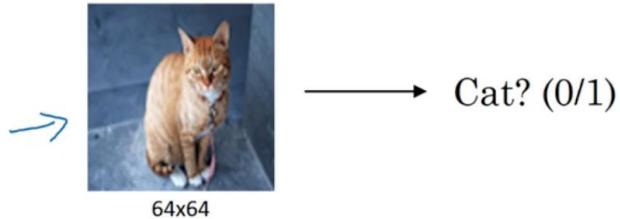


# Computer Vision Problems

Image Classification



Neural Style Transfer

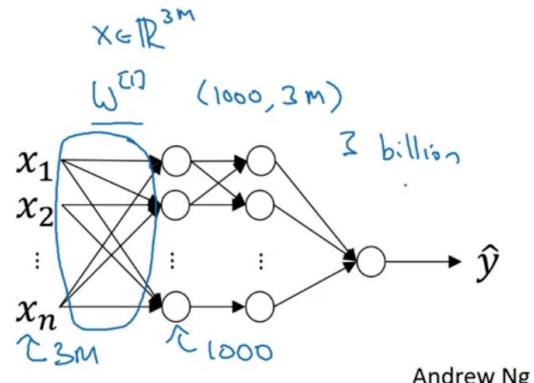
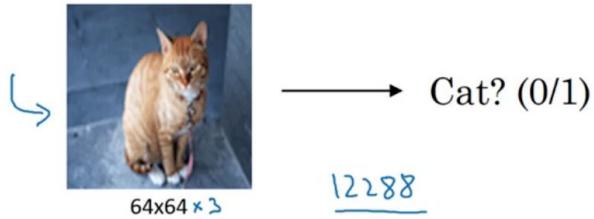


Object detection



Andrew Ng

## Deep Learning on large images



Andrew Ng

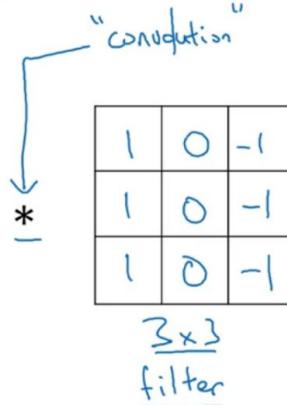
Small images are ok - but what about larger images? Need a way to manage 3b parameters!

## Vertical edge detection

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 3 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$$

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

6x6



=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

4x4

python: conv-forward  
tensorflow: tf.nn.conv2d  
keras: Conv2D

Andrew Ng

Element wise multiply all parts of the matrix with numbers in the image to get the convolution outputs.

BRILLIANT! Why does kernel actually work? Convolution actually masks on various shades of color and then tries to find the region of space where edge might occur - here in the middle!

## Vertical edge detection

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array} \quad * \quad
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} \quad = \quad
 \begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}
 \\
 \text{6x6} \qquad \qquad \qquad \text{3x3} \qquad \qquad \qquad \text{4x4}
 \end{array}$$

↓      ↓      ↓

↑      ↑      ↑

Andrew Ng

## Vertical edge detection examples

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array} \quad * \quad
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} \quad = \quad
 \begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}
 \\
 \text{6x6} \qquad \qquad \qquad \text{3x3} \qquad \qquad \qquad \text{4x4}
 \end{array}$$

→      ↓      ↓      ↓

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 0 & 0 & 0 & 10 & 10 & 10 \\ \hline
 \end{array} \quad * \quad
 \begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array} \quad = \quad
 \begin{array}{|c|c|c|c|} \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 0 & -30 & -30 & 0 \\ \hline
 \end{array}
 \\
 \text{6x6} \qquad \qquad \qquad \text{3x3} \qquad \qquad \qquad \text{4x4}
 \end{array}$$

→      ↓      ↓      ↓

Andrew Ng

# Vertical and Horizontal Edge Detection

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

Vertical

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

Horizontal

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

$6 \times 6$

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

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

Andrew Ng

If you go to larger say  $1000 \times 1000$  images - you won't see the transitions actually  $(10, -10)$

## Learning to detect edges

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

↑

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

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

Sobel filter

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

Scharr filter

convolution

$$\begin{array}{|c|c|c|} \hline w_1 & w_2 & w_3 \\ \hline w_4 & w_5 & w_6 \\ \hline w_7 & w_8 & w_9 \\ \hline \end{array}$$

$3 \times 3$

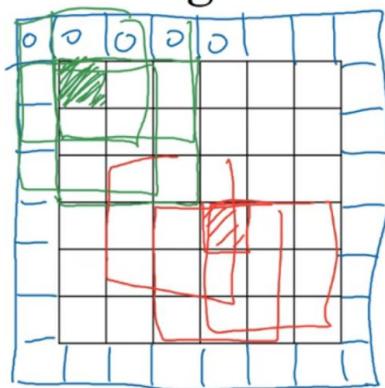
$$\begin{array}{c} 45^\circ \\ 70^\circ \\ 73^\circ \end{array}$$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

Andrew Ng

Use deep learning to learn the weights instead of hand coding them! You can even detect random edges like at  $73$  degrees etc. Whatever the model chooses to learn using backprop. Convolution is still the key though.

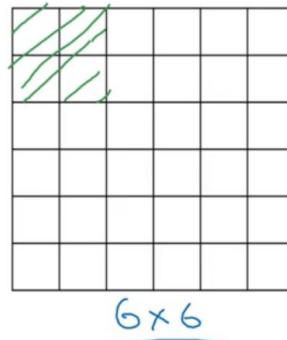
# Padding



- shrinks output  
- throws away info from edge

$$\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array}$$

$3 \times 3$   
 $f \times f$



$$\frac{6 \times 6}{n \times n} \rightarrow 8 \times 8$$

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

$$6-3+1=4$$

$$P = \text{padding} = 1$$

$$\xrightarrow{\quad \quad \quad} \underline{\underline{4 \times 4}}$$

$$\frac{n+2p-f+1}{6+2-3+1} \times \frac{n+2p-f-1}{n+2-3+1} = 6 \times 6$$

Andrew Ng

Padding ---> keeps same image dimension which is important for deep networks, otherwise you will lose info! Also preserves info about corners and keeps them.

## Valid and Same convolutions

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

$$\text{"Valid": } \begin{array}{c} n \times n \\ 6 \times 6 \end{array} * \begin{array}{c} f \times f \\ 3 \times 3 \end{array} \rightarrow \frac{n-f+1}{4 \times 4} \times \frac{n-f+1}{4 \times 4}$$

**"Same":** Pad so that output size is the same as the input size.

$$\begin{array}{l} n+2p-f+1 \times n+2p-f+1 \\ n+2p-f+1 = n \Rightarrow p = \frac{f-1}{2} \\ 3 \times 3 \quad p = \frac{3-1}{2} = 1 \end{array}$$

$f$  is usually odd  
 $1 \times 1$   
 $3 \times 3$   
 $5 \times 5$   
 $7 \times 7$

$$p=2$$

Andrew Ng

Odd number  $f$ : symmetric padding instead of asymmetric and also has central pixel

## Strided convolution

Diagram illustrating strided convolution:

- Input Matrix:**  $7 \times 7$
- Filter:**  $3 \times 3$
- Stride:**  $s=2$
- Output Matrix:**  $3 \times 3$

The output values are labeled as  $\lfloor z \rfloor = \text{floor}(z)$ .

$$\begin{matrix} n \times n \\ \text{padding } p \end{matrix} * \begin{matrix} f \times f \\ \text{stride } s \\ s=2 \end{matrix}$$

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

$$\frac{7+0-3}{2} + 1 = \frac{4}{2} + 1 = 3$$

Andrew Ng

## Summary of convolutions

$$n \times n \text{ image} \quad f \times f \text{ filter}$$

$$\text{padding } p \quad \text{stride } s$$

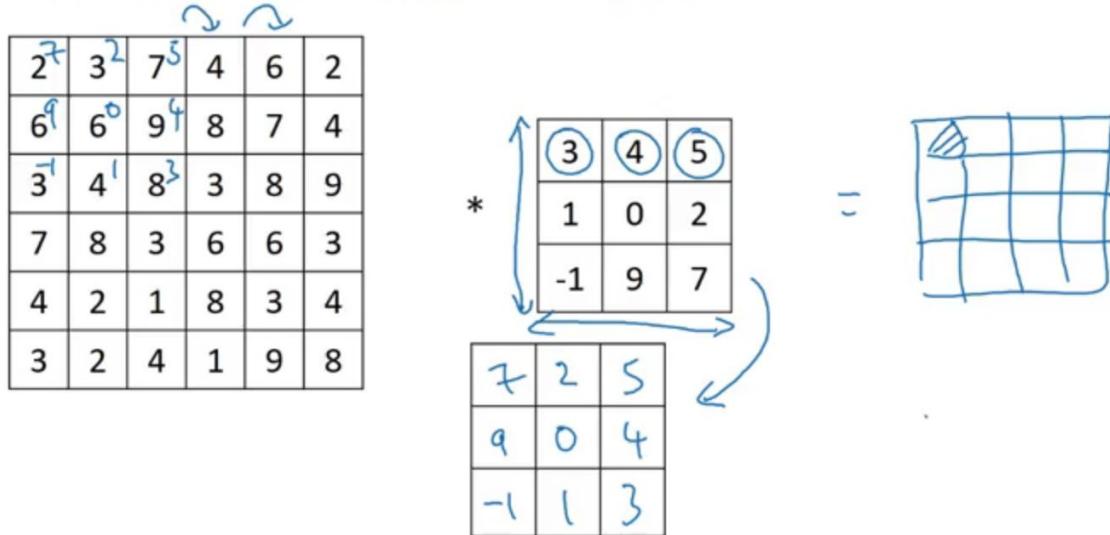
$\text{Output}_0$

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

OUTPUT SIZE with floor if not an integer

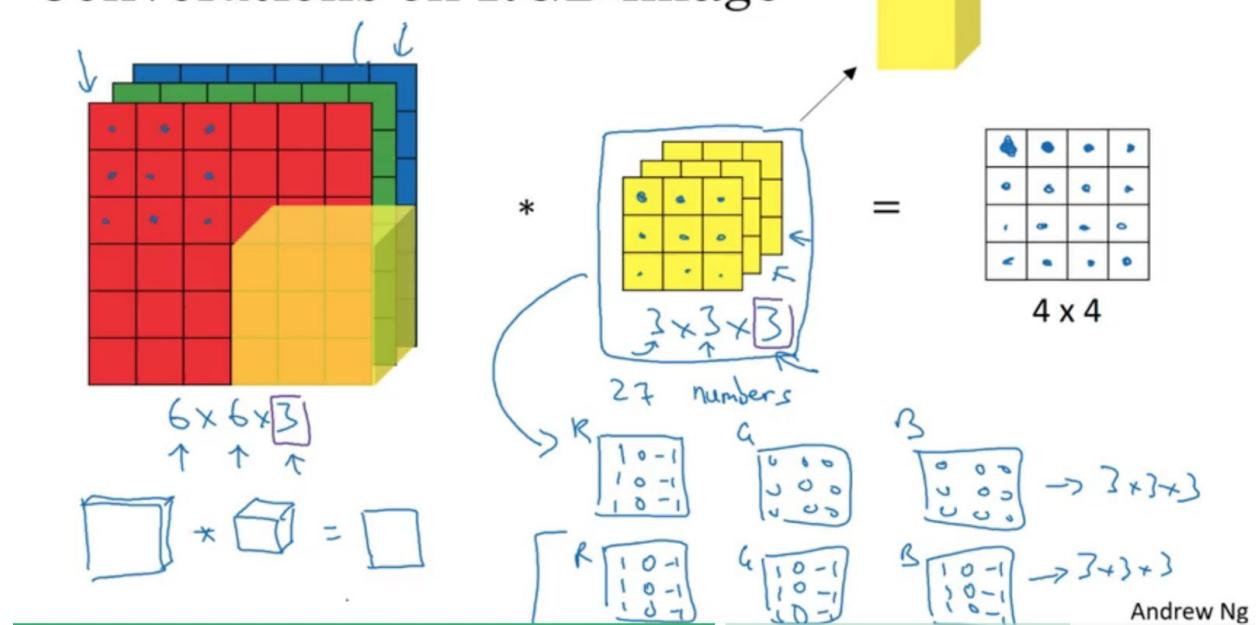
# Technical note on cross-correlation vs. convolution

Convolution in math textbook:



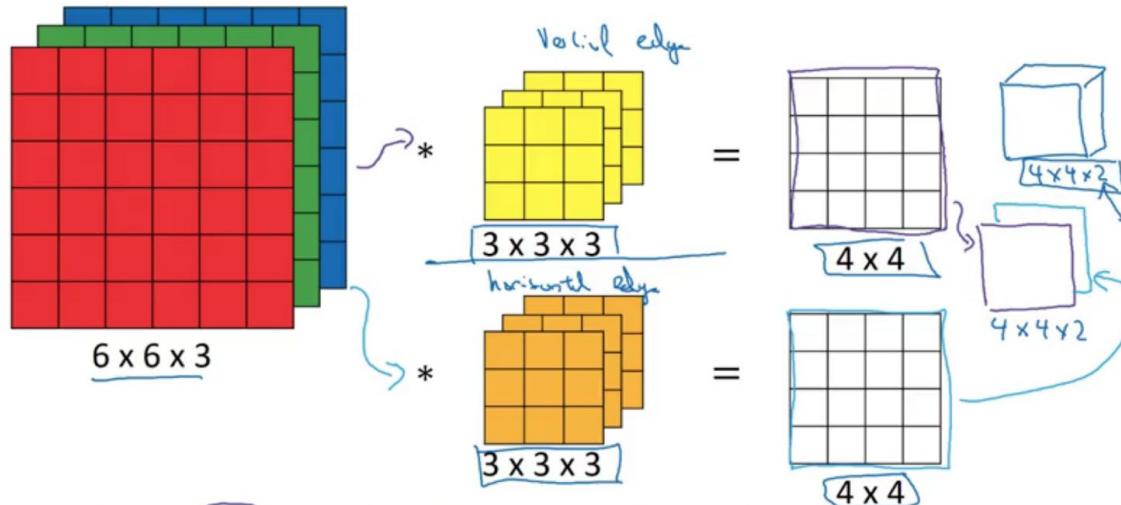
## CONVOLUTIONS OVER A VOLUME

Convolutions on RGB image



Channels in image and channels in filter is exactly the same! To detect only red edges you can make 1,1,1,0,0,0,-1,-1,-1 and rest all 0s and if you dont care which color edge it is then all 3 filter dimensions will have this same value.

## Multiple filters

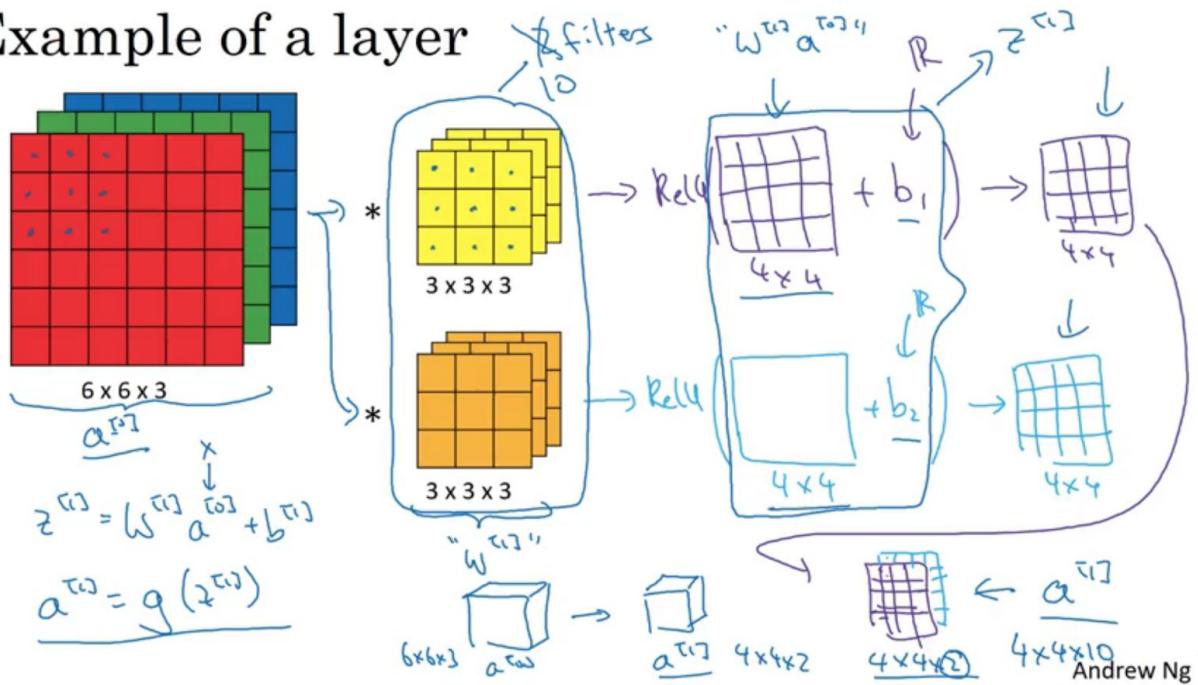


Summary:  $n \times n \times n_c$   $\times f \times f \times n_c$   $\rightarrow \frac{n-f+1}{4} \times \frac{n-f+1}{4} \times n'_c$  #filters

$6 \times 6 \times 3$   $\times 3 \times 3 \times 3$   $\rightarrow 4 \times 4 \times 2$  #filters

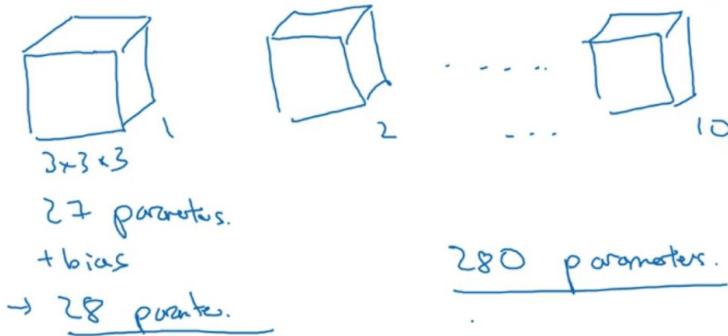
Andrew Ng

## Example of a layer



# Number of parameters in one layer

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



Andrew Ng

Even if image is  $10000 \times 10000$  parameters are still only 280! NO OVERFITTING DUE TO THIS IN CNNs!

## Summary of notation

If layer  $\underline{l}$  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  $l$ .

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

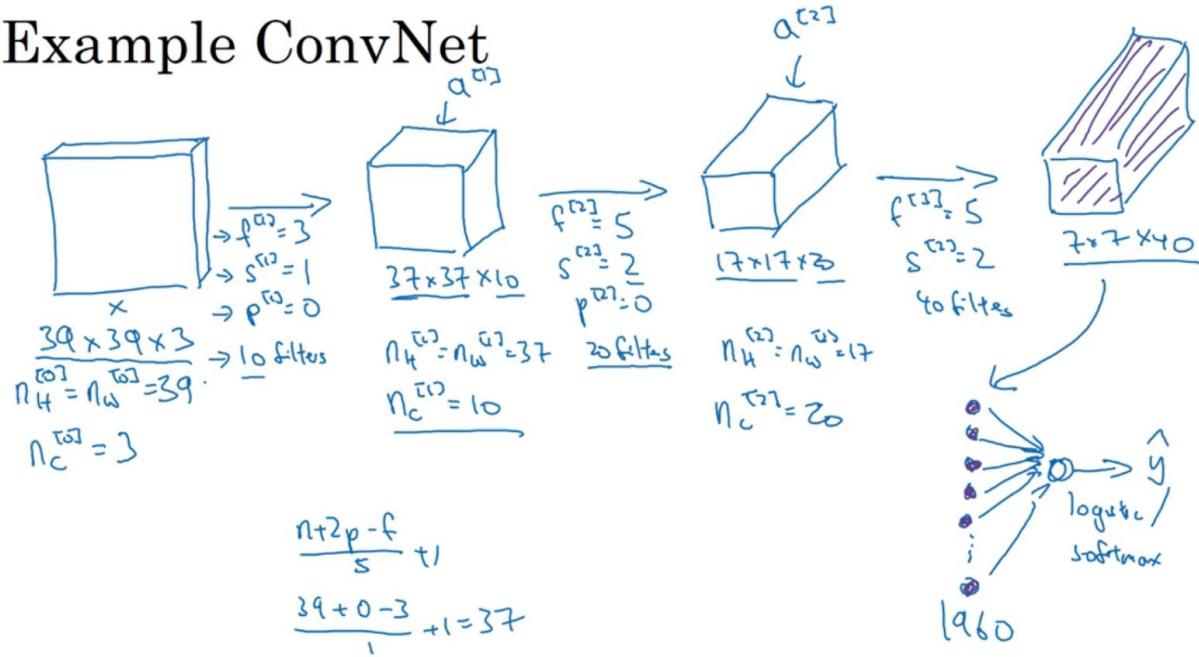
Output:  $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$  ←  $n_c^{[l]}$

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

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

Andrew Ng

## Example ConvNet

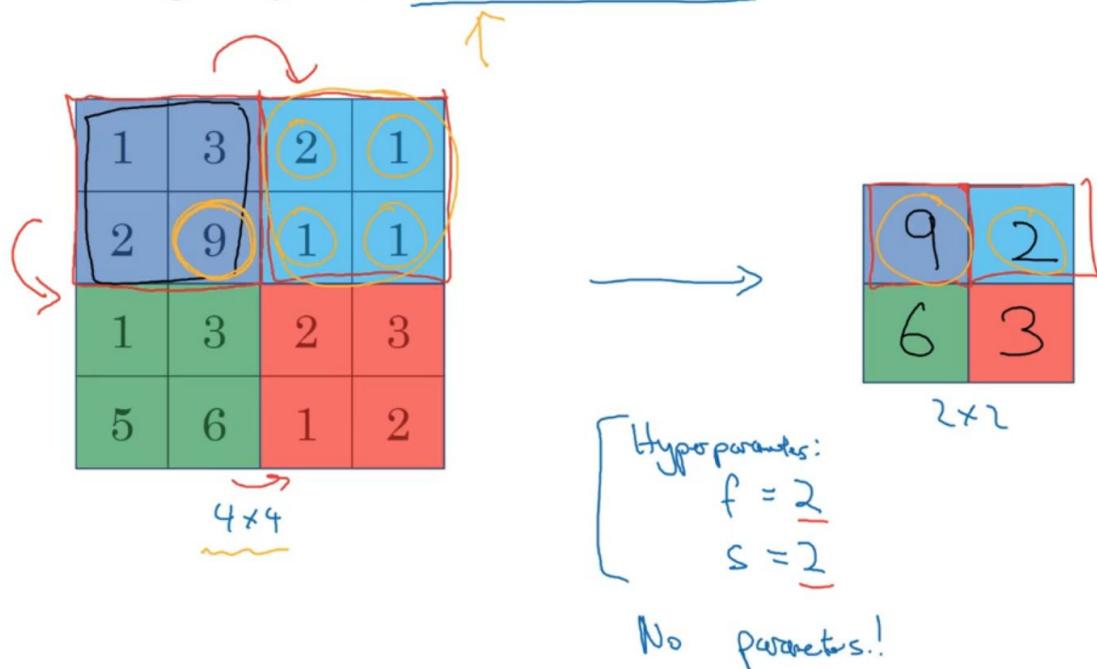


Andrew Ng

Types of layer in a convolutional network:

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

## Pooling layer: Max pooling



## Summary of pooling

Hyperparameters:

$f$  : filter size

$f = 2, s = 2$

$s$  : stride

$f = 3, s = 2$

Max or average pooling

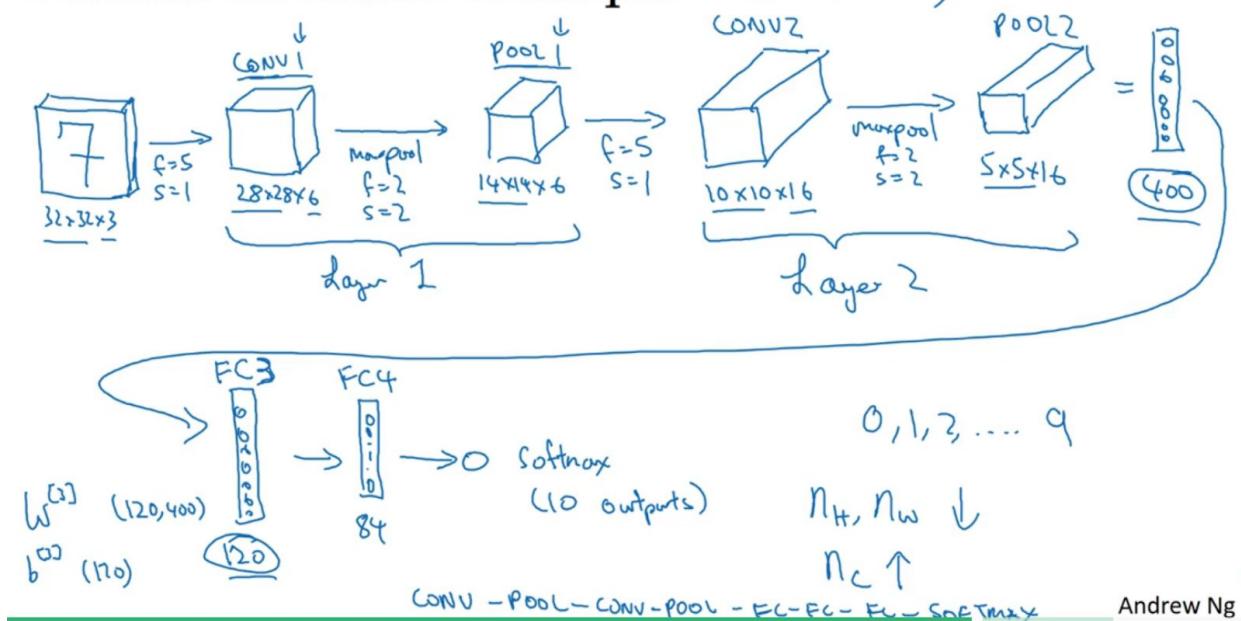
$\Rightarrow p$ : padding

No parameters to learn!

$$n_H \times n_w \times n_c$$

$$\downarrow \\ \left\lfloor \frac{n_H - f + 1}{s} \right\rfloor \times \left\lfloor \frac{n_w - f}{s} + 1 \right\rfloor \\ \times n_c$$

# Neural network example



CONV → POOL → CONV → POOL → FC → FC → FC → SOFTMAX

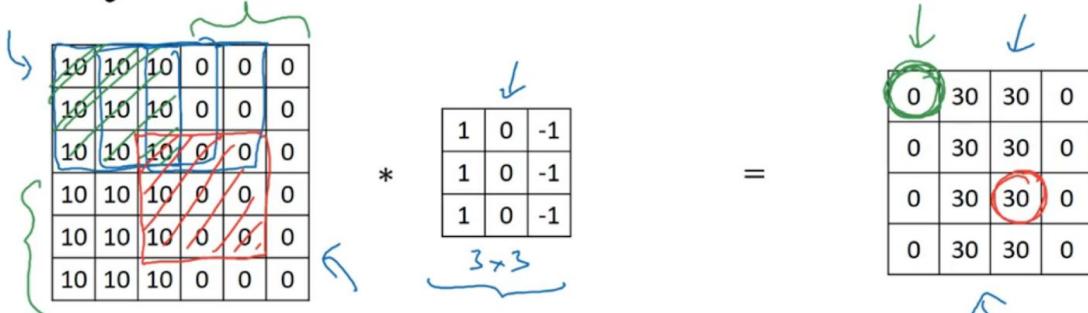
# Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	$\sim 3,072$ $a^{(0)}$	0
CONV1 (f=5, s=1)	(28,28,8)	6,272	208 ↙
POOL1	(14,14,8)	1,568	0 ↙
CONV2 (f=5, s=1)	(10,10,16)	1,600	416 ↙
POOL2	(5,5,16)	400	0 ↙
FC3	(120,1)	120	48,001 ↘
FC4	(84,1)	84	10,081 ↘
Softmax	(10,1)	10	841

Activation size should gradually decrease -- not too sharply! Slowly reduce it. Most parameters are in FC!

WHY CONVOLUTIONS? Sparsity, parameter sharing and translation invariance. Cat is a cat even if shifted by a few pixels!

## Why convolutions



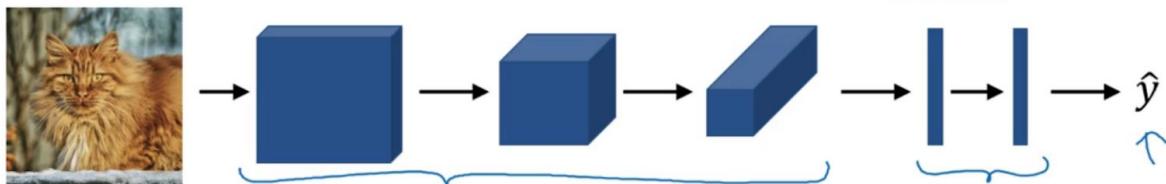
**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

→ **Sparsity of connections:** In each layer, each output value depends only on a small number of inputs.

## Putting it together

Training set  $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ .

$w, b$



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

Use gradient descent to optimize parameters to reduce  $J$

Reduce cost function with the weights and biases and build a CNN!

## WEEK 2 (CASE STUDY OF NETWORKS)

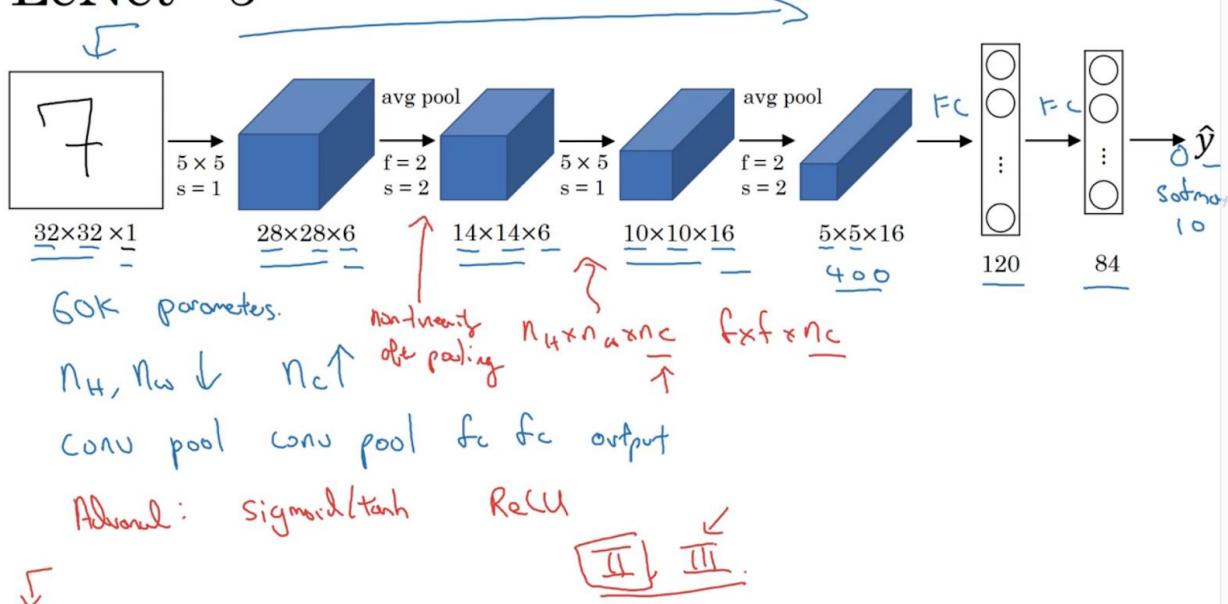
## Classic networks:

- LeNet-5 ←
- AlexNet ←
- VGG ←

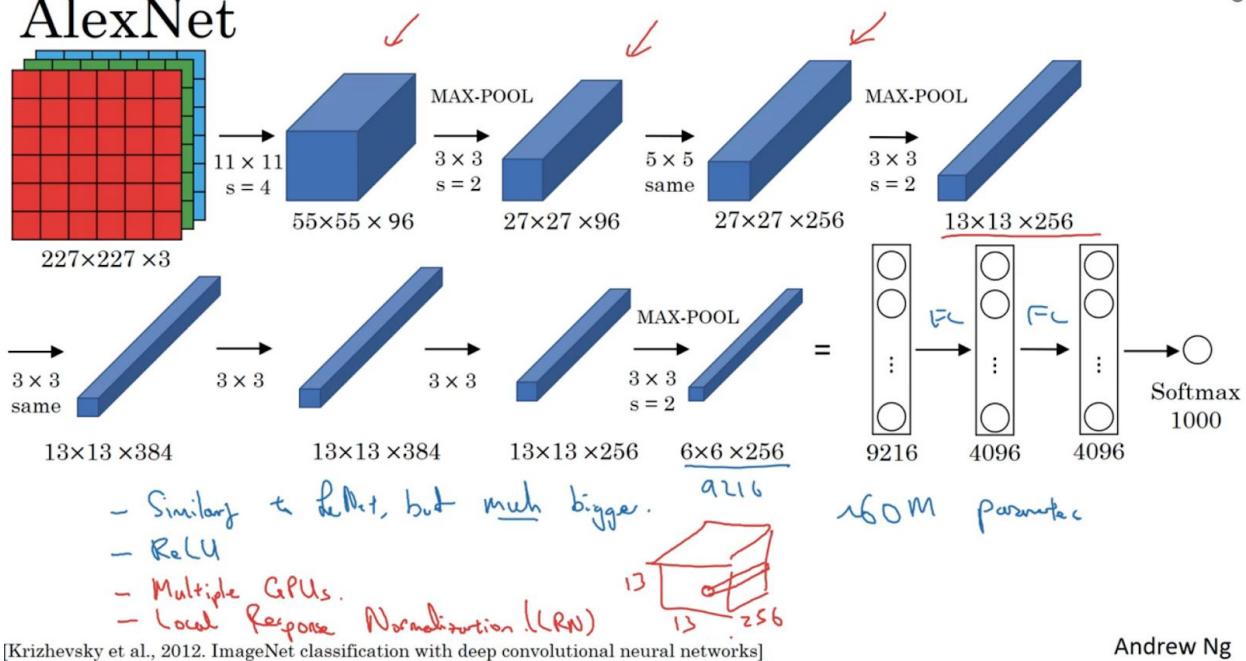
ResNet (152)

Inception

## LeNet - 5



# AlexNet

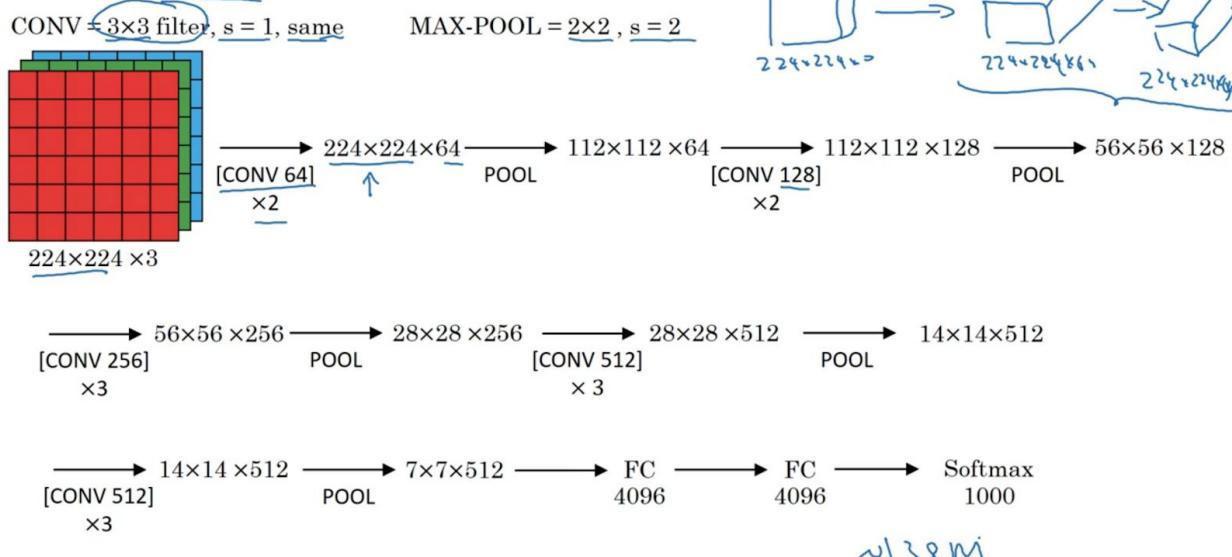


[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

Andrew Ng

LRN is not used today as much. Very influential paper in general though!  
 EASIER PAPER TO READ!

## VGG - 16



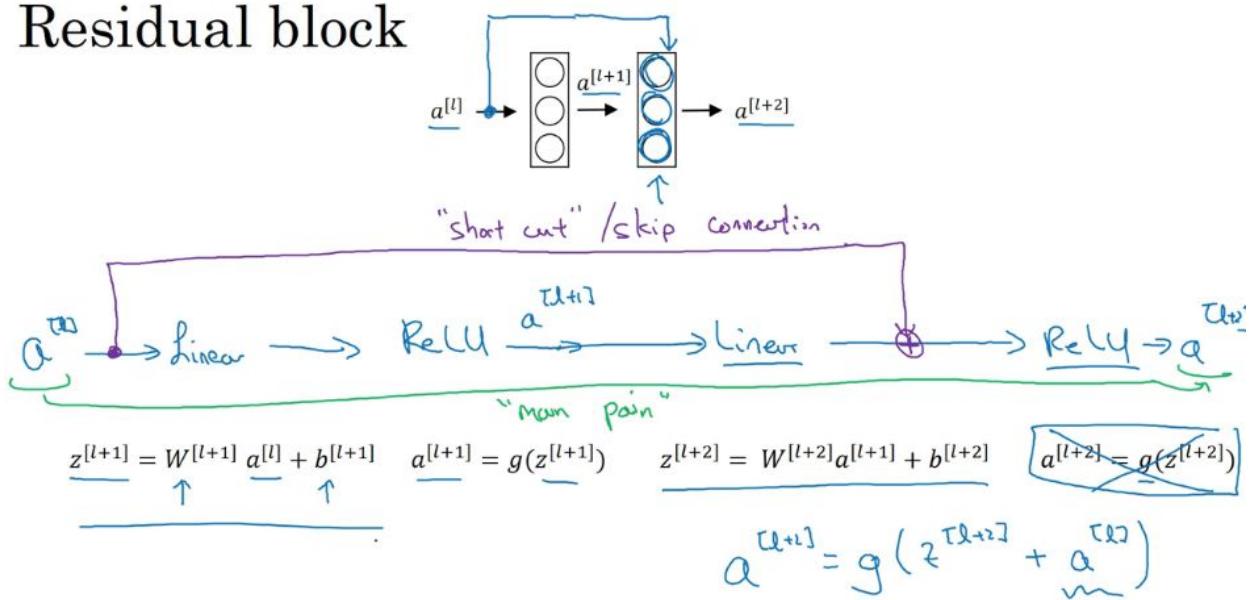
[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]

Andrew Ng

Pattern of height width goes down but channels go up → systematic and attractive!  
 Start from reading Alexnet paper to VGG net paper

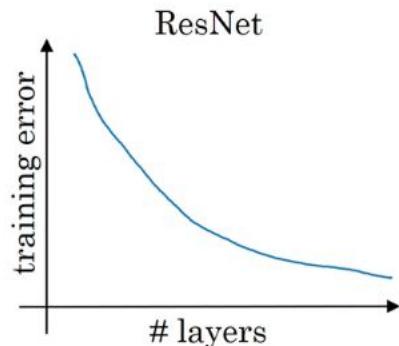
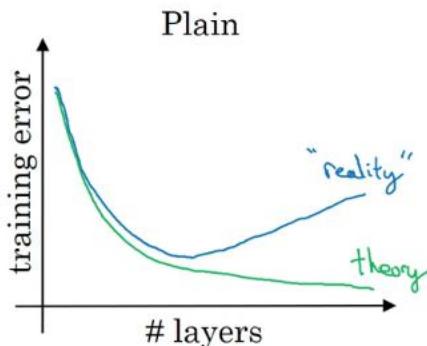
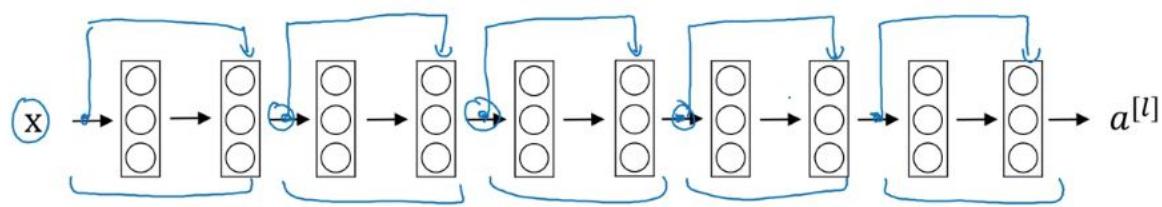
## RESNETS

### Residual block



Skipping over and feed into much deeper layers - can train even deeper ResNet!

### Residual Network

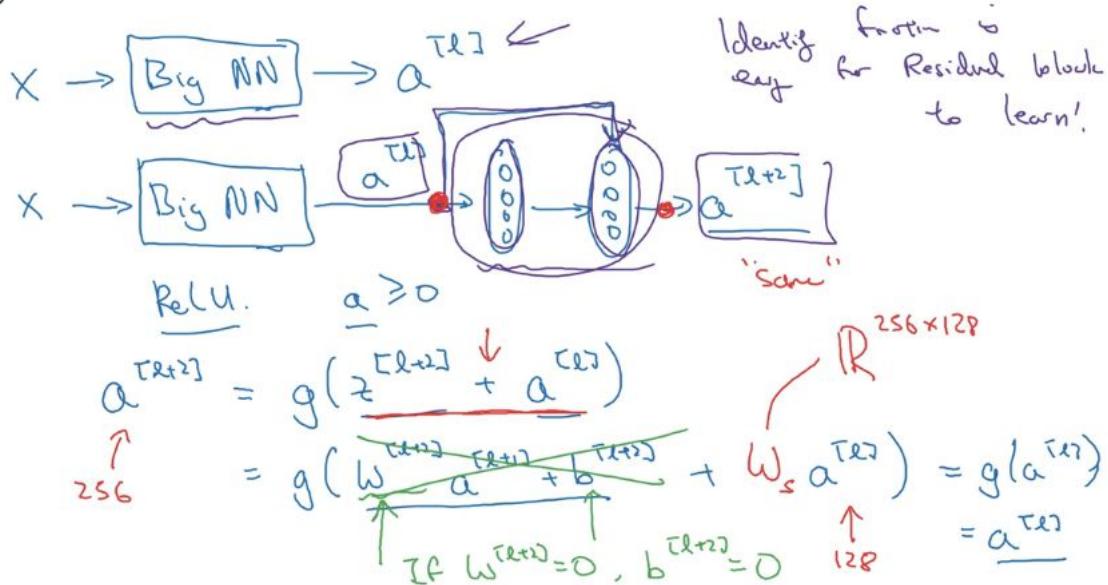


Residual networks work because -- its easy to learn the identity function. Basically if there is no additional complexity to be learnt, it can just feed  $a[l]$  to  $a[l+2]$  but if there is an additional complexity - BONUS!

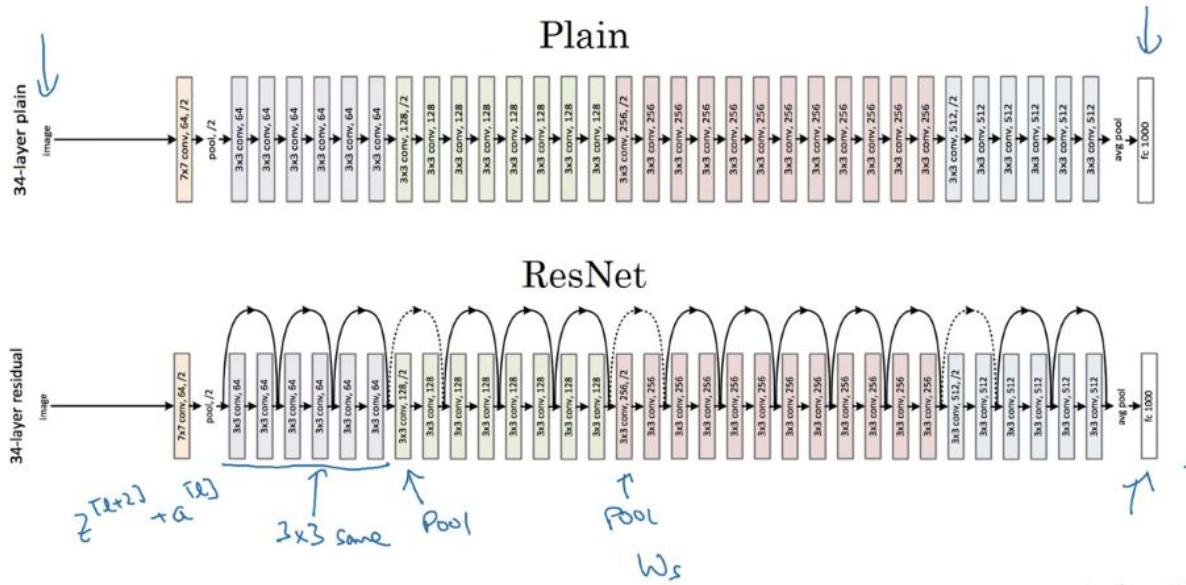
Thats why generally deeper and deeper layers are only good to certain extent, after which your error increases. But with resnet - you have a backup option of identity being moved forward so you can choose more complexity or choose no complexity.

Assume same dimensions of Z and A -- therefore SAME convolutions, but if not you can just do padding using Ws.

## Why do residual networks work?



# ResNet



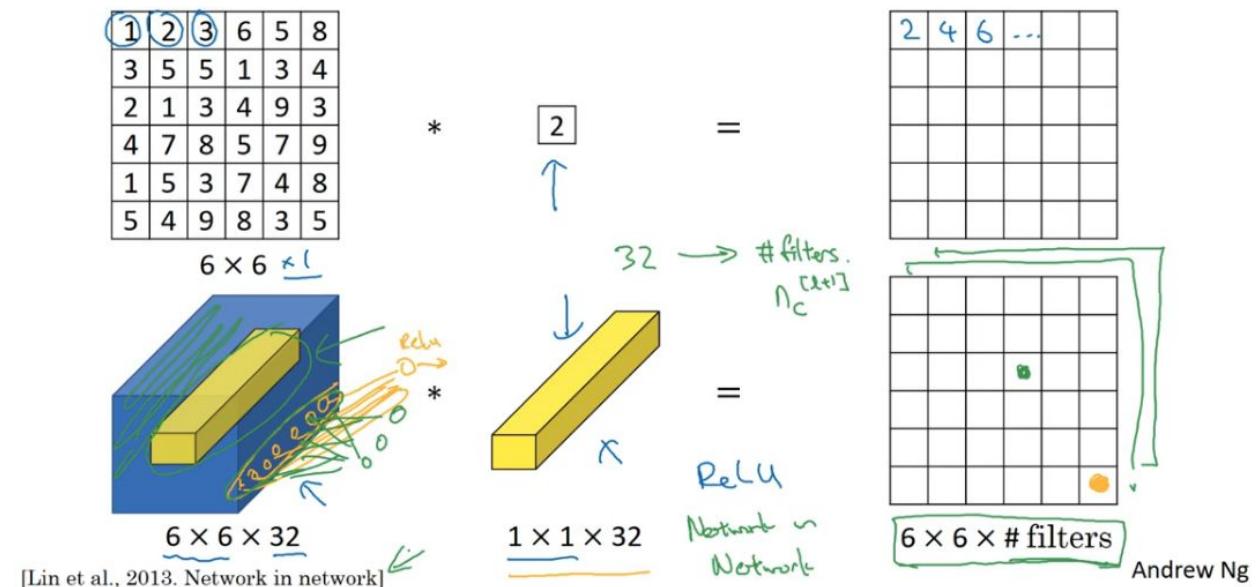
He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

## NETWORK IN NETWORK

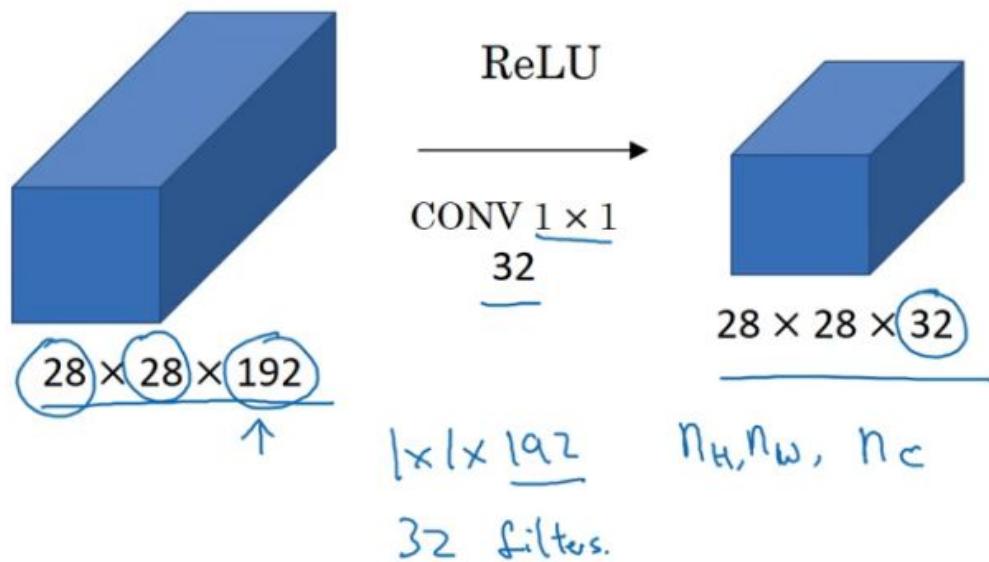
1x1 convolution : If number of channels is too much , you can use 1x1 conv. For shrinking height and width - use max pooling!

## Why does a $1 \times 1$ convolution do?



[Lin et al., 2013. Network in network]

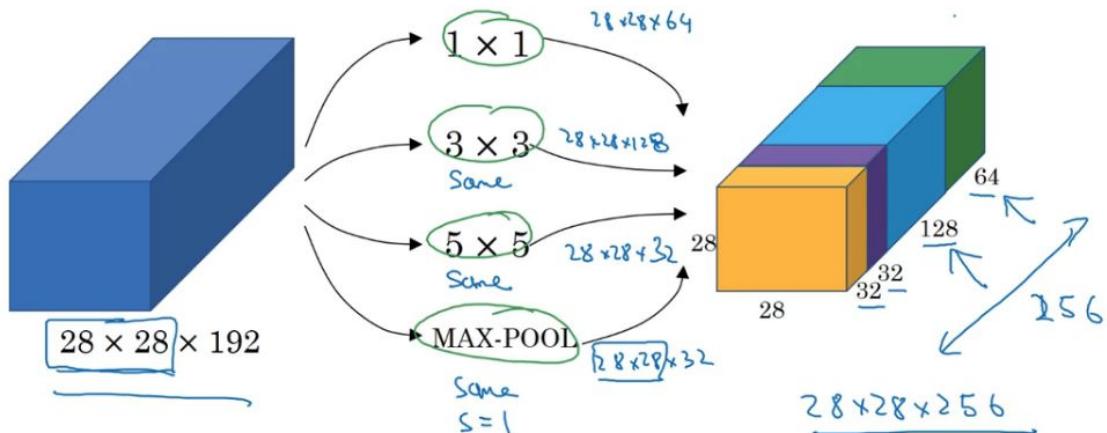
# Using $1 \times 1$ convolutions



## INCEPTION NETWORK

DO ALL YOU MIGHT WANT! If you are not sure - just stack them all!  
For last max pool ,there is  $1 \times 1$  conv at the end to reduce the channels.

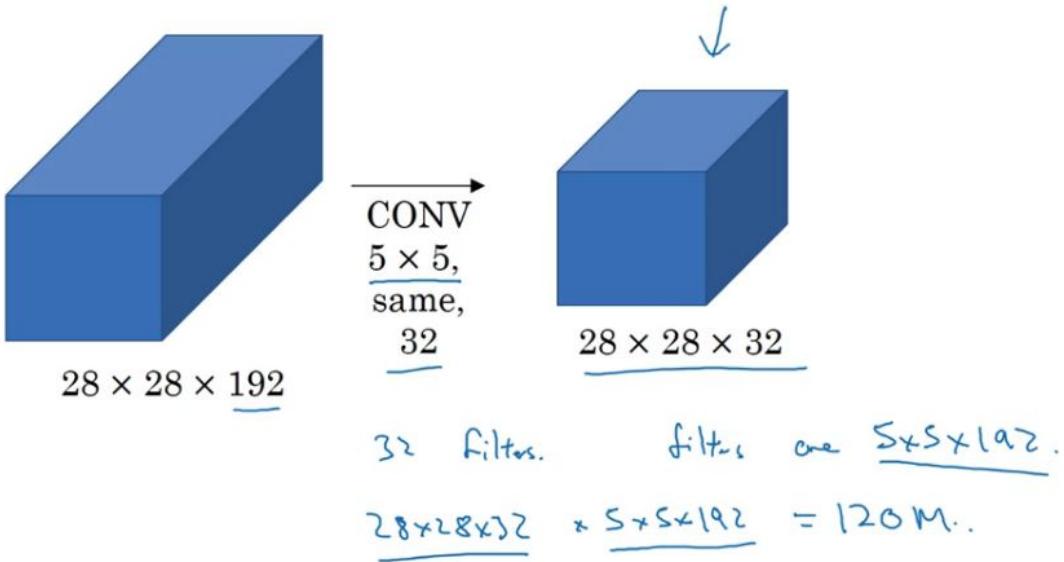
## Motivation for inception network



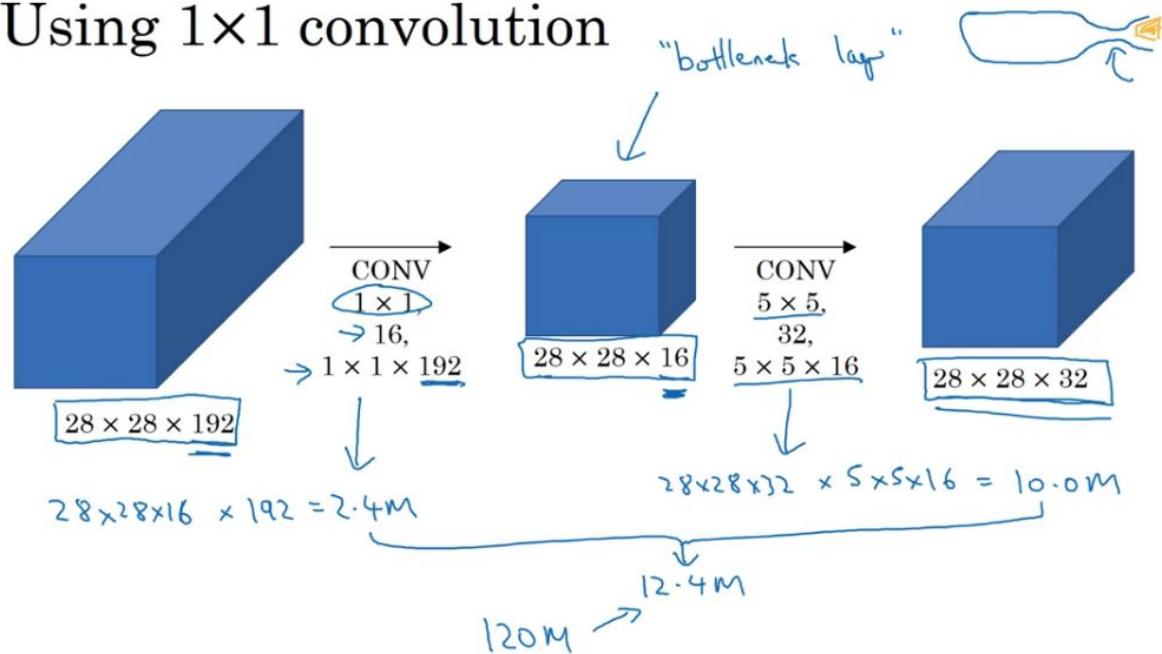
[Szegedy et al. 2014. Going deeper with convolutions]

Andrew Ng

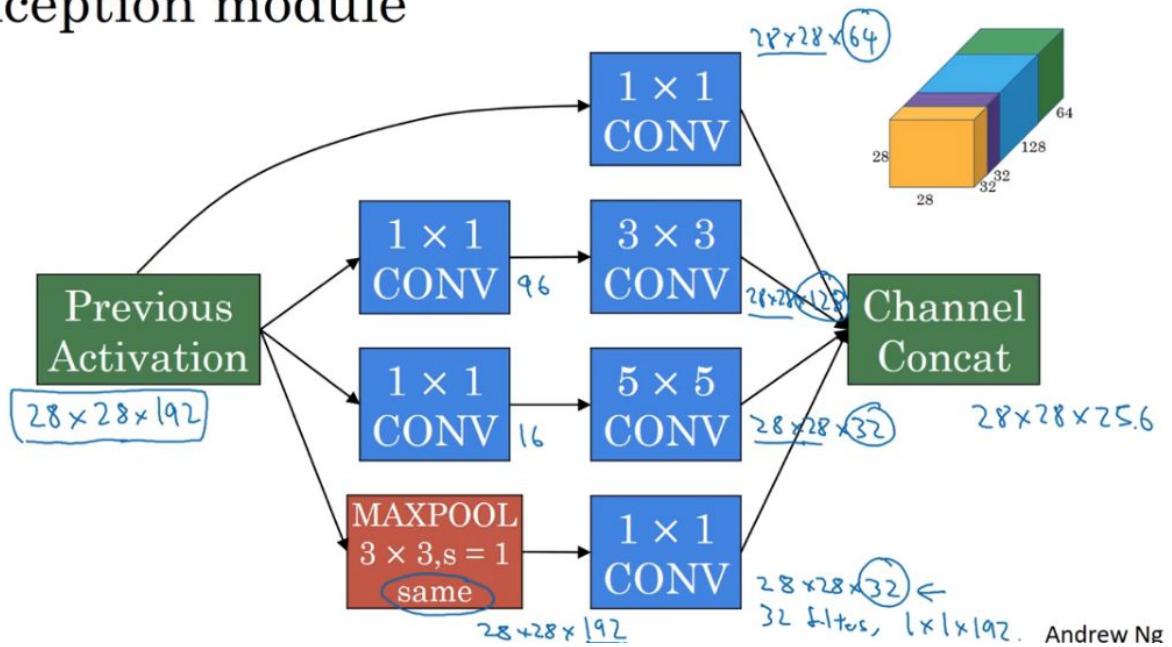
## The problem of computational cost



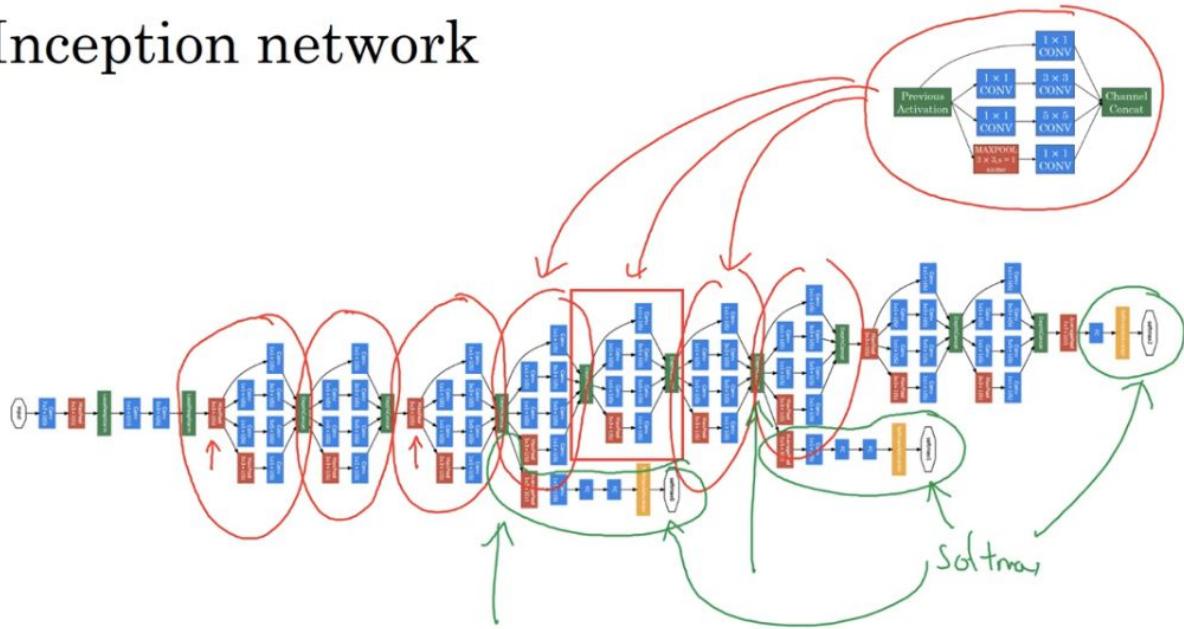
# Using $1 \times 1$ convolution



## Inception module



# Inception network



there are a very few softmax layers even in between - to have regularization. To make sure even the hidden data is not bad for predicting the output.

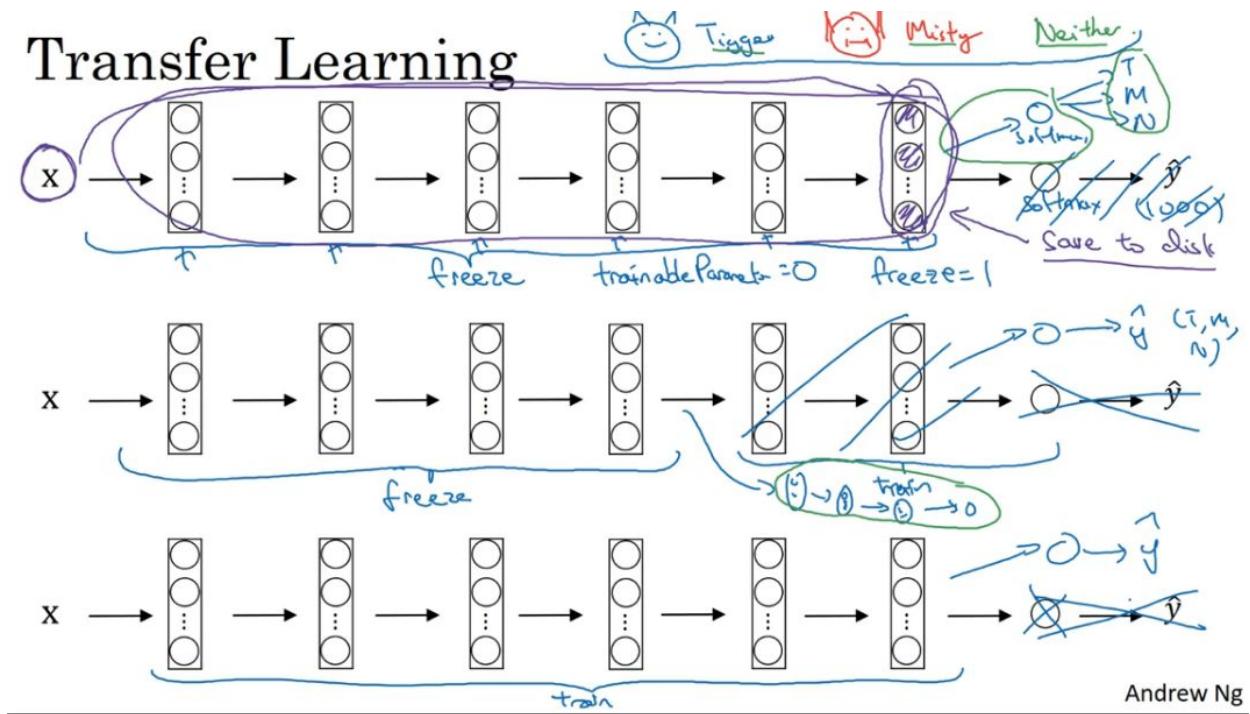
There are ideas of inception network combined with Res nets!

## PRACTICAL ADVICE

### TRANSFER LEARNING:

Trainable parameter = 0 -- tensorflow and keras might have it to freeze earlier layers and only add your own softmax function in the end.

# Transfer Learning



More data you have → more layers you can train! If you have insanely high data, you can use weights and architecture as initialization to your problem and train the whole network!  
 People have spent weeks and months training some models, good to re-use them!

## DATA AUGMENTATION:

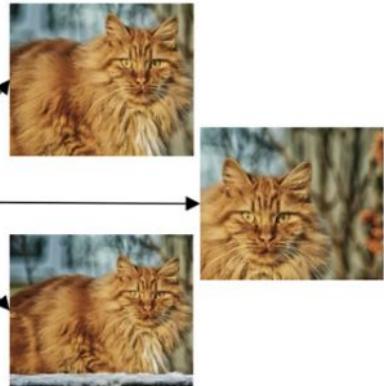
you can train! If you have insanely high data, you can use weights and architecture

# Common augmentation method

Mirroring



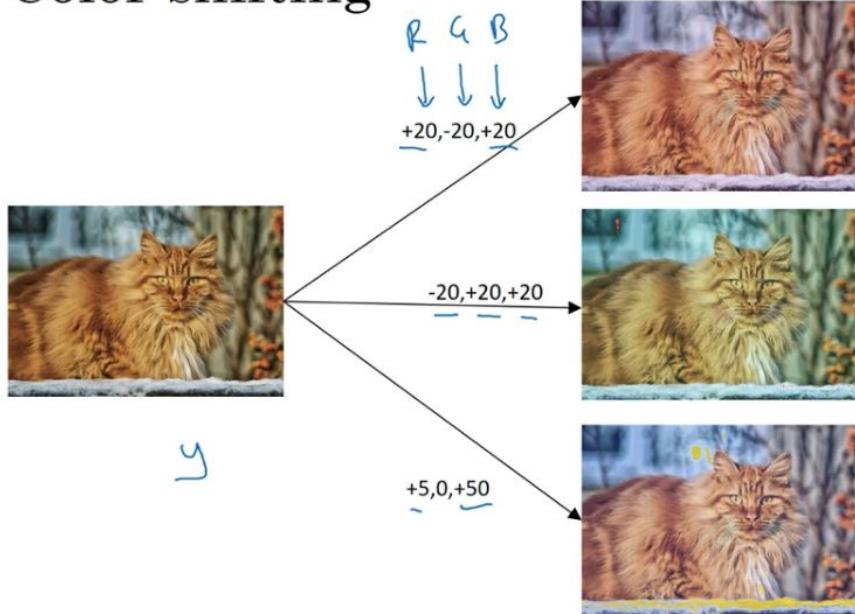
Random Cropping



Δ

Rotation, Shearing, Local Warping etc!

## Color shifting



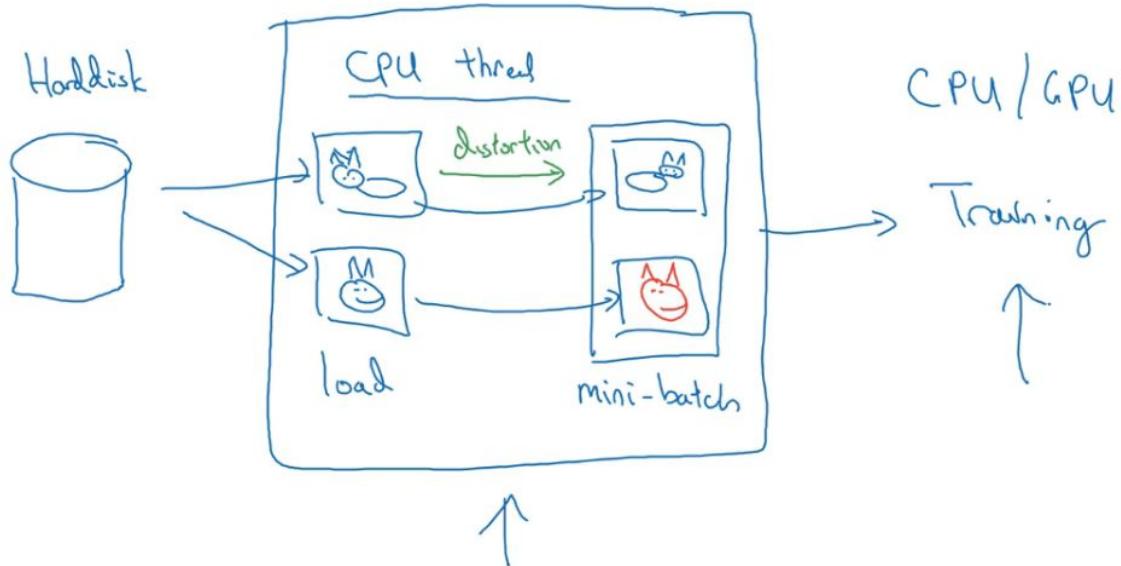
Advanced:  
PCA  
[ml-class.org](http://ml-class.org)

[ AlexNet paper  
[ "PCA color  
augmentation."  
R B G

Andrew Ng

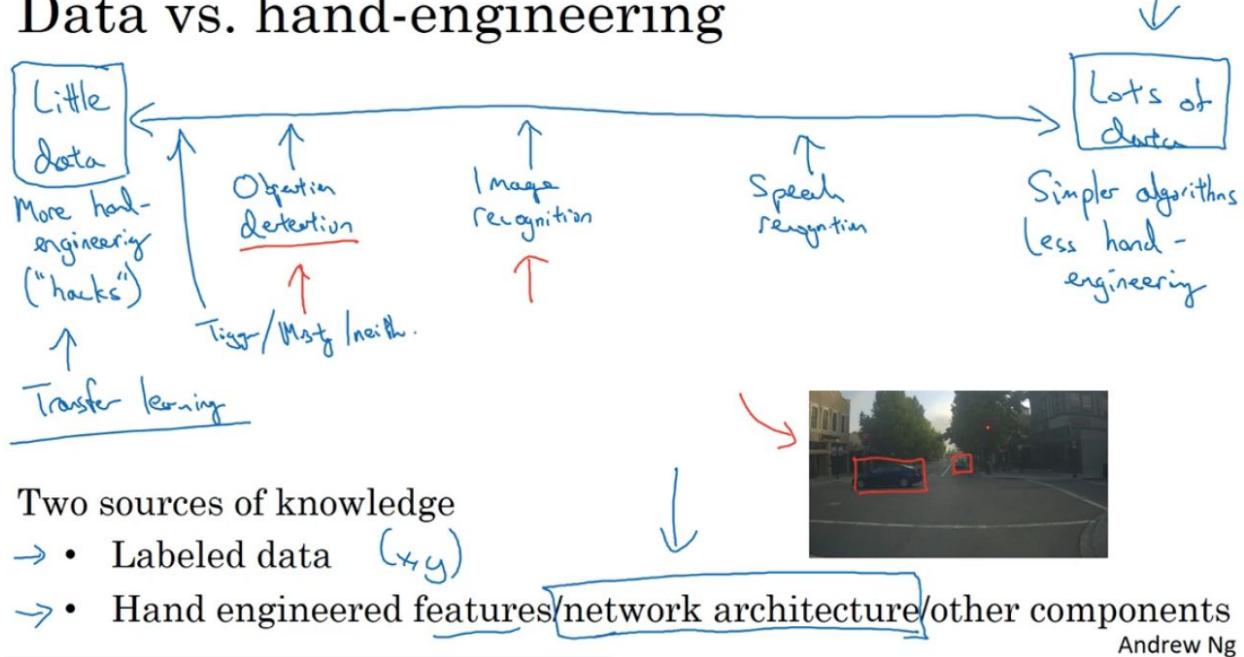
You can sample the increase or decrease of R G B pixels based on PCA \_ so you if you have lower G tint, your new images will continue to have lower G tint. You will not have very very high G or something.

# Implementing distortions during training



STATE OF COMPUTER VISION!

## Data vs. hand-engineering



# Tips for doing well on benchmarks/winning competitions

Ensembling

3-15 networks

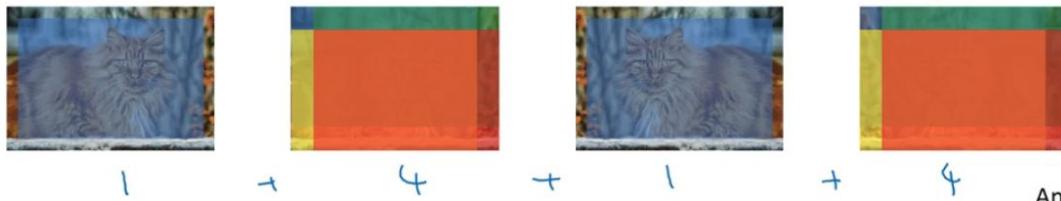


- Train several networks independently and average their outputs

Multi-crop at test time

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

10-crop



Can't use them at production for customers as it's more expensive! But great for winning competitions or benchmarks!

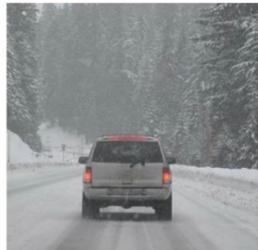
## Use open source code

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset

## OBJECT LOCALIZATION

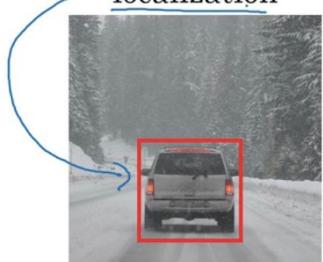
# What are localization and detection?

Image classification



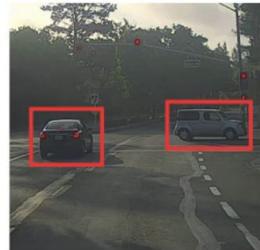
"Car"  
1 object

Classification with localization



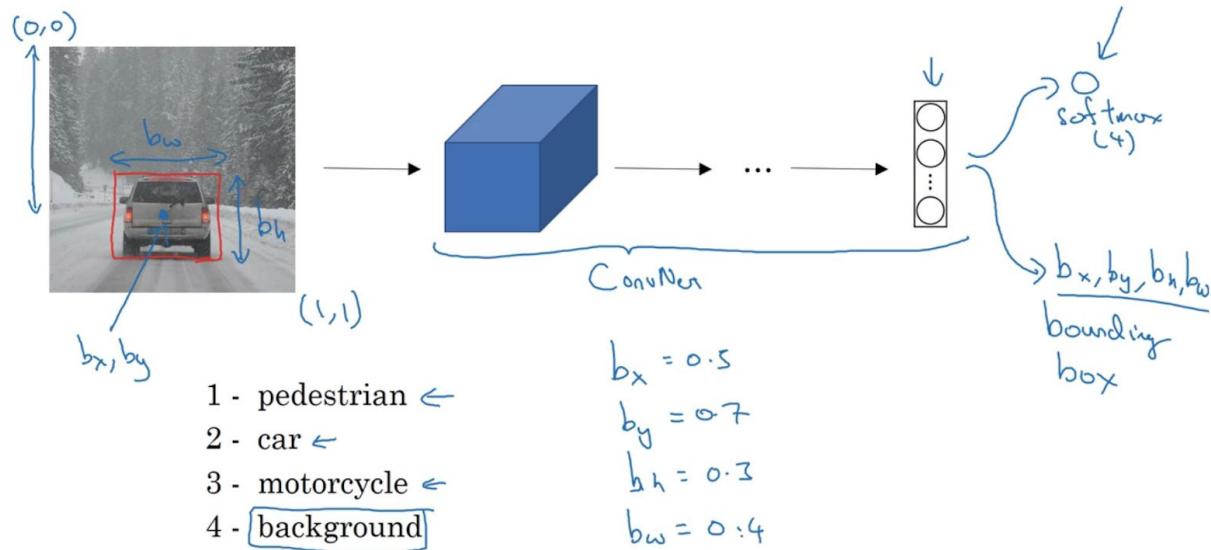
"Car"

Detection



multiple  
objects

## Classification with localization

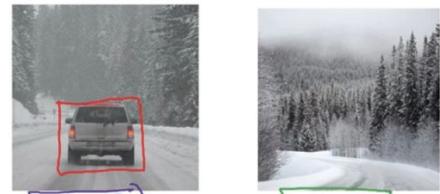


# Defining the target label $y$

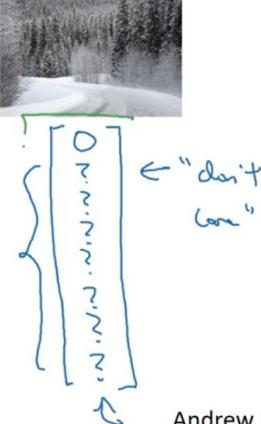
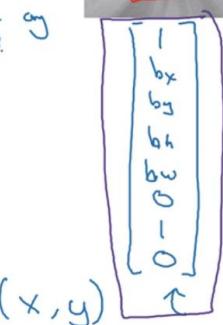
- 1 - pedestrian
- 2 - car  $\leftarrow$
- 3 - motorcycle
- 4 - background  $\leftarrow$

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

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



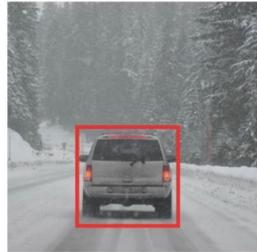
$x =$



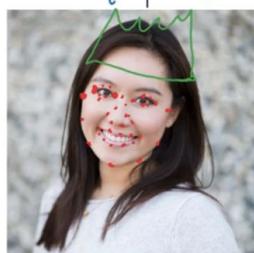
Andrew Ng

If object exists, then all square or log-loss errors for all 8 elements for bounding boxes as well as if object exists or not etc. If object is not there, then rest of the elements are dont care, therefore loss is only with  $P_c$  ( $P_c \rightarrow$  Does object exist in the image?)

## Landmark detection



$b_x, b_y, b_h, b_w$



$l_{1x}, l_{1y},$   
 $l_{2x}, l_{2y},$   
 $l_{3x}, l_{3y},$   
 $l_{4x}, l_{4y},$   
 $\vdots$   
 $l_{64x}, l_{64y}$

$\left\{ \begin{array}{c} x, y \\ \vdots \\ x, y \end{array} \right\}$

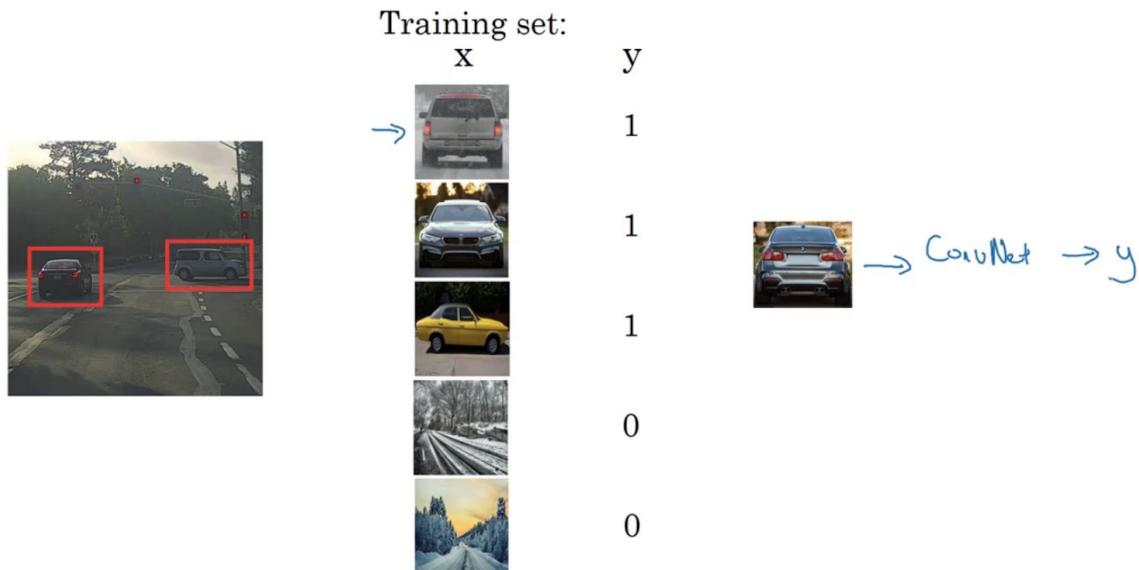


$l_{1x}, l_{1y},$   
 $\vdots$   
 $l_{32x}, l_{32y}$

Andrew Ng

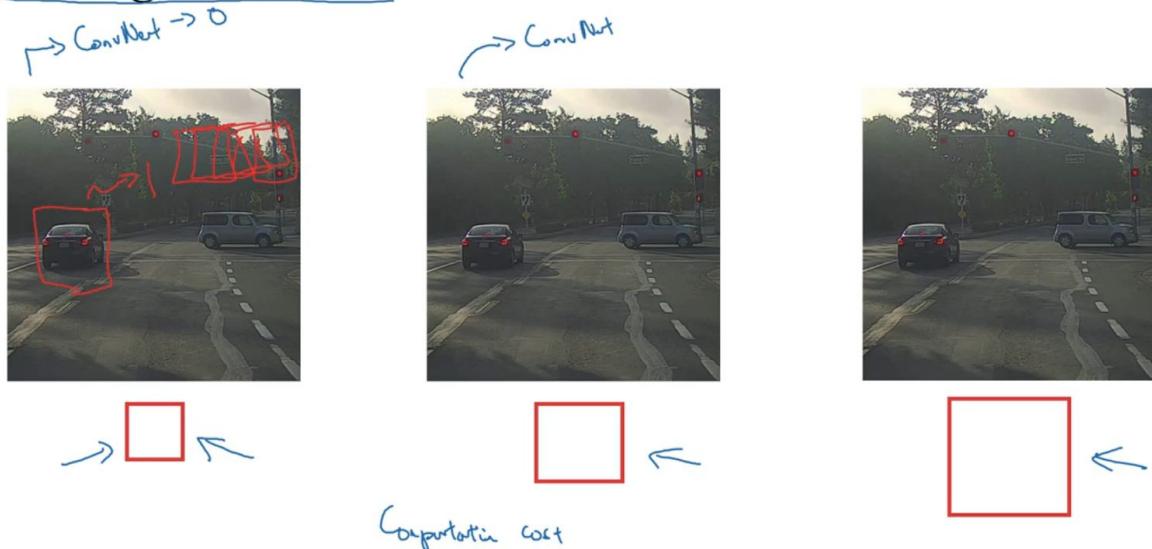
SnapChat AR filters for creating a crown etc needs landmarks on face to identify where various face things appear.

## Car detection example



Just train a convnet first on a closely cropped image of car. Then use sliding window to identify all possible cars in an image.

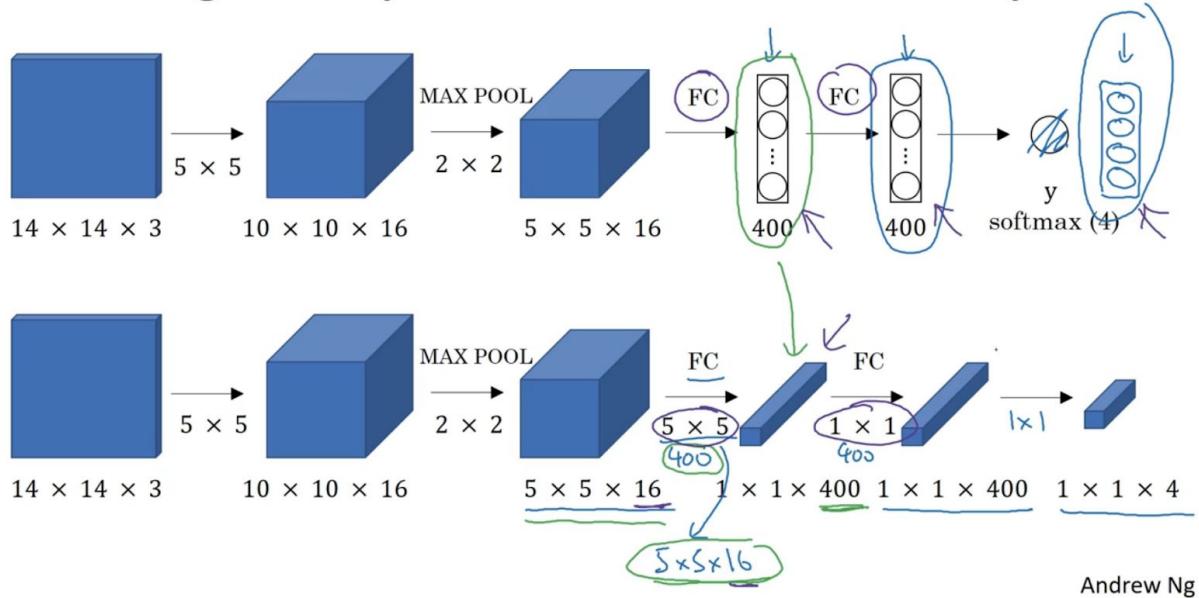
## Sliding windows detection



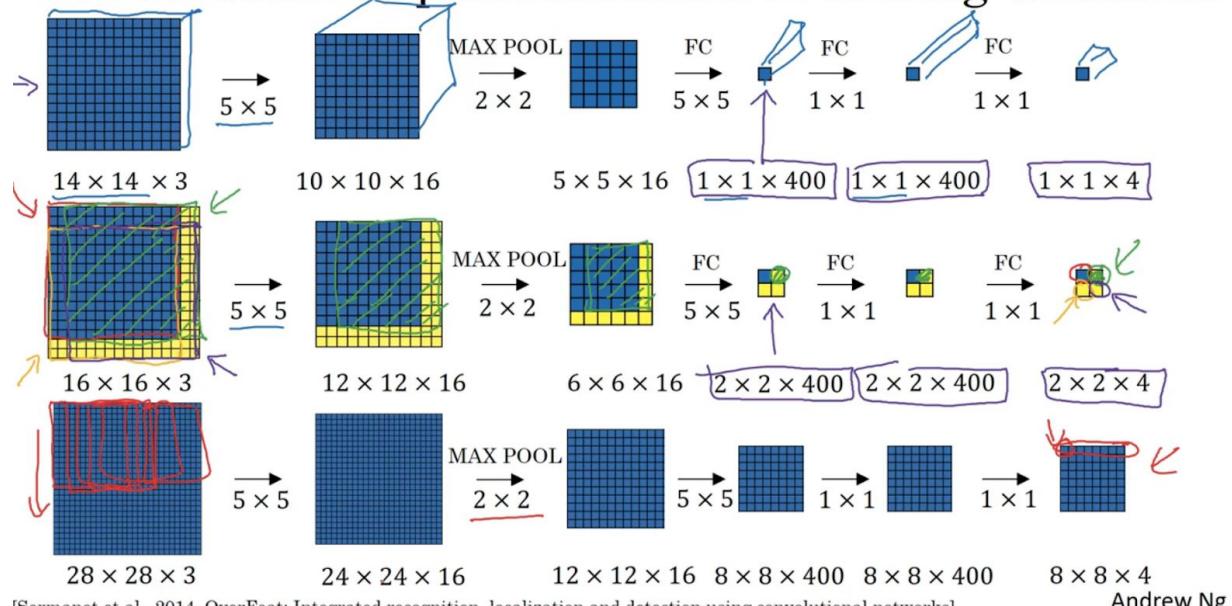
COMPUTATION COST is very high. Previously linear classifiers was ok to implement so many times, but now neural networks are more expensive.

Converting FC to Convolutional layers:

## Turning FC layer into convolutional layers



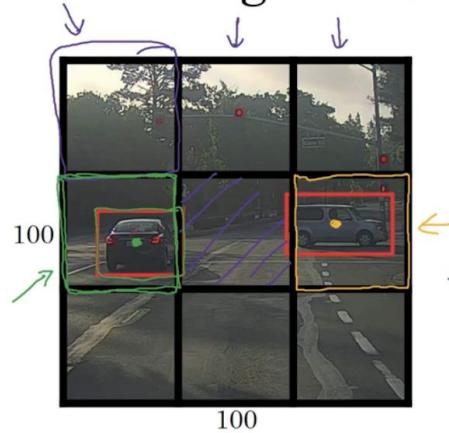
## Convolution implementation of sliding windows



Instead of repeating computation several times, you can just run on the whole image, and now instead of seeing FC as flat layer but see it as CNN. Weakness is that this method can't catch bounding box very well.

Just use the last 4 boxes and softmax output and see which output for every box.

# YOLO algorithm



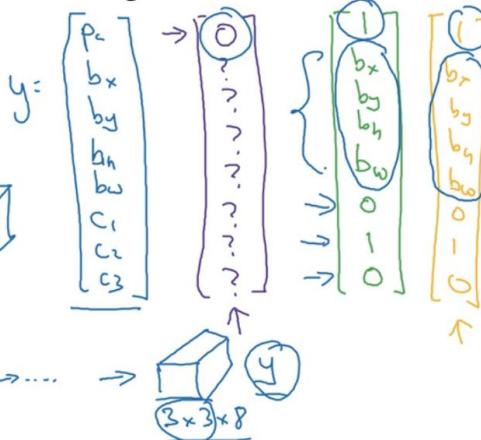
Target output:

$$3 \times 3 \times 8$$

$$\rightarrow X \rightarrow \text{WN} \rightarrow \text{MAX POOL} \rightarrow \dots$$

Labels for training

For each grid cell:

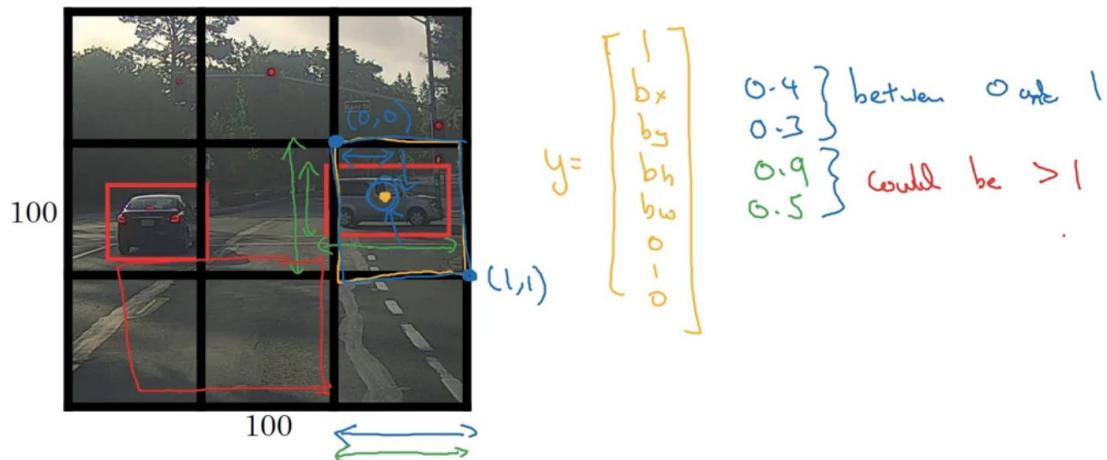


[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

See where the mid point of the object lies in a grid. In real life the grid is 19x19. Its not running it 9 times, but its conv implementation where you re-use parameters.

## Specify the bounding boxes



YOU ONLY LOOK ONCE or Yolo algo is really popular for object detection.

# Evaluating object localization



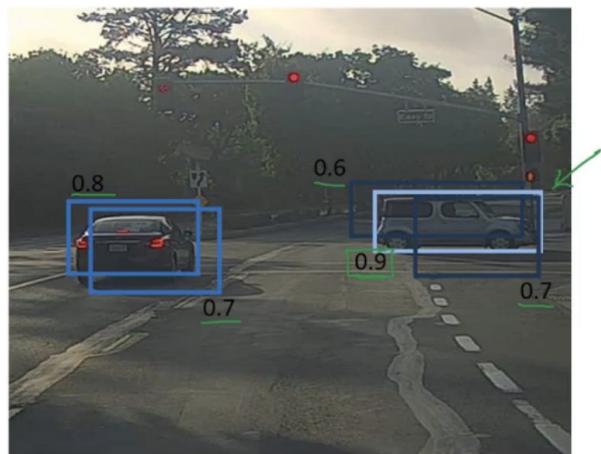
$$\text{Intersection over Union (IoU)}$$
$$= \frac{\text{Size of } \cap}{\text{Size of } \cup}$$

"Correct" if  $\underline{\text{IoU} \geq 0.5}$  ←  
0.6 ←

More generally, IoU is a measure of the overlap between two bounding boxes.

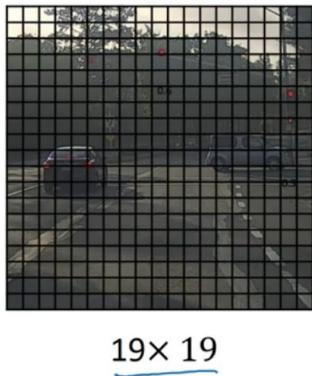
Andrew Ng

## Non-max suppression example

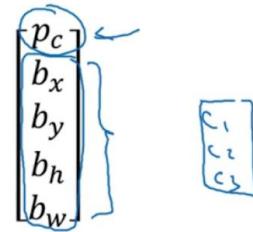


Remove boxes with non max probabilities, its non max suppression. Use IOU to find closely overlapped boxes and remove the non max ones and keep only the max one.

# Non-max suppression algorithm



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

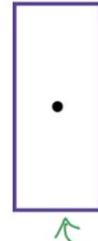
You need to do non max suppression for each classes, each for pedestrian, car etc.

## Overlapping objects:



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

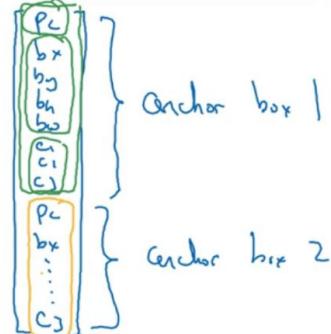
Anchor box 1:



Anchor box 2:



$$y =$$



[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

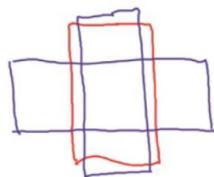
Andrew Ng

# Anchor box algorithm

Previously:

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

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



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.

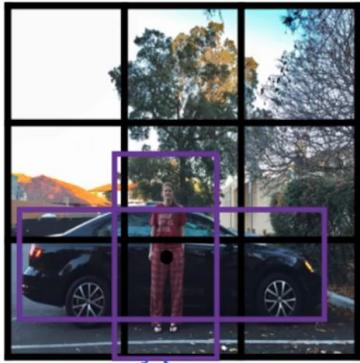
(grid cell, anchor box)

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

Andrew Ng

2 anchor boxes for each grid, find which anchor box is most IOU with the ground truth.

## Anchor box example



Anchor box 1:    Anchor box 2:



y =

$p_c$	1	car only?
$b_x$	bx	?
$b_y$	by	?
$b_h$	bh	?
$b_w$	bw	?
$c_1$	1	?
$c_2$	0	?
$c_3$	0	?
$p_c$	1	?
$b_x$	bx	1
$b_y$	by	bx
$b_h$	bh	by
$b_w$	bw	bh
$c_1$	0	bw
$c_2$	1	0
$c_3$	0	0

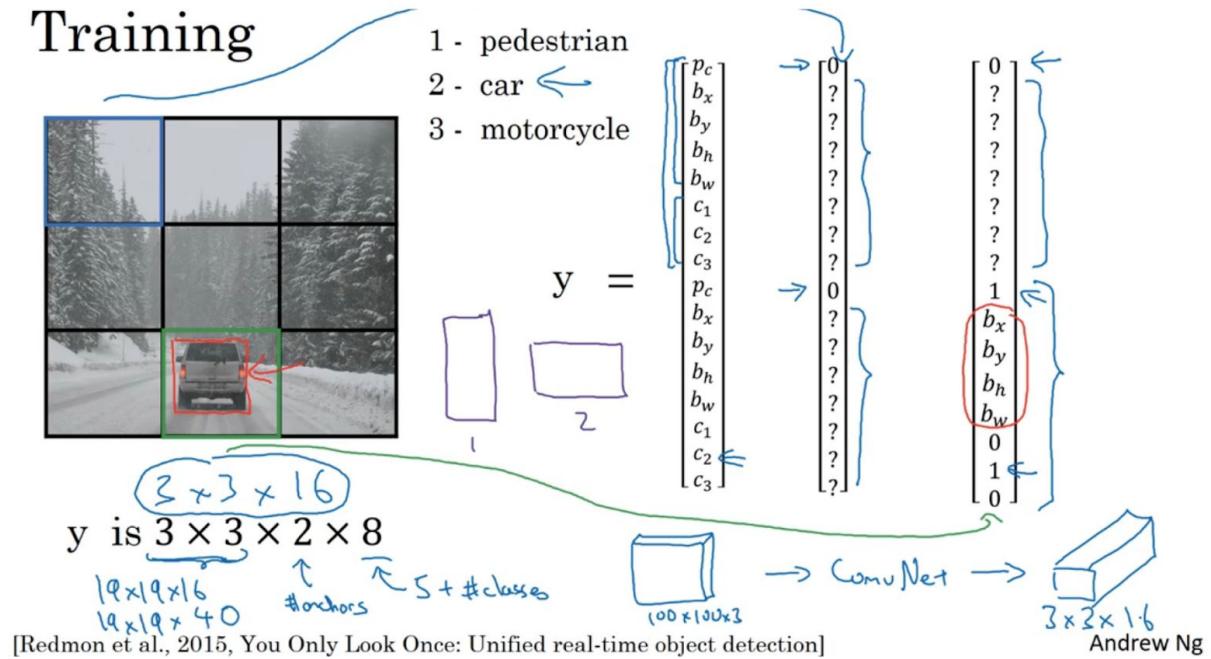
Andrew Ng

Can't handle more than 2 objects in one grid --- but usually 19x19 doesn't have overlapping objects in same grid.

People choose 5 or 10 anchor boxes.

YOLO ALGORITHM BELOW

## Training



## Outputting the non-max suppressed outputs



- For each grid cell, 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.

Another approach instead of YOLO: Find regions where to run sliding window CNN

## Region proposal: R-CNN



[Girshik et. al, 2013, Rich feature hierarchies for accurate object detection and semantic segmentation] Andrew Ng

## Faster algorithms

→ R-CNN: Propose regions. Classify proposed regions one at a time. Output label + bounding box. ↩

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.

[Girshik et. al, 2013. Rich feature hierarchies for accurate object detection and semantic segmentation]

[Girshik, 2015. Fast R-CNN]

[Ren et. al, 2016. Faster R-CNN: Towards real-time object detection with region proposal networks] Andrew Ng

## WEEK 4: 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

99%

99.9

## → Recognition

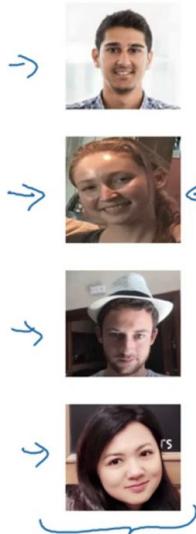
- Has a database of K persons
- Get an input image
- Output ID if the image is any of the K persons (or “not recognized”)

1:K

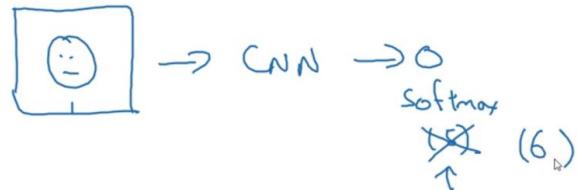
K=100 ←

Easier to verify if a person is what he is claiming to be, harder to do recognition. 99% accuracy is ok in verification, but recognition is not good at all with 99% if say there are 100 people, it might get 1 person wrong.

# One-shot learning



Learning from one example to recognize the person again

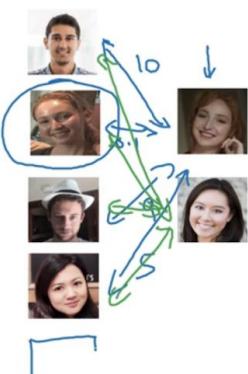


You only have one image for every employee.

# Learning a “similarity” function

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

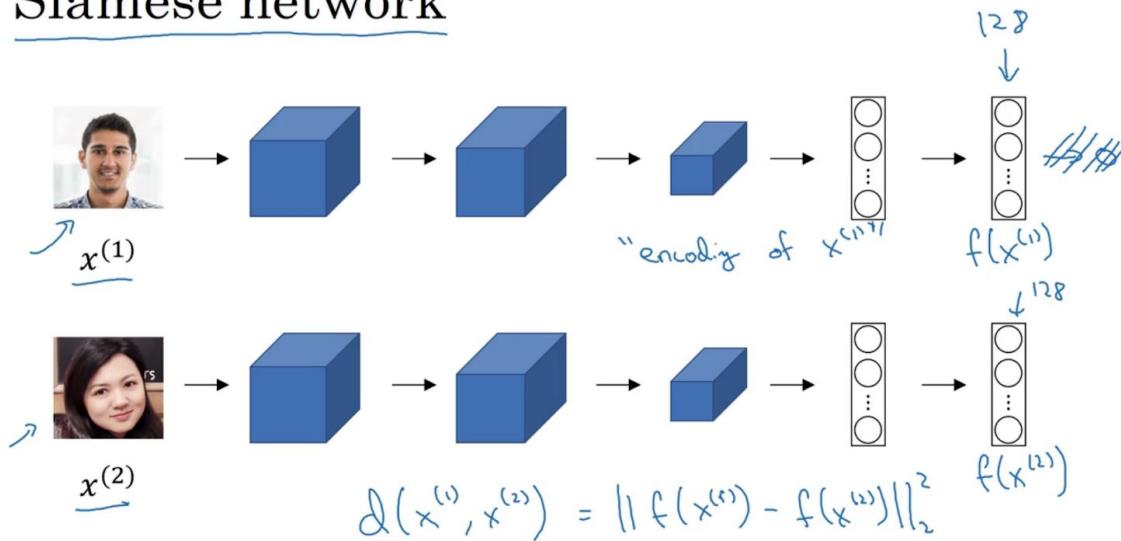
If $d(\text{img1}, \text{img2}) \leq \tau$	"same"	}	Verification.
$> \tau$	"different"		



$$d(\text{img1}, \text{img2})$$

You can add a new employee easily if you can find a similarity function of two images.

## Siamese network

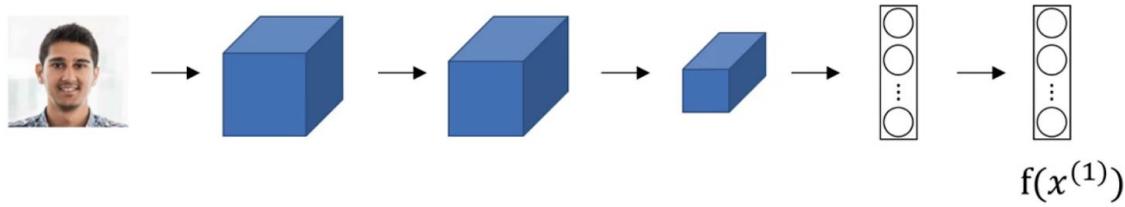


[Taigman et. al., 2014. DeepFace closing the gap to human level performance]

Andrew Ng

Take a norm of the two 128d vectors and see how similar are two images.

# Goal of learning



Parameters of NN define an encoding  $f(x^{(i)})$

Learn parameters so that:

If  $x^{(i)}, x^{(j)}$  are the same person,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is small.

If  $x^{(i)}, x^{(j)}$  are different persons,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is large.

## TRIPLET LOSS

### Learning Objective



$$\text{Want: } \underbrace{\frac{\|f(A) - f(P)\|^2}{d(A, P)}}_{\geq 0} + \underline{\alpha} \leq \underbrace{\frac{\|f(A) - f(N)\|^2}{d(A, N)}}_{\geq 0}$$

$$\frac{\|f(A) - f(P)\|^2}{\underline{0}} - \frac{\|f(A) - f(N)\|^2}{\underline{0}} + \underline{\alpha} \leq \underline{\alpha} \quad \text{Margin} \quad f(\text{img}) = \vec{0}$$

[Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering]

Andrew Ng

## Loss function

Given 3 images:  $A, P, N$ :

$$L(A, P, N) = \max \left( \frac{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \lambda}{\epsilon \rho > 0}, 0 \right)$$

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

Training set: 10k pictures of 1k persons

Atleast while training need pairs of anchors and positive examples. Once trained can use it for one shot learning case with only 1 picture for an employee, but while training definitely needs more than 1 image per person.

## Choosing the triplets $A, P, N$

During training, if  $A, P, N$  are chosen randomly,  
 $d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

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

Choose triplets that're "hard" to train on.

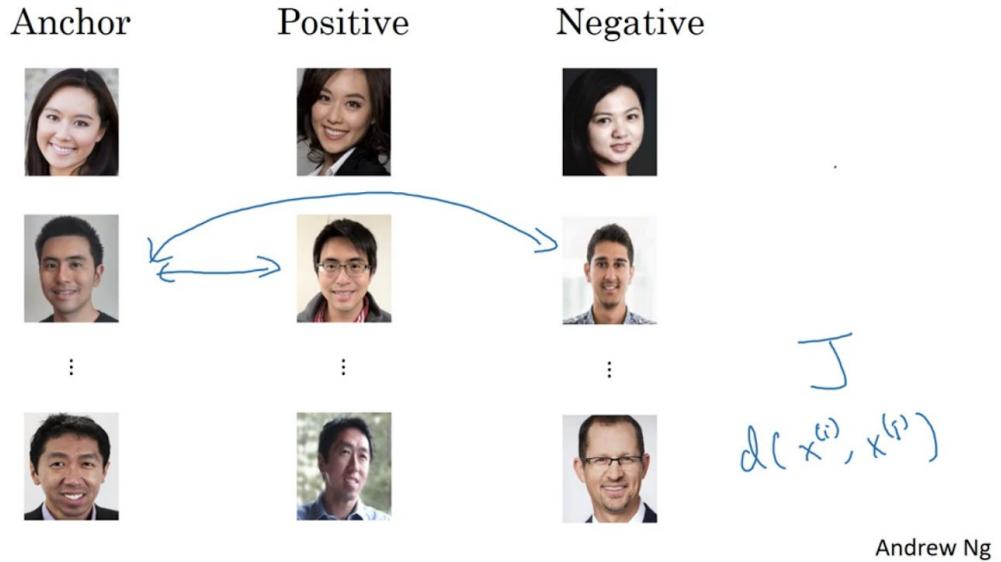
$$\begin{aligned} \lambda \approx & d(A, P) \\ \frac{\lambda}{d(A, P)} & \approx \frac{d(A, N)}{d(A, P)} \end{aligned}$$

[Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering]

Andrew Ng

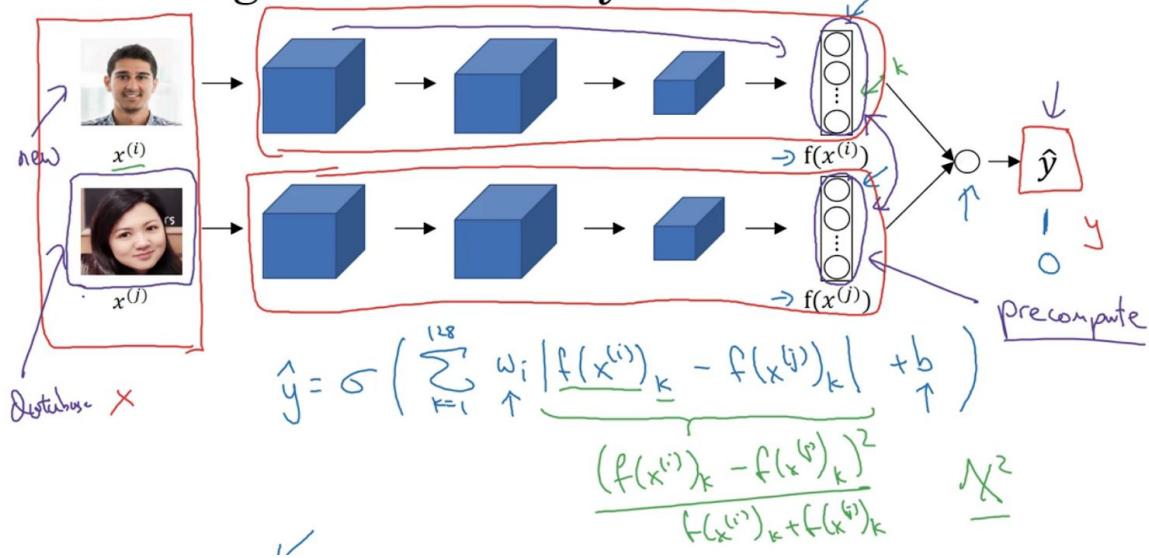
Have to choose triplets that are hard to train.

# Training set using triplet loss



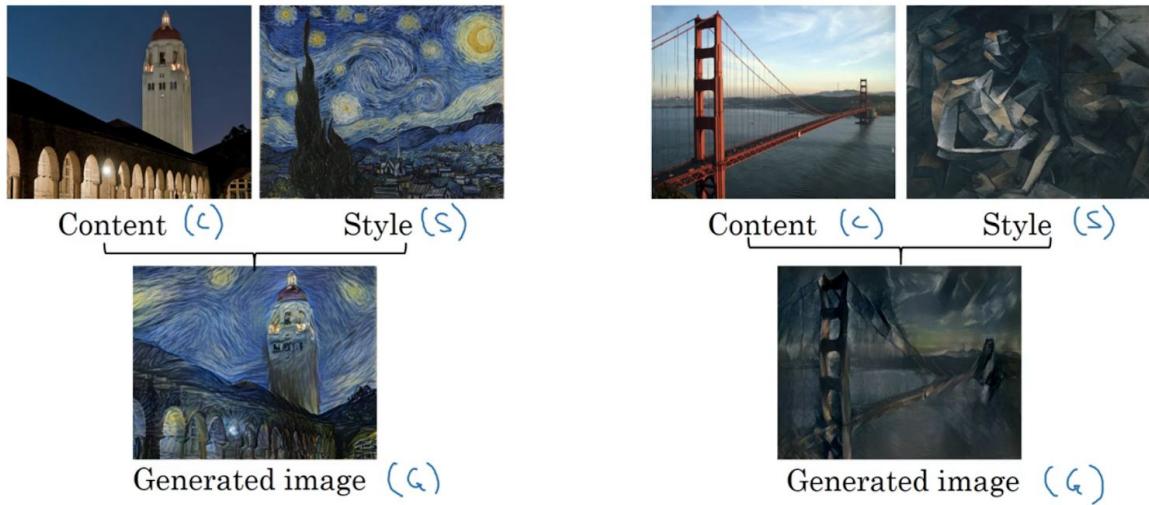
You can also use logistic unit for saying 1 if pictures are similar else 0. You can also keep pre-computed values of employees in database - so you don't lose time.

## Learning the similarity function



## NEURAL STYLE TRANSFER

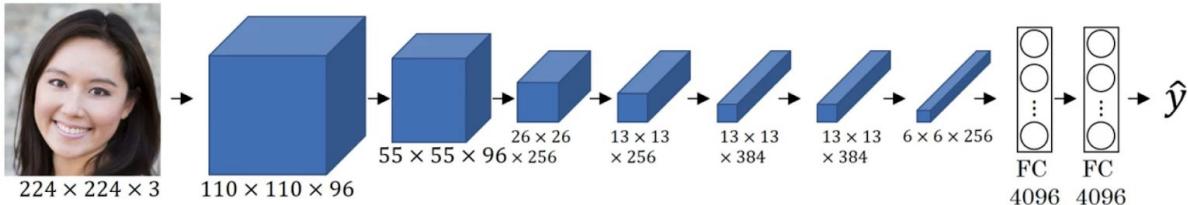
# Neural style transfer



[Images generated by Justin Johnson]

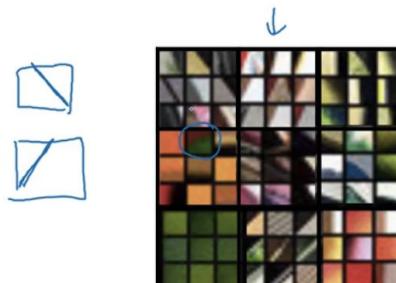
Andrew Ng

## Visualizing what a deep network is learning



Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

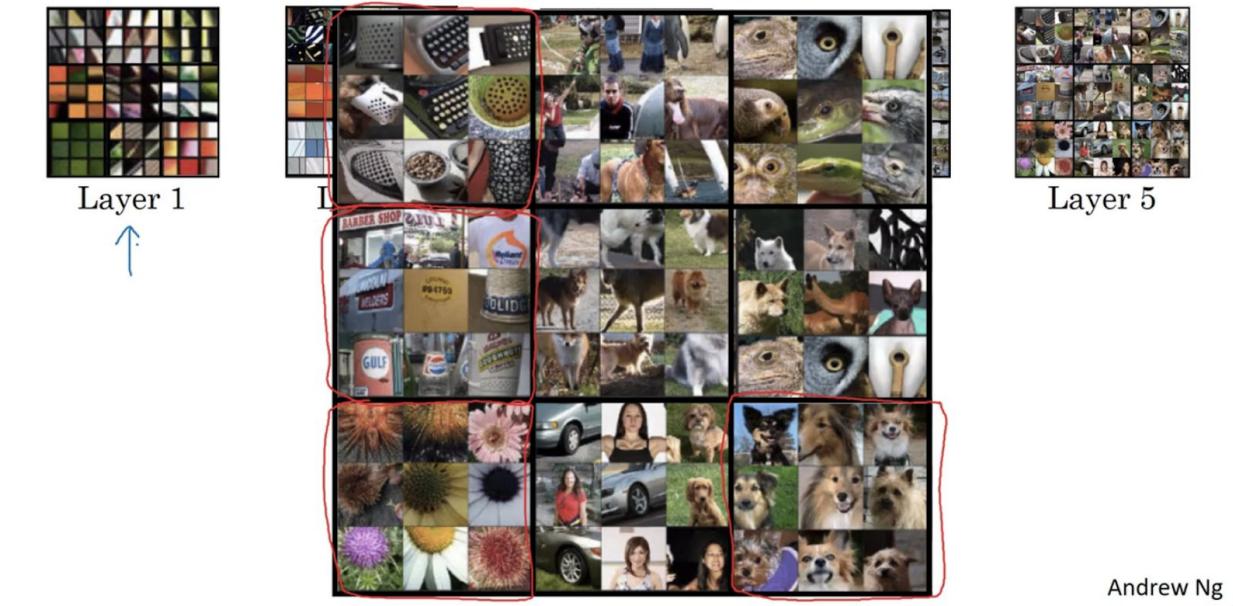
Repeat for other units.



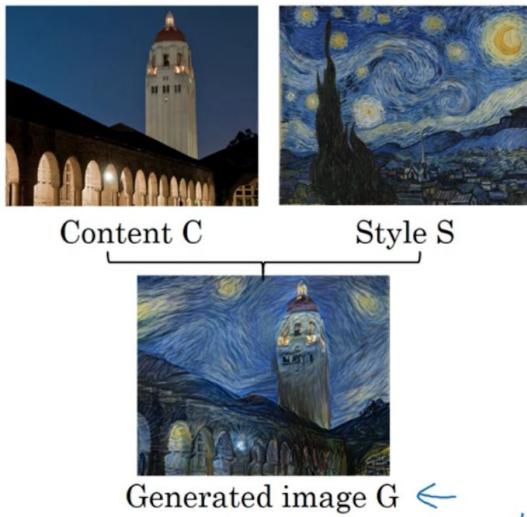
[Zeiler and Fergus., 2013, Visualizing and understanding convolutional networks]

Andrew Ng

## Visualizing deep layers: Layer 5



## Neural style transfer cost function



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

[Gatys et al., 2015. A neural algorithm of artistic style. Images on slide generated by Justin Johnson] Andrew Ng

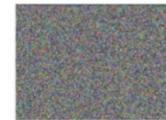
# Find the generated image G

1. Initiate G randomly

$$G: \underbrace{100 \times 100}_{\text{RGB}} \times 3$$

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

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



[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

## Content cost function

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

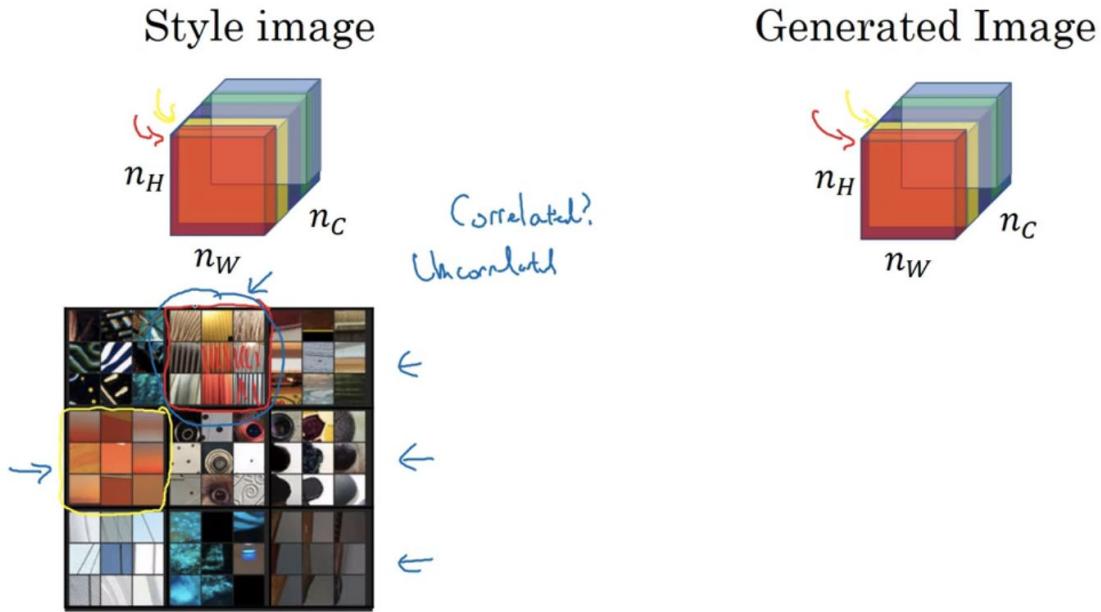
- Say you use hidden layer  $l$  to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let  $a^{[l](C)}$  and  $a^{[l](G)}$  be the activation of layer  $l$  on the images
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content

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

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

# Intuition about style of an image



How co-related are activations across channels? This will help get style cost function.

## Style matrix

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

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

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

$\downarrow \quad \downarrow \quad \downarrow$   
 $n_c$   
 $G_{kk'}^{[l](s)}$   
 $k, k' = 1, \dots, n_c^{[l]}$

"Gram matrix"

$$\begin{aligned} J_{\text{style}}^{[l]}(S, G) &= \frac{1}{C} \| G^{[l](s)} - G^{[l](G)} \|_F^2 \\ &= \frac{1}{(2n_H n_W n_C)^2} \sum_k \sum_{k'} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](G)})^2 \end{aligned}$$

## Style cost function

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

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

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

## Style cost function

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

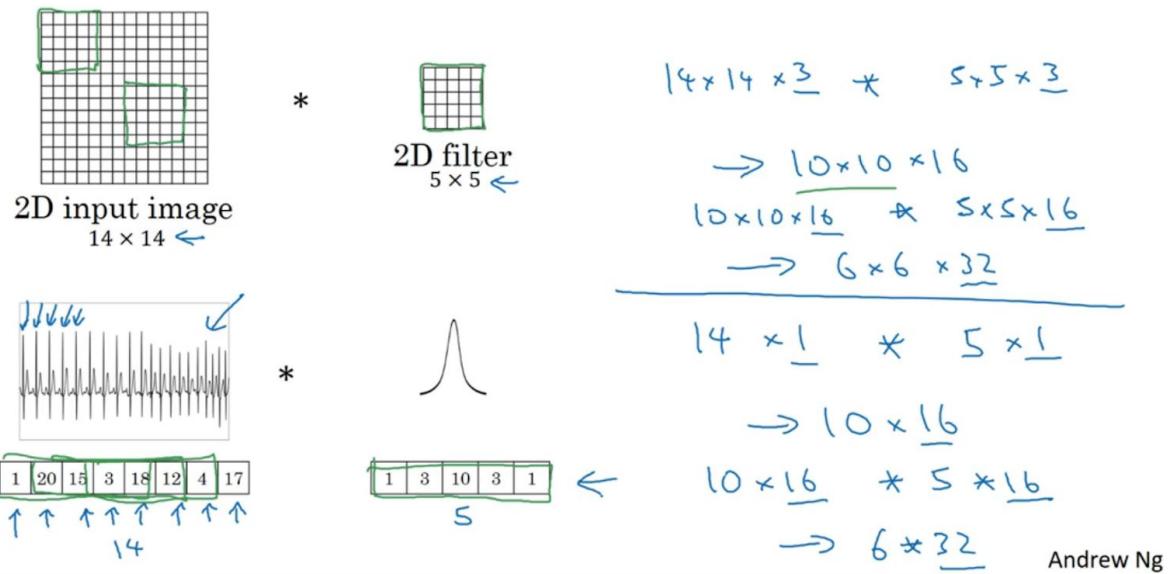
$$J_{style}^{[l]}(S, G) = \frac{1}{(2n_H^{[l]} n_W^{[l]} n_C^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})$$

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

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

G

## Convolutions in 2D and 1D



Sequence models for 1D data → More like RNN but can also use CNN for 1D data

## 3D convolution

