04/12/12

# Object Category Detection: Parts-based Models

Computer Vision CS 543 / ECE 549 University of Illinois

Derek Hoiem

## Administrative stuff

- HW 5 due Tues
- I'm out of town Sun to Tues night
  Unresponsive on Monday
  - Amin Sadeghi will teach Tues
- Final project posters presented May 4 (1:30-4:30pm)
  - Final paper: due May 8

## Goal: Detect all instances of objects

Cars



Faces





Cats

## Last class: sliding window detection





## Object model: last class

- Statistical Template in Bounding Box
  - Object is some (x,y,w,h) in image
  - Features defined wrt bounding box coordinates



Image



**Template Visualization** 

Images from Felzenszwalb

# Last class: statistical template

 Object model = log linear model of parts at fixed positions



## When do statistical templates make sense?



#### Caltech 101 Average Object Images

## **Object models: this class**

- Articulated parts model
  - Object is configuration of parts
  - Each part is detectable





# Deformable objects



Images from Caltech-256

## **Deformable objects**

















































Images from D. Ramanan's dataset

Slide Credit: Duan Tran

# **Compositional objects**



## Parts-based Models

Define object by collection of parts modeled by

- 1. Appearance
- 2. Spatial configuration







Slide credit: Rob Fergus

• One extreme: fixed template



• Another extreme: bag of words



• Star-shaped model

![](_page_14_Figure_2.jpeg)

• Star-shaped model

![](_page_15_Picture_2.jpeg)

Tree-shaped model

• Many others...

![](_page_17_Picture_2.jpeg)

Csurka '04 Vasconcelos '00

![](_page_17_Picture_4.jpeg)

- b) Star shape
- Leibe et al. '04, '08 Crandall et al. '05 Fergus et al. '05

![](_page_17_Figure_7.jpeg)

Crandall et al. '05

![](_page_17_Picture_9.jpeg)

Felzenszwalb & Huttenlocher '05

![](_page_17_Figure_11.jpeg)

Bouchard & Triggs '05

![](_page_17_Picture_13.jpeg)

![](_page_17_Picture_14.jpeg)

g) Sparse flexible model

Carneiro & Lowe '06

#### from [Carneiro & Lowe, ECCV'06]

# Today's class

- 1. Star-shaped model
  - Example: ISM
    - <u>Leibe et al. 2004, 2008</u>

![](_page_18_Figure_4.jpeg)

- 2. Tree-shaped model
  - Example: Pictorial structures
    - Felzenszwalb Huttenlocher 2005

![](_page_18_Picture_8.jpeg)

# ISM: Implicit Shape Model

#### Training overview

- Start with bounding boxes and (ideally) segmentations of objects
- Extract local features (e.g., patches or SIFT) at interest points on objects
- Cluster features to create codebook
- Record relative bounding box and segmentation for each codeword

![](_page_19_Figure_6.jpeg)

# ISM: Implicit Shape Model

#### **Testing overview**

- Extract interest points in test image
- Softly match to codebook entries
- Each matched codeword votes for object bounding box
- Compute modes of votes using mean-shift
- Check which codewords voted for modes

![](_page_20_Figure_7.jpeg)

## **Codebook Representation**

- Extraction of local object features
  - Interest Points (e.g. Harris detector)
  - Sparse representation of the object appearance

![](_page_21_Picture_4.jpeg)

- Collect features from whole training set
- Example:

![](_page_21_Picture_7.jpeg)

## **Agglomerative Clustering**

- Algorithm (Average-Link)
  - 1. Start with each patch as a cluster of its own
  - 2. Repeatedly merge the two most similar clusters X and Y, where the similarity between two clusters is defined as the average similarity between their members

$$sim(X,Y) = \frac{1}{NM} \sum_{i=1}^{N} \sum_{j=1}^{M} sim(x^{(i)}, y^{(j)})$$

3. Until  $sim(X,Y) < \theta$ 

- Commonly used similarity measures
  - Normalized correlation
  - > Euclidean distances

## **Appearance Codebook**

![](_page_23_Picture_1.jpeg)

#### • Clustering Results

- Visual similarity preserved
- > Wheel parts, window corners, fenders, ...
- > Store cluster centers as Appearance Codebook

## **Voting with Local Features**

• For every feature, store possible "occurrences"

![](_page_24_Picture_2.jpeg)

Record relative size and scale of object

• For new image, let the matched features vote for possible object positions

![](_page_24_Picture_5.jpeg)

![](_page_24_Picture_6.jpeg)

 $x_{vote} = x_{img} - x_{occ}(s_{img}/s_{occ})$  $y_{vote} = y_{img} - y_{occ}(s_{img}/s_{occ})$  $s_{vote} = (s_{img}/s_{occ}).$ 

K. Grauman, B. Leibe

## Implicit Shape Model - Recognition

![](_page_25_Figure_1.jpeg)

[Leibe04, Leibe08]

## **Scale Voting: Efficient Computation**

![](_page_26_Figure_1.jpeg)

- Mean-Shift formulation for refinement
  - Scale-adaptive balloon density estimator

$$\hat{p}(o_n, x) = \frac{1}{V_b} \sum_k \sum_j p(o_n, x_j | f_k, \ell_k) K(\frac{x - x_j}{b})$$

## Implicit Shape Model - Recognition

Matched Codebook

![](_page_27_Figure_1.jpeg)

![](_page_27_Picture_2.jpeg)

![](_page_27_Figure_3.jpeg)

![](_page_27_Picture_4.jpeg)

<sup>[</sup>Leibe04, Leibe08]

![](_page_28_Picture_1.jpeg)

#### Original image

![](_page_29_Picture_1.jpeg)

#### **Interest points**

![](_page_30_Picture_1.jpeg)

#### Matched patches

![](_page_31_Picture_1.jpeg)

**Prob. Votes** 

![](_page_32_Picture_1.jpeg)

#### 1<sup>st</sup> hypothesis

![](_page_33_Picture_1.jpeg)

#### 2<sup>nd</sup> hypothesis

![](_page_34_Picture_1.jpeg)

#### 3<sup>rd</sup> hypothesis

## **ISM: Detection Results**

- Qualitative Performance
  - Robust to clutter, occlusion, noise, low contrast

![](_page_35_Picture_3.jpeg)

# Beyond bounding boxes

![](_page_36_Picture_1.jpeg)

Backprojected codewords can vote:

- Pixel segmentation
- Part layout
- Pose
- Depth values

![](_page_36_Picture_7.jpeg)

![](_page_36_Picture_8.jpeg)

## Segmentation: Probabilistic Formulation

![](_page_37_Picture_1.jpeg)

![](_page_37_Picture_2.jpeg)

• Influence of patch on object hypothesis (vote weight)

$$p(f,\ell|o_n,x) = \frac{\sum_i p(o_n,x|C_i)p(C_i|f)p(f,\ell)}{p(o_n,x)}$$

• Backprojection to features f and pixels p:  $p(\mathbf{p} = figure \mid o_n, x) = \sum_{\mathbf{p} \in (f, \ell)} p(\mathbf{p} = figure \mid f, \ell, o_n, x) p(f, \ell \mid o_n, x)$ Segmentation Influence on object hypothesis

[Leibe04, Leibe08] K. Grauman, B. Leibe

## **ISM - Top-Down Segmentation**

![](_page_38_Figure_1.jpeg)

#### **Example Results: Motorbikes**

![](_page_39_Picture_1.jpeg)

## **Example Results: Chairs**

![](_page_40_Picture_1.jpeg)

#### Dining room chairs

![](_page_40_Picture_3.jpeg)

![](_page_40_Picture_4.jpeg)

Office chairs

![](_page_40_Picture_6.jpeg)

![](_page_40_Picture_7.jpeg)

![](_page_40_Picture_8.jpeg)

![](_page_40_Picture_9.jpeg)

### **Inferring Other Information: Part Labels**

Training

![](_page_41_Picture_2.jpeg)

Test

![](_page_41_Picture_4.jpeg)

Output

![](_page_41_Picture_6.jpeg)

![](_page_41_Picture_7.jpeg)

![](_page_41_Picture_8.jpeg)

![](_page_41_Picture_9.jpeg)

![](_page_41_Picture_10.jpeg)

![](_page_41_Picture_11.jpeg)

![](_page_41_Picture_12.jpeg)

#### [Thomas07]

## Inferring Other Information: Part Labels (2)

![](_page_42_Picture_1.jpeg)

![](_page_42_Picture_2.jpeg)

## **Inferring Other Information: Depth Maps**

Test image

![](_page_43_Picture_2.jpeg)

![](_page_43_Picture_3.jpeg)

![](_page_43_Picture_4.jpeg)

![](_page_43_Picture_5.jpeg)

![](_page_43_Picture_6.jpeg)

![](_page_43_Picture_7.jpeg)

![](_page_43_Picture_8.jpeg)

![](_page_43_Picture_9.jpeg)

![](_page_43_Picture_10.jpeg)

![](_page_43_Picture_11.jpeg)

![](_page_43_Picture_12.jpeg)

![](_page_43_Picture_13.jpeg)

#### [Thomas07]

# 2 minute break

- Comparing a sliding window/part detector
  - What are the advantages of the part detector?
  - What are the advantages of the sliding window detector?

## Tree-shaped model

![](_page_45_Picture_1.jpeg)

## **Pictorial Structures Model**

![](_page_46_Picture_1.jpeg)

#### Felzenszwalb and Huttenlocher 2005

## **Pictorial Structures Model**

![](_page_47_Picture_1.jpeg)

$$P(L|I,\theta) \propto \left(\prod_{i=1}^{n} p(I|l_i, u_i) \prod_{(v_i, v_j) \in E} p(l_i, l_j | c_{ij})\right)$$
  
Appearance likelihood Geometry likelihood

# Modeling the Appearance

- Any appearance model could be used
  - HOG Templates, etc.
  - Here: rectangles fit to background subtracted binary map
- Can train appearance models independently (easy, not as good) or jointly (more complicated but better)

$$P(L|I,\theta) \propto \left(\prod_{i=1}^{n} p(I|l_{i}, u_{i}) \prod_{(v_{i}, v_{j}) \in E} p(l_{i}, l_{j}|c_{ij})\right)$$
  
Appearance likelihood Geometry likelihood

## Part representation

Background subtraction

![](_page_49_Picture_2.jpeg)

![](_page_49_Figure_3.jpeg)

## Pictorial structures model

Optimization is tricky but can be efficient

$$L^* = \arg\min_{L} \left( \sum_{i=1}^n m_i(l_i) + \sum_{(v_i, v_j) \in E} d_{ij}(l_i, l_j) \right) \xrightarrow{(\bullet)}_{\mathbb{Z}_{2} \cup \mathbb{Z}_{2}} v_l$$

 $Best_2(l_1) = \min_{l_2} \left[ m_2(l_2) + d_{12}(l_1, l_2) \right]$ 

- Remove v<sub>2</sub>, and repeat with smaller tree, until only a single part
- For k parts, n locations per part, this has complexity of O(kn<sup>2</sup>), but can be solved in ~O(nk) using generalized distance transform

## **Distance Transform**

 For each pixel p, how far away is the nearest pixel q of set S

$$-f(p) = \min_{q \in G} d(p,q)$$

- G is often the set of edge pixels

![](_page_51_Figure_4.jpeg)

## **Distance Transform - Applications**

- Set distances e.g. Hausdorff Distance
- Image processing e.g. Blurring
- Robotics Motion Planning
- Alignment
  - Edge images
  - Motion tracks
  - Audio warping
- Deformable Part Models

## Generalized Distance Transform

- Original form:  $f(p) = \min_{q \in G} d(p,q)$
- General form:  $f(p) = \min_{q \in [1,N]} m(q) + d(p,q)$
- For many deformation costs,  $O(N^2) \rightarrow O(N)$ Quadratic  $d(p,q) = \alpha(p-q)^2 + \beta(p-q)$ Abs Diff  $d(p,q) = \alpha|p-q|$ Min Composition  $d(p,q) = \min(d_1(p,q), d_2(p,q))$ Bounded  $d_{\tau}(p,q) = \begin{cases} d(p,q) & : |p-q| < \tau \\ \infty & : |p-q| \ge \tau \end{cases}$

## Results for person matching

![](_page_54_Picture_1.jpeg)

## Results for person matching

![](_page_55_Picture_1.jpeg)

## Enhanced pictorial structures

#### EICHNER, FERRARI: BETTER APPEARANCE MODELS FOR PICTORIAL STRUCTURES 9

![](_page_56_Picture_2.jpeg)

BMVC 2009

# **Deformable Latent Parts Model**

#### Useful parts discovered during training

Detections

![](_page_57_Picture_3.jpeg)

![](_page_57_Picture_4.jpeg)

![](_page_57_Picture_5.jpeg)

![](_page_57_Picture_6.jpeg)

Template Visualization

![](_page_57_Picture_8.jpeg)

![](_page_57_Picture_9.jpeg)

![](_page_57_Picture_10.jpeg)

root filters coarse resolution

part filters finer resolution

deformation models

Felzenszwalb et al. 2008

# Things to remember

- Rather than searching for whole object, can locate "parts" that vote for object
  - Better encoding of spatial variation
- These parts can vote for other things too
- Models can be broken down into part appearance and spatial configuration
  - Wide variety of models
- Efficient optimization can be tricky but usually possible

![](_page_58_Picture_7.jpeg)

![](_page_58_Picture_8.jpeg)

## Next classes

- Tues: Object tracking with Kalman Filters
  - Presented by Amin Sadeghi
  - HW 5 is due

- Thurs: Action Recognition
  - Presented by me