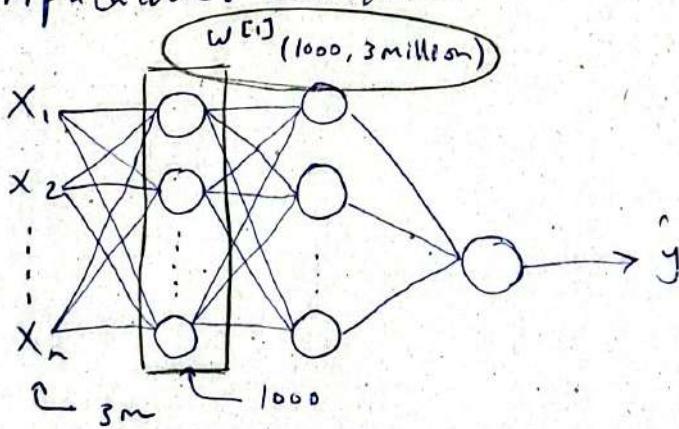


## Course 4. → CNNs

### Foundations of convolutional neural net works :-

- \* When we train a network on 64 or  $32 \times 32$  images this goes fine because if there is a 100 hidden units in the first hidden layer that makes the weights for this layer  $(100, 32 \times 32 \times 3) = (100, 3072)$  which are 30 thousand parameters which are actually not that bad compared to a high resolution image; say  $1000 \times 1000$ ; assuming that we're gonna have 1000 hidden units then the slope of w will be  $(1000, 1000 \times 1000 \times 3) = (1000, 3\text{million})$  resulting in 3 billion parameters which is very expensive in matters of computations and time.



→ To solve this problem we'll use the convolution operation

Edge detection

- \* in the first course when we were trying to understand what's going on under the hood of the neural network we said that in the first layers; the neural network is sort of detecting edges (Lines); maybe vertical or horizontal lines in the image but how does the neural network detect those edges?

## vertical edge detection

→ Let's assume that we have a  $6 \times 6$  grayscale image

$$\begin{array}{|c|c|c|c|c|c|} \hline 3 & 0 & 1 & 2 & 7 & 4 \\ \hline 1 & 5 & 8 & 9 & 3 & 1 \\ \hline 2 & 7 & 2 & 5 & 1 & 3 \\ \hline 0 & 1 & 3 & 1 & 7 & 8 \\ \hline 4 & 2 & 1 & 6 & 2 & 8 \\ \hline 2 & 4 & 5 & 2 & 3 & 9 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline -5 & -4 & 0 & 8 \\ \hline -10 & -2 & 2 & 3 \\ \hline 0 & -2 & -7 & -7 \\ \hline -3 & -2 & -3 & -16 \\ \hline \end{array}$$

→ We've gonna take this  $6 \times 6$  image and convolve it with a  $3 \times 3$  filter (kernel), the result is a  $4 \times 4$  matrix; here's how we can compute it:-

\* We'll take a  $3 \times 3$  window from our image and element wise multiply it with our filter as follows:-

$$\begin{array}{|c|c|c|} \hline 3 & 0 & 1 \\ \hline 1 & 5 & 8 \\ \hline 2 & 7 & 2 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = 3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$$

and so  $\rightarrow -5$  is our first element in the  $4 \times 4$  matrix.

→ we shift the  $3 \times 3$  window one step to the right and continue.

→ after we reach the end to the right we start from the beginning & we shift one step down and start from the left again.

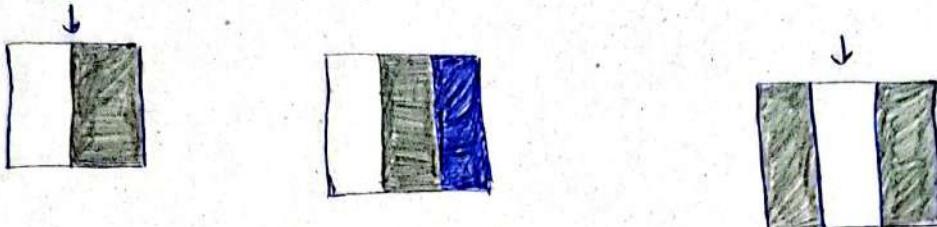
→ we repeat the process until we reach our final window

$$\begin{array}{|c|c|c|} \hline 1 & 7 & 8 \\ \hline 6 & 2 & 8 \\ \hline 2 & 3 & 9 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = 1 + 6 + 2 + 0 + 0 + 0 - 8 - 8 - 9 = -16$$

\* So a  $6 \times 6$  matrix convolved with a  $3 \times 3$  filter gives us a  $4 \times 4$  matrix.

\* So how is this convolution detecting images?

↳ if we take another example a  $6 \times 6$  image ~~that has only 2 colors~~ and convolve it by a  $3 \times 3$  filter we can see why

$$\begin{array}{|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline 1 & 1 & 1 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 3 & 3 & 0 \\ \hline 0 & 3 & 3 & 0 \\ \hline 0 & 3 & 3 & 0 \\ \hline 0 & 3 & 3 & 0 \\ \hline \end{array}$$


So as drawn in these images; the convolution process done by convolving the  $3 \times 3$  filter with the image that had two colored regions and one edge between them now has the edge highlighted so it's basically like image processing.

\* now; to why the region in the middle in the  $4 \times 4$  image is thick? it's because the image that we used which is  $6 \times 6$  is very small but when using larger images such as  $1000 \times 1000$  images; the filter does a good job on them.

\* one intuition to take away from edge detection:-

↳ a vertical edge is the output of convolving images with filters ~~for brighter regions~~.

## Filter

↳ a lot of ideas in the rise of deep learning have been introduced for the values used in the filter; if you're using the filter for vertical edge detection then our previous filter might be suitable but if we're trying to detect something else such as horizontal edges then using another filter might be a better idea

This is a filter that focuses on horizontal edges rather than vertical

1	1	1
0	0	0
-1	-1	-1

↳ other ideas were introduced regarding the values in the filter itself such as:

1	0	-1
2	0.	-2
1	0	-1

Sobel Filter

3	0	-3
10	0	-10
3	0	-3

Scharr Filter

\* We can also consider the values in the filter as a learnable parameter so instead of hardcoding the values in the filter we can consider them to be learnable weights

→ Instead of using values that were predetermined by computer vision researchers; maybe we can set these values to learnable weights

w <sub>1</sub>	w <sub>2</sub>	w <sub>3</sub>
w <sub>4</sub>	w <sub>5</sub>	w <sub>6</sub>
w <sub>7</sub>	w <sub>8</sub>	w <sub>9</sub>

→ Back prop allows the filter to detect more than just vertical or horizontal edges (i.e. learnable parameters (e.g.))

## Padding

→ When we're convolving images with the filters we lose some information from the border pixels since the filter sees them less, especially the corner pixels as the filter only sees them once.

\* One way to deal with this loss of information is to use padding; padding is the process of adding some pixels all around the pixels of the borders of the image.

0	0	0	0
0	3	-2	
0	1	0	
0	2	5	
0			

→ thanks to the padded area now the  $3 \times 3$  filter processes the pixel with the value of 3 four times instead of just once.

\* Remember that the goal of convolutions is to detect images in the picture not shrink the information from the images so it doesn't matter if the output of the convolution is the same size of the image.

## Types of padding

→ the two main types of padding are valid padding and same padding; in the valid convolution we don't pad valid:  $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$   
 $6 \times 6 * 3 \times 3 \rightarrow 1 \times 4$

same: pad so that the output is the same size as the input

$$\text{same: (for } P=1) \rightarrow n+2P-f+1 \times n+2P-f+1, n+2P-f+1=n$$

$P$  is the padding size;  $f$  is almost always odd.

$$P = \frac{f-1}{2}$$

## Strided Convolutions

\*      =

3	9	4
1	0	2
-1	0	3

91	100	83
69	51	127
49	72	74

Stride = 2

→ Striding convolutions is the process of skipping rows and columns; specifically  $s$  number of rows and columns when we're convolving the input with the filter such as:

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 7 \\ \hline 6 & 6 & 9 \\ \hline 3 & 4 & 8 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 3 & 4 & 1 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} = 91$$

now skip 2 columns ↗

$$\begin{array}{|c|c|c|} \hline 7 & 4 & 6 \\ \hline 9 & 8 & 7 \\ \hline 8 & 3 & 8 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 3 & 4 & 1 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} = 100$$

and so on until we hit the border  
↳ now we stride 2 rows down:-

$$\begin{array}{|c|c|c|} \hline 3 & 4 & 8 \\ \hline 7 & 8 & 3 \\ \hline 9 & 2 & 1 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 3 & 4 & 1 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} = 69$$

and so on until we're done.

→ in the previous example we convolved a  $7 \times 7$  image with a  $3 \times 3$  filter with stride = 2 and got a  $3 \times 3$  output.

\* The formula for this is:-

$$n \times n * f \times f \rightarrow \frac{n+2p-f}{s} + 1 \times \frac{n+2p-f}{s} + 1$$

↳ padding  $\rightarrow p$

↳ stride  $\rightarrow s$

↳ in our example we used  $p=0$  and  $s=2$

so the output shape should be  $\frac{7+0-3}{2} + 1 = \frac{4}{2} + 1 = 3$

which is the correct shape that we got.

\* if  $\frac{n+2p-f}{s} + 1$  is not an integer we can round it down to the nearest integer.

\* if input image is  $8 \times 8$  and we use same padding having a filter  $3 \times 3$  and stride = 2 we get:-

$$p = \frac{f-1}{2} = \frac{3-1}{2} = 1, \quad s = 2$$

$$\frac{n+2p-f}{s} + 1 = \frac{8+2-3}{2} + 1 = \frac{7}{2} + 1 = [4.5] = \boxed{4}$$

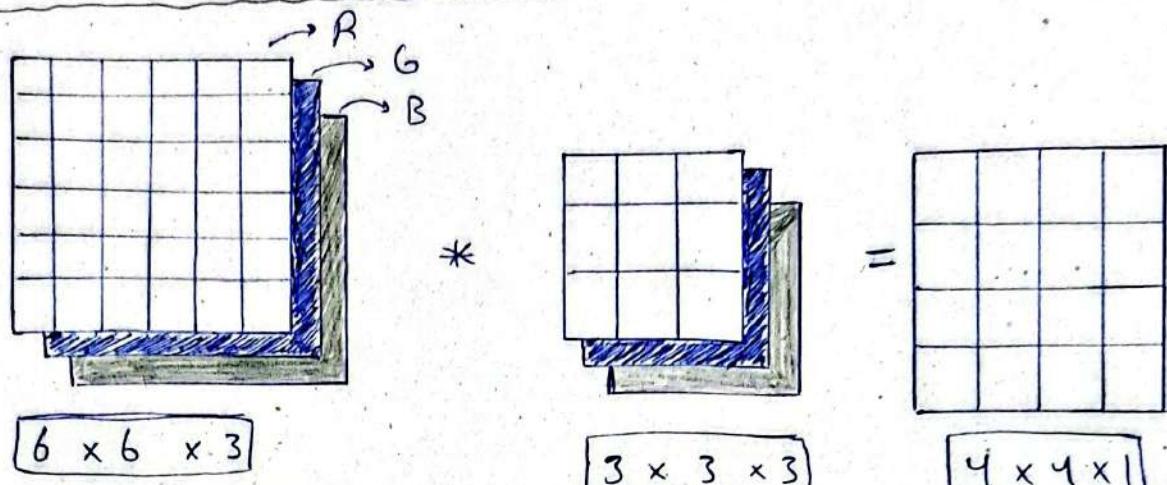
→ the way rounding down is implemented in the stride is that if the final stride doesn't apply fully to the filter (isn't the same size); you remove the whole computation completely

7	6	1	X
0	2	5	X
-3	4	2	X

3	4	4
1	0	2
-1	0	3

→ removed computation because the shape of the 2 matrices isn't the same.

## Convolutions on RGB images:-



↳ 3-channel RGB image

$6 \times 6 \times 3$

↳ 3x3x3 Filter

$3 \times 3 \times 3$

\* The number of channels in the image needs to match the number of channels in the filter.

→ the way that we calculate the convolution of the image with the filter is that we apply the  $3 \times 3 \times 3$  to the first  $3 \times 3 \times 3$  cube in the image this results in the multiplication of 27 numbers, we add the results and this is our first output.

↳ we continue to do so on as we did before even with padding and strides.

\* having a  $3 \times 3 \times 3$  filter allows us to detect edges in the way that we want so if we want to detect edges in the red. channels we can set the filter to be:

R	1	0	-1
	1	0	-1
	1	0	-1

G	0	0	0
	0	0	0
	0	0	0

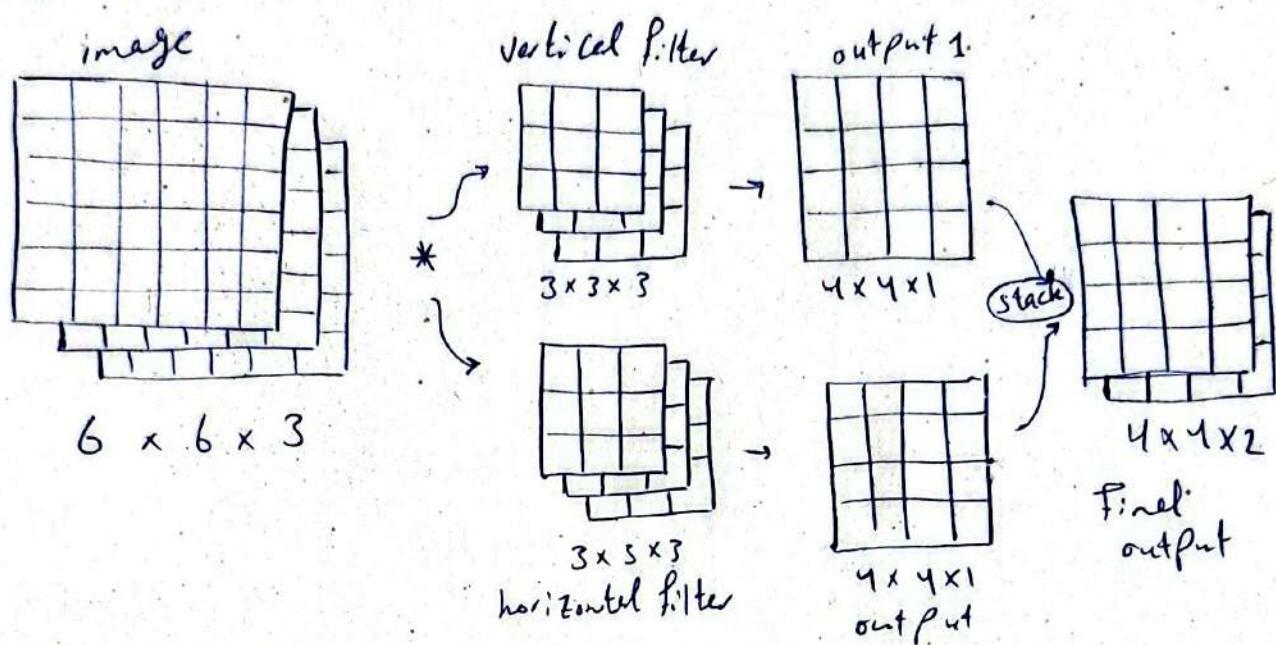
0	0	0
0	0	0
0	0	0

→ just like we wanted now the filter detects edges in the red channel only.

↳ those are the three channels

## \*using multiple filters:-

↳ now we've seen how to use filters on 3-channel images so what if we want to use multiple filters; for example a filter for vertical edges and another filter for horizontal images?



→ The process is just as seen in the drawing, first we convolve the image with the first filter and then convolve it with the second filter. and in the end we stack the two  $4 \times 4 \times 1$  outputs to be  $4 \times 4 \times 2$ .

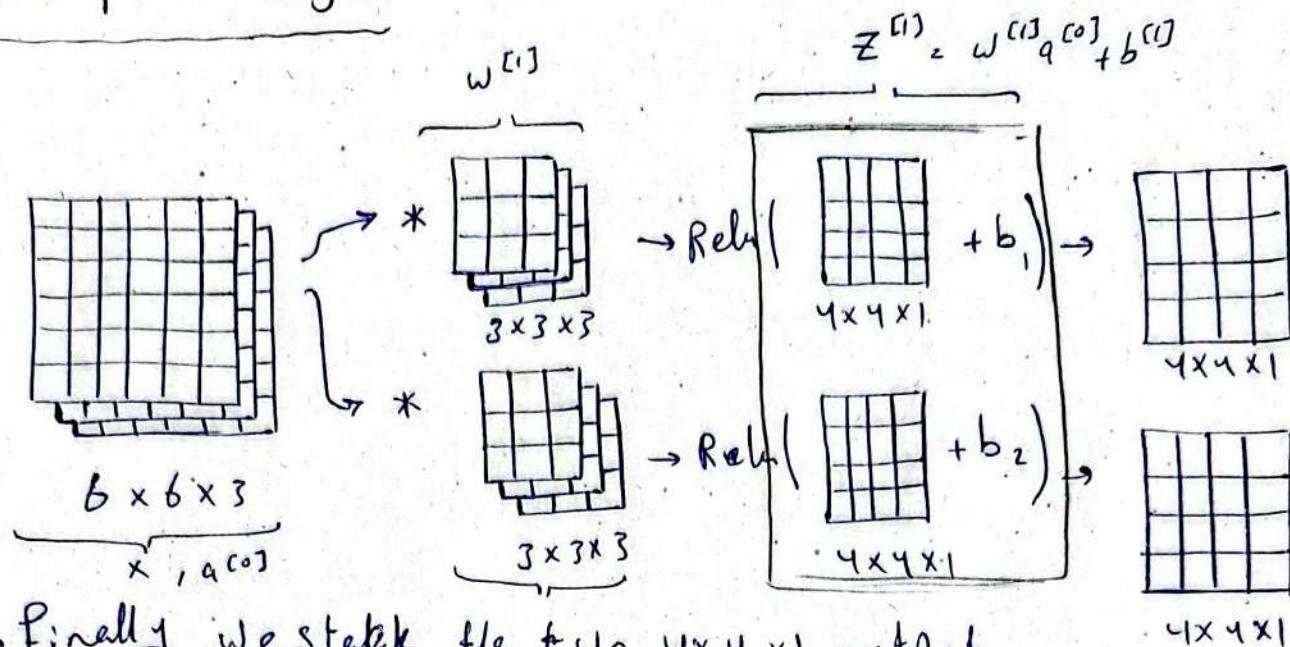
summary:-

$$n \times n \times \boxed{nc} * f \times f \times \boxed{nc} \rightarrow n-f+1 * n-f+1 \times \underline{\boxed{nc}}$$

↳ assuming  $s=1, p=0$

\* in the past example we used two filters to detect 2 features, we can add as much filters as we want to detect as features as we want.

Example of a layer:-



↳ Finally we stack the two  $4 \times 4 \times 1$  output to make the  $4 \times 4 \times 2$  output of the whole layer.

\* So just like neural networks; we take some input multiply it with weights (filters) and add some bias and pass it through some non-linear function such as ReLU and after stacking the two filter outputs now we have what we used to call  $a^{[1]}$

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

↳ probably ReLU.

\* In the previous example the output is  $4 \times 4 \times 2$  because we only had 2 filters but if we had 10 filters then the output would have been  $4 \times 4 \times 10$  so the channels of the output depend on the number of filters we use.

\* no matter how big the input image is; the number of parameters in a CNN layer depends only on the size of the filter and the number of filters.

\* If Layer L is a Convolutional layer:-

$f^{[l]}$  = filter size

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_c^{[l]}$  = number of filters

Input:  $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output:  $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$* n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$$* n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

Each filter :-

$$f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$$

Activations:-

$$a^{[l]} \rightarrow n_H^{[l]}, n_W^{[l]}, n_C^{[l]}$$

$$A^{[l]} = m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

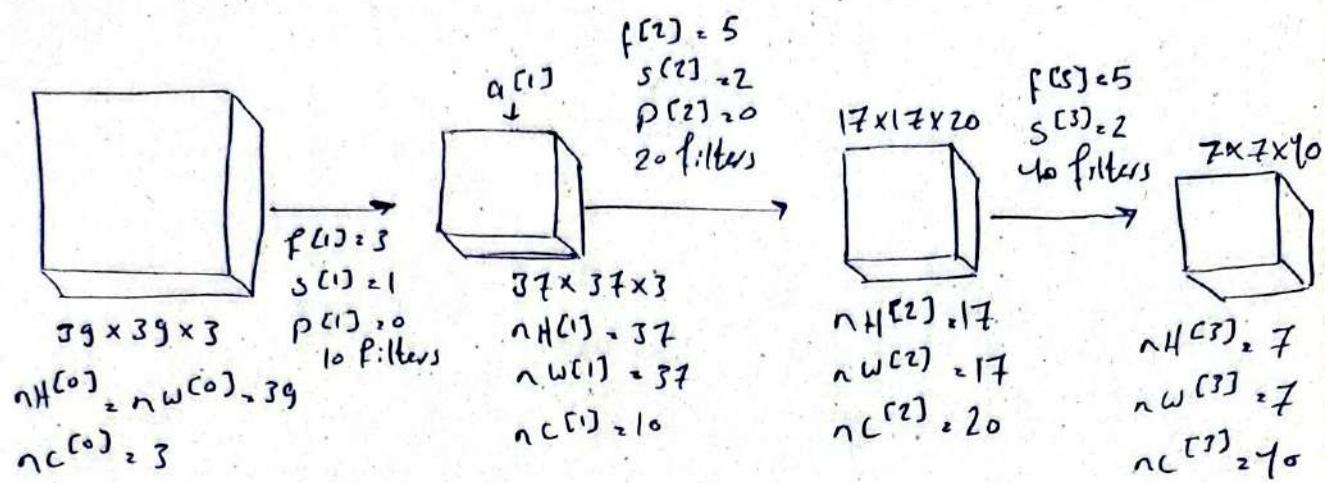
Weights:-

$$w \rightarrow f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times \underline{\underline{n_C^{[l]}}} \quad \text{number of filters in layer } l$$

$$b \rightarrow n_C^{[l]}$$

$\underline{\underline{\text{number of filters in layer } l}}$ .

## Example Conv.Nets:



→ now what's commonly done in the next step is to flatten the final  $7 \times 7 \times 10$  output that we have to 1960 elements and after that we can have fully connected layer or dropout or whatever we want to do before feeding the output to the final (sigmoid or softmax) function for classification.

## Pooling Layers

→ In this Max pooling we started with a filter size 2 and stride=2 but instead of convolving the matrix with the filter like we did with convolutional layers here we just take the max of the filtered area (the window); so the hyperparameters of the pooling layer are ,  $f$  (filter size) and stride.

→ no parameters to learn

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

$\downarrow$   
 $4 \times 4$   
 $f=2$   
 $s=2$

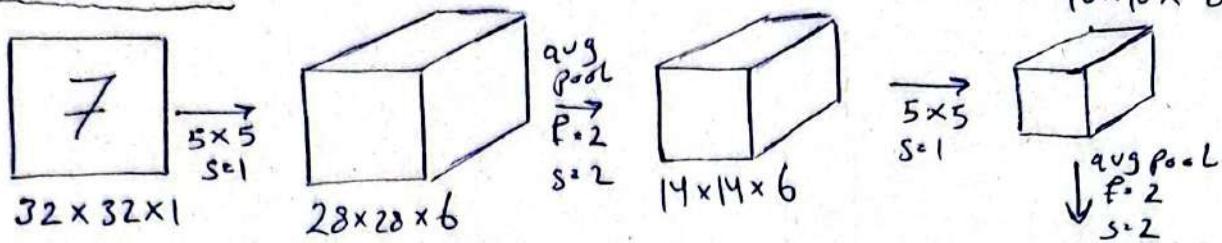
9	2
6	3

$2 \times 2$

$$\left\lfloor \frac{nH-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{nW-f}{s} + 1 \right\rfloor \times nC$$

## Week 2:- Deep Convolutional models

### LeNet-5



→ The LeNet-5 was a neural network  $5 \times 5 \times 16$  used to classify hand written digits, it's made of 2 convolutional layers followed by max pooling layers; back then, it was more common to use Averagepooling instead of Maxpooling.

→ The LeNet-5 is considered a small neural net by since it has  $\approx 60k$  parameters only.

$$\text{Conv-1} \rightarrow 6 \text{ filters } (5 \times 5) = (25 + \text{bias}) \times 6 = 156 \text{ params}$$

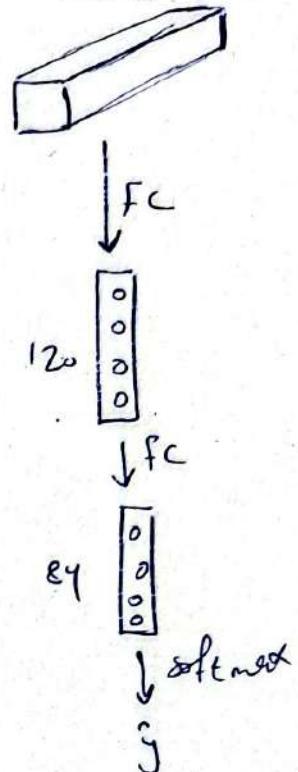
$$\text{Conv-2} \rightarrow 16 \text{ filters } (5 \times 5) = (25 + \text{bias}) \times 16 = 416 \times 2416$$

$$\text{FC-1} \rightarrow 400 \text{ inputs, } 120 \text{ hidden units} \rightarrow 48,000 + 120 = 48120$$

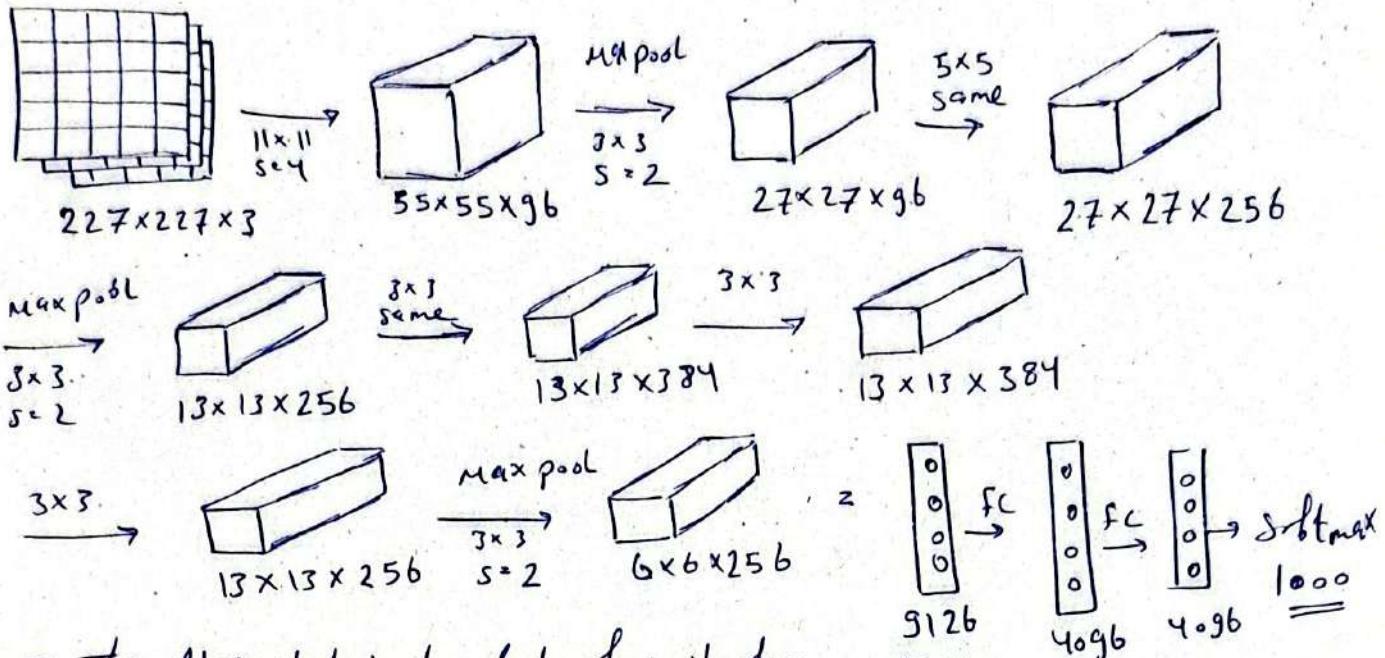
$$\text{FC-2} \rightarrow 120 \text{ inputs, } 84 \text{ hidden units} \rightarrow 10080 + 84 = 10164$$

$$\text{softmax} \rightarrow 84 \text{ inputs, } 10 \text{ units} \rightarrow 840 + 10 = 850$$

$$\rightarrow \text{total parameters for LeNet-5} = 51706 \approx 60k$$



## ALeX Net

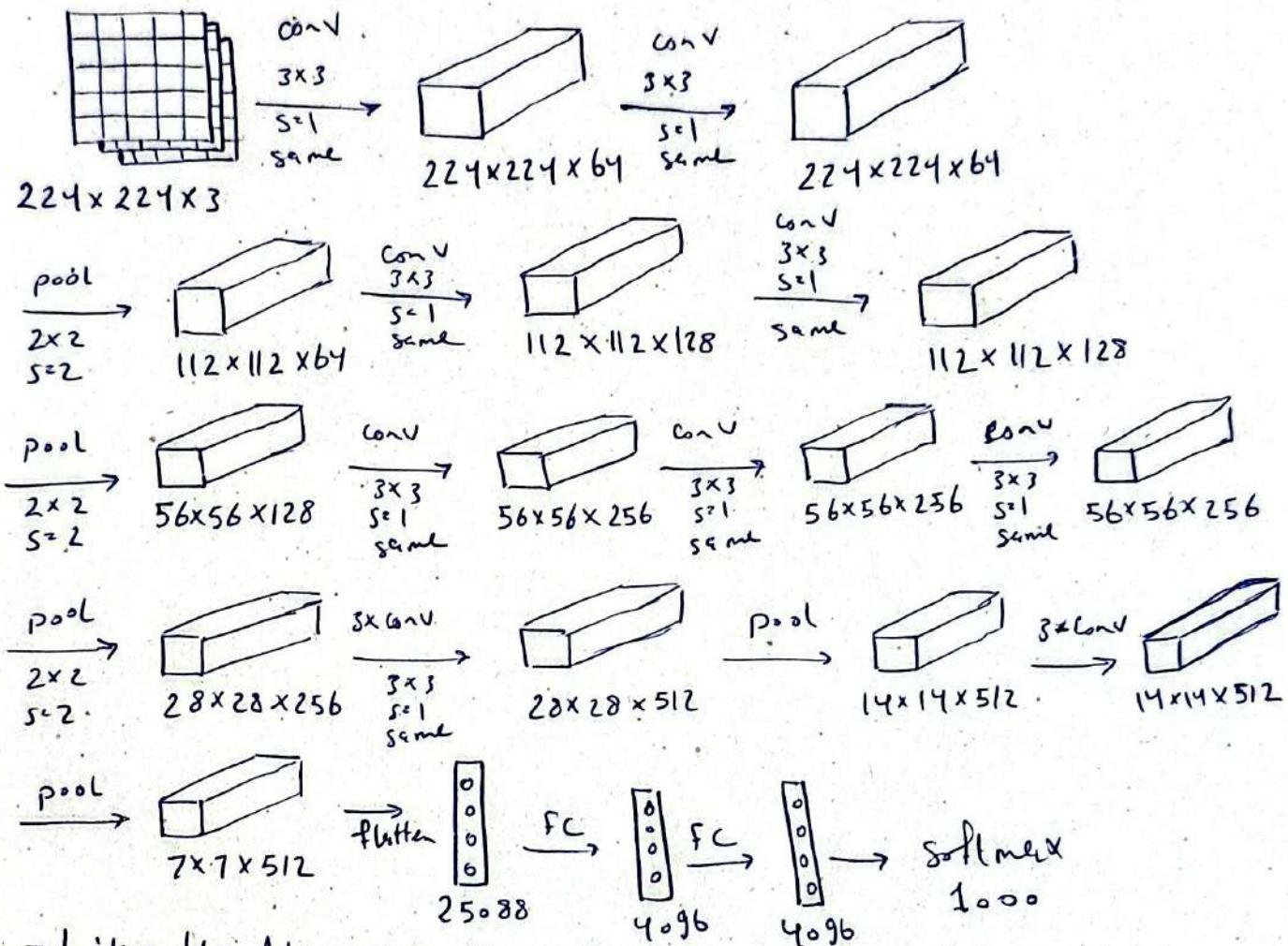


→ The AlexNet had a lot of similarities to the LeNet-5 neural net work but it was much bigger. The LeNet-5 had around 60k parameters, but the AlexNet has around 60 million parameters which is almost 1000 times bigger than the LeNet-5.

→ Another aspect why this was better than the LeNet-5 is that it uses ReLU (rectified Linear unit) activation function instead of the sigmoid function like in the LeNet-5 but back in 1998 when the paper for the LeNet-5 was written, ReLU was not used commonly yet.

## VGG-16

→ In this neural net the researchers decided to go with a much simpler network in their approach compared to something like AlexNet; they used Conv layers with  $3 \times 3$  filters, stride = 1 and used same padding, in the Max pooling layers they used  $2 \times 2$  filters with stride = 2 which is very simple.



→ Like the AlexNet, this also is a pretty big neural net with around 138 million parameters, it uses simpler convolutions but is way deeper than the AlexNet.

→ It's called the VGG-16 because it has 16 layers.

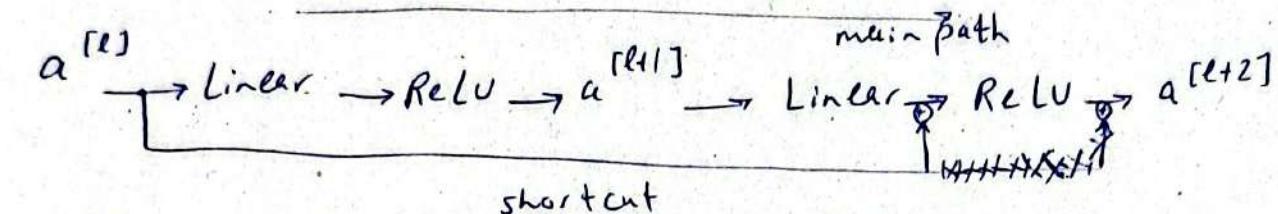
→ The main down side of this neural net is that it is a pretty big neural net even by modern standards.

## ResNets

→ very deep neural net works are harder to train because of problems like vanishing and exploding gradients so in this case study we'll see how skip connections are implemented.

\* skip connections allows us to take the activations from one layer and feed it to a layer much deeper in the neural net.

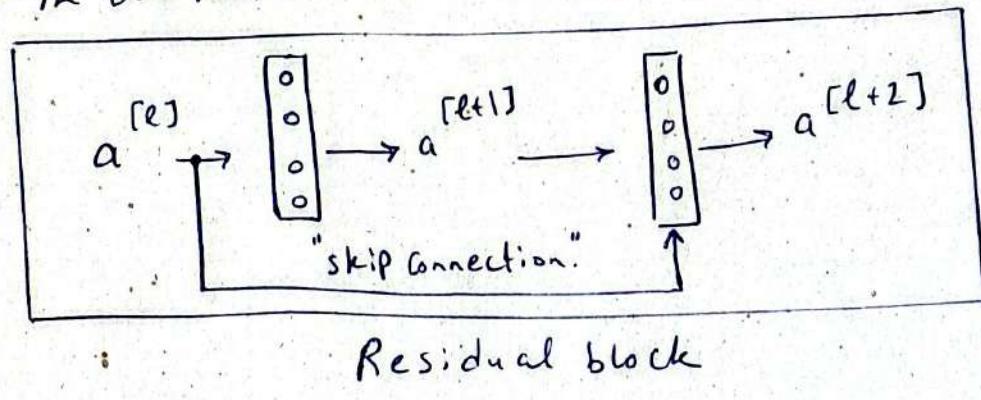
→ ResNets are built on something called a Residual block;



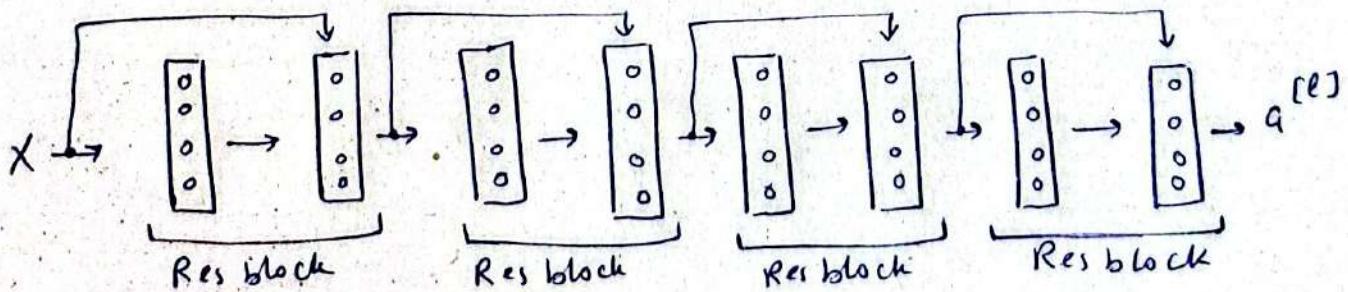
$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]} \rightarrow a^{[l+1]} = g(z^{[l+1]})$$

$$z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]} \rightarrow a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

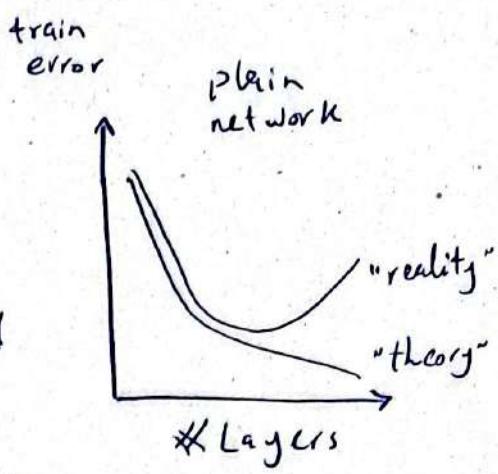
\* so what happened here is that we took  $a^{[l]}$  and we directly fed it or added it to another activation much deeper in the network.



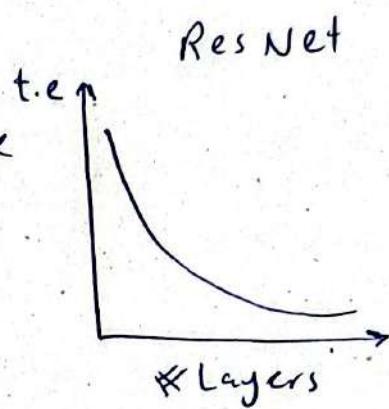
## ResNet



\* so in theory; training a neural network or like what the authors of the Res Net's paper called it; "plain network" (meaning it does not have any residual blocks); is supposed to have lower training error as we increase the number of layers but in reality; the ~~no~~ training error actually increases.



→ while in plain networks the training error increases as you increase the number of layers; this does not happen when we're using Res Nets because Res Nets handle the problem of exploding and vanishing gradients much better than the plain neural networks.



### Why Res. Nets work so well?

\* complicated; it just works..

What does a  $1 \times 1$  convolution do?

→ if you think about it in terms of a having a picture with only one channel then the  $1 \times 1$  convolution is only multiplying the pixel value with the weight of the filter

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline \end{array} \quad * \quad \boxed{2} \quad = \quad \begin{array}{|c|c|c|} \hline 2 & 4 & 6 \\ \hline \end{array}$$

$5 \times 5 \times 1$

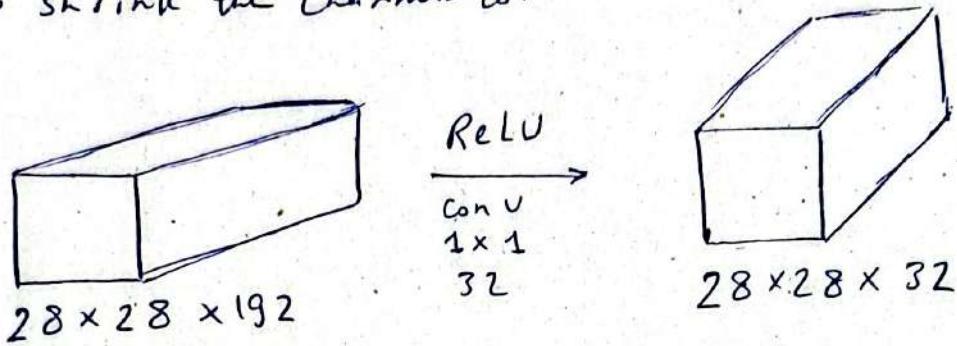
\* In reality the  $1 \times 1$  convolutions are not used with images that have only 1 channel but rather with matrices that can have larger number of channels

$$\begin{array}{c} \text{3D Volume} \\ 6 \times 6 \times 32 \end{array} \quad * \quad \begin{array}{c} \text{Filter} \\ 1 \times 1 \times 32 \end{array} \quad = \quad \begin{array}{|c|c|c|} \hline \text{6 filters} \\ \hline \end{array}$$

→ so actually in this case the value that we're computing is being computed across the whole channels of the  $6 \times 6 \times 32$  volume not just multiplication of two numbers and in a single step of computing the output volume you can imagine it as a fully connected layer between 32 and 32 hidden units you can imagine it as applying ReLU to some values multiplied by the weights

→ This is sometimes called "Net work in Net work".

\* a useful use of the  $1 \times 1$  convolutions is that if you have for instance an input volume that has grown too big in number of channels and you want to decrease the number of channels then you can use  $1 \times 1$  convolutions with the number of filters that you want to shrink the channels to.

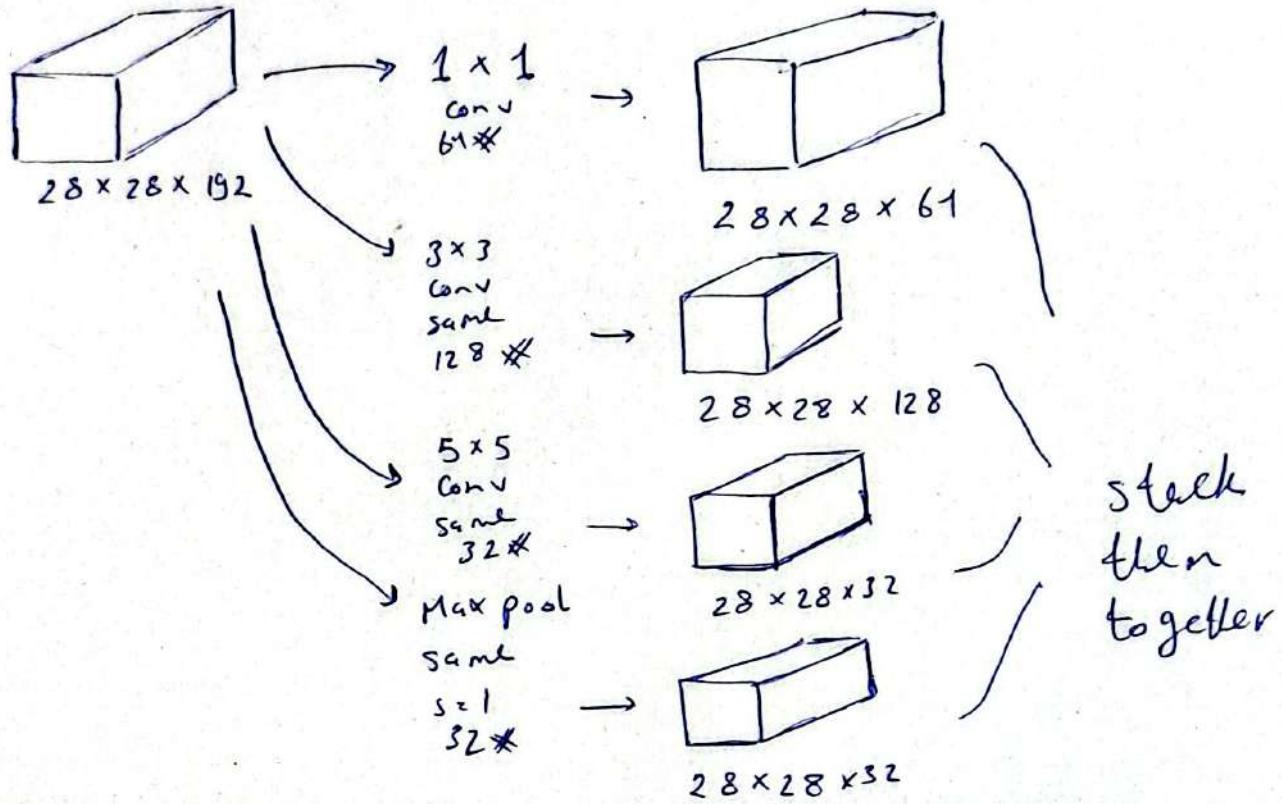


- with each  $1 \times 1$  filter having 192 channels as the input volume.
- so the idea of having  $1 \times 1$  convolutions allows you to shrink the number of channels allowing you to save on computations without shrinking the actual height or width of the input volume.
- even if we don't use the  $1 \times 1$  convolutions to shrink the number of channels it still plays a role as it applies ReLU or nonlinearity to its inputs so it doesn't have to be used only to shrink channels.

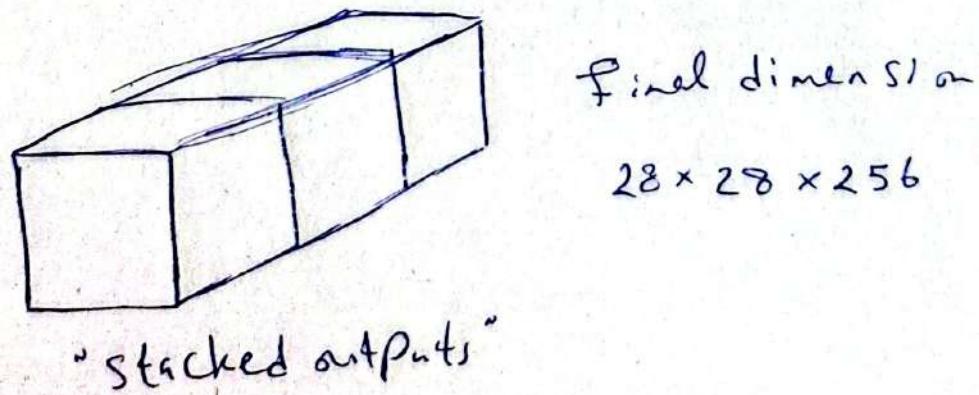
## Inception network

→ the idea behind the inception neural net is that it allows you to do more than one convolution; even maxpooling in the same computational process, so how does that exactly work?

- let's say you have a  $28 \times 28 \times 192$  input volume



• and then we have an output volume of the stacked volumes from each process and we can do that because we paid attention to the dimensions of each output volume

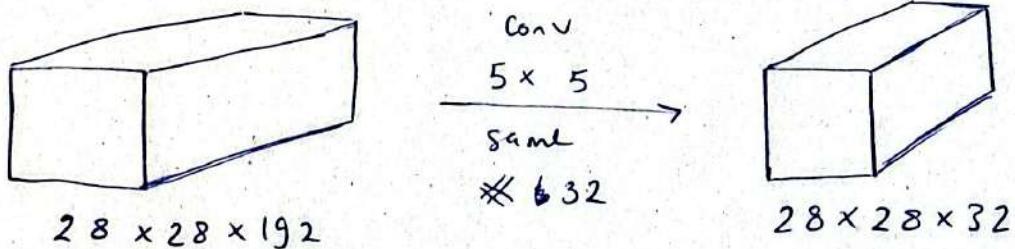


so instead of having to pick a single computation for each layer like a certain kernel size with a certain number of filters for a single conv layer; you can do them all and have multiple convolutions and even max pooling like we said in a single "Inception module" and the network will learn whatever parameters it wants to use for whatever combinations of these filter sizes that we have.

The problem with Inception layers:-

↪ Of course as it's predicted the inception layer or module has a problem which is the Computational Cost.

- For just the  $5 \times 5$  output block that we wanted to stack in the previous illustration the computational cost is :-



↪ 32 filters with each being  $5 \times 5 \times 192 = 4800$  values

these values will be multiplied by  $28 \times 28 \times 32 = 25088$

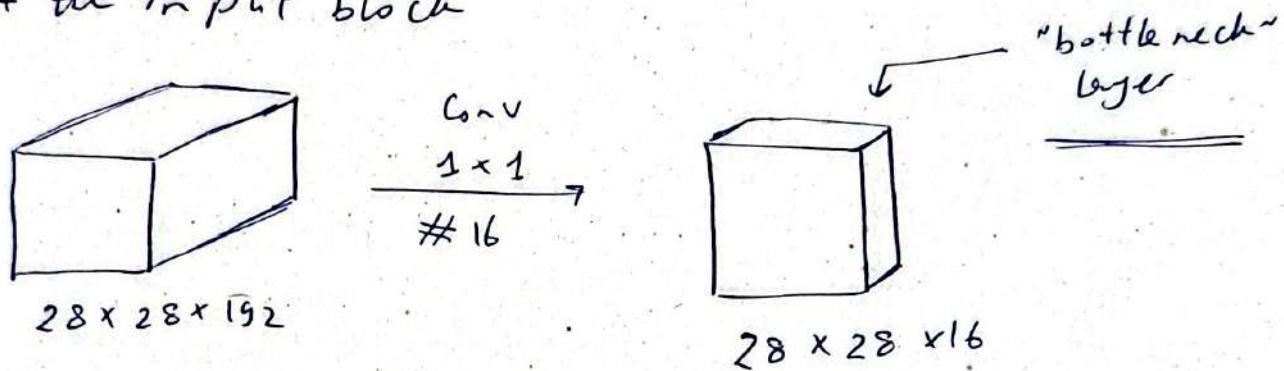
- so the full number of computations for this step will

$$\text{be } 25088 \times 5 \times 5 \times 32 \approx 120 \text{ million}$$

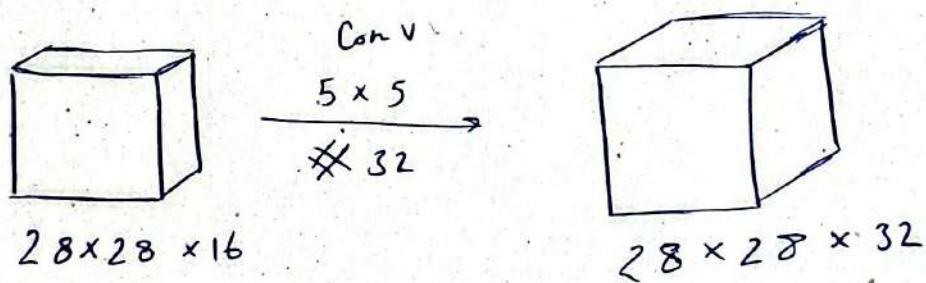
✓  
number of computations  
for each output block      ↓  
                                  number of computations  
                                  for each value in each  
                                  Filter

↳ while 120 million computations on a modern Computer is possible; it is still very computationally expensive.

- Here comes the role of  $1 \times 1$  convolutions as it can decrease the number of computations needed for this step by nearly a factor of 10; making them  $\approx 12$  million computations only instead of 120 million.
- First we use  $1 \times 1$  convolutions to reduce the ~~then~~ volume of the input block



- Then we can use the  $28 \times 28 \times 16$  output volume and run the  $5 \times 5$  convolution on it.



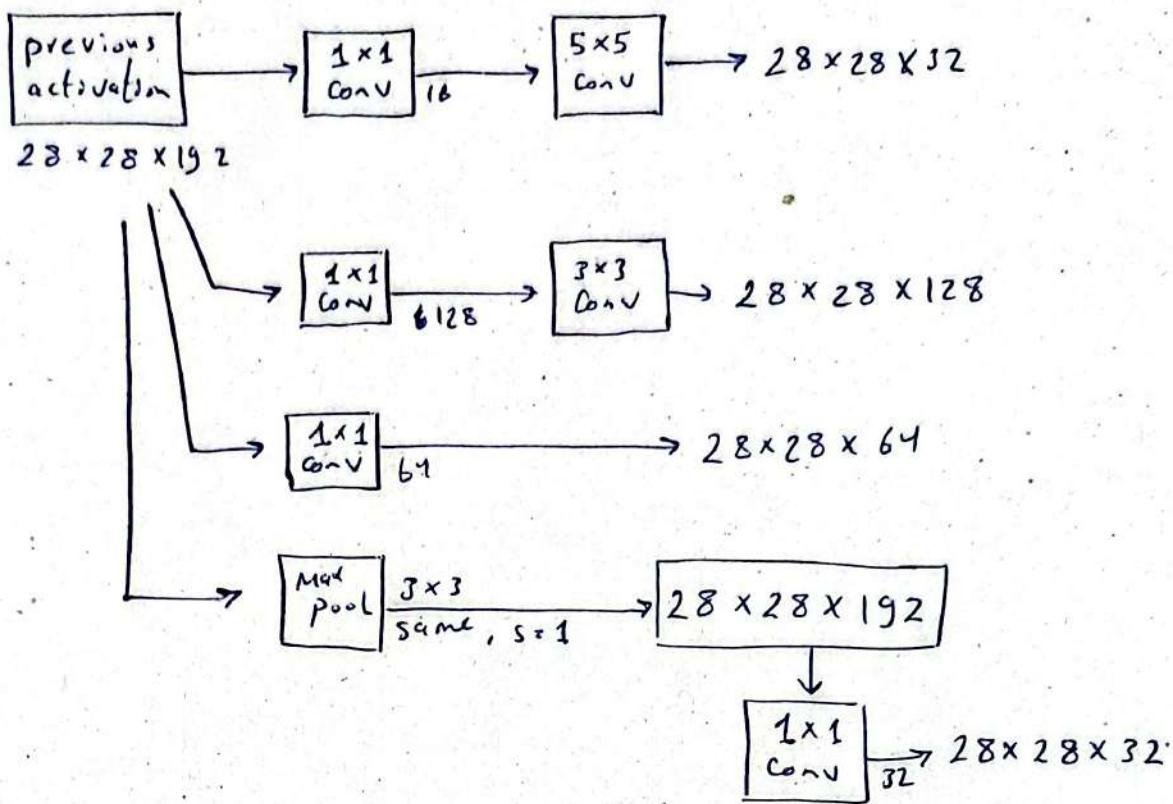
- Computational cost for the past step:-

$$\rightarrow \text{Conv 1} \rightarrow 1 \times 1 \times 192 \times 28 \times 28 \times 16 = 2.4 \text{ million}$$

$$\rightarrow \text{Conv 2} \rightarrow 28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10.0 \text{ Million}$$

total  $\approx 12$  million

## Inception network



and then we concatenate all the output volumes

resulting in a  $28 \times 28 \times 256$  output volume.

→ This was the inception module; what the inception network does is put a lot of inception modules together.

- one famous inception network is the googlenet neural network developed by google research team.

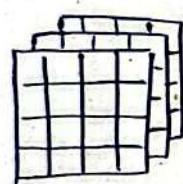
## Mobile Net

↳ Mobile net is another foundational neural network architecture used in low compute environments such as mobile phones.

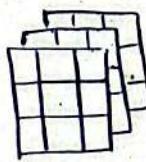
↳ The previous neural networks that we learned about are actually computationally expensive and if we want a neural net to work on a low compute environment such as mobile phones then we need another neural network architecture that can perform better.

∴ Normal vs depth wise separable convolutions.

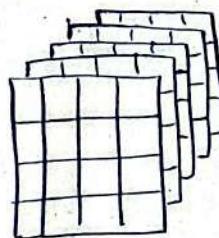
↳ normal convolutions :-



$$b \times 6 \times 3 \\ n \times n \times nc$$



$$3 \times 3 \times 3 \\ f \times f \times nc \\ padding = 0 \\ stride = 1$$



$$4 \times 4 \times \frac{\text{number of filters}}{nout \times nout} \quad (5)$$

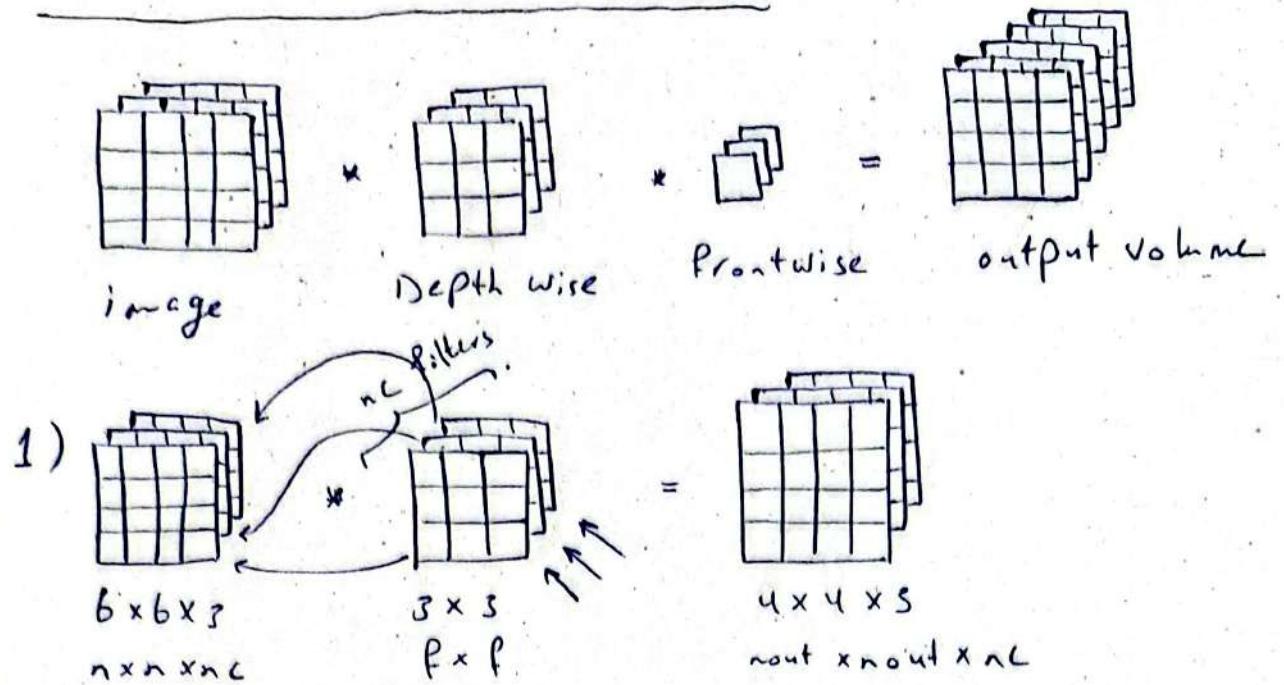
→ Computational Cost :-

Cost = \* filter params \* filter positions \* # of filters

$$= 3 \times 3 \times 3 * 4 \times 4 * 5$$

$$= 2160$$

## Depth wise separable convolutions



The way that we apply these convolutions is that we take each channel of the filter and apply it to the corresponding channel in the input volume; not like we did in the normal convolutions where we applied the whole filter by the whole slice of the image.

↳ By convolving each channel of the filter with each channel of the image we get 3 three 2-dimensional matrices then we stack these together to get our output volume.

↳ in the previous process, we got each value of the output matrix by multiplying only 9 values from the filter by 9 values from the corresponding image channel; much less than we used to do in normal convolutions.

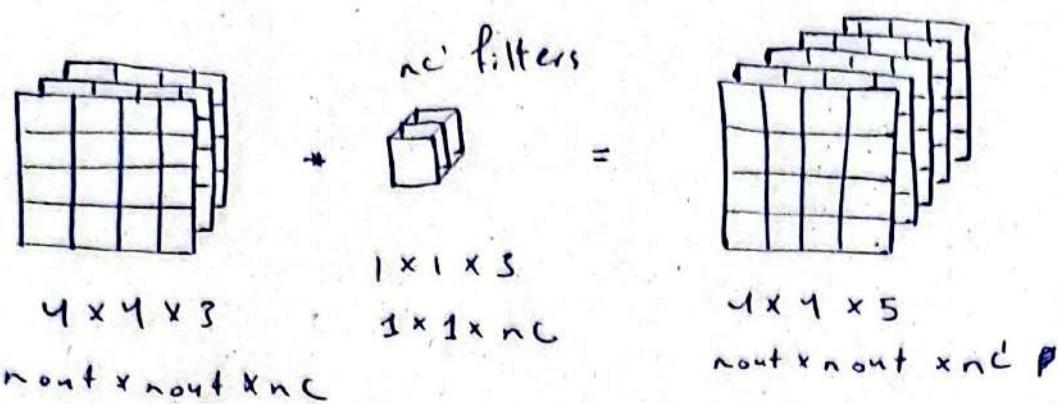
Cost =  $\#$  filter params  $\times$   $\#$  filter positions  $\times$   $\#$  of filters

$$= 3 \times 3 \times 4 \times 4 \times 3$$

$$= 432$$

which is much less than normal convolutions

↳ We're not done yet as this was the depth wise part of the convolution only; now we still have to do the front wise convolution.



→ Now with this step we would do as we did with normal convolutions, we have a  $4 \times 4 \times 3$  input volume and to get the same output as the  $4 \times 4 \times 5$  normal convolution that we did before we have to convolve our input with 5 filters that are  $1 \times 1 \times 3$ .

↳ Computational cost for this step:-

$$\text{Cost} = \frac{1 \times 1 \times 3}{\text{filter params}} \times \frac{4 \times 4}{\text{Filter position}} \times \frac{5}{\text{number of filters}}$$

$$= 240 \text{ multiplications}$$

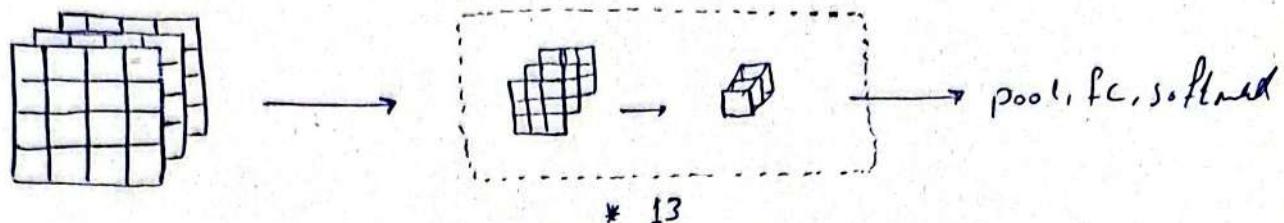
Now the total number of computations needed by this depth wise separable convolution to get to the same output volume as the normal convolution is

$$= 432 + 240 = \underline{\underline{672}} \text{ multiplications.}$$

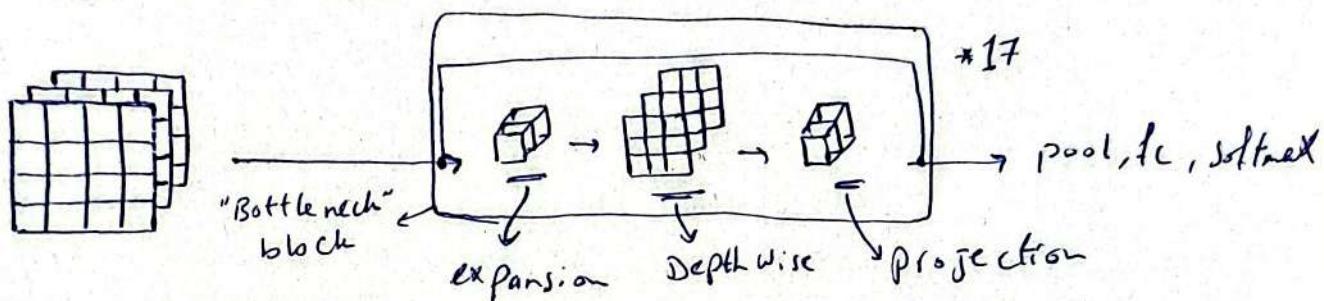
→ which is much less than the 2160 that was needed by the normal convolution.

→ The general formula for relative computation used to normal convolutions =  $\frac{1}{n_{in}} + \frac{1}{f^2} = \frac{1}{5} + \frac{1}{3^2} = 0.51$

## MobileNet Architectures.



→ In MobileNet Version 1 they fed the input volume to 13 blocks of Depth wise separable convolutions and then applied some pooling followed by some fully connected layers and finally a softmax function for classification.

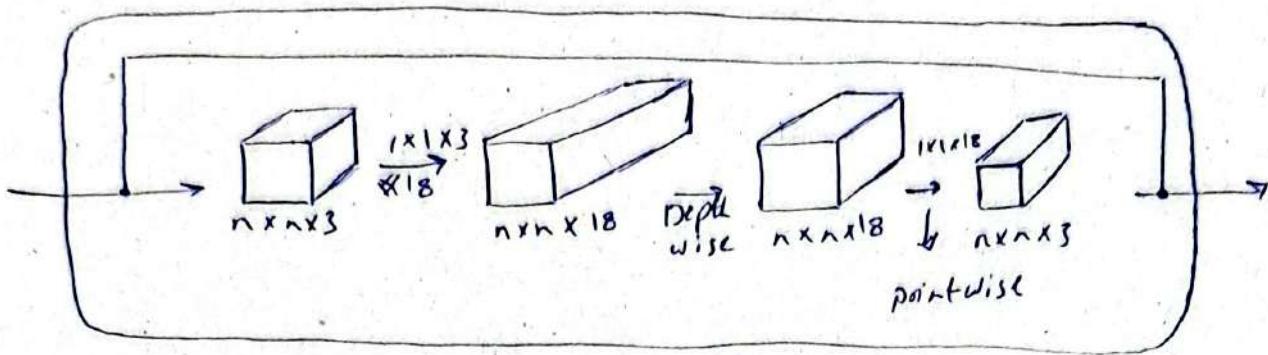


→ In MobileNet Version 2 they applied the concept of skip connections aka (Residual Nets) to the input volume and combined it with the Depth wise separable convolutions; they first fed the input volume to an expansion layer then to the normal Depth wise layer followed by the front wise layer (here called projection layer).

→ Using the Residual Connection allows the gradients to propagate more efficiently.

→ Just like MobileNet v1 ; MobileNet v2 fed the input volume not to 1 but to 17 consecutive residual blocks and then did some pooling followed by fully connected layers then softmax.

## MobileNet V2 bottleneck



→ so what happened in this Residual block or bottle neck

block that we made is:-

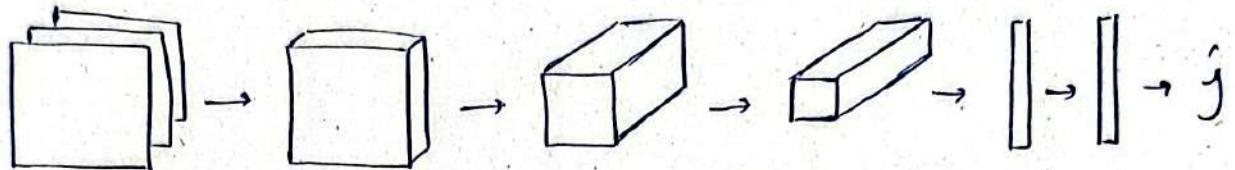
- we fed the input volume to convolution layer that uses  $1 \times 1 \times 3$  filters, specifically 18 of them and by using the  $1 \times 1 \times 3$  we maintained the height and width but with 18 filters we increased the number of channels to 18 hence why it's called an expansion layer.
- we then took the  $n \times n \times 18$  output volume and applied a Depth wise convolution to it which will maintain the same height and width. and also the same number of channels because it convolves each filter with the corresponding channel.
- we then fed the  $n \times n \times 18$  output volume that we got from the depth wise convolution and fed it to a Point wise (front wise) convolutional layer, which is a typical convolution with filters of size  $1 \times 1 \times n_c$ , specifically with 3 filters resulting in an output volume of  $n \times n \times 3$ .

→ The expansion Layer expands the information from the input volume which allows the model to learn better which is good over all for the performance.

→ But since we're doing all this in low compute environments we have to project the output volume down again to save up on computations.

## EfficientNet

→ MobileNet v1 and v2 gave us a way to implement a neural network that is more computationally efficient but there is a way to scale or computations either up if we have more computational resources or down if we're constrained by less resources; this is EfficientNet.



→ Let's assume that we have a "Baseline neural network" in which we have a image of resolution  $r$  and depth of the network  $d$ , and the layers have a certain width  $w$ .

∴ The three things that you can do to scale things up or down are

↳ use a higher or lower resolution ( $r$ ).

↳ Make the network more or less deep ( $d$ ):

↳ Make the layers more or less wide ( $w$ ).

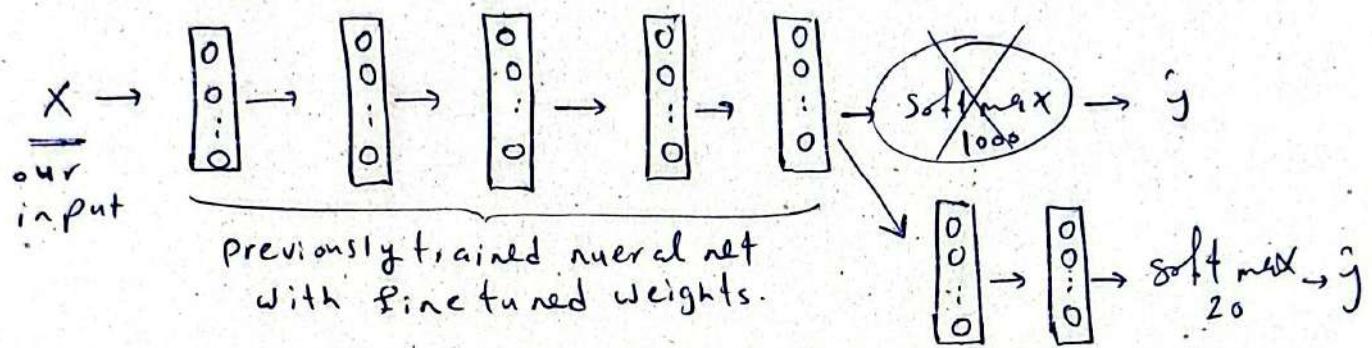
→ so the three things that we can control are resolution, depth and width of the network.

• Now the question is, if we want to scale up  $r$ ,  $d$  and  $w$  what's the rate at which we scale them up?

↳ a good way to consider this is to look up some of the open source implementations of EfficientNet, this will help us decide the rate at which we chose to scale up to.

## Transfer Learning

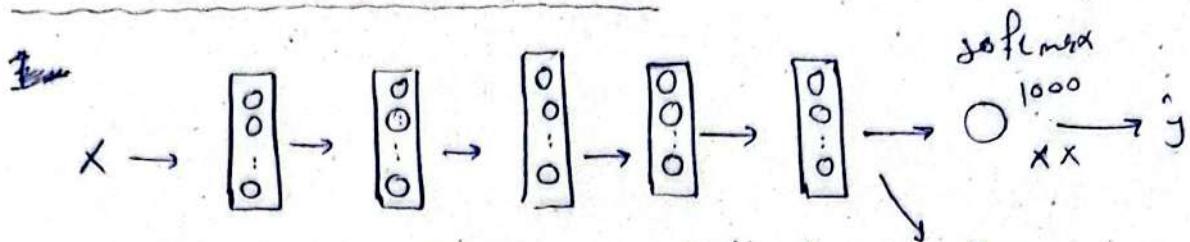
↳ Transfer learning is the process of transferring knowledge or information from a model that was previously trained and fine tuned to our model that we're trying to train, we do this by using the weights of these trained models as our initialization weights and of course this is much much better than randomly initializing the weights.



→ So it's just like getting rid of the output layer of a previously trained model and adding a couple of new layers so the model can adapt to the new data which are our data set and then we can have a new softmax with number of hidden units = our data set classes

→ Transfer learning also helps when we have a small data set so we don't have a lot of training examples to train the model on, in this case we can just use the knowledge gained by a previously trained model and adapt this knowledge to our new data set by fine tuning the weight a bit on our data.

## approaches to transfer learning

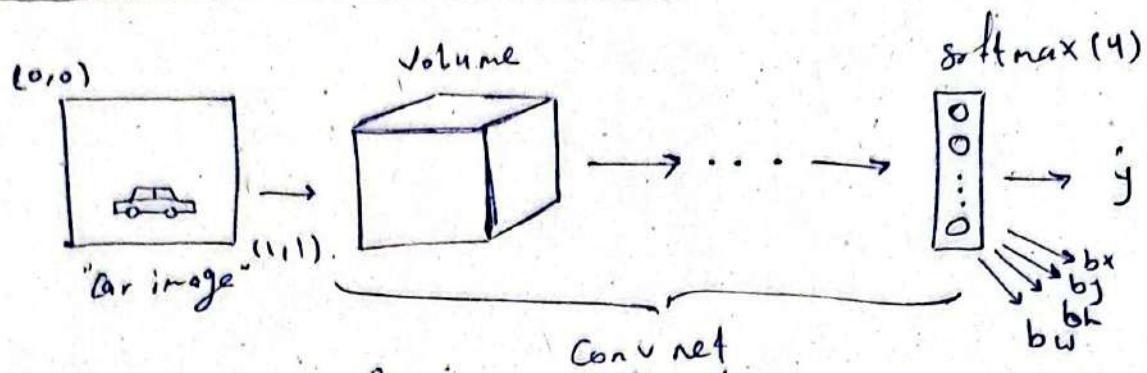


→ the first thing that we need to do in transfer Learning is to download a pretrained neural network like image net, this pretrained neural net will have its own architecture and pretrained weights and based on the dataset in hand we'll do one of the following options.

- 1 - if we have a very small dataset we can just ~~use~~ feed our training examples (the whole dataset) to the pre trained neural net without changing a thing; actually we'll freeze the weights in the layers and just feed the data to the neural net and classify the images using the neural net's weights that we downloaded.
- 2 - if we have a bigger dataset; we can freeze some of the earlier layers and just train some of the layers at the end of the neural net; so we'll freeze the weights in some layers but retrain and update the weights in some of the layers at the end, or we can just remove some of the last layers and create new ones with our own new hidden units.
- 3 - if we have a lot of data we can just use the neural net's architecture and the pretrained weights and use them for initialization and retrain the whole network on our data with the weights that we already have as initial values and use the same or almost same architectures.

## Week 31-

### Classification with Localizations



→ We've seen before the classification problems where we had a picture on which we did some convolutions turning it to a vector of features in the end and we fed that vector of features to a softmax function and that softmax layer predicts the class in the image; this week we'll add to that a bounding box around the object that we're trying to classify so now we're given ~~and~~ an image and our goal is not only to classify the image but also to localize it first in the image.

→ The classes that our softmax function has to predict may be:-

1- pedestrian

2- car

3- Motor cycle

4- back ground (picture with no target objects)

- So far our "Bounding box" we'll need our model to output some parameters associated with the location of the object in the image; these parameters are:-

1-  $b_x$  → mid point of the Bounding box

2-  $b_y$  →

3-  $b_h$  → Height of the box

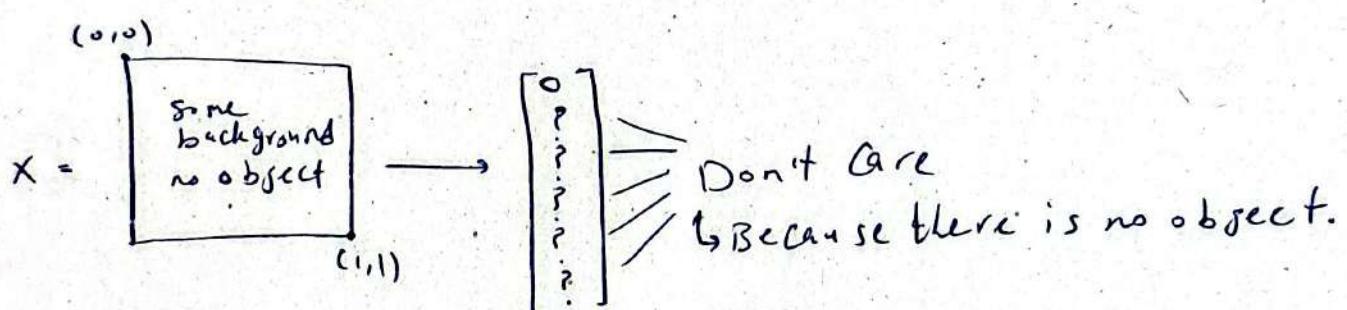
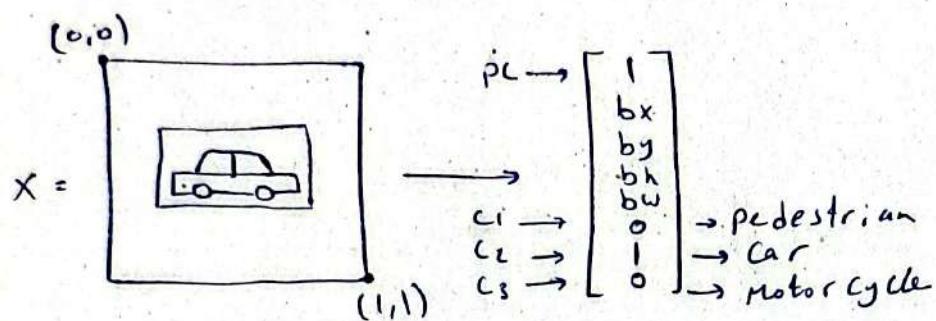
4-  $b_w$  → width of the box

→ with the new parameters, the output  $y$  will change to include  $bx$ ,  $by$ ,  $bh$  and  $bw$  + class labels.

$y = \begin{bmatrix} pc \\ bx \\ by \\ bh \\ bw \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$

- is there an object?
- Bounding Box
- if  $pc$  is true then which object is it?

\* with the problem that we're dealing with; we assume that there is only one object in the image.



Loss Function:-

→ if we're using squared error

if  $pc = 1$

$$L(g, y) = \{(g_1 - y_1)^2 + (g_2 - y_2)^2 + \dots + (g_8 - y_8)^2\}$$

if  $pc = 0$

↳ there are 8 components of the output vector

$$L(g, y) = \{(g_1 - y_1)\}$$

→ we said we'll use squared error just for the simplicity of the example but in reality we'll probably use log likelihood loss for  $c_1, c_2, c_3$ ; squared error for the rest.

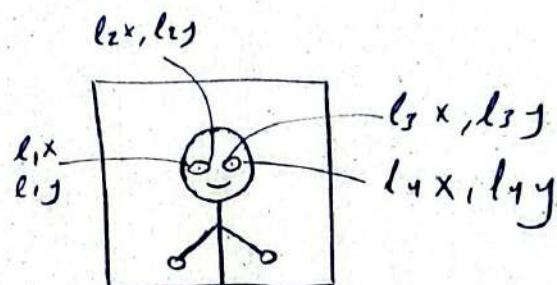
## Land mark detection.

→ a Land mark is a  $(x, y)$  location in the image at which you want the neural network to point.

↳ May be in Face recognition problem this Land mark could be the corner of someone's eye.

\* In general Cases you can have a neural network just output these Land mark in an image; for example, for the Face detection problem you can have the neural network output  $l_1 x, l_1 y, l_2 x, l_2 y$ .

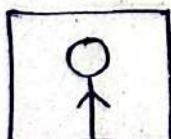
$l_1 x, l_1 y \rightarrow$  left corner  
of the  
first eye



and so on.

→ Generalizing this idea; you might want the neural network to output a set of numbers that point to some landmarks in the image; like a punch of points for the eyes; another punch for the mouth ... etc.

\* Of course for a neural network to be trained to do such things you have to have Labeled data of images with Labeled Land marks.



→ Conv net →

$\begin{bmatrix} PL \\ l_{1x} \\ l_{1y} \\ \vdots \\ l_{4x} \\ l_{4y} \end{bmatrix}$  → is there a face?  
↔ positions for 6<sup>4</sup> Land marks in the image.

→ These Landmarks can be used to classify the emotions or the facial expressions in the image such as smiling or frowning and so on.

→ Some apps use these Landmarks to add filters to images like Instagram and Snapchat.

## Object detection.

### Sliding window detection

→ In the sliding window detection you first get some closely cropped car images [car icon] no surroundings; just a small cropped picture of the car and train a classification Conv net on these images having the network output 1 for Car and 0 for no car.

→ Then you would input it to this Convnet small squares from the image you want to detect an image in and get your output; You slide the window to the right and continue the process.

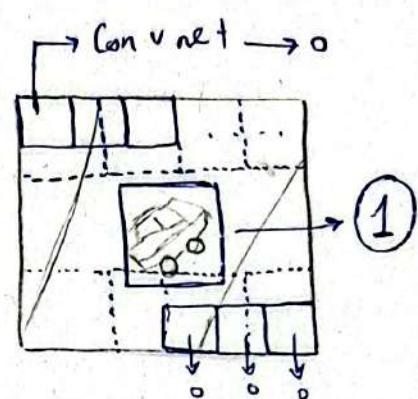
→ Having done the previous process once; we now take larger regions or larger squares and feed them to the Convnet.

→ Then you might do it even a third time using larger regions and also feed them to the Convnet.

→ Hopefully by doing this you'll finally be able to predict if there's a car in the image or not.

→ Of course the disadvantage of sliding window detection is the computational cost because you're taking all these square regions and feeding them to the Convnet and each of these processes require some computations.

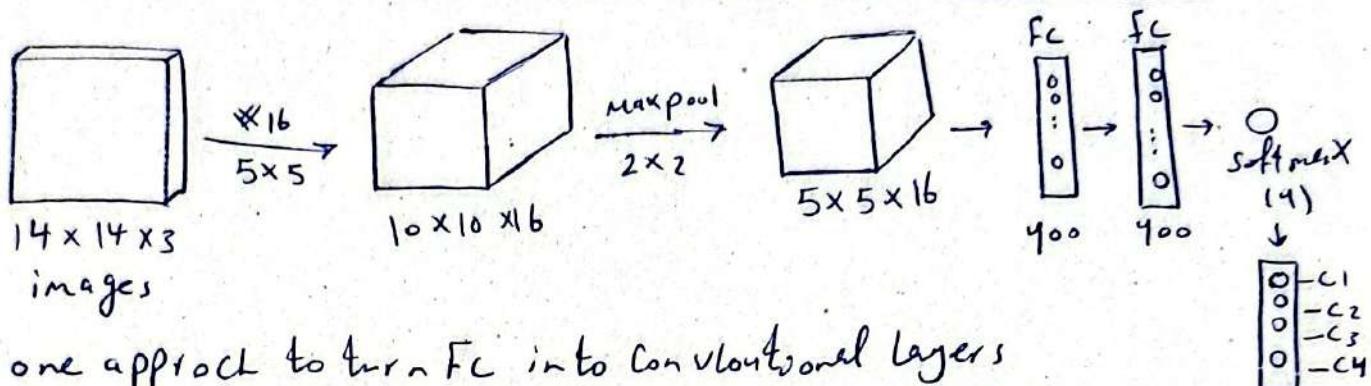
→ If we use big stride (skipping some pixels between each window) this of course would decrease the number of windows but also it would probably hurt performance.



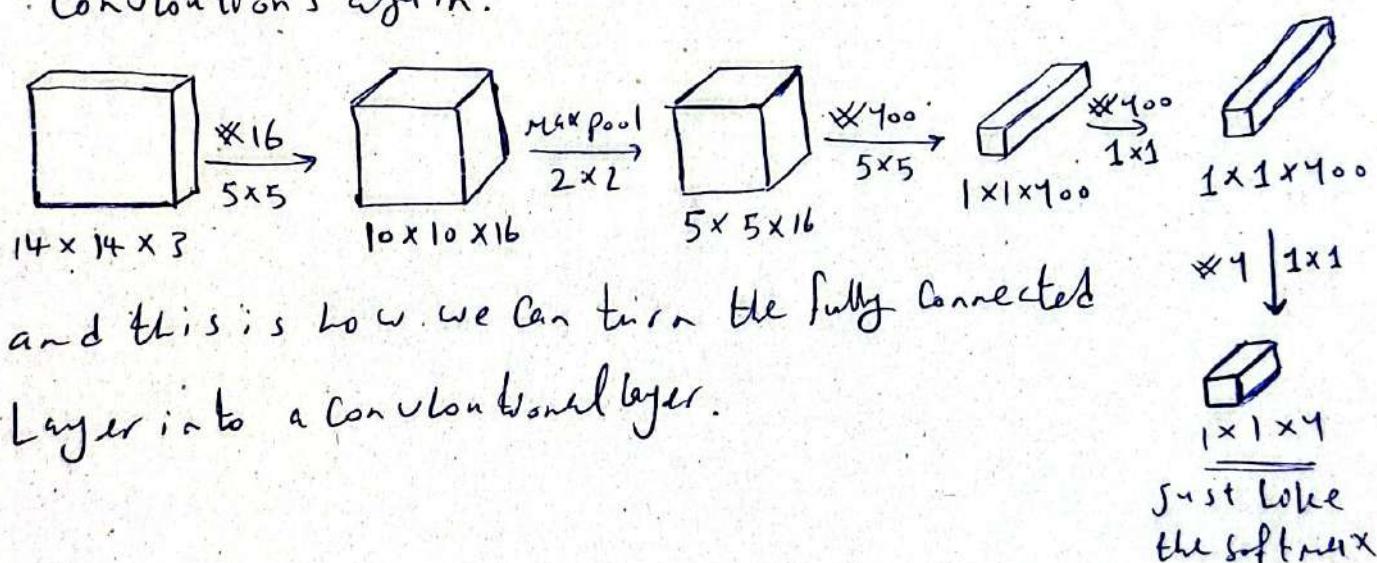
## Convolutional implementation of sliding windows

→ Since sliding window algorithm is computationally expensive when dealing with it and feeding each image to the conv net; an implementation of convolutional sliding window was made to reduce computations.

### Turning Fully Connected Layers into convolutional layers



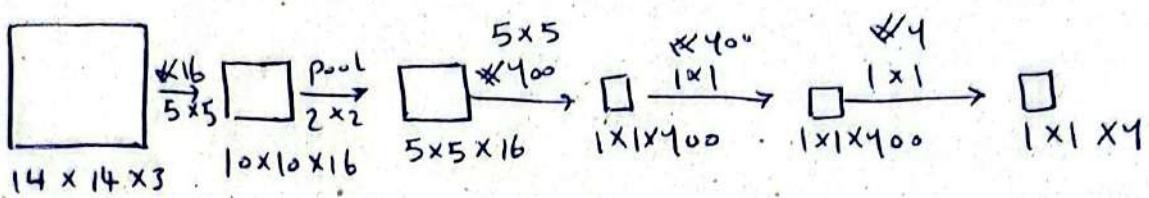
→ one approach to turn FC into convolutional layers is instead of flattening the output volume of the last convolutional layer, we'll just turn it into  $1 \times 1 \times n_c$  output volume and for the next FC we'll do  $1 \times 1 \times n_c$  convolutions again.



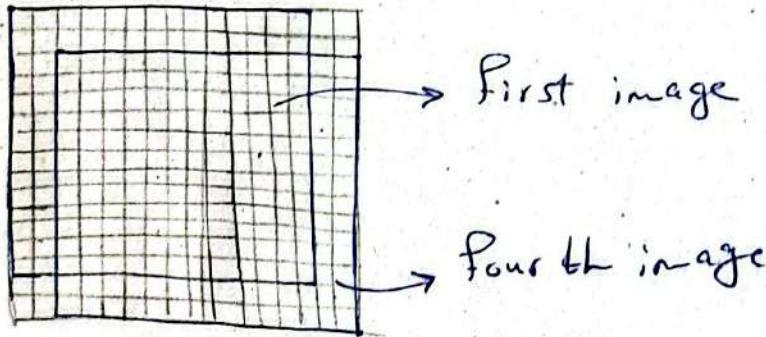
- What if the image that we're inputting to the Convnet is bigger than the image size the convnet was trained on?  
↳ We use sliding window.

\* But since the normal sliding window is like we said before computationally expensive we'll have to do it using convolutions.

→ Let's say we have  $16 \times 16 \times 3$  images and the Convnet was trained on  $14 \times 14 \times 3$  images:-

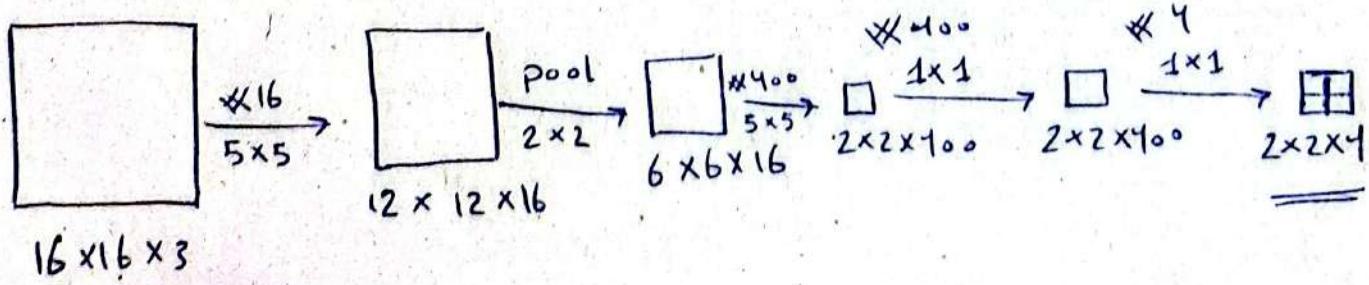


→ Now since our image is  $16 \times 16 \times 3$  we can't just feed it to the Convnet; we'll have to use  $14 \times 14$  windows and a stride of 1; this will result in having 4 windows



so instead of having to feed multiple smaller crops of the same image to the Convnet, you just feed the Convnet the whole image and have multiple output vectors one for each crop

→ so instead of feeding 4 different images to the Convnet we'll just feed it the  $16 \times 16 \times 3$  image to the convolutional layers and have our output be  $2 \times 2 \times 4$  with each vector corresponding to one of the four images



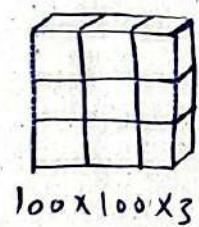
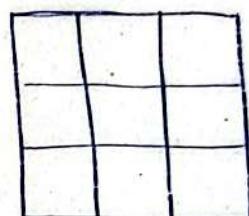
## Bounding Box predictions.

### Yolo algorithm

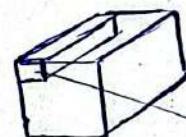
→ Yolo stands for "you only look once" and it's an algorithm that utilizes convolutional sliding windows to predict Bounding boxes around objects in images.

\* Instead of having to take multiple crops of images and feeding them to the CNN; now you can only look at the input image once as a whole and the algorithm will utilize convolutional sliding windows on them which will include partial crops to be convoluted with filters.

→ First we layout an  $5 \times 5$  grid on the image, like the illustration and then do some convolutions on the image and make sure that the output shape is a  ~~$5 \times 5 \times n_c$~~   $3 \times 3 \times 8$  image grid vector volume; each of the g vectors now has some values that correspond to one grid cell.



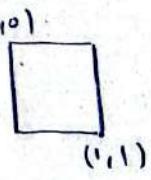
→ Conv net → maxpool → ... →



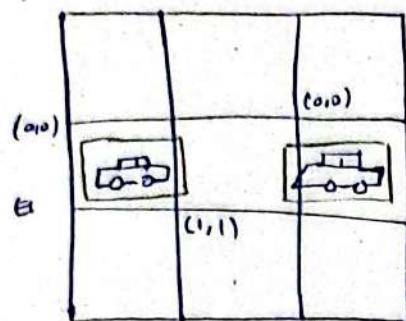
→ our output shape consists of g vectors each of them have 8 channels; these channels correspond to some values

$v_i = \begin{bmatrix} p_c \\ bx \\ by \\ bh \\ bw \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$  → is there an object in the grid cell  
 → Bounding Box  
 → Class 1?  
 → Class 2?  
 → Class 3?

\* each grid cell have coordinates from  $(0,0)$  to  $(1,1)$  so  $bx$  and  $by$  are  $0 < bx, by < 1$   
 → the bounding box can be larger than the grid cell so  $bh, bw$  can be greater than 1.



## specifying the Bounding Box



→ Coordinates for  $b_x$  and  $b_y$  are relative to the upper left and lower right corners of the cell grid as they are the  $(0,0)$  and  $(1,1)$  coordinates for the object.

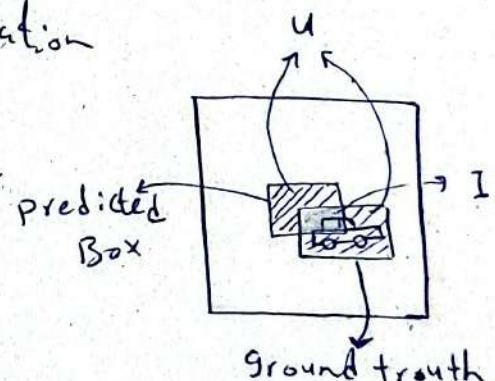
→ So  $b_x$  and  $b_y$  can never be 1 as they point to the center of each object unless the center of some object is at some corner  $0 < b_x, b_y < 1$  unless in a corner.

→  $b_h$  and  $b_w$  are calculated relative to the height and width of the grid cell but since the box can be larger than the grid cell; they can be greater than 1.

## Intersection over union:-

→ Intersection over union is a evaluation method for object localization; it computes the size of the intersection and divides it by the size of the union

$$\text{IoU} = \frac{\text{size of } \square}{\text{size of } \blacksquare}$$



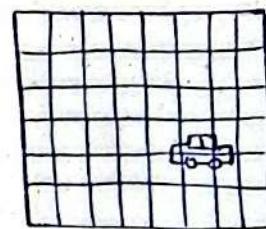
→ By convention; a lot of computer vision algorithms would judge a Bounding Box Prediction if  $\text{IoU} \geq 0.5$ .

→ If the predicted and the ground truth Boxes overlaped perfectly; the Intersection over union would equal 1.

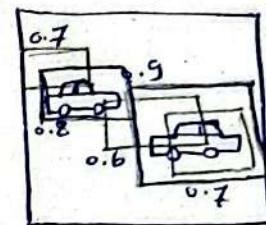
→ IoU is more generally a measure of the overlap between 2 Bounding Boxes.

## Non Max suppression :-

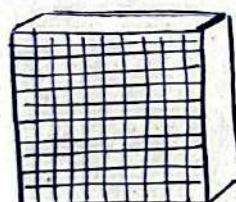
→ When we lay out the  $5 \times 5$  grid on the image each of the grid cells has an indicator whether the cell has an object or not and of course there can be a lot of cases where more than one grid cell has  $pc$  (the indicator of the object's existence) equal 1 indicating that the center of the object is in that grid cells. Non max suppression helps clean out these multiple predicted Bounding Boxes around an object to 1 Box.



→ Each of the Bounding Boxes has a parameter which tells us the confidence of whether an object exists in the Box or not; non max suppression first looks at these values, it takes the highest confidence value and assumes that the object exists in this Box so it suppresses all the other Boxes that have a high IoU with that Box and then looks at the next highest confidence and suppresses the other Boxes that overlap with it and so on.



→ All the non suppressed Boxes are our final predictions.

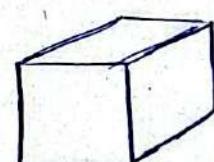


→ Conv net → Max pool →

input. image → Discard all boxes with  
19x19 grid  $pc < 0.6$

→ pick the Box with the  
highest  $pc$  and output  
it as a prediction

etc...etc → Discard any remaining Box with high IoU  
with the previous Box.



in this example  
we don't care about  
the class  
of the obj

$19 \times 19 \times 5$  of the  
output volume

## Anchor Boxes

→ one of the problems of object detection as we've seen it so far is that a grid cell can only detect one object but what if two object's centers are in the same grid cell?

→ what if we have 2 objects overlapping in the same grid?

\* If we go about this problem with our current implementation of YOLO

then we won't even be able to tell that

there are 2 objects in the image because the vector associated with this grid cell has only 1 PC.

→ To solve this problem; Anchor Boxes were introduced as a solution; Anchor Boxes are predefined boxes or shapes in the  $y$  vector; so instead of having just 1 PC and one of each  $b_x, b_y, b_L, b_W, c_1, c_2, \dots$  now we'll repeat these values again in the vector.

$$y = \left[ \begin{array}{c} p_c \\ b_x \\ b_y \\ b_L \\ b_W \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ \vdots \end{array} \right] \quad \left. \begin{array}{l} \text{Values associated with Anchor Box 1} \\ \text{Values associated with Anchor Box 2} \end{array} \right\}$$

→ The Anchor Box that gets assigned to an object is the one that has a higher IoU (overlap) with the ground truth Bounding Box.

obj1 → (gridcell, AnchorBox)  
↳ it gets assigned.

## Training with Yolo

→ The main idea of the Yolo algorithm is that we have an input image and some "Encoded" Labels, these encoded labels match our desired output volume shape which is a number of  $y$  predicted vectors each corresponding to one grid cell of the  $s \times s$  grid that we laid out on top of the input image, the number of elements (channels) of that vector can vary based on the number of anchor boxes that we specify for each grid cell.

→ This is a  $4 \times 4$  image on which we laid a  $4 \times 4$  grid, we need to encode the labels for this image so we can use it for training.

\* Each grid cell will have some values associated with a vector  $y$  corresponding to a number of anchor boxes; 2 in this example

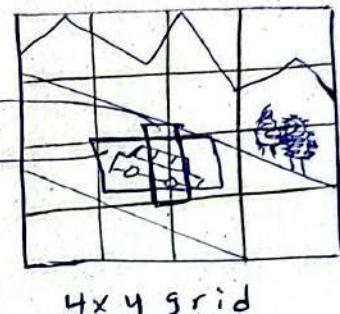
→ The values for each vector are as follows

→  $p_c \rightarrow$  is there an object in the bounding box

→  $b_x, b_y \rightarrow$  coordinates for the center of the object

→  $b_h, b_w \rightarrow$  height and width for the box

→  $c_1, c_2, c_3 \rightarrow$  which class is the object



$4 \times 4$  grid

$p_c$	/
$b_x$	/
$b_y$	-
$b_h$	-
$b_w$	x
$c_1$	-
$c_2$	o
$c_3$	o
$p_c$	/
$b_x$	/
$b_y$	-
$b_h$	-
$b_w$	x
$c_1$	-
$c_2$	-
$c_3$	o

→ Of course we have to encode these labels ourselves for each image in order to be able to use it in the training phase.

→ How do we encode these labels?

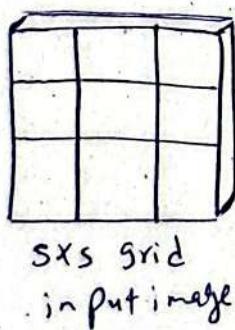
- \* For an input image like the illustration;
- We have 1 object which is one of our target classes (car, pedestrian, motorcycle)

→ So since we have 1 object; this object will have only one center point associated with a certain grid cell; in this example it's grid cell 10.

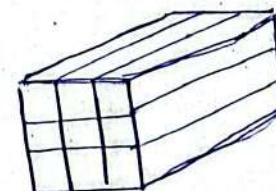
→ Only grid cell 10 would have  $pc=1$  and we will call it the ground truth, all other cells would have  $pc=0$  so the rest of the encoded vector of labels are values that we don't care about.

→ Grid cell 10 now has  $pc=1$ ,  $bx$ ,  $by$ ,  $bw$ ,  $bh$  have assigned values for the first anchor box and also the second anchor box resulting in a 16 channel vector.

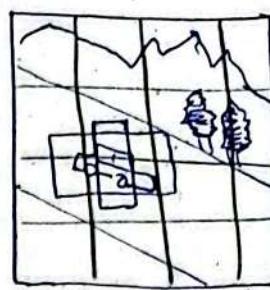
→ Now we can train our neural net using the input image + the  $S \times S$  output volume (labels)



→ Conv net → Maxpool →



$S \times S \times [\text{Anchor Boxes Params}]$

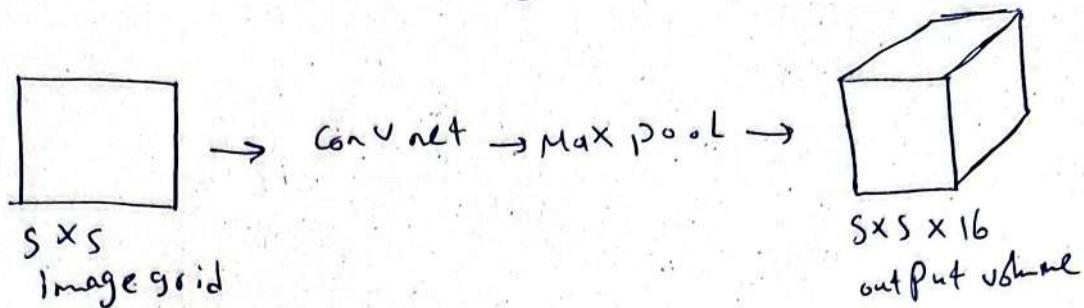


1
0.5
0.4
0.6
0.9
0
1
0
1
0.2
0.4
0.6
0.8
0
1
0

→ Each of the  $S \times S$  output volume vectors corresponds to a grid cell.

## Evaluation with Yolo

→ After being done with the training phase; now the Yolo neural net is ready to be tested.



→ We first take our input image and layout a  $S \times S$  grid on it; feed it to the Conv net and get our output volume.

\* Now we can evaluate our output volume using the IOU and non Max suppression first to drop one of the 2 anchor boxes that we have , using a threshold for the intersection over union between the predicted output volume and the ground truth encoded label.

For grid cell 10

1
0.5
0.4
0.6
0.9
0
1
0
1
0.2
0.4
0.6
0.8
0.9
0

ground truth

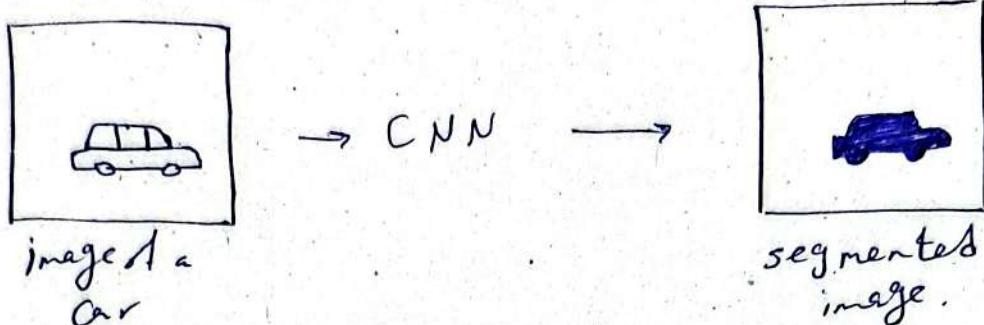
1
0.7
0.8
0.2
0.6
0
1
0
1
0.5
0.5
0.6
0.8
0
1

predicted

→ Now we can backpropagate using the loss function to improve the performance of the Neural Network.

## Semantic Segmentation

→ semantic segmentation is the process of inputting an image and outputting a segmented image that corresponds to the classes in the input image.



→ So how does that work?

↳ to train a segmentation CNN you have to have the training data as labeled segmented images and these labeled images will have a pixel value for each element in the image so if there is a car, a motor cycle and background then the car might have a pixel value of 1, the motor cycle 2 and the background 0.

→ Here the car has pixel value 1, the trees pixel value 2, the road pixel value 3 and everything else which is the background has pixel value 0.



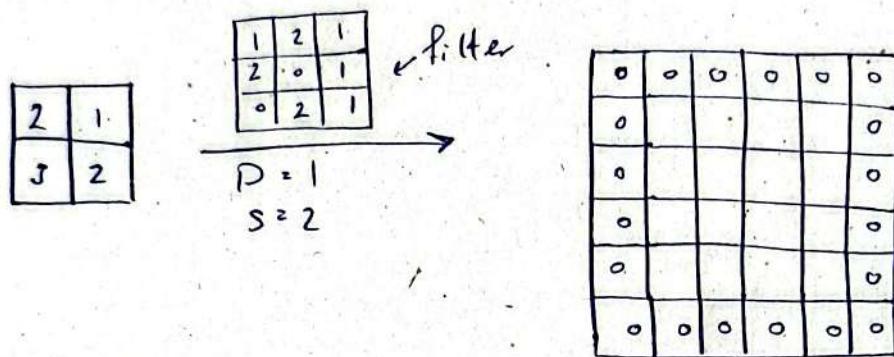
→ Again we find ourselves in front of a new problem which is that the output of the CNN has to be the same size as the input image, before we learned about  $1 \times 1$  convolutions which allowed us to maintain the same image size but here it would be computationally very expensive.

→ The solution to this problem is that we have to be able to down sample the image to decrease the number of operations per layer which lowers the computational cost and up sample the image again to have an output which is the same size as the input.

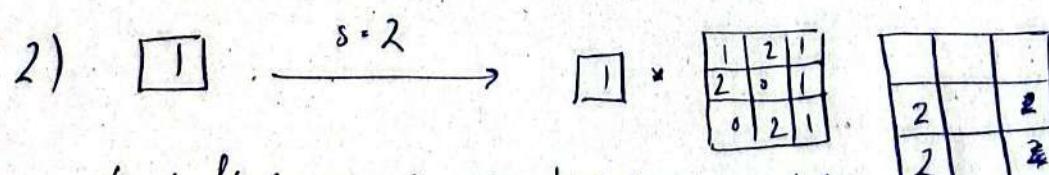
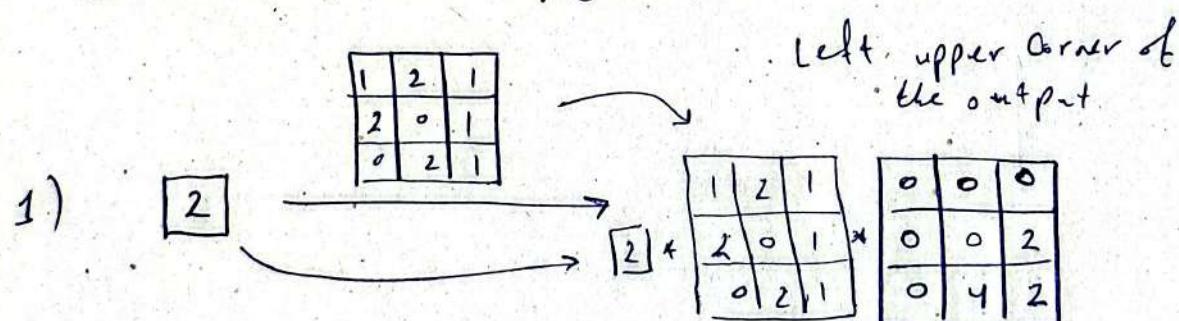
### Transpose Convolutions

→ Transpose Convolutions is the solution to our problem since it reverses the convolution process.

\* How is transpose convolution done?



→ Instead of placing the filter on the input; we would place it on the output this time and multiplying its values with the corresponding value in the input and go on and add the values of the multiplied numbers to get them in the output



The left lower corner has some numbers which are values produced by the previous computation we'll add the new values to them.

2)  $\begin{bmatrix} 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & 1 \\ 0 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 3 & 0 & 0 \\ 2 & 2 \\ 2 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 4 & 0 & 1 \\ 2 & 2 & 1 \end{bmatrix}$

3)  $\begin{bmatrix} 3 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & 1 \\ 0 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 4 & 2 \\ 0 & 0 \\ 0 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 10 & 5 \\ 0 & 0 & 3 \\ 0 & 6 & 3 \end{bmatrix}$

4)  $\begin{bmatrix} 2 \end{bmatrix} \rightarrow \begin{bmatrix} 2 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & 1 \\ 0 & 2 & 1 \end{bmatrix} + \begin{bmatrix} 5 & 2 & 1 \\ 3 \\ 3 \end{bmatrix} \rightarrow \begin{bmatrix} 7 & 6 & 3 \\ 7 & 0 & 2 \\ 3 & 4 & 2 \end{bmatrix}$

so the output will be →

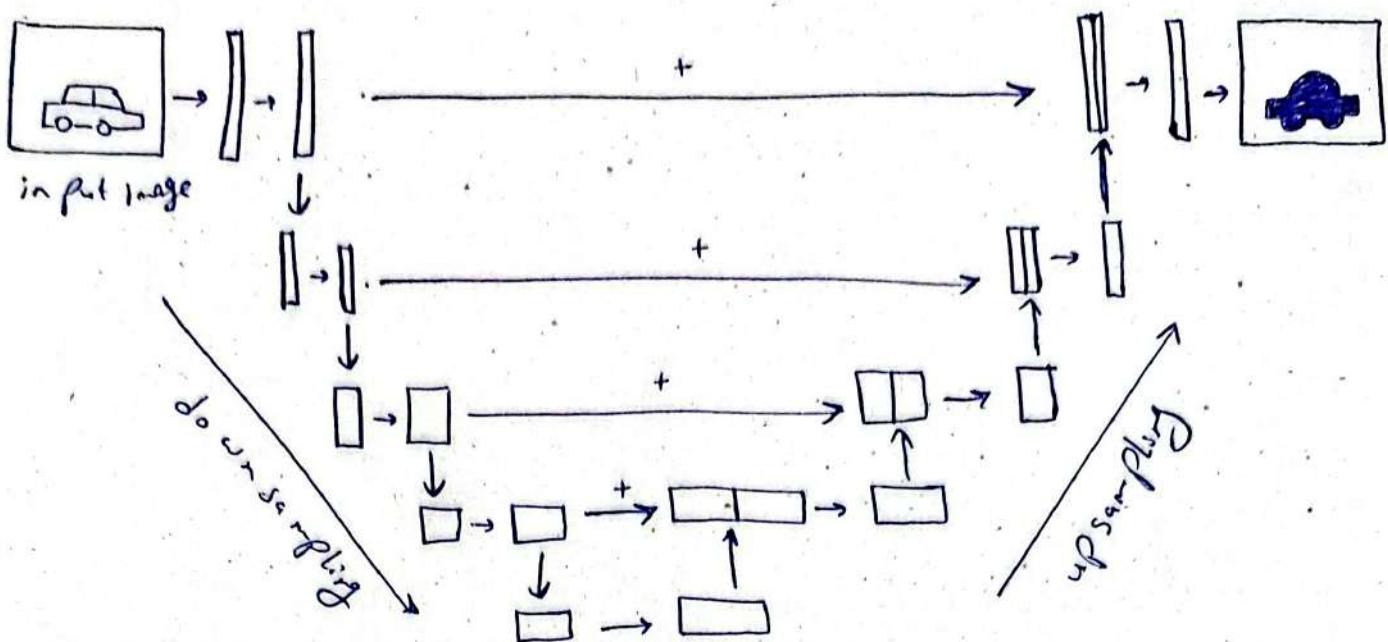
0	0	0	0	0	0
0	0	4	0	1	0
0	10	7	6	3	0
0	0	7	0	2	0
0	6	3	4	2	0
0	0	0	0	0	0

There are multiple other ways to take a small input and turn it to a big output so other ways like the

transpose convolutions can be used but transpose convolutions is just as effective, if not even more.

→ Next we'll use the idea of transpose convolutions in implementing the UNet which is a segmentation CNN based neural network.

## U-Net



- The U-Net is an image segmentation CNN that utilizes transpose convolutions to upscale the output volume back to the input volume shape.
- The U-Net also uses the concept of ResNets, adding each output volume in down sampling to the corresponding upsampled output volume with the same shape; this makes it more robust to problems like exploding and vanishing gradients; it also helps the model to learn the identity function of the inputs.
- It's called a U-Net because it looks like the letter U.
- When the U-Net was first implemented it was intended to be an application that helps with segmenting medical images but it turned out to be useful for a lot of computer vision problems.

## Week 4

### Face Recognition-

→ Face Recognition is the problem of ~~regal~~ recognizing faces and not just checking if a face matches an image; it's a technology that is potentially capable of matching a human face from a digital image or a video frame against a "database" of faces.

### Face verification vs face recognition-

Fv →

↳ input image, name / ID

↳ output whether the input image is that of the claimed person.

FR →

↳ Has a database of k persons

↳ Gets an input image

↳ output ID of the image is any of the k persons.

\* Face ~~reg~~ recognition is much harder than face verification because if you have for example a face verification model that has 99% accuracy for matching an image to a single face; it would do much worse in matching a face to a database of images of a 100 people.

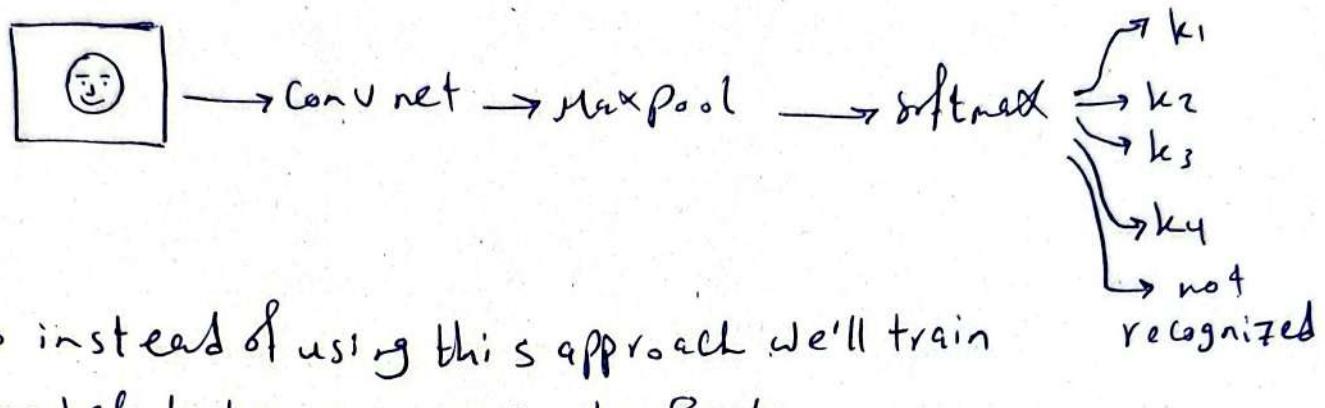
→ one of the challenges of solving the face recognition problem is solving the "one shot learning" problem.

## one shot learning

→ one shot Learning is being able to recognize a face with just one or only few images as training data.

→ Learning from one example to recognize a person again.

\* passing an image of the face to a convnet and then to a softmax function to recognize if it's a face in the database or not is not a good solution because with just 1 image as a training example the model can't learn.



\* so instead of using this approach we'll train a model to learn a similarity function.

$$d(\text{image1}, \text{image2}) = \text{degree of difference between images}$$

$\downarrow$   
difference

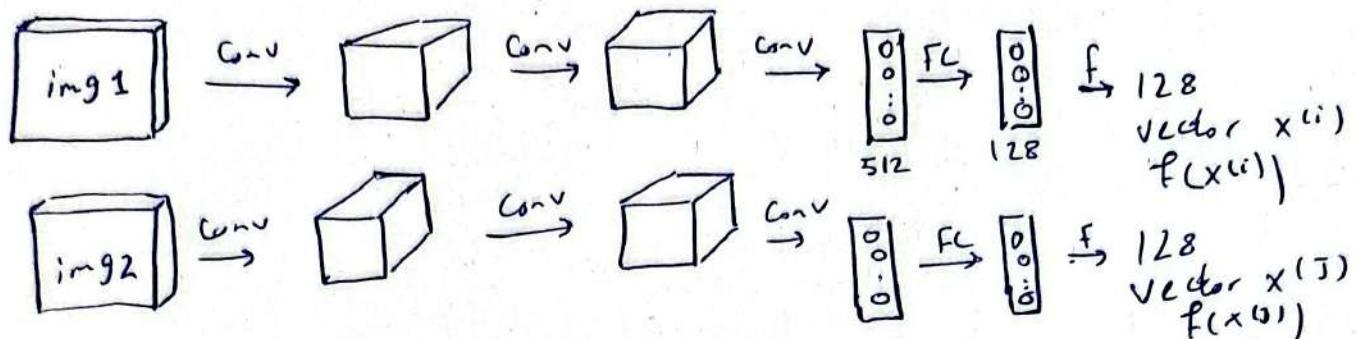
→ if the images are of the same person then you want  $d$  to output a small number (small difference) and if the images are of different people then you want  $d$  to be a large number.

---

if  $d \leq T \rightarrow \text{threshold} \rightarrow \text{"same person"}$   
if  $d \geq T \rightarrow \text{"different people"}$ .

## Siamese network:-

A siamese network take 2 vector representations of 2 images and calculates the difference between these representations. If the difference is large then they are not of the same person and if small then they probably are of the same person.



These 2 - 128 vectors are an encoding of the images; if you believe that this encoding is a good representation of the images then you can just calculate the difference between them using some function.

$$d(x^{(i)}, x^{(j)}) = \|f(x^{(i)}) - f(x^{(j)})\|_2^2$$

Small Large  
 ↓ ↑  
 "same person" "different people"

## Triplet Loss :-

→ In the terminology of the triplet loss function what you're going to do is look at 1 anchor image and a positive image which is an image of the person in the anchor image and again look at the same anchor image but with a negative image which is an image of a different person.

↳ you want to measure the distance between their vector representations so you can have

$$\|f(A) - f(P)\|_2^2 \leq \|f(A) - f(N)\|_2^2$$

So you want:-

$$\underbrace{\|f(A) - f(P)\|_2^2}_{d(A, P)} - \underbrace{\|f(A) - f(N)\|_2^2}_{d(A, N)} \leq 0$$

↓  
Anchor → positive      ↓  
Anchor → negative

But since  $0 - 0 = 0$  referring to the differences between the vector representations; the neural network can learn to set everything to zero to satisfy the loss function so instead of using  $\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 \leq 0$  we will set it to  $\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 \leq -\alpha$  where alpha is a hyperparameter that we can set; this small change prevents the network from setting all the parameters to zero and by convention the new loss is:

$$\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 + \alpha \leq 0$$

→ alpha here is also called a margin which is derived from support vector machines.

---

Loss function :-

$\downarrow$   
Same  
Given 3 images  $\rightarrow A, P, N$   
 $\uparrow$   
Different

→ The triplet loss function is defined on triplets of images.

$$L(A, P, N) \cdot \underbrace{(\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 + \alpha)}$$

we take the max between this  
and zero

$$\text{Loss} = \max(L(A, P, N), 0)$$

→ The effect of taking the max in the equation

$$\text{Loss} = \max(L(A, P, N), 0)$$

is that as long as the  $L(A, P, N)$  is smaller than 0 then the loss is equal to zero which is what we want but if  $L(A, P, N)$  is larger than zero then the loss is equal to it's value.

$$L(A, P, N) = \max(\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 + \alpha, 0)$$

↳ Now the overall cost function is :-

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

\* So if we have a 10k training set of 1k people then we can use the training data to create triplets of  $(A, P, N)$  and then train our algorithm using Gradient descent on the loss function that we defined.

Choosing the triplets  $(A, P, N)$ .

One of the problems of this approach is if  $A, P, N$  were chosen randomly given that  $A$  and  $P$  are images of the same person then the condition that :-

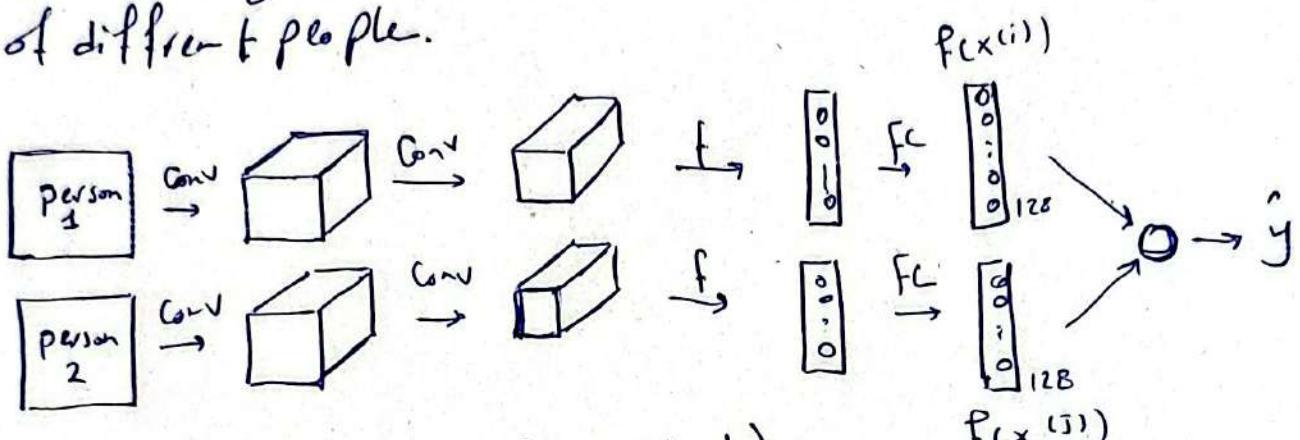
$$d(A, P) + d(A, N) + \alpha \leq 0 \rightarrow \text{is easily satisfied}$$

Because given  $A, P, N$  the difference between  $A$  and  $N$  is likely very large regarding the vector representation.

→ Choose triplets where  $d(A, P) \approx d(A, N)$  because having a small difference increases the efficiency of the algorithm.

## Face verification and Binary classification

- Another way to train a face recognition algorithm is having two neural networks embed images of same or of different people then feed these vector embeddings to a binary classification function; if the output  $\hat{y} = 1$  then the two images are of the same person & not then they are of different people.



$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

*(weights)*      *bias*

↳ so the sigmoid function classifies the difference between the different vector encodings of the images.

→ we first subtract element wise values of the vectors and then feed it to the sigmoid function and get our output  $\hat{y}$ .

$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

\* other ways of measuring the difference could be

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{2}$$

$\frac{f(x^{(i)})_k + f(x^{(j)})_k}{2}$  → This is called the chi-squared similarity or formula  $\chi^2$