detail from the visualization of ResNet-18 // Graphcore

# CMP784

## DEEP LEARNING

Lecture #05 – Convolutional Neural Networks

Erkut Erdem // Hacettepe University // Fall 2021

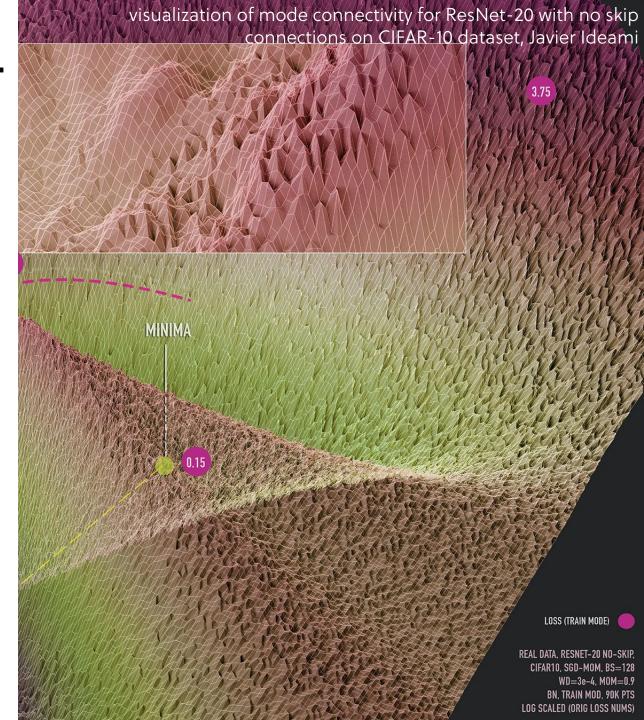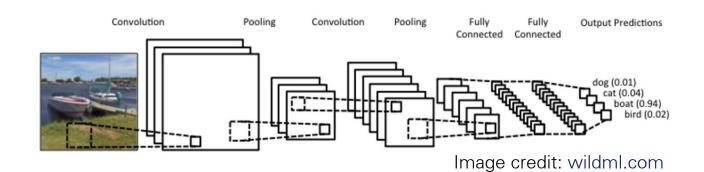# Previously on CMP784

- data preprocessing and normalization

- weight initializations

- ways to improve generalization

- babysitting the learning process

- hyperparameter selection

- optimization

visualization of mode connectivity for ResNet-20 with no skip connections on CIFAR-10 dataset, Javier Ideami

3.75

MINIMA

0.15

LOSS (TRAIN MODE)

REAL DATA, RESNET-20 NO-SKIP,
CIFAR10, SGD-MOM, BS=128
WD=3e-4, MOM=0.9
BN, TRAIN MOD, 90K PTS
LOG SCALED (ORIG LOSS NUMS)

# Breaking news!

- Practical 2 is out!
  - Convolutional Neural Networks
  - Due Wednesday, Nov. 17, 23:59:59



Image credit: wildml.com

- Project proposals is due Nov. 3!
  - about a half page
  - the research topic to be investigated,
  - what data you will use,
  - design overview,
  - a list of key readings.

  **Note:** The project should be done <u>in pairs</u>.

Deadlines in the syllabus are closer than they appear

# Lecture Overview

- convolution layer

- pooling layer

- revolution of depth

- design guidelines

- residual connections

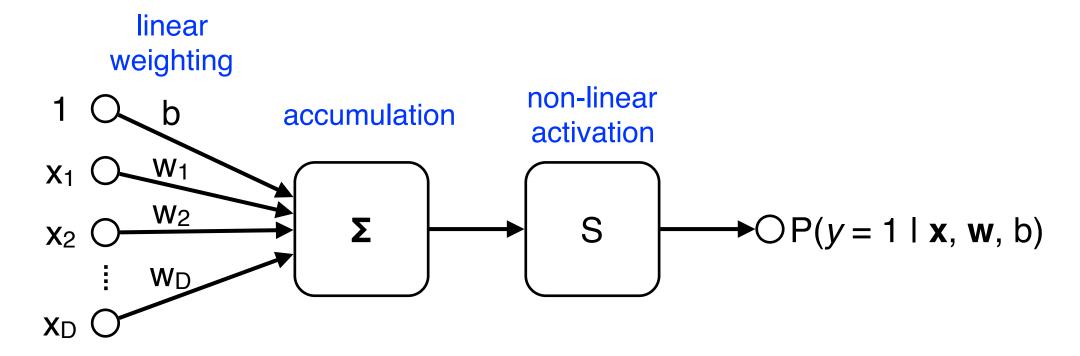- semantic segmentation networks

- object detection networks

**Disclaimer:** Much of the material and slides for this lecture were borrowed from

— Andrea Vedaldi's tutorial on Convolutional Networks for Computer Vision Applications

— Kaiming He's ICML 2016 tutorial on Deep Residual Networks: Deep Learning Gets Way Deeper

— Ross Girshick's talk on The Past, Present, and Future of Object Detection

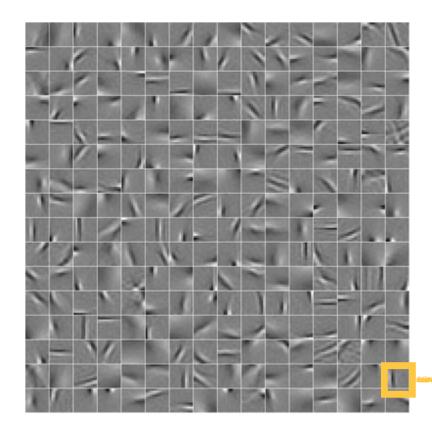— Fei-Fei Li, Andrej Karpathy and Justin Johnson's CS231n class
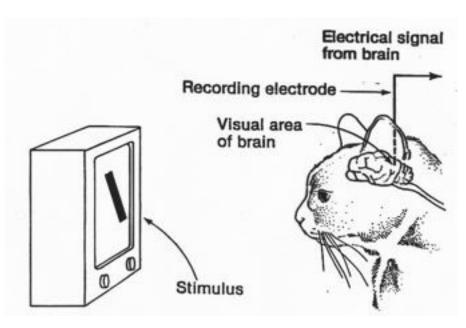
# Perceptron

[Rosenblatt 57]

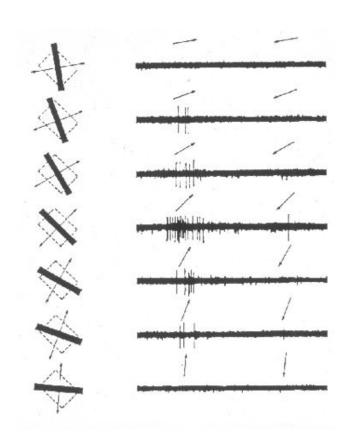- The goal is estimating the posterior probability of the binary label y of a vector **x**:

# Discovery of oriented cells in the visual cortex

[Hubel and Wiesel 59]



oriented filter

# Convolution



- Convolution = Spatial filtering

$$(a \star b)[i,j] = \sum_{i',j'} a[i',j']b[i-i',j-j']$$

vec **y** = [ ] × vec **x**

Banded matrix equivalent to  *F*

**Convolution transpose**

**Transposed**

- Different filters (weights) reveal a different characteristics of the input.



| $*^{1/8}$ | $F$ | | |
|---|---|---|---|
| | 0 | 1 | 0 |
| | 1 | 4 | 1 |
| | 0 | 1 | 0 |

vec **y** = × vec **x**

Transposed matrix

# Convolution



- Convolution = Spatial filtering

$$(a \star b)[i, j] = \sum_{i', j'} a[i', j'] b[i - i', j - j']$$

**vec y** ... **=** [ ... ] **×** **vec x**

Banded matrix equivalent to $F$

**Convolution transpose**

**Transposed**

- Different filters (weights) reveal a different characteristics of the input.



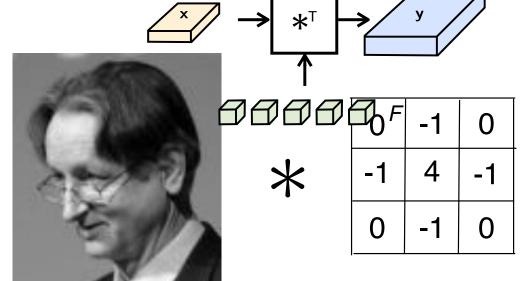| 0 | -1 | 0 |
|---|----|---|
| -1 | 4 | -1 |
| 0 | -1 | 0 |

**vec y** ... **=** ... **×** **vec x**

Transposed matrix

# Convolution



- Convolution = Spatial filtering

$$(a \star b)[i,j] = \sum_{i',j'} a[i',j']b[i-i',j-j']$$

vec **y** = [ banded matrix ] × vec **x**

Banded matrix equivalent to *F*

**Convolution transpose**

- Different filters (weights) reveal a different characteristics of the input.



| 1 | 0 | -1 |
|---|---|---|
| 2 | 0 | -2 |
| 1 | 0 | -1 |

**Transposed**

vec **y** = × vec **x**

Transposed matrix

12

# Convolutional Neural Networks in a Nutshell

- A neural network model that consists of a sequence of local & translation invariant layers
    - Many identical copies of the same neuron: Weight/parameter sharing
    - Hierarchical feature learning



A. Krizhevsky, I. Sutskever, and G. E. Hinton. **Imagenet classification with deep convolutional neural networks**. In NIPS 2012.

# A bit of history

- Neocognitron model by Fukushima (1980)

- The first convolutional neural network (CNN) model

- so-called "sandwich" architecture
  - simple cells act like filters
  - complex cells perform pooling

- Difficult to train
  - No backpropagation yet

# A bit of history

- LeNet-5 model



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

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.

# A bit of history

- AlexNet model



A. Krizhevsky, I. Sutskever, and G. E. Hinton. **Imagenet classification with deep convolutional neural networks**. In NIPS 2012.

# Convolutional Neural Network



A. Krizhevsky, I. Sutskever, and G. E. Hinton. **Imagenet classification with deep convolutional neural networks**. In NIPS 2012.

# Convolutional layer

- Learn a filter bank (a set of filters) once
- Use them over the input data to extract features

$$\mathbf{y} = F * \mathbf{x} + b$$

input data $\mathbf{x}$      filter bank $F$      output data $\mathbf{y}$

# Data = 3D Tensors

- There is a vector of feature channels (e.g. RGB) at each spatial location (pixel).

# Convolutions with 3D Filters

- Each filter acts on multiple input channels

  – **Local**
    Filters look locally

  – **Translation invariant**
    Filters act the same everywhere

# Convolutional Layer

32x32x3 input

32

32

3

5x5x3 filter

Convolve the filter with the input i.e. "slide over the image spatially, computing dot products"

# Convolutional Layer

<span style="color:red">32x32x3 input</span>
<span style="color:blue">5x5x3 filter</span>

32

32

3

1 number:
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the input
(i.e. 5*5*3 = 75-dimensional dot product + bias)

# Convolutional Layer



32x32x3 input
5x5x3 filter

activation map

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Convolutional Layer

consider a second, green filter

32x32x3 input
5x5x3 filter

activation maps

convolve (slide) over all
spatial locations

32

32

3

28

28

1

24

# Convolutional Layer

- Multiple filters produce multiple output channels
- For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

activation maps

32

32

3

Convolutional Layer

28

28

6

We stack these up to get an output of size 28x28x6.

# Spatial Arrangement of Output Volume

32

32

3

28

28

5

- **Depth:** number of filters
- **Stride:** filter step size (when we "slide" it)
- **Padding:** zero-pad the input

Input Volume (+pad 1) (7x7x3)  Filter W0 (3x3x3)  Filter W1 (3x3x3)  Output Volume (3x3x2)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

w0[:,:,0]

| 0 | 0 | -1 |
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
| 0 | 1 | 1 |
| 0 | -1 | 0 |

w1[:,:,0]

| -1 | 0 | 0 |
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
| -1 | 1 | -1 |
| 0 | 1 | -1 |

o[:,:,0]

| -3 | -1 | 4 |
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
| -7 | -11 | -1 |
| -4 | -2 | -4 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |

**Input Volume (+pad 1) (7x7x3)**　　　**Filter W0 (3x3x3)**　　　**Filter W1 (3x3x3)**　　　**Output Volume (3x3x2)**

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|----|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|----|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|----|----|----|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

w1[:,:,0]

| -1 | 0 | 0 |
|----|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|----|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|----|----|----|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

o[:,:,0]

| -3 | -1 | 4 |
|----|----|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|----|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

**Bias b0 (1x1x1)**　　　**Bias b1 (1x1x1)**

b0[:,:,0]

| 1 |
|---|

b1[:,:,0]

| 0 |
|---|

## Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

## Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|----|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|----|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | 1 | -1 |
|----|---|----|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

## Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|----|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|----|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|----|----|----|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

## Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|----|----|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|----|----|----|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

## Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

## Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|---|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|---|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | 1 |
|---|---|---|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|---|---|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|---|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
| 0 | 1 | 1 |
| 0 | -1 | 0 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | 1 |
| -1 | 1 | -1 |
| 0 | 1 | -1 |

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |

Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
| -7 | -11 | -1 |
| -4 | -2 | -4 |

31

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|---|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|---|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|---|---|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|---|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

32

## Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

## Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|---|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

## Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|---|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

## Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|---|---|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|---|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

## Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

## Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|---|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

Bias b0 (1x1x1)

b0[:,:,0]

1

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|---|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

Bias b1 (1x1x1)

b1[:,:,0]

0

Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|---|---|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|---|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 2 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 2 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 2 | 2 | 0 |
| 0 | 1 | 2 | 2 | 2 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter W0 (3x3x3)

w0[:,:,0]

| 0 | 0 | -1 |
|---|---|---|
| -1 | 0 | 0 |
| -1 | -1 | -1 |

w0[:,:,1]

| 1 | 1 | -1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | 1 | 1 |

w0[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | -1 | 0 |

Filter W1 (3x3x3)

w1[:,:,0]

| -1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 0 | 0 | -1 |

w1[:,:,1]

| -1 | 0 | 1 |
|---|---|---|
| 1 | -1 | 1 |
| -1 | 0 | 1 |

w1[:,:,2]

| -1 | -1 | -1 |
|---|---|---|
| -1 | 1 | -1 |
| 0 | 1 | -1 |

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Output Volume (3x3x2)

o[:,:,0]

| -3 | -1 | 4 |
|---|---|---|
| -2 | -7 | -4 |
| 1 | -1 | 1 |

o[:,:,1]

| -7 | 3 | 1 |
|---|---|---|
| -7 | -11 | -1 |
| -4 | -2 | -4 |

# Convolutional layers

- Local receptive field

- Each column of hidden units looks at a different input patch

feature component

features

input image

receptive field

# Effective Receptive Field

Contributing input units to a convolutional filter.

@jimmfleming // fomoro.com

**Input Features**

**7 // 2 Convolution**
Each filter sees 7 input units

strides continue…

**Convolutional Features**

**2 // 2 Max Pool**
Each filter sees 9 input units

**Max Pool Features**

**3 // 1 Convolution**
Each filter sees 17 input units

**Convolutional Features**

Features

Conv1D Filter

Padding or Stride

Receptive Field

[http://fomoro.com/tools/receptive-fields/index.html]

37

# Convolutional layers



32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

28

28

6

# Repeat linear / non-linear operators



32
32
3

CONV,
ReLU
e.g. 6
5x5x3
filters

28
28
6

CONV,
ReLU
e.g. 10
5x5x**6**
filters

24
24
10

CONV,
ReLU

....

# Linear/Non-linear Chains

- The basic blueprint of most architectures
- Stack multiple layers of convolutions



**filtering**      **ReLU**      **filtering
& downsampling**      **ReLU**      …

# Feature Learning

- Hierarchical layer structure allows to learn hierarchical filters (features).

# Feature Learning

• Hierarchical layer structure allows to learn hierarchical filters (features).



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

# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:
- Max pooling, average pooling, etc.

Single depth slice



x

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

max pool with 2x2 filters and stride 2 →

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

y



224x224x64

pool →

112x112x64

224

224

downsampling

112

112

43

# Fully connected layer

- contains neurons that connect to the entire input volume, as in ordinary Neural Networks

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

Details/Retrospectives:
- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

# Convolutional Neural Network Demo

- ConvNetJS demo: training on CIFAR-10

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

# Three Years of Progress
## From AlexNet (2012) to ResNet (2015)

# Revolution of Depth

AlexNet, 8 layers

(ILSVRC 2012)

11x11 conv, 96, /4, pool/2

↓

5x5 conv, 256, pool/2

↓

3x3 conv, 384

↓

3x3 conv, 384

↓

3x3 conv, 256, pool/2

↓

fc, 4096

↓

fc, 4096

↓

fc, 1000



- 5 convolutional layers
- 3 fully connected layers
- ReLU
- End-to-end (no pre-training)
- Data augmentation

# Revolution of Depth

AlexNet, 8 layers

(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers

(ILSVRC 2014)

| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

- Very deep
- Simply deep

INPUT: [224x224x3]        memory: 224*224*3=150K   params: 0
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M   params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M   params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]  memory: 112*112*64=800K   params: 0
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M   params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M   params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]  memory: 56*56*128=400K   params: 0
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]  memory: 28*28*256=200K   params: 0
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]  memory: 14*14*512=100K   params: 0
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]  memory: 7*7*512=25K  params: 0
FC: [1x1x4096]  memory: 4096  params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]  memory: 4096  params: 4096*4096 = 16,777,216
FC: [1x1x1000]  memory: 1000 params: 4096*1000 = 4,096,000

TOTAL memory: 24M * 4 bytes ~= 93MB / image
(only forward! ~*2 for bwd)
TOTAL params: 138M parameters

| ConvNet Configuration | | | |
| --- | --- | --- | --- |
| B | C | D | 19 |
| 13 weight layers | 16 weight layers | 16 weight layers | |
| input (224 × 224 RGB image) | | | |
| conv3-64 | conv3-64 | conv3-64 | co |
| conv3-64 | conv3-64 | conv3-64 | co |
| maxpool | | | |
| conv3-128 | conv3-128 | conv3-128 | co |
| conv3-128 | conv3-128 | conv3-128 | co |
| maxpool | | | |
| conv3-256 | conv3-256 | conv3-256 | co |
| conv3-256 | conv3-256 | conv3-256 | co |
| | conv1-256 | conv3-256 | co |
| | | | co |
| maxpool | | | |
| conv3-512 | conv3-512 | conv3-512 | co |
| conv3-512 | conv3-512 | conv3-512 | co |
| | conv1-512 | conv3-512 | co |
| | | | co |
| maxpool | | | |
| conv3-512 | conv3-512 | conv3-512 | co |
| conv3-512 | conv3-512 | conv3-512 | co |
| | conv1-512 | conv3-512 | co |
| | | | co |
| maxpool | | | |
| FC-4096 | | | |
| FC-4096 | | | |
| FC-1000 | | | |
| soft-max | | | |

# VGG-16 Net

INPUT: [224x224x3]        memory: 224*224*3=150K   params: 0          (not counting biases)
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M   params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M   params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]  memory: 112*112*64=800K   params: 0
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M   params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128]  memory: 112*112*128=1.6M   params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]  memory: 56*56*128=400K   params: 0
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256]  memory: 56*56*256=800K   params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]  memory: 28*28*256=200K   params: 0
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512]  memory: 28*28*512=400K   params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]  memory: 14*14*512=100K   params: 0
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]  memory: 14*14*512=100K   params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]  memory: 7*7*512=25K  params: 0
FC: [1x1x4096]  memory: 4096  params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]  memory: 4096  params: 4096*4096 = 16,777,216
FC: [1x1x1000]  memory: 1000 params: 4096*1000 = 4,096,000

Note:

Most memory is in early CONV

Most params are in late FC

TOTAL memory: 24M * 4 bytes ~= 93MB / image
(only forward! ~*2 for bwd)
TOTAL params: 138M parameters

VGG-16 Net

# Revolution of Depth

AlexNet, 8 layers

(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers

(ILSVRC 2014)

| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

- Very deep
- Simply deep

GoogLeNet,
22 layers

(ILSVRC 2014)

- Branching
- Bootleneck
- Skip connection

# GoogLeNet

Inception module

# GoogLeNet

| type | patch size/ stride | output size | depth | #1×1 | #3×3 reduce | #3×3 | #5×5 reduce | #5×5 | pool proj | params | ops |
|---|---|---|---|---|---|---|---|---|---|---|---|
| convolution | 7×7/2 | 112×112×64 | 1 | | | | | | | 2.7K | 34M |
| max pool | 3×3/2 | 56×56×64 | 0 | | | | | | | | |
| convolution | 3×3/1 | 56×56×192 | 2 | | 64 | 192 | | | | 112K | 360M |
| max pool | 3×3/2 | 28×28×192 | 0 | | | | | | | | |
| inception (3a) | | 28×28×256 | 2 | 64 | 96 | 128 | 16 | 32 | 32 | 159K | 128M |
| inception (3b) | | 28×28×480 | 2 | 128 | 128 | 192 | 32 | 96 | 64 | 380K | 304M |
| max pool | 3×3/2 | 14×14×480 | 0 | | | | | | | | |
| inception (4a) | | 14×14×512 | 2 | 192 | 96 | 208 | 16 | 48 | 64 | 364K | 73M |
| inception (4b) | | 14×14×512 | 2 | 160 | 112 | 224 | 24 | 64 | 64 | 437K | 88M |
| inception (4c) | | 14×14×512 | 2 | 128 | 128 | 256 | 24 | 64 | 64 | 463K | 100M |
| inception (4d) | | 14×14×528 | 2 | 112 | 144 | 288 | 32 | 64 | 64 | 580K | 119M |
| inception (4e) | | 14×14×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 840K | 170M |
| max pool | 3×3/2 | 7×7×832 | 0 | | | | | | | | |
| inception (5a) | | 7×7×832 | 2 | 256 | 160 | 320 | 32 | 128 | 128 | 1072K | 54M |
| inception (5b) | | 7×7×1024 | 2 | 384 | 192 | 384 | 48 | 128 | 128 | 1388K | 71M |
| avg pool | 7×7/1 | 1×1×1024 | 0 | | | | | | | | |
| dropout (40%) | | 1×1×1024 | 0 | | | | | | | | |
| linear | | 1×1×1000 | 1 | | | | | | | 1000K | 1M |
| softmax | | 1×1×1000 | 0 | | | | | | | | |

Fun features:

- Only 5 million params!
(Removes FC layers completely)

**Compared to AlexNet:**
- 12X less params
- 2x more compute
- 6.67% (vs. 16.4%)

# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

ResNet,
152 layers
(ILSVRC 2015)

# Revolution of Depth

AlexNet, 8 layers

(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers

(ILSVRC 2014)

| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

ResNet,
152 layers

(ILSVRC 2015)

| 7x7 conv, 64, /2, pool/2 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x1 conv, 64 |
| 3x3 conv, 64 |
| 1x1 conv, 256 |
| 1x2 conv, 128, /2 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |
| 1x1 conv, 512 |
| 1x1 conv, 128 |
| 3x3 conv, 128 |

56

# Residual Net (ResNet)

## Plain Net

Any two
stacked layers

$x$

weight layer

relu

weight layer

relu

$H(x)$

## Residual Net

$x$

weight layer

relu

$F(x)$

weight layer

identity

$x$

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

relu

# How deep is enough?

GoogLeNet (2014)

VGG-VD-16 (2014)

VGG-M (2013)

AlexNet (2012)

ResNet 50 (2015)

ResNet 152 (2015)



16 convolutional layers

50 convolutional layers

152 convolutional layers

Krizhevsky, I. Sutskever, and G. E. Hinton. *ImageNet classification with deep convolutional neural networks*. In Proc. NIPS, 2012.

C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. *Going deeper with convolutions*. In Proc. CVPR, 2015.

K. Simonyan and A. Zisserman. *Very deep convolutional networks for large-scale image recognition*. In Proc. ICLR, 2015.

K. He, X. Zhang, S. Ren, and J. Sun. *Deep residual learning for image recognition*. In Proc. CVPR, 2016.

58

# How deep is enough?

- 3 × more accurate in 3 years

# Speed

- 5 × slower



- **Remark:** 101 ResNet layers same size/speed as 16 VGG-VD layers
- **Reason:** Far fewer feature channels (quadratic speed/space gain)
- **Moral:** Optimize your architecture

# Model Size

- Num. of parameters is about the same



- **Remark:** 101 ResNet layers same size/speed as 16 VGG-VD layers
- **Reason:** Far fewer feature channels (quadratic speed/space gain)
- **Moral:** Optimize your architecture

# ResNeXt: Both Wider and Deeper



Inception:
heterogeneous multi-branch

ResNeXt:
uniform multi-branch

- shortcut, bottleneck, and multi-branch
- Better accuracy (when having the same FLOPs/#params as ResNet)

Saining Xie et al. "Aggregated Residual Transformations for Deep Neural Networks". CVPR 2017

# DenseNet

- 201 layers, 20M parameters

- Densely connected blocks
- Alleviates vanishing gradient
- Strengthens feature propagation
- Encourages feature reuse



G. Huang, Z. Liu, L. van der Maaten, & K. Weinberger, **"Densely connected convolutional networks"**, CVPR 2017

# Design Guidelines

# Design Guidelines

features



image

## Guideline 1: Avoid tight bottlenecks

- **From bottom to top**
  - The spatial resolution H×W decreases
  - The number of channels C increases

- **Guideline**
  - Avoid tight information bottleneck
  - Decrease the data volume H × W × C slowly

K. Simonyan and A. Zisserman. **Very deep convolutional networks for large-scale image recognition**. In ICLR 2015.

C. Szegedy, V. Vanhoucke, S. Ioffe, and J. Shlens. **Rethinking the inception architecture for computer vision**. In CVPR 2016.

# Receptive Field

"neuron"

neuron's
receptive field

## Must be large enough

- **Receptive field of a neuron**
  - The image region influencing a neuron
  - Anything happening outside is invisible to the neuron

- **Importance**
  - Large image structures cannot be detected by neurons with small receptive fields

- **Enlarging the receptive field**
  - Large filters
  - Chains of small filters

# Design Guidelines

**Guideline 2:** Prefer small filter chains



**One big filter bank**

**Two smaller filter banks**

prefer

$5 \times 5$ filters
+ ReLU

$3 \times 3$ filters
+ ReLU

$3 \times 3$ filters
+ ReLU

- **Remark:** 101 ResNet layers same size/speed as 16 VGG-VD layers
- **Reason:** Far fewer feature channels (quadratic speed/space gain)
- **Moral:** Optimize your architecture

# Design Guidelines

**Guideline 3:**

Keep the number of channels at bay

$H \times W \times C$

x

$*$

y

$F$

$H_f \times W_f \times C \times K$

$C$ = num. input channels

$K$ = num. output channels

**Num. of operations**

$$\frac{H \times H_f}{\text{stride}} \times \frac{W \times W_f}{\text{stride}} \times C \times K$$

**Num. of parameters**

$$H_f \times W_f \times C \times K$$

complexity $\propto C \times K$

# Design Guidelines

**Guideline 4:**

Less computations with filter groups



**M filters**

**G groups of M/G filters**

consider instead

**split channels**    **filter groups**    **put back**

Did we see this before?

$$\text{complexity} \propto (C \times K) / G$$

# AlexNet



A. Krizhevsky, I. Sutskever, and G. E. Hinton. **Imagenet classification with deep convolutional neural networks**. In NIPS 2012.

# Design Guidelines

**Guideline 4:**

Less computations with filter groups

**Full filters**

$$\mathbf{y} = [\; C \times K \;] \times \mathbf{x}$$

y      F      x

complexity: $C \times K$

**Group-sparse filters**

$$\mathbf{y} = \begin{bmatrix} & 0 & 0 \\ 0 & & 0 \\ 0 & 0 & \end{bmatrix} \times \mathbf{x}$$

y      F      x

complexity: $C \times K / G$

**Groups** = filters, seen as a matrix, have a "block" structure

# Design Guidelines

**Guideline 5:**

Low-rank decompositions



**filter bank**
$3 \times 3 \times C \times K$

decompose
spatially

decompose
channels

**vertical**
$1 \times 3 \times C \times K$

$*$

**horizontal**
$3 \times 1 \times K \times K$

$*$

**vertical**
$1 \times 3 \times K \times K$

**groups**
$3 \times 3 \times C/G \times K/G$

$*$

**"network in network"**
$1 \times 1 \times K \times K$

Make sure to mix the information

# Design Guidelines

**Guideline 6:**

Dilated Convolutions

7x7



49 coefficients
18 degrees of freedom

=

3x3

5x5

| a | 0 | b | 0 | c |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| d | 0 | e | 0 | f |
| 0 | 0 | 0 | 0 | 0 |
| g | 0 | h | 0 | i |

○

25 coefficients
9 degrees of freedom

**Exponential expansion of the receptive field without loss of resolution**

# A Closer Look to Residual Learning

# Residual Learning

$$\mathbf{x}_{n+5} = \mathbf{x}_n + (\phi_{\mathrm{ReLU}} \circ \phi_* \circ \phi_{\mathrm{ReLU}} \circ \phi_*)(\mathbf{x}_n)$$

**Fixed identity // learned residual**

$$\mathbf{x}_{n+5} = \mathbf{x}_n + (\phi_{\mathrm{ReLU}} \circ \phi_* \circ \phi_{\mathrm{ReLU}} \circ \phi_*)(\mathbf{x}_n)$$

identity          residual

K. He, X. Zhang, S. Ren, and J. Sun. **Deep residual learning for image recognition**. In CVPR 2016.

# Residual Learning



CIFAR-10

ImageNet-1000

56-layer
44-layer
32-layer
20-layer

solid: test/val
dashed: train

34-layer
18-layer

- "Overly deep" plain nets have **higher training error**

- A general phenomenon, observed in many datasets

- This is optimization issue, deeper models are harder to optimize

# Residual Learning

- Richer solution space

- A deeper model should not have **higher training error**

- A solution by construction:
  - original layers: copied from a
  - learned shallower model
  - extra layers: set as identity
  - at least the same training error

"extra"
layers

**Shallower model (18 layers):**

- 7x7 conv, 64, /2
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 128, /2
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 256, /2
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 512, /2
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- fc 1000

**Deeper counterpart (34 layers):**

- 7x7 conv, 64, /2
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 64
- 3x3 conv, 128, /2
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 128
- 3x3 conv, 256, /2
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 512, /2
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- fc 1000

# Residual Learning



- The loss surface of a 56-layer net using the CIFAR-10 dataset, both without (left) and with (right) residual connections.

Hao Li et al., "Visualizing the Loss Landscape of Neural Nets". ICLR 2018

# Transfer Learning with Convolutional Neural Networks

# Beyond CNNs

- Do features extracted from the CNN generalize other tasks and datasets?
  - Donahue et al. (2013), Chatfield et al. (2014), Razavian et al. (2014), Yosinski et al. (2014), etc.

- CNN activations as deep features
- Finetuning CNNs

# CNN activations as deep features

- CNNs discover effective representations. Why not to use them?



92

# CNN activations as deep features

- CNNs discover effective representations. Why not to use them?



Layer 1 Filters (Gabor and color blobs)

# CNN activations as deep features

- CNNs discover effective representations. Why not to use them?



Layer 1 Filters (Gabor and color blobs)

Layer 2

Layer 5

Zeiler et al., 2014

# CNN activations as deep features

- CNNs discover effective representations. Why not



Layer 1 Filters (Gabor and color blobs)

Layer 2

Layer 5

Windsor tie: 0.998959

Windsor tie: 0.992462

Last Layer

Zeiler et al., 2014

Nguyen et al., 2014

# CNNs as deep features

- CNNs discover effective representations. Why not to use them?



t-SNE feature visualizations on the ILSVRC-2012

Legend:
- structure, construction
- covering
- commodity, trade good, good
- conveyance, transport
- invertebrate
- bird
- hunting dog

LLC          GIST          Conv-1 activations          Conv-6 activations

Donahue et al. **DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition**, 2014

# Transfer Learning with CNNs

- A CNN trained on a (large enough) dataset generalizes to other visual tasks



A. Joulin, L.J.P. van der Maaten, A. Jabri, and N. Vasilache **Learning visual features from Large Weakly supervised Data.**
ECCV 2016

# Transfer Learning with CNNs

- Keep layers 1-7 of our ImageNet-trained model fixed
- Train a new softmax classifier on top using the training images of the new dataset.

1. Train on Imagenet

2. Small dataset: feature extractor

Freeze these

Train this

3. Medium dataset: finetuning

more data = retrain more of the network (or all of it)

Freeze these

tip: use only ~1/10th of the original learning rate in finetuning top layer, and ~1/100th on intermediate layers

Train this

# How transferable are features in CNN networks?

- Divide ImageNet into man-made objects A (449 classes) and natural objects B (551 classes)
- The transferability of features decreases as the distance between the base task and target task increases



99

# How transferable are features in CNN networks?

- An open research problem

A. Zamir, A. Sax, W. Shen, L. Guibas, J. Malik, S. Savarese. **Taskonomy: Disentangling Task Transfer Learning**. CVPR 2018.   100

# Semantic Segmentation



Sky

Building

Window

Door

# Semantic Image Segmentation

- Label individual pixels





input = image

$c_1$ → $c_2$ → $c_3$ → $c_4$ → $c_5$ → $f_6$ → $f_7$ → $f_8$

**convolutional**　　　　**fully-connected**

output = image

# Convolutional Layers

- Local receptive field

feature component

features

input image

receptive field

# Fully Connected Layers

• Global receptive field

class predictions

fully-connected

fully-connected

fully-connected

# Convolutional vs. Fully Connected

- Comparing the receptive fields

**Downsampling filters**

Responses are spatially selective, can be used to localize things.

**Upsampling filters**

Responses are global, do not characterize well position

Which one is more useful for pixel level labelling?

# Fully-Connected Layer = Large Filter



$K$

$\mathbf{w}^{(k)}$

$W \times H \times C$

$=$

$1 \times 1 \times K$

$*$

$\mathrm{F}^{(k)}$

$W \times H \times C \times K$

# Fully-Convolutional Neural Networks

class predictions

# Fully-Convolutional Neural Networks



- **Dense evaluation**
  - Apply the whole network convolutional
  - Estimates a vector of class probabilities at each pixel

- **Downsampling**
  - In practice most network downsample the data fast
  - The output is very low resolution (e.g. 1/32 of original)

# Upsampling The Resolution

- Interpolating filter

**Downsampling filters**

**Upsampling filters**



Upsampling filters allow to increase the resolution of the output

Very useful to get full-resolution segmentation results

# Deconvolution Layer

- Or convolution transpose

**Convolution**



**As matrix multiplication**

$$\text{vec } \mathbf{y} = [\quad] \times \text{vec } \mathbf{x}$$

Banded matrix equivalent to $F$

**Convolution transpose**



**Transposed**

$$\text{vec } \mathbf{y} = \times \text{vec } \mathbf{x}$$

Transposed matrix

110

# Deconvolution Layer

- Or convolution transpose

**Convolution**



**As matrix multiplication**

vec **y** = [ ] × vec **x**

Banded matrix equivalent to *F*

**Convolution transpose**



**Transposed**

vec **y** = × vec **x**

Transposed matrix

# Deconvolution Layer

- Or convolution transpose

**Convolution**



$F$

**As matrix multiplication**

vec $\mathbf{y}$ $=$ $\begin{bmatrix} \\ \\ \end{bmatrix}$ $\times$ vec $\mathbf{x}$

Banded matrix equivalent to $F$

**Convolution transpose**



$*^T$

$F$

**Transposed**

vec $\mathbf{y}$ $=$ $\times$ vec $\mathbf{x}$

Transposed matrix

# U-Architectures

- Image to image

net

skip
layers

input image

# U-Architectures

- Image to image



**input image**

**segmentation mask
(output image)**

# U-Architectures

- Image to image



net

net

net

net

net

**skip layers**

**input image**

**segmentation mask (output image)**

# U-Architectures

- Several variants: FCN, U-arch, deco

J. Long, E. Shelhamer, and T. Darrell. **Fully convolutional models for semantic segmentation.** In CVPR 2015
H. Noh, S. Hong, and B. Han. **Learning deconvolution network for semantic segmentation.** In ICCV 2015
O. Ronneberger, P. Fischer, and T. Brox. **U-net: Convolutional networks for biomedical image segmentation.** In MICCAI 2015

116

# Object Detection

# MS COCO Dataset Images

# MS COCO
## Annotations

- 80 different categories

# MS COCO
Dataset Images
+
Annotations

# COCO Object Detection Average Precision (%)

- Area under a detector's precision-recall curve, averaged over…
  - Object categories
  - True positive overlap requirement (IoU from 0.5 to 0.95; see below)

boxes



IoU = 0.5          IoU = 0.7          IoU = 0.9

masks



Ground truth          IoU = 0.55          IoU = 0.70          IoU = 0.91

Figure credits: Dollár and Zitnick (top), Krähenbühl and Kulton (bottom)

# More than one "stage" (≈proposal based; but doesn't require proposals)

classification of reduced output space elements

Rol transformation

aeroplane? no.

:

person? yes.

:

tvmonitor? no.

Cascade-like reduction in output space

Input image

Object / region proposals

Deep Learning region classifier

Region classification, box regression

# One stage

Direct classification
Of all output space elements



1. Resize image.
2. Run convolutional network.
3. Non-max suppression.

"You only look once»
"Single shot"

Redmond et al. You Only Look Once:
Unified Real-time Object Detection. In CVPR 2016

123

# COCO Object Detection Average Precision (%)

Past
(best circa
2012)

5

DPM
(Pre DL)

Felzenszwalb, Girshick, McAllester, Ramanan. Object Detection with Discriminatively Trained Part Based Models. PAMI 2010.

# COCO Object Detection Average Precision (%)

Past
(best circa
2012)

Early
2015

15

5

} Movement to
DL methods

DPM
(Pre DL)

Fast R-CNN
(AlexNet)

Girshick. Fast R-CNN. ICCV 2015.

# COCO Object Detection Average Precision (%)

Past
(best circa
2012)

Early
2015

19

15

5

DPM
(Pre DL)

Fast R-CNN
(AlexNet)

Fast R-CNN
(VGG-16)

Girshick. Fast R-CNN. ICCV 2015.

# COCO Object Detection Average Precision (%)

Past
(best circa
2012)

Early
2015



| DPM (Pre DL) | Fast R-CNN (AlexNet) | Fast R-CNN (VGG-16) | Faster R-CNN (VGG-16) |
|:---:|:---:|:---:|:---:|
| 5 | 15 | 19 | 29 |

Ren, He, Girshick, Sun. Faster R-CNN: Towards Real-Time Object Detection. NIPS 2015.

# COCO Object Detection Average Precision (%)



Past (best circa 2012)

Early 2015

| 5 | 15 | 19 | 29 | 36 |

DPM (Pre DL) — Fast R-CNN (AlexNet) — Fast R-CNN (VGG-16) — Faster R-CNN (VGG-16) — Faster R-CNN (ResNet-50)

Ren, He, Girshick, Sun. Faster R-CNN: Towards Real-Time Object Detection. NIPS 2015.

128

# COCO Object Detection Average Precision (%)



Past
(best circa
2012)

Early
2015

DPM
(Pre DL)

Fast R-CNN
(AlexNet)

Fast R-CNN
(VGG-16)

Faster R-CNN
(VGG-16)

Faster R-CNN
(ResNet-50)

Faster R-CNN
(R-101-FPN)

5    15    19    29    36    39

Lin et al. Feature Pyramid Networks. CVPR 2017.

# COCO Object Detection Average Precision (%)



Past (best circa 2012) — DPM (Pre DL): 5

Early 2015 — Fast R-CNN (AlexNet): 15

Fast R-CNN (VGG-16): 19

Faster R-CNN (VGG-16): 29

Faster R-CNN (ResNet-50): 36

Faster R-CNN (R-101-FPN): 39

2017 — Mask R-CNN (X-152-FPN): 46

He, Gkioxari, Dollár, Girshick. Mask R-CNN. ICCV 2017.

# COCO Object Detection Average Precision (%)



Past (best circa 2012)   Early 2015   **2.5 years**   2017

Progress within DL methods

| | | | | | | |
|---|---|---|---|---|---|---|
| 5 | 15 | 19 | 29 | 36 | 39 | 46 |
| DPM (Pre DL) | Fast R-CNN (AlexNet) | Fast R-CNN (VGG-16) | **Faster R-CNN (VGG-16)** | Faster R-CNN (ResNet-50) | Faster R-CNN (R-101-FPN) | Mask R-CNN (X-152-FPN) |

# "Slow" R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$



$I:$

Selective search, Edge Boxes, MCG, ...

① 

②

Crop & warp

ConvNet($r_i$)

③

1-vs-rest SVMs

④

Box regressor

⑤

132

# "Slow" R-CNN



Per-image computation

Per-region computation for each $r_i \in r(I)$

$I$:

Selective search, Edge Boxes, MCG, ...

①

②

Crop & warp

③

ConvNet($r_i$)

④

1-vs-rest SVMs

⑤

Box regressor

**Very heavy per-region computation**
**E.g., 2000 full network evaluations**

133

# "Slow" R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$



$I$:

Identity

Crop & warp

ConvNet($r_i$)

Selective search, Edge Boxes, MCG, ...

1-vs-rest SVMs

Box regressor

# Generalized R-CNN Approach to Detection

Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = f(I)$

$I:$

$r(I)$

$g(f_I, r_i)$

$h(g_i)$

Classification

$\vdots$

Box regressor

# Fast R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$



$I:$

FCN($I$)

Selective search,
Edge Boxes,
MCG, …

RoIPool

MLP

Softmax clf.

Box regressor

**Lightweight per-region computation**

136

# Fast R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$



FCN($I$)

$I$:

Selective search, Edge Boxes, MCG, ...

RoIPool

MLP

Softmax clf.

Box regressor

# Whole-image FCN

- Use **any standard ConvNet** as the "**backbone architecture**"
  – AlexNet, VGG, ResNet, Inception, Inception-ResNet, ResNeXt, DenseNet, …
  – Use the first N layers with spatial extent (e.g., up to "conv5")

$I$:

FCN($I$)

Example feature map dimensions:
(512, H/16, W/16)

# Fast R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$

# RoIPool (on each Proposal)

Transform **arbitrary size proposal** into a **fixed-dimensional** representation (e.g., 2x2)

$f_I = \text{FCN}(I)$

**Proposal Region of Interest (RoI)**

(Variable size RoI)

# RoIPool (on each Proposal)

Transform **arbitrary size proposal** into a
**fixed-dimensional** representation (e.g., 2x2)

$f_I$ = FCN($I$)

Snapped RoI

Proposal
Region of
Interest
(RoI)

(Variable size RoI)

# RoIPool (on each Proposal)

Transform **arbitrary size proposal** into a **fixed-dimensional** representation (e.g., 2x2)

$f_I = \text{FCN}(I)$

**Proposal Region of Interest (RoI)**

**Snapped RoI**

(Variable size RoI)

RoIPool transform

(Fixed dimensional representation)

Feature value is **max** over input cells

# Fast R-CNN



Per-image computation

Per-region computation for each $r_i \in r(I)$

$I$:

FCN($I$)

Selective search, Edge Boxes, MCG, ...

RoIPool

MLP

Softmax clf.

Box regressor

Region proposals have very poor recall
(ok for PASCAL VOC, major bottleneck for COCO)
Also, they can be slow

# Faster R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = \text{FCN}(I)$

$I:$

$\text{RPN}(f_I)$

RoIPool

MLP

Softmax clf.

Box regressor

Learned proposals
Sharing computation with whole-image network

# Region Proposal Network (RPN)

Proposals = sliding window object/not-object classifier + box regression
**inside the same network**



$f_I$ = FCN($I$)

2k scores

classification fc

4k coordinates

regression fc

(Shared over FPN levels)

k anchor boxes

256-d

intermediate fc

sliding window

conv feature map

Anchors are prototypical object boxes

# Mask R-CNN

Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = \text{FCN}(I)$

$I$:

$\text{RPN}(f_I)$

RoIAlign

MLP

Softmax clf.

Box regressor

Mask FCN

# Mask R-CNN



Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = \text{FCN}(I)$

$I:$

$\text{RPN}(f_I)$

RoIAlign

MLP

Softmax clf.

Box regressor

Mask FCN

# RoIAlign (on each Proposal)

Smoothly transform RoI features into
a fixed-dimensional representation (e.g., 2x2)

Grid of bilinear
interpolation points

$f_I$ = FCN($I$)

Proposal
RoI from
RPN

(Variable size RoI)

# RoIAlign (on each Proposal)

Smoothly transform RoI features into
a fixed-dimensional representation (e.g., 2x2)

$f_I = \text{FCN}(I)$

Proposal RoI from RPN

Grid of bilinear interpolation points

(Variable size RoI)

RoIAlign transform

(Fixed dimensional representation)

Feature value is average of interpolated values on grid

# Compare to RoIPool

Preserve alignment or not?

+20% relative at high IoU

| | align? | bilinear? | agg. | AP | AP$_{50}$ | AP$_{75}$ |
|---|---|---|---|---|---|---|
| *RoIPool* [12] | | | max | 26.9 | 48.8 | 26.4 |
| *RoIWarp* [10] | | ✓ | max | 27.2 | 49.2 | 27.1 |
| | | ✓ | ave | 27.1 | 48.9 | 27.1 |
| *RoIAlign* | ✓ | ✓ | max | **30.2** | **51.0** | **31.8** |
| | ✓ | ✓ | ave | **30.3** | **51.2** | **31.5** |

(c) **RoIAlign** (ResNet-50-C4): Mask results with various RoI layers. Our RoIAlign layer improves AP by ~3 points and AP$_{75}$ by ~5 points. Using proper alignment is the only factor that contributes to the large gap between RoI layers.

# Compare to RoIPool

Quantization breaks pixel-to-pixel alignment



Original RoI

RoIPool coordinate quantization

Snapped RoI

# Instance Segmentation

# Mask R-CNN



Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = \text{FCN}(I)$

$I:$

$\text{RPN}(f_I)$

RoIAlign

MLP

Softmax clf.

Box regressor

Mask FCN

153

# Mask Head (on each Proposal)

- Task specific heads for …
  - Object classification
  - Bounding box detection
  - Instance mask prediction



Standard Fast/er R-CNN head

RoIAlign transformed features

RoIAlign

RoI   7×7 ×256   1024   1024   class   box

# Mask Head (on each Proposal)

- Task specific heads for …
  - Object classification
  - Bounding box detection
  - Instance mask prediction



RoIAlign transformed features

Per-proposal FCN predicts instance masks

# Mask R-CNN: Extension to 2D Human Pose



Per-image computation

Per-region computation for each $r_i \in r(I)$

$f_I = \text{FCN}(I)$

$I:$

$\text{RPN}(f_I)$

RoIAlign

MLP

Softmax clf.

Box regressor

Mask FCN

Pose FCN

# Pose Head



7×7 ×256 → 1024 → 1024 → class

box

RoI

14×14 ×256 ×4 → 14×14 ×256 → 28×28 ×256 → 28×28 ~~×80~~ x17

RoI

~~mask~~ keypoints

(Not shown: Head architecture is slightly different for keypoints)

0.94   nose 1.00   left_eye 1.00   right_eye 0.98   left_ear 0.98
right_ear 0.93   left_shoulder 0.97   right_shoulder 1.00   left_elbow 0.41   right_elbow 0.99
left_wrist 0.91   right_wrist 0.97   left_hip 0.96   right_hip 0.97   left_knee 0.99
right_knee 0.99   left_ankle 0.91   right_ankle 0.98

17 keypoint "mask" predictions shown as heatmaps with OKS scores from argmax positions

- Add keypoint head (28x28x17)
- Predict one "mask" for each keypoint
- Softmax over spatial locations (encodes one keypoint per mask "prior")

157

# Mask R-CNN: Training

- Same as "image centric" Fast/er R-CNN training

- But with <span style="color:red">training targets for masks</span>

# Example Mask Training Targets



Image with training proposal     28x28 mask target     Image with training proposal     28x28 mask target

# Mask R-CNN: Inference

1. **Perform Faster R-CNN inference**
   - Run backbone FCN
   - Generate proposals with RPN
   - Score the proposals with clf. head
   - Refine proposals with box regressor
   - Apply NMS and take the top K (= 100, e.g.)

2. **Run RoIAlign and mask head on top-$K$ refined, post-NMS boxes**
   - Fast (only compute masks for top-K detections)
   - Improves accuracy (uses refined detection boxes, not proposals)

# Mask Prediction

28x28 soft prediction from Mask R-CNN
(enlarged)



Soft prediction **resampled to image coordinates**
(bilinear and bicubic interpolation work equally well)



Final prediction (threshold at 0.5)



Validation image with box detection shown in red

# Mask Prediction



Validation image with box detection shown in red

28x28 soft prediction from Mask R-CNN
(enlarged)



Soft prediction **resampled to image coordinates**
(bilinear and bicubic interpolation work equally well)



Final prediction (threshold at 0.5)

Quantization breaks pixel-to-pixel alignment

Original RoI

RoIPool coordinate quantization

Snapped RoI

163

# Mask Prediction

28x28 soft prediction



Resized soft prediction    Final mask



Validation image with box detection shown in red

# Mask Prediction



28x28 soft prediction

Resized Soft prediction

Final mask

Validation image with box detection shown in red

# Is Object Detection Solved?

- Obviously no; there are <span style="color:red">frequently silly errors</span>

- But it is getting frustratingly good

- The errors are often reasonable

- The bottlenecks are raw recognition and "reasoning"

177

# Addressing other tasks...

# Addressing other tasks...



image

CNN

features

224x224x3

7x7x512

A block of compute with a
few million parameters.

# Addressing other tasks...

# Addressing other tasks...

image
224x224x3

CNN

features
7x7x512

this part changes from task to task

**predicted thing**

**desired thing**

A block of compute with a few million parameters.

# Image Classification

**thing** = a vector of probabilities for different classes



224x224x3

7x7x512

fully connected layer

e.g. vector of 1000 numbers giving probabilities for different classes.

# Segmentation



image        class "map"



image

224x224x3

CNN

features

7x7x512

**deconv layers**

**224x224x20 array of class probabilities at each pixel.**

# Localization



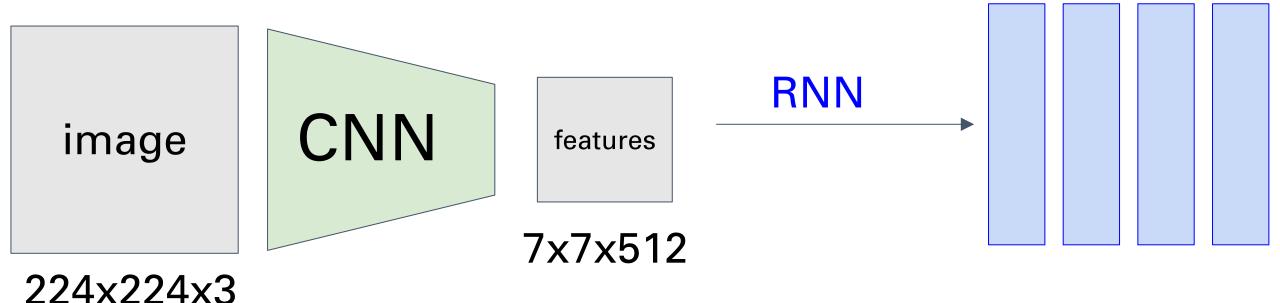image
224x224x3

CNN

features
7x7x512

fully connected layer

Class probabilities (as before)

4 numbers:
- X coord
- Y coord
- Width
- Height

# Image Captioning



A person on a beach flying a kite.
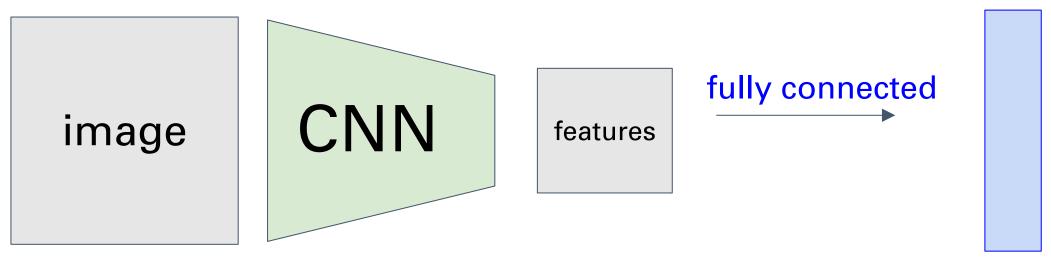
image
224x224x3

CNN

features
7x7x512

RNN

A sequence of 10,000-dimensional vectors giving probabilities of different words in the caption.

# Reinforcement Learning



Mnih et al. 2015



image

CNN

features

fully connected
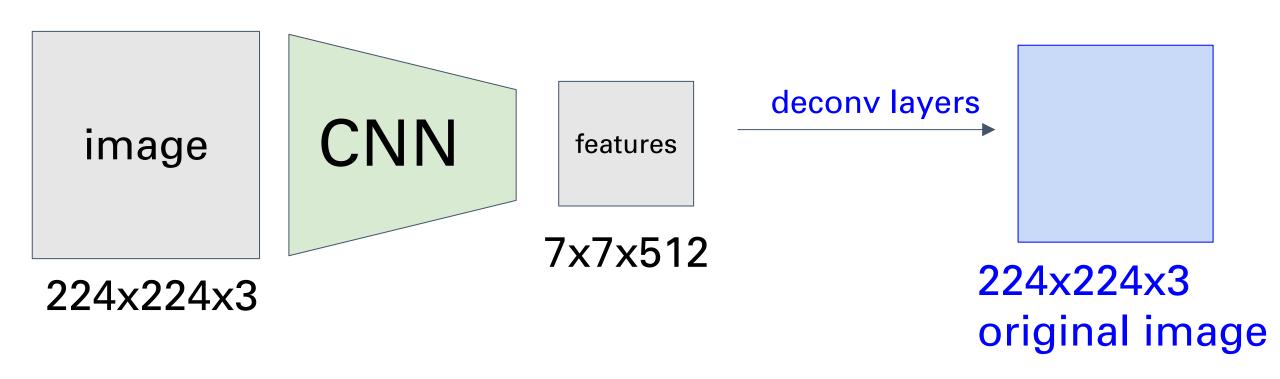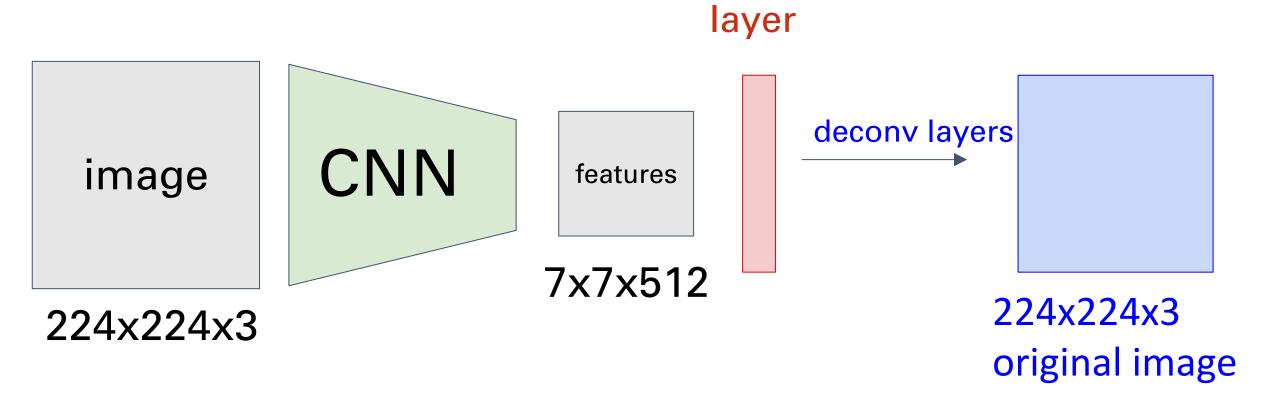
160x210x3

e.g. vector of 8 numbers giving probability of wanting to take any of the 8 possible ATARI actions.

# Autoencoders

image
224x224x3
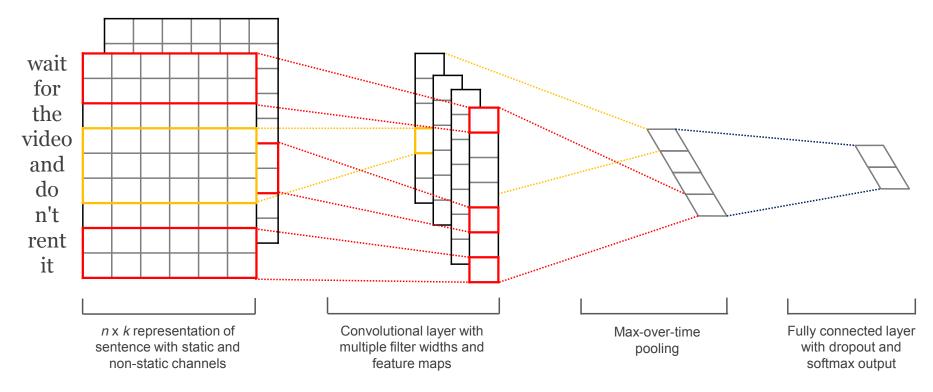
CNN

features
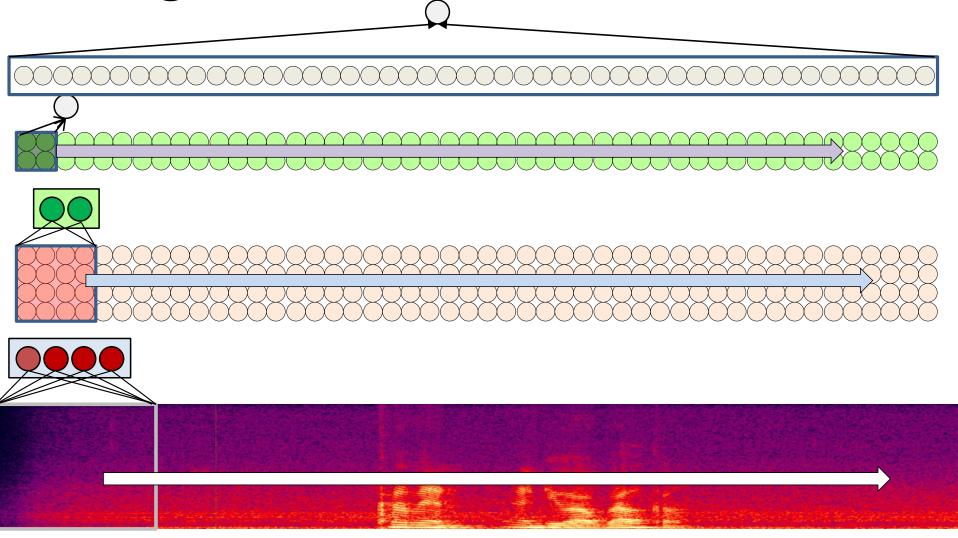7x7x512

deconv layers

224x224x3
original image

# Variational Autoencoders



reparameterization layer

image

CNN

features

7x7x512

deconv layers

224x224x3

224x224x3
original image

[Kingma et al.], [Rezende et al.], [Salimans et al.]

# Addressing other tasks...



n x k representation of sentence with static and non-static channels | Convolutional layer with multiple filter widths and feature maps | Max-over-time pooling | Fully connected layer with dropout and softmax output

- 1D convolution ≈ Time Delay Neural Networks (Waibel et al. 1989, Collobert and Weston 2011)

- Two main paradigms:

  - **Context window modeling:** For tagging, etc. get the surrounding context before tagging
  - **Sentence modeling:** Do convolution to extract n-grams, pooling to combine over whole sentence

Figure credit: Yoon Kim

# Addressing other tasks...



- **CNNs for audio processing:** MFCC features + Time Delay Neural Networks

# **Next lecture:**
# Understanding and Visualizing ConvNets