

# Unsupervised Learning of Models for Object Class Recognition

CS 395T  
Object Recognition

Tuyen Huynh  
Prateek Jain

\* Slides taken from R. Fergus, P. Perona and A. Zisserman

## Goal

- Recognition of object categories
- Unassisted learning

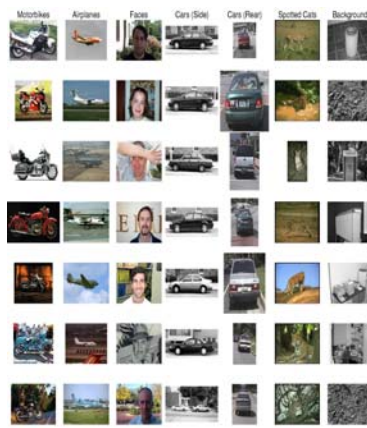


## Some object categories

Learn from examples

Difficulties:

- Size variation
- Background clutter
- Occlusion
- Intra-class variation



## Problems

- Three problems
    - Segmentation of training images
    - Part selection
    - Estimation of model parameters
- Provide framework to solve these problem automatically

## Related Work

- Hierarchical model from edge elements
  - Statistical model from shape space densities
  - Active appearance models
  - Gradient descent on a deformation energy function
- Require some kinds of labeled in the training images

## Approach

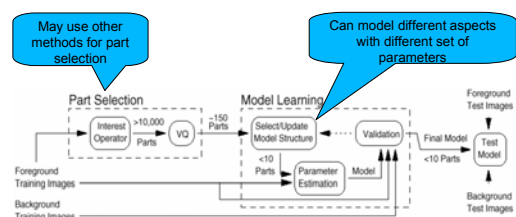
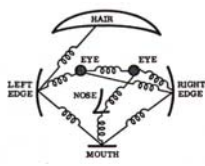
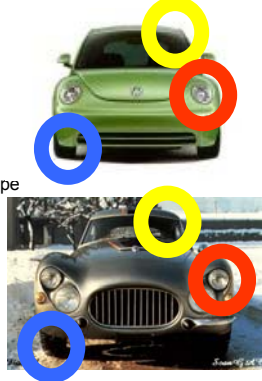


Fig. 2. Block diagram of our method. "Foreground images" are images containing the target objects in cluttered background. "Background images" contain background only.

## Model: Constellation of Parts



Object composition by parts and shape



## Part selection I

- ✧ Interest operator by Forstner
  - ✧ corners
  - ✧ intersection of 2+ lines
  - ✧ center points of circular patterns
- ✧ Vector quantization by k-means clustering
  - ✧ Retain clusters with at least 10 patterns
- ✧ Remove patterns which are similar to others
- ✧ Use greedy search to find the most informative parts

## Part selection I



Fig. 3. Points of interest (left) identified on a training image of a human face in cluttered background using Forstner's method. Crosses denote corner-type patterns while circles mark circle-type patterns. A sample of the patterns obtained using k-means clustering of small image patches is shown for faces (center) and cars (right). The car images were high-pass filtered before the part selection process. The total number of patterns selected were 81 for faces and 80 for cars.

## A generative object model I

- ✧ Assume T different types of parts
- ✧ Observable data

$$X^0 = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1N_1} \\ x_{21} & x_{22} & \dots & x_{2N_2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{T1} & x_{T2} & \dots & x_{TN_T} \end{pmatrix}$$

2-dimensional vector

Row t = All locations of part type t in the image

- ✧ Hidden/Missing data:
  - ✧ Vector  $\mathbf{h} = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T)$  where  $\mathbf{h}_i \in \{0, 1, 2, \dots, N_i\}$
  - ✧ Vector  $\mathbf{x}^m$ : positions of missing parts

## A generative object model I

- ✧ Joint probability density:  $P(X^0, \mathbf{x}^m, \mathbf{h})$
  - ✧ Two auxiliary variables:
    - ✧ Binary vector  $\mathbf{b} = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_T)$  where  $\mathbf{b}_i = 1$  if  $\mathbf{h}_i > 0$
    - ✧ Vector  $\mathbf{n} = (\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_T)$  where  $\mathbf{n}_i$  is the number of background candidates in row i of  $X^0$
- $$\rightarrow P(X^0, \mathbf{x}^m, \mathbf{h}) = P(X^0, \mathbf{x}^m, \mathbf{h}, \mathbf{n}, \mathbf{b})$$
- $$= P(X^0, \mathbf{x}^m | \mathbf{h}, \mathbf{n}) p(\mathbf{h} | \mathbf{n}, \mathbf{b}) p(\mathbf{n}) p(\mathbf{b})$$

## A generative model I

$$p(\mathbf{n}) = \prod_{t=1}^T \frac{1}{n_t!} (M_t)^{n_t} e^{-M_t}$$

$M_t$ : the average number of background of type t per image

$$p(\mathbf{h} | \mathbf{n}, \mathbf{b}) = \begin{cases} \frac{1}{\prod_{j=1}^T N_j^{h_j}} & \mathbf{h} \in \mathcal{H}(\mathbf{b}, \mathbf{n}) \\ 0 & \text{other } \mathbf{h} \end{cases}$$

$\mathcal{H}(\mathbf{b}, \mathbf{n})$ : the set of all hypotheses consistent with  $\mathbf{b}$  and  $\mathbf{n}$

$N_i$ : the total number of detections of the type of part i

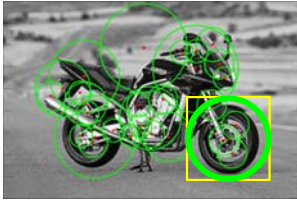
$$p(X^0, \mathbf{x}^m | \mathbf{h}, \mathbf{n}) = p_{bg}(z) p_{fg}(x_{bg})$$

$z^i = (x_{bg}^i, \mathbf{x}^m)$ : the coordinates of all foreground detections. Modeled by Gaussian( $\mu, \Sigma$ )

$$p_{bg}(x_{bg}) = \prod_{t=1}^T \frac{1}{A^{n_t}}$$

$x_{bg}$ : the coordinates of all background detections. Modeled by a uniform density.

## Part Selection II



- Find regions within image
- Use Kadir and Brady's salient region operator [IJCV '01]

### Location

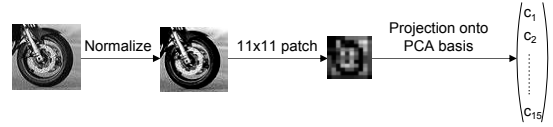
(x,y) coords. of region center

### Scale

Diameter of region (pixels)

- Uses gray-scale input
- Finds maxima in entropy over scale and location

## Appearance



- Normalize all the images to same size
- Representation in low-dimensional vector space, hence speeds up computation
- Ignores noise, hence makes algorithm robust

## Model II

$$p(\mathbf{X}, \mathbf{S}, \mathbf{A} | \theta) = \sum_{\mathbf{h} \in H} p(\mathbf{X}, \mathbf{S}, \mathbf{A}, \mathbf{h} | \theta) = \sum_{\mathbf{h} \in H} \underbrace{p(\mathbf{A} | \mathbf{X}, \mathbf{S}, \mathbf{h}, \theta)}_{\text{Appearance}} \underbrace{p(\mathbf{X} | \mathbf{S}, \mathbf{h}, \theta)}_{\text{Shape}} \underbrace{p(\mathbf{S} | \mathbf{h}, \theta)}_{\text{Rel. Scale}} \underbrace{p(\mathbf{h} | \theta)}_{\text{Other}}$$

- Appearance is highest level abstraction
- Occlusion is lowest level abstraction
- $\mathbf{h}$  – Hypothesis vector, mapping a feature to each part
- If  $h(p)=0$ , part  $p$  is occluded

## Model II

### • Appearance

- Each part has Gaussian Density in Appearance Space
- Independent for each part
- Gaussian density for background
- Feature not in hypothesis calculated under background density
- $d_p=0$  if part  $\hat{\theta}_{pp}^{a1, \dots, aP, d} = \{c_p, V_p\}$   $\hat{\theta}_{ba}^{a1, \dots, aP, d} = \{c_{bg}, V_{bg}\}$
- Parameters:

$$\frac{p(\mathbf{A} | \mathbf{X}, \mathbf{S}, \mathbf{h}, \theta)}{p(\mathbf{A} | \mathbf{X}, \mathbf{S}, \mathbf{h}, \theta_{bg})} = \prod_{p=1}^P \left( \frac{G(\mathbf{A}(h_p) | c_p, V_p)}{G(\mathbf{A}(h_p) | c_{bg}, V_{bg})} \right)^{d_p}$$

## Model II

### • Shape

- Same as Weber et al.'s model
- Joint Gaussian density for foreground features
- Uniform density for background features
- Parameters:  $\theta^{shape} = \{\mu, \Sigma\}$

$$\frac{p(\mathbf{X} | \mathbf{S}, \mathbf{h}, \theta)}{p(\mathbf{X} | \mathbf{S}, \mathbf{h}, \theta_{bg})} = G(\mathbf{X}(\mathbf{h}) | \mu, \Sigma) \alpha^f$$

## Model II

### • Relative Scale

- Scale with reference to a fixed frame
- Modeled by Gaussian density
- Parts assumed independent of each other
- Parameters  $\theta^{scale} = \{t_p, \tilde{U}_p\}$

$$\frac{p(\mathbf{S} | \mathbf{h}, \theta)}{p(\mathbf{S} | \mathbf{h}, \theta_{bg})} = \prod_{p=1}^P G(\mathbf{S}(h_p) | t_p, \tilde{U}_p)^{d_p} r^f$$

## Model II

### ✧ Occlusion

- ✧ Same as Weber et al.'s model
- ✧ Number of features detected modeled using Poisson distribution
- ✧ Probability table for all possible occlusion patterns is a parameter

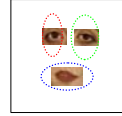
$$\frac{p(\mathbf{h}|\theta)}{p(\mathbf{h}|\theta_{bg})} = \frac{pPois(n|M)}{pPois(N|M)} \frac{1}{nCr(N, f)} p(d|\theta)$$

## Generative probabilistic model

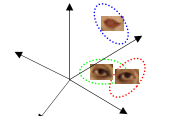
### Foreground model

based on Burl, Weber et al. [ECCV '98, '00]

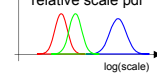
#### Gaussian shape pdf



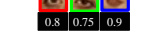
#### Gaussian part appearance pdf



#### Gaussian relative scale pdf



#### Prob. of detection

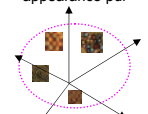


### Clutter model

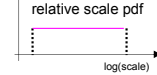
#### Uniform shape pdf



#### Gaussian background appearance pdf



#### Uniform relative scale pdf



#### Poisson pdf on # detections

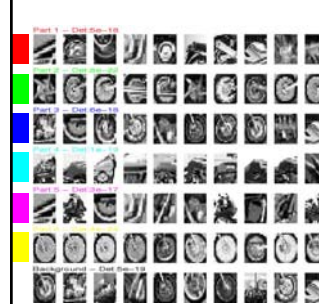
## Classification

$$\begin{aligned} R &= \frac{p(\text{Object}|\mathbf{X}, \mathbf{S}, \mathbf{A})}{p(\text{No object}|\mathbf{X}, \mathbf{S}, \mathbf{A})} \\ &= \frac{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\text{Object}) p(\text{Object})}{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\text{No object}) p(\text{No object})} \\ &\approx \frac{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\theta) p(\text{Object})}{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\theta_{bg}) p(\text{No object})} \end{aligned}$$

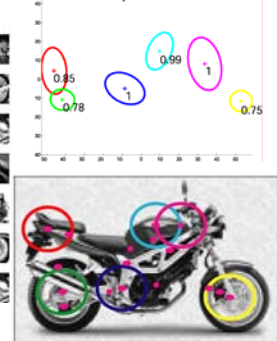
- ✧ Priors estimated using training set
- ✧ Object iff  $R > \text{threshold}$

## Motorbikes

### Samples from appearance model



### Shape model



## Learning

- Task: Estimation of model parameters
- Chicken and Egg type problem, since we initially know neither:
  - Model parameters
  - Assignment of regions to foreground / background
- Let the assignments be a hidden variable and use EM algorithm to learn them and the model parameters



## Learning procedure

- Find regions & their location, scale & appearance
- Initialize model parameters
- Use EM and iterate to convergence:
  - E-step: Compute assignments for which regions are foreground / background
  - M-step: Update model parameters
- Trying to maximize likelihood – consistency in shape & appearance



# Experiments

## Experimental procedure

Two series of experiments:

- Fixed-scale model - Objects the same size (manual normalization)
- Scale-invariant model - Objects between 100 and 550 pixels in width

### Datasets

#### Training

- 50% images
- No identification of object within image

Motorbikes



Airplanes



Frontal Faces



#### Testing

- 50% images
- Simple object present/absent test

Cars (Side)



Cars (Rear)

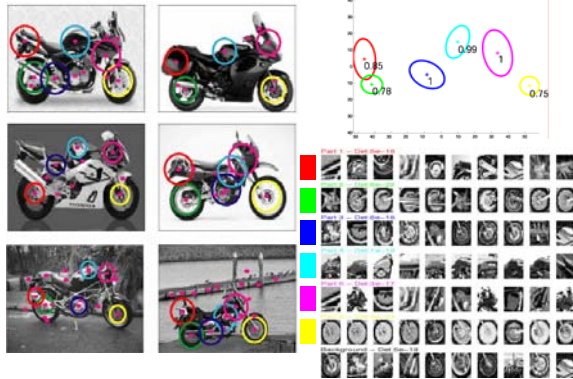


Spotted cats

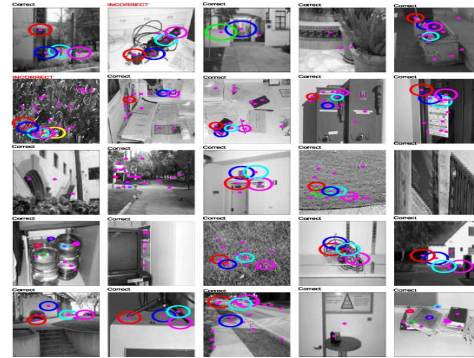


## Motorbikes

Shape model

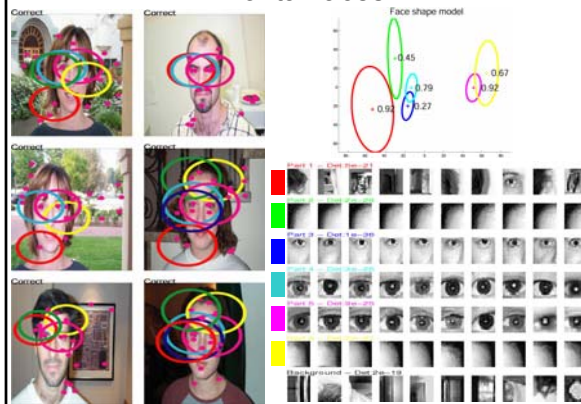


## Background images evaluated with motorbike model



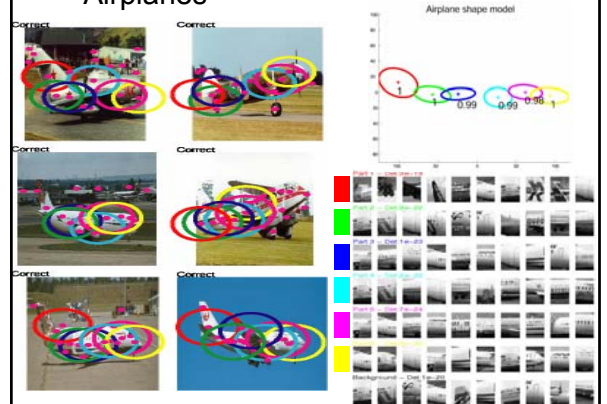
## Frontal faces

Face shape model

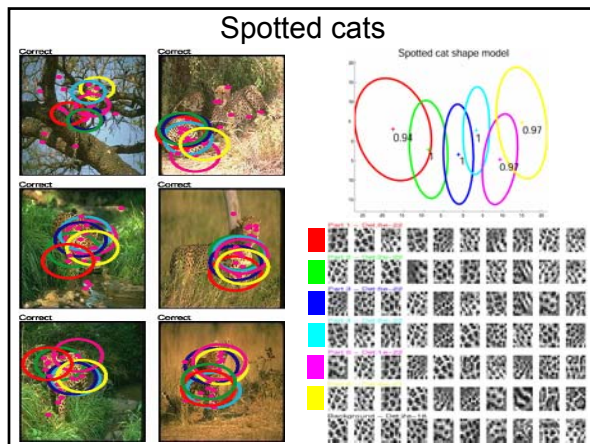


## Airplanes

Airplane shape model





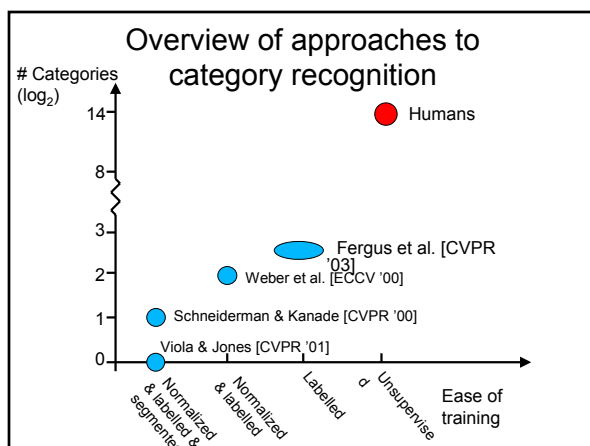
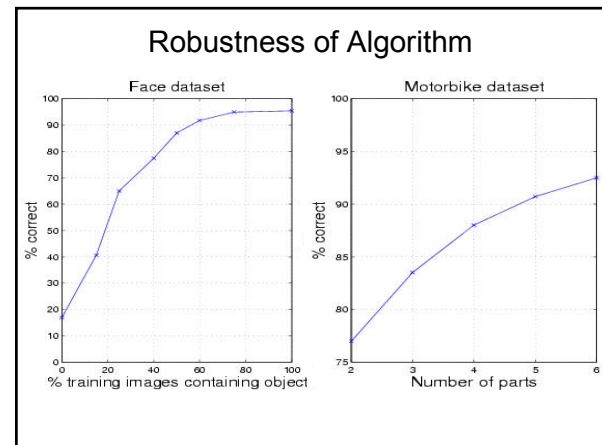
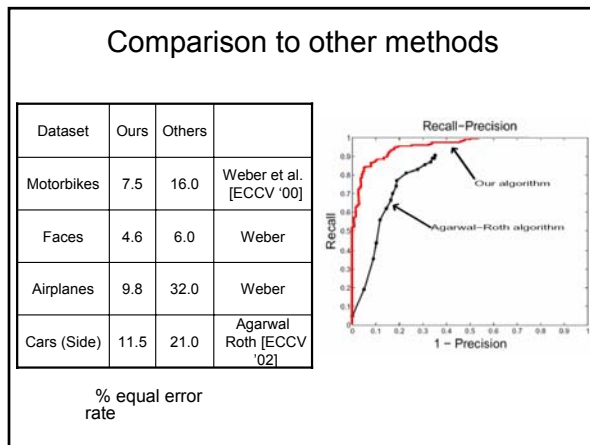


### Summary of results

Dataset	Fixed scale experiment	Scale invariant experiment
Motorbikes	7.5	6.7
Faces	4.6	4.6
Airplanes	9.8	7.0
Cars (Rear)	15.2	9.7
Spotted cats	10.0	10.0

% equal error rate

Note: Within each series, same settings used for all datasets



- ### Summary
- Comprehensive probabilistic model for object classes
  - Learn appearance, shape, relative scale, occlusion etc. simultaneously in scale and translation invariant manner
  - Same algorithm gives  $\leq 10\%$  error across 5 diverse datasets with identical settings
- ### Limitations $\rightarrow$ future work
- Very reliant on region detector  
Different part types (e.g. edge curves)
  - Only learns a single viewpoint  
Use mixture models
  - Need lots of images to learn  
Bayesian learning - fewer images [ICCV '03 (Fei Fei, Fergus, Perona)]
  - Need more thorough testing