

Object Detection consists of two sub tasks (i) Object Localization (ii) Object Classification.

For implementation we need to break it into three steps instead -

[A] Region Proposal Generation. Each of them called Region of Interest (RoI).

[B] Extracting feature maps for each region proposal. Lets call it Fmap RoI.

[C] Classifying each Fmap RoI as object classes.

Here step [B] is necessary because the images or crop from a image requires to be broken down to feature space which we are calling feature maps. Each of the convolutional filter can be responsible for extracting different image features. Like convolutional filters responsible for edge filtering can result in edge images kind of the the canny detector output.

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

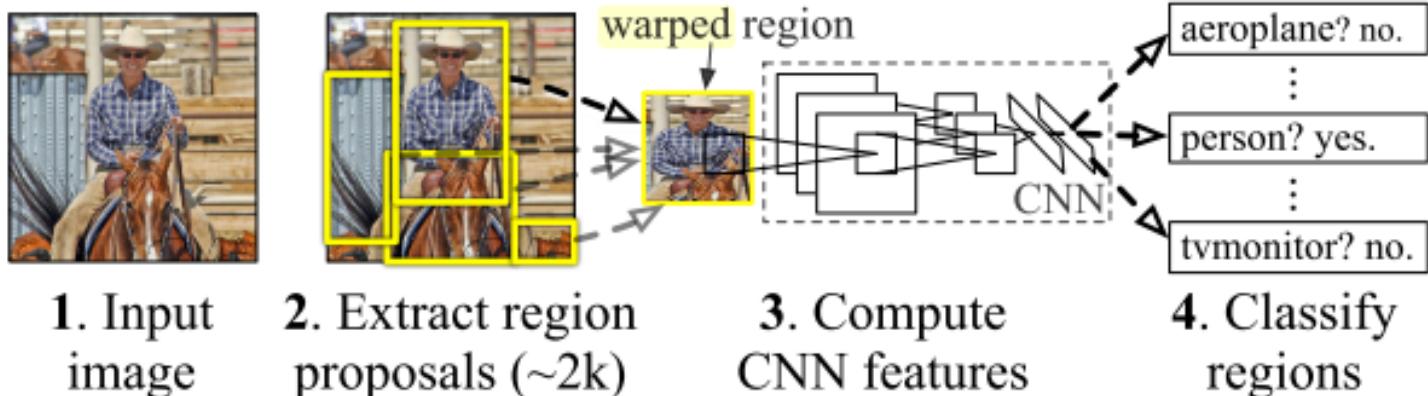


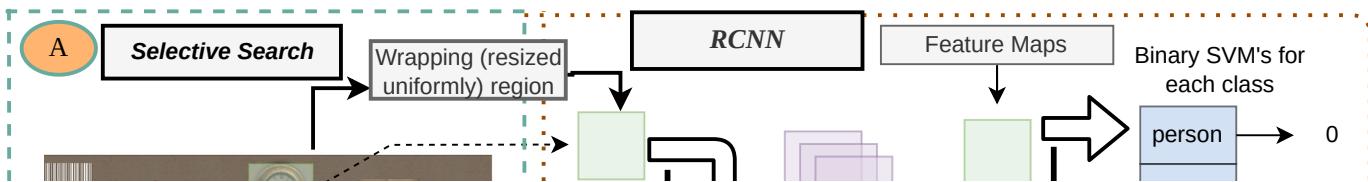
Figure : From RCNN paper

### The Algorithm :

- [A] Selective Search with 2k proposals per image. [Objectness[1], Selective Search[2]]
- [A] Wrapping the images [resizing with bilinear/nearest sampling should work here]
- [B] Using a CNN based classifier network for feature extraction
  - > Pre training with ICSLV2012 images with 200 classes lr 0.01 on Caffe implementation of a CNN based architecture.
  - > Fine tuning with lr 0.001 with wrapped region proposal with N + 1 class. N is no of object class 1 is background
  - > Oversampling positive classes 32 positive sample + 96 neg sample per mini-batch
  - > During finetuning a region proposal with more than 0.50 IoU with a positive sample was taken as positive.
  - > During Feature Extraction the last layer is removed and features from an internal dense layer is collected. i.e 1024 or 2048 features
- [C] Using binary classifier for per class classification from the extracted features
  - > One classifier is used for one class
  - > Hard negative mining [4] was very useful for fast convergence. It was possible to achieve results after just one pass over the data
- [\*] Fine tuning of localization was done by bbox regression using Ridge regressor with lambda = 1000 [detail below]

Regarding using the outputs of the classifiers after fine tuning instead of additionally using SVM :

Quote from the paper's appendix B : "We tried this and found that performance on VOC 2007 dropped from 54.2% to 50.9% mAP. This performance drop likely arises from a combination of several factors including that the definition of positive examples used in fine-tuning not emphasize precise localization and the softmax classifier was trained on randomly sampled negative examples rather than on the sub "hard negatives" used for SVM training."





res

1.

g does  
subset of

Ground Truth - Gx , Gy, Gw, Gh

Predictions - Px, Py, Pw, Pz

$$G_x = P_w d_x(P) + P_x$$

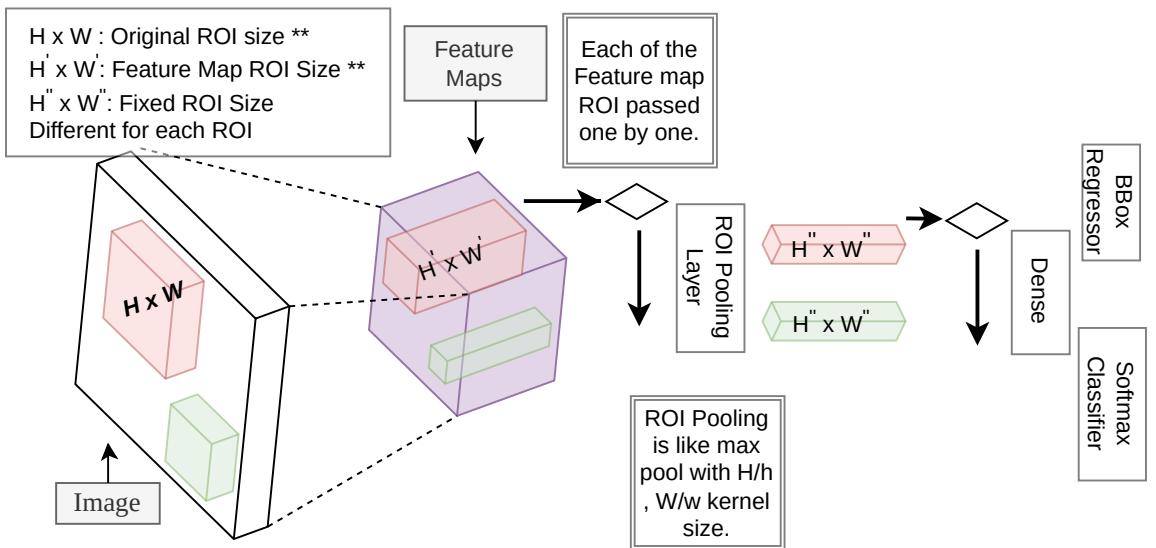
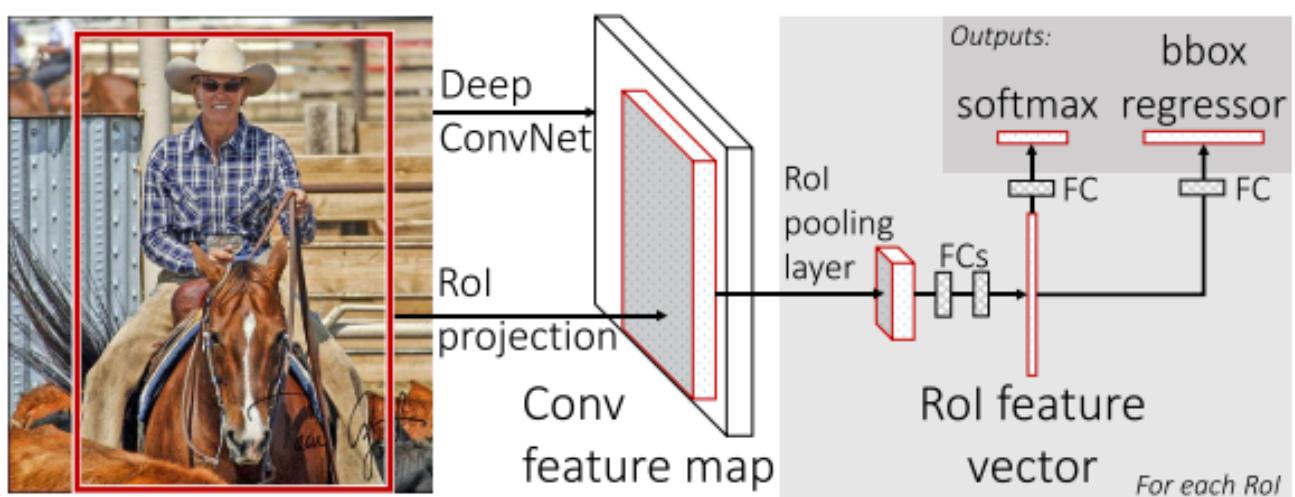
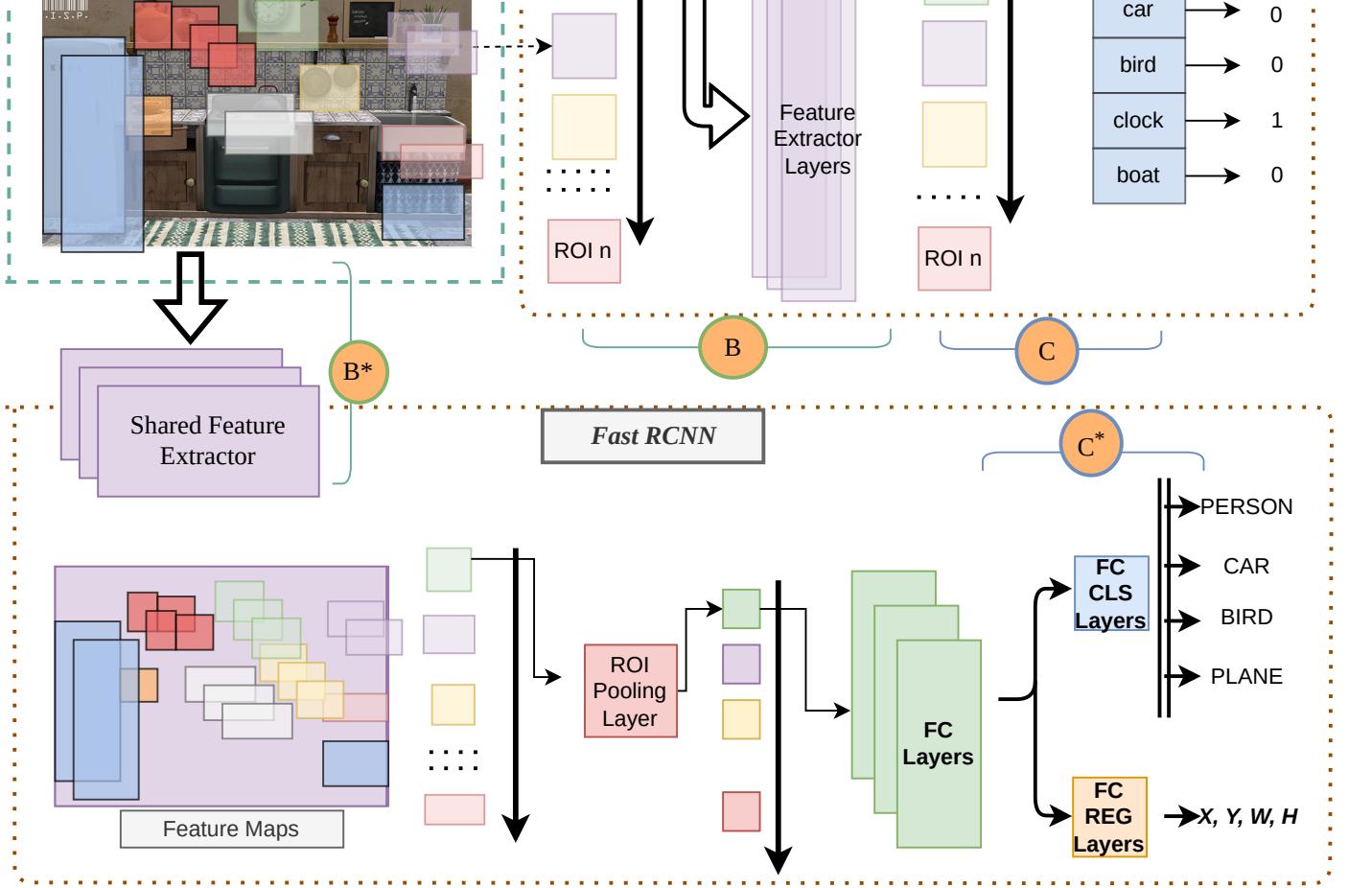
$$G_y = P_h d_y(P) + P_y$$

$$G_w = P_w \exp(d_m(P))$$











$$G_h = P_h \exp(d_h(P)).$$

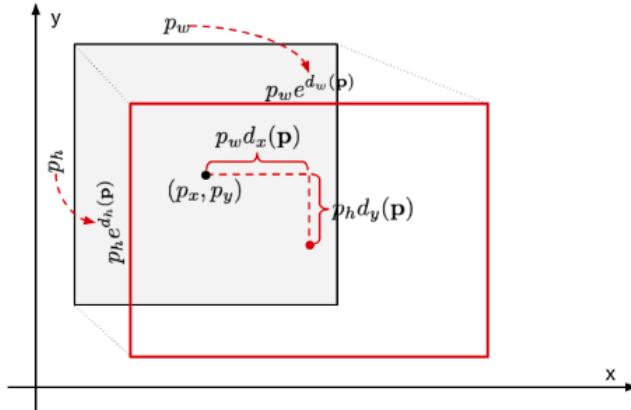
$$t_x = (G_x - P_x)/P_w$$

$$t_y = (G_y - P_y)/P_h$$

$$t_w = \log(G_w/P_h)$$

$$t_h = \log(G_h/P_h).$$

$$(t_i - d_i(P))^2 + \lambda * \|w\|$$



: from <https://lilianweng.github.io/posts/2017-12-31-object-recognition-part-3/>







\*\* Getting rid of one by RoI input to the feature extractor.

\*\*

### The Algorithm :

1.[A] Slective Search with 2k proposals per image. [Objectness[1] , Selective Search[2] ]

2.[B] Using a CNN based classifier network for feature extraction

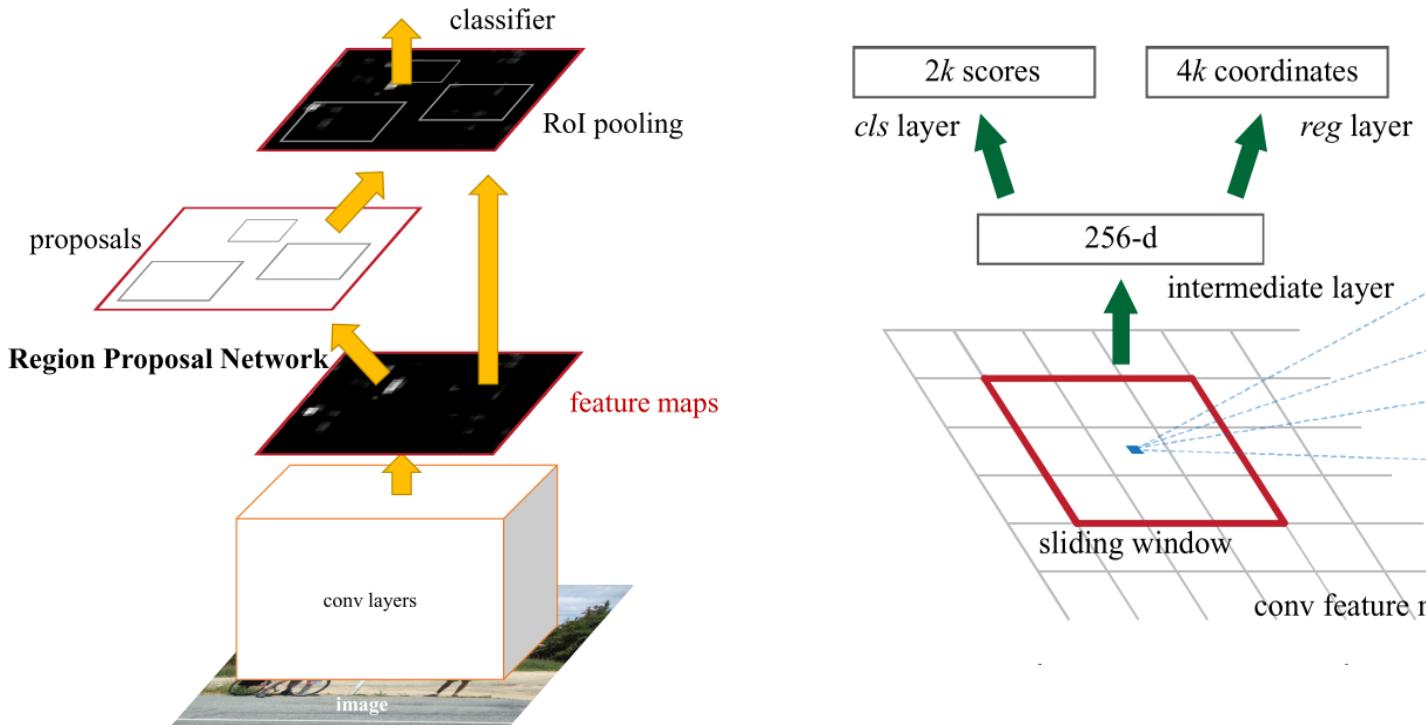
-> Pre training with ImageNet data on VGG16 (and two other smaller nets) images with 200 classes lr 0.01

-> Fine tuning with lr 0.001 with wrapped region proposal with  $N + 1$  class. N is no of object class 1 is background

-> Oversampling positive classes 32 positive sample + 96 neg sample per mini-batch

3.[B] The last max pool layer is replaced with a RoI Pooling layer for extracting fixed sized feature maps for every Fmap RoI so that it can layer (H=7, W=7 for VGG 16)

4. [C] The last fully connected layer is replaced with 2 sibling dense layer (i) one for classificatiob (ii) anther for bbox regression.



\*\* The selective search is replaced with Region Proposal Network that mostly shares weights with Feature extractor.

\*\* The region proposals are further broken down to 3 different scales with 3 different aspect ratio called anchor boxes.

### The Alorithm :

1.[A] Region Proposal Network proposes possible box coordinates and objectness confidence at each pixel location at a sliding window

2.[B] At each location box coordinates are predicted at 9 possible aspect ratio, so at each pixel location we can get 9 boxes with different

confidences instead of just one. It will make possible to detect two overlapping objects of different aspect ratio or scale. Please look at the

predictions on the image on top right. 3 aspect ration for big boxes , 3 for medium, 3 for small.

2.[B] Using a CNN based classifier network for feature extraction

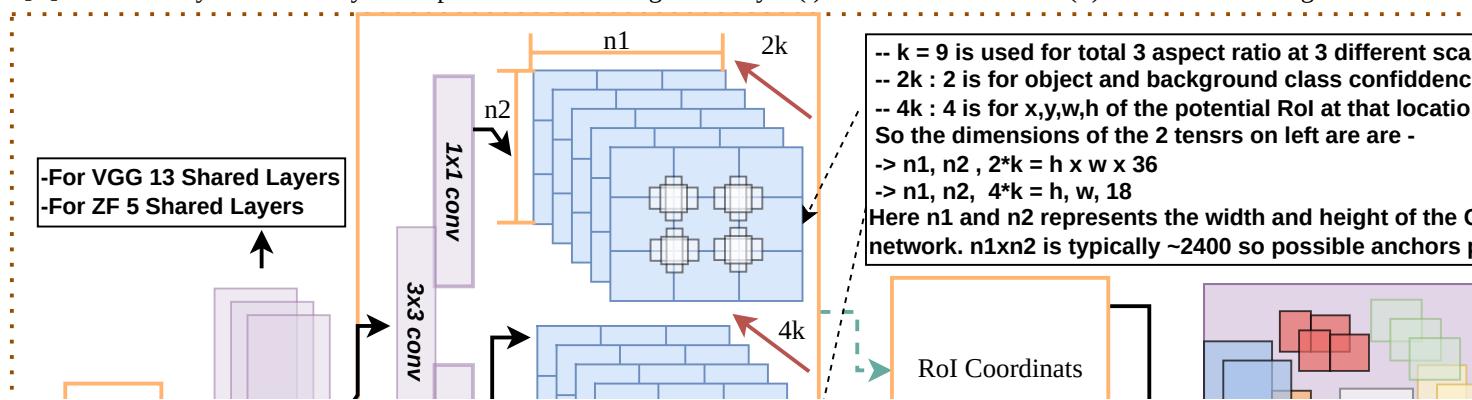
-> Pre training with ImageNet data on VGG16 (and two other smaller nets) images with 200 classes lr 0.01

-> Fine tuning with lr 0.001 with wrapped region proposal with  $N + 1$  class. N is no of object class 1 is background

-> Oversampling positive classes 32 positive sample + 96 neg sample per mini-batch

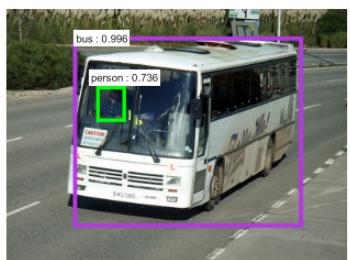
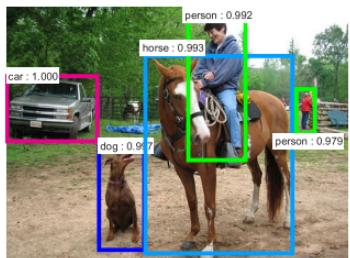
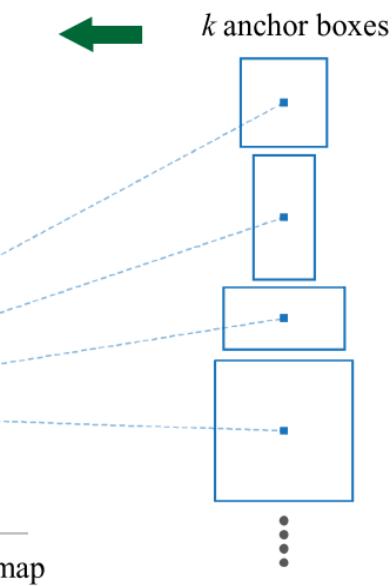
3.[B] The last max pool layer is replaced with a RoI Pooling layer for extracting fixed sized feature maps for every Fmap RoI so that it can into the first dense layer (H=7, W=7 for VGG 16)

4. [C] The last fully connected layer is replaced with 2 sibling dense layer (i) one for classificatiob (ii) anther for bbox regression.





can fit into the first dense



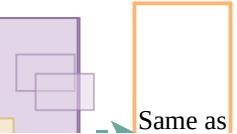
manner.

e

an fit

les.  
e score  
n.

CNN output of Region Proposal  
per image 9\*2400

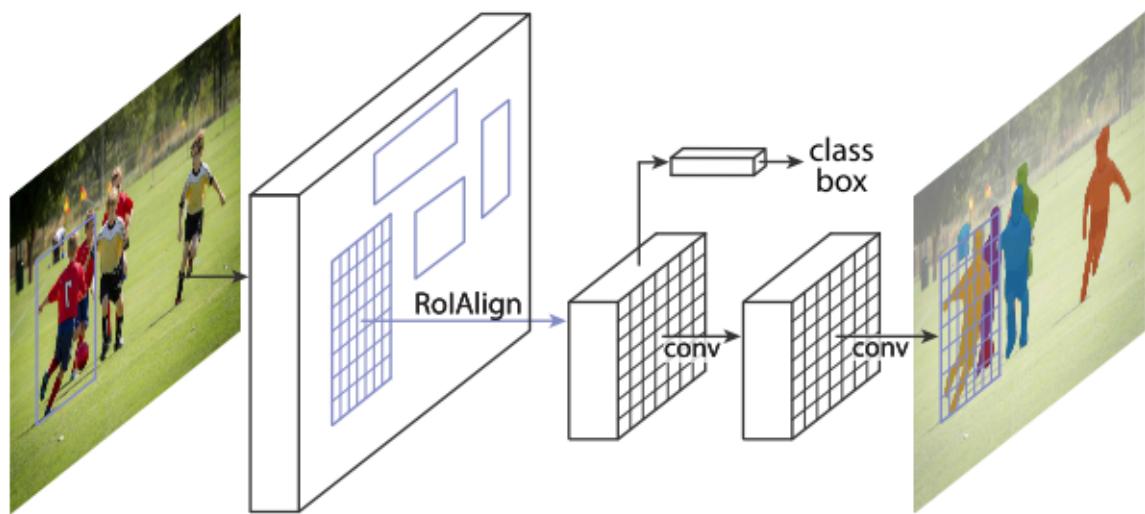
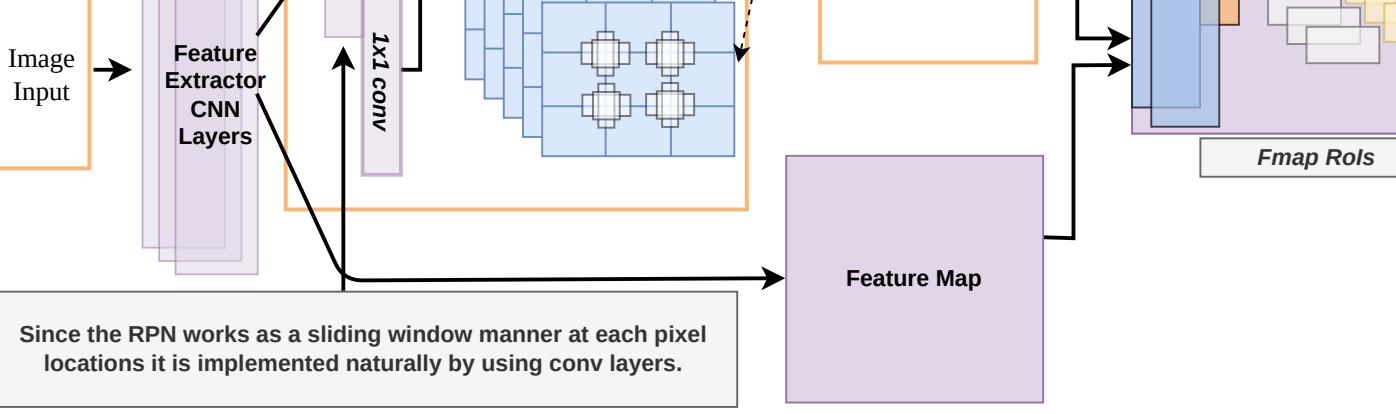


Same as

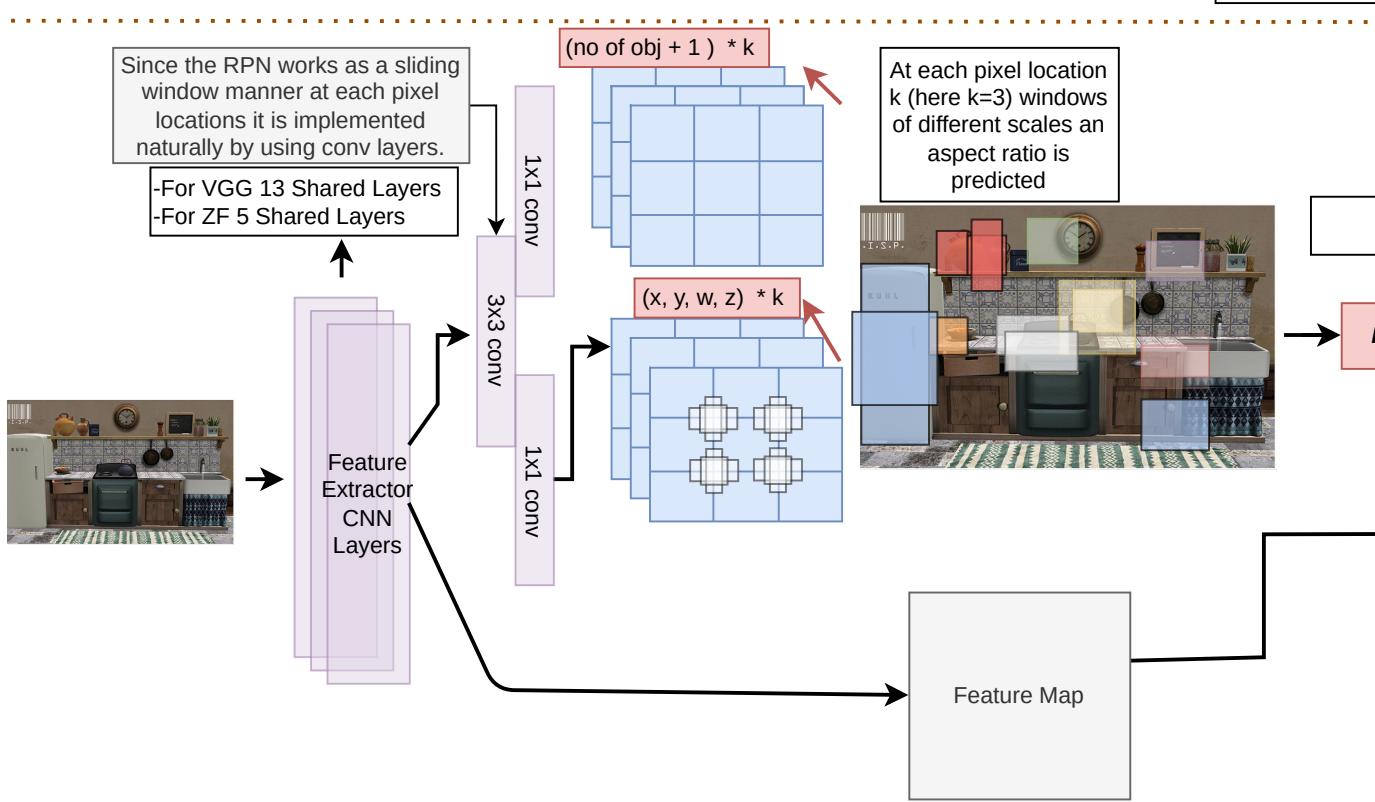




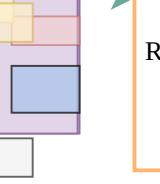




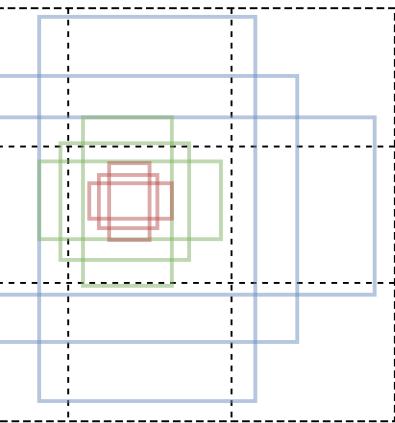
Anchor was used  
The used 800 anchors  
same 9 boxes are  
capable of predicting  
-> For Multiclass  
fully connected  
-> For Faster R-CNN





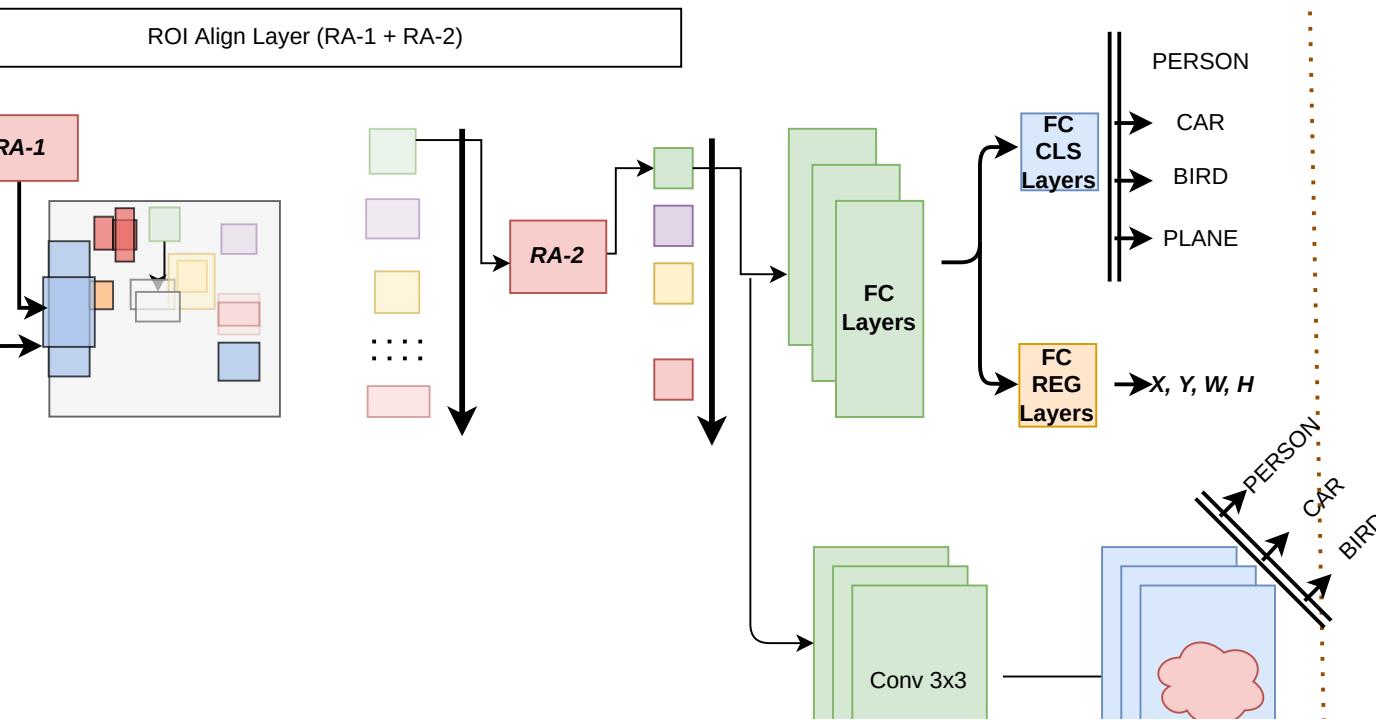


Fast  
RCNN



possible anchor box at one pixel locations

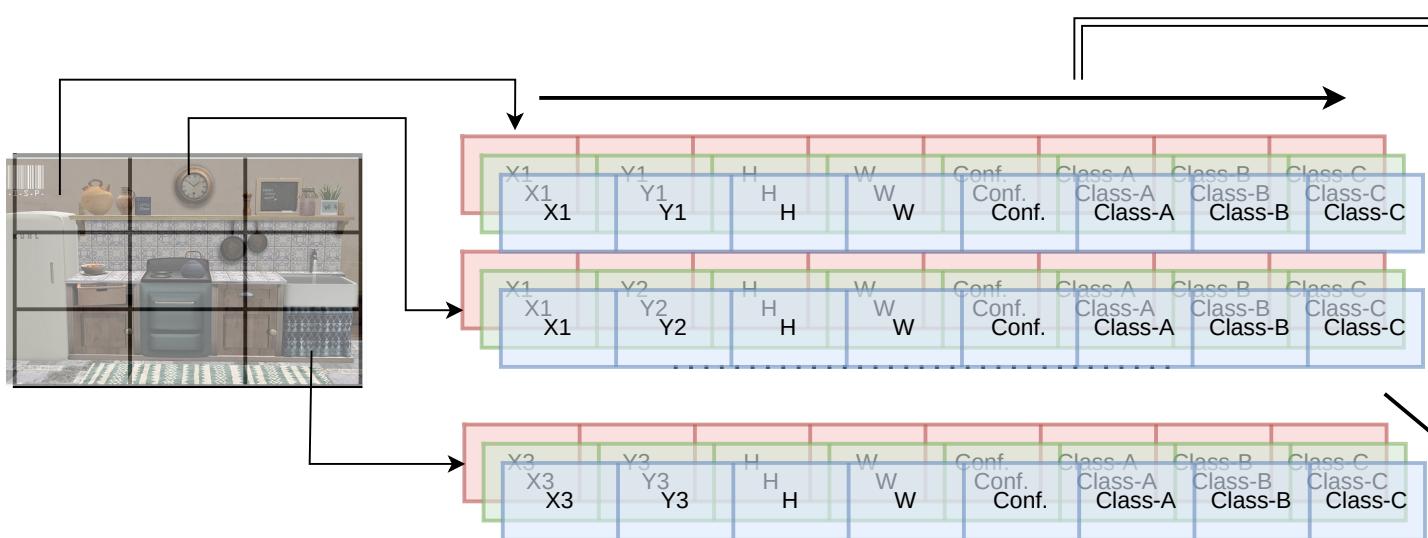
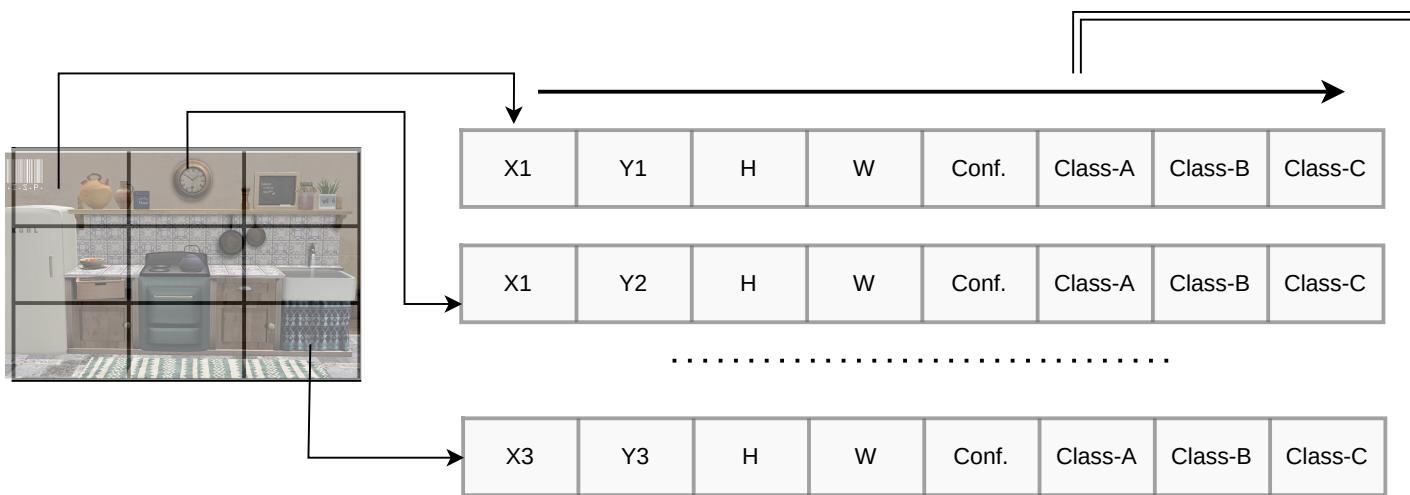
ed by Multibox[] but they were not translation invariant.  
anchor boxes different locations of the image, faster rcnn used  
at every locations  
reducing ( $w \cdot h = \sim 2400$ )  $\sim 2400 \times 9$  RoI per image at maximum.  
set the shape of tensor for ROI prediction was  $4+1 \times 800$  dim.  
selected layer  
RCNN it reduced to  $4+2 \times 9$  convolutional output layer











**V 1.0 : GRID\_W , GRID\_H , (B\*5 + Classes)**

Predicted Tensor Grid\_Wx Grid\_HxN

Step :1

Step :2

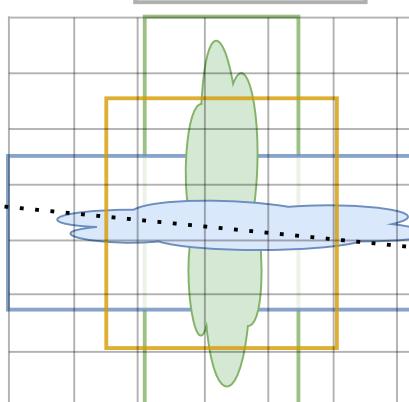
Ground Truth Tensor

CLASS C  
CLASS B  
CLASS A  
CONFIDENCE

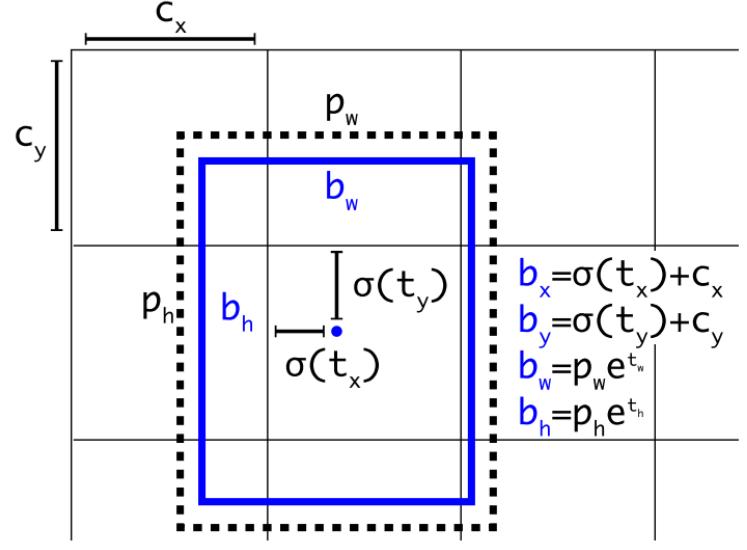
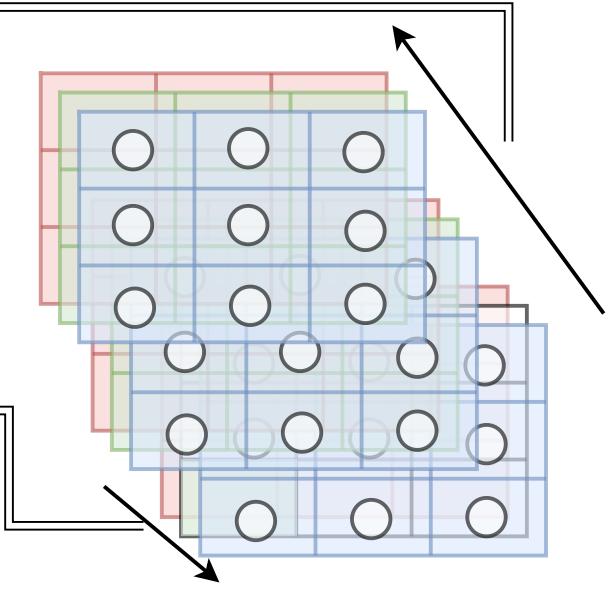
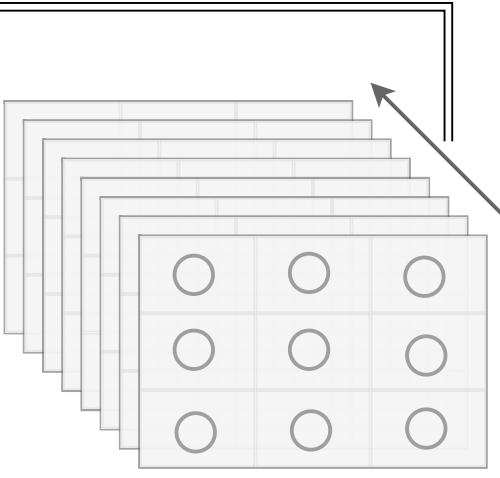
GRID W  
Y  
H  
X

**YOLO**

Anchor boxes

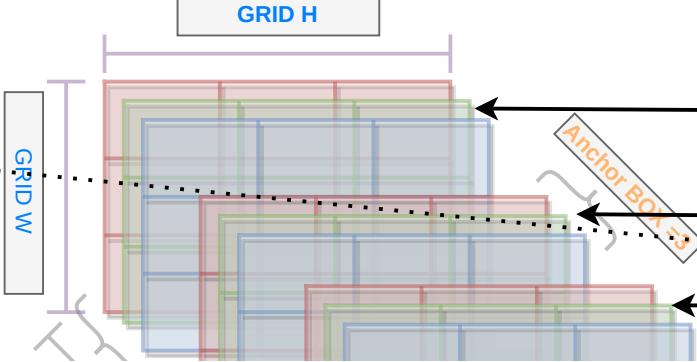






V 2.0 & 3.0 :GRID\_W , GRID\_H , BOX , ( 5 + Classes )

Predicted Tensor Grid\_W x Grid\_H x Box x N



Step :3

$P(\text{Class}) = P(\text{Class} | \text{Object}) * P(\text{Object})$

At each grid center of the anc box find max probability

For x in Grid  
For y in Grid

For i in Box

$P(\text{Class C})$

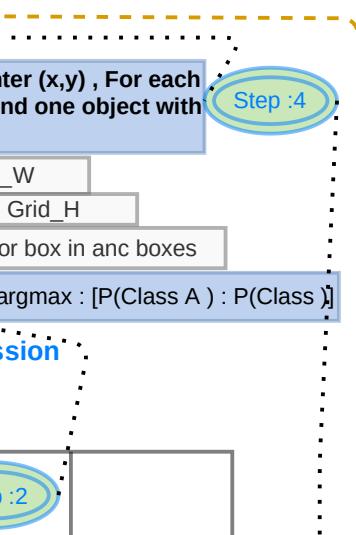
$P(\text{Class B})$

$P(\text{Class A})$

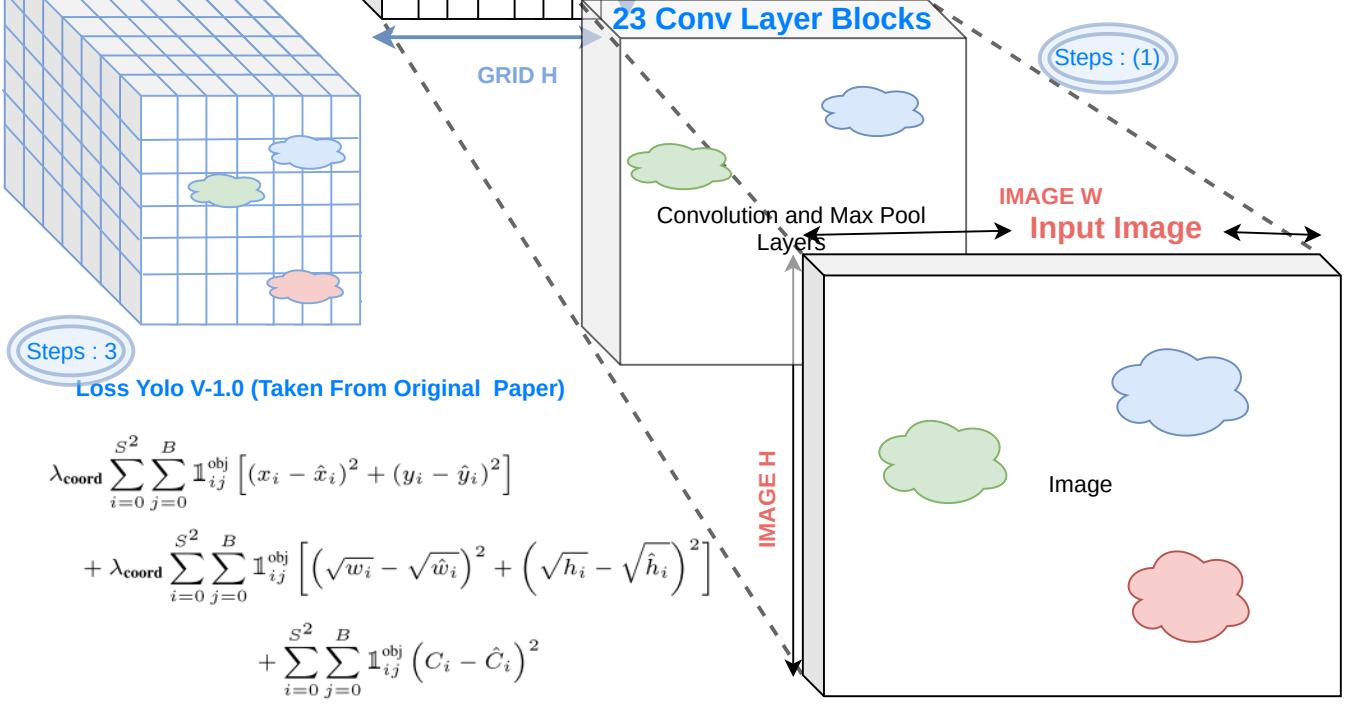
a Non Max Suppression  
b Removing low confidence box

Step









Steps : 3

**Loss Yolo V-1.0 (Taken From Original Paper)**

$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
 \end{aligned}$$

Designed By - Irfan Mohammad Al Hasib  
email - irfanhasib.me@gmail.com

Some figures are inspired from the original YOLO papers

**Train**  
1. Cr  
2. Ge  
3. Ca  
**Final Infer**  
1. Pr  
2. Re  
3. P  
4. Fo  
5. N  
6. Ca



