# Computer Vision – Lecture 16

## Deep Learning Applications

11.01.2017

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

leibe@vision.rwth-aachen.de

Computer Vision WS 16/17

---

## Announcements

- **Seminar registration period starts on Friday**
  - We will offer a lab course in the summer semester "Deep Robot Learning"
  - Topic: Deep reinforcement learning for robot control
    - Either UAV or grasping robot
  - If you're interested, you can register at http://www.graphics.rwth-aachen.de/apse
  - Registration period: 13.01.2016 – 29.01.2016

  - *Quick poll: Who would be interested in that?*

Computer Vision WS 16/17

2

---

## 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 16/17

3

---

## Recap: Convolutional Neural Networks



- **Neural network with specialized connectivity structure**
  - Stack multiple stages of feature extractors
  - Higher stages compute more global, more invariant features
  - Classification layer at the end

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86(11): 2278-2324, 1998.

Slide credit: Svetlana Lazebnik

B. Leibe

Computer Vision WS 16/17

4

---

## Recap: CNN Structure

- **Feed-forward feature extraction**
  1. Convolve input with learned filters
  2. Non-linearity
  3. Spatial pooling
  4. (Normalization)
- **Supervised training of convolutional filters by back-propagating classification error**

Feature maps
⇧
Normalization
⇧
Spatial pooling
⇧
Non-linearity
⇧
Convolution (Learned)
⇧
Input Image

Slide credit: Svetlana Lazebnik

B. Leibe

Computer Vision WS 16/17

5

---

## Recap: Intuition of CNNs



- **Convolutional net**
  - Share the same parameters across different locations
  - Convolutions with learned kernels

- **Learn *multiple* filters**
  - E.g. 1000×1000 image
    100 filters
    10×10 filter size
  ⇒ only 10k parameters

- **Result: Response map**
  - size: 1000×1000×100
  - Only memory, not params!

Slide adapted from Marc'Aurelio Ranzato

B. Leibe

Image source: Yann LeCun

Computer Vision WS 16/17

6

## Recap: Convolution Layers



Naming convention:
HEIGHT / WIDTH / DEPTH

- **All Neural Net activations arranged in 3 dimensions**
  - Multiple neurons all looking at the same input region, stacked in depth
  - Form a single [1×1×depth] depth column in output volume.

Computer Vision WS 16/17

Slide credit: FeiFei Li, Andrej Karpathy    B. Leibe    7

---

## Recap: Activation Maps



one filter = one depth slice (or activation map)

**5×5 filters**

**Activation maps**

Each activation map is a depth slice through the output volume.

Computer Vision WS 16/17

Slide adapted from FeiFei Li, Andrej Karpathy    B. Leibe    8

---

## Recap: Pooling Layers



max pool with 2x2 filters and stride 2

- **Effect:**
  - Make the representation smaller without losing too much information
  - Achieve robustness to translations

Computer Vision WS 16/17

Slide adapted from FeiFei Li, Andrej Karpathy    B. Leibe    9

---

## Recap: Effect of Multiple Convolution Layers



Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Computer Vision WS 16/17

Slide credit: Yann LeCun    B. Leibe    10

---

## Recap: AlexNet (2012)



- **Similar framework as LeNet, but**
  - Bigger model (7 hidden layers, 650k units, 60M parameters)
  - More data ($10^6$ images instead of $10^3$)
  - GPU implementation
  - Better regularization and up-to-date tricks for training (Dropout)

A. Krizhevsky, I. Sutskever, and G. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012.

Computer Vision WS 16/17

11

Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012

---

## Recap: VGGNet (2014/15)

- **Main ideas**
  - Deeper network
  - Stacked convolutional layers with smaller filters (+ nonlinearity)
  - Detailed evaluation of all components

- **Results**
  - Improved ILSVRC top-5 error rate to 6.7%.



Computer Vision WS 16/17

B. Leibe    12

Image source: Simonyan & Zisserman

## Recap: GoogLeNet (2014)

- **Ideas:**
  - Learn features at multiple scales
  - Modular structure



**Inception module** **+ copies**

**Convolution**
**Pooling**
**Softmax**
**Other**

**Auxiliary classification outputs for training the lower layers (deprecated)**

(b) Inception module with dimension reductions

*Computer Vision WS 16/17*

B. Leibe

13

Image source: Szegedy et al.

---

## Recap: Residual Networks

AlexNet, 8 layers (ILSVRC 2012)     VGG, 19 layers (ILSVRC 2014)     ResNet, 152 layers (ILSVRC 2015)

- **Core component**
  - Skip connections bypassing each layer
  - Better propagation of gradients to the deeper layers
  - This makes it possible to train (much) deeper networks.



$x$

weight layer

$F(x)$    relu

weight layer

$H(x) = F(x) + x$    relu

*Computer Vision WS 16/17*

B. Leibe

14

---

## Transfer Learning with CNNs



**1. Train on ImageNet**

**2. If small dataset: fix all weights (treat CNN as fixed feature extractor), retrain only the classifier**

**I.e., swap the Softmax layer at the end**

*Computer Vision WS 16/17*

Slide credit: Andrej Karpathy

B. Leibe

15

---

## Transfer Learning with CNNs



**1. Train on ImageNet**

**3. If you have medium sized dataset, "finetune" instead: use the old weights as initialization, train the full network or only some of the higher layers.**

**Retrain bigger portion of the network**

*Computer Vision WS 16/17*

Slide credit: Andrej Karpathy

B. Leibe

16

---

## Topics of This Lecture

- **Object Detection with CNNs**
  - R-CNN
  - Fast R-CNN
  - Faster R-CNN

- **Semantic Image Segmentation**

- **Human Pose Estimation**

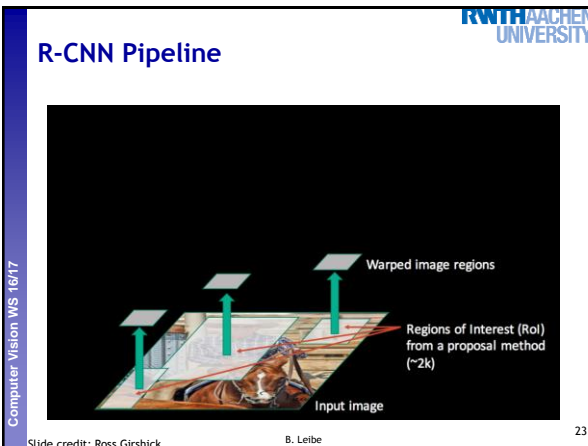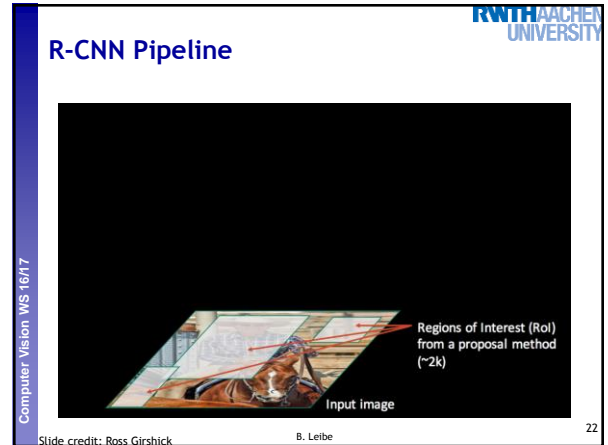- **Face/Person Identification**
  - DeepFace
  - FaceNet

*Computer Vision WS 16/17*

B. Leibe

17

---

## The Learned Features are Generic



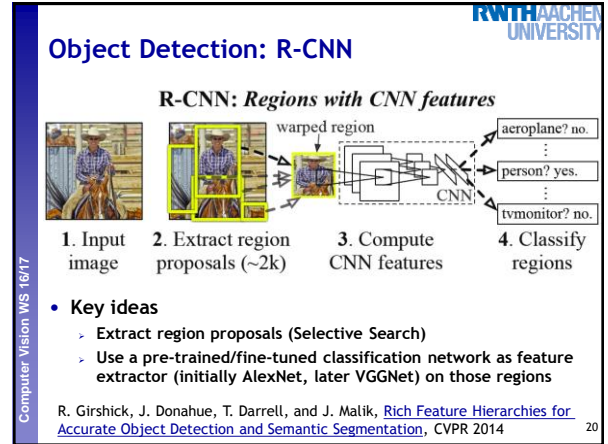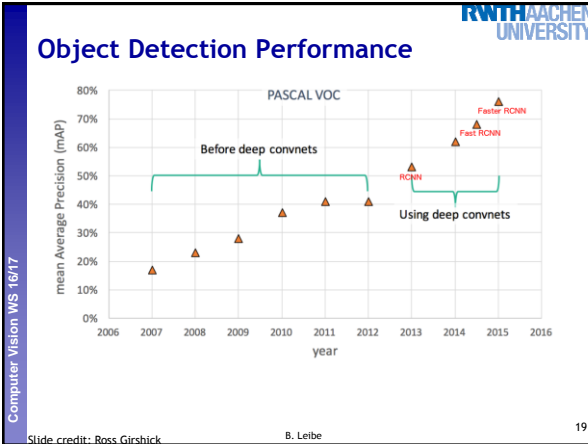state of the art level (pre-CNN)

- **Experiment: feature transfer**
  - Train network on ImageNet
  - Chop off last layer and train classification layer on CalTech256
  - ⇒ State of the art accuracy already with only 6 training images

*Computer Vision WS 16/17*

B. Leibe

18

Image source: M. Zeiler, R. Fergus

## Object Detection Performance

PASCAL VOC

Before deep convnets

Using deep convnets

Faster RCNN
Fast RCNN
RCNN

mean Average Precision (mAP)

year

Computer Vision WS 16/17

Slide credit: Ross Girshick

B. Leibe

19

---

## Object Detection: R-CNN

**R-CNN: *Regions with CNN features***

warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

- **Key ideas**
  - **Extract region proposals (Selective Search)**
  - **Use a pre-trained/fine-tuned classification network as feature extractor (initially AlexNet, later VGGNet) on those regions**

R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014

Computer Vision WS 16/17

20

---

## R-CNN Pipeline

Input image

Computer Vision WS 16/17

Slide credit: Ross Girshick

B. Leibe

21

---

## R-CNN Pipeline

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Computer Vision WS 16/17

Slide credit: Ross Girshick

B. Leibe

22

---

## R-CNN Pipeline

Warped image regions

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Computer Vision WS 16/17

Slide credit: Ross Girshick

B. Leibe

23

---

## R-CNN Pipeline

ConvNet
ConvNet
ConvNet

Forward each region through ConvNet

Warped image regions

Regions of Interest (RoI) from a proposal method (~2k)

Input image

Computer Vision WS 16/17

Slide credit: Ross Girshick

B. Leibe

24

---

## R-CNN Pipeline

Slide credit: Ross Girshick — B. Leibe — 25

## R-CNN Pipeline

Slide credit: Ross Girshick — B. Leibe — 26

## Classification



Input Image — Region Proposals — Feature Extraction — Classification

- **Linear model with class-dependent weights**
  - Linear SVM

$$f_c(x_{fc7}) = w_c^T x_{fc7}$$

  - where
    - $x_{fc7}$ = features from the network (fully-connected layer 7)
    - $c$ = object class

Slide credit: Ross Girshick, Kaustav Kundu — B. Leibe — 27

## Bounding Box Regressors

- **Prediction of the 2D box**
  - Necessary, since the proposal region might not fully coincide with the (annotated) object bounding box
  - Perform regression for location $(x^*, y^*)$, width $w^*$ and height $h^*$

$$\frac{x^* - x}{w} = w_{c,x}^T x_{pool5}$$
$$\frac{y^* - y}{h} = w_{c,y}^T x_{pool5}$$
$$\ln \frac{w^*}{w} = w_{c,w}^T x_{pool5}$$
$$\ln \frac{h^*}{h} = w_{c,w}^T x_{pool5}$$

  - Where $x_{pool5}$ are the features from the pool5 layer of the network.

Slide credit: Ross Girshick, Kaustav Kundu — B. Leibe — 28

## Problems with R-CNN

- **Ad hoc training objectives**
  - Fine tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-box regressors (squared loss)

- **Training (3 days) and testing (47s per image) is slow.**
  - Many separate applications of region CNNs

- **Takes a lot of disk space**
  - Need to store all precomputed CNN features for training the classifiers
  - Easily 200GB of data



Slide credit: Ross Girshick — B. Leibe — 29

## Fast R-CNN

- **Forward Pass**



Slide credit: Ross Girshick — B. Leibe — 30

## Fast R-CNN

- **Forward Pass**



Slide credit: Ross Girshick

B. Leibe

31

## Fast R-CNN

- **Forward Pass**



Slide credit: Ross Girshick

B. Leibe

32

## Fast R-CNN Training

- **Backward Pass**



Slide credit: Ross Girshick

B. Leibe

33

## Region Proposal Networks (RPN)

- **Idea**
  - Remove dependence on external region proposal algorithm.
  - Instead, infer region proposals from same CNN.
  - ⇒ Feature sharing
  - ⇒ Object detection in a single pass becomes possible.

- **Faster R-CNN = Fast R-CNN + RPN**



Slide credit: Ross Girshick

34

## Faster R-CNN

- **One network, four losses**
  - Joint training



Slide credit: Ross Girshick

35

## Faster R-CNN (based on ResNets)



K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

B. Leibe

36

## Slide 37

### Faster R-CNN (based on ResNets)



Computer Vision WS 16/17

K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

B. Leibe

37

## Slide 38

### Summary

- **Object Detection**
  - Find a variable number of objects by classifying image regions
  - Before CNNs: dense multiscale sliding window (HoG, DPM)
  - Avoid dense sliding window with region proposals
  - R-CNN: Selective Search + CNN classification / regression
  - Fast R-CNN: Swap order of convolutions and region extraction
  - Faster R-CNN: Compute region proposals within the network
  - Deeper networks do better

Computer Vision WS 16/17

38

## Slide 39

### Topics of This Lecture

- **Object Detection with CNNs**
  - R-CNN
  - Fast R-CNN
  - Faster R-CNN
- **Semantic Image Segmentation**
- **Human Pose Estimation**
- **Face/Person Identification**
  - DeepFace
  - FaceNet

Computer Vision WS 16/17

B. Leibe

39

## Slide 40

### Semantic Image Segmentation



- **Perform pixel-wise prediction task**
  - Usually done using **Fully Convolutional Networks** (FCNs)
    - All operations formulated as convolutions
    - Advantage: can process arbitrarily sized images

Computer Vision WS 16/17

40

Image source: Long, Shelhamer, Darrell

## Slide 41

### CNNs vs. FCNs

- **CNN**
- **FCN**



- **Intuition**
  - Think of FCNs as performing a sliding-window classification, producing a heatmap of output scores for each class

Computer Vision WS 16/17

41

Image source: Long, Shelhamer, Darrell

## Slide 42

### Semantic Image Segmentation



- **Encoder-Decoder Architecture**
  - Problem: FCN output has low resolution
  - Solution: perform upsampling to get back to desired resolution
  - Use skip connections to preserve higher-resolution information

Computer Vision WS 16/17

42

Image source: Newell et al.

## Other Tasks: Semantic Segmentation



[Farabet et al. ICML 2012, PAMI 2013]

Computer Vision WS 16/17

B. Leibe

43

## Semantic Segmentation



[Pohlen, Hermans, Mathias, Leibe, arXiv 2016]

• **More recent results**
  ➢ **Based on an extension of ResNets**

Computer Vision WS 16/17

## Topics of This Lecture

• **Object Detection with CNNs**
  ➢ R-CNN
  ➢ Fast R-CNN
  ➢ Faster R-CNN

• **Semantic Image Segmentation**

• **Human Pose Estimation**

• Face/Person Identification
  ➢ DeepFace
  ➢ FaceNet

Computer Vision WS 16/17

B. Leibe

45

## FCNs for Human Pose Estimation

• **Input data**



Image    Keypoints    Labels

• **Task setup**
  ➢ **Annotate images with keypoints for skeleton joints**
  ➢ **Define a target disk around each keypoint with radius r**
  ➢ **Set the ground-truth label to 1 within each such disk**
  ➢ **Infer heatmaps for the joints as in semantic segmentation**

Computer Vision WS 16/17

Slide adapted from Georgia Gkioxari

46

## Heat Map Predictions from FCN



Test Image   Right Ankle   Right Knee   Right Hip   Right Wrist   Right Elbow   Right Shoulder

Computer Vision WS 16/17

Slide adapted from Georgia Gkioxari

47

## Example Results: Human Pose Estimation



[Rafi, Gall, Leibe, BMVC 2016]

Computer Vision WS 16/17

48

## Slide 49

**Topics of This Lecture**

- Object Detection with CNNs
  - R-CNN
  - Fast R-CNN
  - Faster R-CNN
- Semantic Image Segmentation
- Human Pose Estimation
- **Face/Person Identification**
  - **DeepFace**
  - **FaceNet**

Computer Vision WS 16/17

B. Leibe

49

## Slide 50

**Other Tasks: Face Verification**



Y. Taigman, M. Yang, M. Ranzato, L. Wolf, DeepFace: Closing the Gap to Human-Level Performance in Face Verification, CVPR 2014

Slide credit: Svetlana Lazebnik

50

## Slide 51

**Discriminative Face Embeddings**

- **Learning an embedding using a Triplet Loss Network**
  - **Present the network with triplets of examples**

Negative    Anchor    Positive



  - Apply triplet loss to learn an embedding $f(\cdot)$ that groups the positive example closer to the anchor than the negative one.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$

⇒ Used with great success in Google's FaceNet face recognition

B. Leibe

51

## Slide 52

**Vector Arithmetics in Embedding Space**

- **Learned embeddings often preserve linear regularities between concepts**
  - **Analogy questions can be answered through simple algebraic operations with the vector representation of words.**
  - **E.g., vec("King") – vec("Man") + vec("Woman") ≈ vec("Queen")**
  - **E.g.,**



smiling woman    neutral woman    neutral man    smiling man

B. Leibe

[Mikolov, NIPS 2013], [Radford, ICLR 2016]

52

## Slide 53

**Commercial Recognition Services**

- **E.g., clarifai**

Try it out with your own media

Upload an image or video file under 100mb or give us a direct link to a file on the web.

Paste a url here...    ENGLISH ▼

USE THE URL    CHOOSE A FILE INSTEAD

*By using the demo you agree to our terms of service

- **Be careful when taking test images from Google Search**
  - Chances are they may have been seen in the training set...

B. Leibe

Image source: clarifai.com

53

## Slide 54

**Commercial Recognition Services**



clarifai

B. Leibe

Image source: clarifai.com

54

## References and Further Reading

- RCNN and related ideas:
  - Girshick et al., Region-based Convolutional Networks for Accurate Object Detection and Semantic Segmentation, PAMI, 2014.
  - Zhu et al., segDeepM: Exploiting Segmentation and Context in Deep Neural Networks for Object Detection, 2015.
- Fast RCNN and related ideas:
  - He et al., Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition, 2014.
  - Girshick, Ross, Fast R-CNN, 2015.
- Faster RCNN and related ideas:
  - Szegedy et al., Scalable, High-Quality Object Detection, 2014.
  - Ren et al., Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, 2015.

10