# CS 558: Computer Vision
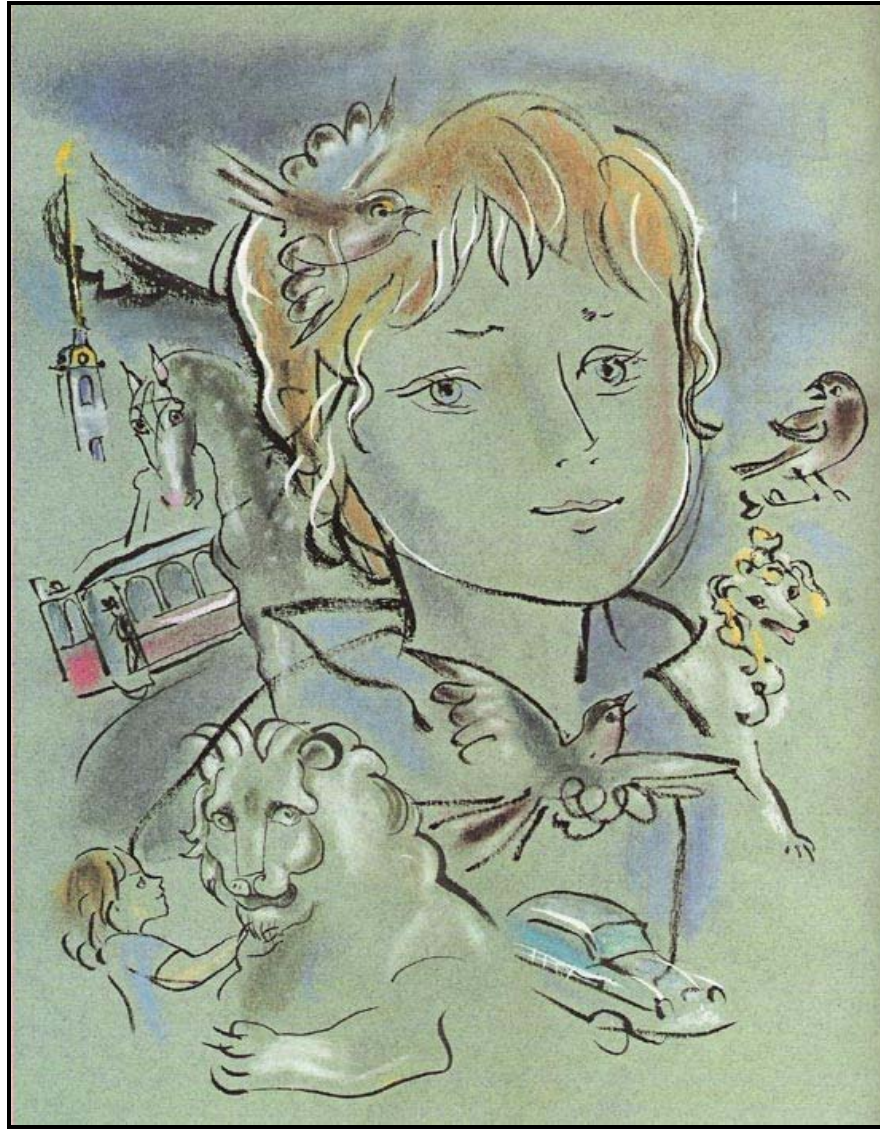# 9th Set of Notes

Instructor: Philippos Mordohai
Webpage: www.cs.stevens.edu/~mordohai
E-mail: Philippos.Mordohai@stevens.edu
Office: Lieb 215

# Introduction to object recognition

By Svetlana Lazebnik



Slides adapted from Fei-Fei Li, Rob Fergus, Antonio Torralba, and Jean Ponce

# Overview

- Basic recognition tasks
- A machine learning approach
  - Example features
  - Example classifiers
  - Levels of supervision
  - Datasets
- Current trends and advanced recognition tasks

# Specific recognition tasks

# Scene categorization



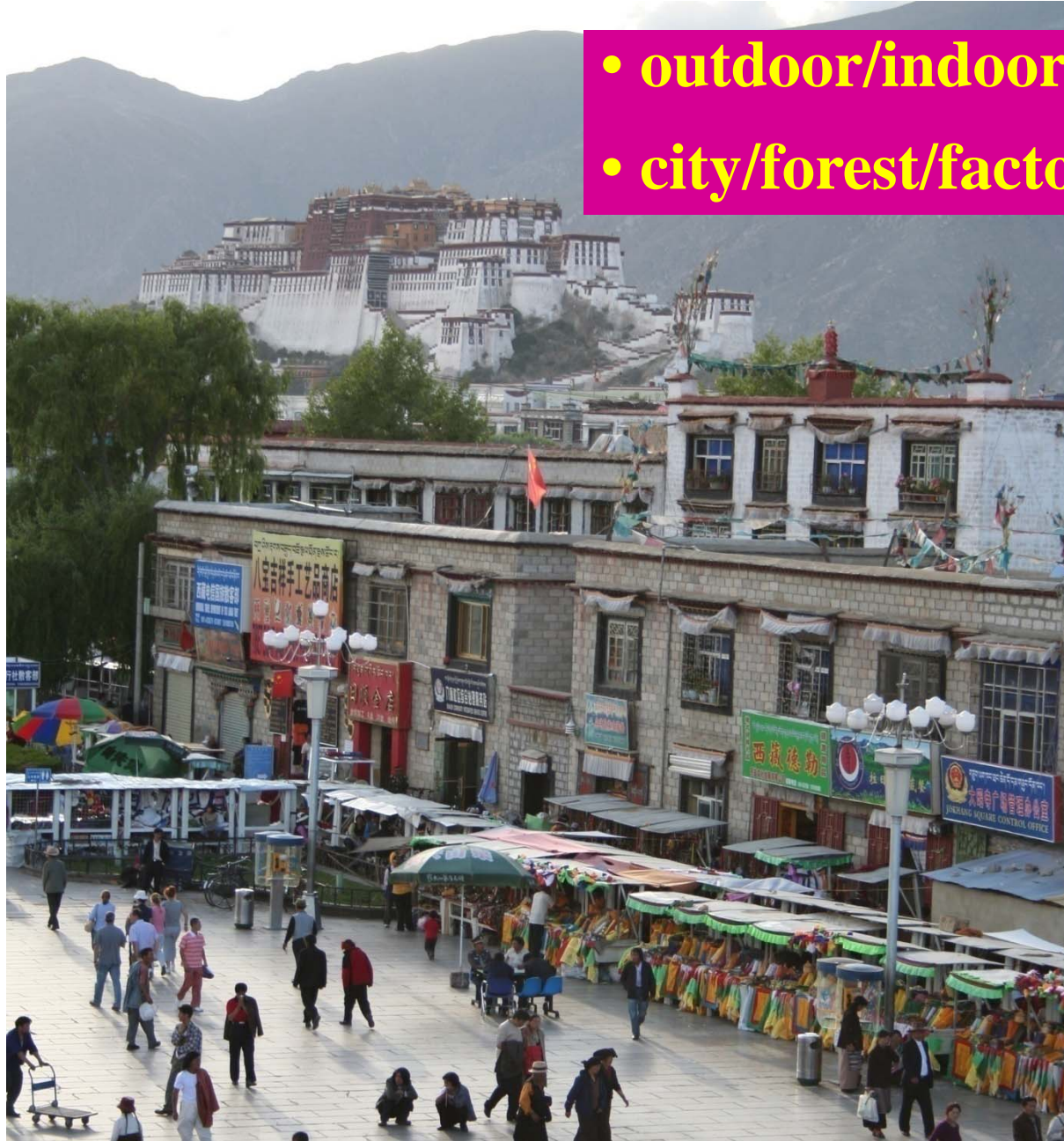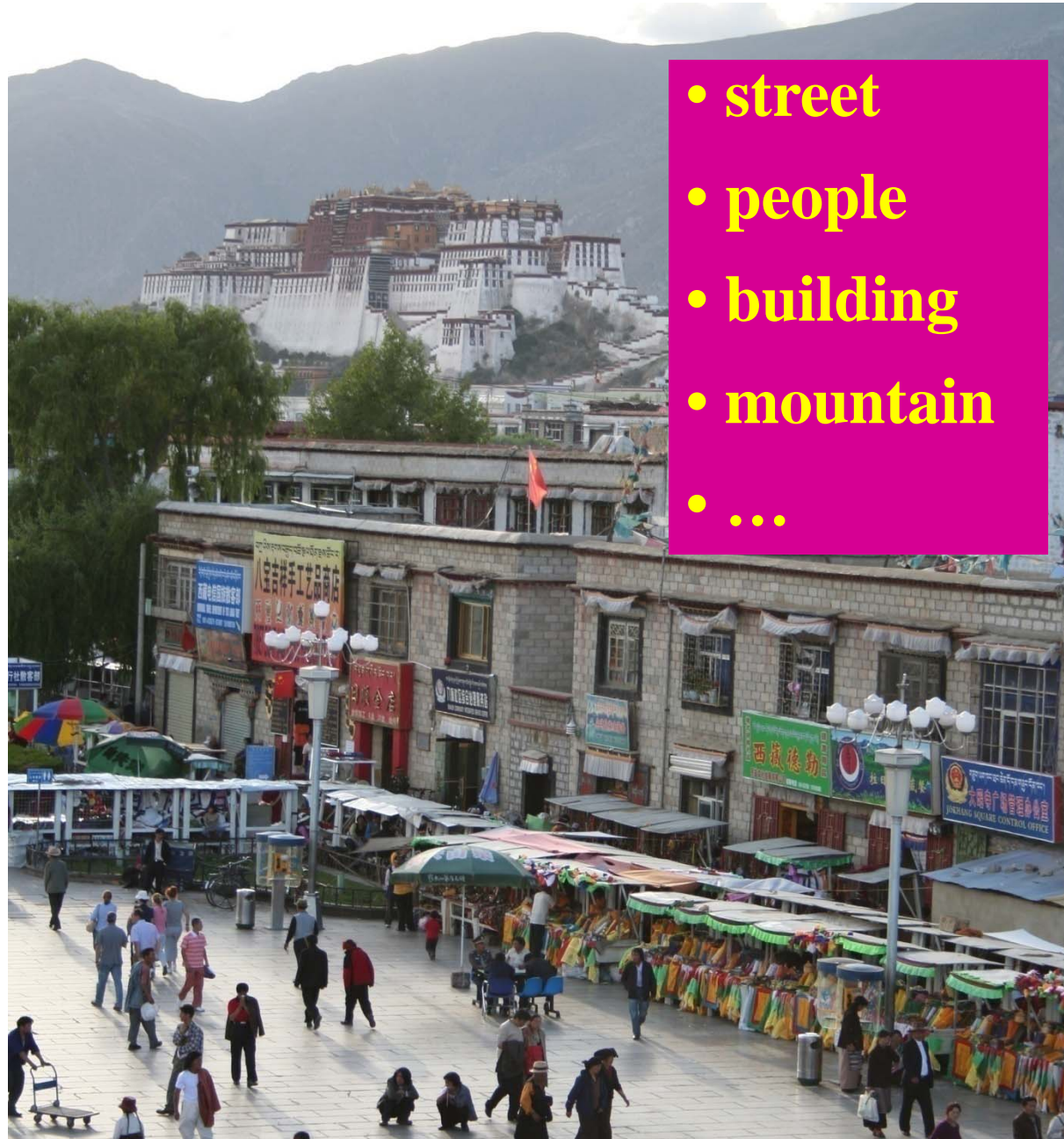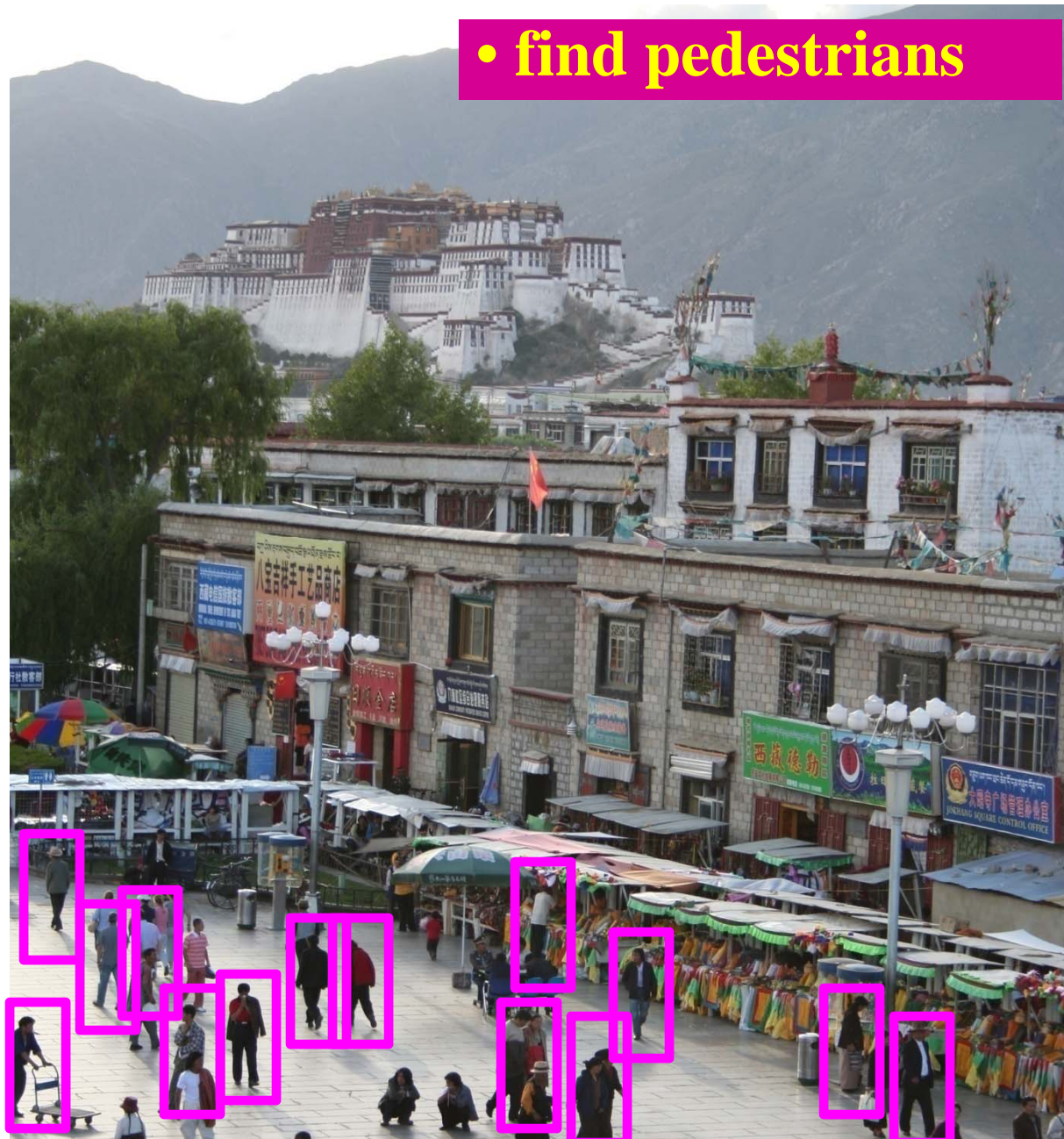- **outdoor/indoor**
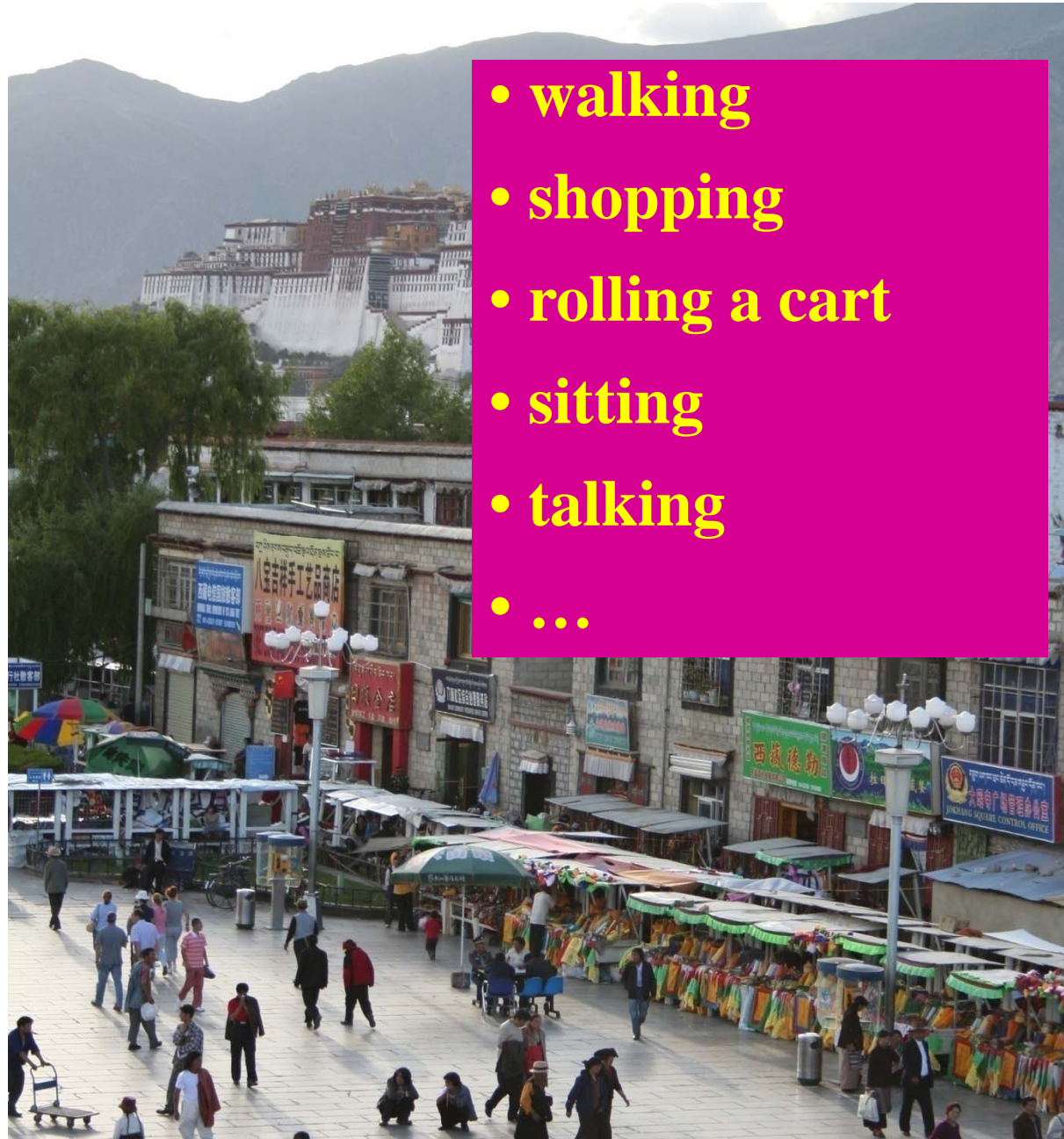- **city/forest/factory/etc.**

# Image annotation/tagging



- **street**
- **people**
- **building**
- **mountain**
- **…**

# Object detection



- **find pedestrians**

# Activity recognition



- **walking**
- **shopping**
- **rolling a cart**
- **sitting**
- **talking**
- **…**

# Image parsing

# Image understanding?

# How many visual object categories are there?



~10,000 to 30,000

Biederman 1987

~10,000 to 30,000

# OBJECTS

## ANIMALS

### …..

### VERTEBRATE

#### MAMMALS

TAPIR

BOAR

#### BIRDS

GROUSE

## PLANTS

## INANIMATE

### NATURAL

### MAN-MADE

CAMERA

# Recognition: A machine learning approach

# The machine learning framework

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

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

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

$$f(\text{}) = \text{"cow"}$$
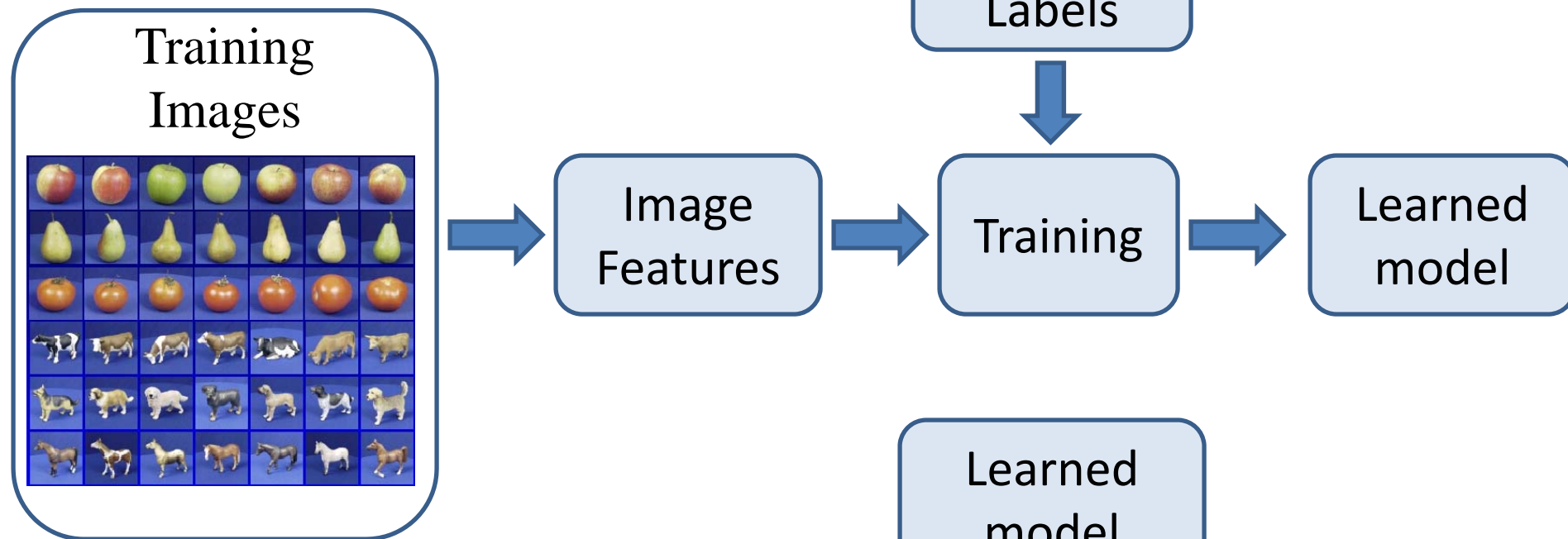
# The machine learning framework

$$y = f(x)$$

output    prediction function    Image feature

- **Training:** given a *training set* of labeled examples $\{(x_1, y_1), \ldots, (x_N, y_N)\}$, estimate the prediction function $f$ by minimizing the prediction error on the training set
- **Testing:** apply $f$ to a never before seen *test example* $x$ and output the predicted value $y = f(x)$

# Steps



**Training**

Training Images

Training Labels

Image Features → Training → Learned model

**Testing**

Test Image → Image Features → Prediction

Learned model

Slide credit: D. Hoiem

# Generalization



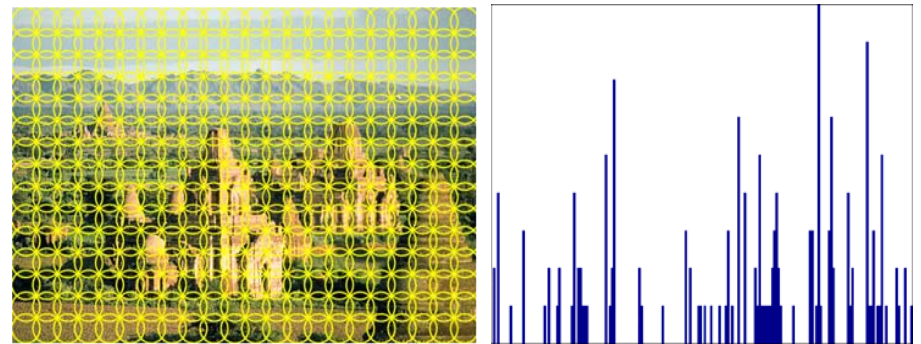Training set (labels known)

Test set (labels unknown)

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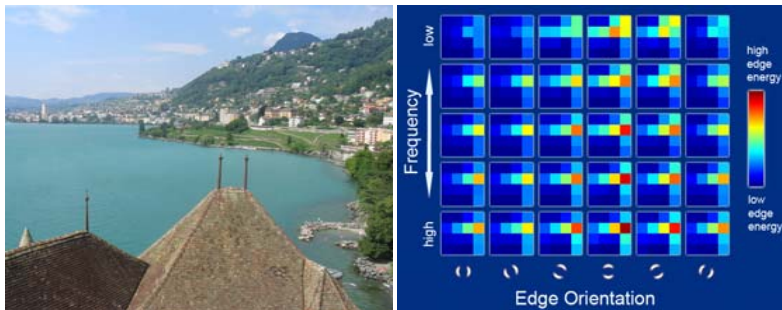
# Popular global image features

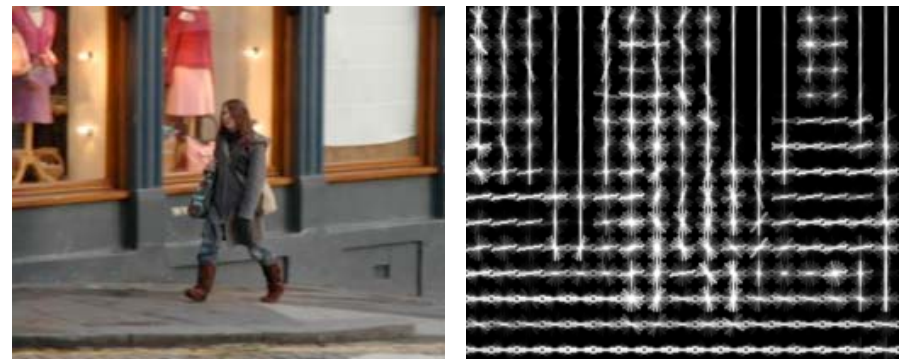- Raw pixels (and simple functions of raw pixels)
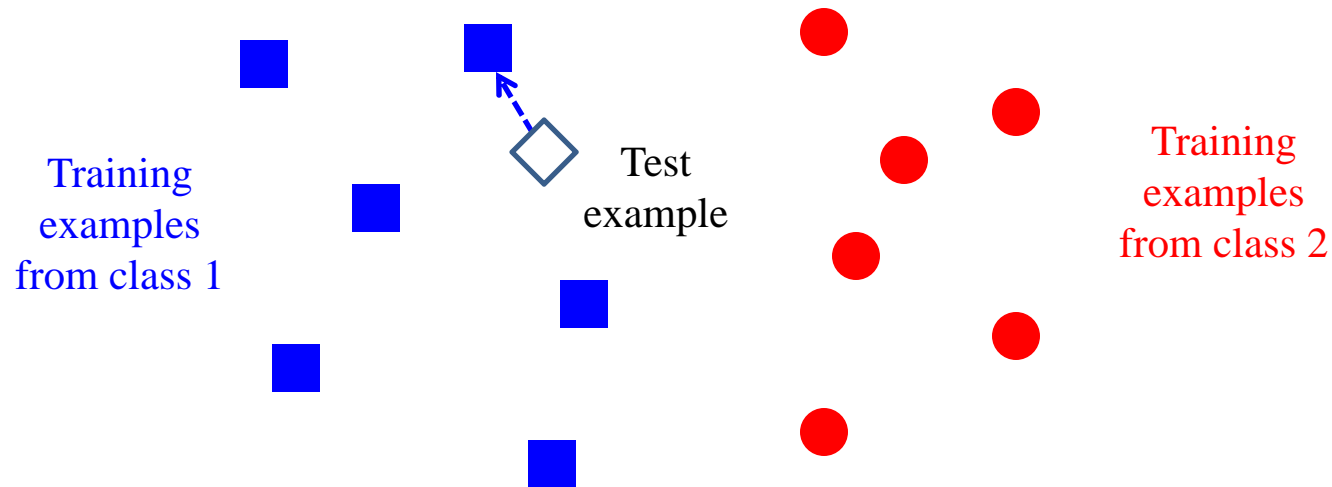


- Histograms, bags of features



- GIST descriptors [Oliva and Torralba, 2001]



- Histograms of oriented gradients (HOG) [Dalal and Triggs, 2005]

# Classifiers: Nearest neighbor



Training examples from class 1

Test example

Training examples from class 2
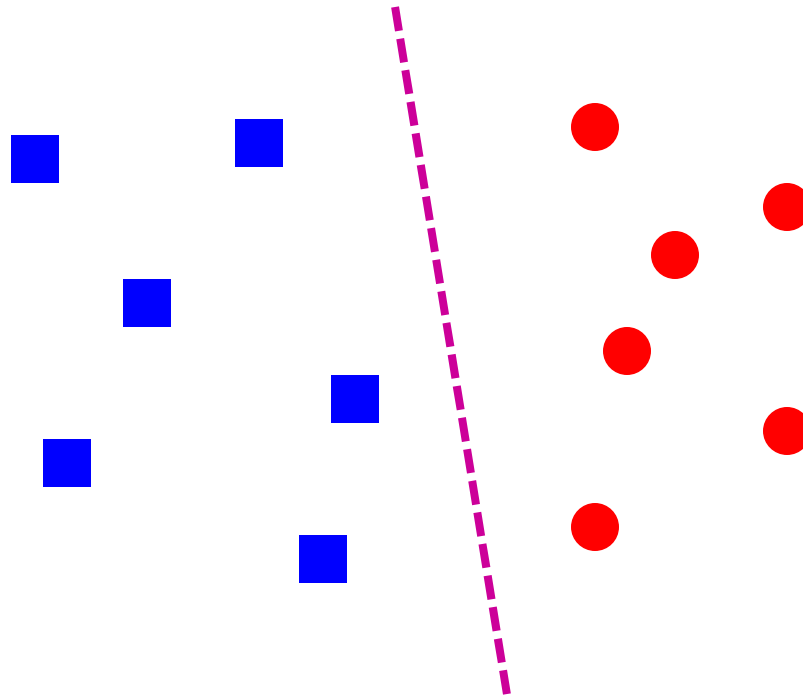
f(**x**) = label of the training example nearest to **x**

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

# Classifiers: Linear
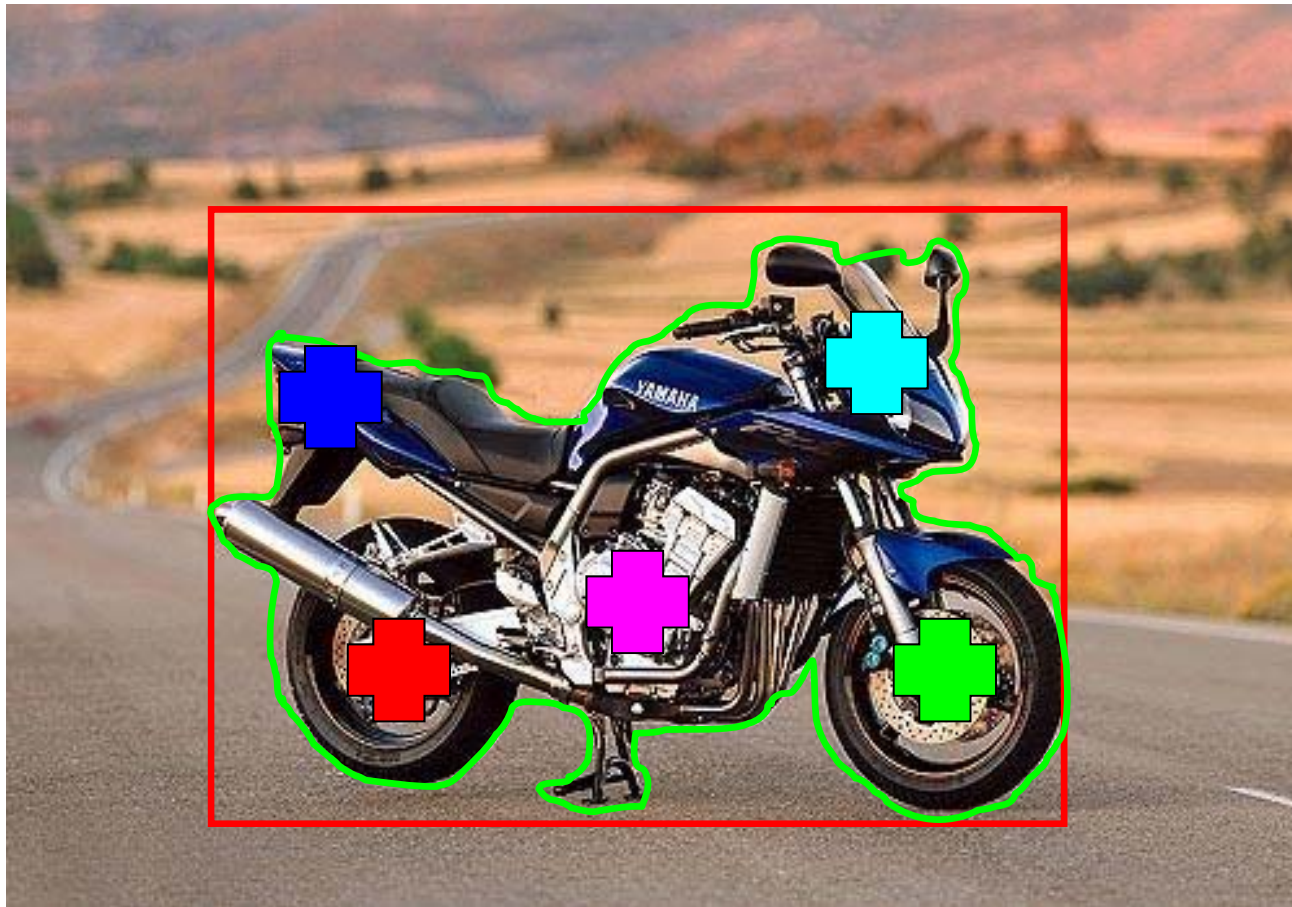


- Find a *linear function* to separate the classes:

$$f(x) = sgn(w \cdot x + b)$$

# Recognition task and supervision

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

Contains a motorbike

# Spectrum of supervision



Less                            More

Unsupervised        "Weakly" supervised        Fully supervised

Definition depends on task

# Datasets

- **Circa 2001:** five categories, hundreds of images per category

- **Circa 2004:** 101 categories

- **Today:** tens of thousands of categories, millions of images

# Caltech 101 & 256

http://www.vision.caltech.edu/Image_Datasets/Caltech101/
http://www.vision.caltech.edu/Image_Datasets/Caltech256/



Fei-Fei, Fergus, Perona, 2004



Griffin, Holub, Perona, 2007

# Caltech-101: Intra-class variability

# ImageNet

IM**A**GENET

14,197,122 images, 21841 synsets indexed

Explore  Download  Challenge<sup>New</sup>  People  Publication  About

Not logged in. Login | Signup

**ImageNet** is an image database organized according to the WordNet hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. Currently we have an average of over five hundred images per node. We hope ImageNet will become a useful resource for researchers, educators, students and all of you who share our passion for pictures.
Click here to learn more about ImageNet, Click here to join the ImageNet mailing list.

SEARCH

What do these images have in common? *Find out!*

The ImageNet Challenge 2013 is announced!

© 2013 Stanford Vision Lab, Stanford University, Princeton University  support@image-net.org  Copyright infringement

# The PASCAL Visual Object Classes Challenge (2005-2012)

http://pascallin.ecs.soton.ac.uk/challenges/VOC/

- **Challenge classes:**

  *Person:* person
  *Animal:* bird, cat, cow, dog, horse, sheep
  *Vehicle:* aeroplane, bicycle, boat, bus, car, motorbike, train
  *Indoor:* bottle, chair, dining table, potted plant, sofa, tv/monitor

- **Dataset size (by 2012):**

  11.5K training/validation images, 27K bounding boxes, 7K segmentations

# PASCAL competitions

- **Classification:** For each of the twenty classes, predicting presence/absence of an example of that class in the test image

- **Detection:** Predicting the bounding box and label of each object from the twenty target classes in the test image
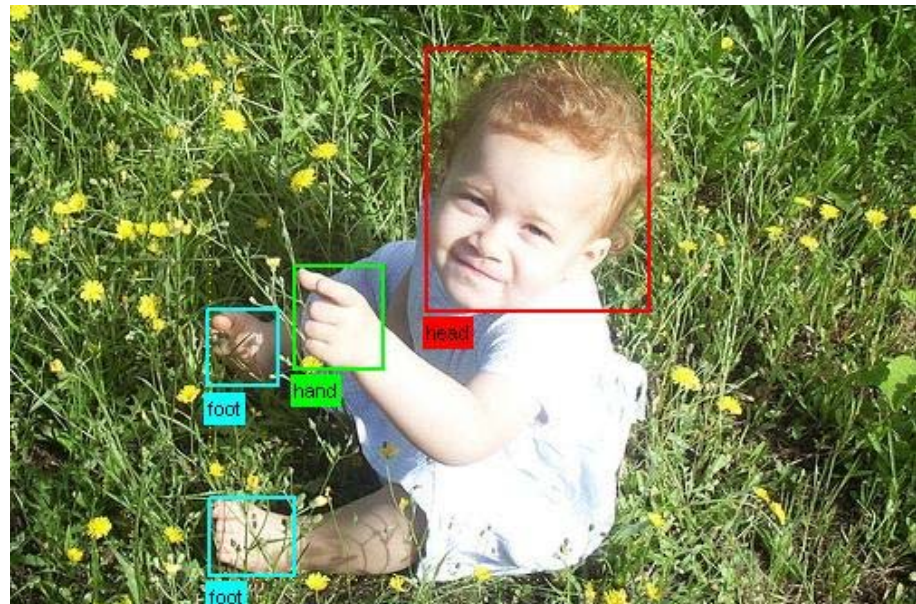
# PASCAL competitions

- **Segmentation:** Generating pixel-wise segmentations giving the class of the object visible at each pixel, or "background" otherwise
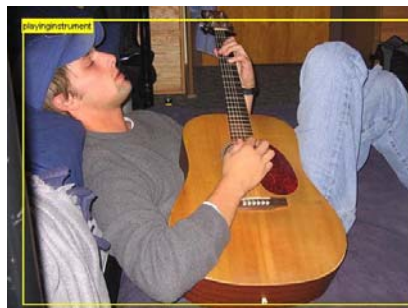


- **Person layout:** Predicting the bounding box and label of each part of a person (head, hands, feet)

# PASCAL competitions

- Action classification (10 action classes)

# LabelMe Dataset

http://labelme.csail.mit.edu/



Russell, Torralba, Murphy, Freeman, 2008

# SUN dataset

~900 scene categories (~400 well-sampled), 130K images



J. Xiao, J. Hays, K. Ehinger, A. Oliva, and A. Torralba, "SUN Database: Large-scale Scene Recognition from Abbey to Zoo," CVPR 2010

http://groups.csail.mit.edu/vision/SUN/

# Fine-grained recognition

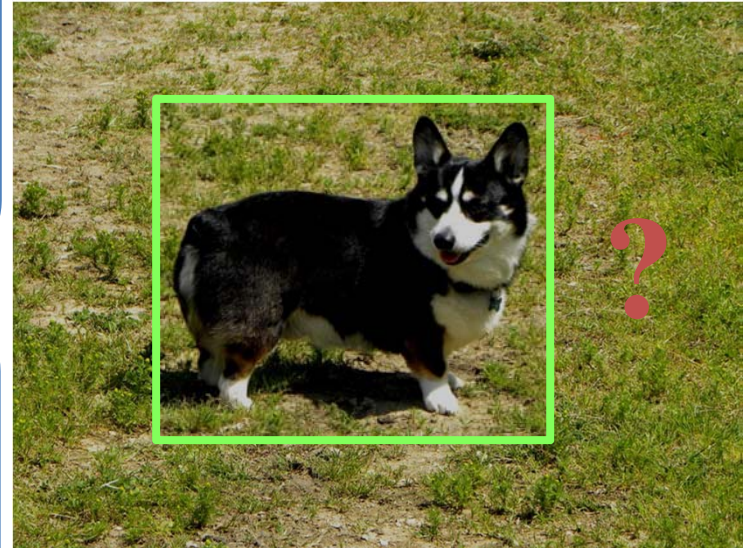# Fine-grained recognition
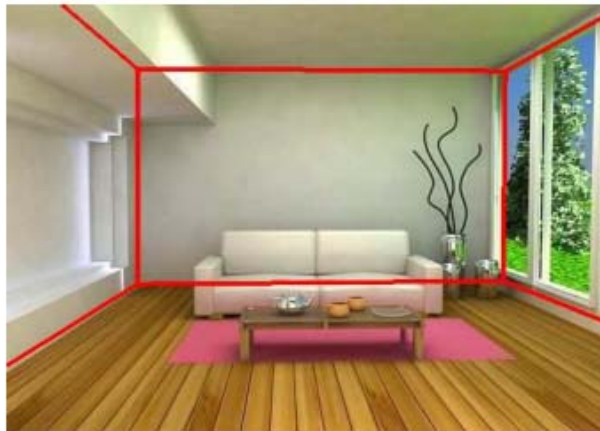


Cardigan Welsh Corgi
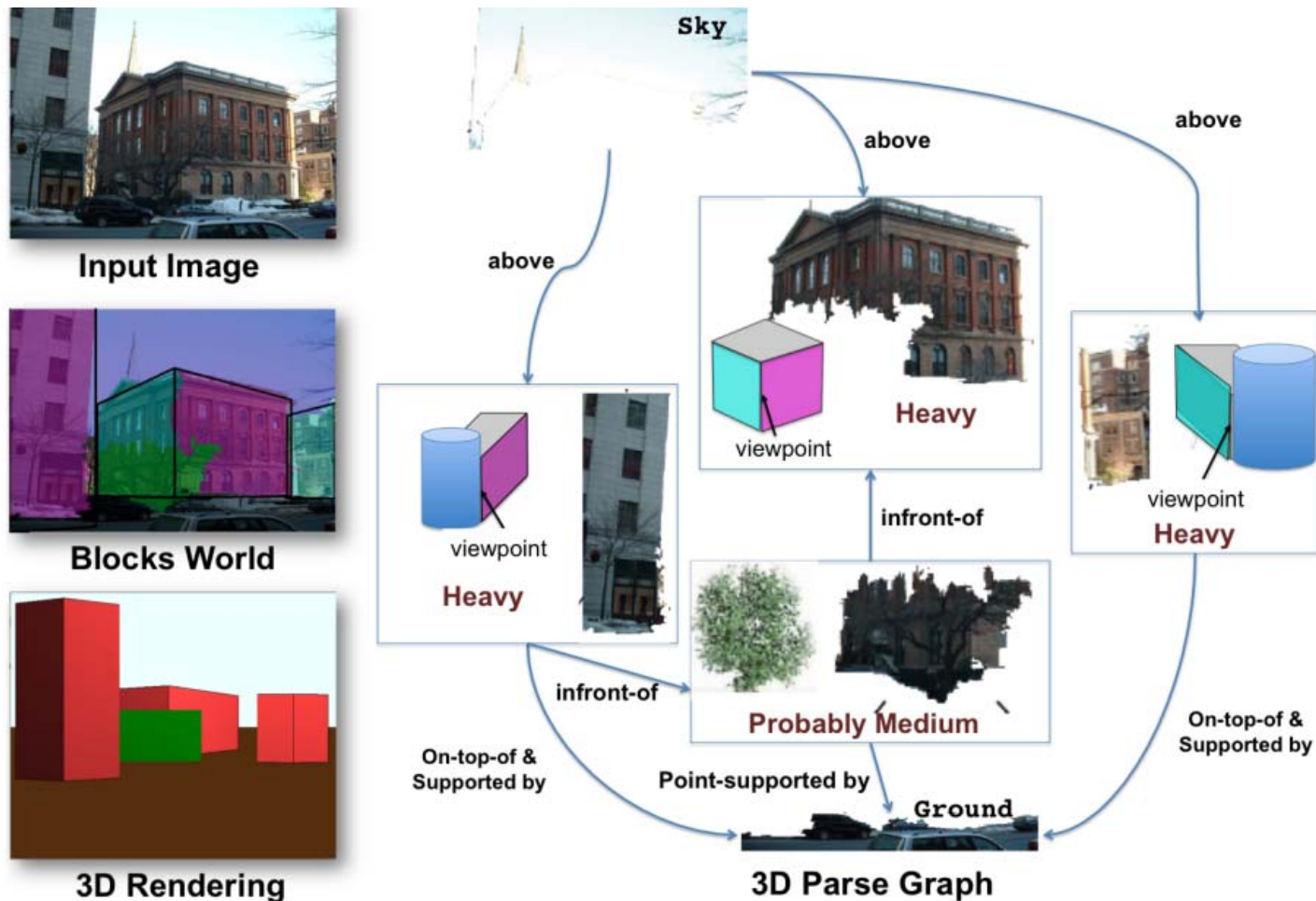
Pembroke Welsh Corgi

**What breed is this dog?**

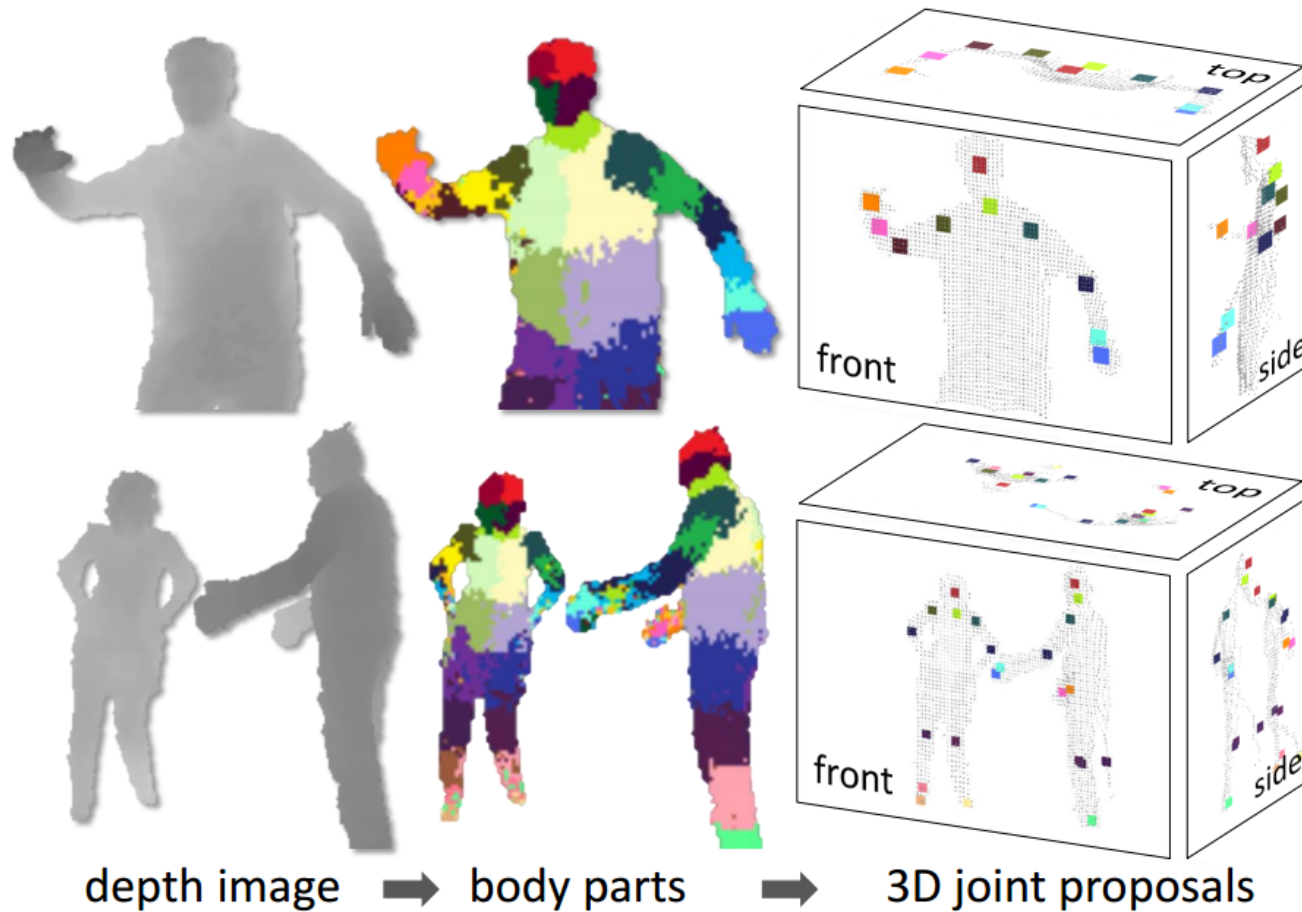**Key: Find the right features.**

# Geometric image interpretation



V. Hedau, D. Hoiem, and D. Forsyth, Recovering the Spatial Layout of Cluttered Rooms, ICCV 2009.

# Geometric image interpretation
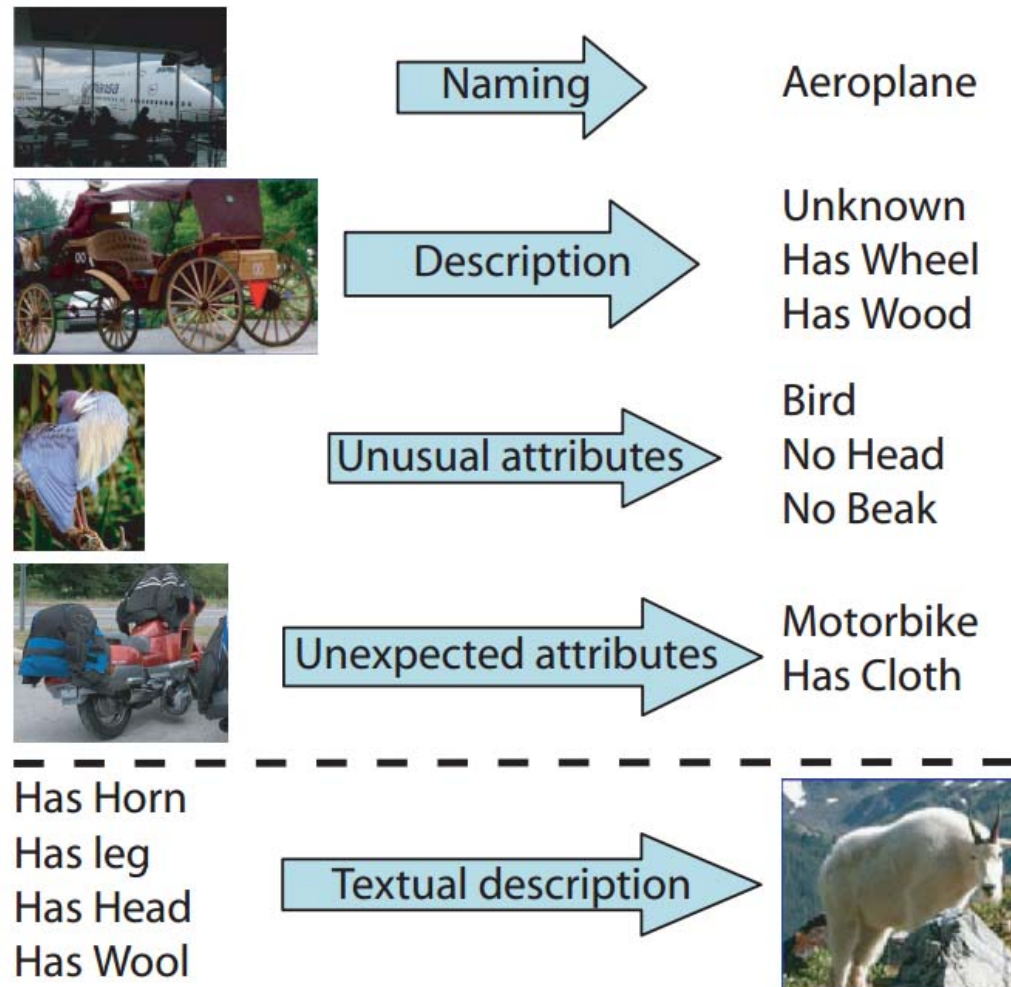


A. Gupta, A. Efros and M. Hebert, **Blocks World Revisited: Image Understanding Using Qualitative Geometry and Mechanics**, ECCV 2010

# Recognition from RGBD Images



depth image ➡ body parts ➡ 3D joint proposals

J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake, **Real-Time Human Pose Recognition in Parts from a Single Depth Image**, CVPR 2011

# Attribute-based recognition



A. Farhadi, I. Endres, D. Hoiem, and D Forsyth, **Describing Objects by their Attributes**, CVPR 2009
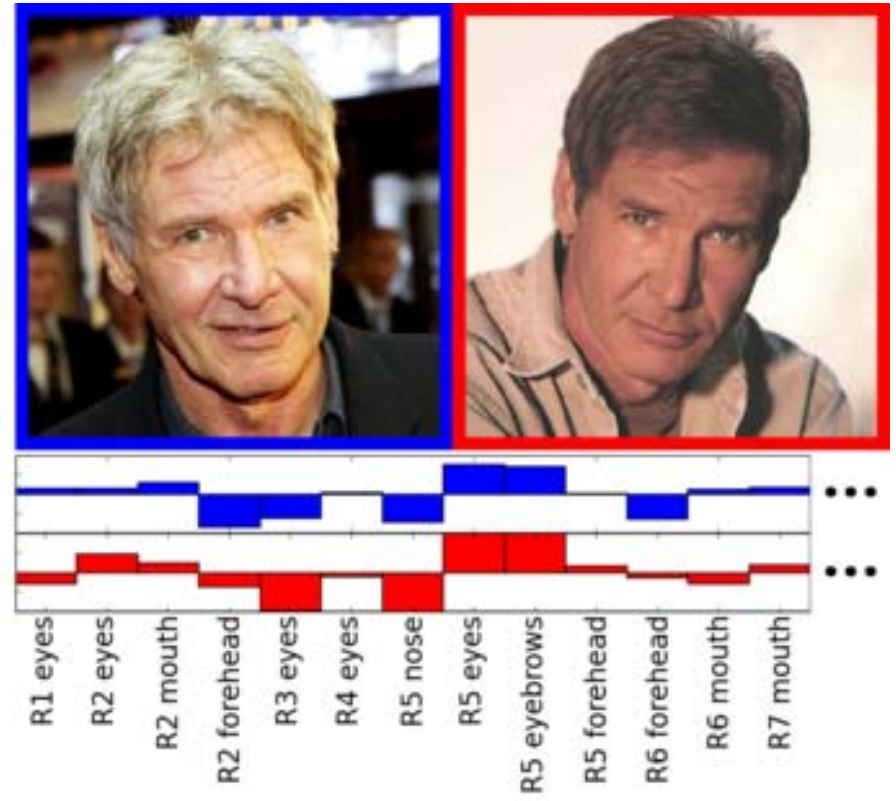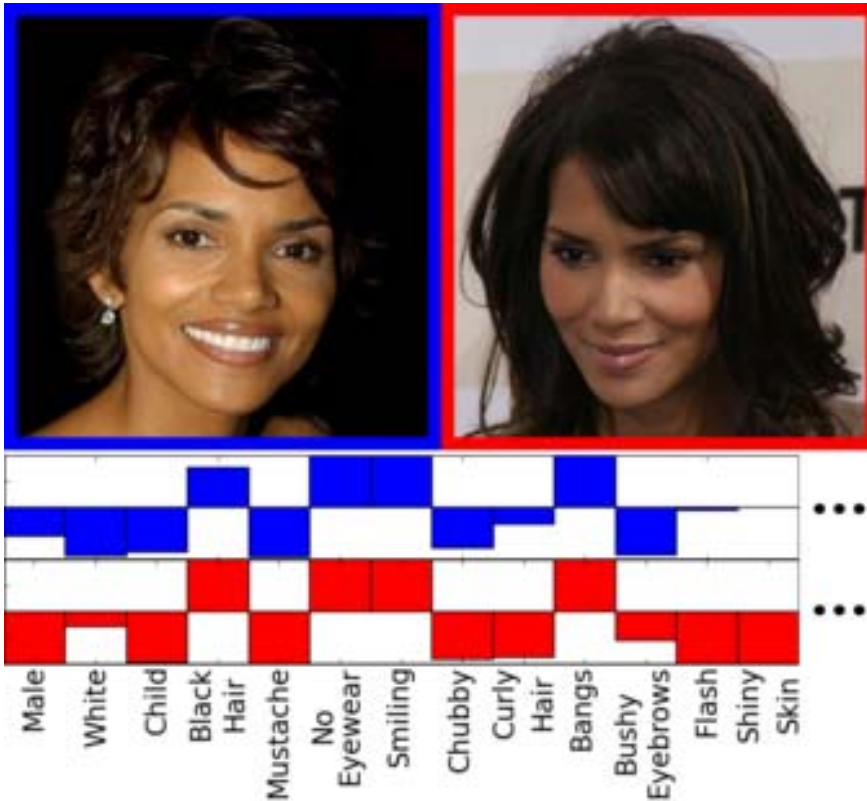
# Attribute-based search



A. Kovashka, D. Parikh and K. Grauman, **WhittleSearch: Image Search with Relative Attribute Feedback**, CVPR 2012

# Face verification



N. Kumar, A. C. Berg, P. N. Belhumeur, and S. K. Nayar, Attribute and Simile Classifiers for Face Verification, ICCV 2009

# Sentence generation from images



This is a photograph of one sky, one road and one bus. The blue sky is above the gray road. The gray road is near the shiny bus. The shiny bus is near the blue sky.

There are two aeroplanes. The first shiny aeroplane is near the second shiny aeroplane.

There are one cow and one sky. The golden cow is by the blue sky.

There are one dining table, one chair and two windows. The wooden dining table is by the wooden chair, and against the first window, and against the second white window. The wooden chair is by the first window, and by the second white window. The first window is by the second white window.

This is a picture of one sky, one road and one sheep. The gray sky is over the gray road. The gray sheep is by the gray road.

Here we see one road, one sky and one bicycle. The road is near the blue sky, and near the colorful bicycle. The colorful bicycle is within the blue sky.

Here we see two persons, one sky and one aeroplane. The first black person is by the blue sky. The blue sky is near the shiny aeroplane. The second black person is by the blue sky. The shiny aeroplane is by the first black person, and by the second black person.

This is a picture of two dogs. The first dog is near the second furry dog.

G. Kulkarni, V. Premraj, S. Dhar, S. Li, Y. Choi, A. Berg, T. Berg, **Baby Talk: Understanding and Generating Simple Image Descriptions**, CVPR 2011