



MIT CSAIL

**6.869: Advances in Computer Vision**

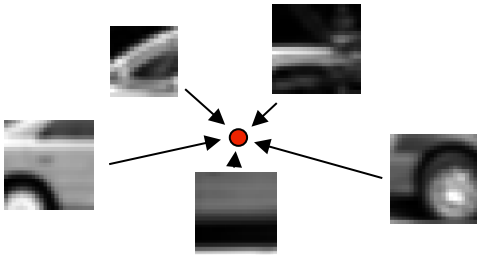
**MIT**  
COMPUTER  
VISION

# Lecture 21

## Object recognition IV

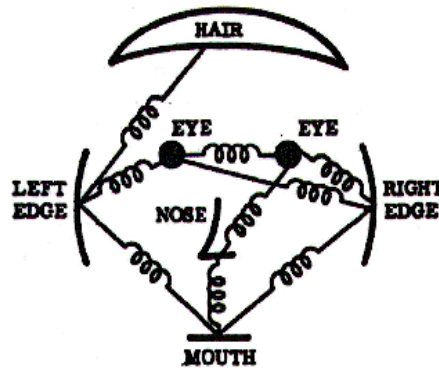
# Structure models

## Voting models



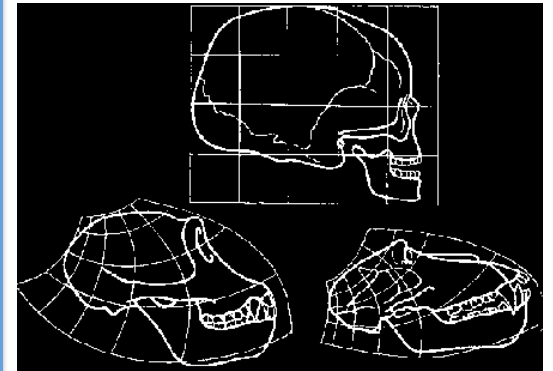
- Many parts ( $>100$ )

## Constellation models



- Few parts ( $\sim 6$ )

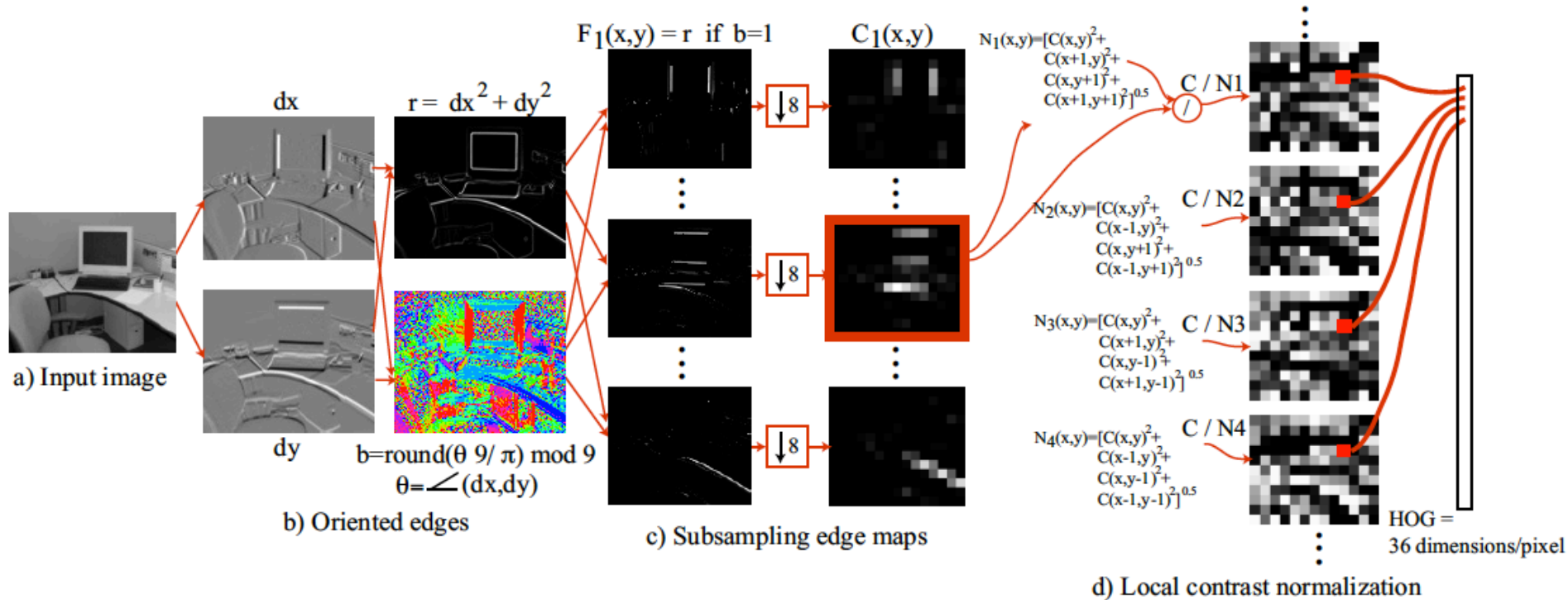
## Deformable models



- No parts



# HOG

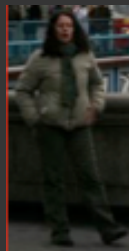


# Scanning-window templates

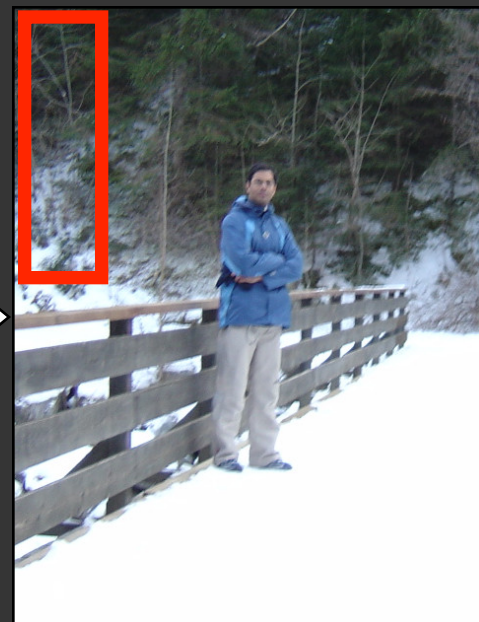
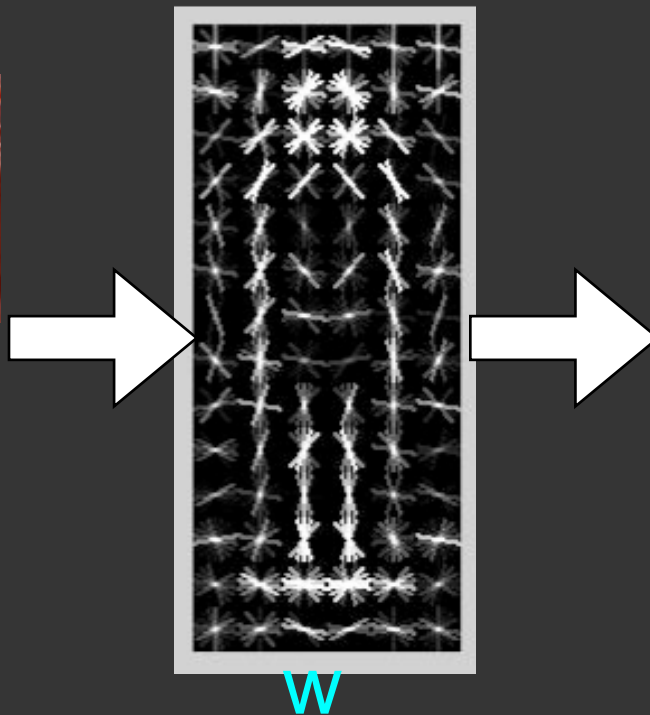
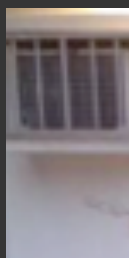
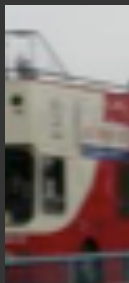
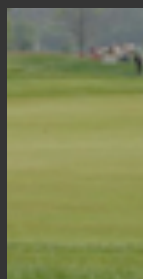
Dalal and Triggs CVPR05 (HOG)

Papageorgiou and Poggio ICIP99 (wavelets)

pos



neg



$w$  = weights for orientation and spatial bins

$$w \cdot x > 0$$



Train with a linear classifier (perceptron, logistic regression, SVMs...)

# Object Detection with Discriminatively Trained Part Based Models

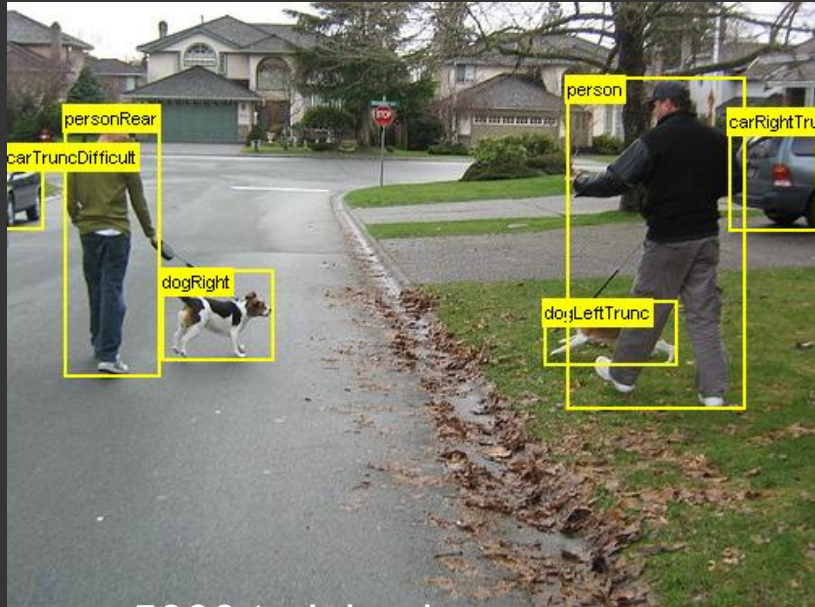
Pedro F. Felzenszwalb, Ross B. Girshick, David McAllester and Deva Ramanan

**Abstract**—We describe an object detection system based on mixtures of multiscale deformable part models. Our system is able to represent highly variable object classes and achieves state-of-the-art results in the PASCAL object detection challenges. While deformable part models have become quite popular, their value had not been demonstrated on difficult benchmarks such as the PASCAL datasets. Our system relies on new methods for discriminative training with partially labeled data. We combine a margin-sensitive approach for data-mining hard negative examples with a formalism we call *latent SVM*. A latent SVM is a reformulation of MI-SVM in terms of latent variables. A latent SVM is semi-convex and the training problem becomes convex once latent information is specified for the positive examples. This leads to an iterative training algorithm that alternates between fixing latent values for positive examples and optimizing the latent SVM objective function.

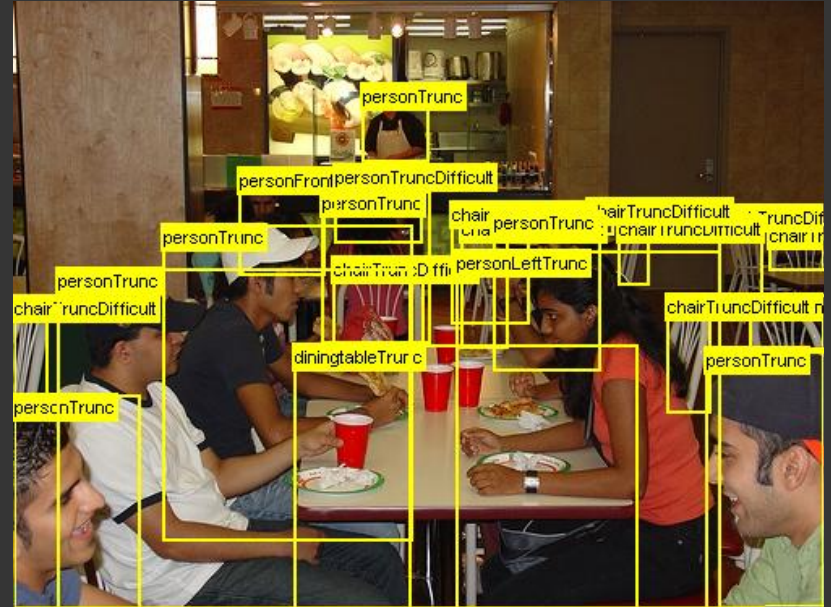
**Index Terms**—Object Recognition, Deformable Models, Pictorial Structures, Discriminative Training, Latent SVM



# PASCAL Visual Object Challenge



5000 training images



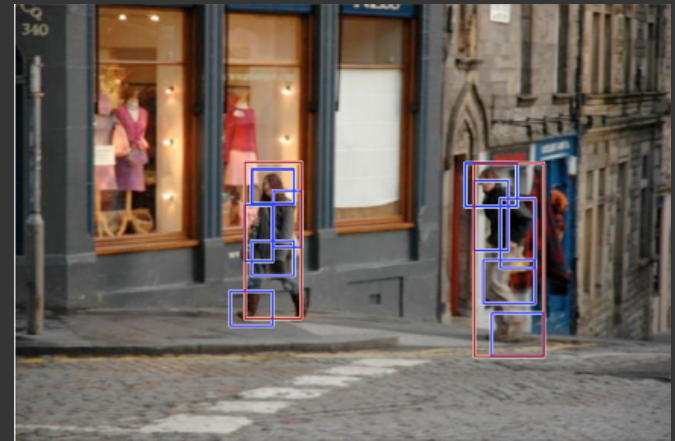
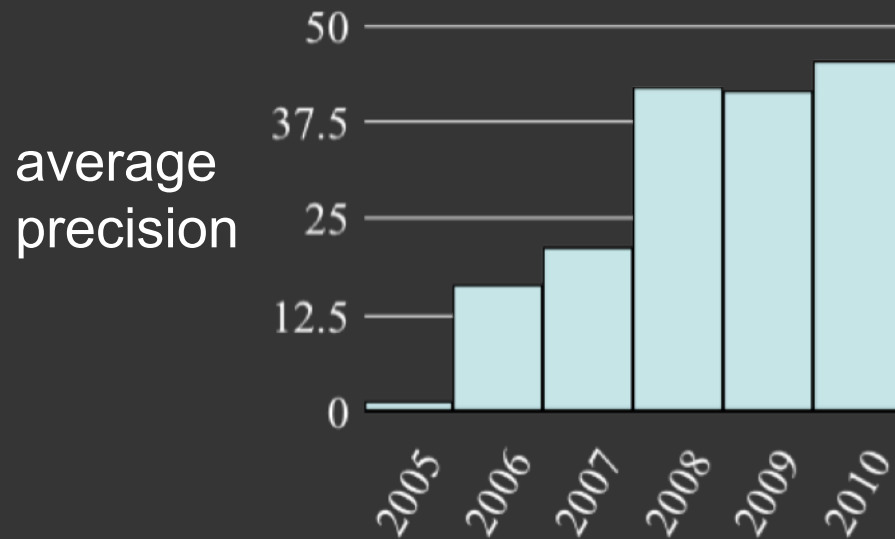
5000 testing images

20 everyday object categories

aeroplane bike bird boat bottle bus car cat chair cow table  
dog horse motorbike person plant sheep sofa train tv



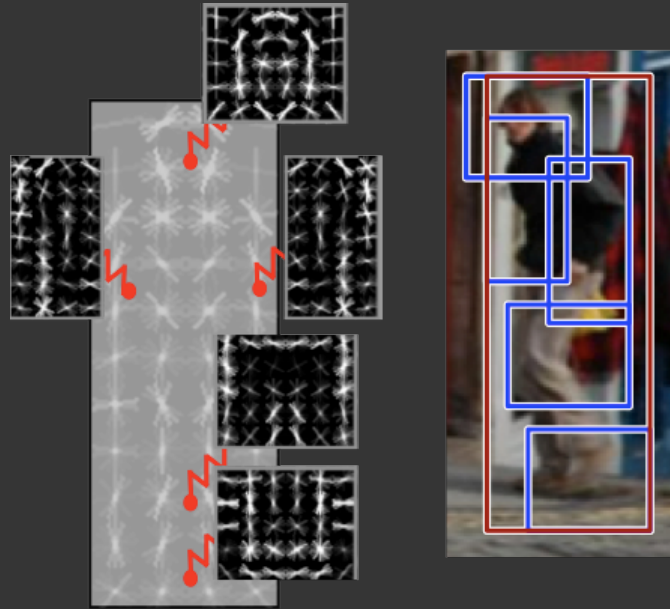
# 5 years of PASCAL people detection



1% to 45% in 5 years

Discriminative mixtures of star models 2007-2010 Felzenszwalb,  
McAllester, Ramanan CVPR 2008  
Felzenszwalb, Girshick, McAllester, and Ramanan PAMI 2009

# Deformable part models



Model encodes **local appearance** + **pairwise geometry**

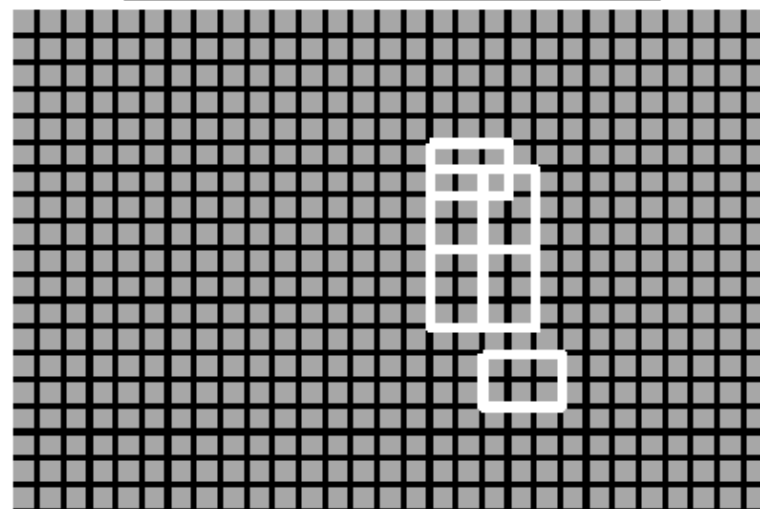
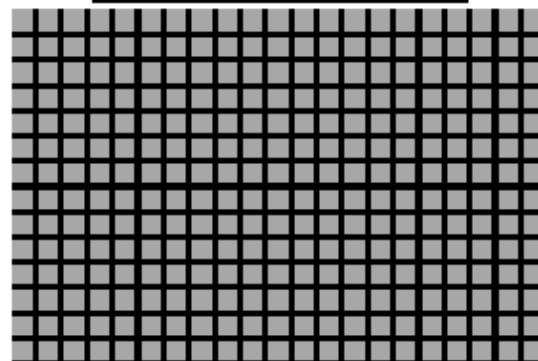
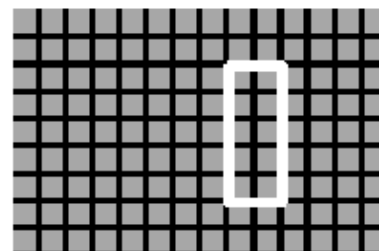
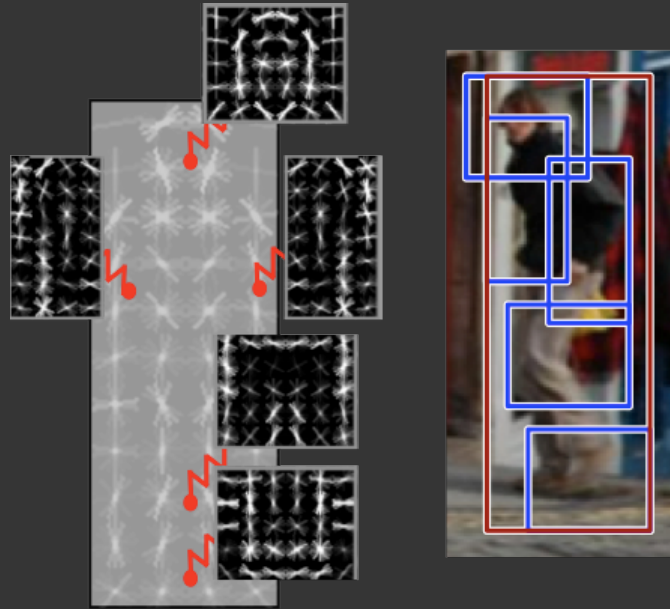


Image pyramid

Feature pyramid

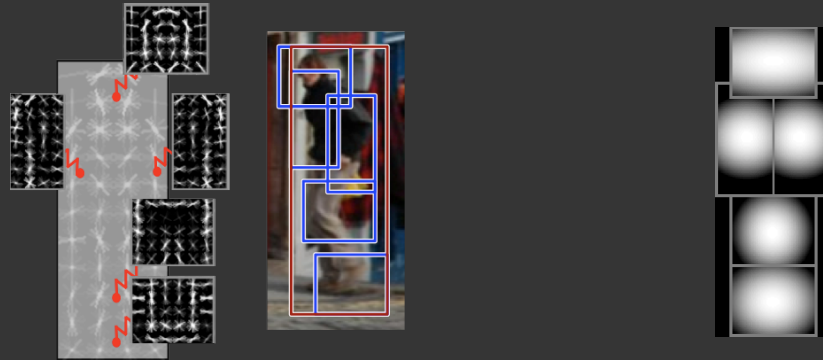
# Deformable part models



Model encodes **local appearance** + **pairwise geometry**



# Scoring function



$$\text{score}(x, z) = \sum_i w_i \phi(x, z_i) + \sum_{i,j} w_{ij} \Psi(z_i, z_j)$$

$x$  = image  
 $z_i = (x_i, y_i)$   
 $z = \{z_1, z_2, \dots\}$

part template  
scores

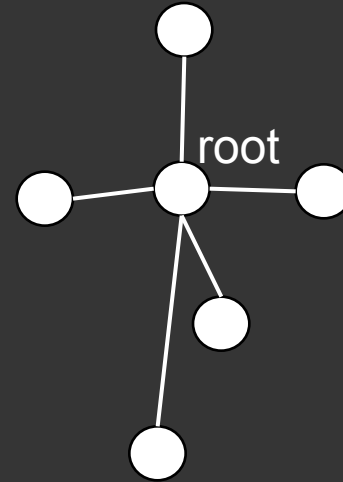
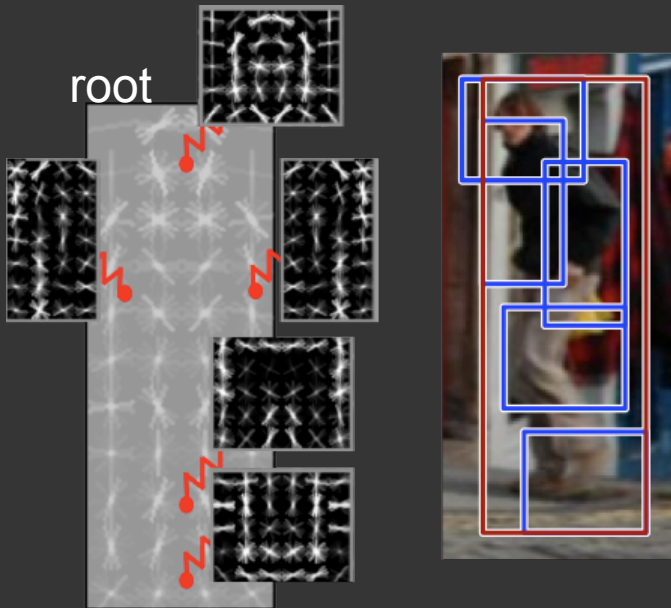
spring deformation model

Score is linear in local templates  $w_i$  and spring parameters  $w_{ij}$

$$\text{score}(x, z) = w \cdot \Phi(x, z)$$

# Inference: $\max_z \text{score}(x, z)$

Felzenszwalb & Huttenlocher 05



Star model: the location of the root filter is the anchor point  
Given the root location, all part locations are independent

# Classification



$$f_w(x) > 0$$

$$f_w(x) = w \cdot \Phi(x)$$



# Latent-variable classification



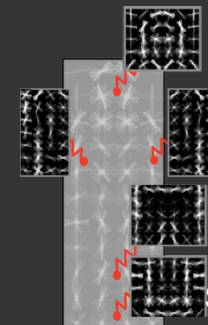
$$f_w(x) = w \cdot \Phi(x)$$

$$f_w(x) > 0$$



$$f_w(x) = \max_z S(x, z)$$

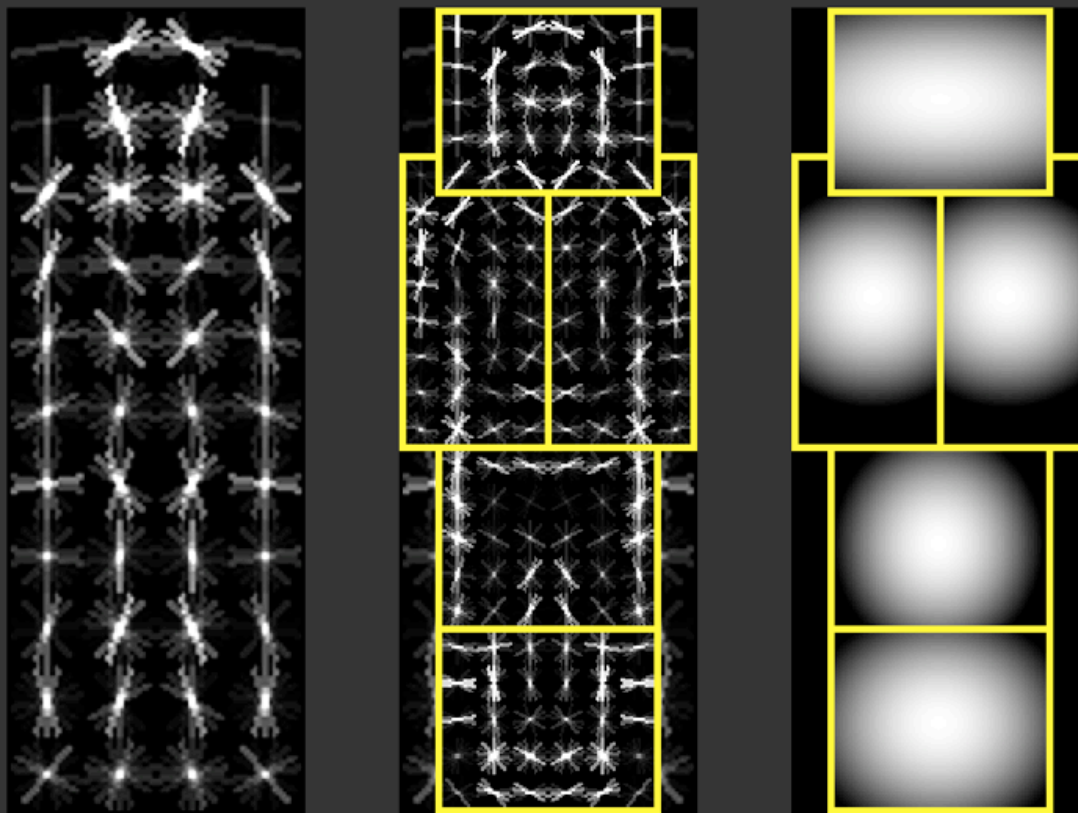
$$= \max_z w \cdot \Phi(x, z)$$



# Learning Initialization

Learn root filter with SVM

Initialize part filters to regions in  
root filter with lots of energy



# Coordinate descent

1) Given positive part locations, learn  $w$  with a convex program

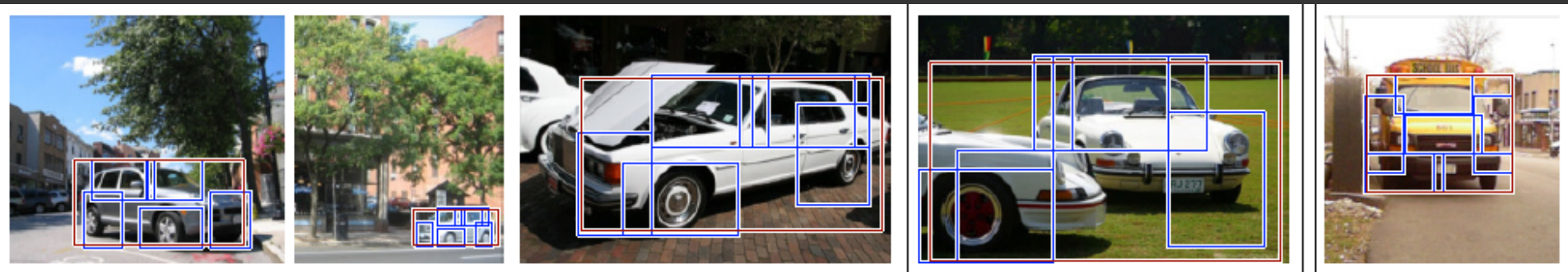
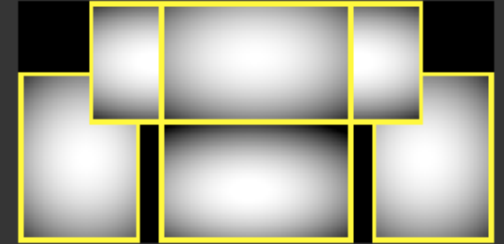
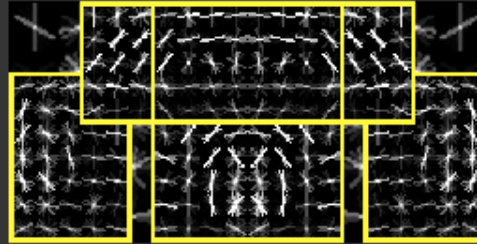
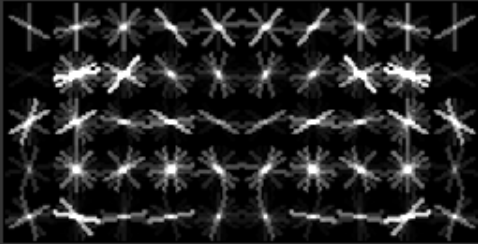
$$w = \underset{w}{\operatorname{argmin}} L(w) \quad \text{with fixed} \quad \{z_n : n \in \text{pos}\}$$

2) Given  $w$ , estimate part locations on positives

$$z_n = \underset{z}{\operatorname{argmax}} w \cdot \Phi(x_n, z) \quad \forall n \in \text{pos}$$

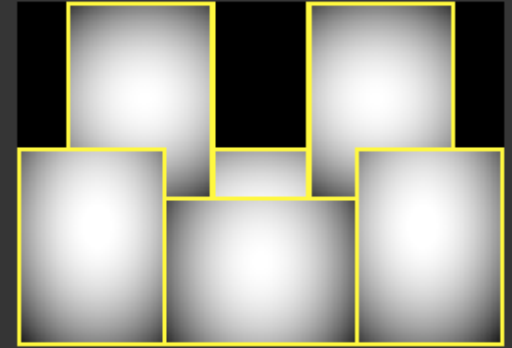
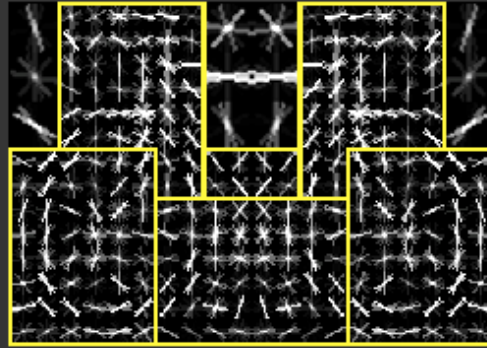
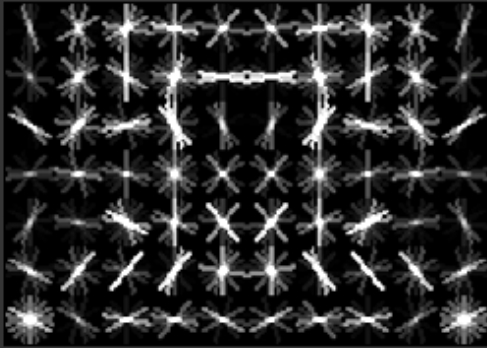
The above steps perform coordinate descent on a joint loss

# Example models



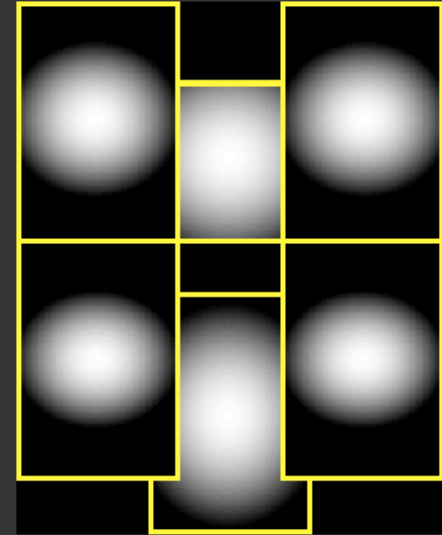
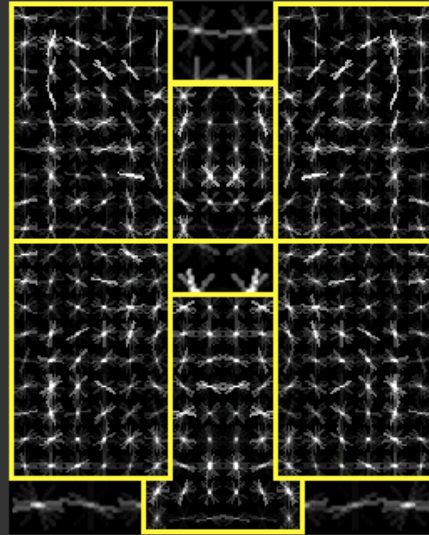


# Example models

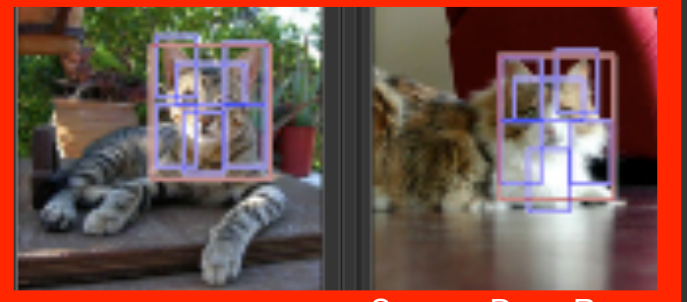
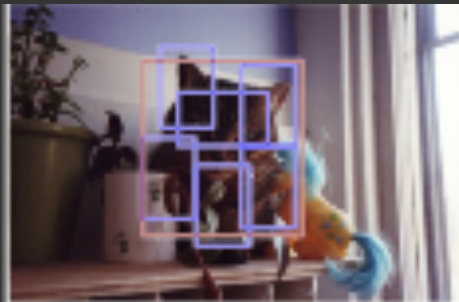
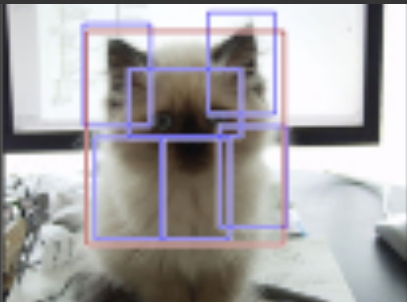




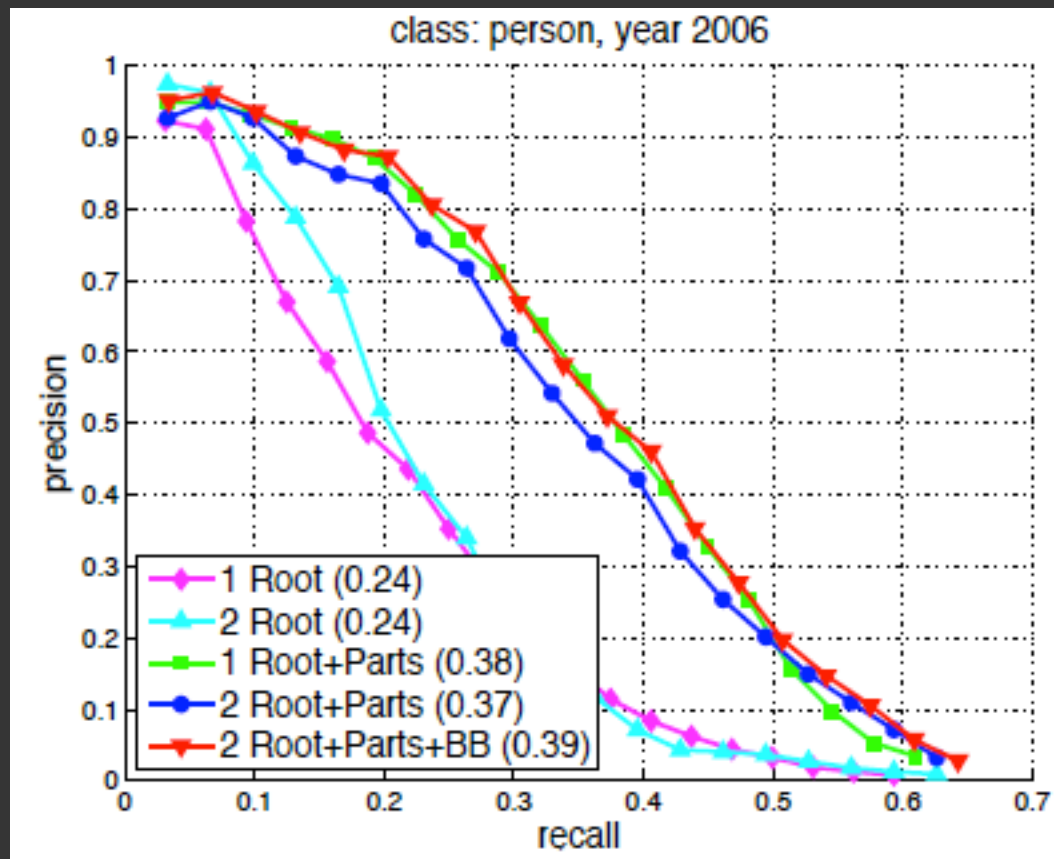
# Example models



False positive due to imprecise bounding box



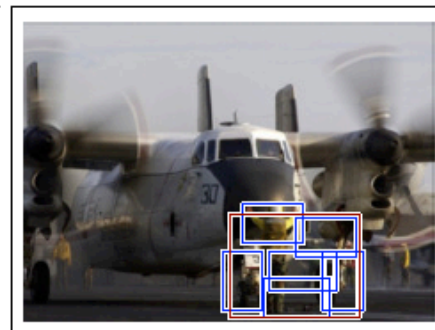
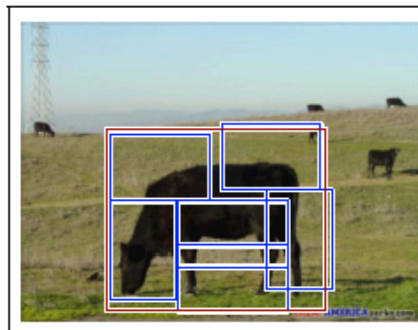
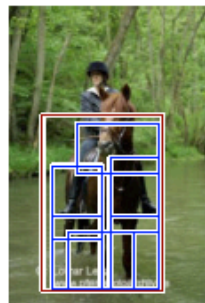
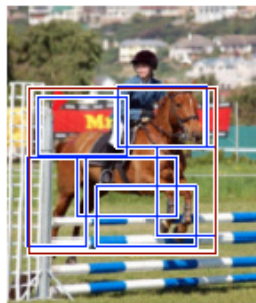
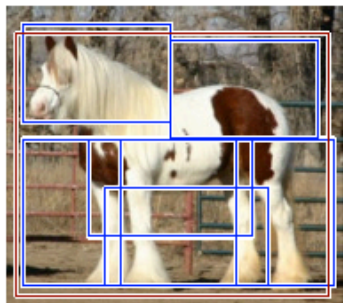
Source: Deva Ramanan



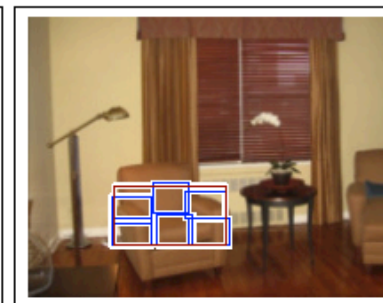
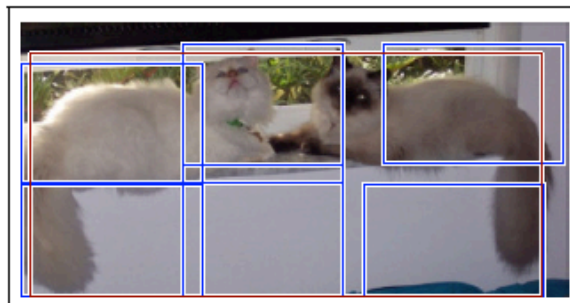
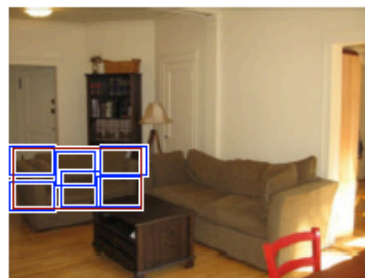
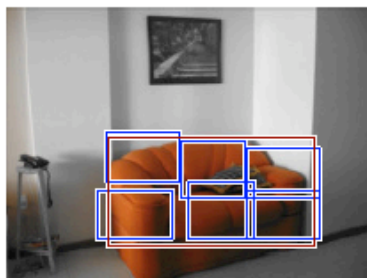
Other tricks:

- Mining hard negative examples
- Noisy annotations

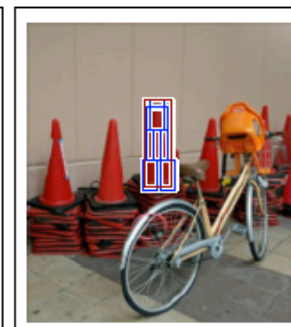
horse



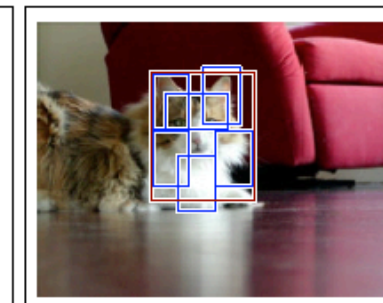
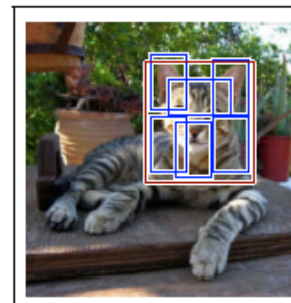
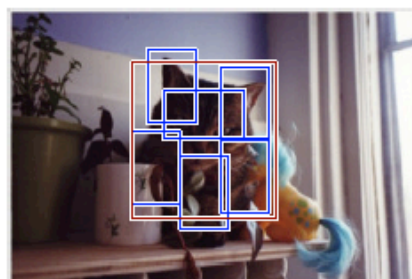
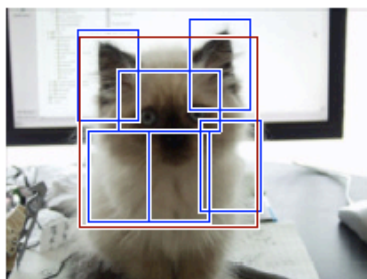
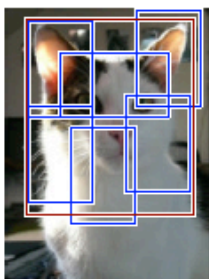
sofa



bottle

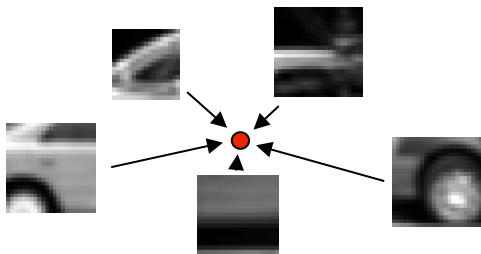


cat



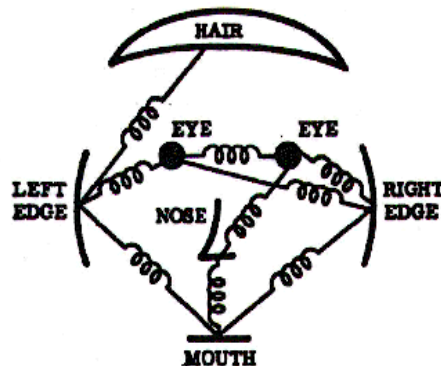
# Structure models

## Voting models



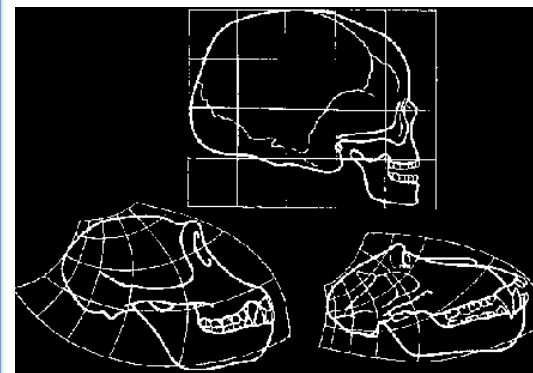
- Many parts ( $>100$ )

## Constellation models



- Few parts ( $\sim 6$ )

## Deformable models



- No parts



# ON GROWTH AND FORM

The Complete Revised Edition



D'Arcy Wentworth Thompson

to the lines of our new curved ordinates. In like manner, the still more bizarre outlines of other fishes of the same family of Chaetodonts will be found to correspond to very slight modifications of

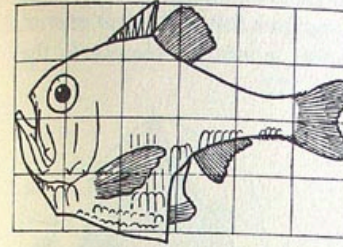


Fig. 146. *Argyropelecus olfersi*.

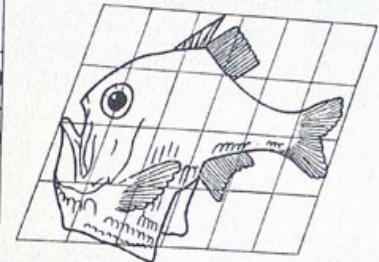


Fig. 147. *Sternoptyx diaphana*.

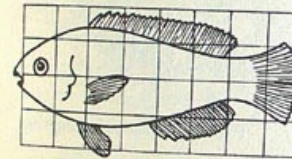


Fig. 148. *Scarus* sp.

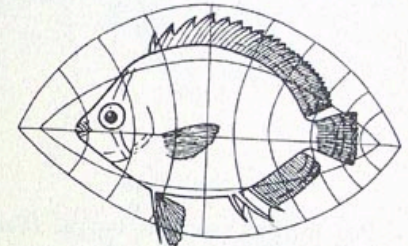


Fig. 149. *Pomacanthus*.

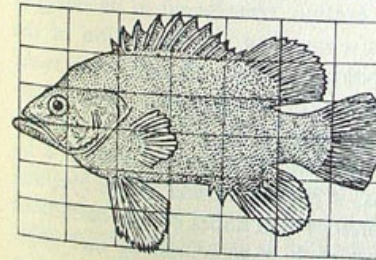


Fig. 150. *Polyprion*.

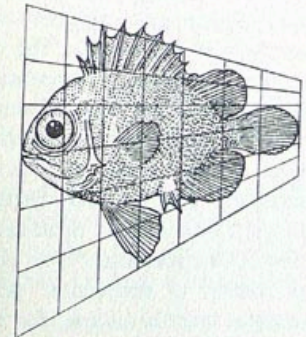


Fig. 151. *Pseudopriacanthus altus*.

similar co-ordinates; in other words, to small variations in the values of the constants of the coaxial curves.

In Figs. 150-153 I have represented another series of Acanthopterygian fishes, not very distantly related to the foregoing. If we

From wikipedia: Perhaps the most famous part of the work is chapter XVII, "The Comparison of Related Forms," where Thompson explored the degree to which differences in the forms of related animals could be described by means of relatively simple mathematical transformations.

# Shape Matching and Object Recognition Using Shape Contexts

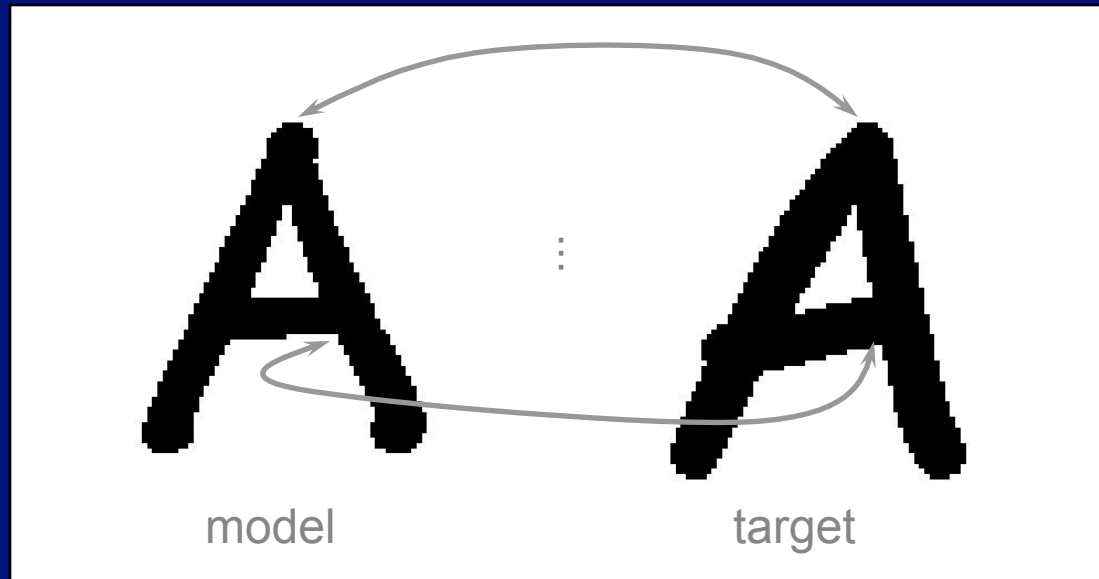
Serge Belongie, *Member, IEEE*, Jitendra Malik, *Member, IEEE*, and Jan Puzicha

**Abstract**—We present a novel approach to measuring similarity between shapes and exploit it for object recognition. In our framework, the measurement of similarity is preceded by 1) solving for correspondences between points on the two shapes, 2) using the correspondences to estimate an aligning transform. In order to solve the correspondence problem, we attach a descriptor, the *shape context*, to each point. The shape context at a reference point captures the distribution of the remaining points relative to it, thus offering a globally discriminative characterization. Corresponding points on two similar shapes will have similar shape contexts, enabling us to solve for correspondences as an optimal assignment problem. Given the point correspondences, we estimate the transformation that best aligns the two shapes; regularized thin-plate splines provide a flexible class of transformation maps for this purpose. The dissimilarity between the two shapes is computed as a sum of matching errors between corresponding points, together with a term measuring the magnitude of the aligning transform. We treat recognition in a nearest-neighbor classification framework as the problem of finding the stored prototype shape that is maximally similar to that in the image. Results are presented for silhouettes, trademarks, handwritten digits, and the COIL data set.

**Index Terms**—Shape, object recognition, digit recognition, correspondence problem, MPEG7, image registration, deformable templates.

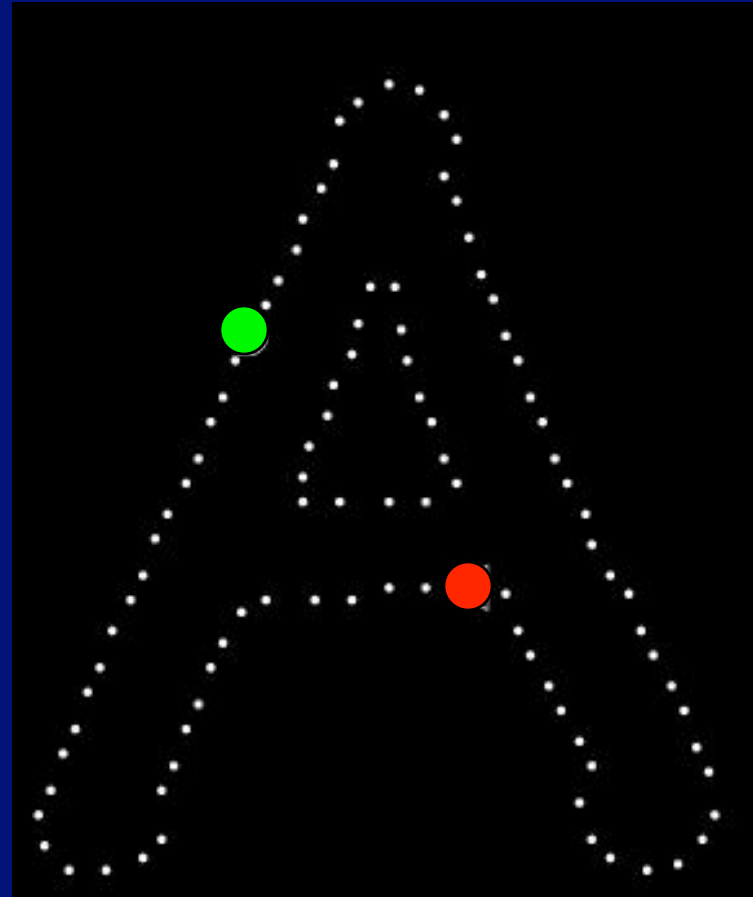
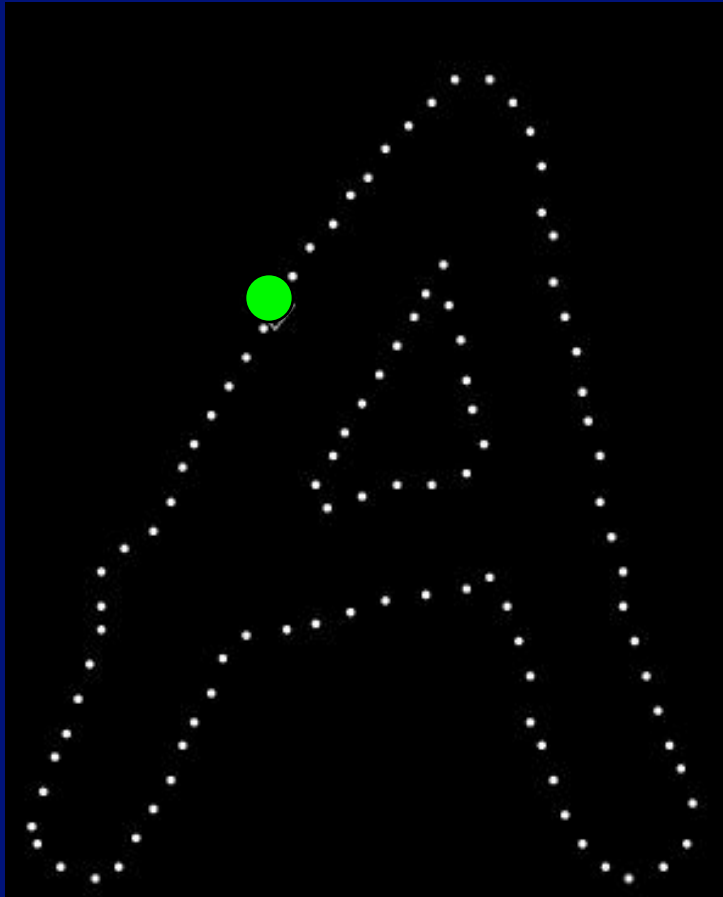


# Matching Framework



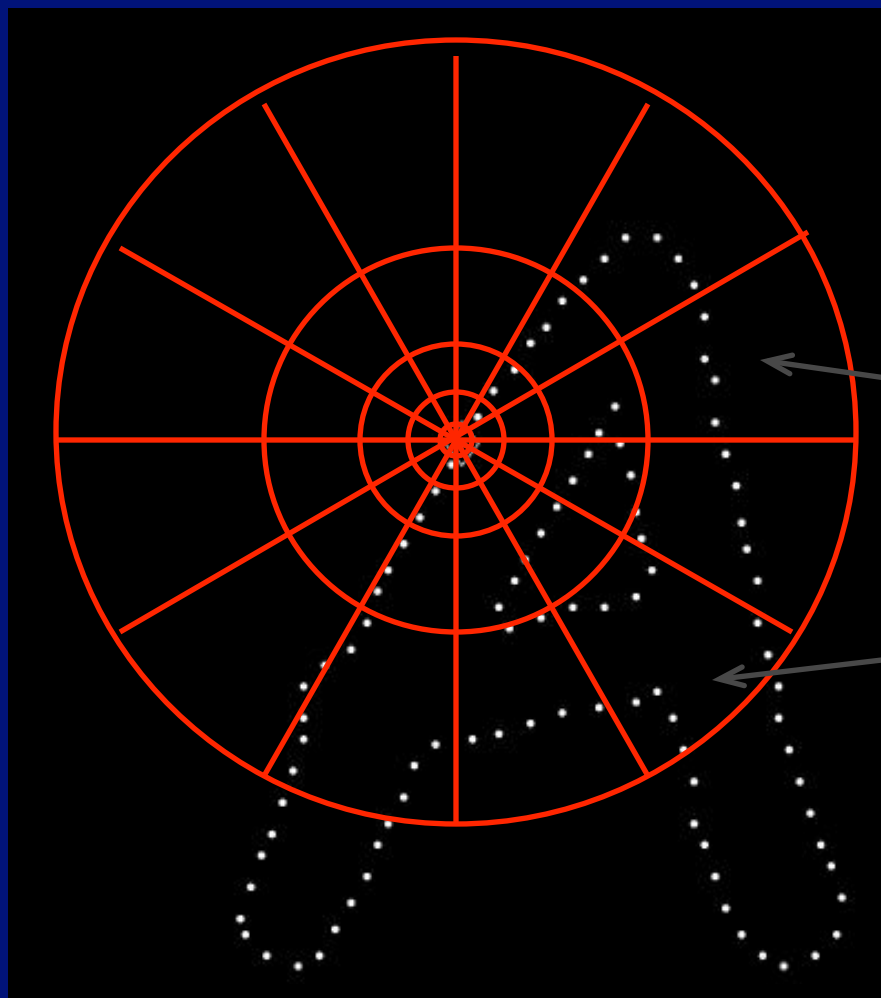
- Find correspondences between points on shape
- Fast pruning
- Estimate transformation & measure similarity

# Comparing Pointsets





# Shape Context



Count the number of points inside each bin, e.g.:

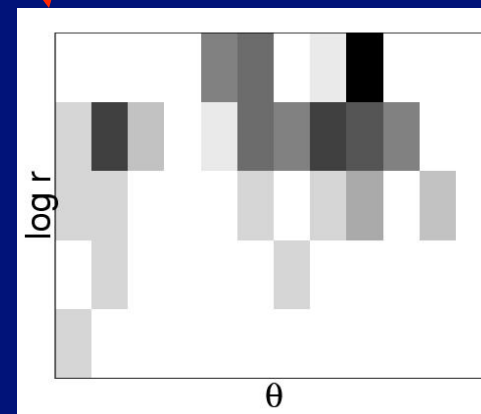
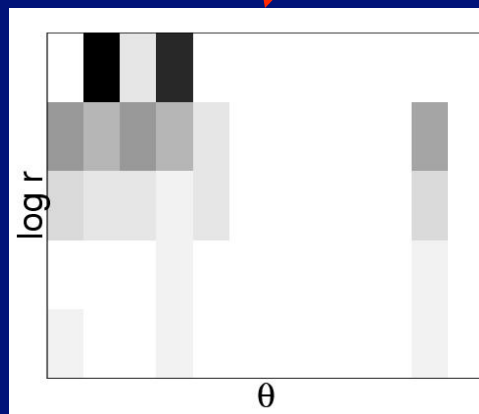
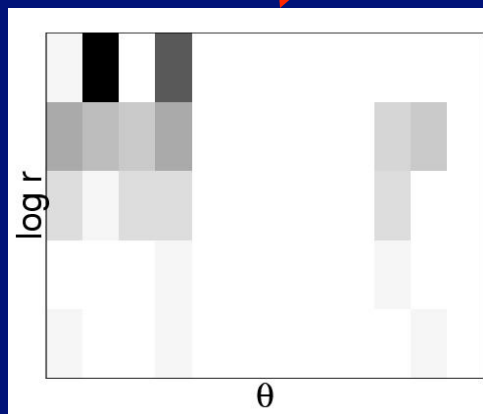
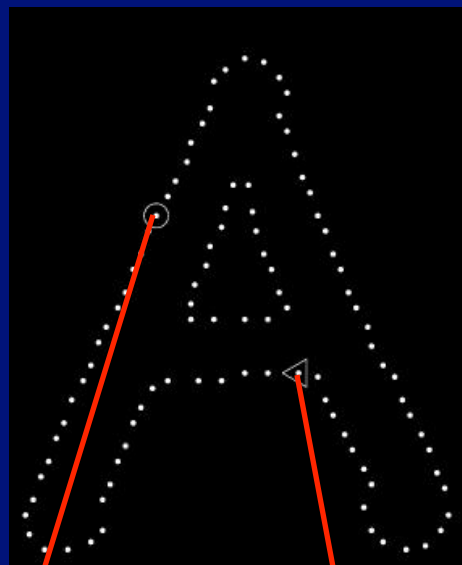
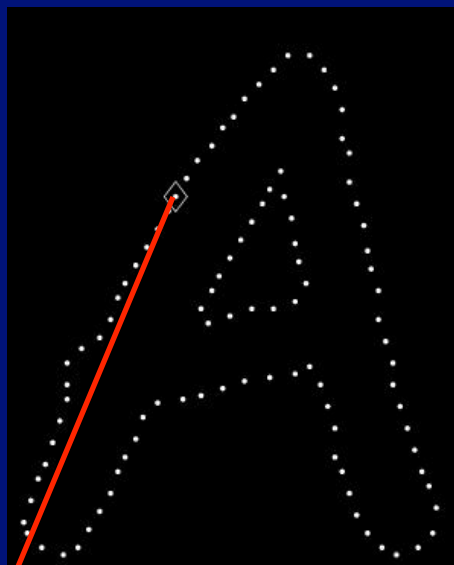
Count = 4

⋮

Count = 10

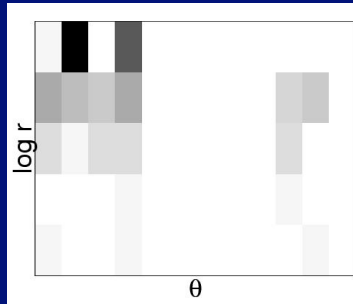
- ✦ Compact representation of distribution of points relative to each point

# Shape Context

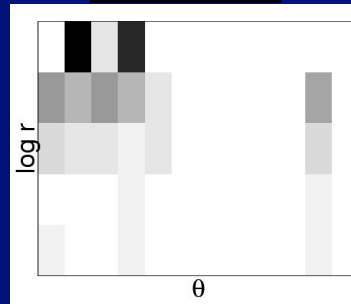


# Comparing Shape Contexts

$h_i(k)$

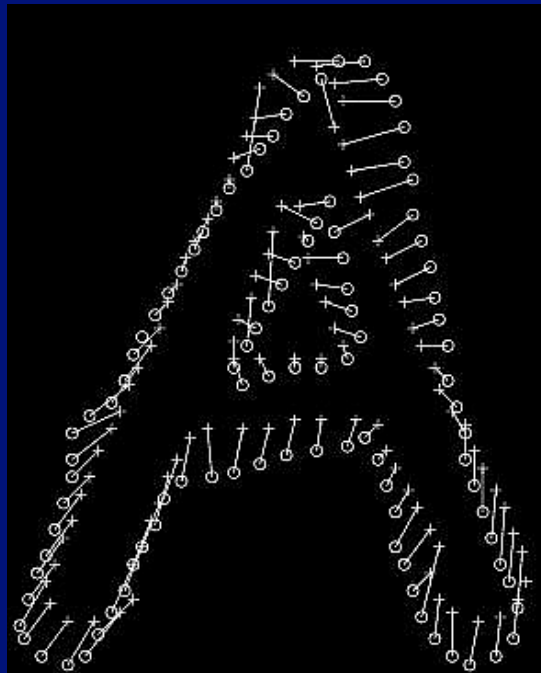


$h_j(k)$



Compute matching costs using Chi Squared distance:

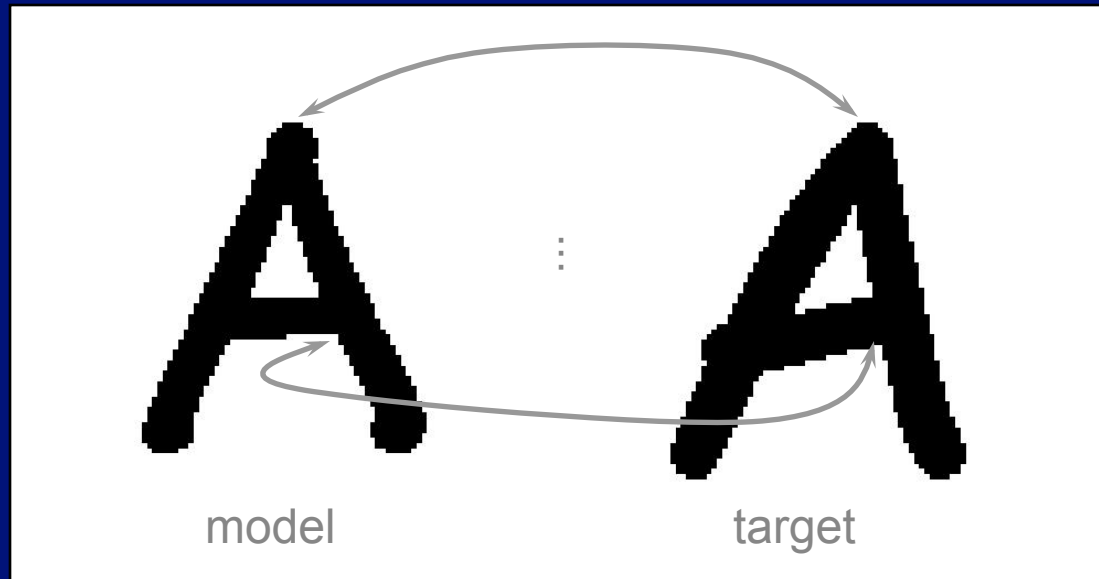
$$C_{ij} = \frac{1}{2} \sum_{k=1}^K \frac{[h_i(k) - h_j(k)]^2}{h_i(k) + h_j(k)}$$



Recover correspondences by solving linear assignment problem with costs  $C_{ij}$

[Jonker & Volgenant 1987]

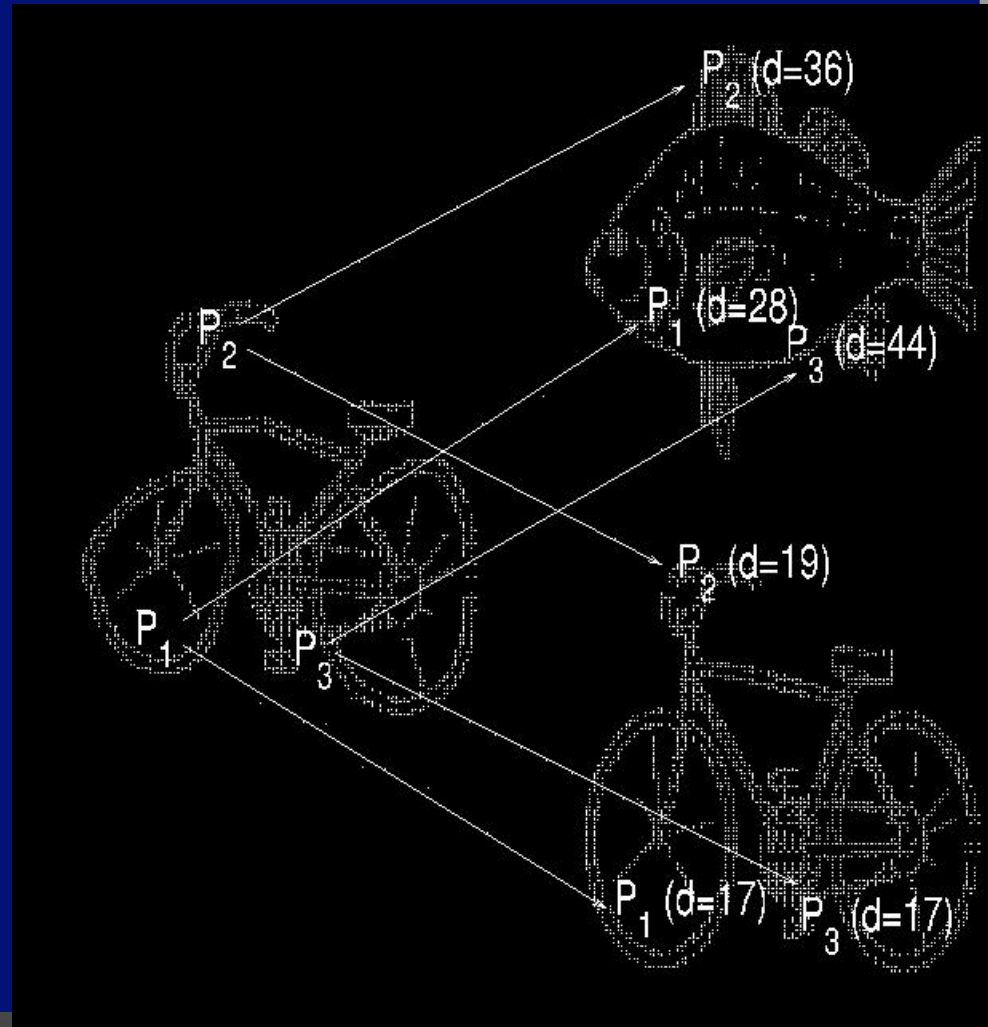
# Matching Framework



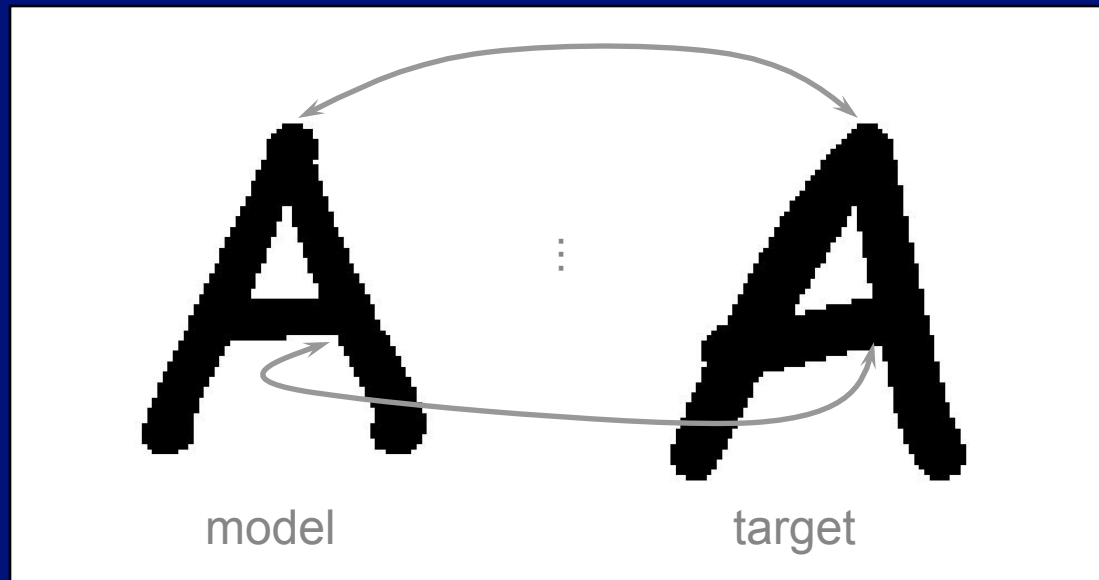
- Find correspondences between points on shape
- **Fast pruning**
- Estimate transformation & measure similarity

# Fast pruning

- Find best match for the shape context at only a few random points and add up cost

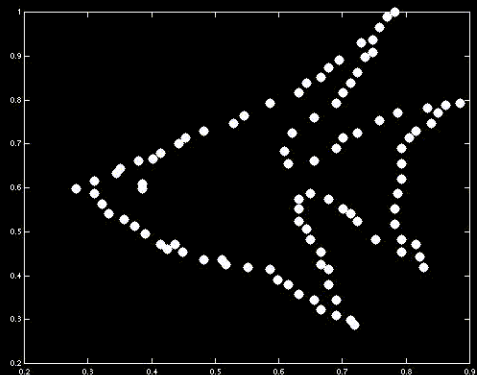


# Matching Framework

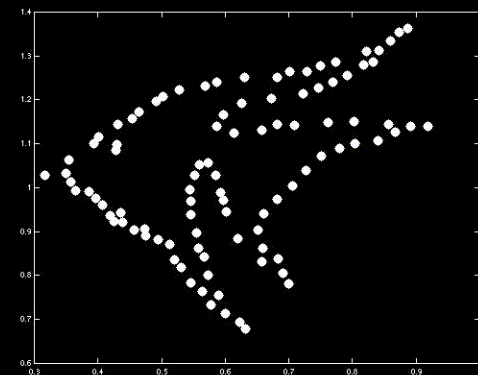


- Find correspondences between points on shape
- Fast pruning
- Estimate transformation & measure similarity

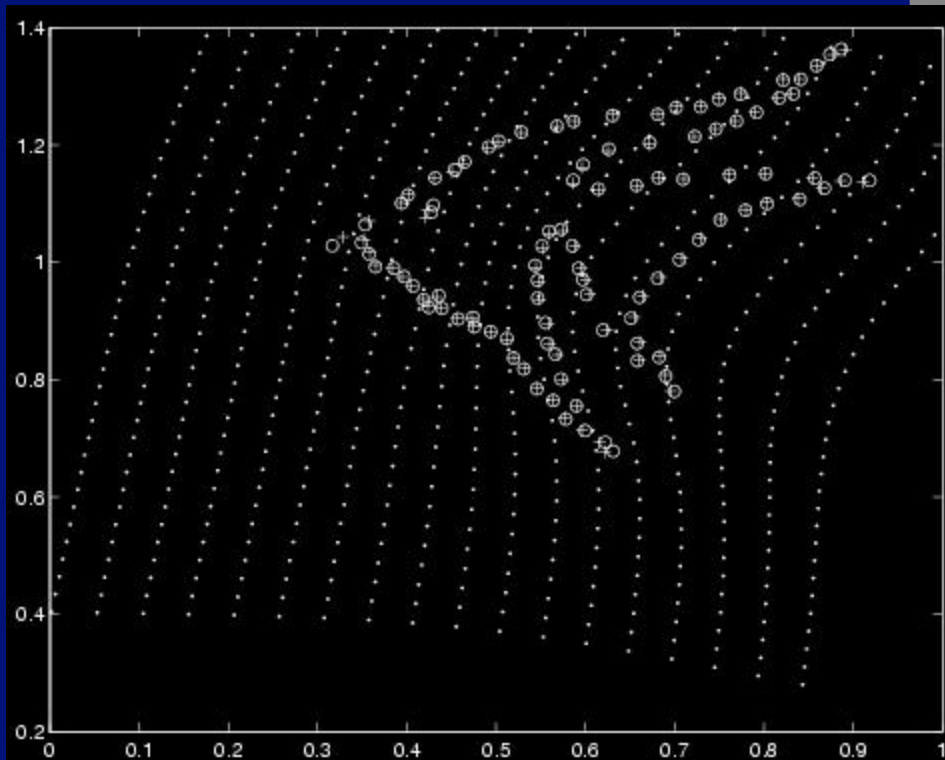
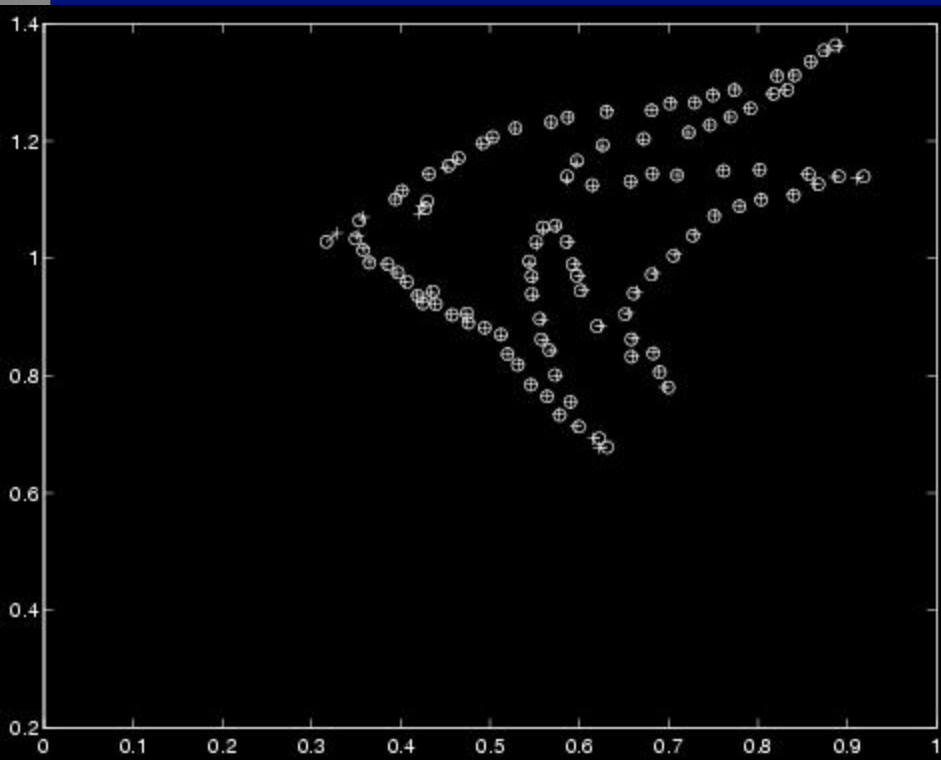
# Matching Example



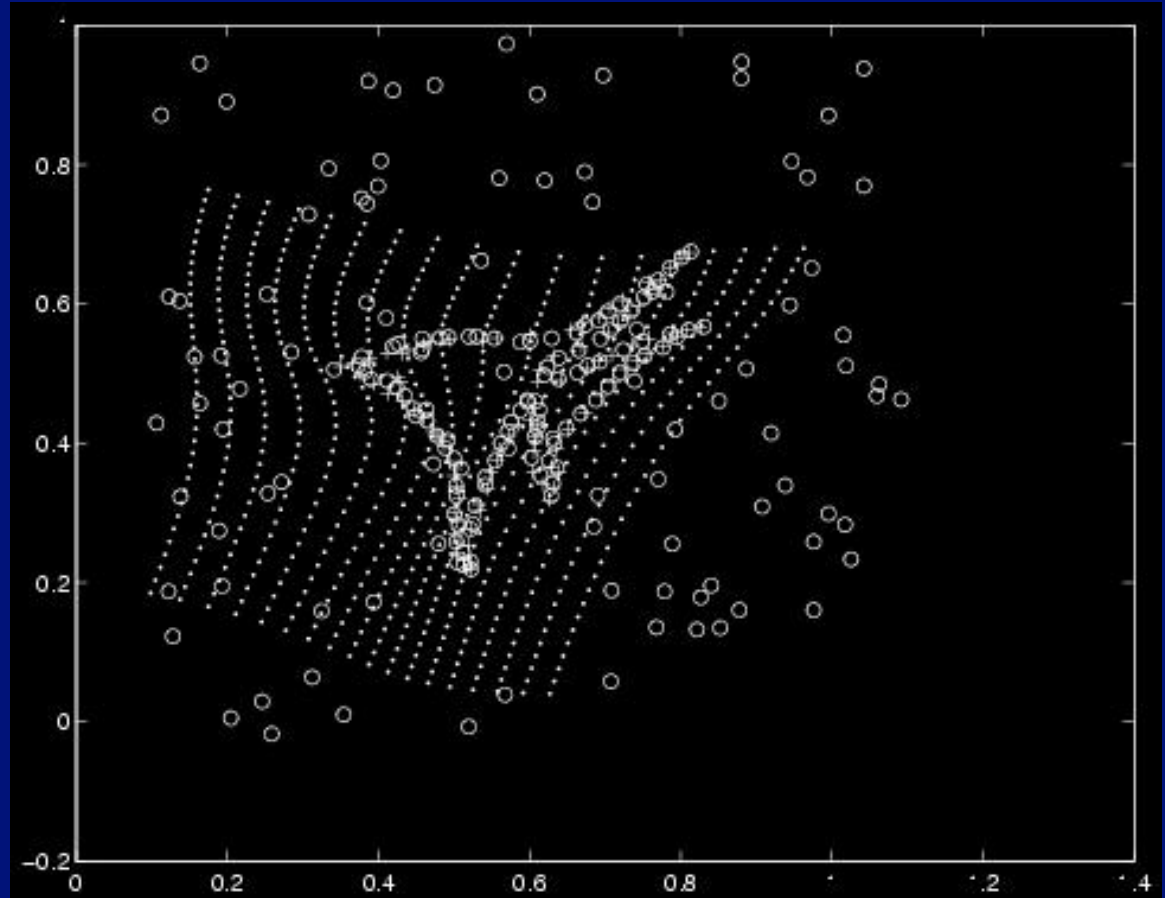
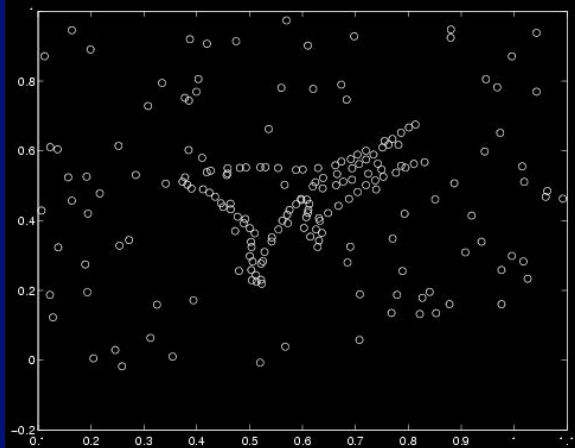
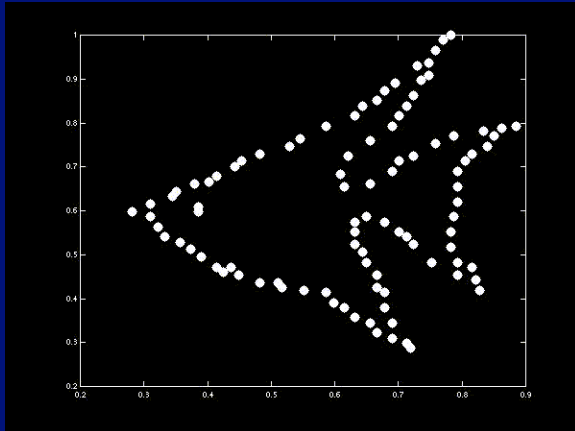
model



target



# Outlier Test Example





# The spaces of faces is not convex



The average of two faces is ...

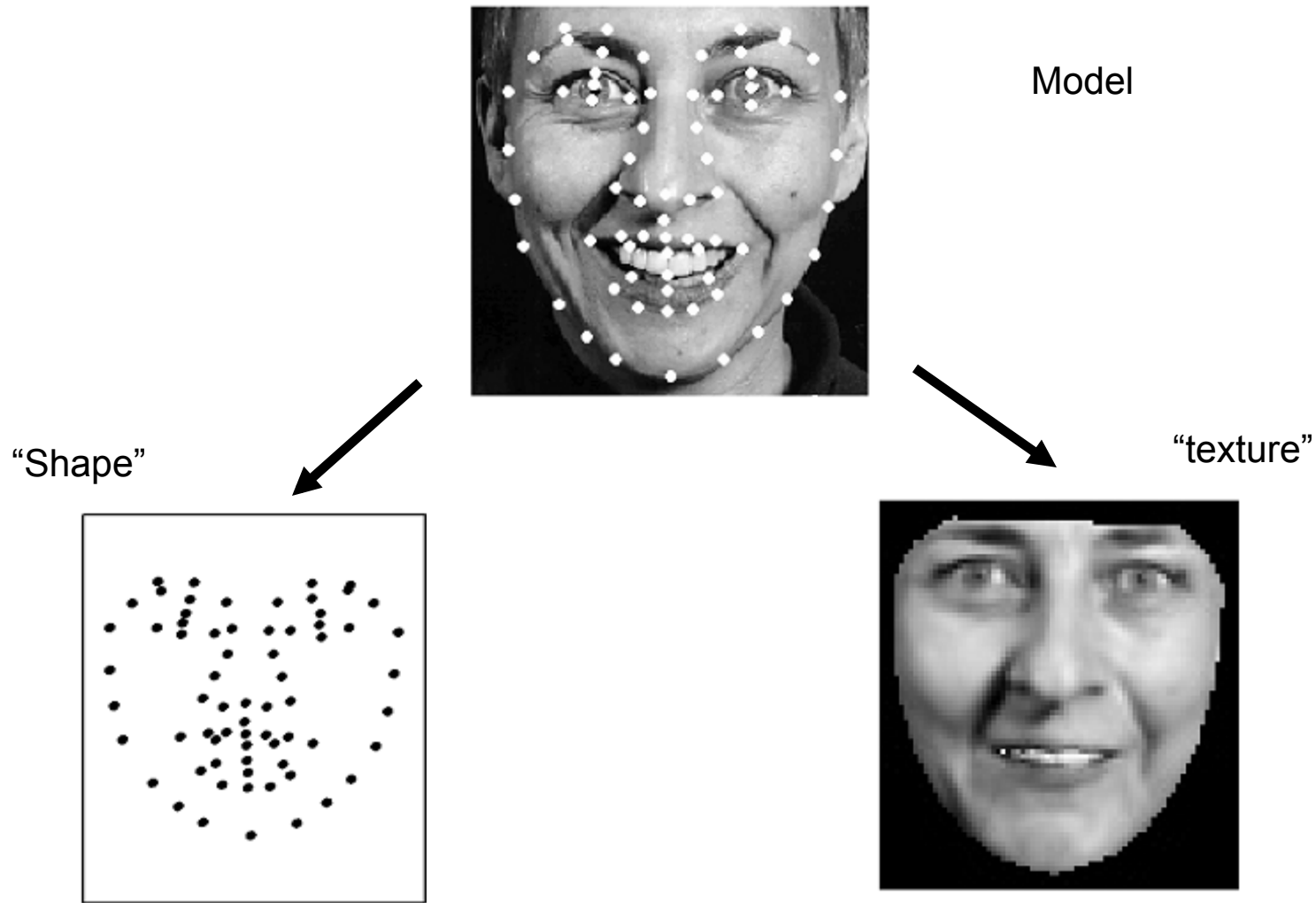
# The spaces of faces is not convex



The average of two faces is not another face

# A shape-texture face model

---

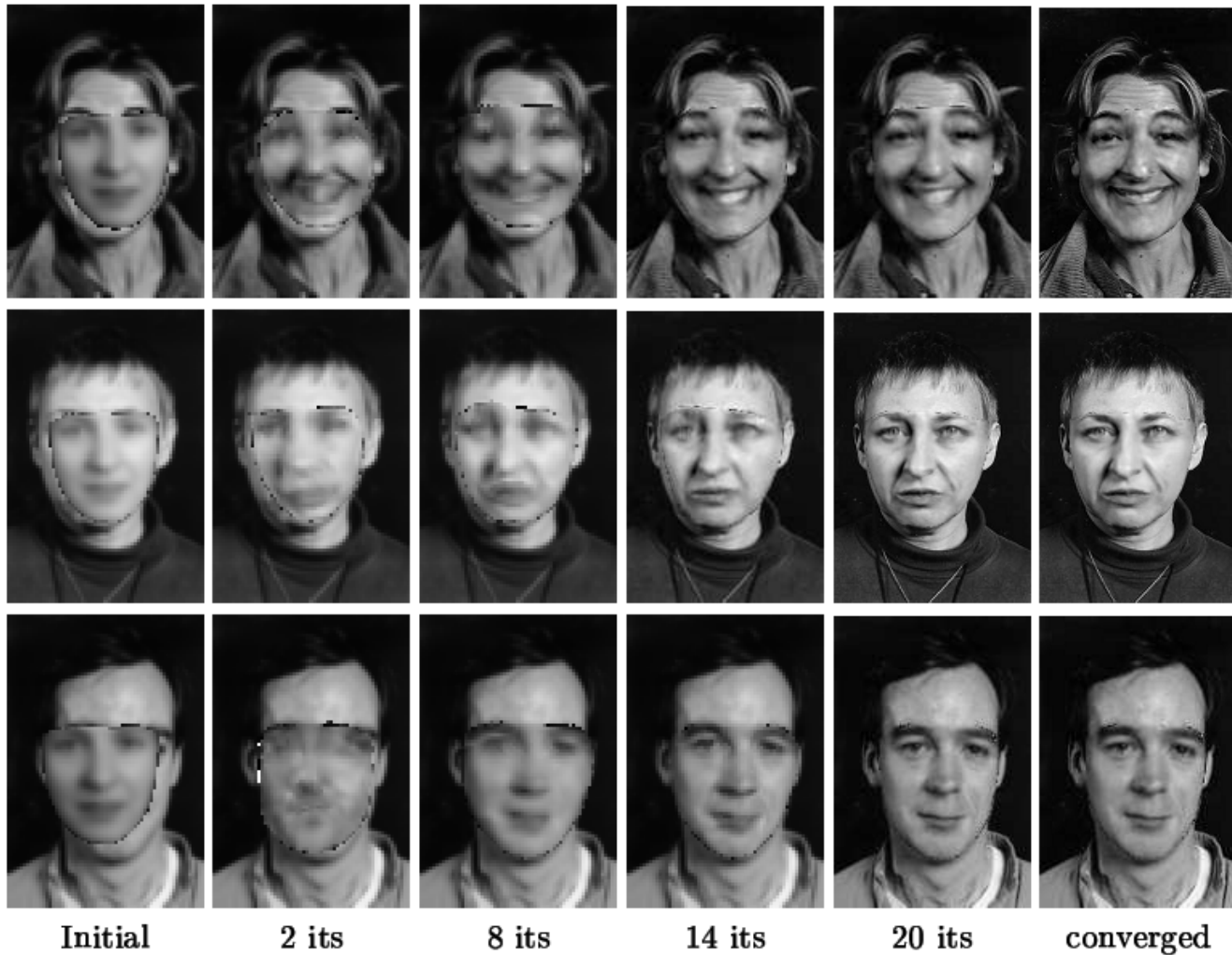


Cootes, Edwards, and Taylor, “Active Appearance Models”, ECCV 1998

Blanz, V. and Vetter, T., A morphable model for the synthesis of 3D faces, 1999

# Active Appearance Model Search (Results)

---



# Essence of the Idea: Recognition by Synthesis

Explain a new example in terms of the model parameters



Initial



3 its



8 its



11 its



Converged



Original

# Enhancing gender

---



more same **original** androgynous more opposite



# Changing age

---

Face becomes  
“rounder” and “more  
textured” and “grayer”

original



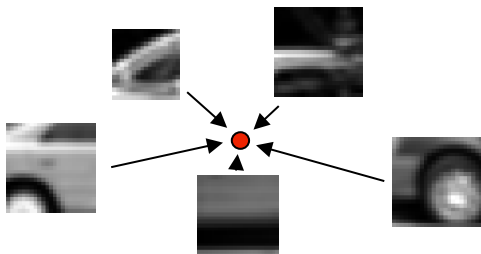
shape

color

both

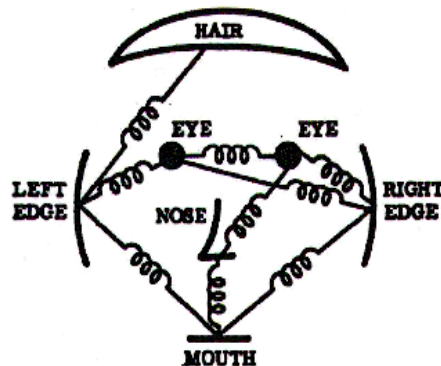
# Structure models

## Voting models



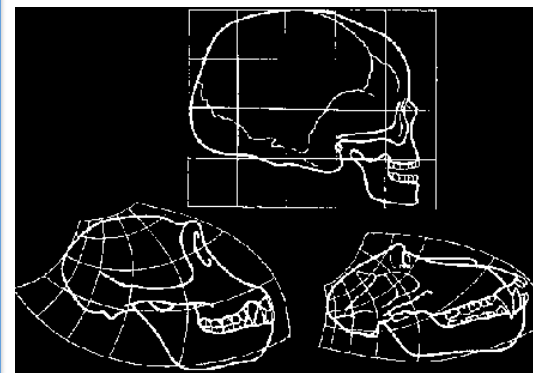
- Many parts ( $>100$ )

## Constellation models



- Few parts ( $\sim 6$ )

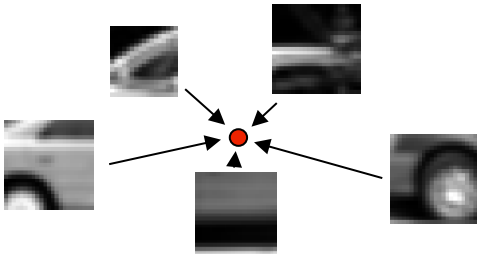
## Deformable models



- No parts

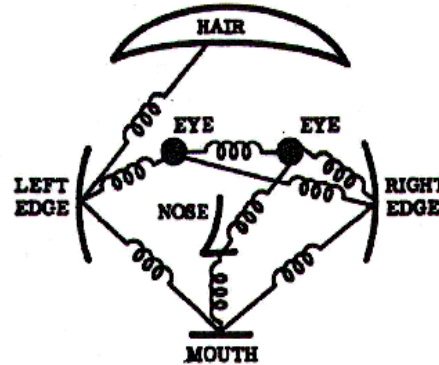
# Structure models

## Voting models



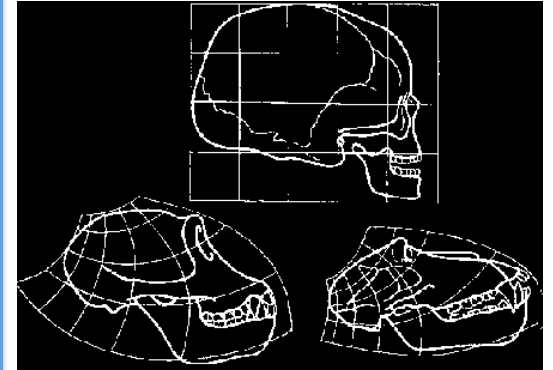
- Many parts ( $>100$ )

## Constellation models



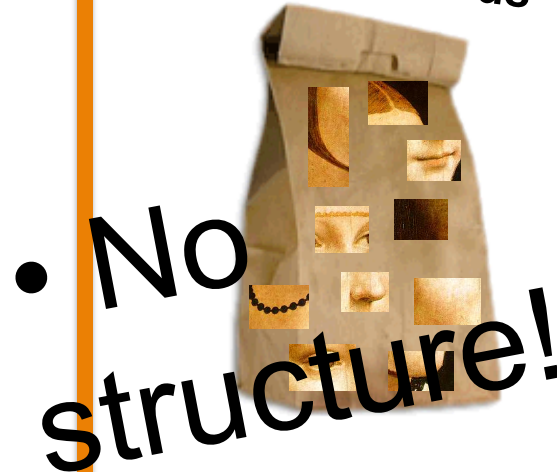
- Few parts ( $\sim 6$ )

## Deformable models



- No parts

## Bag of words



Object



Bag of 'words'



# Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes.

For a long time, the retinal image was considered as a visual centers in the brain.

Hubel and Wiesel discovered that the eye, cell, optical nerve, image perception, more complex, following the various cells of the cortex, Hubel and Wiesel demonstrate that the *message about the image falling on the retina undergoes a point-by-point analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.*

sensory, brain,  
visual, perception,  
retinal, cerebral cortex,  
eye, cell, optical  
nerve, image  
Hubel, Wiesel

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$580bn in 2004.

The increase in exports will also annoy the US, which has long complained about China's trade surplus. China's government has agreed to a deal with the US that the yuan is to be allowed to rise against the dollar. The government also needs to increase demand so that the yuan can be used in the country. China has also been the target of the US yuan against the dollar. The US has permitted it to trade within a narrow band but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

China, trade,  
surplus, commerce,  
exports, imports, US,  
yuan, bank, domestic,  
foreign, increase,  
trade, value

# Related works

- Early “bag of words” models: mostly texture recognition
  - Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003;
- Hierarchical Bayesian models for documents (pLSA, LDA, etc.)
  - Hoffman 1999; Blei, Ng & Jordan, 2004; Teh, Jordan, Beal & Blei, 2004
- Object categorization
  - Csurka, Bray, Dance & Fan, 2004; Sivic, Russell, Efros, Freeman & Zisserman, 2005; Sudderth, Torralba, Freeman & Willsky, 2005;
- Natural scene categorization
  - Vogel & Schiele, 2004; Fei-Fei & Perona, 2005; Bosch, Zisserman & Munoz, 2006



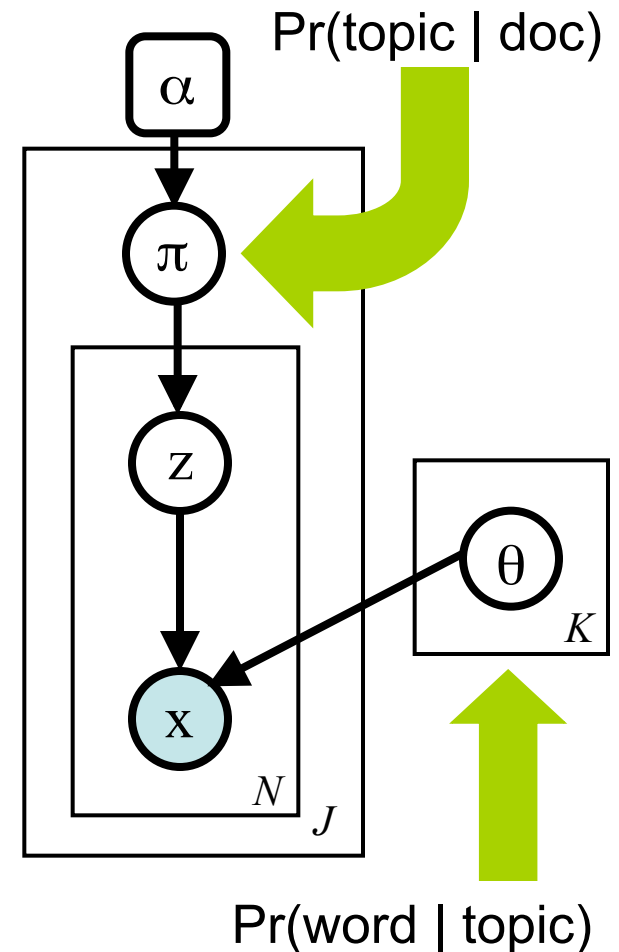
# Discovering topics in text collections

## Text document

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

## Discovered topics

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI



Latent Dirichlet Allocation (LDA)  
Blei, Ng, & Jordan, JMLR 2003

# Visual analogy

document - image

word - visual word

topics - objects

# Demo


A demonstration of bag-of-words classifiers - Microsoft Internet Explorer provided by Insight Broadband

File Edit View Favorites Tools Help

Back Forward Stop Reload Home Search Favorites RSS Print Mail News Groups

Address <http://people.csail.mit.edu/fergus/iccv2005/bagwords.html>

Google Search 100 blocked Check AutoLink AutoF



## Two bag-of-words classifiers

ICCV 2005 short courses on  
[Recognizing and Learning Object Categories](#)

A simple approach to classifying images is to treat them as a collection of regions, describing only their appearance and ignoring their location. This approach has been successfully used in the text community for analyzing documents and are known as "bag-of-words" models, since each document is represented by a distribution over fixed vocabulary(s). Using such a representation, methods such as probabilistic latent semantic analysis (pLSA) [1] (LDA) [2] are able to extract coherent topics within document collections in an unsupervised manner.

Recently, Fei-Fei et al. [3] and Sivic et al. [4] have applied such methods to the visual domain. The demo code implements pLSA, including a Naïve Bayes classifier. For comparison, a Naïve Bayes classifier is also provided which requires labelled training data, unlike pLSA.

The code consists of Matlab scripts (which should run under both Windows and Linux) and a couple of 32-bit Linux binaries for doing image representation. Hence the whole system will need to be run on Linux. The code is for teaching/research purposes only. If you find a bug, please email [fergus@csail.mit.edu](mailto:fergus@csail.mit.edu) where csail point mit point edu.

---

## Download

[Download](#) the code and datasets (32 Mbytes)

---

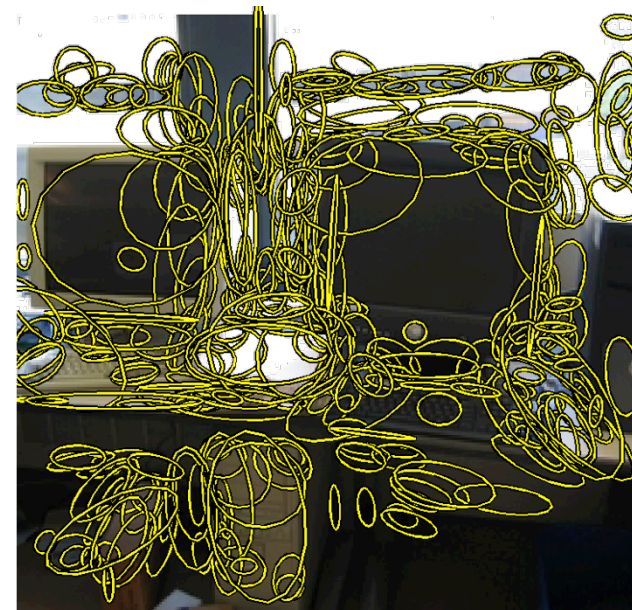
## Operation of code

To run the demos:

start Microsoft Outlook We... 未名空间(mitbbs.co... A demonstration of b... ICCV2005

# From Images to Features

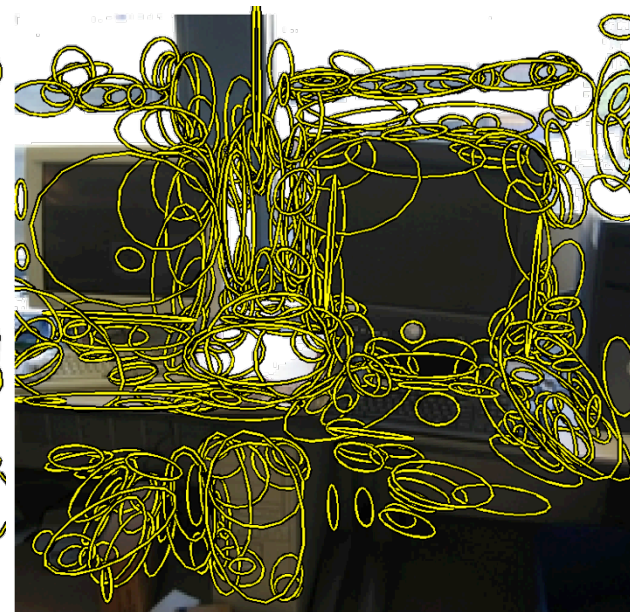
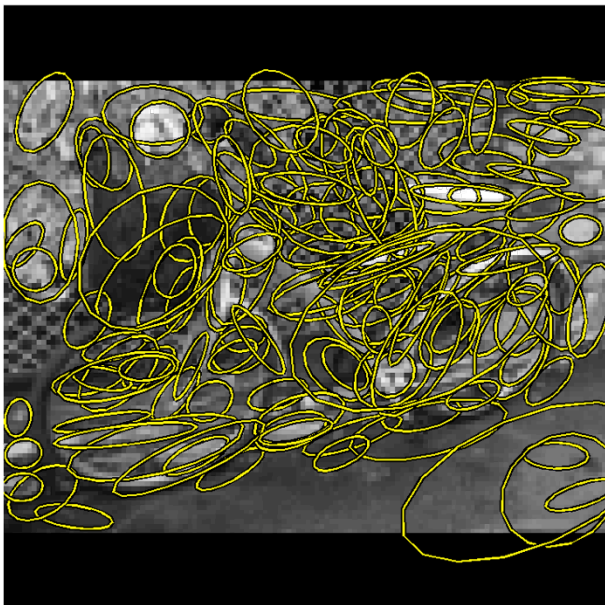
- Pixels are very sensitive to changes in lighting & pose
- Instead represent image as *affine covariant regions*:
  - Harris affine invariant regions (corners & edges)
  - Maximally stable extremal regions (segmentation)



Software provided by  
*Oxford Visual Geometry Group*

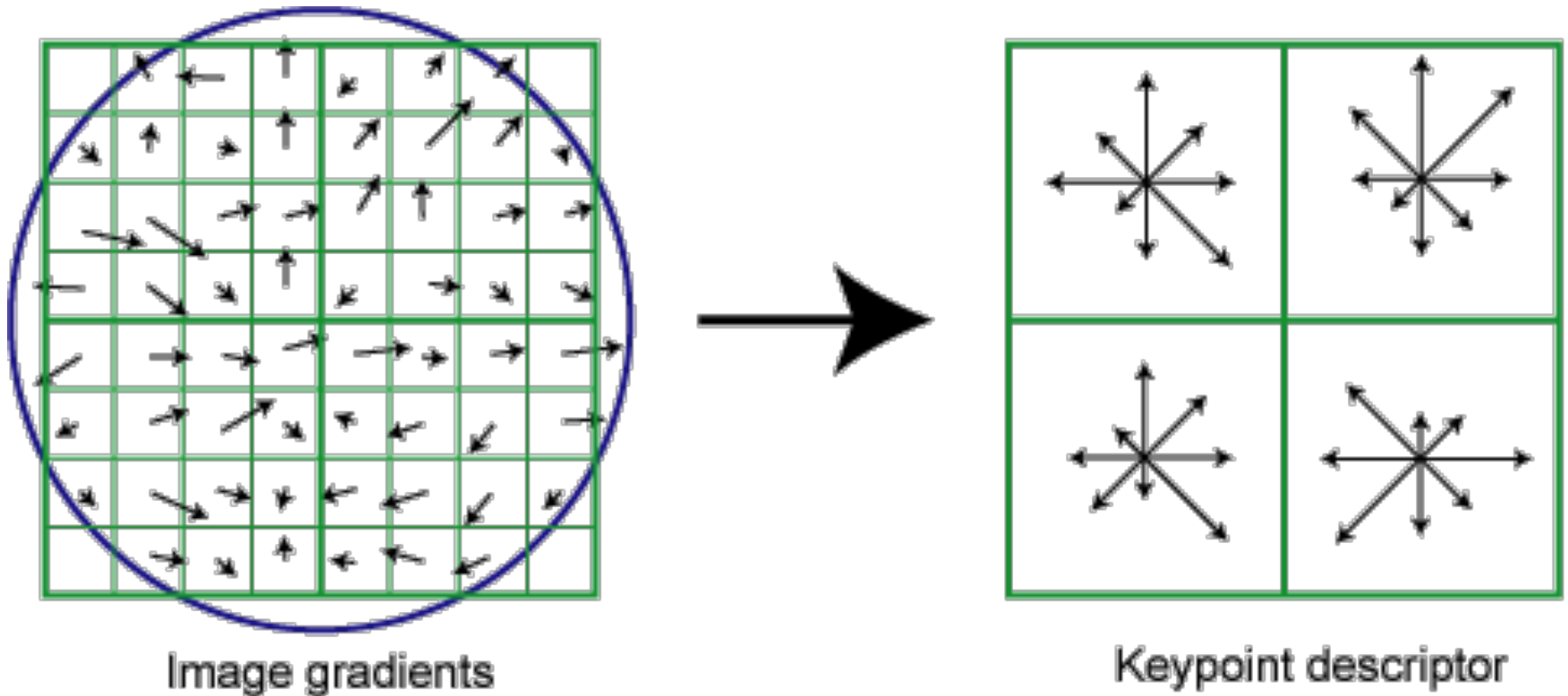


# Sample Detected Features



# Describing Feature Appearance

- **SIFT**: Scale Invariant Feature Transform
- Normalized histogram of orientation energy in each affinely adapted region (128-dim.)



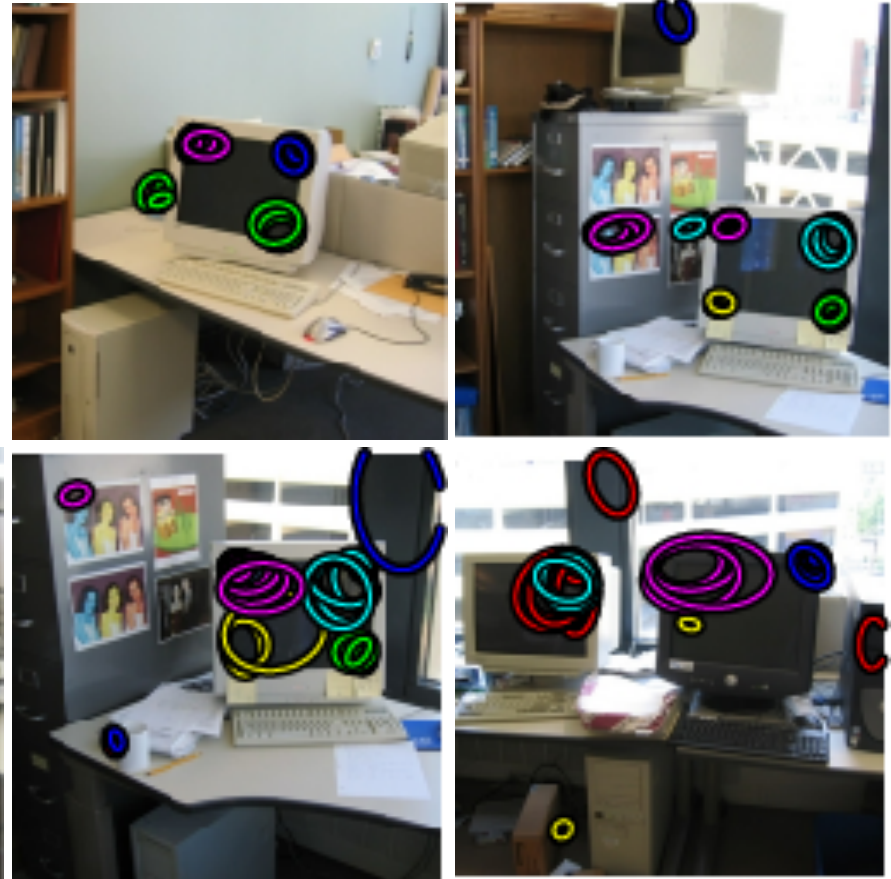


# A Discrete Feature Vocabulary

- Using all training images, build a dictionary via K-means clustering (~1000 words)
- Map each SIFT descriptor to nearest word

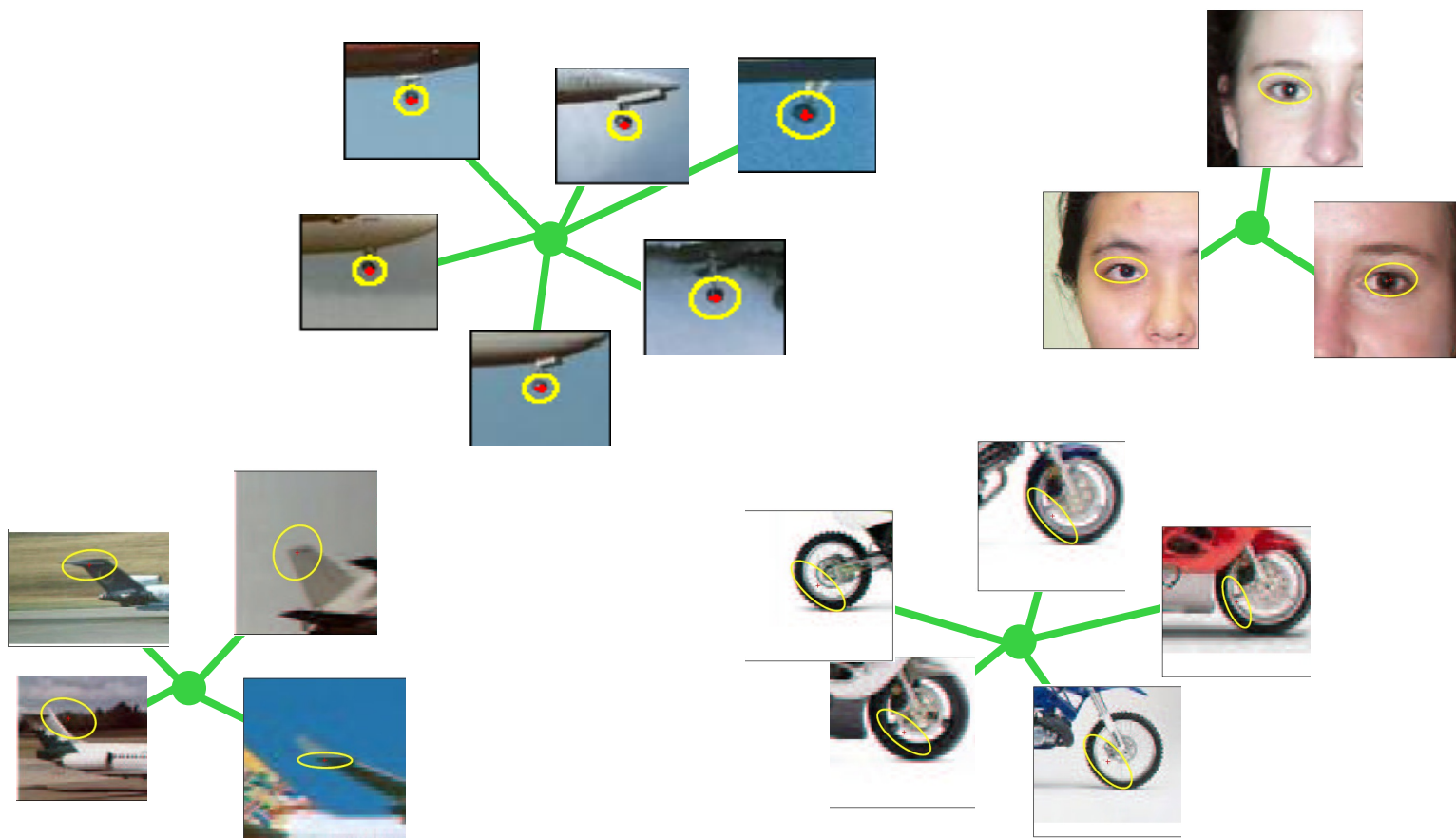
$w_{ji}$   $\longrightarrow$  appearance of feature  $i$  in image  $j$

$y_{ji}$   $\longrightarrow$  2D position of feature  $i$  in image  $j$



# Form dictionary

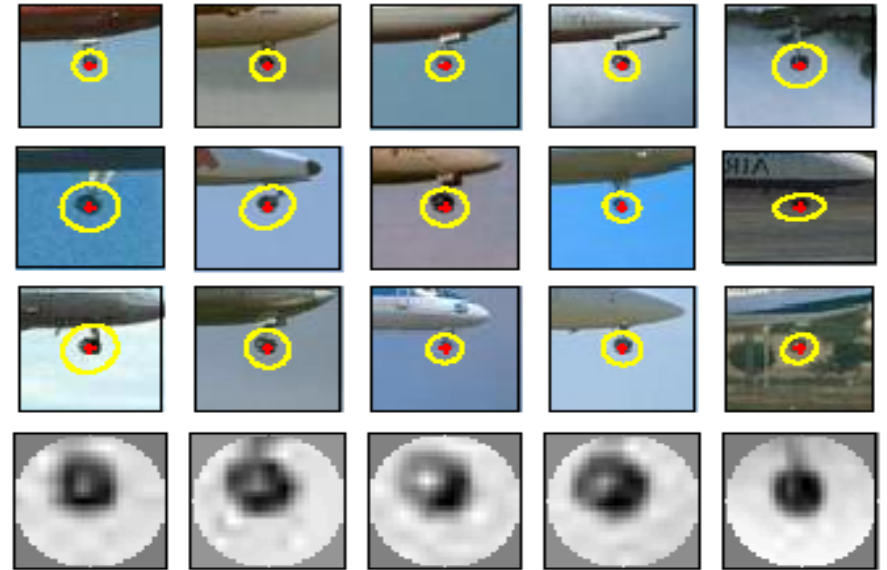
Build visual vocabulary by k-means clustering  
SIFT descriptors (K~2,000)



# Example regions assigned to the same dictionary cluster



Cluster 1



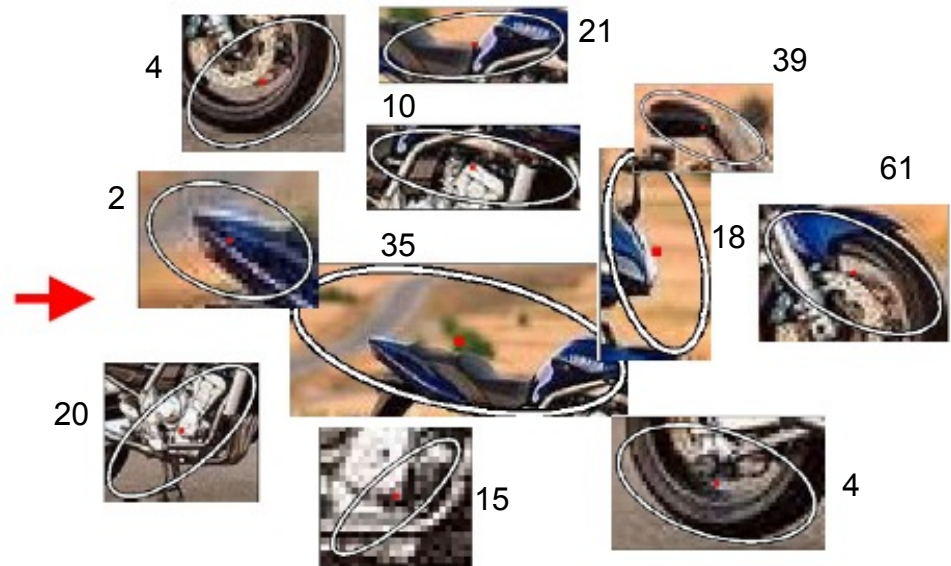
Cluster 2

# Representing an image with visual words

Sivic & Zisserman '03



Interest regions



Visual words



# System overview



Input image

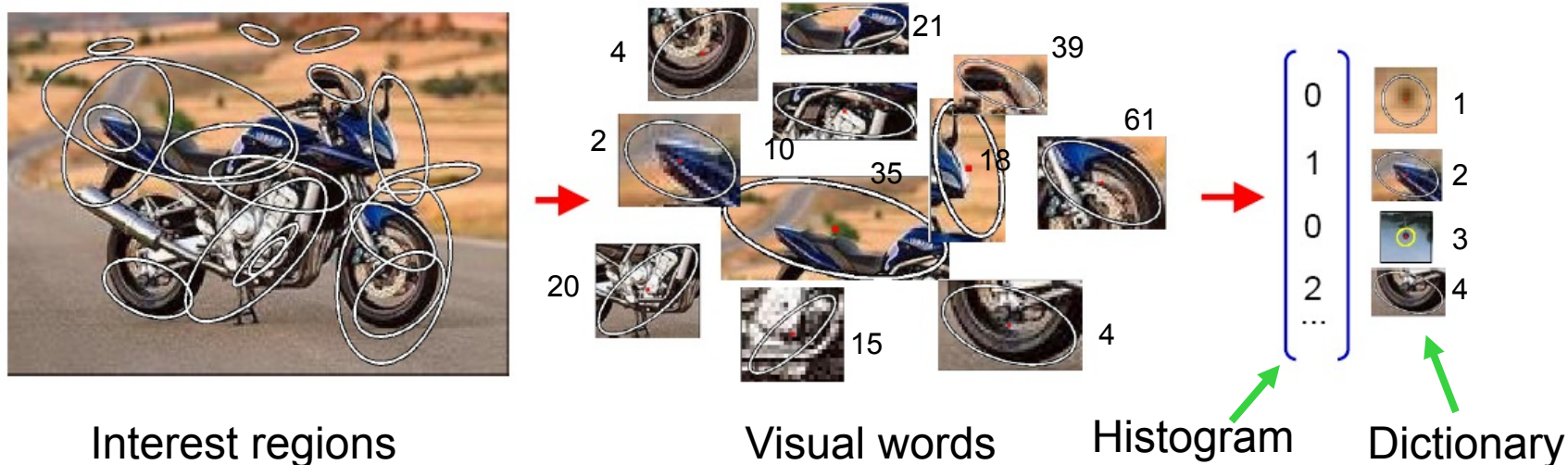


Compute visual words



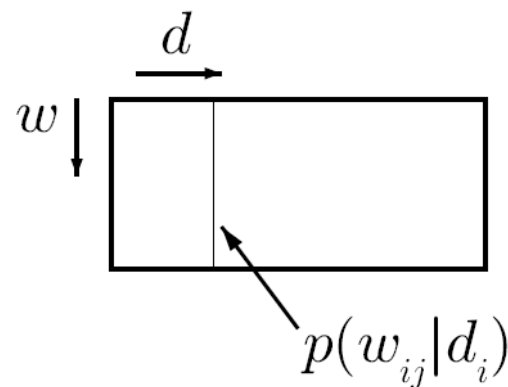
Discover visual topics

# Bag of words



Stack visual word histograms  
as columns in matrix

Throw away spatial information!





# Documents collection

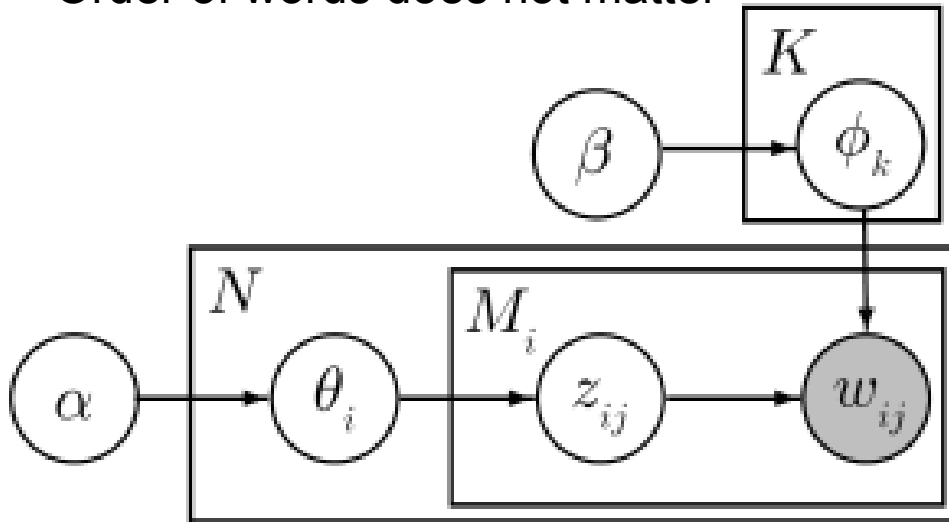
Co-occurrence table:

i	Number of times word $i$ appears on document/image $j$
	j

# Latent Dirichlet Allocation (LDA)

Blei, et al. 2003

- LDA model assumes exchangeability
- Order of words does not matter



$w_{ij}$  - words

$z_{ij}$  - topic assignments

$\theta_i$  - topic mixing weights

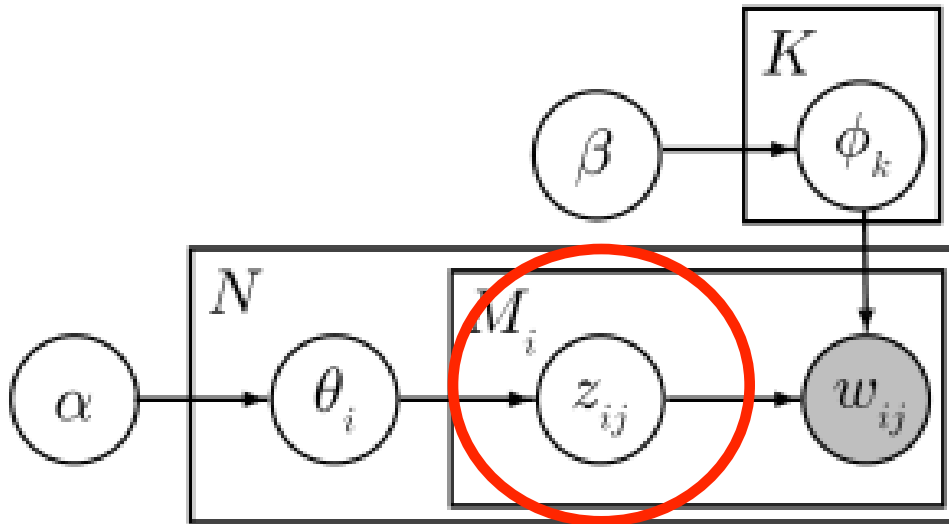
$\Phi_k$  - word mixing weights

$$z_{ij}|\theta_i \sim \theta_i \quad \theta_i|\alpha \sim \text{Dirichlet}(\alpha)$$

$$w_{ij}|z_{ij} = k, \phi \sim \phi_k \quad \phi_k|\beta \sim \text{Dirichlet}(\beta)$$

$$p(w_{ij}) \propto \sum_{k=1}^K p(w_{ij}|z_{ij} = k, \phi_k) p(z_{ij} = k|\theta_i)$$

# Inference



$w_{ij}$  - words

$z_{ij}$  - topic assignments

$\theta_i$  - topic mixing weights

$\phi_k$  - word mixing weights

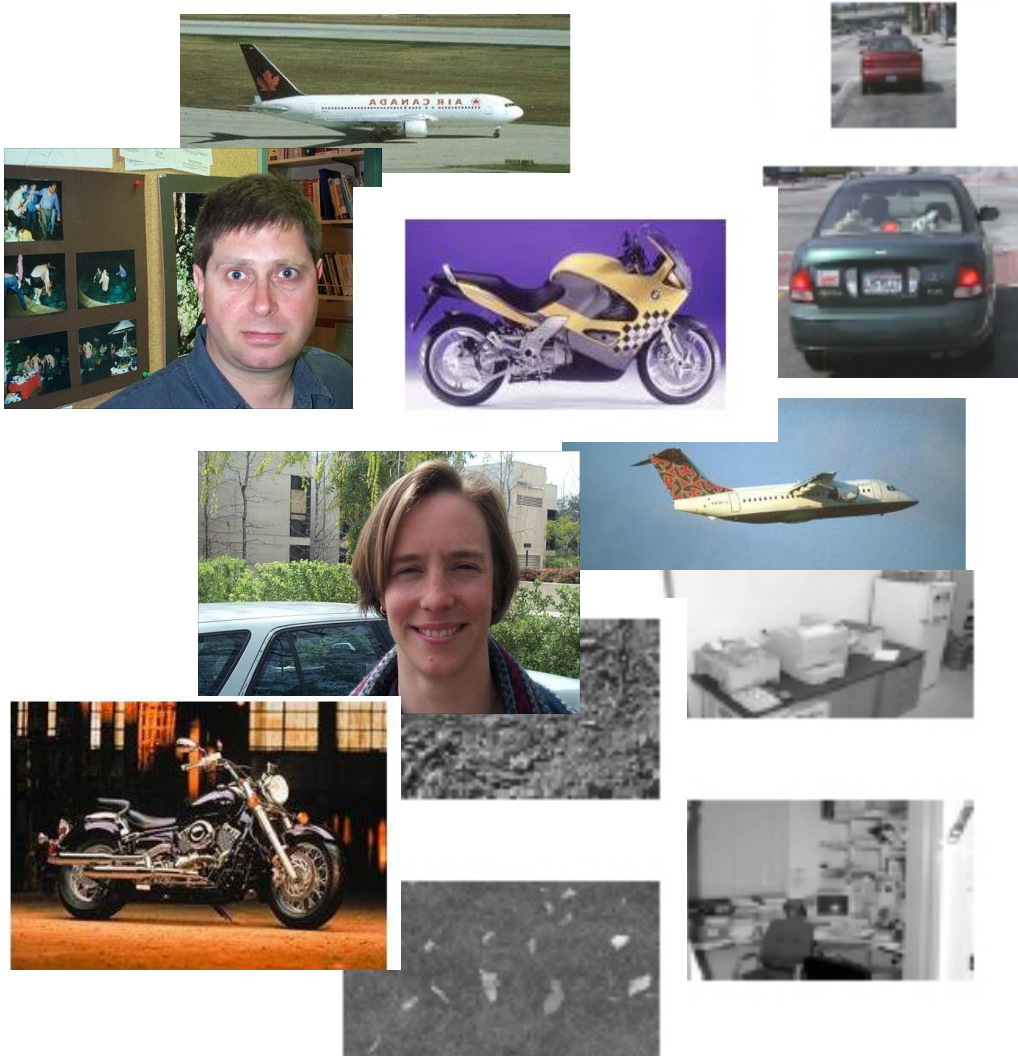
Use Gibbs sampler to sample topic assignments

[Griffiths & Steyvers 2004]

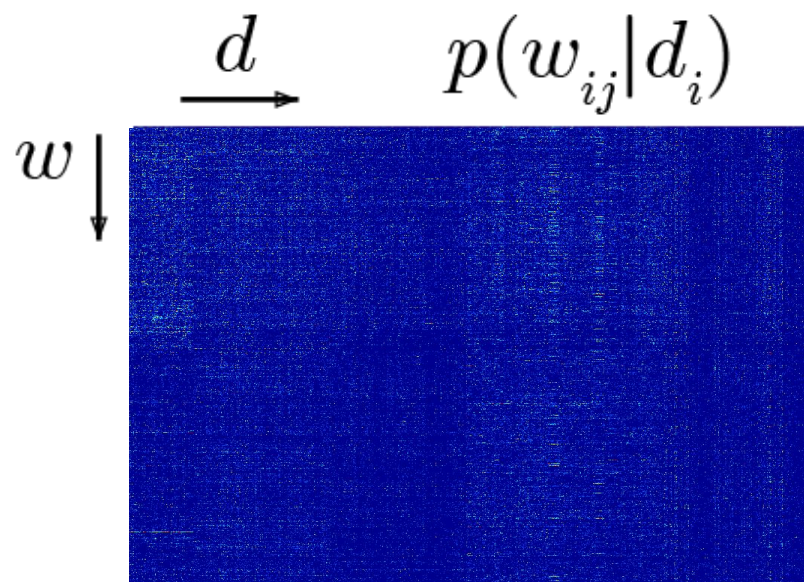
$$z_{ij} \sim p(z_{ij} = k | w_{ij} = v, w_{\setminus(ij)}, z_{\setminus(ij)}, \alpha, \beta)$$

- Only need to maintain counts of topic assignments
- Sampler typically converges in less than 50 iterations
- Run time is less than an hour

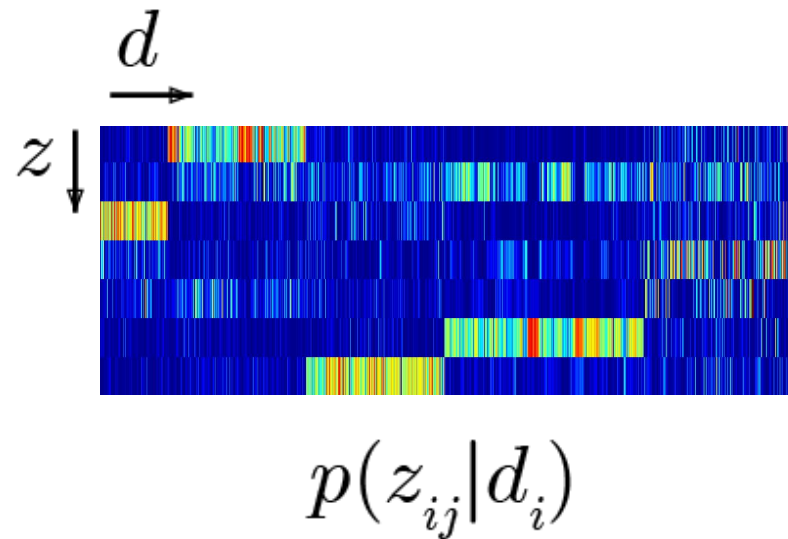
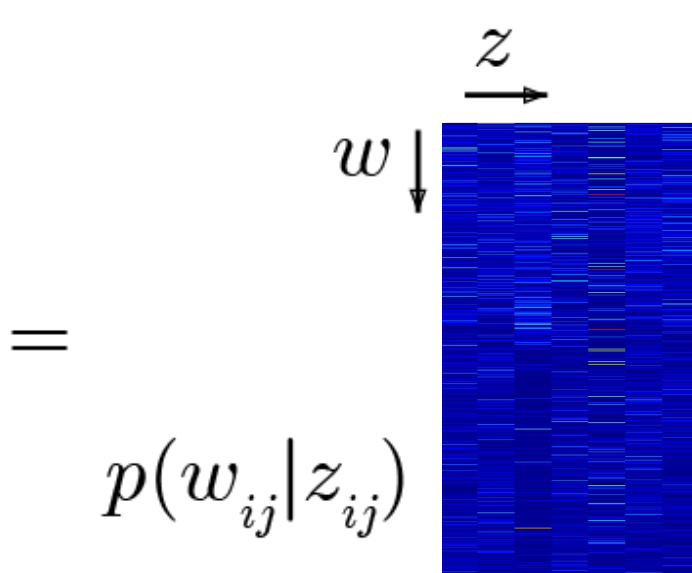
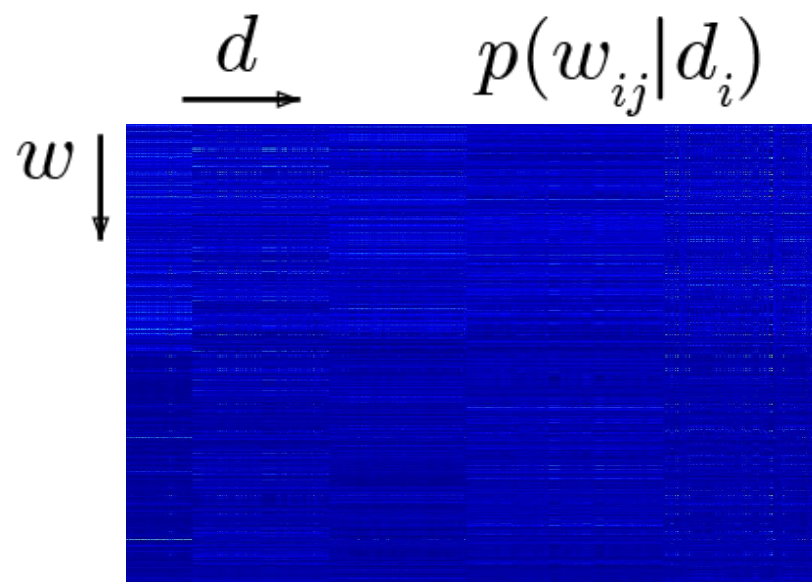
# Apply to Caltech 4 + background images



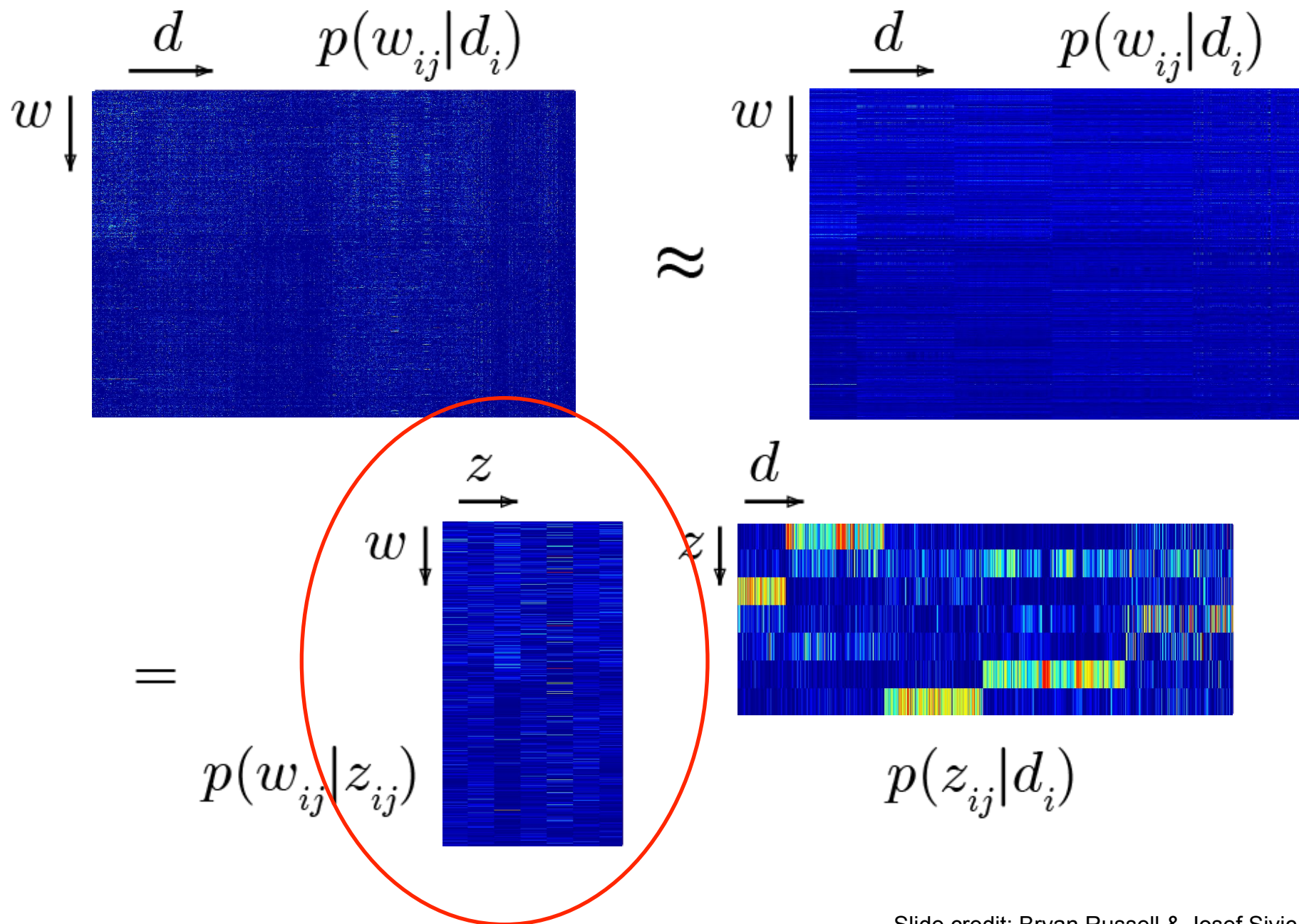
Faces	435
Motorbikes	800
Airplanes	800
Cars (rear)	1155
Background	900
Total:	4090



$\approx$



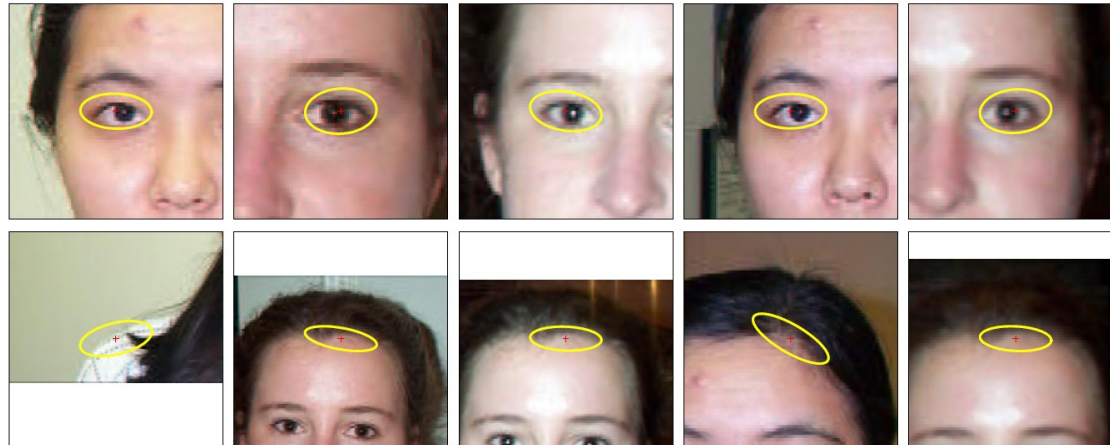




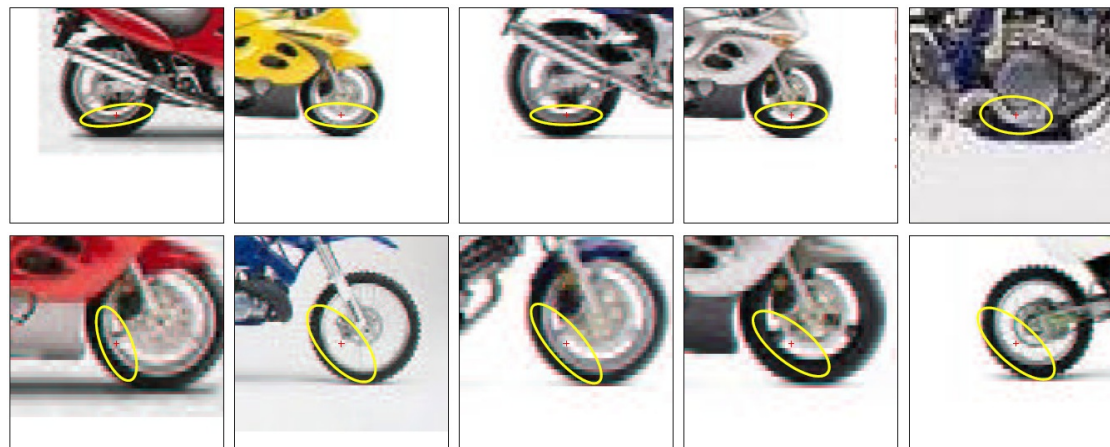


# Most likely words given topic

Topic 1



Topic 2



# Most likely words given topic

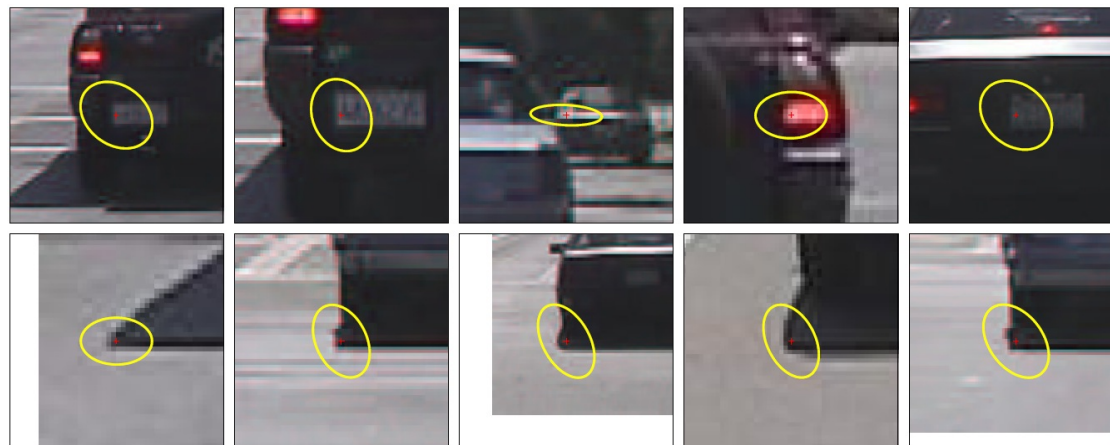
Topic 3



Word 1

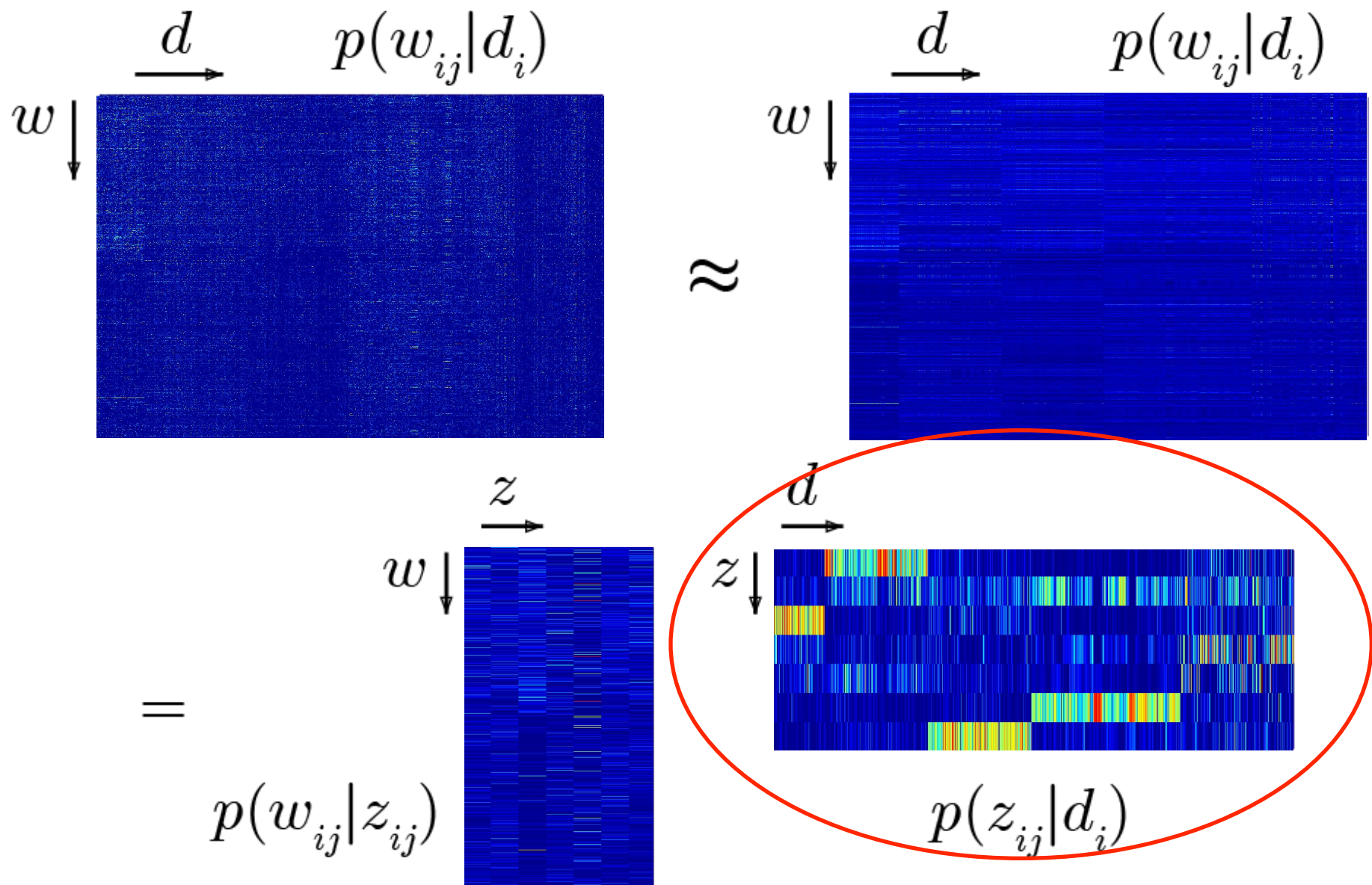
Word 2

Topic 4



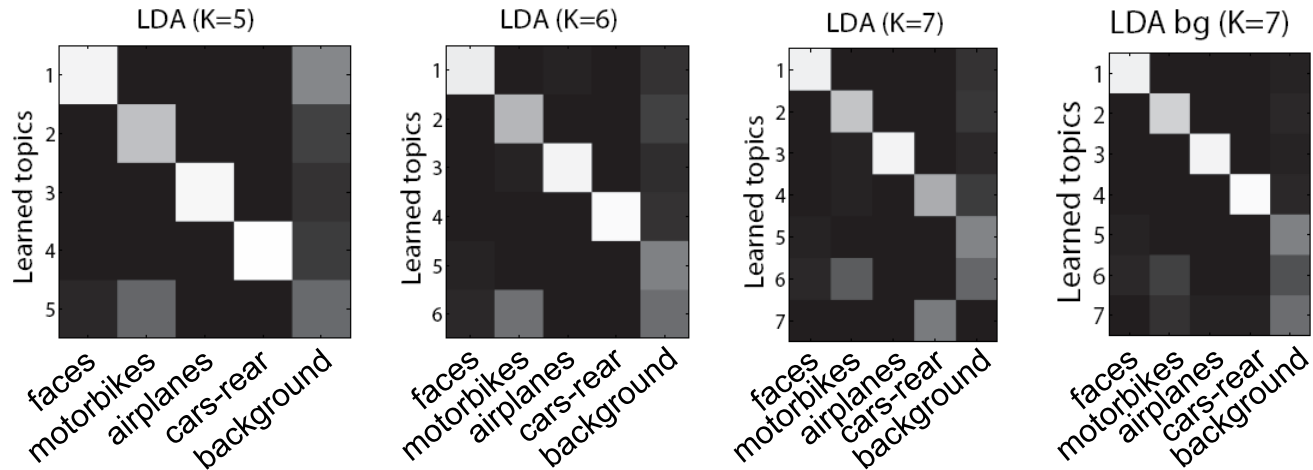
Word 1

Word 2



# Image clustering

Confusion matrices:

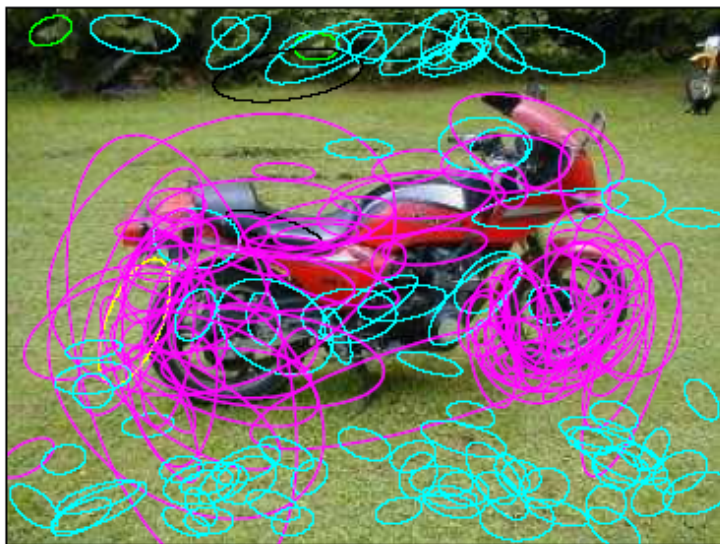


Average confusion:

Expt.	Categories	T	LDA		pLSA		KM baseline	
			%	#	%	#	%	#
(1)	4	4	97	86	98	70	72	908
(2)	4 + bg	5	78	931	78	931	56	1820
(2)*	4 + bg	6	84	656	76	1072	—	—
(2)*	4 + bg	7	78	1007	83	768	—	—
(2)*	4 + bg-fxd	7	90	330	93	238	—	—



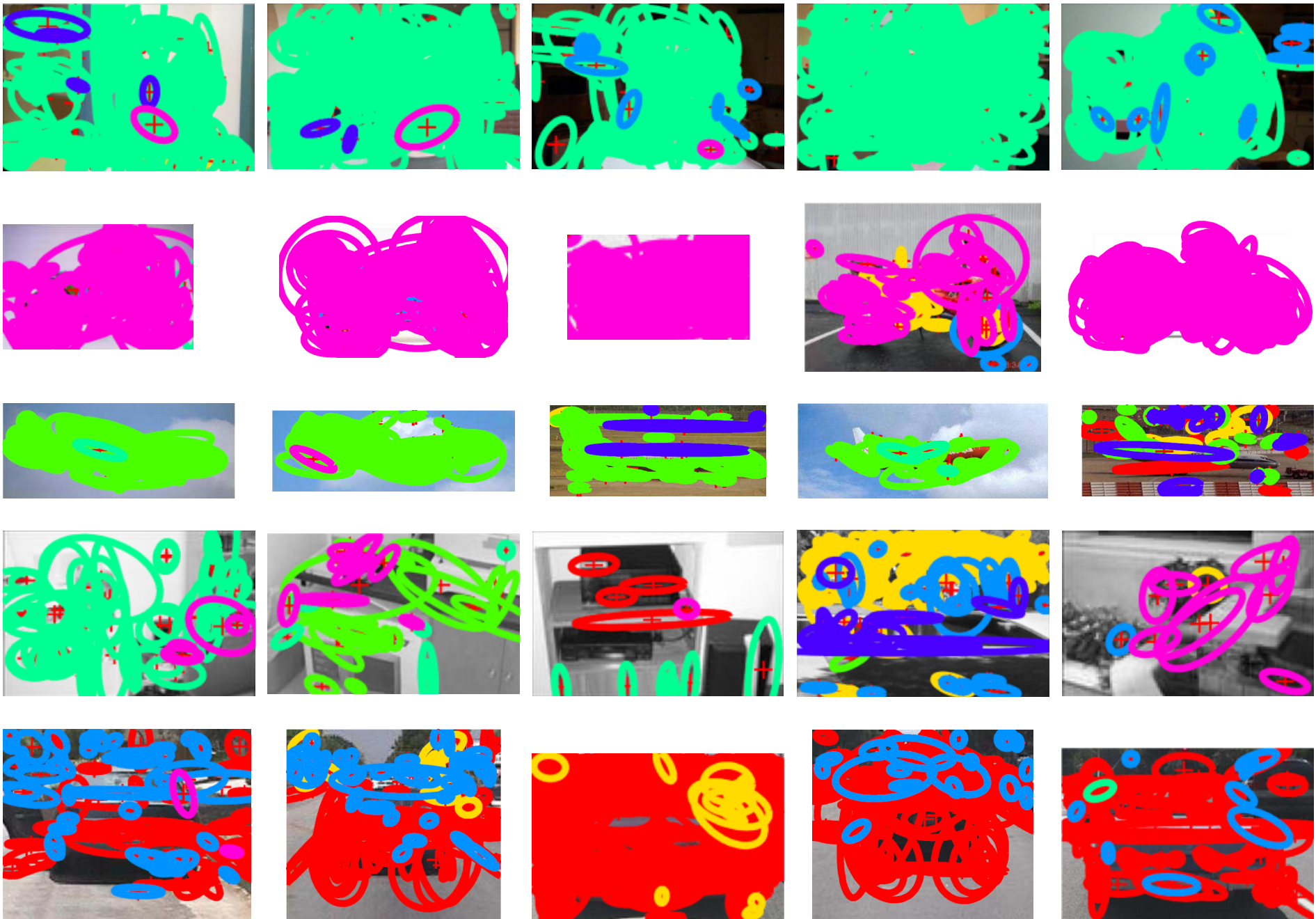
# Image as a mixture of topics (objects)





Slide credit: Bryan Russell & Josef Sivic





Slide credit: Bryan Russell & Josef Sivic

# Beyond single classes

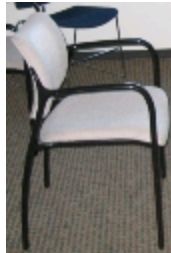
- Multiclass
- Multiview
- Datasets

# Beyond single classes

- **Multiclass**
- Multiview
- Datasets

# Shared features

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



...

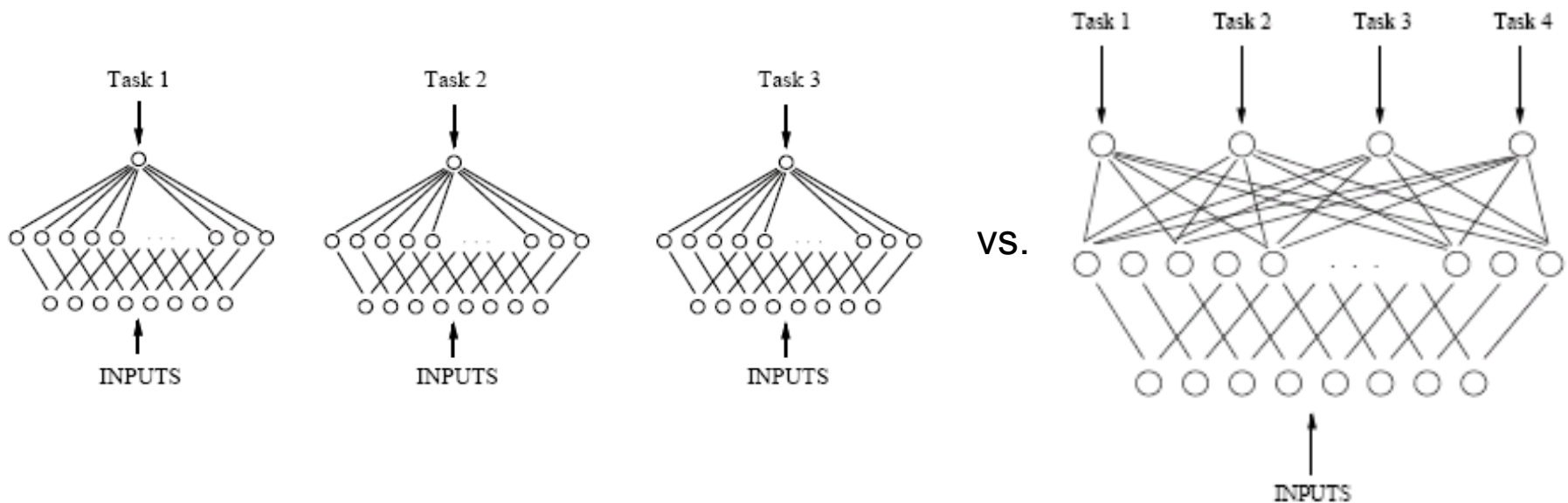


- Can we transfer knowledge from one object to another?
- Are the shared properties interesting by themselves?

# Multitask learning

**R. Caruana. Multitask Learning. ML 1997**

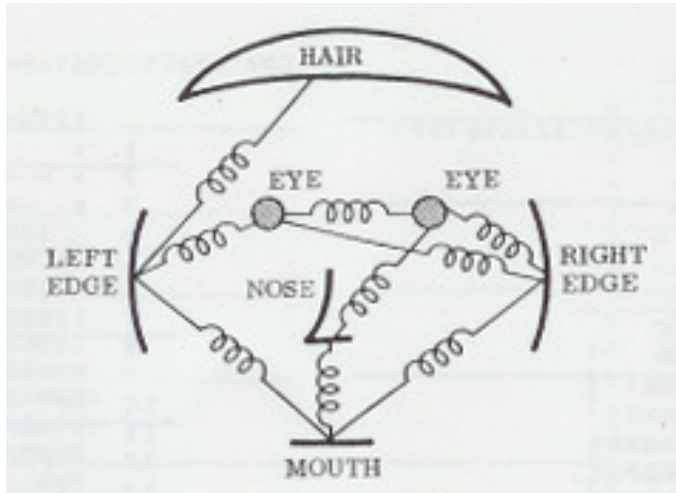
“MTL improves generalization by leveraging the domain-specific information contained in the training signals of *related* tasks. It does this by training tasks in parallel while using a shared representation”.



Sejnowski & Rosenberg 1986; Hinton 1986; Le Cun et al. 1989; Suddarth & Kergosien 1990; Pratt et al. 1991; Sharkey & Sharkey 1992; ...

# Sharing in constellation models

(next Wednesday)



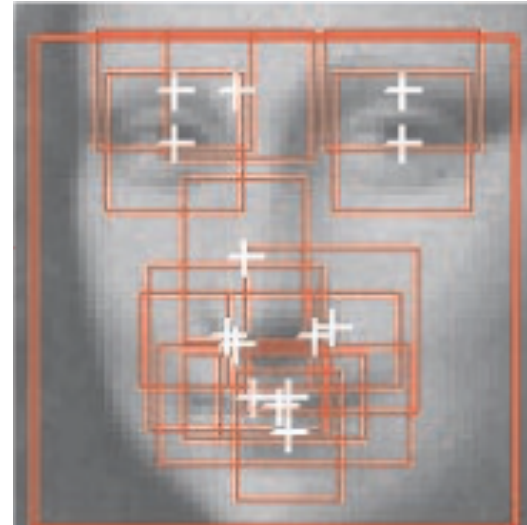
## Pictorial Structures

*Fischler & Elschlager, IEEE Trans. Comp. 1973*



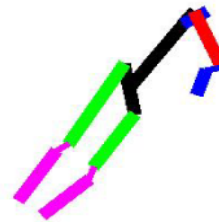
## Constellation Model

*Fergus, Perona, & Zisserman, CVPR 2003*



## SVM Detectors

*Heisele, Poggio, et. al., NIPS 2001*



## Model-Guided Segmentation

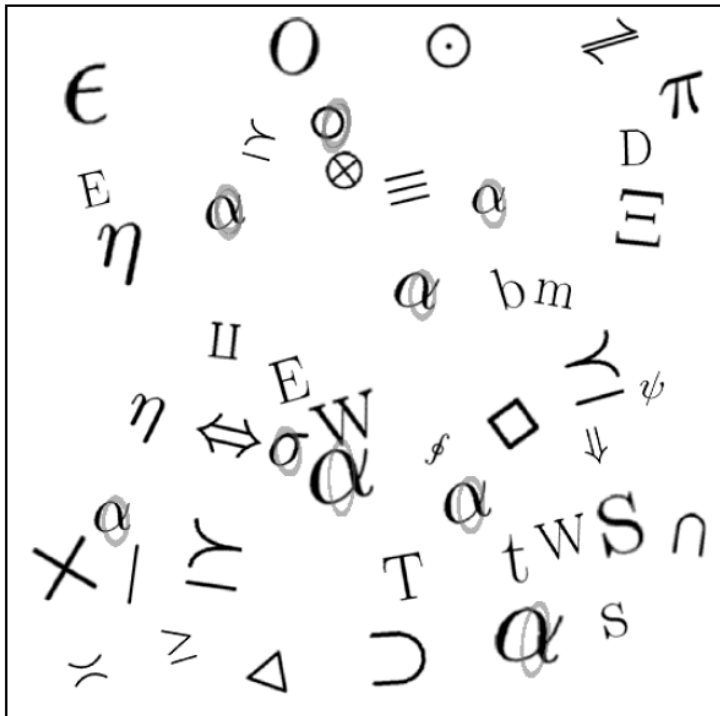
*Mori, Ren, Efros, & Malik, CVPR 2004*



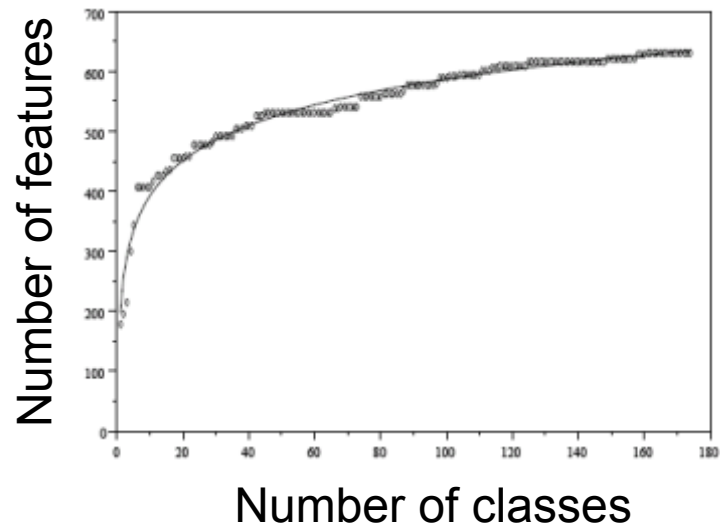
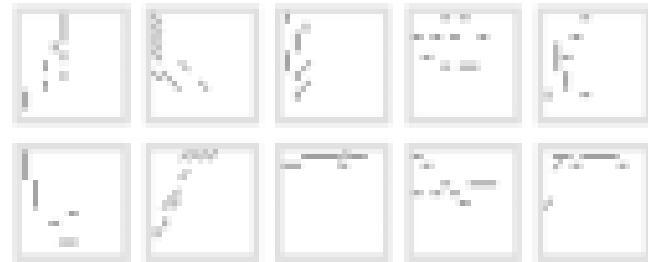
# Reusable Parts

Krempf, 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.

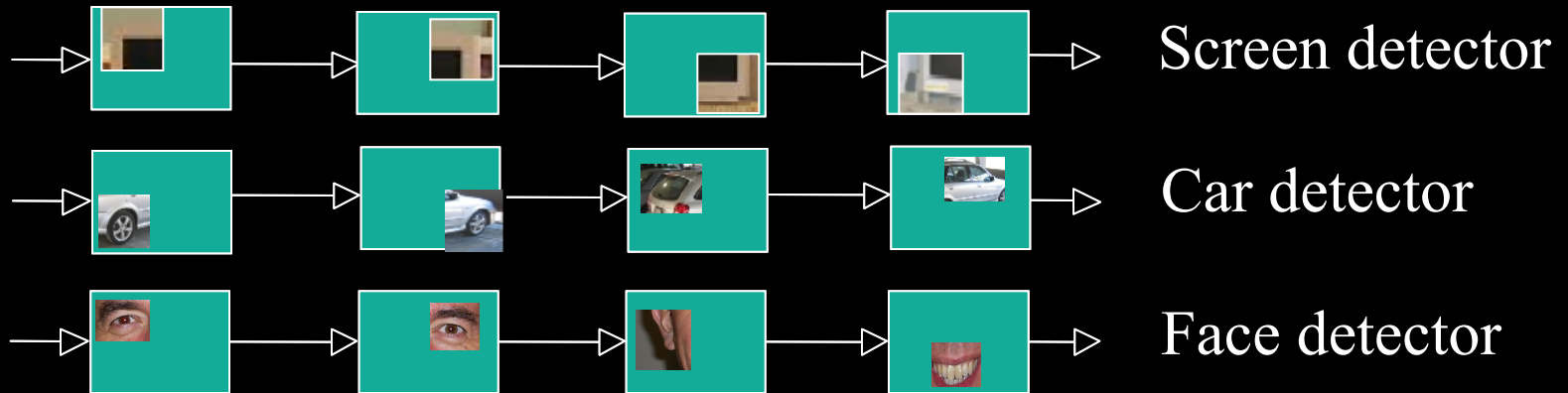


Examples of reused parts

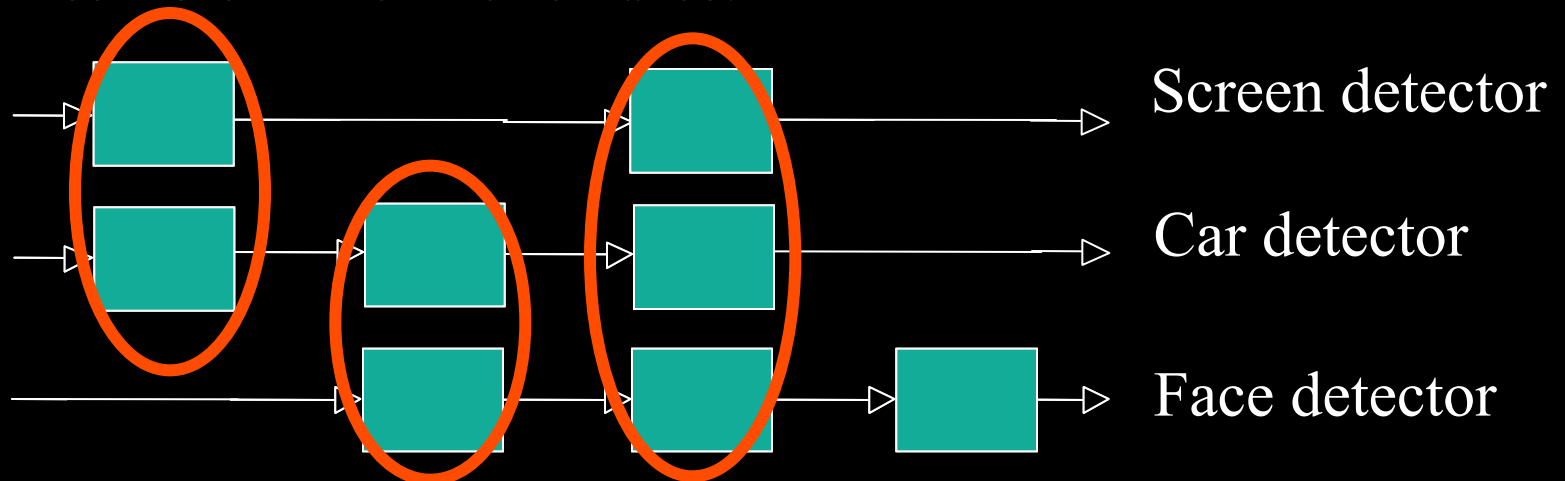


# Additive models and boosting

- Independent binary classifiers:

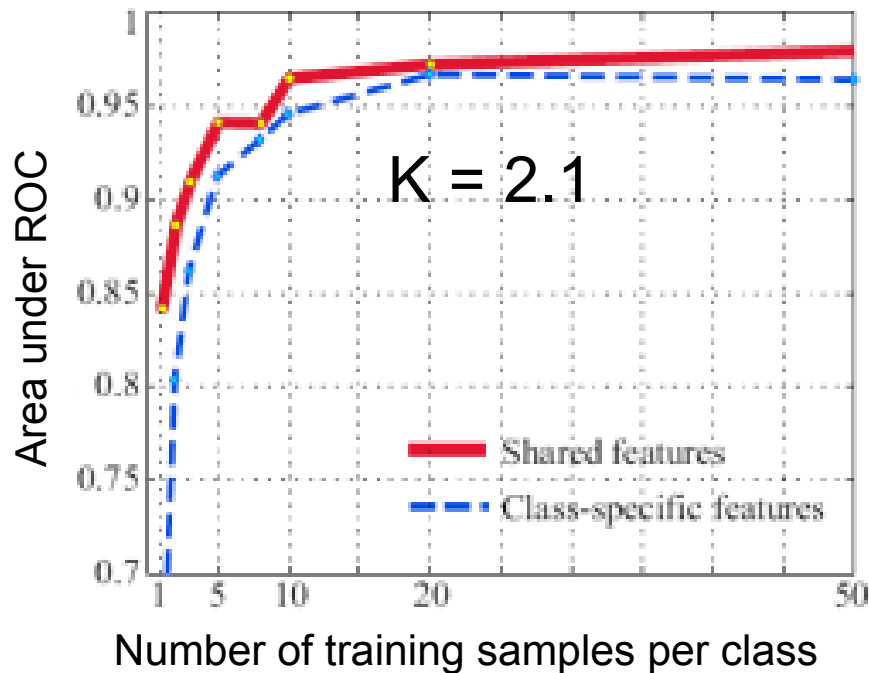
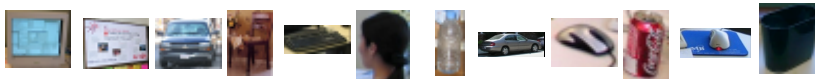


- Binary classifiers that share features:

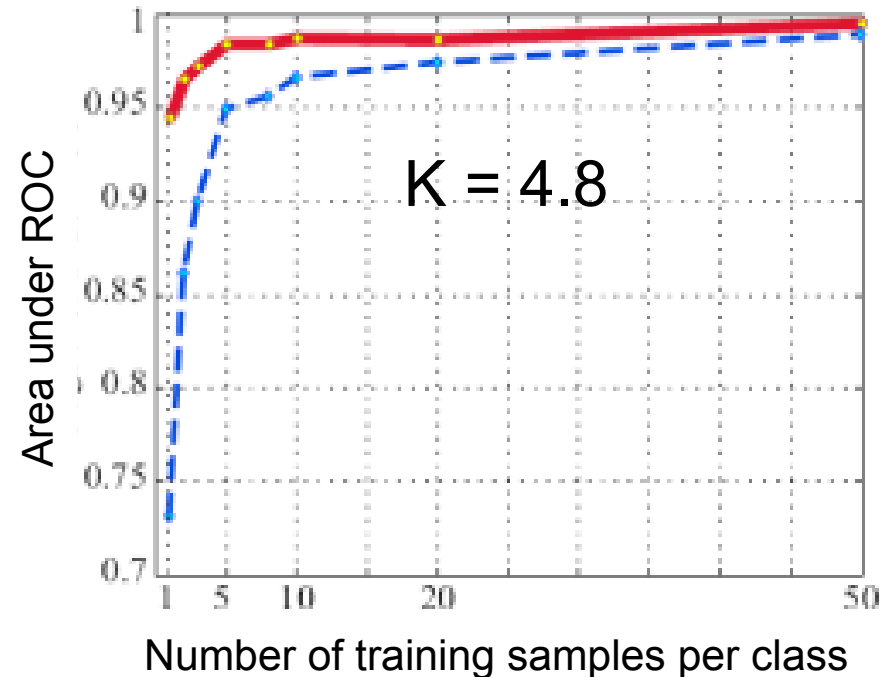


# Generalization as a function of object similarities

12 unrelated object classes



12 viewpoints



# 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.

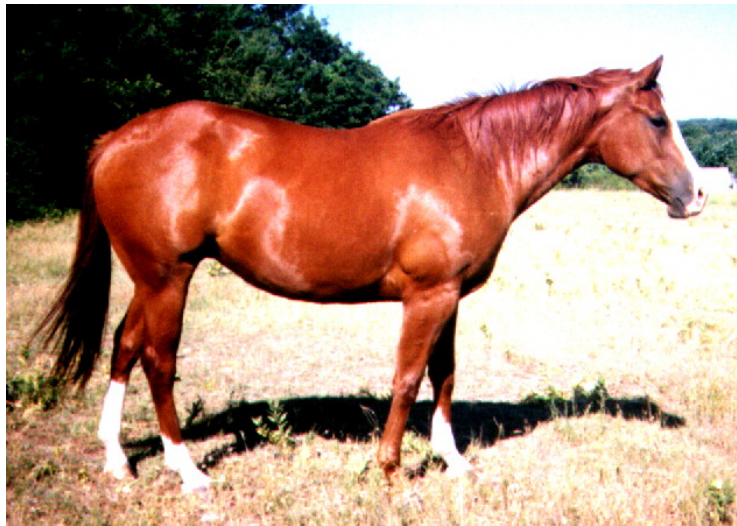


# Class experiment

**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



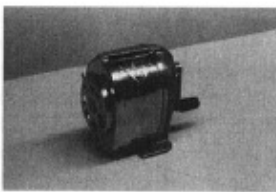
CLOCK



TELEPHONE



HOUSE



PENCIL SHARPENER



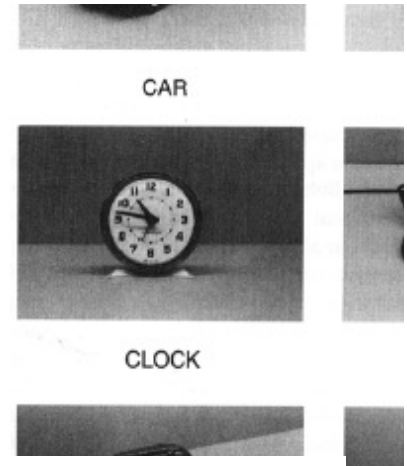
SHOE



IRON

# Canonical Viewpoint

Clocks are preferred as purely frontal

[Search Images](#)[Search the Web](#)[Advanced Image Search](#)  
[Preferences](#)

Moderate SafeSearch is on

Images Showing:

Results 1 - 18 of about 38,300,000 for

Related searches: [cartoon clock](#) [clock clipart](#) [alarm clock](#) [clock face](#)



clock character

359 x 344 - 4k - gif

[school.discoveryeducation.com](http://school.discoveryeducation.com)



Wind-up alarm clocks have been

...

346 x 510 - 22k - jpg

[electronics.howstuffworks.com](http://electronics.howstuffworks.com)



Artistic Clock And Wall Clock

360 x 360 - 18k - jpg

[www.global-b2b-network.com](http://www.global-b2b-network.com)



... mechanical clock

screensaver.

640 x 480 - 53k - jpg

[davinciautomata.wordpress.com](http://davinciautomata.wordpress.com)



If it is 3 o'clock and we add 5 ...

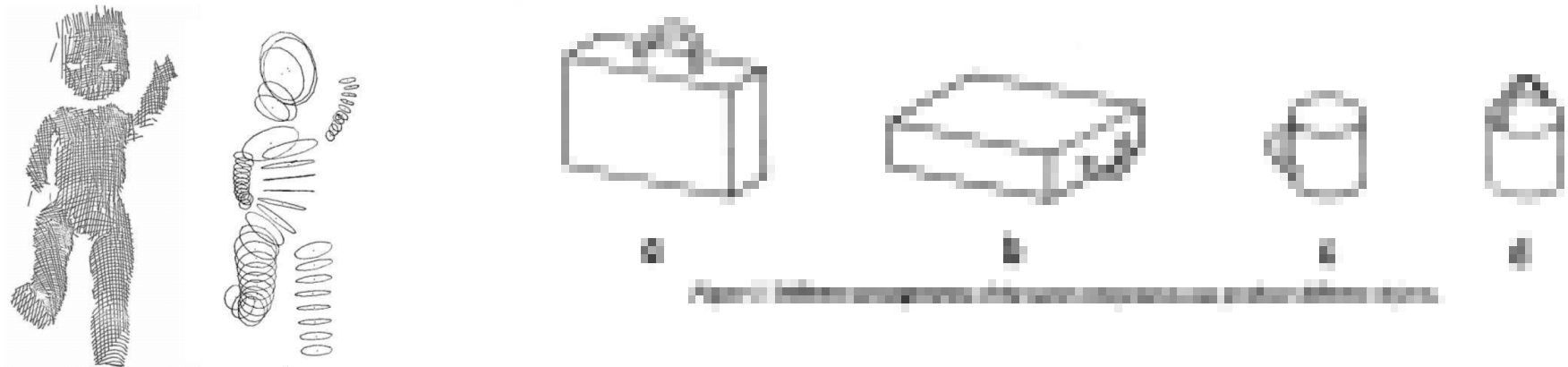
305 x 319 - 4k - gif

[www-math.cudenver.edu](http://www-math.cudenver.edu)

[ [More from](#)  
[www-math.cudenver.edu](http://www-math.cudenver.edu) ]

# 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

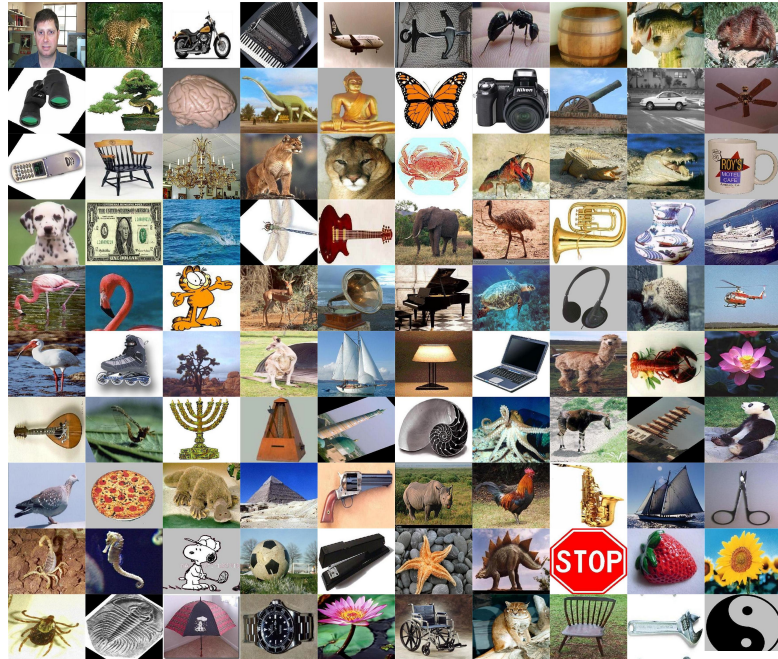
*Animal:* bird, cat, cow, dog, horse, sheep

*Vehicle:* aeroplane, bicycle, boat, bus, car, motorbike, train

*Indoor:* bottle, chair, dining table, potted plant, sofa, tv/monitor



# Caltech 101 and 256

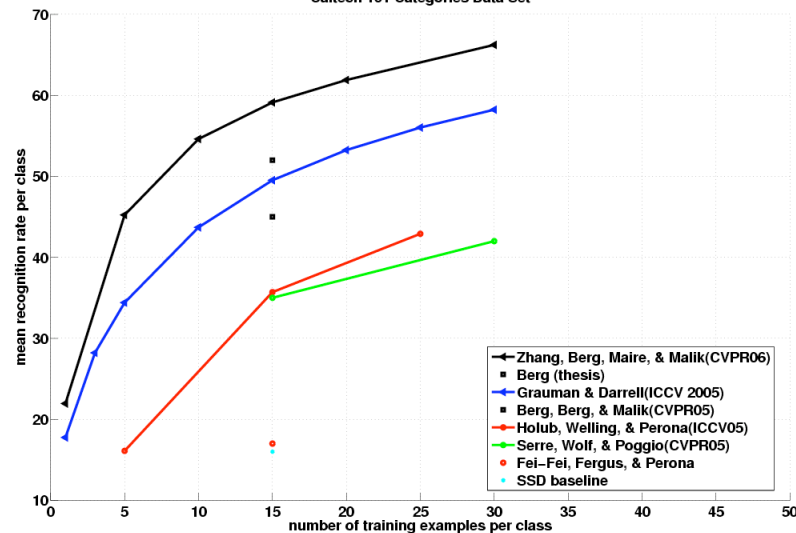


Fei-Fei, Fergus, Perona, 2004



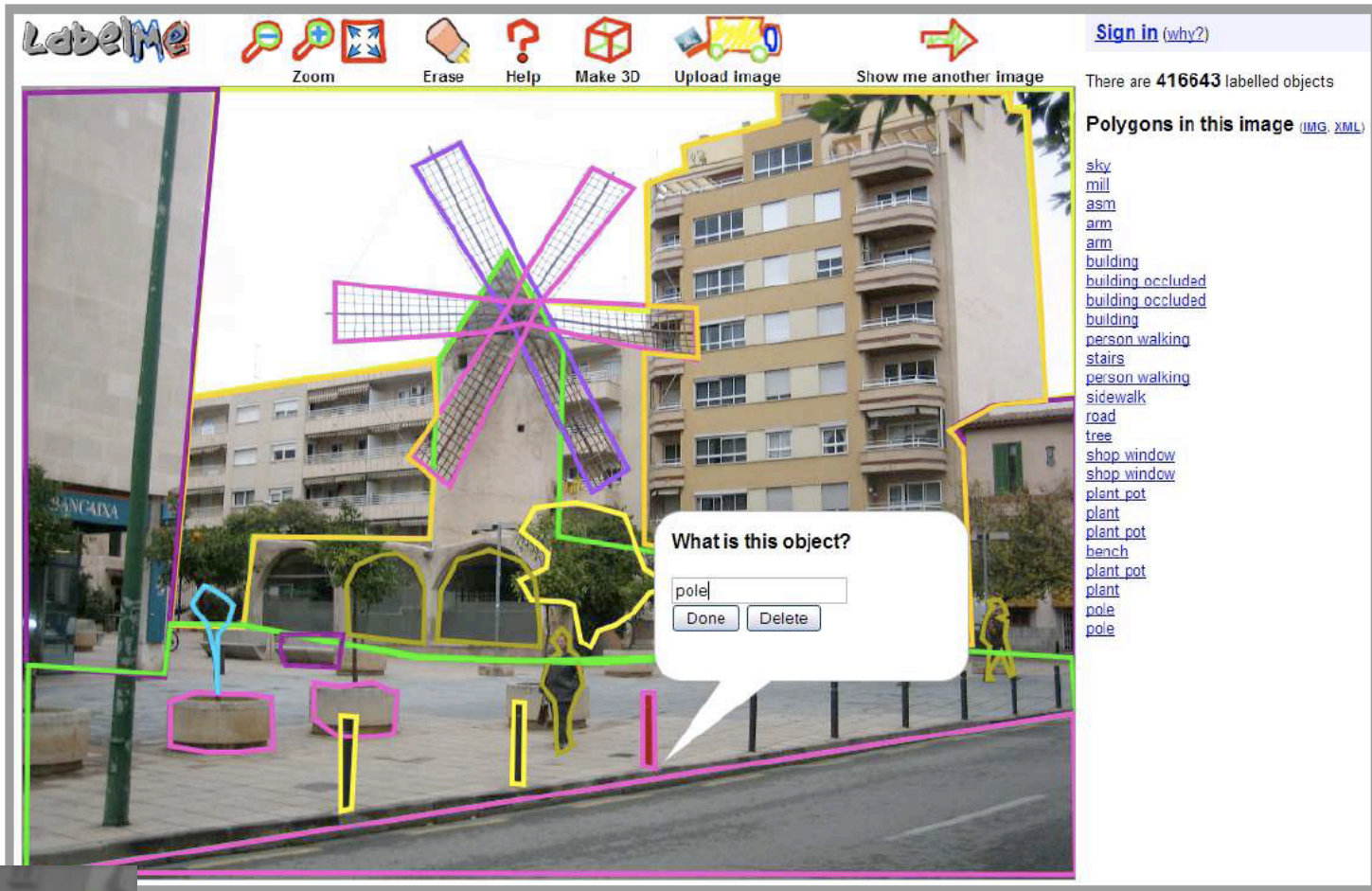
Griffin, Holub, Perona, 2007

Caltech 101 Categories Data Set





# LabelMe



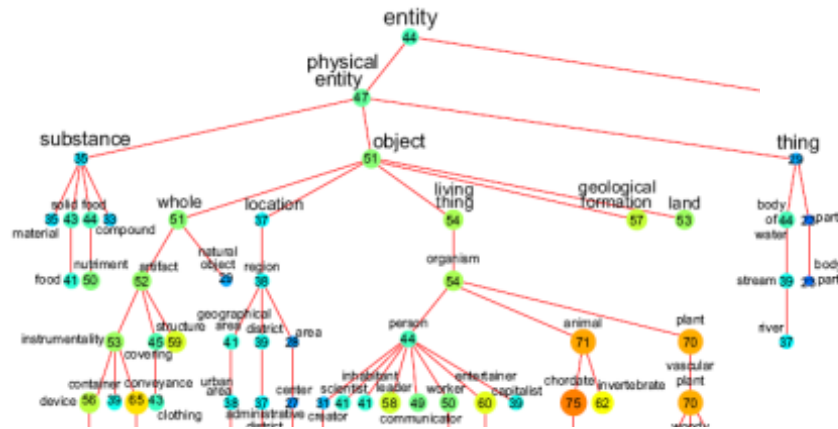
**Tool went online July 1st, 2005**  
**530,000 object annotations collected**

Labelme.csail.mit.edu

# 80.000.000 images

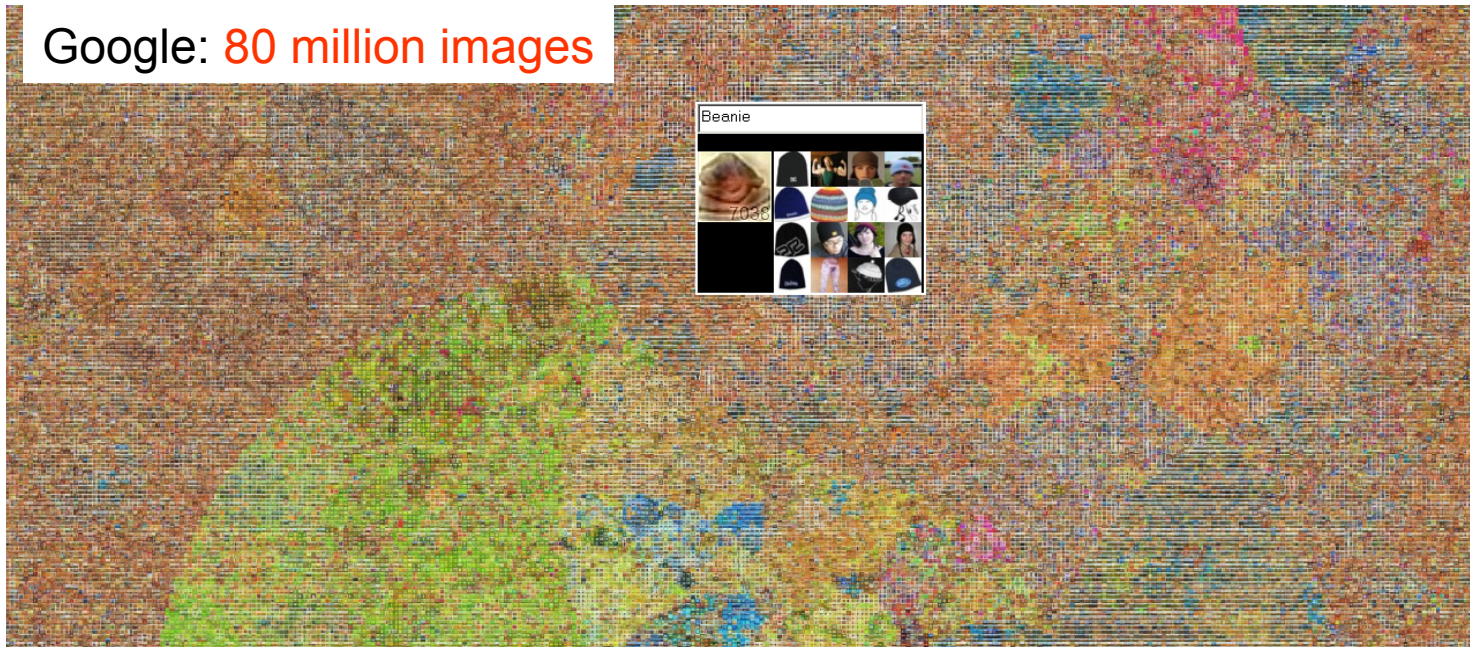
75.000 non-abstract nouns from WordNet

7 Online image search engines



And after 1 year downloading images

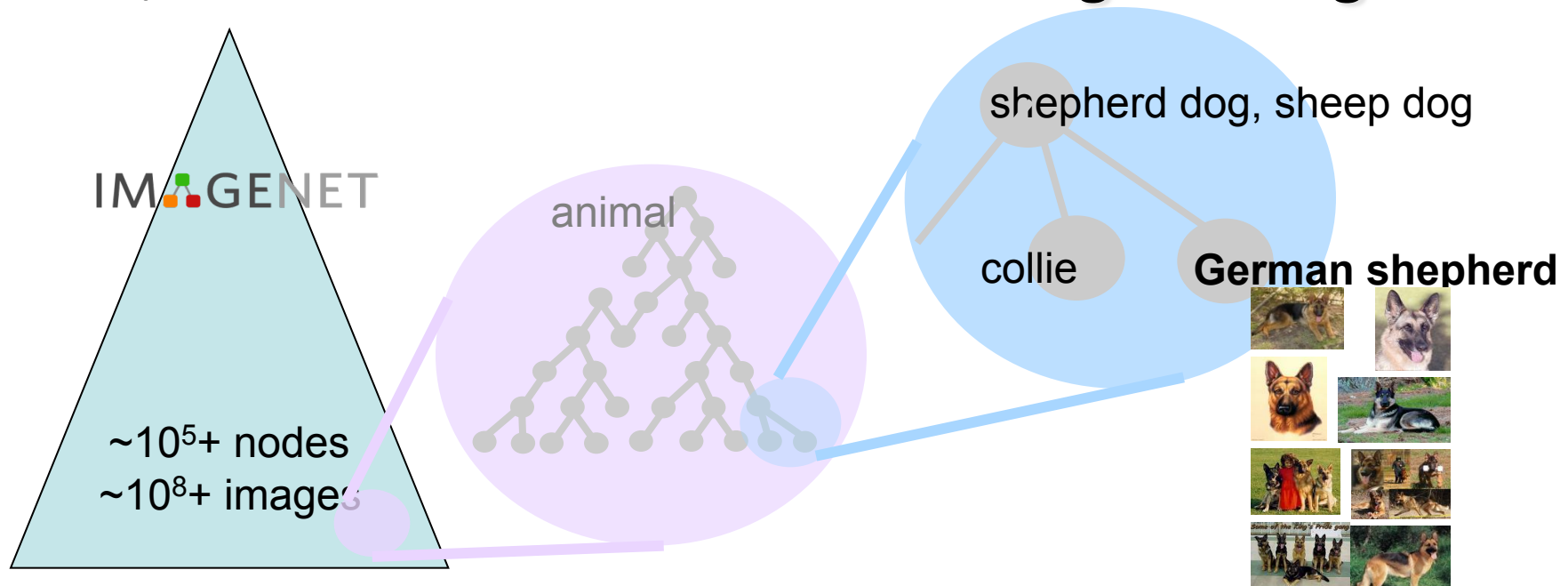
Google: 80 million images





# IMAGENET

- 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**





mug

Search

SafeSearch moderate

About 10,100,000 results (0.09 seconds)

Advanced search

Google mugs

# Dataset biases

59¢ Logo Coffee Mugs

[www.DiscountMugs.com](http://www.DiscountMugs.com) Lead Free & Dishwasher Safe. Save 40-50%. No Catch. Factory Direct !

Custom Mugs On Sale

[www.Vistaprint.com](http://www.Vistaprint.com) Order Now & Save 50% On Custom Mugs No Minimums. Upload Photos & Logos.

Promotional Mugs from 69¢

[www.4imprint.com/Mugs](http://www.4imprint.com/Mugs) Huge Selection of Style Colors- Buy 72 Mugs @ \$1.35 ea-24hr Service

Sponsored

Related searches: [white mug](#) [coffee mug](#) [mug root beer](#) [mug shot](#)



**Representational**  
500 × 429 - 91k - jpg  
[eagereyes.org](#)  
[Find similar images](#)



**Ceramic Happy Face**  
300 × 300 - 77k - jpg  
[larose.com](http://larose.com)  
[Find similar images](#)



**Here I go then, trying**  
600 × 600 - 35k - jpg  
[beeper.wordpress.com](http://beeper.wordpress.com)  
[Find similar images](#)



**The Chalk Mug »**  
304 × 314 - 17k - jpg  
[coolest-gadgets.com](http://coolest-gadgets.com)  
[Find similar images](#)



**mug**  
300 × 279 - 54k - jpg  
[reynosawatch.org](http://reynosawatch.org)



**Bring your own**  
500 × 451 - 15k - jpg  
[cookstownunited.ca](http://cookstownunited.ca)  
[Find similar images](#)



**ceramic mug**  
980 × 1024 - 30k - jpg  
[diytrade.com](http://diytrade.com)



**Dual Purpose Drinking**  
490 × 428 - 16k - jpg  
[freshhome.com](http://freshhome.com)  
[Find similar images](#)



**This coffee mug.**  
300 × 300 - 22k - jp  
[gizmodo.com](http://gizmodo.com)  
[Find similar images](#)



**Back to Ceramic**  
400 × 400 - 8k - jpg  
[freshpromotions.com.au](http://freshpromotions.com.au)  
[Find similar images](#)



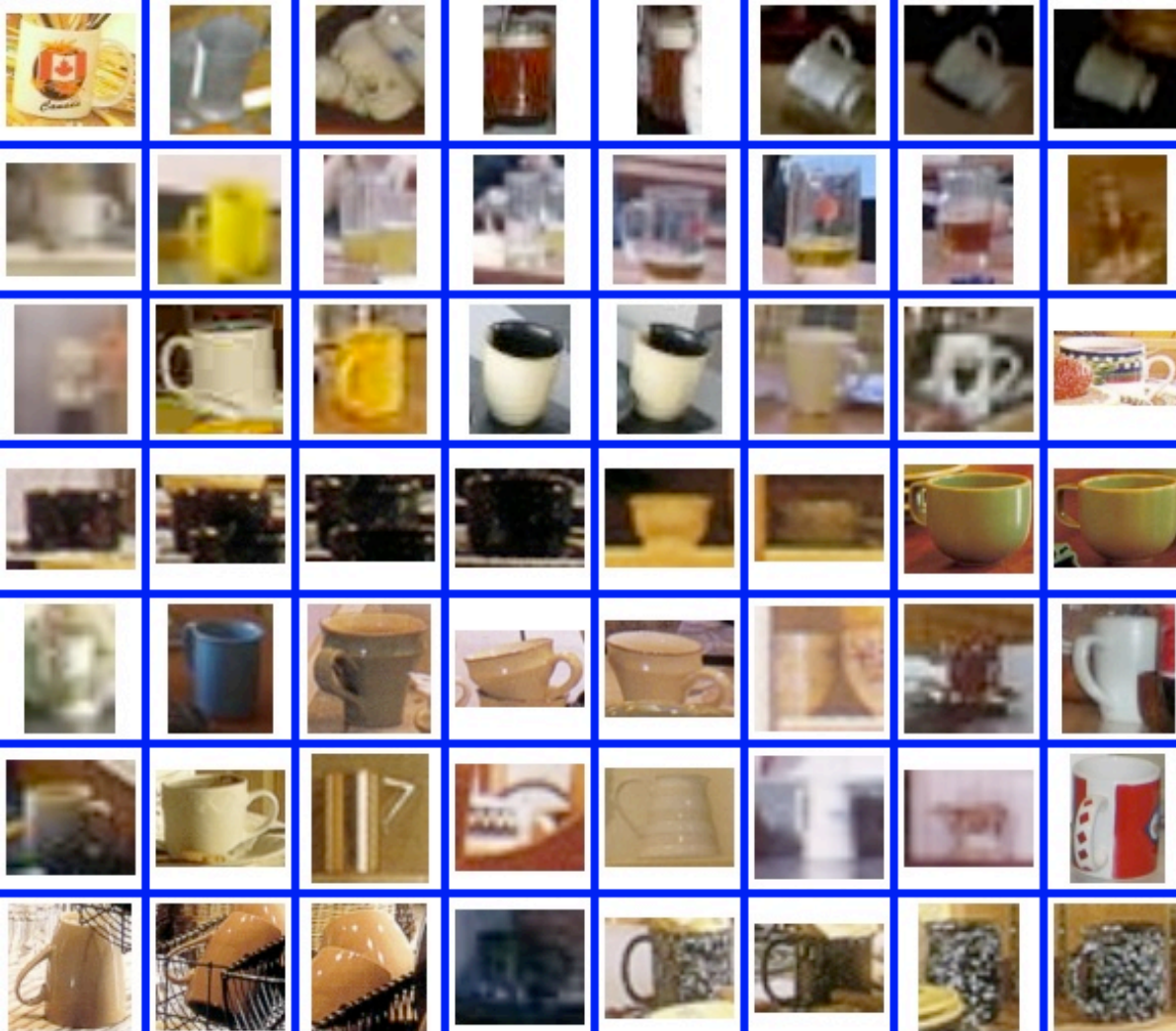
**Coffee Mug as a**  
303 × 301 - 10k - jpg  
[dustbowl.wordpress.com](http://dustbowl.wordpress.com)  
[Find similar images](#)



**SASS Life Member**  
300 × 302 - 6k - jpg  
[sassnet.com](http://sassnet.com)



**personalized coffee**  
400 × 343 - 15k - jp  
[walyou.com](http://walyou.com)  
[Find similar images](#)



Mugs from LabelMe

# Dataset biases

