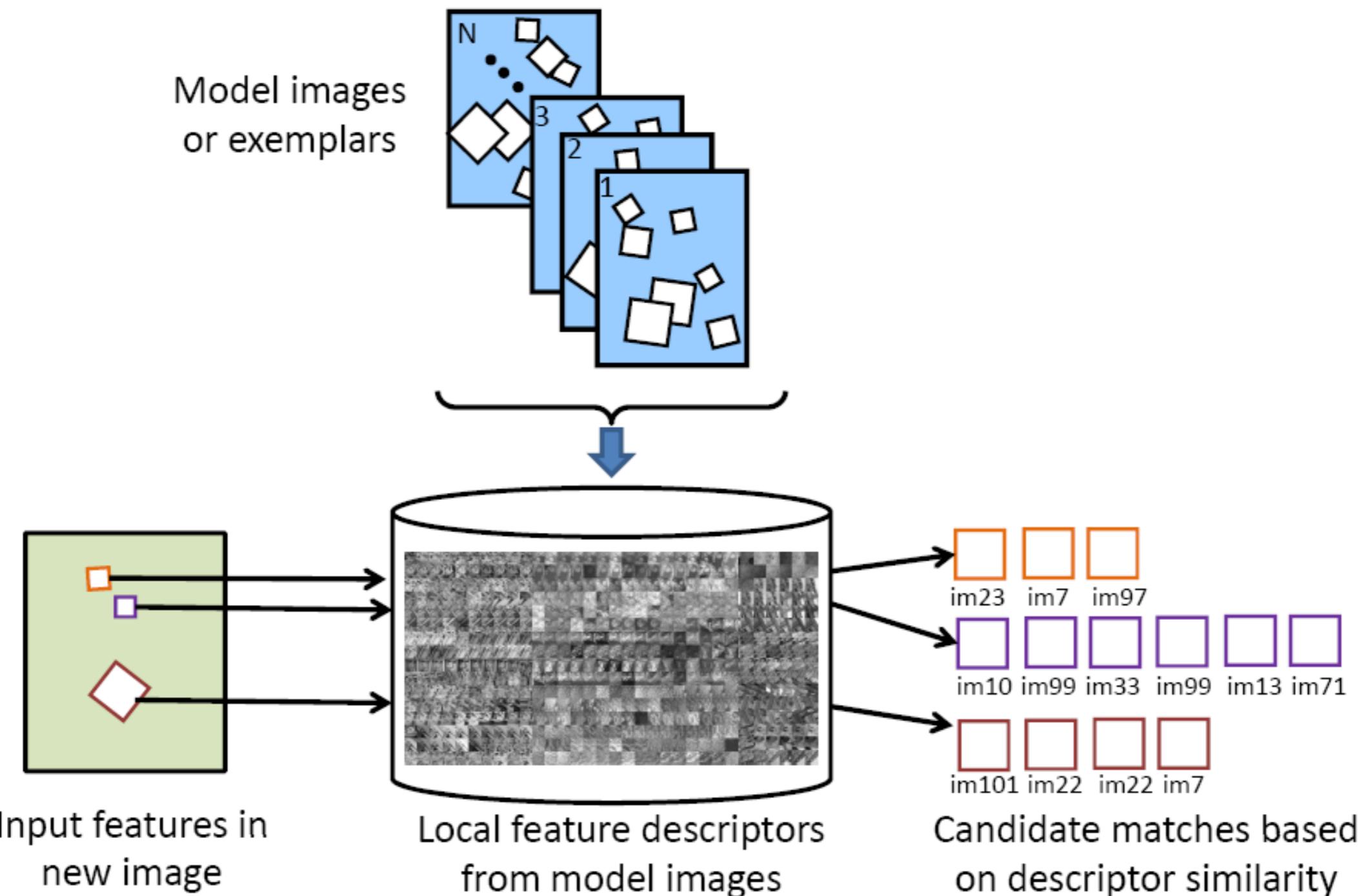


Image credit: K. Grauman and B. Leibe

---

# History of ideas in recognition

---

- ❖ 1960s – early 1990s: the geometric era
- ❖ 1990s: appearance-based models
- ❖ Mid-1990s: sliding window approaches
- ❖ Late 1990s: local features
- ❖ Early 2000s: parts-and-shape models

# Parts-and-shape models

- ❖ Model:
  - ❖ Object as a set of parts
  - ❖ Relative locations between parts
  - ❖ Appearance of part

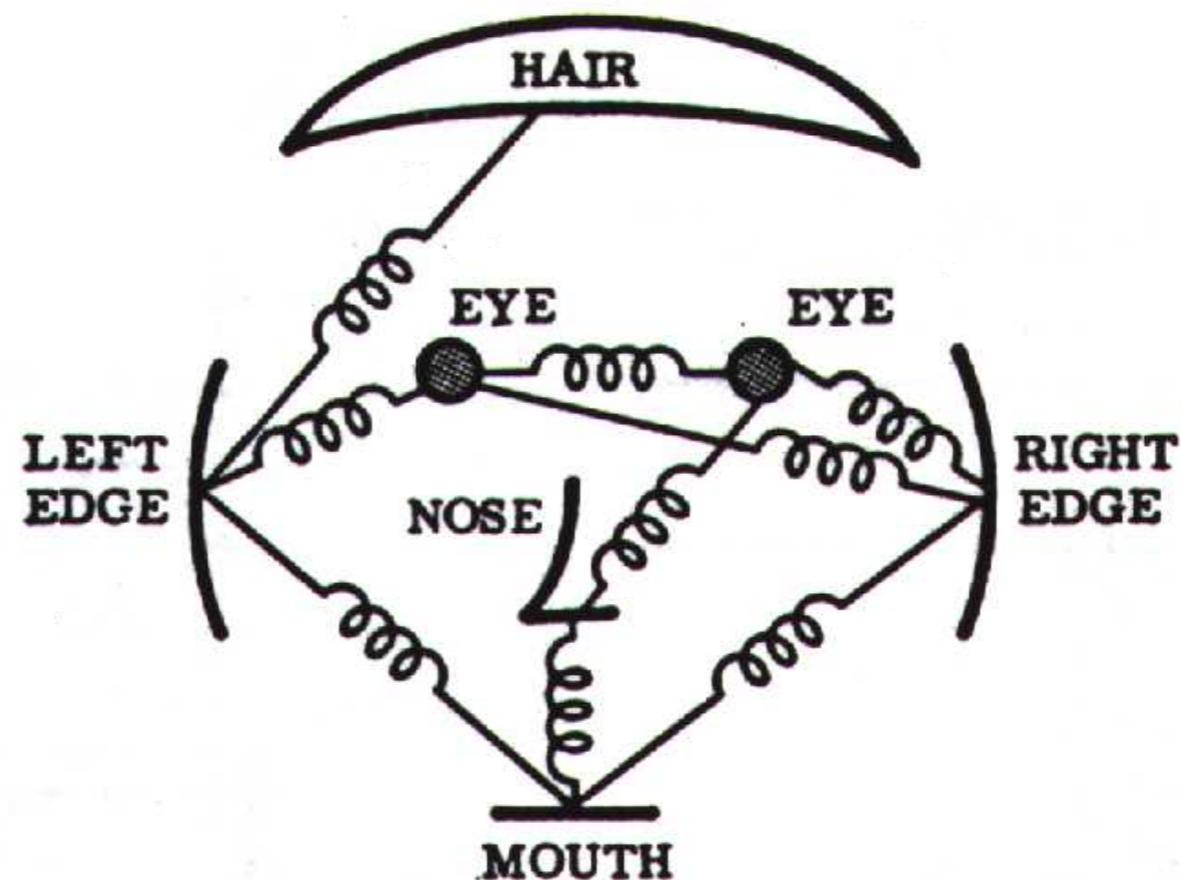
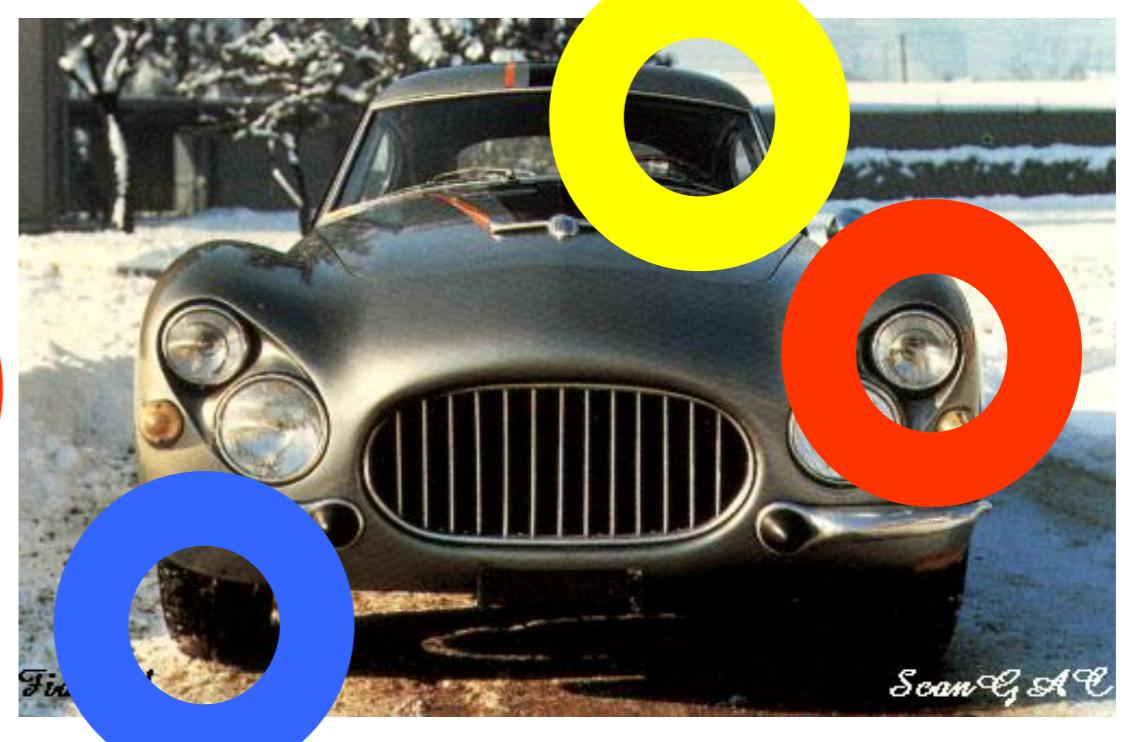


Figure from [Fischler & Elschlager 73]

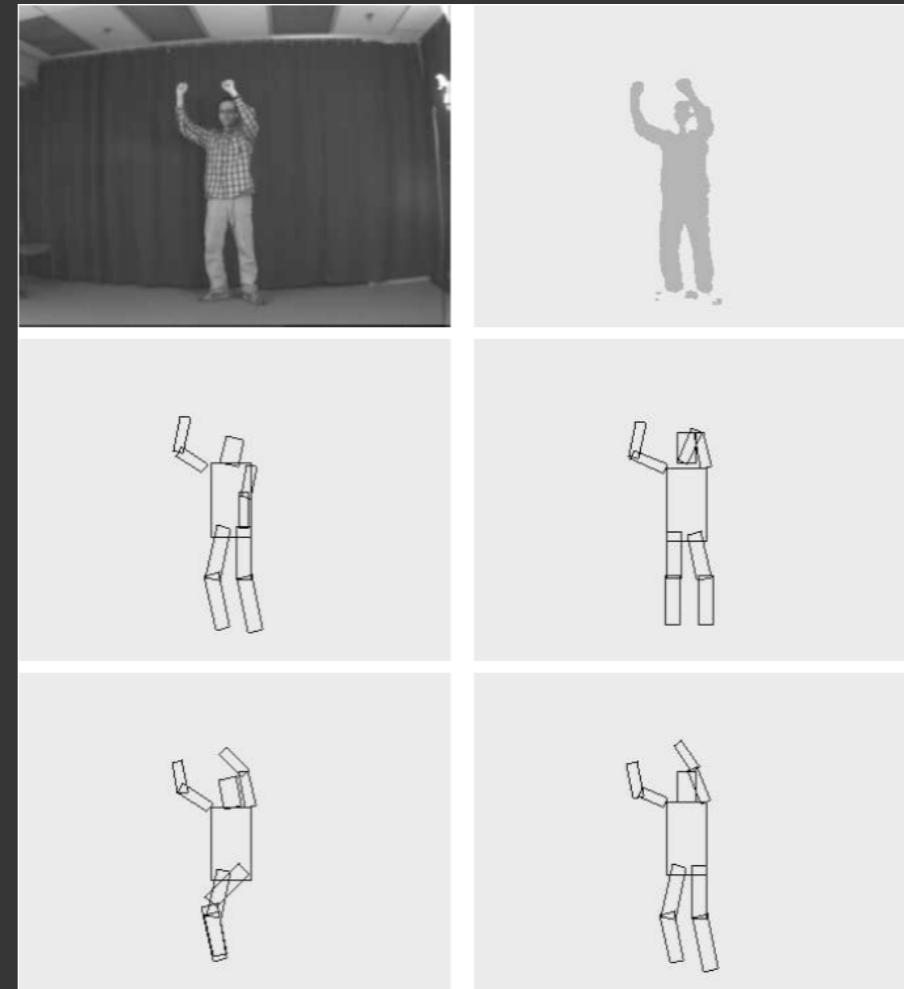
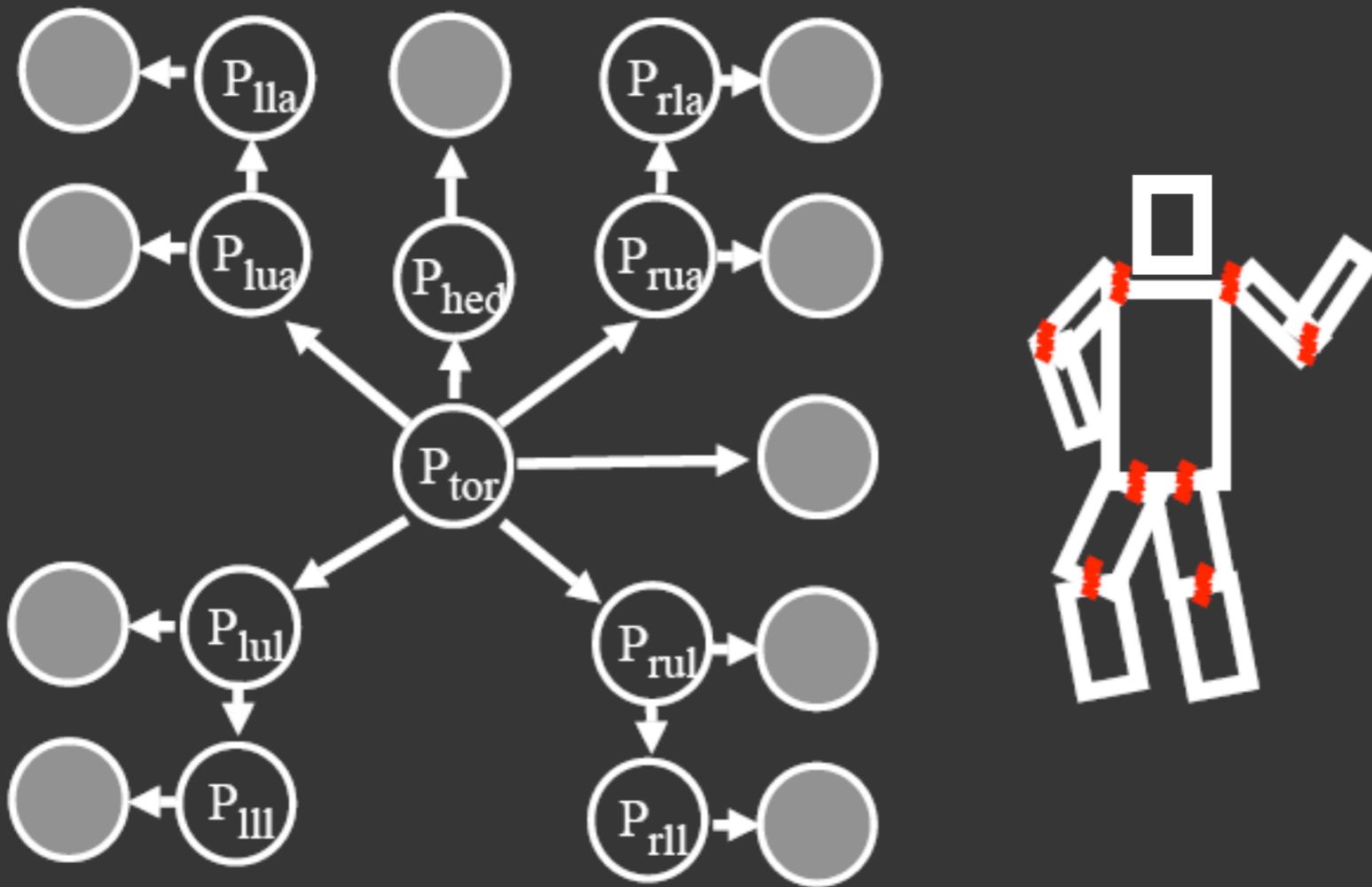
# Constellation models



Weber, Welling & Perona (2000), Fergus, Perona & Zisserman (2003)

# Pictorial structure model

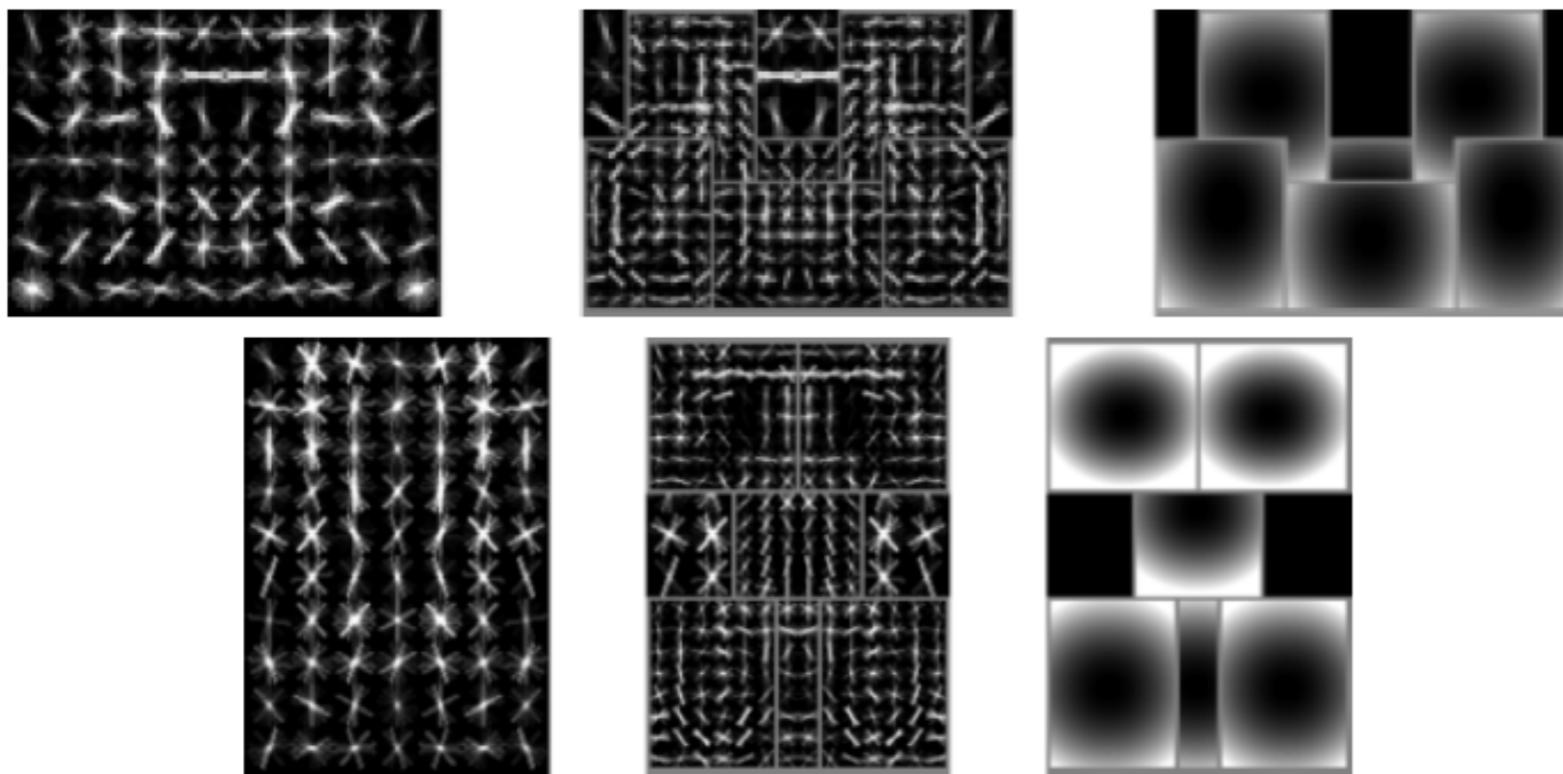
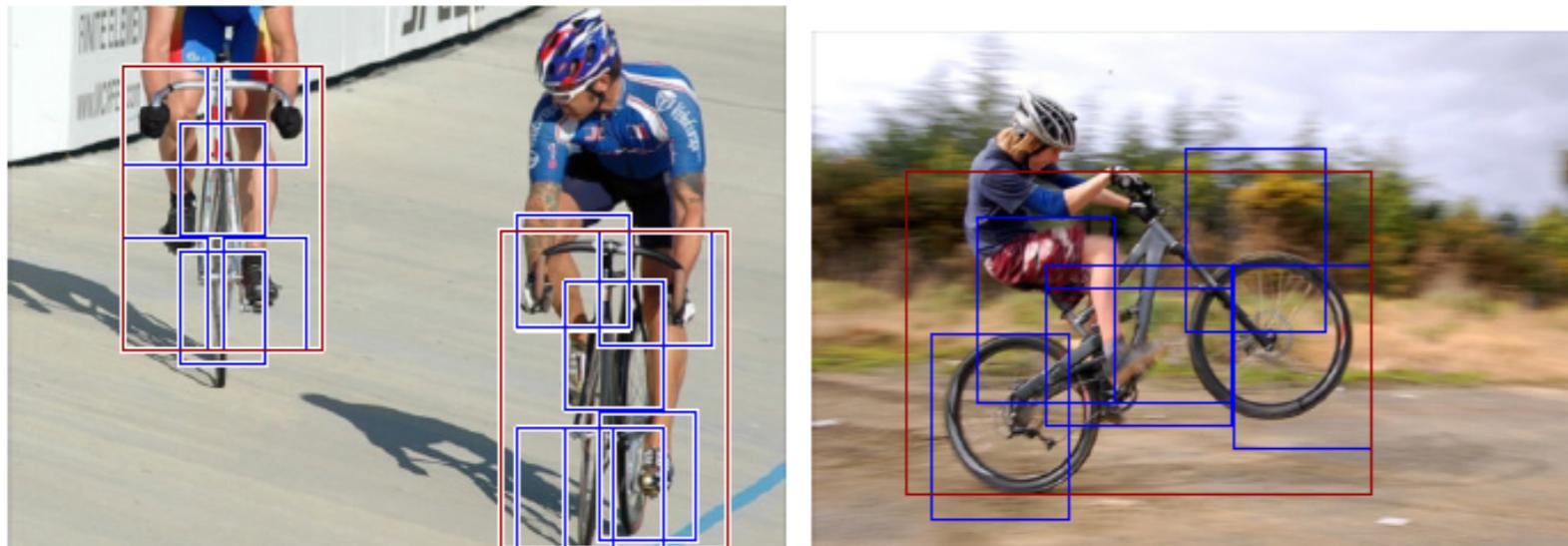
Fischler and Elschlager(73), Felzenszwalb and Huttenlocher(00)



$$\Pr(P_{\text{tor}}, P_{\text{arm}}, \dots | \text{Im}) \propto \prod_{i,j} \Pr(P_i | P_j) \prod_i \Pr(\text{Im}(P_i))$$

↑  
part geometry      ↙  
part appearance

# Discriminatively trained part-based models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, PAMI 2009,

**“Object Detection with Discriminatively Trained Part-Based Models”**

# History of ideas in recognition

- ❖ 1960s – early 1990s: the geometric era No digital cameras!  
Slow compute!
- ❖ 1990s: appearance-based models Slow compute!
- ❖ Mid-1990s: sliding window approaches
- ❖ Late 1990s: local features
- ❖ Early 2000s: parts-and-shape models
- ❖ Mid-2000s: bags of features Early GPU compute.

# History of ideas in recognition

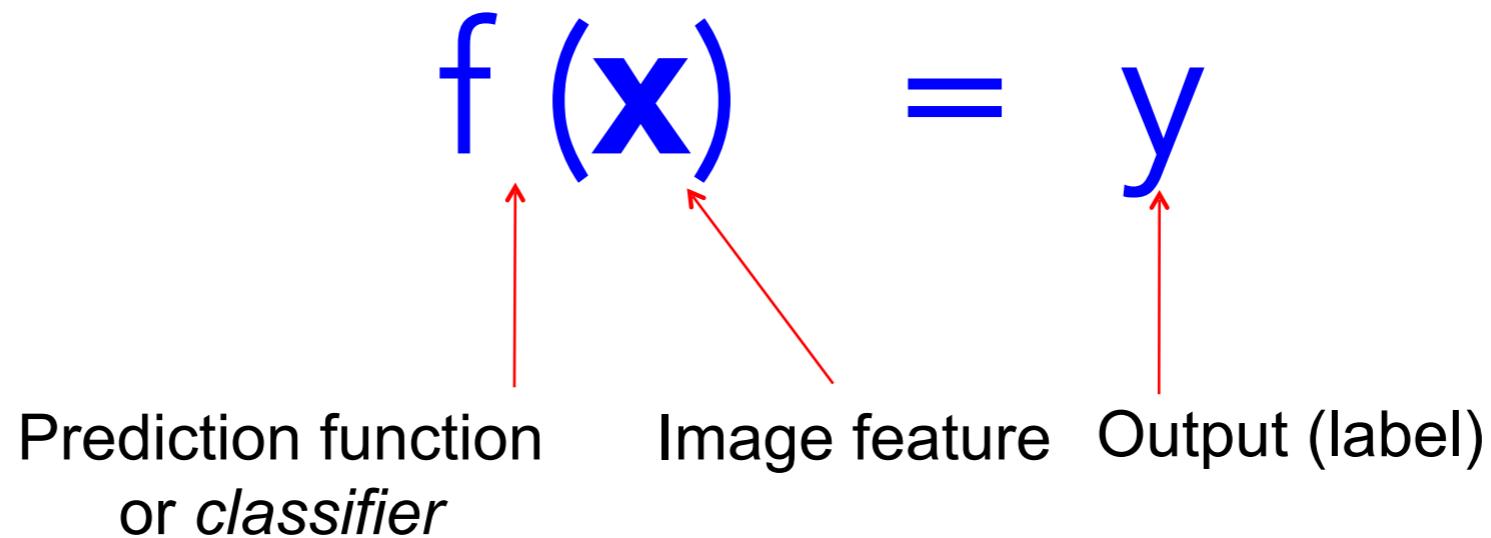
- ❖ 1960s – early 1990s: the geometric era No digital cameras!  
Slow compute!
- ❖ 1990s: appearance-based models Slow compute!
- ❖ Mid-1990s: sliding window approaches
- ❖ Late 1990s: local features
- ❖ Early 2000s: parts-and-shape models
- ❖ Mid-2000s: bags of features (next!) Early GPU compute.
- ❖ *Present trends:*  
Combined local and global methods,  
context, deep learning GPU/cloud compute.

# The machine learning framework

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

$$f(\text{apple}) = \text{"apple"}$$
$$f(\text{tomato}) = \text{"tomato"}$$
$$f(\text{cow}) = \text{"cow"}$$

# The machine learning framework



- ❖ **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 unseen *test example*  $\mathbf{x}_u$  and output the predicted value  $y_u = f(\mathbf{x}_u)$  to *classify*  $\mathbf{x}_u$ .

# Dataset

Training  
Images



Validation  
Images



Testing  
Images



- Train classifier

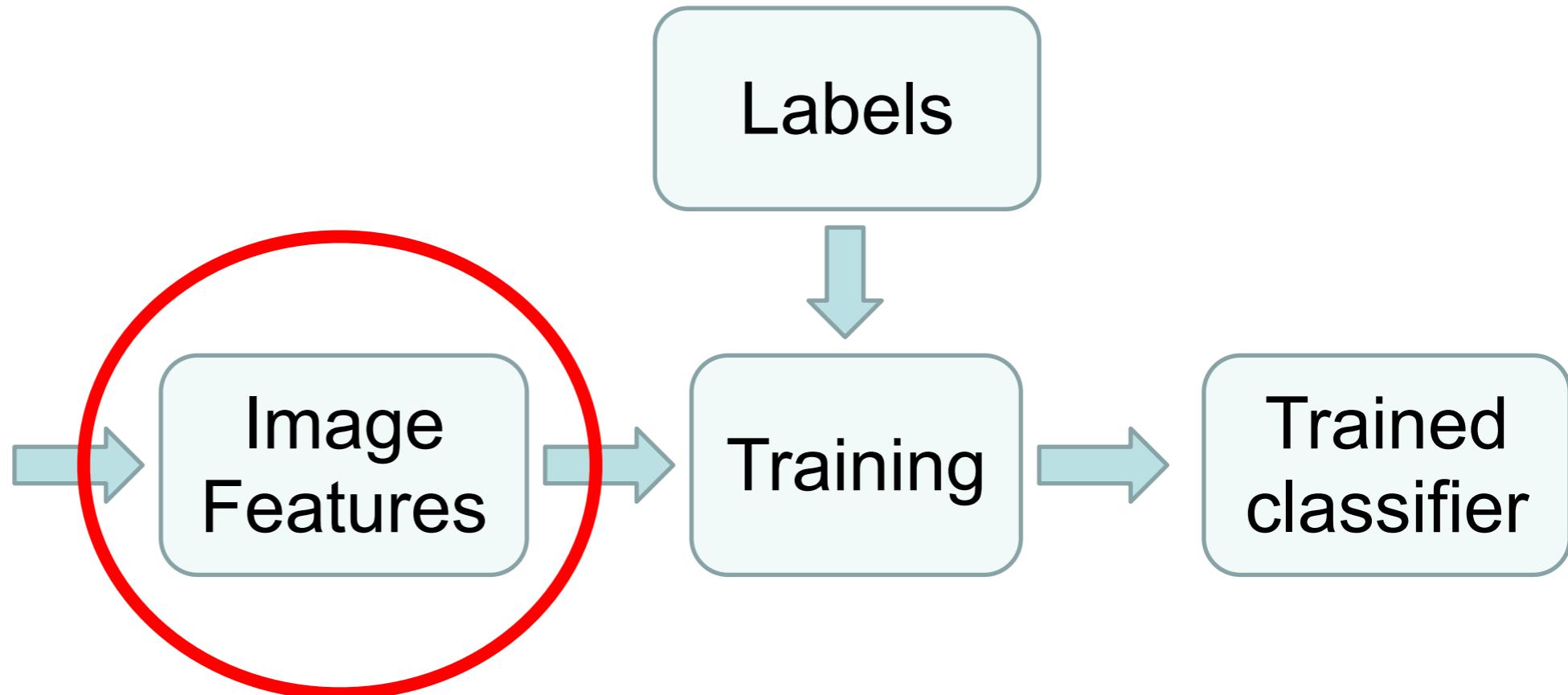
- Measure error
- Tune model hyperparameters

- Secret labels
- Measure error

*Random train/validate splits = cross validation*

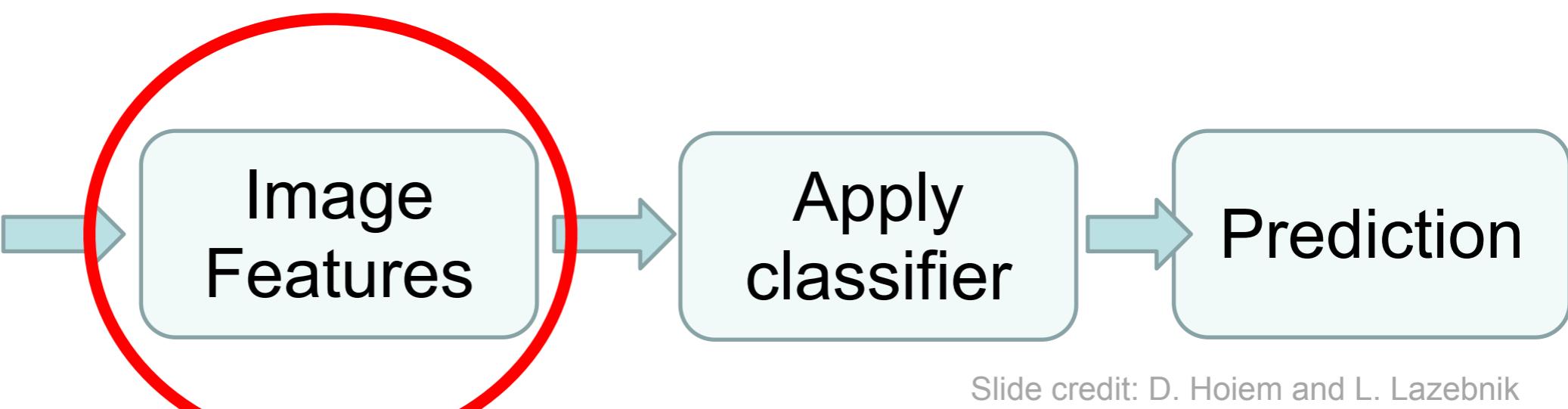
# Training

Images



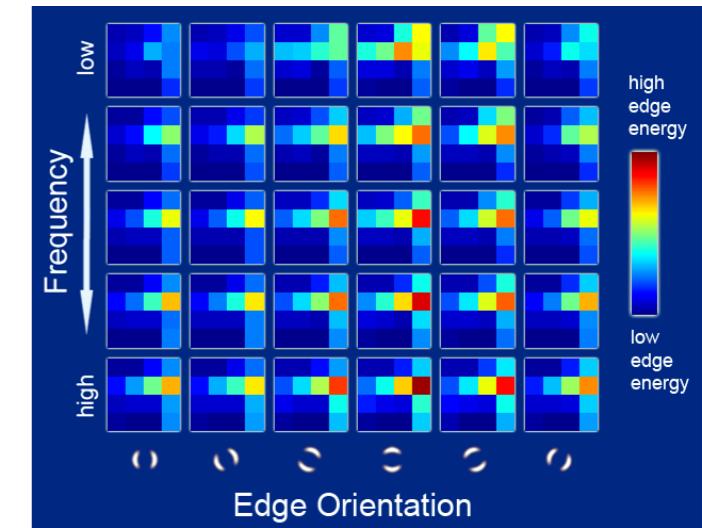
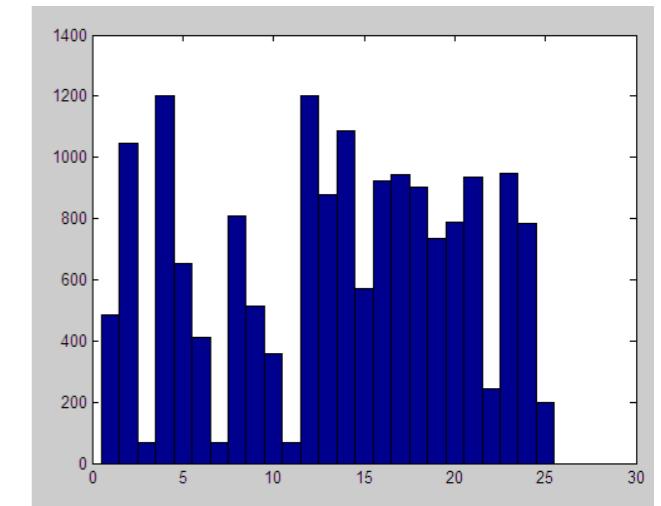
# Testing

Image  
*not in*  
training set

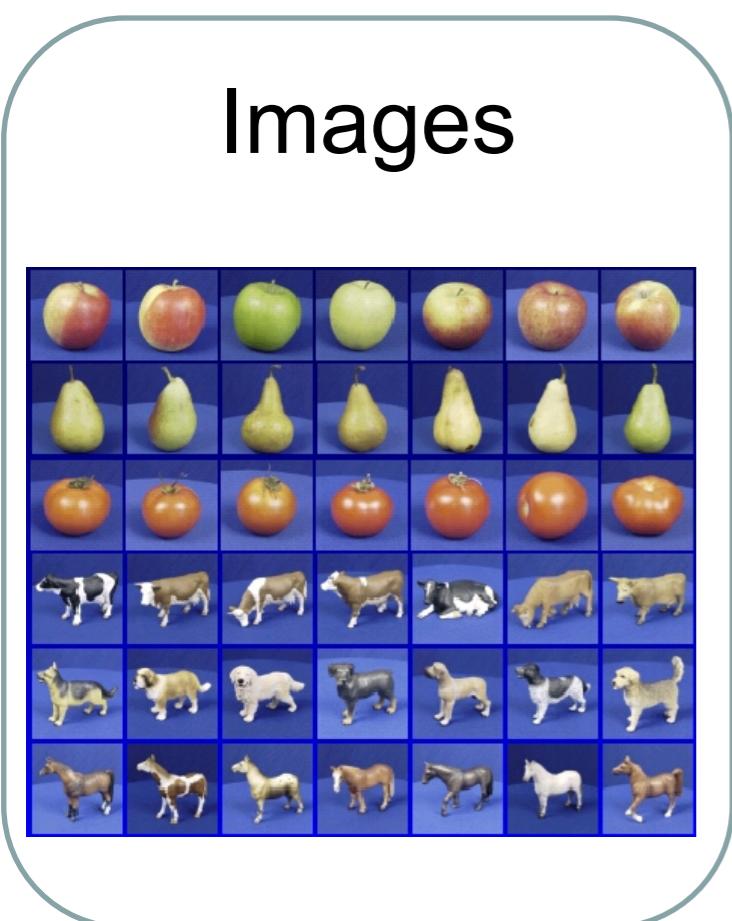


# Features

- ❖ Raw pixels
- ❖ Histograms
- ❖ Templates
- ❖ SIFT descriptors
  - ❖ GIST
  - ❖ ORB
  - ❖ HOG....



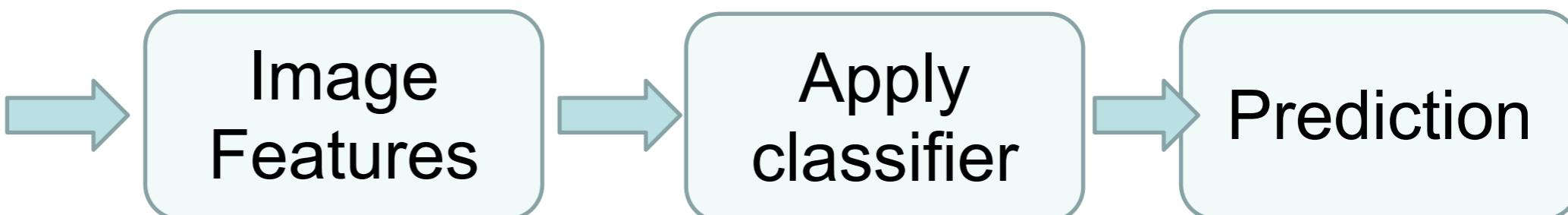
# Training



# Testing

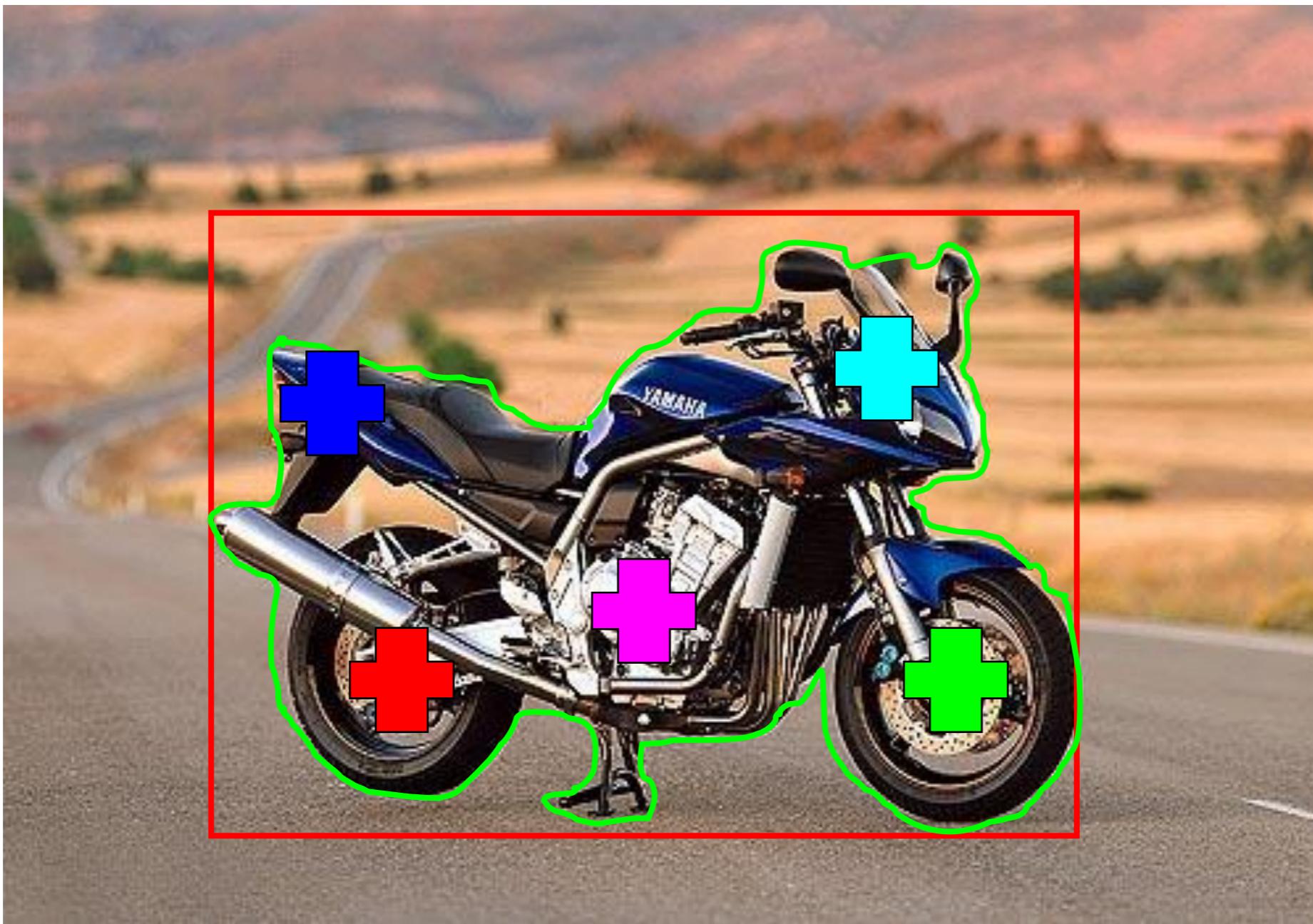
Image  
*not in*  
training set

A single image of an apple.



# Recognition task and supervision

- ❖ Images in the training set must be annotated with the “correct answer” that the model is expected to produce



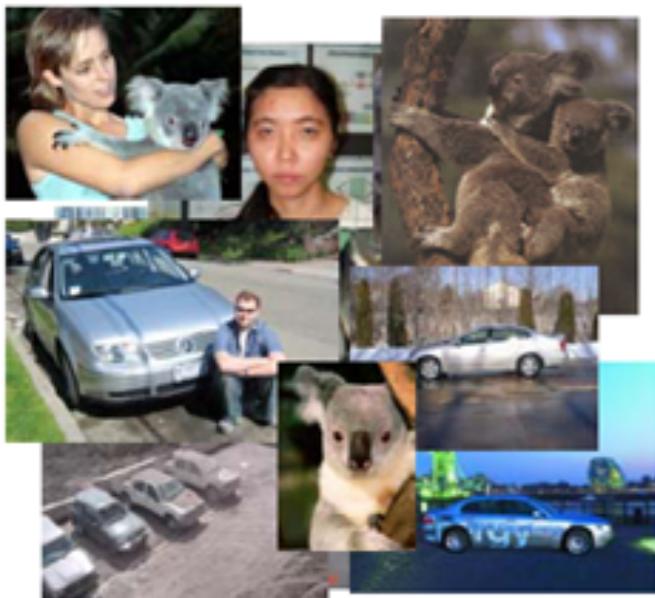
# Spectrum of supervision

Less

More



E.G., ImageNet

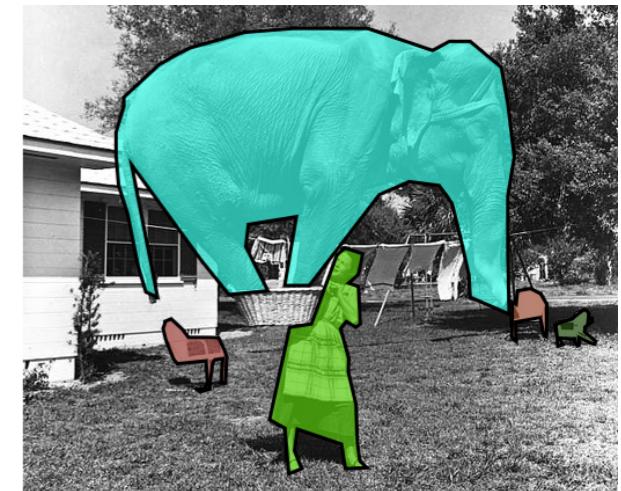


Unsupervised



“Weakly” supervised

E.G., MS Coco

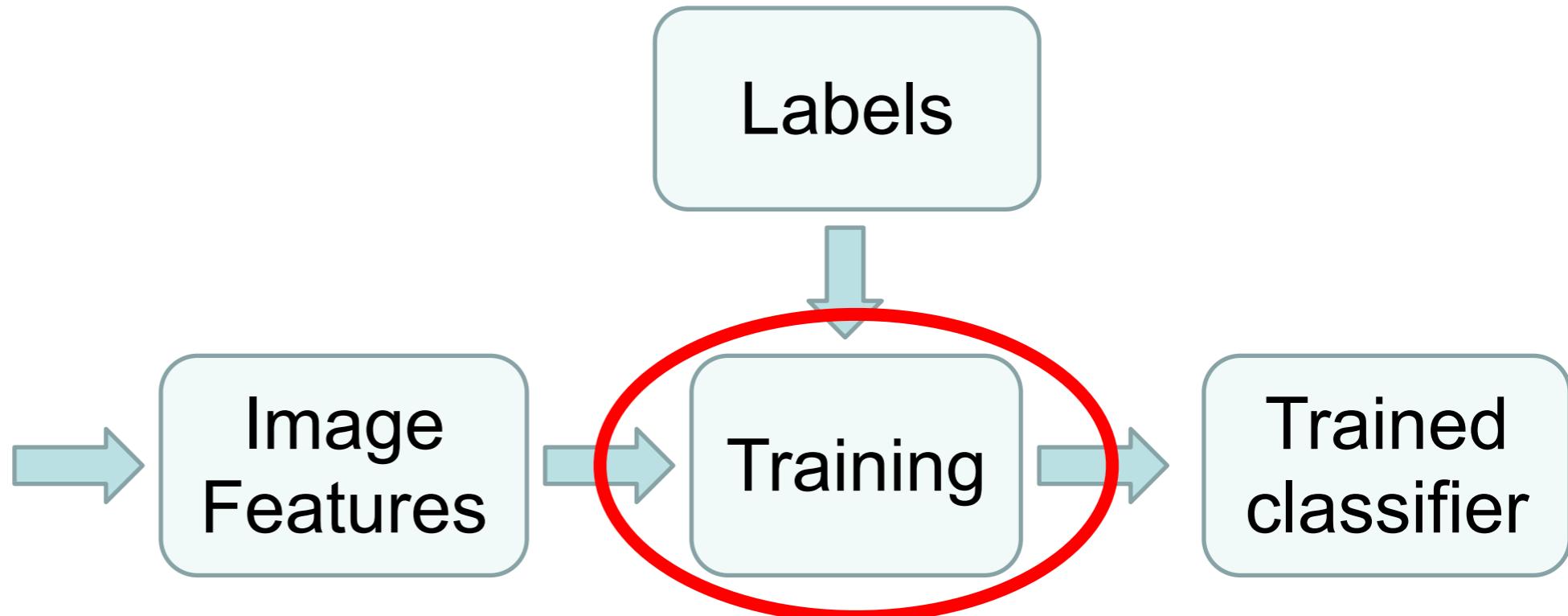
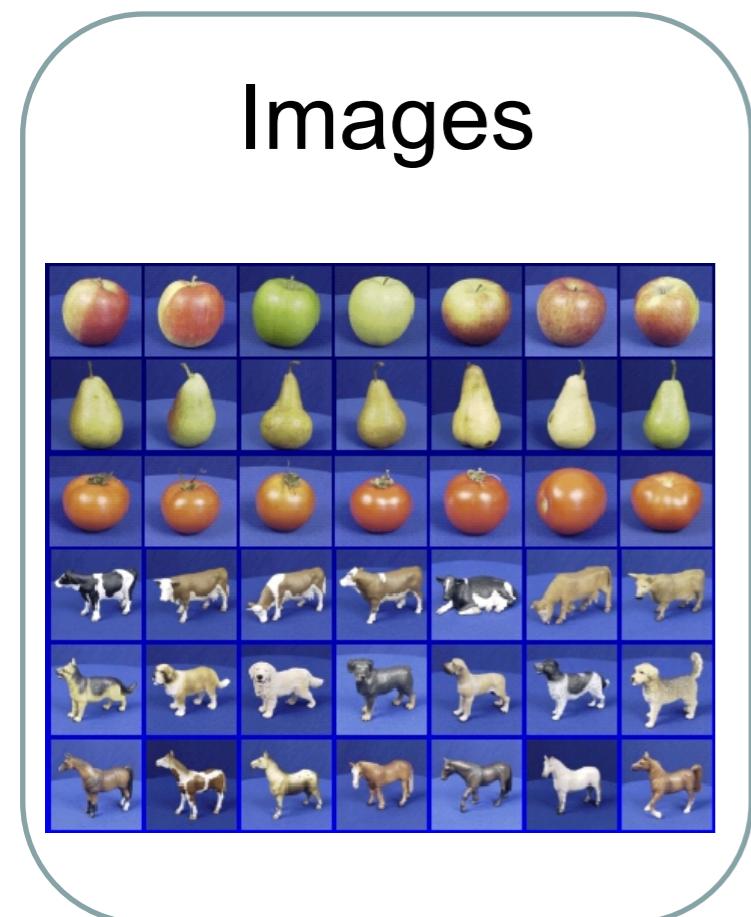


Fully supervised

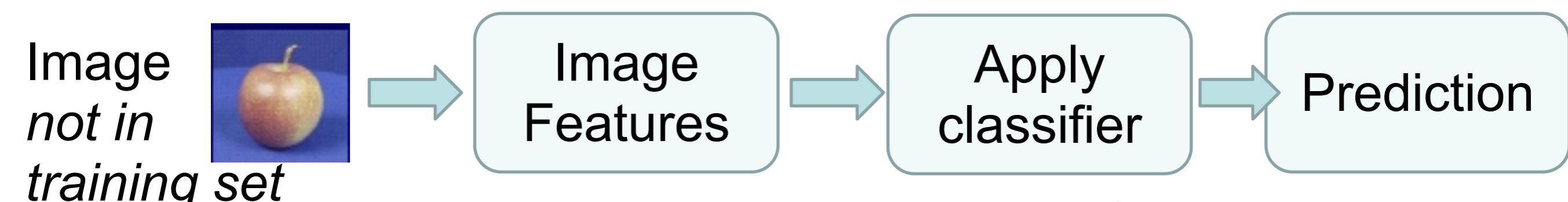
Fuzzy; definition depends on task

‘Semi-supervised’: small partial labeling

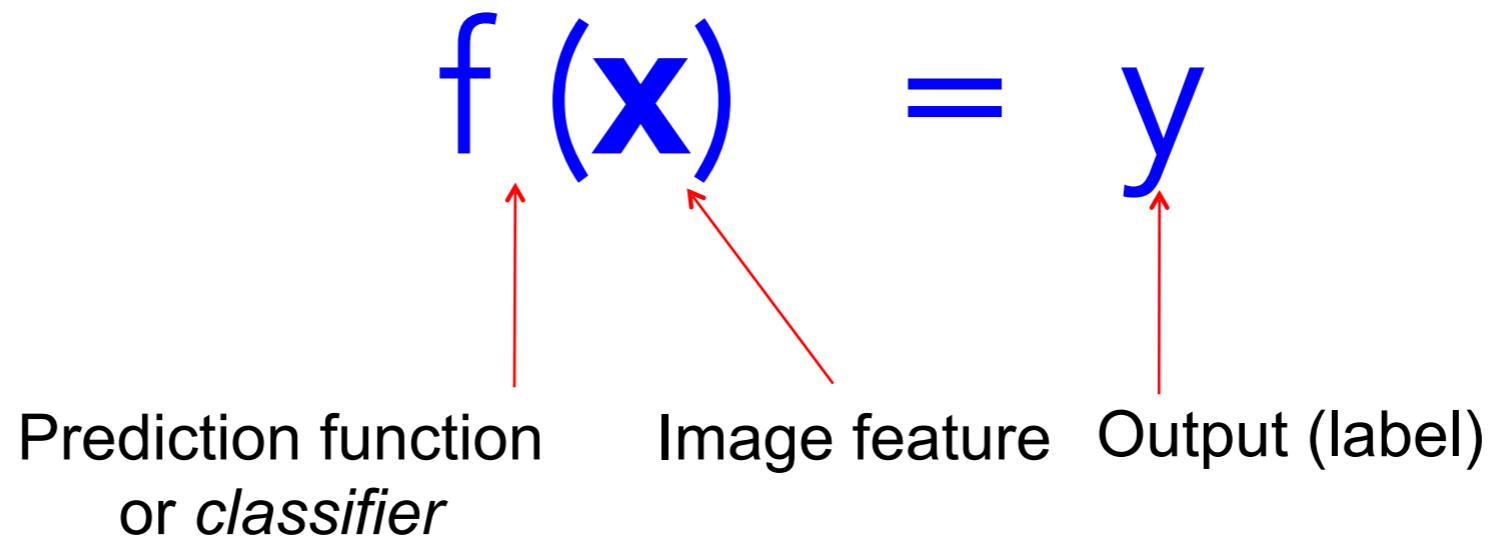
# Training



# Testing



# The machine learning framework



- ❖ **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 unseen *test example*  $\mathbf{x}_u$  and output the predicted value  $y_u = f(\mathbf{x}_u)$  to *classify*  $\mathbf{x}_u$ .

---

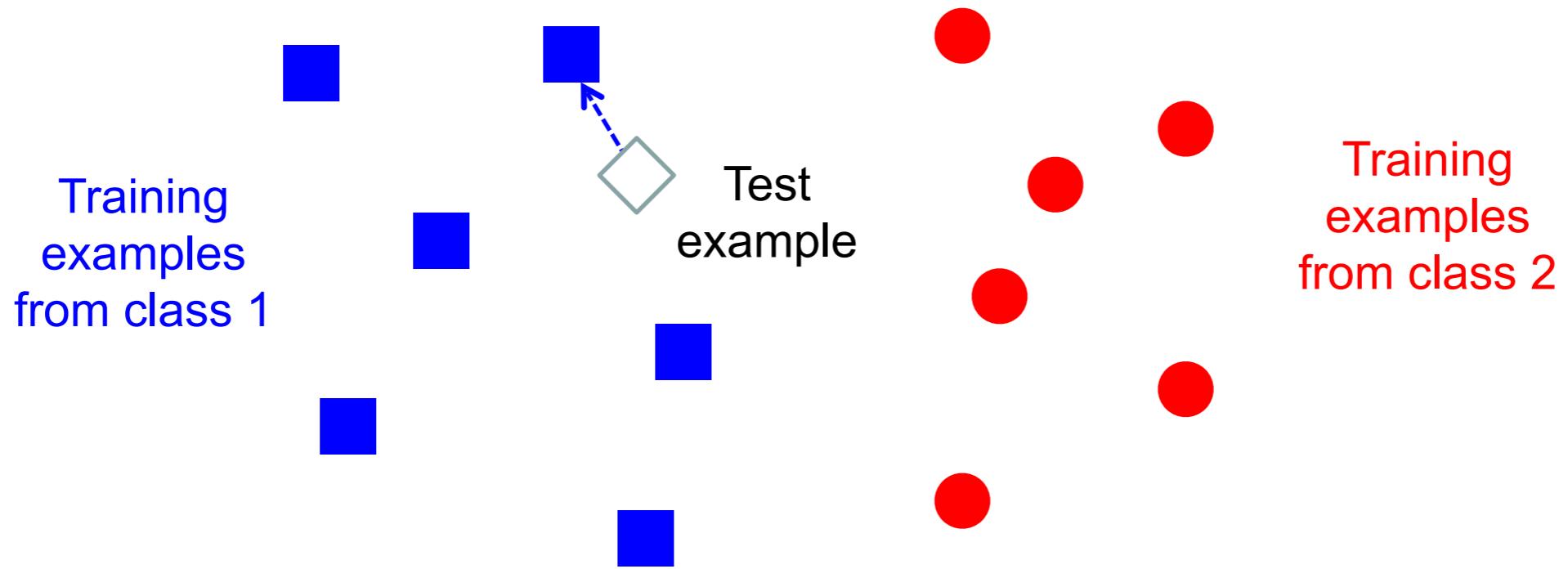
# Classification

---

Assign **x** to one of two (or more) classes.

A decision rule divides input space into *decision regions* separated by *decision boundaries* – literally boundaries in the space of the features.

# Classifiers: Nearest neighbor



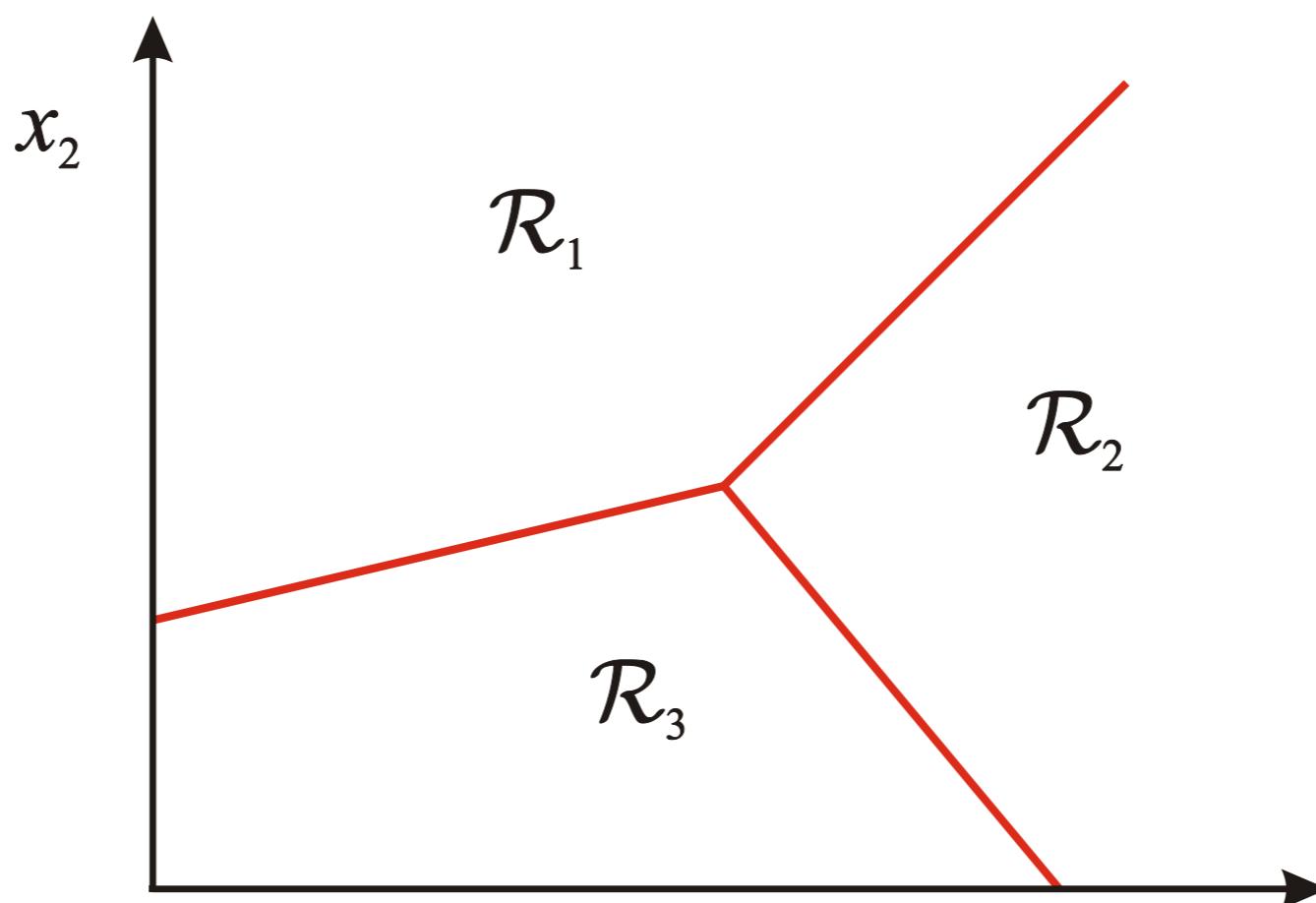
$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

- ❖ All we need is a distance function for our inputs
- ❖ No training required!

# Classification

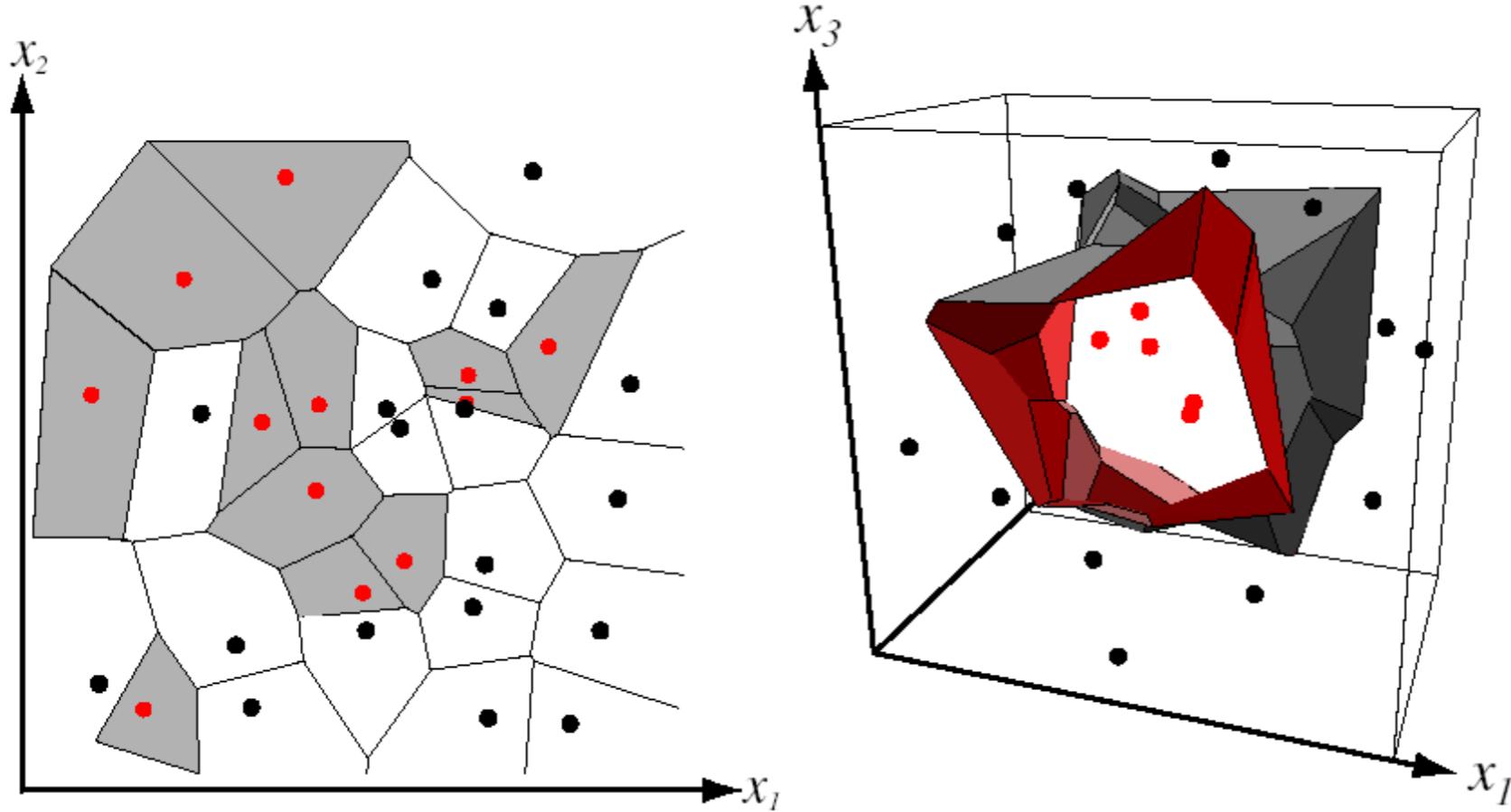
Assign  $\mathbf{x}$  to one of two (or more) classes.

A decision rule divides input space into *decision regions* separated by *decision boundaries* – literally boundaries in the space of the features.



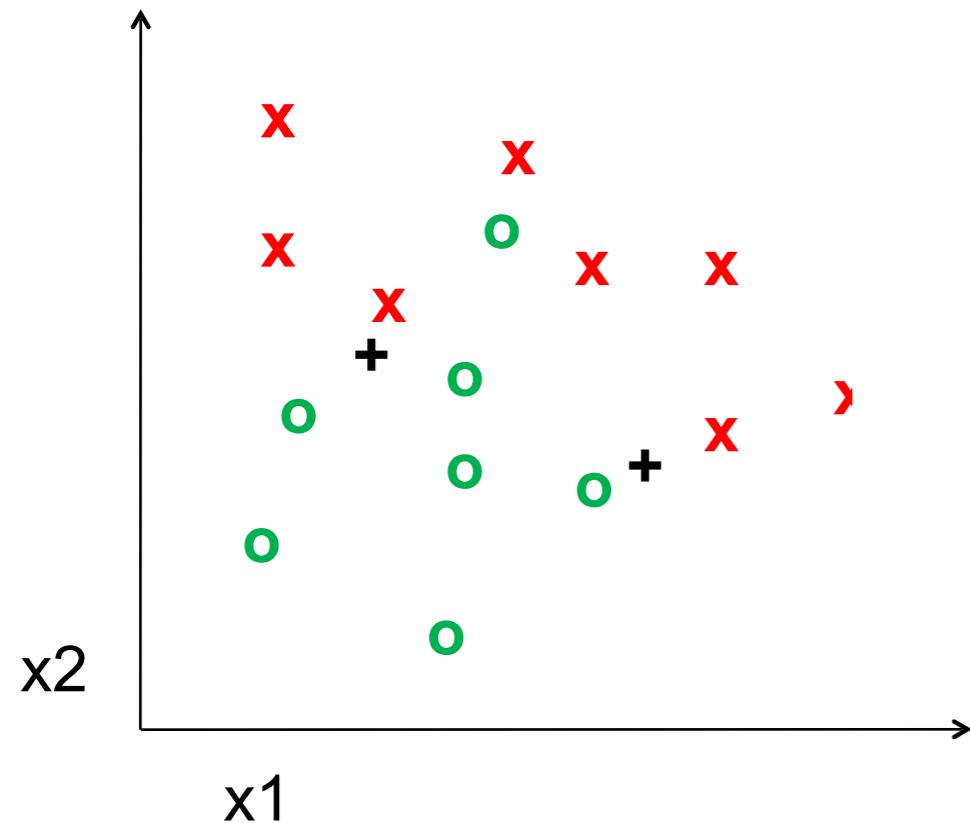
# Decision boundary for Nearest Neighbor Classifier

Divides input space into *decision regions* separated by *decision boundaries* – *Voronoi*.

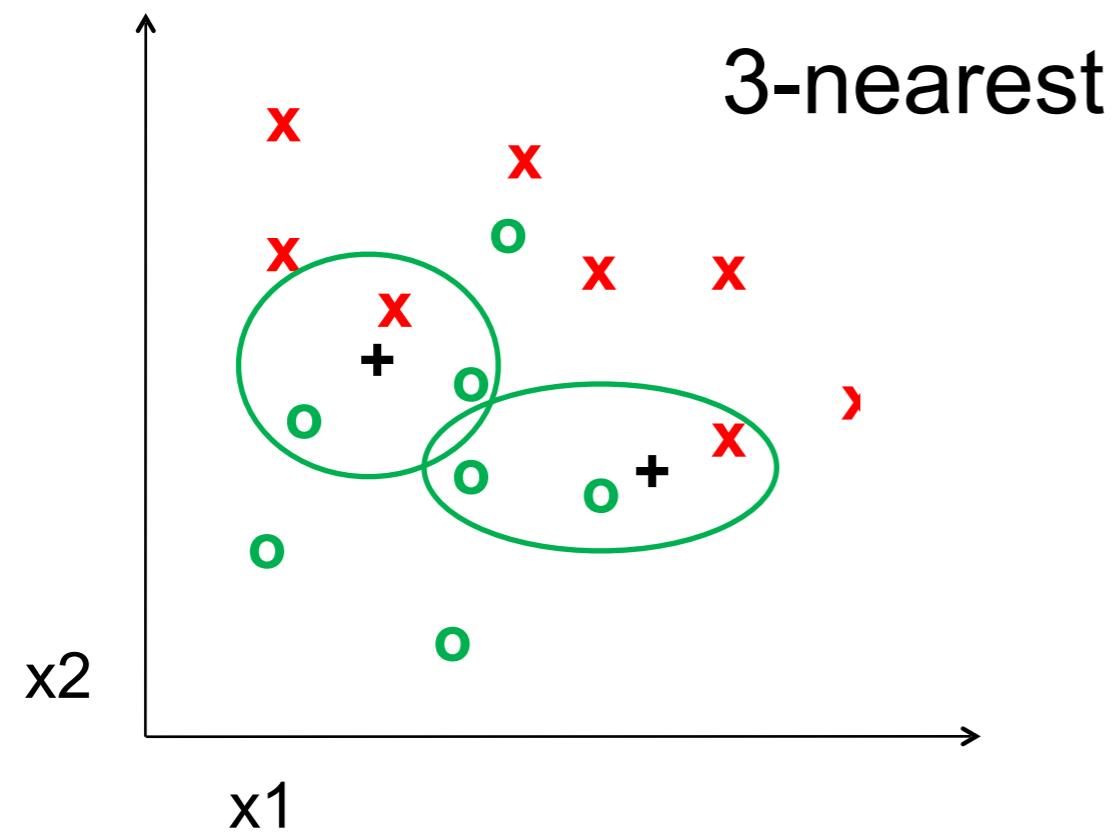


Voronoi partitioning  
of feature space  
for two-category  
2D and 3D data

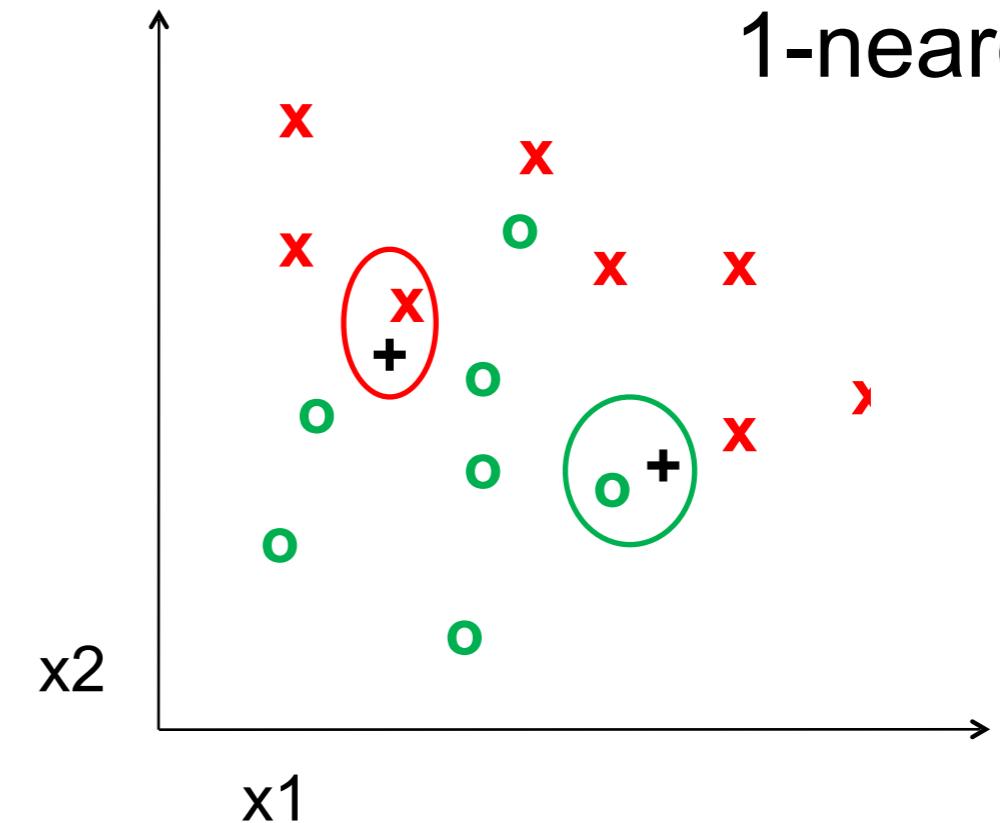
# k-nearest neighbor



1-nearest



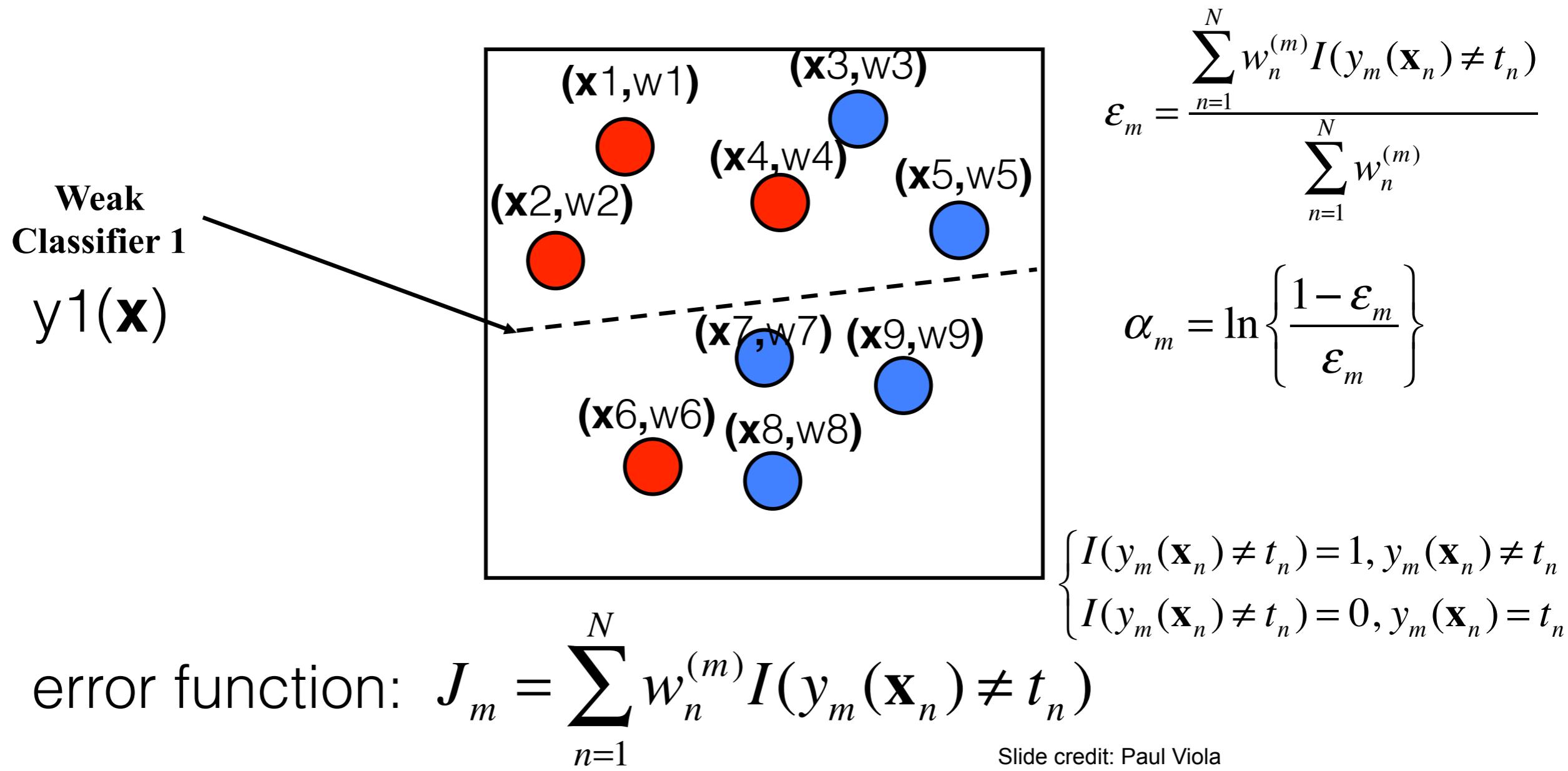
3-nearest



5-nearest

# Discriminative classifiers

## - Boosting

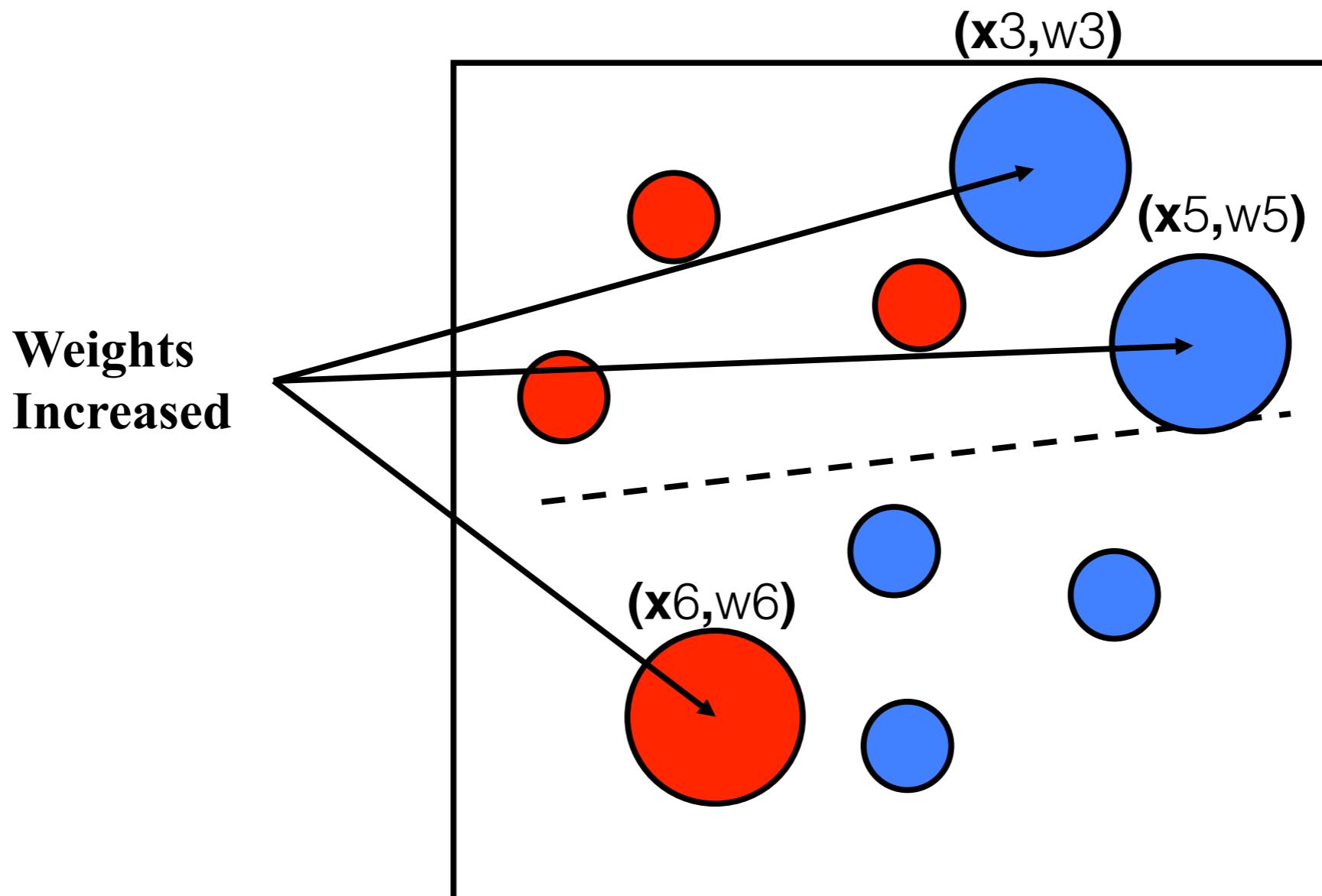


error function:  $J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)$

Slide credit: Paul Viola

# Discriminative classifiers

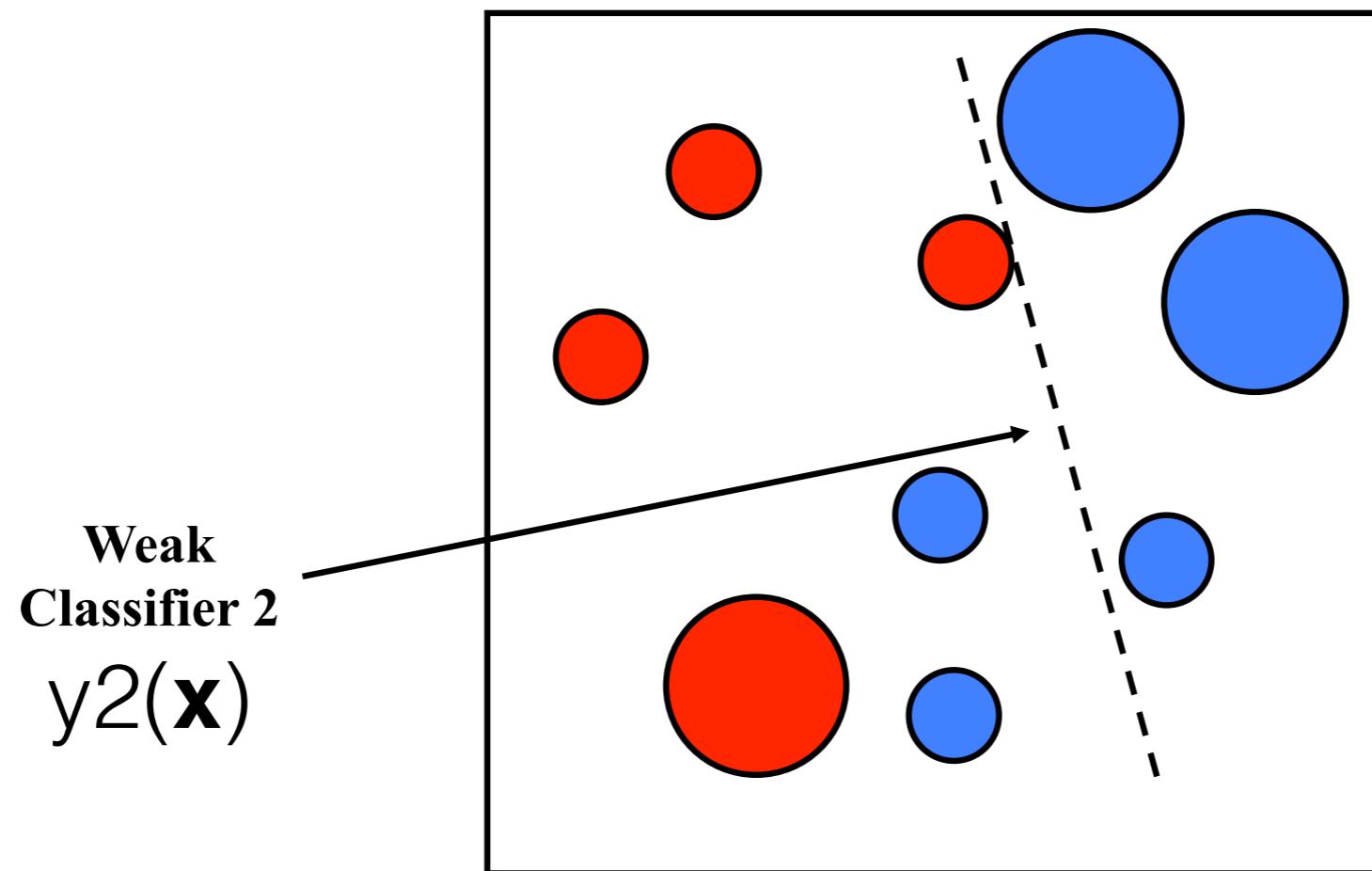
## - Boosting



$$w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(\mathbf{x}_n) \neq t_n)\}$$

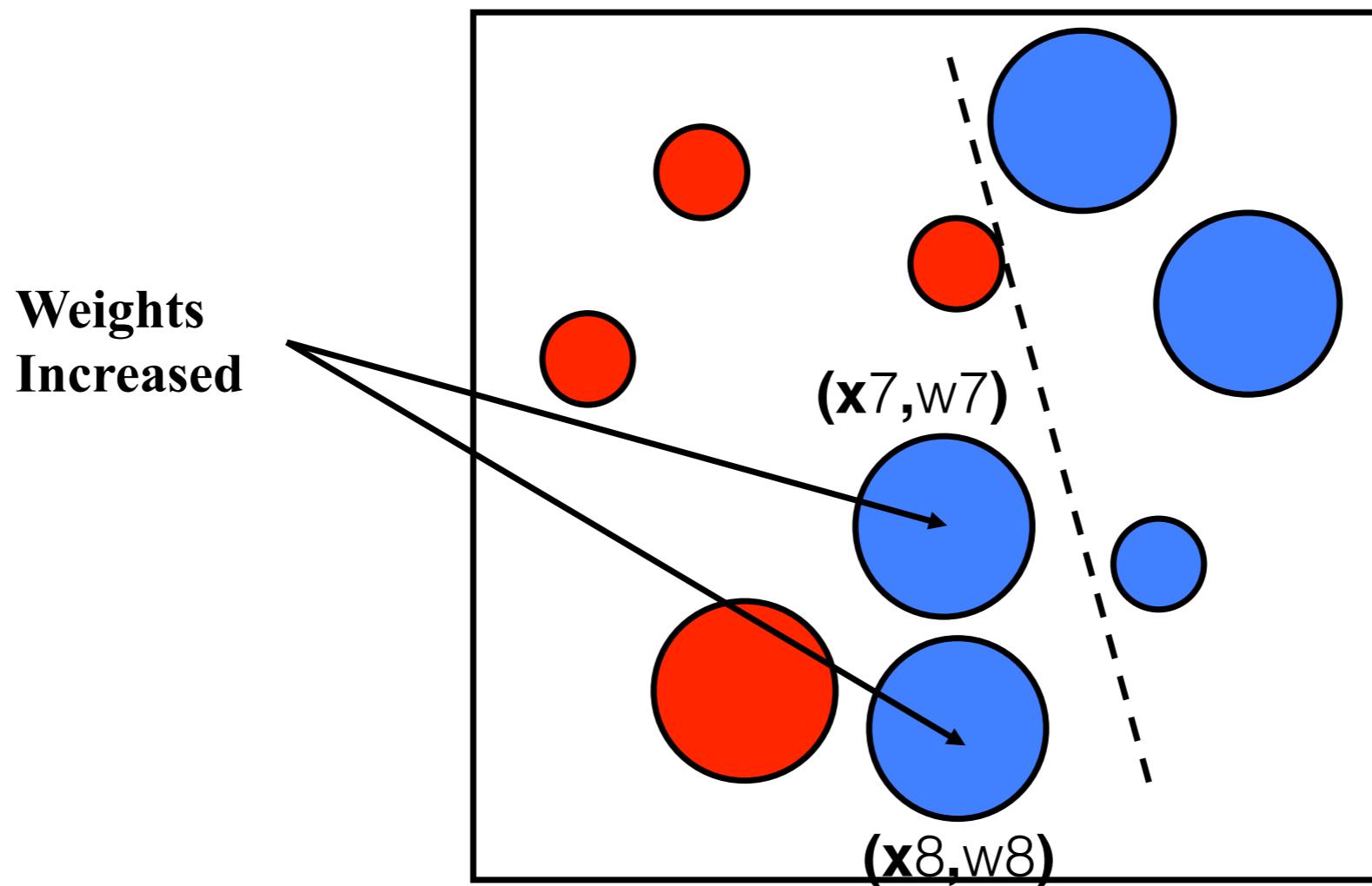
# Discriminative classifiers

## - Boosting



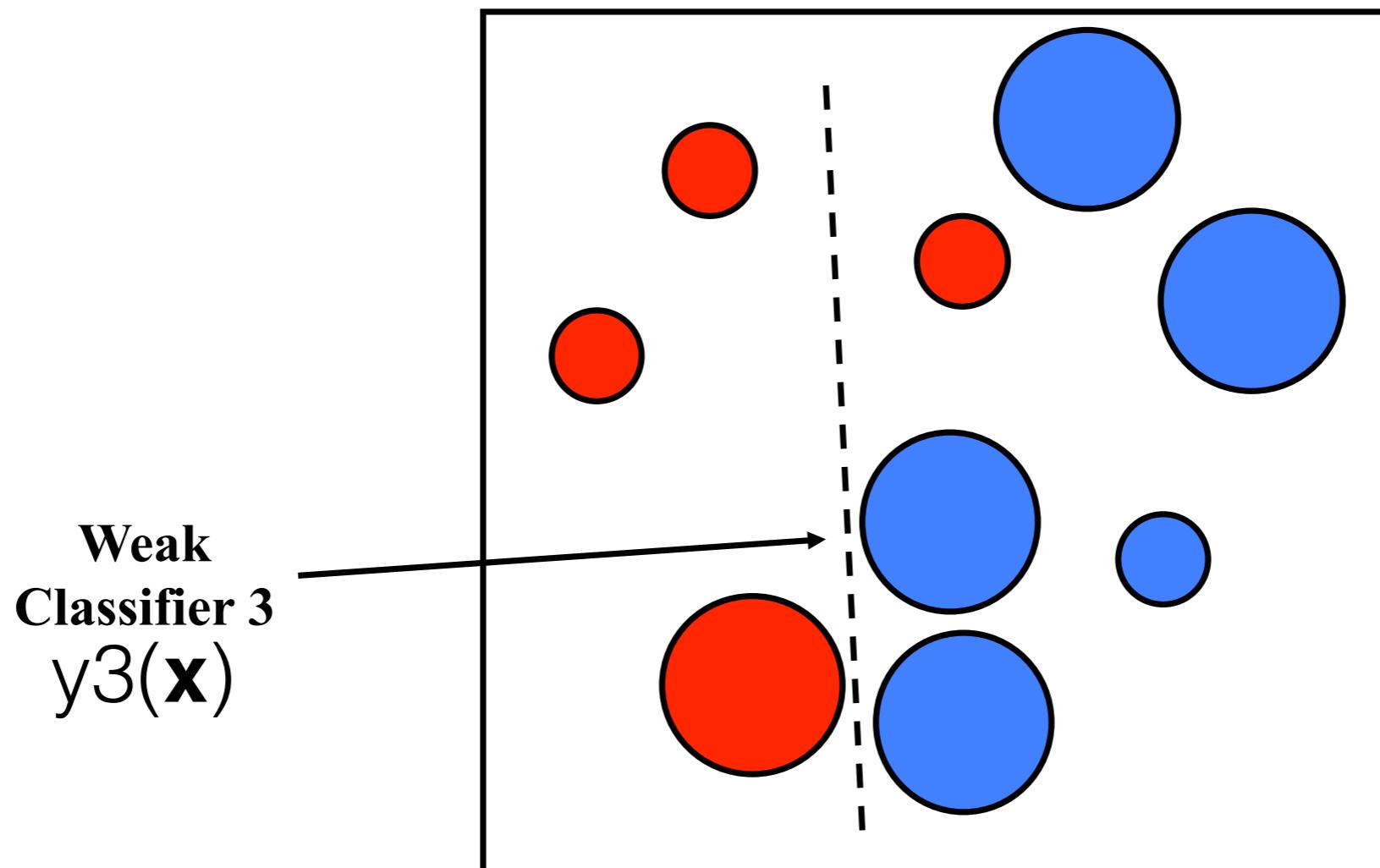
# Discriminative classifiers

- Boosting



# Discriminative classifiers

## - Boosting

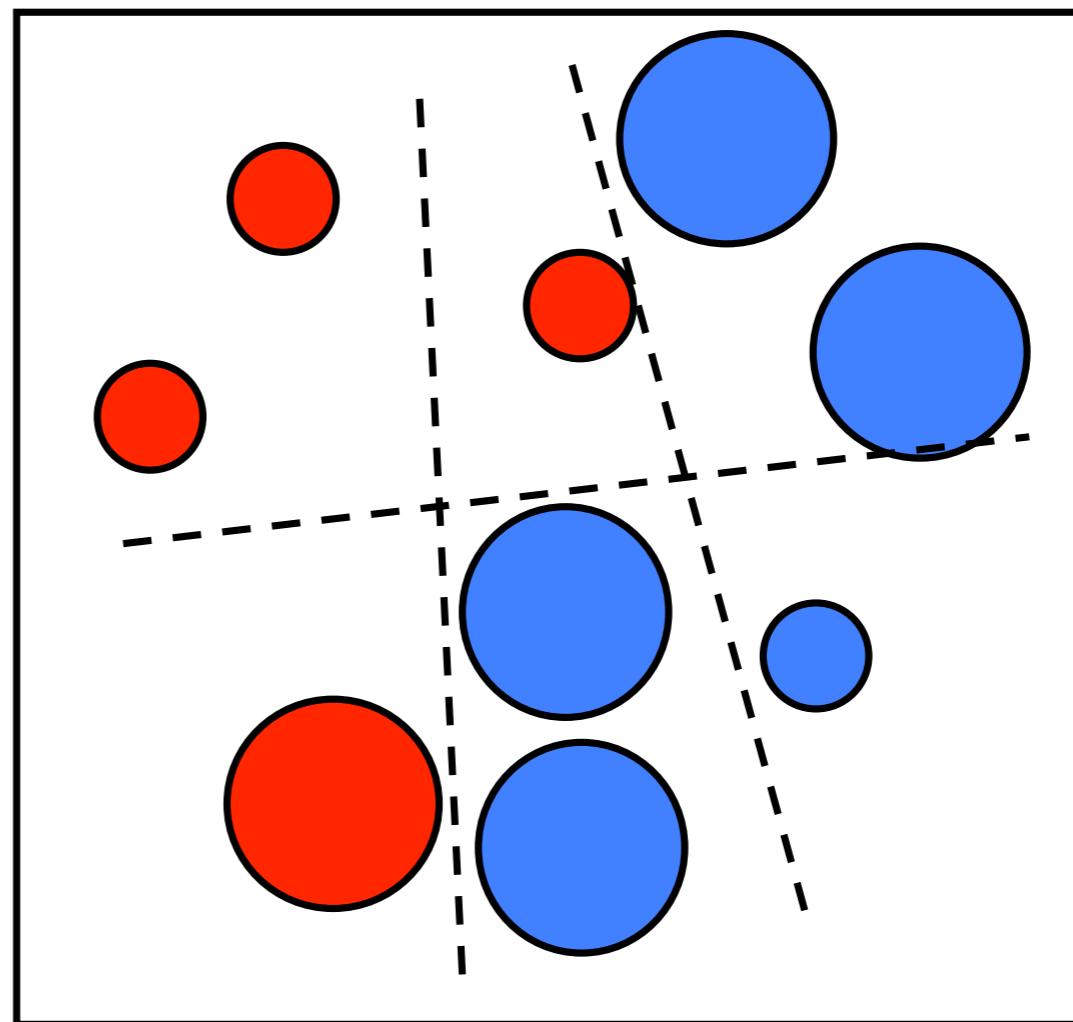


# Discriminative classifiers

## - Boosting

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

$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m y_m(\mathbf{x}) \right)$$



# Boosting: training

- Initially, weight each training example equally
- In each boosting round:
  - Find the weak learner that achieves the lowest weighted training error
  - Raise weights of training examples misclassified by current weak learner
- Compute final classifier as linear combination of all weak learners (weight of each learner is directly proportional to its accuracy)
- Exact formulas for re-weighting and combining weak learners depend on the particular boosting scheme (e.g., AdaBoost)

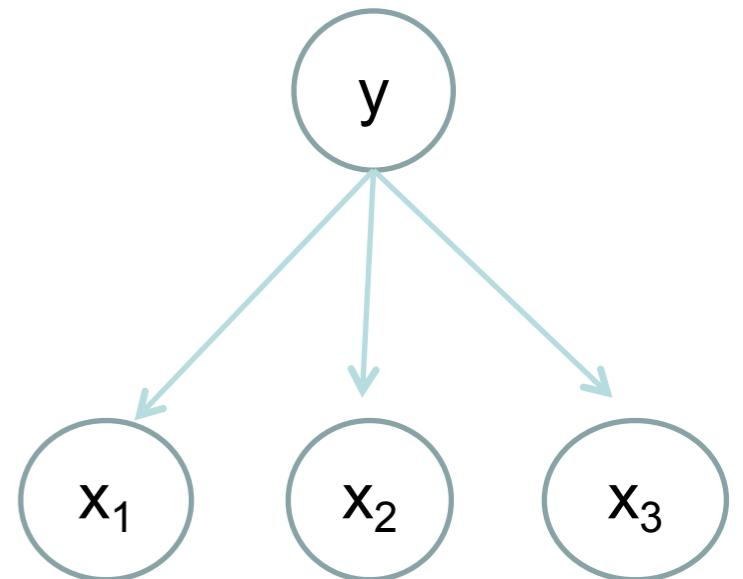
# Boosting: pros and cons

- Advantages of boosting
  - Integrates classification with feature selection
  - Complexity of training is linear in the number of training examples
  - Flexibility in the choice of weak learners, boosting scheme
  - Testing is fast
  - Easy to implement
- Disadvantages
  - Needs many training examples
  - Often found not to work as well as an alternative discriminative classifier, support vector machine (SVM)
    - especially for many-class problems

# Classifier: Naïve Bayes

$$p(C_k \mid x_1, \dots, x_n) = \frac{1}{Z} p(C_k) \prod_{i=1}^n p(x_i \mid C_k)$$

- ❖ Conditional probability model over *classes*  $C_k$



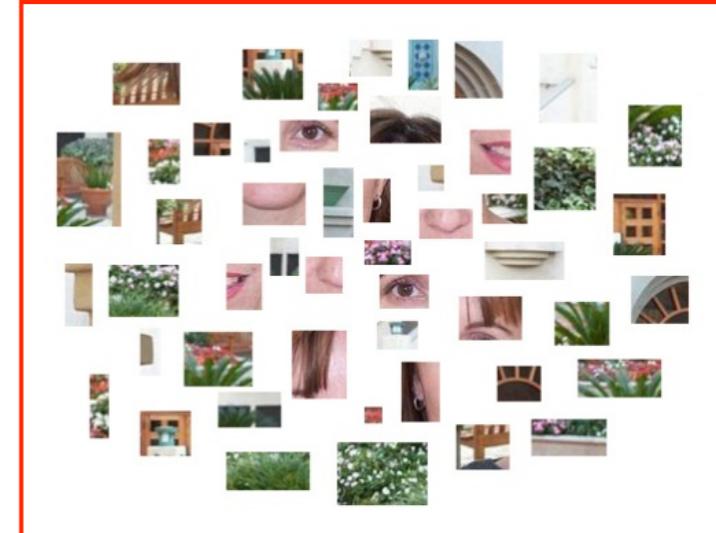
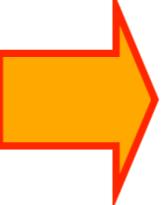
Classifier:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i \mid C_k).$$

# Naïve Bayes model for classification



$N$  patches



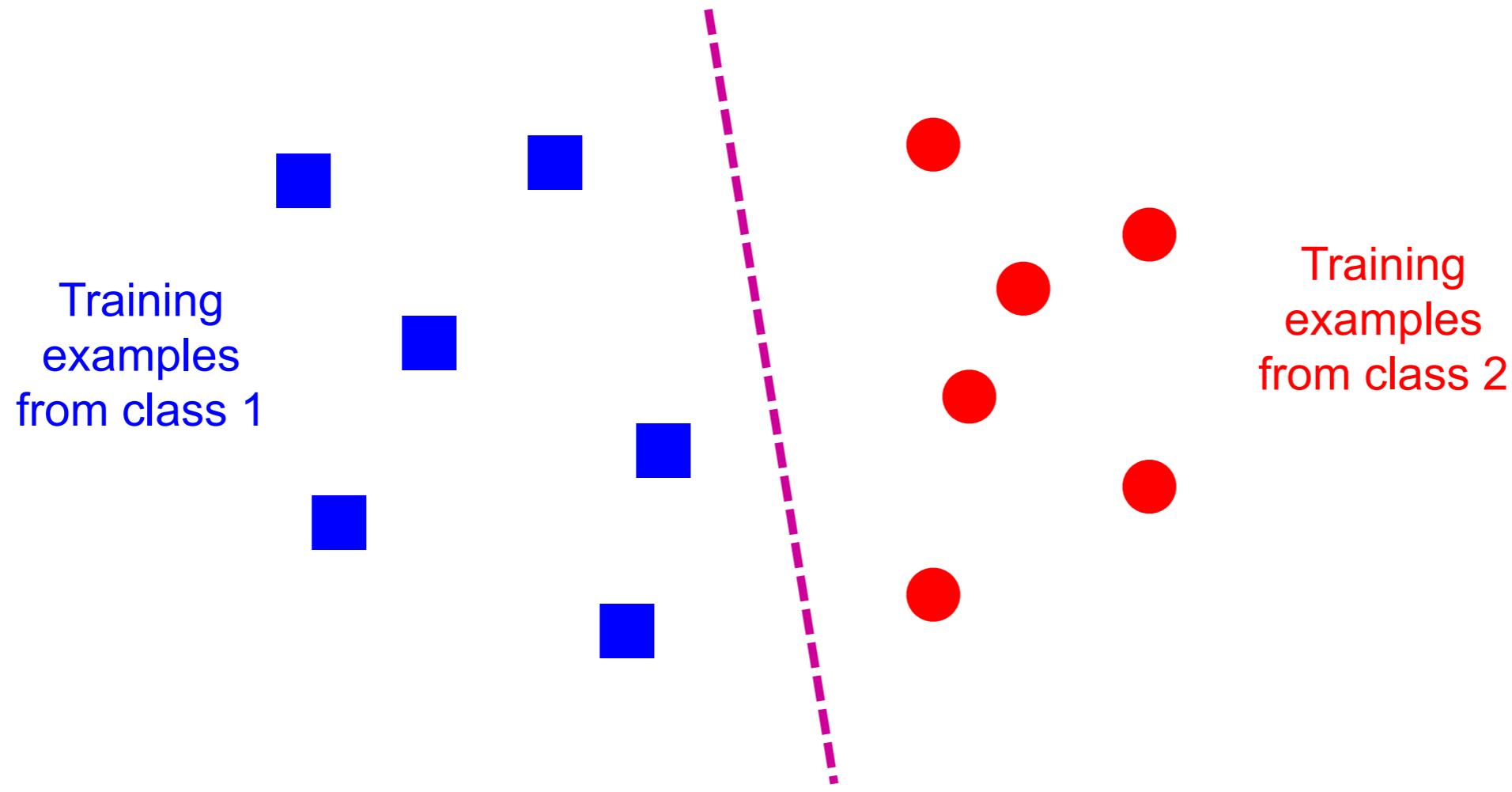
$$c^* = \arg \max_c p(c | w) \propto p(c)p(w|c) = p(c) \prod_{n=1}^N p(w_n | c)$$

Object class decision

Prior prob. of the object classes

Image likelihood given the class

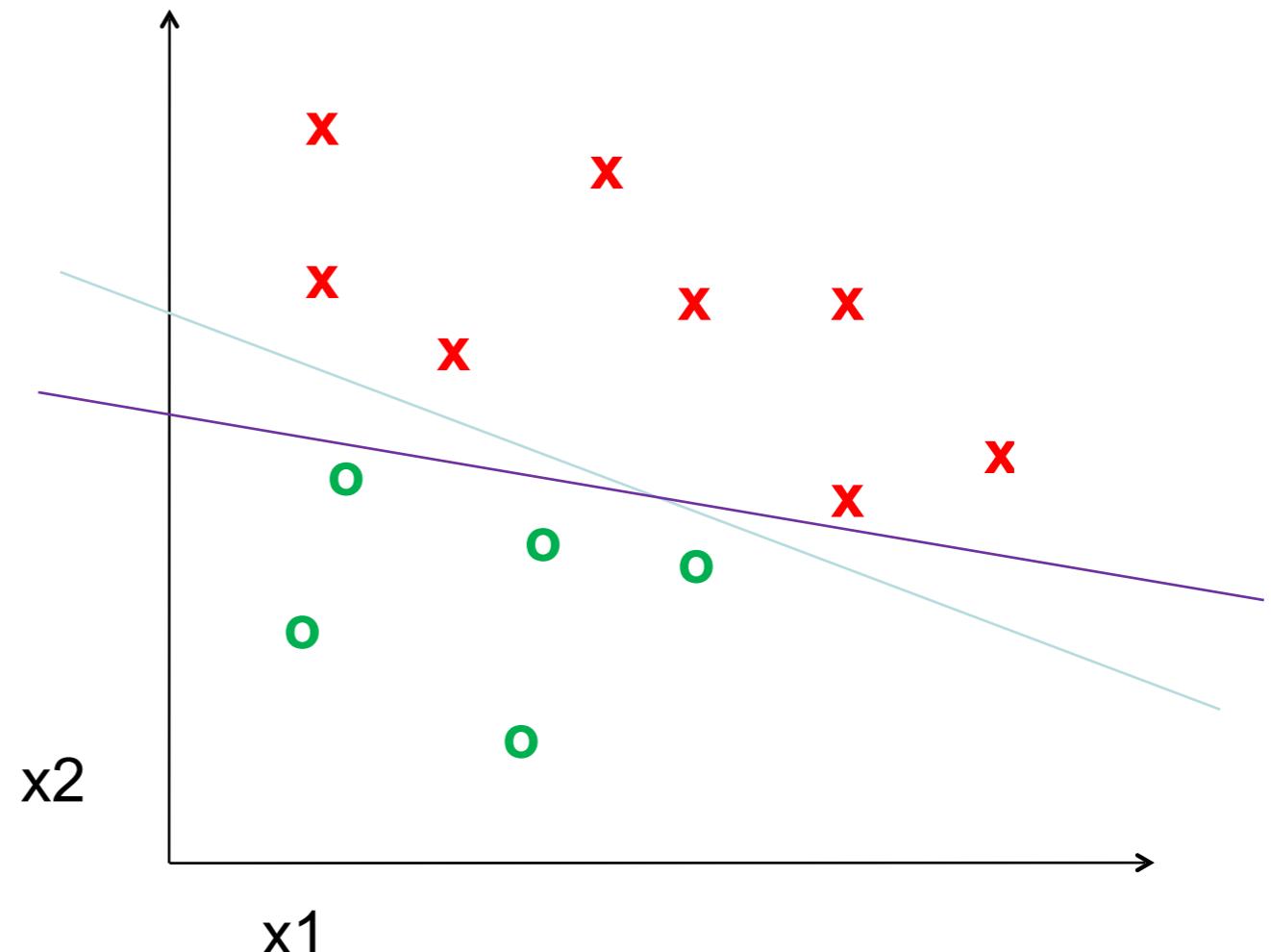
# Classifiers: Linear



Find a *linear function* to separate the classes

# Classifiers: Linear SVM

Find a *linear function*  
to separate the classes:  
 $f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$



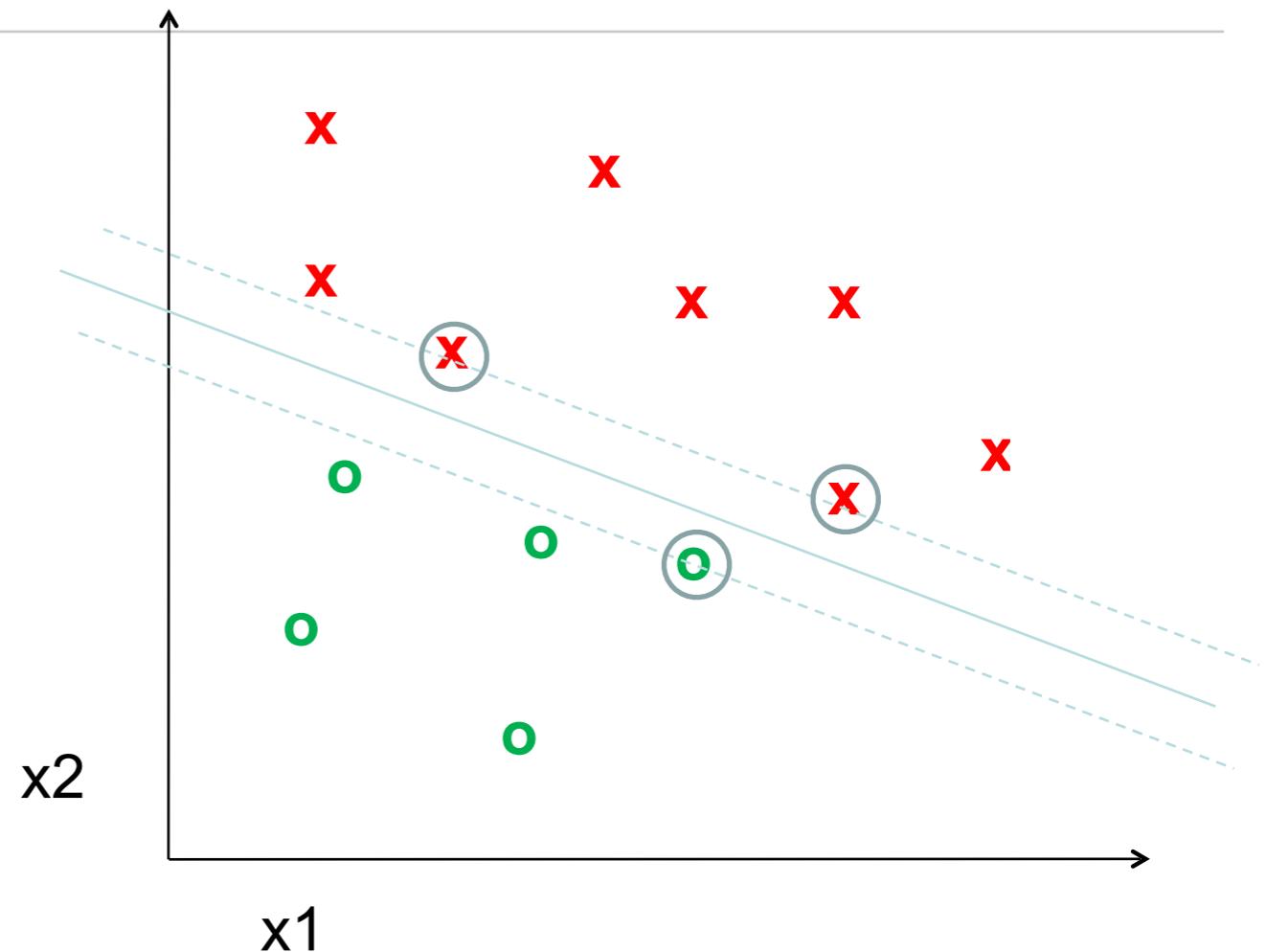
# Classifiers: Linear SVM

Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

How?

$\mathbf{X}$  = all data points



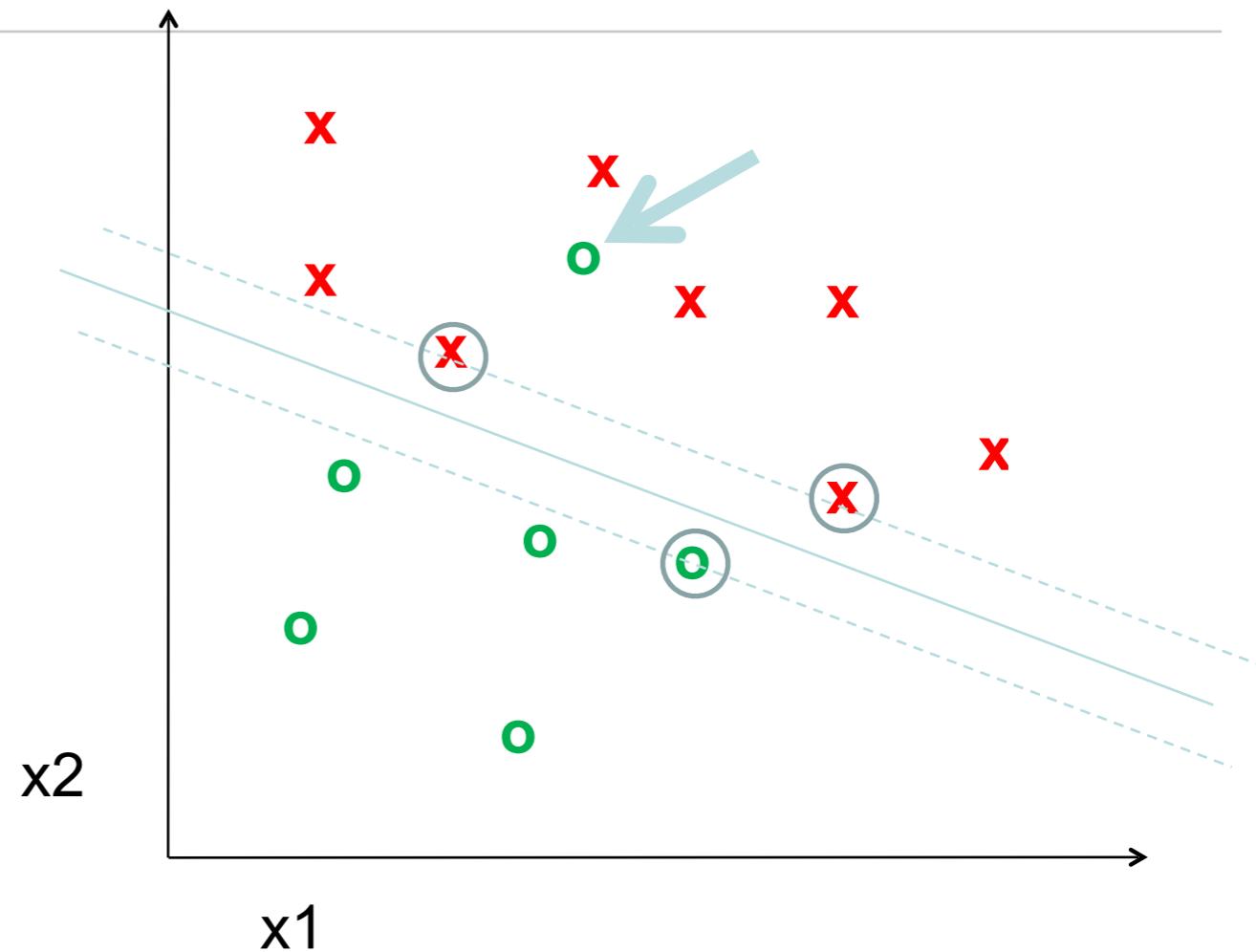
Define *hyperplane*  $t\mathbf{X} - b = 0$ , where  $t$  is tangent to hyperplane.

Minimize  $\|t\|$  s.t.  $t\mathbf{X} - b$  produces correct label for all  $\mathbf{X}$

# Classifiers: Linear SVM

Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

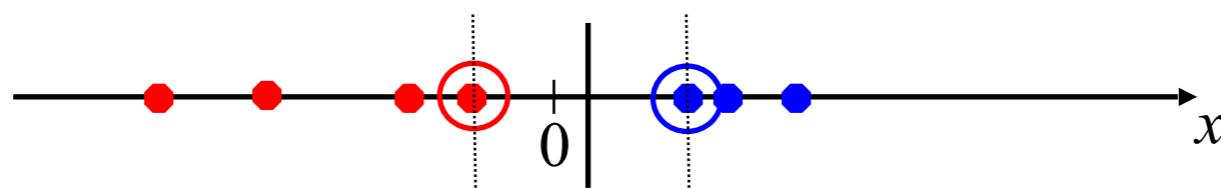


What if my data are not linearly separable?

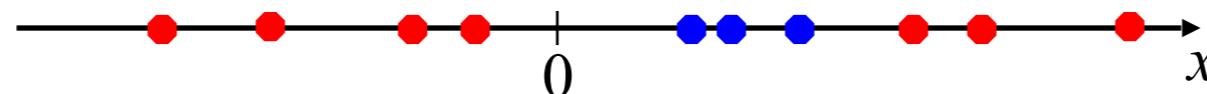
Introduce flexible ‘hinge’ loss (or ‘soft-margin’)

# Nonlinear SVMs

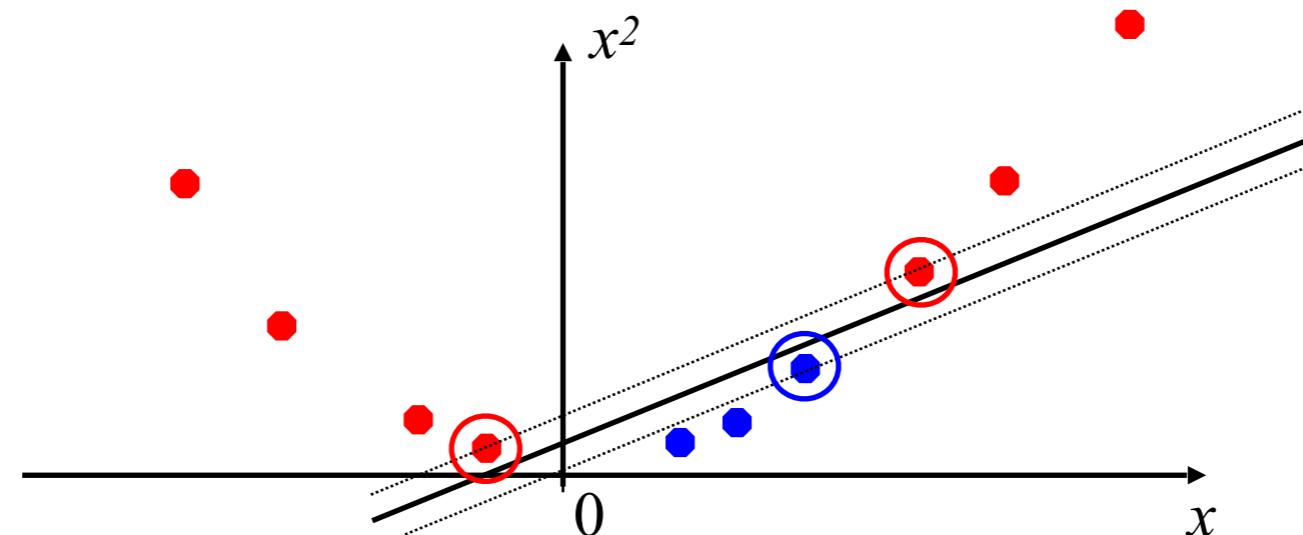
- ❖ Datasets that are linearly separable work out great:



- ❖ But what if the dataset is just too hard?

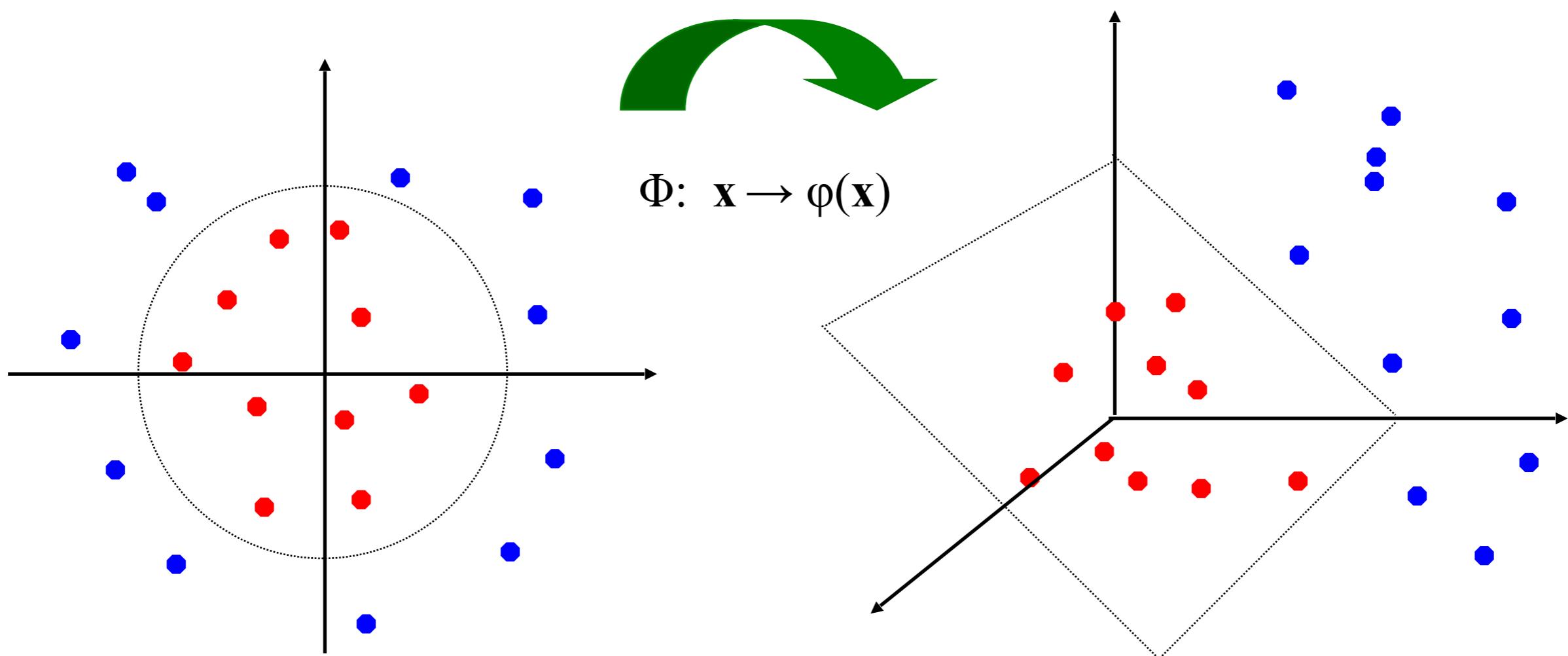


- ❖ We can map it to a higher-dimensional space:



# Nonlinear SVMs

Map the original input space to some higher-dimensional feature space where the training set is separable:



# Nonlinear SVMs

*The kernel trick:* instead of explicitly computing the lifting transformation  $\varphi(\mathbf{x})$ , define a kernel function  $K$  such that:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

This gives a *non-linear* decision boundary in the original feature space:

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

But...we only transformed the distance function  $K$ !

Common kernel function: Radial basis function kernel

# What about multi-class SVMs?

Unfortunately, there is no “definitive” multi-class SVM.

In practice, we combine multiple two-class SVMs

One vs. others

- ❖ Training: learn an SVM for each class vs. the others
- ❖ Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value

One vs. one

- ❖ Training: learn an SVM for each pair of classes
- ❖ Testing: each learned SVM “votes” for a class to assign to the test example

---

# SVMs: Pros and cons

---

- ❖ Pros
  - ❖ Many publicly available SVM packages:  
<http://www.kernel-machines.org/software>
  - ❖ Kernel-based framework is very powerful, flexible
  - ❖ SVMs work very well in practice, even with very small training sample sizes
- ❖ Cons
  - ❖ No “direct” multi-class SVM, must combine two-class SVMs
  - ❖ Computation, memory
    - ❖ During training time, must compute matrix of kernel values for every pair of examples
    - ❖ Learning can take a very long time for large-scale problems

---

# Ideals for a classification algorithm

---

- ❖ Objective function: encodes the right loss for the problem
- ❖ Parameterization: takes advantage of the structure of the problem
- ❖ Regularization: good priors on the parameters
- ❖ Training algorithm: can find parameters that maximize objective on training set
- ❖ Inference algorithm: can solve for labels that maximize objective function for a test example

# Two ways to think about classifiers

---

1. What is the objective?  
What are the parameters?  
How are the parameters learned?  
How is the learning regularized?  
How is inference performed?
  
2. How is the data modeled?  
How is similarity defined?  
What is the shape of the boundary?

# Training

Training  
Images



Image  
Features

Training  
Labels

Training

Learned  
classifier

# Testing



Test Image

Image  
Features

Apply  
classifier

Prediction

**Features and distance measures**

*define visual similarity.*

**Training labels**

*dictate that examples are the same or different.*

**Classifiers**

*learn weights (or parameters) of features and distance measures...*

*so that visual similarity predicts label similarity.*

# Generative vs. Discriminative Classifiers

## Discriminative Models

- Learn to directly predict the labels from the data
- Often, assume a simple boundary (e.g., linear)
- Examples
  - Logistic regression
  - SVM
  - Boosted decision trees
- Often easier to predict a label from the data than to model the data

## Generative Models

- Represent both the data and the labels
- Often, makes use of conditional independence and priors
- Examples
  - Naïve Bayes classifier
  - Bayesian network
- Models of data may apply to future prediction problems

---

# Making decisions about data

---

- ❖ 3 important design decisions:
  - 1) What data do I use?
  - 2) How do I represent my data (what feature)?
  - 3) What classifier / regressor / machine learning tool do I use?
- ❖ These are in decreasing order of importance
- ❖ Deep learning addresses 2 and 3 simultaneously (and blurs the boundary between them).
- ❖ You can take the representation from deep learning and use it with any classifier.

---

# Many classifiers to choose from...

---

- ❖ K-nearest neighbor
- ❖ SVM
- ❖ Naïve Bayes
- ❖ Bayesian network
- ❖ Logistic regression
- ❖ Randomized Forests
- ❖ Boosted Decision Trees
- ❖ Restricted Boltzmann Machines
- ❖ Neural networks
- ❖ Deep Convolutional Network
- ❖ ...

Claim:

*It is more important to have more or better labeled data than to use a different supervised learning technique.*

“The Unreasonable Effectiveness of Data” - Norvig

# What to remember about classifiers

---

- ❖ No free lunch: machine learning algorithms are tools, not dogmas
- ❖ Try simple classifiers first
- ❖ Better to have smart features and simple classifiers than simple features and smart classifiers
- ❖ Use increasingly powerful classifiers with more training data (bias-variance tradeoff)