

RWTH AACHEN  
UNIVERSITY

# Machine Learning – Lecture 16

## Convolutional Neural Networks II

18.12.2019

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

Machine Learning Winter '19

RWTH AACHEN  
UNIVERSITY

## Course Outline

- Fundamentals
  - Bayes Decision Theory
  - Probability Density Estimation
- Classification Approaches
  - Linear Discriminants
  - Support Vector Machines
  - Ensemble Methods & Boosting
  - Random Forests
- Deep Learning
  - Foundations
  - Convolutional Neural Networks
  - Recurrent Neural Networks

B. Leibe

RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

- Recap: CNNs
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNets
- Visualizing CNNs
  - Visualizing CNN features
  - Visualizing responses
  - Visualizing learned structures
- Applications

B. Leibe

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.

Slide credit: Svetlana Lazebnik

B. Leibe

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
  - ⇒ only 10k parameters
- Result: Response map
  - size: 1000×1000×100
  - Only memory, not params!

Slide adapted from Marc'Aurelio Ranzato

B. Leibe

RWTH AACHEN  
UNIVERSITY

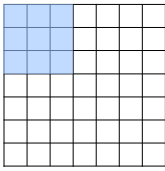
## Recap: Convolution Layers

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

Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

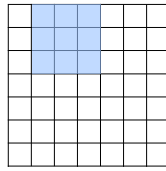
### Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1

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

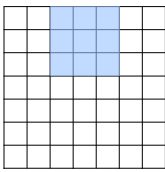
### Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1

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

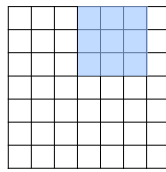
### Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1

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

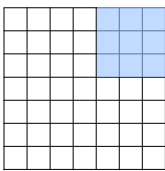
### Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1

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

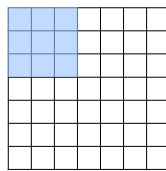
### Convolution Layers



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

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

### Convolution Layers



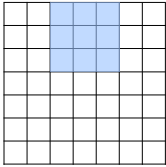
Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?

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

RWTH AACHEN UNIVERSITY

## Convolution Layers



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

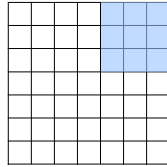
What about stride 2?

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

Machine Learning Winter '19 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 15

RWTH AACHEN UNIVERSITY

## Convolution Layers



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

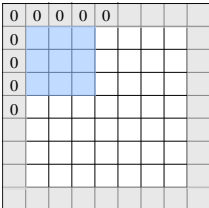
What about stride 2?  
 $\Rightarrow 3 \times 3$  output

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

Machine Learning Winter '19 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 16

RWTH AACHEN UNIVERSITY

## Convolution Layers



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

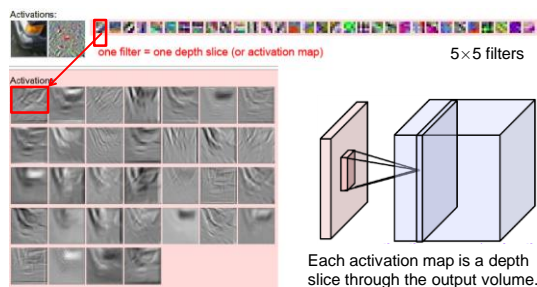
What about stride 2?  
 $\Rightarrow 3 \times 3$  output

- 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.

Machine Learning Winter '19 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 17

RWTH AACHEN UNIVERSITY

## Activation Maps of Convolutional Filters



Activations:  
 one filter = one depth slice (or activation map)  
 $5 \times 5$  filters

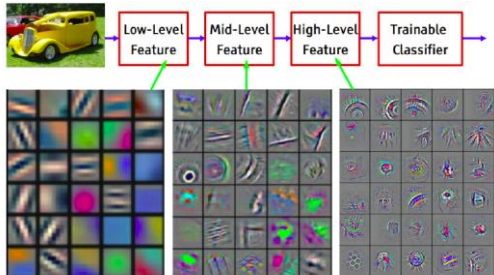
Activation maps

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

Machine Learning Winter '19 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 18

RWTH AACHEN UNIVERSITY

## Effect of Multiple Convolution Layers

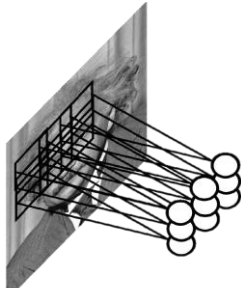


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Machine Learning Winter '19 | Slide credit: Yann LeCun | B. Leibe | 19

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?

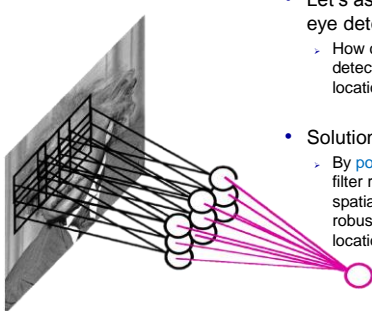
Machine Learning Winter '19 | Slide adapted from Marc'Aurelio Ranzato | B. Leibe | Image source: Yann LeCun | 20

Machine Learning Winter '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.

21

Machine Learning Winter '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.

24

Machine Learning Winter '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.

25

Machine Learning Winter '19

## CNNs: Implication for Back-Propagation

RWTH AACHEN UNIVERSITY

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

Machine Learning Winter '19

B. Leibe.

26

Machine Learning Winter '19

## Topics of This Lecture

RWTH AACHEN UNIVERSITY

- Recap: CNNs
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNet
- Visualizing CNNs
  - Visualizing CNN features
  - Visualizing responses
  - Visualizing learned structures
- Applications

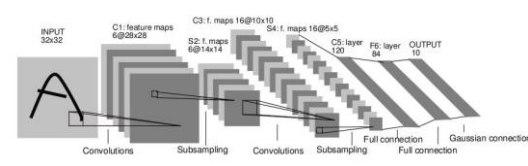
B. Leibe.

27

Machine Learning Winter '19

## 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.

28

Machine Learning Winter '19

## 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

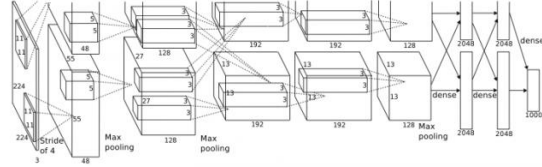


[Deng et al., CVPR'09]

B. Leibe 29

Machine Learning Winter '19

## CNN Architectures: AlexNet (2012)



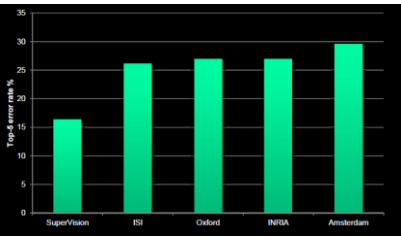
- Similar framework as LeNet, but
  - Bigger model (7 hidden layers, 650k units, 60M parameters)
  - More data ( $10^6$  images instead of  $10^3$ )
  - GPU implementation
  - Better regularization and up-to-date tricks for training (Dropout)

A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.

Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012. 30

Machine Learning Winter '19

## ILSVRC 2012 Results



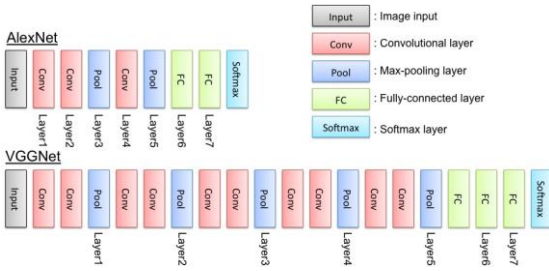
Team	Top-5 Error Rate (%)
SuperVision	~16.4
ISI	~26.2
Oxford	~26.2
INRIA	~26.2
Amsterdam	~26.2

- 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 31

Machine Learning Winter '19

## CNN Architectures: VGGNet (2014/15)



Legend:

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

AlexNet: Input → Conv Layer1 → Conv Layer2 → Pool Layer3 → Conv Layer4 → Conv Layer5 → FC Layer6 → FC Layer7 → Softmax

VGGNet: Input → Conv Layer1 → Conv Layer2 → Pool Layer3 → Conv Layer4 → Conv Layer5 → Pool Layer6 → Conv Layer7 → Pool Layer8 → FC Layer9 → FC Layer10 → FC Layer11 → Softmax

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

Image source: Hirokatsu Kataoka 33

Machine Learning Winter '19

## 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	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)				
conv-3-64	conv-3-64	conv-3-64	conv-3-64	conv-3-64
maxpool				
conv-3-128	conv-3-128	conv-3-128	conv-3-128	conv-3-128
maxpool				
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-1-256	conv-3-256
maxpool				
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-1-512	conv-3-512
maxpool				
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
maxpool				
FC-4096				
FC-4096				
FC-1000				
soft-max				

Mainly used

B. Leibe 34

Machine Learning Winter '19

## Comparison: AlexNet vs. VGGNet

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

B. Leibe 35

RWTH AACHEN UNIVERSITY

## CNN Architectures: GoogLeNet (2014/2015)

(a) Inception module, naïve version      (b) Inception module with dimension reductions

- 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, CVPR'15, 2015.

B. Leibe 36

RWTH AACHEN UNIVERSITY

## GoogLeNet Visualization

Inception module + copies

Auxiliary classification outputs for training the lower layers (deprecated)

Convolution  
Pooling  
Softmax  
Other

B. Leibe 37

RWTH AACHEN UNIVERSITY

## 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
Zeller & Fergus (Zeller & Fergus, 2013) (6 nets)	36.0	14.7	14.8
Zeller & Fergus (Zeller & 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/>

B. Leibe 38

RWTH AACHEN UNIVERSITY

## Newer Developments: Residual Networks

AlexNet, 8 layers (ILSVRC 2012)

VGG, 19 layers (ILSVRC 2014)

GoogLeNet, 22 layers (ILSVRC 2014)

B. Leibe 39

RWTH AACHEN UNIVERSITY

## Newer Developments: 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
  - We'll analyze this mechanism in more detail later...

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

B. Leibe 40

RWTH AACHEN UNIVERSITY

## ImageNet Performance

152 layers

22 layers    19 layers    11.7    16.4    25.8    28.2

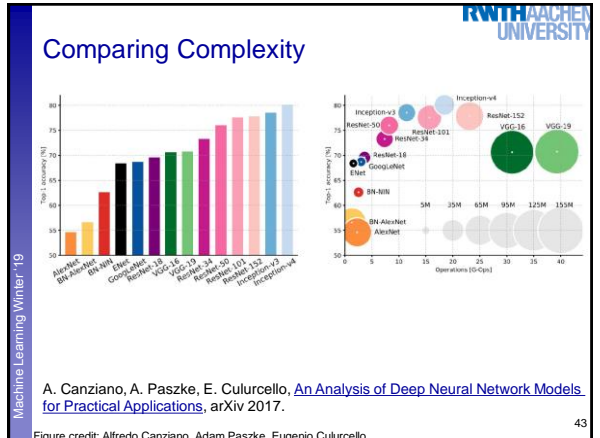
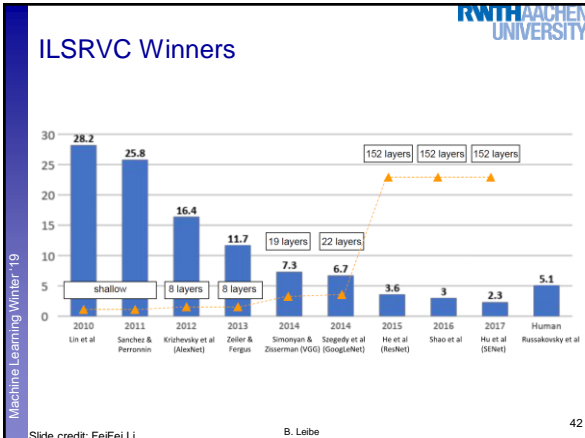
8 layers    8 layers    shallow

3.57    6.7    7.3

ILSVRC'15 ResNet    ILSVRC'14 GoogLeNet    ILSVRC'14 VGG    ILSVRC'13    ILSVRC'12 AlexNet    ILSVRC'11    ILSVRC'10

ImageNet Classification top-5 error (%)

B. Leibe 41



### Understanding the ILSRVC Challenge

- Imagine the scope of the problem!
  - 1000 categories
  - 1.2M training images
  - 50k validation images
- This means...
  - Speaking out the list of category names at 1 word/s... takes 15mins.
  - Watching a slideshow of the validation images at 2s/image... takes a full day (24h+).
  - Watching a slideshow of the training images at 2s/image... takes a full month.

B. Leibe

B. Leibe

### More Finegrained Classes

B. Leibe

Image source: O. Russakovsky et al

### Quirks and Limitations of the Data Set

- Generated from WordNet ontology
  - Some animal categories are overrepresented
  - E.g., 120 subcategories of dog breeds

⇒ 6.7% top-5 error looks all the more impressive

B. Leibe

Image source: A. Krizhevsky

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Recap: CNNs
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNets
- Visualizing CNNs
  - Visualizing CNN features
  - Visualizing responses
  - Visualizing learned structures
- Applications

Machine Learning Winter '19  
B. Leibe 48

RWTH AACHEN UNIVERSITY

## Visualizing CNNs

Machine Learning Winter '19  
Image source: M. Zeiler, R. Fergus 49

RWTH AACHEN UNIVERSITY

## Visualizing CNNs

reconstruction of image patches from that unit (indicates aspect of patches which unit is sensitive to)

top 9 image patches that cause maximal activation in layer 2 unit

M. Zeiler, R. Fergus, [Visualizing and Understanding Convolutional Neural Networks](#), ECCV 2014.

Machine Learning Winter '19  
B. Leibe  
Image source: M. Zeiler, R. Fergus 50

RWTH AACHEN UNIVERSITY

## Visualizing CNNs

Machine Learning Winter '19  
B. Leibe  
Image source: M. Zeiler, R. Fergus 51

RWTH AACHEN UNIVERSITY

## Visualizing CNNs

Machine Learning Winter '19  
B. Leibe  
Image source: M. Zeiler, R. Fergus 52

RWTH AACHEN UNIVERSITY

## What Does the Network React To?

- Occlusion Experiment
  - Mask part of the image with an occluding square.
  - Monitor the output


Machine Learning Winter '19  
B. Leibe 53



Machine Learning Winter '19

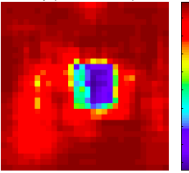
### What Does the Network React To?

Input image




True Label: Pomeranian

$p(\text{True class})$



Most probable class

- Pomeranian
- Tennis ball
- Keeshond
- Beleisene



Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

54

Machine Learning Winter '19

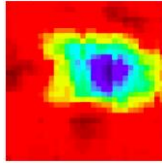
### What Does the Network React To?

Input image

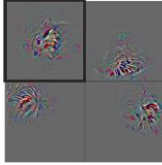


True Label: Pomeranian

Total activation in most active 5<sup>th</sup> layer feature map



Other activations from the same feature map.



Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

55

Machine Learning Winter '19

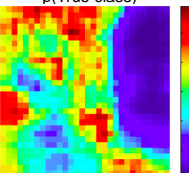
### What Does the Network React To?

Input image



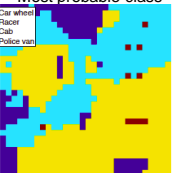
True Label: Car Wheel

$p(\text{True class})$



Most probable class

- Car wheel
- Racer
- Cab
- Police van



Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

56

Machine Learning Winter '19

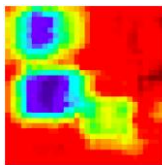
### What Does the Network React To?

Input image




True Label: Car Wheel

Total activation in most active 5<sup>th</sup> layer feature map



Other activations from the same feature map.



Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

57

Machine Learning Winter '19

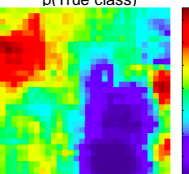
### What Does the Network React To?

Input image



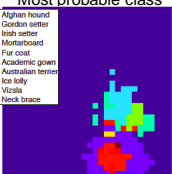
True Label: Afghan Hound

$p(\text{True class})$



Most probable class

- Afghan hound
- Cardigan setter
- Irish setter
- Montarboard
- Fur coat
- Academic gown
- Australian terrier
- Ice lolly
- Vizsla
- Neck brace



Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

58

Machine Learning Winter '19

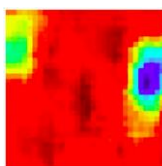
### What Does the Network React To?

Input image




True Label: Afghan Hound

Total activation in most active 5<sup>th</sup> layer feature map



Other activations from the same feature map.



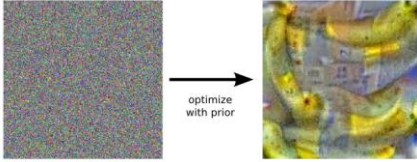
Slide credit: Svetlana Lazebnik, Rob Fergus

Image source: M. Zeiler, R. Fergus

59

Machine Learning Winter '19

## Inceptionism: Dreaming ConvNets



- Idea
  - Start with a random noise image.
  - Enhance the input image such as to enforce a particular response (e.g., banana).
  - Combine with prior constraint that image should have similar statistics as natural images.
- ⇒ Network hallucinates characteristics of the learned class.

<http://googleresearch.blogspot.de/2015/06/inceptionism-going-deeper-into-neural.html>

RWTH AACHEN UNIVERSITY

Machine Learning Winter '19

## Inceptionism: Dreaming ConvNets

- Results



<http://googleresearch.blogspot.de/2015/07/deepdream-code-example-for-visualizing.html>

RWTH AACHEN UNIVERSITY

Machine Learning Winter '19

## Inceptionism: Dreaming ConvNets



<https://www.youtube.com/watch?v=IREsx-xWQ0g>

RWTH AACHEN UNIVERSITY

Machine Learning Winter '19

## Topics of This Lecture

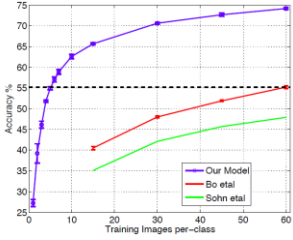
- Recap: CNNs
- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNets
- Visualizing CNNs
  - Visualizing CNN features
  - Visualizing responses
  - Visualizing learned structures
- Applications

B. Leibe

RWTH AACHEN UNIVERSITY

Machine Learning Winter '19

## The Learned Features are Generic



- 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

B. Leibe

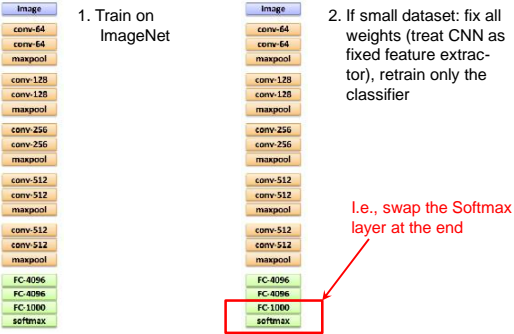
Image source: M. Zeiler, B. Fergus

RWTH AACHEN UNIVERSITY

Machine Learning Winter '19

## 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



B. Leibe

Slide credit: Andrei Karpathy

RWTH AACHEN UNIVERSITY

### Transfer Learning with CNNs

1. Train on ImageNet

- conv-64
- conv-64
- maxpool
- conv-128
- conv-128
- maxpool
- conv-256
- conv-256
- maxpool
- conv-512
- conv-512
- maxpool
- conv-512
- conv-512
- maxpool
- FC-4096
- FC-4096
- FC-1000
- softmax

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.

- conv-64
- conv-64
- maxpool
- conv-128
- conv-128
- maxpool
- conv-256
- conv-256
- maxpool
- conv-512
- conv-512
- maxpool
- conv-512
- conv-512
- maxpool
- FC-4096
- FC-4096
- FC-1000
- softmax

Retrain bigger portion of the network

Slide credit: Andrei Karpathy. B. Leibe. 66

### Other Tasks: Detection

#### 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

Machine Learning Winter '19. 67

### 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.
  - mAP improved to >70%

Slide credit: Ross Girshick. 68

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

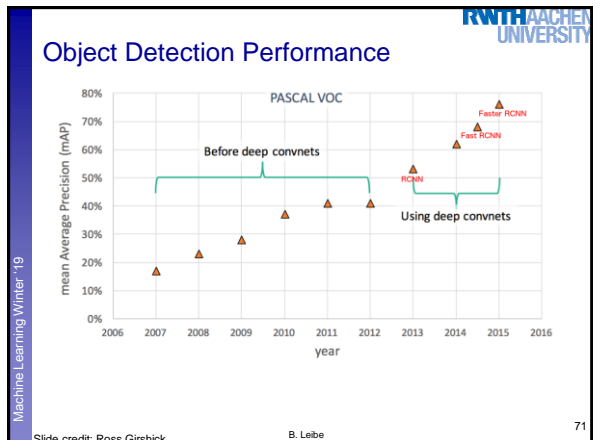
K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

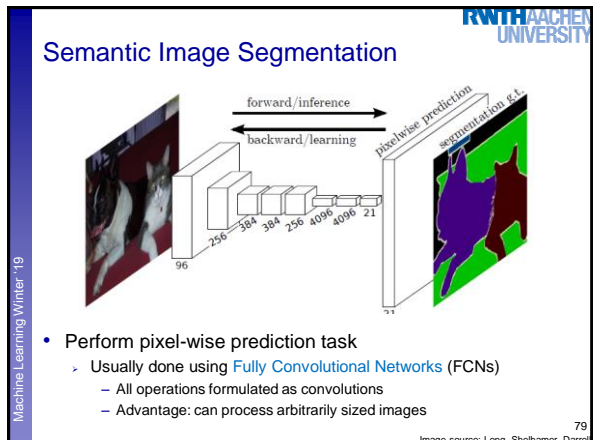
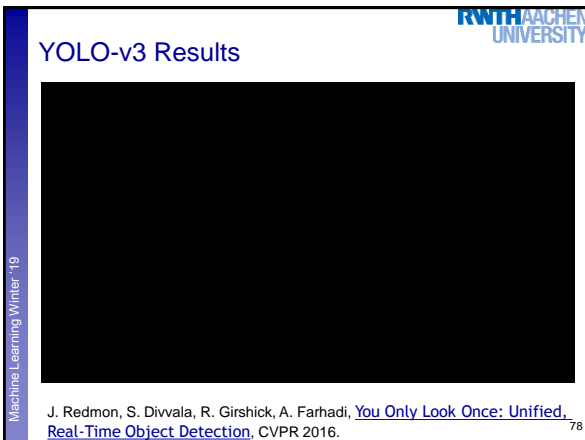
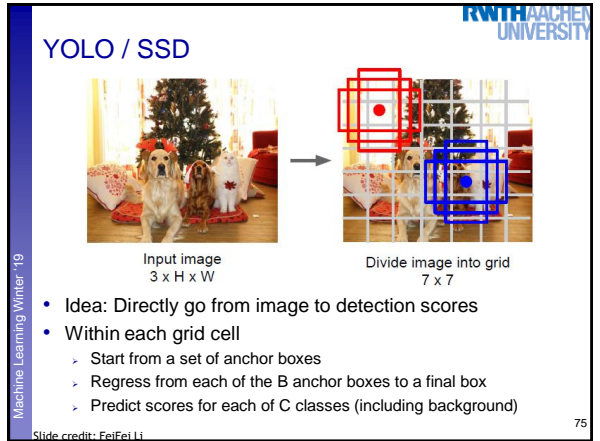
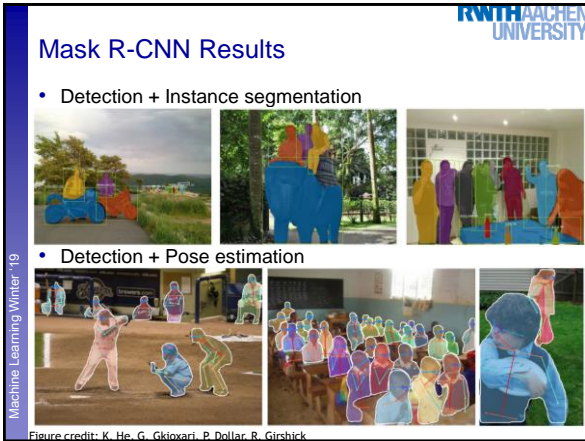
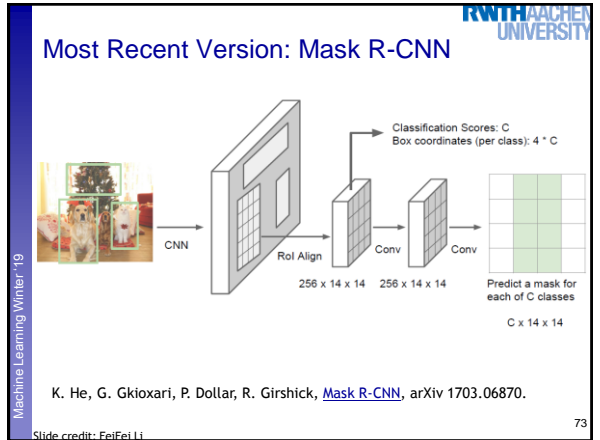
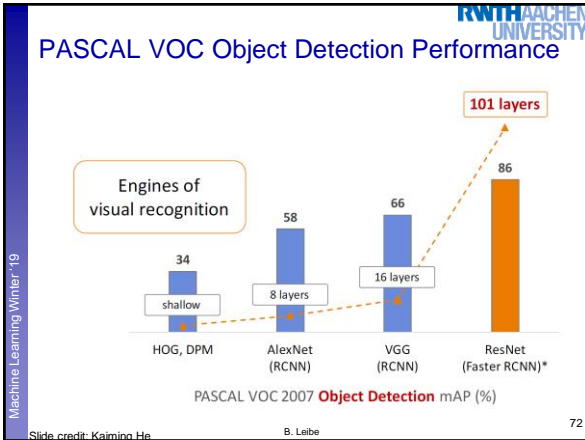
Machine Learning Winter '19. B. Leibe. 69

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

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

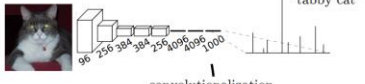
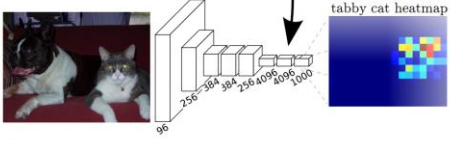
Machine Learning Winter '19. B. Leibe. 70





Machine Learning Winter '19

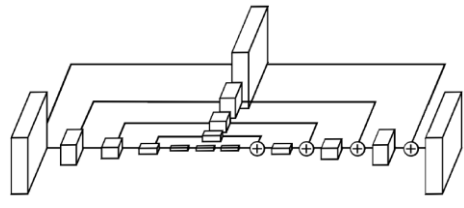
## CNNs vs. FCNs

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

80  
Image source: Long, Shelhamer, Darrell

Machine Learning Winter '19

## 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

81  
Image source: Newell et al.

Machine Learning Winter '19

## Semantic Segmentation

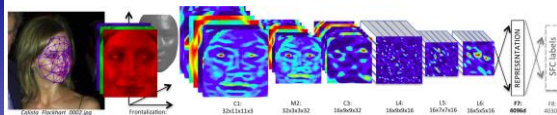
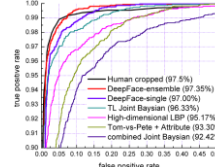


- Current state-of-the-art
  - Based on an extension of ResNets

82  
[Pohlen, Hermans, Mathias, Leibe, CVPR 2017]

Machine Learning Winter '19


## Other Tasks: Face Verification

83  
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

Machine Learning Winter '19

## 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...

84  
B. Leibe  
Image source: clarifai.com

Machine Learning Winter '19

## 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.

85  
B. Leibe

## References and Further Reading

- ResNets
  - K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.
  - A. Veit, M. Wilber, S. Belongie, [Residual Networks Behave Like Ensembles of Relatively Shallow Networks](#), NIPS 2016.

## 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.

## 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.