

---

# Image Classification and Object detection

---

Øyvind Meinich-Bache  
Associate Professor II

# Classification vs detection

---

Image classification



«car»

Classification with localization



«car»

Object detection



Multiple objects

1 object

# Image Classification

---



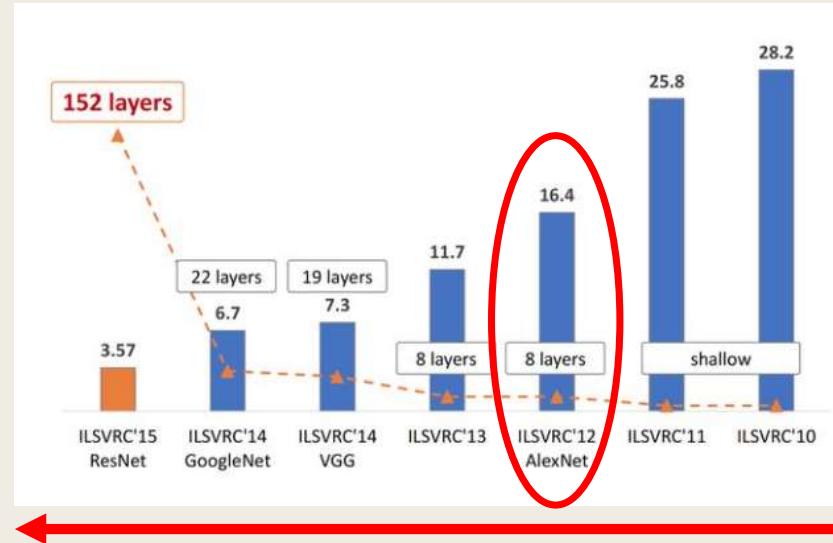
All images in the dataset are only assigned to **one** class – «car»

# Image classification history

## ImageNet Large Scale Visual Recognition Challenge

2012 winner- AlexNet

- First «deep» model
- Utilized GPU



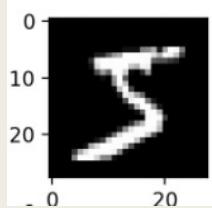
Today: State-of-the-art - Deep convolutional Neural Networks



# Image Classification with MLP

---

0	1	2	3	4	5	6	7	8	9
0	1	<b>2</b>	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



MNIST dataset: <http://yann.lecun.com/exdb/mnist/>

# Image Classification with CNN

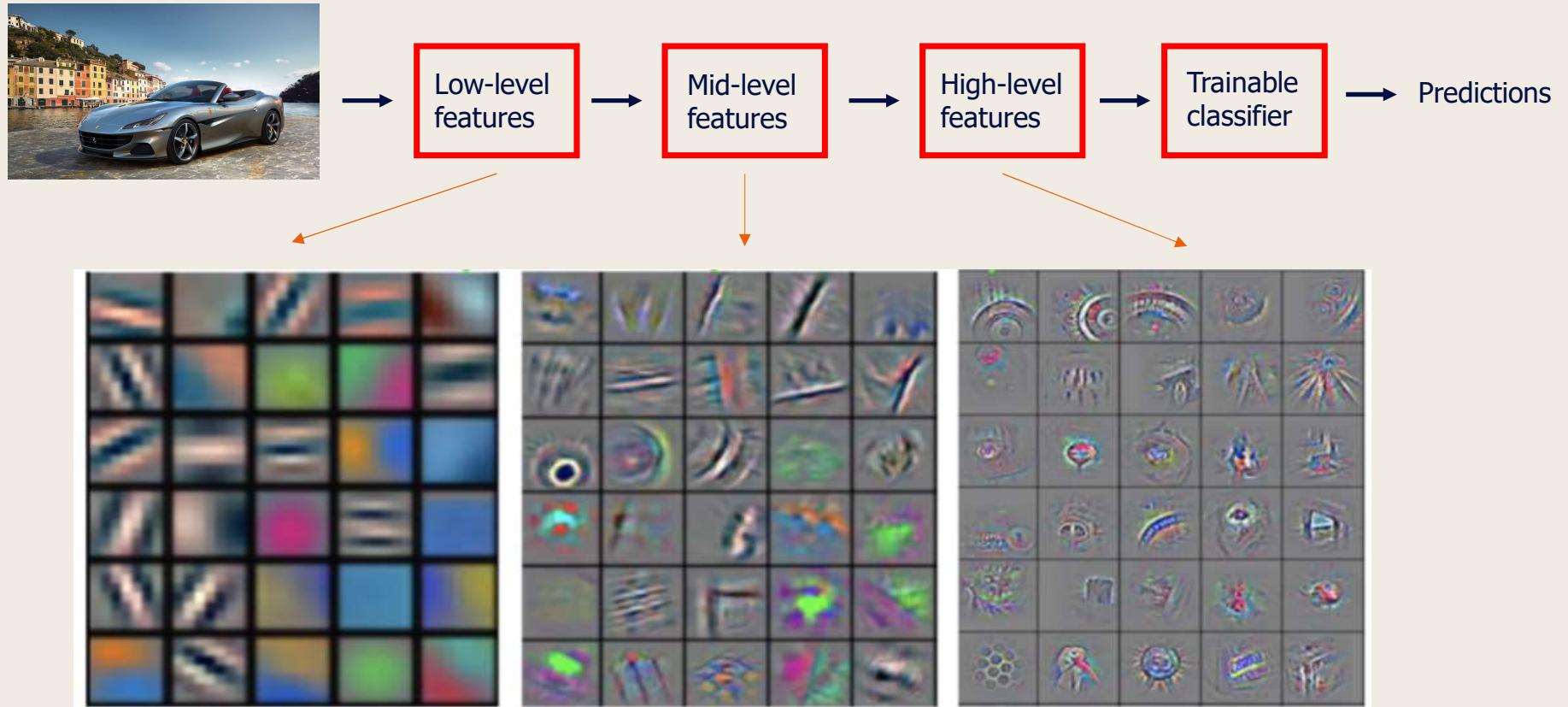
---



CNNs:

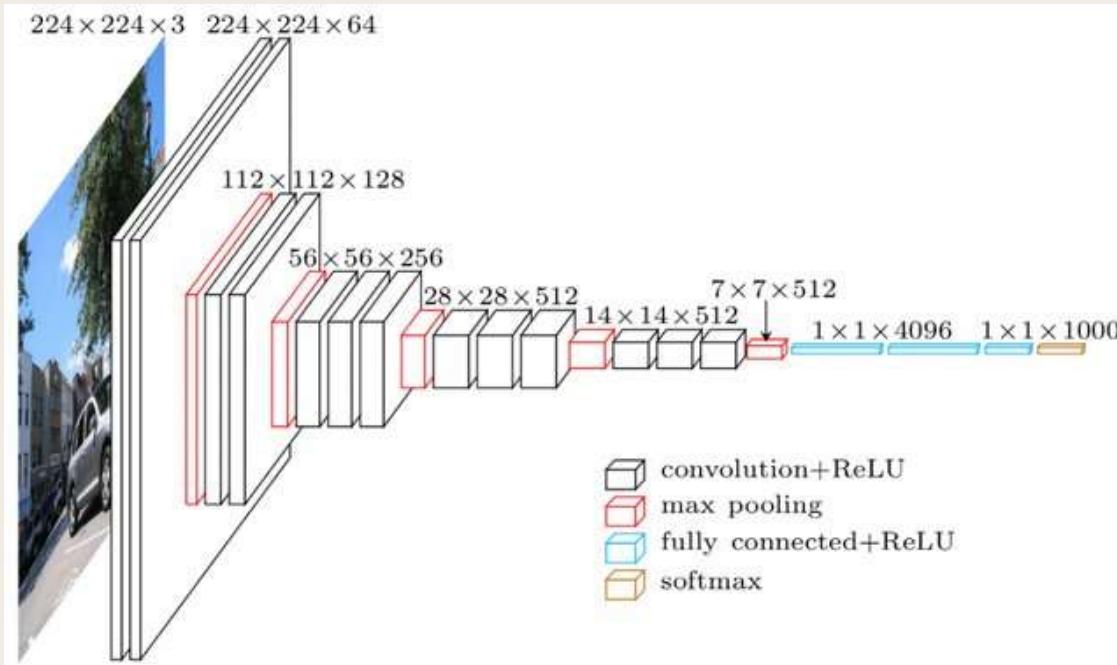
- Good at data with spatial information – filters learn spatial features
- Parameter/weight sharing  
-> efficient

# Image Classification with CNN



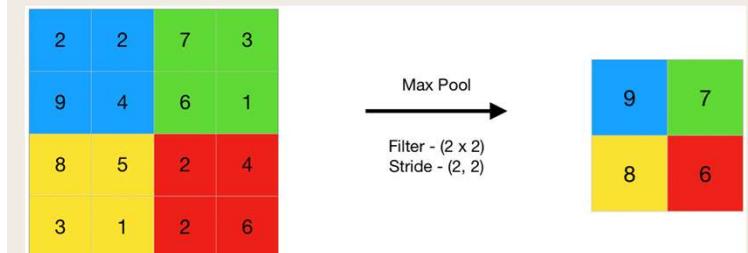
# Image Classification – network structures

## VGG-16 Network



Uses smaller receptive fields in early layers,  $3 \times 3$  with stride of 1 compared to AlexNet which use  $11 \times 11$  and stride of 4.

Winner of ImageNet 2014.

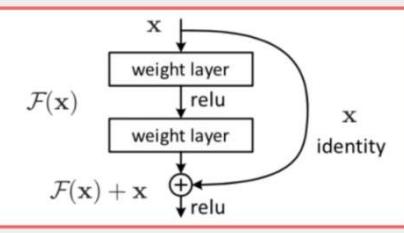
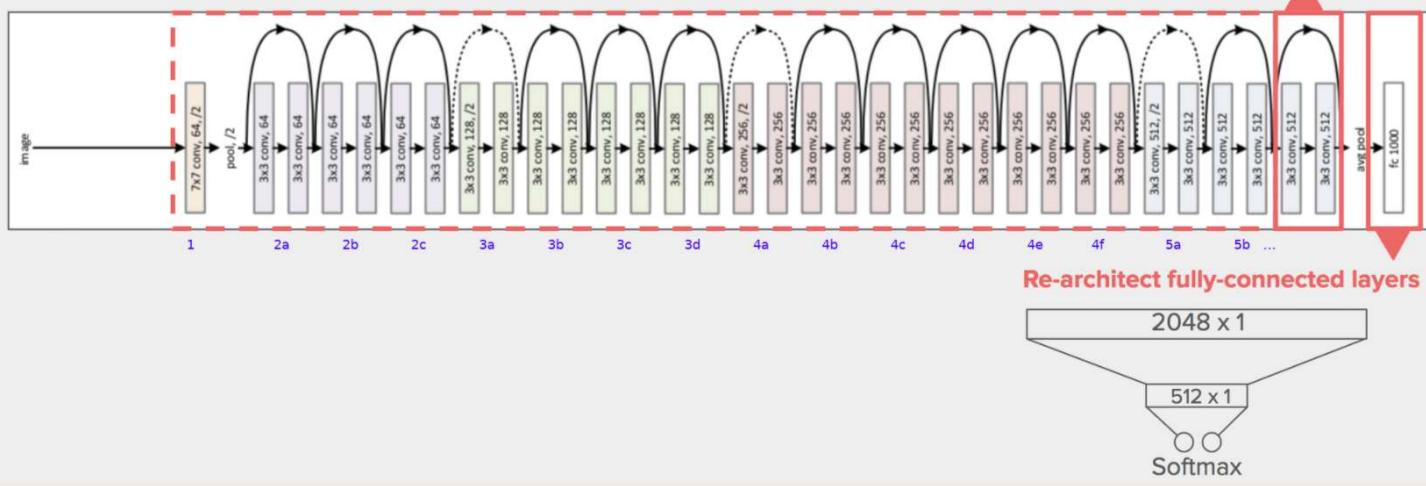


Source: <https://towardsdatascience.com/step-by-step-vgg16-implementation-in-keras-for-beginners-a833c686ae6c>

# Image Classification – network structures

## Retrain ResNet50

ResNet50 Diagram

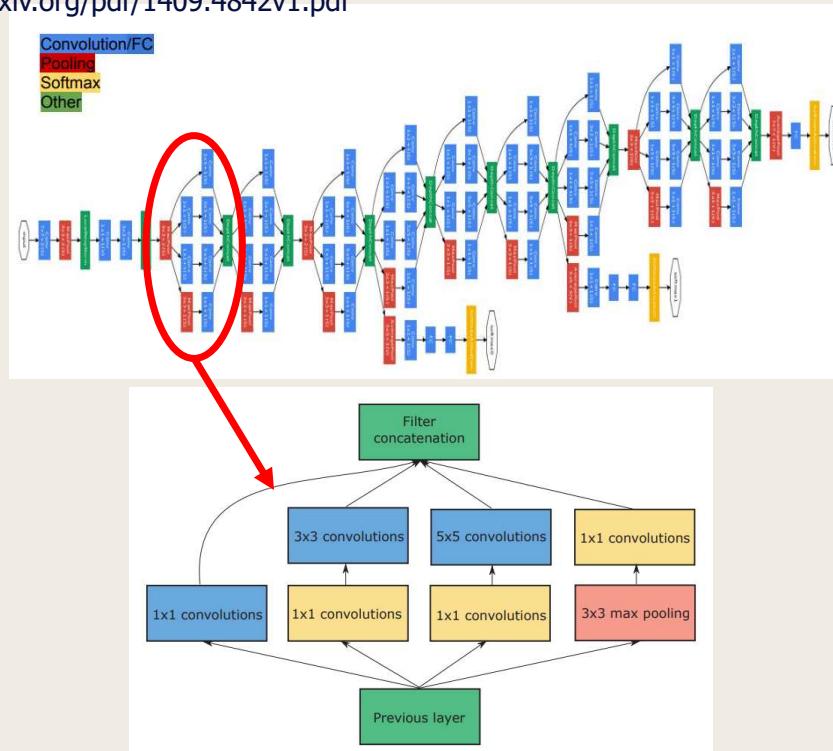


- ResNet50 have 50 layers
- Often used as starting point for transfer learning for image classification problems
- Also a 152 layer version of ResNet (very deep CNN) - ImageNet winner of 2015
- Handle vanishing gradient problem during training using the «skip connection architecture»

# Image Classification – network structures

## Inception Network

Source: <https://arxiv.org/pdf/1409.4842v1.pdf>



Too deep CNNs:

- vanishing/exploding gradients
- Computational expensive

Inception:

- Tackle these problems by letting filters of different size operate on the same level
- More «wider» and less «deeper»

# Image Classification – network structures

---

VGG: <https://arxiv.org/pdf/1409.1556.pdf>

Resnet: <https://arxiv.org/pdf/1512.03385.pdf>

Inception: <https://arxiv.org/pdf/1409.4842.pdf>

Optional material:

ViT Networks:

Scaling vision transformers: <https://arxiv.org/pdf/2106.04560v1.pdf>

Image Classification on ImageNet:

<https://paperswithcode.com/sota/image-classification-on-imagenet>

# Image Classification with localization



- Classes:
- 1 - Car
  - 2 - Building
  - 3 – Person
  - 4 - Background

→



CNN

Output layer  
Softmax

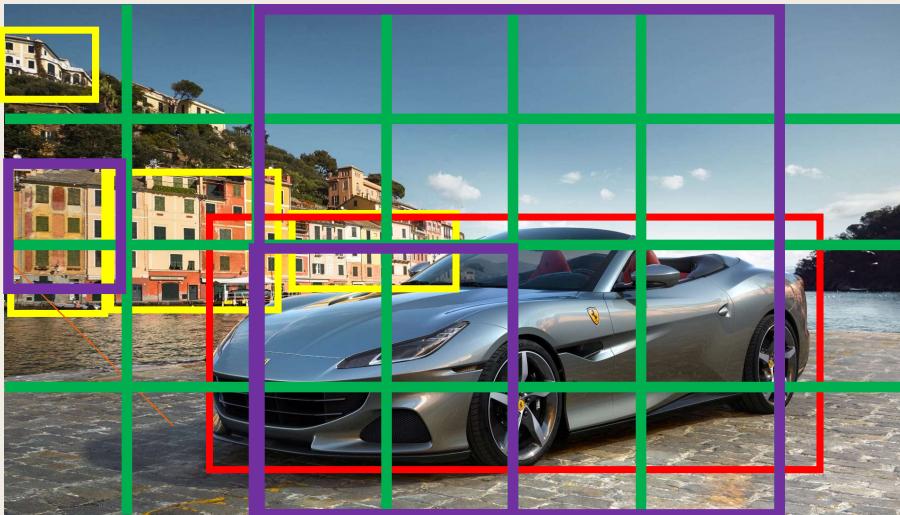


$pc$   
 $bx$   
 $by$   
 $bh$   
 $bw$   
 $c1,$   
 $1$   
 $c2$   
 $c3$   
 $c4$   
 $0$   
 $0$   
 $0$

$\begin{bmatrix} 1 \\ bx \\ by \\ bh \\ bw \\ c1, \\ 1 \\ 0 \\ c2 \\ c3 \\ 0 \\ c4 \\ 0 \end{bmatrix}$

= Y

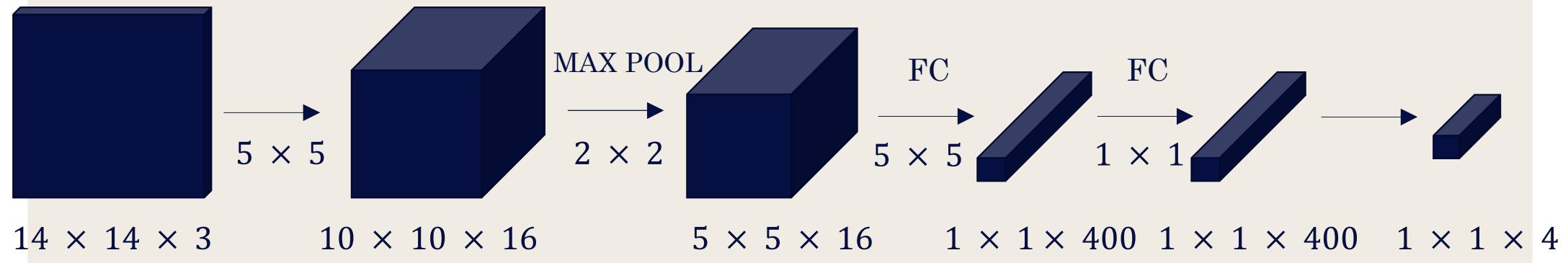
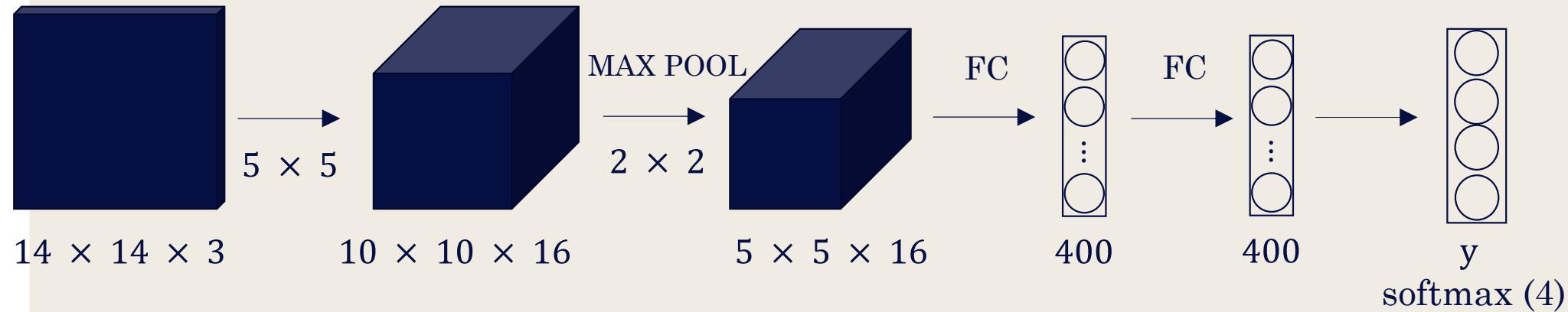
# Object Detection – sliding window



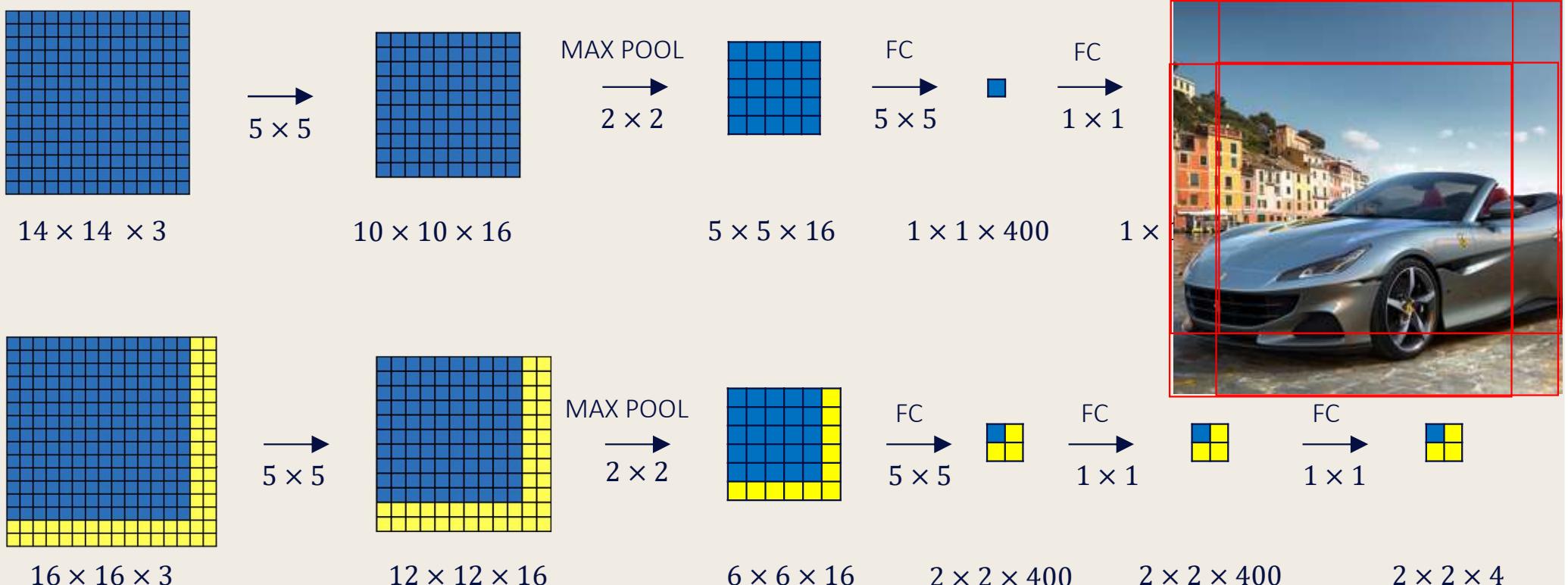
Grid split – image classification in each region

- Objects can cover more than one grid region
    - Need to do a «sliding window» approach
    - Would also need to analyse different sized grid regions
- > not an efficient approach!

# Convolution implementation of sliding windows



## Object Detection – Convolution implementation of sliding window



# Object Detection

---



«car»  
«car»  
«car»  
«car»

# Object Detection

---



«car»

«car»

→  $P_c = 0.7$

→  $P_c = 0.9$

$P_c$  = Prediction probability

If multiple bounding boxes of same object – keep only the one with the highest prediction probability

How?

# Object Detection - IoU

---



Intersection over union (IoU)

- measurement of overlap between two bounding boxes

$$\text{IoU} = \frac{\text{Size of intersection}}{\text{Size of union}}$$

# Object Detection - Non-maximum suppression



«car»

«car»

→  $P_c = 0.8$

→  $P_c = 0.9$

1: discard all boxes with,  $P_c < \text{Threshold}$

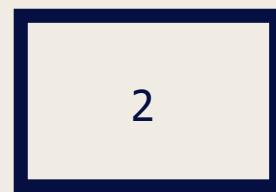
If there are remaining boxes:

2: Pick the box with the largest  $P_c$  and output that as a prediction

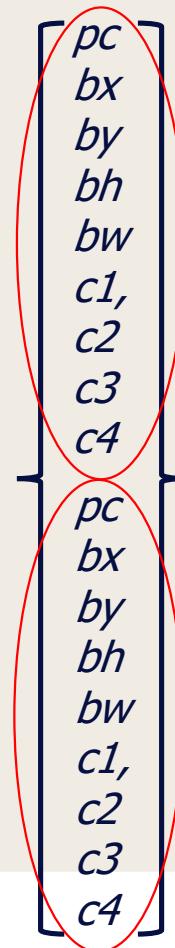
3: discard any remaining boxes with  $\text{IoU} > 0.5$

# Object Detection – Anchor boxes

Region vector in output tensor



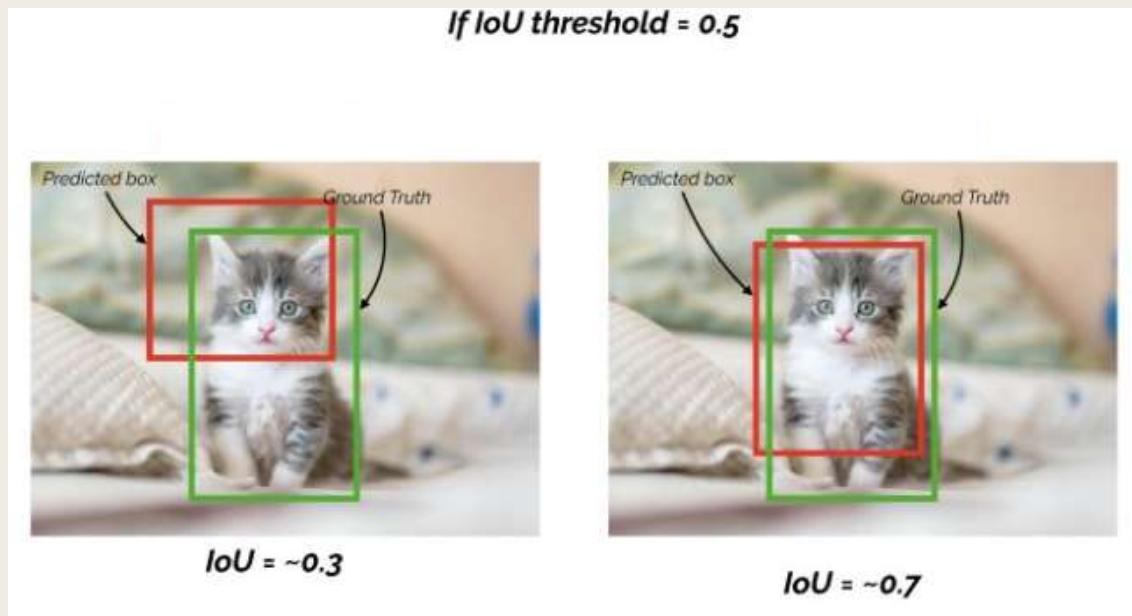
$Y =$



Make it possible to predict multiple object that has centerpoints in the same region

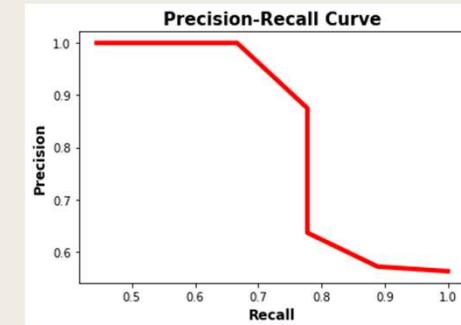
# Object Detection – mean Average Precision (mAP)

Object detector performance metric. Estimates an Average Precision for each class and then average over all classes to get the mAP

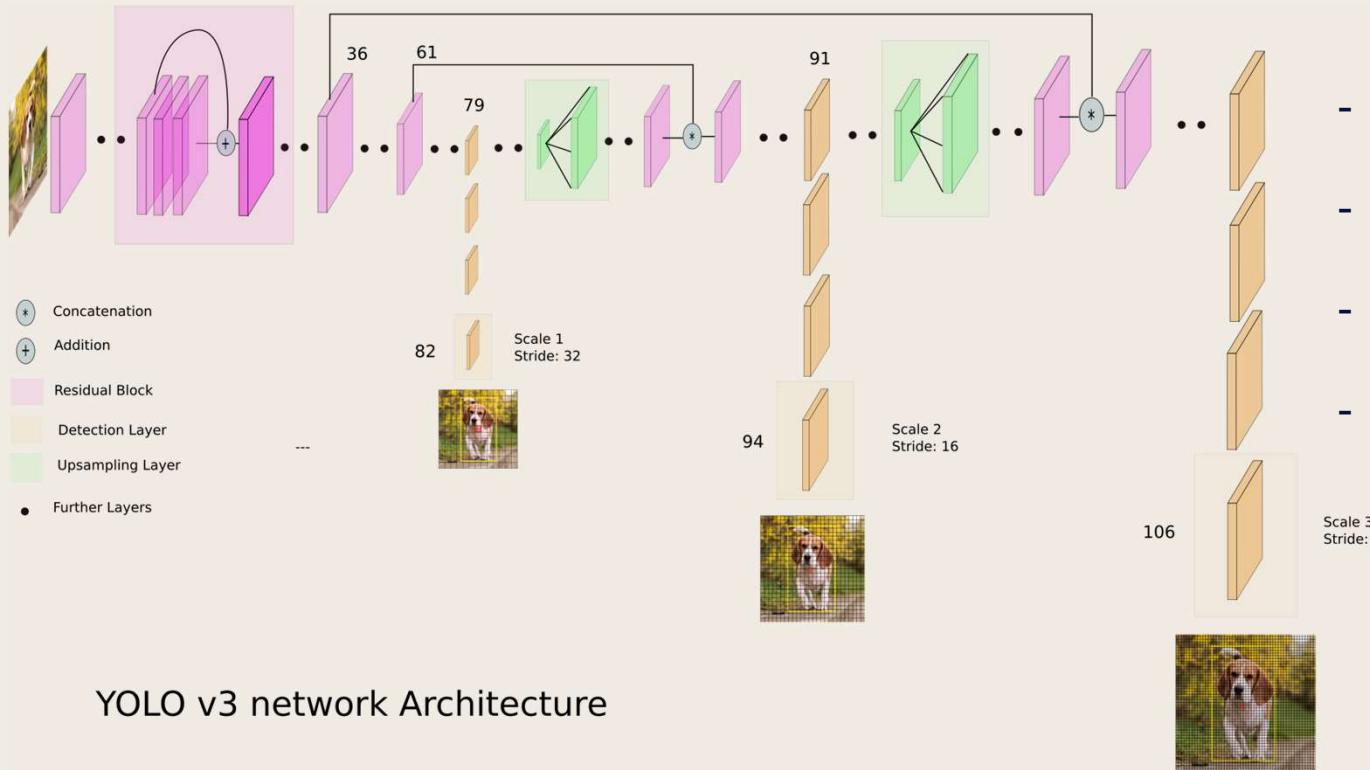


$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$



# Object Detection – YOLOv3



- 9 anchor boxes, 3 for each scale
- Scales - 13x13, 26x26 and 52x52
- -> Output 10 647 predicted boxes
- Non-maximum suppression

One-stage  
object detector

Source: <https://towardsdatascience.com/yolo-v3-object-detection-53fb7d3bfe6b>

# Object Detection – one-stage detectors

- YOLO
  - v3 : <https://arxiv.org/pdf/1804.02767.pdf>
  - v4 : <https://arxiv.org/pdf/2004.10934.pdf>
  - PP-YOLO : <https://arxiv.org/pdf/2007.12099.pdf>
- RetinaNet
  - <https://arxiv.org/pdf/1708.02002.pdf>



# Object Detection – two-stage detectors

Two-stage object detectors:

- First step: region proposal,  
Second step: classification based on feature extraction from regions
- Normally not end-to end trainable – trained in two steps
- Highest accuracy, but slower. Performance (fps) drops.



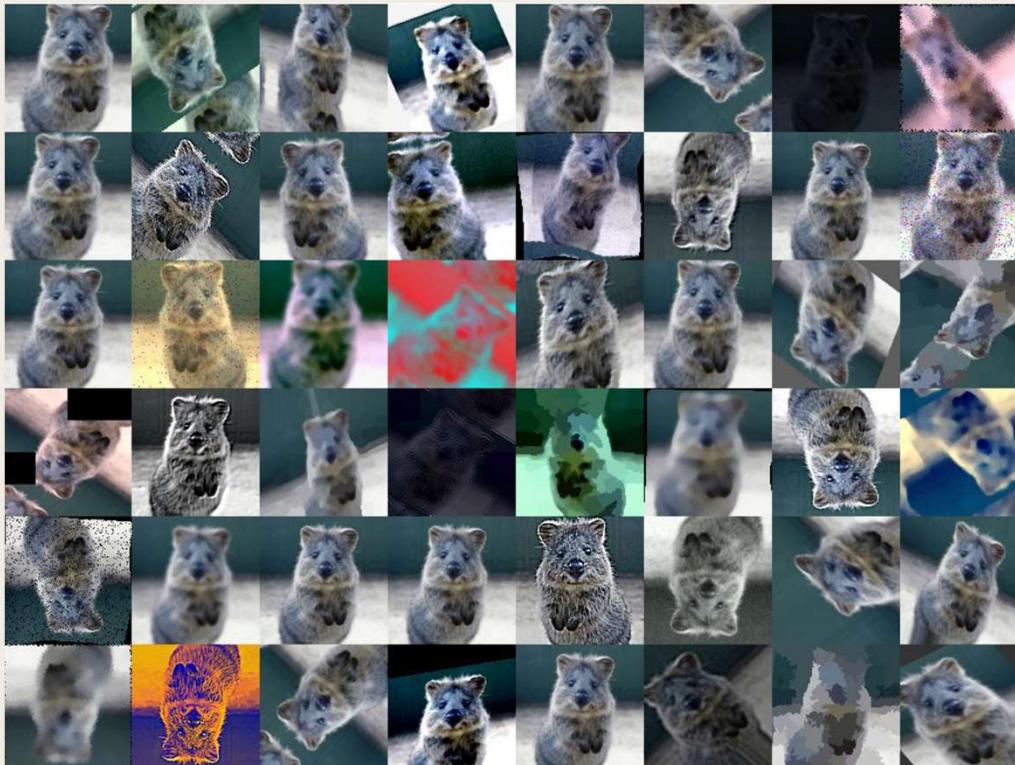
# Object Detection - two-stage detectors

---

- Popular 2-stage object detectors:
  - Region-Based Convolutional Neural Networks (R-CNN) –
    - <https://arxiv.org/pdf/1311.2524v5.pdf>
  - Faster R-CNN,
    - <https://arxiv.org/pdf/1506.01497.pdf>
  - Mask R-CNN
    - <https://arxiv.org/pdf/1703.06870.pdf>

# Data Augmentation

---



- Increase variations and dataset size
- Rotation, cropping, blurring, noise, color adjustment etc

Source: <https://github.com/aleju/imgaug>

# Training, Validation, Test data

## Car detector

- Need to come from the same distributions!



Training/val data

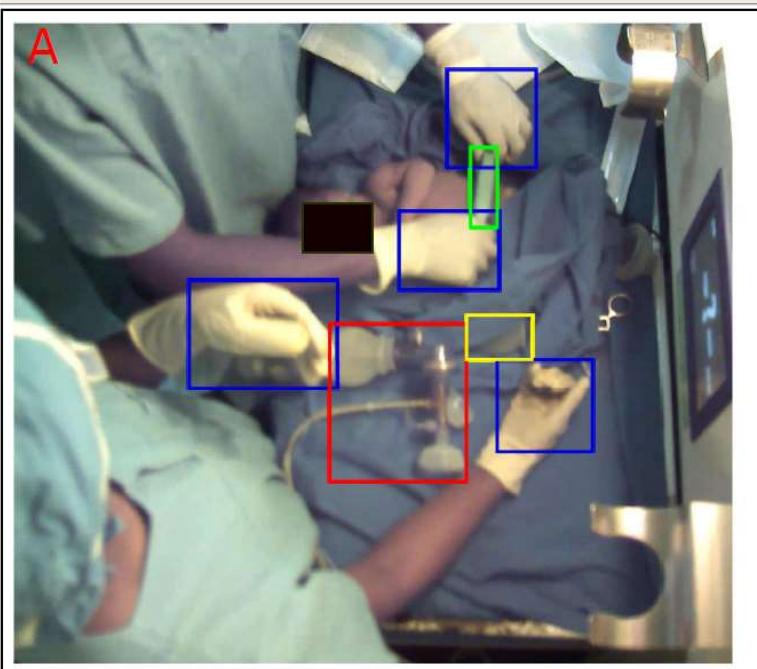


Test data



# Object detection in Newborn Resuscitation Videos

---



Source: <https://pubmed.ncbi.nlm.nih.gov/31247581/>