# Machine Learning – Lecture 14

## Convolutional Neural Networks II

07.01.2019

Bastian Leibe
RWTH Aachen
http://www.vision.rwth-aachen.de

leibe@vision.rwth-aachen.de

Machine Learning Winter '18

---

## 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

2

---

## 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

3

---

## 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

4

---

## 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

Image source: Yann LeCun

5

---

## Recap: 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 [1×1×depth] depth column in output volume.

Slide credit: FeiFei Li, Andrei Karpathy

B. Leibe

6

## Recap: Activation Maps



5×5 filters

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

Activation maps

B. Leibe

7

## Recap: Pooling Layers



Single depth slice

max pool with 2x2 filters and stride 2

- Effect:
  - Make the representation smaller without losing too much information
  - Achieve robustness to translations
  - Pooling happens independently across each slice, preserving the number of slices

B. Leibe

8

## Topics of This Lecture

- Recap: CNNs

- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNet

- Visualizing CNNs
  - Visualizing CNN features
  - Visualizing responses
  - Visualizing learned structures

- Applications

B. Leibe

9

## 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.

B. Leibe

10

## 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

11

## 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.

12

## 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

B. Leibe

13

---

## CNN Architectures: VGGNet (2014/15)



K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015

B. Leibe

15

Image source: Hirokatsu Kataoka

---

## 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%.



B. Leibe

16

Image source: Simonyan & Zisserman

---

## 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

17

---

## 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

18

---

## GoogLeNet Visualization



Inception module  + copies

**Convolution**
**Pooling**
**Softmax**
**Other**

Auxiliary classification outputs for training the lower layers (deprecated)

B. Leibe

19

3

## Slide 20

### 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/

Image source: Simonyan & Zisserman

## Slide 21

### Newer Developments: Residual Networks



AlexNet, 8 layers (ILSVRC 2012)    VGG, 19 layers (ILSVRC 2014)    GoogLeNet, 22 layers (ILSVRC 2014)

## Slide 22

### 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…

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

## Slide 23

### ImageNet Performance



152 layers

| 3.57 | 6.7 | 7.3 | 11.7 | 16.4 | 25.8 | 28.2 |
| ILSVRC'15 ResNet | ILSVRC'14 GoogleNet | ILSVRC'14 VGG | ILSVRC'13 | ILSVRC'12 AlexNet | ILSVRC'11 | ILSVRC'10 |

ImageNet Classification top-5 error (%)

## Slide 24

### ILSRVC Winners

## Slide 25

### Comparing Complexity



A. Canziano, A. Paszke, E. Culurcello, An Analysis of Deep Neural Network Models for Practical Applications, arXiv 2017.

Figure credit: Alfredo Canziano, Adam Paszke, Eugenio Culurcello

## Understanding the ILSVRC Challenge

**RWTH**AACHEN UNIVERSITY

- Imagine the scope of the problem!
  - 1000 categories
  - 1.2M training images
  - 50k validation images

**IM☆GENET**

- 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

26

---

**RWTH**AACHEN UNIVERSITY

American alligator, American black bear, American chameleon, American coot, American egret, American lobster, American Staffordshire terrier, amphibian, analog clock, anemone fish, Angora, ant, apiary, Appenzeller, apron, Arabian camel, Arctic fox, armadillo, artichoke, ashcan, assault rifle, Australian terrier, axolotl, baboon, backpack, badger, bagel, bakery, balance beam, bald eagle, balloon, ballplayer, ballpoint, banana, Band Aid, banded gecko, banjo, bannister, barbell, barber chair, barbershop, barn, barn spider, barometer, barracouta, barrel, barrow, baseball, basenji, basketball, basset, bassinet, bassoon, bath towel, bathing cap, bathtub, beach wagon, beacon, beagle, beaker, bearskin, beaver, Bedlington terrier, bee, bee eater, beer bottle, beer glass, bell cote, bell pepper, Bernese mountain dog, bib, bicycle-built-for-two, bighorn, bikini, binder, binoculars, birdhouse, bison, bittern, black and gold garden spider, black grouse, black stork, black swan, black widow, black-and-tan coonhound, black-footed ferret, Blenheim spaniel, bloodhound, bluetick, boa constrictor, boathouse, bobsled, bolete, bolo tie, bonnet, book jacket, bookcase, bookshop, Border collie, Border terrier, borzoi, Boston bull, bottlecap, Bouvier des Flandres, bow, bow tie, box turtle, boxer, Brabancon griffon, brain coral, brambling, brass, brassiere, breakwater, breastplate, briard, Brittany spaniel, broccoli, broom, brown bear, bubble, bucket, buckeye, buckle, bulbul, bull mastiff, bullet train, bulletproof vest, bullfrog, burrito, bustard, butcher shop, butternut squash, cab, cabbage butterfly, cairn, caldron, can opener, candle, cannon, canoe, capuchin, car mirror, carwheel, carbonara, Cardigan, cardigan, cardoon, carousel, carpenter's kit, carton, cash machine, cassette, cassette player, castle, catamaran, cauliflower, CD player, cello, cellular telephone, centipede, chain, chain mail, chain saw, chainlink fence, chambered nautilus, cheeseburger, cheetah, Chesapeake Bay retriever, chest, chickadee, chiffonier, Chihuahua, chime, chimpanzee, china cabinet, chiton, chocolate sauce, chow, Christmas stocking, church, cicada, cinema, cleaver, cliff, cliff dwelling, cloak, clog, clumber, cock, cocker spaniel, cockroach, cocktail shaker, coffee mug, coffeepot, coho, coil, collie, colobus, combination lock, comic book, common iguana, common newt, computer keyboard, conch, confectionery, consomme, container ship, convertible, coral fungus, coral reef, corkscrew, corn, cornet, coucal, cougar, cowboy boot, cowboy hat, coyote, cradle, crane, crane, crash helmet, crate, crayfish, crib, cricket, Crock Pot, croquet ball, crossword puzzle, crutch, cucumber, cuirass, cup, curly-coated retriever, custard apple, daisy, dalmatian, dam, damselfly, Dandie Dinmont, desk, desktop computer, dhole, dial telephone, diamondback, diaper, digital clock, digital watch, dingo, dining table, dishrag, dishwasher, disk brake, Doberman, dock, dogsled, dome, doormat, dough, dowitcher, dragonfly, drake, drilling platform, drum, drumstick, dugong, dumbbell, dung beetle, Dungeness crab, Dutch oven, ear, earthstar, echidna, eel, eft, eggnog, Egyptian cat, electric fan, electric guitar, electric locomotive, electric ray, English foxhound, English setter, English springer, entertainment center, EntleBucher, envelope, Eskimo dog, espresso, espresso maker, European fire salamander, European gallinule, face powder, feather boa, fiddler crab, fig, file, fire engine, fire screen, fireboat, flagpole, flamingo, flat-coated retriever, flatworm, flute, fly, folding chair, football helmet, forklift, fountain, fountain pen, four-poster, fox squirrel, freight car, French bulldog, French horn, French loaf, frilled lizard, frying pan, fur coat, gar, garbage truck, garden spider, garter snake, gas pump, gasmask, gazelle, German shepherd, German short-haired pointer, geyser, giant panda, giant schnauzer, gibbon, Gila monster, go-kart, goblet, golden retriever, goldfinch, goldfish, golf ball, golfcart, gondola, gong, goose, Gordon setter, gorilla, gown, grand piano, Granny Smith, grasshopper, Great Dane, great grey owl, Great Pyrenees, great white shark

27

---

## More Finegrained Classes

**RWTH**AACHEN UNIVERSITY



B. Leibe

Image source: O. Russakovsky et al.

28

---

## Quirks and Limitations of the Data Set

**RWTH**AACHEN UNIVERSITY



- 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. Karpathy

29

---

## Topics of This Lecture

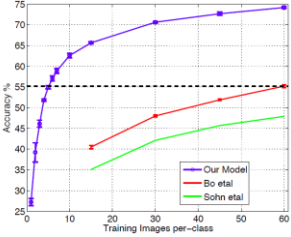**RWTH**AACHEN UNIVERSITY

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

B. Leibe

30

---

## Visualizing CNNs

**RWTH**AACHEN UNIVERSITY



DeconvNet / ConvNet

Image source: M. Zeiler, R. Fergus

31

---

5

Visualizing CNNs

Layer 1

Layer 2

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.

Slide credit: Richard Turner

B. Leibe

32

Image source: M. Zeiler, R. Fergus

---

Visualizing CNNs

Layer 3

B. Leibe

33

Image source: M. Zeiler, R. Fergus

---

Visualizing CNNs

Layer 4

Layer 5

B. Leibe

34

Image source: M. Zeiler, R. Fergus

---

What Does the Network React To?

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

B. Leibe

35

---

What Does the Network React To?

Input image

True Label: Pomeranian

p(True class)

Most probable class

Pomeranian
Tennis ball
Keeshond
Pekinese

Slide credit: Svetlana Lazebnik, Rob Fergus

36

Image source: M. Zeiler, R. Fergus

---

What Does the Network React To?

Input image

True Label: Pomeranian

Total activation in most active 5th layer feature map

Other activations from the same feature map.

Slide credit: Svetlana Lazebnik, Rob Fergus

37

Image source: M. Zeiler, R. Fergus

6

## Inceptionism: Dreaming ConvNets



optimize with prior

- 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

## Inceptionism: Dreaming ConvNets

- Results



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

## Inceptionism: Dreaming ConvNets



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

44

---

## 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

45

---

## 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

Image source: M. Zeiler, R. Fergus

46

---

## 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

I.e., swap the Softmax layer at the end

Slide credit: Andrej Karpathy

B. Leibe

47

---

## Transfer Learning with CNNs



1. Train on ImageNet

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.

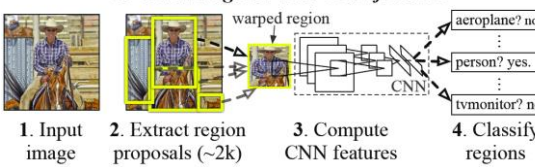Retrain bigger portion of the network

Slide credit: Andrej Karpathy

B. Leibe

48

---

## Other Tasks: Detection



**R-CNN: Regions with CNN features**

warped region

aeroplane? no.

person? yes.

tvmonitor? no.

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

49

---

8

## 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

50

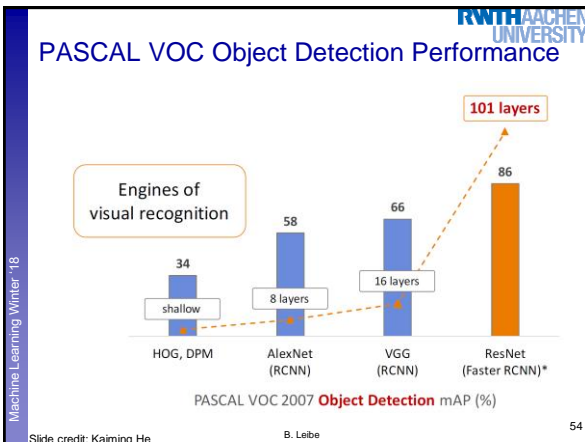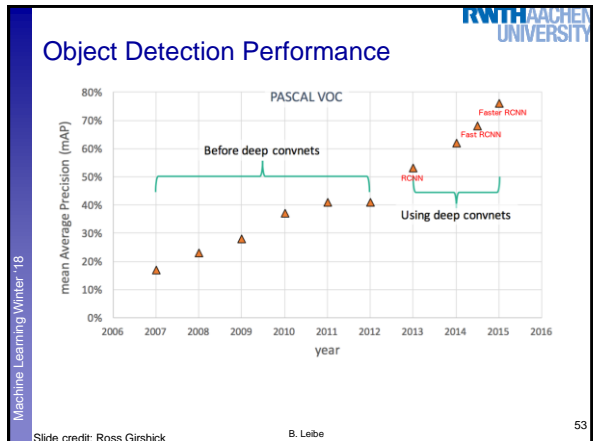

## Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

B. Leibe

51



## Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

B. Leibe

52



## Object Detection Performance

Slide credit: Ross Girshick

B. Leibe

53



## PASCAL VOC Object Detection Performance

Slide credit: Kaiming He

B. Leibe

54



## Most Recent Version: Mask R-CNN

K. He, G. Gkioxari, P. Dollar, R. Girshick, Mask R-CNN, arXiv 1703.06870.

Slide credit: FeiFei Li

55

9

## Mask R-CNN Results

- Detection + Instance segmentation



- Detection + Pose estimation



Figure credit: K. He, G. Gkioxari, P. Dollar, R. Girshick

## YOLO / SSD



Input image
3 x H x W

Divide image into grid
7 x 7

- Idea: Directly go from image to detection scores
- Within each grid cell
  - Start from a set of anchor boxes
  - Regress from each of the B anchor boxes to a final box
  - Predict scores for each of C classes (including background)

Slide credit: FeiFei Li

57

## YOLO-v3 Results

J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016.

60

## Semantic Image Segmentation



- Perform pixel-wise prediction task
  - Usually done using Fully Convolutional Networks (FCNs)
    – All operations formulated as convolutions
    – Advantage: can process arbitrarily sized images

61

Image source: Long, Shelhamer, Darrell

## CNNs vs. FCNs

- CNN



"tabby cat"

- FCN

convolutionalization

tabby cat heatmap

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

62

Image source: Long, Shelhamer, Darrell

## Semantic Image Segmentation



- Encoder-Decoder Architecture
  - Problem: FCN output has low resolution
  - Solution: perform upsampling to get back to desired resolution
  - Use skip connections to preserve higher-resolution information
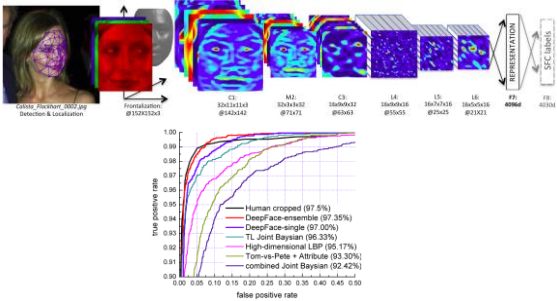
63

Image source: Newell et al.

## Semantic Segmentation



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

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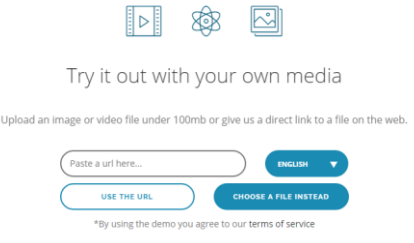
---

## 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

65

---

## 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...

B. Leibe

66

Image source: clarifai.com

---

## 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

67

---

## 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.

B. Leibe

68

---

## References: Computer Vision Tasks

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

B. Leibe

69

# 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.

Machine Learning Winter '18

B. Leibe

70

12