

RWTH AACHEN
UNIVERSITY

Advanced Machine Learning Lecture 15

Convolutional Neural Networks III

12.01.2017

Bastian Leibe
RWTH Aachen
<http://www.vision.rwth-aachen.de/>
leibe@vision.rwth-aachen.de

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

Announcement

- Lecture evaluation
 - Please fill out the evaluation forms...

B. Leibe

2

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

This Lecture: *Advanced Machine Learning*

- Regression Approaches
 - Linear Regression
 - Regularization (Ridge, Lasso)
 - Kernels (Kernel Ridge Regression)
 - Gaussian Processes
- Approximate Inference
 - Sampling Approaches
 - MCMC
- Deep Learning
 - Linear Discriminants
 - Neural Networks
 - Backpropagation & Optimization
 - CNNs, RNNs, ResNets, etc.

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

Topics of This Lecture

- Recap: CNN Architectures
- Residual Networks
- Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

B. Leibe

4

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

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.

B. Leibe

5

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

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.

B. Leibe

6

Advanced Machine Learning Winter'16

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%.

ConvNet Configuration				
A	A+LRN	B	C	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)				
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool				
FC-4096				
FC-4096				
FC-1000				
soft-max				
Mainly used				

B. Leibe 7

Recap: GoogLeNet (2014)

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

B. Leibe 8

Recap: Visualizing CNNs

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

B. Leibe 10

Topics of This Lecture

- Recap: CNN Architectures
- Residual Networks
- Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

B. Leibe 11

Newest Development: Residual Networks

B. Leibe 12

Newest Development: Residual Networks

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

B. Leibe 13

RWTH AACHEN UNIVERSITY

Spectrum of Depth

shallow ← → deeper

5 layers: easy
 >10 layers: initialization, Batch Normalization
 >30 layers: skip connections
 >100 layers: identity skip connections
 >1000 layers: ?

Advanced Machine Learning Winter'16

Slide credit: Kaiming He

B. Leibe

14

RWTH AACHEN UNIVERSITY

Spectrum of Depth

shallow ← → deeper

5 layers: easy
 >10 layers: initialization, Batch Normalization
 >30 layers: skip connections
 >100 layers: identity skip connections
 >1000 layers: ?

- Deeper models are more powerful
 - > But training them is harder.
 - > Main problem: getting the gradients back to the early layers
 - > The deeper the network, the more effort is required for this.

Advanced Machine Learning Winter'16

Slide adapted from Kaiming He

B. Leibe

15

RWTH AACHEN UNIVERSITY

Initialization

22-layer ReLU net: good init converges faster

30-layer ReLU net: good init is able to converge

- Importance of proper initialization (Recall Lecture 11)
 - > Glorot initialization for tanh nonlinearities
 - > He initialization for ReLU nonlinearities
 - ⇒ For deep networks, this really makes a difference!

Advanced Machine Learning Winter'16

Slide credit: Kaiming He

B. Leibe

16

RWTH AACHEN UNIVERSITY

Batch Normalization

accuracy

iter.

- Effect of batch normalization
 - > Greatly improved speed of convergence

Advanced Machine Learning Winter'16

Image source: Ioffe & Szegedy

B. Leibe

17

RWTH AACHEN UNIVERSITY

Going Deeper

- Checklist
 - > Initialization ok
 - > Batch normalization ok
 - > Are we now set?
 - Is learning better networks now as simple as stacking more layers?

Advanced Machine Learning Winter'16

Slide credit: Kaiming He

B. Leibe

18

RWTH AACHEN UNIVERSITY

Simply Stacking Layers?

train error (%)

test error (%)

56-layer

20-layer

- Experiment going deeper
 - > Plain nets: stacking 3x3 convolution layers
 - ⇒ 56-layer net has higher training error than 20-layer net

Advanced Machine Learning Winter'16

Slide credit: Kaiming He

B. Leibe

19

RWTH AACHEN UNIVERSITY

Simply Stacking Layers?

CIFAR-10: 56-layer, 44-layer, 32-layer, 20-layer
 ImageNet-1000: 34-layer, 18-layer

solid: test/val
 dashed: train

- General observation
 - Overly deep networks have higher training error
 - A general phenomenon, observed in many training sets

Slide credit: Kaiming He
 B. Leibe
 20

RWTH AACHEN UNIVERSITY

Why Is That???

- A deeper model should not have higher training error!
 - Richer solution space should allow it to find better solutions
- Solution by construction
 - Copy the original layers from a learned shallower model
 - Set the extra layers as identity
 - Such a network should achieve at least the same low training error.
- Reason: Optimization difficulties
 - Solvers cannot find the solution when going deeper...

Slide credit: Kaiming He
 B. Leibe
 21

RWTH AACHEN UNIVERSITY

Deep Residual Learning

- Plain net
 - any two stacked layers
 - weight layer → relu
 - weight layer → relu
 - $H(x)$

$H(x)$ is any desired mapping
 Hope the 2 weight layers fit $H(x)$

Slide credit: Kaiming He
 B. Leibe
 22

RWTH AACHEN UNIVERSITY

Deep Residual Learning

- Residual net
 - weight layer → relu
 - weight layer → relu
 - identity shortcut
 - $F(x)$
 - $H(x) = F(x) + x$

$H(x)$ is any desired mapping
 Hope the 2 weight layers fit $H(x)$
 Hope the 2 weight layers fit $F(x)$
 Let $H(x) = F(x) + x$

Slide credit: Kaiming He
 B. Leibe
 23

RWTH AACHEN UNIVERSITY

Deep Residual Learning

- $F(x)$ is a residual mapping w.r.t. identity
 - weight layer → relu
 - weight layer → relu
 - identity shortcut
 - $F(x)$
 - $H(x) = F(x) + x$

If identity were optimal, it is easy to set weights as 0
 If optimal mapping is closer to identity, it is easier to find small fluctuations
 Further advantage: direct path for the gradient to flow to the previous stages

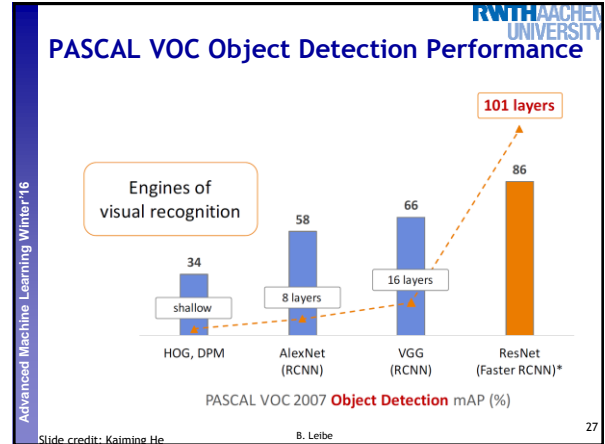
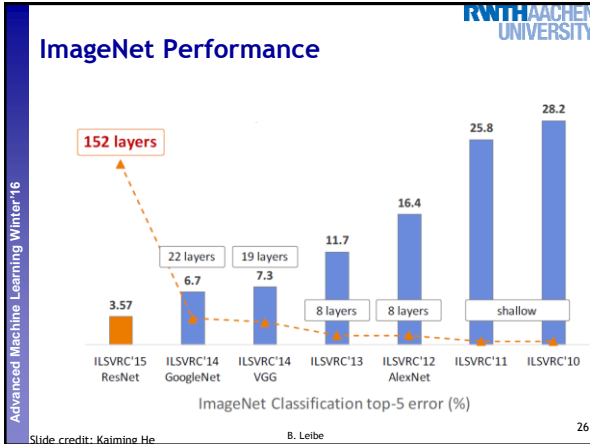
Slide credit: Kaiming He
 B. Leibe
 24

RWTH AACHEN UNIVERSITY

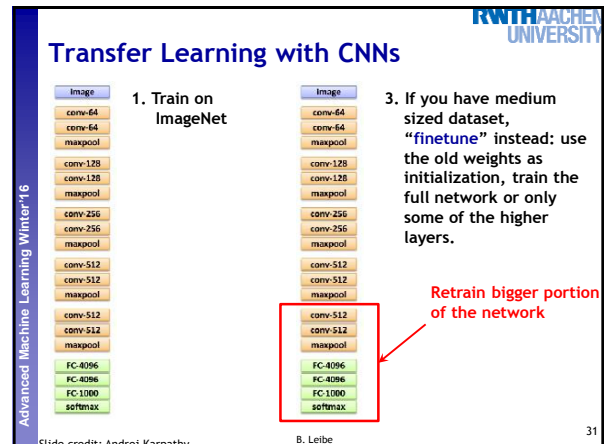
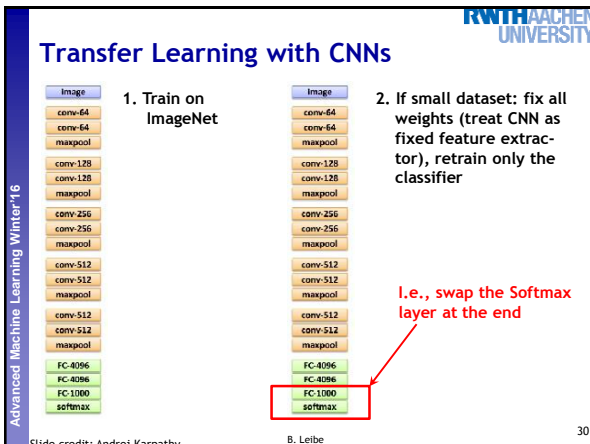
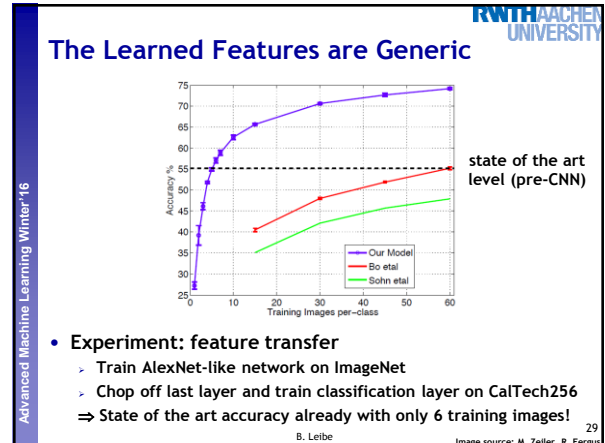
Network Design

- Simple, VGG-style design
 - (Almost) all 3x3 convolutions
 - Spatial size / 2 ⇒ #filters · 2 (same complexity per layer)
 - Batch normalization
 - ⇒ Simple design, just deep.

Slide credit: Kaiming He
 B. Leibe
 25



- ### Topics of This Lecture
- Recap: CNN Architectures
 - Residual Networks
 - Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification
- Slide credit: B. Leibe



RWTH AACHEN UNIVERSITY

Other Tasks: Object Detection

R-CNN: Regions with CNN features

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

RWTH AACHEN UNIVERSITY

Object Detection: R-CNN

R-CNN: Regions with CNN features

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

- Results on PASCAL VOC Detection benchmark
 - Pre-CNN state of the art: 35.1% mAP [Uijlings et al., 2013]
 - 33.4% mAP DPM
 - R-CNN: 53.7% mAP

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

RWTH AACHEN UNIVERSITY

Most Recent Version: Faster R-CNN

- One network, four losses
 - Remove dependence on external region proposal algorithm.
 - Instead, infer region proposals from same CNN.
 - Feature sharing
 - Joint training
 - ⇒ Object detection in a single pass becomes possible.

Slide credit: Ross Girshick

RWTH AACHEN UNIVERSITY

Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

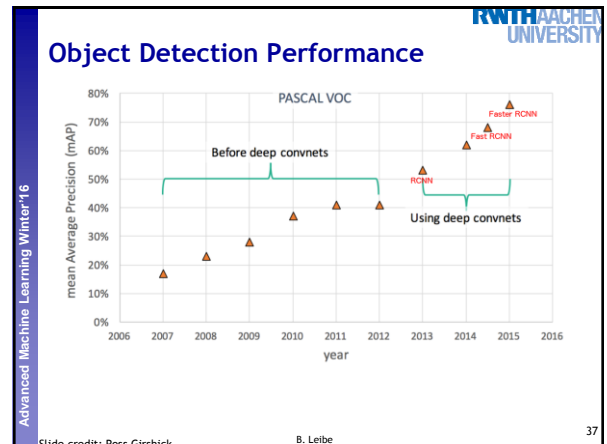
B. Leibe

RWTH AACHEN UNIVERSITY

Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

B. Leibe



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

Image source: Long, Shelhamer, Darrell

CNNs vs. FCNs

- CNN
 - Output: "tabby cat"
- FCN
 - convolutionalization
 - Output: tabby cat heatmap
- Intuition
 - Think of FCNs as performing a sliding-window classification, producing a heatmap of output scores for each class

Image source: Long, Shelhamer, Darrell

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

Image source: Newell et al.

Semantic Segmentation

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

- More recent results
 - Based on an extension of ResNets

Other Tasks: Face Identification

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

References: Computer Vision Tasks

- Object Detection
 - R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014.
 - S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NIPS 2015.
 - J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You Only Look Once: Unified Real-Time Object Detection, CVPR 2016.
 - W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A.C. Berg, SSD: Single Shot Multi Box Detector, ECCV 2016.

B. Leibe

References: Computer Vision Tasks

- **Semantic Segmentation**

- J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation, CVPR 2015.
- H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid Scene Parsing Network, arXiv 1612.01105, 2016.