

## Course 4

\* Kernel for vertical edge :  $\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$

Kernel for horizontal edge :  $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

Sobel filter :  $\begin{bmatrix} 1 & 0 & 1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$

Schless filter :  $\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$

$$\text{Ex: } \begin{matrix} 6 \times 6 \\ \text{image} \end{matrix} \times \begin{matrix} 3 \times 3 \\ \text{filter} \end{matrix} = \begin{matrix} 4 \times 4 \\ \text{image} \end{matrix}$$

$$\begin{matrix} n \times n \\ \text{image} \end{matrix} \times \begin{matrix} f \times f \\ \text{filter} \end{matrix} = \begin{matrix} (n-f+1) \times (n-f+1) \\ \text{image} \end{matrix}$$

Image shrinks on applying convolution filters.

Also, edge pixels ~~are~~ are not used as much as the non-edge pixels.

To counter this, 1 pixel padding is applied to the original image.

This increases weight or usage of edge pixels and retains the same size output as original image.

Valid convolution :  $n \times n \times f \times f = n-f+1 \times n-f+1$

Same convolution :  $n+2p \times f \times f = n \times n \times n+2p$

$$R = R + 2p - f + 1 \times n \times n \times n$$

$$\therefore p = \frac{f-1}{2} \quad p \rightarrow \text{padding}$$

filters are generally odd size square matrix.

\* Strided Convolution: jumps 's' pixels when the filter is moved

$$\text{Output size} : \left[ \frac{n+2p-f+1}{s} \right] \times \left[ \frac{n+2p-f+1}{s} \right]$$

$\lfloor x \rfloor$  = floor function

\* Convolution over Volume.

$$6 \times 6 \times 3 * 3 \times 3 \times 3 = 4 \times 4 \times 1$$

$\downarrow \quad \downarrow \quad \downarrow$  channels       $\downarrow \quad \downarrow \quad \downarrow$  channels       $\downarrow$   
 $h \quad w$        $h \quad w$        $h \quad w$

2D image

\* Notation

$f^{[C]}$  → filter layer

$n_c^{[C]}$  → number of filters

$p^{[C]}$  → padding

$s^{[C]}$  → stride

Input:  $n_h^{[C]} \times n_w^{[C]} \times n_c^{[C]}$

Output:  $n_h^{[C]} \times n_w^{[C]} \times n_c^{[C]}$

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

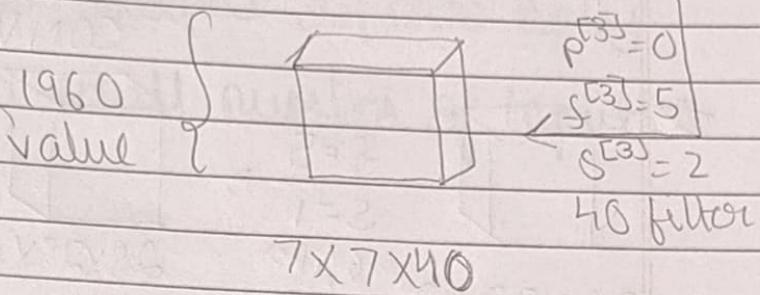
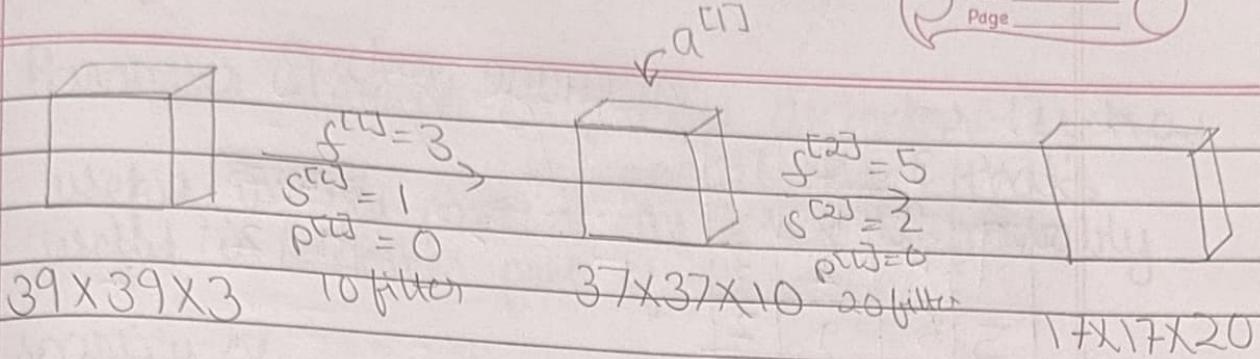
$$n_w^{[C]} = \left\lfloor \frac{n_w^{[C]} + 2p^{[C]} - f^{[C]}}{s^{[C]}} + 1 \right\rfloor$$

each filter →  $f^{[C]} \times f^{[C]} \times n_c^{[C]}$

$$a^{[C]} = n_h^{[C]} \times n_w^{[C]} \times n_c^{[C]}$$

$$\text{weight: } f^{[C]} \times f^{[C]} \times n_c^{[C]} \times n_c^{[C]}$$

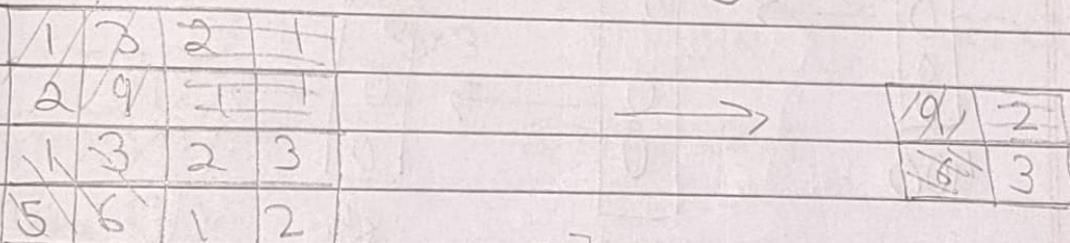
$$\text{bias: } n_c^{[C]}$$



generally,

$n, w \downarrow$  } as we go  
 $c \uparrow$  } deeper  
in the CNN

## \* Pooling Layer : Max pooling



$$\begin{aligned} f^{[1]} &= 2 \\ s^{[1]} &= 2 \end{aligned}$$

↓  
hyperparameters  
fixed

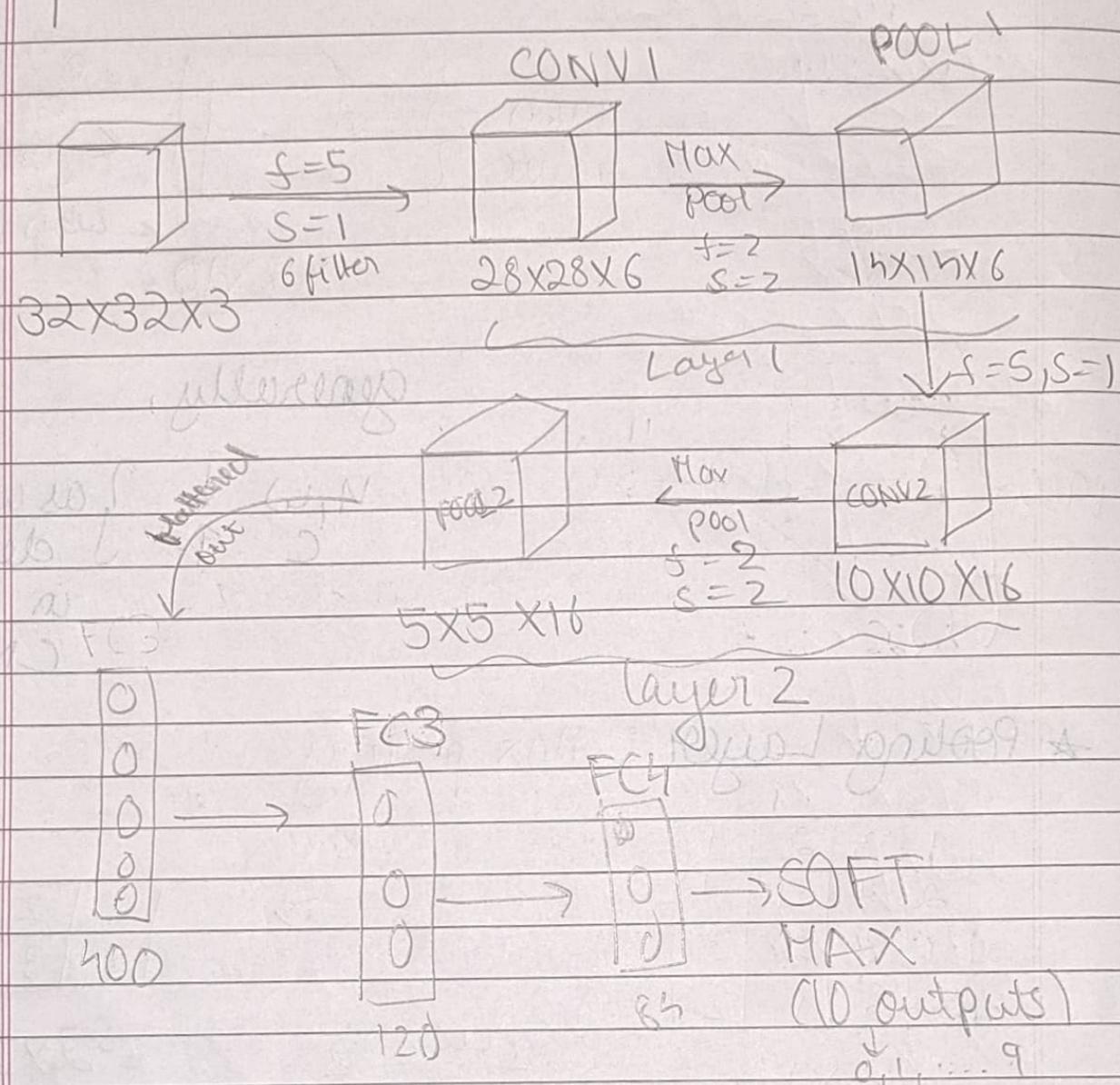
no learning

## average pooling

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

→

3	15	F:25
4		2



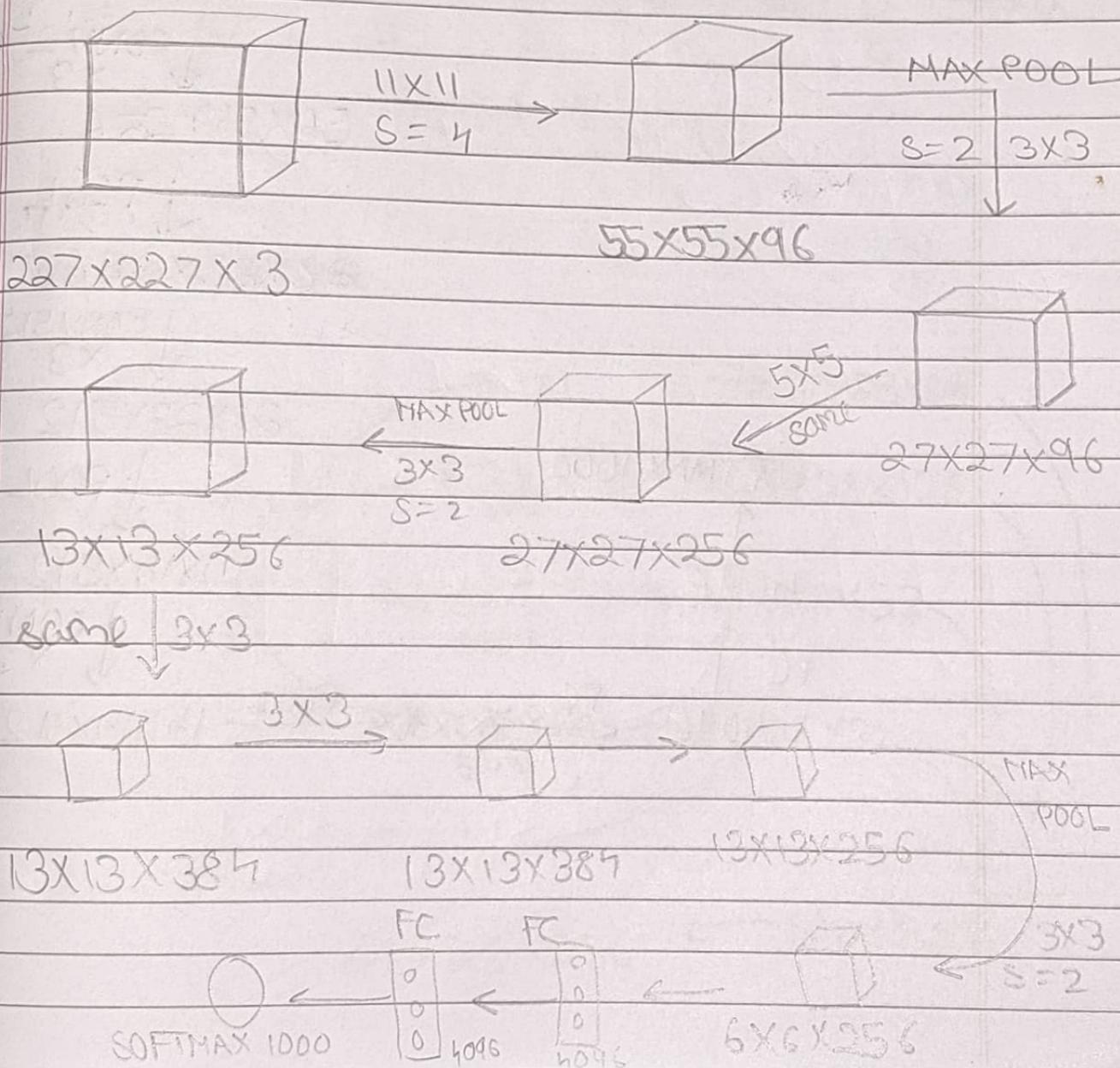
$n_h, n_w \downarrow$        $n_c \uparrow$

→

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

**Sparcity of connection:** each output depends only on small number of inputs

## \* Alexnet



~~60~~ ~~VGG-16~~: CONV =  $3 \times 3$ , S = 1, same  
~~POOL~~ =  $2 \times 2$ , S = 2

$224 \times 224 \times 3$   $\xrightarrow[\times 2]{\text{CONV} 64}$   $224 \times 224 \times 64$   
 $\downarrow \text{POOL}$

$112 \times 112 \times 64$   $\xrightarrow{\text{CONV} 128}$   
 $\downarrow \times 2$

$112 \times 112 \times 128$

$\downarrow \text{POOL}$

$56 \times 56 \times 128$

$\downarrow \text{CONV} 256$   
 $\times 3$

$56 \times 56 \times 256$

$\downarrow \text{POOL}$

$28 \times 28 \times 256$

$\downarrow \text{CONV} 512$   
 $\times 3$

$28 \times 28 \times 512$

$\downarrow \text{POOL}$

SOFTMAX 1000

$\uparrow$   
 4096

$\uparrow$   
 FC

4096

FC

$\xleftarrow{\text{POOL}}$   $7 \times 7 \times 512$

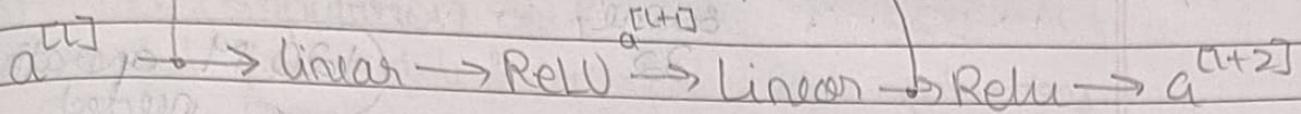
$14 \times 14 \times 512$

$\downarrow \text{CONV} 512$   
 $\times 3$

$14 \times 14 \times 512$

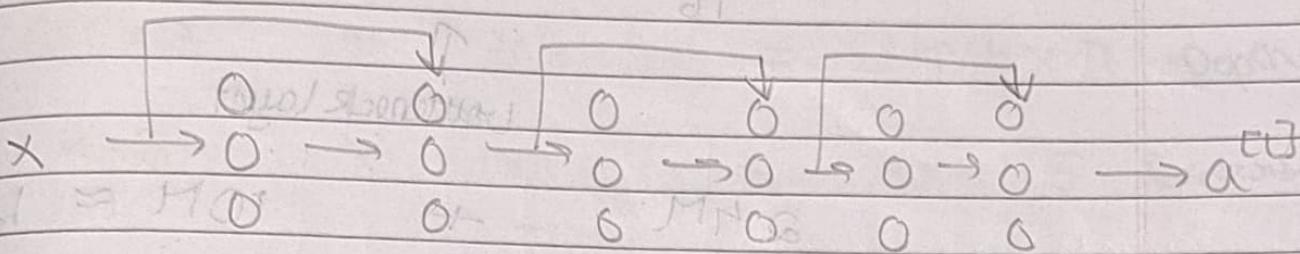
## \* Residual Block

Shortcut / skip connection



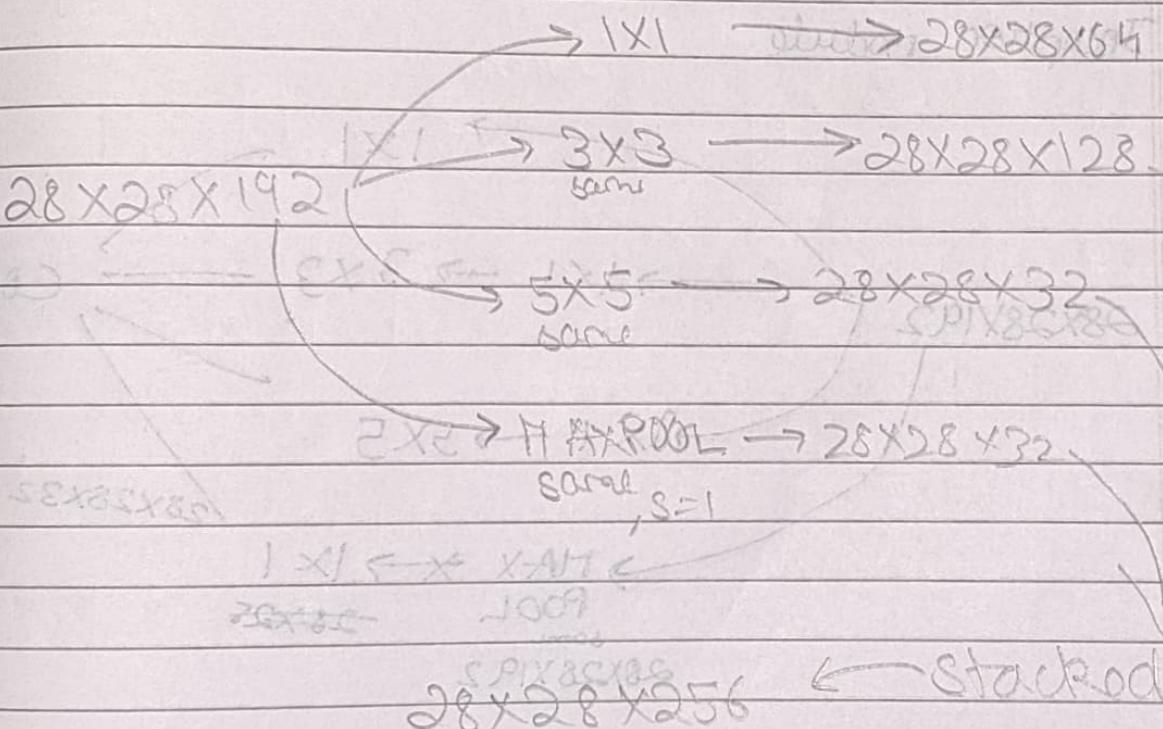
$$\therefore a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

## \* Resnet



ResNet Example.

## \* Inception network



\*  $28 \times 28 \times 192 \xrightarrow[\text{same, } 32]{\text{CONV } 5 \times 5} 28 \times 28 \times 32$

$120 \text{ M operations needed}$

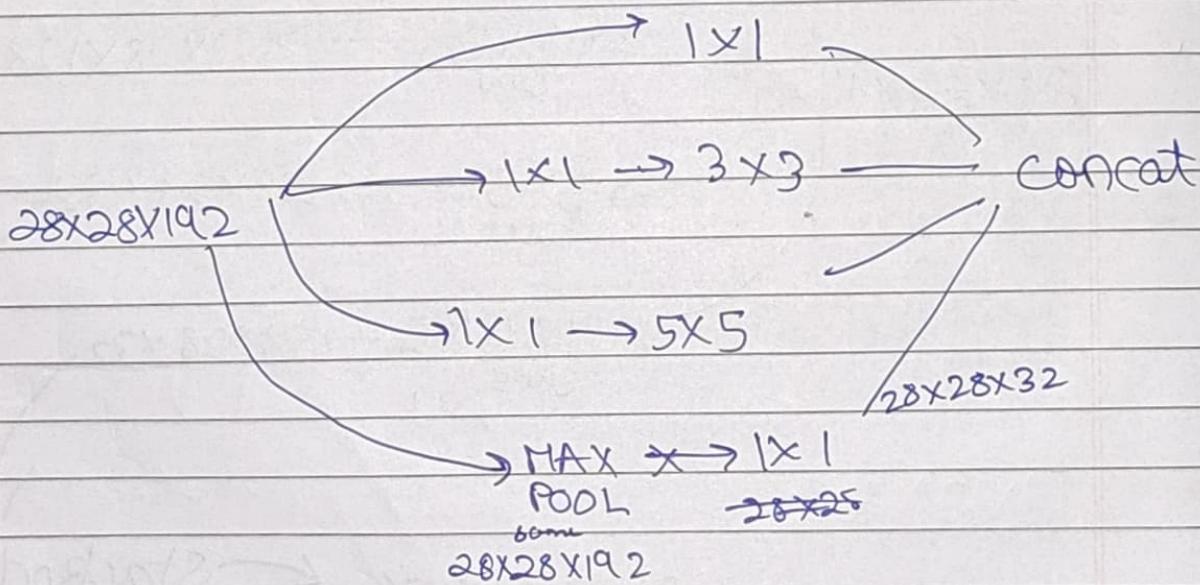
\*  $28 \times 28 \times 192 \xrightarrow[\text{1x1 } 16]{\text{CONV}} 28 \times 28 \times 16 \xrightarrow[\text{32 } 32]{\text{CONV } 5 \times 5} 28 \times 28 \times 32$

$\uparrow$   
bottleneck layer

$2.4 \text{ M} + 10 \text{ M} \approx 12.4 \text{ M}$

Output is of same dimension but the computation is reduced by a factor of 10

#### \* Inception module



## \* Mobile Nets

UN 1013

↳ low computation cost

useful for mobile & embedded vision applications

## # Depthwise Convolution

$$\text{I: } 6 \times 6 \times 3 \times 3 \rightarrow 4 \times 4 \times 3 \text{ Depthwise}$$

$$\text{II: } h \times w \times 3 \times 1 \times 1 \times 3 = h \times w \times 5 \text{ Pointwise}$$

5 filters

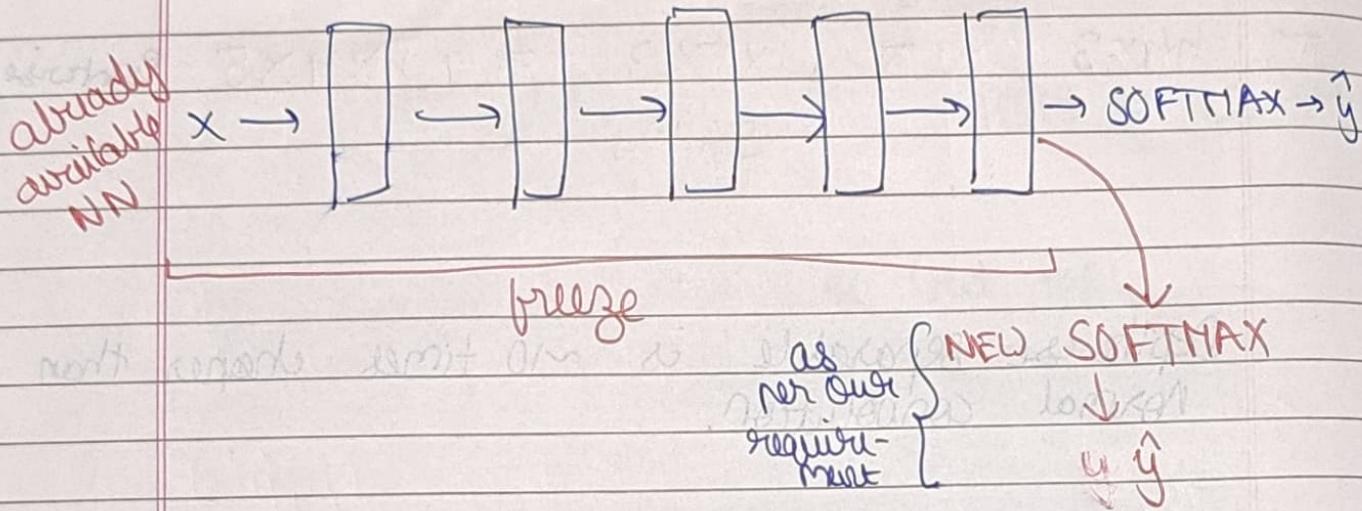
Depthwise separable is ~10 times cheaper than normal convolution.

## \* Efficient NN

→ best NN with limited computational power

Can be →  $r$  resolution of image  
done →  $d$  depth of NN  
using →  $w$  width of NN (No of neurons)

## \* Transfer Learning



Now the data → less layer freeze and train on the remaining layers

# Output of the frozen layers can be computed and saved to disk to enhance the speed of training.

# If lots of data then use the weight only as initialization.

## \* Data Augmentation

- Mirroring image }
  - Random cropping }
  - rotation }
  - shearing }
  - local wrapping }
  - Color shifting → can RGB value of all pixel
- not used as much

## \* # Ensembleing → average outputs of multiple NN

# Multicrop → 1 image → 10 version  $\xrightarrow{\text{NN}} \text{TL, TR, BL, BR}$

$$\begin{aligned} & \text{TL} \\ & \text{TR} \\ & \text{BL} \\ & \text{BR} \\ & \text{center crop} \\ & \text{left crop} \\ & \text{right crop} \\ & \text{top crop} \\ & \text{bottom crop} \\ & \text{left top crop} \\ & \text{right top crop} \\ & \text{left bottom crop} \\ & \text{right bottom crop} \end{aligned}$$

$$\frac{(\mu_1, \sigma_1) + (\mu_2, \sigma_2)}{2} = (\bar{\mu}, \bar{\sigma})$$

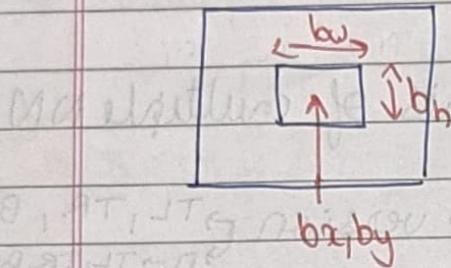
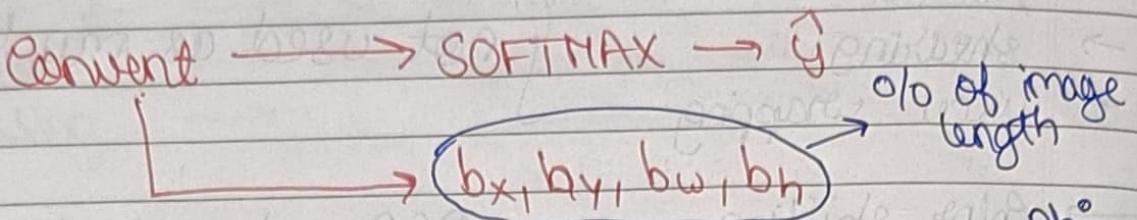
Post - tempo  
PSI patient

Low risk

downgrade  
(green)

## \* Object Detection

Classification with localisation

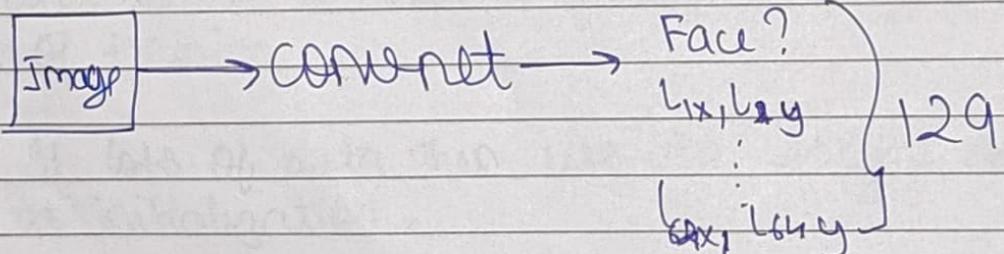


$$\hat{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \text{Obj hai kirdi?}$$

bounding box classes

$y_1 = 0 \rightarrow$  no obj → other parameters don't matter

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



→ landmark  
 (point)

detection

## \* Sliding Window Algorithm

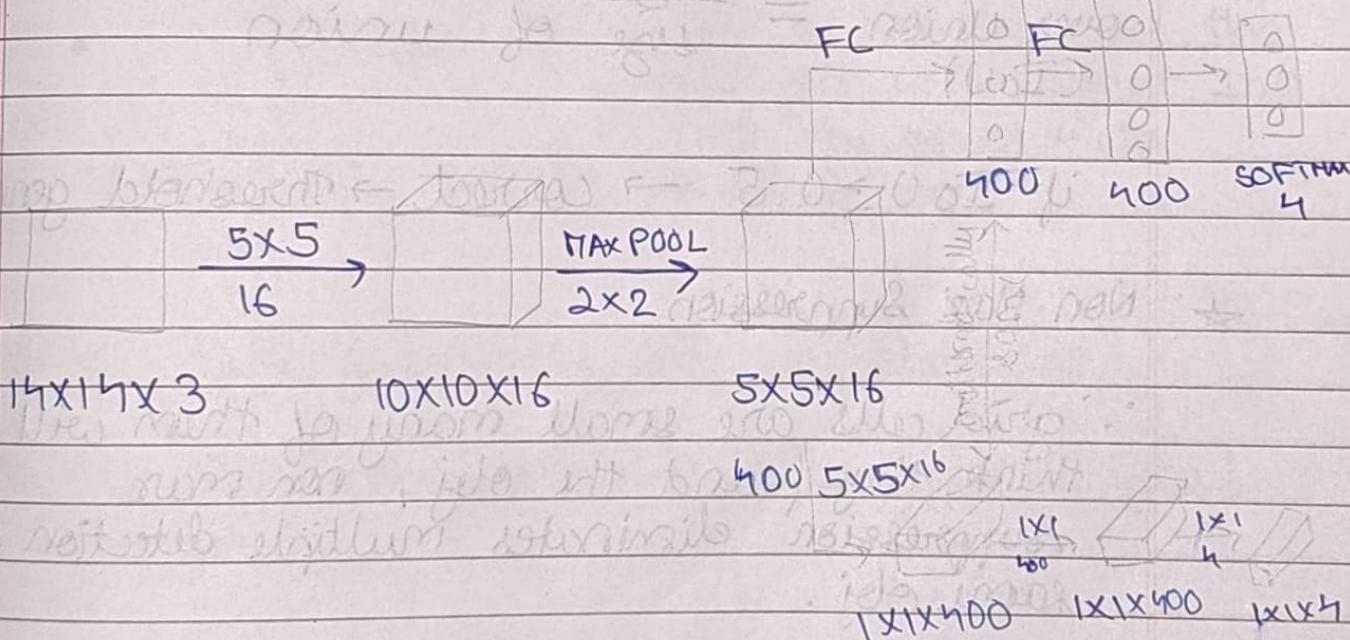
SI : Train a convnet using closely cropped car images i.e. image mainly is car

SII : Just like filters in convolution, take parts of images and feed to the CNN.

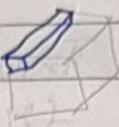
The window part of image is made bigger and bigger and resized to the size of input CNN requires.

# It is very expensive to run.

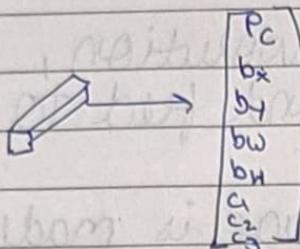
## \* Turning FC Layer into convolution layer.



## ★ YOLO Algorithm

multiple conv layer → and max pool → 

$100 \times 100 \times 3 \rightarrow 3 \times 3 \times 8$



# Object is assigned to only one of the grid cells even if it spans multiple grid cells.

$b_x, b_y, b_w, b_h$  are specified relative to grid cell  
 $(0,1)$  would be  $> 1$

## ★ Evaluating object localization

$$\text{Intersection over Union} = \frac{\text{size of intersection}}{\text{size of union}}$$

if  $\text{IoU} \geq 0.5 \rightarrow \text{correct} \rightarrow \text{threshold generally}$

## ★ Non Max Suppression

∴ grid cells are small many of them will think they found the obj, non max suppression eliminates multiple detection of same obj.

It takes the most confident prediction and eliminate the other prediction which have high IoU with the most confident prediction.

Also,  $P_c \leq 0.6$  are discarded anyway

remaining boxes having  $\text{IoU} \geq 0.5$  with most confident prediction are eliminated.

### \* Anchor box

define many anchor boxes and make y output  $P_c, b_x, b_y, b_w, b_h, c_1, c_2, c_3$  corresponding to each anchor box.

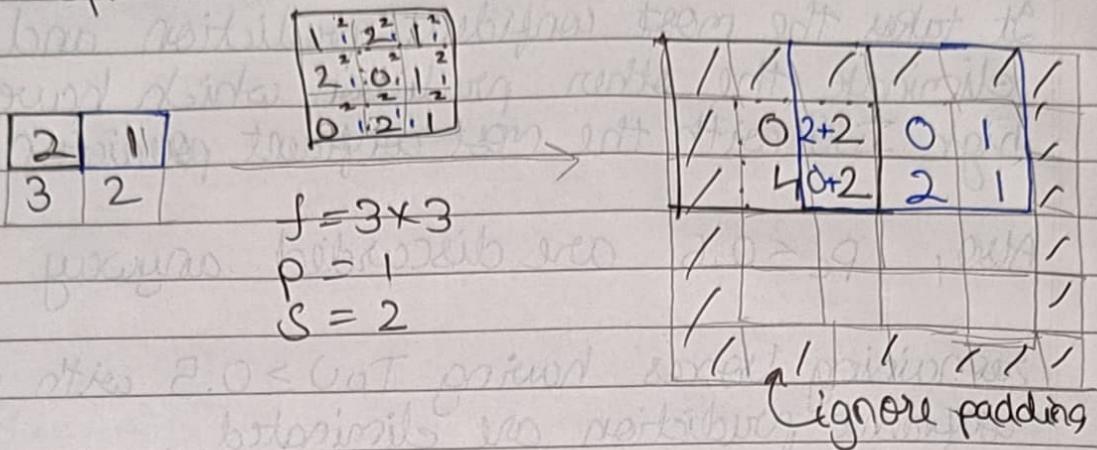
$$y = \begin{bmatrix} P_c \\ b_x \\ \vdots \\ c_3 \\ P_c \\ b_x \\ \vdots \\ c_3 \end{bmatrix} \left\{ \begin{array}{l} \text{Anchor 1} \\ \text{Anchor 2} \end{array} \right.$$

each obj in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU with the ground truth bounding box.

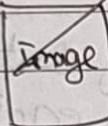
# low prob to bounding box are eliminated

# non max suppression used to generate final predictions.

## \* Transpose Convolution.



## \* U-Net Architecture



gradienții sunt în joc  
 Vor fi săptăni să se  
 trimit și la gradienții totul  
 să se săptăneze să se  
 săptăneze să se săptăneze  
 să se săptăneze să se săptăneze

$$\frac{2^7}{2^6} = 8$$

$$\frac{2^6}{2^5} = 8$$

$$\frac{2^5}{2^4} = 8$$

$$\frac{2^4}{2^3} = 8$$

$$\frac{2^3}{2^2} = 8$$

$$\frac{2^2}{2^1} = 8$$

$$\frac{2^1}{2^0} = 8$$

gradenții sunt săptănezi să se săptăneze

Lor săptănezi să se săptăneze să se săptăneze

## \* Face verification vs Face Recognition.

Verify: 1:1

Output is whether input is that of the claimed person.

Recognition: 1:k

To check against database of k persons. Output is ID of the person.

## \* Oneshot learning

↳ learn from a single picture

CNN don't work well bcz of small train set

So, a similarity function is used.

$d(1, 2)$  = degree of diff between 1 & 2

if  $d(1, 2)$  is small  $\rightarrow$  same person

if  $d(1, 2)$  is large  $\rightarrow$  diff person

$$(u, A)b \geq (q, A)b$$

$$0 \geq (u, A)b - (q, A)b$$

$$0 \geq x + (u, A)b - (q, A)b$$

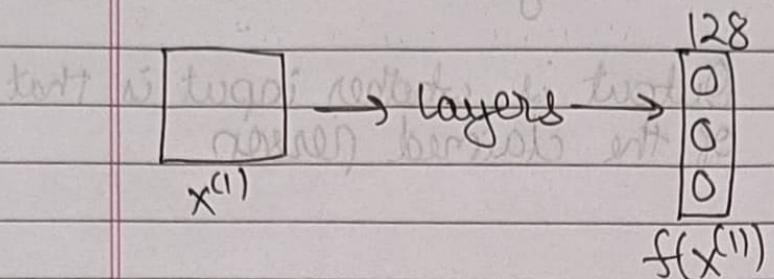
signum

sign of  $(u, A)$  treasury of

the set working more

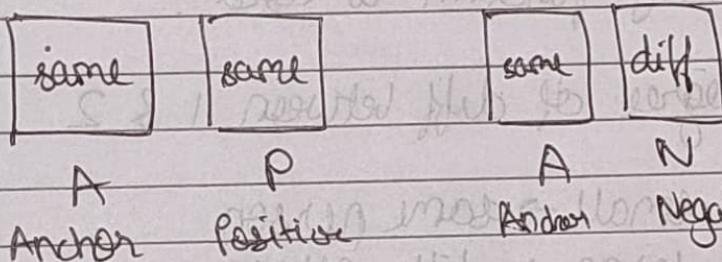
- steps;

## \* Siamese network



$$d(\cancel{x^{(1)}}, x^{(2)}) = \|f(x^{(1)}, x^{(2)})\|_2$$

## \* Triplet Loss.



$$d(A, P) \leq d(A, N)$$

$$d(A, P) - d(A, N) \leq 0$$

$$d(A, P) - d(A, N) + \alpha \leq 0$$

↓  
margin

α generally  
0.2

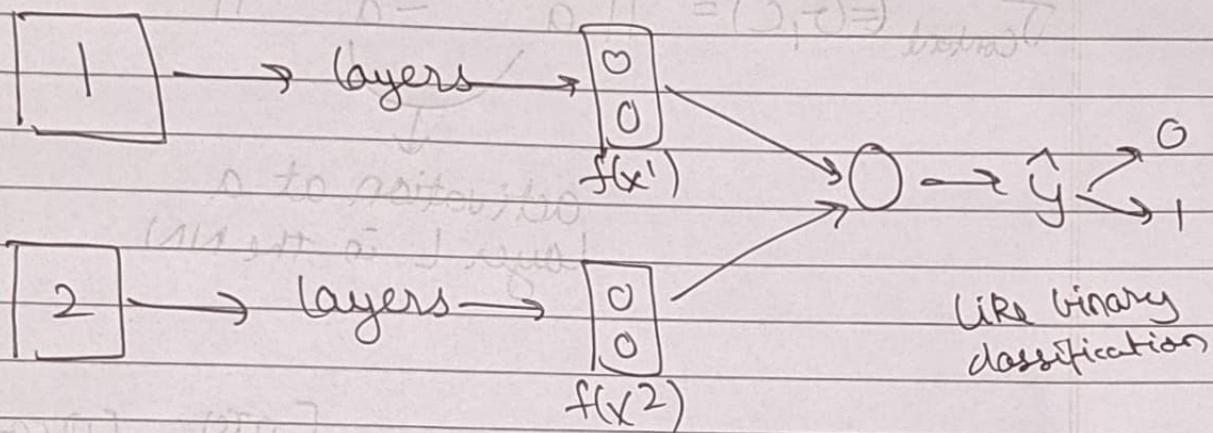
to prevent NN to give  
same output to all  
inputs.

Loss func:  $L(A, P, N) = \max(d(A, P) - d(A, N) + \alpha, 0)$

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

( $\Rightarrow$ )  $J$  minimum of training loss  $\Rightarrow$   $J$  minimum of training error

# Choose triplets that are hard to train bcz if random the eqn is easily satisfied



$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^1)_k - f(x^2)_k| + b \right)$$

## \* Neural style transfer

Content + Style = Generated image

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

SI: Randomly initialize a image of req dimension

SIT: Use gradient descent to minimize  $J(G)$ .

Content cost function:

$$J_{\text{content}}(G, C) = \|a^{(L)(C)} - a^{(L)(G)}\|^2$$

activation at a layer L in the NN

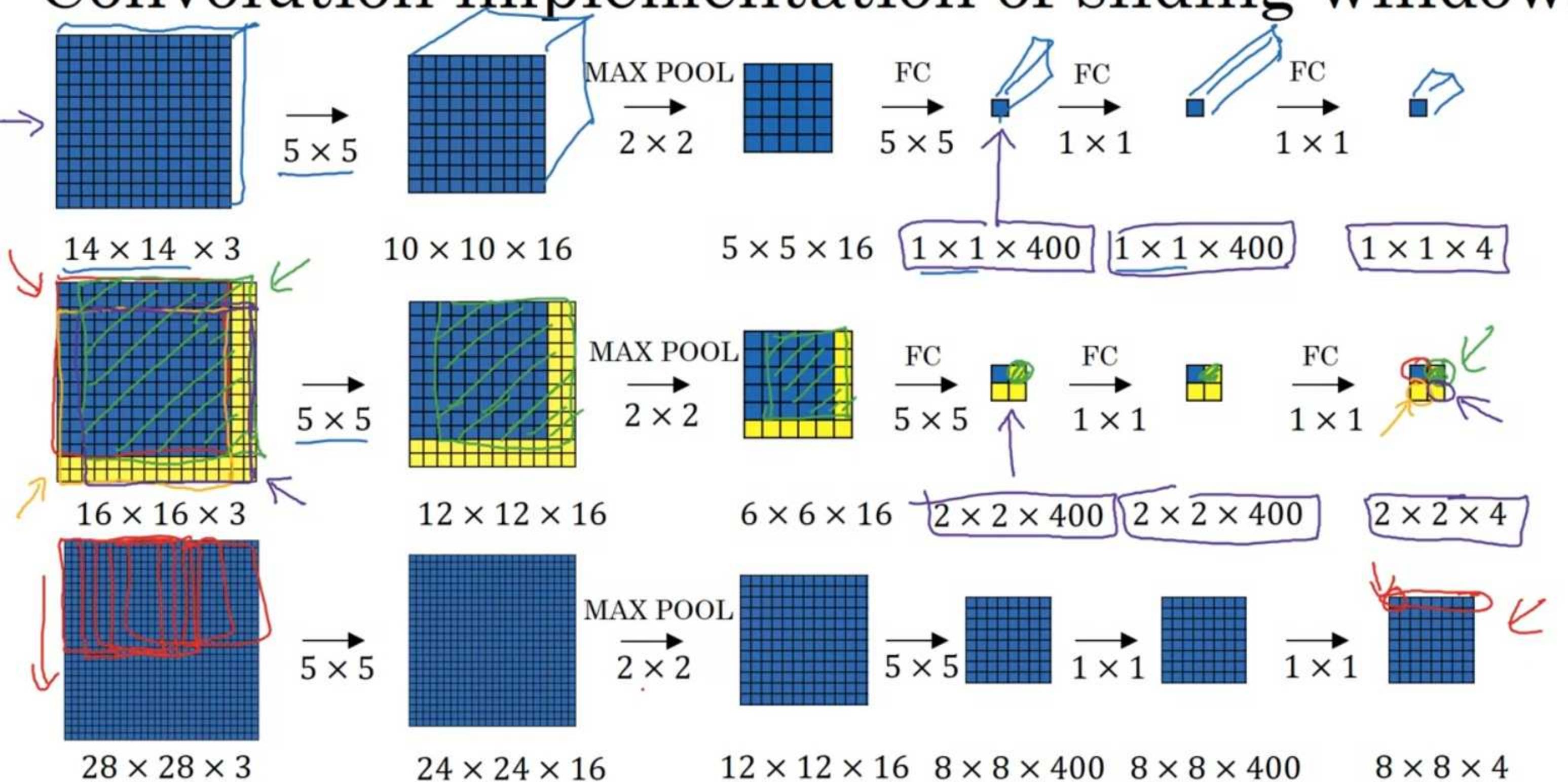
$$J_{\text{style}}(G, G) = \frac{1}{(2n_H n_W n_C)^2} \sum_k \sum_k [G_{kk'}^{(S)} - G_{kk'}^{(G)}]^2$$

$$G_{kk'}^{(S)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} a_{ijk} a_{ijk'}$$

$$G_{kk'}^{(G)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} F_{ij}(G) a_{ijk} a_{ijk'}$$

$$[(\partial, 2), \partial, T]q + [(-\partial), \partial, T]x = (-\partial)T$$

# Convolution implementation of sliding windows



# U-Net

