## Solved, Half-Solved and Unsolved Problems in Visual Recognition

Jitendra Malik University of California at Berkeley

| "Possibly outdoor fight. Twe seem to scene, maybe a fam. I could not fam. I fam. I fam. I fam. I could not fell for sure."     "Some kind of game or fight. Twe groups of two men. One in the foreground was getting a fast in the on grass. Another on grass. Another sure."     "Some kind of game or fight. Twe groups of two men. One in the foreground was getting a fast in the on grass. Another is scene, "     "Some kind of game or fight. Twe groups of two men. One in the fore outset, we still a fast in the scene of the scene."     Some kind of game or fight. Twe groups of two men. One in the fore outset, we still a fast in the scene of the scene."     Some kind of game or fight. Twe groups of two men. One in the fore outset, we still a fast in the scene of the scene."     Some kind of game or fight. Twe groups of two men. One in the fore outset, we still a fast in the fore outset with this of a game, tough game though, more like rupby tha fore a back for the scene." | Image shown to subjects                                    | 40ms  | 80ms   | 107ms  | 500ms   |
|---|--|---|--|--|---|
| Pierre 2. Human authings and sting on other holds some in an interest house for different and outration   |  | "Possibly<br>outdoor<br>scene,<br>maybe a<br>farm. I<br>could not<br>tell for<br>sure." | "There<br>seem to<br>be two<br>people in<br>the<br>center of<br>the<br>scene." | " People playing<br>rugby. Two<br>persons in close<br>contact, wreatling,<br>on grass. Another<br>man more distant.<br>Goal in sight." | "Some kind of game or fight. Two<br>groups of two men. One in the<br>foreground was getting a fix in the<br>face. Outdoors, because I see grass<br>and maybe lines on the grass? That<br>is why I think of a game, rough<br>game though, more like rugby than<br>football because they weren't in<br>pads and hemets" |
| Figure 2. Human subjects reporting on what ne sue saw in an image shown for different presentation durations (PD=27, 40, 67, 80, 107, 500ms). From Fei-Fei and Perona [26].   | Figure 2. Human subjects re<br>durations (PD=27, 40, 67, 8 | eporting on v<br>0, 107, 500r   | what he/she<br>ns). From F   | saw in an image sh<br>ei-Fei and Perona  | own for different presentation<br>26].  |









| mage shown to subjects      | 40ms<br>"Possibly<br>outdoor<br>scene,<br>maybe a<br>farm. I<br>could not<br>tell for | 80ms<br>"There<br>seem to<br>be two<br>people in<br>the<br>center of<br>the | 107ms<br>"People playing<br>rugby. Two<br>persons in close<br>contost, wreotling,<br>on grass. Another<br>man more distant.<br>Goal in sight." | 500ms<br>"Some kind of game or fight. Two<br>groups of two men. One in the<br>foreground was getting a fist in the<br>foreground was getting a fist in the<br>foreground was getting a fist on the<br>foreground was getting a fist<br>and maybe lines on the grass? That<br>is why I think of a game, rough<br>game though, more like roughy than |
|-----------------------------|---|---|--|--|
| Figure 2. Human subjects ro | eporting on v   | what he/she   | saw in an image sh   | pads and helmets"<br>own for different presentation<br>261   |

#### We need to identify

- Objects
- Agents
- Relationships among objects with objects, objects with agents, agents with agents ...
- Events and Actions

Berkeley

Computer Vision Group



# A brief history of computer vision .. Those who cannot remember the past are condemned to repeat it -George Santayana

#### Fifty years of computer vision 1963-2013

- 1960s: Beginnings in artificial intelligence, image processing and pattern recognition
- 1970s: Foundational work on image formation: Horn, Koenderink, Longuet-Higgins ...
- 1980s: Vision as applied mathematics: geometry, multi-scale analysis, probabilistic modeling, control theory, optimization
- 1990s: Geometric analysis largely completed, vision meets graphics, statistical learning approaches resurface
- 2000s: Significant advances in visual recognition, range of practical applications

UC Berkeley

**Object recognition in computer vision** 

- Recognition as Pose Estimation
- Recognition as Description using Volumetric primitives
- Recognition as Pattern Classification
- Recognition as Deformable Matching

University of Ca Berkeley

Computer Vision Group

Computer Vision Group

### **Recognition as Pose Estimation: Object as a set of points in 3D** • Roberts (1963), Faugeras & Hebert (1983), Huttenlocher & Ullman (1987) Variants - Geometric Hashing : Lamdan & Wolfson (1988) - Pose Clustering : Stockman (1987), Olson (1994) - Linear Combination of Views: Basri & Ullman (1991)

Berkeley

Computer Vision Group

Berkeley

#### **Recognition as Fitting Volumetric Primitives: Object as a hierarchy of simple shapes**

- Binford (1971), Marr & Nishihara (1978), Biederman(1987)
- Discredited as an approach for recognition in general, it has retained appeal for analyzing images of people

Computer Vision Group

The Stick Figure Ideal



#### **Recognition as Statistical Pattern Classification: Object as a feature vector**

- Optical Character Recognition studied as far back as the 1950s. Recent years focus on handwritten digit classification and face detection.
- Some examples:
  - Neural networks: Neocognitron (Fukushima, 1980, 1988) , Convolution Neural Networks (LeCun et al), C2 Features (Serre, Wolf & Poggio 2005)
  - Support Vector Machines (various)

  - Decision Trees (Amit, Geman, & Wilder, 1997)
    Boosted Decision Trees (Viola & Jones, 2001)

University of C Berkeley

Computer Vision Group

#### **Recognition as Pictorial Structure Matching: Object as a configuration of feature points**

- Transformations to model shape variation-D'Arcy Wentworth Thompson (1910)
- · Grenander (1970s and later)probabilistic models ontransformations
- Fischler and Elschlager (1973) deformable matching of landmarks, "point masses", in a configuration of "springs" to model deformable templates.
- Von derMalsburg-dynamic link architecture for neural modelling, elastic graph matching for face recognition (1993, 1997)
- Felzenszwalb and Huttenlocher (2000) pictorial structures for aligning human bodies to stick figures using dynamic programming
- Belongie, Malik &Puzicha (2001) use"shape contexts" as point descriptors, and thin plate splines to model deformation

University of C Berkeley

Computer Vision Group

| Handwritten digit recognition<br>(MNIST,USPS)  | EZ-Gimpy Results (N       | Mori & Malik, 2003)    |
|--|---------------------------|------------------------|
|  | • 171 of 192 images corre | ectly identified: 92 % |
| 0123436789   | horse                     | spade                  |
| <ul> <li>LeCun's Convolutional Neural Networks variations (0.8%,<br/>0.6% and 0.4% depending on different ways of virtually<br/>augmenting dataset)</li> </ul> | horse                     | spade                  |
| • SVMs (DeCoste & Scholkopf : 0.6%)  | SIMTAGE                   | Join                   |
| <ul> <li>K-NN based Shape context/TPS matching (Belongie, Malik &amp;<br/>Puzicha: 0.6%)</li> </ul>  | smile                     | join                   |
| • On USPS comparison to humans: 2.5% (Bromley and  | Canvas                    | where                  |

Sackinger, 1991), <sup>University</sup> distance; 2.59% ct. Zhang e Computer Vision Group

## join where

Computer Vision Group



## Multiscale sliding window



UC Berkeley

Paradigm introduced by Rowley, Baluja & Kanade 96 for face detection Viola & Jones 01, Dalal & Triggs 05, Felzenszwalb, McAllester, Ramanan 08

#### Problems with the multi-scale scanning paradigm

#### •Computational complexity

• Not natural for irregularly shaped objects

Computer Vision Group

- Segmentation is delinked
- Context is delinked

UC Berkelev

### Caltech-101 [Fei-Fei et al. 04]

#### 102 classes, 31-300 images/class





















Patches are often far visually, but they are close semantically (Bourdev& Malik, 09; Bourdev et al, 10)

## How do we train a poselet for a given pose configuration?



### **Finding Correspondences**



Given part of a human pose



How do we find a similar pose configuration in the training set?



We use keypoints to annotate the joints, eyes, nose, etc. of people

## Finding Correspondences



Residual Error



## Training poselet classifiers



- 1. Given a seed patch
- 2. Find the closest patch for every other person
- 3. Sort them by residual error
- 4. Threshold them

## Training poselet classifiers



- 1. Given a seed patch
- 2. Find the closest patch for every other person
- 3. Sort them by residual error
- 4. Threshold them
- 5. Use them as positive training examples for a classifier (HOG features, linear SVM)

## How do we find poselets?

Choose thousands of random windows, generate poselet candidates, train linear SVMs



- Select a small set of poselets that are:
  - Individually effective
  - Complementary

## Segmenting people (Brox et al, CVPR 2011)



## Actions in still images ...



#### have characteristic :

- pose and appearance
- interaction with objects and agents

## Some discriminative poselets





















- "Morpho-kinetics" of action (shape and movement of the body)
- Identity of the object/s
- Activity context



10 May 2011

53









| The more you look, the more you see!                        |   |  |  |   |  |  |  |
|---|---|--|--|---|--|--|--|
| Image shown to subjects                                     | 40ms  | 80ms   | 107ms  | 500ms   |  |  |  |
|   | "Possibly<br>outdoor<br>scene,<br>maybe a<br>farm. I<br>could not<br>tell for<br>sure." | "There<br>seem to<br>be two<br>people in<br>the<br>center of<br>the<br>scene." | " People playing<br>rugby. Two<br>persons in close<br>contact, wrestling,<br>on grass. Another<br>man more distant.<br>Goal in sight." | "Some kind of game or fight. Two<br>groups of two men. One in the<br>foreground was getting a fist in the<br>face. Outdoors, because 1 see gross<br>and maybe lines on the grass? That<br>is why 1 think of a game, rough<br>game though, more like rough than<br>football because they weren't in<br>pads and helmets" |  |  |  |
| Figure 2. Human subjects re<br>durations (PD=27, 40, 67, 80 | porting on v<br>0, 107, 500r  | what he/she<br>ns). From F   | saw in an image sh<br>ei-Fei and Perona [  | own for different presentation<br>26].  |  |  |  |
|   |   |  |  |   |  |  |  |

#### So much remains to be done...

- Objects, Scenes, Events
- The semantic gap is to be confronted, not avoided!

UC Berkeley

Computer Vision Group