



CS598:VISUAL INFORMATION RETRIEVAL

Lecture IV: Image Representation:

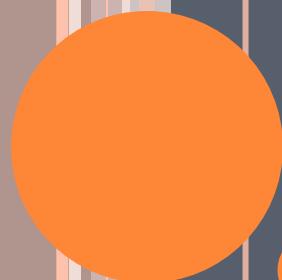
- Semantic and Attribute based Representation



RECAPT OF LECTURE IV

- Histogram of local features
- Bag of words model
- Soft quantization and sparse coding
- Supervector with Gaussian mixture model



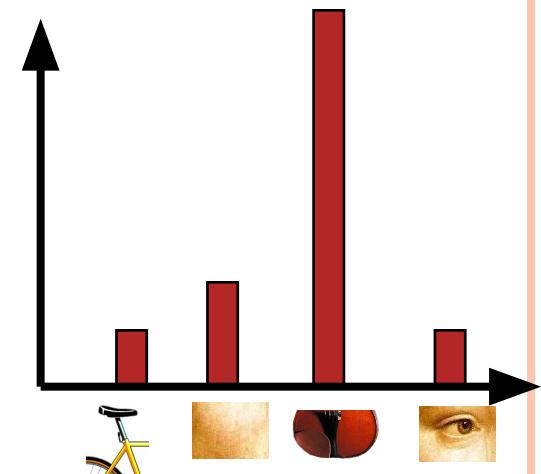
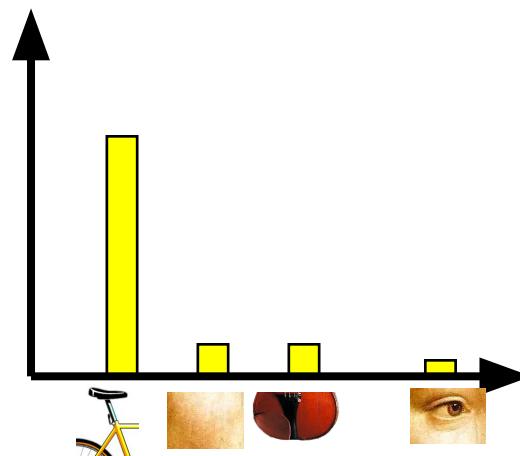
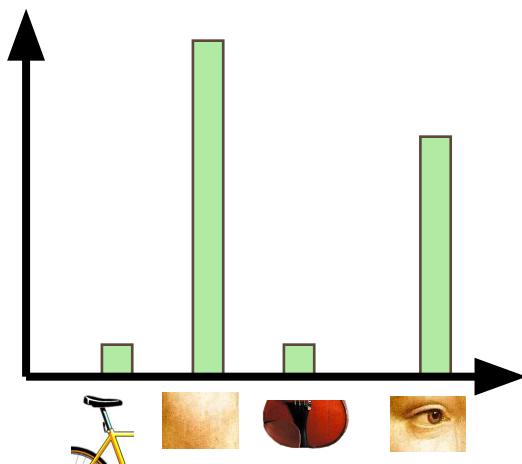


LECTURE V: PART I

Image classification basics

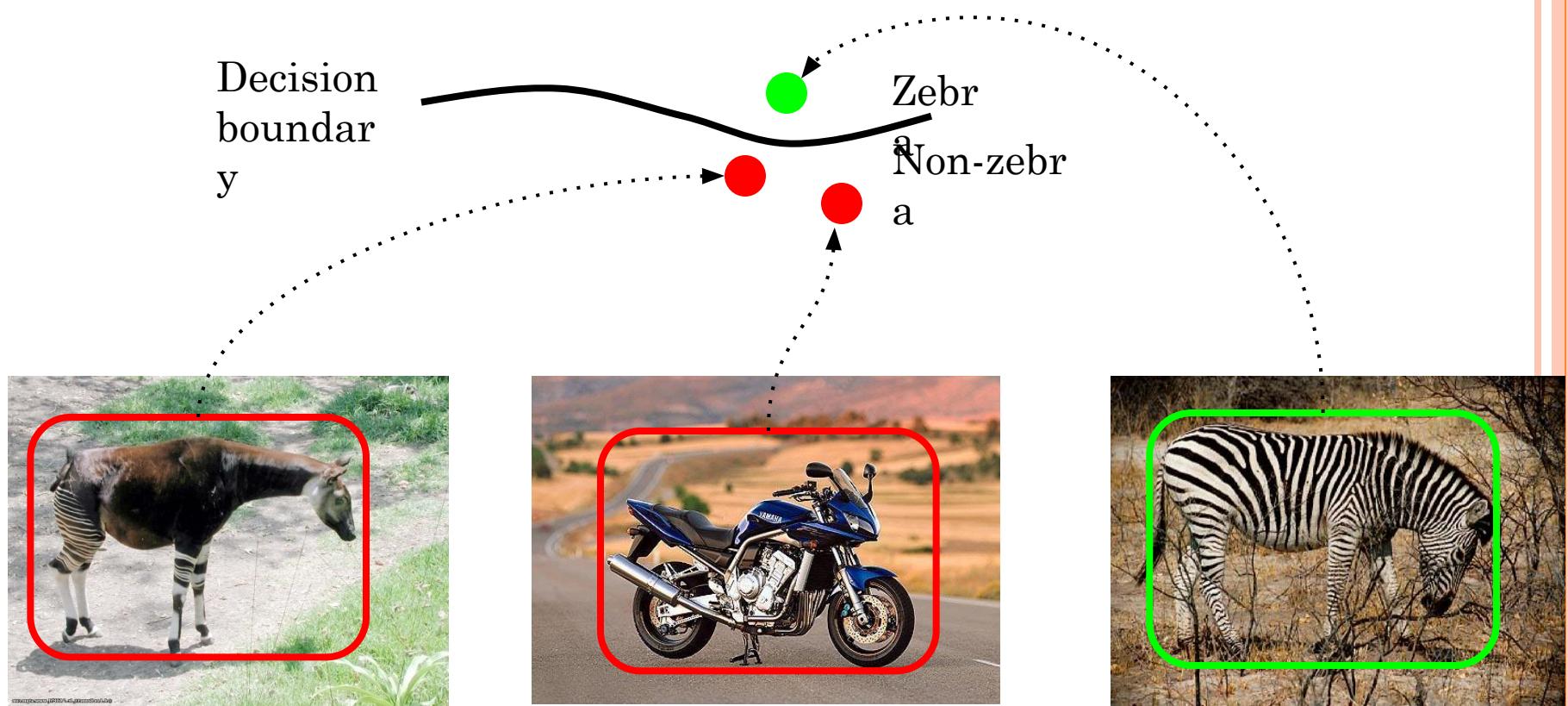
IMAGE CLASSIFICATION

- Given the image representation, e.g., color, texture, shape, local descriptors, of images from different classes, how do we learn a model for distinguishing them?



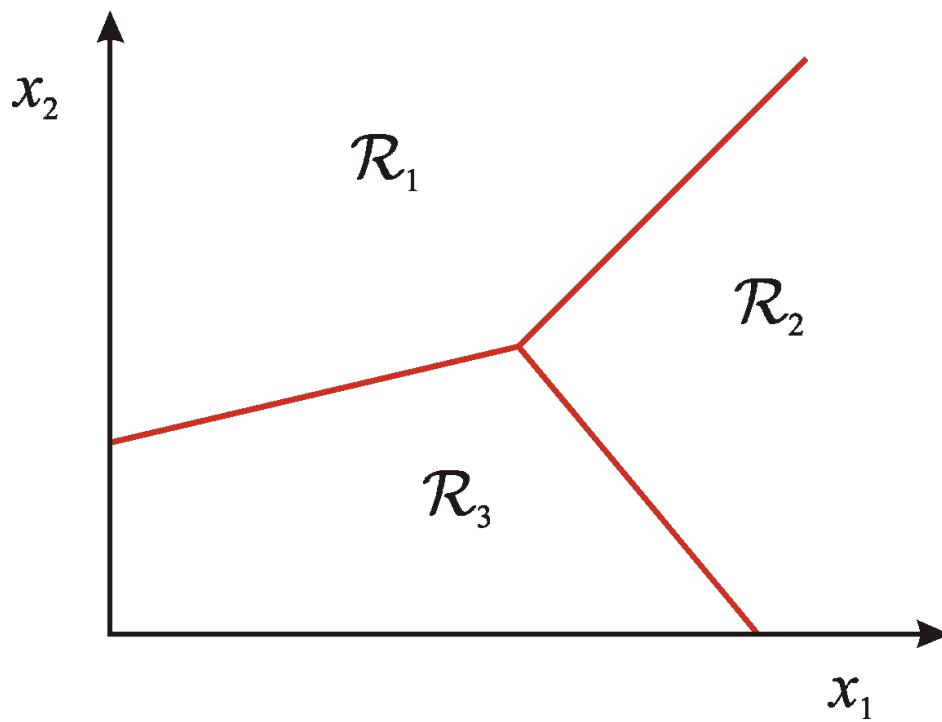
CLASSIFIERS

- Learn a decision rule assigning bag-of-features representations of images to different classes



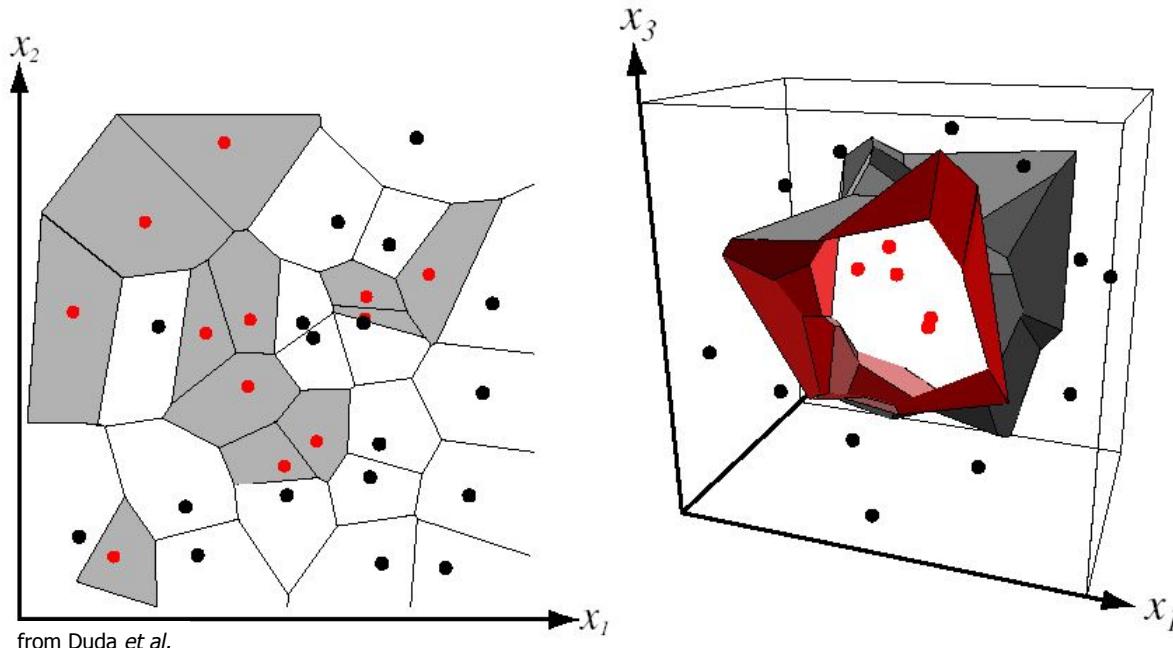
CLASSIFICATION

- Assign input vector to one of two or more classes
- Any decision rule divides input space into *decision regions* separated by *decision boundaries*



NEAREST NEIGHBOR CLASSIFIER

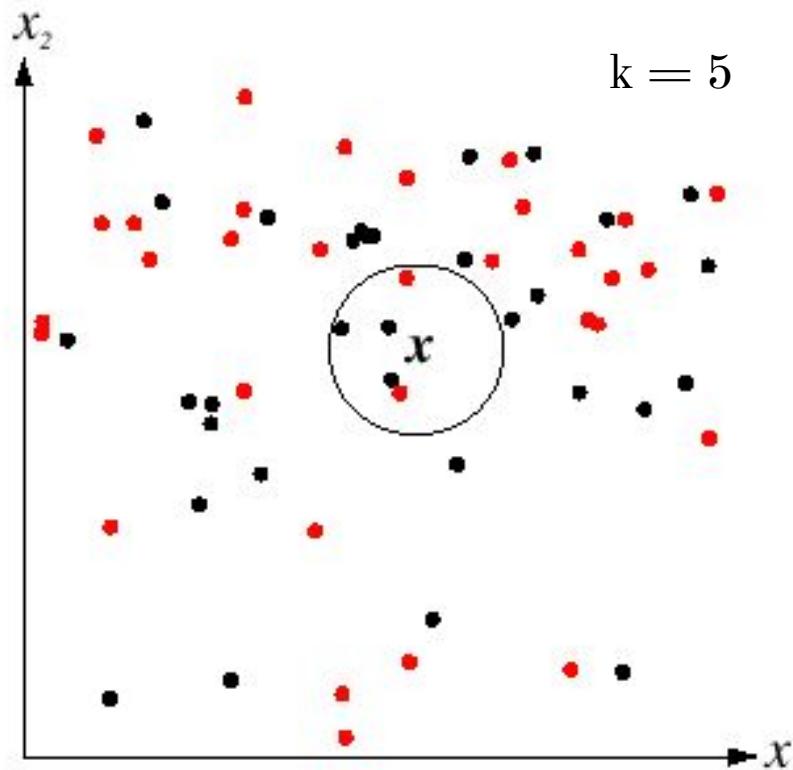
- Assign label of nearest training data point to each test data point



Voronoi partitioning of feature space
for two-category 2D and 3D data

K-Nearest Neighbors

- For a new point, find the k closest points from training data
- Labels of the k points “vote” to classify
- Works well provided there is lots of data and the distance function is good



FUNCTIONS FOR COMPARING HISTOGRAMS

- L1 distance:

$$D(h_1, h_2) = \sum_{i=1}^N |h_1(i) - h_2(i)|$$

- χ^2 distance:

$$D(h_1, h_2) = \sum_{i=1}^N \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$$

- Quadratic distance (*cross-bin distance*):

$$D(h_1, h_2) = \sum_{i,j} A_{ij} (h_1(i) - h_2(j))^2$$

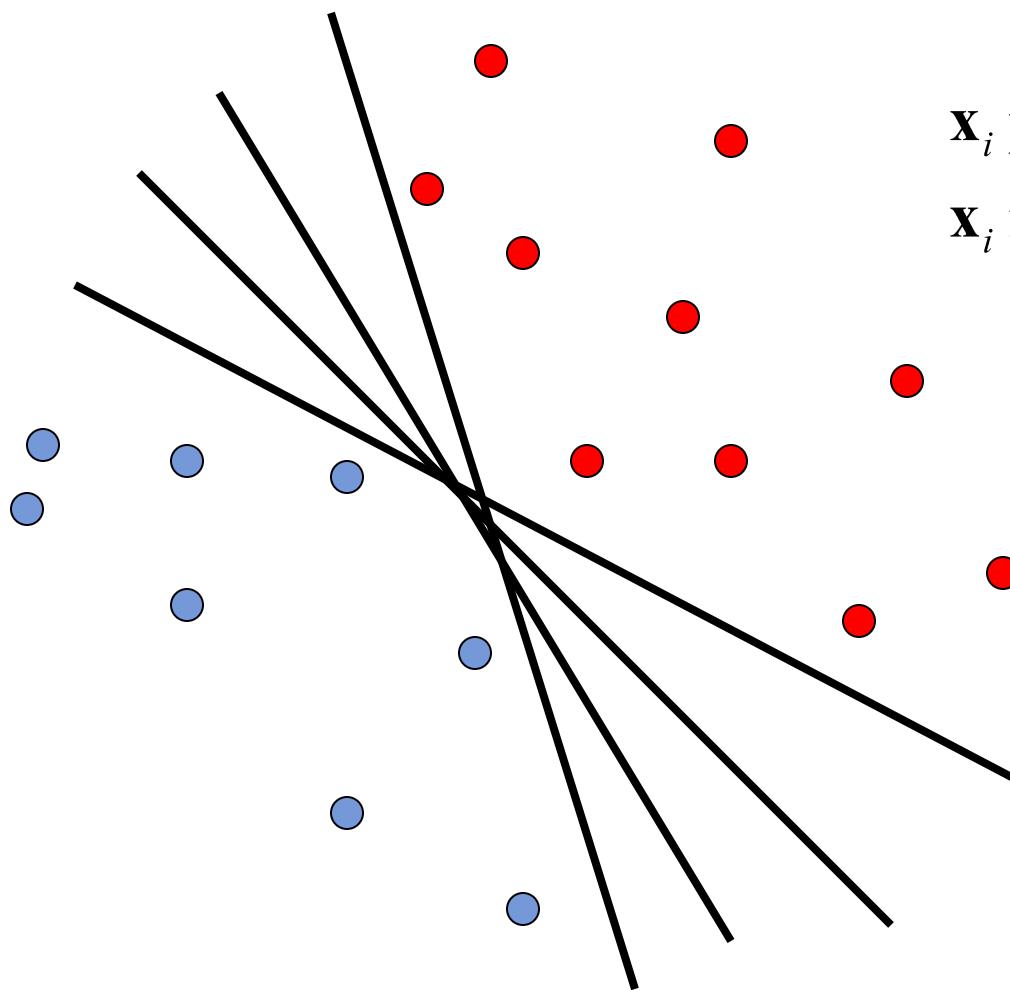
- Histogram intersection (similarity function):

$$I(h_1, h_2) = \sum_{i=1}^N \min(h_1(i), h_2(i))$$



LINEAR CLASSIFIERS

- Find linear function (*hyperplane*) to separate positive and negative examples



\mathbf{x}_i positive : $\mathbf{x}_i \cdot \mathbf{w} + b \geq 0$
 \mathbf{x}_i negative : $\mathbf{x}_i \cdot \mathbf{w} + b < 0$

Which hyperplane
is best?



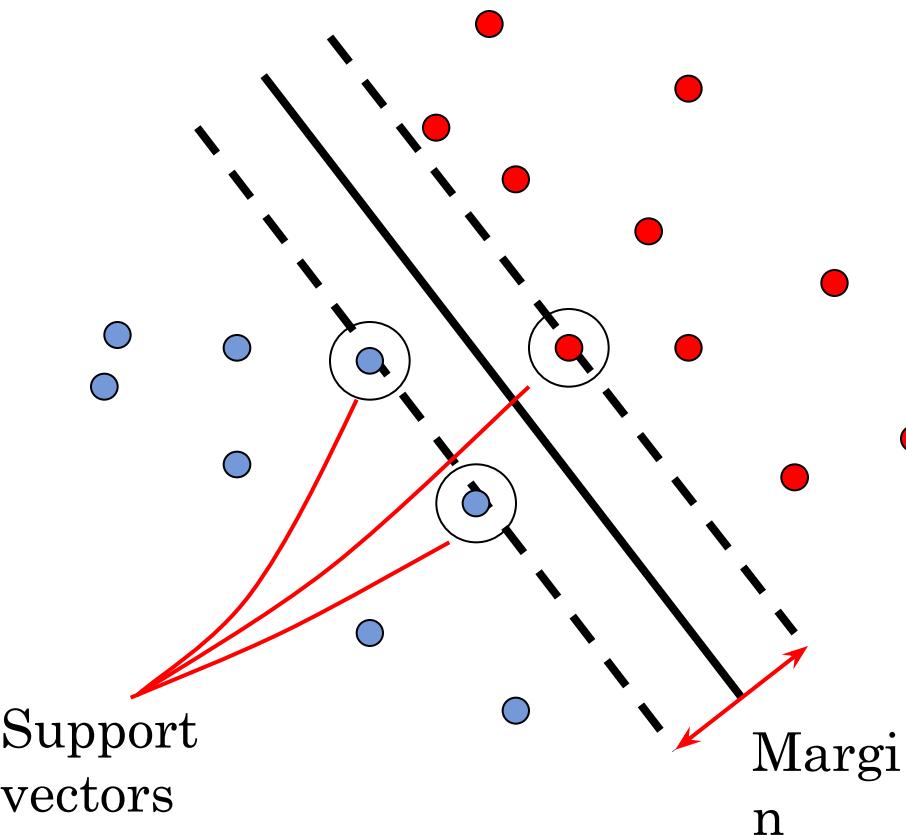
SUPPORT VECTOR MACHINES

- Find hyperplane that maximizes the *margin* between the positive and negative examples



SUPPORT VECTOR MACHINES

- Find hyperplane that maximizes the *margin* between the positive and negative examples



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

For support vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and hyperplane:

Therefore, the margin is $2 / \|\mathbf{w}\|$



FINDING THE MAXIMUM MARGIN HYPERPLANE

1. Maximize margin $2/\|\mathbf{w}\|$
2. Correctly classify all training data:

\mathbf{x}_i positive ($y_i = 1$): $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$

\mathbf{x}_i negative ($y_i = -1$): $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

□ *Quadratic optimization problem:*

□ Minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

Subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$



FINDING THE MAXIMUM MARGIN HYPERPLANE

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$



FINDING THE MAXIMUM MARGIN HYPERPLANE

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$ for any support vector
- Classification function (decision boundary):

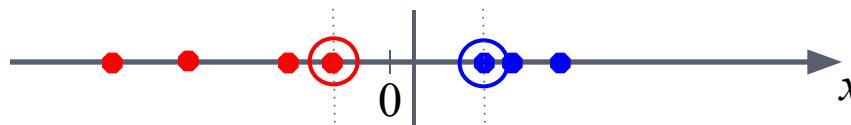
$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
- Solving the optimization problem also involves computing the inner products $\mathbf{x}_i \cdot \mathbf{x}_j$ between all pairs of training points



NONLINEAR SVMs

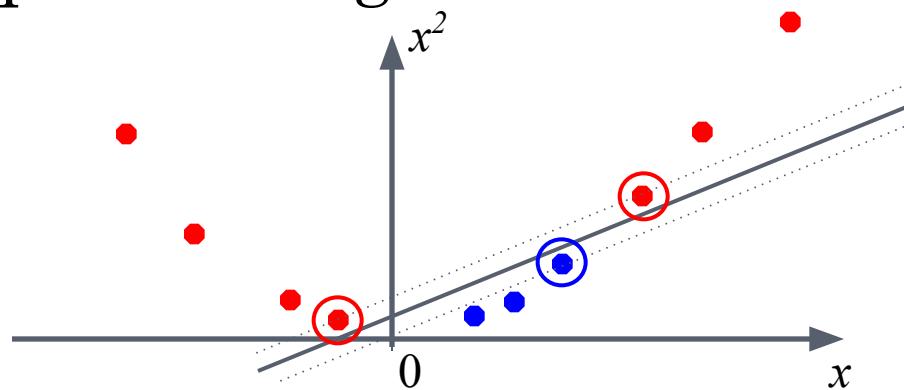
- Datasets that are linearly separable work out great:



- But what if the dataset is just too hard?

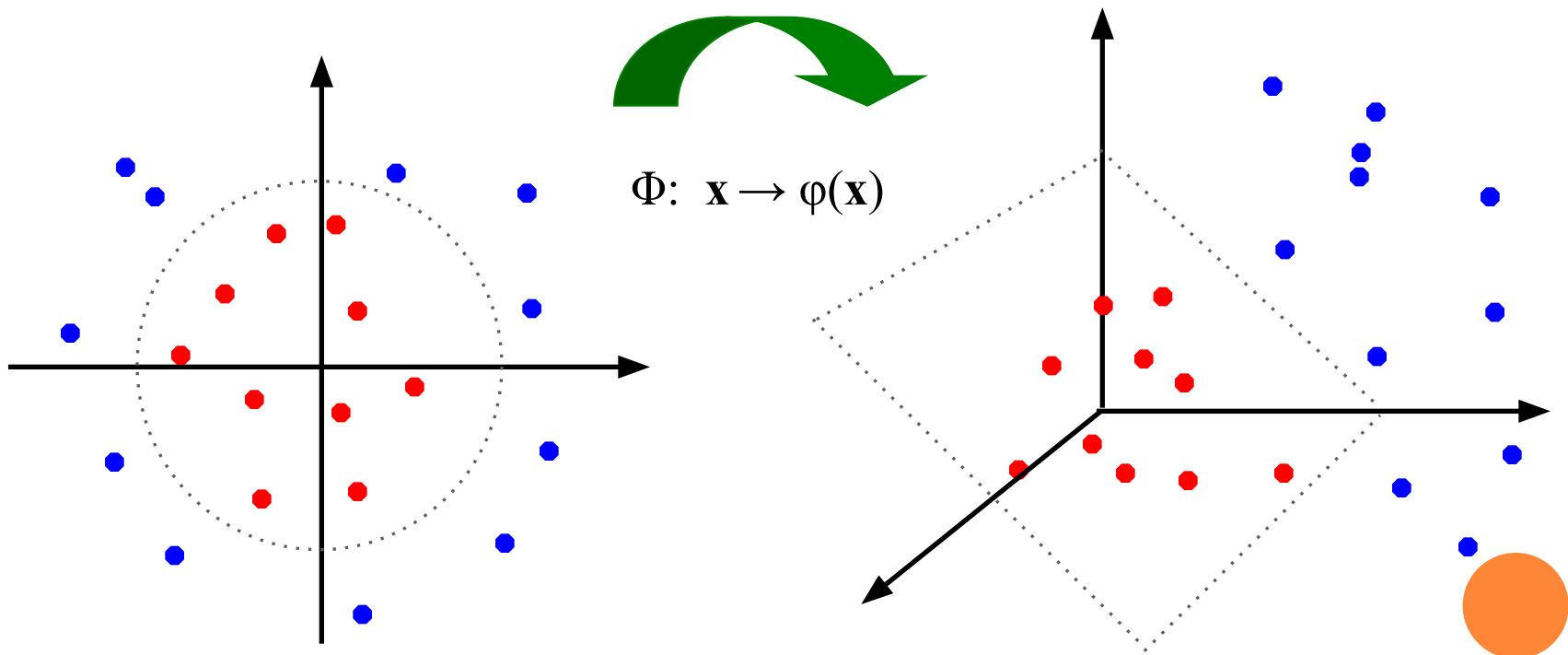


- We can map it to a higher-dimensional space:



NONLINEAR SVMs

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



NONLINEAR SVMs

- *The kernel trick:* instead of explicitly computing the lifting transformation $\varphi(\mathbf{x})$, define a kernel function K such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

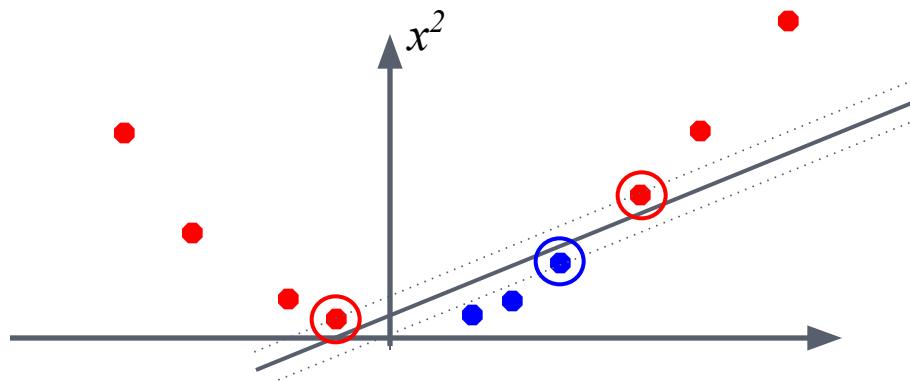
- (to be valid, the kernel function must satisfy *Mercer's condition*)
- This gives a nonlinear decision boundary in the original feature space:

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$



NONLINEAR KERNEL: EXAMPLE

- Consider the mapping $\varphi(x) = (x, x^2)$



$$\varphi(x) \cdot \varphi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2 y^2$$

$$K(x, y) = xy + x^2 y^2$$



KERNELS FOR BAGS OF FEATURES

- Histogram intersection kernel:

$$I(h_1, h_2) = \sum_{i=1}^N \min(h_1(i), h_2(i))$$

- Generalized Gaussian kernel:

$$K(h_1, h_2) = \exp\left(-\frac{1}{A} D(h_1, h_2)^2\right)$$

- D can be L1 distance, Euclidean distance, χ^2 distance, etc.

SUMMARY: SVMs FOR IMAGE CLASSIFICATION

1. Pick an image representation (in our case, bag of features)
2. Pick a kernel function for that representation
3. Compute the matrix of kernel values between every pair of training examples
4. Feed the kernel matrix into your favorite SVM solver to obtain support vectors and weights
5. At test time: compute kernel values for your test example and each support vector, and combine them with the learned weights to get the value of the decision function



WHAT ABOUT MULTI-CLASS SVMs?

- Unfortunately, there is no “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two-class SVMs
- One vs. others
 - Training: learn an SVM for each class vs. the others
 - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- One vs. one
 - Training: learn an SVM for each pair of classes
 - Testing: each learned SVM “votes” for a class to assign to the test example



SVMs: PROS AND CONS

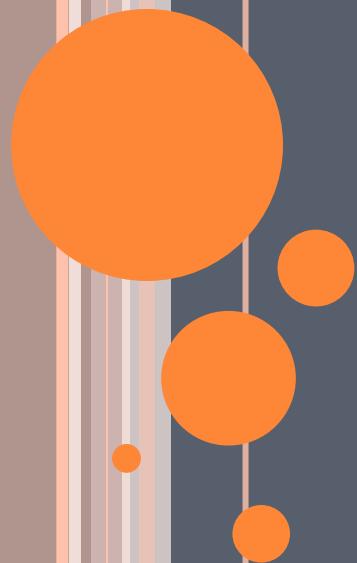
- Pros
 - Many publicly available SVM packages:
<http://www.kernel-machines.org/software>
 - Kernel-based framework is very powerful, flexible
 - SVMs work very well in practice, even with very small training sample sizes
- Cons
 - No “direct” multi-class SVM, must combine two-class SVMs
 - Computation, memory
 - During training time, must compute matrix of kernel values for every pair of examples
 - Learning can take a very long time for large-scale problems



SUMMARY: CLASSIFIERS

- Nearest-neighbor and k-nearest-neighbor classifiers
 - L1 distance, χ^2 distance, quadratic distance, histogram intersection
- Support vector machines
 - Linear classifiers
 - Margin maximization
 - The kernel trick
 - Kernel functions: histogram intersection, generalized Gaussian, pyramid match
 - Multi-class
- Of course, there are many other classifiers out there
 - Neural networks, boosting, decision trees, ...



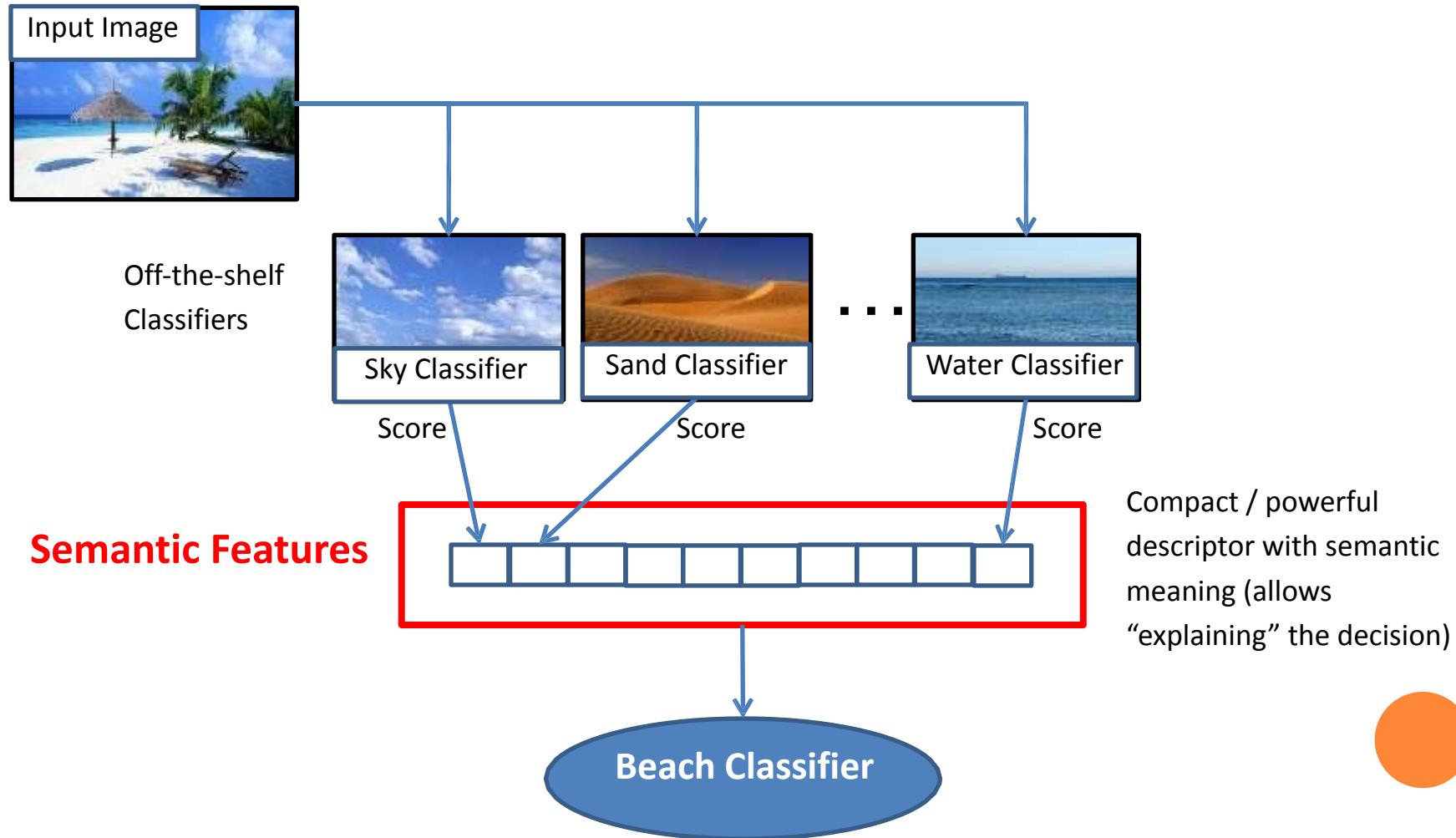


LECTURE IV: PART II

Semantic and Attribute Features

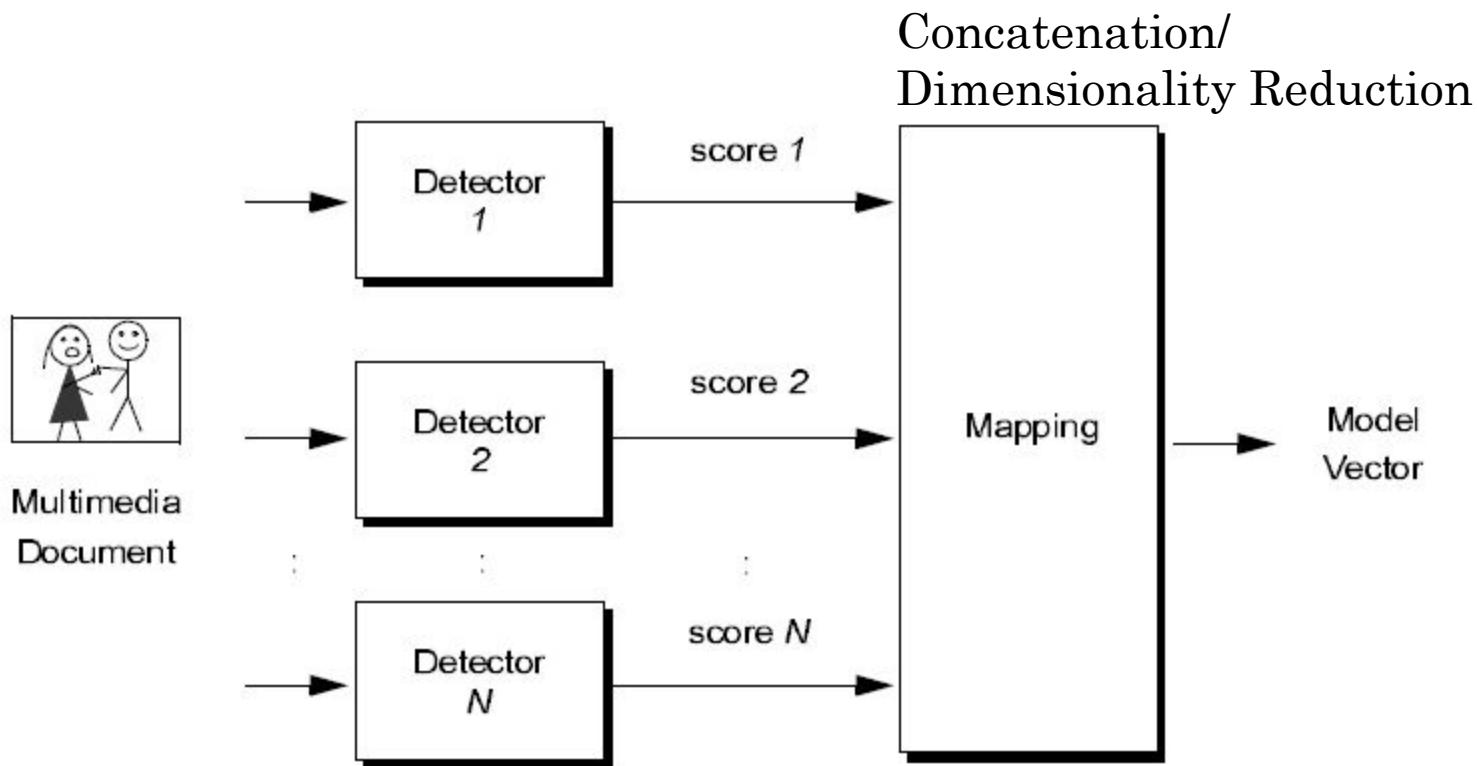
SEMANTIC FEATURES

- Use the scores of semantic classifiers as high-level features



FRAME LEVEL SEMANTIC FEATURES (1)

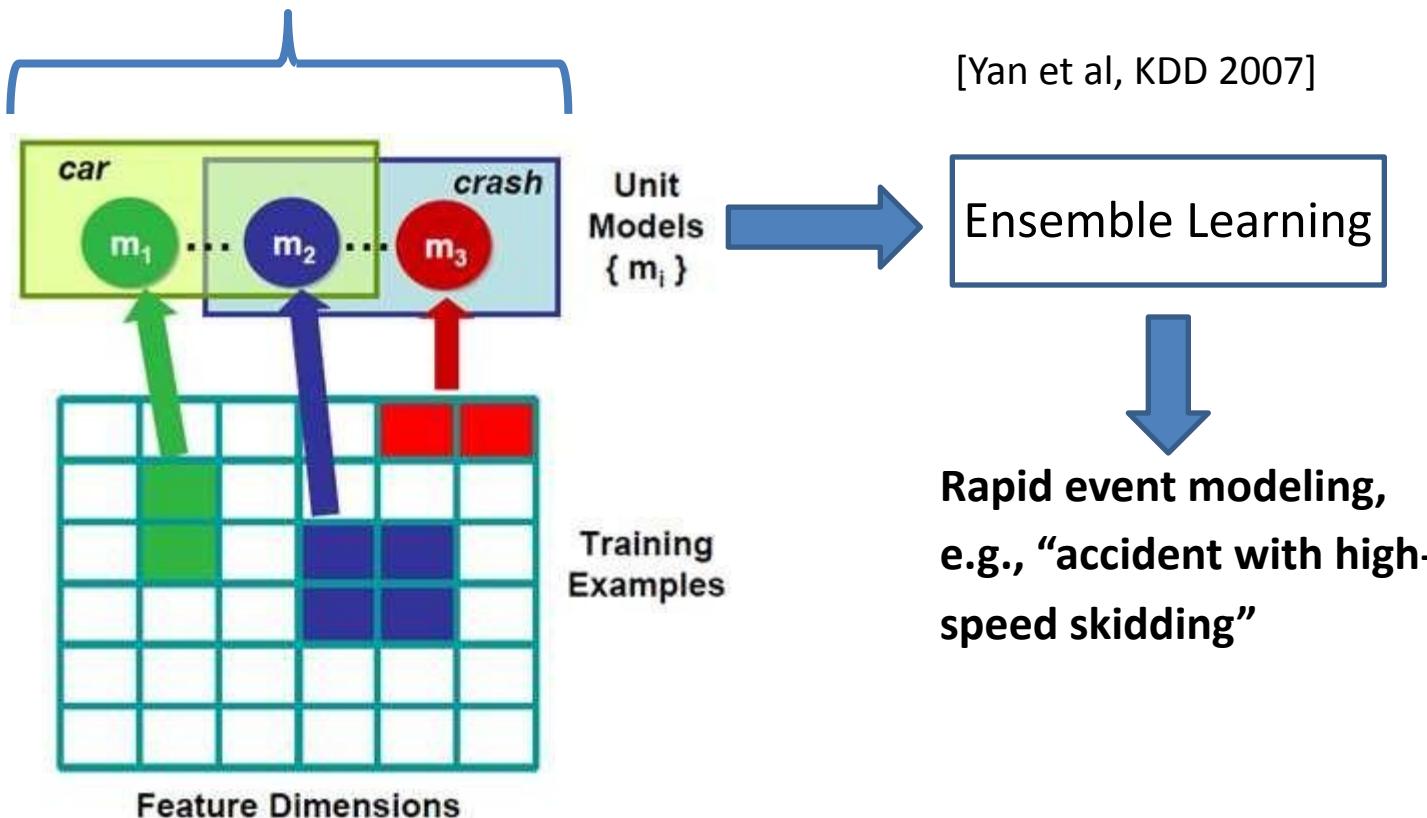
- Early IBM work from multimedia community
 - [Smith et al., ICME2003]



FRAME LEVEL SEMANTIC FEATURES (2)

- IMARS: IBM Multimedia Analysis and Retrieval System

Discriminative semantic basis

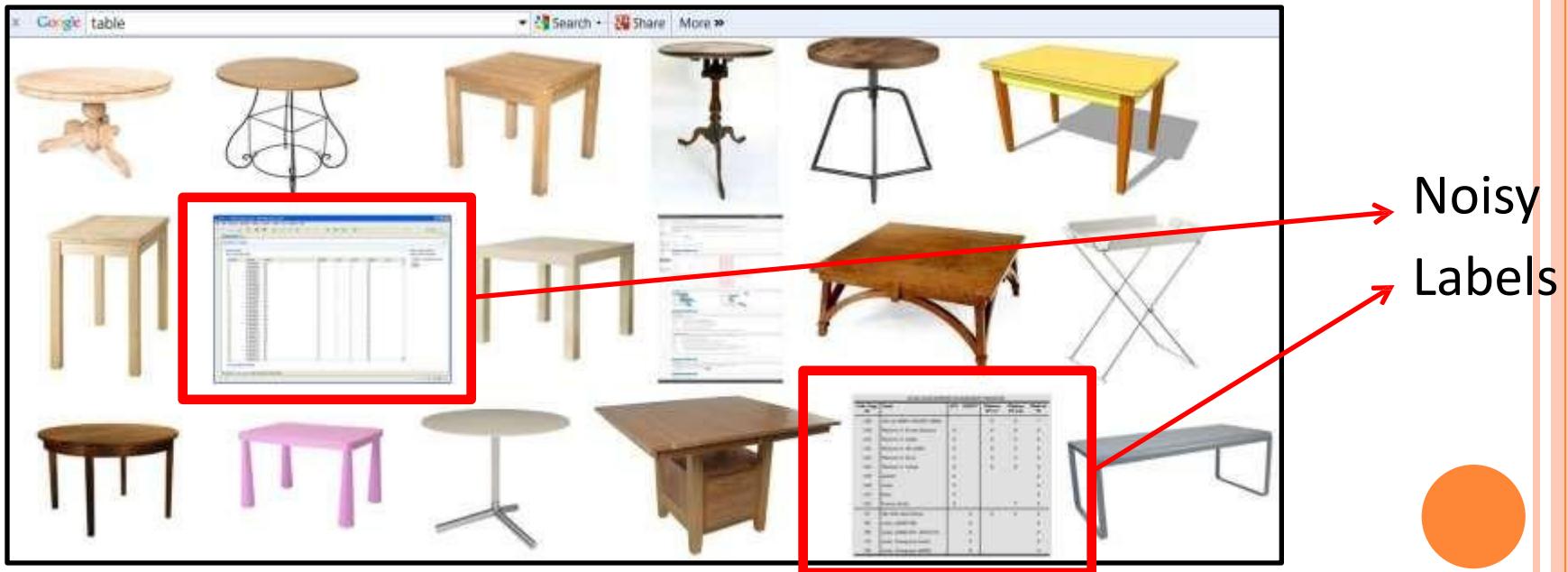


FRAME LEVEL SEMANTIC FEATURES (3)

- CLASSME: descriptor is formed by concatenating the outputs of weakly trained classifiers with noisy labels

[Torresani et al, ECCV2010]

Images used to train the “table” classeme (from Google image search)



FRAME LEVEL SEMANTIC FEATURES (3)

□ CLASSMES

New category	Highly weighted classemes				
cowboy-hat	helmet	sports_track	cake_pan	collectible	muffin_pan
duck	bomber_plane	body_of_water	swimmer	walking	straight
elk	figure_skater	bull_male_herd_animal	cattle	gravesite	dead_body
frisbee	watercraft_surface	scsi_cable	alarm_clock	hindu	serving_tray
trilobite-101	convex_thing	mined_area	cdplayer	roasting_pan	western_hemisphere_person
wheelbarrow	taking_care_of_something	baggage_porter	canopy_closure_open	rowing_shell	container_pressure_barrier



Compact and Efficient Descriptor, useful for large-scale classification

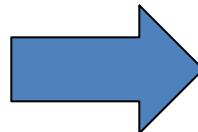


Features are not really semantic!

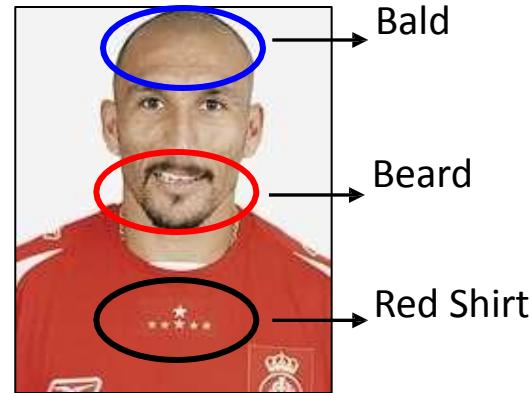


SEMANTIC ATTRIBUTES

Naming



Describing



- Modifiers rather than (or in addition to) nouns
- Semantic properties that are shared among objects
- Attributes are category independent and trasferrable



PEOPLE SEARCH IN SURVEILLANCE VIDEOS

- Traditional approaches: face recognition (naming)
 - Confronted by lighting, pose, and low SNR in surveillance
- Attribute based people search (describing)
 - Semantic attributes based search
 - Example:
 - Show me all bald people at the 42nd street station last month with dark skin, wearing sunglasses, wearing a red jacket

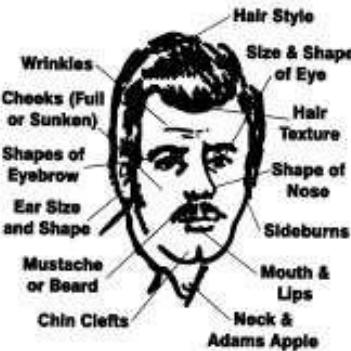


PEOPLE SEARCH IN SURVEILLANCE VIDEOS

PENNSYLVANIA CAPITOL POLICE SUSPECT DESCRIPTION

SEX	RACE	AGE	HEIGHT	WEIGHT	TYPE OF WEAPON	
HAIR/FACIAL HAIR						HAT (color, type)
GLASSES (type)						TIE
TATTOOS						COAT
COMPLEXION						SHIRT
SCARS/MARKS						PANTS/SHOES
HARRISBURG EMERGENCY DIAL 1-911						
POLICE	FIRE	MEDICAL	DON'T HANG UP			
NON-EMERGENCY 717-787-3199						
AUTO MAKE MODEL, COLOR		LICENSE NUMBER	DIRECTION OF ESCAPE	TIME OF DEPARTURE		

FACIAL APPEARANCE

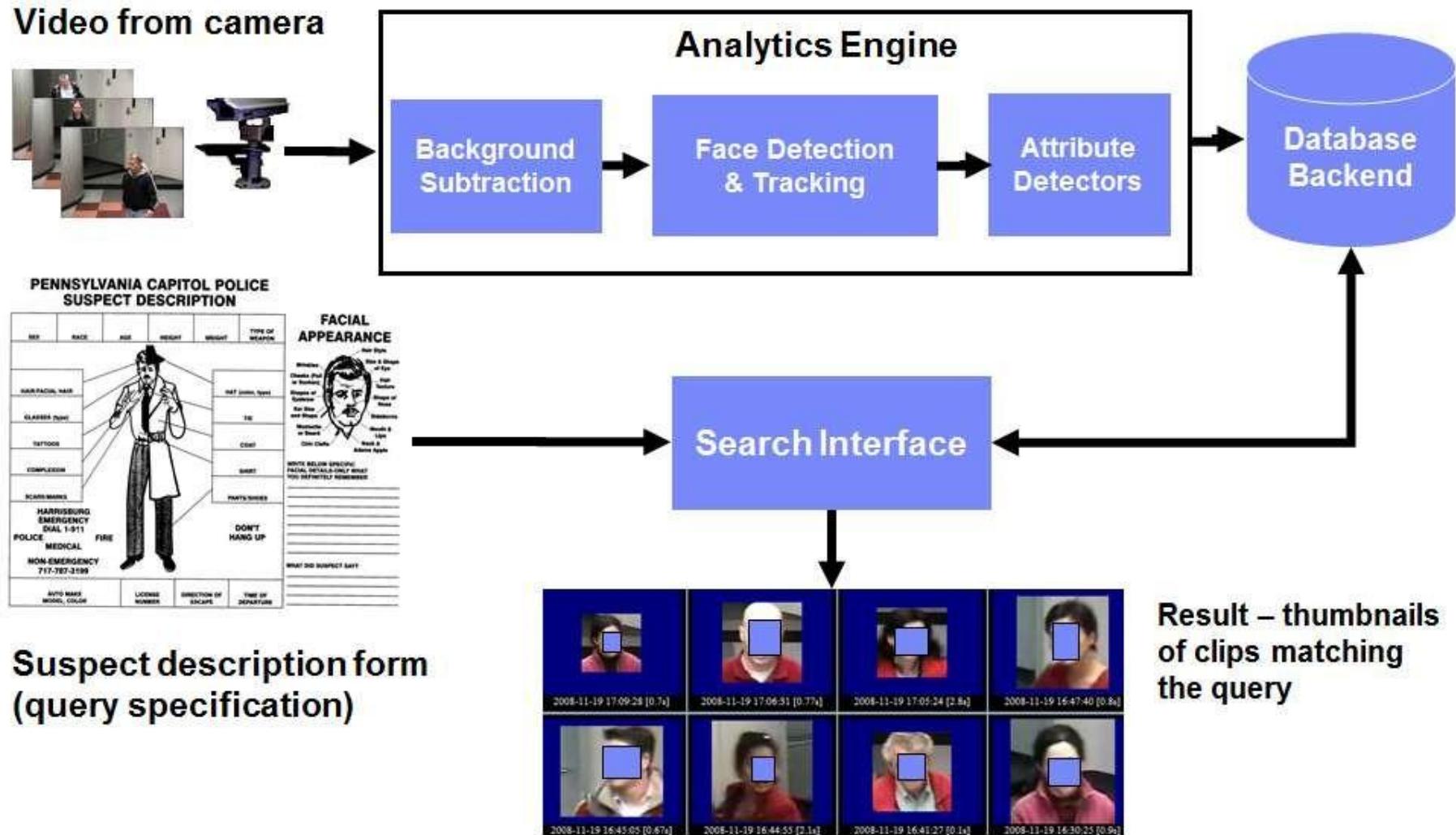


WRITE BELOW SPECIFIC
FACIAL DETAILS-ONLY WHAT
YOU DEFINITELY REMEMBER

WHAT DID SUSPECT SAY?



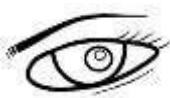
PEOPLE SEARCH IN SURVEILLANCE VIDEOS



PEOPLE SEARCH IN SURVEILLANCE VIDEOS



People Search based on textual descriptions - It does not require training images for the target suspect.



Robustness: attribute detectors are trained using lots of training images covering different lighting conditions, pose variation, etc.



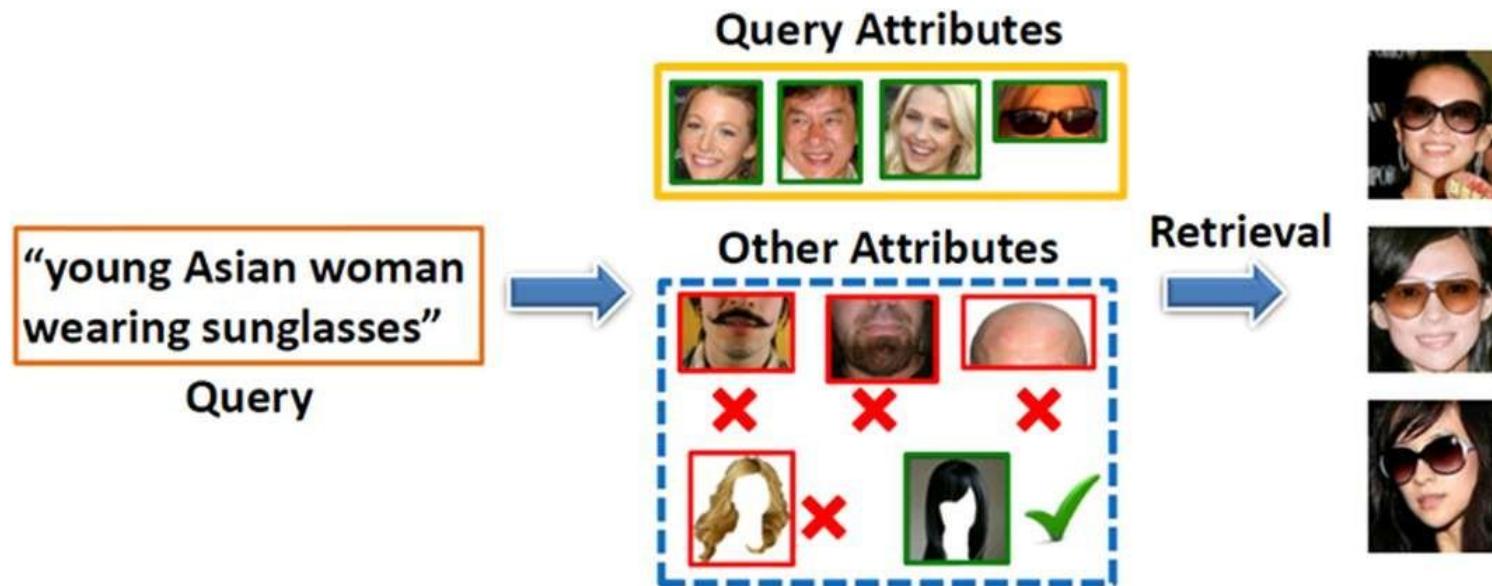
Works well in low-resolution imagery (typical in video surveillance scenarios)



PEOPLE SEARCH IN SURVEILLANCE VIDEOS

Modeling attribute correlations

[Siddiquie et al., “Image Ranking and Retrieval Based on Multi-Attribute Queries”, CVPR 2011]



Attribute-Based Classification



ATTRIBUTE-BASED CLASSIFICATION

Recognition of Unseen Classes (Zero-Shot Learning)

[Lampert et al., Learning To Detect Unseen Object Classes by Between-Class Attribute Transfer, CVPR 2009]

otter

black:	yes
white:	no
brown:	yes
stripes:	no
water:	yes
eats fish:	yes



(1). Train semantic attribute classifiers

polar bear

black:	no
white:	yes
brown:	no
stripes:	no
water:	yes
eats fish:	yes



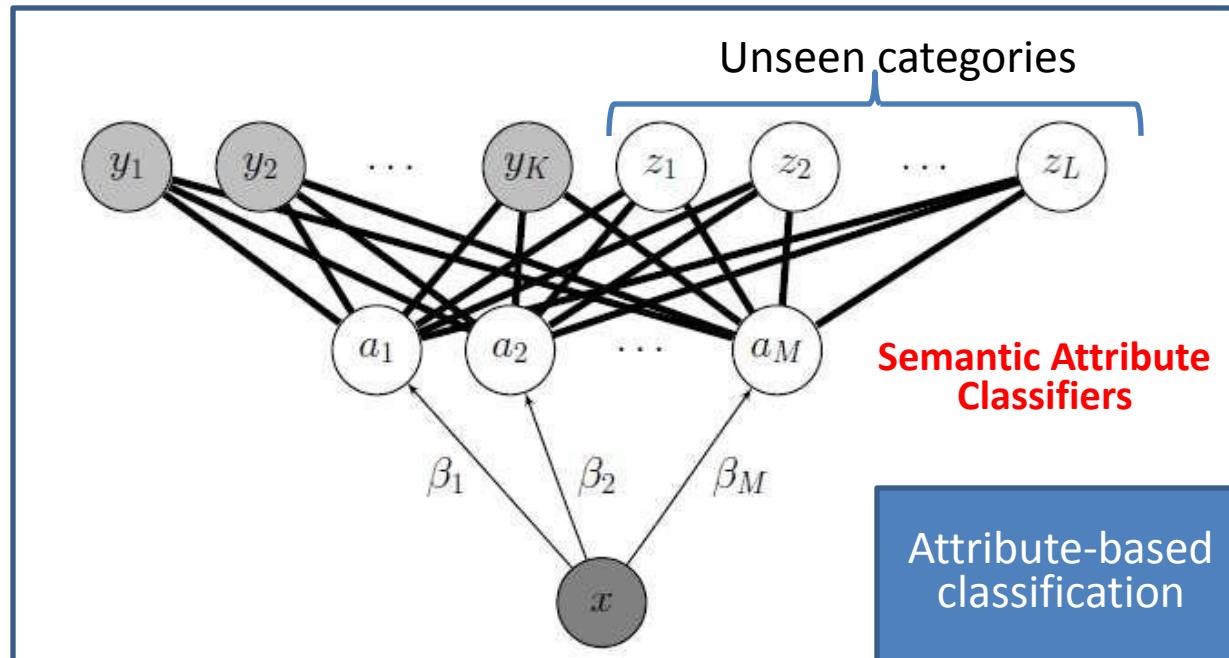
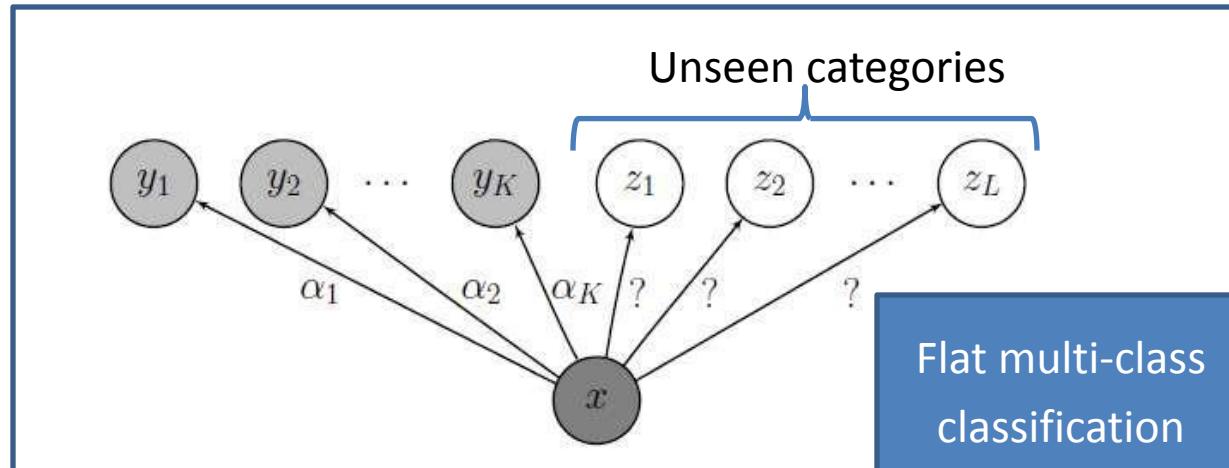
(2). Obtain a classifier for an unseen object (no training samples) by just specifying which attributes it has

zebra

black:	yes
white:	yes
brown:	no
stripes:	yes
water:	no
eats fish:	no

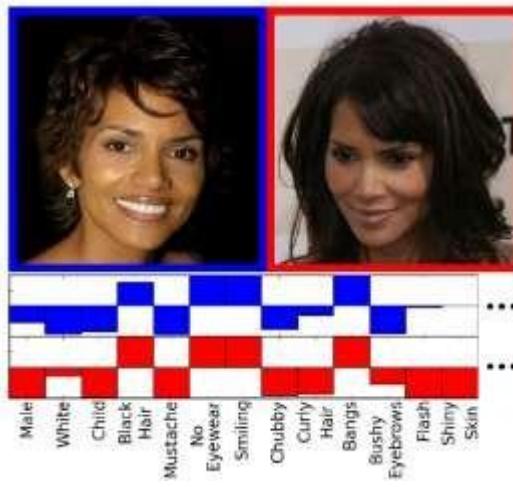


ATTRIBUTE-BASED CLASSIFICATION



ATTRIBUTE-BASED CLASSIFICATION

Face verification [Kumar et al, ICCV 2009]



Action recognition [Liu et al, CVPR2011]

Two action recognition examples with their descriptions and feature tables.

Naming: Walking

Description	Indoor related	Outdoor related	Translation motion	Arm pendulum-like motion	Torso up-down motion	Torso twist	Having stick-like tool
	Yes	Yes	Yes	Yes	No	No	No

Naming: Golf-Swinging

Description	Indoor related	Outdoor related	Translation motion	Arm pendulum-like motion	Torso up-down motion	Torso twist	Having stick-like tool
	No	Yes	No	No	No	Yes	Yes

Animal Recognition

[Lampert et al, CVPR 2009]



Person Re-identification
[Layne et al, BMVC 2012]



Bird Categorization [Farrell et al, ICCV 2011]



Many more! Significant growth in the past few years

ATTRIBUTE-BASED CLASSIFICATION



Note: Several recent methods use the term “attributes” to refer to non-semantic model outputs

In this case attributes are just mid-level features, like PCA, hidden layers in neural nets, ... (non-interpretable splits)



ATTRIBUTE-BASED CLASSIFICATION

<http://rogerioferis.com/VisualRecognitionAndSearch/Resources.html>

Datasets

Attributes

- [Animals with Attributes](#) – 30,475 images of 50 animals classes with 6 pre-extracted feature representations for each image.
- [aYahoo and aPascal](#) – Attribute annotations for images collected from Yahoo and Pascal VOC 2008.
- [FaceTracer](#) – 15,000 faces annotated with 10 attributes and fiducial points.
- [PubFig](#) – 58,797 face images of 200 people with 73 attribute classifier outputs.
- [LFW](#) – 13,233 face images of 5,749 people with 73 attribute classifier outputs.
- [Human Attributes](#) – 8,000 people with annotated attributes. Check also this [link](#) for another dataset of human attributes.
- [SUN Attribute Database](#) – Large-scale scene attribute database with a taxonomy of 102 attributes.
- [ImageNet Attributes](#) – Variety of attribute labels for the ImageNet dataset.
- [Relative attributes](#) – Data for OSR and a subset of PubFig datasets. Check also this [link](#) for the WhittleSearch data.
- [Attribute Discovery Dataset](#) – Images of shopping categories associated with textual descriptions.
- [Caltech-UCSD Birds Dataset](#) – Hundreds of bird categories with annotated parts and attributes.



Attributes for Fine-Grained Categorization

Easy for Humans



Hard for Humans









- Main page
- Contents
- Featured content
- Current events
- Random article

search

interaction

- About Wikipedia
- Community portal
- Recent changes
- Contact Wikipedia
- Donate to Wikipedia
- Help

toolbox

- What links here
- Related changes
- Upload file
- Special pages
- Printable version

Cere

From Wikipedia, the free encyclopedia

The **cere** (from the Latin *cera*: wax) is a waxy skin-like area at the top of the beaks of certain birds. Hawks, pigeons, and falcons are among the birds that have cerea. The word 'cere' is derived from the Latin word 'cera'. The two are not identical. The cere is located at the top of the beak of certain dimorphic birds, and also has a different texture.

Contents [hide]

- 1 Physical characteristics
- 2 Role in respiration
- 3 Role in indication of reproductive cycle
- 4 References
- 5 See also

Physical characteristics

The cere is located at the top of the beak of certain birds. The colour of the cere may vary between species, and also depends on the sex of the bird. In falcons, the opening of the nostrils (nostrilia). The shape of the cere is often used to identify different species. In falcons, the opening of the nostrils is located at the top of the beak.

FINE-GRAINED CATEGORIZATION

Visipedia (<http://visipedia.org/>)

- Machines collaborating with humans to organize visual knowledge, connecting text to images, images to text, and images to images
- Easy annotation interface for experts (powered by computer vision)

Visual Query: Fine-grained Bird Categorization

The image shows a hand holding a silver flip phone. The phone's screen displays a photograph of a blue jay perched on a branch. To the right of the phone is a screenshot of a Wikipedia page for "Blue Jay". The page header includes the Wikipedia logo and navigation links for "article", "discussion", "edit this page", and "history". The main content area is titled "Blue Jay" and describes it as a passerine bird native to North America. It mentions its adaptability, aggressiveness, and omnivorous nature. Below the main text is a "Contents" section with links to "Description", "Vocalizations", and "Distribution and habitat". To the right of the main content area is a large image of a blue jay perched on a snow-covered branch.

Picture credit: Serge Belongie

FINE-GRAINED CATEGORIZATION

African



Is it an African or Indian Elephant?



Indian



Example-based Fine-Grained Categorization is Hard!!

FINE-GRAINED CATEGORIZATION

African



Larger
Ears

Is it an African or Indian Elephant?



Indian



Smaller Ears



Visual distinction of subordinate categories may be quite subtle,
usually based on **Attributes**

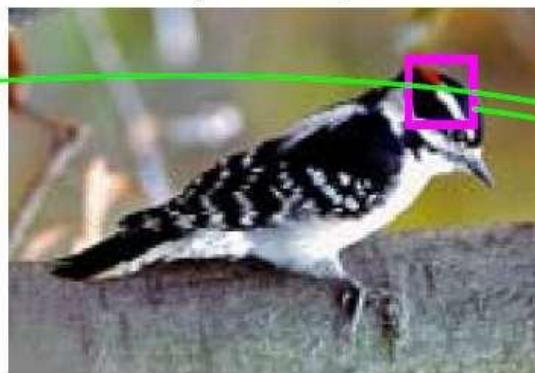
FINE-GRAINED CATEGORIZATION

- Standard classification methods may not be suitable because the variation between classes is small ...

Three toed woodpecker

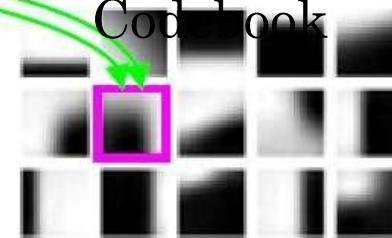


Downy woodpecker



[B. Yao, CVPR
2012]

Cod book

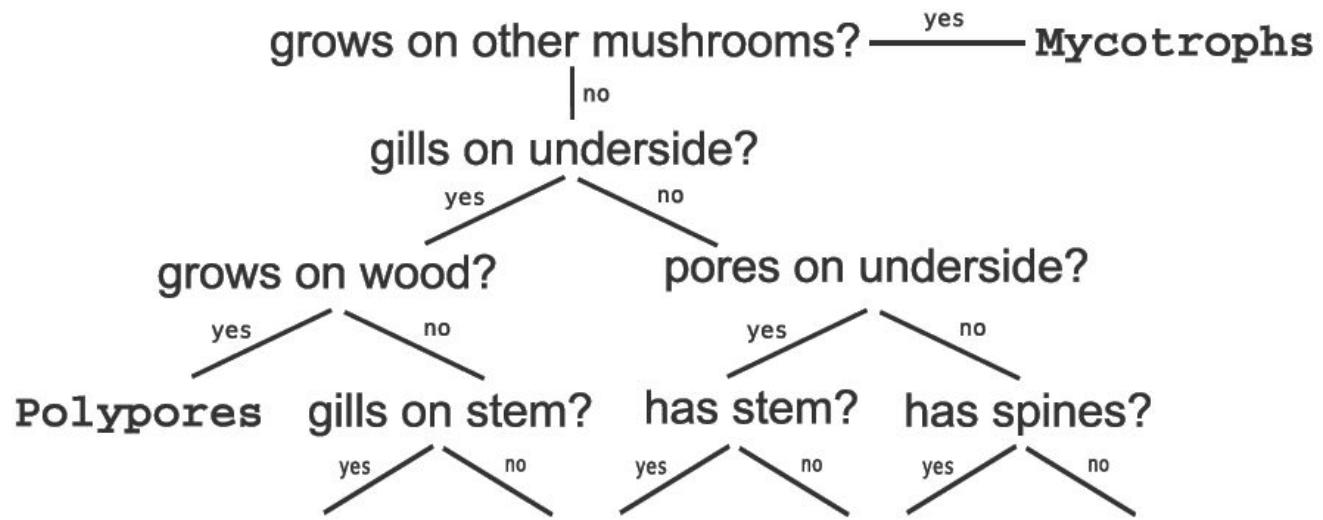
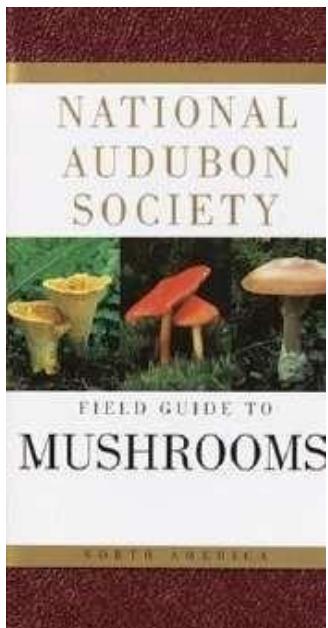


- ... and intra-class variation is still high.



FINE-GRAINED CATEGORIZATION

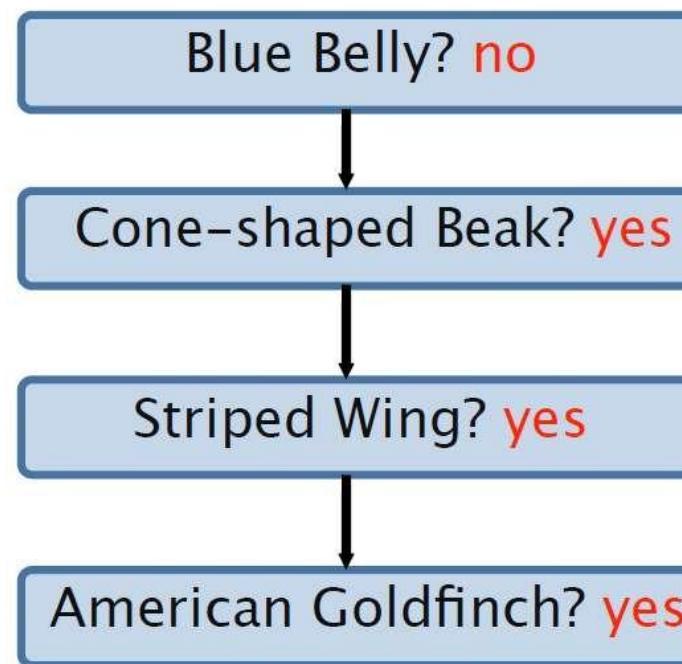
- Humans rely on field guides!
- Field guides usually refer to parts and attributes of the object



FINE-GRAINED CATEGORIZATION

[Branson et al, Visual Recognition with Humans in the Loop, ECCV 2010]

Visual 20 Questions Game



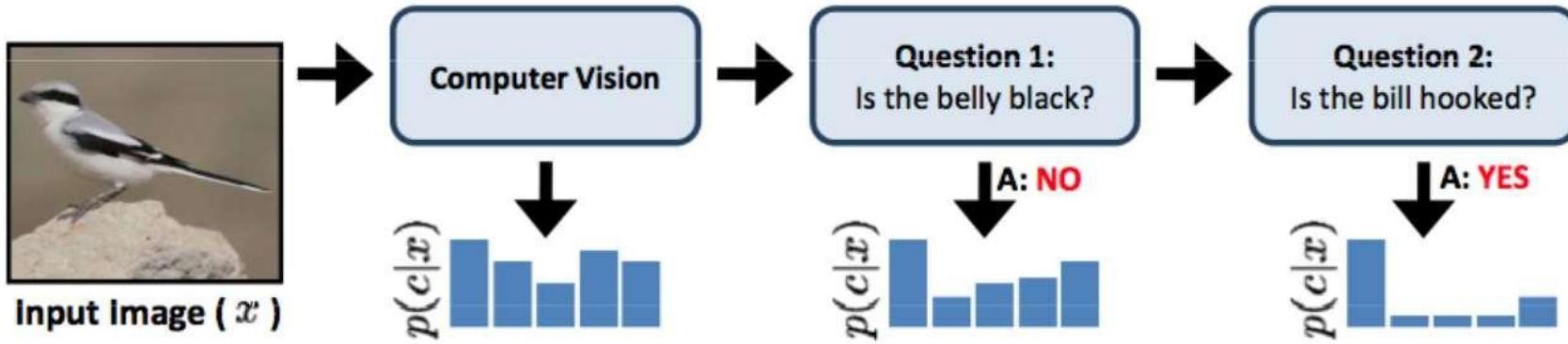
Hard classification problems can be turned into a sequence of easy ones



FINE-GRAINED CATEGORIZATION

[Branson et al, Visual Recognition with Humans in the Loop, ECCV 2010]

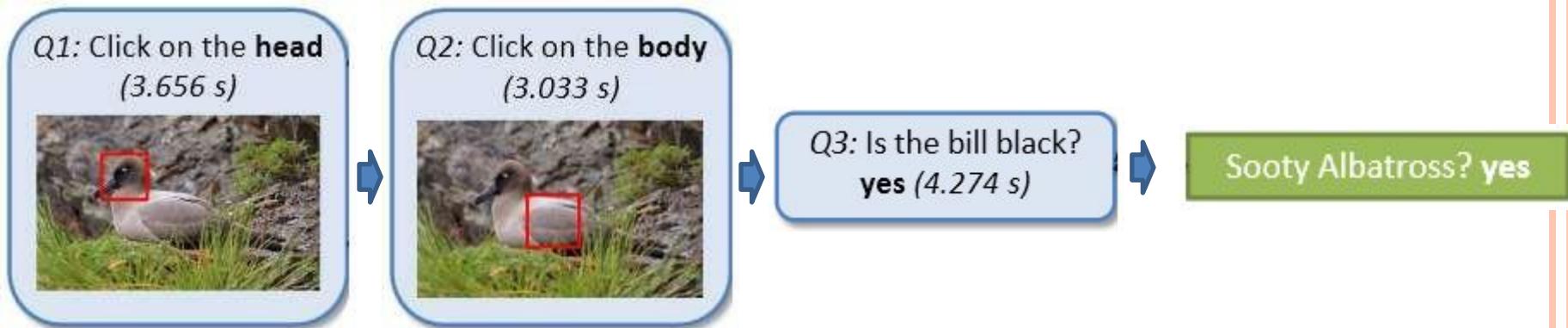
- Computer vision reduces the amount of human-interaction (minimizes the number of questions)



FINE-GRAINED CATEGORIZATION

[Wah et al, Multiclass Recognition and Part Localization with Humans in the Loop, ICCV 2011]

- Localized part and attribute detectors.
- Questions include asking the user to localize parts.



FINE-GRAINED CATEGORIZATION

- ❑ <http://www.vision.caltech.edu/visipedia/CUB-200-2011.html>

Caltech-UCSD Birds-200-2011

Browse



Click here to browse the dataset.

Details

Caltech-UCSD Birds-200-2011 (CUB-200-2011) is an extended version of the CUB-200 dataset, with roughly double the number of images per class and new part location annotations. For detailed information about the dataset, please see the technical report linked below.

- **Number of categories:** 200
- **Number of images:** 11,788
- **Annotations per image:** 15 Part Locations, 312 Binary Attributes, 1 Bounding Box

Some related datasets are Caltech-256, the Oxford Flower Dataset, and Animals with Attributes. More datasets are available at the Caltech Vision Dataset Archive.



FINE-GRAINED CATEGORIZATION

Video Demo:

http://www.youtube.com/watch?v=_ReKVqnDXzA



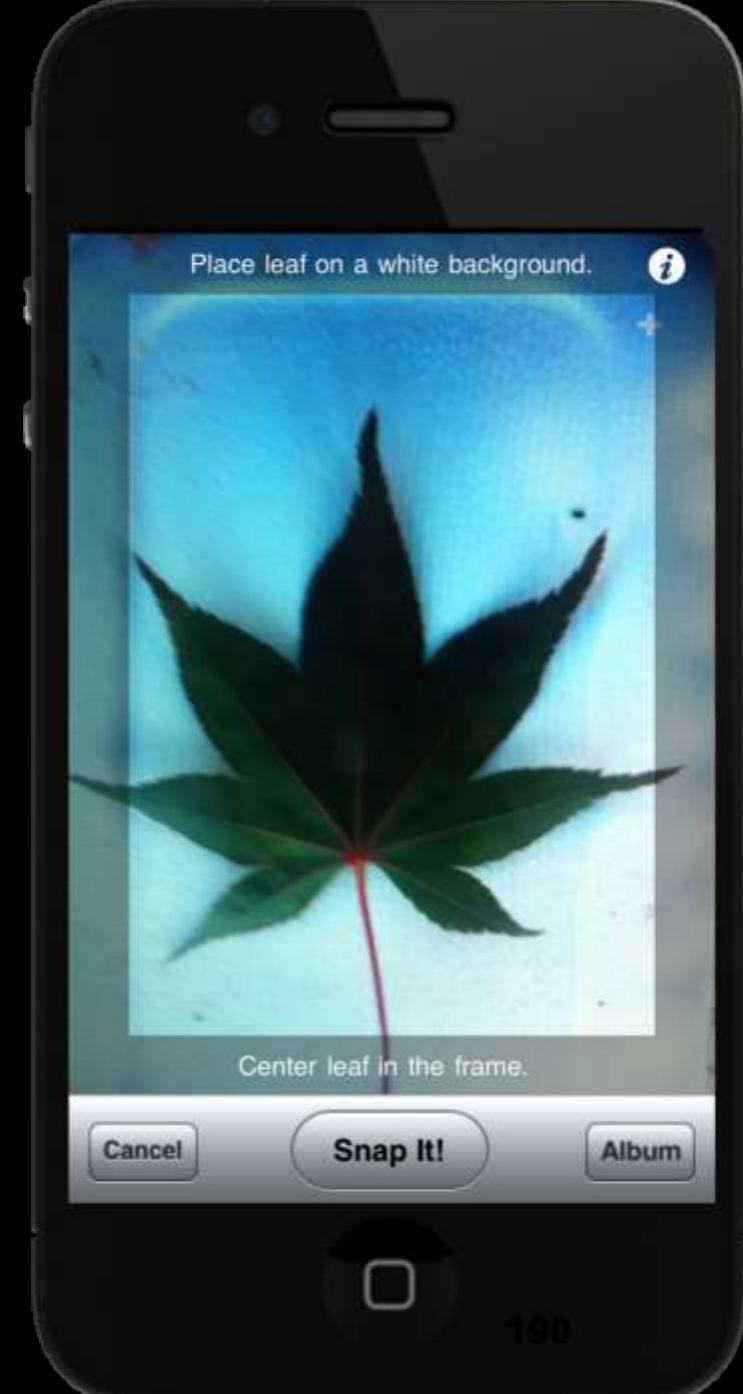


Like a normal field guide...

- that you can search and sort
- and with visual recognition



See N. Kumar et al,
"Leafsnap: A Computer
Vision System for
Automatic Plant Species
Identification, ECCV 2012



leafsnap



Available on the
App Store



Available on the iPad
App Store

- Nearly 1 million downloads
 - 40k new users per month
 - 100k active users
- 1.7 million images taken
 - 100k new images/month
 - 100k users with > 5 images
- Users from all over the world
- Botanists, educators, kids, hobbyists, photographers, ...



FINE-GRAINED CATEGORIZATION



Check the fine-grained visual categorization workshop:
<http://www.fgvc.org/>

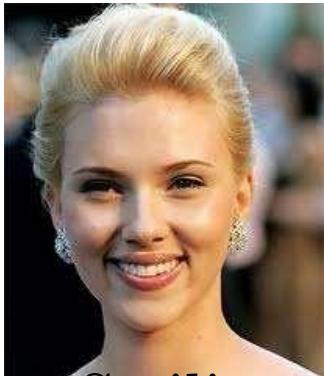


Relative Attributes



RELATIVE ATTRIBUTES

[Parikh & Grauman, Relative Attributes, ICCV 2011]



Smiling



???



No smiling
t



Natural



???



Not natural



LEARNING RELATIVE ATTRIBUTES

For each attribute a_m , e.g., “openness”

Supervision consists of:

$$O_m: \left\{ \left(\begin{matrix} \text{[Image of a narrow street]} \\ \text{[Image of a dense city]} \end{matrix} \right) \succ, \dots \right\},$$

Ordered pairs

$$S_m: \left\{ \left(\begin{matrix} \text{[Image of a beach]} \\ \text{[Image of a field]} \end{matrix} \right) \sim, \dots \right\}$$

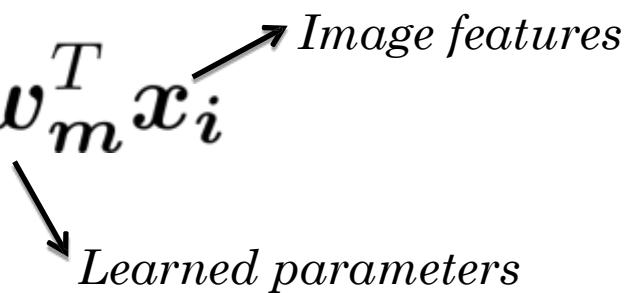
Similar pairs



LEARNING RELATIVE ATTRIBUTES

Learn a ranking function

$$r_m(\mathbf{x}_i) = \mathbf{w}_m^T \mathbf{x}_i$$



Learned parameters

Image features

that best satisfies the constraints:

$$\forall (i, j) \in O_m : \mathbf{w}_m^T \mathbf{x}_i > \mathbf{w}_m^T \mathbf{x}_j$$

$$\forall (i, j) \in S_m : \mathbf{w}_m^T \mathbf{x}_i = \mathbf{w}_m^T \mathbf{x}_j$$



LEARNING RELATIVE ATTRIBUTES

Max-margin learning to rank formulation

$$\begin{aligned} \min \quad & \left(\frac{1}{2} \|\mathbf{w}_m^T\|_2^2 + C \left(\sum \xi_{ij}^2 + \sum \gamma_{ij}^2 \right) \right) \\ \text{s.t.} \quad & \mathbf{w}_m^T (\mathbf{x}_i - \mathbf{x}_j) \geq 1 - \xi_{ij}, \forall (i, j) \in O_m \\ & |\mathbf{w}_m^T (\mathbf{x}_i - \mathbf{x}_j)| \leq \gamma_{ij}, \forall (i, j) \in S_m \\ & \xi_{ij} \geq 0; \gamma_{ij} \geq 0 \end{aligned}$$

Based on [Joachims 2002]

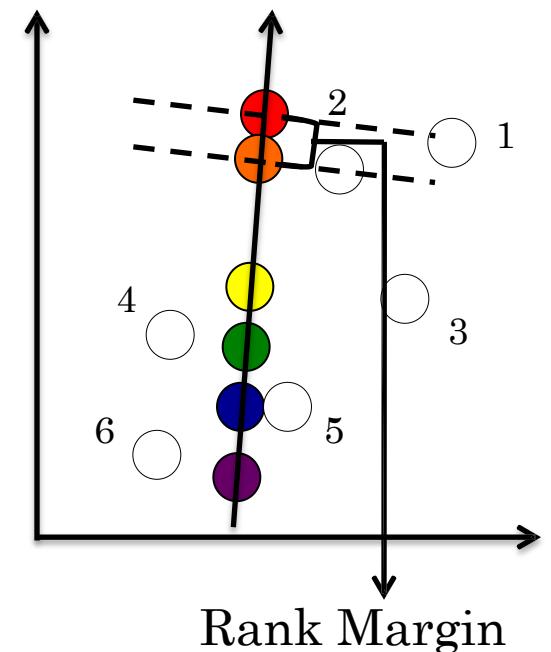


Image \rightarrow Relative Attribute Score



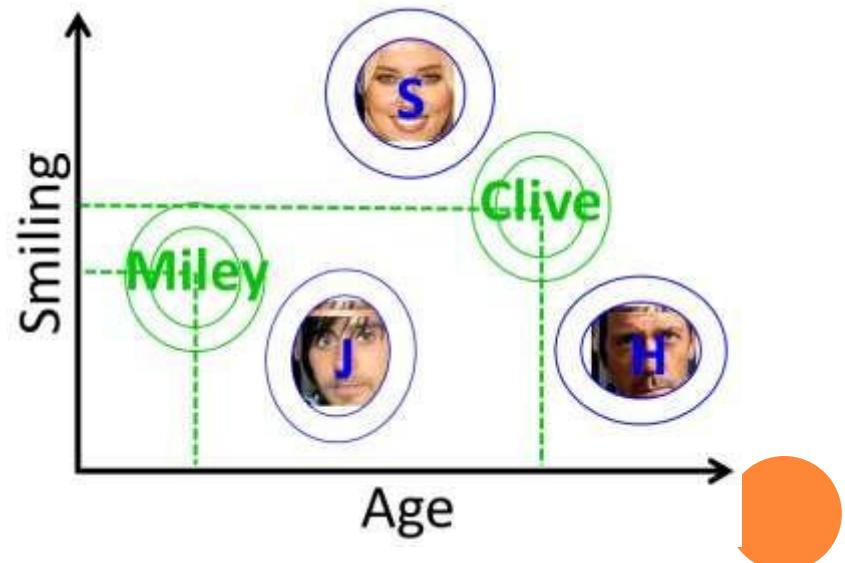
RELATIVE ZERO-SHOT LEARNING

- Each image is converted into a vector of relative attribute scores indicating the strength of each attribute
- A Gaussian distribution for each category is built in the relative attribute space. The distribution of unseen categories is estimated based on the specified constraints and the distributions of seen categories
- Max-likelihood is then used for classification

Age: Hugh \searrow Clive \searrow Scarlett
Jared \searrow Miley

Smiling: Miley \searrow Jared

Blue: Seen class Green: Unseen class



RELATIVE IMAGE DESCRIPTION

Binary:

Not natural

Not open

Has perspective



Relative :

More natural than tallbuilding
Less natural than forest

More open than tallbuilding
Less open than coast

Has more perspective than tallbuilding



WHITTLE SEARCH

Query: "I want a bright,
open shoe that is short
on the leg."



Round 1

More open than



Round 2

Selected feedback

More bright in color than
Less ornaments than



Round 2

Less high at the heel than



Round 3

More formal than
More bright in color than
Higher at the heel than

More open than



[Kovashka, Parikh, & Grauman, CVPR 2012]



SUMMARY

Semantic attribute classifiers can be useful for:

- + Describing images of unknown objects [Farhadi et al, CVPR 2009]
- + Recognizing unseen classes [Lampert et al, CVPR 2009]
- + Reducing dataset bias (trained across classes)
- + Effective object search in surveillance videos [Vaquero et al, WACV 2009]
- + Compact descriptors / Efficient image retrieval [Douze et al, CVPR 2011]
- + Fine-grained object categorization [Wah et al, ICCV 2011]
- + Face verification [Kumar et al, 2009], Action recognition [Liu et al, CVPR 2011], Person re-identification [Layne et al, BMVC 2012] and other classification tasks.
- + Other applications, such as sentence generation from images [Kulkarni et al, CVPR 2011], image aesthetics prediction [Dhar et al CVPR 2011], ...