

RWTH AACHEN
UNIVERSITY

Computer Vision – Lecture 11

Deep Learning II

28.05.2019

Computer Vision Summer'19

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

RWTH AACHEN
UNIVERSITY

Course Outline

- Image Processing Basics
- Segmentation & Grouping
- Object Recognition & Categorization
 - Sliding Window based Object Detection
- Local Features & Matching
 - Local Features – Detection and Description
 - Recognition with Local Features
- **Deep Learning**
 - Convolutional Neural Networks
- 3D Reconstruction

2

RWTH AACHEN
UNIVERSITY

Topics of This Lecture

- **Recap: Convolutional Neural Networks**
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- Background: Deep Learning
 - Recap from ML lecture
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
 - ResNet

Computer Vision Summer'19

3

RWTH AACHEN
UNIVERSITY

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

Computer Vision Summer'19

4

RWTH AACHEN
UNIVERSITY

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
 - ⇒ 10k parameters
- Result: Response map
 - size: $1000 \times 1000 \times 100$
 - Only memory, not params!

Computer Vision Summer'19

5

RWTH AACHEN
UNIVERSITY

Recap: Convolution Layers


- All Neural Net activations arranged in 3 dimensions
 - Multiple neurons all looking at the same input region, stacked in depth
 - Form a single $[1 \times 1 \times \text{depth}]$ depth column in output volume.


Computer Vision Summer'19

6

RWTH AACHEN UNIVERSITY

Recap: Activation Maps

Activations:  one filter = one depth slice (or activation map) 5 × 5 filters

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

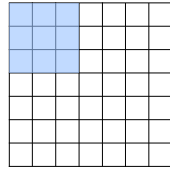
Activation maps

Computer Vision Summer'19 7

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
7 × 7 input
assume 3 × 3 connectivity
stride 1

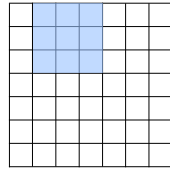
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 8

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
7 × 7 input
assume 3 × 3 connectivity
stride 1

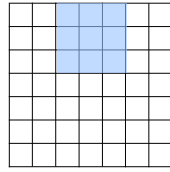
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 9

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
7 × 7 input
assume 3 × 3 connectivity
stride 1

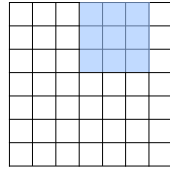
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 10

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
7 × 7 input
assume 3 × 3 connectivity
stride 1

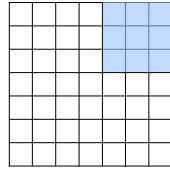
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 11

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
7 × 7 input
assume 3 × 3 connectivity
stride 1
⇒ 5 × 5 output

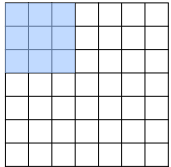
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 12

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
 7×7 input
 assume 3×3 connectivity
 stride 1
 $\Rightarrow 5 \times 5$ output

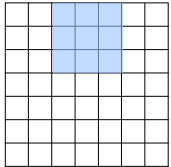
What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 | Slide credit: FeiFei Li, Andrej Karpathy | B. Leibe | 13

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
 7×7 input
 assume 3×3 connectivity
 stride 1
 $\Rightarrow 5 \times 5$ output

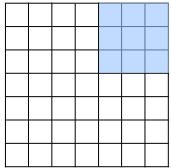
What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 | Slide credit: FeiFei Li, Andrej Karpathy | B. Leibe | 14

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
 7×7 input
 assume 3×3 connectivity
 stride 1
 $\Rightarrow 5 \times 5$ output

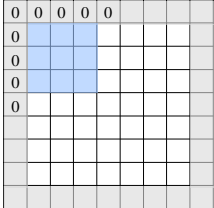
What about stride 2?
 $\Rightarrow 3 \times 3$ output

- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 | Slide credit: FeiFei Li, Andrej Karpathy | B. Leibe | 15

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
 7×7 input
 assume 3×3 connectivity
 stride 1
 $\Rightarrow 5 \times 5$ output

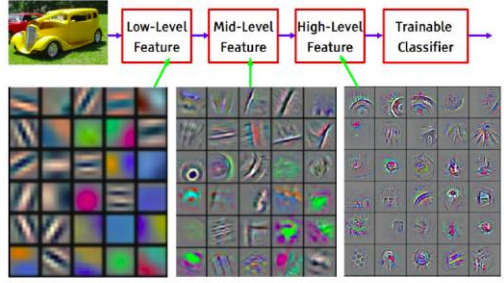
What about stride 2?
 $\Rightarrow 3 \times 3$ output

- Replicate this column of hidden neurons across space, with some **stride**.
- In practice, common to zero-pad the border.
 - Preserves the size of the input spatially.

Computer Vision Summer'19 | Slide credit: FeiFei Li, Andrej Karpathy | B. Leibe | 16

RWTH AACHEN UNIVERSITY

Effect of Multiple Convolution Layers



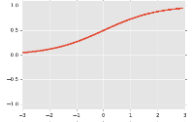
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

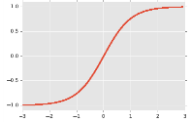
Computer Vision Summer'19 | Slide credit: Yann LeCun | B. Leibe | 19

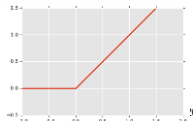
RWTH AACHEN UNIVERSITY

Commonly Used Nonlinearities

- Sigmoid

$$g(a) = \sigma(a) = \frac{1}{1 + \exp\{-a\}}$$

- Hyperbolic tangent

$$g(a) = \tanh(a) = 2\sigma(2a) - 1$$

- Rectified linear unit (ReLU)

$$g(a) = \max\{0, a\}$$


Preferred option for deep networks

Computer Vision Summer'19 | B. Leibe | 20

Computer Vision Summer'19

Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- Let's assume the filter is an eye detector
 - How can we make the detection robust to the exact location of the eye?

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 21

Computer Vision Summer'19

Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- Let's assume the filter is an eye detector
 - How can we make the detection robust to the exact location of the eye?
- Solution:
 - By **pooling** (e.g., max or avg) filter responses at different spatial locations, we gain robustness to the exact spatial location of features.

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 22

Computer Vision Summer'19

Max Pooling

RWTH AACHEN UNIVERSITY

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

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

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 23

Computer Vision Summer'19

Max Pooling

RWTH AACHEN UNIVERSITY

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

- Note
 - Pooling happens independently across each slice, preserving the number of slices.

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 24

Computer Vision Summer'19

Compare: SIFT Descriptor

RWTH AACHEN UNIVERSITY

Image Pixels

Apply oriented filters [Lowe [ICV 2004]]

Spatial pool (Sum)

Normalize to unit length

Feature Vector

Slide credit: Svetlana Lazebnik. B. Leibe. 25

Computer Vision Summer'19

Compare: Spatial Pyramid Matching

RWTH AACHEN UNIVERSITY

SIFT features

Filter with Visual Words [Lazebnik, Schmid, Ponce [CVPR 2006]]

Take max VW response (L-inf normalization)

Multi-scale spatial pool (Sum)

Global image descriptor

Slide credit: Svetlana Lazebnik. B. Leibe. 26

RWTH AACHEN UNIVERSITY

Topics of This Lecture

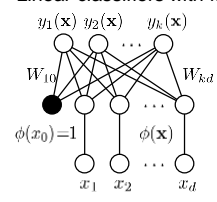
- Recap: Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- Background: Deep Learning
 - Recap from ML lecture
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
 - ResNet

27

RWTH AACHEN UNIVERSITY

Recap: Generalized Linear Discriminants

- Linear classifiers with fixed feature transformation



Output layer

Weights

Feature layer

Mapping (fixed)

Input layer
- Outputs
 - Linear outputs
$$y_k(\mathbf{x}) = \sum_{i=0}^d W_{ki} \phi(x_i)$$
- Logistic outputs

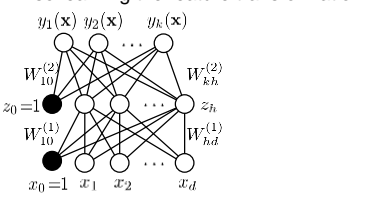
$$y_k(\mathbf{x}) = \sigma \left(\sum_{i=0}^d W_{ki} \phi(x_i) \right)$$

28

RWTH AACHEN UNIVERSITY

Recap: Multi-Layer Perceptrons

- Also learning the feature transformation



Output layer

Hidden layer

Mapping (learned!)

Input layer
- Output

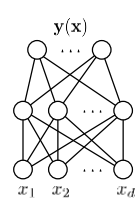
$$y_k(\mathbf{x}) = g^{(2)} \left(\sum_{i=0}^h W_{ki}^{(2)} g^{(1)} \left(\sum_{j=0}^d W_{ij}^{(1)} x_j \right) \right)$$

29

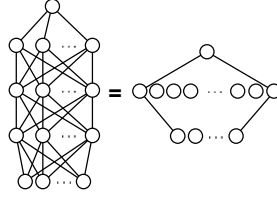
RWTH AACHEN UNIVERSITY

Two Important Remarks

- Why are hierarchical multi-layered models attractive?



An MLP with 1 hidden layer can implement any function (universal approximator)



However, if the function is deep, a very large hidden layer may be required.

30

RWTH AACHEN UNIVERSITY

Two Important Remarks

- Nonlinearities are essential for deep models

$$y_k(\mathbf{x}) = g^{(2)} \left(\sum_{i=0}^h W_{ki}^{(2)} g^{(1)} \left(\sum_{j=0}^d W_{ij}^{(1)} x_j \right) \right)$$
 - If we leave out the nonlinearity $g^{(1)}(\cdot)$, the two layers collapse into a single linear function

$$\tilde{W} = W^{(1)}W^{(2)}$$

⇒ The nonlinearities make multi-layer representation more powerful!

31

RWTH AACHEN UNIVERSITY

Recap: Supervised Learning

- Given
 - Training data set $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ with target labels $\mathbf{t} = (t_1, \dots, t_N)^T$.
- Solve an optimization problem
 - Set up an error function

$$E(\mathbf{W}) = \sum_n L(t_n, y(\mathbf{x}_n; \mathbf{W})) + \lambda \Omega(\mathbf{W})$$

with a loss $L(\cdot)$ and a regularizer $\Omega(\cdot)$.
 - E.g., $L(t, y(\mathbf{x}; \mathbf{W})) = \sum_n (y(\mathbf{x}_n; \mathbf{W}) - t_n)^2$ L₂ loss
 - $\Omega(\mathbf{W}) = \|\mathbf{W}\|_F^2$ L₂ regularizer ("weight decay")

⇒ Update each weight $W_{ij}^{(k)}$ in the direction of the gradient $\frac{\partial L(\mathbf{W})}{\partial W_{ij}^{(k)}}$

32

RWTH AACHEN UNIVERSITY

Recap: Loss Functions

- We can now also apply other loss functions
 - L2 loss $L(t, y(\mathbf{x})) = \sum_n (y(\mathbf{x}_n) - t_n)^2$ \Rightarrow Least-squares regression
 - L1 loss: $L(t, y(\mathbf{x})) = \sum_n |y(\mathbf{x}_n) - t_n|$ \Rightarrow Median regression
 - Cross-entropy loss $L(t, y(\mathbf{x})) = -\sum_n \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$ \Rightarrow Logistic regression
 - Hinge loss $L(t, y(\mathbf{x})) = \sum_n [1 - t_n y(\mathbf{x}_n)]_+$ \Rightarrow SVM classification
 - Softmax loss $L(t, y(\mathbf{x})) = -\sum_n \sum_k \left\{ \mathbb{1}(t_n = k) \ln \frac{\exp(y_k(\mathbf{x}))}{\sum_j \exp(y_j(\mathbf{x}))} \right\}$ \Rightarrow Multi-class probabilistic classification

33

RWTH AACHEN UNIVERSITY

Recap: Obtaining the Gradients

- Approach: Incremental Analytical Differentiation
 - Idea: Compute the gradients layer by layer.
 - Each layer below builds upon the results of the layer above.
 - The gradient is propagated backwards through the layers.
 - \Rightarrow Backpropagation algorithm

34

RWTH AACHEN UNIVERSITY

Recap: Backpropagation Algorithm

- More general formulation (used in deep learning packages)
 - Convert the network into a computational graph.
 - Perform reverse-mode-differentiation this graph
 - Each new layer/module just needs to specify how it affects the
 - forward pass $\mathbf{y} = \text{module.fprop}(\mathbf{x})$
 - backward pass $\frac{\partial E}{\partial \mathbf{x}} = \text{module.bprop}(\frac{\partial E}{\partial \mathbf{y}})$

\Rightarrow Very general framework, *any differentiable layer* can be used.

35

RWTH AACHEN UNIVERSITY

Recap: Supervised Learning

- Two main steps
 - Computing the gradients for each weight
 - Adjusting the weights in the direction of the gradient
- Gradient Descent: Basic update equation

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \Big|_{\mathbf{w}^{(\tau)}}$$
- Important considerations
 - On what data do we want to apply this? \Rightarrow Minibatches
 - How should we choose the step size η (the learning rate)?
 - More advanced optimizers (Momentum, RMSProp, Adam, ...)

36

RWTH AACHEN UNIVERSITY

Practical Considerations

- Vanishing gradients problem
 - In multilayer nets, gradients need to be propagated through many layers
 - The magnitudes of the gradients are often very different for the different layers, especially if the initial weights are small.
 - Gradients can get very small in the early layers of deep nets.
- When designing deep networks, we need to make sure gradients can be propagated throughout the network
 - By restricting the network depth (shallow networks are easier)
 - By very careful implementation (*numerics matter!*)
 - By choosing suitable nonlinearities (e.g., ReLU)
 - By performing proper initialization (Glorot, He)

37

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Recap: Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- Background: Deep Learning
 - Recap from ML lecture
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
 - ResNet

38

RWTH AACHEN UNIVERSITY

CNN Architectures: LeNet (1998)

- Early convolutional architecture
 - 2 Convolutional layers, 2 pooling layers
 - Fully-connected NN layers for classification
 - Successfully used for handwritten digit recognition (MNIST)

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

39

RWTH AACHEN UNIVERSITY

ImageNet Challenge 2012

- ImageNet
 - ~14M labeled internet images
 - 20k classes
 - Human labels via Amazon Mechanical Turk
- Challenge (ILSVRC)
 - 1.2 million training images
 - 1000 classes
 - Goal: Predict ground-truth class within top-5 responses
 - Currently one of the top benchmarks in Computer Vision

[Deng et al., CVPR'09]

40

RWTH AACHEN UNIVERSITY

CNN Architectures: 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.

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

41

RWTH AACHEN UNIVERSITY

ILSVRC 2012 Results

- AlexNet almost halved the error rate
 - 16.4% error (top-5) vs. 26.2% for the next best approach
 - ⇒ A revolution in Computer Vision
 - Acquired by Google in Jan '13, deployed in Google+ in May '13

42

RWTH AACHEN UNIVERSITY

AlexNet Results

43

RWTH AACHEN UNIVERSITY

AlexNet Results

Test image Retrieved images

44

CNN Architectures: VGGNet (2014/15)

Computer Vision Summer'19

K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

B. Leibe 45

CNN Architectures: 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 Configurations					
A	B	C	D	E	
11 weight layers	A-LRN 11 weight layers	13 weight layers	16 weight layers	19 weight layers	
conv-3-64	conv-3-64 LRN conv-3-64	conv-3-64 conv-3-64	conv-3-64 conv-3-64	conv-3-64 conv-3-64	
conv-3-128	conv-3-128 conv-3-128	conv-3-128 conv-3-128	conv-3-128 conv-3-128	conv-3-128 conv-3-128	
conv-3-256	conv-3-256 conv-3-256	conv-3-256 conv-3-256	conv-3-256 conv-3-256	conv-3-256 conv-3-256	
conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	
conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	conv-3-512 conv-3-512	
					Mainly used

Computer Vision Summer'19

B. Leibe 46

Comparison: AlexNet vs. VGGNet

- Receptive fields in the first layer
 - AlexNet: 11×11 , stride 4
 - Zeiler & Fergus: 7×7 , stride 2
 - VGGNet: 3×3 , stride 1
- Why that?
 - If you stack three 3×3 on top of another 3×3 layer, you effectively get a 5×5 receptive field.
 - With three 3×3 layers, the receptive field is already 7×7 .
 - But much fewer parameters: $3 \cdot 3^2 = 27$ instead of $7^2 = 49$.
 - In addition, non-linearities in-between 3×3 layers for additional discriminativity.

Computer Vision Summer'19

B. Leibe 47

CNN Architectures: GoogLeNet (2014)

- Main ideas
 - "Inception" module as modular component
 - Learns filters at several scales within each module

C. Szegedy, W. Liu, Y. Jia, et al, [Going Deeper with Convolutions](#), arXiv:1409.4842, 2014.

Computer Vision Summer'19

B. Leibe 48

GoogLeNet Visualization

Computer Vision Summer'19

B. Leibe 49

Results on ILSVRC

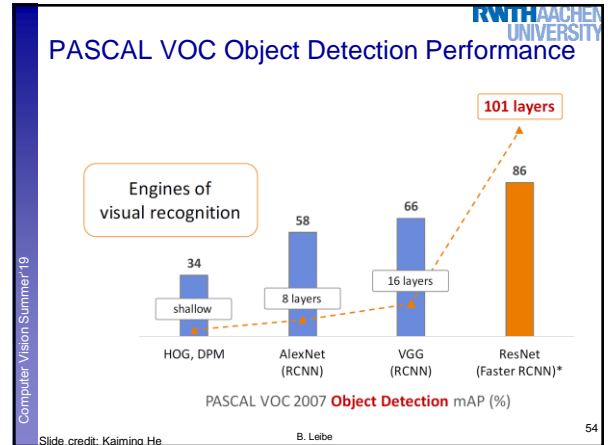
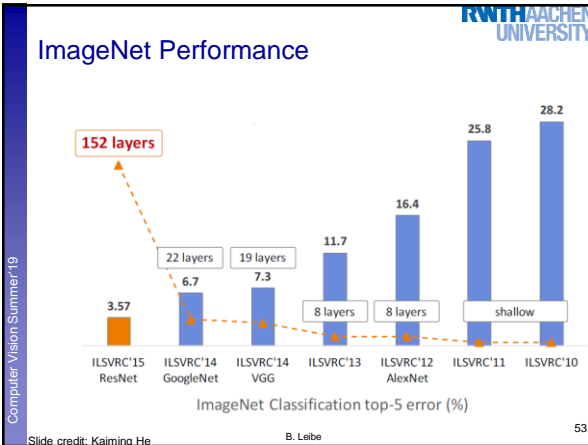
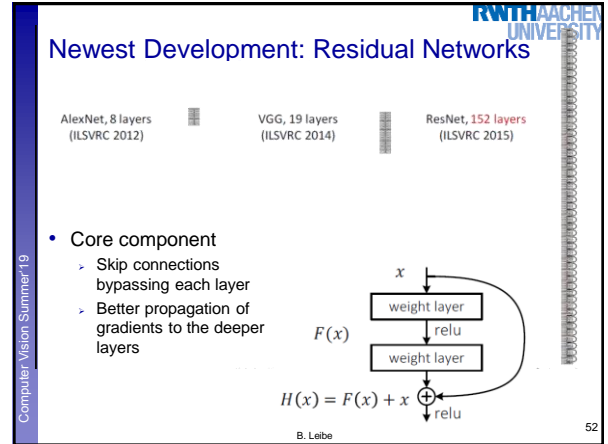
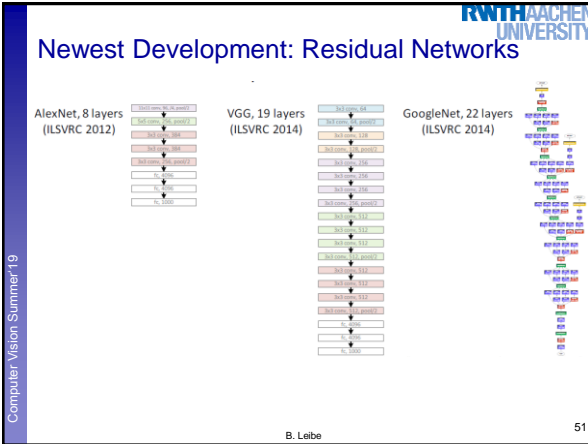
Method	top-1 val. error (%)	top-5 val. error (%)	top-5 test error (%)
VGG (2 nets, multi-crop & dense eval.)	23.7	6.8	6.8
VGG (1 net, multi-crop & dense eval.)	24.4	7.1	7.0
VGG (ILSVRC submission, 7 nets, dense eval.)	24.7	7.5	7.3
GoogLeNet (Szegedy et al., 2014) (1 net)	-	-	7.9
GoogLeNet (Szegedy et al., 2014) (7 nets)	-	-	6.7
MSRA (He et al., 2014) (11 nets)	-	-	8.1
MSRA (He et al., 2014) (1 net)	27.9	9.1	9.1
Clarifai (Russakovsky et al., 2014) (multiple nets)	-	-	11.7
Clarifai (Russakovsky et al., 2014) (1 net)	-	-	12.5
Zeiler & Fergus (Zeiler & Fergus, 2013) (6 nets)	36.0	14.7	14.8
Zeiler & Fergus (Zeiler & Fergus, 2013) (1 net)	37.5	16.0	16.1
OverFeat (Sermanet et al., 2014) (7 nets)	34.0	13.2	13.6
OverFeat (Sermanet et al., 2014) (1 net)	35.7	14.2	-
Krizhevsky et al. (Krizhevsky et al., 2012) (5 nets)	38.1	16.4	16.4
Krizhevsky et al. (Krizhevsky et al., 2012) (1 net)	40.7	18.2	-

- VGGNet and GoogLeNet perform at similar level
 - Comparison: human performance ~5% [Karpathy]

<http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/>

Computer Vision Summer'19

B. Leibe 50



RWTH AACHEN UNIVERSITY

References and Further Reading

- More information on Deep Learning and CNNs can be found in Chapters 6 and 9 of the Goodfellow & Bengio book

I. Goodfellow, Y. Bengio, A. Courville
Deep Learning
MIT Press, 2016
<http://www.deeplearningbook.org/>

B. Leibe

55