

Convolutional Neural Networks

– Part 3 Object Detection

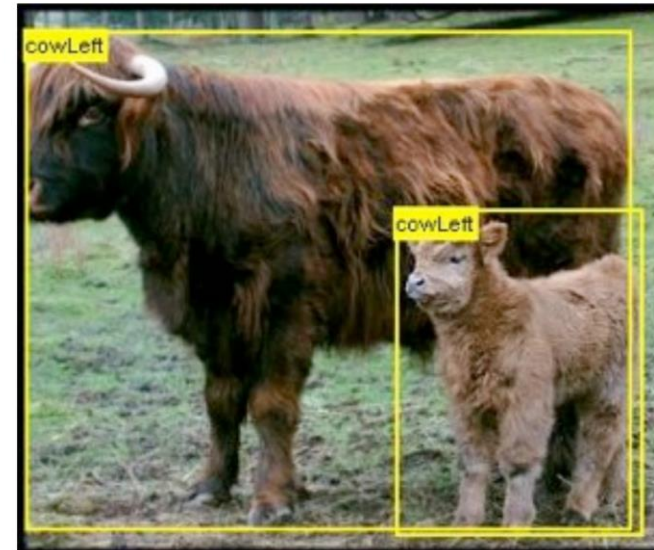
ECE 449

Object Detection

- Before deep learning
 - Deformable part models (DPM)
- Multi-stage
 - R-CNN
 - Fast R-CNN (one-stage)
 - Faster R-CNN
 - Mask R-CNN
- One-stage
 - You only look once (YOLO)
 - Single-shot detector (SSD)
 - CenterNet (anchor free)
- Non-maximum suppression (NMS)

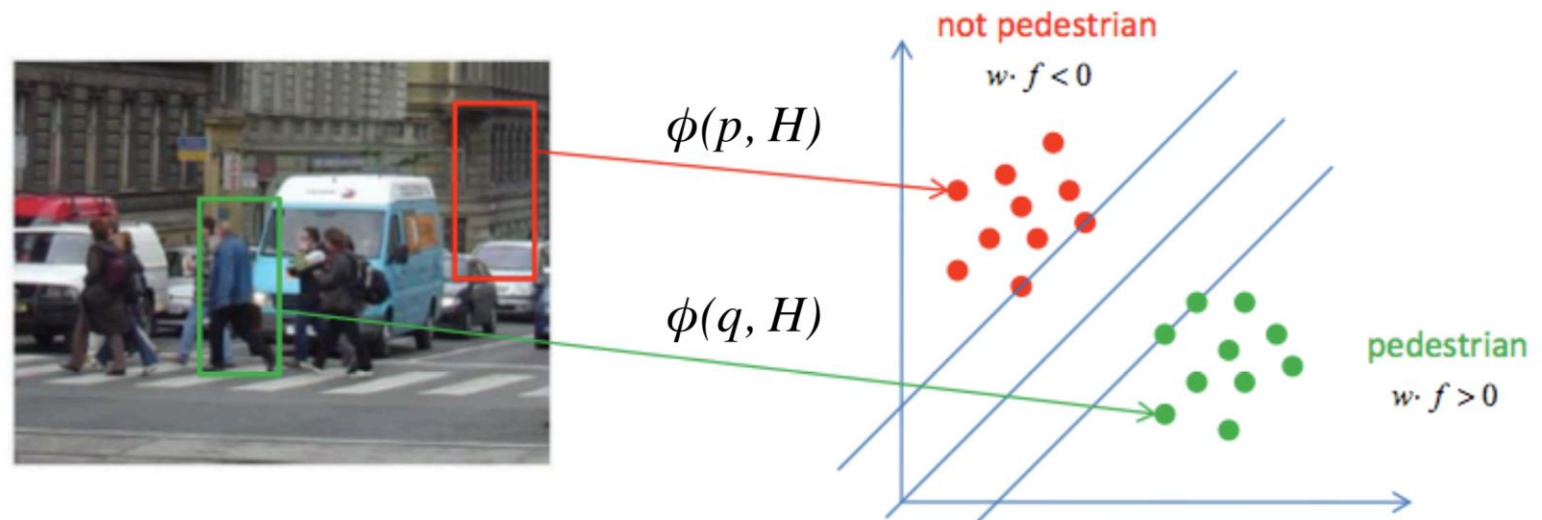
Object Detection

- PASCAL Challenge
 - Objects from 20 categories
 - person, car, bicycle, bus, airplane, sheep, cow, table, ...
 - Objects are annotated with bounding boxes



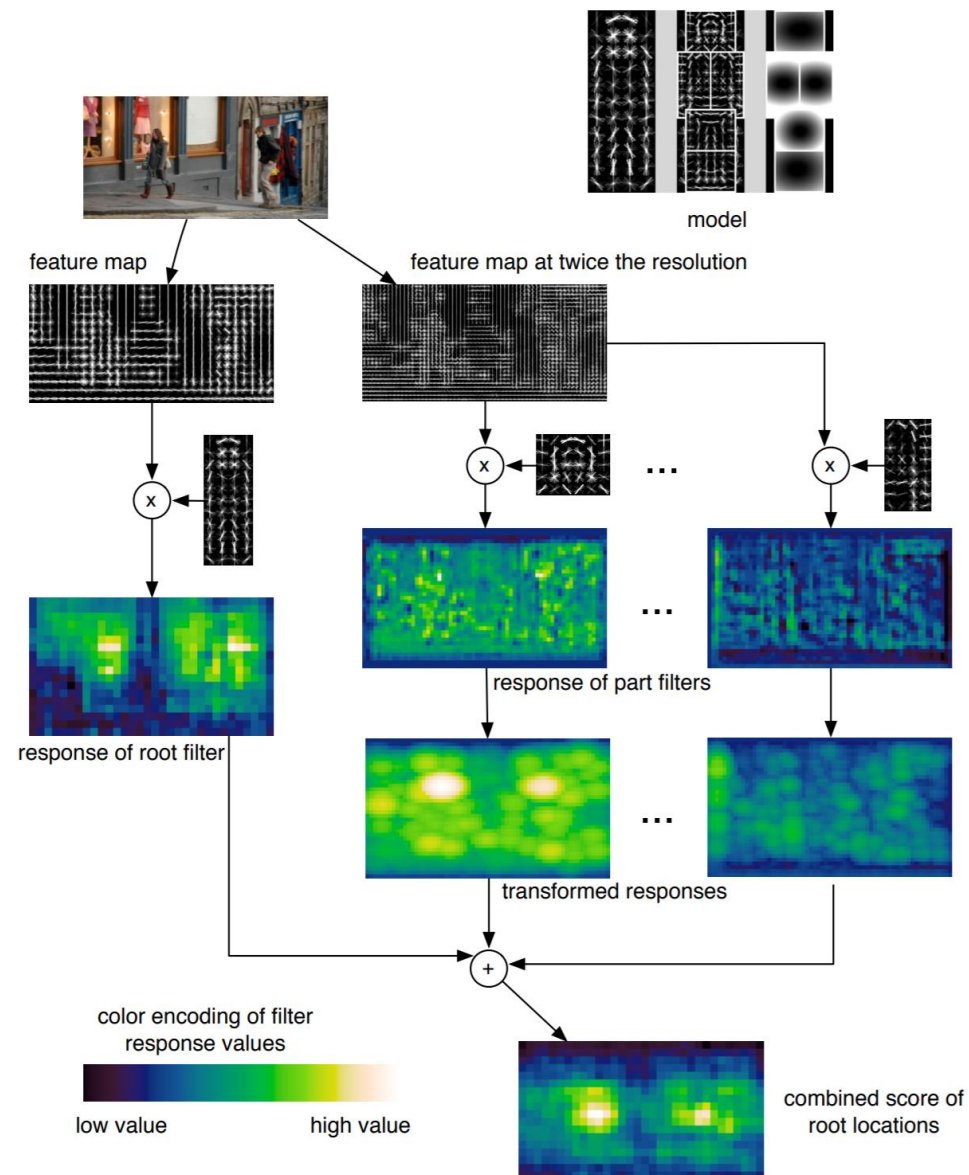
Deformable Part Models (DPM)

- Strategy
 - Sliding window approach, reduces object detection to binary classification
- Feature
 - Histogram of Gradient (HOG)
- Classifier
 - Linear SVMs



DPM

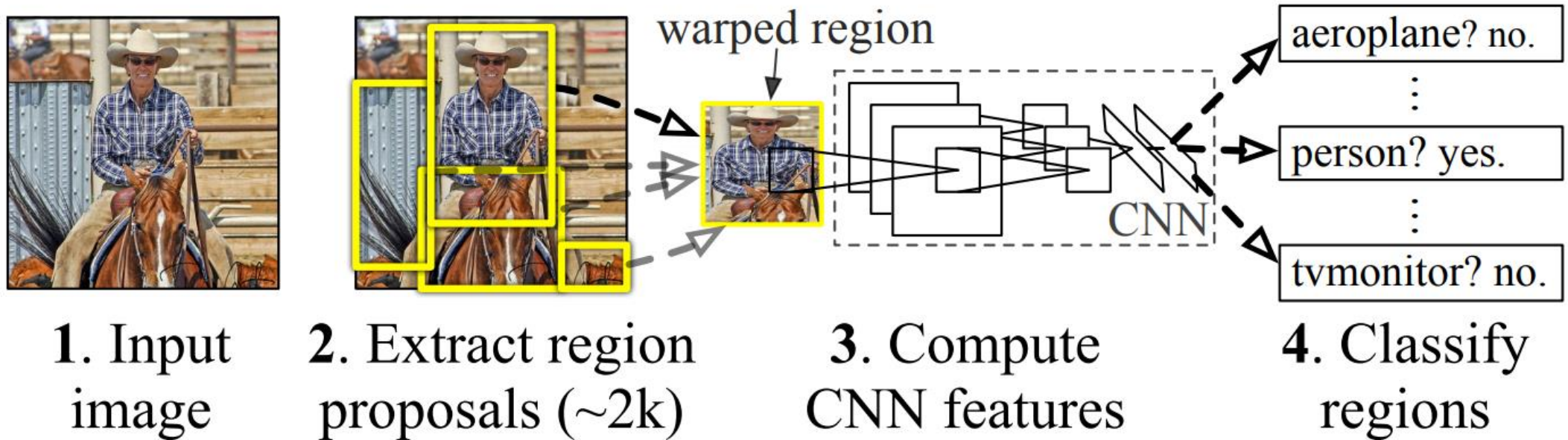
- Root filter + part filters



R-CNN

- Classify region proposals

R-CNN: *Regions with CNN features*



SVMs

R-CNN

- Generate region proposals
 - Selective search

Algorithm 1: Hierarchical Grouping Algorithm

Input: (colour) image

Output: Set of object location hypotheses L

Obtain initial regions $R = \{r_1, \dots, r_n\}$ using Felzenszwalb and Huttenlocher (2004) Initialise similarity set $S = \emptyset$;

foreach Neighbouring region pair (r_i, r_j) **do**

 Calculate similarity $s(r_i, r_j)$;

$S = S \cup s(r_i, r_j)$;

while $S \neq \emptyset$ **do**

 Get highest similarity $s(r_i, r_j) = \max(S)$;

 Merge corresponding regions $r_t = r_i \cup r_j$;

 Remove similarities regarding r_i : $S = S \setminus s(r_i, r_*)$;

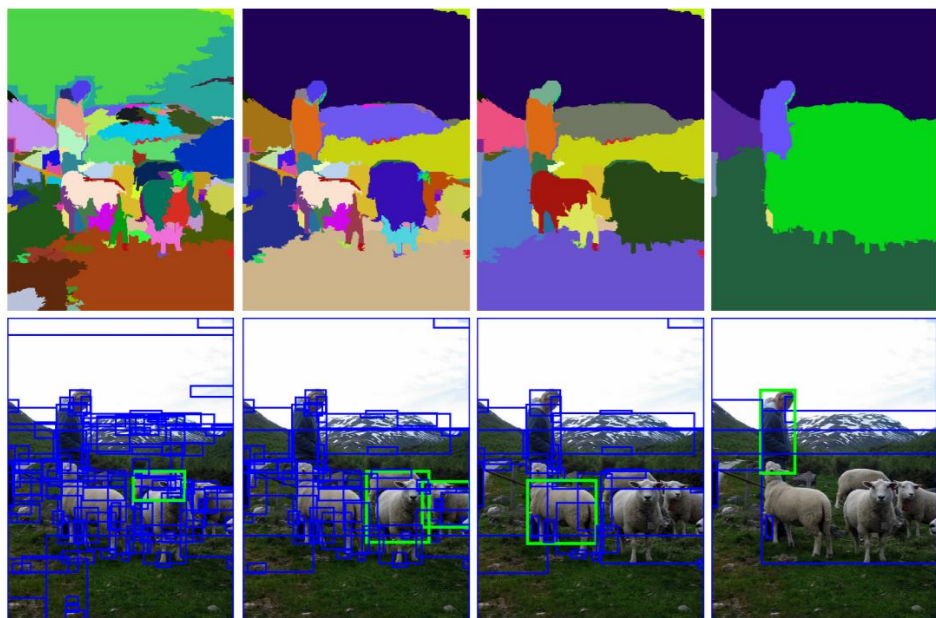
 Remove similarities regarding r_j : $S = S \setminus s(r_*, r_j)$;

 Calculate similarity set S_t between r_t and its neighbours;

$S = S \cup S_t$;

$R = R \cup r_t$;

Extract object location boxes L from all regions in R ;



(a)

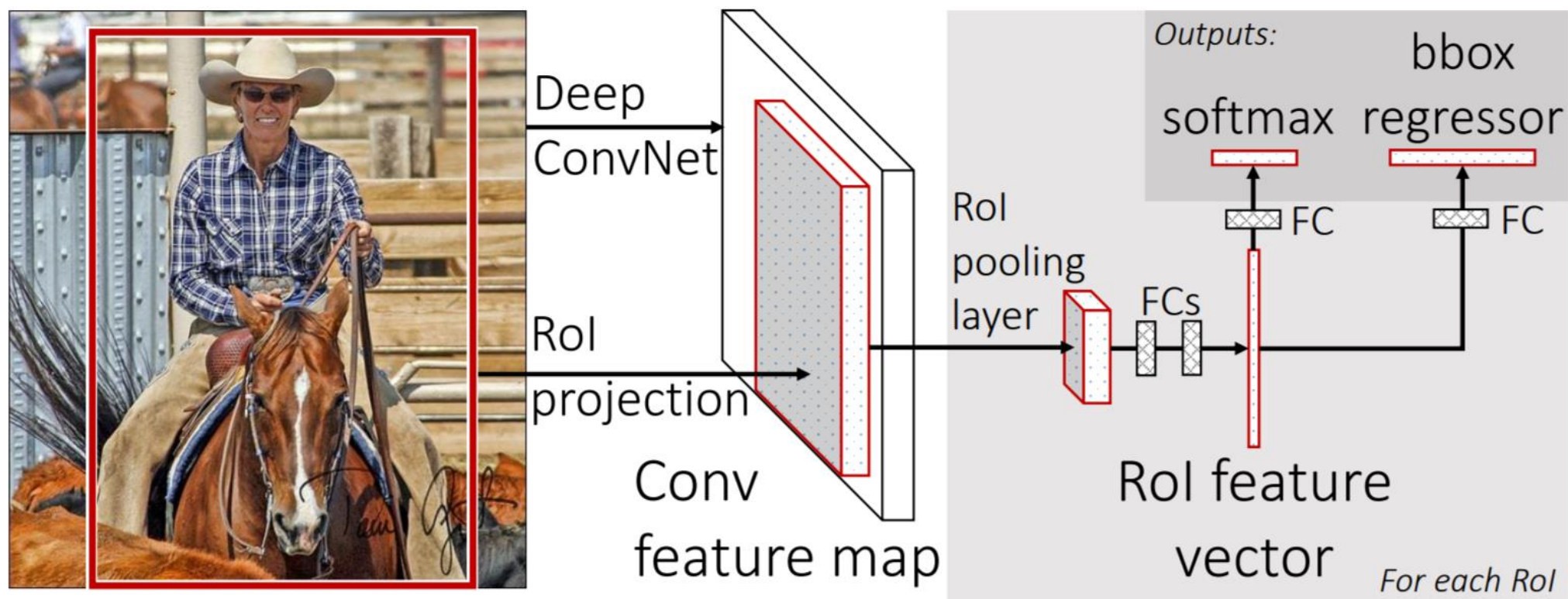


(b)

R-CNN

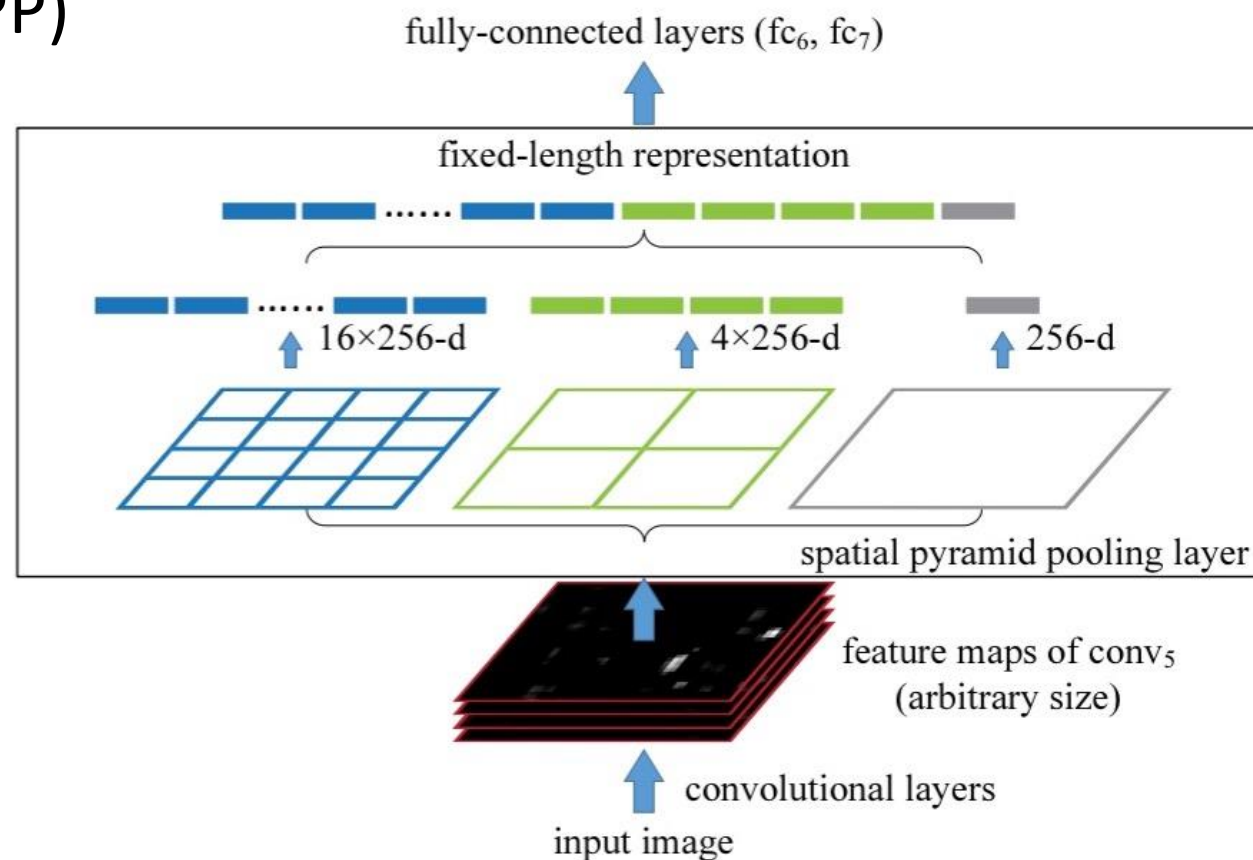
- Limitations
 - Training is a multi-stage pipeline
 - Training is expensive in space and time
 - Object detection is slow

Fast R-CNN



Fast R-CNN

- Spatial pyramid pooling (SPP)
- One-stage training

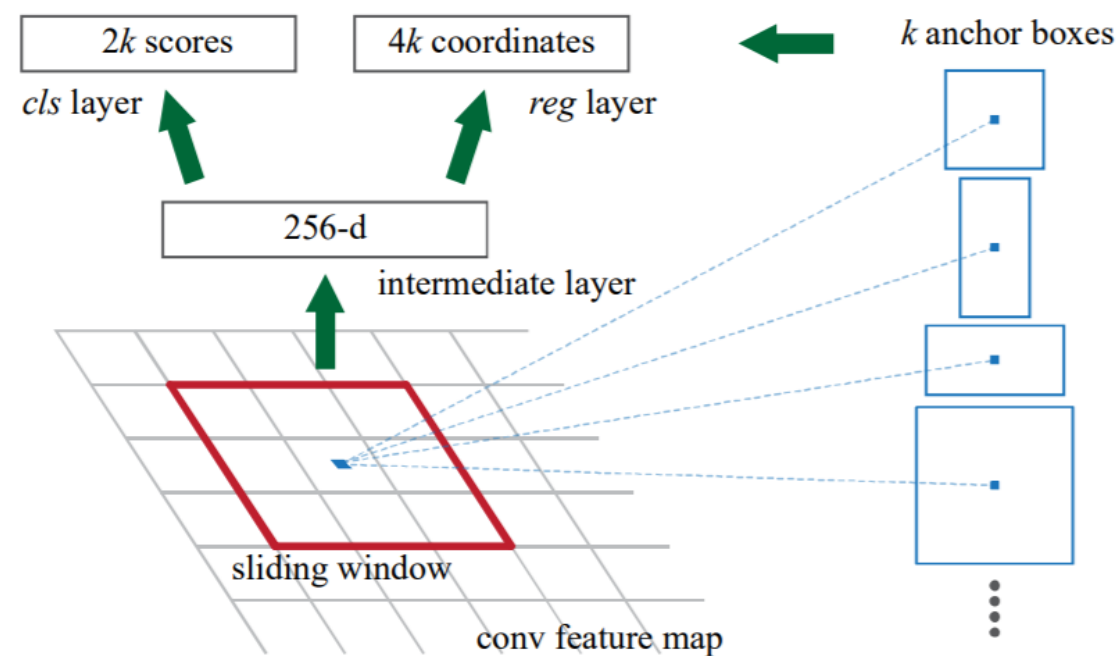
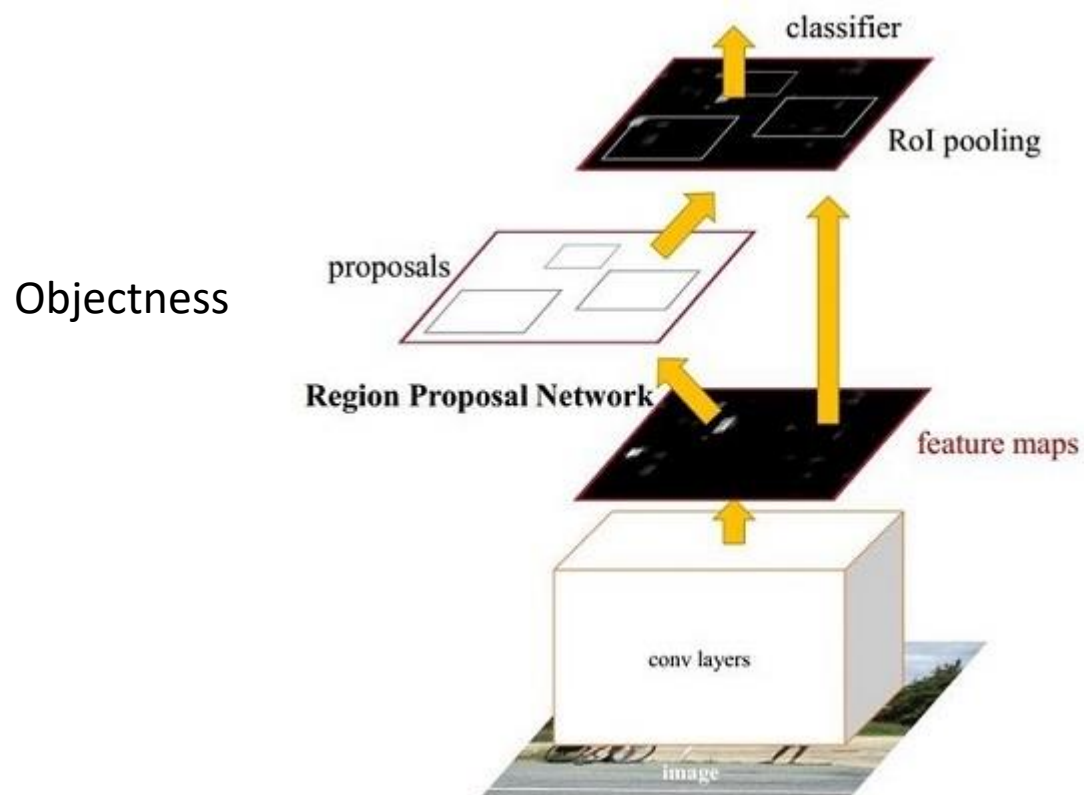


Fast R-CNN

- Fast R-CNN vs R-CNN
 - One-stage
 - No SVMs
 - No warp / resize
 - Fast
- Limitations
 - Region proposals are not efficient

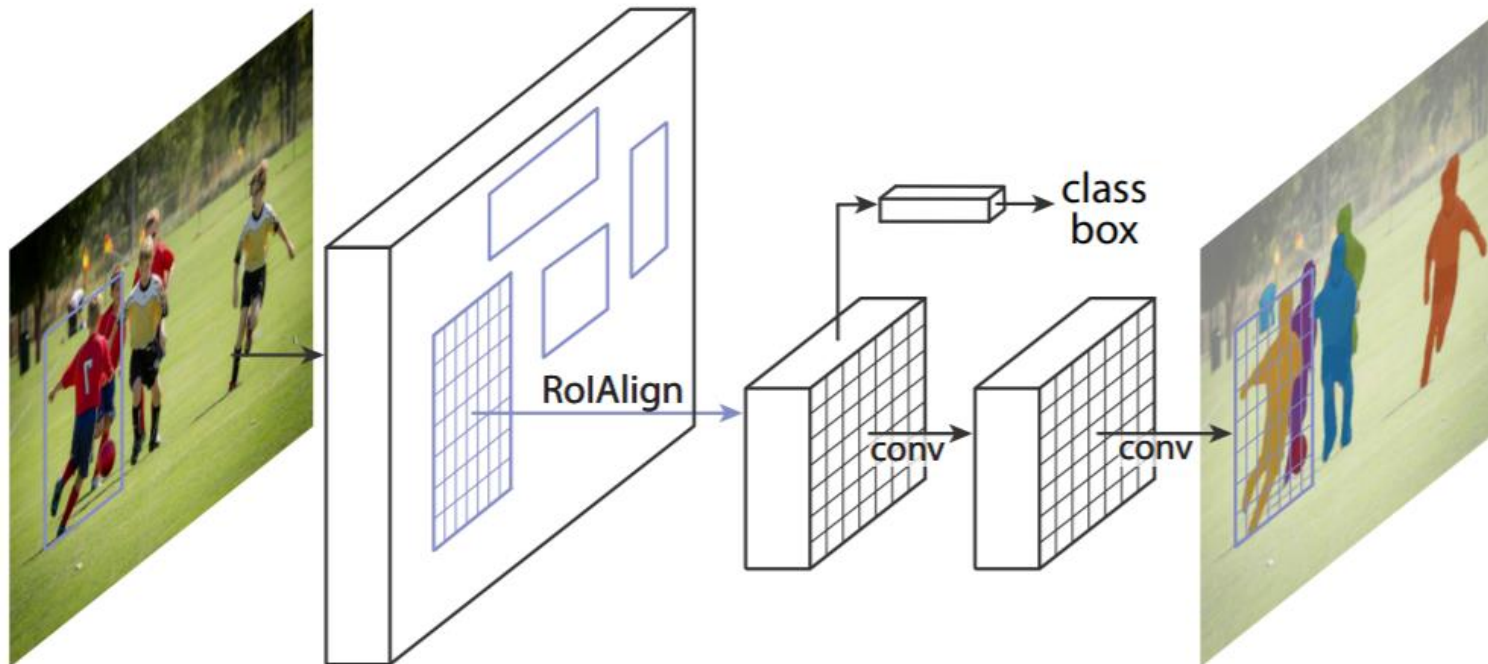
Faster R-CNN

- Region proposal network (RPN)



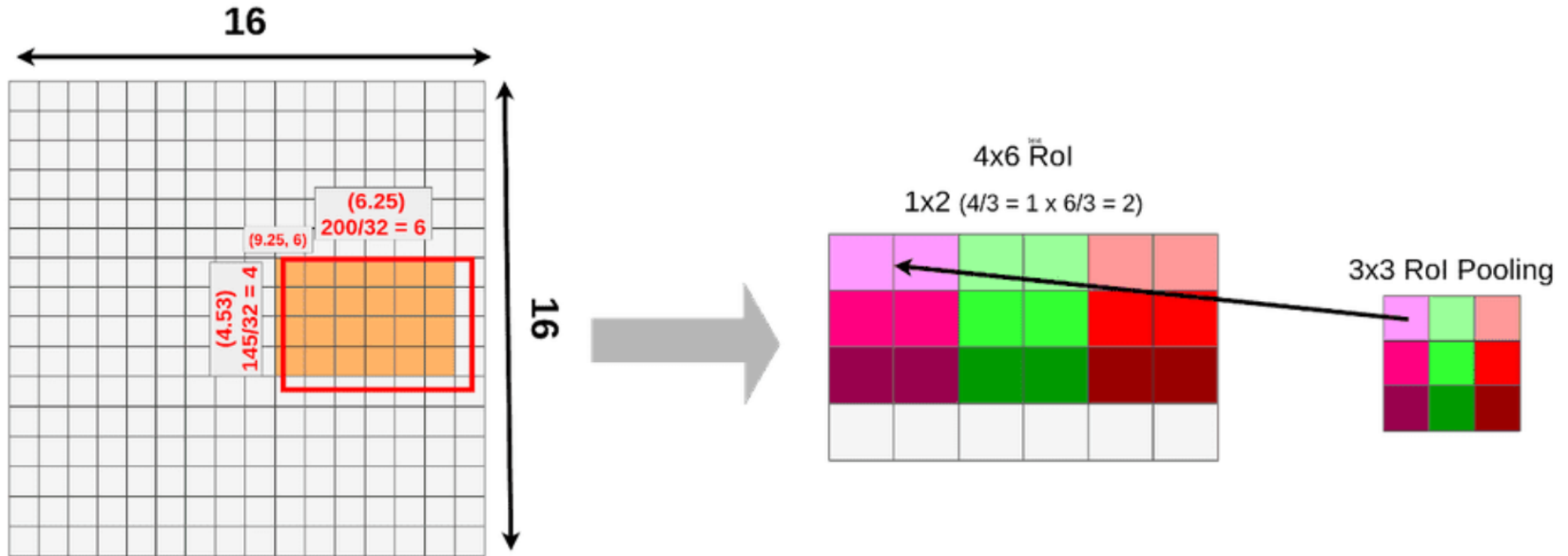
Mask R-CNN

- Mask prediction
- ROI align



Mask R-CNN

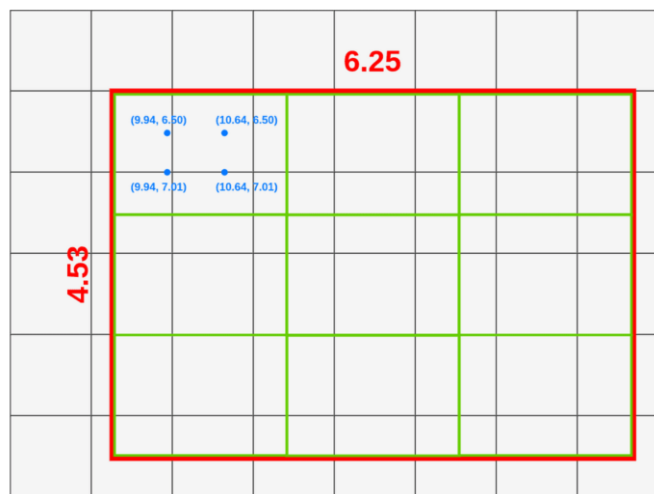
- **ROI pooling** vs ROI align



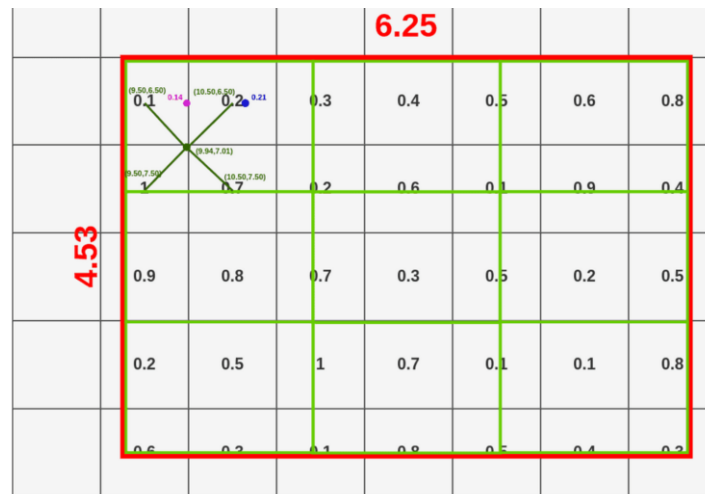
Quantization two times

Mask R-CNN

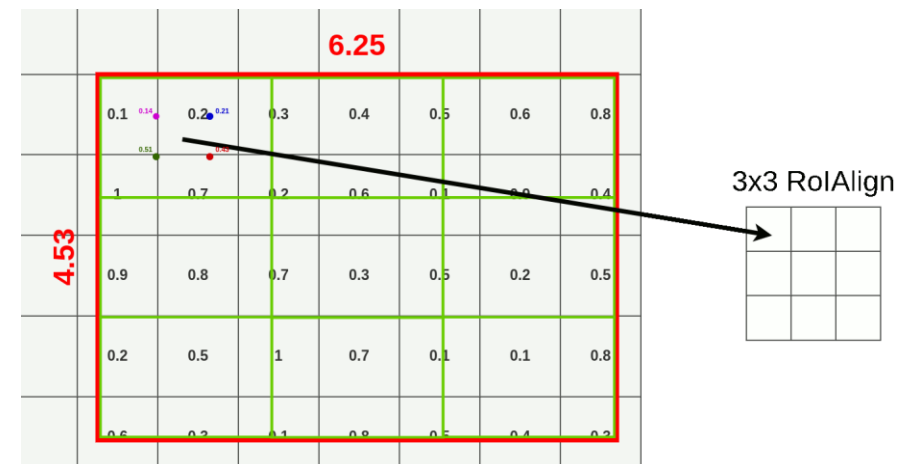
- ROI pooling vs **ROI align**



Sampling

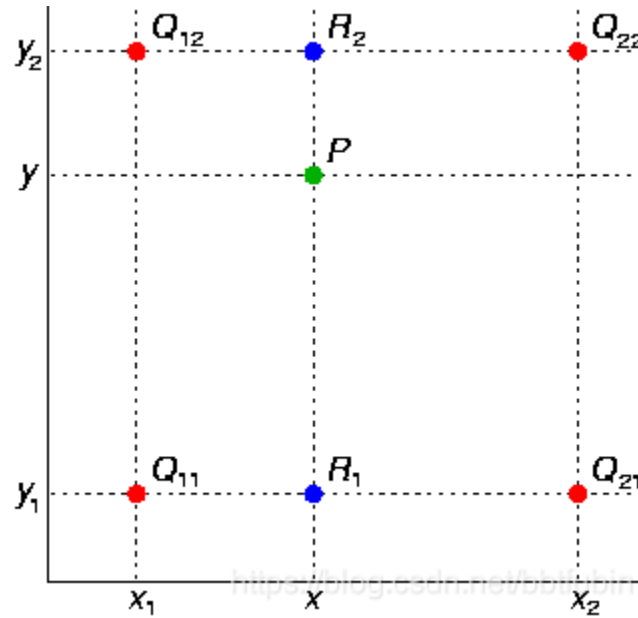


Bilinear interpolation



Pooling

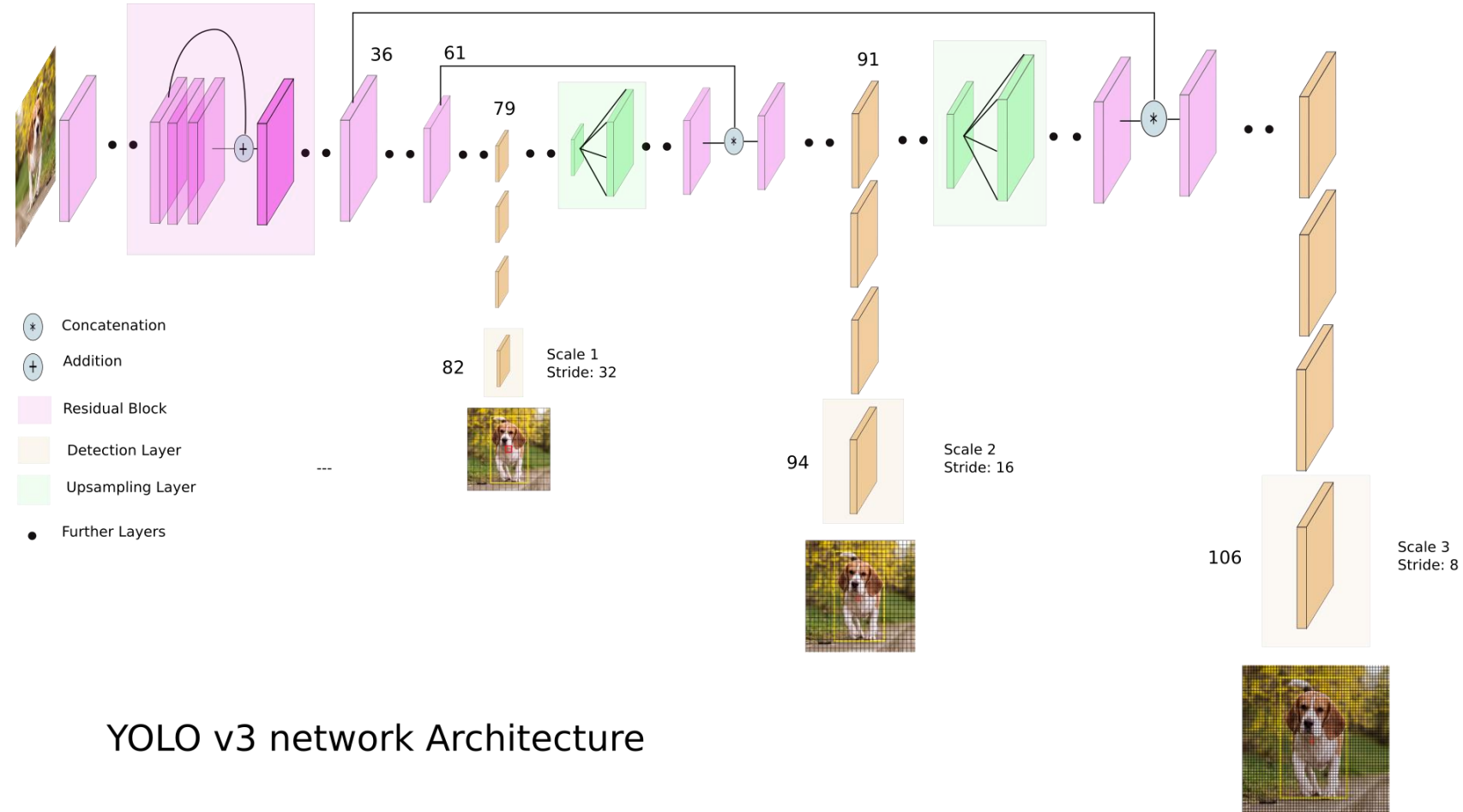
Bilinear Interpolation



$$f(x, y) \approx \frac{f(Q_{11})}{(x_2 - x_1)(y_2 - y_1)}(x_2 - x)(y_2 - y) + \frac{f(Q_{21})}{(x_2 - x_1)(y_2 - y_1)}(x - x_1)(y_2 - y) \\ + \frac{f(Q_{12})}{(x_2 - x_1)(y_2 - y_1)}(x_2 - x)(y - y_1) + \frac{f(Q_{22})}{(x_2 - x_1)(y_2 - y_1)}(x - x_1)(y - y_1).$$

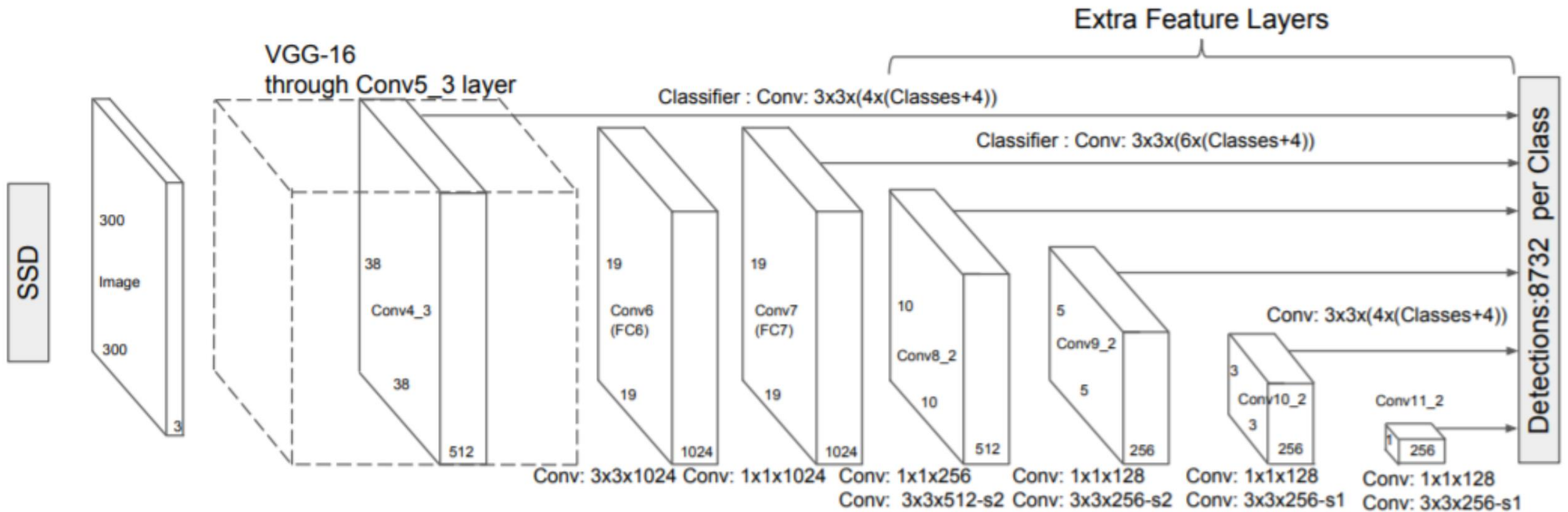
YOLO

- Example
 - Scale 1: 7×7 grid
 - Each grid: $(1 + \text{class} + 4) \times k$
- Objectness Bbox params Class number Anchors



SSD

- Multibox and multi-scale



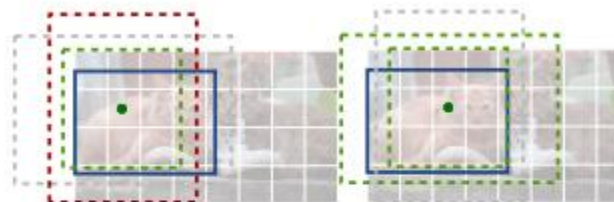
CenterNet

- Anchor vs anchor-free

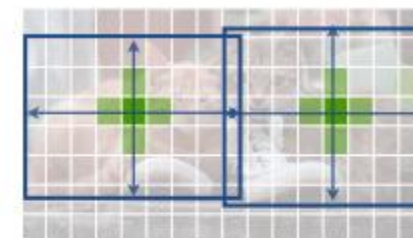
$$L_k = \frac{-1}{N} \sum_{xyc} \begin{cases} (1 - \hat{Y}_{xyc})^\alpha \log(\hat{Y}_{xyc}) & \text{if } Y_{xyc} = 1 \\ (1 - Y_{xyc})^\beta (\hat{Y}_{xyc})^\alpha \log(1 - \hat{Y}_{xyc}) & \text{otherwise} \end{cases}$$

$$L_{off} = \frac{1}{N} \sum_p \left| \hat{O}_{\tilde{p}} - \left(\frac{p}{R} - \tilde{p} \right) \right|$$

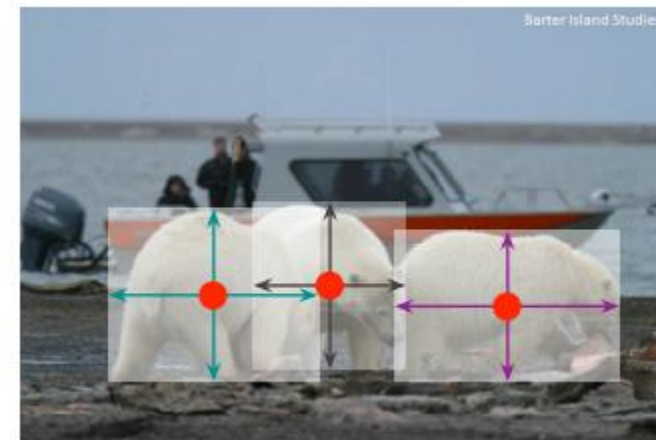
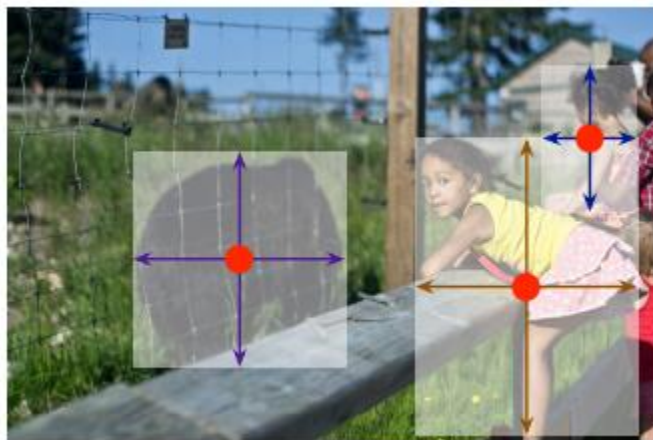
$$L_{size} = \frac{1}{N} \sum_{k=1}^N \left| \hat{S}_{p_k} - s_k \right|$$



(a) Standard anchor based detection. Anchors count as **positive** with an overlap $IoU > 0.7$ to any **object**, **negative** with an overlap $IoU < 0.3$, or are **ignored** otherwise.

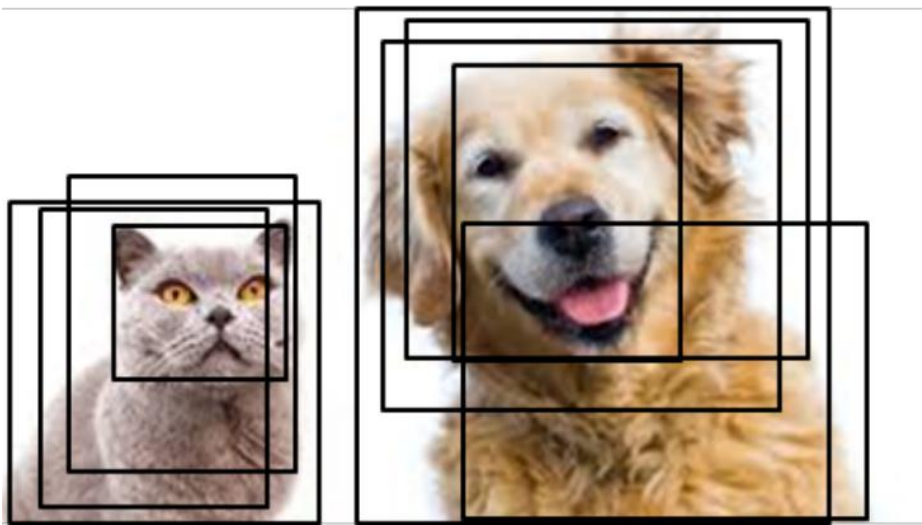


(b) Center point based detection. The **center pixel** is assigned to the **object**. Nearby points have a reduced negative loss. Object size is regressed.



NMS

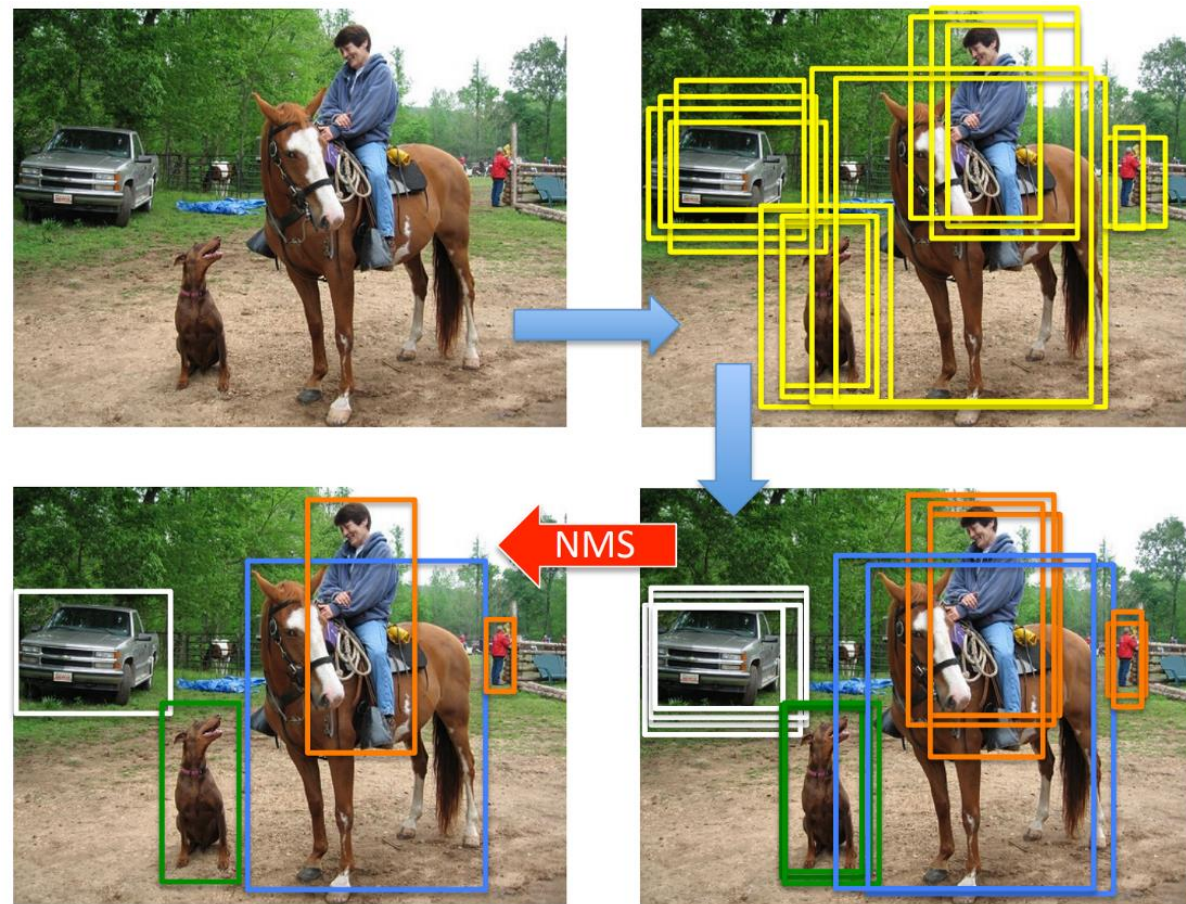
- What left?



NMS

Algorithm 1 Non-Max Suppression

```
1: procedure NMS( $B, c$ )
2:    $B_{nms} \leftarrow \emptyset$ 
3:   for  $b_i \in B$  do
4:      $discard \leftarrow \text{False}$ 
5:     for  $b_j \in B$  do
6:       if  $\text{same}(b_i, b_j) > \lambda_{nms}$  then
7:         if  $\text{score}(c, b_j) > \text{score}(c, b_i)$  then
8:            $discard \leftarrow \text{True}$ 
9:       if not  $discard$  then
10:         $B_{nms} \leftarrow B_{nms} \cup b_i$ 
11:   return  $B_{nms}$ 
```



Improve Your Detector Performance

- Add extra training data
- Use appropriate augmentation strategies
- Change the backbone architecture
- Modify the classification and regression methods
- Use proper loss functions
- ...