# Computer Vision – Lecture 12

## Deep Learning III

**17.06.2019**

Bastian Leibe

Visual Computing Institute
RWTH Aachen University
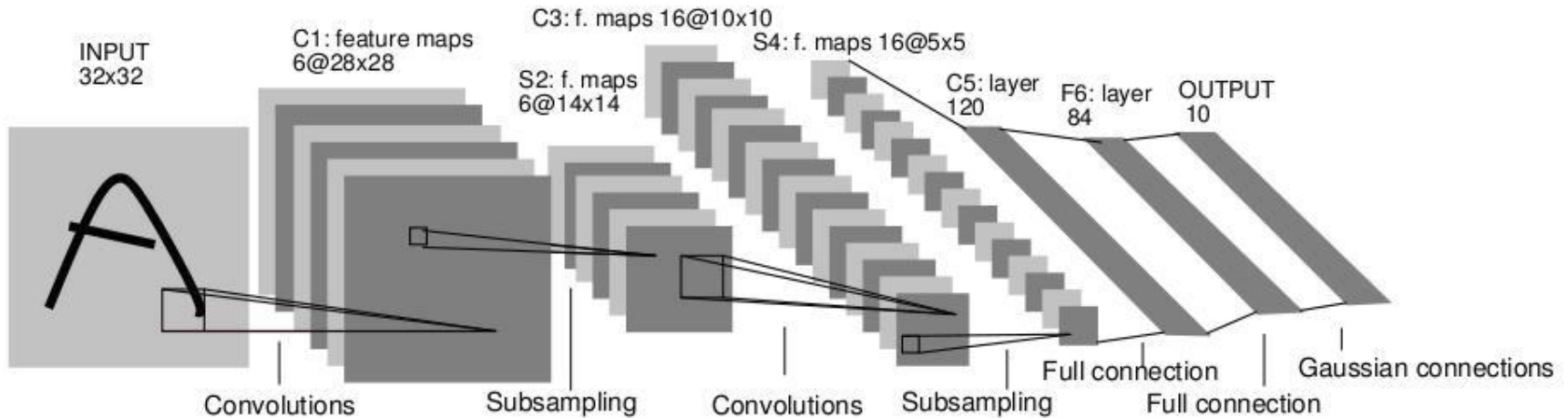http://www.vision.rwth-aachen.de/

leibe@vision.rwth-aachen.de

# Course Outline

- Image Processing Basics

- Segmentation & Grouping

- Object Recognition & Categorization
  - Sliding Window based Object Detection

- Local Features & Matching

- Deep Learning
  - Convolutional Neural Networks (CNNs)
  - Deep Learning Background
  - CNNs for Object Detection
  - CNNs for Semantic Segmentation
  - CNNs for Matching

- 3D Reconstruction

# Topics of This Lecture

- **CNN Architectures**
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNet

- **CNNs for Object Detection**
  - R-CNN
  - Fast R-CNN
  - Faster R-CNN
  - Mask R-CNN
  - YOLO / SSD

B. Leibe

# CNN Architectures: LeNet (1998)

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 — Full connection — Gaussian connections

- **Early convolutional architecture**
  - 2 Convolutional layers, 2 pooling layers
  - Fully-connected NN layers for classification
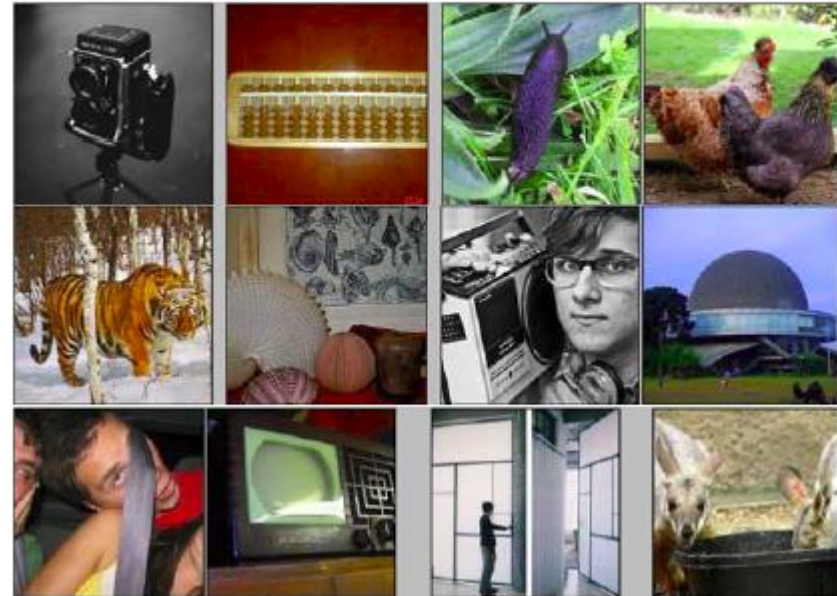  - Successfully used for handwritten digit recognition (MNIST)

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.

Slide credit: Svetlana Lazebnik

B. Leibe

# ImageNet Challenge 2012

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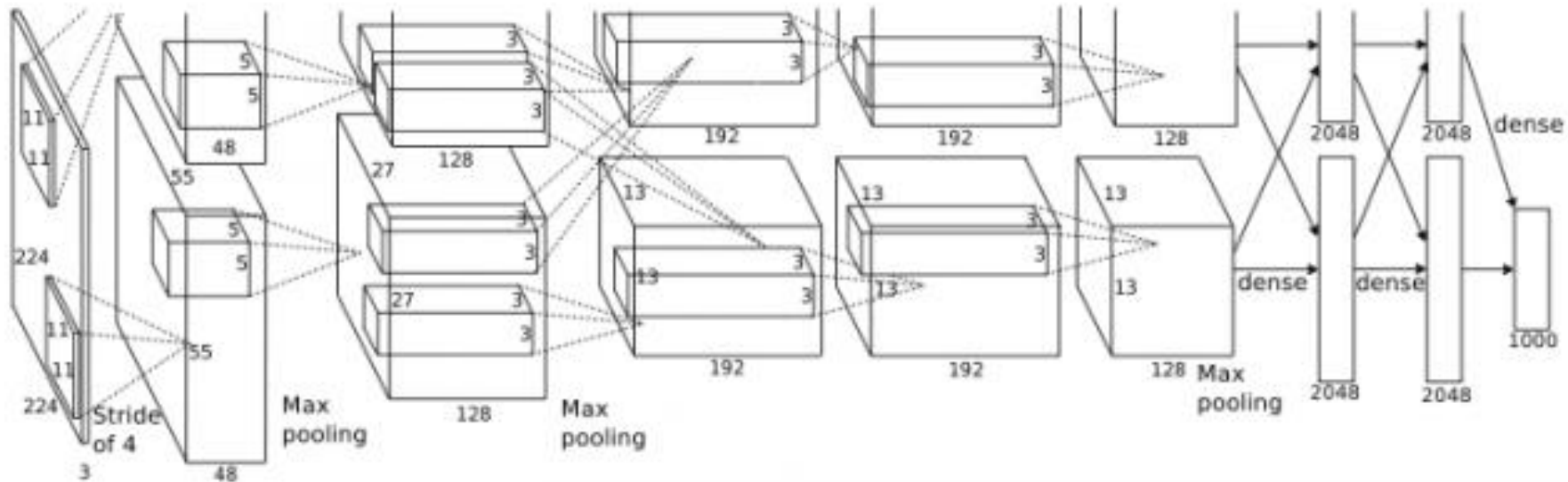  - Human labels via Amazon Mechanical Turk

- Challenge (ILSVRC)
  - 1.2 million training images
  - 1000 classes
  - Goal: Predict ground-truth class within top-5 responses
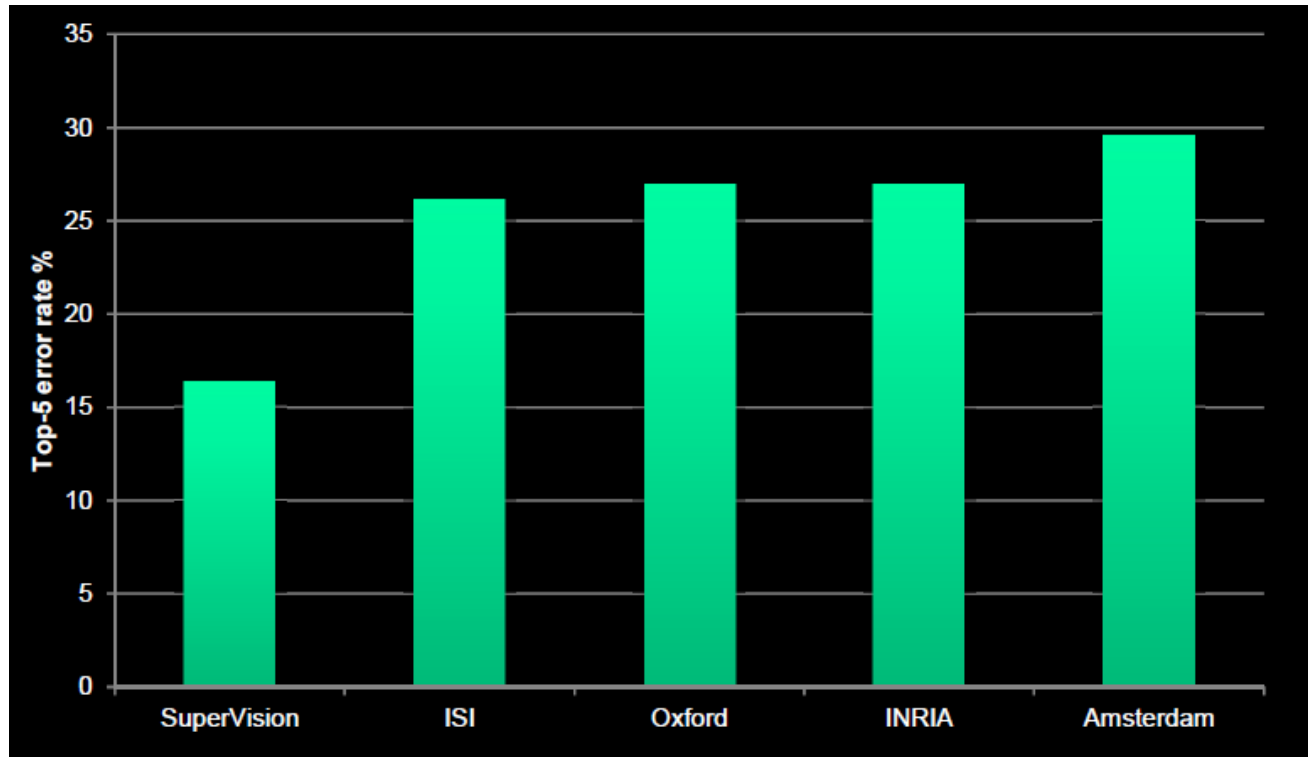  - Currently one of the top benchmarks in Computer Vision

[Deng et al., CVPR'09]

B. Leibe

5

# CNN Architectures: AlexNet (2012)



- ## Similar framework as LeNet, but
  - Bigger model (7 hidden layers, 650k units, 60M parameters)
  - More data ($10^6$ images instead of $10^3$)
  - GPU implementation
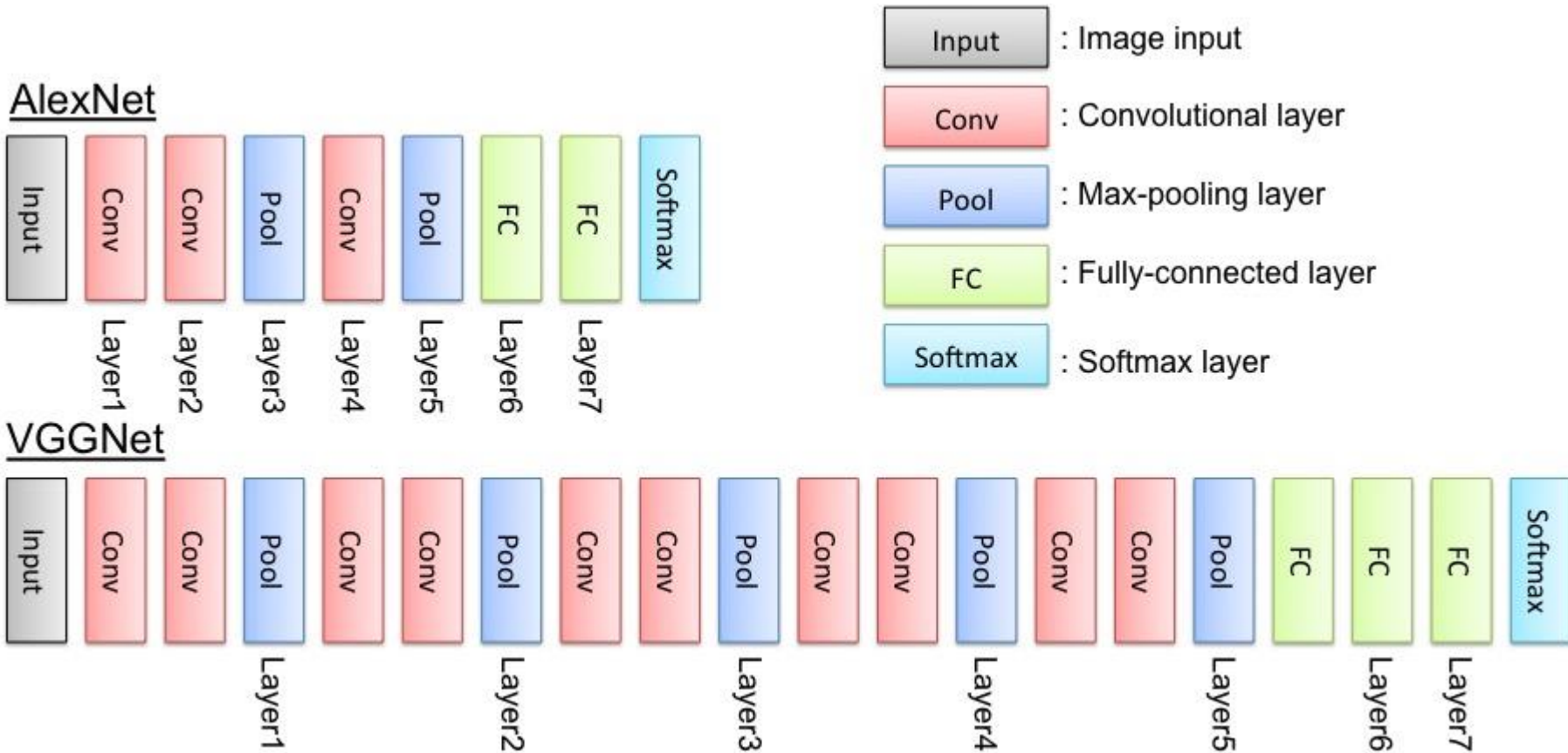  - Better regularization and up-to-date tricks for training (Dropout)

A. Krizhevsky, I. Sutskever, and G. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012.

Image source: A. Krizhevsky, I. Sutskever and G.E. Hinton, NIPS 2012

# ILSVRC 2012 Results



- **AlexNet almost halved the error rate**
  - ➢ 16.4% error (top-5) vs. 26.2% for the next best approach
  - ⇒ A revolution in Computer Vision
  - ➢ Acquired by Google in Jan '13, deployed in Google+ in May '13

B. Leibe

# CNN Architectures: VGGNet (2014/15)



K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015

Computer Vision Summer'19

B. Leibe

Image source: Hirokatsu Kataoka

# CNN Architectures: VGGNet (2014/15)

- ## Main ideas
  - ➢ Deeper network
  - ➢ Stacked convolutional layers with smaller filters (+ nonlinearity)
  - ➢ Detailed evaluation of all components
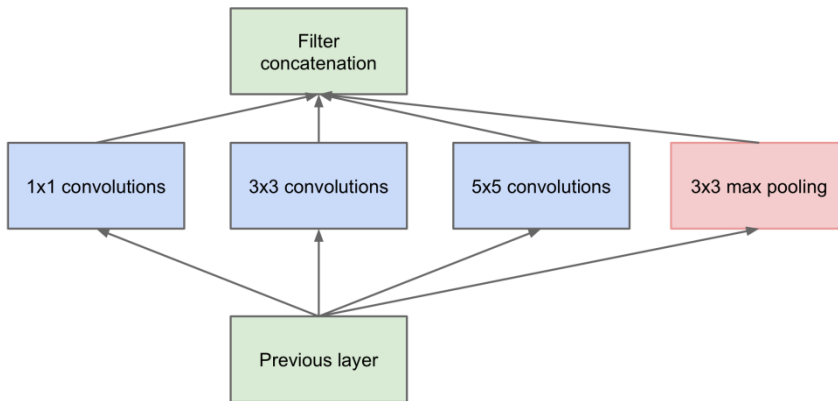
- ## Results
  - ➢ Improved ILSVRC top-5 error rate to 6.7%.
  - ➢ 138M parameters (VGG16), most of those in the FC layers (102M)

| 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 | | | | | |

Mainly used

B. Leibe

Image source: Simonyan & Zisserman

# Comparison: AlexNet vs. VGGNet

- **Receptive fields in the first layer**
  - AlexNet:        $11 \times 11$, stride 4
  - Zeiler & Fergus:  $7 \times 7$, stride 2
  - VGGNet:        $3 \times 3$, stride 1

- **Why that?**
  - If you stack a $3 \times 3$ layer on top of another $3 \times 3$ layer, you effectively get a $5 \times 5$ receptive field.
  - With three $3 \times 3$ layers, the receptive field is already $7 \times 7$.
  - But much fewer parameters: $3 \cdot 3^2 = 27$ instead of $7^2 = 49$.
  - In addition, non-linearities in-between $3 \times 3$ layers for additional discriminativity.

B. Leibe

# CNN Architectures: GoogLeNet (2014)
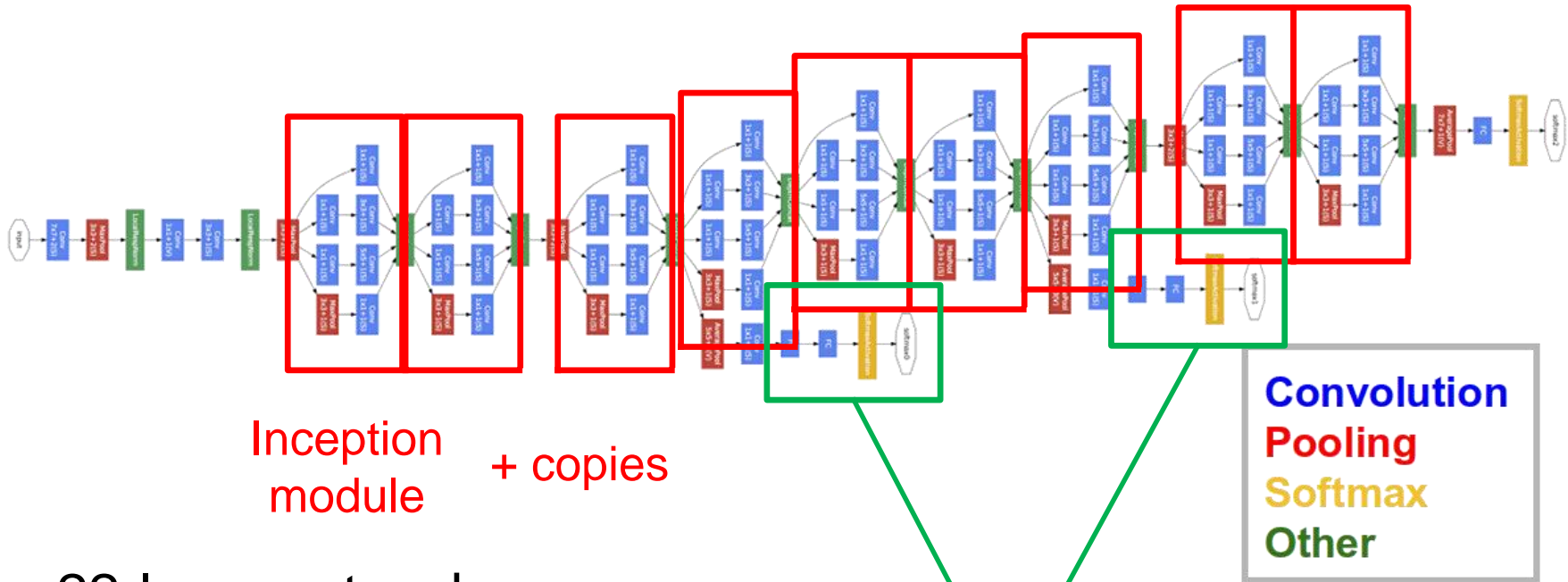


(a) Inception module, naïve version

(b) Inception module with dimension reductions

- **Main ideas**
  - ➢ "Inception" module as modular component
  - ➢ Learns filters at several scales within each module
  - ➢ 1x1 convolutions ("bottleneck layers") for dimensionality reduction

    C. Szegedy, W. Liu, Y. Jia, et al, Going Deeper with Convolutions, arXiv:1409.4842, 2014.

13

# GoogLeNet Visualization



Inception module

+ copies

Auxiliary classification outputs for training the lower layers (deprecated)

**Convolution**
**Pooling**
**Softmax**
**Other**

- 22-layer network
  - No FC layers
  - Only 5M parameters
  - ILSVRC'14 winner with 6.7% top-5 error

B. Leibe

# Results on ILSVRC

| Method | top-1 val. error (%) | top-5 val. error (%) | top-5 test error (%) |
|---|---|---|---|
| VGG (2 nets, multi-crop & dense eval.) | **23.7** | **6.8** | **6.8** |
| VGG (1 net, multi-crop & dense eval.) | 24.4 | 7.1 | 7.0 |
| VGG (ILSVRC submission, 7 nets, dense eval.) | 24.7 | 7.5 | 7.3 |
| GoogLeNet (Szegedy et al., 2014) (1 net) | - | 7.9 | |
| GoogLeNet (Szegedy et al., 2014) (7 nets) | - | **6.7** | |
| MSRA (He et al., 2014) (11 nets) | - | - | 8.1 |
| MSRA (He et al., 2014) (1 net) | 27.9 | 9.1 | 9.1 |
| Clarifai (Russakovsky et al., 2014) (multiple nets) | - | - | 11.7 |
| Clarifai (Russakovsky et al., 2014) (1 net) | - | - | 12.5 |
| Zeiler & Fergus (Zeiler & Fergus, 2013) (6 nets) | 36.0 | 14.7 | 14.8 |
| Zeiler & Fergus (Zeiler & Fergus, 2013) (1 net) | 37.5 | 16.0 | 16.1 |
| OverFeat (Sermanet et al., 2014) (7 nets) | 34.0 | 13.2 | 13.6 |
| OverFeat (Sermanet et al., 2014) (1 net) | 35.7 | 14.2 | - |
| Krizhevsky et al. (Krizhevsky et al., 2012) (5 nets) | 38.1 | 16.4 | 16.4 |
| Krizhevsky et al. (Krizhevsky et al., 2012) (1 net) | 40.7 | 18.2 | - |

- **VGGNet and GoogLeNet perform at similar level**
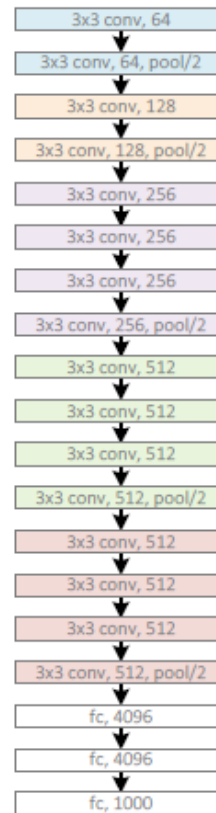  - Comparison: human performance ~5% [Karpathy]

http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/

B. Leibe

Image source: Simonyan & Zisserman

# Residual Networks

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 |

GoogleNet, 22 layers
(ILSVRC 2014)

B. Leibe

16

# Residual Networks

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

ResNet, 152 layers
(ILSVRC 2015)

- Core component
  - Skip connections bypassing each layer
  - Better propagation of gradients to the deeper layers

$F(x)$

$x$

weight layer

relu

weight layer

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

relu

B. Leibe

# ILSRVC Winners

Slide credit: FeiFei Li

B. Leibe

# PASCAL VOC Object Detection Performance



**101 layers**

**Engines of visual recognition**

86

66

58

34

16 layers

8 layers

shallow

HOG, DPM | AlexNet (RCNN) | VGG (RCNN) | ResNet (Faster RCNN)*

PASCAL VOC 2007 **Object Detection** mAP (%)

Slide credit: Kaiming He

B. Leibe

# Comparing Complexity



A. Canziano, A. Paszke, E. Culurcello, An Analysis of Deep Neural Network Models for Practical Applications, arXiv 2017.

Computer Vision Summer'19

Figure credit: Alfredo Canziano, Adam Paszke, Eugenio Culurcello

# The Learned Features are Generic



state of the art
level (pre-CNN)

- Experiment: feature transfer
  - Train AlexNet-like network on ImageNet
  - Chop off last layer and train classification layer on CalTech256
  - ⇒ State of the art accuracy already with only 6 training images!

22

B. Leibe

Image source: M. Zeiler, R. Fergus

# Transfer Learning with CNNs



1. Train on ImageNet

2. If small dataset: fix all weights (treat CNN as fixed feature extractor), retrain only the classifier

I.e., swap the Softmax layer at the end

Slide credit: Andrej Karpathy

B. Leibe

# Transfer Learning with CNNs

1. Train on ImageNet

3. If you have medium sized dataset, "finetune" instead: use the old weights as initialization, train the full network or only some of the higher layers.

Retrain bigger portion of the network

Slide credit: Andrej Karpathy

B. Leibe

24

# Topics of This Lecture

- CNN Architectures
  - LeNet
  - AlexNet
  - VGGNet
  - GoogLeNet
  - ResNet

- CNNs for Object Detection
  - R-CNN
  - Fast R-CNN
  - Faster R-CNN
  - Mask R-CNN
  - YOLO / SSD

B. Leibe

# Object Detection: R-CNN

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



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

- Results on PASCAL VOC Detection benchmark
  - Pre-CNN state of the art:  35.1% mAP  [Uijlings et al., 2013]
    33.4% mAP  DPM
  - R-CNN:  53.7% mAP

R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014
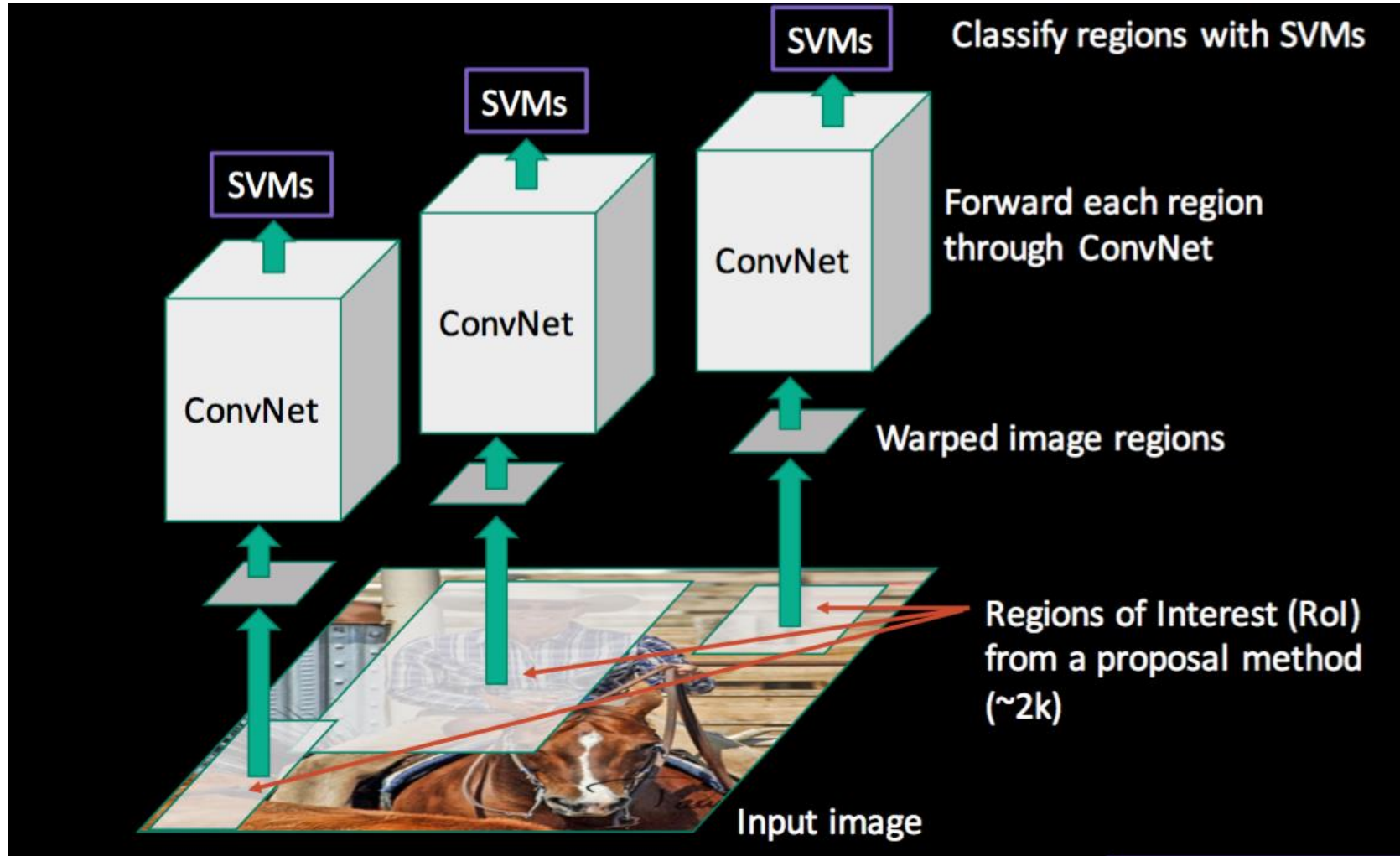
26

# R-CNN Pipeline



Input image

Slide credit: Ross Girshick

B. Leibe

# R-CNN Pipeline



Regions of Interest (RoI) from a proposal method (~2k)

Input image

Slide credit: Ross Girshick

B. Leibe

# R-CNN Pipeline

Slide credit: Ross Girshick

B. Leibe

# R-CNN Pipeline

Slide credit: Ross Girshick

B. Leibe

# R-CNN Pipeline

Slide credit: Ross Girshick

B. Leibe

# R-CNN Pipeline

Slide credit: Ross Girshick

B. Leibe

33

# Classification



| Input Image | Region Proposals | Feature Extraction | Classification |

- ## Linear model with class-dependent weights
  - ➢ Linear SVM

$$f_c(x_{fc7}) = w_c^T x_{fc7}$$

  - ➢ where
    - $x_{fc7}$ = features from the network (fully-connected layer 7)
    - $c$ = object class

Slide credit: Ross Girshick, Kaustav Kundu          B. Leibe

# Bounding Box Regressors

- Prediction of the 2D box

  - ➢ Necessary, since the proposal region might not fully coincide with the (annotated) object bounding box

  - ➢ Perform regression for location $(x^*, y^*)$, width $w^*$ and height $h^*$

$$\frac{x^* - x}{w} = w_{c,x}^T x_{pool5}$$

$$\frac{y^* - y}{h} = w_{c,y}^T x_{pool5}$$

$$\ln \frac{w^*}{w} = w_{c,w}^T x_{pool5}$$

$$\ln \frac{h^*}{h} = w_{c,w}^T x_{pool5}$$

  - ➢ Where $x_{pool5}$ are the features from the pool5 layer of the network.

Slide credit: Ross Girshick, Kaustav Kundu          B. Leibe

# Problems with R-CNN

- ## Ad hoc training objectives

  - ➢ Fine tune network with softmax classifier (log loss)

  - ➢ Train post-hoc linear SVMs (hinge loss)

  - ➢ Train post-hoc bounding-box regressors (squared loss)

- ## Training (3 days) and testing (47s per image) is slow.

  - ➢ Many separate applications of region CNNs

- ## Takes a lot of disk space

  - ➢ Need to store all precomputed CNN features for training the classifiers
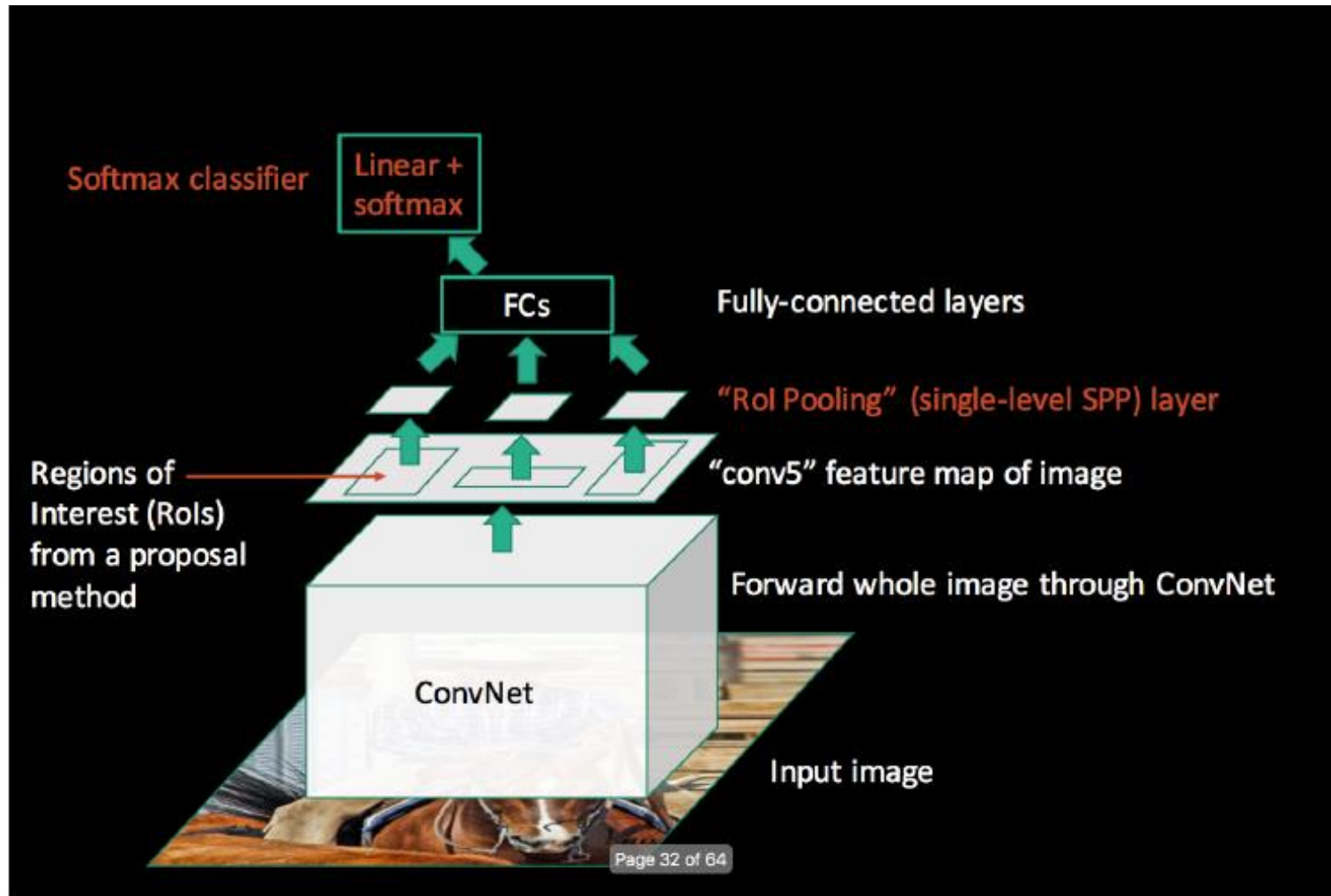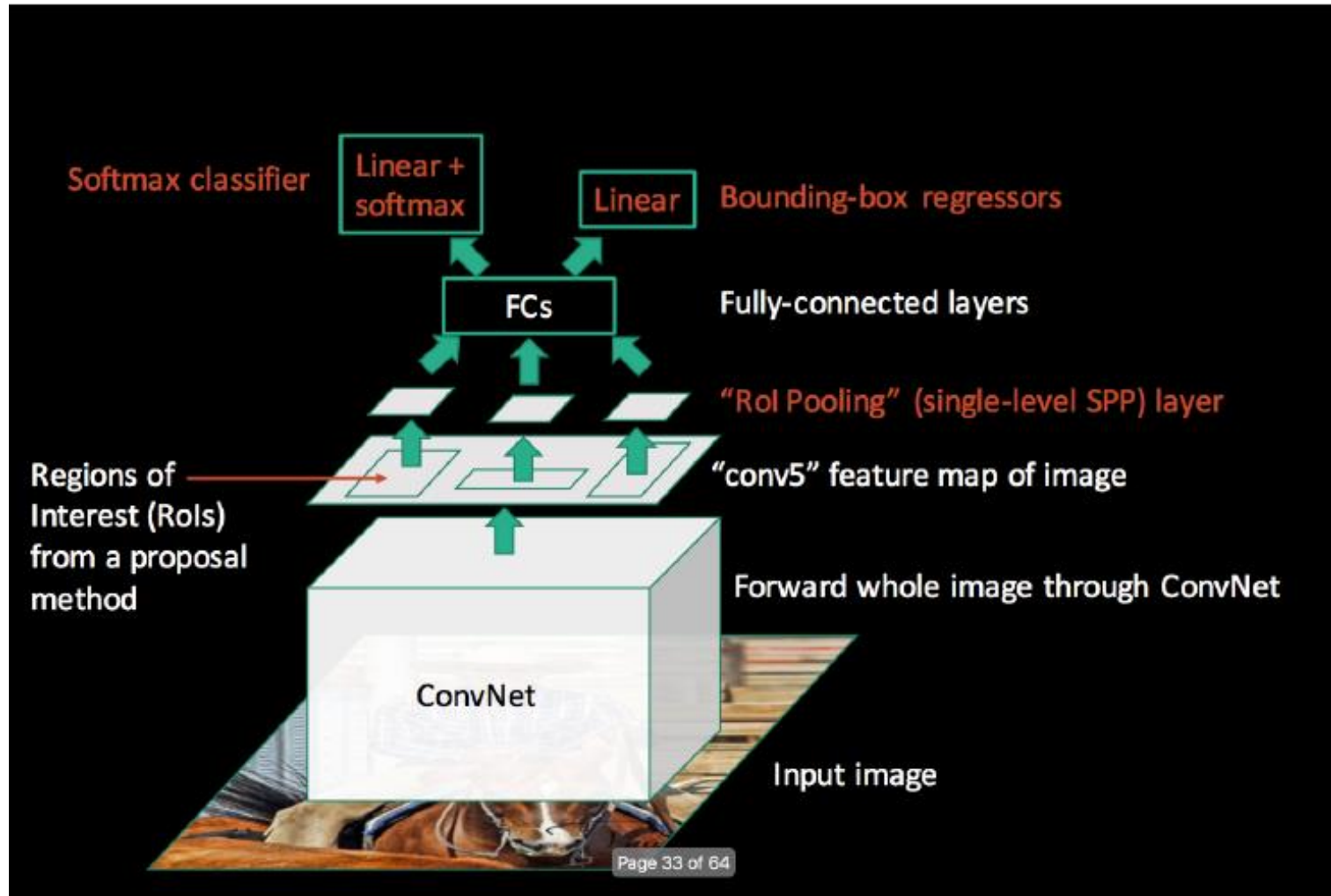
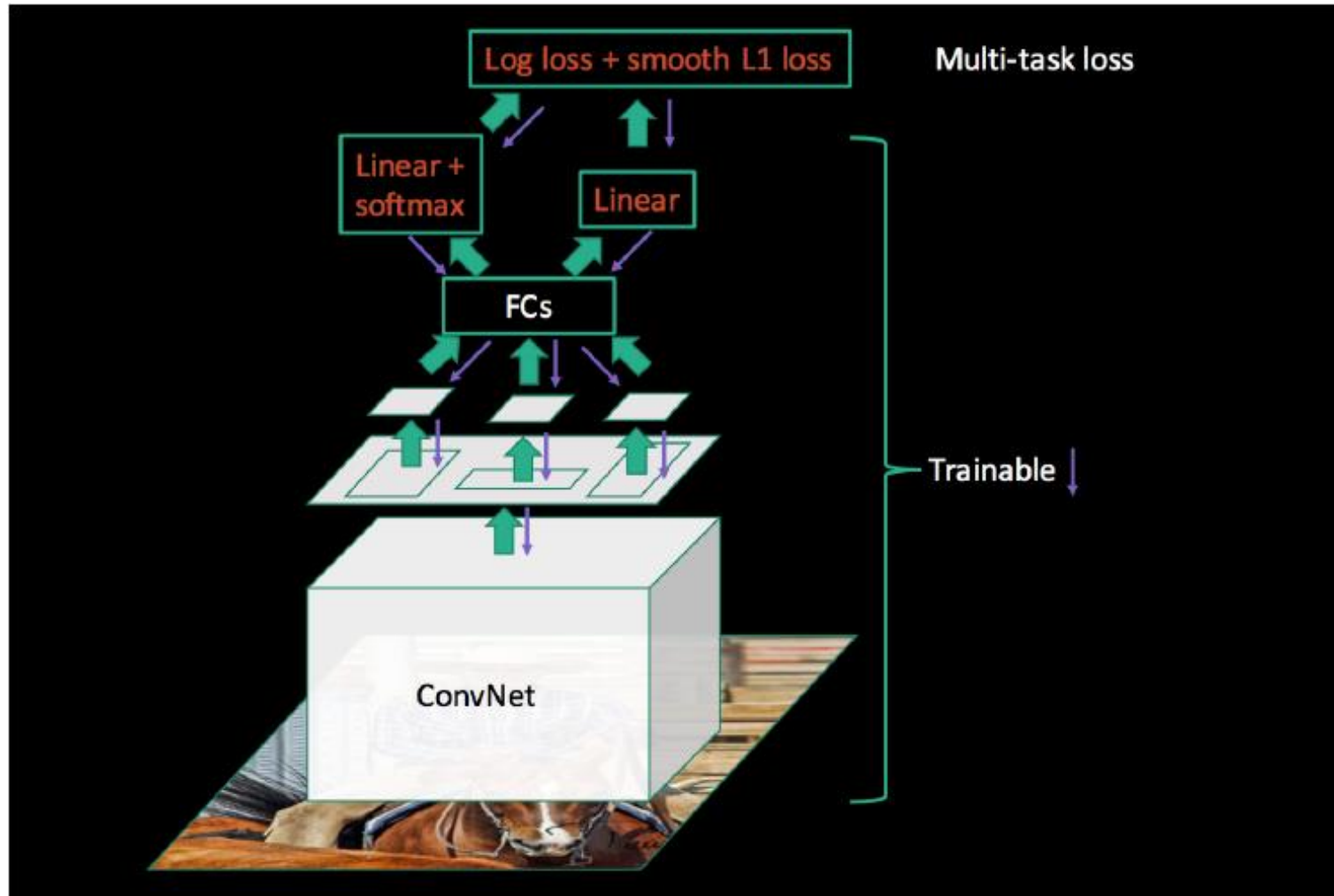  - ➢ Easily 200GB of data

B. Leibe

# Fast R-CNN

- Forward Pass

Slide credit: Ross Girshick

B. Leibe

# Fast R-CNN

- Forward Pass

Slide credit: Ross Girshick

B. Leibe

# Fast R-CNN

- Forward Pass

B. Leibe

# Fast R-CNN Training

- Backward Pass

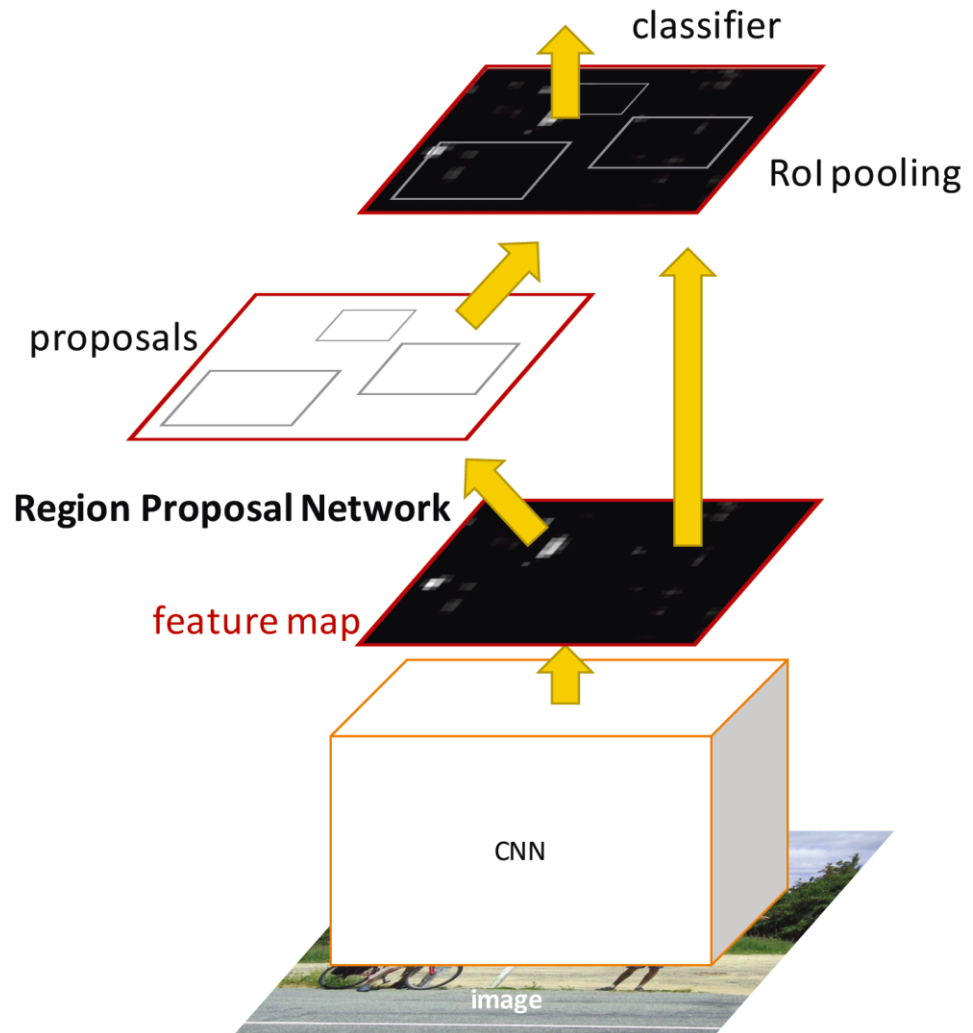Slide credit: Ross Girshick

B. Leibe

# Region Proposal Networks (RPN)

- Idea
  - Remove dependence on external region proposal algorithm.
  - Instead, infer region proposals from same CNN.
  - $\Rightarrow$ Feature sharing
  - $\Rightarrow$ Object detection in a single pass becomes possible.

- Faster R-CNN = Fast R-CNN + RPN



classifier

RoI pooling

proposals

**Region Proposal Network**

feature map

CNN

image

42

# Faster R-CNN

- One network, four losses
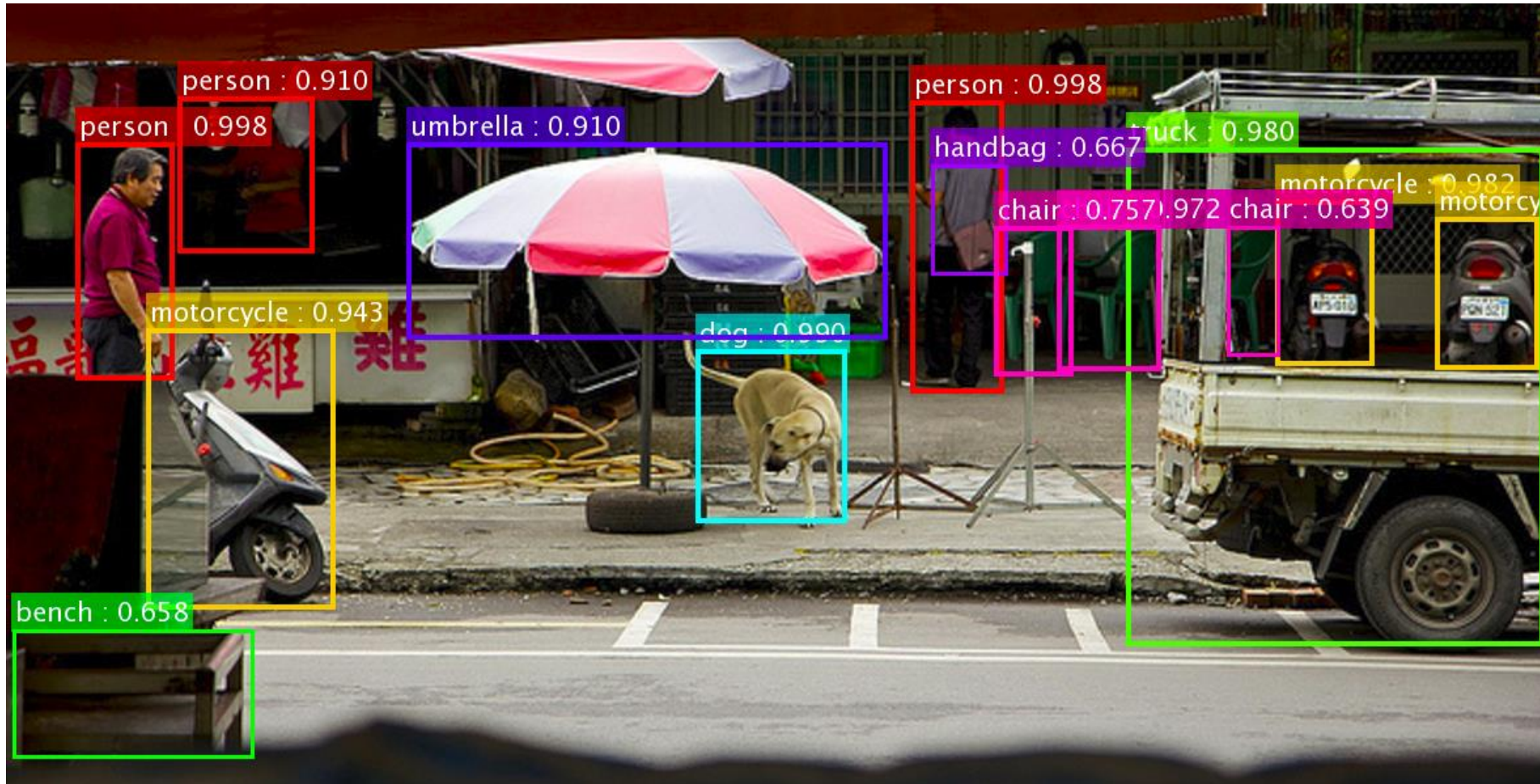  - ➢ Joint training

Slide credit: Ross Girshick
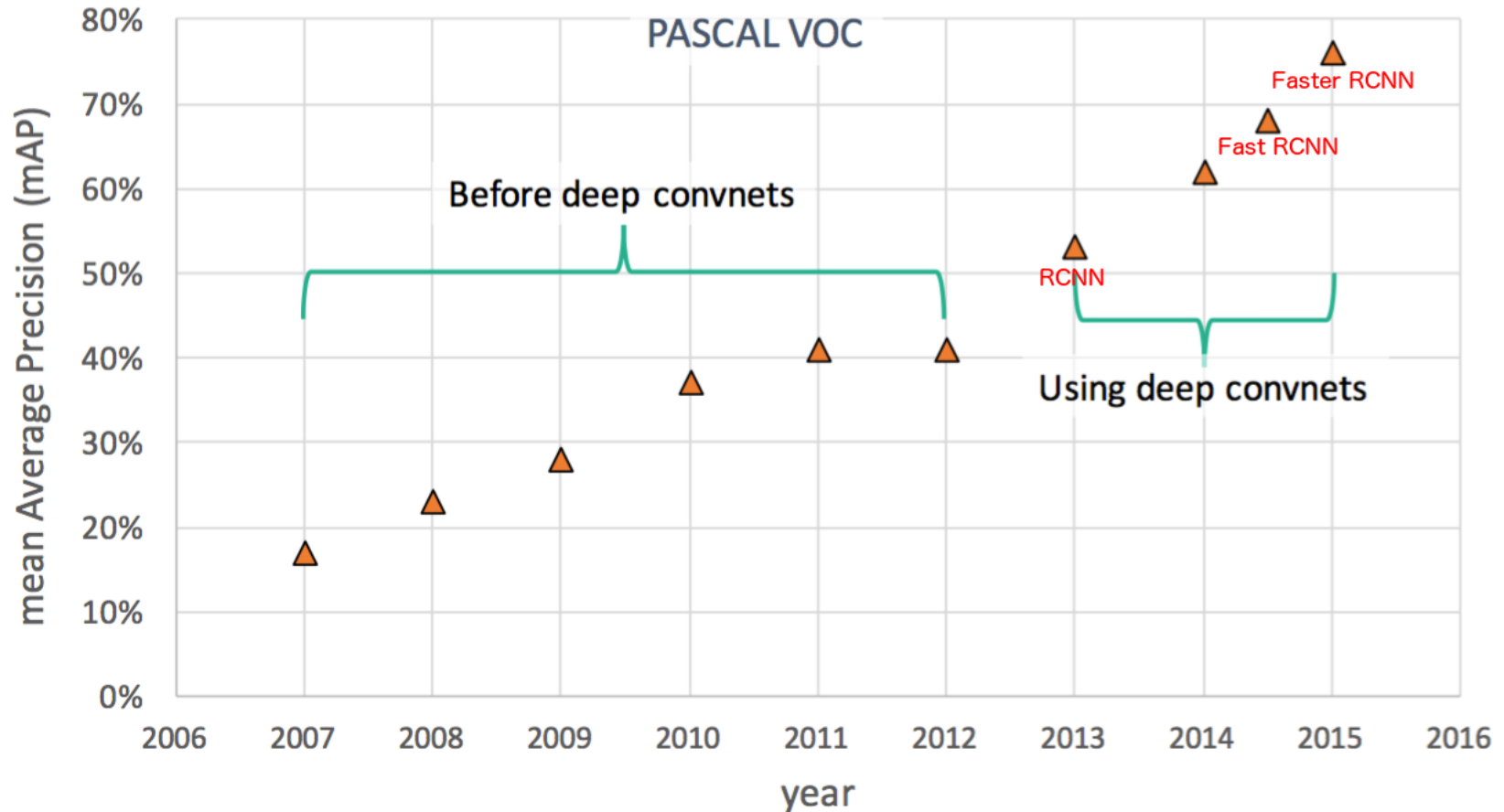
# Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.
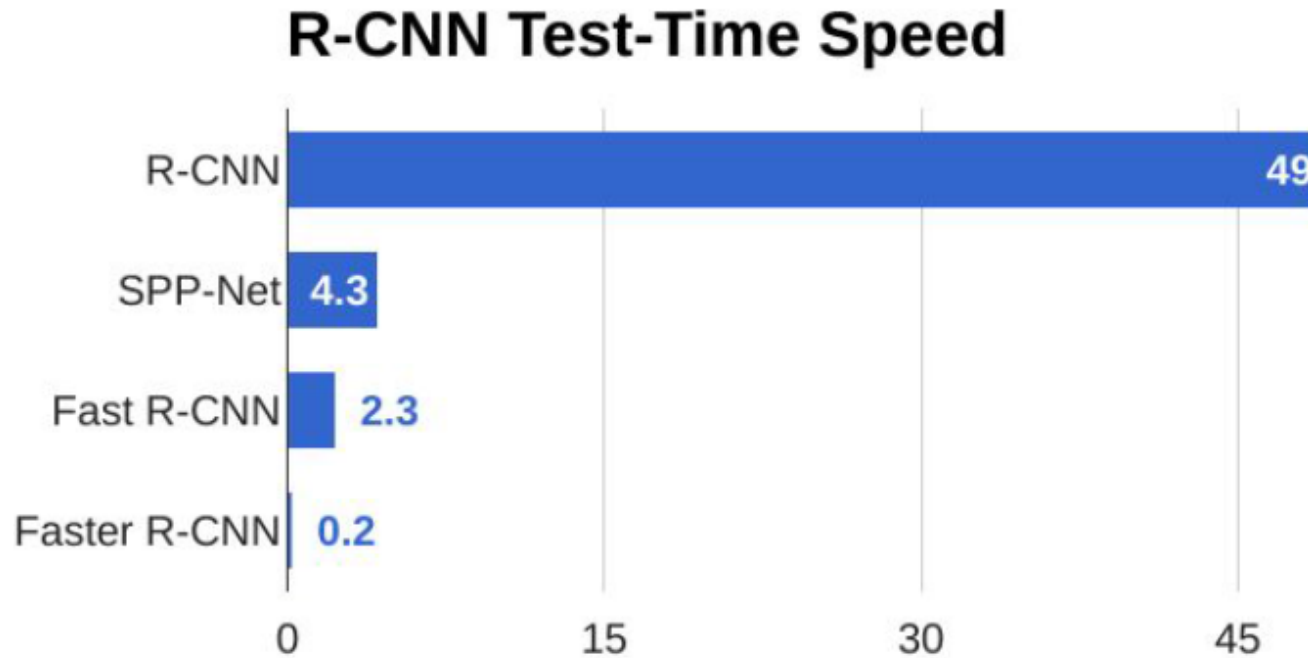
B. Leibe

# Faster R-CNN (based on ResNets)



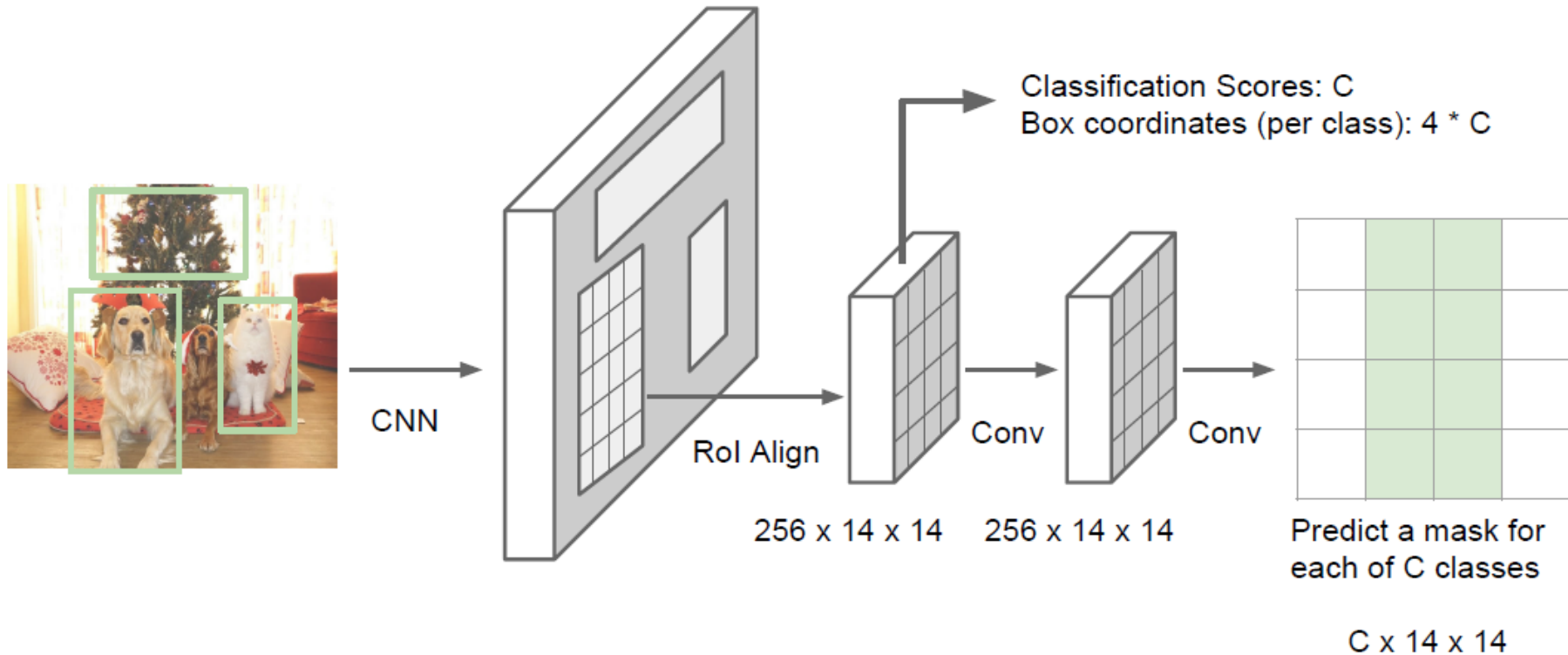K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.
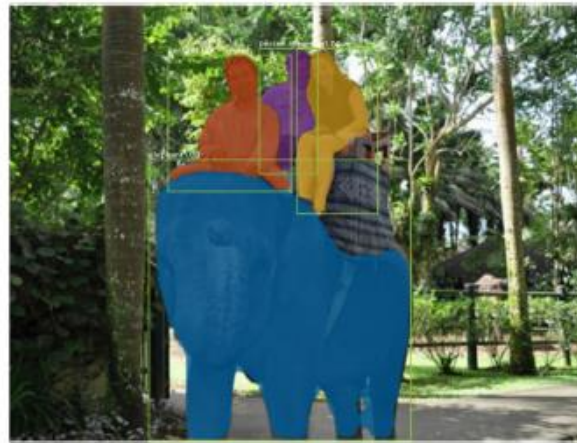
B. Leibe

# Object Detection Performance

Slide credit: Ross Girshick

B. Leibe

# Runtime Comparison



**R-CNN Test-Time Speed**

| | |
|---|---|
| R-CNN | 49 |
| SPP-Net | 4.3 |
| Fast R-CNN | 2.3 |
| Faster R-CNN | 0.2 |

Slide credit: FeiFei Li

B. Leibe

Computer Vision Summer'19

# Most Recent Version: Mask R-CNN



K. He, G. Gkioxari, P. Dollar, R. Girshick, Mask R-CNN, arXiv 1703.06870.

Computer Vision Summer'19

# Mask R-CNN Results

- Detection + Instance segmentation
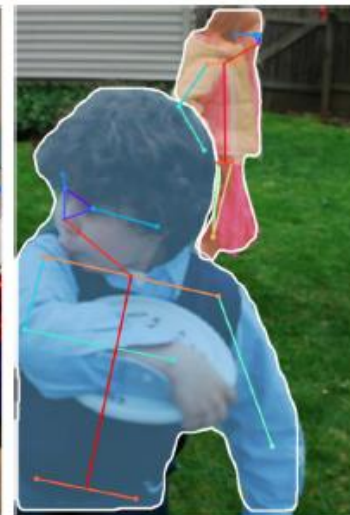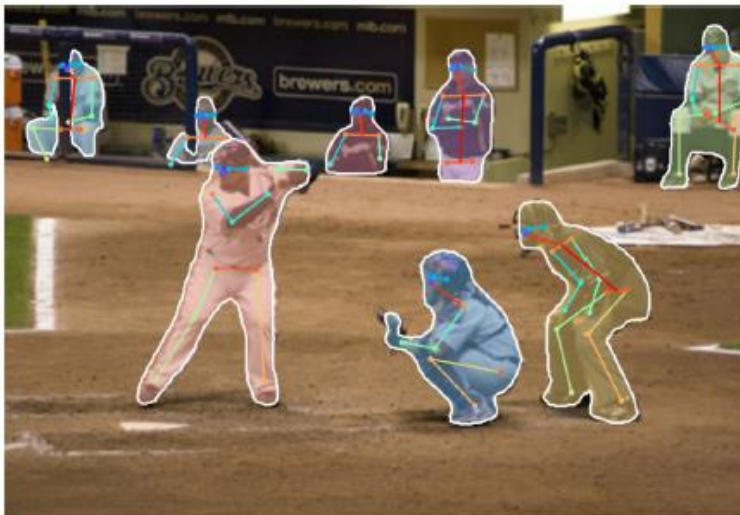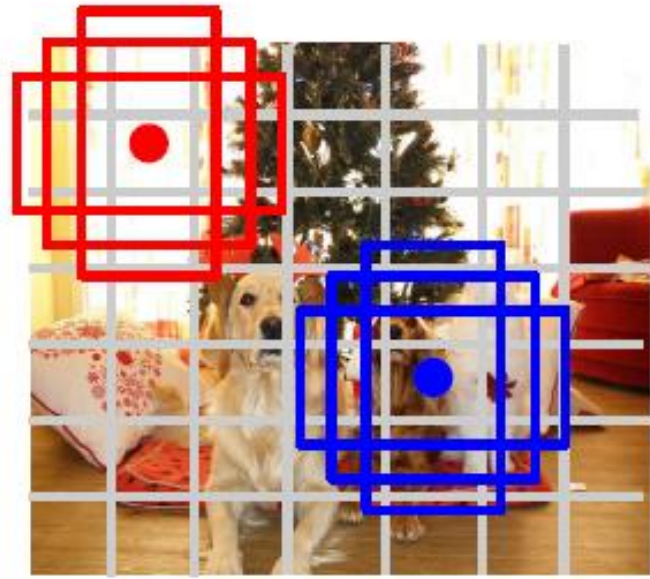


- Detection + Pose estimation



Figure credit: K. He, G. Gkioxari, P. Dollar, R. Girshick
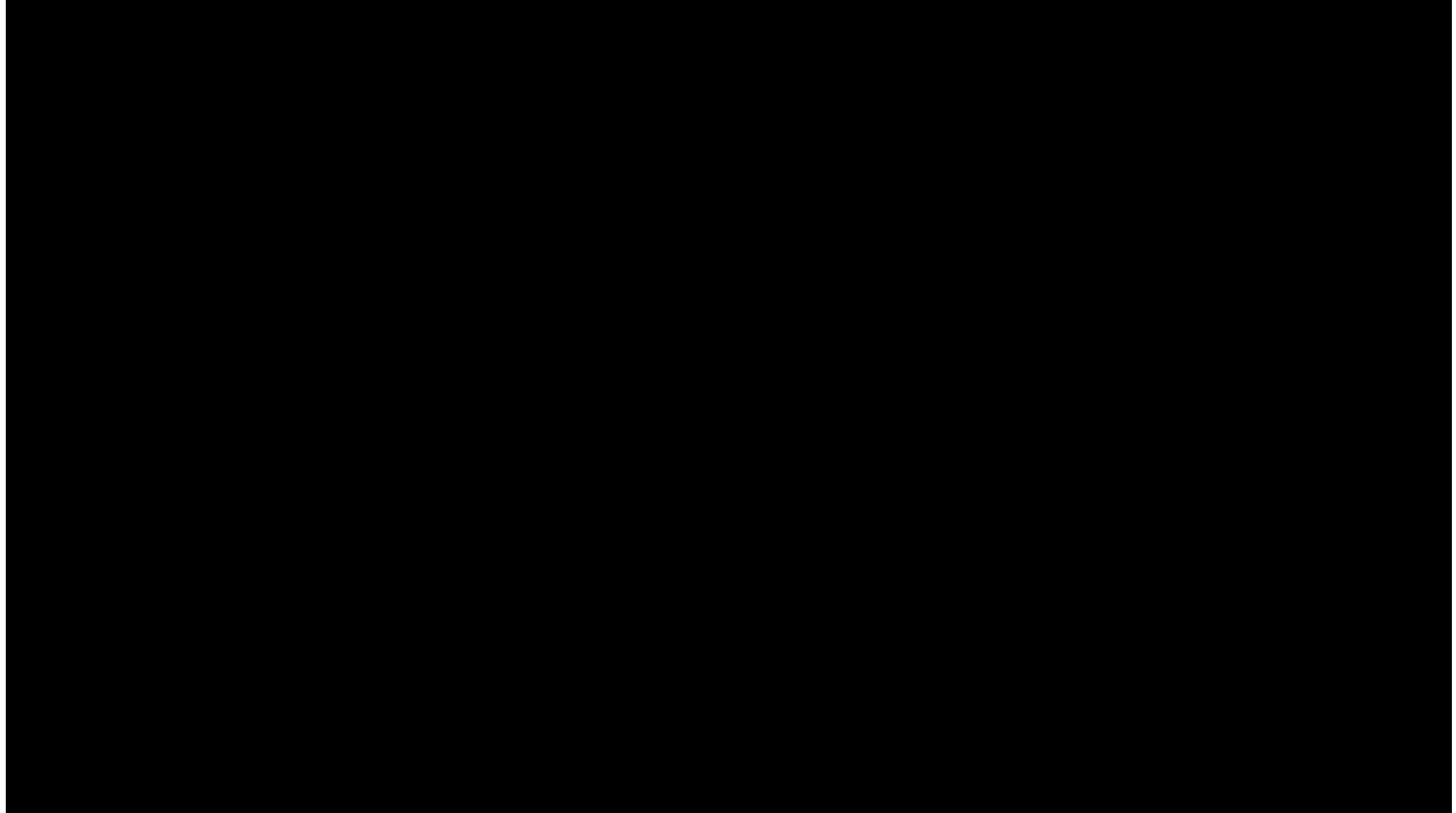
# YOLO / SSD



Input image
3 x H x W

Divide image into grid
7 x 7

- Idea: Directly go from image to detection scores
- Within each grid cell
  - Start from a set of anchor boxes
  - Regress from each of the B anchor boxes to a final box
  - Predict scores for each of C classes (including background)

50

# YOLO-v3 Results

J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016.

# Summary

- ## Object Detection
  - ➢ Find a variable number of objects by classifying image regions
  - ➢ Before CNNs: dense multiscale sliding window (HoG, DPM)

- ## Region proposal based detectors
  - ➢ Idea: Avoid dense sliding window with region proposals
  - ➢ R-CNN: Selective Search + CNN classification / regression
  - ➢ Fast R-CNN: Swap order of convolutions and region extraction
  - ➢ Faster R-CNN: Compute region proposals within the network
  - ➢ Mask R-CNN: Detection + instance segmentation + pose estimation

- ## Anchor box based detectors
  - ➢ Idea: Perform detection in a single step using grid of anchor boxes
  - ➢ YOLO, YOLO-v2, YOLO-v3
  - ➢ SSD

# References and Further Reading

- **LeNet**
  - Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86(11): 2278–2324, 1998.

- **AlexNet**
  - A. Krizhevsky, I. Sutskever, and G. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012.

- **VGGNet**
  - K. Simonyan, A. Zisserman, Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015

- **GoogLeNet**
  - C. Szegedy, W. Liu, Y. Jia, et al, Going Deeper with Convolutions, arXiv:1409.4842, 2014.

B. Leibe

# References and Further Reading

- ResNet
    - K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, CVPR 2016.

# References: Computer Vision Tasks

- Object Detection

  - R. Girshick, J. Donahue, T. Darrell, J. Malik, <u>Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation</u>, CVPR 2014.

  - S. Ren, K. He, R. Girshick, J. Sun, <u>Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks</u>, NIPS 2015.

  - K. He, G. Gkioxari, P. Dollar, R. Girshick, <u>Mask R-CNN</u>, ICCV 2017.

  - J. Redmon, S. Divvala, R. Girshick, A. Farhadi, <u>You Only Look Once: Unified, Real-Time Object Detection</u>, CVPR 2016

  - W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C-Y. Fu, A.C. Berg, <u>SSD: Single Shot Multi Box Detector</u>, ECCV 2016.

B. Leibe

Computer Vision Summer'19