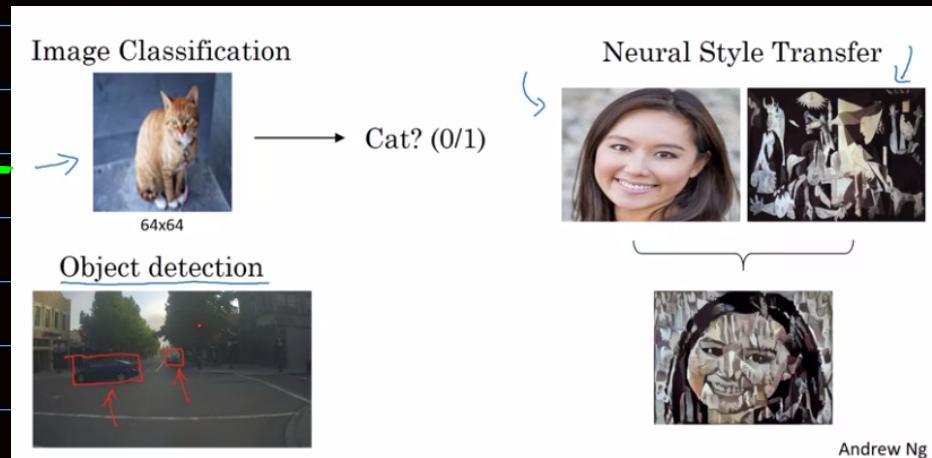


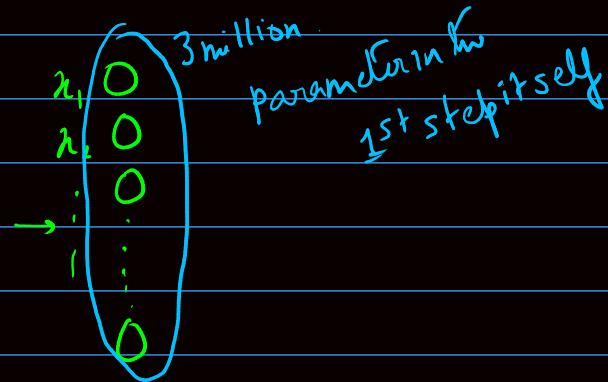
# Convolutional Neural Networks

Week 1

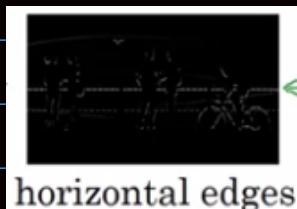
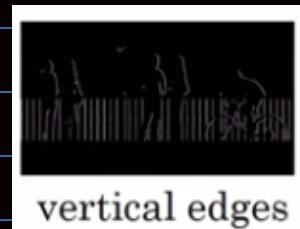
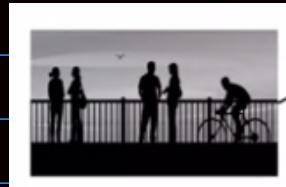
Examples : Image Classification , Object Detection , Neural Style Transfer



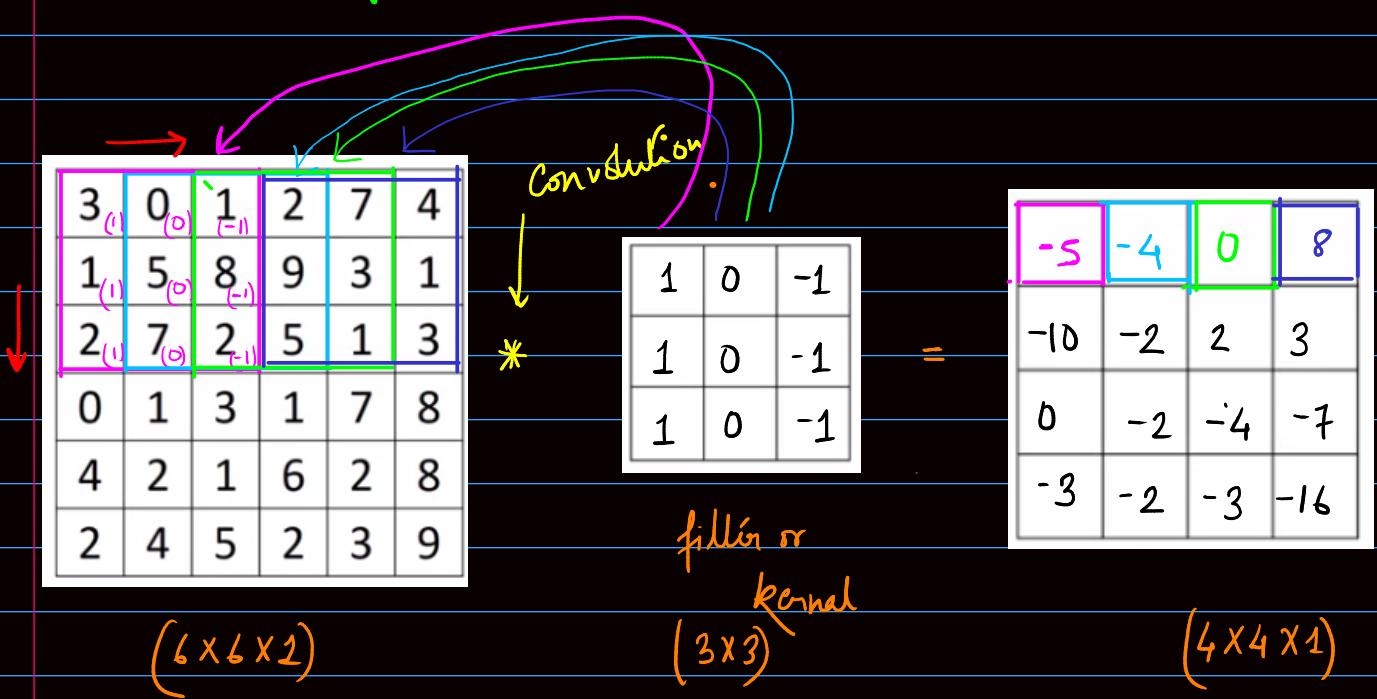
If we use  $1000 \times 1000 \times 3$  pixels then



Edge Detection



## Vertical Edge Detection



10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

(6x6)

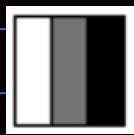
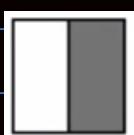
1	0	-1
1	0	-1
1	0	-1

\*

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

(3x3)

(4x4)



Since we are using very small sized image the edges seems very thick but when we use very large image the results does a pretty good job

## Verticle Edge Detection Examples.

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 \end{array} \\
 *
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} \\
 = \\
 \begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}
 \end{array}$$
  

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 \end{array} \\
 *
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} \\
 = \\
 \begin{array}{|c|c|c|c|} \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 \end{array}
 \end{array}$$

Andrew Ng

$$\begin{array}{cc}
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} &
 \begin{array}{|c|c|c|} \hline
 1 & 1 & 1 \\ \hline
 0 & 0 & 0 \\ \hline
 -1 & -1 & -1 \\ \hline
 \end{array} \\
 \text{Vertical} & \text{Horizontal}
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 & 10 \\ \hline
 \end{array} \\
 *
 \begin{array}{|c|c|c|} \hline
 1 & 1 & 1 \\ \hline
 0 & 0 & 0 \\ \hline
 -1 & -1 & -1 \\ \hline
 \end{array} \\
 = \\
 \begin{array}{|c|c|c|c|} \hline
 0 & 0 & 0 & 0 \\ \hline
 30 & 10 & -10 & -30 \\ \hline
 30 & 10 & -10 & -30 \\ \hline
 0 & 0 & 0 & 0 \\ \hline
 \end{array}
 \end{array}$$

→ Other filters

$$\begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 2 & 0 & -2 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}$$

Sobel filter

$$\begin{array}{|c|c|c|} \hline
 3 & 0 & -3 \\ \hline
 10 & 0 & -10 \\ \hline
 3 & 0 & -3 \\ \hline
 \end{array}$$

Scharr filter

But in neural networks, kernel filters are not handpicked, instead they are treated as parameters in a network which is automatically figured out by the network.

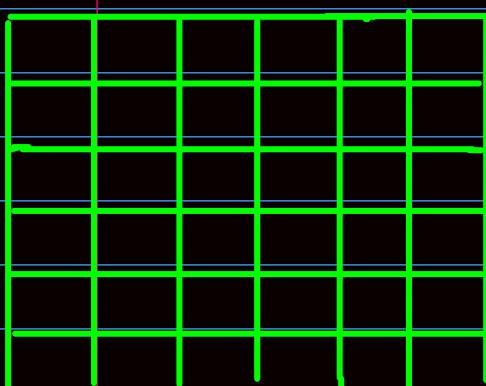
$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

With this apart from vertical & horizontal parameters we can detect  $45^\circ$ ,  $70^\circ$ ,  $73^\circ$ ,  $60^\circ$  edges as well even more robustly.

### PADDING

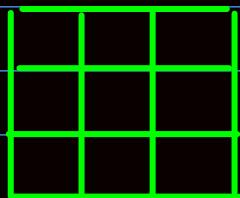
✓ with Padding

$$(n + 2p - f + 1) \times (n + 2p - f + 1)$$



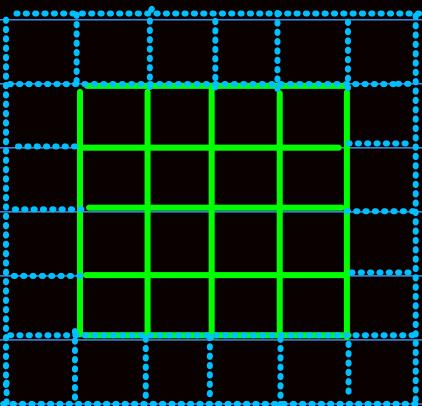
$(n \times n)$

\*



$(f \times f)$

=



$(n - f + 1) \times (n - f + 1)$

green block only ↑

\* The off size will go on shrinking as we apply more filters after the other.

\* The edge pixels will be used much lesser during computation than those in the middle.

\* Padding is used (.....) in the figure to get the same sized.

image & pad with 0 values (generally)  
Here Padding = 1

### Valid & Same Convolution

Valid :  $(n \times n) * (f \times f) \rightarrow (n-f+1) \times (n-f+1) \rightarrow$  No padding

Same : Pad so that the off size = ifp size

$$(n + 2p - f + 1) = n \\ p = \frac{(f-1)}{2}$$

Now

$$(n \times n) * (f \times f) = (n + 2p - f + 1) \times (n + 2p - f + 1)$$

\* f is usually odd. ( $3 \times 3, 5 \times 5, 7 \times 7, \dots$ )

### STRIDED CONVOLUTIONS

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	5	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

$7 \times 7$

\*

3	4	4
1	0	2
-1	0	3

3	4	4
1	0	2
-1	0	3


$3 \times 3$

[floor]

$3 \times 3$

$\left\lfloor \frac{n+2p-f+1}{s} \right\rfloor \times \left\lfloor \frac{n+2p-f+1}{s} \right\rfloor$

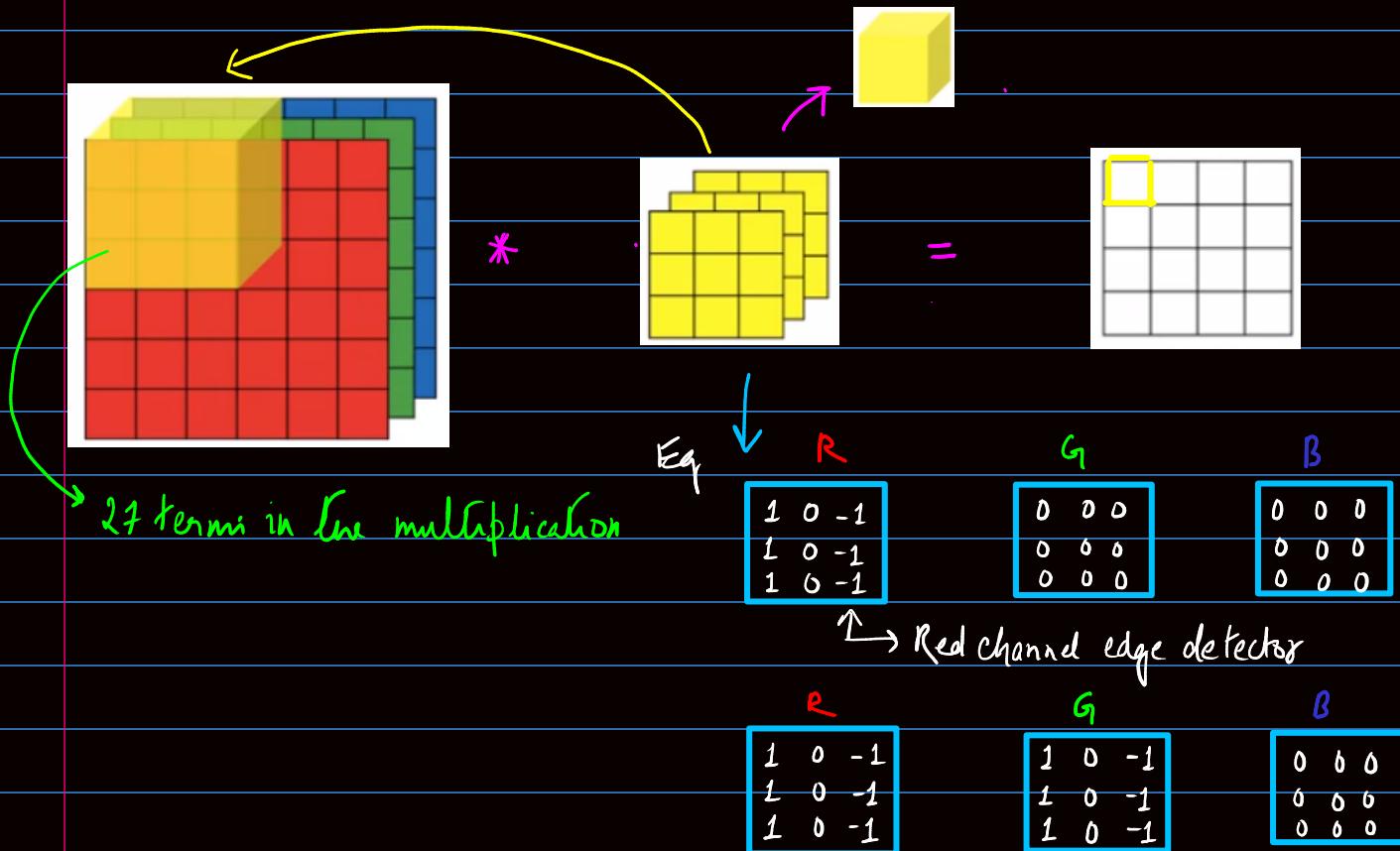
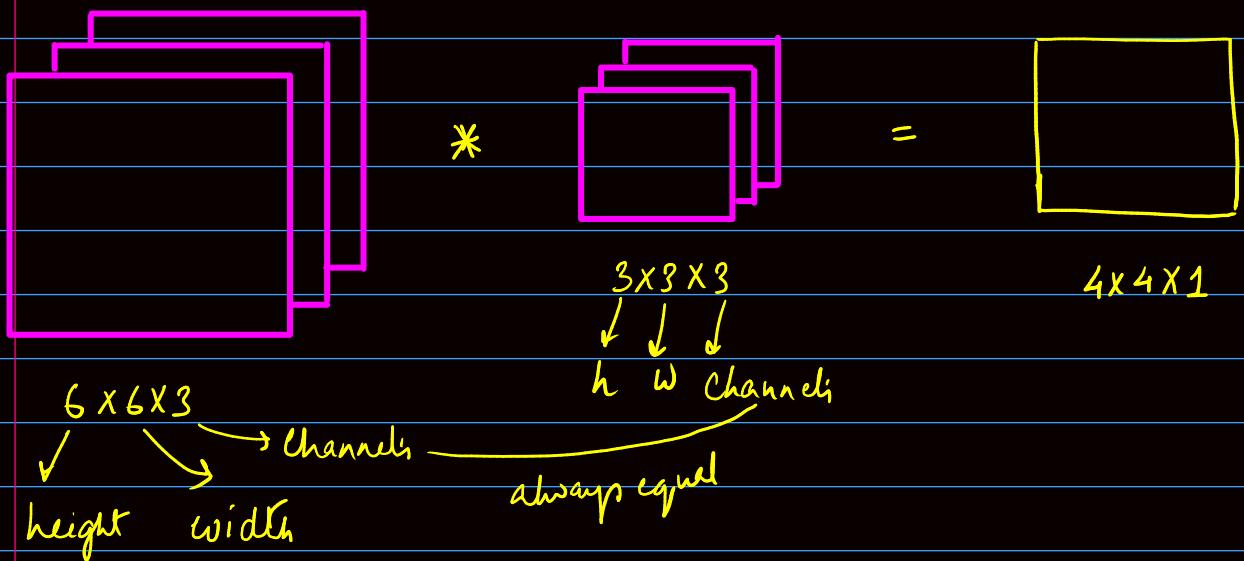
Stride = 2 [Scan by jumping 2 steps at a time]

$$\text{Eq: } \left( \frac{7-0-3}{2} + 1 = 3 \right)$$

\* Caution : Don't follow the math & SS textbook for convolution formula.

In ML actually convolution is cross-convolution.

## CONVOLUTIONS OVER VOLUMES



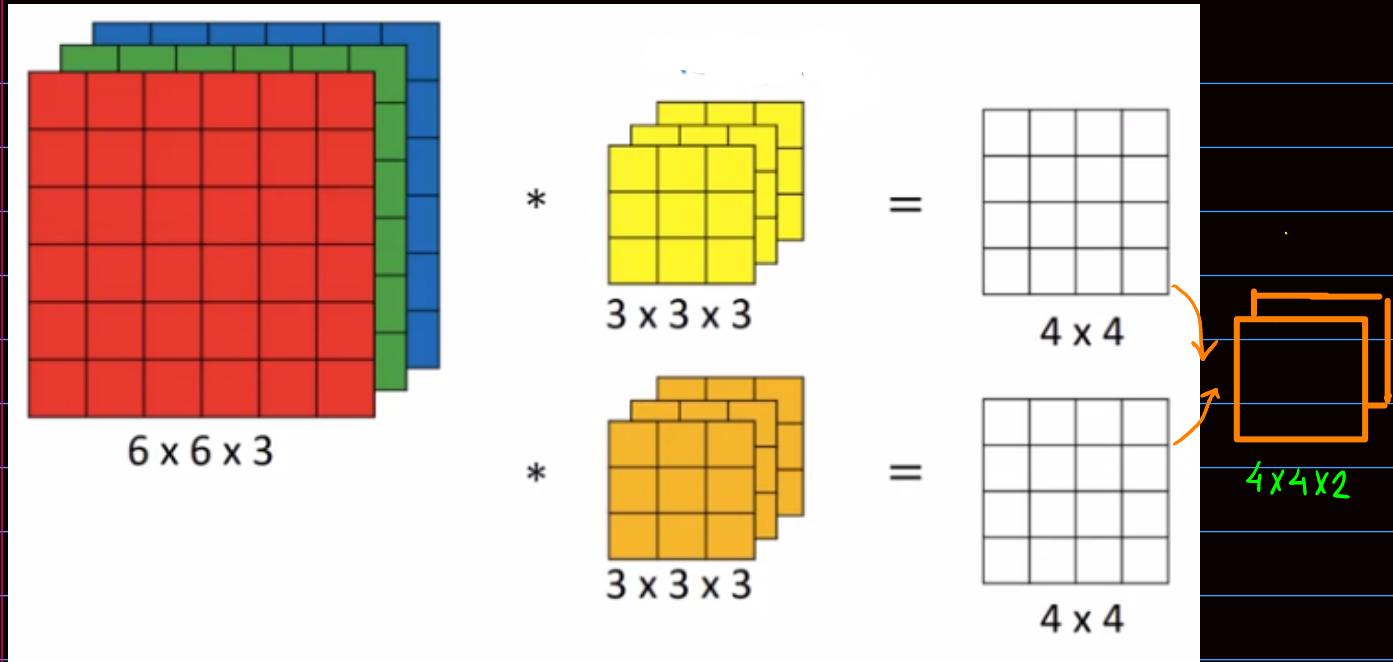
Multiple filters

filter used  
m.s.

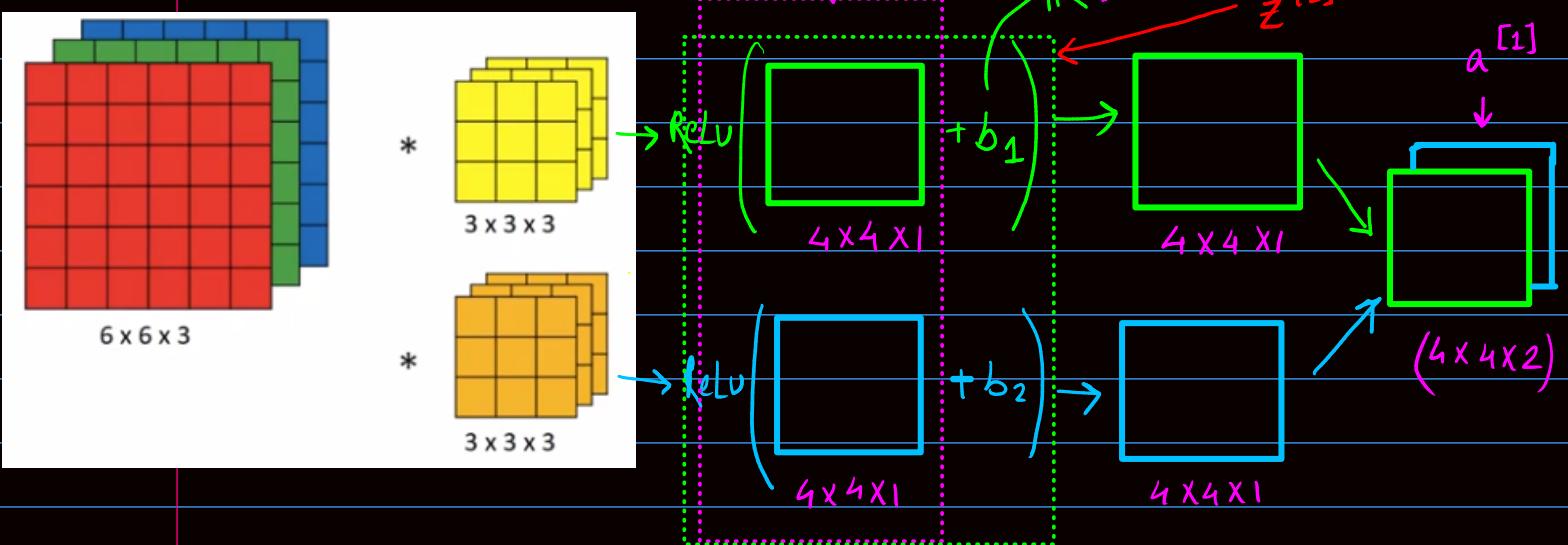
$$n \times n \times n_c \times (f \times f \times n_c) \rightarrow (n - f + 1) \times (n - f + 1) \times n'_c$$

$(f \times f \times n_c)^2$

$6 \times 6 \times 3 \quad 3 \times 3 \times 3 \times 2 \quad 4 \times 4 \times 2$



### Layer Example.



Q If you have 10 filter that are  $3 \times 3 \times 3$  in one layer of a neural network how many parameters does that layer have?

Ans:  $10 \times 3 \times 3 \times 3 + 10 \cdot \underbrace{bias}_{b_{10}} = 280$  parameters

\* Notation

- ①  $f^{[l]}$  = filter size
- ②  $b^{[l]}$  = padding

$$③ S^{[l]} = \text{Stride}$$

$$④ \text{Input} \rightarrow n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$$

$$⑤ \text{Output} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

$$\text{for } H \quad n^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{S^{[l]}} + 1 \right\rfloor$$

$$⑥ n_C^{[l]} \rightarrow \text{no. of filters}$$

$$⑦ \text{Each filter} n = f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$$

$$⑧ \text{Activation} : a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

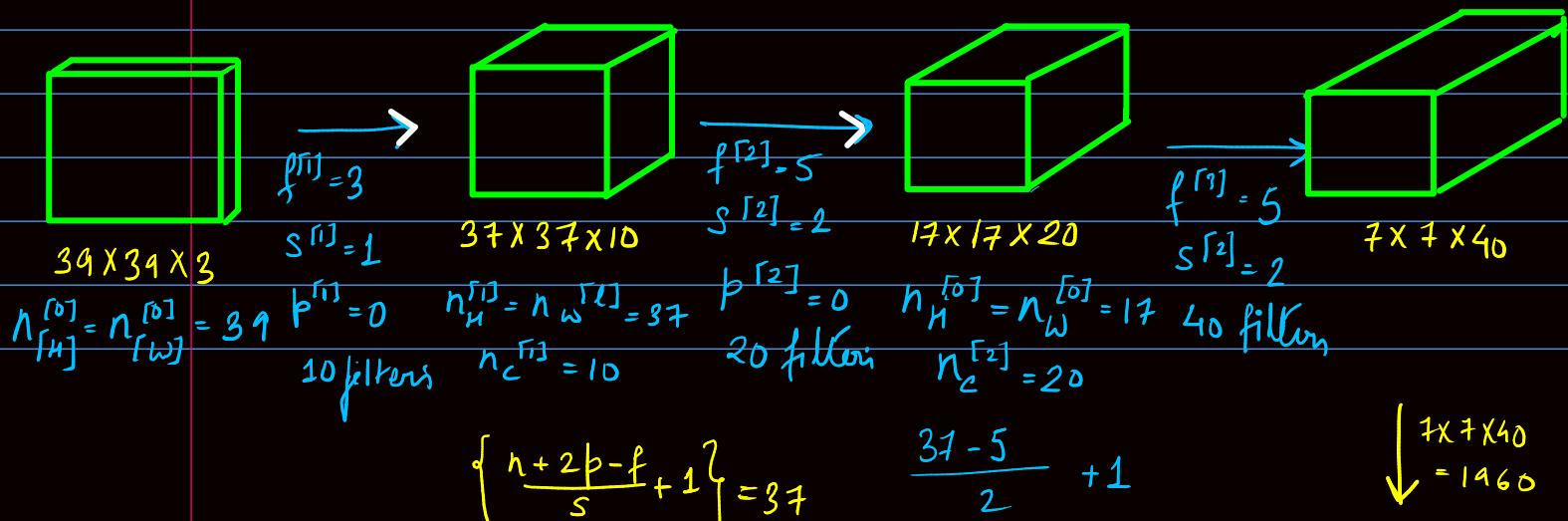
$$A = m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

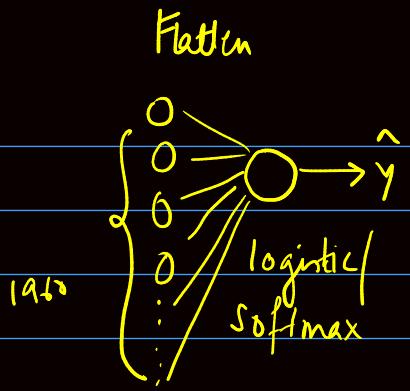
$$⑨ \text{Weights: } f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$$

↓ total no. of filters

$$⑩ \text{Bias} : n_C^{[l]} \rightarrow (1, 1, 1, n_C^{[l]})$$

\* Simple Convolutional Network example:



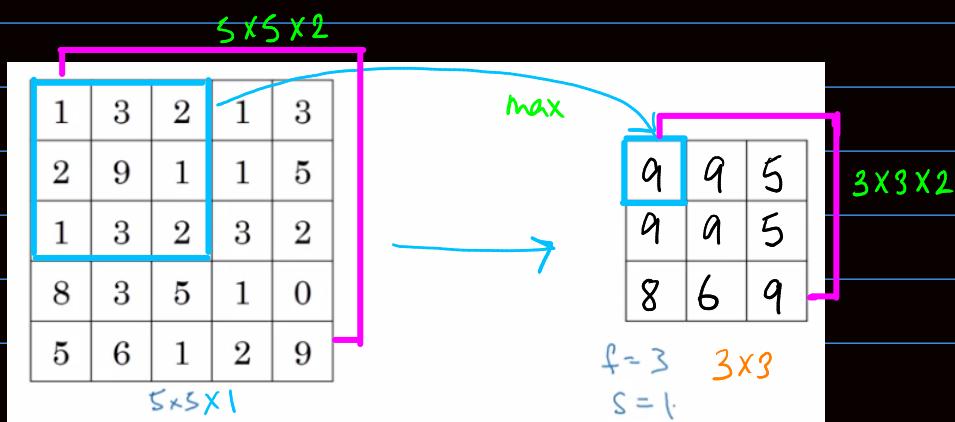
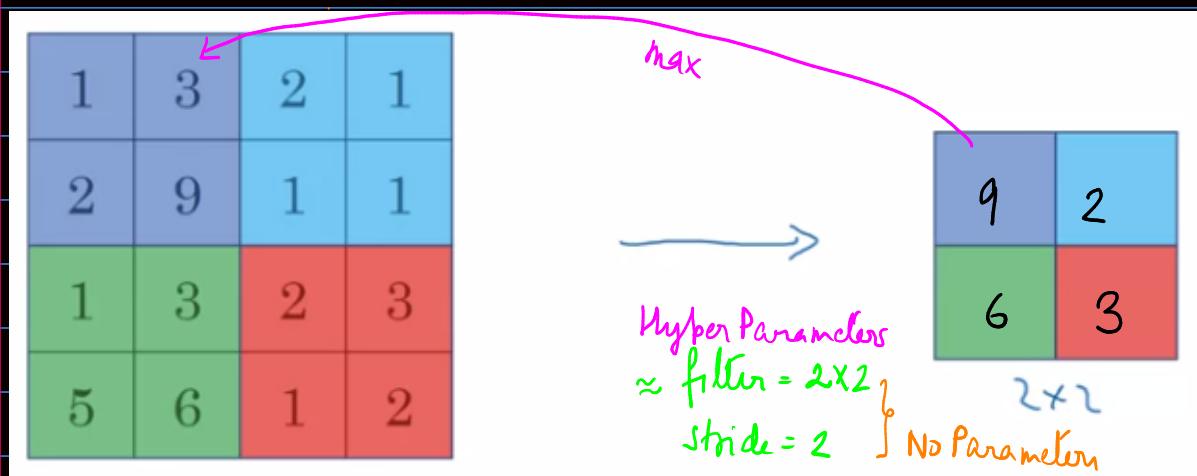


## Types of Convolutional Network

- ① Convolution (conv)
- ② Pooling (pool)
- ③ Fully Connected (FC)

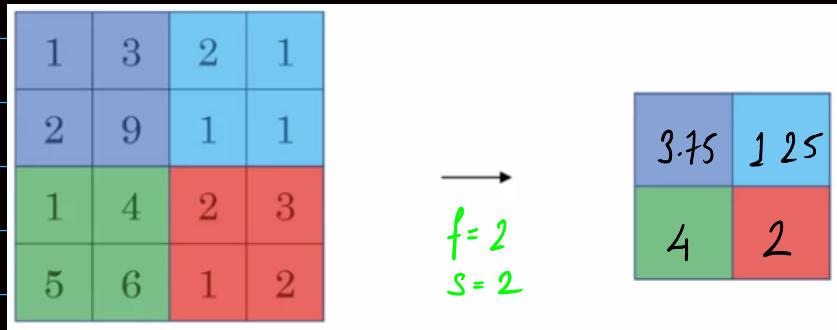
### POOLING LAYERS

#### ① Max Pooling



\* Max Pooling is done independently on each of the channels i.e  $5 \times 5 \times n_c \rightarrow 3 \times 3 \times n_c$

## ② Average Pooling



\* Max Pooling is most commonly used.

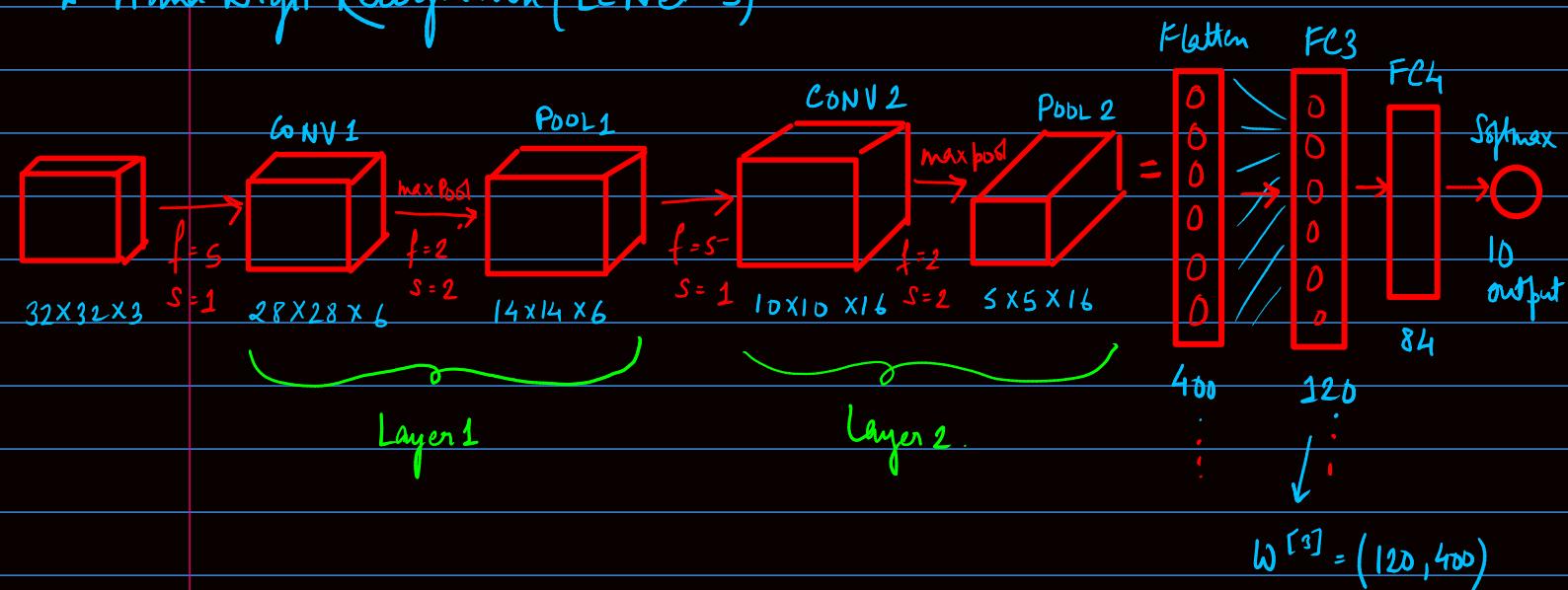
Hyperparameters Used:  $f = \text{filter size}$ . { common  $f=2, s=2$   
 $s = \text{stride}$   $f=3, s=1$

Max or Average pooling

\* No Learnable Parameters in Pooling Layer

## CNN Example

\* Hand Digit Recognition (LeNet-5)



\* As we get deeper  $n_h \& n_w \downarrow$  &  $n_c \uparrow$

\* CONV - POOL - CONV - POOL - FC - FC - SOFTMAX { Common Pattern }

## On the basis of the LeNet:

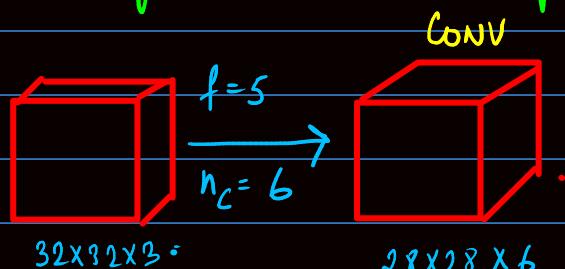
	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	— 3,072 $\alpha^{[3]}$	0
CONV1 (f=5, s=1)	(28,28,8)	6,272	608
POOL1	(14,14,8)	1,568	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	3216
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48120
FC4	(84,1)	84	10164
Softmax	(10,1)	10	850

Inference: ① Pooling layers have no parameters

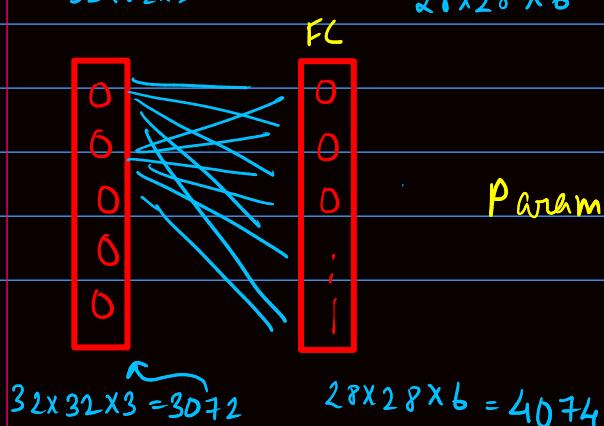
② CONV layers have very less parameters as compared to FC layers

③ Activation size goes down gradually  
 { If it drops very quickly it is not a great performance

Q Why convolutions are useful?



$$\text{Parameters} = 5 \times 5 \times 6 + 1 = 156, \text{ Very Small}$$



$$\begin{aligned} \text{Parameters} &= 3072 \times 4094 \\ &= 14 \text{ million in one layer and a low pixel image} \end{aligned}$$

∴ Convolutional layers has much lesser parameters

Reason:

$$\begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array} * \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} = \begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}$$

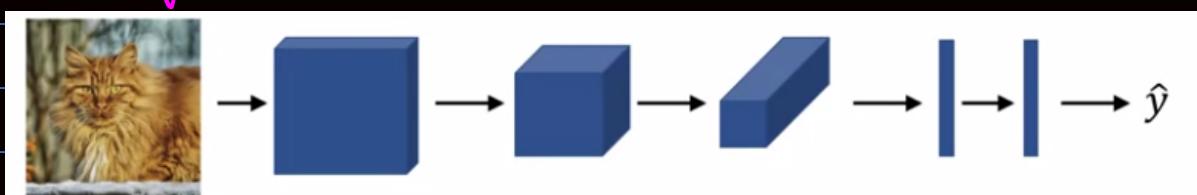
i.e same filter is used for all the time all around the image ∴ A parameter are shared

Parameter Sharing: A feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image

Sparsity of Connections: In each layer, each output value depends on small no. of input's

Translation Invariance is captured very well in CNN.

Training Set:  $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

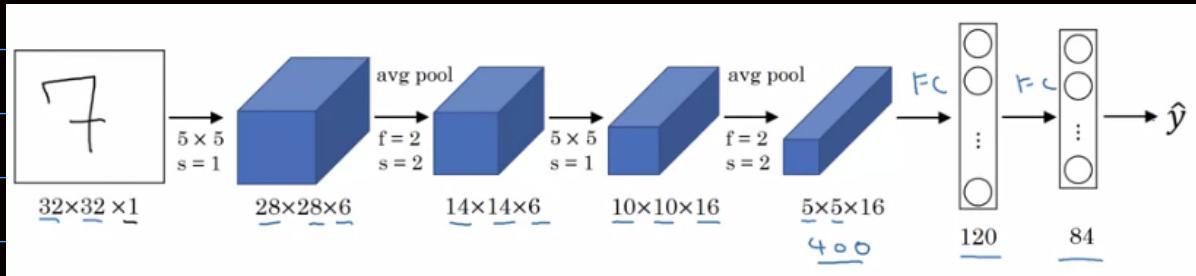


$$Cost J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

\* Use gradient descent to optimize J.

## Week 2

### ① LeNet-5



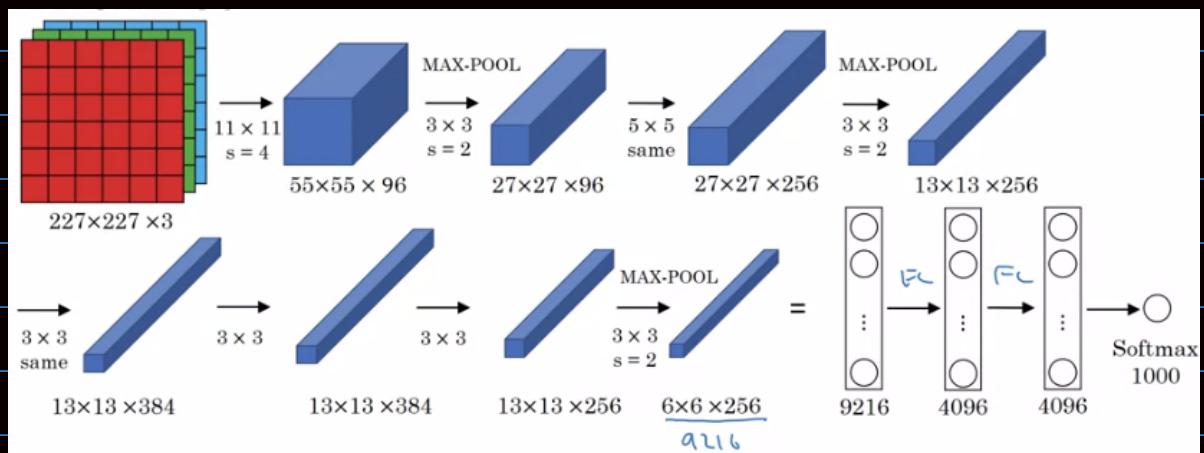
\* 60K Parameters

\*  $n_H, n_W \downarrow$

\* CONV - POOL - CONV - POOL - FC - FC - OUTPUT

{ Original Paper → used sigmoid / tanh  
→ used non-linearity after pooling }

### ② ALEX Net.



\* Around 60 M Parameters

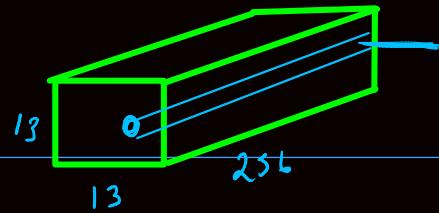
\* Used multiple GPU's

\* ReLU

\* Used Local Response Normalization Layer (LRN)

\* Not used in practise

## LRN

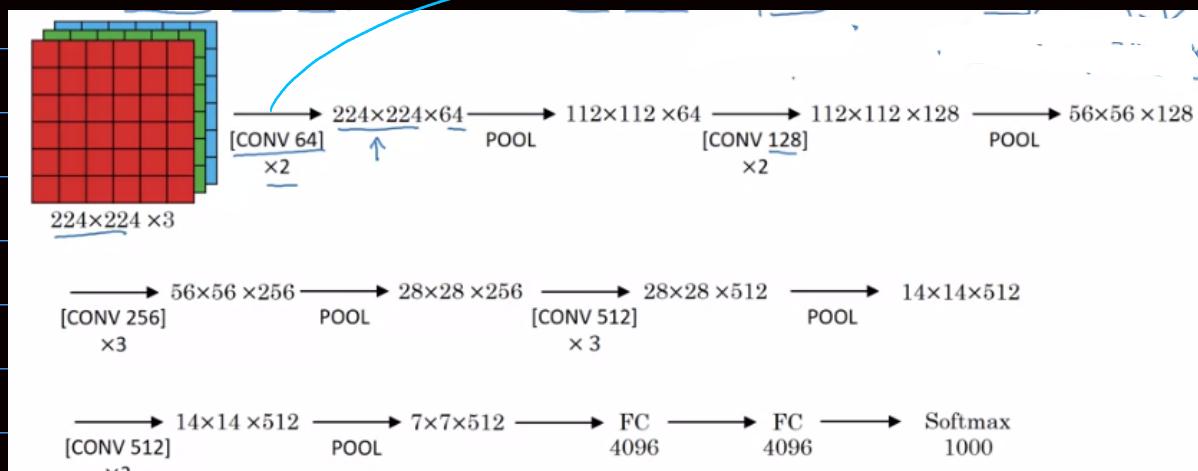


Normalize a position in height & weight  
face & normalize across the  
depth

## ③ VGG-16 (16 steps)

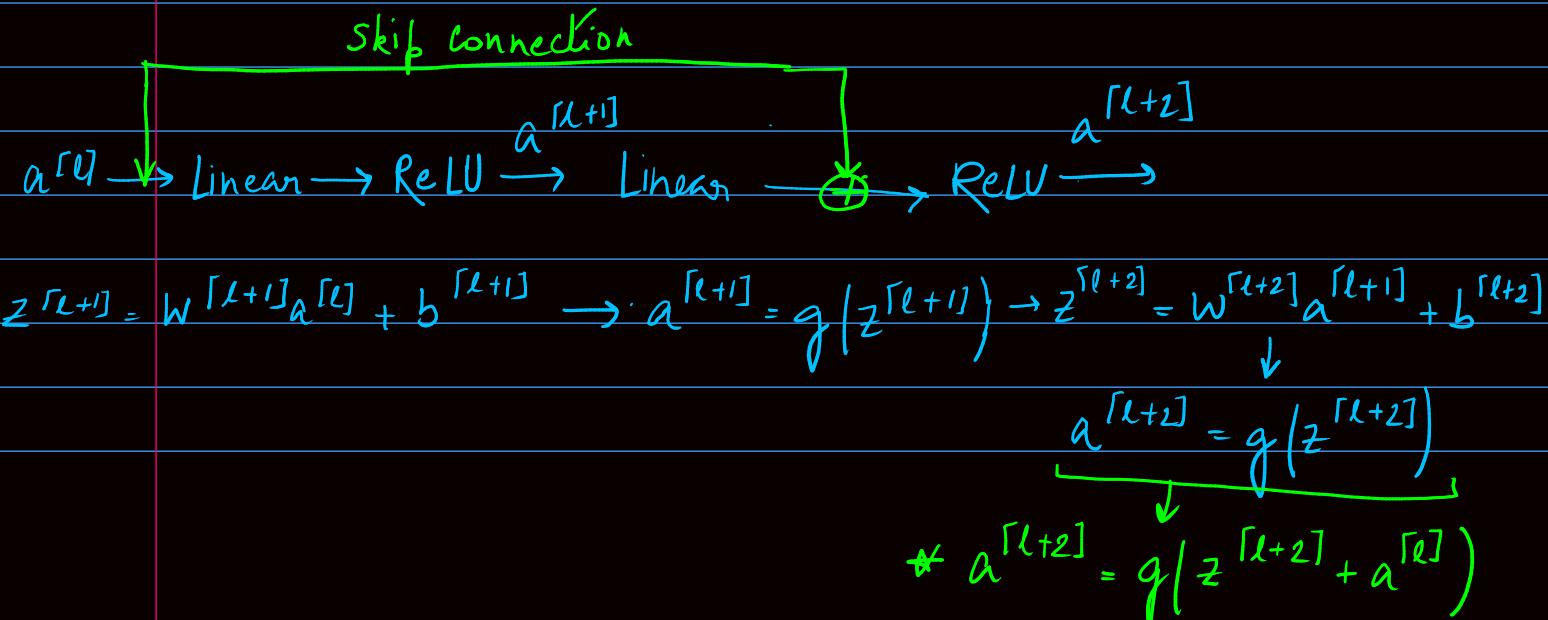
- \* Conv =  $3 \times 3$  filter,  $s=1$ , same.
- \* Max Pool =  $2 \times 2$ ,  $s=2$

Means 2 conv 2D layers with 64 filters

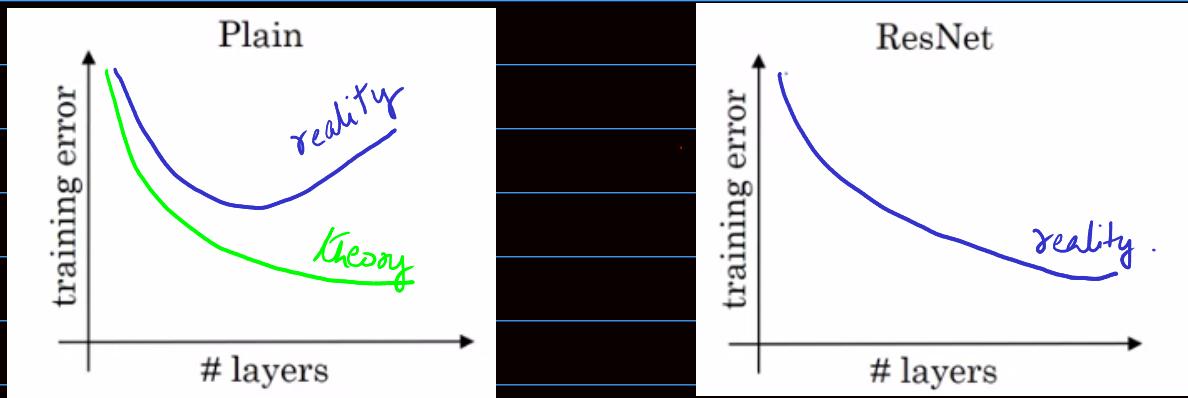
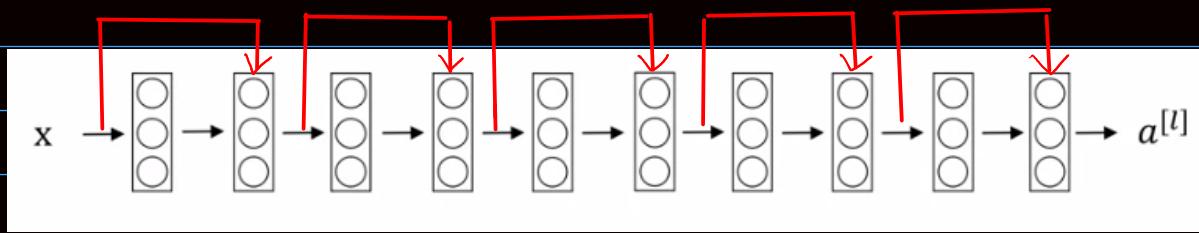


- \* 138 M Parameters
- \*  $n_h, n_w \downarrow$  but  $n_c \uparrow$

## RESNETS

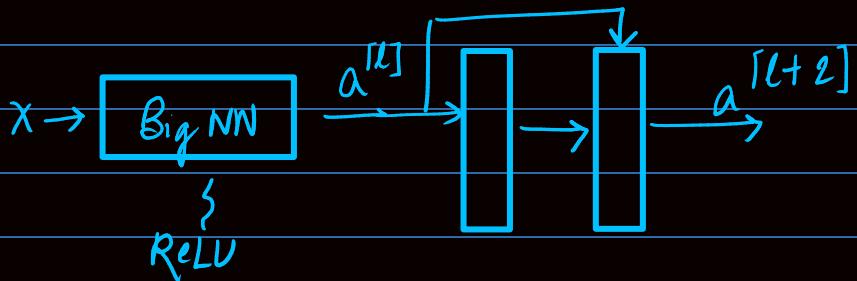


\* It is useful when we build very very deep network.



Why RESNETS work so well?

$$x \rightarrow \boxed{\text{Big NN}} \rightarrow a^{[l]}$$



$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

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

*(After some cut decay)*

$$\text{if } w^{[l+2]} = 0 = b^{[l+2]}$$

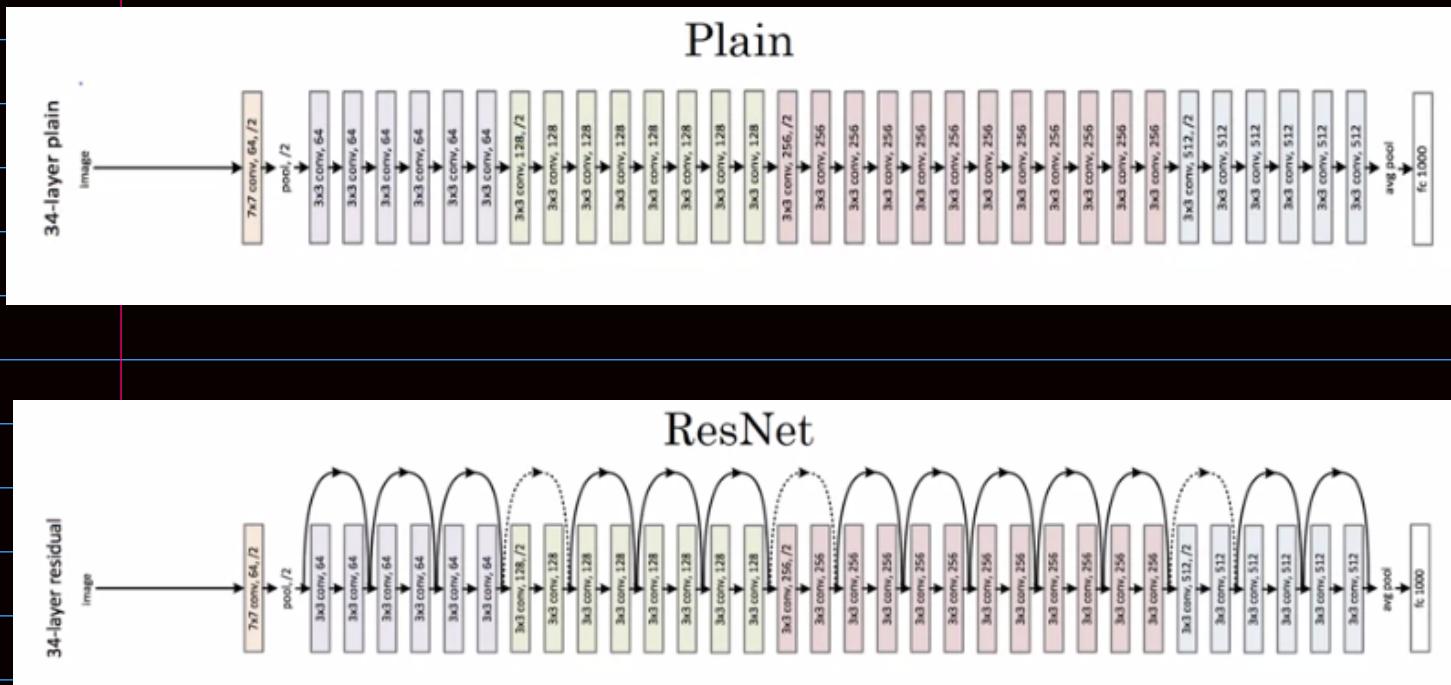
Identity fun<sup>n</sup> → easy for R-Block to learn

- \* Adding the 2 extra FC layers doesn't hurt the network at all
- \* Here dimension of  $Z^{[L+2]} = a^{[L]}$ . Thus in all of Resnet we prefer using same padding

or using same padding  
or

$$a^{[l+2]} = g\left(z^{[l+2]} + w_s a^{[l+2]}\right)$$

$256 \times 128$



## Network in Networks, & 1x1 convolutions.

1	2	3	6	5	8
3	5	5	1	3	4
2	1	3	4	9	3
4	7	8	5	7	9
1	5	3	7	4	8

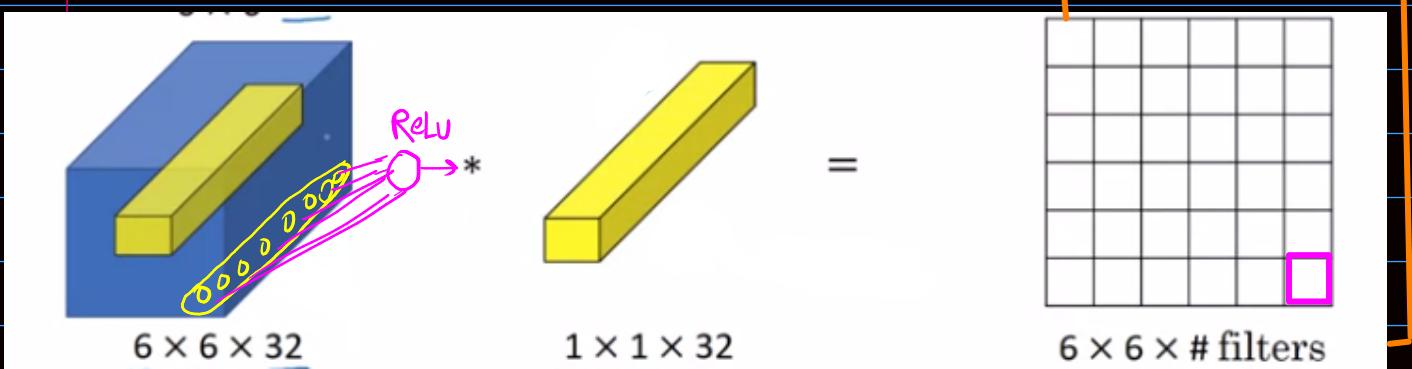
$6 \times 6 \times 1$

\*

2

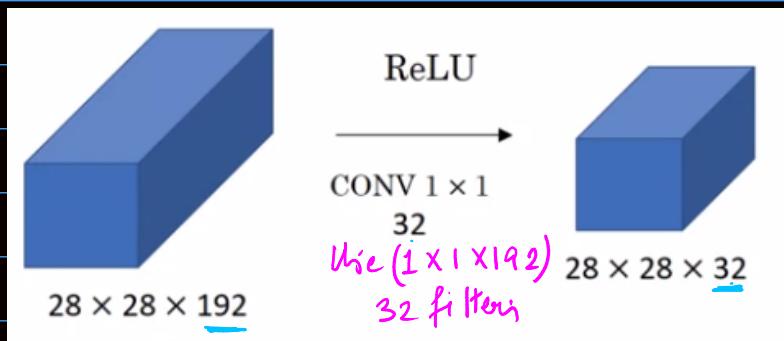
2

2	4	6	12.



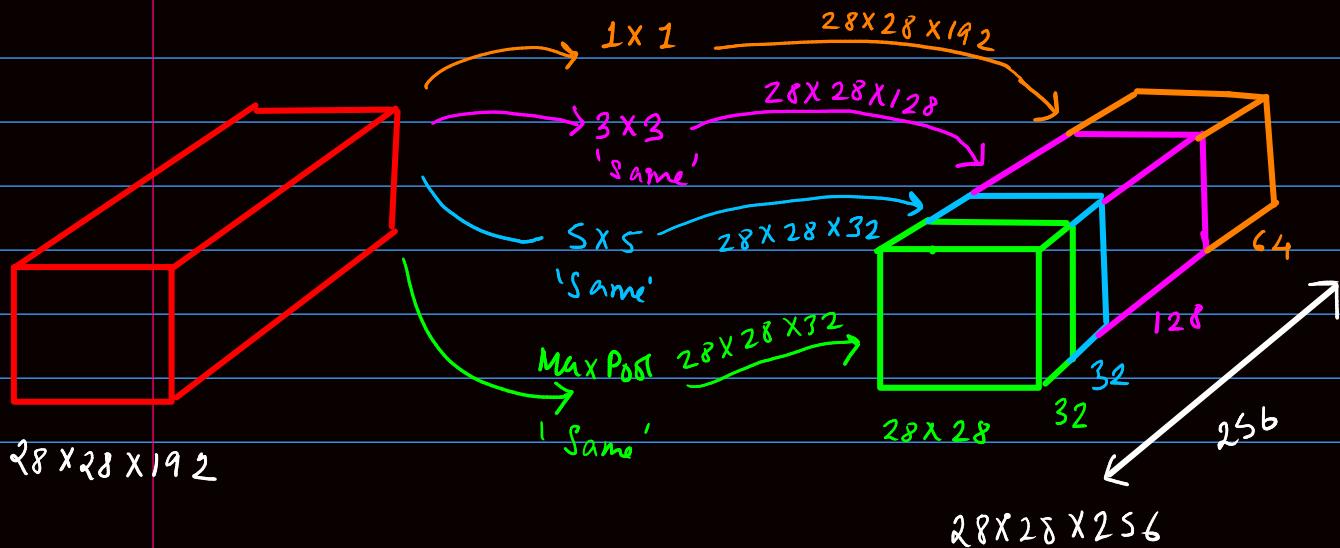
This is called 1x1 convolution or Network in Networks

Example :



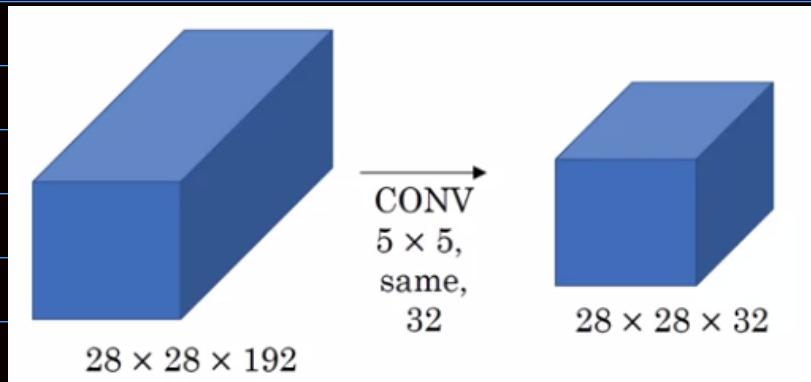
- \* Helps to shrink the filter size effectively & also save on computation
- \* It also adds non-linearity (ReLU) to the n/w

### Inception Network



Hence instead of picking one filter sizes or pool we want we can pick from all  $g_i$  concatenate the outputs  $g_i$ , let the network learn what parameters it wants to use or which filter sizes it wants

\* Problem with computational cost

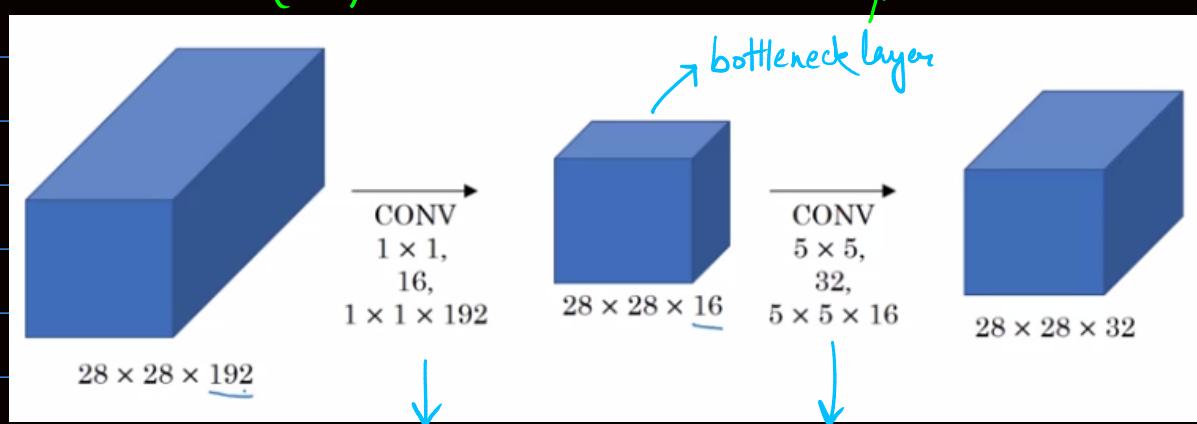


Parameters  $\in 5 \times 5 \times 192 \times 32$

$$\text{No. of Multiplications} = 28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120M$$

$\uparrow$   
Very Expensive Operation

\* Instead use  $(1 \times 1)$  convolution.

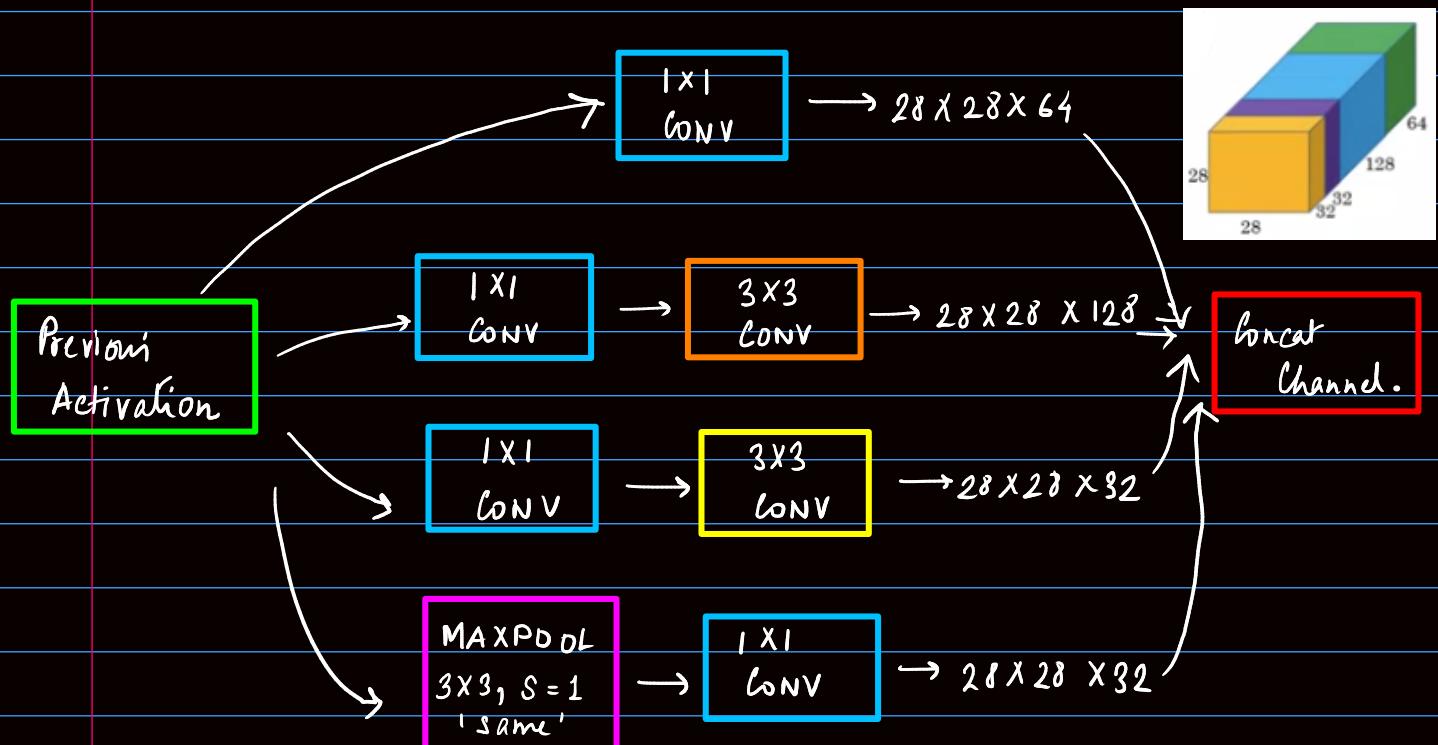


$$28 \times 28 \times 16 \times 192 = 2.4M$$

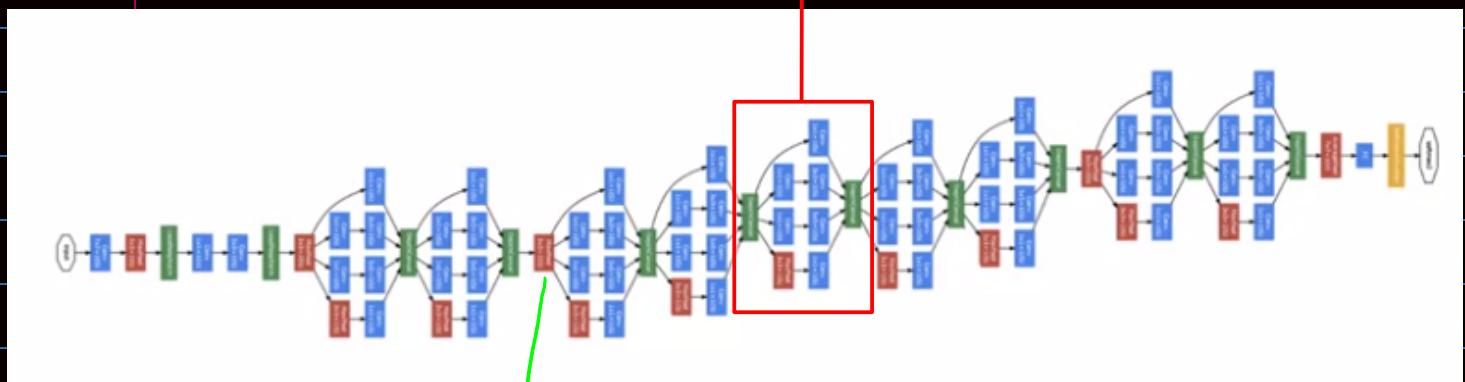
$$28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10.0M$$

$$\therefore 2.4M + 10M = 12.4M \ll 120M$$

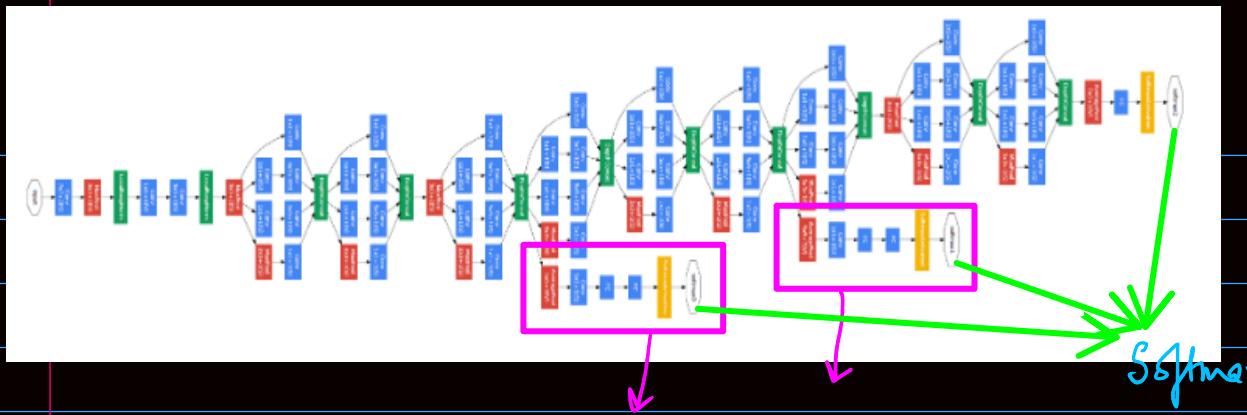
\* Shrinking down the layers doesn't hurt the performance of the model but saves a lot of computation, unless we are using bottleneck layers



\* Inception Module as a whole.



Inception Network is a list of these blocks



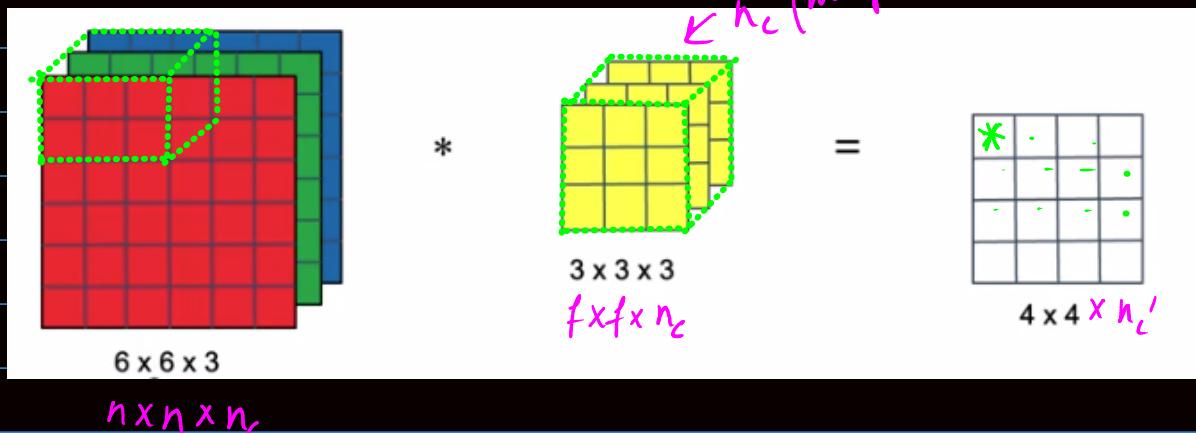
In the original paper we have these side blocks which are classifiers (consisting of a few FC layers & a softmax layer) to predict the output at different stages of the network.

\* Also called GoogleNet

### MobileNet

- \* Low computational cost at deployment
- \* Useful for mobile & embedded vision application .
- \* Normal vs Depthwise Separable Convolutions .

### Normal Convolution

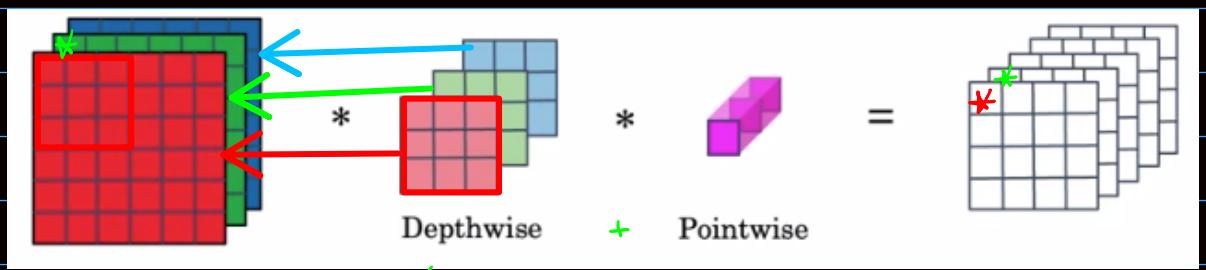


### Computational Cost

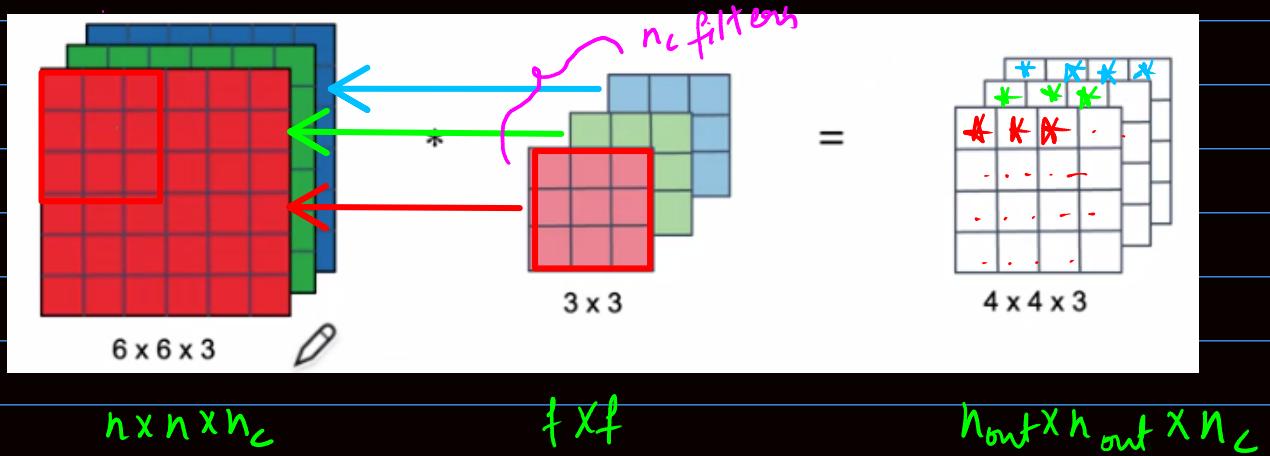
$$= \# \text{filter.params} \times \# \text{filter.positions} \times \# \text{of filters}$$

$$(3 \times 3 \times 3) \times (4 \times 4) \times 5 = 2160$$

# Depthwise Separable Convolution



## ① Depthwise Convolution



## Computational Cost

$$= \# \text{filter params} \times \# \text{filter positions} \times \# \text{of filters}$$

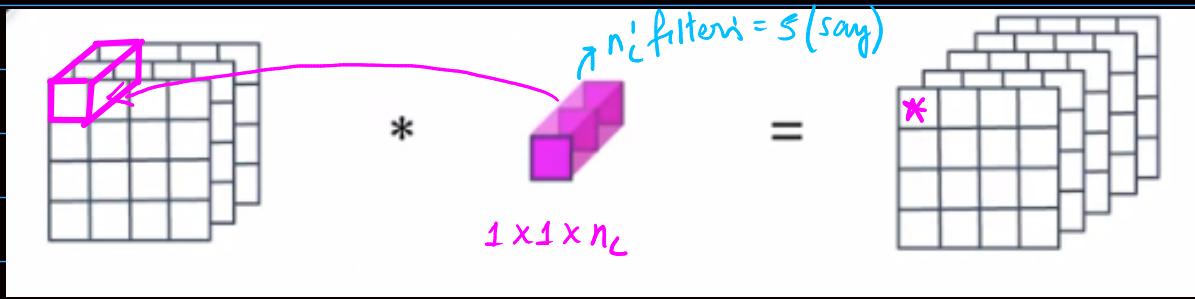
$3 \times 3$

$4 \times 4$

$3$

$= 432 - ①$

## ② Pointwise Convolution



$h_{out} \times h_{out} \times n_c$

$n_{out} \times n_{out} \times n_c'$

### Computational Cost

= #filter.params  $\times$  #filter positions  $\times$  #of filter

$$1 \times 1 \times 3 \quad \times \quad 4 \times 4 \quad \times \quad 5 \quad = \quad 240 \quad - \textcircled{1}$$

### Cost Summary

Cost of normal convolution = 2160

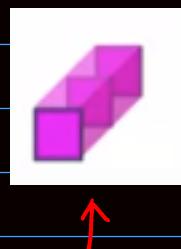
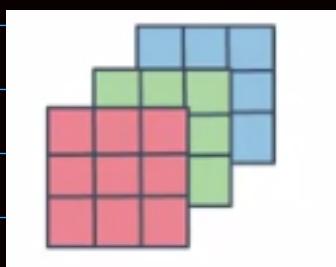
Cost of depthwise separable convolution =  $432 + 240$   
 $= 672 \ll 2160$

Generally,

$$\frac{\text{Cost of normal convolution}}{\text{Cost of depthwise separable convolution}} = \frac{1}{n_c'} + \frac{1}{f^2}$$

In reality  $\approx$  10 times cheaper

\*\* For notes purpose (Notation)

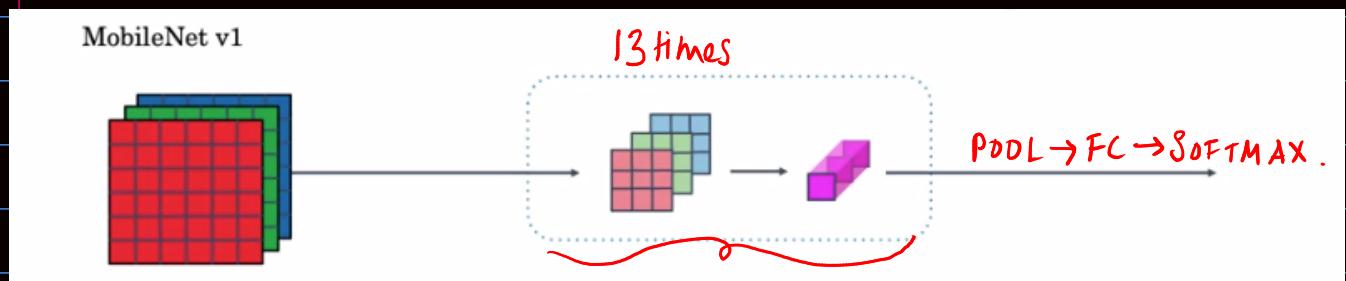


↑  
Denotes Depthwise convolution in general  
(not necessarily three channels)

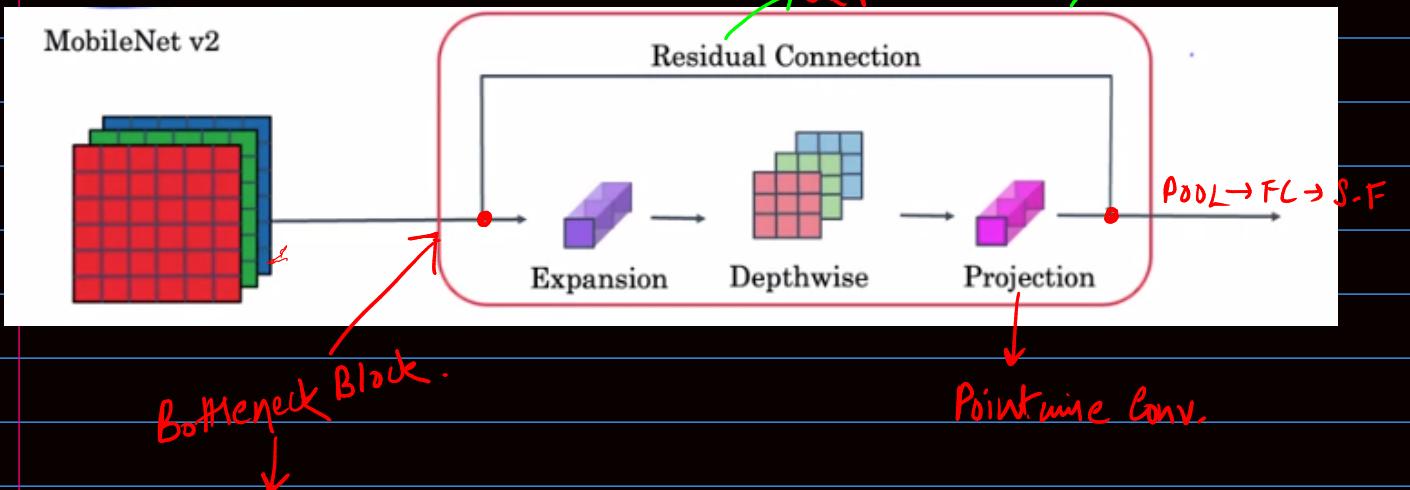
↑  
Denotes pointwise convolution in general (not necessarily 3 channels)

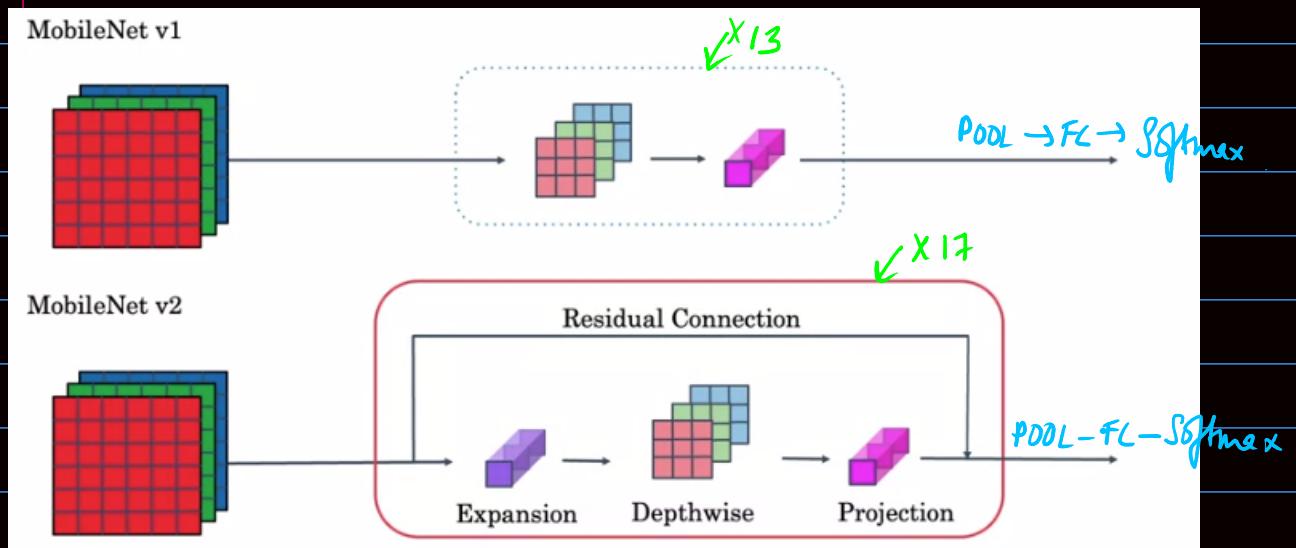
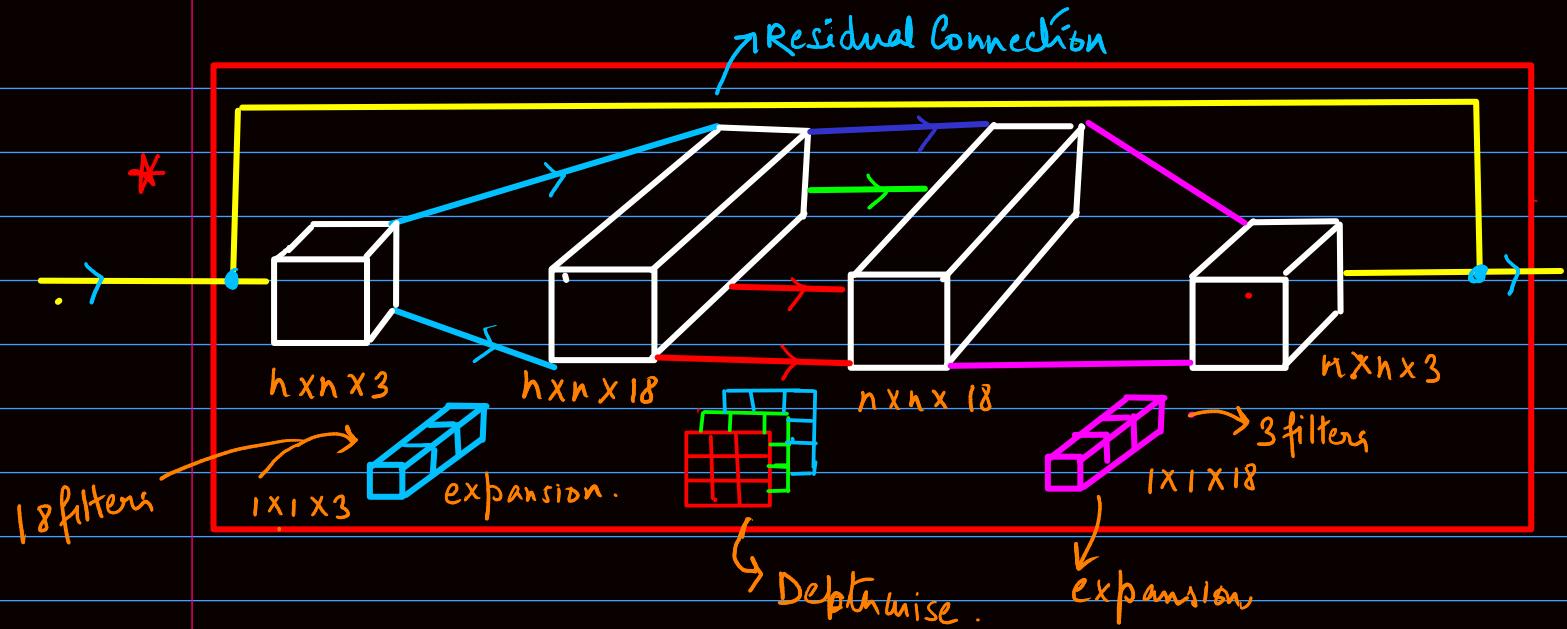
## MobileNet Architecture

①

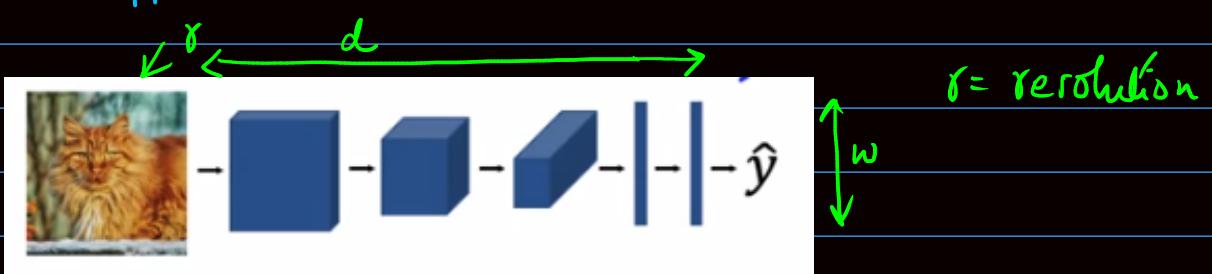


②





## Efficient Net



To scale the components up we can:

- ①  $\gamma \uparrow$
- ②  $d \uparrow$
- ③  $w \uparrow$

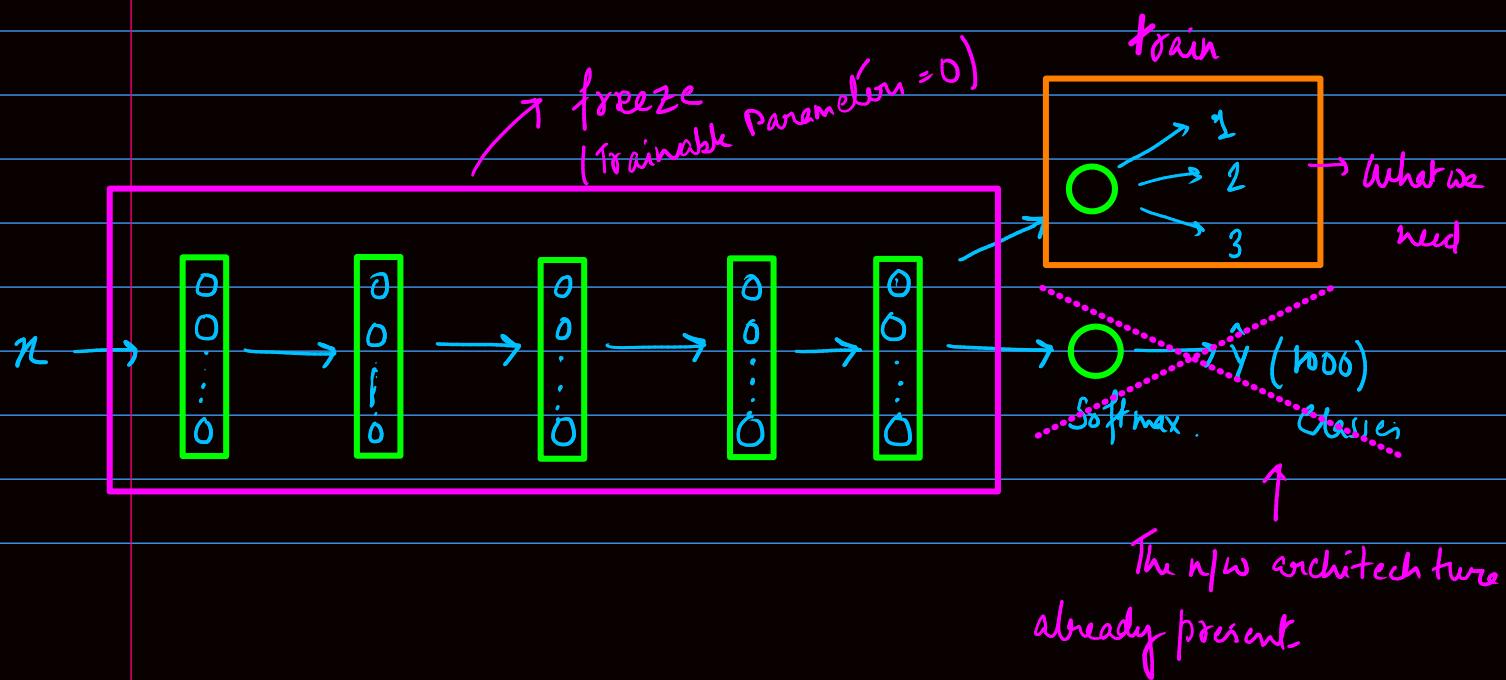
q What is the best tradeoff betn  $\gamma$ ,  $d$  &  $w$  to be most computationally efficient?

Ans: With efficient net we can scale up or down depending on resources we are working on.

### Practical Implementation

- \* Using Pretrained networks for similar tasks saves a lot of time
- \* Use Transfer learning to compensate for the small datasets  
(Transfer learning discussed in the previous course)

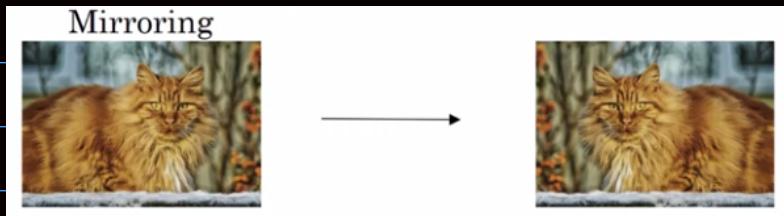
Eg.



- \* If we have a slightly significant dataset then we can try freezing few layers
  - \* dataset size ↑ , freezed layer ↓
- \* We can also download a model & just use the initialization & train the entire model from scratch.

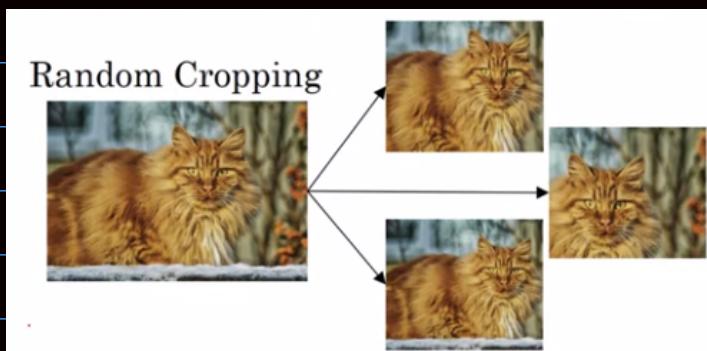
## Data Augmentation

1. Mirroring



2.

Random Cropping



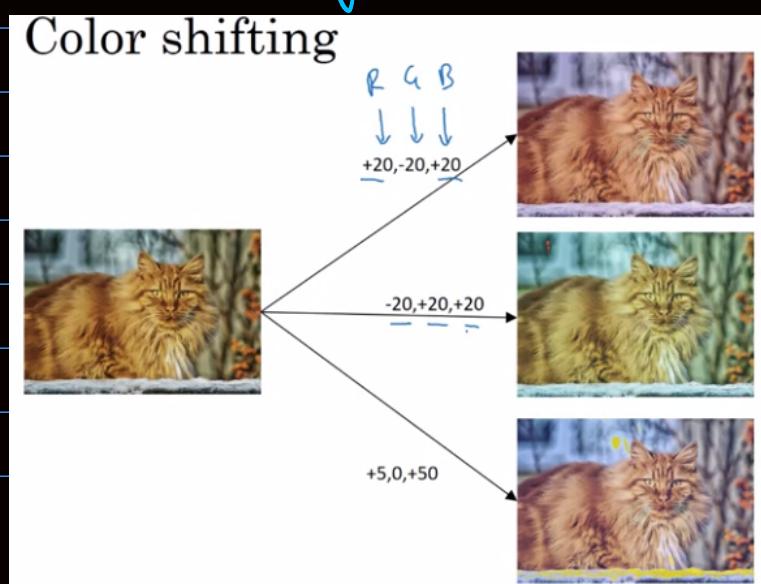
3. Rotation

4. Shearing

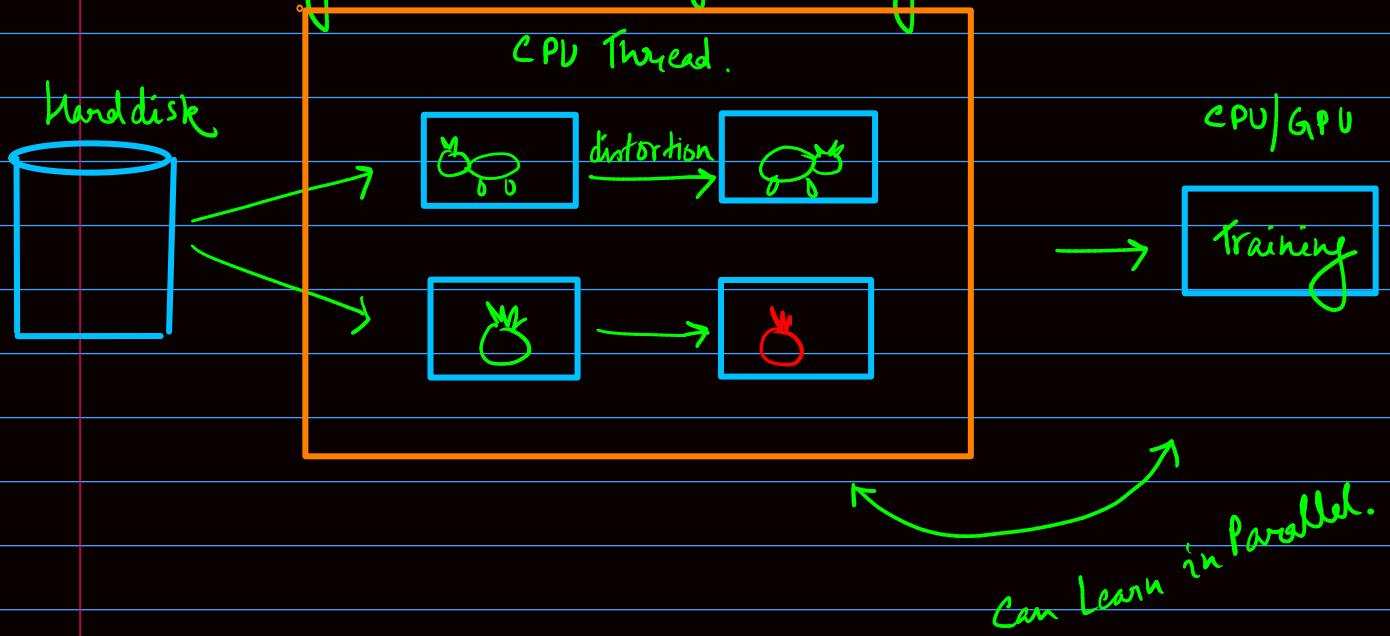
5. Local Warping

6.

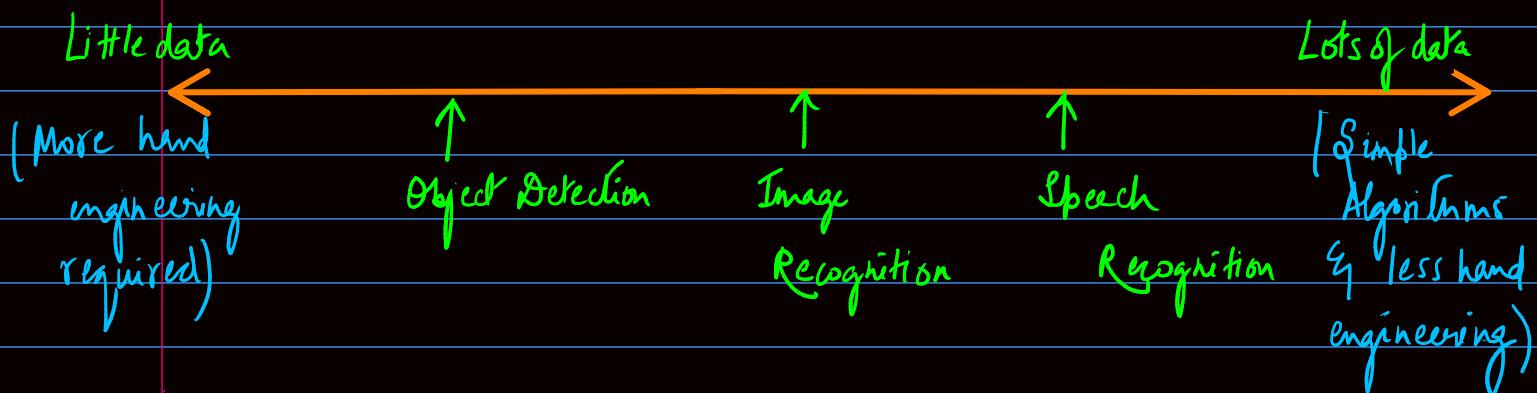
Color shifting



## \* Implementing Distortion during training



## Scale of Computer Vision



Two sources of knowledge.

- Labelled data.
- Hand engineered features/network architecture of other components

\* CV has more hand engineered data due to the problems of less data readily available.

\* But for less data we have transfer learning to rescue

Tips for doing well on benchmarks | winning competitions

→ Ensembling

- Train several networks independently & average their opp's.  
(3-15 n/w's)

→ Not preferred in production.

→ Multicrop at test time

- Run classifier on multiple versions of test images & average results



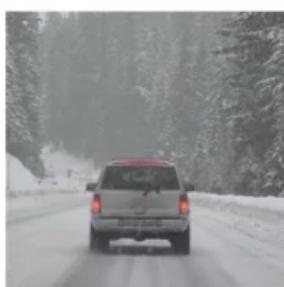
→ Not preferred in production

- \* Use architecture of networks published in literature
- \* Use open source implementation if possible.
- \* Use pretrained models & finetune on your dataset

Week 3:

## Object Localization

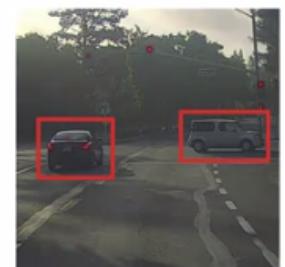
Image classification



Classification with localization



Detection

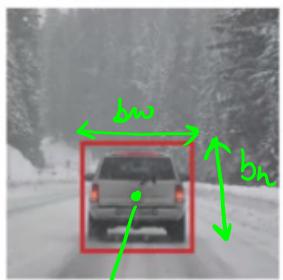


one central object

multiple objects

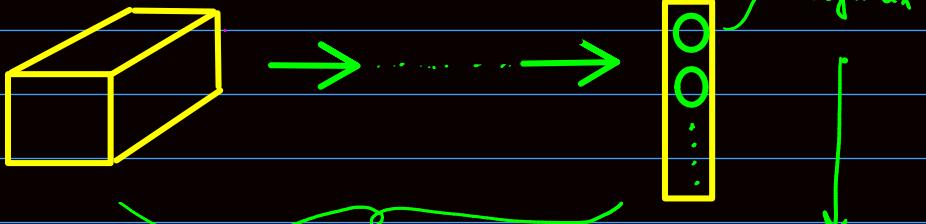
(0,0)

Classification with  
localization



$(b_x, b_y)$

(1,1)



ConvNet

need to output

- ① 4 classes (say)
- ②  $b_x, b_y, b_h, b_w$

Eq.  $b_x = 0.5, b_y = 0.7, b_h = 0.3,$   
 $b_w = 0.4$

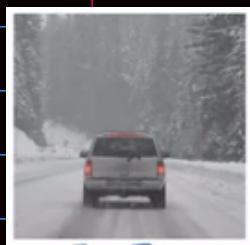
Let the classes be, ① Pedestrian ( $c_1$ )

② Car ( $c_2$ )

③ Motorcycle ( $c_3$ )

④ Background

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \text{Is there any object (1/0)}$$



Eq

$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

don't care.

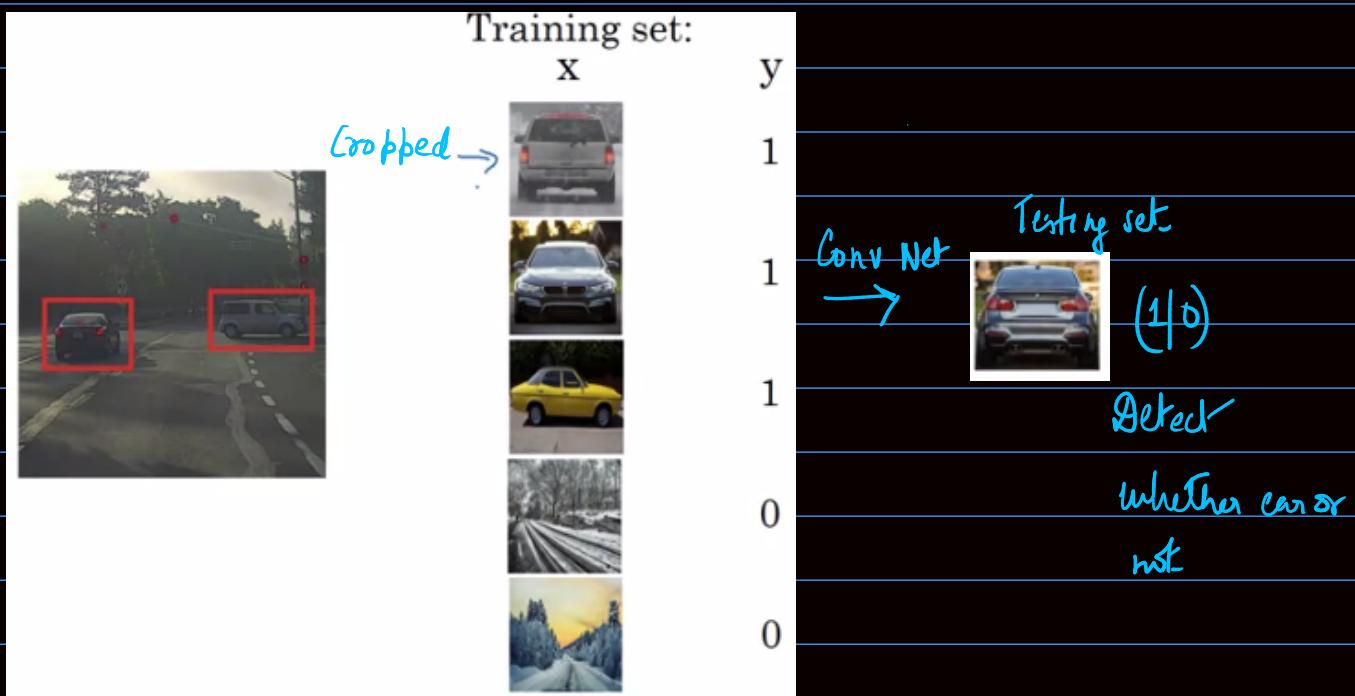
$\hookrightarrow$  Loss function

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

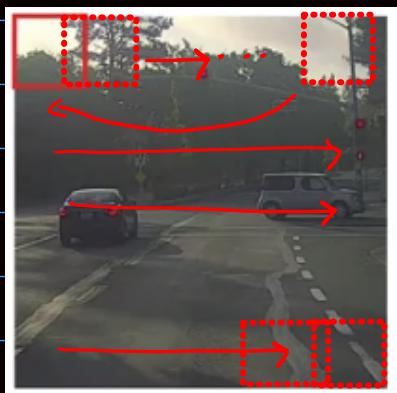
## Landmark Detection



## Object Detection



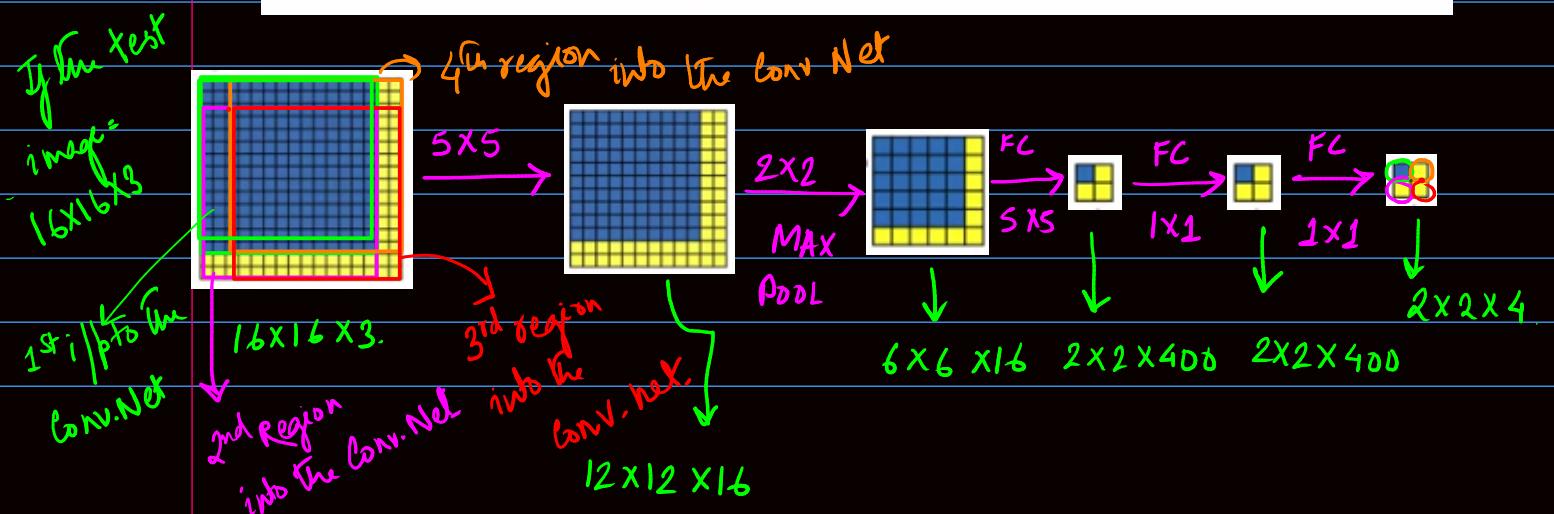
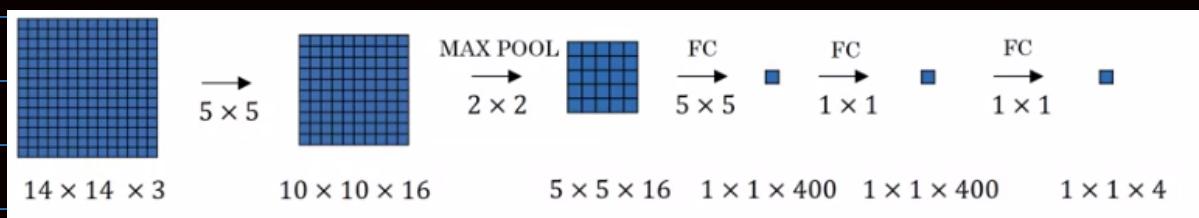
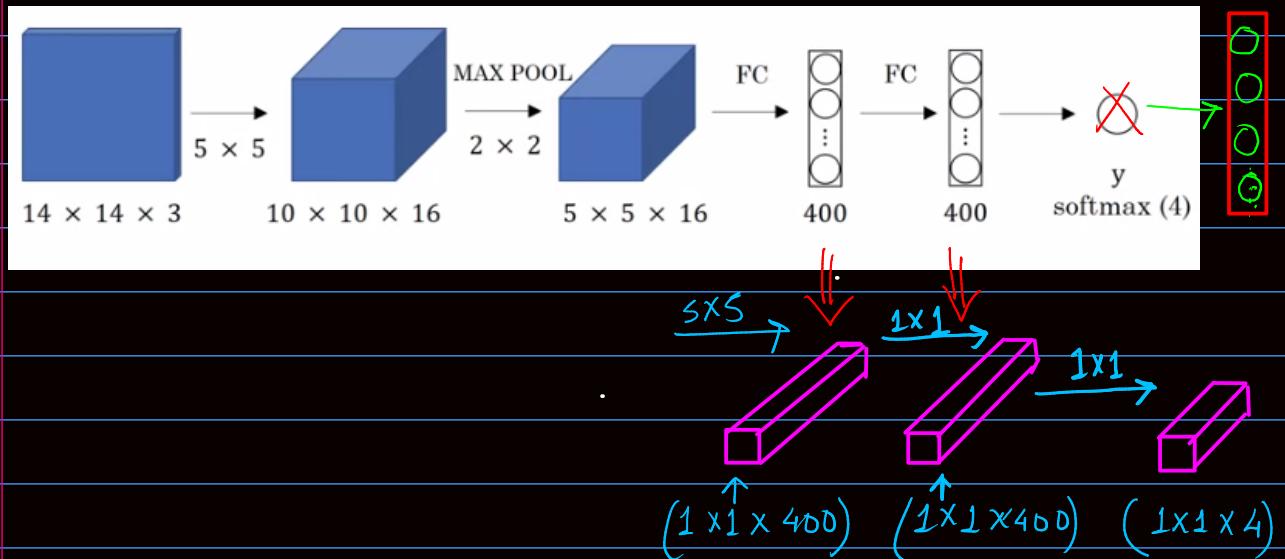
## Sliding Windows Detection



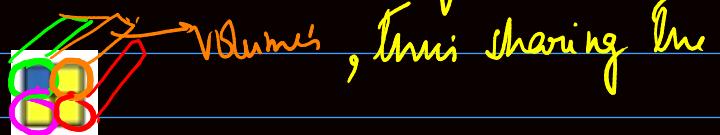
Windows  
Slide a kernel of varying sizes & strides all over the image

\* Disadvantages: longer computation as we are cropping a lot throughout the image and much slow

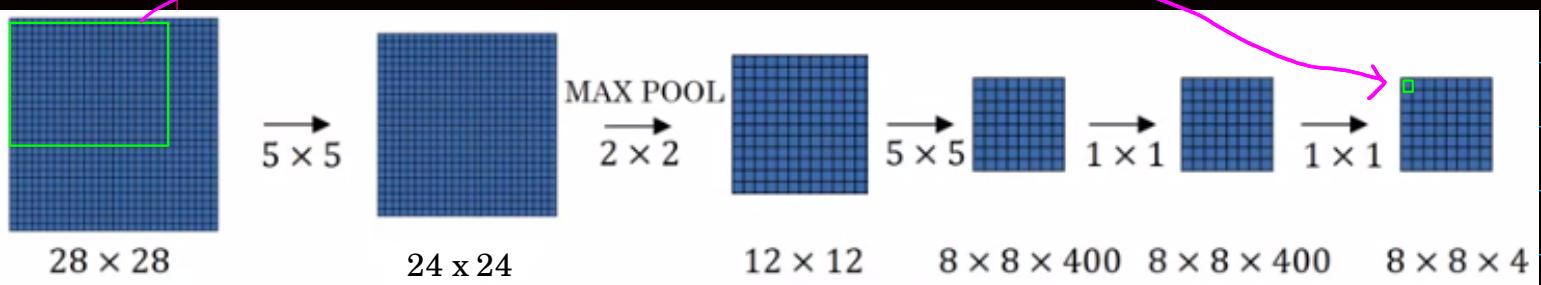
## Convolutional Implementation of Sliding Window



\* So individually passing all the four regions into the ConvNet is expensive rather we can pass the entire net & the output volume gives us the representation of all the four nets

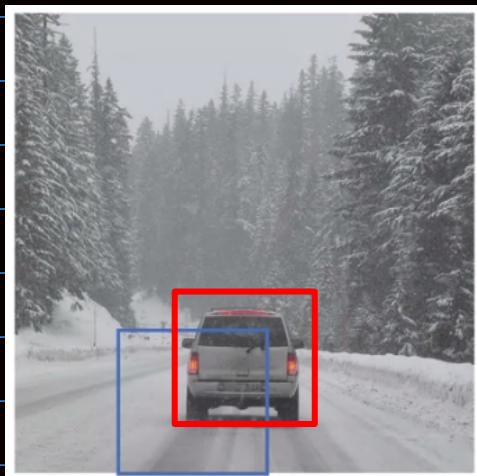


### Example:



∴ Instead of many forward. pass for each separate sliding window we do all the predictions together.

### Bounding Box Predictions



→ None of the bounding box fits the object well.

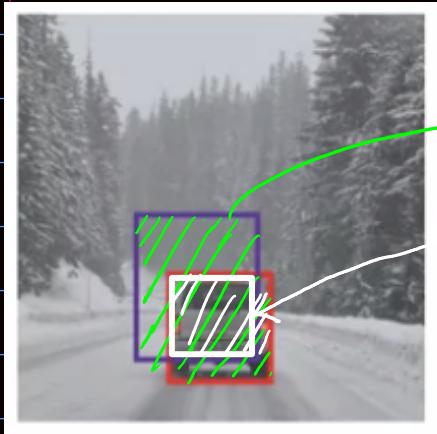
→ Ground Truth.

\* With normal window slides we might not get the bounding box even (ground truth)



\* Or in this type of example with different dimensional bounding box, this normal convolutional sliding becomes inefficient.

## Intersection over Union (IOU)



Union

Intersection

$$IOU = \frac{\text{Size of } \cap}{\text{Size of } \cup}$$

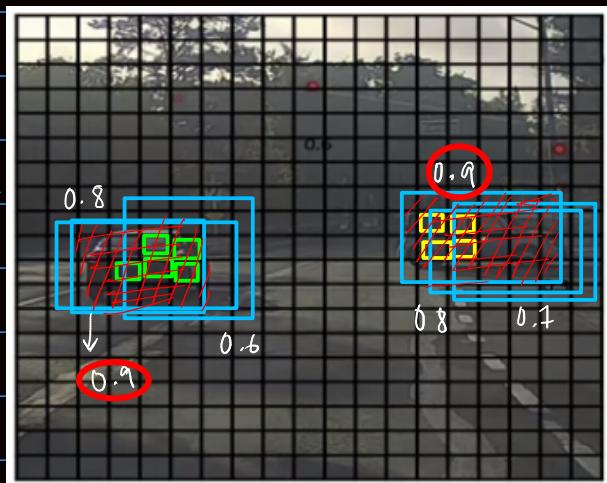
"Correct" if  $IOU \geq 0.5$

Higher the IOU more accurate is the bounding box

IOU is the measure of the overlap between two bounding boxes.

## Non-max-Suppression

It is a way to make sure that the algorithm detects the object only once



→ Multiple detections per car.  
( $P_c \rightarrow \text{Prob}$ )  
↳ only max prob will remain &  
others will get suppressed.

The max. prob bounding boxes will remain



19x19

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

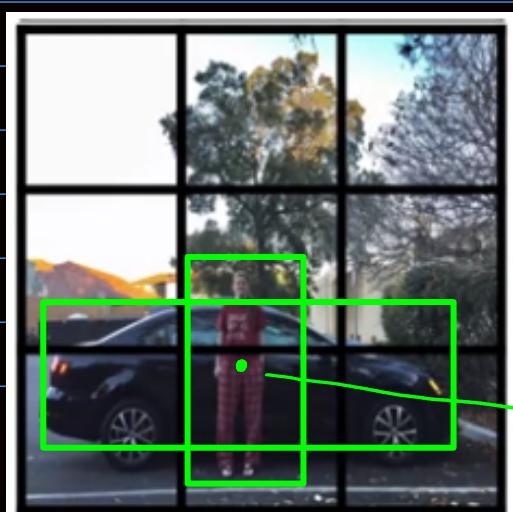
$\left. \begin{array}{l} C_1, C_2 \text{ are absent since} \\ \text{we have only one object} \end{array} \right\}$

### Algorithm

- ① Describe the output prediction i.e.  $y'$
- ② Discard all boxes with  $p_c \leq 0.6$
- ③ While there are any remaining boxes —
  - Pick the box with largest  $p_c$  & output that as prediction
  - Discard any remaining box with  $IoU \geq 0.5$  with the box o/p in the previous step.

\* If more than one class we apply non-max suppression to all the objects separately.

### Anchor Boxes



$$y_1 = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

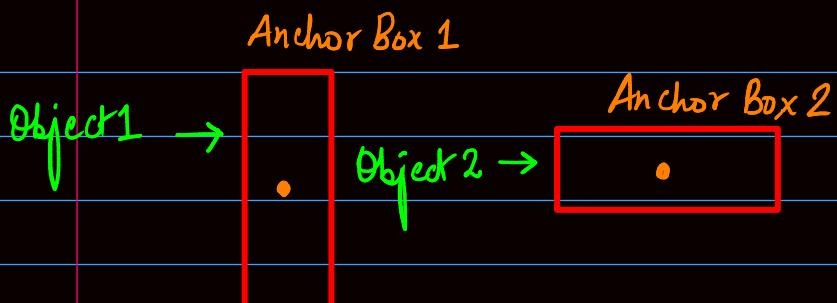
(Say)

$C_1 \rightarrow \text{Car.}$

$C_2 \rightarrow \text{Person}$

$C_3 \rightarrow \text{Cycle.}$

Same centroid for both the objects  
(Person & Car)



$$y = \begin{bmatrix} p_c \\ b_x \\ \vdots \\ c_1 \\ p_c \\ b_x \\ \vdots \\ c_1 \\ \vdots \end{bmatrix} \quad \left. \begin{array}{l} \text{anchor box 1} \\ \text{anchor box 2} \end{array} \right\}$$

Each object in training image is assigned to grid cell that contain object midpoint & anchor box for the grid cell with highest IoU.

(grid cell, anchor box)

Output box:  $y = 3 \times 3 \times 2 \times 8$

↓  
2 Anchor 8D  
Box anchor

## Anchor box example



Anchor box 1: [●]      Anchor box 2: [●]



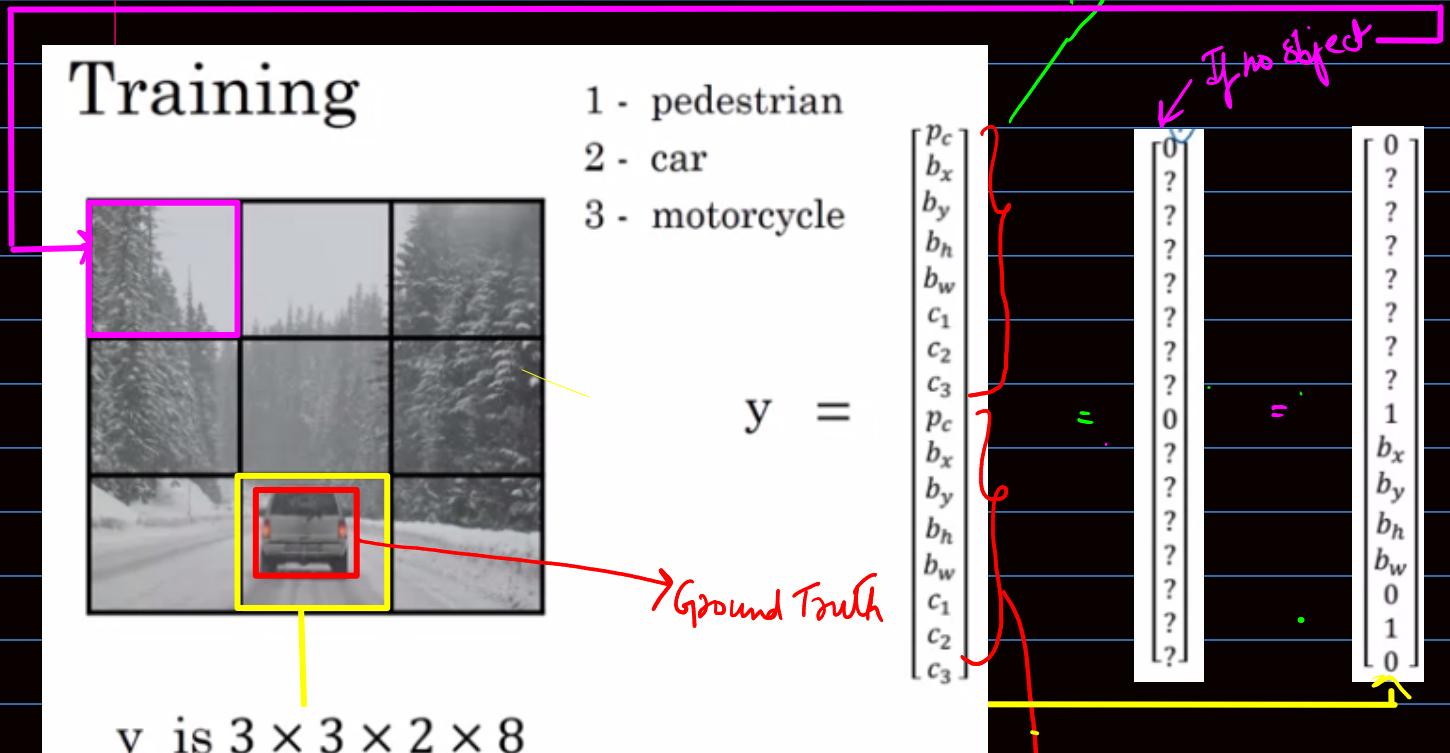
2 Object	Only 1 obj (car)
1	0
$b_x$	$b_x$
$b_y$	$b_y$
$b_h$	$b_h$
$b_w$	$b_w$
$c_1$	0
$c_2$	1
$c_3$	0
$p_c$	0
$b_x$	1
$b_y$	$b_y$
$b_h$	$b_h$
$b_w$	$b_w$
$c_1$	0
$c_2$	1
$c_3$	0

$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$

\* Anchor Box Algorithm does not work well when,  
 ① 3 objects & 2 Anchor Boxes

② 2 objects & 1 Anchor Box

## YOLO ALGORITHM



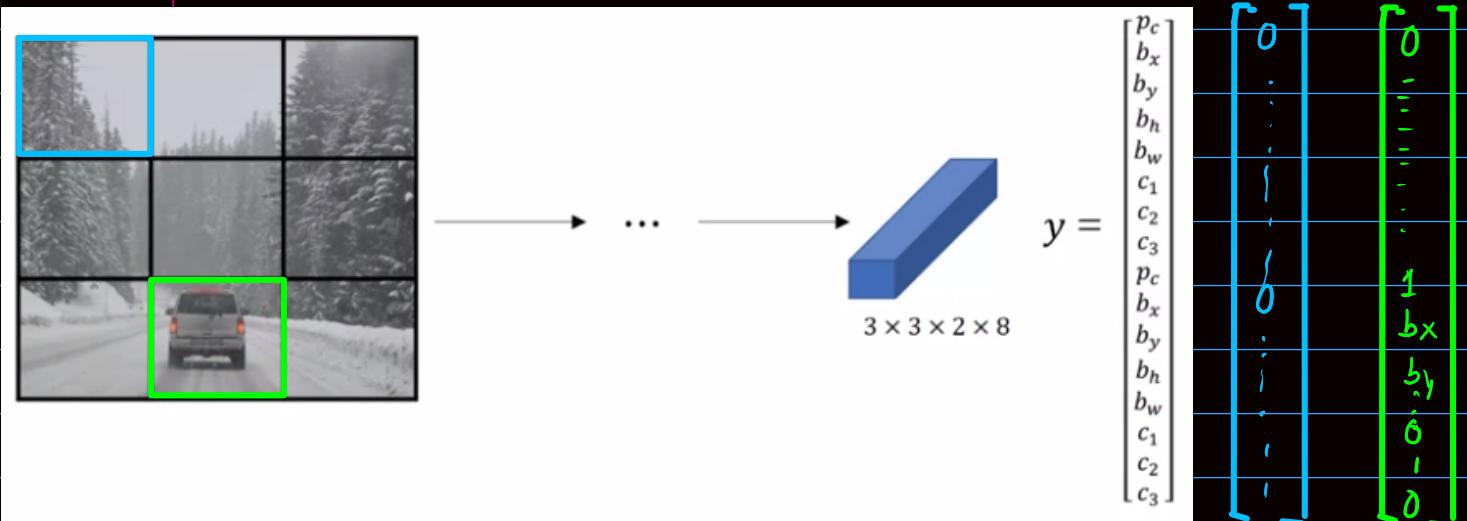
$y$

Anchor Box 1  $\rightarrow$



Anchor Box  $\rightarrow$





Outputting the non-max suppressed outputs.

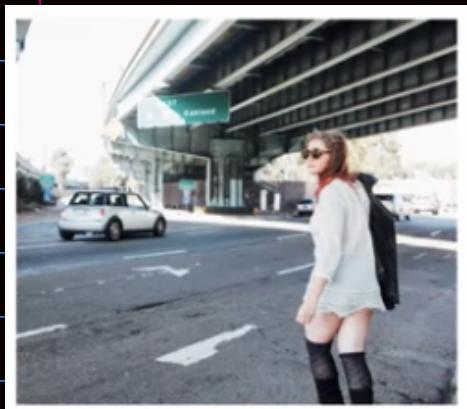


## Region Proposals (R-CNN's)

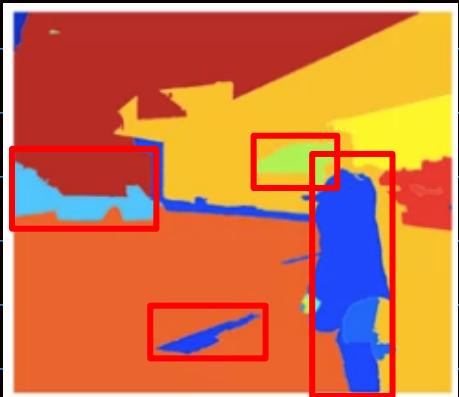


This algorithm chooses some particular regions which makes sense to pass through the conv Net.

Rather than running the classifier through all windows we select some specific windows



Segmentation →



- ① Propose Regions
- ② Classify Proposed Regions one at a time
- ③ Find the accurate Output label  
Bounding Box

Selects the blob regions  
(Say) find 2000 blobs &  
running the classifier on only  
those blobs.

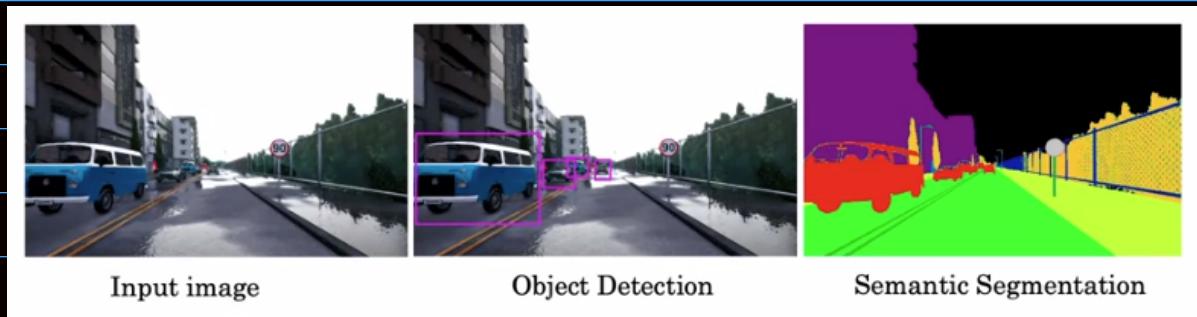
Drawbacks: The algorithm is quite slow.

Fast-(R-CNN) → Propose Regions Use Convolution implementation of  
sliding windows to classify all the proposed regions

Faster R-CNN: Use convolutional network to propose regions

\* Andrew prefers YOLO

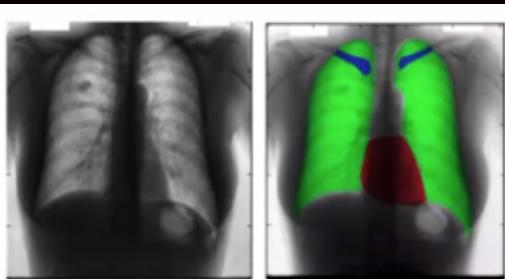
### Semantic Segmentation with U-Net



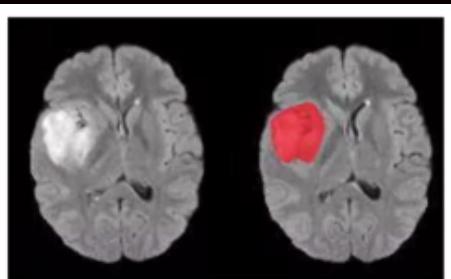
↑  
Indicates every single pixels

Uses: Self Driving can engg. to know about the drivable surface.

### Other uses:

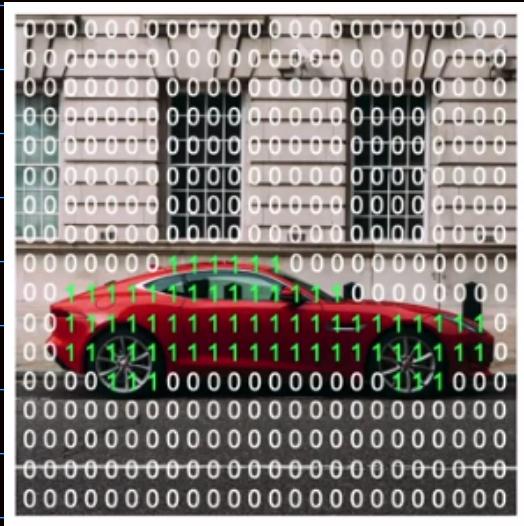


Chest X-Ray

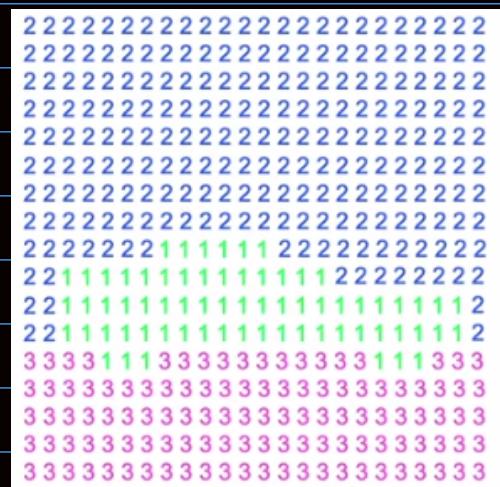
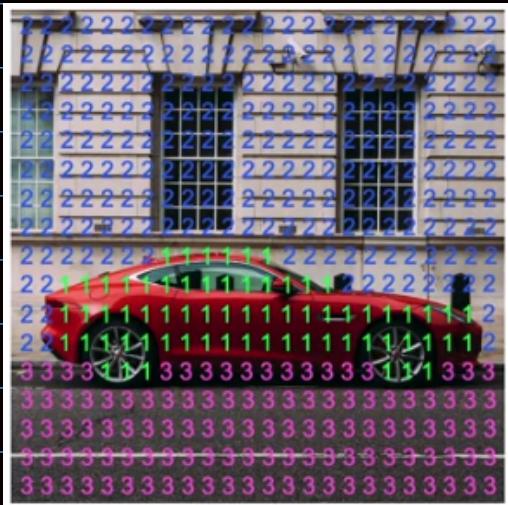


Brain MRI

## Per Pixel Class Labels



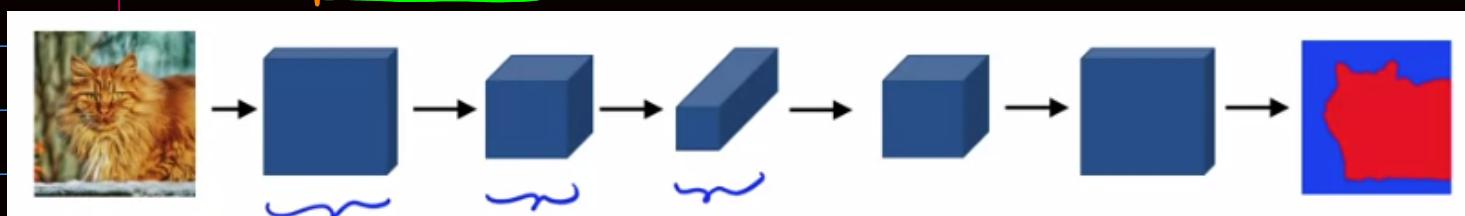
$1 \rightarrow$  can ,  $0 \rightarrow$  not can .



1: Car, 2: Building,  
3: Road

## Segmentation Map.

## UNet Implementation:



## Transpose Convolution

$$\begin{matrix} \text{2x2} \\ \boxed{\begin{matrix} & & \\ & & \end{matrix}} \end{matrix} * \begin{matrix} \text{3x3} \\ \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \end{matrix} = \begin{matrix} \text{4x4} \\ \boxed{\begin{matrix} & & & \\ & & & \\ & & & \\ & & & \end{matrix}} \end{matrix}$$

Upsampling ↑

$$\begin{matrix} 2 & 1 \\ 3 & 2 \end{matrix}$$

2x2

$$\begin{matrix} 1 & 2 & 1 \\ 2 & 0 & 1 \\ 2 & 2 & 1 \\ 0 & 2 & 1 \end{matrix}$$

$$\begin{matrix} / & / & / & / & / & / & / \\ 0 & 2+2 & 0 & 1 \\ / & 4 & 2+0 & 2 & 1 \\ \dots & & & & \end{matrix}$$

and so on . . .

$$\text{filter} = f \times f = 3 \times 3 \text{ (say)}$$

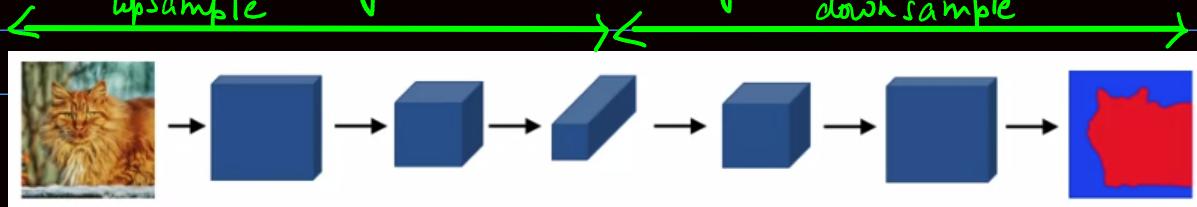
$$\text{padding} = 1$$

$$\text{stride} = 2$$

$$\begin{matrix} 0 & 4 & 0 & 1 \\ 10 & 7 & 6 & 3 \\ 0 & 7 & 0 & 2 \\ 6 & 3 & 4 & 2 \end{matrix}$$

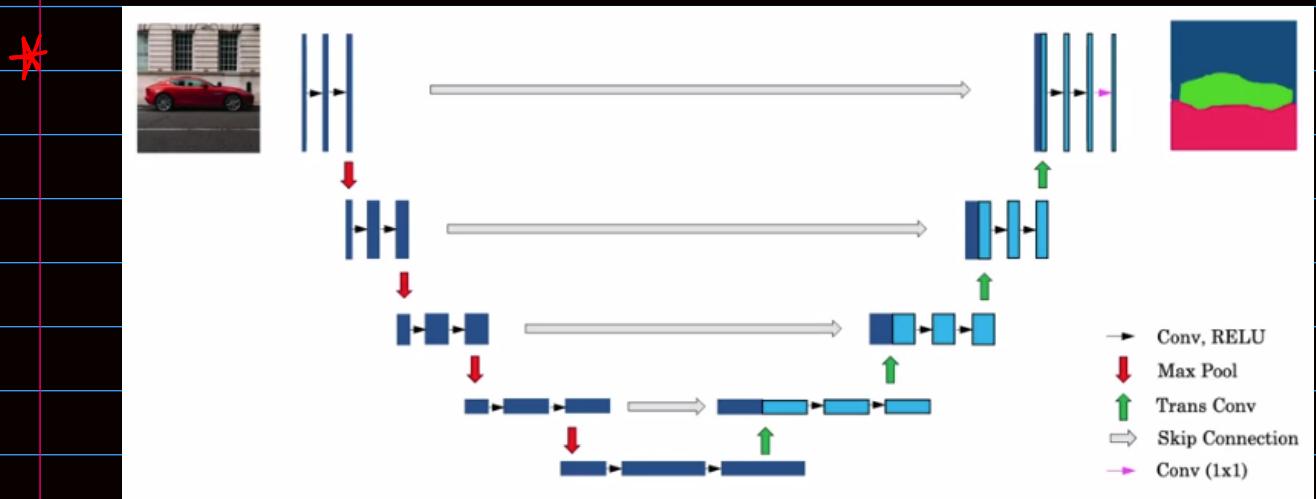
## U-Net Architecture

Deep Learning for Semantic Segmentation



Skip Connections are used in U-Net to get the low level information which was downsampled to get the more complex high level features

## Architecture



$\rightarrow$  Conv, RELU  $\downarrow$  MaxPool  $\uparrow$  Trans Conv  $\Rightarrow$  Skip Connections  $\rightarrow$  Conv = 1x1

## Week 4

### Face Recognition

#### Verification

I/p image, name ID

o/p whether the input image is that of the claimed person

#### Recognition

Has a database of K persons

Get an i/p image.

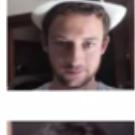
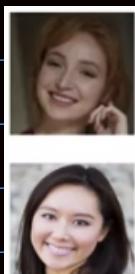
b/p ID if the image is any of the K Persons

### One Shot Learning Problem

Learning from one example to recognise the person again



X



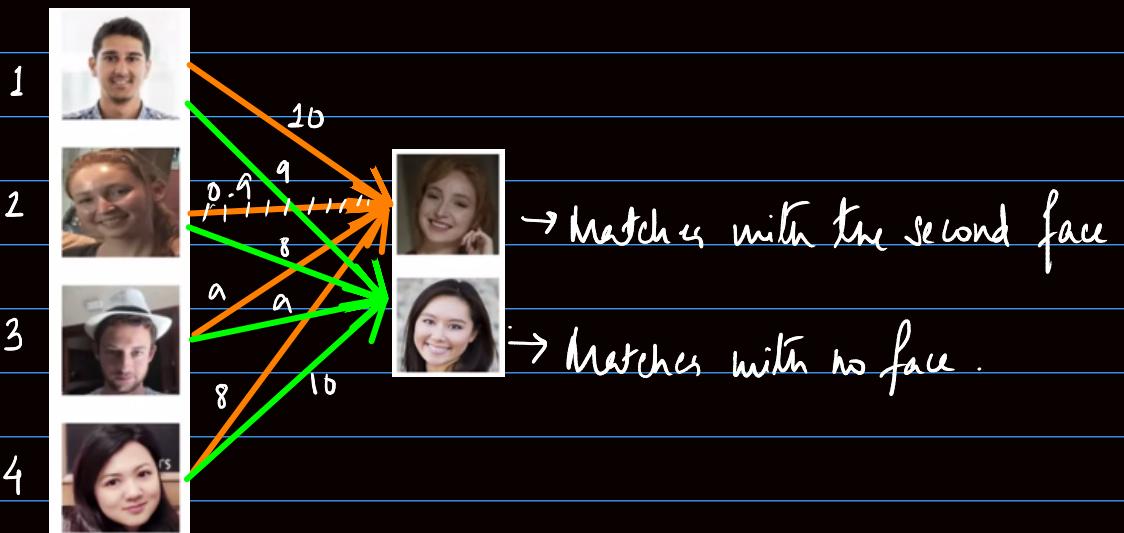
## Learning a 'similarity' function

$d(\text{img}_1, \text{img}_2)$  = degree of difference between images

$d(\text{img}_1, \text{img}_2) \leq 1 \xrightarrow{\text{some threshold}} \text{'same'}$

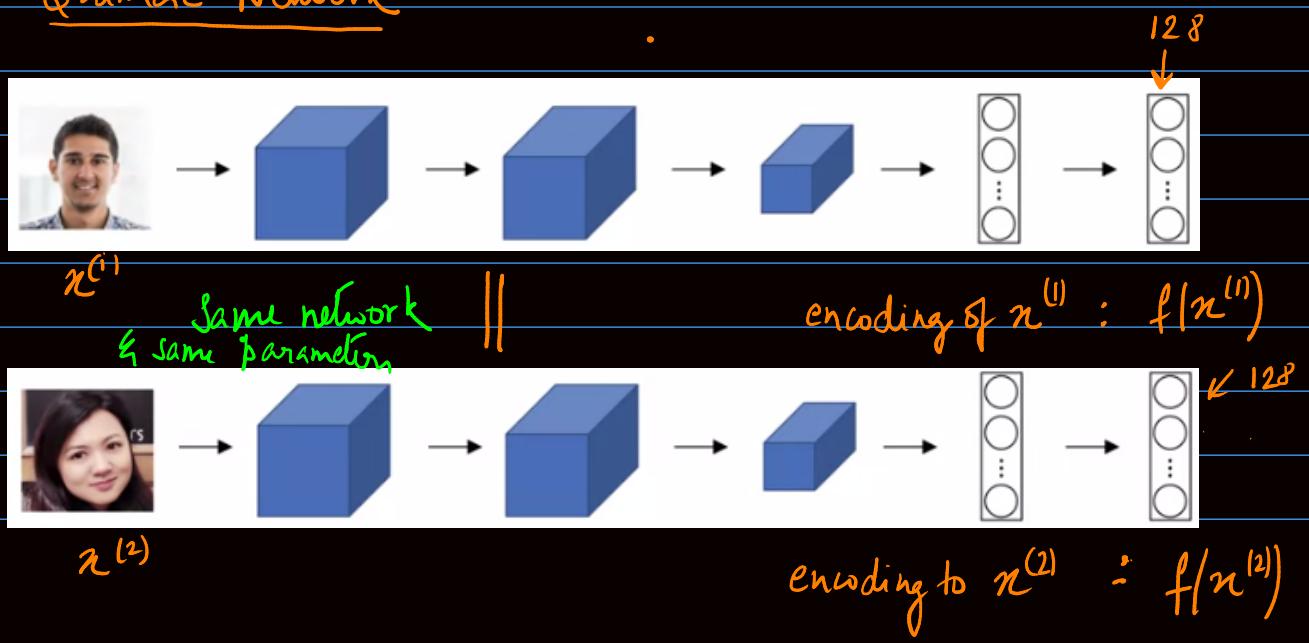
$> 1 \rightarrow \text{'different'}$

face Verification Problem



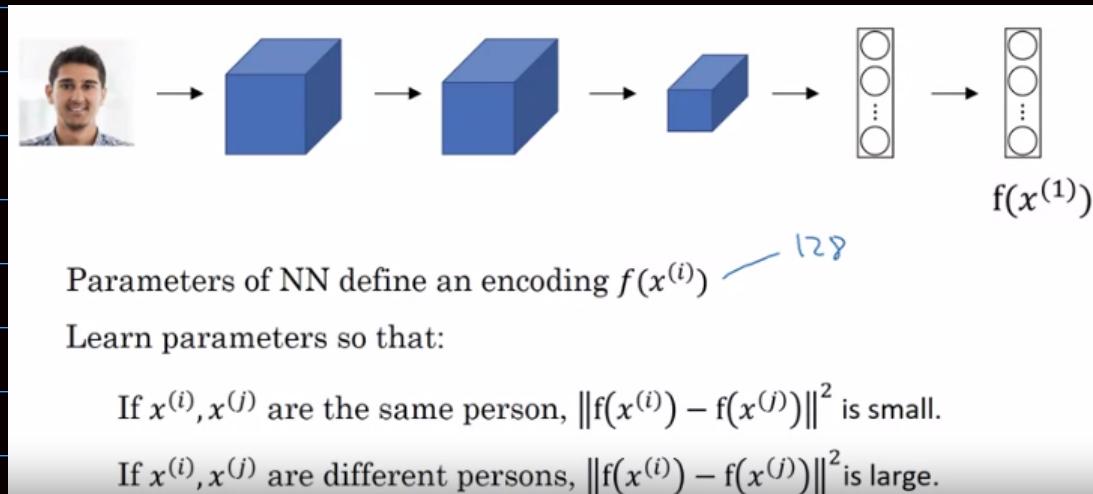
∴ In one shot learning adding extra images to the dataset will not effect the face verification.

## Siamese Network



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_F^2$$

This way of calling two NN's for 2 different  $x^{(i)}$ 's & then comparing them → Siamese network.



## Triplet Loss



Anchor  
A

Positive  
P

Anchor  
A

Negative  
N

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

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

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

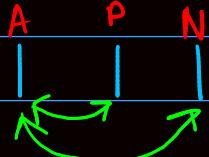
margin (pushes AP & AN gap further)

## Loss Function

Given 3 images  $\overbrace{A, P, N}$

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

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



e.g. Training Set of 10K Pictures  $\in$  1K persons

(A P)  $\rightarrow$  We need multiple pictures of the same person  $\therefore$  1K persons  $\in$  10K images

Choosing the triplets  $A, P \notin N$

1. During training if  $A, P \notin N$  are chosen randomly

$$\begin{aligned} d(A, P) + \alpha &\leq d(A, N) \text{ is easily satisfied} \\ \|f(A) - f(P)\|^2 + \alpha &\leq \|f(A) - f(N)\|^2 \end{aligned}$$

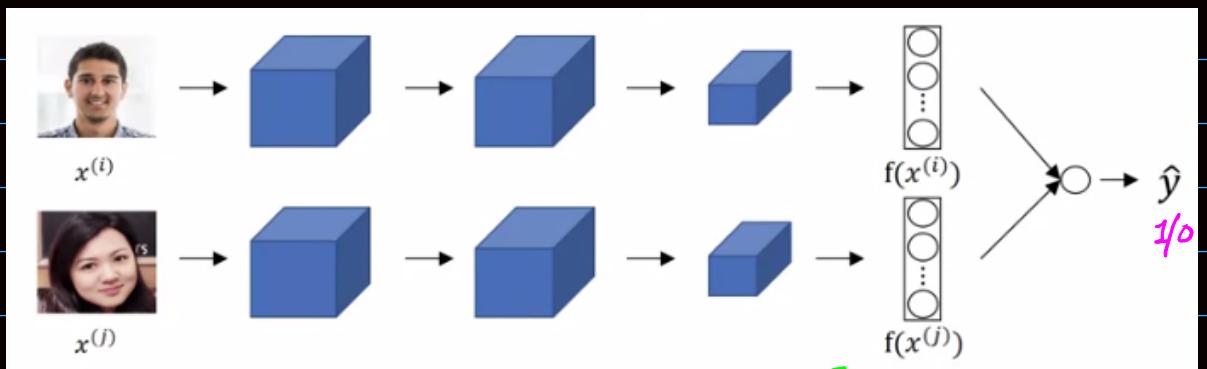
2. Choose triplets that are hard to train on.

$$\text{s.t. } \underline{d(A, P)} \approx \underline{d(A, N)}$$

so that to satisfy.

$$d(A, P) + \alpha \leq d(A, N) \text{ the classifier finds it hard}$$

# Face Verification as a Binary Classification Problem



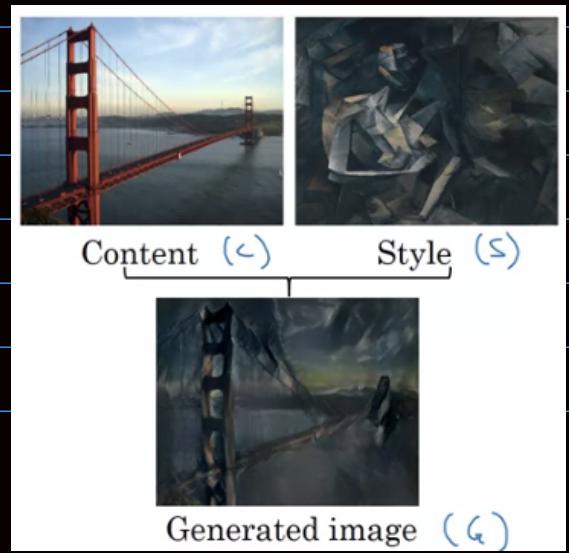
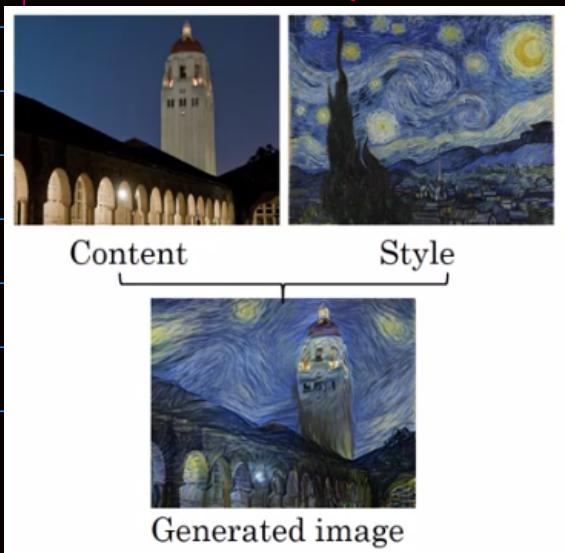
$$\hat{y} = \alpha \left( \sum_{k=1}^{128} w_k \left| f(x^{(i)})_k - f(x^{(j)})_k \right| + b \right)$$

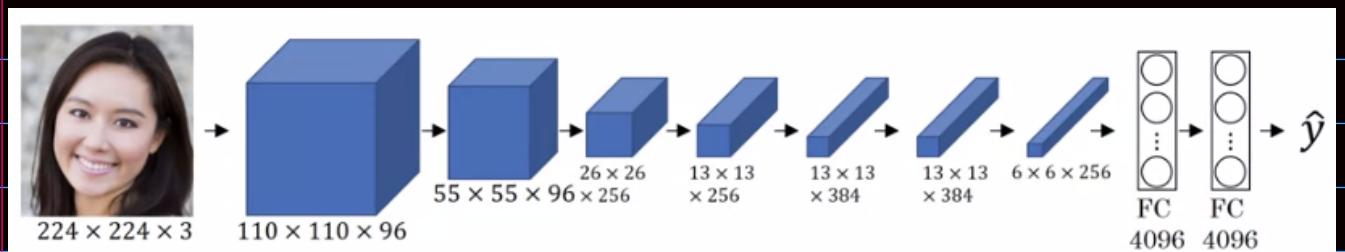
element wise diff betn the two encodings

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

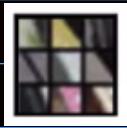
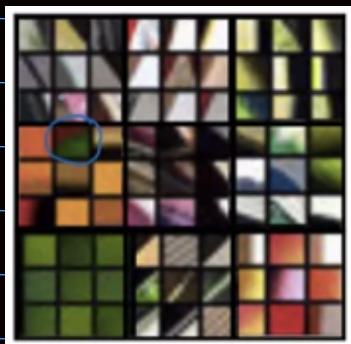
\* If a new person is included in the dataset (not instead of computing in the entire Siamese Network), we can use the upper conv. net layers to directly get the encoding & then comparing it with the precomputed encodings

## Neural Style Transfer

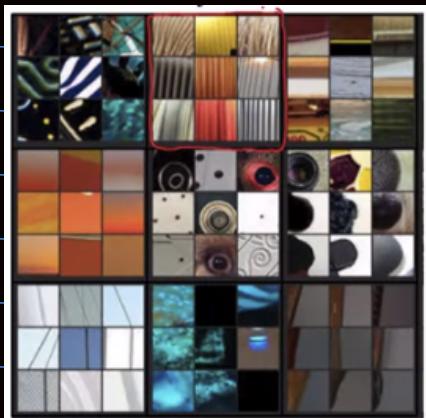




Pick a unit in layer 1. Find the nine image patches that maximizes the unit's activation  
Repeat for other units.



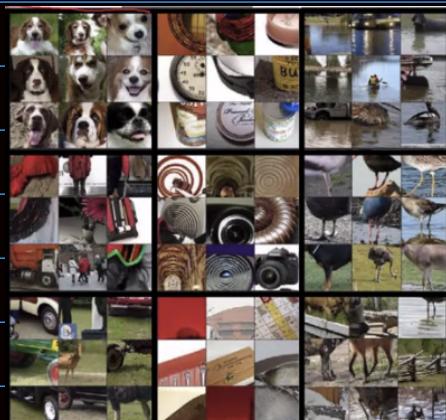
← Layer 1 (These are the 9 patches that cause Layer 1 to be highly activated!)



Layer 2



Layer 3

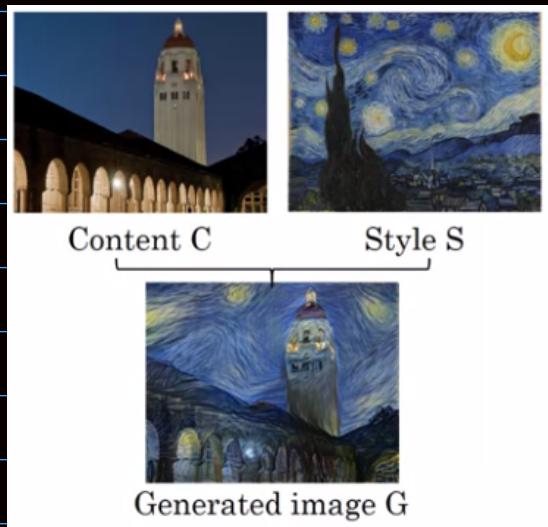


Layer 4



Layer 5

## Cost Function (for Neural Style Transfer)



$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

hyper parameters

Find The Generated Image  $G_1$ ,

1. Initiate  $G_1$  randomly.

$$G_1: 100 \times 100 \times 3$$

2. Use gradient descent to minimize  $J(G)$

$$G = G_1 - \frac{d}{dG} J(G)$$



→ Random Noise



+



→



## Content Loss Function

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

① Say we use hidden layer  $l$  to compute content cost.

{  $l \rightarrow$  neither too shallow nor too deep }

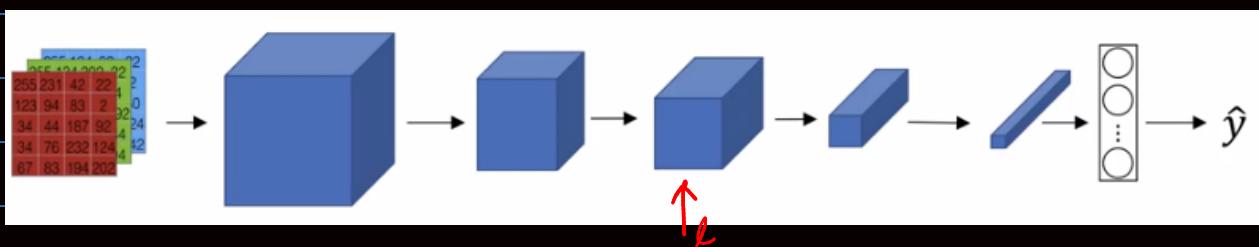
② Use Pretrained Conv. Net. (Eg VGG Net)

③ Let  $a^{[l]}(c)$  and  $a^{[l]}(G)$  be the activation of layer  $l$  on the images

④ If  $a^{[l]}(c)$  and  $a^{[l]}(G)$  are similar, both images have similar content.

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[l]}(c) - a^{[l]}(G) \|^2$$

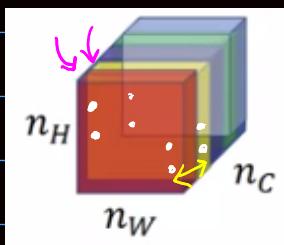
## Meaning of "style" of an image



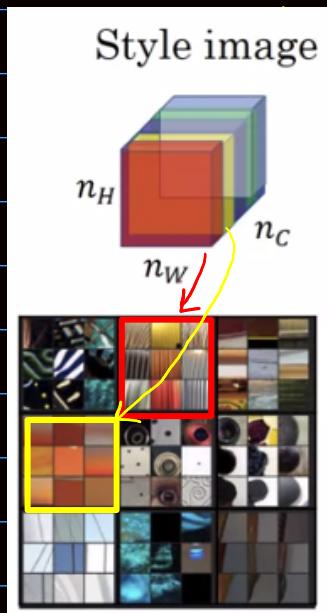
→ Say you are using layer  $l$ 's activation to measure 'style'

→ Define style as correlation between activations across channels

$$n_c = 5 \text{ (say)}$$



How correlated are the activation across different channels?



If these two channels are  
Correlated = whenever part of the image  
will have texture  will  
have have texture  also  
otherwise not correlated

- \* This degree of correlation gives us an idea about how often thin features occur together or don't occur.
- \* Used to compare the style of the generated image.

### STYLE MATRIX

Let  $a_{i,j,k}^{[l]}$  = activation at  $(i, j, k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$  no. of channels

$G_{KK'}^{[l]}, (k=1, 2, \dots, n_c^{[l]})$  = Will measure how correlated are activations in channel  $K$  to activations in channel  $K'$

$$G_{KK'}^{[L](s)} = \sum_i^{n_H^{[L]}} \sum_j^{n_W^{[L]}} a_{ijk}^{[L](s)} a_{ijK'}^{[L](s)}$$

$$G_{KK'}^{[L](G)} = \sum_i^{n_H^{[L]}} \sum_j^{n_W^{[L]}} a_{ijk}^{[L](G)} a_{ijK'}^{[L](G)}$$

If correlated  $G_{KK'}^{[L]} = \text{large}$ .

else,  $G_{KK'}^{[L]} = \text{small}$ .

\* 'G' in linear algebra is 'gram matrix'

$$\mathcal{J}_{\text{Style}}^{[L]}(s, G) = \left\| G^{[L](s)} - G^{[L](G)} \right\|^2$$

$$= \frac{1}{2n_H^{[L]} n_W^{[L]} n_C^{[L]}} \sum_K \sum_{K'} \left( G_{KK'}^{[L](s)} - G_{KK'}^{[L](G)} \right)^2$$

$\nearrow$   
Normalization  
constant

If we use Style cost fun^ from multiple layers

Therefore,  
Overall Cost Function

$$\mathcal{J}_{\text{Style}}(s, G) = \sum_L \gamma^{[L]} \mathcal{J}_{\text{Style}}^{[L]}(s, G).$$

\* Take both high level & low level layers  
into a/c

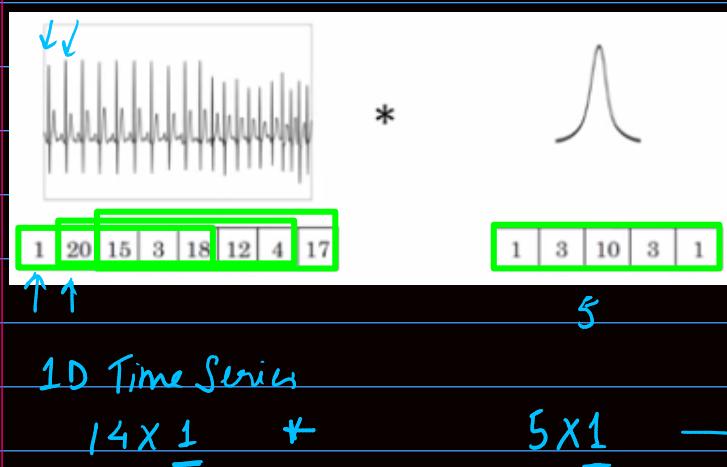
$$\mathcal{J}(G) = \alpha \mathcal{J}_{\text{Content}}(C, G) + \beta \mathcal{J}_{\text{Style}}(s, G)$$

$G$  to minimize  $\mathcal{J}(G)$

## 1D & 3D Generalization

### Convolutions

2D Input Image \* 2D Filter → Result  
 $14 \times 14 \times 3$        $5 \times 5 \times 3$        $10 \times 10 \times 16$



3D Data. → E.g. CT Scan, Movie Data

