

Inception

(102)

The inception network was an important milestone in the development of CNN classifier. Prior to its inception, most popular CNNs just stacked convolution layers deeper and deeper, hoping to get better performance.

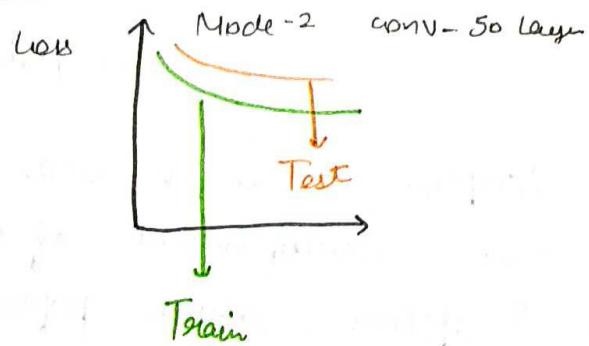
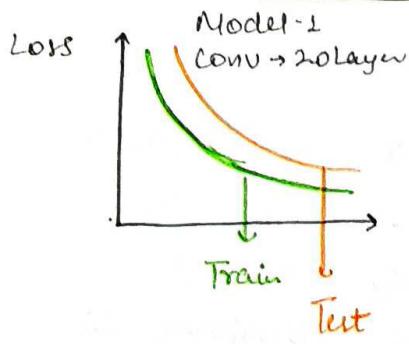
The Inception network on the other hand, was complex (carefully engineered). It used a lot of tricks to push performance; both in term of speed and accuracy. Its constant evolution lead to the creation of several versions of the network. The popular versions are as follows:

Inception V1

Inception V2 and Inception V3

Inception V4 and Inception-ResNet

Each version is an iterative improvement over the previous one. Understanding the upgrades can help us to build custom classifiers that are optimized both in speed and accuracy. Also



At a certain conv layer \rightarrow loss decrease

After that increase the layers \rightarrow loss increase

conv layers increasing \rightarrow loss decreasing

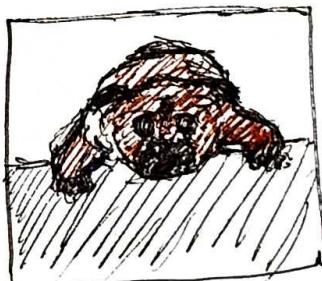
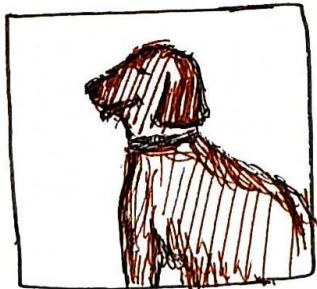
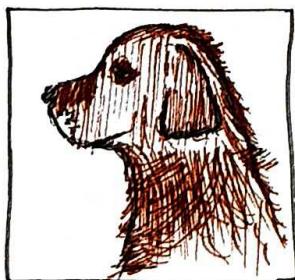
conv layers increasing \rightarrow loss steady

conv layers increasing \rightarrow loss increasing

Inception V1

The Premise:

Salient parts in the image can have extremely large variation in size. For instance, an image with a dog can be either of the following, as shown below. The area occupied by the dog is different in each image.

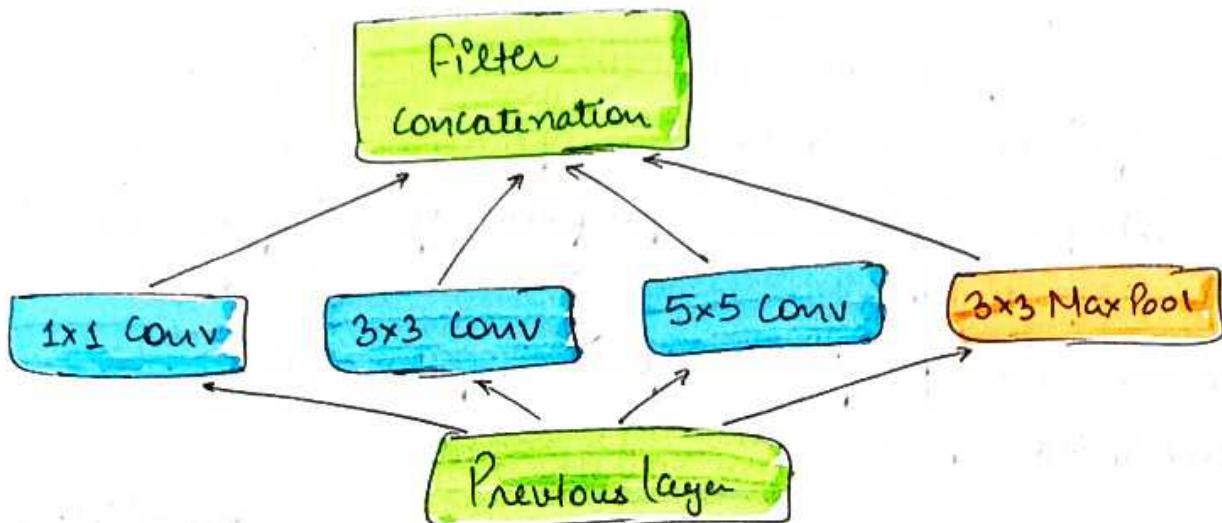


- Because of this huge variation in the location of the information choosing the right kernel size for the convolution operation becomes tough. A larger Kernel is preferred for information that is more distributed more globally and a smaller kernel is preferred for information that is distributed more locally.
- Very deep networks are prone to overfitting. It also hard to pass gradient update through the entire network.
- Naively stacking large convolutional operation is computationally expensive.

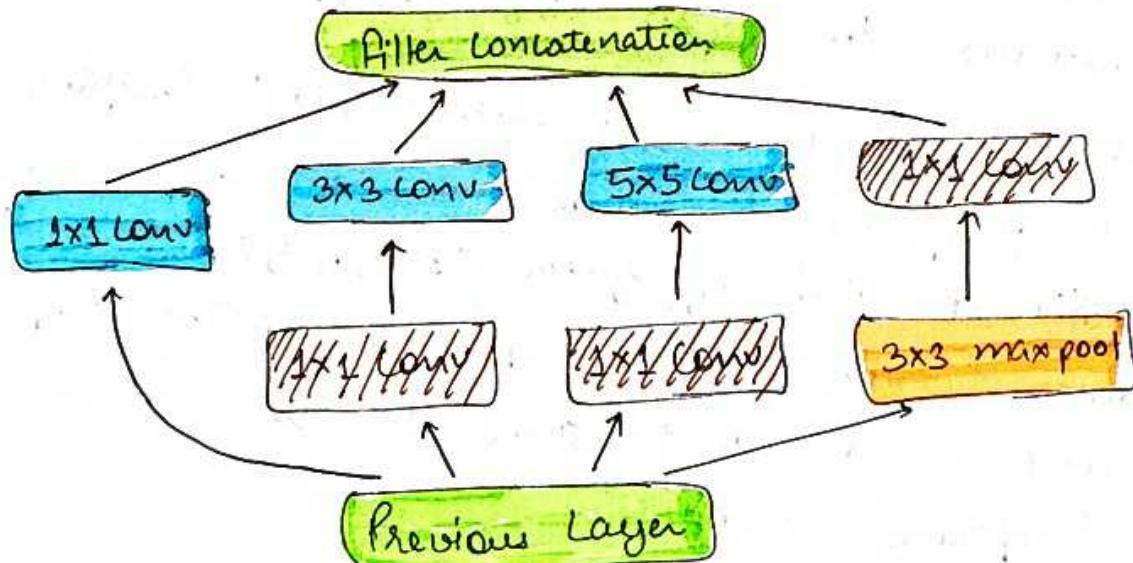
The Solution?

why not have filters with multiple sizes operate on the same level? The network essentially would get a bit "wider" rather than "deeper". The authors designed the Inception module to reflect the same.

The below image is the "naive" Inception module. It performs convolution on an input, with 3 different sizes of filters ($1 \times 1, 3 \times 3, 5 \times 5$). Additionally, max pooling is also performed. The outputs are concatenated and sent to the next Inception module.



As stated before, deep neural networks are computationally expensive. To make it cheaper, the authors limit the number of input channels by adding an extra 1×1 convolutional before the 3×3 and 5×5 conv. Though adding an extra operation may seem counterintuitive, 1×1 conv are far more cheaper than 5×5 conv and the reduced number of input channel also help. Do not that however, the 1×1 conv is introduced after the max pooling layer, rather than before.

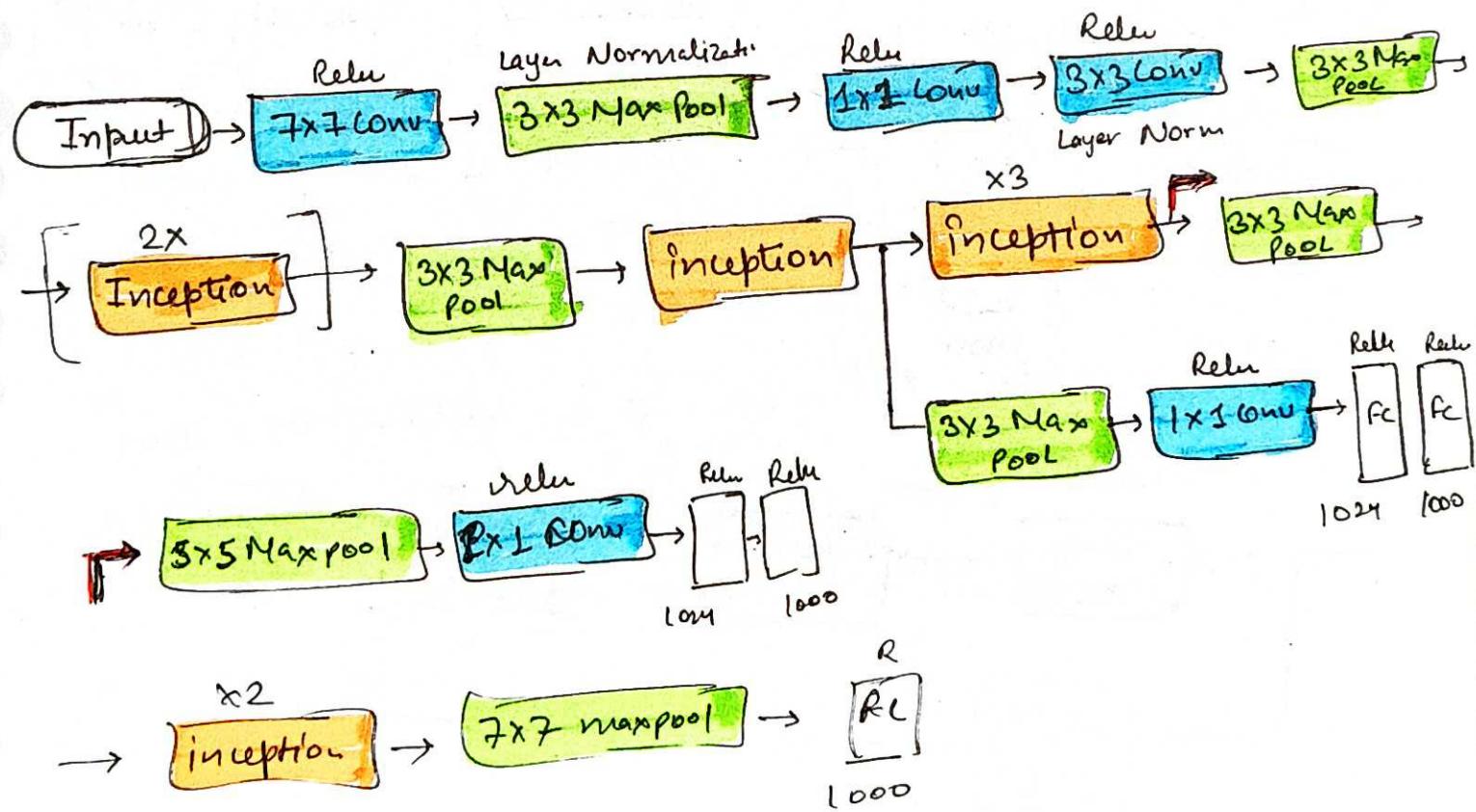


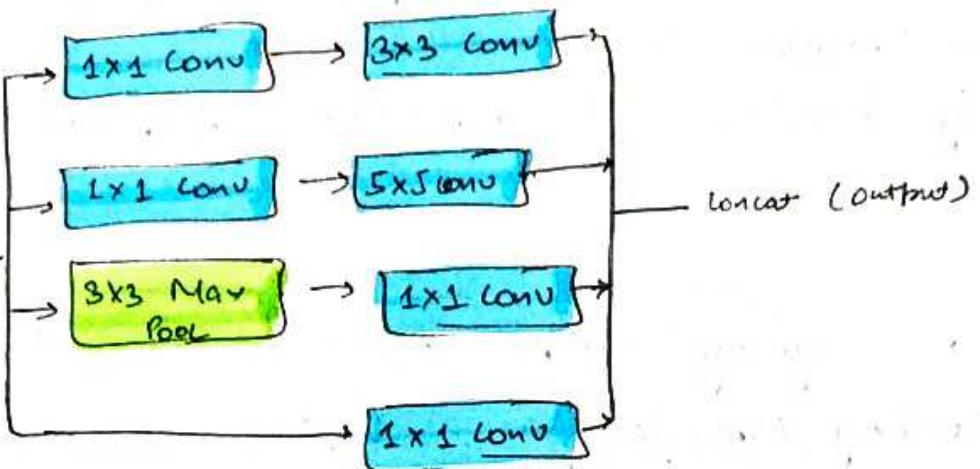
* Can also apply $\text{stride} = 2 \iff \text{Conv}(1 \times 1)$
 ↳ work same as low in Inception

Using the dimension reduced Inception module (104) a neural network architecture was built. This was popularly known as googleNet (Inception v3)

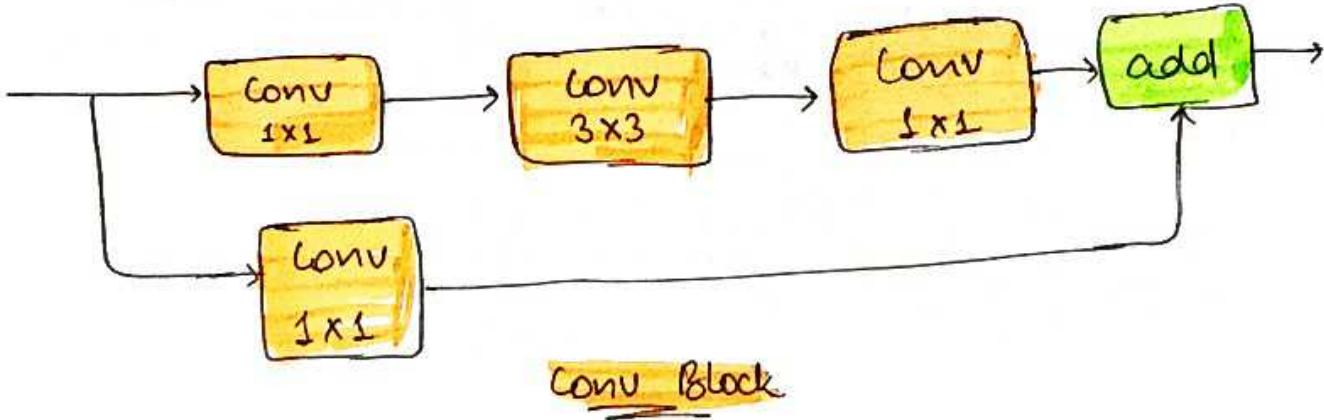
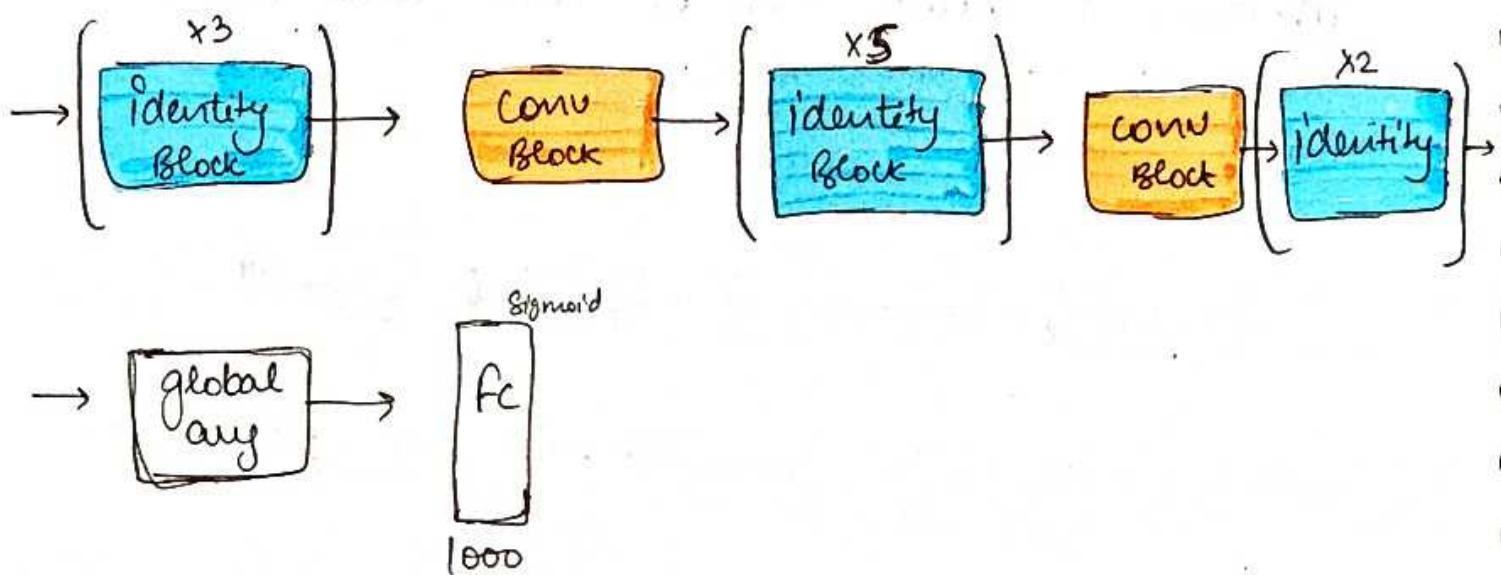
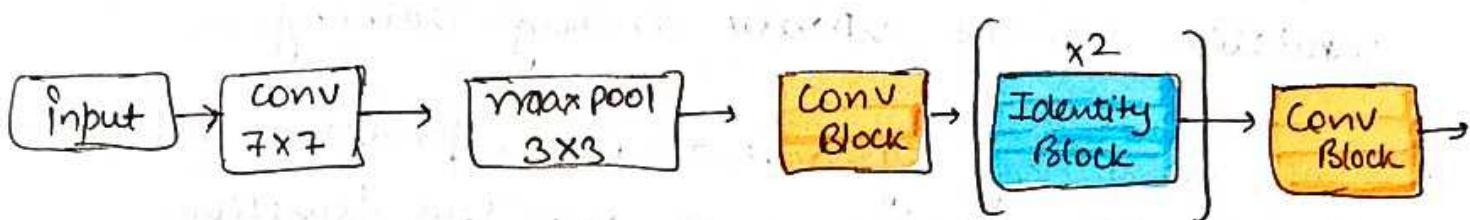
Inception v3 → going deeper has a caveat:
Exploding / Vanishing gradients:

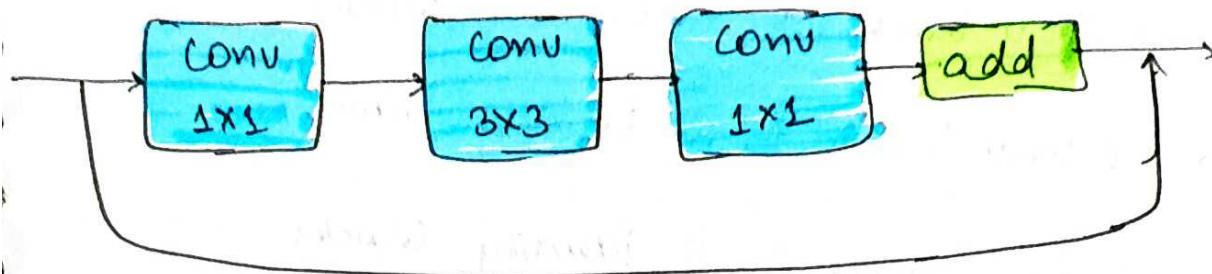
1. The exploding gradient is a problem when large error gradients accumulate and result in unstable weight updates during training.
2. The vanishing gradient is a problem when the partial derivative of the loss function approaches a value close to zero and the network couldn't train.





Resnet





Identity block

The degradation problem is addressed by introducing bottleneck residual blocks. There are 2 kind of residual blocks:

Identity block: consists of 3 convolution layer with 1×1 , 3×3 , and 1×1 kernel sizes, all of which are equipped with BN. The ReLU activation function is applied to the first two layers, while the input of the identity block is added to the last layer before applying ReLU.

convolution block: same as identity block but the input of the convolution block is first passed through a convolution layer with 1×1 kernel size and BN before adding to the last convolution layer of the main series.

Notice that both residual blocks have 3 layer. In total resnet -50 has 26 million parameters in 50 layers:

1 conv layer with BN then ReLU is applied

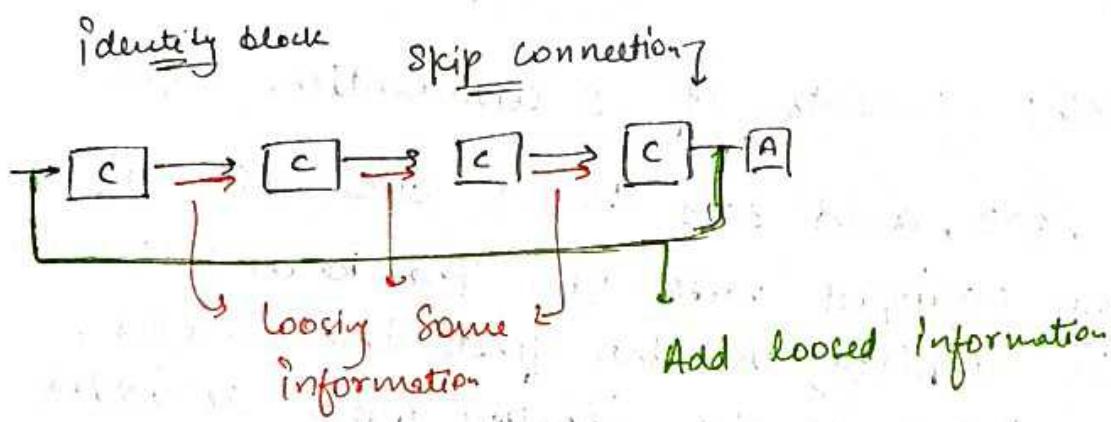
9 layers \rightarrow 1 conv block & 2 Identity blocks

12 layers \rightarrow 1 conv block & 3 Identity blocks

18 layer \rightarrow 1 conv block & 5 Identity blocks

9 layer \rightarrow 1 conv block & 2 Identity blocks

1 fully connected layer with softmax



CNN Architecture	Default input	Default output	No. of layers	No. of Param	New addition per
LeNet - 5	32x32x1	10	5	60k	conv layer
AlexNet	224x224x3	1000	8	60M	local Response norm
VGG-16	224x224x3	1000	16	138M	Very deep but still single thread
Inception v1	224x224x3	1000	22	7M	Auxiliary classif 2 inception Module
ResNet	224x224 X 3	1000	50	26M	Batch Norm 2 Residual block

Resnet solve the problem of vanishing gradients, which occurs in deep neural networks. As layers increase, the gradients become too small during backpropagation, making it hard for the network to learn effectively. Resnet introduces skip connections (or residual connections) that allow the model to bypass some layers, enabling better gradient flow and improving training in deep architecture.

Resnet Architecture

The ResNet architecture is built around the concept of residual learning, where layers learn the residual (or difference) betn the input and the output and the output rather than learning the full transformation. This is achieved using skip connection.

Here's a breakdown of the ResNet architecture:

1. Input layer: Initial convolutional layer with a large filter (e.g 7×7) and stride 2, followed by max pooling.
2. Residual Block:
 - Each block contains a few convolutional layers (typically two or three).
 - A skip connection (shortcut) is added, bypassing these layers and directly connecting the input to the output of the block.
 - The output of the convolutional layers is added to the input before passing through the next activation function.

The residual block can be either:

→ Identity block: If the input and output dimension are the same.

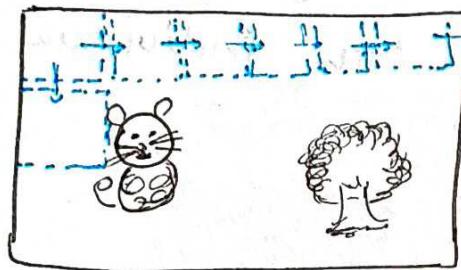
→ Convolutional blocks: If dimensions are the same.

3. Stacked Residual Blocks: Multiple residual blocks are stacked, increasing the depth of the network without losing gradient flow.

4. Global Average Pooling: After all residual blocks, a global average pooling layer reduces the spatial dimensions.

5. Fully connected layer: Finally, a fully connected layer produces the classification output.

Object Detection → R-CNN
↳ Fast RCNN
↳ Faster RCNN
↳ YOLO



filters are moving and finding images. (Basic approach)
Iteration of this process is ↑↑ more.
This process is called exhaustive search.

RCNN → Region based CNN

* Regions

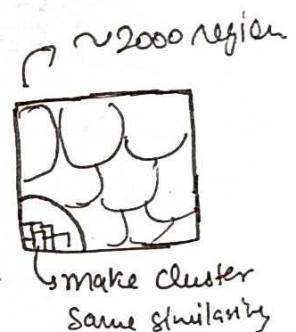
↳ Exhaustive Search (Basic)

↳ Selective Search → greedy algorithm →

↳ Region proposal Method

used in RCNN is selective search

~ 2000^{propose} region.



Region Proposal

The first stage of the R-CNN pipeline is the generation of "region proposals" or regions in an image that could belong to a particular object. The authors use the selective search algo. The selective search algo works by generating sub-segmentations of the image that could belong to one object.

based on color, texture, size and shape - and iteratively combining similar regions to form objects. This gives 'object proposal' of different scales. Note the R-CNN pipeline is agnostic to the region proposal algorithm. The authors use the selective search algo to generate 2000 category independent region proposals (usually indicated by rectangular regions or bounding boxes) for each individual images.



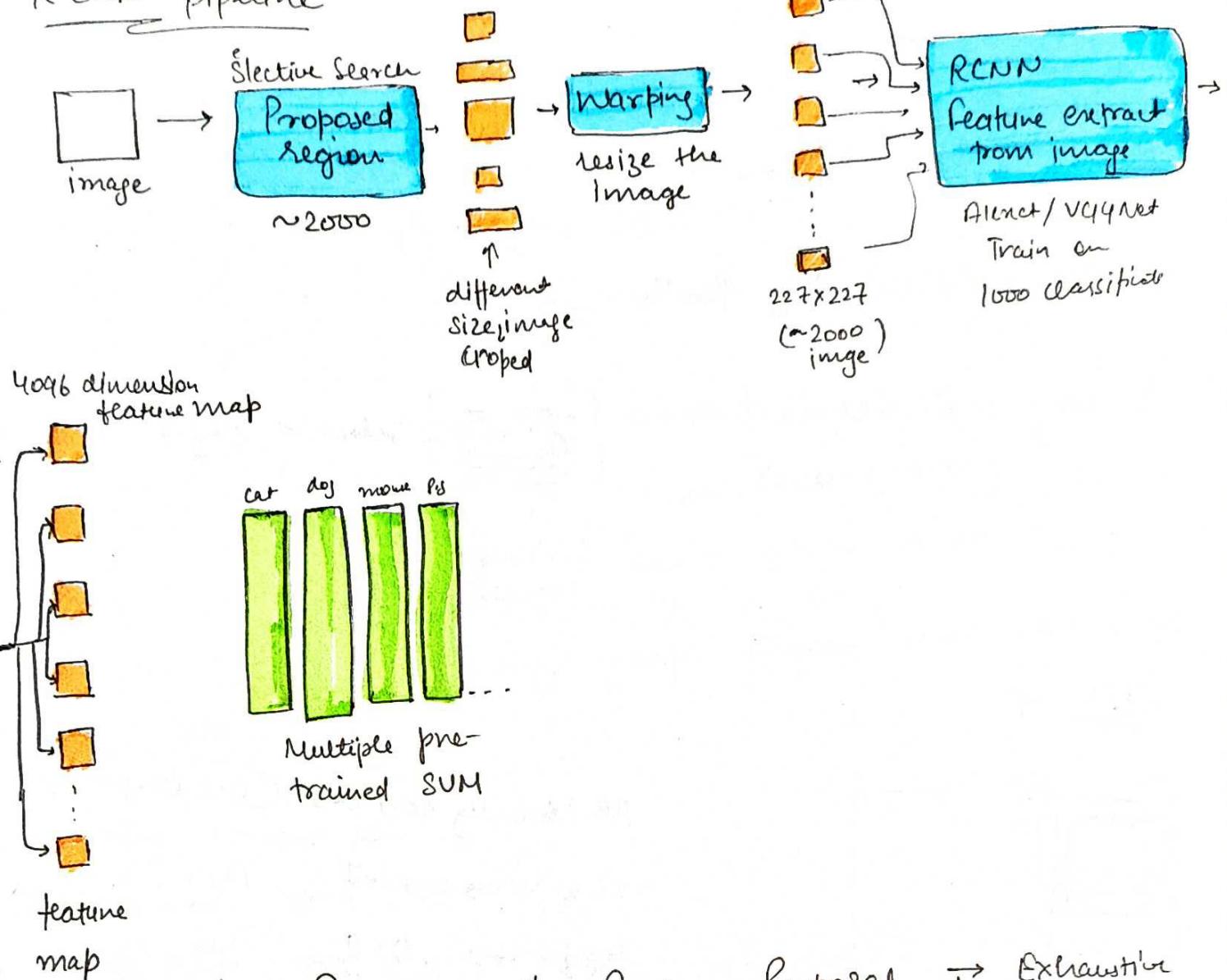
R-CNN Pipeline Stages

Stage 1: Feature extraction from the Region Proposals

Stage 2: SVM for object classification

Stage 3: Bounding box regression

R-CNN pipeline



RCNN → (i) Regions Proposal → Exhaustive
selective search

(ii) Warping and resize

(iii) Pretrained CNN → Alex Net
↳ VGG Net

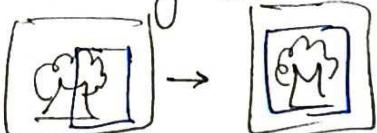
(iv) Pretrained → SVN_{cat}, SVN_{dog}, SVN_{person}

(v) Clean Up → Maybe same image
come many time

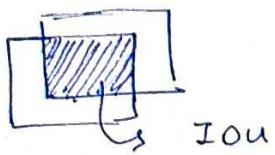


Same Image → 2 time

(vi) Bounding Box → fine tune

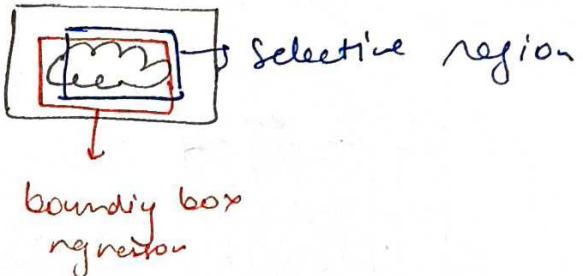


IOU → Intersection of union

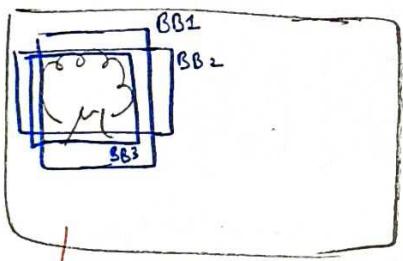


Metric → Mean Avg Position

Bounding box regression →

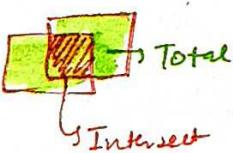


Clean up



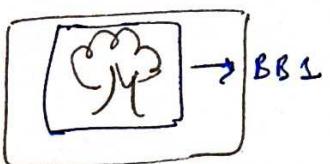
BB (bounding box) 1 → select tree
BB 2 → 0.7 tree
BB 3 → 0.8 tree

Find IOU →



If Total divide by intersect is greater than 50 then

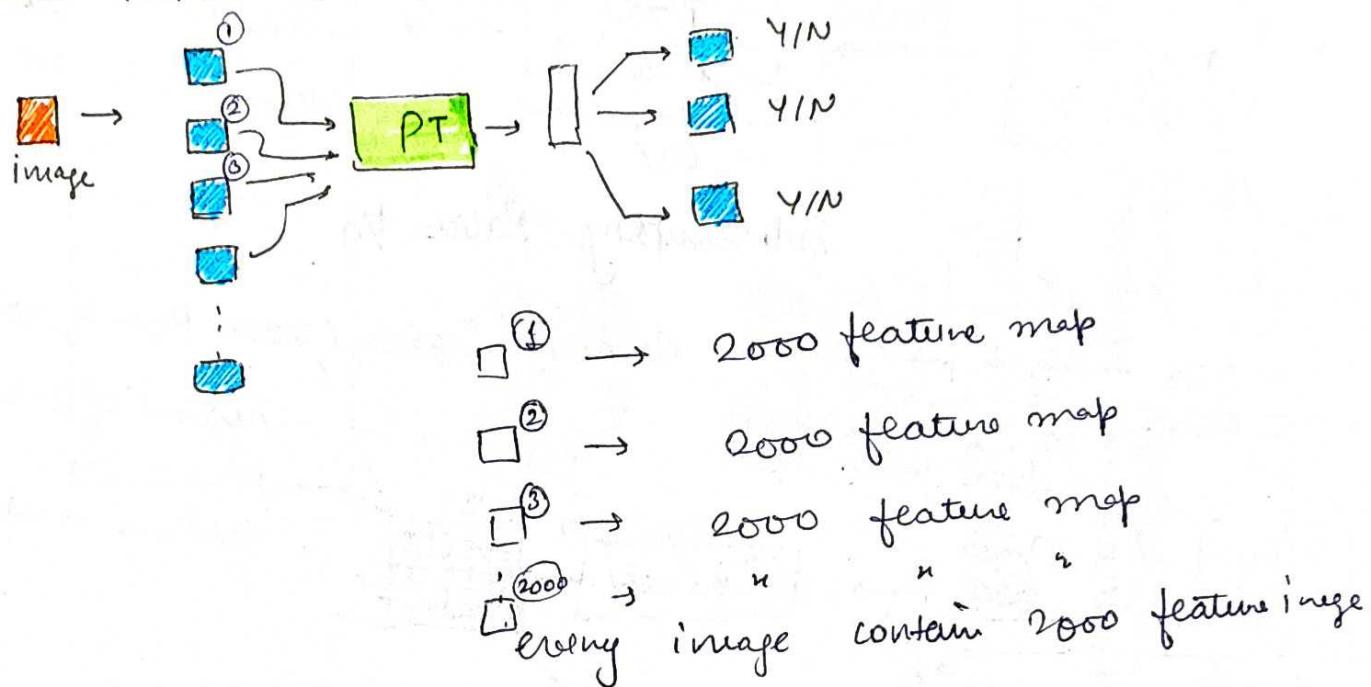
we drop BB2 and BB3



Fast RCNN

(109)

In RCNN (2014)



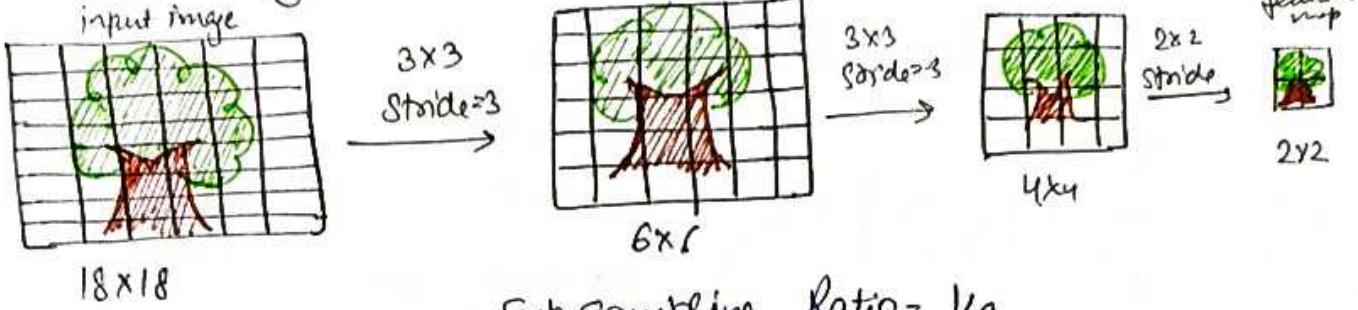
Problem:

- ① Very computation
- ② Extreme time consuming
 - one image took → 47 second
- ③ It has a complex multi stage training pipeline.

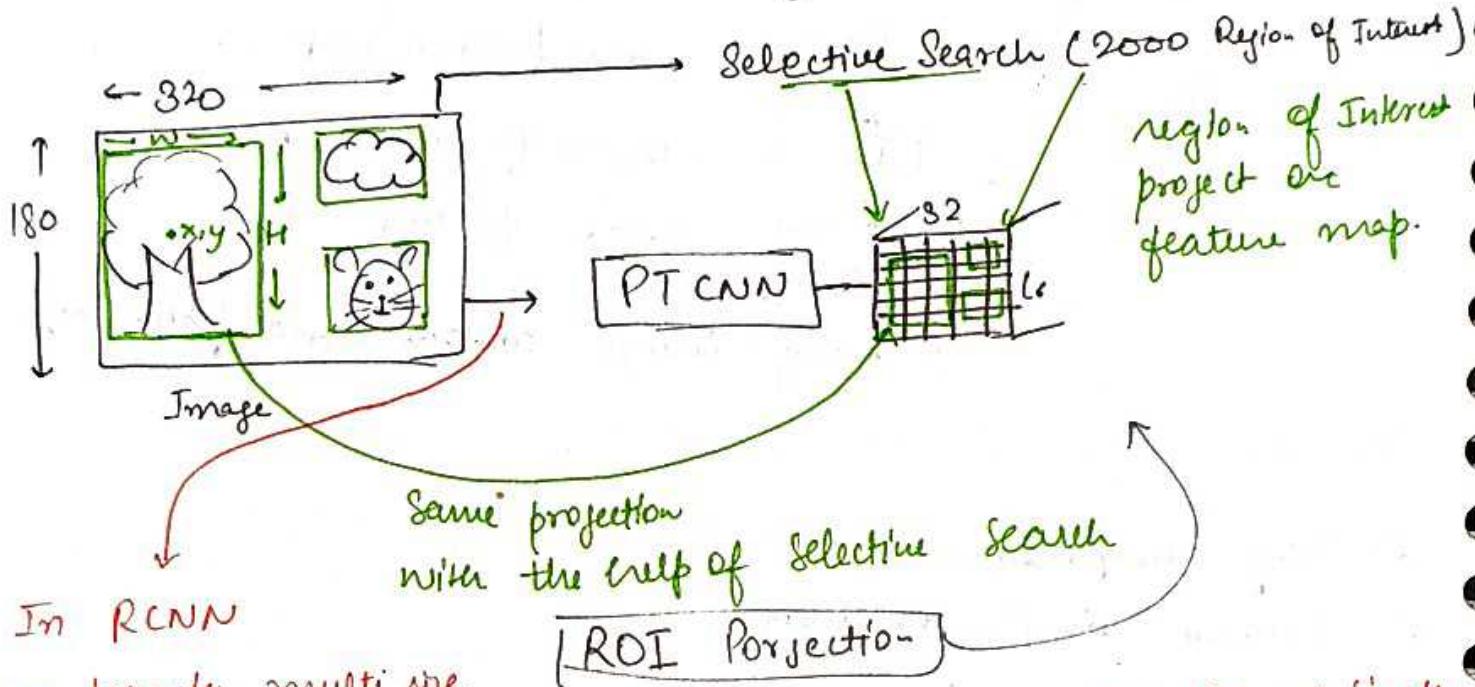
Fast - R-CNN

1. Sub Sampling
2. ROI Protection (ROI \rightarrow Region of Interest)
3. ROI Pooling

Sub Sampling

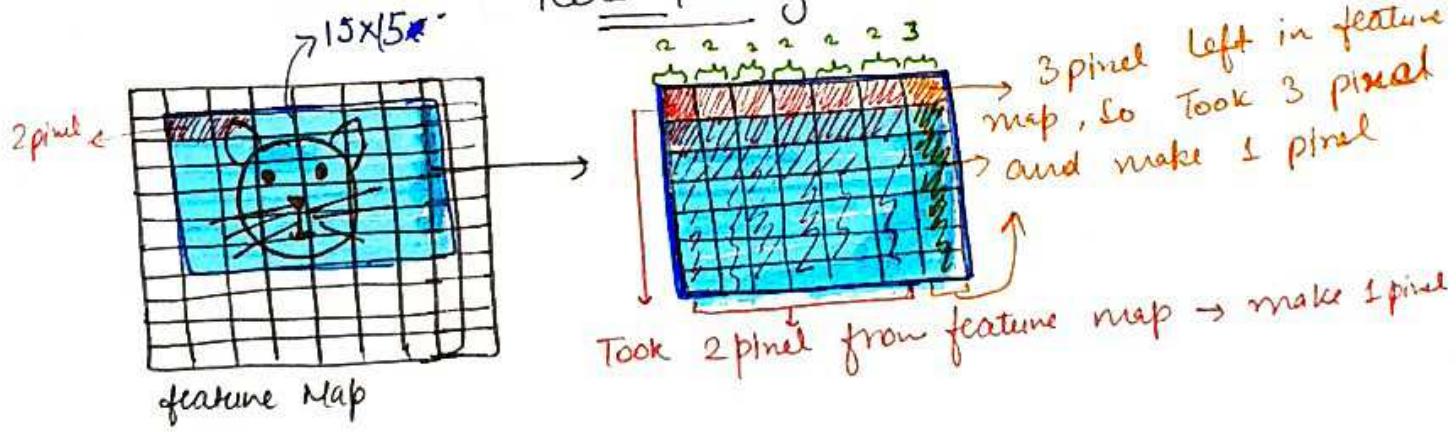


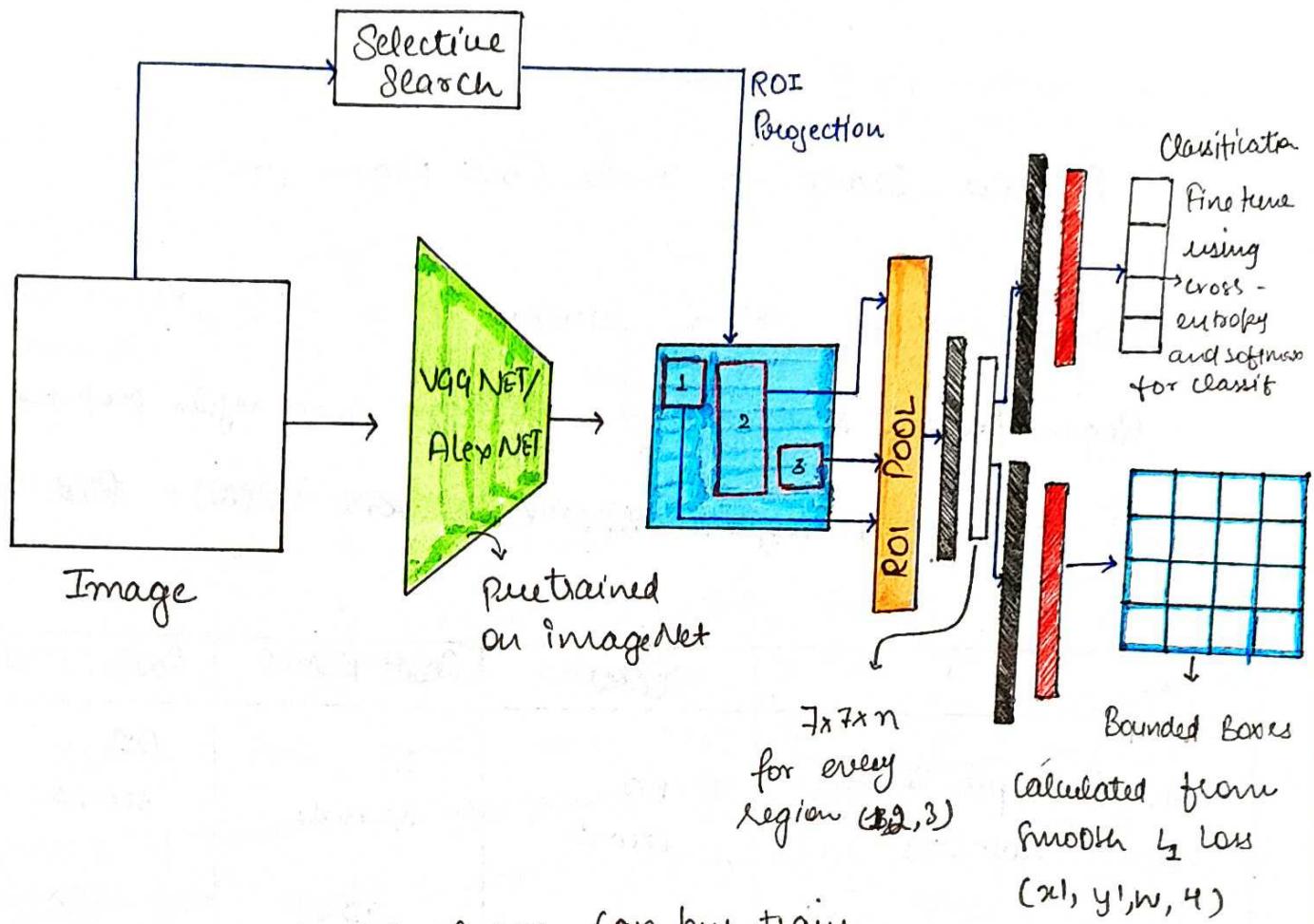
Sub Sampling Ratio = $1/9$



ROI Pooling And PTCNN helps to reduce dimension and selective search help to project ROI on feature map. ROI size also reduce accordingly to feature map.

ROI Pooling





In classification \rightarrow we have car, bus, train
 So Total class is 4 because one class is for background

In Bounded box \rightarrow find L_1 loss where a bounded box is ^{and initial} perfect or not.

	RCNN	Fast RCNN
Training Time <u>(Speedup)</u>	84 hours 1x	9.5 hours 8.8x
Test time per image <u>(Speedup)</u>	47 second 1x	0.32 seconds 146x
Test time per image with selective search <u>(Speedup)</u>	50 second 1x	2 second 25x

Faster R-CNN

Problems with Fast R-CNN

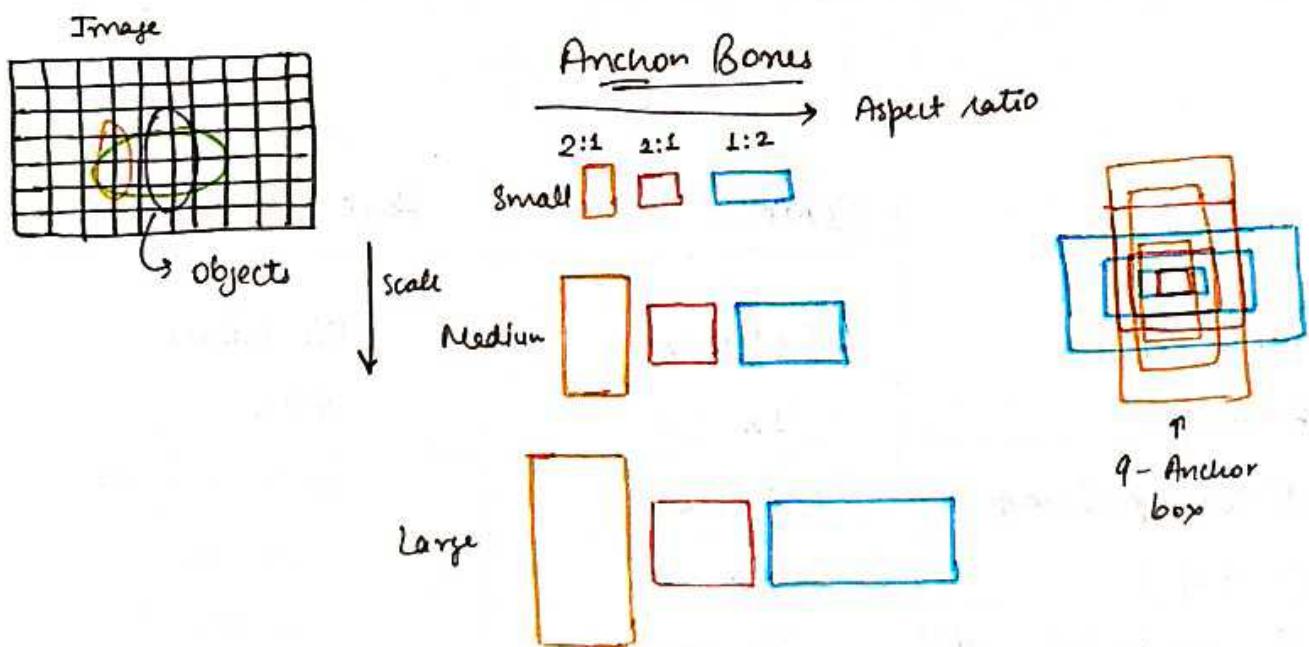
Selective Search → create 2000 Region proposed

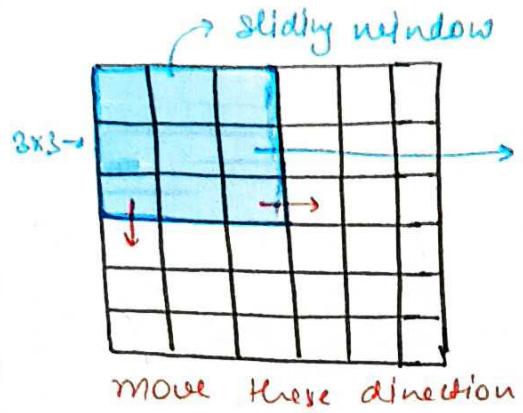
In Faster R-CNN → replace Selective Search with Region Proposal Network

Region Proposal Network → less than 2000 region proposed

Faster R-CNN → Region Proposal Network (RPN) + Fast R-CNN

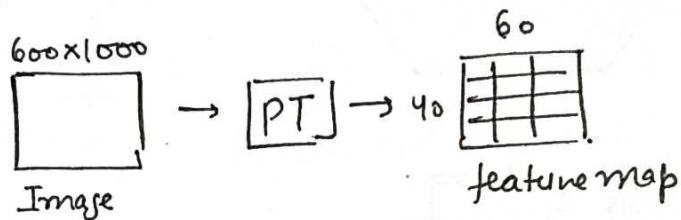
	RCNN	Fast RCNN	Faster RCNN
Test time per image with Proposals (Speedup)	50 seconds 1x 66.0	2 seconds 25x 66.9	0.2 seconds 250x 66.9
mAP (Pascal VOC 07)			





9 region of Interest for one sliding window.

No. of slide ↑↑ = region of Interest ↑↑



600x1000, we would have a feature map of size - 40x60. And using a diff sliding window at each position for all the 40x60 values in the feature Map, we end up having $40 \times 60 \times 9 = \sim 20000$ proposals. Compared to Selective Search which give just 2000 proposals, we have almost 10 times more proposals. This would be computationally more expensive and also would have more false positives.

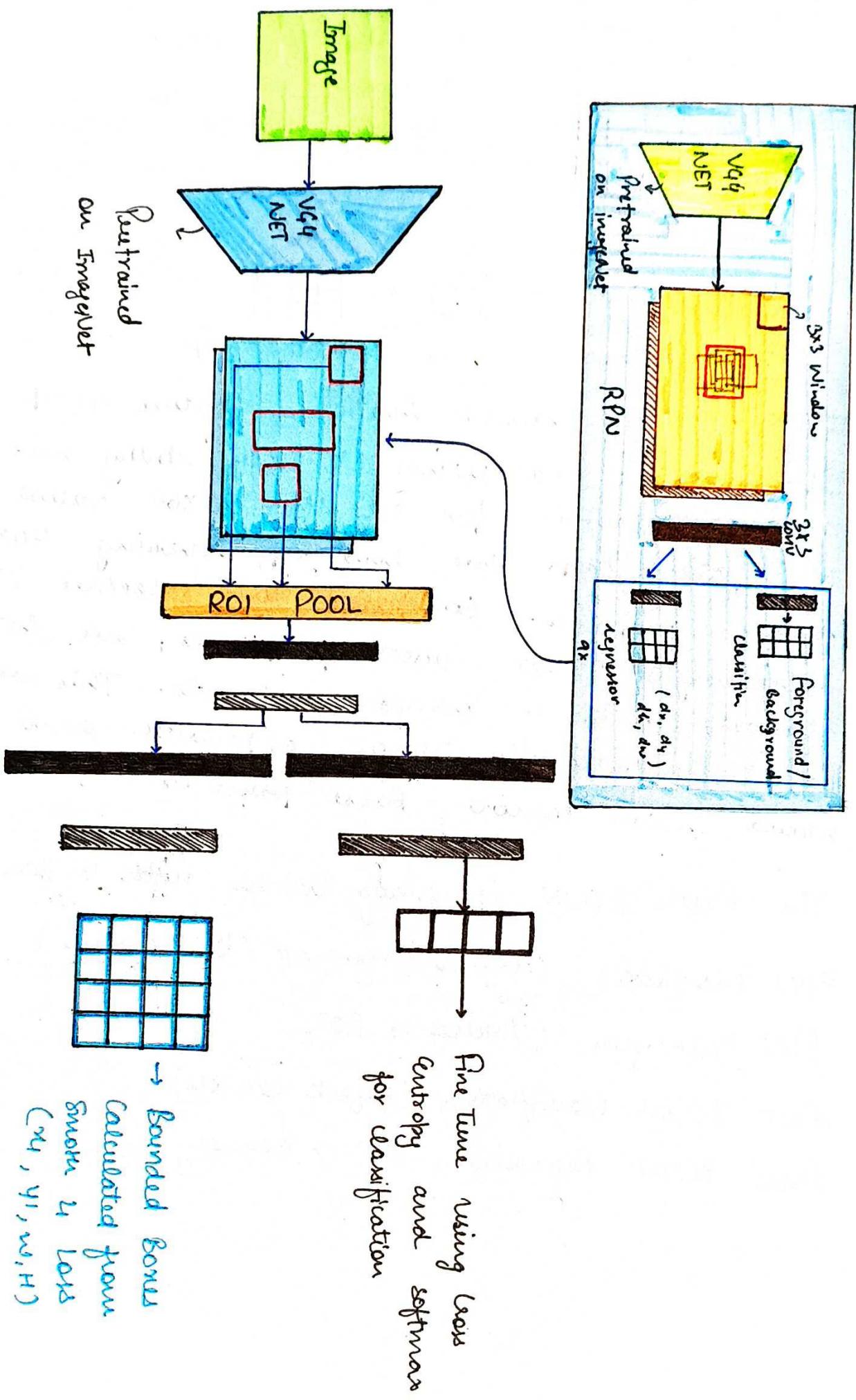
The Faster R-CNN is jointly trained with 4 losses
RPN classification (Object foreground / background)

RPN Regression (Anchor \rightarrow ROI)

Fast RCNN Classification (Object classes)

Fast RCNN Regression (ROI \rightarrow Bounding Box)

Faster R-CNN

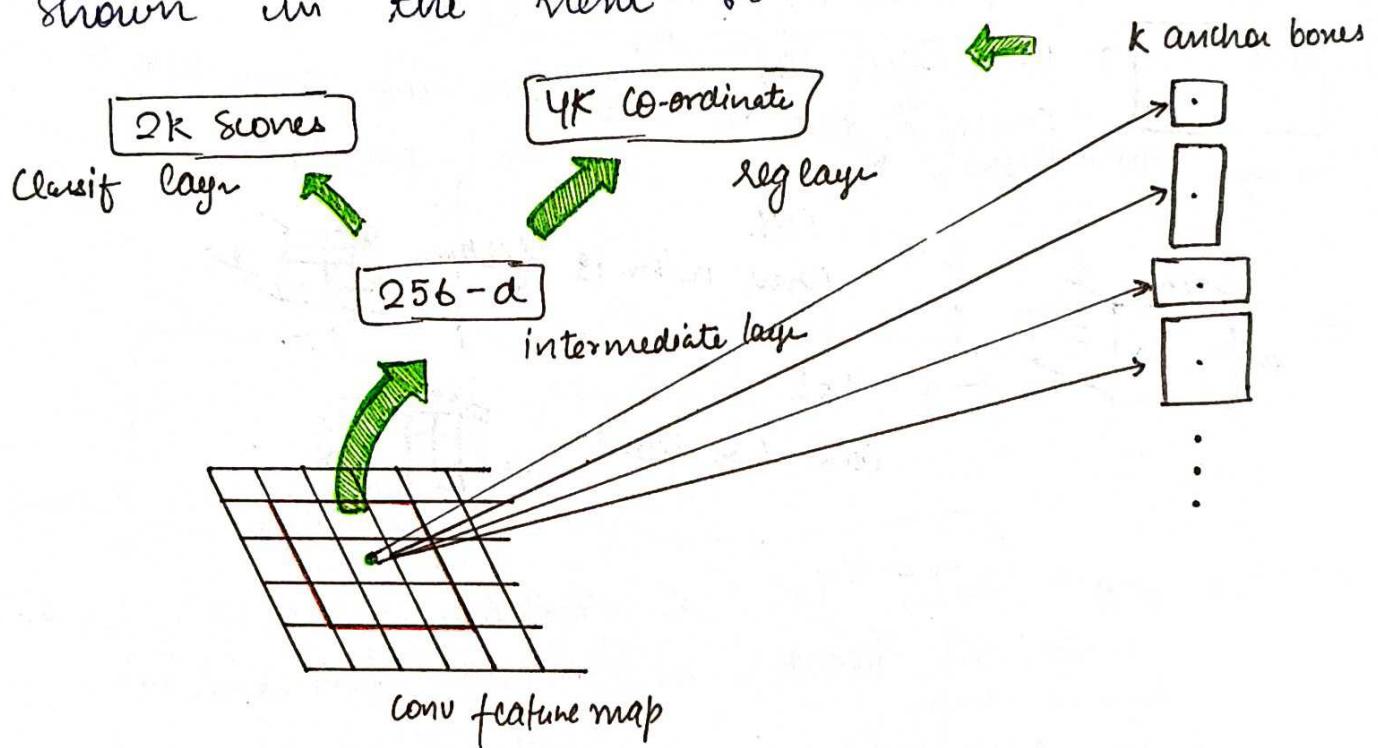


Anchor → depend on

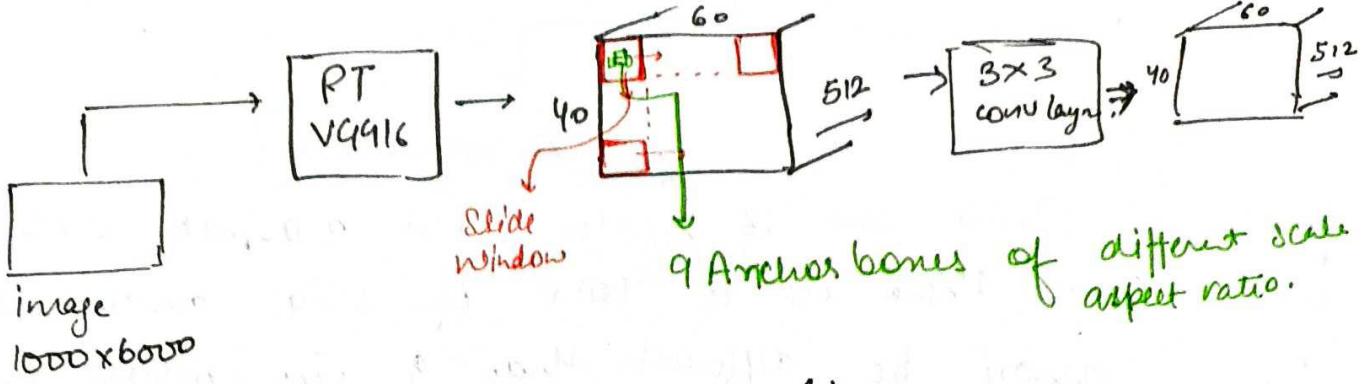
Scale

Aspect Ratio

Generally, there are 3 scales and 3 aspect ratios and thus there is a total of $k=9$ anchor boxes. But k may be different than a in other words, k regions are produced from each region proposal, where each of the k regions varies in either the scale or the aspect ratio. Some of the anchor variations are shown in the next figure.



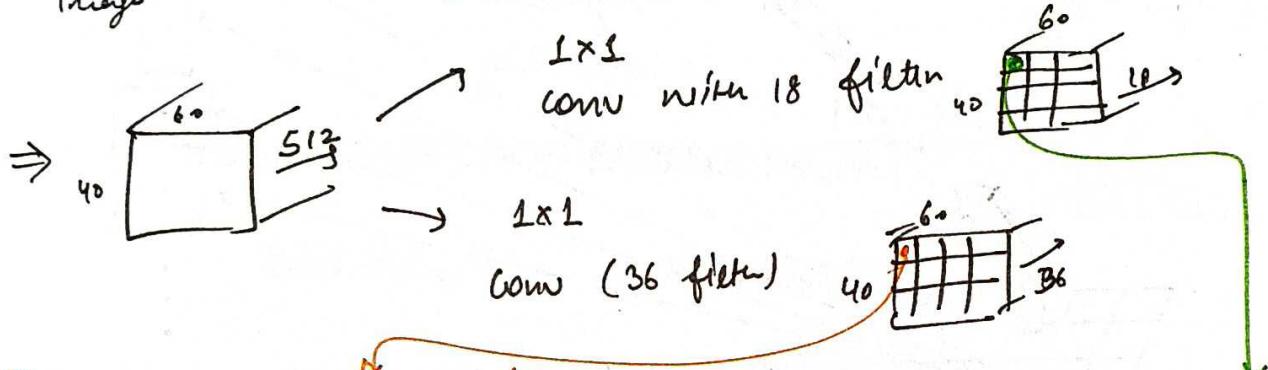
Faster R-CNN Detailed



$60 \times 40 \times 9 = \sim 21k$ crossborder $\sim 6k$
 Non Max Suppression \rightarrow IOU \rightarrow If < 0.5 not accept
 If > 0.5 accept
 Intersection of Union \downarrow confidence

crossborder

image \rightarrow ROI (Region of Interest)
 crossing border. So, decrease from 21k to 6k

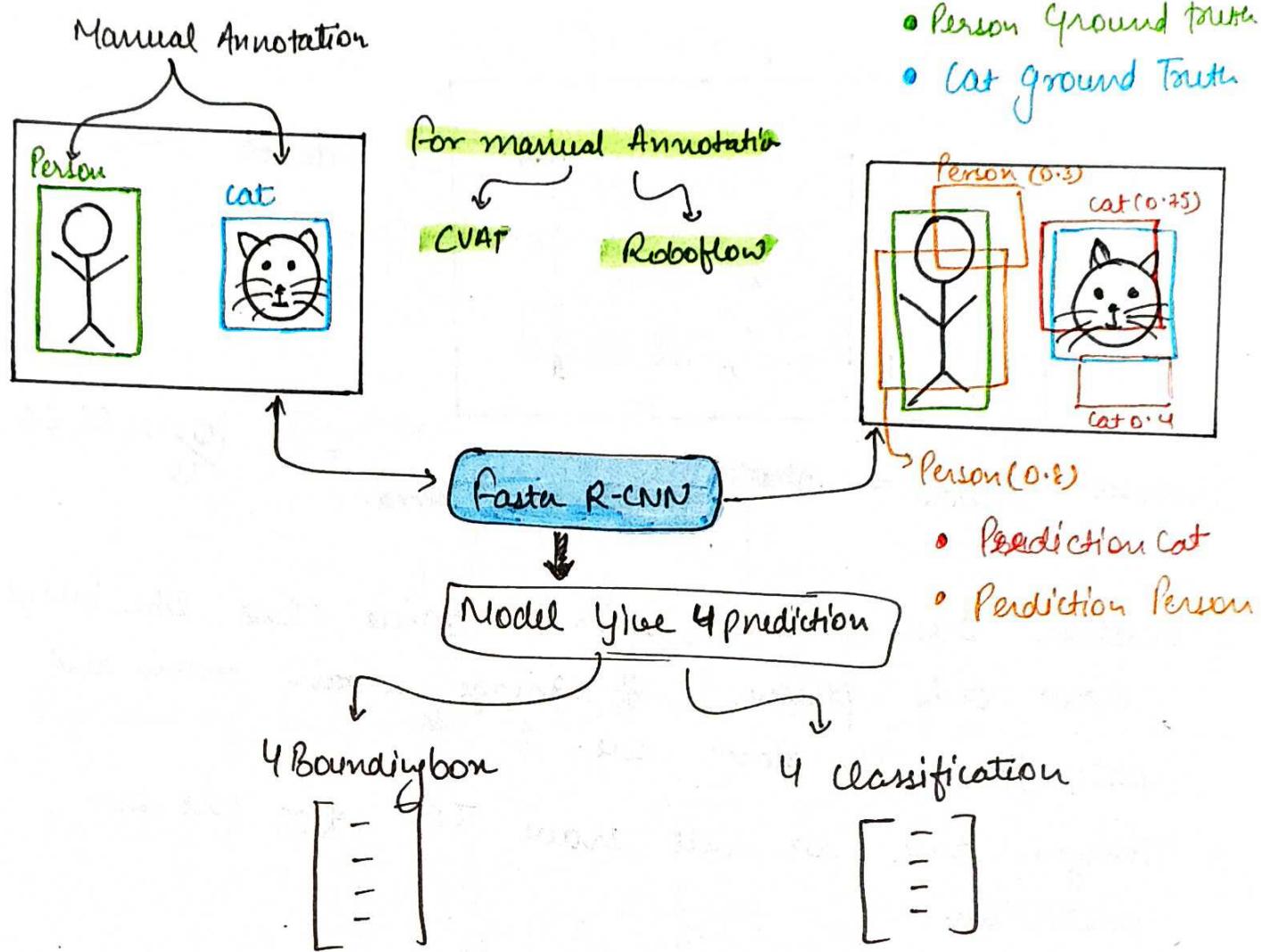


9D boxes, define co-ordinate
 d_x, d_y, d_w, d_h of
 Anchor box

9 Anchor Boxes \rightarrow 1, 2, 3, 4, 5, 6, 7, 8, 9
 background or foreground
 we have prob not 0 or 1

IOU, Non Max Supression And Map

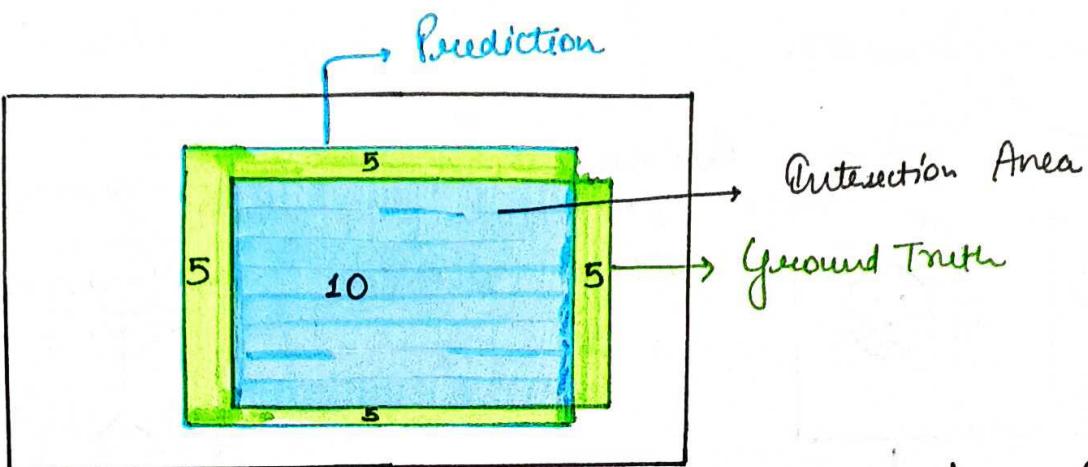
113



Model Performance

- ① We usually detect a lot more objects than we usually have in that image.
 - Selective [confidence vary]
 - Anchor
- ② We need a way to decide → Threshold
Evaluate model → Mean Avg., Precision.
- ③ Because of Point 1, we also need a clean up process → No. of prediction limited and Precision and Recall also increase

Intersection over union (IOU)

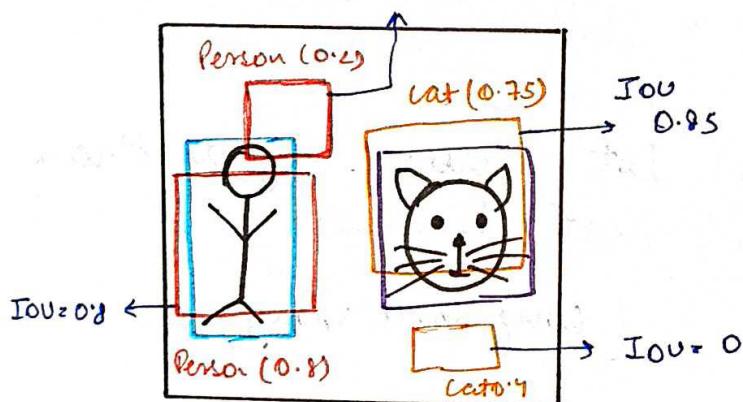


$$\text{Intersection/union} = \frac{\text{blue Area}}{\text{blue+green Area}} = \frac{10}{15} = 66.66\%$$

Intersection over union use → Same class like image have only person. If image contain person and cat then we don't use.

- * Through this we will have IOU for all our predictions.

$$\text{IOU} = 0.25$$

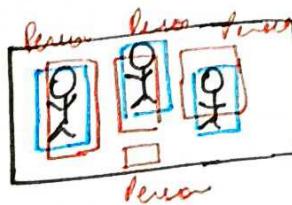


- ground truth: cat
- ground truth: person

Mean Average Precision

114

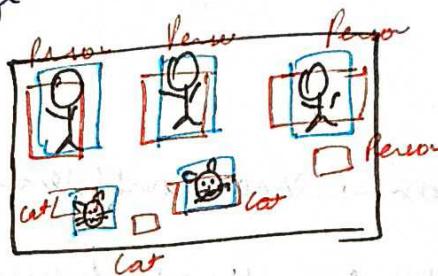
$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$



Object detection class of Person

$$= \frac{3}{3+1} = \frac{3}{4} = 75\%$$

let say, Image \rightarrow both cat and dog



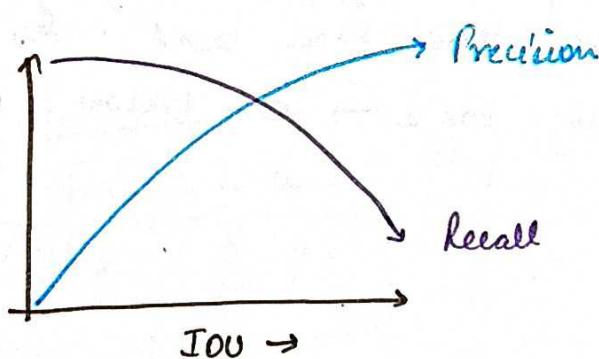
Precision of Person ~~2~~ = $\frac{2}{3} = 66.67\%$

Precision of Cat = $\frac{2}{3}$

Mean Average Precision $\rightarrow \frac{1}{k} \sum_{i=1}^n AP_k$

$$\boxed{\text{Precision} = \frac{\text{TP}}{\text{All detects}}}$$

$$\boxed{\text{Recall} = \frac{\text{TP}}{\text{Actual Ground Truth}}}$$



0.7 0.8 0.9 \rightarrow Prob

New threshold = 0.5

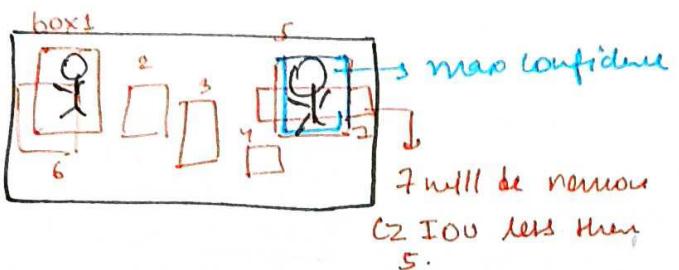
all prediction are above threshold that's why our prediction is true.

Avg Precision (Person) = 0.75 = $\frac{3}{4}$

let say, My threshold = 0.85
now, our two prediction are false

$$\text{Avg Precision} = 0.25 = \frac{1}{4}$$

Non Max Suppression



box	confidence level
1	$\rightarrow 0.9$
2	$\rightarrow 0.4$
3	$\rightarrow 0.7$
4	$\rightarrow 0.8$
5	$\rightarrow 0.95$

Step 1 → To remove all bound boxes with confidence level less than α (where $\alpha = 0.75$)

removed = 2, 3

Person class

↳ Bounding box → max confidence (Box Max)

↳ All bounding box that have $IOU > 0.4$ (any value) with max box will be removed.

→ We remove 7 cz IOU less than 5 bound box. So we cleaned second image of person (only 5 bound box there)

→ Now, same process apply on first person → box 1 have higher IOU than box 6. So, box 6 will remove. Only box 1 → on person (cleaned person 1 too)

let discuss more about Precision and Recall

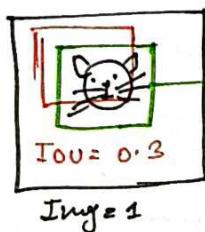
115

$$\text{Precision} = \frac{\text{TP}}{\text{All detection}}$$

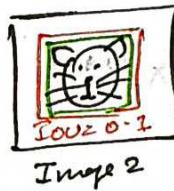
$$\text{Recall} = \frac{\text{TP}}{\text{Actual Ground Truth}}$$

True Positives

↳ Prediction = $\text{IOU} > \text{Threshold}$

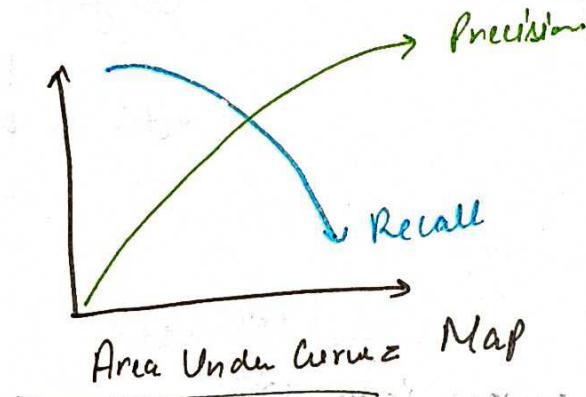


Ground Truth



our Threshold = 0.5

- Image 1 reject Cz $\text{IOU} \leq$ less than threshold ($0.3 < 0.5$)
- Image 2 accepted as TP Cz $\text{IOU} \geq$ greater than threshold ($0.7 \geq 0.5$)



$$\text{Map} = \frac{1}{K} \sum_{k=1}^K \text{AP}_k$$

* We find this for every class