

## Objects and scenes:

# Recognizing Multiple Object Classes

Josef Sivic and Ivan Laptev

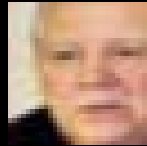
<http://www.di.ens.fr/~josef>

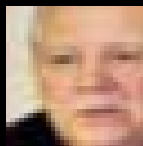
INRIA, WILLOW, ENS/INRIA/CNRS UMR 8548

Laboratoire d'Informatique, Ecole Normale Supérieure, Paris

With slides from: A. Torralba, D. Hoiem, D. Ramanan and others.

# Multiclass object detection

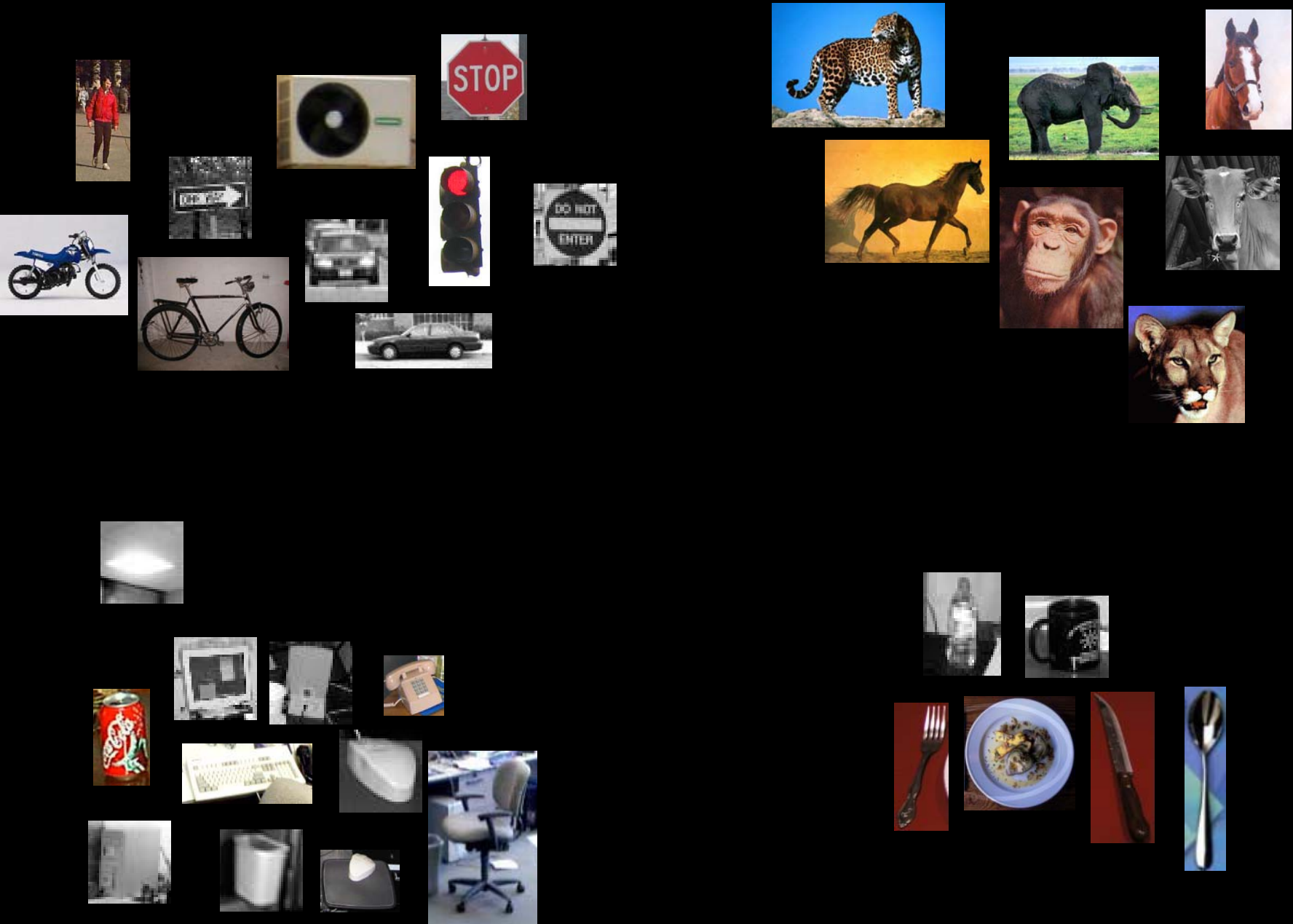




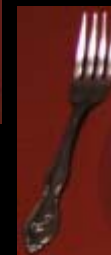
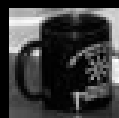
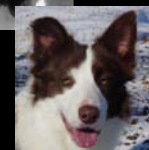
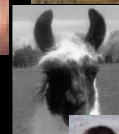
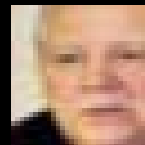




# Context: objects appear in configurations



# Generalization: objects share parts



# How many categories?



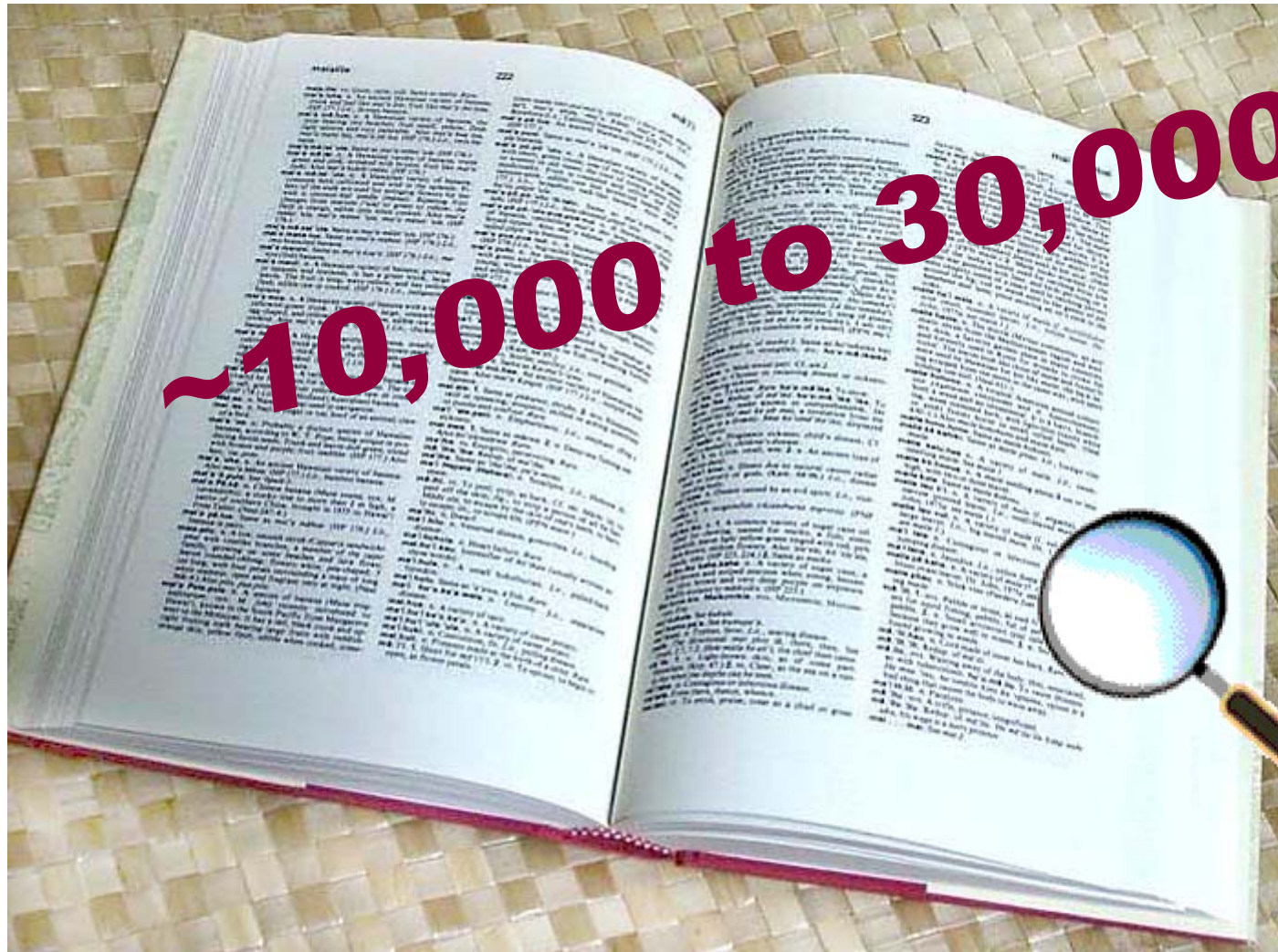
# How many categories?



Slide by Aude Oliva



# How many object categories are there?



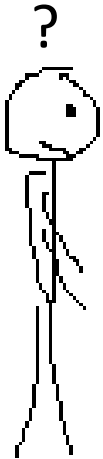
# How many categories?

- Probably this question is not even specific enough to have an answer

# Which level of categorization is the right one?

Car is an object composed of:

a few doors, four wheels (not all visible at all times), a roof, front lights, windshield



If you are thinking in buying a car, you might want to be a bit more specific about your categorization level.

# Entry-level categories

(Jolicoeur, Gluck, Kosslyn 1984)

- Typical member of a basic-level category are categorized at the expected level
- Atypical members tend to be classified at a subordinate level.



A bird



An ostrich



# We do not need to recognize the exact category

A new class can borrow information from similar categories



# So, where is computer vision?

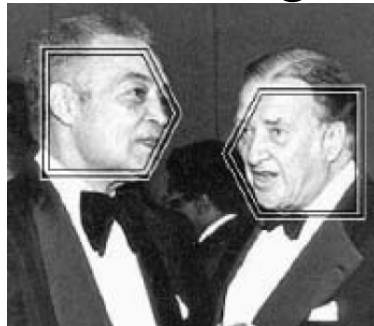
Well...

# Multiclass object detection

## the not so early days

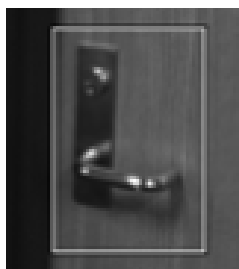
Using a set of independent binary classifiers was a common strategy:

- Viola-Jones extension for dealing with rotations



- two cascades for each view

- Schneiderman-Kanade multiclass object detection



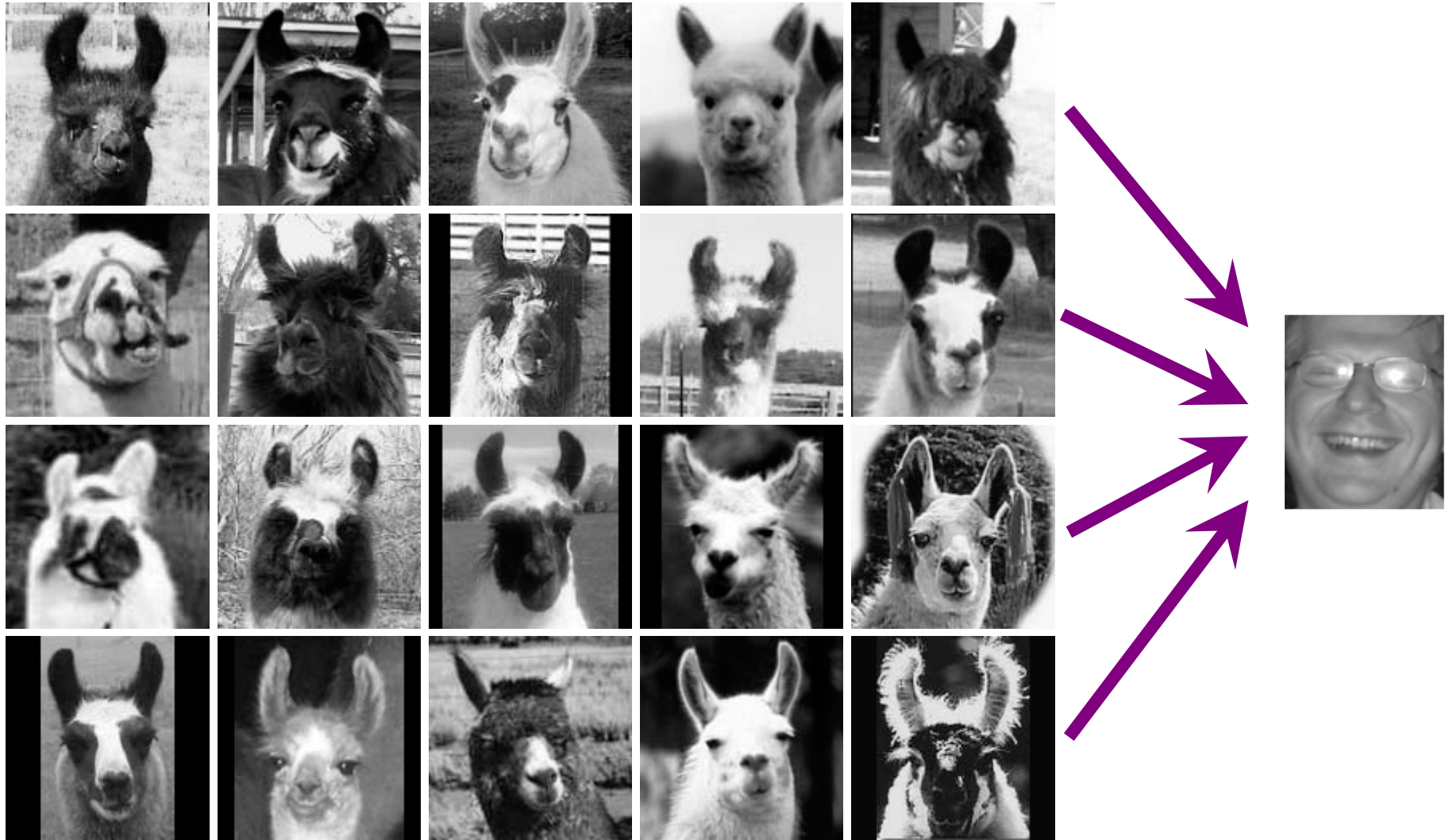
(a) One detector for each class



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

There is nothing wrong with this approach if you have access to lots of training data and you do not care about efficiency.

# Generalizing Across Categories



*Can we transfer knowledge from one object category to another?*

Slide by Erik Sudderth

# Shared features

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



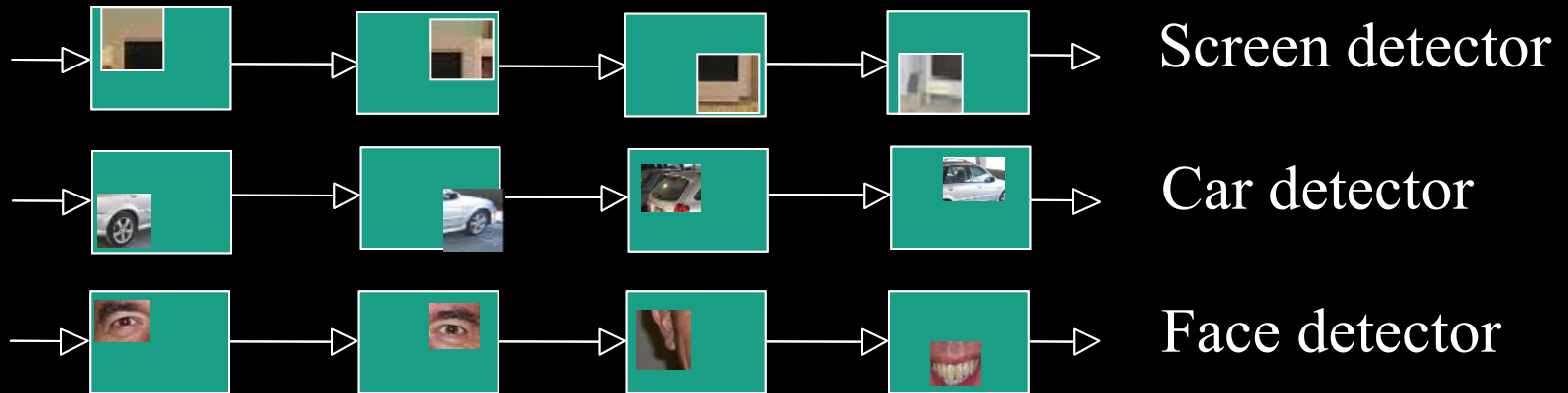
...



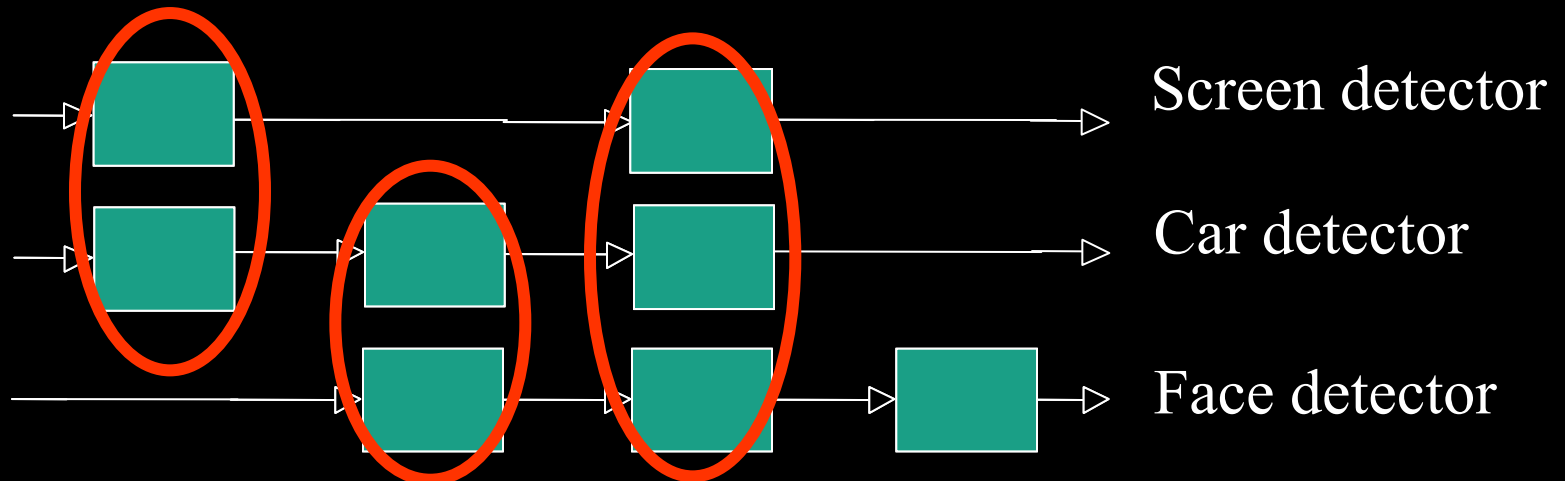
- Can we transfer knowledge from one object to another?
- Are the shared properties interesting by themselves?

# Additive models and boosting

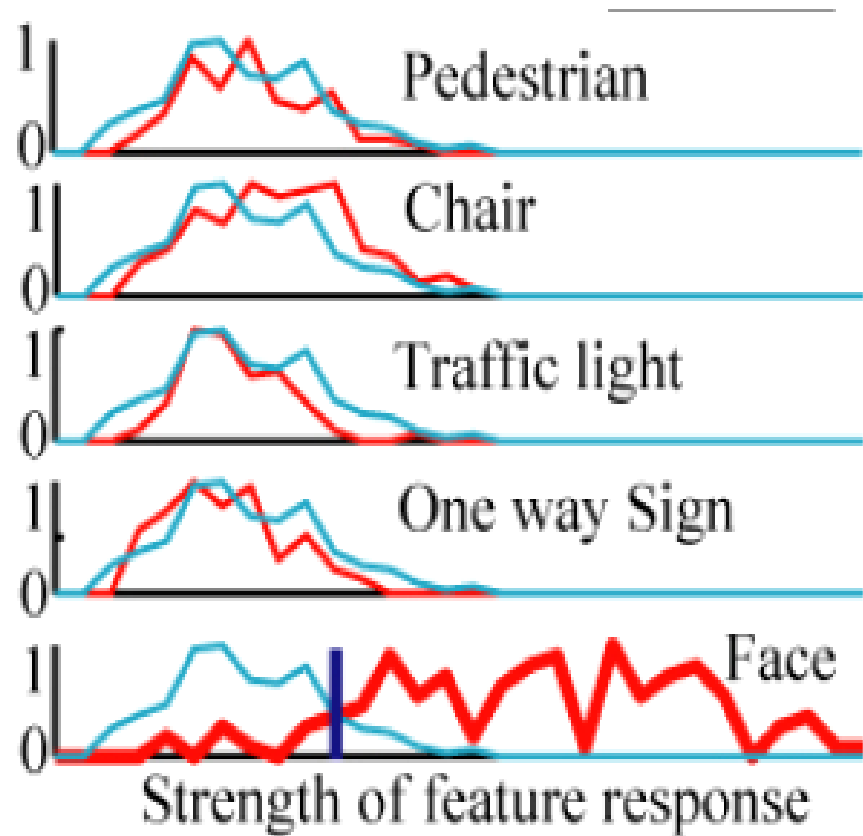
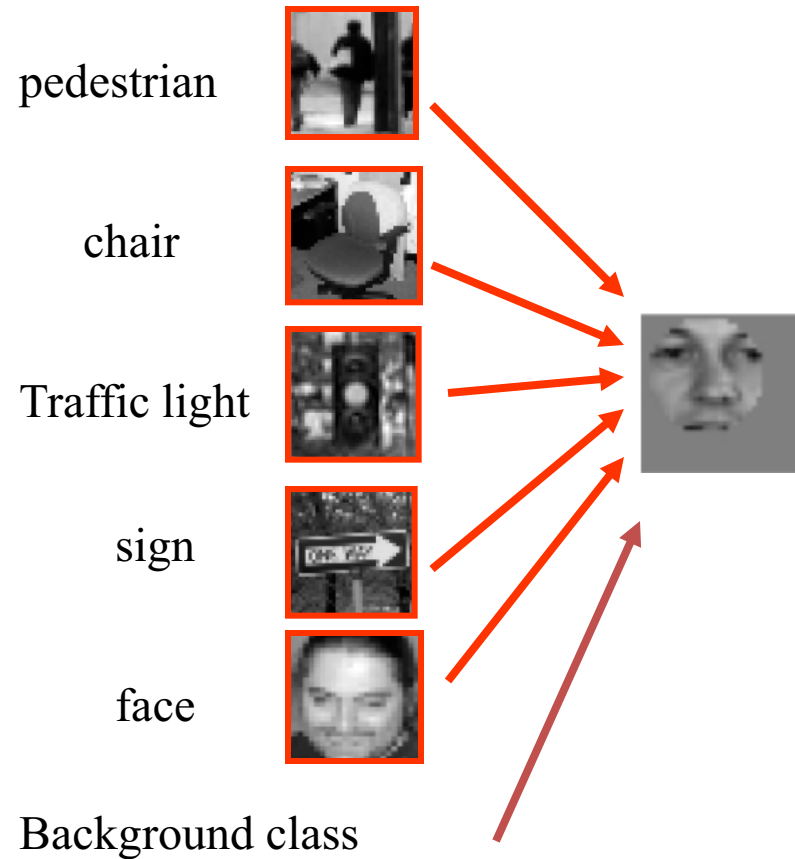
- Independent binary classifiers:



- Binary classifiers that share features:

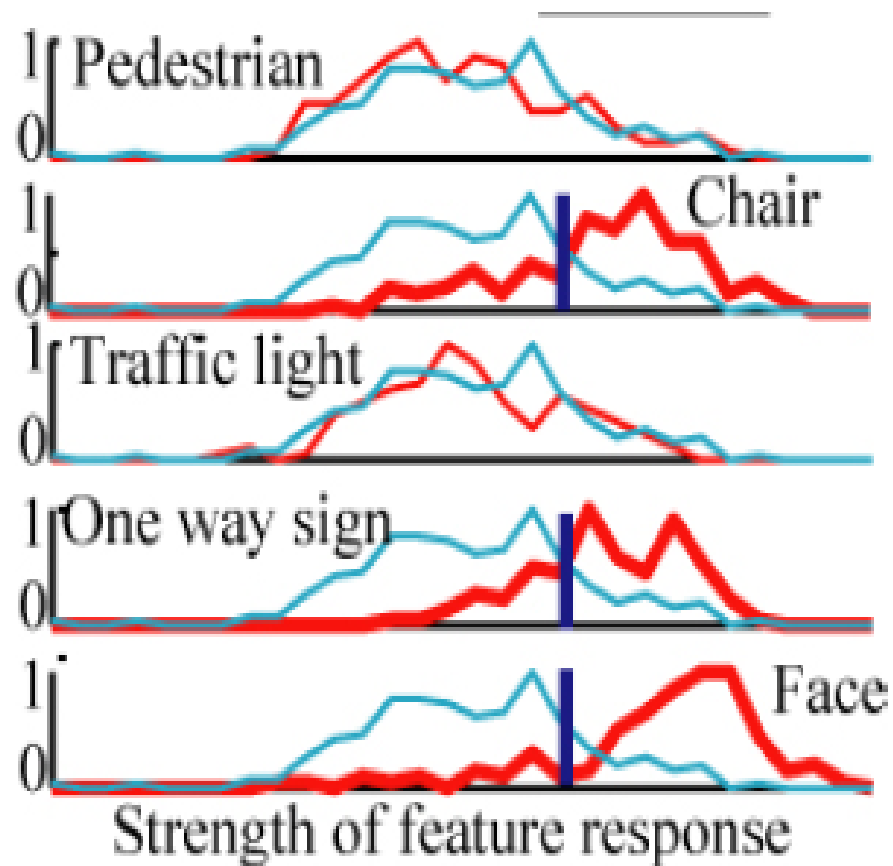
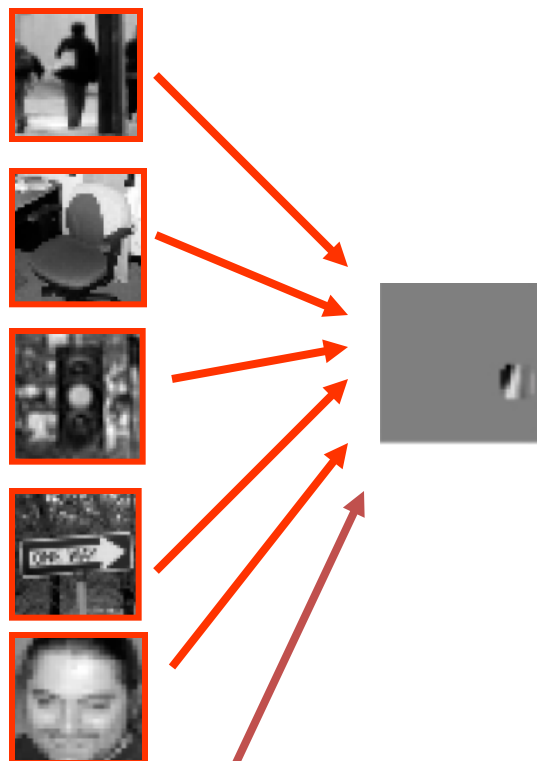


# Specific feature



Non-shared feature: this feature is too specific to faces.

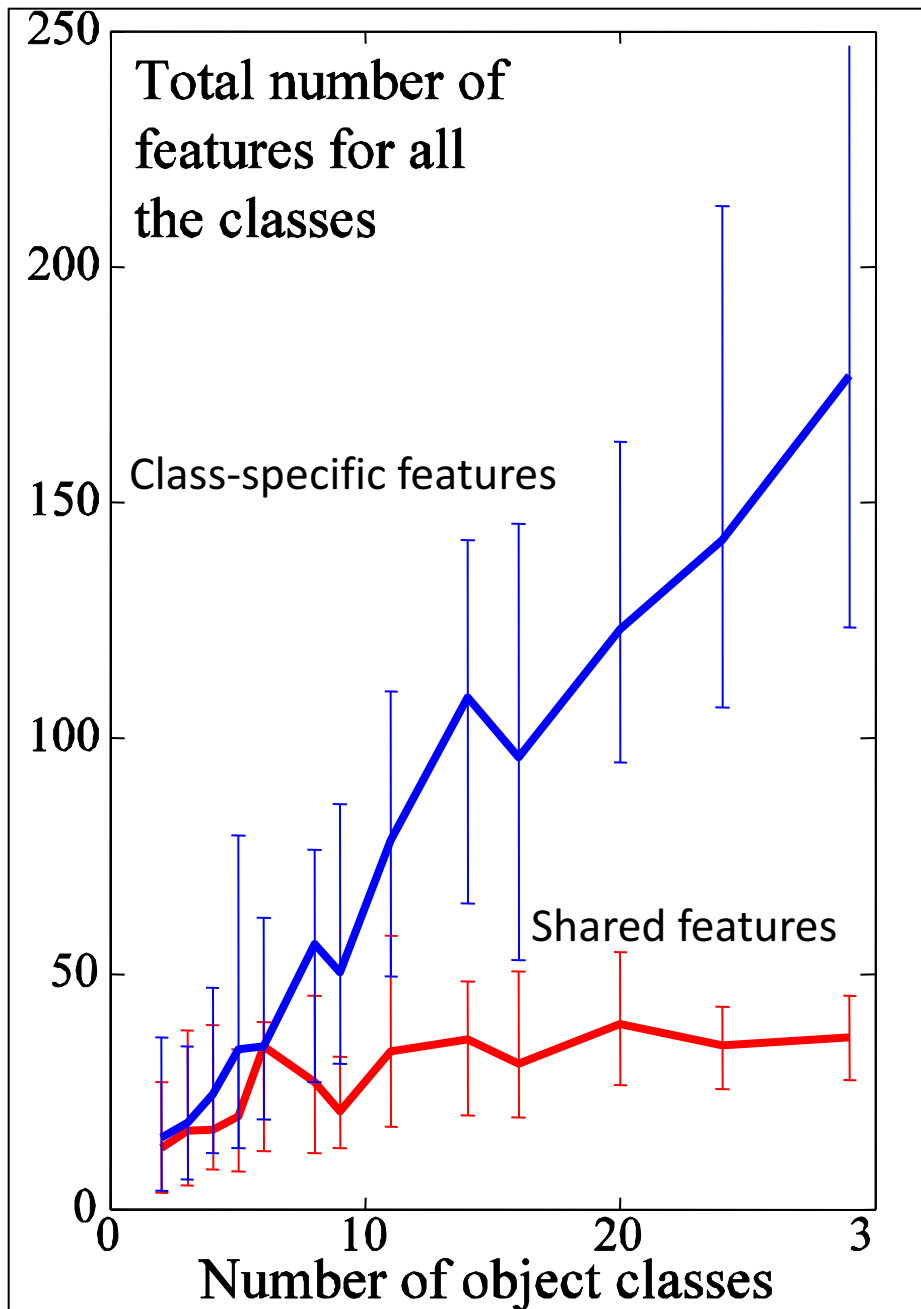
# Shared feature



shared feature





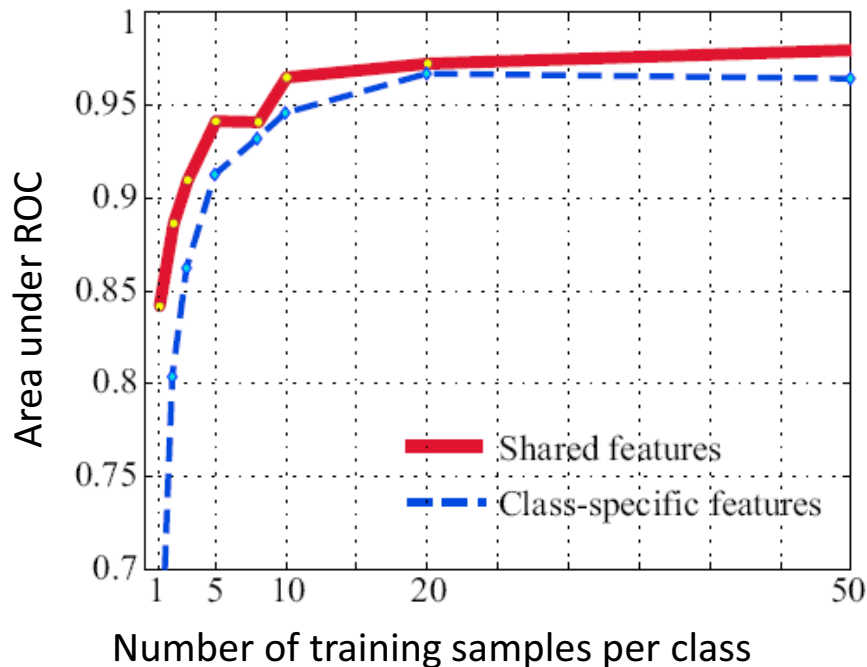
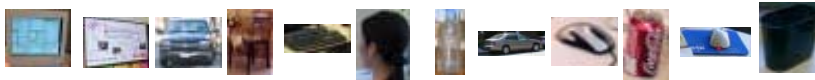


50 training samples/class  
29 object classes  
2000 entries in the dictionary

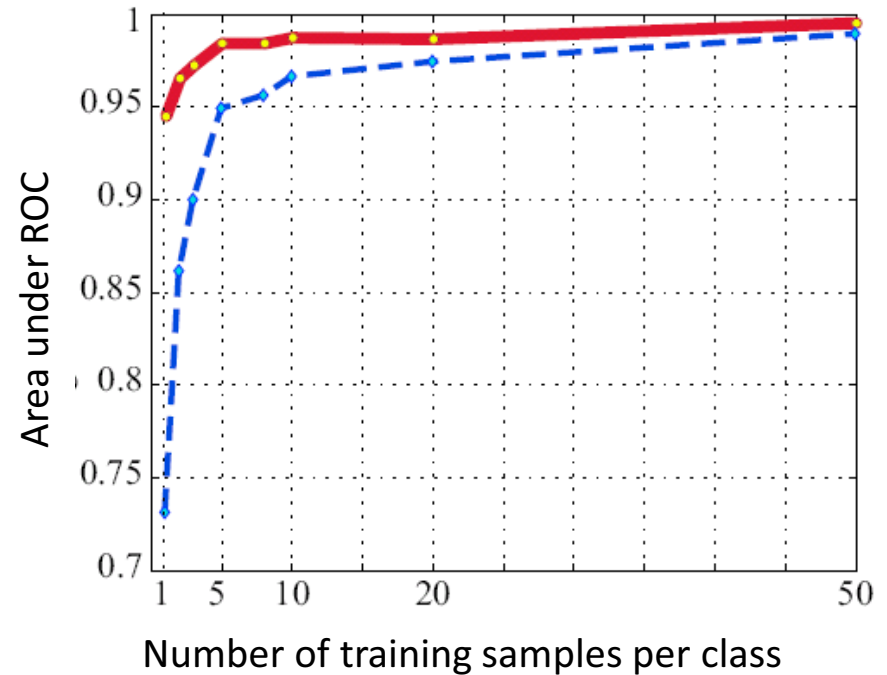
Results averaged on 20 runs

# Generalization as a function of object similarities

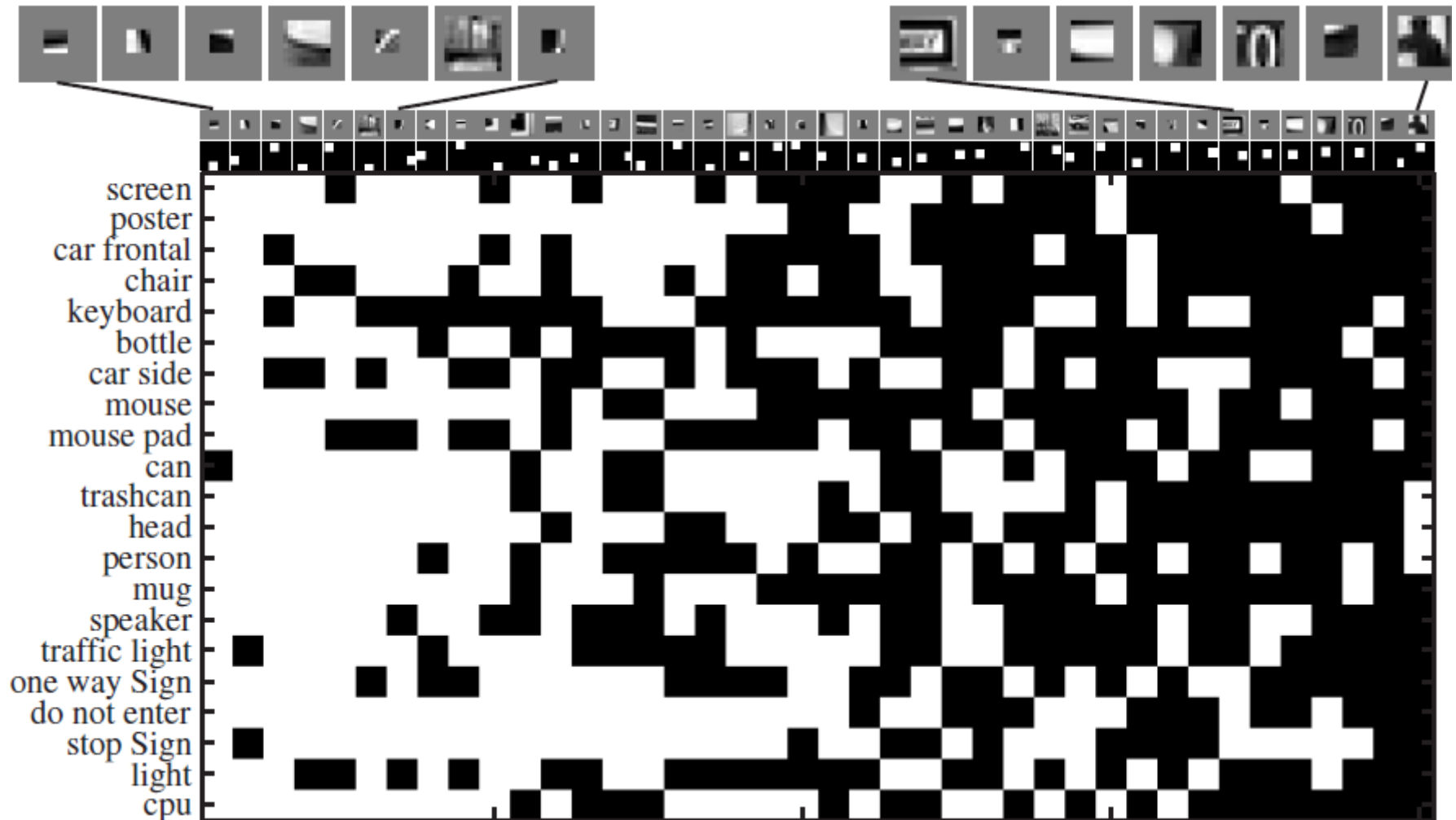
12 unrelated object classes



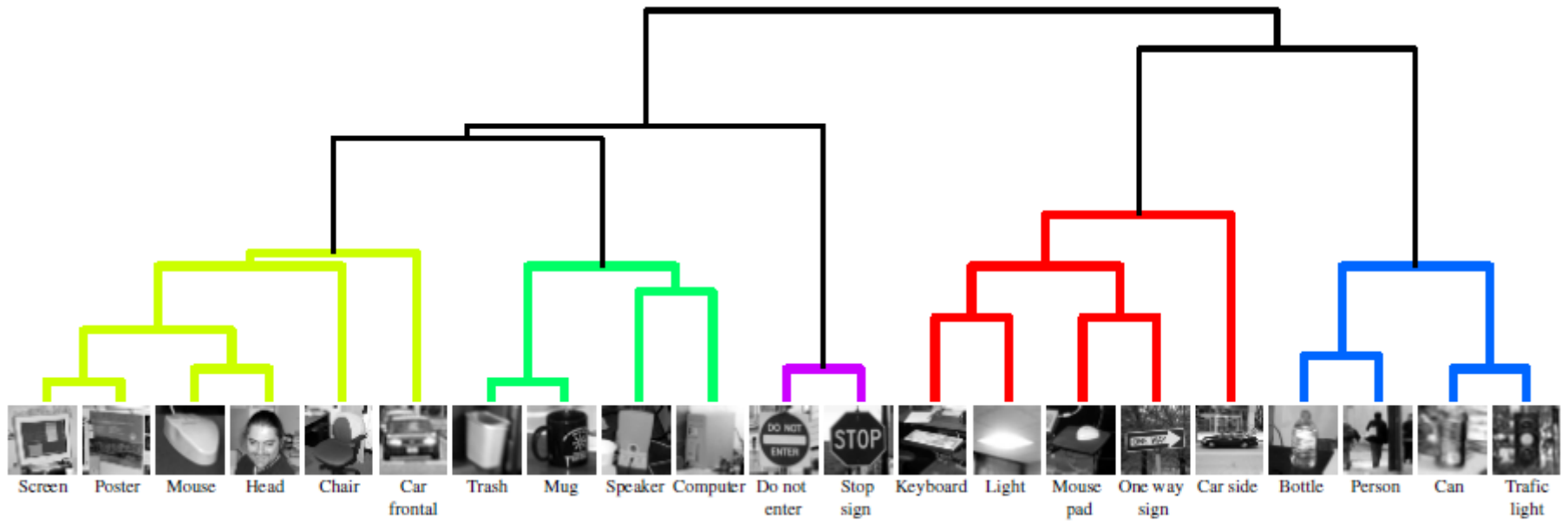
12 viewpoints



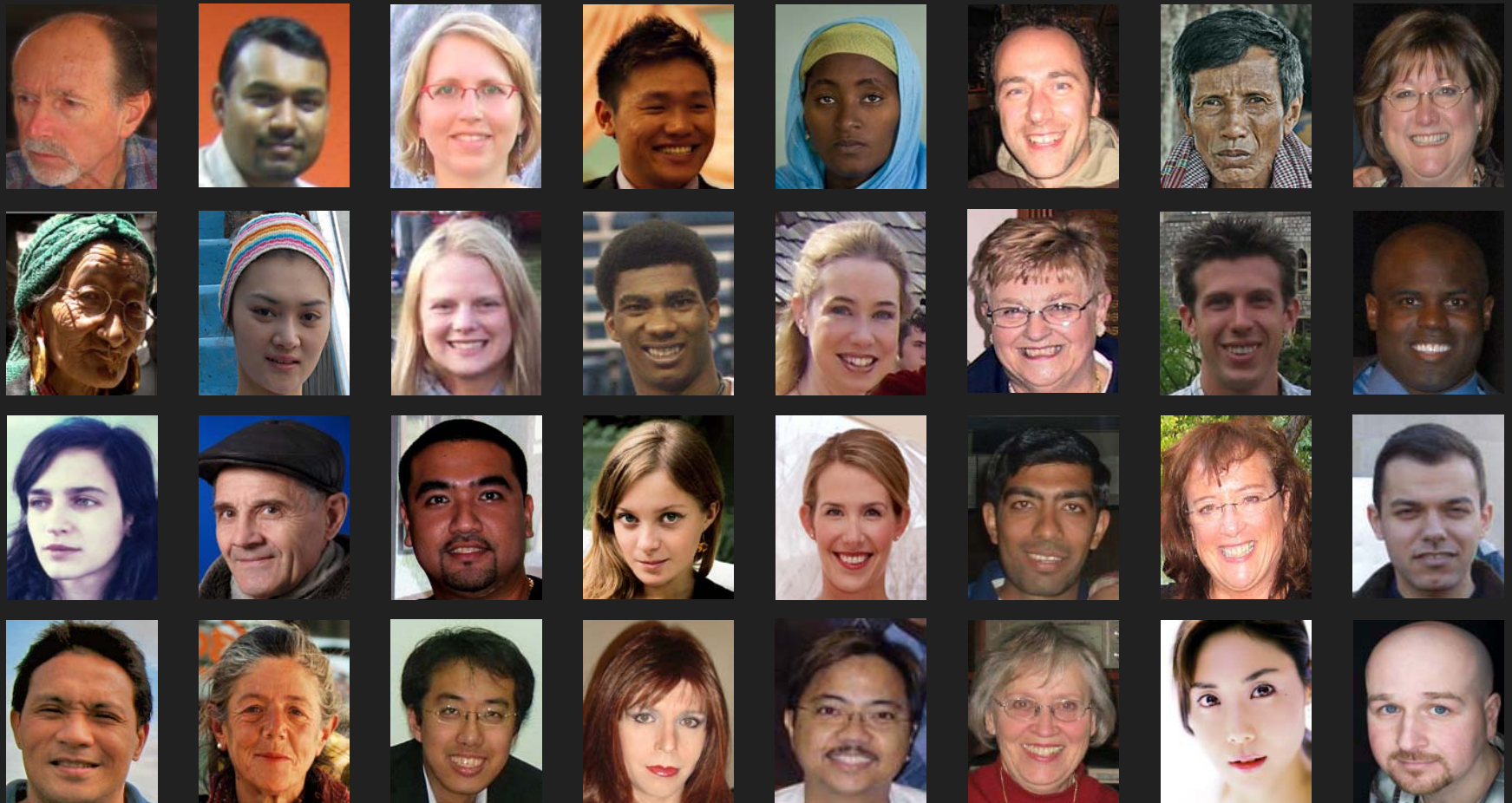
# Generic vs. specific features



# Object clustering according to shared features



# Another multi-class problem: Face recognition



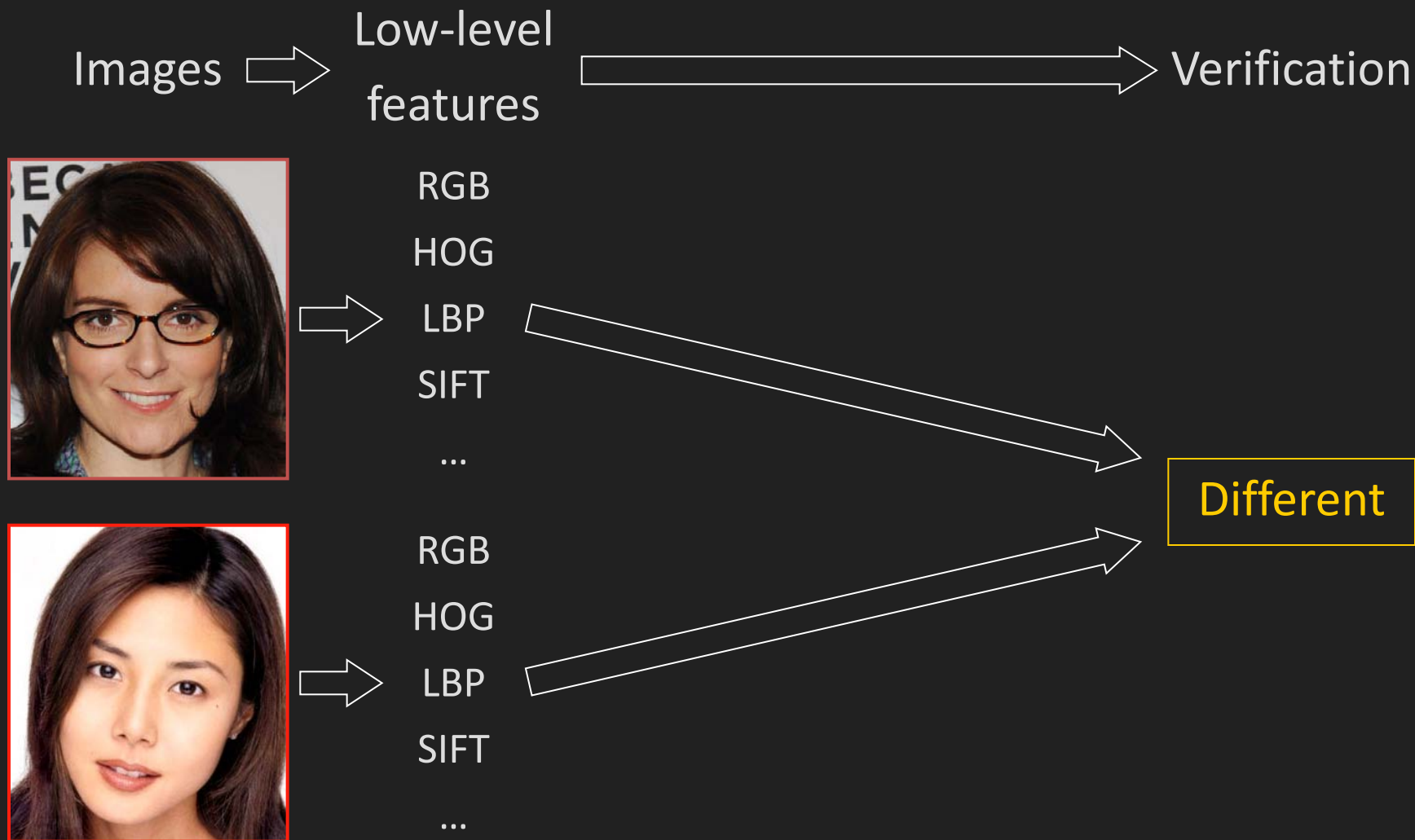
We do not want to learn recognition of each person from scratch!



# Are these images of the same person?

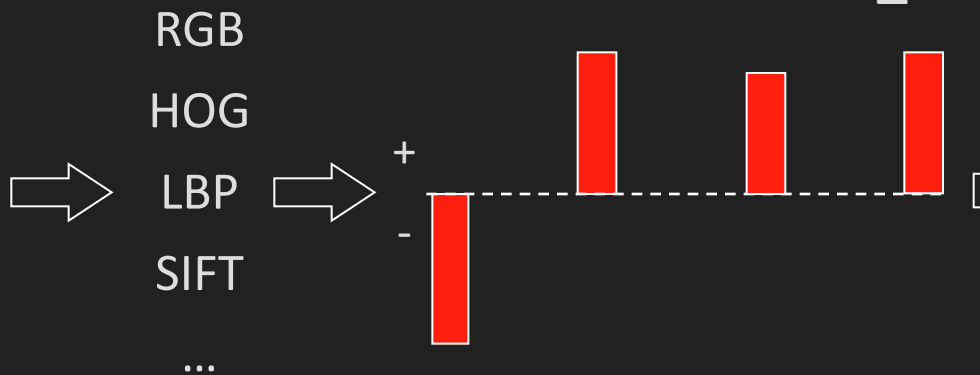
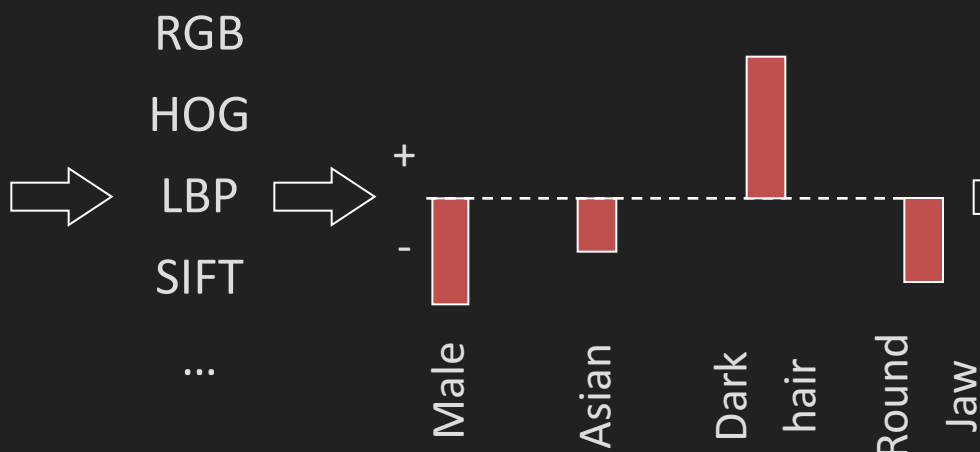


# Prior approaches



# Approach: attributes

Images  $\Rightarrow$  Low-level features  $\Rightarrow$  **Attributes**  $\Rightarrow$  Verification



**Different**



# Attributes can define categories

Female      Caucasian      Middle-aged  
                 Eyeglasses      Dark hair



# Some attributes may be irrelevant

Teeth showing

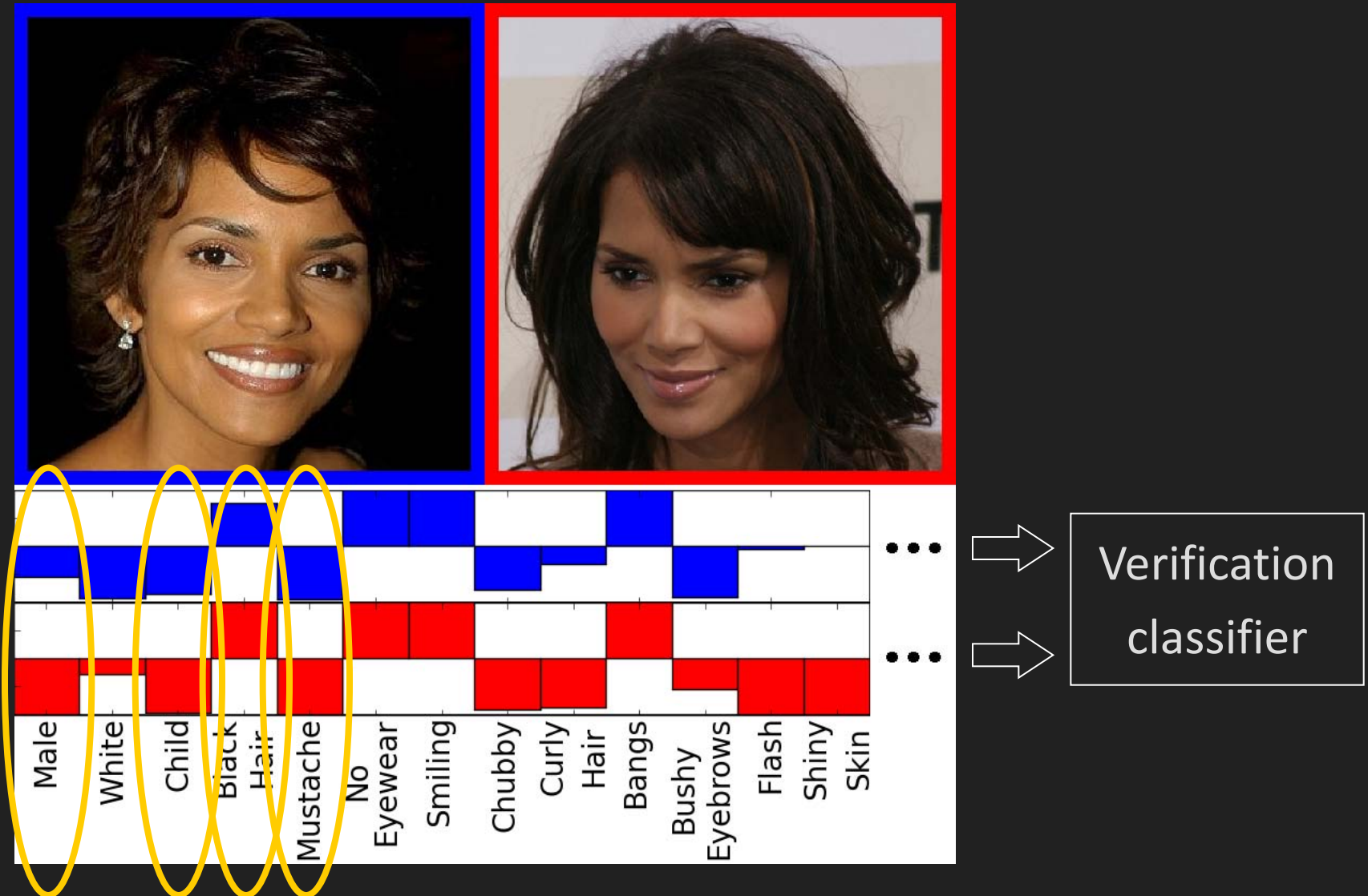
Tilted head

Outside

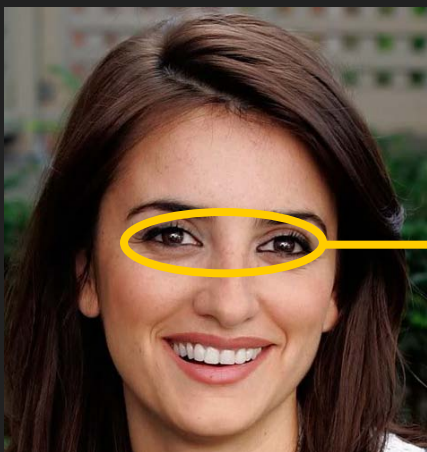




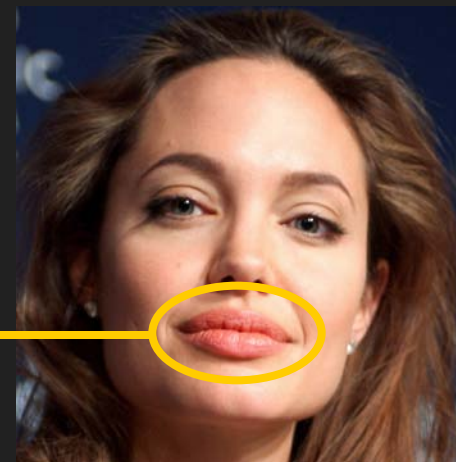
# Using attributes to perform verification



# Describe faces using similes



Penelope Cruz



Angelina Jolie

# Training simile classifiers



Images of Penelope Cruz 's eyes

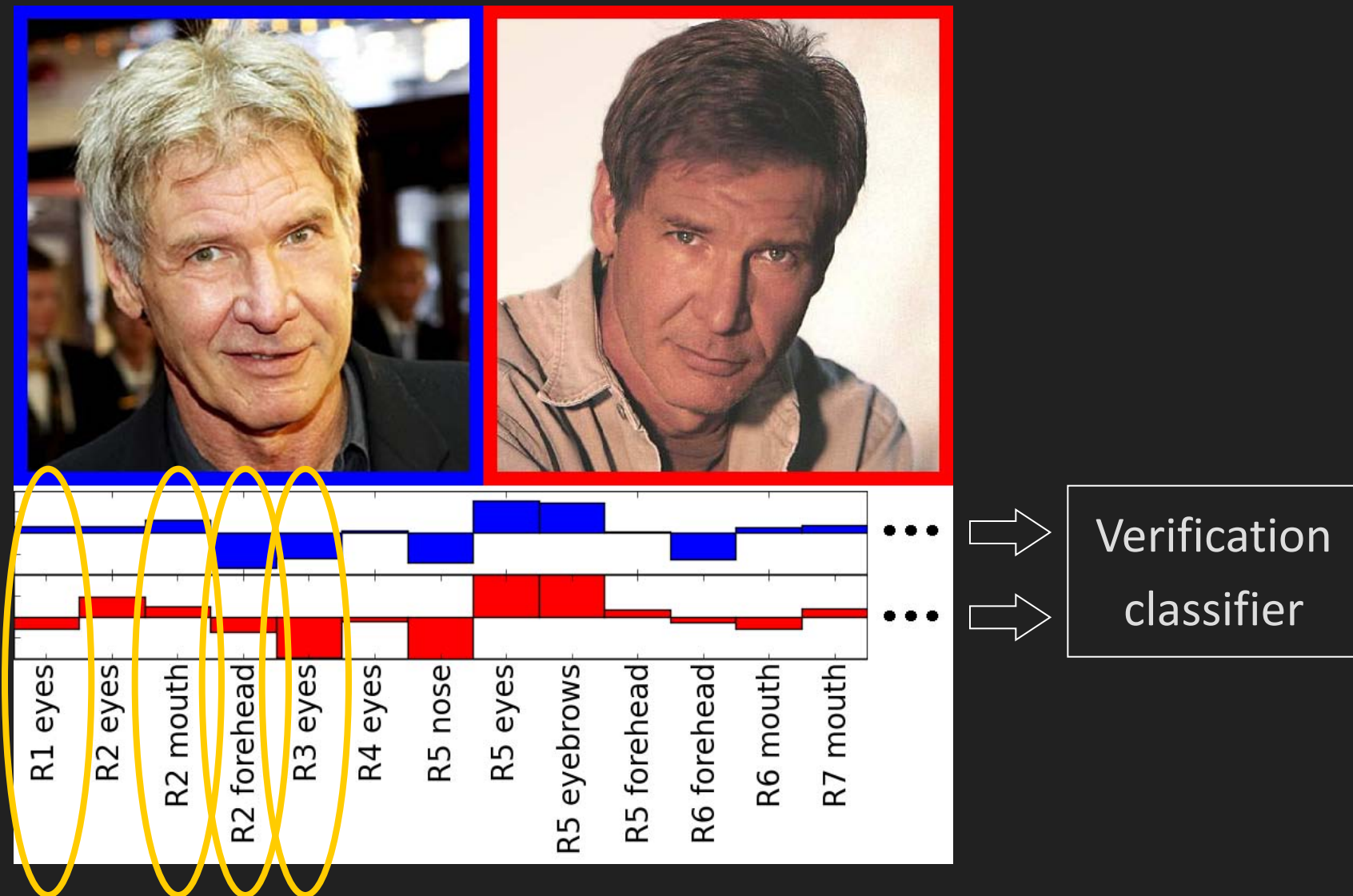
---



Images of other people 's eyes




# Using simile classifiers for verification



# Experimental evaluation

## LFW Image-Restricted Benchmark:

- 6,000 face pairs (3,000 same, 3,000 different)
- 10-fold cross-validation






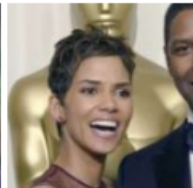
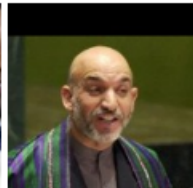

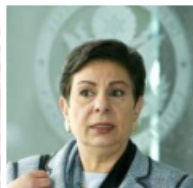

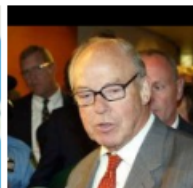

Menu

- LFW Home
- UMass Vision

### Labeled Faces in the Wild

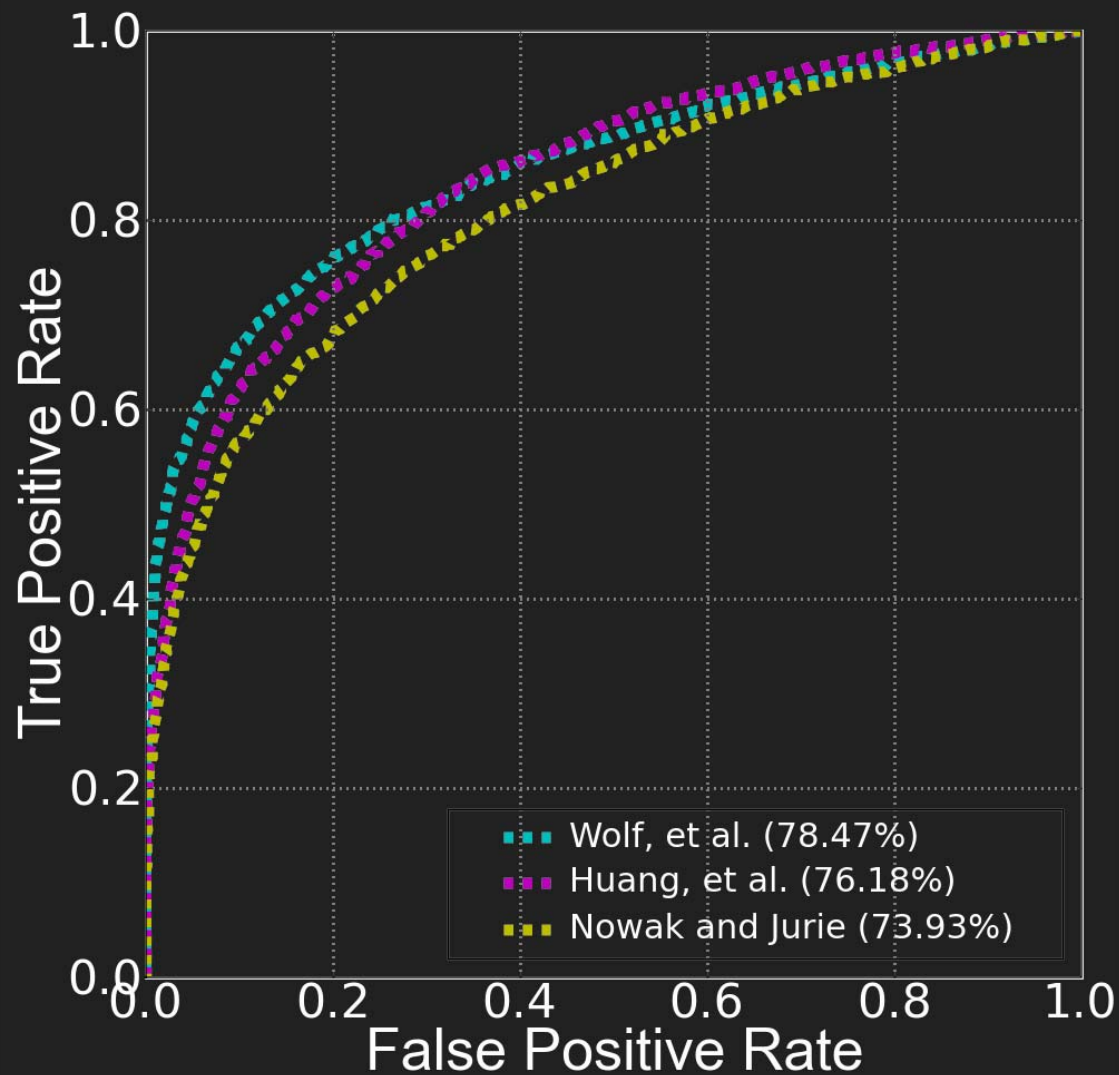
Database by name, non-singleton

[A][B][C][D][E][F][G][H][I][J][K][L][M][N][O][P][Q][R][S][T][U][V][W][X][Y][Z]

 Habib Rizieq (5)	 Hal Gehman (5)	 Hal Sutton (2)	 Halle Berry (16)	 Hamid Karzai (22)
 Hamzah Haz (2)	 Hanan Ashrawi (2)	 Hannah Stockbauer (2)	 Hans Blix (39)	 Hans Eichel (3)

<http://vis-www.cs.umass.edu/lfw>

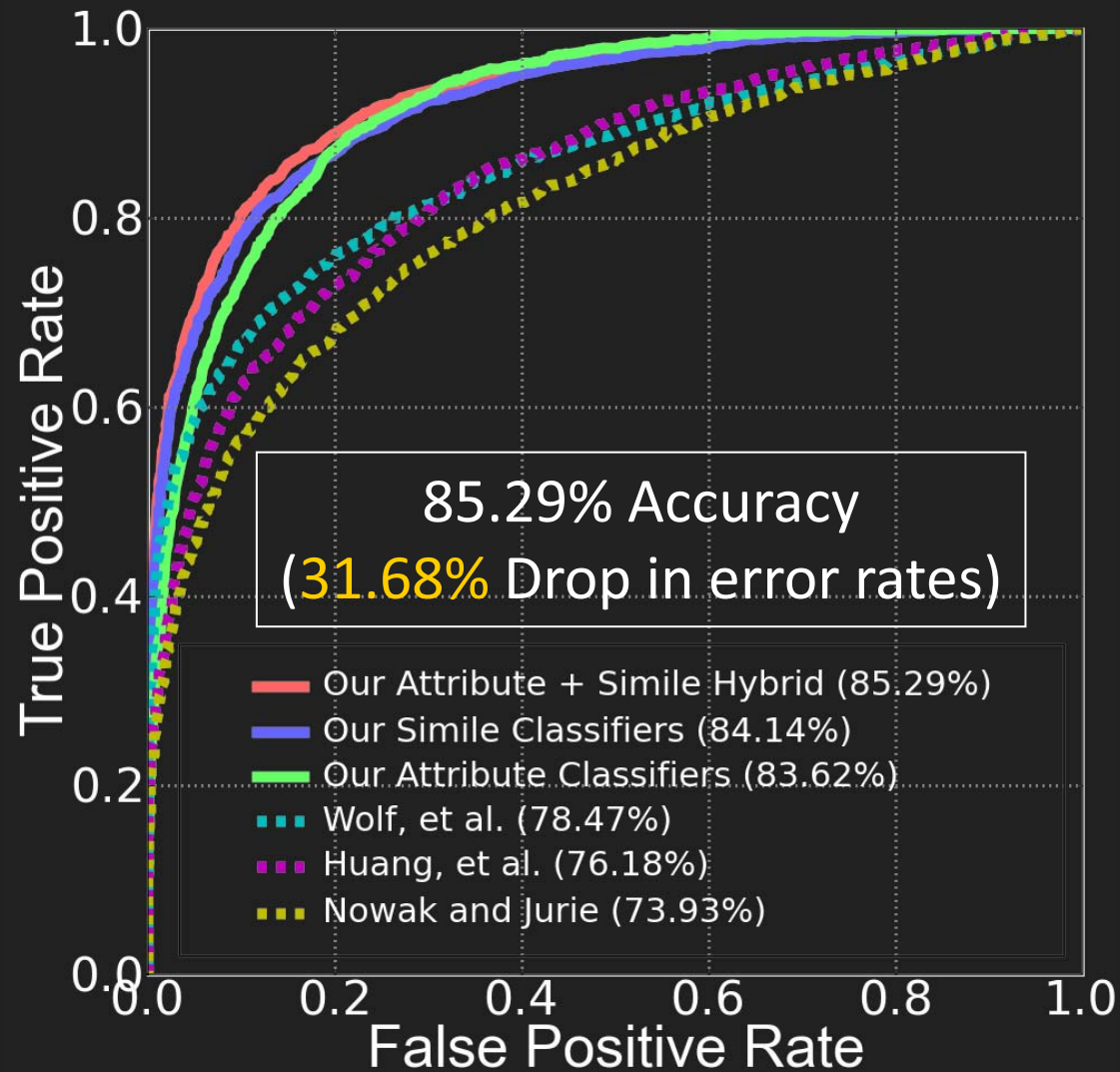
# Previous state-of-the-art on LFW



as of May 2009

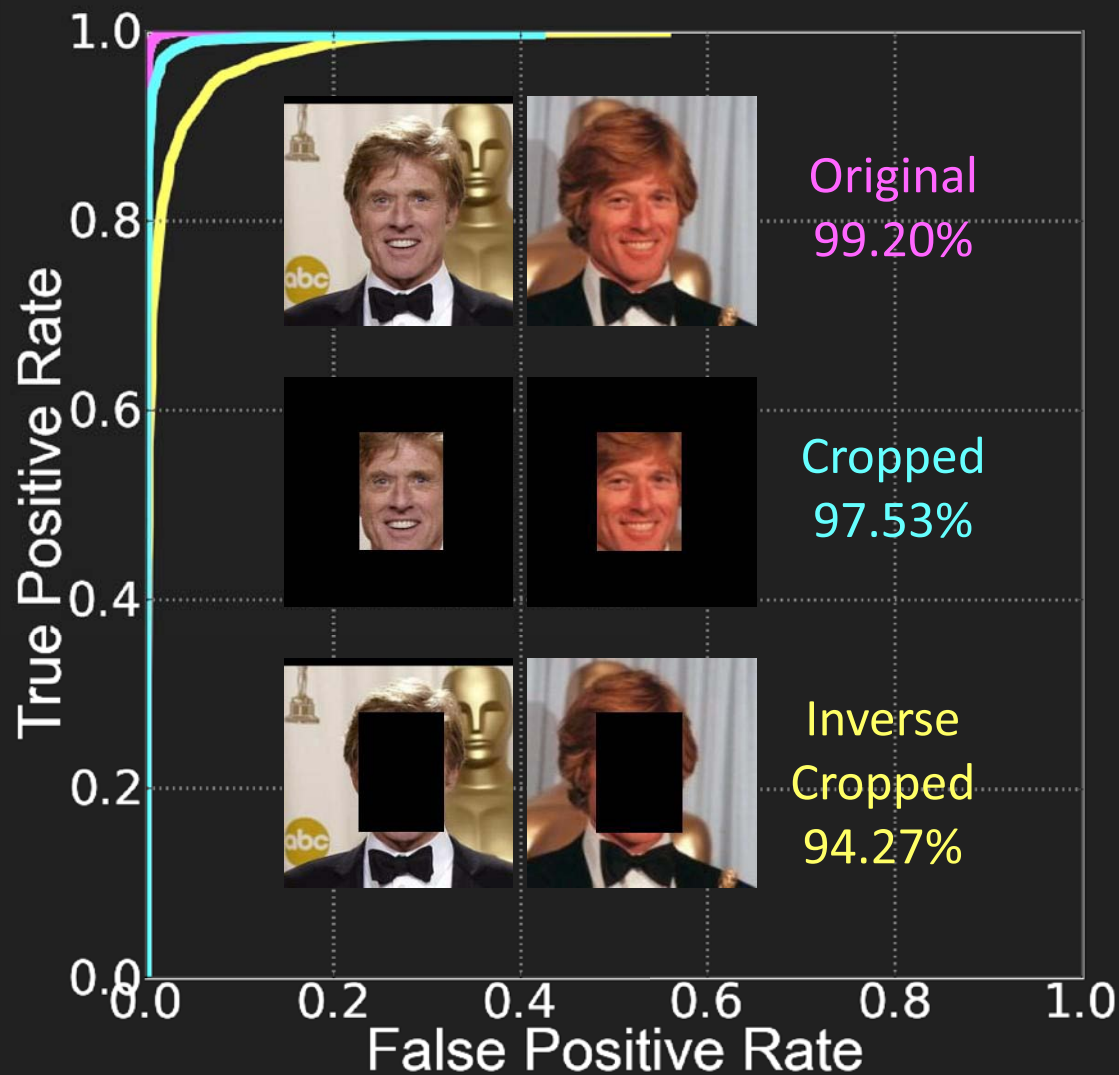


# Kumar et al. 2009 on LFW

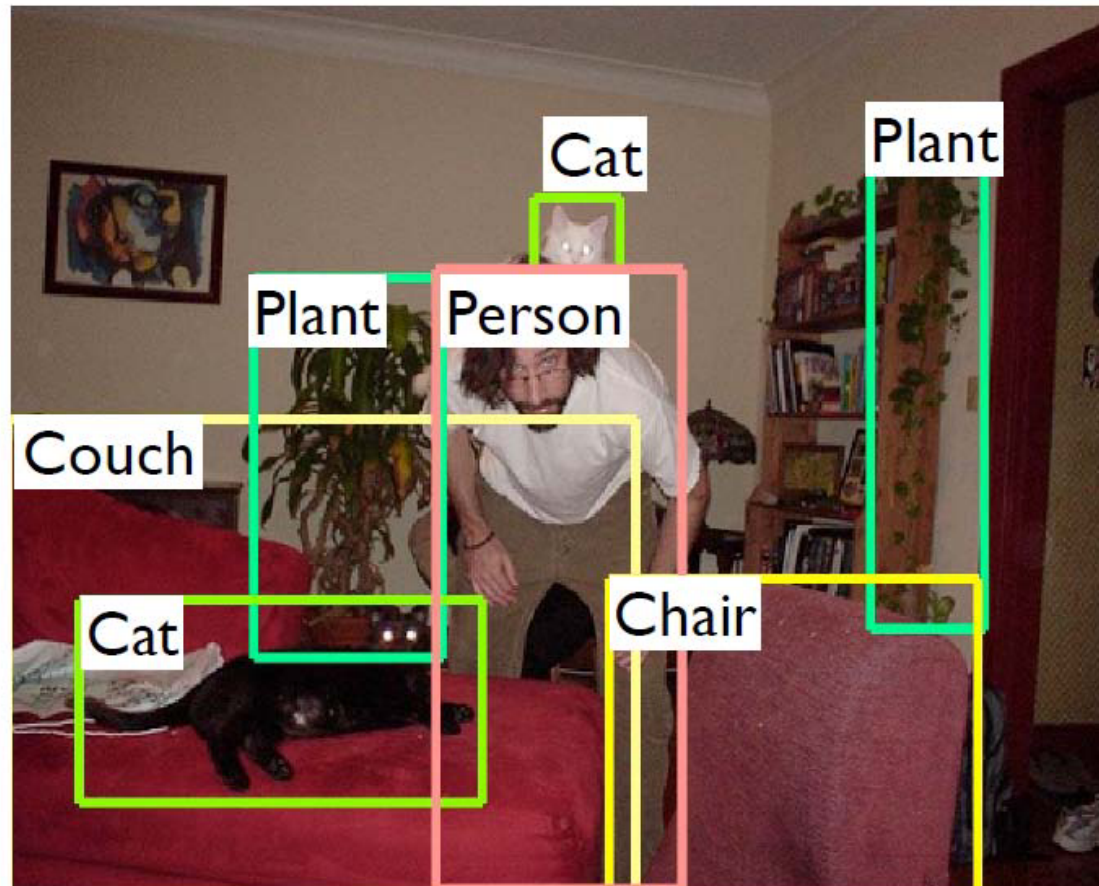


as of May 2009

# Human face verification performance

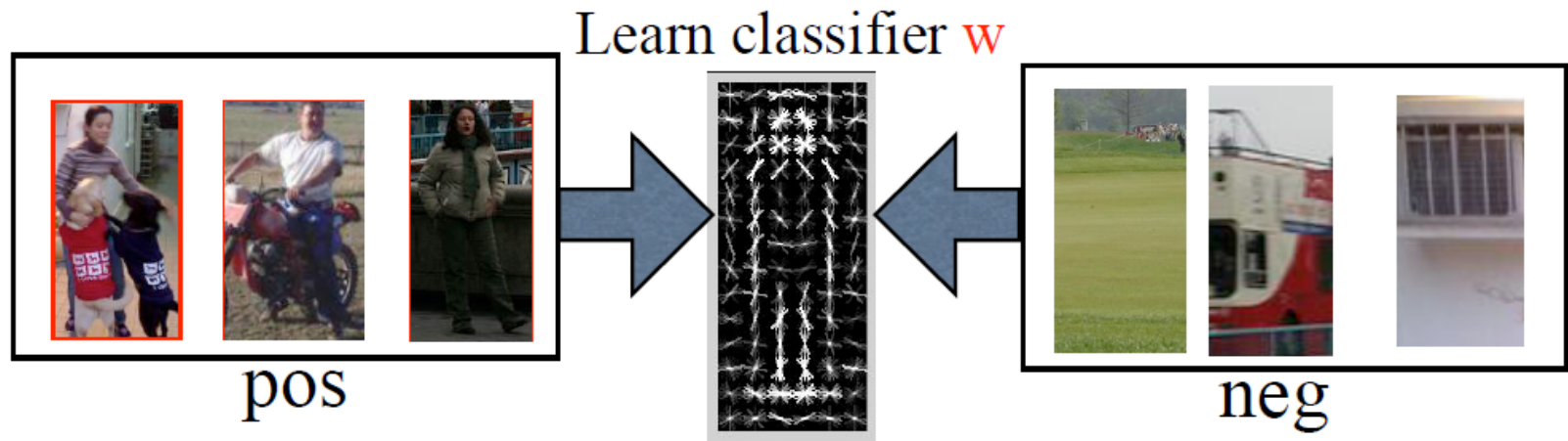


# What about multiple objects in the same image?



## Multiclass object detection

# Scanning-window pattern classification



$$w^T x > 0 ?$$

Face detection

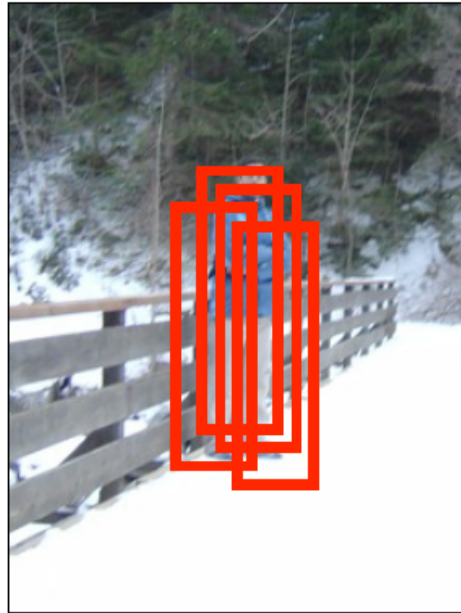
Rowley, Baluja, & Kanade. CVPR 96  
Viola & Jones IJCV 01



Pedestrian detection  
(and other objects)

Oren et al. CVPR 97  
Dalal & Triggs CVPR 05  
Felzenswalb et al. PAMI 09

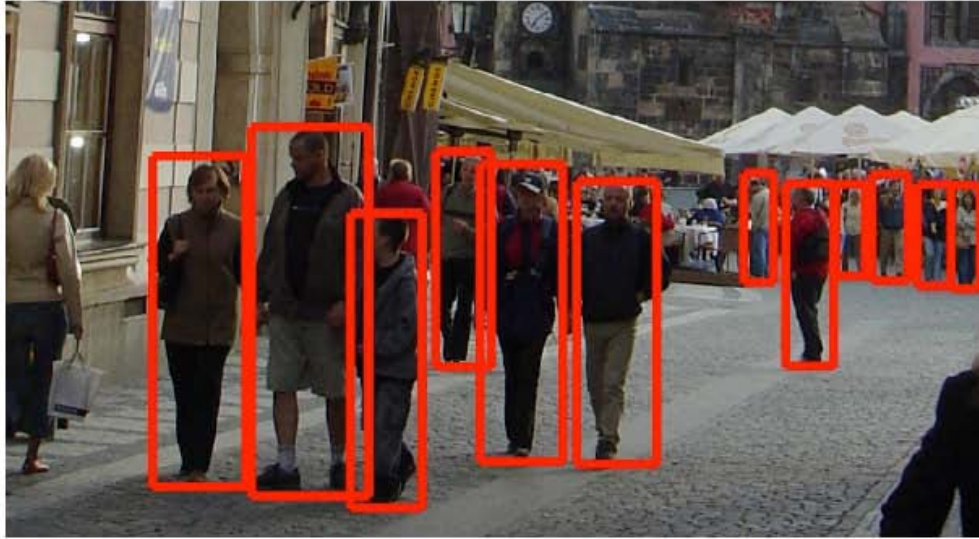
# Non-maxima suppression (NMS)



We need to suppress overlapping detections  
Many heuristics (mode finding, greedy selection)



# NMS in cluttered scenes

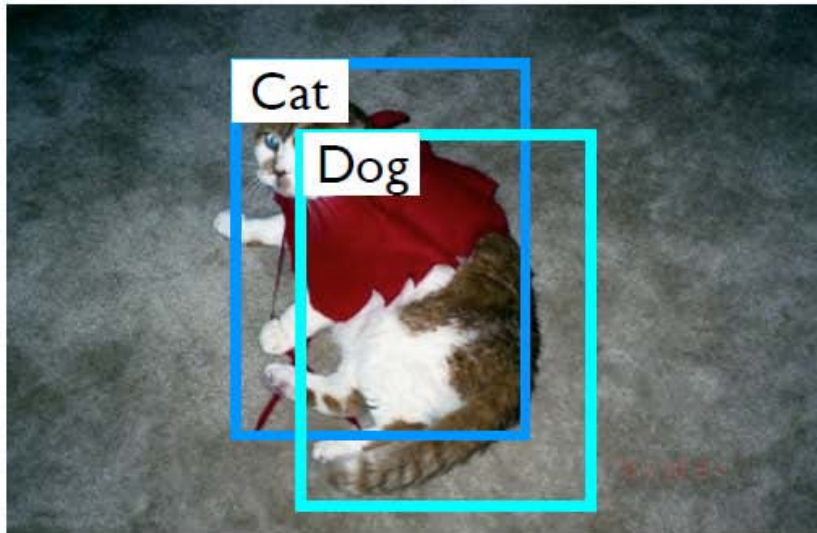


NMS should exploit spatial statistics of objects in real scenes

Is there a principled way to **learn** how to perform NMS?

# Inter-class NMS

Mutual exclusion: two objects cannot occupy the same 3D volume



May not be a strict constraint due to porous or transparent objects

# Taxonomy of spatial interactions

	within-class	between-class
negative	NMS	mutual exclusion
positive	textures of objects	spatial cueing

Most past work focuses on positive interactions,  
heuristically performing NMS & mutual exclusion.

Our contribution: a model for **all** of the above

Our inspiration: Torralba, Murphy, & Freeman NIPS 04

Kumar & Hebert ICCV 05

He, Zemel, & Carreira-Perpinan CVPR 05

Galleguillos, Rabinovich & Belongie CVPR 08

Hoeim, Efros, & Hebert IJCV 08



# Object detection as ....

## Classification



$x = \text{image window}$   
 $y \in \{0,1\}$

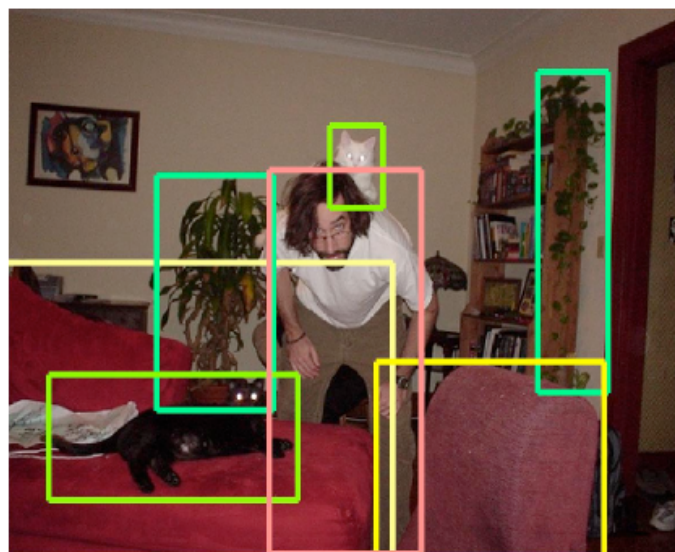
# Object detection as a structured labeling task

Classification



$x$  = image window  
 $y \in \{0,1\}$

Structured, sparse label



$X$  = entire image

$Y = [...4...3...2.7..1..]$

# Global scoring function

$$S_w(X, Y)$$



$$X = \{x_i\} \quad Y = \{y_i\}$$

$x_i$  = feature vector extracted from  $i^{\text{th}}$  window (e.g. HOG)

$y_i$  = class label (0..K) for  $i^{\text{th}}$  window

# Global scoring function

$$S_w(X, Y) =$$

$$X = \{x_i\} \quad Y = \{y_i\}$$

$x_i$  = feature vector  
extracted from  $i^{\text{th}}$   
window (e.g. HOG)

$y_i$  = class label (0..K)

$$\sum_i w_{y_i}^T x_i$$

$w_{y_i}$  = template for class  $y_i$

$$\begin{matrix} w_{car} \\ w_{bus} \\ \vdots \end{matrix}$$

sum of per-window classifier scores

# Global scoring function

$$S_w(X, Y) =$$

$$X = \{x_i\} \quad Y = \{y_i\}$$

$x_i$  = feature vector  
extracted from  $i^{\text{th}}$   
window (e.g. HOG)

$y_i$  = class label (0..K)

$$\sum_i w_{y_i}^T x_i$$

$w_{y_i}$  = template for class  $y_i$

$w_{car}$

$w_{bus}$

$\vdots$

sum of per-window classifier scores

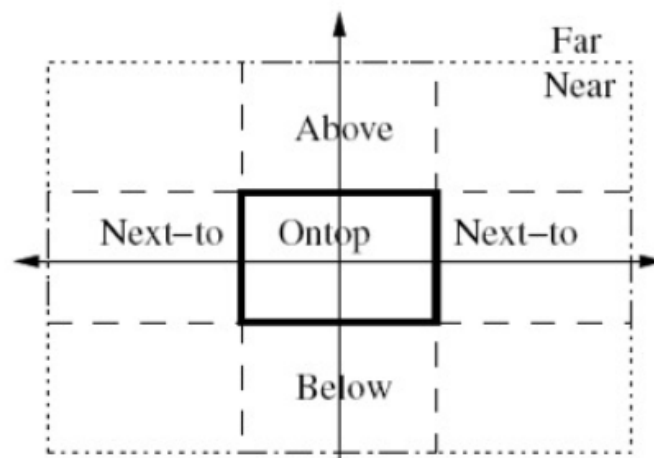
+

$$\sum_{i,j} w_{y_i y_j}^T d_{ij}$$

$d_{ij}$  = spatial context  
descriptor for  
window  $i$  and  $j$

$w_{y_i y_j}$  = spatial interaction  
model for  
class  $y_i$  &  $y_j$

sum of pairwise window-label interactions



# Inference

$$S_w(X, Y) = \sum_i w_{y_i}^T x_i + \sum_{i,j} w_{y_i, y_j}^T d_{ij}$$

$$L(X) = \operatorname{argmax}_Y S_w(X, Y)$$

Looks like an MRF - can we use standard inference techniques?

Our model is not sub-modular

Sub-modular interactions: neighboring labels should be similar

NMS interactions: neighboring labels should be **different**




# Greedy inference

$$L(X) = \operatorname{argmax}_Y S(X, Y) \quad S(X, Y) = \sum_i w_{y_i}^T x_i + \sum_{i,j} w_{y_i, y_j}^T d_{ij}$$

Analogous to common NMS schemes

(1) Initialize all labels to bg

Initialize per-window scores with local template

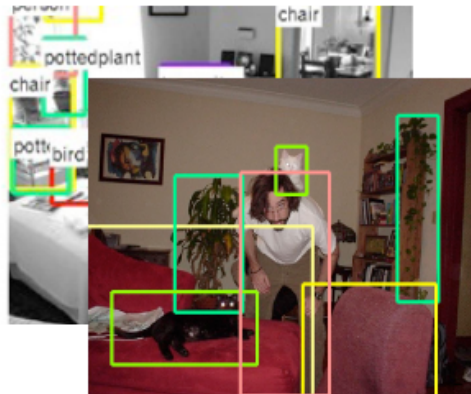
- 
- (2) Select highest scoring un-instanced window
- (3) Instance it and add pairwise contribution to remaining windows
- (4) Stop when remaining windows score  $< 0$

**Effectiveness:** Greedy solution close to optimal in practice

(See Numhauser et al. 78 for theoretical arguments)

# Learning

Training data consists of pairs of  $\{X_n, Y_n\}$



$$S_w(X, Y) = \sum_i w_{y_i}^T x_i + \sum_{i,j} w_{y_i, y_j}^T d_{ij}$$

$$S_w(X, Y) = w^T \Psi(X, Y)$$

# Learning with SVMs

$$\operatorname{argmin}_w \frac{1}{2} w^T w$$

$$\text{s.t.} \quad \forall n, H_n \neq Y_n \quad w^T \Psi(X_n, Y_n) - w^T \Psi(X_n, H_n) \geq 1$$

“Find a small  $w$  such that for each image, score of true label  $Y_n$  dominates all other hypothesized labels  $H_n$  by at least 1 unit”

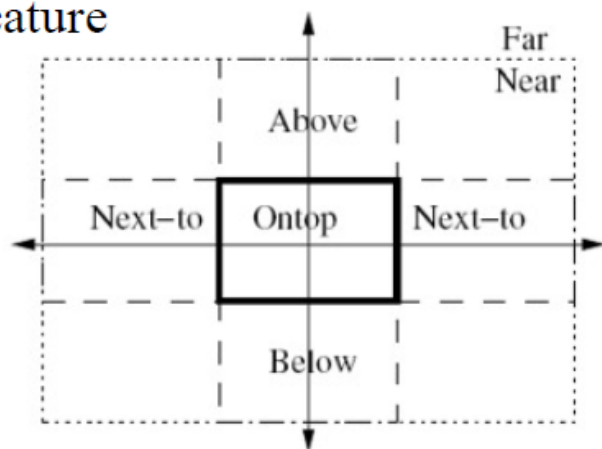


Only a tiny fraction of exponential number of constraints are necessary  
(i.e., support vectors)

Structured Prediction  
Tsochantaridis et al. ICML 04

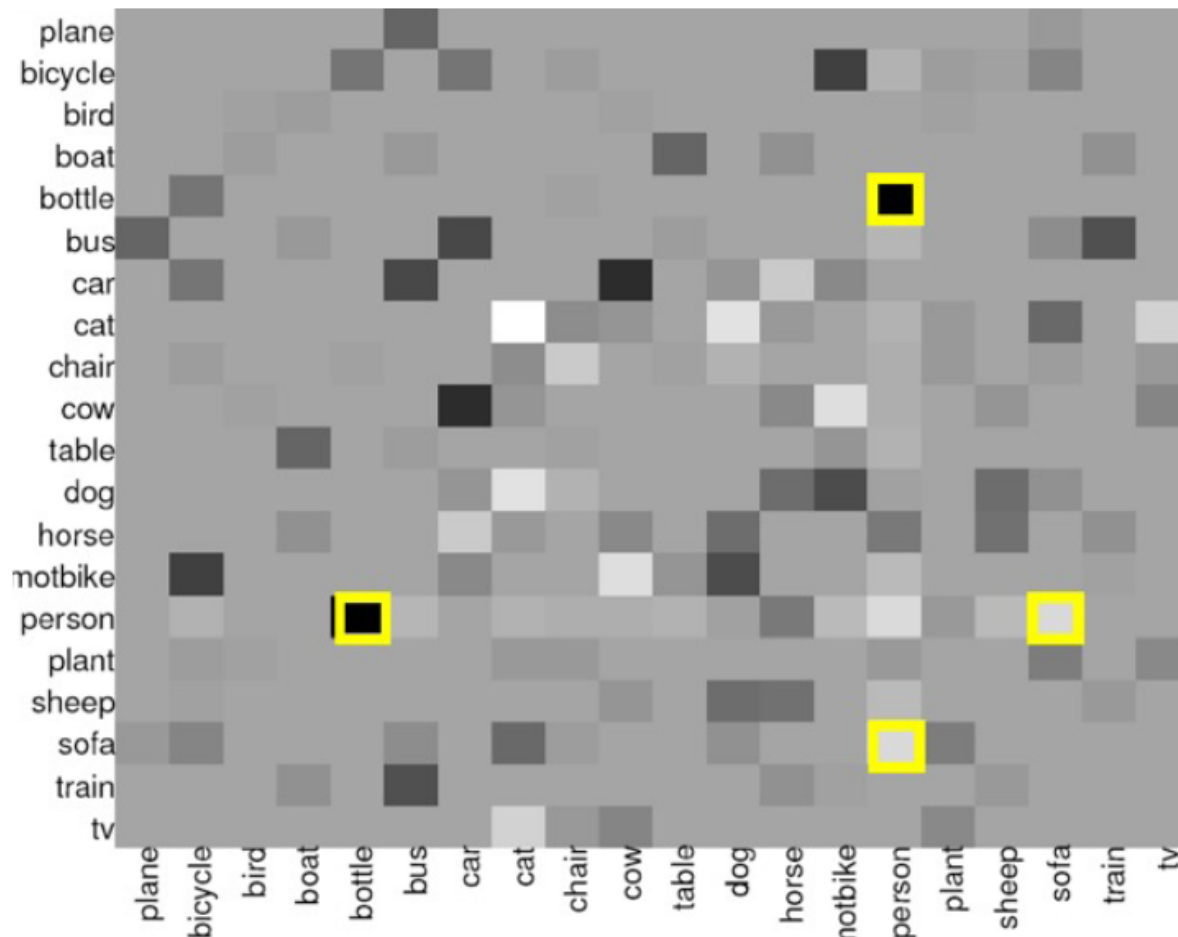
# Experiments

- 1) We use PASCAL 2007 training and test data  
20 classes, 5000 training images, 5000 test images
- 2) Baseline: Felzenswalb et al. PAMI 09 (with default NMS)
- 3) Local feature = [score of baseline detector 1]  
(We learn bias and offset for each local detector)
- 4) Pairwise feature



+ 50% overlap feature

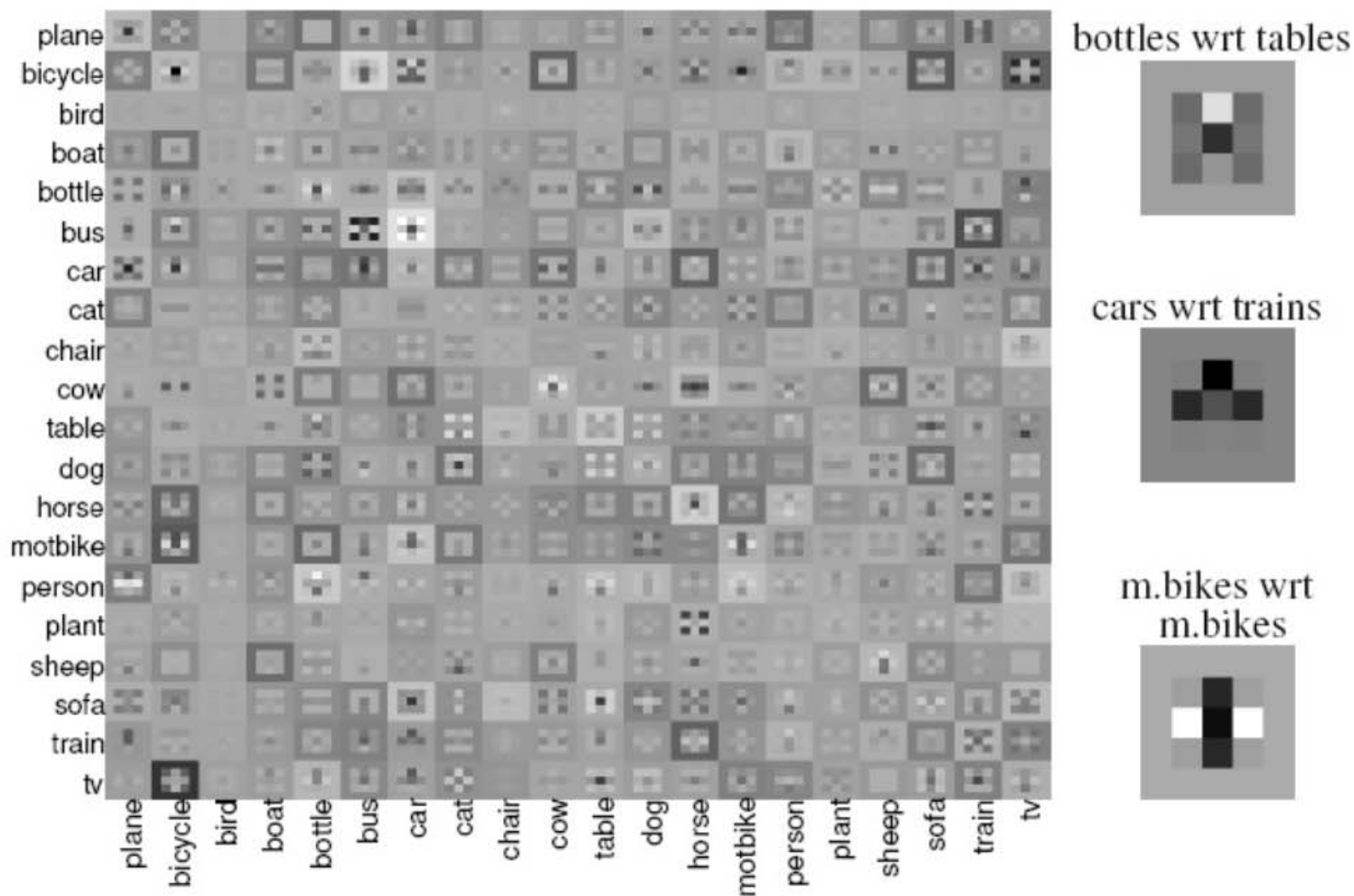
# Overlap feature in pairwise potential



Mutual exclusion can be subtle

Parameters are trained with knowledge of local detectors

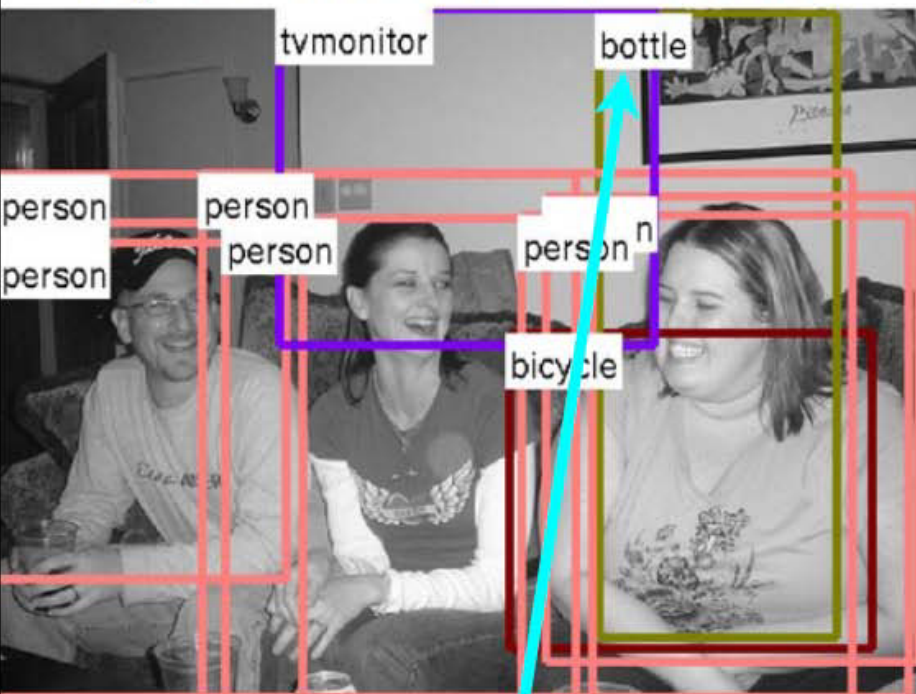
# Remaining pairwise potentials





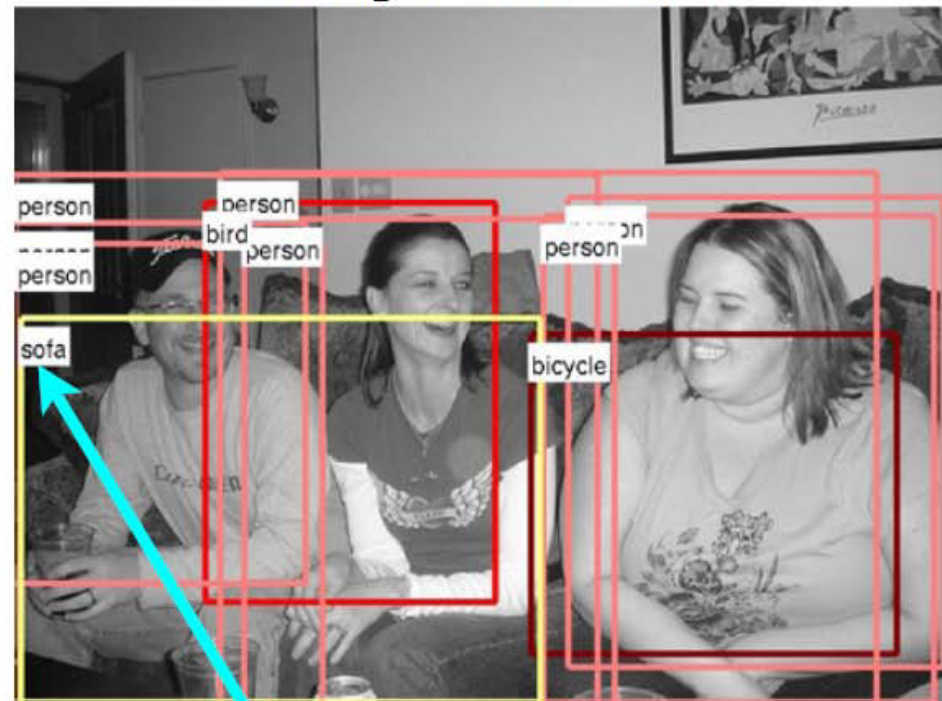
# Results

## Top 10 detections for baseline



Inhibit  
overlapping people & bottles  
because local detectors confuse them

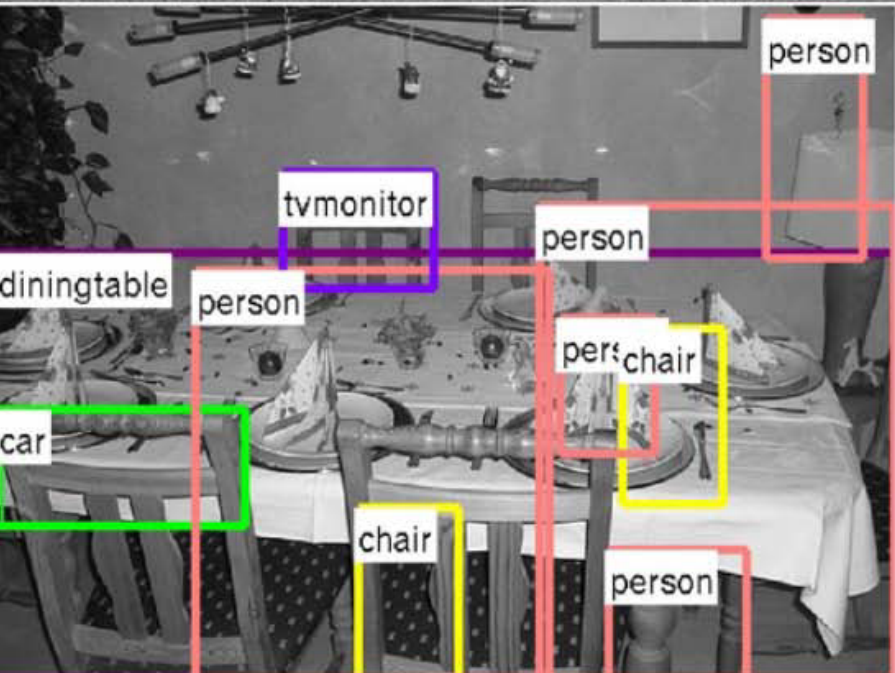
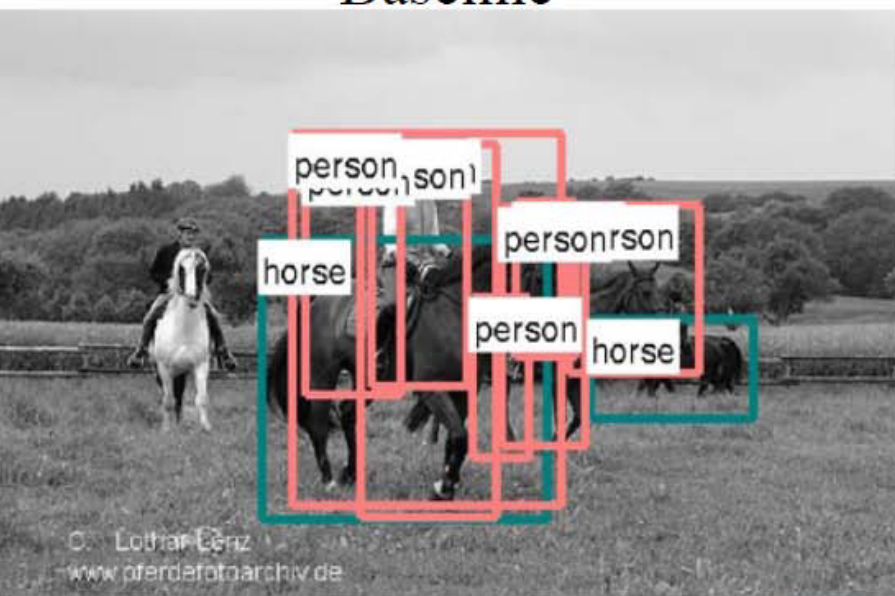
## Our top 10 detections



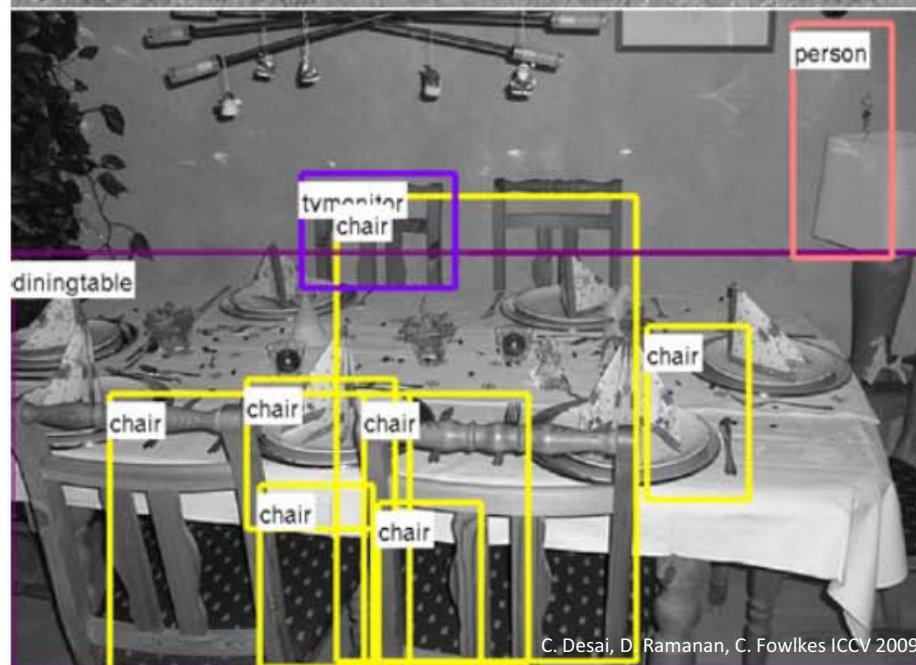
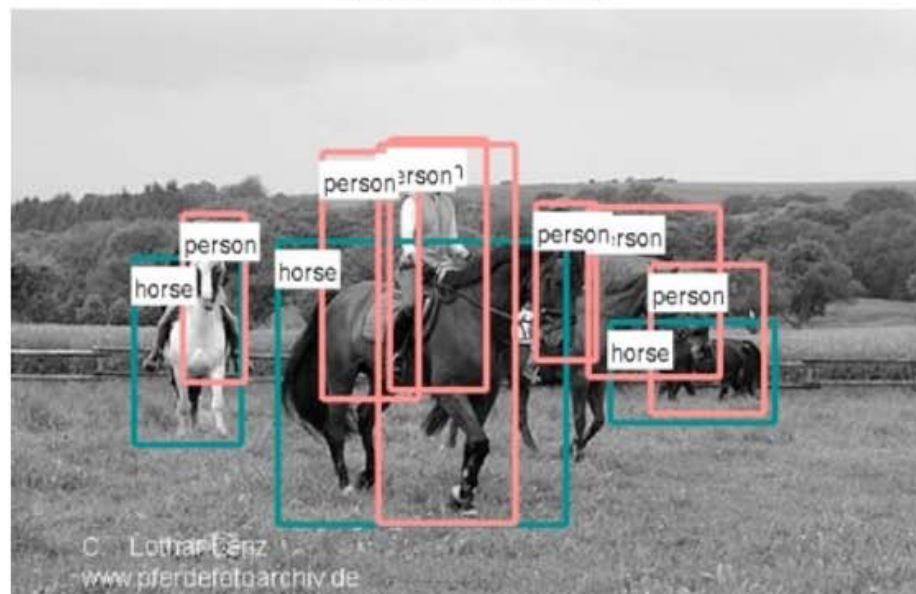
Favor  
overlapping people & sofas  
because people sit on sofas

# Results

## Baseline



## Our model



Default NMS heuristics

Default heuristics don't work for Mutual Exclusion

Winning  
PASCAL07  
score

Felzenszwalb et al.  
PAMI 09  
code

Mutual  
Exclusion

Our  
model

plane	.262	0.278	0.270	<b>0.288</b>
bike	.409	0.559	0.444	<b>0.562</b>
bird	<b>.098</b>	0.014	0.015	0.032
boat	.094	<b>0.146</b>	0.125	0.142
bottle	.214	0.257	0.185	<b>0.294</b>
bus	<b>.393</b>	0.381	0.299	0.387
car	.432	0.470	0.466	<b>0.487</b>
cat	<b>.240</b>	0.151	0.133	0.124
chair	.128	<b>0.163</b>	0.145	0.160
cow	.140	0.167	0.109	<b>0.177</b>
table	.098	0.228	0.191	<b>0.240</b>
dog	<b>.162</b>	0.111	0.091	0.117
horse	.335	0.438	0.371	<b>0.450</b>
motbike	.375	0.373	0.325	<b>0.394</b>
person	.221	0.352	0.342	<b>0.355</b>
plant	.120	0.140	0.091	<b>0.152</b>
sheep	<b>.175</b>	0.169	0.091	0.161
sofa	.147	0.193	0.188	<b>0.201</b>
train	.334	0.319	0.318	<b>0.342</b>
TV	.289	<b>0.373</b>	0.359	0.354

Our model outperforms Felzenszwalb et al.'s baseline for most classes

Yet, another multi-object detection problem





## Density-aware person detection and tracking in crowds

M. Rodriguez, I. Laptev, J. Sivic, J.-Y. Audibert



# Motivation

- Recognize crowd events
- Predict future [potentially dangerous] events
  - ▣ Detect and track individual people



# Problem

- As the density of people in a scene increases:
  - ▣ The accuracy of current detection and tracking methods degrades



Increasing person density

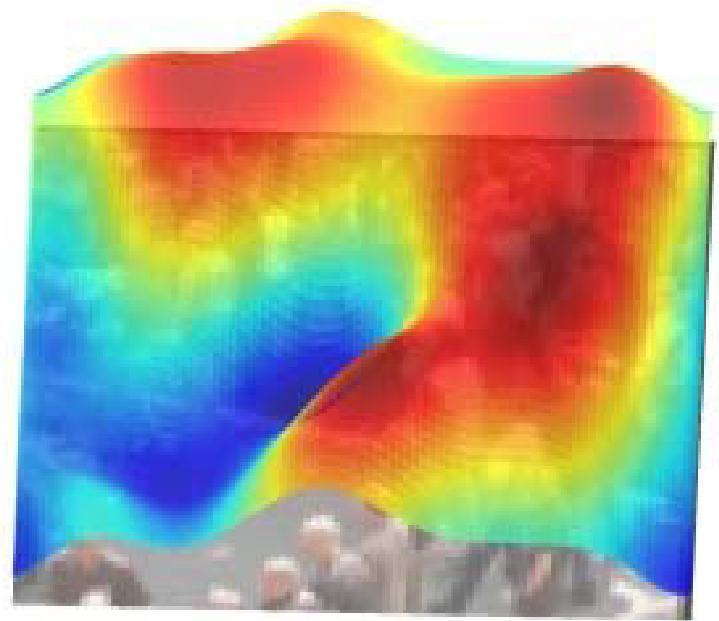


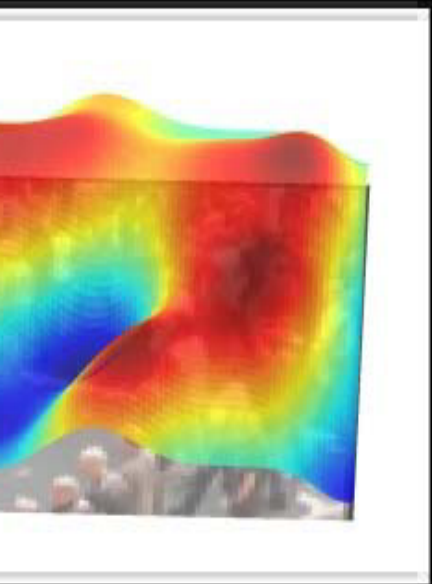
# Detection and Tracking



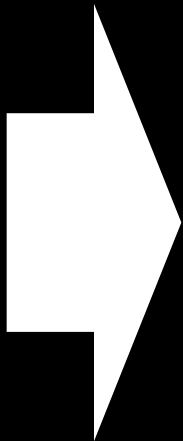
Typical state-of-the-art detection and tracking

# Density Estimation





Density estimate



Improved Detection and Tracking





# Crowd Model

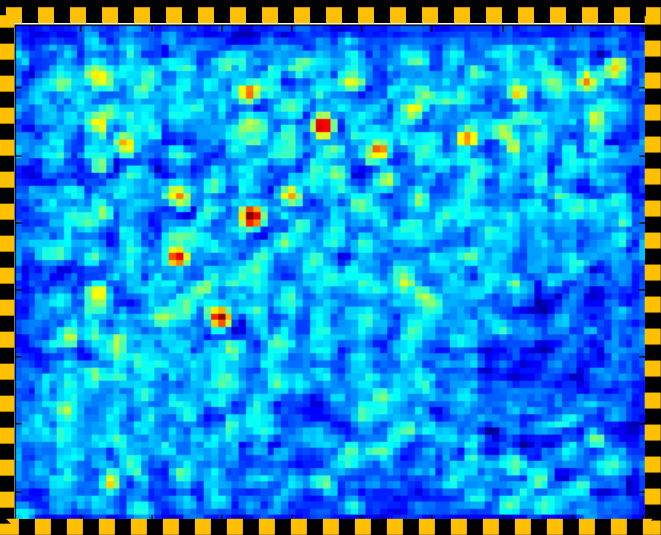
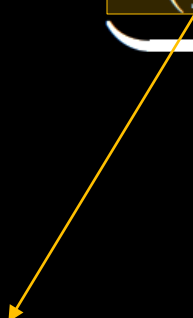
Energy Formulation



$$\min_{\mathbf{x} \in \{0,1\}^N} \underbrace{-s(\mathbf{p})^\top \mathbf{x}}_{E_S} + \underbrace{\mathbf{x}^\top W \mathbf{x}}_{E_P} + \underbrace{\alpha \|D(\mathbf{p}) - A\mathbf{x}\|_2^2}_{E_D}$$

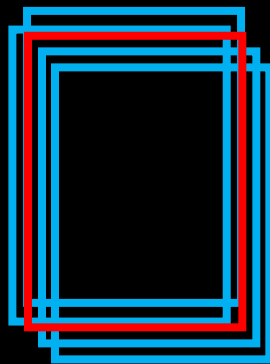
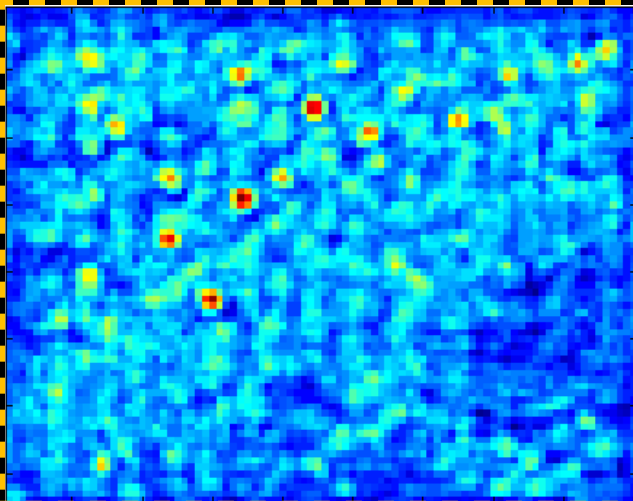
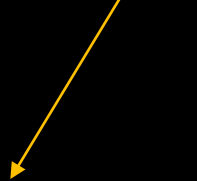


$$\min_{\mathbf{x} \in \{0,1\}^N} \underbrace{-s(\mathbf{p})^\top \mathbf{x}}_{E_S} + \underbrace{\mathbf{x}^\top W \mathbf{x}}_{E_P} + \underbrace{\alpha \|D(\mathbf{p}) - A\mathbf{x}\|_2^2}_{E_D}$$



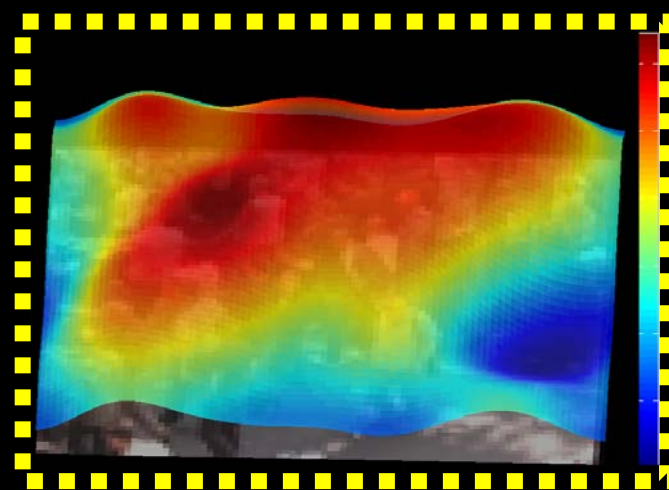
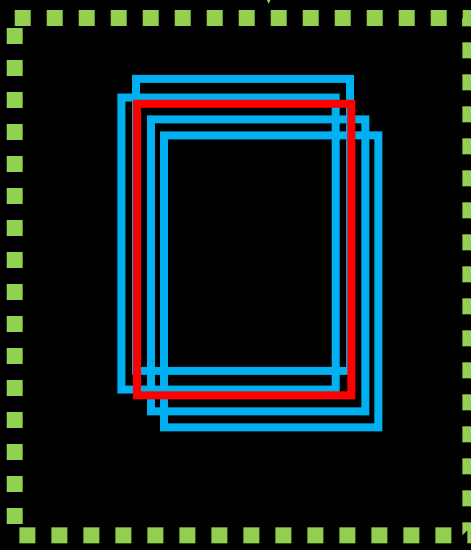
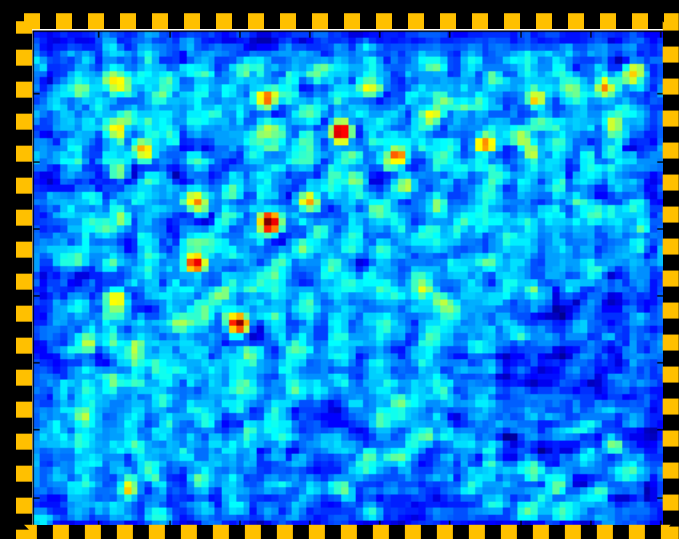


$$\min_{\mathbf{x} \in \{0,1\}^N} \underbrace{-s(\mathbf{p})^\top \mathbf{x}}_{E_S} + \underbrace{\mathbf{x}^\top W \mathbf{x}}_{E_P} + \underbrace{\alpha \|D(\mathbf{p}) - A\mathbf{x}\|_2^2}_{E_D}$$





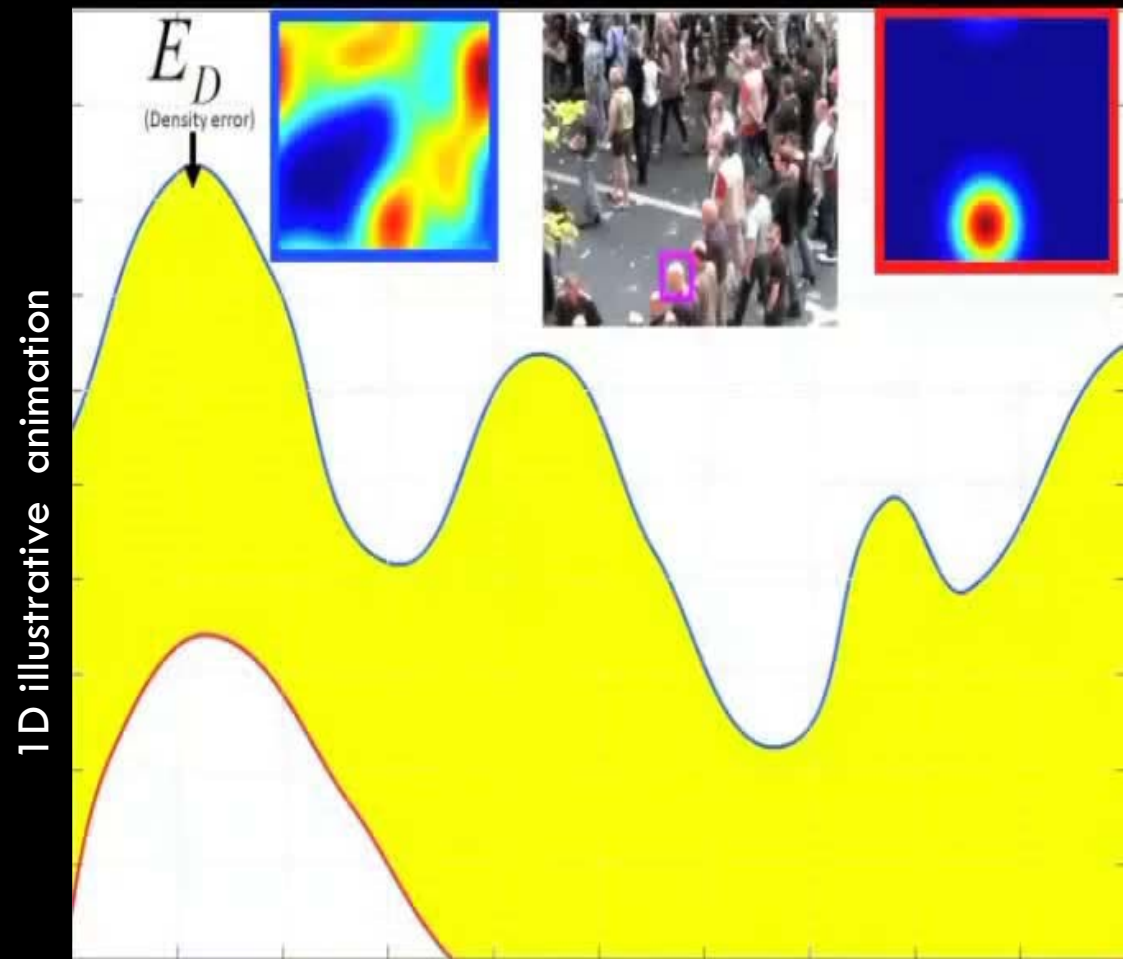
$$\min_{\mathbf{x} \in \{0,1\}^N} \underbrace{-s(\mathbf{p})^\top \mathbf{x}}_{E_S} + \underbrace{\mathbf{x}^\top W \mathbf{x}}_{E_P} + \underbrace{\alpha \|D(\mathbf{p}) - A\mathbf{x}\|_2^2}_{E_D}$$





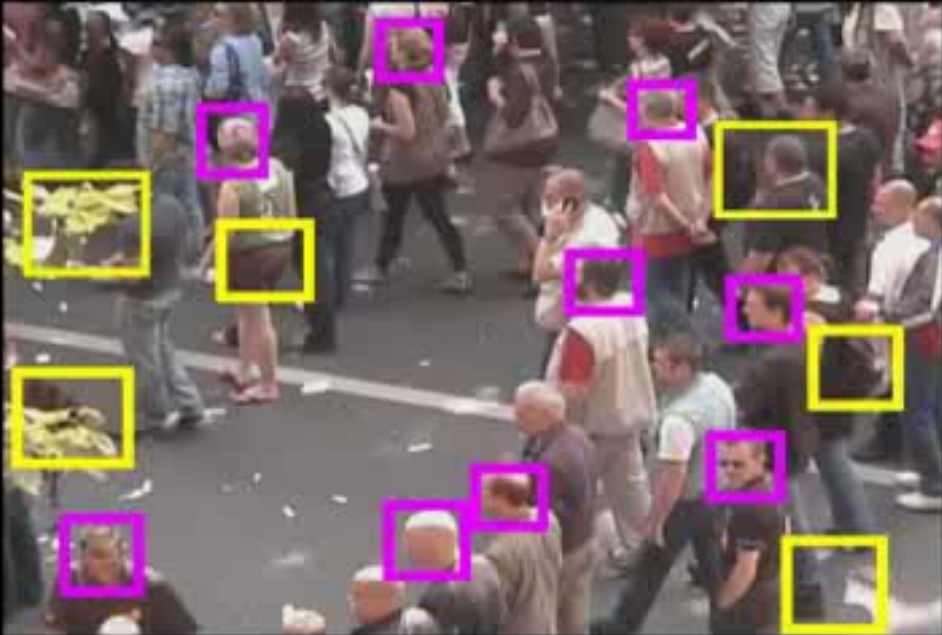
# Optimization

- NP-hard problem
- We adopt a greedy search procedure

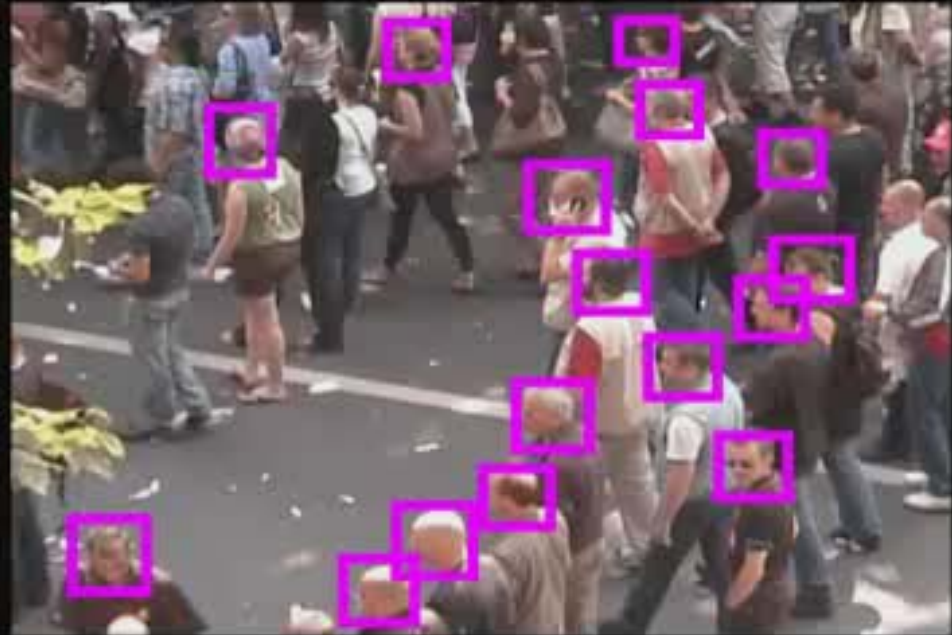




# Results and Evaluation



Baseline



Density-aware

Incorrect  
track



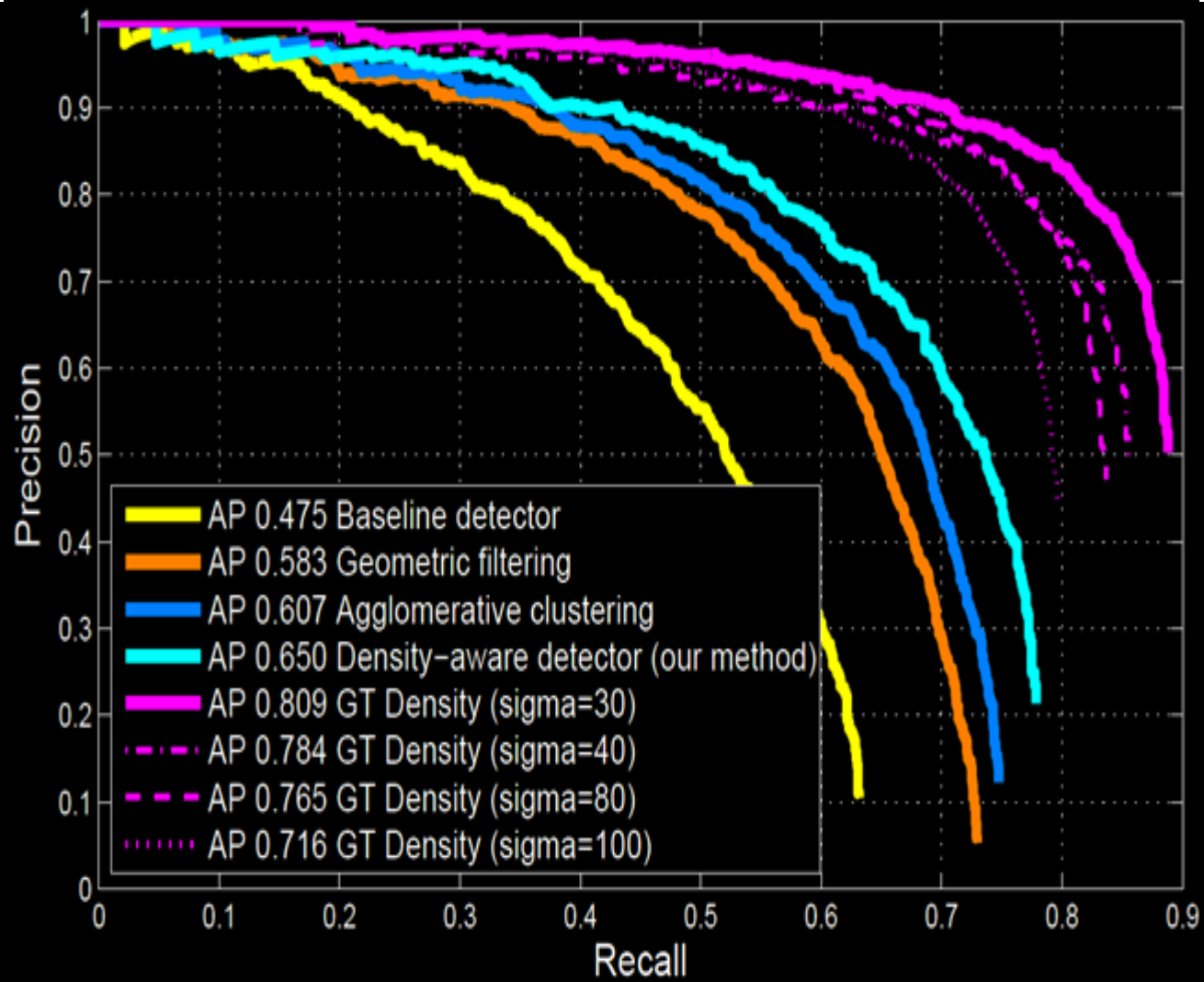
Baseline



Density-aware

Incorrect  
track

# Detection Evaluation





# Innrernship Topic 1

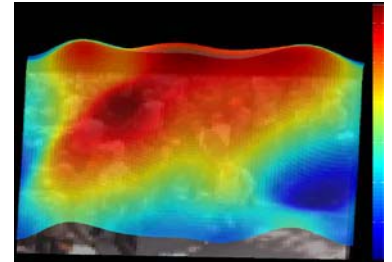
## Person Detection and Tracking in Crowds

Willow team, advisors: Ivan Laptev and Josef Sivic

<http://www.di.ens.fr/willow/teaching/recvis11/internships>

### Goals:

- Improve person detection and tracking with *better person density estimation*.
- Leverage space-time constraints with *density-aware track detection*.
- Address new research problems such as *person detection and tracking in low resolution video*



### Potential outcomes:

- Impact on the very active and challenging research domain
- Conference/journal publication and the start of a PhD thesis



# What to do about The Object That Cannot Be Named?



Photo by Daniel J. Cox

Slides by Derek Hoiem  
Computer Science Department  
University of Illinois at Urbana- Champaign

A. Farhadi, I. Endres, and D. Hoiem 2010



# A failure/success story



# Dealing with inevitable failure

Failure in categorization should not mean  
failure in recognition

# What to do about the **Object That Cannot Be Named?**





# Example

## Assisted Driving



# Example

## Security

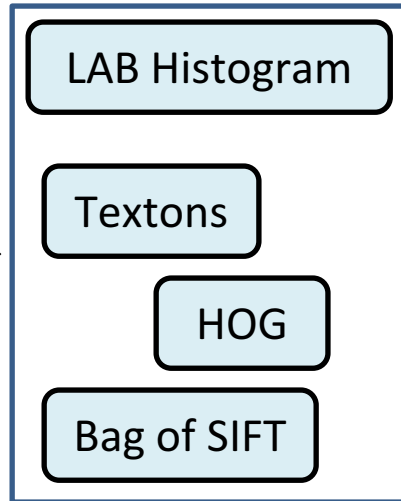


# Current View of Recognition

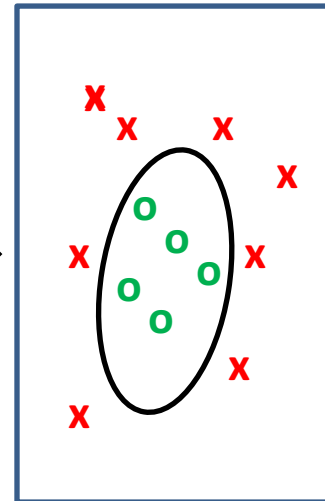
Training Examples



Appearance Representation



Appearance Model



Object Representation

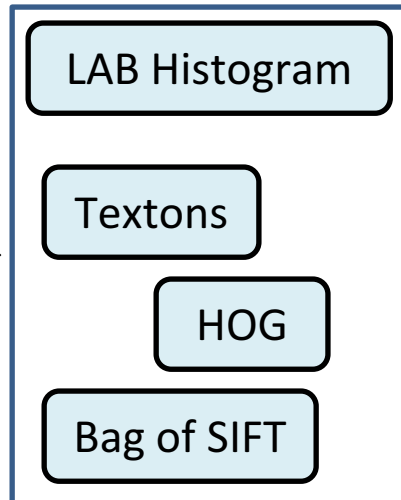


# Current View of Recognition

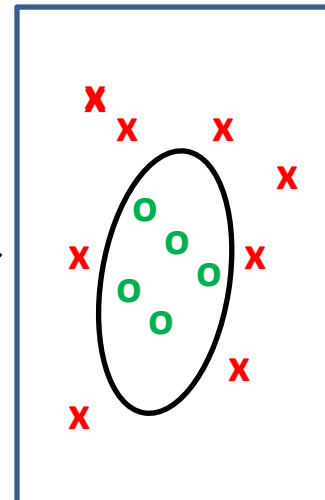
Training Examples



Appearance Representation



Appearance Model



Object Representation

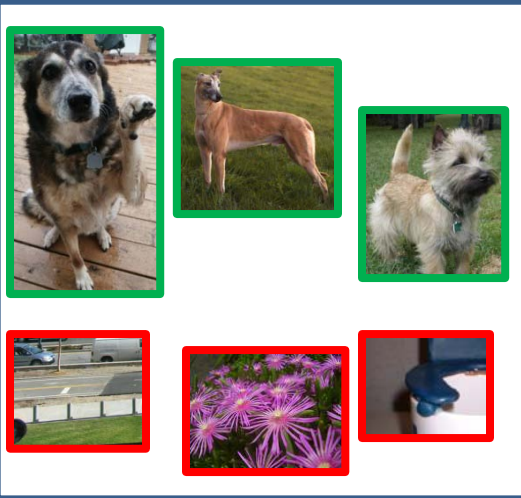


**Lots of effort – fancy stuff**

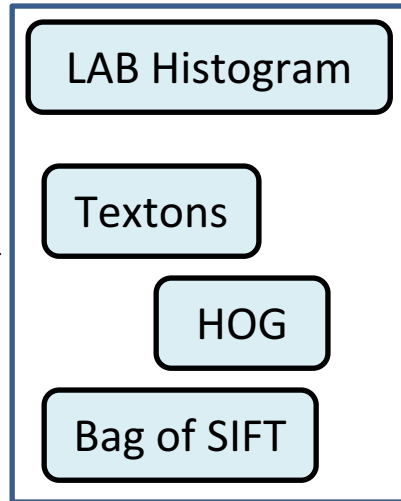


# Current View of Recognition

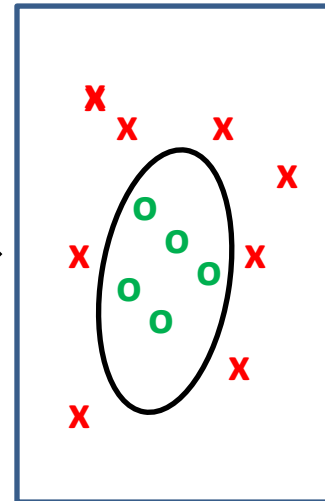
Training Examples



Appearance Representation



Appearance Model



Object Representation



**Not much changed**

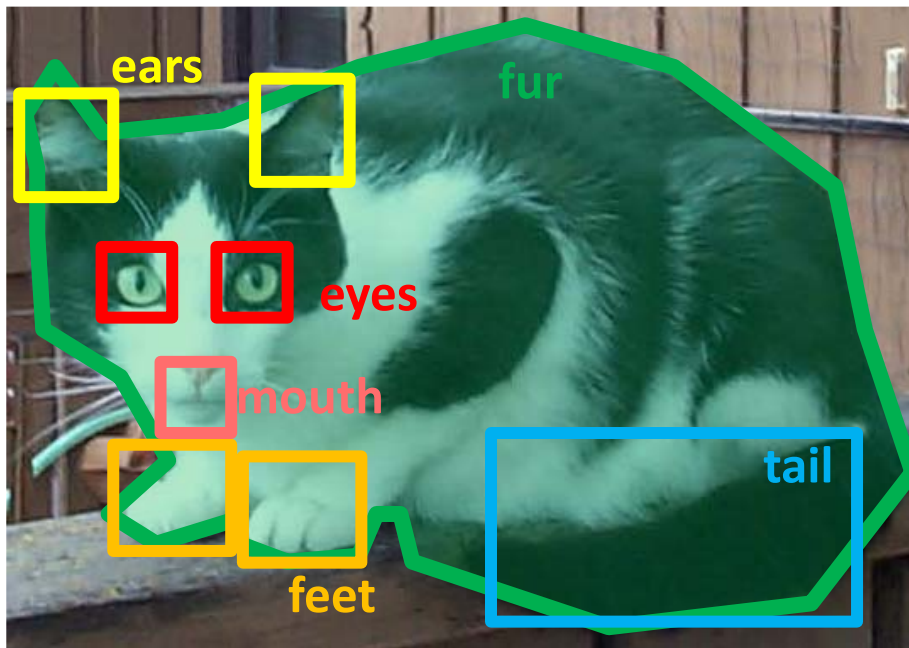
# Category-based representation

- Limited description and prediction
- No generalization to objects outside of learned categories
- Provides little guidance for learning

**So what would make a better representation?**

# Attribute-based Representation

- Properties that we want to describe or predict
- Shared across basic categories
- Made explicit through supervision



## Multiple Categories

animal, land animal, ..., cat

## Viewpoint/pose

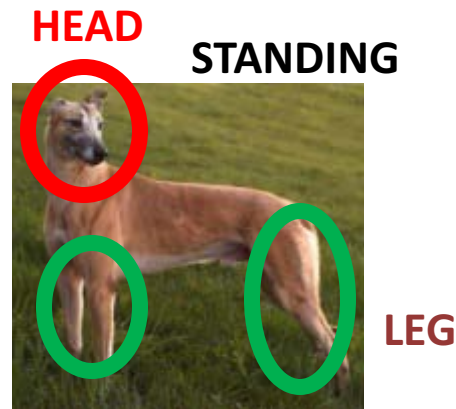
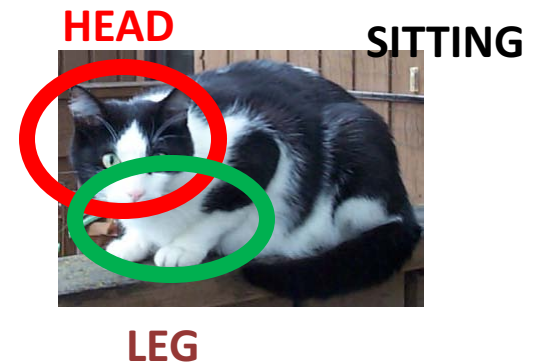
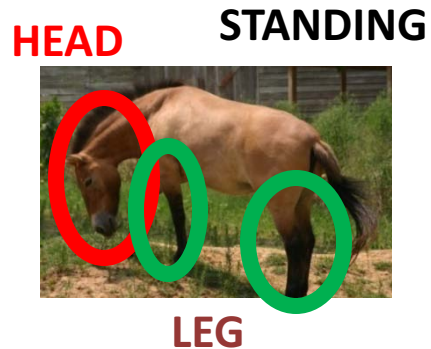
lying down, left side, facing camera

## Function

fast runner, climb trees, eat small animals, jump high, household pet, scratch

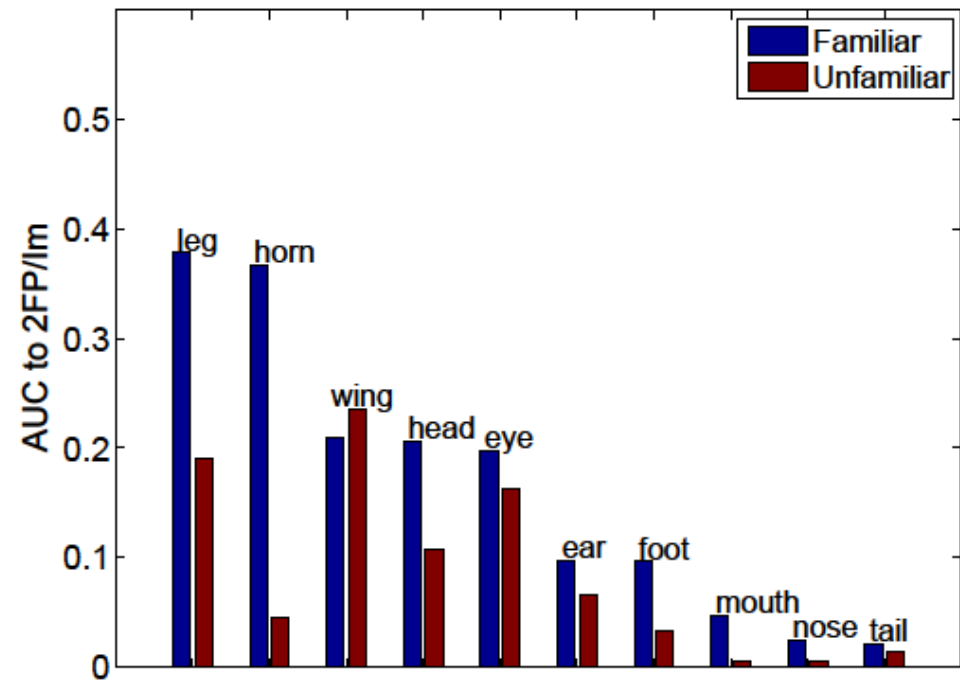
# Advantages of supervised attributes

- Provides correspondence for objects from different categories

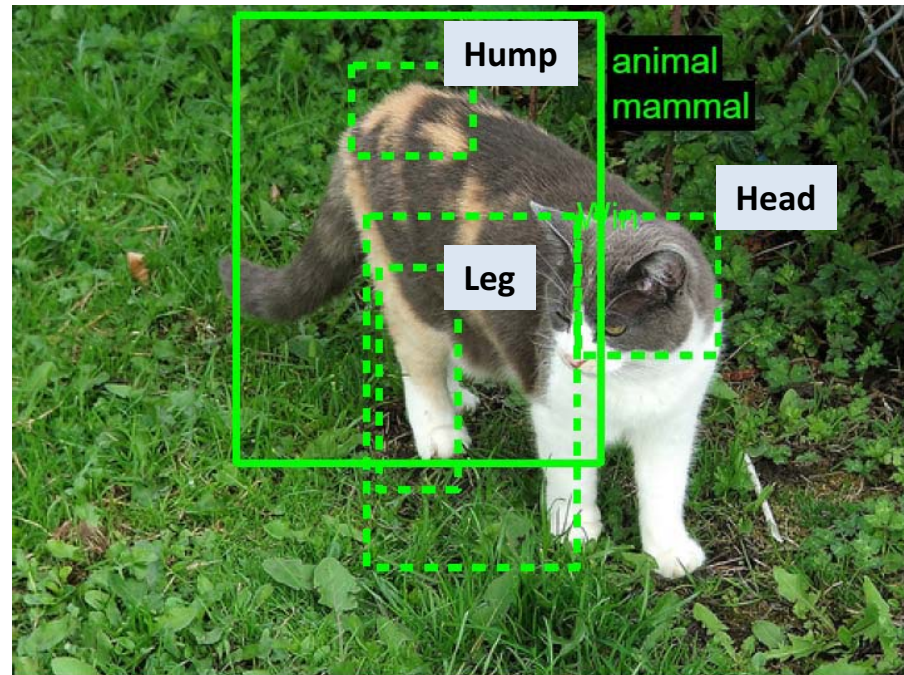


# Result: Part detectors can generalize across categories

Animal Parts



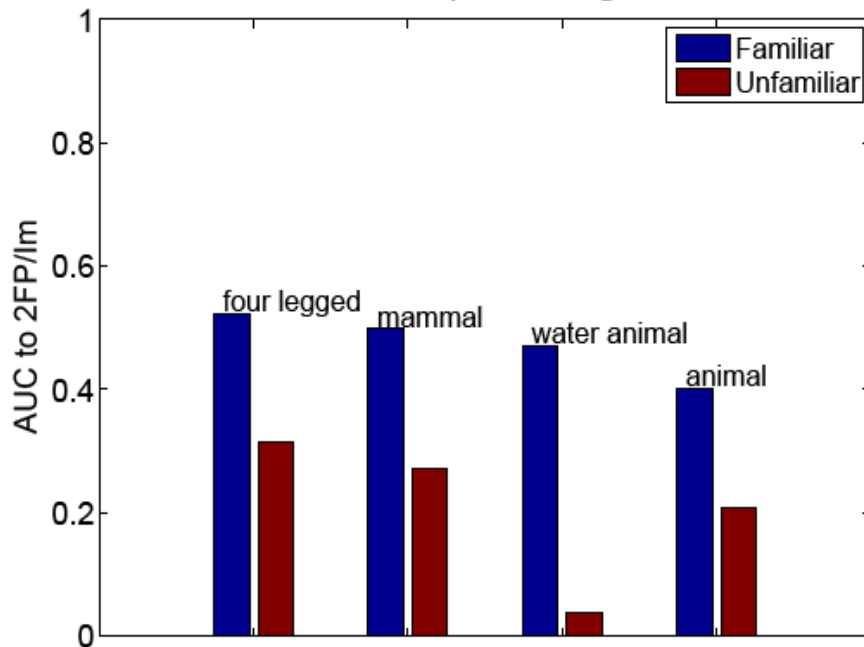
Part Detections for Novel Object



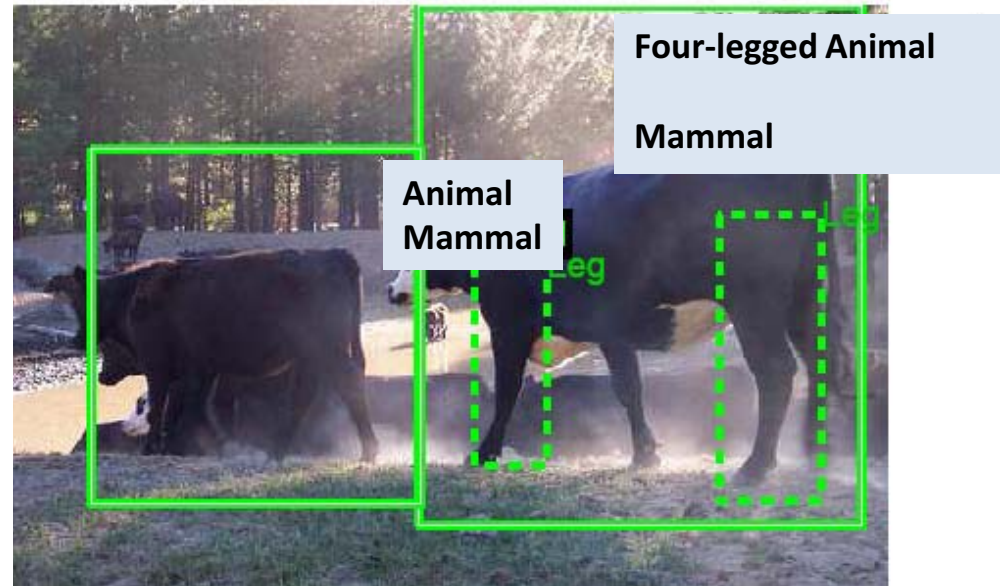


# Result: Broad category detectors can generalize across basic categories

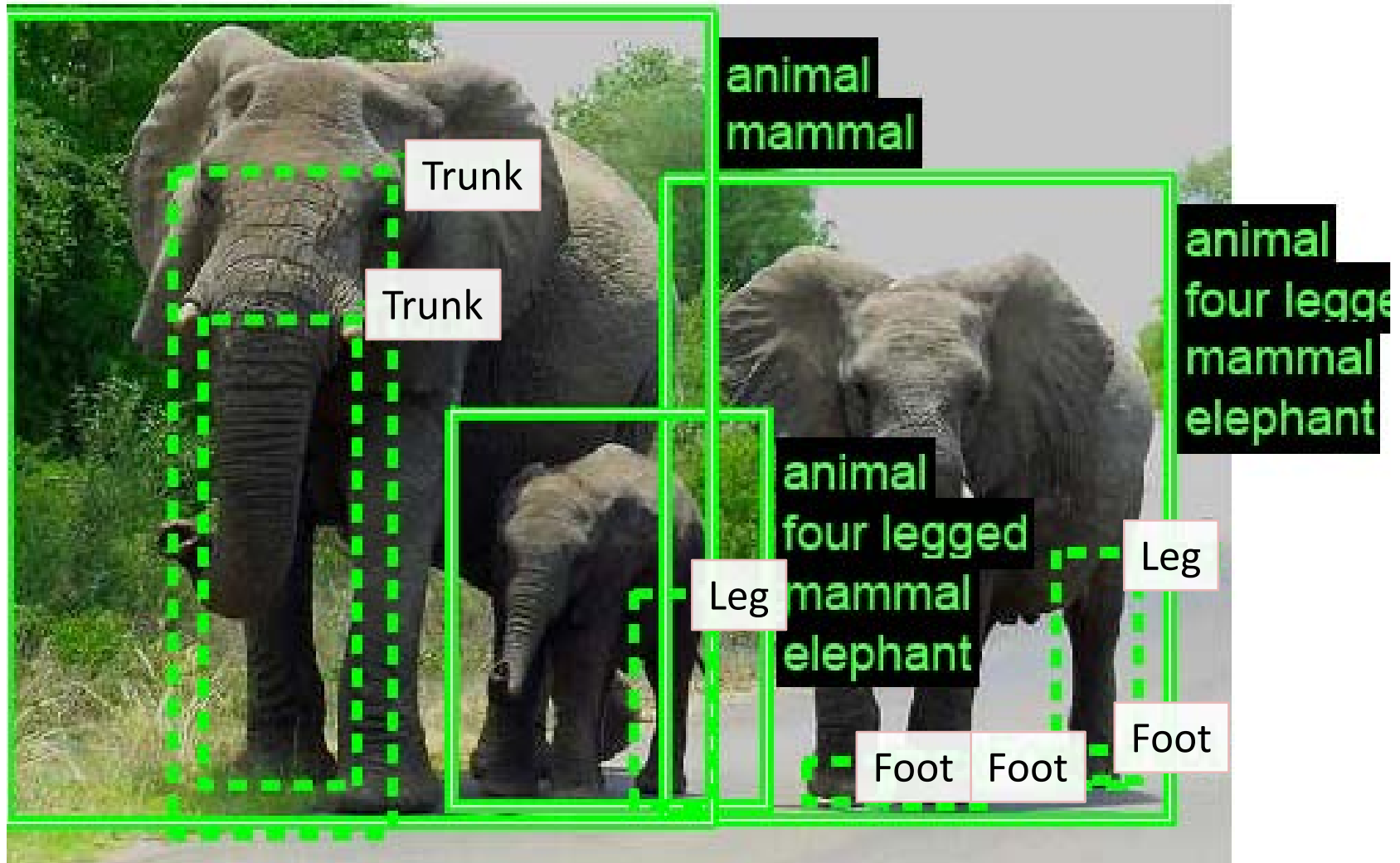
Animal SuperCategories



Category Detections for Novel Object

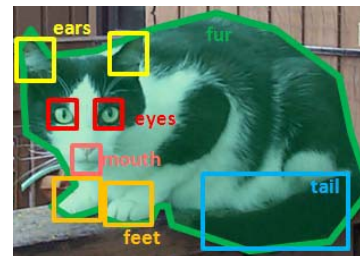
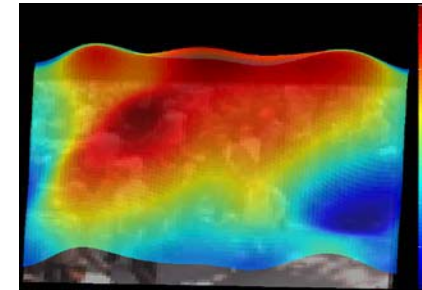
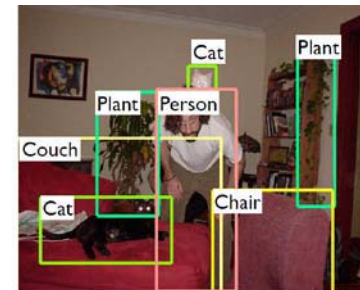


Result: We can better find and describe objects from familiar categories



# What we have seen so far?

- Objects in the context of scenes.
- Objects in relation with each other
- Objects defined by parts and attributes



Is this the end of the story?



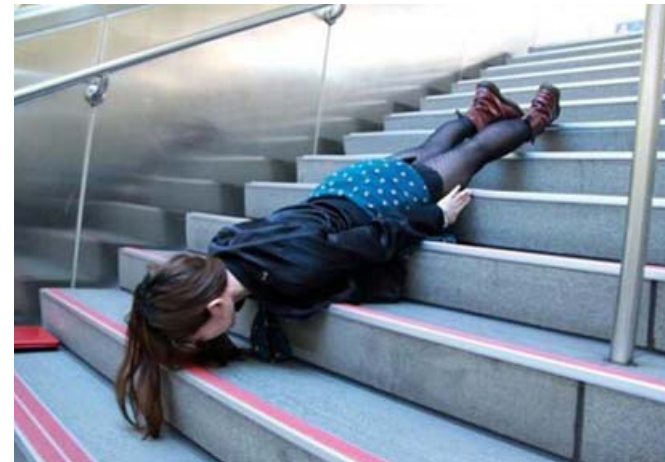
cat

woman

trash bin

2011-03-18 15:54:33

It seems important to recognize object use

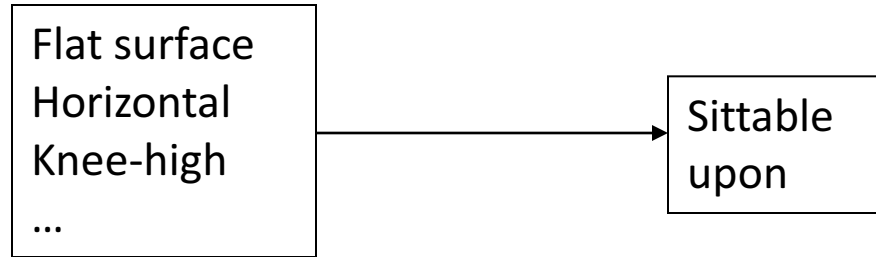


How to reason about typical / non-typical object use?



# The perception of function

- Direct perception (affordances): Gibson (70s-80s)



- Mediated perception (Categorization)

