

# Advanced Machine Learning Lecture 16

## Convolutional Neural Networks II

14.01.2016

Bastian Leibe

RWTH Aachen

<http://www.vision.rwth-aachen.de/>

[leibe@vision.rwth-aachen.de](mailto:leibe@vision.rwth-aachen.de)

# Announcements

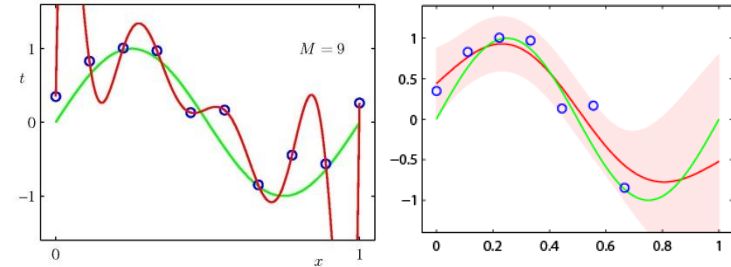
- **Lecture evaluation**
  - Please fill out the evaluation forms.

# This Lecture: *Advanced Machine Learning*

## • Regression Approaches

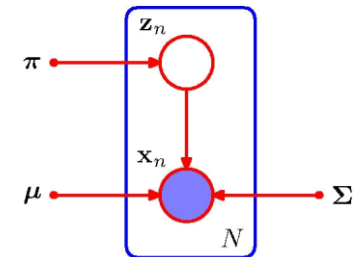
- Linear Regression
- Regularization (Ridge, Lasso)
- Gaussian Processes

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



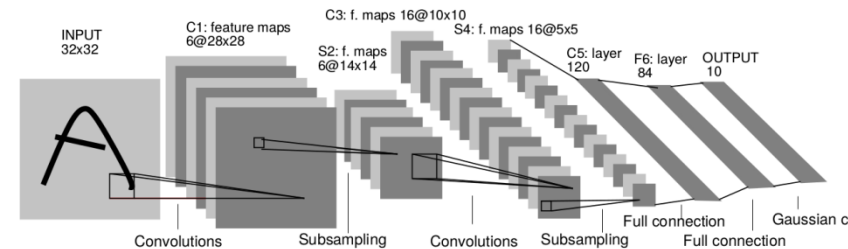
## • Learning with Latent Variables

- Prob. Distributions & Approx. Inference
- Mixture Models
- EM and Generalizations



## • Deep Learning

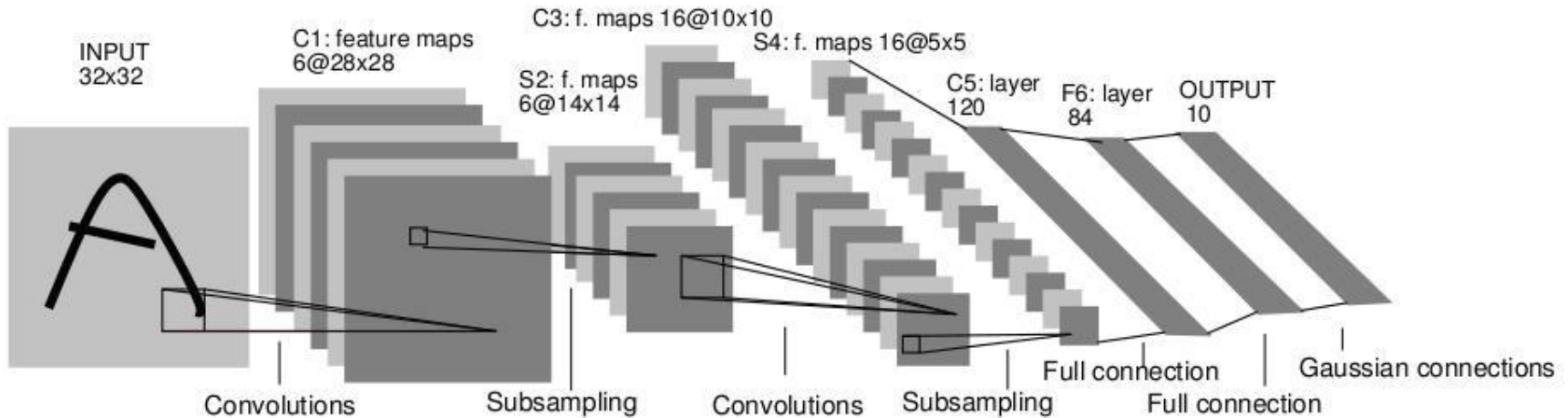
- Linear Discriminants
- Neural Networks
- Backpropagation & Optimization
- CNNs, RNNs, RBMs, etc.



# Topics of This Lecture

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

# 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.

# Recap: Intuition of CNNs

- Convolutional net

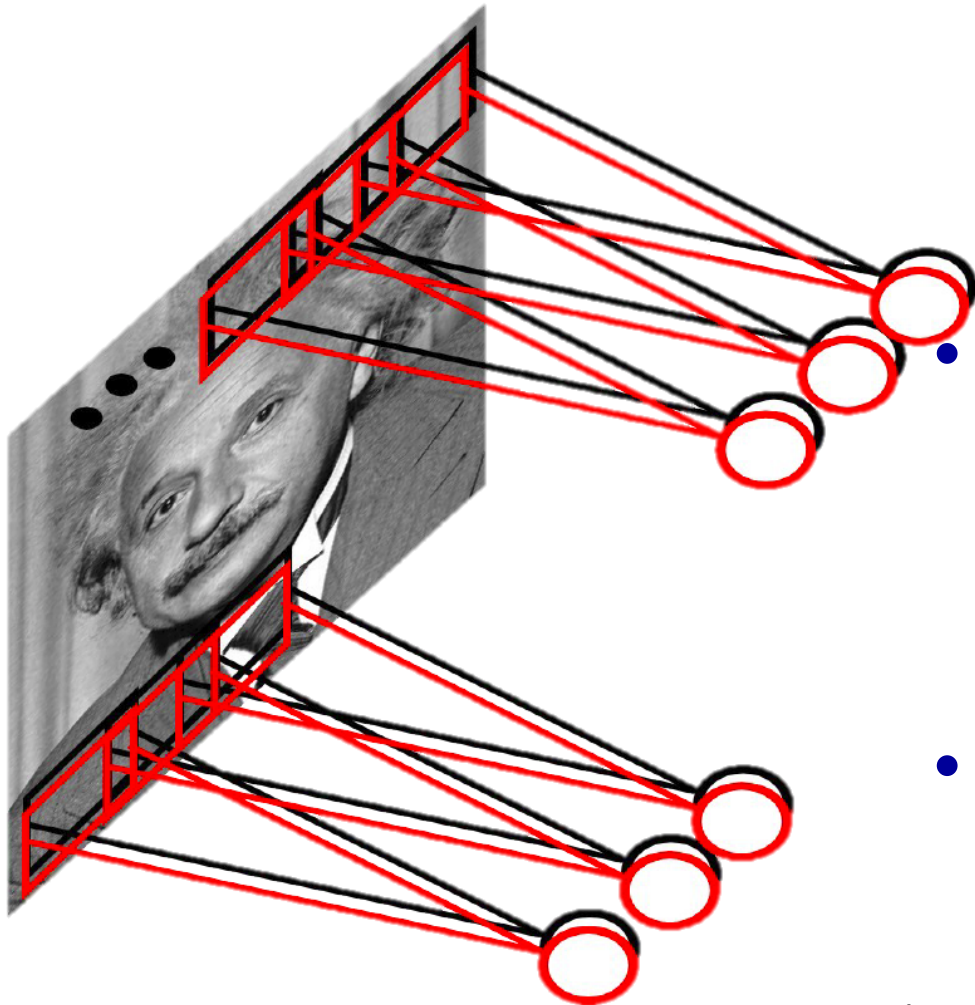
- Share the same parameters across different locations
- Convolutions with learned kernels

- Learn *multiple* filters

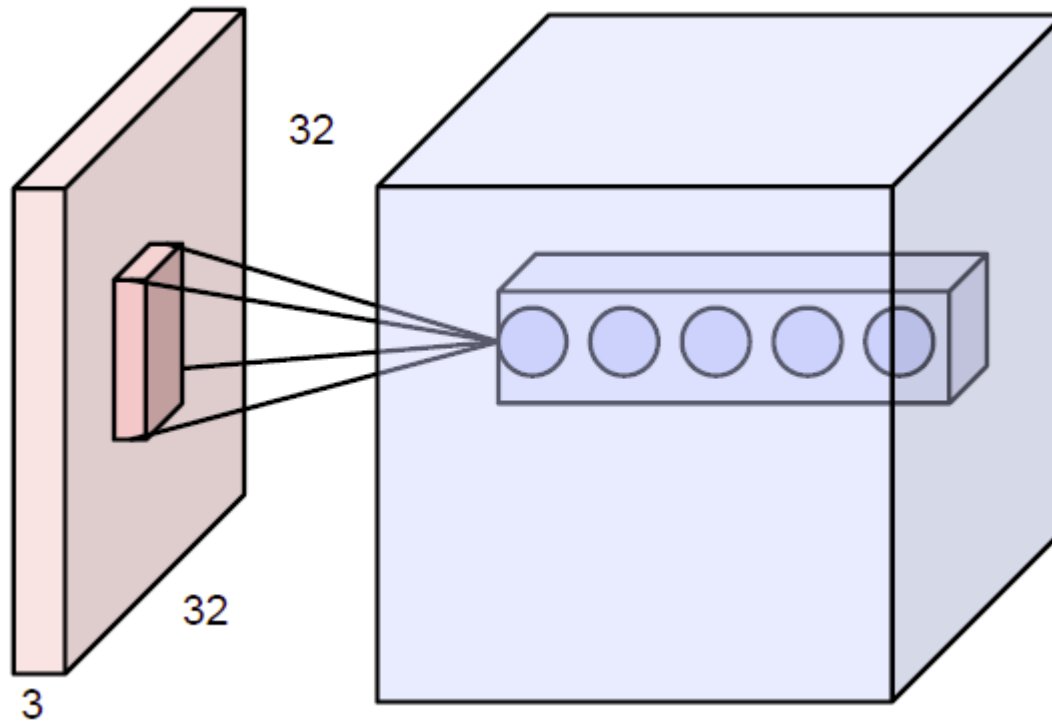
- E.g.  $1000 \times 1000$  image  
100 filters  
 $10 \times 10$  filter size  
⇒ only 10k parameters

- Result: Response map

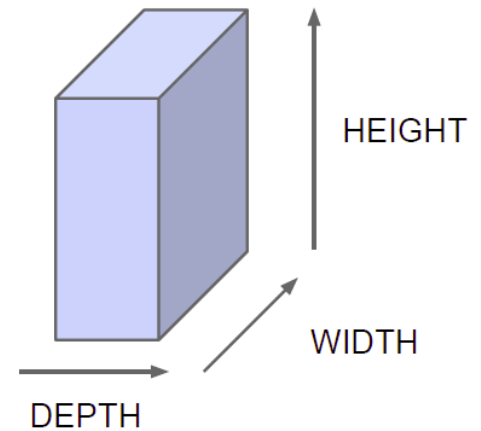
- size:  $1000 \times 1000 \times 100$
- Only memory, not params!



# Recap: Convolution Layers



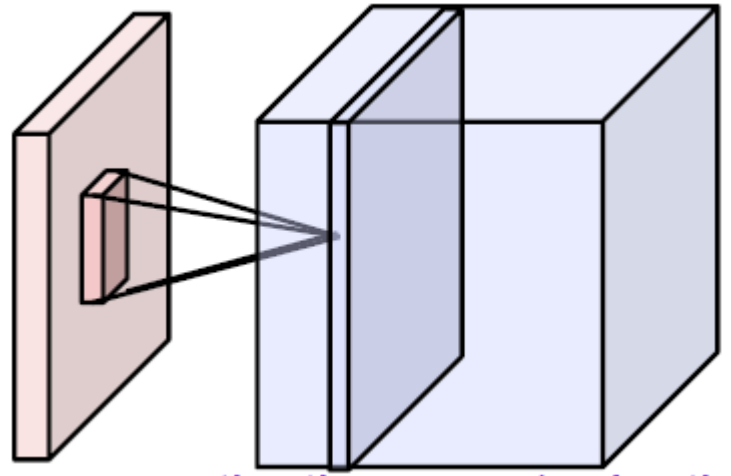
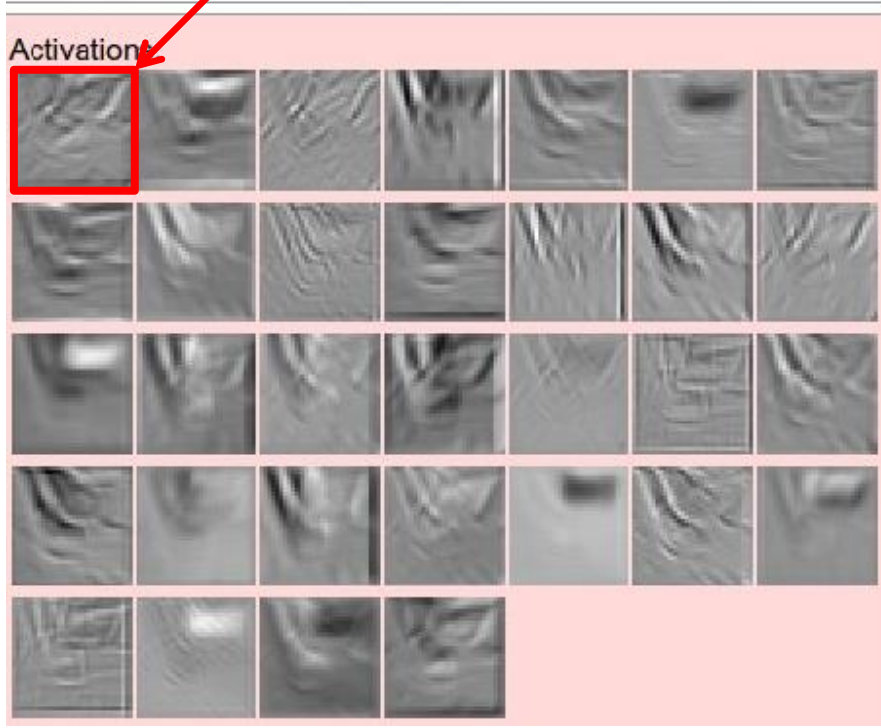
Naming convention:



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

# Recap: Activation Maps

Activations:

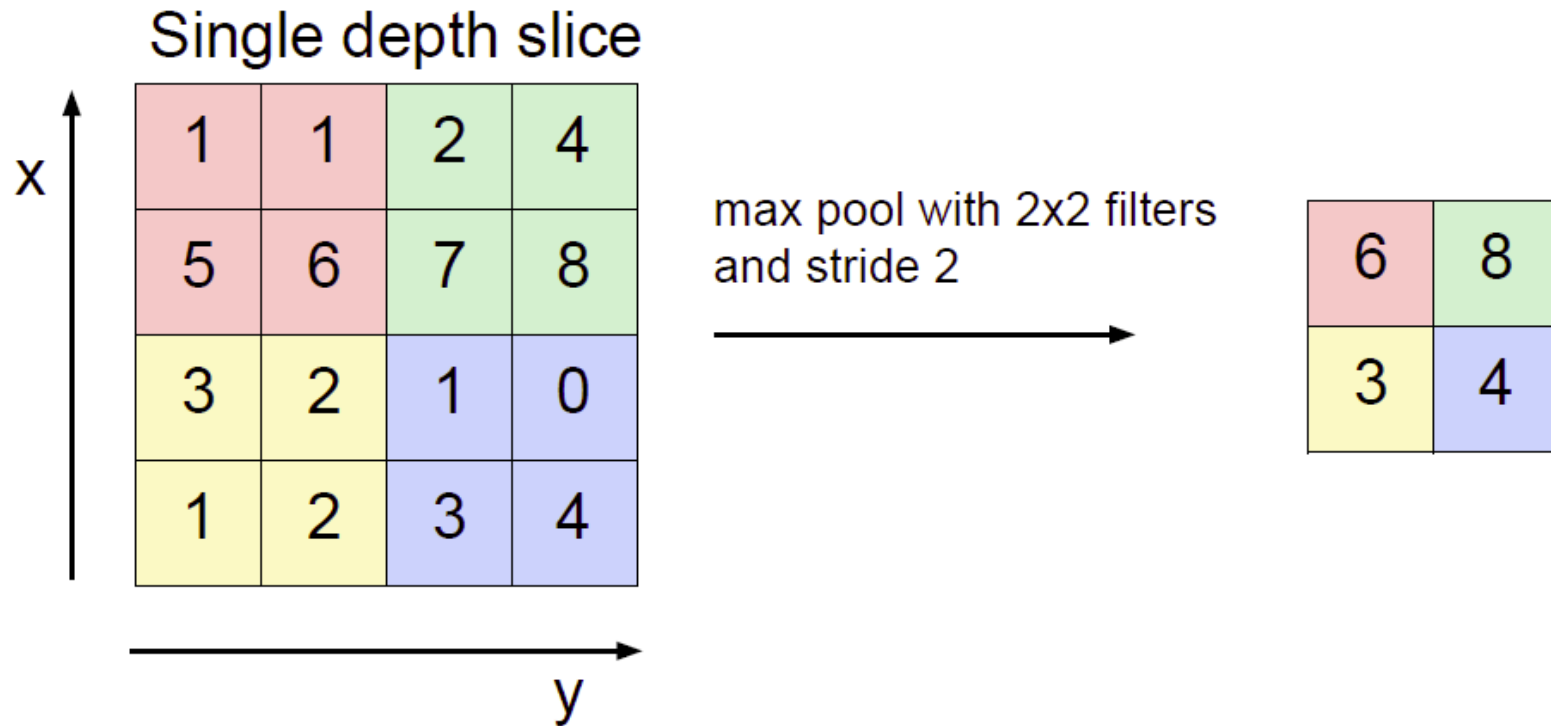


Each activation map is a depth slice through the output volume.

Activation maps



# Recap: Pooling Layers



- **Effect:**

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

# Topics of This Lecture

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

# Recap: ImageNet Challenge 2012

- ImageNet

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

IM  GENET

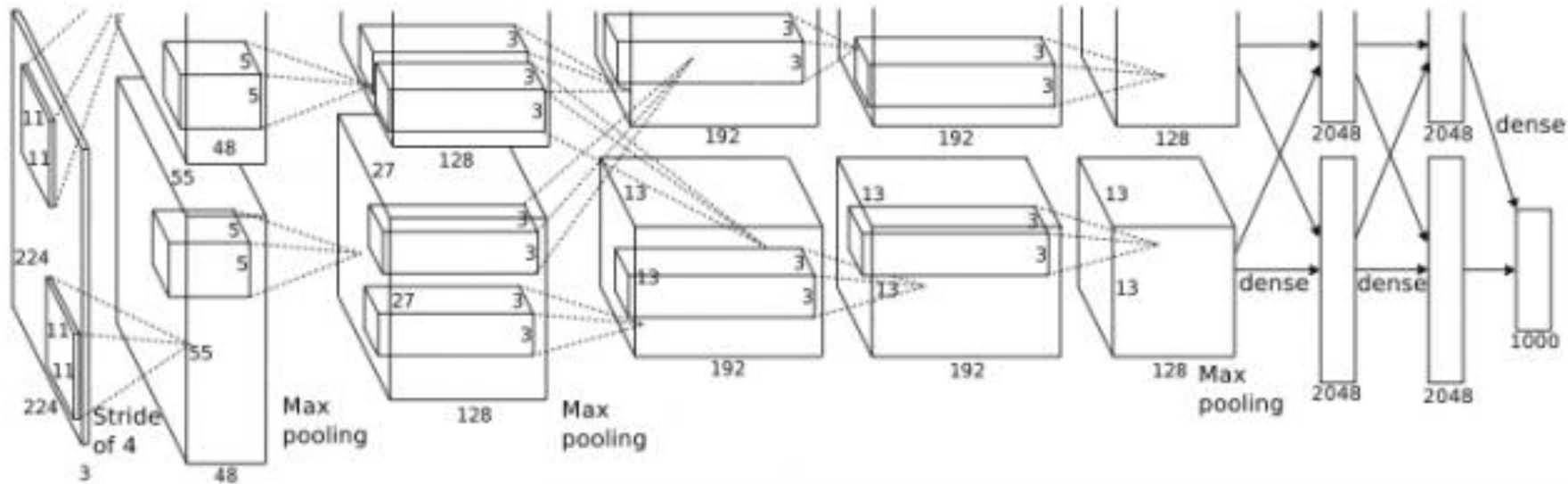


- Challenge (ILSVRC)

- 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]

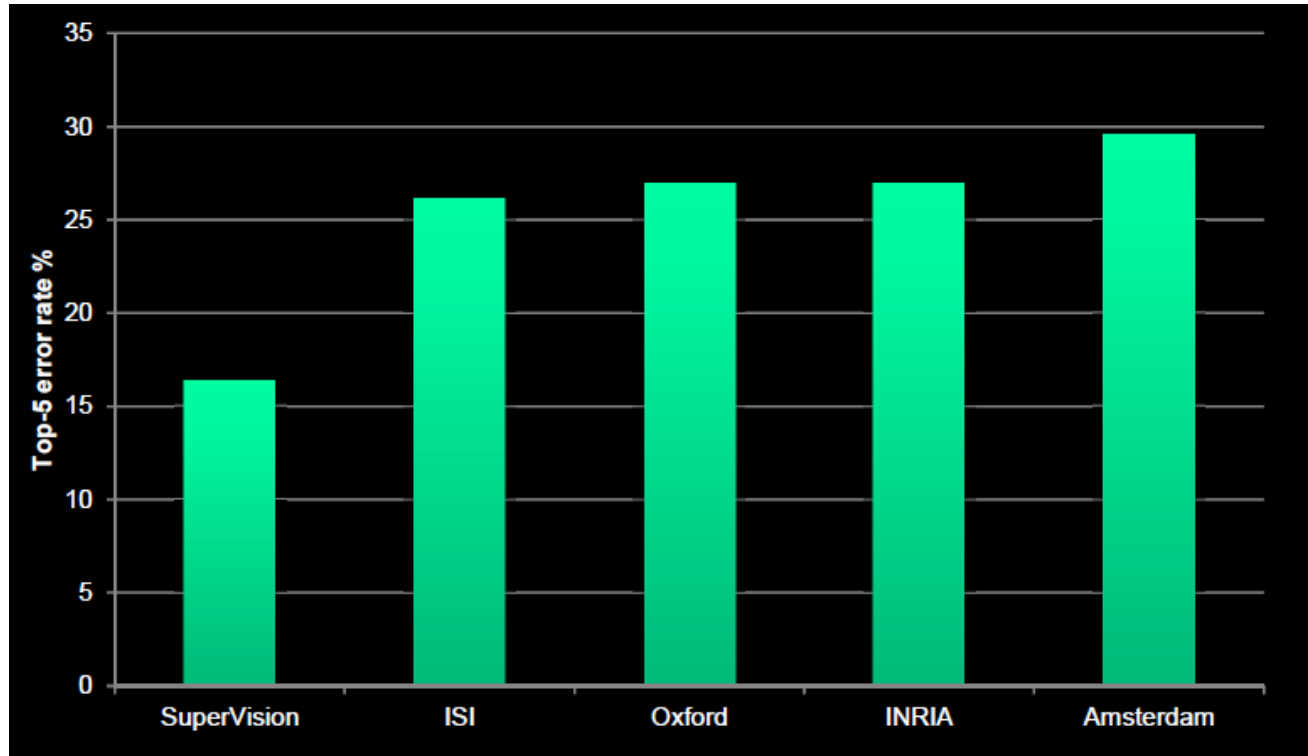
# CNN Architectures: AlexNet (2012)



- **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.

# 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

# CNN Architectures: VGGNet (2015)

- Main ideas

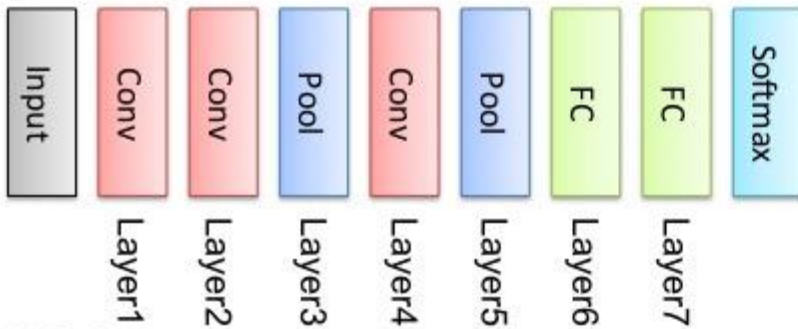
- Deeper network
- Stacked convolutional layers with smaller filters (+ nonlinearity)
- Detailed evaluation of all components

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 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	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 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

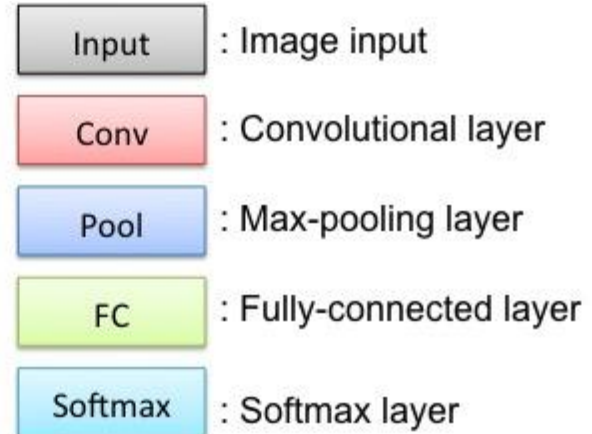
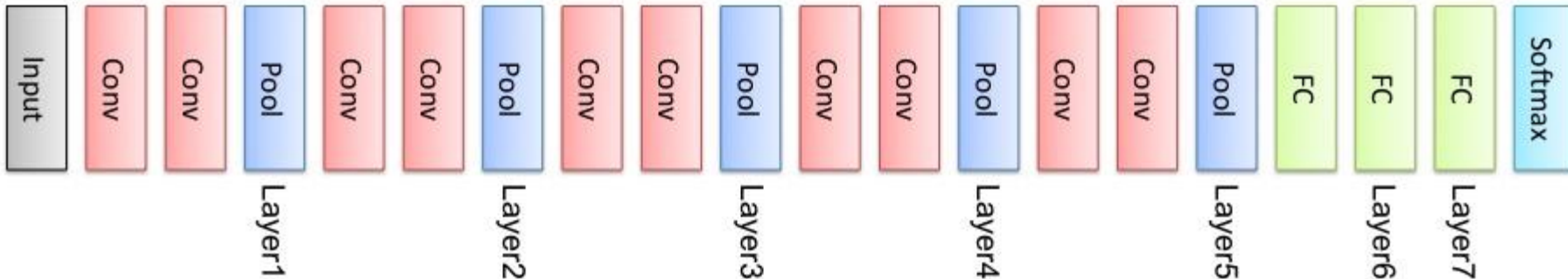
Mainly used

# Comparison: AlexNet vs. VGGNet

## AlexNet



## VGGNet



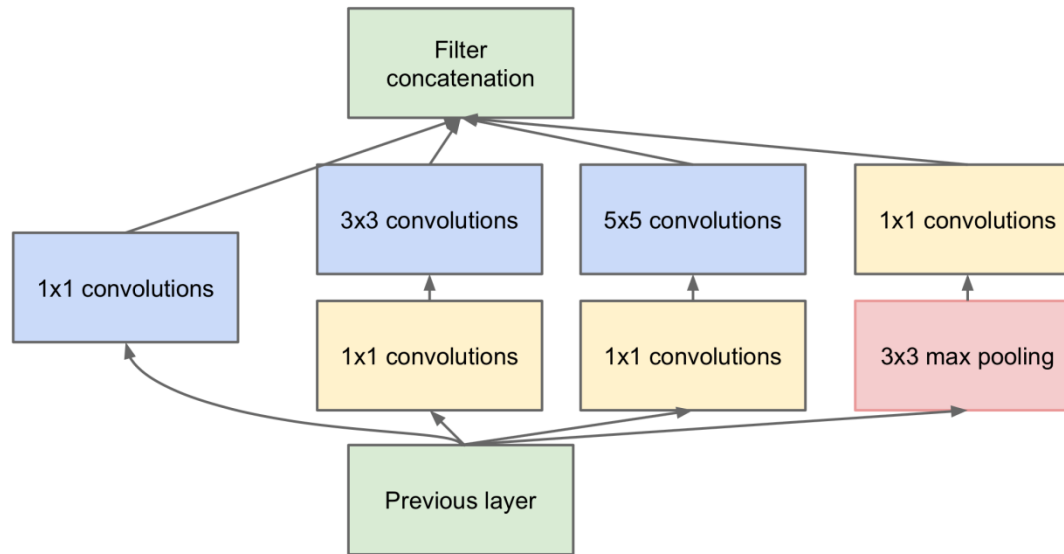
K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

# 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 three  $3 \times 3$  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.



# CNN Architectures: GoogLeNet (2014)



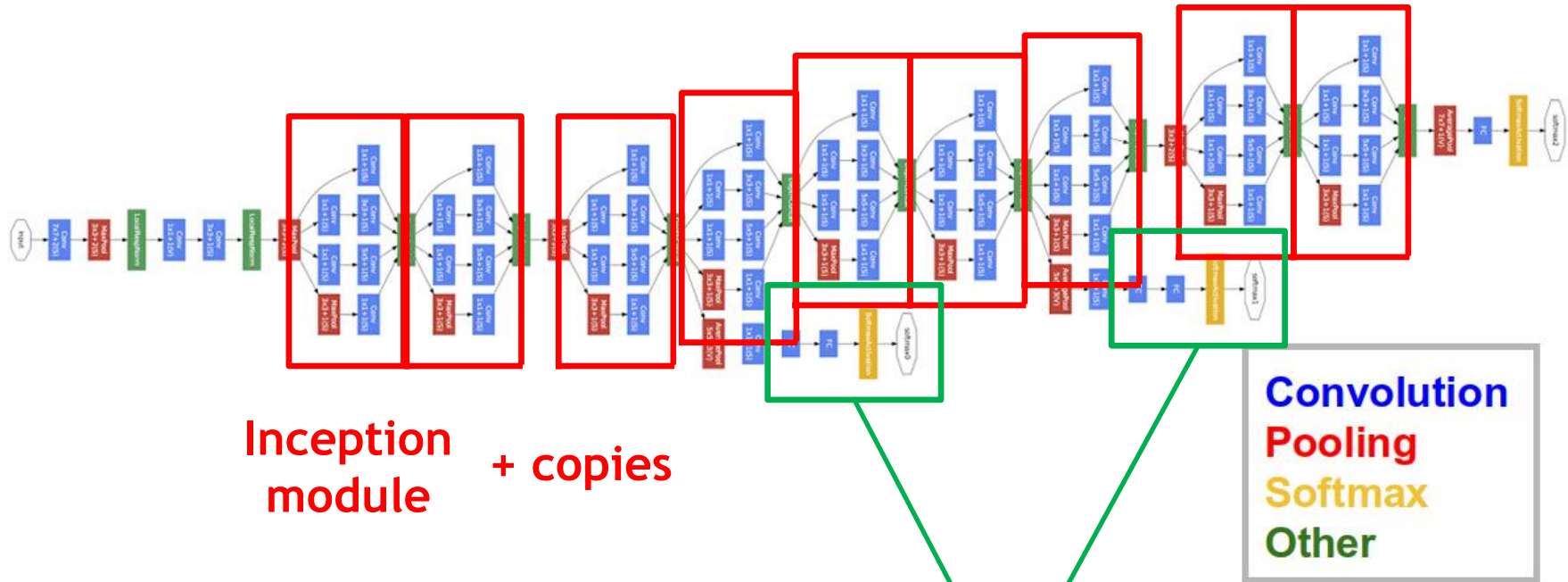
Inception module with dimension reductions

- **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.

# GoogLeNet Visualization



Inception module + copies

Auxiliary classification outputs for training the lower layers (deprecated)

Convolution  
Pooling  
Softmax  
Other

# Results on ILSVRC

Method	top-1 val. error (%)	top-5 val. error (%)	top-5 test error (%)
VGG (2 nets, multi-crop & dense eval.)	<b>23.7</b>	<b>6.8</b>	<b>6.8</b>
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)	-	<b>6.7</b>	
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/>

# Understanding the ILSVRC Challenge

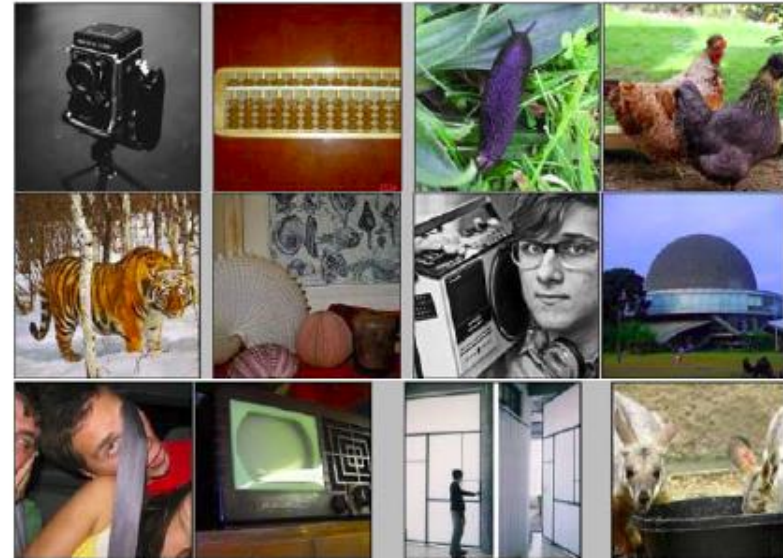
- Imagine the scope of the problem!

- 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.

# IMAGENET



r, Alredale, airliner, airship, albatross, alligator lizard, alp, altar, ambulance, 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, 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, chain-machine, Shetland sheepdog, shield, Shih-Tzu, shoe shop, shoji, shopping basket, shopping cart, shovel, shower cap, shower curtain, siamang, Siamese cat, Siberian husky, sidewinder, silky terrier, ski, ski mask, skunk, sleeping bag, slide rule, sliding door, slot, sloth bear, slug, snail, snorkel, snow leopard, snowmobile, snowplow, soap dispenser, soccer ball, sock, soft-coated wheaten terrier, solar dish, sombrero, sorrel, soup bowl, space bar, space heater, space shuttle, spaghetti squash, spatula, speedboat, spider monkey, spider web, spindle, spiny lobster, spoonbill, sports car, spotlight, spotted salamander, squirrel monkey, Staffordshire bullterrier, stage, standard poodle, standard schnauzer, starfish, steam locomotive, steel arch bridge, steel drum, stethoscope, stingray, stinkhorn, stole, stone wall, stopwatch, stove, strainer, strawberry, street sign, stretcher, studio couch, stupa, sturgeon, submarine, suit, sulphur butterfly, sulphur-crested cockatoo, sundial, sunglass, sunglasses, sunscreen, suspension bridge, Sussex spaniel, swab, sweatshirt, swimming trunks, swing, switch, syringe, tabby, table lamp, tailed frog, tank, tape player, tarantula, teapot, teddy, television, tench, tennis ball, terrapin, thatch, theater curtain, thimble, three-toed sloth, thresher, throne, thunder snake, Tibetan mastiff, Tibetan terrier, tick, tiger, tiger beetle, tiger cat, tiger shark, tile roof, timber wolf, titi, toaster, tobacco shop, toilet seat, toilet tissue, torch, totem pole, toucan, tow truck, toy poodle, toy terrier, toyshop, tractor, traffic light, trailer truck, tray, tree frog, trench coat, triceratops, tricycle, trifle, trilobite, trimaran, tripod, triumphal arch, trolleybus, trombone, tub, turnstile, tusker, typewriter keyboard, umbrella, unicycle, upright, vacuum, valley, vase, vault, velvet, vending machine, vestment, viaduct, vine snake, violin, vizsla, volcano, volleyball, vulture, waffle iron, Walker hound, walking stick, wall clock, wallaby, wallet, wardrobe, warplane, warthog, washbasin, washer, water bottle, water buffalo, water jug, water ouzel, water snake, water tower, weasel, web site, weevil, Weimaraner, Welsh springer spaniel, West Highland white terrier, whippet, whiptail, whiskey jug,

# More Finegrained Classes

PASCAL

birds



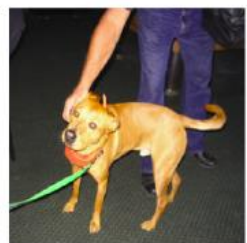
bird

cats



cat

dogs



dog

ILSVRC



flamingo



cock



ruffed grouse



quail



partridge

...



Egyptian cat



Persian cat



Siamese cat

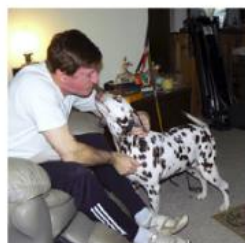


tabby



lynx

...



dalmatian



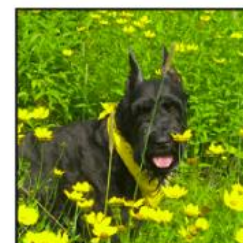
keeshond



miniature schnauzer



standard schnauzer



giant schnauzer

...

# 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

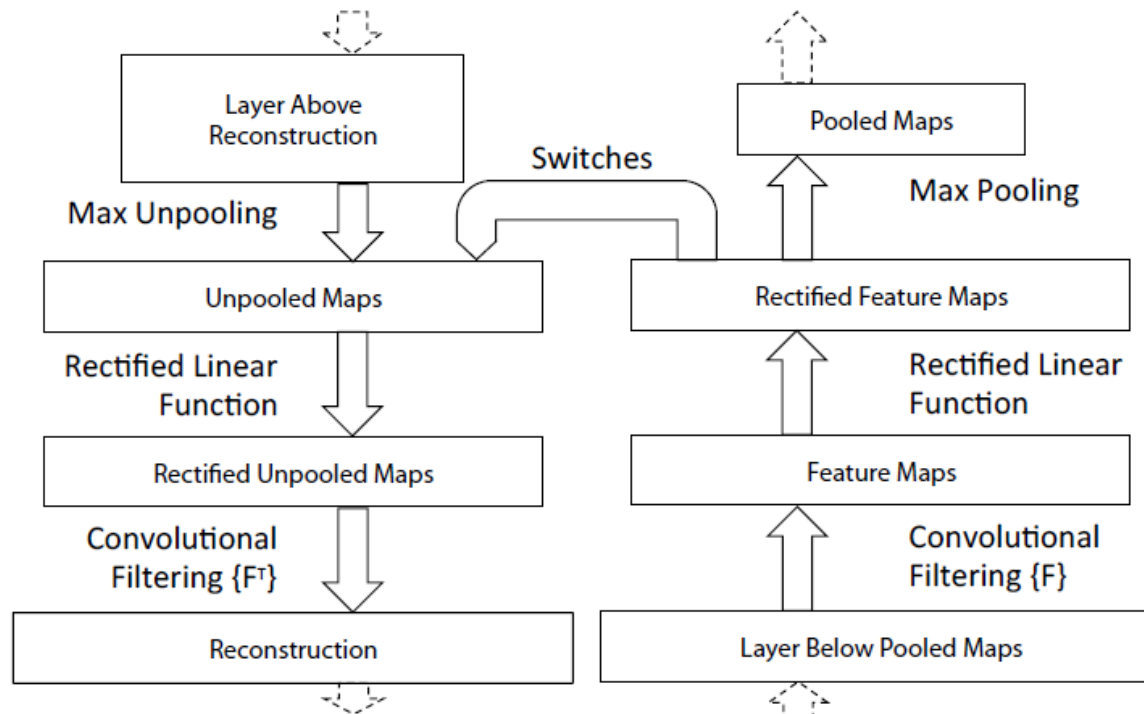
# Topics of This Lecture

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

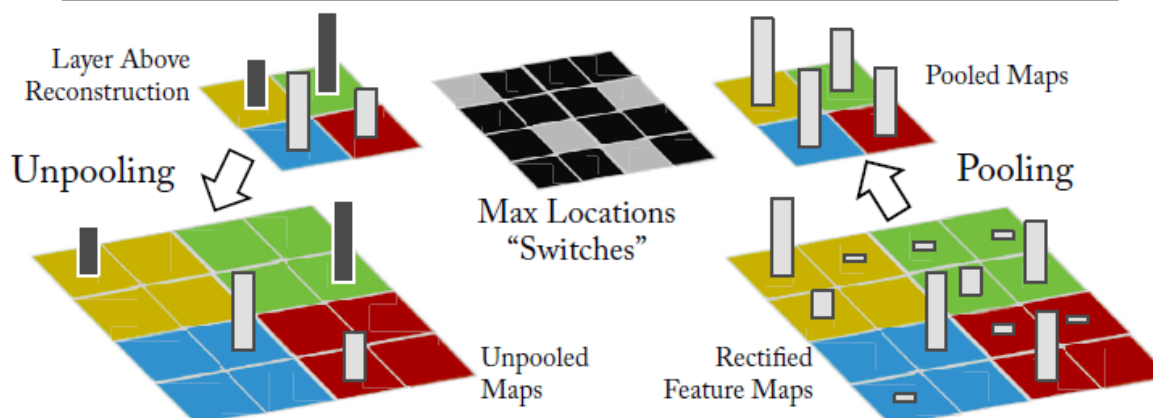


# Visualizing CNNs

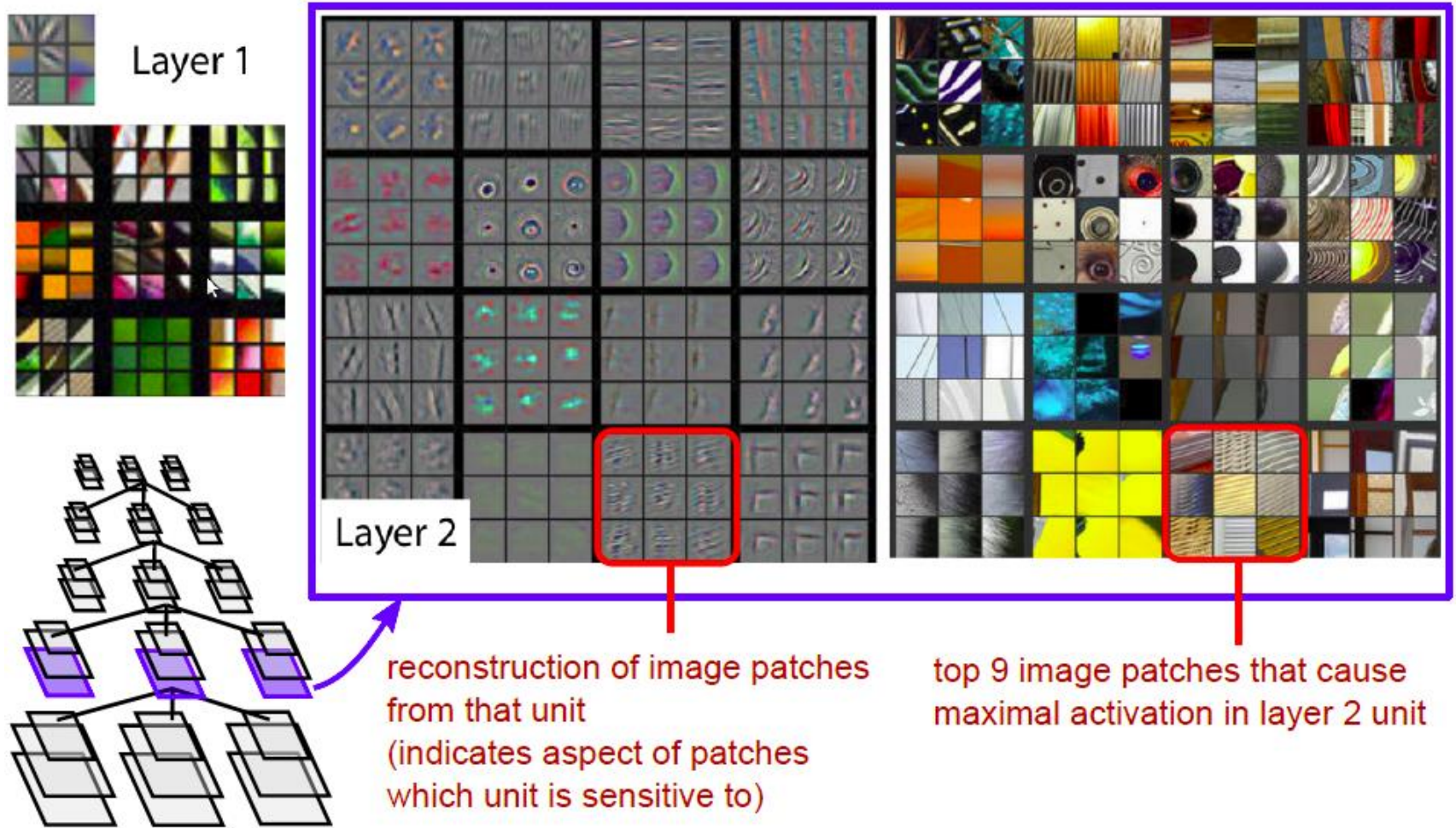
DeconvNet



ConvNet

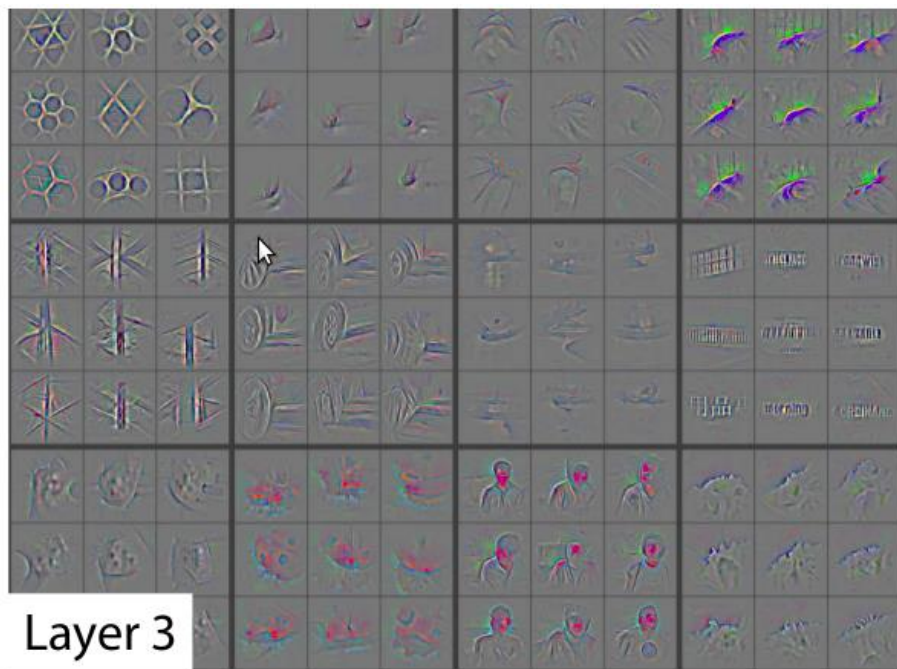


# Visualizing CNNs

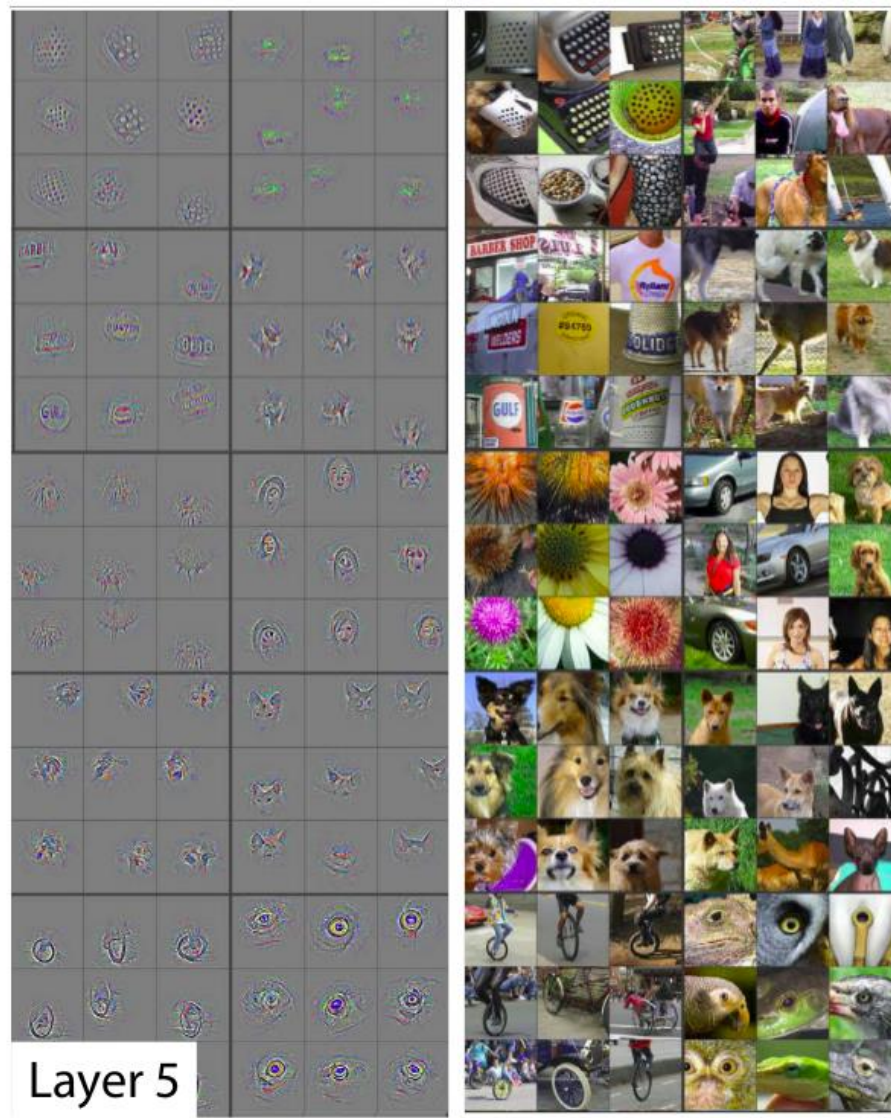
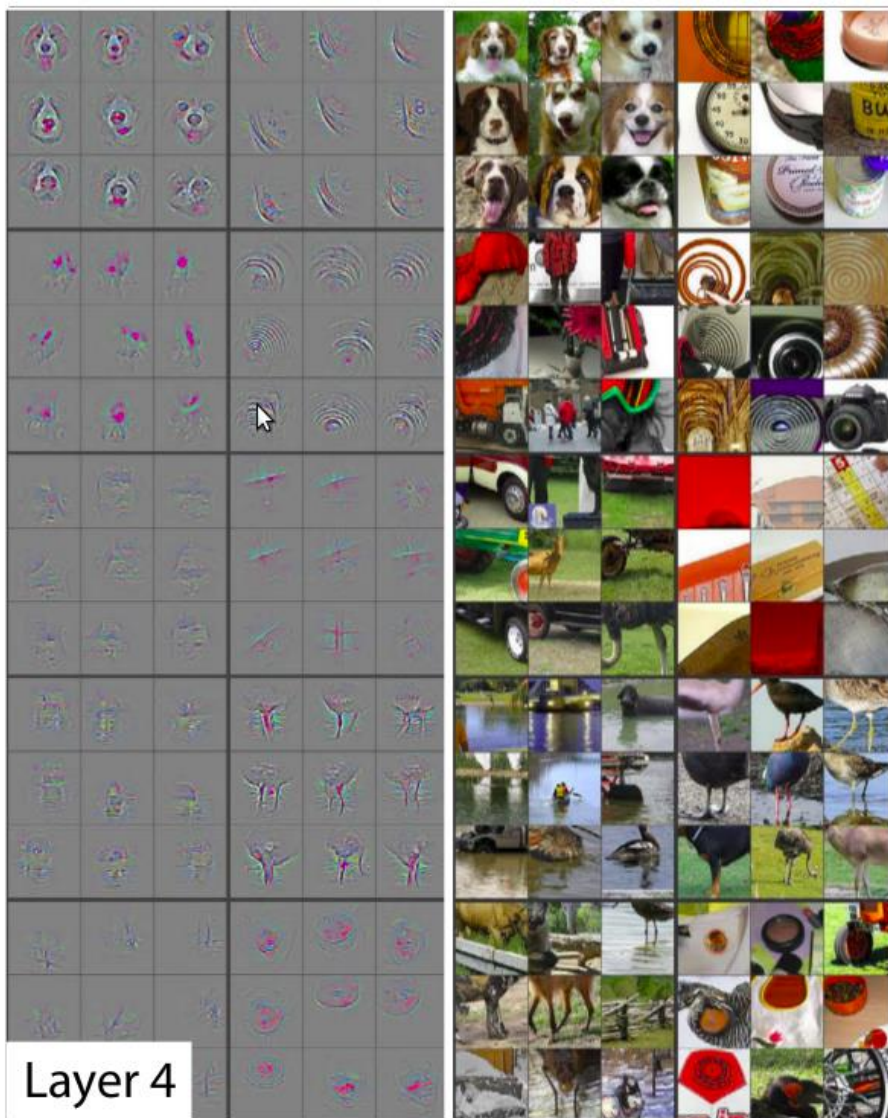


M. Zeiler, R. Fergus, [Visualizing and Understanding Convolutional Neural Networks](#), ECCV 2014.

# Visualizing CNNs



# Visualizing CNNs



# What Does the Network React To?

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

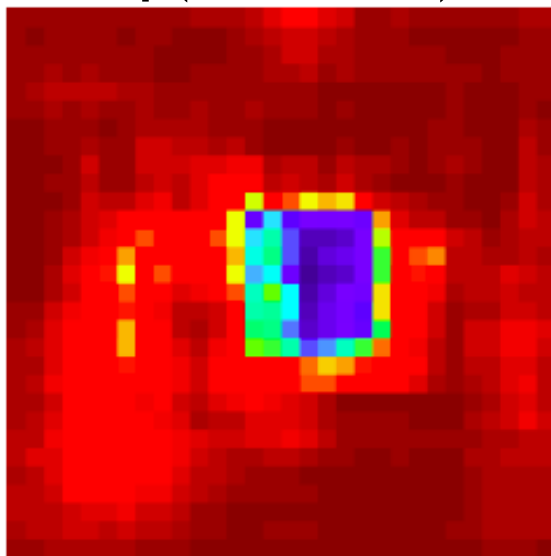


# What Does the Network React To?

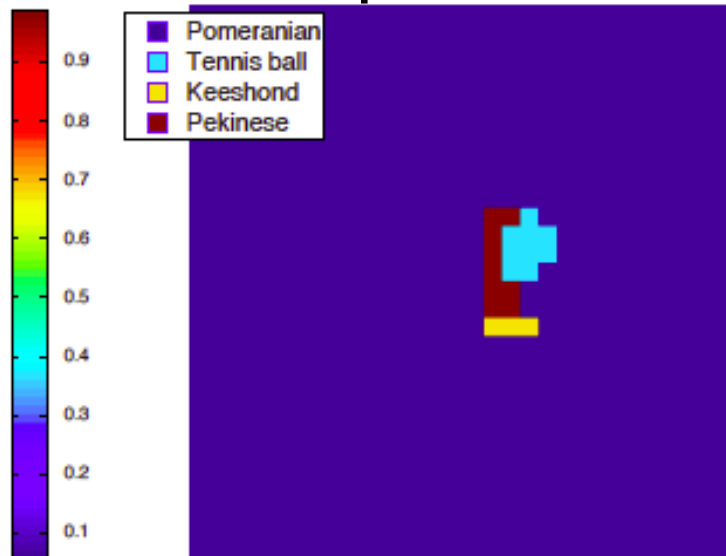
Input image



$p(\text{True class})$



Most probable class

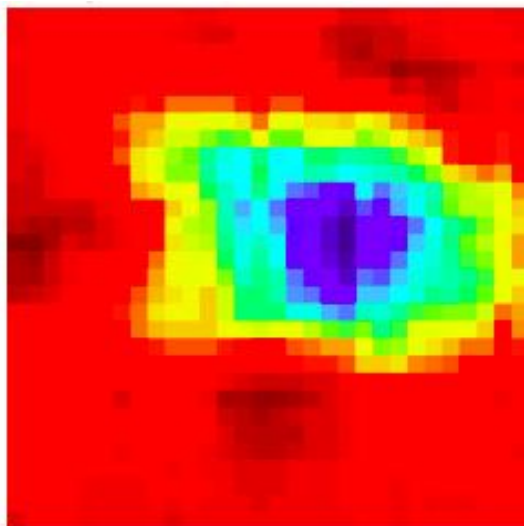


# What Does the Network React To?

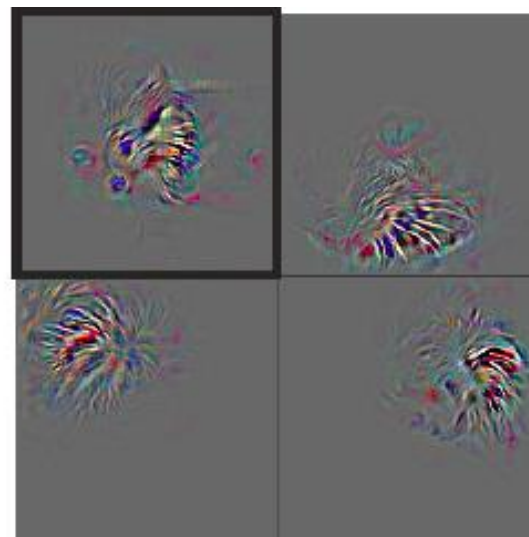
Input image



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



Other activations from the same feature map.

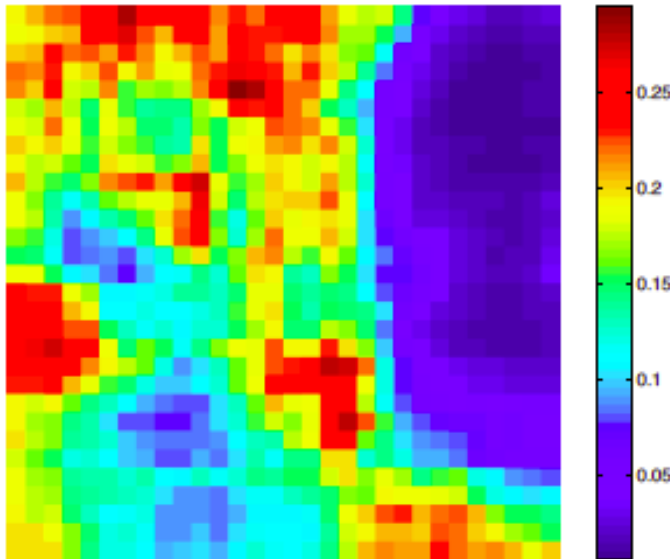


# What Does the Network React To?

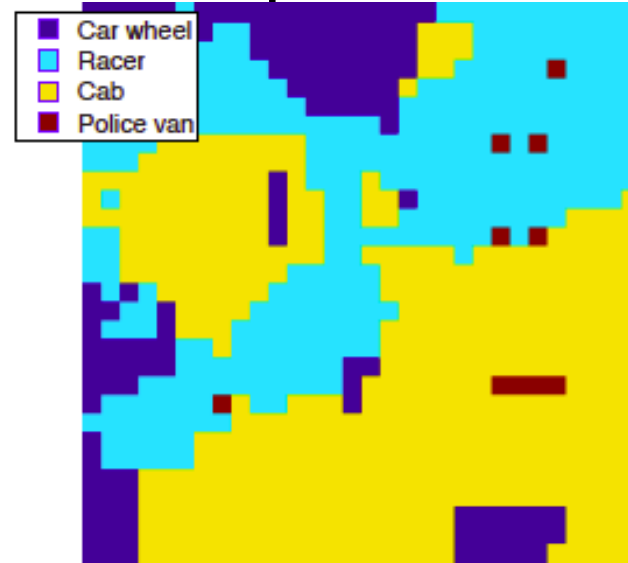
Input image



$p(\text{True class})$



Most probable class



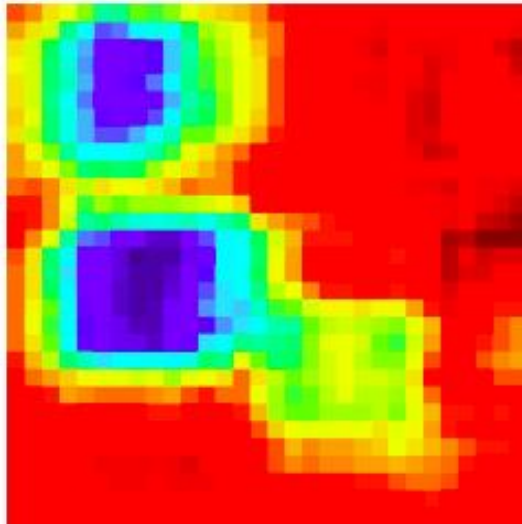


# What Does the Network React To?

Input image



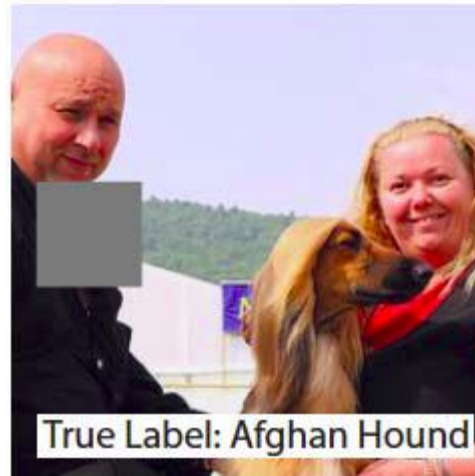
Total activation in most active 5<sup>th</sup> layer feature map



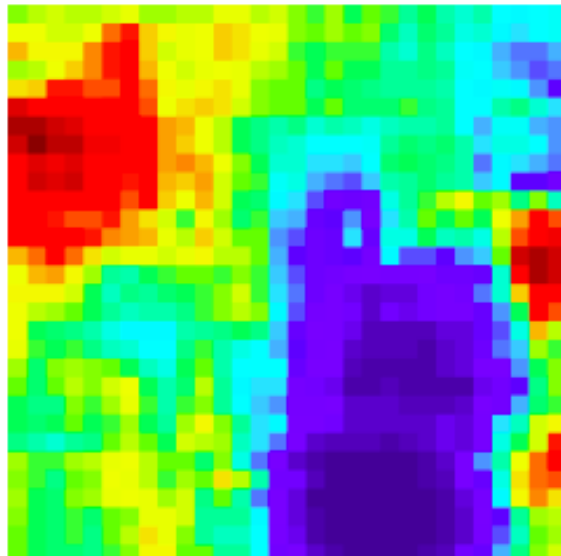
Other activations from the same feature map.

# What Does the Network React To?

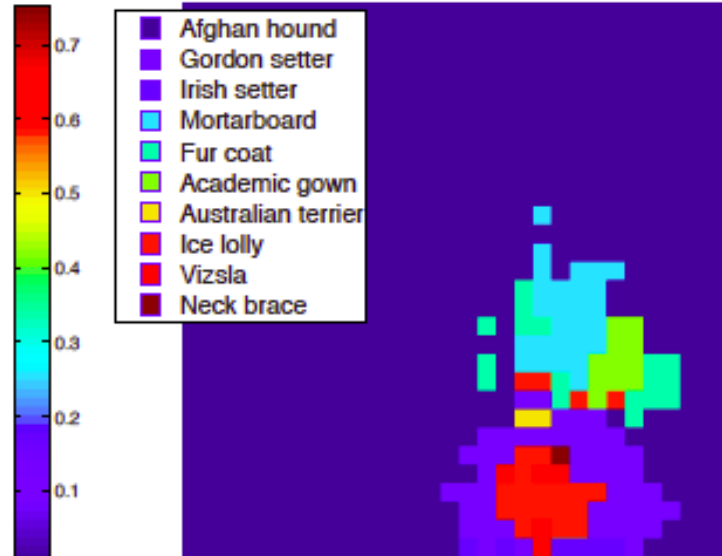
Input image



$p(\text{True class})$

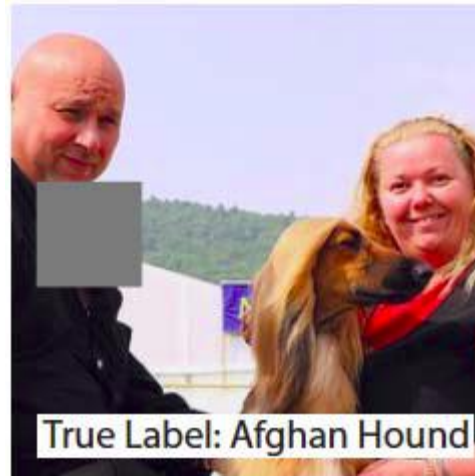


Most probable class

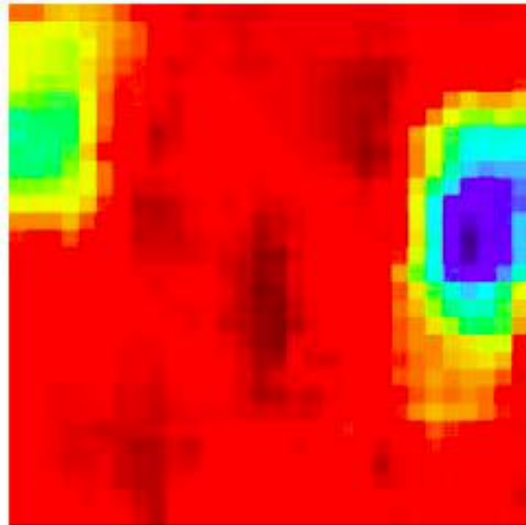


# What Does the Network React To?

Input image



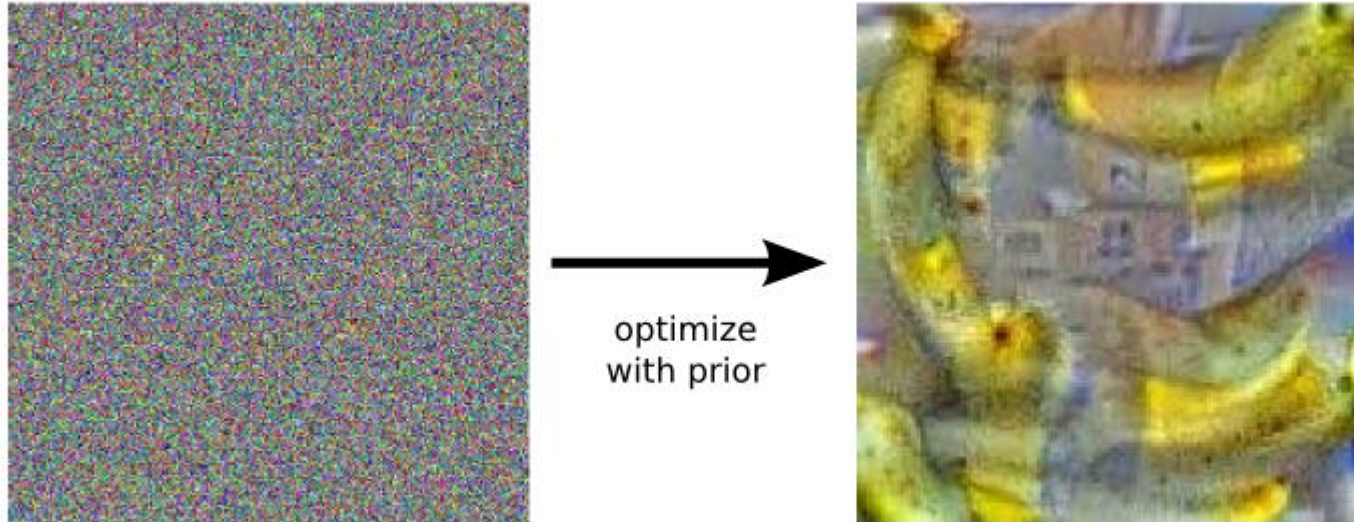
Total activation in most active 5<sup>th</sup> layer feature map



Other activations from the same feature map.



# Inceptionism: Dreaming ConvNets



- **Idea**

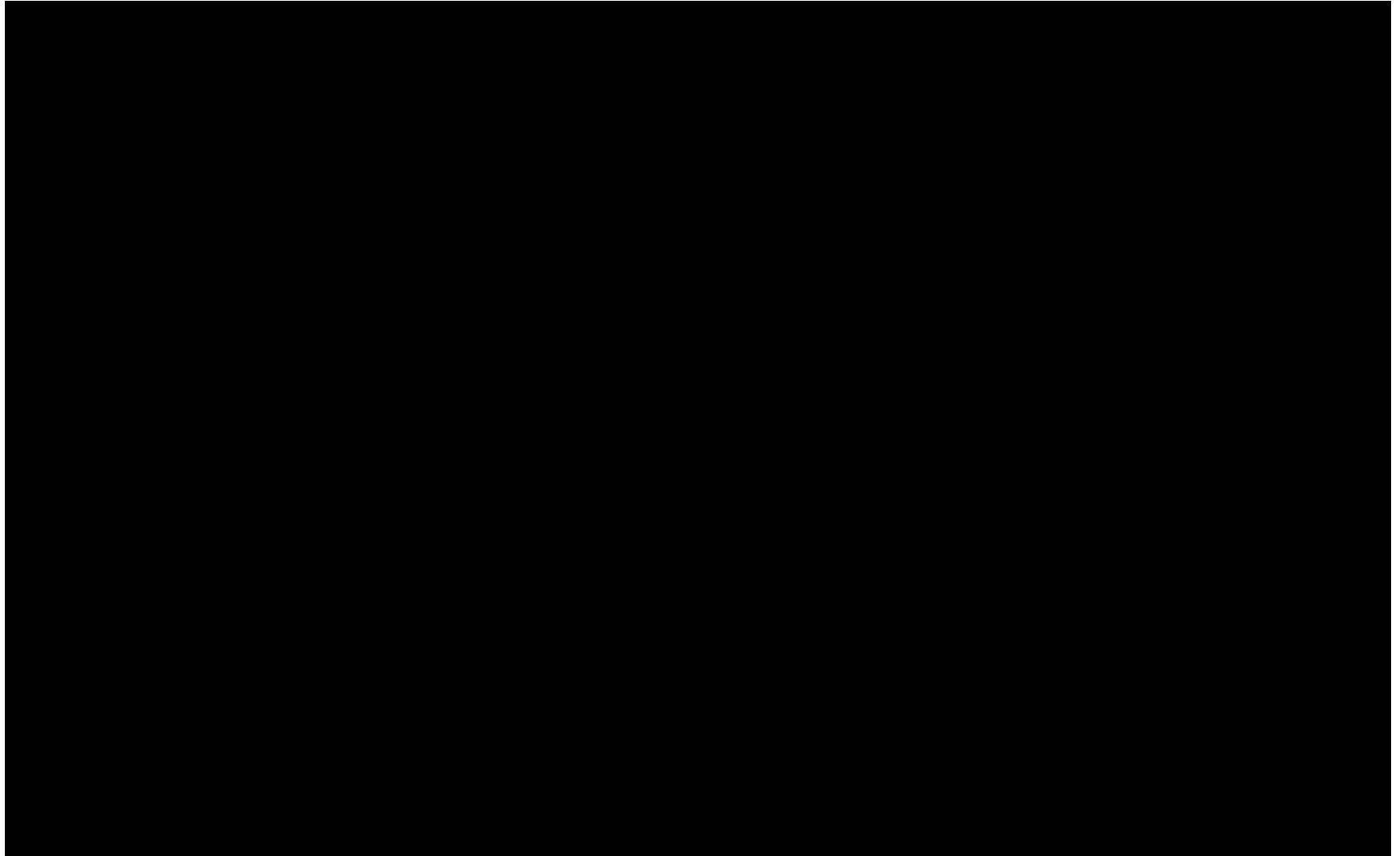
- Start with a random noise image.
  - Enhance the input image such as to enforce a particular response (e.g., banana).
  - Combine with prior constraint that image should have similar statistics as natural images.
- ⇒ Network hallucinates characteristics of the learned class.

# Inceptionism: Dreaming ConvNets

- Results



# Inceptionism: Dreaming ConvNets

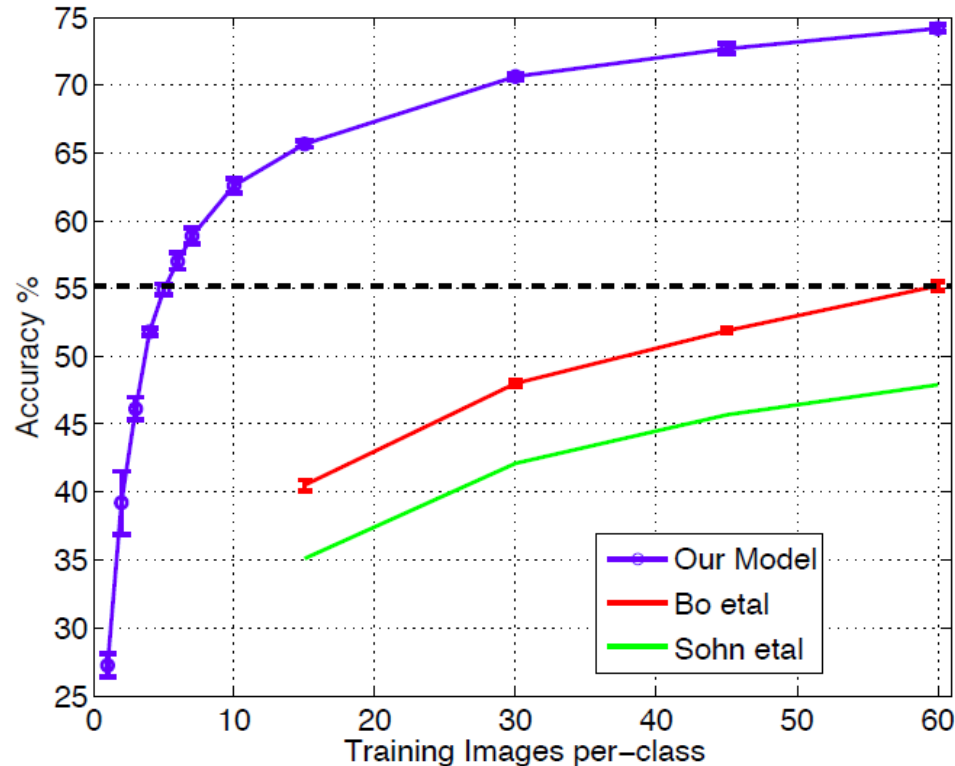


<https://www.youtube.com/watch?v=IREsx-xWQ0g>

# Topics of This Lecture

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

# The Learned Features are Generic



state of the art  
level (pre-CNN)

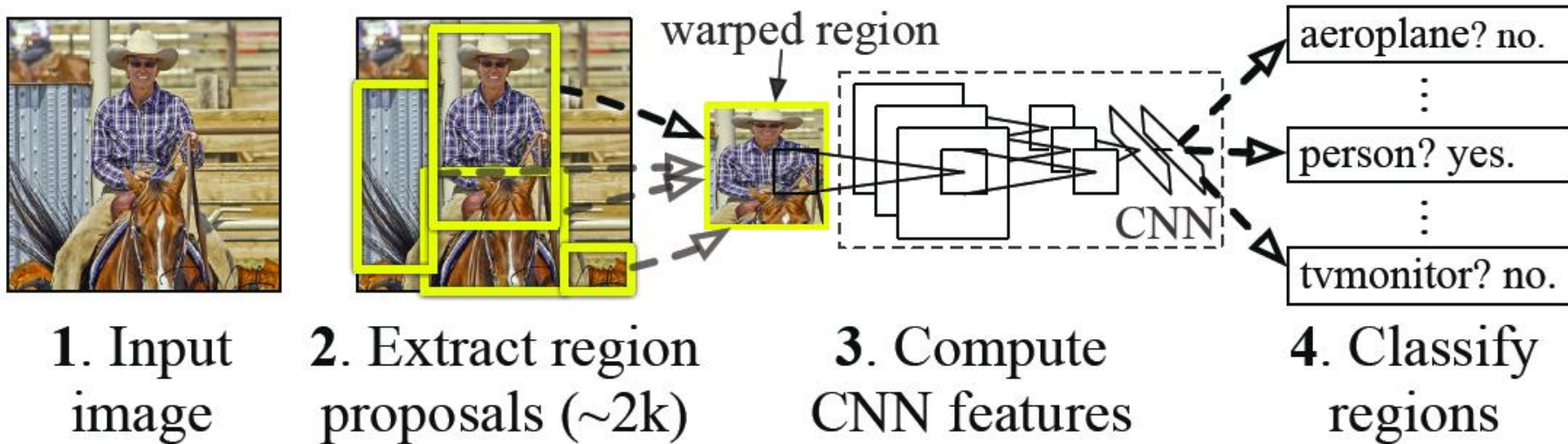
- **Experiment: feature transfer**

- 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



# Other Tasks: Detection

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

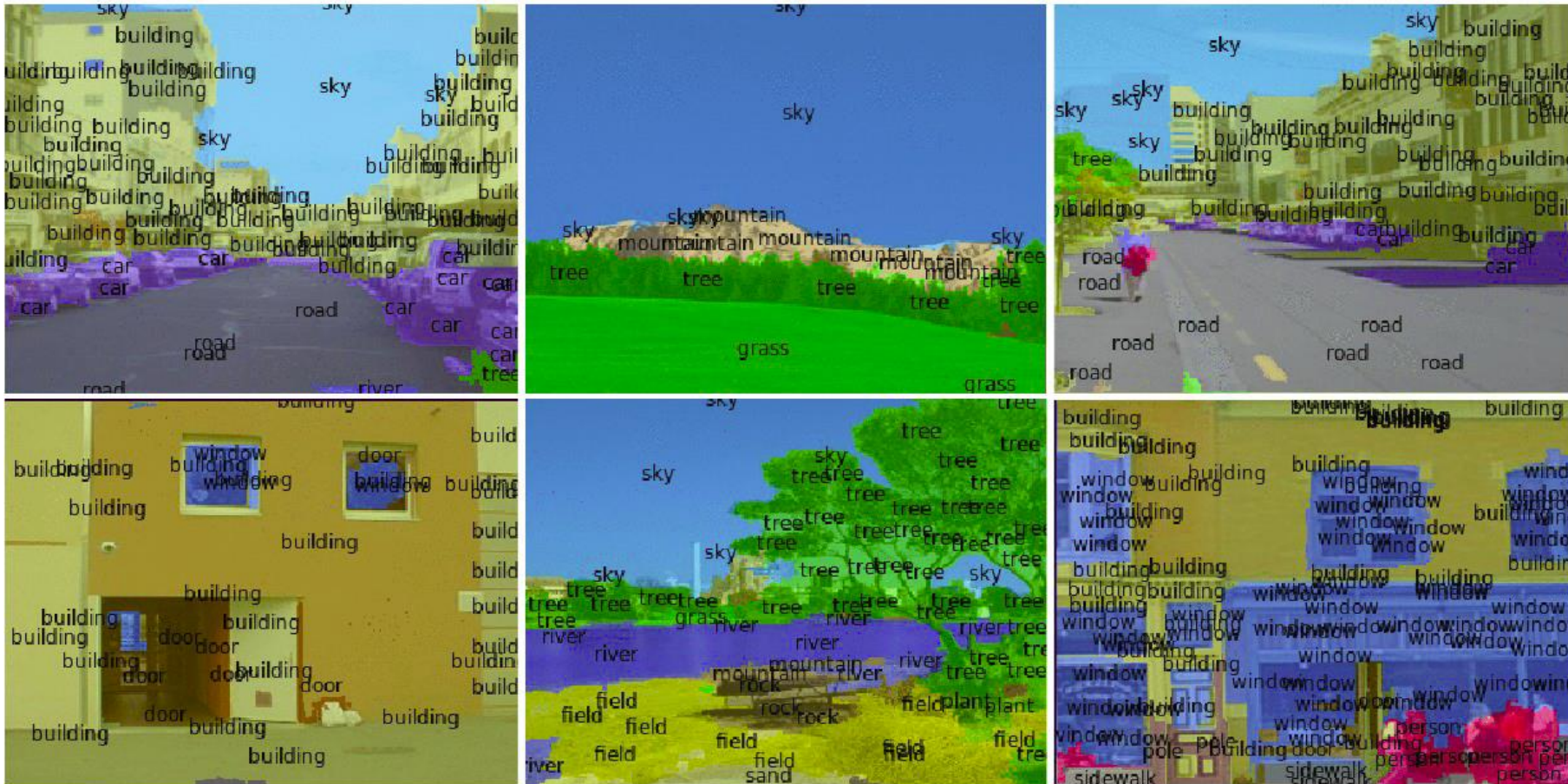


- **Results on PASCAL VOC Detection benchmark**

- **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

# Other Tasks: Semantic Segmentation



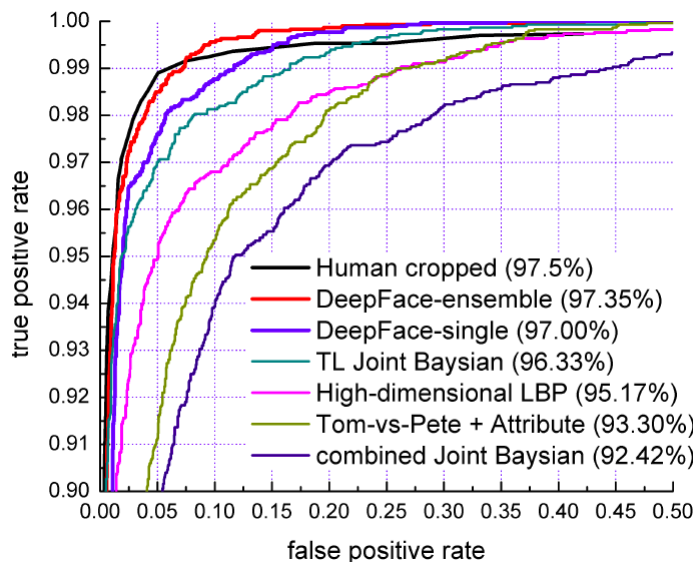
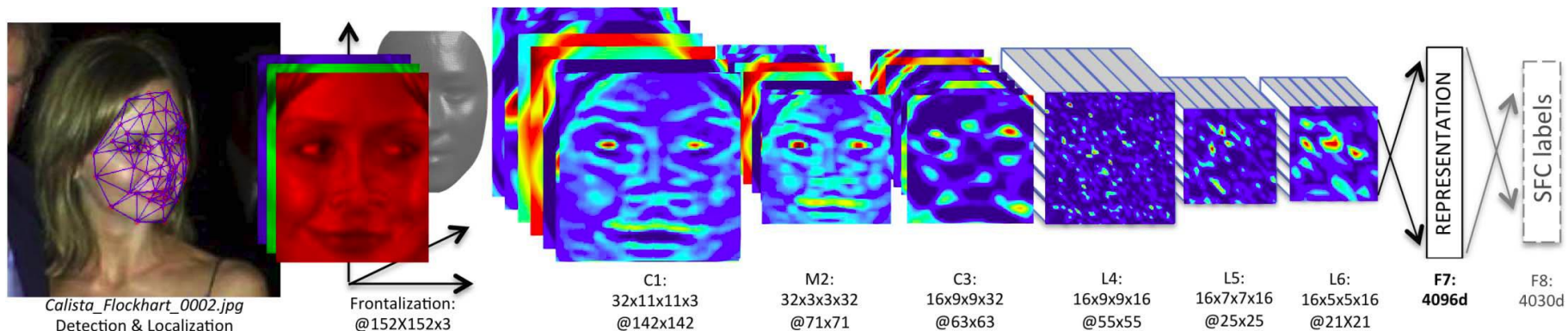
[Farabet et al. ICML 2012, PAMI 2013]

# Other Tasks: Semantic Segmentation



[Farabet et al. ICML 2012, PAMI 2013]

# 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

# 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 ▼

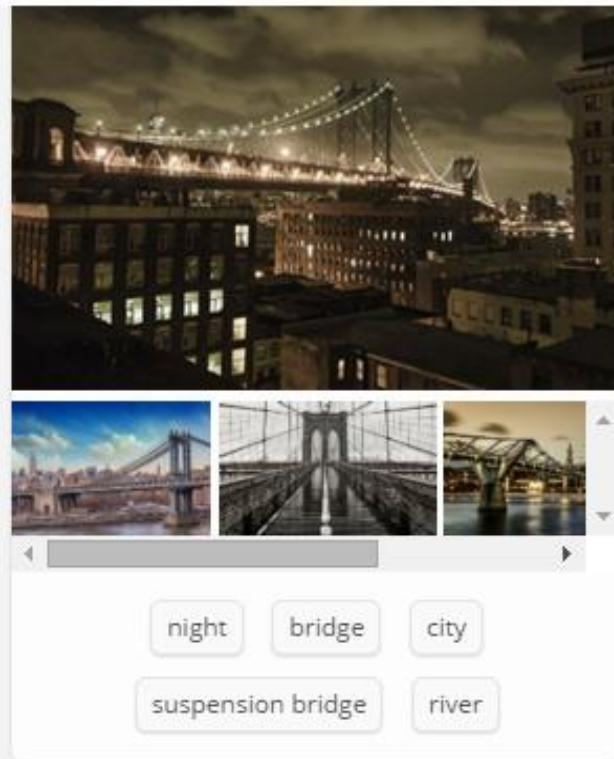
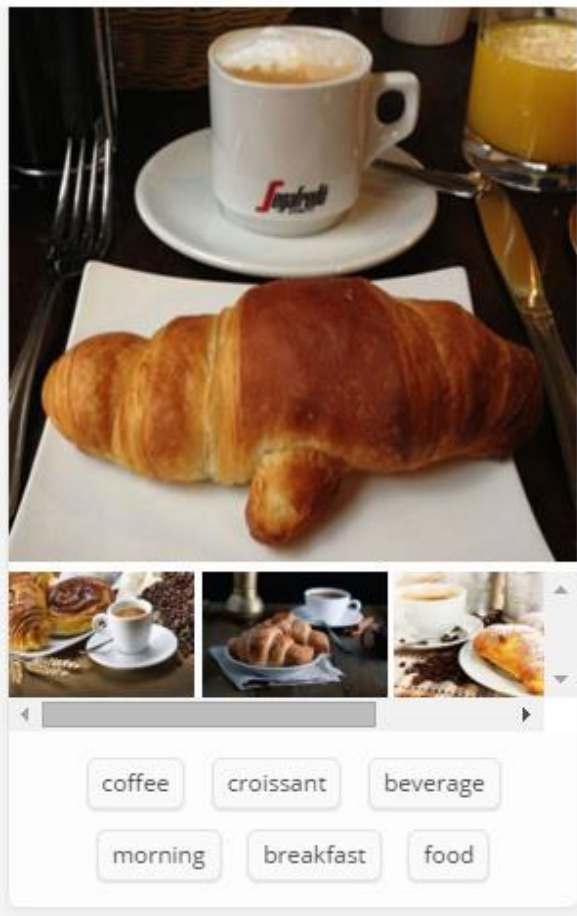
USE THE URL

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...

# Commercial Recognition Services



clarifai

# 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.

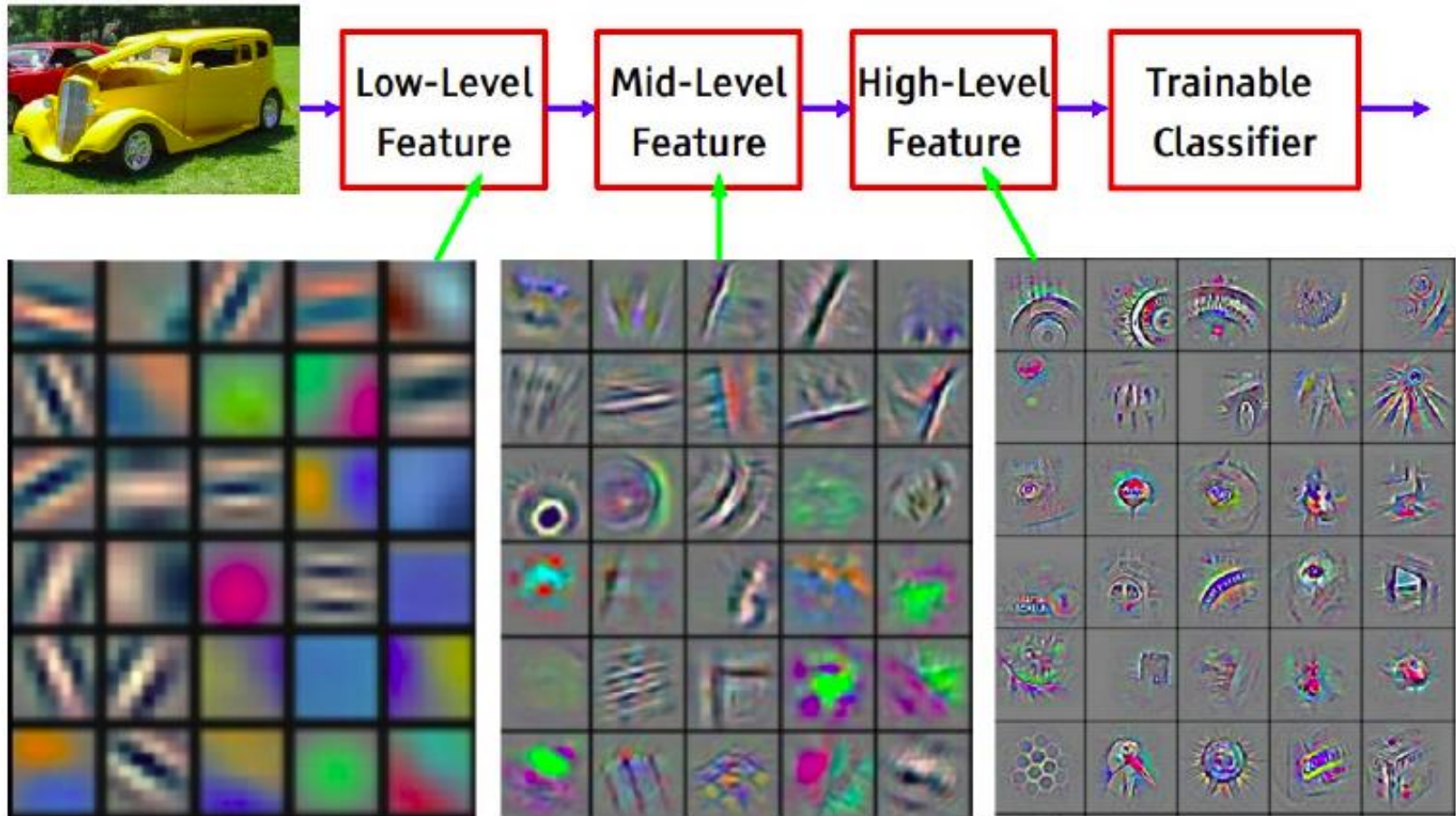
- **VGGNet**

- 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.

# Effect of Multiple Convolution Layers



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