The seal of Stanford University is visible in the background on the left side of the slide. It features a red circular border with the text 'STANFORD UNIVERSITY' at the top and '1891' at the bottom. Inside the circle is a red tree (El Palo Alto) on a hill, with the words 'FREIHEIT WEHT' (Liberty rings) above it.

Lecture 16: Object recognition: Part-based generative models

Professor Fei-Fei Li
Stanford Vision Lab

What we will learn today?

- Introduction
- Constellation model
 - Weakly supervised training
 - One-shot learning
- (Problem Set 4 (Q1))

Challenges: intra-class variation



Usual Challenges:

Variability due to:

- View point
- Illumination
- Occlusions

Basic issues

- **Representation**

- 2D Bag of Words (BoW) models;
- Part-based models;
- Multi-view models;

- **Learning**

- Generative & Discriminative BoW models
- Generative models
- Probabilistic Hough voting

- **Recognition**

- Classification with BoW
- Classification with Part-based models

Basic issues

- **Representation**

- 2D Bag of Words (BoW) models;
- Part-based models;
- Multi-view models;

- **Learning**

- Generative & Discriminative BoW models
- Generative models
- Probabilistic Hough voting

- **Recognition**

- Classification with BoW
- Classification with Part-based models

Basic issues

- **Representation**

- 2D Bag of Words (BoW) models;
- Part-based models;
- Multi-view models (Lecture #19);

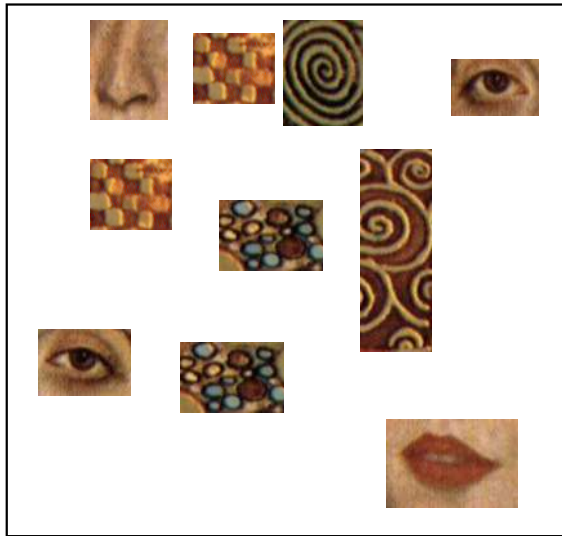
- **Learning**

- Generative & Discriminative BoW models
- Generative models
- Probabilistic Hough voting

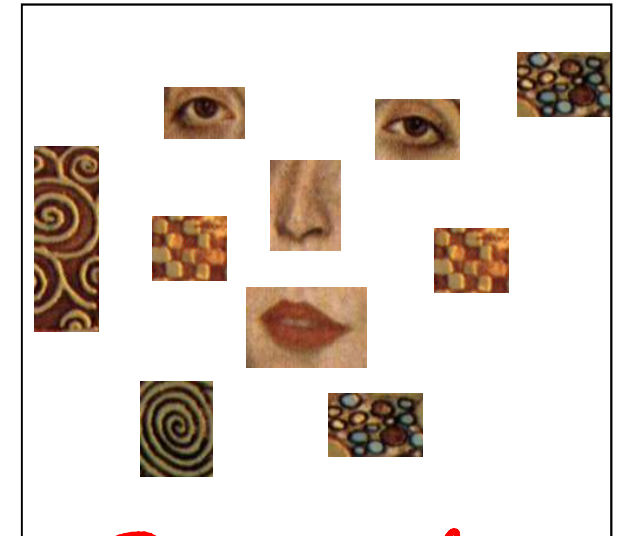
- **Recognition**

- Classification with BoW
- Classification with Part-based models

Problem with bag-of-words



BoW



Part-based

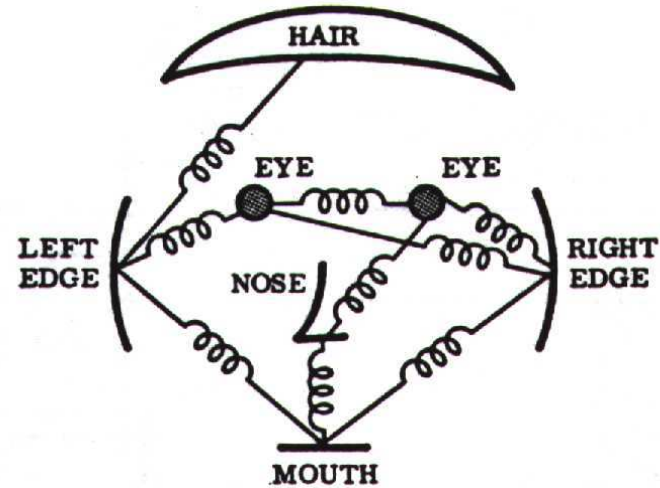
- All have equal probability for bag-of-words methods
- Location information is important

Model: Parts and Structure

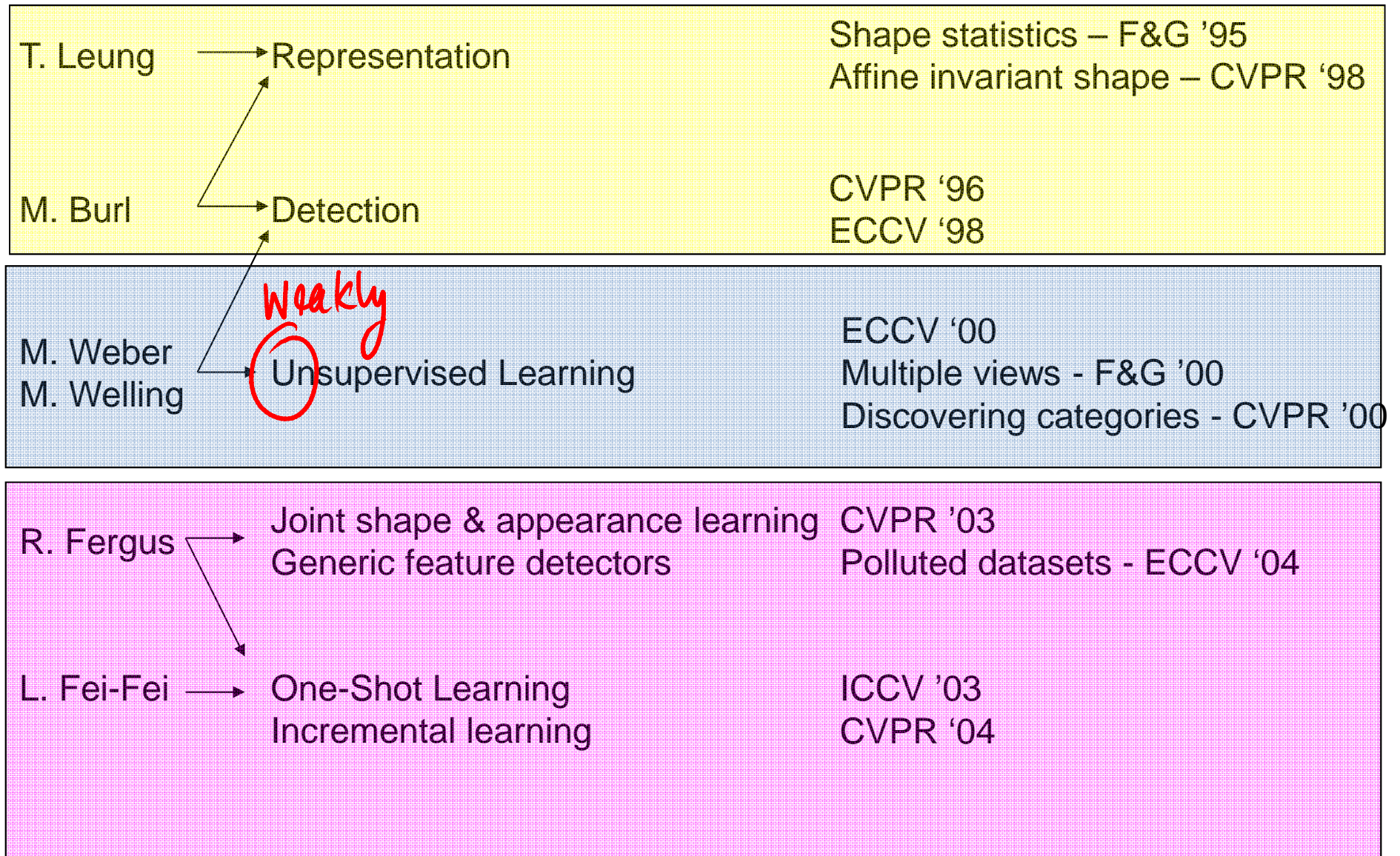


Parts and Structure Literature

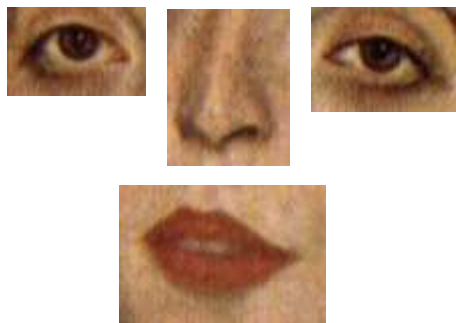
- Fischler & Elschlager 1973
- Yuille '91
- Brunelli & Poggio '93
- Lades, v.d. Malsburg et al. '93
- Cootes, Lanitis, Taylor et al. '95
- Amit & Geman '95, '99
- **et al. Perona '95, '96, '98, '00, '03**
- Huttenlocher et al. '00
- Agarwal & Roth '02
- etc...



The Constellation Model



Deformations



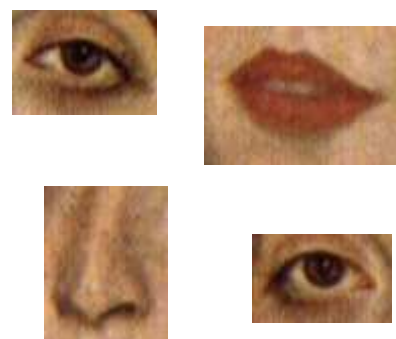
A



B



C



D

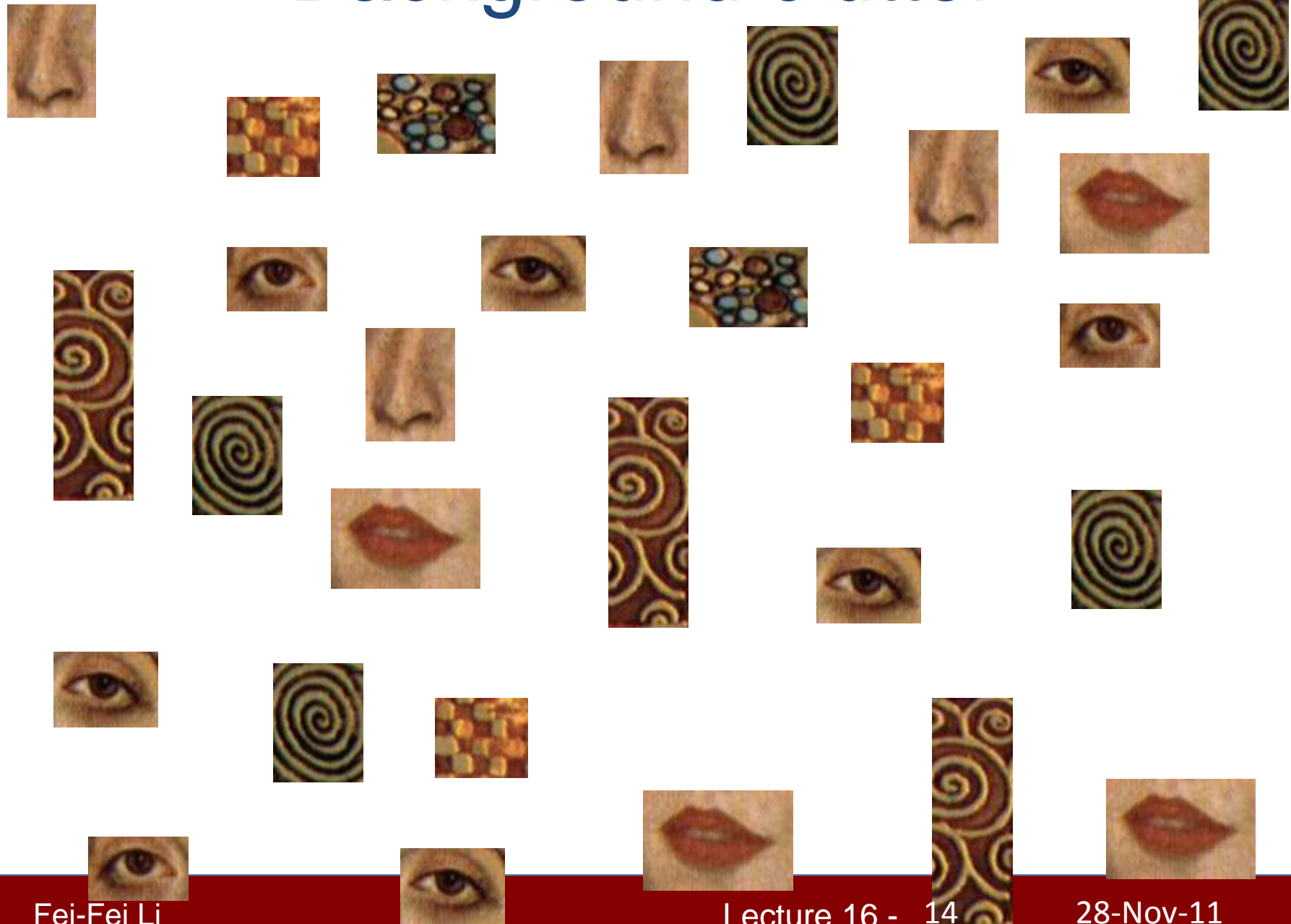
Presence / Absence of Features



occlusion



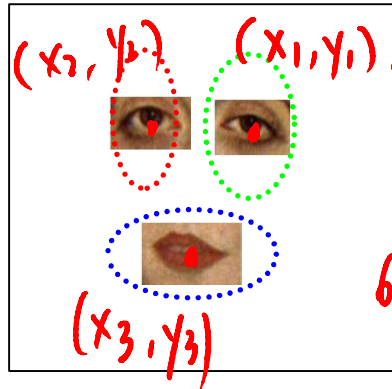
Background clutter



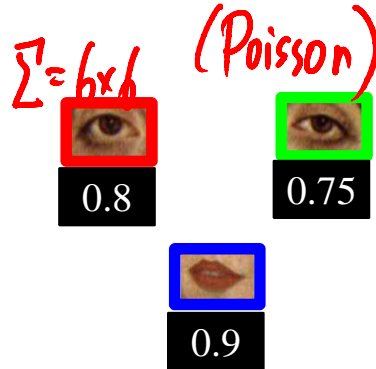
Generative probabilistic model

Foreground model (3-part) ← Clutter model

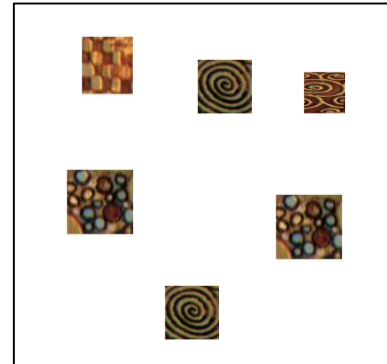
Gaussian shape pdf



Prob. of detection



Uniform shape pdf



detections

$$p_{\text{Poisson}}(N_1/\lambda_1)$$

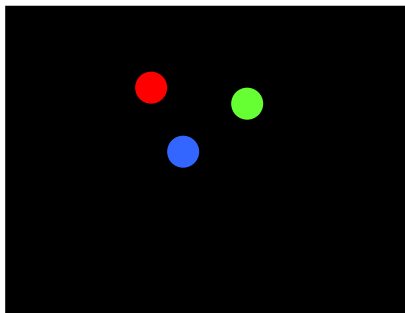
$$p_{\text{Poisson}}(N_2/\lambda_2)$$

$$p_{\text{Poisson}}(N_3/\lambda_3)$$

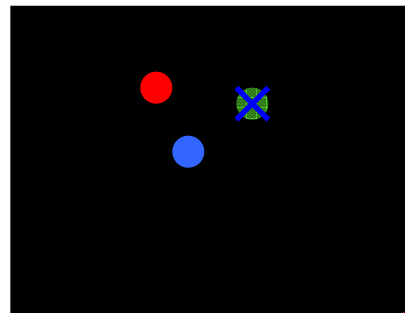
Assumptions: (a) Clutter independent of foreground detections
(b) Clutter detections independent of each other

Example

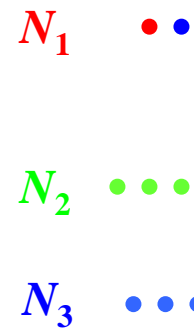
1. Object Part Positions



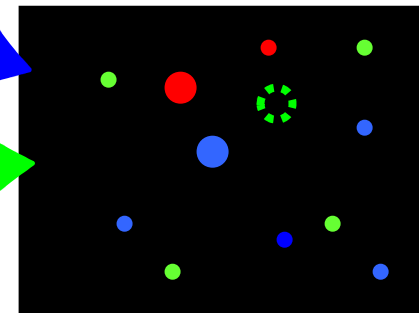
2. Part Absence



3a. N false detect



3b. Position f. detect

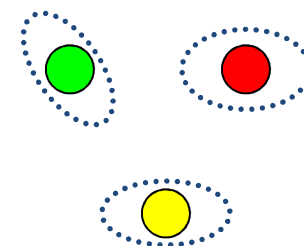


Learning Models 'Manually'



Goal: μ, Σ

- Obtain set of training images
- Choose parts
- Label parts by hand, train detectors
- Learn model from labeled parts

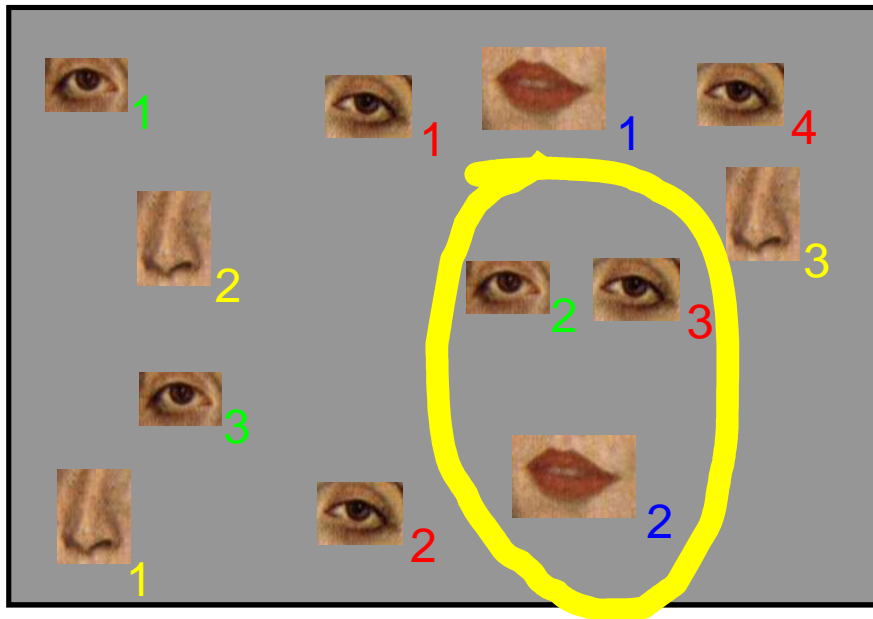


$$\prod_{i=1}^{100} \frac{1}{|A|} \quad A: \text{area of the image}$$

Recognition

$$p(fg|O, h) = \text{Gaussian}\left(\begin{bmatrix} x_1 \\ y_1 \\ x_2 \\ y_2 \\ x_3 \\ y_3 \end{bmatrix} \middle| \begin{matrix} \mu^* \\ \Sigma^* \end{matrix}\right)$$

1. Run part detectors exhaustively over image



hypothesis list

$$h = \begin{pmatrix} 0 \dots N_1 \\ 0 \dots N_2 \\ 0 \dots N_3 \\ 0 \dots N_4 \end{pmatrix}$$

$$\text{e.g. } h = \begin{pmatrix} 2 \\ 3 \\ 0 \\ 2 \end{pmatrix}$$

2. Try different combinations of detections in model
 - Allow detections to be missing (occlusion)

3. Pick hypothesis which maximizes:

$$\text{Threshold} < \frac{p(\text{Data} | \text{Object Hyp})}{p(\text{Data} | \text{Clutter, Hyp})}$$

$$= \frac{p(fg|O) \cdot p(bg|O)}{p(bg|O)}$$

learned μ, Σ
 3 patches 97 patches uniform
 100 patches

4. If ratio is above threshold then, instance detected

So far.....

- Representation
 - Joint model of part locations
 - Ability to deal with background clutter and occlusions
- Learning
 - Manual construction of part detectors
 - Estimate parameters of shape density
- Recognition
 - Run part detectors over image
 - Try combinations of features in model
 - Use efficient search techniques to make fast

100 det. (100)
3-part (3)

Image: label is given (i.e. a "face" image)

Object: no label is given (i.e. don't know where each part is)

Unsupervised Learning

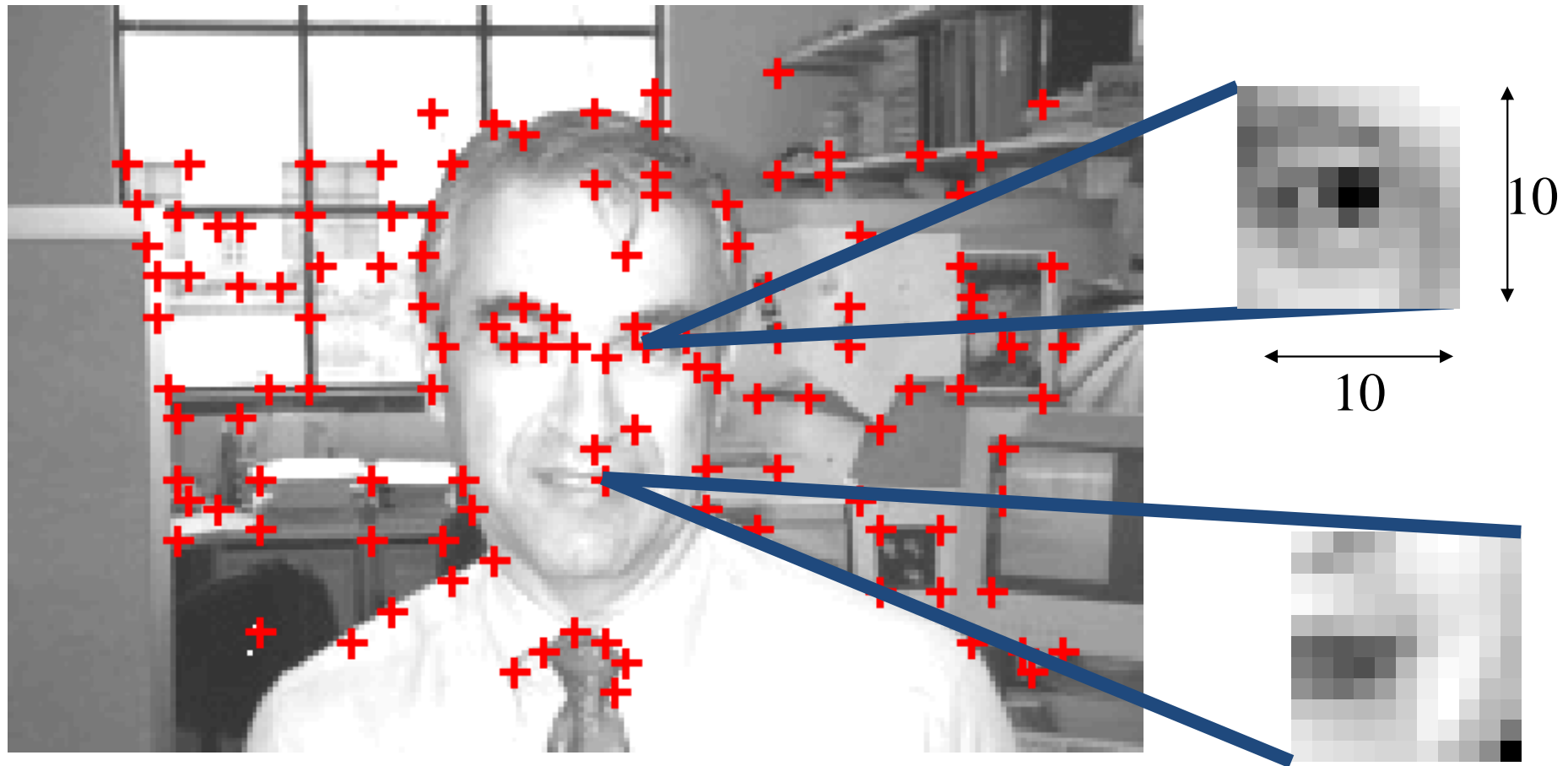
Weber & Welling et. al.

(Semi) Unsupervised learning



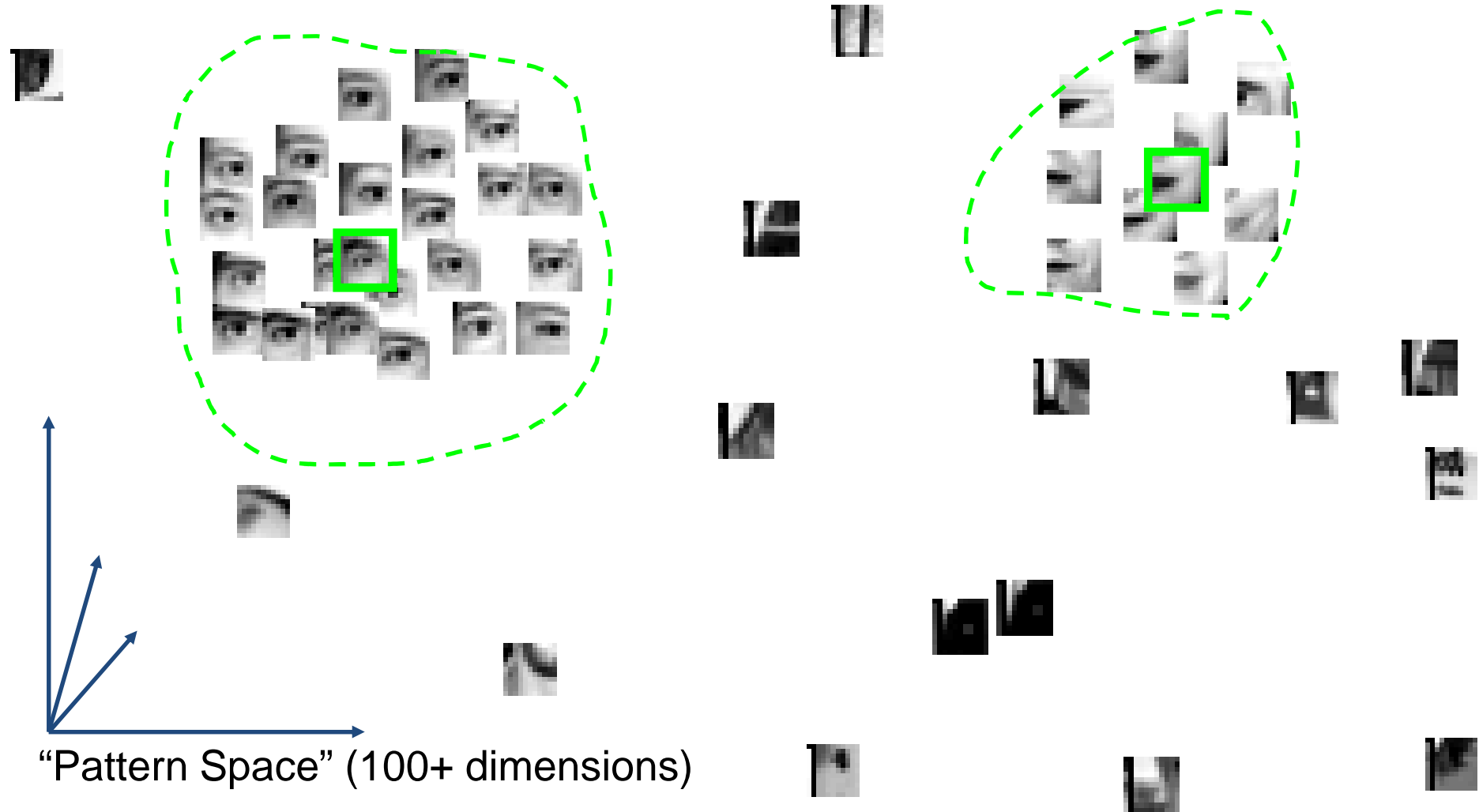
- Know if image contains object or not
- But no segmentation of object or manual selection of features

Unsupervised detector training - 1



- Highly textured neighborhoods are selected automatically
- produces 100-1000 patterns per image

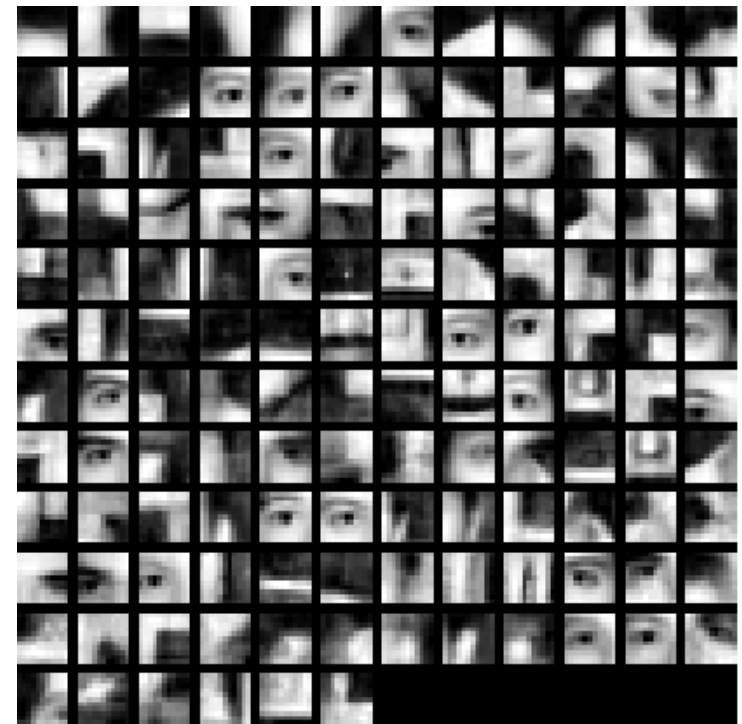
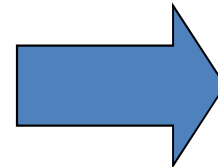
Unsupervised detector training - 2



Unsupervised detector training - 3



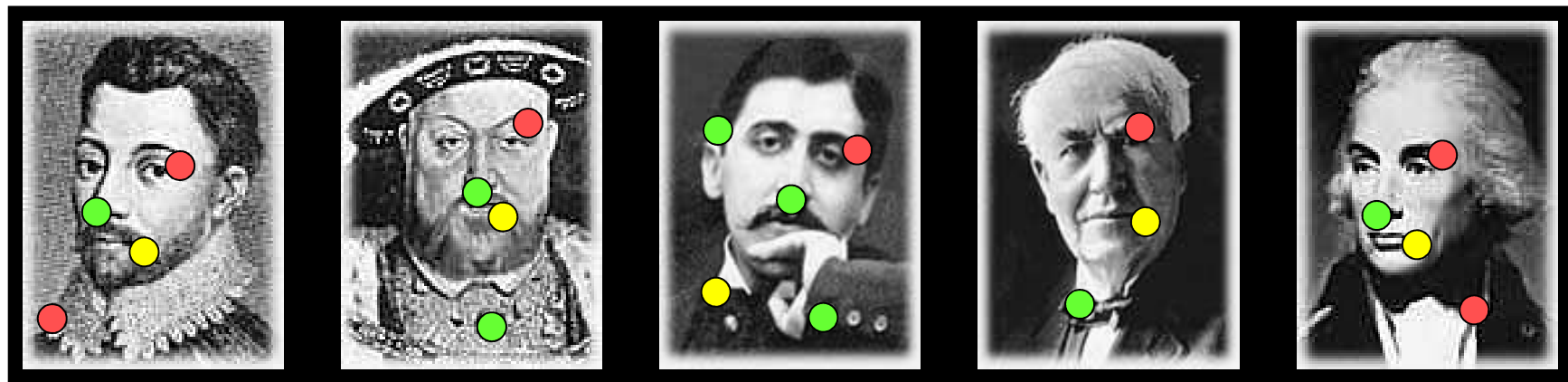
100-1000 images



~100 detectors

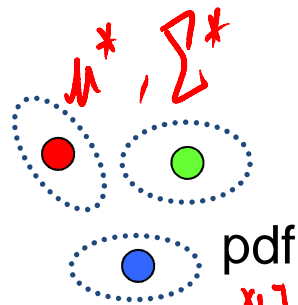
Learning

- Take training images. Pick set of detectors. Apply detectors.
- 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



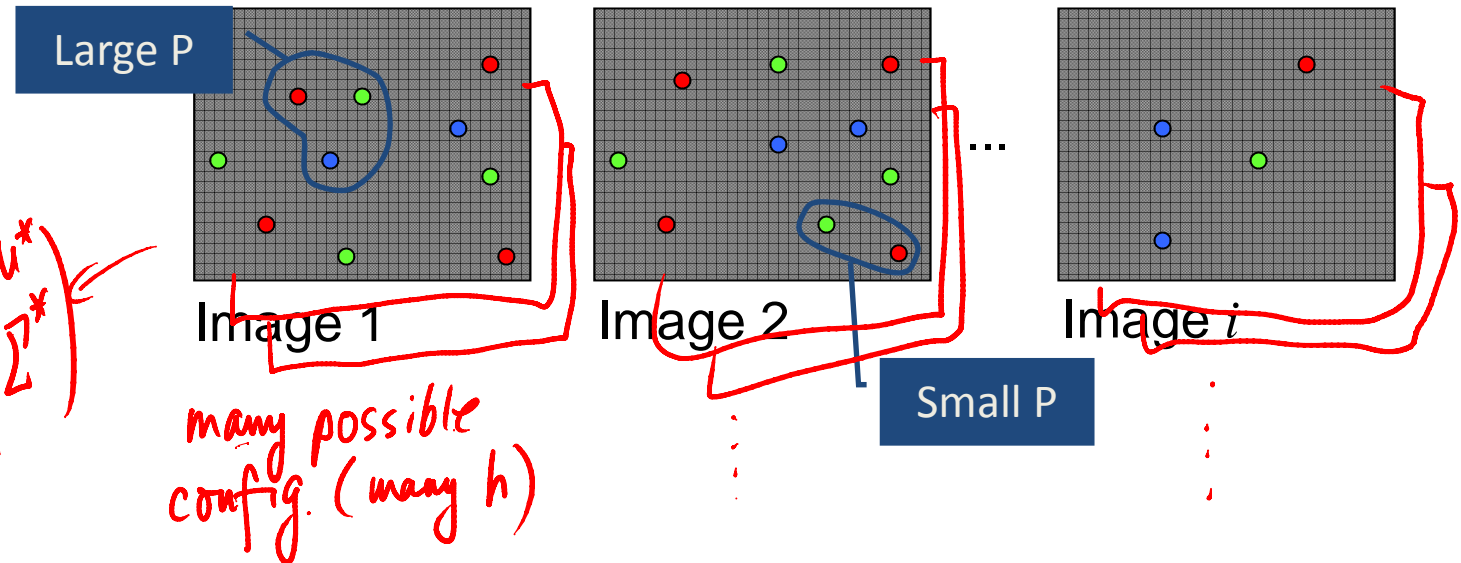
ML using EM

1. Current estimate



$$p = p\left(\begin{bmatrix} x_1 \\ y_1 \\ x_2 \\ y_2 \\ x_3 \\ y_3 \end{bmatrix} \middle| \begin{matrix} \mu^* \\ \Sigma^* \end{matrix}\right)$$

2. Assign probabilities to constellations



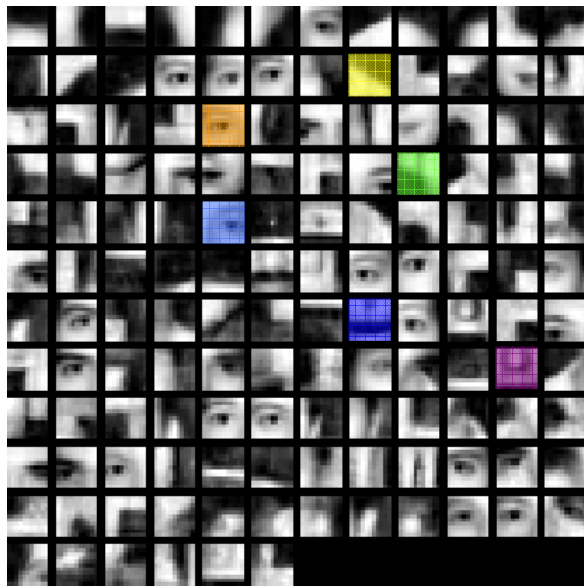
3. Use probabilities as weights to re-estimate parameters. Example: μ

$$\text{Large P} \times \begin{matrix} \bullet \\ \bullet \end{matrix} + \text{Small P} \times \begin{matrix} \bullet \\ \bullet \end{matrix} + \dots = \begin{matrix} \bullet \\ \bullet \end{matrix}$$

new estimate of μ

Detector Selection

- Try out different combinations of detectors
(Greedy search)

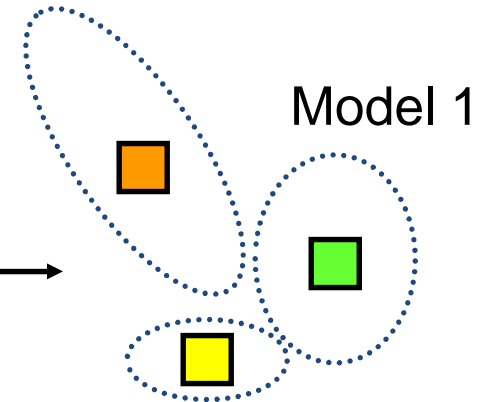


Detectors (≈ 100)

Choice 1

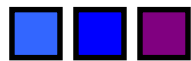


Parameter
Estimation

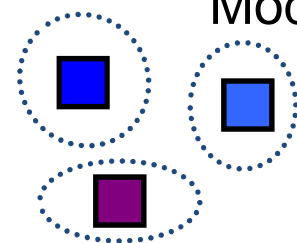


Model 1

Choice 2



Parameter
Estimation



Model 2

•
•
•

Predict / measure model performance
(validation set or directly from model)

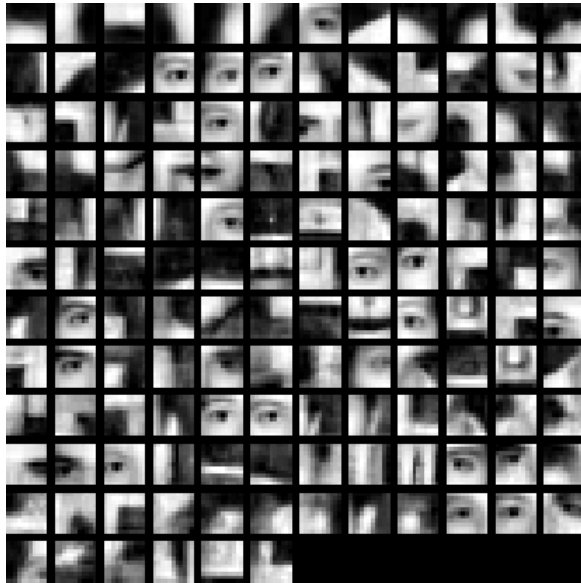
Frontal Views of Faces



- 200 Images (100 training, 100 testing)
- 30 people, different for training and testing

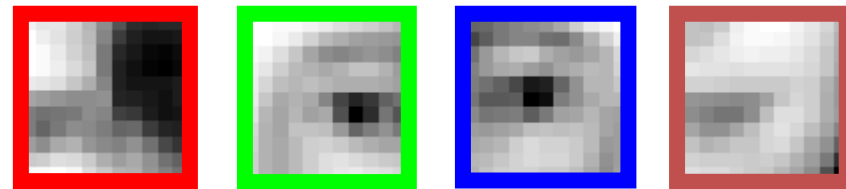
Learned face model

Pre-selected Parts

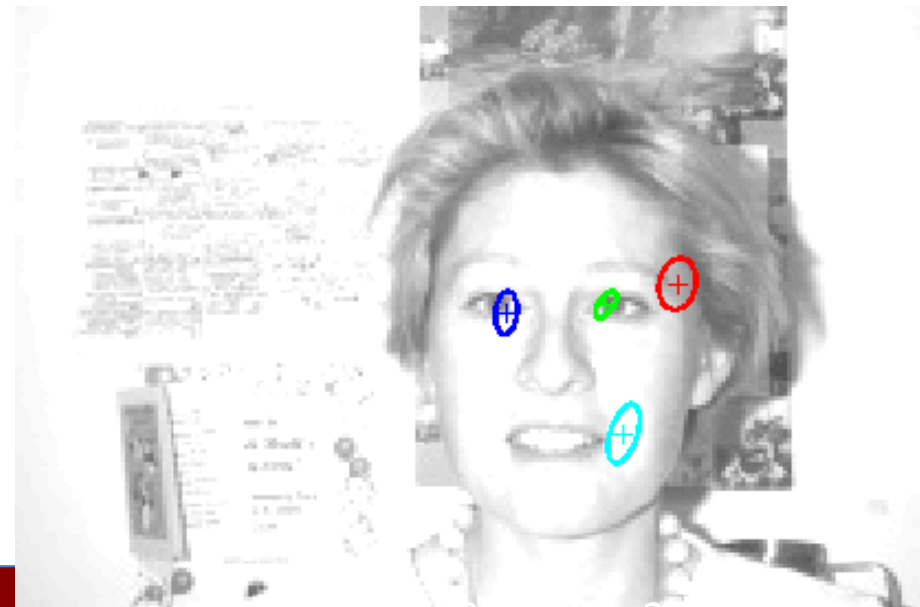


Test Error: 6% (4 Parts)

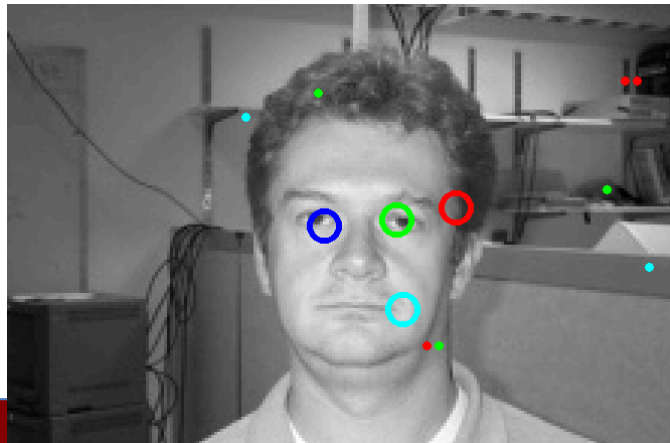
Parts in Model



Model Foreground pdf



Sample Detection



Fei-Fei Li

Lecture 16 - 28

28-NOV-11

Face images



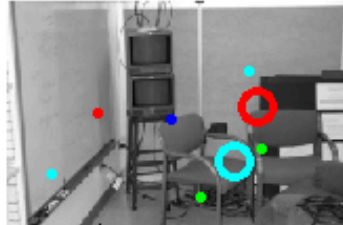
28 Nov 11

Background images

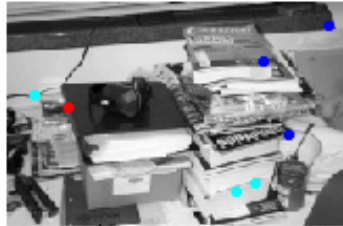
incorrect



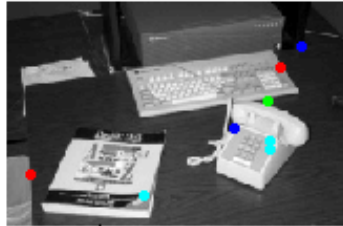
incorrect



correct



correct



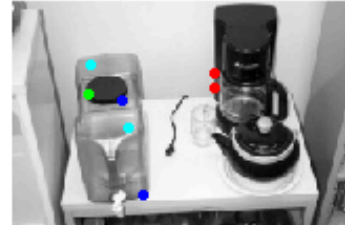
correct



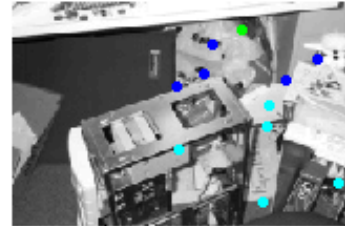
correct



correct



correct



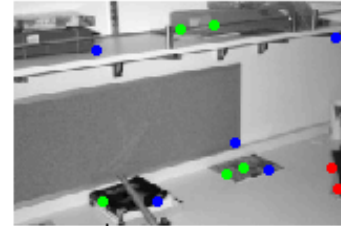
correct



correct



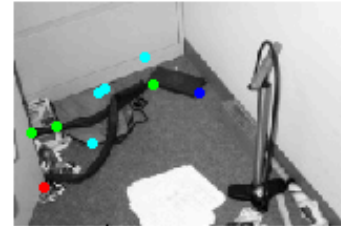
correct



correct



correct



correct



correct



correct



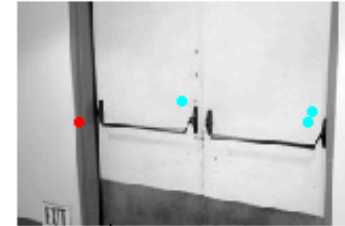
correct



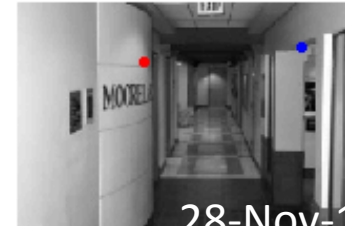
correct



correct



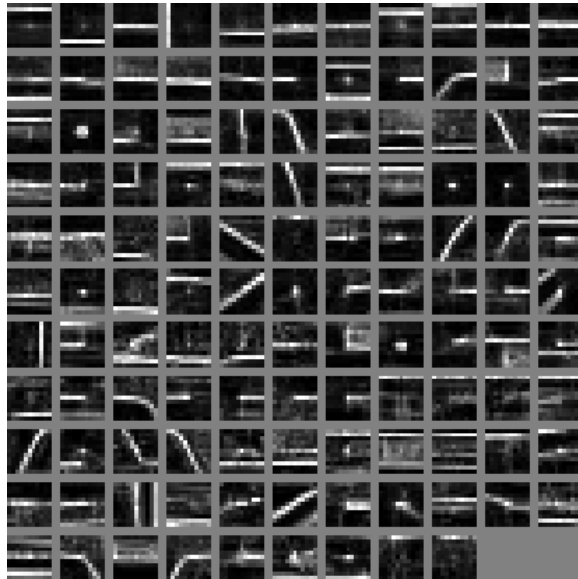
correct



28-Nov-1

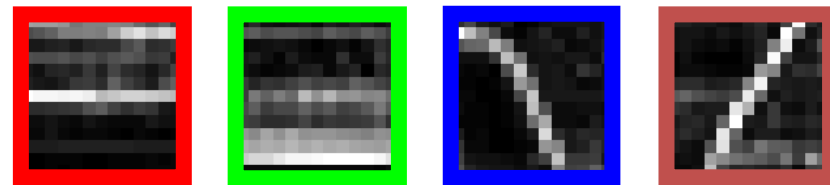
Car from Rear

Preselected Parts

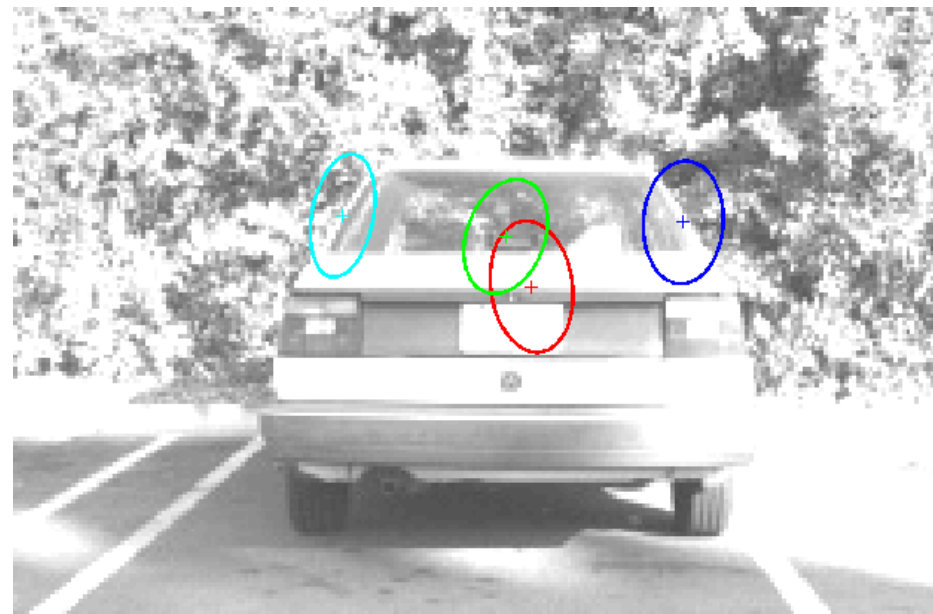


Test Error: 13% (5 Parts)

Parts in Model



Model Foreground pdf



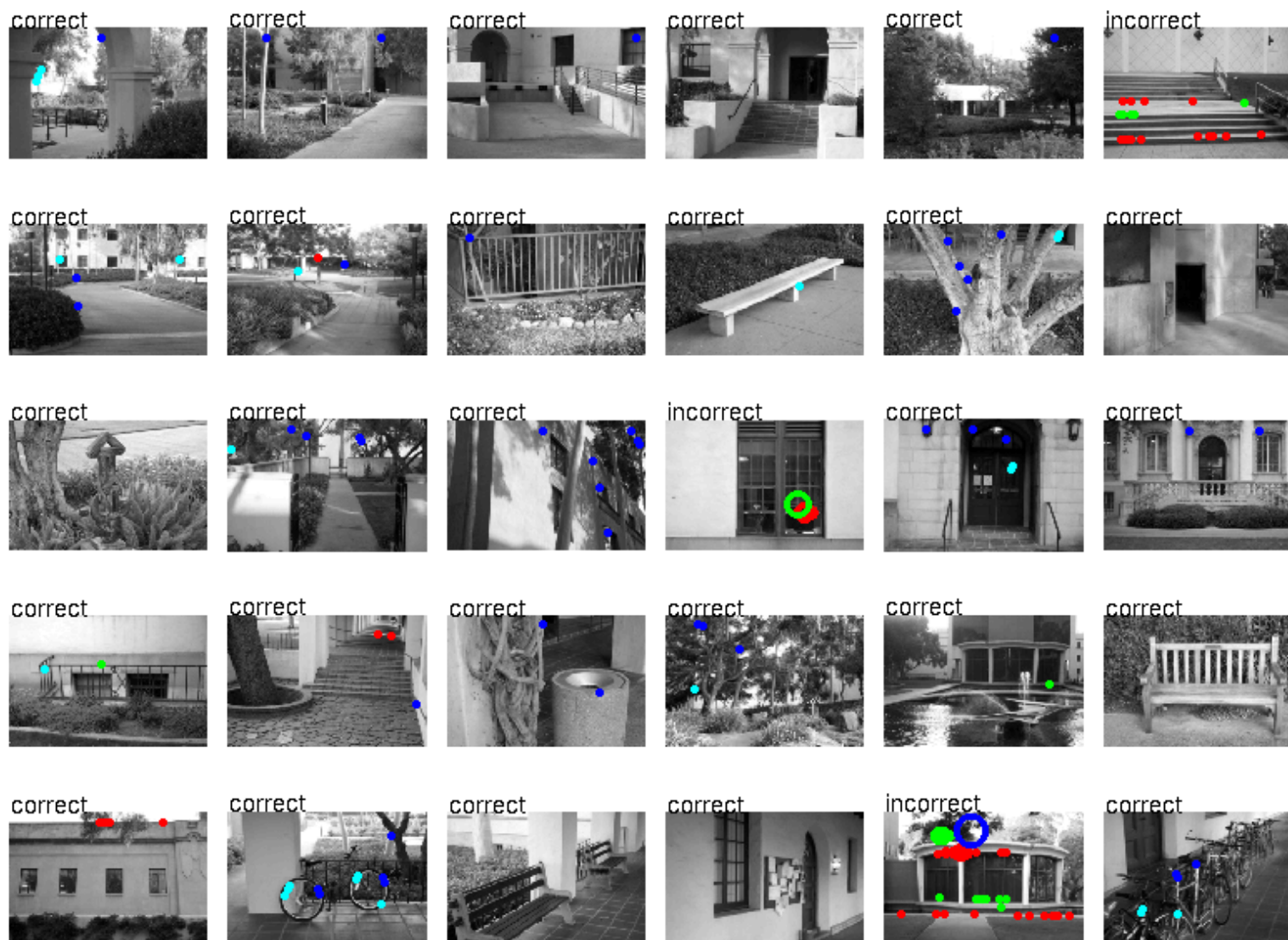
Sample Detection



Detections of Cars



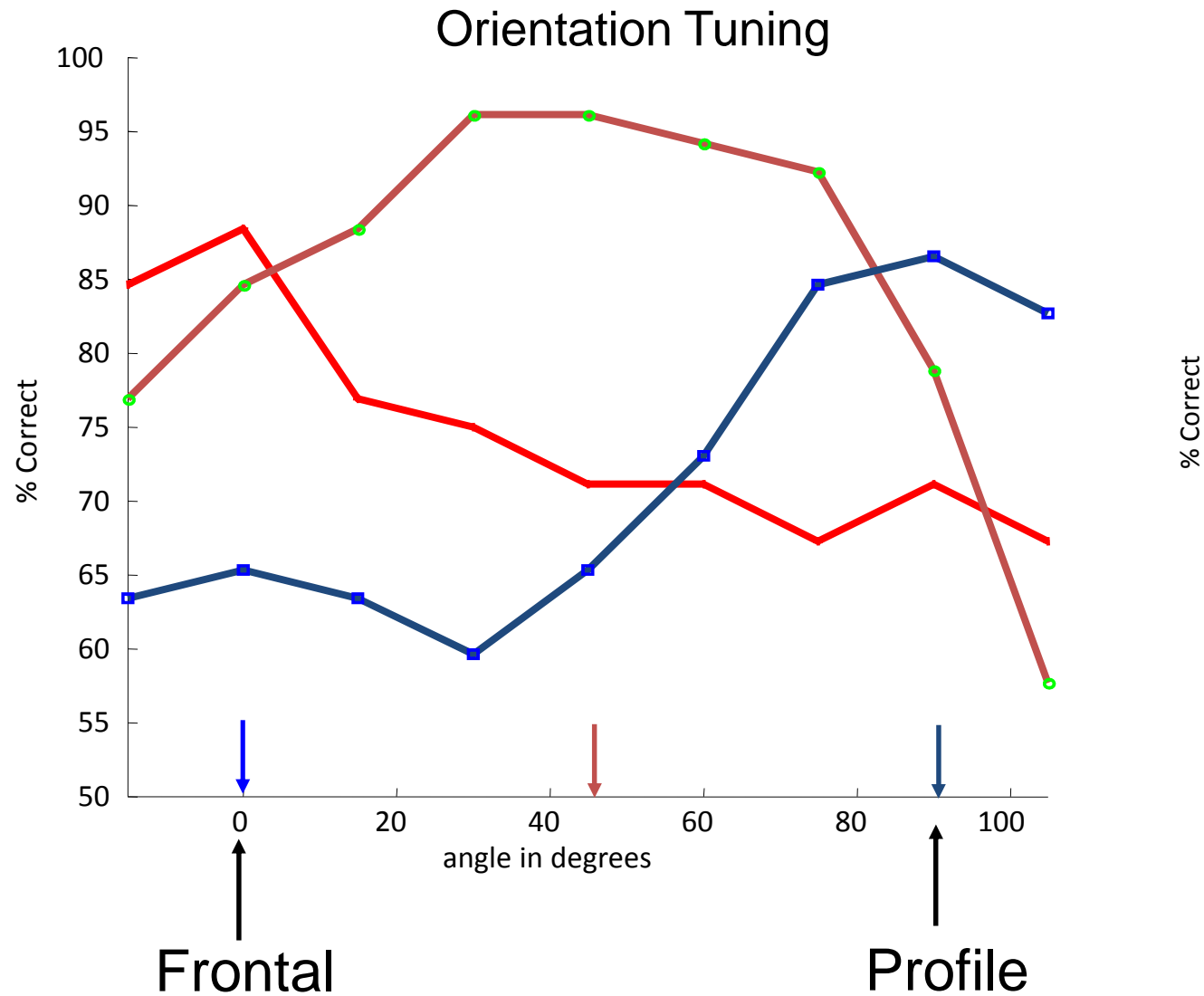
Background Images



3D Object recognition – Multiple mixture components



3D Orientation Tuning



So far (2).....

- Representation
 - Multiple mixture components for different viewpoints
- Learning
 - Now semi-unsupervised
 - Automatic construction and selection of part detectors
 - Estimation of parameters using EM
- Recognition
 - As before
- Issues:
 - Learning is slow (many combinations of detectors)
 - Appearance learnt first, then shape

Issues

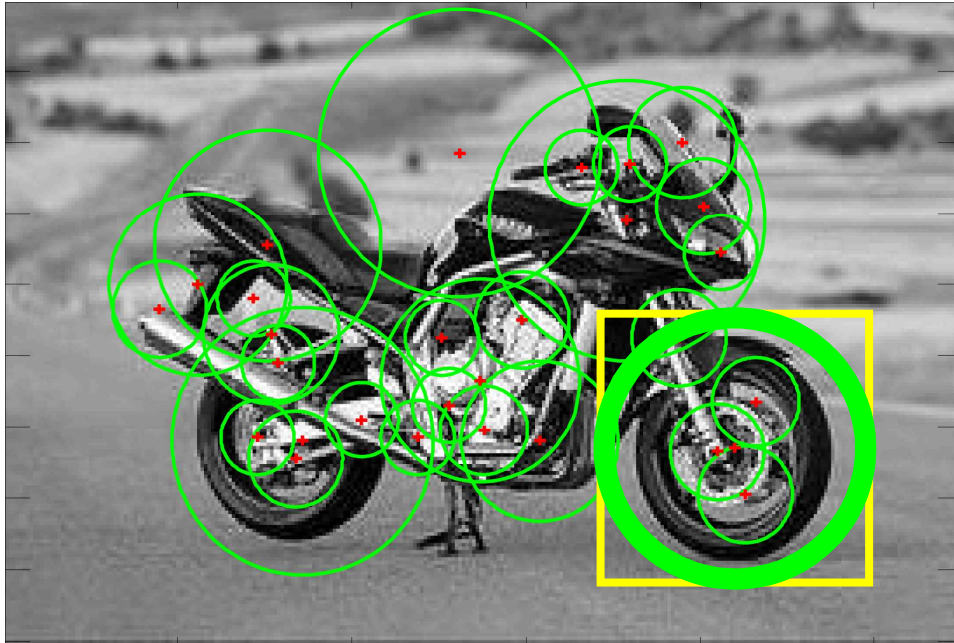
- Speed of learning
 - Slow (many combinations of detectors)
- Appearance learnt first, then shape
 - Difficult to learn part that has stable location but variable appearance
 - Each detector is used as a cross-correlation filter, giving a hard definition of the part's appearance
- Would like a fully probabilistic representation of the object

Object categorization

Fergus et. al.

CVPR '03, IJCV '06

Detection & Representation of regions



Appearance

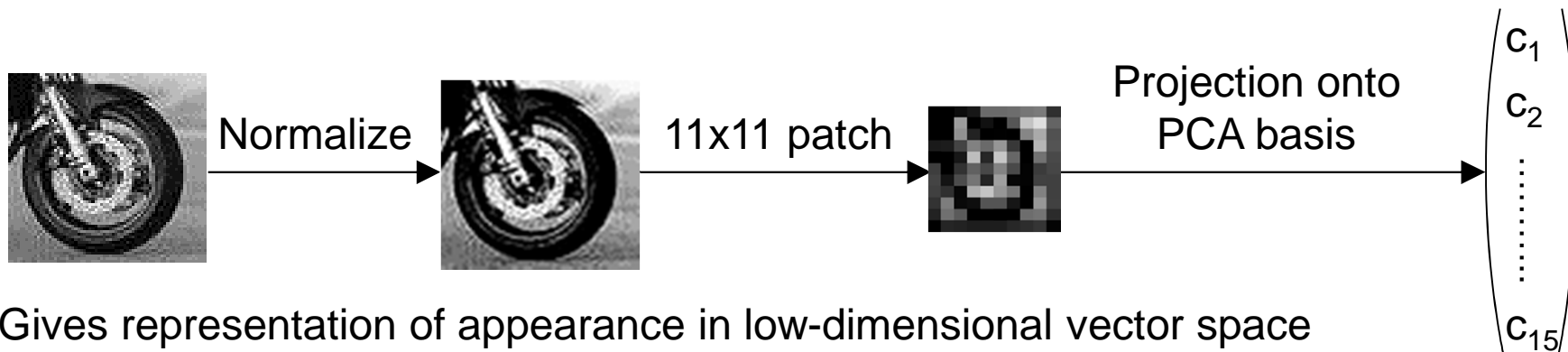
- Find regions within image
- Use salient region operator (Kadir & Brady 01)

Location

(x,y) coords. of region centre

Scale

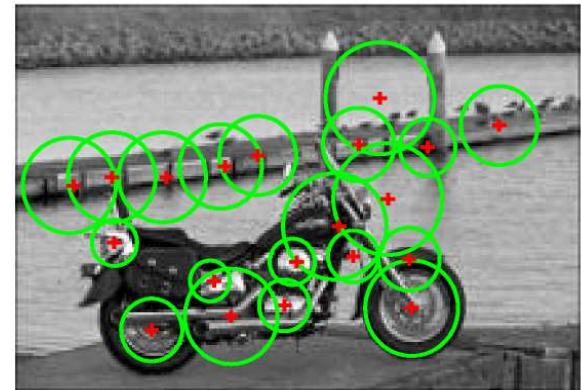
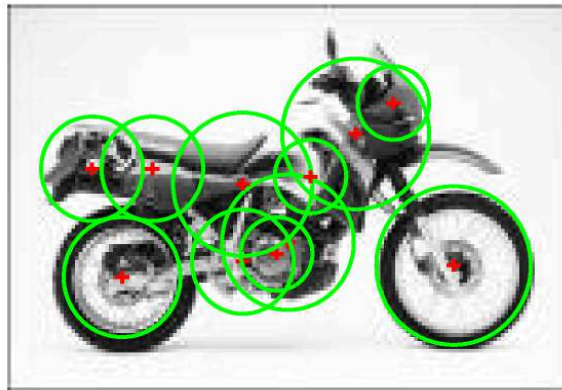
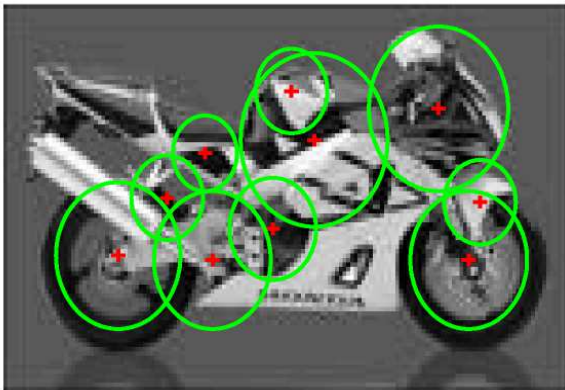
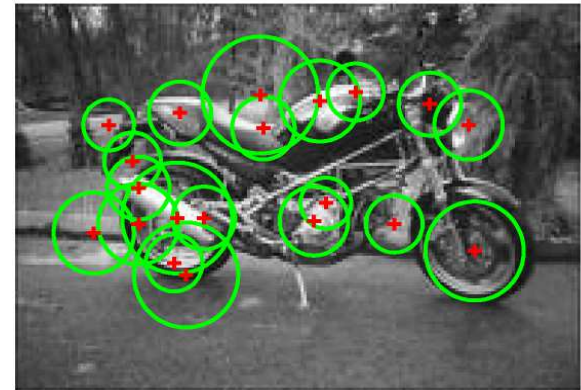
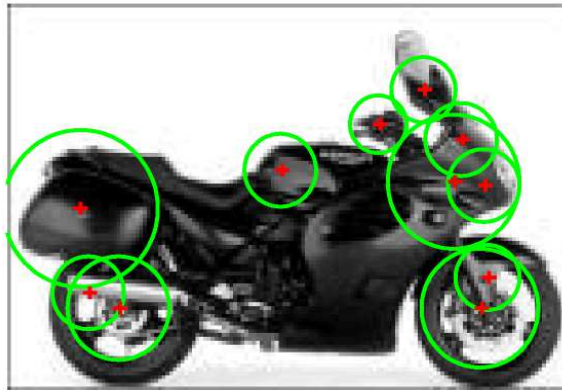
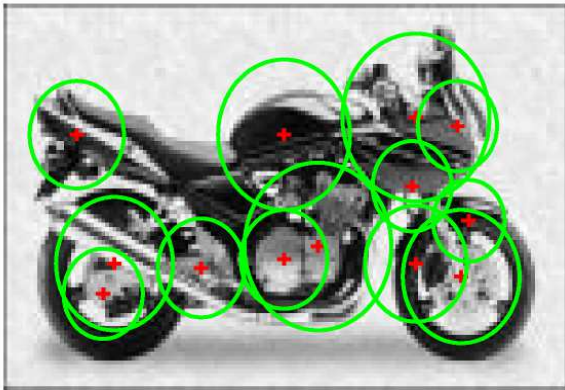
Radius of region (pixels)



Gives representation of appearance in low-dimensional vector space

Motorbikes example

- Kadir & Brady saliency region detector

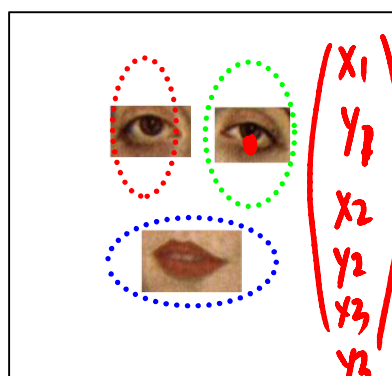


Generative probabilistic model (2)

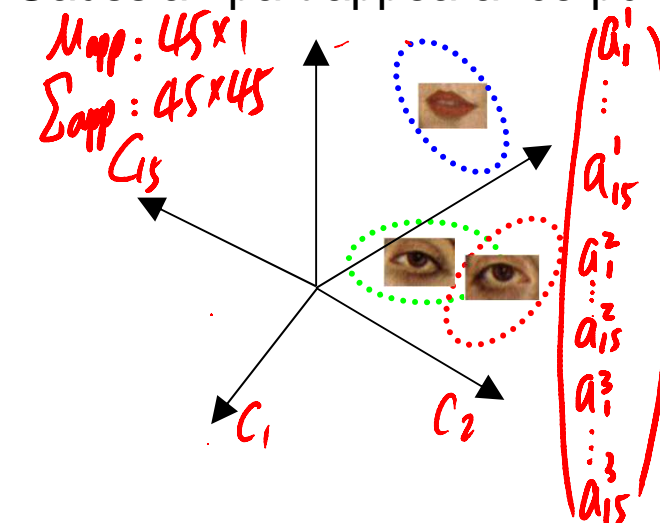
Foreground model

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

Gaussian shape pdf

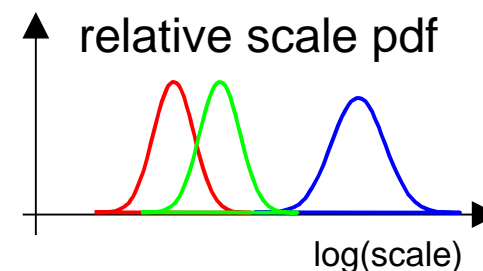


Gaussian part appearance pdf

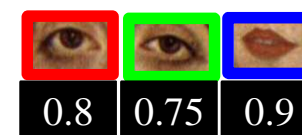


Gaussian

relative scale pdf

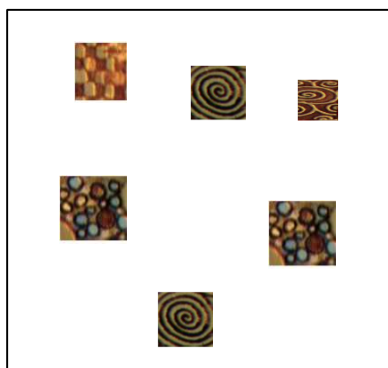


Prob. of detection

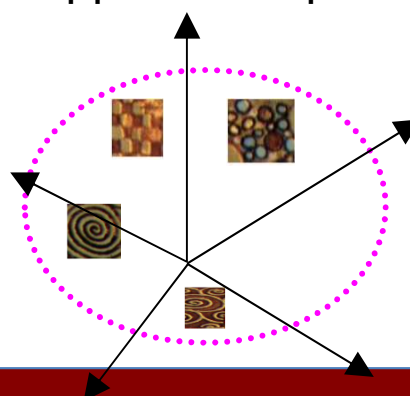


Clutter model

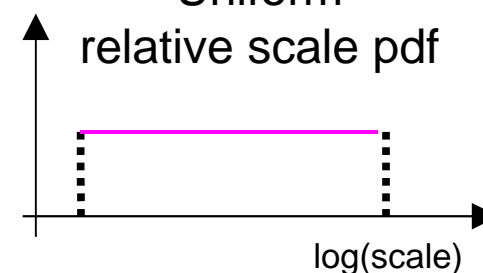
Uniform shape pdf



Gaussian background appearance pdf



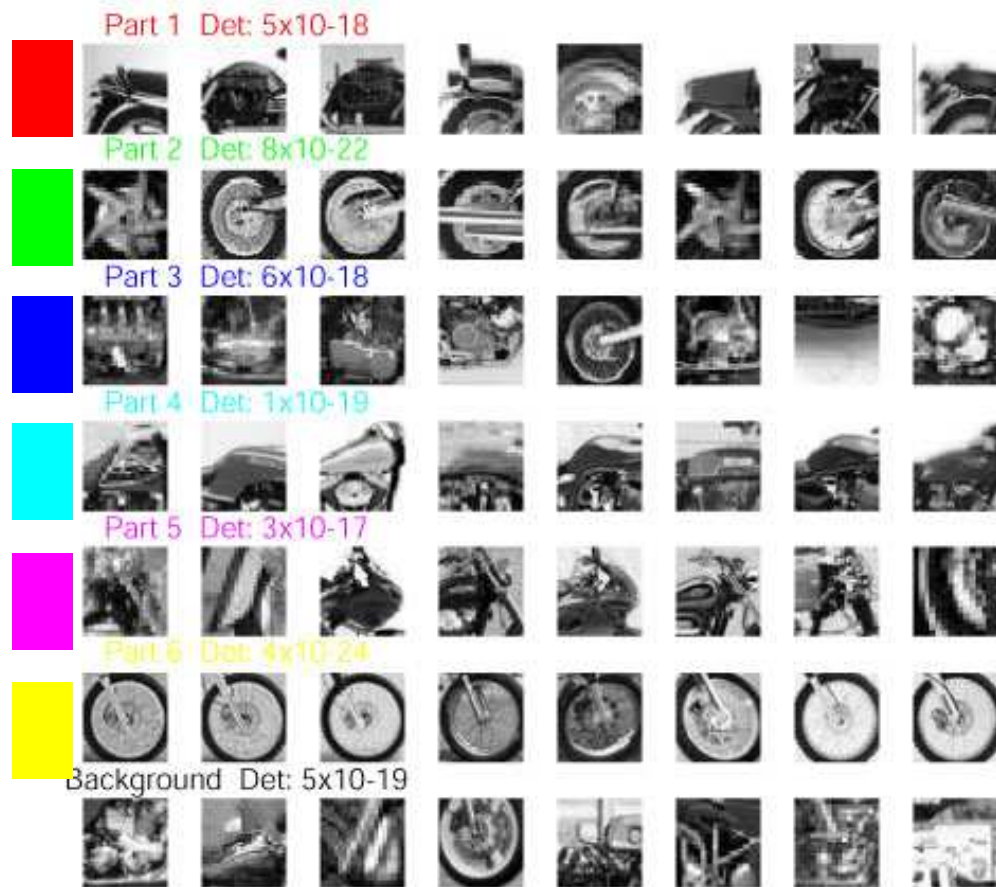
Uniform



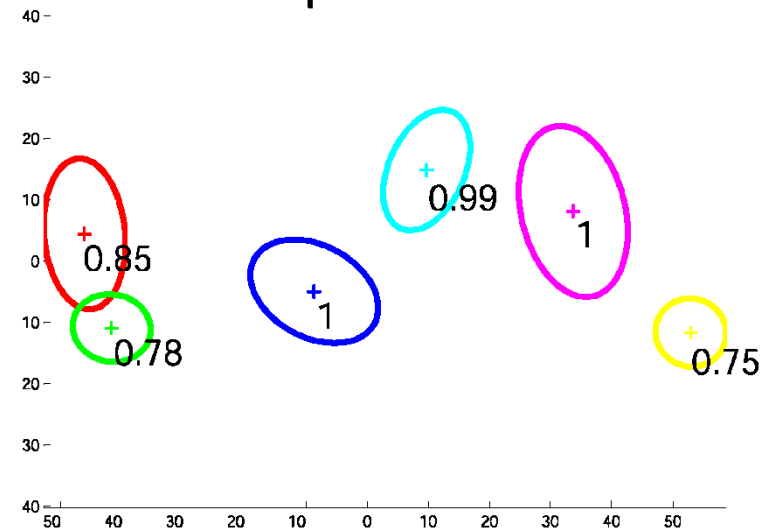
Poisson pdf on # detections

Motorbikes

Samples from appearance model

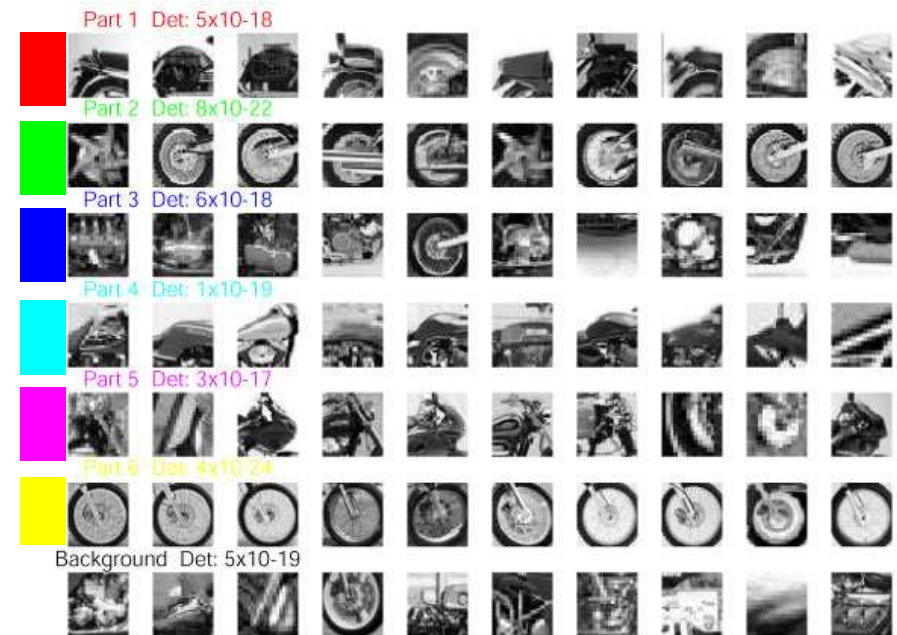
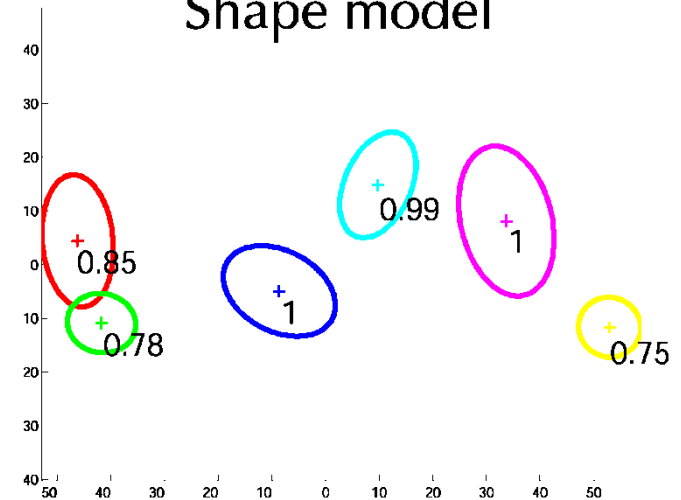
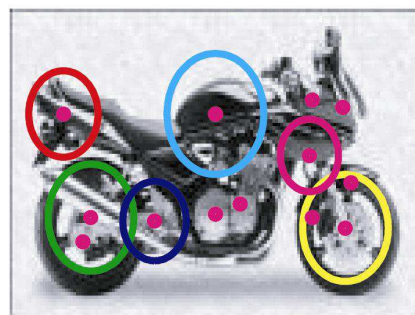
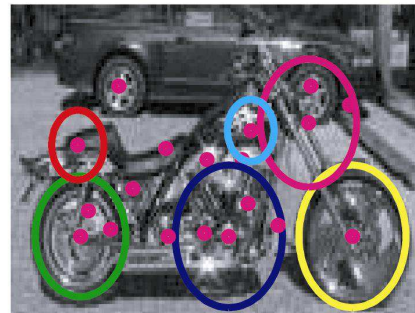
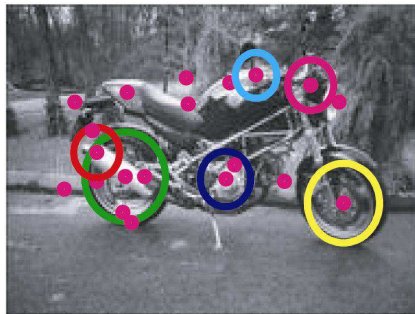
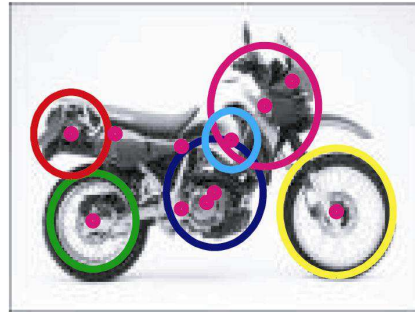
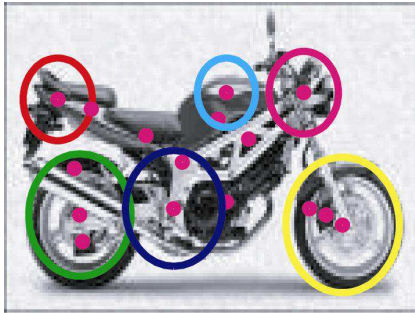


Shape model



Recognized Motorbikes

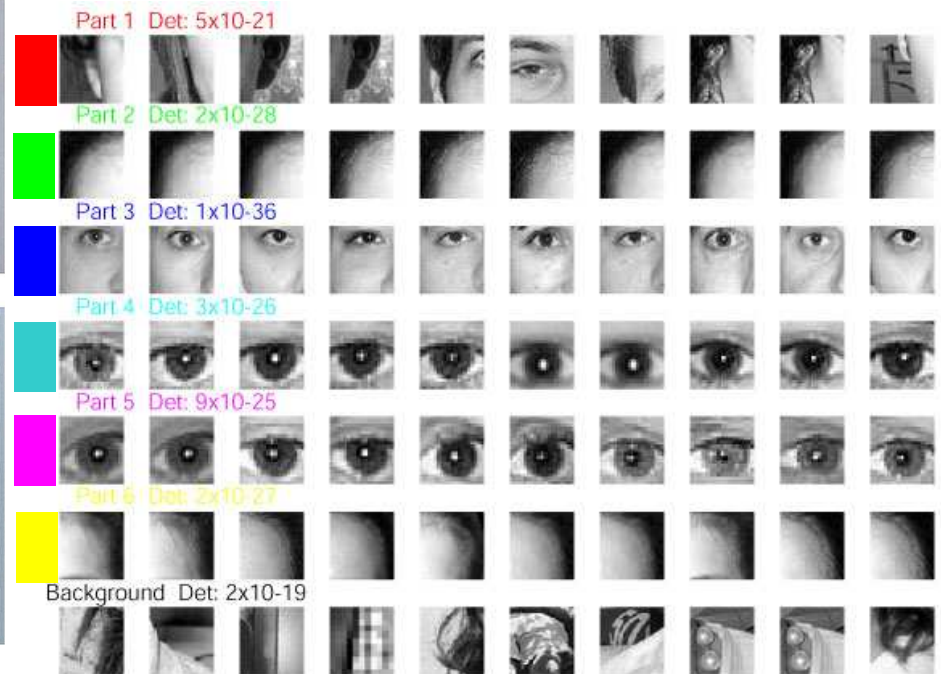
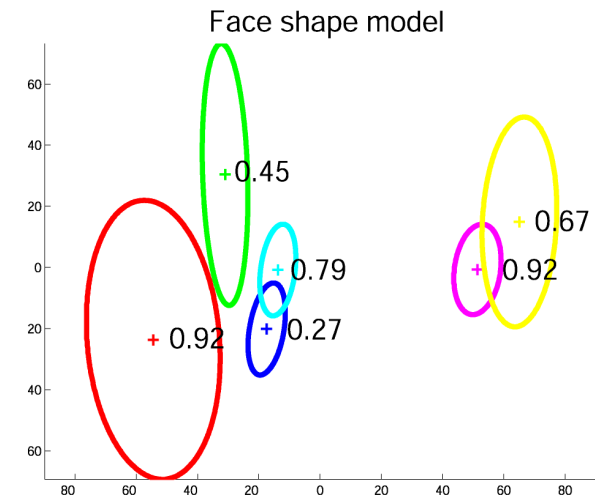
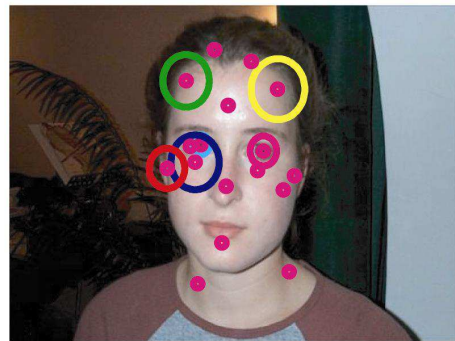
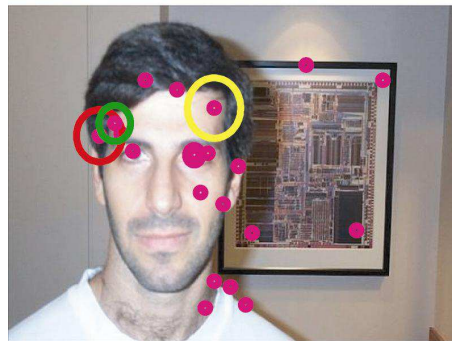
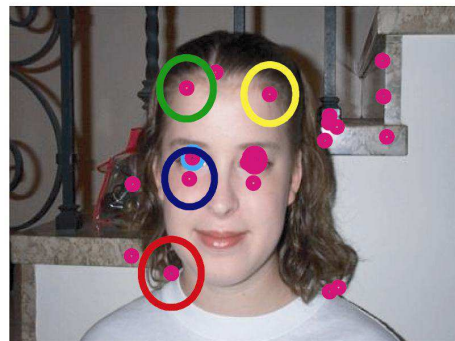
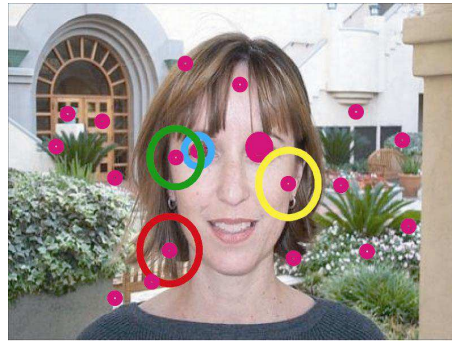
Shape model



Background images evaluated with motorbike model



Frontal faces

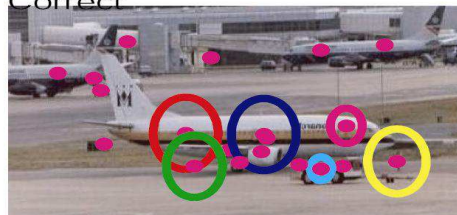


Airplanes

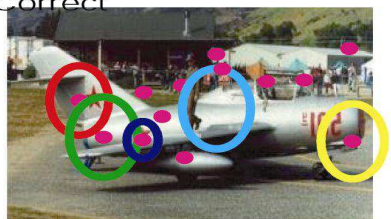
INCORRECT



Correct



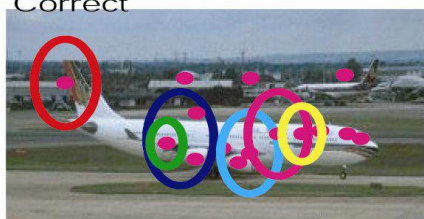
Correct



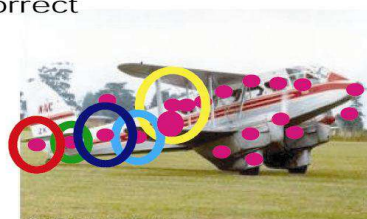
Correct



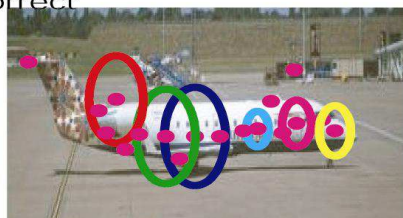
Correct



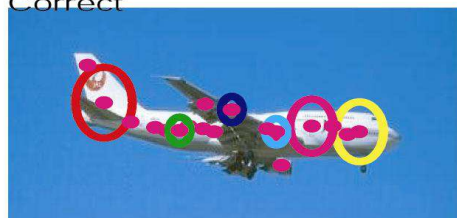
Correct



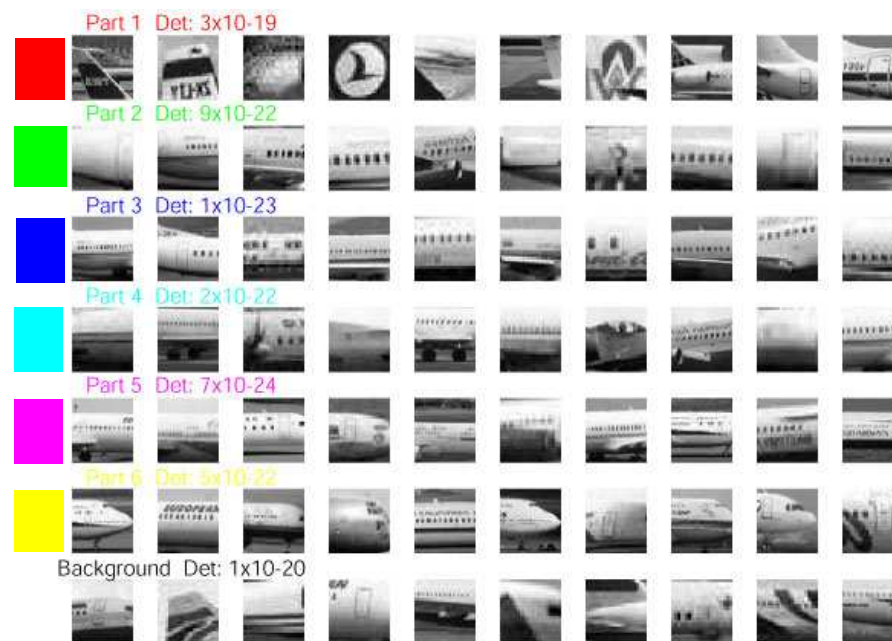
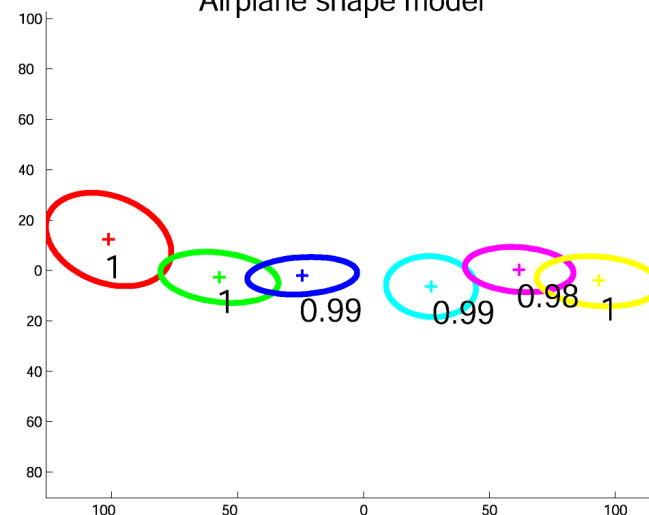
Correct



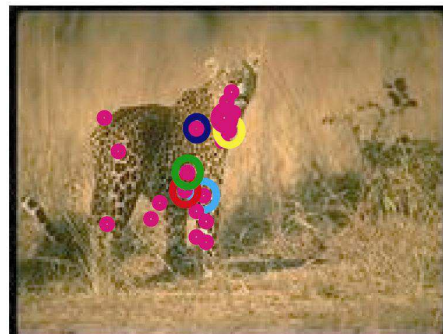
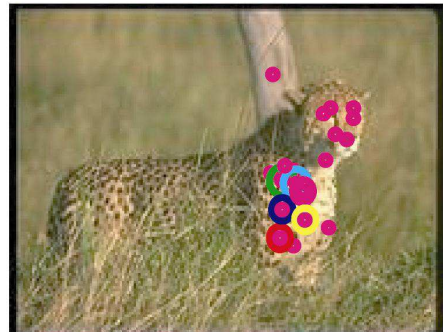
Correct



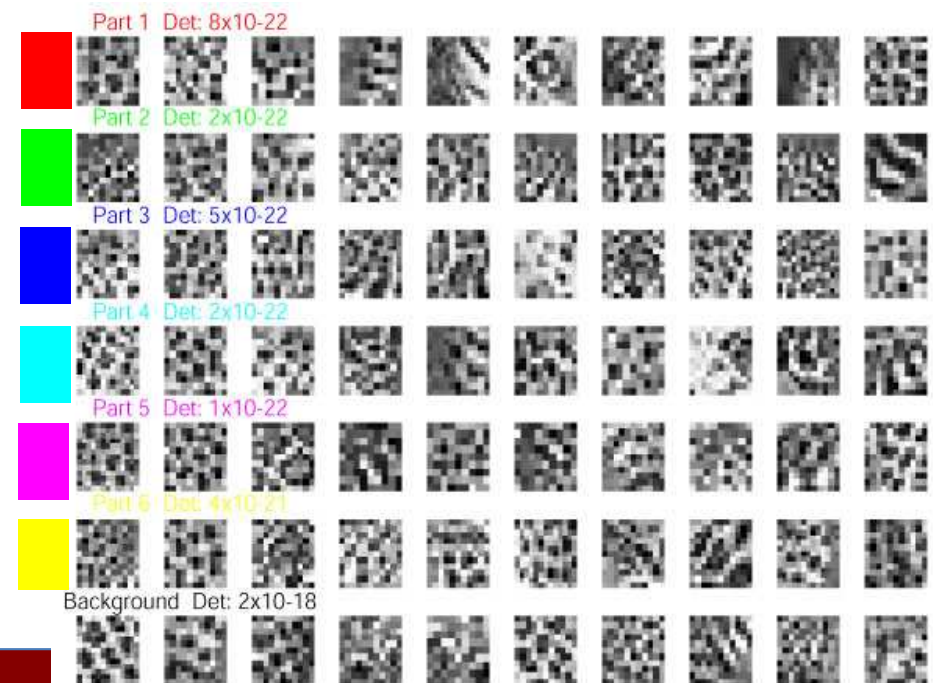
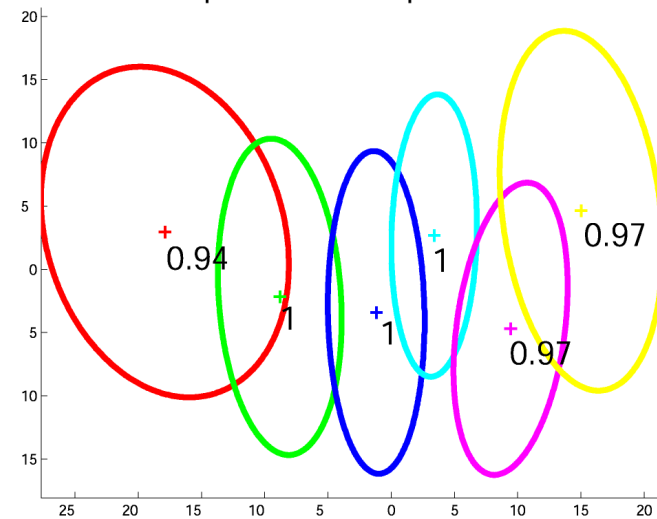
Airplane shape model



Spotted cats



Spotted cat 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

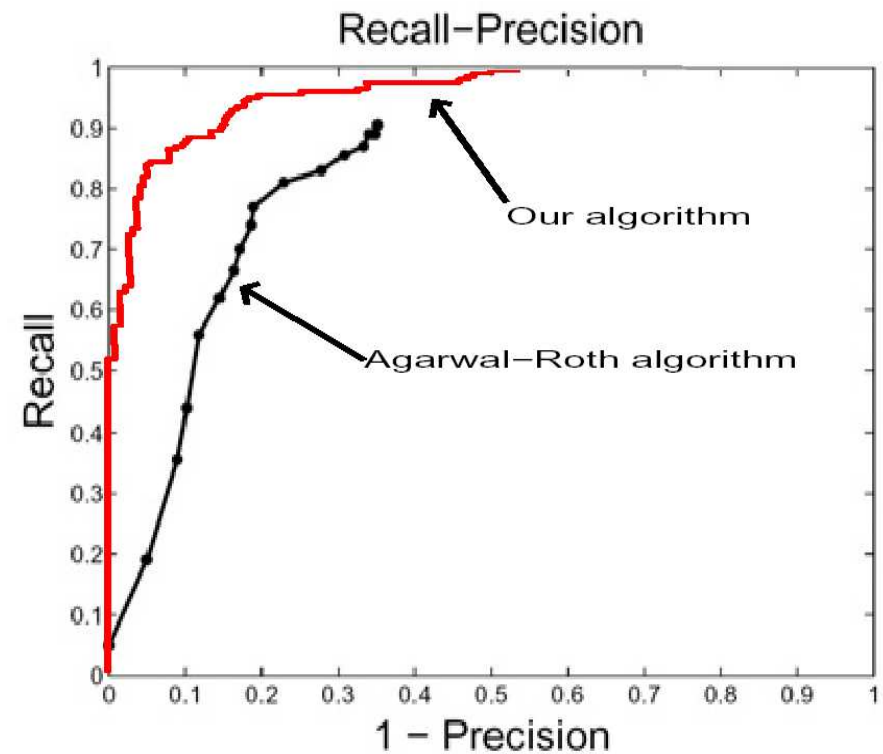
% equal error rate

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

Comparison to other methods

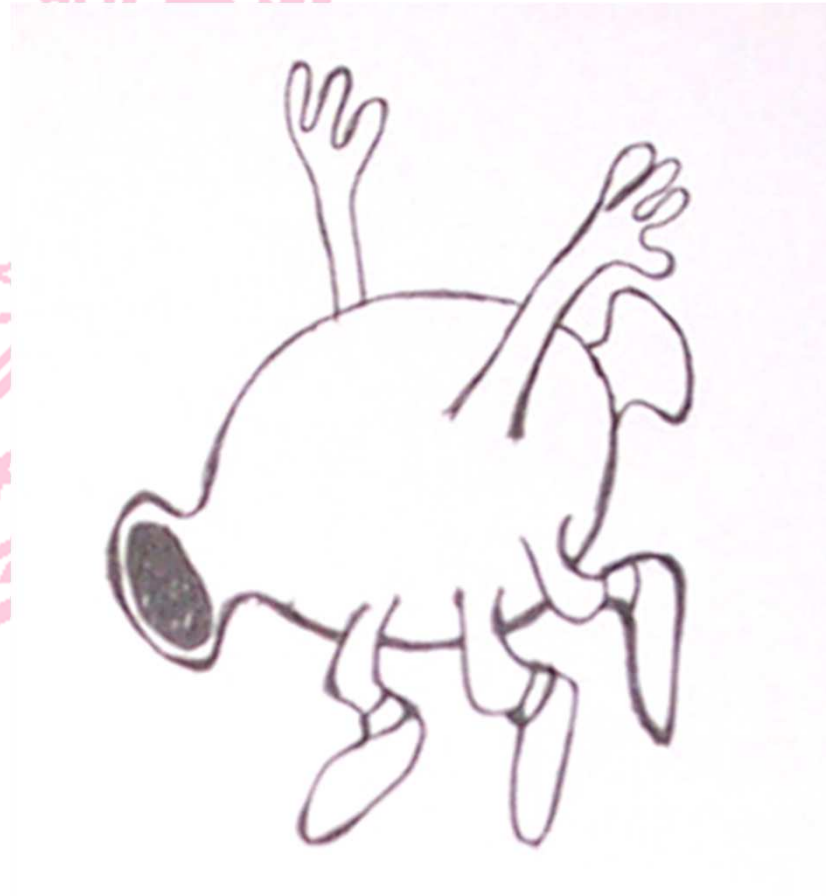
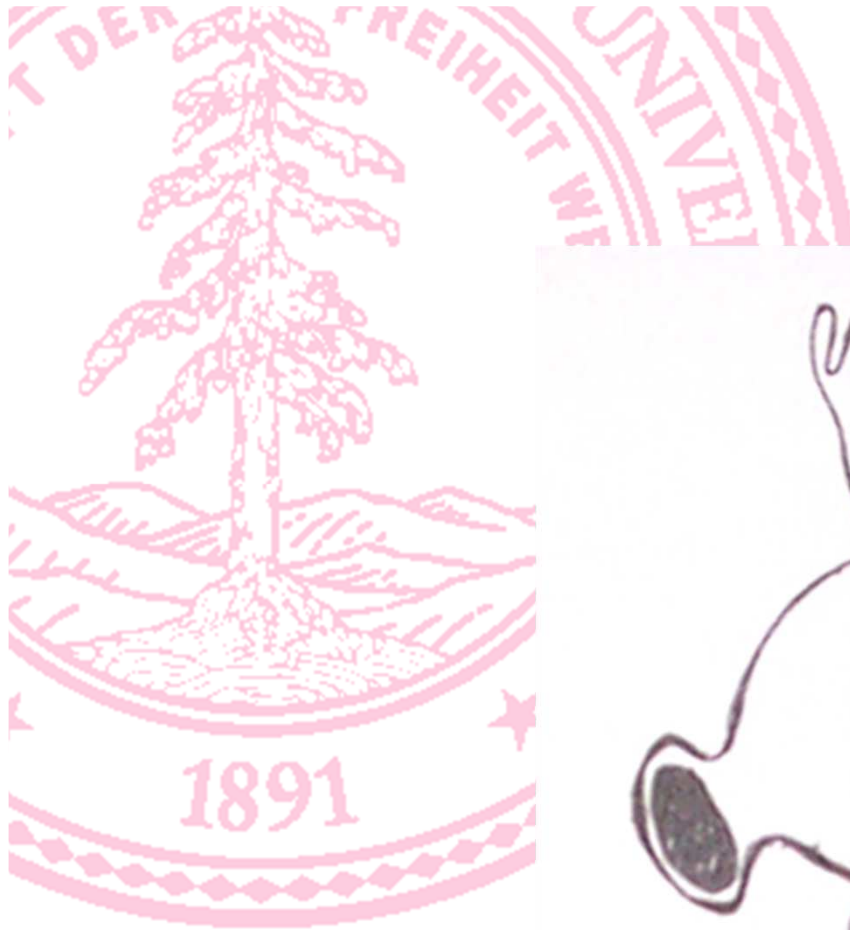
Dataset	Ours	Others	
Motorbikes	7.5	16.0	Weber et al. [ECCV '00]
Faces	4.6	6.0	Weber
Airplanes	9.8	32.0	Weber
Cars (Side)	11.5	21.0	Agarwal Roth [ECCV '02]

% equal error rate



Why this design?

- Generic features seem to well in finding consistent parts of the object
- Some categories perform badly – different feature types needed
- Why PCA representation?
 - Tried ICA, FLD, Oriented filter responses etc.
 - But PCA worked best
- Fully probabilistic representation lets us use tools from machine learning community



S. Savarese, 2003



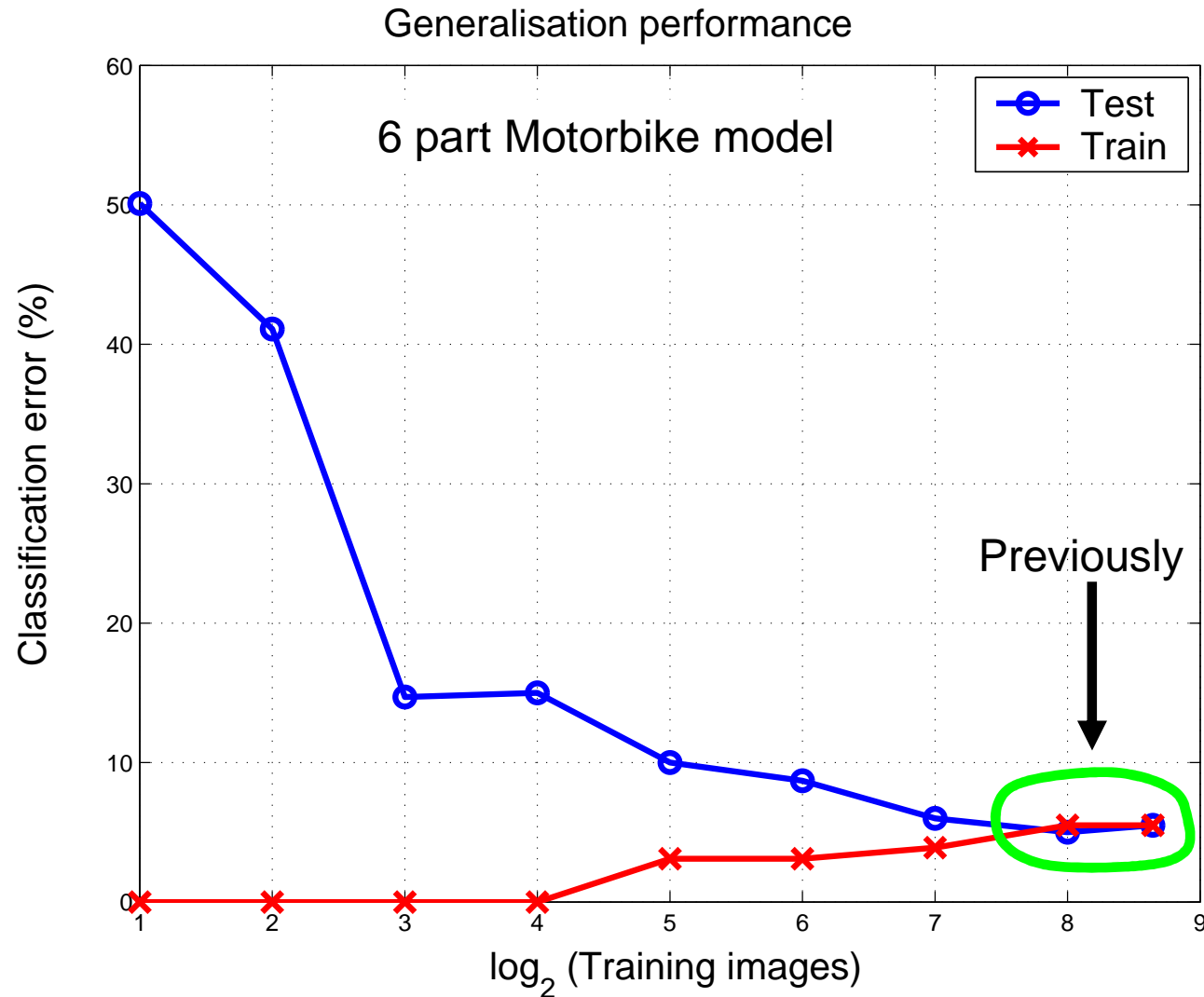
One-Shot learning

Fei-Fei et. al.

ICCV '03, PAMI '06

Algorithm	Training Examples	Categories
Burl, et al. Weber, et al. Fergus, et al.	200 ~ 400	Faces, Motorbikes, Spotted cats, Airplanes, Cars
Viola et al.	~10,000	Faces
Schneiderman, et al.	~2,000	Faces, Cars
Rowley et al.	~500	Faces

Number of training examples

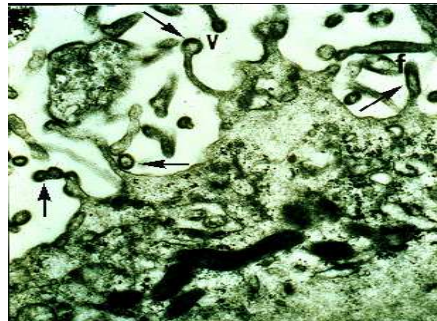


How do we do better than what statisticians have told us?

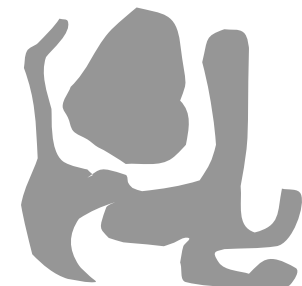
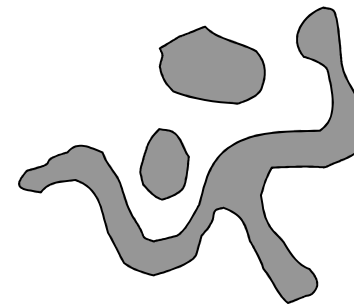
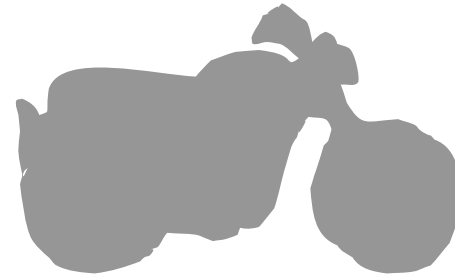
- Intuition 1: use **Prior** information
- Intuition 2: make best use of training information

Prior knowledge: means

Appearance



Shape



likely

unlikely

Bayesian framework

$P(\text{object} \mid \text{test}, \text{train})$ vs. $P(\text{clutter} \mid \text{test}, \text{train})$

Bayes Rule

$p(\text{test} \mid \text{object}, \text{train}) p(\text{object})$

Expansion by parametrization

$\int p(\text{test} \mid \theta, \text{object}) p(\theta \mid \text{object}, \text{train}) d\theta$

Bayesian framework

$P(\text{object} \mid \text{test}, \text{train})$ vs. $P(\text{clutter} \mid \text{test}, \text{train})$

Bayes Rule

$p(\text{test} \mid \text{object}, \text{train}) p(\text{object})$

Expansion by parametrization

$\int p(\text{test} \mid \theta, \text{object}) p(\theta \mid \text{object}, \text{train}) d\theta$

Previous Work: $\delta(\theta^{\text{ML}})$

Bayesian framework

$P(\text{object} \mid \text{test, train})$ vs. $P(\text{clutter} \mid \text{test, train})$

Bayes Rule

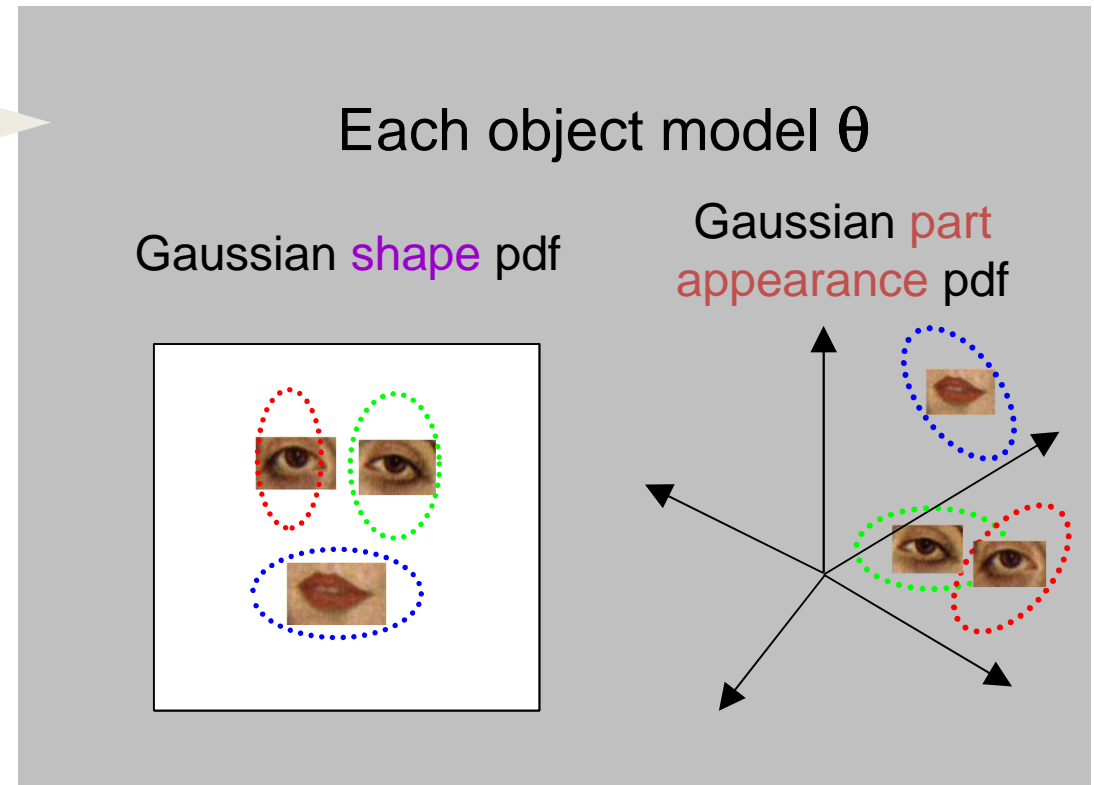
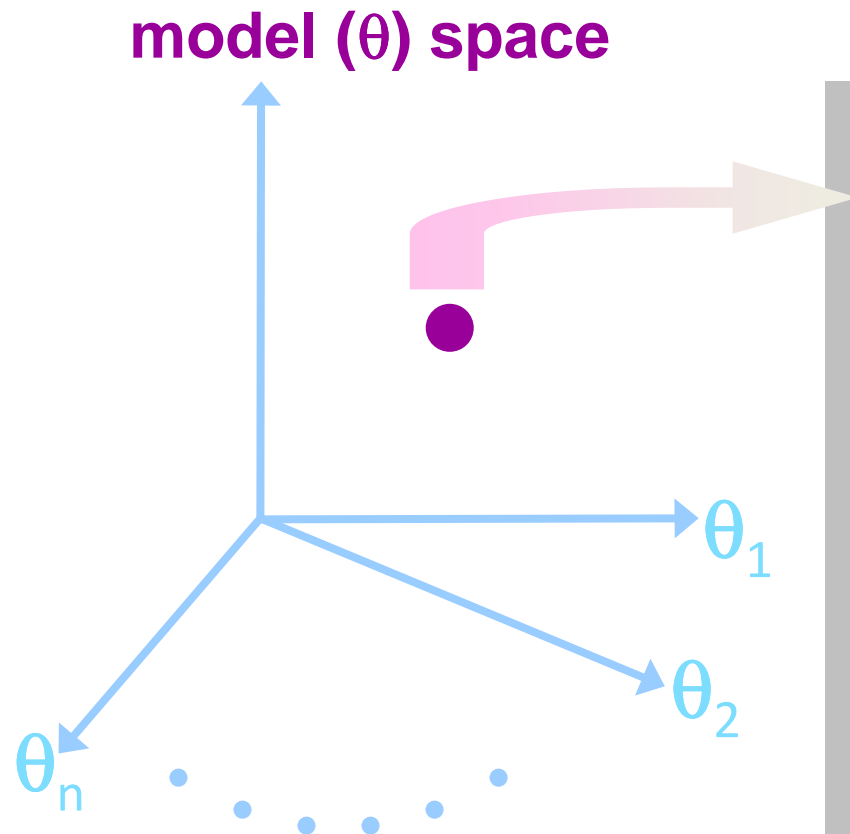
$$p(\text{test} \mid \text{object, train}) p(\text{object})$$

Expansion by parametrization

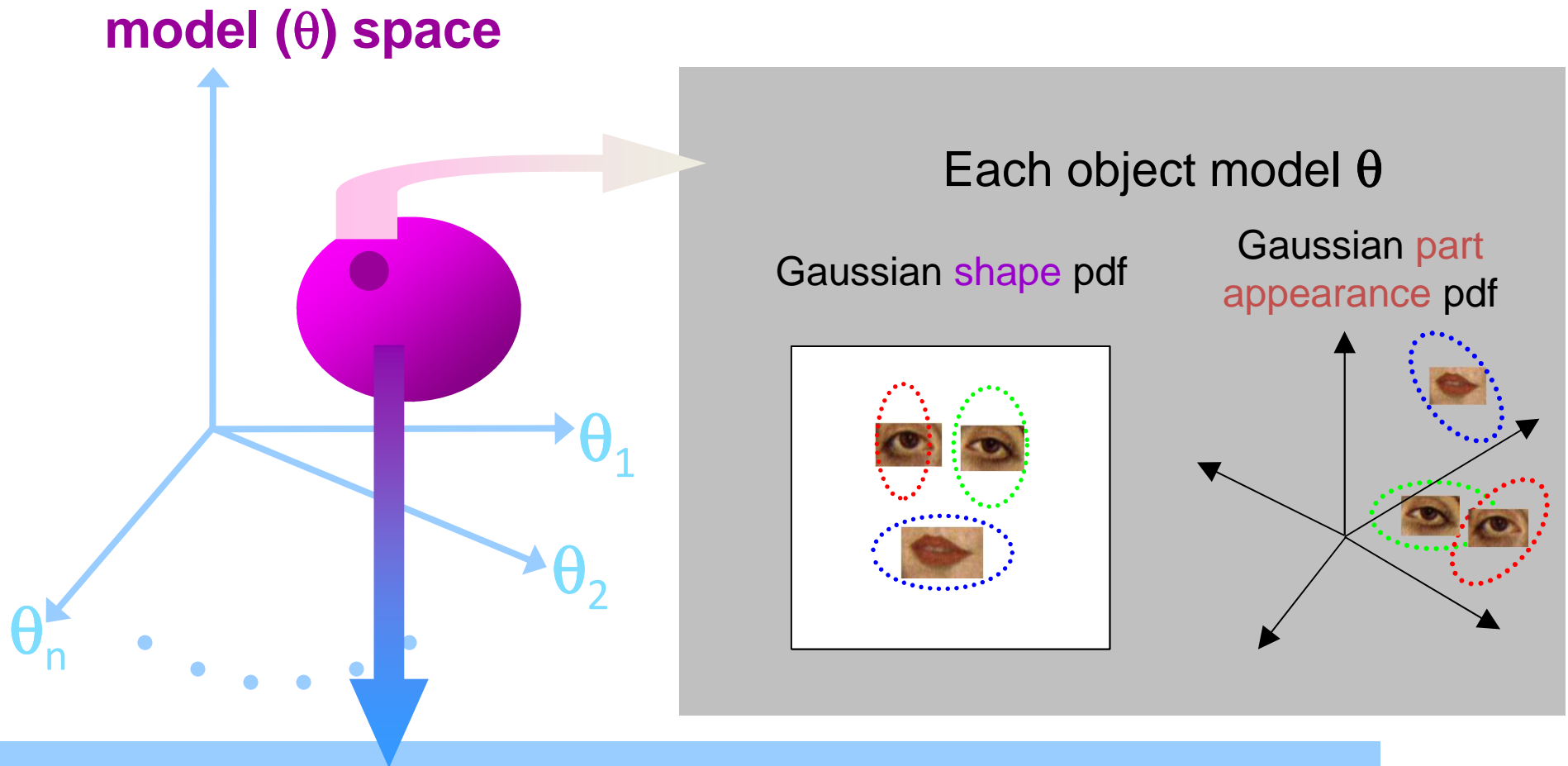
$$\int p(\text{test} \mid \theta, \text{object}) p(\theta \mid \text{object, train}) d\theta$$

One-Shot learning: $p(\text{train} \mid \theta, \text{object}) p(\theta)$

Model Structure



Model Structure

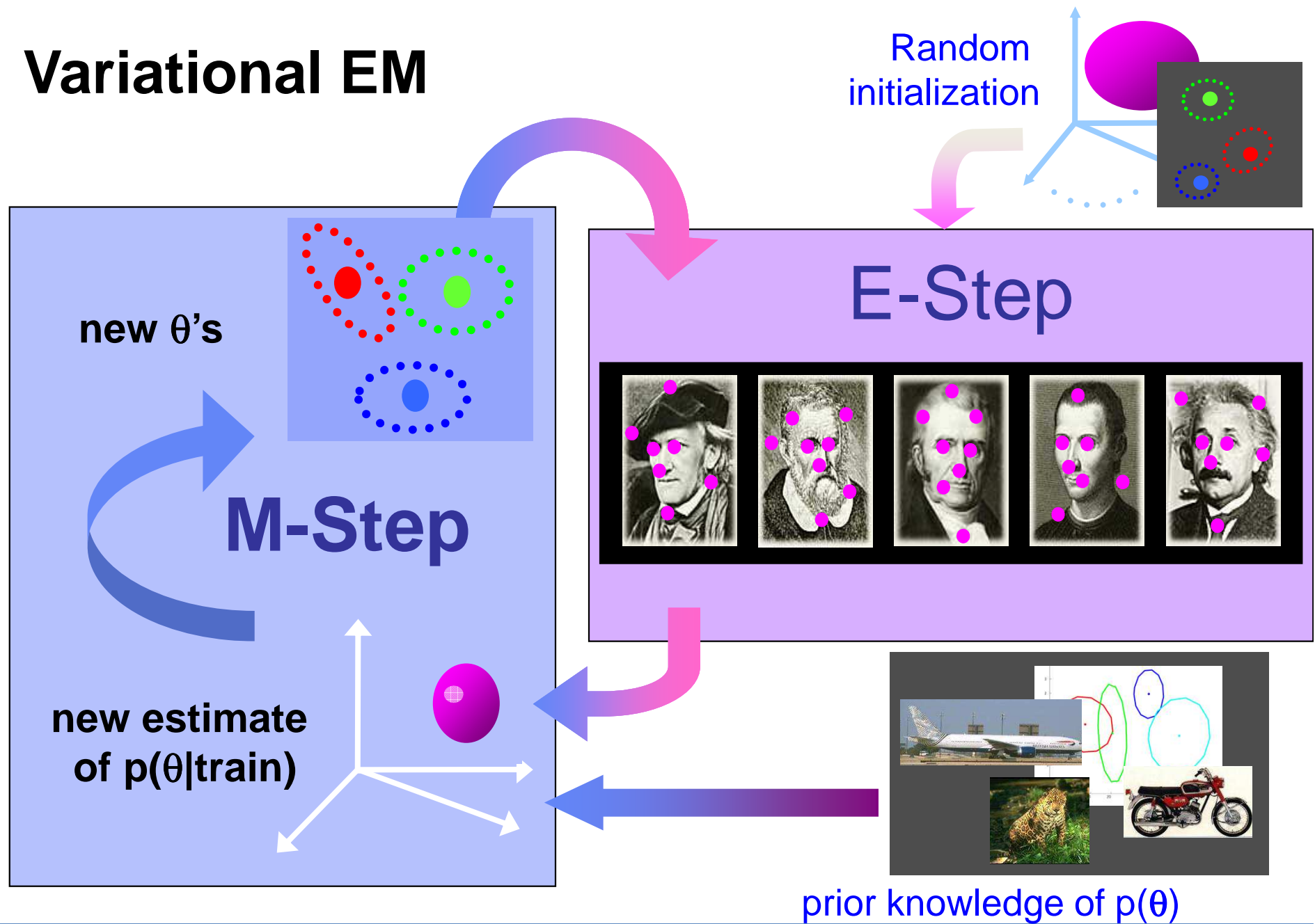


Learning Model Distribution

$$p(\theta | \text{object, train}) \propto p(\text{train} | \theta, \text{object}) p(\theta)$$

- use **Prior** information
- Bayesian learning
 - marginalize over theta
- ❖ **Variational EM** (Attias, Hinton, Minka, etc.)

Variational EM



Experiments

Training:

1- 6 randomly
drawn images

Testing:

50 fg/ 50 bg images
object present/absent

Datasets



faces



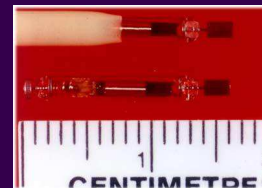
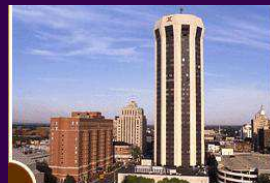
airplanes



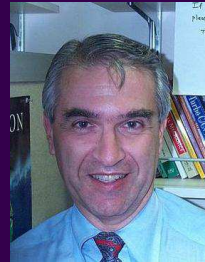
spotted cats



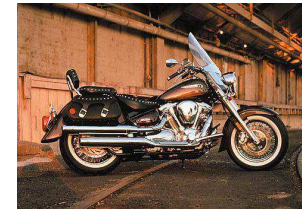
motorbikes



Faces



Motorbikes



Airplanes



Spotted cats



66

28-Nov-11

Experiments: obtaining priors



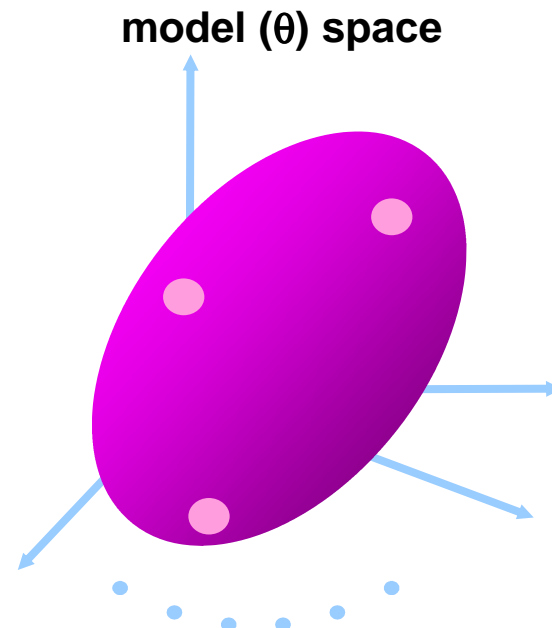
airplanes



spotted cats



motorbikes



faces

Experiments: obtaining priors



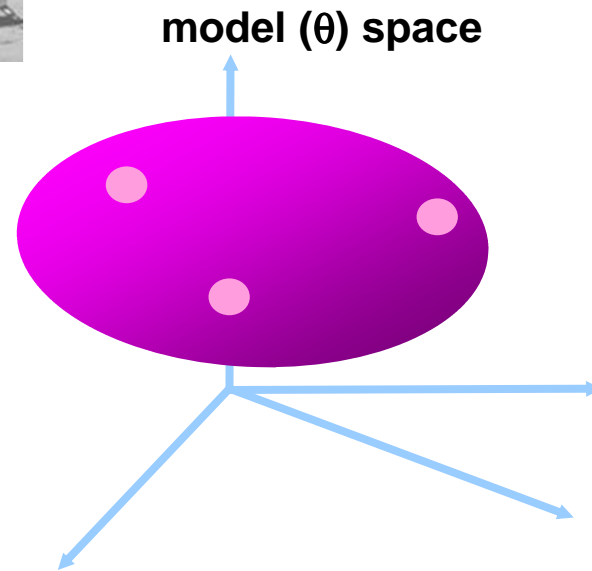
airplanes



faces

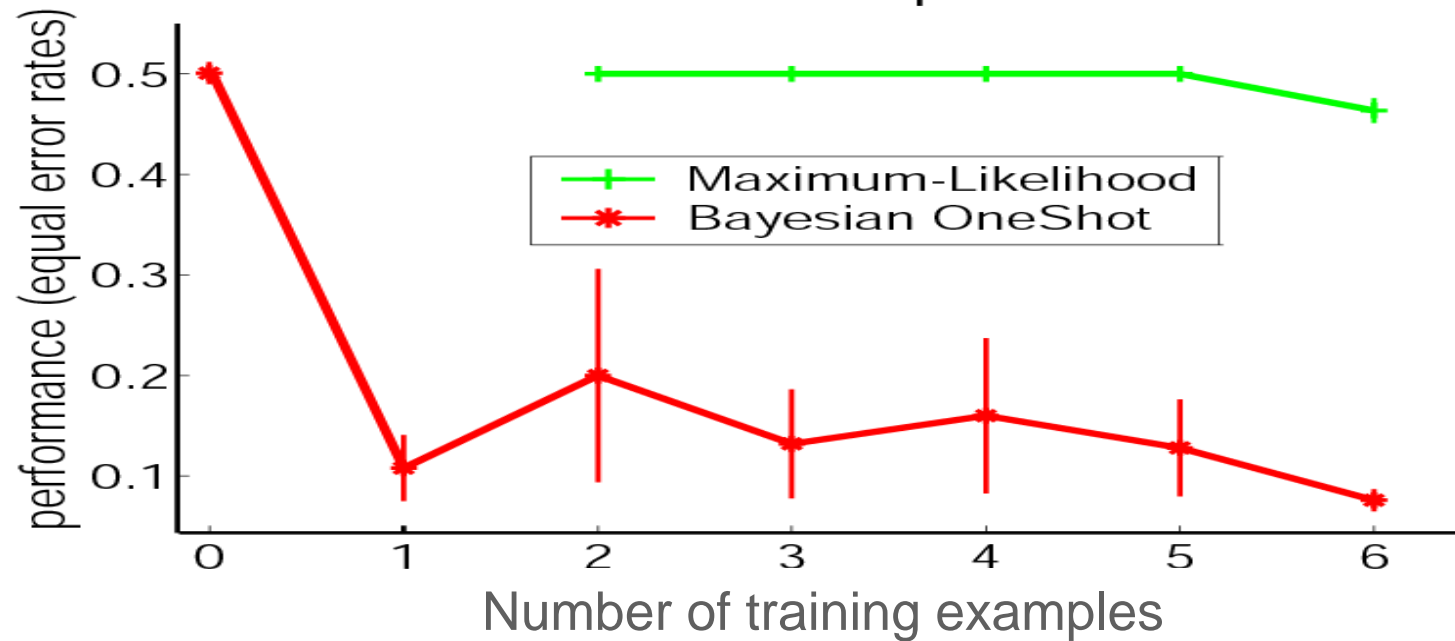


motorbikes



spotted cats

Performance comparison



Part 1



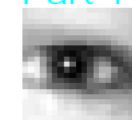
Part 2



Part 3



Part 4



Correct



Correct



Correct



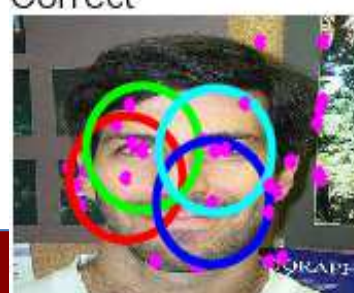
Correct



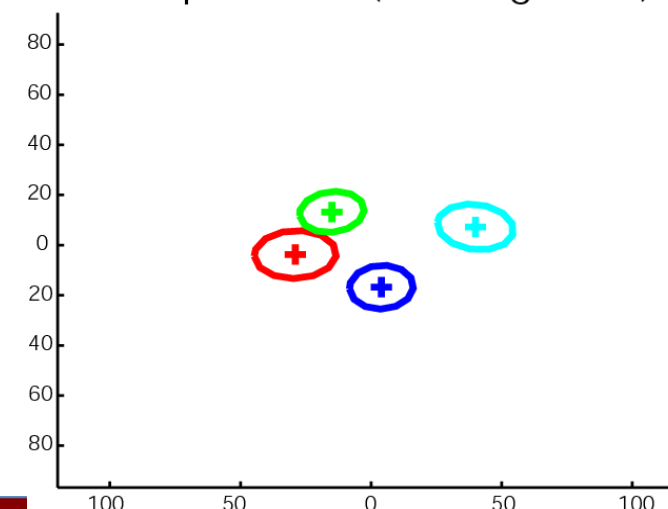
INCORRECT



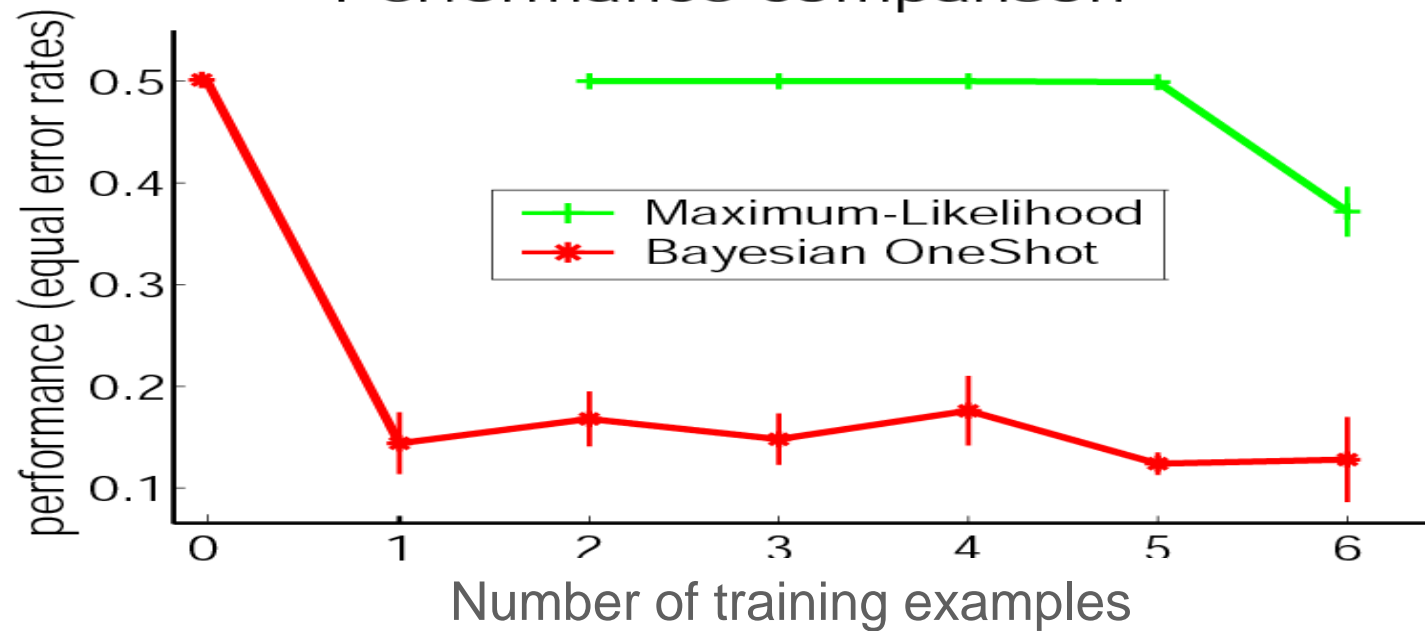
Correct



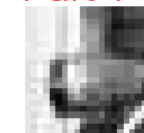
Shape Model (Training # = 1)



Performance comparison



Part 1



Part 2



Part 3



Part 4



Correct



INCORRECT



Correct



INCORRECT



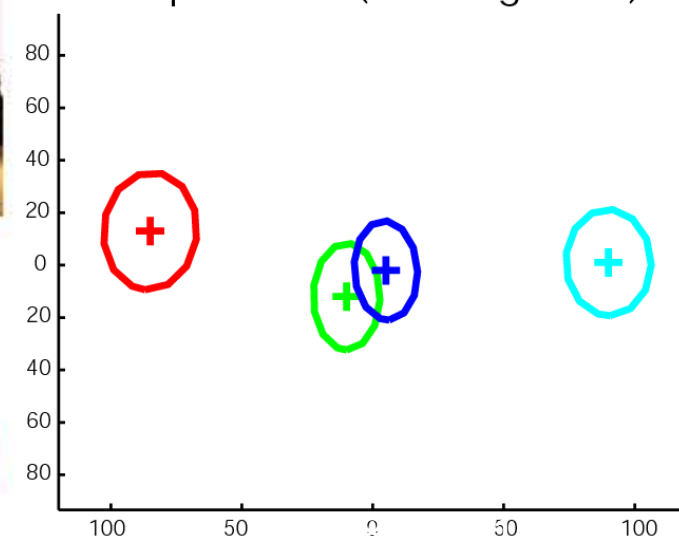
Correct



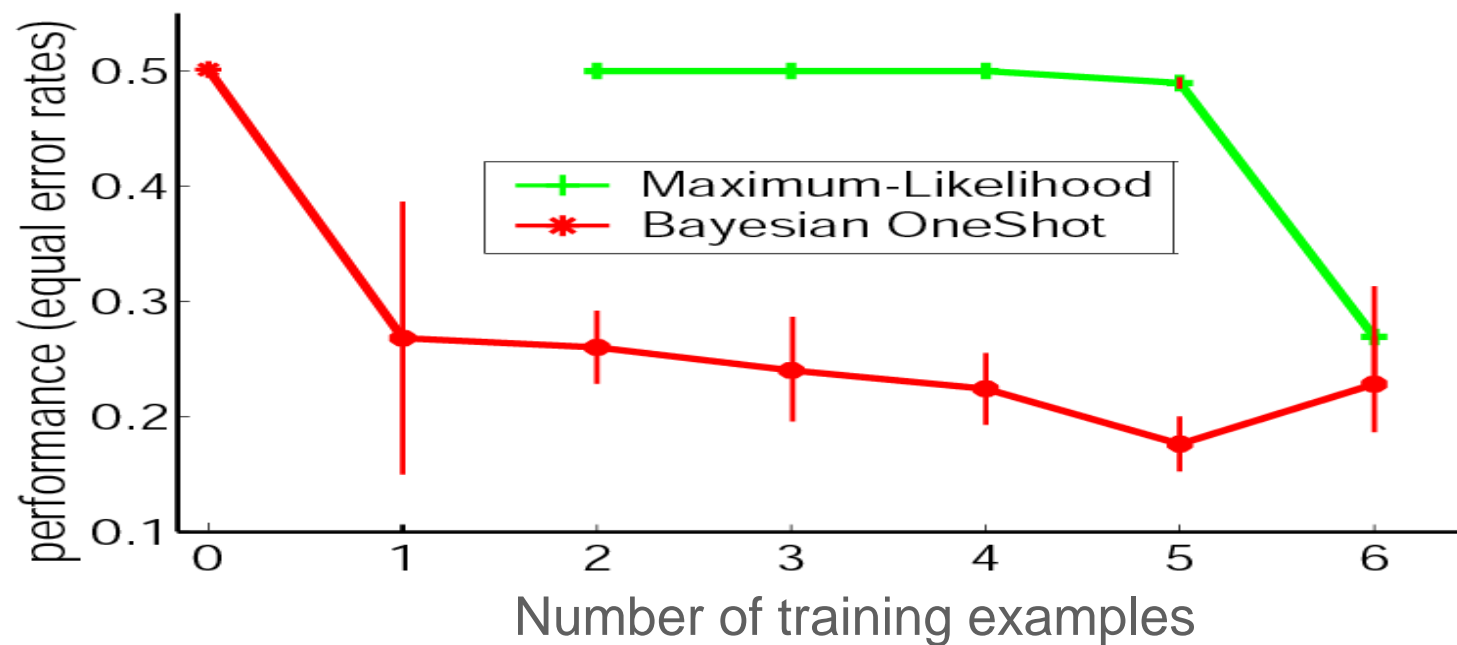
Correct



Shape Model (Training # = 1)



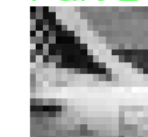
Performance comparison



Part 1



Part 2



Part 3



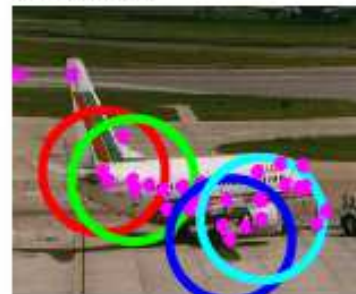
Part 4



Correct



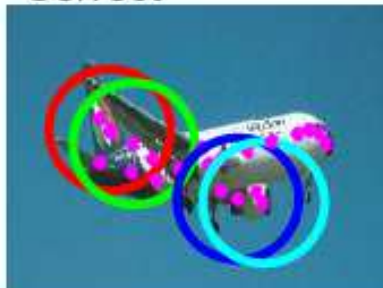
Correct



INCORRECT



Correct



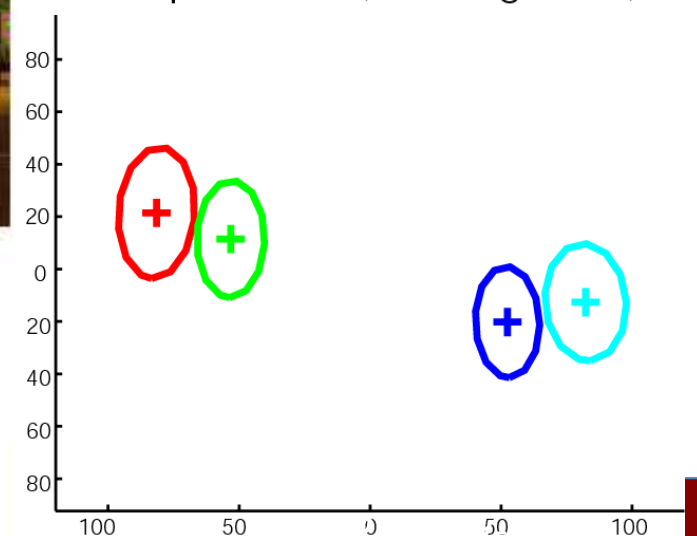
Correct



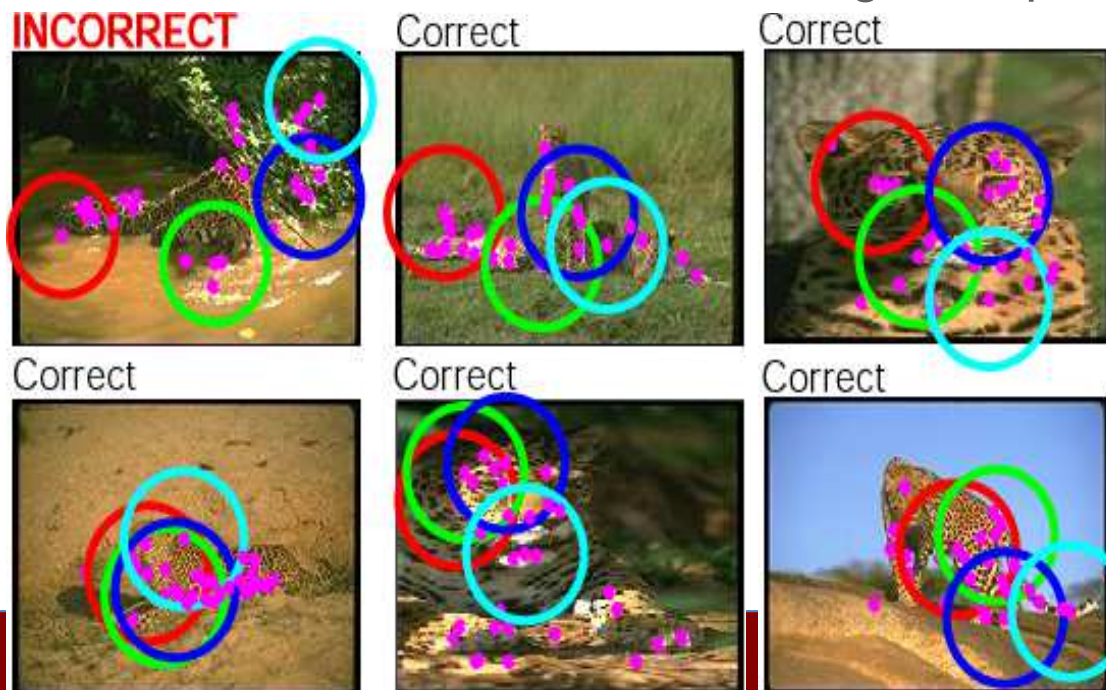
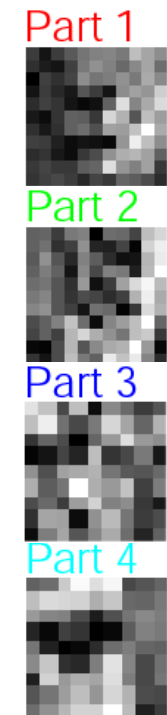
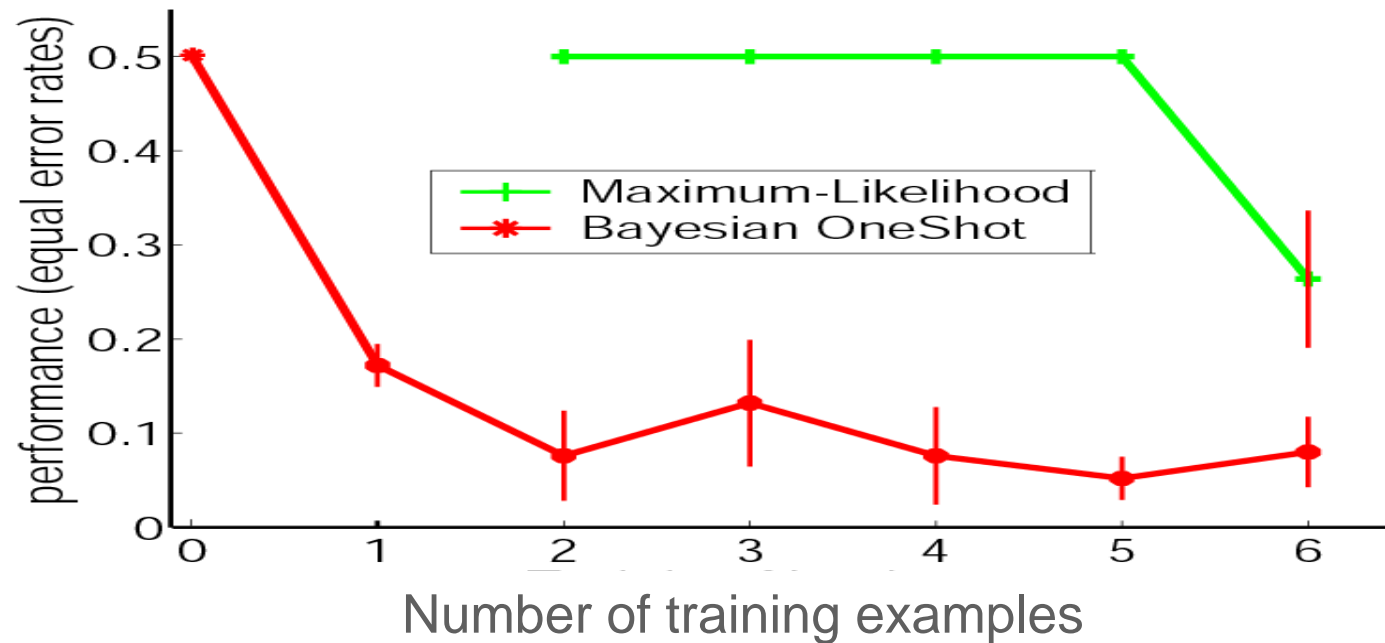
Correct



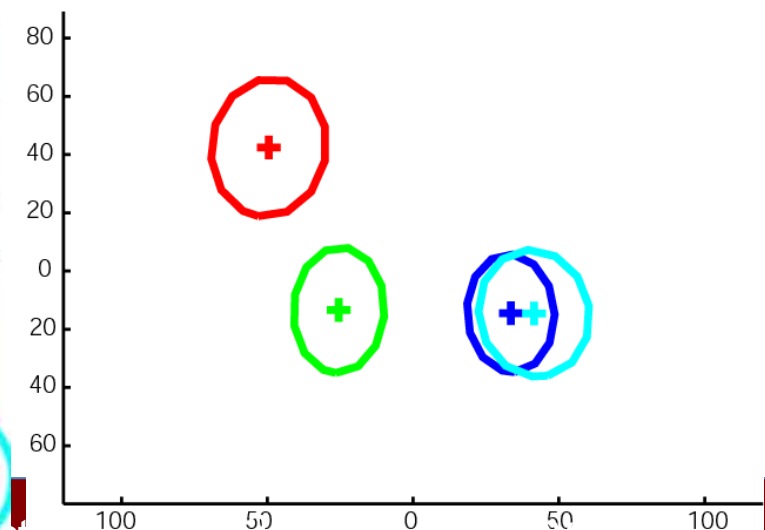
Shape Model (Training # = 1)



Performance comparison



Shape Model (Training # = 1)



Algorithm	Training Examples	Categories	Results(error)
Burl, et al. Weber, et al. Fergus, et al.	200 ~ 400	Faces, Motorbikes, Spotted cats, Airplanes, Cars	5.6 - 10 %
Viola et al.	~10,000	Faces	7-21%
Schneiderman, et al.	~2,000	Faces, Cars	5.6 – 17%
Rowley et al.	~500	Faces	7.5 – 24.1%
Bayesian One-Shot	1 ~ 5	Faces, Motorbikes, Spotted cats, Airplanes	8 – 15 %

What we have learned today?

- Introduction
- Constellation model
 - Weakly supervised training
 - One-shot learning
- (Problem Set 4 (Q1))