# Recap: Viola/Jones detector

- Rectangle features

- Integral images for fast computation

- Boosting for feature selection

- Attentional cascade for fast rejection of negative windows

# Project 3

- I have office hours today
- Let's talk more about scene recognition

# SUN Database: Large-scale Scene Categorization and Detection

Jianxiong Xiao, James Hays[†], Krista A. Ehinger, Aude Oliva, Antonio Torralba

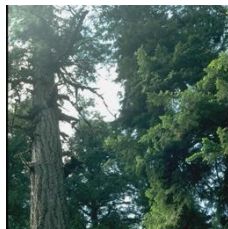Massachusetts Institute of Technology
[†] Brown University

# Scene Categorization

## Oliva and Torralba, 2001



| Coast | Forest | Highway | Inside City | Mountain | Open Country | Street | Tall Building |

## Fei Fei and Perona, 2005

+



| Bedroom | Kitchen | Living Room | Office | Suburb |

## Lazebnik, Schmid, and Ponce, 2006

+



| Industrial | Store |

# 15 Scene Database

# How many object categories are there?



~10,000 to 30,000

Biederman 1987

abbey

airplane cabin

airport terminal

**apple orchard**

assembly hall

bakery

car factory

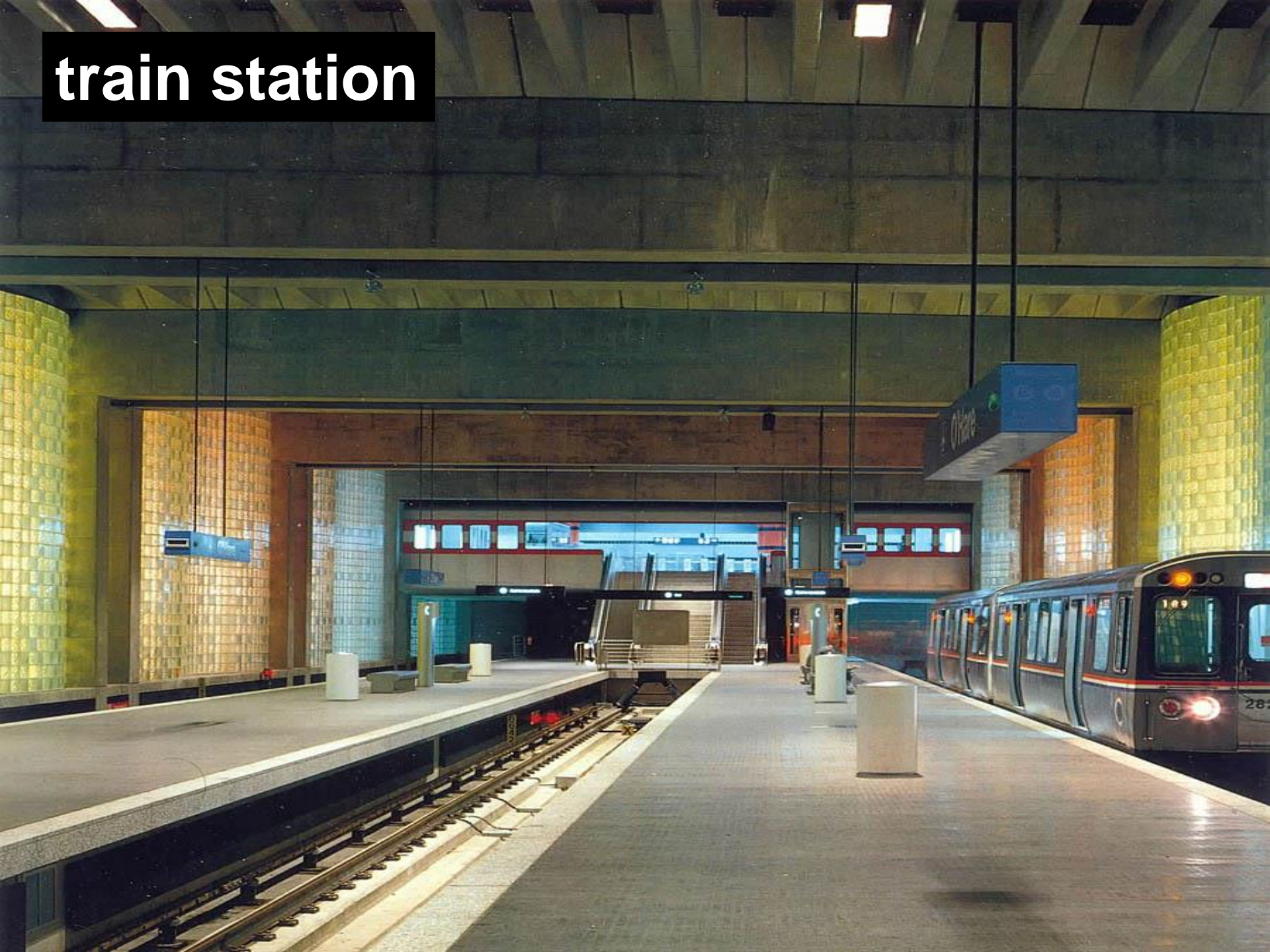cockpit
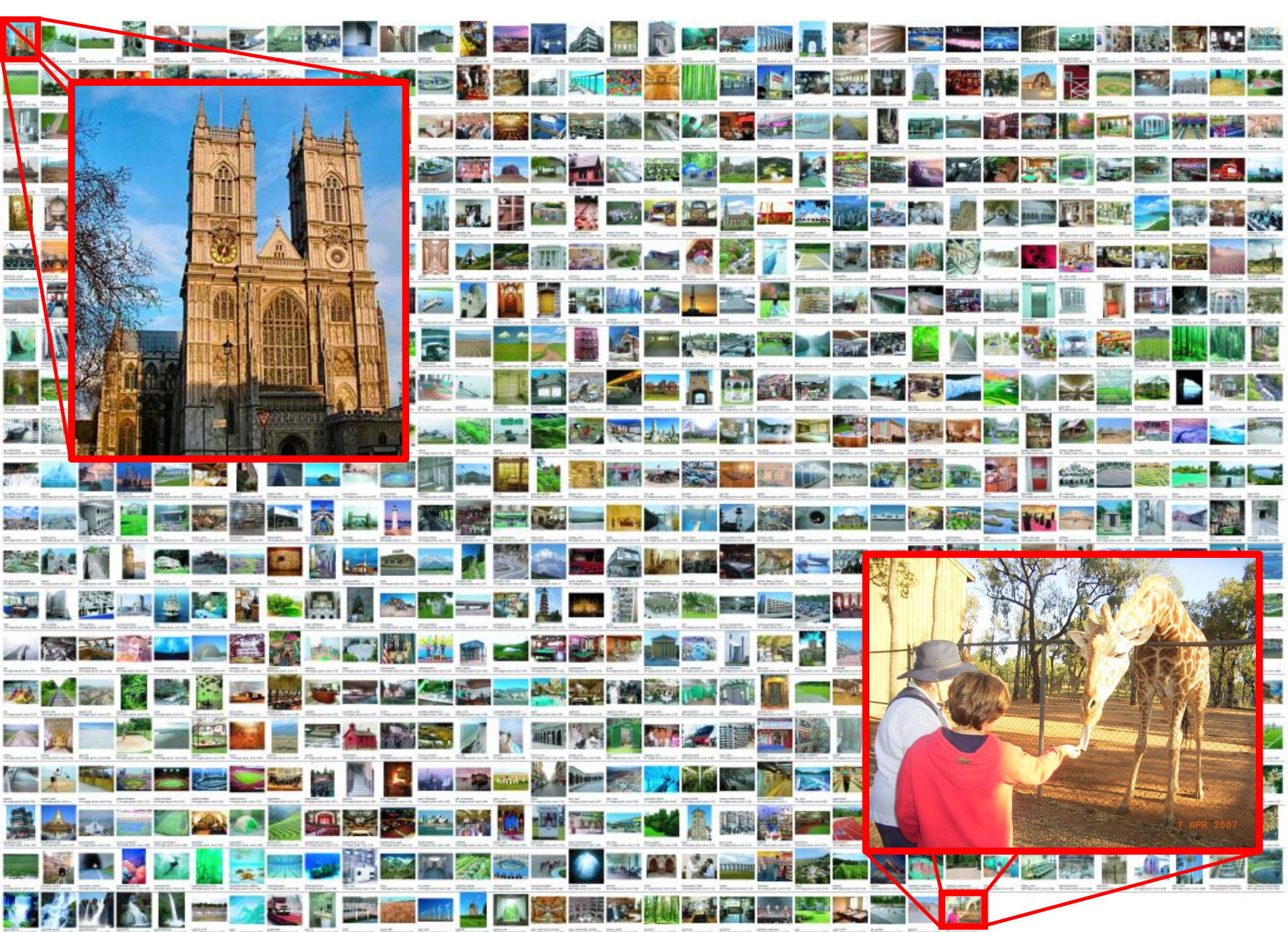
construction site

food court

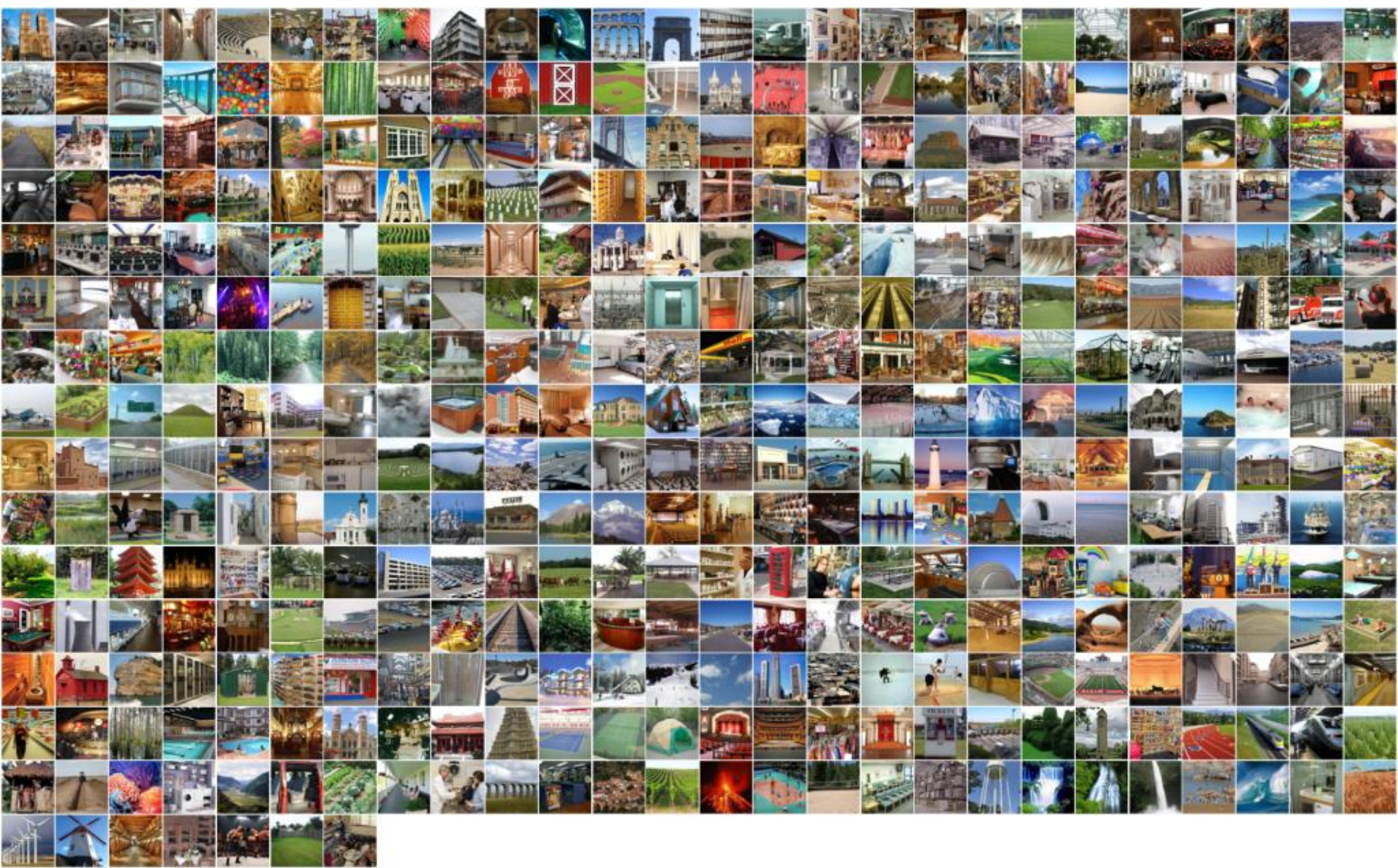interior car

lounge

stream

train station

# 397 Well-sampled Categories

# Evaluating Human Scene Classification



?

"Good worker" Accuracy     98%     90%     68%

bathroom(100%)

beauty salon(100%)

bedroom(100%)

bullring(100%)

playground(100%)

phone booth(100%)

greenhouse outdoor(100%)

podium outdoor(100%)

tennis court outdoor(100%)

wind farm(100%)

veterinarians office(100%)

riding arena(100%)

## Scene category

Inn  (0%)



Bayou  (0%)



Basilica  (0%)



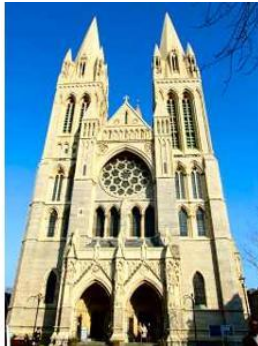## Most confusing categories

Restaurant patio (44%)



Chalet (19%)



River (67%)



Coast (8%)



Cathedral(29%)



Courthouse (21%)

# Conclusion: humans can do it

- The SUN database is reasonably consistent and differentiable -- even with a huge number of very specific categories, humans get it right 2/3rds of the time *with no training.*

- We also have a good benchmark for computational methods.

## How do we classify scenes?
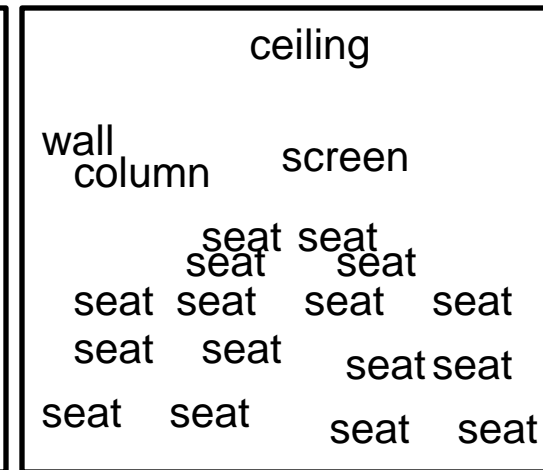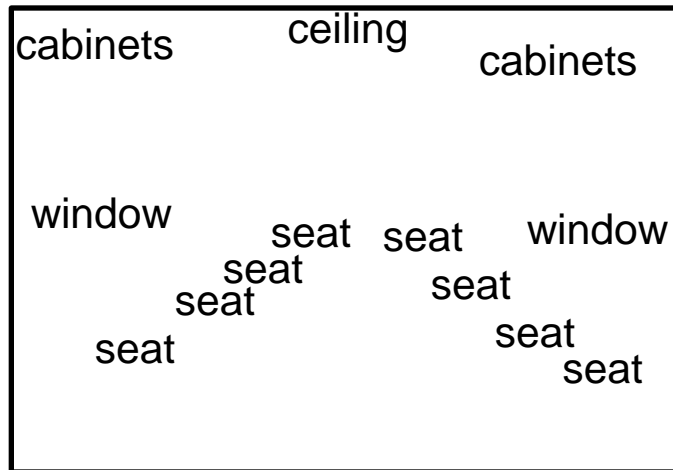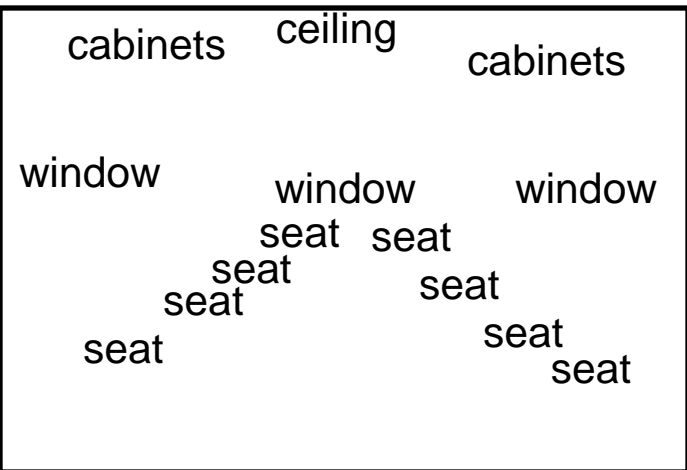
# How do we classify scenes?



| | | |
|---|---|---|
| Ceiling | Ceiling | |
| Light | Lamp | wall |
| | | painting |
| Door  Door | Painting  mirror | |
| Door | mirror | wall    Lamp |
| Wall   Door    Wall   Door | wall | |
| | | phone |
| | Fireplace | Bed   alarm |
| Floor | armchair    armchair | |
| | | Side-table |
| | Coffee table | carpet |

Different objects, different spatial layout

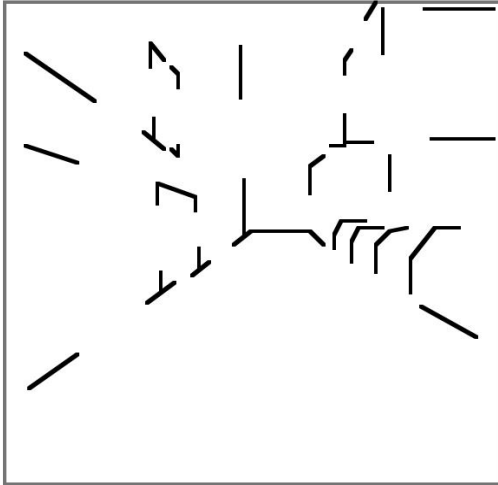# Which are the important elements?



Similar objects, and similar spatial layout

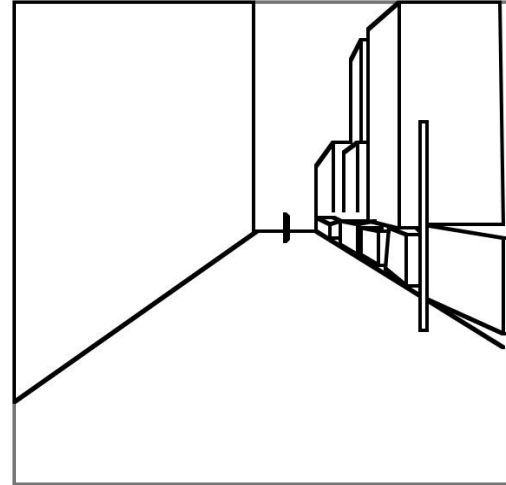Different lighting, different materials, different "stuff"

# Scene emergent features

"Recognition via features that are not those of individual objects but "emerge" as objects are brought into relation to each other to form a scene." – Biederman 81



Biederman, 1981

Suggestive edges and junctions



Biederman, 1981

Simple geometric forms



Bruner and Potter, 1969

Blobs



Oliva and Torralba, 2001

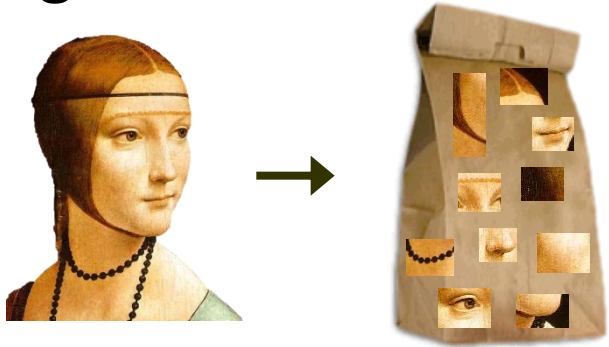Textures

# Global Image Descriptors

- Tiny images (Torralba et al, 2008)
- Color histograms
- Self-similarity (Shechtman and Irani, 2007)
- Geometric class layout (Hoiem et al, 2005)
- Geometry-specific histograms (Lalonde et al, 2007)
- Dense and Sparse SIFT histograms
- Berkeley texton histograms (Martin et al, 2001)
- HoG 2x2 spatial pyramids
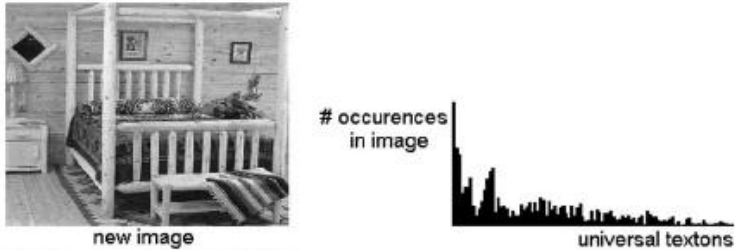- Gist scene descriptor (Oliva and Torralba, 2008)

Texture Features

# Global Texture Descriptors

## Bag of words



Sivic et. al., ICCV 2005
Fei-Fei and Perona, CVPR 2005
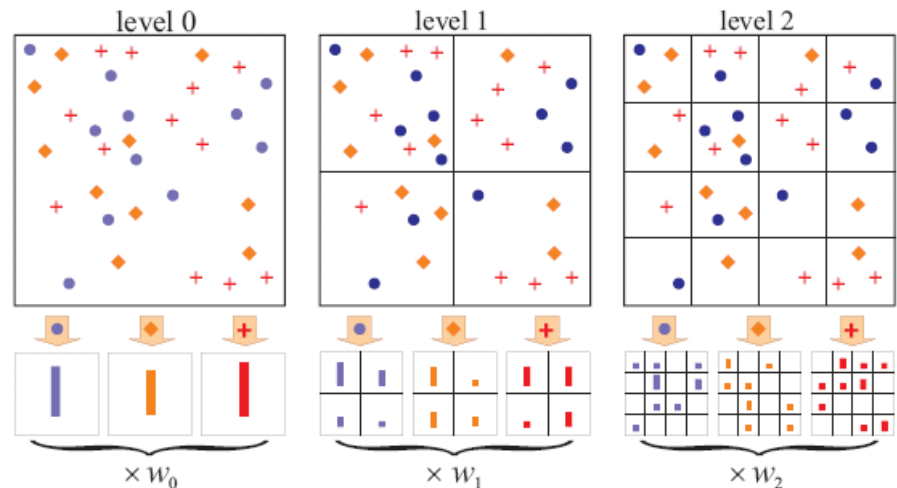
## Non localized textons



Walker, Malik. Vision Research 2004

...

## Spatially organized textures



M. Gorkani, R. Picard, ICPR 1994
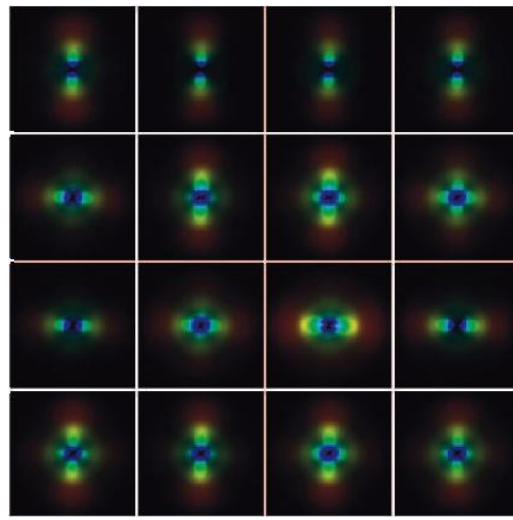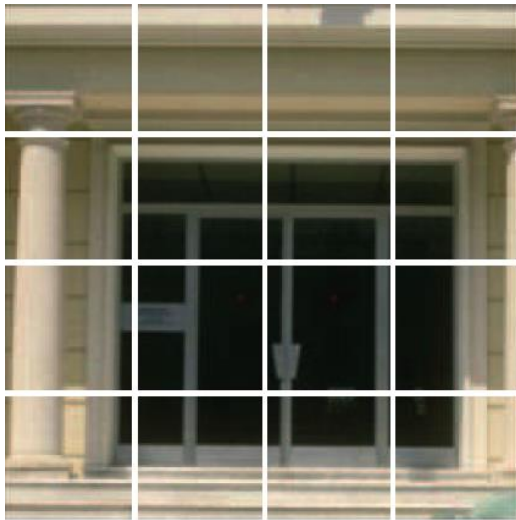A. Oliva, A. Torralba, IJCV 2001



level 0        level 1        level 2

$\times w_0$        $\times w_1$        $\times w_2$

S. Lazebnik, et al, CVPR 2006        ...

R. Datta, D. Joshi, J. Li, and J. Z. Wang, **Image Retrieval: Ideas, Influences, and Trends of the New Age**, *ACM Computing Surveys*, vol. 40, no. 2, pp. 5:1-60, 2008.
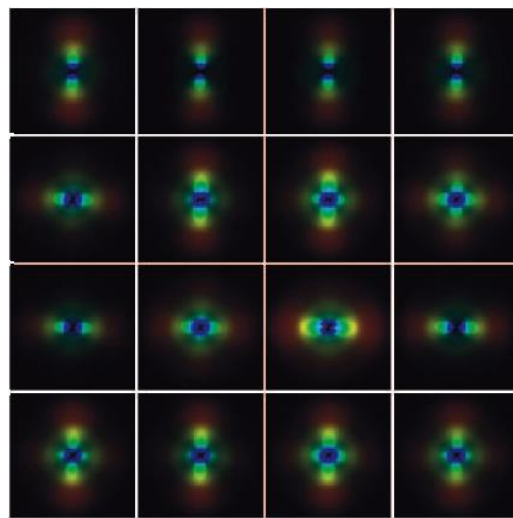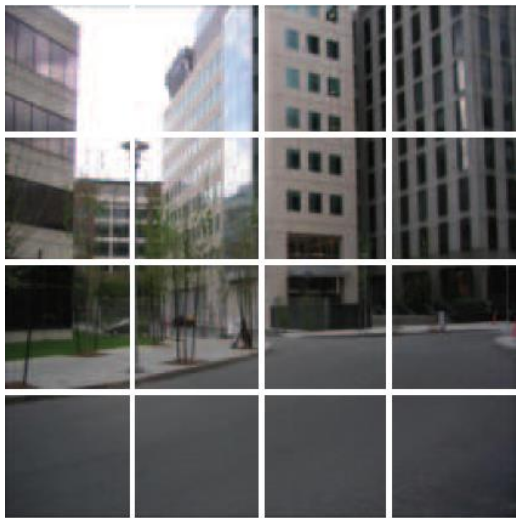
# Gist descriptor

Oliva and Torralba, 2001



- Apply oriented Gabor filters over different scales
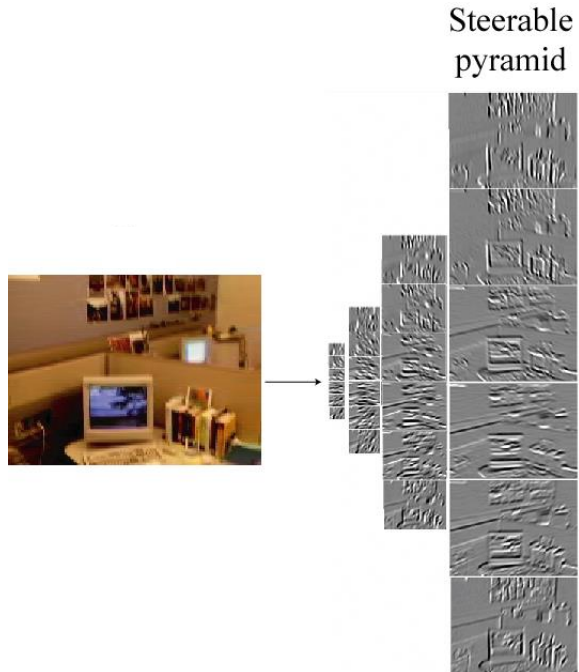- Average filter energy in each bin

8   orientations
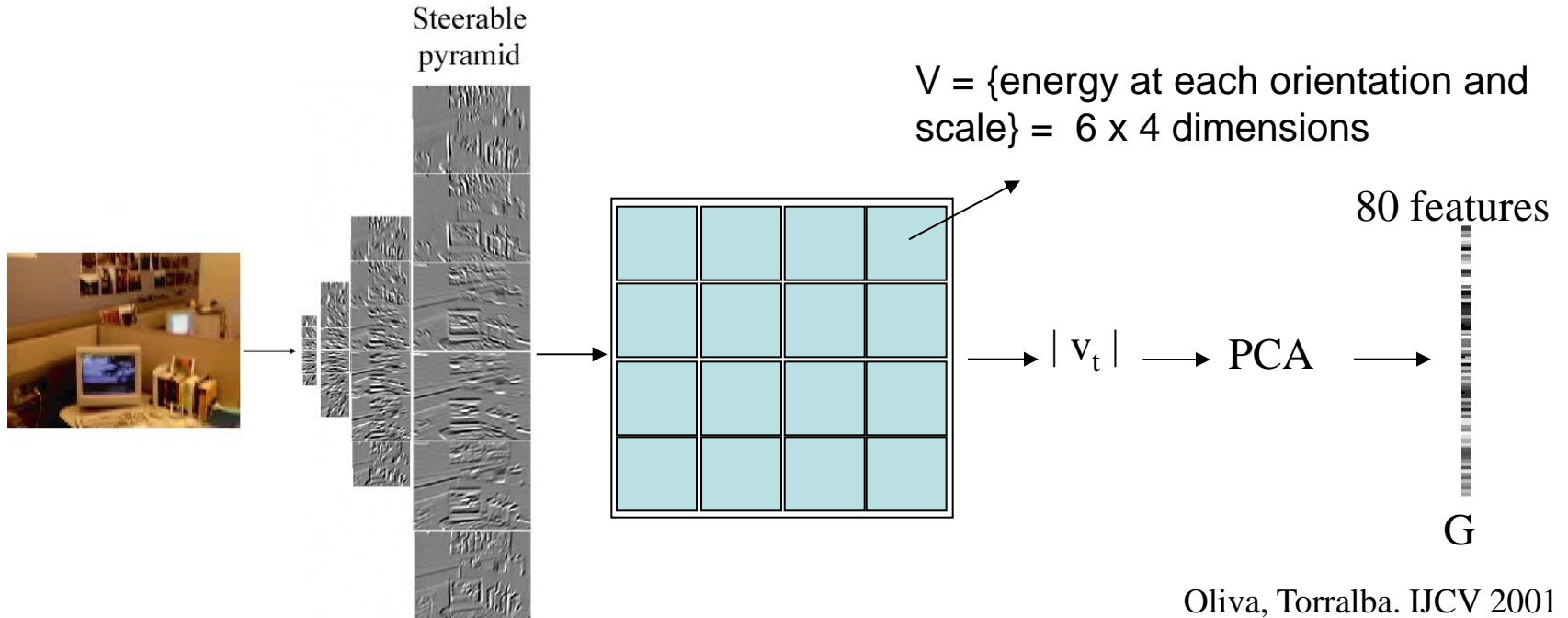4   scales
<u>x 16</u>   bins
512   dimensions

Similar to SIFT (Lowe 1999) applied to the entire image

M. Gorkani, R. Picard, ICPR 1994; Walker, Malik. Vision Research 2004;  Vogel et al. 2004;
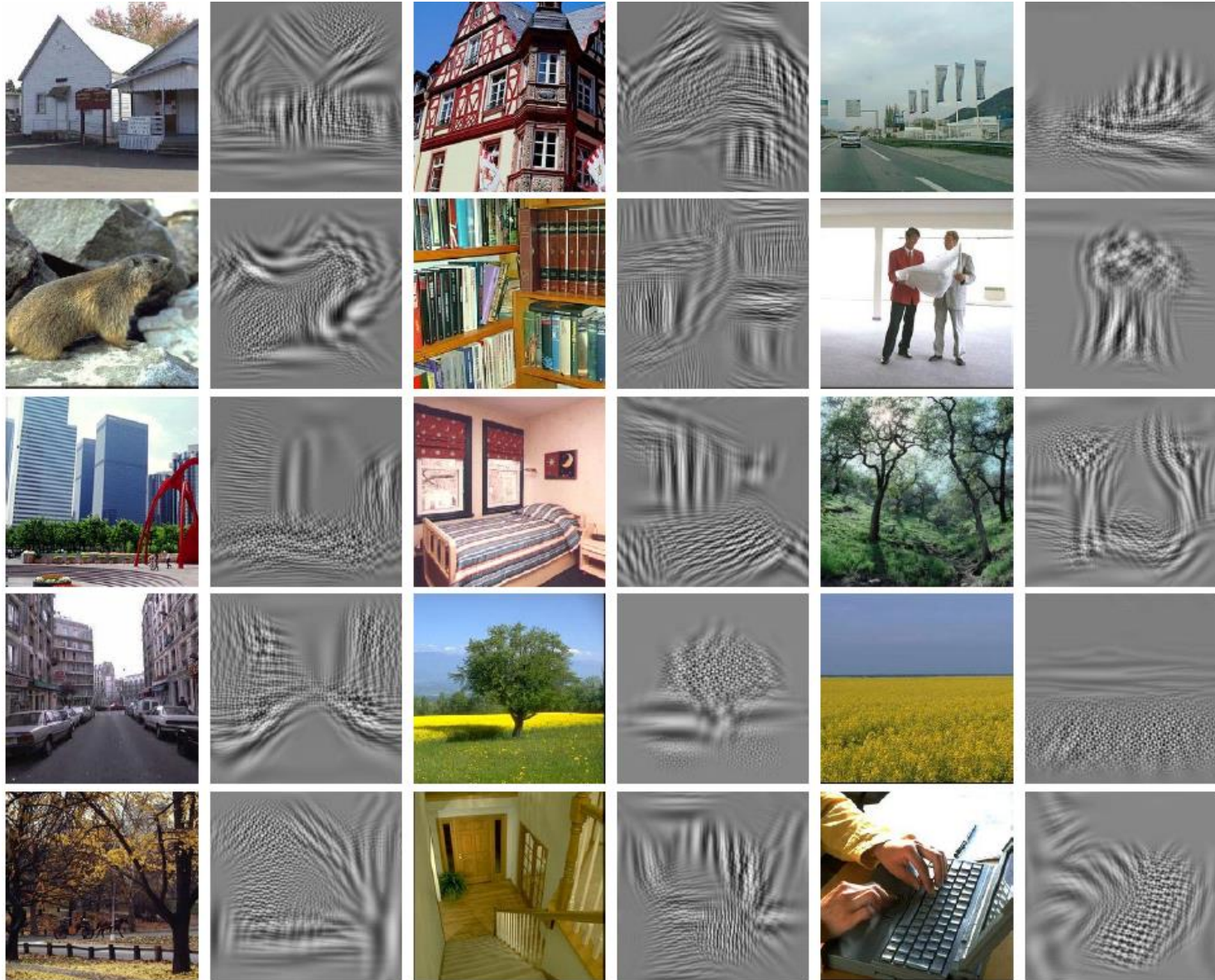Fei-Fei and Perona, CVPR 2005; S. Lazebnik, et al, CVPR 2006; …
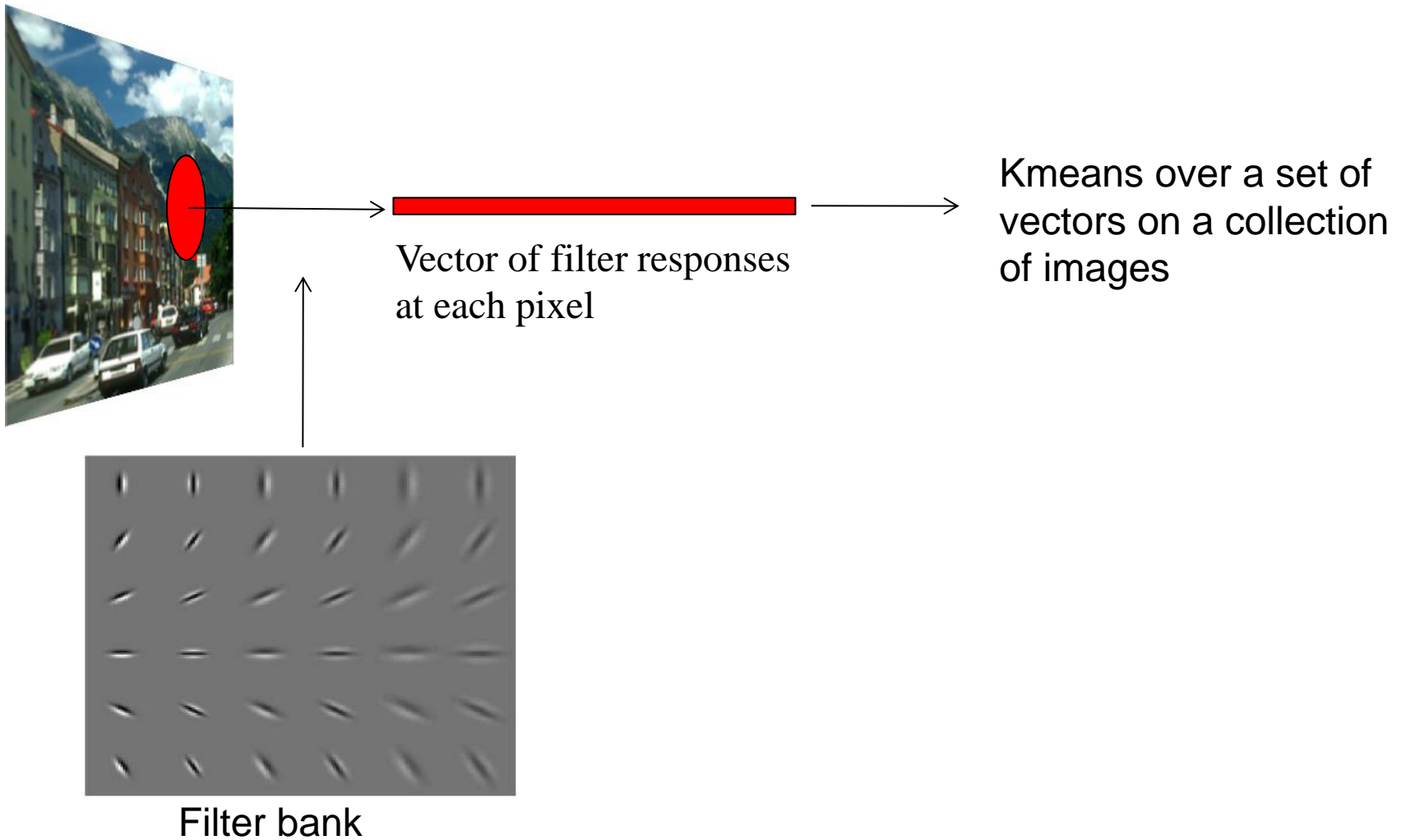
# Gist descriptor



Steerable
pyramid

# Gist descriptor



Steerable pyramid

V = {energy at each orientation and scale} = 6 x 4 dimensions

80 features

$$|v_t| \longrightarrow PCA \longrightarrow$$

G

Oliva, Torralba. IJCV 2001

# Example visual gists



Global features (I) ~ global features (I')

Oliva & Torralba (2001)

# Textons



Vector of filter responses
at each pixel

Kmeans over a set of
vectors on a collection
of images

Filter bank
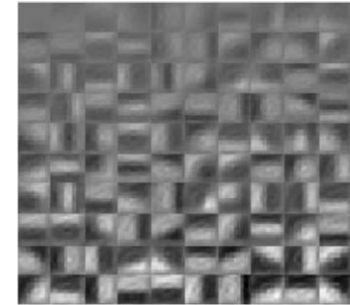
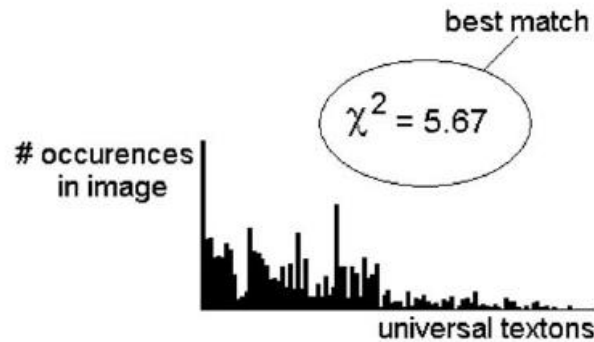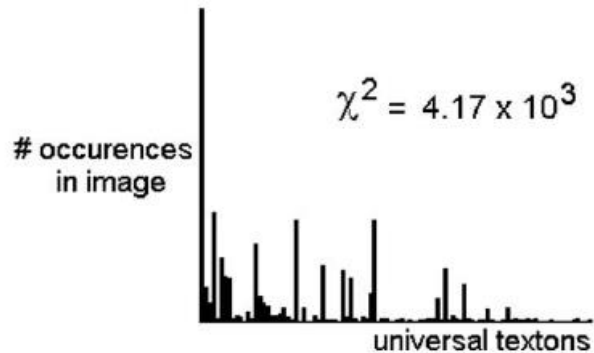Malik, Belongie, Shi, Leung, 1999

# Textons



Filter bank

K-means (100 clusters)

Malik, Belongie, Shi, Leung, 1999

label = bedroom

best match

$\chi^2 = 5.67$

\# occurences in image

universal textons

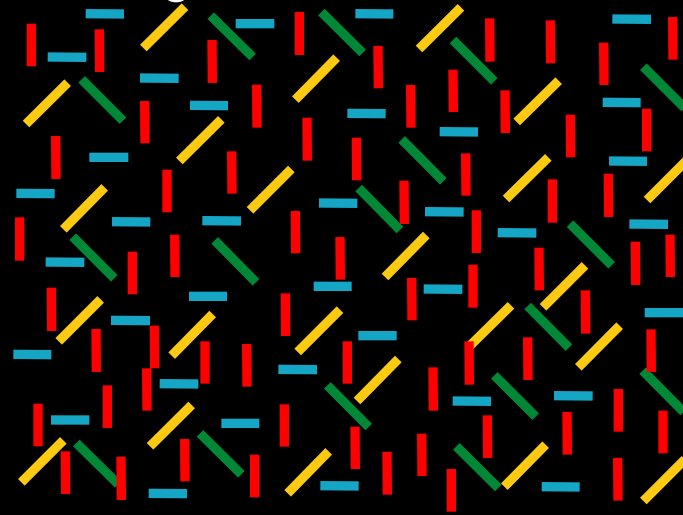label = beach

$\chi^2 = 4.17 \times 10^3$

\# occurences in image

universal textons

Walker, Malik, 2004

# Bag of words

## Bag of words model



65 17 23 36

## Spatially organized textures



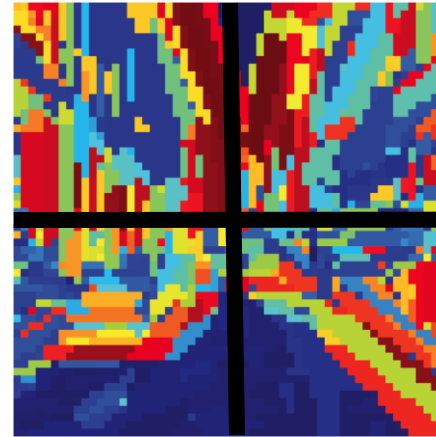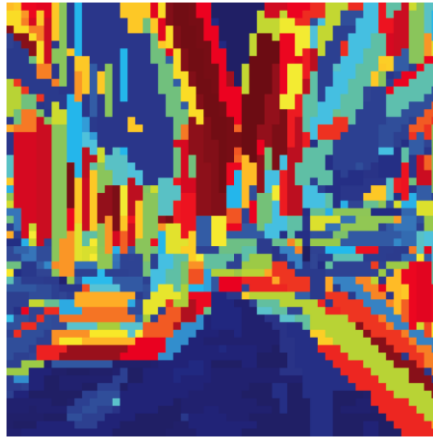| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 8 | 0 | 0 | 0 | 2 | 0 | 0 | 7 | 0 | 4 | 0 |
| 20 | 0 | 0 | 0 | 11 | 1 | 0 | 2 | 14 | 0 | 3 | 3 |
| 3 | 0 | 12 | 4 | 0 | 0 | 4 | 16 | 3 | 6 | 0 | 11 |

# Bag of words & spatial pyramid matching

Sivic, Zisserman, 2003. Visual words = Kmeans of SIFT descriptors



S. Lazebnik, et al, CVPR 2006
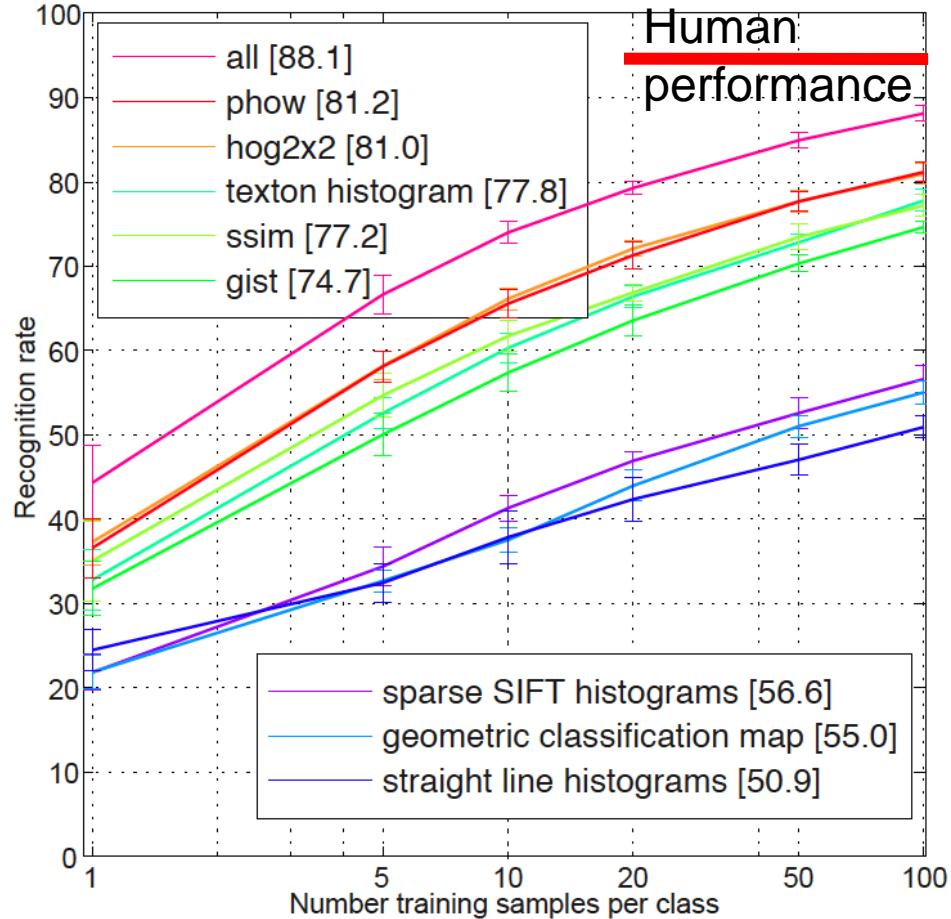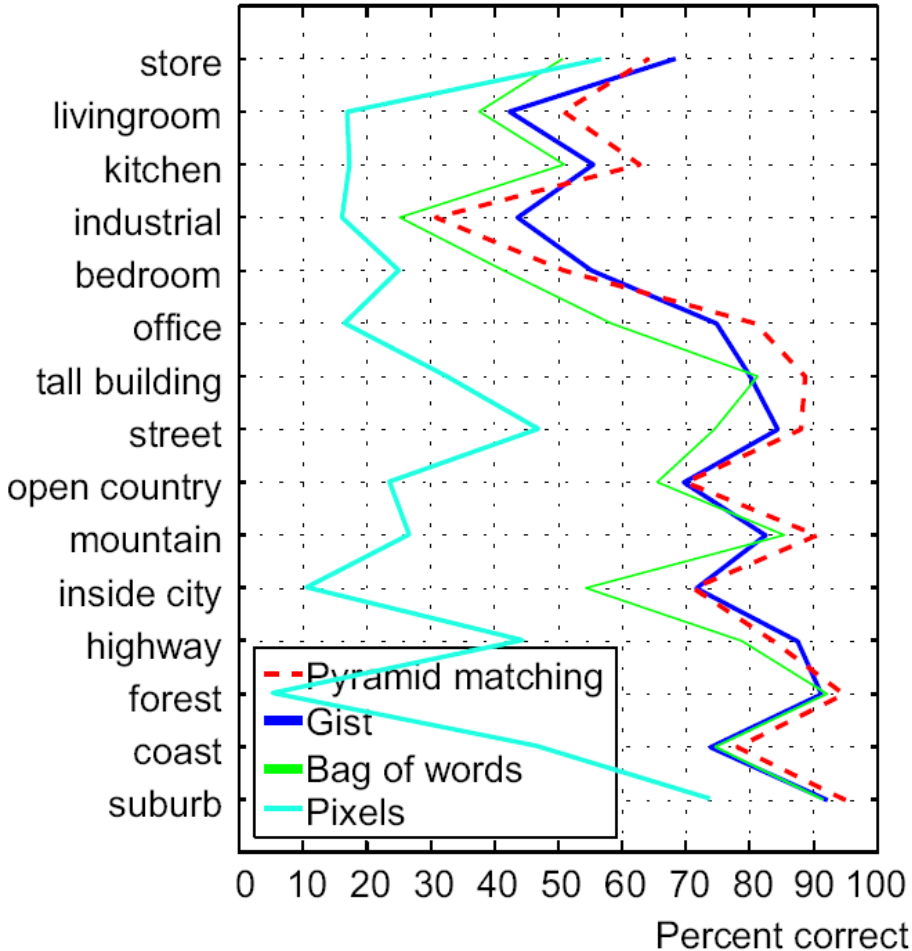
# Learning Scene Categorization



Forest path
Vs.
all

Living - room
Vs.
all

# Scene recognition

100 training samples per class

SVM classifier in both cases

# Feature Accuracy

Legend:
- all [38.0]
- hog2x2 [27.2]
- geometry texton histograms [23.5]
- ssim [22.5]
- dense SIFT [21.5]
- lbp [18.0]
- texton histogram [17.6]
- gist [16.3]
- all (1NN) [13.0]
- lbphf [12.8]
- sparse SIFT histograms [11.5]
- geometry color histograms [9.1]
- color histograms [8.2]
- geometric classification map [6.0]
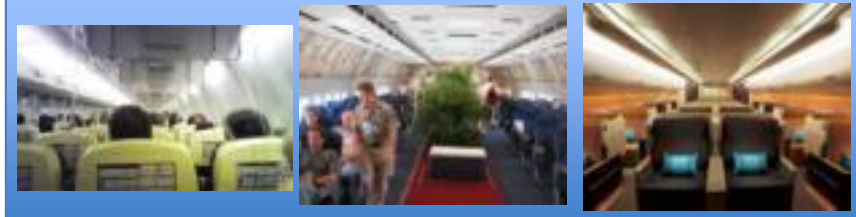- straight line histograms [5.7]
- tiny image [5.5]

Classifier: 1-vs-all SVM with histogram intersection, chi squared, or RBF kernel.

# A look into the results

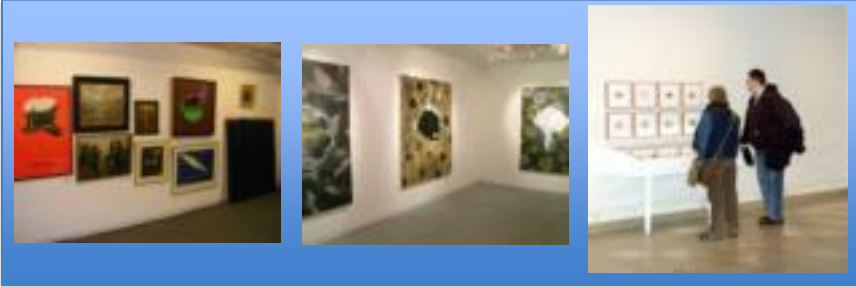## Airplane cabin (64%)



Van interior    Discotheque    Toyshop

## Art gallery (38%)



Iceberg    Hotel room    Kitchenette

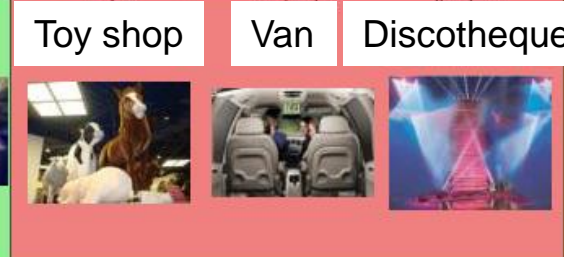All the results available on the web                ...

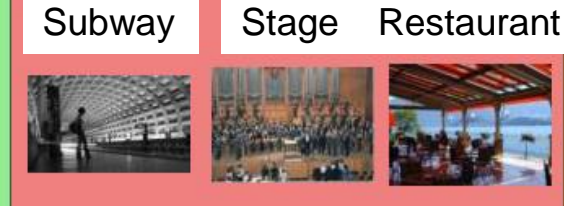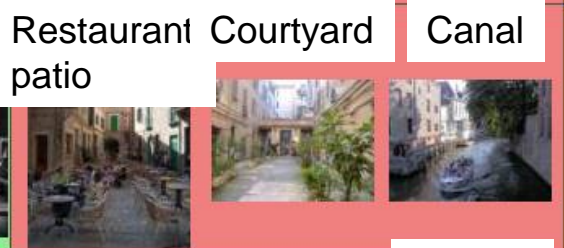|  | Training images | Correct classifications | Miss-classifications |
|---|---|---|---|

**Abbey** — Miss-classifications: Monastery, Cathedral, Castle

**Airplane cabin** — Miss-classifications: Toy shop, Van, Discotheque

**Airport terminal** — Miss-classifications: Subway, Stage, Restaurant

**Alley** — Miss-classifications: Restaurant patio, Courtyard, Canal

**Amphitheater** — Miss-classifications: Harbor, Coast, Athletic field

Xiao, Hays, Ehinger, Oliva, Torralba; maybe 2010

|  | limousine interior (95% vs 80%) | riding arena (100% vs 90%) | sauna (96% vs 95%) | skatepark (96% vs 90%) | subway interior (96% vs 80%) |

**Humans good Comp. good**
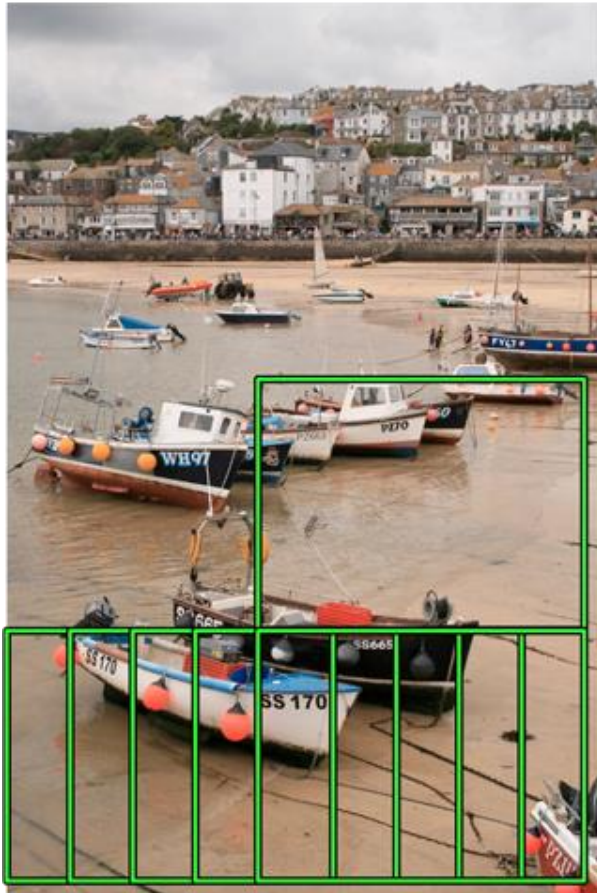
**Humans bad Comp. bad**

**Human good Comp. bad**

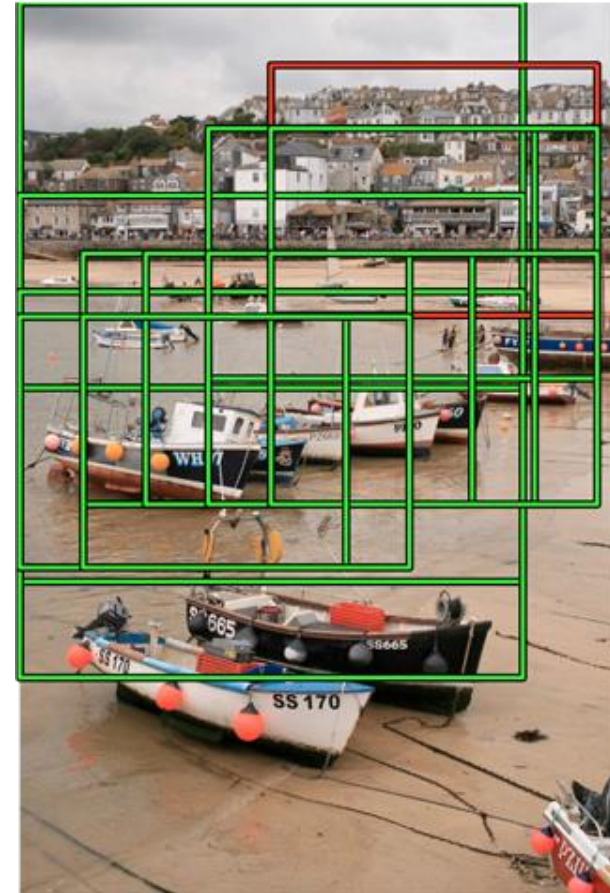**Human bad Comp. good**

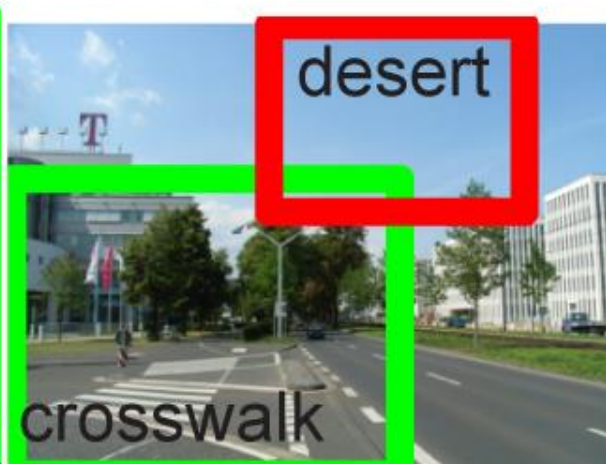# Local Scene Detection



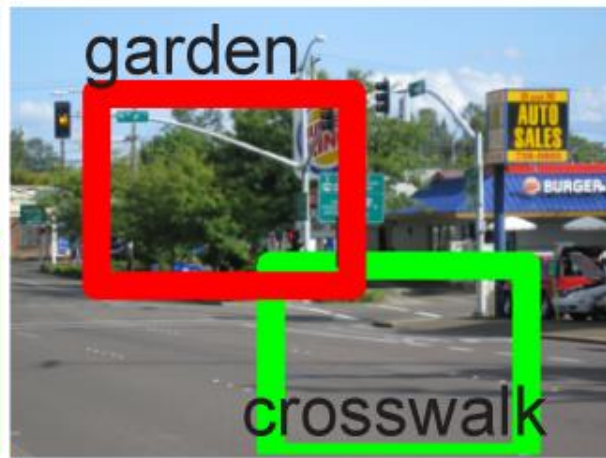beach detections     village detections     harbor detections

# Confident Subscene Detections

Database and source code available at
  http://groups.csail.mit.edu/vision/SUN/

Additional details available:

**SUN Database: Large-scale Scene Recognition from Abbey to Zoo.** Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, Antonio Torralba. *CVPR 2010.*

# How do we do better than 40%?

- Deep learning gets about the same performance

- Fisher vector encoding gets up to 47.2%