



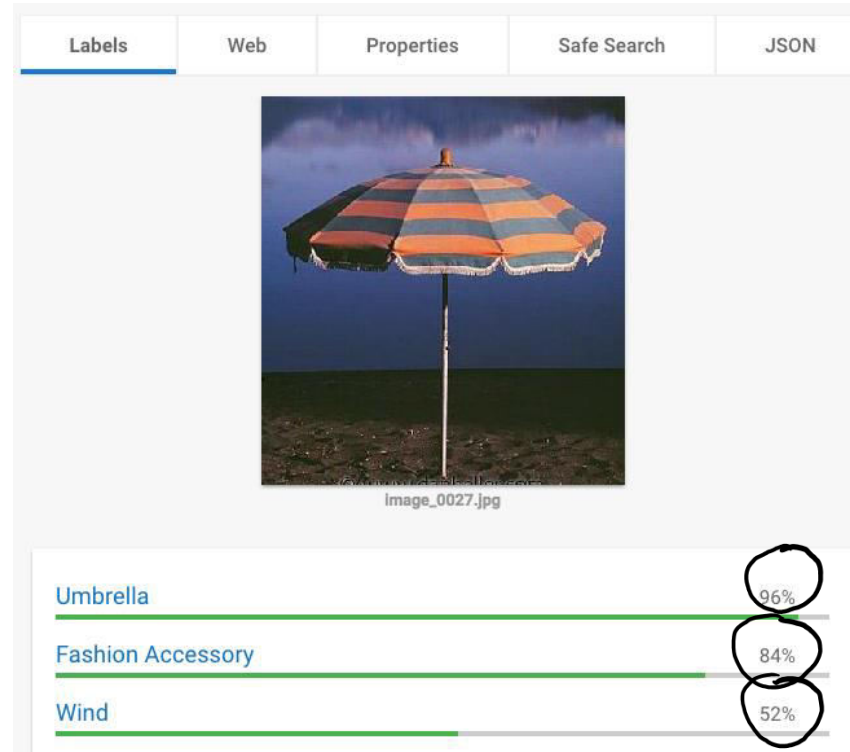
# Video 11.1

## Kostas Daniilidis

# By object recognition we might mean...

Image classification:

Input is an image  
and output is a set of  
class labels with  
probabilities



<https://cloud.google.com/vision/docs/drag-and-drop>

# .. Or here:

clarifai

[Demo](#)

[Solutions](#) ▾

[Pricing](#)

[Developer API](#)

[Resources](#) ▾

[Get your free API key](#)

[Talk to us](#)



Clarifai Demo

[Configure](#)

GENERAL MODEL

PREDICTED CONCEPT

PROBABILITY

vehicle	0.987
police	0.978
road	0.957
car	0.952
action	0.936
transportation system	0.936
people	0.933
street	0.931

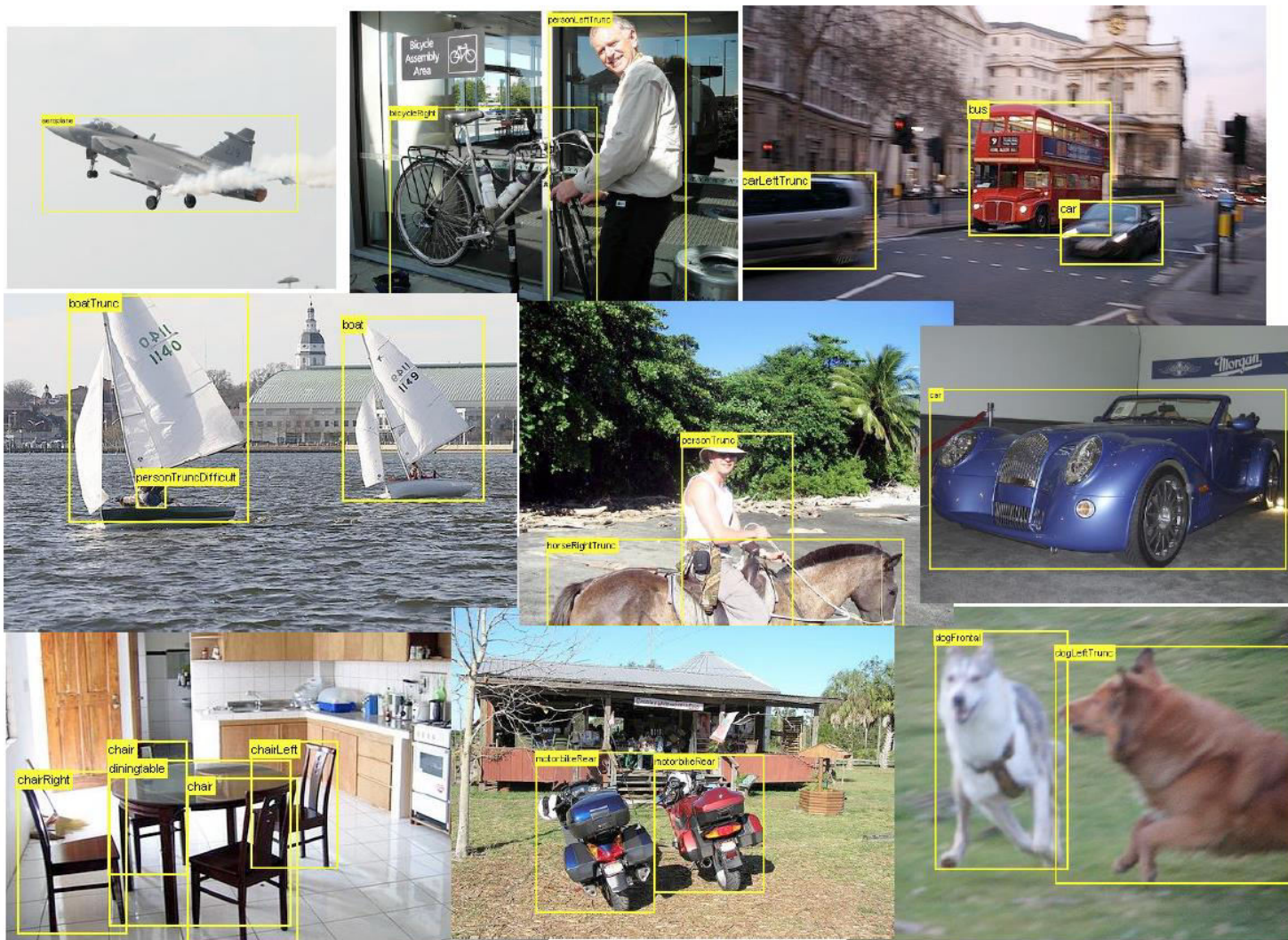
# But object recognition usually means...

## Object Detection and Localization:

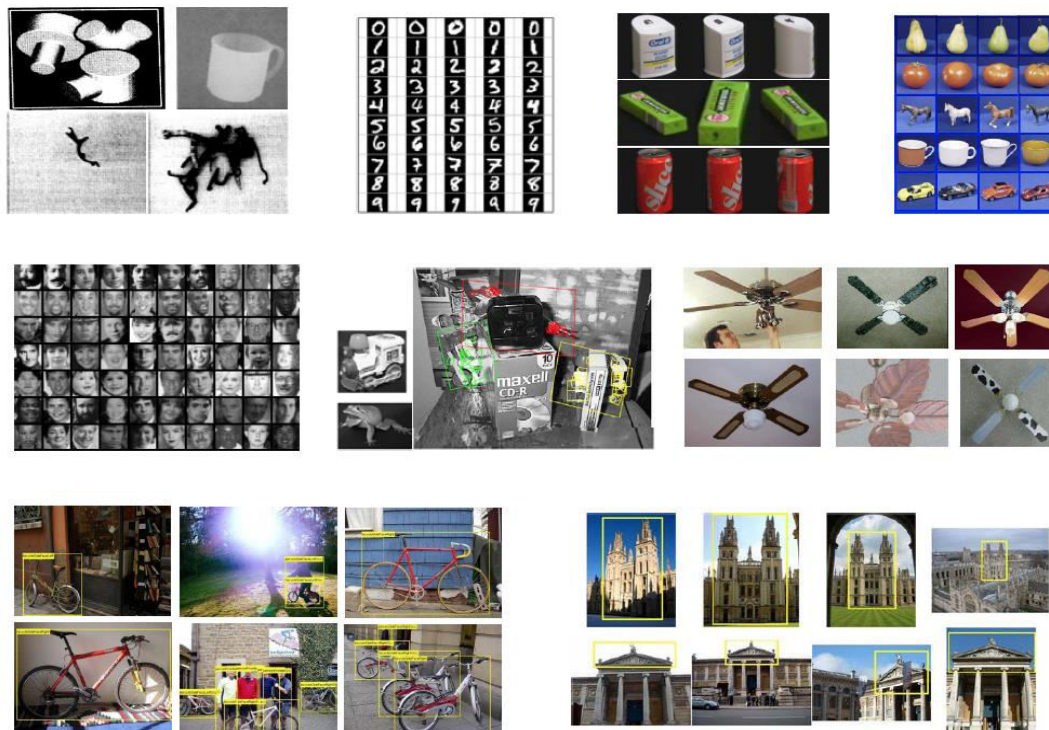
Find where is a **car** in the image?







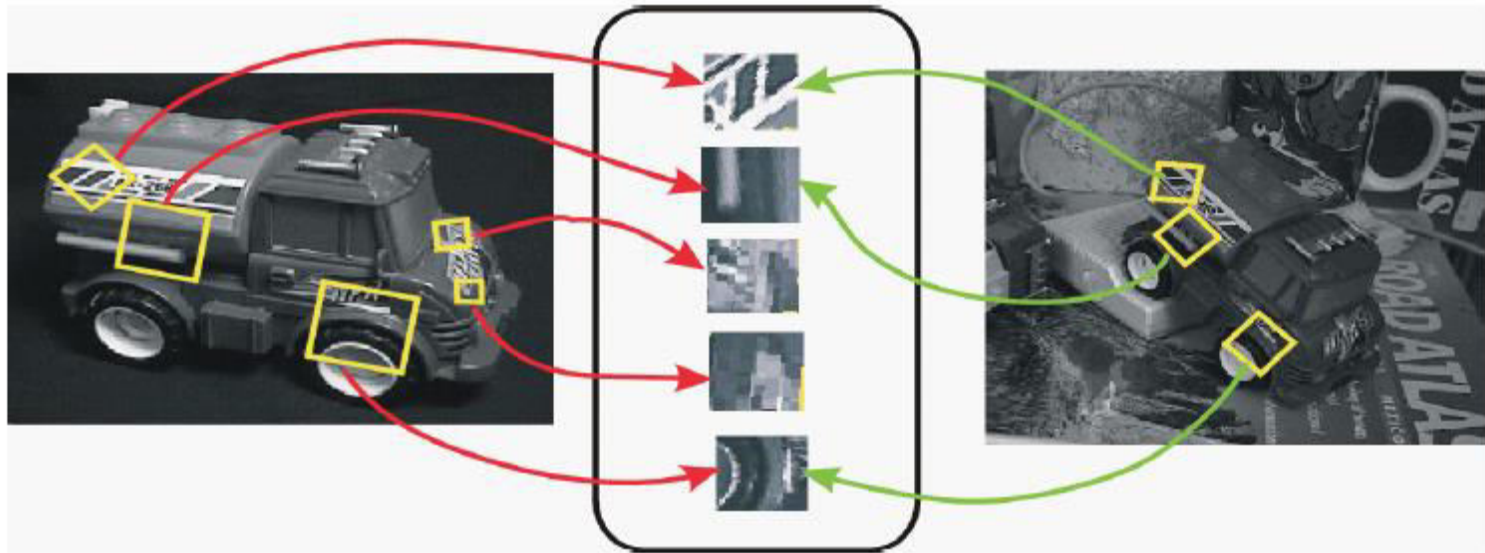
# History (Grauman and Leibe, 2011)



# Recognition ingredients

- Object **features** (for example, HOG) that is resistant to geometric and photometric nuisances.
- Learning (**training**) the class model (SVM) to represent intra-class variations.
- **Testing** on an input image and output all bounding boxes for a specific class (Sliding window)

## Instance recognition by matching (Lowe, 2004)



Match SIFT Features and verify geometric consistency





# Video 11.2

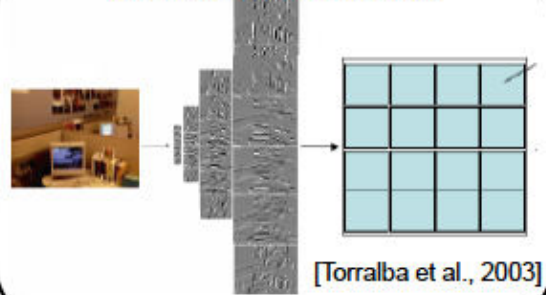
## Kostas Daniilidis

# Holistic approaches (Grauman and Leibe, 2011)

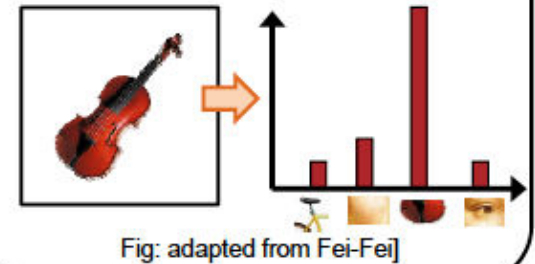
## Subspace-based, e.g., Eigenfaces



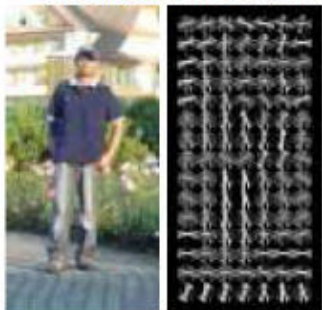
## GIST scene descriptor



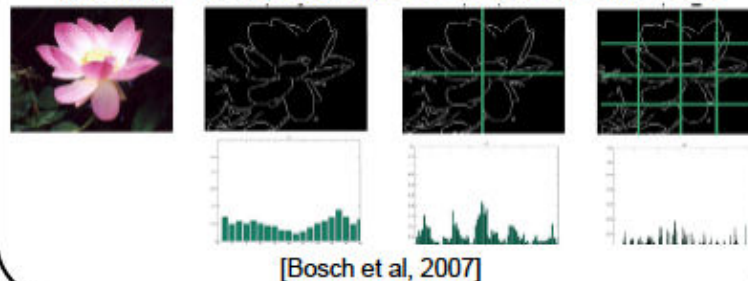
## Bag-of-words



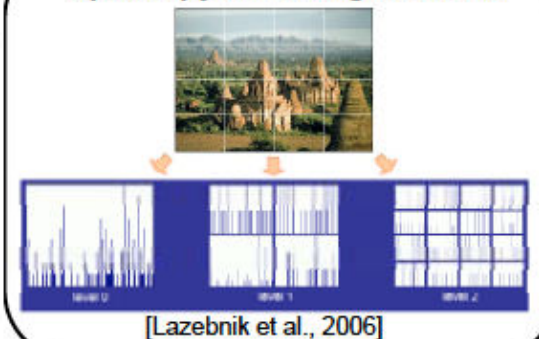
## Histograms of Oriented Gradients



## Pyramid of Histograms of Oriented Gradients

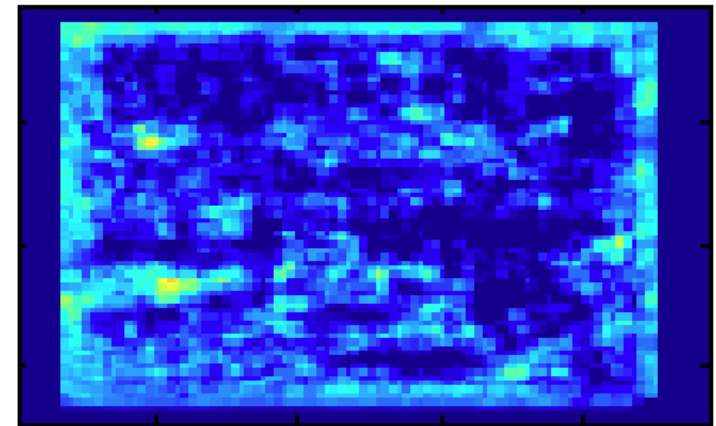
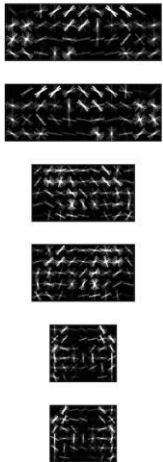


## Spatial pyramid bag-of-words



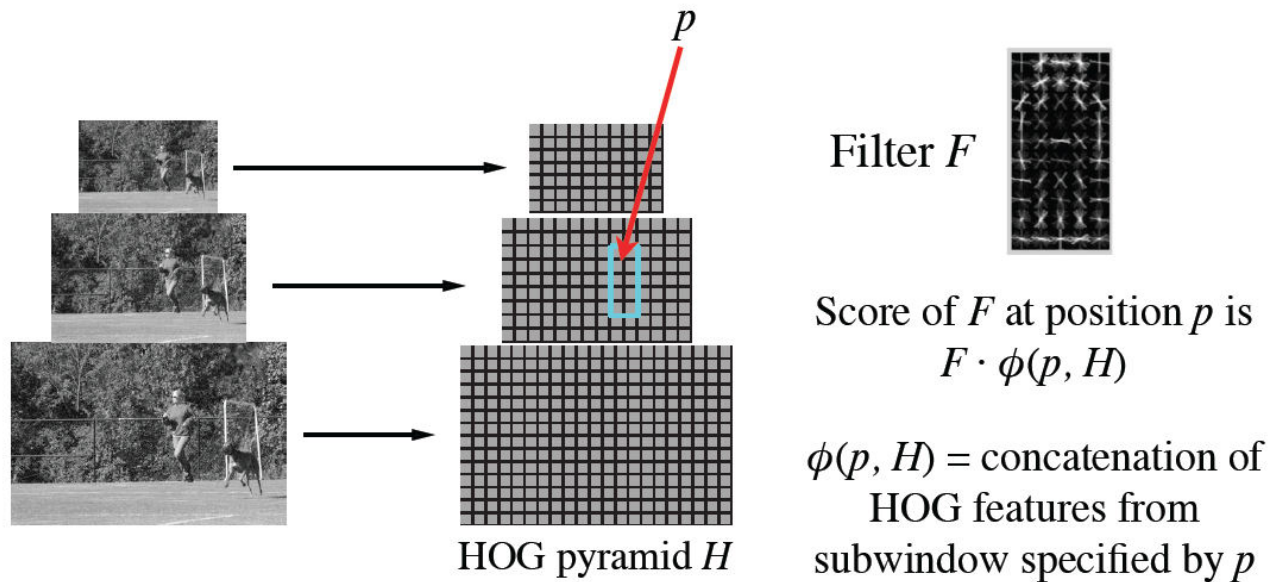
# Class Recognition: Sliding windows

- Learn a HOG representative of a car over multiple scales
- What happens if we **slide** this window over an input image?

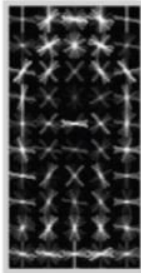
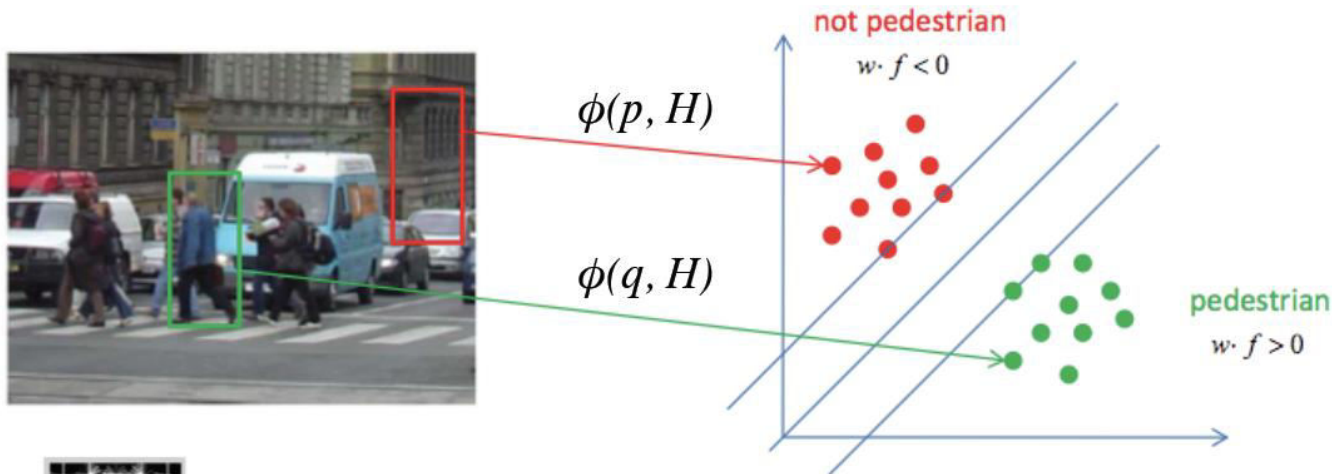


# HOG filter

- Array of weights for features in subwindow of HOG pyramid
- Score is dot product of filter and feature vector



# Dalal & Triggs: HOG + linear SVMs



Typical form of  
a model

There is much more background than objects

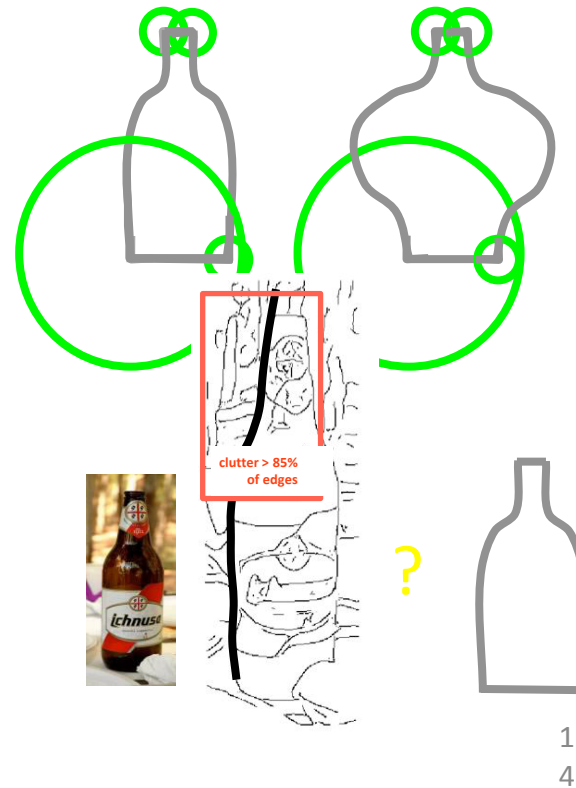
Start with random negatives and repeat:

- 1) Train a model
- 2) Harvest false positives to define “hard negatives”

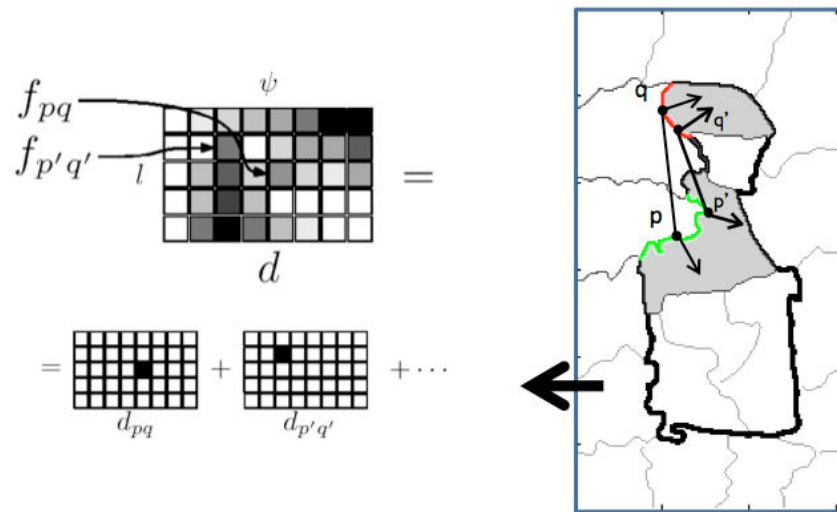


# Holistic / Global Representation

- + Gestaltism: shape is perceived as a whole:
  - not merely sum of parts
  - global dependencies



# Localization through Shape Verification

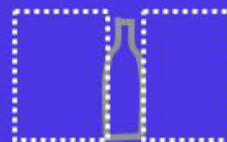


# Holistic approaches need perceptual grouping

Bottom-up perceptual grouping:  
+ Strong prior on possible shapes.

- Grouping is not always repeatable / stable.

Is there a bottle?



No edge  
grouping

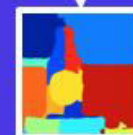


Contour and  
Region grouping

Are segmentations the same?



≠



# Chordigram

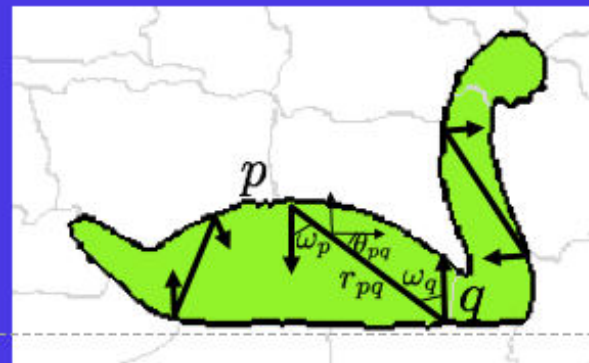
**Idea:** capture all global dependences among edges.

**Chord:** pair of boundary edges  $(p, q)$ .

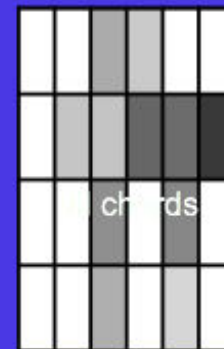
**Chord feature:**  $f_{pq} = (r_{pq}, \theta_{pq}, \omega_p, \omega_q)$

**Chordigram:** histogram of chord features:

$$ch_k = \#\{(p, q) | f_{pq} \in bin(k)\}$$

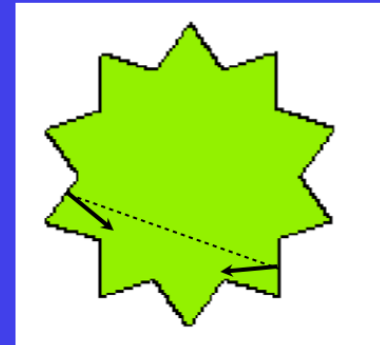
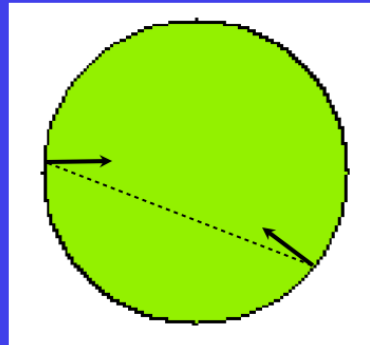


↓ binning



# Chord Features

Chord normals  
capture local finer  
shape information:



Chord diagram over  
only two features:

Chord normals  $\omega \rightarrow$



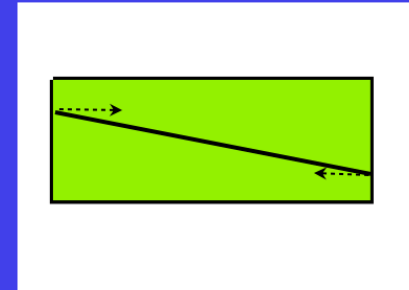
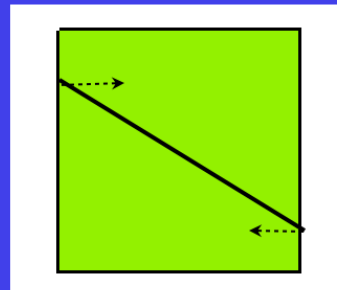
Chord length and  
orientation  $r, \theta \rightarrow$





# Chord Features

Chord length and orientation capture global coarse shape information:



Chord diagram over only two features:

Chord normals  $\omega \rightarrow$



Chord length and orientation  $r, \theta \rightarrow$



# Gestalt Properties of the Chordigram

## Global:

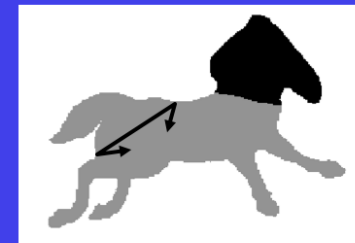
short as well as long range chords.



ch(centaur torso)

## Holistic:

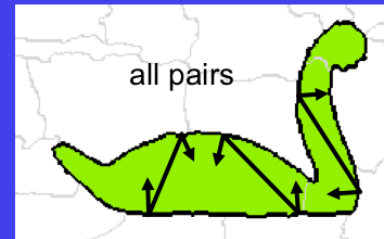
the chords of an edge are affected by all object parts.



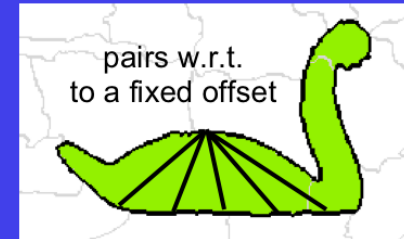
ch(horse torso)

# Transformation Properties of the Chordigram

**Translation invariance:**  
a chord captures only relative location.

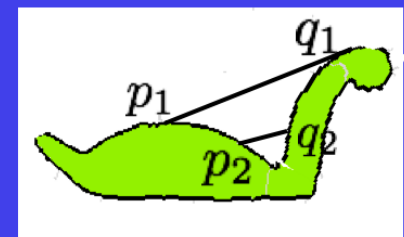
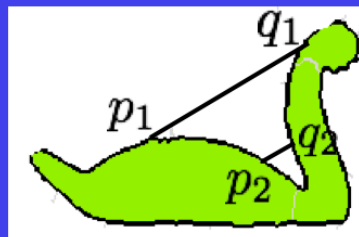


chordigram

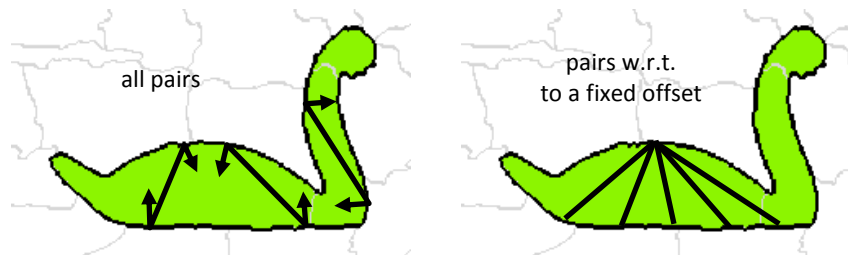


shape context

**Robustness to shape variation:**  
lengths are quantized uniformly in log space.



# Chord diagram vs Shape Context

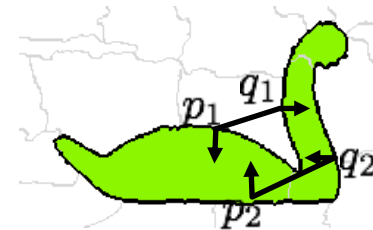


Translation invariance:	Yes	No
Global support:	Yes	No
Global configuration:	No	Yes
Notion of parts:	No	Yes

# Figure/Ground Organization

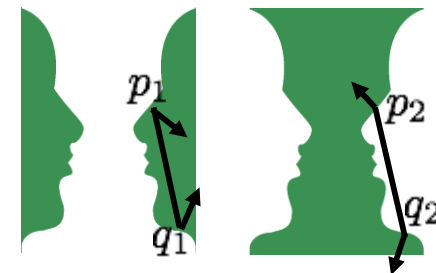
Chord normals point  
towards the object  
interior.

$$\text{chord}(p_1, q_1) \neq \text{chord}(p_2, q_2)$$

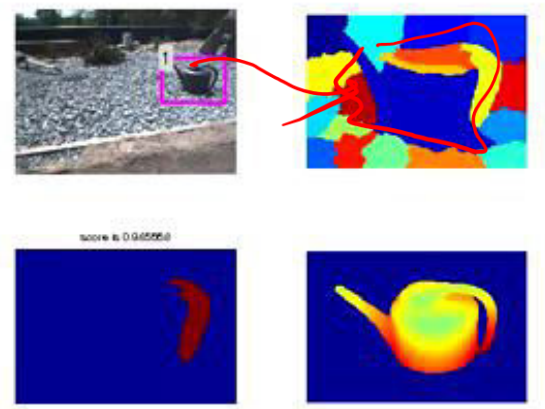
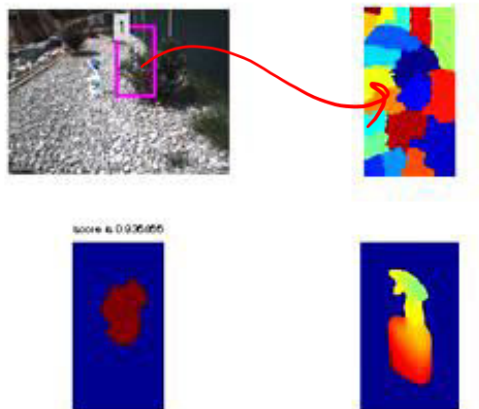
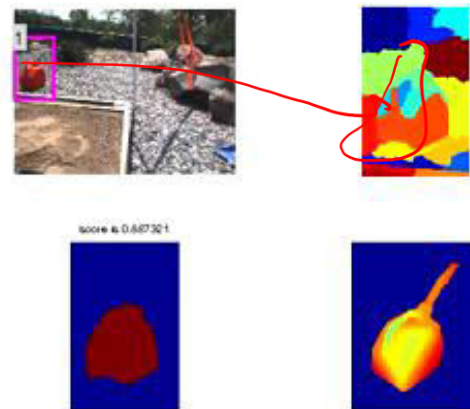
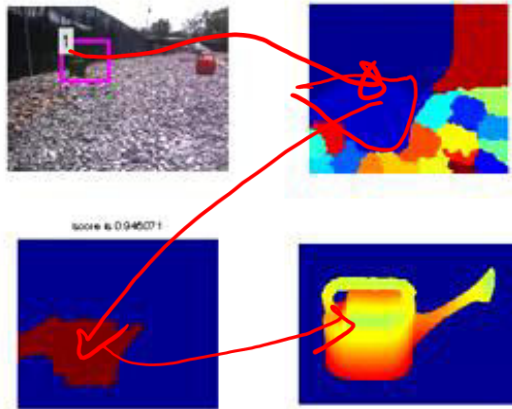


## Contour and interior descriptor:

More discriminative than pure  
contour descriptors.  
Encode figure/ground  
organization.









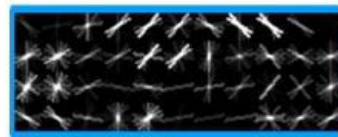
# Video 11.3

## Kostas Daniilidis

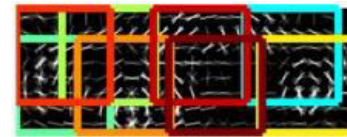
# Deformable Part Models

- DPM consists a Root and several Parts
  - Root represents holistic object shape
  - Part captures detailed part appearance

A car DPM model



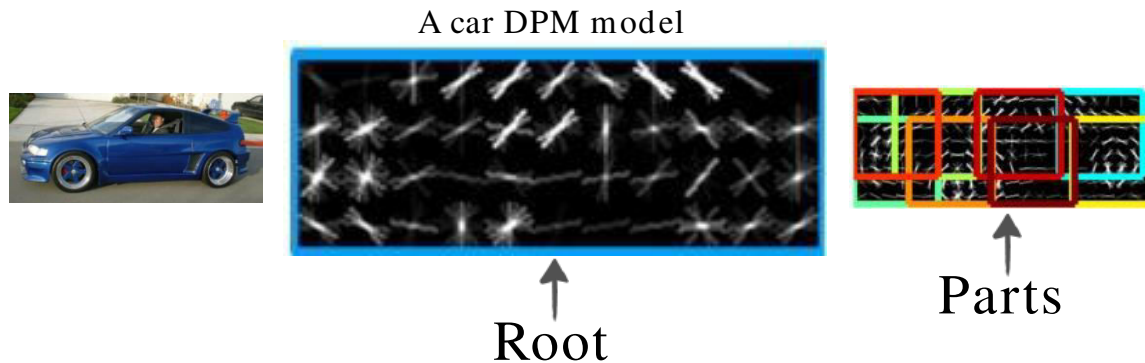
↑  
Root



↑  
Parts

# Deformable Part Models

- DPM consists a Root and several Parts
  - Root represents holistic object shape
  - Part captures detailed part appearance

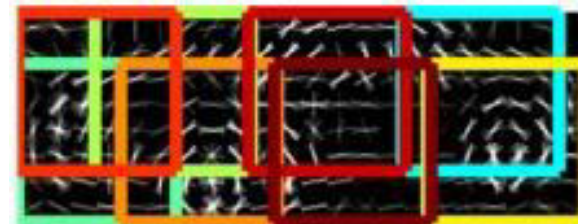


# Deformable Part Models

- Flexible part configuration
  - Parts are allowed to move around the anchor points (default position)



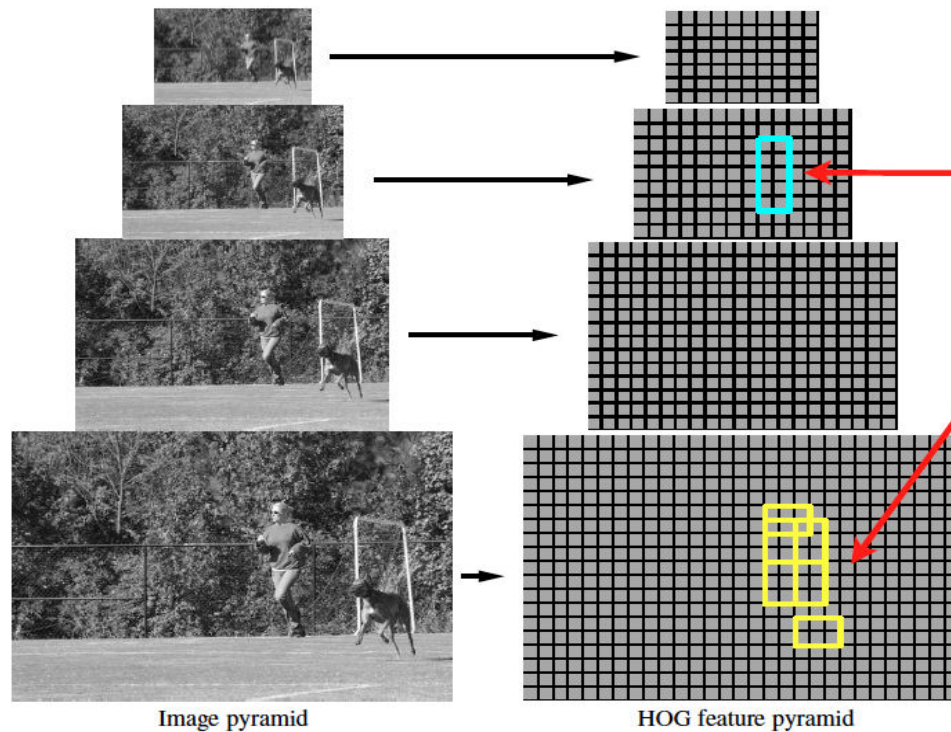
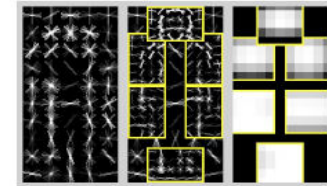
↑  
Deformation



↑  
Parts



# Object hypothesis



$$z = (p_0, \dots, p_n)$$

$p_0$  : location of root

$p_1, \dots, p_n$  : location of parts

Score is sum of filter  
scores minus  
deformation costs




Multiscale model captures features at two-resolutions

# Deformable Part Models

- Detector response at a given root location

$$x = (r, c, l)$$

$$F'_0 \cdot \phi(H, x) + \sum_{k=1}^n \max_{x_k} \left( F'_k \cdot \phi(H, x_k) - d_k \cdot \phi_d(\delta_k) \right) + b$$

  
Root
  
Parts
  
bias

# Score of a hypothesis

$$\text{score}(p_0, \dots, p_n) = \sum_{i=0}^n F_i \cdot \phi(H, p_i) - \sum_{i=1}^n d_i \cdot (dx_i^2, dy_i^2)$$

“data term”

$\uparrow$

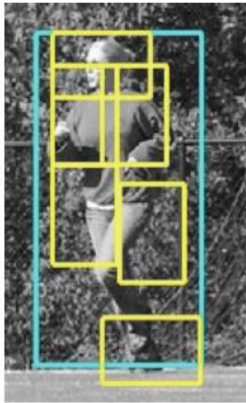
filters

“spatial prior”

$\uparrow$

displacements

deformation parameters



$$\text{score}(z) = \beta \cdot \Psi(H, z)$$

$\uparrow$

concatenation filters and  
deformation parameters

$\uparrow$

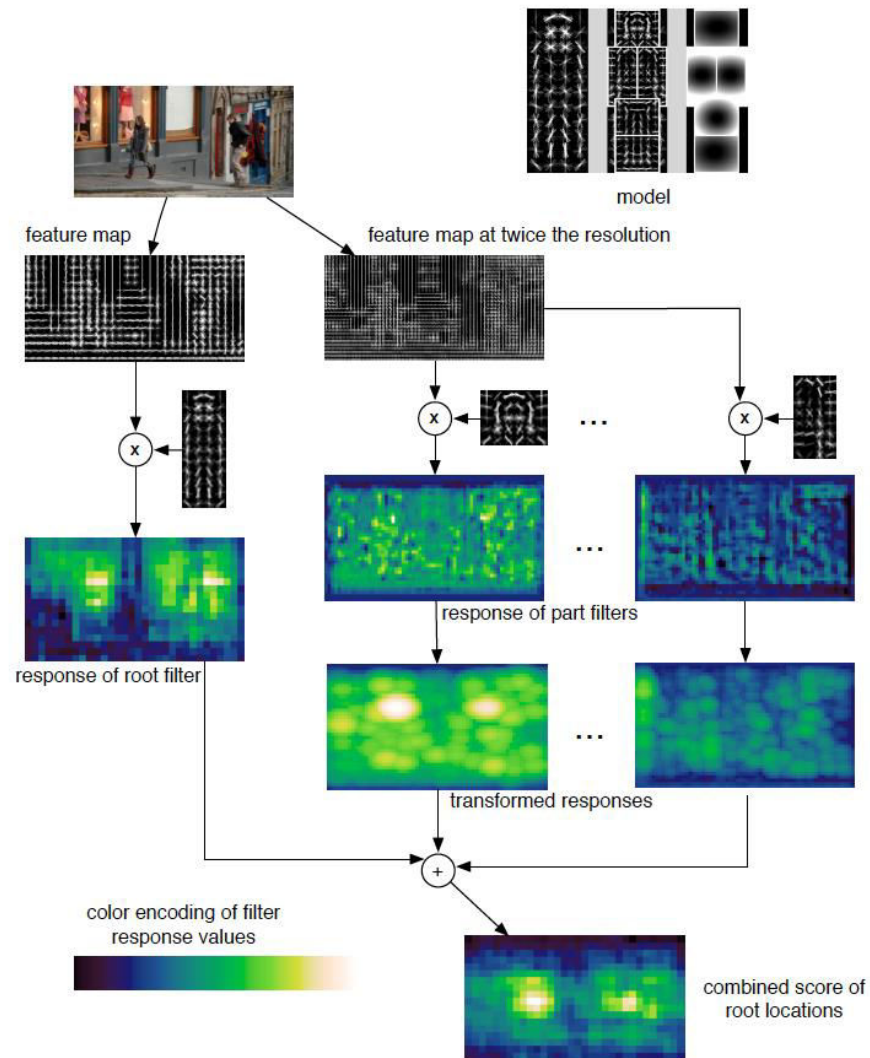
concatenation of HOG  
features and part  
displacement features

# Matching

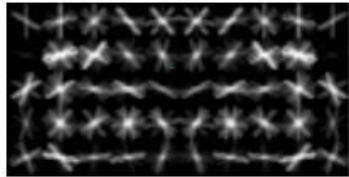
- Define an overall score for each root location
  - Based on best placement of parts

$$\text{score}(p_0) = \max_{p_1, \dots, p_n} \text{score}(p_0, \dots, p_n).$$

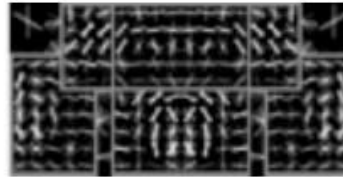
- High scoring root locations define detections
  - “sliding window approach”
- Efficient computation: dynamic programming + generalized distance transforms (max-convolution)



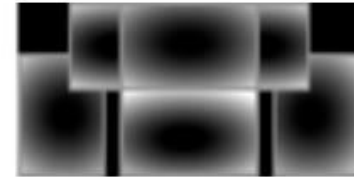
# Car model



root filters  
coarse resolution

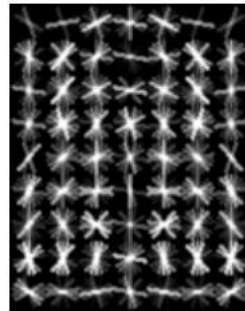
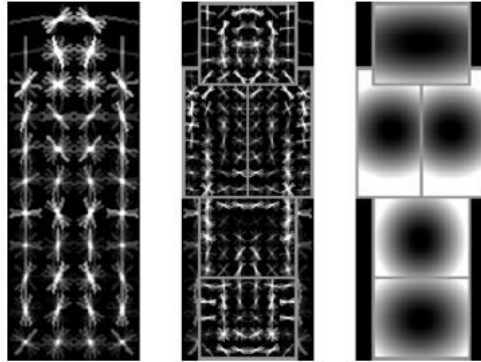


part filters  
finer resolution

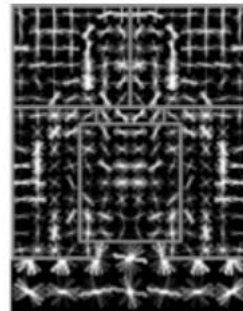


deformation  
models

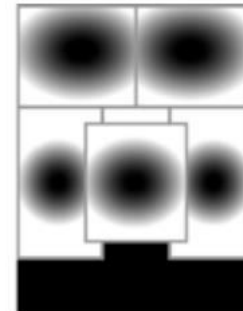
# Person model



root filters  
coarse resolution



part filters  
finer resolution



deformation  
models



# Precision/Recall results on Bicycles 2008

