

## 26 - Types of neural networks used in supervised learning

### Supervised Learning

Input(x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1.....1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

## 31 - Backpropagation working 2

If the "W" value is changed to 4, notice the errors:

Input	Desired Output	Model Output (W=3)	Absolute Error	Square Error	Model Output (W=4)	Square Error
0	0	0	0	0	0	0
1	2	3	1	1	4	4
2	4	6	2	4	8	16

Backpropagation trains a neural network by assigning random weights to the algorithms and analyzing where the error in the system increases. When errors occur, the difference between the model output and the actual output is calculated. Once calculated, a different weight is assigned and the system is run again, to see if the error is minimized. If the error is not minimized, then an update of parameters is required

To update parameters, weights and biases are adjusted. Biases are located after weights and are in a different layer of a network, always being assigned the value of 1. After the parameters are updated, the process is run again. Once the error is at a minimum, the model is ready to start predicting.

In looking at the diagram below, if "W" also known as weight is changed, then the error of the system goes up, if "W" is changed into a smaller number the error goes down. Once the error is as close to zero as possible that weight is set as the parameter and the model can start predicting.

## 31 - computational graph

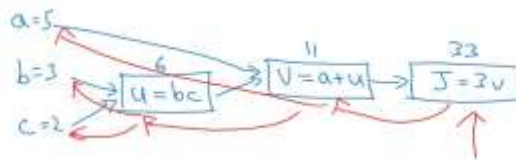
### Computation Graph

$$J(a,b,c) = J(a+bc) = J(5+3*2) = 33$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



## 34 - vectorization 3

## Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = w^T x^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$z = [z^{(1)} \ z^{(2)} \ z^{(3)}] = w^T X + [b \ b \ b]$$

$$z = np.dot(w.T, X) + b$$

$$A = [a^{(1)} \ a^{(2)} \ a^{(3)}] = \sigma(z)$$

Andrew Ng

## 31 - Backpropagation working

## How Backpropagation Algorithms Work

Note the output of the model when the "W" value is 3:

Input	Desired Output	Model Output (W=3)
0	0	0
1	2	3
2	4	6

Notice the difference between the actual output and the desired output:

Input	Desired Output	Model Output (W=3)	Absolute Error	Square Error
0	0	0	0	0
1	2	3	1	1
2	4	6	2	4

## 34 - vectorization

$$z = w^T x + b$$

Non-vectorized:

$$z = 0$$

for i in range(n):

$$z += w[i] * x[i]$$

$$z += b$$

Vectorized:

$$z = np.dot(w, x) + b$$

GPU } SIMD - single instruction multiple data

CPU }

## 34 - vectorization 2

```

a = np.random.rand(1000000)
b = np.random.rand(1000000)

tic = time.time()
c = np.dot(a,b)
toc = time.time()

print(c)
print("Vectorized version:" + str(1000*(toc-tic)) + "ms")

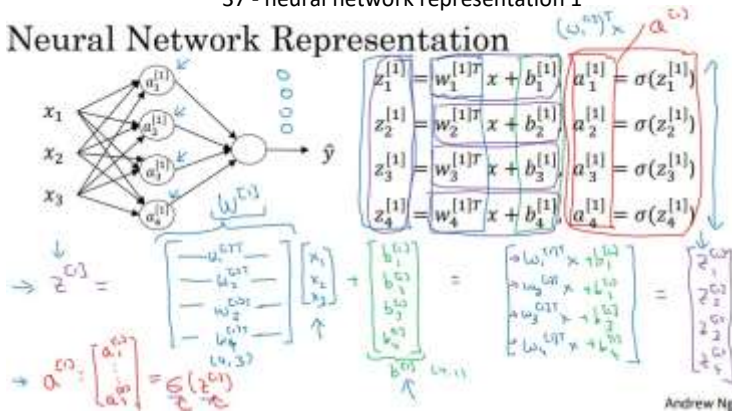
c = 0
tic = time.time()
for i in range(1000000):
    c += a[i]*b[i]
toc = time.time()

print(c)
print("For loop:" + str(1000*(toc-tic)) + "ms")

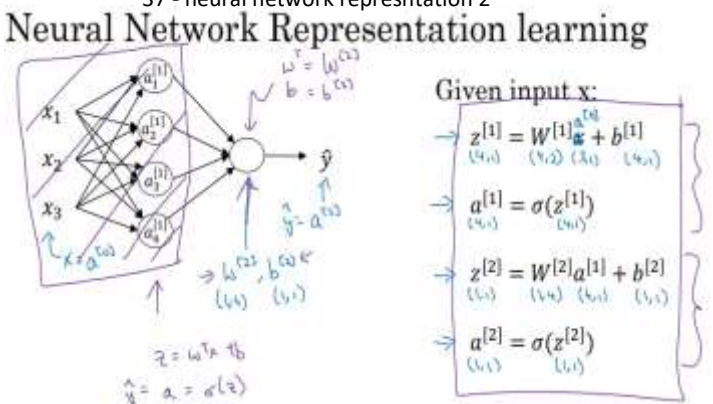
250286.989866
Vectorized version:1.5027523040771484ms
250286.989866
For loop:474.29513931274414ms

```

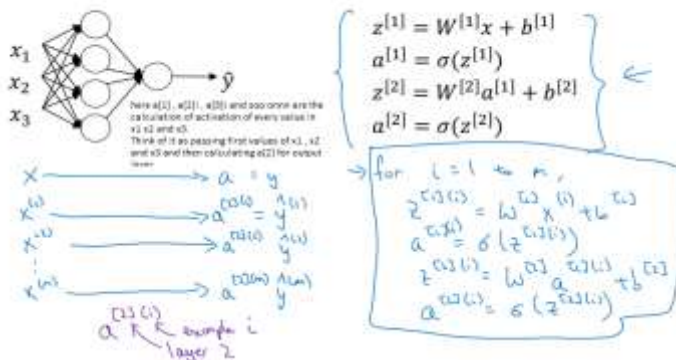
## Neural Network Representation



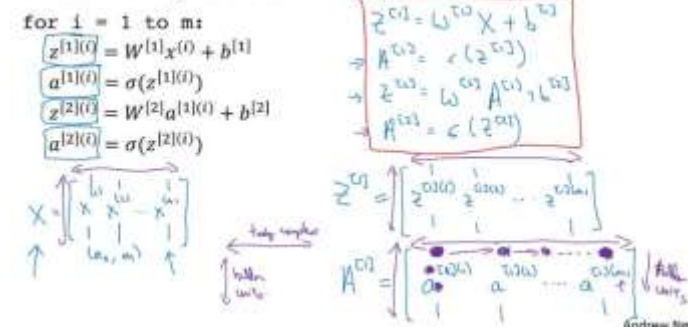
## Neural Network Representation learning



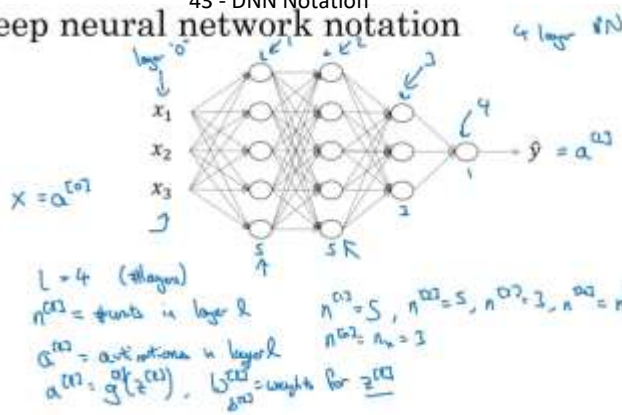
## Vectorizing across multiple examples



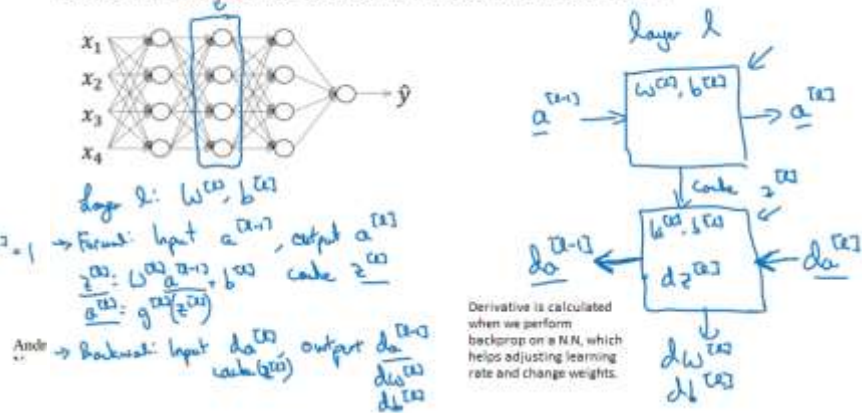
## Vectorizing across multiple examples



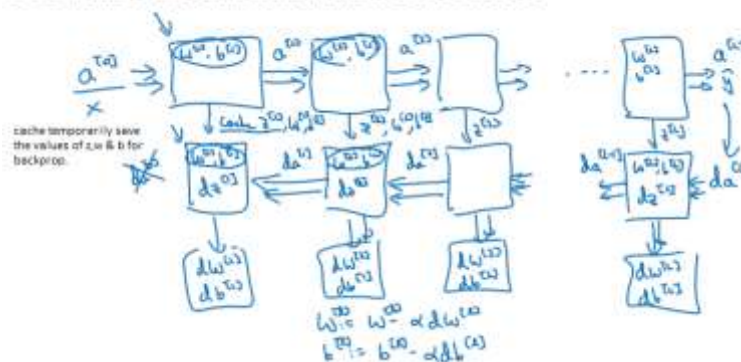
## Deep neural network notation



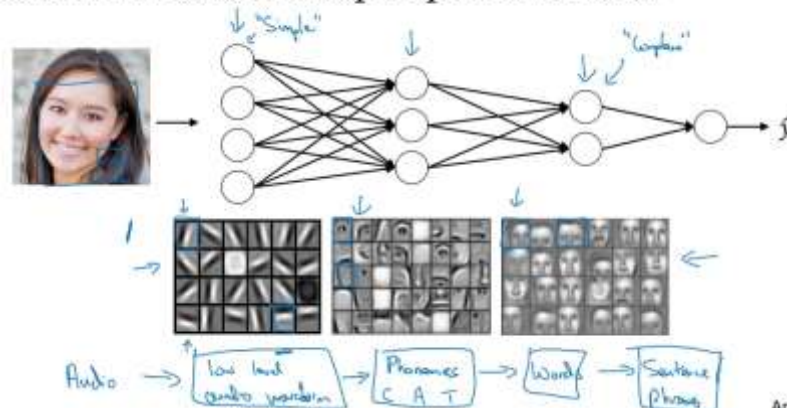
## Forward and backward functions



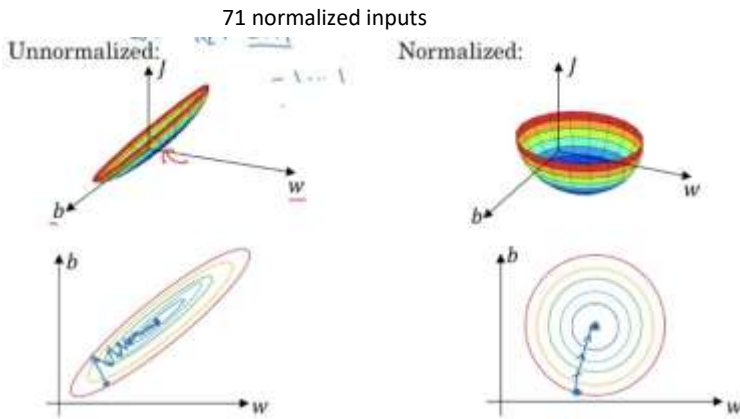
## Forward and backward functions



## Intuition about deep representation



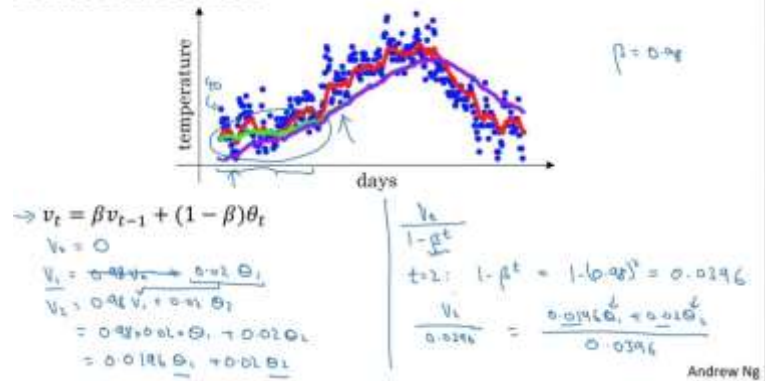




78 - gradient descent momentum. 1

## Bias correction

78 - bias correction in EMA



78 - gradient descent momentum

## Implementation details

$$v_{dw} = 0, v_{db} = 0$$

On iteration  $t$ :

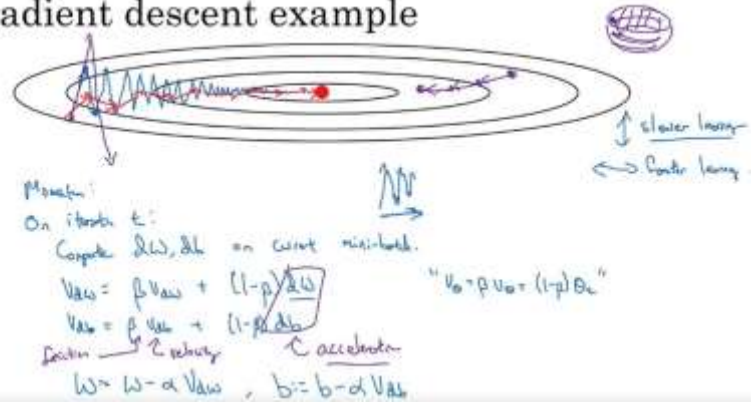
Compute  $dW, db$  on the current mini-batch

$$\begin{aligned} \rightarrow v_{dw} &= \beta v_{dw} + (1 - \beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) db \end{aligned} \quad \left| \quad v_{dw} = \beta v_{dw} + dW \right.$$

$$W = W - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta$   $\beta = 0.9$

## Gradient descent example



79 - Adam 1

## Hyperparameters choice:

- $\rightarrow \alpha$ : needs to be tune
- $\rightarrow \beta_1$ : 0.9  $\rightarrow (dW)$
- $\rightarrow \beta_2$ : 0.999  $\rightarrow (dW^2)$
- $\rightarrow \epsilon$ :  $10^{-8}$

Adam: Adaptive moment estimation

79 - Adam

## Adam optimization algorithm

$$v_{dw} = 0, s_{dw} = 0, v_{db} = 0, s_{db} = 0$$

On iteration  $t$ :

Compute  $dW, db$  using current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db \leftarrow \text{"moment"} \beta_1$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dW^2, \quad s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 \leftarrow \text{"RMSprop"} \beta_2$$

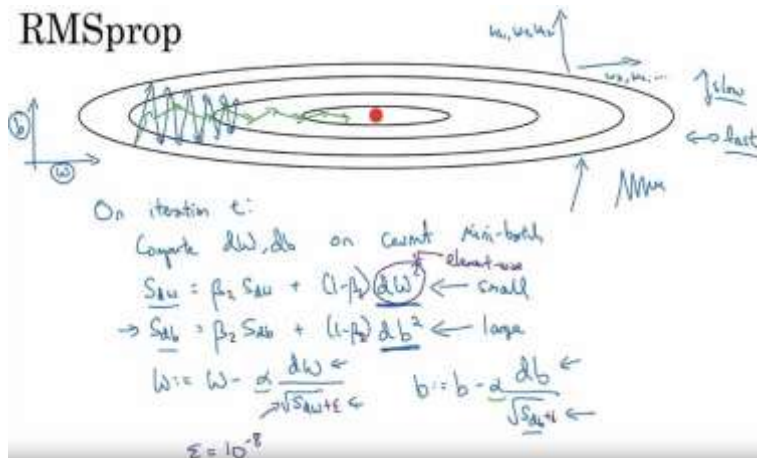
$$v_{dw}^{correct} = v_{dw} / (1 - \beta_1^t), \quad v_{db}^{correct} = v_{db} / (1 - \beta_1^t)$$

$$s_{dw}^{correct} = s_{dw} / (1 - \beta_2^t), \quad s_{db}^{correct} = s_{db} / (1 - \beta_2^t)$$

$$W = W - \alpha \frac{v_{dw}^{correct}}{\sqrt{s_{dw}^{correct} + \epsilon}}, \quad b = b - \alpha \frac{v_{db}^{correct}}{\sqrt{s_{db}^{correct} + \epsilon}}$$

79 - RMSProp

## RMSprop



80 - learning rate decay 1

## Learning rate decay

1 epoch = 1 pass through data

$$\alpha = \frac{\alpha_0}{1 + \text{decay-rate} \times \text{epoch-number}}$$

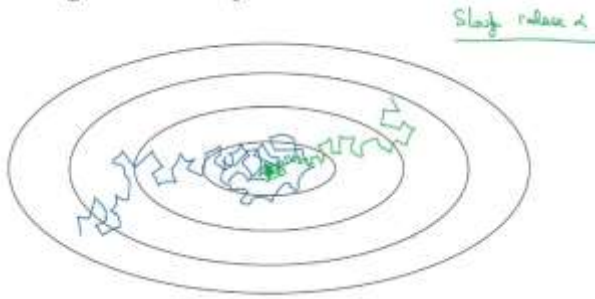
Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04
$\vdots$	$\vdots$

$$\alpha_0 = 0.2$$

$$\text{decay-rate} = 1$$



## Learning rate decay



## Picking hyperparameters at random

$$\rightarrow n^{T1} = 50, \dots, 100$$



$$\rightarrow \# \text{ layers } L: 2 - 4$$

$$2, 3, 4$$

## Hyperparameters

$$\beta = 0.9 \dots 0.999$$

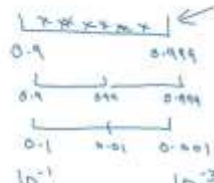
$$\downarrow$$

$$1 - \beta = 0.1 \dots 0.001$$

$$\beta: 0.999 \rightarrow 0.9995 \quad \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995 \quad \sim 1000$$

$$\frac{1}{1-\beta}$$

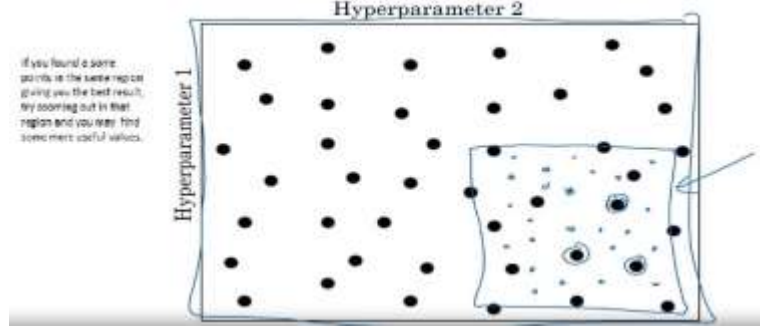


$$r \in [-3, -1]$$

$$1 - \beta = 10^{-5}$$

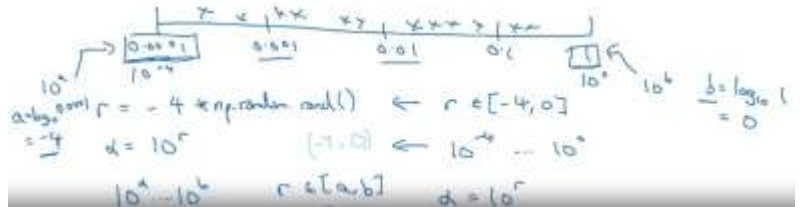
$$\beta = 1 - 10^{-5}$$

## Coarse to fine

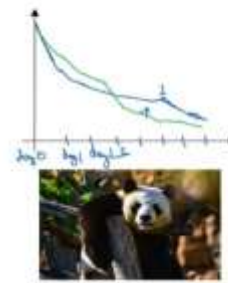


## Appropriate scale for hyperparameters

$$\alpha = 0.0001 \dots 1$$

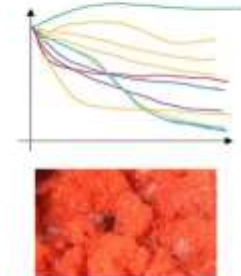


## Babysitting one model



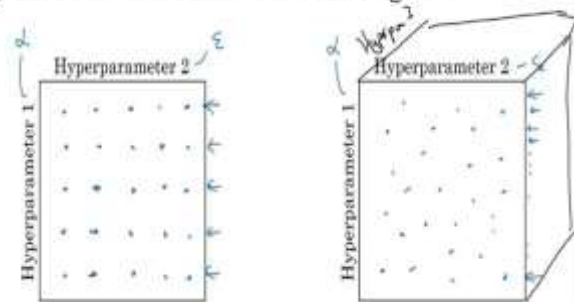
Panda

## Training many models in parallel

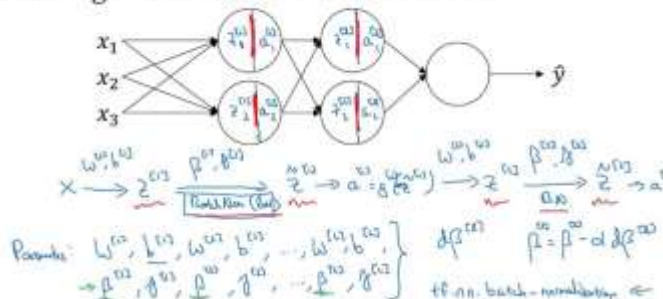


Caviar

## Try random values: Don't use a grid



## Adding Batch Norm to a network



## Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values  $z^{(i)}$  within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

## Implementing gradient descent

for  $t = 1 \dots \text{num\_Mini\_Batches}$   
 Compute Partial pop on  $X^{(t)}$ .  
 In each hidden layer, use BN to rep  $z^{(i)}$  with  $\hat{z}^{(i)}$ .  
 Use backprop to compute  $\frac{\partial L}{\partial W^{(i)}}$ ,  $\frac{\partial L}{\partial b^{(i)}}$ ,  $\frac{\partial L}{\partial z^{(i)}}$ .  
 Update params  $W^{(i)} := W^{(i)} - \alpha \frac{\partial L}{\partial W^{(i)}}$ ,  $b^{(i)} := b^{(i)} - \alpha \frac{\partial L}{\partial b^{(i)}}$ ,  $z^{(i)} := z^{(i)} - \alpha \frac{\partial L}{\partial z^{(i)}}$ .  
 Works w/ mom, RMSprop, Adam.



## 82 - implementation of batch norm

### Implementing Batch Norm

Given some intermediate values in NN  $z^{(L)} \dots z^{(L)}$

$$\mu = \frac{1}{m} \sum z^{(L)}$$

$$\sigma^2 = \frac{1}{m} \sum (z^{(L)} - \mu)^2$$

$$z_{\text{norm}}^{(L)} = \frac{z^{(L)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

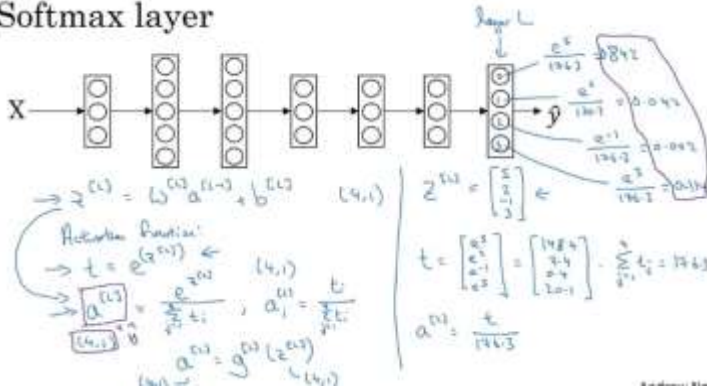
$$\hat{z}^{(L)} = \gamma z_{\text{norm}}^{(L)} + \beta$$

Learnable parameters  $\gamma, \beta$

Use  $\frac{1}{\sqrt{2}}$  instead of  $\frac{1}{\sqrt{2}}$

## 83 - Softmax

### Softmax layer



## 88 - dev and test for inappropriate images

### Another example

Algorithm A: 3% error  
Algorithm B: 5% error



If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.

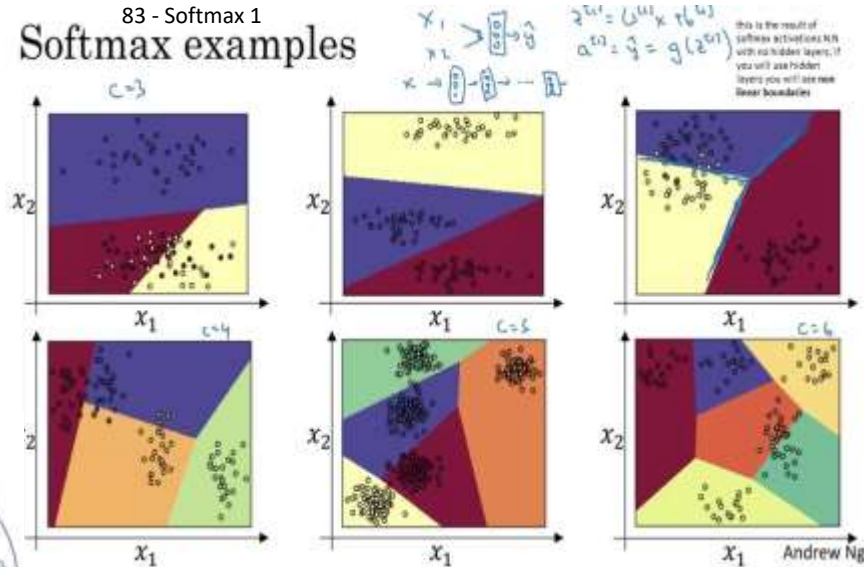
## 89 - assumptions related to Supervised Learning

### The two fundamental assumptions of supervised learning

- You can fit the training set pretty well.  $\approx$  Bias
- The training set performance generalizes pretty well to the dev/test set.  $\approx$  Variance

## 83 - Softmax 1

### Softmax examples



## 86 - cat classification example

### Another cat classification example

Classifier	Accuracy	Running time
A	90%	80ms
B	92%	95ms
C	95%	1,500ms

$$\text{Cost} = \text{accuracy} - 0.5 \times \text{Running Time}$$

Maximize accuracy  
Subject to Running Time  $\leq 100 \text{ ms}$

N metric: 1 optimizing  
N-1 satisfactory

Wabanda / Wager words  
Alma, OK Google  
Hey Sir, Nicholas  
12 43 百度

Accuracy  
False positive

Maximize accuracy  
s.t.  $\leq 1$  false positive  
every 24 hours

## 88 - human level performance

### Why compare to human-level performance

Humans are quite good at a lot of tasks. So long as ML is worse than humans, you can:

- Get labeled data from humans.  $(x, y)$
- Gain insight from manual error analysis: Why did a person get this right?
- Better analysis of bias/variance.

## 89 - bias-variance w human level error 1

### Summary of bias/variance with human-level performance

Human-level error  
(proxy for Bayes error)

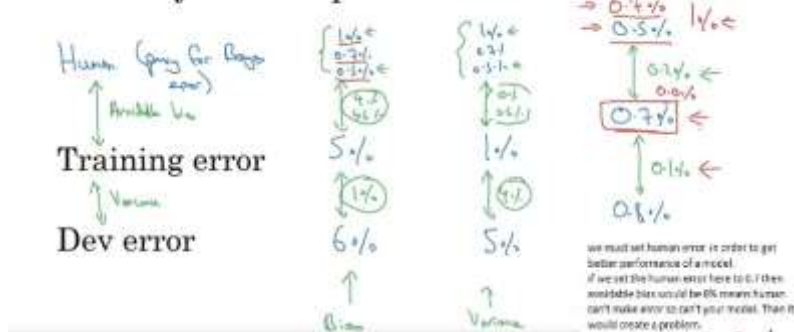
Training error

Dev error

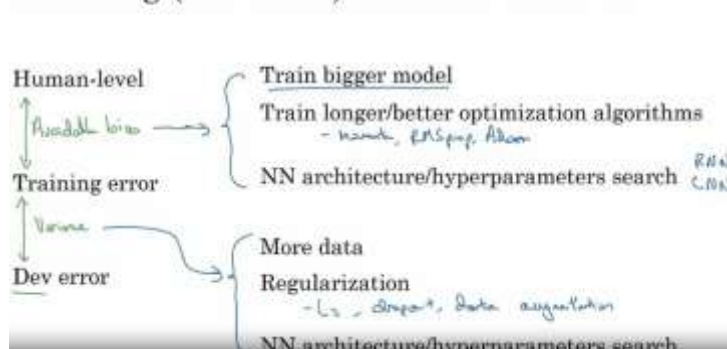
"Bias"  $\rightarrow$  "Variance"

Variance here defines how much of information is model giving off.

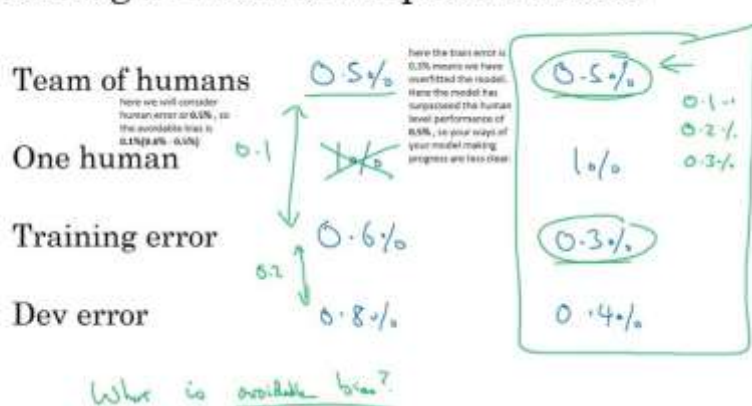
## 89 - bias-variance w human level error 2



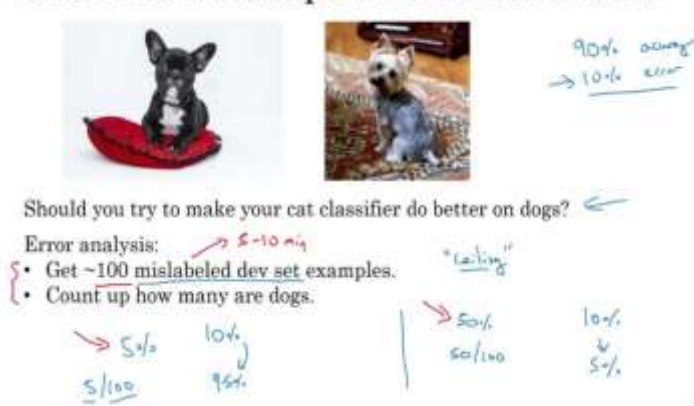
## 89 - reducing avoidable bias and variance



## 89 - surpassing human level performance



## 90 - error analysis



### 90 - error analysis 1

Image	Dog	Great Cats	Blurry	Instagram	Comments
1	✓			✓	Pitbull
2			✓	✓	
3		✓	✓		Rainy day at zoo
⋮	⋮	⋮	⋮	⋮	
% of total	8%	43%	61%	12%	

Key to is make a table and put all the pic name in the column you and on the left put all the number of images you are counting manually.  
 We can see that we have to majorly work on blurry & Great Cats

### 91 - incorrectly labeled data

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeler missed cat in background
99					
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	

Overall dev set error: 10%  
 Errors due incorrect labels: 0.6%  
 Errors due to other causes: 9.4%

Handwritten calculations:  
 2%  
 0.6%  
 1.4%  
 2.1%  
 1.9%

In this case the error is 2% so solving incorrect labeling will definitely help improving model's performance

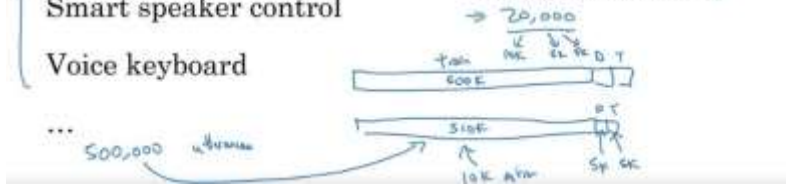
Goal of dev set is to help you select between two classifiers A & B.



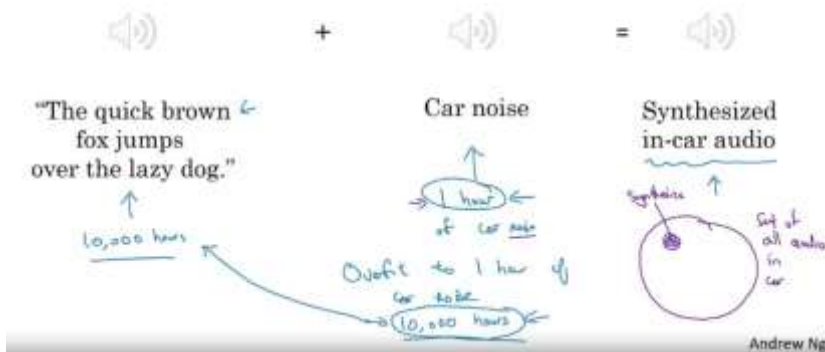
## Speech recognition example

Dev/test

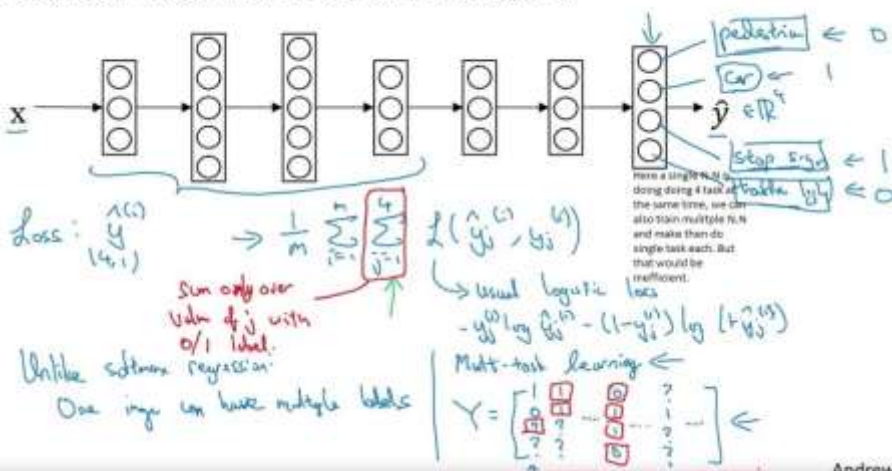
Speech activated }  
rearview mirror }



## Artificial data synthesis

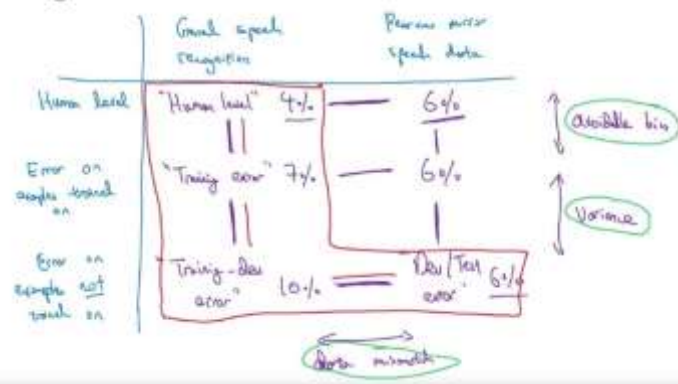


## Neural network architecture



Andrew

### More general formulation



### Simplified autonomous driving example



Pedestrians  $y^{(1)}$   
 Cars  $0$   
 Stop signs  $1$   
 Traffic light  $1$   
 ...  $0$   
 ...  $1$

$(4, n)$

$$= \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(n)} \\ y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(n)} \\ y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(n)} \\ y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(n)} \end{bmatrix}$$

$(4, n)$

### 99 - Multitask Learning 3

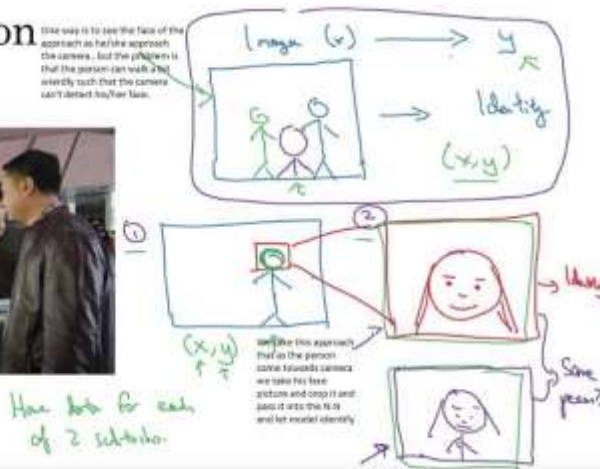
## When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
  - Usually: Amount of data you have for each task is quite similar.
- 
- The diagram illustrates multitask learning with shared lower-level features. On the left, a vertical stack shows task A with 1,000,000 data points and task B with 1,000 data points, connected by a downward arrow. On the right, a diagram shows two overlapping ellipses representing shared lower-level features. The top ellipse contains A<sub>1</sub>, 1,000 and A<sub>2</sub>, 1,000. The bottom ellipse contains A<sub>n</sub>, 1,000. A red arrow points from the top ellipse to the bottom one. A bracket on the right side of the ellipses is labeled 99,000.
- Can train a big enough neural network to do well on all the tasks.

## Face recognition



[Image courtesy of Baidu]



## Pros and cons of end-to-end deep learning

Pros:

- Let the data speak  $x \rightarrow y \rightarrow p$
- Less hand-designing of components needed

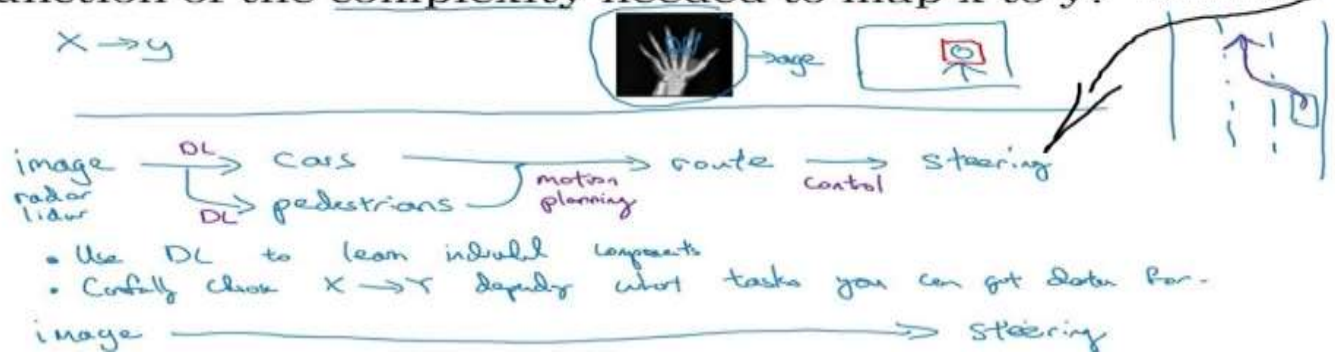
**Cons:**

- May need large amount of data  $X \rightarrow y$
- Excludes potentially useful hand-designed components Data      Hand-design

Date                      Homework

# Applying end-to-end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map  $x \rightarrow y$ ?



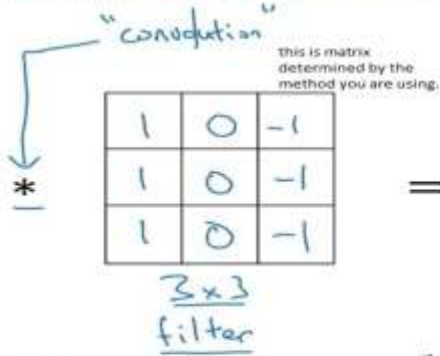
103 - Edge detection In C.N.N

## Vertical edge detection

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

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

If we have no padding then stride of 1 is applied (means we are not skipping any pixel).  
The Shape of the resultant image will be  $(6-1+1) \times (6-1+1)$

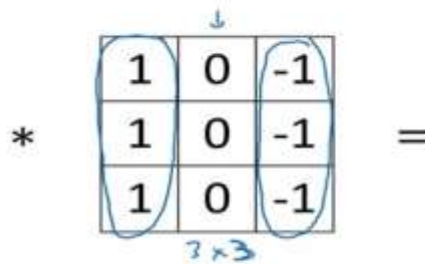
python: conv-forward  
tensorflow: tf.nn.conv2d  
keras: Conv2D

103 - Edge detection In C.N.N 1

## 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



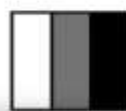
=

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

4x4

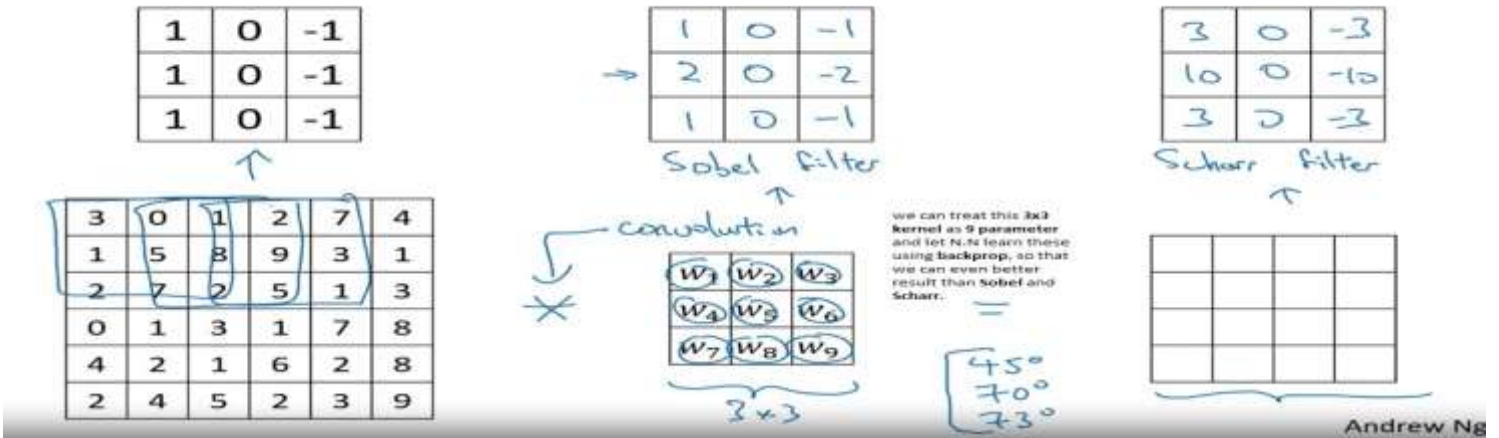


\*

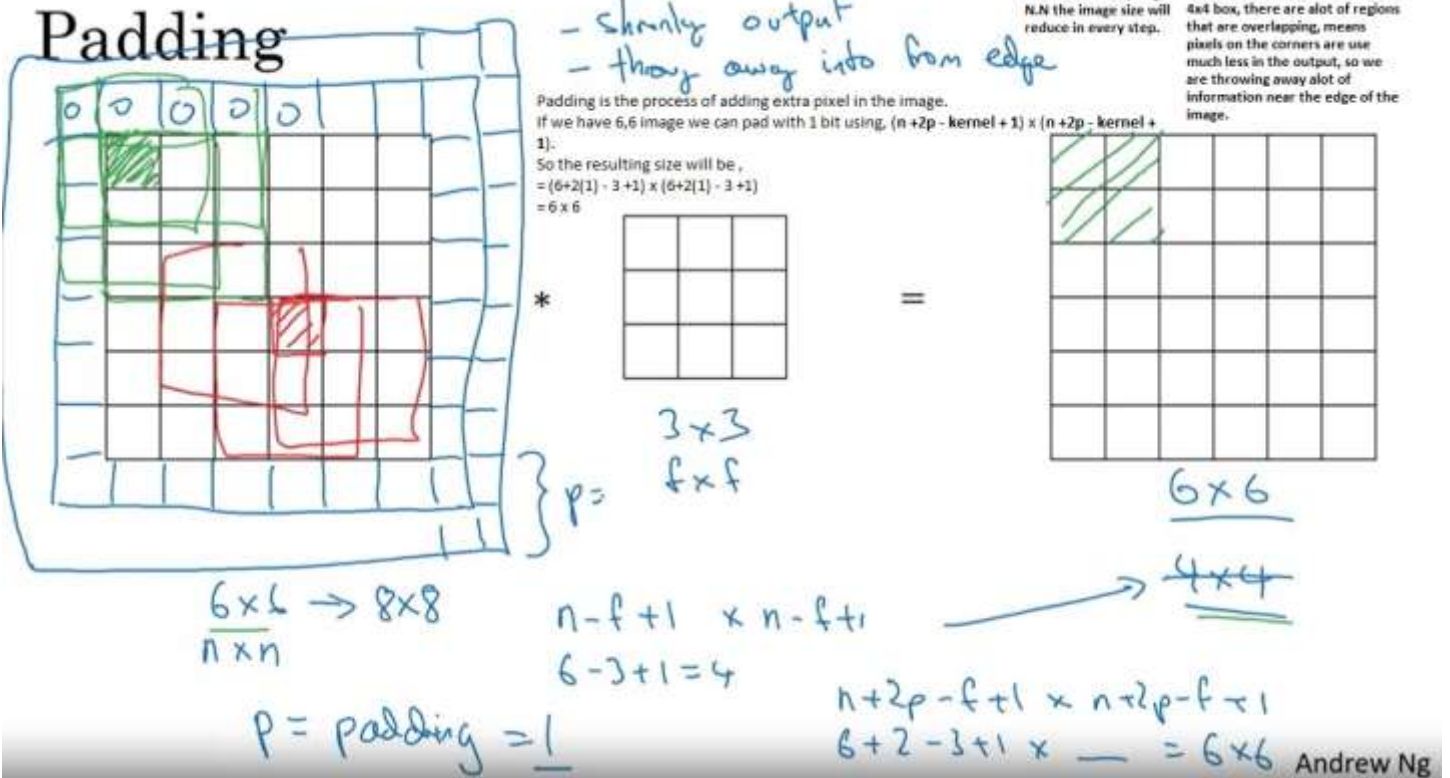




# Learning to detect edges



## 104 - padding



## 104 - padding 1

### Valid and Same convolutions

no padding

“Valid”:  $n \times n$   $\times$   $f \times f$   $\rightarrow$   $n - f + 1 \times n - f + 1$   
 $6 \times 6$   $\times$   $3 \times 3$   $\rightarrow$   $4 \times 4$

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

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

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

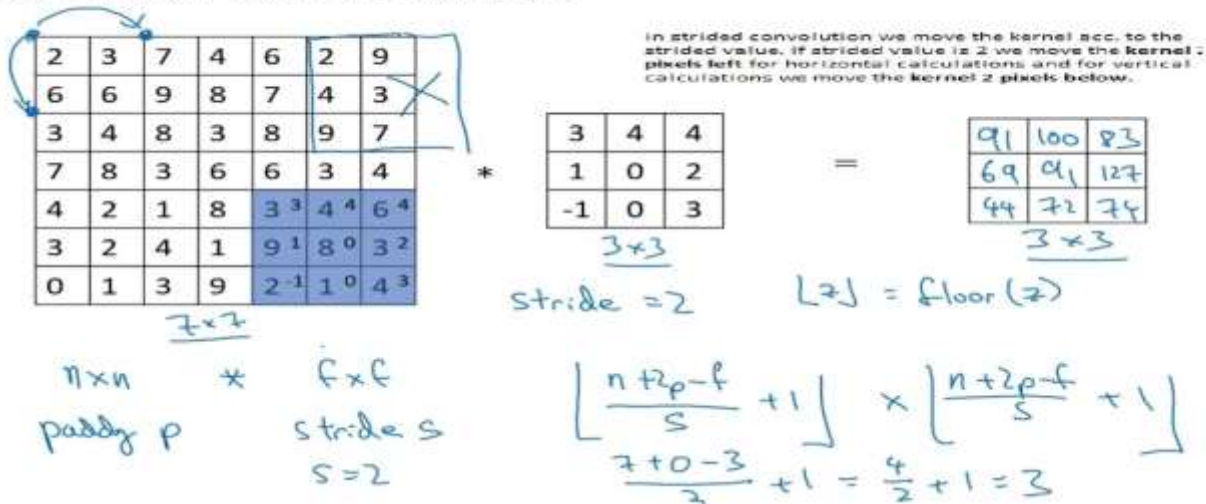
$$3 \times 3 \quad p = \frac{3-1}{2} = 1 \quad \begin{matrix} 5 \times 5 \\ f=5 \end{matrix} \quad p=2$$

f is usually odd

Andrew Ng

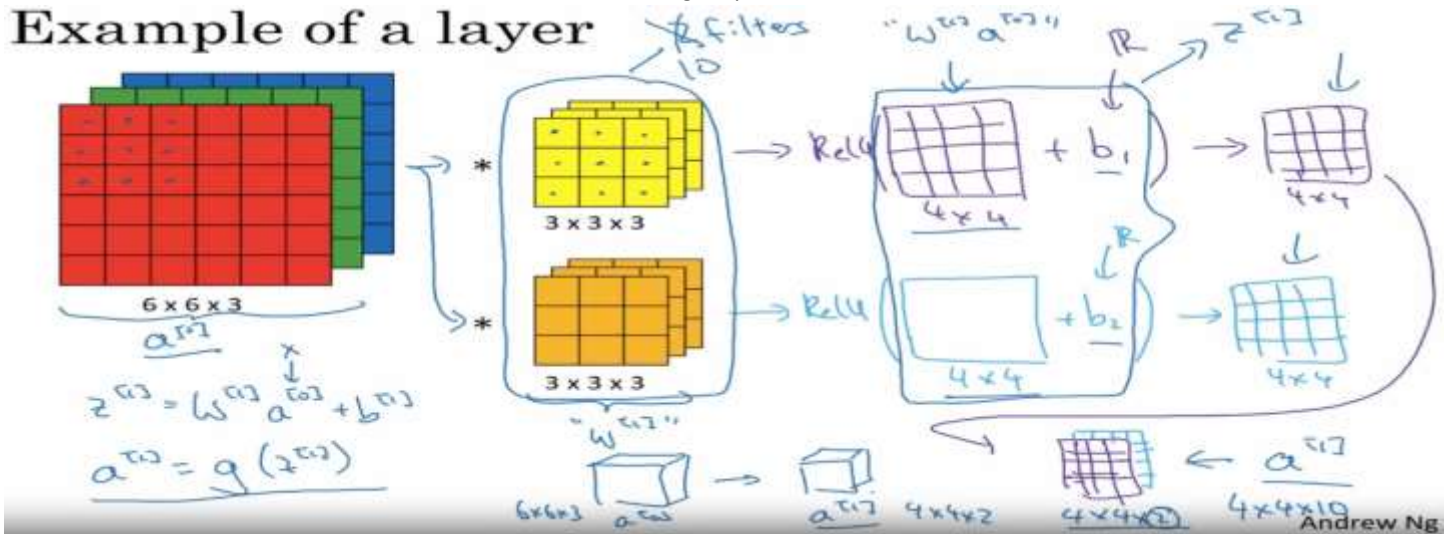
Strided convolution<sup>10</sup>

## 104 - strided convolutions





## Example of a layer



## 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

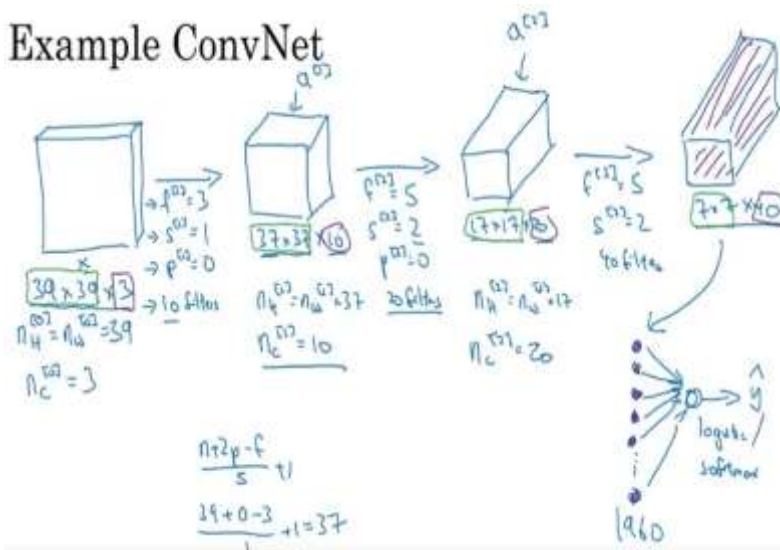
→ Each filter is:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$   
 Activations:  $a^{[l-1]} \rightarrow n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$   
 Weights:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$   
 bias:  $n_c^{[l]} = (1, 1, 1, n_c^{[l]})$  ← #filters in layer  $l$ .

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

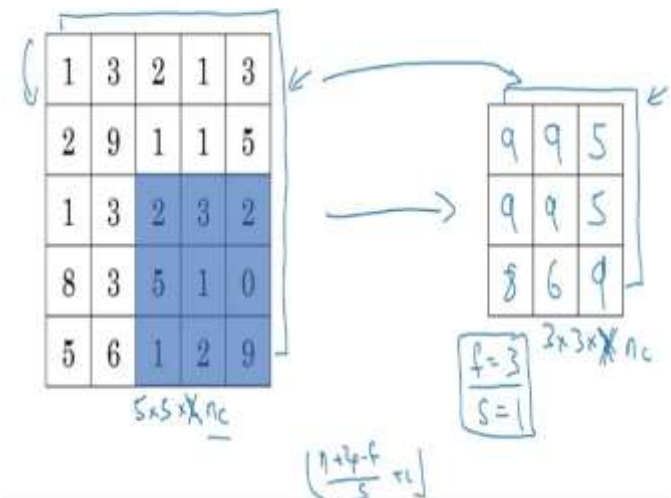
$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$
  

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

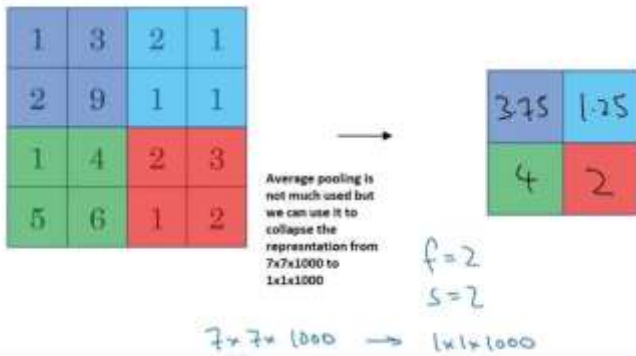
## Example ConvNet



## Pooling layer: Max pooling



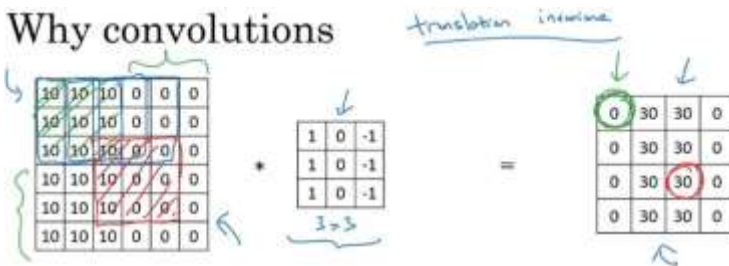
## Pooling layer: Average pooling



## Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	— 3,072 $\alpha^{(0)}$	0
CONV1 ( $f=5, s=1$ )	(28,28,6)	4,704	456 ←
POOL1	(14,14,6)	1,176	0 ←
CONV2 ( $f=5, s=1$ )	(10,10,16)	1,600	2,416 ←
POOL2	(5,5,16)	400	0 ←
FC3	(120,1)	120	48,120
FC4	(84,1)	84	10,164
Softmax	(10,1)	10	850

## Why convolutions

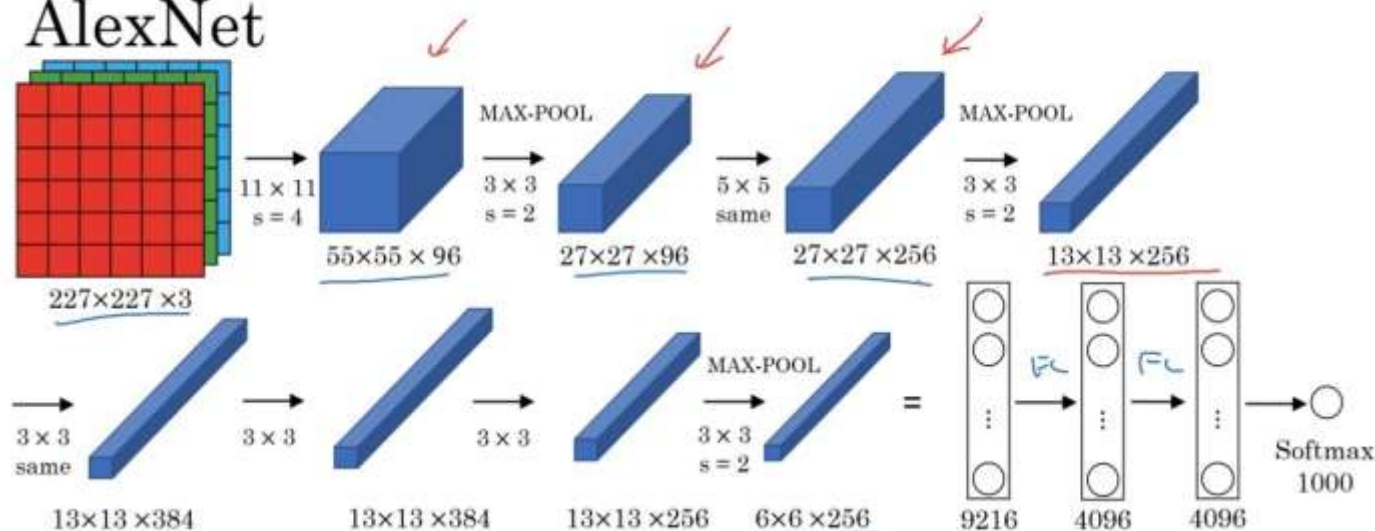


**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

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



# AlexNet



- Similar to LeNet, but much bigger.

- ReLU

- Multiple GPUs.

- Local Response Normalization (LRN)



Local response normalization picks a pixel and apply normalization across all the channels, but it does not affect that much.

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

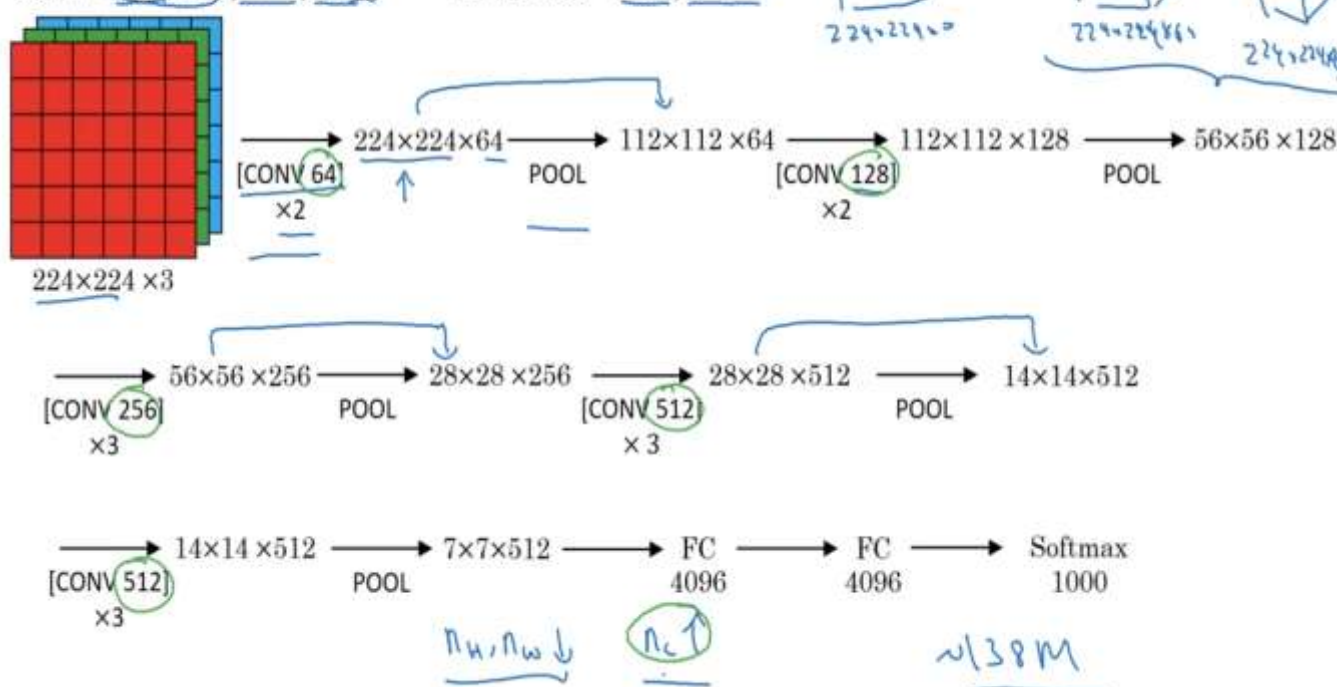
Andrew Ng

## 115 - VGG-16

# VGG - 16

CONV = 3x3 filter, s=1, same

MAX-POOL = 2x2, s=2

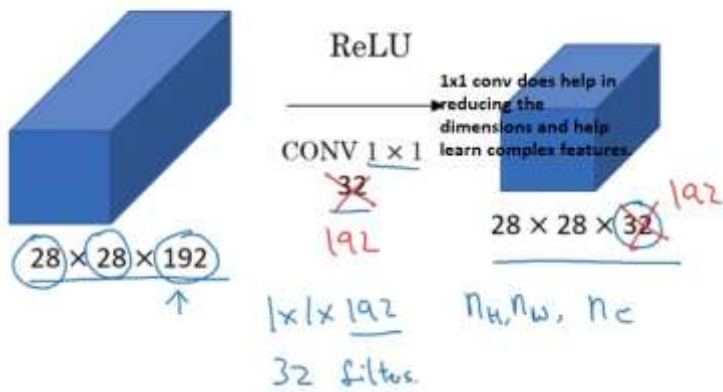


[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]

Andrew Ng

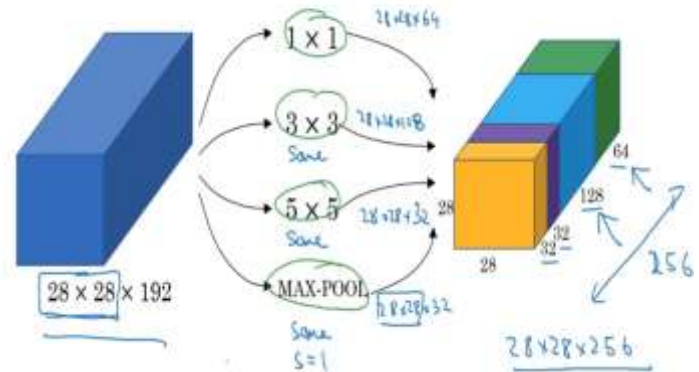
119

## Using 1x1 convolutions



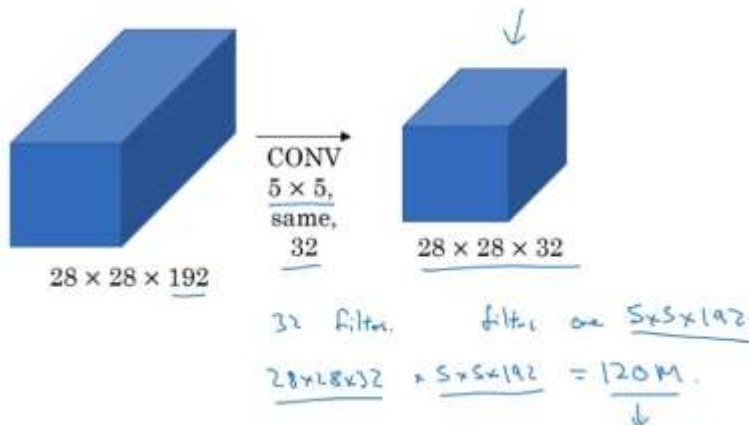
121.1

## Motivation for inception network



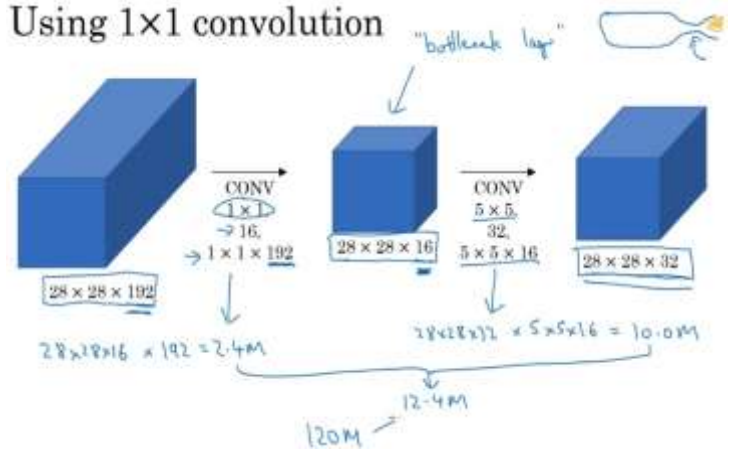
121.2

## The problem of computational cost



121.3

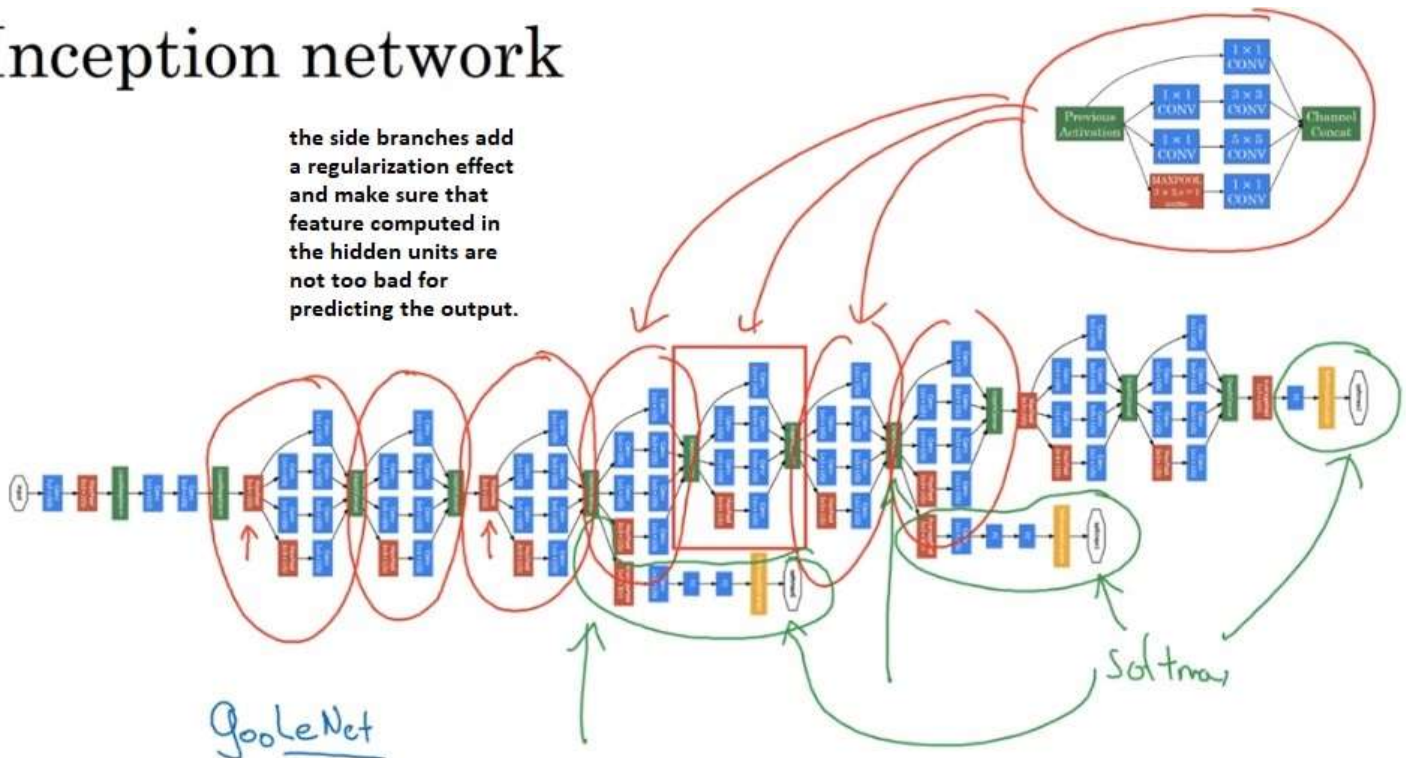
## Using 1x1 convolution



122.1

## Inception network

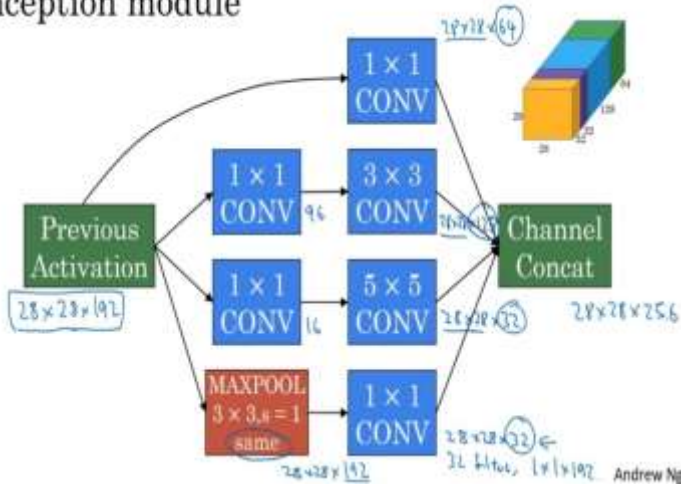
the side branches add a regularization effect and make sure that feature computed in the hidden units are not too bad for predicting the output.





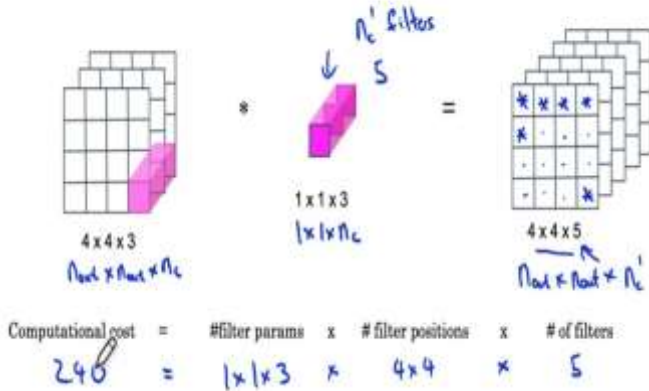
# Inception module

122

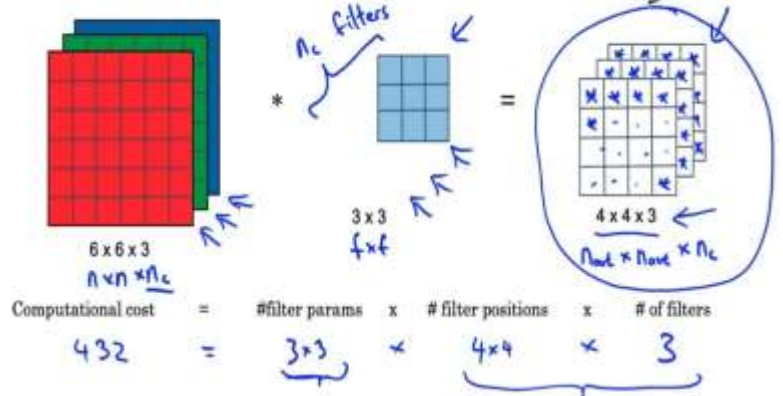


123.2

## Pointwise Convolution



## Depthwise Convolution



123.3

## Cost Summary

Cost of normal convolution  $2160$

Cost of depthwise separable convolution

$$\text{depthwise} + \text{pointwise} \\ 432 + 240 = 672$$

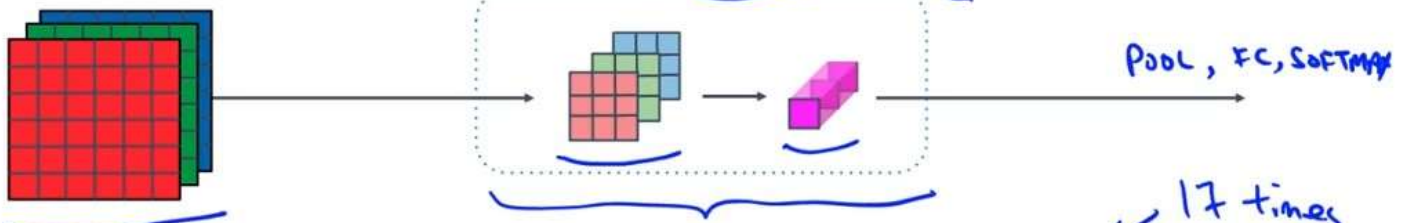
$$\frac{672}{2160} = 0.31 \leftarrow$$

$$= \frac{1}{N_c} + \frac{1}{f^2} \\ \frac{1}{5} + \frac{1}{9} \\ = \frac{1}{512} + \frac{1}{3^2 \times 8} \\ \sim 10 \text{ times cheaper}$$

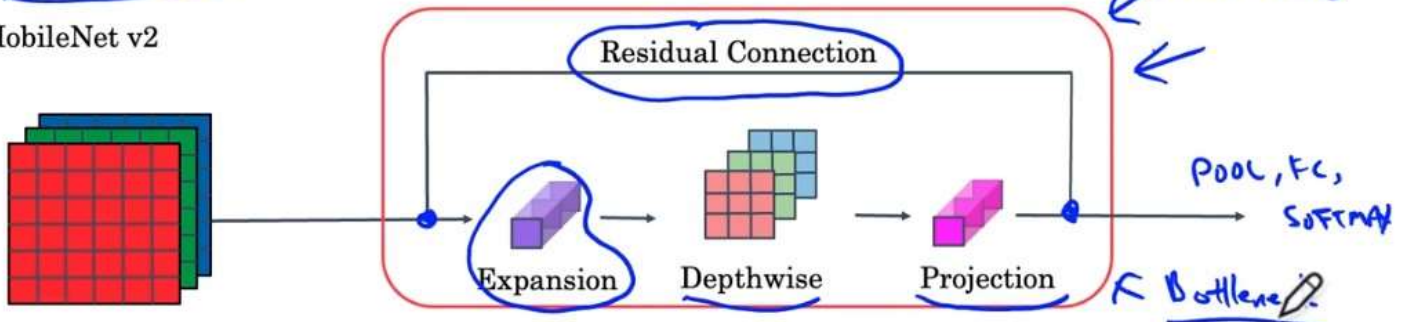
124.1

# MobileNet

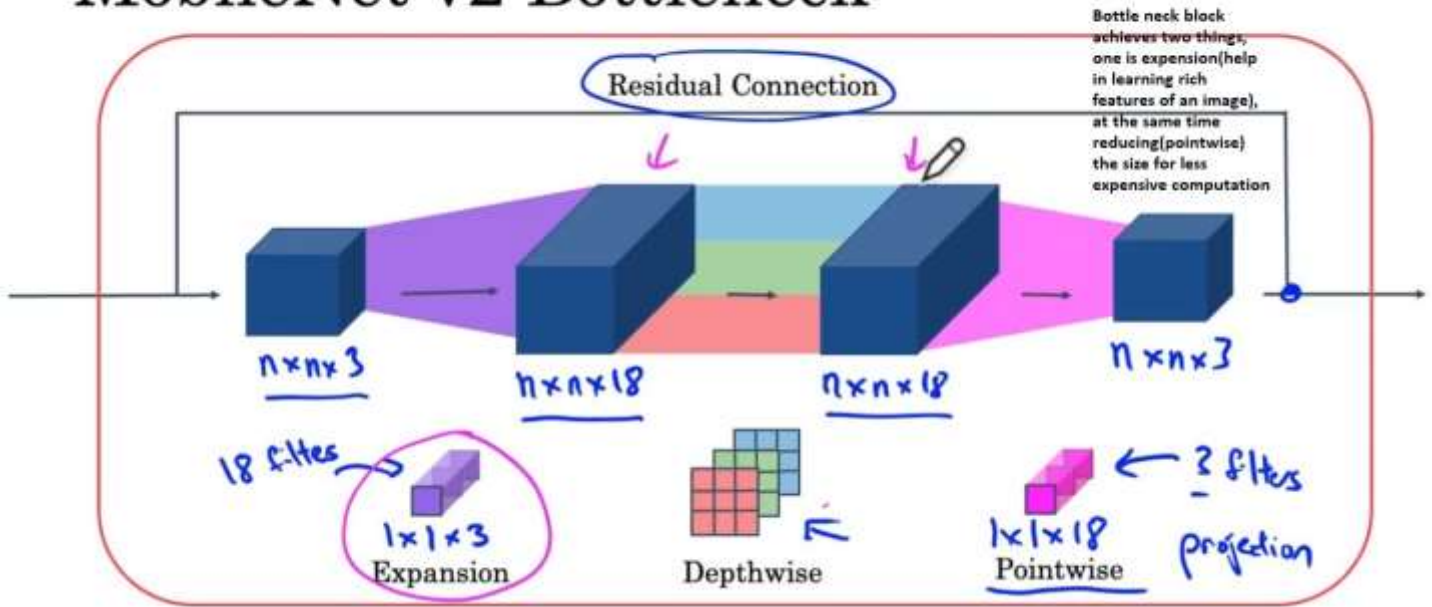
MobileNet v1



MobileNet v2

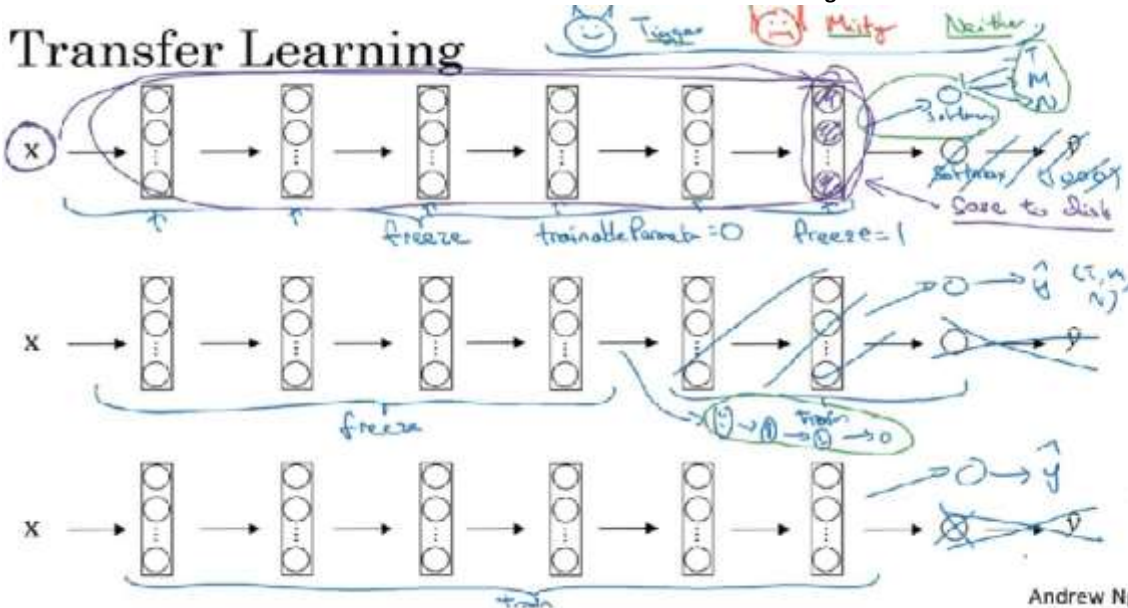


# MobileNet v2 Bottleneck



126 - transfer learning

## Transfer Learning



if u hv small dataset then download trained model from github(or imagenet) freeze, weighted layers and just remove output layer and apply softmax with your new output

if u hv a bit large dataset then download trained model from github(or imagenet) freeze initial, weighted layers and train later layers and just remove output layer and apply softmax with your new output

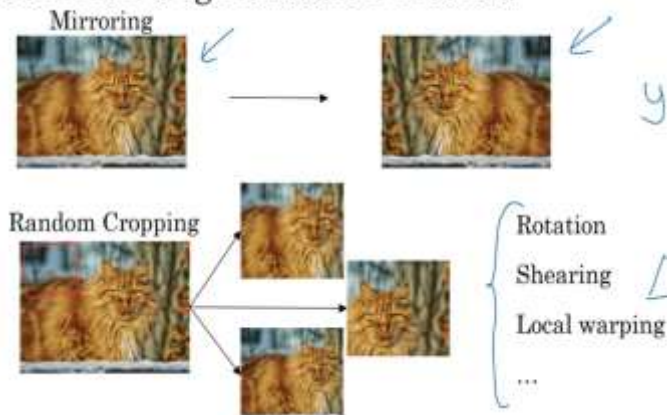
if u hv a very large dataset then download trained model from github(or imagenet) use all the weights as initialization and just remove output layer and apply softmax with your new output

Andrew Ng

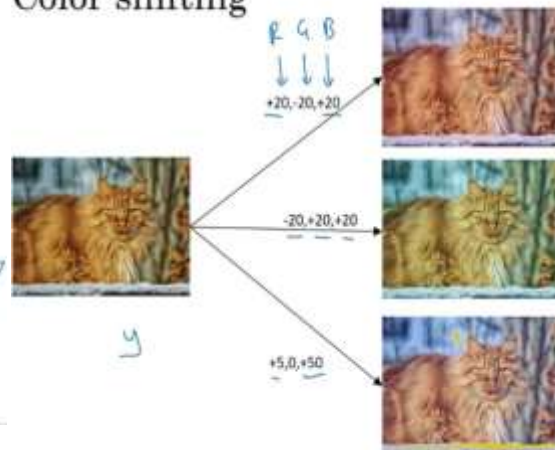
127.1

127.2

## Common augmentation method



## Color shifting



Advanced:

PCA

n2-class-org

AlexNet paper

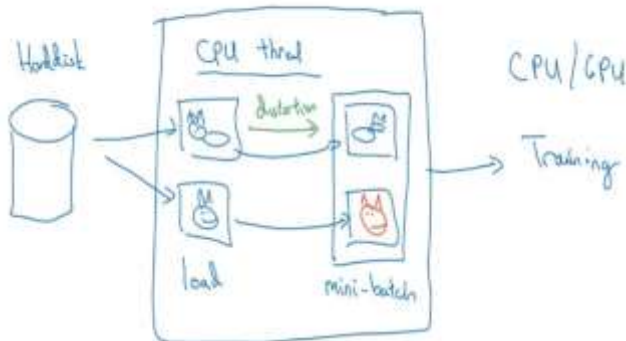
"PCA color augmentation"

R B G

R B G



## Implementing distortions during training



## Landmark detection



$b_x, b_y, b_h, b_w$

we can label the points in order to identify the landmarks, whether the landmarks belong to face or legs or arms etc



Facial Landmark Detection is a computer vision task that involves detecting and localizing specific points or landmarks on a face, such as the eyes, nose, mouth, and chin. The goal is to accurately identify these landmarks in images or videos of faces in real-time and use them for various applications, such as face recognition, facial expression analysis, and head pose estimation.

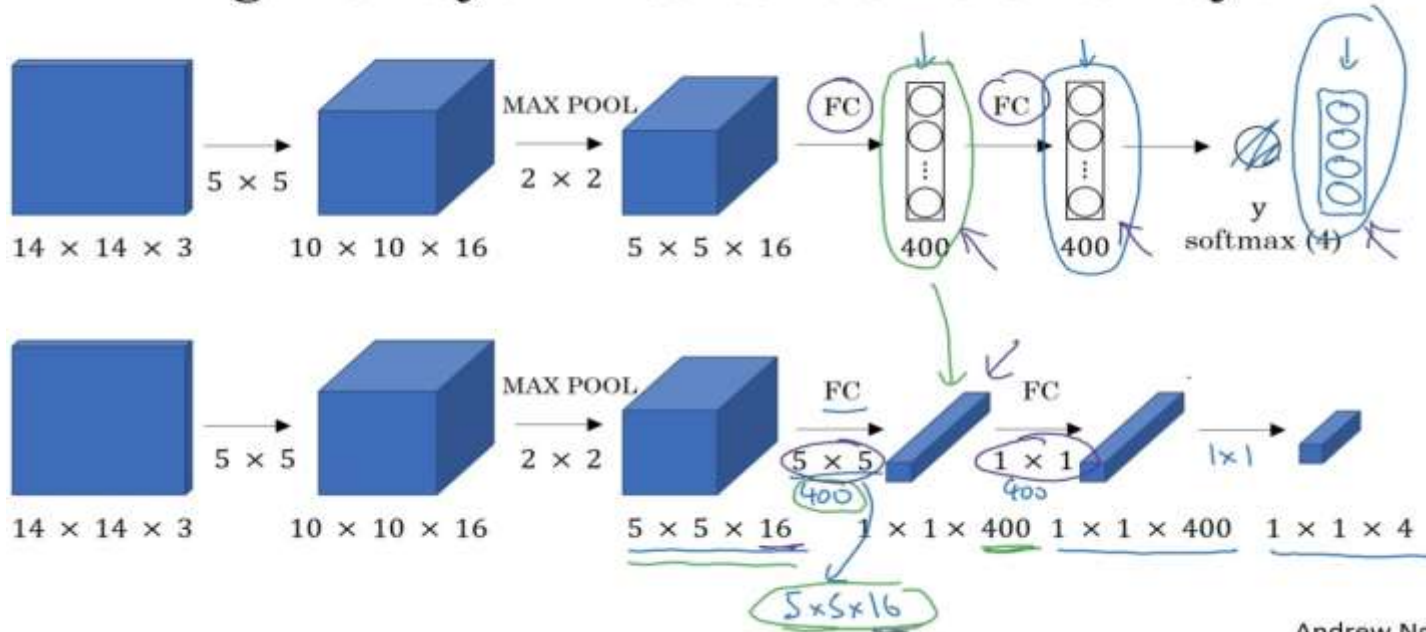


$l_{1x}, l_{1y},$   
 $l_{2x}, l_{2y},$   
 $l_{3x}, l_{3y},$   
 $l_{4x}, l_{4y},$   
 $\vdots$   
 $l_{64}, l_{64y}$

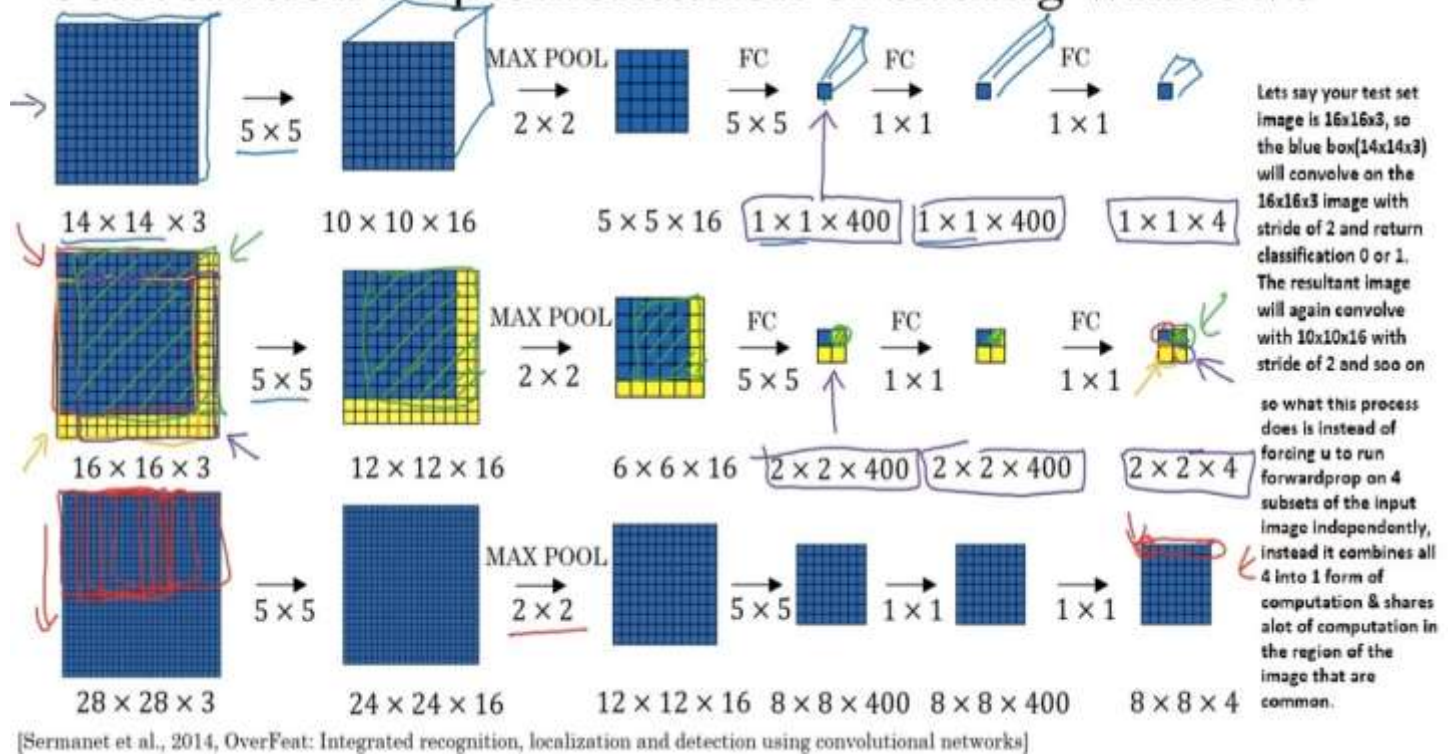
$x, y$

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

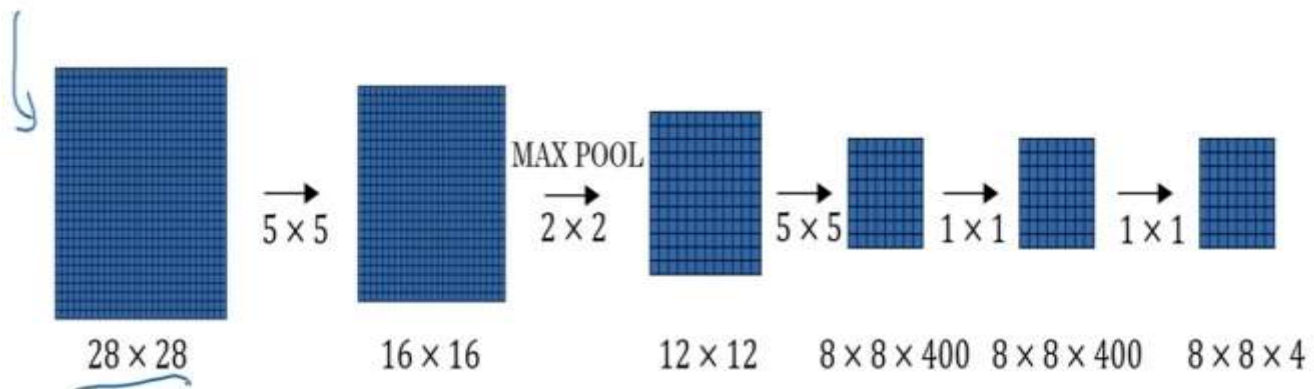
## Turning FC layer into convolutional layers



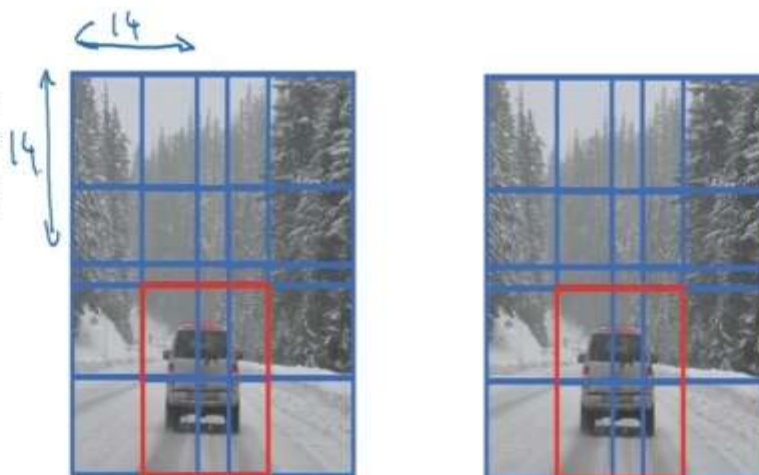
# Convolution implementation of sliding windows



# Convolution implementation of sliding windows



Previously what we did was take, let's say a kernel of  $14 \times 4$  and slide this onto the picture one by one, and that makes the computation expensive and we observed that the majority of the computation caused by moving kernel was duplicative.

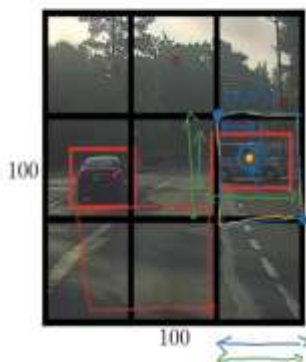


Now with the help of convNet on 132.2 we can directly pass the image in the network and the picture gets detected.



## Specify the bounding boxes

133.2



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

0.4 } between 0 and 1  
0.3 }  
0.9 } could be > 1  
0.5 }

Here's how YOLO works in more detail: 133.3

**Input image:** YOLO takes an input image of any size and divides it into a grid of cells.

**Anchor boxes:** Each grid cell predicts a fixed number of anchor boxes, which are predefined boxes of different shapes and sizes. Each anchor box is represented by its center coordinates, width, height, and a confidence score, which represents the probability that the box contains an object.

**Class probabilities:** For each anchor box, YOLO predicts class probabilities for each object class that the model has been trained to detect. These probabilities represent the probability that the object inside the box belongs to a particular class.

**Non-max suppression:** After prediction, YOLO applies non-maximum suppression (NMS) to remove overlapping bounding boxes and keep only the most confident predictions.

**Output:** The final output of YOLO is a set of bounding boxes and class probabilities for each object detected in the input image.

## Evaluating object localization

134.1



Intersection over Union (IoU)

$$= \frac{\text{Size of } \text{Intersection}}{\text{Size of } \text{Union}}$$

IOU gives accurate bounding box, if IOU is 1 then the bounding box is 100% accurate.

"Correct" if  $\text{IoU} \geq 0.5$

0.6

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

## Non-max suppression algorithm

134.2



19x19

Each output prediction is:

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

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

## Overlapping objects:

135.1

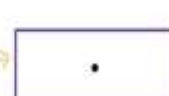


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



Anchor box 2:



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

Anchor box 1

Anchor box 2

Andrew Ng

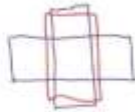
Redmon et al., 2015, You Only Look Once: Unified real-time object detection

## 135.2 Anchor box algorithm

Previously:

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

Output  $y$ :  
 $3 \times 3 \times 9$



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.

Output 4  
 $3 \times 3 \times 1$   
 $3 \times 3 \times 1$

(good cell, only  
box)

136.1

### Anchor box example



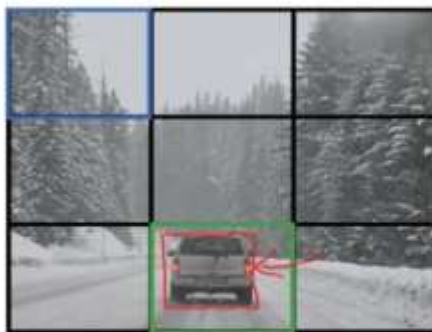
Anchor box 1:    Anchor box 2:



**135.3**

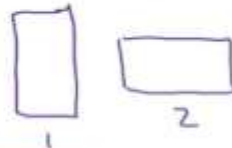
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_r \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\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

## Training



- 1 - pedestrian  
2 - car ←  
3 - motorcycle

$y =$



y is  $3 \times 3 \times 2 \times 8$

$$19 \times 19 \times 16$$
$$19 \times 19 \times 40$$

anchors

↖  $S + \# \text{ classes}$



→ Comu Ne-



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

Andrew Ng

## Making predictions



this  $3 \times 3 \times 2 \times 3$  will be calculated for each grid and we will get this  $\longrightarrow$  matrix against each grid.

 $3 \times 3 \times 2 \times 8$ 

$y =$



# Outputting the non-max suppressed outputs

136.3



- For each grid cell, get 2 predicted bounding boxes.

Each grid will have 2 anchor boxes of different sizes.



- Get rid of low probability predictions.



- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

Here's how RCNN works in more detail:

136.4

**Region proposal:** The first step in RCNN is to generate a set of region proposals using a separate algorithm, such as Selective Search. Selective Search combines low-level and high-level image features to identify regions that are likely to contain objects.

**Feature extraction:** For each proposed region, RCNN uses a deep convolutional neural network (CNN) to extract a fixed-length feature vector. This feature vector captures the salient features of the proposed region.

**Classification:** The feature vector is then fed into a set of fully connected layers, which are trained to classify the object in the proposed region into one of several classes. These classes typically correspond to different object categories, such as "car", "person", "dog", etc.

**Bounding box regression:** In addition to classifying the object, RCNN also predicts the object's bounding box coordinates (i.e., the object's location and size within the proposed region) using a separate set of fully connected layers.

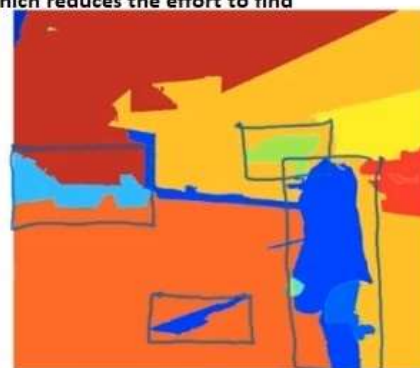
**Non-maximum suppression:** After classification and bounding box regression, RCNN applies non-maximum suppression (NMS) to remove overlapping bounding boxes and keep only the most confident predictions.

**Output:** The final output of RCNN is a set of bounding boxes and class probabilities for each object detected in the input image.

136.5

## Region proposal: R-CNN

Convolving an image is an expensive task, so what R-CNN does is it applies the segmentation algorithm and creates multiple segments of an image, which makes it easier to identify objects since it makes multiple blobs which reduces the effort to find object



Segmentation algorithm

~2,000

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

## Faster algorithms

136.6

137.1

→ R-CNN:

Propose regions. Classify proposed regions one at a time. Output label + bounding box.

Proposal generation: Like other region-based methods, Fast R-CNN relies on a separate algorithm to generate region proposals, such as Selective Search, which can be time-consuming.

Fast R-CNN:

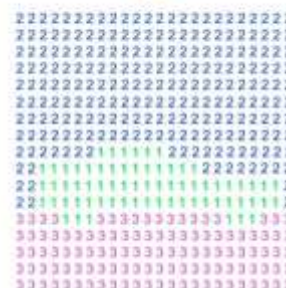
Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions.

Faster R-CNN: Use convolutional network to propose regions.

## Per-pixel class labels



- Car
- Building
- Road



Segmentation Map

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

[Girshik, 2015, Fast R-CNN]

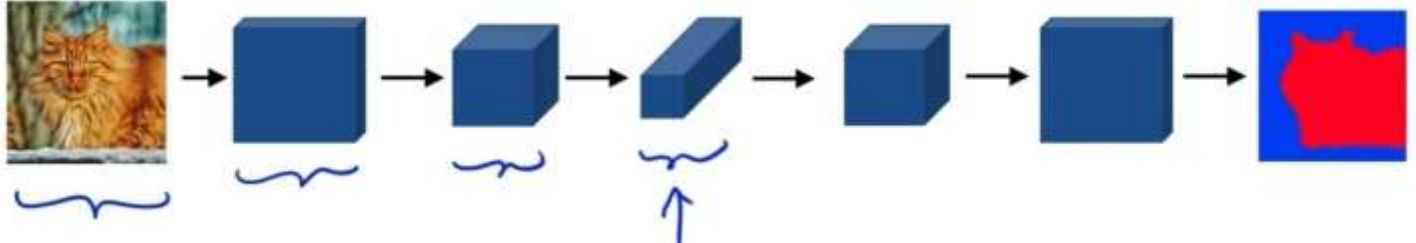
[Ren et. al, 2016, Faster R-CNN: Towards real-time object detection with region proposal networks]

Andrew Ng

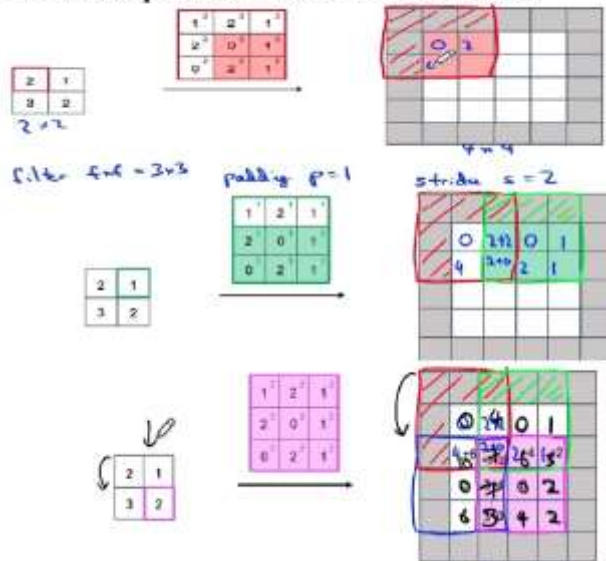
# Deep Learning for Semantic Segmentation

So when we pass the image through N.N the normally the image dimensions gets smaller but here the dimensions get bigger (means as we go deeper into the N.N the height and width will go back and # of channels will decrease) so they can gradually blow it back-up to a full-size image, which is a size we want for the output.

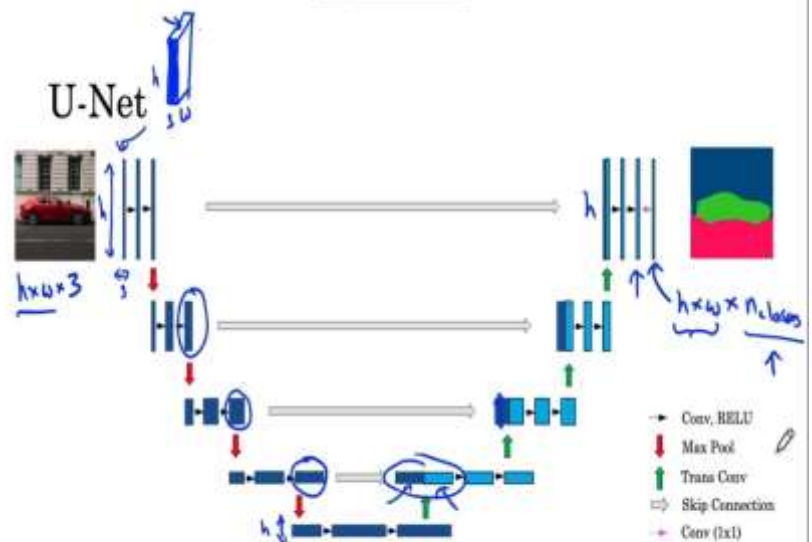
All this is done using transpose convolution



## Transpose Convolution

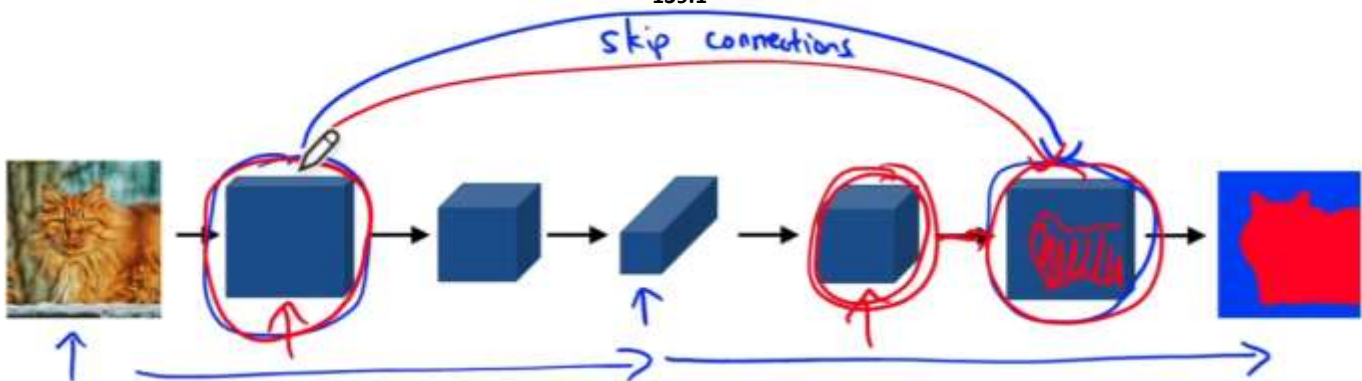


## U-Net



[Ronneberger et al., 2015, U-Net: Convolutional Networks for Biomedical Image Segmentation]

Andrew Ng



By reducing the channels we again make picture bigger using transpose convolution and by skipping connection we are making sure that N.N knows where in the original image the cat located.



144

## 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.

$$\|C(A)-C(P)\|^2 + \alpha \leq \|C(A)-C(N)\|^2$$

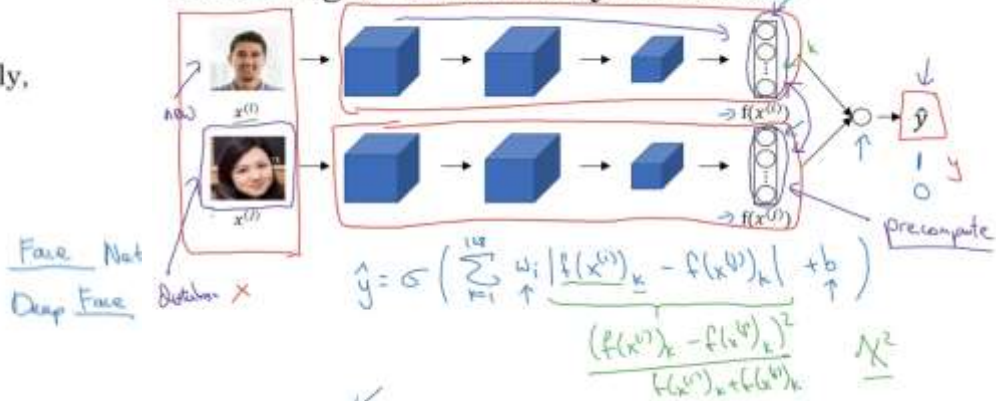
Choose triplets that're "hard" to train on.

$$\frac{d(A,P)}{d(A,N)} + \alpha \leq \frac{d(A,N)}{d(A,N)}$$

Face Not  
Deep Face

## Learning the similarity function

145.1



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

Andre

145.2

## Face verification supervised learning

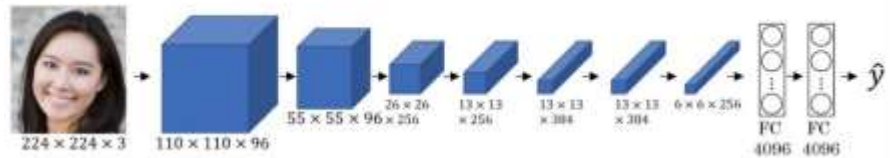


"Same"

"Different"

## Visualizing what a deep network is learning

146.1



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

Repeat for other units.



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

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

Andrew Ng

146.2

## Visualizing deep layers



Layer 1



Layer 2



Layer 3



Layer 4



Layer 5

## Content cost function

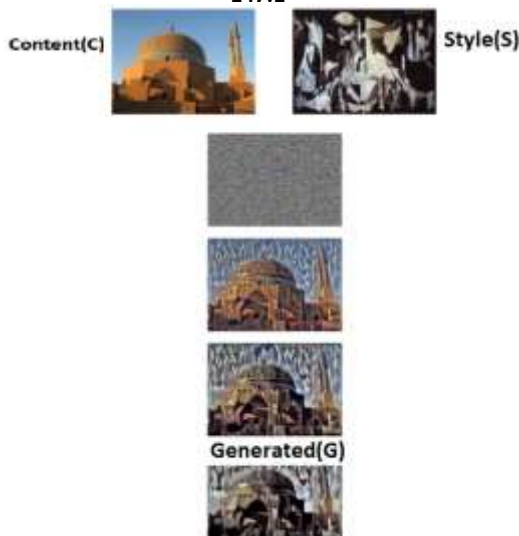
147.2

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

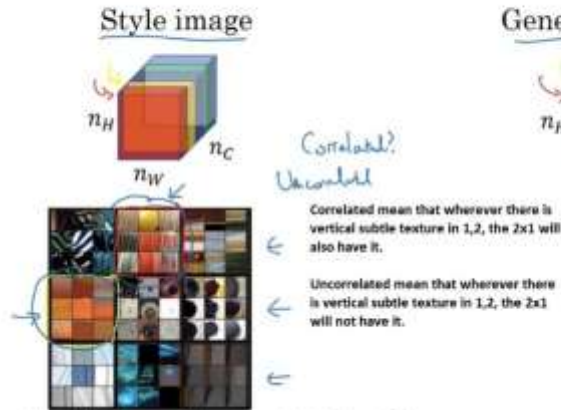
- Say you use hidden layer  $l$  to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let  $a^{[l](C)}$  and  $a^{[l](G)}$  be the activation of layer  $l$  on the images
- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar, both images have similar content

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

If  $J(C, G)$  is greater, this represents that the generated image will be more similar to content image, and if  $J(S, G)$  is greater then this represents that the generated image will be more similar to Style image, the hyperparameters are there to incentivize these  $J(C, G)$  and  $J(S, G)$ .



# Intuition about style of an image



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

## Style cost function

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

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

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

## Style matrix

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

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

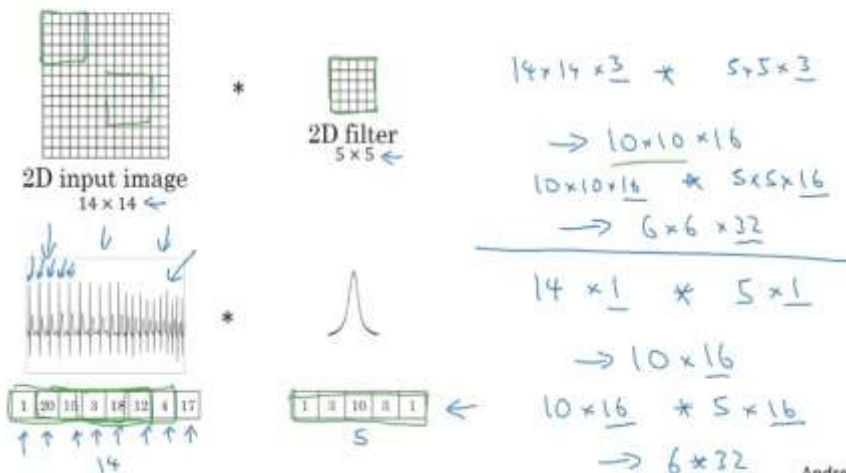
this is done for both style(S) and generated(G) image.

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

Gram matrix  
Calculates how correlated the activations are in channel k and k' then compare them.

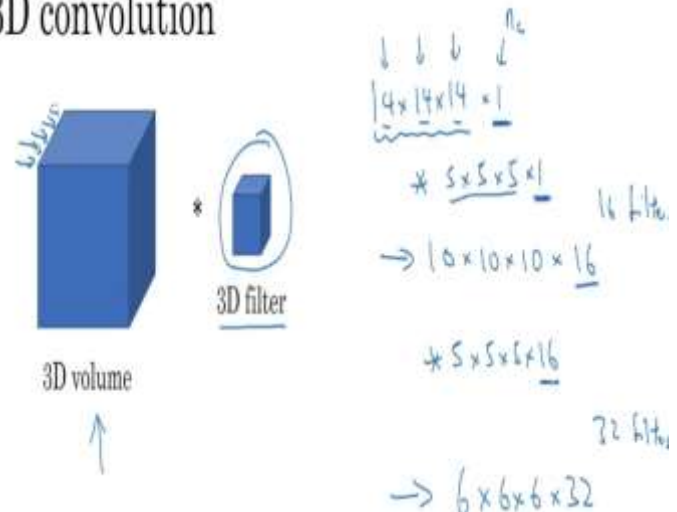
$$J_{style}^{[l]}(S, G) = \frac{1}{(2n_H^{[l]}n_W^{[l]}n_C^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l]}(S) - G_{kk'}^{[l]}(G))^2$$

## Convolutions in 2D and 1D



Andres

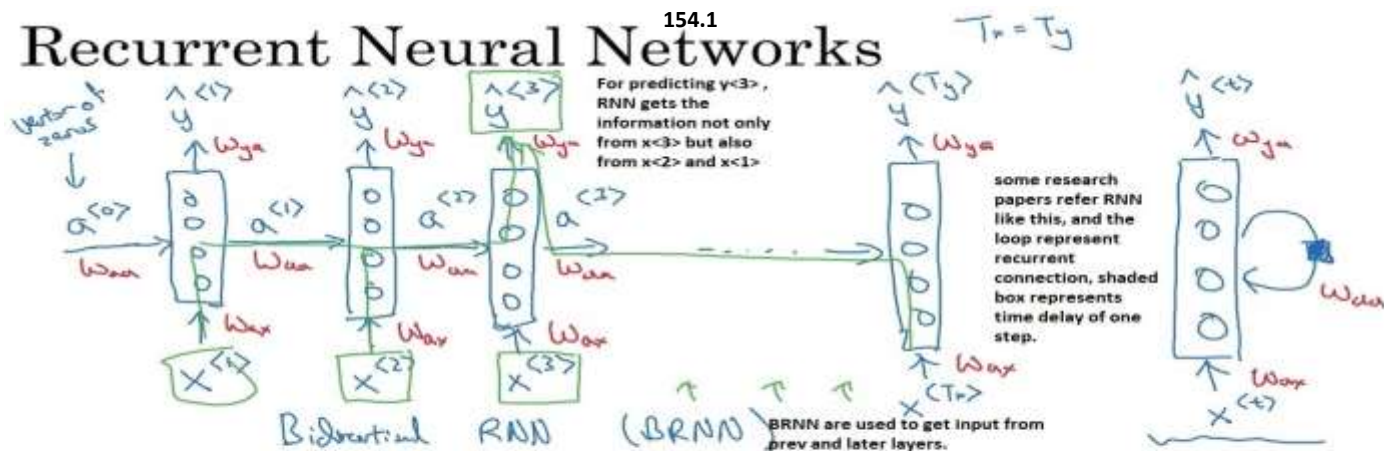
## 3D convolution





# Recurrent Neural Networks

154.1



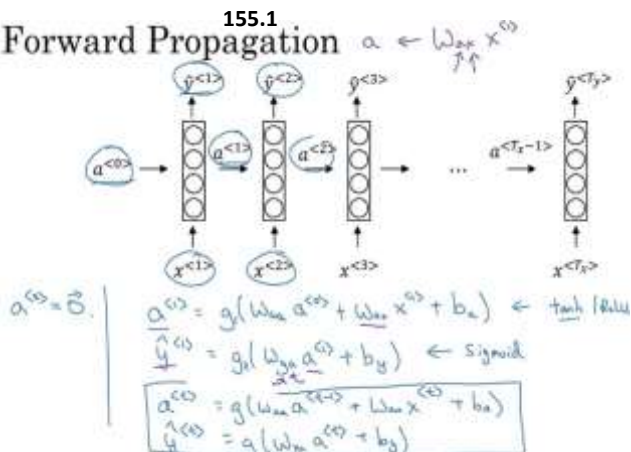
Downside is that RNN only uses the info that is in the earlier in the seq to predict  $y^{(3)}$ , it does not use info from  $x^{(4)}$ ,  $x^{(5)}$  and so on. Below we can see given the first three words it not possible to predict whether teddy is part of a person's name. In the first eg it is, in the 2nd it is not.

He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

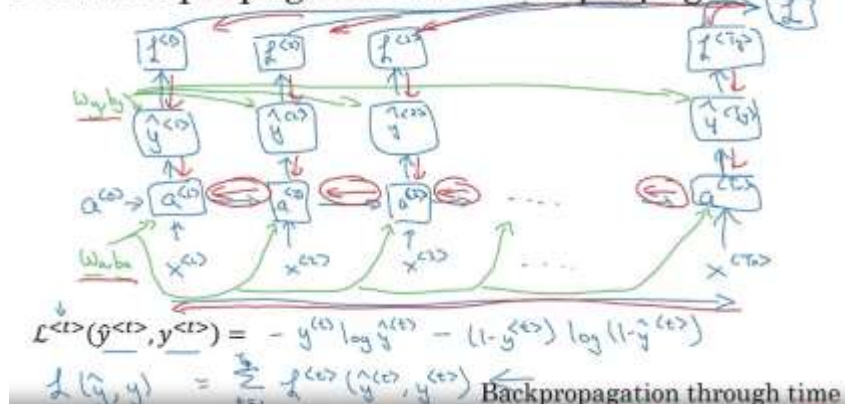
## Forward Propagation

155.1



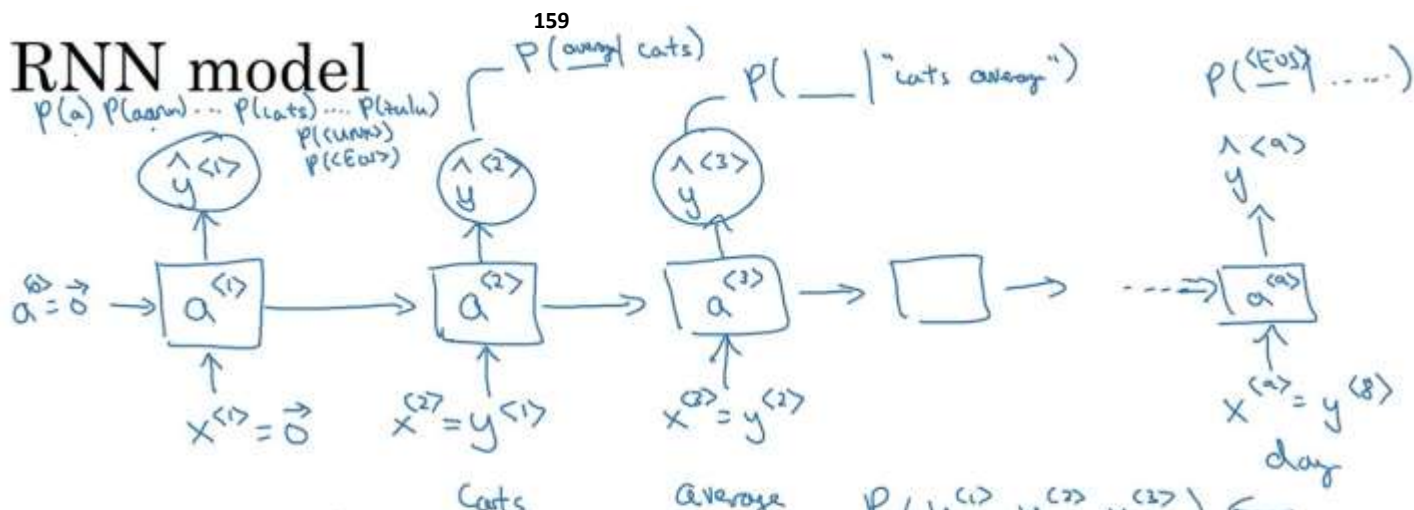
## Forward propagation and backpropagation

155.3



## RNN model

159



Cats average 15 hours of sleep a day. <EOS>

$$\mathcal{L}(\hat{y}^{(t)}, y^{(t)}) = -\sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

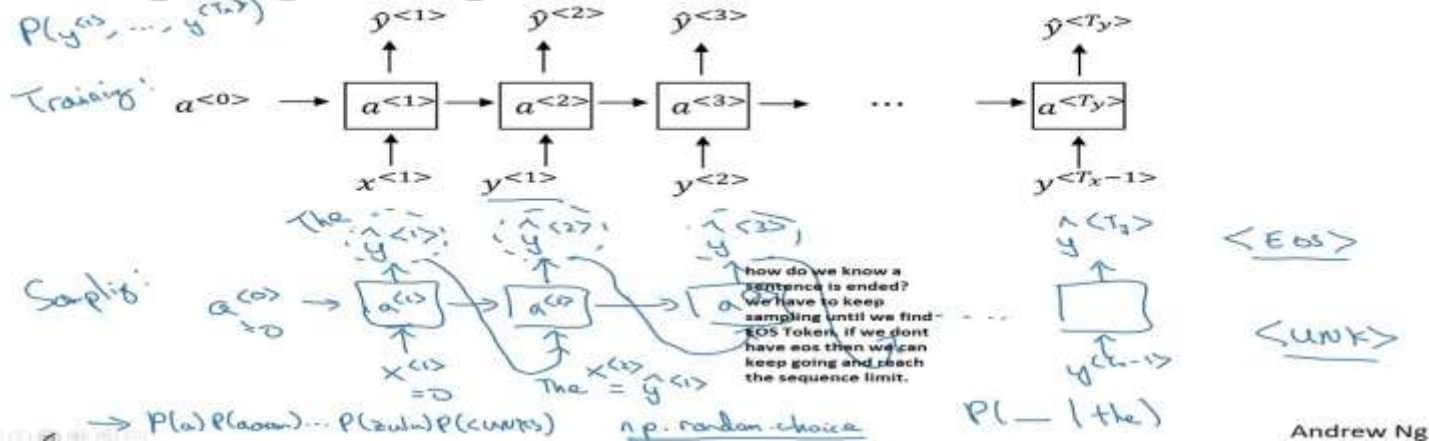
$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$P(y^{(1)}, y^{(2)}, y^{(3)}) \leftarrow$$

$$= \frac{P(y^{(1)}) P(y^{(2)} | y^{(1)})}{P(y^{(3)} | y^{(1)}, y^{(2)})}$$

# Sampling a sequence from a trained RNN

160.1



160.2

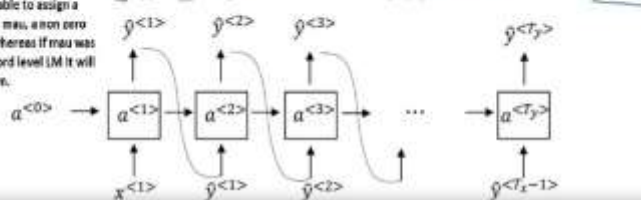
## Character-level language model

Disadvantage of character level LM is we end up with much more longer sequences.

→ Vocabulary = [a, aaron, ..., zulu, <UNK>]

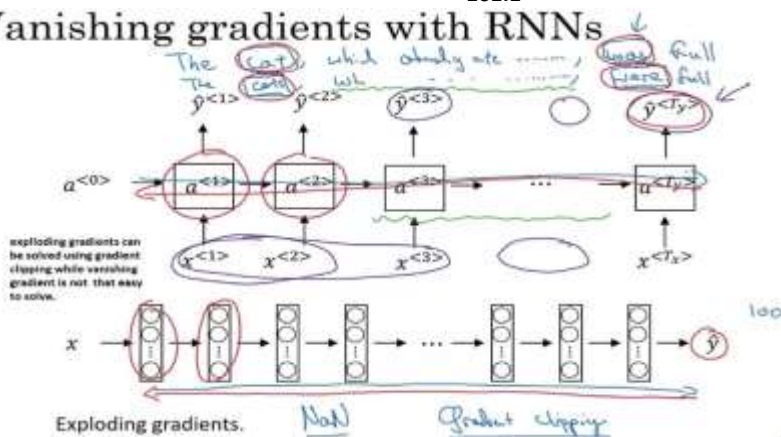
→ Vocabulary = [a, b, c, ..., z, ., ,, ;, : , , - , \_ , ' , " , <UNK>]

In this model we don't need to worry about unknown words, this model is able to assign a sequence like mau, a non zero probability, whereas if mau was not in your word level LM it will be an unknown.



161.1

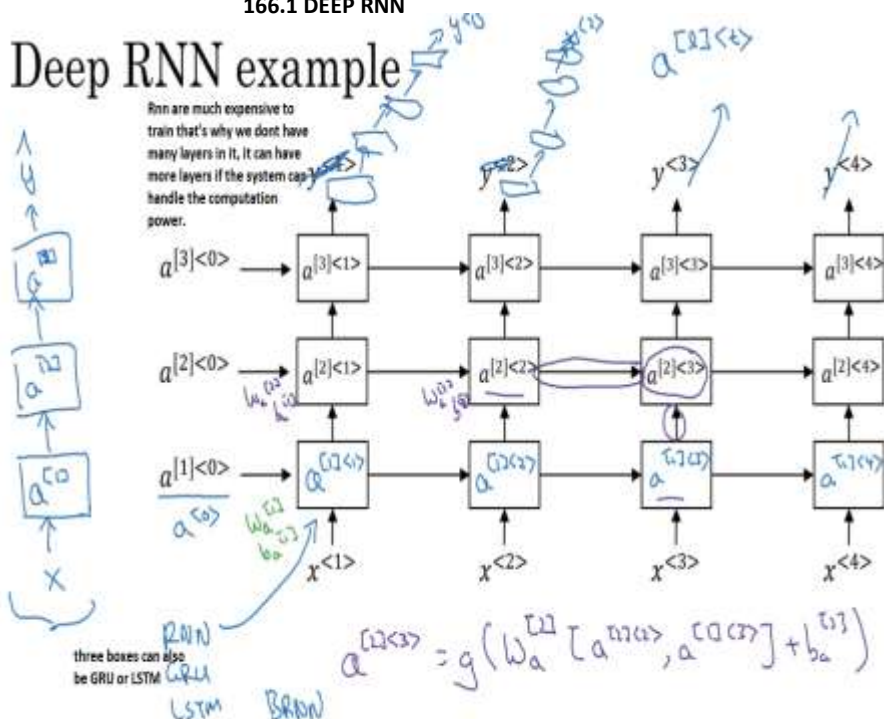
## Vanishing gradients with RNNs



166.1 DEEP RNN

## Deep RNN example

Rnn are much expensive to train that's why we don't have many layers in it, it can have more layers if the system can handle the computation power.



166BRNN

Andrew Ng

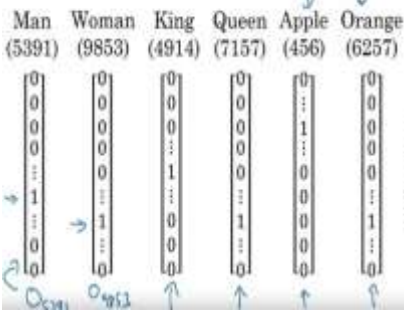


167.1

## Word representation

$V = [a, aaron, ..., zulu, <UNK>]$

1-hot representation



I want a glass of orange juice.  
I want a glass of apple juice.

A model doesn't know that apple and orange juice are similar, because if you take the product of any two vectors, it will be 0. So it does not know that somehow apple and orange are similar category.

167.2

## Featurized representation: word embedding

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

In this way we assign different features to different words based on numbers, upon processing this whole vectorized feature representation, the model can see that in the food category orange and apple are correlated.

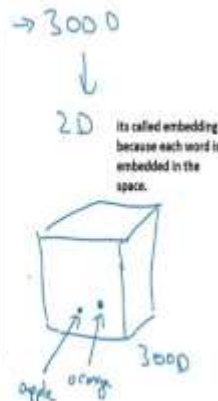
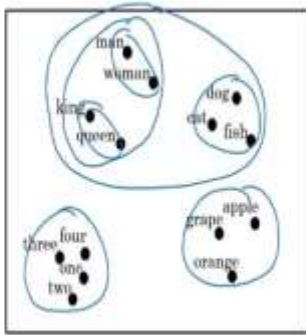
I want a glass of orange juice.  
I want a glass of apple juice.

Andrew Ng

## Visualizing word embeddings

167.3

The t-SNE algorithm calculates a similarity measure between pairs of instances in the high dimensional space and in the low dimensional space. It then tries to optimize these two similarity measures using a cost function.



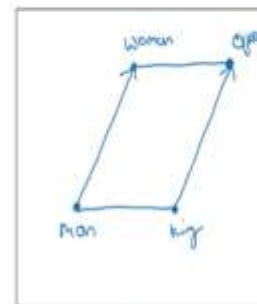
168.1

## Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)  
(Or download pre-trained embedding online.)
2. Transfer embedding to new task with smaller training set. (say, 100k words)  
→ 10,000 → 300
3. Optional: Continue to finetune the word embeddings with new data.

169.1

## Analogy using word vectors



similarity function is used in order to find the word for analogy, some research papers get the accuracy b/w 80 to 75%. Most commonly used similarity function is cosine similarity.

$$e_{man} - e_{woman} \approx e_{king} - e_{queen}$$

300D

$$\text{Find word w: } \arg \max_w \sin(e_w, e_{king} - e_{man} + e_{woman})$$

30-75%

Andrew Ng

168.2

## Analogy

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

we can use model learn embeddings using word embeddings by subtracting the respective vectors, just like we did with man - woman and king - queen. The goal is to balance the eqn, if the eqn gets balanced, it's a sign that we have to use that word.

$e_{man} - e_{woman} \approx e_{king} - e_{queen}$

Man → Woman    King → ? Queen

$e_{man} - e_{woman} \approx e_{king} - e_{queen}$

[Mikolov et. al., 2013. Linguistic regularities in continuous space word representations]

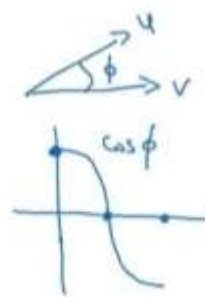
Andrew Ng

# Cosine similarity

169.2

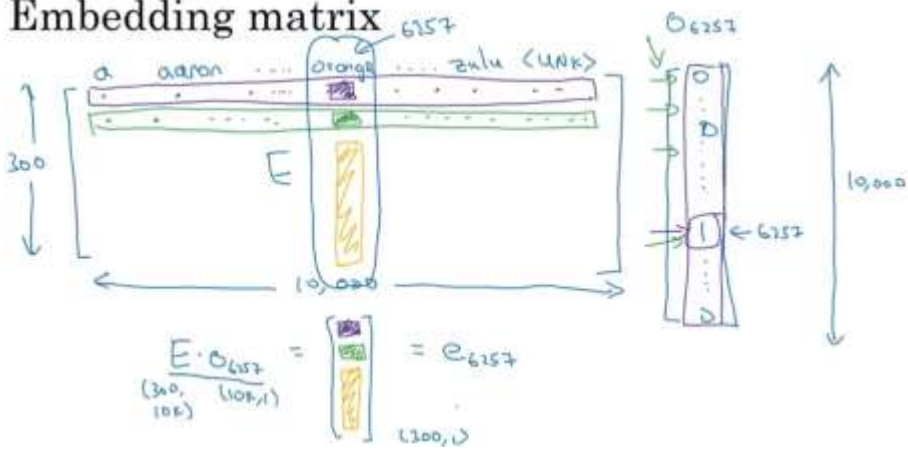
$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

it's called cosine because of the cosine angle btw two vectors u and v



$\|u - v\|^2$   
we can also use squared difference similarity, but it is more like dissimilarity but it works.

## Embedding matrix



170.2

## Other context/target pairs

I want a glass of orange juice to go along with my cereal.

Context: Last 4 words

4 words on left & right

Last 1 word

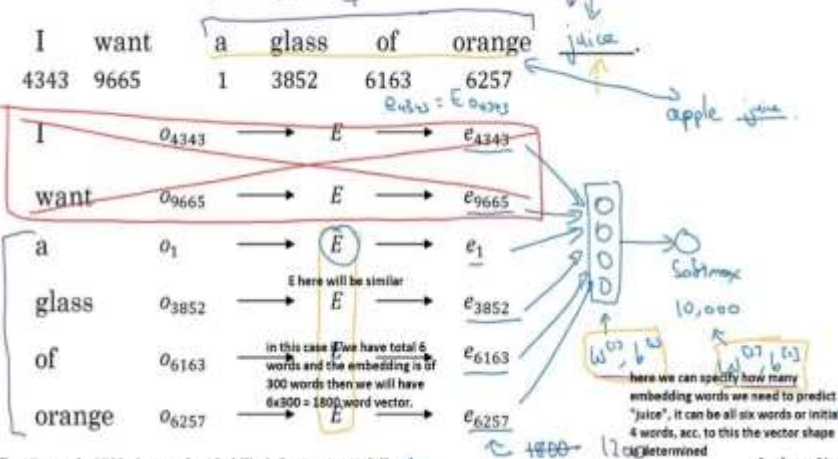
Nearby 1 word

We can specify orange to identify juice. We can identify one nearby word, such as glass, which refers to juice if we think about it.

a glass of orange juice to go along with my cereal.  
orange ?  
glass ?  
to go along with

## Neural language model

170.1



[Bengio et. al., 2009, A neural probabilistic language model]

Andrew Ng

## Model

178.1

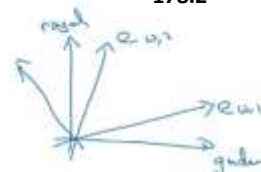
$$\text{Minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} \ell(x_{ij}) (\theta_i^T e_j + b_i + b'_j - \log x_{ij})^2$$

weight term  
 $\ell(x_{ij}) = 0$  if  $x_{ij} = 0$   
 $0 \log 0 = 0$   
this is, is, a, a, ...  
Deriv  
 $e_u^{(R)} = \frac{e_u + \theta_u}{2}$

## A note on the featurization view of word embeddings

178.2

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)
Gender	-1	1	-0.95	0.97
Royal	0.01	0.02	0.93	0.95
Age	0.03	0.02	0.70	0.69
Food	0.09	0.01	0.02	0.01



$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\theta_i^T e_j + b_i + b'_j - \log x_{ij})^2$$

Andrew Ng

## Sentiment classification problem

179.1

$x \rightarrow y$

The dessert is excellent. ★★★★★

Service was quite slow. ★★☆☆☆

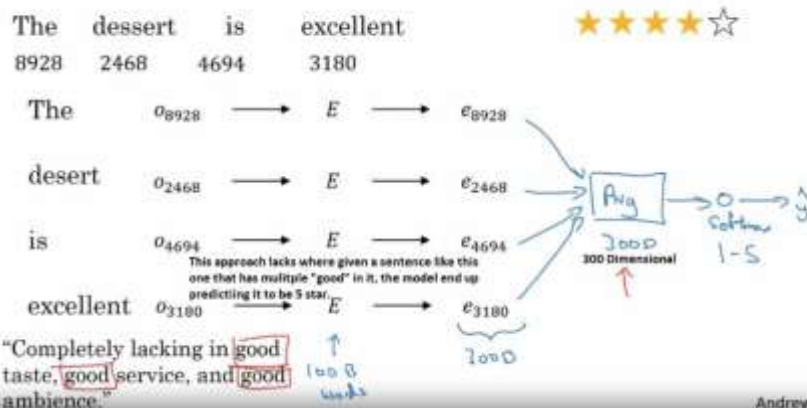
Good for a quick meal, but nothing special. ★★★☆☆

Completely lacking in good taste, good service, and good ambience. ★☆☆☆☆

10,000 → 100,000 words

## Simple sentiment classification model

179.2

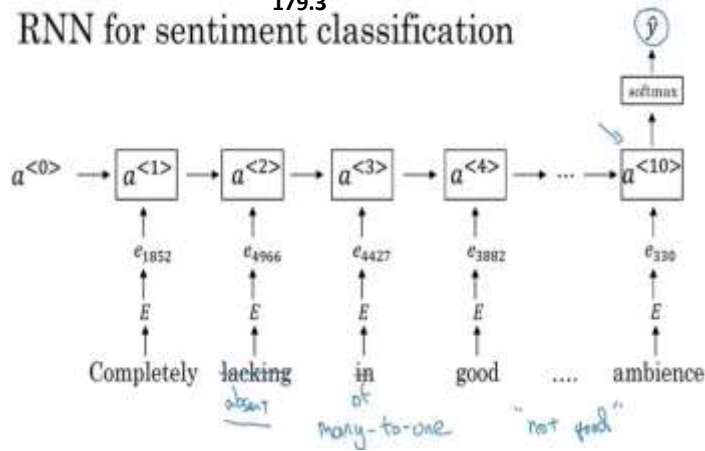


Andrew



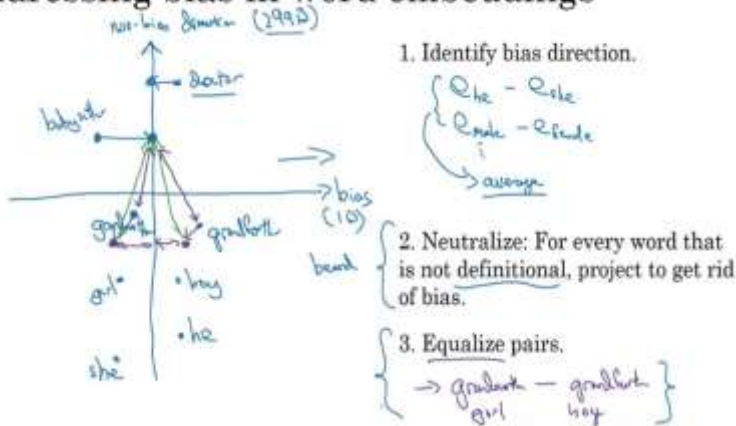
179.3

## RNN for sentiment classification



180.2

## Addressing bias in word embeddings



180.1

## The problem of bias in word embeddings

Man:Woman as King:Queen

Man:Computer\_Programmer as Woman:Homemaker ✗

Father:Doctor as Mother:Nurse ✗

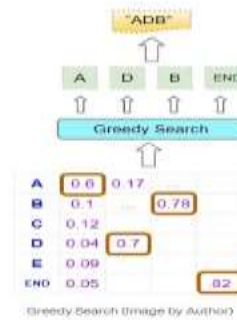
Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.

Dobnik et al., 2016, Man is to computer programmer as woman is to homemaker? Debiasing word embeddings Andrew Ng

## Greedy Search

184.1

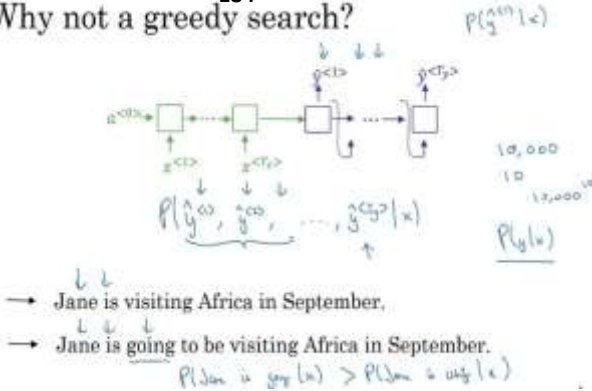
A fairly obvious way is to simply take the word that has the highest probability at each position and predict that. It is quick to compute and easy to understand, and often does produce the correct result.



In fact, Greedy Search is so easy to understand, that we don't need to spend more time explaining it. But can we do better?

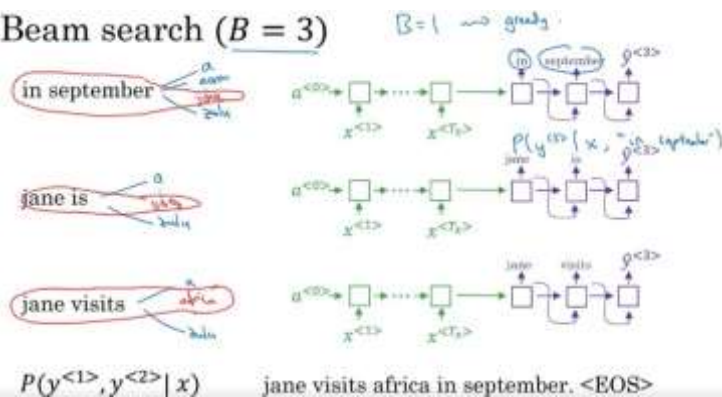
184

## Why not a greedy search?



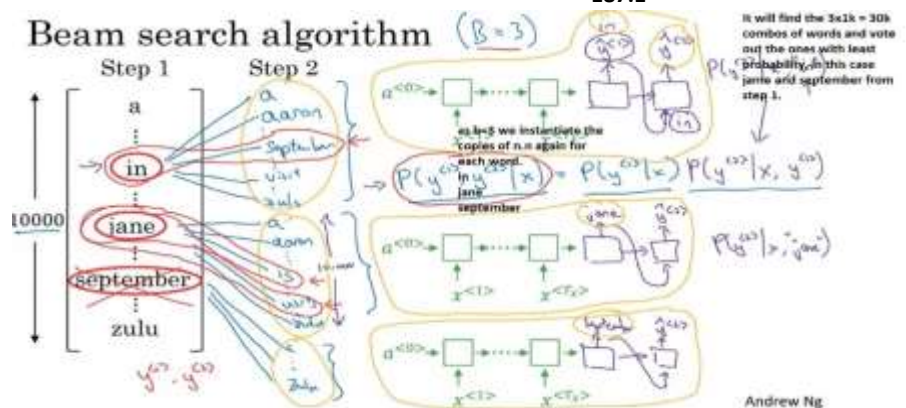
187.2

## Beam search (B = 3)



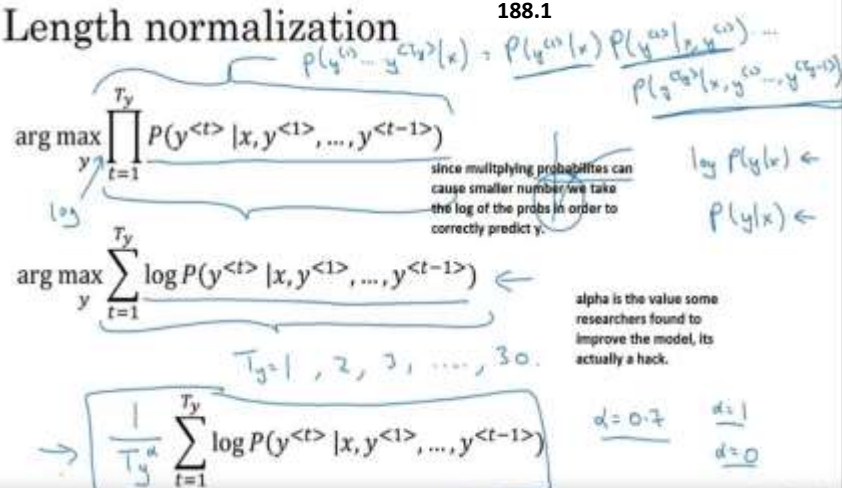
187.1

## Beam search algorithm (B=3)



188.1

## Length normalization



## Beam search discussion

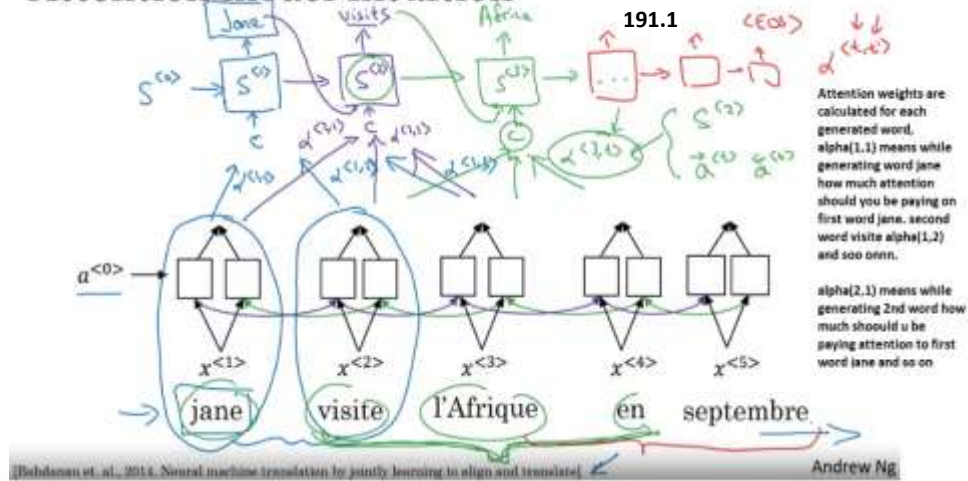
Beam width B?

|→| 10, 100, 1000, → 2000

Large B: better result, slower  
Small B: worse result, faster

Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but is not guaranteed to find exact maximum for  $\arg \max_y P(y|x)$ .

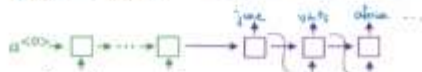
## Attention model intuition



188

## Example

Jane visite l'Afrique en septembre.

Human: Jane visits Africa in September. ( $y^*$ )Algorithm: Jane visited Africa last September. ( $\hat{y}$ )RNN computes  $P(\hat{y}|x) \leq P(y^*|x)$ 

## Error analysis on beam search

Human: Jane visits Africa in September. ( $y^*$ )Algorithm: Jane visited Africa last September. ( $\hat{y}$ )Case 1:  $P(y^*|x) > P(\hat{y}|x)$  ←  $\text{as } P(y^*|x) > P(\hat{y}|x)$ Beam search chose  $\hat{y}$ . But  $y^*$  attains higher  $P(y|x)$ .

Conclusion: Beam search is at fault.

Case 2:  $P(y^*|x) \leq P(\hat{y}|x)$  ← $y^*$  is a better translation than  $\hat{y}$ . But RNN predicted  $P(y^*|x) < P(\hat{y}|x)$ .

Conclusion: RNN model is at fault.

## Error analysis process

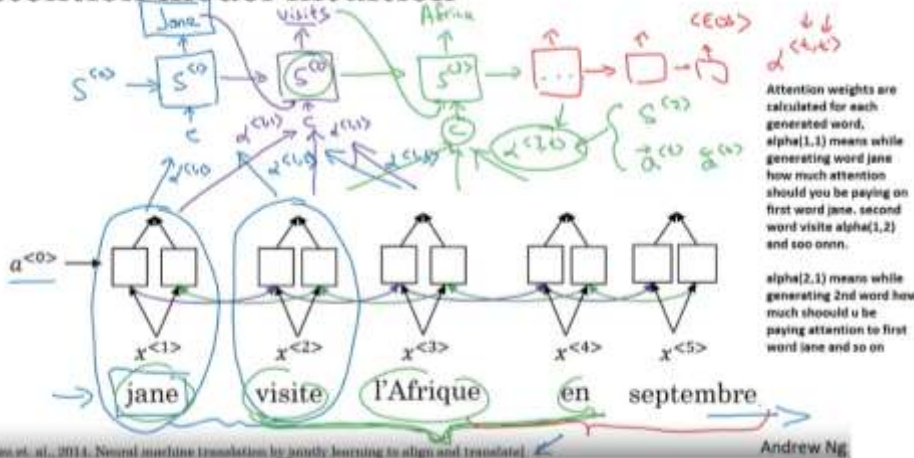
Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	$2 \times 10^{-10}$	$1 \times 10^{-10}$	(B)
...	...	...	...	(R)
				(R)
				(R)
				(R)

If most errors are by beam, then we can try increasing b value and tune other hyperparameters.

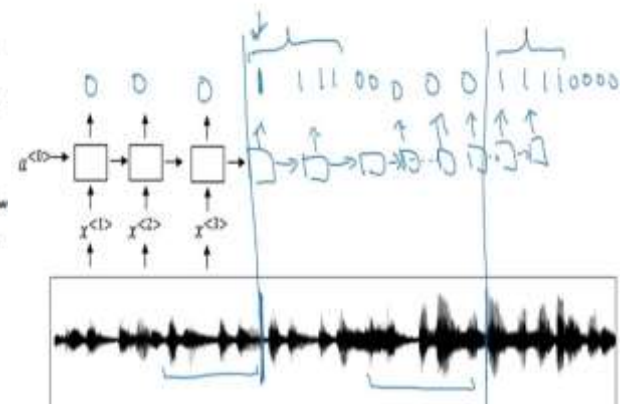
If RNN is making more errors then we can add regularization or get more training data and so on.

Figures out what fraction of errors are "due to" beam search vs. RNN model

## Attention model intuition

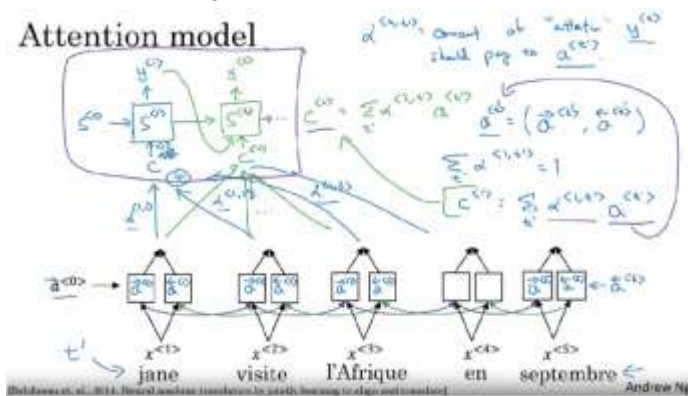
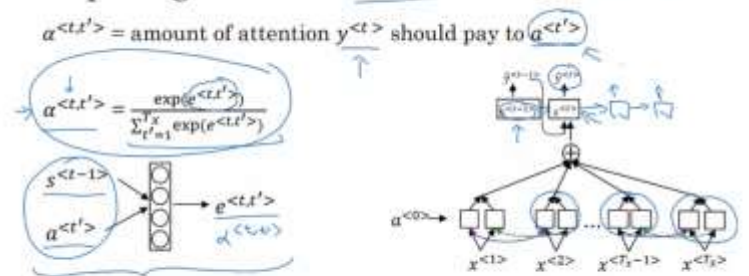


## Trigger word detection algorithm





## Attention model

Computing attention  $\alpha^{<t,t'>}$ 

Bahdanau et al., 2014. Neural machine translation by jointly learning to align and translate.  
 Xu et al., 2015. Show, attend and tell: Neural image caption generation with visual attention.

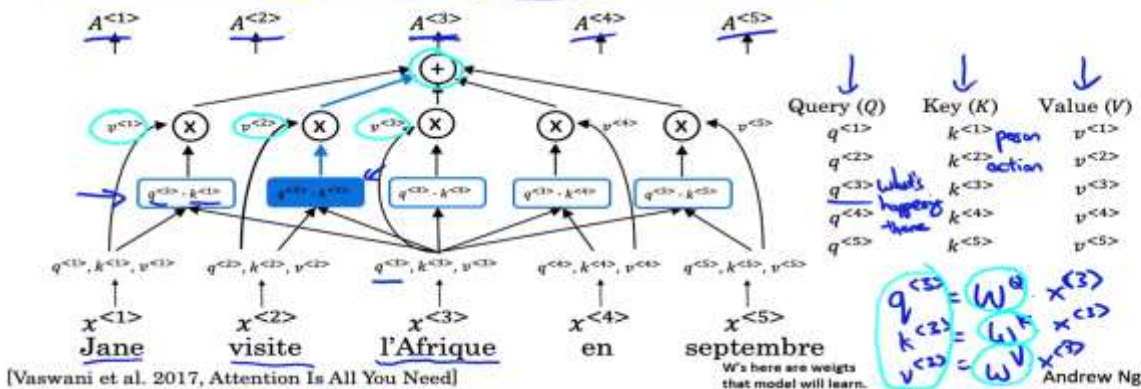
Andrew Ng

## Self-Attention

WHEN COMPUTING ALL THE 5 WORDS IS DONE, WE CAN REPRESENT THIS USING THIS FUNCTION.

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k^{<i>})}{\sum_j \exp(q \cdot k^{<j>})} v^{<i>}$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



[Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

## Multi-Head Attention

"head"

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) W_o$$

$$\text{head}_i = \text{Attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

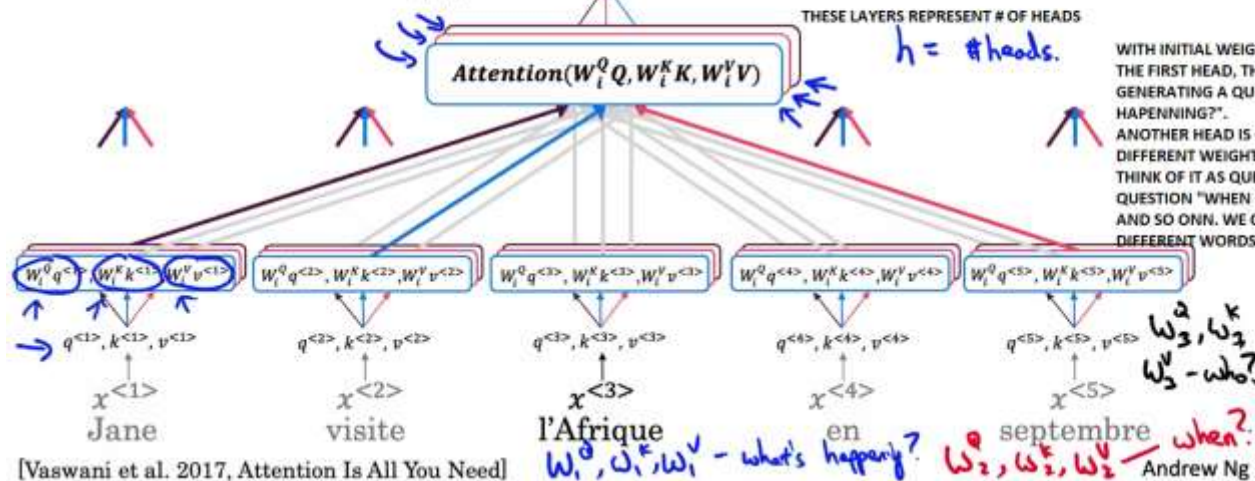
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

HERE WE CAN CONCATENATE AND MULTIPLY BY  $W_o$ , IT IS JUST A CONCEPTUAL WAY, BEHIND THE BARS IT IS COMPUTING PARALLEL.

THESE LAYERS REPRESENT # OF HEADS

 $h = \# \text{heads.}$ 

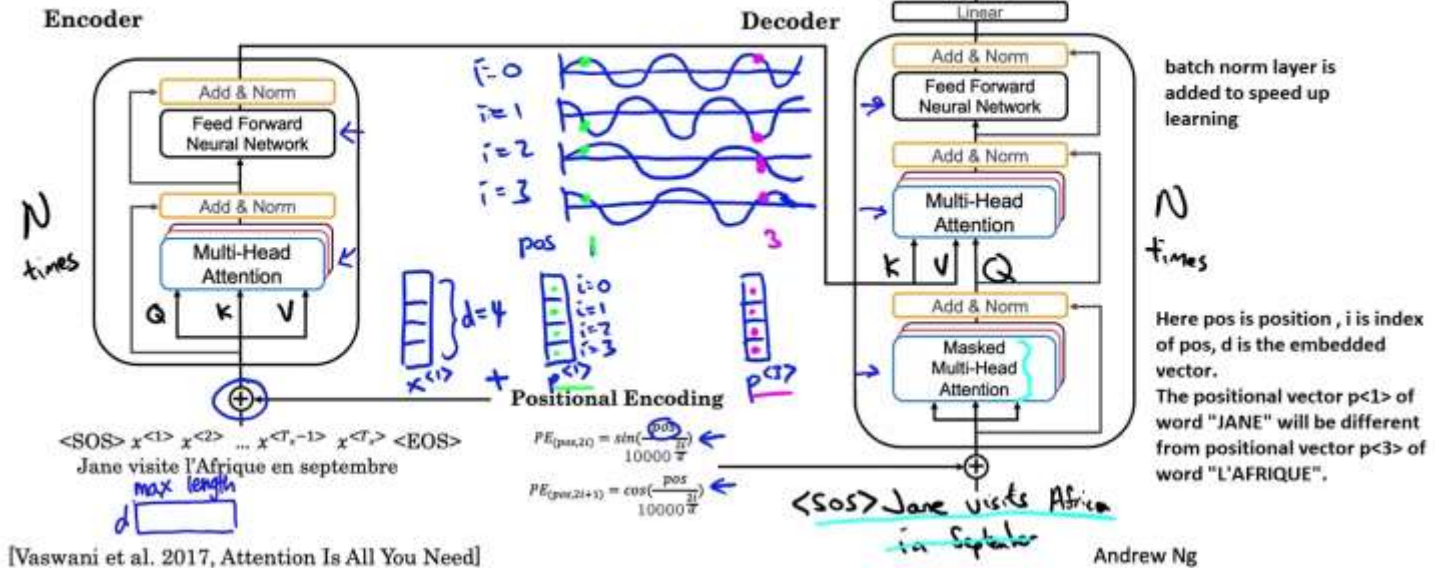
WITH INITIAL WEIGHTS WE WILL CALCULATE THE FIRST HEAD, THINK OF IT AS QUERY GENERATING A QUESTION "WHAT'S HAPPENING?". ANOTHER HEAD IS CALCULATED USING DIFFERENT WEIGHTS LEARNED BY MODEL, THINK OF IT AS QUERY GENERATING A QUESTION "WHEN IS SMTH HAPPENING?" AND SO ONN. WE CALCULATE THIS FOR DIFFERENT WORDS.



[Vaswani et al. 2017, Attention Is All You Need]

Andrew Ng

# Transformer Details



# Transformer

