# Advanced Machine Learning Lecture 16

## Convolutional Neural Networks II

### 22.12.2016

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

leibe@vision.rwth-aachen.de

*Advanced Machine Learning Winter'16*

---

## This Lecture: *Advanced Machine Learning*

- **Regression Approaches**
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
  - Gaussian Processes

$$f : \mathcal{X} \to \mathbb{R}$$

- **Approximate Inference**
  - Sampling Approaches
  - MCMC

- **Deep Learning**
  - Linear Discriminants
  - Neural Networks
  - Backpropagation & Optimization
  - **CNNs**, RNNs, ResNets, etc.

B. Leibe

*Advanced Machine Learning Winter'16*

---

## 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

4

*Advanced Machine Learning Winter'16*

---

## 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

5

*Advanced Machine Learning Winter'16*

---

## 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

6

Image source: Yann LeCun

*Advanced Machine Learning Winter'16*

---

## Recap: Convolution Layers



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
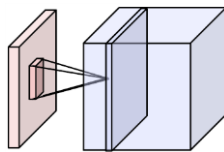  - Form a single [1×1×depth] depth column in output volume.

Slide credit: FeiFei Li, Andrej Karpathy

B. Leibe

7

*Advanced Machine Learning Winter'16*

## Recap: Activation Maps

Activations:

one filter = one depth slice (or activation map)

**5×5 filters**

Activations:

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

**Activation maps**

Slide adapted from FeiFei Li, Andrej Karpathy    B. Leibe    8

---

## Recap: Pooling Layers

Single depth slice

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

x

max pool with 2x2 filters and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

y

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

Slide adapted from FeiFei Li, Andrej Karpathy    B. Leibe    9

---

## Topics of This Lecture

- Recap: CNNs

- **CNN Architectures**
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet

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

- Applications

B. Leibe    10

---

## Recap: ImageNet Challenge 2012

- **ImageNet**
  - ~14M labeled internet images
  - 20k classes
  - Human labels via Amazon Mechanical Turk

**IMAGENET**
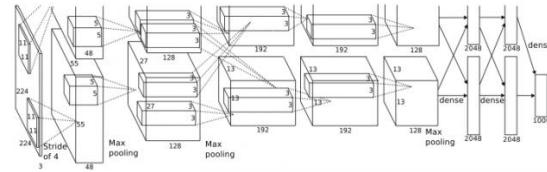
- **Challenge (ILSVRC)**
  - 1.2 million training images
  - 1000 classes
  - Goal: Predict ground-truth class within top-5 responses
  - Currently one of the top benchmarks in Computer Vision

**[Deng et al., CVPR'09]**

B. Leibe    12
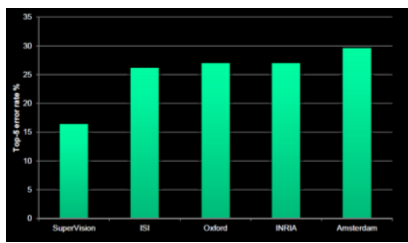
---

## 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.    13

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

---

## 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    14

## CNN Architectures: VGGNet (2014/15)



**AlexNet**

Input — Conv — Conv — Pool — Conv — Pool — Conv — Conv — FC — FC — Softmax

Layer1 Layer2 Layer3 Layer4 Layer5 Layer6 Layer7

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

**VGGNet**

Input — Conv — Conv — Pool — Conv — Conv — Pool — Conv — Conv — Pool — Conv — Conv — Pool — Conv — Conv — Pool — FC — FC — FC — Softmax

Layer1 Layer2 Layer3 Layer4 Layer5 Layer6 Layer7

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

B. Leibe

Image source: Hirokatsu Kataoka

16

---

## 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 |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 LRN | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 conv1-256 | conv3-256 conv3-256 conv3-256 | conv3-256 conv3-256 conv3-256 conv3-256 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | 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

Image source: Simonyan & Zisserman

17

---

## Comparison: AlexNet vs. VGGNet

- **Receptive fields in the first layer**
  - AlexNet: 11×11, stride 4
  - Zeiler & Fergus: 7×7, stride 2
  - VGGNet: 3×3, stride 1

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

B. Leibe

18

---

## CNN Architectures: GoogLeNet (2014)



(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.

B. Leibe

19

---

## GoogLeNet Visualization



Inception module

+ copies

Convolution
Pooling
Softmax
Other

Auxiliary classification outputs for training the lower layers (deprecated)

B. Leibe

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/

B. Leibe

Image source: Simonyan & Zisserman

21

## Newest Development: Residual Networks

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

GoogleNet, 22 layers
(ILSVRC 2014)



Advanced Machine Learning Winter'16

B. Leibe

22

---

## Newest Development: 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…**

$$x$$

weight layer

$$F(x)$$ relu

weight layer

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

Advanced Machine Learning Winter'16

B. Leibe

23

---

## ImageNet Performance

**152 layers**

3.57

6.7 — 22 layers

7.3 — 19 layers

11.7 — 8 layers

16.4 — 8 layers

25.8 — shallow

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 (%)

Advanced Machine Learning Winter'16

B. Leibe

24

---

## Understanding the ILSVRC Challenge

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

IMAGENET

- **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.**

Advanced Machine Learning Winter'16

B. Leibe

25

---

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, 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, car wheel, 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, flatcoated 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,

Advanced Machine Learning Winter'16

26

---

## More Finegrained Classes

PASCAL                  ILSVRC

birds

bird | flamingo | cock | ruffed grouse | quail | partridge | ...

cats

cat | Egyptian cat | Persian cat | Siamese cat | tabby | lynx | ...

dogs

dog | dalmatian | keeshond | miniature schnauzer | standard schnauzer | giant schnauzer | ...

Advanced Machine Learning Winter'16

B. Leibe

27

Image source: O. Russakovsky et al.

---

4

## 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
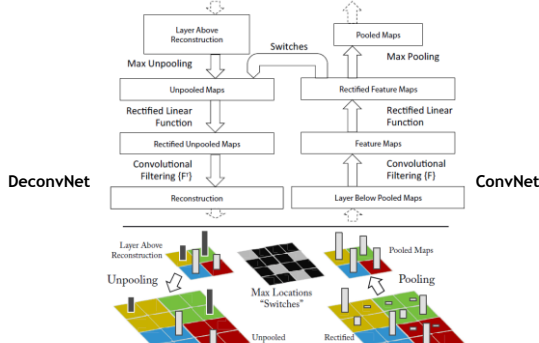Image source: A. Karpathy

---

## Topics of This Lecture

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

B. Leibe

---

## Visualizing CNNs



DeconvNet · ConvNet

Image source: M. Zeiler, R. Fergus

---

## 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.

Slide credit: Richard Turner
B. Leibe
Image source: M. Zeiler, R. Fergus

---

## Visualizing CNNs



Layer 3

B. Leibe
Image source: M. Zeiler, R. Fergus

---

## Visualizing CNNs



Layer 4 · Layer 5

B. Leibe
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

---

## What Does the Network React To?

**Input image**

True Label: Pomeranian

**p(True class)**     **Most probable class**

- Pomeranian
- Tennis ball
- Keeshond
- Pekinese

---

## 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.**

---

## What Does the Network React To?

**Input image**

True Label: Car Wheel

**p(True class)**     **Most probable class**

- Car wheel
- Racer
- Cab
- Police van

---

## What Does the Network React To?

**Input image**

True Label: Car Wheel

**Total activation in most active 5th layer feature map**     **Other activations from the same feature map.**

---

## What Does the Network React To?

**Input image**

True Label: Afghan Hound

**p(True class)**     **Most probable class**

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

## What Does the Network React To?

Input image

True Label: Afghan Hound

Total activa-
tion in most
active 5th
layer feature
map

Other activa-
tions from the
same feature
map.

Advanced Machine Learning Winter'16

40

---

## 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.**

Advanced Machine Learning Winter'16

41

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

---

## Inceptionism: Dreaming ConvNets

- **Results**



Advanced Machine Learning Winter'16

42

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

---

## Inceptionism: Dreaming ConvNets

Advanced Machine Learning Winter'16

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

43

---

## Topics of This Lecture

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

Advanced Machine Learning Winter'16

B. Leibe

44

---

## 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**

Advanced Machine Learning Winter'16

B. Leibe

45

## Other Tasks: Detection

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



warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

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

46

---

## Faster R-CNN (based on ResNets)



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

B. Leibe

47

---

## Faster R-CNN (based on ResNets)



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

B. Leibe

48

---

## Other Tasks: Semantic Segmentation



[Farabet et al. ICML 2012, PAMI 2013]

B. Leibe

49

---

## Semantic Segmentation



[Pohlen, Hermans, Mathias, Leibe, arXiv 2016]

- **More recent results**
  - **Based on an extension of ResNets**

---

## 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

52

Slide credit: Svetlana Lazebnik

## Commercial Recognition Services



**clarifai**

*Advanced Machine Learning Winter'16*

B. Leibe

54

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.

*Advanced Machine Learning Winter'16*

B. Leibe

55

## References and Further Reading
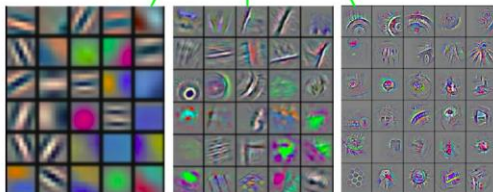
- **ResNet**
  - K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

*Advanced Machine Learning Winter'16*

B. Leibe

56

## Effect of Multiple Convolution Layers



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

Slide credit: Yann LeCun

*Advanced Machine Learning Winter'16*

B. Leibe

60