Lecture 9
# Convolutional neural networks

Information Systems
(Machine Learning)

Andrey Filchenkov

22.11.2018

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- Convolutional neural networks
- Deconvolution and visualization of neurons
- Architecture overview
- Computer vision problems

- The presentation is prepared with materials of D. Polykovsky and K. Khrabrov "Neural networks in machine learning"
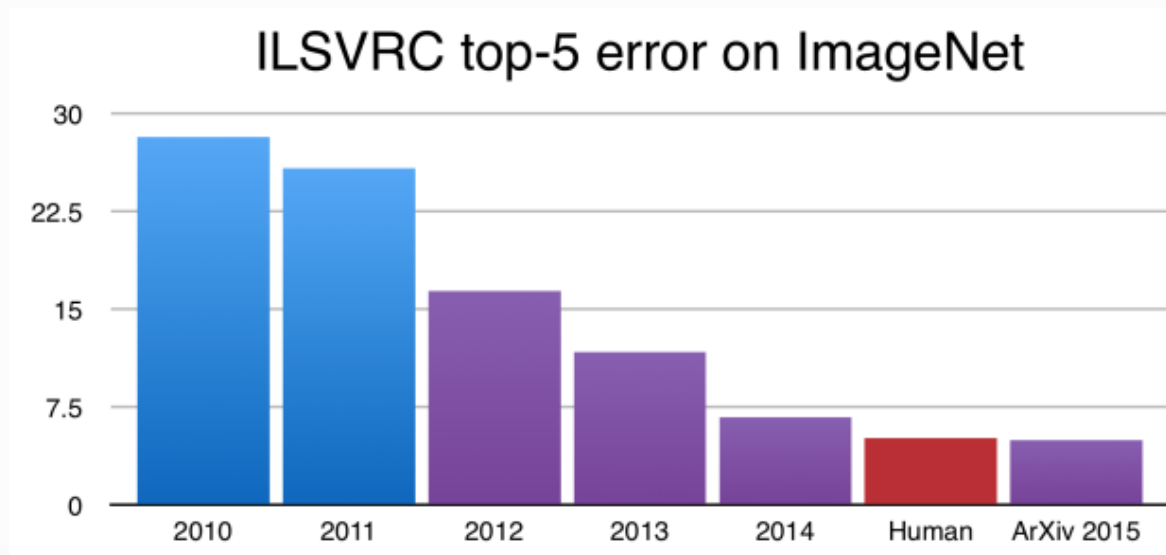- Slides are available online: **goo.gl/BspjhF**

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- Convolutional neural networks
- Deconvolution and visualization of neurons
- Architecture overview
- Computer vision problems

# Today history (reminder)

2012 Hinton, Krizhevsky, and Sutskever suggest Dropout

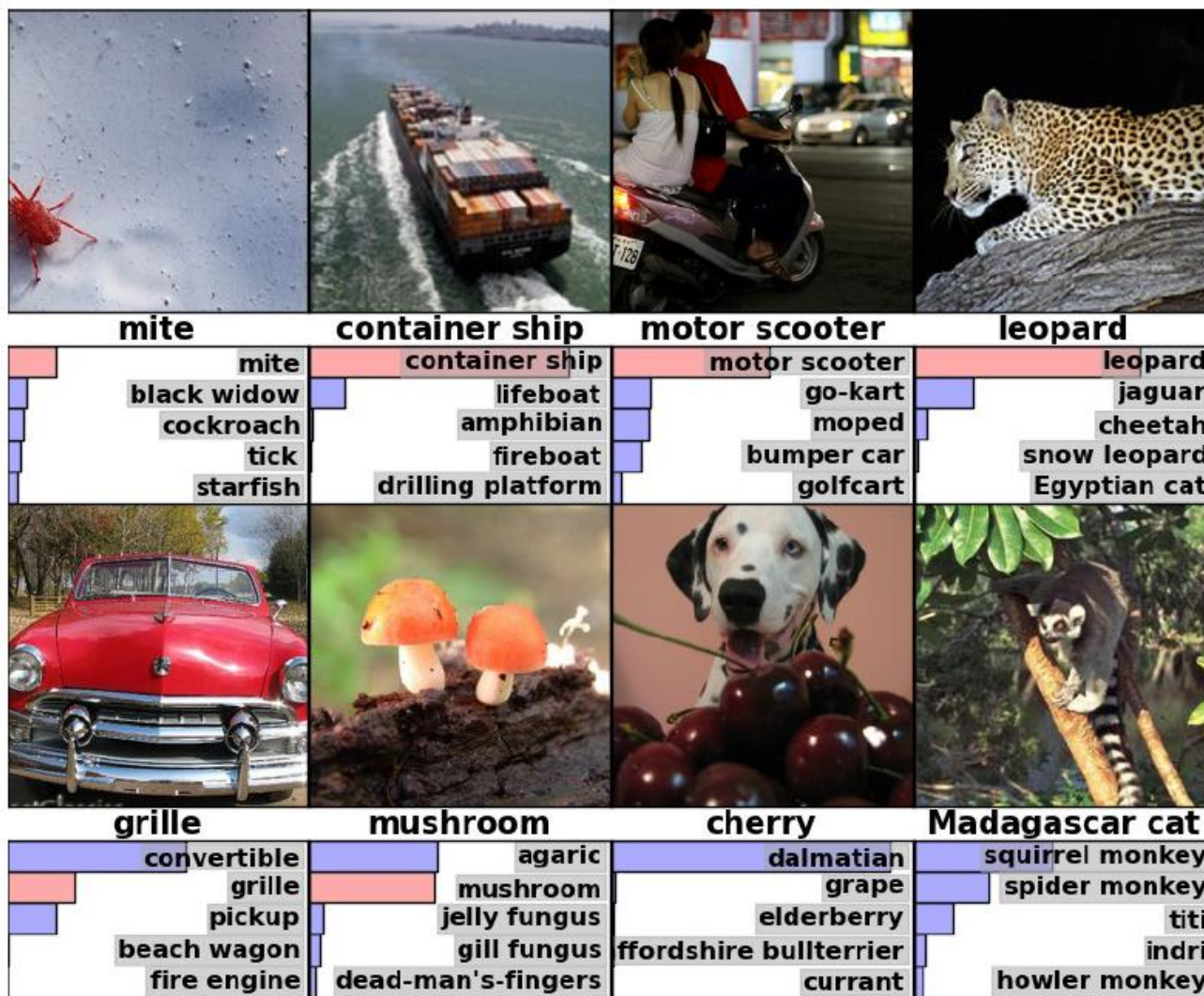2012 They win ImageNet (and two less known competitions). Deep learning era begins.

ILSVRC top-5 error on ImageNet

# Imagenet Challenge



- 1000 images per class
- 1000 classes
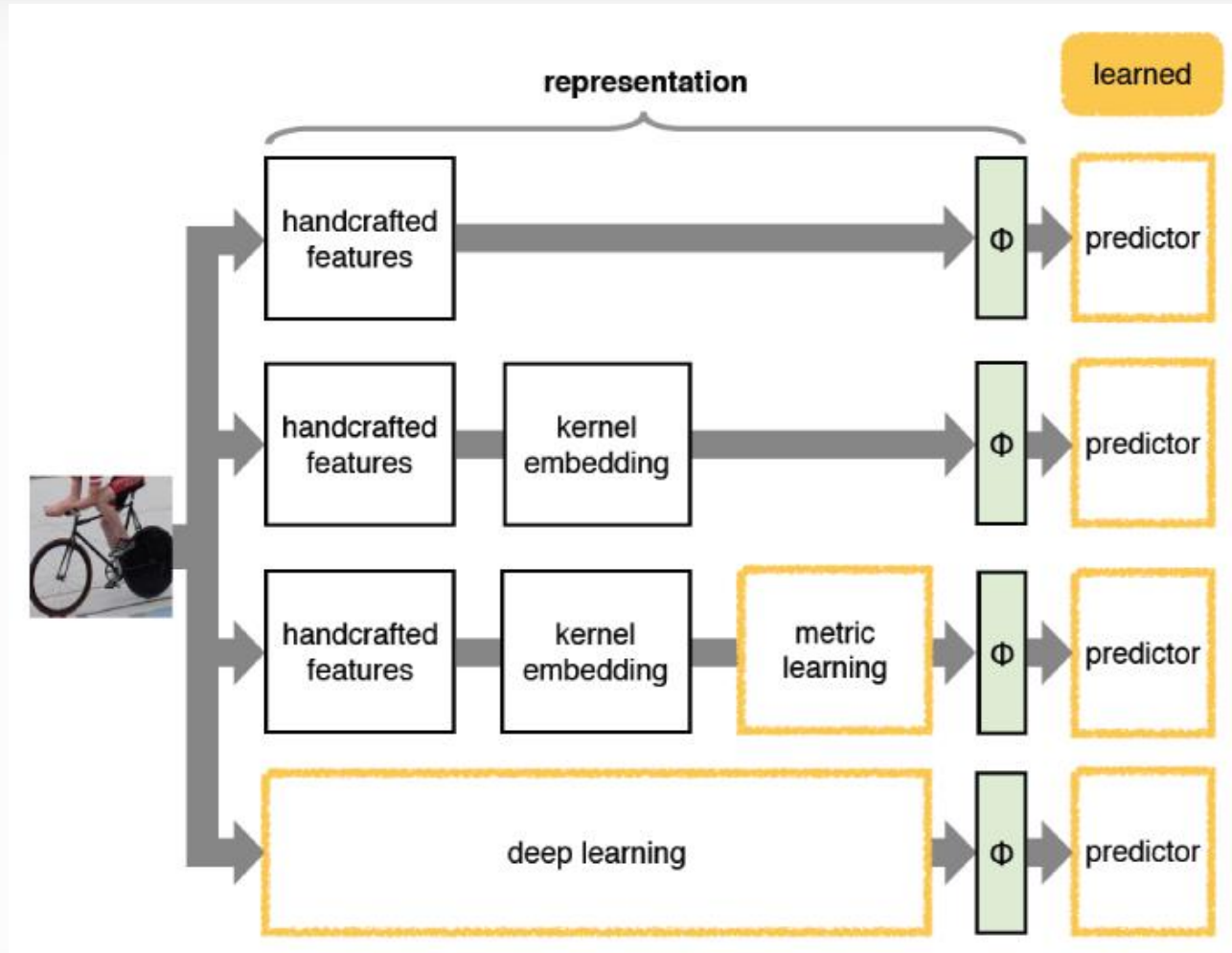- Today, 14 mln images

# Examples of images

# Lecture plan

- Brief overview of ImageNet
- **Earlier approaches in computer vision**
- Convolutional neural networks
- Deconvolution and visualization of neurons
- Architecture overview
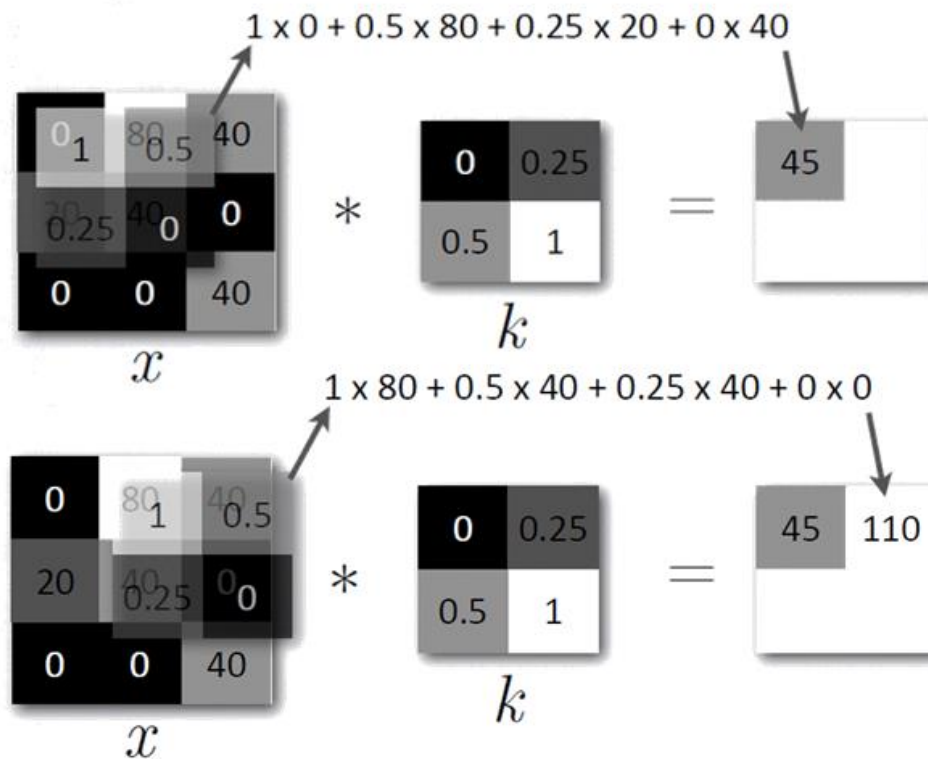- Computer vision problems
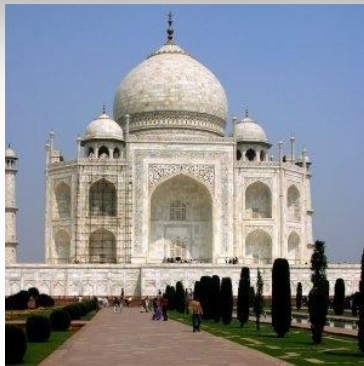
# Short history of computer vision

# Core concepts

- **Local perception**: each neuron sees a small part of the object. Use kernels (filters) to capture 1-D or 2-D structure of objects. For instance, capture all pixel neighbors for an image.
- **Weight sharing:** use small and the same sets of kernels for all objects, this leads to reduction of number of adjusting parameters in comparison with MLP
- **Subsampling/pooling**: use dimensionality reduction for images in order to provide invariance to scale

# Discrete kernel

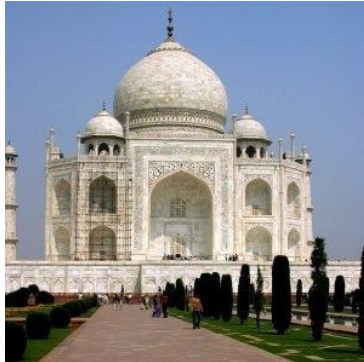$$(x * k)_{ij} = \sum_{pq} x_{i+p,j+q} \; k_{r-p,r-q}$$

1 x 0 + 0.5 x 80 + 0.25 x 20 + 0 x 40

| 0 1 | 80 0.5 | 40 |
|---|---|---|
| 0.25 | 40 0 | 0 |
| 0 | 0 | 40 |

$x$

| 0 | 0.25 |
|---|---|
| 0.5 | 1 |

$k$

$*$

$=$

| 45 | |
|---|---|
| | |

1 x 80 + 0.5 x 40 + 0.25 x 40 + 0 x 0

| 0 | 80 1 | 40 0.5 |
|---|---|---|
| 20 | 40 0.25 | 0 0 |
| 0 | 0 | 40 |

$x$

| 0 | 0.25 |
|---|---|
| 0.5 | 1 |

$k$

$*$

$=$

| 45 | 110 |
|---|---|
| | |

# What kernels can do?



$*$ $\begin{array}{|c|c|c|c|c|} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$ $=$

blur

$*$ $\begin{array}{|c|c|c|} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{array}$ $=$

edge detection

$*$ $\begin{array}{|c|c|c|c|c|} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & -1 & 5 & -1 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array}$ $=$
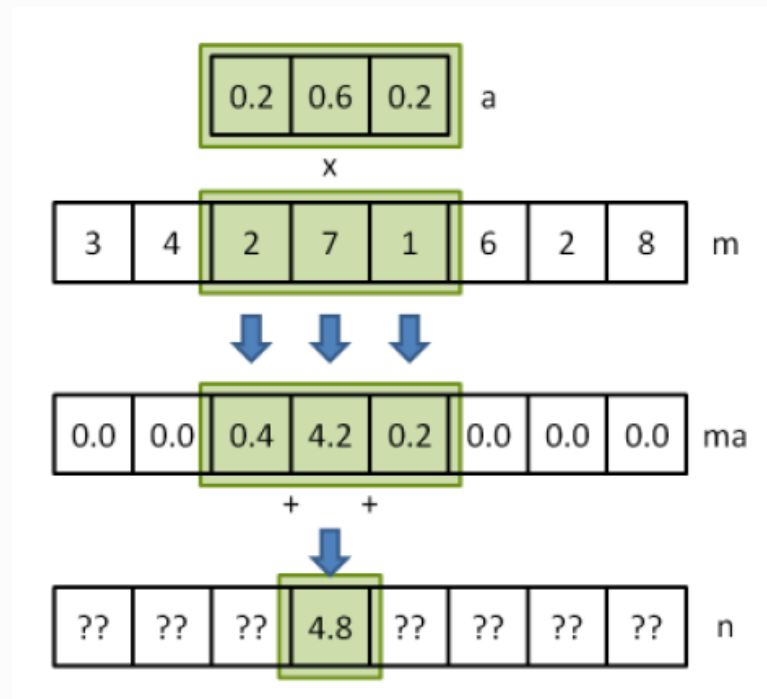
sharpen

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- **Convolutional neural networks**
- Deconvolution and visualization of neurons
- Architecture overview
- Computer vision problems

# Convolution

**Convolution of array $m$ with kernel $a$** is an array $ma[k] = \sum_{i=-w}^{w} m[k+i]\, a[-i]$

# Convolution properties

- Associative property
- Commutative property
- Linearity

# Padding

Zero shift

| 0 | 0 | **A** | **B** | **C** | 0 | 0 |
|---|---|-------|-------|-------|---|---|

Border extension

| A | A | **A** | **B** | **C** | C | C |
|---|---|-------|-------|-------|---|---|

Mirror shift

| B | A | **A** | **B** | **C** | C | B |
|---|---|-------|-------|-------|---|---|

| C | B | **A** | **B** | **C** | B | A |
|---|---|-------|-------|-------|---|---|

Cyclic shift

| B | C | **A** | **B** | **C** | A | B |
|---|---|-------|-------|-------|---|---|

# 2-D convolution



Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.

$(4 \times 0)$
$(0 \times 0)$
$(0 \times 0)$
$(0 \times 0)$
$(0 \times 1)$
$(0 \times 1)$
$(0 \times 0)$
$(0 \times 1)$
$+ (-4 \times 2)$
$-8$

Source pixel

Convolution kernel (emboss)

New pixel value (destination pixel)

# Convolutional tensors

# Pooling

- Pooling is used to reduce dimensionality

## Single depth slice

| 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 |

# VGG-16 conceptual scheme



224 x 224 x 3   224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096   1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

# VGG-16 technical scheme

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- Convolutional neural networks
- **Deconvolution and visualization of neurons**
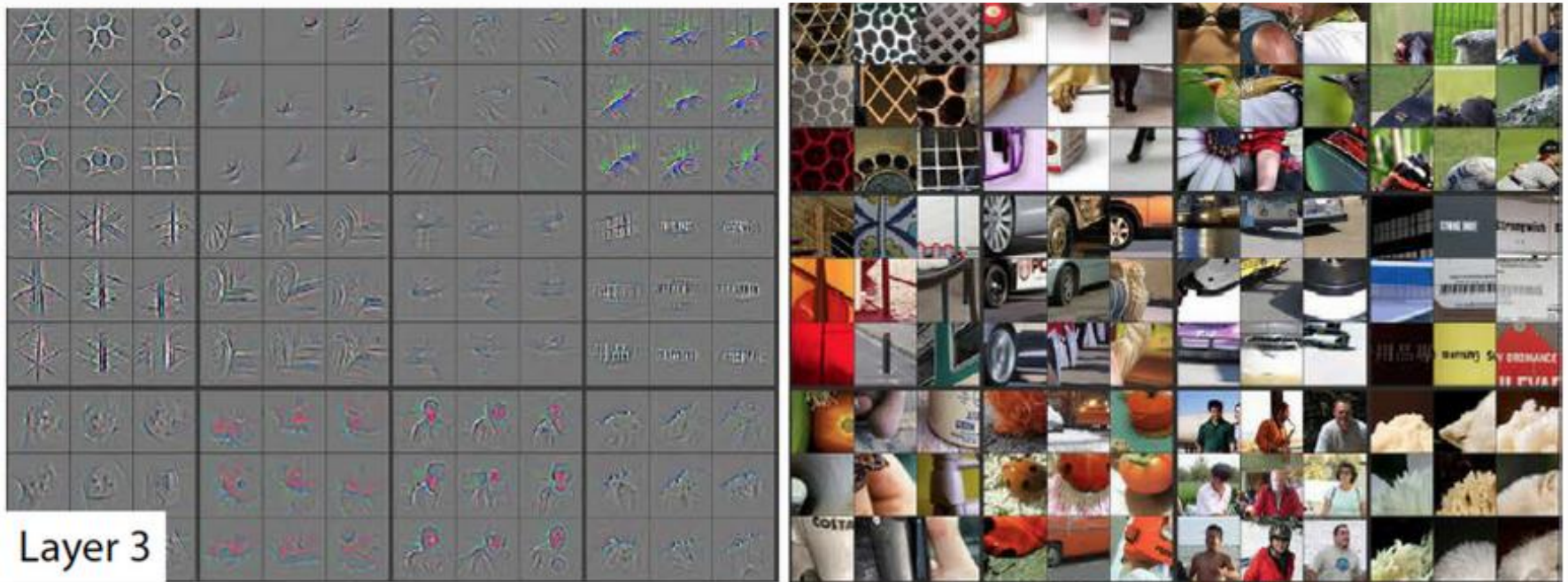- Architecture overview
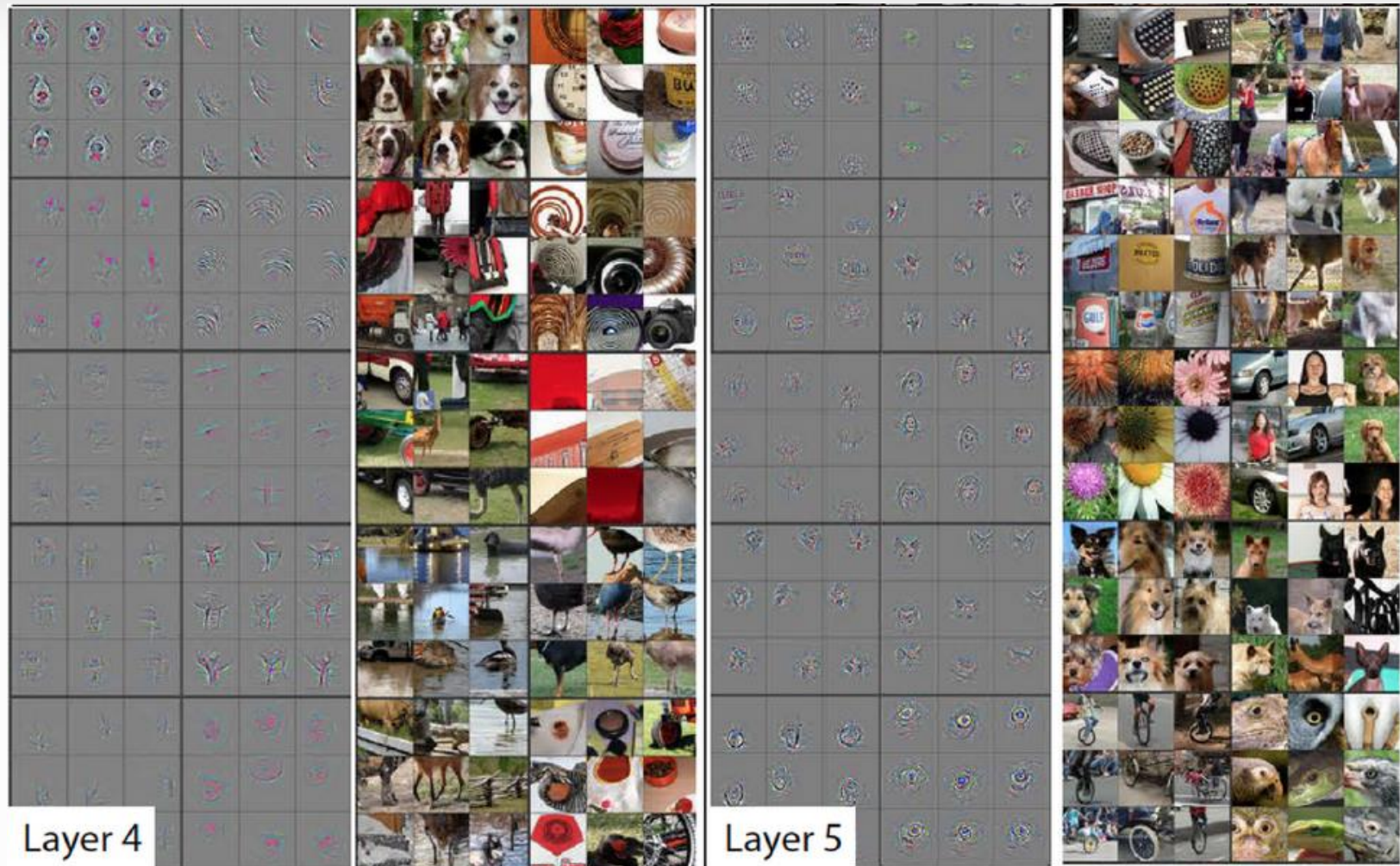- Computer vision problems

# Deconvolution neural network

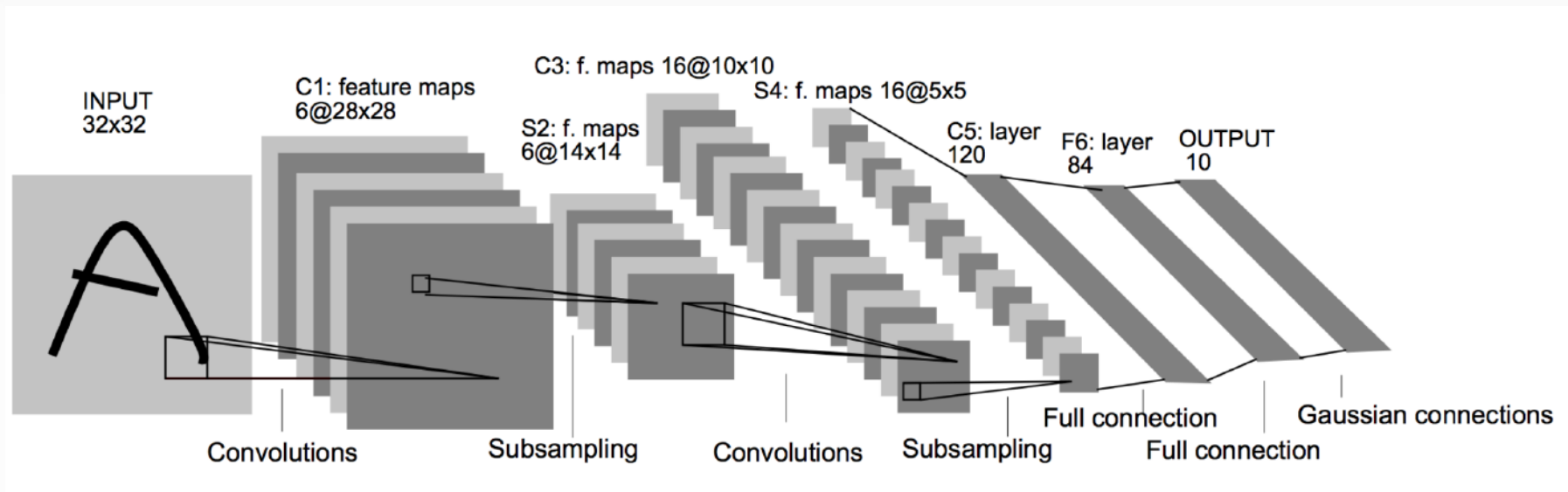# Visualization of neuron activation

# Visualization of neuron activation



Layer 3

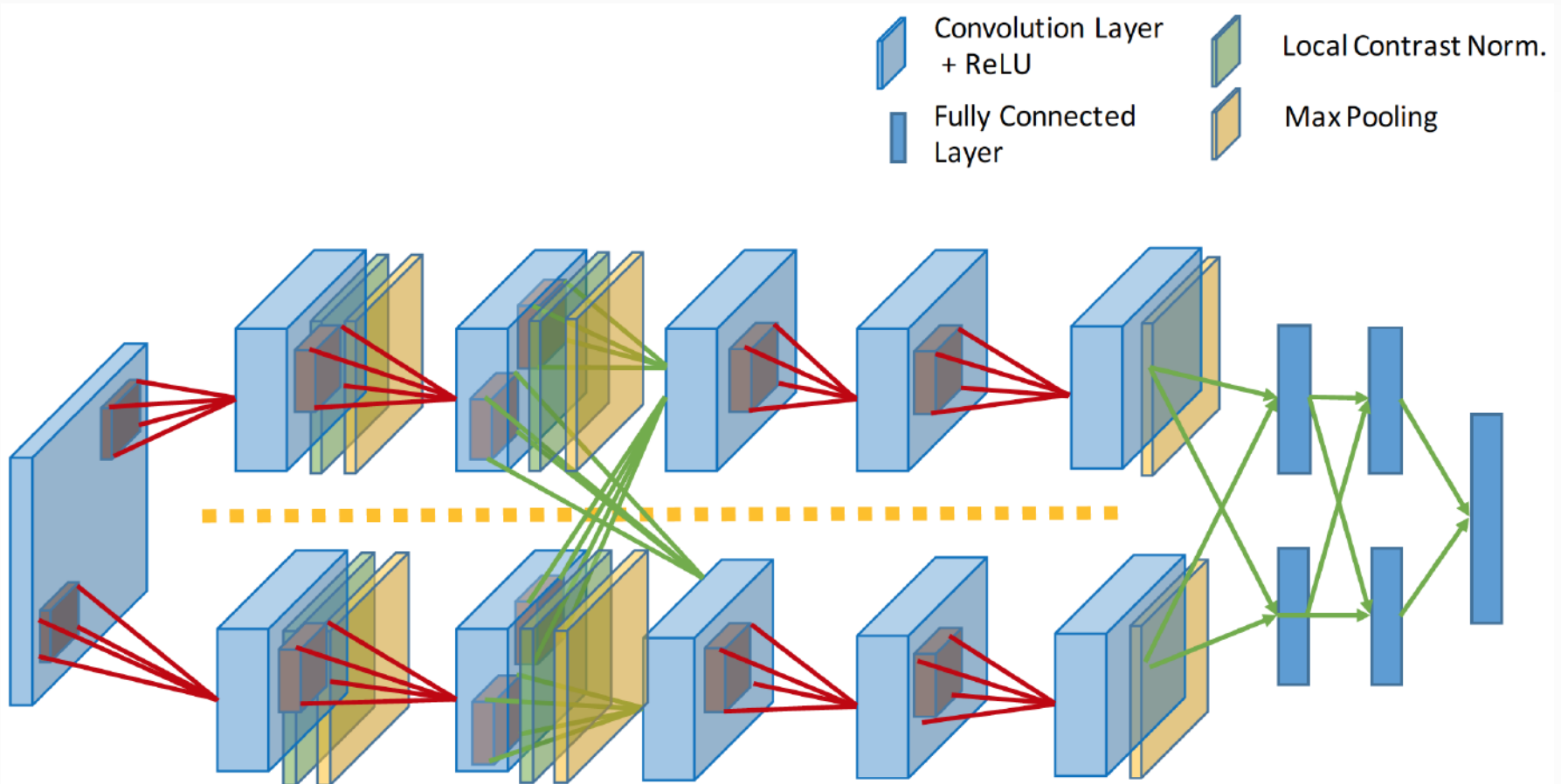# Visualization of neuron activation



Layer 4

Layer 5

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- Convolutional neural networks
- Deconvolution and visualization of neurons
- **Architecture overview**
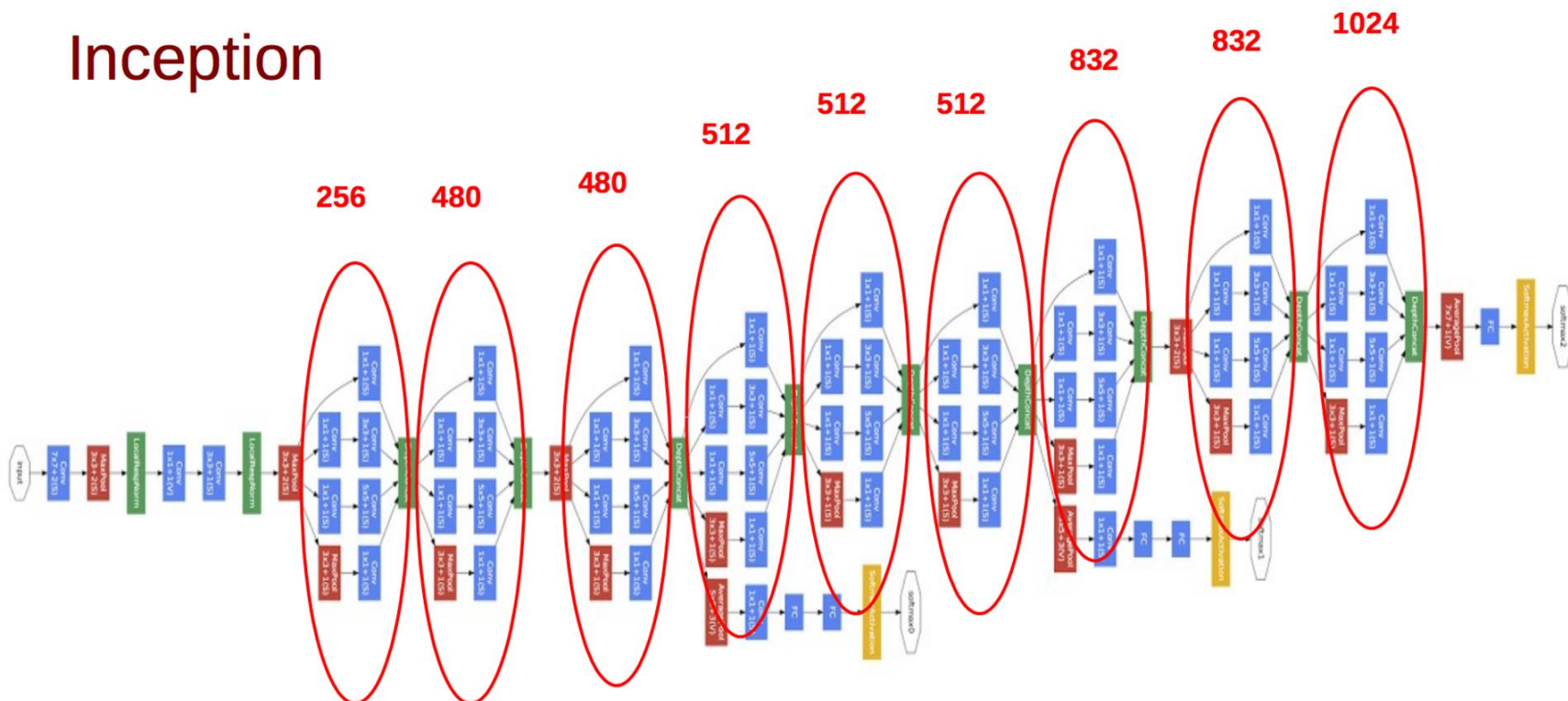- Computer vision problems

# LeNet

# AlexNet

# VGG-16

| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

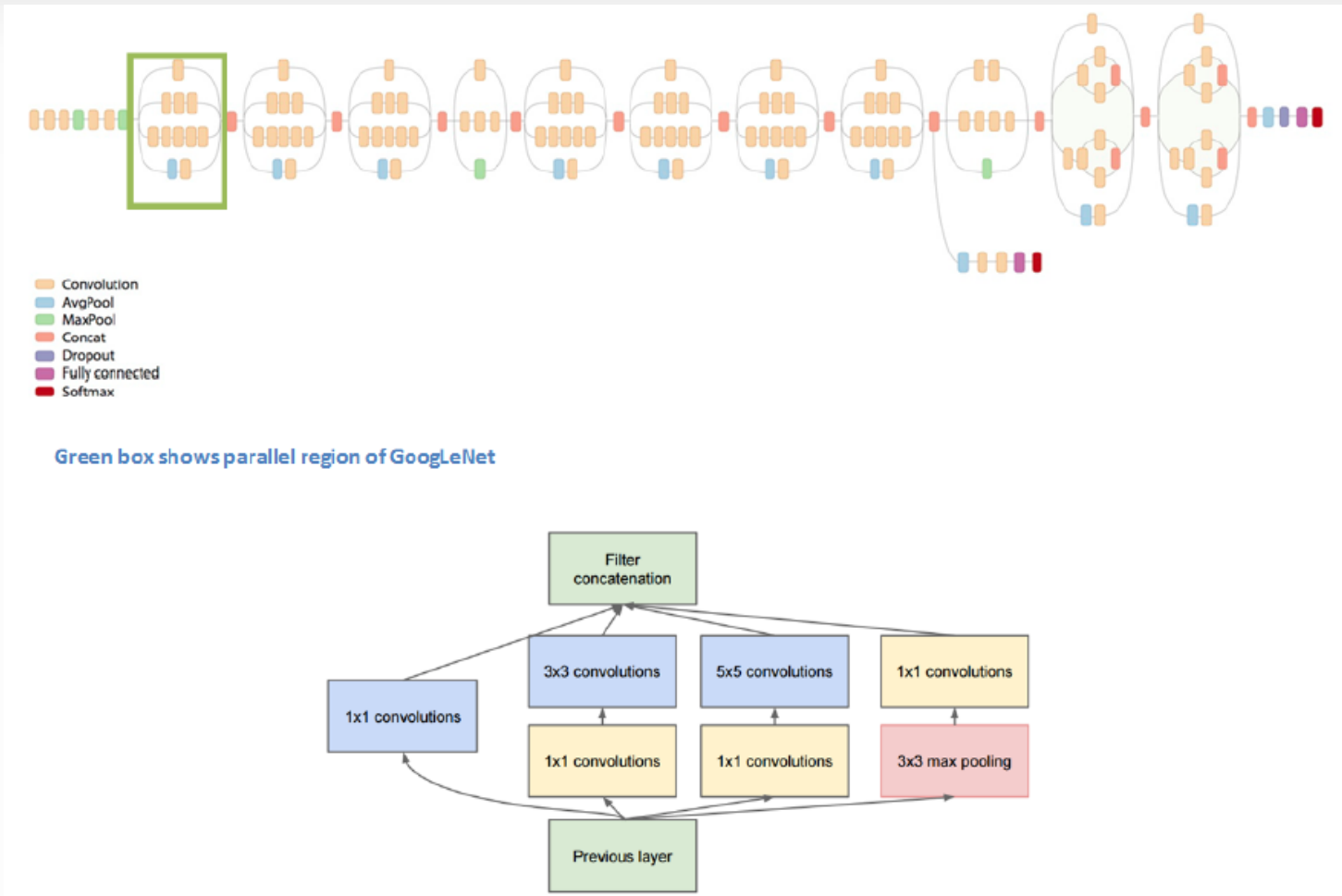IS (Machine learning). Lecture 9. Convolutional. 22.11.2018.

# Inception



Width of **inception modules** ranges from 256 filters (in early modules) to 1024 in top inception modules.

# Inception aka GoogleNet



Convolution
AvgPool
MaxPool
Concat
Dropout
Fully connected
Softmax

Green box shows parallel region of GoogLeNet
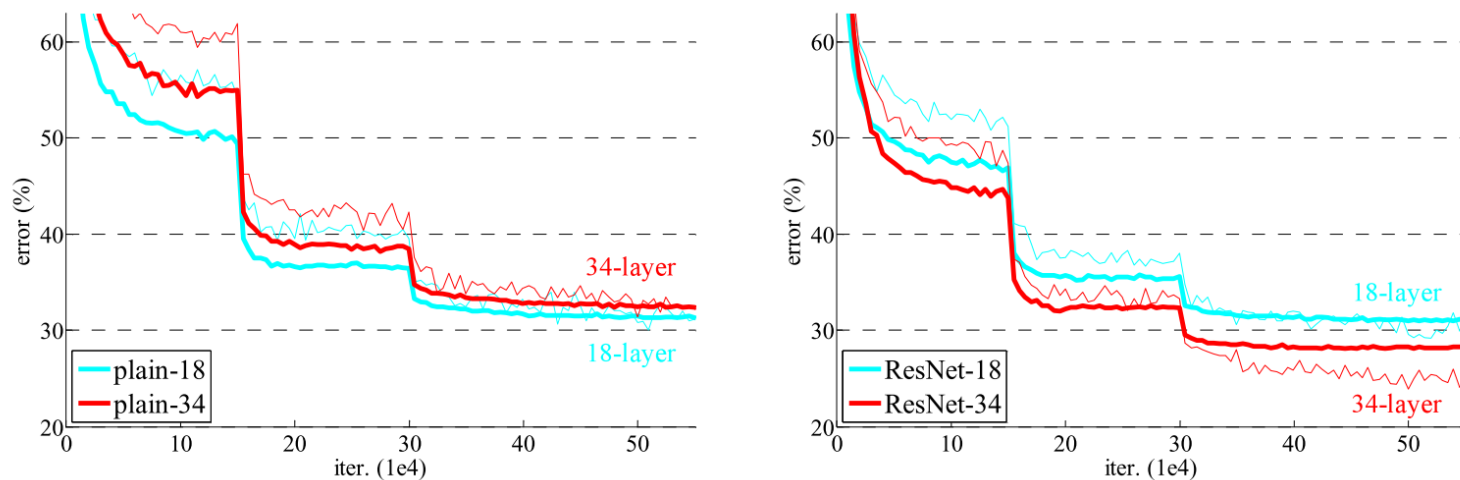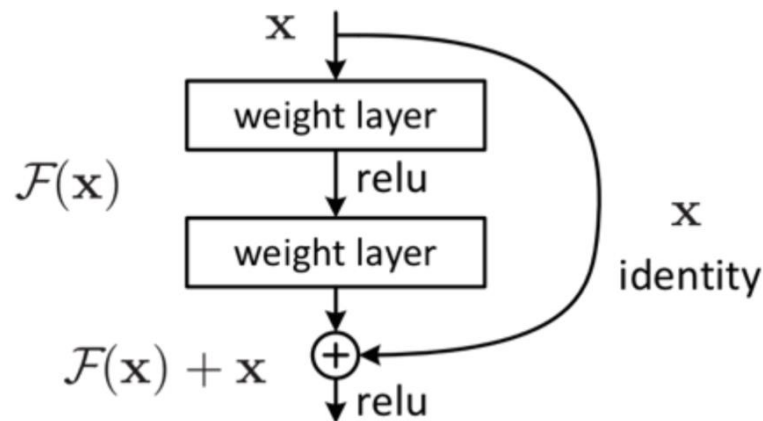
# ResNet (1/2)

- Additional layers not always help



Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

- Adding skip layers may help



Figure 2. Residual learning: a building block.

$\mathcal{H}(\mathbf{x})$ is the true function we want to learn

Let's pretend we want to learn

$$\mathcal{F}(\mathbf{x}) := \mathcal{H}(\mathbf{x}) - \mathbf{x}$$

instead.

The original function is then

$$\mathcal{F}(\mathbf{x}) + \mathbf{x}$$

# Network comparison

# Lecture plan

- Brief overview of ImageNet
- Earlier approaches in computer vision
- Convolutional neural networks
- Deconvolution and visualization of neurons
- Architecture overview
- Computer vision problems

# CV tasks
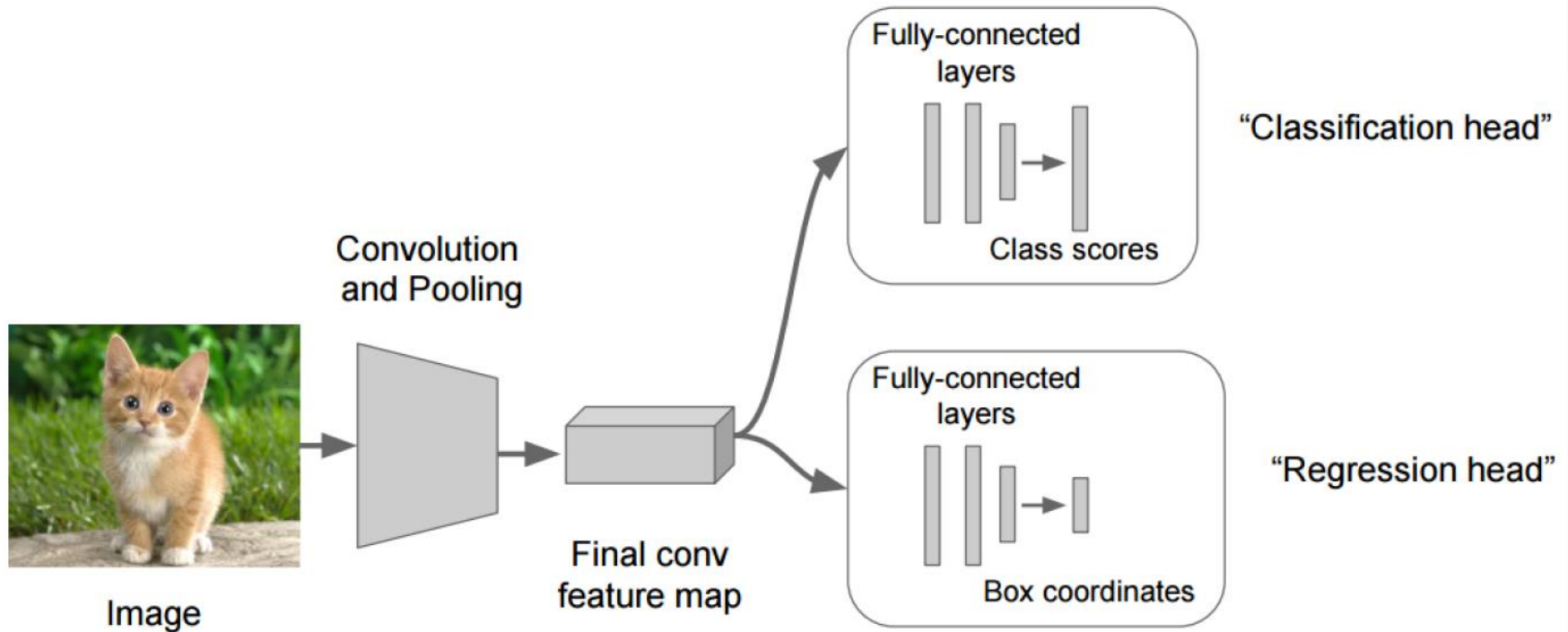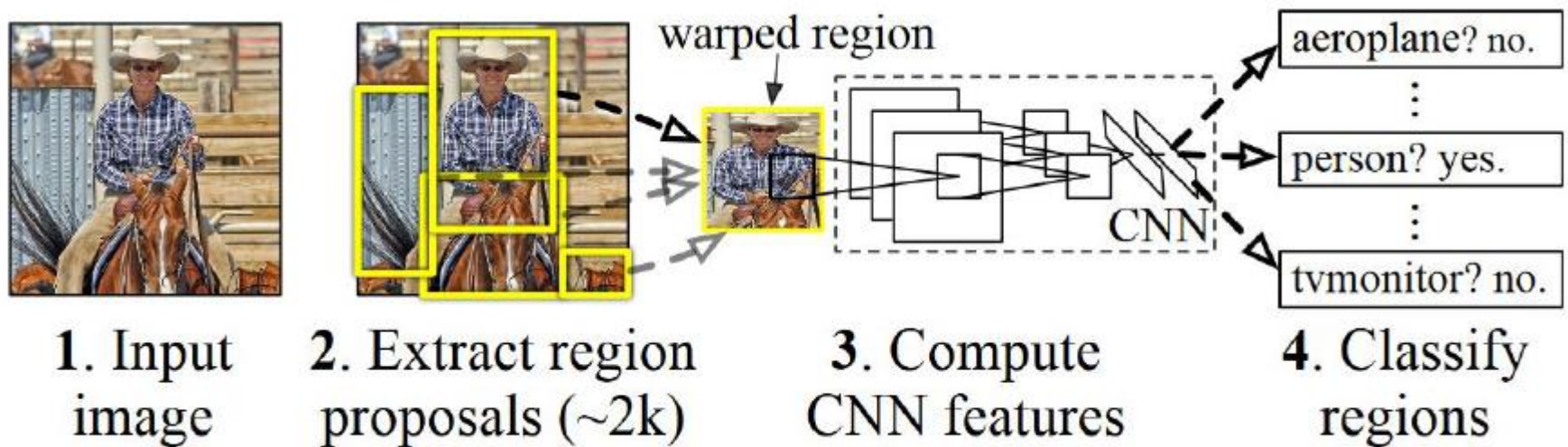
# Semantic segmentation

# U-Net

# Object localization

# Object detection. Pascal VOC

# Detection via R-CNN



1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

warped region

aeroplane? no.
person? yes.
tvmonitor? no.

CNN

# Detection via YOLO