# Computer Vision – Lecture 16

## Part-based Models for Object Categorization

### 08.01.2015

**Bastian Leibe**

**RWTH Aachen**

http://www.vision.rwth-aachen.de

leibe@vision.rwth-aachen.de

# 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**
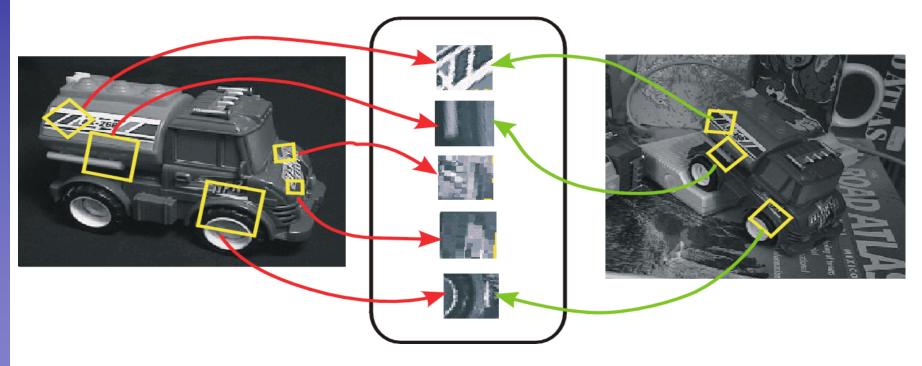- **3D Reconstruction**
- **Optical Flow**

# 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

B. Leibe

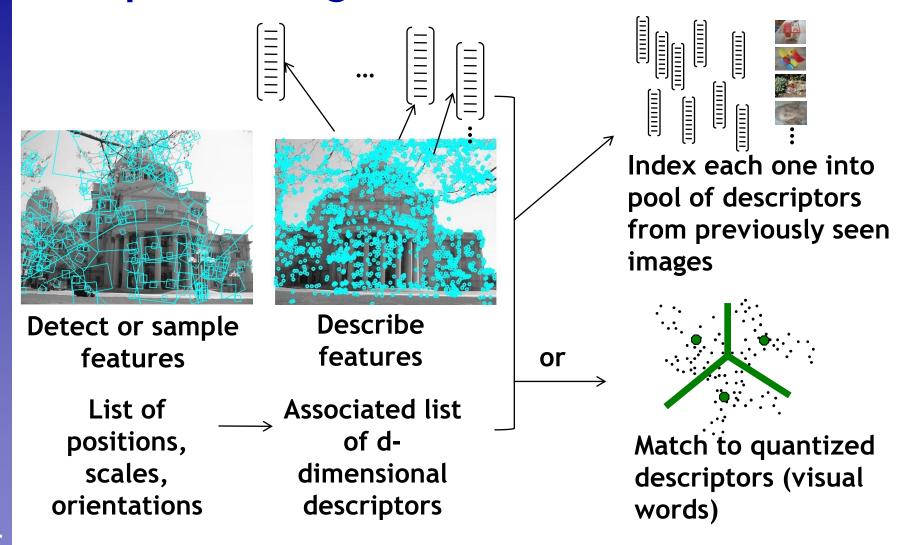# Recap: Recognition with Local Features

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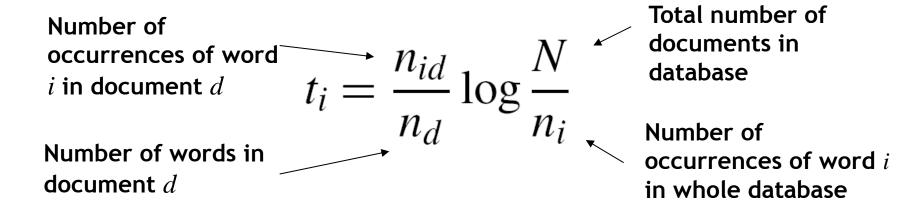


**Local Features,
e.g. SIFT**

B. Leibe

# Recap: Indexing features

**Detect or sample features**

**Describe features**

**Index each one into pool of descriptors from previously seen images**

List of positions, scales, orientations

→

Associated list of d-dimensional descriptors

**or**

**Match to quantized descriptors (visual words)**

⇒ *Shortlist of possibly matching images + feature correspondences*

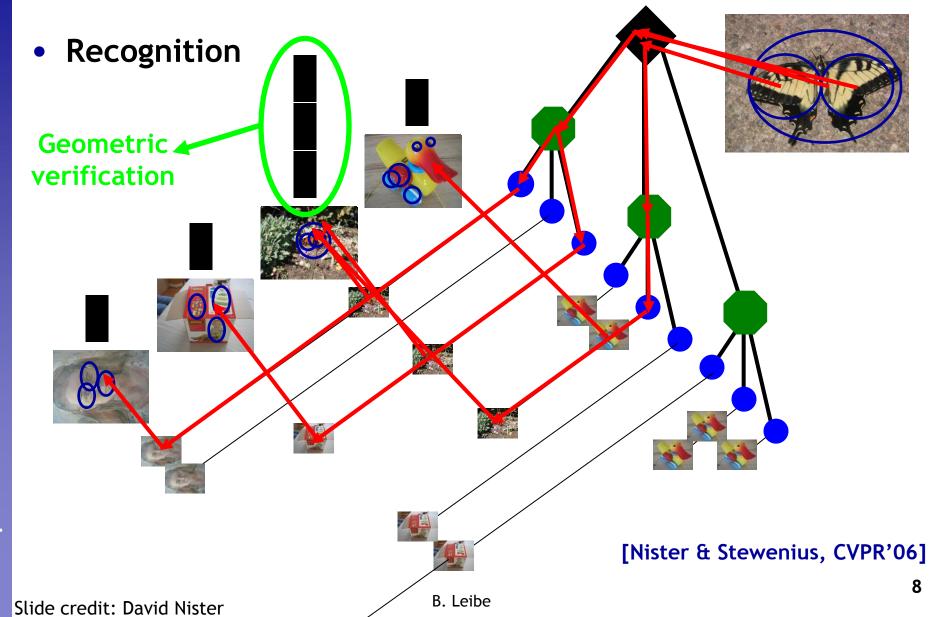Slide credit: Kristen Grauman

B. Leibe

# 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

Slide credit: Kristen Grauman

B. Leibe

# Recap: Fast Indexing with Vocabulary Trees

- **Recognition**

Geometric verification

Slide credit: David Nister

B. Leibe

8

# Application for Content Based Img Retrieval

- **What if query of interest is a portion of a frame?**



Visually defined query

"Groundhog Day" [Rammis, 1993]

"Find this clock"

"Find this place"

Slide credit: Andrew Zisserman

B. Leibe

[Sivic & Zisserman, ICCV'03]

9

# Video Google System

1. **Collect all words within query region**
2. **Inverted file index to find relevant frames**
3. **Compare word counts**
4. **Spatial verification**

**Sivic & Zisserman, ICCV 2003**

- **Demo online at :**
  http://www.robots.ox.ac.uk/~vgg/research/vgoogle/index.html



**Query region**

**Retrieved frames**

Slide credit: Kristen Grauman

B. Leibe

10

# Collecting Words Within a Query Region

- **Example: Friends**



**Query region:**
**pull out only the SIFT**
**descriptors whose**
**positions are within the**
**polygon**

B. Leibe

11

# Example Results



**Query**

B. Leibe

12

# More Results



**Query**



**Retrieved shots**

13

Slide credit: Kristen Grauman

B. Leibe

# 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

$$At = b \qquad\qquad\qquad Ah = 0$$

- **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**

B. Leibe

# Applications: Aachen Tourist Guide

B. Leibe

# Applications: Fast Image Registration



B. Leibe

16

# Applications: Mobile Augmented Reality



**D. Wagner, G. Reitmayr, A. Mulloni, T. Drummond, D. Schmalstieg,**
**Pose Tracking from Natural Features on Mobile Phones**. In *ISMAR 2008*.

B. Leibe

# 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

# 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**

Computer Vision WS 14/15

# Part-Based Models

- **Fischler & Elschlager 1973**

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

B. Leibe

# Different Connectivity Structures



$\mathcal{O}(N)$

a) Bag of visual words

Csurka et al. '04
Vasconcelos et al. '00

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

b) Constellation

Fergus et al. '03
Fei-Fei et al. '03

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

c) Star shape

Leibe et al. '04, '08
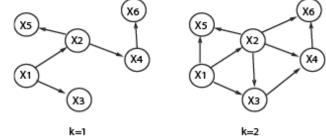Crandall et al. '05
Fergus et al. '05

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

d) Tree

Felzenszwalb &
Huttenlocher '05

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

e) k-fan (k = 2)

Crandall et al. '05

Center
Part
Subpart

f) Hierarchy

Bouchard & Triggs '05

k=1    k=2

g) Sparse flexible model

Carneiro & Lowe '06

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

# 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

# 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 ~~~~ brain from our eyes. For ~~~~ ought that the re~~~~ point by ~~~~ brain; t~~~~ screen ~~~~ in the ~~~~ discov~~~~ know th~~~~ perceptio~~~~ considerab~~~~ of events. By fo~~~~ ses along their path ~~~~ vers of the optical cortex, Hubel and ~~~~ have been able to demonstrate tha ~~~~ *message about the image falling o~~~~ retina undergoes a step-wise analysi~~~~ system of nerve cells stored in colum~~~~ 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 create~~~~ 30% jump in exports t~~~~ th a 18% rise i~~~~ ures are likel~~~~ has lo~~~~ unfair~~~~ under~~~~ surplu~~~~ only on~~~~ Zhou Xia~~~~ needed to ~~~~ demand so m~~~~ the country. China ino~~~~ the yuan against the dollar by 2.1% ~~~~ and permitted it to trade within a ~~~~ band, but the US wants the yuan to ~~~~ allowed to trade freely. However, Beij~~~~ has made it clear that it will take its tin~~~~ and tread carefully before allowing the yuan to rise further in value.
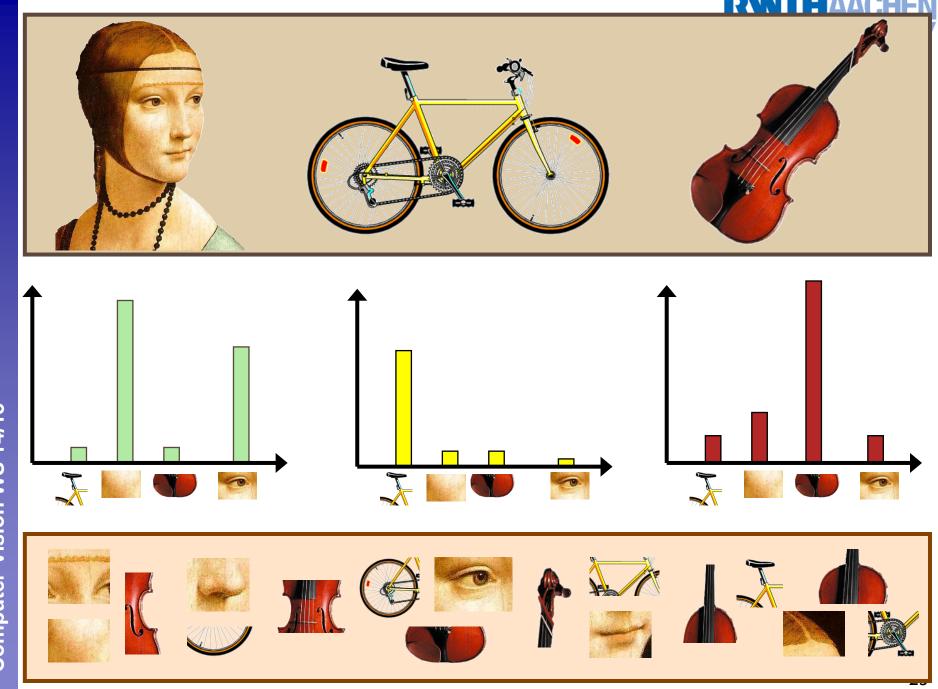
**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**
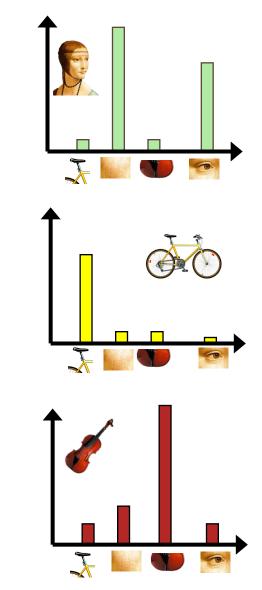
**Object** → **Bag of 'words'**

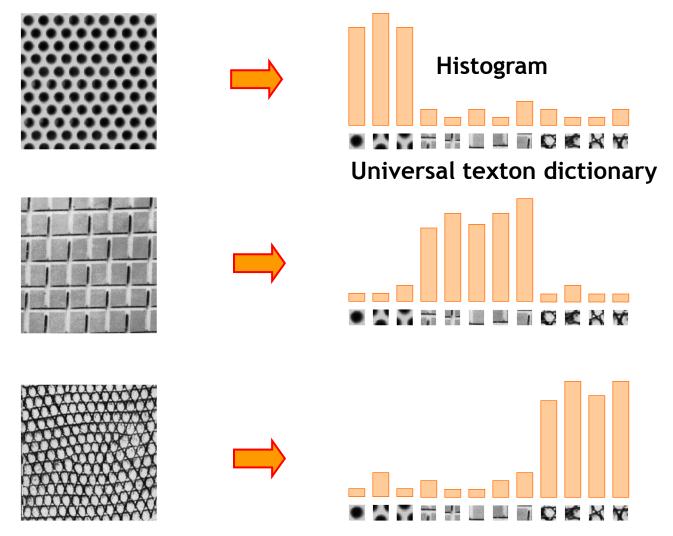Source: ICCV 2005 short course, Li Fei-Fei

# Bags of Visual Words



- **Summarize entire image based on its distribution (histogram) of word occurrences.**

- **Analogous to bag of words representation commonly used for documents.**

- **Main difference to text retrieval: visual words are not given a priori, but obtained through clustering (e.g., using k-means)**

B. Leibe

Computer Vision WS 14/15

# Similarly, Bags-of-Textons for Texture Repr.



Histogram

Universal texton dictionary

**Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003**

Slide credit: Svetlana Lazebnik

# Comparing Bags of Words

- We build up histograms of word activations, so any histogram comparison measure can be used here.

- E.g. we can rank frames by normalized scalar product between their (possibly weighted) occurrence counts

  ➢ *Nearest neighbor* search for similar images.

[1  8  1  4]'  •  [5  1  1  0]

$$sim(d_j, q) = \frac{\vec{d_j} \bullet \vec{q}}{|\vec{d_j}| \times |\vec{q}|}$$

$$= \frac{\sum_{i=1}^{t} w_{i,j} \times w_{i,q}}{\sqrt{\sum_{i=1}^{t} w_{i,j}^2} \times \sqrt{\sum_{j=1}^{t} w_{i,q}^2}}$$

$\vec{d}_j$          $\vec{q}$

B. Leibe

28

Slide credit: Kristen Grauman

# Learning/Recognition with BoW Histograms

- **Bag of words representation makes 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.**
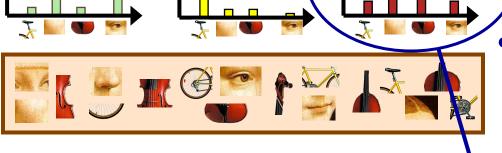
B. Leibe

29

# 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

Slide adapted from Kristen Grauman

B. Leibe

# BoW for Object Categorization



{face, flowers, building}

- **Works pretty well for image-level classification**

Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)

B. Leibe

Computer Vision WS 14/15

# BoW for Object Categorization

## Caltech6 dataset



| class | bag of features | bag of features | Parts-and-shape model |
|---|---|---|---|
| | Zhang et al. (2005) | Willamowski et al. (2004) | Fergus et al. (2003) |
| airplanes | **98.8** | 97.1 | 90.2 |
| cars (rear) | 98.3 | **98.6** | 90.3 |
| cars (side) | **95.0** | 87.3 | 88.5 |
| faces | **100** | 99.3 | 96.4 |
| motorbikes | **98.5** | 98.0 | 92.5 |
| spotted cats | **97.0** | — | 90.0 |

- **Good performance for pure classification (object present/absent)**
  - **Better than more elaborate part-based models with spatial constraints...**
  - **What could be possible reasons why?**

Slide credit: Svetlana Lazebnik

B. Leibe

# 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?*

Slide adapted from Bill Freeman

B. Leibe

# Spatial Pyramid Representation

- **Representation in-between orderless BoW and global appearance**
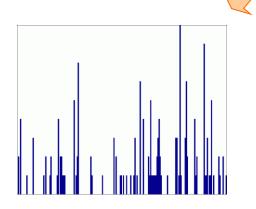
Slide credit: Svetlana Lazebnik          B. Leibe          [Lazebnik, Schmid & Ponce, CVPR'06]

# Spatial Pyramid Representation

- **Representation in-between orderless BoW and global appearance**

Slide credit: Svetlana Lazebnik

B. Leibe

[Lazebnik, Schmid & Ponce, CVPR'06]

# Spatial Pyramid Representation

- **Representation in-between orderless BoW and global appearance**



Slide credit: Svetlana Lazebnik

B. Leibe

[Lazebnik, Schmid & Ponce, CVPR'06]

37

# Summary: Bag-of-Words

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

- ## <u>Cons:</u>
    - ➢ Basic model ignores geometry – must verify afterwards, or encode via features.
    - ➢ Background and foreground mixed when bag covers whole image
    - ➢ Interest points or sampling: no guarantee to capture object-level parts.
    - ➢ Optimal vocabulary formation remains unclear.

B. Leibe

Slide credit: Kristen Grauman

# 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

# 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

B. Leibe

40

# 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**
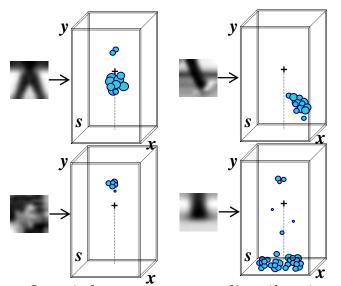
B. Leibe, A. Leonardis, and B. Schiele, Robust Object Detection with Interleaved Categorization and Segmentation, International Journal of Computer Vision, Vol. 77(1-3), 2008.
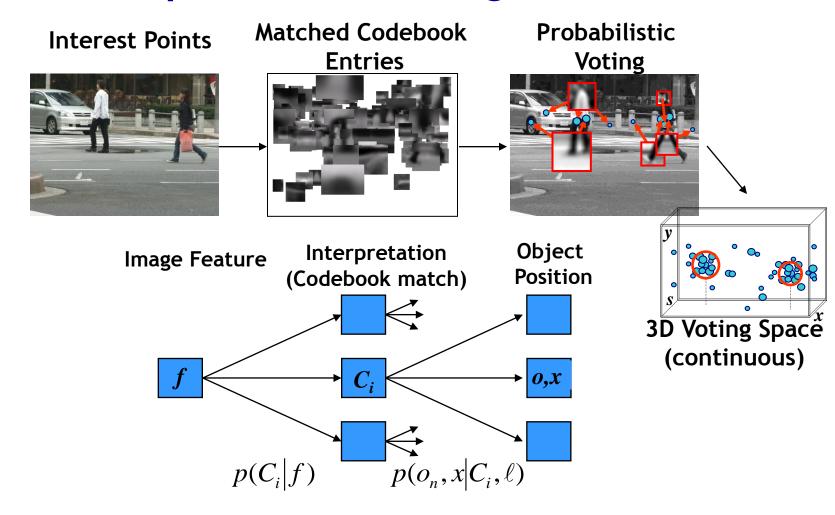
# Implicit Shape Model: Basic Idea

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



**Test image**

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

# 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**

B. Leibe

# Implicit Shape Model - Recognition

**Interest Points**

**Matched Codebook Entries**

**Probabilistic Voting**



**3D Voting Space (continuous)**

**Image Feature**

**Interpretation (Codebook match)**

**Object Position**



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

**Probabilistic vote weighting**

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

Computer Vision WS 14/15

# Implicit Shape Model - Recognition

**Interest Points**

**Matched Codebook Entries**

**Probabilistic Voting**



*y*

*s*

*x*

**3D Voting Space (continuous)**

**Backprojected Hypotheses**

**Backprojection of Maxima**

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

# Example: Results on Cows



## Original image

B. Leibe

# Example: Results on Cows



**Interest points**

B. Leibe

# Example: Results on Cows



## Matched patches

B. Leibe

# Example: Results on Cows



**Prob. Votes**

B. Leibe

# Example: Results on Cows



**1st hypothesis**

K. Grauman, B. Leibe

# Example: Results on Cows

**2nd hypothesis**

B. Leibe

51

# Example: Results on Cows



**3rd hypothesis**

B. Leibe

# Scale Invariant Voting

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



- **Generate scale votes**
  - ➢ Scale as 3$^{rd}$ 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

$s$

$y$

$x$

53

B. Leibe

# Detection Results

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

B. Leibe

# Detections Using Ground Plane Constraints



180°

150°     150°

90°

30°     30°

0°

**Battery of 5
ISM detectors
for different
car views**

**left camera
1175 frames**

57

B. Leibe     [Leibe, Cornelis, Cornelis, Van Gool, CVPR'07]

# 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)

B. Leibe

[Mikolajczyk, Leibe, Schiele, CVPR'06]

# Sometimes, Rotation Invariance Is Needed...

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

B. Leibe

# Implicit Shape Model – Segmentation

**Local Features**

**Matched Codebook Entries**

**Probabilistic Voting**

**3D Voting Space (continuous)**

**Segmentation**

**Backproject Meta-information**

**Pixel Contributions**

**Backprojected Hypotheses**

**Backprojection of Maxima**

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

# Example Results: Motorbikes

B. Leibe    [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

B. Leibe

# 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

# 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
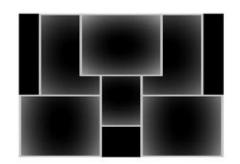
B. Leibe

64

# Deformable Part-based Models



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

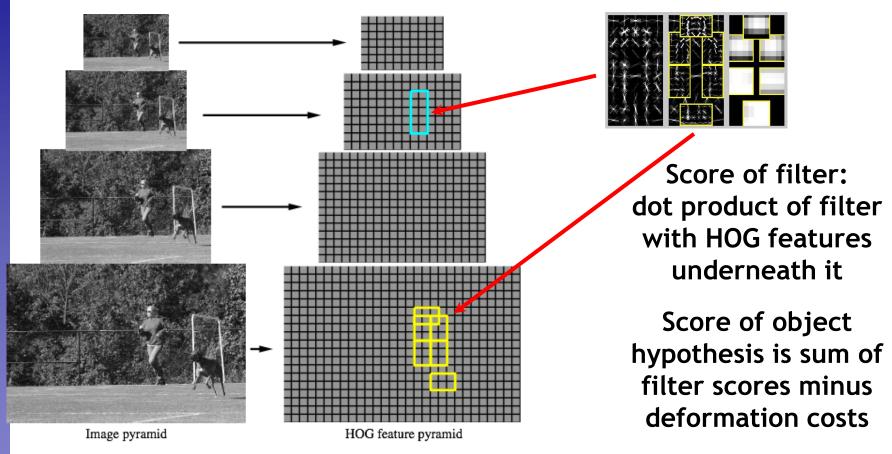65

# 2-Component Bicycle Model



**Root filters
coarse resolution**

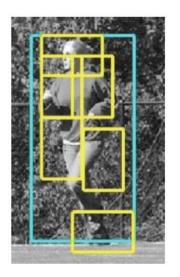**Part filters
finer resolution**

**Deformation
models**

Slide credit: Pedro Felzenszwalb

B. Leibe

# Object Hypothesis



Image pyramid

HOG feature pyramid

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

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

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

Slide credit: Pedro Felzenszwalb

B. Leibe

# Score of a Hypothesis



$$\text{score}(p_0, \ldots, p_n) = \underbrace{\sum_{i=0}^{n} F_i \cdot \phi(H, p_i)}_{\substack{\text{"data term"} \\ \text{filters}}} - \underbrace{\sum_{i=1}^{n} d_i \cdot (dx_i^2, dy_i^2)}_{\substack{\text{"spatial prior"} \\ \text{displacements} \\ \text{deformation parameters}}}$$

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

concatenation filters and deformation parameters

concatenation of HOG features and part displacement features

Slide credit: Pedro Felzenszwalb

B. Leibe

68

# 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**
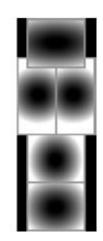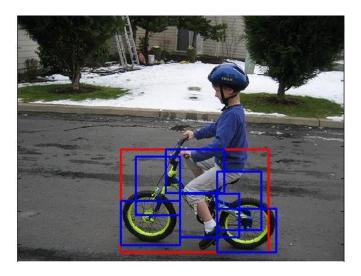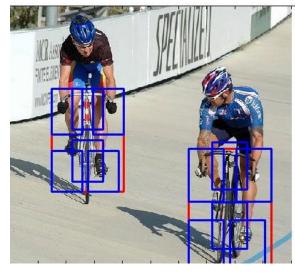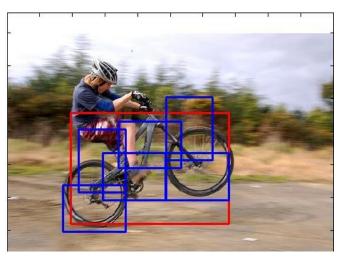
Slide credit: Pedro Felzenszwalb

B. Leibe

# Results: Persons



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

Slide credit: Pedro Felzenszwalb
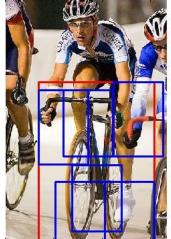
B. Leibe

# Results: Bicycles

Slide adapted from Trevor Darrell

B. Leibe

Computer Vision WS 14/15

# False Positives

- **Bicycles**

B. Leibe

# Results: Cats



**High-scoring true positives**

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

Slide credit: Pedro Felzenszwalb

# You Can Try It At Home…

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

$\Rightarrow$ **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

# 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