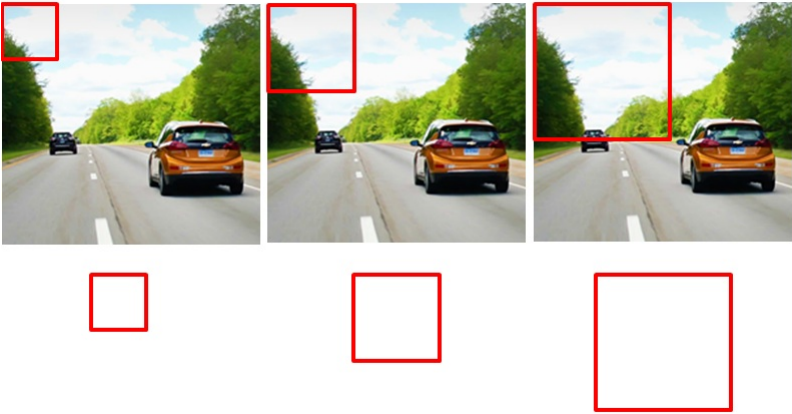
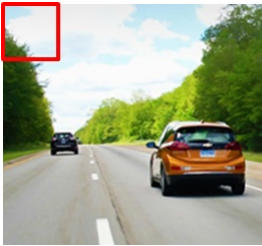


Object Detection

Sliding window approach



RCNN

Image classification such as cats vs dogs has been a problem on which a considerable accuracy has been achieved but object detection and localization is one area which didn't have satisfactory results by the time this paper is written. This paper proposes a CNN architecture to deal with object detection and localization problem and achieves a relative improvement of more than 30% on the measure mAP(Mean average precision) compared with previous works on the task.

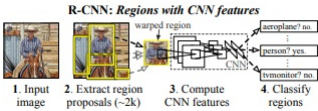
Object detection with R-CNN :

3 steps are followed in the process of object detection

- 1. Generation of a substantial number of region proposals where an object could be present
- 2. CNN to generate a 4096 dimensional vector from every object proposal
- 3. Class specific SVM trained to classify the vector

Region Proposals :-

Selective search is used to generate the region proposals. Selective search works by grouping adjacent pixels either by color, texture etc. and keeping them as one region. The process is repeated iteratively to generate a set of proposals each of which are warped into size 227x227 to use as input for CNN. Warping has negative effects since the object can look completely different after warping, this issue is later resolved in Fast R-CNN.



- It first takes an image as input:



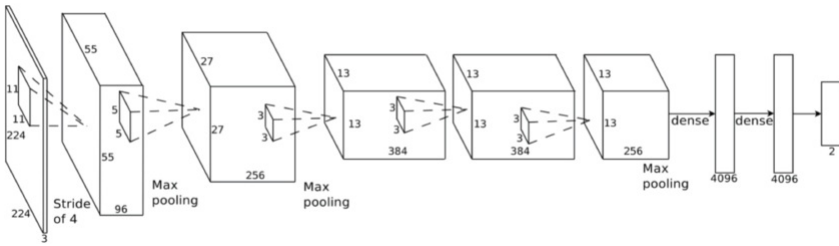
- Then, it generates initial sub-segmentations so that we have multiple regions from this image:



- The technique then combines the similar regions to form a larger region (based on color similarity, texture similarity, size similarity, and shape compatibility);
- Finally, these regions then produce the final object locations (Region of Interest).

Feature Extraction :-

Alexnet, a popular architecture which performed great on Imagenet is used for CNN.



Alexnet

Features are computed by passing the image through 5 convolutional layers and 2 fully connected layers. The output of last fully connected layer is taken as feature vector.

In order to converge fast and to work with less amount of data the CNN is first trained on a classification task on large auxiliary dataset ILSVRC2012 classification. The final layer of 1000 outputs in Alexnet is replaced by N+1 classification outputs where N is the number of classes in dataset and +1 for background.

All region proposals with IOU>0.5 are treated as +ve samples and the rest as -ve. IOU is a simple metric used to compare similarity between 2 spaces. IOU stands for Intersection over Union and does the same as the name i.e intersection of region proposal with ground truth divided by union of both. 32 +ve windows and 96 -ve windows(background) are combined in a minibatch for training. +ve windows are oversampled since they are rare compared to background.

Object Classification :-

Once the feature vector is obtained by passing every proposal through CNN, a SVM specific for each class is trained for classification. For classification training with SVM images with IOU<0.3 are considered as background. The *positives* for that class are simply the features from the ground truth bounding boxes itself. All other proposals (IoU overlap greater than 0.3, but not a ground truth bounding box) are *ignored for the purpose of training the SVM*. Hard negative mining is used to increase the speed of convergence. Hard negative mining is the process of finding false positives predicted by the network and sampling them as negatives in the dataset to improve the convergence speed.

Results :-

VOC 2010 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
DPM v5 [20] ¹	49.2	53.8	13.1	15.3	35.5	53.4	49.7	27.0	17.2	28.8	14.7	17.8	46.4	51.2	47.7	10.8	34.2	20.7	43.8	38.3	33.4
UVA [39]	56.2	42.4	15.3	12.6	21.8	49.3	36.8	46.1	12.9	32.1	30.0	36.5	43.5	52.9	32.9	15.3	41.1	31.8	47.0	44.8	35.1
Regionlets [41]	65.0	48.9	25.9	24.6	24.5	56.1	54.5	51.2	17.0	28.9	30.2	35.8	40.2	55.7	43.5	14.3	43.9	32.6	54.0	45.9	39.7
SegDPM [18] ¹	61.4	53.4	25.6	25.2	35.5	51.7	50.6	50.8	19.3	33.8	26.8	40.4	48.3	54.4	47.1	14.8	38.7	35.0	52.8	43.1	40.4
R-CNN	67.1	64.1	46.7	32.0	30.5	56.4	57.2	65.9	27.0	47.3	40.9	66.6	57.8	65.9	53.6	26.7	56.5	38.1	52.8	50.2	50.2
R-CNN BB	71.8	65.8	53.0	36.8	35.9	59.7	60.0	69.9	27.9	50.6	41.4	70.0	62.0	69.0	58.1	29.5	59.4	39.3	61.2	52.4	53.7

Table 1: Detection average precision (%) on VOC 2010 test. R-CNN is most directly comparable to UVA and Regionlets since all methods use selective search region proposals. Bounding-box regression (BB) is described in Section C. At publication time, SegDPM was the top-performer on the PASCAL VOC leaderboard. ¹DPM and SegDPM use context rescoring not used by the other methods.

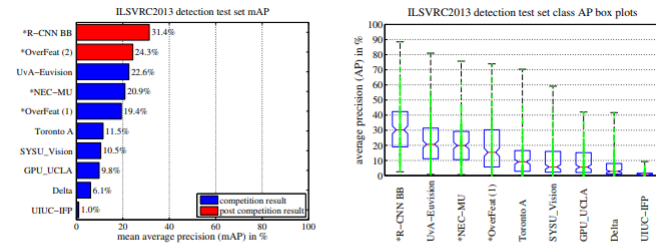


Figure 3: (Left) Mean average precision on the ILSVRC2013 detection test set. Methods preceded by * use outside training data (images and labels from the ILSVRC classification dataset in all cases). **(Right)** Box plots for the 200 average precision values per method. A box plot for the post-competition OverFeat result is not shown because per-class APs are not yet available (per-class APs for R-CNN are in Table 8 and also included in the tech report source uploaded to arXiv.org; see R-CNN-ILSVRC2013-APs.txt). The red line marks the median AP, the box bottom and top are the 25th and 75th percentiles. The whiskers extend to the min and max AP of each method. Each AP is plotted as a green dot over the whiskers (best viewed digitally with zoom).

Figure 3: (Left) Mean average precision on the ILSVRC2013 detection test set. Methods preceded by * use outside training data (images and labels from the ILSVRC classification dataset in all cases). (Right) Box plots for the 200 average precision values per method. A box plot for the post-competition OverFeat result is not shown because per-class APs are not yet available (per-class APs for R-CNN are in Table 8 and also included in the tech report source uploaded to arXiv.org; see R-CNN-ILSVRC2013-APs.txt). The red line marks the median AP, the box bottom and top are the 25th and 75th percentiles. The whiskers extend to the min and max AP of each method. Each AP is plotted as a green dot over the whiskers (best viewed digitally with zoom).

VOC 2007 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
R-CNN pool _s	51.8	60.2	36.4	27.8	25.2	52.8	60.6	49.2	18.3	47.8	44.3	40.8	56.6	58.7	42.4	25.4	46.1	36.7	51.3	55.7	44.2
R-CNN f ₆₄	59.3	61.8	43.1	34.0	25.1	53.1	60.6	52.8	21.7	47.8	42.7	47.8	52.5	58.5	44.6	25.6	48.3	34.0	53.1	58.0	46.2
R-CNN f ₆₄	57.6	57.9	38.5	31.8	23.7	51.2	58.9	51.4	20.0	50.5	40.9	46.0	51.6	55.9	43.3	23.3	48.1	35.3	51.0	57.4	44.7
R-CNN FT pool _s	58.2	63.3	37.9	27.6	26.1	54.1	66.9	51.4	26.7	55.5	43.4	43.1	57.7	59.0	45.8	28.1	50.8	40.6	53.1	56.4	47.3
R-CNN FT f ₆₄	63.5	66.0	47.9	37.7	29.9	62.5	70.2	60.2	32.0	57.9	47.0	53.5	60.1	64.2	52.2	31.3	55.0	50.0	57.7	63.0	53.1
R-CNN FT f ₆₄	64.2	69.7	50.0	41.9	32.0	62.6	71.0	60.7	32.7	58.5	46.5	56.1	60.6	66.8	54.2	31.5	52.8	48.9	57.9	64.7	54.2
R-CNN FT f ₆₄ BB	68.1	72.8	56.8	43.0	36.8	66.3	74.2	67.6	34.4	63.5	54.5	61.2	69.1	68.6	58.7	33.4	62.9	51.1	62.5	64.8	58.5
DPM v5 [20]	33.2	60.3	10.2	16.1	27.3	54.3	58.2	23.0	20.0	24.1	26.7	12.7	58.1	48.2	43.2	12.0	21.1	36.1	46.0	43.5	33.7
DPM ST [28]	23.8	58.2	10.5	8.5	27.1	50.4	52.0	7.3	19.2	22.8	18.1	8.0	55.9	44.8	32.4	13.3	15.9	22.8	46.2	44.9	29.1
DPM HSC [31]	32.2	58.3	11.5	16.3	30.6	49.9	54.8	23.5	21.5	27.7	34.0	13.7	58.1	51.6	39.9	12.4	23.5	34.4	47.4	45.2	34.3

Table 2: Detection average precision (%) on VOC 2007 test. Rows 1-3 show R-CNN performance without fine-tuning. Rows 4-6 show results for the CNN pre-trained on ILSVRC 2012 and then fine-tuned (FT) on VOC 2007 trainval. Row 7 includes a simple bounding-box regression (BB) stage that reduces localization errors (Section C). Rows 8-10 present DPM methods as a strong baseline. The first uses only HOG, while the next two use different feature learning approaches to augment or replace HOG.

VOC 2007 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
R-CNN-Net	64.2	69.7	50.0	41.9	32.0	62.6	71.0	60.7	32.7	58.5	46.5	56.1	60.6	66.8	54.2	31.5	52.8	48.9	57.9	64.7	54.2
R-CNN-Net BB	68.1	72.8	56.8	43.0	36.8	66.3	74.2	67.6	34.4	63.5	54.5	61.2	69.1	68.6	58.7	33.4	62.9	51.1	62.5	64.8	58.5
R-CNN O-Net	71.6	73.5	58.1	42.2	39.4	70.7	76.0	74.5	38.7	71.0	56.9	74.5	67.9	69.6	59.3	35.7	62.1	64.0	66.5	71.2	62.2
R-CNN O-Net BB	73.4	77.0	63.4	45.4	44.6	75.1	78.1	79.8	40.5	73.7	62.2	79.4	78.1	73.1	64.2	35.6	66.8	67.2	70.4	71.1	66.0

Table 3: Detection average precision (%) on VOC 2007 test for two different CNN architectures. The first two rows are results from Table 2 using Krizhevsky et al.'s architecture (T-Net). Rows three and four use the recently proposed 16-layer architecture from Simonyan and Zisserman (O-Net) [43].

Visualizing learned features :-

A simple method is proposed to visualize the features in CNN layers. To find what all images the feature map is identifying the activations for all region proposals are found and sorted, non-maximum suppression is performed which ignores duplicate proposals and finds the best proposal and the top regions for each are seen.

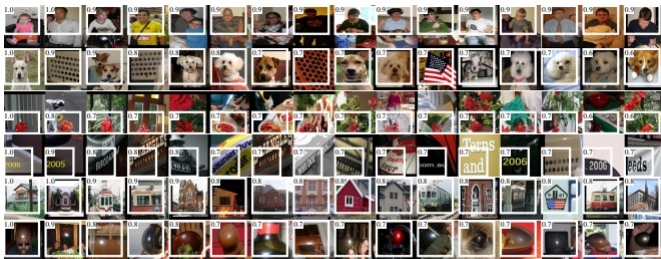


Figure 4: Top regions for six pool units. Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).

Bounding-box regression:-

After each proposal is classified with a class specific SVM, the co-ordinates of the bounding box are predicted using a class specific regressor. Inputs are center coordinates of region proposal, width and height of the region proposal and outputs are predicted center coordinates, width and height. The mapping is done by below 4 equations where P's correspond to region proposal's variables and G's are for predicted ground truth variables.

$$\hat{G}_x = P_w d_x(P) + P_x \quad (1)$$

$$\hat{G}_y = P_h d_y(P) + P_y \quad (2)$$

$$\hat{G}_w = P_w \exp(d_w(P)) \quad (3)$$

$$\hat{G}_h = P_h \exp(d_h(P)). \quad (4)$$

All values of d in above 4 equations are calculated by using pool 5 features from above CNN using the below equation

$$\mathbf{w}_* = \underset{\mathbf{w}_*}{\operatorname{argmin}} \sum_i^N (t_i^* - \hat{\mathbf{w}}_*^T \phi_i(P^i))^2 + \lambda \|\mathbf{w}_*\|^2. \quad (5)$$

Weights w are learnt to minimise the difference. The ground truth targets t are calculated using below equations

$$t_x = (G_x - P_x) / P_w \quad (6)$$

$$t_y = (G_y - P_y) / P_h \quad (7)$$

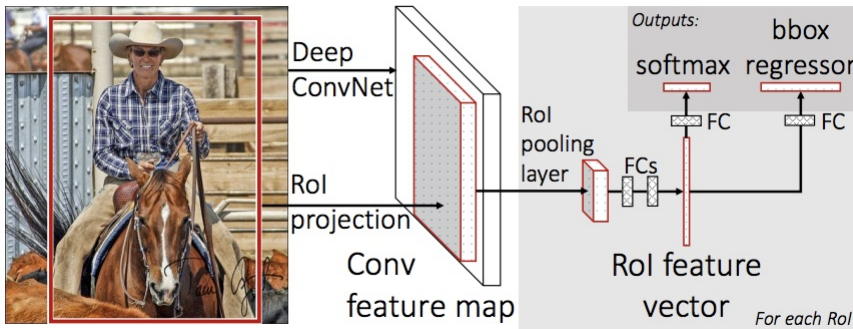
$$t_w = \log(G_w / P_w) \quad (8)$$

$$t_h = \log(G_h / P_h). \quad (9)$$

Problems with R-CNN

- It still takes a huge amount of time to train the network as you would have to classify 2000 region proposals per image.
- It cannot be implemented real time as it takes around 47 seconds for each test image.
- The selective search algorithm is a fixed algorithm. Therefore, no learning is happening at that stage. This could lead to the generation of bad candidate region proposals.

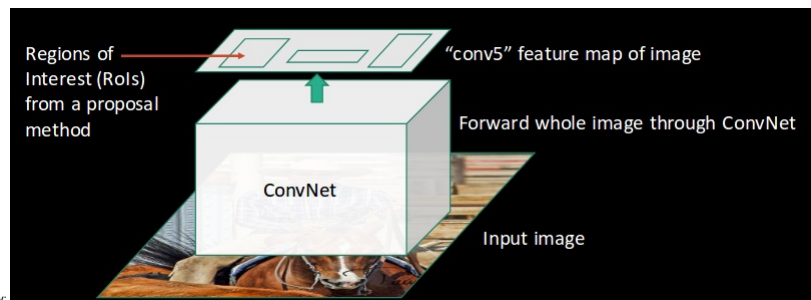
Fast R-CNN



The same author of the previous paper(R-CNN) solved some of the drawbacks of R-CNN to build a faster object detection algorithm and it was called Fast R-CNN. The approach is similar to the R-CNN algorithm. But, instead of feeding the region proposals to the CNN, we feed the input image to the CNN to generate a convolutional feature map. From the convolutional feature map, we identify the region of proposals and warp them into squares and by using a RoI pooling layer we reshape them into a fixed size so that it can be fed into a fully connected layer. From the RoI feature vector, we use a softmax layer to predict the class of the proposed region and also the offset values for the bounding box.

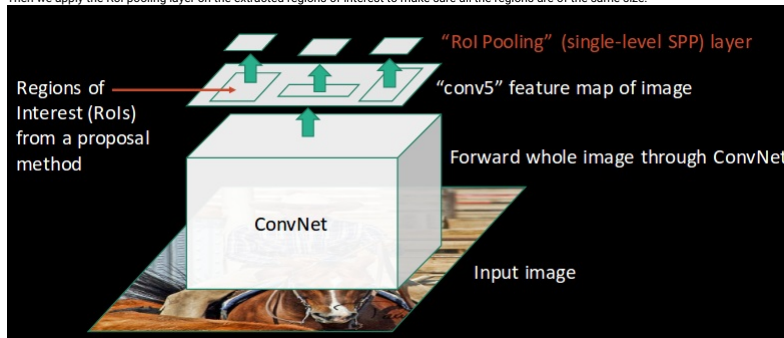


- We follow the now well-known step of taking an image as input:

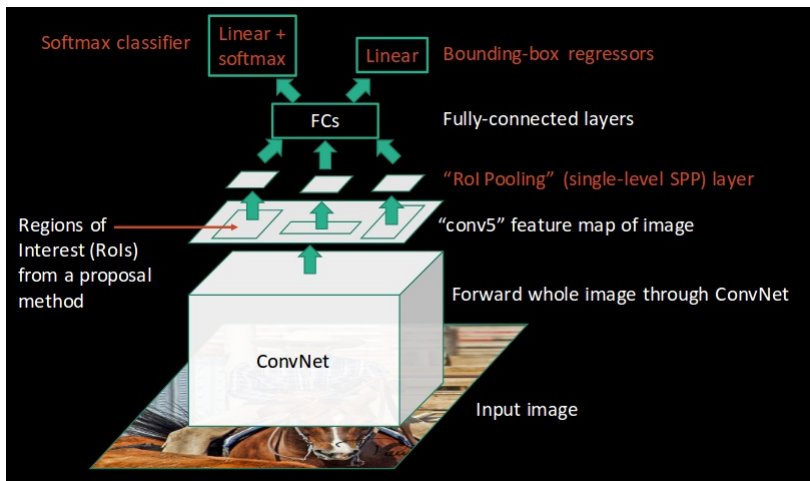


- This image is passed to a ConvNet which returns the region of interests accordingly:

- Then we apply the RoI pooling layer on the extracted regions of interest to make sure all the regions are of the same size:

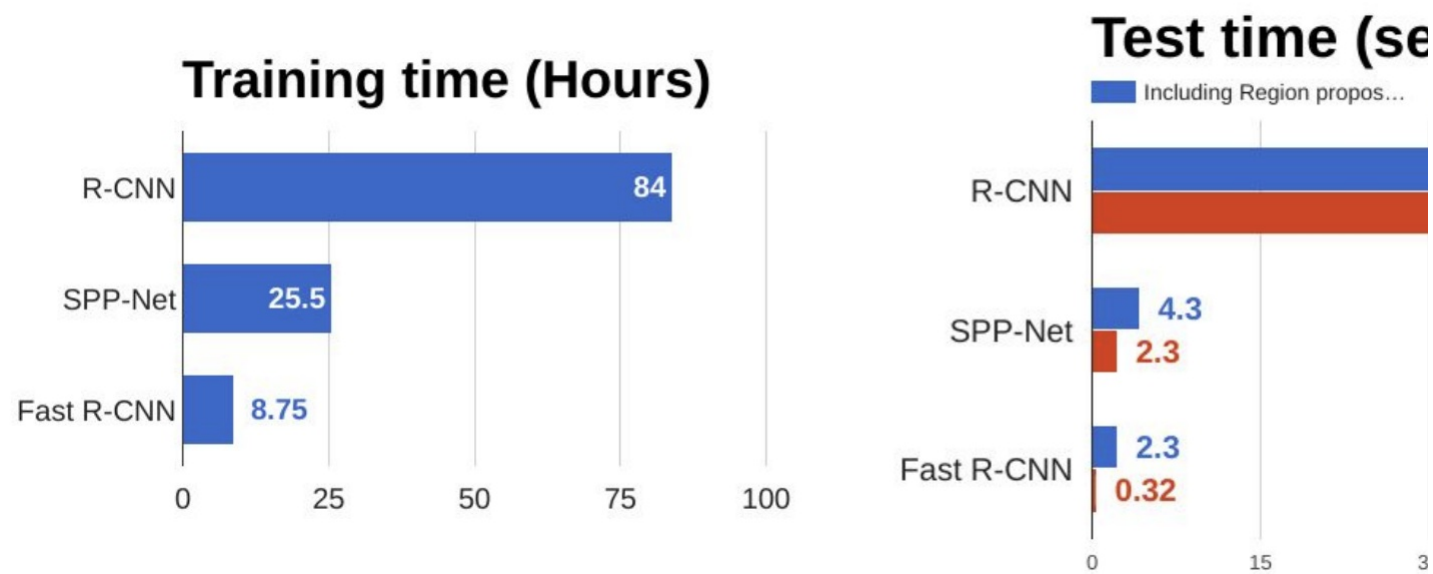


- Finally, these regions are passed on to a fully connected network which classifies them, as well as returns the bounding boxes using softmax and linear regression layers simultaneously:



This is how Fast RCNN resolves two major issues of RCNN, i.e., passing one instead of 2,000 regions per image to the ConvNet, and using one instead of three different models for extracting features, classification and generating bounding boxes.

The reason "Fast R-CNN" is faster than R-CNN is because you don't have to feed 2000 region proposals to the convolutional neural network every time. Instead, the convolution operation is done only once per image and a feature map is generated from it.



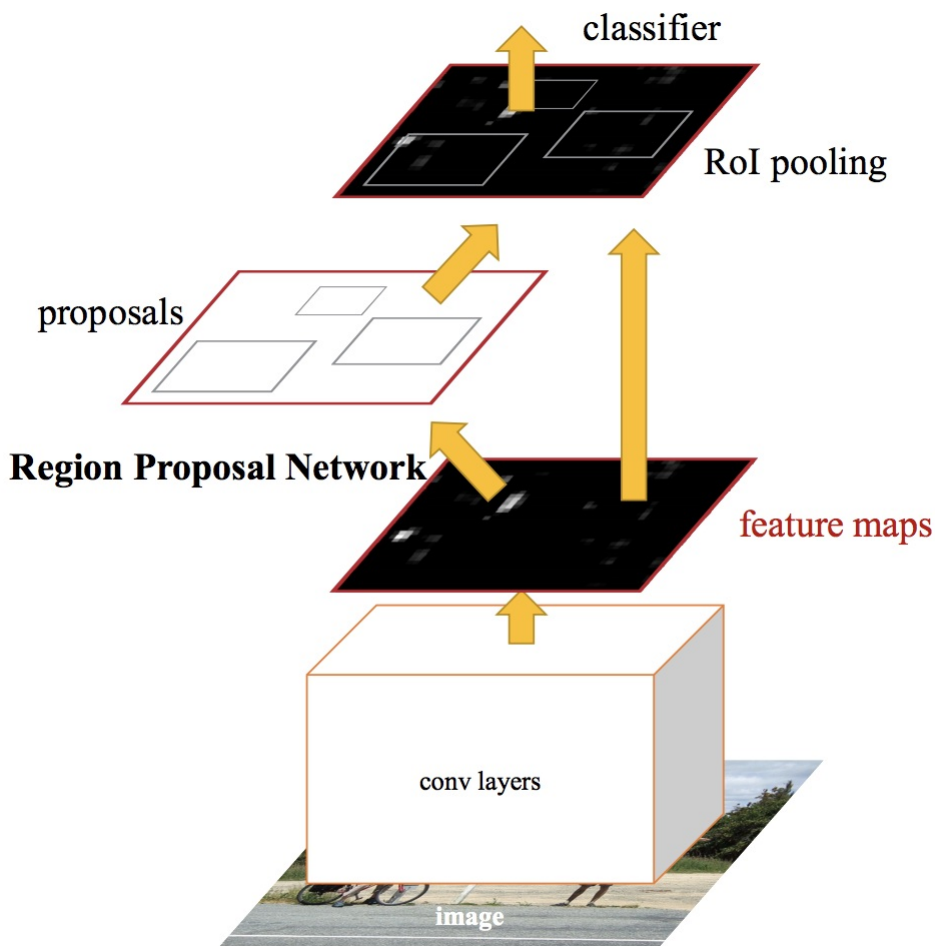
Comparison of object detection algorithms

From the above graphs, you can infer that Fast R-CNN is significantly faster in training and testing sessions over R-CNN. When you look at the performance of Fast R-CNN during testing time, including region proposals slows down the algorithm significantly when compared to not using region proposals. Therefore, region proposals become bottlenecks in Fast R-CNN algorithm affecting its performance.

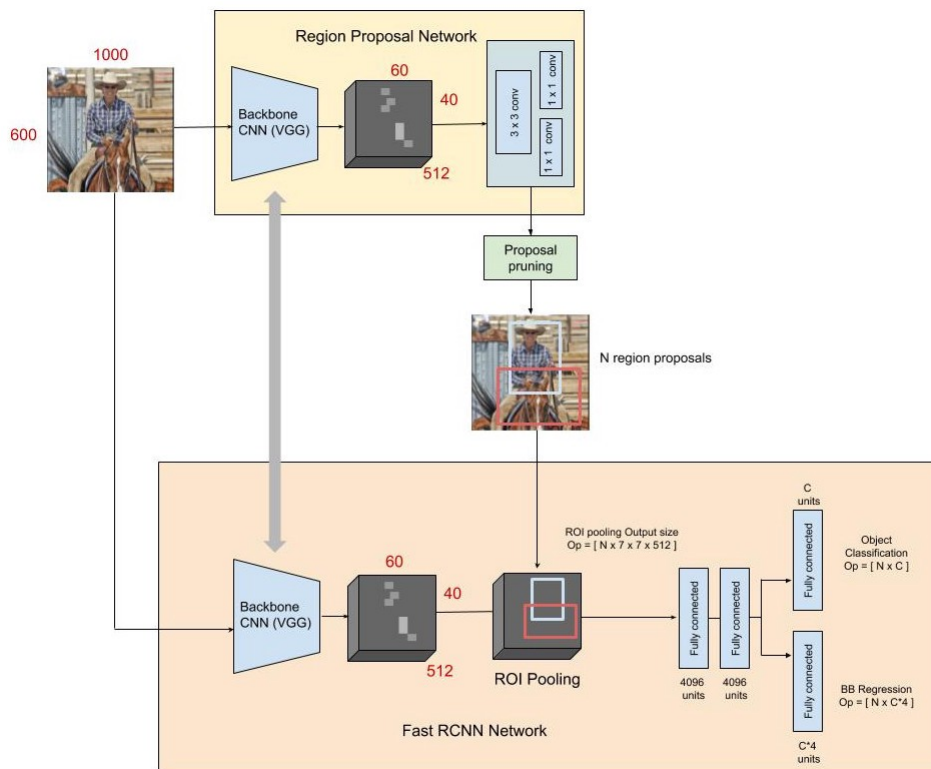
-
-
-
-

Faster R-CNN

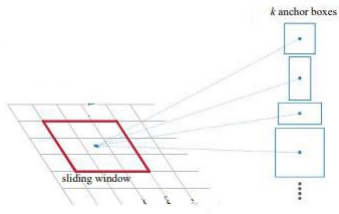
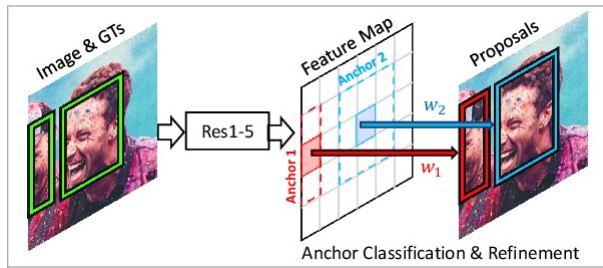
Both of the above algorithms(R-CNN & Fast R-CNN) uses selective search to find out the region proposals. Selective search is a slow and time-consuming process affecting the performance of the network. Therefore the author came up with an object detection algorithm that eliminates the selective search algorithm and lets the network learn the region proposals.



Similar to Fast R-CNN, the image is provided as an input to a convolutional network which provides a convolutional feature map. Instead of using selective search algorithm on the feature map to identify the region proposals, a separate network is used to predict the region proposals. The predicted region proposals are then reshaped using a RoI pooling layer which is then used to classify the image within the proposed region and predict the offset values for the bounding boxes.



RPN consists of a backbone network giving feature maps. Each point in the $M \times N$ Feature Map is found in the corresponding position in the original image. For each projected position, k priori bounding boxes of different sizes and ratios are set.

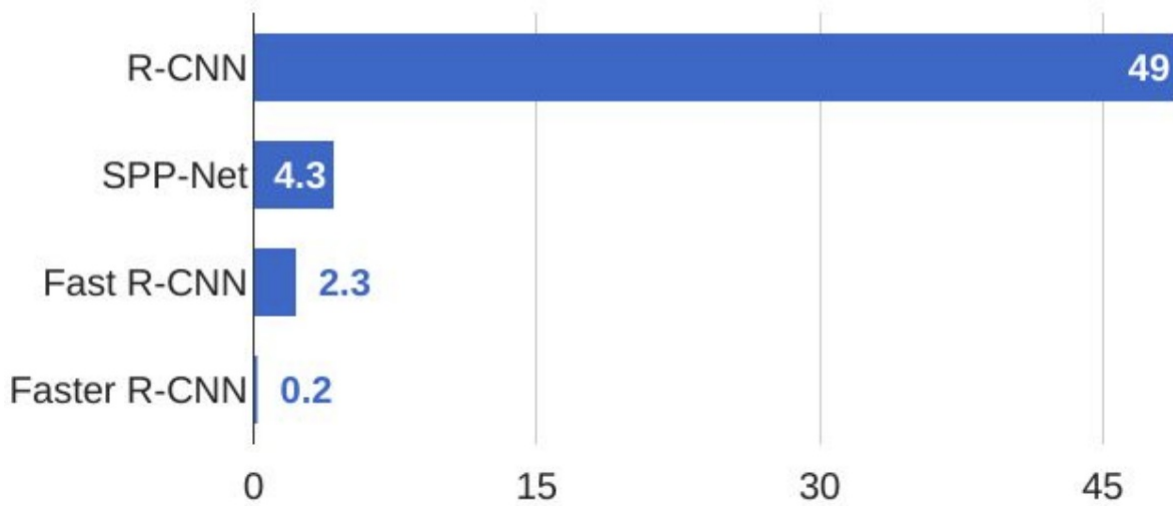


RPN uses a two-class classification, which only distinguishes the background from the object, but does not predict the class of the object.

The a priori box is matched with the ground-truth box, if a priori bounding box has an **IoU value** (intersection over union) with the ground-truth bounding box greater than **0.7** it is considered as a positive sample (belonging to an object).

For those a priori boxes whose **IoU value** is lower than 0.3 with any ground-truth box, it is considered as a negative sample.

R-CNN Test-Time Speed



Comparison of test-time speed of object detection algorithms

From the above graph, you can see that Faster R-CNN is much faster than its predecessors. Therefore, it can even be used for real-time object detection.