

April 7, 2023

DSML: Computer Vision.

## Object localization and Detection - 2

Class starts  
@ 9:05 pm.



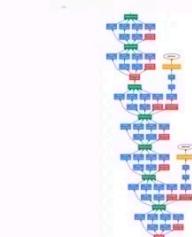
What normal people see  
when they walk on street



What Computer Vision  
folks see



**WHO WOULD WIN?**



STATE OF THE ART  
NEURAL NETWORK



ONE NOISY BOI

## Recap:

### \* Bounding boxes

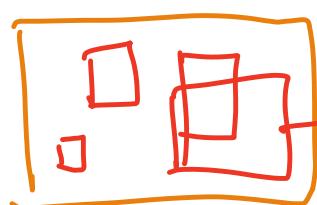
1]  $(x_c, y_c), (x_e, y_e)$       2]  $(x_c, y_c, w, h)$ .

### \* Output of an object detection network.

0)  $c$  - class of the object detected,  $\rightarrow$  Classification ✓

(2-5) B.Box co-ordinates  $\cdot \rightarrow$  Regression. ✓

### \* How to deal with multiple bounding box proposals?



$\rightarrow$  20,000 guesses.

+ 1 network to evaluate them all!!

↓  
20

### \* Non-max suppression (NMS).

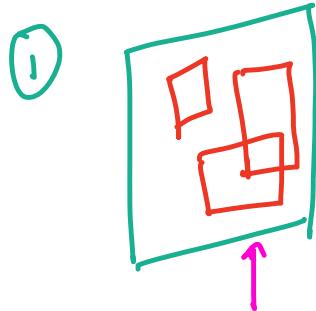
IoU - Intersection over Unions to check if boxes are representing the same object or not.

## Agenda :

- \* Two stage methods - R-CNN family. ✓  
→ Slow, but high accuracy.
- \* Single stage methods - YOLO family.  
→ Faster, slightly lower accuracy.
- \* Real time application - Self driving cars.

## R-CNN family :

### \* Naive :



①

↓  
Maybe no  
CNN  
here.

② Sliding window  
to get region  
proposals.

CNN.

③ Apply Image  
classification  
to all region  
proposals.

↓ O/P

④ Confidence for  
each proposal.

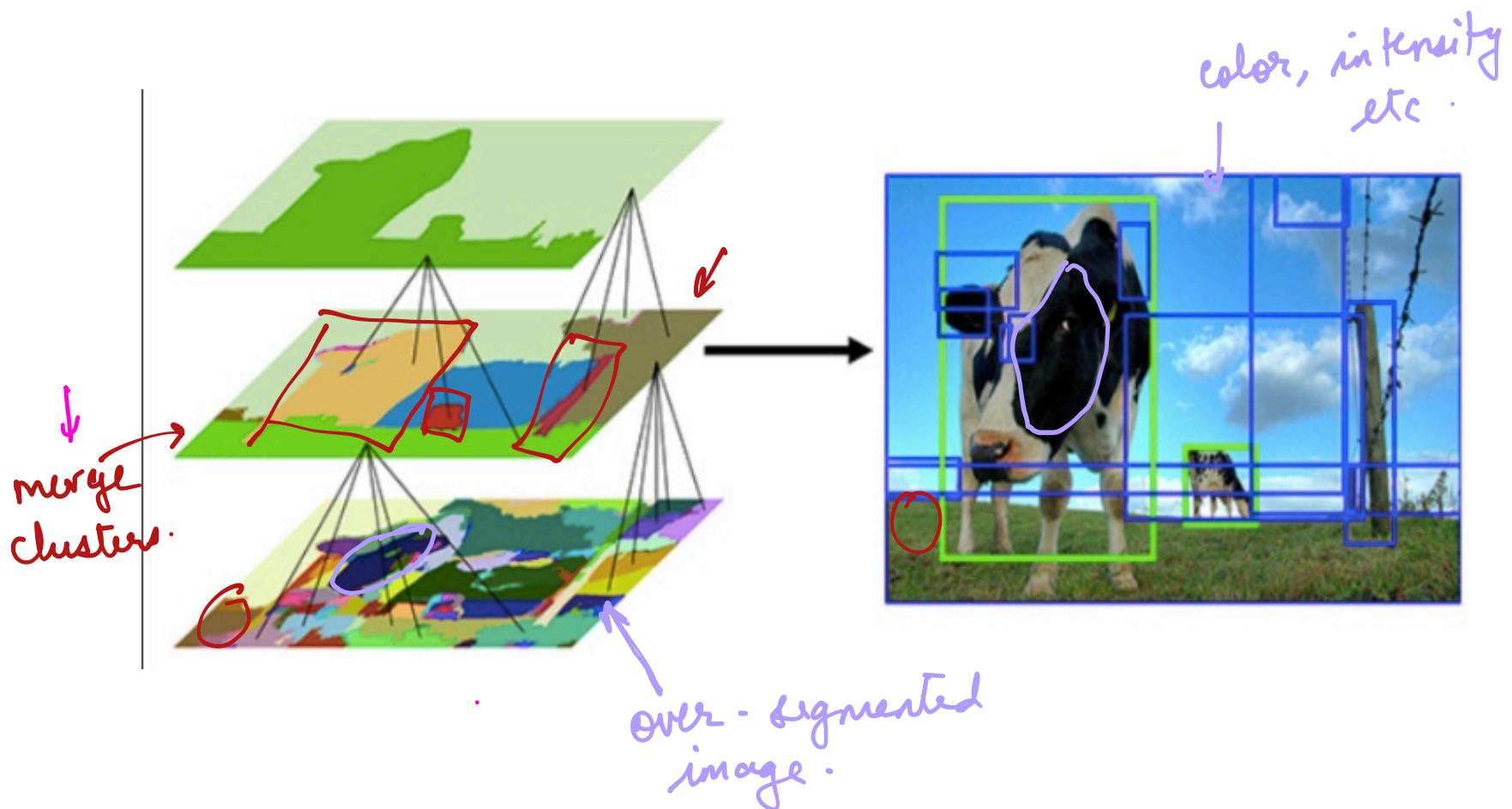
CNN.

⑤ Only keep ←  
high confidence  
bounding boxes.

R-CNN - 2014 →  
(Imagenet 2013  
detection dataset).

② and instead of using all  
possible shifts & scales of window,  
use selective search to extract proposals

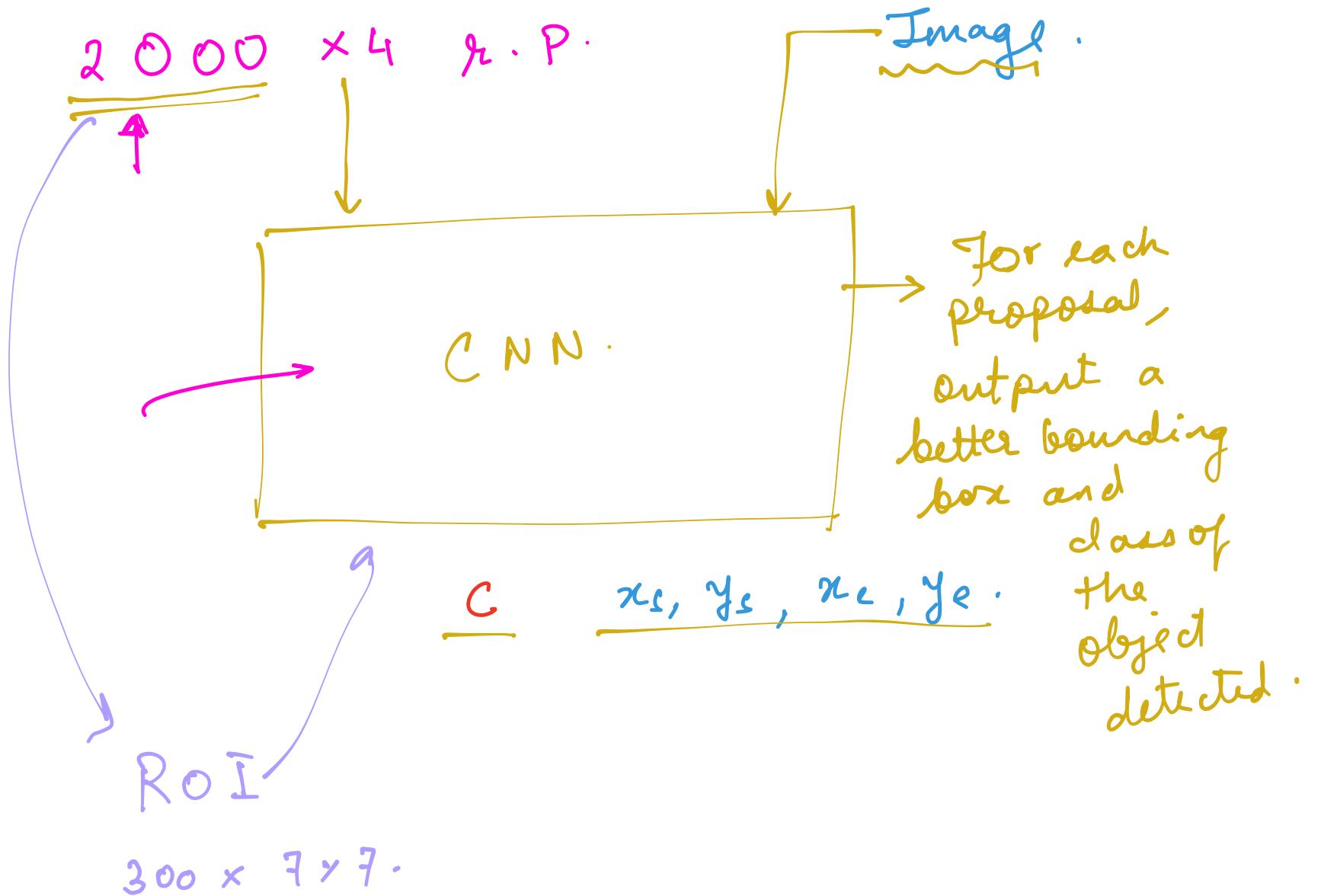
Selective search instead of all possibilities:



"Region proposal algorithm"  $\approx$  2000 B.B.

Important because it is much much faster than step 2 in the naive algorithm.

$$(2000 \times 4)$$



External region proposals method (Selective Search Algorithm)	✓
	R-CNN
Test time per image	<u>50 seconds</u>
Speed up	<u>1x</u>
mAP (VOC 2007)	66.0% ✗

Mean Average Precision

- Not practical at all !!
- (a) Clustering / S.S. takes time.
- (b) We have to run a CNN 2000 times per image !!
- major bottleneck

## Fast R-CNN

- \* The bottleneck seems to be 1000 proposals per image.
- \* Solution: Add a ROI pooling algorithm.  
CNN architecture.

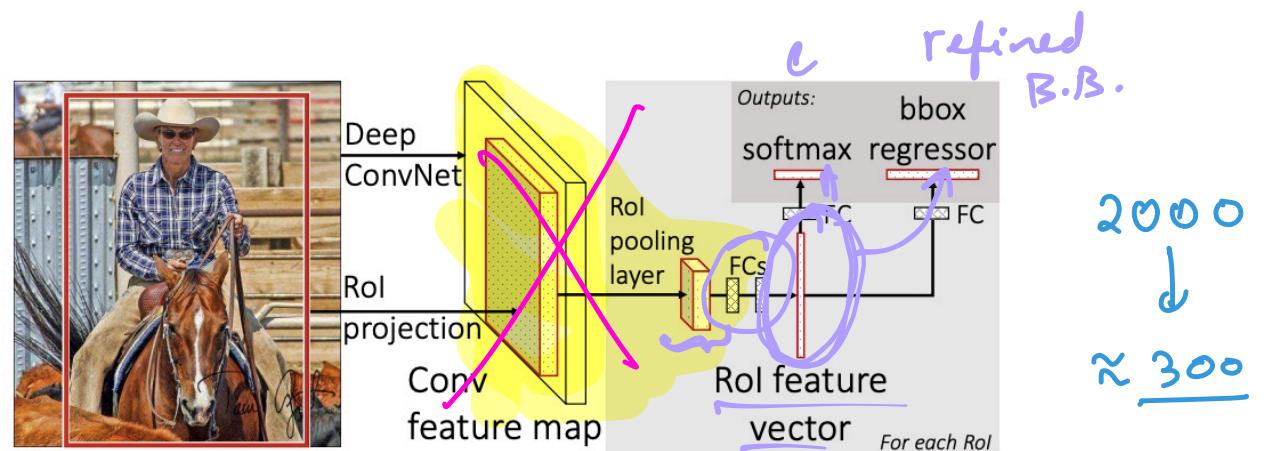
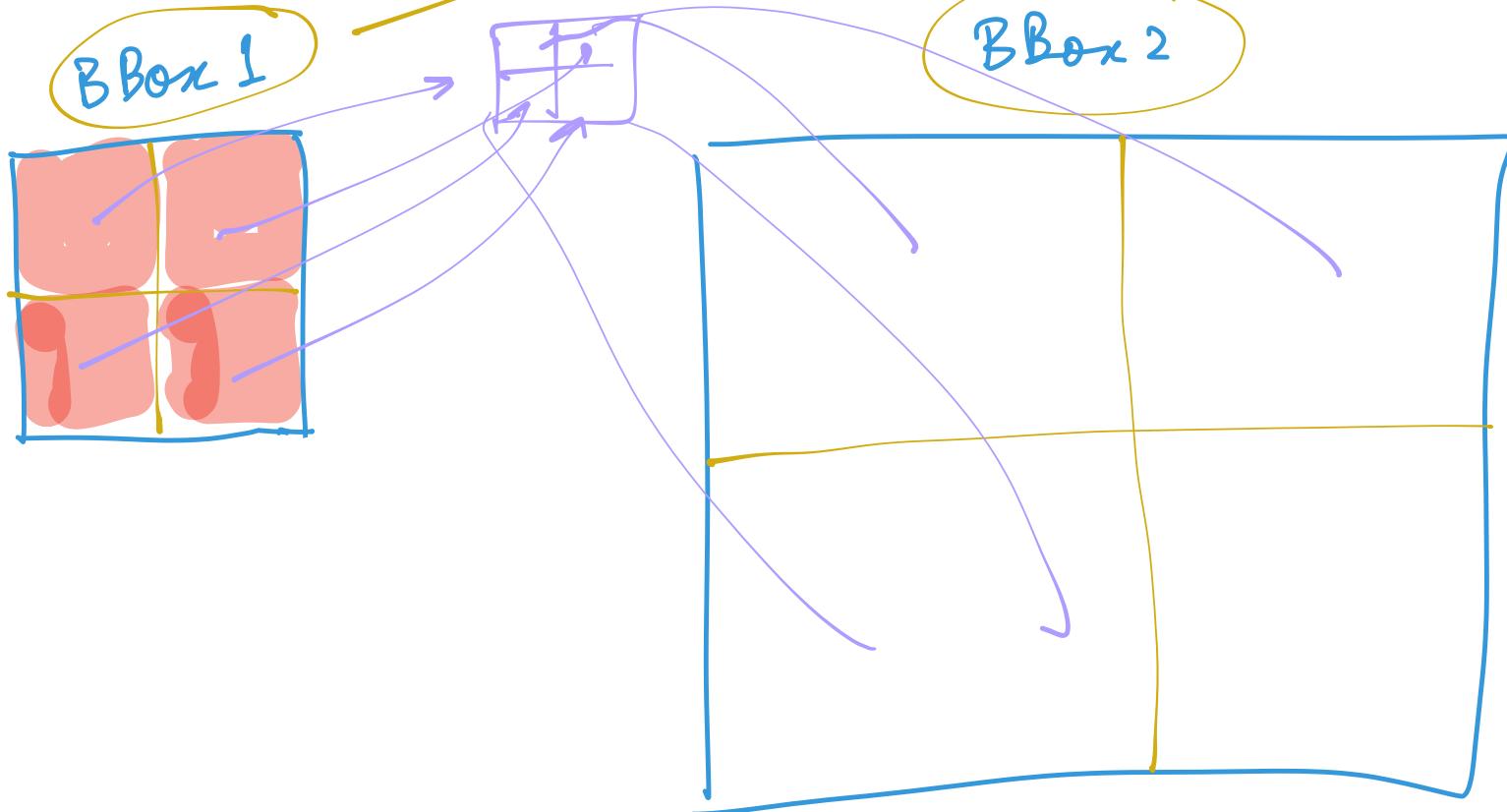


Figure 1. Fast R-CNN architecture. An input image and multiple regions of interest (RoIs) are input into a fully convolutional network. Each ROI is pooled into a fixed-size feature map and then mapped to a feature vector by fully connected layers (FCs). The network has two output vectors per ROI: softmax probabilities and per-class bounding-box regression offsets. The architecture is trained end-to-end with a multi-task loss.

$2000 \times 4$

ROI proposals.

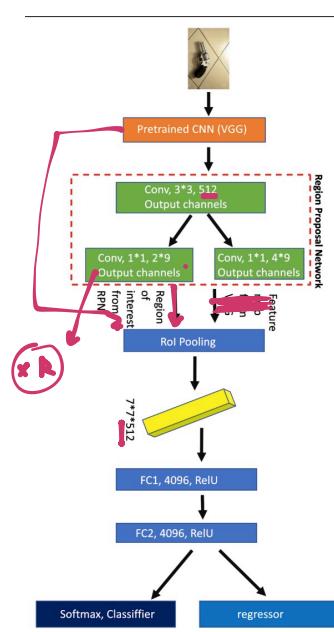
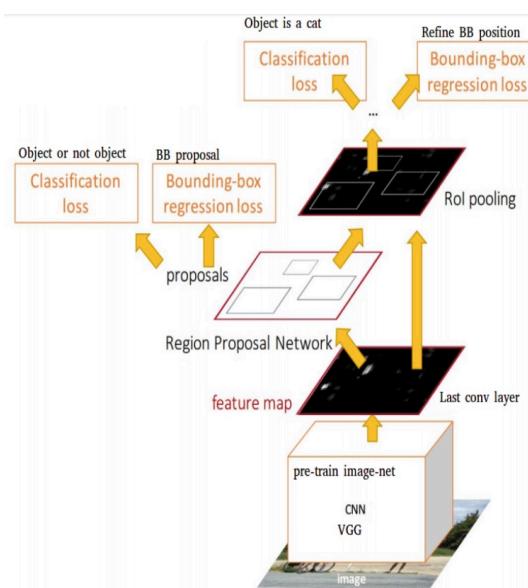
$2 \times 2$ .



External region proposals method (Selective Search Algorithm)	✓	✓
	<b>R-CNN</b>	<b>Fast R-CNN</b>
Test time per image	50 seconds	2 seconds
Speed-up	1x	25x
mAP (VOC 2007)	66.0%	66.9%

## Faster R-CNN paper:

- These guys introduced a CNN which does the job of the I.S./Clustering algorithm.
- RPN - region proposal network.



External region proposals method (Selective Search Algorithm)	✓	✓	✗
	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image	50 seconds	2 seconds	0.2 seconds
Speed-up	1x	25x	250x
mAP (VOC 2007)	66.0%	66.9%	66.9%

1 — 0.2 seconds.

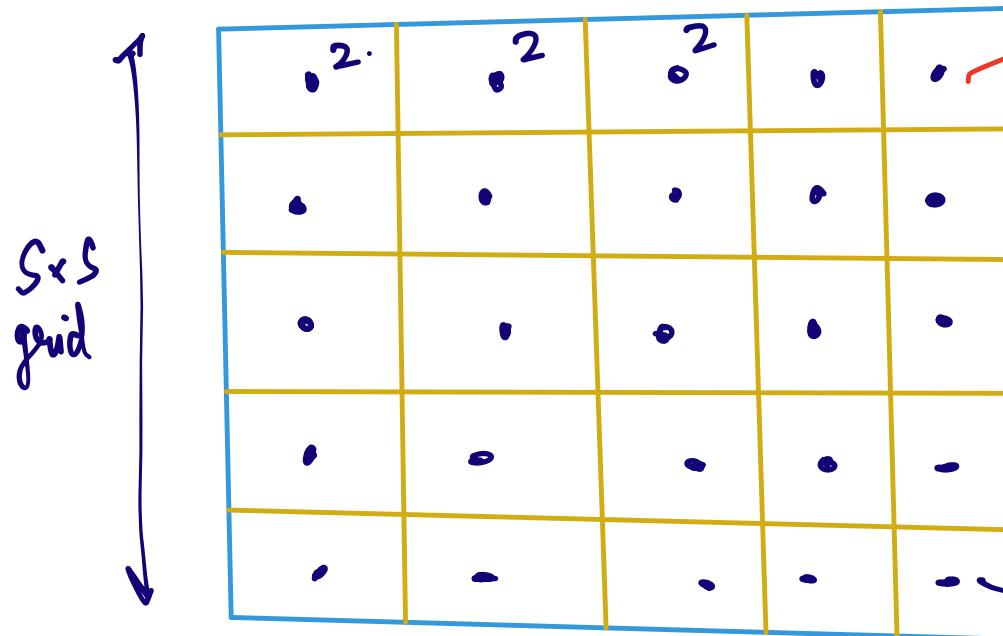
Human visual system

? — 1 second?

$$\frac{1}{0.2} = \text{S}$$

30 - 60 fps

YOLO family: → You only look once!



R-CNN - look for proposals which can be anywhere in this image.

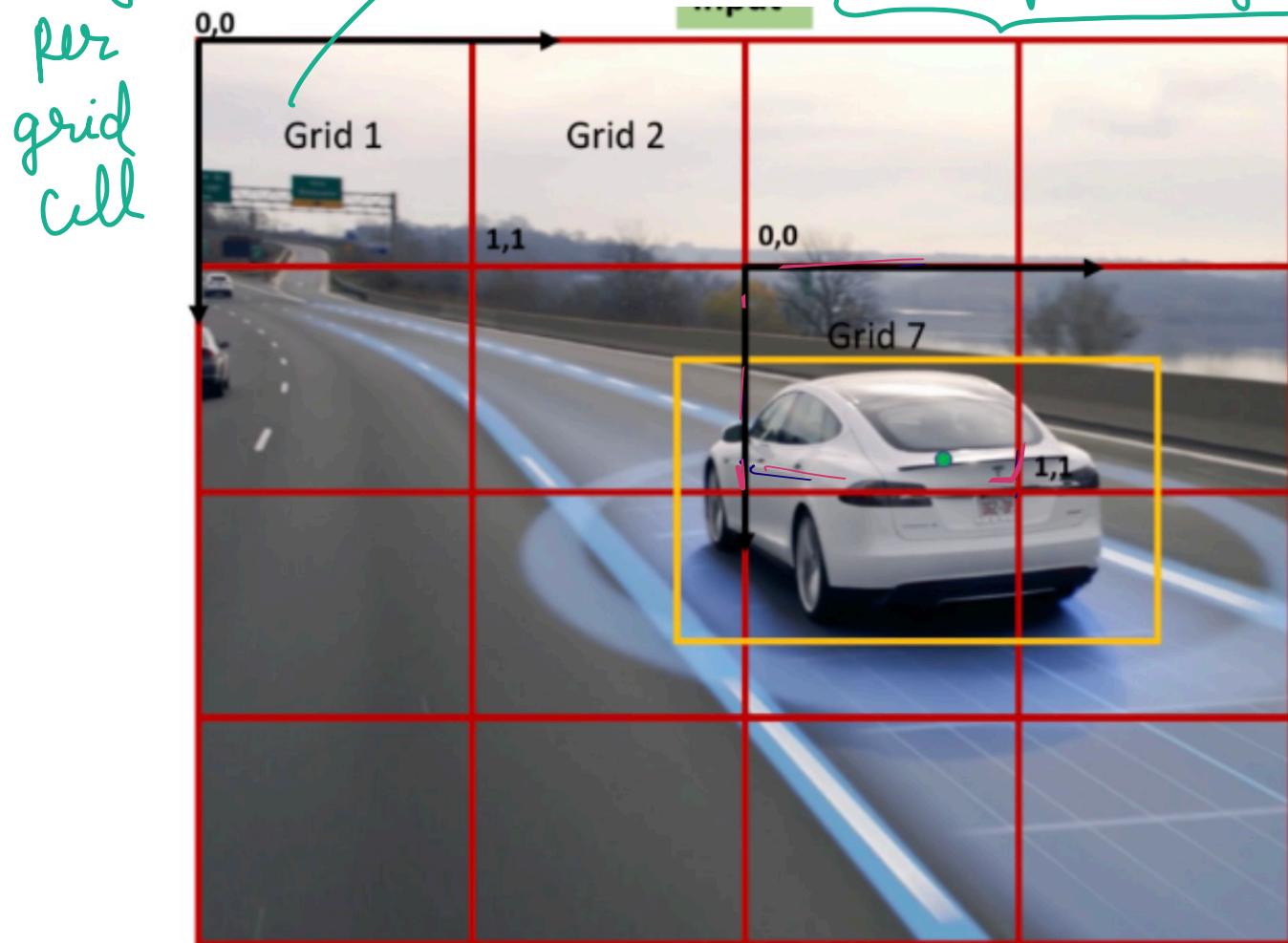
YOLO family: 2 proposals per cell.

# YOLO v3 in easy steps :

## YOLO in easy steps:

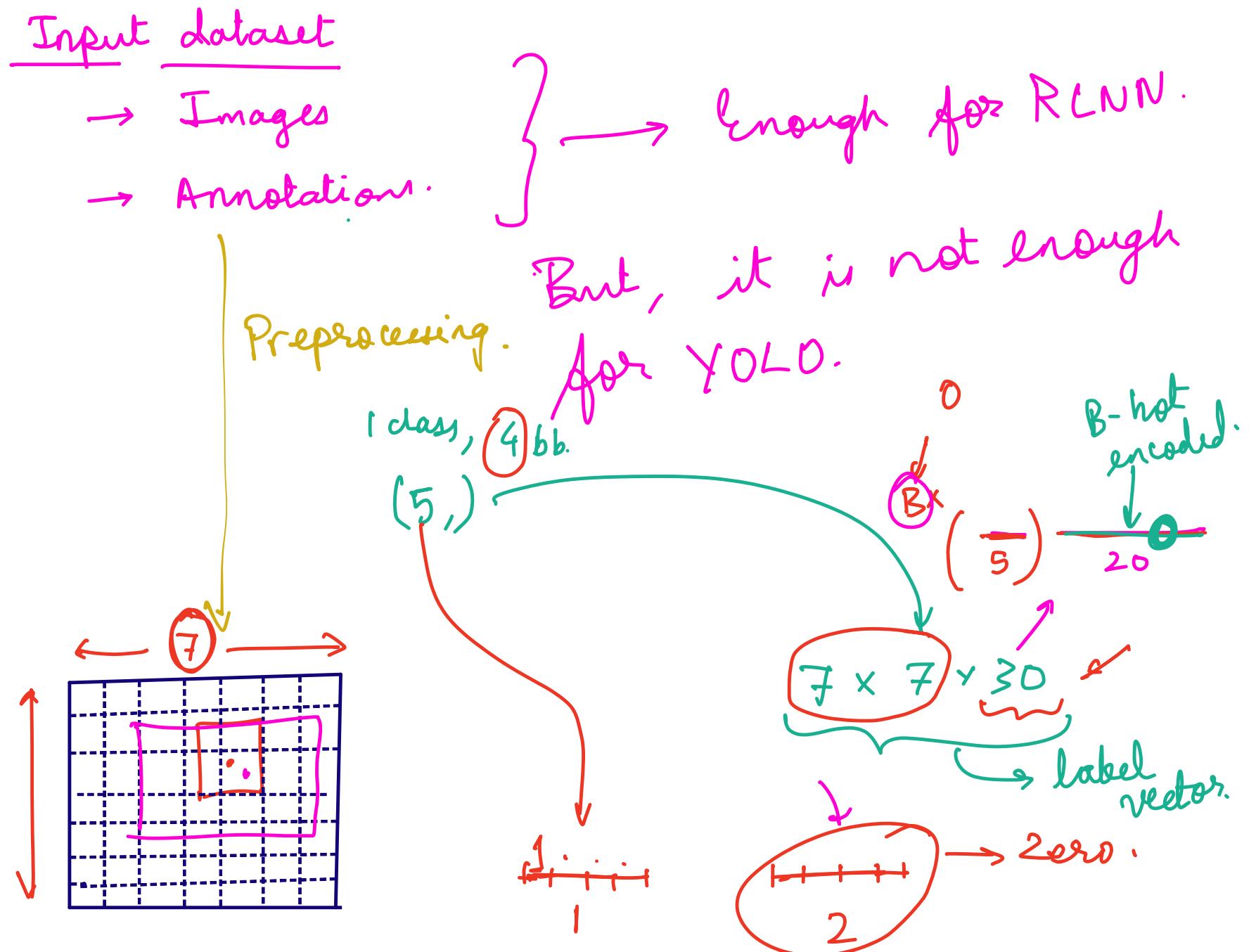
1. Divide the image into multiple grids. For illustration, I have drawn 4x4 grids in above figure.
2. Label the training data as shown in the above figure.
  - If C is number of unique objects in our data,  $s \times s$  is number of grids into which we split our image, then our output vector will be of length  $s \times s \times (C+5)$ .
  - For e.g. in above case, our target vector is  $4 \times 4 \times (3+5)$  as we divided our images into  $4 \times 4$  grids and are training for 3 unique objects: Car, Light and Pedestrian.
3. Make one deep convolutional neural net with loss function as error between output activations and label vector.
  - Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.
4. Keep in mind that the label for object being present in a grid cell (P.Object) is determined by the presence of object's centroid in that grid.
  - This is important to not allow one object to be counted multiple times in different grids.

$B \rightarrow$  No. of region proposals per grid cell → CNN → ① If object is there →  
 ② 4 co-ordinates →  
 ③ Class of the object. } Each proposal for that grid .

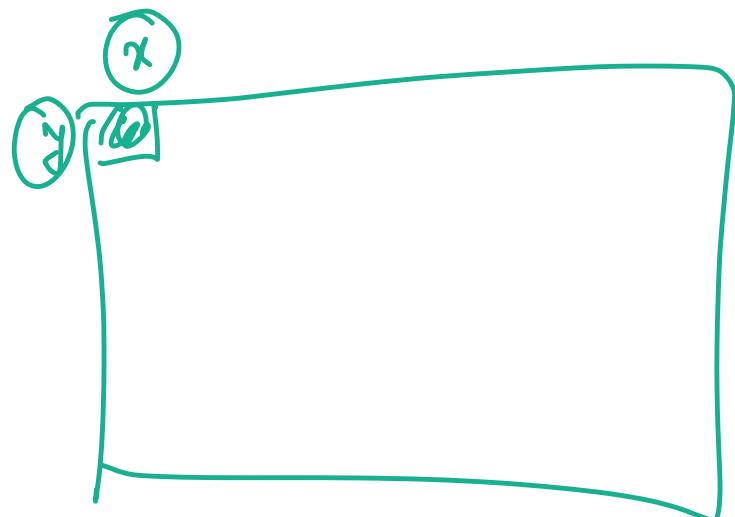


$$1 \rightarrow \frac{B \times (5 + C)}{\text{No. of classes present.}} \rightarrow \text{no. of classes present.}$$

B = No. of region proposals per grid cell  
 C = No. of classes present.  
 5 = No. of coordinates per object

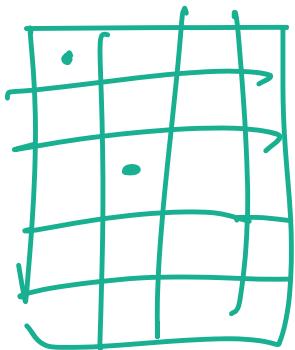


$$R(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix}$$



$$R(\theta) \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_{\text{new}} \\ y_{\text{new}} \end{bmatrix}$$



$$q = k$$

$$q \times w \times h = \frac{300}{100}$$
$$w \times h = \frac{300}{q \times 3}$$