# Lecture 10: Object detection

#### COS 429: Computer Vision



#### Last time





- 1. Object detection evaluation
- 2. Deformable parts model
- 3. Challenges in object detection

# Recognition: classification vs detection

#### Classification



- Image-level label
- Doesn't require/assume object position in the image (blessing and curse)
- Frequently relies on context
- Doesn't require counting
- Doesn't require delineating multiple instances

#### Detection



- Box-level label
- Box tight around the object instance (blessing and curse)
- Requires counting and delineating nearby instances
- Requires finding all instances
- May require non-max suppression

#### Annotation costs

#### Draw a tight bounding box around the moped



#### Annotation costs

#### Draw a tight bounding box around the moped



#### This took 14.5 seconds

(7 sec [Jain&Grauman ICCV'13], 10.2 sec [Russakovsky, Li, Fei-Fei CVPR'15], 25.5 sec [Su, Deng, Fei-Fei AAAIW'12])

# Datasets drive computer vision progress



Dataset scale and complexity

#### PASCAL VOC benchmark

The **PASCAL** Visual Object Classes Homepage



#### The PASCAL VOC project:

- · Provides standardised image data sets for object class recognition
- · Provides a common set of tools for accessing the data sets and annotations
- · Enables evaluation and comparison of different methods
- Ran challenges evaluating performance on object class recognition (from 2005-2012, now finished)

#### **Pascal VOC data sets**

Data sets from the VOC challenges are available through the challenge links below, and evalution of new methods on these data sets can be achieved through the <u>PASCAL VOC Evaluation Server</u>. The evaluation server will remain active even though the challenges have now finished.

#### News

- Nov-2014: A new feature for the <u>Leaderboards</u> of the PASCAL VOC evaluation server has been added, indicating if the differences between a selected submission and others are statistically significant or not.
- May-2014: A new paper covering the 2008-12 years of the challenge, and lessons learnt, is now available:

The PASCAL Visual Object Classes Challenge: A Retrospective Everingham, M., Eslami, S. M. A., Van Gool, L., Williams, C. K. I., Winn, J. and Zisserman, A. International Journal of Computer Vision, 111(1), 98-136, 2015 Bibtex source | Abstract | PDF

#### PASCAL VOC benchmark



- *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

## Object detection evaluation: average precision



- give a "car" score to each window

$$Recall = \frac{NumTruePositives}{NumPositives}$$
$$Precision = \frac{NumTruePositives}{NumPredictions}$$



#### Threshold for Correct Detection



# Object detection evaluation

<u>All</u> instances of <u>all</u> target object classes expected to be localized on <u>all</u> test images



- Algorithm outputs a list of bounding box detections with confidences
- A detection is considered correct if overlap with ground truth is big enough
  - duplicate detections are penalized
- Evaluated by average precision (AP) per object class
- Overall evaluated usually by mAP
  - In competitions, also by number of classes won

Everingham, Van Gool, Williams, Winn and Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# Object detection is a collection of problems

#### Intra-class Variation for "Airplane"

Occlusion









Shape











Viewpoint



+









Credit: Derek Hoiem

# Object detection is a collection of problems

#### **Confusing Distractors for "Airplane"**

#### Background











Similar Categories











Dissimilar Categories











Localization Error











Credit: Derek Hoiem

## PASCAL VOC challenge



- 20 categories
- Annual (2005-2012) classification, detection, segmentation, ... challenges

#### Deformable Parts Model

## Machine learning for object detection

What features do we use?
 – intensity, color, gradient information, BOW, …

 Which machine learning methods?
 – k-nearest neighbors, boosting, least squares, SVMs, …

• What hacks do we need to get things working?

## Review: Person detection via Dalal Triggs

1. Represent each example with a single, fixed HoG template



2. Learn a single linear detector



N. Dalal and B. Triggs, <u>Histograms of Oriented Gradients for Human Detection</u>, CVPR 2005

Credit: Noah Snavely

#### Sliding window detection



- Compute HOG of the whole image at multiple resolutions
- Score every subwindow of the feature pyramid
- Apply non-maxima suppression

#### Detection



number of locations *p* ~ 250,000 per image test set has ~ 5000 images >> 1.3x10<sup>9</sup> windows to classify typically only ~ 1,000 true positive locations

Extremely unbalanced binary classification

#### Dalal&Triggs on PASCAL VOC 2007



#### Part-based models

- Parts local appearance templates
- "Springs" spatial connections between parts (geom. prior)



Image: [Felzenszwalb and Huttenlocher 05]

#### Part-based models

- Local appearance is easier to model than the global appearance
  - Training data shared across deformations
  - "part" can be local or global depending on resolution
- Generalizes to previously unseen configurations



#### Part configuration score function



## Part configuration score function

- Objective: maximize score over p<sub>1</sub>,...,p<sub>n</sub>
- $h^n$  configurations! (h = |P|, about 250,000)
- Dynamic programming

Object Detection with Discriminatively Trained Part Based Models Pedro F. Felzenszwalb, Ross B. Girshick, David McAllester and Deva Ramanan. PAMI 2010

#### Star-structured deformable part models



#### Dalal-Triggs + parts





- - Add parts to the Dalal & Triggs detector
    - HOG features
    - Linear filters / sliding-window detector
    - Discriminative training

#### Sliding window DPM score function



#### DPM detection in a slide



# What are the parts?





# Clustering by viewpoints (aspect ratios as proxy)



General philosophy: enrich models to better represent the data

## DPM with mixture models

Data driven: aspect, occlusion modes, subclasses





<u>Person detection</u> Without parts: AP = 0.12 Parts but no mixtures: AP = 0.27 Parts + mixtures: AP = 0.36

## Discriminatively trained deformable parts model

#### Discriminatively trained deformable part models

Version 5 (Sept. 5, 2012)





#### Introduction

Over the past few years we have developed a complete learning-based system for detecting and localizing objects in images. Our system represents objects using mixtures of deformable part models. These models are trained using a discriminative method that only requires bounding boxes for the objects in an image. The approach leads to efficient object detectors that achieve state of the art results on the PASCAL and INRIA person datasets.

At a high level our system can be characterized by the combination of

- 1. Strong low-level features based on histograms of oriented gradients (HOG)
- 2. Efficient matching algorithms for deformable part-based models (pictorial structures)
- 3. Discriminative learning with latent variables (latent SVM)

This work was awarded the PASCAL VOC "Lifetime Achievement" Prize in 2010.

Code: <u>http://www.rossgirshick.info/latent/</u>

Slides: <u>http://vision.stanford.edu/teaching/cs231b\_spring1213/slides/dpm-slides-ross-girshick.pdf</u> Paper: Object Detection with Discriminatively Trained Part Based Models Pedro F. Felzenszwalb, Ross B. Girshick, David McAllester and Deva Ramanan. PAMI 2010

## Where would DPM succeed and fail?

- *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

aero	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
33.2	60.3	10.2	16.1	27.3	54.3	58.2	23.0	20.0	24.1	26.7	12.7	58.1	48.2	43.2	12.0	21.1	36.1	46.0	43.5	33.7

#### How would DPM compare to SIFT+SPM?

#### (not a perfect comparison by any means, but an attempt)

	aero										dining			motor		potted				tv/
	plane	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	bike	person	plant	sheep	sofa	train 1	monitor
UOC_OXFORD_DPM_MKL	59.6	54.5	21.9	21.6	32.1	52.5	<mark>49.3</mark>	40.8	19.1	<mark>35.2</mark>	28.9	37.2	50.9	49.9	46.1	15.6	39.3	35.6	48.9	42.8
NEC STANFORD OCP	65.1	46.8	25.0	24.6	16.0	51.0	44.9	51.5	13.0	26.6	31.0	40.2	39.7	51.5	32.8	12.6	35.7	33.5	48.0	44.8

#### UOC\_OXFORD\_DPM\_MKL

This method is similar to last year DPM-MKL entry. We updated several aspects of the implementation (e.g. the type of features).

#### NEC\_STANFORD\_OCP

Object-centric pooling (OCP) is a method which represents a bounding box by pooling the coded low-level descriptors on the foreground and background separately and then concatenating them (Russakovsky et al. ECCV 2012). This method exploits powerful classification features that have been developed in the past years. In this system, we used DHOG and LBP as low-level descriptors. We developed a discriminative LCC coding scheme in addition to traditional LCC coding. We make use of candidate bounding boxes (van de Sande et al. ICCV 2011).

http://host.robots.ox.ac.uk/pascal/VOC/voc2012/results/index.html



# Average precision evaluation is not enough

- Average Precision (AP)
  - Good summary statistic for quick comparison
  - Not a good driver of research

	aero	bike	bird	boat	bottle	bus
a) base	.290	.546	.006	.134	.262	.394
b) BB	.287	.551	.006	.145	.265	.397
c) context	.328	.568	.025	.168	.285	.397

Typical evaluation through comparison of AP numbers

- Need tools to determine
  - where detectors fail
  - potential impact of particular improvements



#### Tool for object detection analysis

#### **Diagnosing Error in Object Detectors**

<u>Derek Hoiem</u> and Qieyun Dai and Yodsawalai Chodpathumwan <u>Computer Vision Group</u> <u>Department of Computer Science</u> <u>University of Illinois at Urbana-Champaign</u>

#### **Overview**

This work provides a set of tools for analyzing object detector performance.

**Note:** (11/12/14) The summary plots (e.g., "animal" or "vehicle") for displayDetectionTrend were computed incorrectly. The revised code is now in the .tar.gz file, but the pdfs have not been updated. Thanks to Shaoqing Ren for noticing the bug and providing the fix. Another method displayDetectionTrend2.m is also provided, which averages across tic marks to summarize several categories.

#### **Downloads**

The following resources are available:

- An updated version (v2) of the code/annotations: [src/data (84.5MB)]
- Description of updates: [pdf]
- Examples of automatic analysis reports: [dpm v4.pdf] [vedaldi2009.pdf] [cnn7 bb.pdf]
- Original version (v1) of the code/annotations (in case you have trouble with the new version): [src/data (69MB)]

#### **Publications**

#### **Diagnosing Error in Object Detectors**

Derek Hoiem, Yodsawalai Chodpathumwan, and Qieyun Dai *ECCV*, 2012. [pdf] [slides]

http://dhoiem.web.engr.illinois.edu/projects/detectionAnalysis/

## Analysis of object characteristics



Level of occlusion: 2 (moderate) Parts visible: bike body, handlebars, wheel Parts not visible: seat View: side visible (front, top, etc., not visible) Area: 3233 pixels Aspect Ratio (w/h): 1.24



http://dhoiem.web.engr.illinois.edu/projects/detectionAnalysis/ Credit: Derek Hoiem

## Top false positives: Airplane (DPM)



# Top false positives: Dog (DPM)



#### Summary of False Positive Analysis



http://dhoiem.web.engr.illinois.edu/projects/detectionAnalysis/ Credit: Derek Hoiem

# Challenging false positives



**Incorrect Localization** 



Right



Wrong



Wrong





Dog Model

Challenges of Confusion with Similar Categories

http://dhoiem.web.engr.illinois.edu/projects/detectionAnalysis/ Credit: Derek Hoiem

## ImageNet detection challenge

# Datasets drive computer vision progress



Dataset scale and complexity

# IMAGEN

#### **22K** categories and **15M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate

- - Tree
    Artifact
  - Flower
- Food
- Materials
- Plants
   Structures
  - - Tools
  - Appliances
    - Structures

- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activity

#### www.image-net.org

Deng et al. 2009, Russakovsky et al. 2015

#### ILSVRC object detection task

#### Allows evaluation of generic object detection in cluttered scenes at scale



Person Car Motorcycle Helmet

#### ILSVRC object detection data









# ILSVRC object detection data

#### **Comparative scale**

		PASCAL VOC 2012	ILSVRC 2014			
Number of	object classes	20	200			
Training	Num images	5717	456567			
Training	Num objects	13609	478807			
Validation	Num images	5823	20121			
validation	Num objects	13841	55502			
Testing	Num images	10991	40152			
resurig	Num objects					

#### Comparative statistics (on validation set)

	PASCAL VOC 2012	ILSVRC 2013
Average image resolution	469x387 pixels	482x415 pixels
Average object classes per image	1.521	1.534
Average object instances per image	2.711	2.758
Average object scale (bounding box area as fraction of image area)	0.207	0.170

# ImageNet challenge: participation and performance



## Easiest and hardest categories



Olga Russakovsky\*, Jia Deng\*, et al. ImageNet Large Scale Visual Recognition Challenge. IJCV, 2015.

## What are the remaining challenges?





1

Number of concepts

#### PASCAL VOC

[Everingham et al. IJCV10]



20

#### ImageNet challenge

[Russakovsky et al. IJCV15]



200

## Impact of object scale on detection accuracy

(ImageNet challenge 2013-2015 winning object detection entries)



## Impact of object scale on detection accuracy

(ImageNet challenge 2013-2015 winning object detection entries)



## Impact of object scale on detection accuracy

(ImageNet challenge 2013-2015 winning object detection entries)



#### ILSVRC data, challenge, algorithms

http://image-net.org/challenges/LSVRC/

 Olga Russakovsky\*, Jia Deng\*, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg and Li Fei-Fei. (\* = equal contribution) ImageNet Large Scale Visual Recognition Challenge. IJCV, 2015.
 https://arxiv.org/abs/1409.0575

We'll come back to this in the deep learning section of the course

## Summary: key concepts

- 1. Object detection evaluation
- 2. Deformable parts model
- 3. Challenges in object detection

# Next class: segmentation

