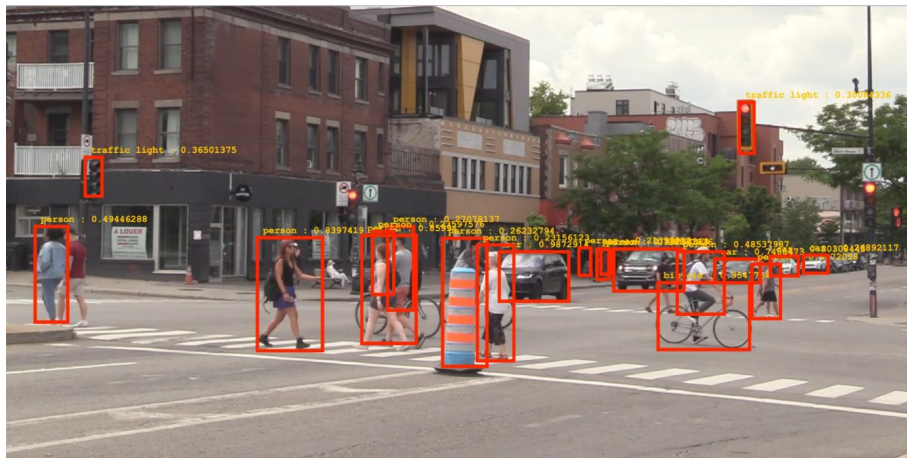


Object Detection and Segmentation

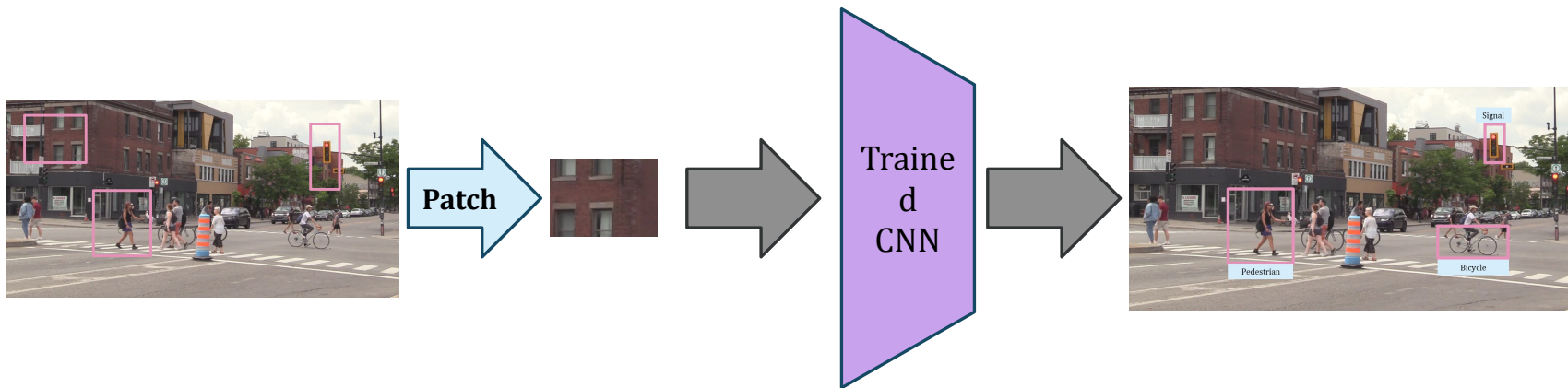
What is Object Detection?

- Input: An Image
- Output: A set of detected objects, each with
 - A class label (from a set of pre-defined class labels)
 - A bounding box
 - May be of the form $[x_{left-top}, y_{left-top}, height, width]$



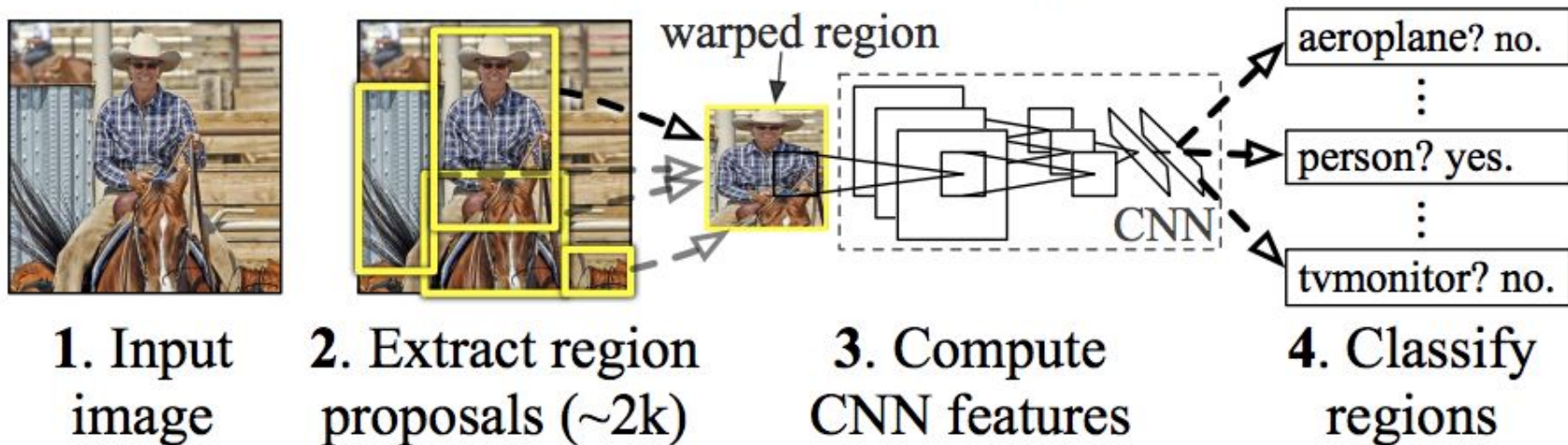
Object Detection using CNNs: A Simple Approach

- Pass the different sized patches through the trained CNN



R-CNN

R-CNN: *Regions with CNN features*



Identifying Potential Regions

- R-CNN can use different methods to identify potential regions
 - Selective search
 - Objectness
 - Category-independent object proposals, etc.
- In the original paper, the authors use selective search

Region Proposal - Selective Search



Region Proposal - Selective Search



Input Image



**After Initial
Segmentation**



**After few
Iterations**



**After many
iterations**

Region Proposal - Selective Search



Input Image



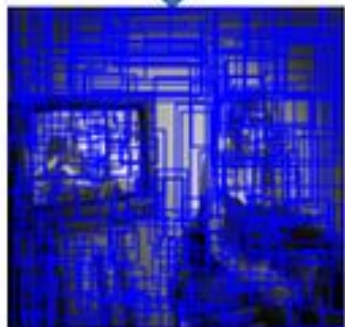
After Initial
Segmentation



After few
Iterations



After many
iterations



Region Proposal - Selective Search

- Color Similarity
- Texture Similarity
- Size Similarity
- Shape Compatibility

Selective Search

- Step 1: Use graph-based image segmentation to create initial regions
- Step 2: Extract features from each region that represents the characteristics of the region
- Step 3: Iteratively group the regions together based on the similarity of features
 - Step 3.1: Two most similar neighbour regions are grouped together
 - Step 3.2: The features of the grouped region is calculated
 - Step 3.3: The similarity of the grouped regions with its neighbours are calculated
 - Step 3.4: Go to step 3.1 if there is a pair of neighbouring regions having a feature similarity score more than a threshold. Stop, otherwise.



Feature of the regions could be colour distribution, texture, or other features, such as HOG

Feature Extraction

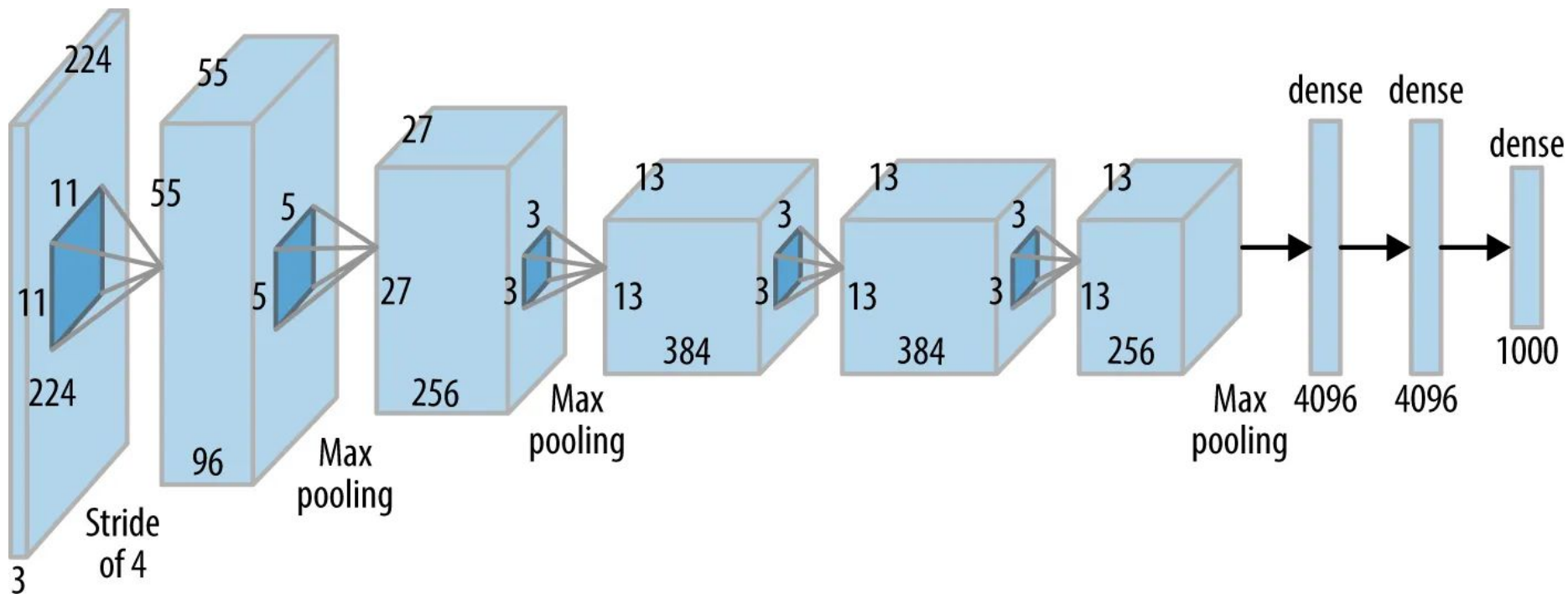


Warping



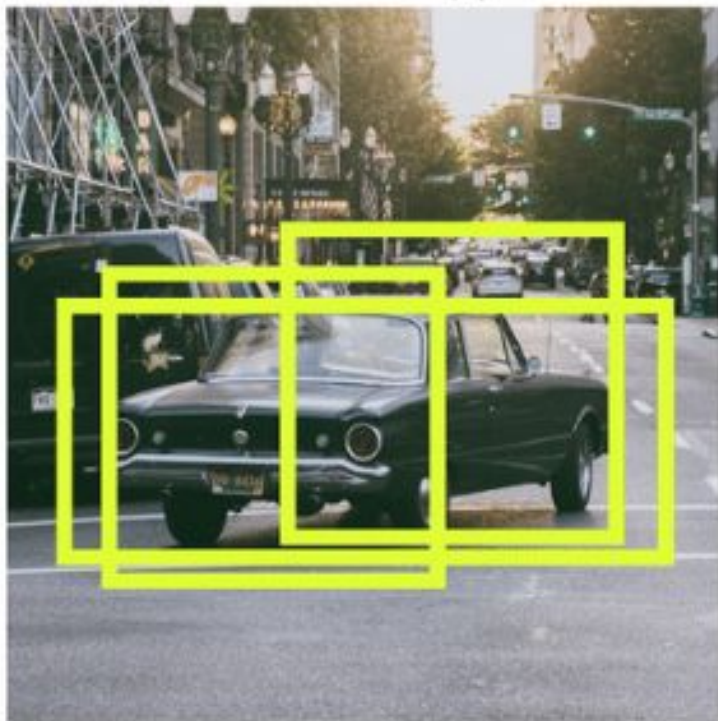
AlexNet

AlexNet



Non-Maximal Suppression

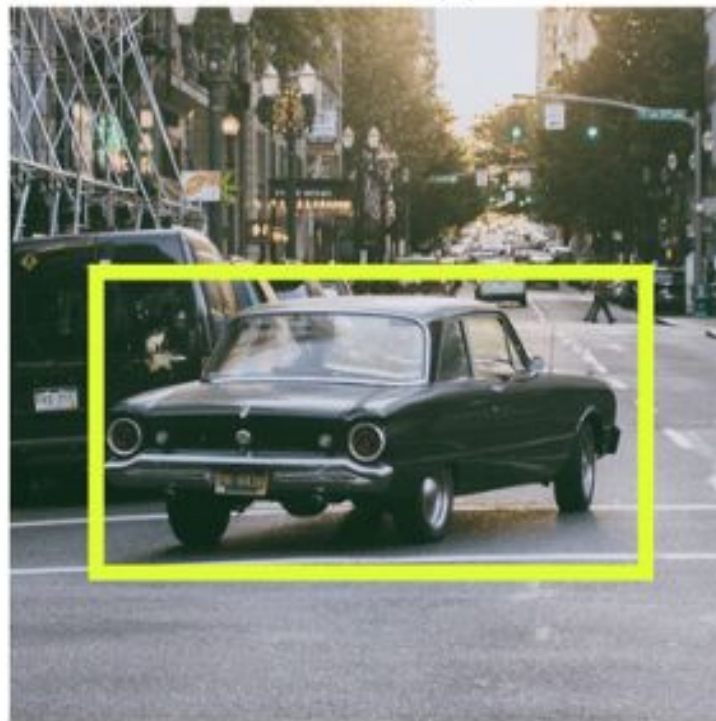
Before non-max suppression



Non-Max
Suppression



After non-max suppression



Improving the Performance: Non Maximal Suppression (NMS)

- Take the most confident bounding boxes among all the bounding boxes



Improving the Performance: Non Maximal Suppression (NMS)

- Take the most confident bounding boxes among all the bounding boxes



Improving the Performance: Non Maximal Suppression (NMS)

- Find out all the bounding boxes that overlap with the most confident box by a certain amount



Improving the Performance: Non Maximal Suppression (NMS)

- Find out all the bounding boxes that overlap with the most confident box by a certain amount
- Remove all such boxes



Improving the Performance: Non Maximal Suppression (NMS)

- Find the next most confident box

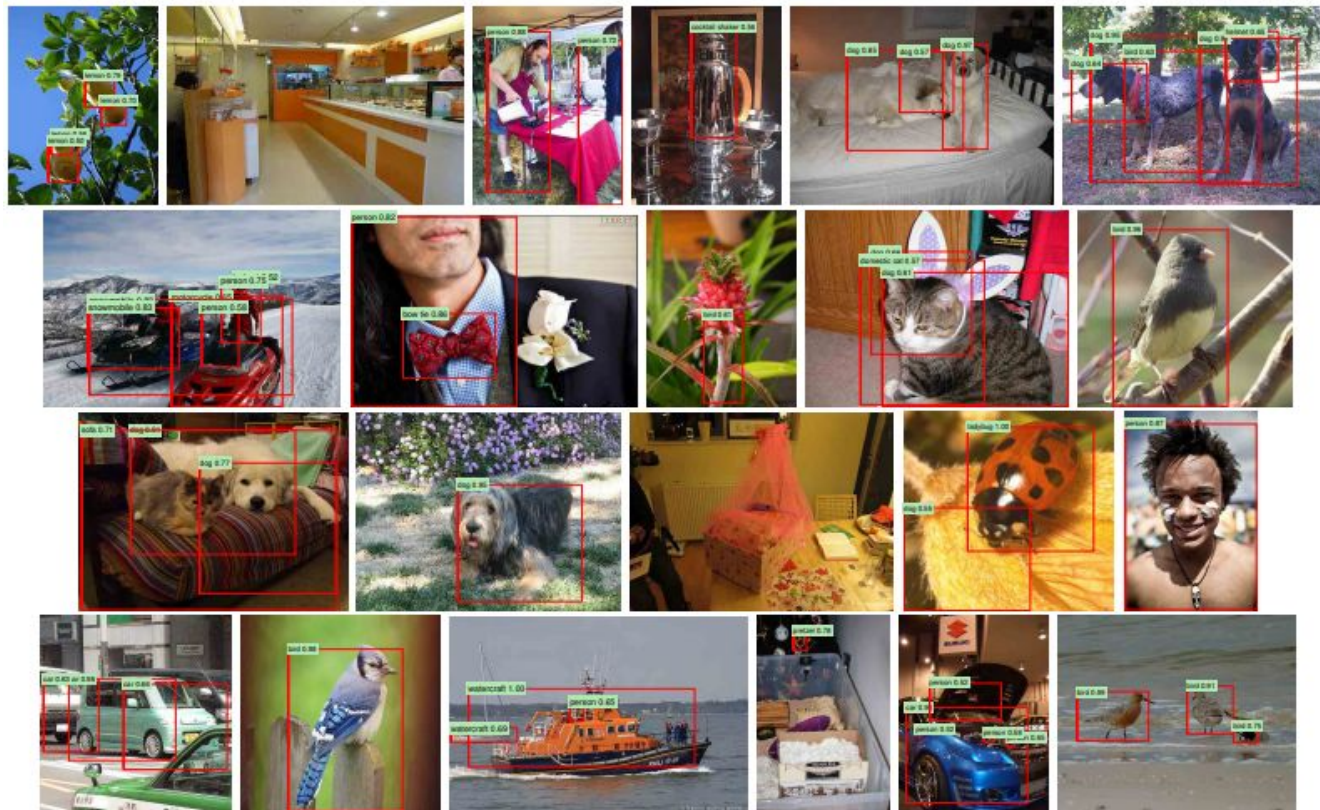


Improving the Performance: Non Maximal Suppression (NMS)

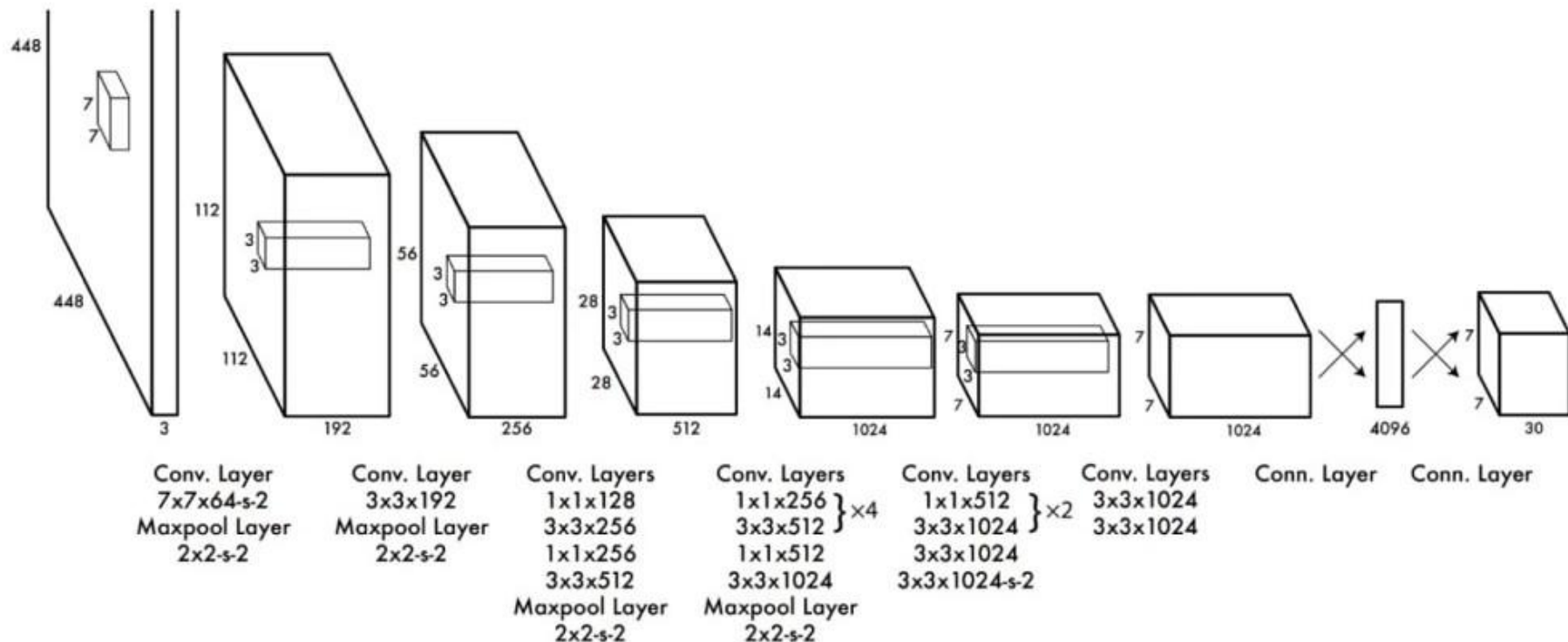
- Find the next most confident box
- Remove the overlapping boxes



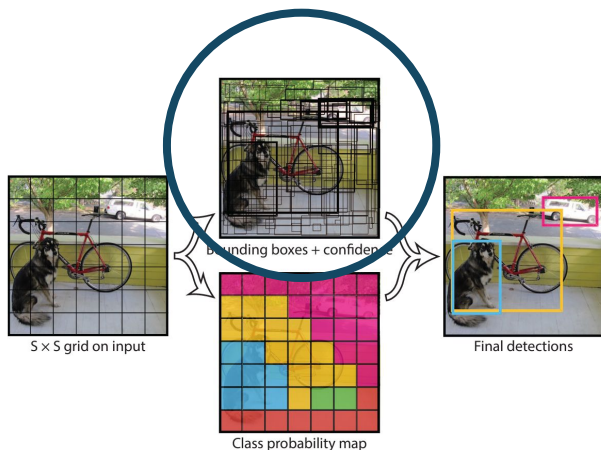
Examples



YOLO - You Only Look Once



YOLO: You Only Look Once



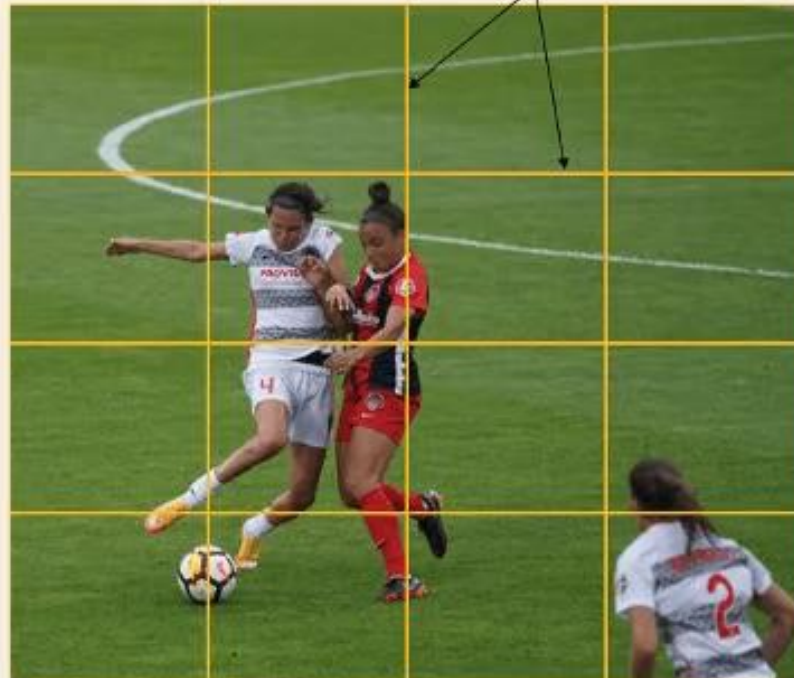
- Input image divided into $S \times S$ grid
- For each grid cell, the model predicts B bounding boxes
- Each bounding box has 5 predictions
 - x, y, w, h
 - x, y : represent the center of the box relative to the bounds of the grid cell
 - w, h : predicted relative to the whole image
- Confidence score: reflect how confident the model is that the box contains an object and also how accurate it thinks the box is that it predicts
 - $P(object) \times IOU_{pred}^{gt}$
 - If no object exists, confidence score should be zero

Residual Blocks

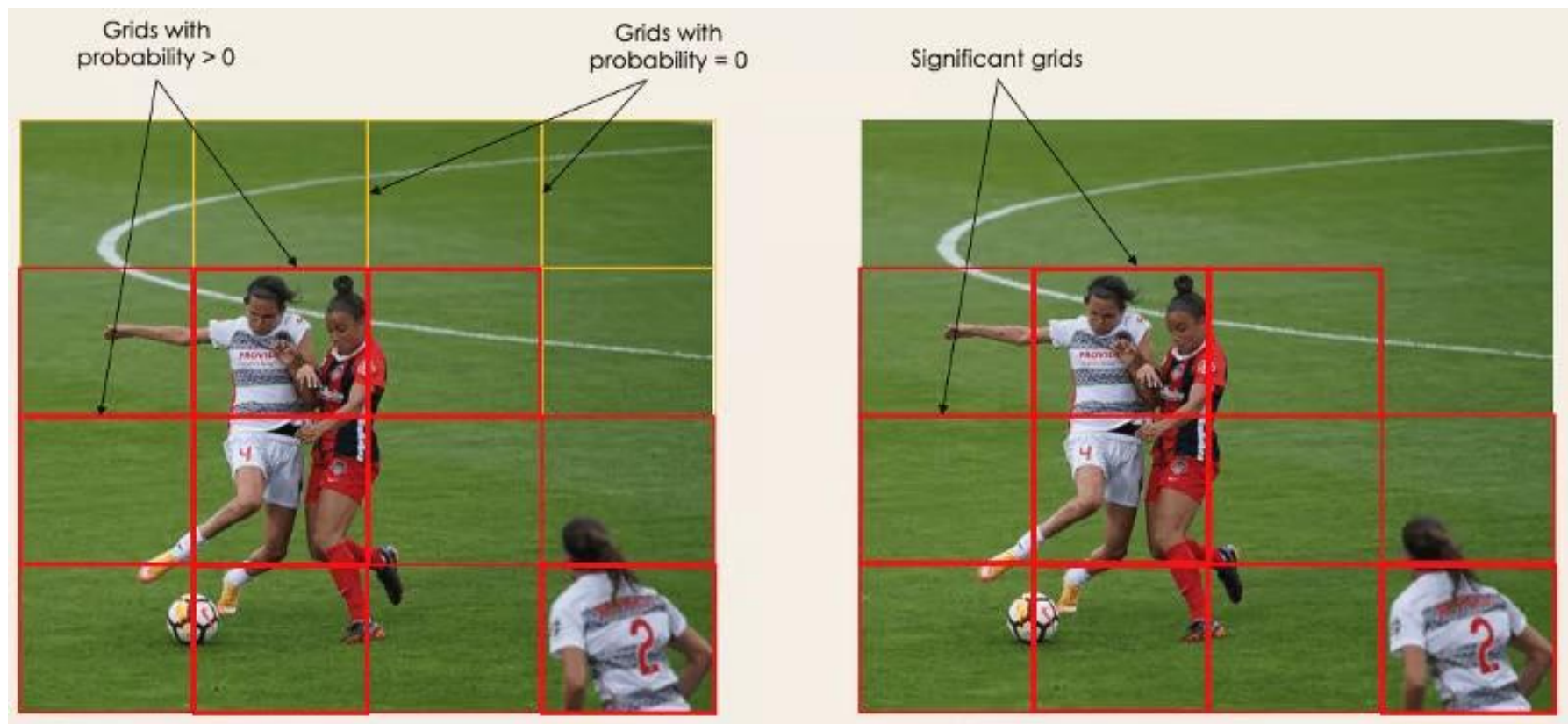
Original input Image



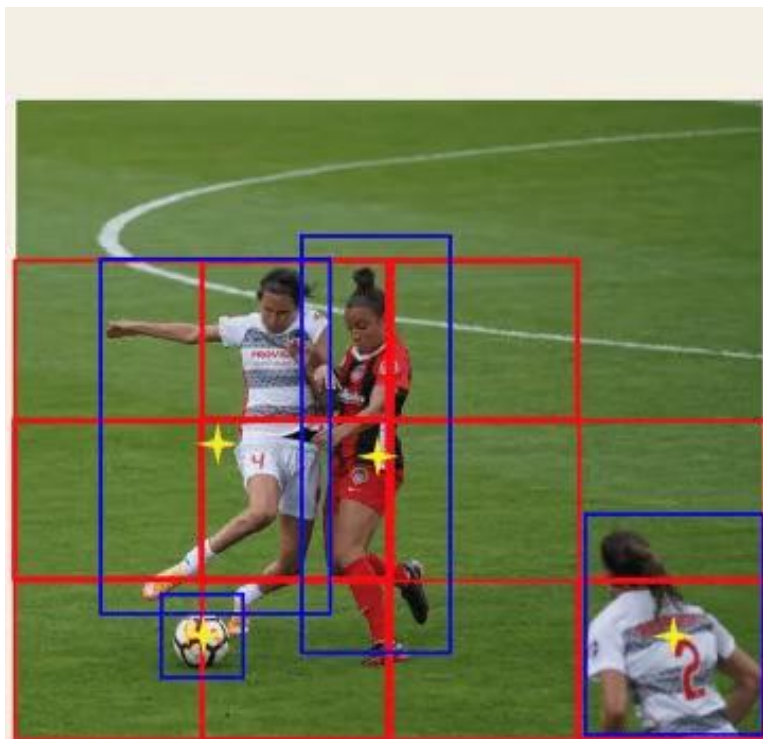
4x4 grid cells



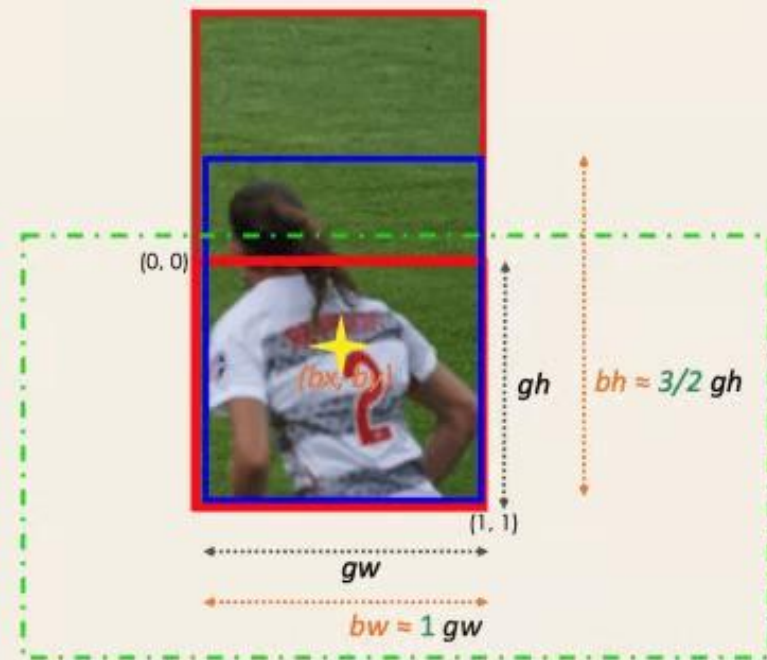
Bounding box regression



Bounding box regression



✦ Bounding box centers



From the previous info we can have for e.g.

$Y = [1, bx, by, 3/2, 1, c1, c2]$

- First 1 means 100% of object presence

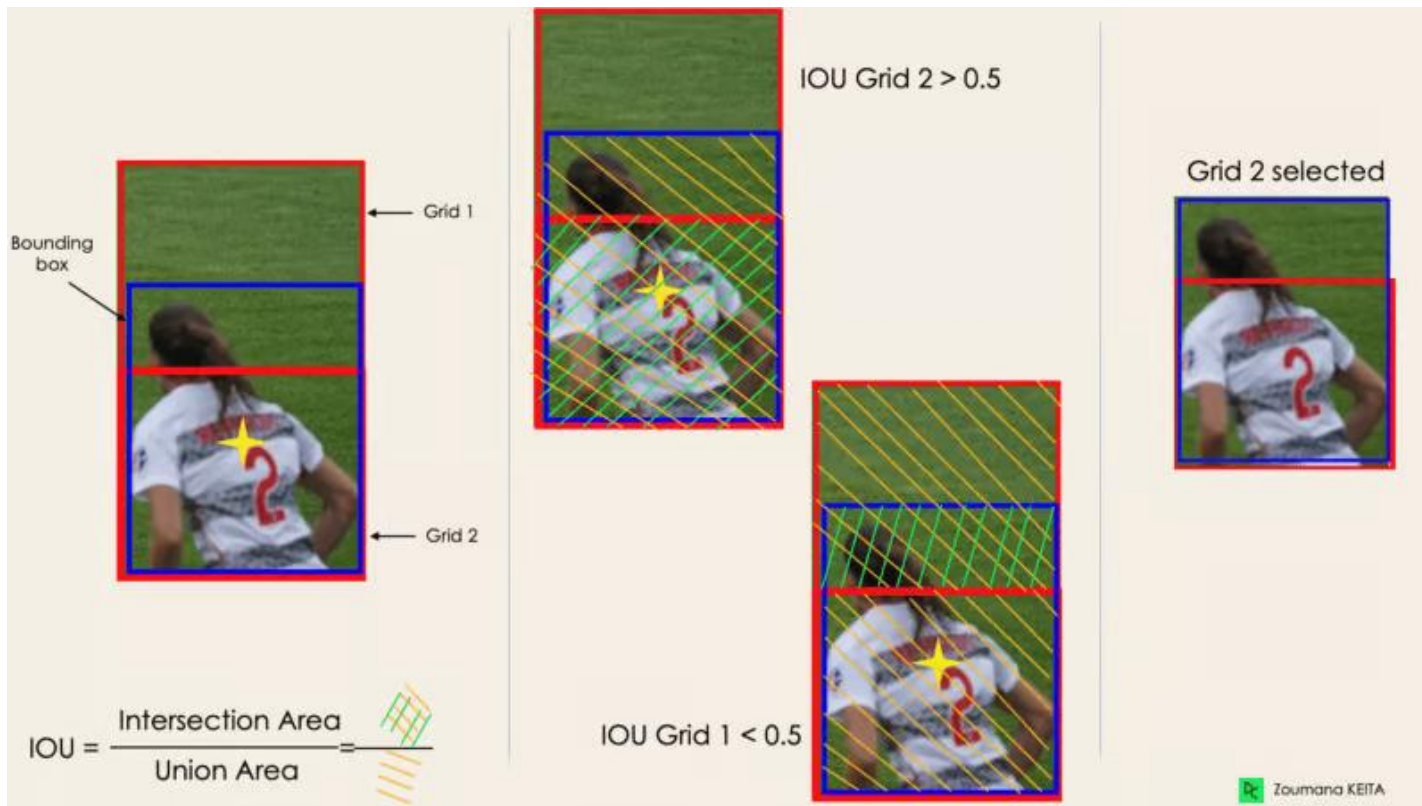
- gh, gw : height & width of the grid

- $0 \leq bx \leq 1$

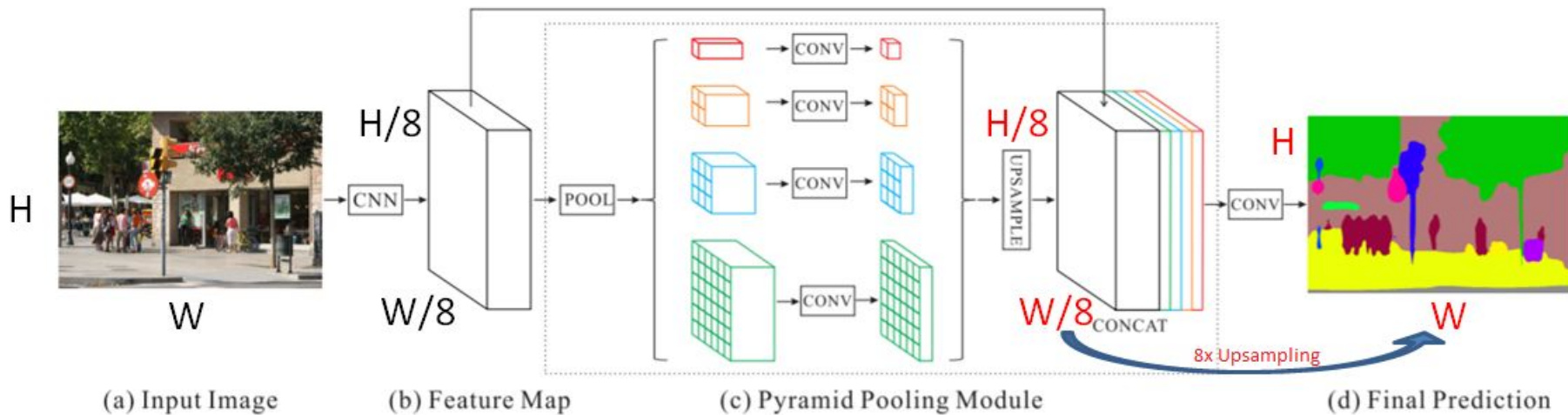
- $0 \leq by \leq 1$

- bh and bw can be more than 1

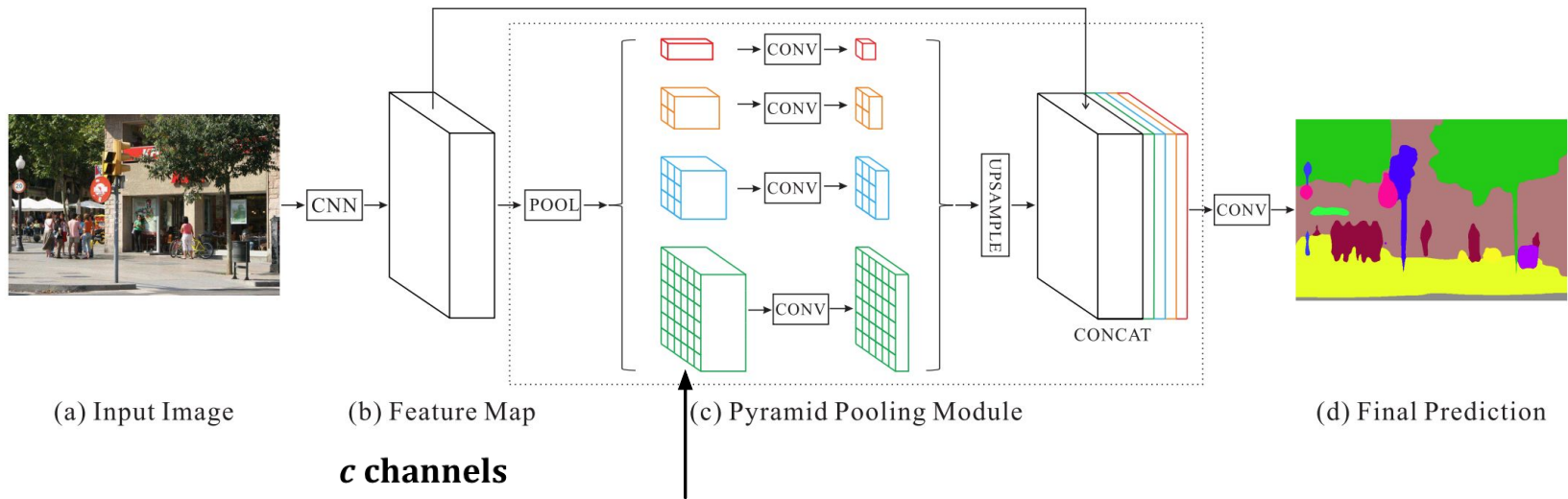
Intersection Over Unions or IOU



PSPNet



PSP Net: Pyramid Scene Parsing Network



Pooling at different scales

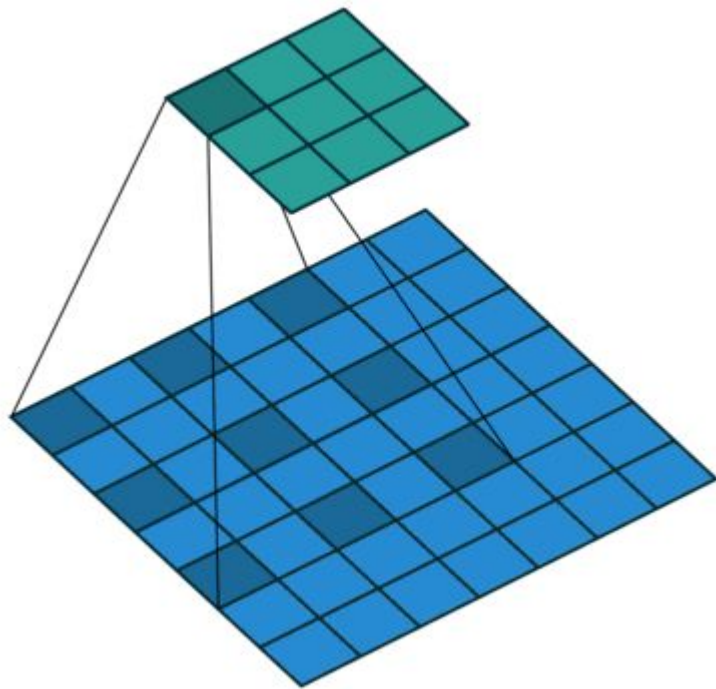
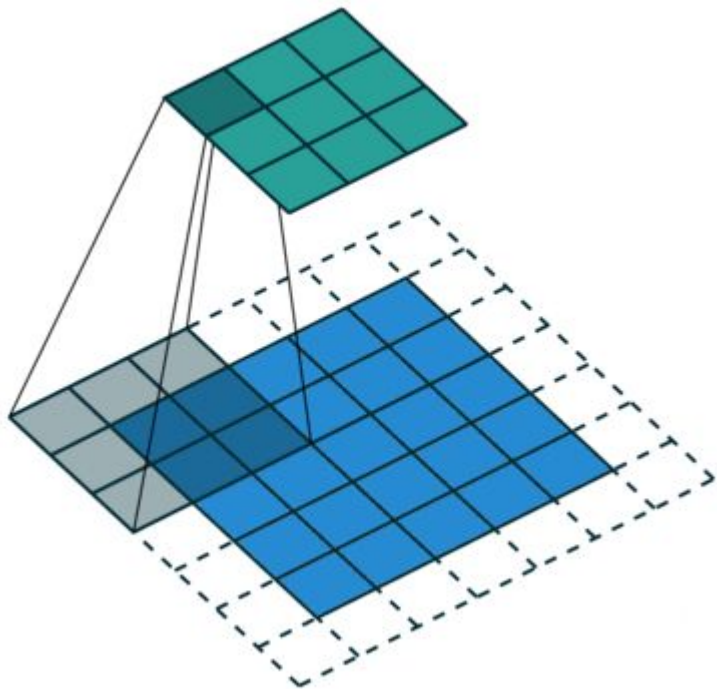
Global average pooling ($1 \times 1 \times c$ output maps)

Pooling that results in $2 \times 2 \times c$ output maps

Pooling that results in $3 \times 3 \times c$ output maps

Pooling that results in $6 \times 6 \times c$ output maps

Dilated Convolutions



SegNet

