# CNN

Tuesday, September 28, 2021     6:20 PM



n - f + 1  = features size,
where n - matrix
        f  - filter size

**Padding**



As the filter is smoothing through image, convolution happens. The convolution basically happens more in central part.
So if image has features on side we have to do padding around image to slide feature towards center.
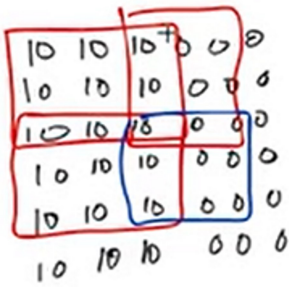Then New formula will be
 n + 2p - f + 1 where p - padding

Valid padding  - without padding
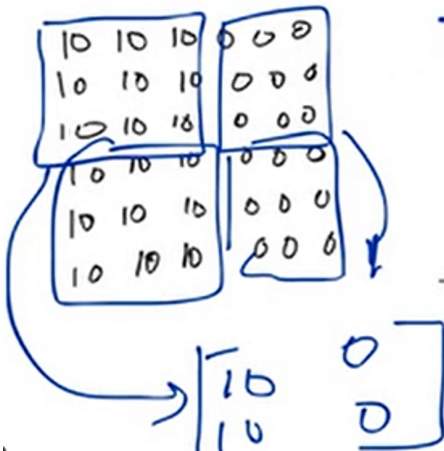Same padding - with padding

**Striding**

## Striding



If the image is huge dense or large megapixel, the convolution takes time and results in more training time. In Striding, increase step while convolution instead of take 1 step. It reduces the image size very high

Updated formula will be

$$\frac{n + 2p - f}{n} + 1, \quad \text{where n - no of striding}$$

**Pooling** :
It is completely a feature extractor. It always take the highest pixel value.



-> If our filter is too big, it will sometimes considers two images as one. So that why our filter should be smaller and we also pass additional features in network>
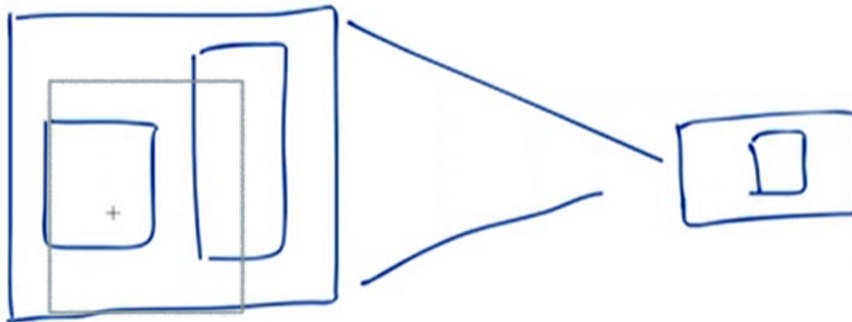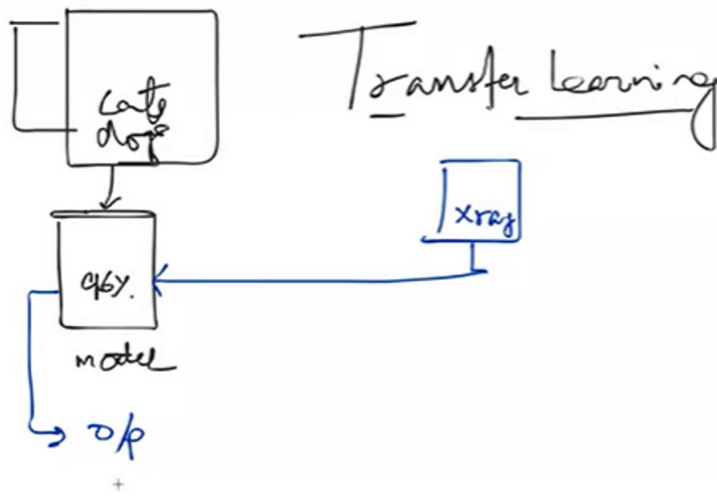


Image Preprocessing
  1. resize image
  2. Gray scaling

3. Smoothing or blurring / Threshold
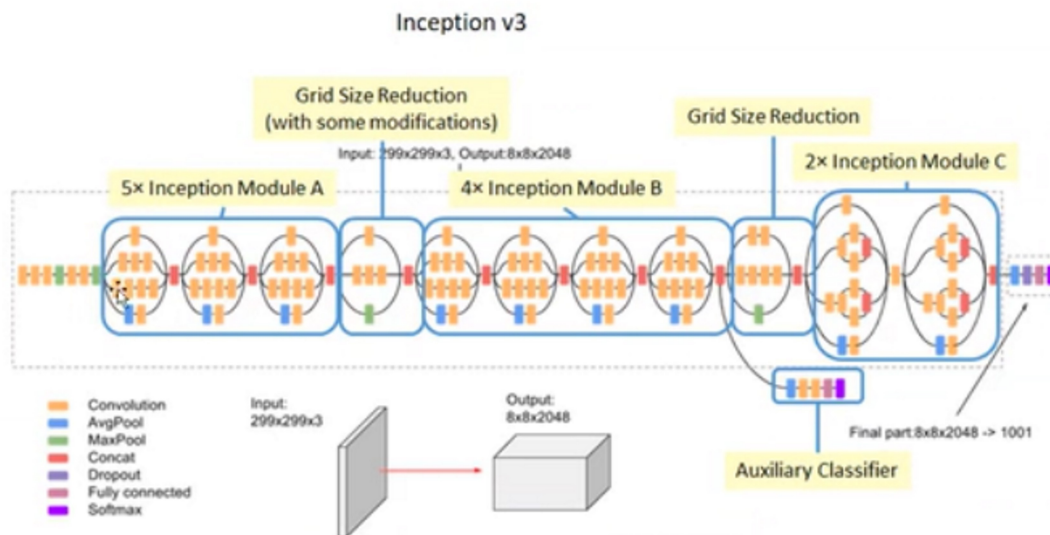4. You can dilate or erode
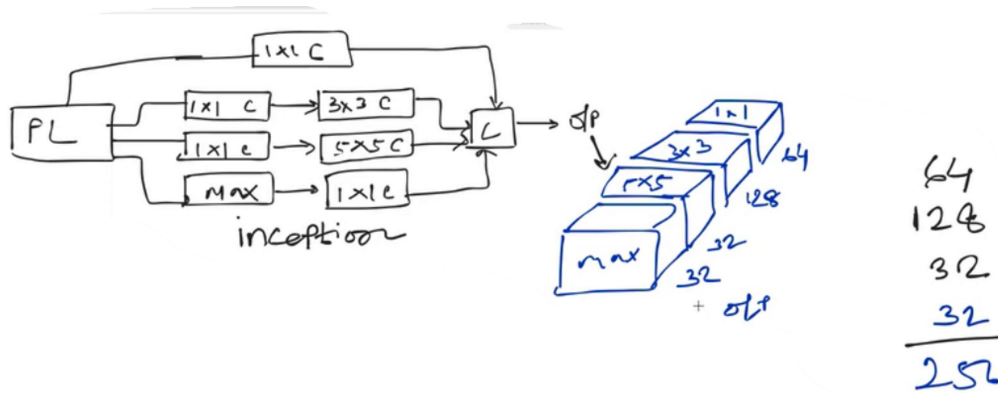
Transfer Learning Technique:



It uses pre-trained model to predict the output with different data. If required, we train the model, but the training time is less. The Hyperparameters are already tuned, we don't need to do that, it will saves time and resources for training. Some of the pre-trained Models are Lenet, Alexnet, vgg16, vgg19 etc.

Savoula Threshold : It reduces noise to 255 or white.
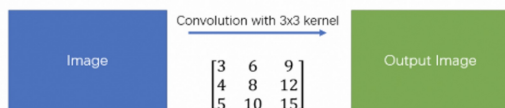
Inception Model by Google :

Instead of using normal 256 filters, it uses different range & sizes of filters in each module. Each module is connected by other.
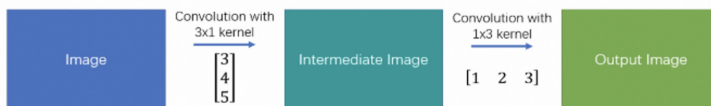It saves memory and time. It is used in devices having low resources.

Separable Convolution :
1. Spatial Separable Convention :
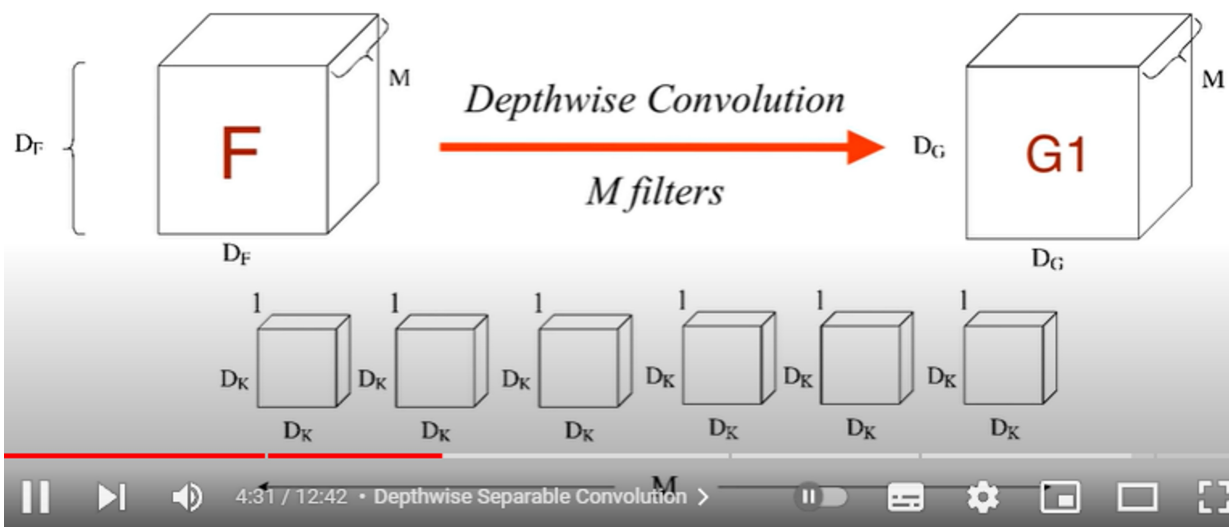   It breaks the kernel. Its not use in Inception
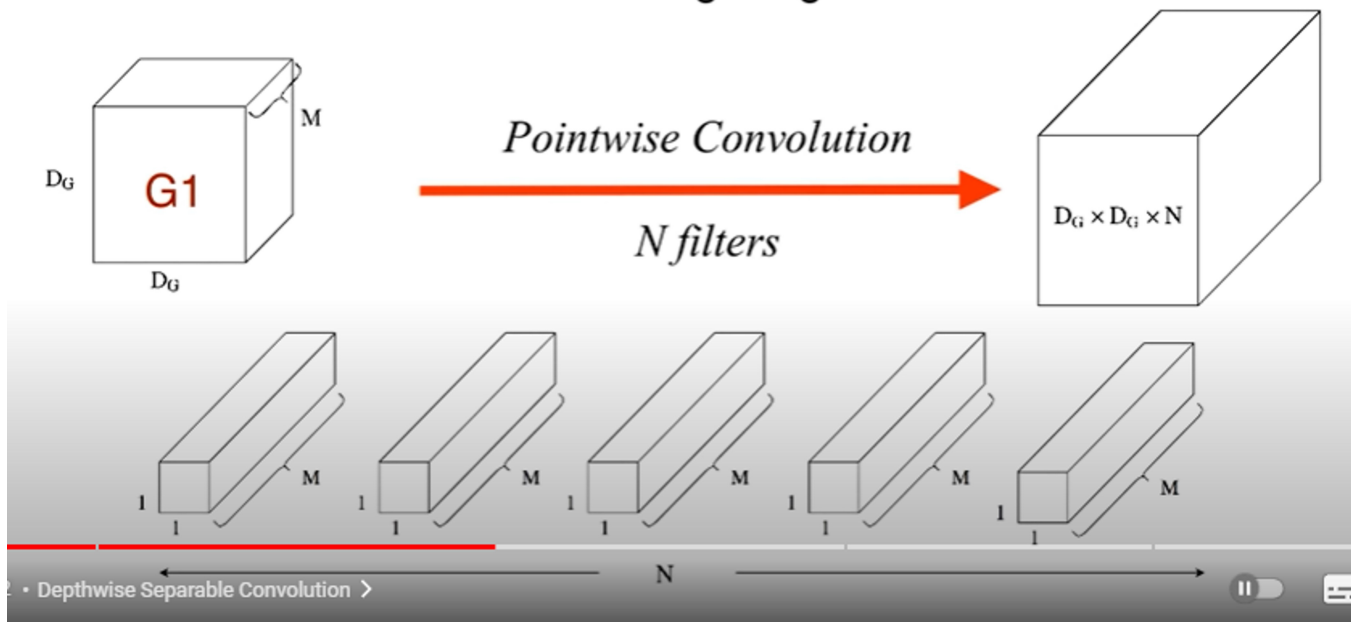


2. Depth Separable Convolution :
   It has two parts :

## 1. Depthwise Convolution: Filtering Stage
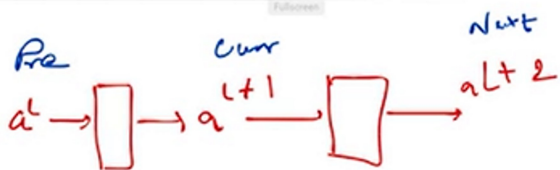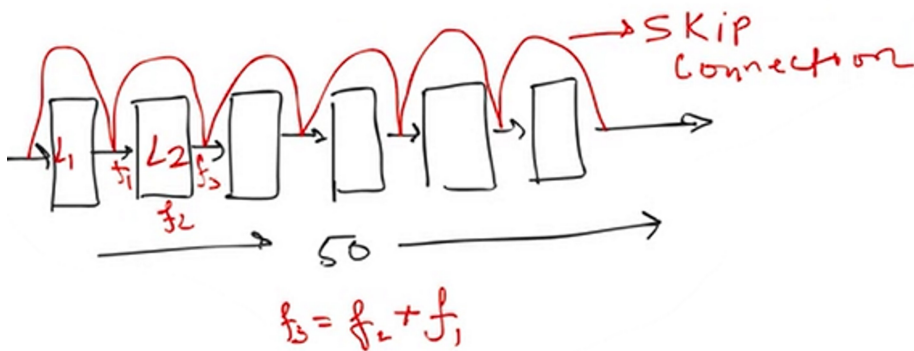
## 2. Pointwise Convolution: Filtering Stage

Resnet :
If we have many layers of CNN like 50 layers,
1. Image/ feature will vanish
2. Vanishing Gradient.

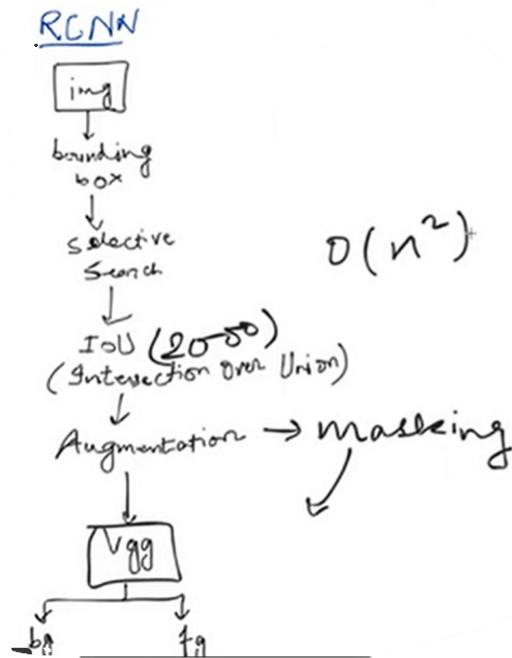For this we use Resnet, it skips connection.



$$f_3 = f_2 + f_1$$



$$a^{l+1} = ReLu(w^{l+1} \cdot a^l + b^{l+1})$$
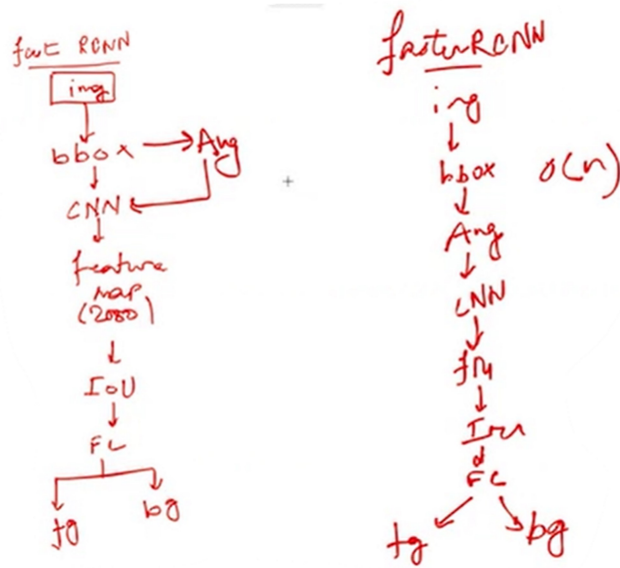$$a^{l+2} = Relu(N^{l+2} \cdot a^{l+1} + b^{l+2})$$

$$f_3 = \text{Relu}\left(W^{L+1} \cdot a^L + b^{L+1}\right) + \text{Relu}\left(W^L \cdot x + b^L\right)$$
$$\underbrace{\qquad}_{f_1} \qquad \underbrace{\qquad}_{f_2}$$

Object Detection :
In CNN, we can't segment out object or annotated object. Segmentation is detection of an object in an area. That's why we go with RCNN (Region based CNN).

RCNN



The limitation of RCNN is basically there are only 2000 objects can be detected because the time complexity is very high 0(n2).



There is fast RCNN in which after bbox and Augmentation, we pass image through CNN for feature extraction.
We use depth separable convolution in fast RCNN which takes less resources and time.
If we using annotated images, then why limit fast RCNN to 2000. We can use as many features as we want.

GANS :

GAN

Latent
space q

CNN — G

Generator

Real
features

+

CNN-D

Discriminator

loss