

Computer Vision with CNN

Behnia - Heydari

Amirkabir University of Technology

May 21, 2019

Overview

- 1 Motivation
- 2 ImageNet Challenge
- 3 From Cat's Brain to ResNet (Classification)
- 4 Localization
- 5 Object Detection

Cameras Everywhere

Smartphones

- Exploding number of sensors vs. humans



ConvNets Everywhere

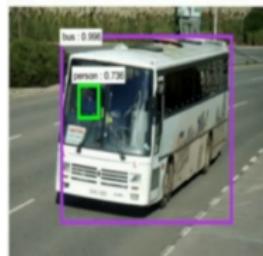
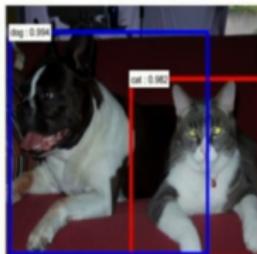
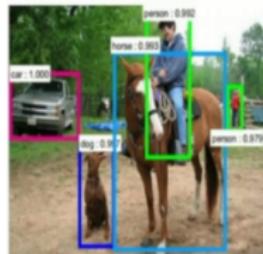


Figure credit: Channing Ran, Elmira Ha, Bruce Gürsoy, Jian Sun, 2015. Reproduced with permission.

Figure credit: Clement Farabet, 2013

ConvNets Everywhere



self-driving cars

Photo by Lane McIntosh. Copyright CS231n 2017.



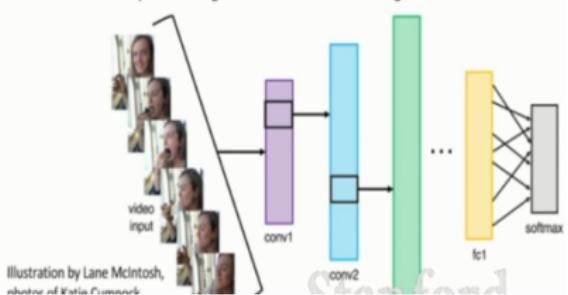
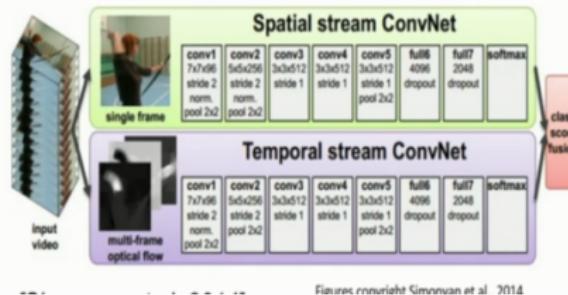
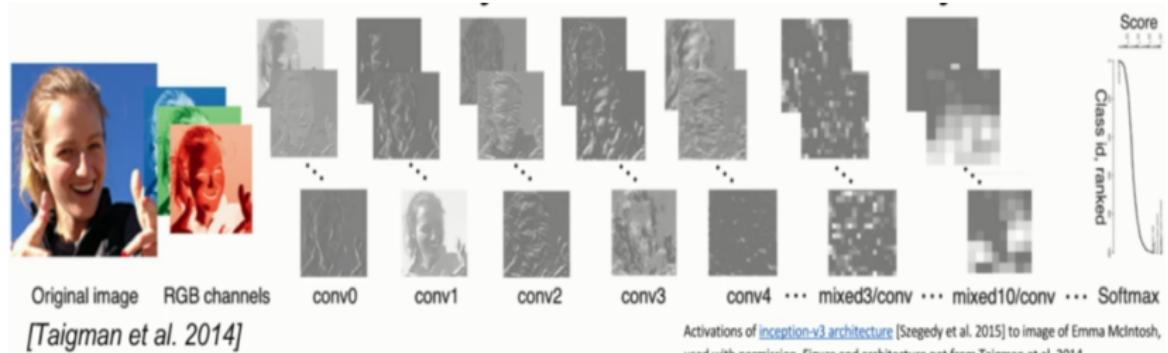
[This image](#) by GBPublic_PR is
licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

(these are the GPUs on rye01.stanford.edu)

Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

ConvNets Everywhere

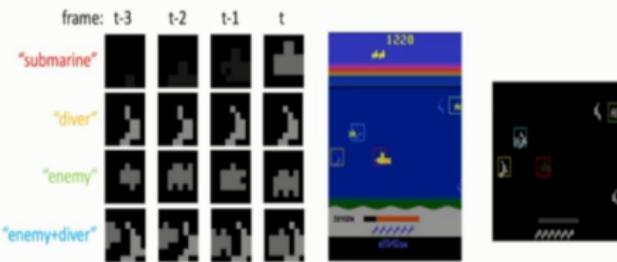


ConvNets Everywhere



Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.

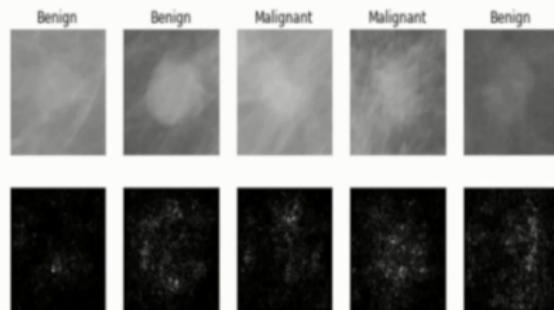
[Toshev, Szegedy 2014]



[Guo et al. 2014]

Figures copyright Xiaoxiao Guo, Satinder Singh, Hon-Jak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reprinted with permission.

ConvNets Everywhere

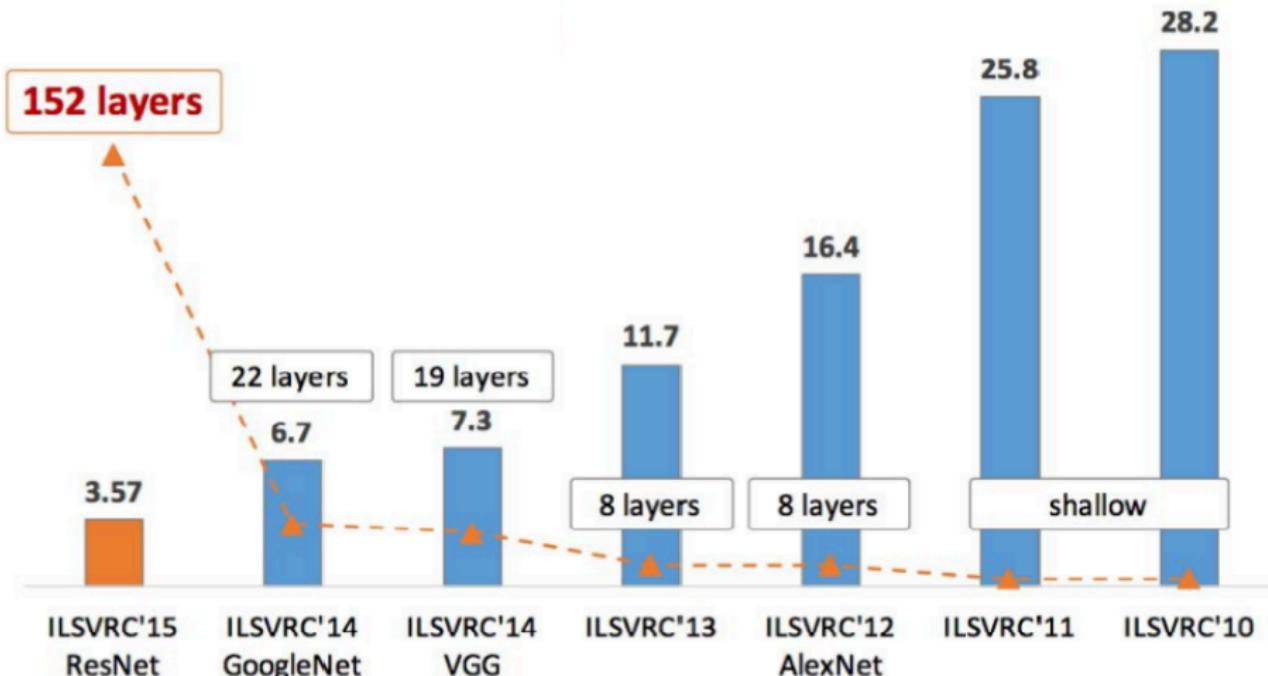


[Sermanet et al. 2011]
Photos by Lane McIntosh,
Copyright CS231n 2017.

Neat Dataset



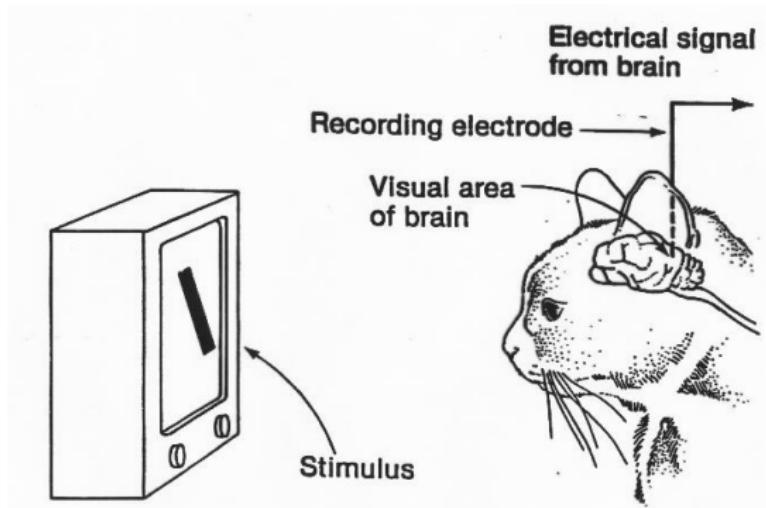
Role of CNN



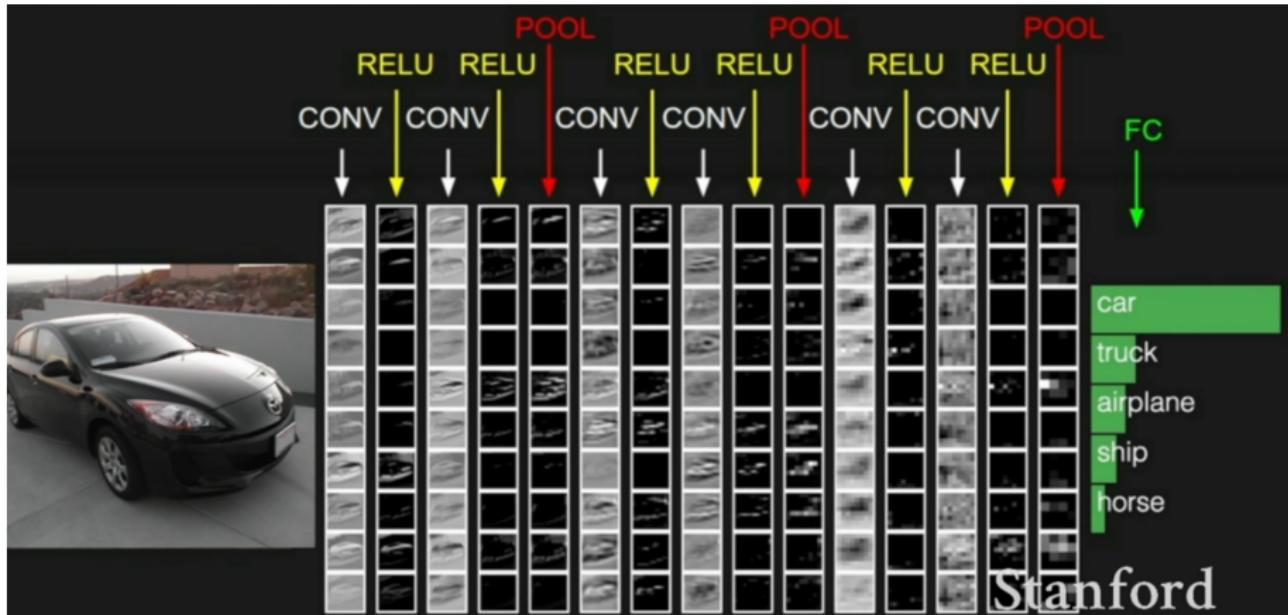
Cat's Brain

Types of cells:

- simple cells
 - complex cells
 - hypercomplex cells



CNN



Convolution

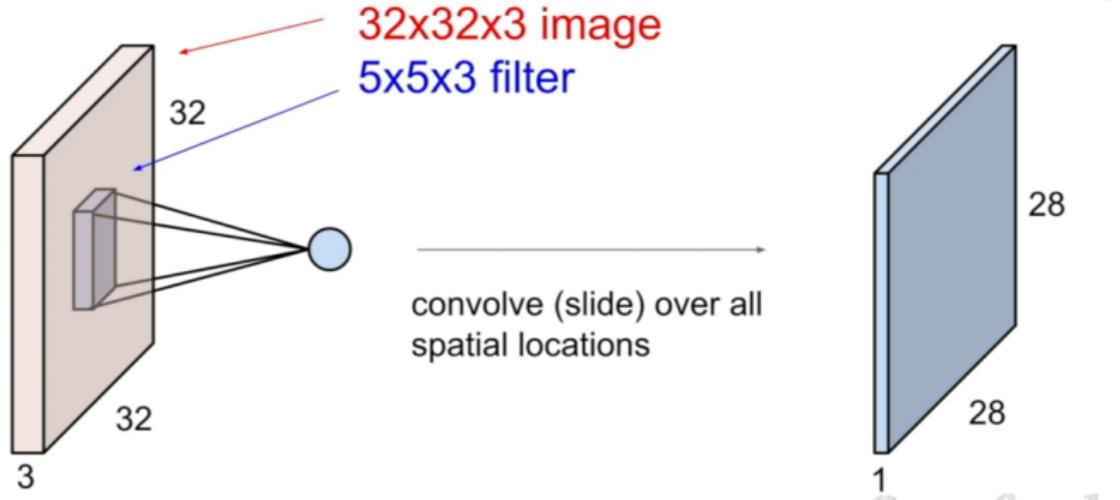


Figure: Convolution Layer

Complexity of Features

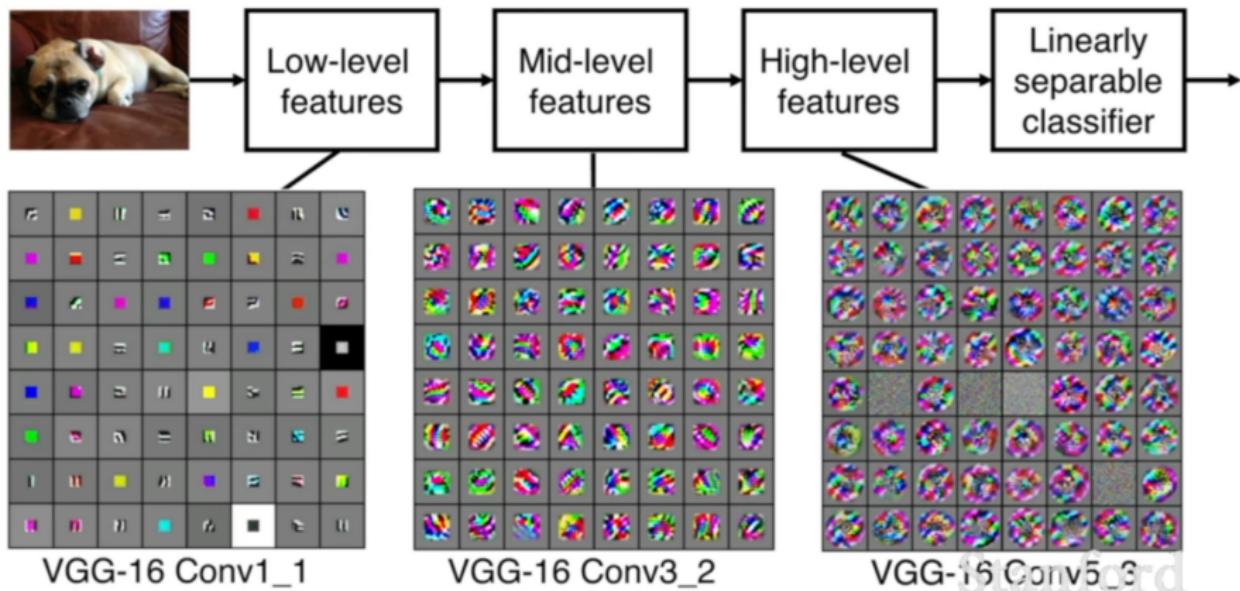


Figure: complexity of features in each layer

Complexity of Features

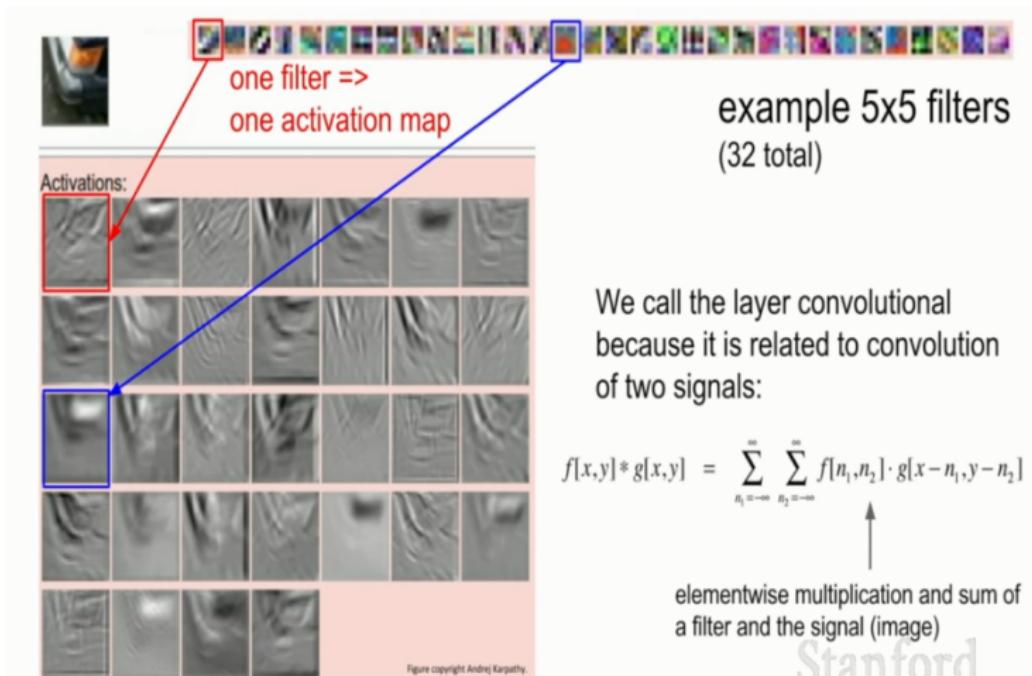


Figure: complexity of each layer output

Sliding a Filter

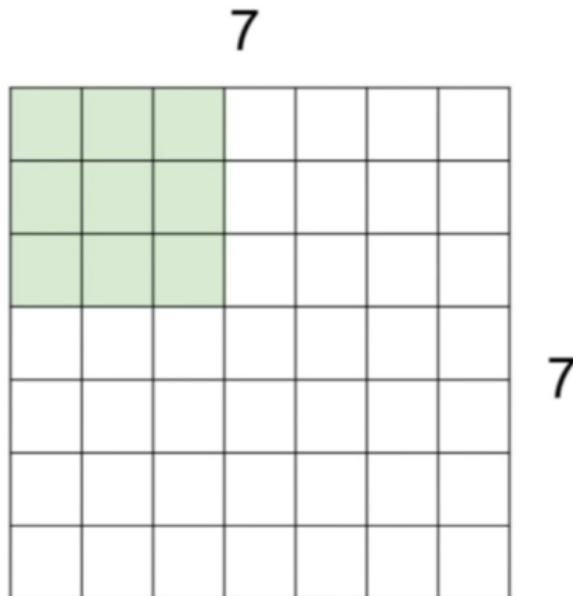


Figure: filtering example

Sliding a Filter

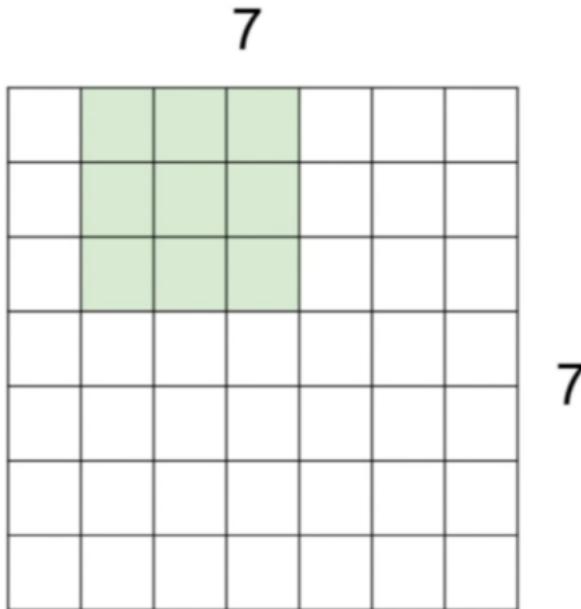


Figure: filtering example

Sliding a Filter

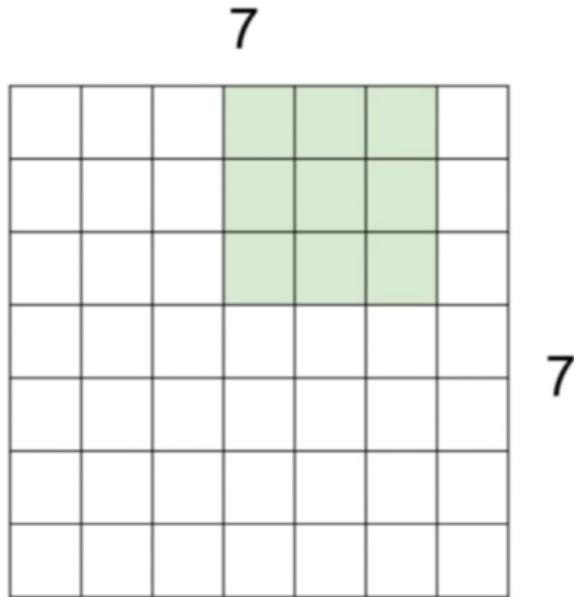


Figure: filtering example

Sliding a Filter

What if using a 3×3 filter with stride 3 ?!

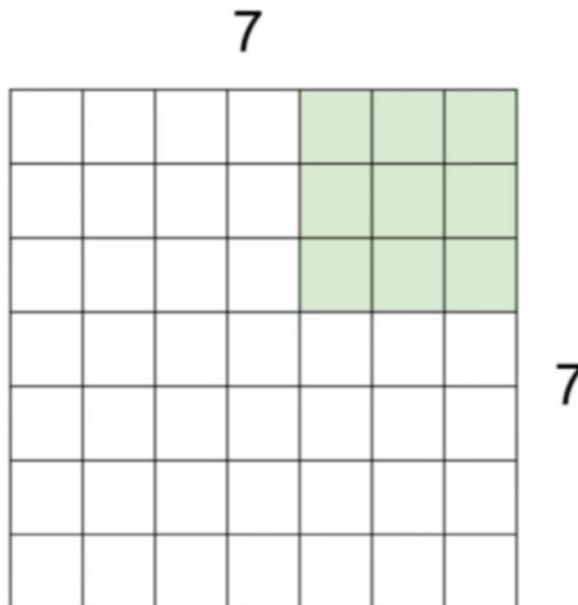


Figure: filtering example

Padding

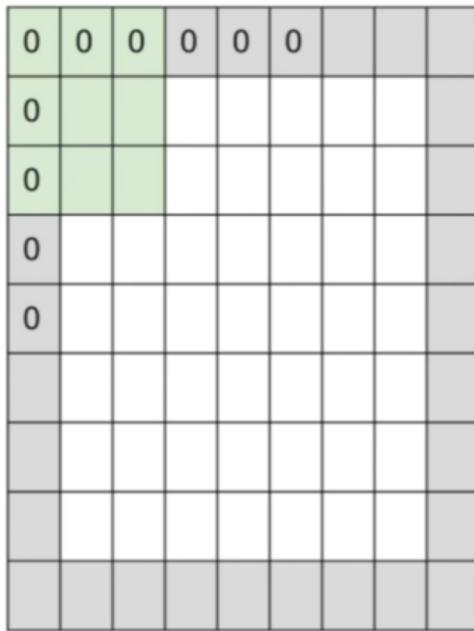


Figure: Padding: control output size

Summary

- Accept a volume of size $N \times N$
- requires 4 hyper parameters:
 - Filter's spatial extent F
 - Filter's Stride S
 - Amount of zero padding P
 - Number of kernels ($\times depth \times$)
- Produces a volume of size $M \times M$
 - $M = \frac{(N-F+2p)}{S} + 1$
 - Number of parameters : $(F \times F \times D + b) \times k$

CNN

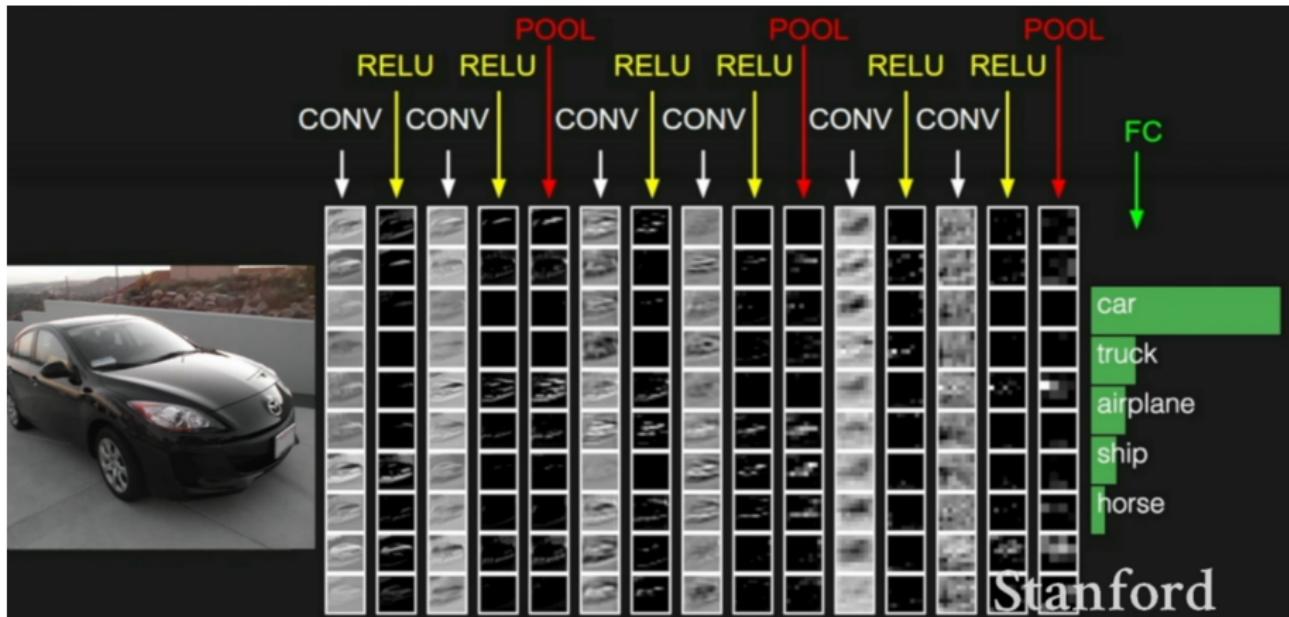


Figure: Fully Connected Layer

Pooling

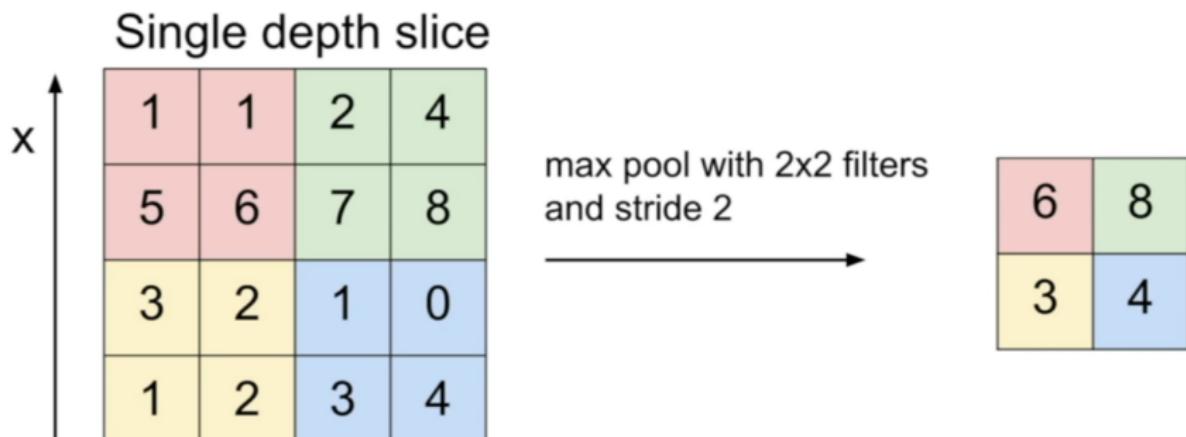


Figure: Max-pooling

CNN

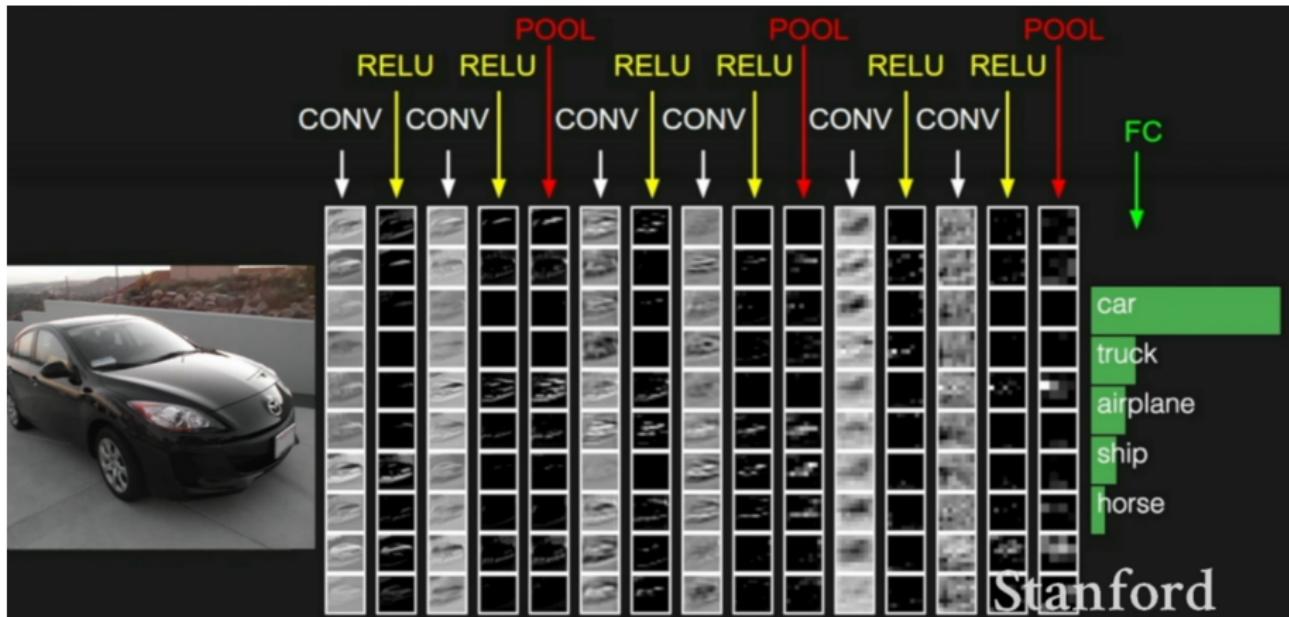


Figure: Fully Connected Layer

Comparing

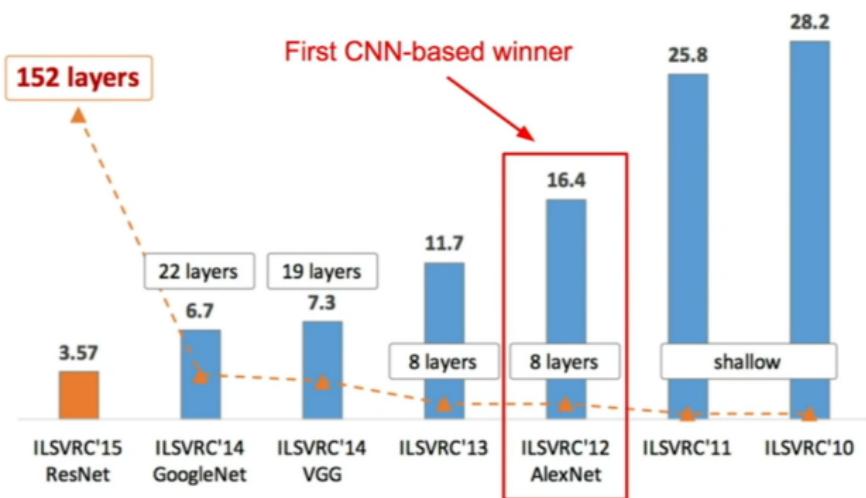


Figure: Comparing Structures

AlexNet

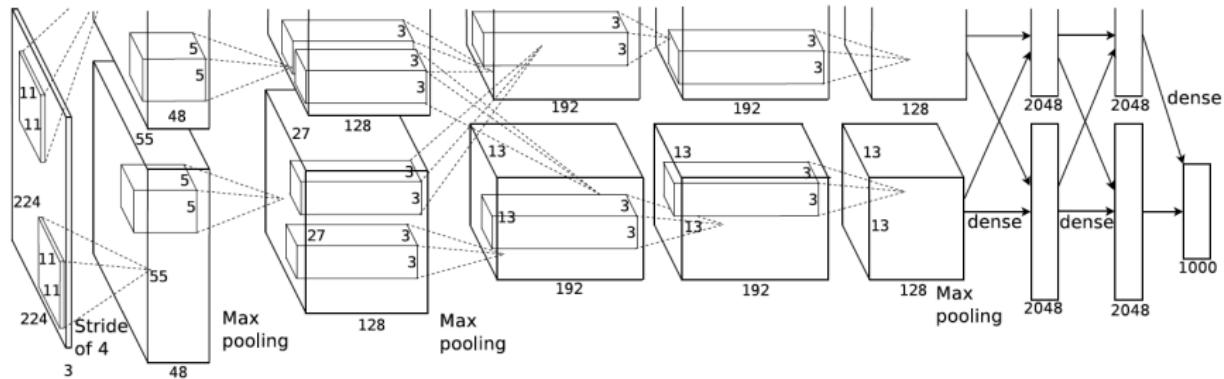
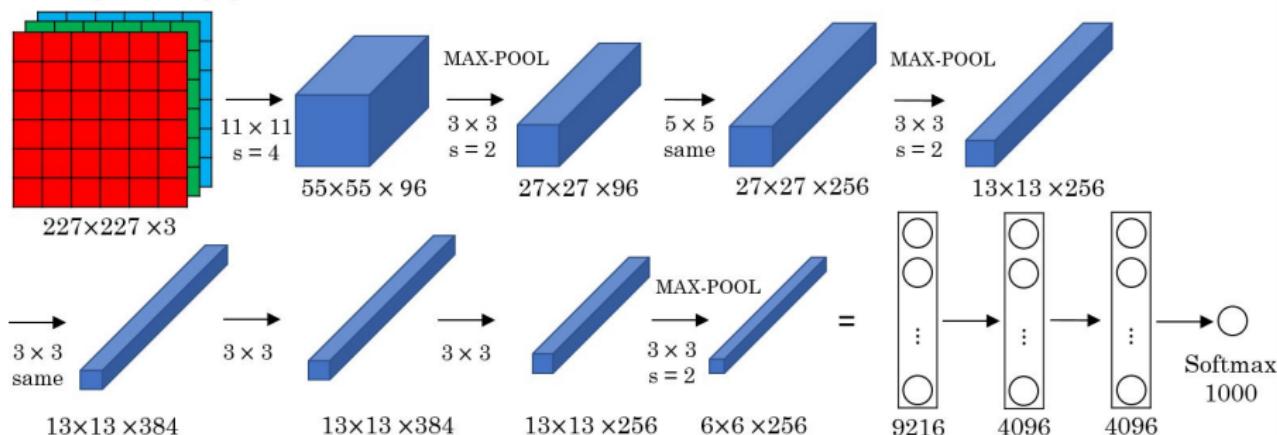


Figure: AlexNet architecture

AlexNet

AlexNet



Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks

Andrew Ng

Figure: AlexNet architecture

Comparing

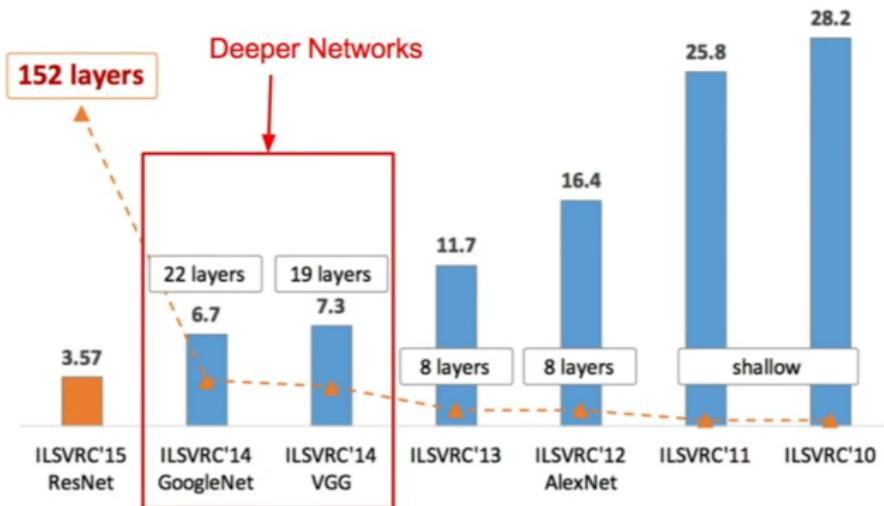
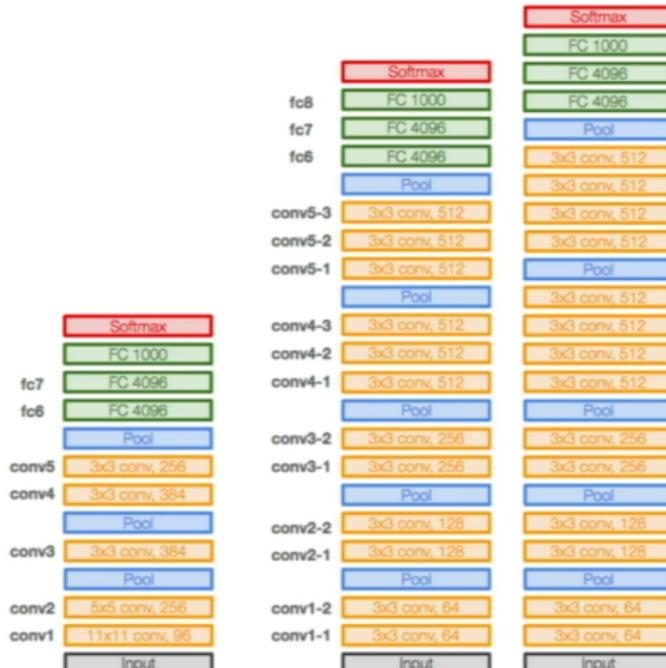


Figure: Comparing Structures

Deeper Networks, Smaller Filters



AlexNet

VGG16

VGG19

Comparing

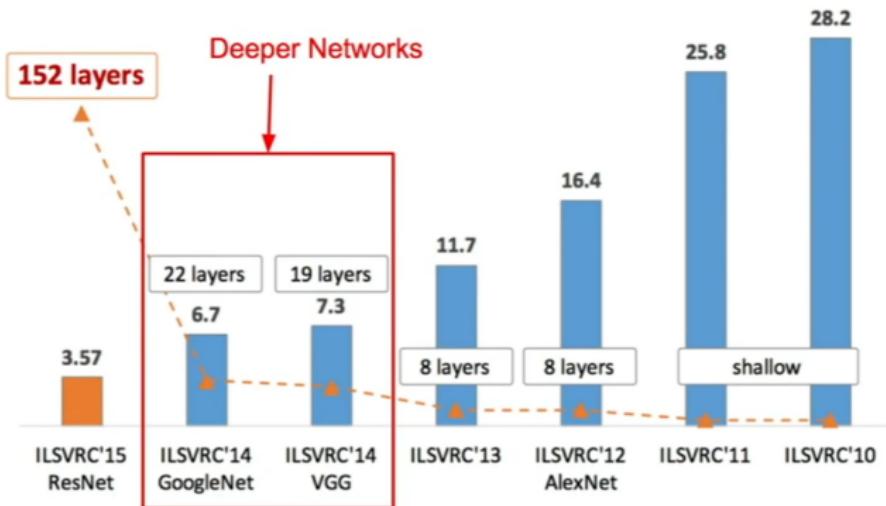
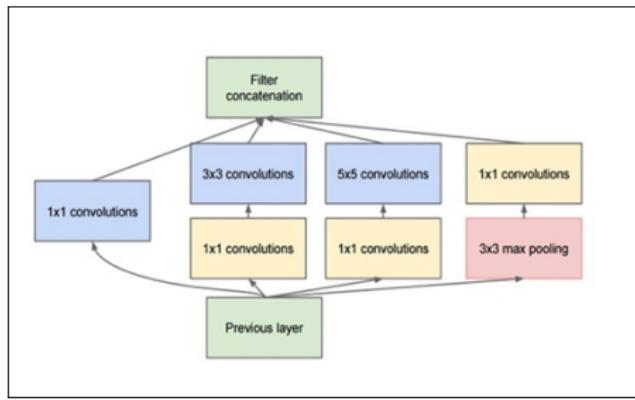
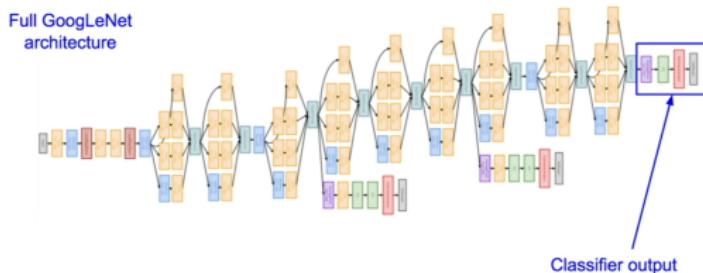


Figure: Comparing Structures

GoogLeNet

Make deeper networks computationally inexpensive



Comparing

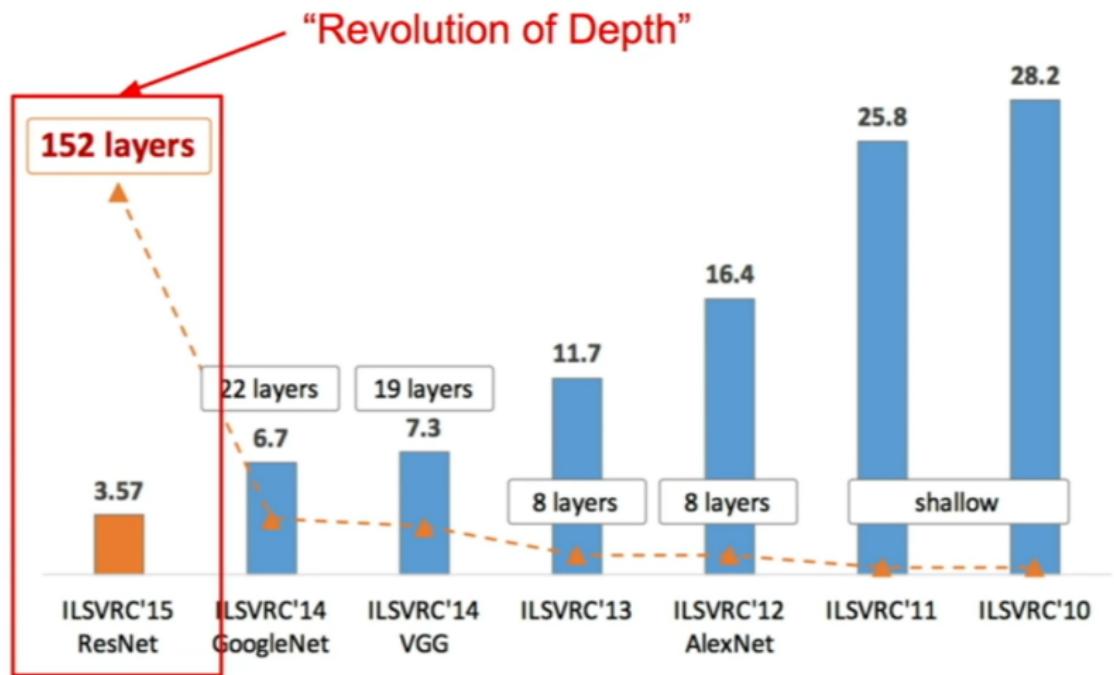
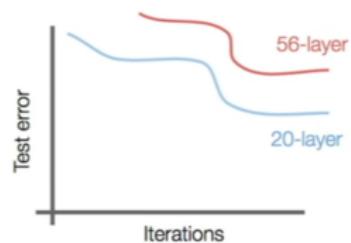
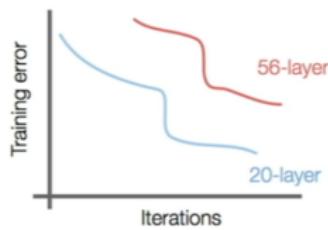


Figure: Comparing Structures

What happens when we continue stacking deeper layers on a “plain” convolutional neural network?



ResNet

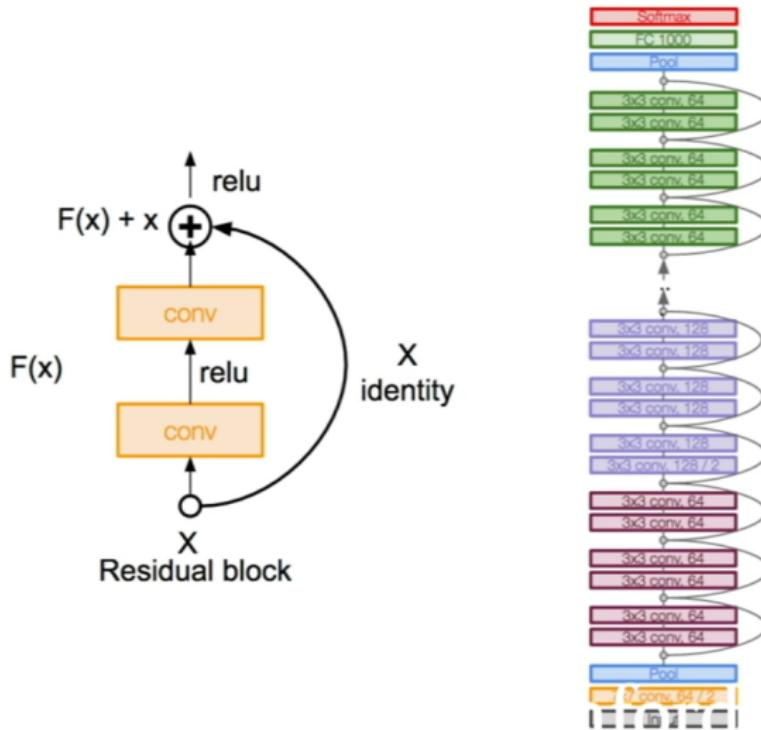


Figure: ResNet Architecture

Comparing

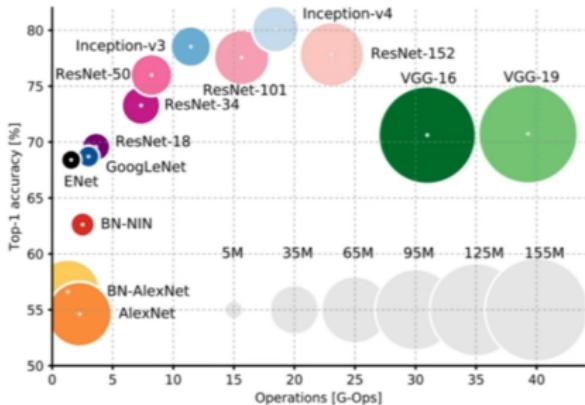
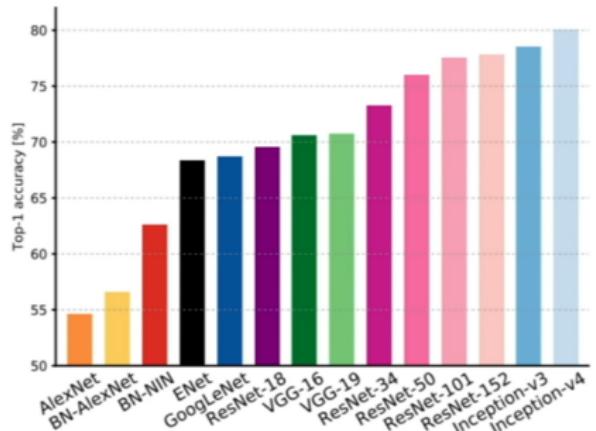


Figure: Comparison

Example

ConvNetJS CIFAR-10 demo

Description

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to parse the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).

Training Stats

pause

Forward time per example: 12ms
Backprop time per example: 28ms
Classification loss: 1.38449
L2 Weight decay loss: 0.04078
Training accuracy: 0.51
Validation accuracy: 0.43
Examples seen: 54991

Loss:



change

Momentum:

change

Batch size:

change

Weight decay:

change

clear graph

save network snapshot as JSON

init network from JSON snapshot

load a pretrained network (achieves ~80% accuracy)

Instantiate a Network and Trainer

```
layer_defs = [{}]
layer_defs.push(type:'input', out_sz:32, out_ty:32, out_depth:3);
layer_defs.push(type:'conv', sz:5, filters:16, stride:1, pad:2, activation:'relu');
layer_defs.push(type:'pool', sz:2, stride:2);
layer_defs.push(type:'conv', sz:5, filters:32, stride:1, pad:2, activation:'relu');
layer_defs.push(type:'pool', sz:2, stride:2);
layer_defs.push(type:'conv', sz:5, filters:20, stride:1, pad:2, activation:'relu');
layer_defs.push(type:'pool', sz:2, stride:2);
layer_defs.push(type:'softmax', num_classes:10);
```

```
net = new convnetjs.Net();
net.addLayers(layer_defs);
|
```

Computer Vision Tasks

Classification + Localization



GRASS, CAT,
TREE, SKY

No objects, just pixels



CAT

Single Object



DOG, DOG, CAT

Multiple Object

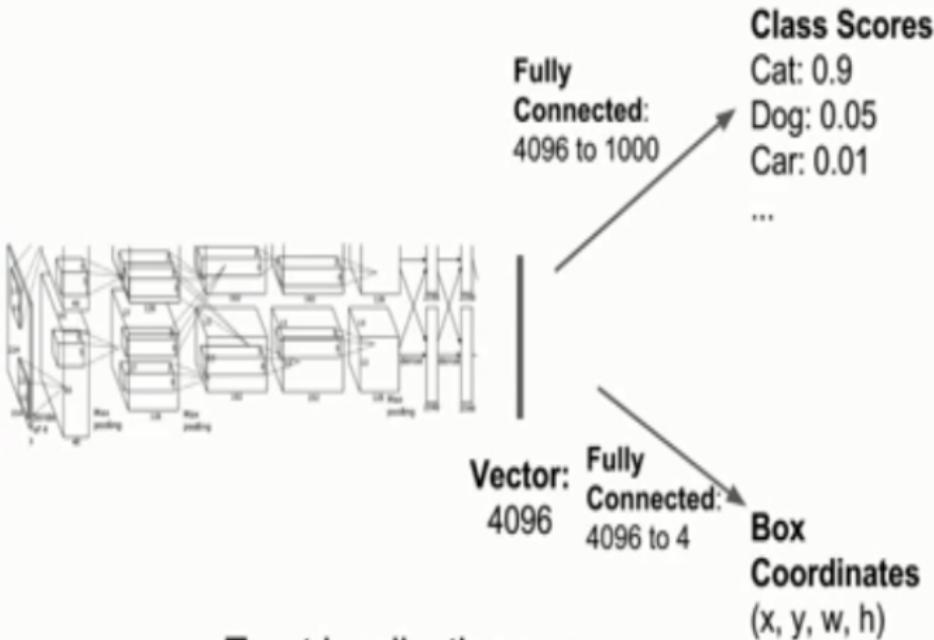


DOG, DOG, CAT

Classification + Localization

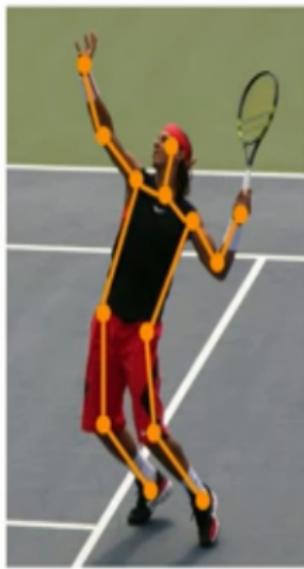


This image is CC0 public domain



Treat localization as a regression problem!

Example



This image is licensed under CC-BY 2.0.

Represent pose as a set of 14 joint positions:

- Left / right foot
- Left / right knee
- Left / right hip
- Left / right shoulder
- Left / right elbow
- Left / right hand
- Neck
- Head top

Figure: Body pose estimation

Fancier Task

Object Detection



GRASS, CAT,
TREE, SKY

No objects, just pixels



CAT

Single Object



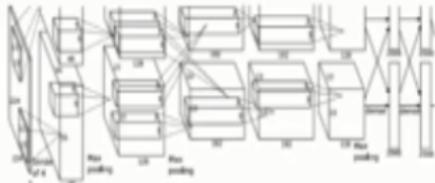
DOG, DOG, CAT

Multiple Object

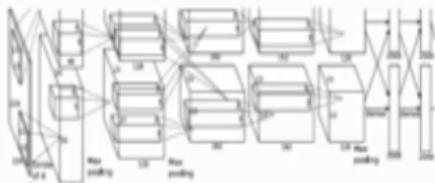


DOG, DOG, CAT

Object Detection



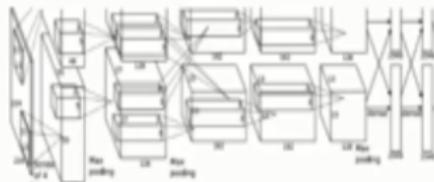
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



DUCK: (x, y, w, h)

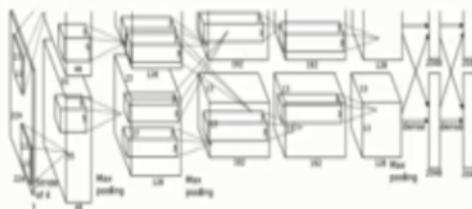
DUCK: (x, y, w, h)

...

Figure: New approach is needed

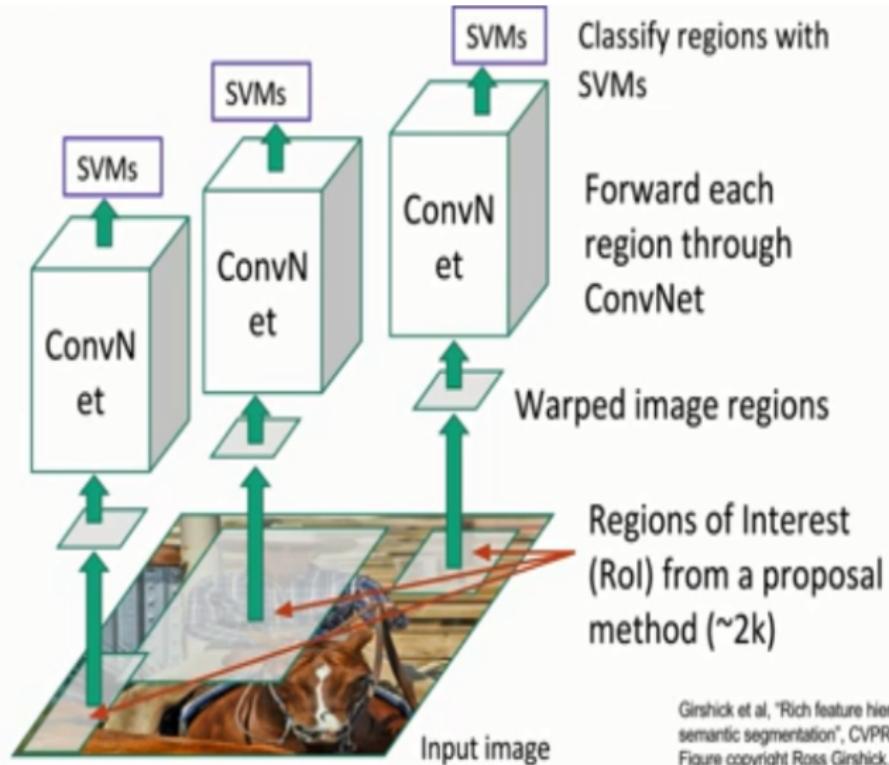
Object Detection

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

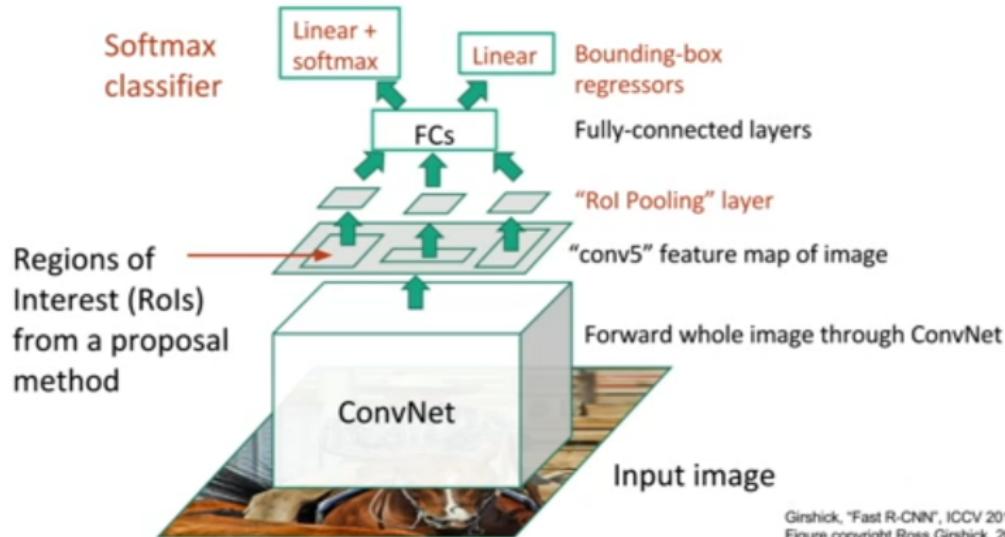
RCNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [http://csail.mit.edu/%7Egirshic/fcns.html](#) (FCC BY-NC-ND license).

Fast RCNN

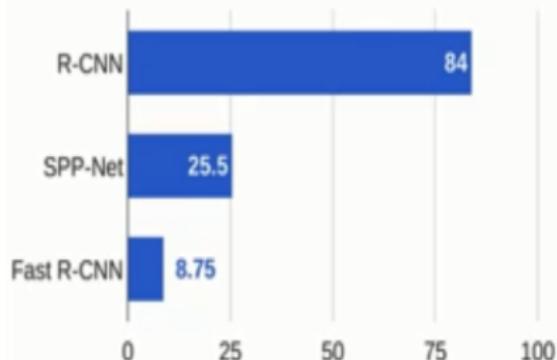
Fast R-CNN



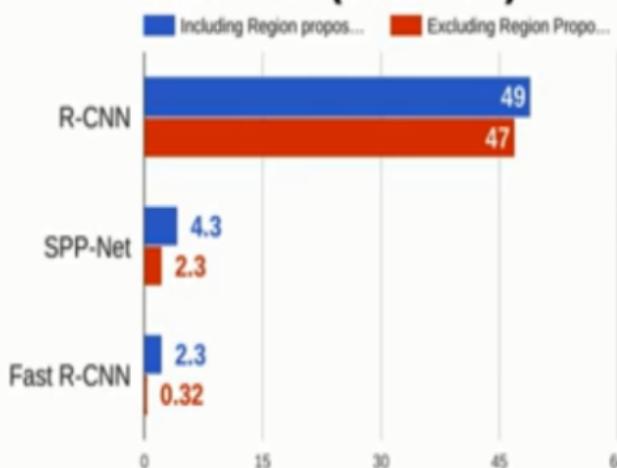
Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015. <http://csail.mit.edu/people/girshick/pubs.html>

Comparing

Training time (Hours)



Test time (seconds)



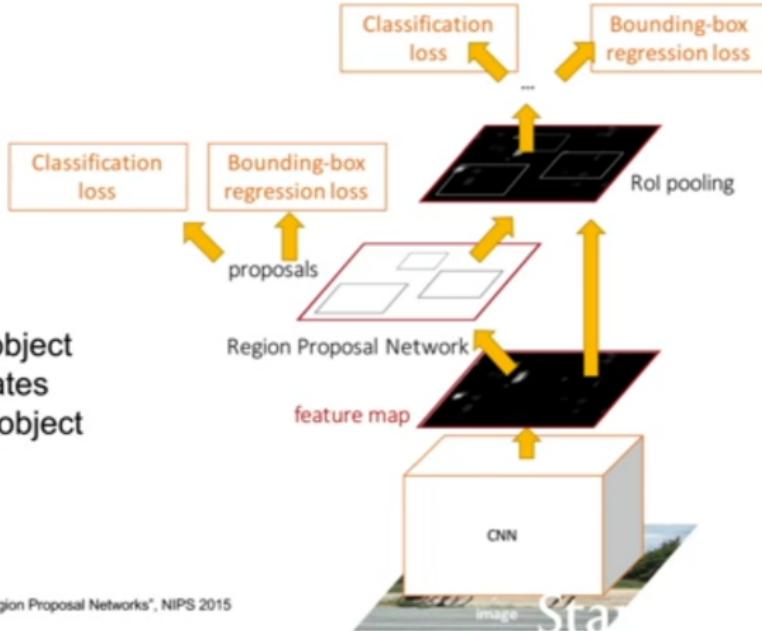
Even Faster

Faster R-CNN: Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

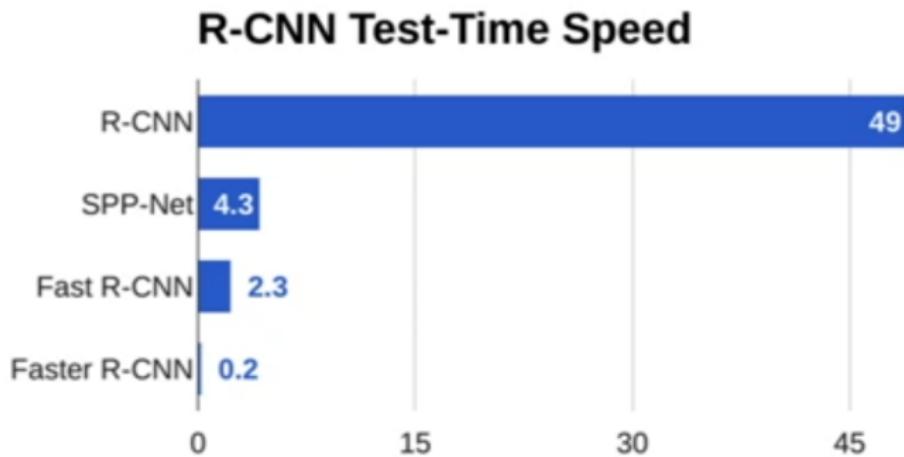
Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

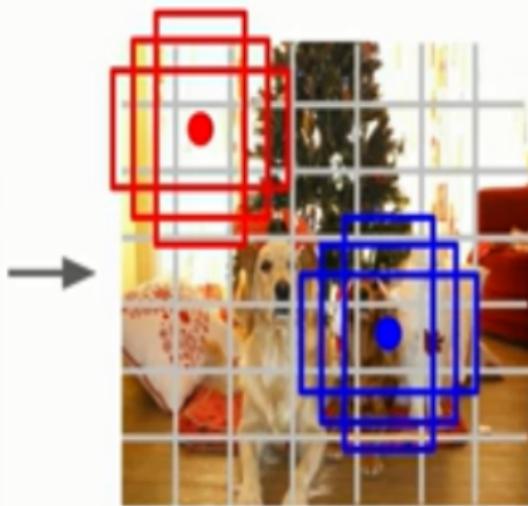
Faster R-CNN: Make CNN do proposals!



Yolo



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

References

References:

- Krizhevsky, Alex & Sutskever, Ilya & E. Hinton, Geoffrey. (2012). ImageNet Classification with Deep Convolutional Neural Networks. *Neural Information Processing Systems*. 25. 10.1145/3065386.
- Very Deep Convolutional Networks for Large-Scale Image Recognition Karen Simonyan, Andrew Zisserman (Submitted on 4 Sep 2014 (v1), last revised 10 Apr 2015 (this version, v6))
- Going Deeper with Convolutions Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich (Submitted on 17 Sep 2014)
- Stanford CS231n 2017

Any Question?