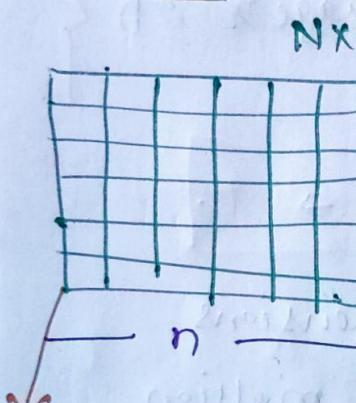


# CNN (Convolution Neural Networks)



- It is used in text / images / visualization.
- It is particularly well-suited for image recognition and processing tasks.
- It is made up of multiple layers.

## Image



$N \times N$  (pixels)

Gray scale  
↓  
Black & White

$6 \times 6$

$[0 - 255]$

tells about

intensity of image

Color word  $\rightarrow \{6 \times 6 \times 6\}$

$\left. \begin{array}{c} R, G, B \\ Red \\ Green \\ Blue \end{array} \right\}$

$[0 - 255]$

Convolution means we use some filters to on given image (gray scale or RGB) and traverse filters to retrieve some features from images.

\* Convolution  $\rightarrow$  retrieve features from image by applying filters on image

## Q Why Convolutions ?

- Ans → Spatial invariance or less features →  
→ the Spatial features of a 2D image are lost when it is flattened to a 1D vector but, before feeding an image to the hidden layers of an MLP, we must flatten the image matrix to a 1D vector, ~~as we~~ this implies that all of the image's 2D information is discarded.

## (2) Increase in Parameters Issue →

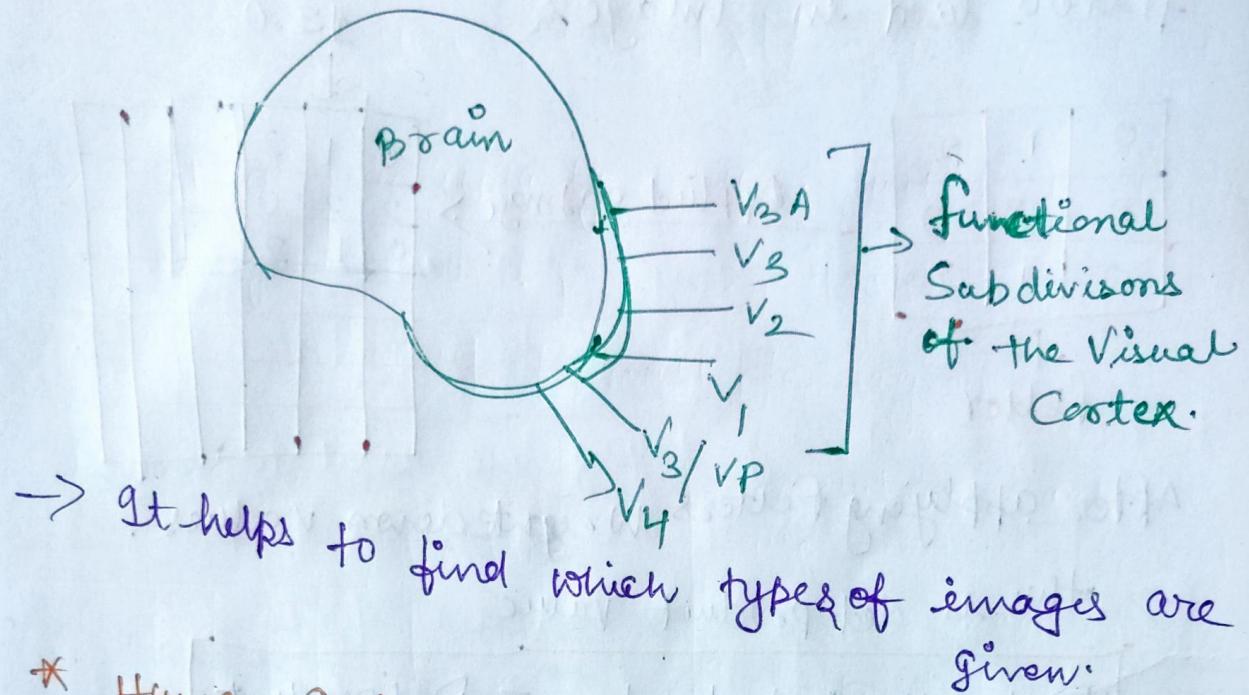
- If we have an image with dimensions  $1000 \times 1000$ , it will yield 1 million parameters for each node in the first hidden layer. So, ~~the~~ flattened image is necessary.

## (3) Local Connected Neural Net →

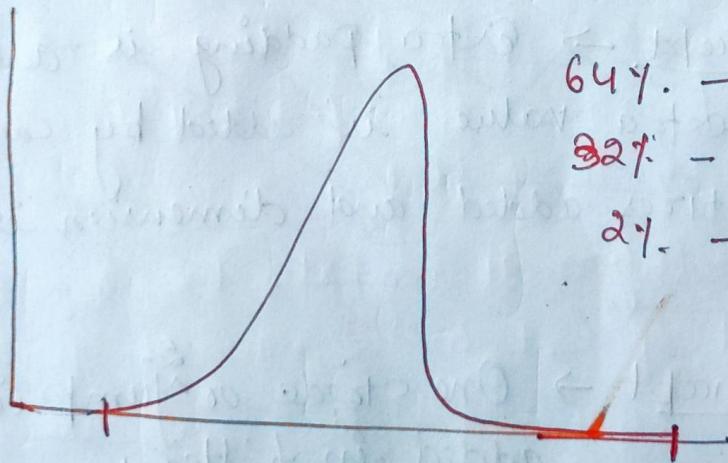
It is used in image recognition to focus on parts of images.

Flatten → Used to flatten images to 1D vector.

## \* Human Brain Visual Cortex Processing →



## \* Human Eye Color Sensitivity →



Cones are  
64% - Red sensitivity  
32% - Green sensitive  
2% - blue sensitive

## \* Visualizing the Process →

- ① Simple Convolution → By using filters select features from images.  
→ Gradient used in backpropagation also here.

\* Matrix calculations → filters used in images.

0	1	2
2	2	0
0	1	2

filter

Applied on Image

12.0				
0	1	2		
2	2	0		
0	1	2		

After applying filters weighted down values.

→ this is dot product value

Filter \* Input value = Dot product.

write

\* Padding concept → Extra padding is added so, one extra value is added by convolution so, extra features added and dimension is preserved.

\* Stride Concept → One stride or jump is extra added in filter.

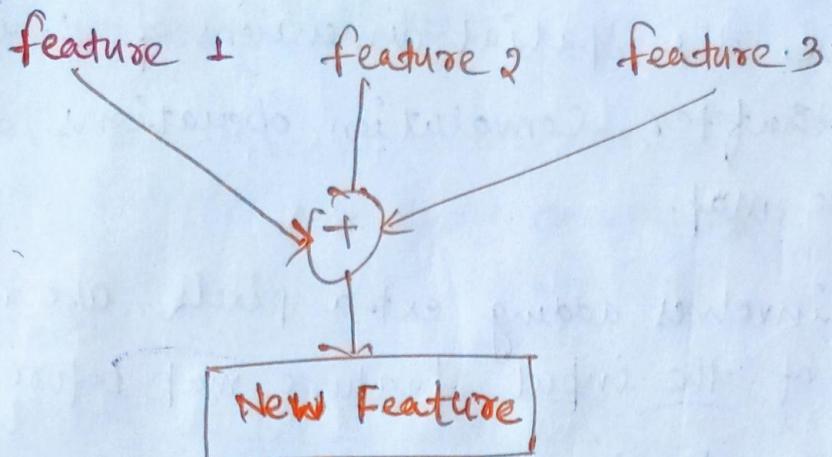
so,  $2 \times 2$  convolution features are preserved.

2-jump.

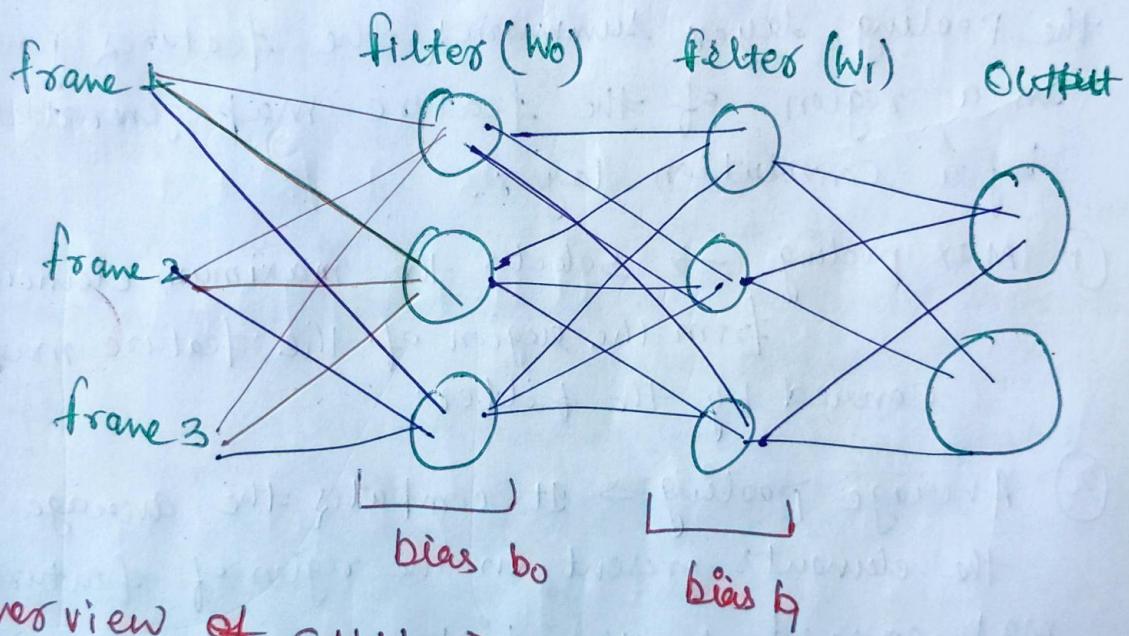
\* Feature Accumulation →

filters 1 → filters 2 → filters 3  
↓            ↓            ↓  
features      feature 2      feature 3

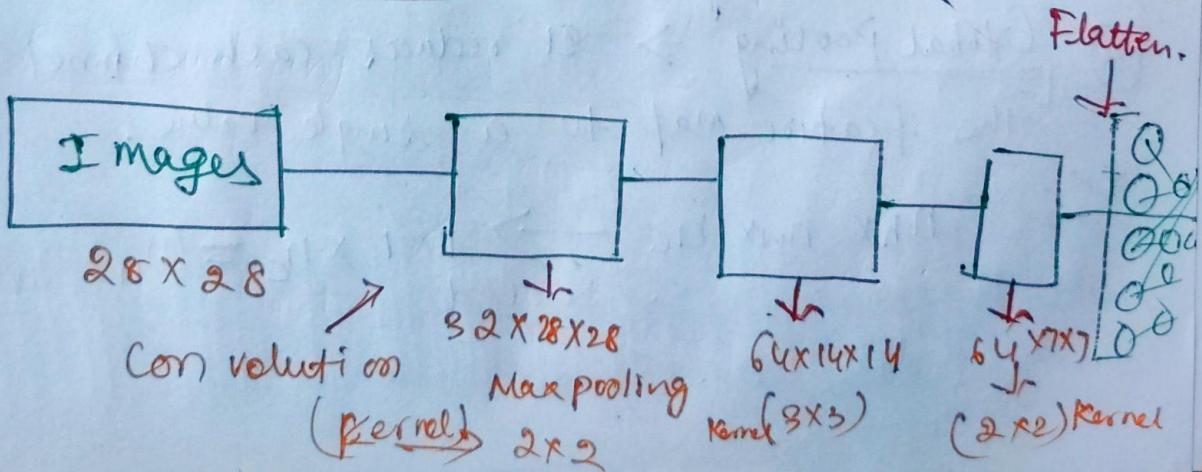
## \* Feature Aggregation →



## Convolution operation →



## Overview of CNN →



\* Padding  $\Rightarrow$  It is a technique used to preserve the spatial dimensions of the input image after convolution operations on a feature map.

- Padding involves adding extra pixels around the border of the input feature map before convolution.

\* Pooling  $\Rightarrow$  Pooling layers are used to reduce the dimensions of the feature maps.

The Pooling layer summarizes the features present in a region of the feature map generated by a convolution layer.

- ① Max pooling  $\rightarrow$  Selects the maximum element from the region of the feature map covered by the filter.
- ② Average Pooling  $\rightarrow$  It computes the average of the elements present in the region of feature map covered by the filter.
- ③ Global Pooling  $\rightarrow$  It reduces each channel in the feature map to a single value.

Ex-  
 $\rightarrow n_h \times n_w \times n_c \rightarrow 1 \times 1 \times n_c = (n_h \times n_w)$   
Selected

## Advantages of pooling →

- ① Dimensionality reduction
- ② Translation invariance → position of an object in the image does not affect the classification result.
- ③ Feature selection

## Disadvantage of pooling ⇒

- ① Information loss
- ② Over-smoothing
- ③ Hyperparameter tuning
- ④ Stride in CNN

→ It is a parameter of the neural network's filter that modifies the amount of movement over the image or video.  
 → the no. of pixels turning to the input matrix is known as the strides.

$$\text{no of strides} = \text{no. of filters}$$

## Important formulas

### ① Input dimensions (IMAGE) ⇒

Width ( $w$ ) → Width of image

Height ( $h$ ) → Height of image

channels ( $C$ ) → the no. of color planes in the input image (1 for grayscale  
3 for RGB)

## Convolutional Layer $\rightarrow$

Kernel size (K) - width & height of the convolutional kernel or filter.

Padding (P)  $\rightarrow$  the no. of zeros added to the border of the input image to preserve spatial dimensions.

Stride (S) : the step size at which the kernel moves across the input image.

No. of filters (N)  $\rightarrow$  no. of convolutional filters applied in the layer.

$$\text{Width (w\_out)} : ((w-K+2P)/s) + 1$$

$$\text{Height (h\_out)} : ((h-K+2P)/s) + 1$$

$$\text{Channels (c\_out)} : N$$

## Pooling layer $\rightarrow$

Pool size (P) : width & height of the pooling

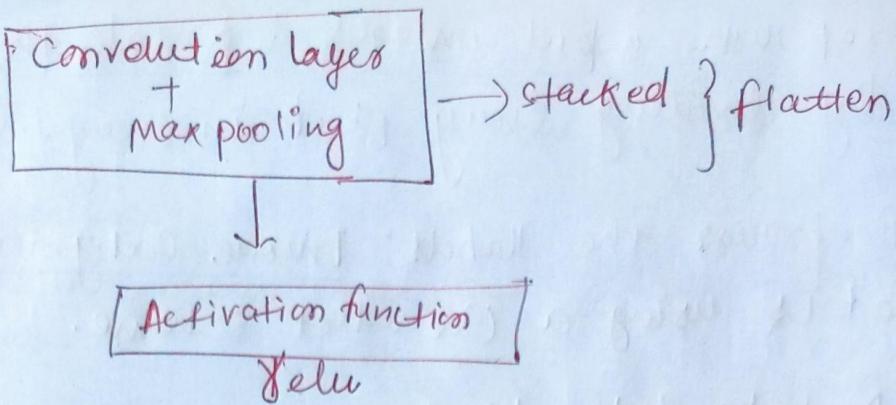
Stride (S) : - the step size at which the Pooling window moves across the input.

## Output Dimension $\Rightarrow$

$$(w\_out) : ((w-P)/s) + 1$$

$$(h\_out) : ((h-P)/s) + 1$$

$$(c\_out) : C \text{ (remains the same)}$$



\* ~~Kernel~~  $\Rightarrow$

\* Image Kernel  $\rightarrow$

A image kernel is a small matrix used to apply effects like the ones you might find in Photoshop or Gimp, such as blurring, sharpening, outlining or embossing.

\* Practical Intaition with code in CNN

① Load MNIST Database

MNIST is one of the most famous datasets in the field of machine learning.

- It has 70,000 images of hand-written digits.
- Very straight forward to download.
- Images dimensions are  $28 \times 28$
- Grayscale Images.

② Visualize the first six Training Images.

③ View an image in more detail

(4) Preprocess input images : Rescale the images  
by dividing every pixel in every image by 255

(5) Preprocess the labels: Encode Categorical Integers  
Labels using a One-hot scheme.

(6) Reshape data to fit our CNN (and input shape)

(7) Define the Model Architecture  $\Rightarrow$

filters - no. of filters

Kernels - no. specifying both height & width of  
the (square) convolution window.

Stride - by default set to 1.

Padding - by default - Valid

Activation - (relu)

(8) Compile Model

(9) Train the Model

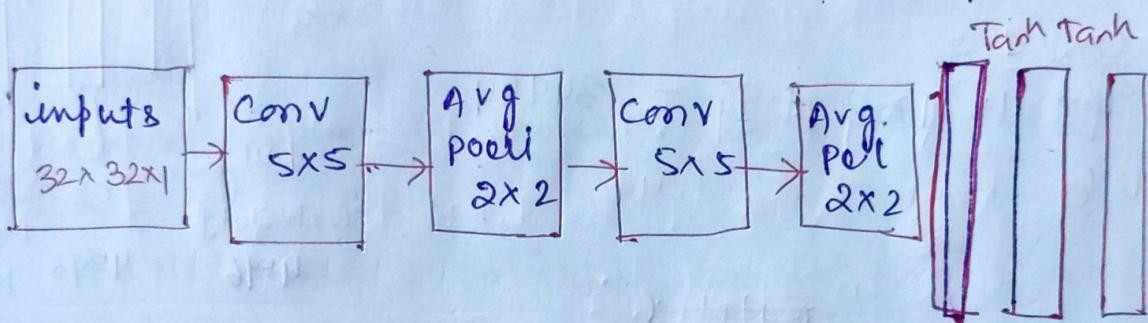
(10) Load the model with the Best classification  
Accuracy on the validation set.

(11) Calculate the Classification Accuracy on the Test set.

## CNN-4 LENET

LENET → LENET-5 is one of the simplest architectures.

- It has 2 convolutional and 3 fully-connected layers.
- It is the most widely known CNN architecture.
- It was introduced in 1998 and is widely used for handwritten method digit recognition.
- LENET-5, architecture has 60,000 parameters.



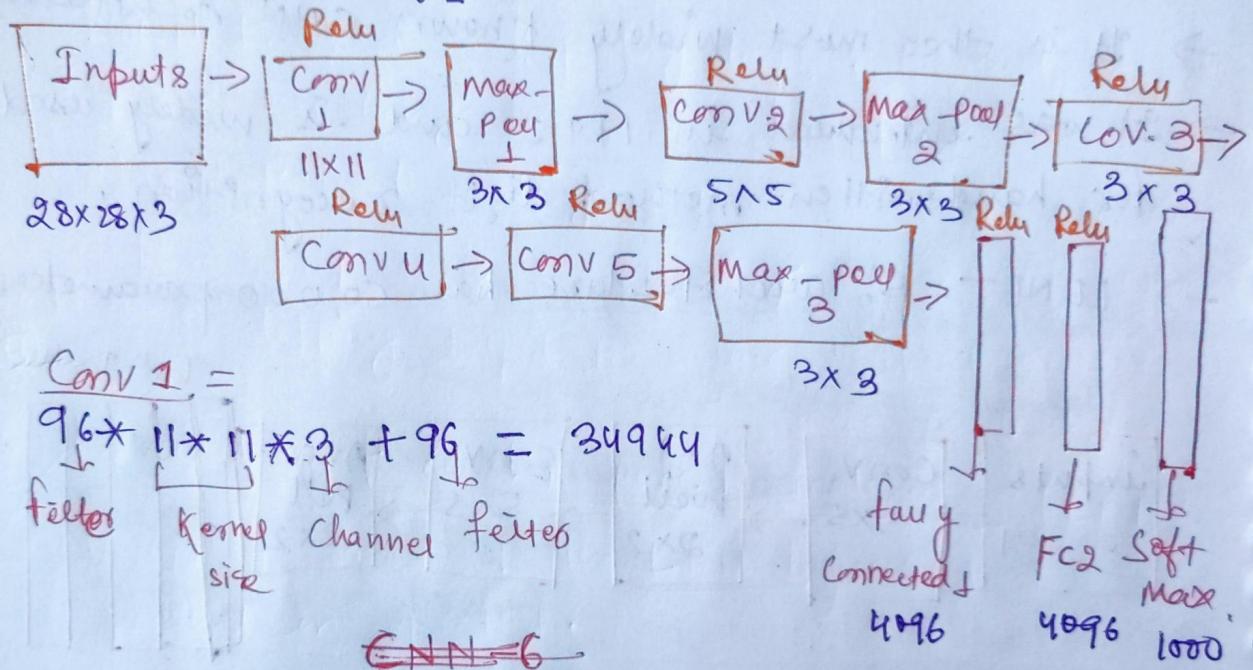
→ The LeNet-5 has the ability to process higher one-resolution images that require larger and more CNN Convolutional layers.

ALEXNET → [2012]

- It consists of 8 layers, including 5 conv. layers and 3 fully connected layers.
- It uses traditional stacked Convolutional Layers with max-pooling in between.
- It allows for the extraction of complex features from images.

Activation function → ReLU used.

- Activation function, and dropout Regularization, which enhances the model's ability to capture non-linear relationships within the data.
- Batch size is of 128



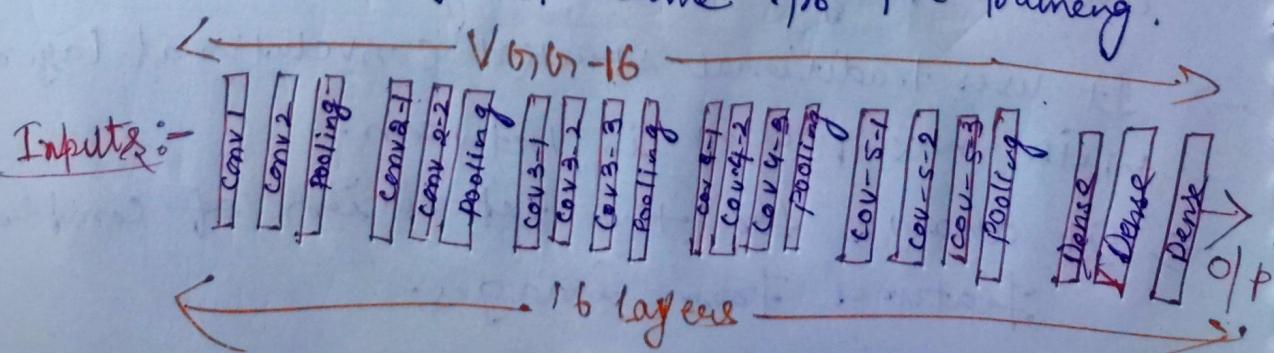
CNN = 6

\* VGG NET \* [Visual Geometry Group]

- VGG 16 [16 - no. of layers)  
→ VGG 19 [19 - no. of layers)

With trainable Parameters

In 2015. → Image net came for pre training.



## GnVGG NET

- All convolution layers were made to have kernel sizes equal to  $3 \times 3$ .
- CNN is very deep as 16 layers for large scale image recognition.
- LRN - removed.
- Max pooling →  $2 \times 2$  and stride is (2) so, no overlapping occurred.

\*

### \* Some key changes made in VGGNet from AlexNet

#### (1) Deeper Architecture ⇒

VGGNet introduced a much deeper architecture compared to AlexNet.

→ AlexNet has 8 layers, VGGNet has 16 or 19 layers. The increased depth allowed VGGNet to capture more complex features and patterns.

#### (2) Smaller filter sizes ⇒

VGGNet utilized small  $3 \times 3$  conv. filters throughout the network, whereas AlexNet used a mixture of  $3 \times 3$  and  $5 \times 5$  filters.

The use of small filters for a more localized feature representation and enabled deeper networks without a significant increase in parameters.

#### (3) Uniform structures ⇒

→ VGGNet maintained a uniform structure with a stack of conv. layers followed by max pooling layers.

- VGGNet helped to reduce spatial resolution more gradually.
- the uniform made it easier to design and understand the network architecture.
- ⇒ ④ NO Overlapping Pooling ⇒
- VGGNet used max pooling layers with a stride of 2, resulting in non-overlapping pooling regions.
- AlexNet used pooling regions with a stride of 2 but with Overlapping regions of size  $3 \times 3$ .  
Non-overlapping helped to reduce spatial resolution more gradually.
- ⑤ Elimination of Local Response Normalization (LRN)
- VGGNet removed the LRN layers that were present in AlexNet.  
LRN is initially introduced in AlexNet as a technique to enhance local contrast.
- However, subsequent research found that LRN has limited impact on performance and was computationally expensive.
- ⑥ Smaller Fully Connected Layers ⇒
- VGGNet used smaller fully connected layers compared to AlexNet.
- VGGNet used 4096 neurons in the last two fully connected layers reduced the no. of neurons.

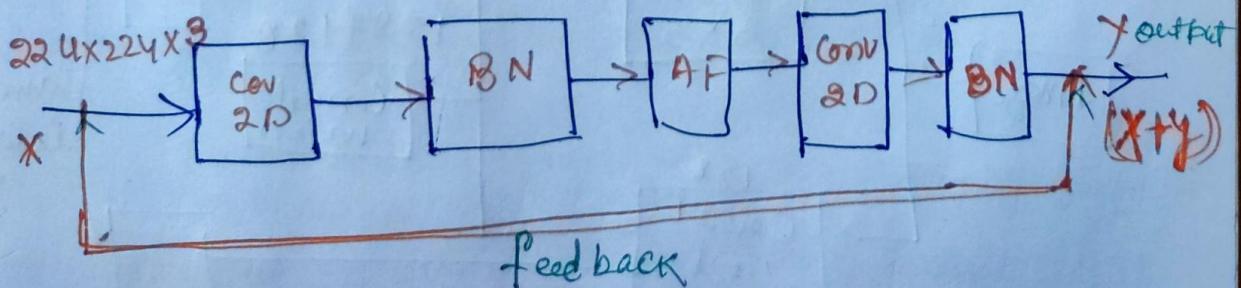
- to 1000 in the output layers.
- the ~~change~~ change helped to reduce the no. of parameters and computational complexity.
- ⇒ if padding is same then conv. NN's dimension can't reduced.

## RESNET \*

[ Residual Networks ]

- In this network, we use a technique called Skip connections.
- Skipping
- the skip connection connects activations of a layer to further layers by skipping some layers in between.

Residual Block (U.S.P) → used to solve Vanishing/exploding gradient problem.



This Residual block was the main U.S.P for the success of RESNET

# Inception

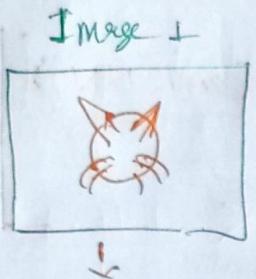
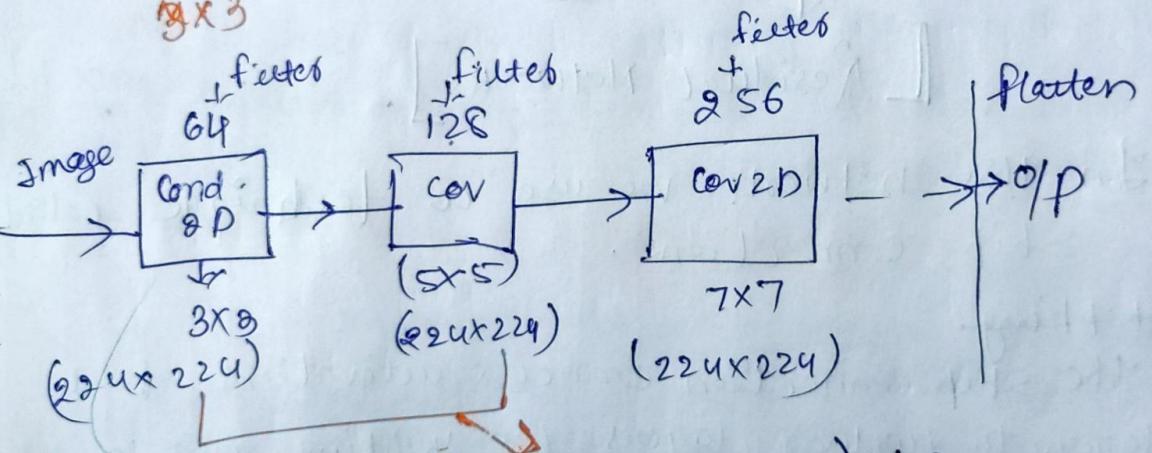


Image 2



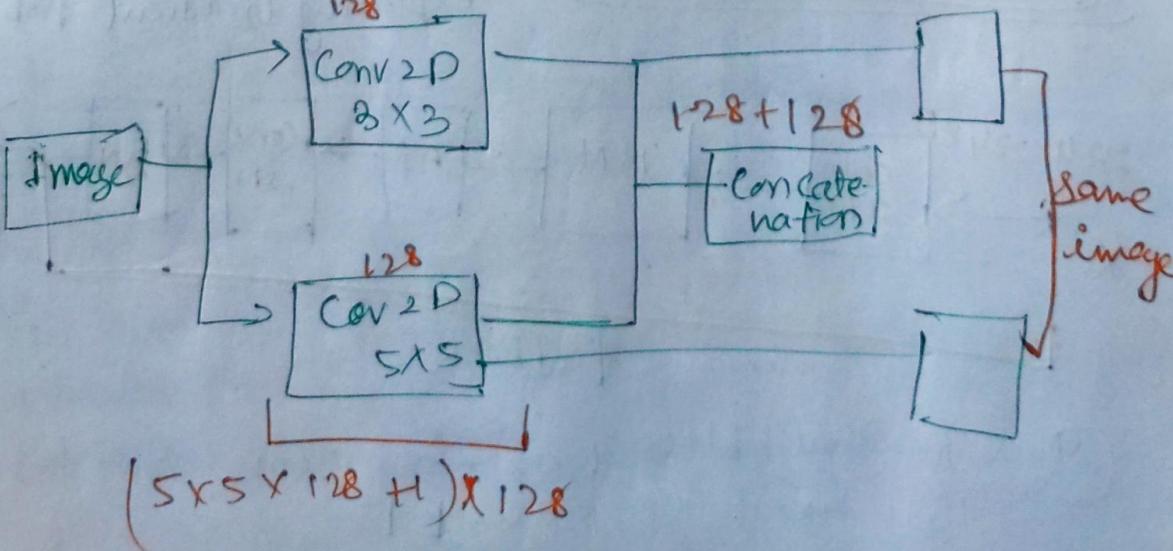
Feature map

$3 \times 3$



Parameters :  $-(3 \times 3 \times 64 + 1) \times 128$

kernel bias filters



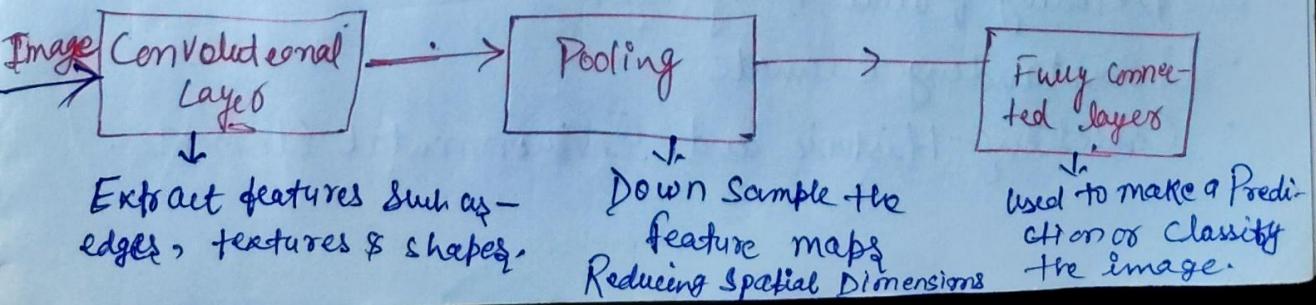
- In short, this parallel approach of applying on different dimensions of filters to our filter map in images and concatenating them to get the desired output with the same dimension.
- \* Less no. of trainable parameters.
- \* Multiple filter sizes.

### features of Inception :-

- ① It adds non-linearity in the architecture.
- ② No. of Parameters remains same.
- ③ Inception Architecture creates multiple sizes of images which helps in feature engineering.

### CNN (Convolutional Neural Network) :-

- It is well-suited for image recognition and processing tasks.
  - It consists of multiple layers :-
- i) Convolutional layer
  - ii) Pooling layer
  - iii) Fully connected layers.



→ Output of convolutional layers is passed to Pooling layers and then output of pooling layers are passed to ~~to~~ fully connected layers for prediction or classify the image.

### Types of CNN Models ⇒

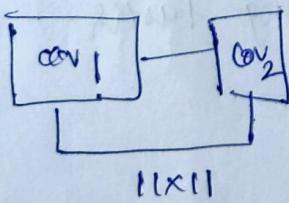
① LeNet → It was specifically designed for handwritten digit recognition, and was one of the first successful CNNs for image recognition.

→ It was trained on the MNIST dataset, which consists of 70,000 images of handwritten digits, and was able to achieve high recognition accuracy.

② AlexNet ⇒

→ It was the first CNN to win the ImageNet Large Scale Visual Recognition challenge (ILSVRC), a major image recognition competition.

→  $\text{Conv2D} = 5$ , Pooling layer = 3, fully connected layers = 3



→ first 2 conv layers uses  $11 \times 11$  size of Kernel.

### Applications ⇒

- Decoding facial Recognition
- Understanding Climate
- Collecting Historic and Environmental Elements.

## R-CNN (Region-based CNN)

[ Used for Object Detection ]

### Object Detection :-

- RCNN (Naïve Approach)
- Fast RCNN
- Faster RCNN
- YOLO
- SSD

\* Exhaustive Search → Search types of images or no. of images.

\* Selective Search → Search iteratively  
steps → It is greedy algo.

Step 1 → Cropped images

Methods to extract Cropped images:-

- ① Exhaustive Search ② Selective Search.

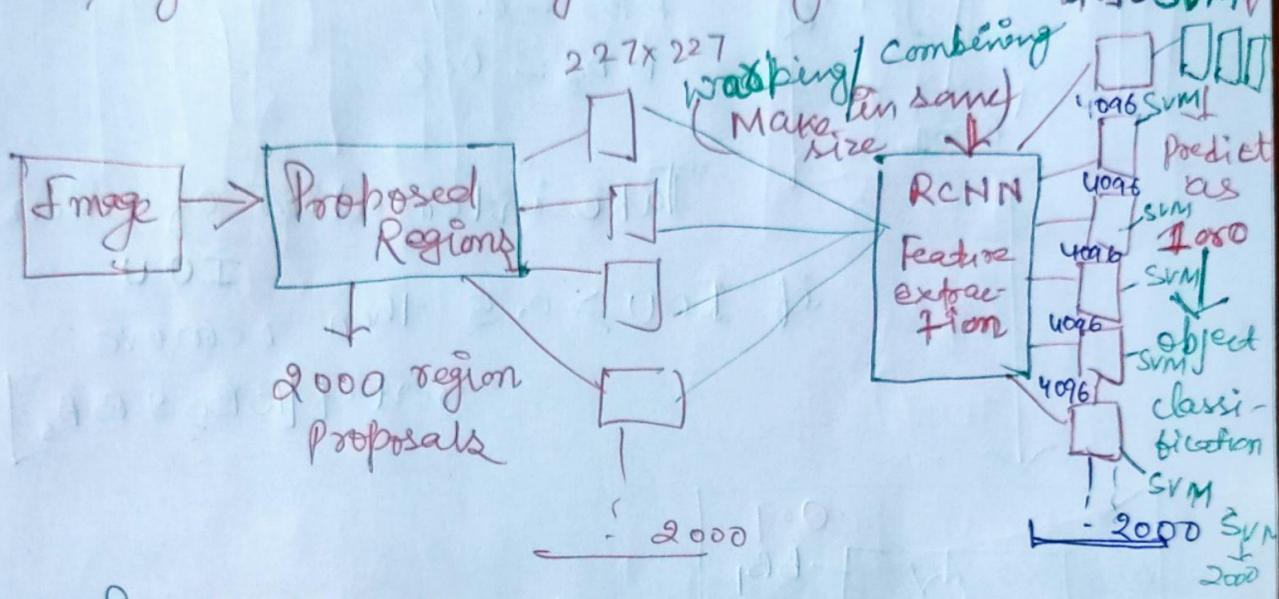
\* Regions extracted from R-CNN.

① Selective Search → Region Proposal Method  
↳ Used in RCNN is selective search around 2000 regions proposals.

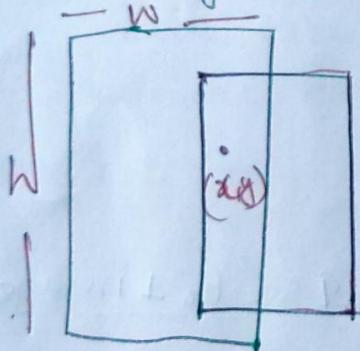
→ This algorithm works by generating sub-segmentations of the image that could belong to one object - based on color, texture, size and shape - and iteratively combining similar regions to form objects.

## RCNN Pipeline $\Rightarrow$

- ① stage 1 - Feature extraction from Region Proposals.
- stage 2 - SVM for object classification
- stage 3 - Bounding box regression



## Bounding Box Regression $\Rightarrow$

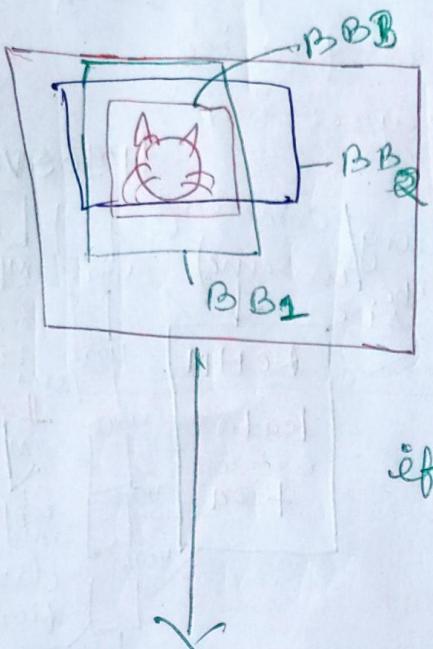


$(dx, dy, dw, dh) \rightarrow$  Regress on these points for fine tuning

## \* Steps of RCNN $\Rightarrow$

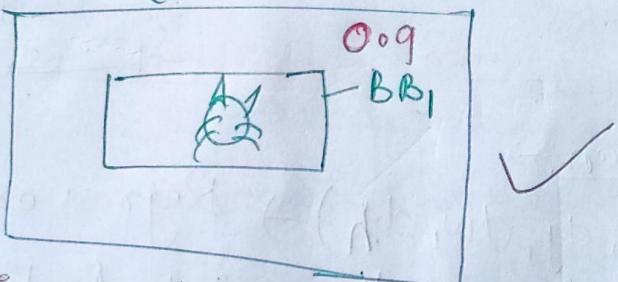
- ① Region Proposals [by exhaustive search]
- ② Warping and Resize
- ③ Pre-trained CNN [AlexNet, VGGNet used]
  - $\hookrightarrow$  for fine tuning on category
- ④ Pre-trained SVM for every image.
- ⑤ Clean up images  $\rightarrow$  [NON MAX Suppression]
- ⑥ Bound Box Regression

Non-Max Suppression  $\Rightarrow$  For clean-up or fine-tuning for best accuracy of Model.



$$\begin{aligned}BB_1 &= 0.9 \text{ [Accuracy]} \\BB_2 &= 0.7 \\BB_3 &= 0.8\end{aligned}\quad \left.\right\}$$

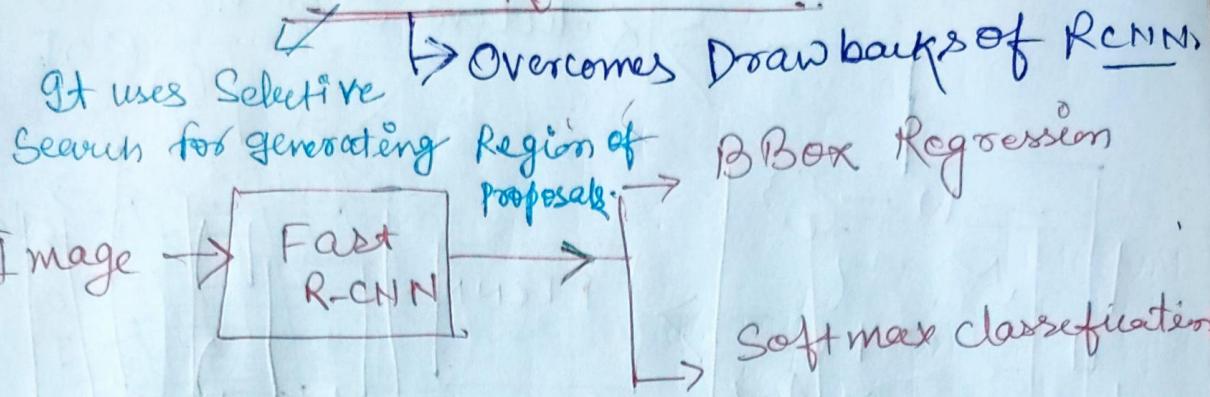
Take intersection of union  
 $\text{IOU}$   
if  $\text{IOU} > 0.5$  then remove  
rest of  $[BB_2 \cup BB_3]$



Disadvantage of RCNN  $\Rightarrow$

- ① Very computationally expensive,
- ② Extremely time consuming (49 see for 1 image)
- ③ We need to pretrain [CNN, SVM]
- ④ huge amount of time, storage & computation power used.

## Fast R-CNN

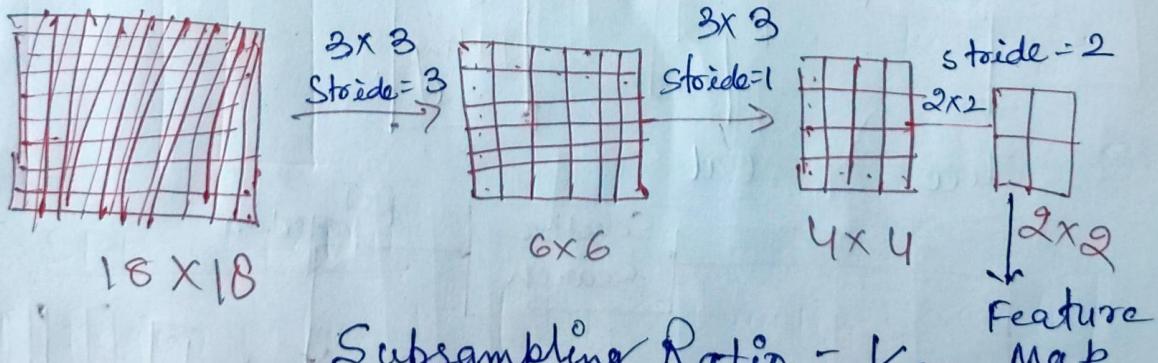


### \* Fast-CNN Features →

- ① Subsampling }
- ② ROI Projection }
- ③ ROI Pooling }

↑  
Non Max Suppression

① Subsampling ⇒ A method that reduces data size by selecting a subset of the original data.



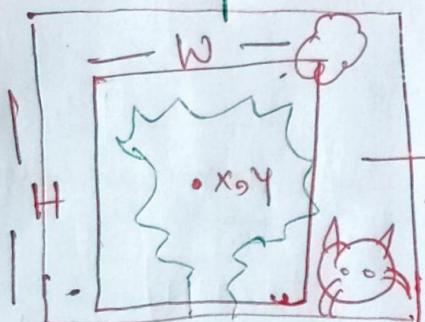
$$\text{Subsampling Ratio} = \frac{1}{3}$$

### ② ROI (Region of Interest) ⇒

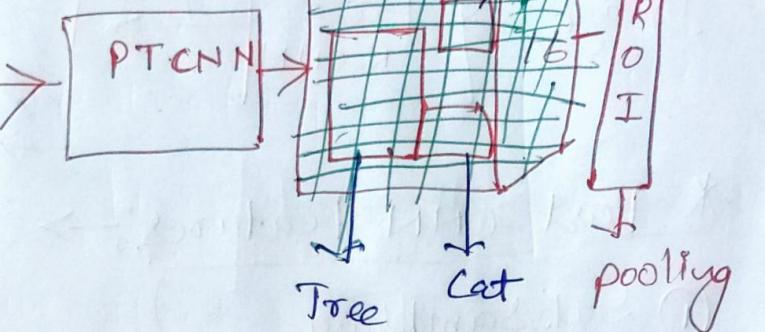
→ It extracts the feature vector of fixed length for each feature map.

\* Fast R-CNN → It runs 1 time instead of 2000 times.

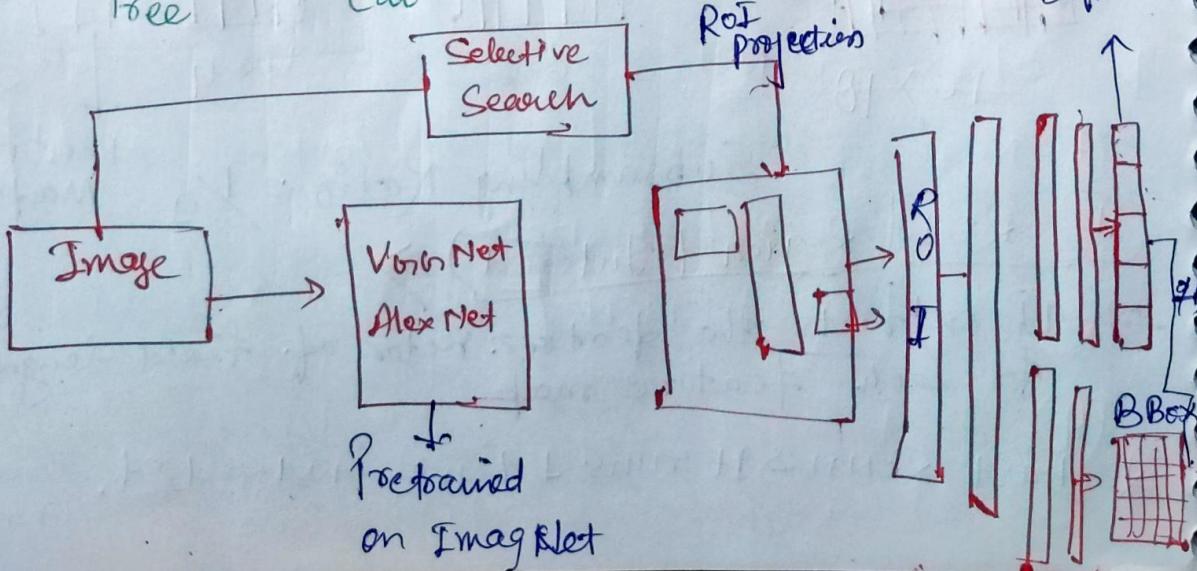
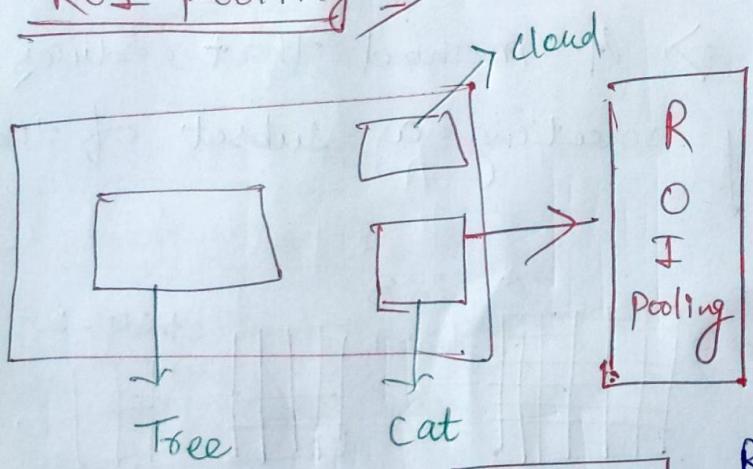
## Selective Search (2000)



- (1) Tree
- (2) cat
- (3) cloud



## ROI Pooling $\Rightarrow$

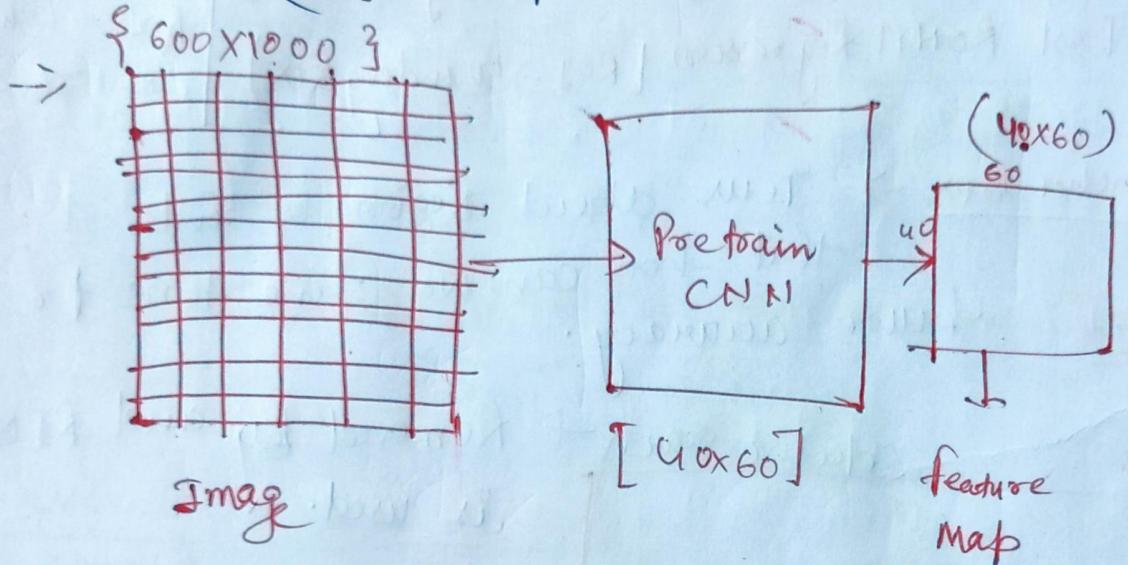


~~Faster RCNN  $\rightarrow$  RPN used to generate ROI.~~

[RPN + Fast RCNN]



(Region Proposal Network)



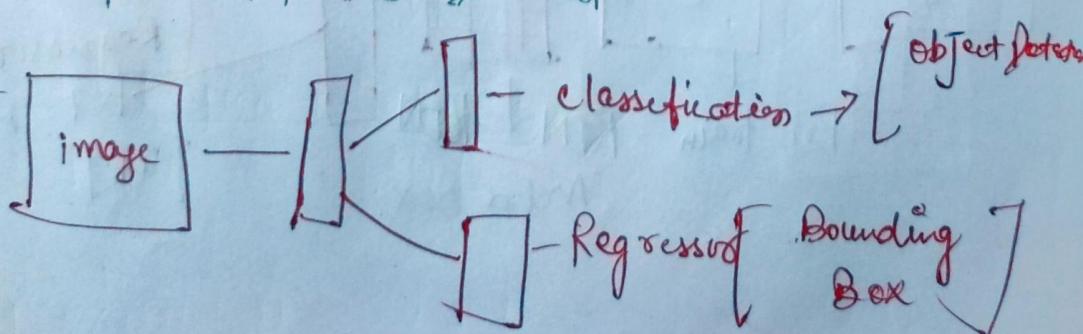
$$40 \times 60 \times 9 \geq 20000$$



no. of sliding  
windows

→ Comparing with my Anchor Box Result, with selective search it gives 10 times more ROI's.

so, this problem is solved



\* the Faster R-CNN has 4 losses: -

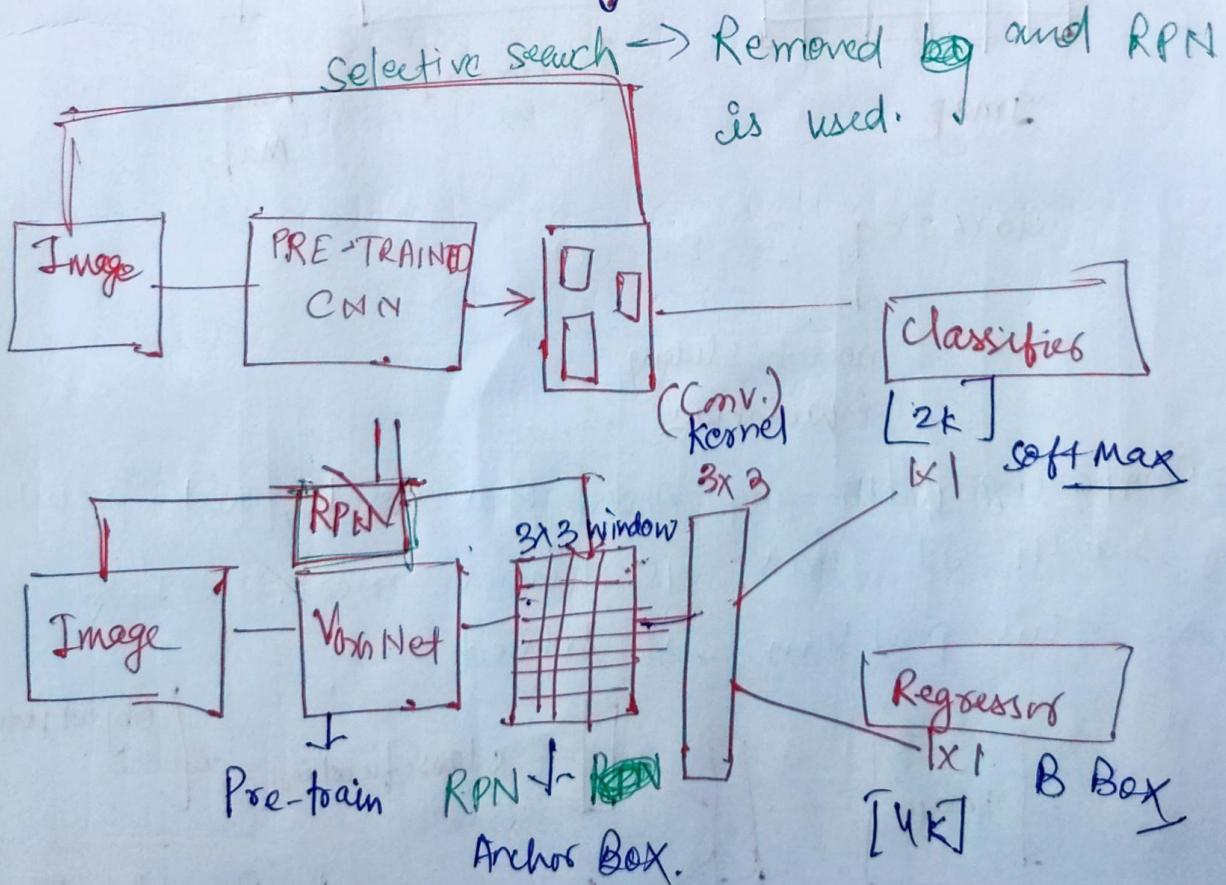
(i) RPN Classifier (Object foreground /  $\rightarrow$  Binary classifier background)

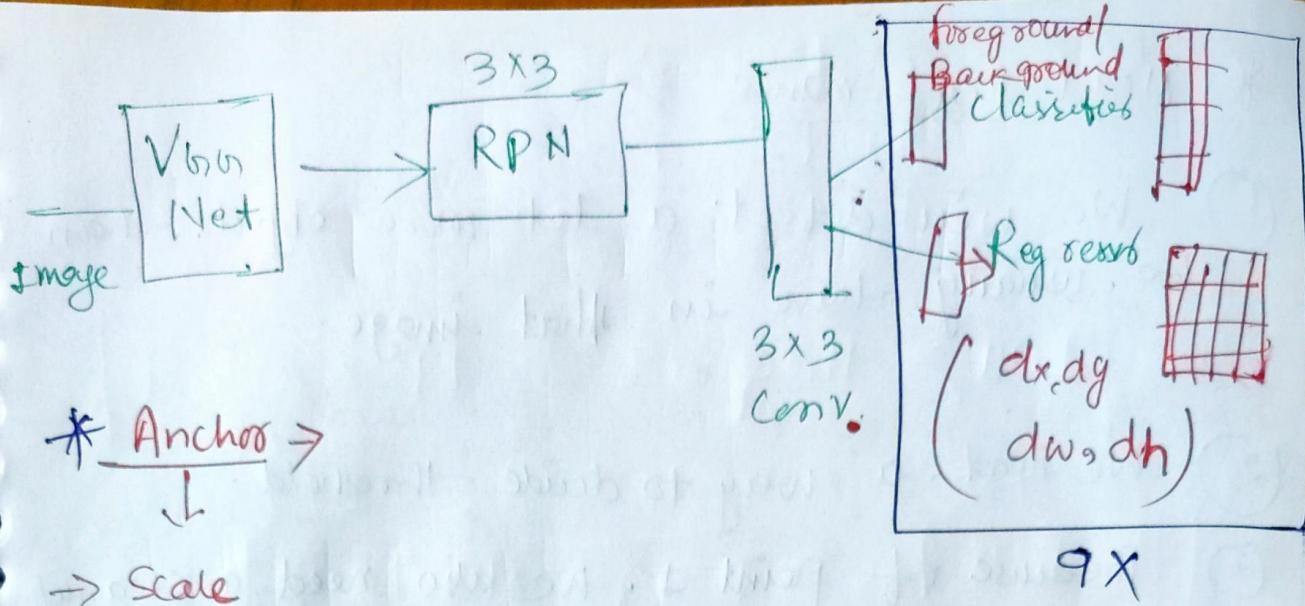
(ii) RPN Regression (Anchor  $\rightarrow$  ROI)

(iii) Fast R-CNN Classification (Object classes)

(iv) Fast R-CNN Regression [ROI  $\rightarrow$  Bounding Box] object class  
ROI = Boundary  
dx, dy, dw, dh

Anchor Box  $\rightarrow$  Tells about Region of interest that how can we fine tune for better accuracy.



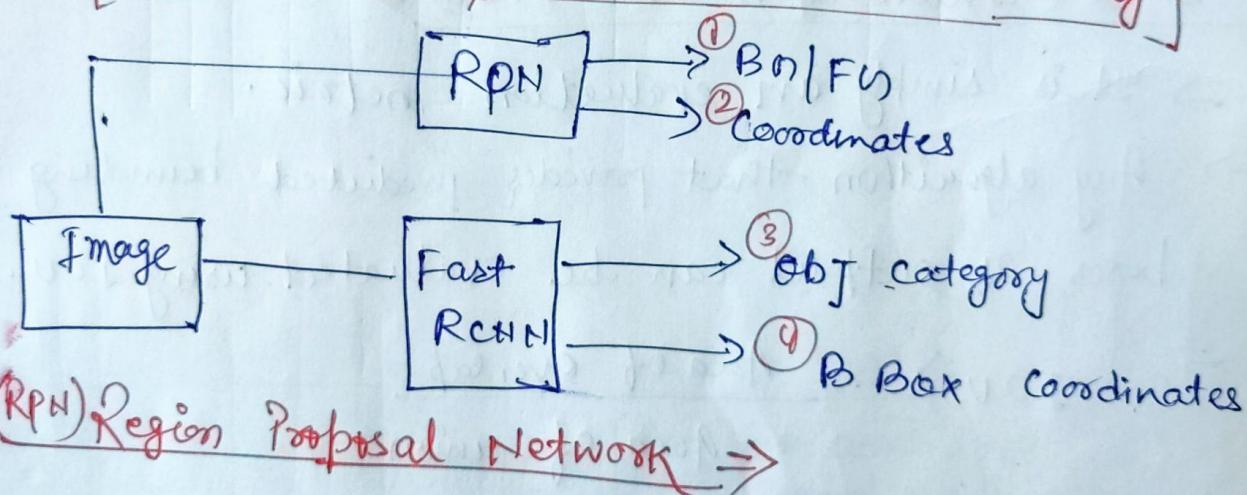


\* Anchors →  
↓

→ Scale

→ Aspect Ratio

[ Faster-Rcnn Detailed Understanding ]



## \* Model Performance $\Rightarrow$

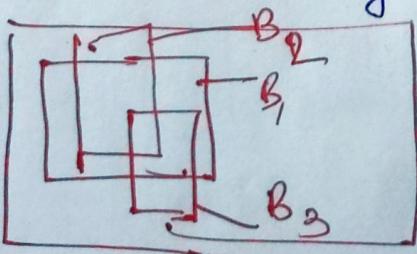
- (1) We will detect a lot more objects than we usually have in that image.
- (2) We need a way to decide threshold.
- (3) Because of Point 1, we also need a cleanup process, by this predictions will be limited.

## IoU (Intersection over Union) $\Rightarrow$

- $\rightarrow$  It is simply an evaluation metric.
- $\rightarrow$  Any algorithm that provides predicted bounding boxes as output can be evaluated using IoU.

$$\leftarrow \text{IoU} = \frac{\text{Area of overlap}}{\text{Area of union}}$$

Area of overlap = Between predicted bounding box and the ground-Truth bounding box.



$$\text{IoU} = B_1 \cap B_2 \cap B_3$$

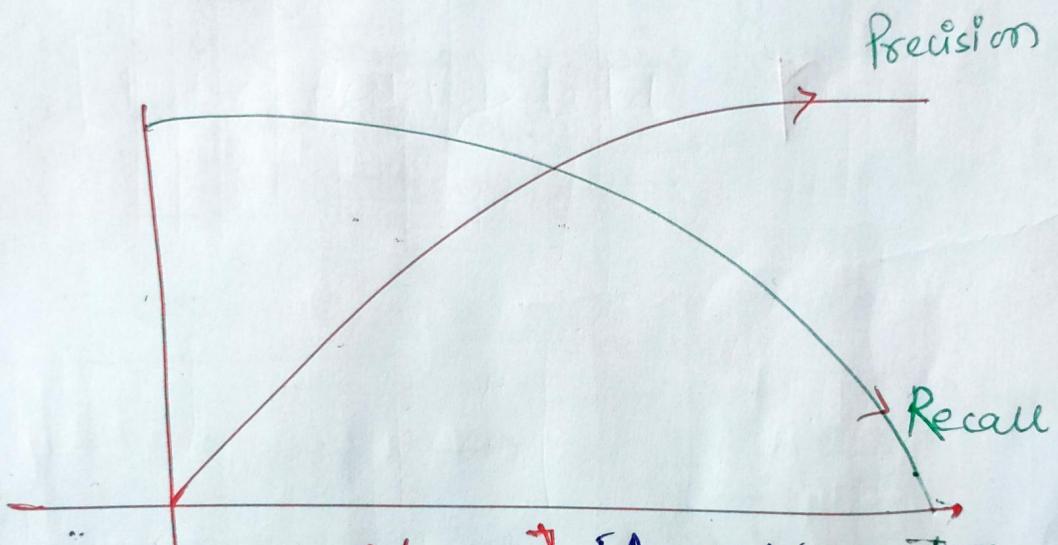
\* Mean Average Precision  $\rightarrow$  It is a metric used to evaluate object detection model.

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

It varies  $[0-1]$

$\rightarrow$  IOU, Recall, Precision, confusion matrix etc.

$$\text{MAP} = \sum_{i=1}^n \text{AP}_i$$



$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

[Area under curve IOU = Mean Avg. Precision]

Non-Max Suppression  $\Rightarrow$

- ① Remove all B Box with confidence level  $< 0.75$
- ② Find B Box with max confidence level.
- ③ All B Box have  $\text{IOU} > 0.5$   
then remove all B Box except large value B Box.

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

Actual Ground  
Truth of class

$\rightarrow$   $\boxed{\text{IoU} > \text{threshold} \rightarrow \text{TP considered}}$

$\rightarrow$   $\boxed{\text{If IoU } \uparrow \text{ True} \rightarrow \text{Precision also True}}$