

# Image Classification and Object Recognition

---

Selim Aksoy

Department of Computer Engineering

Bilkent University

saksoy@cs.bilkent.edu.tr

# Image classification

---

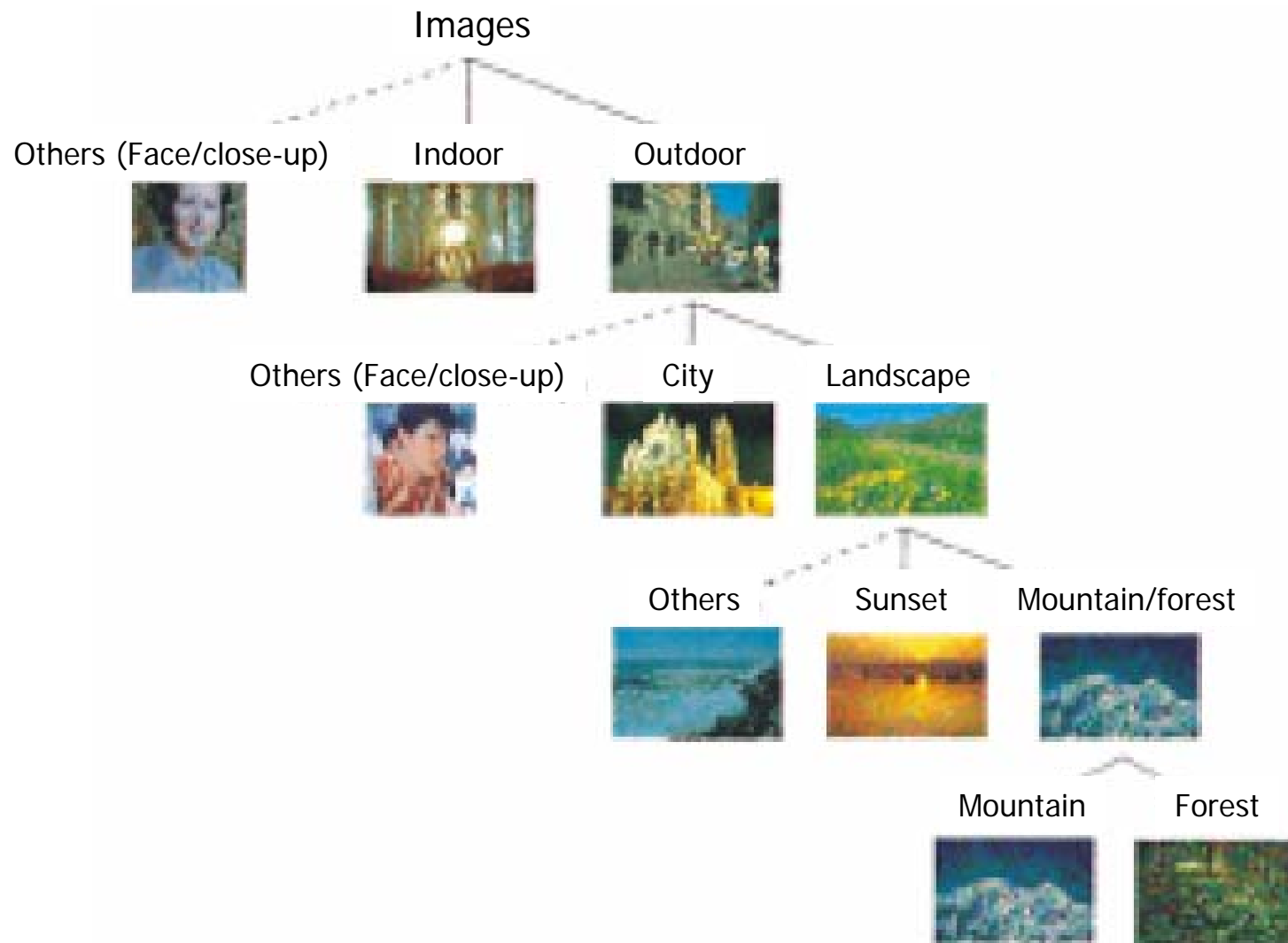
- Image (scene) classification is a fundamental problem in image understanding.
- Automatic techniques for associating scenes with semantic labels have a high potential for improving the performance of other computer vision applications such as
  - browsing (natural grouping of images instead of clusters based only on low-level features),
  - retrieval (filtering images in archives based on content), and
  - object recognition (the probability of an unknown object/region that exhibits several local features of a ship actually being a ship can be increased if the scene context is known to be a coast with high confidence but can be decreased if no water related context is dominant in that scene).

# Image classification

---

- The image classification problem has two critical components: **representing** images and **learning** models for semantic categories using these representations.
- Early work used low-level global features extracted from the whole image or from a fixed spatial layout.
- More recent approaches exploit local statistics in images using patches extracted by interest point detectors.
- Other configurations that use regions and their spatial relationships are also proposed.

# Hierarchical image classification



Hierarchy of 11 scene categories (Vailaya et al., "Image classification for content-based indexing," IEEE Trans. Image Processing, 2001).



# Hierarchical image classification

---

- Image representation:
  - Mean and std. dev. of LUV values in 10x10 blocks for indoor/outdoor classification.
  - Edge direction histograms for city/landscape classification.
  - Histograms of HSV and LUV values for sunset/mountain/forest classification.
- Classification:
  - Class-conditional density estimation using vector quantization.
  - Bayesian classification.

# Hierarchical image classification

TABLE III  
ACCURACIES (IN PERCENT) FOR INDOOR/OUTDOOR CLASSIFICATION USING  
COLOR MOMENTS; TEST SET 1 AND TEST SET 2 ARE INDEPENDENT TEST SETS

Test Data	Database Size	Accuracy (%)
Training Set	2,541	94.2
Test Set 1	2,540	88.2
Test Set 2	1,850	88.7
Entire Database	6,931	90.5

TABLE IV  
CLASSIFICATION ACCURACIES (IN PERCENT) FOR CITY/LANDSCAPE CLASSIFICATION; THE FEATURES ARE ABBREVIATED AS FOLLOWS: EDGE DIRECTION  
HISTOGRAM (EDH), EDGE DIRECTION COHERENCE VECTOR (EDCV), COLOR HISTOGRAM (CH), AND COLOR COHERENCE VECTOR (CCV)

Test Data	EDH	EDCV	CH	CCV	EDH & CH	EDH & CCV	EDCV & CH	EDCV & CCV
Training Set	94.7	97.0	83.7	83.5	94.8	95.4	96.4	96.9
Test Set	92.0	92.9	75.4	76.0	92.5	92.8	93.4	93.8
Entire Database	93.4	95.0	79.6	79.8	93.7	94.1	94.9	95.3

# Hierarchical image classification

TABLE V

CLASSIFICATION ACCURACIES (IN PERCENT) FOR SUNSET/FOREST/MOUNTAIN CLASSIFICATION; *SPM* STANDS FOR “SPATIAL COLOR MOMENTS”

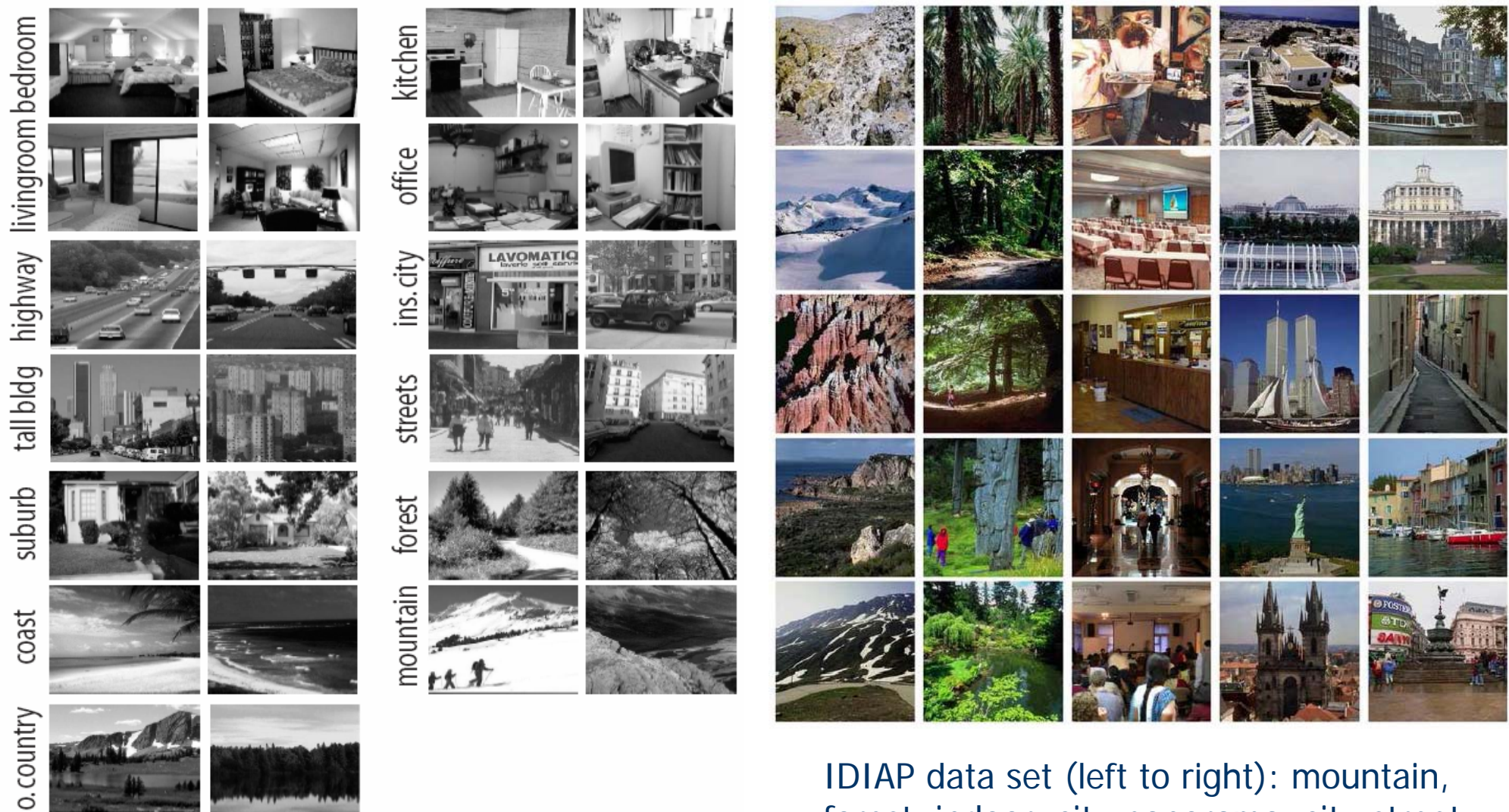
Test Data	EDH	EDCV	CH	CCV	SPM	EDH & CH	EDH & CCV	EDCV & CH	EDCV & CCV
Training Set	88.3	88.3	96.2	99.2	98.9	95.9	96.6	95.5	97.0
Test Set	86.3	89.0	89.7	93.9	93.9	90.1	95.4	90.5	95.1
Entire Database	87.4	88.7	93.0	96.6	96.4	93.0	96.0	93.0	96.1

TABLE VI

CLASSIFICATION ACCURACIES (IN PERCENT) FOR FOREST/MOUNTAIN CLASSIFICATION

Test Data	EDH	EDCV	CH	CCV	SPM	EDH & CH	EDH & CCV	EDCV & CH	EDCV & CCV
Training Set	83.4	78.1	92.0	98.9	98.4	94.1	98.4	93.6	98.4
Test Set	87.1	77.2	91.4	91.9	93.6	93.0	92.5	93.5	91.9
Entire Database	85.3	77.7	91.7	95.5	96.0	93.6	95.5	93.6	95.2

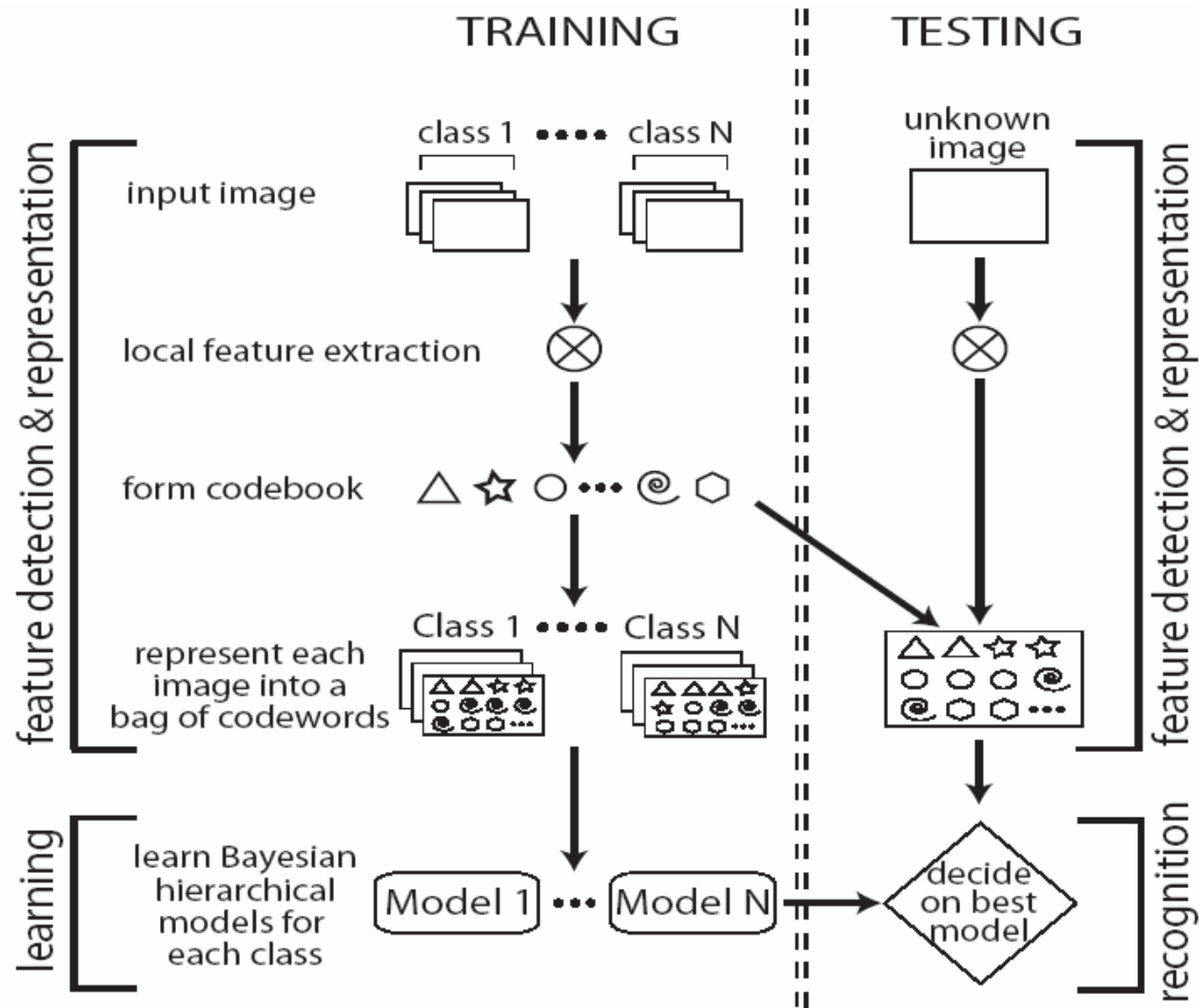
# Image classification using bag-of-words



IDIAP data set (left to right): mountain, forest, indoor, city-panorama, city-street.

Caltech data set: 13 natural scene categories.

# Image classification using bag-of-words

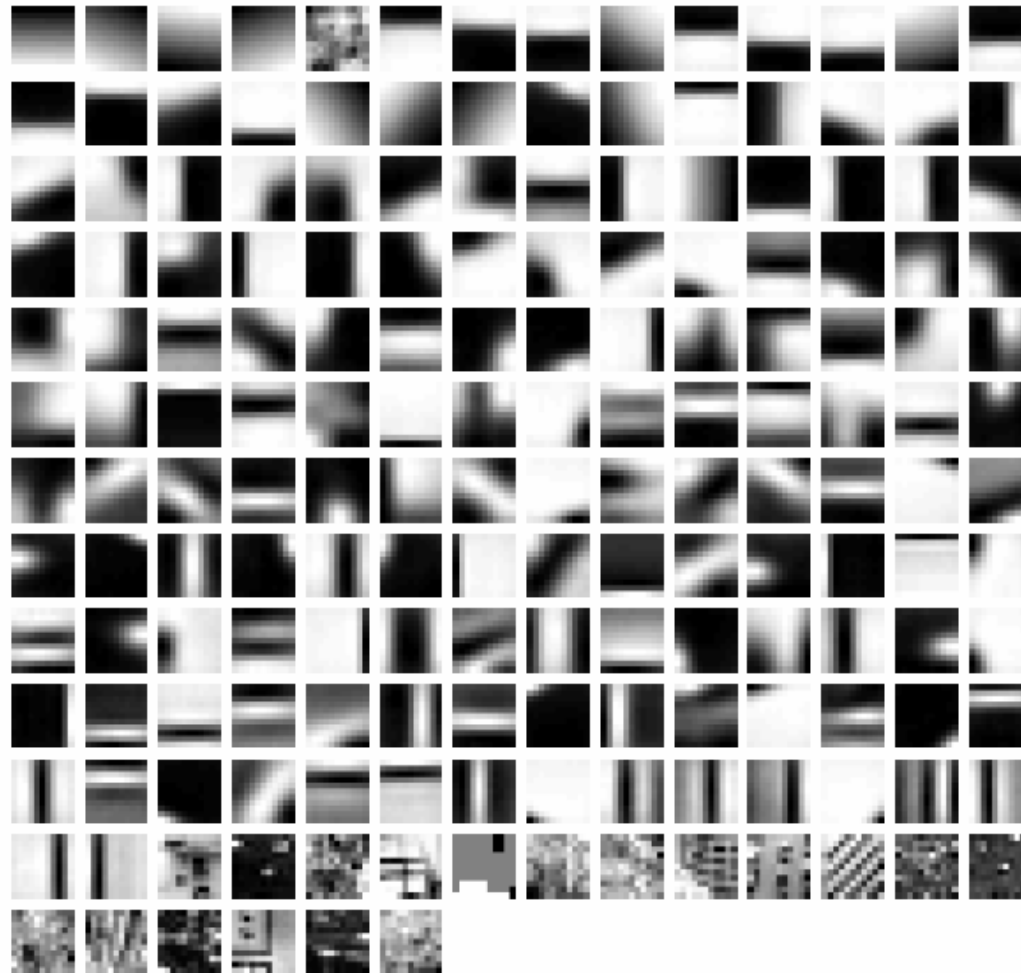


Flowchart from Fei-Fei Li, Pietro Perona, "A Bayesian hierarchical model for learning natural scene categories," IEEE CVPR, 2005.

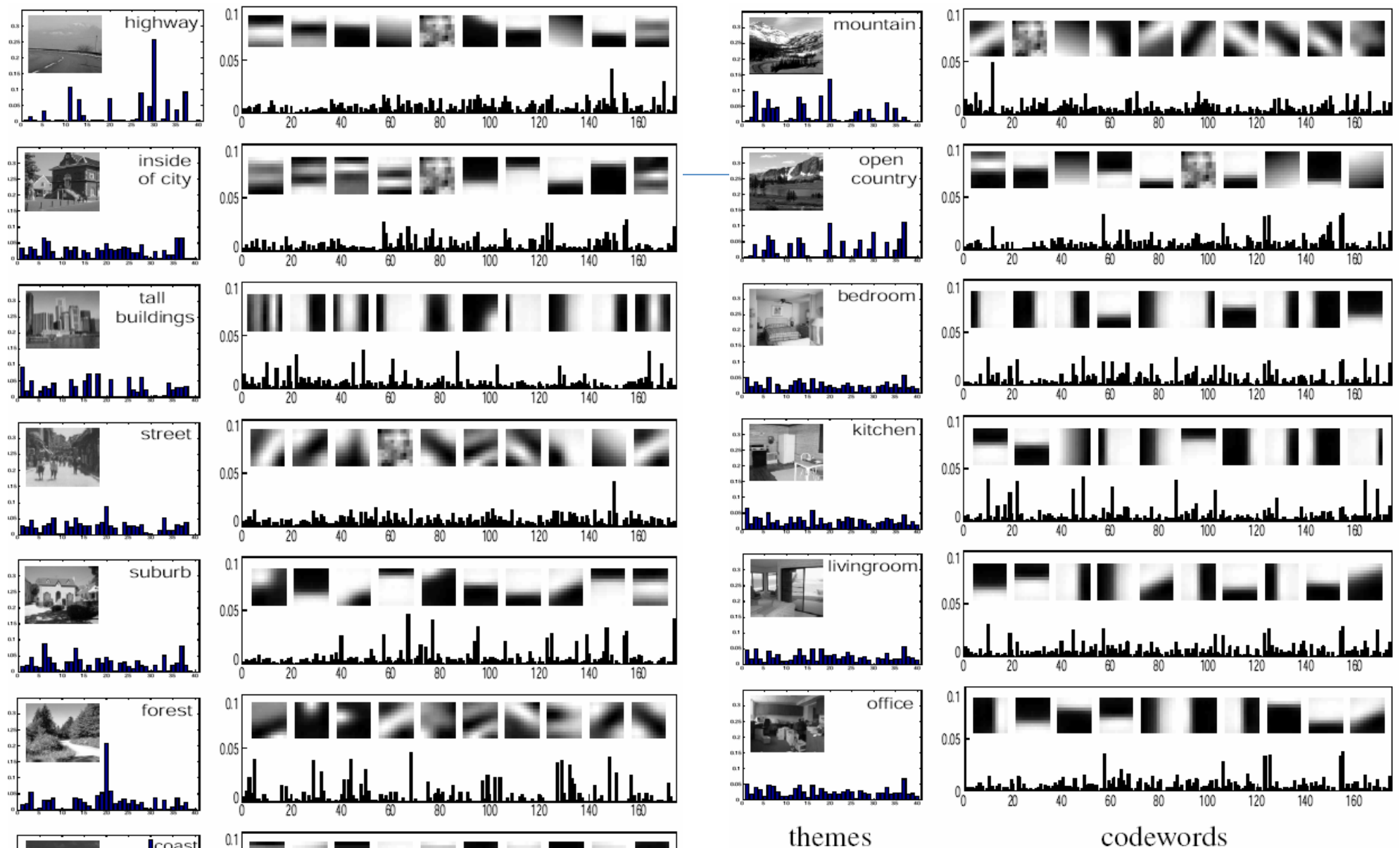


# Image classification using bag-of-words

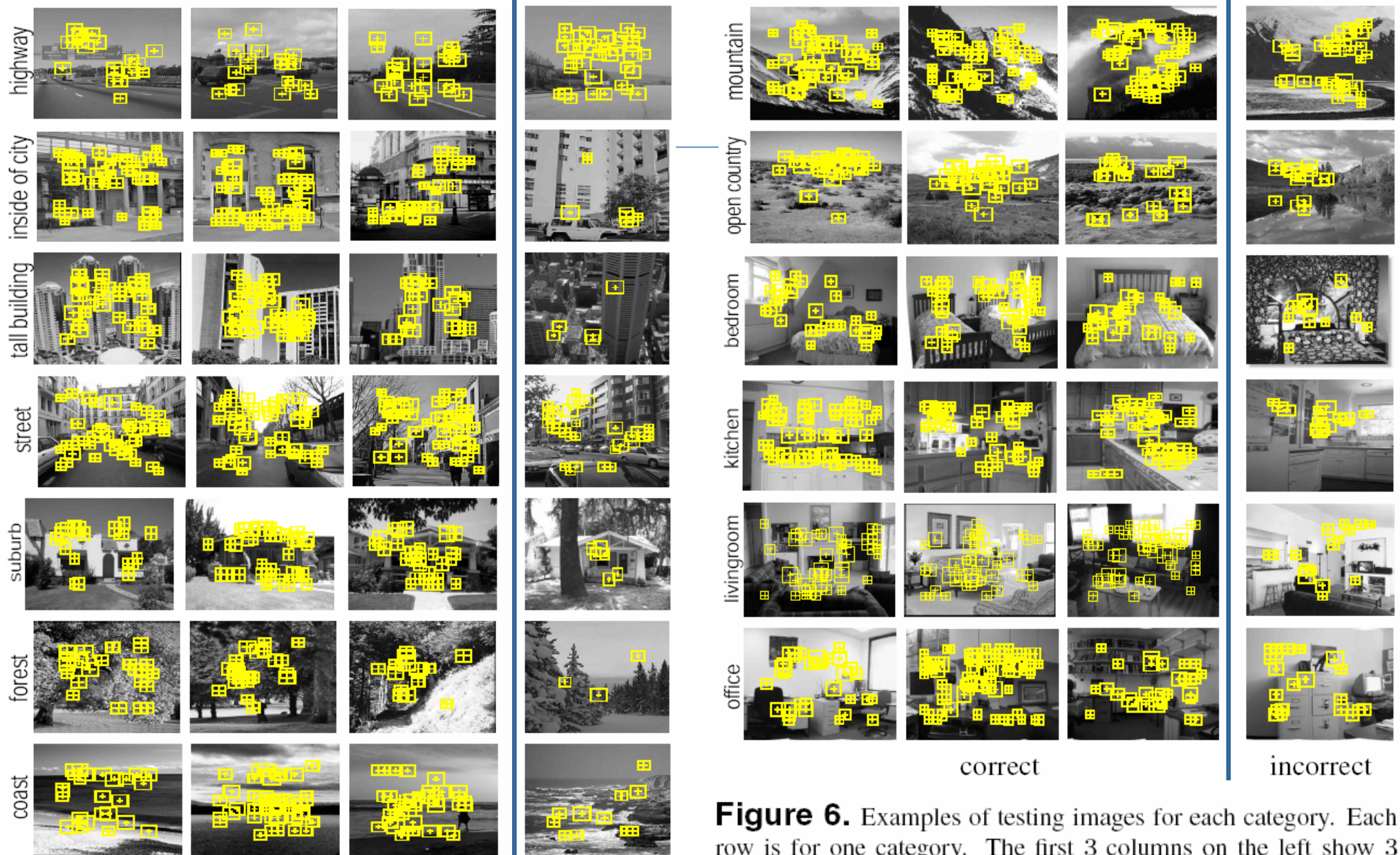
---



A codebook obtained from 650 training examples from 13 categories.  
Image patches are detected by a sliding grid and random sampling of scales.



**Figure 5.** Internal structure of the models learnt for each category. Each row represents one category. The left panel shows the distribution of the 40 intermediate themes. The right panel shows the distribution of codewords as well as the appearance of 10 codewords selected from the top 20 most likely codewords for this category model.



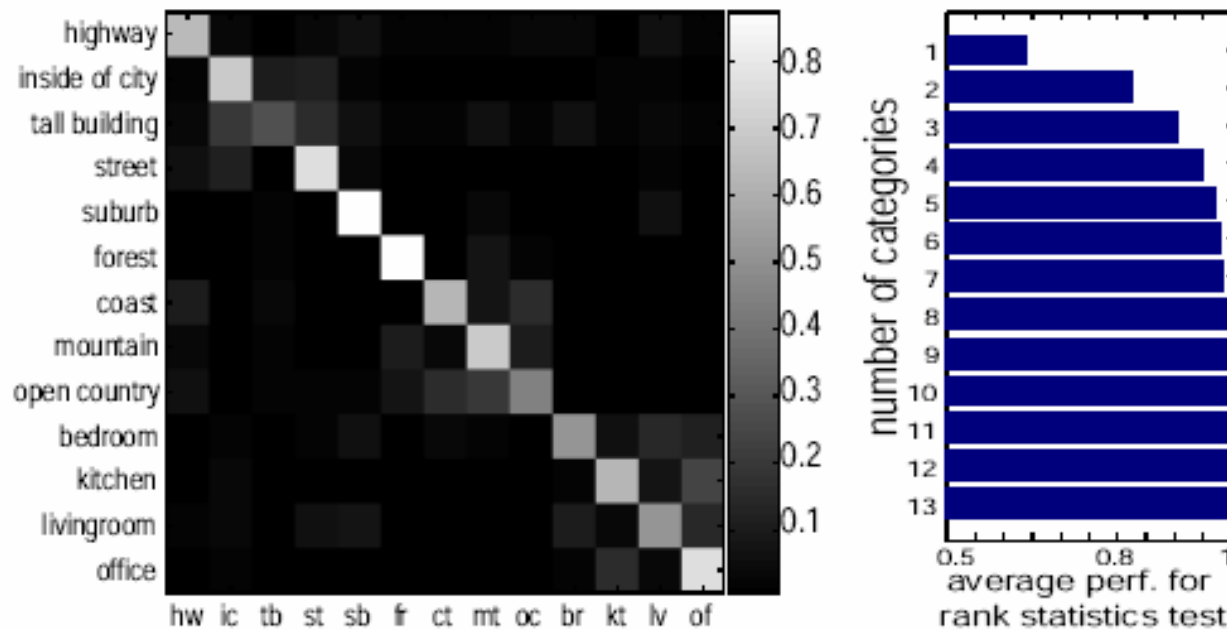
correct

incorrect

**Figure 6.** Examples of testing images for each category. Each row is for one category. The first 3 columns on the left show 3 examples of correctly recognized images, the last column on the right shows an example of incorrectly recognized image. Super-imposed on each image, we show samples of patches that belong to the most significant set of codewords given the category model.

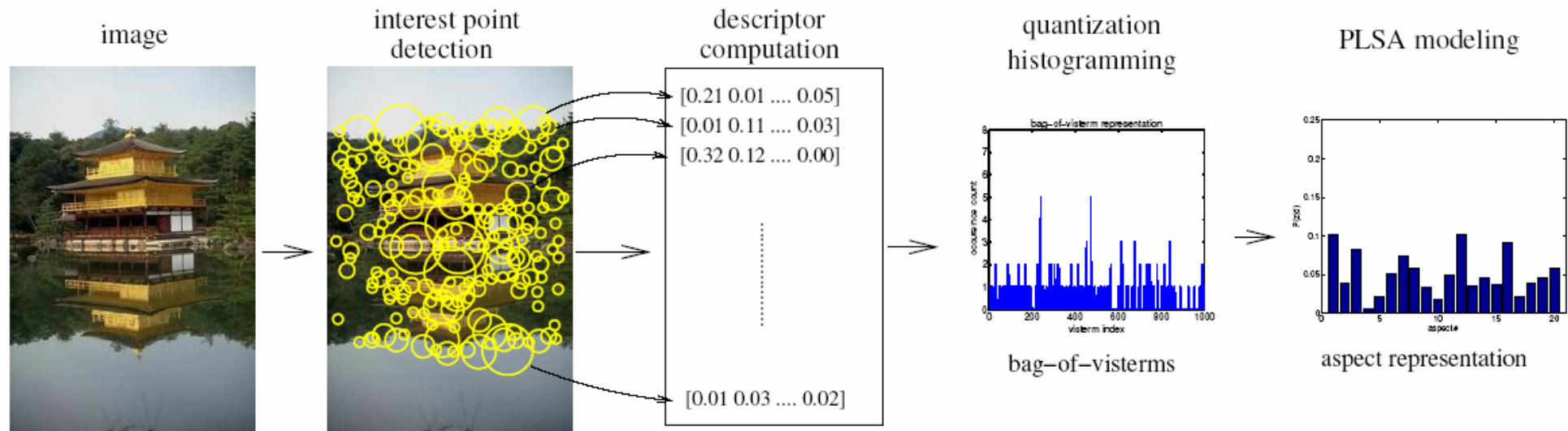


# Image classification using bag-of-words



**Figure 7. Left Panel.** Confusion table of Theme Model 1 using 100 training and 50 test examples from each category, the grid detector and patch based representation. The average performance is 64.0%. **Right Panel.** Rank statistics of the confusion table, which shows the probability of a test scene correctly belong to one of the top  $N$  most probable categories.  $N$  ranges from 1 to 13.

# Image classification using bag-of-words



Flowchart from Quelhas et al., "A thousand words in a scene," IEEE Trans. PAMI, 2007.

- Probabilistic Latent Semantic Analysis (PLSA) is used to learn aspect models to capture co-occurrences of visterms (visual terms).
- Bag-of-visterns representation or the aspect parameters are given as input to Support Vector Machines for classification.

# Image classification using bag-of-words

Total class. error					11.1 (0.8)
Gr. Truth	Classification (%)			Class. Error (%)	# of images
	indoor	city	land.		
indoor	89.7	9.0	1.3	10.3	2777
city	14.5	74.8	10.7	25.2	2505
landscape	1.2	2.0	96.8	3.1	4175

TABLE III

CONFUSION MATRIX FOR THE THREE-CLASS CLASSIFICATION PROBLEM, USING VOCABULARY  $V_{1000}$ .

Total class. error rate: 20.8 (2.1) (Baseline: 30.1 (1.1))							
	m.	f.	i.	c.-p.	c.-s.	error (%)	# of images
mount.	85.8	8.6	2.5	0.5	2.6	14.2	590
forest	8.9	80.3	1.6	2.4	6.7	19.7	492
indoor	0.4	0	91.1	0.4	8.1	8.9	2777
city-pan.	3.5	1.8	8.0	46.9	39.8	53.1	549
city-str.	2.0	2.2	20.8	6.0	68.9	31.1	1957

TABLE V

CLASSIFICATION RATE AND CONFUSION MATRIX FOR THE FIVE-CLASS, USING BOV AND VOCABULARY  $V_{1000}$ .

Total class. error					11.9(1.0)
	indoor	city	land.	class error(%)	# images
indoor	86.6	11.8	1.6	13.4	2777
city	14.8	75.4	9.8	24.5	2505
land.	1.3	1.9	96.8	3.1	4175

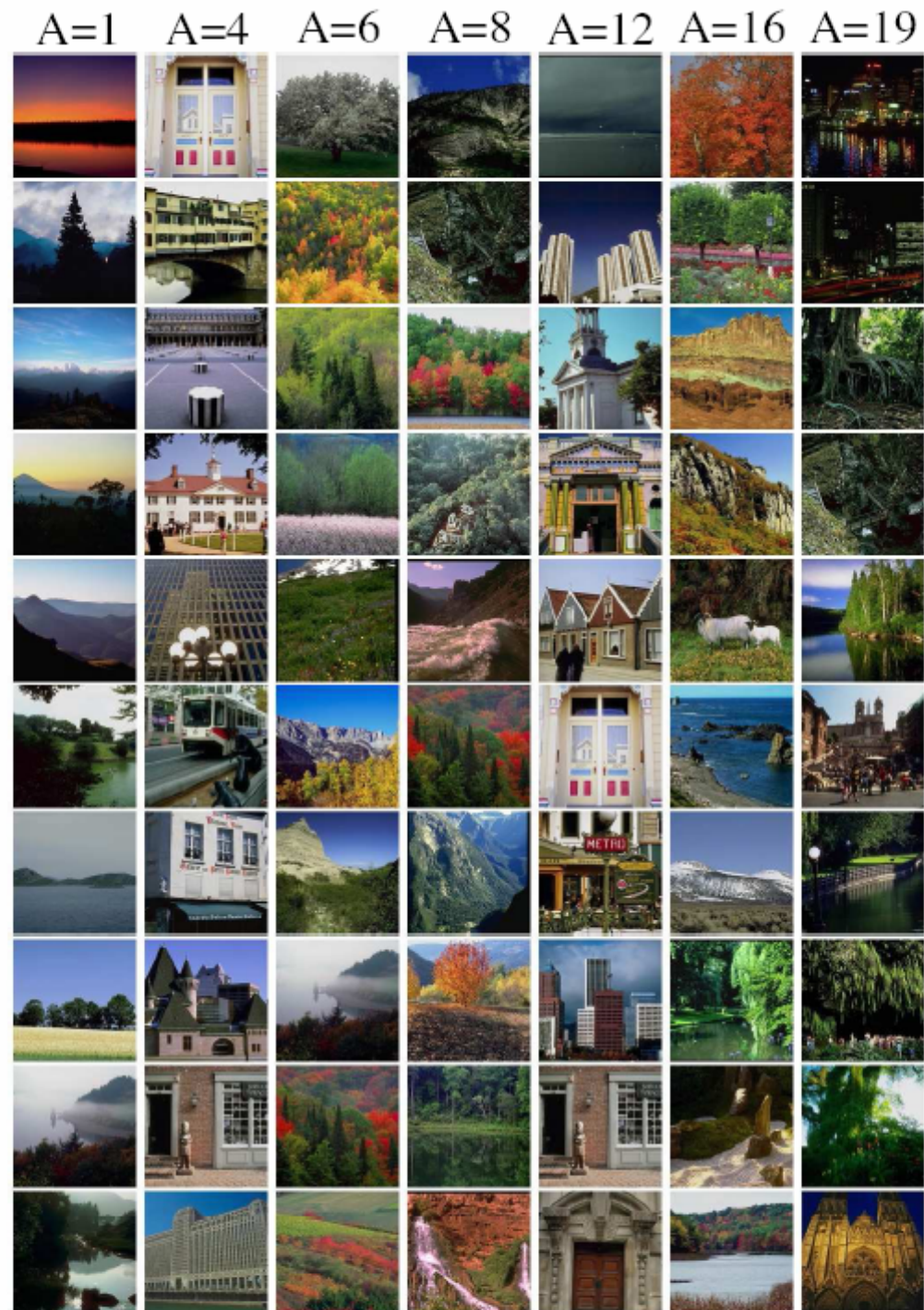
TABLE VIII

CLASSIFICATION ERROR AND CONFUSION MATRIX FOR THE THREE-CLASS PROBLEM USING PLSA, WITH  $V_{1000}$  AND 60 ASPECTS.

Total error rate (BOV: 20.8 (2.1), Baseline: 30.1 (1.1))						
	m.	f.	i.	c.-p.	c.-s.	error (%)
mountain	85.5	12.2	0.8	0.3	1.2	14.5
forest	12.8	78.3	0.8	0.4	7.7	21.7
indoor	0.3	0.1	88.9	0.2	10.5	11.1
city-pan.	3.6	4.9	8.8	12.6	70.1	87.4
city-str.	1.6	1.4	20.4	1.7	74.9	25.1

TABLE X

CLASSIFICATION ERROR AND CONFUSION MATRIX FOR THE FIVE-CLASS PROBLEM USING PLSA-O WITH 60 ASPECTS.





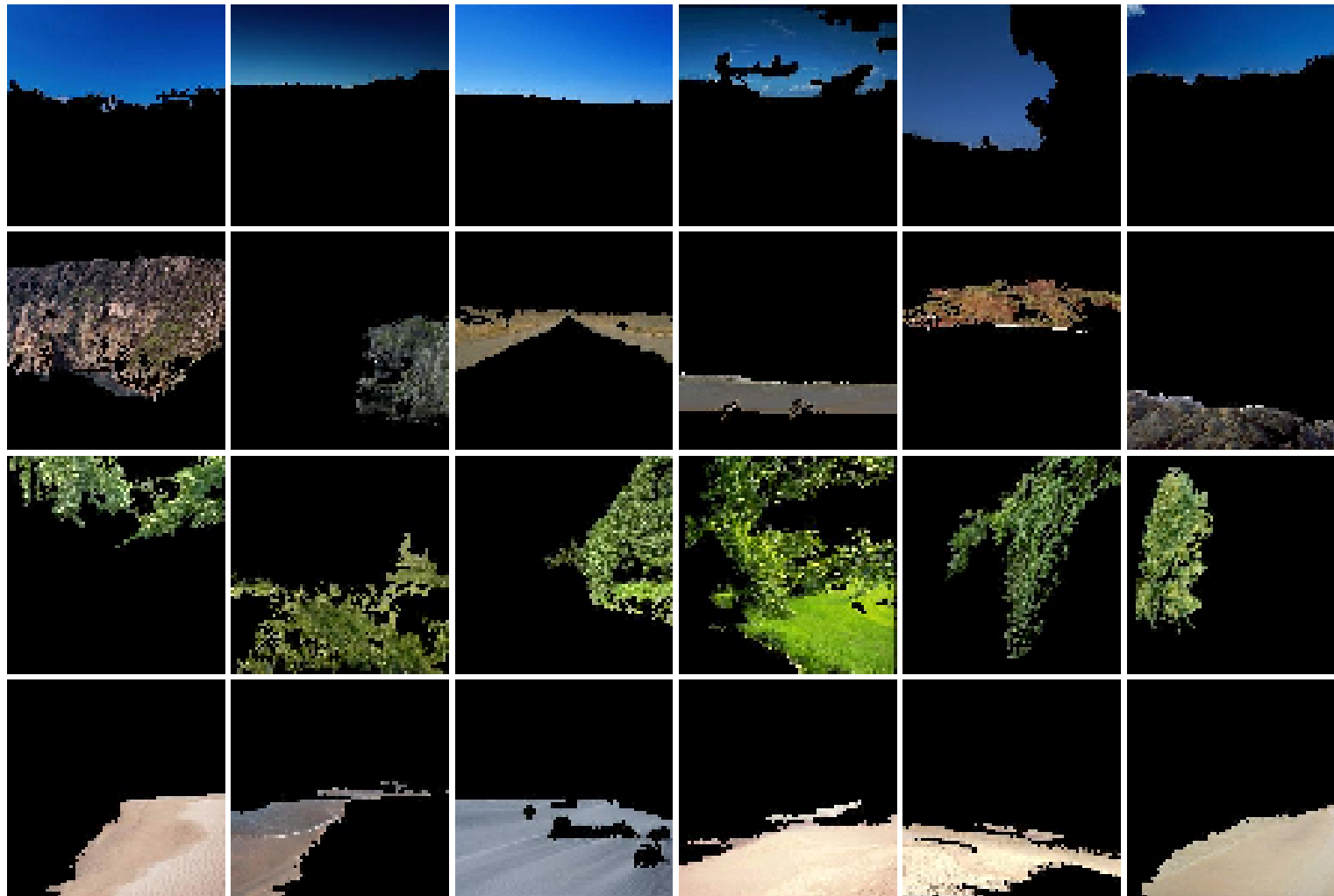
# Image classification using bag-of-regions

---

- D. Gökalp, S. Aksoy, "Scene classification using bag-of-regions representations," IEEE CVPR, Beyond Patches Workshop, 2007.
  - Region segmentation
  - Region clustering → region codebook
  - Above-below spatial relationships → region pairs
  - Statistical region selection: identify region types that
    - are frequently found in a particular class of scenes but rarely exist in other classes, and
    - consistently occur together in the same class of scenes.
  - Bayesian scene classification using
    - bag of individual regions,
    - bag of region pairs.

# Image classification using bag-of-regions

---



Examples for region clusters.  
Each row represents a different cluster.

# Image classification using bag-of-regions

Table 3. Confusion matrix for the bag of individual regions representation after region selection.

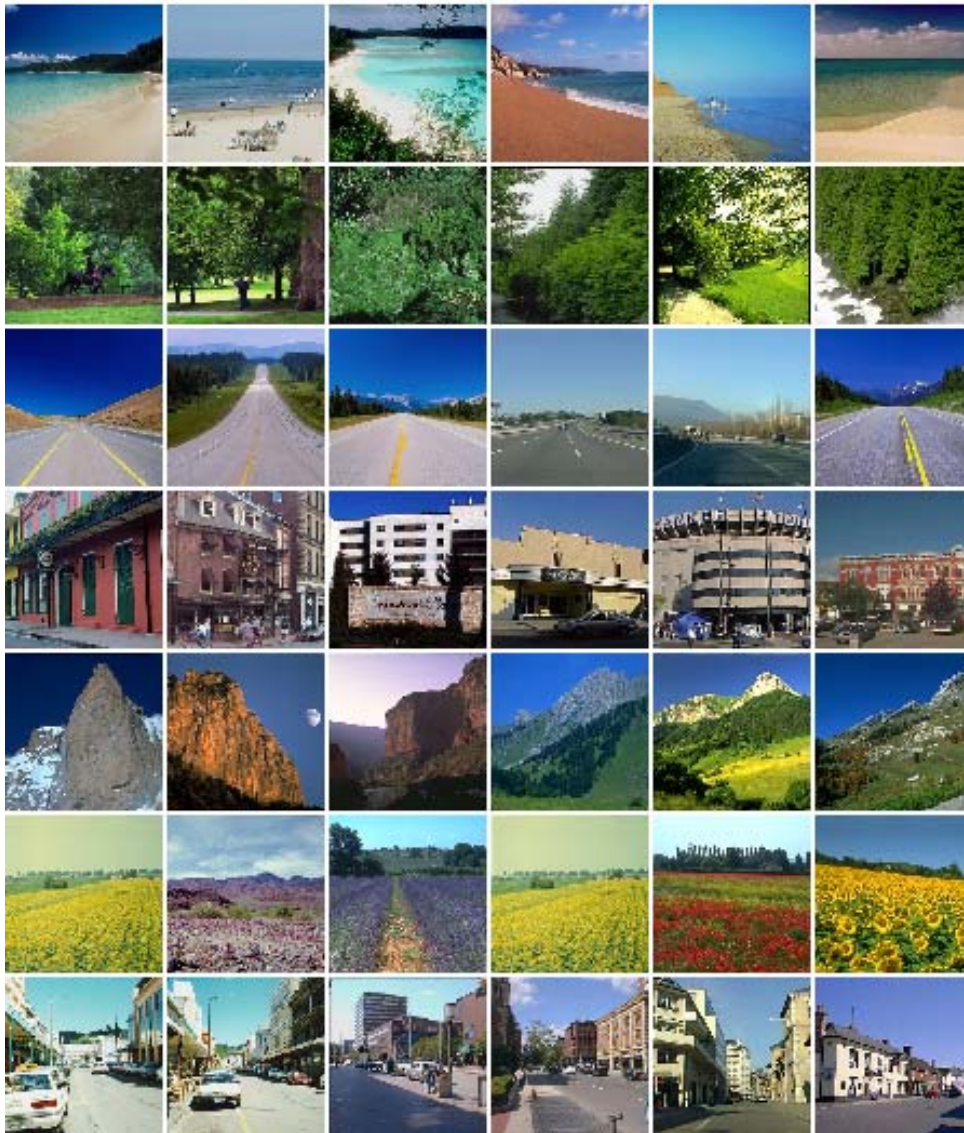
		Assigned							Total	% Agree
		coast	forest	highway	insidecity	mountain	opencountry	street		
True	coast	38	2	2	1	3	4	0	50	76.00
	forest	4	36	0	0	7	2	1	50	72.00
	highway	2	2	32	6	0	2	6	50	64.00
	insidecity	3	1	12	22	2	0	10	50	44.00
	mountain	2	3	5	0	32	6	2	50	64.00
	opencountry	9	8	3	1	14	14	1	50	28.00
	street	0	0	9	6	2	6	27	50	54.00
Total		58	52	63	36	60	34	47	350	57.43

Table 4. Confusion matrix for the bag of region pairs representation after region selection.

		Assigned							Total	% Agree
		coast	forest	highway	insidecity	mountain	opencountry	street		
True	coast	42	0	0	1	3	4	0	50	84.00
	forest	1	38	0	2	4	4	1	50	76.00
	highway	1	1	31	4	2	2	9	50	62.00
	insidecity	3	4	12	19	1	1	10	50	38.00
	mountain	1	5	0	0	40	3	1	50	80.00
	opencountry	8	5	1	2	9	25	0	50	50.00
	street	2	1	8	12	2	3	22	50	44.00
Total		58	54	52	40	61	42	43	350	62.00



# Image classification using bag-of-regions



Examples for correctly classified scenes.



Examples for wrongly classified scenes.



# Image classification using factor graphs

- Boutell et al., "Scene Parsing Using Region-Based Generative Models," IEEE Trans. Multimedia, 2007.

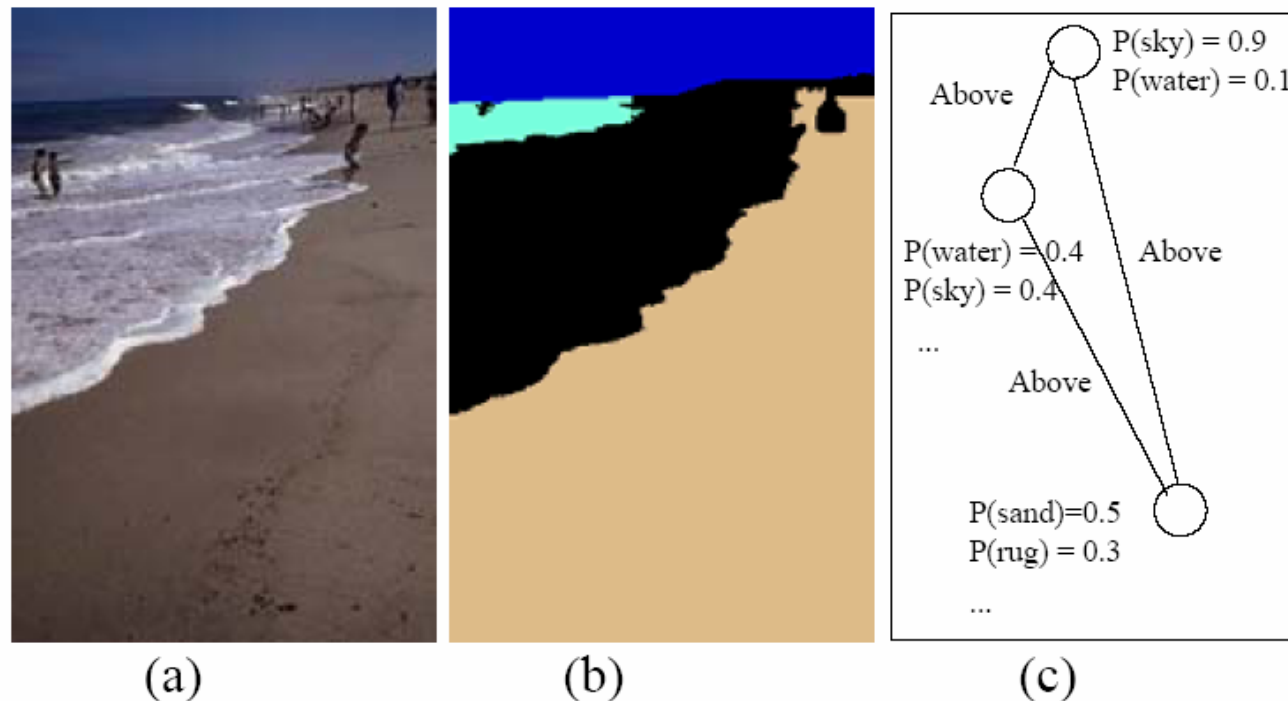


Figure 1. (a) A beach scene. (b) Its manually-labeled materials. The true configuration includes *sky above water*, *water above sand*, and *sky above sand*. (c) The underlying graph showing detector results and spatial relations.

# Recognizing and Learning Object Categories

---

Li Fei-Fei, UIUC

Rob Fergus, MIT

Antonio Torralba, MIT



TM

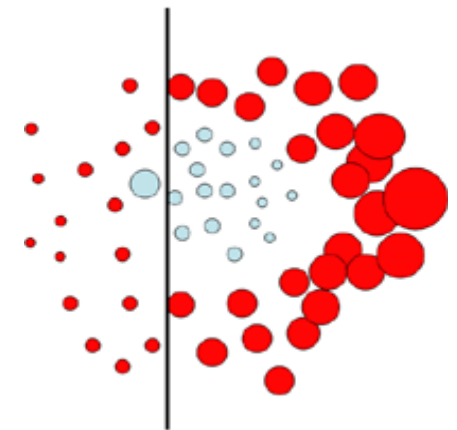
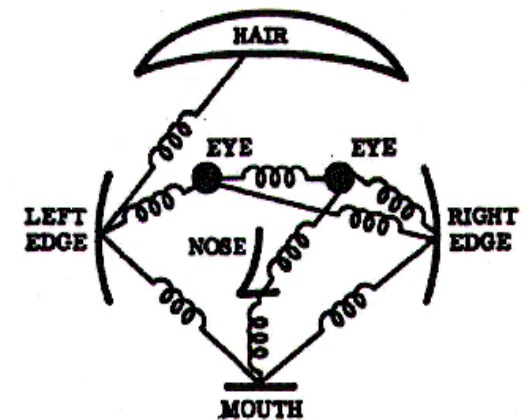
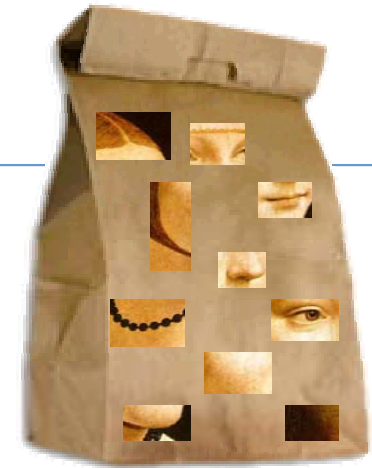
ILLINOIS

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN



# Agenda

- Introduction
- Bag of words models
- Part-based models
- Discriminative methods
- Segmentation and recognition
- Conclusions



**ob·ject**   [Pronunciation Key](#) (ŏb'jĕkt, -jĕkt')  
 n.

1. Something that can be perceived by one or more of the senses, especially sight or touch; a perceptible object.
2. A focus of attention, thought, or action: *an object of contemplation*.
3. The purpose or goal of a specific action or effort: *the object of the game*.
4. Grammar.
  - a. A noun, pronoun, or noun phrase that receives or is affected by the action of a verb within a sentence.
  - b. A noun or substantive governed by a preposition.
5. Philosophy. Something intelligible or perceptible by the mind.
6. Computer Science. A discrete item that can be selected and maneuvered, such as an onscreen graphic. In object-oriented programming, objects include data and the procedures necessary to operate on that data.

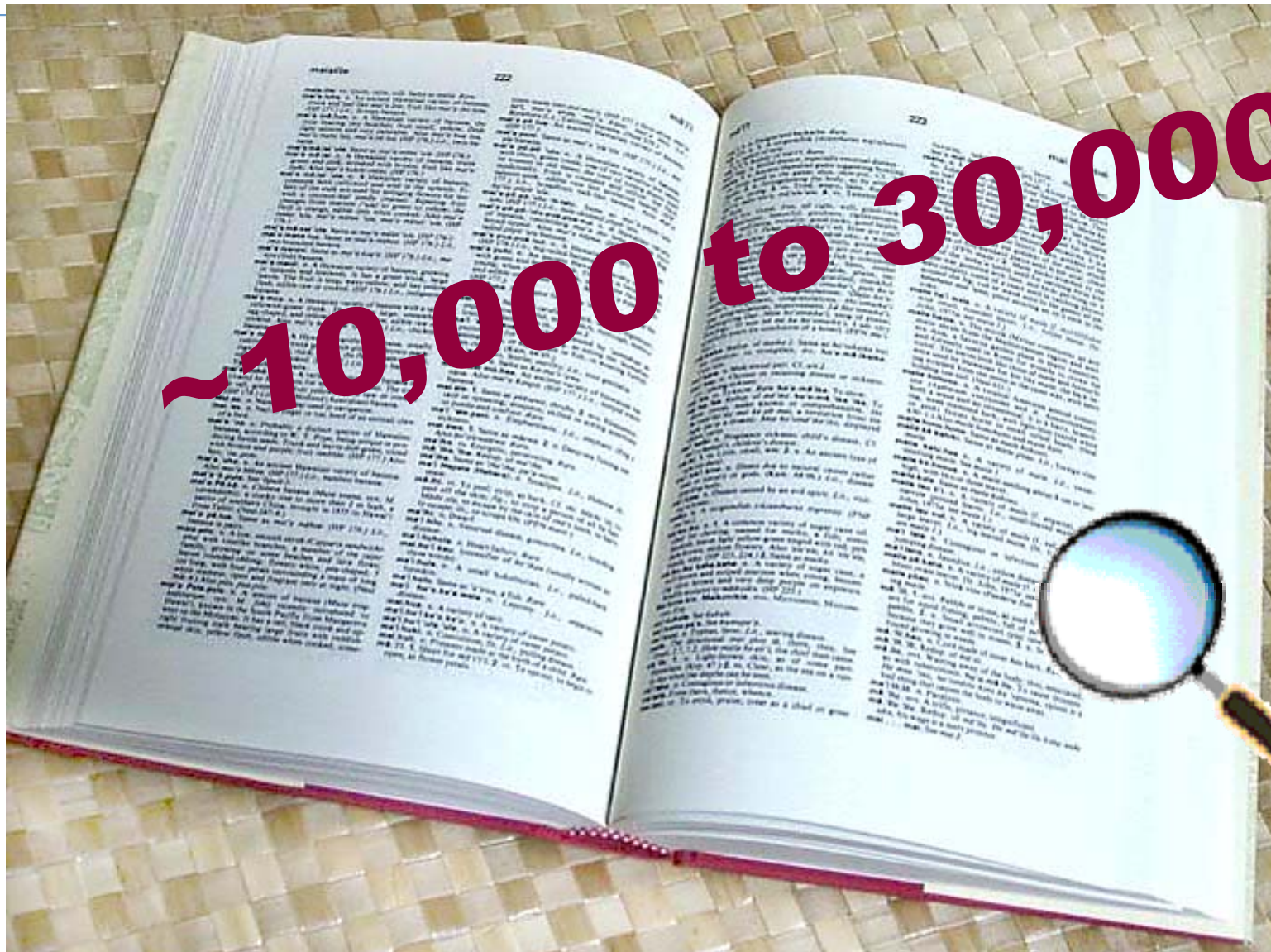
perceptible

vision

material  
thing



# How many object categories are there?



Biederman 1987

So what does object recognition involve?





Verification: is that a bus?



Detection: are there cars?

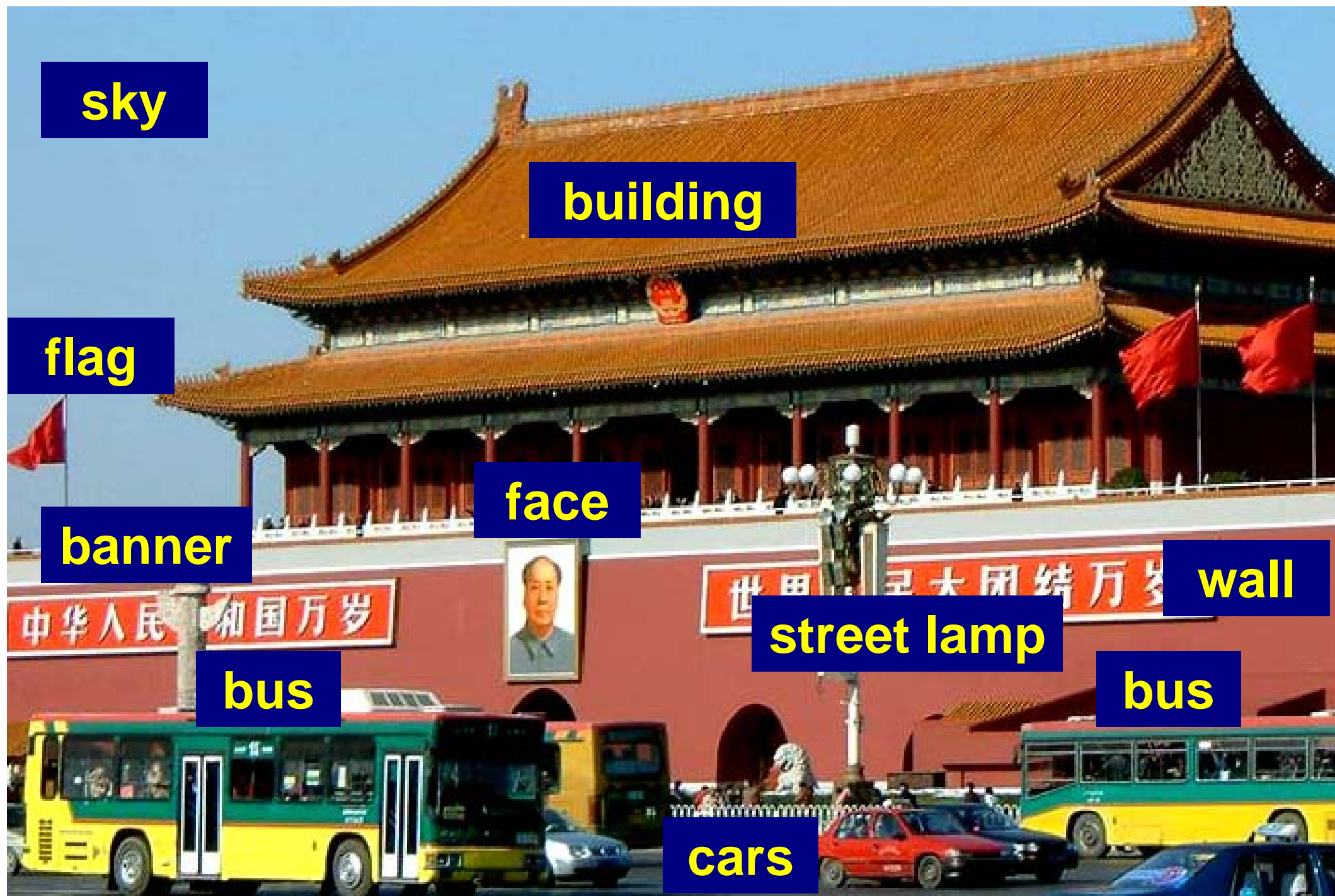




Identification: is that a picture of Mao?

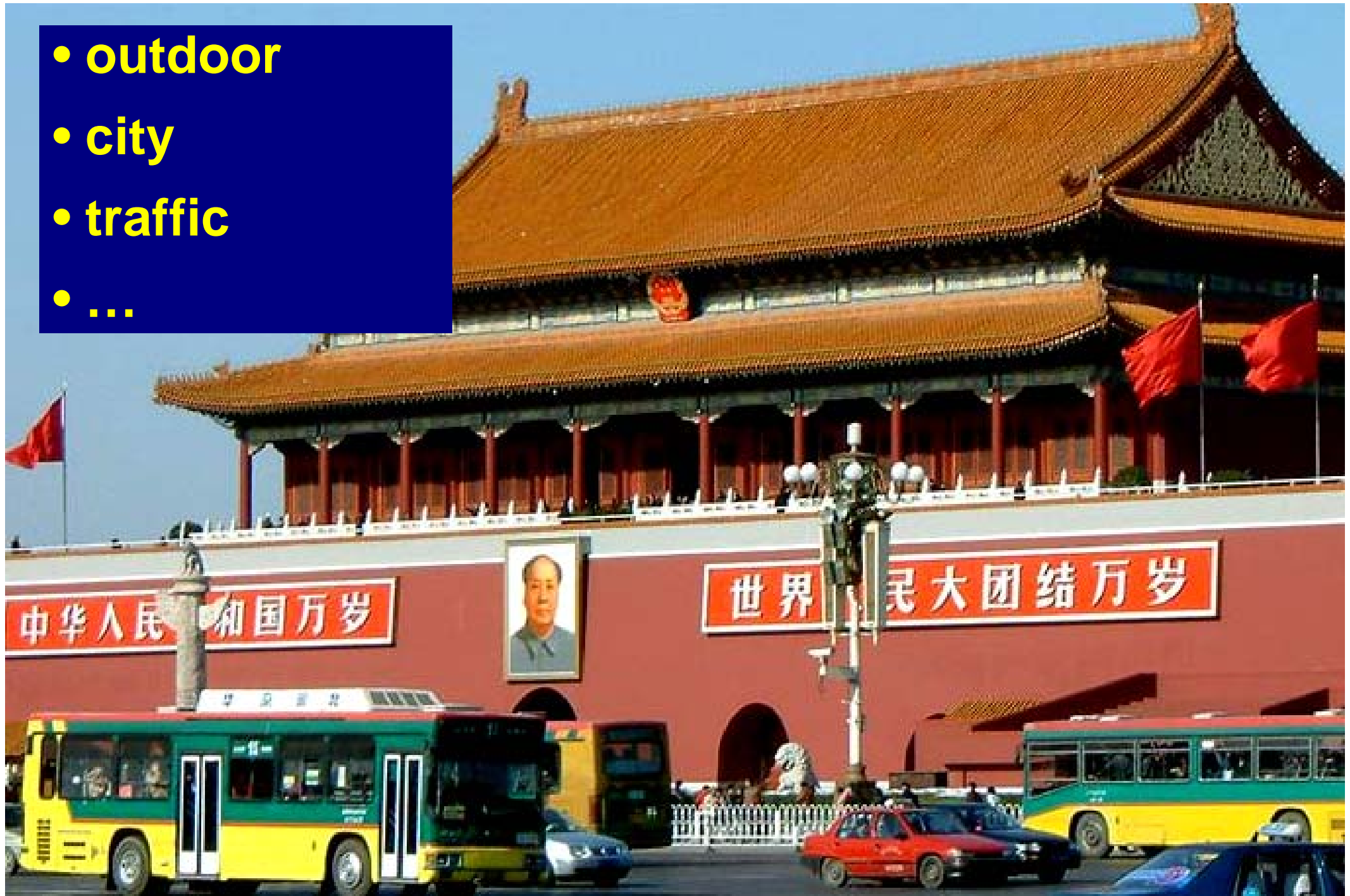


# Object categorization



# Scene and context categorization

- outdoor
- city
- traffic
- ...



# Challenges 1: view point variation

---



Michelangelo 1475-1564

## Challenges 2: illumination



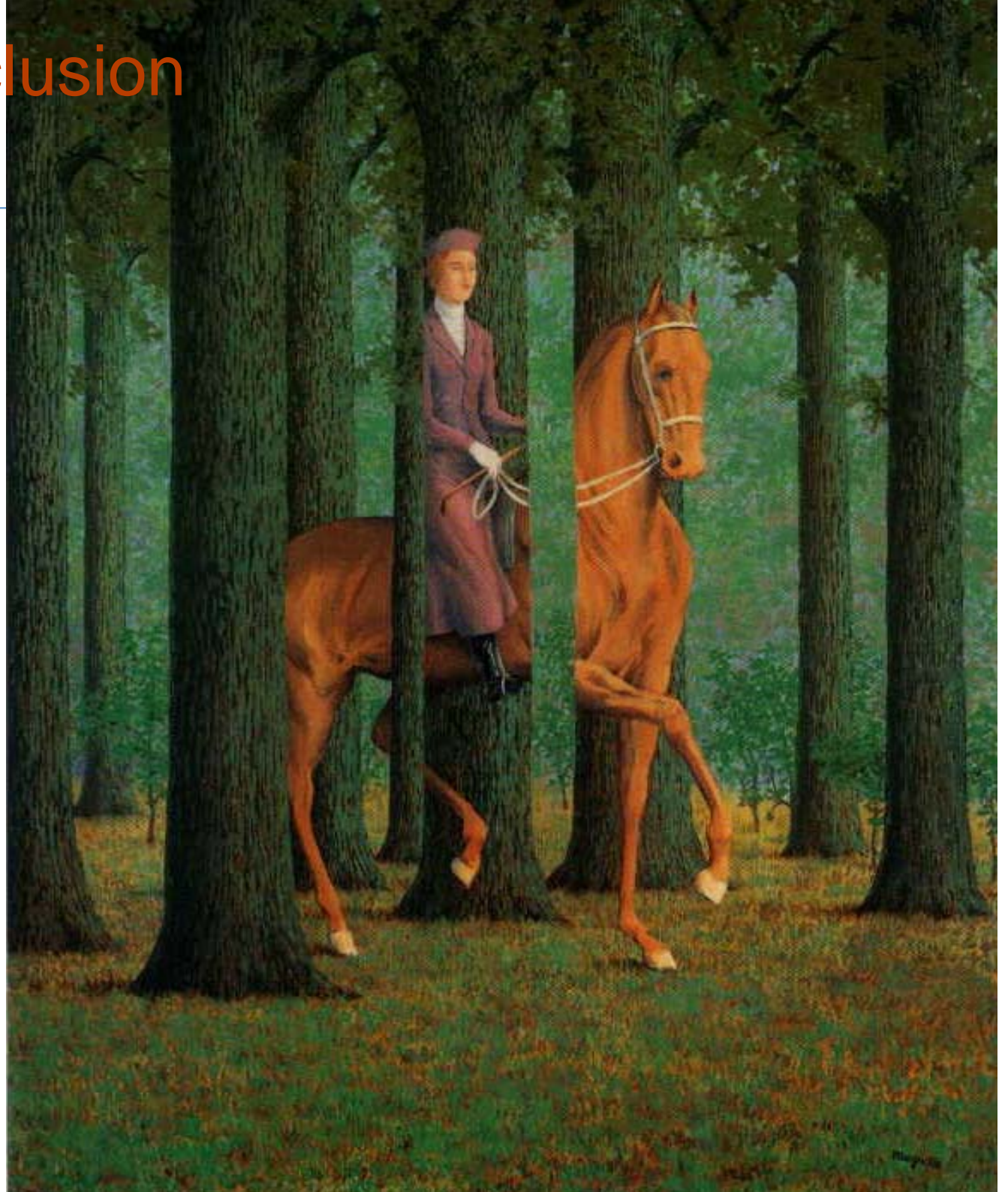
slide credit: S. Ullman



# Challenges 3: occlusion

---

Magritte, 1957



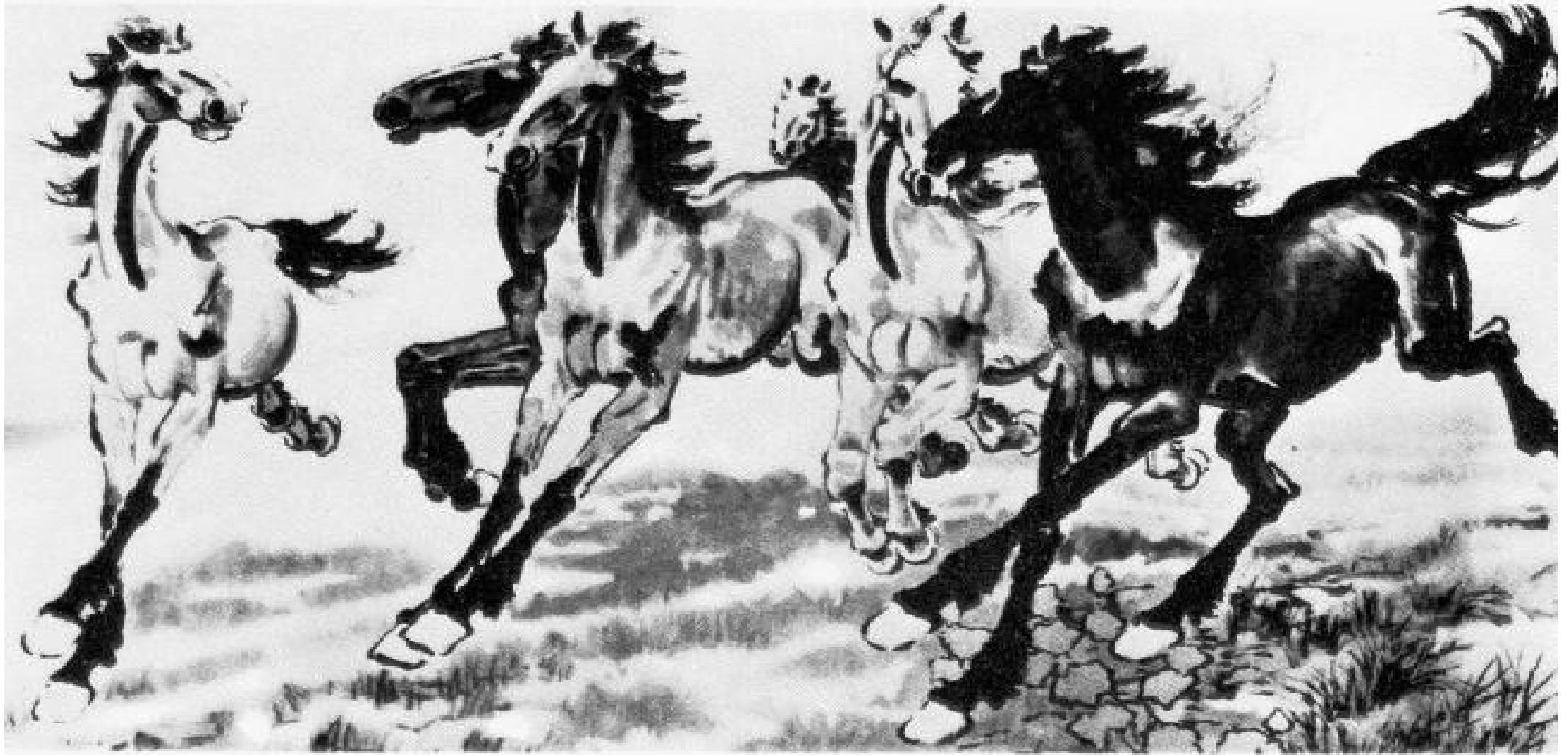
# Challenges 4: scale

---



# Challenges 5: deformation

---



Xu, Beihong 1943



# Challenges 6: background clutter

---



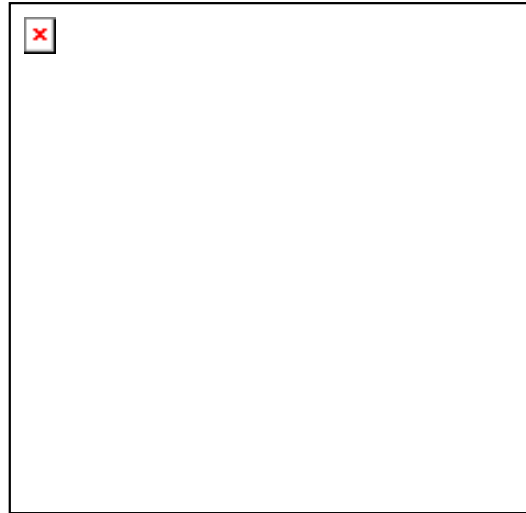
Klimt, 1913

# History: single object recognition





# Challenges 7: intra-class variation



# History: early object categorization

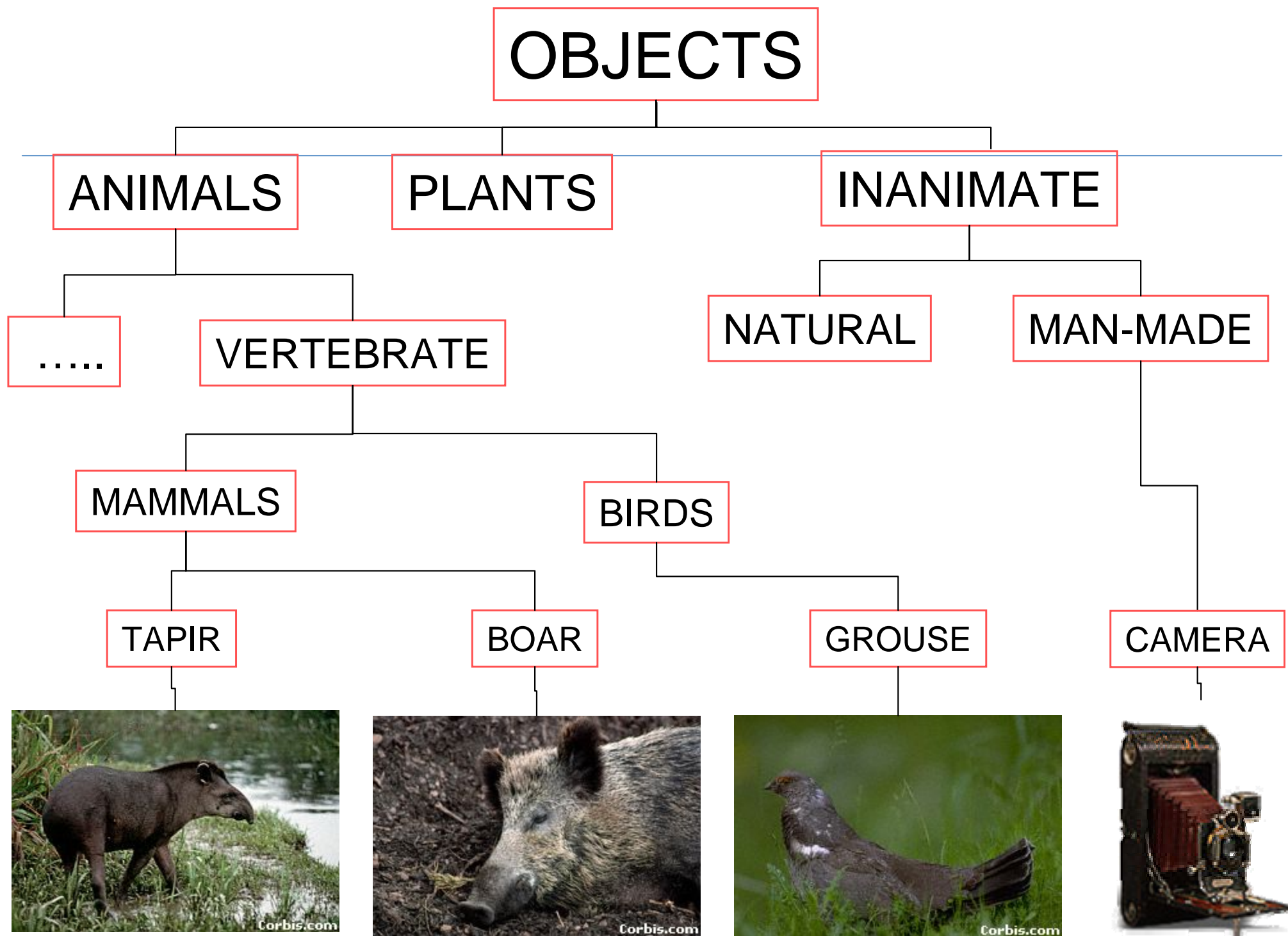


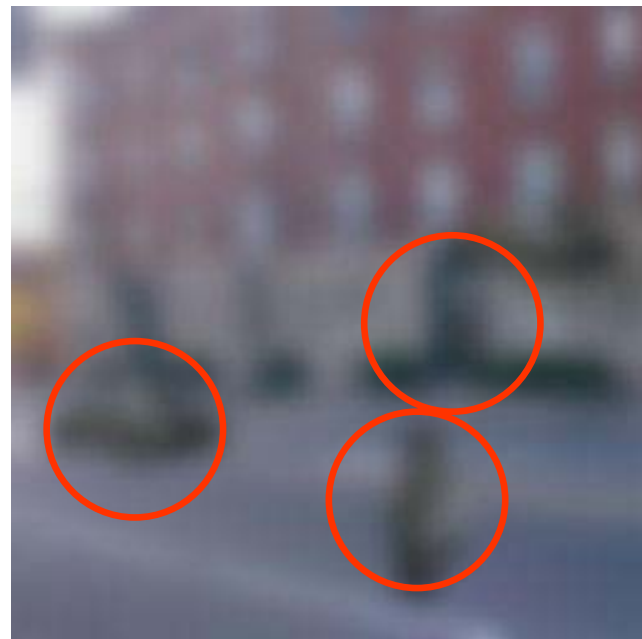
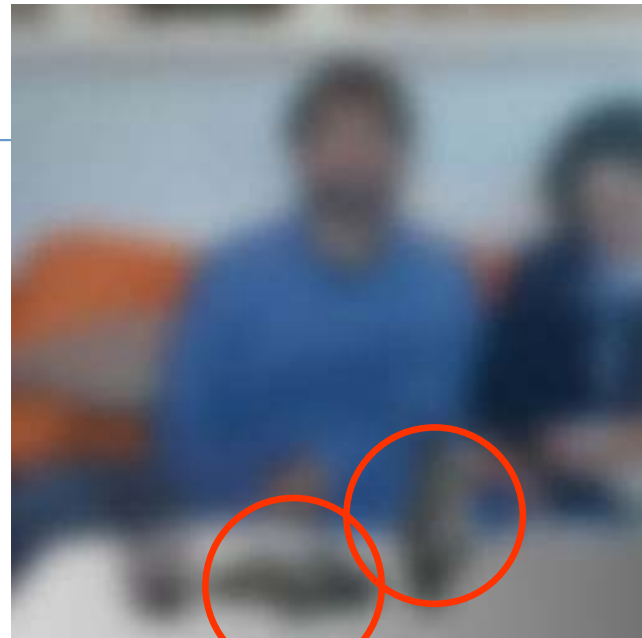
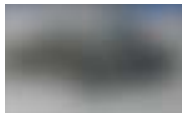
1 7 9 6  
7 8 6 3  
2 1 7 9 7 1 2  
4 8 1 9 0 1 8  
7 6 1 8 6 4 1 5 0 0  
7 5 9 2 6 5 8 1 9 7  
2 2 2 2 2 3 4 4 8 0  
0 2 3 8 0 7 3 8 5 7  
0 1 4 6 4 6 0 2 4 3  
7 1 2 8 7 6 9 8 6 1



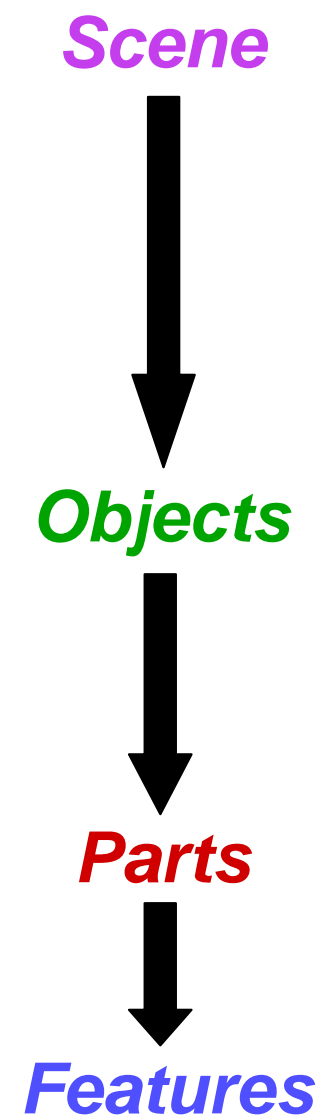
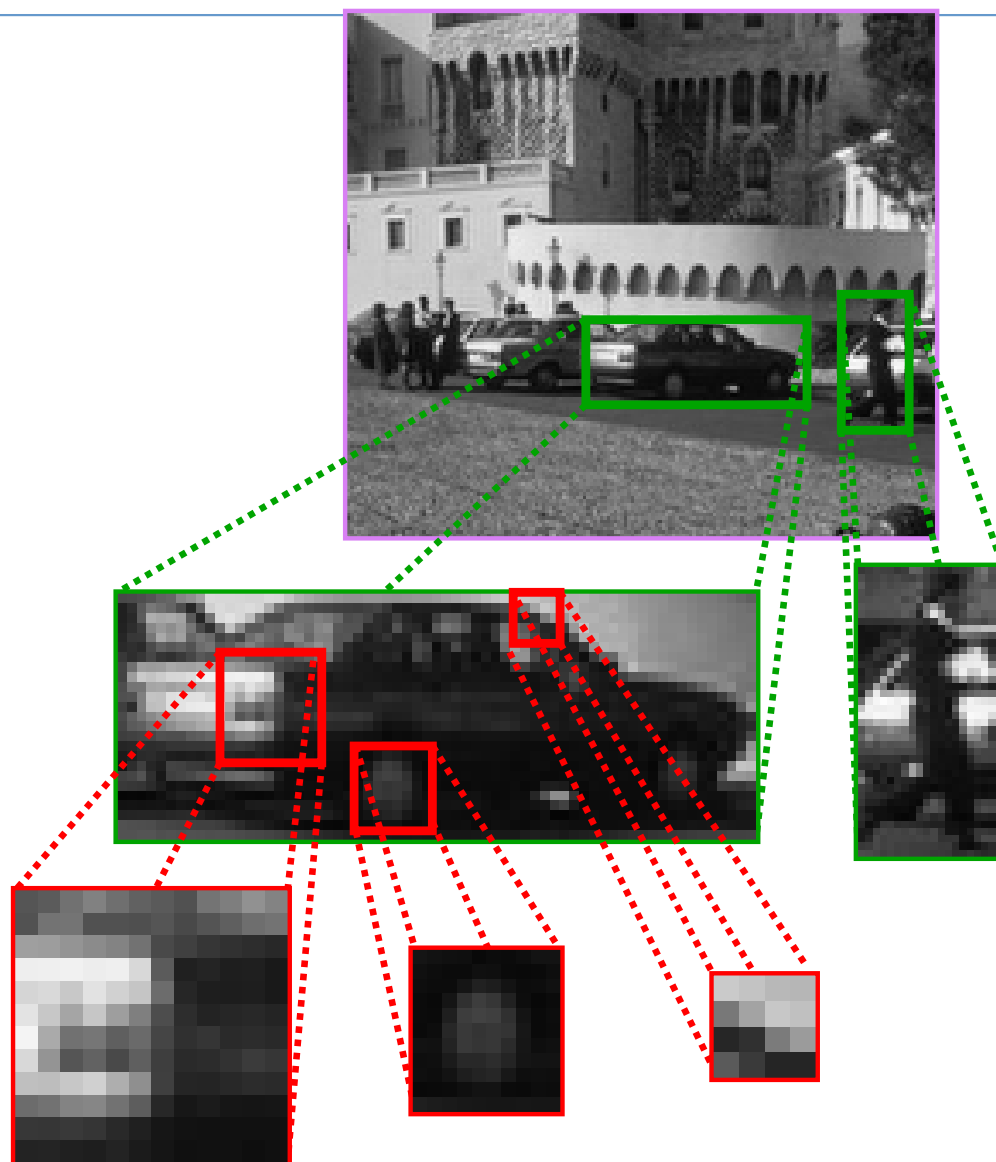
**~10,000 to 30,000**







# Scenes, Objects, and Parts



E. Sudderth, A. Torralba, W. Freeman, A. Willsky. ICCV 2005.



# Object categorization: the statistical viewpoint



$$p(\text{zebra} \mid \text{image})$$

vs.

$$p(\text{no zebra} \mid \text{image})$$

## ■ Bayes rule:

$$\underbrace{\frac{p(\text{zebra} \mid \text{image})}{p(\text{no zebra} \mid \text{image})}}_{\text{posterior ratio}} = \underbrace{\frac{p(\text{image} \mid \text{zebra})}{p(\text{image} \mid \text{no zebra})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{zebra})}{p(\text{no zebra})}}_{\text{prior ratio}}$$

# Object categorization: the statistical viewpoint

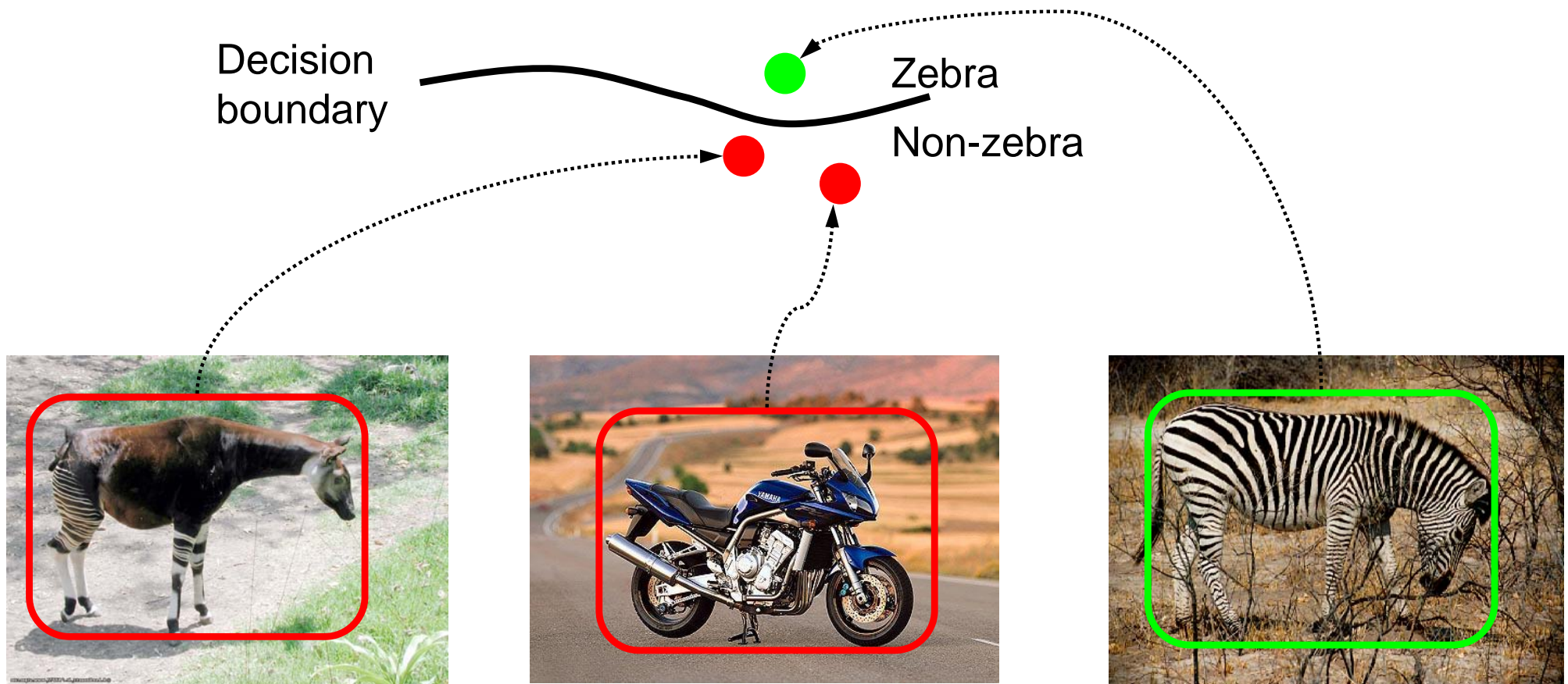
---

$$\underbrace{\frac{p(\text{zebra} | \text{image})}{p(\text{no zebra} | \text{image})}}_{\text{posterior ratio}} = \underbrace{\frac{p(\text{image} | \text{zebra})}{p(\text{image} | \text{no zebra})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{zebra})}{p(\text{no zebra})}}_{\text{prior ratio}}$$

- Discriminative methods model posterior
- Generative methods model likelihood and prior

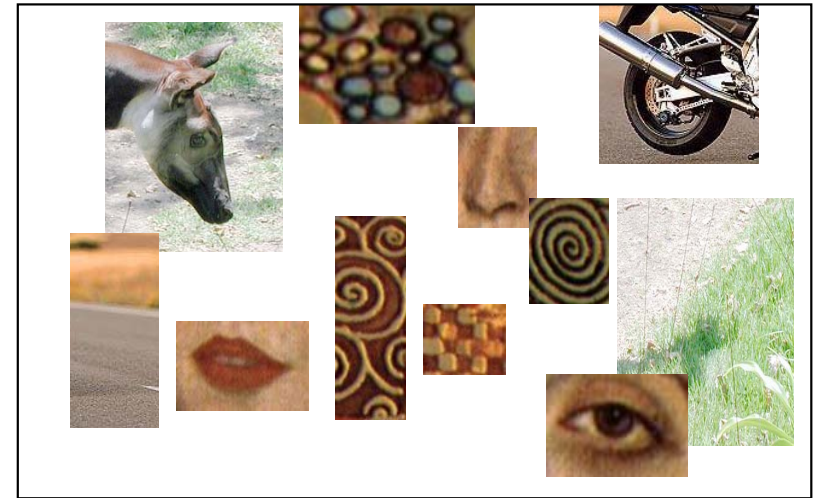
# Discriminative

- Direct modeling of  $\frac{p(\text{zebra} | \text{image})}{p(\text{no zebra} | \text{image})}$



# Generative

- Model  $p(\text{image} | \text{zebra})$  and  $p(\text{image} | \text{no zebra})$



$p(image  zebra)$	$p(image  no\ zebra)$
Low	Middle
High	Middle $\rightarrow$ Low

© 2008, Selim Aksoy



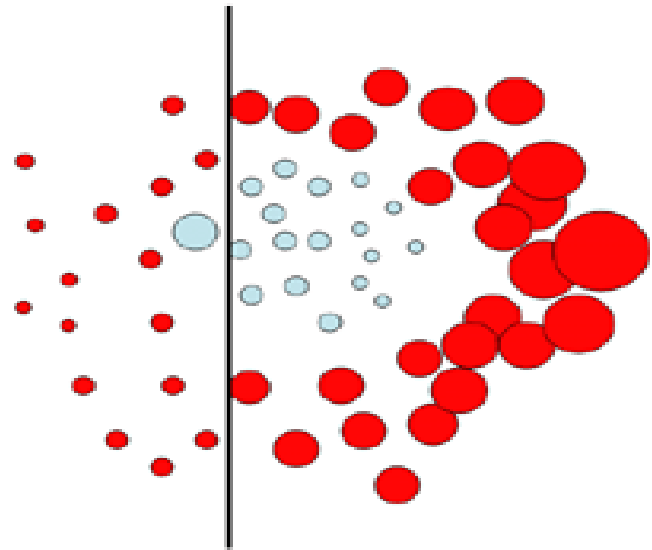
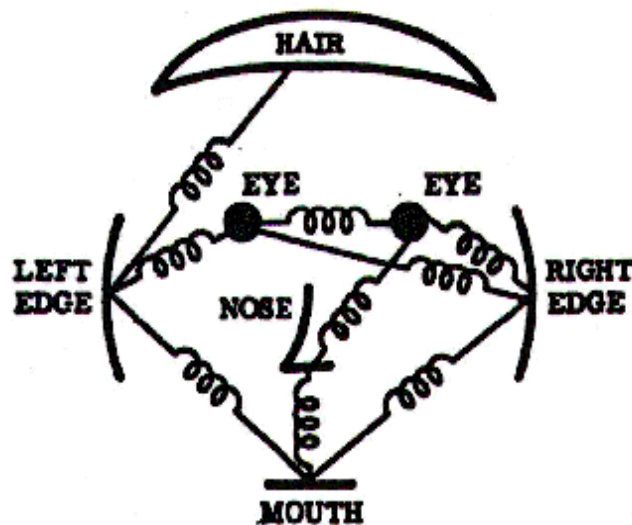
# Three main issues

---

- Representation
  - How to represent an object category
- Learning
  - How to form the classifier, given training data
- Recognition
  - How the classifier is to be used on novel data

# Representation

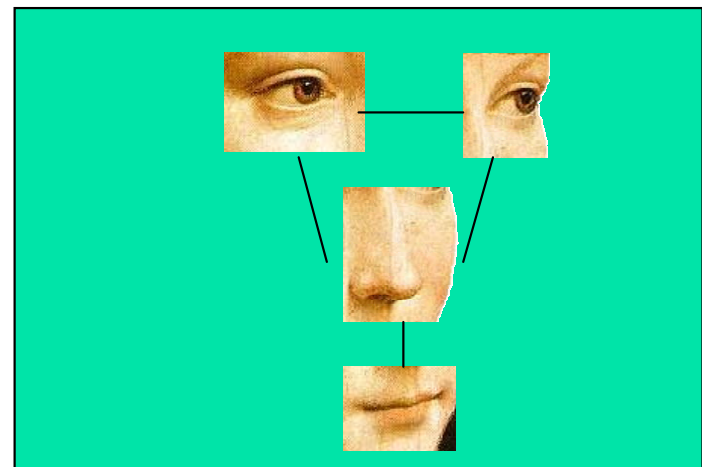
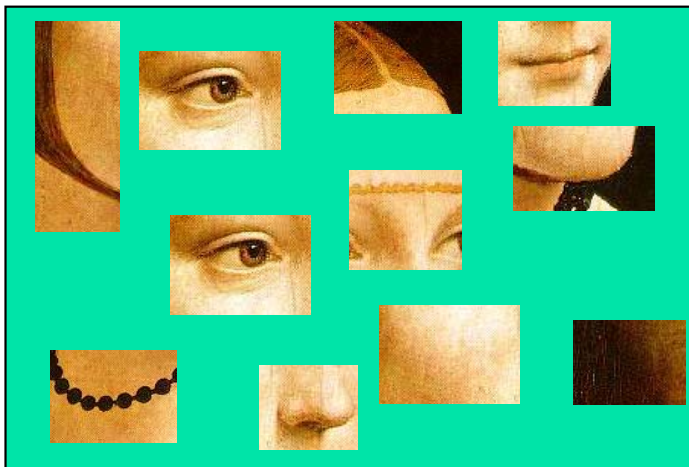
- Generative / discriminative / hybrid



# Representation

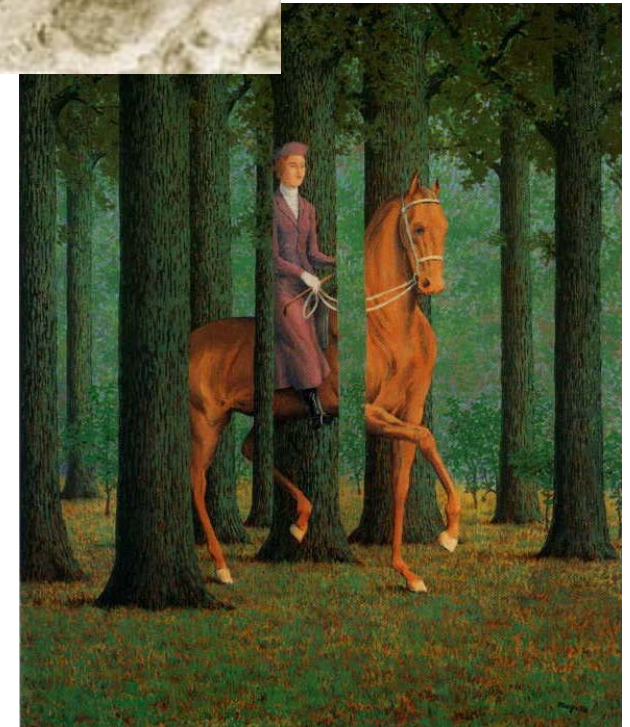
---

- Generative / discriminative / hybrid
- Appearance only or location and appearance



# Representation

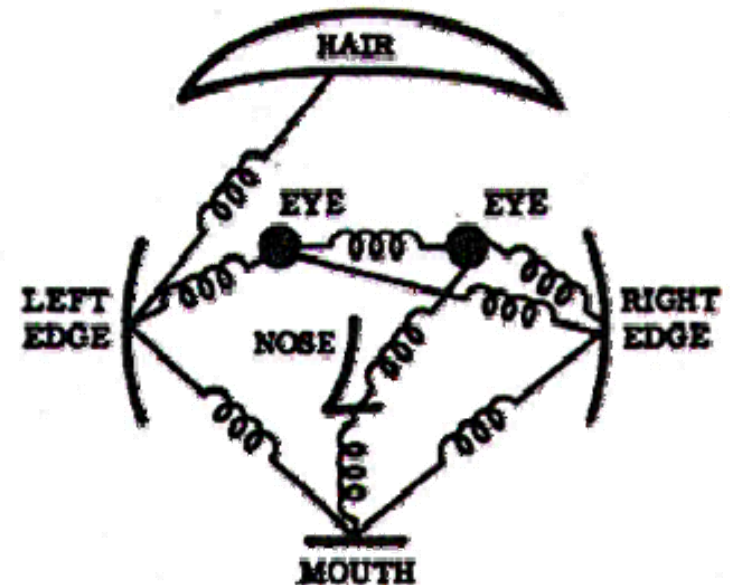
- Generative / discriminative / hybrid
- Appearance only or location and appearance
- Invariances
  - View point
  - Illumination
  - Occlusion
  - Scale
  - Deformation
  - Clutter
  - etc.





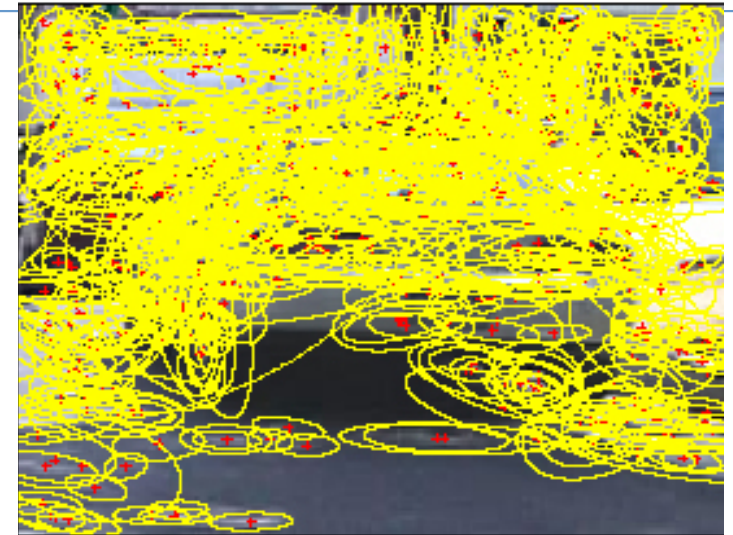
# Representation

- Generative / discriminative / hybrid
- Appearance only or location and appearance
- invariances
- Part-based or global w/sub-window



# Representation

- Generative / discriminative / hybrid
- Appearance only or location and appearance
- invariances
- Parts or global w/sub-window
- Use set of features or each pixel in image



# Learning

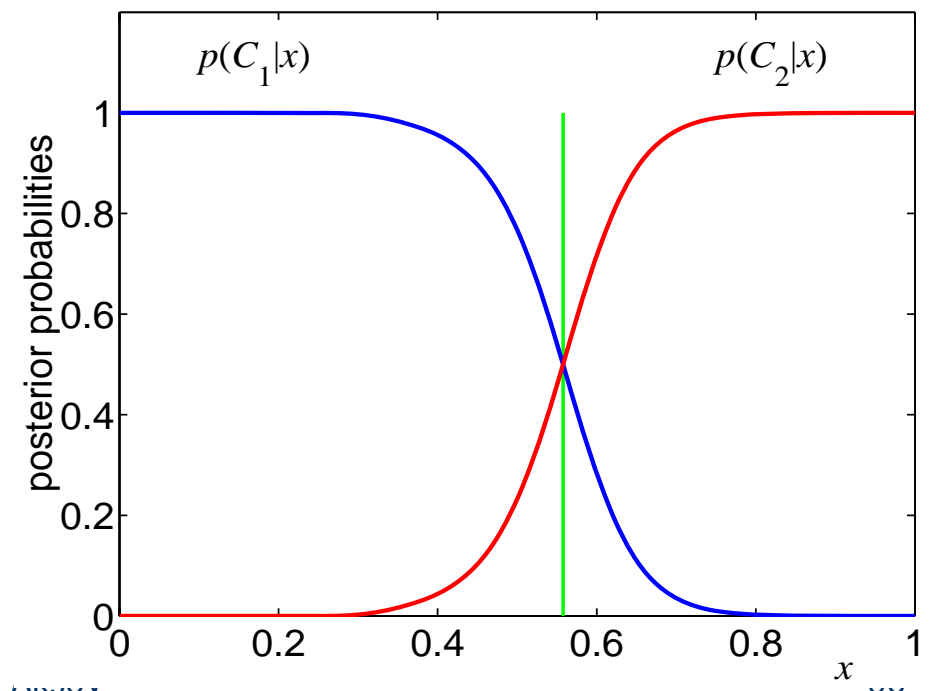
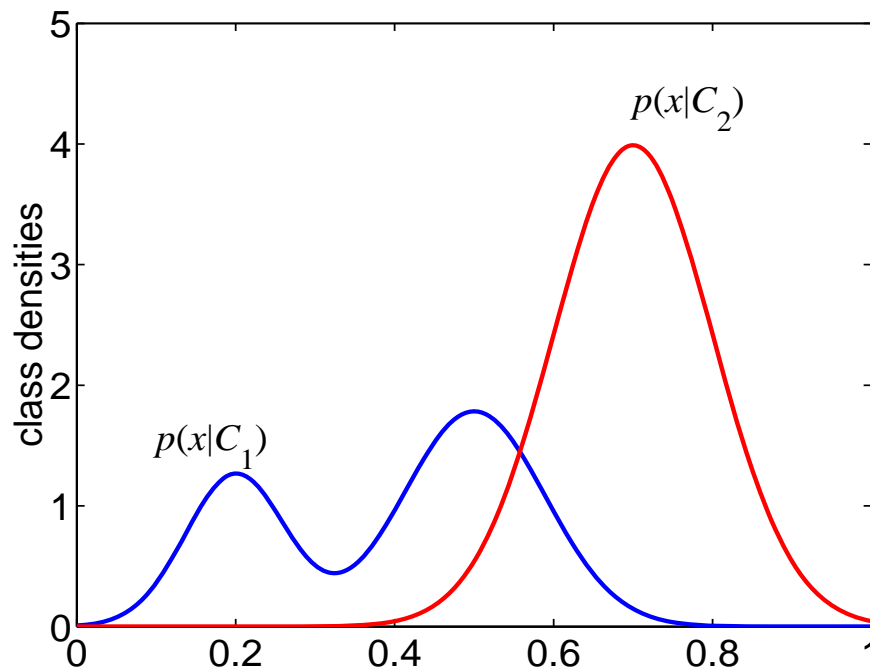
---

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning



# Learning

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning)
- Methods of training: generative vs. discriminative



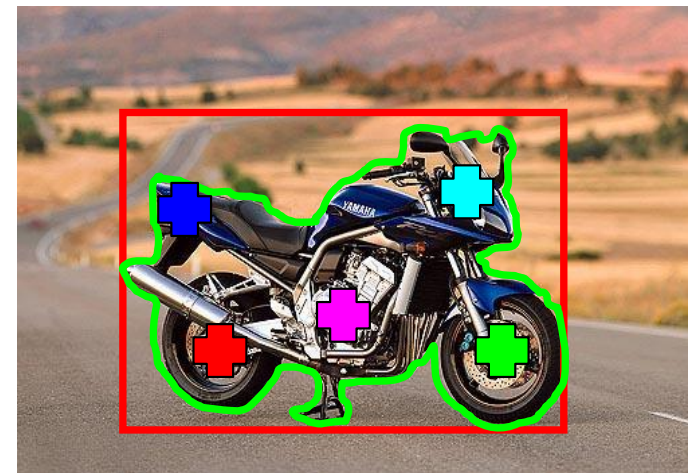


# Learning

---

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning)
- What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels

Contains a motorbike



# Learning

---

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning)
- What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels
- Batch/incremental (on category and image level; user-feedback )

# Learning

---

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning)
- What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels
- Batch/incremental (on category and image level; user-feedback )
- Training images:
  - Issue of overfitting
  - Negative images for discriminative methods

# Learning

---

- Unclear how to model categories, so we learn what distinguishes them rather than manually specify the difference -- hence current interest in machine learning)
- What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels
- Batch/incremental (on category and image level; user-feedback )
- Training images:
  - Issue of overfitting
  - Negative images for discriminative methods
- Priors



# Recognition

---

- Scale / orientation range to search over
- Speed



# Object Class Recognition using Images of Abstract Regions

---

Yi Li, Jeff A. Bilmes, and Linda G. Shapiro  
Department of Computer Science and Engineering  
Department of Electrical Engineering  
University of Washington

# Problem Statement

**Given:** Some images and their corresponding descriptions



{trees, grass, cherry trees}



{cheetah, trunk}



{mountains, sky}



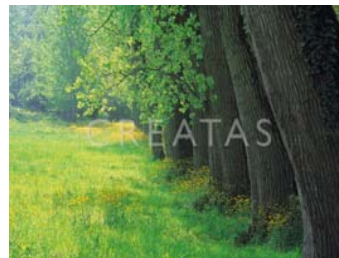
{beach, sky, trees, water}

...

**To solve:** What object classes are present in new images



?



?



?



?

...

# Abstract Regions

---

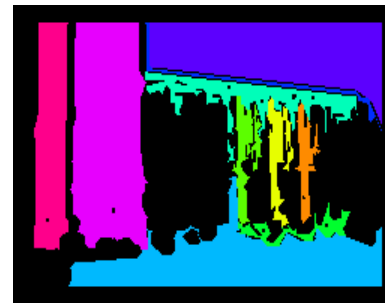
Original Images



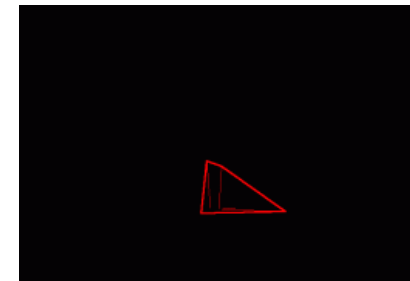
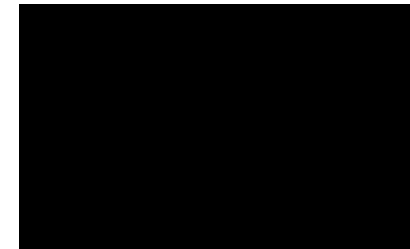
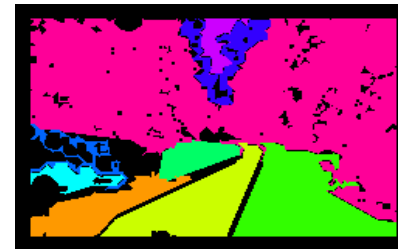
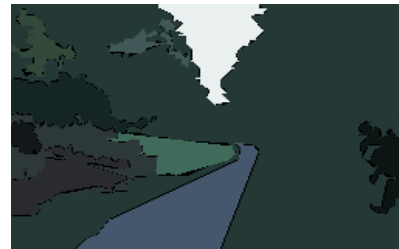
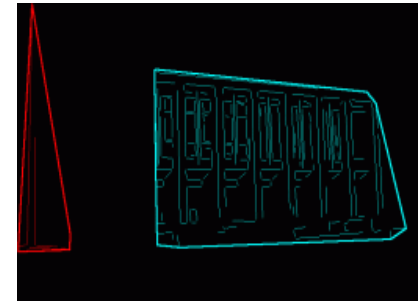
Color Regions



Texture Regions

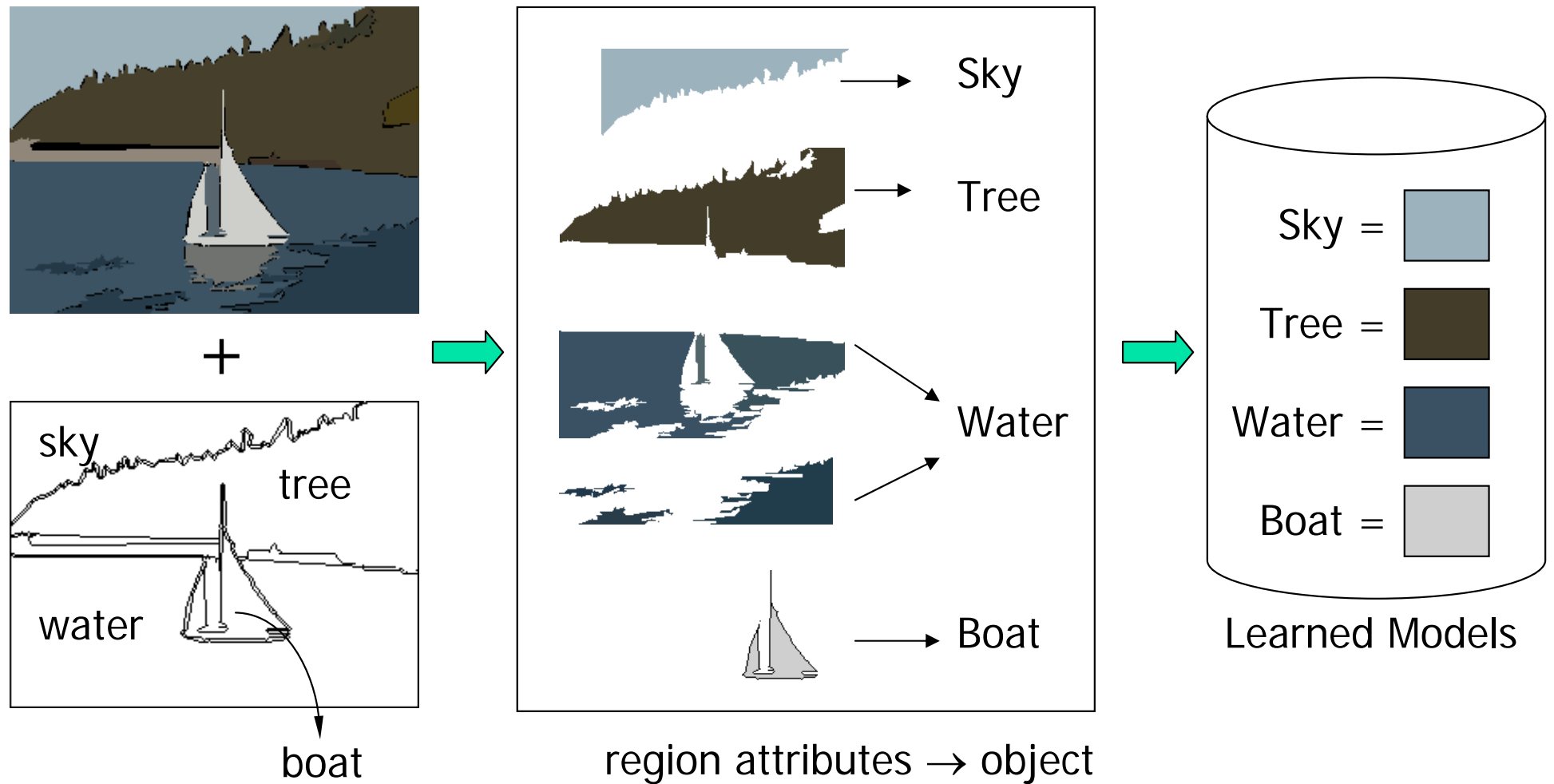


Line Clusters





# Object Model Learning (Ideal)



# Our Scenario: Abstract Regions

Multiple segmentations whose regions are not labeled;  
a list of labels is provided for each training image.

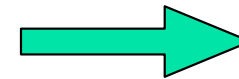
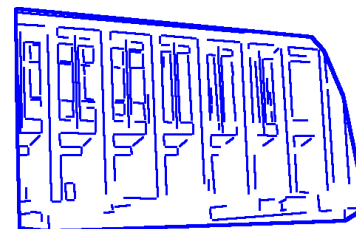
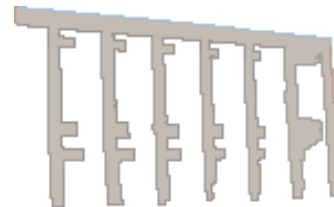
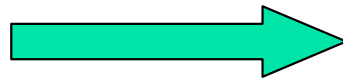
image



labels

{sky, building}

various different  
segmentations



region  
attributes  
from several  
different  
types of  
regions

# Object Model Learning

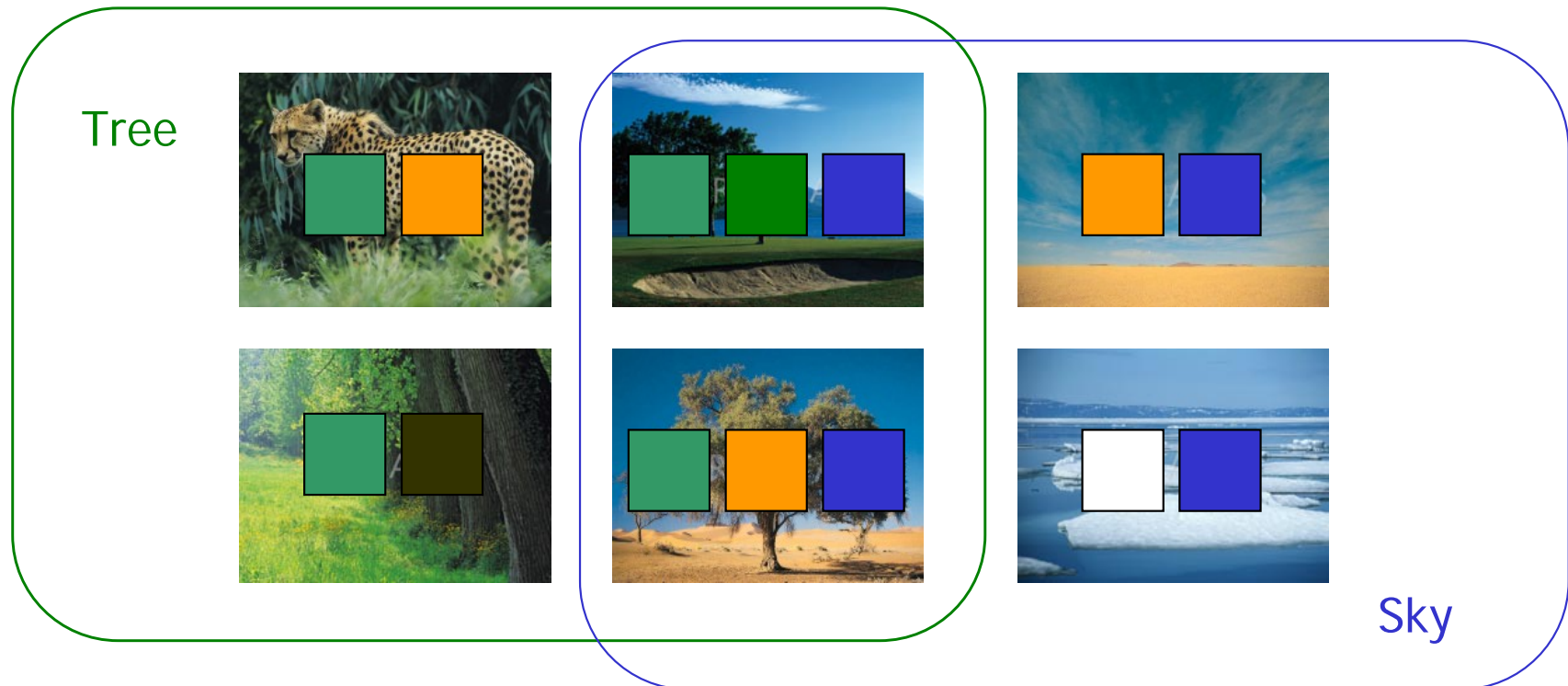
---

## Assumptions

- The feature distribution of each object within a region is a Gaussian;
- Each image is a set of regions; each region can be modeled as a mixture of multivariate Gaussian distributions.

# Model Initial Estimation

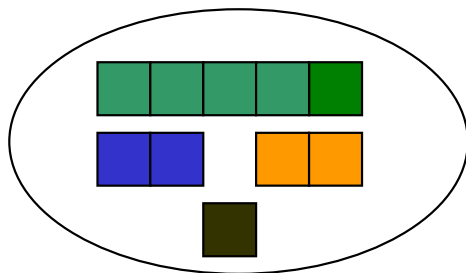
- Estimate the initial model of an object using all the region features from all images that contain the object



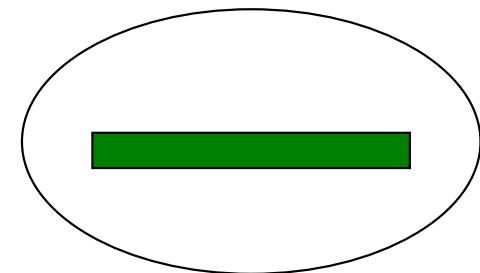


# Expectation-Maximization

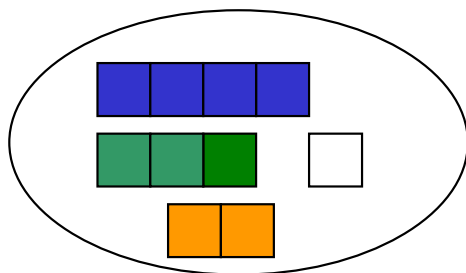
Initial Model for “trees”



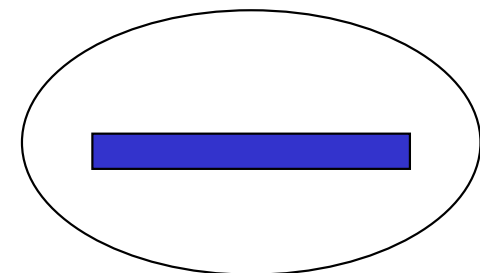
Final Model for “trees”



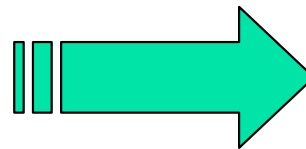
Initial Model for “sky”



Final Model for “sky”

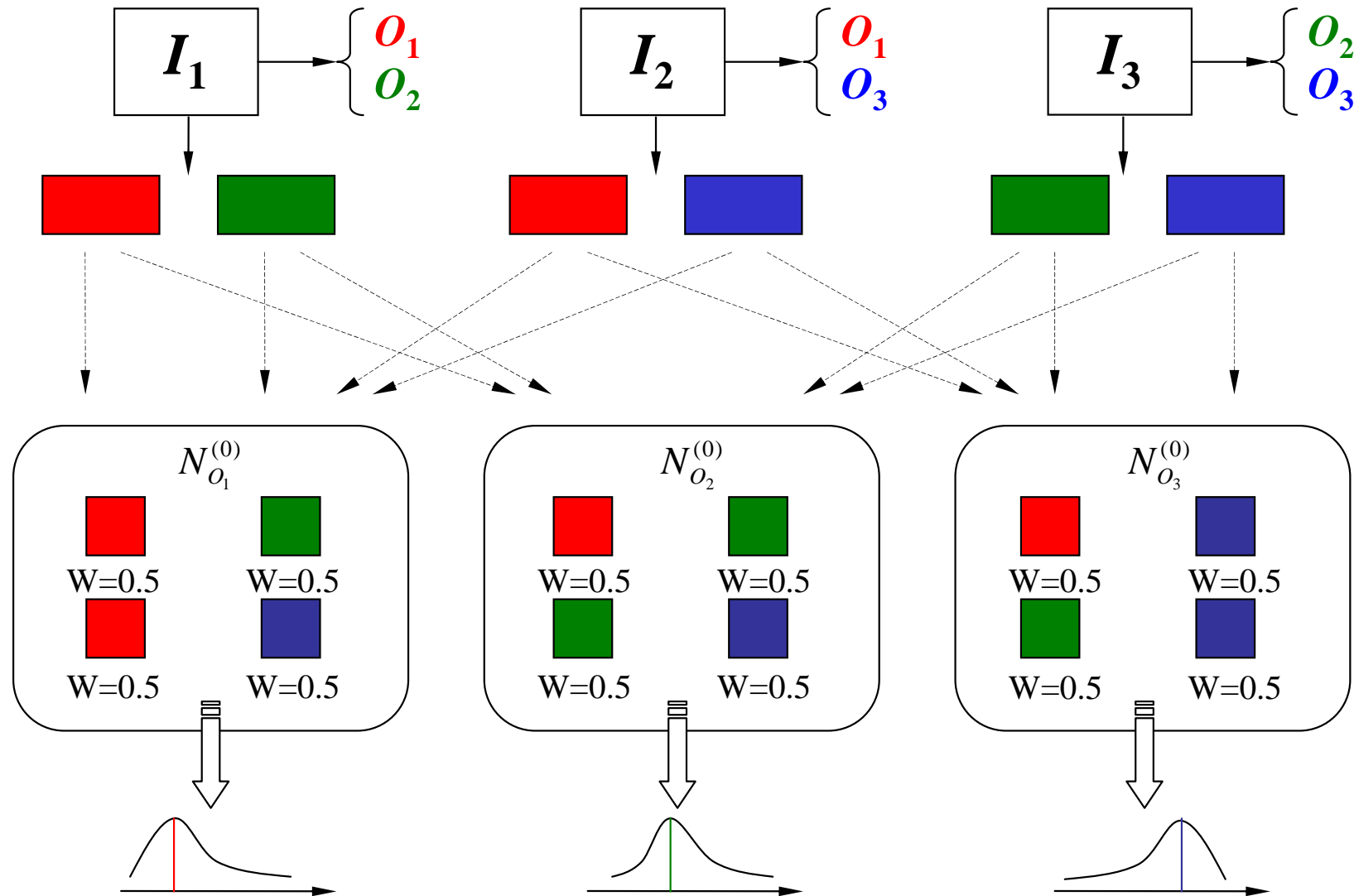


EM

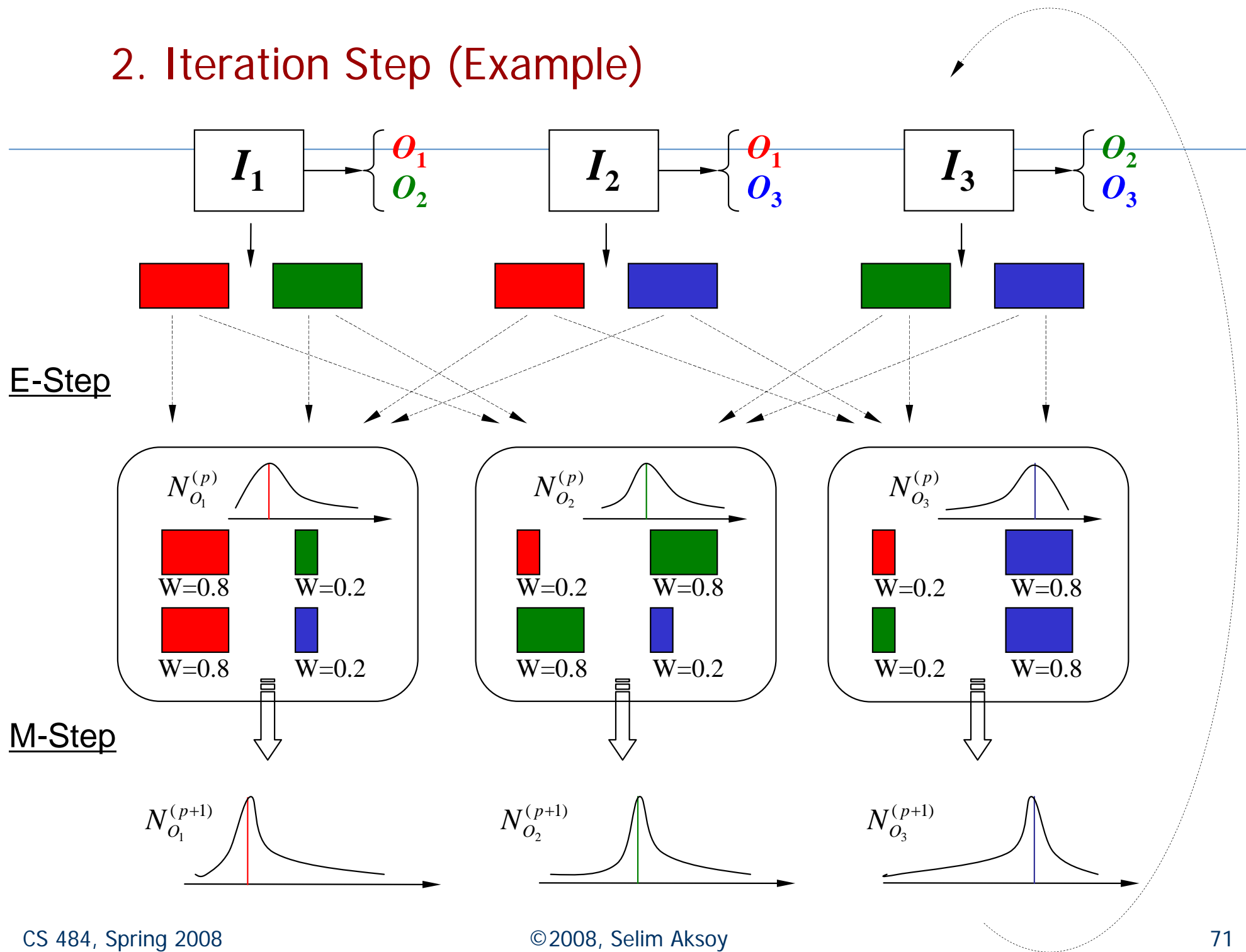


# 1. Initialization Step (Example)

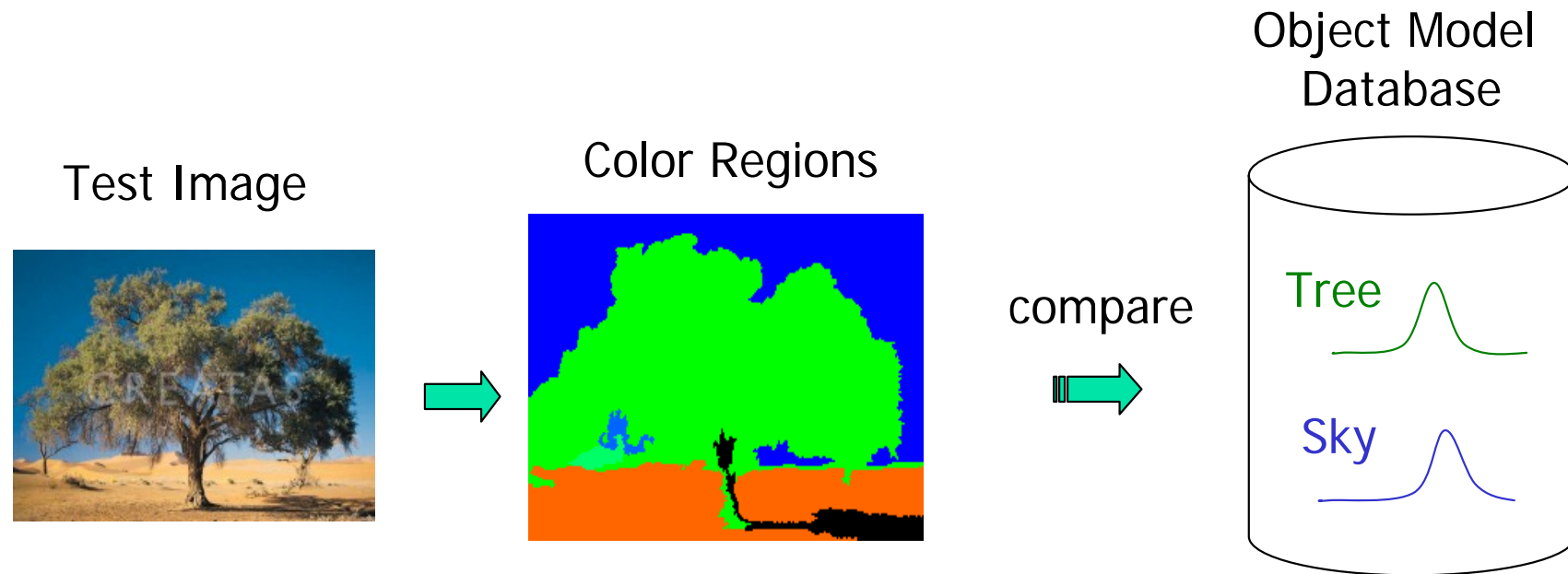
Image & description



## 2. Iteration Step (Example)



# Recognition



To calculate  $p(\text{tree} / \text{image})$

$$p(\text{tree} / \text{image}) = f \left( \begin{array}{l} p(\text{tree} / \text{blue}) \\ p(\text{tree} / \text{green}) \\ p(\text{tree} / \text{orange}) \\ p(\text{tree} / \text{black}) \end{array} \right)$$

$$p(o | F_I^a) = \prod_{r^a \in F_I^a} (p(o | r^a))$$



# Combining different abstract regions

---

- Treat the different types of regions **independently** and combine at the time of classification.

$$p(o | \{F_I^a\}) = \prod_a p(o | F_I^a)$$

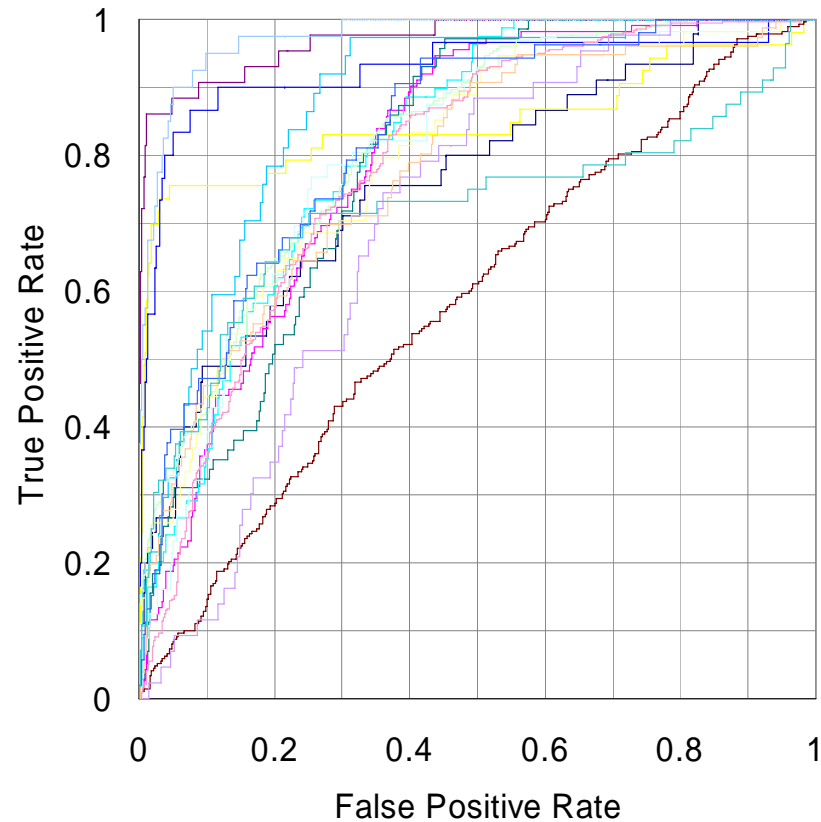
- Form **intersections** of the different types of regions, creating smaller regions that have both color and texture properties for classification.

# Experiments (on 860 images)

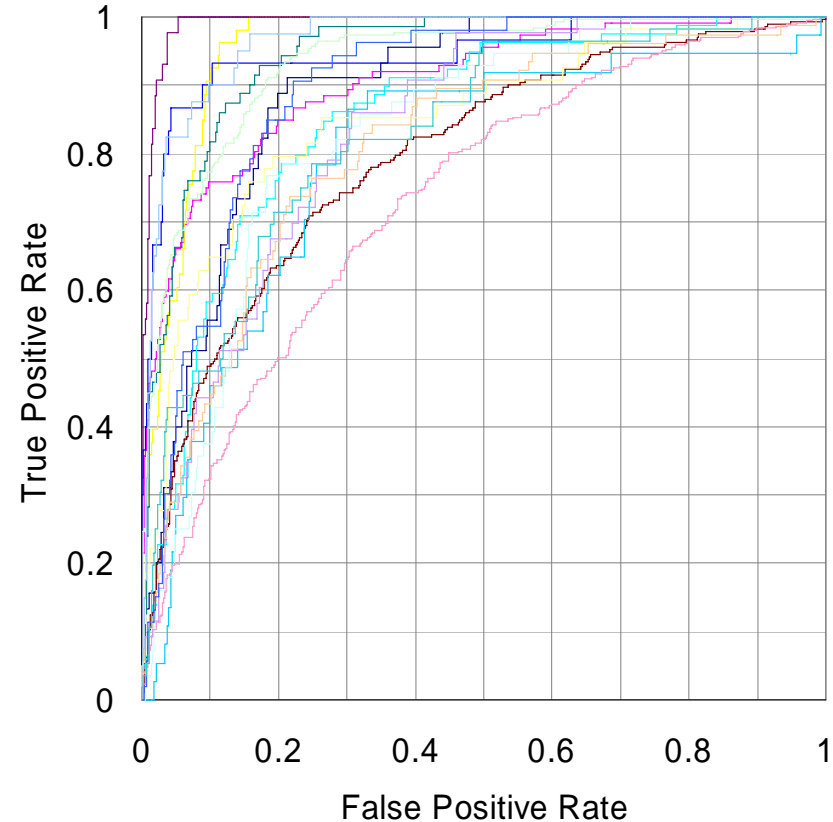
---

- 18 keywords: mountains (30), orangutan (37), track (40), tree trunk (43), football field (43), beach (45), prairie grass (53), cherry tree (53), snow (54), zebra (56), polar bear (56), lion (71), water (76), chimpanzee (79), cheetah (112), sky (259), grass (272), tree (361).
- A set of cross-validation experiments (80% as training set and the other 20% as test set)
- The poorest results are on object classes “tree,” “grass,” and “water,” each of which has a high variance; a single Gaussian model is insufficient.

# ROC Charts



Independent Treatment of  
Color and Texture

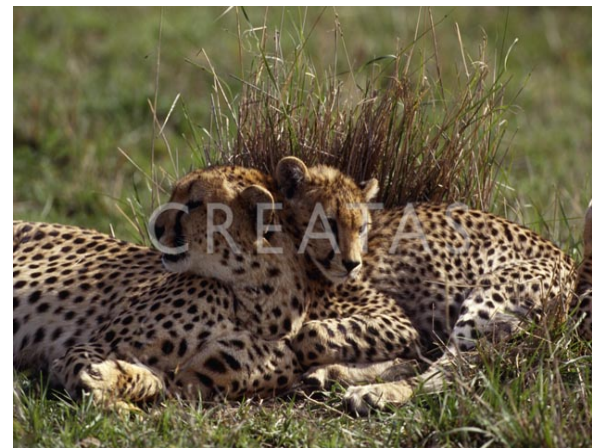
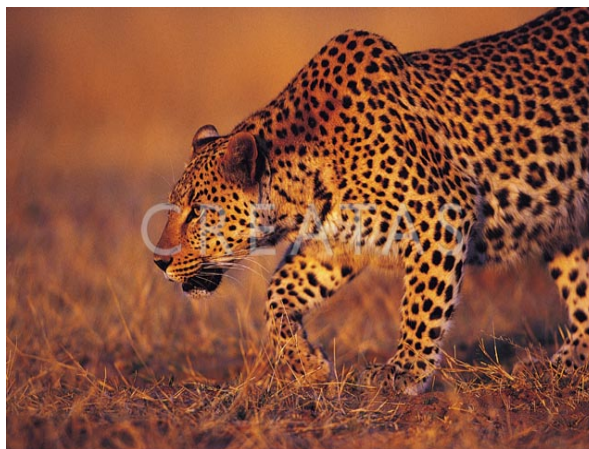


Using Intersections of  
Color and Texture Regions

# Sample Results

---

cheetah





# Sample Results (Cont.)

---

grass





# Sample Results (Cont.)

---

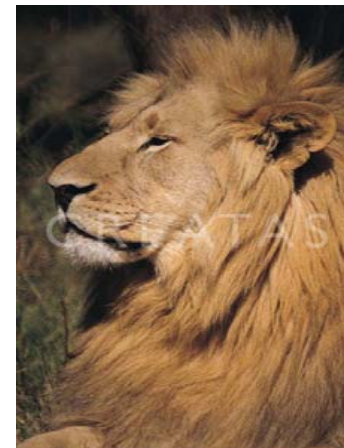
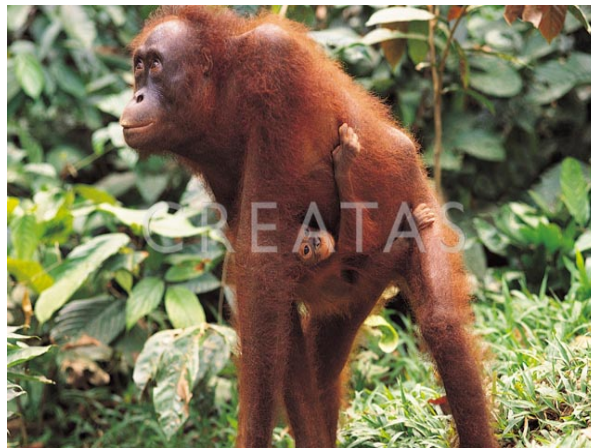
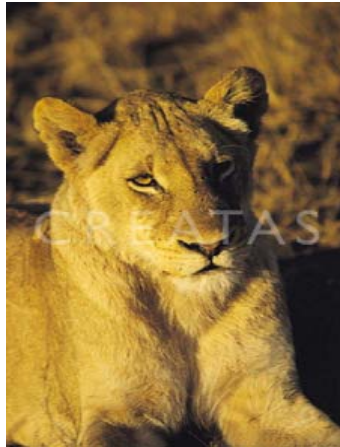
cherry tree



# Sample Results (Cont.)

---

lion





# Summary

---

- Designed a set of abstract region features: color, texture, structure, . . .
- Developed a new semi-supervised EM-like algorithm to recognize object classes in color photographic images of outdoor scenes; tested on 860 images.
- Compared two different methods of combining different types of abstract regions. The intersection method had a higher performance

# Groundtruth Data Set

---

- UW Ground truth database (1224 images)
- 31 elementary object categories: *river* (30), *beach* (31), *bridge* (33), *track* (35), *pole* (38), *football field* (41), *frozen lake* (42), *lantern* (42), *husky stadium* (44), *hill* (49), *cherry tree* (54), *car* (60), *boat* (67), *stone* (70), *ground* (81), *flower* (85), *lake* (86), *sidewalk* (88), *street* (96), *snow* (98), *cloud* (119), *rock* (122), *house* (175), *bush* (178), *mountain* (231), *water* (290), *building* (316), *grass* (322), *people* (344), *tree* (589), *sky* (659)
- 20 high-level concepts: *Asian city*, *Australia*, *Barcelona*, *campus*, *Cannon Beach*, *Columbia Gorge*, *European city*, *Geneva*, *Green Lake*, *Greenland*, *Indonesia*, *indoor*, *Iran*, *Italy*, *Japan*, *park*, *San Juans*, *spring flowers*, *Swiss mountains*, and *Yellowstone*.





*beach, sky, tree, water*



*people, street, tree*



*building, grass, people,  
sidewalk, sky, tree*



*building, bush, sky,  
tree, water*



*flower, house, people,  
pole, sidewalk, sky*



*flower, grass, house,  
pole, sky, street, tree*



*building, flower, sky,  
tree, water*



*boat, rock, sky,  
tree, water*



*building, car, people, tree*



*car, people, sky*



*boat, house, water*



*building* 82

# Groundtruth Data Set:

## ROC Scores

---

<i>street</i>	60.4	<i>tree</i>	80.8	<i>stone</i>	87.1	<i>columbia gorge</i>	94.5
<i>people</i>	68.0	<i>bush</i>	81.0	<i>hill</i>	87.4	<i>green lake</i>	94.9
<i>rock</i>	73.5	<i>flower</i>	81.1	<i>mountain</i>	88.3	<i>italy</i>	95.1
<i>sky</i>	74.1	<i>iran</i>	82.2	<i>beach</i>	89.0	<i>swiss moutains</i>	95.7
<i>ground</i>	74.3	<i>bridge</i>	82.7	<i>snow</i>	92.0	<i>sanjuans</i>	96.5
<i>river</i>	74.7	<i>car</i>	82.9	<i>lake</i>	92.8	<i>cherry tree</i>	96.9
<i>grass</i>	74.9	<i>pole</i>	83.3	<i>frozen lake</i>	92.8	<i>indoor</i>	97.0
<i>building</i>	75.4	<i>yellowstone</i>	83.7	<i>japan</i>	92.9	<i>greenland</i>	98.7
<i>cloud</i>	75.4	<i>water</i>	83.9	<i>campus</i>	92.9	<i>cannon beach</i>	99.2
<i>boat</i>	76.8	<i>indonesia</i>	84.3	<i>barcelona</i>	92.9	<i>track</i>	99.6
<i>lantern</i>	78.1	<i>sidewalk</i>	85.7	<i>geneva</i>	93.3	<i>football field</i>	99.8
<i>australia</i>	79.7	<i>asian city</i>	86.7	<i>park</i>	94.0	<i>husky stadium</i>	100.0
<i>house</i>	80.1	<i>european city</i>	87.0	<i>spring flowers</i>	94.4		



# Groundtruth Data Set: Top Results

*Asian city*



*Cannon beach*



*Italy*



*park*



# Groundtruth Data Set: Top Results

---

*sky*



*spring flowers*



*tree*



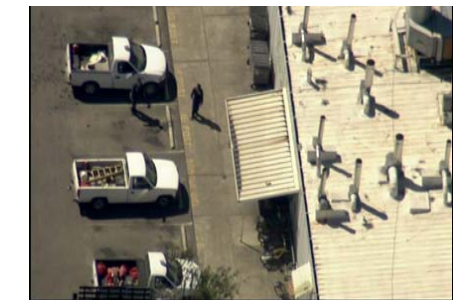
*water*



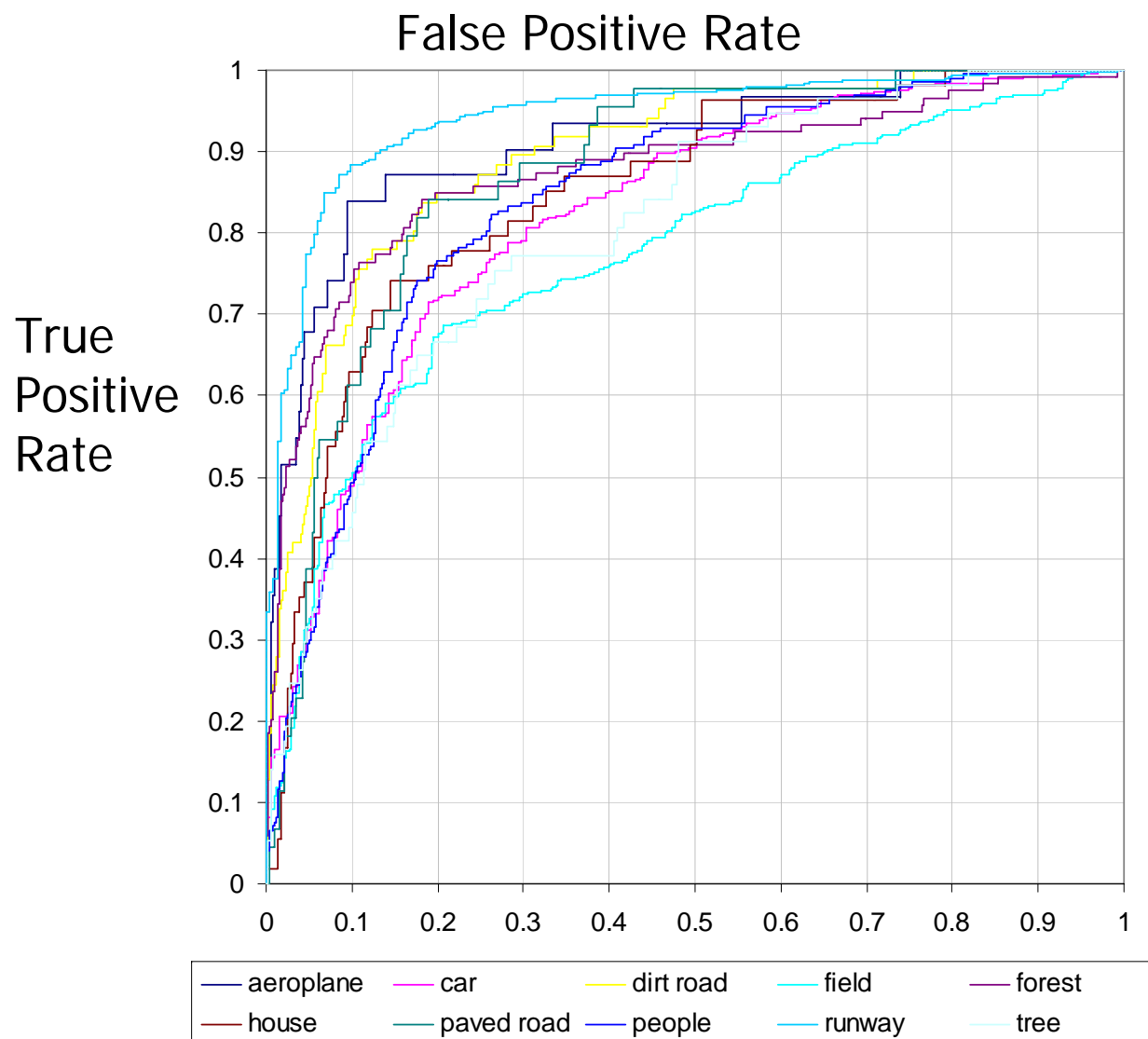


# VACE Test Image Set (828 images and 10 object classes): from Boeing, VIVID, and NGA videos

---



# Experiments: ROC Curves



field	77.5
tree	80.6
car	82.3
people	83.9
house	84.9
paved road	87.5
forest	87.6
dirt road	89.5
airplane	91.1
runway	94.4

# Objects detected in frames



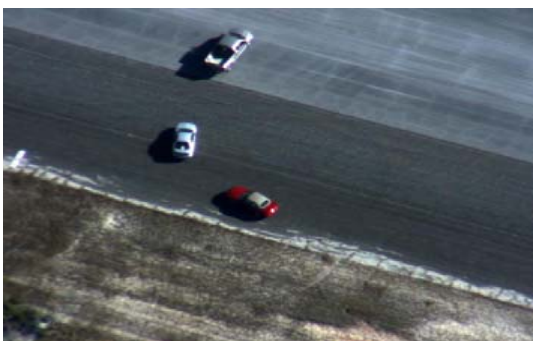
**forest(94.37) house(64.09)**  
car(46.5) dirt road(23.44) paved  
road(4.77) tree(2.29) airplane(1.47)  
runway(0.03) field(0.02) people(0)



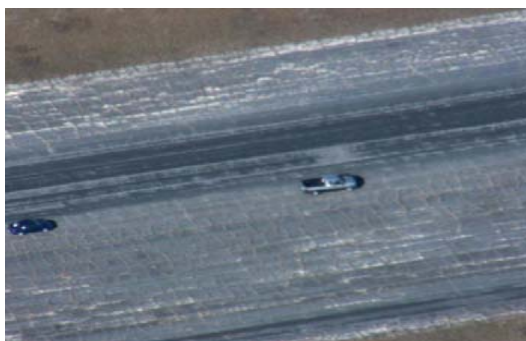
**runway(99.98) field(98.66) car(96.24)**  
people(10.04) airplane(2.74) paved  
road(2.39) forest(0.82) house(0.48) dirt  
road(0.41) tree(0)



**car(94.3) dirt road(91.7) field(16.17)**  
tree(14.23) paved road(5.34) airplane(5.17)  
people(3.91) forest(0.53) house(0.47)  
runway(0.41)



**runway(100) car(99.23) field(98.07)** dirt  
road(92.1) house(85.24) tree(19.43)  
paved road(5.77) airplane(3.56)  
forest(2.85) people(0.07)



**runway(99.98) car(99.84) field(99.27)**  
paved road(18.28) people(13.13)  
tree(8.71) airplane(7.94) forest(1.67)  
house(0.14) dirt road(0.08)



**car(97.92) forest(94.2) paved road(85)**  
**dirt road(72.94)** tree(68.84)  
airplane(39.13) house(33.17)  
people(12.97) field(2.38) runway(0.04)