

RWTH AACHEN
UNIVERSITY

Computer Vision - Lecture 16

Deep Learning for Object Categorization

14.01.2016

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

Computer Vision WS 15/16

RWTH AACHEN
UNIVERSITY

Announcements

- Seminar registration period starts today
 - We will offer a seminar in the summer semester "Current Topics in Computer Vision and Machine Learning"
 - Block seminar, presentations at beginning of semester break
 - If you're interested, you can register at <http://www.graphics.rwth-aachen.de/apse>
 - Registration period: 14.01.2016 - 27.01.2016
 - *Quick poll: Who would be interested in that?*

2

Computer Vision WS 15/16

RWTH AACHEN
UNIVERSITY

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

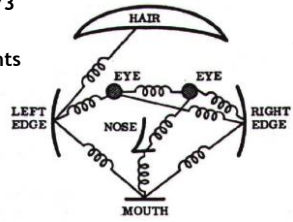
3

Computer Vision WS 15/16

RWTH AACHEN
UNIVERSITY

Recap: Part-Based Models

- Fischler & Elschlager 1973
- Model has two components
 - parts (2D image fragments)
 - structure (configuration of parts)

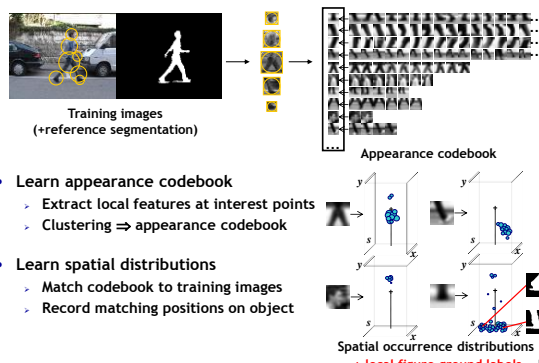


4

Computer Vision WS 15/16

RWTH AACHEN
UNIVERSITY

Recap: Implicit Shape Model - Representation



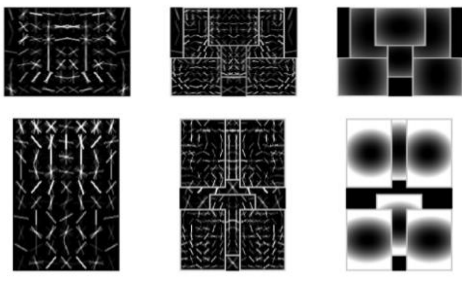
- Learn appearance codebook
 - Extract local features at interest points
 - Clustering ⇒ appearance codebook
- Learn spatial distributions
 - Match codebook to training images
 - Record matching positions on object

5

Computer Vision WS 15/16

RWTH AACHEN
UNIVERSITY

Recap: Deformable Part-Based Model



Root filters
coarse resolution

Part filters
finer resolution

Deformation models

6

Computer Vision WS 15/16

RWTH AACHEN UNIVERSITY

Recap: Object Hypothesis

Score of filter: dot product of filter with HOG features underneath it

Score of object hypothesis is sum of filter scores minus deformation costs

- Multiscale model captures features at two resolutions

Computer Vision WS 15/16 | Slide credit: Pedro Felzenszwalb | B. Leibe | 7

RWTH AACHEN UNIVERSITY

Recap: Score of a Hypothesis

$$\text{score}(p_0, \dots, p_n) = \sum_{i=0}^n F_i \cdot \phi(H, p_i) - \sum_{i=1}^n d_i \cdot (dx_i^2, dy_i^2)$$

“data term” (filters) “spatial prior” (displacements)

$$\text{score}(z) = \beta \cdot \Psi(H, z)$$

concatenation filters and deformation parameters concatenation of HOG features and part displacement features

Computer Vision WS 15/16 | Slide credit: Pedro Felzenszwalb | B. Leibe | 8

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Deep Learning
 - Motivation
- Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Applications

Computer Vision WS 15/16 | B. Leibe | 9

RWTH AACHEN UNIVERSITY

We've finally got there!

Deep Learning

Computer Vision WS 15/16 | B. Leibe | 10

RWTH AACHEN UNIVERSITY

Traditional Recognition Approach

- Characteristics
 - Features are not learned, but engineered
 - Trainable classifier is often generic (e.g., SVM)
 - ⇒ Many successes in 2000-2010.

Computer Vision WS 15/16 | Slide credit: Svetlana Lazebnik | B. Leibe | 11

RWTH AACHEN UNIVERSITY

Traditional Recognition Approach

- Features are key to recent progress in recognition
 - Multitude of hand-designed features currently in use
 - SIFT, HOG,
 - ⇒ Where next? Better classifiers? Or keep building more features?

Computer Vision WS 15/16 | Slide credit: Svetlana Lazebnik | 12

What About Learning the Features?

- Learn a *feature hierarchy* all the way from pixels to classifier
 - Each layer extracts features from the output of previous layer
 - Train all layers jointly

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik B. Leibe 13

“Shallow” vs. “Deep” Architectures

Traditional recognition: “Shallow” architecture

Deep learning: “Deep” architecture

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik B. Leibe 14

Background: Perceptrons

Input

Weights

Output: $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

Sigmoid function

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik 15

Inspiration: Neuron Cells

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik, Rob Fergus 16

Background: Multi-Layer Neural Networks

- Nonlinear classifier**
 - Training:** find network weights \mathbf{w} to minimize the error between true training labels t_n and estimated labels $f_{\mathbf{w}}(x_n)$:

$$E(\mathbf{W}) = \sum_n L(t_n, f(\mathbf{x}_n; \mathbf{W}))$$
 - Minimization can be done by gradient descent provided f is differentiable
 - Training method: **back-propagation.**

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik B. Leibe 17

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

feature hierarchy

- hyper-complex cells
- complex cells
- simple cells

- high level
- mid level
- low level

Computer Vision WS 15/16

Slide credit: Svetlana Lazebnik, Rob Fergus B. Leibe 18

RWTH AACHEN UNIVERSITY

Convolutional Neural Networks (CNN, ConvNet)

INPUT 32x32 C1: feature maps 6@28x28 C3: feature maps 16@10x10 S4: feature maps 16@5x5 C5: feature maps 120 FC: layer 84 OUTPUT 10

Convolutions Subsampling Convolutions Subsampling Full connection Full connection Gaussian connections

- 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.

Computer Vision WS 15/16 Slide credit: Svetlana Lazebnik B. Leibe 19

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Deep Learning
 - Motivation
- Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Applications

Computer Vision WS 15/16 B. Leibe 20

RWTH AACHEN UNIVERSITY

Convolutional Networks: Structure

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

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

Computer Vision WS 15/16 Slide credit: Svetlana Lazebnik B. Leibe 21

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

Computer Vision WS 15/16 Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann LeCun 22

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

Computer Vision WS 15/16 Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann LeCun 23

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition

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

Computer Vision WS 15/16 Slide adapted from Marc'Aurelio Ranzato B. Leibe Image source: Yann LeCun 24

Computer Vision WS 15/16

Convolutional Networks: Intuition

- Convolutional net
 - Share the same parameters across different locations
 - Convolutions with learned kernels
- Learn *multiple* filters
 - E.g. 1000x1000 image
 - 100 filters
 - 10x10 filter size
 - ⇒ 10k parameters
- Result: Response map
 - size: 1000x1000x100
 - Only memory, not params!

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 25

Computer Vision WS 15/16

Important Conceptual Shift

- Before
- Now:

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe. 26

Computer Vision WS 15/16

Convolution Layers

Example image: 32x32x3 volume

Before: Full connectivity
32x32x3 weights

Now: Local connectivity
One neuron connects to, e.g., 5x5x3 region.
⇒ Only 5x5x3 shared weights.

- Note: Connectivity is
 - Local in space (5x5 inside 32x32)
 - But full in depth (all 3 depth channels)

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 27

Computer Vision WS 15/16

Convolution Layers

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. 28

Computer Vision WS 15/16

Convolution Layers

Naming convention:

- All Neural Net activations arranged in 3 dimensions
 - Multiple neurons all looking at the same input region, stacked in depth
 - Form a single [1x1xdepth] depth column in output volume.

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe. 29

Computer Vision WS 15/16

Convolution Layers

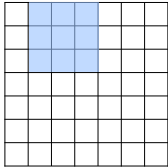
Example:
7x7 input
assume 3x3 connectivity
stride 1

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

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe. 30

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

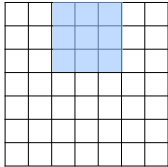
Computer Vision WS 15/16

31

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

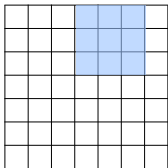
Computer Vision WS 15/16

32

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

RWTH AACHEN
UNIVERSITY

Convolution Layers



Example:
7x7 input
assume 3x3 connectivity
stride 1

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

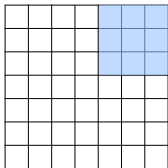
Computer Vision WS 15/16

33

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

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**.

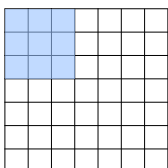
Computer Vision WS 15/16

34

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

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**.

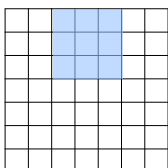
Computer Vision WS 15/16

35

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

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**.

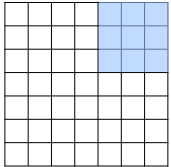
Computer Vision WS 15/16

36

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?
 $\Rightarrow 3 \times 3$ output

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

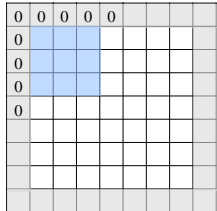
37

Computer Vision WS 15/16

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?
 $\Rightarrow 3 \times 3$ 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.

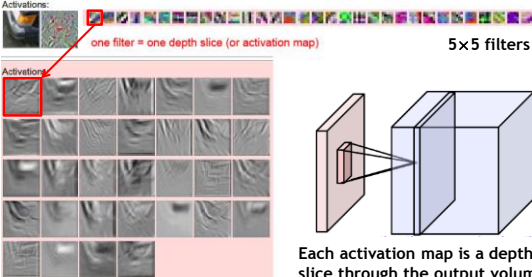
38

Computer Vision WS 15/16

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

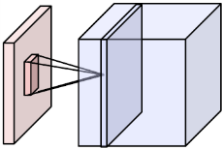
RWTH AACHEN UNIVERSITY

Activation Maps of Convolutional Filters



one filter = one depth slice (or activation map)

5x5 filters



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

Activation maps

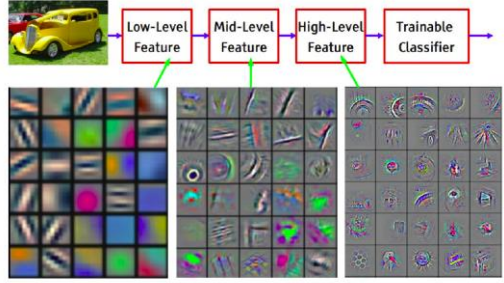
39

Computer Vision WS 15/16

Slide adapted from FeiFei Li, Andrei Karpathy B. Leibe

RWTH AACHEN UNIVERSITY

Effect of Multiple Convolution Layers



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

40

Computer Vision WS 15/16

Slide credit: Yann LeCun B. Leibe

RWTH AACHEN UNIVERSITY

Commonly Used Nonlinearities

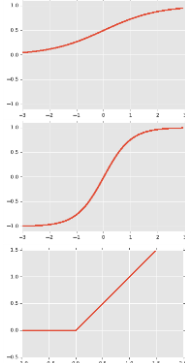
- Sigmoid**

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

$$g(a) = \tanh(a) = 2\sigma(2a) - 1$$
- Rectified linear unit (ReLU)**

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

Currently, preferred option



41

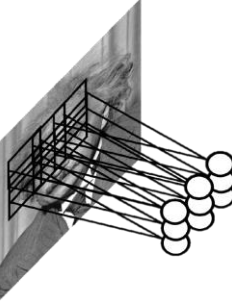
Computer Vision WS 15/16

B. Leibe

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?



42

Computer Vision WS 15/16

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

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?
- Solution:
 - By **pooling** (e.g., max or avg) filter responses at different spatial locations, we gain robustness to the exact spatial location of features.

43

Computer Vision WS 15/16 | Slide adapted from Marc'Aurelio Ranzato | B. Leibe | Image source: Yann Lecun

RWTH AACHEN UNIVERSITY

Max Pooling

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

44

Computer Vision WS 15/16 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe

RWTH AACHEN UNIVERSITY

Max Pooling

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.

45

Computer Vision WS 15/16 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe

RWTH AACHEN UNIVERSITY

Compare: SIFT Descriptor

Image Pixels → Apply oriented filters (Lowe [IJCV 2004])

Spatial pool (Sum)

Normalize to unit length → Feature Vector

46

Computer Vision WS 15/16 | Slide credit: Svetlana Lazebnik | B. Leibe

RWTH AACHEN UNIVERSITY

Compare: Spatial Pyramid Matching

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

47

Computer Vision WS 15/16 | Slide credit: Svetlana Lazebnik | B. Leibe

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Deep Learning
 - Motivation
- Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Applications

46

Computer Vision WS 15/16 | B. Leibe

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

49

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]

50

RWTH AACHEN UNIVERSITY

CNN Architectures: AlexNet (2012)

- Similar framework as LeNet, but
 - Bigger model (7 hidden layers, 650k units, 60M parameters)
 - More data (10⁶ images instead of 10³)
 - 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

51

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

52

RWTH AACHEN UNIVERSITY

AlexNet Results

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

53

RWTH AACHEN UNIVERSITY

AlexNet Results

Test image Retrieved images

54

Computer Vision WS 15/16

Topics of This Lecture

- Deep Learning
 - Motivation
- Convolutional Neural Networks
 - Convolutional Layers
 - Pooling Layers
 - Nonlinearities
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Applications

B. Leibe 61

Computer Vision WS 15/16

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

B. Leibe 62

Computer Vision WS 15/16

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

B. Leibe 63

Computer Vision WS 15/16

Other Tasks: Semantic Segmentation

[Farabet et al. ICML 2012, PAMI 2013]

B. Leibe 64

Computer Vision WS 15/16

Other Tasks: Semantic Segmentation

[Farabet et al. ICML 2012, PAMI 2013]

B. Leibe 65

Computer Vision WS 15/16

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

B. Leibe 66

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

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

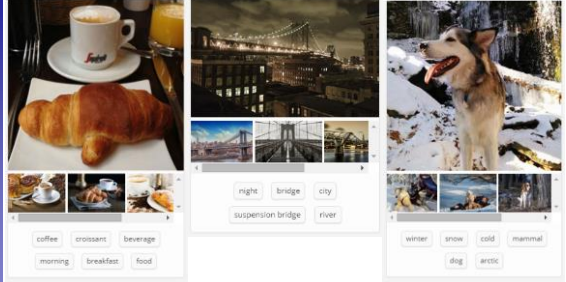
*By using the demo you agree to our terms of service

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

67 B. Leibe Image source: clarifai.com

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

Commercial Recognition Services



68 B. Leibe Image source: clarifai.com

Computer Vision WS 15/16 RWTH AACHEN UNIVERSITY

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.

69 B. Leibe