

RWTH AACHEN
UNIVERSITY

Advanced Machine Learning Lecture 13

Convolutional Neural Networks

15.12.2016

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

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

- Tricks of the Trade
 - Recap
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

Recap: Choosing the Right Learning Rate

- Convergence of Gradient Descent
 - Simple 1D example

$$W^{(\tau-1)} = W^{(\tau)} - \eta \frac{dE(W)}{dW}$$
 - What is the optimal learning rate η_{opt} ?
 - If E is quadratic, the optimal learning rate is given by the inverse of the Hessian

$$\eta_{opt} = \left(\frac{d^2 E(W^{(\tau)})}{dW^2} \right)^{-1}$$
 - Advanced optimization techniques try to approximate the Hessian by a simplified form.
 - If we exceed the optimal learning rate, bad things happen!

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

Recap: Advanced Optimization Techniques

- Momentum
 - Instead of using the gradient to change the *position* of the weight "particle", use it to change the *velocity*.
 - Effect: dampen oscillations in directions of high curvature
 - Nesterov-Momentum: Small variation in the implementation
- RMS-Prop
 - Separate learning rate for each weight: Divide the gradient by a running average of its recent magnitude.
- AdaGrad
- AdaDelta
- Adam

Some more recent techniques, work better for some problems. Try them.

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN
UNIVERSITY

Trick: Patience

- Saddle points dominate in high-dimensional spaces!

⇒ Learning often doesn't get stuck, you just may have to wait...

B. Leibe

Advanced Machine Learning Winter'16

RWTH AACHEN UNIVERSITY

Recap: Reducing the Learning Rate

- Final improvement step after convergence is reached
 - Reduce learning rate by a factor of 10.
 - Continue training for a few epochs.
 - Do this 1-3 times, then stop training.
- Effect
 - Turning down the learning rate will reduce the random fluctuations in the error due to different gradients on different minibatches.
- Be careful: Do not turn down the learning rate too soon!**
 - Further progress will be much slower after that.

7

Slide adapted from Geoff Hinton. B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Tricks of the Trade
 - Recap
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

8

B. Leibe

RWTH AACHEN UNIVERSITY

Neural Networks for Computer Vision

- How should we approach vision problems?

- Architectural considerations
 - Input is 2D ⇒ 2D layers of units
 - No pre-segmentation ⇒ Need robustness to misalignments
 - Vision is hierarchical ⇒ Hierarchical multi-layered structure
 - Vision is difficult ⇒ Network should be deep

9

B. Leibe

RWTH AACHEN UNIVERSITY

Why Hierarchical Multi-Layered Models?

- Motivation 1: Visual scenes are hierarchically organized

10

Slide adapted from Richard Turner. B. Leibe

RWTH AACHEN UNIVERSITY

Why Hierarchical Multi-Layered Models?

- Motivation 2: *Biological vision* is hierarchical, too

11

Slide adapted from Richard Turner. B. Leibe

Inferotemporal cortex
V4: different textures
V1: simple and complex cells
Photoreceptors, retina

RWTH AACHEN UNIVERSITY

Inspiration: Neuron Cells

12

Slide credit: Svetlana Lazebnik, Rob Fergus. B. Leibe

RWTH AACHEN UNIVERSITY

Hubel/Wiesel Architecture

- D. Hubel, T. Wiesel (1959, 1962, Nobel Prize 1981)
 - Visual cortex consists of a hierarchy of *simple*, *complex*, and *hyper-complex* cells

Hubel & Wiesel

topographical mapping

featural hierarchy

hyper-complex cells
complex cells
simple cells

high level
 mid level
 low level

Slide credit: Svetlana Lazebnik, Rob Fergus. B. Leibe 13

RWTH AACHEN UNIVERSITY

Why Hierarchical Multi-Layered Models?

- Motivation 3: Shallow architectures are inefficient at representing complex functions

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.

Slide adapted from Richard Turner. B. Leibe 14

RWTH AACHEN UNIVERSITY

What's Wrong With Standard Neural Networks?

- Complexity analysis
 - How many parameters does this network have?

$$|\theta| = 3D^2 + D$$
 - For a small 32x32 image

$$|\theta| = 3 \cdot 32^4 + 32^2 \approx 3 \cdot 10^6$$
- Consequences
 - Hard to train
 - Need to initialize carefully
 - Convolutional nets reduce the number of parameters!

Slide adapted from Richard Turner. B. Leibe 15

RWTH AACHEN UNIVERSITY

Convolutional Neural Networks (CNN, ConvNet)

- 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 16

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition

- Fully connected network
 - E.g. 1000x1000 image
 - 1M hidden units
 - ⇒ 1T parameters!
- Ideas to improve this
 - Spatial correlation is local

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann LeCun 17

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition

- Locally connected net
 - E.g. 1000x1000 image
 - 1M hidden units
 - 10x10 receptive fields
 - ⇒ 100M parameters!
- Ideas to improve this
 - Spatial correlation is local
 - Want translation invariance

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann LeCun 18

Advanced Machine Learning Winter'16

Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

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

Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann LeCun 19

Advanced Machine Learning Winter'16

Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- 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!

Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann LeCun 20

Advanced Machine Learning Winter'16

Important Conceptual Shift

RWTH AACHEN UNIVERSITY

- Before
- Now:

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe 21

Advanced Machine Learning Winter'16

Convolution Layers

RWTH AACHEN UNIVERSITY

Example image: $32 \times 32 \times 3$ volume

Before: Full connectivity $32 \times 32 \times 3$ weights

Now: Local connectivity
One neuron connects to, e.g., $5 \times 5 \times 3$ region.
⇒ Only $5 \times 5 \times 3$ shared weights.

- Note: Connectivity is
 - Local in space (5×5 inside 32×32)
 - But full in depth (all 3 depth channels)

Slide adapted from FeiFei Li, Andrei Karpathy B. Leibe 22

Advanced Machine Learning Winter'16

Convolution Layers

RWTH AACHEN UNIVERSITY

before: "hidden layer of 200 neurons"
now: "output volume of depth 200"

- All Neural Net activations arranged in 3 dimensions
 - Multiple neurons all looking at the same input region, stacked in depth

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

Advanced Machine Learning Winter'16

Convolution Layers

RWTH AACHEN UNIVERSITY

Naming convention:

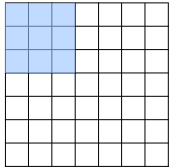
HEIGHT
WIDTH
DEPTH

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

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe 24

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

Advanced Machine Learning Winter'16

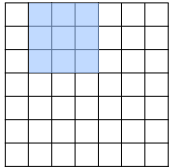
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

26

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

Advanced Machine Learning Winter'16

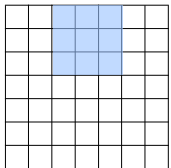
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

27

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

Advanced Machine Learning Winter'16

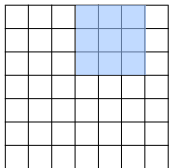
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

28

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

Advanced Machine Learning Winter'16

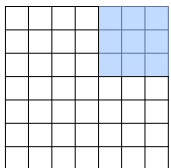
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

29

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1
⇒ 5x5 output

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

Advanced Machine Learning Winter'16

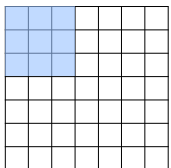
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

30

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1
⇒ 5x5 output

What about stride 2?

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

Advanced Machine Learning Winter'16

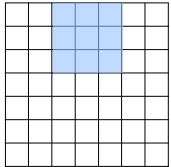
Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

31

RWTH AACHEN UNIVERSITY

Convolution Layers



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

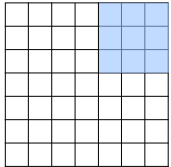
What about stride 2?

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

Advanced Machine Learning Winter'16 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 32

RWTH AACHEN UNIVERSITY

Convolution Layers



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

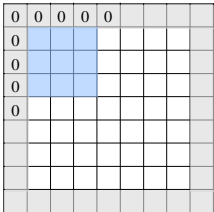
What about stride 2?
 $\Rightarrow 3 \times 3$ output

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

Advanced Machine Learning Winter'16 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 33

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.

Advanced Machine Learning Winter'16 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 34

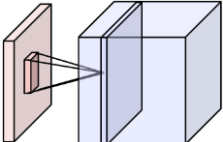
RWTH AACHEN UNIVERSITY

Activation Maps of Convolutional Filters



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

5x5 filters




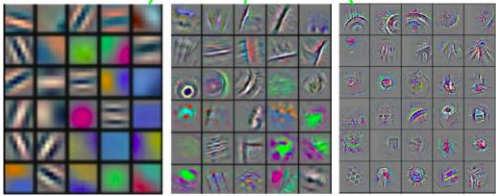
Each activation map is a depth slice through the output volume.

Activation maps

Advanced Machine Learning Winter'16 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 35

RWTH AACHEN UNIVERSITY

Effect of Multiple Convolution Layers

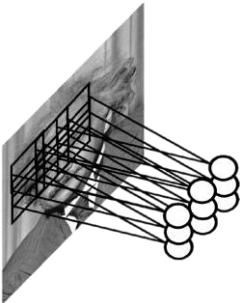



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Advanced Machine Learning Winter'16 | Slide credit: Yann LeCun | B. Leibe | 36

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition



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

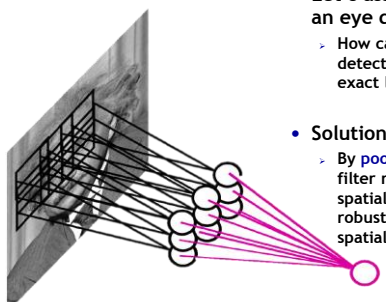
Advanced Machine Learning Winter'16 | Slide adapted from Marc'Aurelio Ranzato | B. Leibe | Image source: Yann LeCun | 37

Advanced Machine Learning Winter'16

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. 38

Advanced Machine Learning Winter'16

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. 39

Advanced Machine Learning Winter'16

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. 40

Advanced Machine Learning Winter'16

CNNs: Implication for Back-Propagation

RWTH AACHEN UNIVERSITY

- Convolutional layers
 - Filter weights are shared between locations
 - ⇒ Gradients are added for each filter location.

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 41

Advanced Machine Learning Winter'16

Topics of This Lecture

RWTH AACHEN UNIVERSITY

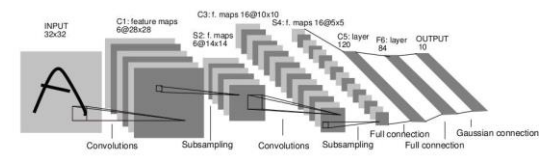
- Tricks of the Trade
 - Recap
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

B. Leibe. 42

Advanced Machine Learning Winter'16

CNN Architectures: LeNet (1998)

RWTH AACHEN UNIVERSITY



- 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. B. Leibe. 43

Advanced Machine Learning Winter'16

ImageNet Challenge 2012

RWTH AACHEN UNIVERSITY

- 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



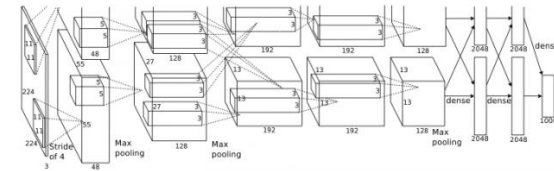
[Deng et al., CVPR'09]

B. Leibe 44

Advanced Machine Learning Winter'16

CNN Architectures: AlexNet (2012)

RWTH AACHEN UNIVERSITY



- 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.

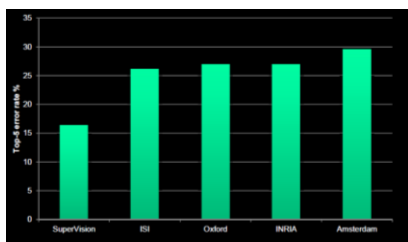
Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012

B. Leibe 45

Advanced Machine Learning Winter'16

ILSVRC 2012 Results

RWTH AACHEN UNIVERSITY



Team	Top-5 error rate %
SuperVision	~16.4
ISI	~26.2
Oxford	~26.2
NIPS	~26.2
Amsterdam	~26.2

- 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

B. Leibe 46

Advanced Machine Learning Winter'16

AlexNet Results

RWTH AACHEN UNIVERSITY

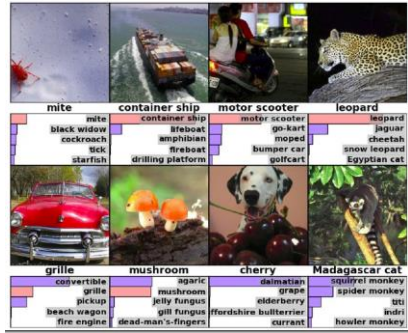



Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012

B. Leibe 48

Advanced Machine Learning Winter'16

AlexNet Results

RWTH AACHEN UNIVERSITY



Test image Retrieved images

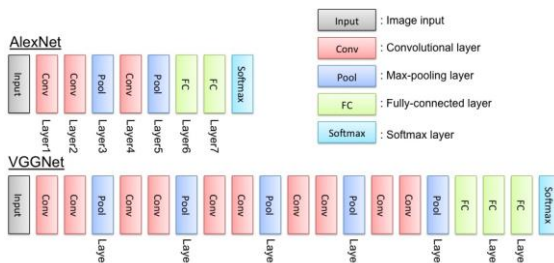
Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012

B. Leibe 49

Advanced Machine Learning Winter'16

CNN Architectures: VGGNet (2014/15)

RWTH AACHEN UNIVERSITY



Legend:

- Input: Image input
- Conv: Convolutional layer
- Pool: Max-pooling layer
- FC: Fully-connected layer
- Softmax: Softmax layer

AlexNet: Input → Conv → Pool → Conv → Pool → Conv → Pool → FC → Softmax

VGGNet: Input → Conv → Pool → Conv → Pool → Conv → Pool → Conv → Pool → Conv → Pool → Conv → Pool → FC → FC → FC → Softmax

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

Image source: Hirokatsu Kobatake

B. Leibe 50

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%.

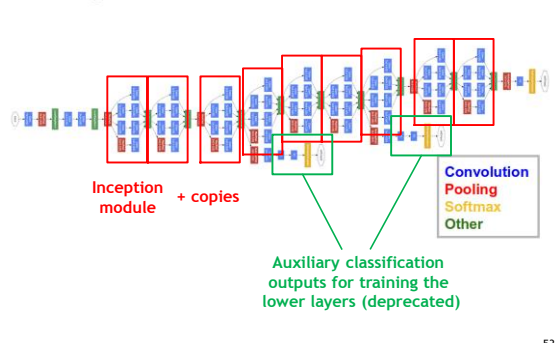
ConvNet Configuration

A	A-LRN	B	C	D	E
input: 224 × 224 RGB image					
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
conv3-64	conv3-64 LRN	conv3-64 conv3-64	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	conv3-128 conv3-128
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 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

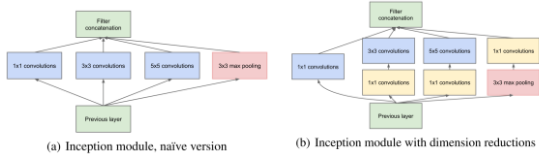
Mainly used

B. Leibe 51
Image source: Simonyan & Zisserman

GoogLeNet Visualization



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.

B. Leibe 52

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	-

B. Leibe 54
Image source: Simonyan & Zisserman

References and Further Reading

- LeNet
 - 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.
 - AlexNet
 - A. Krizhevsky, I. Sutskever, and G. Hinton, *ImageNet Classification with Deep Convolutional Neural Networks*, NIPS 2012.
 - VGGNet
 - K. Simonyan, A. Zisserman, *Very Deep Convolutional Networks for Large-Scale Image Recognition*, ICLR 2015
 - GoogLeNet
 - C. Szegedy, W. Liu, Y. Jia, et al., *Going Deeper with Convolutions*, arXiv:1409.4842, 2014.
- B. Leibe 58

References

- ReLU
 - X. Glorot, A. Bordes, Y. Bengio, *Deep sparse rectifier neural networks*, AISTATS 2011.
 - Batch Normalization
 - S. Ioffe, C. Szegedy, *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*, arXiv 1502.03167, 2015.
- B. Leibe 59