

# Introduction to Recognition

Computer Vision

CS 543 / ECE 549

University of Illinois

# Outline

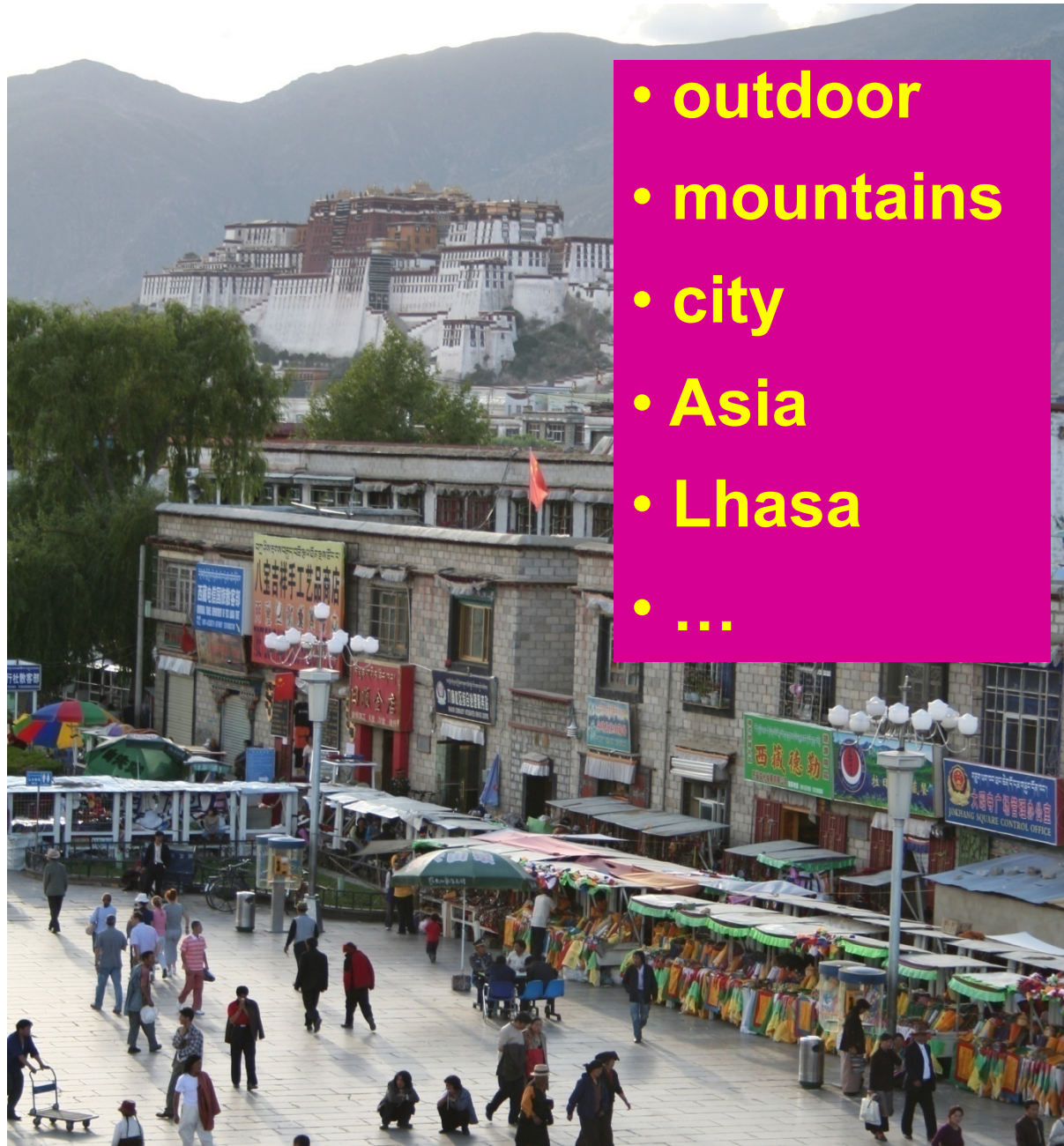
- Overview of image and region categorization
  - Task description
  - What is a category
- Example of spatial pyramids bag-of-words scene categorizer
- Key concepts: features and classification
- Deep convolutional neural networks (CNNs)

# Common recognition tasks



Adapted from  
Fei-Fei Li

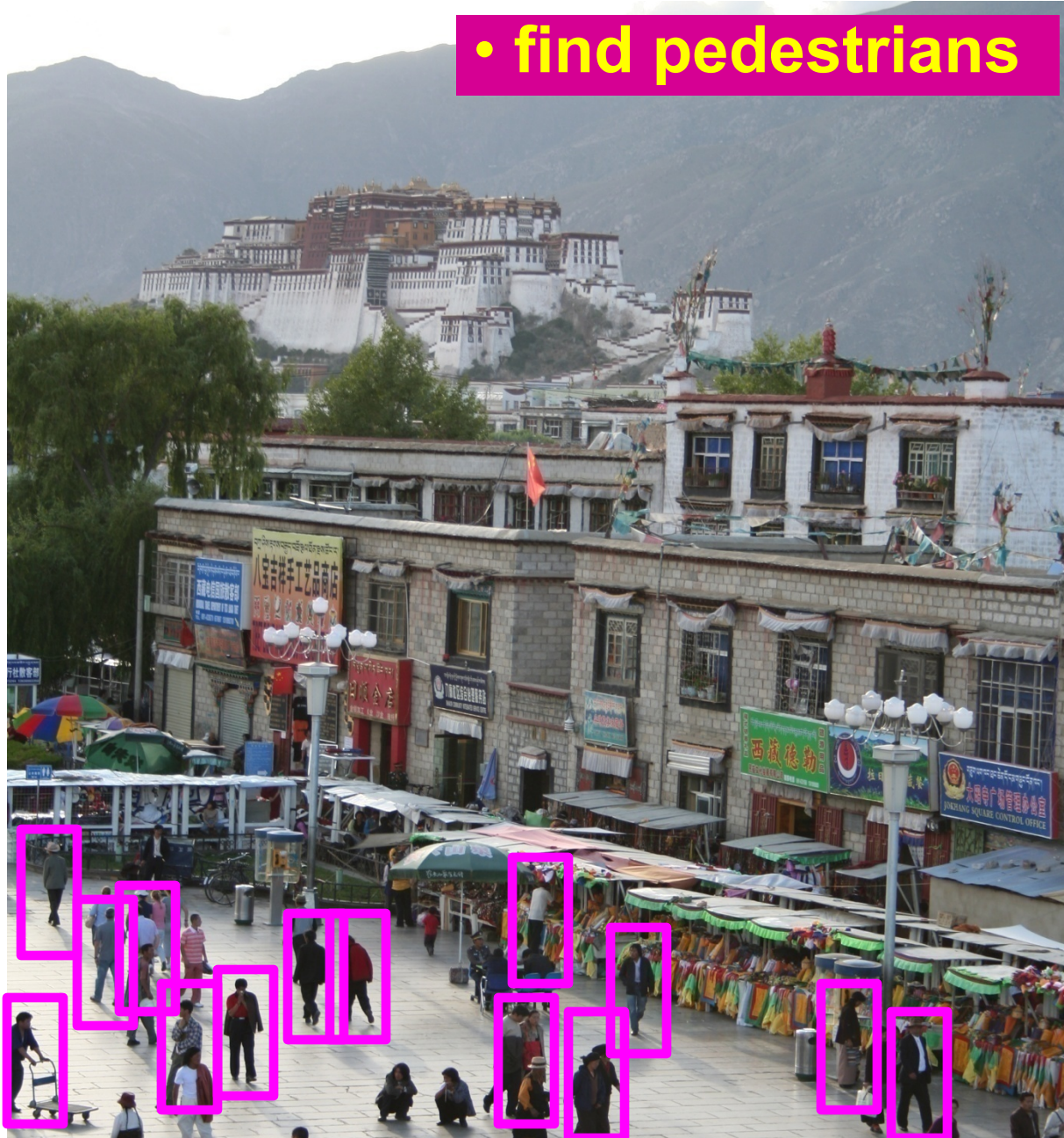
# Image classification and tagging



- outdoor
- mountains
- city
- Asia
- Lhasa
- ...

# Object detection

- find pedestrians



# Activity recognition



- walking
- shopping
- rolling a cart
- sitting
- talking
- ...



# Semantic segmentation



Adapted from  
Fei-Fei Li

# Semantic segmentation

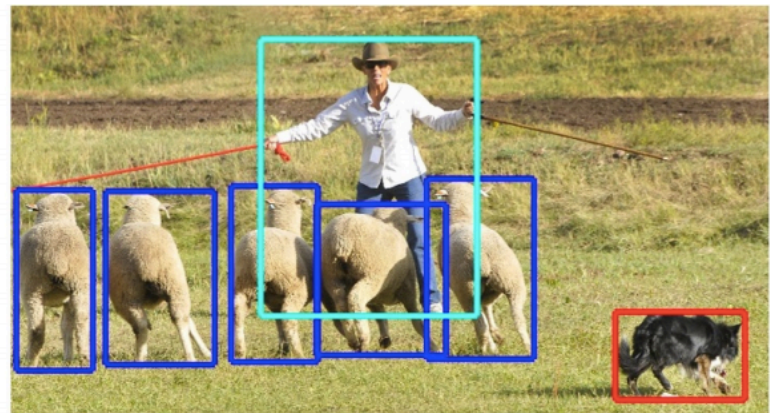




# Detection, semantic segmentation, instance segmentation



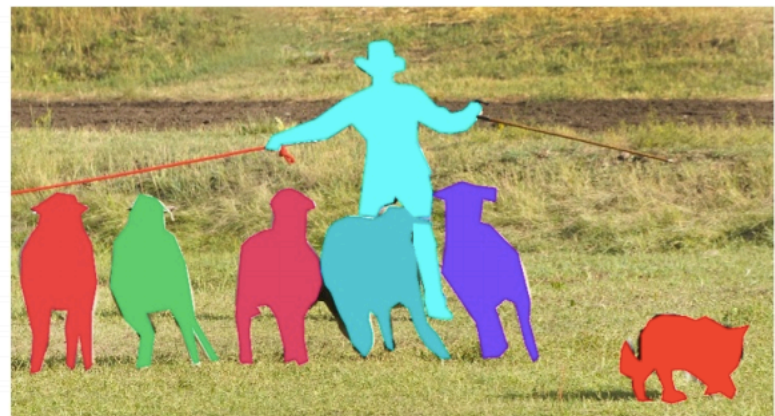
image classification



object detection

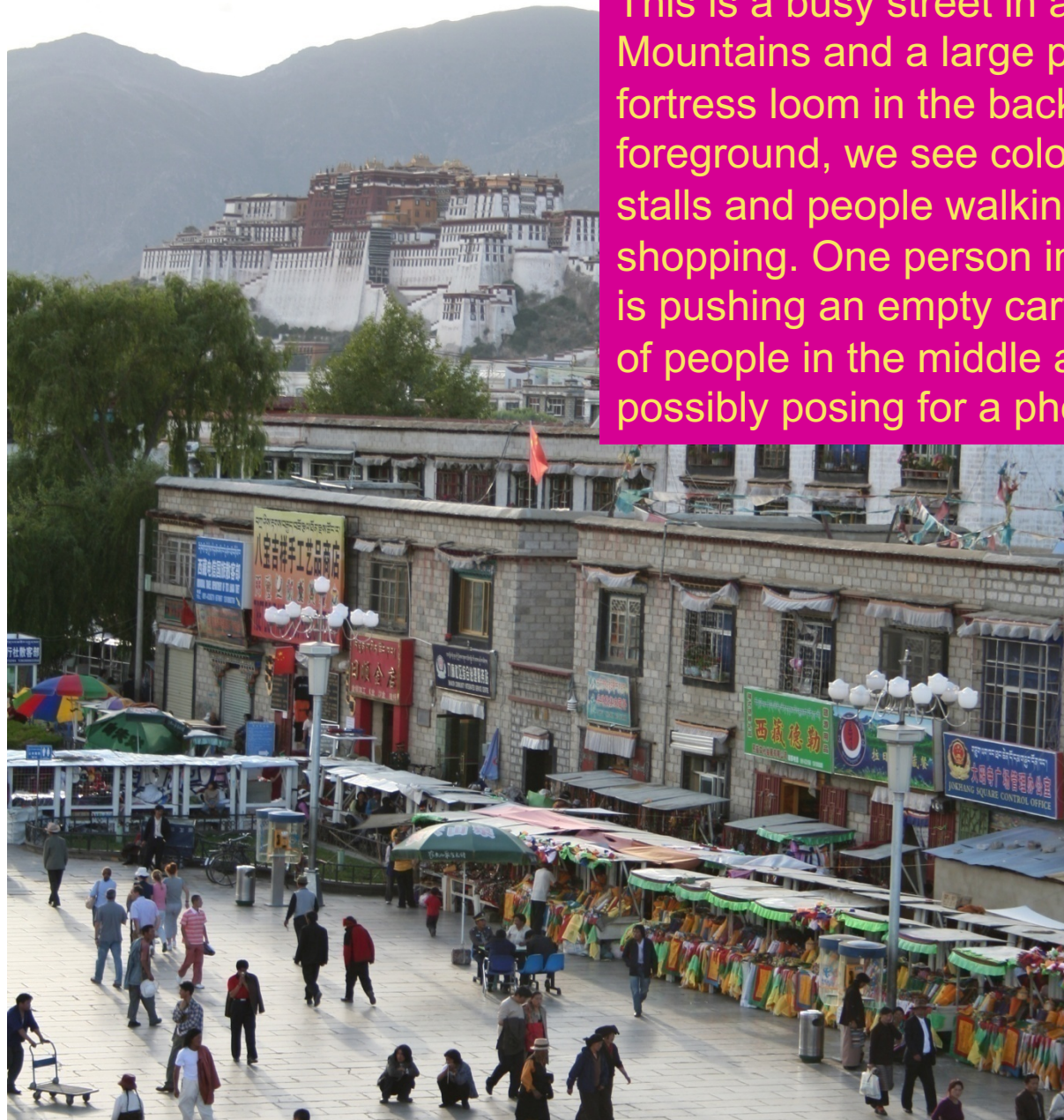


semantic segmentation



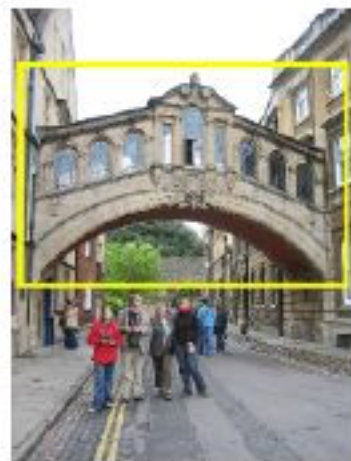
instance segmentation

# Image description



This is a busy street in an Asian city. Mountains and a large palace or fortress loom in the background. In the foreground, we see colorful souvenir stalls and people walking around and shopping. One person in the lower left is pushing an empty cart, and a couple of people in the middle are sitting, possibly posing for a photograph.

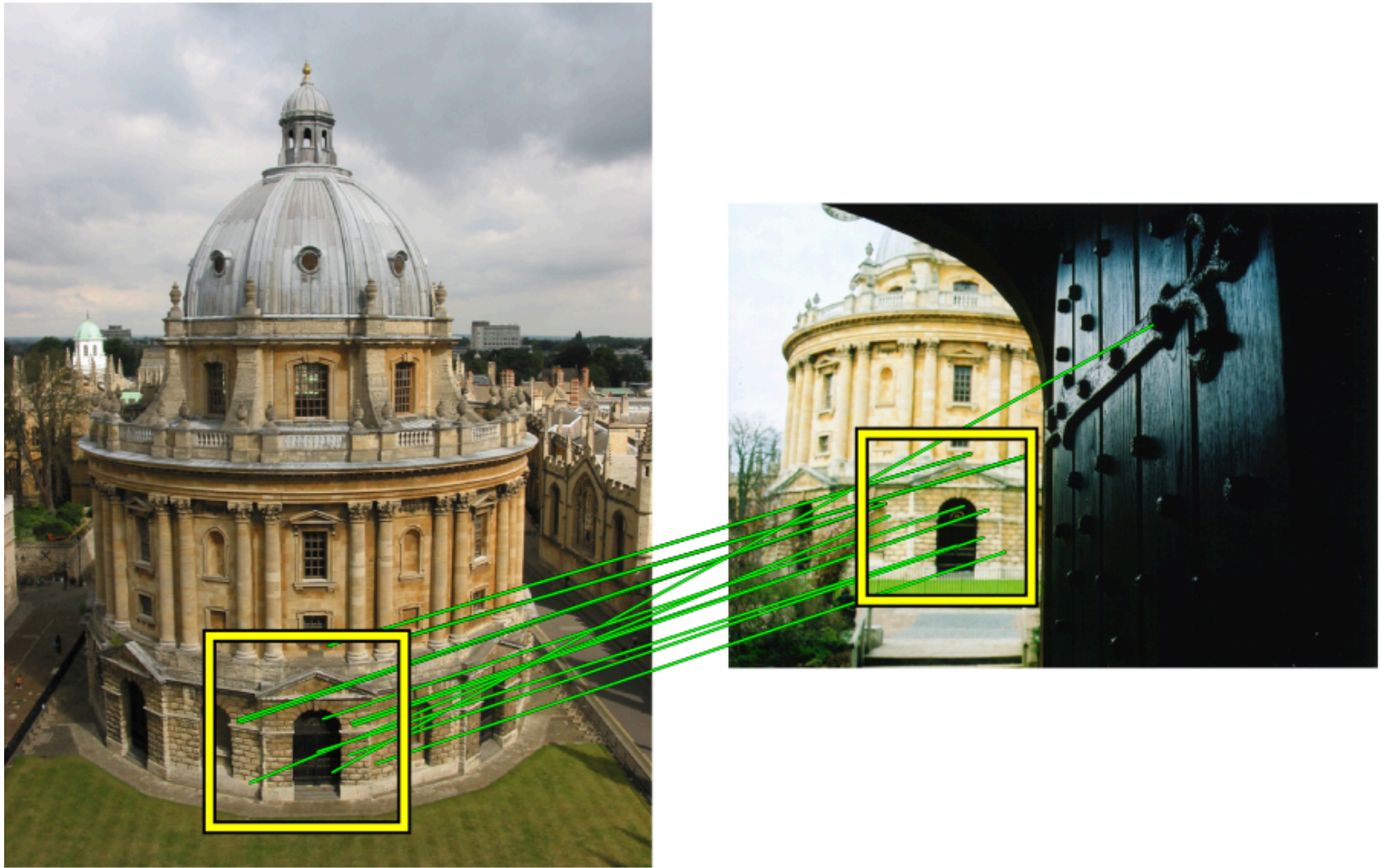
# Recognition as 3D Matching



Find these landmarks

...In these images

# Recognition as 3D Matching



Recognizing solid objects by alignment with an image. Huttenlocher and Ullman IJCV 1990.

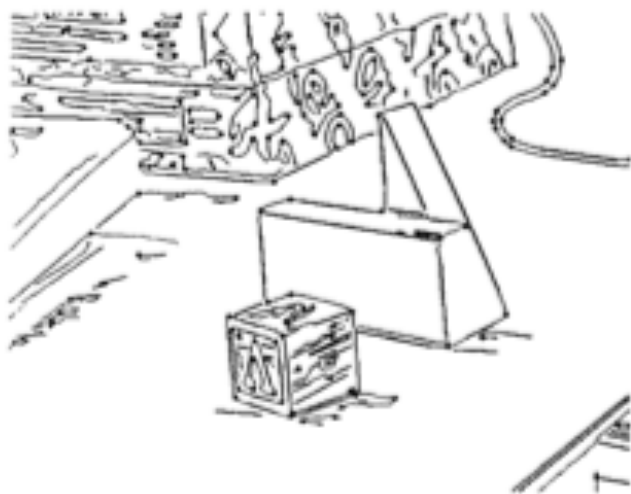
# Recognition as 3D Matching



a)



b)



c)



d)

“Instance”  
Recognition  
“Category-level”  
Recognition

Fig. 8. The output of the recognizer: (a) grey-level image input, (b) Canny edges, (c) edge segments, (d) recovered instances.

Recognizing solid objects by alignment with an image. Huttenlocher and Ullman IJCV 1990.

# Theory of categorization

How do we determine if something is a member of a particular category?

- Definitional approach
- Prototype approach
- Exemplar approach

# Definitional approach: classical view of categories

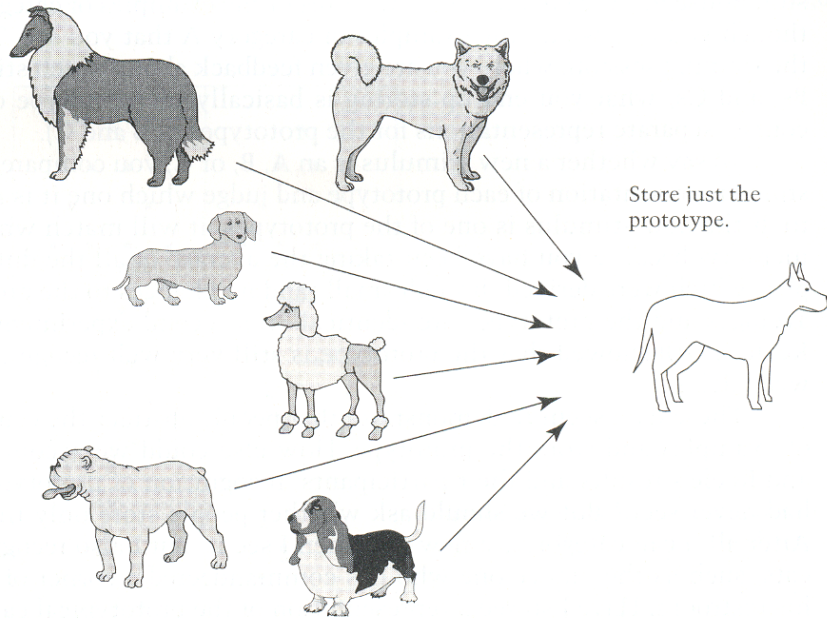
- Plato & Aristotle
  - Categories are defined by a list of properties shared by all elements in a category
  - Category membership is binary
  - Every member in the category is equal



Aristotle by Francesco Hayez

## The Categories (Aristotle)

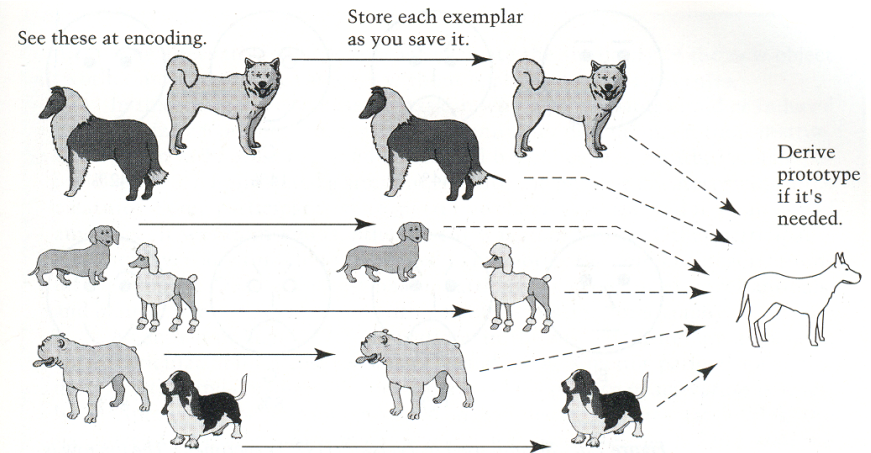
# Prototype Model



**Figure 7.3.** Schematic of the prototype model. Although many exemplars are seen, only the prototype is stored. The prototype is updated continually to incorporate more experience with new exemplars.

Category judgments are made by comparing a new exemplar to the prototype.

# Exemplars Model



**Figure 7.4.** Schematic of the exemplar model. As each exemplar is seen, it is encoded into memory. A prototype is abstracted only when it is needed, for example, when a new exemplar must be categorized.

Category judgments are made by comparing a new exemplar to all the old exemplars of a category or to the exemplar that is the most appropriate

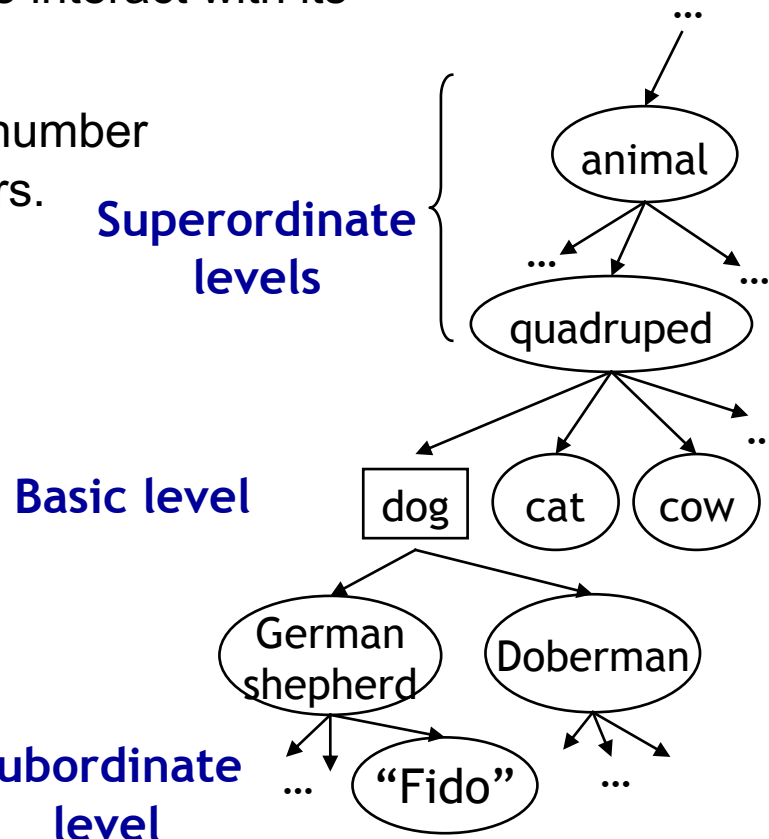
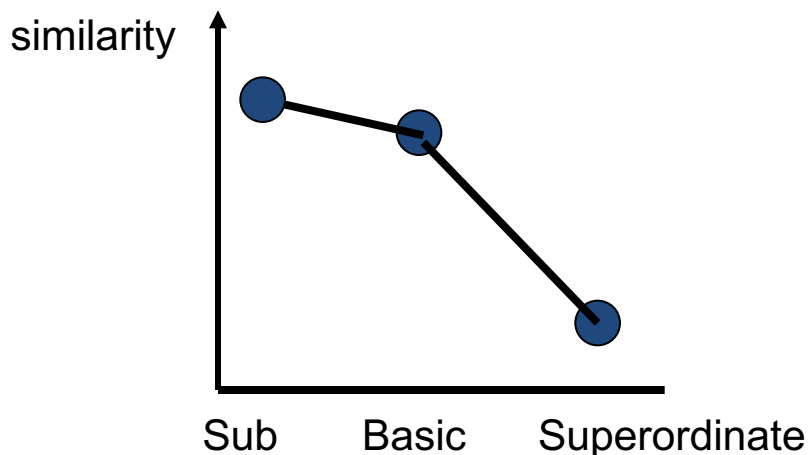


# Levels of categorization [Rosch 70s]



Definition of Basic Level:

- **Similar shape:** Basic level categories are the highest-level category for which their members have similar shapes.
- **Similar motor interactions:** ... for which people interact with its members using similar motor sequences.
- **Common attributes:** ... there are a significant number of attributes in common between pairs of members.



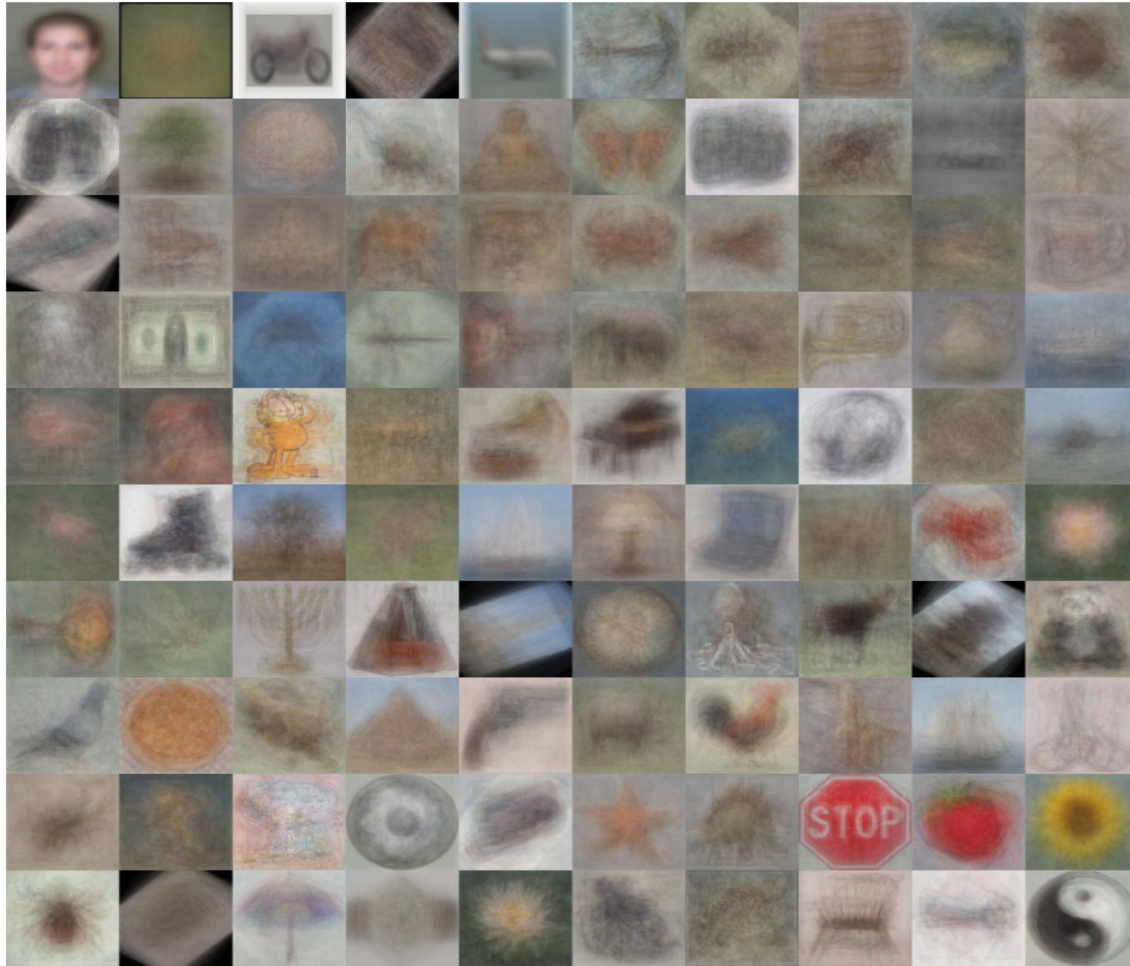
# Image categorization

- Cat vs Dog



# Image categorization

- Object recognition



Caltech 101 Average Object Images

# Image categorization

- Fine-grained recognition



Generalist



Insect catching



Grain eating



Coniferous-seed eating



Nectar feeding



Chiseling



Dip netting



Surface skimming



Scything



Probing



Aerial fishing



Pursuit fishing



Scavenging



Raptorial

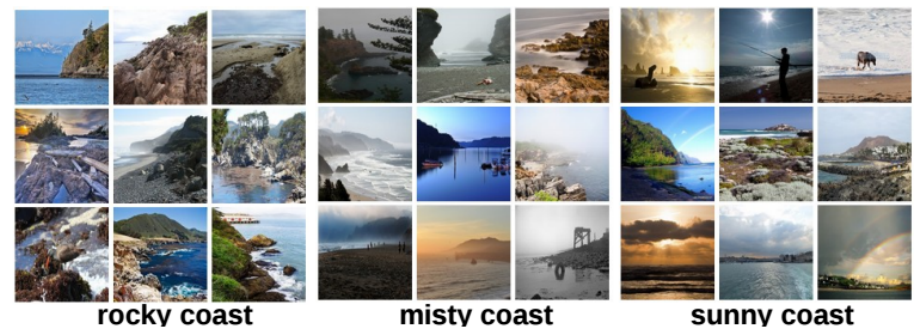


Filter feeding

[Visipedia Project](#)

# Image categorization

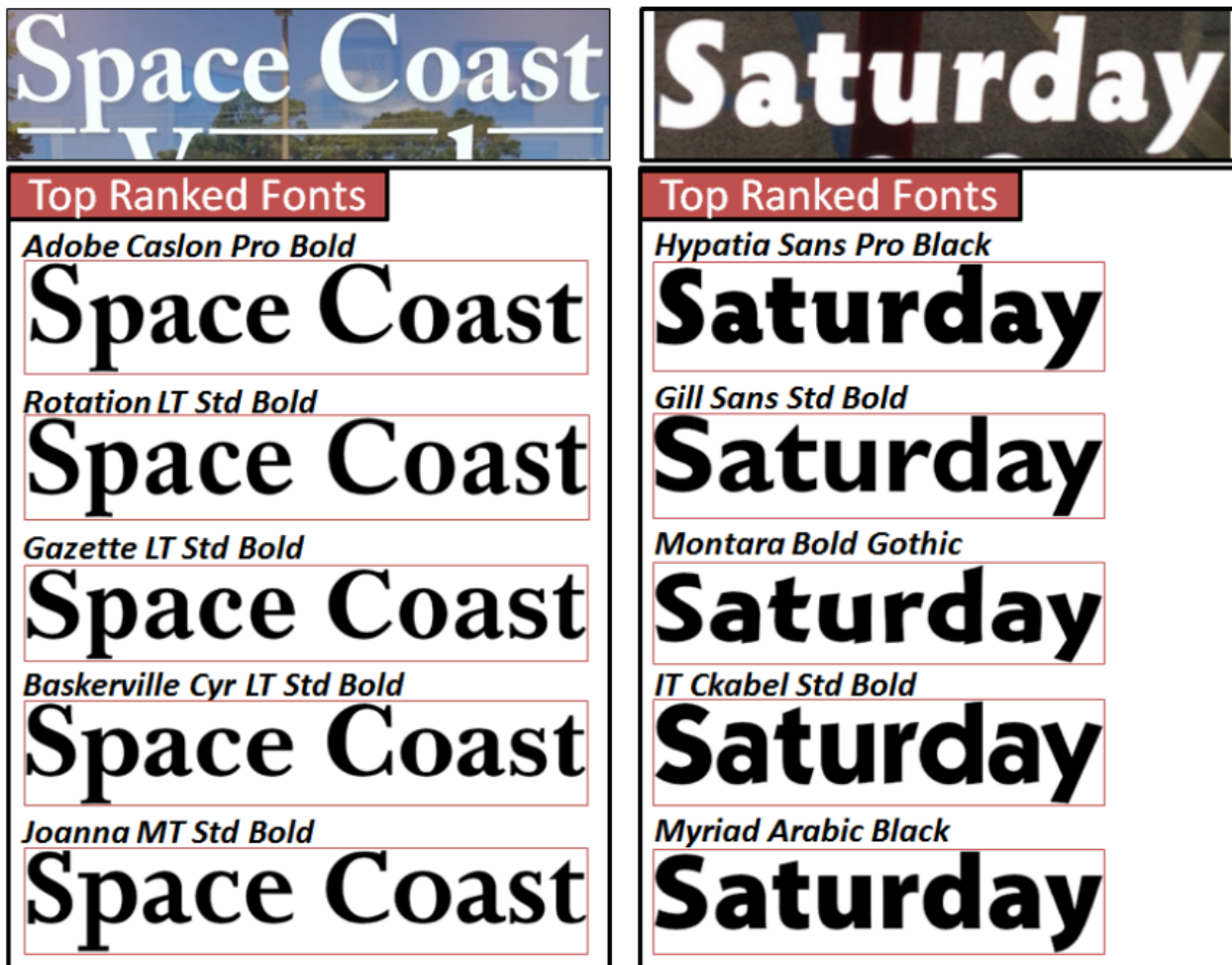
- Place recognition



Places Database [[Zhou et al. NIPS 2014](#)]

# Image categorization

- Visual font recognition

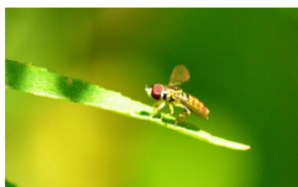


# Image categorization

- Image style recognition



HDR



Macro



Baroque



Roccoco



Vintage



Noir



Northern Renaissance



Cubism



Minimal



Hazy



Impressionism



Post-Impressionism



Long Exposure



Romantic



Abs. Expressionism



Color Field Painting

Flickr Style: 80K images covering 20 styles.

Wikipaintings: 85K images for 25 art genres.

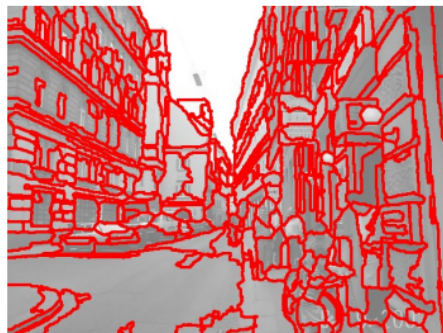
[[Karayev et al. BMVC 2014](#)]

# Region categorization

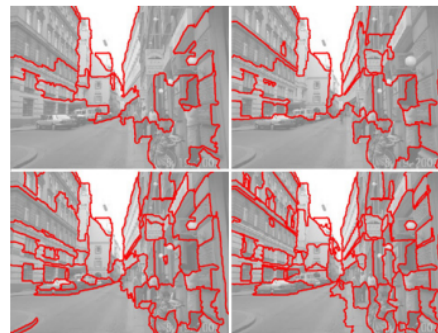
- Layout prediction



Input



Superpixels



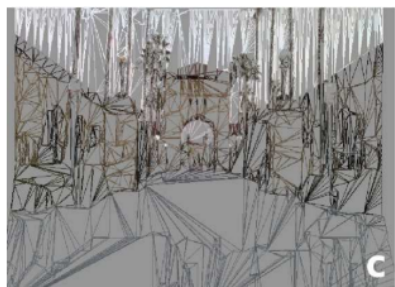
Multiple Segmentations



Surface Layout

Assign regions to orientation

Geometric context [[Hoiem et al. IJCV 2007](#)]



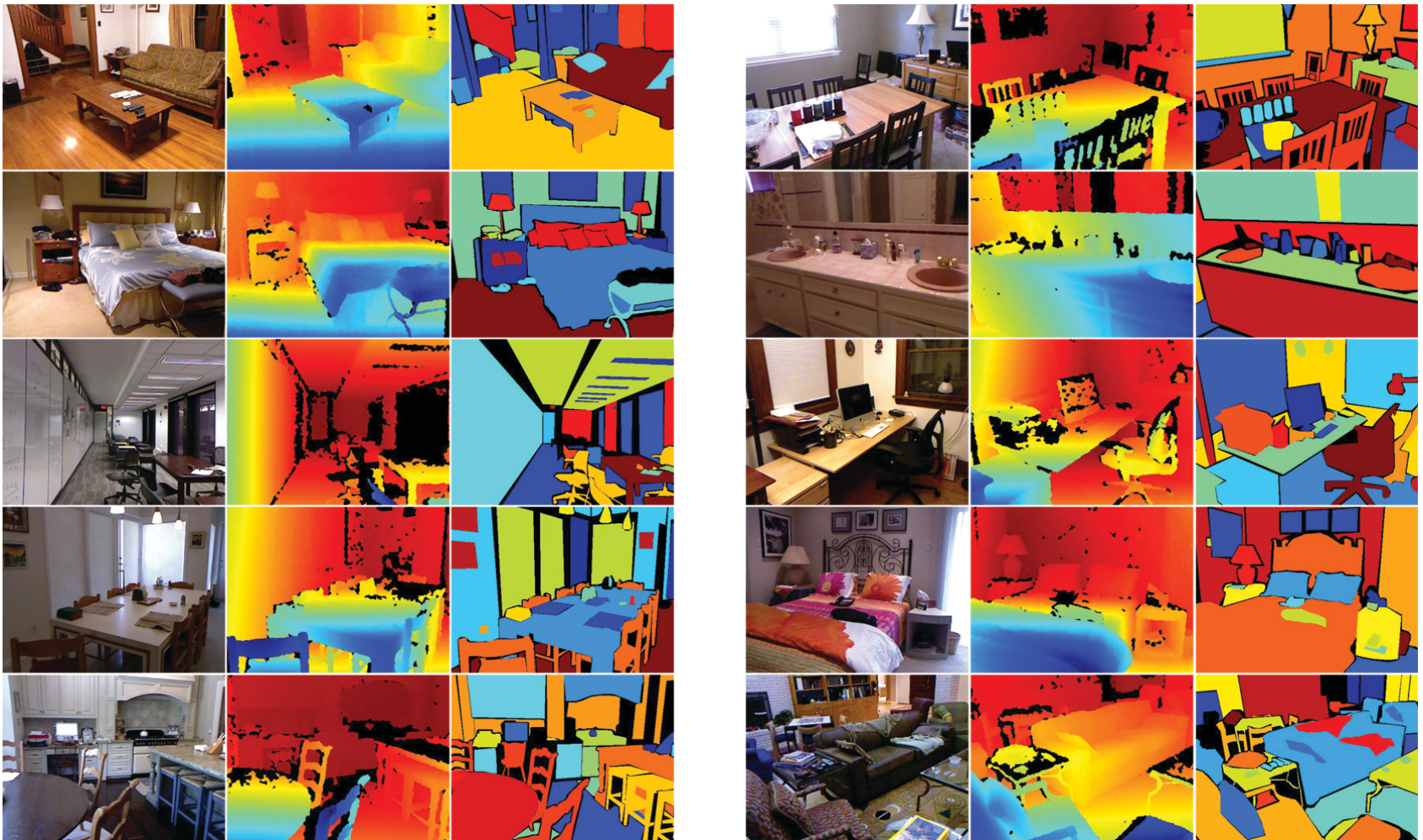
Assign regions to depth

Make3D [[Saxena et al. PAMI 2008](#)]



# Region categorization

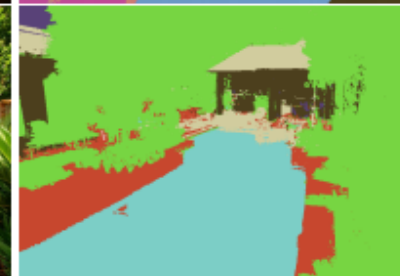
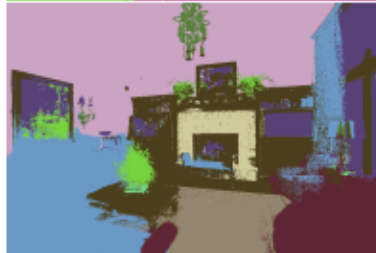
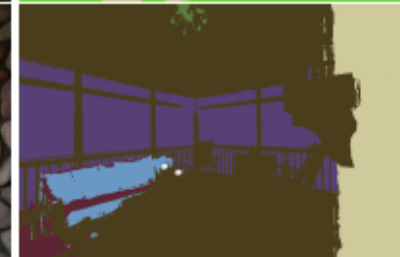
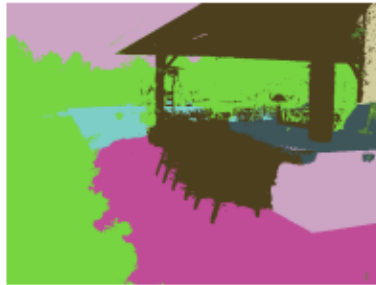
- Semantic segmentation from RGBD images



[Silberman et al. ECCV 2012]

# Region categorization

- Material recognition

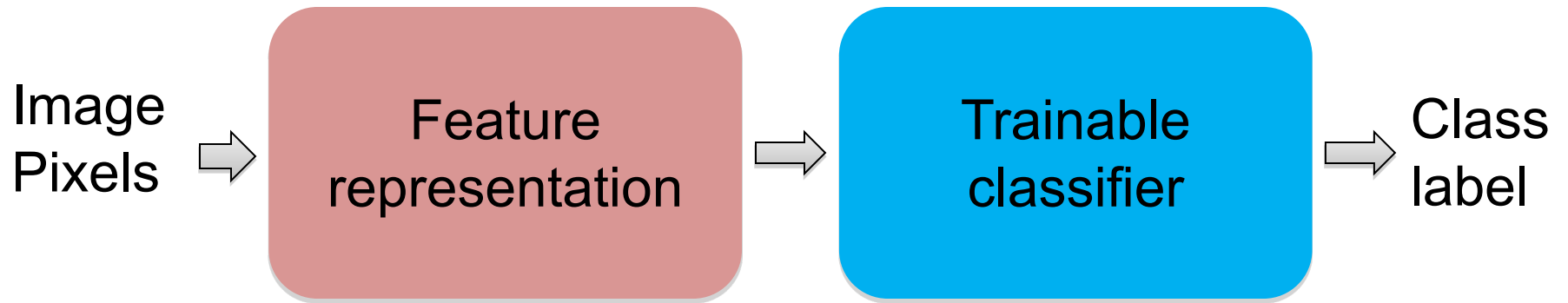


[Bell et al. CVPR 2015]

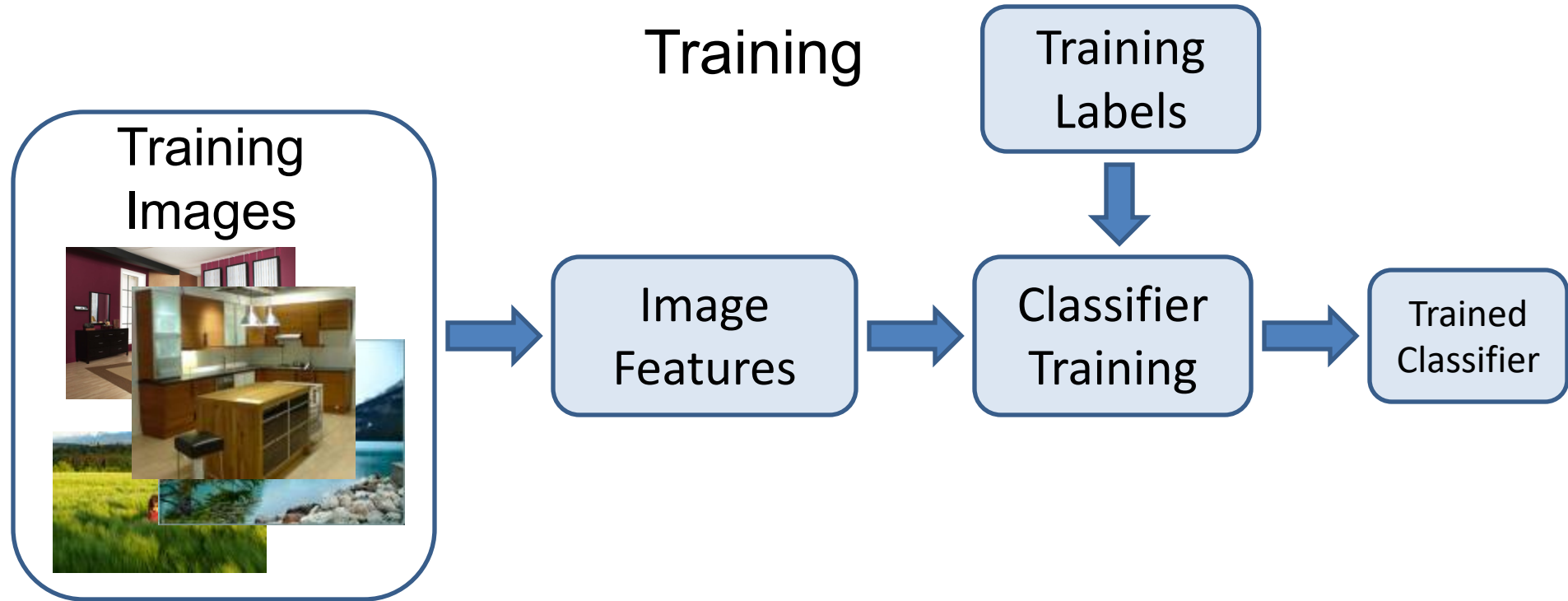
# Many vision problems involve categorization

- Image: Classify as indoor/outdoor, which room, what objects are there, etc.
- Object Detection: classify location (bounding box or region) as object or non-object
- Semantic Segmentation: classify pixel into an object, material, part, etc.
- Action Recognition: classify a frame or sequence into an action type
- ...

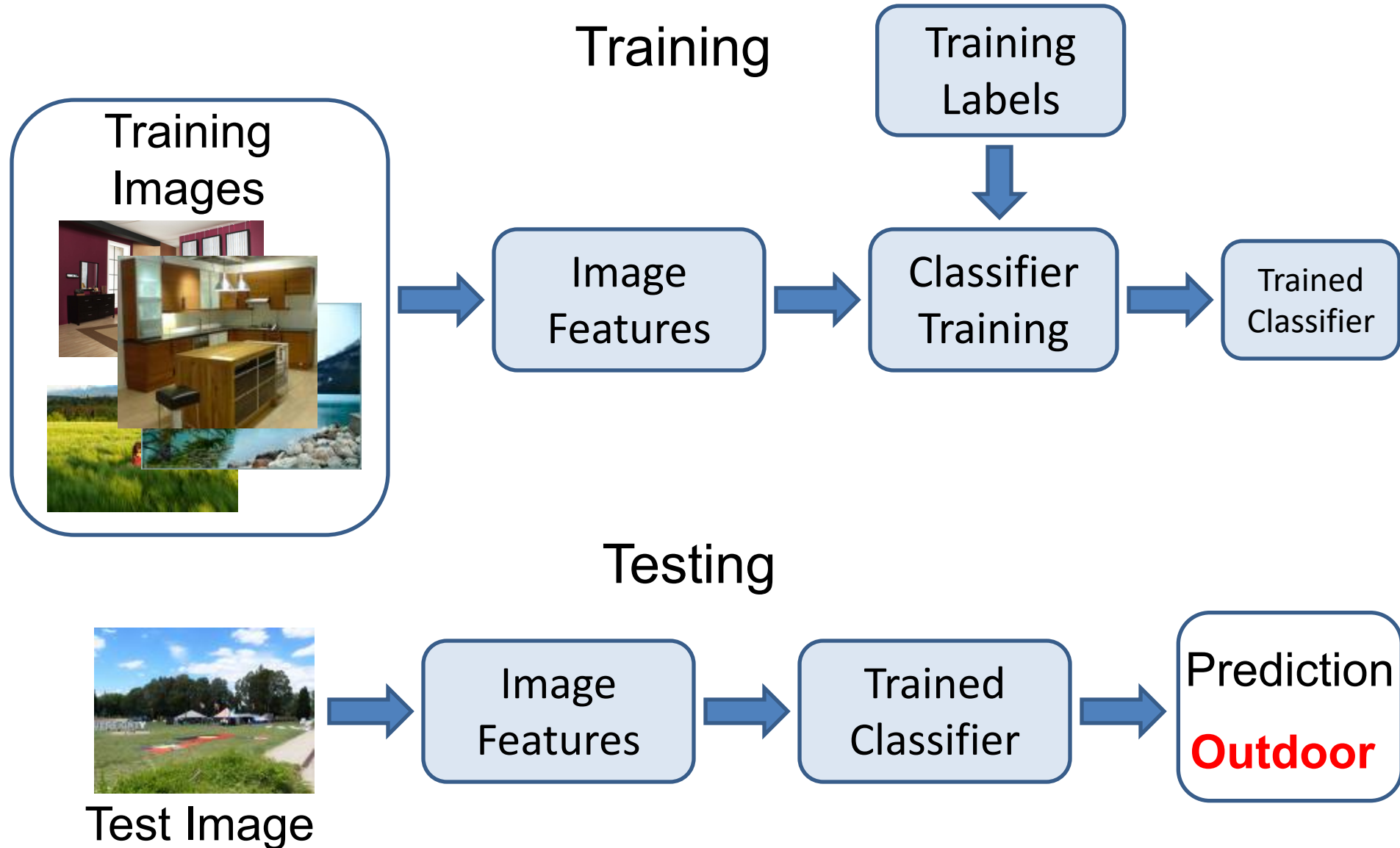
# “Classic” recognition pipeline



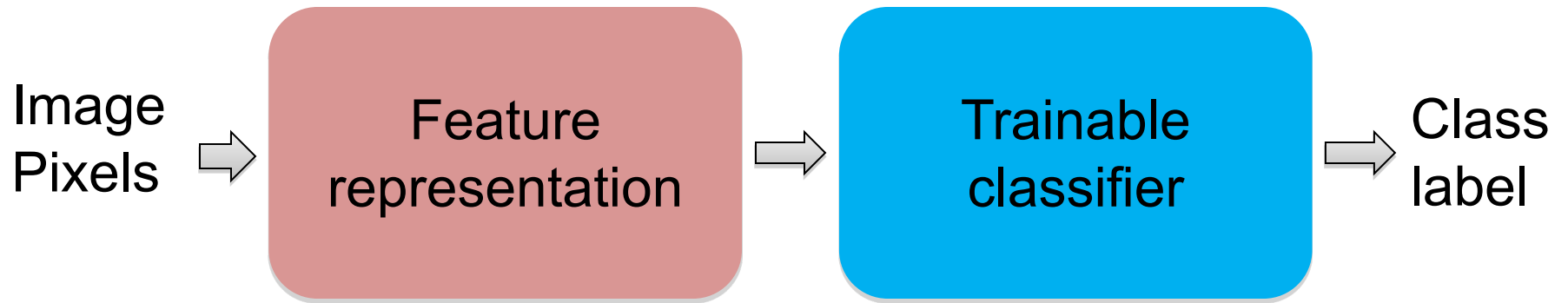
# Training phase



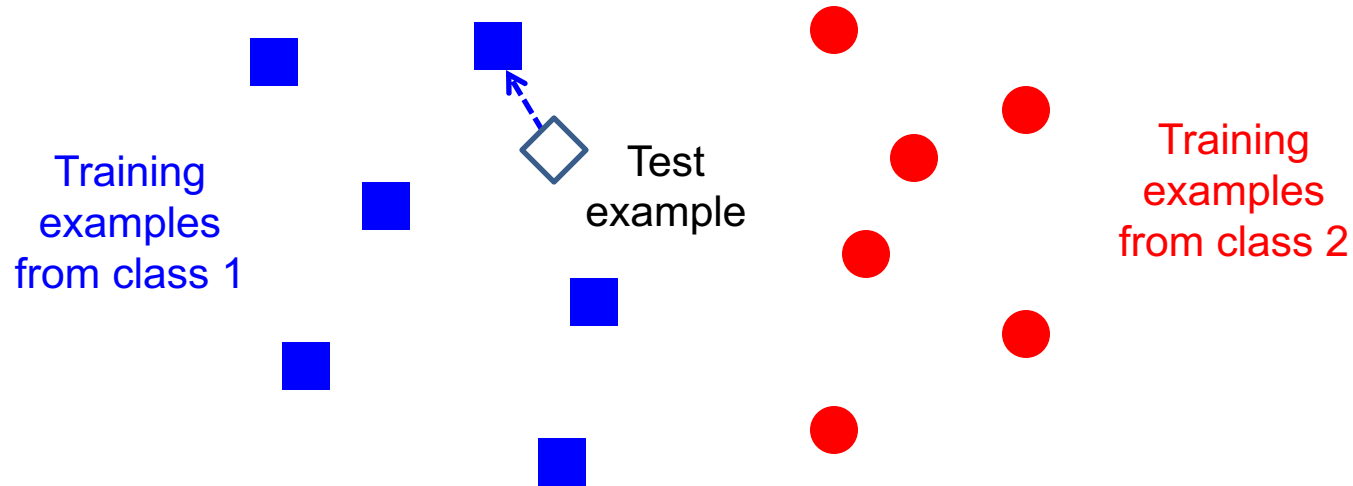
# Testing phase



# “Classic” recognition pipeline



# Classifiers: Nearest neighbor



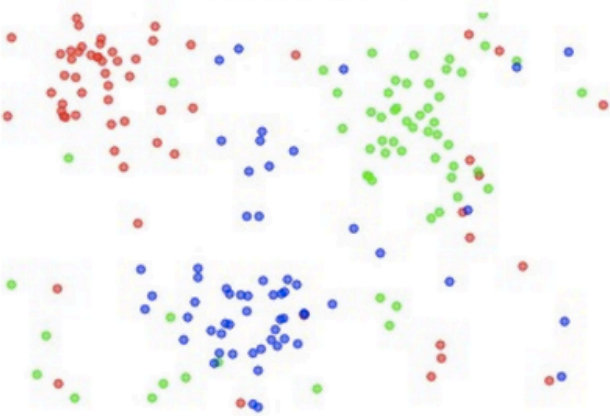
$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

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

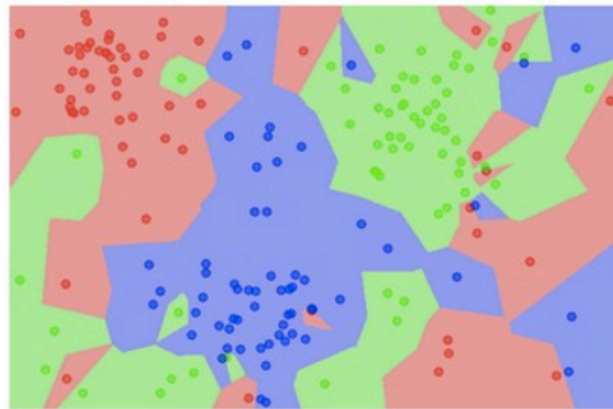


# K-nearest neighbor classifier

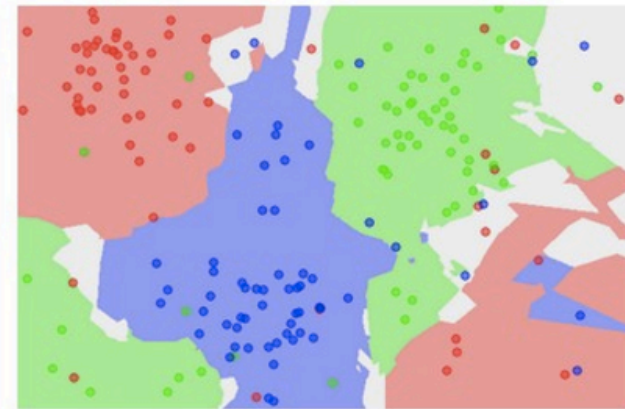
the data



NN classifier



5-NN classifier



- Which classifier is more robust to *outliers*?

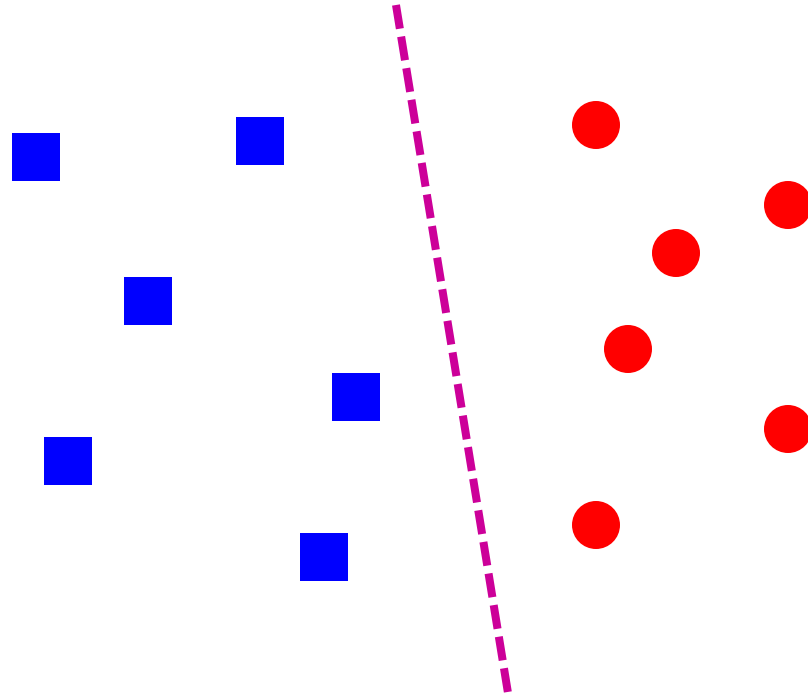
# K-nearest neighbor classifier



Left: Example images from the [CIFAR-10 dataset](#). Right: first column shows a few test images and next to each we show the top 10 nearest neighbors in the training set according to pixel-wise difference.

Credit: Andrej Karpathy, <http://cs231n.github.io/classification/>

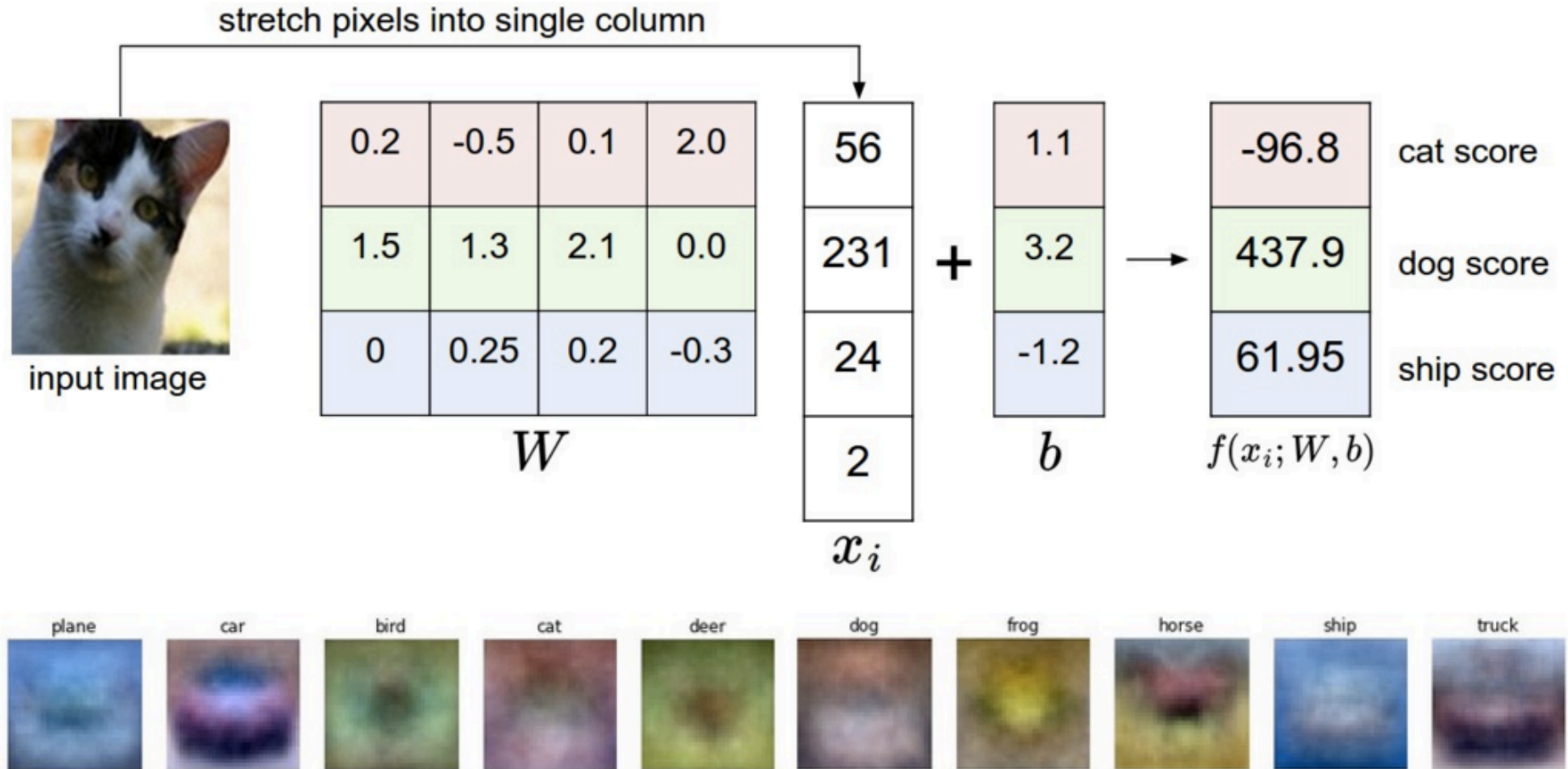
# Linear classifiers



- Find a *linear function* to separate the classes:

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

# Visualizing linear classifiers



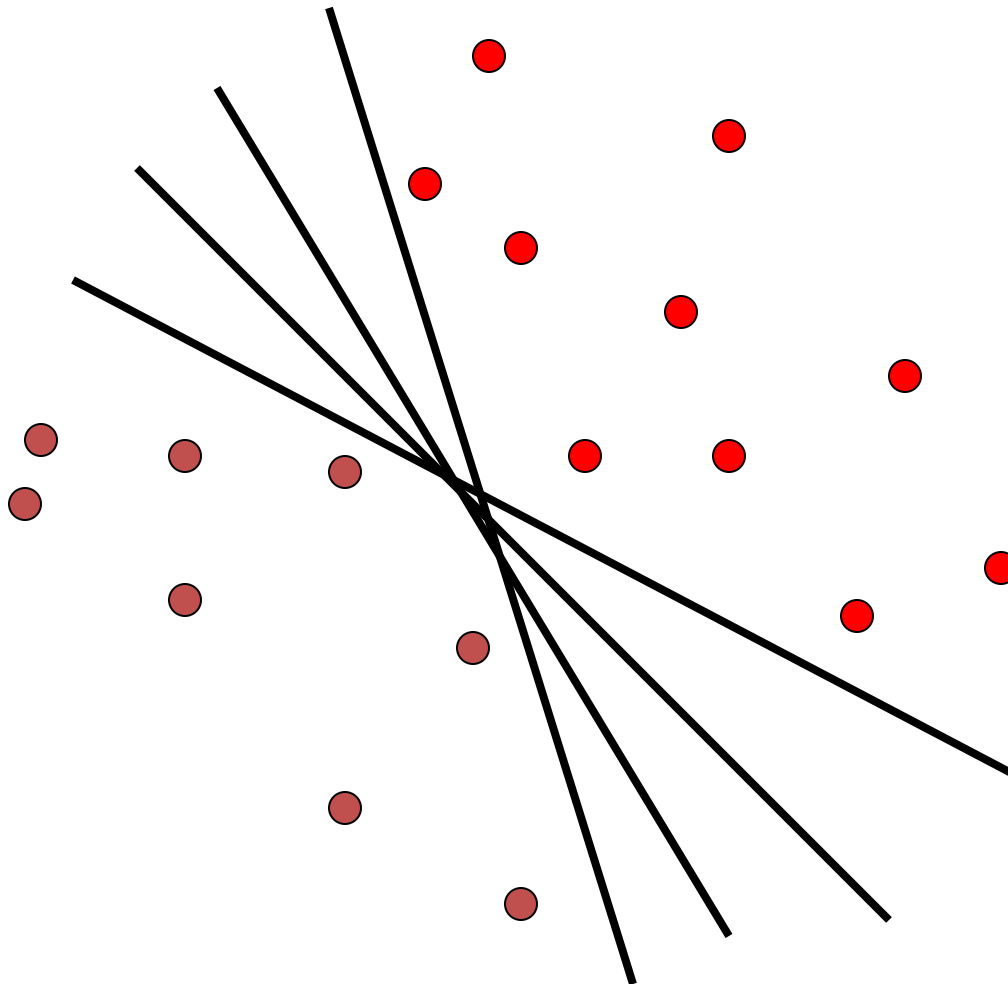
Source: Andrej Karpathy, <http://cs231n.github.io/linear-classify/>

# Nearest neighbor vs. linear classifiers

- **NN pros:**
  - Simple to implement
  - Decision boundaries not necessarily linear
  - Works for any number of classes
  - *Nonparametric* method
- **NN cons:**
  - Need good distance function
  - Slow at test time
- **Linear pros:**
  - Low-dimensional *parametric* representation
  - Very fast at test time
- **Linear cons:**
  - Works for two classes
  - How to train the linear function?
  - What if data is not linearly separable?

# Linear classifiers

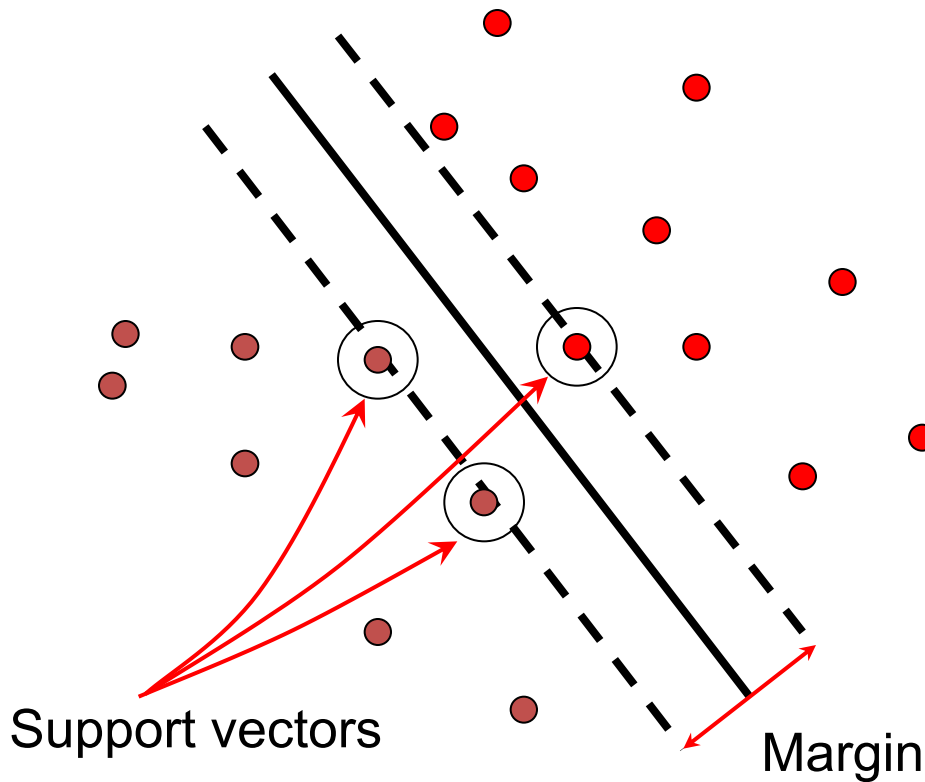
- When the data is linearly separable, there may be more than one separator (hyperplane)



Which separator  
is best?

# Support vector machines

- Find hyperplane that maximizes the *margin* between the positive and negative examples



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

$$\text{For support vectors, } \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$$

$$\text{Distance between point and hyperplane: } \frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

$$\text{Therefore, the margin is } 2 / \|\mathbf{w}\|$$

# Finding the maximum margin hyperplane

1. Maximize margin  $2 / \|\mathbf{w}\|$
2. Correctly classify all training data:
  - $\mathbf{x}_i$  positive ( $y_i = 1$ ):  $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$
  - $\mathbf{x}_i$  negative ( $y_i = -1$ ):  $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

- *Quadratic optimization problem:*

- $$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$



# SVM parameter learning

• Separable data:  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

Maximize  
margin

Classify training data correctly

• Non-separable data:

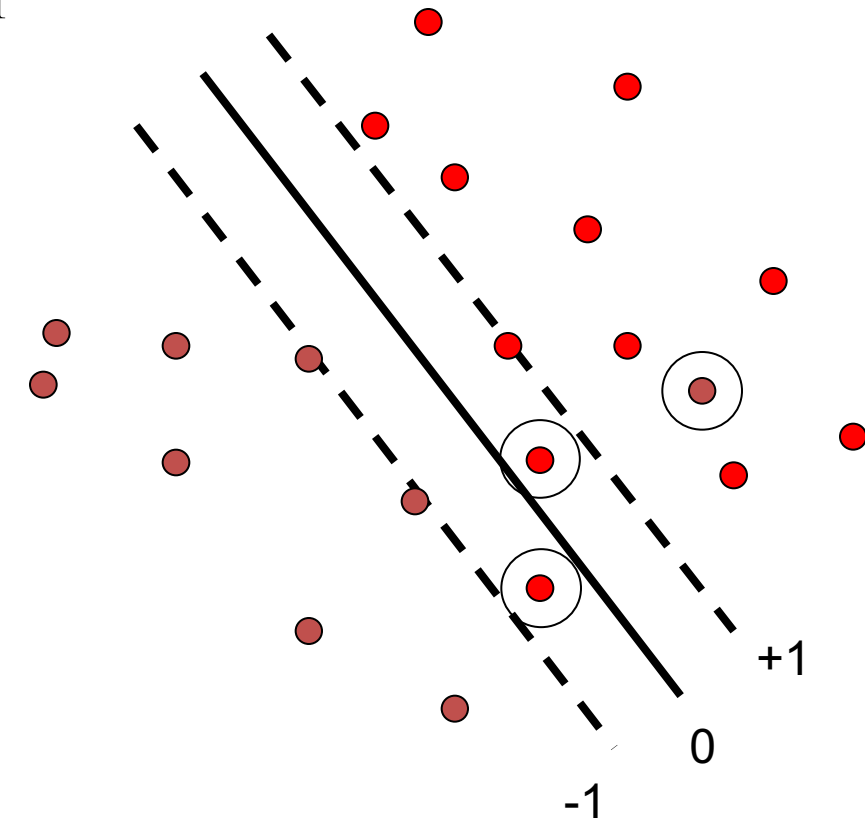
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i (\mathbf{w} \cdot \mathbf{x}_i + b))$$

Maximize  
margin

Minimize classification mistakes

# SVM parameter learning

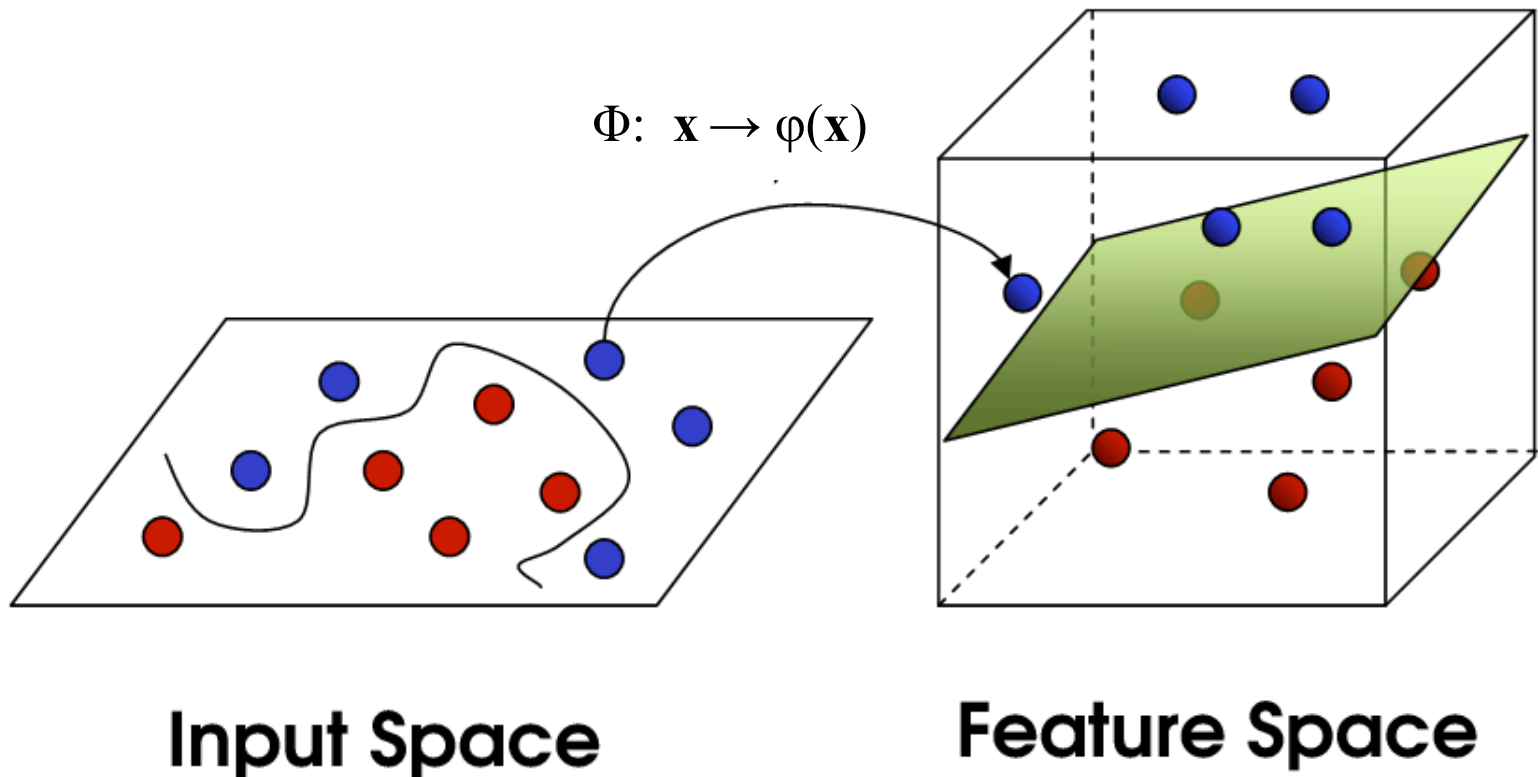
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$



- Demo: <http://cs.stanford.edu/people/karpathy/svmjs/demo>

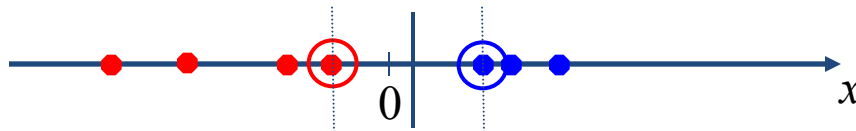
# Nonlinear SVMs

- **General idea:** the original input space can always be mapped to some higher-dimensional feature space where the training



# Nonlinear SVMs

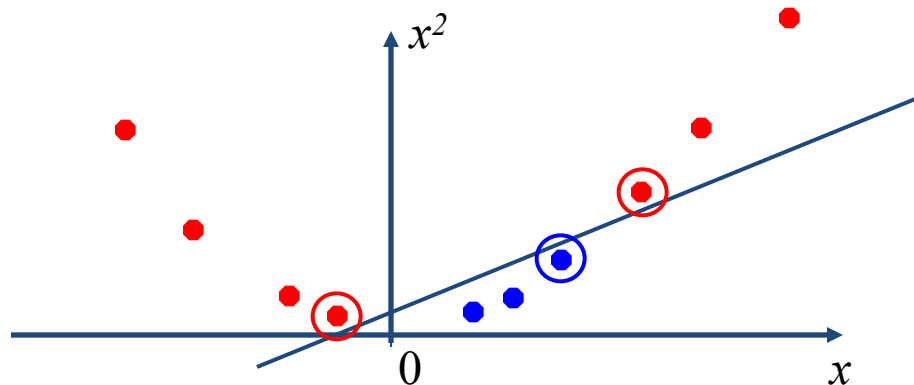
- Linearly separable dataset in 1D:



- Non-separable dataset in 1D:



- We can map the data to a *higher-dimensional space*:



# The kernel trick

- **General idea:** the original input space can always be mapped to some higher-dimensional feature space where the training set is separable
- **The kernel trick:** instead of explicitly computing the lifting transformation  $\varphi(\mathbf{x})$ , define a kernel function  $K$  such that

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

- (to be valid, the kernel function must satisfy *Mercer's condition*)

# The kernel trick

- Linear SVM decision function:

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

learned  
weight

Support  
vector

# The kernel trick

- Linear SVM decision function:

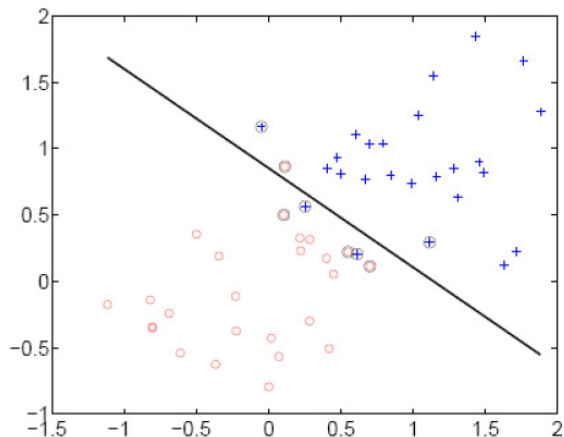
$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

- Kernel SVM decision function:

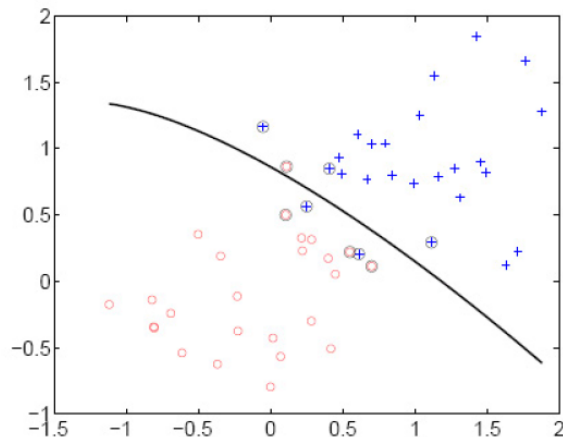
$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

- This gives a nonlinear decision boundary in the original feature space

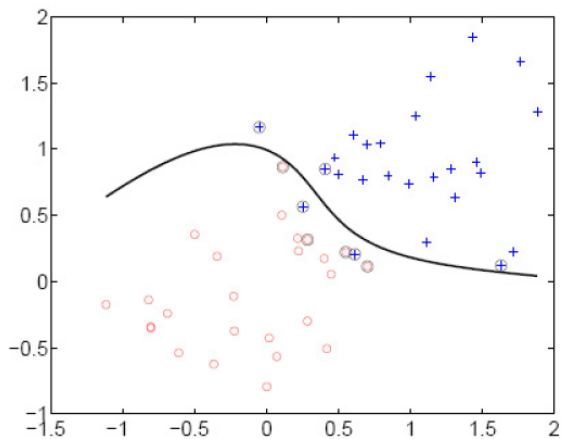
# Polynomial kernel: $K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x} \cdot \mathbf{y})^d$



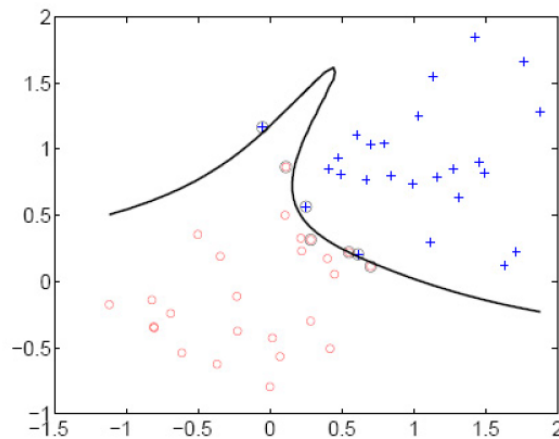
linear



$2^{nd}$  order polynomial



$4^{th}$  order polynomial



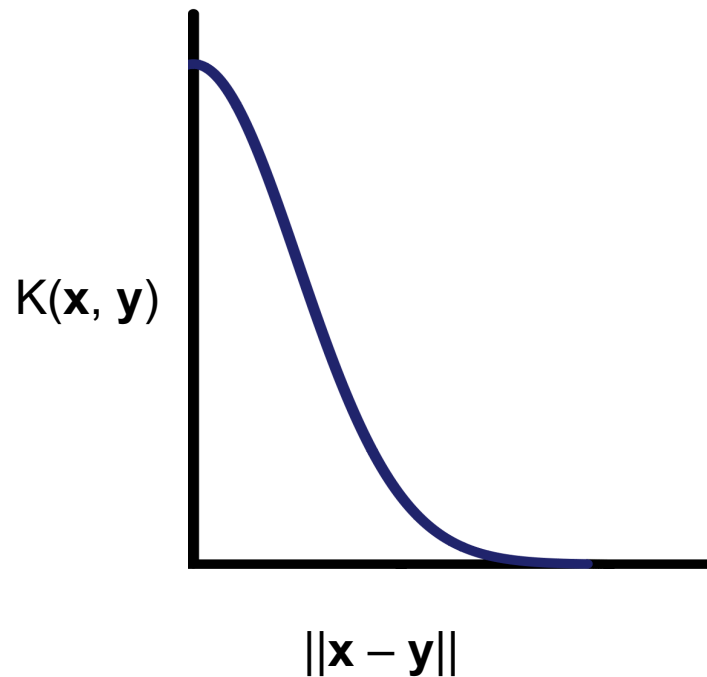
$8^{th}$  order polynomial



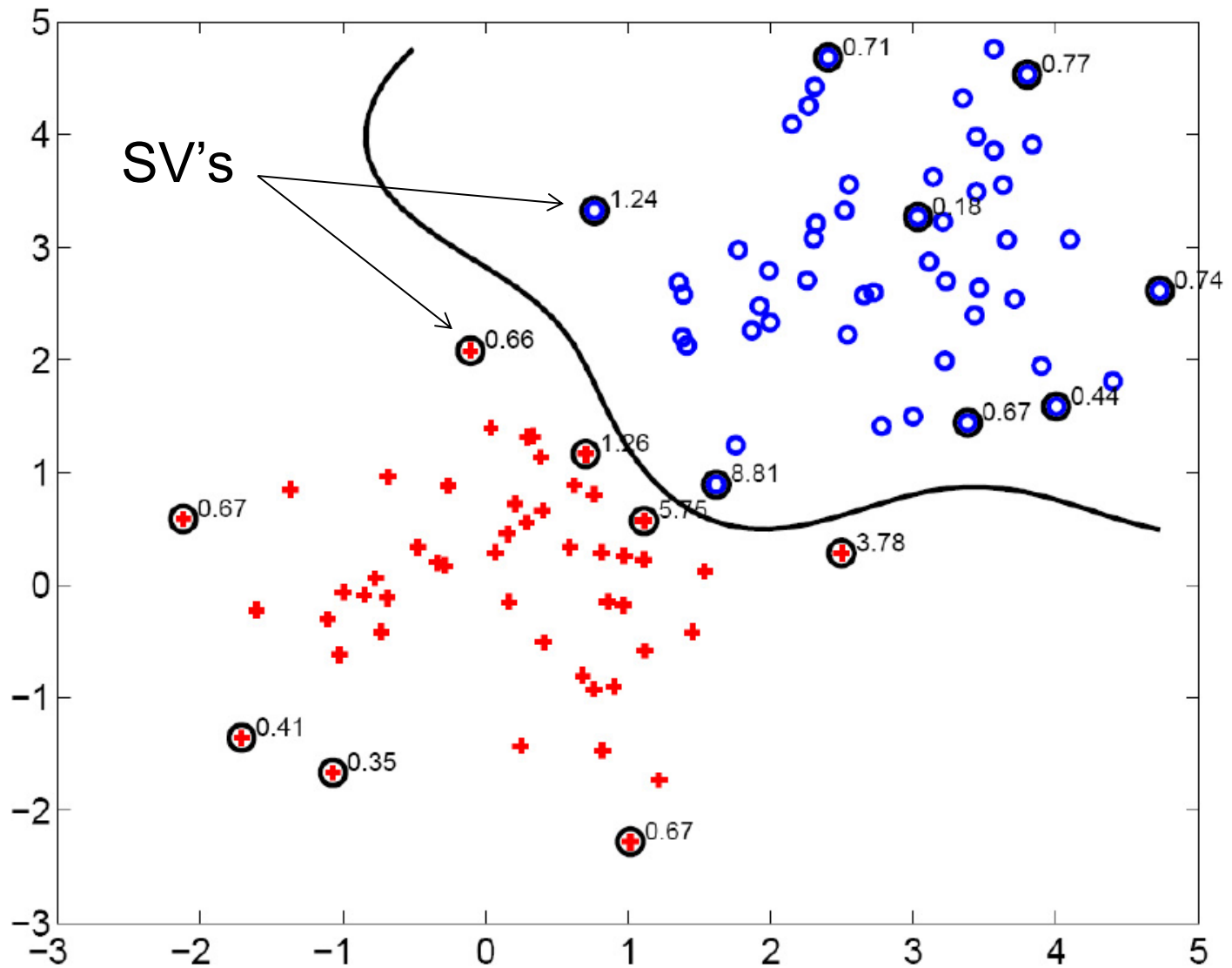
# Gaussian kernel

- Also known as the radial basis function (RBF) kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right)$$



# Gaussian kernel



# SVMs: Pros and cons

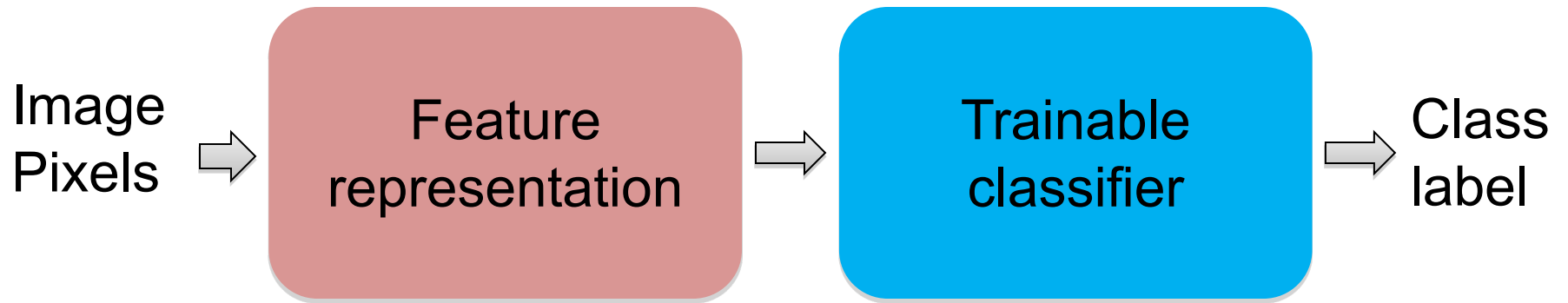
- Pros

- Kernel-based framework is very powerful, flexible
- Training is convex optimization, globally optimal solution can be found
- Amenable to theoretical analysis
- SVMs work very well in practice, even with very small training sample sizes

- Cons

- No “direct” multi-class SVM, must combine two-class SVMs (e.g., with one-vs-others)
- Computation, memory (esp. for nonlinear SVMs)

# “Classic” recognition pipeline



# Building an Image Classifier

Training Data

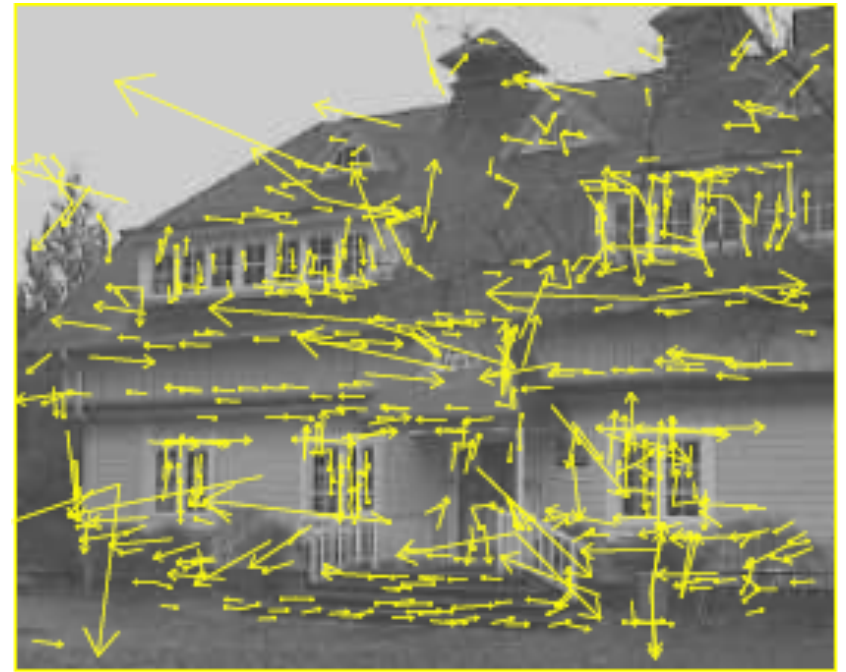


Test Data



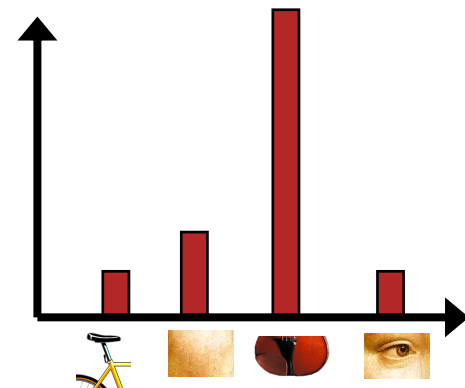
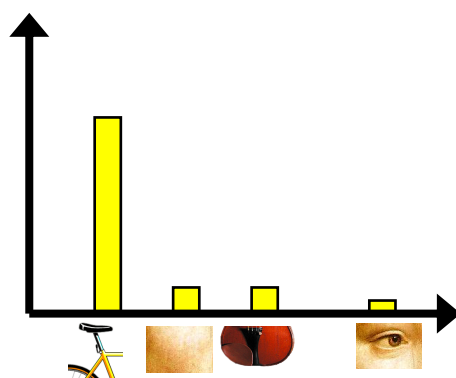
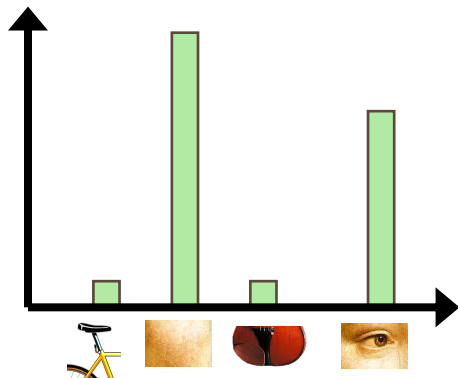
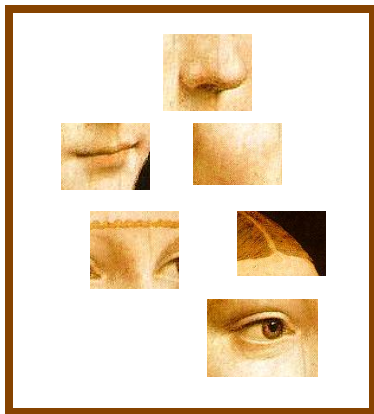
# Building an Image Classifier

- Raw Pixels?
- SIFT features at keypoints?



# Bag of features: Outline

1. Extract local features
2. Learn “visual vocabulary”
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of “visual words”



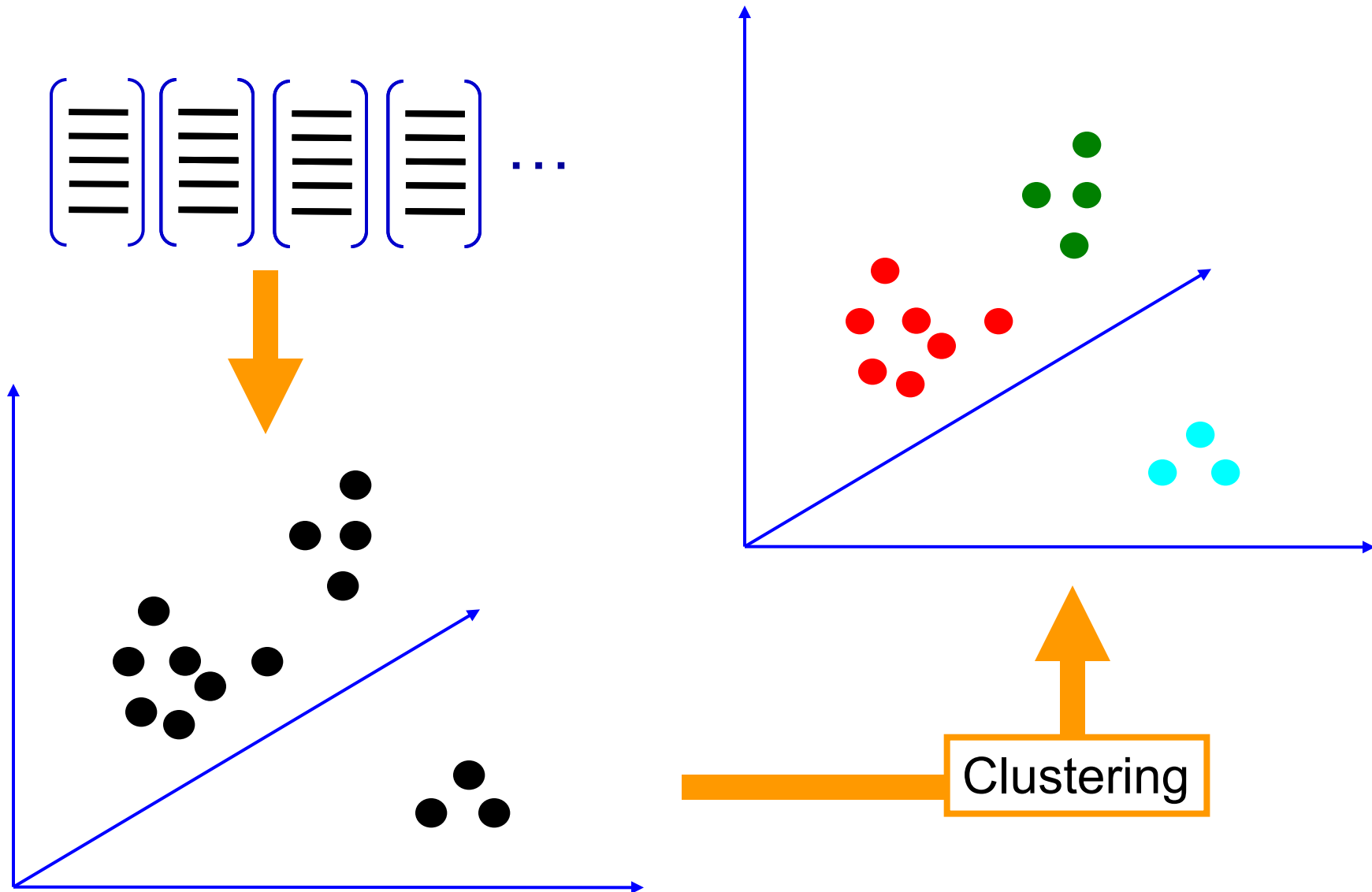
# 1. Local feature extraction

- Sample patches and extract descriptors

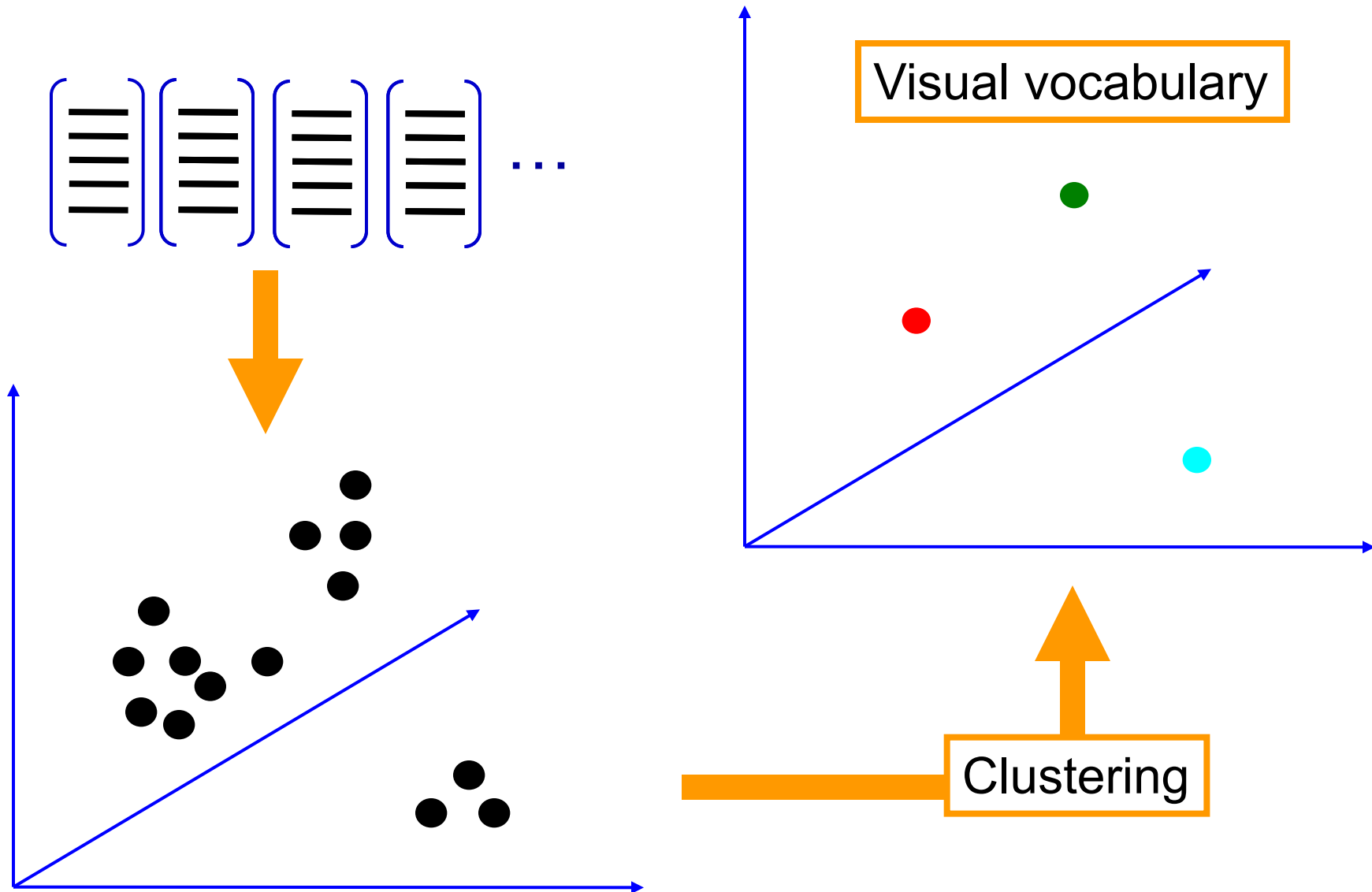




## 2. Learning the visual vocabulary



## 2. Learning the visual vocabulary



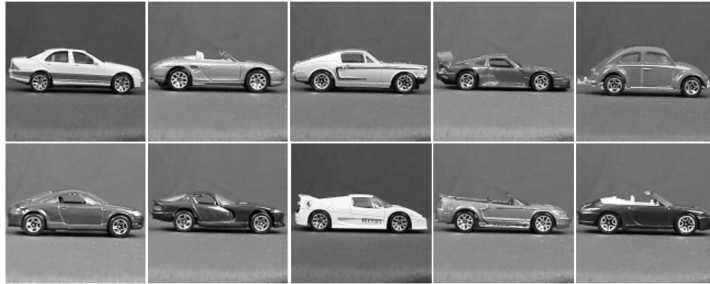
# K-means clustering

- Want to minimize sum of squared Euclidean distances between features  $\mathbf{x}_i$  and their nearest cluster centers  $\mathbf{m}_k$

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\text{point } i \text{ in cluster } k} (\mathbf{x}_i - \mathbf{m}_k)^2$$

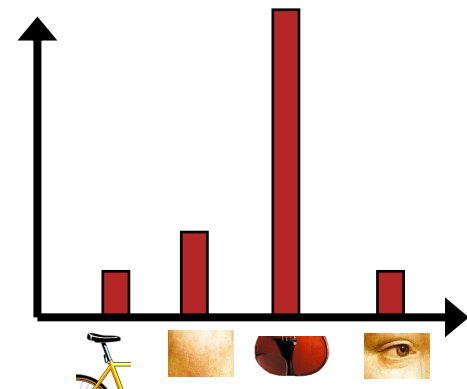
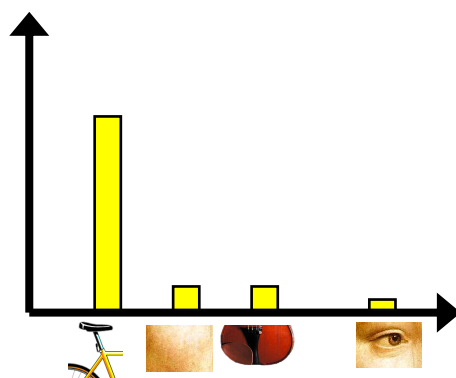
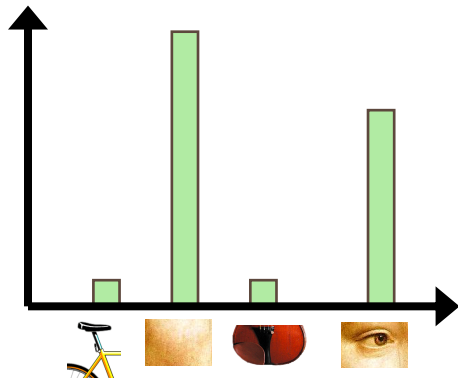
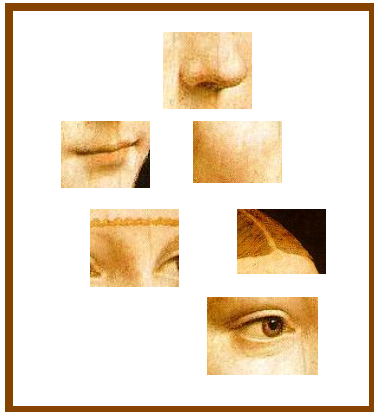
- Algorithm:
- Randomly initialize K cluster centers
- Iterate until convergence:
  - Assign each feature to the nearest center
  - Recompute each cluster center as the mean of all features assigned to it

# Visual vocabularies



# Bag of features: Outline

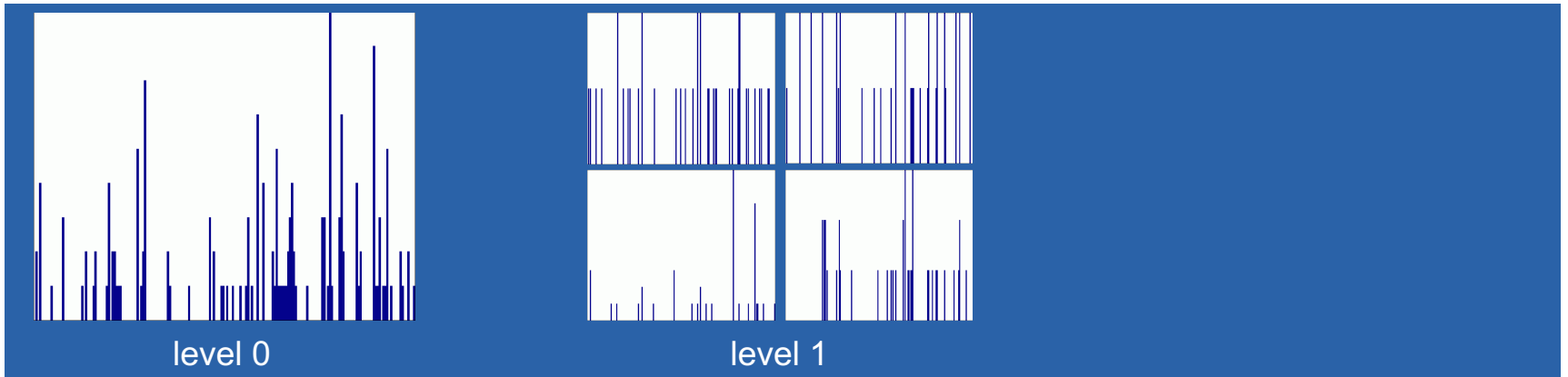
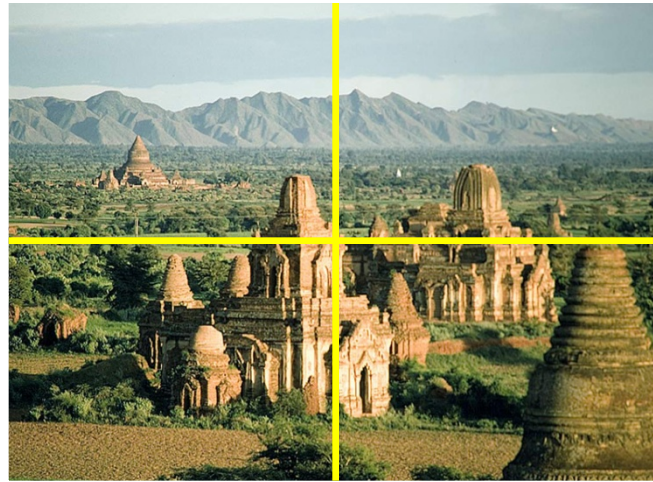
1. Extract local features
2. Learn “visual vocabulary”
3. **Quantize local features using visual vocabulary**
4. **Represent images by frequencies of “visual words”**



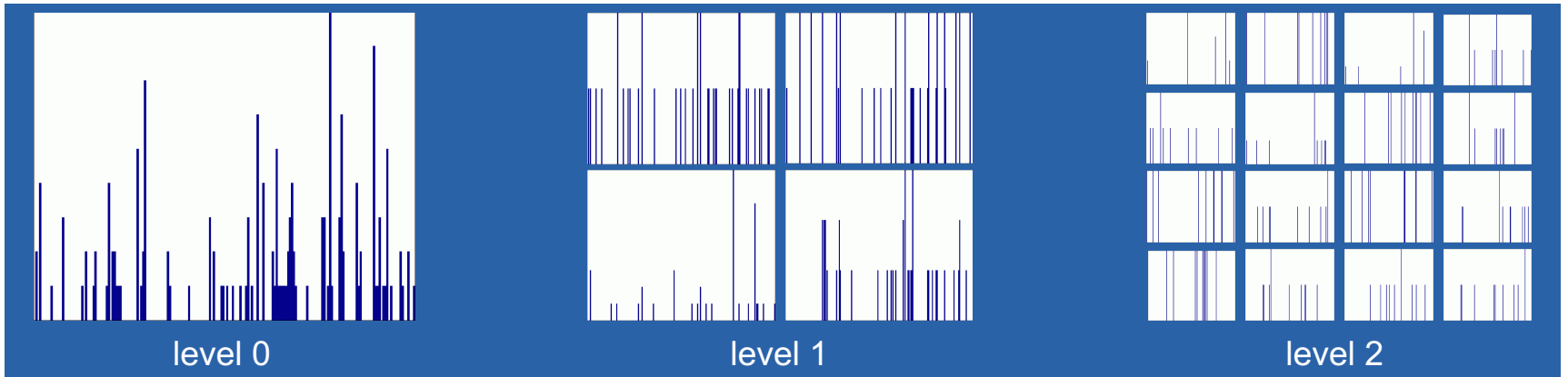
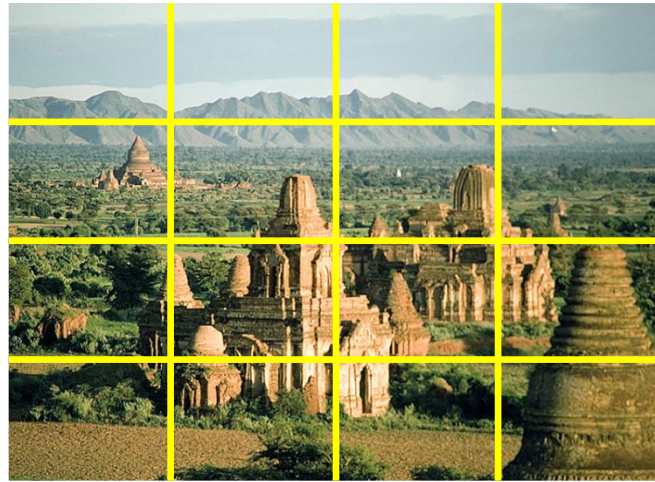
# Spatial pyramids



# Spatial pyramids



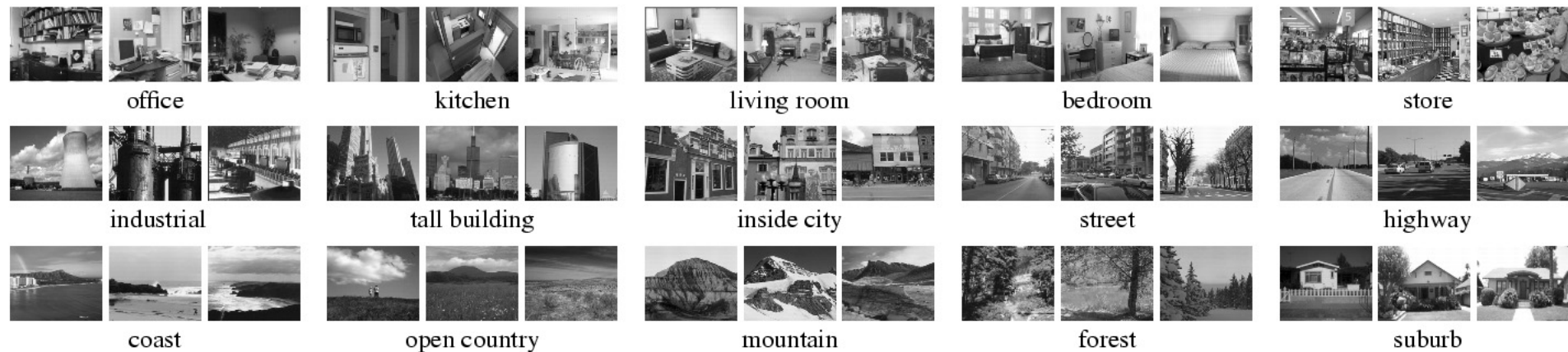
# Spatial pyramids





# Spatial pyramids

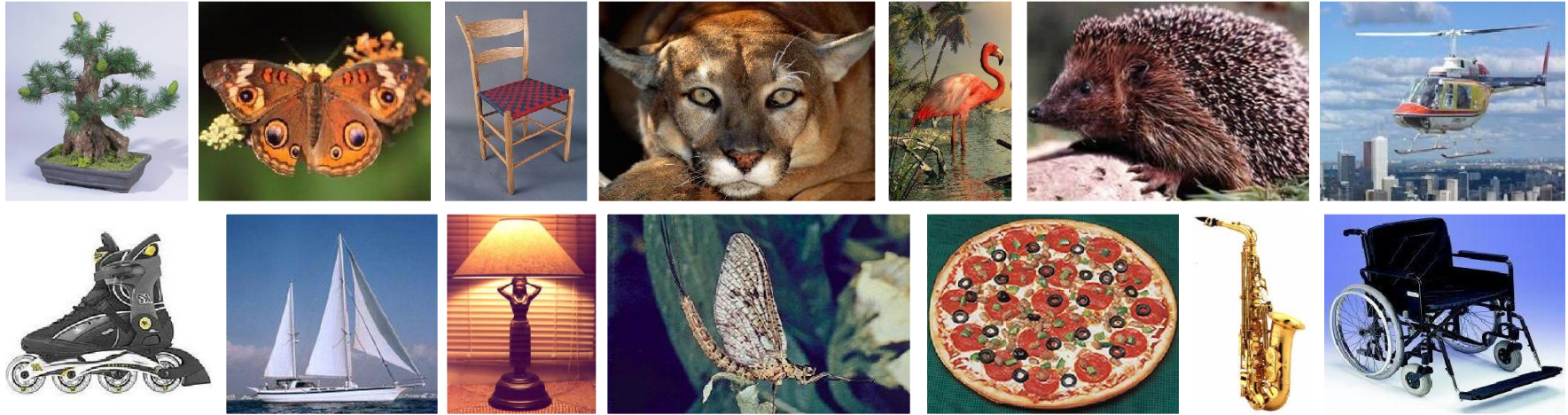
- Scene classification results



Level	Weak features (vocabulary size: 16)		Strong features (vocabulary size: 200)	
	Single-level	Pyramid	Single-level	Pyramid
0 (1 × 1)	45.3 ±0.5		72.2 ±0.6	
1 (2 × 2)	53.6 ±0.3	56.2 ±0.6	77.9 ±0.6	79.0 ±0.5
2 (4 × 4)	61.7 ±0.6	64.7 ±0.7	79.4 ±0.3	<b>81.1 ±0.3</b>
3 (8 × 8)	63.3 ±0.8	<b>66.8 ±0.6</b>	77.2 ±0.4	80.7 ±0.3

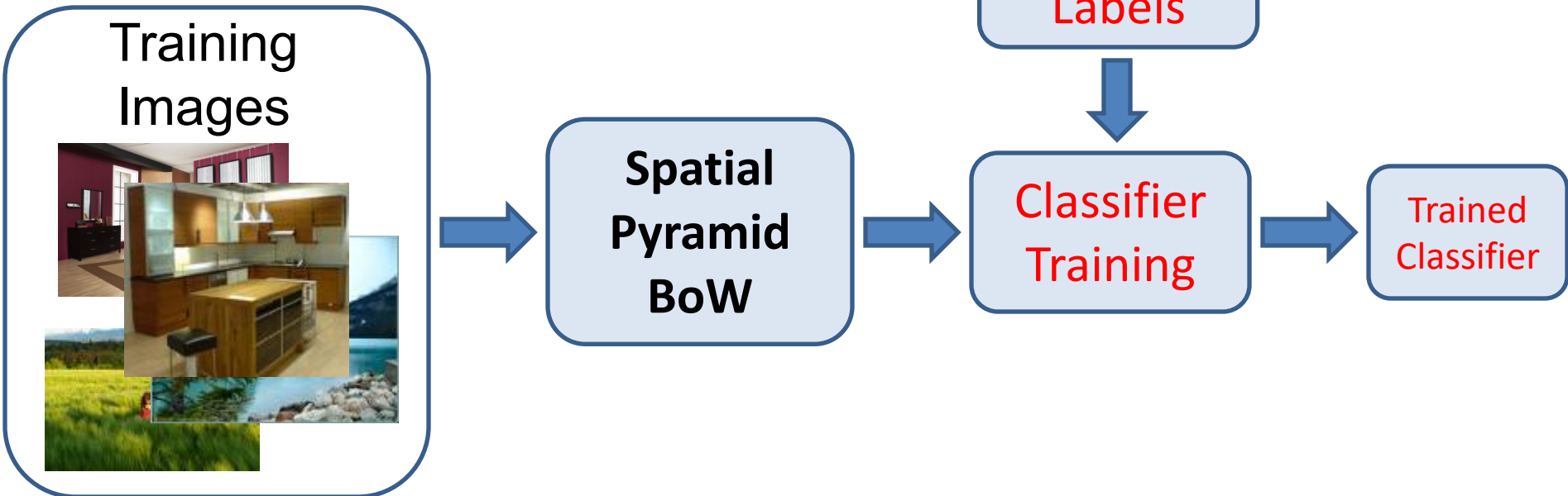
# Spatial pyramids

- Caltech101 classification results

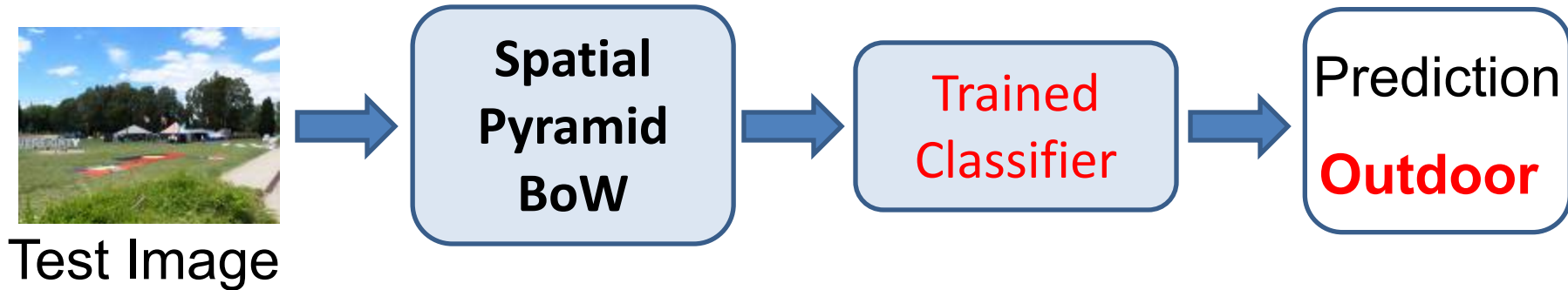


	Weak features (16)		Strong features (200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0	15.5 ±0.9		41.2 ±1.2	
1	31.4 ±1.2	32.8 ±1.3	55.9 ±0.9	57.0 ±0.8
2	47.2 ±1.1	49.3 ±1.4	63.6 ±0.9	<b>64.6</b> ±0.8
3	52.2 ±0.8	<b>54.0</b> ±1.1	60.3 ±0.9	64.6 ±0.7

## Training

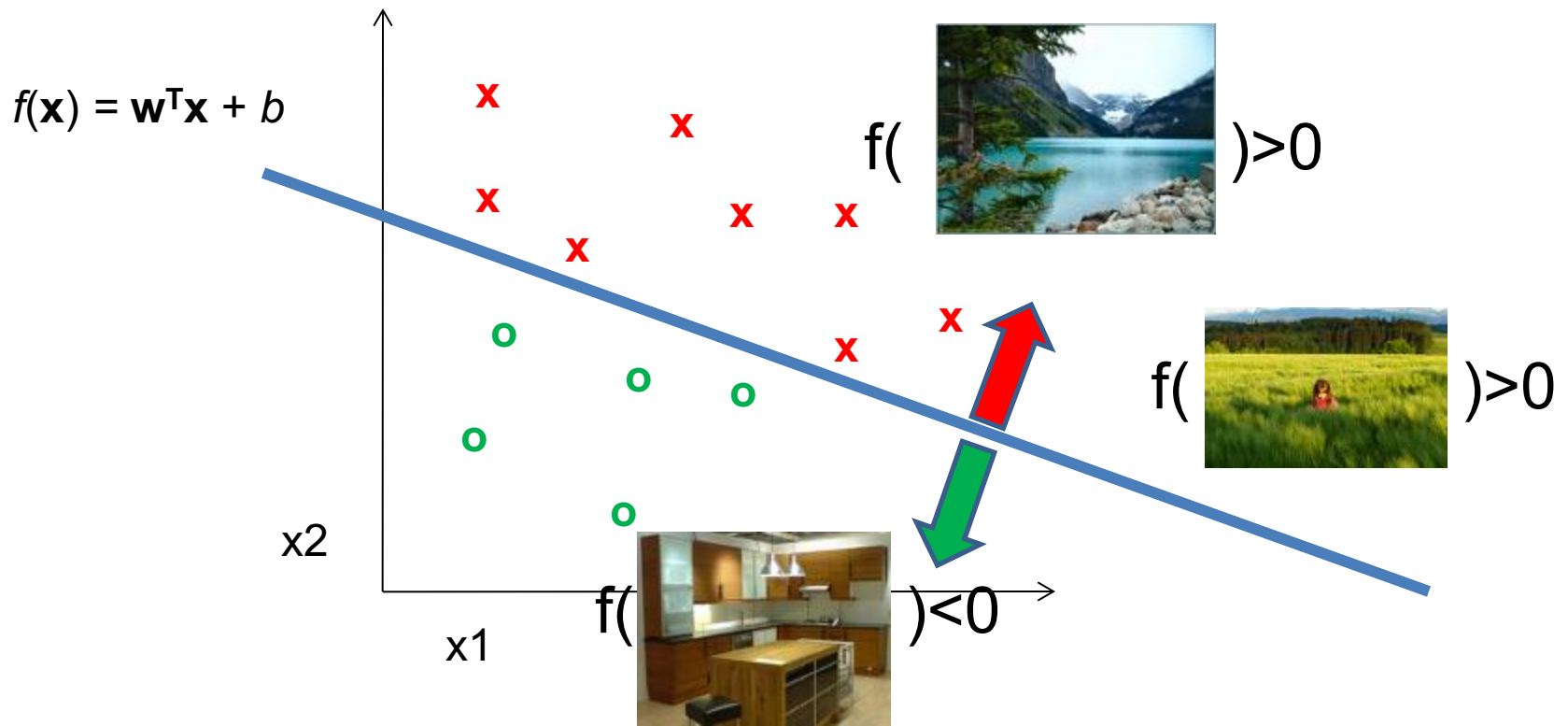


## Testing



# Linear SVM classifier

Find the hyperplane that separate examples of different categories



# Spatial Pyramids Results



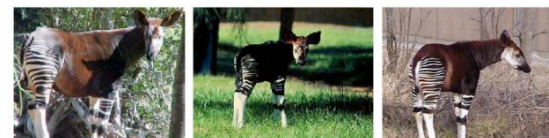
minaret (97.6%)



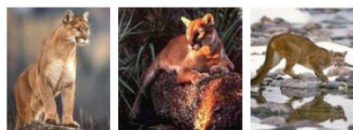
windsor chair (94.6%)



joshua tree (87.9%)



okapi (87.8%)



cougar body (27.6%)



beaver (27.5%)



crocodile (25.0%)



ant (25.0%)

	Weak features		Strong features (200)	
$L$	Single-level	Pyramid	Single-level	Pyramid
0	15.5 $\pm$ 0.9		41.2 $\pm$ 1.2	
1	31.4 $\pm$ 1.2	32.8 $\pm$ 1.3	55.9 $\pm$ 0.9	57.0 $\pm$ 0.8
2	47.2 $\pm$ 1.1	49.3 $\pm$ 1.4	63.6 $\pm$ 0.9	<b>64.6</b> $\pm$ 0.8
3	52.2 $\pm$ 0.8	<b>54.0</b> $\pm$ 1.1	60.3 $\pm$ 0.9	64.6 $\pm$ 0.7

Table 2. Classification results for the Caltech-101 database.

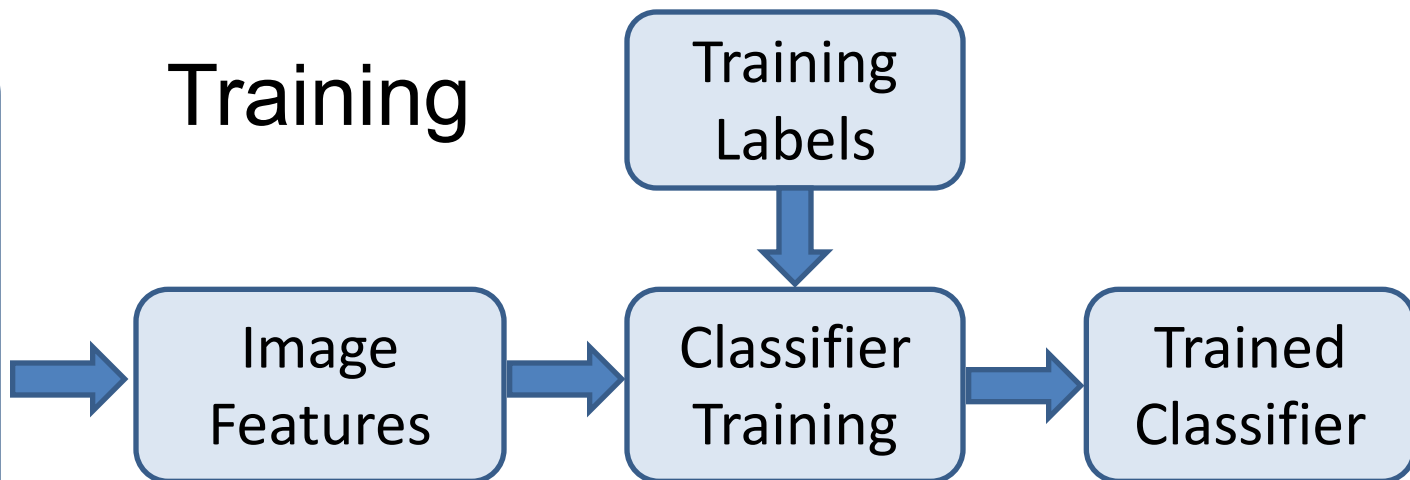
Now let's go back to the core concepts of image categorization in general

# Categorization involves **features** and a classifier

Training Images



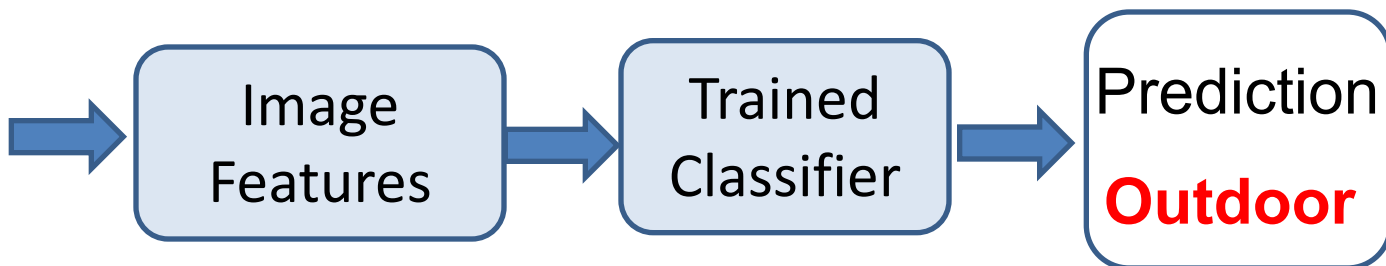
## Training



## Testing



Test Image



# Q: What are good features for...

- recognizing a beach?





# Q: What are good features for...

- recognizing cloth fabric?



# Q: What are good features for...

- recognizing a mug?



# What are the right features?

Depend on what you want to know!

- Object: shape
  - Local shape info, shading, shadows, texture
- Scene : geometric layout
  - linear perspective, gradients, line segments
- Material properties: albedo, feel, hardness
  - Color, texture
- Action: motion
  - Optical flow, tracked points

# General principles of features

- Coverage
  - Ensure that all relevant info is captured
- Concision
  - Minimize number of features without sacrificing coverage
- Directness
  - Ideal features are independently useful for prediction

It's hard to design good features for a specific problem

Many machine learning algorithms solve for both feature and classifier parameters, such as

- Decision trees
- Random forests
- Multilayer neural networks (including CNNs)

# Classifiers

Goal: From labeled training samples, learn parameters of a scoring/decision function that is likely to predict the correct label on test samples

Typical assumptions:

- Training and test samples drawn from same distribution (i.i.d.)
- Training labels are correct

# Many classifiers to choose from

- SVM
- Neural networks
- Naïve Bayes
- Bayesian network
- Logistic regression
- Randomized Forests
- Boosted Decision Trees
- K-nearest neighbor
- RBMs
- Deep networks
- Etc.

Which is the best one?

# Classifiers: three main options

- **Nearest neighbor:** take a vote from K closest neighbors
  - Can learn features or distance measure
- **Linear:** score is linear combination of features
  - SVM, perceptron, naïve bayes, logistic regression
  - Others learn features and then apply linear classifier (e.g. deep network, random forest)
- **Structured prediction:** score an interdependent set of labels
  - E.g. label body part positions
  - Structured SVM, CNN, graphical model algorithms



# Generalization Theory

It's not enough to do well on the training set:  
also should make good predictions for new  
examples

# No Free Lunch Theorem

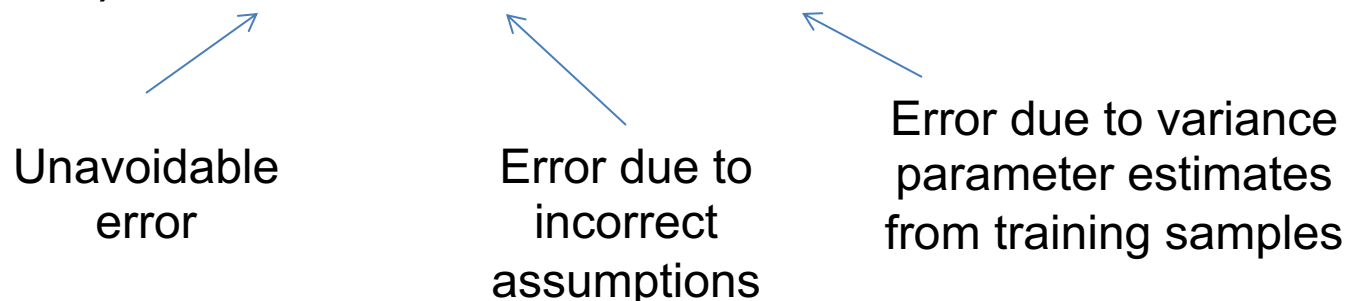
© Original Artist  
Reproduction rights obtainable from  
[www.CartoonStock.com](http://www.CartoonStock.com)



# Bias-Variance Trade-off

$$E(\text{MSE}) = \text{noise}^2 + \text{bias}^2 + \text{variance}$$

Unavoidable  
error



Error due to  
incorrect  
assumptions

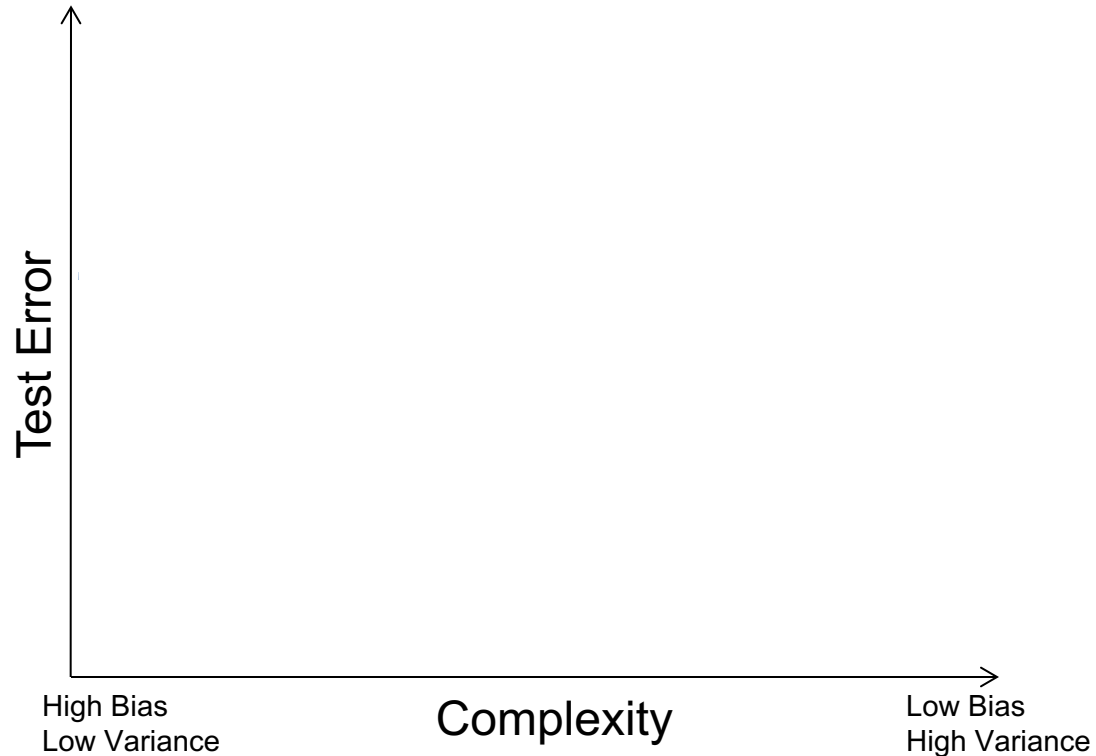
Error due to variance  
parameter estimates  
from training samples

See the following for explanation of bias-variance (also Bishop's "Neural Networks" book):

- <http://www.inf.ed.ac.uk/teaching/courses/mlsc/Notes/Lecture4/BiasVariance.pdf>

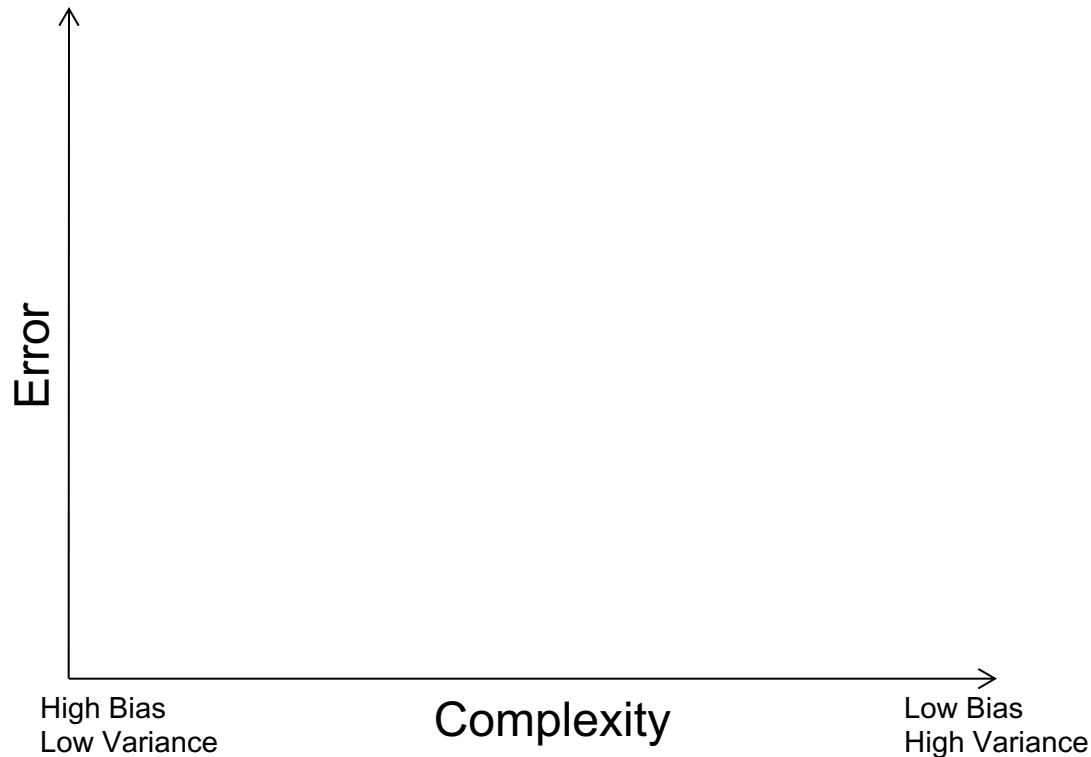
# Bias and Variance

$$\text{Error} = \text{noise}^2 + \text{bias}^2 + \text{variance}$$



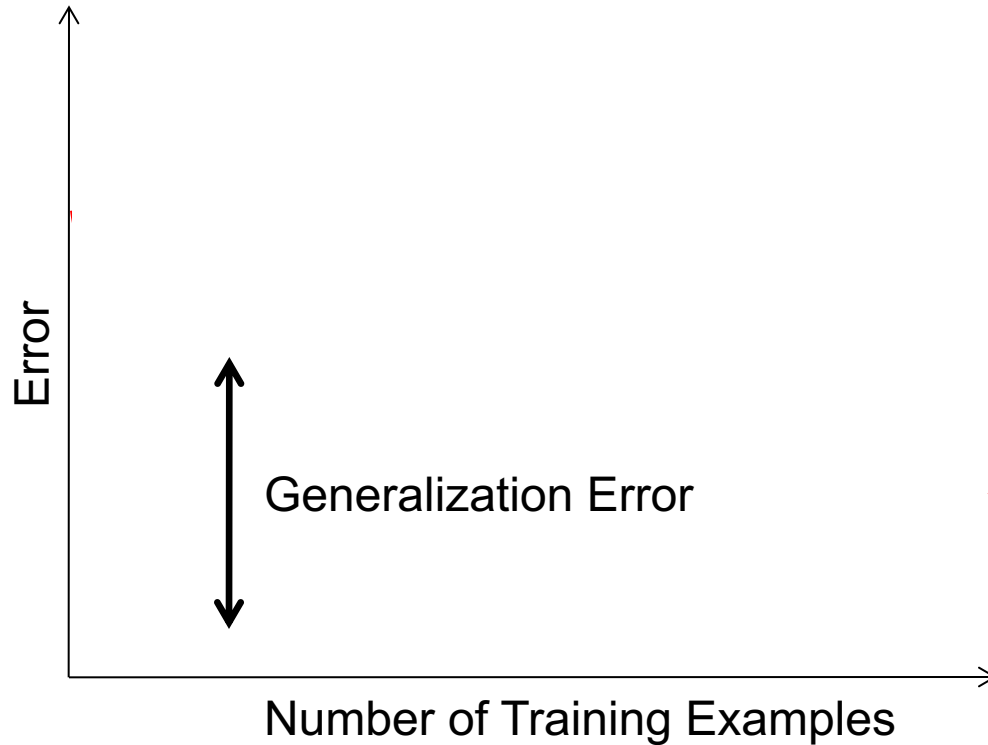
# Choosing the trade-off

- Need validation set
- Validation set is separate from the test set



# Effect of Training Size

Fixed classifier



# Characteristics of vision learning problems

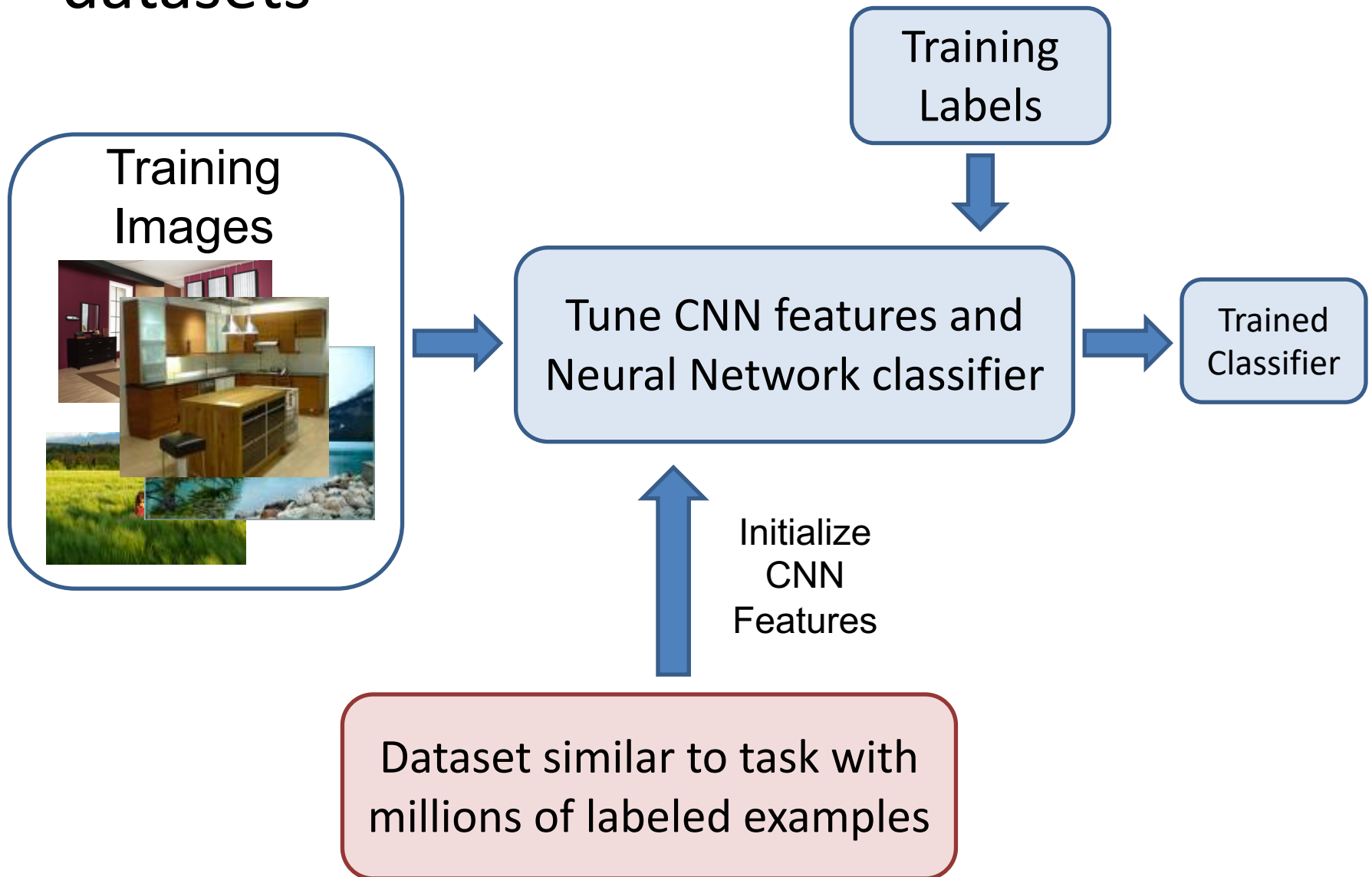
- Lots of continuous features
  - E.g., HOG template may have 1000 features
  - Spatial pyramid may have ~15,000 features
- Imbalanced classes
  - often limited positive examples, practically infinite negative examples
- Difficult prediction tasks

# When a massive training set is available

- Relatively new phenomenon
  - MNIST (handwritten letters) in 1990s, LabelMe in 2000s, ImageNet (object images) in 2009, ...
- Want classifiers with low bias (high variance ok) and reasonably efficient training
- Very complex classifiers with simple features are often effective
  - Random forests
  - Deep convolutional networks



# New training setup with moderate sized datasets



# Practical tips

- Preparing features for linear classifiers
  - Often helps to make zero-mean, unit-dev
  - For non-ordinal features, convert to a set of binary features
- Selecting classifier meta-parameters (e.g., regularization weight)
  - Cross-validation: split data into subsets; train on all but one subset, test on remaining; repeat holding out each subset
    - Leave-one-out, 5-fold, etc.
- Most popular classifiers in vision
  - *SVM*: linear for when fast training/classification is needed; performs well with lots of weak features
  - *Logistic Regression*: outputs a probability; easy to train and apply
  - *Nearest neighbor*: hard to beat if there is tons of data (e.g., character recognition)
  - *Boosted stumps or decision trees*: applies to flexible features, incorporates feature selection, powerful classifiers
  - *Random forests*: outputs probability; good for simple features, tons of data
  - *Deep networks / CNNs*: flexible output; learns features; adapt existing network (which is trained with tons of data) or train new with tons of data
- Always try at least two types of classifiers

# Next class

## Deep convolutional neural networks (CNNs)

- Architecture (Alexnet, VGG, Inception, Resnet)
- Optimizer (stochastic gradient descent, Adam)
- Transfer
- Why it works