## **Project presentations**

Dec 6 and Dec 11

Dec 13 reports due

Assigned days in the website.

Sent us email if there is a mistake in the title/group members or a time conflict.

## Presentations

4 Min + 1 min questions

• Send us presentation. We will run all presentations from the same computer.

## How to give a talk

http://www.cs.berkeley.edu/~messer/Bad\_talk.html

http://www-psych.stanford.edu/~lera/talk.html

## First, some bad news

The more you work on a talk, the better it gets: if you work on it for 1 day, the talk you give will be better than if you had only worked on it for 1 hour. If you work on it for 2 days, it will be better still. 7 days, better yet...

## All talks are important

There are no unimportant talks. There are no big or small audiences.

Prepare each talk with the same enthusiasm.

## How to give a talk

#### **Delivering**:

- Look at the audience! Try not to talk to your laptop or to the screen. Instead, look at the other humans in the room.
- You have to believe in what you present, be confident... even if it only lasts for the time of your presentation.
- Do not be afraid to acknowledge limitations of whatever you are presenting. Limitations are good. They leave job for the people to come. Trying to hide the problems in your work will make the preparation of the talk a lot harder and your self confidence will be hurt.

# Let the audience see your personality

- They want to see you enjoy yourself.
- They want to see what you love about the work.
- People really respond to the human parts of a talk. Those parts help the audience with their difficult task of listening to an hour-long talk on a technical subject. What was easy, what was fun, what was hard about the work?
- Don't be afraid to be yourself and to be quirky.

The different kinds of talks you'll have to give as a researcher

- 2-5 minute talks
- 20 30 minute conference presentations
- 30-60 minute colloquia

## How to give a talk

**Talk organization:** here there are as many theories as there are talks. Here there are some extreme advices:

- 1. Go into details / only big picture
- 2. Go in depth on a single topic / cover as many things as you can
- 3. Be serious (never make jokes, maybe only one) / be funny (it is just another form of theater)
- Corollary: ask people for advice, but at the end, if will be just you and the audience. Chose what fits best your style.
- What everybody agree on is that you have to practice in advance (the less your experience, the more you have to practice). Do it with an audience or without, but practice.

The best advice I got came from Yair Weiss while preparing my job talk:

"just give a good talk"

## How to give the project class talk

Initial conditions:

- I started with a great idea
- It did not work
- The day before the presentation I found 40 papers that already did this work
- Then I also realized that the idea was not so great

How do I present?

Just give a good talk

## Sources on writing technical papers

- How to Get Your SIGGRAPH Paper Rejected, Jim Kajiya, SIGGRAPH 1993 Papers Chair, <u>http://www.siggraph.org/publications/</u> <u>instructions/rejected.html</u>
- Ted Adelson's Informal guidelines for writing a paper, 1991. <u>http://www.ai.mit.edu/courses/6.899/papers/ted.htm</u>
- Notes on technical writing, Don Knuth, 1989.

#### http://www.ai.mit.edu/courses/6.899/papers/knuthAll.pdf

- What's wrong with these equations, David Mermin, Physics Today, Oct., 1989. <u>http://www.ai.mit.edu/courses/6.899/papers/mermin.pdf</u>
- Ten Simple Rules for Mathematical Writing, Dimitri P. Bertsekas <u>http://www.mit.edu:8001/people/dimitrib/Ten\_Rules.html</u>

## Knuth

24. The opening paragraph should be your best paragraph, and its first sentence should be your best sentence. If a paper starts badly, the reader will wince and be resigned to a difficult job of fighting with your prose. Conversely, if the beginning flows smoothly, the reader will be hooked and won't notice occasional lapses in the later parts.

Probably the worst way to start is with a sentence of the form "An x is y." For example,

Bad: An important method for internal sorting is quicksort.

Good: Quicksort is an important method for internal sorting, because ....

Bad: A commonly used data structure is the priority queue.

Good: Priority queues are significant components of the data structures needed for many different applications.

## Knuth on equations

13. Many readers will skim over formulas on their first reading of your exposition. Therefore, your sentences should flow smoothly when all but the simplest formulas are replaced by "blah" or some other grunting noise.

#### The paper impact curve



 $\bigcirc$ 

MIT CSAIL

6.869: Advances in Computer Vision



#### Lecture 22 Scene understanding

## Beyond single classes

- Multiclass
- Multiview
- Datasets

## Beyond single classes

- Multiclass
- Multiview
- Datasets

## Shared features

• Is learning the object class 1000 easier than learning the first?



- Can we transfer knowledge from one object to another?
- Are the shared properties interesting by themselves?

## **Reusable Parts**

Krempp, Geman, & Amit "Sequential Learning of Reusable Parts for Object Detection". TR 2002

Goal: Look for a vocabulary of edges that reduces the number of features.



Examples of reused parts





## Additive models and boosting

• Independent binary classifiers:



• Binary classifiers that share features:



#### Generalization as a function of object similarities



## Beyond single classes

- Multiclass
- Multiview
- Datasets

## **Class experiment**

### **Class experiment**

**Experiment 1:** draw a horse (the entire body, not just the head) in a white piece of paper.

Do not look at your neighbor! You already know how a horse looks like... no need to cheat.

### **Class experiment**

**Experiment 2:** draw a horse (the entire body, not just the head) but this time chose a viewpoint as weird as possible.

## 3D object categorization

Despite we can categorize all three pictures as being views of a horse, the three pictures do not look as being equally typical views of horses. And they do not seem to be recognizable with the same easiness.







## **Canonical Perspective**

**Experiment** (Palmer, Rosch & Chase 81): participants are shown views of an object and are asked to rate "how much each one looked like the objects they depict" (scale; 1=very much like, 7=very unlike)

In a recognition task, reaction time correlated with the ratings.

Canonical views are recognized faster at the entry level.

#### Examples of canonical perspective:







HORSE

PIANO

TEAPOT



CAR



CHAIR

CAMERA



CLOCK



TELEPHONE



HOUSE







PENCIL SHARPENER

SHOE

**IRON** 

From Vision Science, Palmer

## **Canonical Viewpoint**



## **Object representations**

# Explicit 3D models: use volumetric representation. Have an explicit model of the 3D geometry of the object.



Appealing but hard to get it to work...

## **Object representations**

**Implicit 3D models**: matching the input 2D view to view-specific representations.



(b) For cars, classifiers are trained on 8 viewpoints

Not very appealing but somewhat easy to get it to work...

## Beyond single classes

- Multiclass
- Multiview
- Datasets

## The PASCAL Visual Object Classes

In 2007, the twenty object classes that have been selected are:

*Person:* person *Animal:* bird, cat, cow, dog, horse, sheep *Vehicle:* aeroplane, bicycle, boat, bus, car, motorbike, train *Indoor:* bottle, chair, dining table, potted plant, sofa, tv/monitor



M. Everingham, Luc van Gool, C. Williams, J. Winn, A. Zisserman 2007

## Caltech 101 and 256



10<sup>L</sup>

number of training examples per class



Fei–Fei, Fergus, & Perona
SSD baseline

## LabelMe





#### Tool went online July 1st, 2005 530,000 object annotations collected

Labelme.csail.mit.edu

B.C. Russell, A. Torralba, K.P. Murphy, W.T. Freeman, IJCV 2008

#### 80.000.000 images



# IM GENET

- An ontology of images based on WordNet
- ImageNet currently has
  - 13,000+ categories of visual concepts
  - 10 million human-cleaned images (~700im/categ)
  - 1/3+ is released online @ www.image-net.org



Deng, Dong, Socher, Li & Fei-Fei, CVPR 2009
mug

About 10,100,000 results (0.09 seconds)

#### 59¢ Logo Coffee Mugs www.DiscountMugs.com Lead Free & Dishwasher Safe. Save 40-50%. No Catch. Factory Direct !

Custom Mugs On Sale

#### www.Vistaprint.com Order Now & Save 50% On Custom Mugs No Minimums. Upload Photos & Logos.

#### Related searches: white mug coffee mug mug root beer mug shot

Ceramic Happy Face

300 × 300 - 77k - jpg

Find similar images

ceramic mug 980 × 1024 - 30k - jpg

diytrade.com

Coffee Mug as a

303 × 301 - 10k - jpg

Find similar images

dustbowl.wordpress.com

larose.com



Representational 500 × 429 - 91k - jpg eagereyes.org Find similar images

Bring your own

500 × 451 - 15k - jpg

cookstownunited.ca

Find similar images

Back to Ceramic 400 × 400 - 8k - jpg

Find similar images

freshpromotions.com.au



Find similar images



Here I go then, trying 600 × 600 - 35k - jpg beeper.wordpress.com

490 × 428 - 16k - jpg

Find similar images

SASS Life Member

300 × 302 - 6k - jpg

Mugs from LabelMe

sassnet.com

freshome.com





Search

Advanced search



SafeSearch moderate V

www.4imprint.com/Mugs Huge Selection of Style

Colors- Buy 72 Mugs @ \$1.35 ea-24hr Service

Google mugs

Dataset

biases

mug reynosawatch.org



Promotional Mugs from 69¢







personalized coffee 400 × 343 - 15k - jp



walyou.com Find similar i















## Dataset biases



#### Torralba, Efros. Unbiased Look at Dataset Bias. CVPR 2011











# The detector challenge



By looking at the output of a detector on a random set of images, can you guess which object is it trying to detect?

## What object is the detector trying to detect?



By looking at the output of a detector on a random set of images, can you guess which object is it trying to detect?

## What object is the detector trying to detect?



By looking at the output of a detector on a random set of images, can you guess which object is it trying to detect?













### Top 8 out of 4317 images

P. Felzenszwalb, D. McAllester, and D. Ramanan. CVPR, 2008

### Microwave & refrigerator

















### Top 8 out of 4317 images

### What object is hidden behind the red box?





# Objects in context

Torralba, Sinha (2001)



Fink & Perona (2003)

A. eye feature from raw image



B. face feature from raw image



#### Kumar, Hebert (2005)



Carbonetto, de Freitas & Barnard (2004)



Sudderth, Torralba, Wilsky, Freeman (2005)



Heitz and Koller (2008)



Torralba Murphy Freeman (2004)



#### Rabinovich et al (2007)



Hoiem, Efros, Hebert (2005) Horizo Camera Positi **Object Imag** Camera Height 3D Object 3D Obie Object World Height Object Worl

#### Desai, Ramanan, and Fowlkes (2009)



# Increasing the context strength



## Scenes rule over objects



3D percept is driven by the scene, which imposes its ruling to the objects

# Mary Potter (1976)

Mary Potter (1975, 1976) demonstrated that during a rapid sequential visual presentation (100 msec per image), a novel picture is instantly **understood** and observers seem to comprehend a lot of visual information





### **Demo : Rapid image understanding** By Aude Oliva

Instructions: 9 photographs will be shown for half a second each. Your task is to memorize these pictures



















# **Memory Test**

Which of the following pictures have you seen ?

### If you have seen the image clap your hands once

If you have not seen the image do nothing


















#### Have you seen this picture ?





#### Have you seen this picture ?



#### You have seen these pictures



#### You were tested with these pictures



# The gist of the scene

In a glance, we remember the meaning of an image and its global layout but some objects and details are forgotten





# **Scene Categorization**

#### Oliva and Torralba, 2001













Coast

Forest Highway



Mountain

Open Country

Street

Tall Building

#### Fei Fei and Perona, 2005



Bedroom





Living Room





Suburb

Lazebnik, Schmid, and Ponce, 2006



Industrial

Store

15 Scene Database

# Which are the important elements?

Ceiling Liaht	Ceiling Lamp	neinting	wall
Door Door	Painting mirror	painting	
Wall Door Wall Door	wall	wall	Lamp
Floor	Fireplace armchair armchair	Bed	phone alarm
	arrionali		Side-table
	Coffee table		carpet

Different content (i.e. objects), different spatial layout

# Which are the important elements?

cabinets <sup>ceiling</sup> cabinets	cabinets ceiling cabinets	ceiling
window window window seat seat seat seat seat seat seat seat seat seat	window seat seat window seat seat seat seat seat seat seat seat	wall screen seat seat seat seat seat seat seat seat seat seat seat seat seat seat seat seat seat seat seat

Similar objects, and similar spatial layout

Different lighting, different materials, different "stuff"

# What can be an alternative to objects?

## 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



FIG. 8.24. Office, drawn by Robert Mezzanotte.

From "on the semantics of a glance at a scene", Biederman, 1981

#### Examples of scene emergent features



Suggestive edges and junctions







Simple geometric forms



Oliva & Torralba, 2001

Textures ~ Sketch

Blobs

### **Ensemble statistics**

Ariely, 2001, Seeing sets: Representation by statistical properties Chong, Treisman, 2003, Representation of statistical properties Alvarez, Oliva, 2008, 2009, Spatial ensemble statistics



Conclusion: observers had more accurate representation of the mean than of the individual members of the set.

### Global image descriptors

# Global image descriptors

#### Bag of words



Sivic et. al., ICCV 2005 Fei-Fei and Perona, CVPR 2005

#### Non localized textons

i na sua sua su



Walker, Malik. Vision Research 2004

#### Spatially organized textures





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



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



### Gist descriptor



### Example visual gists



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

### **Global features**



"The viewer is presented with a 'potential image', that is, a complex multiplicity of possible images, none of which ever finally resolves".

#### Textons



Filter bank

Malik, Belongie, Shi, Leung, 1999

#### Textons



# **Histogram Intersection**

Histogram intersection

$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = \sum_{j=1}^{\prime} \min(H(\mathbf{X})_j, H(\mathbf{Y})_j)$$

n



Adapted from Kristen Grauman

### SVM

A Support Vector Machine (SVM) learns a classifier with the form:

$$H(x) = \sum_{m=1}^{M} a_m y_m k(x, x_m)$$

Where  $\{x_m, y_m\}$ , for  $m = 1 \dots M$ , are the training data with  $x_m$  being the input feature vector and  $y_m = +1,-1$  the class label.  $k(x, x_m)$  is the kernel and it can be any symmetric function satisfying the Mercer Theorem.

The classification is obtained by thresholding the value of H(x).

There is a large number of possible kernels, each yielding a different family of decision boundaries:

- Linear kernel:  $k(x, x_m) = x^T x_m$
- Radial basis function:  $k(x, x_m) = exp(-|x x_m|^2/\sigma^2)$ .
- Histogram intersection: k(x,x<sub>m</sub>) = sum<sub>i</sub>(min(x(i), x<sub>m</sub>(i)))

# Bag of words



65 17 23 36

2

6 0





# Bag of words & spatial pyramid matching

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



S. Lazebnik, et al, CVPR 2006

### Learning Scene Categorization



#### The 15-scenes benchmark



Oliva & Torralba, 2001 Fei Fei & Perona, 2005 Lazebnik, et al 2006



Office



Skyscrapers









Forest



Bedroom



Living room



Industrial



Street



Highway



Coast



Mountain Open country





Store

### Scene recognition



# **SUN Dataset Project**

We want:

- Large variety of scene categories (we want them all)
- Lots of objects categories
- Multi-object scenes

1. We take all scene words from a dictionary



2. We download images and clean the categories



3. We segment all the images





#### Krista Ehinger





**Jianxiong Xiao** 

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

### 397 Well-sampled Categories



#### Performance with 400 categories



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

#### Training images

Abbey

Airplane cabin

#### Airport terminal

Alley

#### Amphitheater









#### Training images Correct classifications

Abbey

Airplane cabin

Airport terminal

Amphitheater

Alley

11.

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

#### Training images **Correct classifications**

Abbey

Airplane cabin

Airport terminal

Alley

Amphitheater



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

#### Categories or a continuous space?



Check poster by Malisiewicz, Efros
#### Categories or a continuous space?

From the city to the mountains in 10 steps



**Objects in context** 



#### Is local information enough?



#### Is local information even enough?

#### Is local information even enough?







## The system does not care about the scene, but we do...

We know there is a keyboard present in this scene even if we cannot see it clearly.



We know there is no keyboard present in this scene



#### The multiple personalities of a blob









#### The multiple personalities of a blob













ABC



A 13 C 

#### Look-Alikes by Joan Steiner



#### Look-Alikes by Joan Steiner



#### Look-Alikes by Joan Steiner



## The importance of context

- Cognitive psychology
  - Palmer 1975
  - Biederman 1981

- Computer vision
  - Noton and Stark (1971)
  - Hanson and Riseman (1978)
  - Barrow & Tenenbaum (1978)
  - Ohta, kanade, Skai (1978)
  - Haralick (1983)
  - Strat and Fischler (1991)
  - Bobick and Pinhanez (1995)
  - Campbell et al (1997)

Class	Context elements	Operator
SKY	ALWAYS	ABOVE-HORIZON
SKY	SKY-IS-CLEAR ∧ TIME-IS-DAY	BRIGHT
SKY	SKY-IS-CLEAR ∧ TIME-IS-DAY	UNTEXTURED
SKY	SKY-IS-CLEAR ∧ TIME-IS-DAY ∧ RGB-IS-AVAILABLE	BLUE
SKY	SKY-IS-OVERCAST $\land$ TIME-IS-DAY	BRIGHT
SKY	SKY-IS-OVERCAST $\land$ TIME-IS-DAY	UNTEXTURED
SKY	SKY-IS-OVERCAST ∧ TIME-IS-DAY ∧	WHITE
	RGB-IS-AVAILABLE	
SKY	SPARSE-RANGE-IS-AVAILABLE	SPARSE-RANGE-IS-UNDEFINED
SKY	CAMERA-IS-HORIZONTAL	NEAR-TOP
SKY	CAMERA-IS-HORIZONTAL </td <td>ABOVE-SKYLINE</td>	ABOVE-SKYLINE
	CLIQUE-CONTAINS(complete-sky)	
SKY	CLIQUE-CONTAINS(sky)	SIMILAR-INTENSITY
SKY	CLIQUE-CONTAINS(sky)	SIMILAR-TEXTURE
SKY	RGB-IS-AVAILABLE A CLIQUE-CONTAINS(sky)	SIMILAR-COLOR
GROUND	CAMERA-IS-HORIZONTAL	HORIZONTALLY-STRIATED
GROUND	CAMERA-IS-HORIZONTAL	NEAR-BOTTOM
GROUND	SPARSE-RANGE-IS-AVAILABLE	SPARSE-RANGES-FORM-HORIZONT
GROUND	DENSE-RANGE-IS-AVAILABLE	DENSE-RANGES-FORM-HORIZONTA
GROUND	CAMERA-IS-HORIZONTAL	BELOW-SKYLINE
	CLIQUE-CONTAINS(complete-ground)	
GROUND	CAMERA-IS-HORIZONTAL A	BELOW-GEOMETRIC-HORIZON
	CLIQUE-CONTAINS(geometric-horizon) </td <td></td>	
	- CLIQUE-CONTAINS(skyline)	
GROUND	TIME-IS-DAY	DARK

#### **Objects and Scenes**

Stimuli from Hock, Romanski, Galie, and Williams (1978).



Biederman's violations (1981):

- 1. Support (e.g., a floating fire hydrant). The object does not appear to be resting on a surface.
- Interposition (e.g., the background appearing through the hydrant). The objects undergoing this
  violation appear to be transparent or passing through another object.
- 3. Probability (e.g., the hydrant in a kitchen). The object is unlikely to appear in the scene.
- Position (e.g., the fire hydrant on top of a mailbox in a street scene). The object is likely to occur in that scene, but it is unlikely to be in that particular position.
- Size (e.g., the fire hydrant appearing larger than a building). The object appears to be too large or too small relative to the other objects in the scene.

## **CONDOR** system

Strat and Fischler (1991)

Class	Context elements	Operator	
SKY	ALWAYS	ABOVE-HORIZON	
SKY	SKY-IS-CLEAR ^ TIME-IS-DAY	BRIGHT	
SKY	SKY-IS-CLEAR ∧ TIME-IS-DAY	UNTEXTURED	
SKY	SKY-IS-CLEAR ∧ TIME-IS-DAY ∧ RGB-IS-AVAILABLE	BLUE	
SKY	SKY-IS-OVERCAST $\land$ TIME-IS-DAY	BRIGHT	
SKY	SKY-IS-OVERCAST ∧ TIME-IS-DAY	UNTEXTURED	
SKY	SKY-IS-OVERCAST ∧ TIME-IS-DAY ∧	WHITE	
	RGB-IS-AVAILABLE		
SKY	SPARSE-RANGE-IS-AVAILABLE	SPARSE-RANGE-IS-UNDEFINED	
SKY	CAMERA-IS-HORIZONTAL	NEAR-TOP	
SKY	CAMERA-IS-HORIZONTAL $\land$	ABOVE-SKYLINE	
	CLIQUE-CONTAINS(complete-sky)		
SKY	CLIQUE-CONTAINS(sky)	SIMILAR-INTENSITY	
SKY	CLIQUE-CONTAINS(sky)	SIMILAR-TEXTURE	
SKY	RGB-IS-AVAILABLE	SIMILAR-COLOR	
GROUND	CAMERA-IS-HORIZONTAL	HORIZONTALLY-STRIATED	
GROUND	CAMERA-IS-HORIZONTAL	NEAR-BOTTOM	
GROUND	SPARSE-RANGE-IS-AVAILABLE	SPARSE-RANGES-FORM-HORIZONT	
GROUND	DENSE-RANGE-IS-AVAILABLE	DENSE-RANGES-FORM-HORIZONTA	
GROUND	CAMERA-IS-HORIZONTAL A	BELOW-SKYLINE	
	CLIQUE-CONTAINS(complete-ground)		
GROUND	CAMERA-IS-HORIZONTAL A	BELOW-GEOMETRIC-HORIZON	
	CLIQUE-CONTAINS(geometric-horizon) $\land$		
	- CLIQUE-CONTAINS(skyline)		
GROUND	TIME-IS-DAY	DARK	

- Guzman (SEE), 1968
- Noton and Stark 1971
- Hansen & Riseman (VISIONS), 1978
- Barrow & Tenenbaum 1978

- Brooks (ACRONYM), 1979
- Marr, 1982
- Ohta & Kanade, 1978
- Yakimovsky & Feldman, 1973

## An Age of Scene Understanding





(b) Top-down process [Ohta & Kanade 1978]



- Guzman (SEE), 1968
- Noton and Stark 1971
- Hansen & Riseman (VISIONS), 1978
- Barrow & Tenenbaum 1978

- Brooks (ACRONYM), 1979
- Marr, 1982

💬 terriaret ....)

- Ohta & Kanade, 1978
- Yakimovsky & Feldman, 1973















#### Context models





Objects are correlated via the scene



Dependencies among objects

#### **Context models**







Dependencies among objects

#### **Global precedence**

A SPACE IN

Forest Before Trees: The Precedence of Global Features in Visual Perception Navon (1977)



#### Global and local representations



#### Global and local representations



#### An integrated model of Scenes, Objects, and Parts



# Context-based vision system for place and object recognition



- Hidden states = location (63 values)
- Observations =  $v_t^G$  (80 dimensions)
- Transition matrix encodes topology of environment
- Observation model is a mixture of Gaussians centered on prototypes (100 views per place)

## Our mobile rig





Torralba, Murphy, Freeman, Rubin. 2003

#### Place recognition demo



#### Identification and categorization of known places



#### An integrated model of Scenes, Objects, and Parts





Murphy, Torralba, Freeman; NIPS 2003. Torralba, Murphy, Freeman, CACM 2010.

#### Application of object detection for image retrieval

#### **Results using the keyboard detector alone**



#### Application of object detection for image retrieval


### Object retrieval: scene features vs. detector

#### Results using the keyboard detector alone



Results using both the detector and the global scene features





Murphy, Torralba, Freeman; NIPS 2003. Torralba, Murphy, Freeman, CACM 2010.

# Localizing the object



# An integrated model of Scenes, Objects, and Parts





# Predicting object location

### Training set (cars) $Z|g = \sum (A_n g + b_n) W_n(g)$ $\{g^1, Z^1\}$ g(1) $\{g^2, Z^2\}$ Ζ $\{g^3, Z^3\}$ g(2) g(1)

### **Predicting location**















Torralba & Sinha, 2001; Murphy, Torralba, Freeman, 2003; Hoeim, Efros, Hebert. 2006

#### screens







keyboard







car





pedestrian

# An integrated model of Scenes, Objects, and Parts



We train a multiview car detector.





$$p(d | F=1) = N(d | \mu_1, \sigma_1)$$
  
 $p(d | F=0) = N(d | \mu_0, \sigma_0)$ 

# An integrated model of Scenes, Objects, and Parts







a) input image

b) car detector output

c) location priming

c) integrated model output

### Two tasks



### A car out of context ...



### A car out of context ...



### 3d Scene Context



Hoiem, Efros, Hebert ICCV 2005

# 3d Scene Context



Hoiem, Efros, Hebert ICCV 2005

### 3D City Modeling using Cognitive Loops

![](_page_158_Picture_1.jpeg)

Figure 6. Stages of the recognition system: (a) initial detections before and (b) after applying ground plane constraints, (c) temporal integration on reconstructed map, (d) estimated 3D car locations, rendered back into the original image.

N. Cornelis, B. Leibe, K. Cornelis, L. Van Gool. CVPR'06

### **Context models**

![](_page_159_Figure_1.jpeg)

![](_page_159_Figure_2.jpeg)

Objects are correlated via the scene

![](_page_159_Figure_4.jpeg)

Dependencies among objects

![](_page_160_Picture_0.jpeg)

1) Generate candidate objects (run a detector, or segmentation)

M possible object labels N regions

Label:  $c_k = [1...M]$  with k = [1...N]Scores:  $s_k$  = vector length M

2) For each candidate, get a list of possible interpretations with their probabilities

 $p(c_k = m \mid s_k)$ 

 Goal: to assign labels c<sub>k</sub> to each candidate so that they are in contextual agreement. We want to optimize the joint probability of all the labels:

$$p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N)$$

**Goal**: to assign labels c<sub>k</sub> to each candidate so that they are in contextual agreement.

M possible object labels N regions

Label:  $c_k = [1...M]$  with k = [1...N]Scores:  $s_k$  = vector length M

![](_page_161_Picture_3.jpeg)

We want to optimize the joint probability of all the labels:

 $p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N)$ 

Solution 1: Assume objects are independent:

**Computed** building  $p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N) = \prod_{i=1...N} p(c_i = m_i | s_i)$ 

![](_page_161_Picture_8.jpeg)

Independent model

**Problem**: it does not makes use of the correlation between objects in the world. This is fine if the detectors are perfect.

**Goal**: to assign labels c<sub>k</sub> to each candidate so that they are in contextual agreement.

M possible object labels N regions

Label:  $c_k = [1...M]$  with k = [1...N]Scores:  $s_k$  = vector length M

![](_page_162_Picture_3.jpeg)

We want to optimize the joint probability of all the labels:

 $p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N)$ 

Solution 2: Assume objects are fully dependent:

 $p(c_1=m_1,..., c_N=m_N|s_1,..., s_N) =$ 

=

 $\frac{p(s_1,...,s_N|c_1=m_1,...,c_N=m_N) p(c_1=m_1,...,c_N=m_N)}{Z(s_1,...,s_N)}$ 

$$\prod_{i=1...N} p(s_i | c_i = m_i) p(c_1 = m_1, ..., c_N = m_N)$$

 $Z(s_1, \dots, s_N) = \sum_{\text{All } [c_1, \dots, c_N]} \prod_{\text{assignments}} p(s_i | c_i = m_i) p(c_1 = m_1, \dots, c_N = m_N)$ 

 $Z(S_1, \ldots S_N)$ 

**Problem**: learning  $p(c_1=m_1,...,c_N=m_N)$  will need a lot of data. Recognition can be slow.

c3

**Goal**: to assign labels c<sub>k</sub> to each candidate so that they are in contextual agreement.

M possible object labels N regions

Label:  $c_k = [1...M]$  with k = [1...N]Scores:  $s_k$  = vector length M

![](_page_163_Picture_3.jpeg)

We want to optimize the joint probability of all the labels:

 $p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N)$ 

Solution 3: Approximated model of dependencies:

$$p(c_1 = m_1, ..., c_N = m_N | s_1, ..., s_N) = \prod_{i=1...N} p(s_i | c_i = m_i) p(c_1 = m_1, ..., c_N = m_N) Z(s_1, ..., s_N)$$

$$p(c_1=m_1,\ldots,c_N=m_N) = \exp\left(\sum_{i,j=1\ldots N} \Phi(c_i=m_i, c_j=m_j)\right)$$

 $\Phi(c_i=m_i, c_j=m_j) = co-ocurrence matrix on training set (count how many times two objects appear together).$ 

**Problem**: learning  $p(c_1=m_1,...,c_N=m_N)$  will be easier, but recognition may still be slow.

 $\Phi(c_i=m_i, c_j=m_j) = co-ocurrence matrix on training set (count how many times two objects appear together).$ 

#### MSRC training data

![](_page_164_Figure_2.jpeg)

165

A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora and S. Belongie. Objects in Context. ICCV 2007

![](_page_165_Picture_0.jpeg)

A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora and S. Belongie. Objects in Context. ICCV 2007

### Objects in context

Torralba, Sinha (2001)

![](_page_166_Picture_2.jpeg)

Fink & Perona (2003)

A. eye feature from raw image

![](_page_166_Picture_5.jpeg)

B. face feature from raw image

![](_page_166_Picture_9.jpeg)

#### Kumar, Hebert (2005)

![](_page_166_Picture_11.jpeg)

Carbonetto, de Freitas & Barnard (2004)

![](_page_166_Figure_13.jpeg)

Sudderth, Torralba, Wilsky, Freeman (2005)

![](_page_166_Picture_15.jpeg)

Heitz and Koller (2008)

![](_page_166_Picture_17.jpeg)

Torralba Murphy Freeman (2004)

![](_page_166_Figure_19.jpeg)

#### Rabinovich et al (2007)

![](_page_166_Picture_21.jpeg)

Hoiem, Efros, Hebert (2005) Horizo Camera Positi **Object Imag** Camera Height 3D Object 3D Obie Object World Height Object Worl

#### Desai, Ramanan, and Fowlkes (2009)

![](_page_166_Picture_24.jpeg)

# **Object-Object Relationships**

• Fink & Perona (NIPS 03)

Use output of boosting from other objects at previous iterations as input into boosting for this iteration

![](_page_167_Figure_3.jpeg)

Figure 5: A-E. Emerging features of eyes, mouths and faces (presented on windows of raw images for legibility). The windows' scale is defined by the detected object size and by the map mode (local or contextual). C. faces are detected using face detection maps H<sup>Face</sup>, exploiting the fact that faces tend to be horizontally aligned.

# Pixel labeling using MRFs

Enforce consistency between neighboring labels, and between labels and pixels

$$P(L,x) = P(L)P(x|L) = \left[\frac{1}{Z}\prod_{i}\prod_{j\in N_i}\psi_{ij}(L_i,L_j)\right]\left[\prod_{i}P(x_i|L_i)\right]$$

![](_page_168_Picture_3.jpeg)

Carbonetto, de Freitas & Barnard, ECCV'04

# Beyond nearest-neighbor grids

- Most MRF/CRF models assume nearestneighbor graph topology
- This cannot capture long-distance correlations

![](_page_169_Picture_3.jpeg)

![](_page_169_Picture_4.jpeg)

![](_page_169_Picture_5.jpeg)

### Dynamically structured trees

• Each node pick its parents (Storkey& Williams, PAMI'03)

![](_page_170_Picture_2.jpeg)

• 2D SCFGs

(Pollak, Siskind, Harper & Bouman ICASSP'03)

![](_page_170_Picture_5.jpeg)

### **Object-Object Relationships**

Use latent variables to induce long distance correlations between labels in a Conditional Random Field (CRF)

![](_page_171_Figure_2.jpeg)

![](_page_171_Picture_3.jpeg)

![](_page_171_Picture_4.jpeg)

#### He, Zemel & Carreira-Perpinan (04)

# **Object-Object Relationships**

![](_page_172_Picture_1.jpeg)

[Kumar Hebert 2005]

# 3d Scene Context

![](_page_173_Figure_1.jpeg)

### Using stuff to find things

Heitz and Koller, ECCV 2008

In this work, there is not labeling for stuff. Instead, they look for clusters of textures and model how each cluster correlates with the target object.

![](_page_174_Picture_3.jpeg)

# What, where and who? Classifying events by scene and object recognition

![](_page_175_Picture_1.jpeg)

L-J Li & L. Fei-Fei, ICCV 2007

![](_page_176_Picture_0.jpeg)

Slide by Fei-fei

L.-J. Li & L. Fei-Fei ICCV 2007

# Grammars

![](_page_177_Figure_1.jpeg)

[Ohta & Kanade 1978]

![](_page_177_Figure_3.jpeg)

- Guzman (SEE), 1968
- Noton and Stark 1971
- Hansen & Riseman (VISIONS), 1978
- Barrow & Tenenbaum 1978
- Brooks (ACRONYM), 1979
- Marr, 1982
- Yakimovsky & Feldman, 1973

### Grammars for objects and scenes

![](_page_178_Figure_1.jpeg)

S.C. Zhu and D. Mumford. A Stochastic Grammar of Images. Foundations and Trends in Computer Graphics and Vision, 2006.

### Who needs context anyway? We can recognize objects even out of context

![](_page_179_Picture_1.jpeg)

Banksy