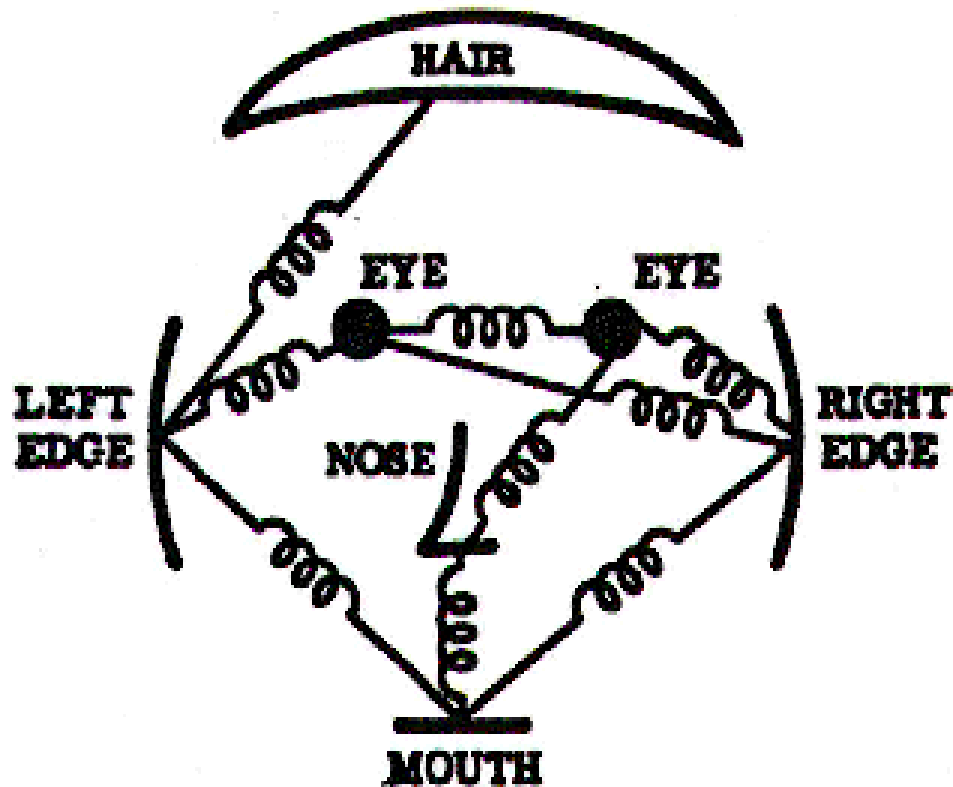


Beyond bags of features: Adding spatial information

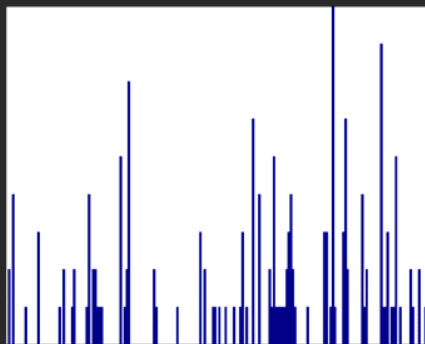


Adding spatial information

- Forming vocabularies from pairs of nearby features – “doublets” or “bigrams”
- Computing bags of features on sub-windows of the whole image
- Using codebooks to vote for object position
- Generative part-based models

Spatial pyramid representation

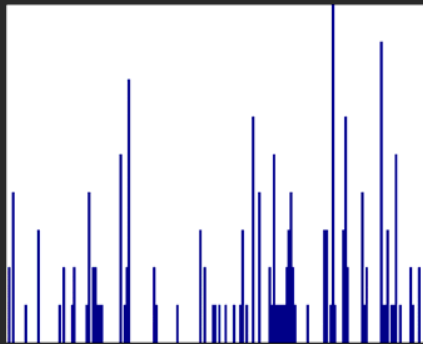
- Extension of a bag of features
- Locally orderless representation at several levels of resolution



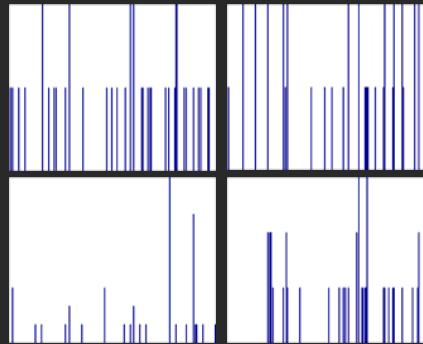
level 0

Spatial pyramid representation

- Extension of a bag of features
- Locally orderless representation at several levels of resolution



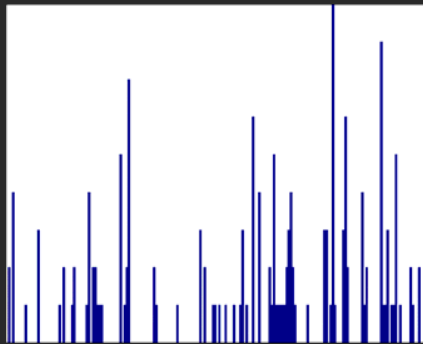
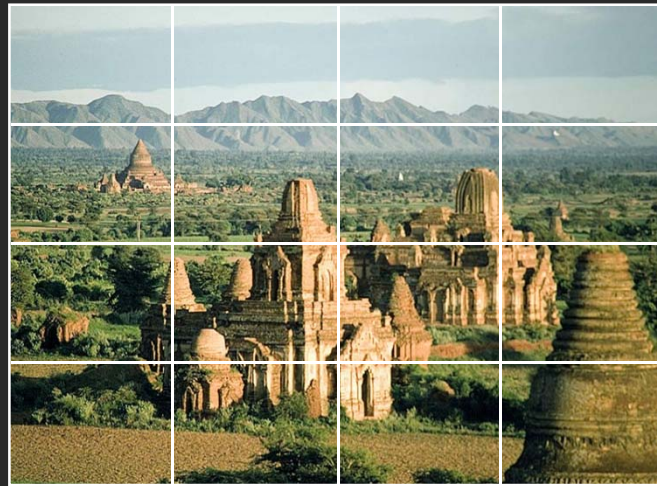
level 0



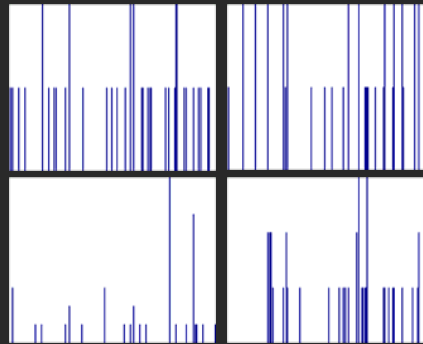
level 1

Spatial pyramid representation

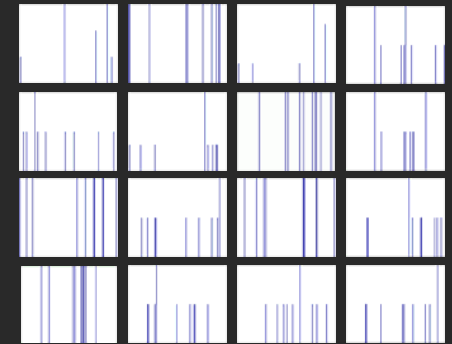
- Extension of a bag of features
- Locally orderless representation at several levels of resolution



level 0

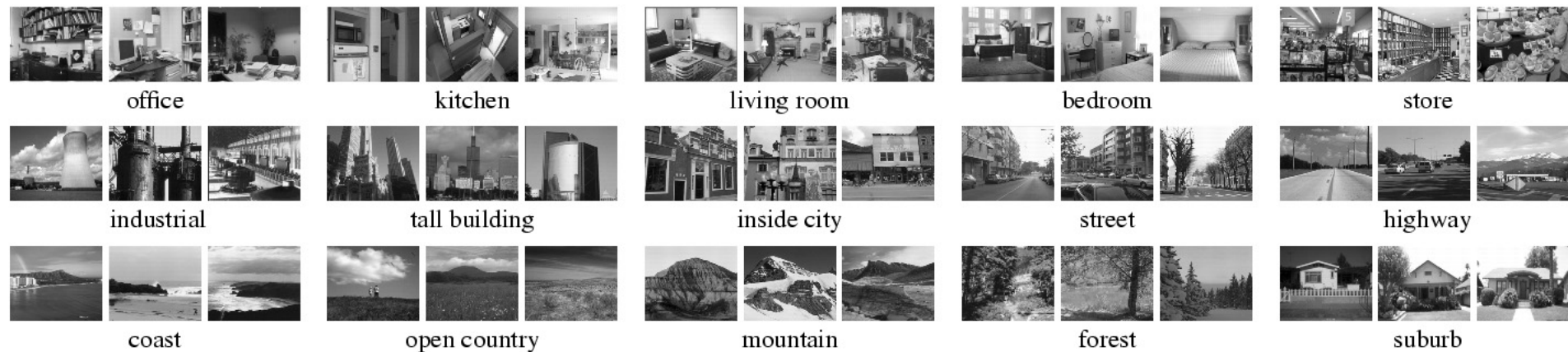


level 1



level 2

Scene category dataset



Multi-class classification results (100 training images per class)

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	81.1 ±0.3
3 (8 × 8)	63.3 ±0.8	66.8 ±0.6	77.2 ±0.4	80.7 ±0.3

Caltech101 dataset

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

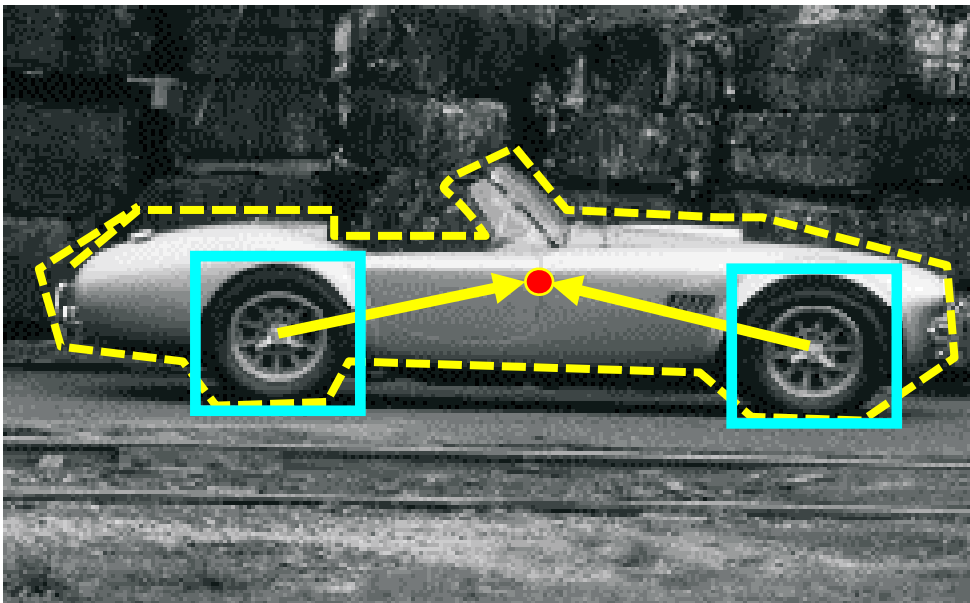


Multi-class classification results (30 training images per class)

	Weak features (16)		Strong features (200)	
Level	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	64.6 \pm 0.8
3	52.2 \pm 0.8	54.0 \pm 1.1	60.3 \pm 0.9	64.6 \pm 0.7

Implicit shape models

- Visual codebook is used to index votes for object position



visual codeword with displacement vectors

training image annotated with object localization info

B. Leibe, A. Leonardis, and B. Schiele, [Combined Object Categorization and Segmentation with an Implicit Shape Model](#), ECCV Workshop on Statistical Learning in Computer Vision 2004

Implicit shape models

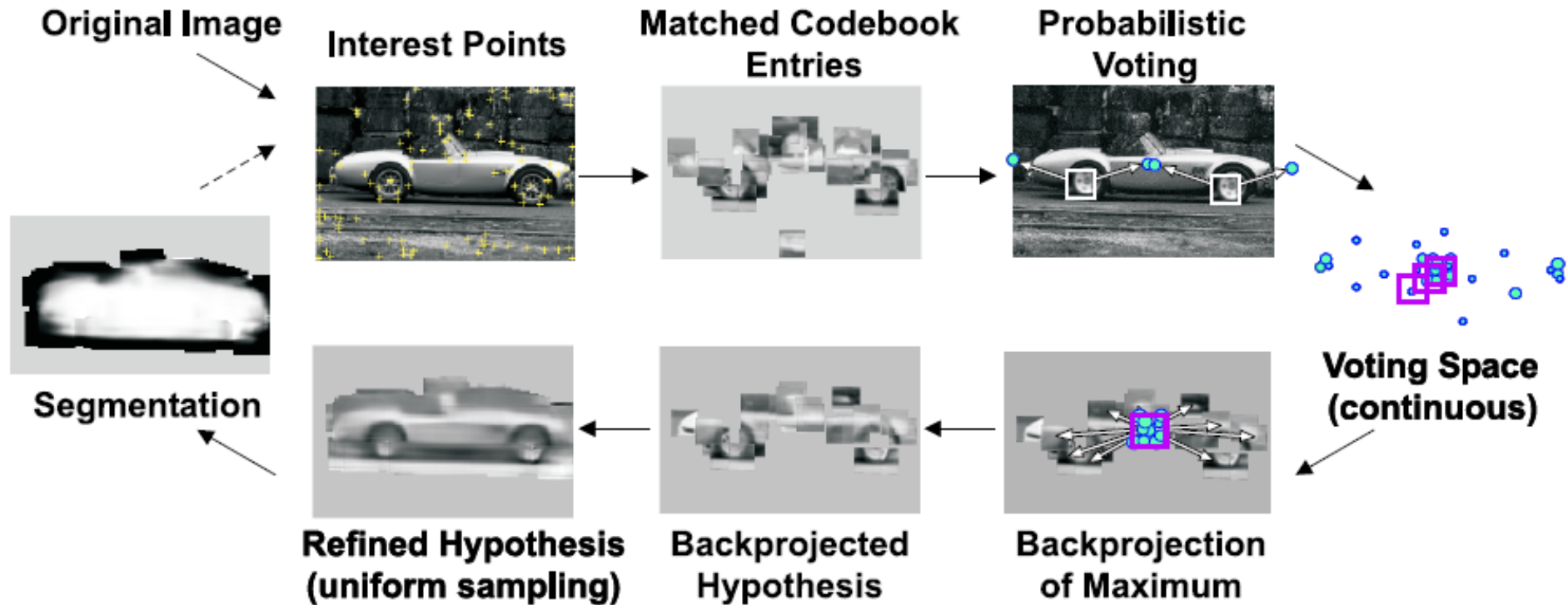
- Visual codebook is used to index votes for object position



test image

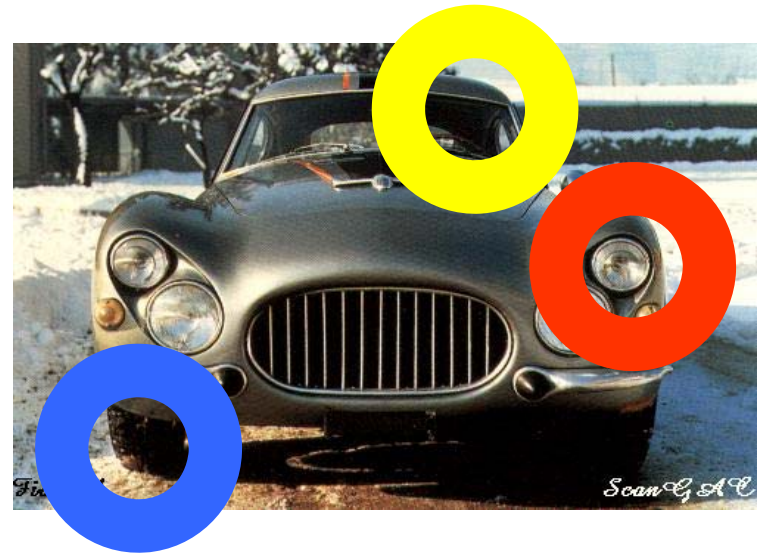
B. Leibe, A. Leonardis, and B. Schiele, [Combined Object Categorization and Segmentation with an Implicit Shape Model](#), ECCV Workshop on Statistical Learning in Computer Vision 2004

Implicit shape models: Details



B. Leibe, A. Leonardis, and B. Schiele, [Combined Object Categorization and Segmentation with an Implicit Shape Model](#), ECCV Workshop on Statistical Learning in Computer Vision 2004

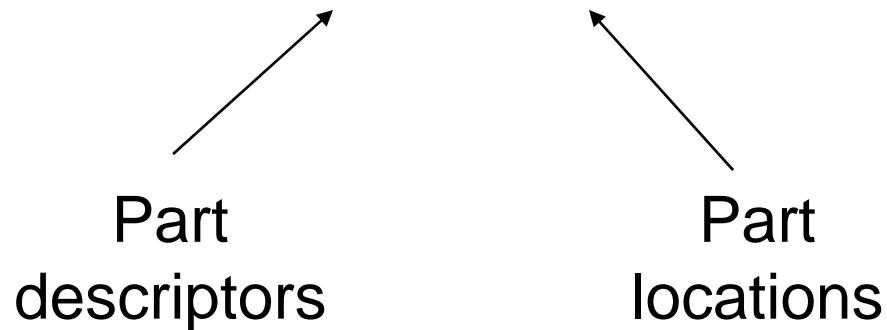
Generative part-based models



R. Fergus, P. Perona and A. Zisserman, [Object Class Recognition by Unsupervised Scale-Invariant Learning](#), CVPR 2003

Probabilistic model

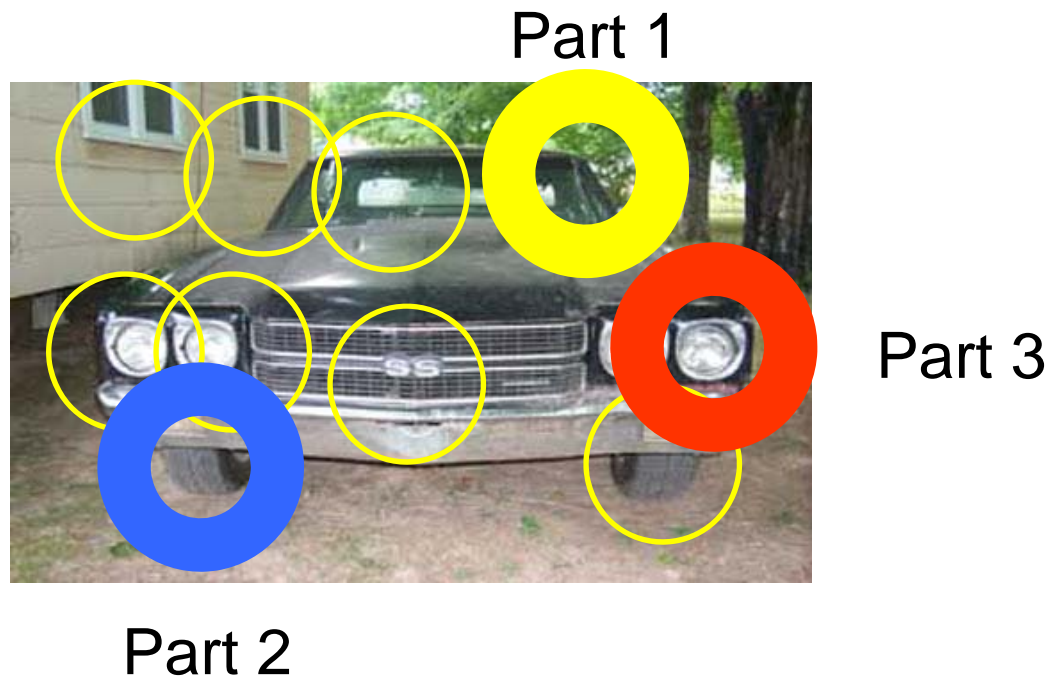
$$P(\text{image} | \text{object}) = P(\text{appearance, shape} | \text{object})$$



Candidate parts

Probabilistic model

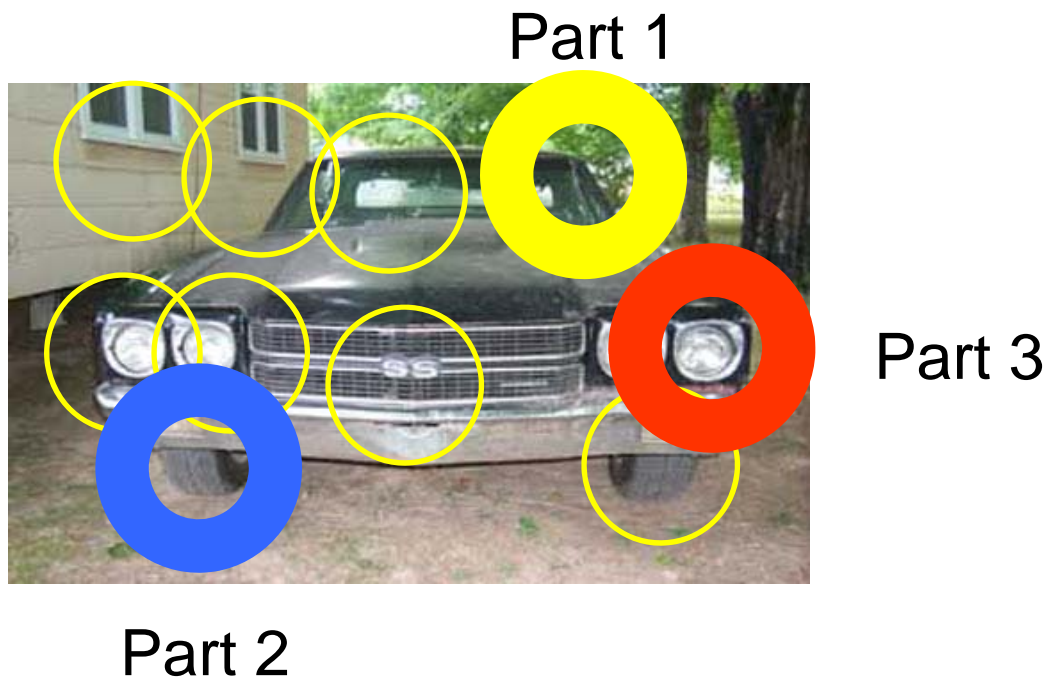
$$P(\text{image} | \text{object}) = P(\text{appearance}, \text{shape} | \text{object})$$



Probabilistic model

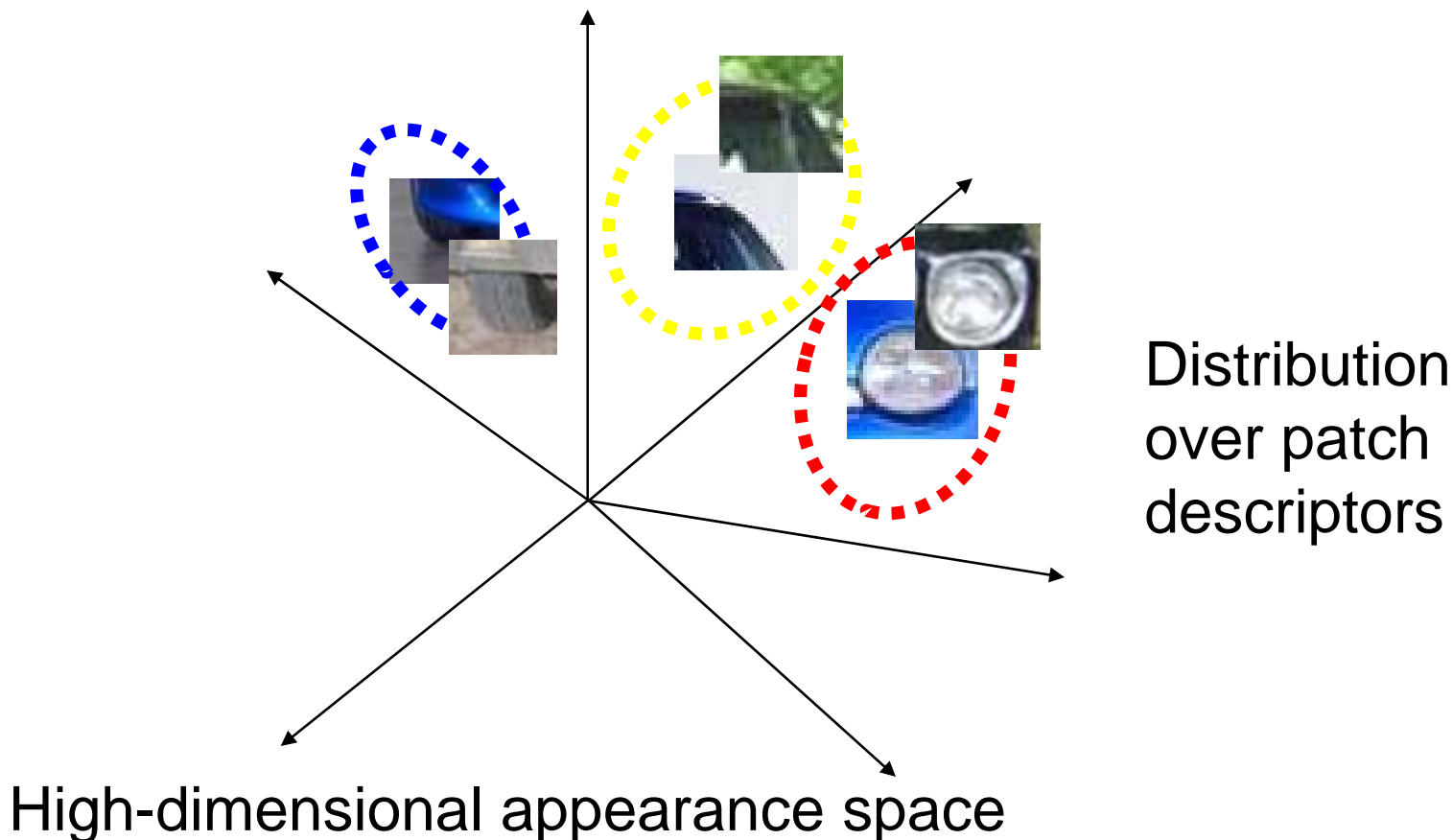
$$P(\text{image} | \text{object}) = P(\text{appearance}, \text{shape} | \text{object})$$
$$= \max_h P(\text{appearance} | h, \text{object}) p(\text{shape} | h, \text{object}) p(h | \text{object})$$

h : assignment of features to parts



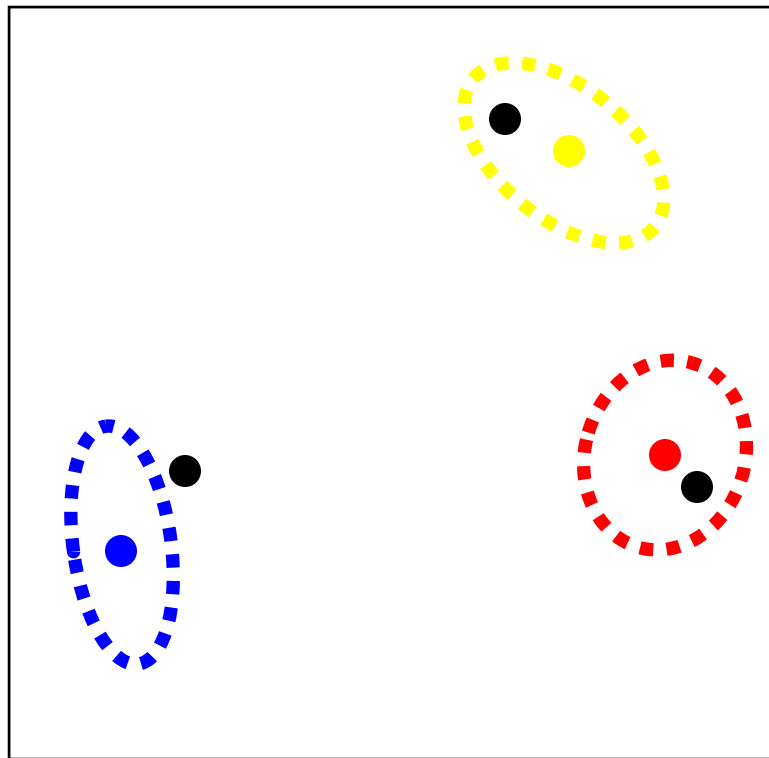
Probabilistic model

$$P(\text{image} \mid \text{object}) = P(\text{appearance}, \text{shape} \mid \text{object})$$
$$= \max_h P(\text{appearance} \mid h, \text{object}) p(\text{shape} \mid h, \text{object}) p(h \mid \text{object})$$



Probabilistic model

$$P(\text{image} \mid \text{object}) = P(\text{appearance}, \text{shape} \mid \text{object})$$
$$= \max_h P(\text{appearance} \mid h, \text{object}) p(\text{shape} \mid h, \text{object}) p(h \mid \text{object})$$

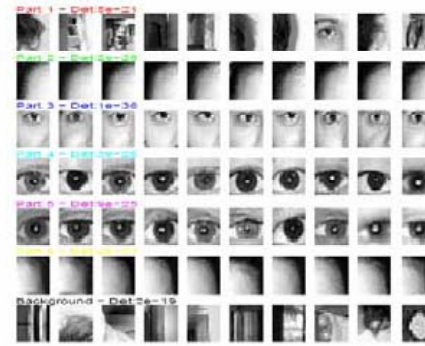
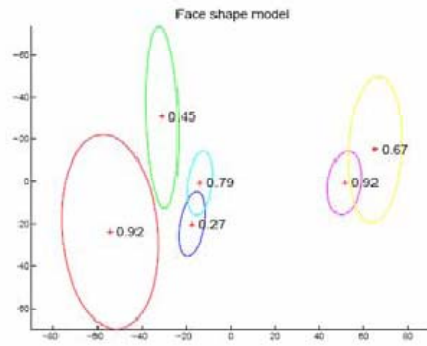


Distribution
over joint
part positions

2D image space

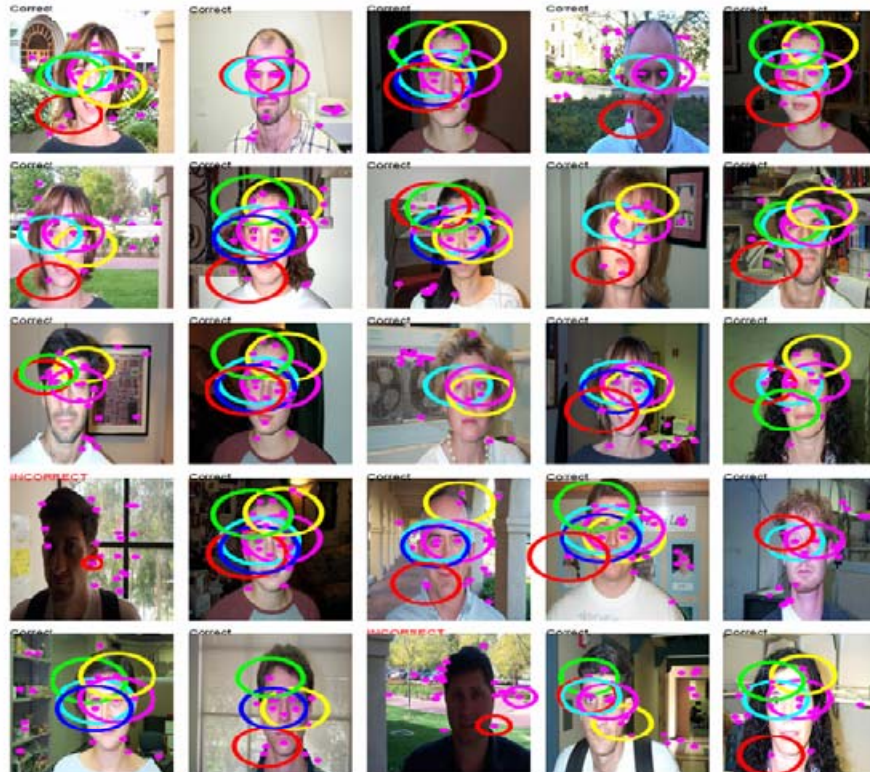
Results: Faces

Face
shape
model

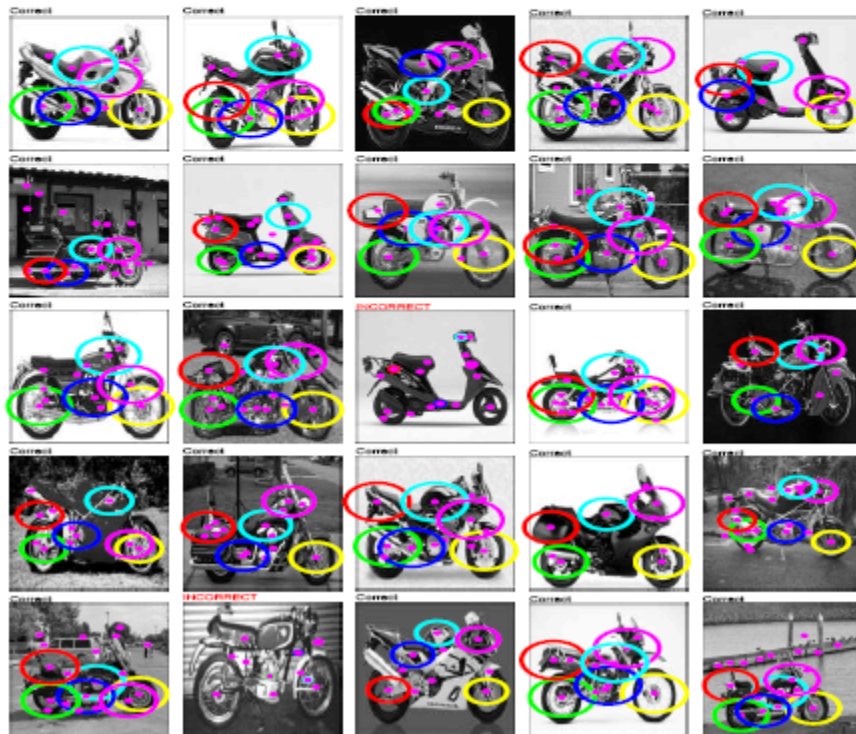
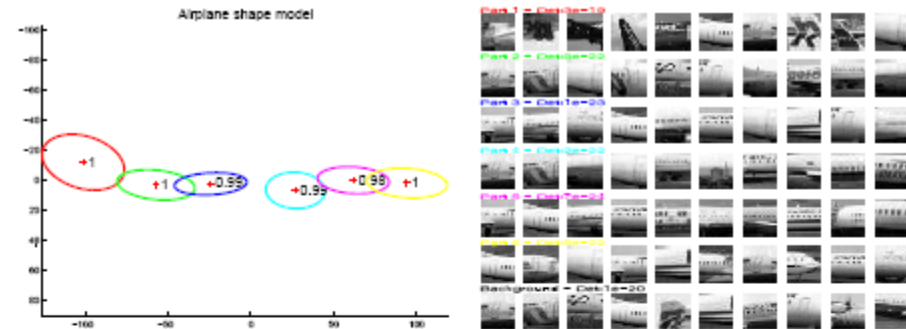
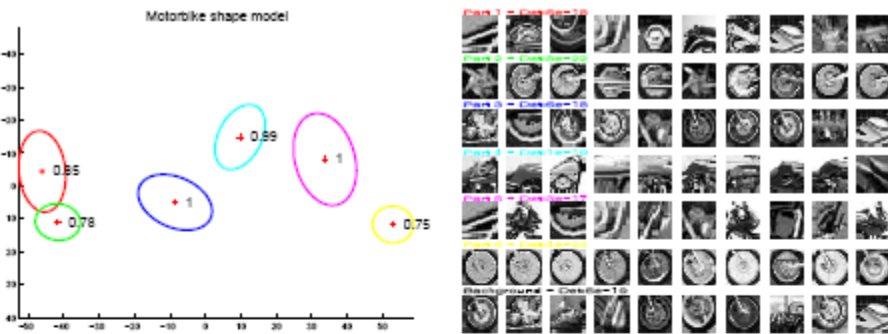


Patch
appearance
model

Recognition
results



Results: Motorbikes and airplanes



Summary: Adding spatial information

- **Doublet vocabularies**
 - Pro: takes co-occurrences into account, some geometric invariance is preserved
 - Con: too many doublet probabilities to estimate
- **Spatial pyramids**
 - Pro: simple extension of a bag of features, works very well
 - Con: no geometric invariance, no object localization
- **Implicit shape models**
 - Pro: can localize object, maintain translation and possibly scale invariance
 - Con: need supervised training data (known object positions and possibly segmentation masks)
- **Generative part-based models**
 - Pro: very nice conceptually, can be learned from unsegmented images
 - Con: combinatorial hypothesis search problem