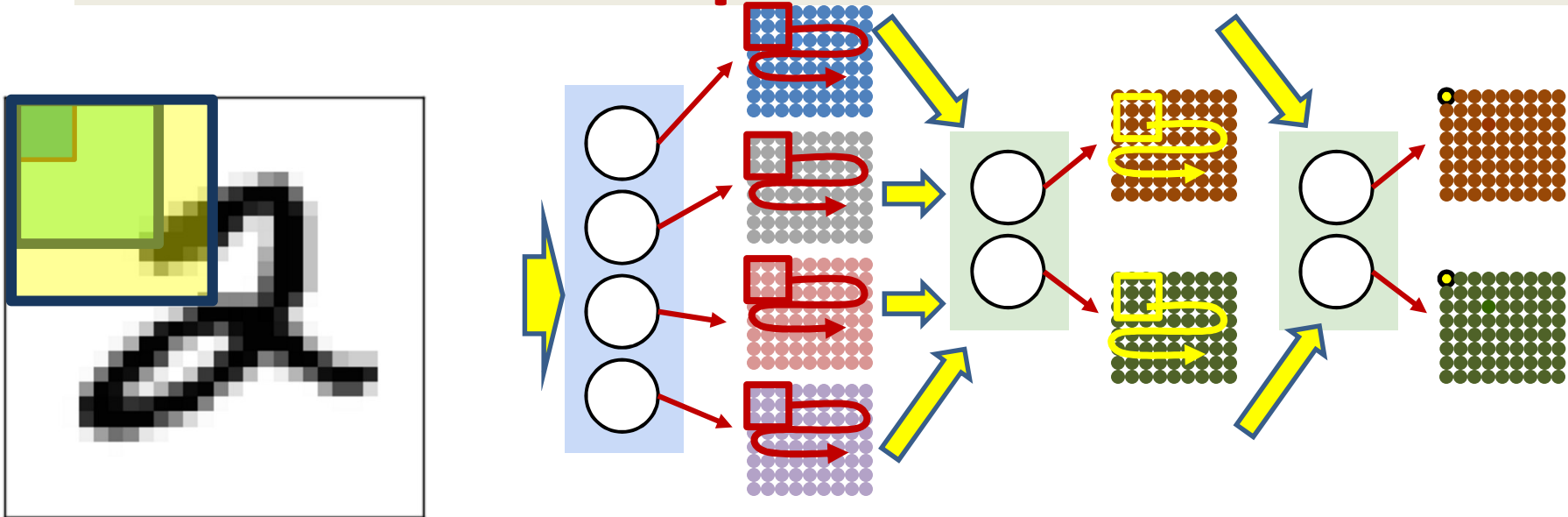# What do the filters learn? Receptive fields



- The pattern in the *input* image that each neuron sees is its "Receptive Field"
- The receptive field for a first layer neurons is simply its arrangement of weights
- For the higher level neurons, the actual receptive field is not immediately obvious and must be *calculated*
  - What patterns in the input do the neurons actually respond to?
  - We estimate it by setting the output of the neuron to 1, and learning the *input* by backpropagation
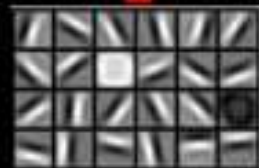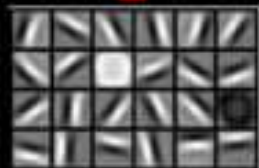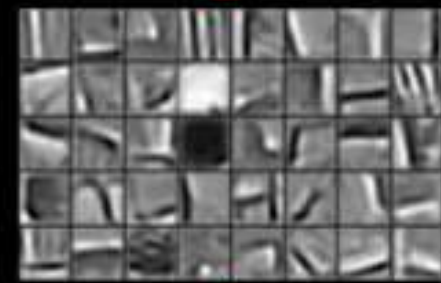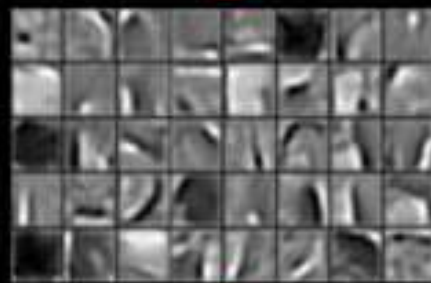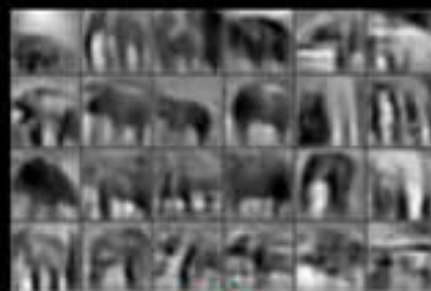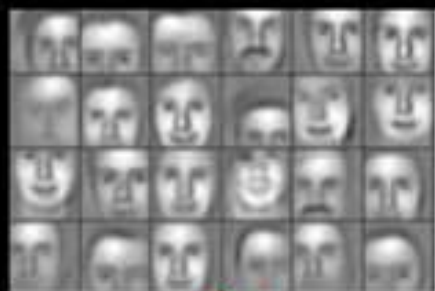
215

# Features learned from training on different object classes.



Faces    Cars    Elephants    Chairs

# Training Issues

- Standard convergence issues
  - Solution: Adam or other momentum-style algorithms
  - Other tricks such as batch normalization

- The number of parameters can quickly become very large
- Insufficient training data to train well
  - Solution: Data augmentation

# Data Augmentation

Original data

Augmented data



- rotation: uniformly chosen random angle between 0° and 360°
- translation: random translation between -10 and 10 pixels
- rescaling: random scaling with scale factor between 1/1.6 and 1.6 (log-uniform)
- flipping: yes or no (bernoulli)
- shearing: random shearing with angle between -20° and 20°
- stretching: random stretching with stretch factor between 1/1.3 and 1.3 (log-uniform)

# Digit classification

# Le-net 5



- Digit recognition on MNIST (32x32 images)
  - **Conv1:** 6 5x5 filters in first conv layer (no zero pad), stride 1
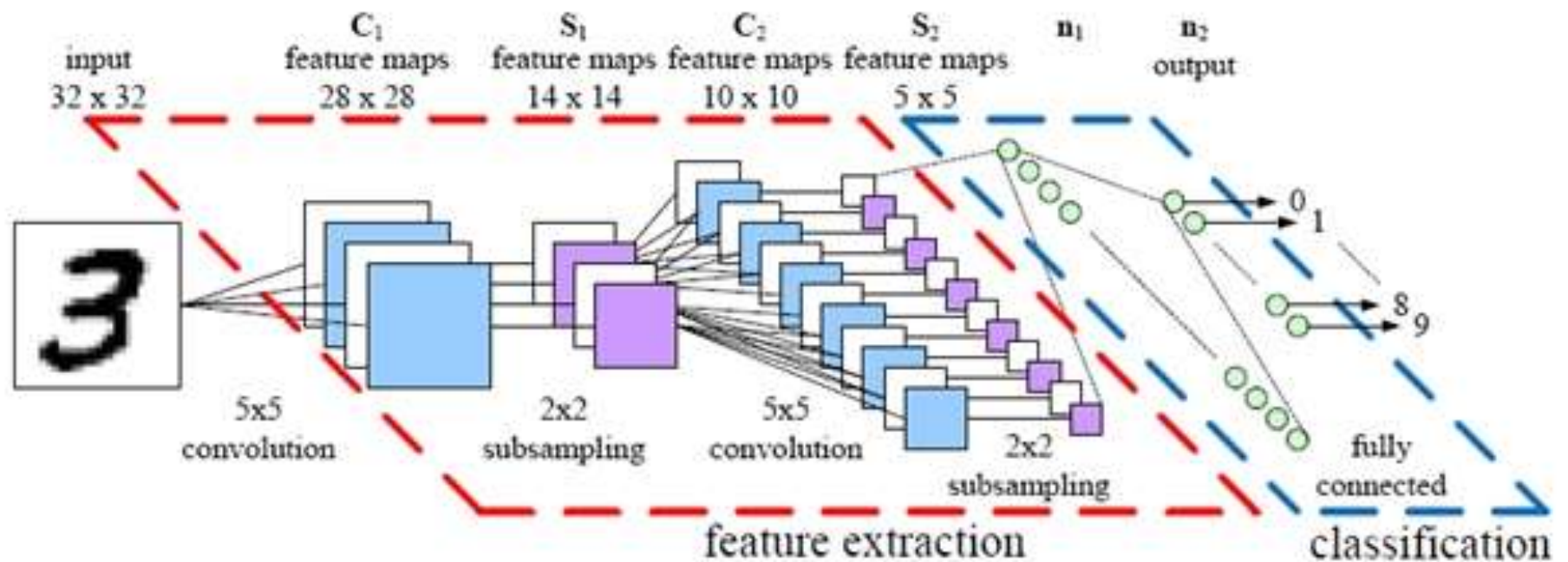    - Result: 6 28x28 maps
  - **Pool1:** 2x2 max pooling, stride 2
    - Result: 6 14x14 maps
  - **Conv2:** 16 5x5 filters in second conv layer, stride 1, no zero pad
    - Result: 16 10x10 maps
  - **Pool2:** 2x2 max pooling with stride 2 for second conv layer
    - Result 16 5x5 maps (400 values in all)
  - **FC:** Final MLP: 3 layers
    - 120 neurons, 84 neurons, and finally 10 output neurons

# Nice visual example

- http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html

# The imagenet task



- **Imagenet Large Scale Visual Recognition Challenge (ILSVRC)**

- http://www.image-net.org/challenges/LSVRC/

- Actual dataset:  Many million images, thousands of categories

- For the evaluations that follow:

  - 1.2 million pictures
  - 1000 categories

# AlexNet

- 1.2 million high-resolution images from ImageNet LSVRC-2010 contest
- 1000 different classes (softmax layer)
- NN configuration
    - NN contains 60 million parameters and 650,000 neurons,
    - 5 convolutional layers, some of which are followed by max-pooling layers
    - 3 fully-connected layers



Krizhevsky, A., Sutskever, I. and Hinton, G. E. "ImageNet Classification with Deep Convolutional Neural Networks" NIPS 2012: Neural Information Processing Systems, Lake Tahoe, Nevada

# Krizhevsky et. al.

- Input: 227x227x3 images
- Conv1:  96 11x11 filters, stride 4, no zeropad
- Pool1: 3x3 filters, stride 2
- "Normalization" layer  [Unnecessary]
- Conv2: 256 5x5 filters, stride 2, zero pad
- Pool2: 3x3,  stride 2
- Normalization layer  [Unnecessary]
- Conv3: 384 3x3,  stride 1, zeropad
- Conv4: 384 3x3, stride 1, zeropad
- Conv5: 256 3x3, stride 1, zeropad
- Pool3: 3x3, stride 2
- FC:  3 layers,
  - 4096 neurons, 4096 neurons, 1000 output neurons

# Alexnet: Total parameters

- 650K neurons

- 60M parameters

- 630M connections



10 patches

- Testing: Multi-crop
  - Classify different shifts of the image and vote over the lot!

# Learning magic in Alexnet

- **Activations were RELU**

  – Made a large difference in convergence

- "Dropout" – 0.5 (in FC layers only)

- *Large amount of data augmentation*

- SGD with mini batch size 128

- Momentum, with momentum factor 0.9

- L2 weight decay 5e-4

- Learning rate: 0.01,  decreased by 10 every time validation accuracy plateaus

- Evaluated using: Validation accuracy


- **Final top-5 error: 18.2% with a single net, 15.4% using an ensemble of 7 networks**

  – **Lowest prior error using conventional classifiers:  > 25%**

# ImageNet

Figure 3: 96 convolutional kernels of size 11×11×3 learned by the first convolutional layer on the 224×224×3 input images. The top 48 kernels were learned on GPU 1 while the bottom 48 kernels were learned on GPU 2. See Section 6.1 for details.
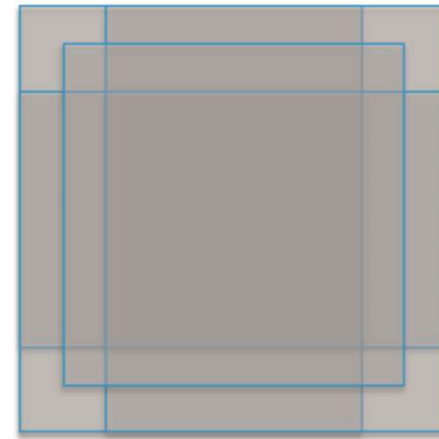


Krizhevsky, A., Sutskever, I. and Hinton, G. E. "ImageNet Classification with Deep Convolutional Neural Networks" NIPS 2012: Neural Information Processing Systems, Lake Tahoe, Nevada
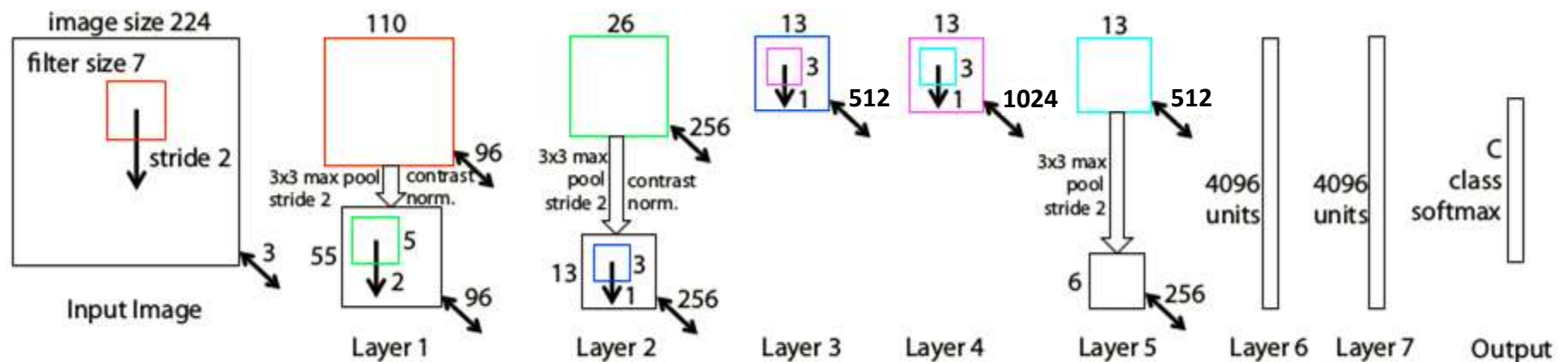
# The net actually *learns* features!



Eight ILSVRC-2010 test images and the five labels considered most probable by our model. The correct label is written under each image, and the probability assigned to the correct label is also shown with a red bar (if it happens to be in the top 5).

Five ILSVRC-2010 test images in the first column. The remaining columns show the six training images that produce feature vectors in the last hidden layer with the smallest Euclidean distance from the feature vector for the test image.

Krizhevsky, A., Sutskever, I. and Hinton, G. E. "ImageNet Classification with Deep Convolutional Neural Networks" NIPS 2012: Neural Information Processing Systems, Lake Tahoe, Nevada
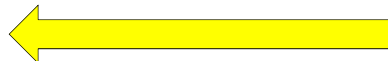
# ZFNet



ZF Net Architecture

- Zeiler and Fergus 2013
- Same as Alexnet except:
  - 7x7 input-layer filters with stride 2
  - 3 conv layers are 512, 1024, 512
  - Error went down from 15.4% → 14.8%
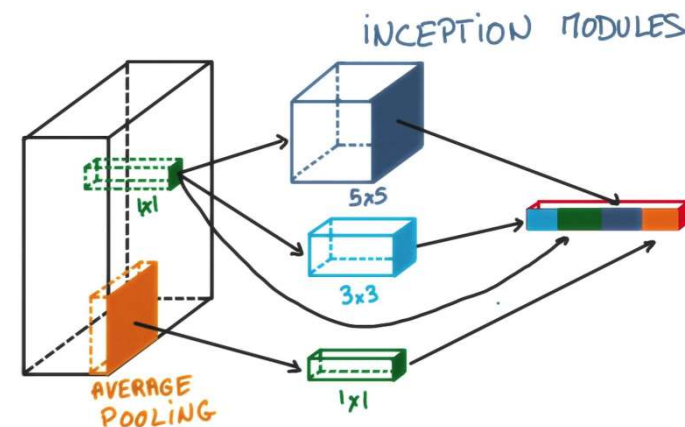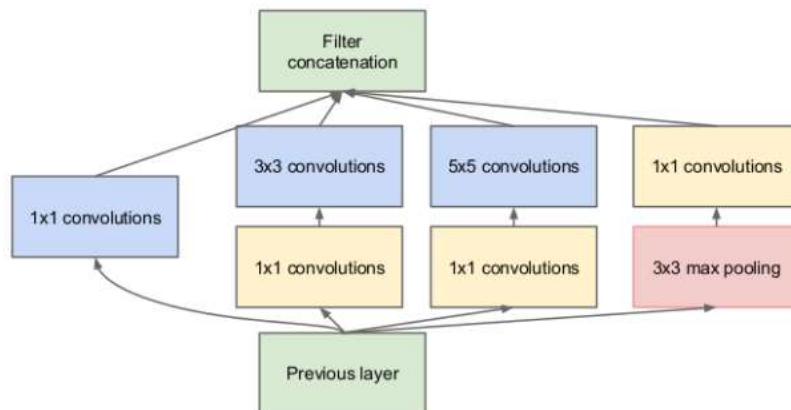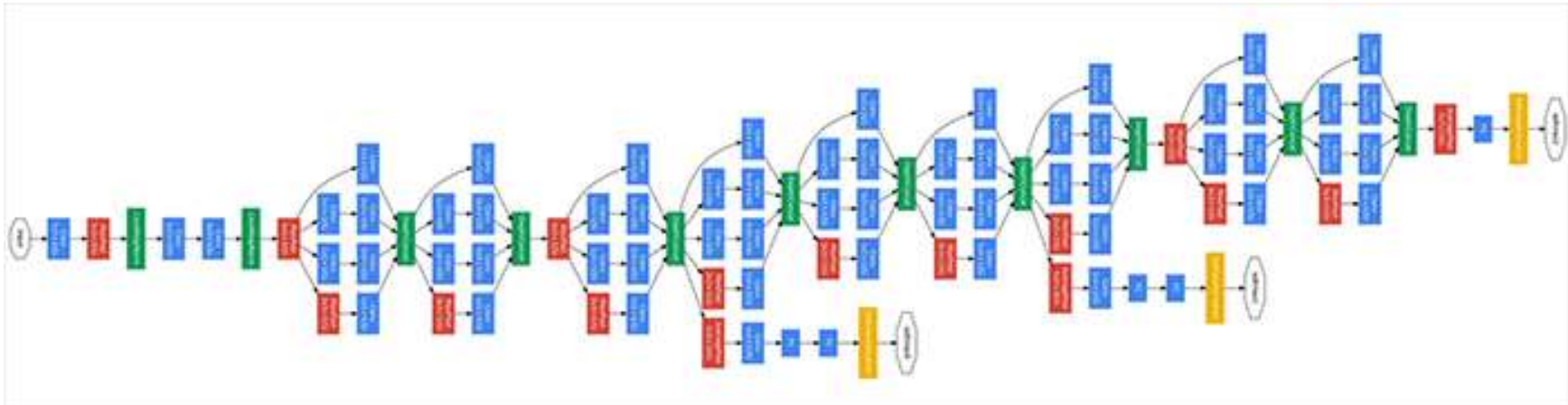    - Combining multiple models as before

231

# VGGNet

- Simonyan and Zisserman, 2014
- *Only* used 3x3 filters, stride 1, pad 1
- *Only* used 2x2 pooling filters, stride 2

- Tried a large number of architectures.
- Finally obtained 7.3% top-5 error using 13 conv layers and 3 FC layers
  - Combining 7 classifiers
  - Subsequent to paper, reduced error to 6.8% using only two classifiers
- Final arch:  64 conv, 64 conv,
  64 pool,
  128 conv, 128 conv,
  128 pool,
  256 conv, 256 conv, 256 conv,
  256 pool,
  512 conv, 512 conv, 512  conv,
  512 pool,
  512 conv, 512 conv, 512  conv,
  512 pool,
  FC with 4096, 4096, 1000
- ~140 million parameters in all!

Madness!

| 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 | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
|  | **LRN** | **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 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
|  |  |  | **conv1-256** | **conv3-256** | conv3-256 |
|  |  |  |  |  | **conv3-256** |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
|  |  |  | **conv1-512** | **conv3-512** | conv3-512 |
|  |  |  |  |  | **conv3-512** |
| maxpool | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
|  |  |  | **conv1-512** | **conv3-512** | conv3-512 |
|  |  |  |  |  | **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

# Googlenet: Inception



- Multiple filter sizes simultaneously
- Details irrelevant;  error → 6.7%
  - Using only 5 million parameters, thanks to average pooling
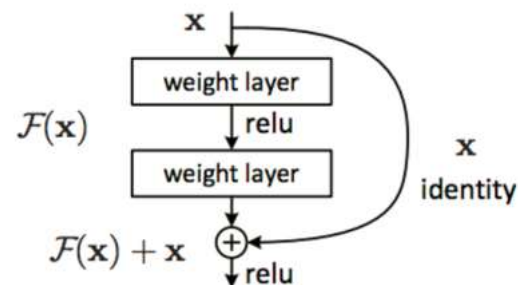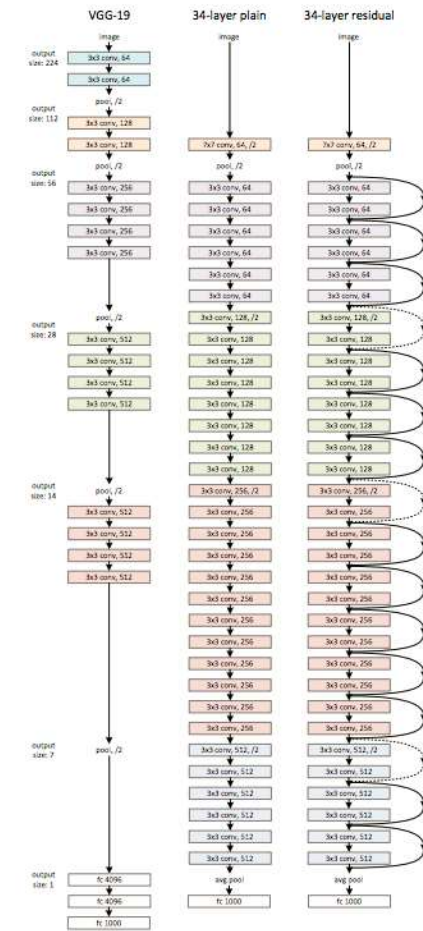
# Imagenet



Figure 2. Residual learning: a building block.

- Resnet: 2015
  - Current top-5 error:  < 3.5%
  - Over 150 layers, with "skip" connections..

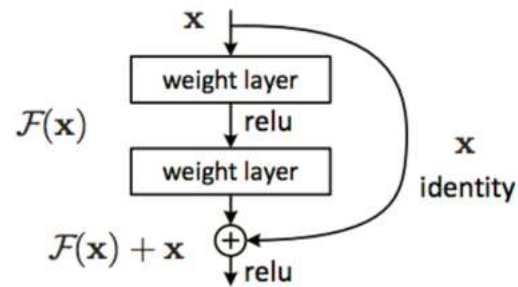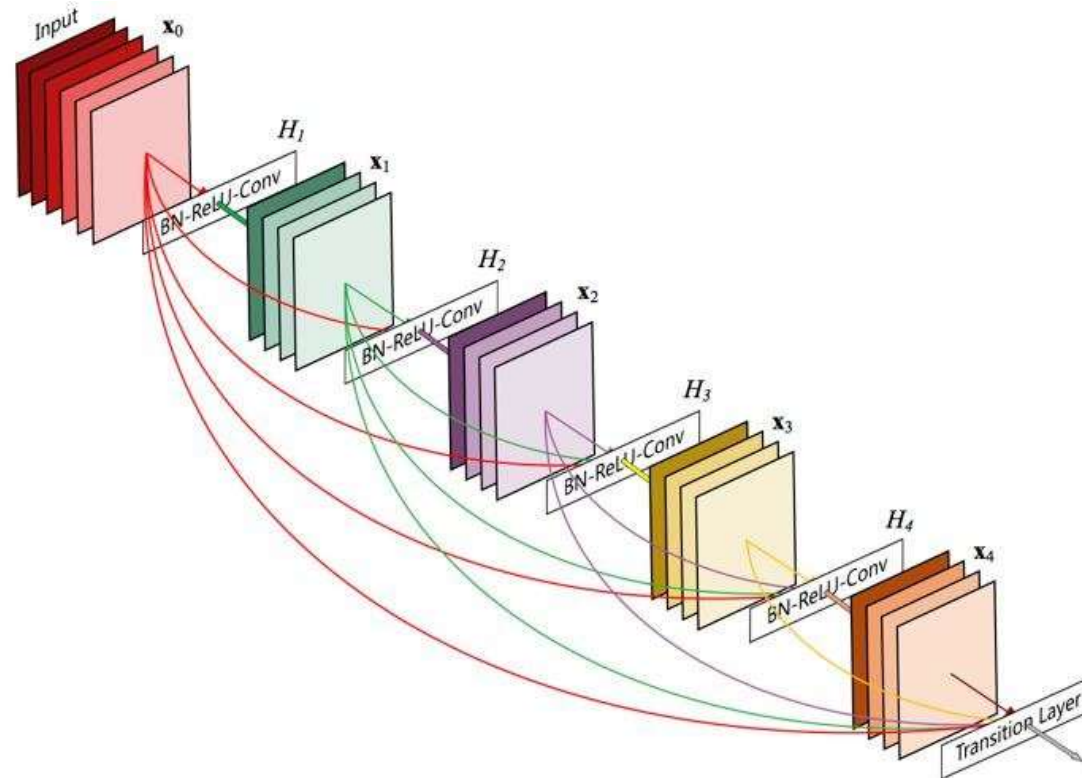# Resnet details for the curious..



Figure 2. Residual learning: a building block.

- Last layer before addition must have the same number of filters as the input to the module

- Batch normalization after each convolution

- SGD + momentum (0.9)

- Learning rate 0.1, divide by 10 (batch norm lets you use larger learning rate)

- Mini batch 256

- Weight decay 1e-5

# Densenet



- All convolutional
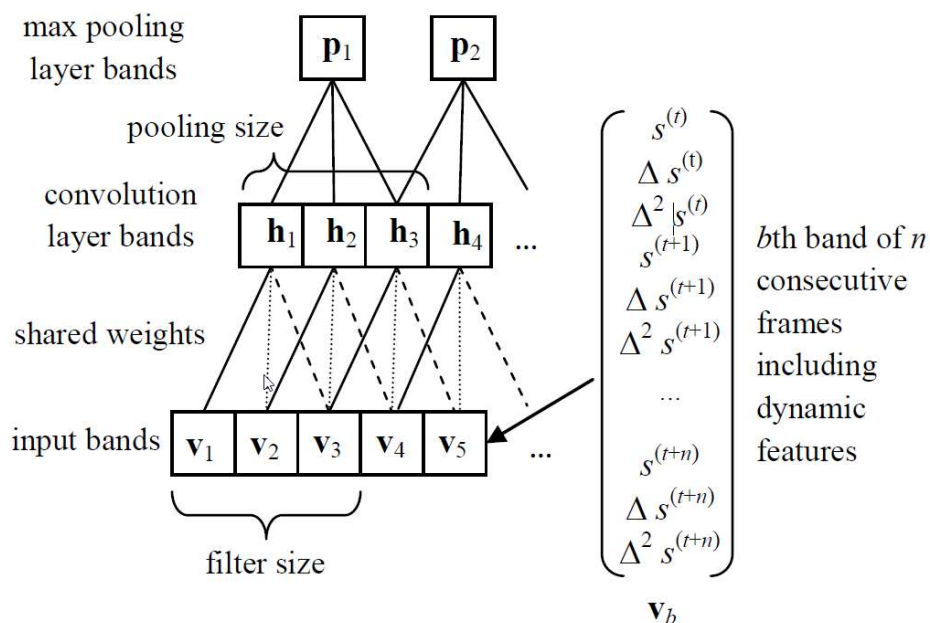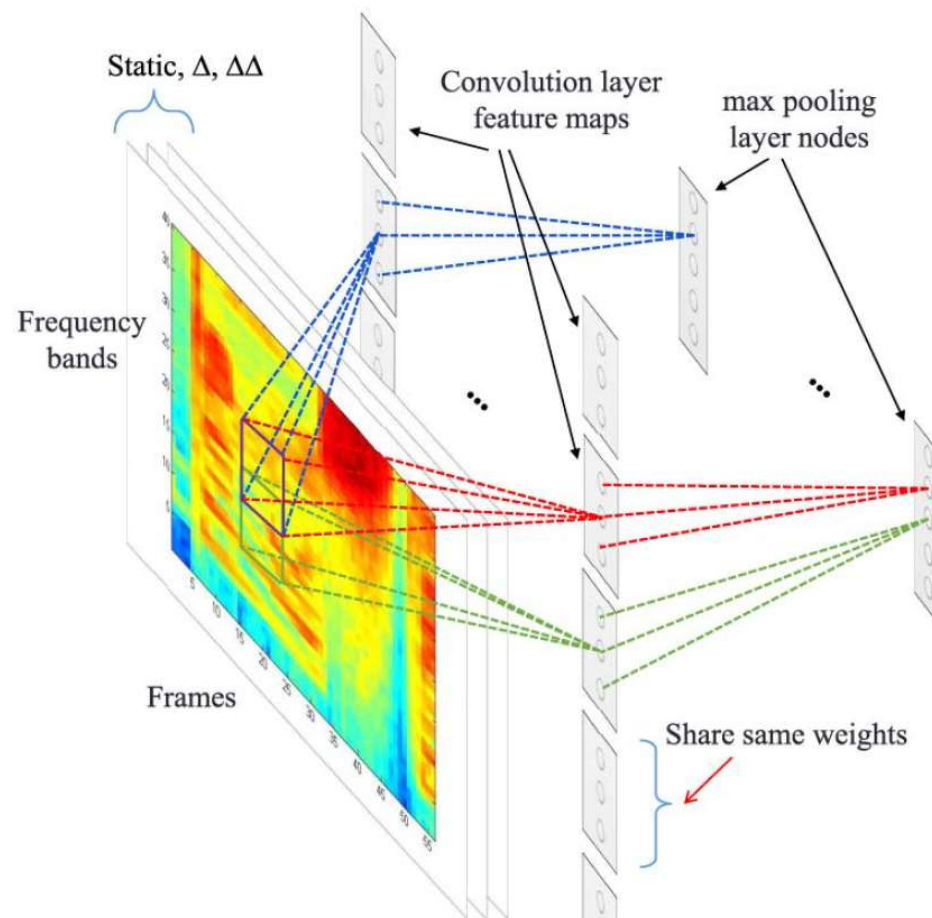- Each layer looks at the <u>union</u> of maps from all previous layers
  - Instead of just the set of maps from the immediately previous layer
- Was state of the art before I went for coffee one day
  - Wasn't when I got back..

# Many many more architectures

- Daily updates on arxiv..

- Many more applications
    - CNNs for speech recognition
    - CNNs for language processing!
    - More on these later..

# CNN for Automatic Speech Recognition

- Convolution over frequencies
- Convolution over time





| Deep Networks | Phone Error Rate |
|---|---|
| DNN (fully connected) | 22.3% |
| CNN-DNN; P=1 | 21.8% |
| CNN-DNN; P=12 | 20.8% |
| CNN-DNN; P=6 (fixed P, optimal) | 20.4% |
| CNN-DNN; P=6 (add dropout) | 19.9% |
| **CNN-DNN; P=1:m (HP, m=12)** | **19.3%** |
| **CNN-DNN; above (add dropout)** | **18.7%** |

*Table 1: TIMIT core test set phone recognition error rate comparisons.*

# CNN-Recap

- Neural network with specialized connectivity structure
- Feed-forward:
    - Convolve input
    - Non-linearity (rectified linear)
    - Pooling (local max)
- Supervised training
- Train convolutional filters by back-propagating error
- Convolution over time



Feature maps

↑

Pooling

↑

Non-linearity

↑

Convolution (Learned)

↑

Input image



x(t)  x(t-1)  x(t-2)  x(t-3)

x(t)



INPUT 32x32

C1: feature maps 6@28x28

S2: f. maps 6@14x14

C3: f. maps 16@10x10

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions    Subsampling    Convolutions    Subsampling    Full connection    Gaussian connections

Full connection