

# Deep machine learning for Computer Vision

Andrii Liubonko  
Samsung R&D Institute Ukraine

# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

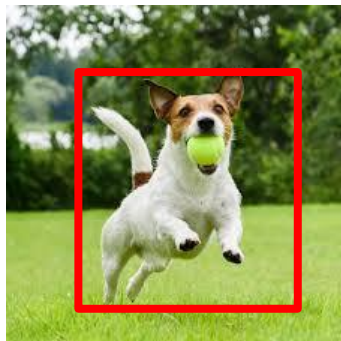
- Faster R-CNN

Single Shot Detector [SSD]

Comparison

Summary

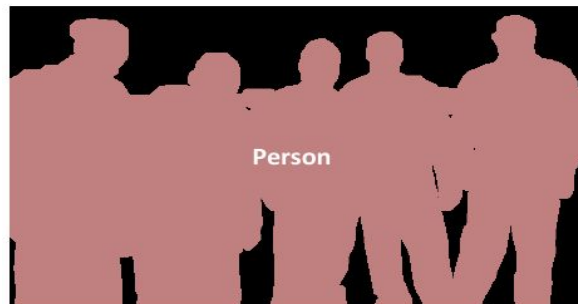
# Visual Perception Problems



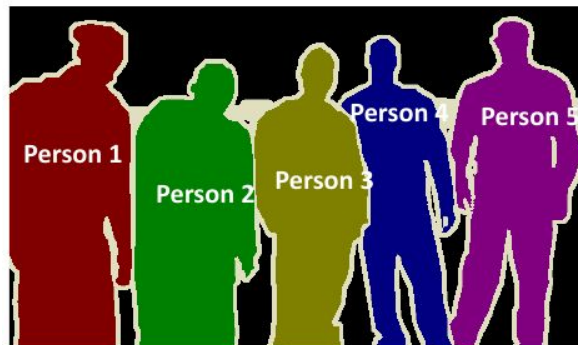
Classification + Localization



Object Detection

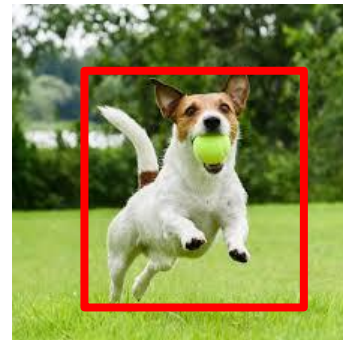


Semantic Segmentation



Instance Segmentation

# Localization

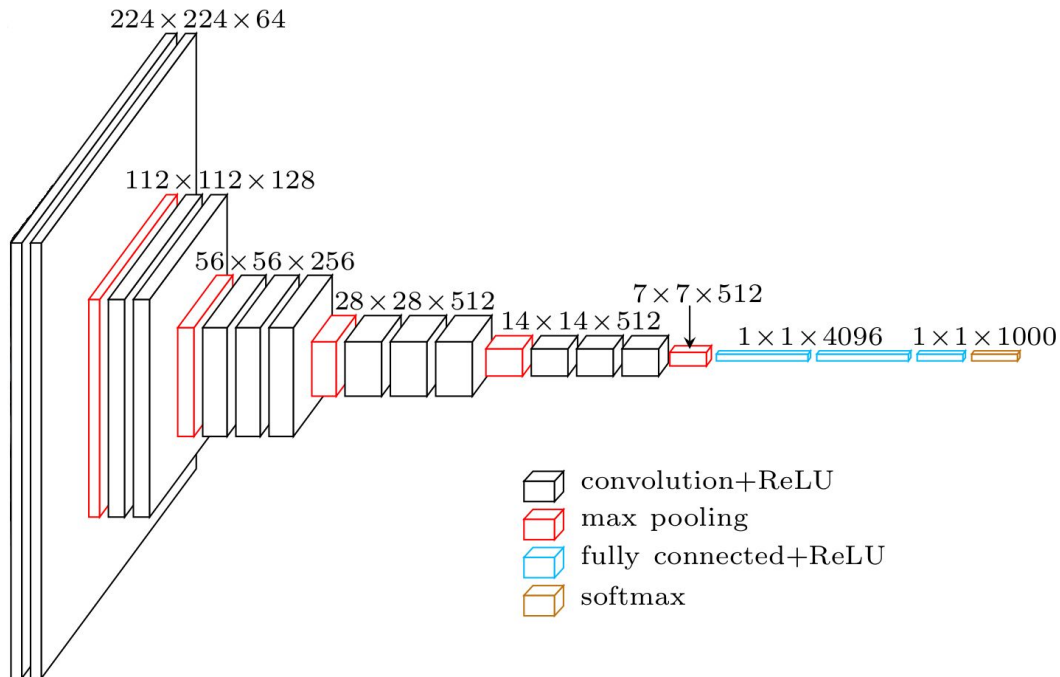


[Score] dog: 0.98  
cow: 0.01

...  
[Bounding Box]  
(x,y,h,w)

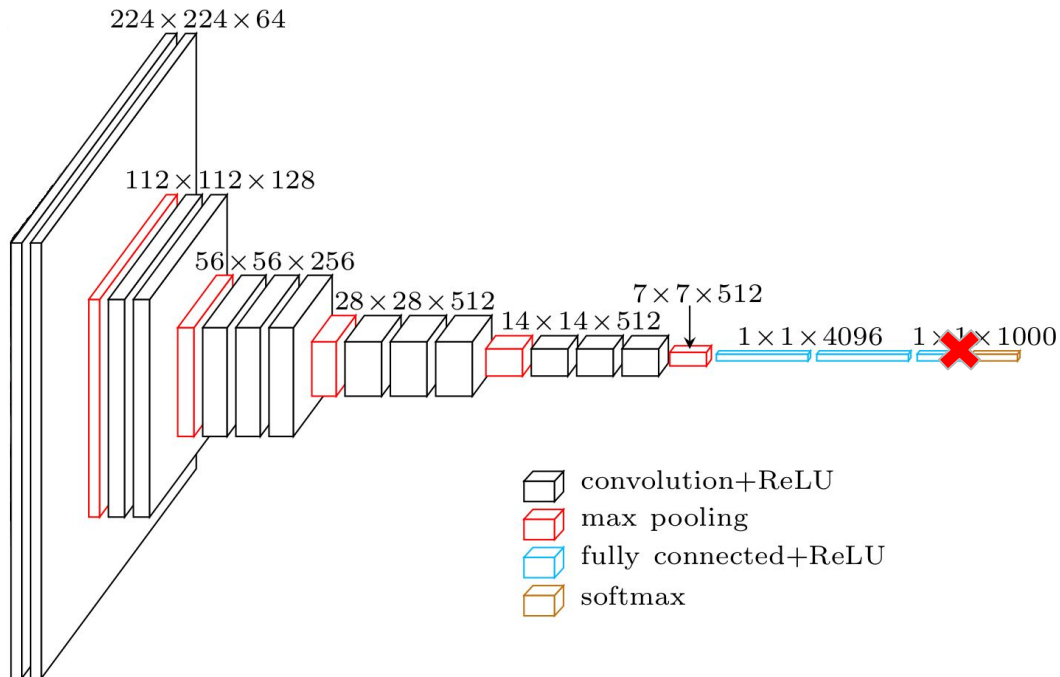
The main assumption is that there is only **one** object in the image

# Image Classification



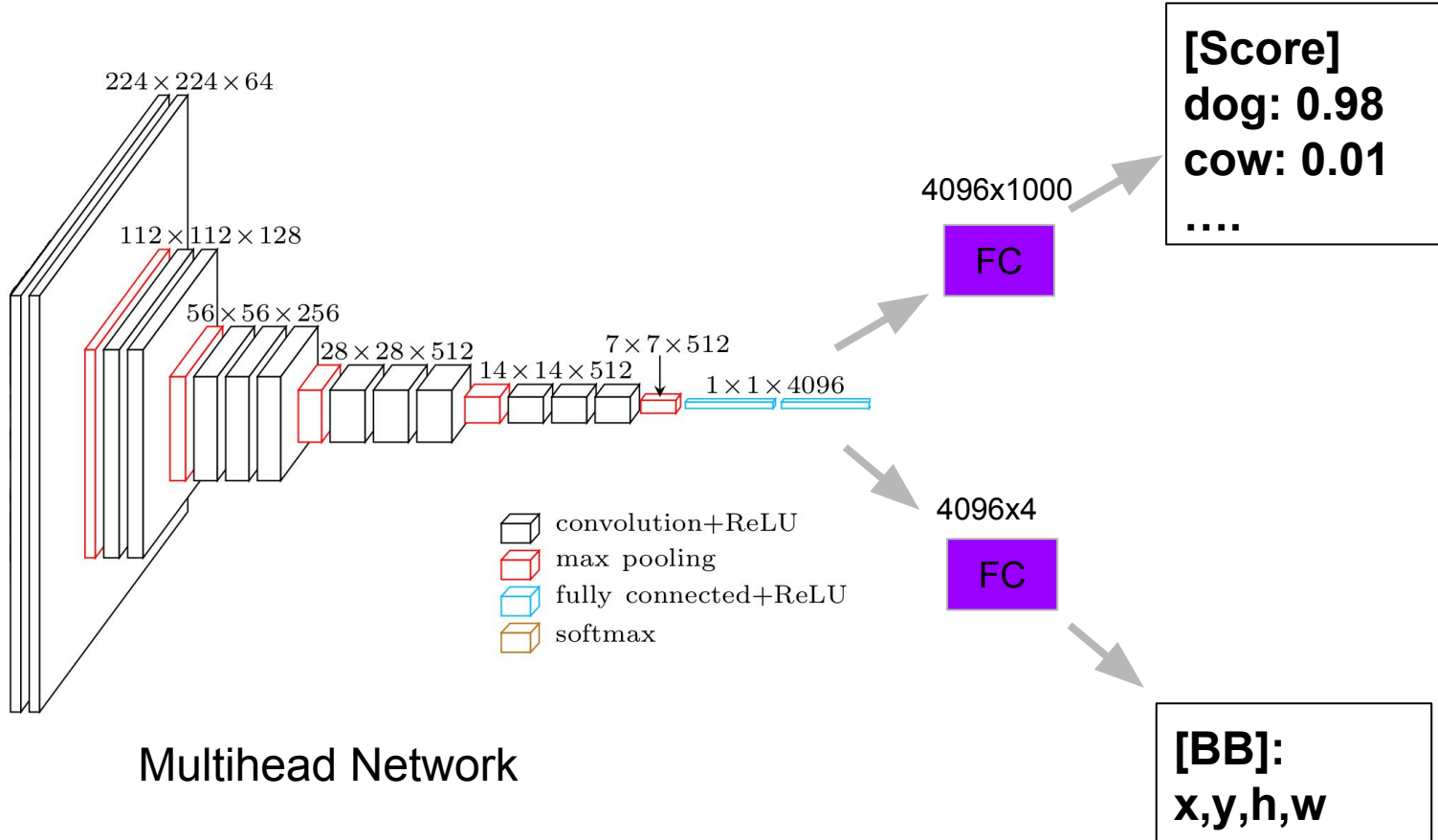
**dog: 0.98**  
**cow: 0.01**  
....

# Image Classification -> Classification + Localization

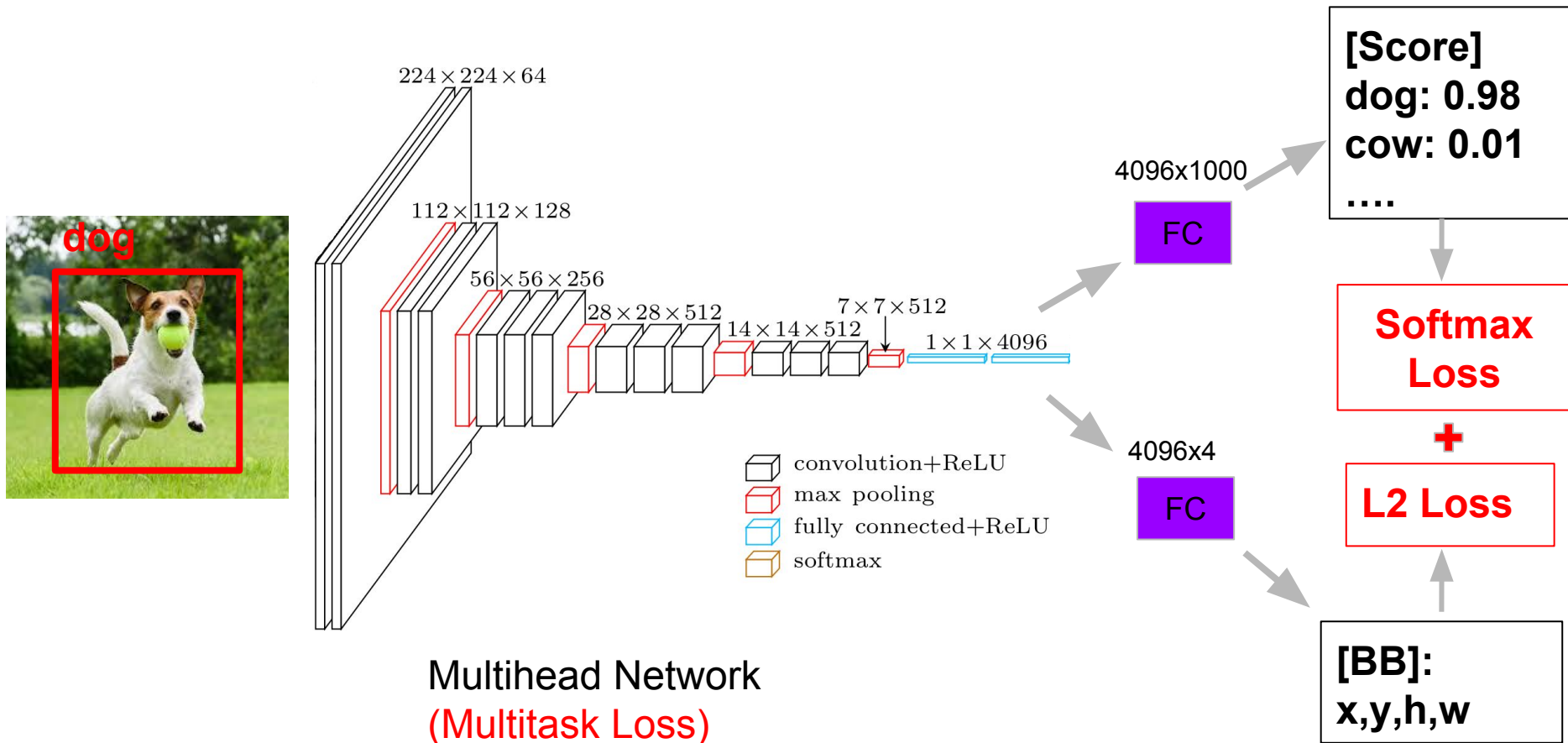


**dog: 0.98**  
**cow: 0.01**  
....

# Image Classification -> Classification + Localization

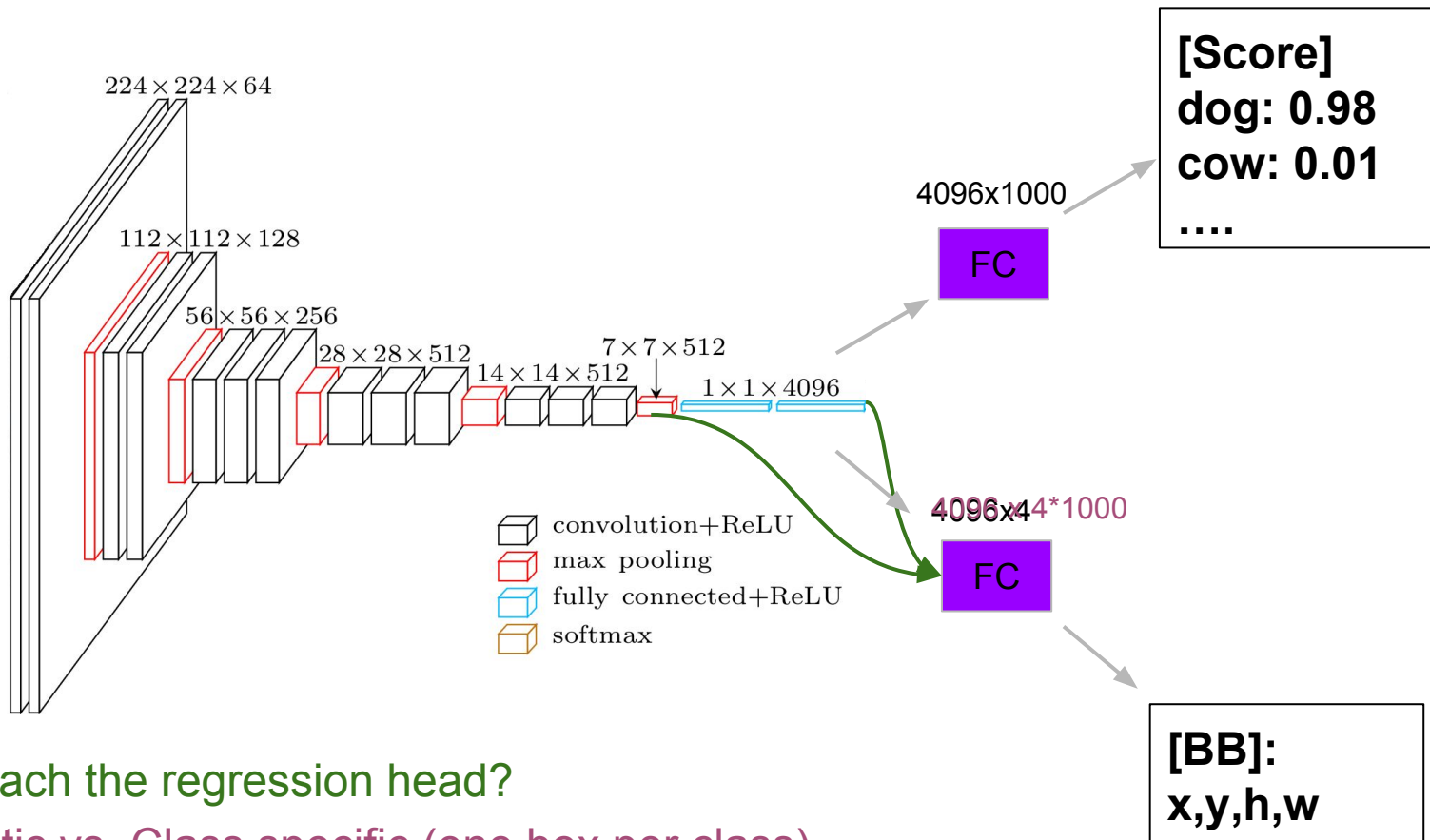


# Localization [Training]





# Localization [Variations]

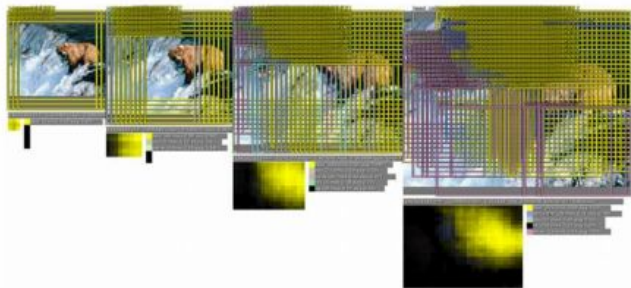


Where to attach the regression head?

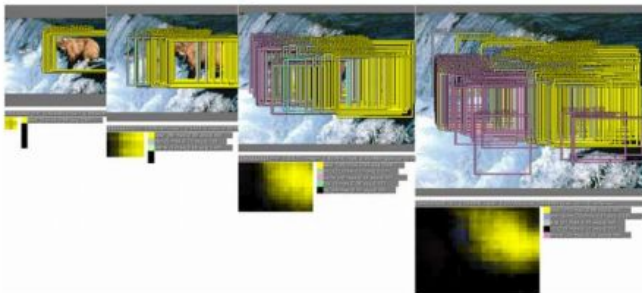
Class agnostic vs. Class specific (one box per class)

# Localization [Sliding Window]

Window positions + score maps



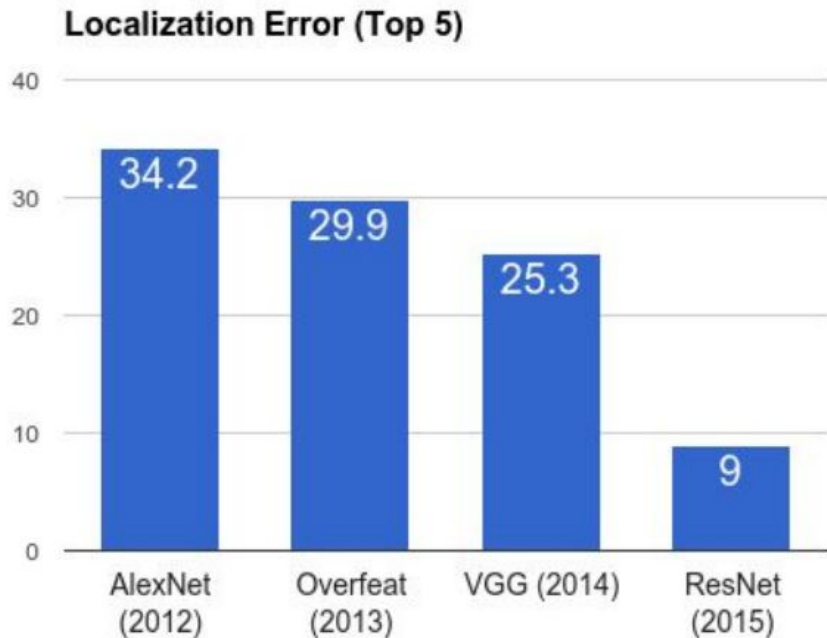
Box regression outputs



Final Predictions



# Localization [Results on ImageNet]



**AlexNet:** Localization method not published

**Overfeat:** Multiscale convolutional regression with box merging

**VGG:** Same as Overfeat, but fewer scales and locations; simpler method, gains all due to deeper features

**ResNet:** Different localization method (RPN) and much deeper features

# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

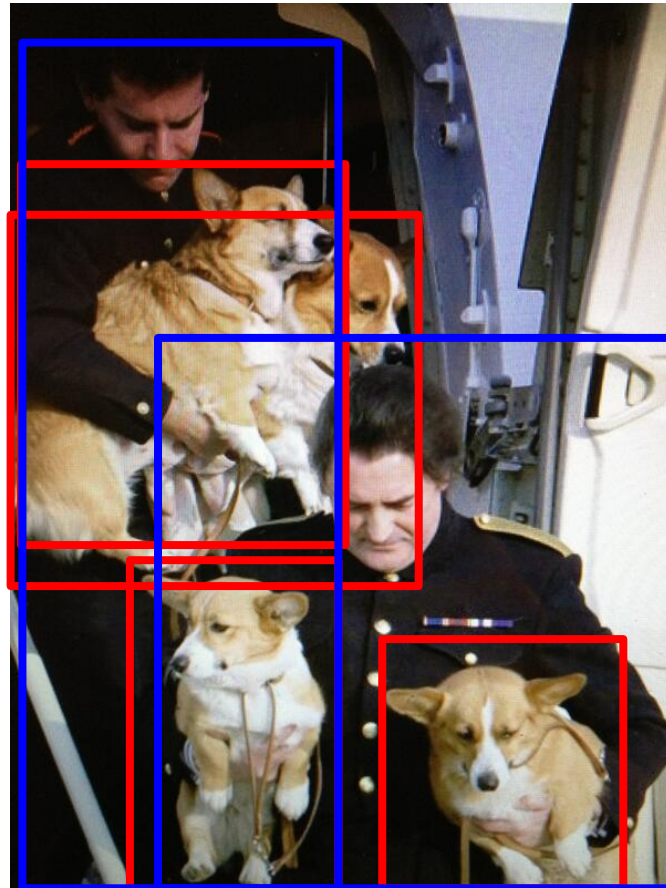
- Faster R-CNN

Single Shot Detector [SSD]

Comparison

Summary

# Object Detection



# Object Detection [Datasets]

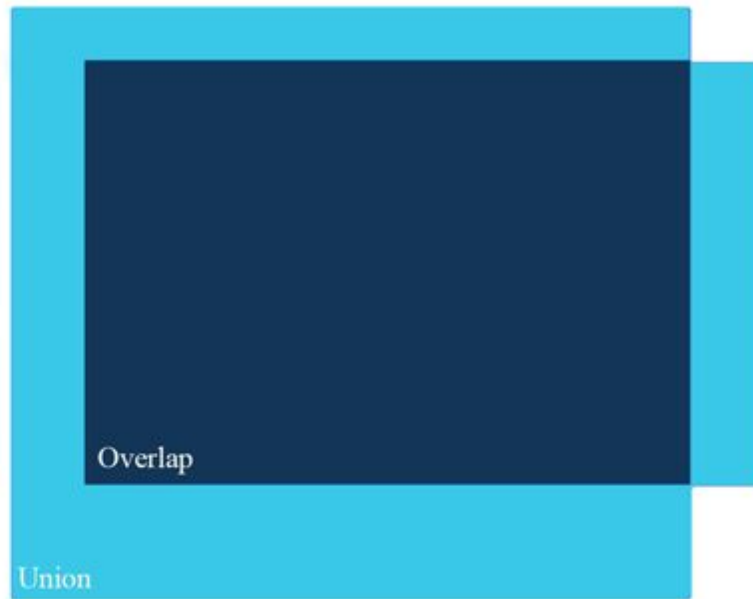
Name	# Images (trainval)	# Classes	Last updated
Open Images Dataset V4	1.74M	600	2018
ImageNet	450k	200	2015
COCO	120K	80	2014
KITTI Vision	7K	3	2014
Pascal VOC	12K	20	2012

# Object Detection [Metrics]

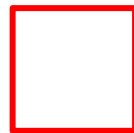
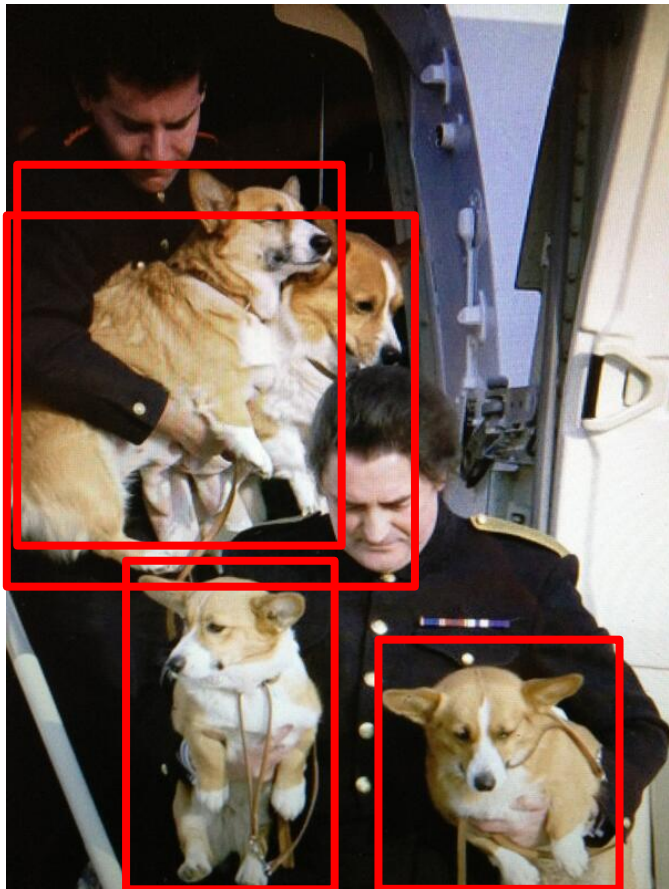
$$IoU = \frac{A \cap B}{A \cup B}$$

A - ground truth (GT)

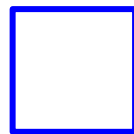
B - detector result (P)



# Object Detection [Metrics]



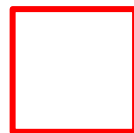
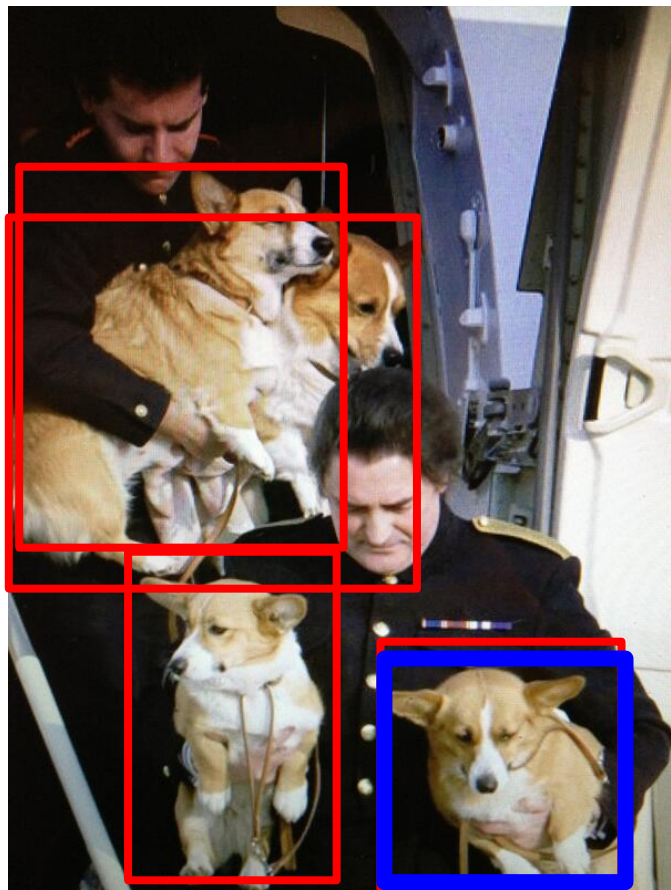
GT



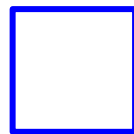
P



# Object Detection [Metrics]



GT

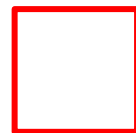
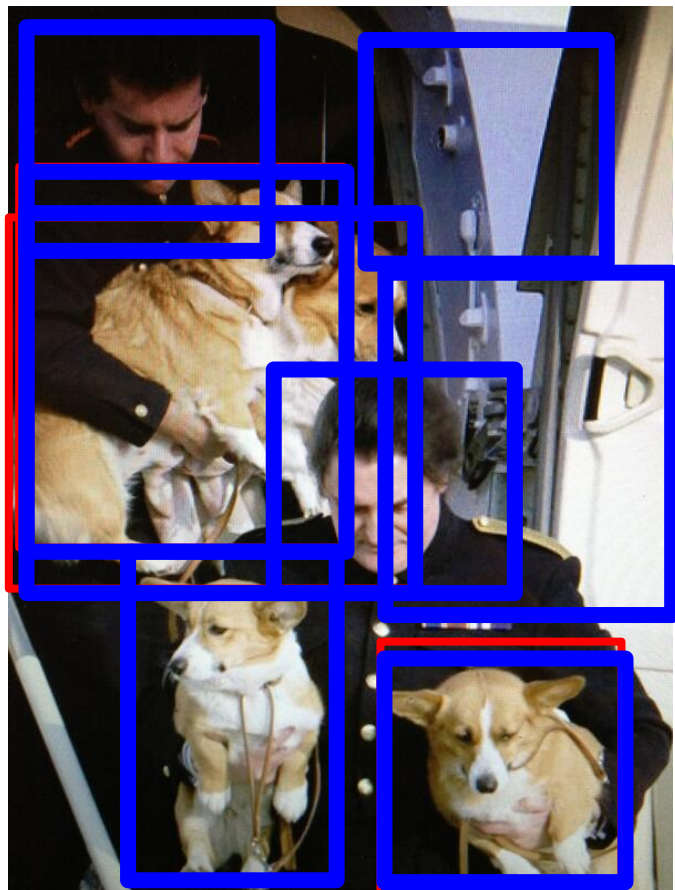


P

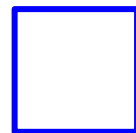
Precision ?

Recall ?

# Object Detection [Metrics]



GT

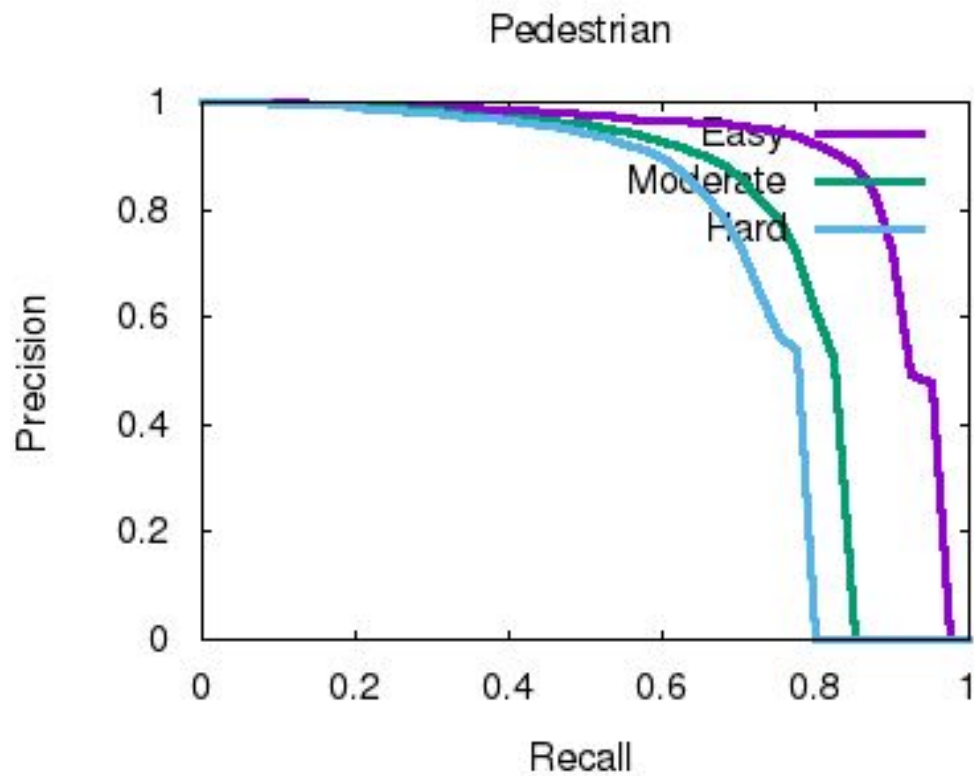


P

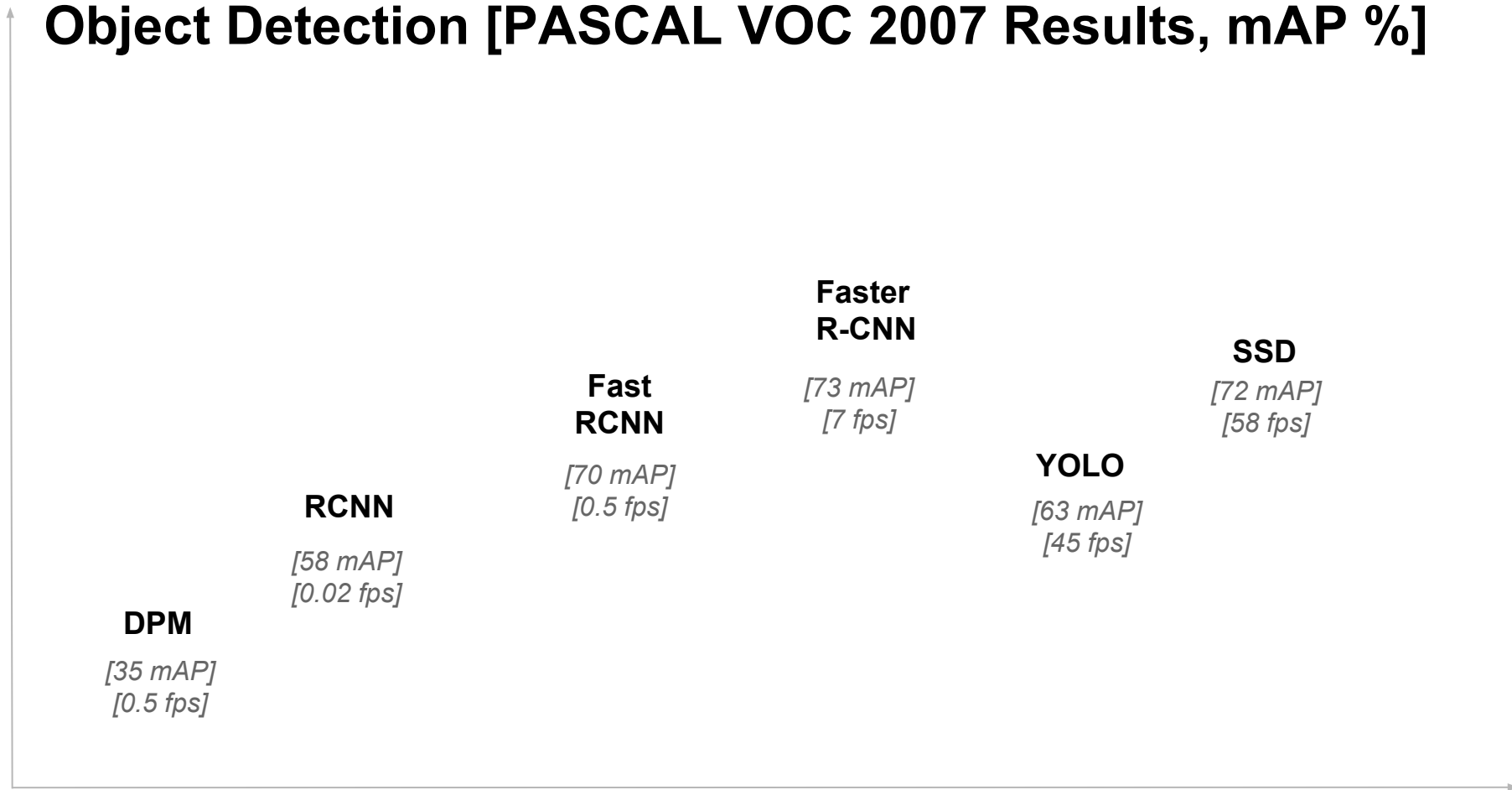
Precision ?

Recall ?

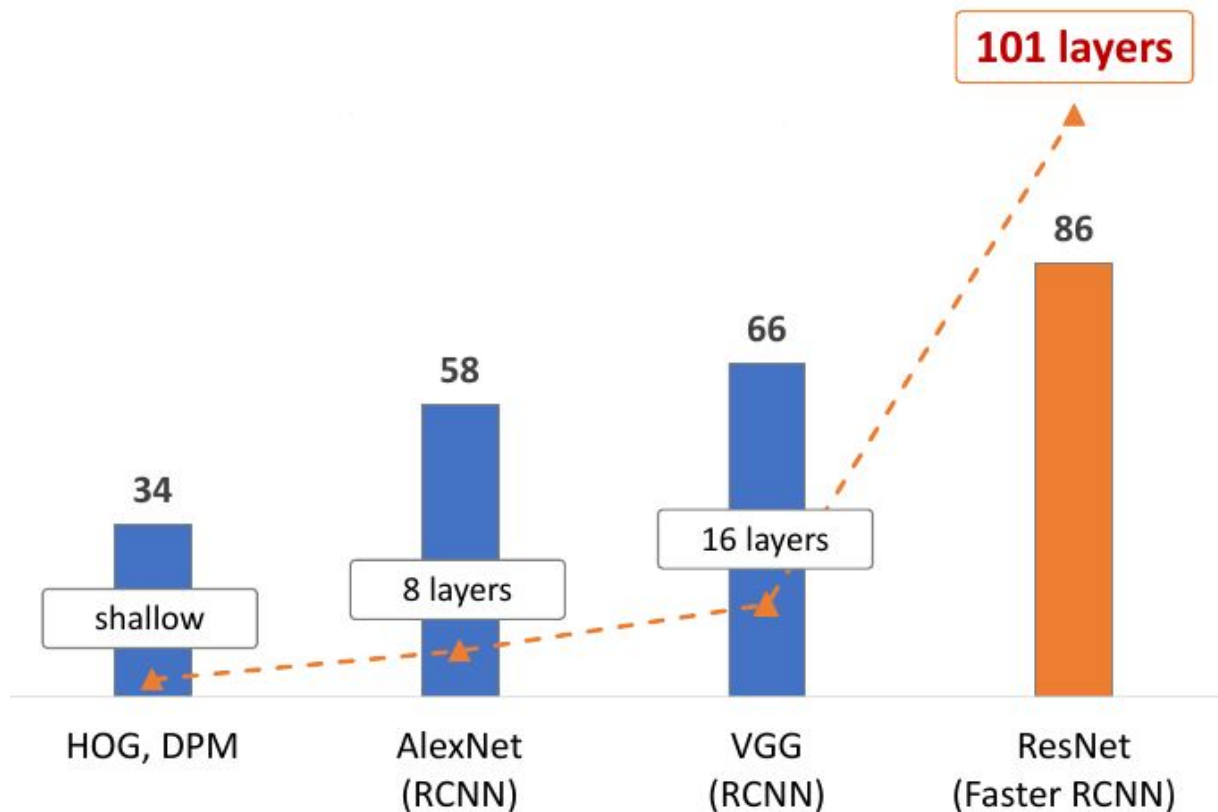
# Object Detection [Metrics]



# Object Detection [PASCAL VOC 2007 Results, mAP %]



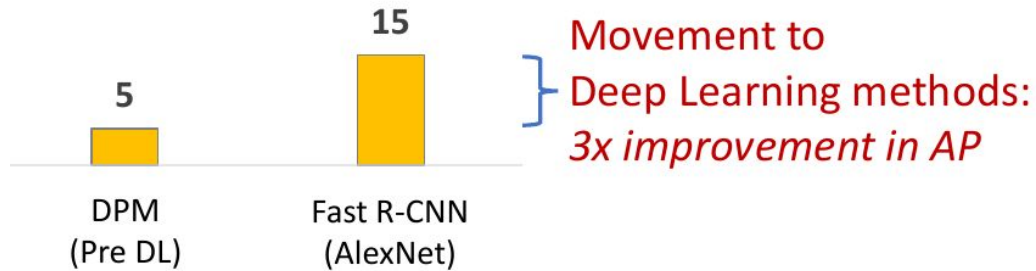
# Object Detection [PASCAL VOC 2007 Results, mAP %]



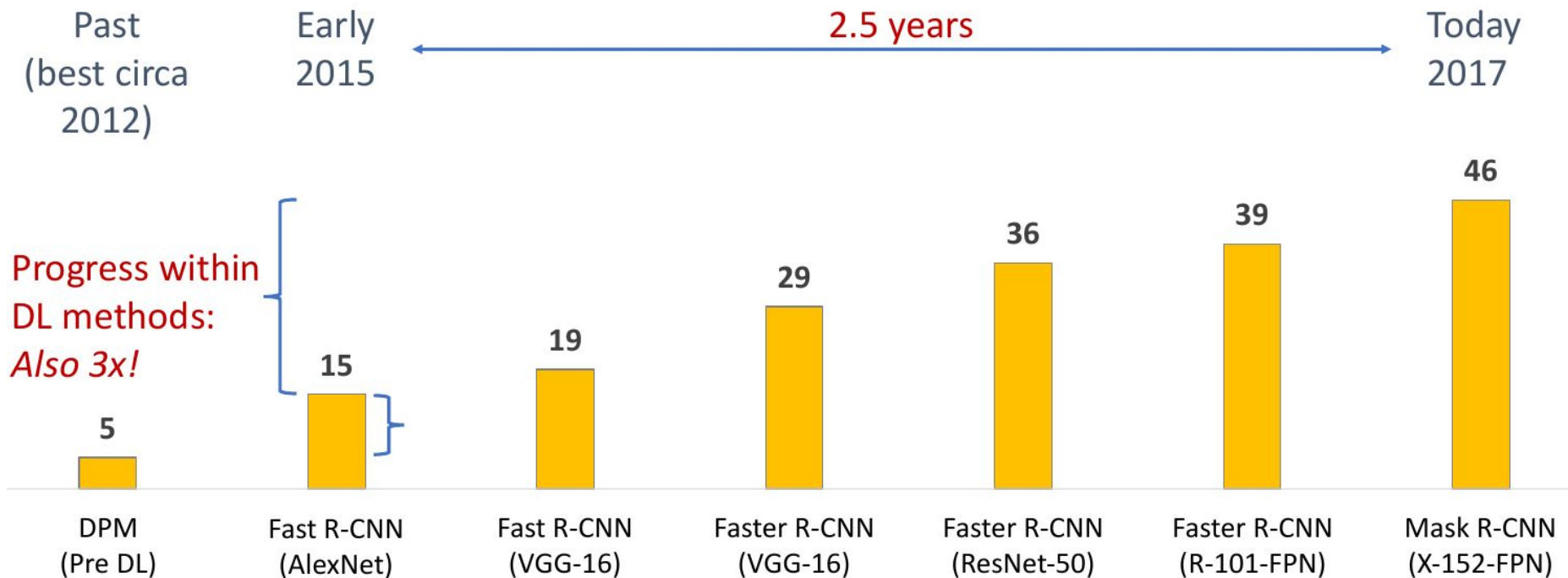
# Object Detection [MS COCO Results, AP %]

Past  
(best circa  
2012)

Early  
2015



# Object Detection [COCO Results, AP %]



# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

- Faster R-CNN

Single Shot Detector [SSD]

Comparison

Summary



Faster

Fast

R-CNN

(Regions with CNN features)



2013

arXiv:1311.2524

Ross Girshick  
[Facebook AI Research (FAIR)]

2015

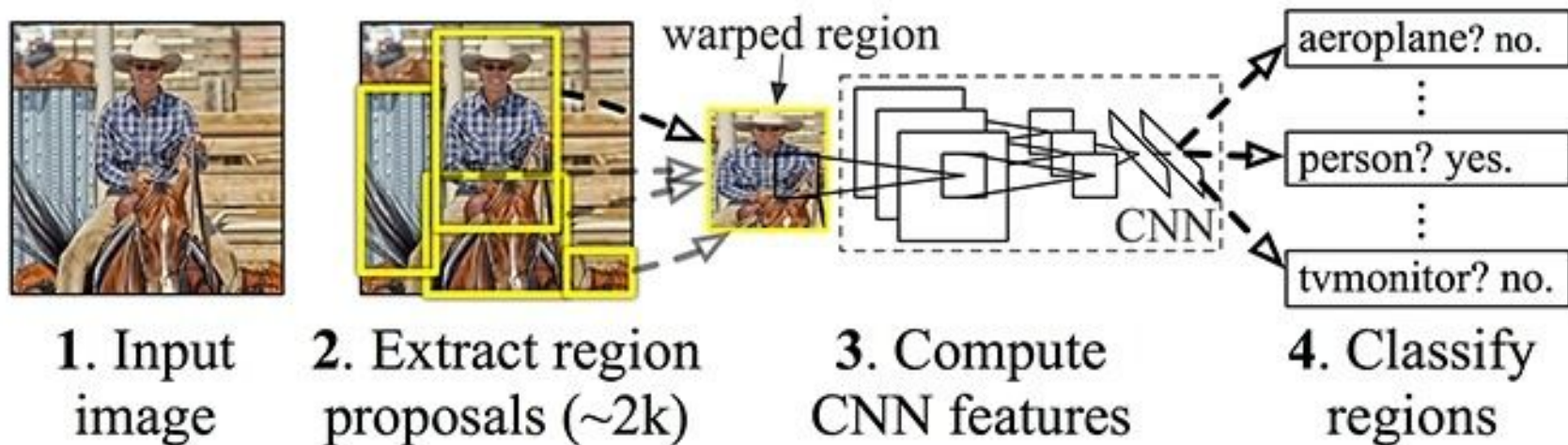
arXiv:1504.08083

<http://www.rossgirshick.info/>

2015

arXiv:1506.01497

# R-CNN: Regions with CNN features

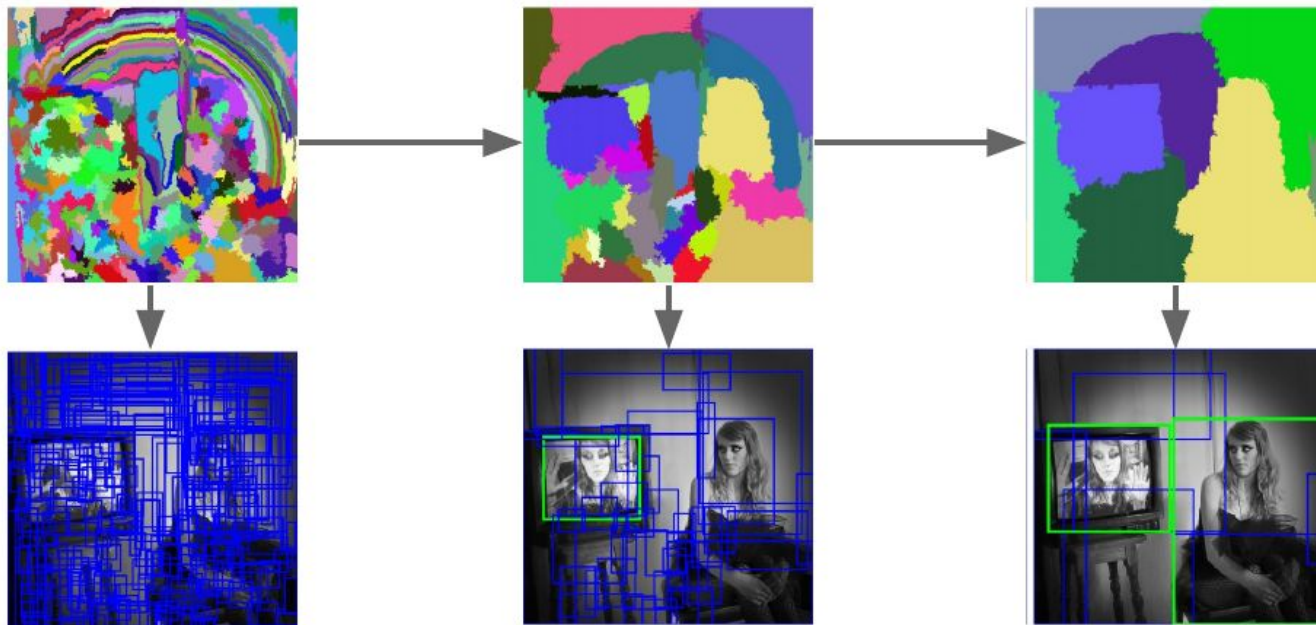


- Extract possible objects using a region proposal method (the most popular one being Selective Search).
- Extract features from each region using a CNN.
- Classify each region with SVMs.

# Object Detection [Region Proposal]

Bottom-up segmentation, merging regions at multiple scales

Convert  
regions  
to boxes



Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

# R-CNN Summary

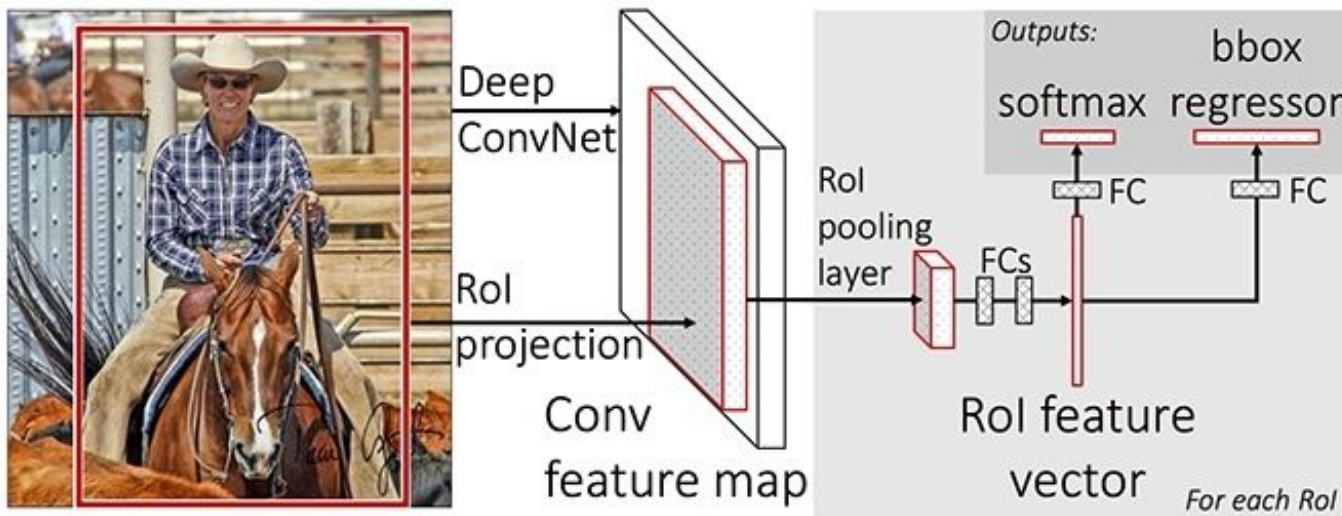
## Pros:

- generalization of CNN features for detection problem;
- almost 50% improvement on the object detection challenge;
- “when labeled training data is scarce, supervised pre-training for an auxiliary task, followed by domain-specific fine-tuning, yields a significant performance boost”

## Cons:

- time & memory inefficiency;
- complex training;
- relied upon external region proposal generator (Selective Search);

# Fast R-CNN



- Use Selective Search to generate object proposals;
- Apply the CNN on the complete image and then used both **Region of Interest (RoI) Pooling** on the feature map;
- Use feed forward network for classification and regression;

# Fast R-CNN Summary

## Pros:

- end-to-end differentiable model, which is easier to train;
- much faster training time (Fast R-CNN trains the very deep VGG16 network 9x faster than R-CNN, is 213x faster at test-time);
- better accuracy than R-CNN

## Cons:

- relied upon external region proposal generator (Selective Search), similar to R-CNN;

# Contents

Object Localization

Object Detection

Datasets

Metrics

Object Detection [RCNN Family]

R-CNN

Fast R-CNN

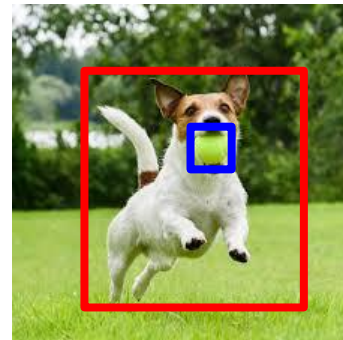
Faster R-CNN

Single Shot Detector [SSD]

Comparison

Summary

# Faster R-CNN (Overall architecture)



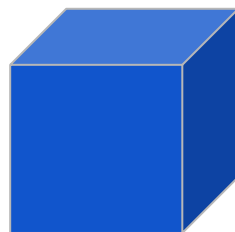
**dog, score: 0.98**  
**ball, score: 0.6**



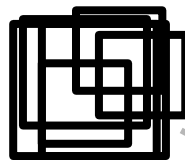
# Faster R-CNN [Overall architecture]



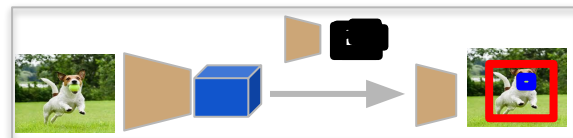
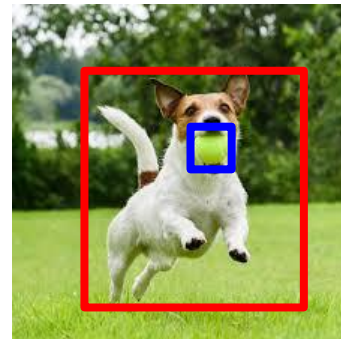
VGG  
up to  
conv5/  
conv5\_1



RPN

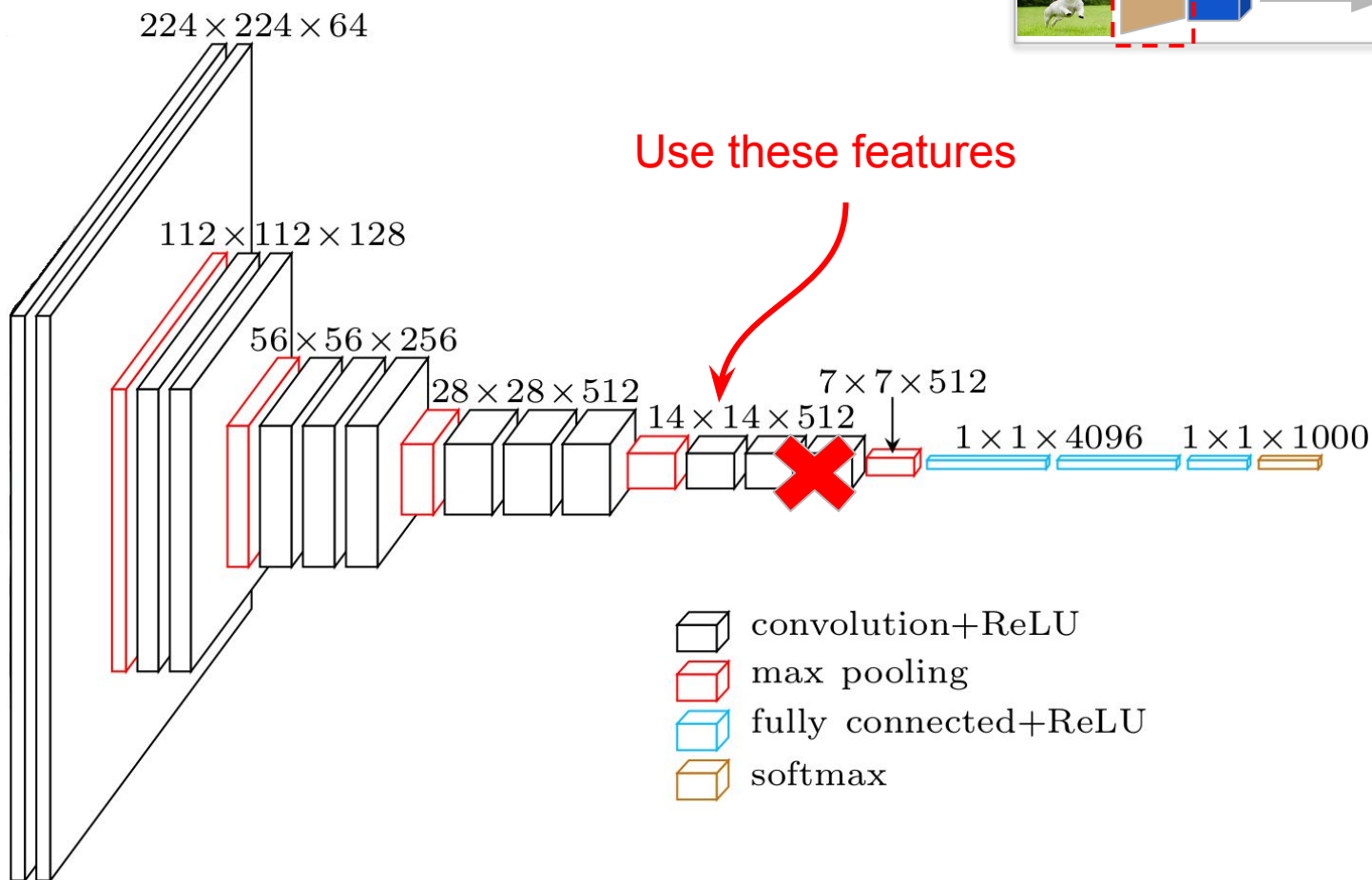
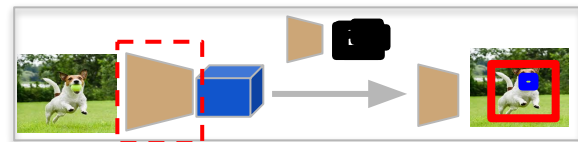


Det.  
&  
class.

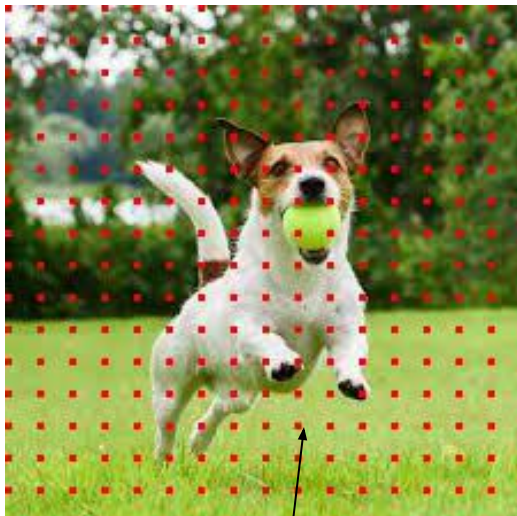
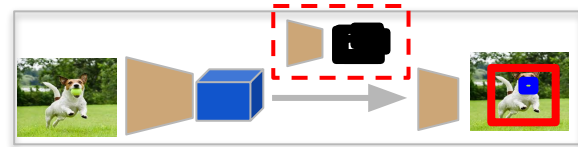


**dog, score: 0.98**  
**ball, score: 0.8**

# Faster R-CNN [Base Network]

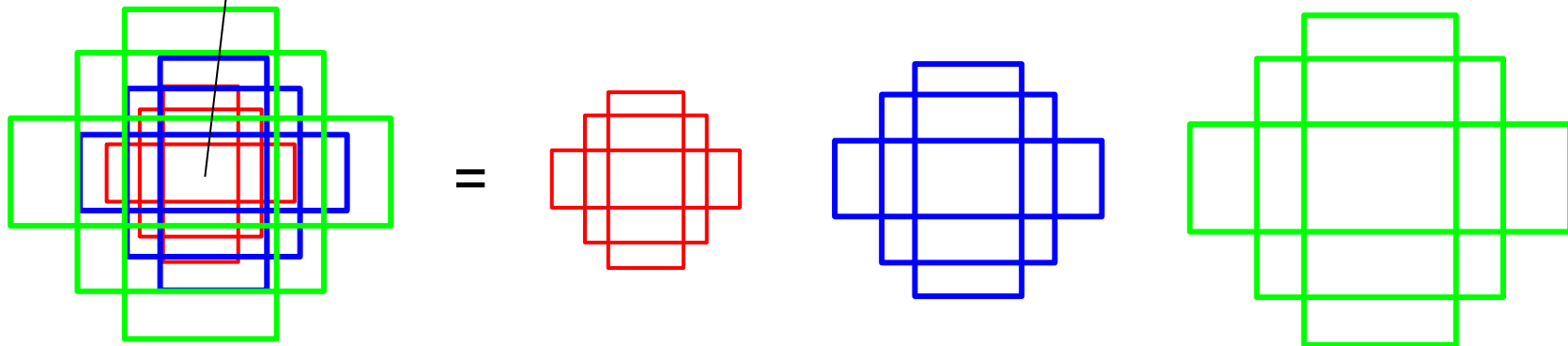


# Faster R-CNN [Anchors]



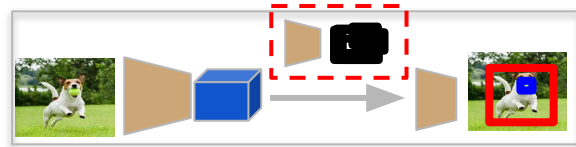
**Anchors** are fixed bounding boxes that are placed throughout the image with different sizes and ratios that are going to be used for reference when first predicting object locations.

- set of sizes (e.g. 64px, 128px, 256px)
  - set of ratios between width and height of boxes (e.g. 0.5, 1, 1.5)
- } = 9



# Faster R-CNN [Region Proposal Network]

RPN takes all the reference boxes (anchors) and outputs a set of good proposals for objects.



$$18 = 2 * 9 = 2 * k$$

14x14x18

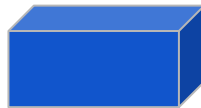


## **Classification layer**

the score of it being background and the score of it being foreground (an actual object)

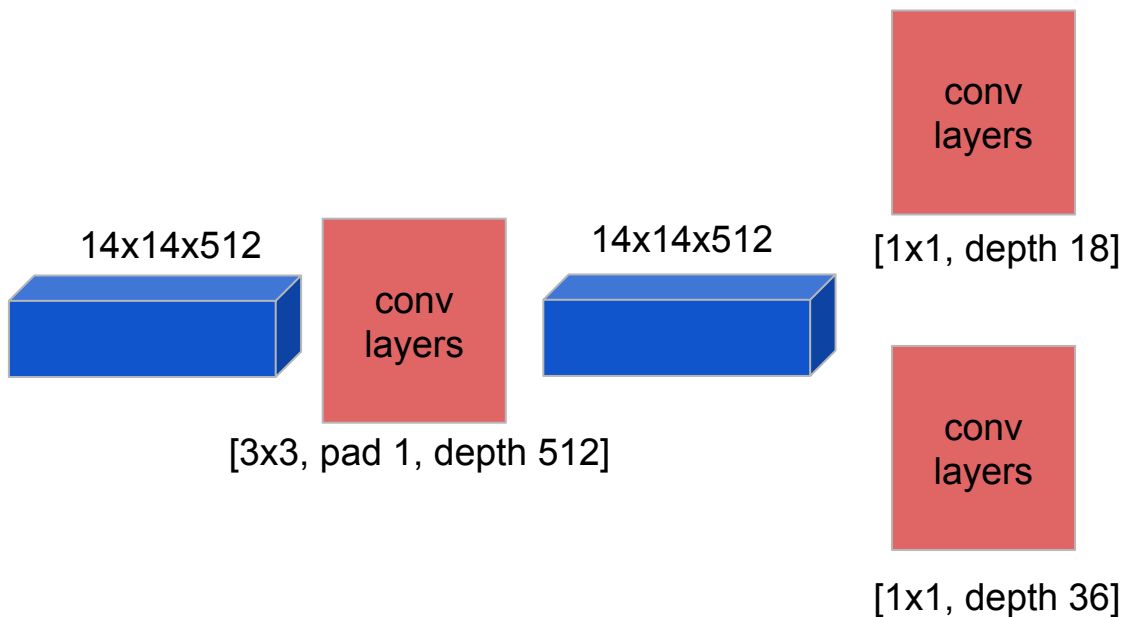
## **Regression layer**

14x14x36

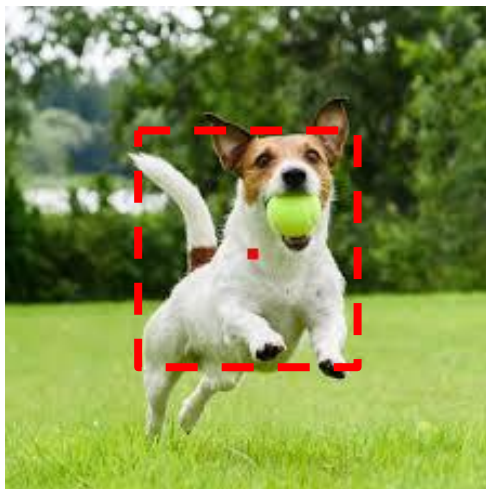
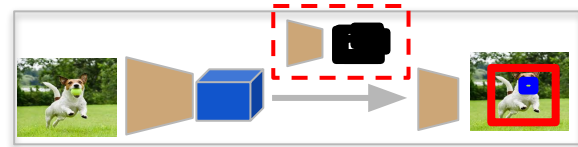


$\Delta x_{center}$ ,  $\Delta y_{center}$ ,  $\Delta w$ ,  $\Delta h$ , which will be applied to the anchors to get the final proposals.

$$36 = 4 * 9 = 4 * k$$



# Faster R-CNN [Region Proposal Network]



RPN output:

score (object) = 0.9

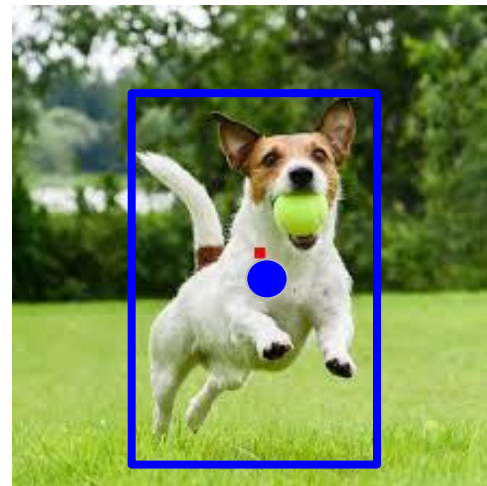
score (background) = 0.1

$\Delta x_{center} = 3$

$\Delta y_{center} = 10$

$\Delta w = 10$

$\Delta h = 40$



default anchor  
position and size:  
 $X_a, Y_a, W_a, H_a$

$$X = \Delta x_{center} * W_a + X_a$$

$$Y = \Delta y_{center} * W_a + Y_a$$

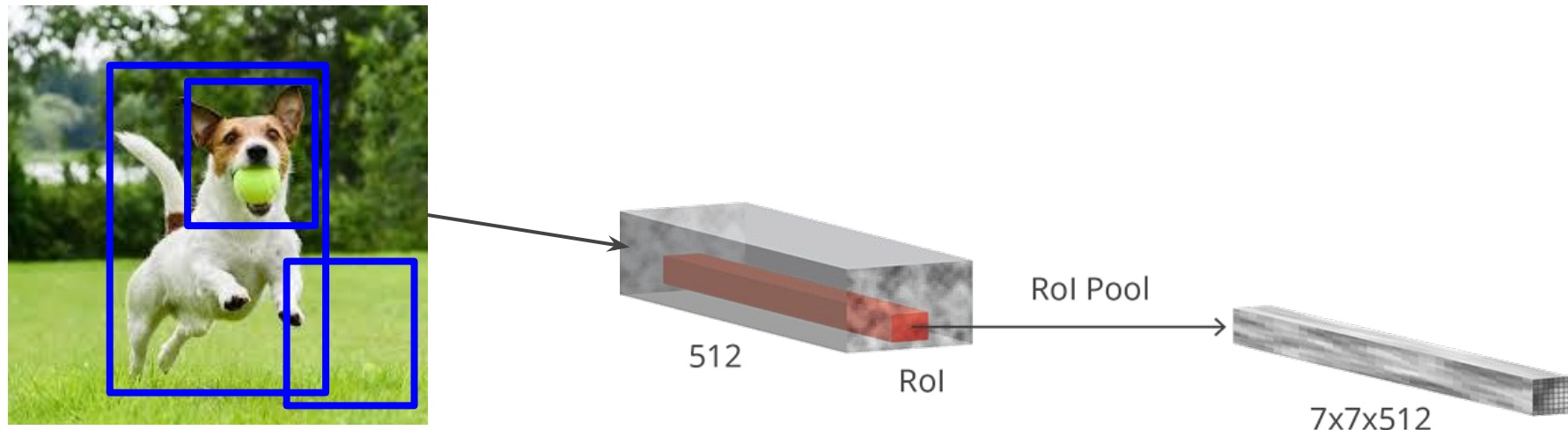
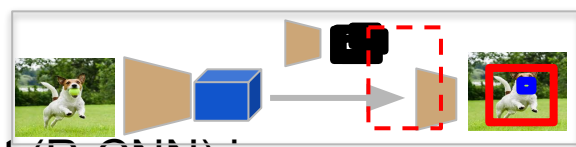
$$W = W_a * \exp(\Delta w)$$

$$H = H_a * \exp(\Delta h)$$

proposal box  
position and size:  
 $X, Y, W, H$

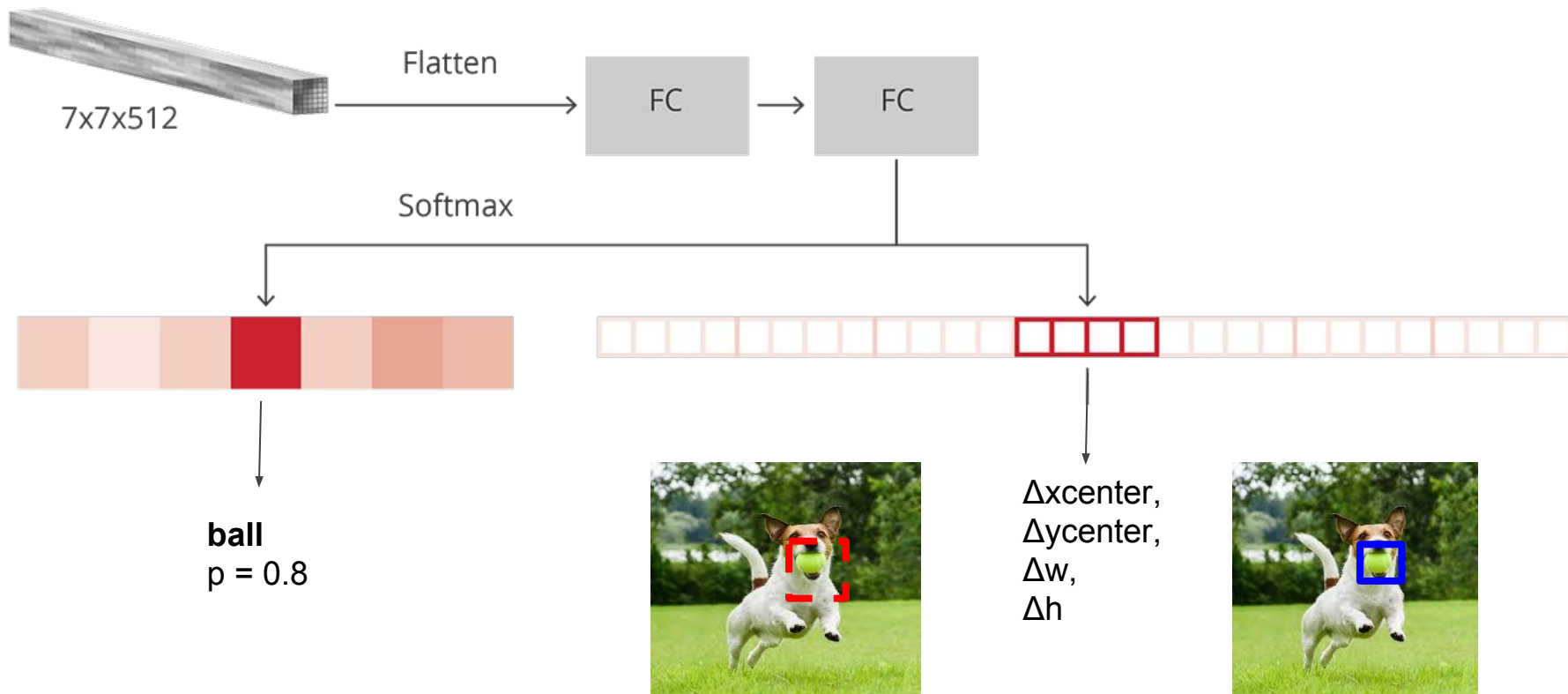
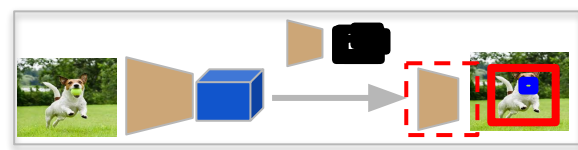
# Faster R-CNN [Region of Interest Pooling]

Fixed size feature maps are needed for the det.&class. part (R-CNN) in order to classify them into a fixed number of classes.

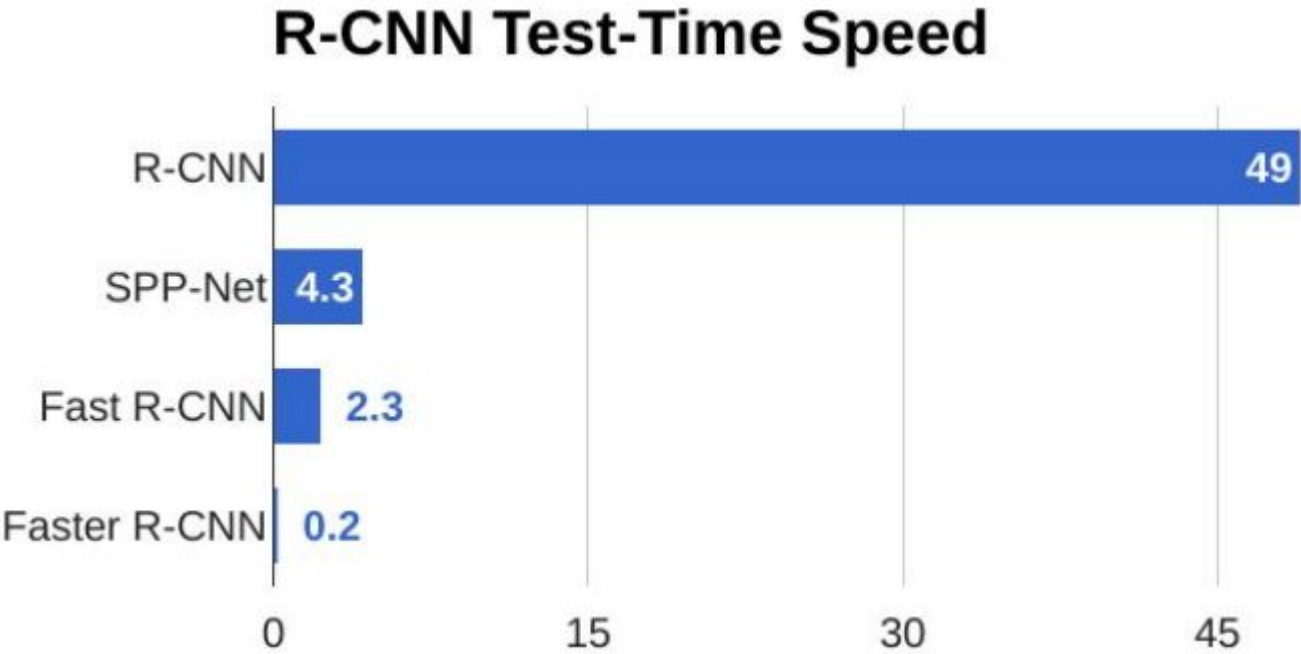


crop the convolutional features using each proposal and then resize to a fixed sized  $14 \times 14 \times D$  using interpolation (usually bilinear). After cropping, max pooling with a  $2 \times 2$  kernel is used to get a final  $7 \times 7 \times D$  for each proposal.

# Faster R-CNN [Region-based CNN]



# Faster R-CNN [Time]





## Faster R-CNN [Training]

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*)$$

mini batch of size 256 — trying to maintain a balanced ratio between foreground and background anchors.

The RPN uses all the anchors selected for the mini batch to calculate the classification loss using binary cross entropy. Then, it uses only those minibatch anchors marked as foreground to calculate the regression loss.

for the regression error, the paper uses Smooth L1 loss. Smooth L1 is basically L1, but when the L1 error is small enough, defined by a certain  $\sigma$ , the error is considered almost correct and the loss diminishes at a faster rate.

# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

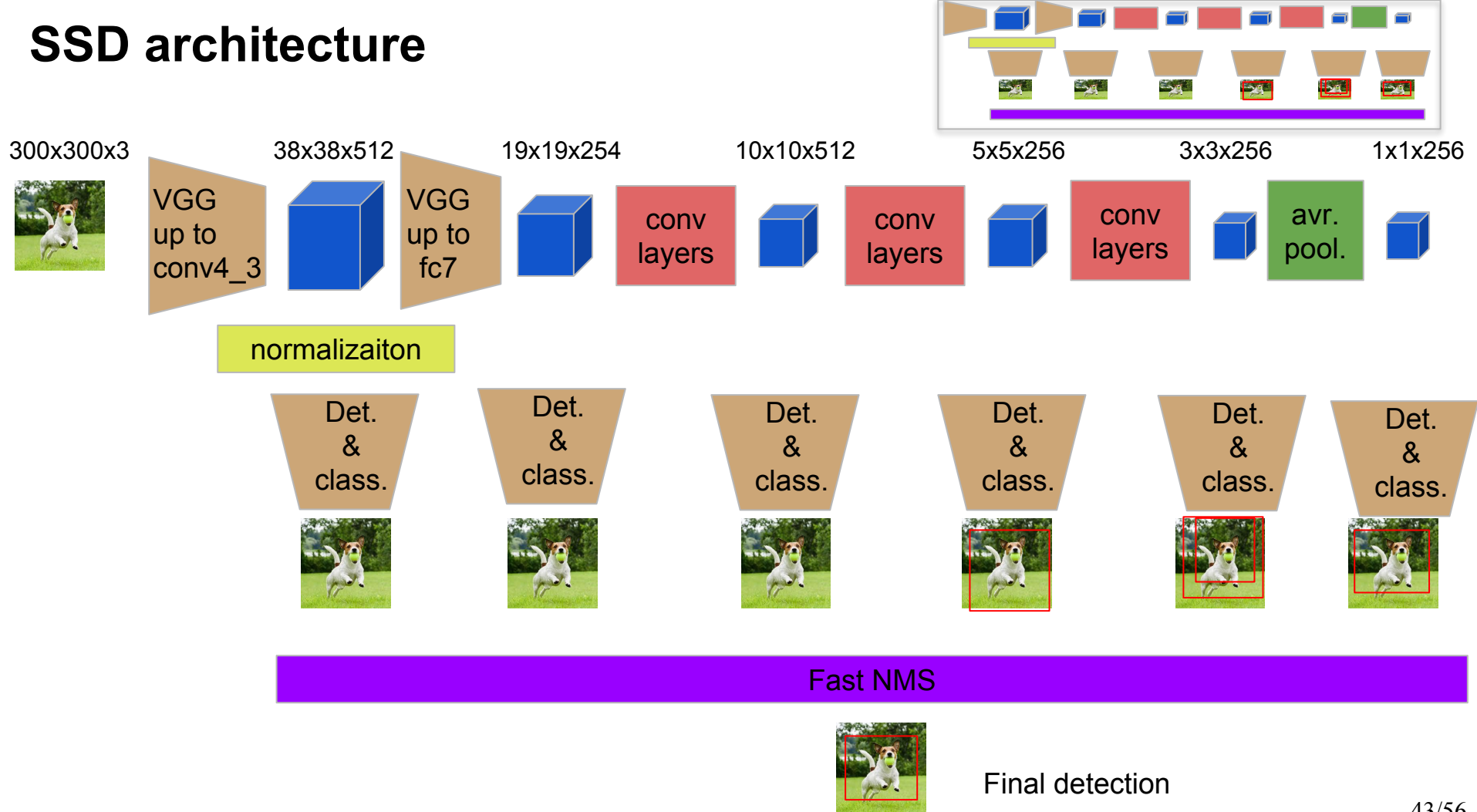
- Faster R-CNN

Single Shot Detector [SSD]

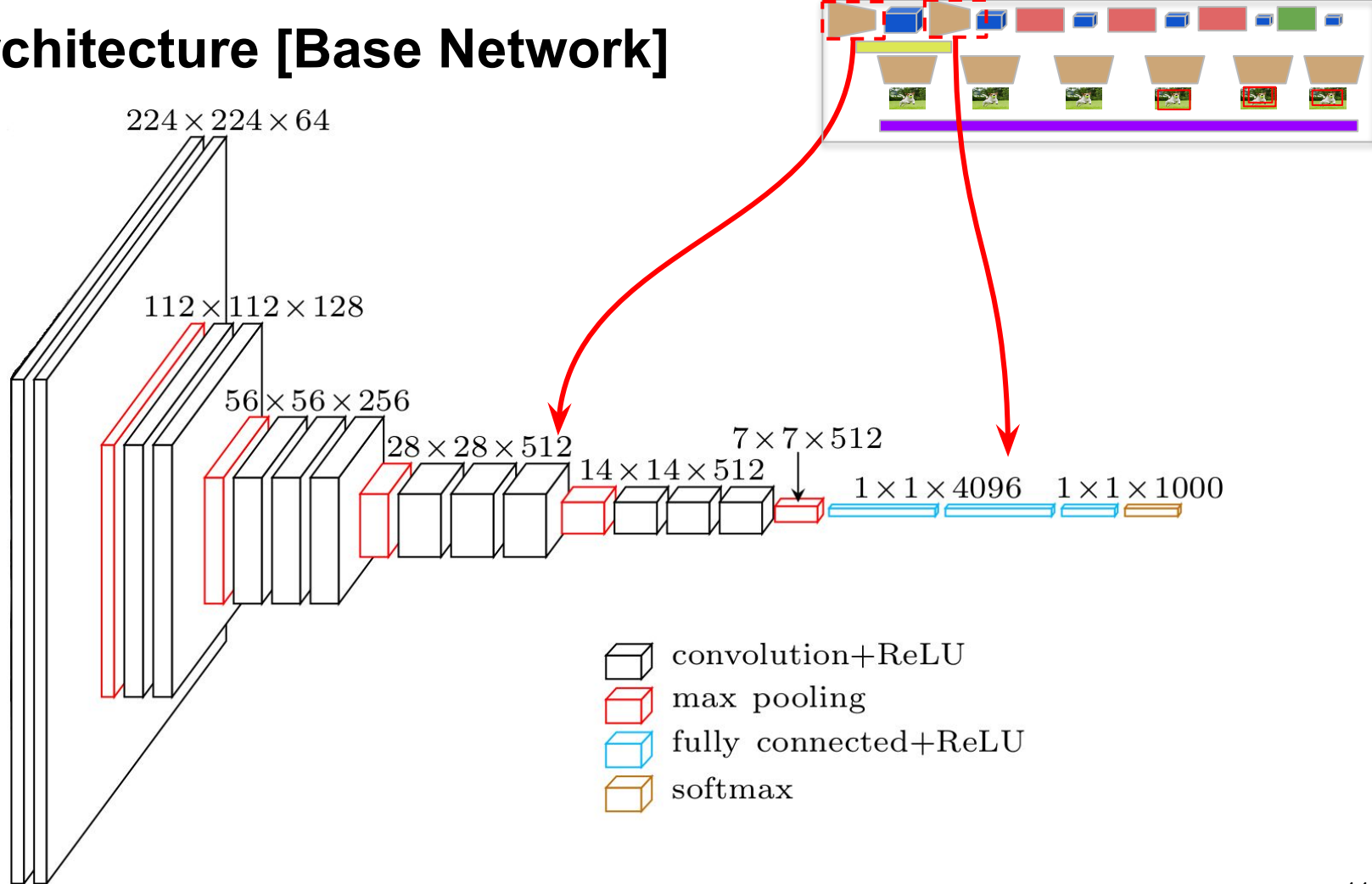
Comparison

Summary

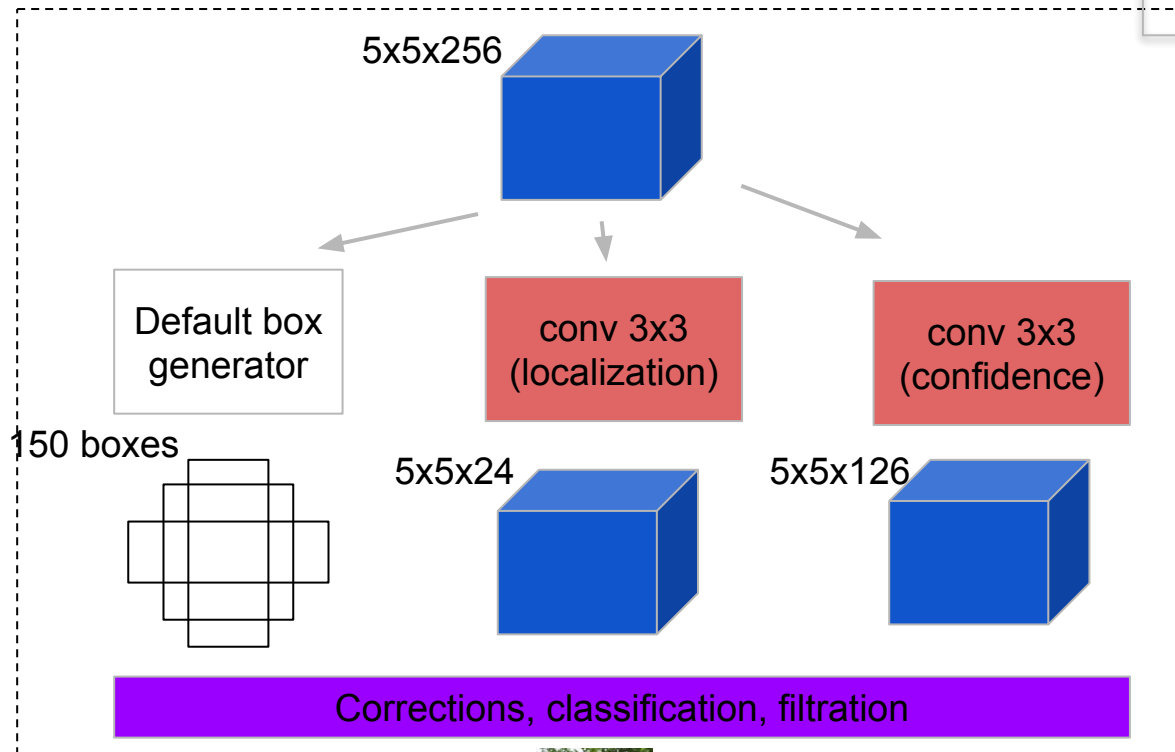
# SSD architecture



# SSD architecture [Base Network]



# SSD architecture [Detector & classifier]

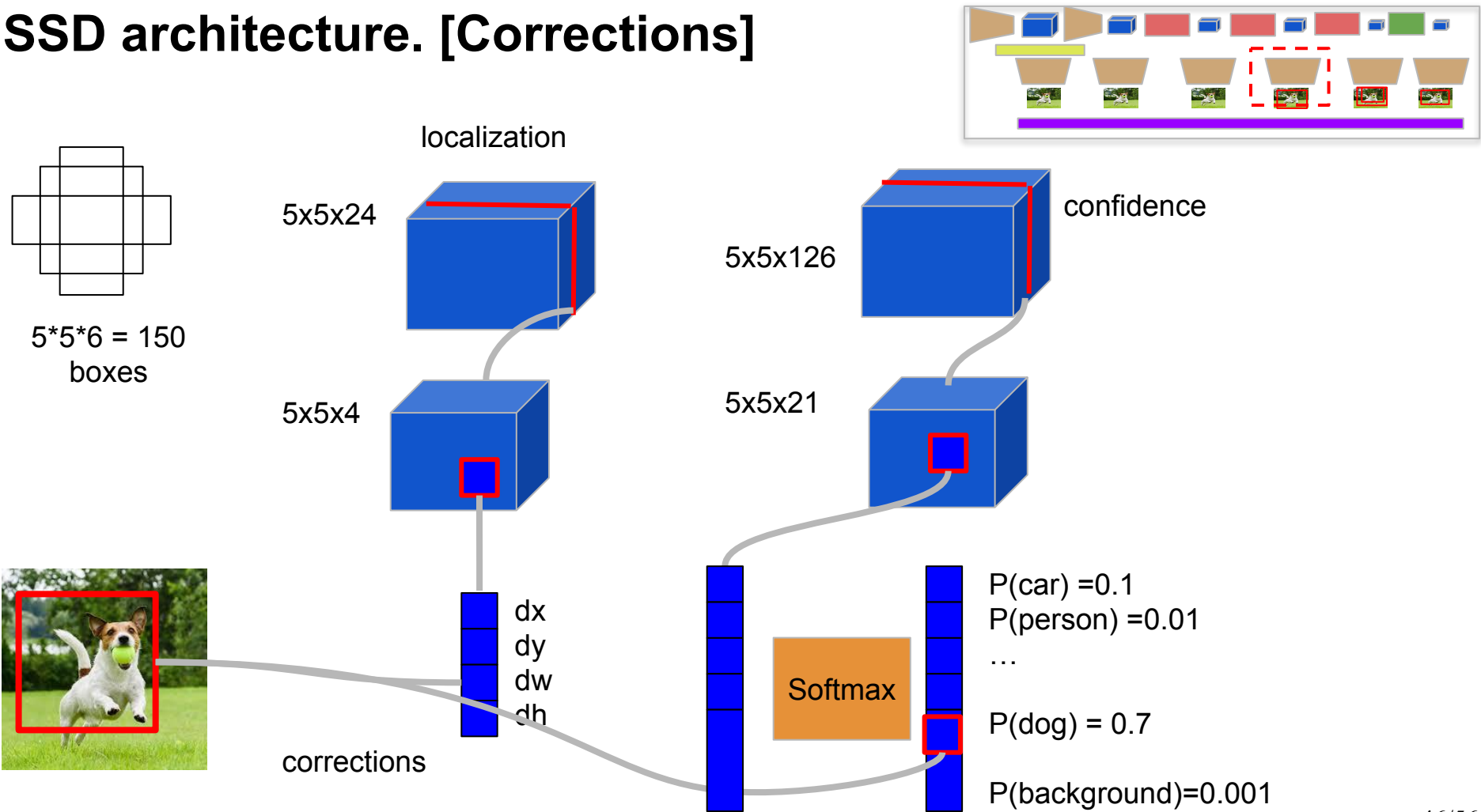


- # default boxes = 6
- 20 + 1 classes
- 4 numbers for BB

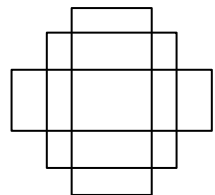
- $126 = 21 \times 6$
- $24 = 4 \times 6$
- $150 = 5 \times 5 \times 6$



# SSD architecture. [Corrections]



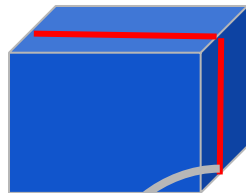
# SSD architecture. [Corrections]



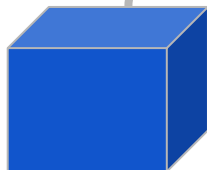
$5 \times 5 \times 6 = 150$   
boxes

localization

$5 \times 5 \times 24$

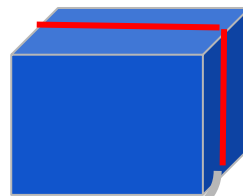


$5 \times 5 \times 4$

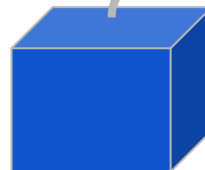


confidence

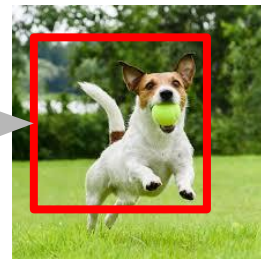
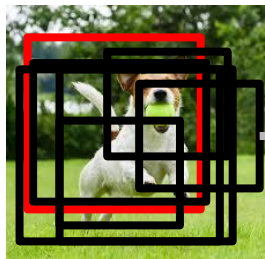
$5 \times 5 \times 126$



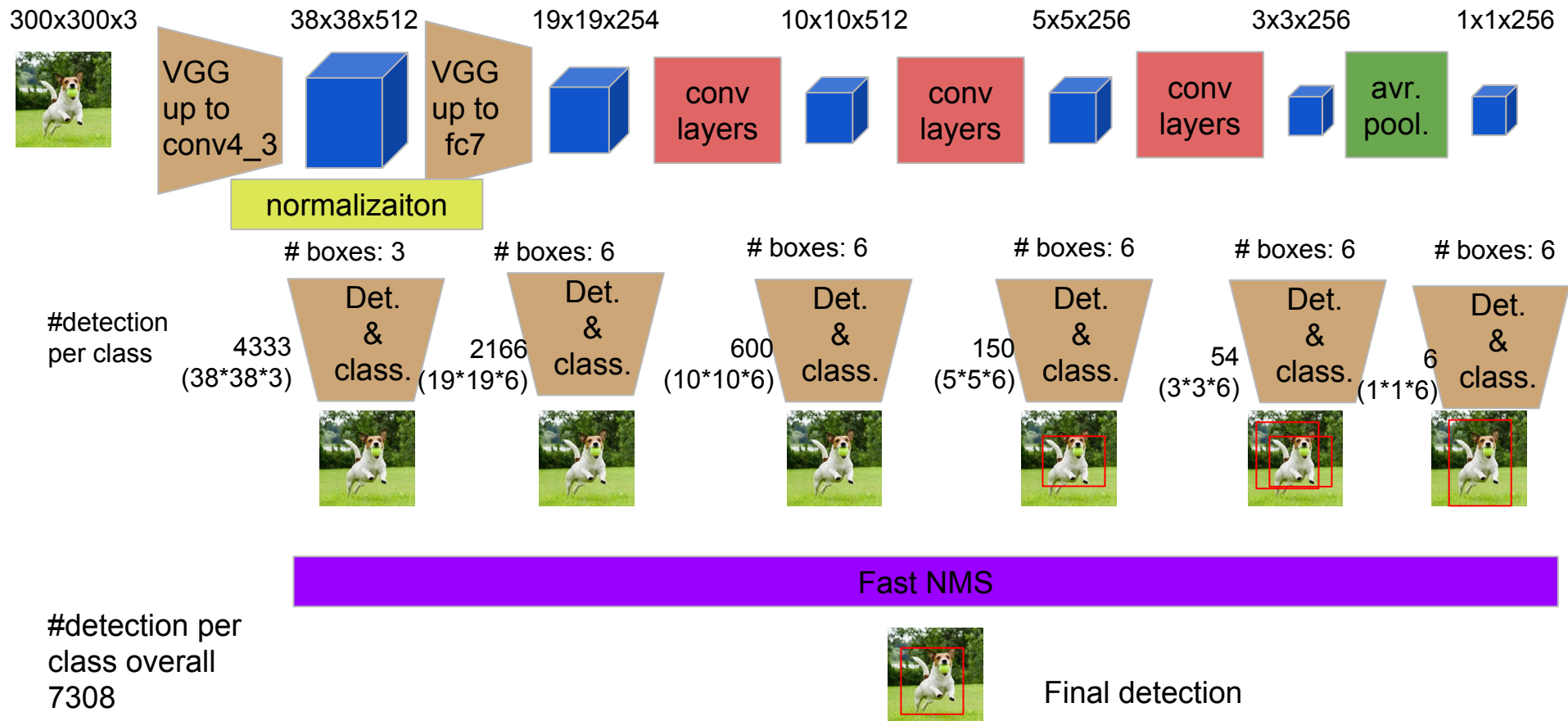
$5 \times 5 \times 21$



Confidence threshold

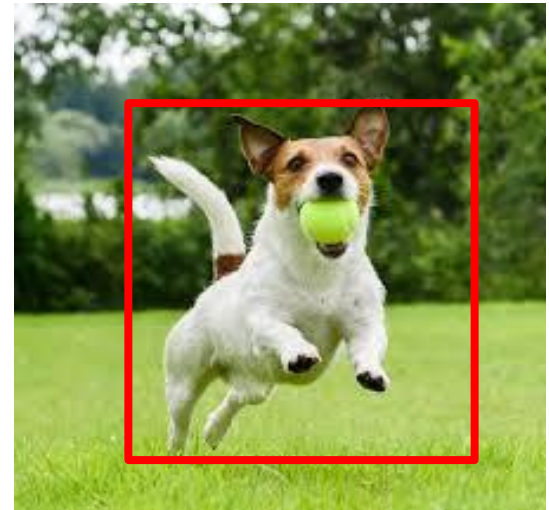
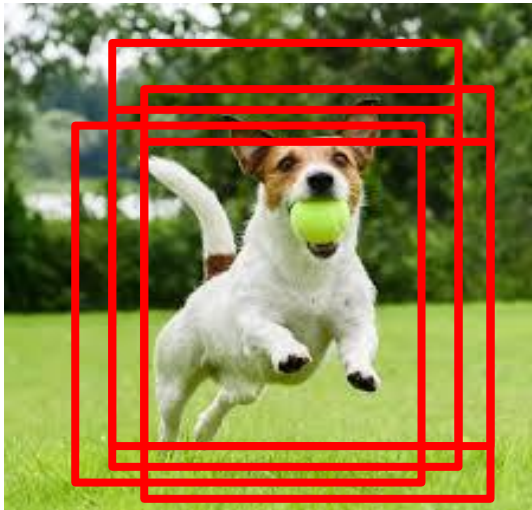


# SSD architecture





# Non-maximum Suppression



# SSD [Training]

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

$$L_{conf}(x, c) = - \sum_{i \in Pos} x_{ij}^p \log(\hat{c}_i^p) - \sum_{i \in Neg} \log(\hat{c}_i^0) \quad \text{where} \quad \hat{c}_i^p = \frac{\exp(c_i^p)}{\sum_p \exp(c_i^p)}$$

$$L_{loc}(x, l, g) = \sum_{i \in Pos} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^k \text{smooth}_{L1}(l_i^m - \hat{g}_j^m)$$

$$\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx}) / d_i^w \quad \hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy}) / d_i^h$$

$$\hat{g}_j^w = \log \left( \frac{g_j^w}{d_i^w} \right) \quad \hat{g}_j^h = \log \left( \frac{g_j^h}{d_i^h} \right)$$

Heavy use of data augmentation.

# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

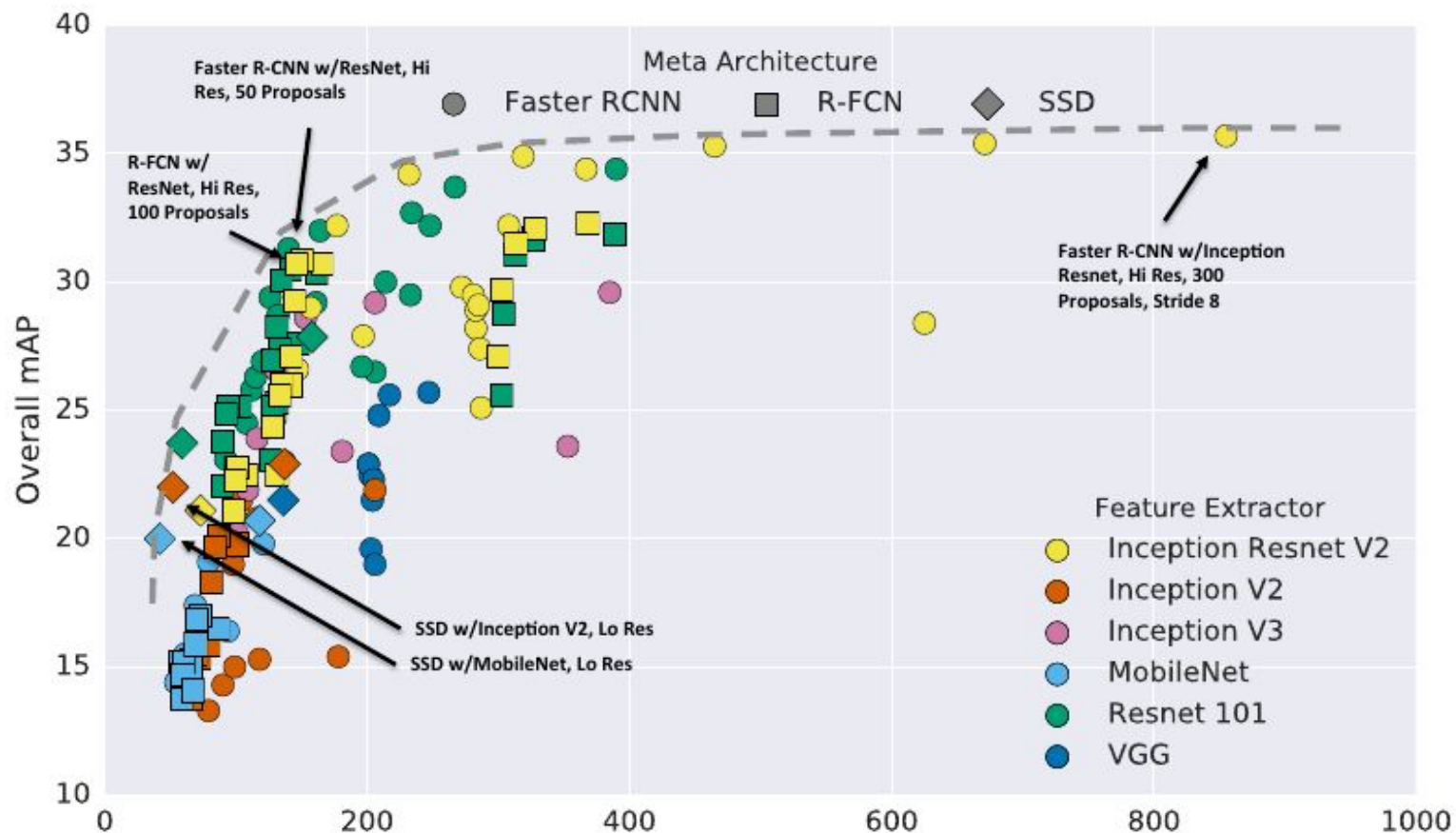
- Faster R-CNN

Single Shot Detector [SSD]

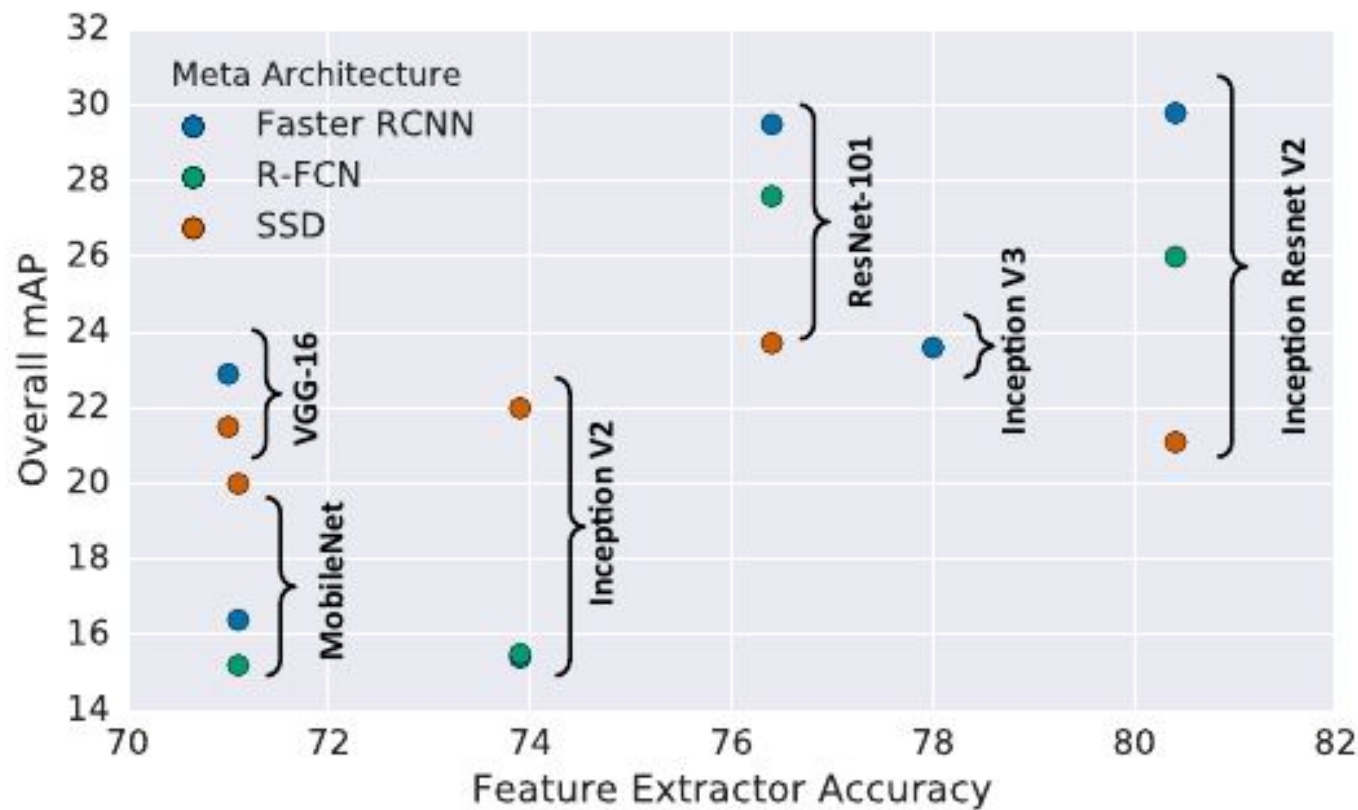
Comparison

Summary

# Object Detection [Comparison]



# Object Detection [Comparison]



# Object Detection [Comparison]

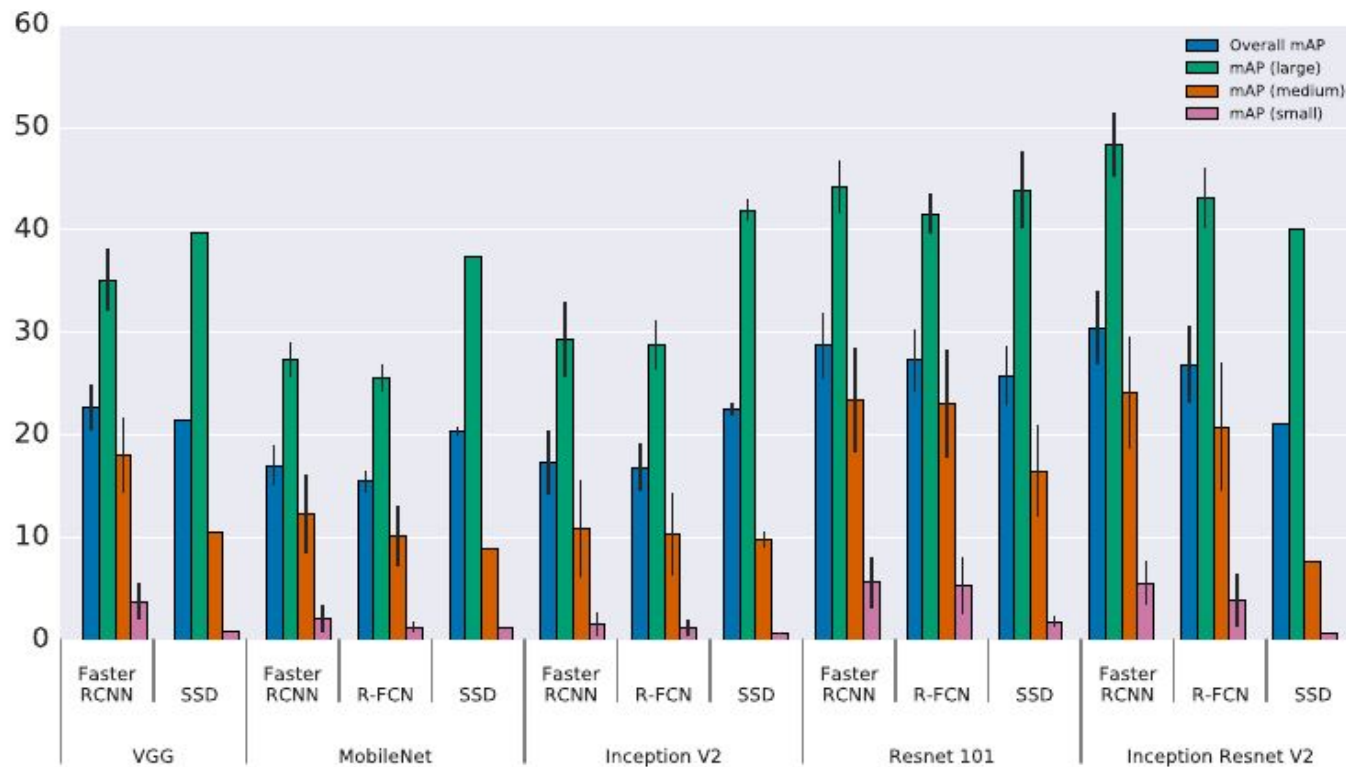


Figure 4: Accuracy stratified by object size, meta-architecture and feature extractor, We fix the image resolution to 300.

# Contents

Object Localization

Object Detection

- Datasets

- Metrics

Object Detection [RCNN Family]

- R-CNN

- Fast R-CNN

- Faster R-CNN

Single Shot Detector [SSD]

Comparison

Summary

# Summary

## Object Localization, Object Detection tasks

### Two stage detector [Faster-RCNN]

- RPN
- End-to-end
- Best accuracy

### One stage detector [SSD architecture]

- Use of multiple convolutional maps to deal with scales;
- More default bounding boxes -> the better the result;
- Comparable to Fast\*-rCNN family but faster;