# Computer Vision – Lecture 15

## Part-based Models for Object Categorization

07.01.2016

Bastian Leibe
RWTH Aachen
http://www.vision.rwth-aachen.de

leibe@vision.rwth-aachen.de

Computer Vision WS 15/16

---

## Course Outline

- **Image Processing Basics**
- **Segmentation & Grouping**
- **Object Recognition**
- **Object Categorization I**
  - ➢ Sliding Window based Object Detection
- **Local Features & Matching**
  - ➢ Local Features – Detection and Description
  - ➢ Recognition with Local Features
  - ➢ Indexing & Visual Vocabularies
- **Object Categorization II**
  - ➢ Bag-of-Words Approaches & Part-based Approaches
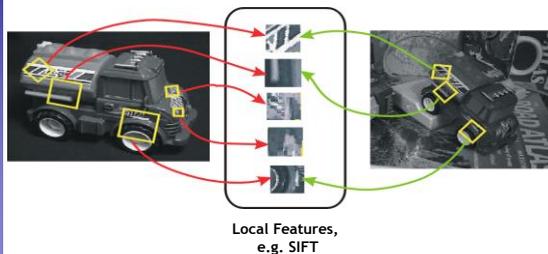  - ➢ Deep Learning Methods
- **3D Reconstruction**

Computer Vision WS 15/16

3

---

## Topics of This Lecture

- **Recap: Specific Object Recognition with Local Features**
  - ➢ Matching & Indexing
  - ➢ Geometric Verification

- **Part-Based Models for Object Categorization**
  - ➢ Structure representations
  - ➢ Different connectivity structures

- **Bag-of-Words Model**
  - ➢ Use for image classification

- **Implicit Shape Model**
  - ➢ Generalized Hough Transform for object category detection

- **Deformable Part-based Model**
  - ➢ Discriminative part-based detection

Computer Vision WS 15/16

B. Leibe

4

---

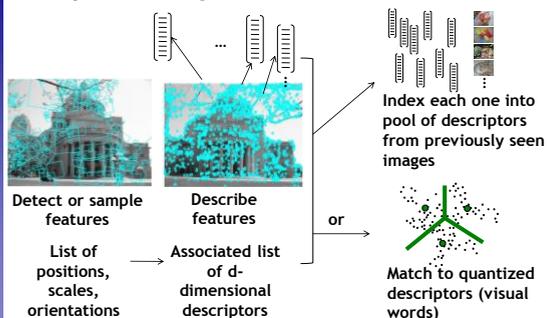## Recap: Recognition with Local Features

- **Image content is transformed into local features that are invariant to translation, rotation, and scale**
- **Goal: Verify if they belong to a consistent configuration**



Local Features,
e.g. SIFT

Computer Vision WS 15/16

Slide credit: David Lowe

B. Leibe

5

---

## Recap: Indexing features



Index each one into
pool of descriptors
from previously seen
images

Detect or sample
features

Describe
features

or

List of
positions,
scales,
orientations

Associated list
of d-
dimensional
descriptors

Match to quantized
descriptors (visual
words)

⇒ *Shortlist of possibly matching images + feature correspondences*

Computer Vision WS 15/16

Slide credit: Kristen Grauman

B. Leibe

6

---

## Extension: *tf-idf* Weighting

- **Term frequency – inverse document frequency**
  - ➢ Describe frame by frequency of each word within it, downweight words that appear often in the database
  - ➢ (Standard weighting for text retrieval)

Number of
occurrences of word
$i$ in document $d$

Total number of
documents in
database

$$t_i = \frac{n_{id}}{n_d} \log \frac{N}{n_i}$$

Number of words in
document $d$

Number of
occurrences of word $i$
in whole database

Computer Vision WS 15/16

Slide credit: Kristen Grauman

B. Leibe

7

## Recap: Fast Indexing with Vocabulary Trees

• **Recognition**

Geometric verification

[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Computer Vision WS 15/16

B. Leibe

8

---

## Recap: Geometric Verification by Alignment

• **Assumption**
  ➢ Known object, rigid transformation compared to model image
  ⇒ *If we can find evidence for such a transformation, we have recognized the object.*

• **You learned methods for**
  ➢ Fitting an *affine transformation* from ≥ 3 correspondences
  ➢ Fitting a *homography* from ≥ 4 correspondences

  Affine: solve a system          Homography: solve a system
  $$At = b$$                        $$Ah = 0$$

• **Correspondences may be noisy and may contain outliers**
  ⇒ Need to use robust methods that can filter out outliers
  ⇒ Use **RANSAC** or the **Generalized Hough Transform**

Computer Vision WS 15/16

B. Leibe

9

---

## Topics of This Lecture

• Recap: Specific Object Recognition with Local Features

• **Part-Based Models for Object Categorization**
  ➢ **Structure representations**
  ➢ **Different connectivity structures**

• Bag-of-Words Model
  ➢ Use for image classification

• Implicit Shape Model
  ➢ Generalized Hough Transform for object category detection

• Deformable Part-based Model
  ➢ Discriminative part-based detection

Computer Vision WS 15/16

B. Leibe

10

---

## Recognition of Object Categories

• **We no longer have exact correspondences…**

• **On a local level, we can still detect similar parts.**

• **Represent objects by their parts**
  ⇒ Bag-of-features

• **How can we improve on this?**
  ➢ Encode structure

Slide credit: Rob Fergus

Computer Vision WS 15/16

11

---

## Part-Based Models

• **Fischler & Elschlager 1973**

• **Model has two components**
  ➢ parts (2D image fragments)
  ➢ structure (configuration of parts)

HAIR
EYE   EYE
LEFT EDGE   RIGHT EDGE
NOSE
MOUTH

Computer Vision WS 15/16

B. Leibe

12

---

## Different Connectivity Structures

$\mathcal{O}(N)$        $\mathcal{O}(N^k)$        $\mathcal{O}(N^2)$        $\mathcal{O}(N^2)$

a) Bag of visual words        b) Constellation        c) Star shape        d) Tree
Csurka et al. '04             Fergus et al. '03       Leibe et al. '04, '08   Felzenszwalb &
Vasconcelos et al. '00        Fei-Fei et al. '03      Crandall et al. '05     Huttenlocher '05
                                                      Fergus et al. '05

$\mathcal{O}(N^3)$

e) k-fan (k = 2)        f) Hierarchy              g) Sparse flexible model
Crandall et al. '05     Bouchard & Triggs '05     Carneiro & Lowe '06

Slide adapted from Rob Fergus        B. Leibe        Image from [Carneiro & Lowe, ECCV'06]

Computer Vision WS 15/16

13

---

2

## Topics of This Lecture

- Recap: Specific Object Recognition with Local Features
- Part-Based Models for Object Categorization
  - Structure representations
  - Different connectivity structures
- **Bag-of-Words Model**
  - **Use for image classification**
- Implicit Shape Model
  - Generalized Hough Transform for object category detection
- Deformable Part-based Model
  - Discriminative part-based detection

---

## Recap: Analogy to Documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach the brain from our eyes. For a long time it was thought that the retinal image was transmitted point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we now know that behind the origin of the visual perception in the brain there is a considerably more complicated course of events. By following the visual impulses along their path to the various cell layers of the optical cortex, Hubel and Wiesel have been able to demonstrate that the *message about the image falling on the retina undergoes a step-wise analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.*

**sensory, brain, visual, perception, retinal, cerebral cortex, eye, cell, optical nerve, image Hubel, Wiesel**

China is forecasting a trade surplus of $90bn (£51bn) to $100bn this year, a threefold increase on 2004's $32bn. The Commerce Ministry said the surplus would be created by a predicted 30% jump in exports to $750bn, compared with a 18% rise in imports to $660bn. The figures are likely to further annoy the US, which has long argued that China's exports are unfair in the US. China's trade surplus, have pledged to measures. Beijing agrees the surplus is too high and now needs domestic demand so more goods stay within the country. China increased the value of the yuan against the dollar by 2.1% in July and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**

---

## Recap: Visual Words

- **Quantize the feature space into "visual words"**
- **Perform matching only to those visual words.**

**Exact feature matching → Match to same visual word**

Figure from Sivic & Zisserman, ICCV 2003

---

## Recap: Bag-of-Word Representations (BoW)

**Object** → **Bag of "words"**

Source: ICCV 2005 short course, Li Fei-Fei

---

## Recap: Categorization with Bags-of-Words

- **Compute the word activation histogram for each image.**
- **Let each such BoW histogram be a feature vector.**
- **Use images from each class to train a classifier (e.g., an SVM).**

**Violins**

B. Leibe · 18

---

## Recap: Advantage of BoW Histograms

- **Bag of words representations make it possible to describe the unordered point set with a single vector (of fixed dimension across image examples).**

- **Provides easy way to use distribution of feature types with various learning algorithms requiring vector input.**

B. Leibe · 19

## Limitations of BoW Representations

- The bag of words removes spatial layout.

- This is both a strength and a weakness.

- *Why a strength?*

- *Why a weakness?*



A  B  C  D

Slide adapted from Bill Freeman
B. Leibe
20

---

## Spatial Pyramid Representation

- Representation in-between orderless BoW and global appearance



Slide credit: Svetlana Lazebnik
B. Leibe
[Lazebnik, Schmid & Ponce, CVPR'06]
21

---

## Spatial Pyramid Representation

- Representation in-between orderless BoW and global appearance



Slide credit: Svetlana Lazebnik
B. Leibe
[Lazebnik, Schmid & Ponce, CVPR'06]
22

---

## Spatial Pyramid Representation

- Representation in-between orderless BoW and global appearance



Slide credit: Svetlana Lazebnik
B. Leibe
[Lazebnik, Schmid & Ponce, CVPR'06]
23

---

## Summary: Bag-of-Words

- **Pros:**
  - Flexible to geometry / deformations / viewpoint
  - Compact summary of image content
  - Provides vector representation for sets
  - Empirically good recognition results in practice

- **Cons:**
  - Basic model ignores geometry – must verify afterwards, or encode via features.
  - Background and foreground mixed when bag covers whole image
  - When using interest points or sampling: no guarantee to capture object-level parts ⇒ Dense sampling is often better.

Slide credit: Kristen Grauman
B. Leibe
24

---

## Topics of This Lecture

- Recap: Specific Object Recognition with Local Features

- Part-Based Models for Object Categorization
  - Structure representations
  - Different connectivity structures

- Bag-of-Words Model
  - Use for image classification

- **Implicit Shape Model**
  - Generalized Hough Transform for object category detection

- Deformable Part-based Model
  - Discriminative part-based detection

B. Leibe
25

4

## Implicit Shape Model (ISM)

- **Basic ideas**
  - Learn an appearance codebook
  - Learn a star-topology structural model
    - Features are considered independent given obj. center

- **Algorithm: probabilistic Gen. Hough Transform**
  - Exact correspondences → Prob. match to object part
  - NN matching → Soft matching
  - Feature location on obj. → Part location distribution
  - Uniform votes → Probabilistic vote weighting
  - Quantized Hough array → Continuous Hough space

Computer Vision WS 15/16

B. Leibe

26

---

## Implicit Shape Model: Basic Idea

- **Visual vocabulary is used to index votes for object position [a visual word = "part"].**

**Training image**

**Visual codeword with displacement vectors**

B. Leibe, A. Leonardis, and B. Schiele, Robust Object Detection with Interleaved Categorization and Segmentation, International Journal of Computer Vision, Vol. 77(1-3), 2008.

Computer Vision WS 15/16

B. Leibe

---

## Implicit Shape Model: Basic Idea

- **Objects are detected as consistent configurations of the observed parts (visual words).**

**Test image**

B. Leibe, A. Leonardis, and B. Schiele, Robust Object Detection with Interleaved Categorization and Segmentation, International Journal of Computer Vision, Vol. 77(1-3), 2008.

Computer Vision WS 15/16

B. Leibe

---

## Implicit Shape Model - Representation

**Training images (+reference segmentation)**

**Appearance codebook**

- **Learn appearance codebook**
  - Extract local features at interest points
  - Agglomerative clustering ⇒ codebook

- **Learn spatial distributions**
  - Match codebook to training images
  - Record matching positions on object

**Spatial occurrence distributions**

Computer Vision WS 15/16

B. Leibe

29

---

## Implicit Shape Model - Recognition

**Interest Points**   **Matched Codebook Entries**   **Probabilistic Voting**

**Image Feature**   **Interpretation (Codebook match)**   **Object Position**

$f$   $C_i$   $o,x$

$p(C_i|f)$   $p(o_n, x | C_i, \ell)$

**Probabilistic vote weighting**

**3D Voting Space (continuous)**

Computer Vision WS 15/16

30

[Leibe, Leonardis, Schiele, SLCV'04; IJCV'08]

---

## Implicit Shape Model - Recognition

**Interest Points**   **Matched Codebook Entries**   **Probabilistic Voting**

**3D Voting Space (continuous)**

**Backprojected Hypotheses**   **Backprojection of Maxima**

Computer Vision WS 15/16

31

[Leibe, Leonardis, Schiele, SLCV'04; IJCV'08]

**Example: Results on Cows**
Original image


**Example: Results on Cows**
Interest points


**Example: Results on Cows**
Matched patches


**Example: Results on Cows**
Prob. Votes


**Example: Results on Cows**
1st hypothesis


**Example: Results on Cows**
2nd hypothesis

## Example: Results on Cows



**3rd hypothesis**

## Scale Invariant Voting

- **Scale-invariant feature selection**
  - Scale-invariant interest regions
  - Extract scale-invariant descriptors
  - Match to appearance codebook

- **Generate scale votes**
  - Scale as 3rd dimension in voting space

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

  - Search for maxima in 3D voting space



Search window

## Detection Results

- **Qualitative Performance**
  - Recognizes different kinds of objects
  - Robust to clutter, occlusion, noise, low contrast

## Detections Using Ground Plane Constraints



Battery of 5
ISM detectors
for different
car views

left camera
1175 frames

## Extension: Rotation-Invariant Detection

- **Polar instead of Cartesian voting scheme**



- **Benefits:**
  - Recognize objects under image-plane rotations
  - Possibility to share parts between articulations.

- **Caveats:**
  - Rotation invariance should only be used when it's really needed. (Also increases false positive detections)

## Sometimes, Rotation Invariance Is Needed…



Figure from [Mikolajczyk et al., CVPR'06]

## Implicit Shape Model – Segmentation

**Local Features**  **Matched Codebook Entries**  **Probabilistic Voting**

Backproject Meta-information

y

x

**3D Voting Space (continuous)**

**Segmentation**

**Pixel Contributions**  **Backprojected Hypotheses**  **Backprojection of Maxima**

Computer Vision WS 15/16

46

[Leibe, Leonardis, Schiele, DAGM'04; IJCV'08]

---

## Example Results: Motorbikes

Computer Vision WS 15/16

B. Leibe

47

[Leibe, Leonardis, Schiele, SLCV'04; IJCV'08]

---

## You Can Try It At Home…

- **Linux source code & binaries available**
  - Including datasets & several pre-trained detectors
  - http://www.vision.rwth-aachen.de/software

Computer Vision WS 15/16

B. Leibe

48

---

## Topics of This Lecture

- Recap: Specific Object Recognition with Local Features
- Part-Based Models for Object Categorization
  - Structure representations
  - Different connectivity structures
- Bag-of-Words Model
  - Use for image classification
- Implicit Shape Model
  - Generalized Hough Transform for object category detection
- **Deformable Part-based Model**
  - **Discriminative part-based detection**

Computer Vision WS 15/16

B. Leibe

49

---

## Starting Point: HOG Sliding-Window Detector

$p$

**Filter $F$**

**Score of $F$ at position $p$ is**
$$F \cdot \phi(p,H)$$

$\phi(p,H)$ = concatenation of HOG features from window specified by $p$.

**HOG pyramid $H$**

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

Computer Vision WS 15/16

B. Leibe

50

---

## Deformable Part-based Models

- **Mixture of deformable part models (pictorial structures)**
- **Each component has global template + deformable parts**
- **Fully trained from bounding boxes alone**

Computer Vision WS 15/16

Slide credit: Pedro Felzenszwalb

B. Leibe

51

---

**2-Component Bicycle Model**

Root filters
coarse resolution

Part filters
finer resolution

Deformation
models

B. Leibe
52

---

**Object Hypothesis**



Score of filter:
dot product of filter
with HOG features
underneath it

Score of object
hypothesis is sum of
filter scores minus
deformation costs

Image pyramid

HOG feature pyramid

• **Multiscale model captures features at two resolutions**

B. Leibe
53

---

**Score of a Hypothesis**

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

"data term"        "spatial prior"

filters            displacements
                   deformation parameters

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

concatenation filters and        concatenation of HOG
deformation parameters           features and part
                                 displacement features

B. Leibe
54

---

**Recognition Model**



$$f_w(x) = w \cdot \Phi(x)$$

$$f_w(x) = \max_z w \cdot \Phi(x, z)$$

• $z$ : **vector of part offsets**
• $\Phi(x, z)$ : **vector of HOG features (from root filter & appropriate part sub-windows) and part offsets**

B. Leibe
55

---

**Results: Persons**



• **Results (after non-maximum suppression)**
  ➢ **~1s to search all scales**

B. Leibe
56

---

**Results: Bicycles**

B. Leibe
57

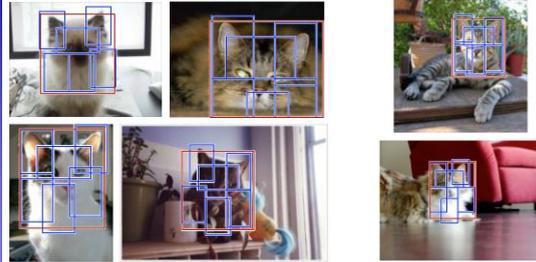9

## False Positives

• **Bicycles**

B. Leibe

58

## Results: Cats



**High-scoring true positives**          **High-scoring false positives (not enough overlap)**

Slide credit: Pedro Felzenszwalb

59

## You Can Try It At Home…

• **Deformable part-based models have been very successful at several recent evaluations.**

⇒ **Currently, state-of-the-art approach in object detection**

• **Source code and models trained on PASCAL 2006, 2007, and 2008 data are available here:**

**http://www.cs.uchicago.edu/~pff/latent**

B. Leibe

60

## References and Further Reading

• **Details about the ISM approach can be found in**
  ➢ *B. Leibe, A. Leonardis, and B. Schiele,* Robust Object Detection with Interleaved Categorization and Segmentation, *International Journal of Computer Vision, Vol. 77(1-3), 2008.*

• **Details about the DPMs can be found in**
  ➢ *P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan,* Object Detection with Discriminatively Trained Part Based Models, *IEEE Trans. PAMI, Vol. 32(9), 2010.*

• **Try the ISM Linux binaries**
  ➢ http://www.vision.ee.ethz.ch/bleibe/code

• **Try the Deformable Part-based Models**
  ➢ http://www.cs.uchicago.edu/~pff/latent