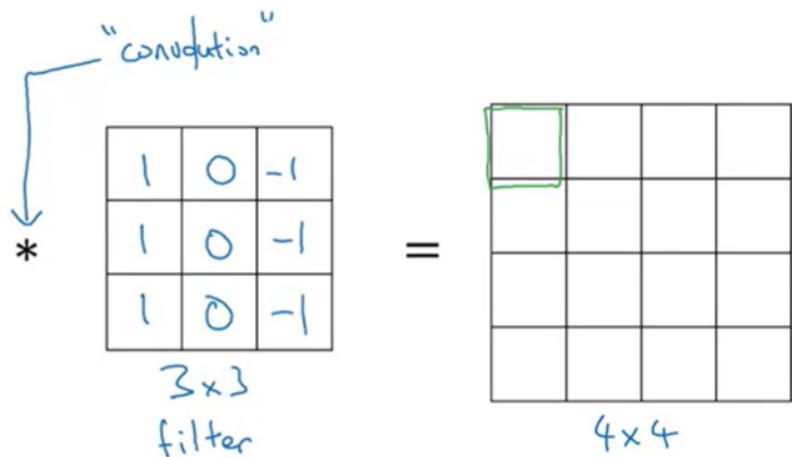


# Edge Detection

One of the applications is Edge Detection. Lets say we have a grayscale input matrix of 6x6 as input. We use convolution operation between this matrix and the filter matrix and a filter 3x3 matrix. We get a 4x4 convolution matrix as output. The calculation is as follows :

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

6x6



We choose the 3x3 section of first matrix, multiply each element with corresponding filter matrix element and find their sum. We fill the sum in the corresponding cell of the convolution matrix.

The filter matrices are :

## Vertical and Horizontal Edge Detection

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

Vertical

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

Horizontal

The intuition of edge detection is as follows :

The diagram illustrates the convolution process for edge detection. It consists of two rows of operations.

**Top Row:**

- An input image of size 6x6 with values [10, 10, 10, 0, 0, 0; ...; 0, 0, 0, 0, 0, 0]. An arrow points to a 3x3 kernel with values [1, 0, -1; 1, 0, -1; 1, 0, -1].
- The multiplication operation is indicated by an asterisk (\*).
- The result is an output image of size 4x4 with values [0, 30, 30, 0; ...; 0, 30, 30, 0]. A blue box highlights the bottom-right 2x2 submatrix [30, 30; 30, 0], which corresponds to the vertical edge in the input.

**Bottom Row:**

- An input image of size 6x6 with values [0, 0, 0, 10, 10, 10; ...; 0, 0, 0, 10, 10, 10]. An arrow points to the same 3x3 kernel.
- The multiplication operation is indicated by an asterisk (\*).
- The result is an output image of size 4x4 with values [0, -30, -30, 0; ...; 0, -30, -30, 0]. A black box highlights the bottom-right 2x2 submatrix [-30, -30; 0, 0], which corresponds to the horizontal edge in the input.

In first image, The 10s show white and 0s show black. There is a sharp vertical edge in the image. After convoluting we get the 4x4 matrix which shows presence of a vertical line in between representing an edge.

On comparing both images we can see that there is a white edge line in first while black edge line in second. This also helps us to find whether there is a light to dark or dark to light transition while going from left to right.

Some popular filters are :

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

Sobel filter

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

Scharr filter

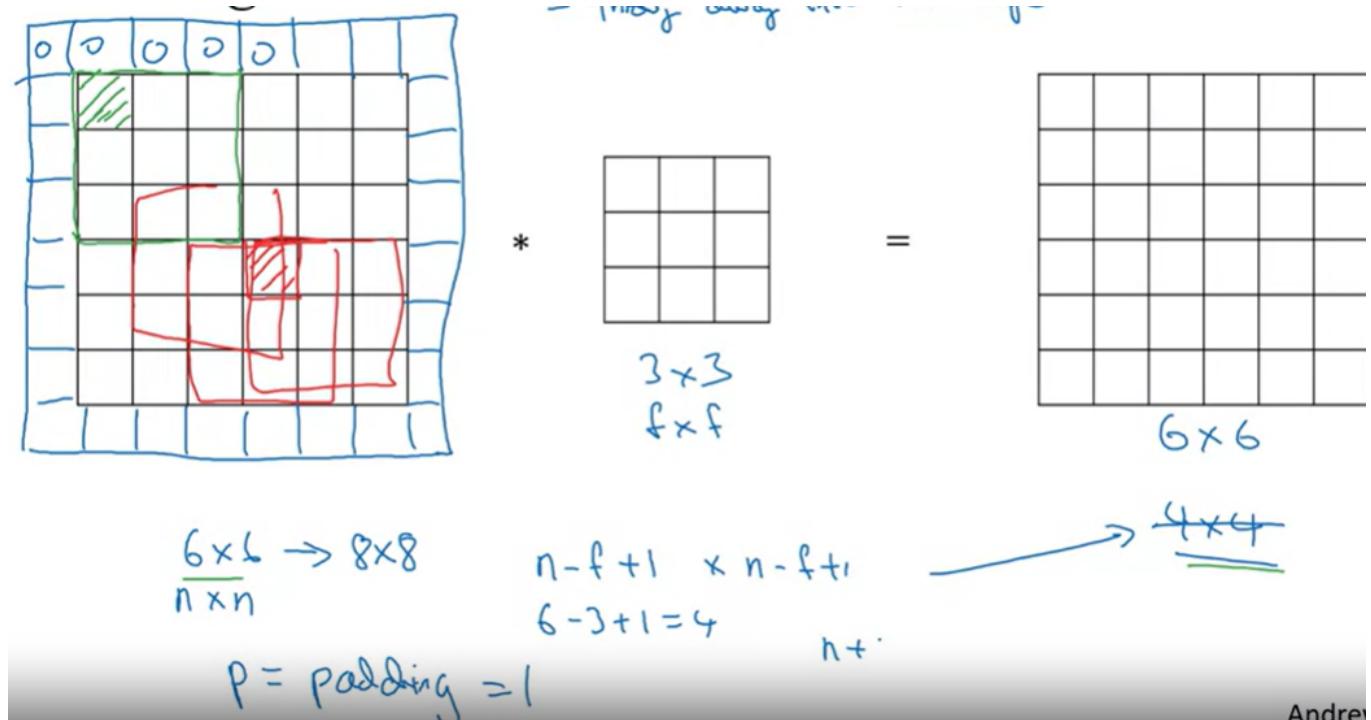
We generally treat the elements of filter as parameters and train them to get a robust model that can detect whatever edges we want.

Function in tensorflow for this purpose is :

"" in keras is :

## Padding

In the above edge detection , there is a shrinking in the above convoluted matrix, so to avoid that we add a layer of pixels around the image , which is called as padding



The dimensions of matrices can be characterised and studied as follows :

# Valid and Same convolutions

$\nearrow n \times \text{padding}$

“Valid”:  $n \times n$   $\times$   $f \times f$   $\rightarrow \frac{n-f+1}{f} \times \frac{n-f+1}{f}$

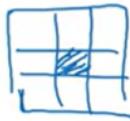
$6 \times 6$   $\times$   $3 \times 3$   $\rightarrow 4 \times 4$

“Same”: Pad so that output size is the same as the input size.

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

$f$  is usually odd

$3 \times 3 \quad p = \frac{3-1}{2} = 1 \quad \begin{matrix} S \times S \\ f \times f \end{matrix} \quad p=2$



Andrew Ng

## Strided Convolution

### Strided convolution

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	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

$\underline{7 \times 7}$

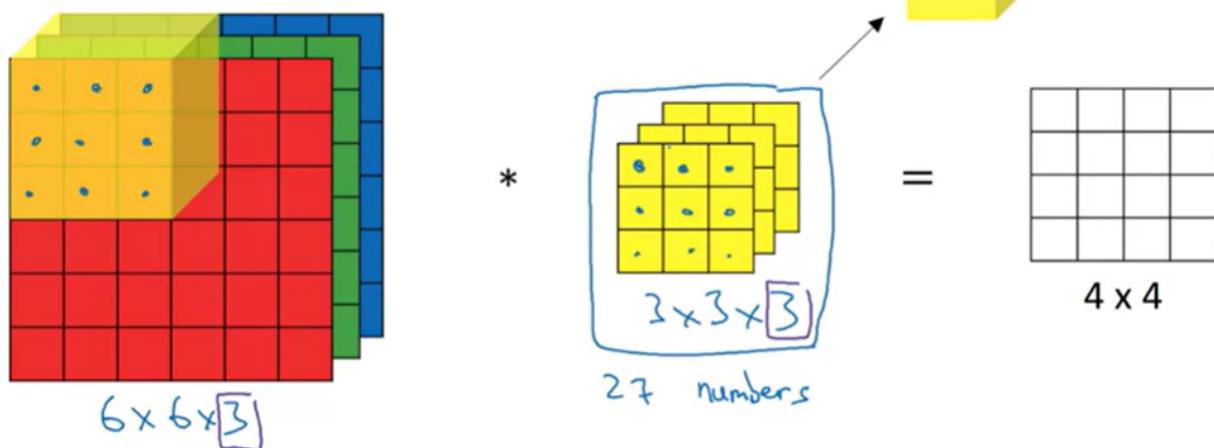
$$\begin{matrix} 3 & 4 & 4 \\ 1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix} \quad * \quad \begin{matrix} 91 & 100 & 83 \\ 69 & 91 & 127 \\ 44 & 72 & 74 \end{matrix} = \begin{matrix} 3 \times 3 \\ 3 \times 3 \end{matrix}$$

stride = 2

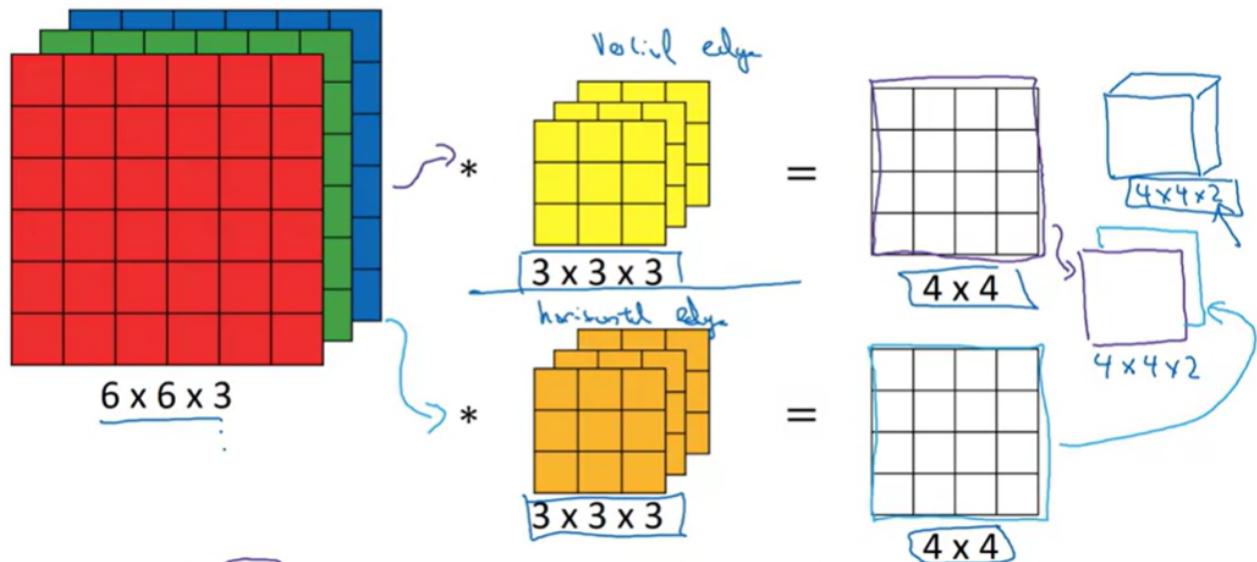
$n \times n$   $\times$   $f \times f$   
padding  $p$  stride  $s$   
 $s=2$

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

# Convolutions on RGB image



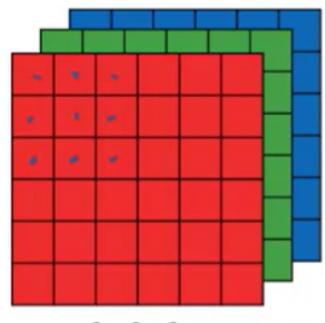
## Multiple filters



$$\text{Summary: } n \times n \times n_c \times f \times f \times n_c \rightarrow \frac{n-f+1}{4} \times \frac{n-f+1}{4} \times n'_c \quad \# \text{filters}$$

Andrew N

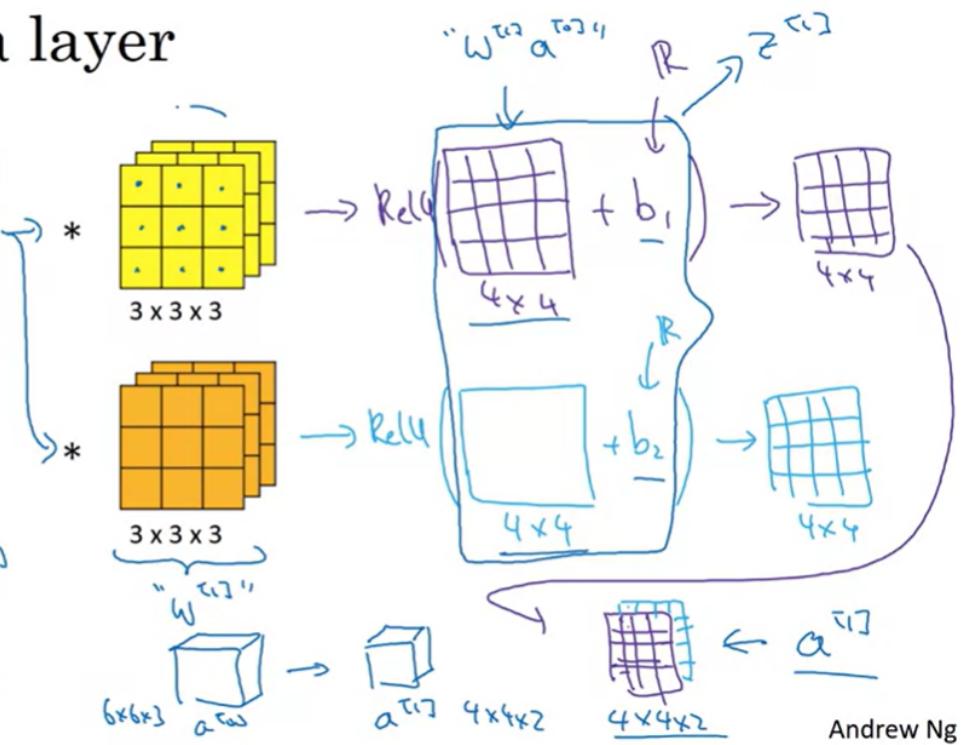
# Example of a layer



$$a^{l-1}$$

$$z^{l-1} = w^{l-1} a^{l-1} + b^{l-1}$$

$$a^l = g(z^l)$$



Andrew Ng

## Summary of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

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

→ Each filter is:  $f^{[l]} \times f^{[l]} \times n_c^{[l]}$

Activations:  $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

Weights:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias:  $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$  ↗ #f: filters in layer l.  $n_c^{[l]} \times n_H^{[l]} \times n_W^{[l]}$

$$\text{Input: } n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]} \leftarrow$$

$$\text{Output: } n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]} \leftarrow$$

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

$$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

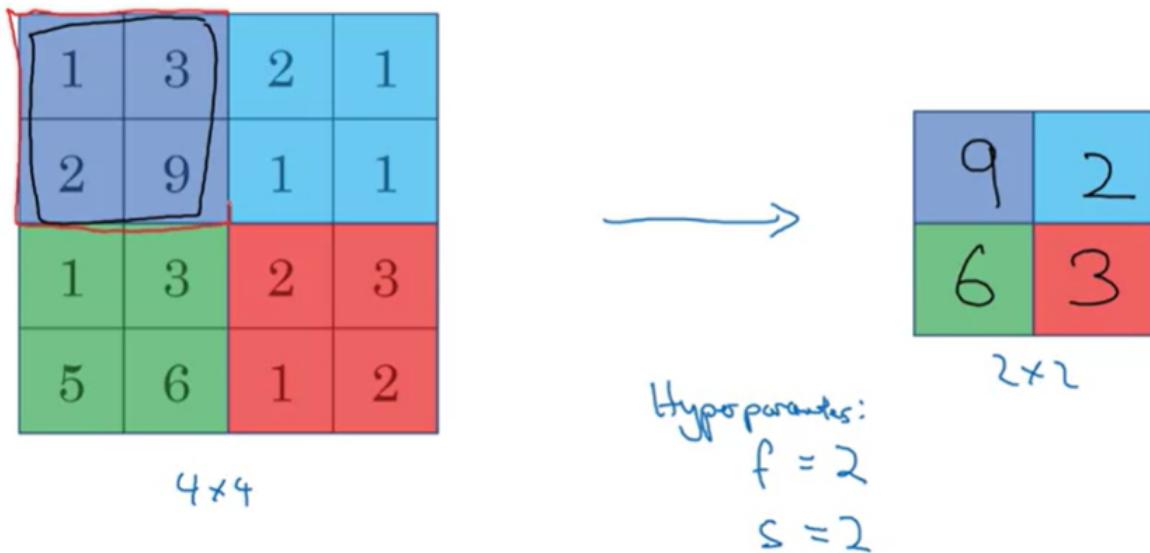
Depth of output volume = number of filters

## Types of layer in a convolutional network:

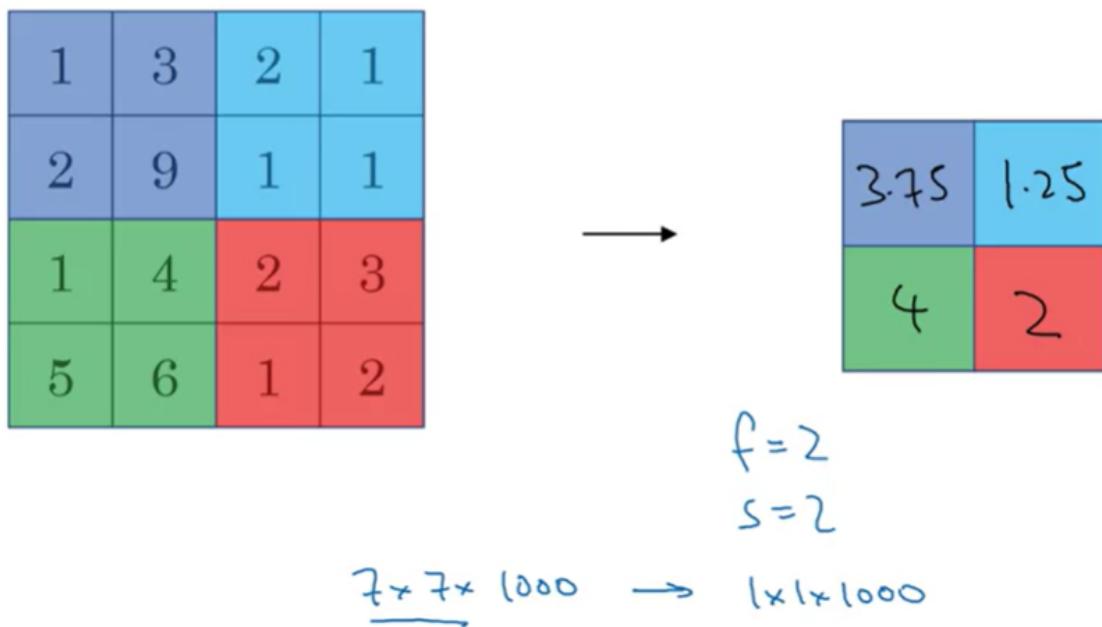
- Convolution (CONV) ←
- Pooling (POOL) ←
- Fully connected (Fc) ←

## Pooling layers

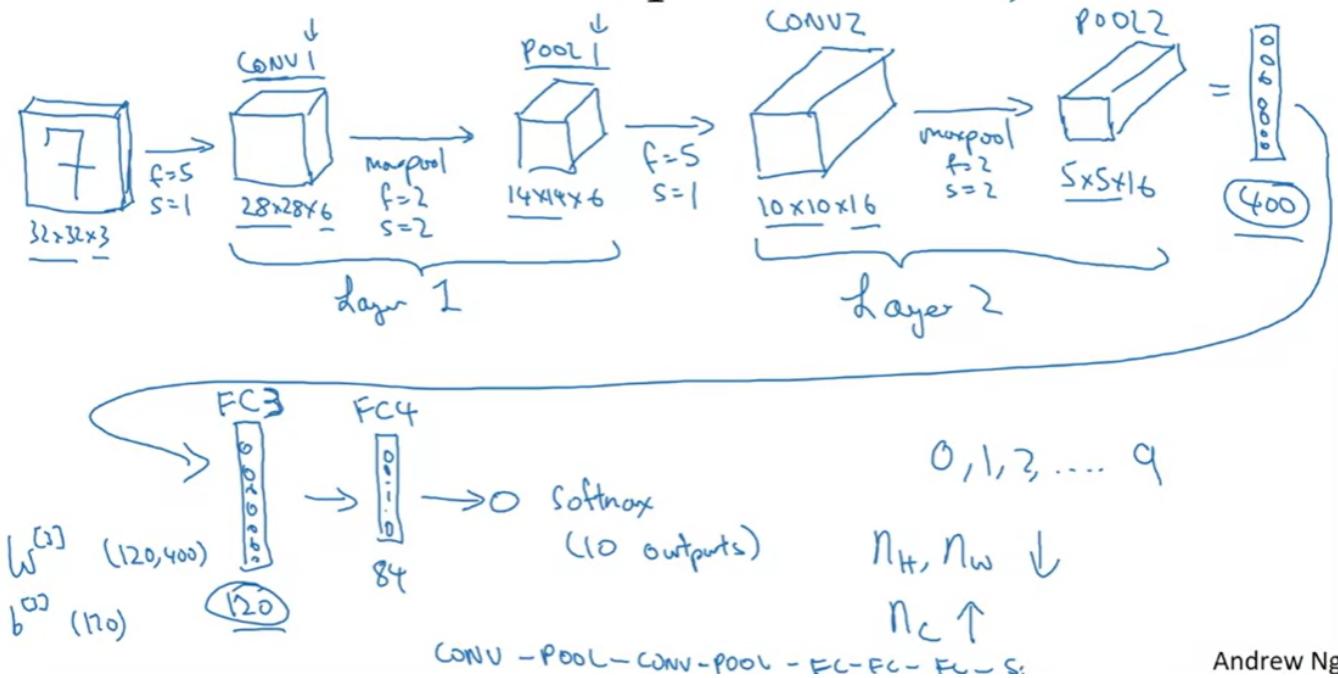
## Pooling layer: Max pooling



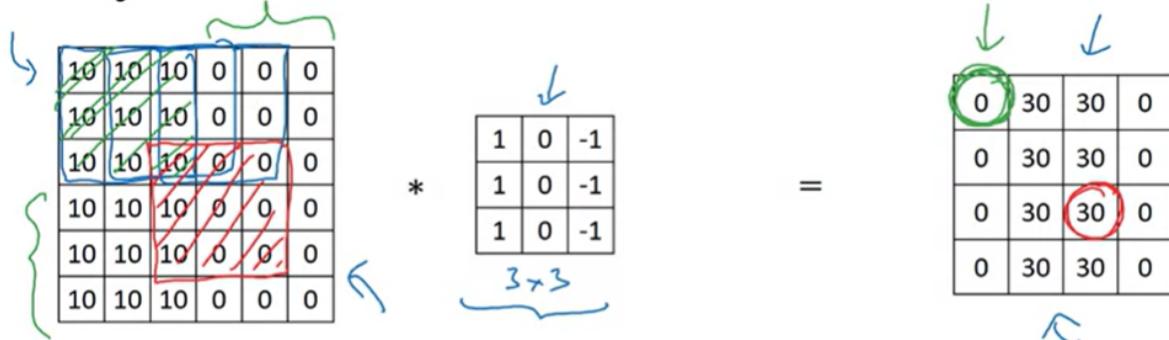
## Pooling layer: Average pooling



# Neural network example



## Why convolutions



**Parameter sharing:** A feature detector (such as a 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 only on a small number of inputs.

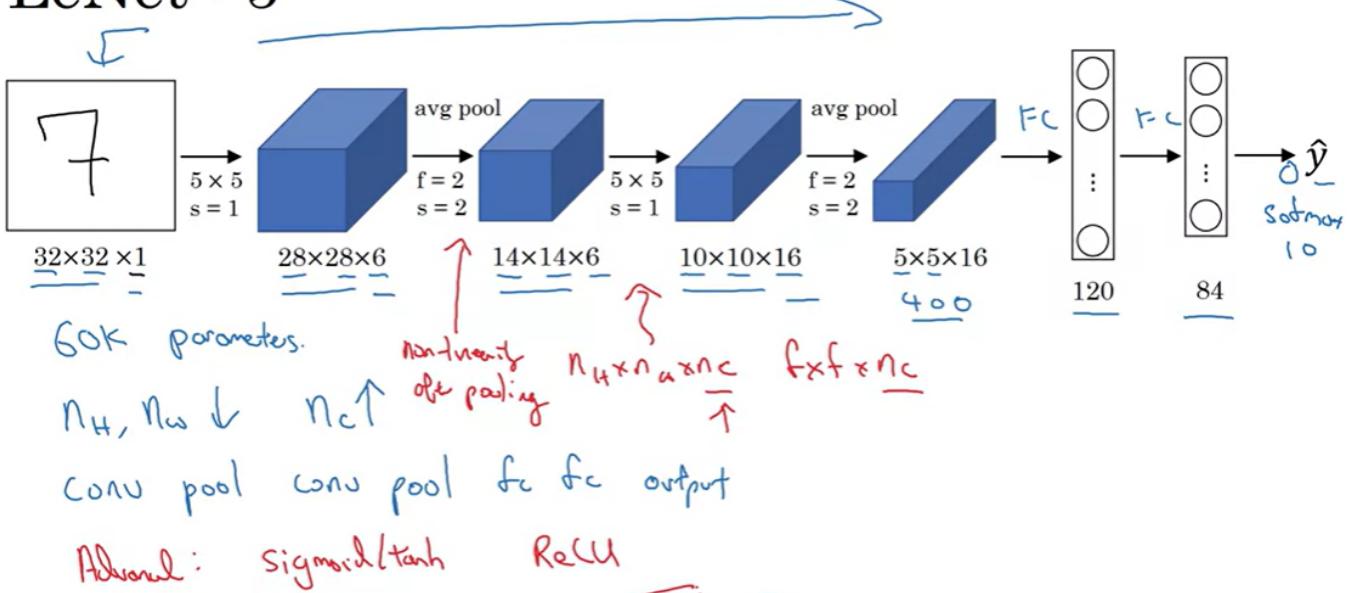
Andrew Ng

## Case Studies

Referring to case studies helps us in formulating our own code.

Classic networks are :

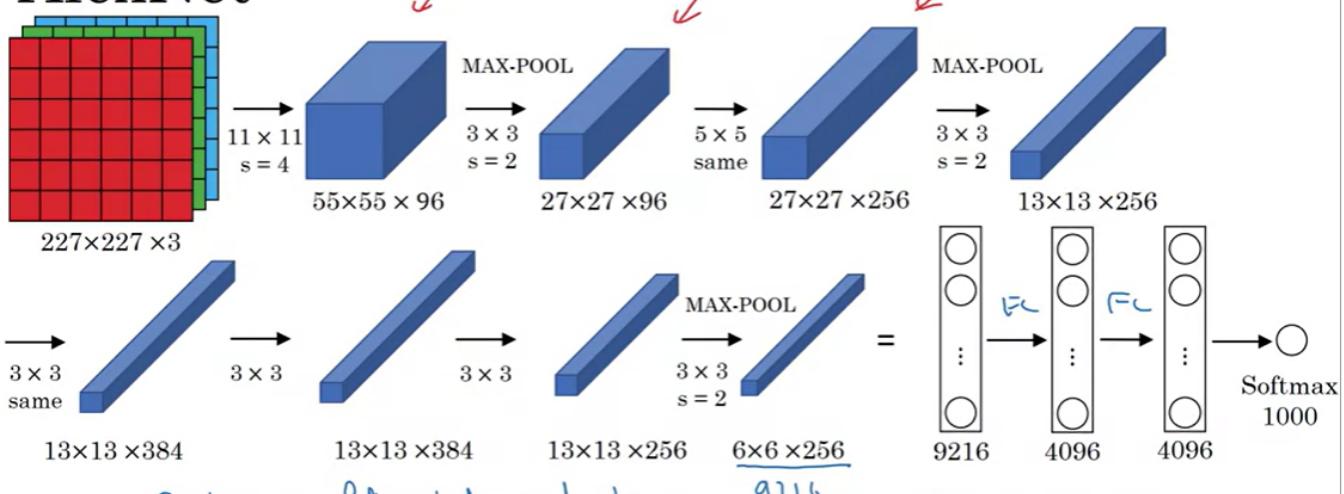
# LeNet - 5



II

III.

## AlexNet

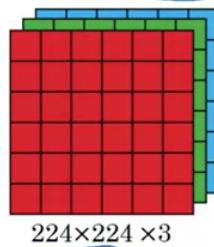


- Similar to LeNet, but much bigger.
- ReLU
- Multiple GPUs.
- Local Response Normalization.

Andrew Ng

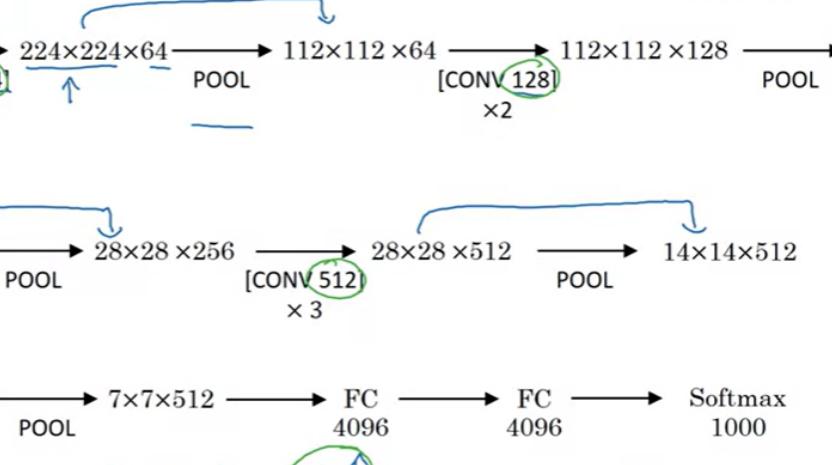
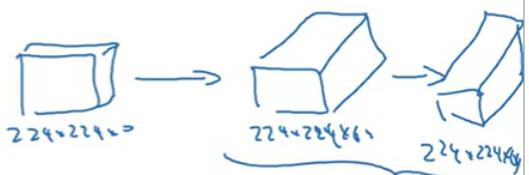
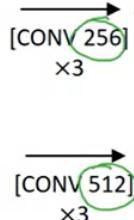
# VGG - 16

CONV = 3x3 filter, s = 1, same



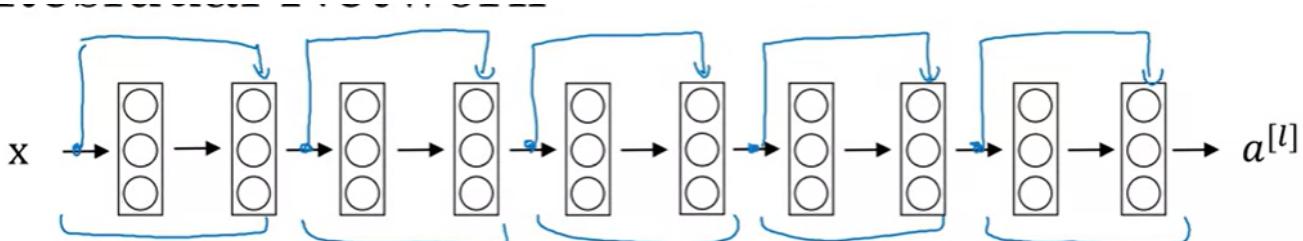
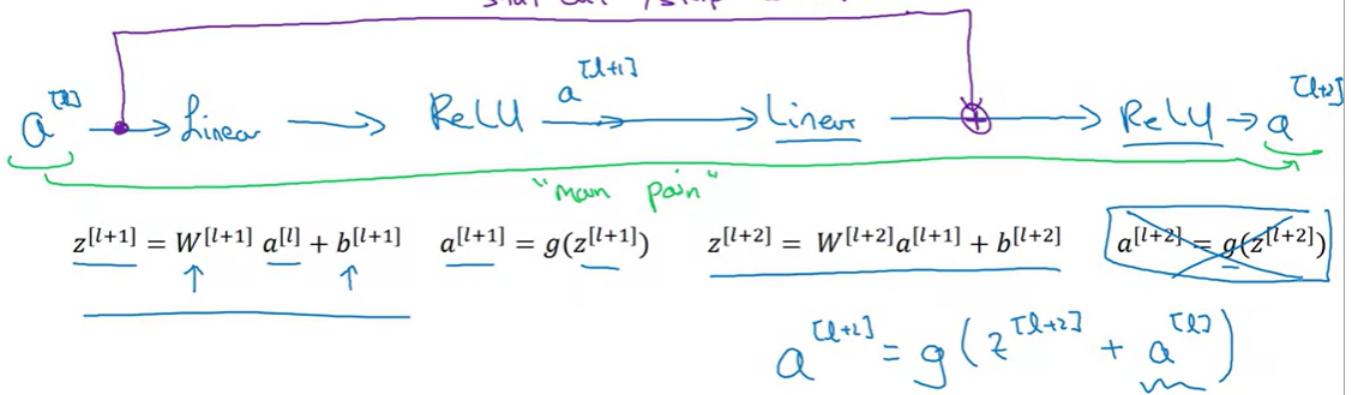
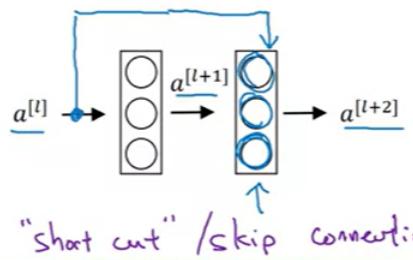
# VGG-19

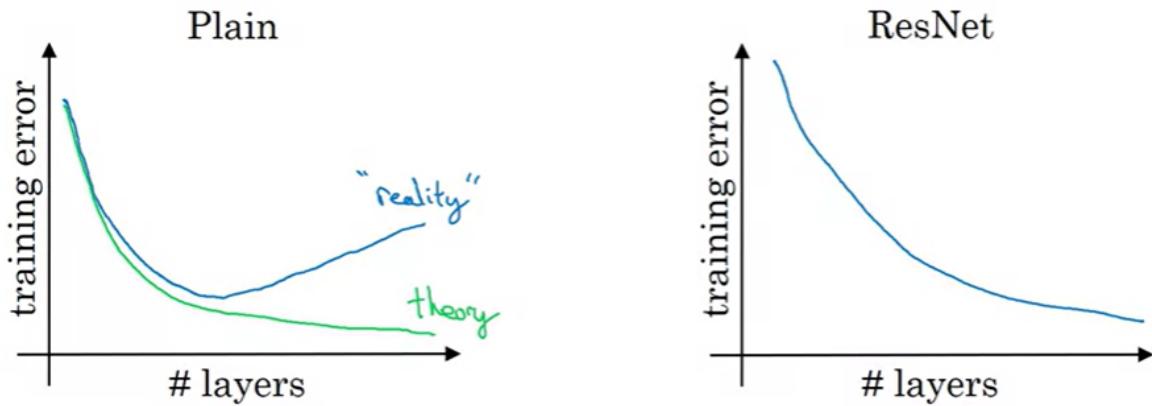
MAX-POOL = 2x2, s = 2



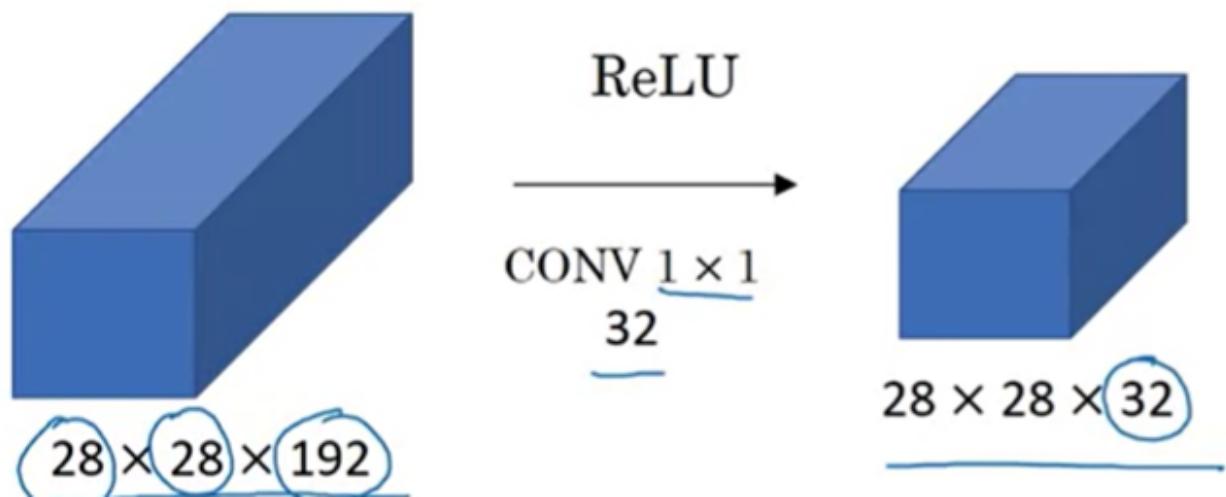
## Residual Networks

### Residual block



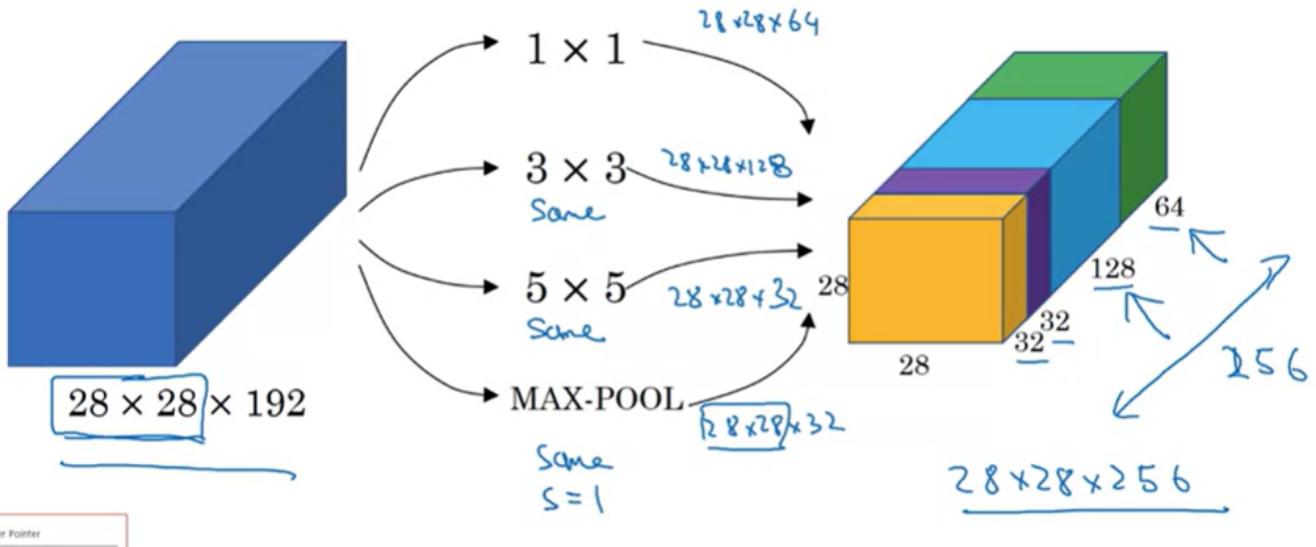


$1 \times 1$  convolutions help to shrink the number of channels if they have gotten too large. It works as follows :



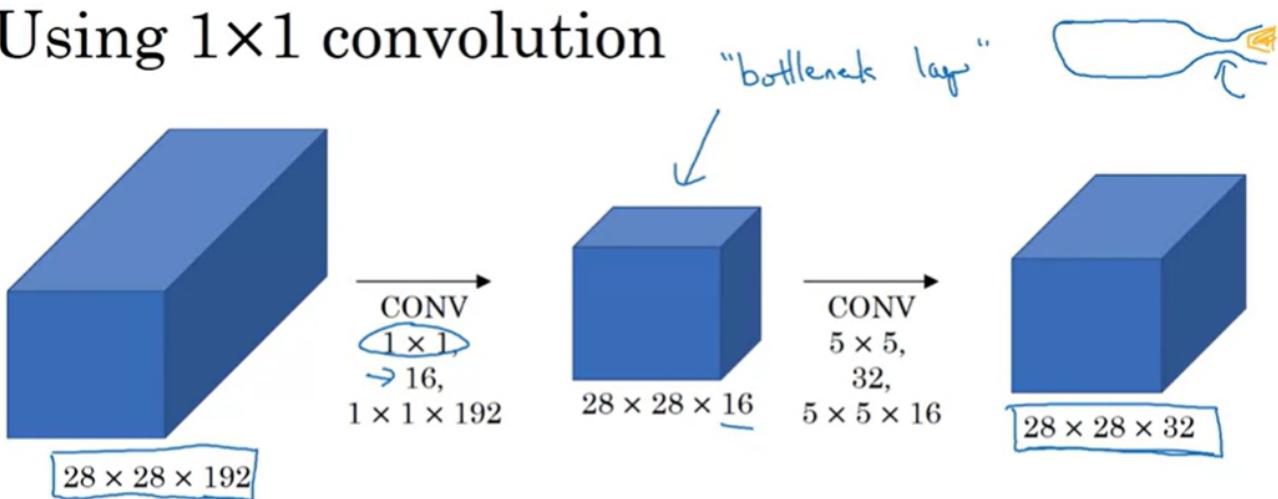
## Inception Network

When we are unsure as to what filters and pooling layers to use, we perform all the possibilities at the same time , condense results into one block and concatinate the results



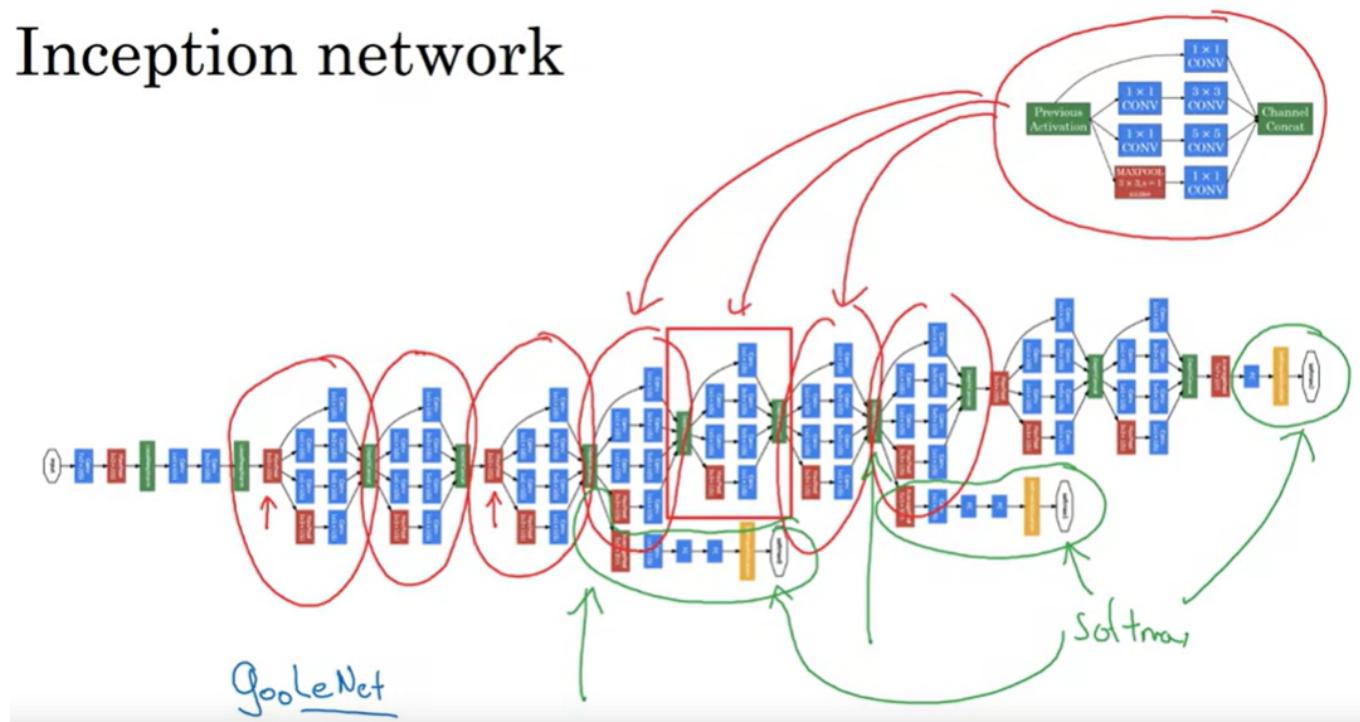
However we run into the problem of high computational cost in for eg  $5 \times 5$  layer above. So we use  $1 \times 1$  convolutions with an intermediate bottleneck layer as shown below to reduce the computational cost

## Using $1 \times 1$ convolution



Example of an inception network

# Inception network



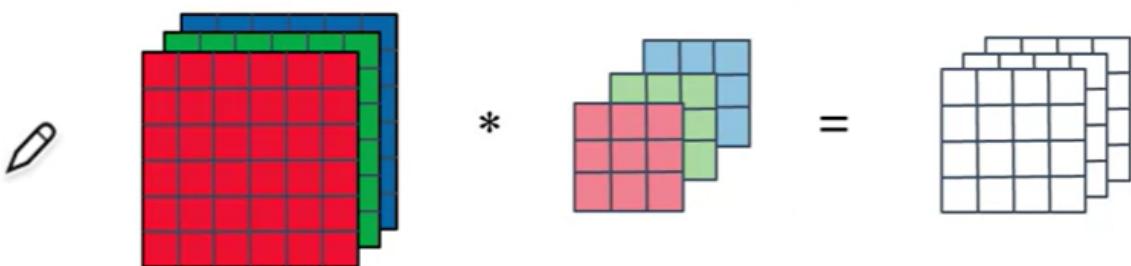
# MobileNets

# Motivation for MobileNets

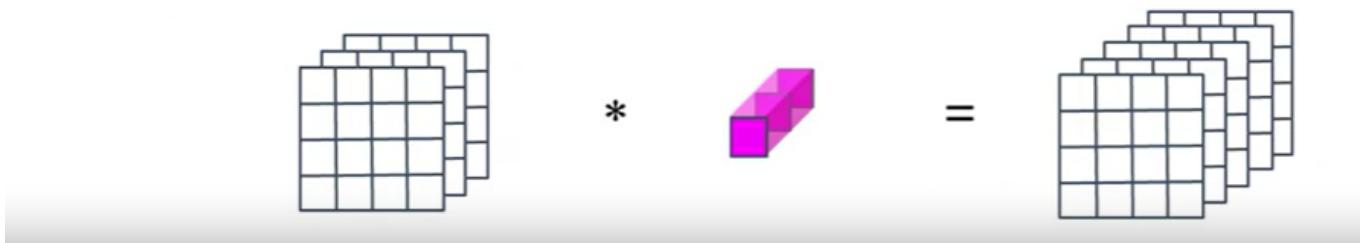
- Low computational cost at deployment
- Useful for mobile and embedded vision applications
- Key idea: Normal vs. depthwise-separable convolutions

## Depthwise Separable Convolution

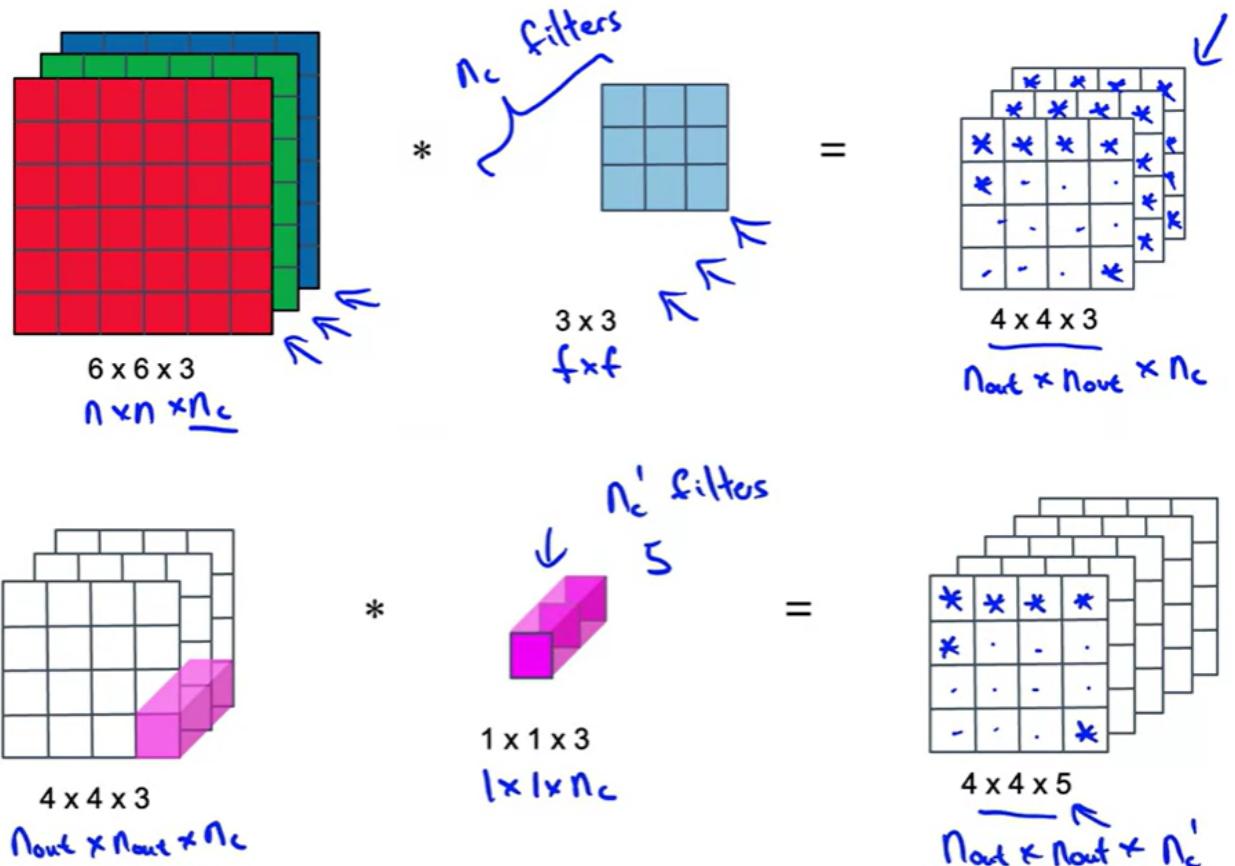
Depthwise Convolution



Pointwise Convolution



# Depthwise Convolution

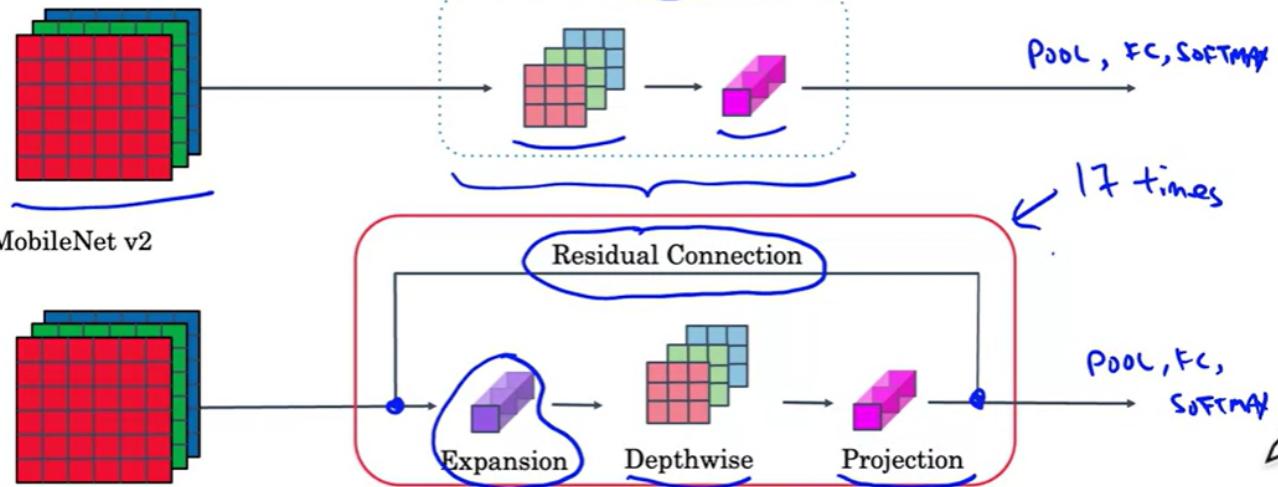


In depthwise convolution we convolve the individual coloured layers with corresponding same coloured filters to get an output . It is observed that this convolution has a much lower cost than normal convolution

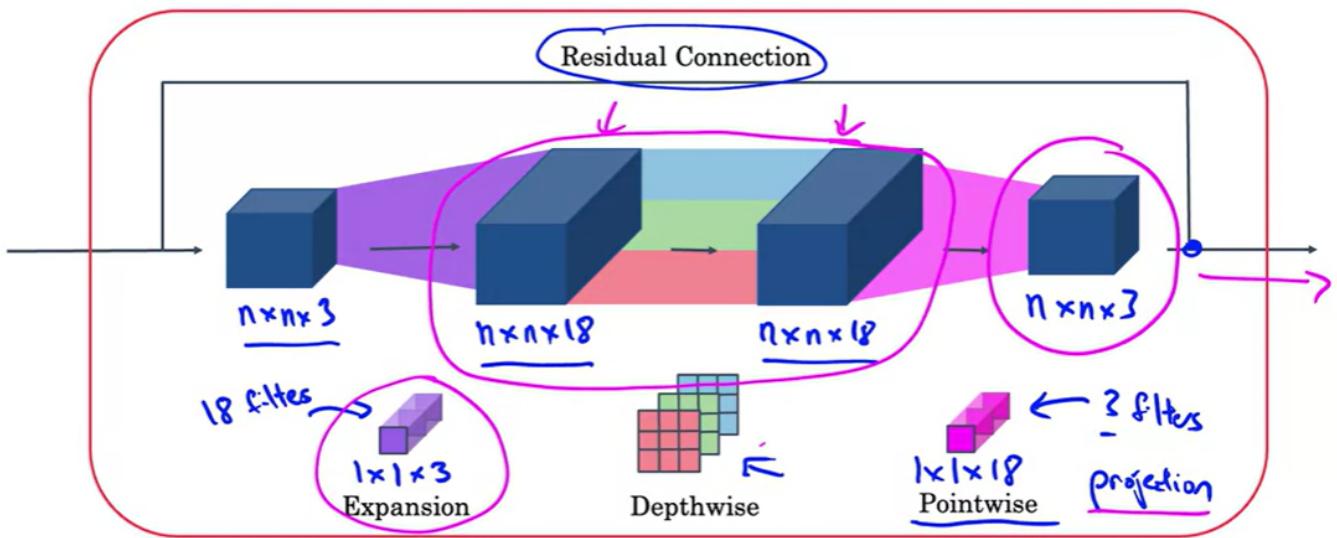
# MobileNet

MobileNet v1

MobileNet v2



## MobileNet v2 Bottleneck



## EfficientNet

Efficiency of a cnn can be improved by

- Increasing resolution of input image
- Increasing depth of the network
- Vary the width of individual layers

## Transfer learning

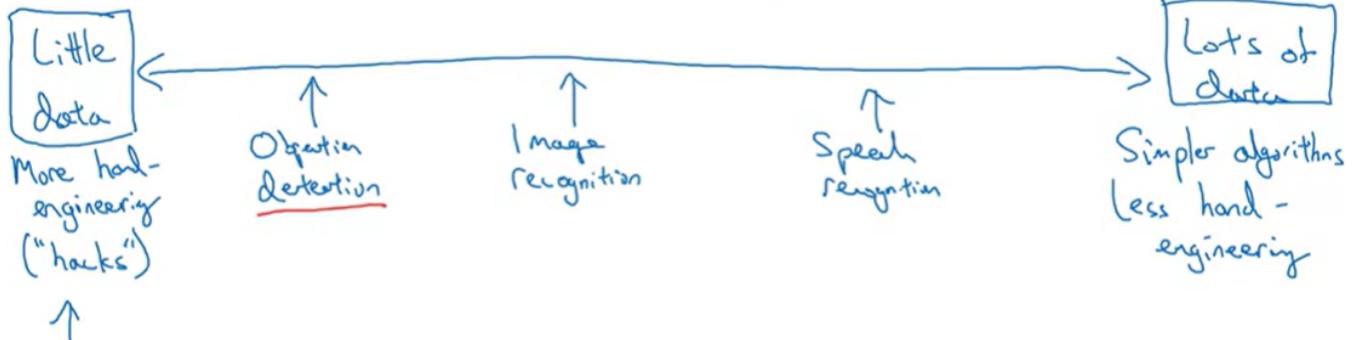
We import pretrained weights of most of layers of network by someone else and freeze those parameters. We only work in the last layers iterating their values. This is generally preferred in

computer vision

## Data Augmentation

Many a times in computer vision we come across a scarcity of training examples. So to tackle this we perform various operations on available set of data. This includes mirroring , random cropping , and using colour distortions(adding or subtracting RGB values)

## Data vs. hand-engineering



## Tips for doing well on benchmarks/winning competitions

### Ensembling

3 - 15 networks

→  $y$

- Train several networks independently and average their outputs

### Multi-crop at test time

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

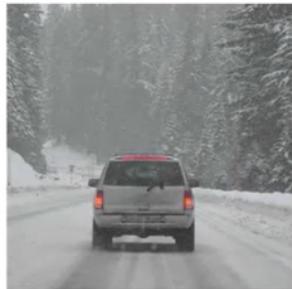
10-crop



It is better to use open-source, pretrained neural networks and finetune according to your requirements

## Object Detection

Image classification



"Car"

1 object

Classification with localization



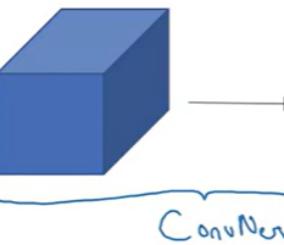
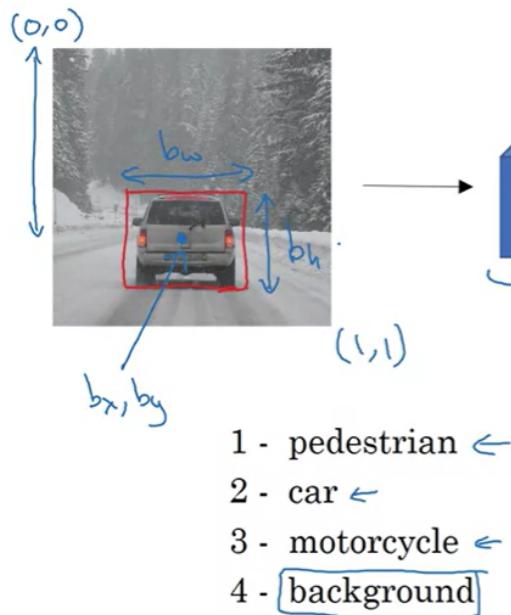
"Car"

Detection

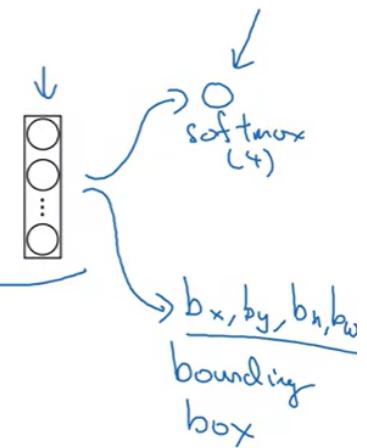


multiple objects

## Classification with localization



$$\begin{aligned} b_x &= 0.5 \\ b_y &= 0.7 \\ b_h &= 0.3 \end{aligned}$$



We can define the centre of object and its dimensions as above.

## Defining the target label $y$

1 - pedestrian  
 2 - car  
 3 - motorcycle  
 4 - background

Need to output  $b_x, b_y, b_h, b_w$ , class label (1-4)



$x =$



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

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

is there obj?

$(x, y) \leftarrow \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$

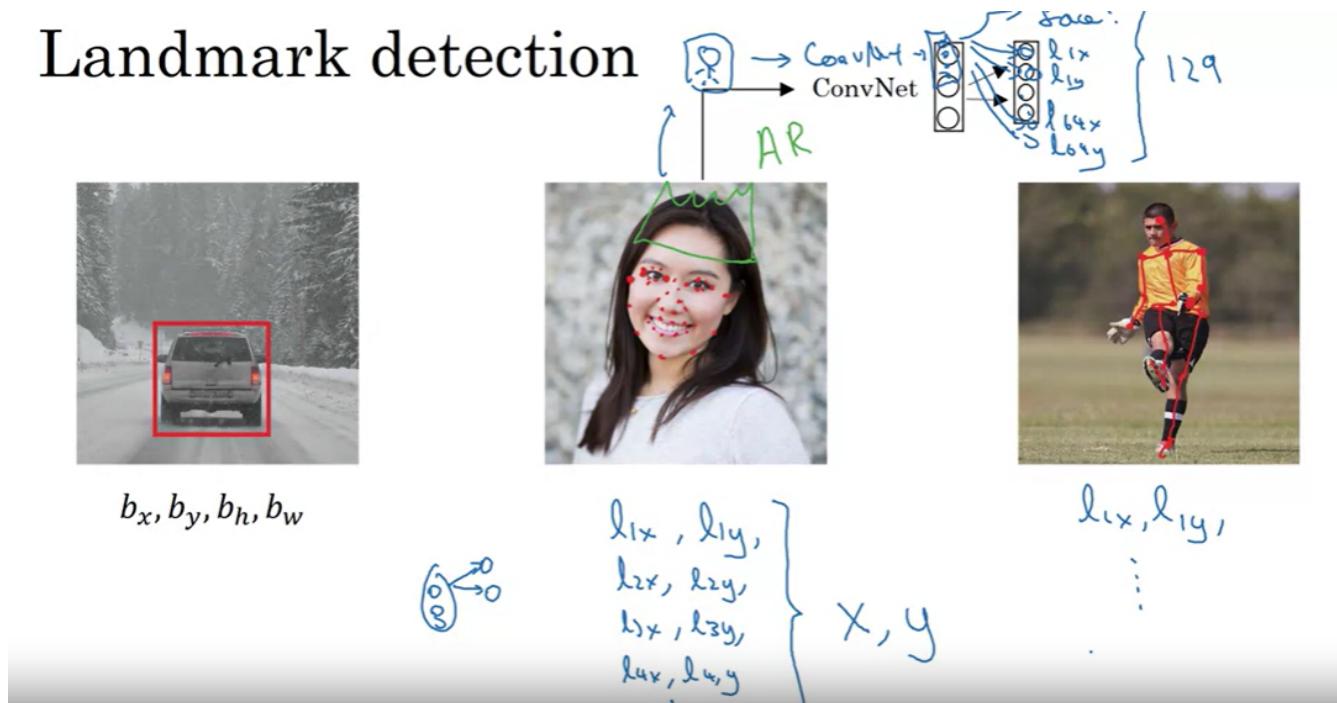
$\left[ \begin{array}{c} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{array} \right]$

← "don't care"

Andrew Ng

We can define the  $y$  vector as above, with  $p_c$  being the probability of an object being in the image. We don't care about other parameters if there is no object, hence the loss function is declared accordingly

## Landmark detection

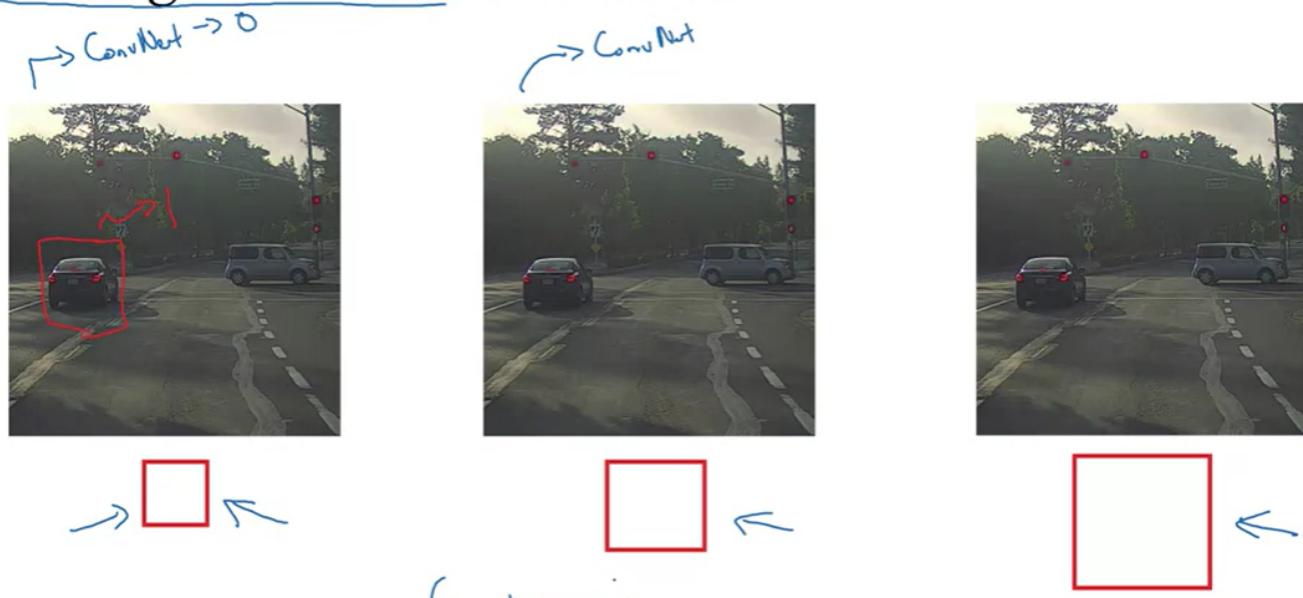


Certain points can be detected . Care has to be taken while labelling

We will first train our model on various car images. Then we use sliding windows method as shown below where we take boxes and pass them over the images and seeing if the pixels

enclosed each time contains a car or not. we will gradually go on increasing size of the boxes

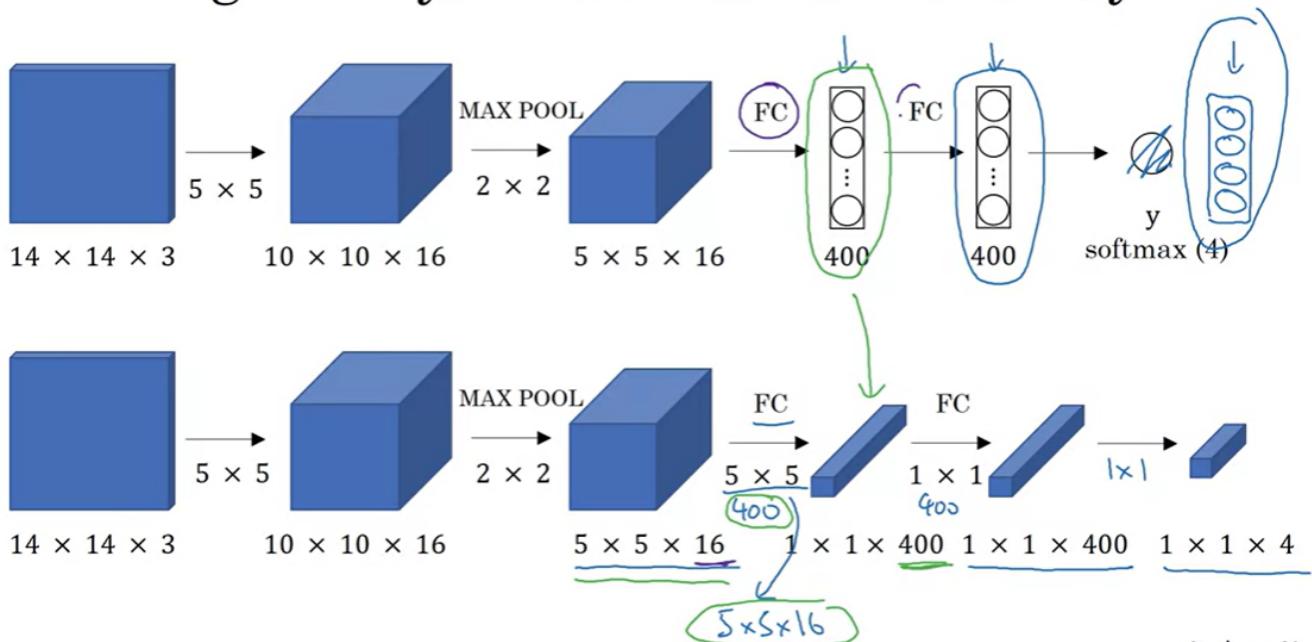
## Sliding windows detection



However if the stride is less then we will get more accuracy but computational cost will also increase

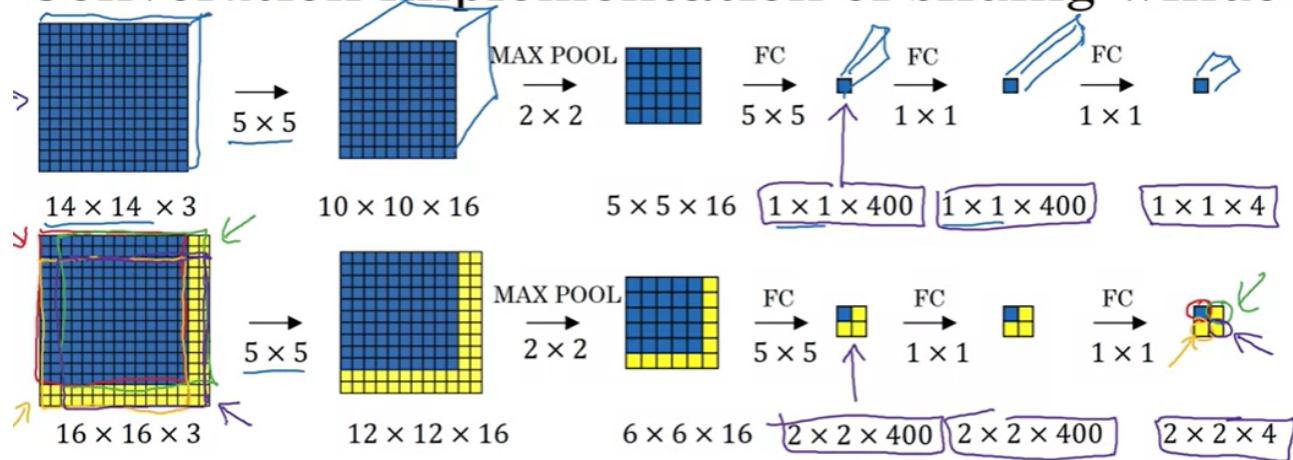
TO implement this in code,

## Turning FC layer into convolutional layers



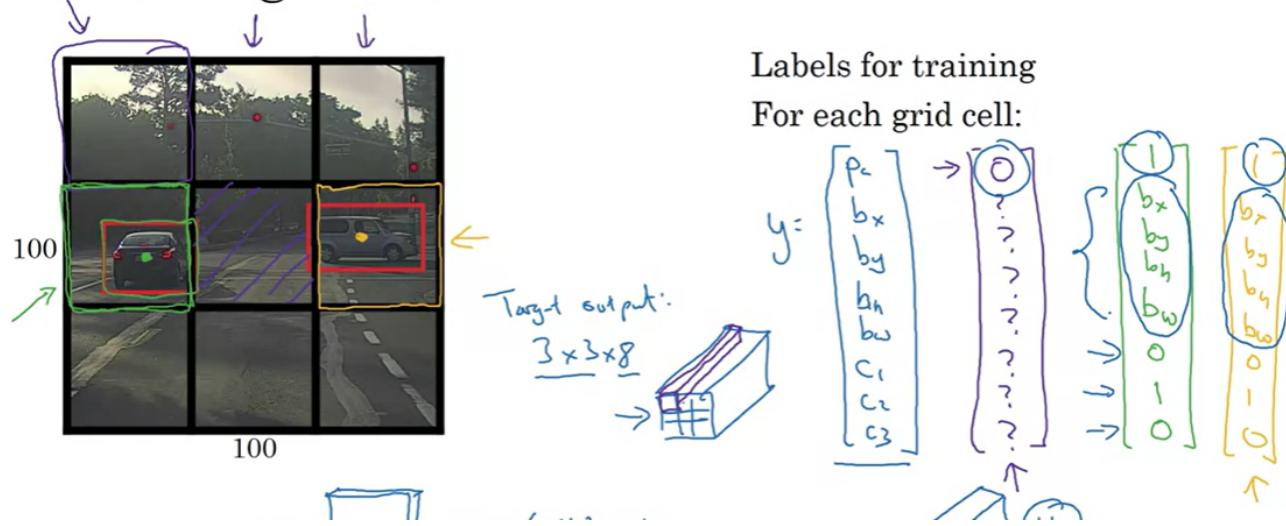
We convert a fc layer to conv layer

## Convolution implementation of sliding windows



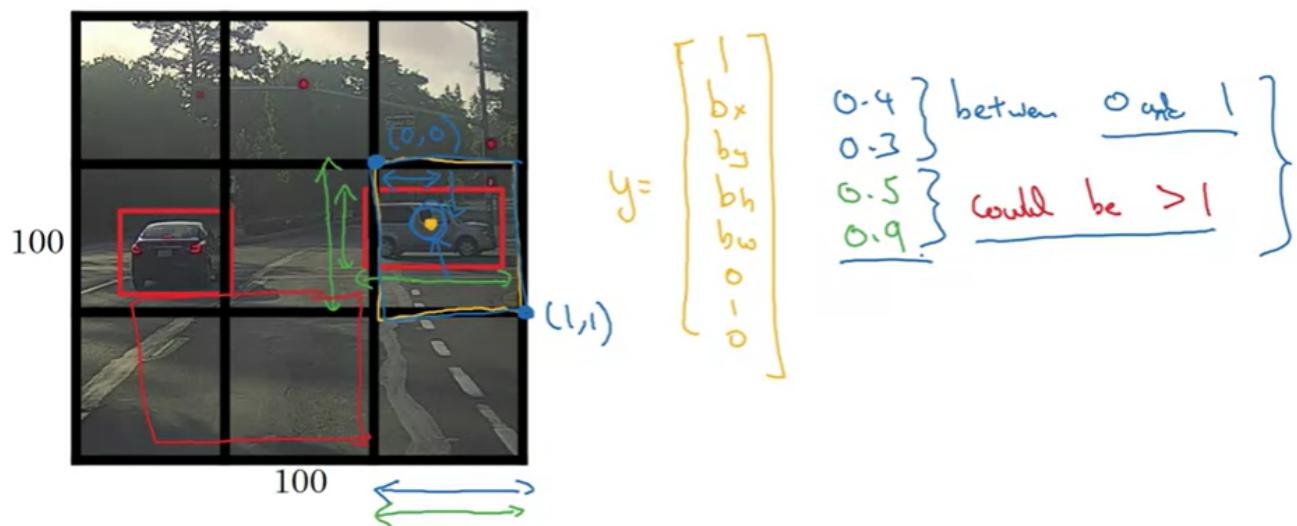
The blue region's output is the final blue top left cell. This helps to compress all the computation into one process. Thus individual sliding window's outputs can be found

## YOLO algorithm



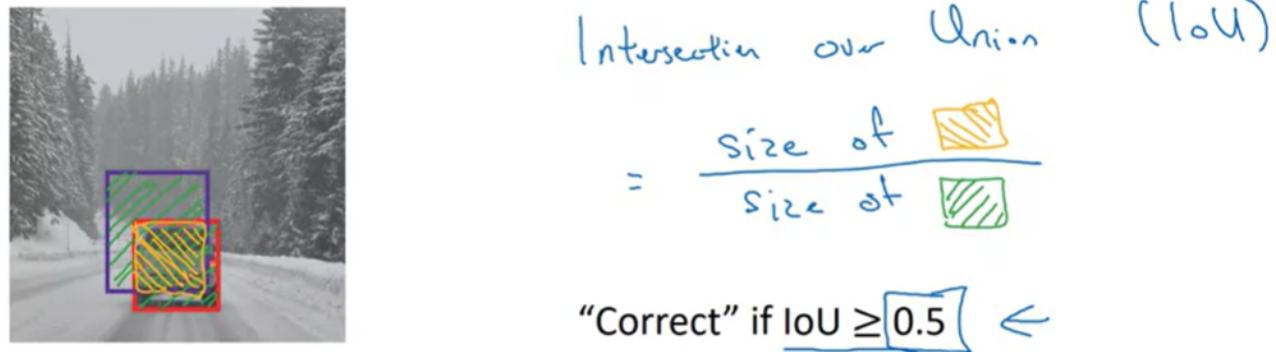
To improve the accuracy and speed, we draw a grid across the input images and perform object detection in each cell. This the output vectors are as shown above in right

# Specify the bounding boxes

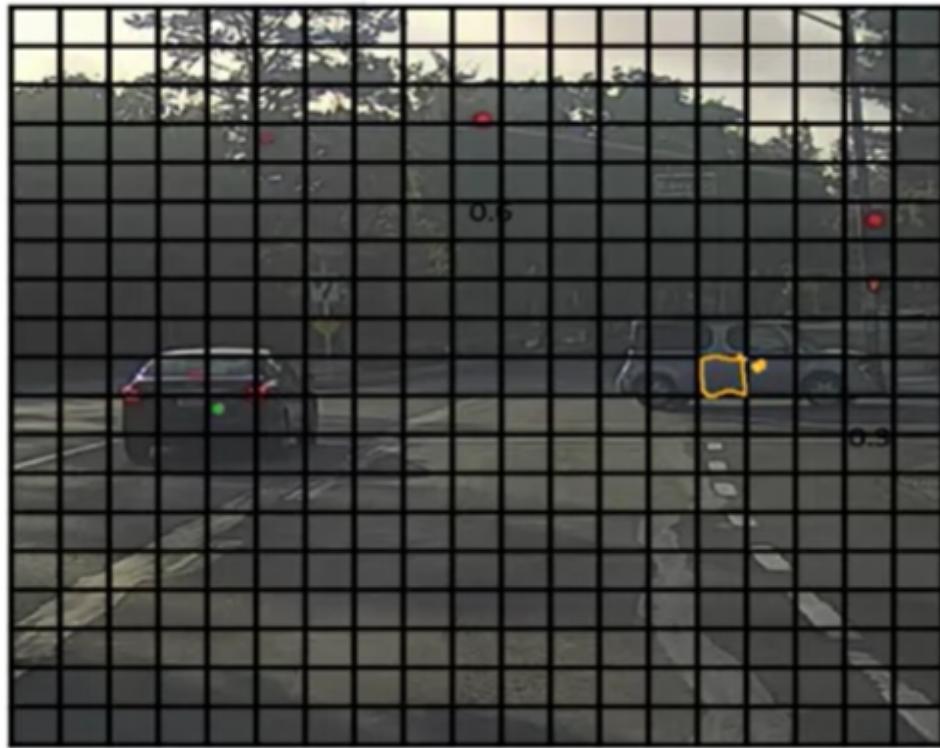


Higher the IoU, more accurate are the bounding boxes

## Evaluating object localization

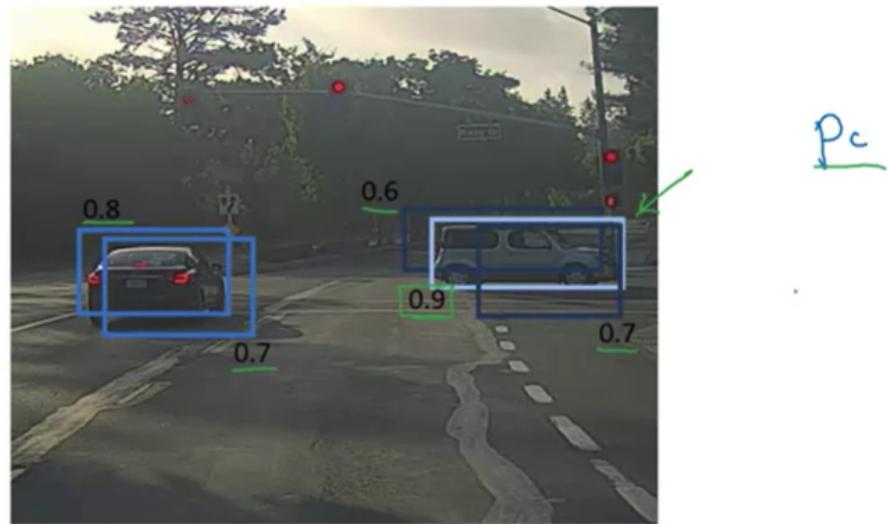


## Non-max suppression



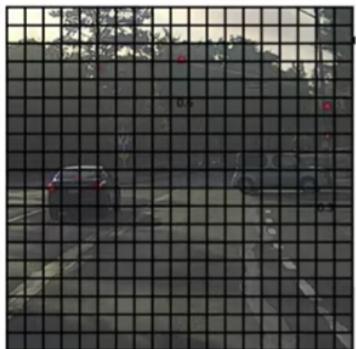
Here all the boxes having car might give multiple detection of car. So what nonmax suppression does is :

## Non-max suppression example



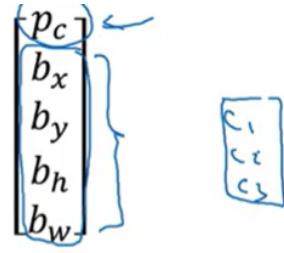
The prediction with highest IoU will be selected.

The algorithm to be followed is as follows :



19 × 19

Each output prediction is:



Discard all boxes with  $p_c \leq 0.6$

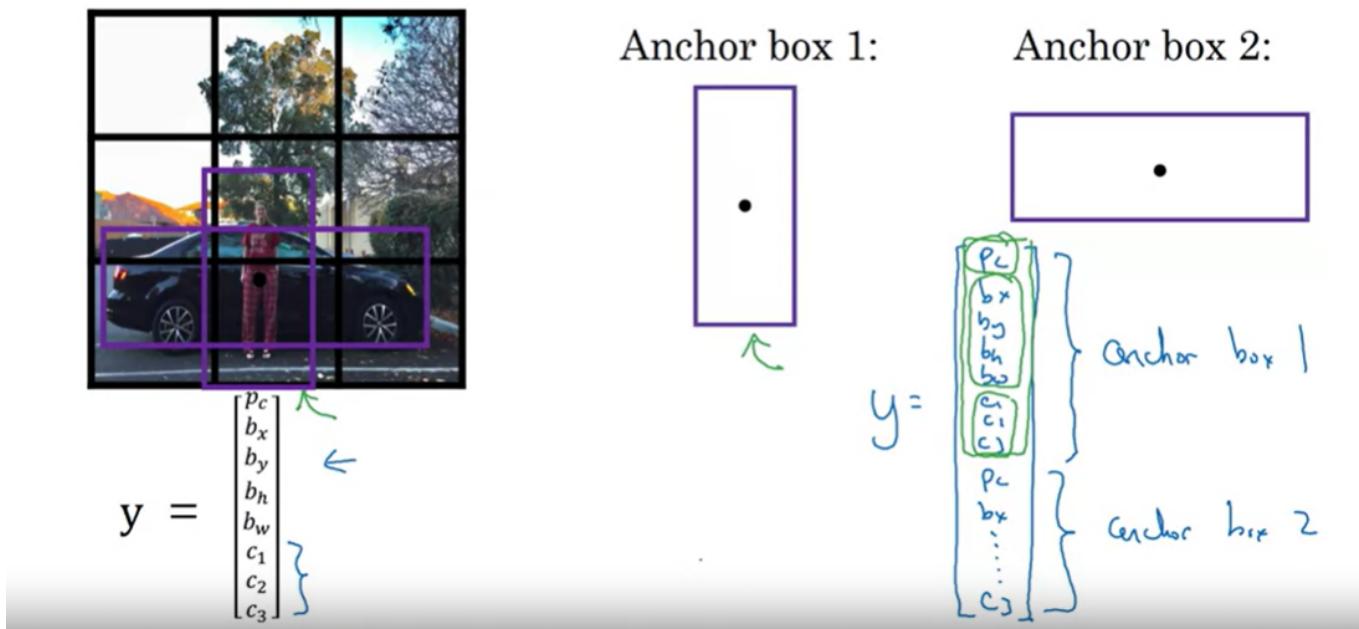
→ While there are any remaining boxes:

- Pick the box with the largest  $p_c$ . Output that as a prediction.
- Discard any remaining box with  $\text{IoU} \geq 0.5$  with the box output in the previous step

Andrew Ng

## Anchor Boxes

# Overlapping objects:

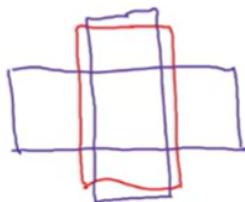


## Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output  $y$ :  
 $3 \times 3 \times 8$



With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

(grid cell, anchor box)

Output  $y$ :  
 $3 \times 3 \times 16$   
 $3 \times 3 \times 2 \times 8$

Andrew N

# Outputting the non-max suppressed outputs

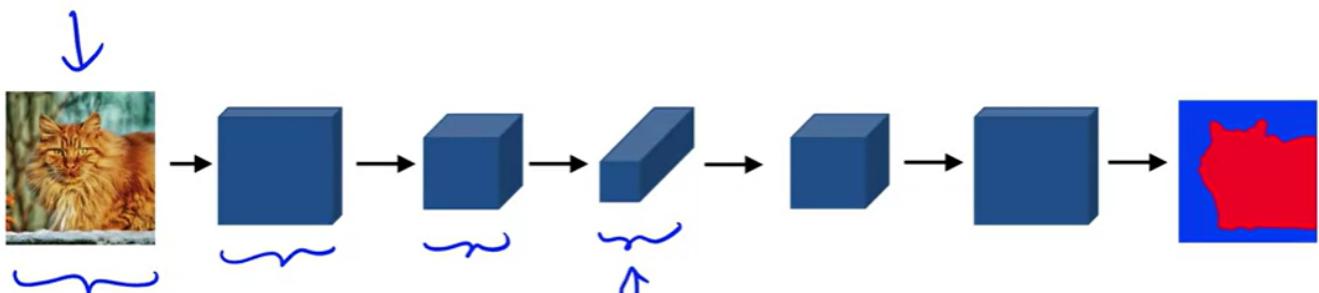


- For each grid call, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

## Semantic Segmentation

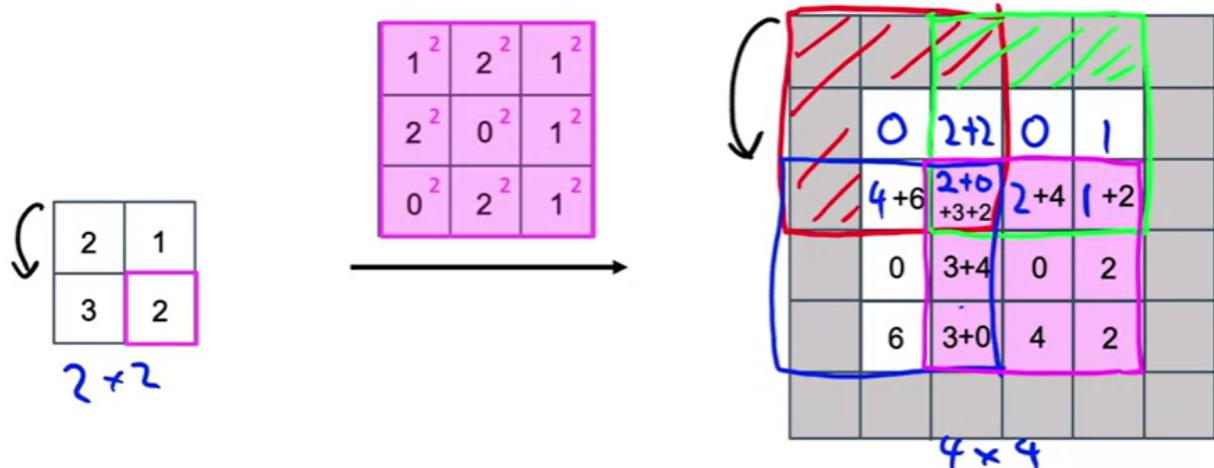
Here we label various pixel regions with different numbers. For e.g. car as 1 , road as 2, sky as 3. This can help to identify for eg. drivable road for self driving cars

Algorithm :



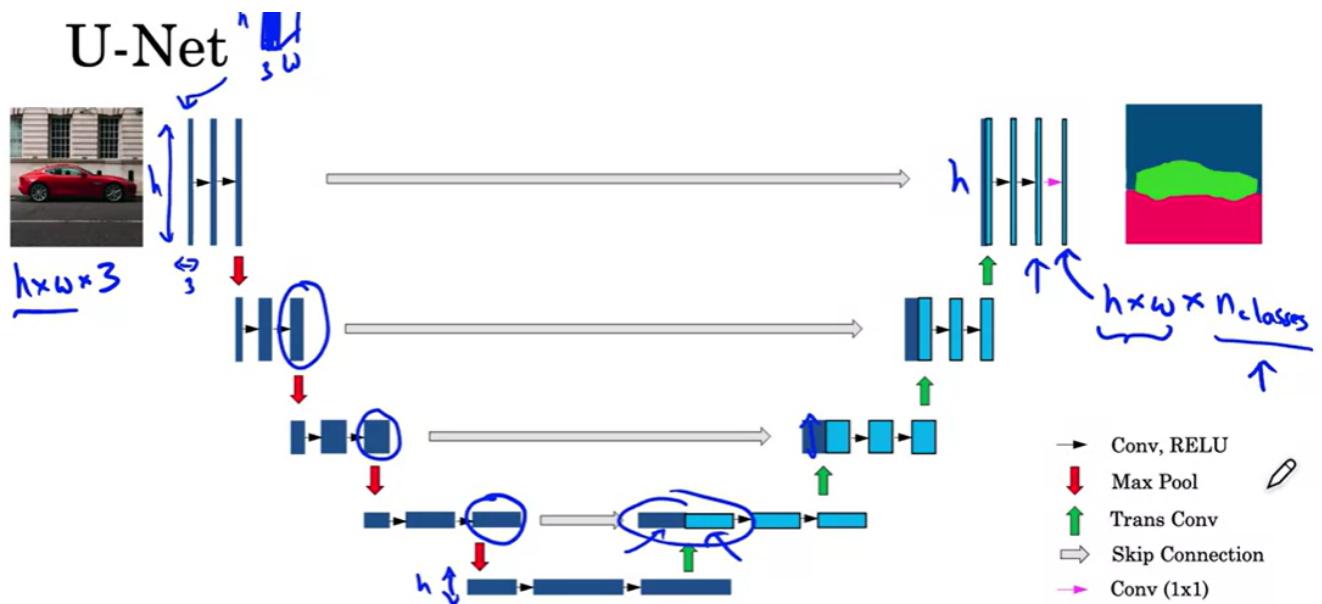
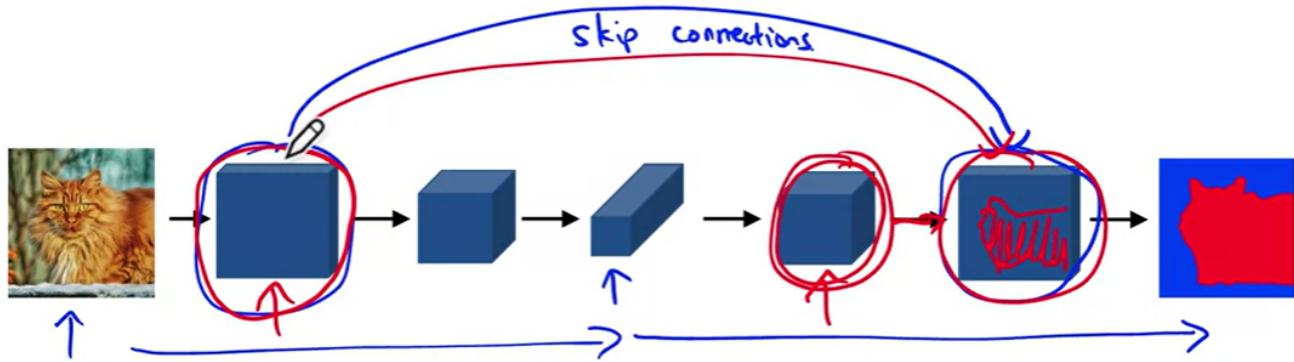
To perform the blowup in above image, we perform transpose convolution

## Transpose Convolution



filter  $f \times f = 3 \times 3$  padding  $p = 1$  stride  $s = 2$

## Deep Learning for Semantic Segmentation



# Face Recognition

## Face verification vs. face recognition

### → Verification

- Input image, name/ID
  - Output whether the input image is that of the claimed person
- $1:1$        $\frac{99\%}{99.9}$

### → Recognition

- Has a database of K persons
  - Get an input image
  - Output ID if the image is any of the K persons (or “not recognized”)
- $1:K$   
 $K=100 \leftarrow$

## One-Shot Learning

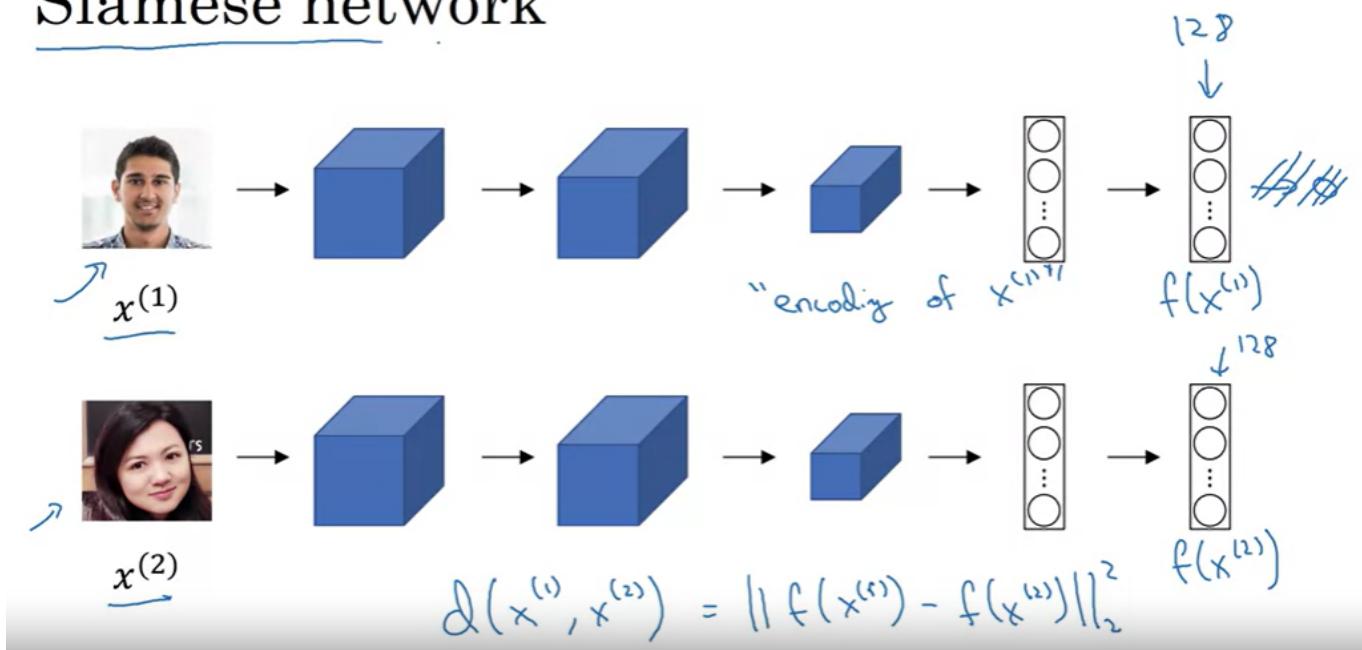
In face recognition we have a very less training set. thus we calculate degree of difference between feeded image and test image, if degree is greater than tau then we say its not the same person

# Learning a “similarity” function

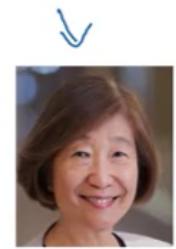
→  $d(\underline{\text{img1}}, \underline{\text{img2}}) = \text{degree of difference between images}$

If  $d(\text{img1}, \text{img2}) \leq \tau$       “same”  
 $> \tau$       “different”      } Verification.

## Siamese network



## Triplet Loss



Anchor  
A

Positive  
P

Anchor  
A

Negative  
N

$$\text{Want: } \underbrace{\|f(A) - f(P)\|^2}_{d(A,P)} + \lambda \leq \underbrace{\|f(A) - f(N)\|^2}_{d(A,N)}$$

## Loss function

Given 3 images  $\overset{\downarrow}{A, P, N}$ :

$$L(A, P, N) = \max \left( \underbrace{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \lambda}_{\geq 0} \right), 0)$$

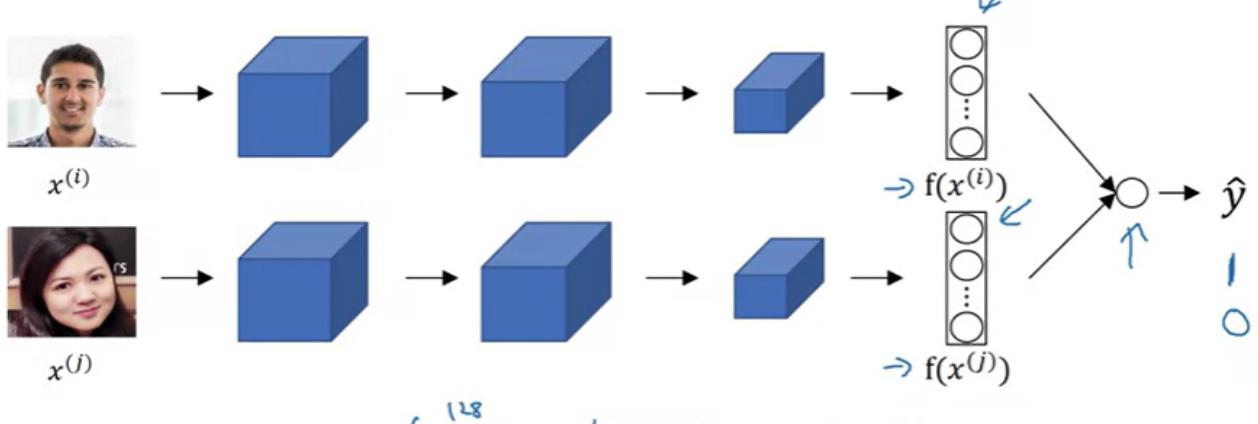
$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

$A, P$

Training set: 10k pictures of 1k persons

Choose hard A,P,N triplets such that it is hard to train on them , otherwise if those triplets are chosen randomly then the condition will be easily satisfied and gradient descent won't do much of a difference

# Learning the similarity function

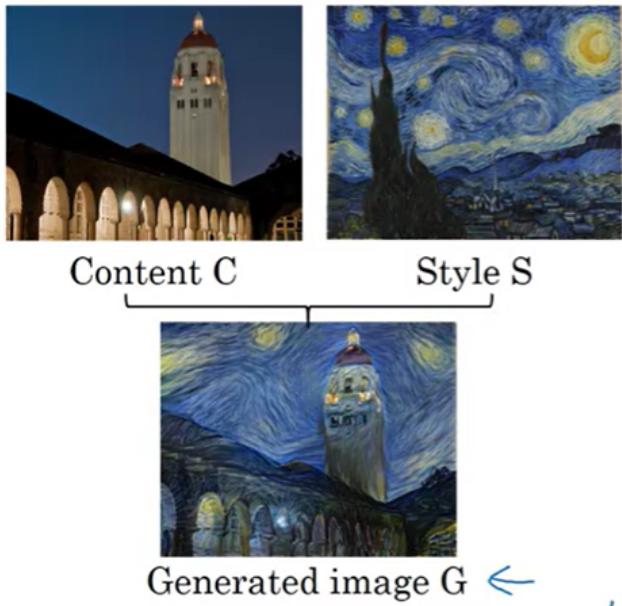


$$y_{\text{hat}} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

## Neural Style Transfer

### Cost Function

# Neural style transfer cost function



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

Find the generated image  $G$

1. Initiate  $G$  randomly

$$G: \underline{100} \times \underline{100} \times 3$$

$\uparrow$   
RGB

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

$$G := G - \frac{\partial}{\partial G} J(G)$$

## Content cost function

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

- Say you use hidden layer  $\underline{l}$  to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let  $\underline{a^{[l](C)}}$  and  $\underline{a^{[l](G)}}$  be the activation of layer  $\underline{l}$  on the images
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content

$$J_{content}(C, G) = \frac{1}{2} \| \underbrace{a^{[l](C)}}_{\uparrow} - \underbrace{a^{[l](G)}}_{\downarrow} \|_2^2$$

## Style matrix

Let  $a_{i,j,k}^{[l]}$  = activation at  $(i, j, k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$

H W C  
↓ ↓ ↓

$n_c$   
 $G_{kk'}^{[l]}$   
 $\sum_{k=1}^{n_c}$

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

$$\rightarrow G_{kk'}^{[l](G)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} a_{ijk}^{[l](G)} a_{ijk'}^{[l](G)}$$

"Gram matrix"

$$J_{style}^{[l]}(S, G) = \| G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)} \|_F^2$$

$$= \frac{1}{2n_H} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2$$

## Style cost function

$$\| G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)} \|_F^2$$

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

$$J_{style}(S, G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S, G)$$

Please note that in the next video at around 8:50 when Andrew wrote down the second formula of matrix, the second factor of the multiplication should be on index  $k'$  (k prime), not  $k$  (otherwise it does not calculate any kind of covariance). This is shown in black below.

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

So the formula should be:

$$G_{kk'}^{[l](G)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} a_{i,j,k}^{[l](G)} a_{i,j,k'}^{[l](G)}.$$

Also at 11:08 the style cost function formula should be the squared difference:

$$(G_S - G_G)^2$$

instead of just the difference:

$$(G_S - G_G).$$

The style cost function should be:

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