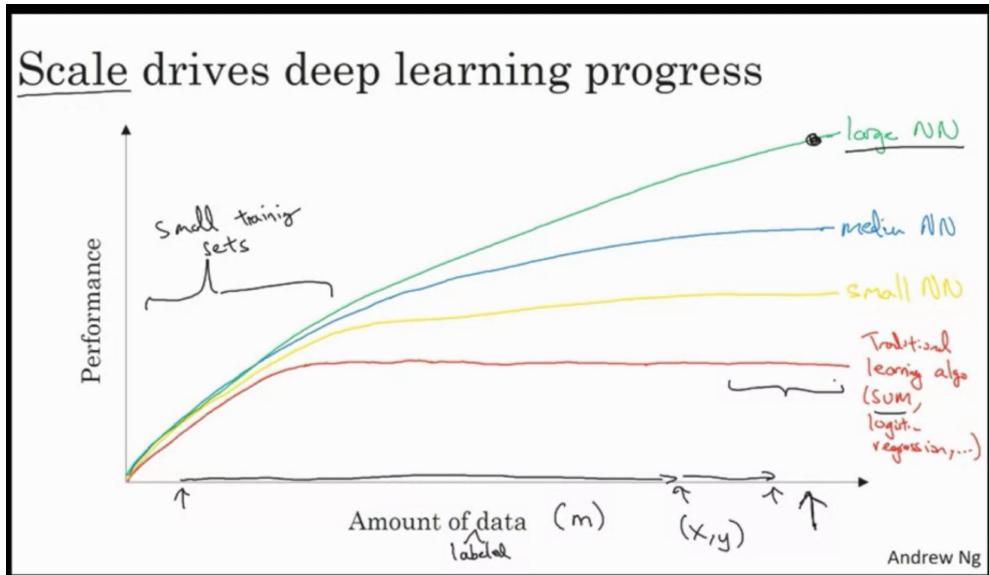


# Deep Learning (Coursera)

## Neural Networks and Deep Learning

Why deep learning take off?



## Course Outline

1. Neural Networks and Deep Learning
2. Improving Deep Neural networks: Hyperparameter tuning, Regularization and Optimization
3. Structuring your Machine Learning project
4. Convolutional Neural Networks
5. Natural Language Processing: Building sequence models

## Logistic Regression as a Neural Network

### Notation

$m$  denote the **number of training examples**

$n$  denote the **number of features**

A single training example:  $(x, y), x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

$m$  training example:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$m_{train} = \#train\ examples, m_{test} = \#test\ examples$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n_x \times m} \text{ (more easier to implement)}$$

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}] \in \mathbb{R}^{1 \times m}$$

## Logistic Regression

Given  $x$ , want  $\hat{y} = P(y=1|x)$ ,  $x \in \mathbb{R}^{n_x}$

Parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

Output:  $\hat{y} = \sigma(w^\top x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$

**Loss function:**

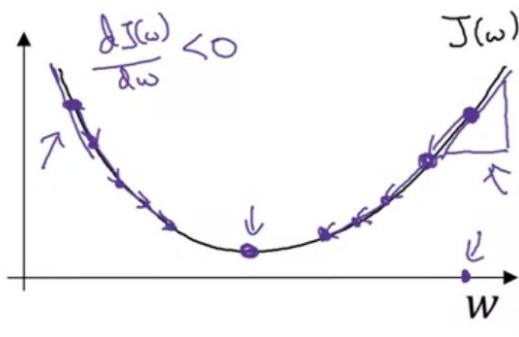
$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

**Cost function:**

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) = -\frac{1}{m} \sum_{i=1}^m [y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)]$$

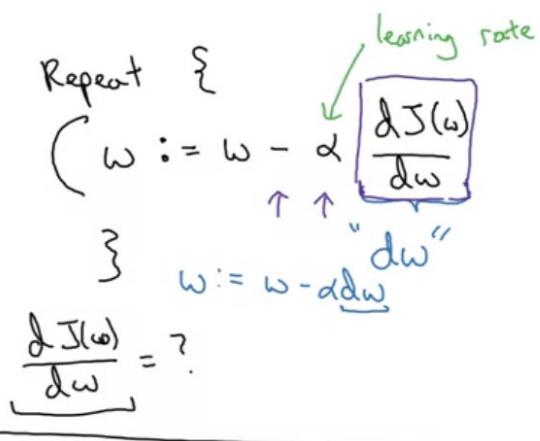
## Gradient Descent

### Gradient Descent



$$J(w, b) \quad w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

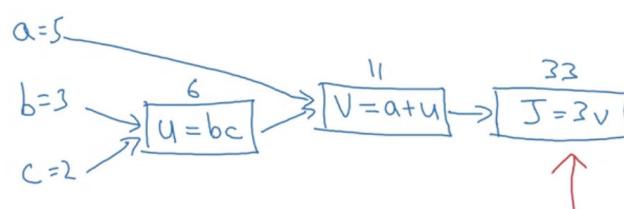
$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$



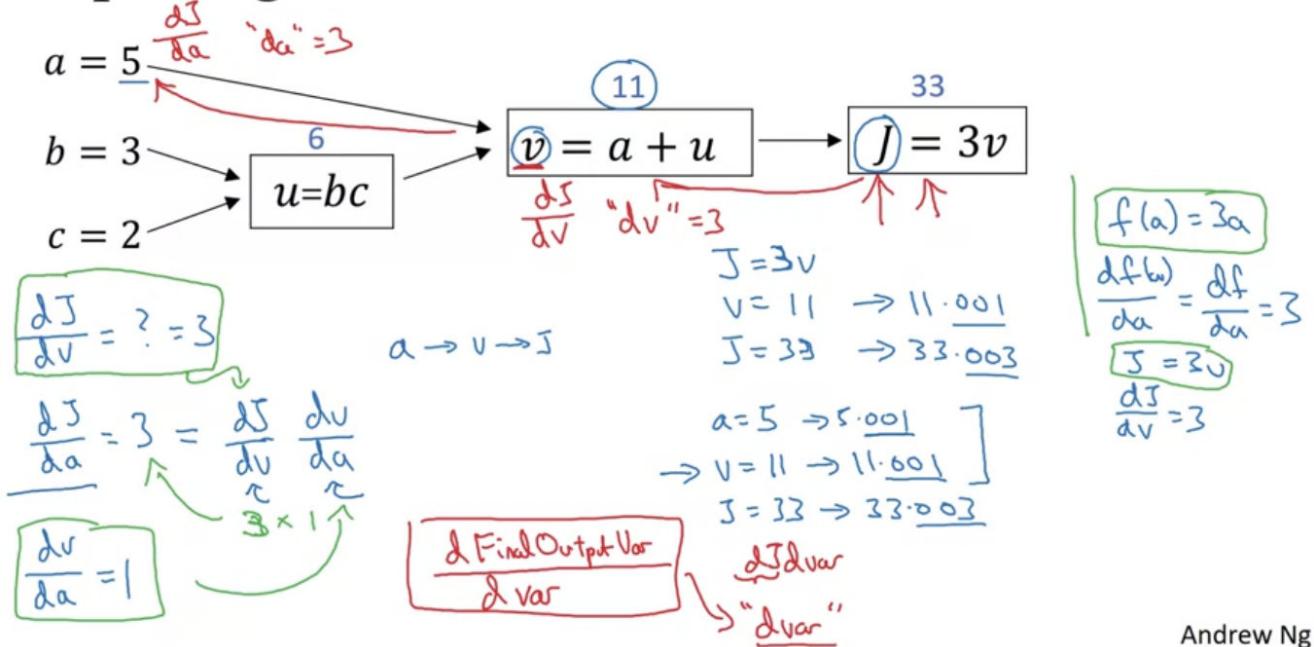
Andrew Ng

## Computation Graph

Forward



# Computing derivatives



$$\frac{dL(a, y)}{dz} = \frac{dL}{da} \frac{da}{dz} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right)a(1-a) = y(a-1) + a(1-y) = a - y$$

Logistic regression on m examples

## Logistic regression on $m$ examples

$$J=0; \underline{dw_1}=0; \underline{dw_2}=0; \underline{db}=0$$

For  $i = 1$  to  $m$

$$z^{(i)} = \omega^T x^{(i)} + b$$

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

$$J_t = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$\underline{dz^{(i)}} = a^{(i)} - y^{(i)}$$

$$\begin{aligned} \underline{dw_1} &+= x_1^{(i)} \underline{dz^{(i)}} \\ \underline{dw_2} &+= x_2^{(i)} \underline{dz^{(i)}} \end{aligned} \quad \left. \right\} n=2$$

$$\underline{db} += \underline{dz^{(i)}}$$

$$J /= m \leftarrow$$

$$\underline{dw_1} /= m; \underline{dw_2} /= m; \underline{db} /= m. \leftarrow$$

$$\underline{dw_1} = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \underline{dw_1}$$

$$w_2 := w_2 - \alpha \underline{dw_2}$$

$$b := b - \alpha \underline{db}$$

Andrew Ng

## Vectorization

# Logistic regression derivatives

$$\begin{aligned}
 J &= 0, \quad \boxed{\cancel{dw_1 = 0, dw_2 = 0}}, \quad db = 0 \quad dw = np.zeros((n_x, 1)) \\
 \rightarrow \text{for } i &= 1 \text{ to } m: \\
 z^{(i)} &= w^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J &+= -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \\
 \cancel{\text{for } j &= 1 \dots n_x:} \\
 dz^{(i)} &= a^{(i)}(1 - a^{(i)}) \\
 \cancel{dw_1 &+= x_1^{(i)} dz^{(i)}} \quad n_x = 2 \quad dw += \cancel{x^{(i)} dz^{(i)}} \\
 \cancel{dw_2 &+= x_2^{(i)} dz^{(i)}} \\
 db &+= dz^{(i)} \\
 J &= J/m, \quad \boxed{dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m} \quad dw /= m
 \end{aligned}$$

## Verctorizing Logistic Regression

$$\begin{aligned}
 X &= \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n_x \times m} \\
 [z^{(1)} & z^{(2)} \dots z^{(m)}] = [w^T x^{(1)} + b & w^T x^{(2)} + b \dots w^T x^{(m)} + b] \\
 &= [w^T x^{(1)} & w^T x^{(2)} \dots w^T x^{(m)}] + [b & b \dots b] \\
 &= [w^T & w^T \dots w^T] \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} + \mathbf{b} \\
 &= \mathbf{w}^T \mathbf{X} + \mathbf{b}
 \end{aligned}$$

```
Z = np.dot(w.T, X) + b # b is real number
```

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

## Gradient computation

$$\begin{aligned}
 dZ &= A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}] \\
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 dw &= \frac{1}{m} X(dZ)^\top
 \end{aligned}$$

# Vectorizing Logistic Regression

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots$$

$$d\bar{z} = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}] \quad \leftarrow$$

$$A = [a^{(1)} \ \dots \ a^{(m)}]. \quad Y = [y^{(1)} \ \dots \ y^{(m)}]$$

$$\rightarrow d\bar{z} = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]$$

$$\begin{aligned} \rightarrow dw &= 0 \\ dw &+ \frac{\partial}{\partial w} X^{(1)} dz^{(1)} \\ dw &+ \frac{\partial}{\partial w} X^{(2)} dz^{(2)} \\ &\vdots \\ dw &/ m \end{aligned}$$

$$\begin{aligned} db &= 0 \\ db &+ \frac{\partial}{\partial b} dz^{(1)} \\ db &+ \frac{\partial}{\partial b} dz^{(2)} \\ &\vdots \\ db &+ \frac{\partial}{\partial b} dz^{(m)} \\ db &/ m \end{aligned}$$

$$\begin{aligned} db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\ &= \frac{1}{m} \text{np.sum}(d\bar{z}) \\ dw &= \frac{1}{m} X d\bar{z}^T \\ &= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\ &= \frac{1}{m} \begin{bmatrix} x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)} \end{bmatrix}_{n \times 1} \end{aligned}$$

Andrew Ng

```

Z = np.dot(w.T, X) + b
A = sigmoid(Z)
dZ = A - Y
dw = 1/m * X * dZ.T
db = 1/m * np.sum(dZ)

w -= lr * dw
b -= lr * db
    
```

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{1}{m} X(A - Y)^T \\ \frac{\partial J}{\partial b} &= \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) \end{aligned}$$

Broadcasting in Python

# General Principle

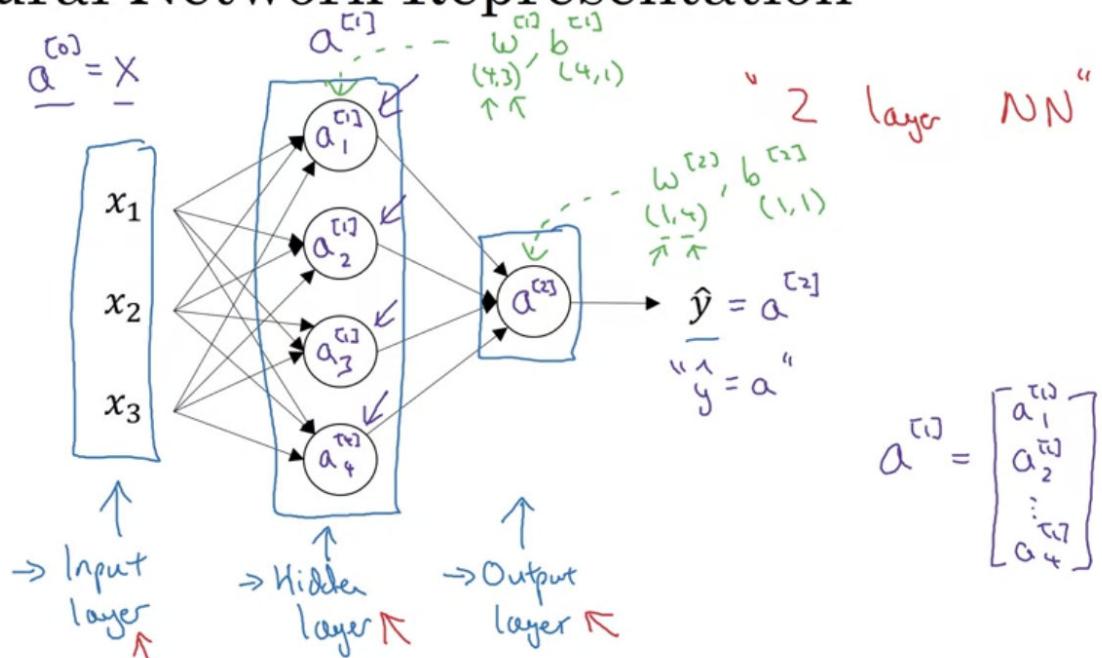
$$\begin{array}{c} (m, n) \\ \text{matrix} \\ \hline + \\ \times \\ / \end{array} \quad \begin{array}{c} (l, n) \\ (m, l) \end{array} \quad \rightsquigarrow (m, n)$$

$$\begin{array}{ccc} (m, 1) & + & R \\ \left[ \begin{smallmatrix} 1 \\ 2 \\ 3 \end{smallmatrix} \right] & + & 100 \\ \left[ \begin{smallmatrix} 1 & 2 & 3 \end{smallmatrix} \right] & + & 100 \end{array} = \begin{array}{c} \left[ \begin{smallmatrix} 101 \\ 102 \\ 103 \end{smallmatrix} \right] \\ = \left[ \begin{smallmatrix} 101 & 102 & 103 \end{smallmatrix} \right] \end{array}$$

Matlab/Octave: `bsxfun`

## Shallow Neural Network

### Neural Network Representation



$$a^{[l]} = \begin{bmatrix} a_1^{[l]} \\ a_2^{[l]} \\ \vdots \\ a_n^{[l]} \end{bmatrix}$$

Andrew Ng

$$a_{i \leftarrow \text{node}}^{[l \leftarrow \text{layer}]}$$

## Vetorization

Single example

$$W^{[l]} = \begin{bmatrix} --(w_1^{[l]})^\top-- \\ --(w_2^{[l]})^\top-- \\ \vdots \\ --(w_n^{[l]})^\top-- \end{bmatrix}$$

$$\begin{aligned} z^{[l]} &= W^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= \sigma(z^{[l]}) \end{aligned}$$

Across the training set

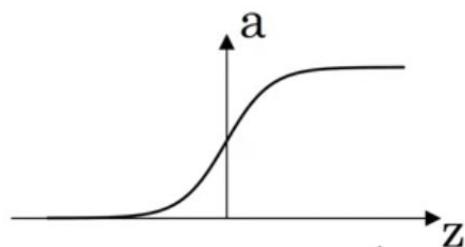
$$\begin{aligned} Z^{[l]} &= \begin{bmatrix} | & | & & | \\ z^{[l](1)} & z^{[l](2)} & \dots & z^{[l](m)} \\ | & | & & | \end{bmatrix} \\ A^{[l]} &= \begin{bmatrix} | & | & & | \\ a^{[l](1)} & a^{[l](2)} & \dots & a^{[l](m)} \\ | & | & & | \end{bmatrix} \\ A^{[l]} &= \sigma(W^{[l]} A^{[l-1]} + b^{[l]}), \quad A^{[0]} = X \end{aligned}$$

## Activation function

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

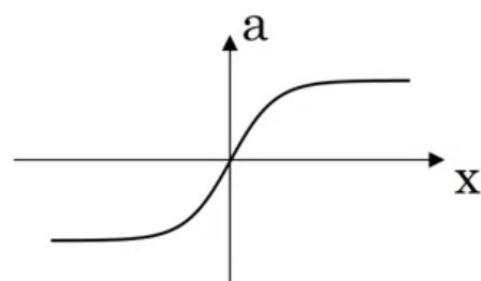
$g$  is the activation function

Sigmoid function (Never use this except for output layer)



$$a = \frac{1}{1 + e^{-z}}$$

Hyperbolic tangent function (tanh)

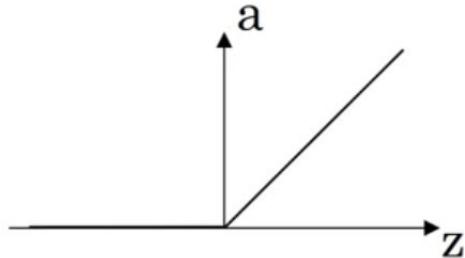


$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

A shifted version of sigmoid function, works **better** than sigmoid.

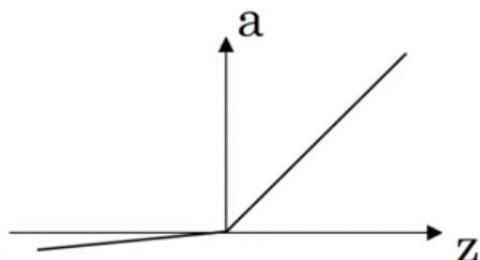
However, for both sigmoid and tanh function, when the z is either very large or very small, the gradient becomes very small.

### Rectified linear unit function (ReLU) -- Default choice



$$a = \max(0, z)$$

### Leaky ReLU



$$a = \max(0.01z, z)$$

### Why do we need non-linear activation function?

If there is no non-linear activation function, there is no hidden layers, because it is just the combination of linear function.

### Derivatives of activation function

#### Sigmoid

$$\frac{d}{dz} g(z) = g(z)(1 - g(z))$$

#### tanh

$$\frac{d}{dz} g(z) = 1 - (\tanh(z))^2 = 1 - (g(z))^2$$

#### ReLU

$$\frac{d}{dz} g(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{otherwise} \end{cases}$$

## Leaky ReLU

$$\frac{d}{dz} g(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{otherwise} \end{cases}$$

The gradient at  $z = 0$  for ReLU and Leaky ReLU is technically not defined. BUT we just take 1.

## Gradient decent for neural networks

### Forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = \sigma(Z^{[l]}), A^{[0]} = X$$

### Backward propagation

$$dZ^{[L]} = A^{[L]} - Y$$

Back propagation:

$$dZ^{[2]} = A^{[2]} - Y \quad \leftarrow \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T} \quad \leftarrow \quad (n^{[2]}) \leftarrow$$

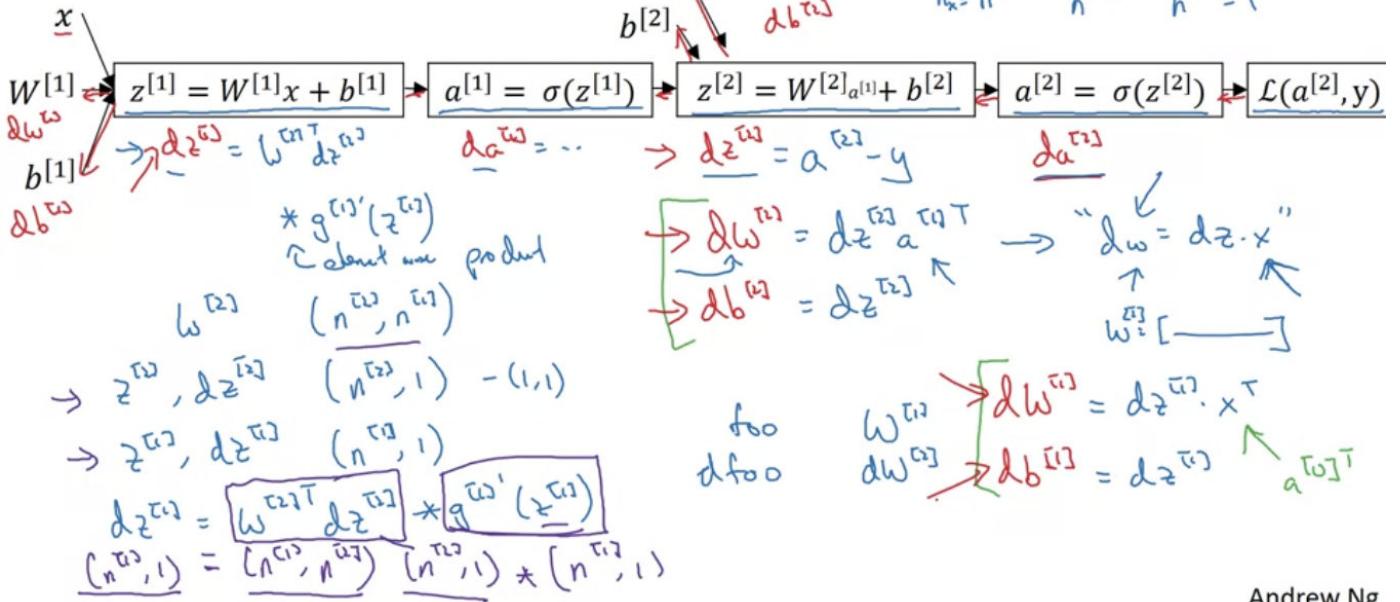
$$db^{[2]} = \frac{1}{m} \underbrace{\text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})}_{\text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims=True})} \quad \leftarrow \quad (n^{[2]}, 1) \leftarrow$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[2]}, m)} \times \underbrace{g^{[1]'}(Z^{[1]})}_{\text{element-wise product}} \quad (n^{[1]}, m)$$

$$dW^{[1]} = \frac{1}{n} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{n} \underbrace{\text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})}_{(n^{[1]}, 1)} \quad (n^{[1]}, 1)$$

# Neural network gradients



Andrew Ng

## Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} X^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{\text{elementwise product}} * \underbrace{g^{[1]'}(Z^{[1]})}_{(n^{[1]}, m)}$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

Andrew Ng

## Random initialization

The purpose of random initialization is to **break symmetry**.

$$w^{[1]} = np.random.randn((2, 2)) * 0.01$$

$$b^{[1]} = np.zeros((2, 1))$$

$b$  do not have symmetric breaking problem.

The general methodology to build a Neural Network is to:

1. Define the neural network structure (# of input units, # of hidden units, etc).
2. Initialize the model's parameters
3. Loop:
  - o Implement forward propagation
  - o Compute loss
  - o Implement backward propagation to get the gradients
  - o Update parameters (gradient descent)

## Deep Neural Network

### Notation

$L$  is the number of layers

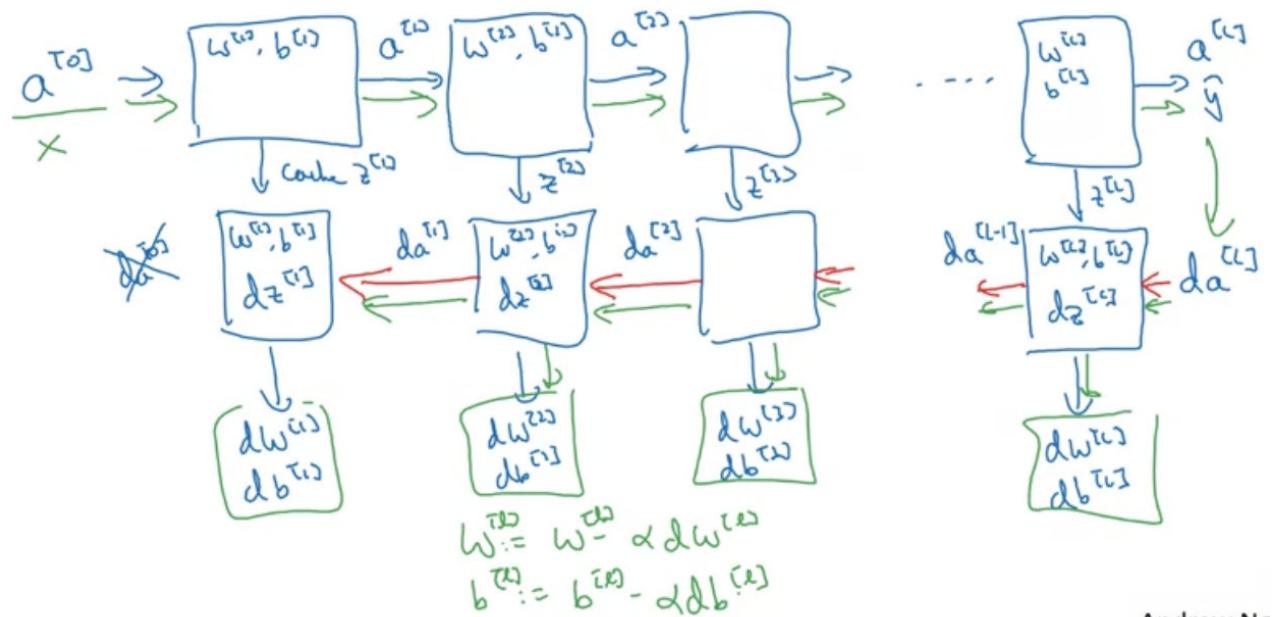
$n^{[l]}$  is the number of units in layer  $l$ ,  $n^{[0]} = n_x$

$a^{[l]}$  is the activations in layer  $l$

$w^{[l]}, b^{[l]}$  are parameters for  $z^{[l]}$

$x = a^{[0]}, \hat{y} = a^{[L]}$

## Forward and backward functions



Andrew Ng

## Forward propagation for layer l

Input  $a^{[l-1]}$

Output  $a^{[l]}$ , cache( $z^{[l]}$ )

$$\begin{aligned} Z^{[l]} &= W^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} &= \sigma(Z^{[l]}), \quad A^{[0]} = X \end{aligned}$$

The dimension of  $W^{[l]}$  and  $b^{[l]}$

$$W^{[l]} \in \mathbb{R}^{n^{[l]} \times n^{[l-1]}}$$

$$b^{[l]} \in \mathbb{R}^{n^{[l]} \times 1}$$

## Backward propagation for layer l

Input  $da^{[l]}$

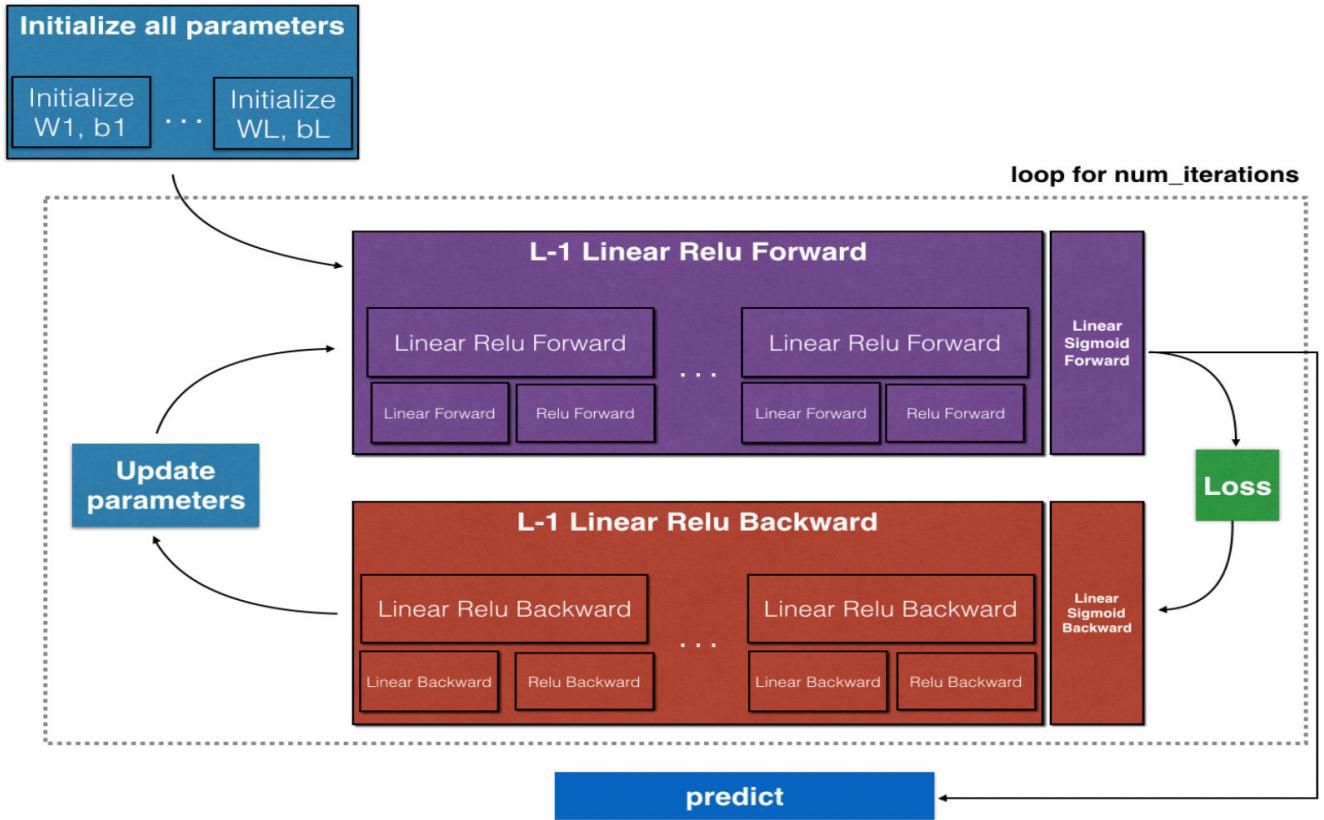
Output  $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$\begin{aligned} dZ^{[l]} &= dA^{[l]} * g^{[l]\prime}(Z^{[l]}) \\ dW^{[l]} &= \frac{1}{m} dZ^{[l]} A^{[l-1]\top} \\ db^{[l]} &= \frac{1}{m} \text{np.sum}(dZ^{[l]}, \text{axis}=1, \text{keepdims=True}) \\ dA^{[l-1]} &= W^{[l]\top} dZ^{[l]} \end{aligned}$$

$$\begin{aligned} dW^{[l]} &= \frac{\partial \mathcal{J}}{\partial W^{[l]}} = \frac{1}{m} dZ^{[l]} A^{[l-1]\top} \\ db^{[l]} &= \frac{\partial \mathcal{J}}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dZ^{[l](i)} \\ dA^{[l-1]} &= \frac{\partial \mathcal{L}}{\partial A^{[l-1]}} = W^{[l]\top} dZ^{[l]} \end{aligned}$$

For final layer

$$dA^{[L]} = \sum_{i=1}^m \left( -\frac{y^{(i)}}{a^{(i)}} + \frac{1 - y^{(i)}}{1 - a^{(i)}} \right)$$



## Hyperparameters

- Learning rate  $\alpha$
- #iterations
- #hidden layer  $L_q$
- #Hidden units  $n^{[l]}$
- choice of activation function
- ...

# Improving Deep Neural Networks: Hyperparameter Tuning, Regularization and Optimization

## Train/dev/test sets

- Small dataset -- 70/30, 60/20/20 traditional ratio.
- large dataset -- set the dev and test set smaller than traditional ratio.

**MAKE SURE** the dev and test sets come from **same distribution**.

## Bias and Variance

## High Bias

- More complex model (bigger network)
- ...

## High Variance

- More data
- Regularization
- ...

## Regularization

### Logistic regression

#### L2 regularization

$$\frac{\lambda}{2m} \|w\|_2^2 = \frac{\lambda}{2m} \sum_{j=1}^{n_x} w_j^2 = \frac{\lambda}{2m} w^\top w$$

L1 regularization ( $w$  will be sparse -- lots of 0 in  $w$ )

$$\frac{\lambda}{m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{m} \|w\|_1$$

## Neural Network

### Frobenius norm

$$\begin{aligned} & \frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|_F^2 \\ & \|W^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (W_{i,j}^{[l]})^2 \end{aligned}$$

$$dW^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} W^{[l]}$$

L2 regularization is also called **weight decay** ( $1 - \frac{\alpha\lambda}{m} < 1$ ).

$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha[(\text{from backprop}) + \frac{\lambda}{m} W^{[l]}] \\ &:= W^{[l]} - \frac{\alpha\lambda}{m} W^{[l]} - \alpha(\text{from backprop}) \\ &:= (1 - \frac{\alpha\lambda}{m}) W^{[l]} - \alpha(\text{from backprop}) \end{aligned}$$

$$J_{regularized} = \underbrace{-\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}) \right)}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{k,j}^{[l]2}}_{\text{L2 regularization cost}} \quad (2)$$

## Dropout regularization

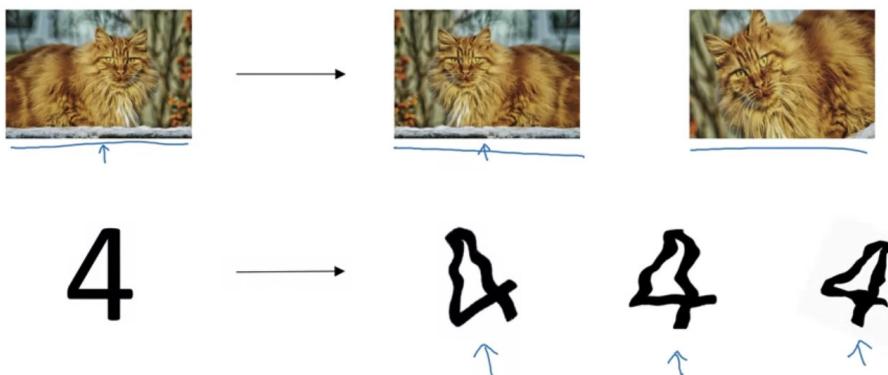
Implementing dropout ("Inverted dropout") -- with l=3 layer

```
d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob  
a3 = np.multiply(a3, d3) # a3 *= d3  
a3 /= keep_prob # Because a3 is reduced by (1 - keep_prob)
```

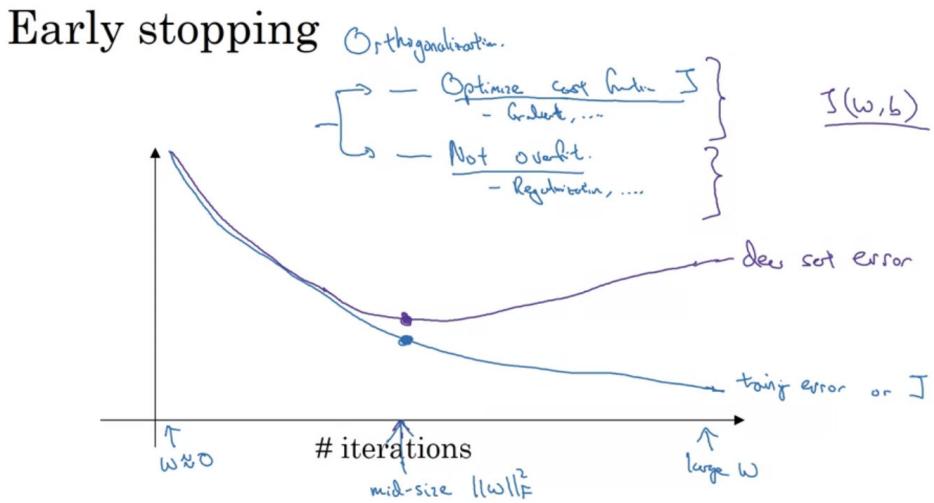
Intuition: Can't rely on any one feature, so have to spread out weights.

## Other regularization methods

- Data augmentation



- Early stopping



## Setting Up Optimization Problem

### Normalizing inputs

1. Move the training set until it has 0 mean.

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$
$$x := x - \mu$$

2. Normalize the variance

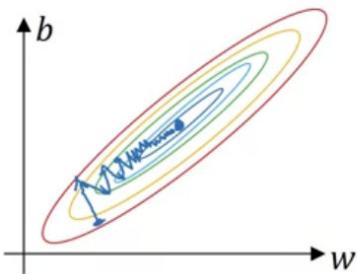
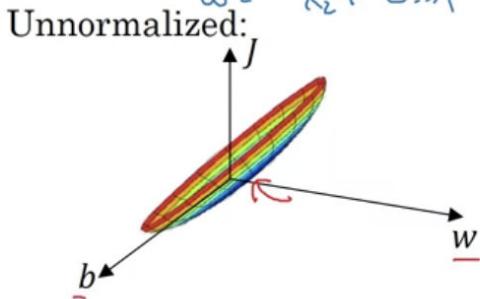
$$\sigma = \sqrt{\frac{1}{m} \sum_{i=1}^m (x^{(i)})^2}$$

$$x := \frac{x}{\sigma}$$

Use same  $\mu$  and  $\sigma$  to normalize the training set and test set.

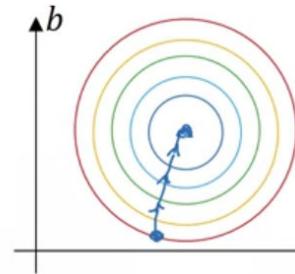
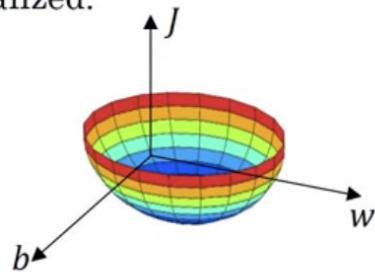
## Why normalize inputs?

$w_1 \quad x_1: 1 \dots \text{load}$   
 $w_2 \quad x_2: 0 \dots 1$



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

Normalized:



Andrew Ng

## Vanishing/exploding gradients

For a very deep network,

$w > 1 \rightarrow$  exploding gradient

$w < 1 \rightarrow$  vanishing gradient

## Weight initialization for deep networks (partial solution for v/e gradients)

ReLU

He Initialization

$$W^{[l]} = \text{np.random.randn(shape)} * \sqrt{\frac{2}{n^{[l-1]}}}$$

tanh

Xavier initialization

$$W^{[l]} = \text{np.random.randn(shape)} * \sqrt{\frac{1}{n^{[l-1]}}}$$

Other version

$$W^{[l]} = \text{np.random.randn(shape)} * \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

## Numerical approximation of gradients

For a very small  $\epsilon$ ,

$$f'(\theta) \approx \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$$

Use this formula for gradient checking.

## Gradient checking for a neural network

Take  $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$  and reshape into a big vector  $\theta$ .

Take  $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$  and reshape into a big vector  $d\theta$ .

For each i:

$$\begin{aligned} d\theta_{approx}[i] &= \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon} \\ &\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \end{aligned}$$

To see whether  $d\theta_{approx} = d\theta$

$$\text{Check } \frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$$

### Notes for grad check:

- Don't use in training - only to debug.
- If algorithm fails grad check, look at components ( $db^{[l]}, dw^{[l]}$ ) to try to identify bug.
- Remember regularization.
- Doesn't work with dropout.
- Run at random initialization; perhaps again after some training.

## Optimization Algorithms

### Mini-batch gradient descent

Mini-batch t:  $X^{\{t\}}, Y^{\{t\}}$

$x^{(i)}$  → the ith example

$z^{[l]}$  → the z value for the L layer of the neural network

$X^{\{t\}}, Y^{\{t\}}$  → mini-batch

# Mini-batch gradient descent

repeat {  
for  $t = 1, \dots, 5000$  {

Forward prop on  $X^{t+1}$ .

$$z^{t+1} = w^{(1)} X^{t+1} + b^{(1)}$$

$$A^{t+1} = g^{(1)}(z^{t+1})$$

$$\vdots$$

$$A^{t+1} = g^{(L)}(z^{t+1})$$

$$\text{Compute cost } J^{t+1} = \frac{1}{m} \sum_{i=1}^m L(A^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot m} \sum_{l=1}^L \|w^{(l)}\|_F^2.$$

Backprop to compute gradients w.r.t  $J^{t+1}$  (using  $(X^{t+1}, Y^{t+1})$ )

$$w^{(l)} := w^{(l)} - \alpha \nabla J^{t+1} \quad b^{(l)} := b^{(l)} - \alpha \nabla b^{(l)}$$

3  
3

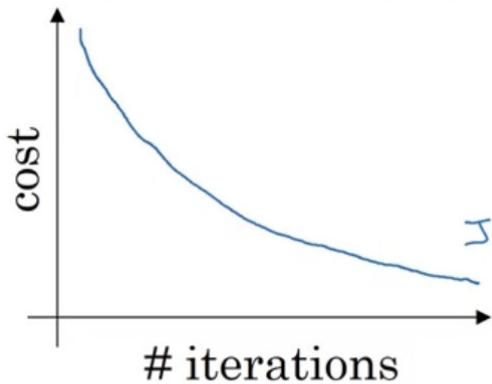
"1 epoch"  
pass through training set.

1 step of gradient descent  
using  $\underline{X^{t+1}}, \underline{Y^{t+1}}$   
(as if  $m=1$ )

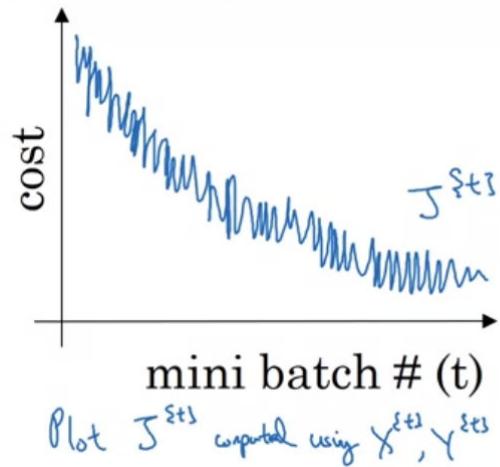
$\underline{X}, \underline{Y}$

Andrew Ng

## Batch gradient descent



## Mini-batch gradient descent



## Choosing the mini-batch size

If mini-batch size =  $m$ : Batch gradient descent.

If mini-batch size = 1: Stochastic gradient descent.

IF small train set ( $m \leq 2000$ ): Use batch GD

Typical mini-batch size: 32, 64, 128, 256, 512

The size is to the power of 2 (the computer will run faster).

Make sure mini-batch fit in CPU/GPU memory.

## Set Up

1. **Shuffle:** Create a shuffled version of the training set ( $X, Y$ ) as shown below. Each column of  $X$  and  $Y$  represents a training example. Note that the random shuffling is done synchronously between  $X$  and  $Y$ . Such that after the shuffling the  $i^{th}$  column of  $X$  is the example corresponding to the  $i^{th}$  label in  $Y$ . The shuffling step ensures that examples will be split randomly into different mini-batches.

$$X = \begin{pmatrix} x_0^{(1)} & x_0^{(2)} & \dots & x_0^{(m-1)} & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m-1)} & x_1^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(1)} & x_{12286}^{(2)} & \dots & x_{12286}^{(m-1)} & x_{12286}^{(m)} \\ x_{12287}^{(1)} & x_{12287}^{(2)} & \dots & x_{12287}^{(m-1)} & x_{12287}^{(m)} \end{pmatrix}$$

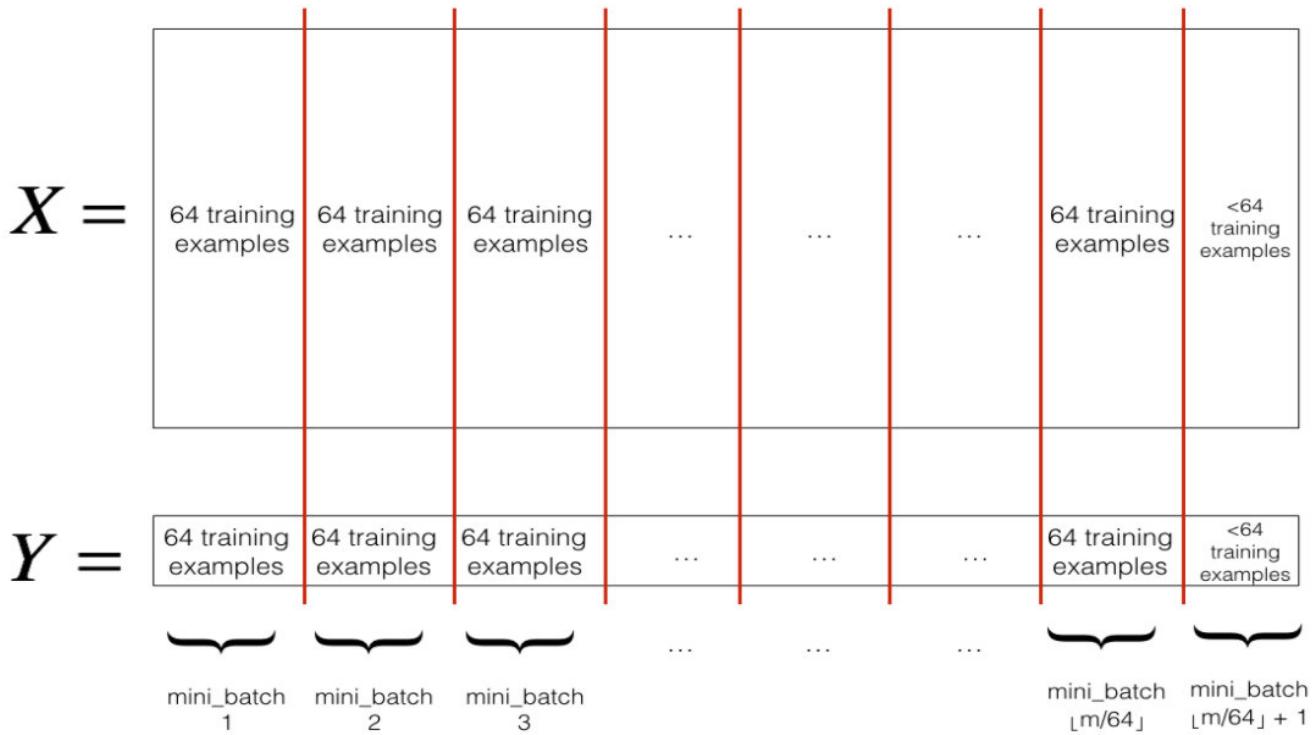
$$Y = \begin{pmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m-1)} & y^{(m)} \end{pmatrix}$$
  

$$X = \begin{pmatrix} x_0^{(5)} & x_0^{(16)} & \dots & x_0^{(2)} & x_0^{(m-1)} \\ x_1^{(5)} & x_1^{(16)} & \dots & x_1^{(2)} & x_1^{(m-1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{12286}^{(5)} & x_{12286}^{(16)} & \dots & x_{12286}^{(2)} & x_{12286}^{(m-1)} \\ x_{12287}^{(5)} & x_{12287}^{(16)} & \dots & x_{12287}^{(2)} & x_{12287}^{(m-1)} \end{pmatrix}$$

$$Y = \begin{pmatrix} y^{(5)} & y^{(16)} & \dots & y^{(2)} & y^{(m-1)} \end{pmatrix}$$

```
permutation = list(np.random.permutation(m)) #m为样本数
shuffled_X = X[:, permutation]
shuffled_Y = Y[:, permutation].reshape((1,m))
```

2. **Partition:** Partition the shuffled ( $X, Y$ ) into mini-batches of size `mini_batch_size` (here 64). Note that the number of training examples is not always divisible by `mini_batch_size`. The last mini batch might be smaller, but you don't need to worry about this. When the final mini-batch is smaller than the full `mini_batch_size`, it will look like this:



### Exponentially weighted averages

$$V_t = \beta V_{t-1} + (1 - \beta) \theta_t$$

$V_t$  is approximately averaging over  $\frac{1}{1-\beta}$  examples

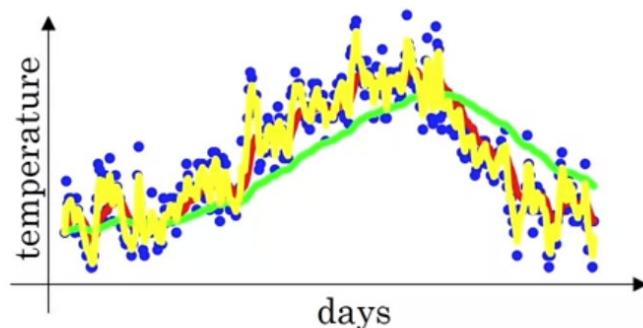
$$(1 - \epsilon)^{\frac{1}{\epsilon}} = \frac{1}{e} \approx 0.35$$

### Exponentially weighted averages Moving

$$V_t = \underline{\beta} \underline{V_{t-1}} + \underline{(1-\beta)} \underline{\theta_t} \leftarrow$$

$\underline{\beta} = 0.9$  :  $\approx 10$  days' temper.  
 $\underline{\beta} = 0.98$  :  $\approx 50$  days  
 $\underline{\beta} = 0.5$  :  $\approx 2$  days

$V_t$  is approximately  
 average over  
 $\rightarrow \approx \frac{1}{1-\beta}$  days'  
 temper.



$$\frac{1}{1-0.98} = 50$$

Andrew Ng

$$\rightarrow V_0 = 0$$

Repeat {

Get next  $\theta_t$

$$V_{\theta} := \beta V_{\theta} + (1-\beta) \theta_t \leftarrow$$

}

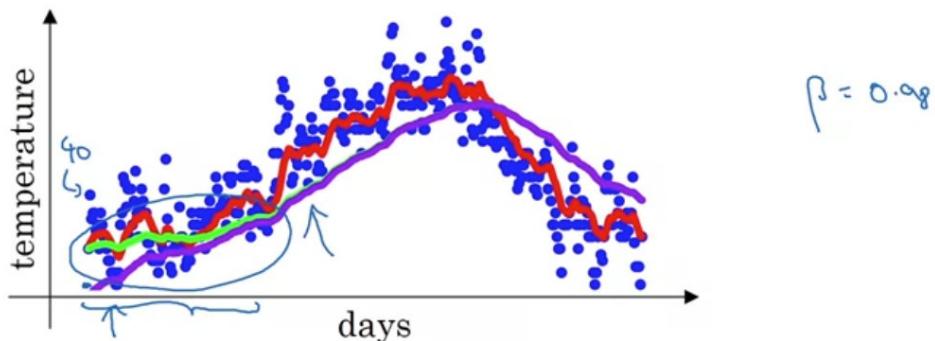
Andrew Ng

The **efficiency** of the exponentially weighted averages is because it the computer can overwrite  $V_{\theta}$  over and over (and only 1 line code), and this requires very small memory. In contrast, the direct computation of average requires very large memory when averaging over many examples.

### Bias correction

$$V_t = \frac{\beta V_{t-1} + (1 - \beta) \theta_t}{1 - \beta^t}$$

## Bias correction



$$\rightarrow v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \underline{\theta_1} + 0.02 \underline{\theta_2}$$

$$\frac{V_t}{1 - \beta^t}$$

$t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396$

$$\frac{V_2}{0.0396} = \frac{0.0196 \underline{\theta_1} + 0.02 \underline{\theta_2}}{0.0396}$$

Andrew Ng

Normally, people just wait for warming up in practice.

### Gradient descent with momentum

Now, implement the parameters update with momentum. The momentum update rule is, for  $l = 1, \dots, L$ :

$$\begin{cases} v_{dW^{[l]}} = \beta v_{dW^{[l]}} + (1 - \beta) dW^{[l]} \\ W^{[l]} = W^{[l]} - \alpha v_{dW^{[l]}} \end{cases}$$

$$\begin{cases} v_{db^{[l]}} = \beta v_{db^{[l]}} + (1 - \beta) db^{[l]} \\ b^{[l]} = b^{[l]} - \alpha v_{db^{[l]}} \end{cases}$$

In one sentence, the basic idea is to **compute an exponentially weighter average of gradients**, and then use that gradient to update weights instead.

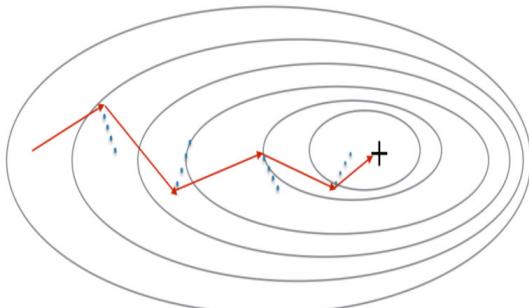
```

VdW, Vdb = zeros(same dimension with W and b)
On iteration t:
    Compute dW, db on current mini-batch
    VdW = beta * VdW + (1 - beta) * dW
    Vdb = beta * Vdb + (1 - beta) * db
    W = W - learning_rate * VdW
    b = b - learning_rate * Vdb

```

The most common value for  $\beta$  is 0.9.

In some literature,  $(1 - \beta)$  is omitted. That is just scaling the  $VdW$  or  $Vdb$  by  $\frac{1}{1-\beta}$  (it will affect the learning rate).



**Figure 3:** The red arrows show the direction taken by one step of mini-batch gradient descent with momentum. The blue points show the direction of the gradient (with respect to the current mini-batch) on each step. Rather than just following the gradient, the gradient is allowed to influence  $v$  and then take a step in the direction of  $v$ .

## RMSProp (Root Mean Square Prop)

```

VdW, Vdb = zeros(same dimension with W and b)
On iteration t:
    Compute dW, db on current mini-batch
    SdW = beta * VdW + (1 - beta) * dW ** 2
    Sdb = beta * Vdb + (1 - beta) * db ** 2
    W = W - learning_rate * dW / (sqrt(SdW + epsilon)) #make sure not 0
    b = b - learning_rate * db / (sqrt(Sdb + epsilon))

```

## Adam optimization algorithm (Adaptive Moment Estimation)

```

VdW = 0, SdW = 0, Vdb = 0, Sdb = 0
On iteration t:
    Compute dW, db on current mini-batch
    VdW = beta_1 * VdW + (1 - beta_1) * dW #'momentum' beta_1
    Vdb = beta_1 * Vdb + (1 - beta_1) * db #'momentum' beta_1
    SdW = beta_2 * VdW + (1 - beta_2) * dW ** 2 #'RMSProp' beta_2
    Sdb = beta_2 * Vdb + (1 - beta_2) * db ** 2 #'RMSProp' beta_2
    VdW_corrected = VdW / (1 - beta_1 ** t)
    Vdb_corrected = Vdb / (1 - beta_1 ** t)
    SdW_corrected = SdW / (1 - beta_2 ** t)
    Sdb_corrected = Sdb / (1 - beta_2 ** t)
    W = W - learning_rate * VdW_corrected / (sqrt(SdW_corrected + epsilon))
    b = b - learning_rate * Vdb_corrected / (sqrt(Sdb_corrected + epsilon))

```

# Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0. \quad V_{db} = 0, S_{db} = 0$$

On iteration  $t$ :

Compute  $\partial W, \partial b$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \partial W, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) \partial b \quad \leftarrow \text{"moment" } \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \partial W^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) \partial b^2 \quad \leftarrow \text{"RMSprop" } \beta_2$$

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

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

$$W := W - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}} \quad b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

How does Adam work?

- It calculates an exponentially weighted average of past gradients, and stores it in variables  $v$  (before bias correction) and  $v^{corrected}$  (with bias correction).
- It calculates an exponentially weighted average of the squares of the past gradients, and stores it in variables  $s$  (before bias correction) and  $s^{corrected}$  (with bias correction).
- It updates parameters in a direction based on combining information from "1" and "2".

The update rule is, for  $l = 1, \dots, L$ :

$$\begin{cases} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1) \frac{\partial \mathcal{J}}{\partial W^{[l]}} \\ v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1 - (\beta_1)^t} \\ s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2) \left( \frac{\partial \mathcal{J}}{\partial W^{[l]}} \right)^2 \\ s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1 - (\beta_2)^t} \\ W^{[l]} = W^{[l]} - \alpha \frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected} + \epsilon}} \end{cases}$$

Hyperparameters choice:

$\alpha$ : needs to be tuned

$\beta_1$ : 0.9 (recommend)  $\rightarrow dw$  (first moment)

$\beta_2$ : 0.999 (recommend)  $\rightarrow dw^2$  (second moment)

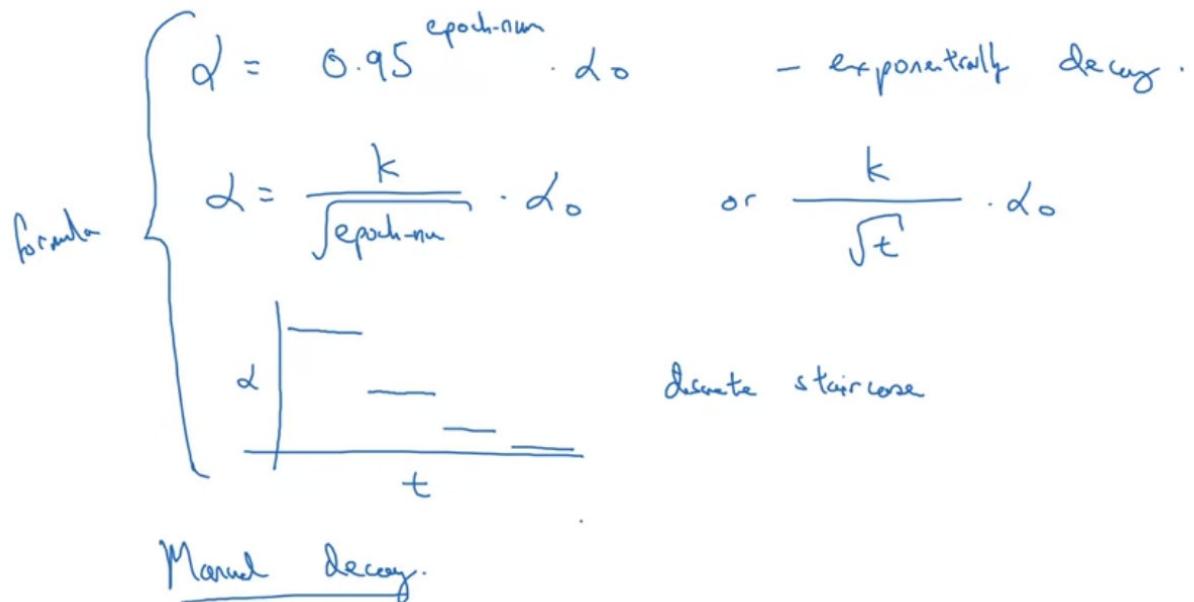
$\epsilon$ :  $10^{-8}$  (recommend)

## Learning rate decay

1 epoch = 1 pass through data

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

## Other learning rate decay methods



Andrew Ng

## Hyperparameter Tuning

(When there are many hyperparameters) Use **random sampling** instead of **grid search**!

### Appropriate scale for hyperparameters

It seems more reasonable to search for hyperparameters on a log scale.

```
r = -4 * np.random.rand() # r in [-4, 0]
alpha = 10 ** r
```

When tuning hyperparameters for exponentially weighted averages, instead of sample  $\beta$ , we sample  $(1 - \beta)$  on a log scale ( $\beta = 1 - 10^r$ ).

## Batch Normalization

Give some intermediate values in NN  $z^{(1)}, \dots, z^{(m)}$

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

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

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta$$

where  $\gamma$  and  $\beta$  here are learnable parameters,  $\gamma$  and  $\beta$  control the mean and variance of  $\tilde{z}^{(i)}$ .

## Implementing gradient descent

for  $t = 1 \dots \text{numMiniBatches}$   
 Compute forward pass on  $X^{[t]}$ .

In each hidden layer, use BN to replace  $\underline{z}^{[t]}$  with  $\hat{z}^{[t]}$ .  
 Use backprop to compute  $\underline{dW}^{[t]}, \underline{d\beta}^{[t]}, \underline{d\gamma}^{[t]}$   
 Update parameters  $\left. \begin{array}{l} W^{[t]} := W^{[t]} - \alpha \underline{dW}^{[t]} \\ \beta^{[t]} := \beta^{[t]} - \alpha \underline{d\beta}^{[t]} \\ \gamma^{[t]} := \dots \end{array} \right\} \leftarrow$

Works w/ momentum, RMSprop, Adam.

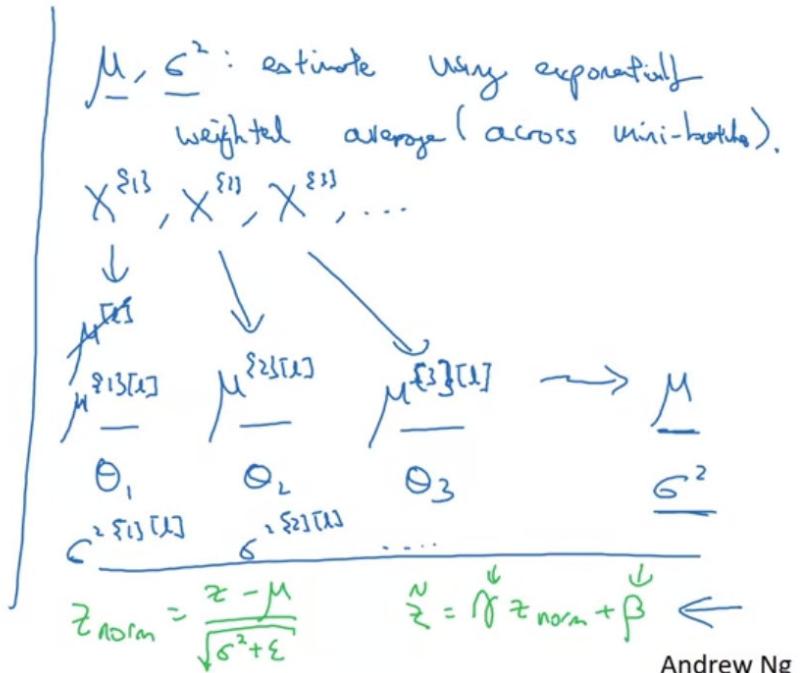
## Batch Norm at test time

$$\rightarrow \boxed{\mu = \frac{1}{m} \sum_i z^{(i)}}$$

$$\rightarrow \boxed{\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2}$$

$$\rightarrow \boxed{z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}}$$

$$\rightarrow \boxed{\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta}$$



## Softmax regression

### Activation function

$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{i=1}^C e^{z^{[l]}}}$$

where  $C$  is the number of classes.

### Loss function

$$L(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

### Back propagation

$$dZ^{[L]} = \hat{y} - y$$

## Structuring Machine Learning Projects

---

### Orthogonalization

Tune exactly one knob to tune the system.

### Satisficing and Optimizing Metric

**maximize** optimizing metric **s.t.** satisficing metric

N metrics: 1 optimizing, N - 1 satisficing.

### Guideline for choosing dev/test set

Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on. And the dev set and test set should come from **same distribution**.

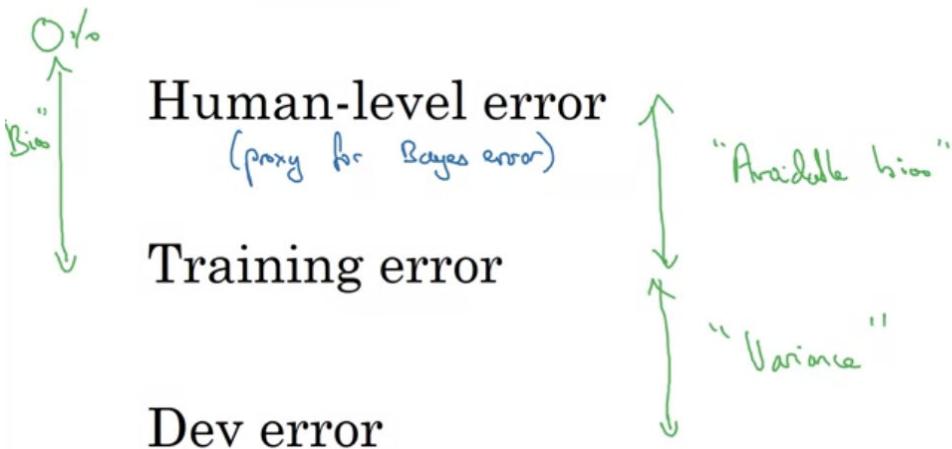
Set the test set to be **big enough to give high confidence** in the overall performance of the system.

### Comparing with human-level error

**Bayes error:** the best possible error any function could achieve. **Human-level error** is an **approximation** of Bayes error.

**Avoidable bias:** the difference between bayes error and training error.

# Summary of bias/variance with human-level performance



## Improving Model Performance

### Two fundamental assumptions of supervised learning

1. You can fit the training set pretty well. **Avoidable bias**
2. The training set performance generalizes pretty well to the dev/test set. **Variance**

#### Avoidable bias

- Train bigger model
- Train longer/better optimization algorithms (momentum, RMSProp, Adam...)
- NN architecture/hyperparameters search

#### Variance

- More data
- Regularization (l2, dropout, data augmentation)
- NN architecture/hyperparameters search

## Error analysis

Find a set of mislabeled examples in dev set, and look at the mislabeled examples for false positives and false negatives. Then count up the number of errors that fall into various different categories.

# Error analysis



| 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%                  |                                   |

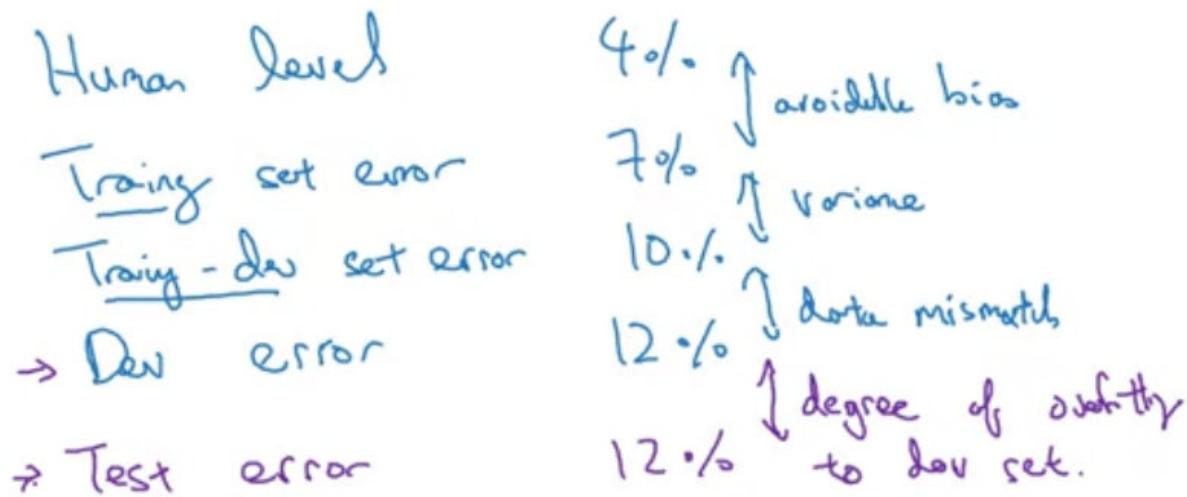
## Clearning Up Incorrectly Labeled Data

DL algoritms are quite **robust** to random errors in the **training set**.

### When correcting incorrect dev/test set examples

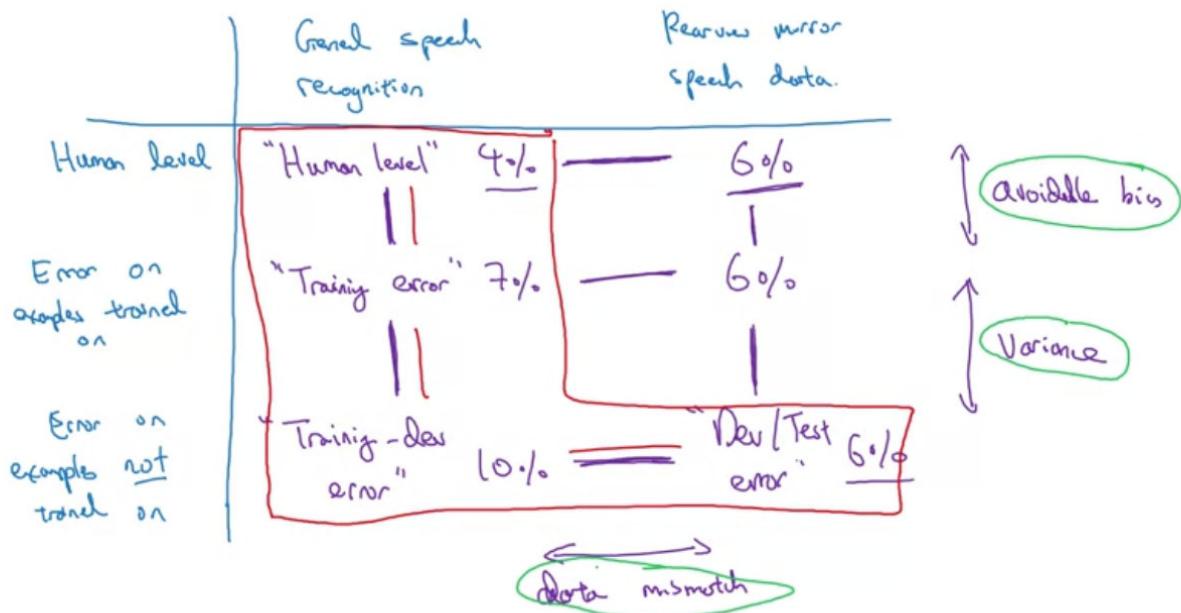
- Apply same process to dev and test sets to make sure they continue to come from the same distribution.
- Consider examining examples the algorithm got right as well as ones it got wrong.
- Train and dev/test data may now come from slightly different distributions.

## Bias/varience on mismatched training and dev/test sets



## More general formulation

Reactive mirror



Andrew Ng

## Addressing Data Mismatch

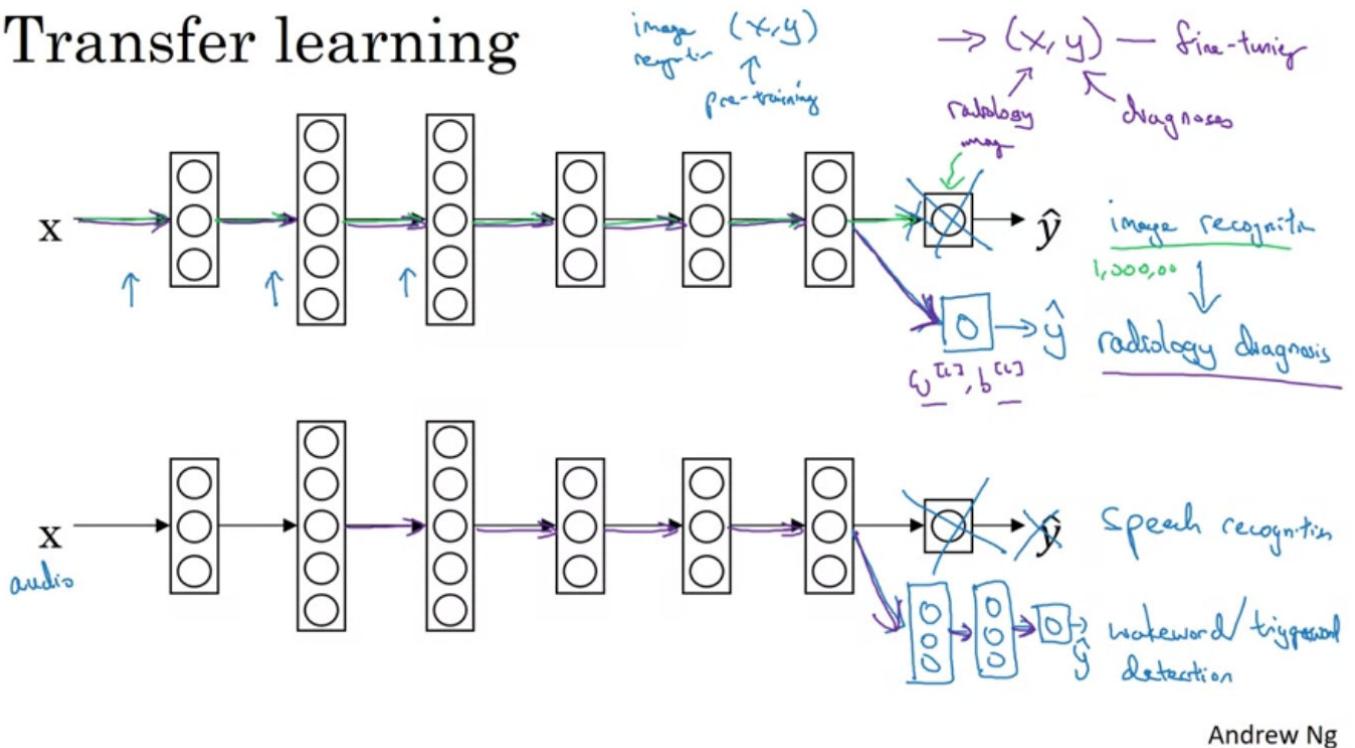
- Carry out **manual error analysis** to try to understand difference between training and dev/test sets.
    - e.g. noisy
  - Make training data more similar; or collect more data similar to dev/test sets.
    - Artificial data synthesis

# Transfer Learning

Transfer learning makes sense when you have **a lot of** data for the problem you are **transferring from**, and relatively **less** data for the problem you are **transferring to**. To be more specific,

- Task A and B have the **same input x**.
  - You have a lot more data for A than B
  - Low level features from A could be helpful for B.

# Transfer learning



Andrew Ng

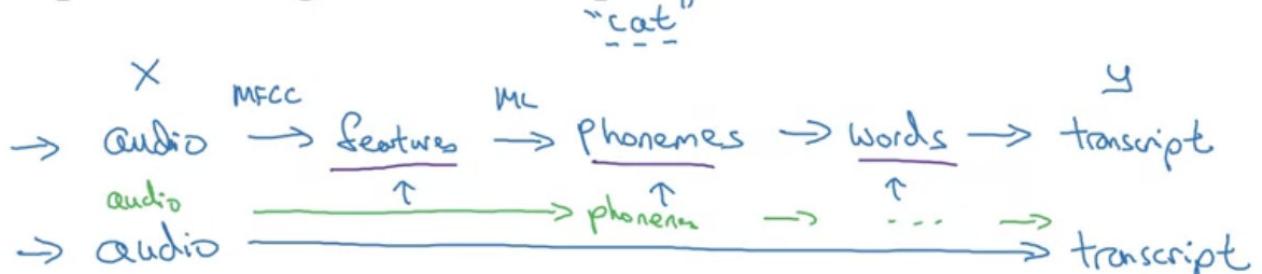
## Multi-task learning

- 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.

## End-to-end Deep Learning

### What is end-to-end learning?

#### Speech recognition example



3,000h  
↑

10,000h  
⋮  
100,000h

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

### Pros

- Let the data speak
- Less hand-designing of components needed

### Cons

- May need large amount of data
- Excludes potentially useful hand-designed components

## Convolutional Neural Networks

### Computer Vision

Image Classification, Object Detection, Neural Style Transfer...

The problem is that, when we use a large image, the input size will be extremely large. So, we need to implement convolution operation.

### Edge Detection Example

#### 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

"convolution"

\*

|   |   |    |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |

3x3  
filter

=

|     |    |    |     |
|-----|----|----|-----|
| -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

# Vertical edge detection examples

|    |    |    |   |   |   |
|----|----|----|---|---|---|
| 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 |



|   |   |   |    |    |    |
|---|---|---|----|----|----|
| 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 | 10 | 10 | 10 |



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

=

|   |    |    |   |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

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

=

|   |     |     |   |
|---|-----|-----|---|
| 0 | -30 | -30 | 0 |
| 0 | -30 | -30 | 0 |
| 0 | -30 | -30 | 0 |
| 0 | -30 | -30 | 0 |

Andrew Ng

|   |   |   |            |            |            |
|---|---|---|------------|------------|------------|
| 1 | 1 | 1 | $\times 1$ | 0          | $\times 1$ |
| 0 | 1 | 1 | $\times 0$ | $\times 1$ | $\times 0$ |
| 0 | 0 | 1 | $\times 1$ | $\times 0$ | $\times 1$ |
| 0 | 0 | 1 | 1          | 1          | 0          |
| 0 | 1 | 1 | 0          | 0          | 0          |

Image

|   |   |   |
|---|---|---|
| 4 | 3 | 4 |
|   |   |   |
|   |   |   |
|   |   |   |

Convolved Feature

## Padding

A  $n \times n$  image convolved by a  $f \times f$  filter, the dimension of output will be  $n - f + 1 \times n - f + 1$ .

There are two major drawbacks:

- The image shrinks every time when applying convolutional operator.
- The pixels on the corners or on the edges are used much less in the output.

In order to fix both of these problems: **Pad** the image with  $p$  (padding amount) before applying convolutional operation. Then the output size becomes  $n + 2p - f + 1 \times n + 2p - f + 1$ .

## Valid and Same convolutions

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

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

$$n + 2p - f + 1 = n$$

$$p = \frac{f - 1}{2}$$

$f$  is usually (almost always) **odd**.

## Strided convolution

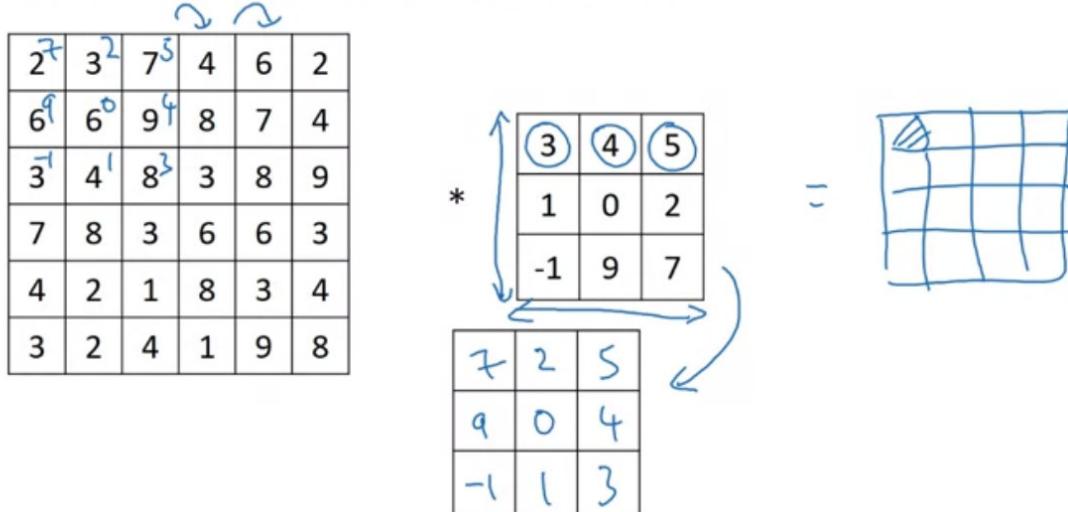
$n \times n$  image,  $f \times f$  filter, padding  $p$ , stride  $s$ :

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

## Technical note on cross-correlation vs. convolution

# Technical note on cross-correlation vs. convolution

## Convolution in math textbook:



Andrew Ng

## Convolutions over volumes

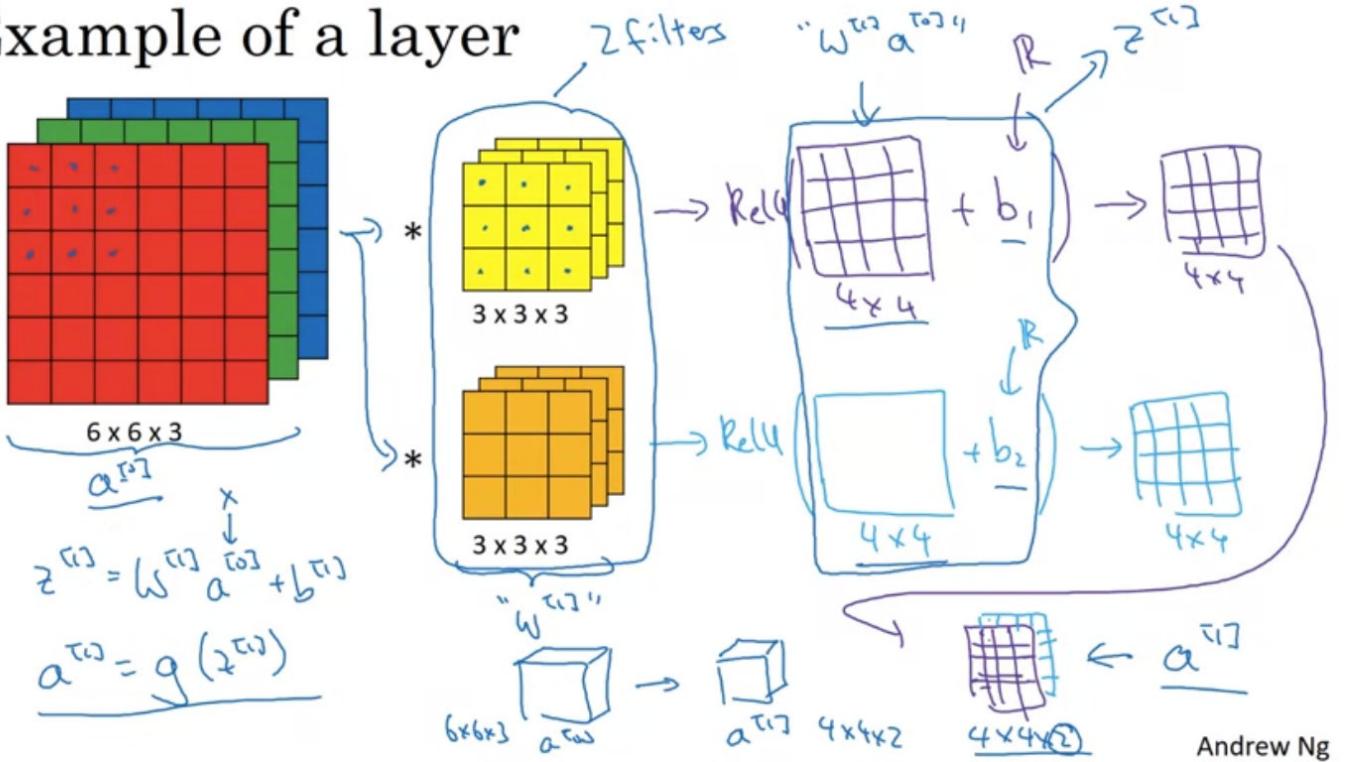
The number of channels of the image is equal to the number of channels of filter. If we use multiple filters, we can stack the outputs together to get the final output.

## Summary (stride = 1)

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

## One Layer of a Convolutional Network

### 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

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

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

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

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

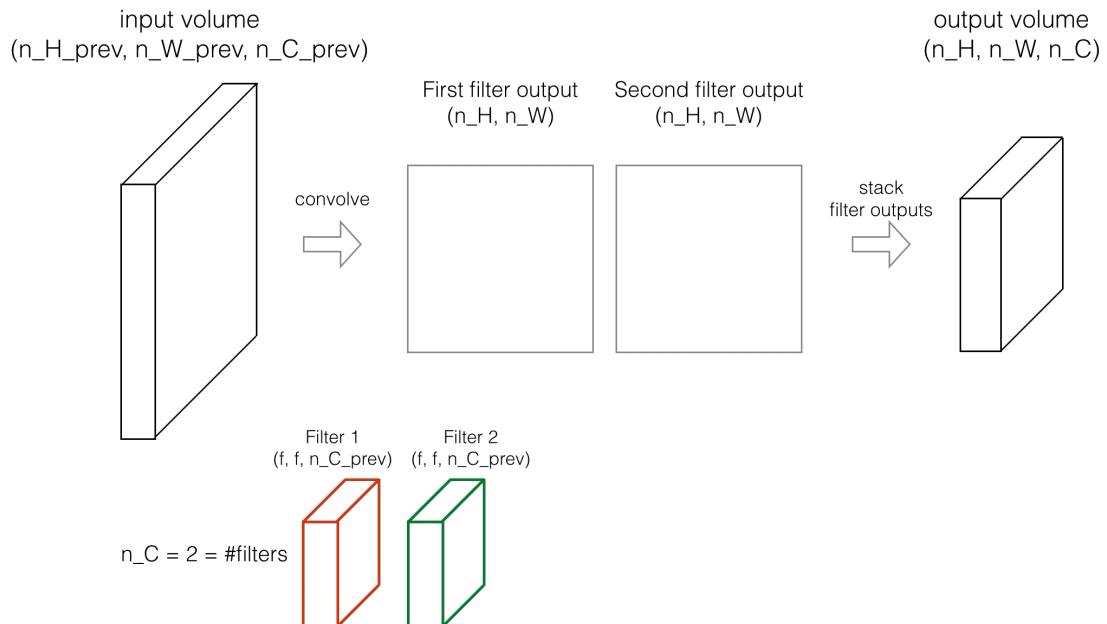
Each filter is:  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

Activations:  $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$ ,  $A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

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

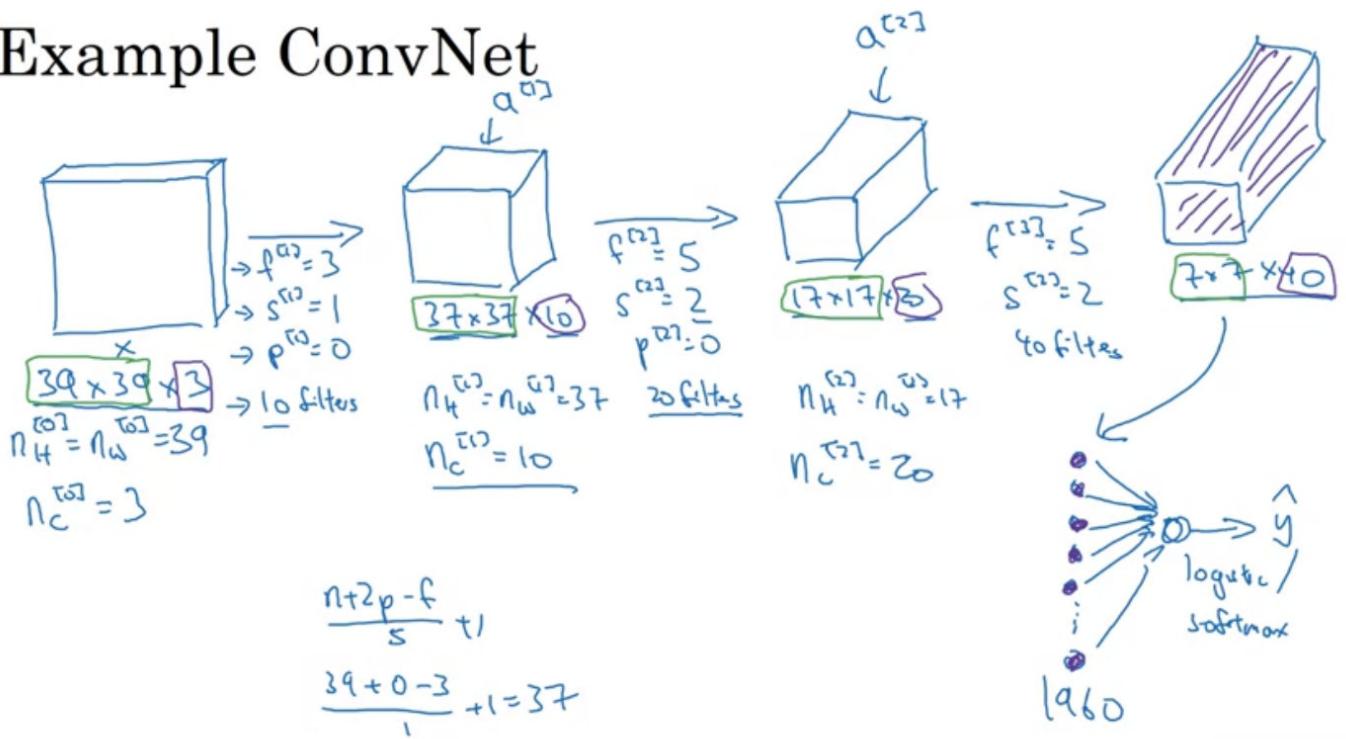
Bias:  $n_C^{[l]} = (1, 1, 1, n_C^{[l]})$

## How do convolutions work?



## Simple ConvNet Example

### Example ConvNet



Andrew Ng

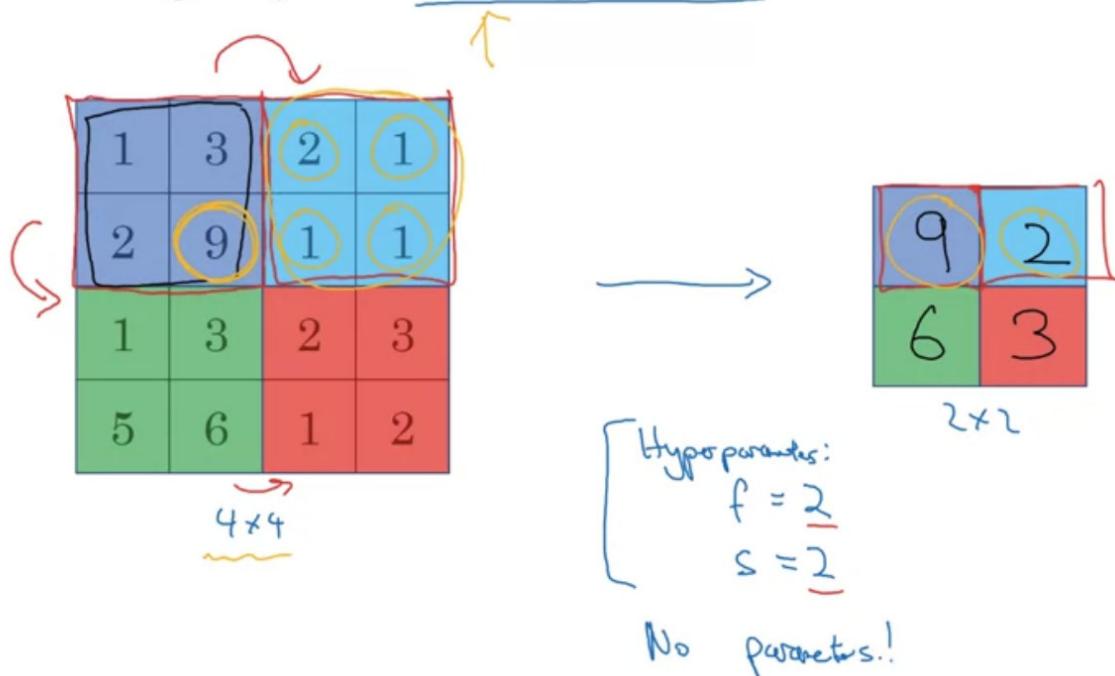
## Types of layer in a convolutional networks:

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

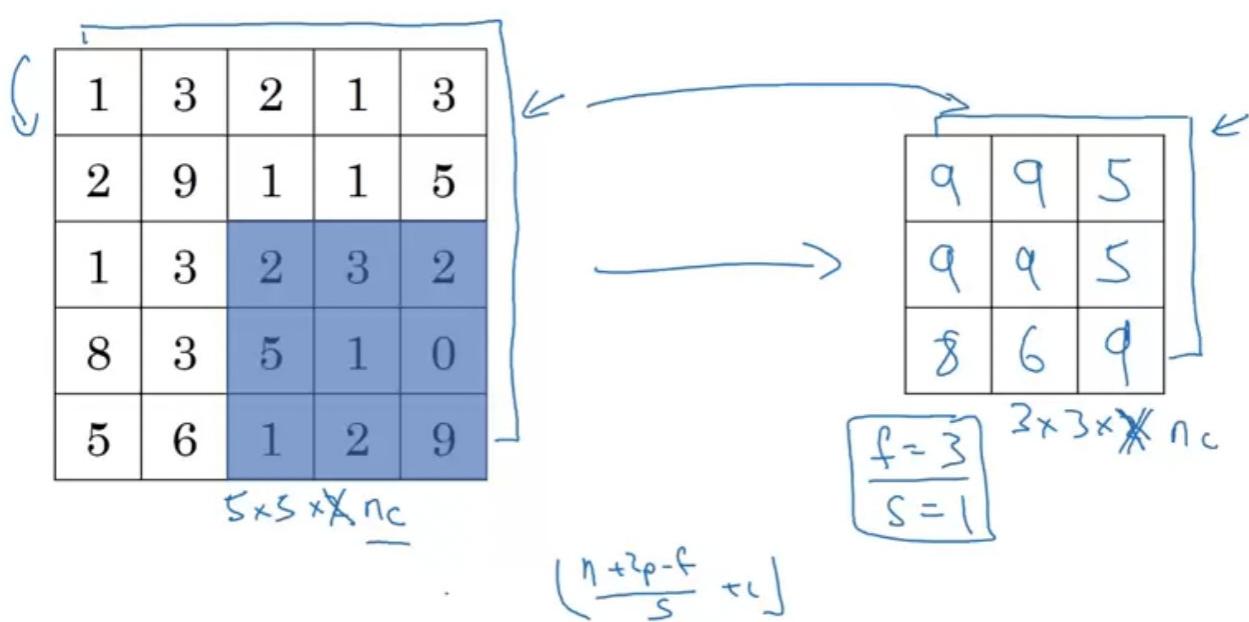
## Pooling Layers

### Max pooling

## Pooling layer: Max pooling



## Pooling layer: Max pooling



The max pooling computation is done **independently** on each of  $n_C$  channels.

## Average pooling (not often used)

Instead of finding maximum, find the average of the "block".

In very deep neural networks, you might use average pooling to collapse the representations.

## Summary of pooling

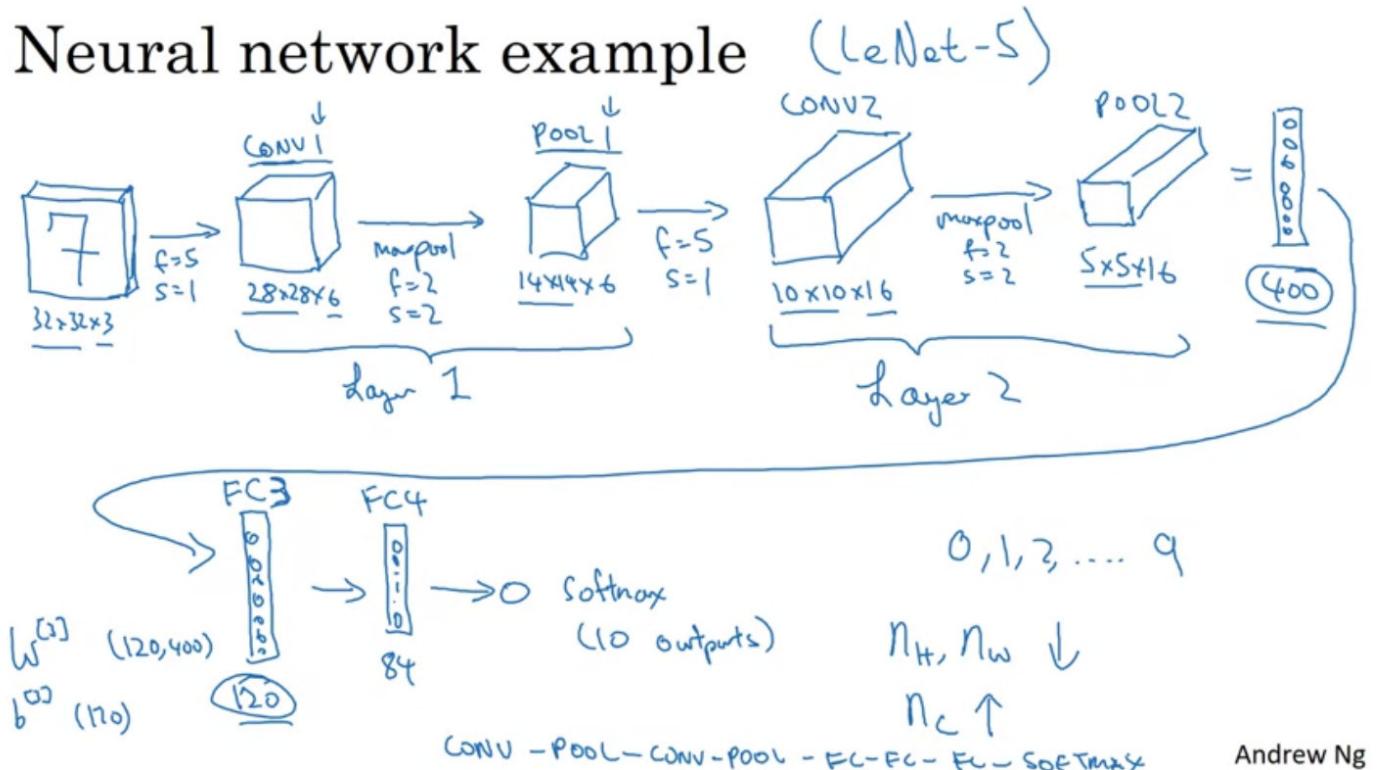
$$n_H \times n_W \times n_C \rightarrow \left\lfloor \frac{n_W^{[l-1]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times \left\lfloor \frac{n_W^{[l-1]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor \times n_C$$

Hyperparameters: normally  $f = 2, s = 2$  (this will half the height and width)

There is **no parameters to learn!**

In fact, pooling layers modify the input by choosing one value out of several values in their input volume. Also, to compute derivatives for the layers that have parameters (Convolutions, Fully-Connected), we **still need to backpropagate the gradient through the Pooling layers**.

## CNN Examples (Similar to LeNet-5)



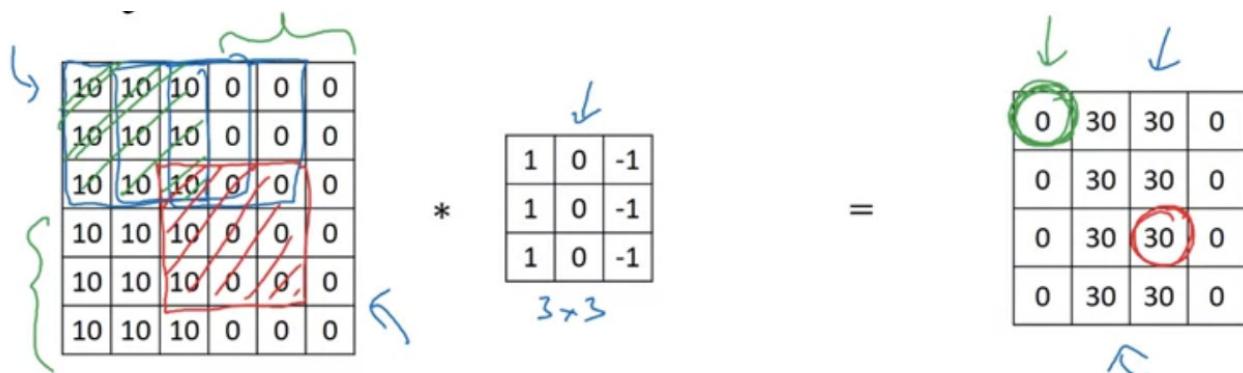
# Neural network example

|                  | Activation shape | Activation Size   | # parameters |
|------------------|------------------|-------------------|--------------|
| Input:           | (32,32,3)        | — 3,072 $a^{(1)}$ | 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          |

## Why Convolutions?

The reasons why ConvNet needs much less parameter):

- **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.



Through these two mechanisms, a neural network has a lot fewer parameters which allows it to be trained with smaller training set and is less prone to be overfitting.

## What you should remember:

- A convolution extracts features from an input image by taking the dot product between the input data and a 2D array of weights (the filter).
- The 2D output of the convolution is called the feature map
- A convolution layer is where the filter slides over the image and computes the dot product
  - This transforms the input volume into an output volume of different size
- Zero padding helps keep more information at the image borders, and is helpful for building deeper networks, because you can build a CONV layer without shrinking the height and width of the volumes

- Pooling layers gradually reduce the height and width of the input by sliding a 2D window over each specified region, then summarizing the features in that region

## Terminology

Window, kernel, filter, pool

The words **kernel** and **filter** are used to refer to the same thing. The word **filter** accounts for the amount of **kernels** that will be used in a single convolution layer. **Pool** is the name of the operation that takes the max or average value of the kernels.

## Convolutional Layer Backward Pass

**Computing dA:**

This is the formula for computing  $dA$  with respect to the cost for a certain filter  $W_c$  and a given training example:

$$dA += \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} W_c \times dZ_{hw} \quad (1)$$

Where  $W_c$  is a filter and  $dZ_{hw}$  is a scalar corresponding to the gradient of the cost with respect to the output of the conv layer Z at the hth row and wth column (corresponding to the dot product taken at the ith stride left and jth stride down). Note that at each time, you multiply the same filter  $W_c$  by a different  $dZ$  when updating  $dA$ . We do so mainly because when computing the forward propagation, each filter is dotted and summed by a different  $a_{slice}$ . Therefore when computing the backprop for  $dA$ , you are just adding the gradients of all the  $a_{slices}$ .

**Computing dW:**

This is the formula for computing  $dW_c$  ( $dW_c$  is the derivative of one filter) with respect to the loss:

$$dW_c += \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} a_{slice} \times dZ_{hw} \quad (2)$$

Where  $a_{slice}$  corresponds to the slice which was used to generate the activation  $Z_{ij}$ . Hence, this ends up giving us the gradient for  $W$  with respect to that slice. Since it is the same  $W$ , we will just add up all such gradients to get  $dW$ .

**Computing db:**

This is the formula for computing  $db$  with respect to the cost for a certain filter  $W_c$ :

$$db = \sum_h \sum_w dZ_{hw} \quad (3)$$

As you have previously seen in basic neural networks,  $db$  is computed by summing  $dZ$ . In this case, you are just summing over all the gradients of the conv output (Z) with respect to the cost.

## Pooling Layer - Backward Pass

Next, let's implement the backward pass for the pooling layer, starting with the MAX-POOL layer. Even though a pooling layer has no parameters for backprop to update, you still need to backpropagate the gradient through the pooling layer in order to compute gradients for layers that came before the pooling layer.

### Max Pooling - Backward Pass

Before jumping into the backpropagation of the pooling layer, you are going to build a helper function called `create_mask_from_window()` which does the following:

$$X = \begin{bmatrix} 1 & 3 \\ 4 & 2 \end{bmatrix} \rightarrow M = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \quad (4)$$

As you can see, this function creates a "mask" matrix which keeps track of where the maximum of the matrix is. True (1) indicates the position of the maximum in X, the other entries are False (0). You'll see later that the backward pass for average pooling is similar to this, but uses a different mask.

**Why keep track of the position of the max?** It's because this is the input value that ultimately influenced the output, and therefore the cost. Backprop is computing gradients with respect to the cost, so anything that influences the ultimate cost should have a non-zero gradient. So, backprop will "propagate" the gradient back to this particular input value that had influenced the cost.

### Average Pooling - Backward Pass

In max pooling, for each input window, all the "influence" on the output came from a single input value--the max. In average pooling, every element of the input window has equal influence on the output. So to implement backprop, you will now implement a helper function that reflects this.

For example if we did average pooling in the forward pass using a 2x2 filter, then the mask you'll use for the backward pass will look like:

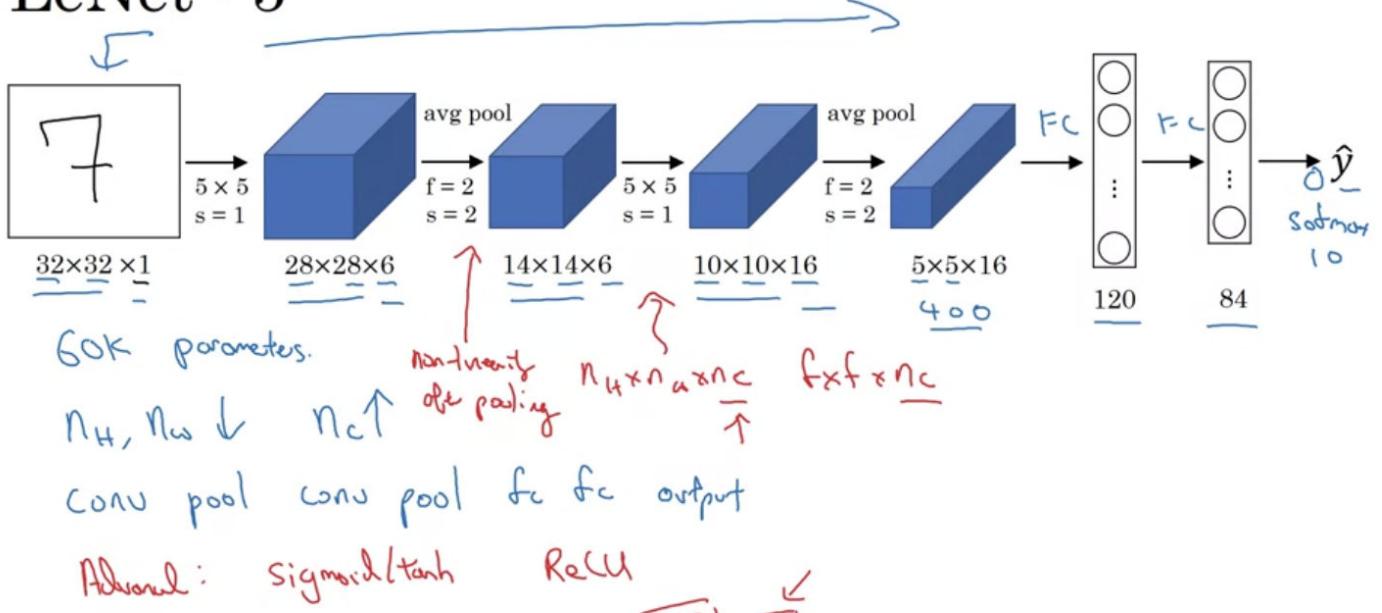
$$dZ = 1 \rightarrow dZ = \begin{bmatrix} 1/4 & 1/4 \\ 1/4 & 1/4 \end{bmatrix} \quad (5)$$

This implies that each position in the  $dZ$  matrix contributes equally to output because in the forward pass, we took an average.

## Case studies

### LeNet - 5

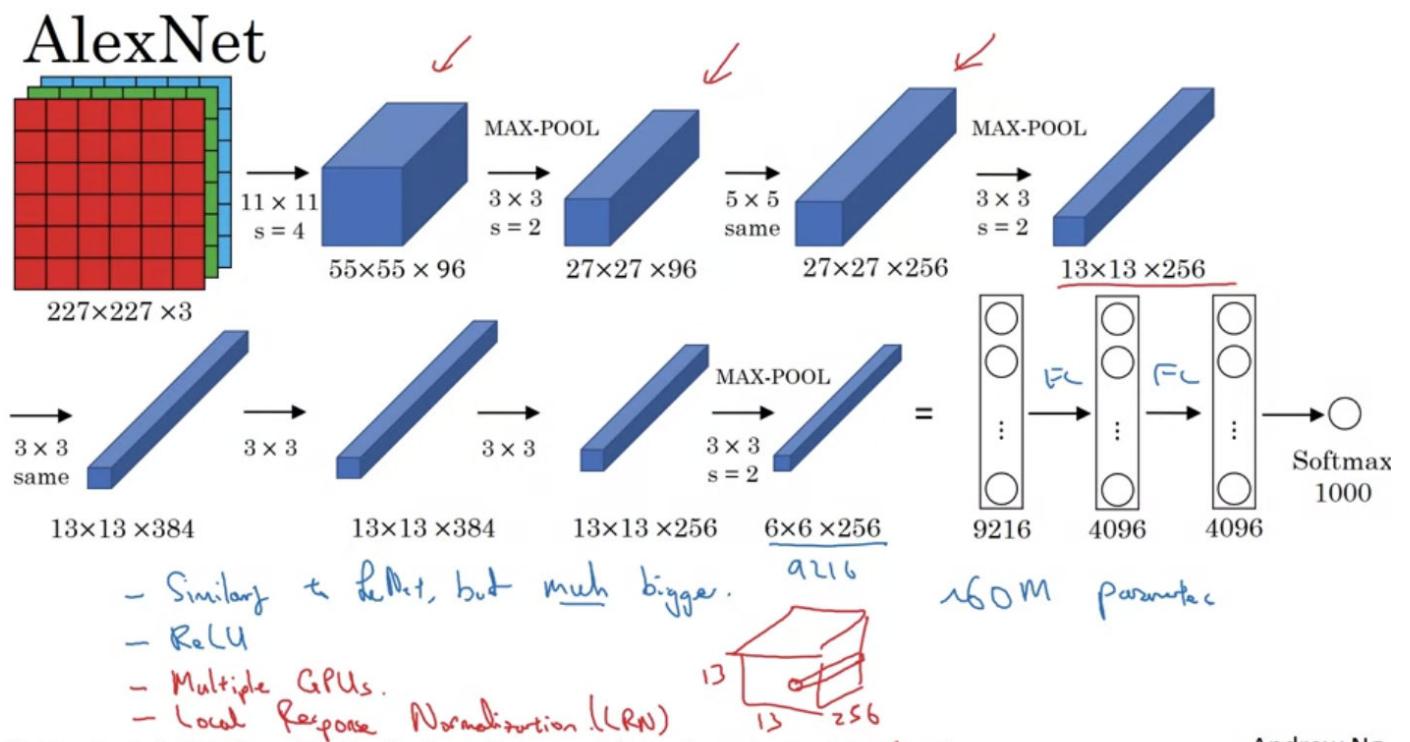
# LeNet - 5



[LeCun et al., 1998. Gradient-based learning applied to document recognition]

Andrew Ng

## AlexNet



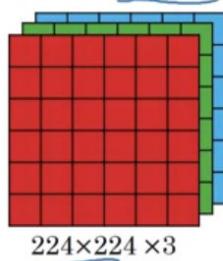
[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

Andrew Ng

## VGG - 16 (16 refers to there are 16 layers with weights)

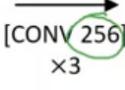
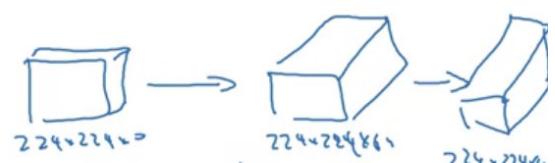
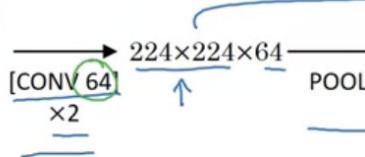
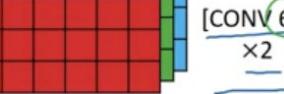
# VGG - 16

CONV = 3x3 filter, s = 1, same

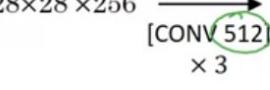


# VGG-19

MAX-POOL = 2x2, s = 2

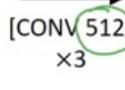


POOL



POOL

14x14x512



POOL

$n_h, n_w \downarrow$

$n_c \uparrow$

$\sim 138M$

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

Andrew Ng

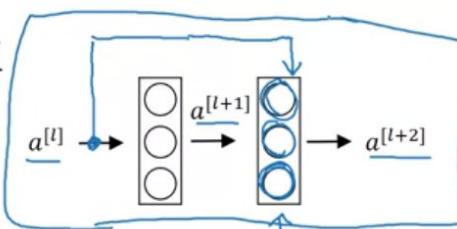
## Residual Networks (ResNets)

Very deep neural networks are difficult to train because of **vanishing** and **exploding** gradient types of problems. The **skip connections** allows you to take activation from one layer and suddenly feed it to another layer even much deeper in the neural network. Using that, you will build **ResNet** which enables you to train very deep networks.

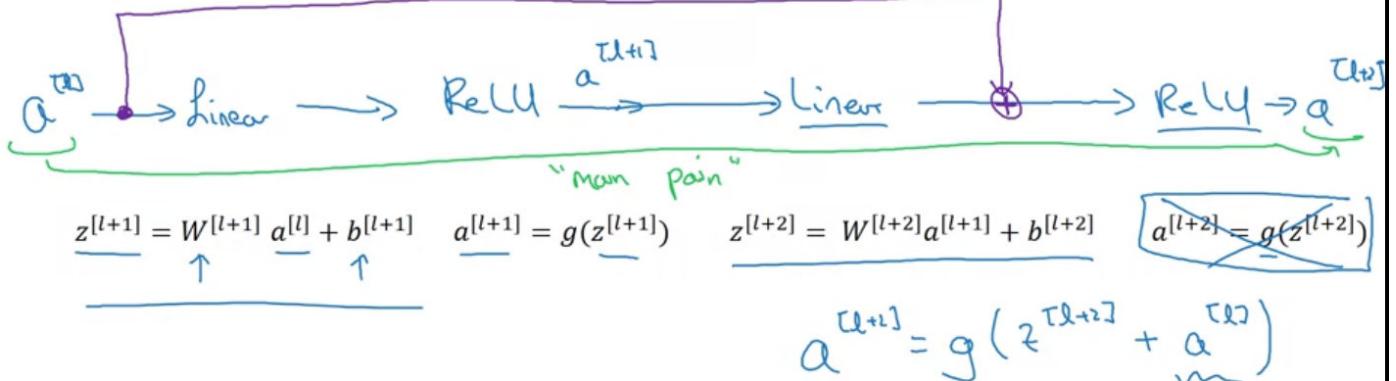
### Residual block (Identity block)

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

## Residual block



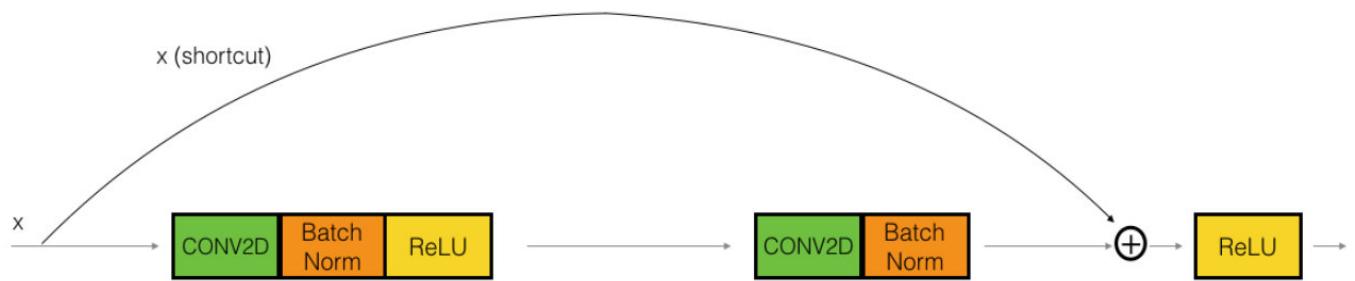
"short cut" / skip connection



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

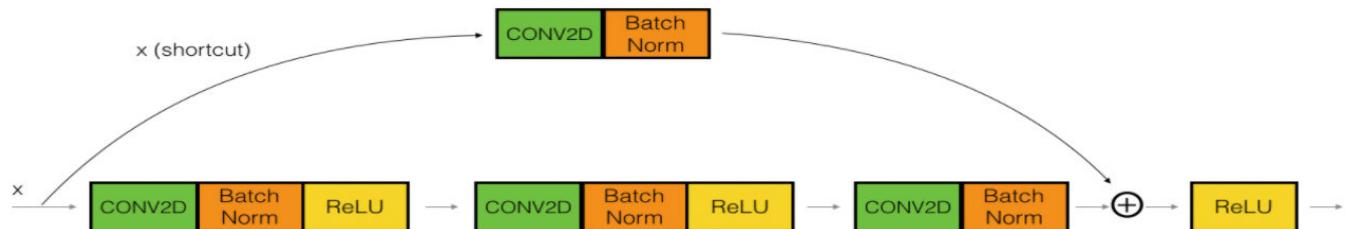
Andrew Ng

The **identity block** is the standard block used in ResNets, and corresponds to the case where the input activation (say  $a^{[l]}$ ) has the **same dimension** as the output activation (say  $a^{[l+2]}$ ).



**Figure 3 : Identity block.** Skip connection "skips over" 2 layers.

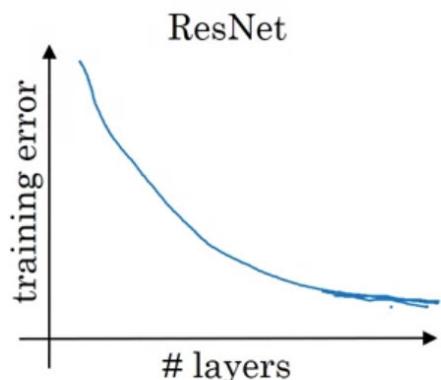
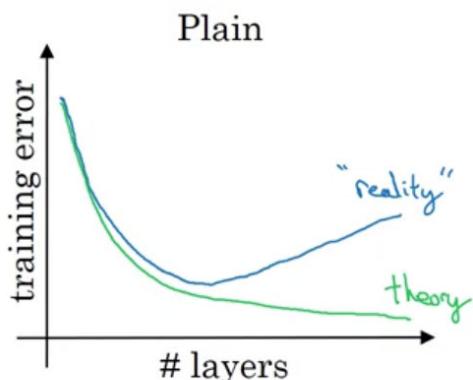
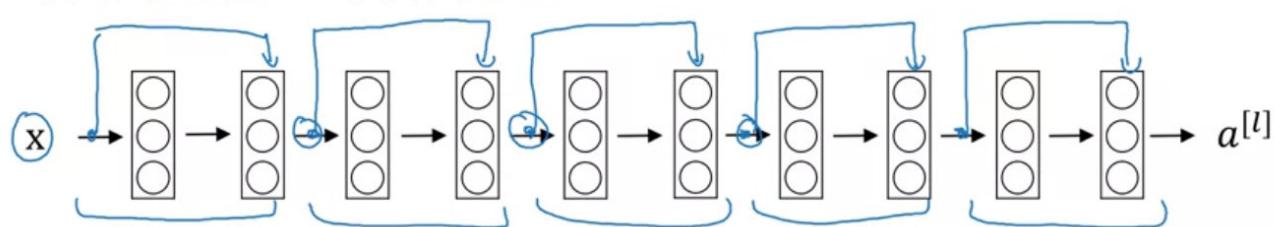
The ResNet "**convolutional block**" is the second block type. You can use this type of block when the input and output dimensions **don't match up**. The difference with the identity block is that there is a CONV2D layer in the shortcut path:



**Figure 4 : Convolutional block**

In theory, having a deeper network should only help. But in practice, having a plain network (no ResNet) is very deep means that your optimization algorithm just has a much harder time training. And so, in reality, your training error gets **worse** if you pick a network that is **too deep**.

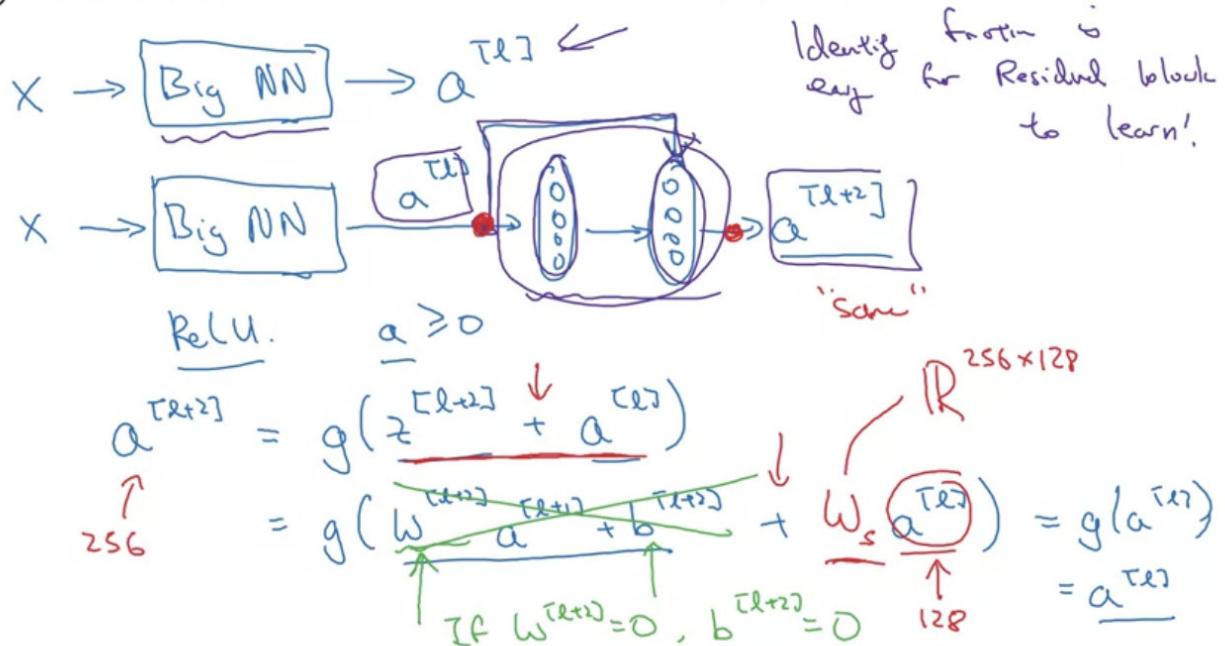
## Residual Network



## Why ResNet works?

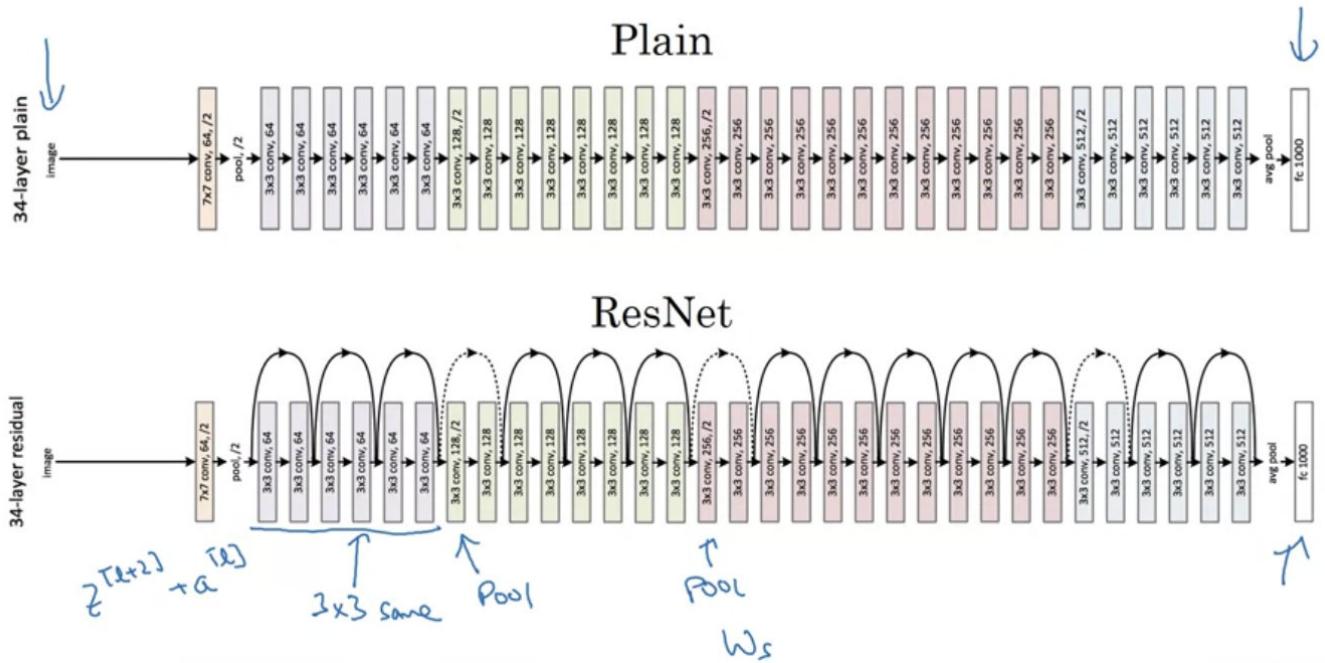
The main reason is that it is so easy for these extra layers to learn the **identity function** that you are kind of guaranteed that **it doesn't hurt performance** and then a lot the time you may be get lucky and then even **helps performance**.

# Why do residual networks work?



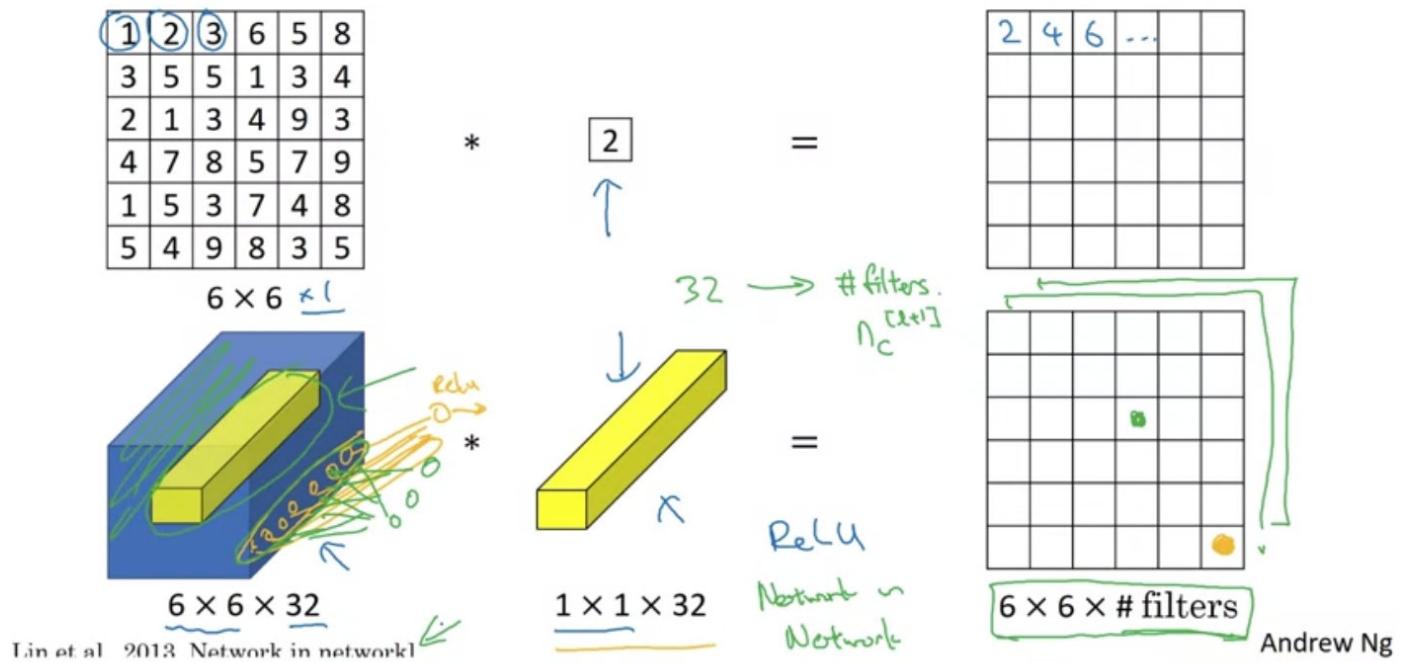
Andrew Ng

ResNet



## Network in Network and $1 \times 1$ convolutions

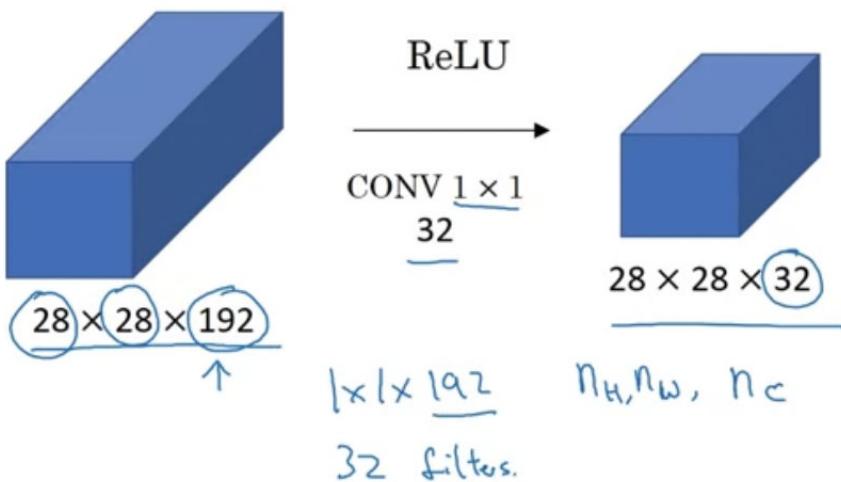
### Why does a $1 \times 1$ convolution do?



The  $1 \times 1$  convolution has nonlinearity, it allows you to learn a more complex function of your network by adding another layer.

The  $1 \times 1$  convolution can **shrink the number of channels**.

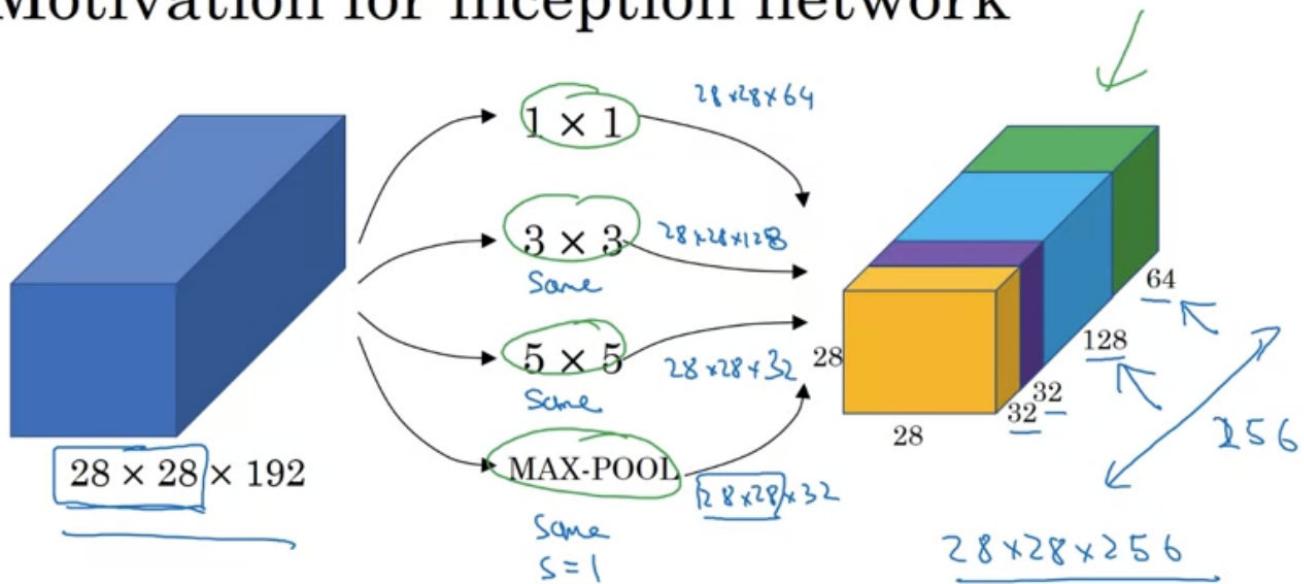
### Using $1 \times 1$ convolutions



## Inception Network

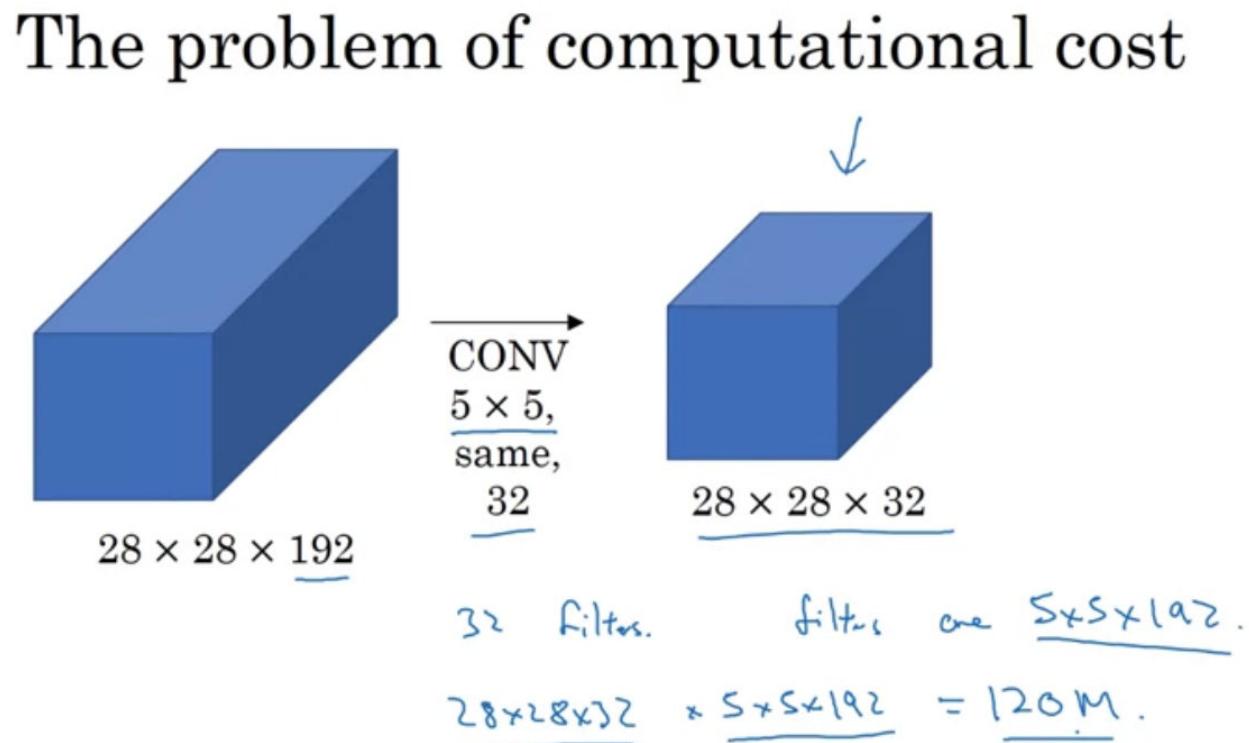
### Motivation for inception network

## Motivation for inception network



The basic idea is that instead of you needing to pick one of these filter sizes or pooling you want and committing to that, you can do them all and just concatenate all the outputs, and let the network learn whatever parameters it wants to use, whatever the combinations of these filter sizes it wants.

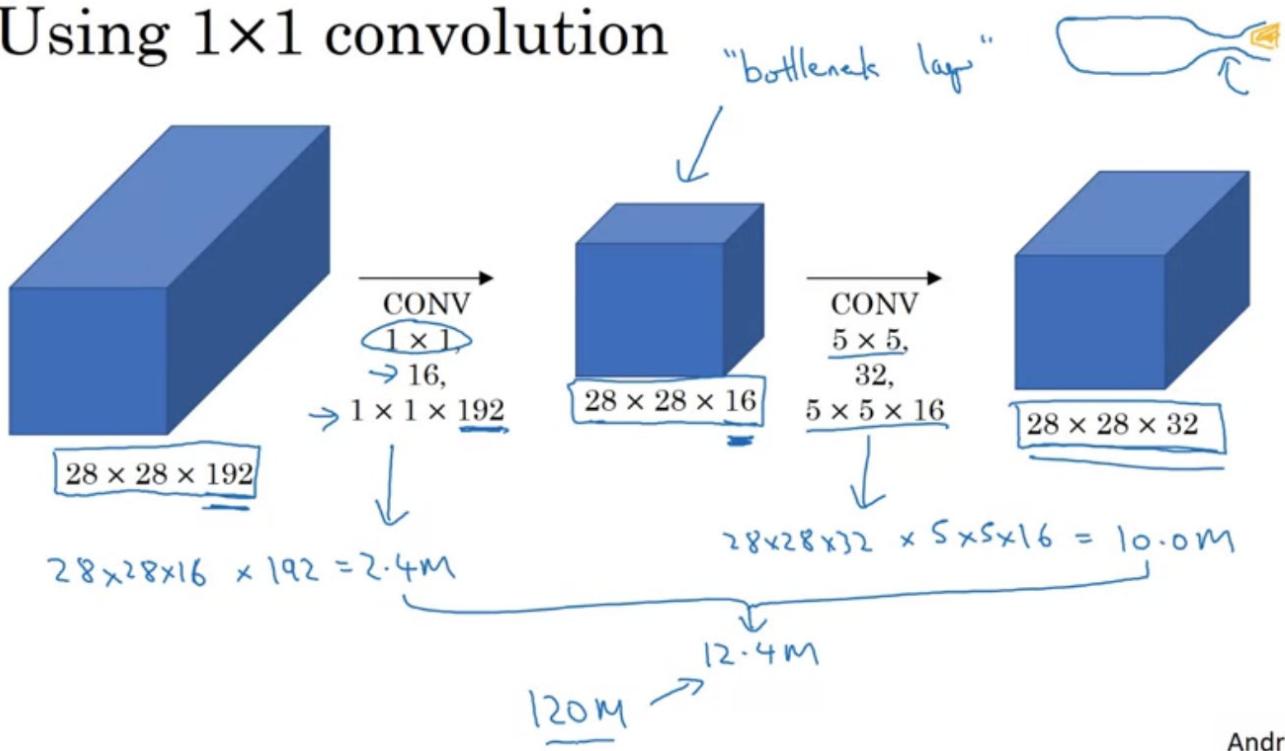
### The problem of computational cost



#multiplies = #multiplies to compute each of output values  $\times$  #output values

## Using $1 \times 1$ convolution

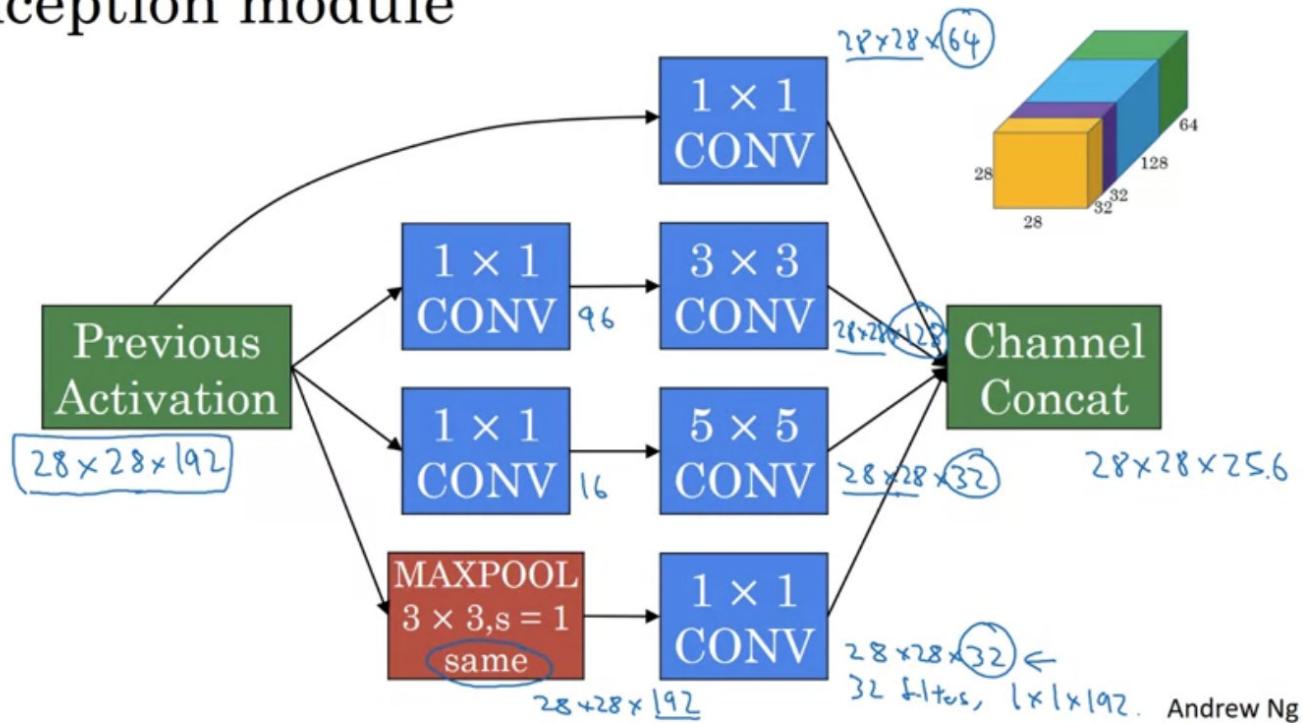
# Using $1 \times 1$ convolution



Andrew Ng

## Inception module

# Inception module



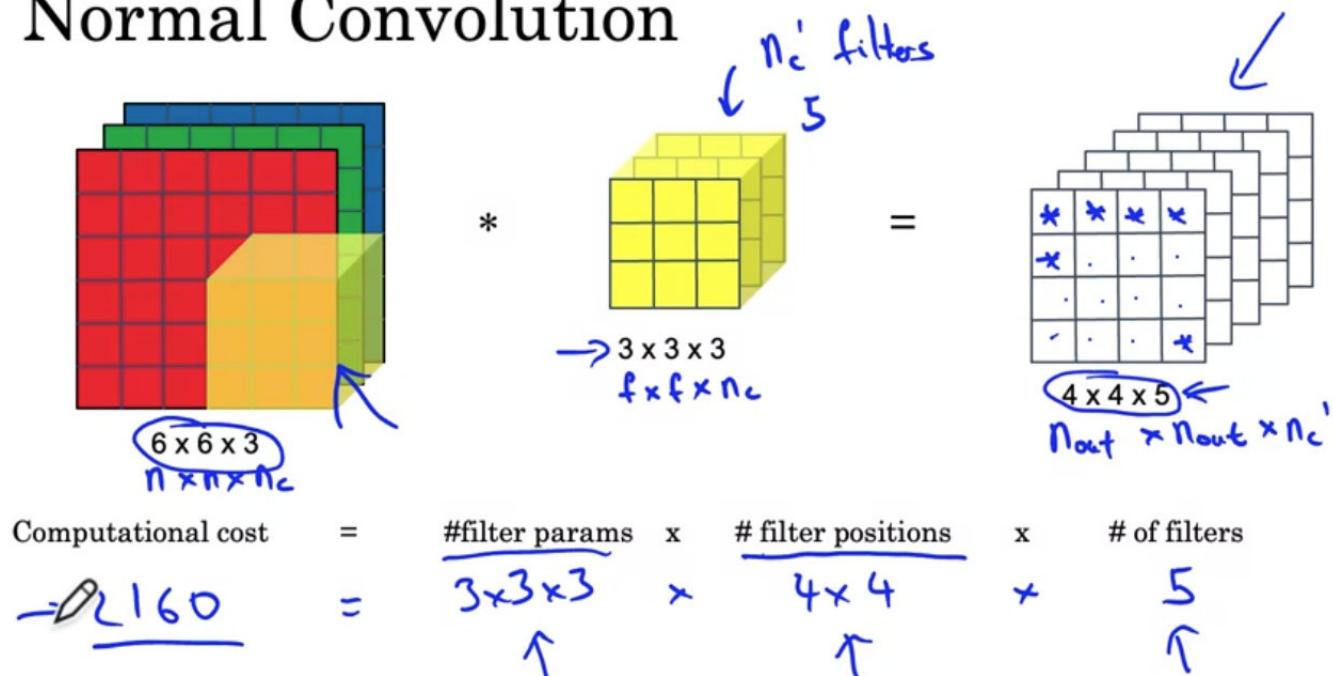
## MobileNet

### Motivation for MobileNets

- Low computational cost at deployment
- Useful for mobile and embedded vision applications

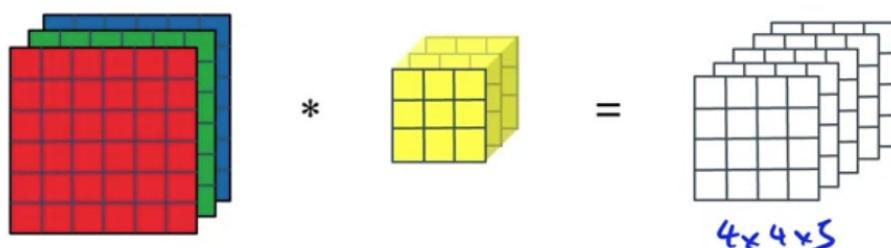
Key idea: Normal vs. depthwise-separable convolutions

## Normal Convolution

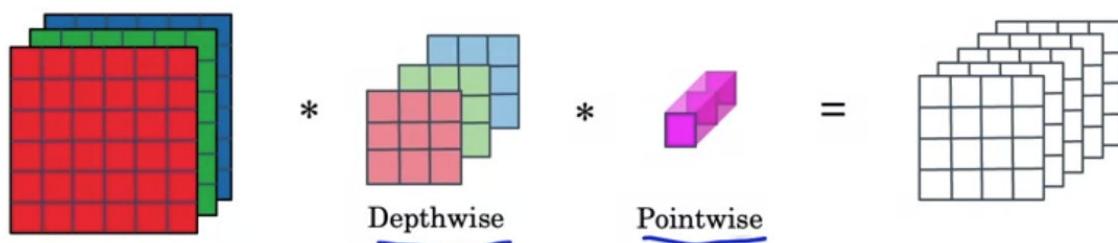


## Depthwise Separable Convolution

Normal Convolution



Depthwise Separable Convolution



Cost of normal convolution: 2160

Cost of depth wise separable convolution:  $432 + 240 = 672$

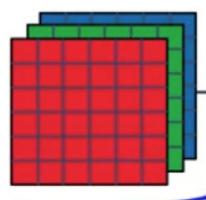
The ratio is equal to

$$ratio = \frac{1}{n'_C} + \frac{1}{f^2}$$

## MobileNet architecture

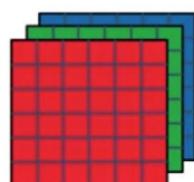
# MobileNet

MobileNet v1

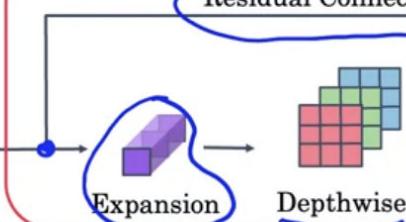


POOL, FC, SOFTMAX

MobileNet v2



Residual Connection



17 times

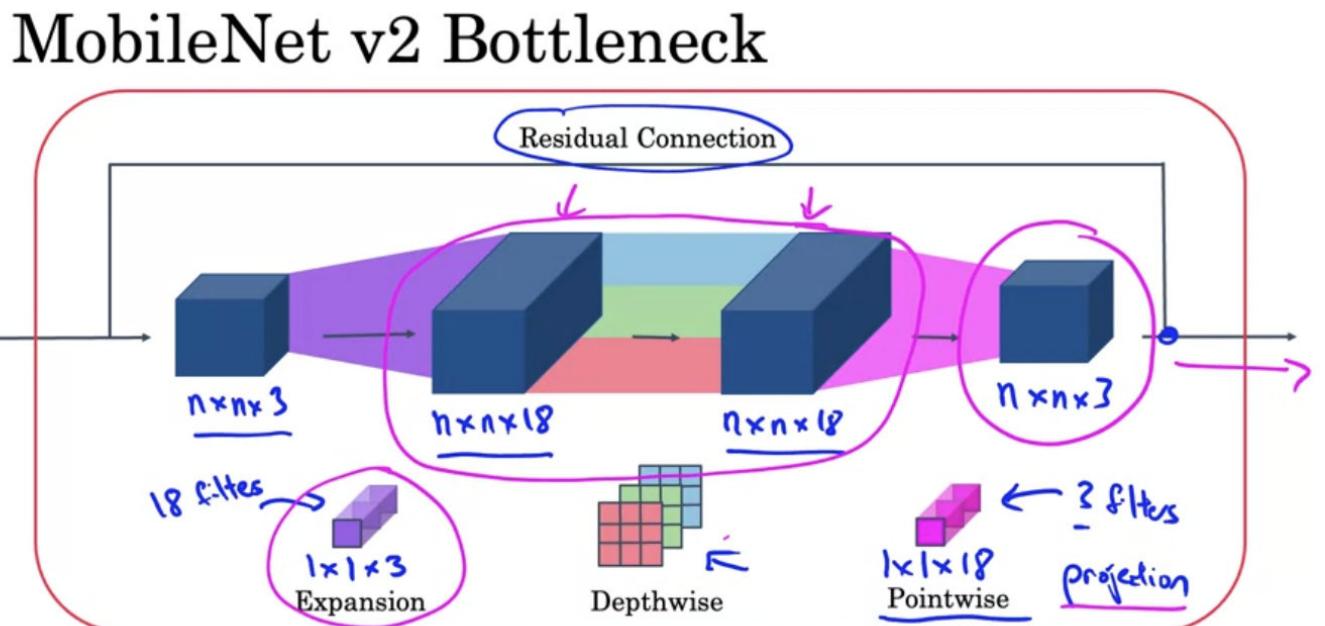
POOL, FC,  
SOFTMAX

Bottleneck

[Sandler et al. 2019, MobileNetV2: Inverted Residuals and Linear Bottlenecks]

Andrew Ng

## MobileNet v2 Bottleneck



## What you should remember:

- MobileNetV2's unique features are:
  - Depthwise separable convolutions that provide lightweight feature filtering and creation

- Input and output bottlenecks that preserve important information on either end of the block
- Depthwise separable convolutions deal with both spatial and depth (number of channels) dimensions

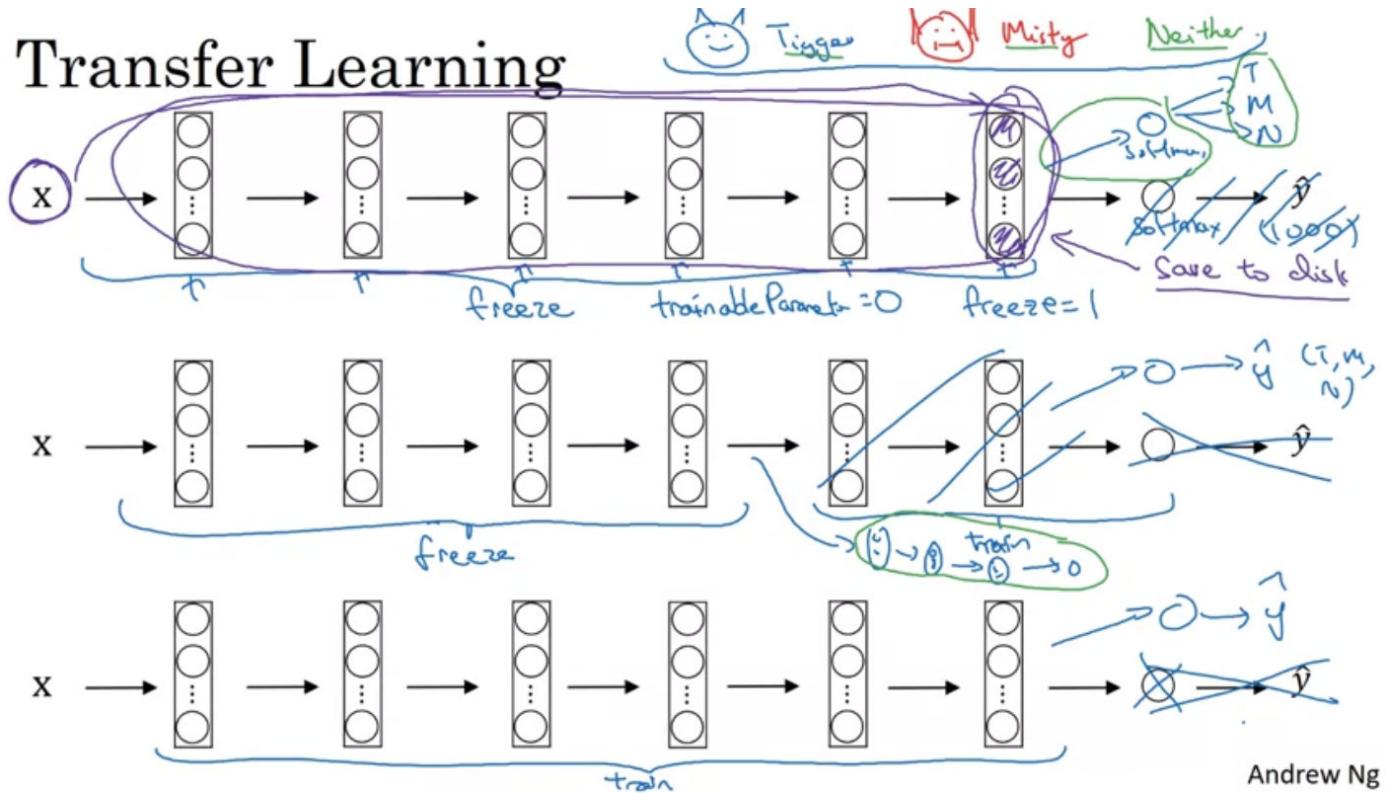
## EfficientNet

Three things you can do to scale things up or down are:

- $r$ : resolution (of images)
- $d$ : depth (of networks)
- $w$ : width (of layers)

The EfficientNet can help you find a way to scale up or down networks.

## Transfer Learning



You could try fine-tuning the model by re-running the optimizer in the last layers to improve accuracy. When you use a smaller learning rate, you take smaller steps to adapt it a little more closely to the new data. In transfer learning, the way you achieve this is by unfreezing the layers at the end of the network, and then re-training your model on the final layers with **a very low learning rate**. Adapting your learning rate to go over these layers in smaller steps can yield more fine details - and higher accuracy.

### What you should remember:

- To adapt the classifier to new data: Delete the top layer, add a new classification layer, and train only on that layer
- When freezing layers, avoid keeping track of statistics (like in the batch normalization layer)
- Fine-tune the final layers of your model to capture high-level details near the end of the network and potentially improve accuracy

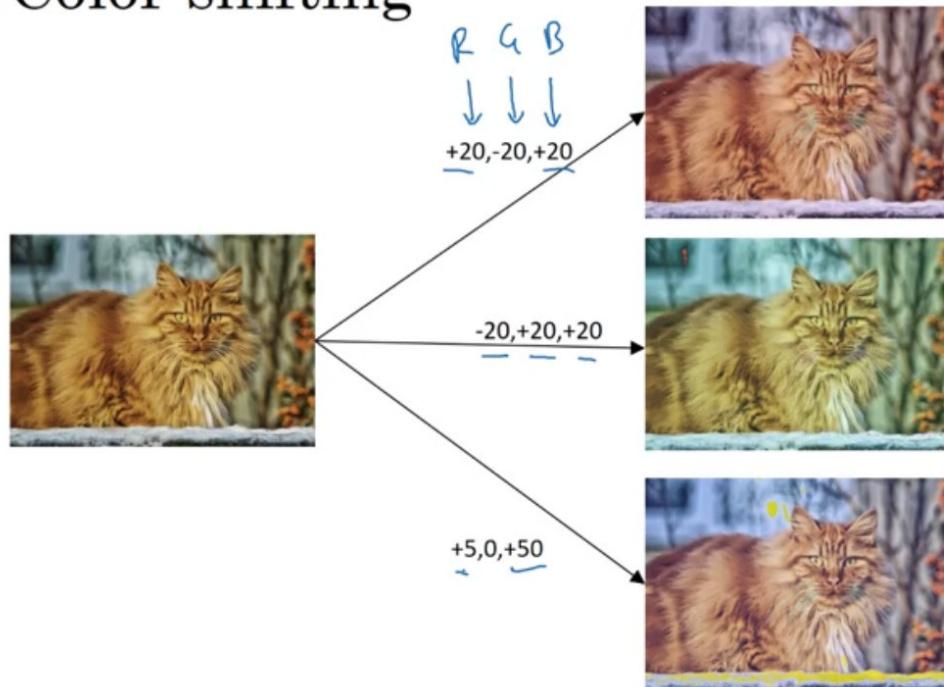
## Data Augmentation

### Common augmentation method

- Mirroring
- Random cropping
- Rotation
- Shearing
- Local warping
- ...

### Color shifting

## Color shifting

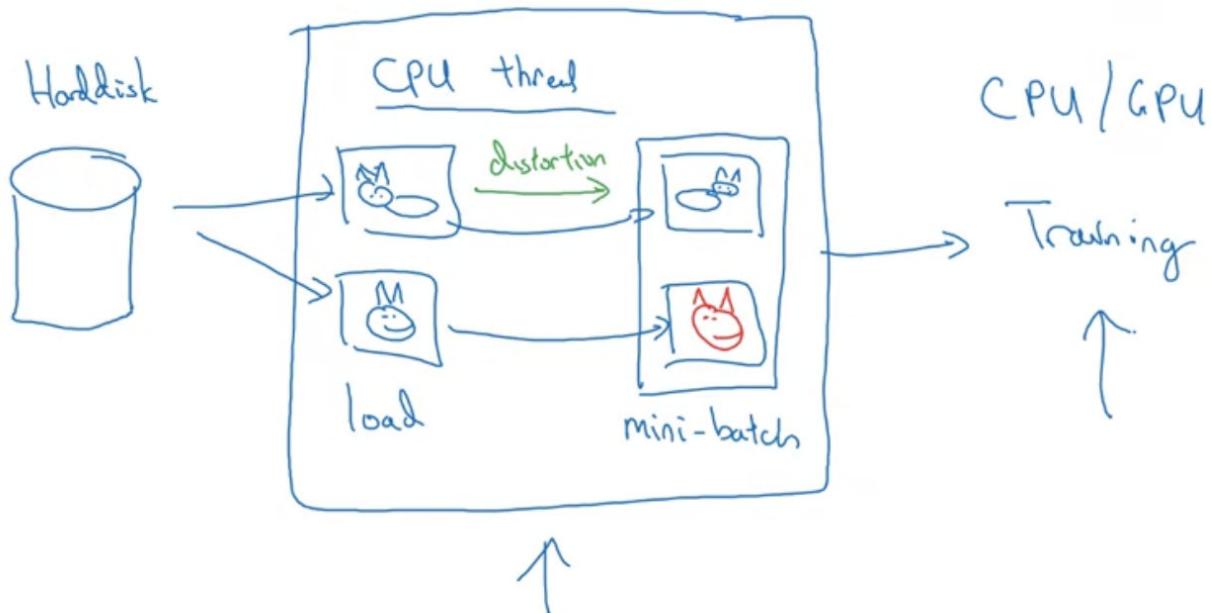


Andrew Ng

This make the learning algorithm more robust to changes in the colors of the images.

"PCA color augmentation" (AlexNet): keep the overall color of the tint the same.

# Implementing distortions during training



Andrew Ng

## Tips for doing well on benchmarks/winning competitions

### [Ensembling

- Train several networks independently and average their outputs

3-15 networks  $\rightarrow \hat{y}$

### [Multi-crop at test time

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

10-crop



1



4



1



4

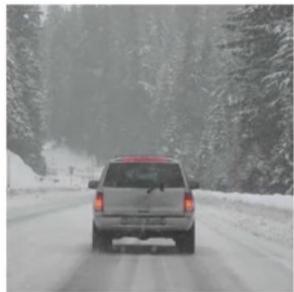
Andrew Ng

## Object Detection

## Object localization

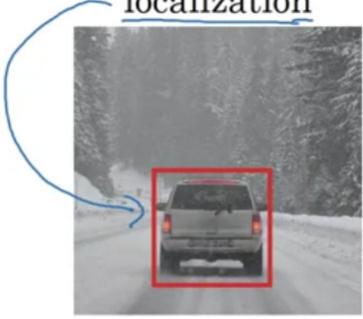
# What are localization and detection?

Image classification



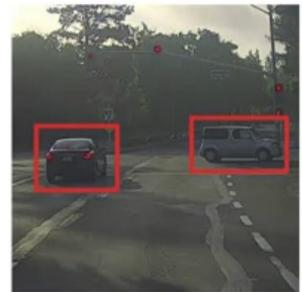
*"Car"*  
↓  
1 object

Classification with localization



*"Car"*

Detection



multiple  
objects

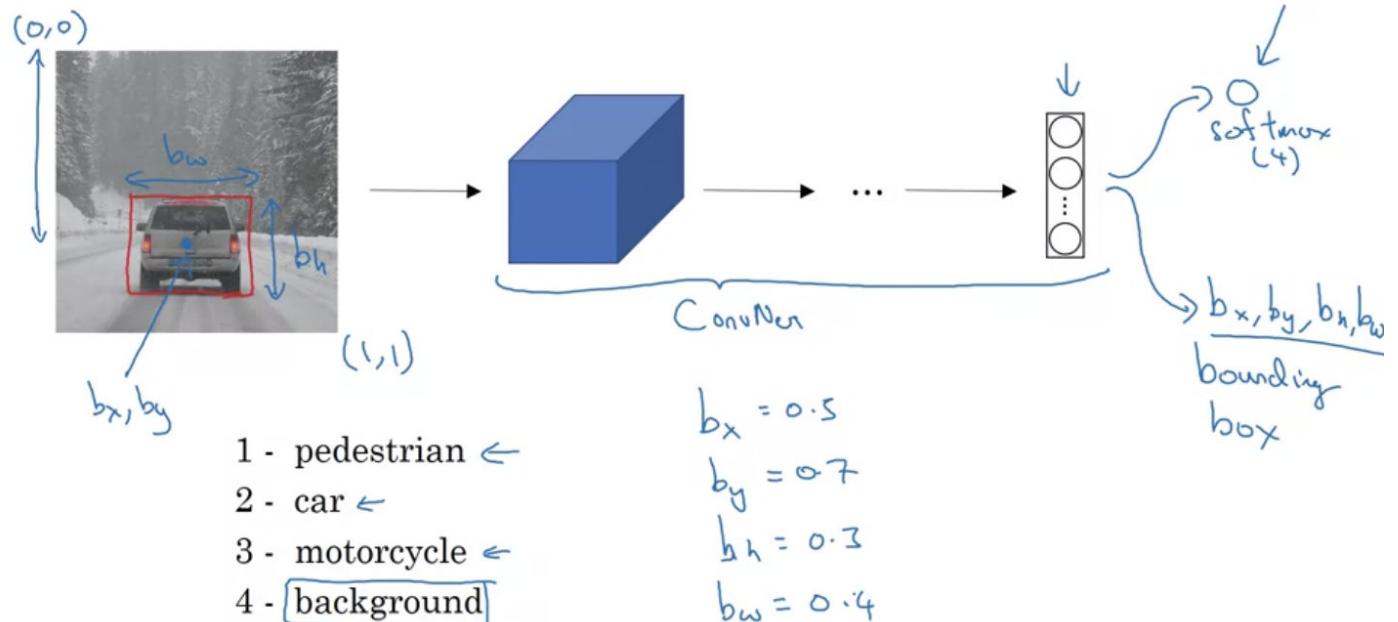
In this section, the **upper left** of the image is  $(0, 0)$ , and at the **lower right** is  $(1, 1)$ .

The bounding box is denoted by:

$$(b_x, b_y, b_h, b_w)$$

where  $b_x, b_y$  are the coordinates of the midpoint,  $b_h, b_w$  are height and weight respectively.

## Classification with localization



Andrew Ng

# Defining the target label $y$

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

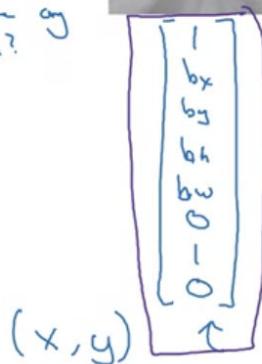
$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

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

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



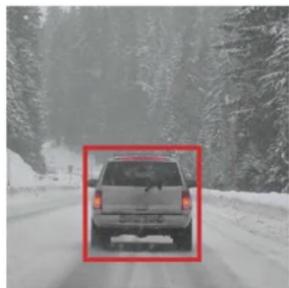
$x =$



Andrew Ng

## Landmark detection

### Landmark detection



$b_x, b_y, b_h, b_w$

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \rightarrow 0$$

$$\left. \begin{array}{l} l_{1x}, l_{1y}, \\ l_{2x}, l_{2y}, \\ l_{3x}, l_{3y}, \\ l_{4x}, l_{4y}, \\ \vdots \\ l_{64x}, l_{64y} \end{array} \right\} X, Y$$

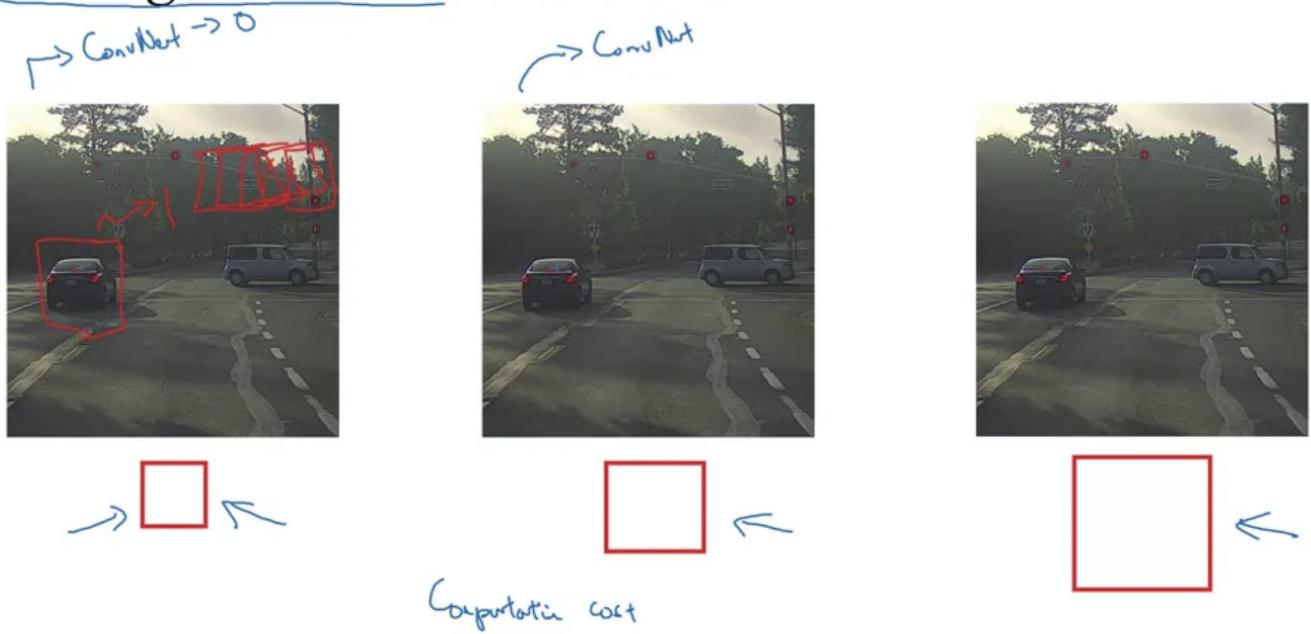
$$\left. \begin{array}{l} l_{1x}, l_{1y}, \\ \vdots \\ l_{32}, l_{32} \end{array} \right.$$

Andrew Ng

## Object detection

Before the raise of deep learning, people used **sliding windows** detection with simple linear classifier over hand-engineered features in order to perform object detection.

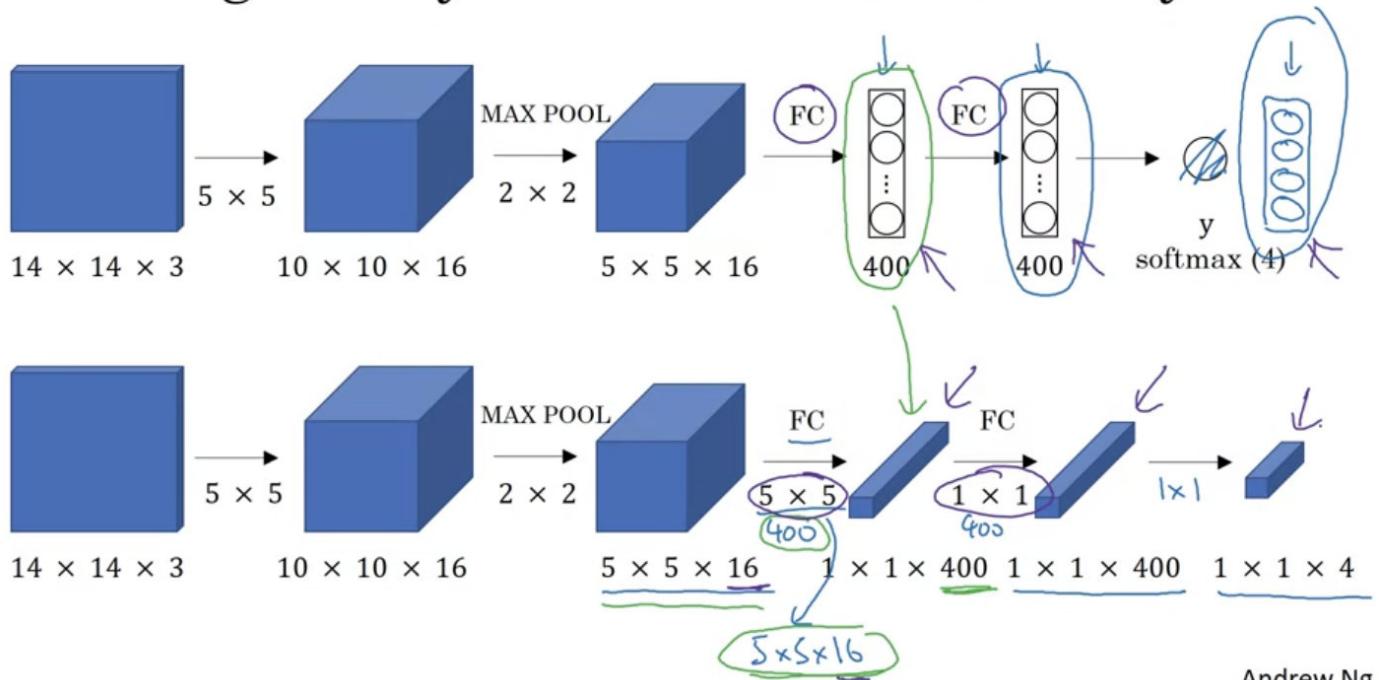
## Sliding windows detection



Andrew Ng

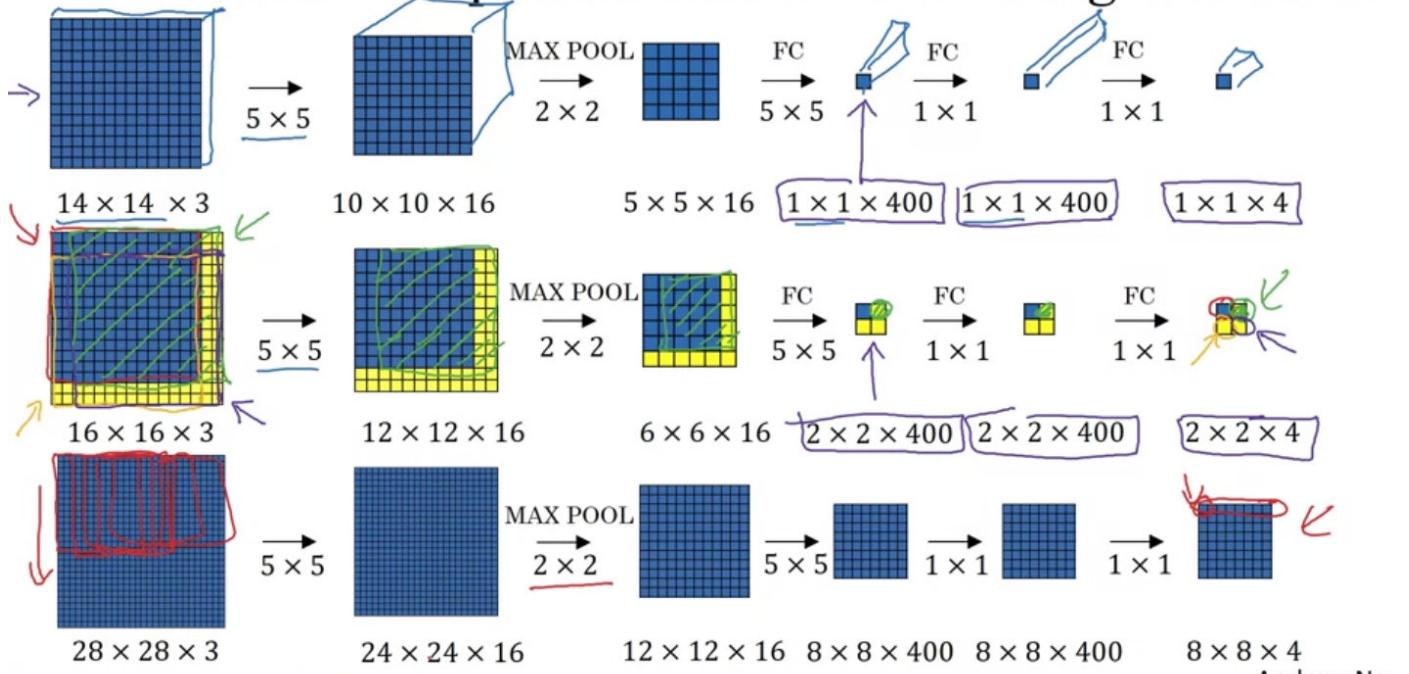
### Convolutional implementation of sliding windows

## Turning FC layer into convolutional layers



Andrew Ng

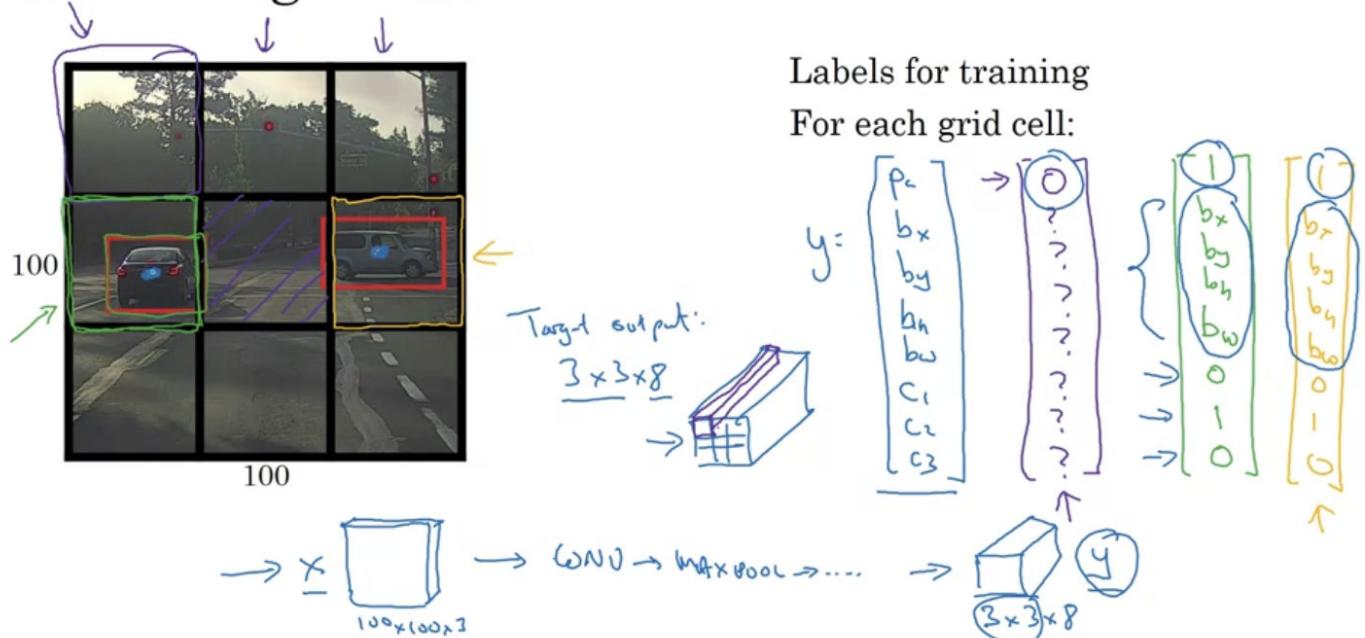
# Convolution implementation of sliding windows



[Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks]

## Bounding box predictions

### YOLO algorithm

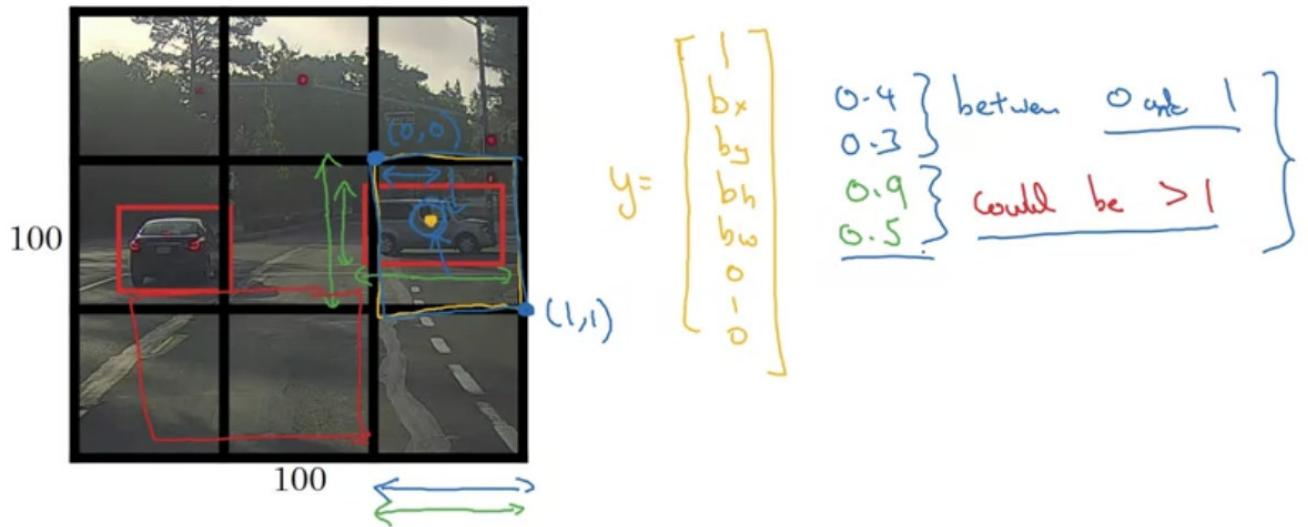


[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

The bounding box for YOLO ( $b_x, b_y, b_h, b_w$ ) is specified relative to the grid cell. And the width and height of the bounding box are specified as **fractions** of the overall width and height of the grid cell.

# Specify the bounding boxes

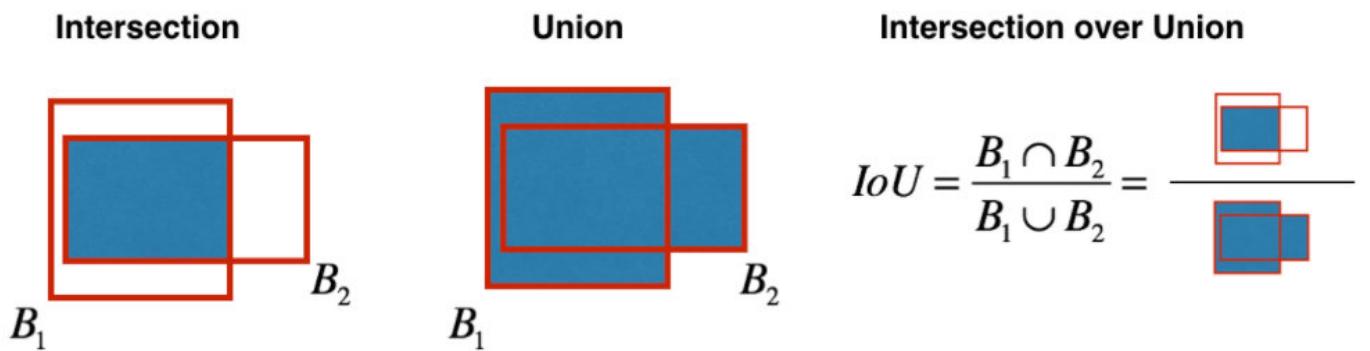


## Intersection over Union (IoU)

$$IoU = \frac{\text{Area of intersection}}{\text{Area of union}}$$

"Correct" if  $IoU \geq 0.5$

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

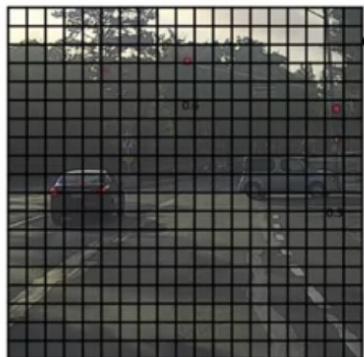


## Non-max suppression

Non-max suppression is a way to make sure the algorithm detects each object **only once**.

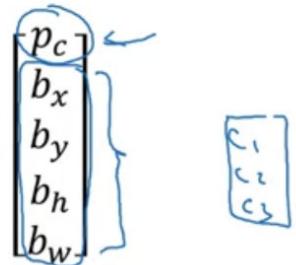
Non-max means that you are going to output your **maximal probabilities classifications**, but suppress the close-by ones that are non-maximal.

# Non-max suppression algorithm



19 × 19

Each output prediction is:



Discard all boxes with  $p_c \leq 0.6$

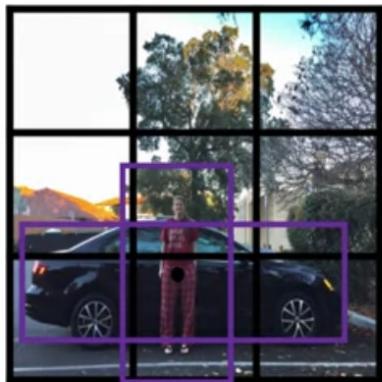
→ While there are any remaining boxes:

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

Andrew Ng

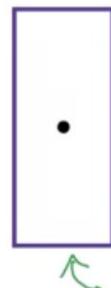
Anchor boxes

Overlapping objects:



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

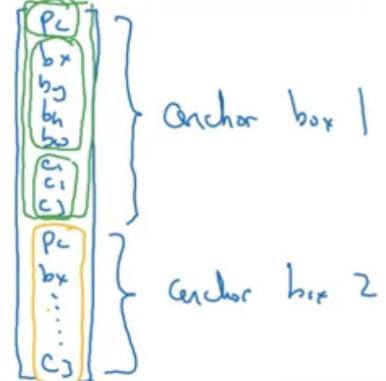
Anchor box 1:



Anchor box 2:



$y =$



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

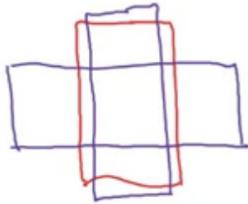
Andrew Ng

# Anchor box algorithm

Previously:

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

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



With two anchor boxes:

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

(grid cell, anchor box)

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

Andrew Ng

Two cases that the algorithm does not handle well:

- If there are **more than two** objects in the same grid cell.
- Two objects are associated with same grid cell but both of them have the **same anchor box shape**.

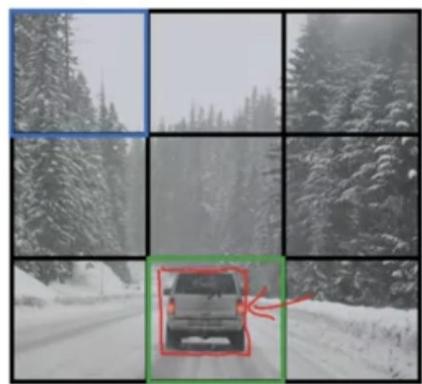
How to choose anchor boxes?

People used to just choose anchor boxes by hand or choose maybe 5 or 10 anchor box shapes that spans a variety of shapes that seems to cover the types of objects you seem to detect. As a much more advanced version is to use a **K-means algorithm** to group together two types of objects shapes you tend to get. And then to use that to select a set of anchor boxes that are most stereotypically representative of the multiple, of the dozens of object classes you're trying to detect.

## YOLO Algorithm

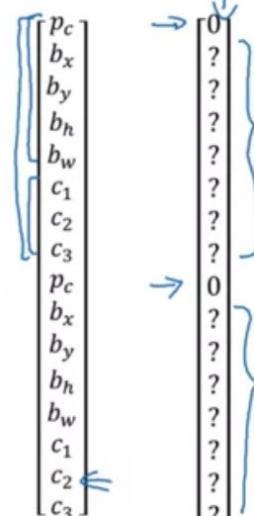
This algorithm "only looks once" at the image in the sense that it requires **only one forward propagation** pass through the network to make predictions.

# Training



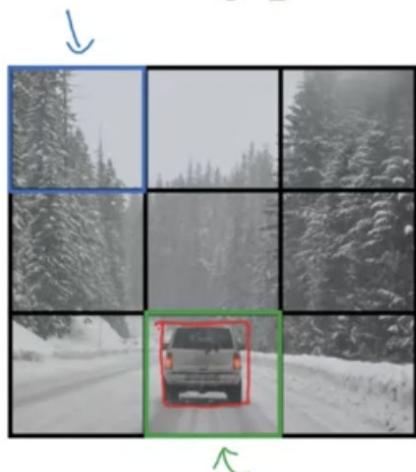
- 1 - pedestrian  
2 - car ←  
3 - motorcycle

$$y = \begin{matrix} 1 \\ 2 \end{matrix}$$



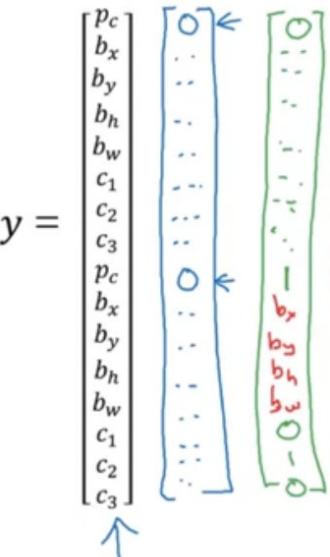
Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

# Making predictions



→ ... →

$$y = \begin{matrix} 3 \times 3 \times 2 \times 8 \end{matrix}$$



# Outputting the non-max suppressed outputs



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

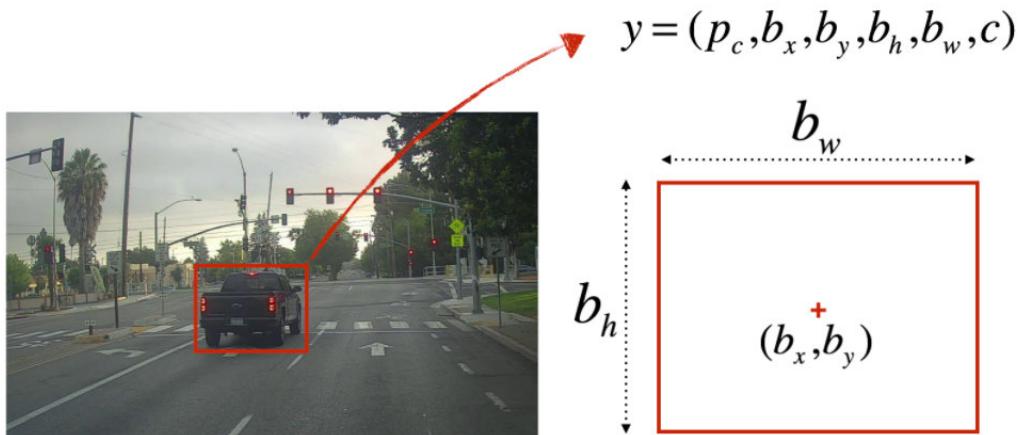
## Programming Assignment Example

### Inputs and outputs

- The **input** is a batch of images, and each image has the shape  $(m, 608, 608, 3)$
- The **output** is a list of bounding boxes along with the recognized classes. Each bounding box is represented by 6 numbers  $(p_c, b_x, b_y, b_h, b_w, c)$  as explained above. If you expand  $c$  into an 80-dimensional vector, each bounding box is then represented by 85 numbers.

### Anchor Boxes

- Anchor boxes are chosen by exploring the training data to choose reasonable height/width ratios that represent the different classes. For this assignment, 5 anchor boxes were chosen for you (to cover the 80 classes), and stored in the file `'./model_data/yolo_anchors.txt'`
- The dimension for anchor boxes is the second to last dimension in the encoding:  $(m, n_H, n_W, anchors, classes)$ .
- The YOLO architecture is: IMAGE  $(m, 608, 608, 3) \rightarrow$  DEEP CNN  $\rightarrow$  ENCODING  $(m, 19, 19, 5, 85)$ .



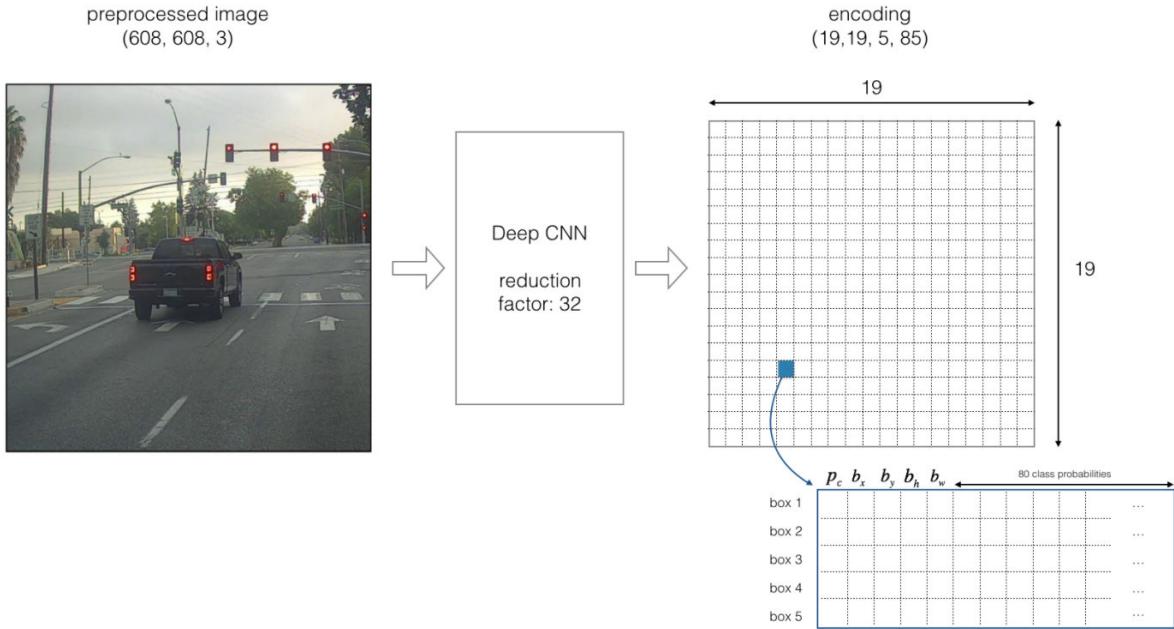
$p_c = 1$  : confidence of an object being present in the bounding box

$c = 3$  : class of the object being detected (here 3 for "car")

**Figure 1:** Definition of a box

## Encoding

Let's look in greater detail at what this encoding represents.

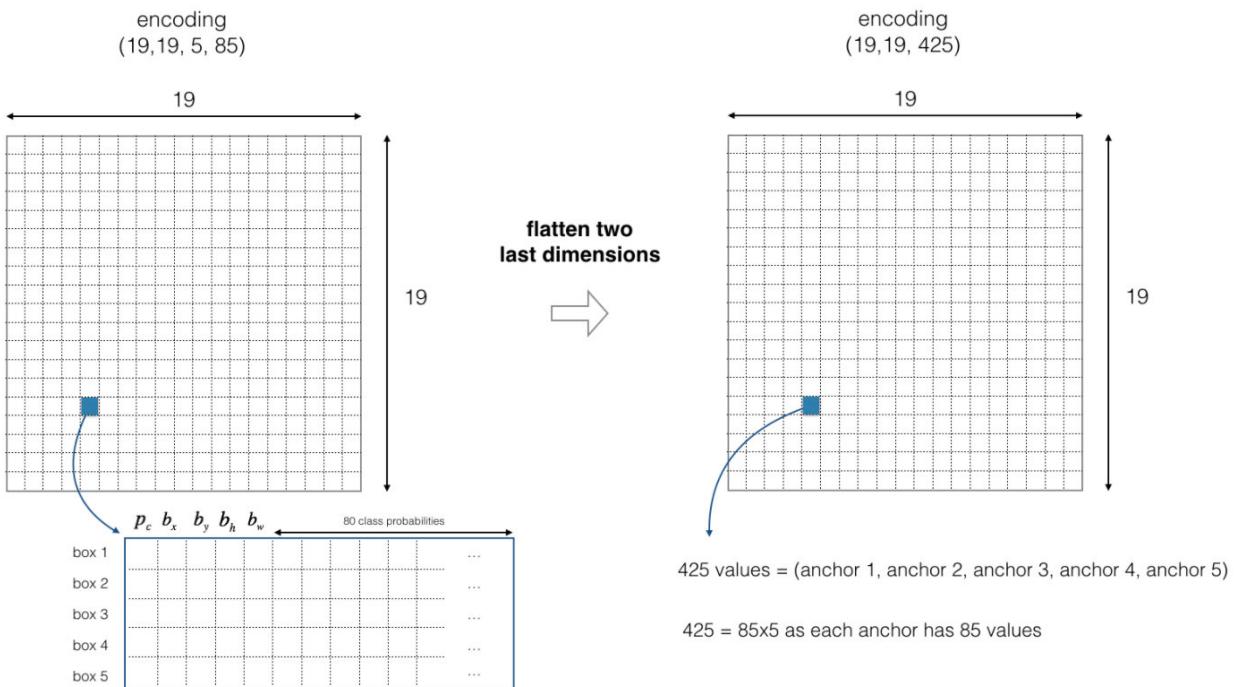


**Figure 2:** Encoding architecture for YOLO

If the center/midpoint of an object falls into a grid cell, that grid cell is responsible for detecting that object.

Since you're using 5 anchor boxes, each of the  $19 \times 19$  cells thus encodes information about 5 boxes. Anchor boxes are defined only by their width and height.

For simplicity, you'll flatten the last two dimensions of the shape  $(19, 19, 5, 85)$  encoding, so the output of the Deep CNN is  $(19, 19, 425)$ .

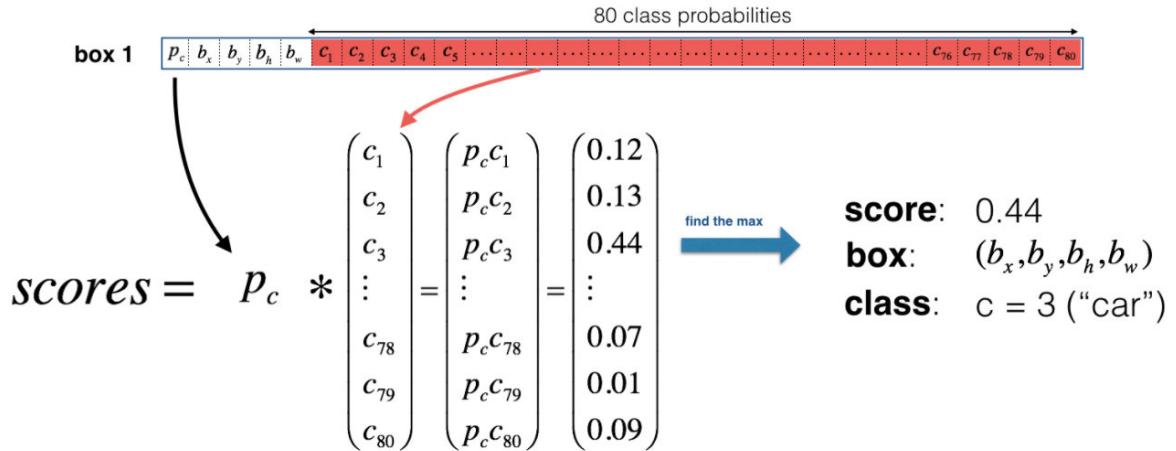


**Figure 3:** Flattening the last two last dimensions

## Class score

Now, for each box (of each cell) you'll compute the following element-wise product and extract a probability that the box contains a certain class.

The class score is  $score_{c,i} = p_c \times c_i$ : the probability that there is an object  $p_c$  times the probability that the object is a certain class  $c_i$ .



the box  $(b_x, b_y, b_h, b_w)$  has detected  $c = 3$  ("car") with probability score: 0.44

**Figure 4:** Find the class detected by each box

## Visualizing classes

Here's one way to visualize what YOLO is predicting on an image:

- For each of the  $19 \times 19$  grid cells, find the maximum of the probability scores (taking a max across the 80 classes, one maximum for each of the 5 anchor boxes).
- Color that grid cell according to what object that grid cell considers the most likely.

Doing this results in this picture:



**Figure 5:** Each one of the  $19 \times 19$  grid cells is colored according to which class has the largest predicted probability in that cell.

Note that this visualization isn't a core part of the YOLO algorithm itself for making predictions; it's just a nice way of visualizing an intermediate result of the algorithm.

## Visualizing bounding boxes

Another way to visualize YOLO's output is to plot the bounding boxes that it outputs. Doing that results in a visualization like this:



**Figure 6:** Each cell gives you 5 boxes. In total, the model predicts:  $19 \times 19 \times 5 = 1805$  boxes just by looking once at the image (one forward pass through the network)! Different colors denote different classes.

## Non-Max suppression

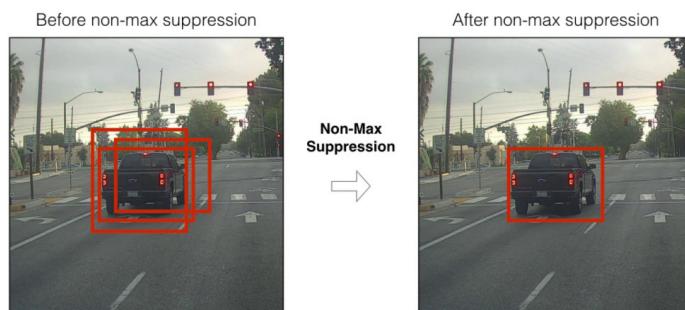
In the figure above, the only boxes plotted are ones for which the model had assigned a high probability, but this is still too many boxes. You'd like to reduce the algorithm's output to a much smaller number of detected objects.

To do so, you'll use **non-max suppression**. Specifically, you'll carry out these steps:

- Get rid of boxes with a low score. Meaning, the box is not very confident about detecting a class, either due to the low probability of any object, or low probability of this particular class.
- Select only one box when several boxes overlap with each other and detect the same object.

The model gives you a total of  $19 \times 19 \times 5 \times 85$  numbers, with each box described by 85 numbers. It's convenient to rearrange the  $(19, 19, 5, 85)$  (or  $(19, 19, 425)$ ) dimensional tensor into the following variables:

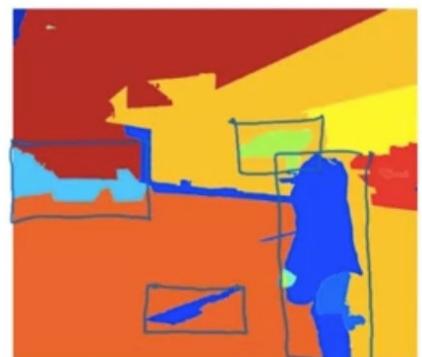
- **box\_confidence** : tensor of shape  $(19, 19, 5, 1)$  containing  $p_c$  (confidence probability that there's some object) for each of the 5 boxes predicted in each of the  $19 \times 19$  cells.
- **boxes** : tensor of shape  $(19, 19, 5, 4)$  containing the midpoint and dimensions  $(b_x, b_y, b_h, b_w)$  for each of the 5 boxes in each cell.
- **box\_class\_probs** : tensor of shape  $(19, 19, 5, 80)$  containing the "class probabilities"  $(c_1, c_2, \dots, c_{80})$  for each of the 80 classes for each of the 5 boxes per cell.



**Figure 7:** In this example, the model has predicted 3 cars, but it's actually 3 predictions of the same car. Running non-max suppression (NMS) will select only the most accurate (highest probability) of the 3 boxes.

## Region proposal: R-CNN

# Region proposal: R-CNN



Segmentation algorithm

$\sim 2,000$

## 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

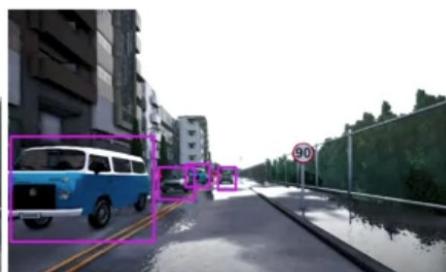
Use convolutional network to propose regions.

## Semantic Segmentation

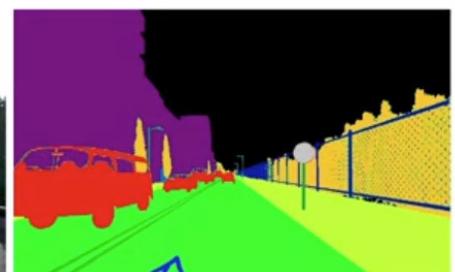
# Object Detection vs. Semantic Segmentation



Input image



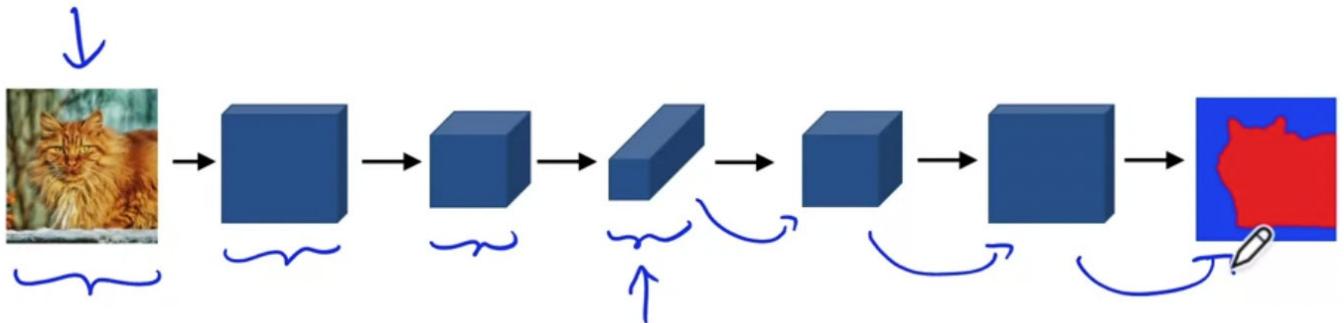
Object Detection



Semantic Segmentation

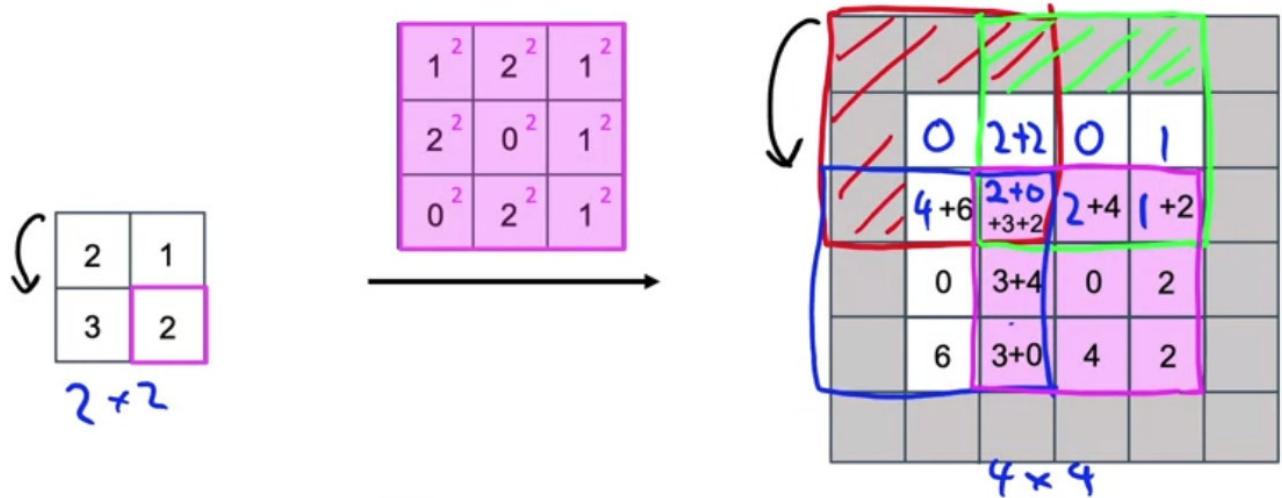
## U-Net

# Deep Learning for Semantic Segmentation



## Transpose Convolutions

### Transpose Convolution



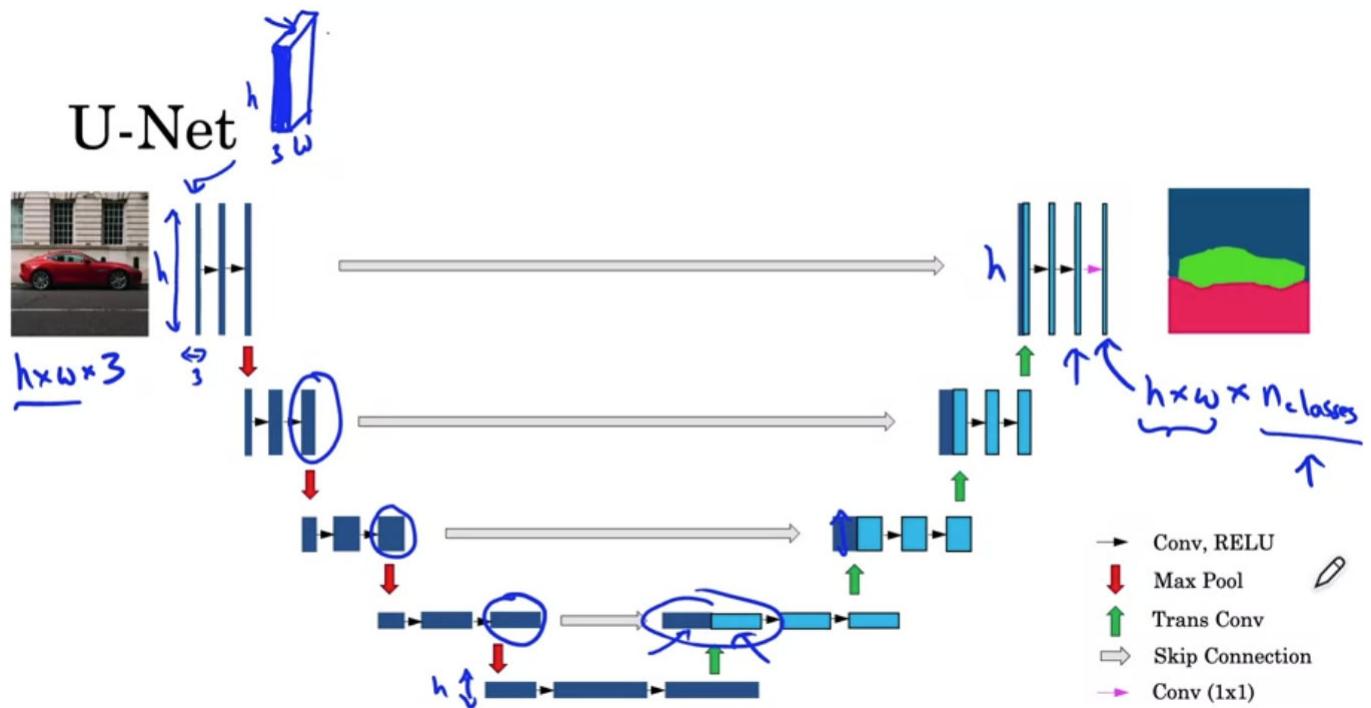
filter  $f \times f = 3 \times 3$    padding  $p = 1$    stride  $s = 2$

## U-Net Architecture

U-Net builds on a previous architecture called the Fully Convolutional Network, or FCN, which replaces the dense layers found in a typical CNN with a transposed convolution layer that upsamples the feature map back to the size of the original input image, while preserving the spatial information. This is necessary because the dense layers destroy spatial information (the "where" of the image), which is an essential part of image segmentation tasks. An added bonus of using transpose convolutions is that the input size no longer needs to be fixed, as it does when dense layers are used.

Unfortunately, the final feature layer of the FCN suffers from information loss due to downsampling too much. It then becomes difficult to upsample after so much information has been lost, causing an output that looks rough.

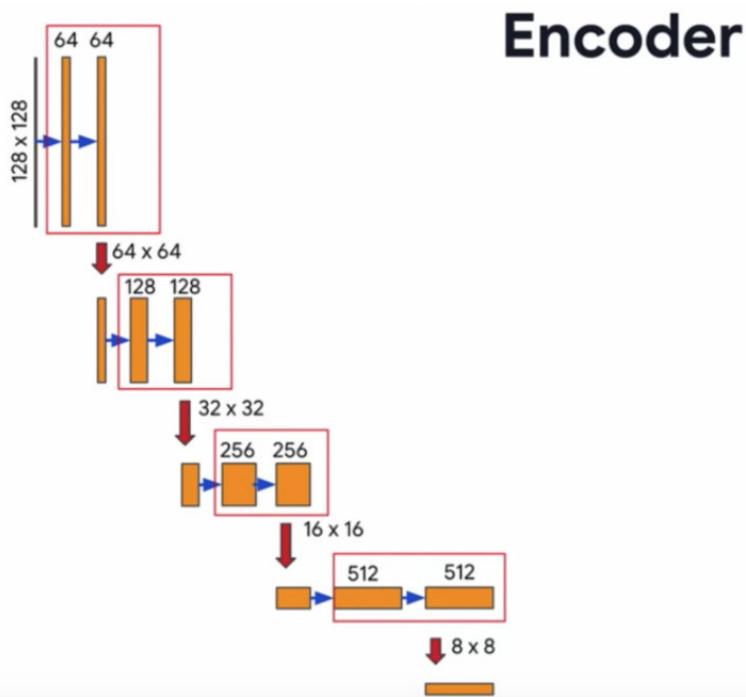
U-Net improves on the FCN, using a somewhat similar design, but differing in some important ways. Instead of one transposed convolution at the end of the network, it uses a matching number of convolutions for downsampling the input image to a feature map, and **transposed convolutions** for upsampling those maps back up to the original input image size. It also adds **skip connections**, to retain information that would otherwise become lost during encoding. Skip connections send information to every upsampling layer in the decoder from the corresponding downsampling layer in the encoder, capturing finer information while also keeping computation low. These help prevent information loss, as well as model overfitting.



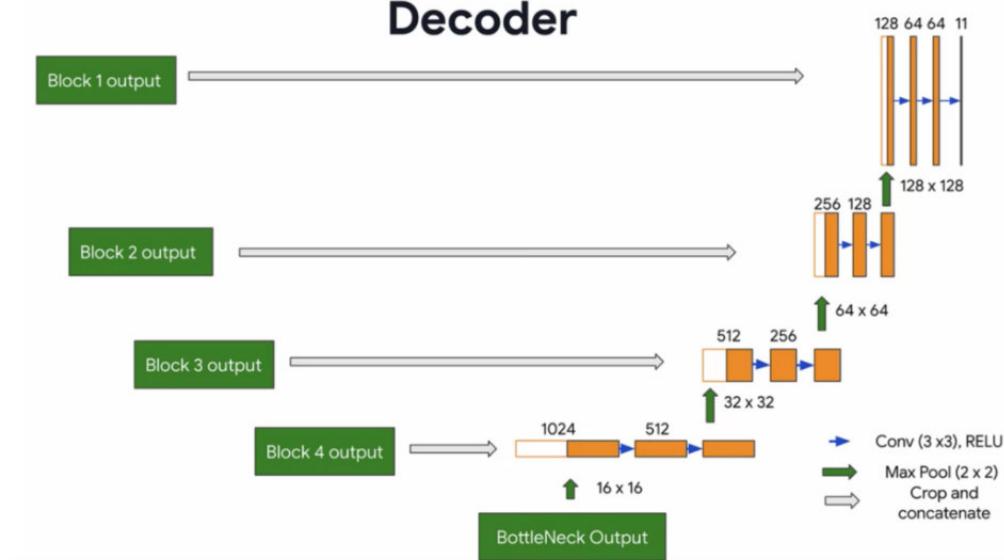
## Encoder

The encoder is a stack of various conv\_blocks:

Each `conv_block()` is composed of 2 Conv2D layers with ReLU activations.



## Decoder



## 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 ←

Andrew Ng

## One-shot learning

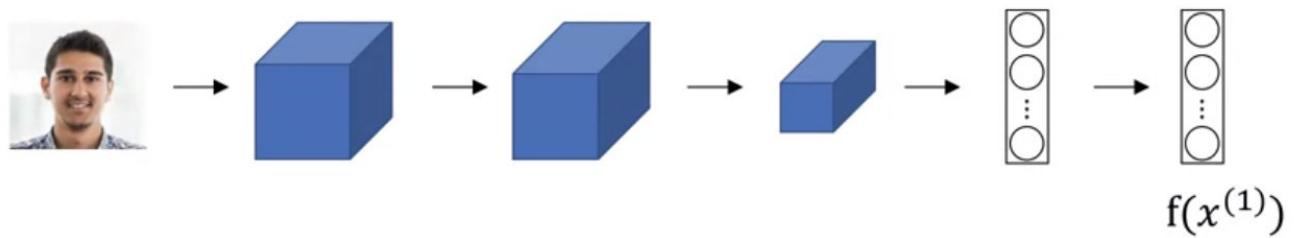
Learning from **one example** to recognize the person again. The method is learn a **similarity function**.

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

If  $d(\text{img1}, \text{img2}) \leq \tau \rightarrow \text{'same'}$ ,  $d(\text{img1}, \text{img2}) > \tau \rightarrow \text{'diff'}$

## Siamese Network

# Goal of learning



Parameters of NN define an encoding  $f(x^{(i)})$  128

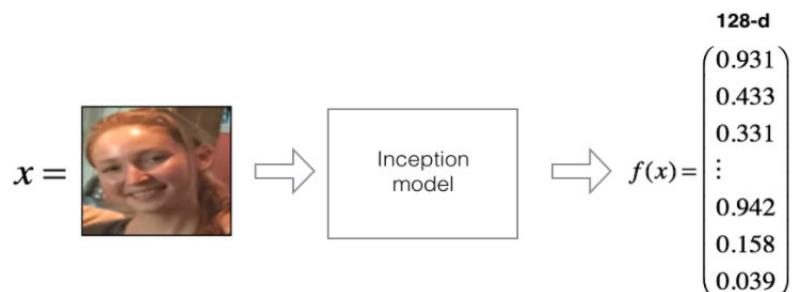
Learn parameters so that:

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

For an image  $x$ , its encoding is denoted as  $f(x)$ , where  $f$  is the function computed by the neural network.



# Learning Objective



Anchor A      Positive P  
 $d(A, P) = 0.5$

Anchor A      Negative N  
 $d(A, N) = 0.55$

Want:  $\frac{\|f(A) - f(P)\|^2}{d(A, P)} + \alpha \leq \frac{\|f(A) - f(N)\|^2}{d(A, N)}$

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

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

Andrew Ng

Training will use triplets of images  $(A, P, N)$ :

- A is an "Anchor" image--a picture of a person.
- P is a "Positive" image--a picture of the same person as the Anchor image.
- N is a "Negative" image--a picture of a different person than the Anchor image.

These triplets are picked from the training dataset.  $(A^{(i)}, P^{(i)}, N^{(i)})$  is used here to denote the  $i$ -th training example.

You'd like to make sure that an image  $A^{(i)}$  of an individual is closer to the Positive  $P^{(i)}$  than to the Negative image  $N^{(i)}$  by at least a margin  $\alpha$ :

$$\|f(A^{(i)}) - f(P^{(i)})\|_2^2 + \alpha < \|f(A^{(i)}) - f(N^{(i)})\|_2^2$$

You would thus like to minimize the following "triplet cost":

$$\mathcal{J} = \sum_{i=1}^m \left[ \underbrace{\|f(A^{(i)}) - f(P^{(i)})\|_2^2}_{(1)} - \underbrace{\|f(A^{(i)}) - f(N^{(i)})\|_2^2}_{(2)} + \alpha \right]_+$$

Here, the notation " $[z]_+$ " is used to denote  $\max(z, 0)$ .

Notes:

- The term (1) is the squared distance between the anchor "A" and the positive "P" for a given triplet; you want this to be small.
- The term (2) is the squared distance between the anchor "A" and the negative "N" for a given triplet, you want this to be relatively large. It has a minus sign preceding it because minimizing the negative of the term is the same as maximizing that term.
- $\alpha$  is called the margin. It's a hyperparameter that you pick manually. e.g.  $\alpha = 0.2$ .

Training set: 10k pictures of 1k persons -- Take 10k pictures and use it to generate the triplets  $(A, P, N)$ .

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

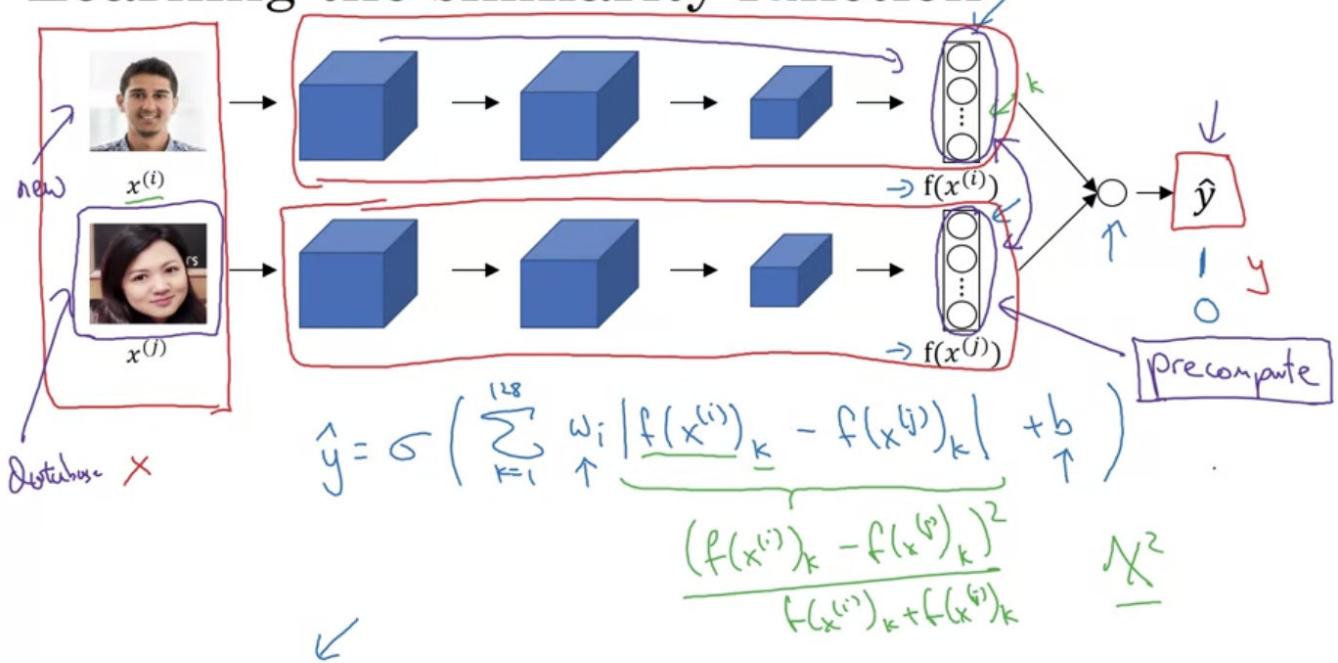
During training, if  $A, P, N$  are chosen **randomly**,  $d(A, P) + \alpha \leq d(A, N)$  is **easily satisfied**.

Choose triplets that are "hard" to train on.

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

## Face Verification and Binary Classification

### Learning the similarity function



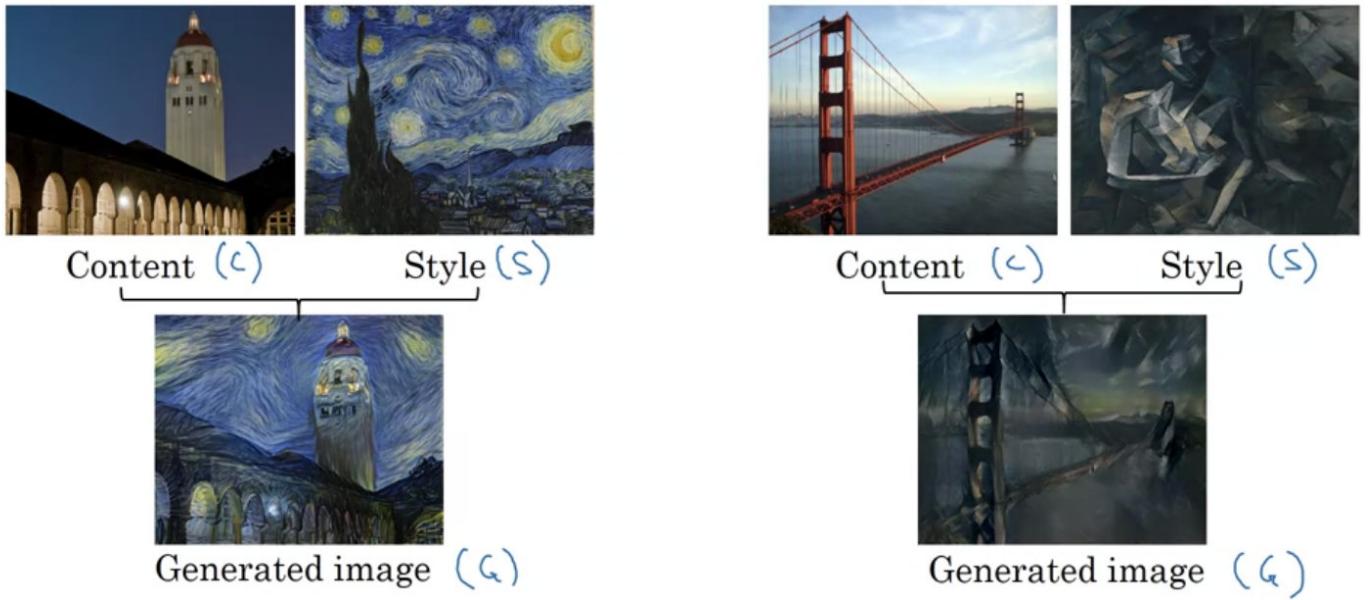
Andrew Ng

#### What you should remember:

- Face verification solves an easier 1:1 matching problem; face recognition addresses a harder 1:K matching problem.
- Triplet loss is an effective loss function for training a neural network to learn an encoding of a face image.
- The same encoding can be used for verification and recognition. Measuring distances between two images' encodings allows you to determine whether they are pictures of the same person.

## Neural Style Transfer

# Neural style transfer



## Cost function

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

Where  $J_{content}(C, G)$  measures how similar is the content of the generated image  $G$ , to the content fo the content image  $C$ ,  $J_{style}(C, G)$  measures how similar is the style of the generated image  $G$ , to the style fo the style image  $S$ .

## Find the generated image $G$

1. Initialize  $G$  randomly.

$$G : 100 \times 100 \times 3$$

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

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

## Content cost function

Recap the cost function:

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

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

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

## Style cost function

Define style as **correlation** between activations across channels.

### Style matrix (Gram matrix)

- The style matrix is also called a "Gram matrix".
- In linear algebra, the Gram matrix  $G$  of a set of vectors  $(v_1, \dots, v_n)$  is the matrix of dot products, whose entries are  $G_{ij} = v_i^T v_j = np. dot(v_i, v_j)$ .
- In other words,  $G_{ij}$  compares how similar  $v_i$  is to  $v_j$ : If they are highly similar, you would expect them to have a large dot product, and thus for  $G_{ij}$  to be large.

Let  $a_{i,j,k}^{[l]} = \text{activation at } (i, j, k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$ . Particularly,  $G_{kk'}^{[l]}$  will measure how correlated are the activations in channel  $k$  compared to the activations in channel  $k'$ .

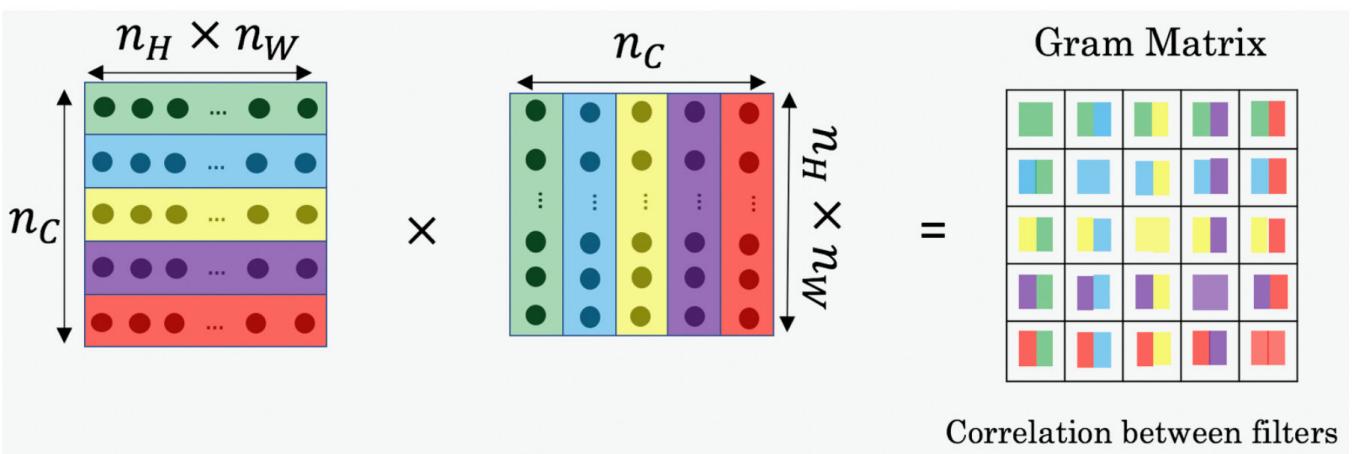
$$G_{kk'}^{[l](S)} = \sum_i^{n_H^{[l]}} \sum_j^{n_W^{[l]}} a_{i,j,k}^{[l](S)} a_{i,j,k'}^{[l](S)}$$

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

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

### Compute Gram matrix $G_{gram}$

You will compute the Style matrix by multiplying the "unrolled" filter matrix with its transpose:



### $G_{(gram)ij}$ : correlation

The result is a matrix of dimension  $(n_C, n_C)$  where  $n_C$  is the number of filters (channels). The value  $G_{(gram)ij}$  measures how similar the activations of filter  $i$  are to the activations of filter  $j$ .

### $G_{(gram),ii}$ : prevalence of patterns or textures

- The diagonal elements  $G_{(gram)ii}$  measure how "active" a filter  $i$  is.
- For example, suppose filter  $i$  is detecting vertical textures in the image. Then  $G_{(gram)ii}$  measures how common vertical textures are in the image as a whole.
- If  $G_{(gram)ii}$  is large, this means that the image has a lot of vertical texture.

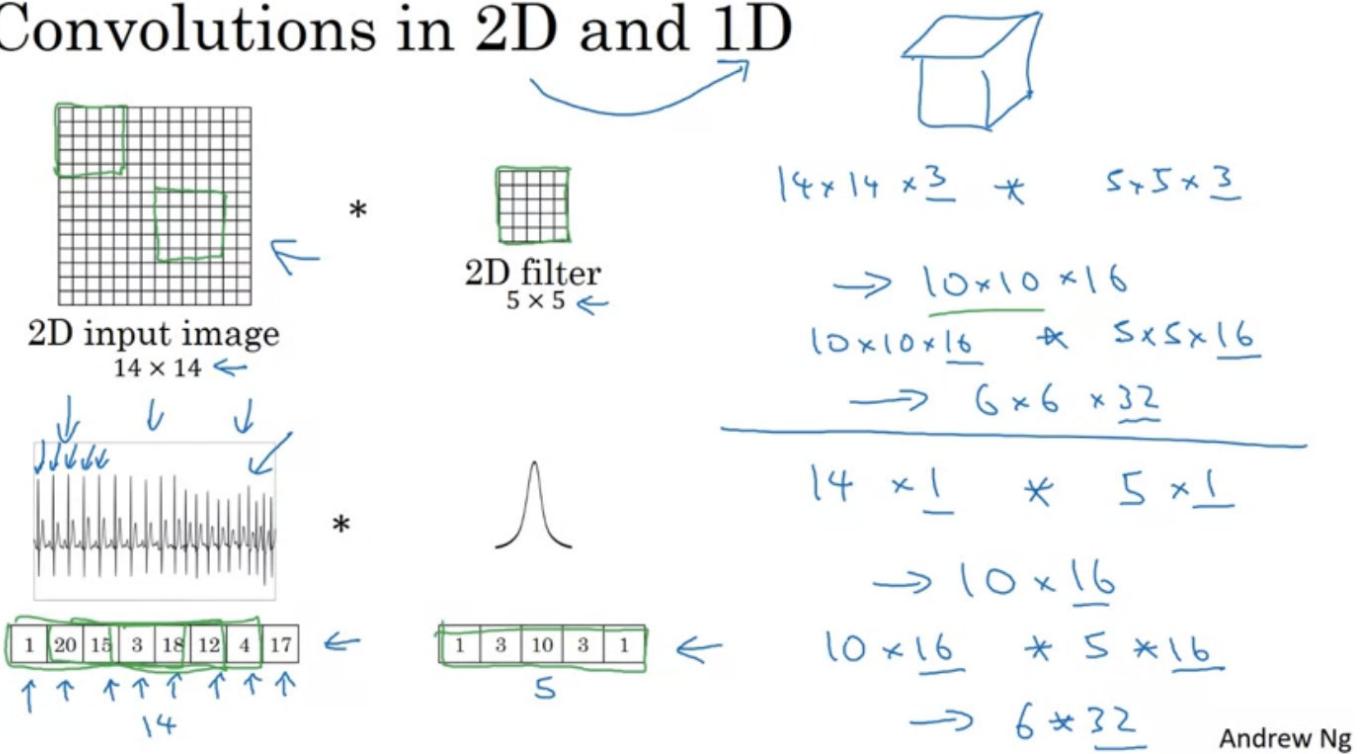
By capturing the prevalence of different types of features ( $G_{(gram)ii}$ ), as well as how much different features occur together ( $G_{(gram)ij}$ ), the Style matrix  $G_{gram}$  measures the style of an image.

Neural style transfer is trained as a supervised learning task in which the goal is to input two images ( $x$ ), and train a network to output a new, synthesized image ( $y$ ).

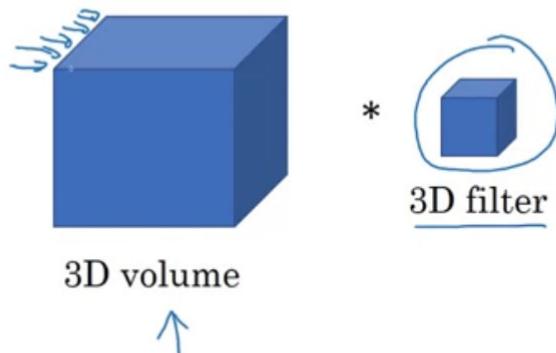
**FALSE:** Neural style transfer is about training on the pixels of an image to make it look artistic, it is **not** learning any parameters.

## 1D and 3D Generalizations

# Convolutions in 2D and 1D



# 3D convolution



$$\begin{array}{c}
 \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\
 \cancel{14 \times 14 \times 14} \times \cancel{1} \quad n_c
 \end{array}$$

$$* \cancel{5 \times 5 \times 5 \times 1} \quad 16 \text{ filters.}$$

$$\rightarrow 10 \times 10 \times 10 \times \underline{16}$$

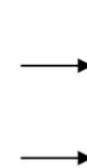
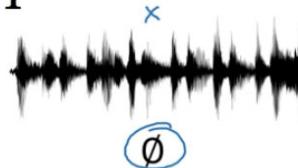
$$* \cancel{5 \times 5 \times 5 \times 16} \quad 32 \text{ filters.}$$

$$\rightarrow 6 \times 6 \times 6 \times 32$$

## Sequence Model

### Examples of sequence data

Speech recognition



"The quick brown fox jumped over the lazy dog."

Music generation



Sentiment classification

"There is nothing to like in this movie."



DNA sequence analysis

AGCCCCTGTGAGGAAC TAG

AGCCCCTGTGAGGAAC TAG

Machine translation

Voulez-vous chanter avec moi?

Do you want to sing with me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter met Hermione Granger.

Yesterday, Harry Potter met Hermione Granger.  
Andrew Ng

## Notation

- Superscript  $[l]$  denotes an object associated with the  $l^{th}$  layer.
- Superscript  $(i)$  denotes an object associated with the  $i^{th}$  example.
- Superscript  $\langle t \rangle$  denotes an object at the  $t^{th}$  time step.
- Subscript  $i$  denotes the  $i^{th}$  entry of a vector.

- $T_x^{(i)}$ : the input sequence length for training example  $i$ .
- $T_y^{(i)}$ : the output sequence length for training example  $i$ .

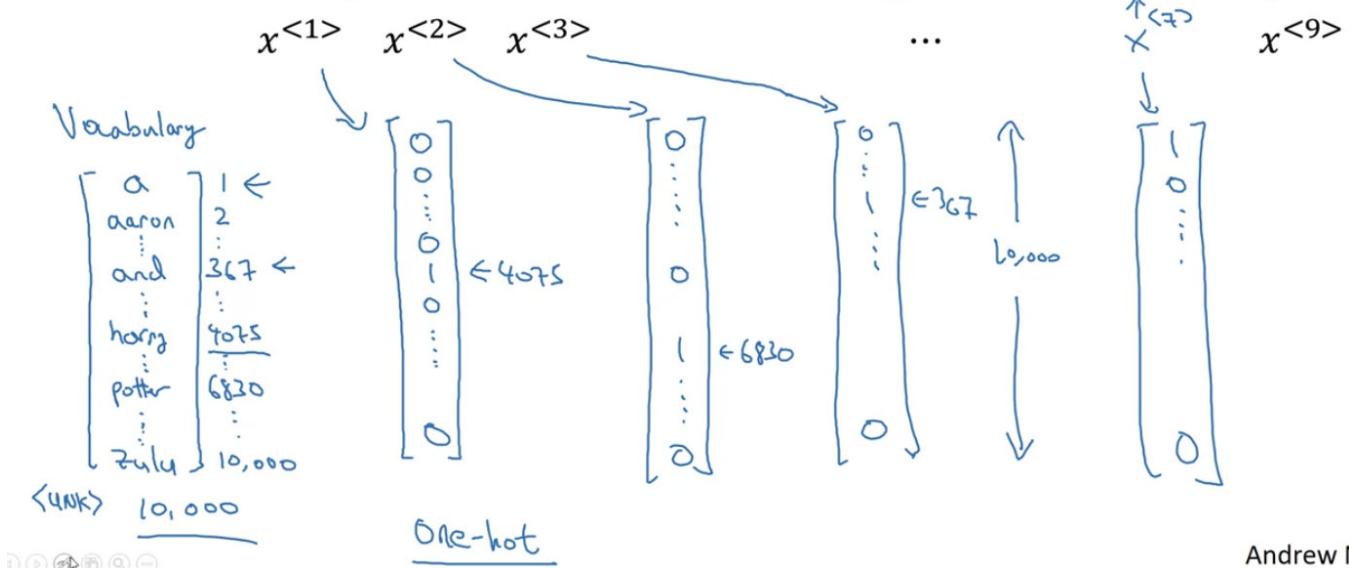
**Example:**

- $a_5^{(2)[3]<4>}$  denotes the activation of the 2nd training example (2), 3rd layer [3], 4th time step  $<4>$ , and 5th entry in the vector.

## Representing words

$$x^{<t>} \rightarrow y^{(x,y)}$$

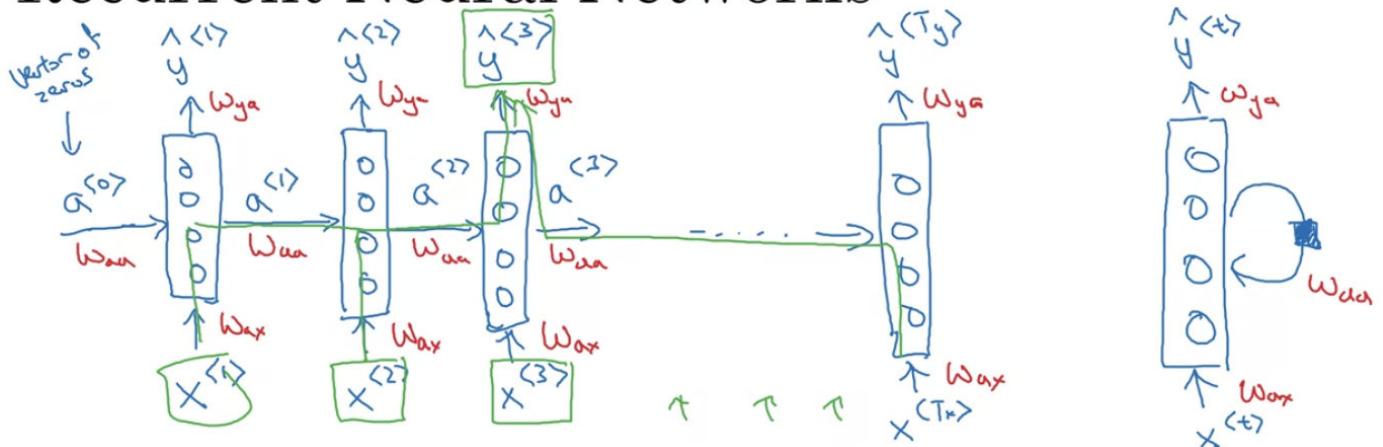
x: Harry Potter and Hermione Granger invented a new spell.



Andrew Ng

## Recurrent Neural Network Model

### Recurrent Neural Networks

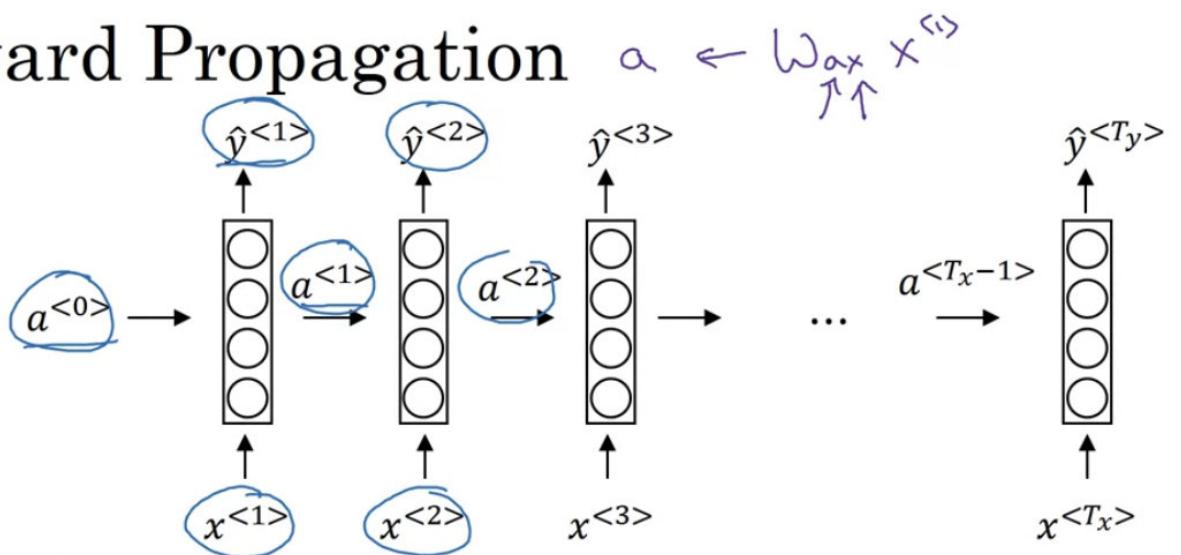


When making the prediction for  $y_3$ , it gets the information **not only from  $x_3$  but also the information from  $x_1$  and  $x_2$** .

One weakness of RNN is that it only uses the information that is **earlier** in the sequence to make prediction.

## Forward Propagation

# Forward Propagation



Initialize  $a^{<0>} = \vec{0}$ .

Then,

$$a^{<1>} = g(W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a) \leftarrow \tanh(\text{ReLU})$$

$$\hat{y}^{<1>} = g(W_{ya}a^{<1>} + b_y) \leftarrow \text{depends on tasks}$$

Generally,

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

## Simplified RNN notation

$$a^{<t>} = g(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

$$\hat{y}^{<t>} = g(W_ya^{<t>} + b_y)$$

where

$$W_a = [W_{aa} | W_{ax}]$$

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

So that,

$$[W_{aa} | W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

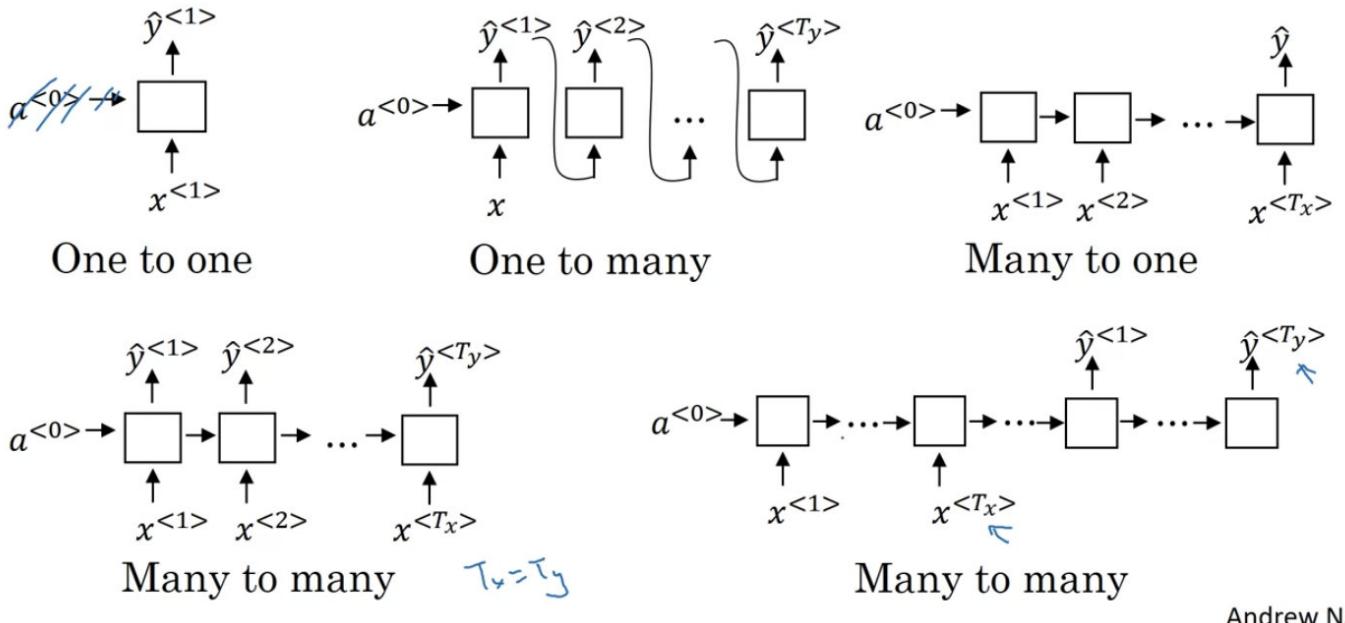
## Backpropagation Through Time

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<y>}) \log(1 - \hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

## Different Types of RNNs

# Summary of RNN types



Andrew Ng

## Language Model and Sequence Generation

# What is language modelling?

Speech recognition

The apple and pair salad.

→ The apple and pear salad.

$$P(\text{The apple and pair salad}) = 3.2 \times 10^{-3}$$

$$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

$$P(\text{sentence}) = ?$$

$$P(y^{<1>} \rightarrow y^{<2>} \rightarrow \dots \rightarrow y^{<T_y>})$$

**The probability of sentence:** the chance that the next sentence you read somewhere out there in the world will be the particular sentence.

The language model represent a sentences as **outputs  $y$**  rather than **inputs  $x$** .

The **basic job** of a language model is to estimate the probability of the particular input sequences of words.

# Language modelling with an RNN

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day.  $\downarrow <\text{EOS}>$

$y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad \dots \quad y^{(8)} \quad y^{(9)}$   
 $x^{(1)} = y^{(t-1)}$

The Egyptian Mau is a bread of cat.  $<\text{EOS}>$   
 $\downarrow <\text{UNK}>$

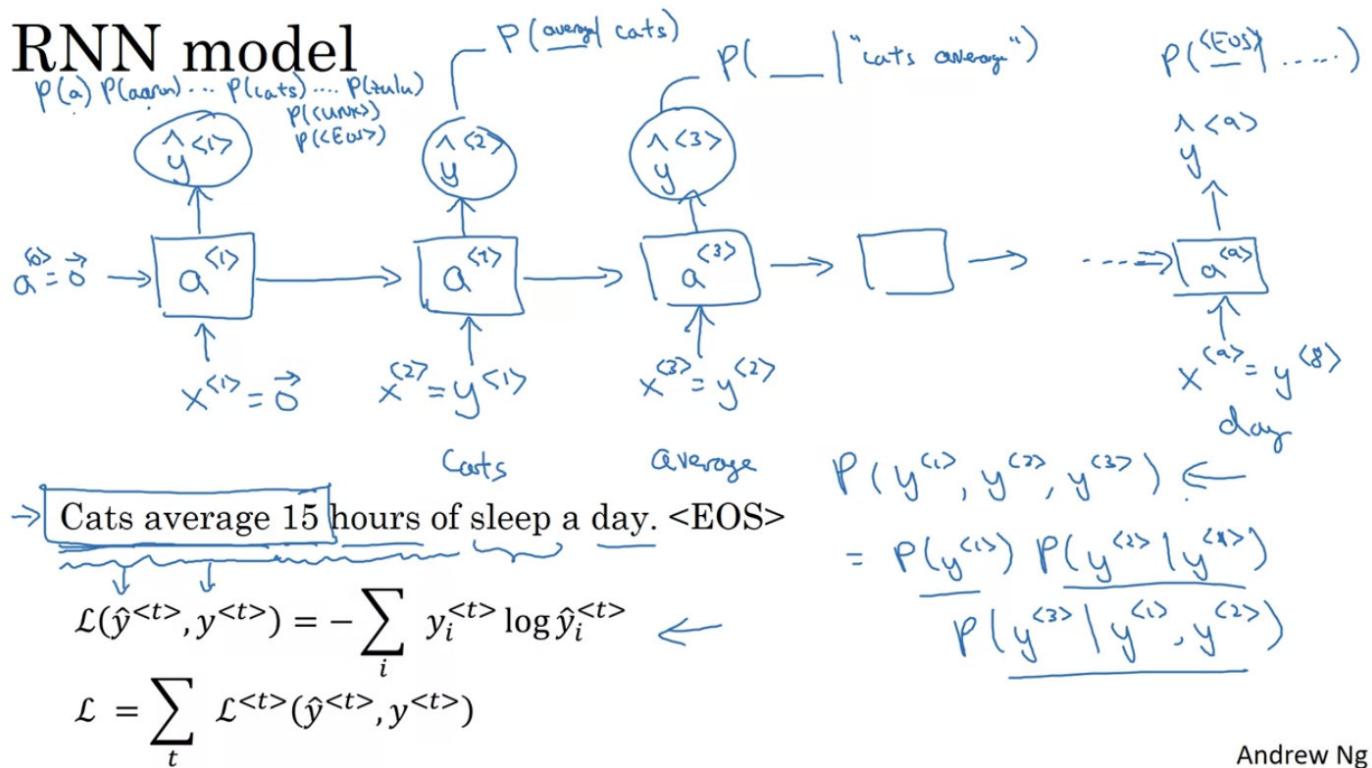
Andrew Ng

**Corpus:** an NLP terminology that just means a large body or a large set of sentences.

**Tokenize:** form a vocabulary and map each of words in the sentence to **one-hot vectors** or **indices** in the vocabulary. By tokenization step, you can decide whether or not the **punctuation** should be token as well.

**<EOS>:** "End Of Sentence", an **optional** extra token that can help to figure out when a sentence ends.

**<UNK>:** a unique token stands for **unknown** words.



Andrew Ng

Each step in the RNN will look at some set of preceding words, such as given the first three words, what is the distribution over the next word? So this learns to predict **one word at a time** going from left to right.

The loss function:

$$L(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>} \leftarrow \text{softmax loss}$$
$$L = \sum_t L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Given a new sentence  $(y^{<1>}, y^{<2>}, y^{<3>})$ , the chance of this entire sentence would be:

$$P(y^{<1>}, y^{<2>}, y^{<3>}) = P(y^{<1>})P(y^{<2>}|y^{<1>})P(y^{<3>}|y^{<2>}, y^{<1>})$$

## Implement RNN

### Dimensions of input $x$

#### Input with $n_x$ number of units

- For a single time step of a single input example,  $x^{(i)(t)}$  is a one-dimensional input vector
- Using language as an example, a language with a 5000-word vocabulary could be one-hot encoded into a vector that has 5000 units. So  $x^{(i)(t)}$  would have the shape (5000,)
- The notation  $n_x$  is used here to denote the number of units in a single time step of a single training example

#### Time steps of size $T_x$

- A recurrent neural network has multiple time steps, which you'll index with  $t$ .
- In the lessons, you saw a single training example  $x^{(i)}$  consisting of multiple time steps  $T_x$ . In this notebook,  $T_x$  will denote the number of timesteps in the longest sequence.

#### Batches of size $m$

- Let's say we have mini-batches, each with 20 training examples
- To benefit from vectorization, you'll stack 20 columns of  $x^{(i)}$  examples
- For example, this tensor has the shape (5000,20,10)
- You'll use  $m$  to denote the number of training examples
- So, the shape of a mini-batch is  $(n_x, m, T_x)$

### Definition of hidden state $a$

- The activation  $a^{(t)}$  that is passed to the RNN from one time step to another is called a "hidden state."

### Dimensions of hidden state $a$

- Similar to the input tensor  $x$ , the hidden state for a single training example is a vector of length  $n_a$
- If you include a mini-batch of  $m$  training examples, the shape of a mini-batch is  $(n_a, m)$
- When you include the time step dimension, the shape of the hidden state is  $(n_a, m, T_x)$
- You'll loop through the time steps with index  $t$ , and work with a 2D slice of the 3D tensor
- This 2D slice is referred to as  $a^{(t)}$

## Dimensions of prediction $\hat{y}$

- Similar to the inputs and hidden states,  $\hat{y}$  is a 3D tensor of shape  $(n_y, m, T_y)$ 
  - $n_y$ : number of units in the vector representing the prediction
  - $m$ : number of examples in a mini-batch
  - $T_y$ : number of time steps in the prediction
- For a single time step  $t$ , a 2D slice  $\hat{y}^{(t)}$  has shape  $(n_y, m)$

## Steps:

- Implement the calculations needed for one time step of the RNN.
- Implement a loop over  $T_x$  time steps in order to process all the inputs, one at a time.

## RNN Cell

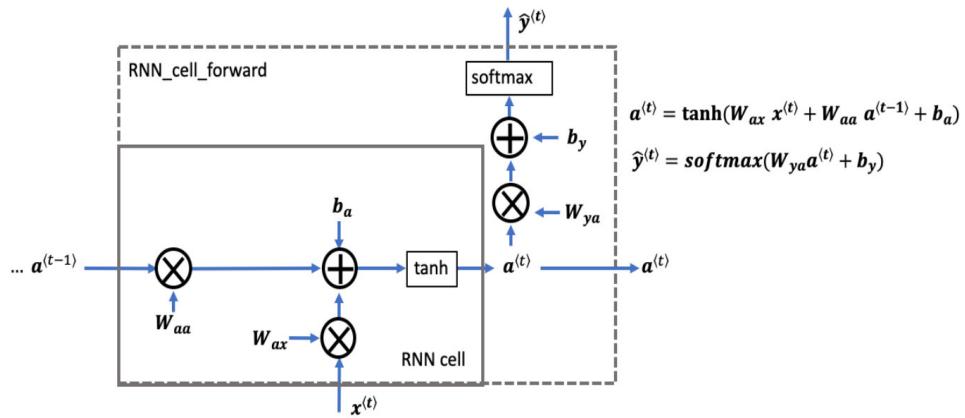


Figure 2: Basic RNN cell. Takes as input  $x^{(t)}$  (current input) and  $a^{(t-1)}$  (previous hidden state containing information from the past), and outputs  $a^{(t)}$  which is given to the next RNN cell and also used to predict  $\hat{y}^{(t)}$

### RNN cell versus RNN\_cell\_forward :

- Note that an RNN cell outputs the hidden state  $a^{(t)}$ .
  - RNN cell** is shown in the figure as the inner box with solid lines
- The function that you'll implement, **rnn\_cell\_forward**, also calculates the prediction  $\hat{y}^{(t)}$ 
  - RNN\_cell\_forward** is shown in the figure as the outer box with dashed lines

## RNN Forward Pass

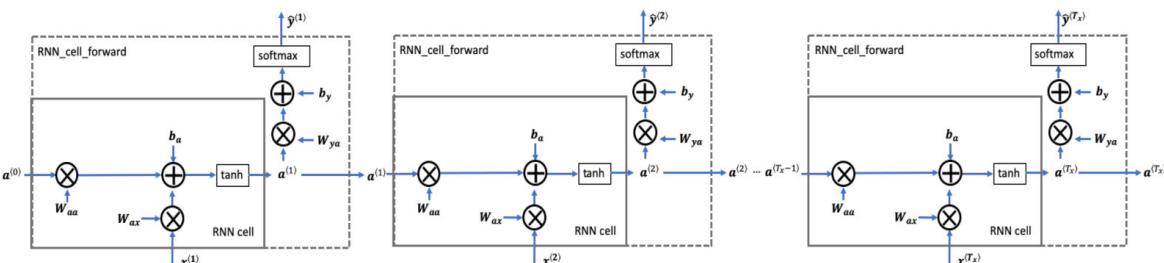


Figure 3: Basic RNN. The input sequence  $x = (x^{(1)}, x^{(2)}, \dots, x^{(T_x)})$  is carried over  $T_x$  time steps. The network outputs  $y = (\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(T_x)})$ .

- A recurrent neural network (RNN) is a repetition of the RNN cell.
  - If your input sequence of data is 10 time steps long, then you will re-use the RNN cell 10 times
- Each cell takes two inputs at each time step:

- $a^{(t-1)}$ : The hidden state from the previous cell
- $x^{(t)}$ : The current time step's input data
- It has two outputs at each time step:
  - A hidden state ( $a^{(t)}$ )
  - A prediction ( $\hat{y}^{(t)}$ )
- The weights and biases ( $W_{aa}, b_a, W_{ax}, b_x$ ) are re-used each time step
  - They are maintained between calls to `rnn_cell_forward` in the 'parameters' dictionary

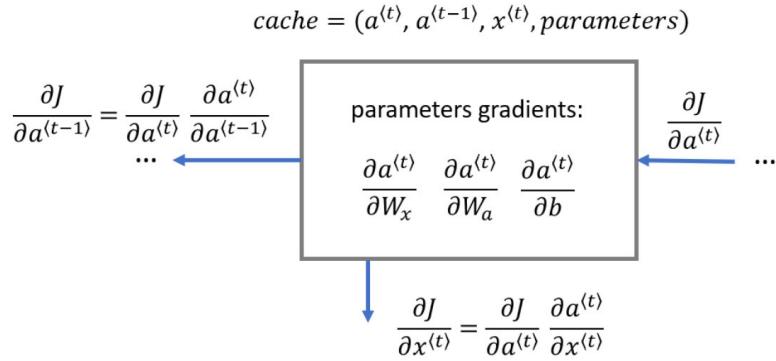
#### Instructions:

- Create a 3D array of zeros,  $a$  of shape  $(n_a, m, T_x)$  that will store all the hidden states computed by the RNN
- Create a 3D array of zeros,  $\hat{y}$ , of shape  $(n_y, m, T_x)$  that will store the predictions
  - Note that in this case,  $T_y = T_x$  (the prediction and input have the same number of time steps)
- Initialize the 2D hidden state `a_next` by setting it equal to the initial hidden state,  $a_0$
- At each time step  $t$ :
  - Get  $x^{(t)}$ , which is a 2D slice of  $x$  for a single time step  $t$ 
    - $x^{(t)}$  has shape  $(n_x, m)$
    - $x$  has shape  $(n_x, m, T_x)$
  - Update the 2D hidden state  $a^{(t)}$ , the prediction  $\hat{y}^{(t)}$  and the cache by running `rnn_cell_forward`
    - $a^{(t)}$  has shape  $(n_a, m)$
  - Store the 2D hidden state in the 3D tensor  $a$ , at the  $t^{th}$  position
    - $a$  has shape  $(n_a, m, T_x)$
  - Store the 2D  $\hat{y}^{(t)}$  prediction in the 3D tensor  $\hat{y}_{pred}$  at the  $t^{th}$  position
    - $\hat{y}^{(t)}$  has shape  $(n_y, m)$
    - $\hat{y}$  has shape  $(n_y, m, T_x)$
  - Append the cache to the list of caches
- Return the 3D tensor  $a$  and  $\hat{y}$ , as well as the list of caches

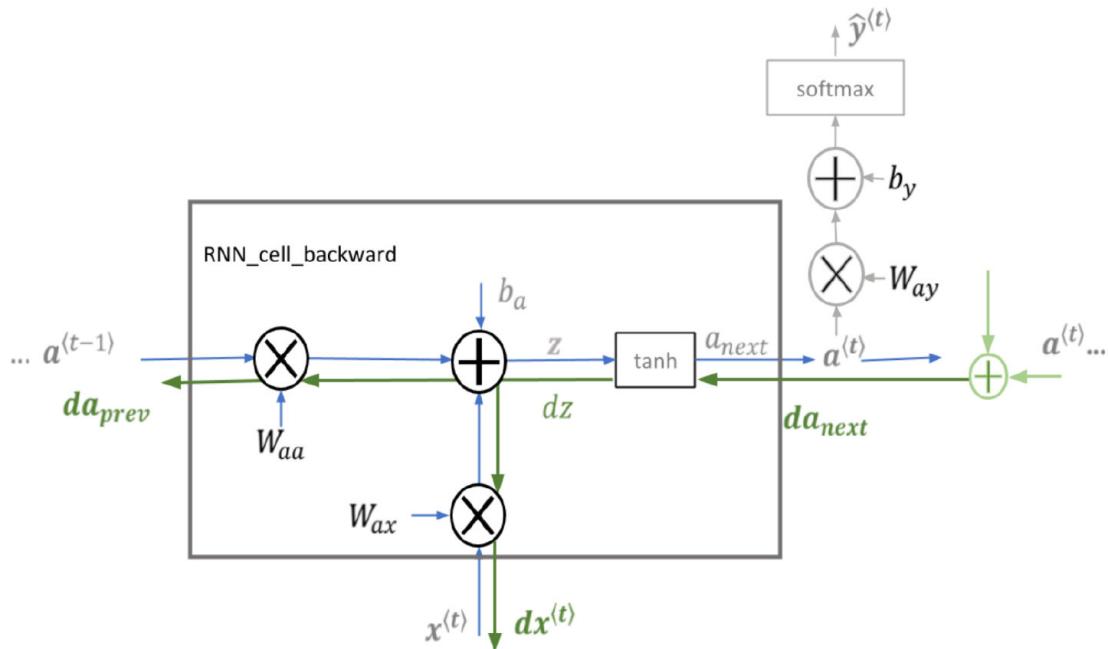
#### What you should remember:

- The recurrent neural network, or RNN, is essentially the repeated use of a single cell.
- A basic RNN reads inputs one at a time, and remembers information through the hidden layer activations (hidden states) that are passed from one time step to the next.
  - The time step dimension determines how many times to re-use the RNN cell
- Each cell takes two inputs at each time step:
  - The hidden state from the previous cell
  - The current time step's input data
- Each cell has two outputs at each time step:
  - A hidden state
  - A prediction

## Basic RNN Backward Pass



**Figure 6:** The RNN cell's backward pass. Just like in a fully-connected neural network, the derivative of the cost function  $J$  backpropagates through the time steps of the RNN by following the chain rule from calculus. Internal to the cell, the chain rule is also used to calculate  $(\frac{\partial J}{\partial W_{ax}}, \frac{\partial J}{\partial W_{aa}}, \frac{\partial J}{\partial b})$  to update the parameters ( $W_{ax}, W_{aa}, b_a$ ). The operation can utilize the cached results from the forward path.



**Figure 7:** This implementation of `rnn\_cell\_backward` does \*\*not\*\* include the output dense layer and softmax which are included in `rnn\_cell\_forward`.  $da_{next}$  is  $\frac{\partial J}{\partial a^{(t)}}$  and includes loss from previous stages and current stage output logic. The addition shown in green will be part of your implementation of `rnn_backward`.

## Equations

To compute `rnn_cell_backward`, you can use the following equations. Here,  $*$  denotes element-wise multiplication while the absence of a symbol indicates matrix multiplication.

$$a^{\langle t \rangle} = \tanh(W_{ax}x^{\langle t \rangle} + W_{aa}a^{\langle t-1 \rangle} + b_a) \quad (-)$$

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x) \quad (-)$$

$$dtanh = da_{next} * (1 - \tanh^2(W_{ax}x^{\langle t \rangle} + W_{aa}a^{\langle t-1 \rangle} + b_a)) \quad (0)$$

$$dW_{ax} = dtanh \cdot x^{\langle t \rangle T} \quad (1)$$

$$dW_{aa} = dtanh \cdot a^{\langle t-1 \rangle T} \quad (2)$$

$$db_a = \sum_{batch} dtanh \quad (3)$$

$$dx^{\langle t \rangle} = {W_{ax}}^T \cdot dtanh \quad (4)$$

$$da_{prev} = {W_{aa}}^T \cdot dtanh \quad (5)$$

$$a^{(t)} = \tanh(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + ba)$$

$$\frac{d\tanh(x)}{dx} = 1 - \tanh^2(x)$$

$$\underline{a^{(t+1)}} = \underline{\tanh(\underline{z^{(t)}})}$$

element-wise

Equation (0)  $d\tanh = da_{next} * (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + ba))$

$$d\tanh = dz = \frac{dJ}{da^{(t+1)}} \cdot \frac{da^{(t+1)}}{dz} = da^{(t+1)} * (1 - \tanh^2(z))$$

= - - - - -

Equation (1)  $dW_{ax} = d\tanh \cdot x^{(t)} T$  Matrix Cookbook

$$dW_{ax} = \frac{dJ}{dz} \cdot \frac{dz}{dW_{ax}} = d\tanh \cdot x^{(t)} T$$

Equation (2)  $dW_{aa} = d\tanh \cdot a^{(t-1)} T$

$$dW_{aa} = \frac{dJ}{dz} \cdot \frac{dz}{dW_{aa}} = d\tanh \cdot a^{(t-1)} T$$

Equation (3)  $dba = \sum_{batch} d\tanh$  vectorization,  $\frac{1}{m}$  is omitted

$$dba = \frac{dJ}{dz} \frac{dz}{dbba} = \sum_{batch} dz = \sum_{batch} d\tanh \quad \underline{db^{(i)}} = \underline{dz^{(i)}} \text{ sum over}$$

Equation (4)  $dx^{(t)} = W_{ax}^T \cdot d\tanh$

$$dx^{(t)} = \frac{dJ}{dz} \frac{dz}{dW_{ax}} = W_{ax}^T \cdot d\tanh \quad \text{order?}$$

Equation (5)  $da_{prev} = W_{aa}^T \cdot d\tanh$  Transpose?

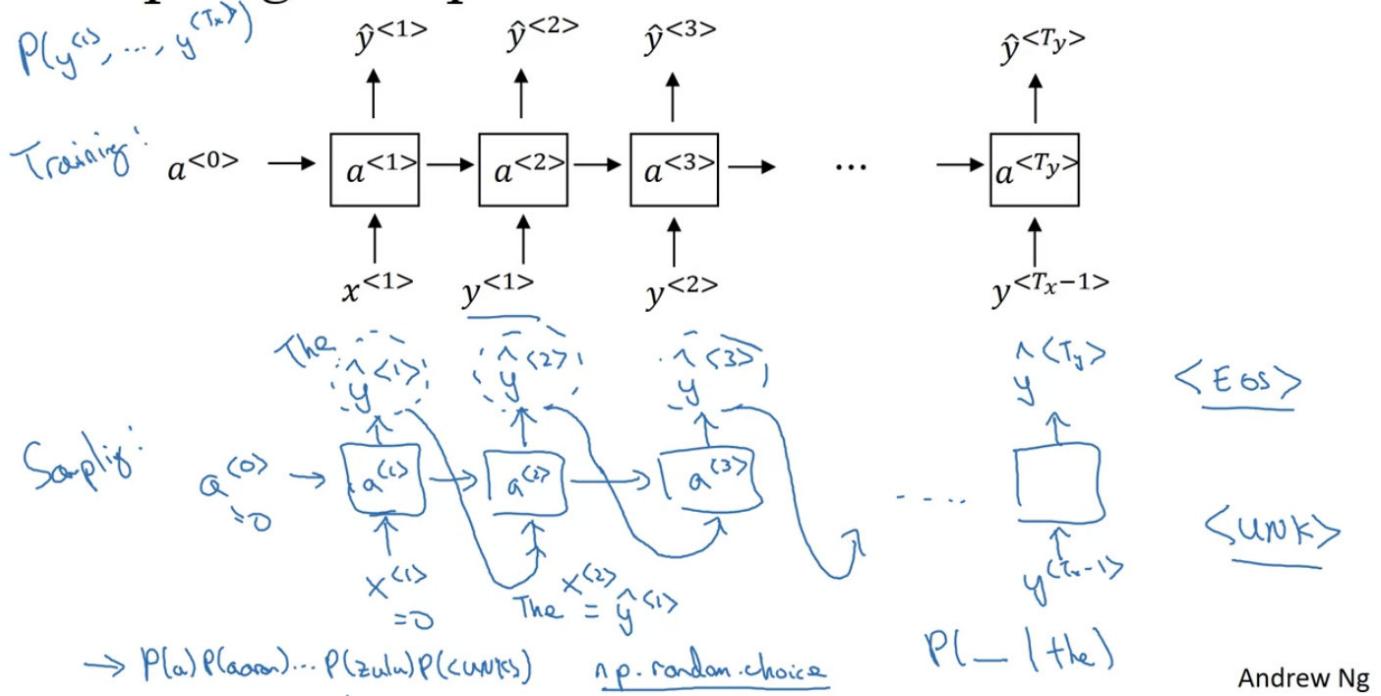
$$da_{prev} = \frac{dJ}{dz} \frac{dz}{dW_{aa}} = W_{aa}^T \cdot d\tanh$$

Check note about Matrix vector derivatives by CS231n:

[vecDerivs.pdf](#)

## Sampling Novel Sequences

# Sampling a sequence from a trained RNN



what you want to do:

1. Randomly sample first word  $\hat{y}^{<1>}$  according to the first softmax distribution.

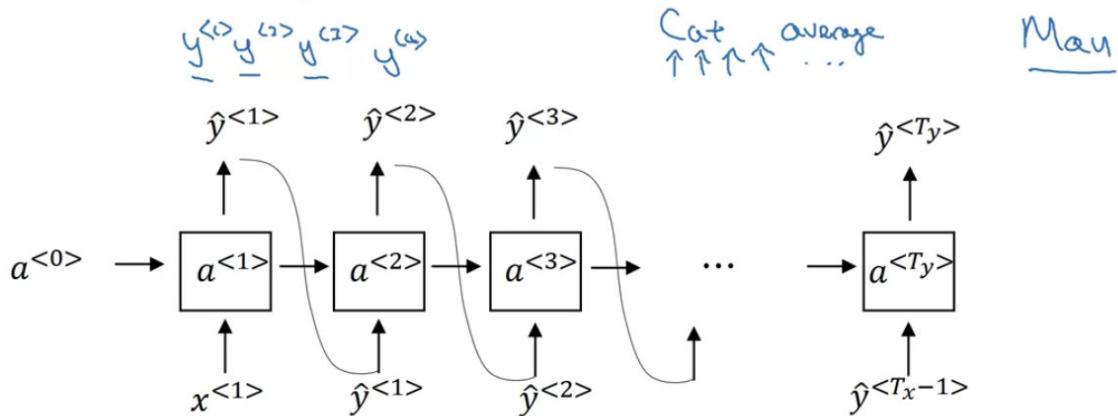
```
np.random.choice() # To sample according to distribution
```

2. Going to next time step, take the  $\hat{y}^{<1>}$  which is just sampled as the input to the next time step and sample  $\hat{y}^{<2>}$ .
3. Repeat 1 → 2 until get the last time step (such as when  $\langle \text{EOS} \rangle$  token is sampled or just decide how many words to sample).

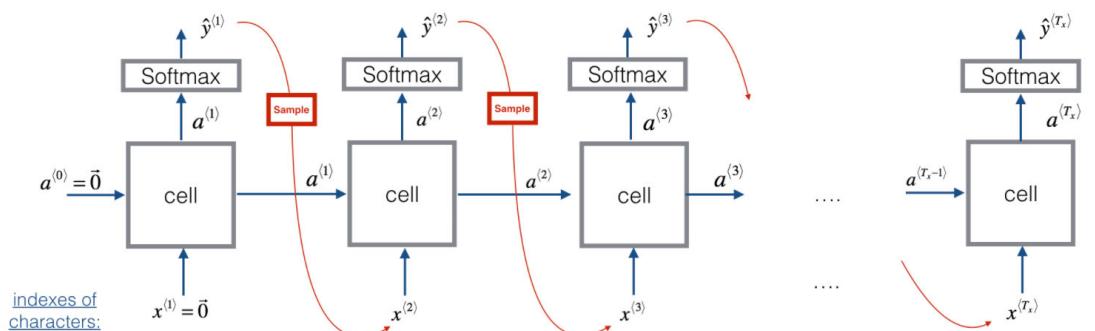
# Character-level language model

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

$\rightarrow \text{Vocabulary} = [a, b, c, \dots, z, \cup, \circ, \cdot, ;, :, 0, \dots, 9, A, \dots, Z]$



Andrew Ng



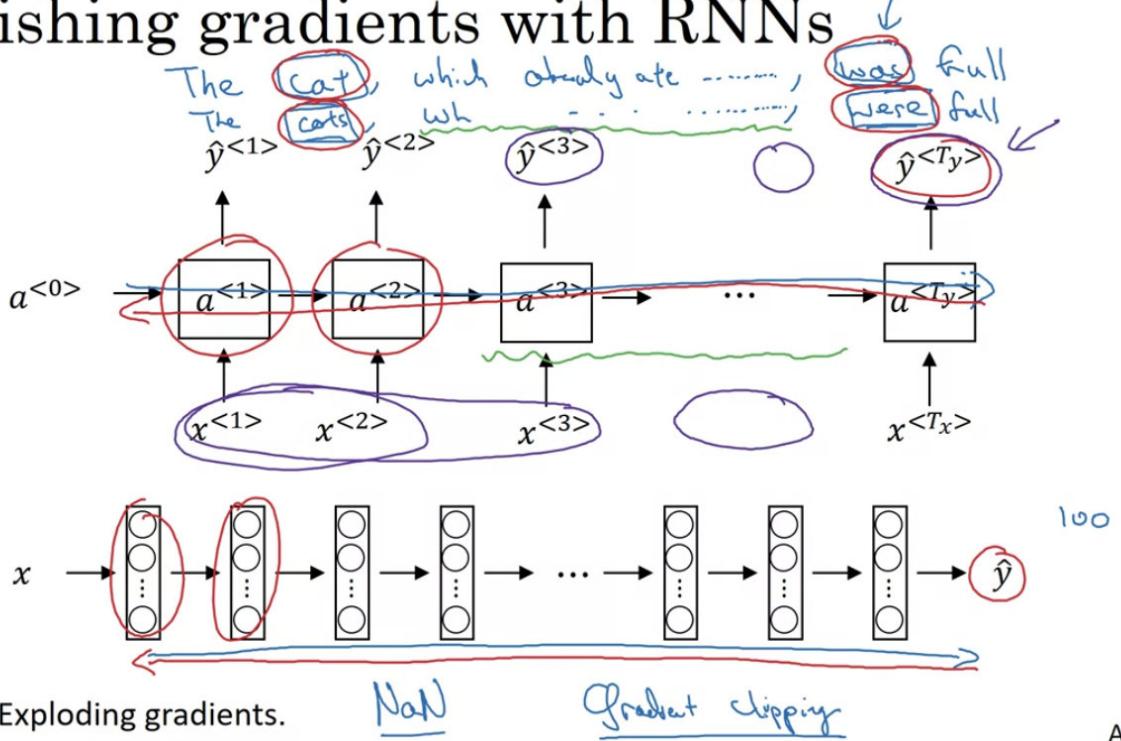
**Figure 3:** In this picture, you can assume the model is already trained. You pass in  $x^{(1)} = \vec{0}$  at the first time-step, and have the network sample one character at a time.

## Pros and cons of Character-level language model

1. Don't have to worry about unknown word tokens.
  2. End up with much longer sequences (computational expensive).

## Vanishing Gradients with RNNs

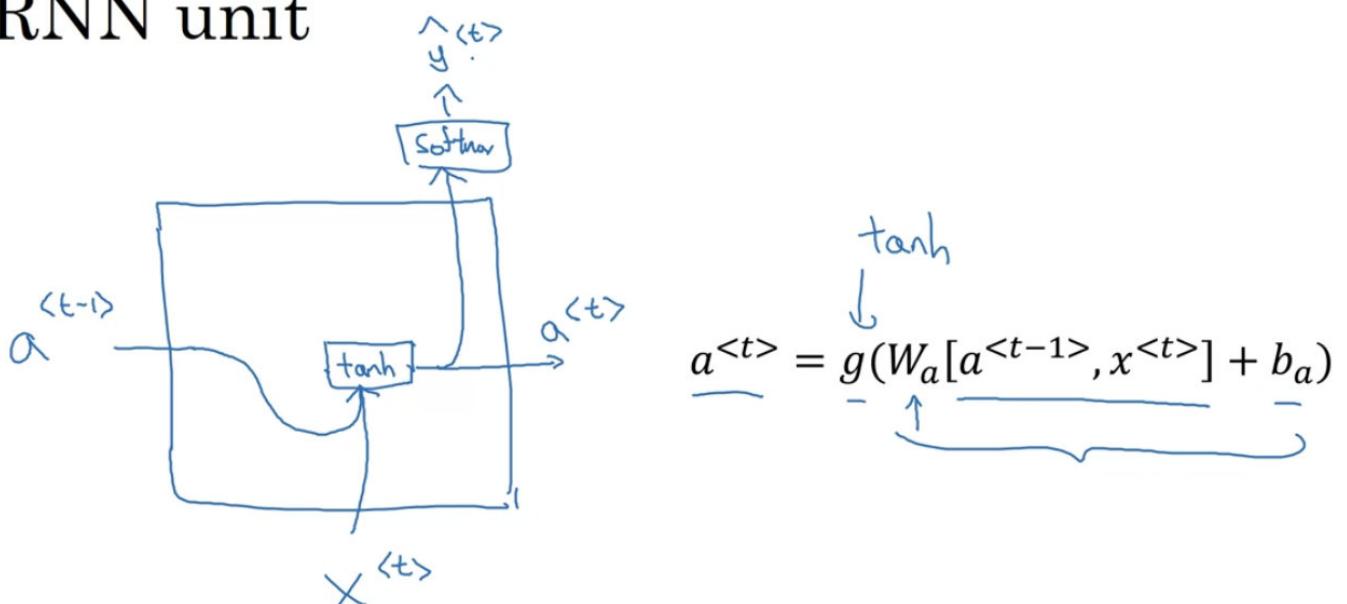
# Vanishing gradients with RNNs



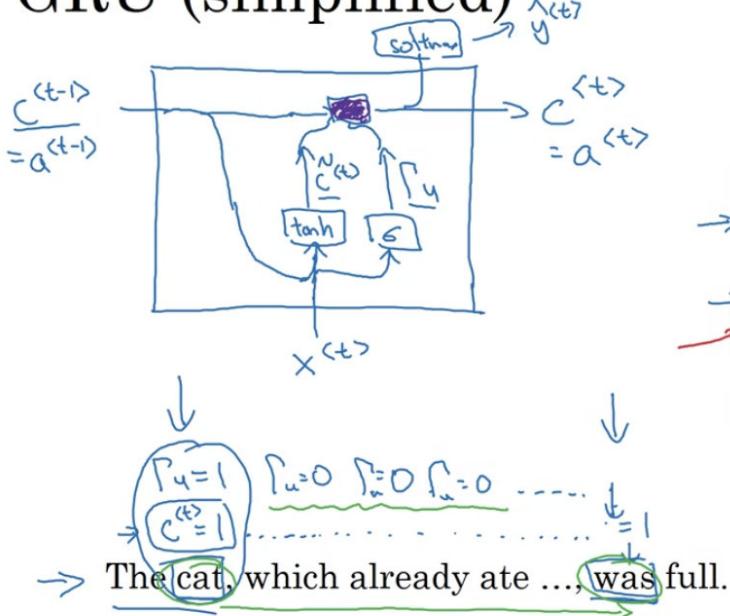
Because of vanishing gradients, the basic RNN model has many **local influences**, meaning that the output  $y^{<t>}$  is mainly influenced by values close to  $y^{<t>}$ , so is difficult for it to be strongly influenced by an input that was **very early** in the sequence.

## Gated Recurrent Unit (GRU)

### RNN unit



# GRU (simplified)



→ The cat, which already ate ... was full.



$C$  = memory cell

$$\rightarrow C^{(t)} = a^{(t)}$$

$$\rightarrow \tilde{C}^{(t)} = \tanh(W_c[c^{(t-1)}, x^{(t)}] + b_c)$$

$$\rightarrow \Gamma_u = \sigma(W_u[c^{(t-1)}, x^{(t)}] + b_u)$$

$$\left\{ \begin{array}{l} C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)} \\ \text{“update”} \\ \text{element-wise} \\ \text{Gate} \end{array} \right.$$

$\Gamma_u = 0.00001$

[Cho et al., 2014. On the properties of neural machine translation: Encoder-decoder approaches]

[Chung et al., 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling]

Andrew Ng

**c:memory cell** will provide a bit of to remember the important information.

The memory cell will have some value  $c^{(t)}$  at time  $t$ . The GRU will output an activation value  $a^{(t)}$  that is equal to  $c^{(t)}$ .

$$c^{(t)} = a^{(t)}$$

At every time step, we are going to **consider** overwriting the memory cell with a **candidate**  $\tilde{c}^{(t)}$  for replacing  $c^{(t)}$ .

$$\tilde{c}^{(t)} = \tanh(W_c[c^{(t-1)}, x^t] + b_c)$$

The key idea of GRU is to have a **update gate**

$$\Gamma_u = \sigma(W_u[c^{(t-1)}, x^t] + b_u)$$

The job of gate  $\Gamma_u$  is to decide when do you update memory cell.

The specific equation for GRU is

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * c^{(t-1)}$$

Because the gate is easily to be set to 0, it is very good at **maintaining** the value for the cell. And because the gate can be so close to 0, it does not suffer from vanishing gradient problems.

**Full GRU**

# Full GRU

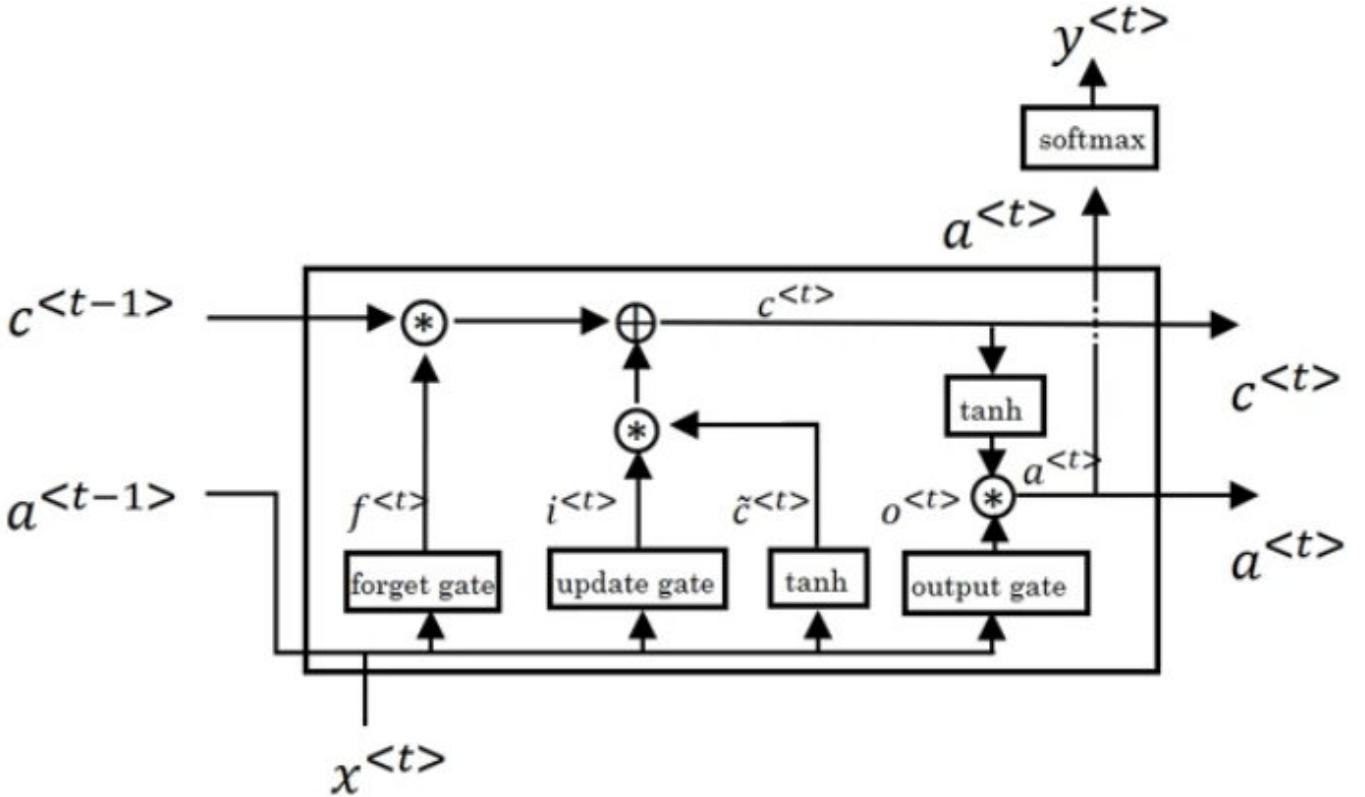
$$\begin{aligned} \tilde{c}^{<t>} &= \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c) \\ u \quad \left\{ \begin{array}{l} \Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ r \quad \left\{ \begin{array}{l} \Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \end{array} \right. \end{array} \right. & \text{LSTM} \\ h \quad c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>} \end{aligned}$$

The cat, which ate already, was full.

For full GRU, there is an additional gate  $\Gamma_r$ .

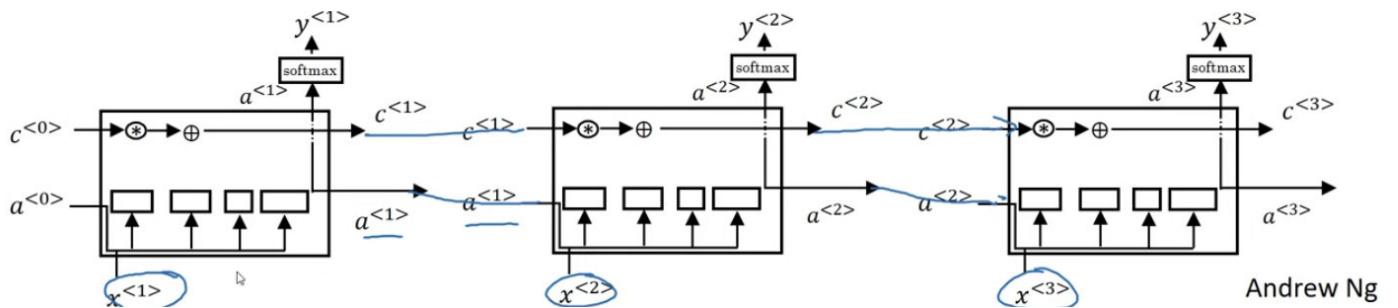
$$\begin{aligned} a^{<t>} &= c^{<t>} \\ \tilde{c}^{<t>} &= \tanh(W_c[\Gamma_r * c^{<t-1>}, x^t] + b_c) \\ \Gamma_r &= \sigma(W_r[c^{<t-1>}, x^t] + b_r) \\ \Gamma_u &= \sigma(W_u[c^{<t-1>}, x^t] + b_u) \\ c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>} \end{aligned}$$

## LSTM (Long Short Term Memory) Unit



Three gates: **update gate**  $\Gamma_u$ , **forget gate**  $\Gamma_f$ , **output gate**  $\Gamma_o$

$$\begin{aligned}
 \tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\
 \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \leftarrow \text{update} \\
 \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \leftarrow \text{forget} \\
 \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \leftarrow \text{output} \\
 c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \\
 a^{<t>} &= \Gamma_o * \tanh(c^{<t>})
 \end{aligned}$$

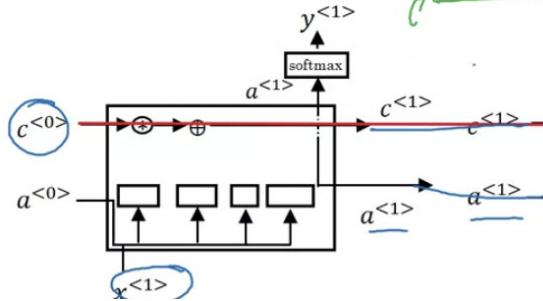


Notice that, there is a line at the top shows that so long as you set the **forget** and the **update** gate appropriately, it is **relatively easy** for the LSTM to have some value  $c_0$  to be passed all the way the right to have  $c_3 = c_0$ . This is why LSTM, as well as the GRU, is very good at memorizing certain values even for a long time.

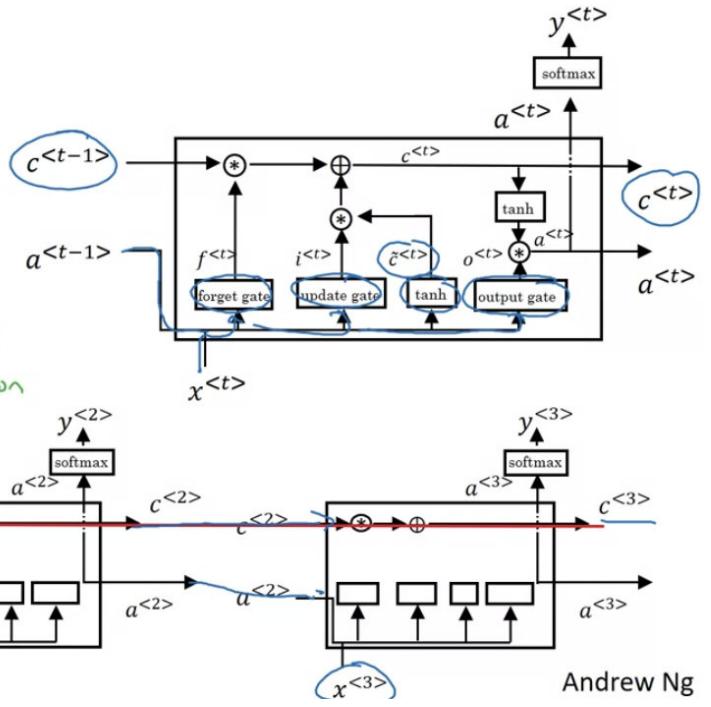
### Peephole LSTM

# LSTM in pictures

$$\begin{aligned}\tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\ \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\ c^{<t>} &= \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \\ a^{<t>} &= \Gamma_o * \tanh c^{<t>}\end{aligned}$$



peephole  
connection

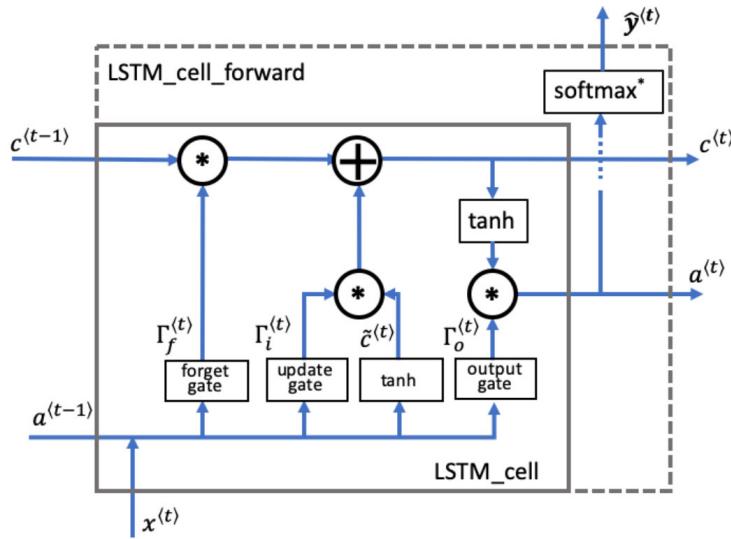


## LSTM and GRU

- The advantage of the GRU is that it is a **simpler** model, so it is actually easier to build a much bigger network.
- LSTM is more **powerful** and more **effective** since it has three gates.

## Details of LSTM

The following figure shows the operations of an LSTM cell:



**Figure 4: LSTM cell.** This tracks and updates a "cell state," or memory variable  $c^{(t)}$  at every time step, which can be different from  $a^{(t)}$ . Note, the softmax includes a dense layer and softmax.

## Forget gate $\Gamma_f$

- Let's assume you are reading words in a piece of text, and plan to use an LSTM to keep track of grammatical structures, such as whether the subject is singular ("puppy") or plural ("puppies").
- If the subject changes its state (from a singular word to a plural word), the memory of the previous state becomes outdated, so you'll "forget" that outdated state.
- The "forget gate" is a tensor containing values between 0 and 1.
  - If a unit in the forget gate has a value close to 0, the LSTM will "forget" the stored state in the corresponding unit of the previous cell state.
  - If a unit in the forget gate has a value close to 1, the LSTM will mostly remember the corresponding value in the stored state.

## Equation

$$\Gamma_f^{(t)} = \sigma(\mathbf{W}_f[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_f) \quad (1)$$

### Explanation of the equation:

- $\mathbf{W}_f$  contains weights that govern the forget gate's behavior.
- The previous time step's hidden state  $a^{(t-1)}$  and current time step's input  $x^{(t)}$  are concatenated together and multiplied by  $\mathbf{W}_f$ .
- A sigmoid function is used to make each of the gate tensor's values  $\Gamma_f^{(t)}$  range from 0 to 1.
- The forget gate  $\Gamma_f^{(t)}$  has the same dimensions as the previous cell state  $c^{(t-1)}$ .
- This means that the two can be multiplied together, element-wise.
- Multiplying the tensors  $\Gamma_f^{(t)} * c^{(t-1)}$  is like applying a mask over the previous cell state.
- If a single value in  $\Gamma_f^{(t)}$  is 0 or close to 0, then the product is close to 0.
  - This **keeps** the information stored in the corresponding unit in  $c^{(t-1)}$  **from** being remembered for the next time step. **forget**
- Similarly, if one value is close to 1, the product is close to the original value in the previous cell state.
  - The LSTM will **keep** the information from the corresponding unit of  $c^{(t-1)}$ , to be used in the next time step. **remember**

## Candidate value $\tilde{c}^{(t)}$

- The candidate value is a tensor containing information from the **current** time step that **may** be stored in the current cell state  $c^{(t)}$ .
- The parts of the candidate value that get passed on depend on the update gate.
- The candidate value is a tensor containing values that range from -1 to 1.
- The tilde "~" is used to differentiate the candidate  $\tilde{c}^{(t)}$  from the cell state  $c^{(t)}$ .

## Equation

$$\tilde{\mathbf{c}}^{(t)} = \tanh(\mathbf{W}_c[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_c) \quad (3)$$

## Explanation of the equation

- The tanh function produces values between -1 and 1.

## Update gate $\Gamma_i$

- You use the update gate to decide what aspects of the candidate  $\tilde{\mathbf{c}}^{(t)}$  to add to the cell state  $c^{(t)}$ .
- The update gate decides what parts of a "candidate" tensor  $\tilde{\mathbf{c}}^{(t)}$  are passed onto the cell state  $c^{(t)}$ .
- The update gate is a tensor containing values between 0 and 1.
  - When a unit in the update gate is close to 1, it allows the value of the candidate  $\tilde{\mathbf{c}}^{(t)}$  to be passed onto the hidden state  $\mathbf{c}^{(t)}$
  - When a unit in the update gate is close to 0, it prevents the corresponding value in the candidate from being passed onto the hidden state.
- Notice that the subscript "i" is used and not "u", to follow the convention used in the literature.

## Equation

$$\Gamma_i^{(t)} = \sigma(\mathbf{W}_i[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_i) \quad (2)$$

## Explanation of the equation

- Similar to the forget gate, here  $\Gamma_i^{(t)}$ , the sigmoid produces values between 0 and 1.
- The update gate is multiplied element-wise with the candidate, and this product ( $\Gamma_i^{(t)} * \tilde{\mathbf{c}}^{(t)}$ ) is used in determining the cell state  $c^{(t)}$ .

## Cell state $\mathbf{c}^{(t)}$

- The cell state is the "memory" that gets passed onto future time steps.
- The new cell state  $\mathbf{c}^{(t)}$  is a **combination** of the previous cell state and the candidate value.

## Equation

$$\mathbf{c}^{(t)} = \Gamma_f^{(t)} * \mathbf{c}^{(t-1)} + \Gamma_i^{(t)} * \tilde{\mathbf{c}}^{(t)} \quad (4)$$

## Explanation of equation

- The previous cell state  $\mathbf{c}^{(t-1)}$  is adjusted (weighted) by the forget gate  $\Gamma_f^{(t)}$
- and the candidate value  $\tilde{\mathbf{c}}^{(t)}$ , adjusted (weighted) by the update gate  $\Gamma_i^{(t)}$

## Output gate $\Gamma_o$

- The output gate decides what gets sent as the prediction (output) of the time step.
- The output gate is like the other gates, in that it contains values that range from 0 to 1.

## Equation

$$\Gamma_o^{(t)} = \sigma(\mathbf{W}_o[\mathbf{a}^{(t-1)}, \mathbf{x}^{(t)}] + \mathbf{b}_o) \quad (5)$$

## Explanation of the equation

- The output gate is determined by the previous hidden state  $\mathbf{a}^{(t-1)}$  and the current input  $\mathbf{x}^{(t)}$
- The sigmoid makes the gate range from 0 to 1.

## Hidden state $\mathbf{a}^{(t)}$

- The hidden state gets passed to the LSTM cell's next time step.
- It is used to determine the three gates ( $\Gamma_f, \Gamma_u, \Gamma_o$ ) of the next time step.
- The hidden state is also used for the prediction  $y^{(t)}$ .

## Equation

$$\mathbf{a}^{(t)} = \Gamma_o^{(t)} * \tanh(\mathbf{c}^{(t)}) \quad (6)$$

## Explanation of equation

- The hidden state  $\mathbf{a}^{(t)}$  is determined by the cell state  $\mathbf{c}^{(t)}$  in combination with the output gate  $\Gamma_o$ .
- The cell state state is passed through the **tanh** function to rescale values between -1 and 1.
- The output gate acts like a "mask" that either preserves the values of  $\tanh(\mathbf{c}^{(t)})$  or keeps those values from being included in the hidden state  $\mathbf{a}^{(t)}$

## Prediction $\mathbf{y}_{pred}^{(t)}$

- The prediction in this use case is a classification, so you'll use a softmax.

## Equation

$$\mathbf{y}_{pred}^{(t)} = \text{softmax}(\mathbf{W}_y \mathbf{a}^{(t)} + \mathbf{b}_y)$$

## Forward Pass for LSTM

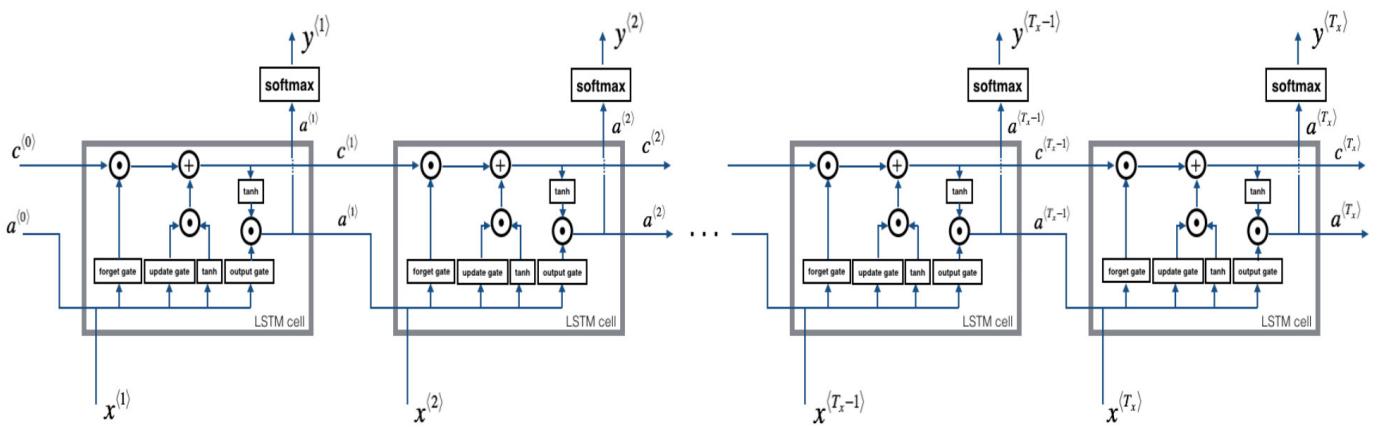


Figure 5: LSTM over multiple time steps.

## Instructions

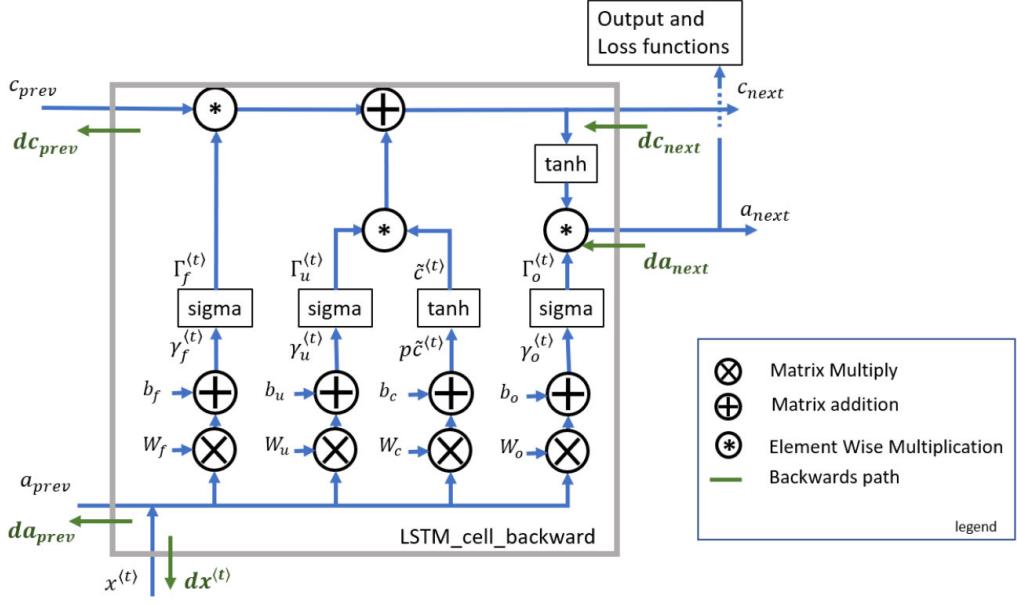
- Get the dimensions  $n_x, n_a, n_y, m, T_x$  from the shape of the variables
- Initialize the 3D tensors  $a, c$  and  $y$

- $a$ : hidden state, shape  $(n_a, m, T_x)$
- $c$ : cell state, shape  $(n_a, m, T_x)$
- $y$ : prediction, shape  $(n_y, m, T_x)$  (Note that  $T_y = T_x$  in this example)
- **Note** Setting one variable equal to the other is a "copy by reference". In other words, don't do ` $c = a$ ', otherwise both these variables point to the same underlying variable.
- Initialize the 2D tensor  $a^{(t)}$ 
  - $a^{(t)}$  stores the hidden state for time step  $t$ .
  - $a^{(0)}$ , the initial hidden state at time step 0, is passed in when calling the function.
  - $a^{(t)}$  and  $a^{(0)}$  represent a single time step, so they both have the shape  $(n_a, m)$
  - Initialize  $a^{(t)}$  by setting it to the initial hidden state ( $a^{(0)}$ ) that is passed into the function.
- Initialize  $c^{(t)}$  with zeros.
  - $c^{(t)}$  represents a single time step, so its shape is  $(n_a, m)$
  - **Note:** create `c_next` as its own variable with its own location in memory. Do not initialize it as a slice of the 3D tensor  $c$ . In other words, don't do `c_next = c[:, :, 0]`.
- For each time step, do the following:
  - From the 3D tensor  $x$ , get a 2D slice  $x^{(t)}$  at time step  $t$
  - Call the `lstm_cell_forward` function that you defined previously, to get the hidden state, cell state, prediction, and cache
  - Store the hidden state, cell state and prediction (the 2D tensors) inside the 3D tensors
  - Append the cache to the list of caches

#### What you should remember:

- An LSTM is similar to an RNN in that they both use hidden states to pass along information, but an LSTM also uses a cell state, which is like a long-term memory, to help deal with the issue of vanishing gradients
- An LSTM cell consists of a cell state, or long-term memory, a hidden state, or short-term memory, along with 3 gates that constantly update the relevancy of its inputs:
  - A **forget** gate, which decides which input units should be remembered and passed along. It's a tensor with values between 0 and 1.
    - If a unit has a value close to 0, the LSTM will "forget" the stored state in the previous cell state.
    - If it has a value close to 1, the LSTM will mostly remember the corresponding value.
  - An **update** gate, again a tensor containing values between 0 and 1. It decides on what information to throw away, and what new information to add.
    - When a unit in the update gate is close to 1, the value of its candidate is passed on to the hidden state.
    - When a unit in the update gate is close to 0, it's prevented from being passed onto the hidden state.
  - And an **output** gate, which decides what gets sent as the output of the time step

## LSTM Backward Pass



**Figure 8:** LSTM Cell Backward. Note the output functions, while part of the `lstm\_cell\_forward`, are not included in `lstm\_cell\_backward`

## Gate Derivatives

Note the location of the gate derivatives ( $\gamma_{..}$ ) between the dense layer and the activation function (see graphic above). This is convenient for computing parameter derivatives in the next step.

$$d\gamma_o^{(t)} = da_{next} * \tanh(c_{next}) * \Gamma_o^{(t)} * \left(1 - \Gamma_o^{(t)}\right) \quad (7)$$

$$dp\tilde{c}^{(t)} = \left( dc_{next} * \Gamma_u^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \Gamma_u^{(t)} * da_{next} \right) * \left(1 - (\tilde{c}^{(t)})^2\right) \quad (8)$$

$$d\gamma_u^{(t)} = \left( dc_{next} * \tilde{c}^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \tilde{c}^{(t)} * da_{next} \right) * \Gamma_u^{(t)} * \left(1 - \Gamma_u^{(t)}\right) \quad (9)$$

$$d\gamma_f^{(t)} = \left( dc_{next} * c_{prev} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * c_{prev} * da_{next} \right) * \Gamma_f^{(t)} * \left(1 - \Gamma_f^{(t)}\right) \quad (10)$$

## 3. Parameter Derivatives

$$dW_f = d\gamma_f^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^T \quad (11)$$

$$dW_u = d\gamma_u^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^T \quad (12)$$

$$dW_c = dp\tilde{c}^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^T \quad (13)$$

$$dW_o = d\gamma_o^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^T \quad (14)$$

To calculate  $db_f, db_u, db_c, db_o$  you just need to sum across all 'm' examples (axis= 1) on  $d\gamma_f^{(t)}, d\gamma_u^{(t)}, dp\tilde{c}^{(t)}, d\gamma_o^{(t)}$  respectively. Note that you should have the `keepdims = True` option.

$$db_f = \sum_{batch} d\gamma_f^{(t)} \quad (15)$$

$$db_u = \sum_{batch} d\gamma_u^{(t)} \quad (16)$$

$$db_c = \sum_{batch} d\gamma_c^{(t)} \quad (17)$$

$$db_o = \sum_{batch} d\gamma_o^{(t)} \quad (18)$$

Finally, you will compute the derivative with respect to the previous hidden state, previous memory state, and input.

$$da_{prev} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\tilde{c}^{(t)} + W_o^T d\gamma_o^{(t)} \quad (19)$$

Here, to account for concatenation, the weights for equations 19 are the first n\_a, (i.e.  $W_f = W_f[:, :n_a]$  etc...)

$$dc_{prev} = dc_{next} * \Gamma_f^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \Gamma_f^{(t)} * da_{next} \quad (20)$$

$$dx^{(t)} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\tilde{c}^{(t)} + W_o^T d\gamma_o^{(t)} \quad (21)$$

where the weights for equation 21 are from n\_a to the end, (i.e.  $W_f = W_f[:, n_a :]$  etc...)

## Bidirectional RNN (BRNN)

# Bidirectional RNN (BRNN)



Andrew Ng

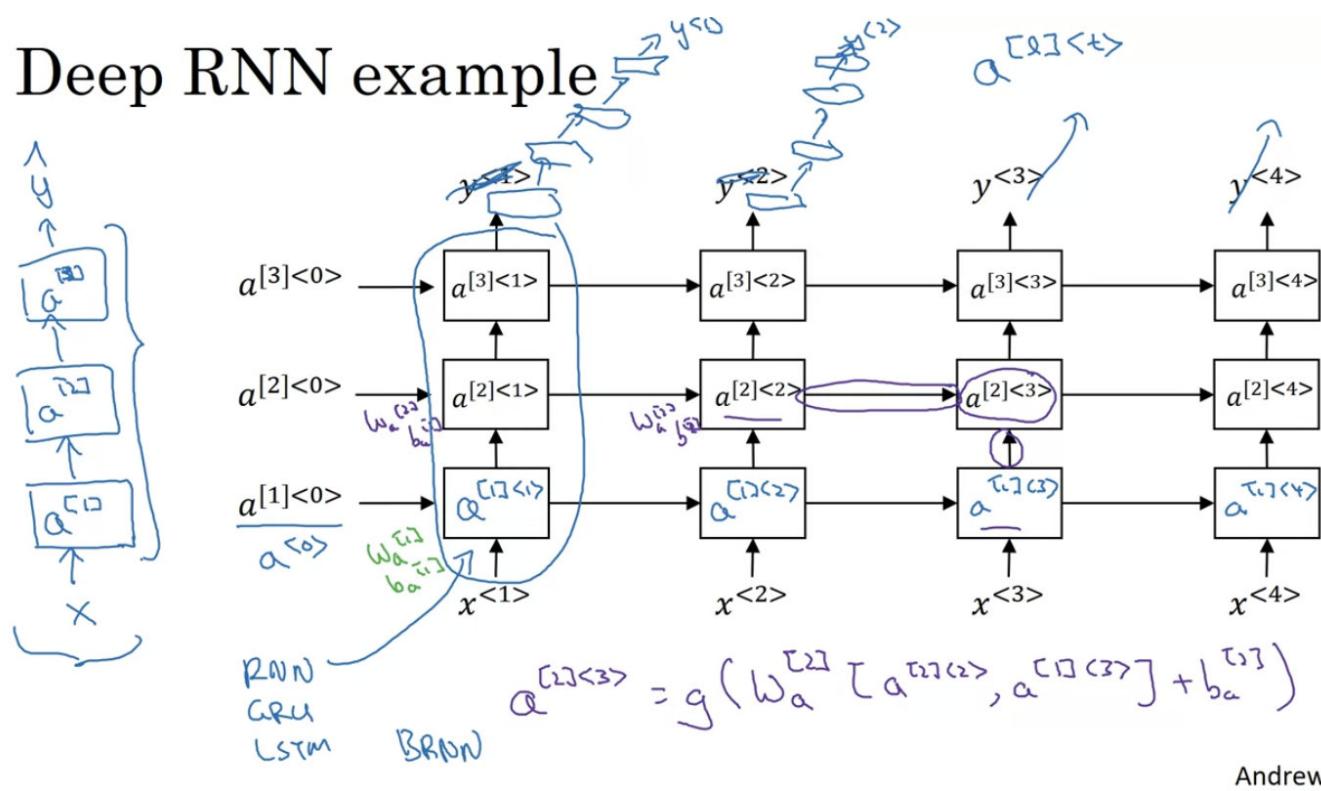
$$\hat{y}^{(t)} = g(W_y[\vec{a}^{(t)}, \vec{a}^{(t)}] + b_y)$$

The blocks can be not just the standard RNN block but they can also be GRU blocks and LSTM blocks.

The disadvantage of the BRNN is that you need the entire sequence of data before you can make predictions anywhere.

## Deep RNNs

### Deep RNN example



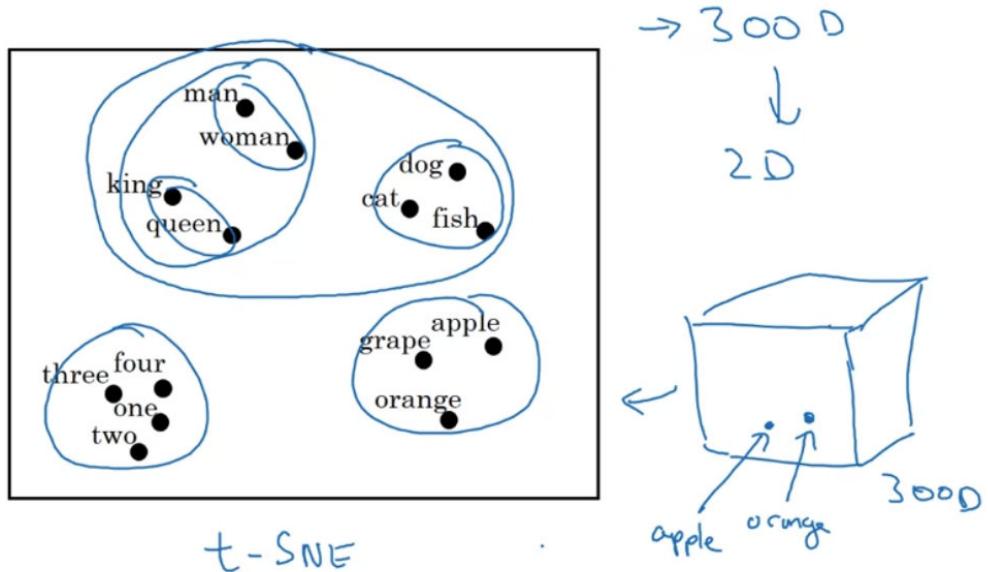
Andrew Ng

## Word Embeddings

Representing words by "features".

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(7157) | Apple<br>(456) | Orange<br>(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             |

# Visualizing word embeddings



## Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (Or download pre-trained embedding online.)
2. Transfer embedding to new task with smaller training set.
3. Optional: Continue to finetune the word embeddings with new data.

The terms **encoding** and **embedding** are used somewhat interchangeably.

## Analogies

One of the most fascinating properties of word embeddings is that they can also help with **analogy reasoning**.

## Analogies

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(7157) | Apple<br>(456) | Orange<br>(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             |

$$\begin{aligned}
 & \text{Man} \rightarrow \text{Woman} \quad \text{vs} \quad \text{King} \rightarrow ? \text{ Queen} \\
 & \underline{\mathbf{e}_{\text{man}}} - \underline{\mathbf{e}_{\text{woman}}} \approx \underline{\mathbf{e}_{\text{King}}} - \underline{\mathbf{e}_{?}} \\
 & \underline{\mathbf{e}_{\text{man}}} - \underline{\mathbf{e}_{\text{woman}}} \approx \underline{\mathbf{e}_{\text{King}}} - \underline{\mathbf{e}_{\text{Queen}}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}
 \end{aligned}$$

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

Find word  $w$ :

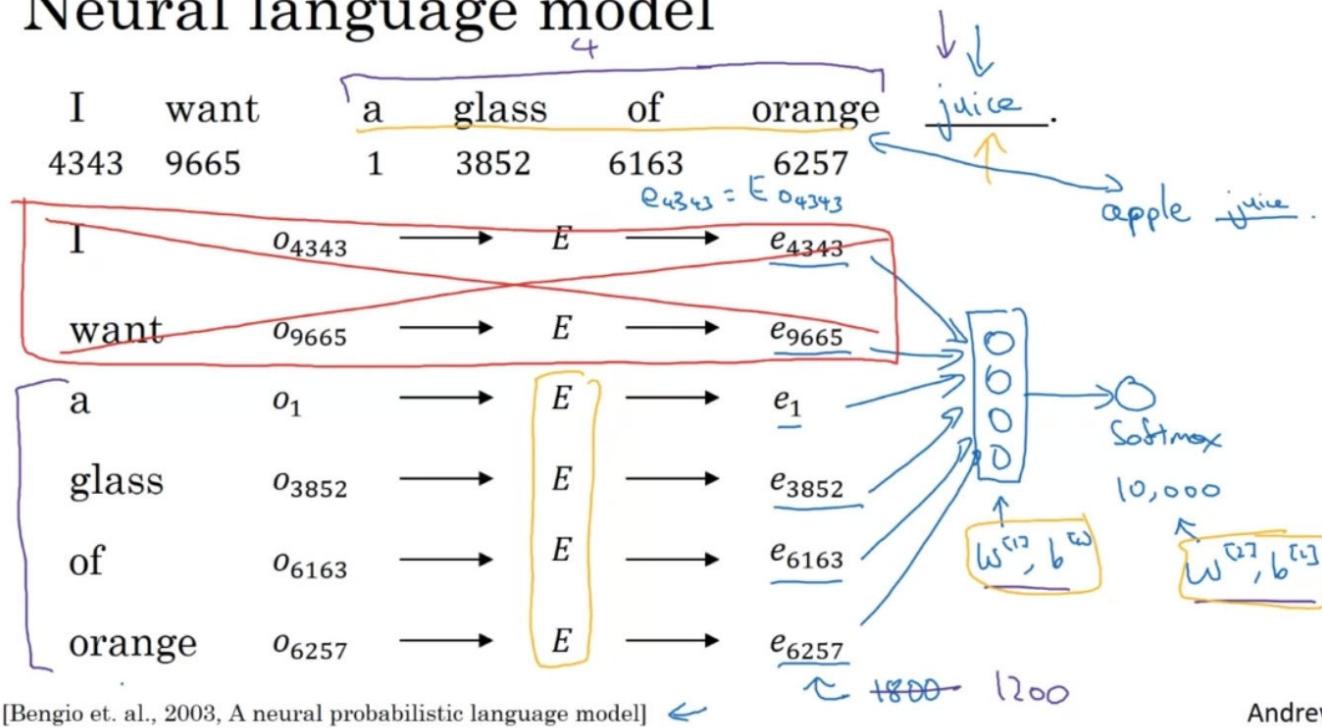
$$\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

## Cosine similarity

$$\text{sim}(u, v) = \frac{u \cdot v}{\|u\|_2 \|v\|_2} = \cos(\theta)$$

- $u \cdot v$  is the dot product (or inner product) of two vectors
- $\|u\|_2$  is the norm (or length) of the vector  $u$
- $\theta$  is the angle between  $u$  and  $v$ .
- The cosine similarity depends on the angle between  $u$  and  $v$ .
  - If  $u$  and  $v$  are very similar, their cosine similarity will be close to 1.
  - If they are dissimilar, the cosine similarity will take a smaller value.

## Neural language model



## Embedding Matrix

$$E = \begin{bmatrix} a & aaron & \dots & orange & \dots & zulu & \langle UNK \rangle \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \in \mathbb{R}^{300 \times 10,000}$$

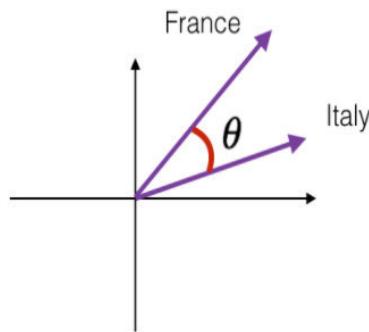
The word orange is the 6257th token in the vocabulary:

$$orange = O_{6257} = \begin{bmatrix} 0 \\ \vdots \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{10,000}$$

$O_{6257}$  is the one-hot vector of the word orange. Notice that:

$$E \cdot O_{6257} = e_{6257} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^{300}$$

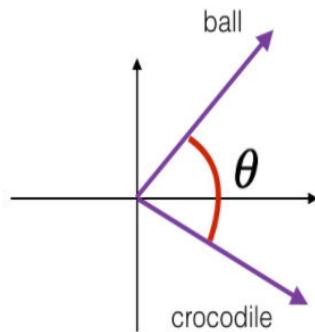
## Learning Word Embeddings



France and Italy are quite similar

$\theta$  is close to  $0^\circ$

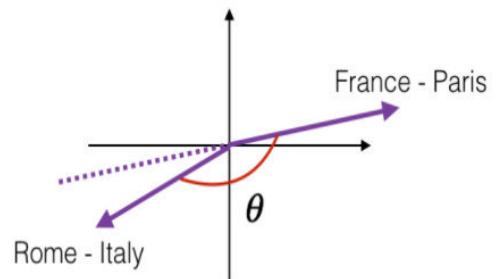
$$\cos(\theta) \approx 1$$



ball and crocodile are not similar

$\theta$  is close to  $90^\circ$

$$\cos(\theta) \approx 0$$

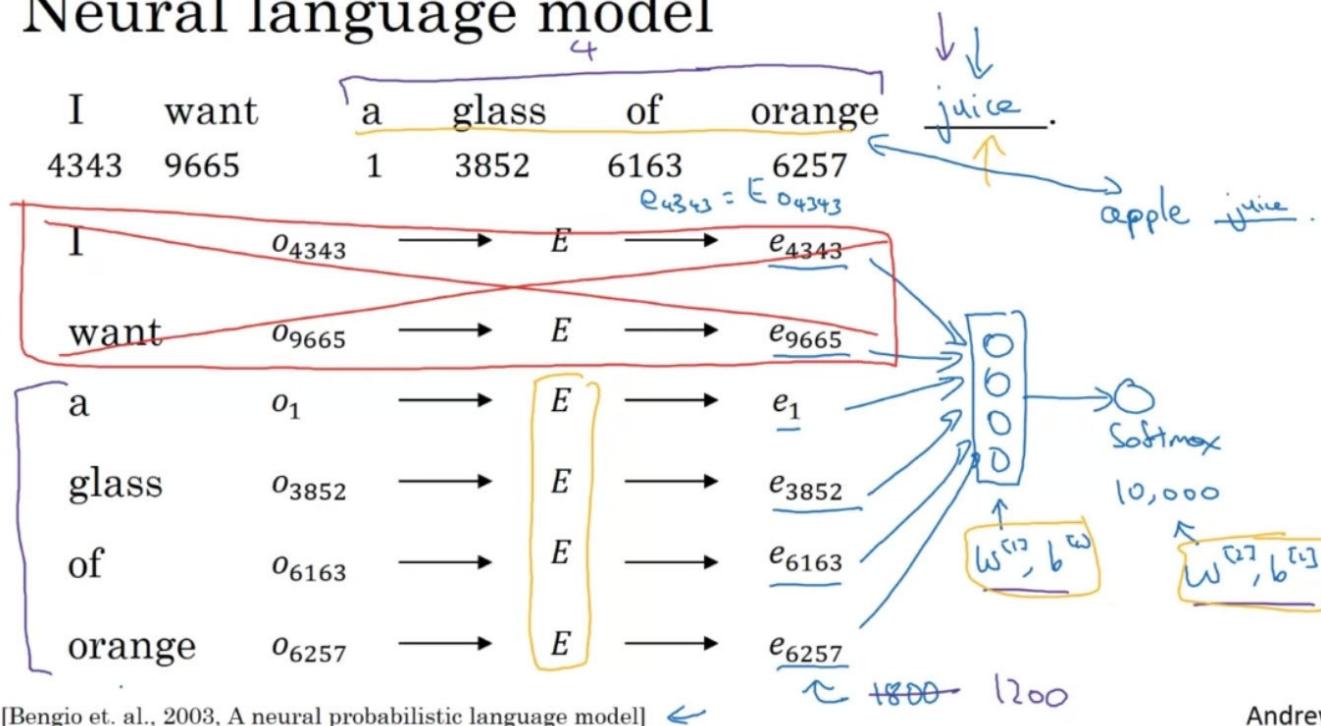


the two vectors are similar but opposite  
the first one encodes (city - country)  
while the second one encodes (country - city)

$\theta$  is close to  $180^\circ$

$$\cos(\theta) \approx -1$$

# Neural language model



If you really want to **build a language model**, it is natural to use **last few words** as the context. But if the main goal is to **learn word embedding**, then you can use **all of these other context**, and they will result in meaningful word embedding as well.

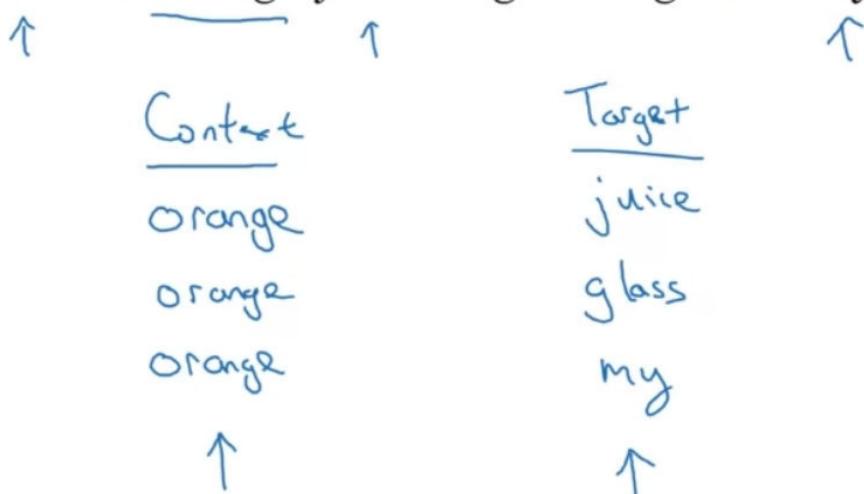
## Word2Vec

### Skip-grams

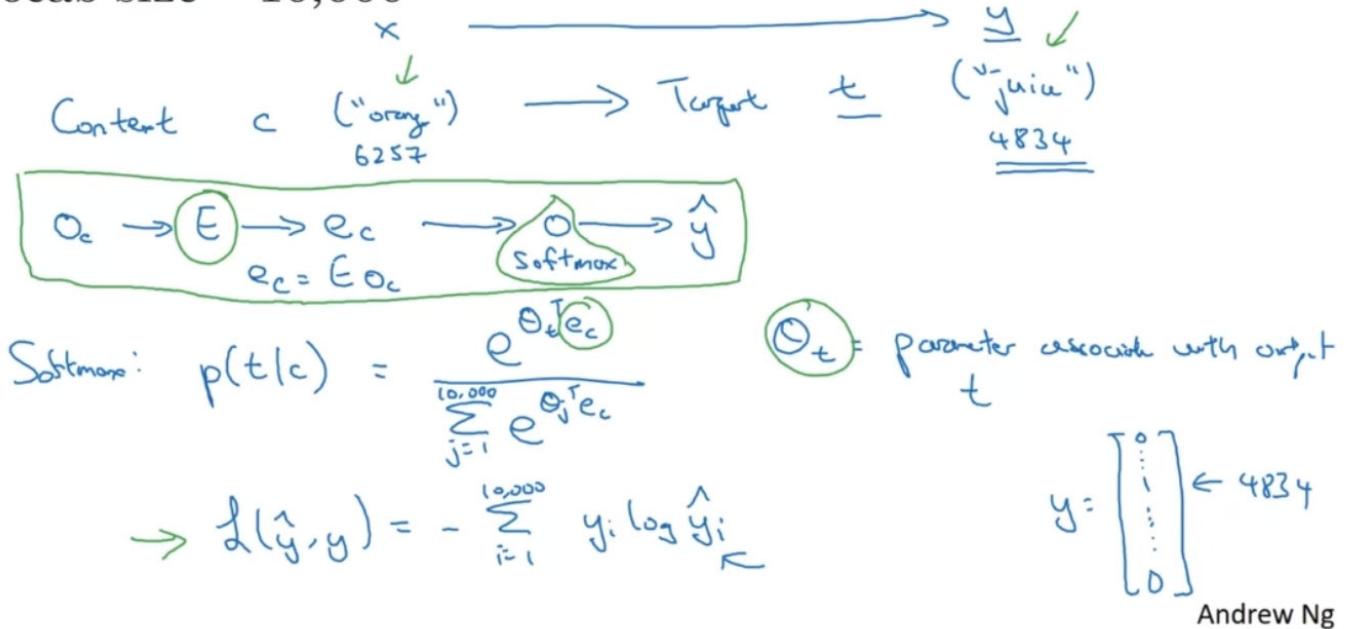
1. Randomly pick a word to be the context word.
2. Randomly pick another word within some window (Say \$pm 5\$ words) of the context word.

## Skip-grams

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



Vocab size = 10,000



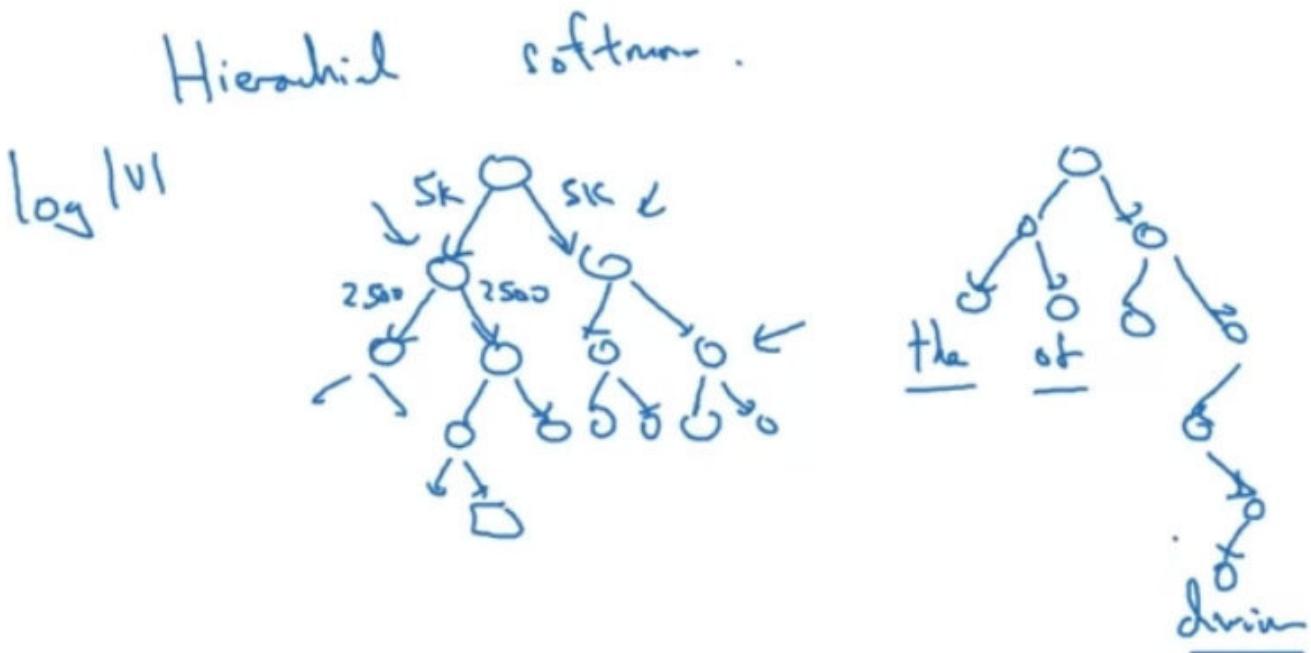
$\theta_t$  is the parameter associated with output  $t$ . It is the chance of the particular word  $t$  being the label.

However, there are a couple problems with using skip-grams.

The primary problem is **computational speed**.

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

To evaluate this probability, you need to carry out a sum over all 10,000 words in the vocabulary. One of the solution is to use **Hierarchical Softmax Classifier**. In practice, the hierarchical softmax classifier can be developed so the **common words** tend to be on **top**, whereas the less common words can be buried much deeper in the tree.



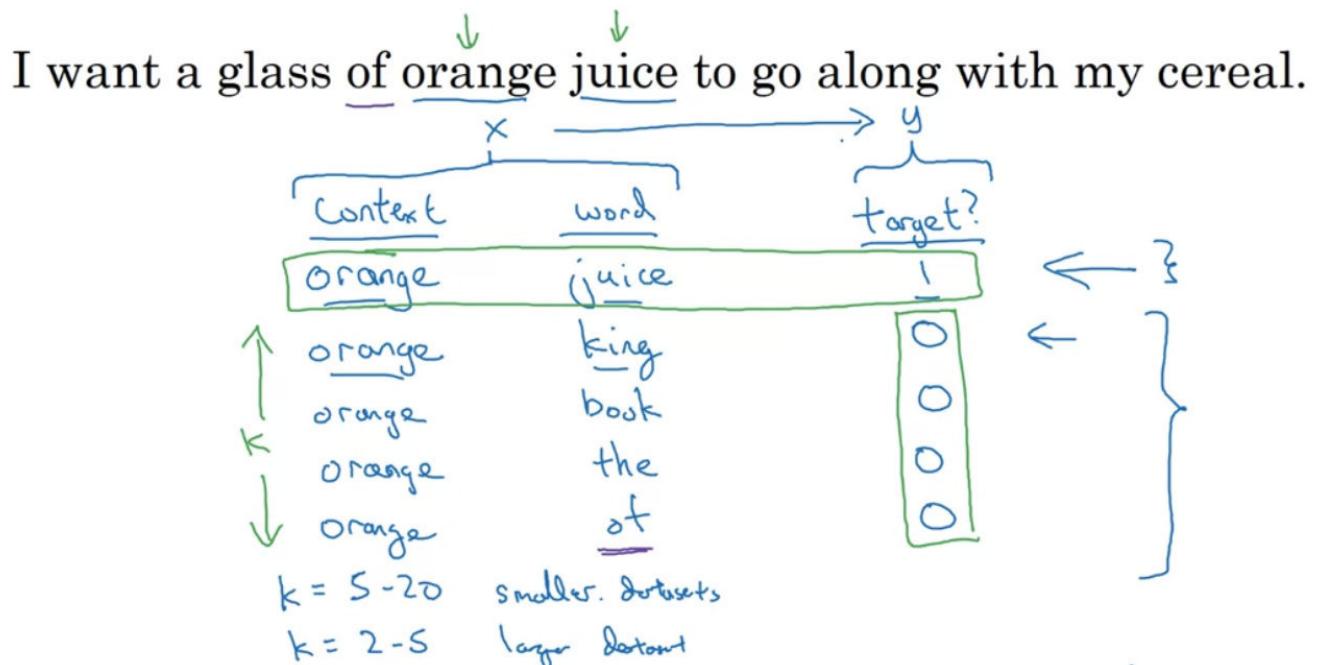
How to sample the context  $c$ ?

In practice the distribution of words pc isn't taken just entirely uniformly at random for the training set purpose, but instead there are different heuristics that you could use in order to balance out something from the common words together with the less common words.

## Negative sampling

To solve the computational speed problems, this algorithm defines a new learning problem.

# Defining a new learning problem



## Model

Softmax:  $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c) \leftarrow$$

Andrew Ng

Think of this as having 10,000 binary logistic regression classifiers, but instead of training all 10,000 of them on every iteration, we are going to train  $k+1$  of them ( $1$  corresponding to actual target word, and  $k$  randomly chosen negative examples).

## Selecting negative examples

(Empirically) Sample according to:

$$p(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

where  $f(w_i)$  is the observed **frequency** of a particular word in the English language or in the training set corpus.

# GloVe Word Vectors

## GloVe (global vectors for word representation)

GloVe (global vectors for word representation)

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

$c, t$

$x_{ij} = \# \text{ times } j \text{ appears in context of } i$

$\begin{matrix} \uparrow & \uparrow \\ c & t \end{matrix}$

[Pennington et. al., 2014. GloVe: Global vectors for word representation] ↗

Andrew Ng

For the purposes of the GloVe algorithm, we define context and target as whether or not the two words appear in close proximity.

Model

minimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) \left( \Theta_i^T e_j + b_i + b_j - \log x_{ij} \right)^2$$

*weight*  
term

$f(x_{ij}) = 0$  if  $x_{ij} = 0$ .      " $0 \log 0 = 0$ "

this, is, of, a, ...

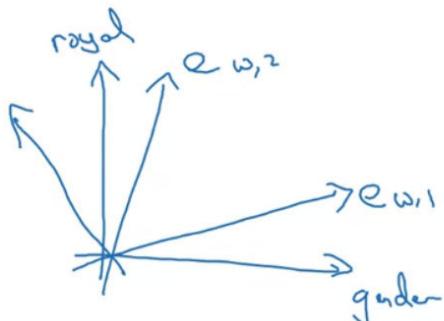
derivation

$\Theta_w^{(\text{final})} = \frac{\Theta_w + \Theta_w}{2}$

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

# A note on the featurization view of word embeddings

|        | Man<br>(5391) | Woman<br>(9853) | King<br>(4914) | Queen<br>(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            | ↙ |



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

$$(A\theta_i)^T (A^T e_j) = \theta_i^T \cancel{A} \cancel{A}^T e_j$$

Andrew Ng

You can't guarantee that the axis used to represent the features will be well-aligned with what might be easily humanly interpretable axis.

## Sentiment Classification

Sentiment classification is the task of looking at a piece of text and telling if someone likes or dislikes the thing they are talking about.

# Sentiment classification problem

$$x \longrightarrow y$$

The dessert is excellent.



Service was quite slow.



Good for a quick meal, but nothing special.



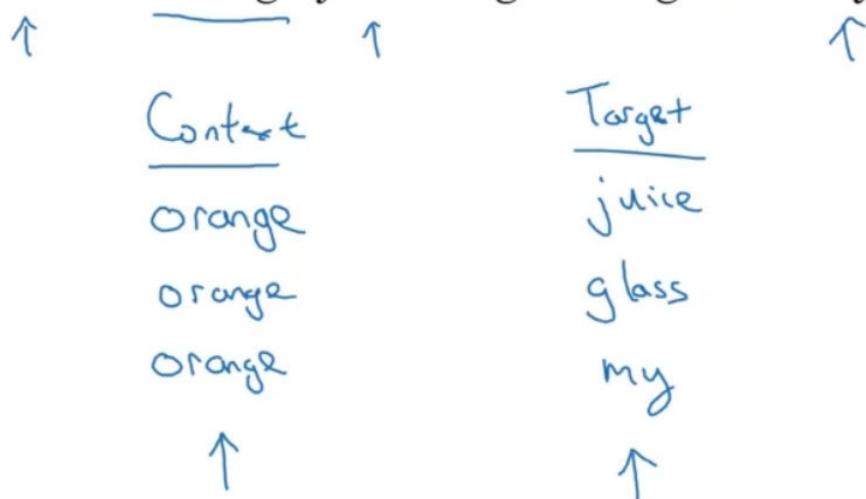
Completely lacking in good taste, good service, and good ambience.



One of the challenges of sentiment classification is you might not have a huge label data set.

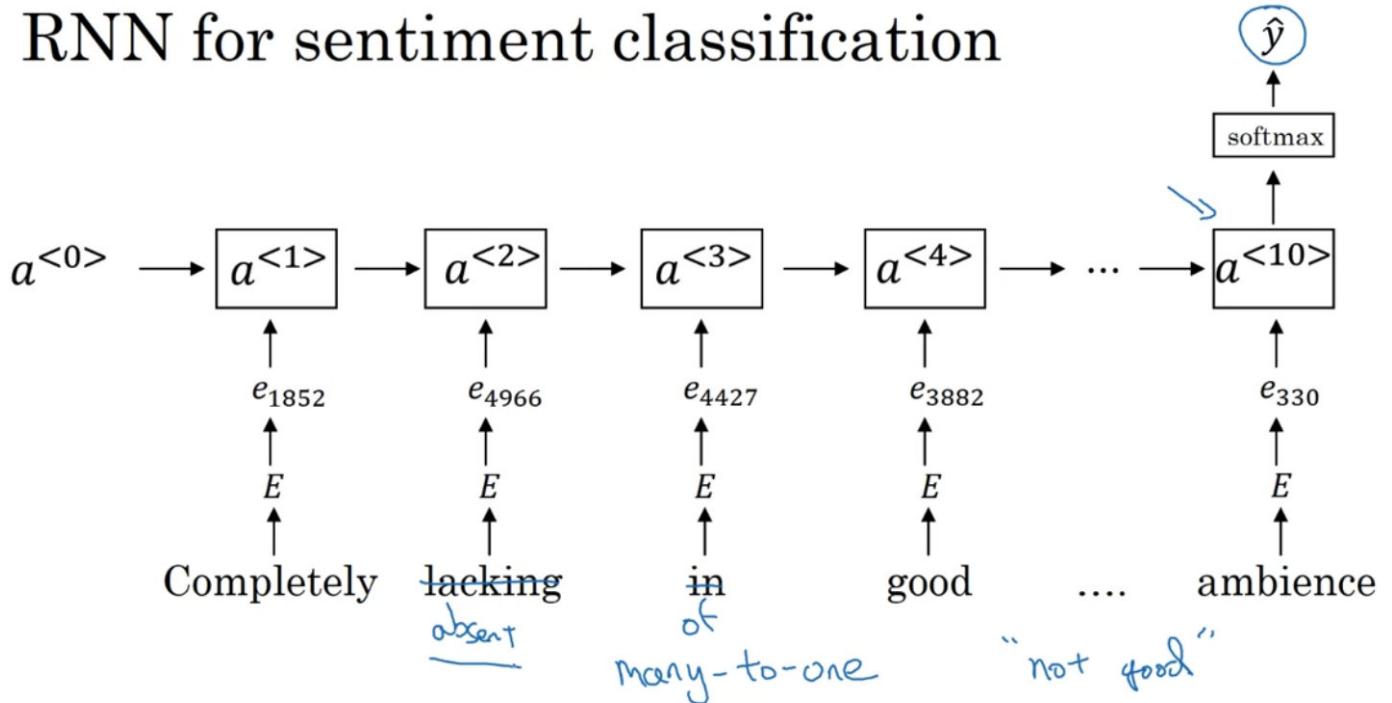
## Skip-grams

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



However, if you use an algorithm like this that ignores word **order** and just sums or averages all of the embeddings for the different words, the result is not convincing for shown example (because there are many "good" in the review, and the algorithm may think this is a good review after averaging the word embedding). So, instead of just summing all of the word embeddings, you can instead use a **RNN** for sentiment classification.

# RNN for sentiment classification



## Debiasing Word Embeddings

### The problem of bias in word embeddings

Man:Woman as King:Queen

Man:Computer\_Programmer as Woman:Homemaker X

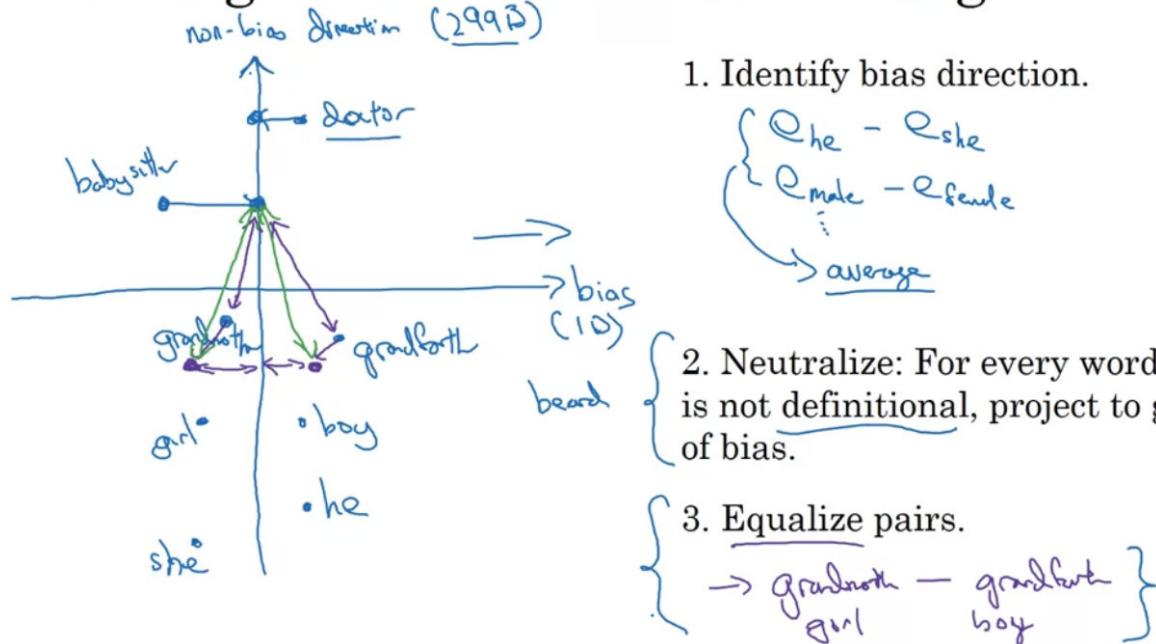
Father:Doctor as Mother:Nurse X

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

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] X

Andrew Ng

# Addressing bias in word embeddings



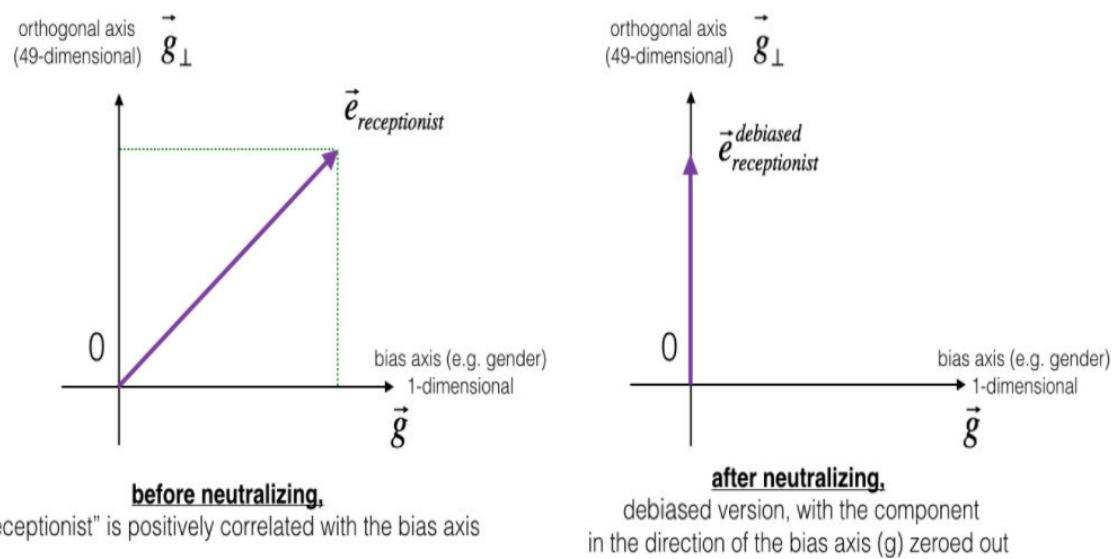
Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings] ↵

Andrew Ng

## Neutralization

The figure below should help you visualize what neutralizing does. If you're using a 50-dimensional word embedding, the 50 dimensional space can be split into two parts: The bias-direction  $g$ , and the remaining 49 dimensions, which is called  $g_{\perp}$  here. In linear algebra, we say that the 49-dimensional  $g_{\perp}$  is perpendicular (or "orthogonal") to  $g$ , meaning it is at 90 degrees to  $g$ . The neutralization step takes a vector such as  $e_{\text{receptionist}}$  and zeros out the component in the direction of  $g$ , giving us  $e_{\text{receptionist}}^{\text{debiased}}$ .

Even though  $g_{\perp}$  is 49-dimensional, given the limitations of what you can draw on a 2D screen, it's illustrated using a 1-dimensional axis below.



**Figure 2:** The word vector for "receptionist" represented before and after applying the neutralize operation.

Given an input embedding  $e$ , you can use the following formulas to compute  $e^{debiased}$ :

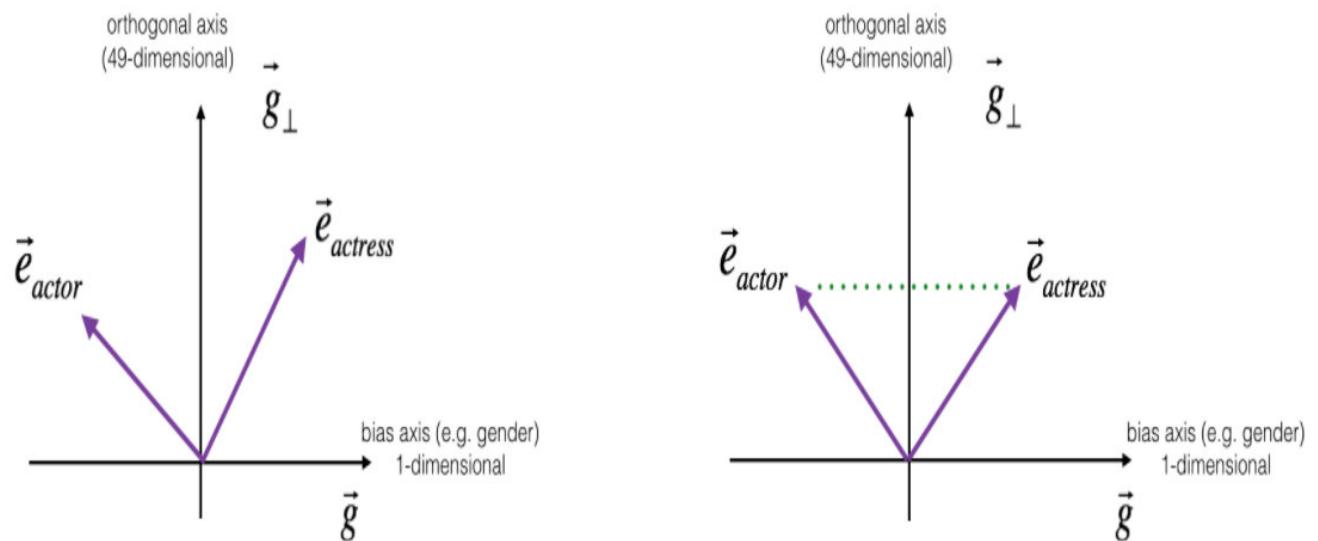
$$e^{bias\_component} = \frac{e \cdot g}{\|g\|_2} * \frac{g}{\|g\|_2} = \frac{e \cdot g}{\|g\|_2^2} * g$$

$$e^{debiased} = e - e^{bias\_component}$$

where  $e^{bias\_component}$  is the projection of  $e$  onto the direction  $g$ .

## Equalization

The key idea behind equalization is to make sure that a particular pair of words are equidistant from the 49-dimensional  $g_{\perp}$ . The equalization step also ensures that the two equalized steps are now the same distance from  $e_{receptionist}^{debiased}$ , or from any other word that has been neutralized. Visually, this is how equalization works:



**before equalizing,**  
"actress" and "actor" differ  
in many ways beyond the  
direction of  $\vec{g}$

**after equalizing,**  
"actress" and "actor" differ  
only in the direction of  $\vec{g}$ , and further  
are equal in distance from  $\vec{g}_{\perp}$

The key equations:

$$\mu = \frac{e_{w1} + e_{w2}}{2}$$

$$\mu_B = \frac{\mu \cdot \text{bias\_axis}}{\|\text{bias\_axis}\|_2^2} * \text{bias\_axis}$$

$$\mu_\perp = \mu - \mu_B$$

$$e_{w1B} = \frac{e_{w1} \cdot \text{bias\_axis}}{\|\text{bias\_axis}\|_2^2} * \text{bias\_axis}$$

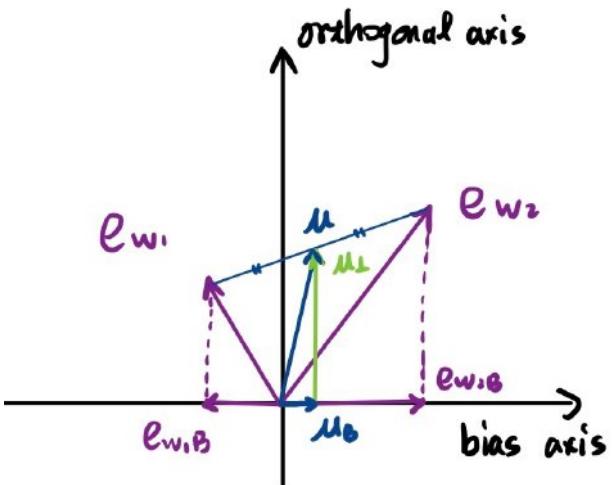
$$e_{w2B} = \frac{e_{w2} \cdot \text{bias\_axis}}{\|\text{bias\_axis}\|_2^2} * \text{bias\_axis}$$

$$e_{w1B}^{corrected} = \sqrt{|1 - \|\mu_\perp\|_2^2|} * \frac{e_{w1B} - \mu_B}{\|(e_{w1} - \mu_\perp) - \mu_B\|_2}$$

$$e_{w2B}^{corrected} = \sqrt{|1 - \|\mu_\perp\|_2^2|} * \frac{e_{w2B} - \mu_B}{\|(e_{w2} - \mu_\perp) - \mu_B\|_2}$$

$$e_1 = e_{w1B}^{corrected} + \mu_\perp$$

$$e_2 = e_{w2B}^{corrected} + \mu_\perp$$



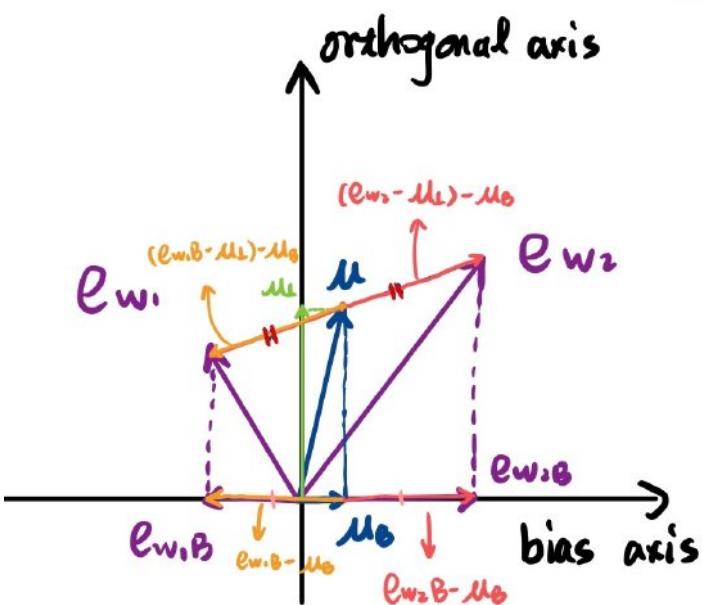
$$M = \frac{e_{w1} + e_{w2}}{2}$$

$$M_B = \frac{M \cdot \text{bias-axis}}{\|\text{bias-axis}\|_2^2} \times \text{bias-axis}$$

$$M_L = M - M_B$$

$$e_{w1B} = \frac{e_{w1} \cdot \text{bias-axis}}{\|\text{bias-axis}\|_2^2} \times \text{bias-axis}$$

$$e_{w2B} = \frac{e_{w2} \cdot \text{bias-axis}}{\|\text{bias-axis}\|_2^2} \times \text{bias-axis}$$

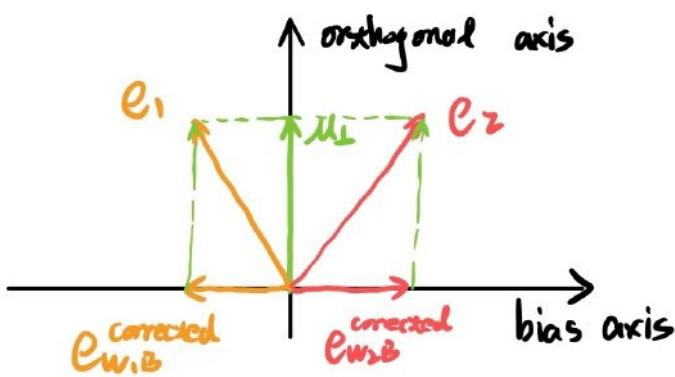


$$e_{w1B}^{\text{corrected}} = \sqrt{1 - \|M_L\|_2^2} \times \frac{e_{w1B} - M_B}{\|(e_{w1B} - M_B) - M_B\|_2}$$

$$e_{w2B}^{\text{corrected}} = \sqrt{1 - \|M_L\|_2^2} \times \frac{e_{w2B} - M_B}{\|(e_{w2B} - M_B) - M_B\|_2}$$

↓

$$\|e_{w1B}^{\text{corrected}}\|_2 = \|e_{w2B}^{\text{corrected}}\|_2$$



$$e_1 = e_{w1B}^{\text{corrected}} + M_L$$

$$e_2 = e_{w2B}^{\text{corrected}} + M_L$$

## Sequence to Sequence Model

### Basic models

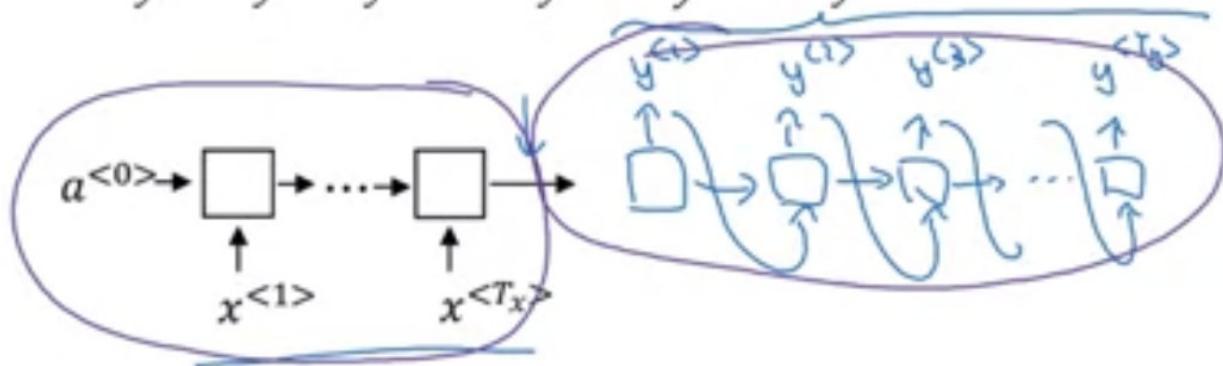
# Sequence to sequence model

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$

Jane visite l'Afrique en septembre

→ Jane is visiting Africa in September.

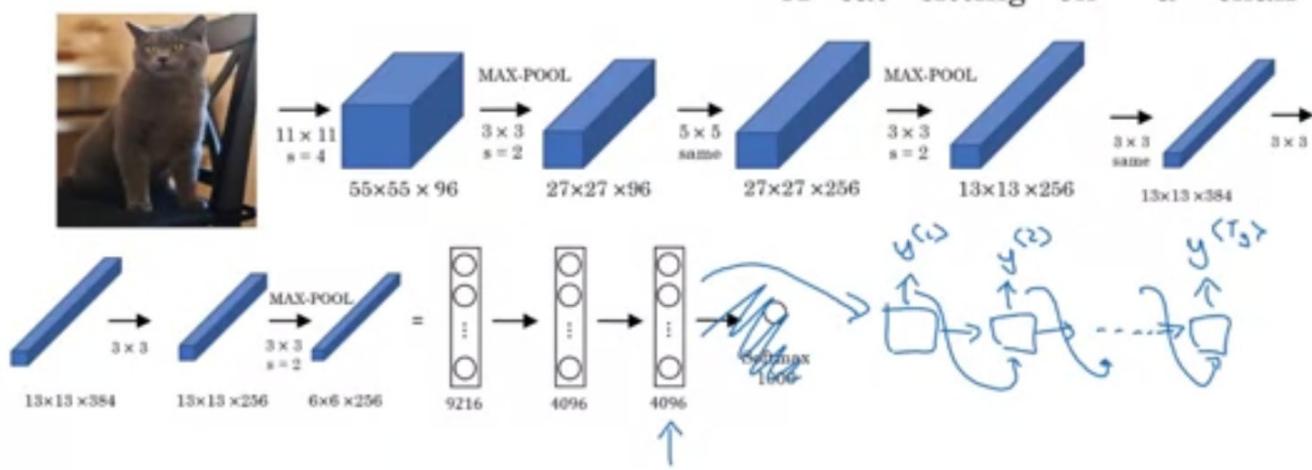
$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$



An encoding network and a decoding network.

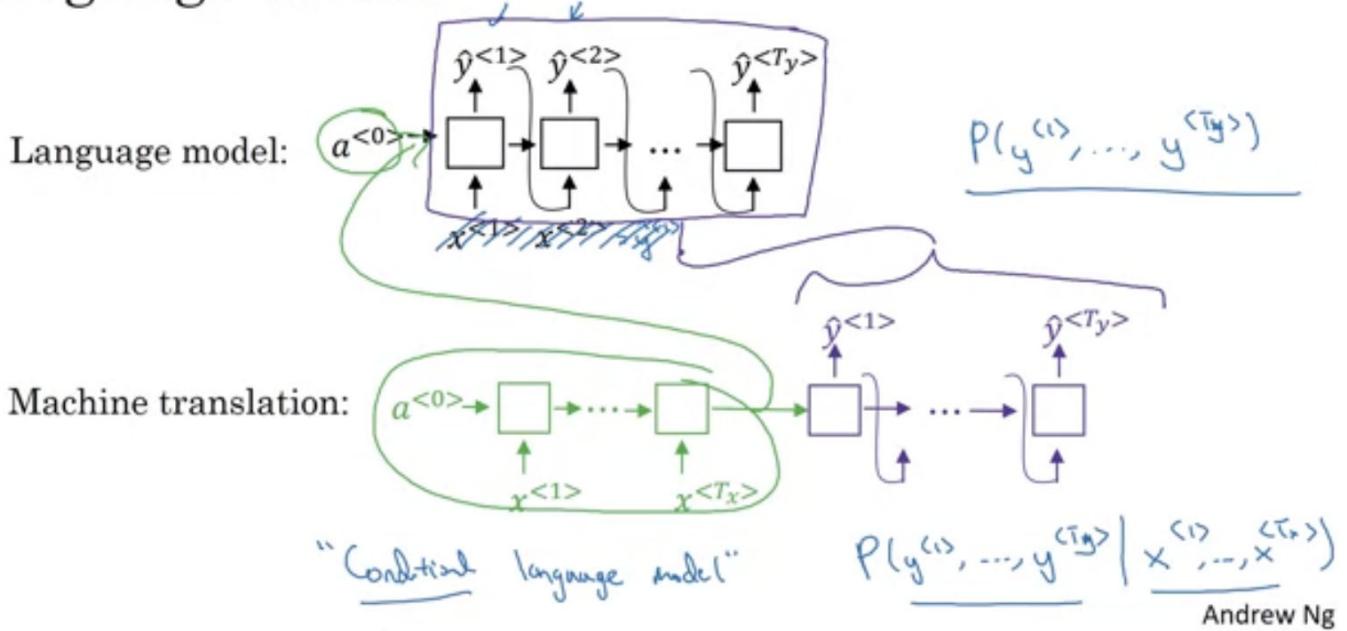
## Image captioning

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>} \quad \dots \quad \{y^{<T>}\}$



Of the key differences is that you do not want to randomly choose (like we did in generating novel text) in translation. You may want most likely translation.

# Machine translation as building a conditional language model



The translation model can be viewed as a **conditional language model**.

## Finding the most likely translation

Jane visite l'Afrique en septembre.

$$P(y^{<1>} \dots, y^{<Ty>} | x)$$

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\arg \max_{y^{<1>} \dots, y^{<Ty>}} P(y^{<1>} \dots, y^{<Ty>} | x)$$

Instead of finding the probability distribution  $P(y^1, \dots, y^{Ty} | x)$ , what you would like is to find the English sentence  $y$  that **maximizes** that conditional probability.

$$\arg \max_{y^{(1)}, \dots, y^{(Ty)}} P(y^{(1)}, \dots, y^{(Ty)} | x)$$

The most common algorithm for doing this is called **beam search**.

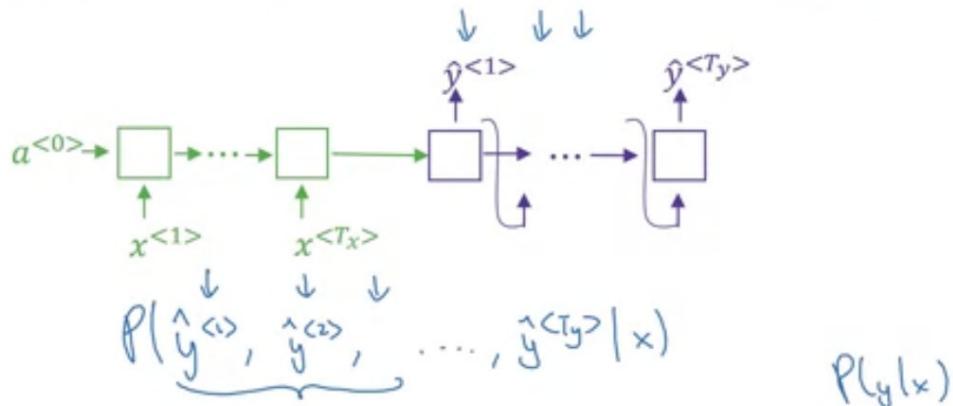
## Beam search

Why not greedy search?

Greedy search is to pick the most likely words every time.

## Why not a greedy search?

$$P(\hat{y}^{(1)} | x)$$



- Jane is visiting Africa in September.  
→ Jane is going to be visiting Africa in September.  
 $P(\text{Jane is going } | x) > P(\text{Jane is visit } | x)$

Andrew Ng

The greedy search may end up resulting in a **less optimal**, or more **verbose** sentence. Moreover, it is not always optimal to just pick one word at a time. If our vocabulary contains 10,000 words, there will be  $10,000^{10}$  possible sentences that are 10 words long in total, it is impossible to rate them all. So, the most common thing to do is use **an approximate search algorithm**. The approximate search algorithm will try to pick the sentence  $y$  that maximizes the conditional probability.

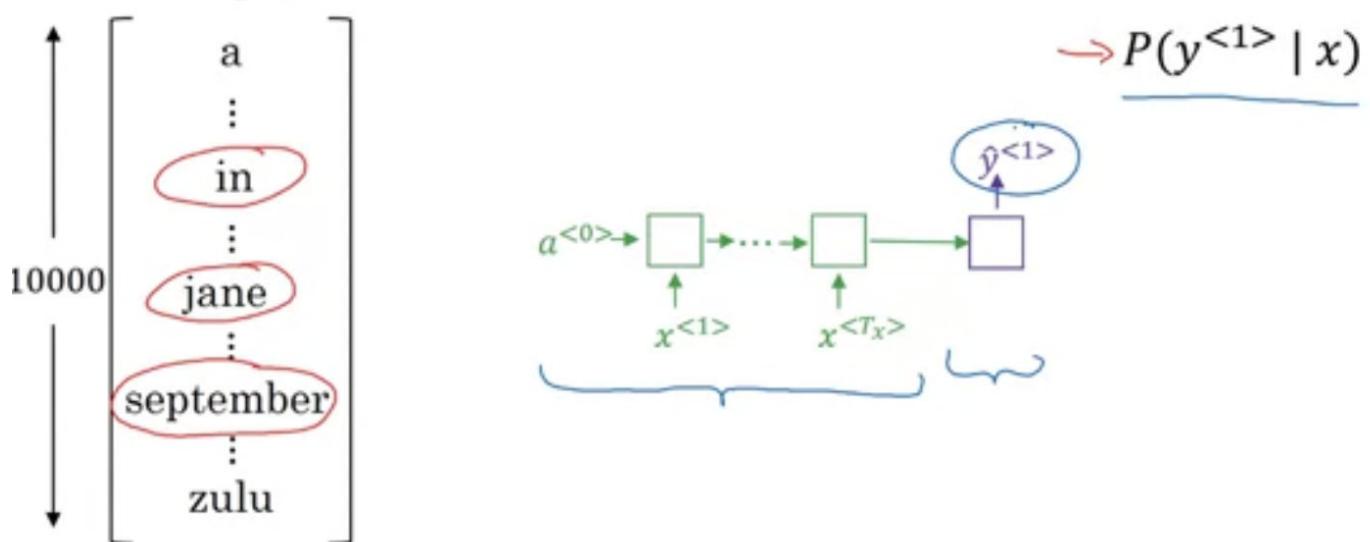
## Beam search algorithm

The beam search algorithm has a parameter called **beam width** denoted by  $B$ . Instead of picking only the one most likely words like greedy search does, beam search can consider **multiple alternatives** according to beam width.

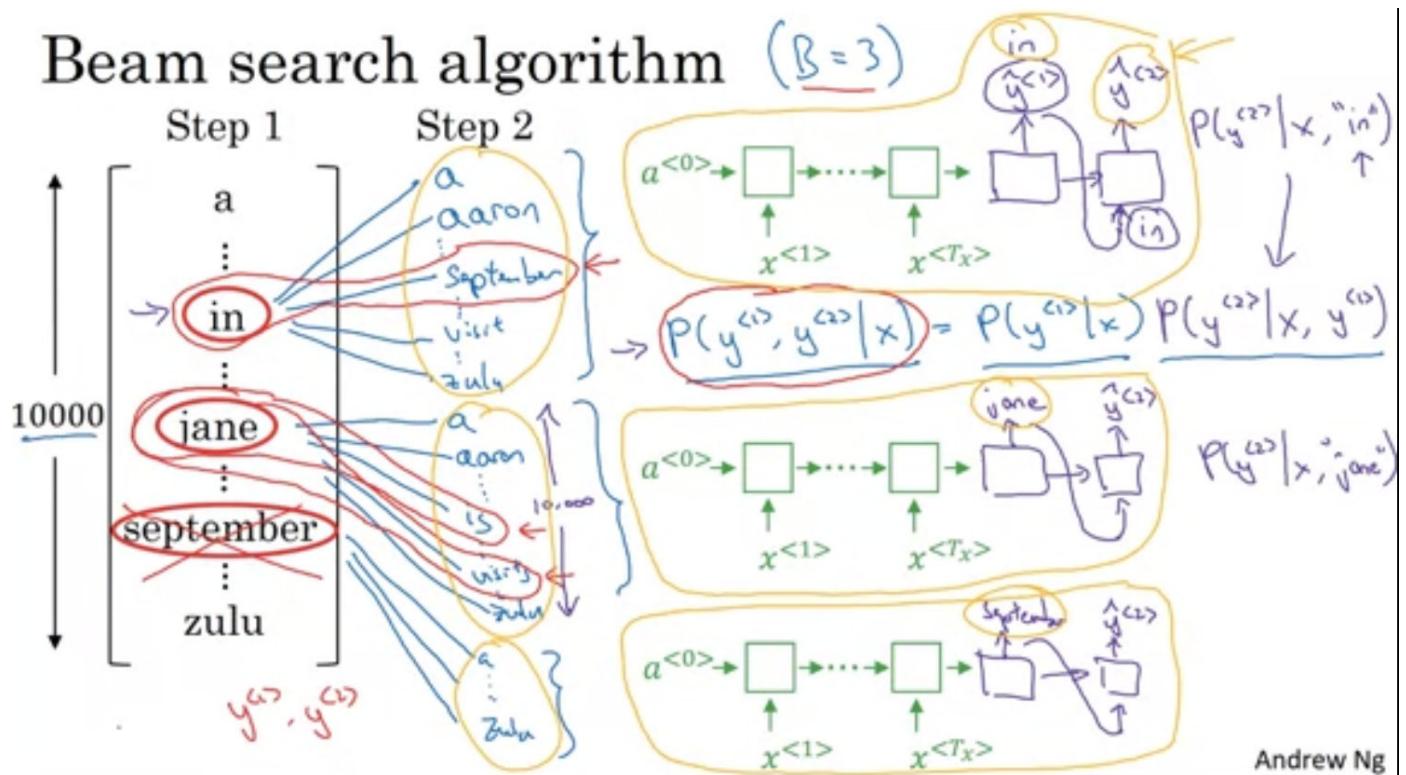
# Beam search algorithm

B = 3 (beam width)

Step 1



# Beam search algorithm (B=3)



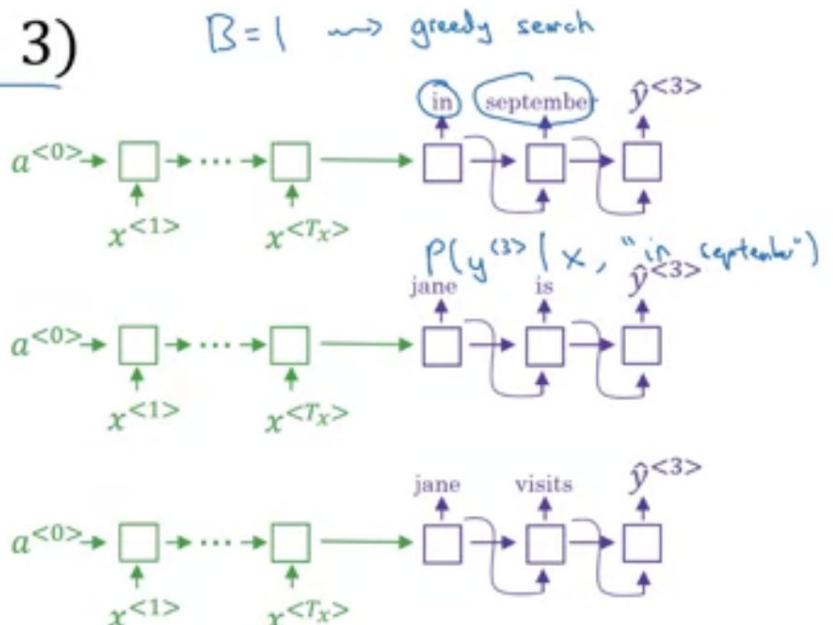
$$P(A, B | C) = \frac{P(A, B, C)}{P(C)} = \frac{P(A, B, C)}{P(B, C)} \frac{P(B, C)}{P(C)} = P(A | B, C) P(B | C)$$

# Beam search ( $B = 3$ )

in september

jane is

jane visits



$$P(y^{<1>} , y^{<2>} | x)$$

jane visits africa in september. <EOS>

Andrew Ng

Notice that if the beam width is set to be equal to one, then the beam search becomes the greedy search algorithm.

## Refinements to beam search

Recall that the objective is to

$$\arg \max_{y^{(1)}, \dots, y^{(T_y)}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

Notice that

$$P(y^{(1)}, \dots, y^{(T_y)} | x) = P(y^{(1)} | x) P(y^{(2)} | x, y^{(1)}) \cdots P(y^{(T_y)} | x, y^{(1)}, \dots, y^{(T_y-1)})$$

So, the objective function is

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

Notice that in the objective function, these probabilities are all numbers less than 1, and multiplying a lot of them will result in a very tiny number, which can result in **numerical underflow**, meaning that it is too small for the floating part representation in your computer to store accurately.

So, in practice, we will take logs and the objective function becomes

$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

## Length normalization

For both of these objective functions, however, if you have a very **long** sentence, the probability of the sentence will be very **low** (multiply lots of number smaller than 1, or sum over many negative numbers). So these objective function have an undesirable affect that they tend to prefer **short** translations.

To make the algorithm works better, we need to add one more term to normalize it:

$$\arg \max_y \frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

This will significantly reduce the penalty for outputting longer translations.

In practice, instead of dividing by  $T_y$ , sometimes we will use a softer approach

$$\arg \max_y \frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

where normally  $\alpha=0.7$ .

## Beam search discussion

### How to choose beam width $B$ ?

- Large  $B$ : better result, but slower
- Small  $B$ : worse result, but faster

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

## Error analysis in beam search

# Example

Jane visite l'Afrique en septembre.

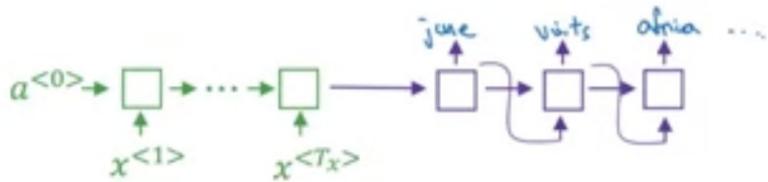
$\rightarrow$  RNN  
 $\rightarrow$  Beam Search

BT

Human: Jane visits Africa in September. ( $y^*$ )

Algorithm: Jane visited Africa last September. ( $\hat{y}$ )  $\leftarrow$

RNN computes  $P(y^*|x) \geq P(\hat{y}|x)$



What you can do is to compute  $P(y^*|x)$  and  $P(\hat{y}|x)$  and see which one is bigger.

# Error analysis on beam search

$$P(y^*|x)$$

$$P(\hat{y}|x)$$

Human: Jane visits Africa in September. ( $y^*$ )

Algorithm: Jane visited Africa last September. ( $\hat{y}$ )

Case 1:  $P(y^*|x) > P(\hat{y}|x)$  ←

$$\arg \max_y P(y|x)$$

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.

Only when a lot of fractions of faults are due to beam search, then you may need to increase the beam width.

## Bleu score

What the Bleu score does is given a machine generated translation, it allows you to automatically compute a score that measures how good is that machine translation.

Bleu stands for bilingual evaluation understudy.

## Evaluating machine translation

French: Le chat est sur le tapis.

Bleu  
bilingual evaluation understudy

→ Reference 1: The cat is on the mat. ← 2 <sup>opposite</sup>

→ Reference 2: There is a cat on the mat. ←

→ MT output: the the the the the the.

$$\text{Precision: } \frac{7}{7}$$

$$\text{Modified precision: } \frac{2}{7} \leftarrow \begin{array}{l} \text{Count}_{\text{clip}}(\text{"the"}) \\ \text{Count}(\text{"the"}) \end{array}$$

# Bleu score on bigrams

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

|         | Count | Count <sub>clip</sub> |  |
|---------|-------|-----------------------|--|
| the cat | 2 ←   | 1 ←                   |  |
| cat the | 1 ←   | 0                     |  |
| cat on  | 1 ←   | 1 ←                   |  |
| on the  | 1 ←   | 1 ←                   |  |
| the mat | 1 ←   | 1 ←                   |  |

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

For unigram (one word), the precision  $P_1$  (1 stands for unigram) is given by:

$$P_1 = \frac{\sum_{\text{unigram} \in \hat{y}} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\text{unigram} \in \hat{y}} \text{Count}(\text{unigram})}$$

For n-gram (n words pair), the precision  $P_n$  is given by:

$$P_n = \frac{\sum_{\text{n-grams} \in \hat{y}} \text{Count}_{\text{clip}}(\text{n-grams})}{\sum_{\text{n-grams} \in \hat{y}} \text{Count}(\text{n-grams})}$$

# Bleu details

$p_n$  = Bleu score on n-grams only

$P_1, P_2, P_3, P_4$

Combined Bleu score:  $BP \exp\left(\frac{1}{4} \sum_{n=1}^4 p_n\right)$

BP = brevity penalty

$$BP = \begin{cases} 1 & \text{if } MT\_output\_length > reference\_output\_length \\ \exp(1 - MT\_output\_length/reference\_output\_length) & \text{otherwise} \\ \exp(1 - reference\_output\_length/MT\_output\_length) \end{cases}$$

[Papineni et. al., 2002. Bleu: A method for automatic evaluation of machine translation]

Andrew Ng

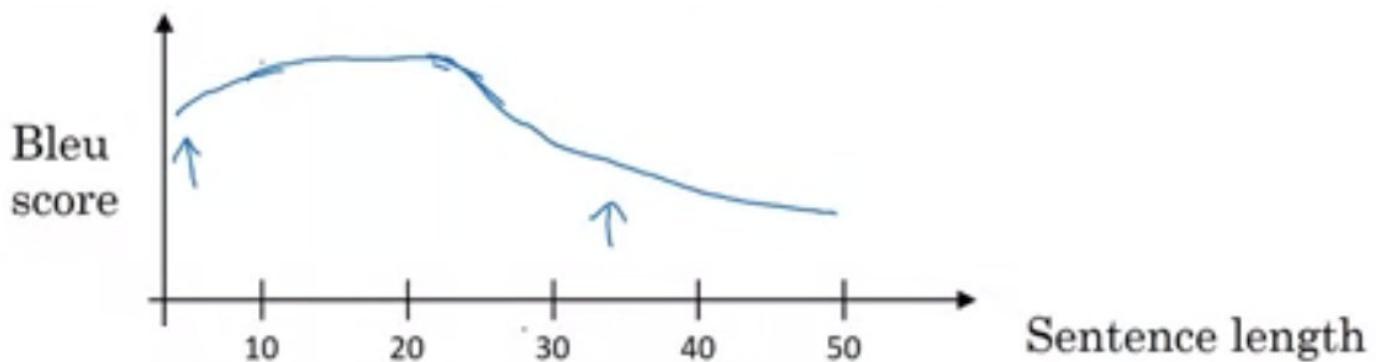
BP stands for **brevity penalty**.

Today, BLEU score is used to evaluate many systems that **generate text**, such as machine translation systems, as well as image captioning systems where you would have a neural network generated image caption. It is a useful single number evaluation metric to use when you want your algorithm to generate text, and you want to see whether it has similar meaning as the reference text written by human.

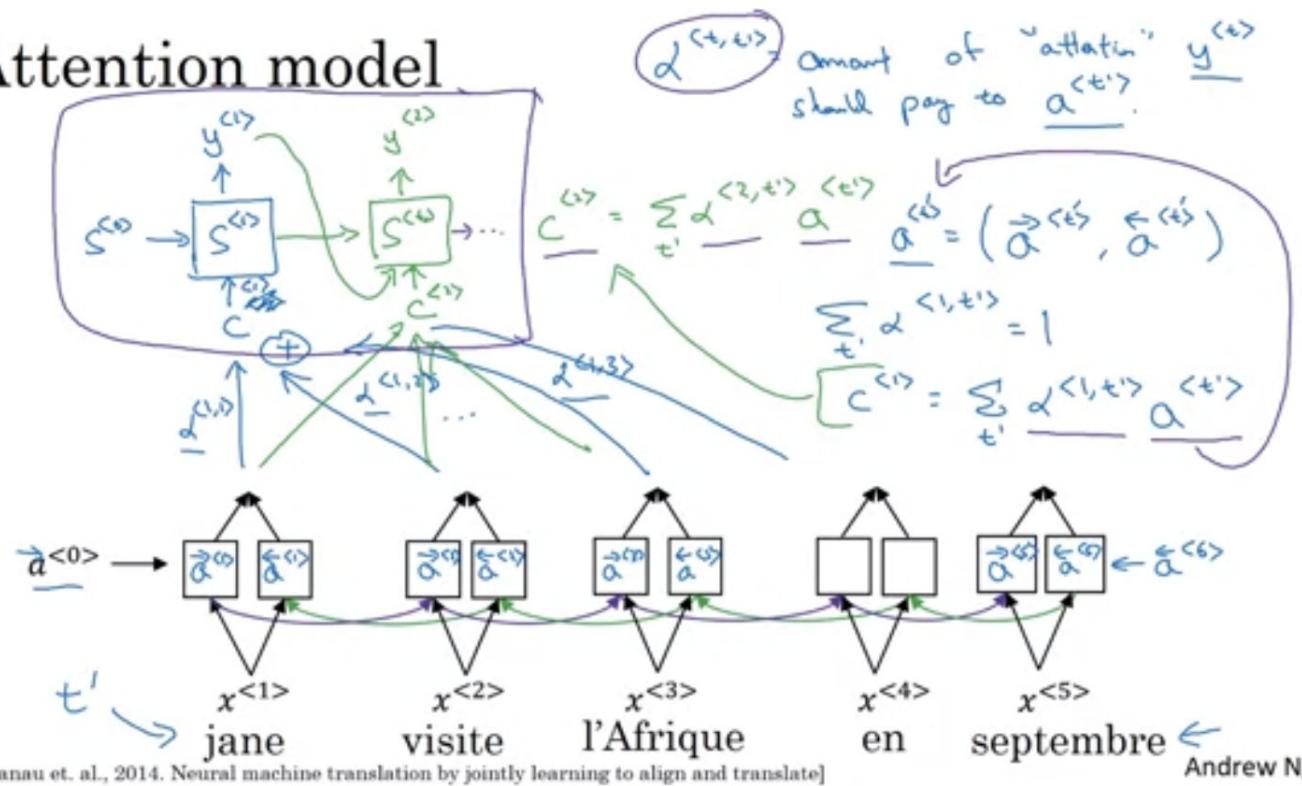
## Attention model

### The problem of long sequences

It is difficult for a network to memorize long sentence.



# Attention model



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

We use

$$a^{(t')} = \begin{bmatrix} \vec{a}^{(t')} \\ \underline{a}^{(t')} \end{bmatrix}$$

to represent feature vector for time step  $t'$ . For generating the translation, we use a single direction RNN with state  $s$ . This single RNN will generate the translation  $y$ . It also has an input called **context**  $c$  for every time step.

$$c^{(t)} = \sum_{t'} \alpha^{(t,t')} a^{(t')}$$

where  $\alpha$  is the **attention parameter**.

$$\alpha^{(t,t')} = \text{amount of attention } y^{(t)} \text{ should pay to } a^{(t')}$$

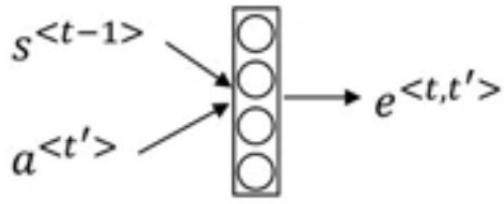
The formula to compute  $\alpha^{(t,t')}$  is

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'=1}^{T_x} \exp(e^{(t,t')})}$$

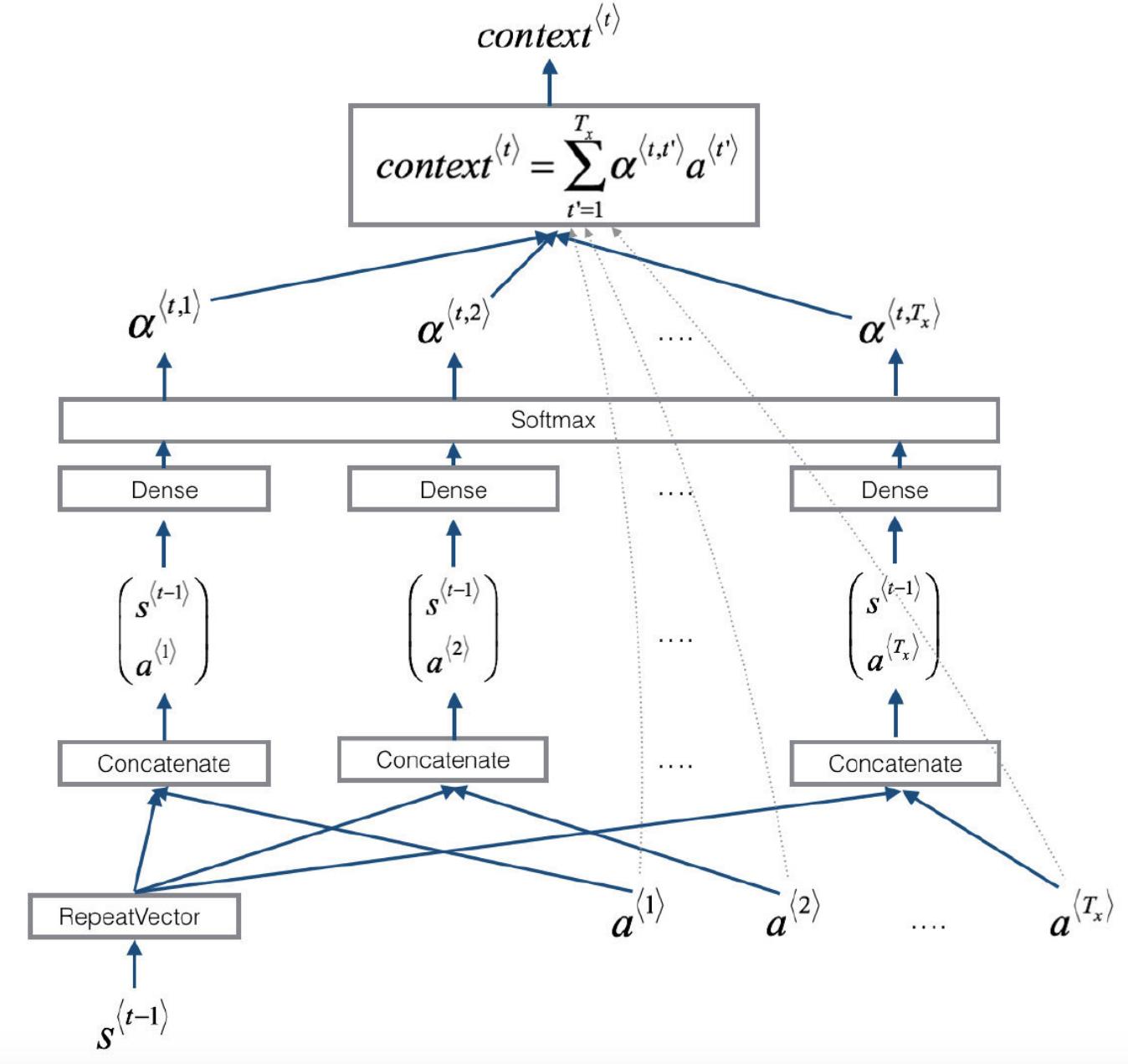
Using the softmax prioritization is just to ensure

$$\sum_{t'} \alpha^{(t,t')} = 1$$

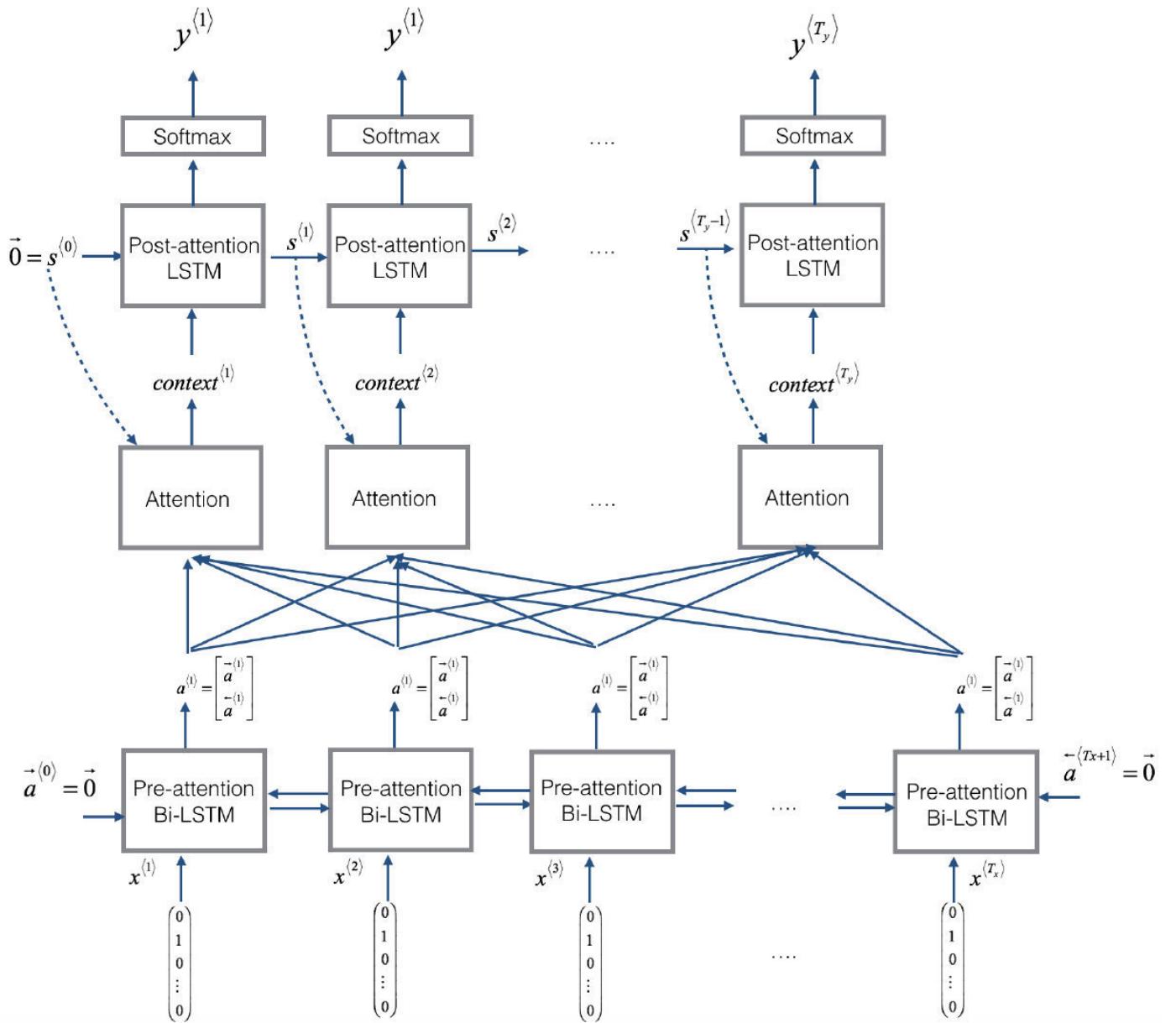
$e$  is called the "**energies**" variable, one way to compute the factors  $e$  is to use a small neural network (usually just one hidden layer) with **softmax** activation as follows:



The intuition is, if you want to decide how much attention to pay to the  $a^{\langle t \rangle}$ , it should depend the most on what is your own hidden state activation from the previous time step, which is  $s^{\langle t-1 \rangle}$ , and also the features of original words, which is  $a^{\langle t \rangle}$ . The following picture shows how the context is generated.



The whole attention model is as follow:



Here are some properties of the model that you may notice:

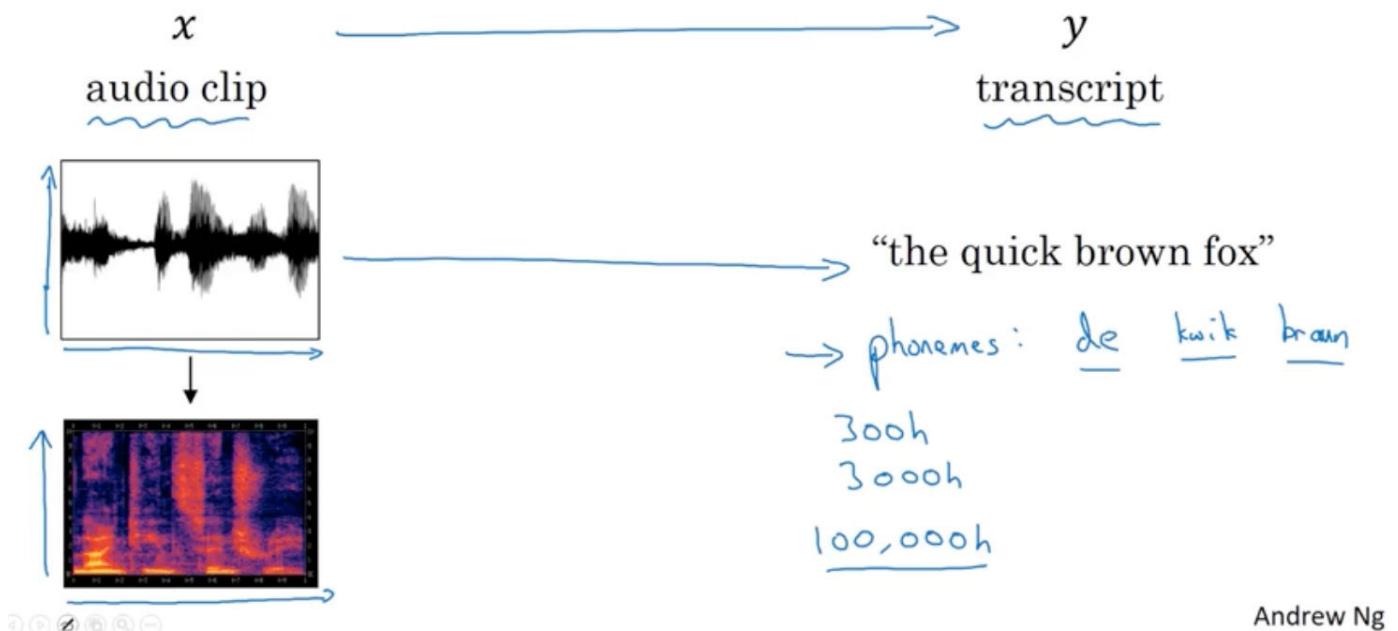
#### Pre-attention and Post-attention LSTMs on both sides of the attention mechanism

- There are two separate LSTMs in this model (see diagram on the left): pre-attention and post-attention LSTMs.
- *Pre-attention* Bi-LSTM
  - The pre-attention Bi-LSTM goes through  $T_x$  time steps
- *Post-attention* LSTM
  - The post-attention LSTM goes through  $T_y$  time steps.
- The post-attention LSTM passes the hidden state  $s^{(t)}$  and cell state  $c^{(t)}$  from one time step to the next.

One **downside** to this algorithm is that it does take **quadratic time** or **quadratic cost** to run this algorithm. Because the total number of attention parameters is  $T_x \times T_y$ .

## Speech Recognition

# Speech recognition problem



Andrew Ng

The first figure plots the audio clip, the horizontal axis is time, and the vertical axis is actually air pressure. The common preprocessing step is to convert the audio clip in to a **spectrogram**. The horizontal axis of the spectrogram is time, and the vertical axis is the frequency, and the intensity of different color shows the amount of energy. It shows "how loud is the sound at different frequency at different time".

### What really is an audio recording?

- A microphone records little variations in air pressure over time, and it is these little variations in air pressure that your ear also perceives as sound.
- You can think of an audio recording as a long list of numbers measuring the little air pressure changes detected by the microphone.
- We will use audio sampled at 44100 Hz (or 44100 Hertz).
  - This means the microphone gives us 44,100 numbers per second.
  - Thus, a 10 second audio clip is represented by 441,000 numbers ( $= 10 \times 44,100$ ).

### Spectrogram

- It is quite difficult to figure out from this "raw" representation of audio whether the word "activate" was said.
- In order to help your sequence model more easily learn to detect trigger words, we will compute a *spectrogram* of the audio.
- The spectrogram tells us how much different frequencies are present in an audio clip at any moment in time.
- If you've ever taken an advanced class on signal processing or on Fourier transforms:
  - A spectrogram is computed by sliding a window over the raw audio signal, and calculating the most active frequencies in each window using a Fourier transform.

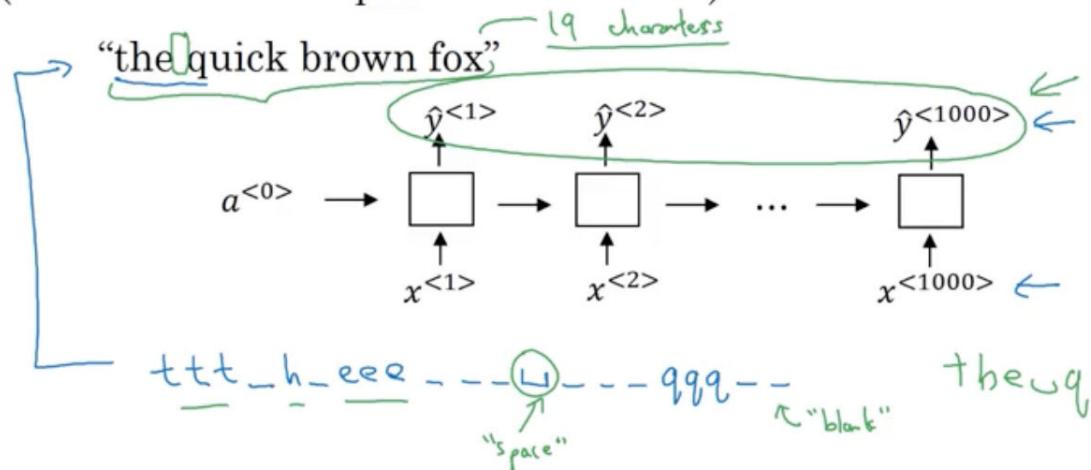
## CTC cost for speech recognition

CTC stands for Connectionist Temporal Classification.

Basic rule: collapse repeated characters not separated by "blank" (not space).

## CTC cost for speech recognition

(Connectionist temporal classification)

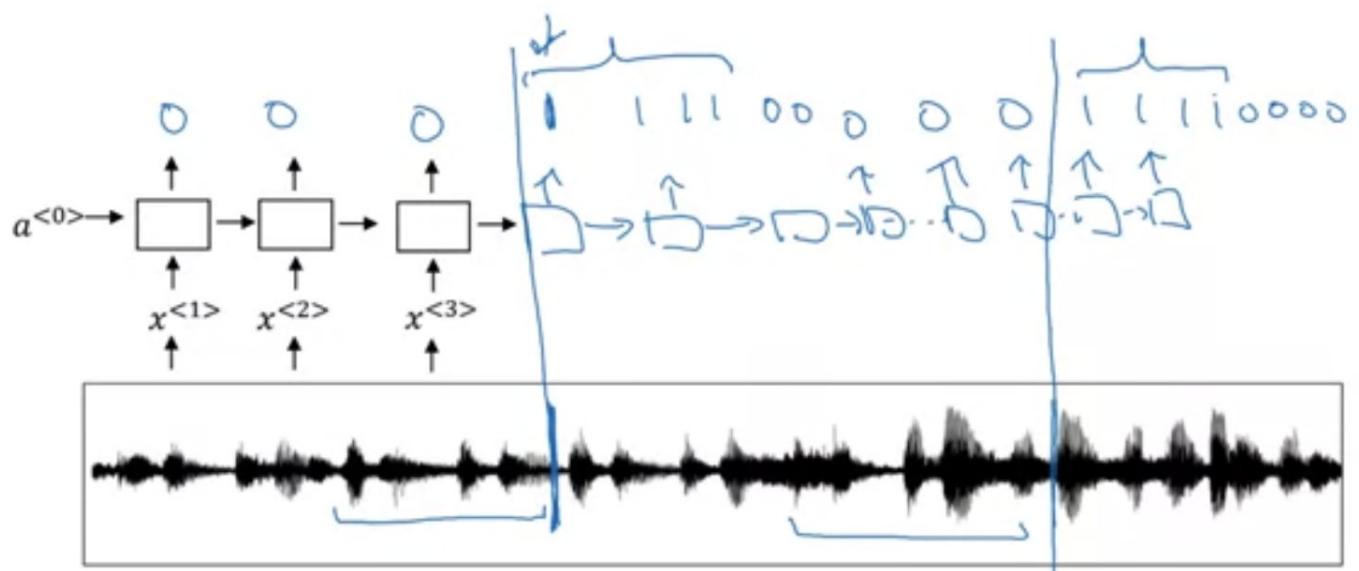


Basic rule: collapse repeated characters not separated by "blank" ↴

[Graves et al., 2006. Connectionist Temporal Classification: Labeling unsegmented sequence data with recurrent neural networks] Andrew Ng

## Trigger word detection

## Trigger word detection algorithm



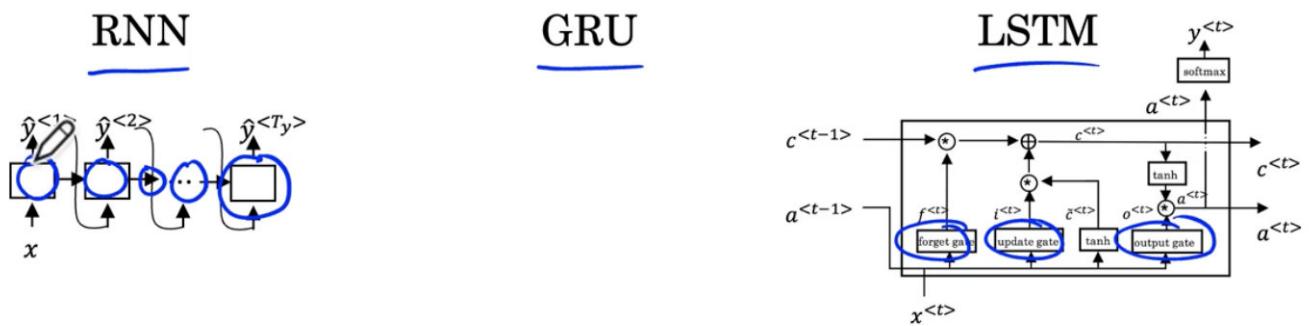
Andrew Ng

# Transformers

## Transformer network motivation

As we move from RNNs to GRU to LSTM, the models became more and more complex. And all of these models are still **sequential models** in that they ingested the input, maybe one word or one token at a time. And so, each unit was like a **bottleneck** to the flow of information. Because to compute the output of the final unit, for example, you have to compute the outputs of all of the unit that come before.

*increased complexity  
sequential*



The transformer architecture allows you to run a lot of these computations for an entire sequence **in parallel**. The major innovation of the transformer architecture is combining the use of **attention based representations** and a **CNN style of processing**. The two key ideas we will go through are:

- Self-attention
- Multi-Head Attention

## Self-attention

Self-attention enable us to use attention with a style more like CNNs.

First, we need to compute an **attention-based representation**:

$$A(q, K, V) = \text{attention-based vector representation of a word}$$

We will compute this for each word to get  $A^{(t)}$ . The intuition of  $A^{(t)}$  is that it will look at the surrounding words to try to figure out the context of this word and find the most appropriate representation for this word.

Recall the RNN attention formula:

$$\alpha^{(t,t')} = \frac{\exp(e^{\langle t,t' \rangle})}{\sum_{t'=1}^{T_x} \exp(e^{\langle t,t' \rangle})}$$

The transformers attention formula is similar:

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

The main difference is that for each word you have three values called

$q^{(t)}$  : query

$k^{(t)}$  : key

$v^{(t)}$  : value

# Self-Attention Intuition

$A(q, K, V)$  = attention-based vector representation of a word  
calculate for each word  $A^{<1>}, \dots, A^{<s>}$

## RNN Attention

$$\alpha^{<t,t'}> = \frac{\exp(e^{<t,t'}>)}{\sum_{t'=1}^T \exp(e^{<t,t'}>)}$$

$\underbrace{\hspace{10em}}_{A^{<3>}}$

$\downarrow \quad \downarrow \quad \downarrow$

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \underline{l'Afrique}$

Jane      visite      l'Afrique

## Transformers Attention

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

$\downarrow \quad \downarrow \quad \downarrow$

$x^{<4>} \quad x^{<5>} \quad q^{<3>} \quad k^{<3>} \quad v^{<3>}$

en      septembre

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

Andrew Ng

Now, let's go through how to get  $A^{(t)}$

First, we will associate each of words with three values, the **query**, **key**, and **value** pairs. If the  $x^{(t)}$  is the word embedding, the way that  $q^{(t)}$  is computed is as a learned matrix:

$$q^{(t)} = W^Q \cdot x^{(t)}$$

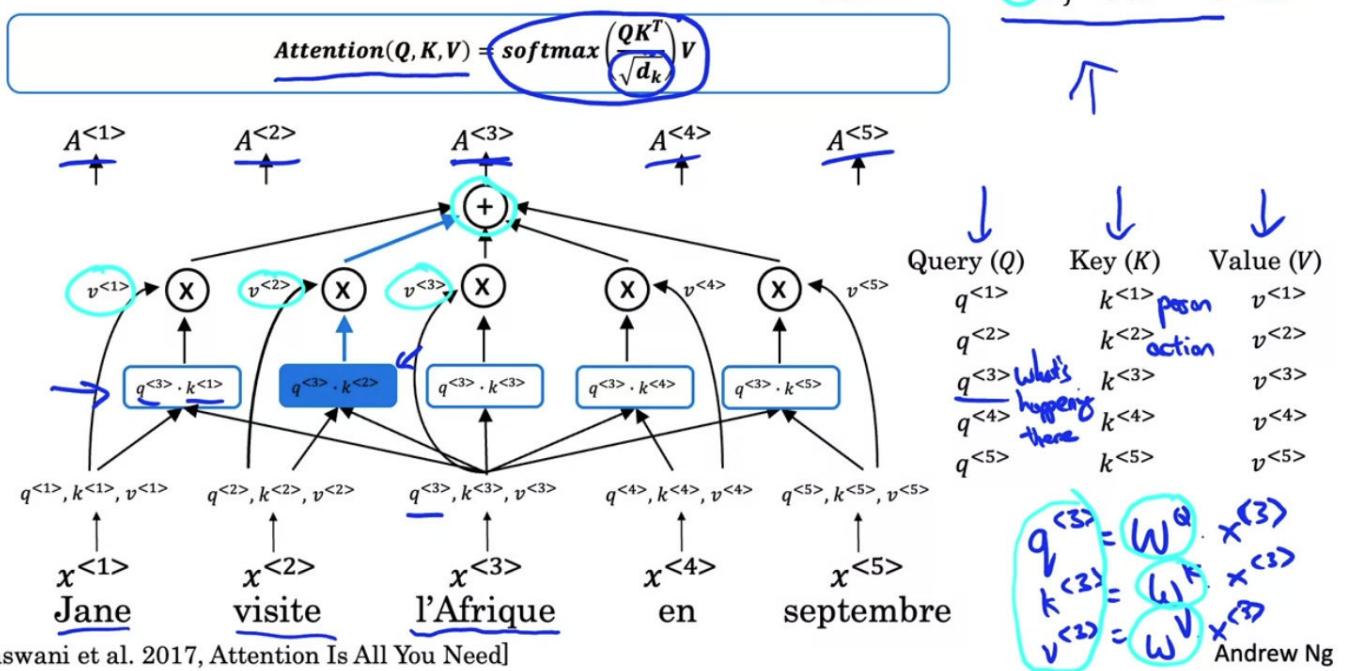
And similarly for the key and value,

$$k^{(t)} = W^K \cdot x^{(t)}$$

$$v^{(t)} = W^V \cdot x^{(t)}$$

These matrices,  $W^Q, W^K, W^V$ , are parameters of this learning algorithm, and they allow you to pull off these query, key, and value vectors for each word.

# Self-Attention



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

One intuition behind the intent of these query, key, value vectors

$q^{(t)}$  may represent a question like, what is happening to the word  $x^{(t)}$  when computing  $A^{(t)}$ ? In the above example,  $q^{(3)}$  may represent what is happening to l'Afrique (Africa) when we compute  $A^{(3)}$ ? Does l'Afrique represent a destination of traveling or the second largest continent in this sentence?

Then, we are going to compute the inner product between  $q^{(t)}$  and  $k^{(1)}$ , and this will tell us how good is an answer word  $x^{(1)}$  to the question what is happening to the word  $x^{(t)}$ ? For example,  $q^{(3)} \cdot k^{(2)}$  will tell us how good is "visite" an answer to the question of what is happening to l'Afrique (Africa), and so on for the other words in the sequence. The goal of this operation is to **pull off the most information that is needed** to help us compute the most useful representation  $A^{(t)}$ .

Next, for the above example, you may find that  $q^{(3)} \cdot k^{(2)}$  gives you largest value, this may suggest that "visite" gives you **the most relevant contexts** for the question of what is happening to l'Afrique (Africa) (view Africa as a destination for a visit).

In summary for intuition:

- $Q$ : interesting questions about the words in a sentence.
- $K$ : qualities of words given a  $Q$ .
- $V$ : specific representations of words given a  $Q$ .

Then we will take softmax over these inner products, and multiply them with  $v^{(t)}$ , which is value for word  $x^{(t)}$ . Finally, we sum them up to get  $A^{(t)}$ .

The key advantage of this representation is that the word of  $x^{(t)}$  is not some fixed word embedding, instead it lets the self-attention mechanism realize that this word is actually associated with another word (Africa is the destination of a visit), and thus compute a richer, more useful representation for this word.

The notation used in literature is like this:

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

- $Q$  is the matrix of queries
- $K$  is the matrix of keys
- $V$  is the matrix of values
- $d_k$  is the dimension of the keys, which is used to scale everything down so the softmax doesn't explode

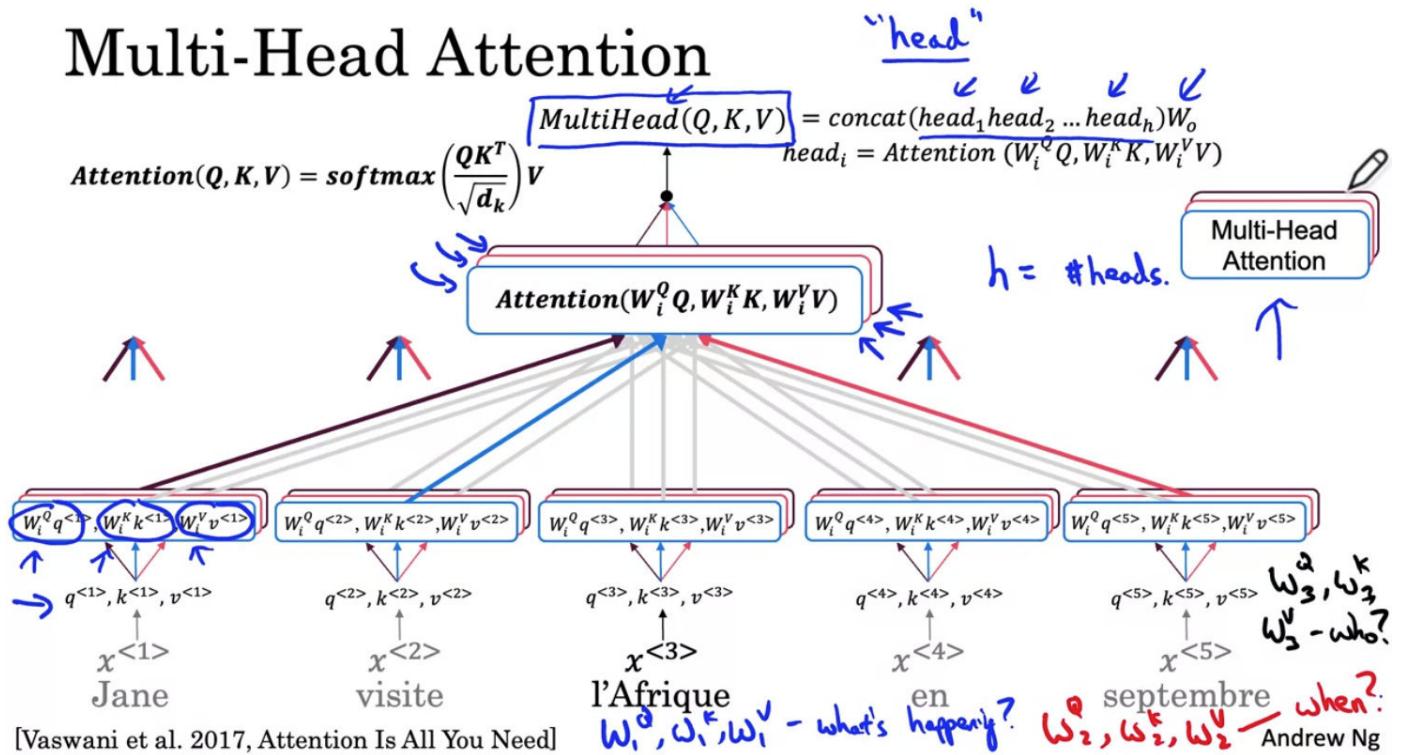
It is the vectorized representation of

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

The term  $\sqrt{d_k}$  in the denominator is just to scale the dot product so it doesn't explode. The another name for this type of attention is called **the scaled dot-product attention**.

## Multi-head attention

Each time we calculate self attention for a sequence is called a **head**.



To calculate the multiple self-attention, we firstly multiply  $q^{(t)}, k^{(t)}, v^{(t)}$  with weight matrices  $W_i^Q, W_i^K, W_i^V$ , the subscript  $i$  represents this is the weight for  $\text{head}_i$ . In literature, we use  $h$  to denote the number of heads. Using  $h = 8$  heads is common in the literature.

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

For the sake of intuition, we can think of  $W_i^Q, W_i^K, W_i^V$  as being learned to help ask and answer different questions ( $i^{th}$  question in particular). For example,  $W_1^Q, W_1^K, W_1^V$  are learned to help with asking and answering "what is happening to the word 'Africa'", and the word "visite" gives the best answer (blue arrow). Similarly,  $W_2^Q, W_2^K, W_2^V$  are learned to help with asking and answering the second question like "when is happening", so the inner product between the September key and l'Afrique query will have the highest value (red arrow). And  $W_3^Q, W_3^K, W_3^V$  may help

with third question "who has something to do with Africa", and this time Jane's value will have the greatest weight in this representation.

After calculating the computation for these heads, the concatenation of these values is used to compute the output of the multi-head attention.

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n)W_0$$

where

$$\text{head}_i = \text{Attention}(W_1^Q Q, W_1^K K, W_1^V V)$$

Recall that

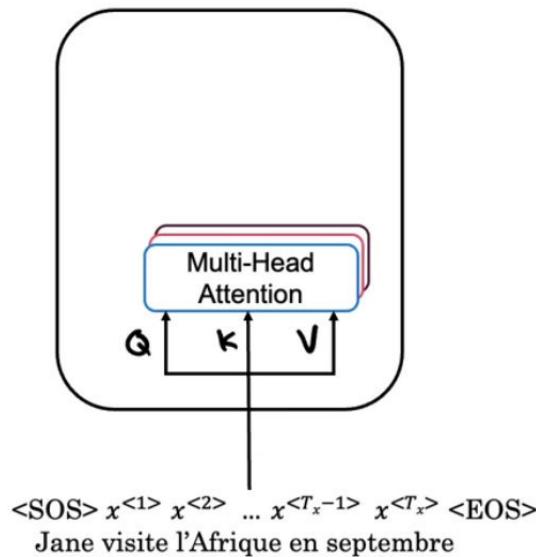
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

In practice, these multi-head attentions are computed in parallel, because no one head's value depends on the value of any other head.

## Transformer network

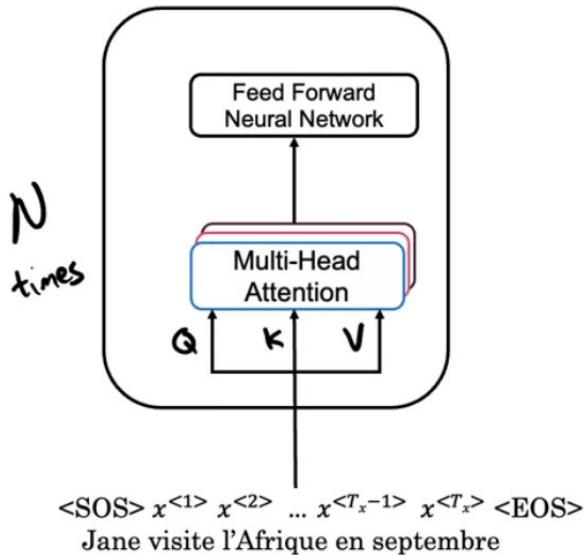
The first step in the transformer is the embeddings of the input sentence get def into an **encoder** block, and the  $Q, K, V$  are computed from the embeddings and weight matrices  $W$ .

### Encoder



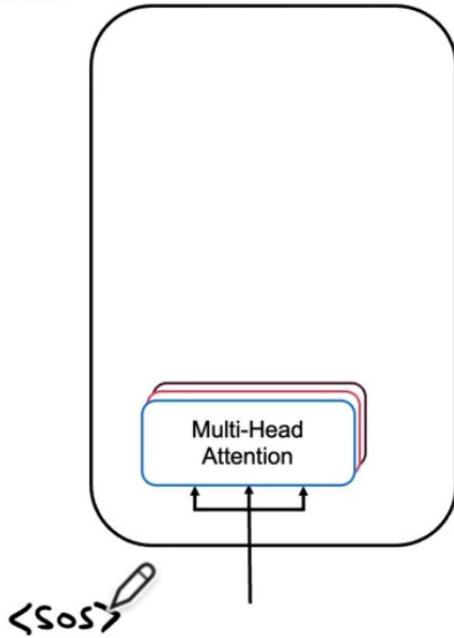
The multi-head attention layer then produces a matrix that can be passed into a feed forward neural network, which helps determine what interesting features there are in the sentence. In the original paper, this encoding block is repeated  $N$  times, typically  $N = 6$ .

## Encoder

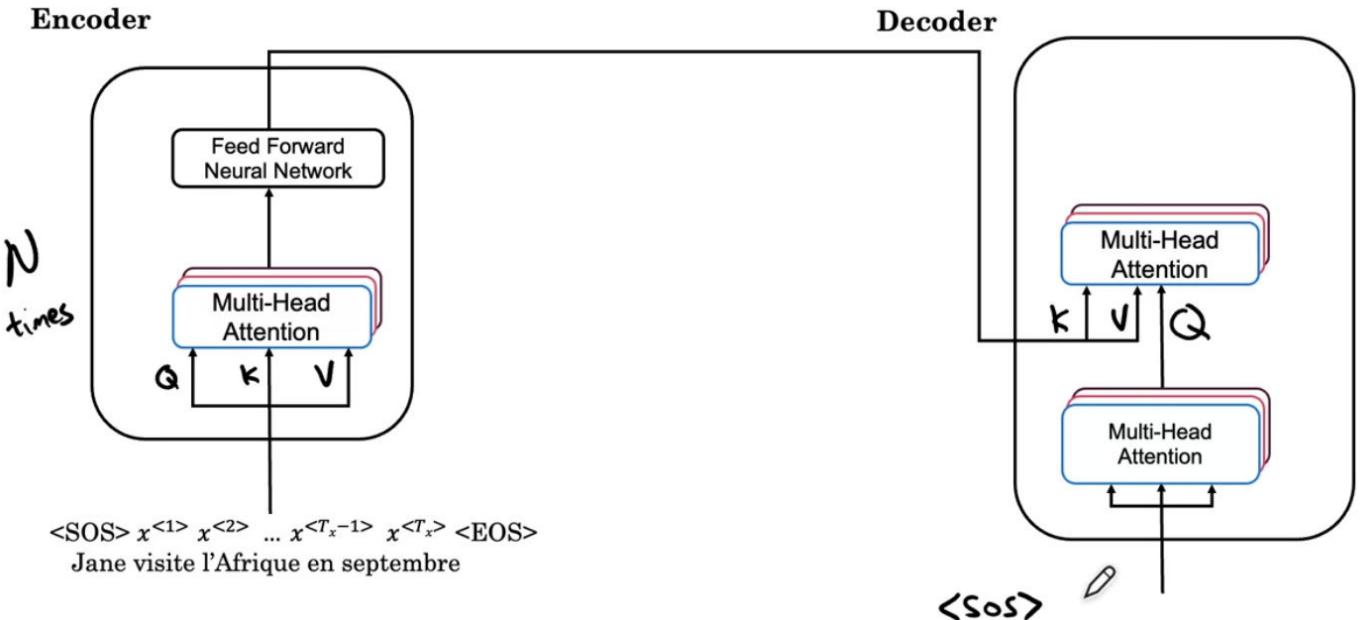


Then the output of the encoder is fed into a **decoder** block. The job of the decoder is to output the English translation. At every step, the decoder block will input the first few words that were generated by decoder and compute  $Q, K, V$  for the first multi-head attention block.

## Decoder

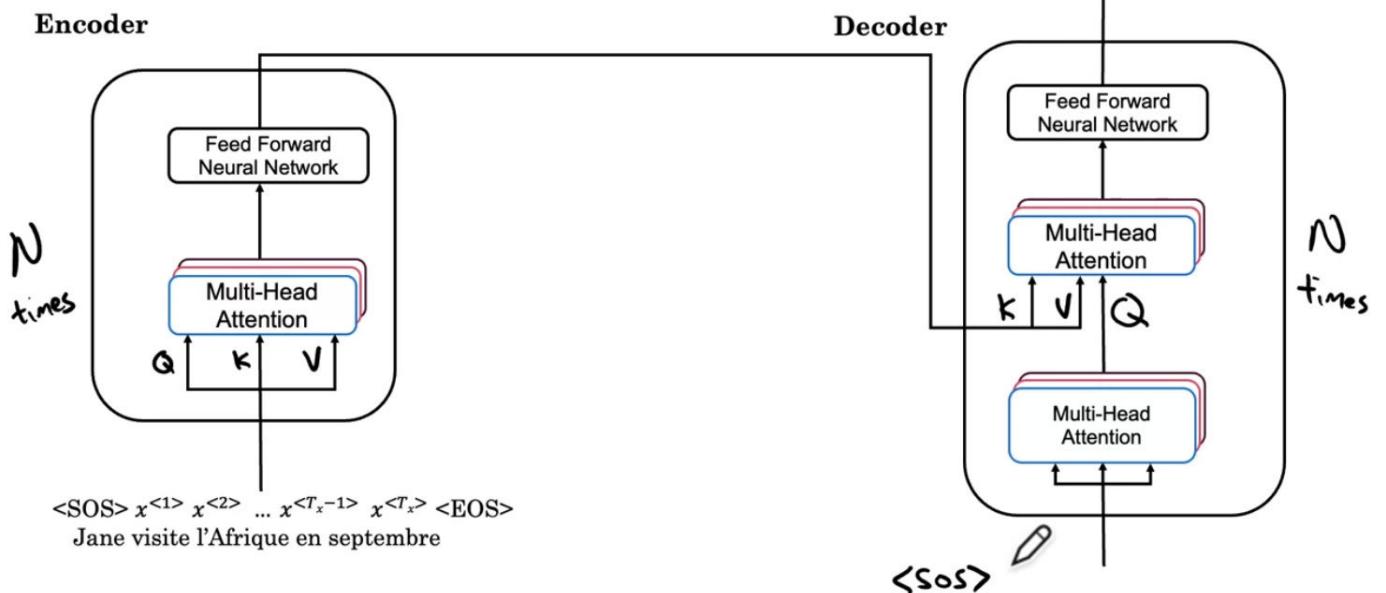


Then the first multi-head attention block will generate  $Q$  matrix for next multi-head attention block, and the matrix  $\$K,V\$$  are derived from the output of the encoder. The intuition for this step is that, the input of the first multi-head attention block is what you have translated of the sentence so far. And this will ask the query to say "what is the next word?", and then we will pull context from  $K$  and  $V$  which are translated from the French to then try to decide the next word.



The outputs of second multi-head attention block will feed to a feed forward neural network, and this decoder block is also going to be repeated  $N$  times. And the job of the feed forward neural network is to predict the next word in the sentence.

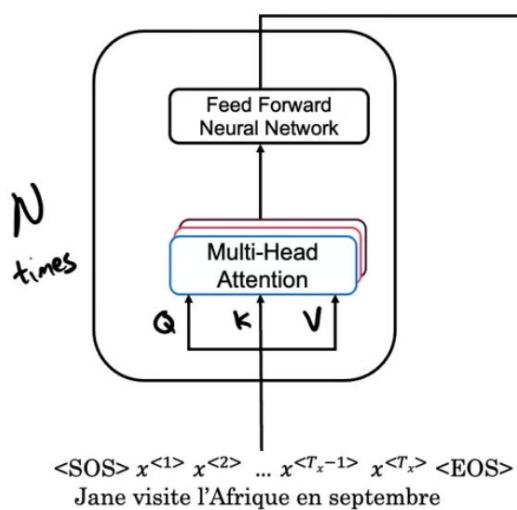
## Transformer



And the generated output will be fed to the input as well and predict the next word.

# Transformer

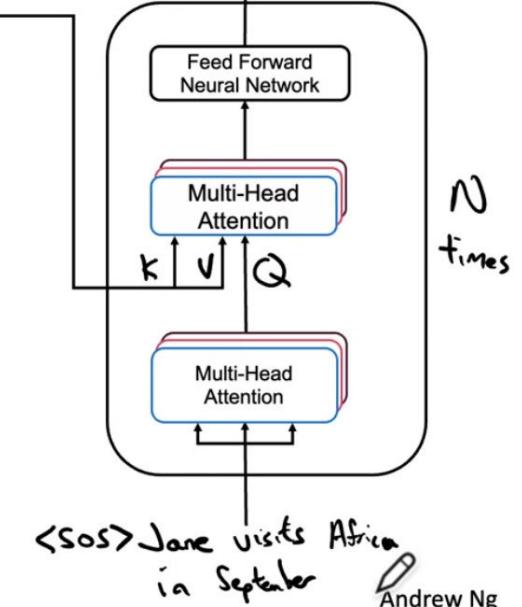
## Encoder



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

<SOS> Jane visits Africa in September <EOS>

## Decoder



This is the **main idea** of transformer.

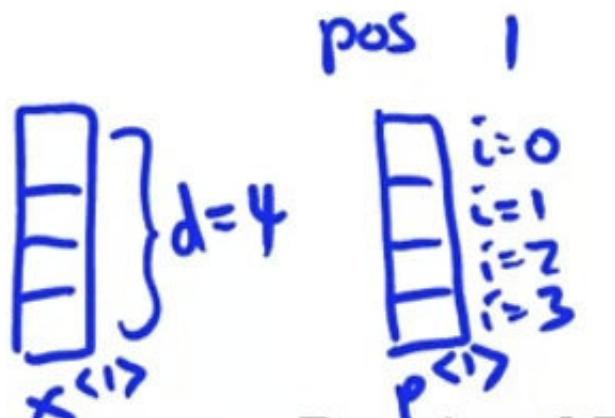
Moreover, there is **positional encoding (PE)** of the input. The position within the sentence can be extremely important to translation. However, when you train a Transformer network, you feed your data into the model all at once. While this dramatically reduces training time, there is no information about the order of your data. This is where positional encoding is useful - you can specifically encode the positions of your inputs and pass them into the network using these sine and cosine formulas:

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

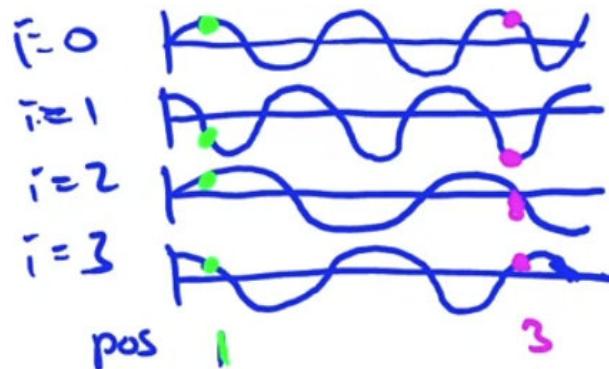
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

Where

- $d$  is the dimension of the word embeddings (how many features)
- $pos$  is the numerical position of the word ( $pos = 1$  for the word Jane).
- $i$  refers to the different dimensions of the encoding
- the position encoding for word  $x^{(t)}$  is  $p^{(t)}$ .

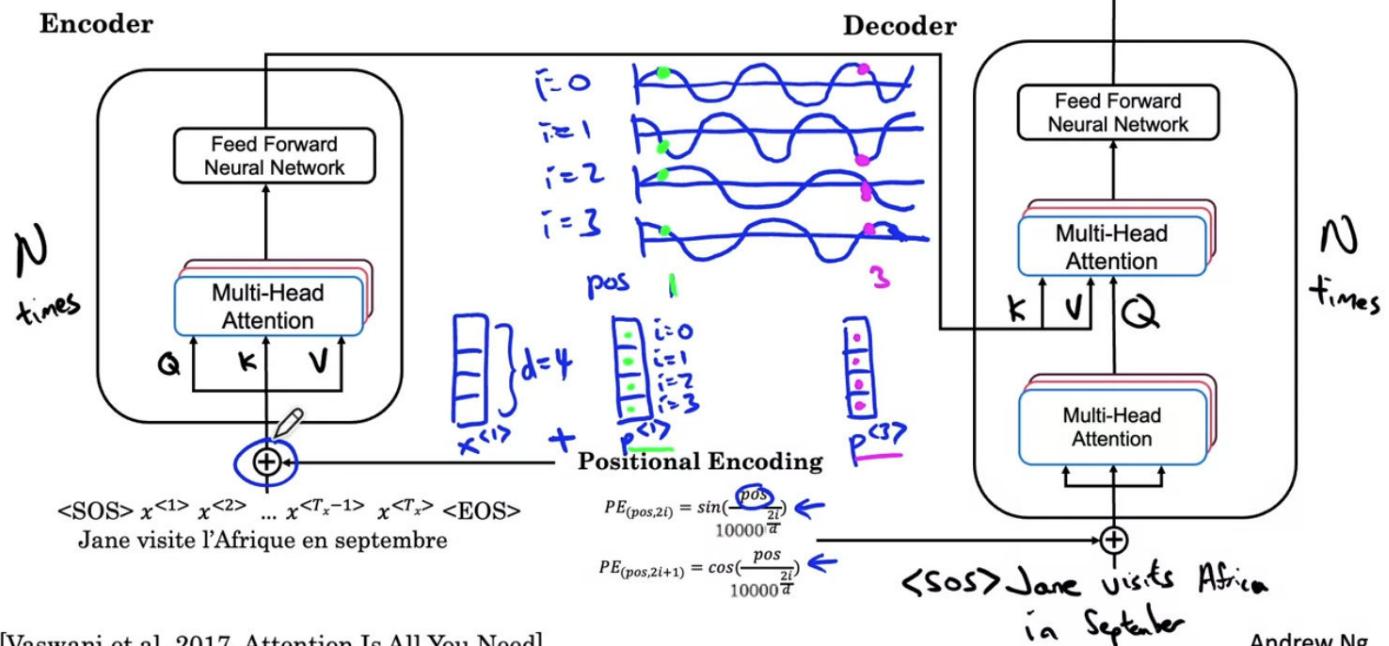


The  $q^{(t)}$  is unique because the properties of sin and cos function:



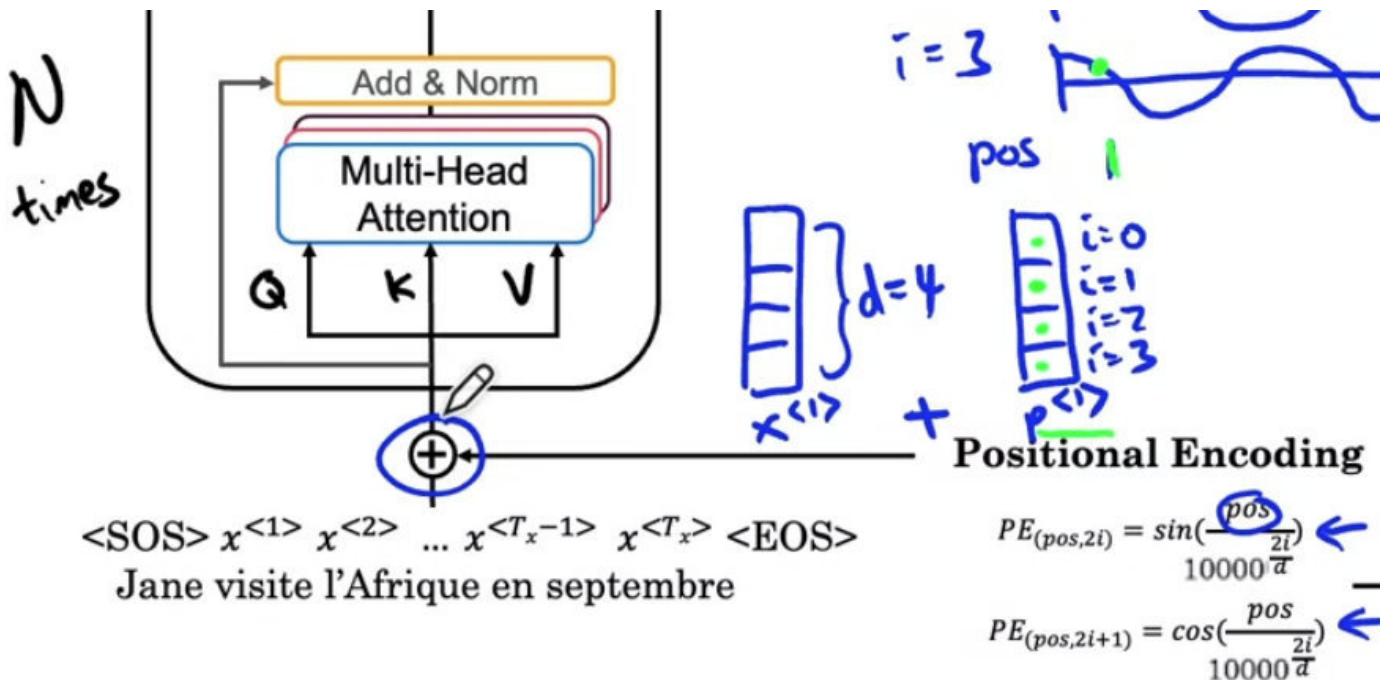
The positional encoding  $p^{(t)}$  is added directly to the  $x^{(t)}$ , so that each of the word vectors is also influenced by the position of the word. The values of the sine and cosine equations are small enough (between -1 and 1) that when you add the positional encoding to a word embedding, the word embedding is not significantly distorted.

## Transformer Details



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

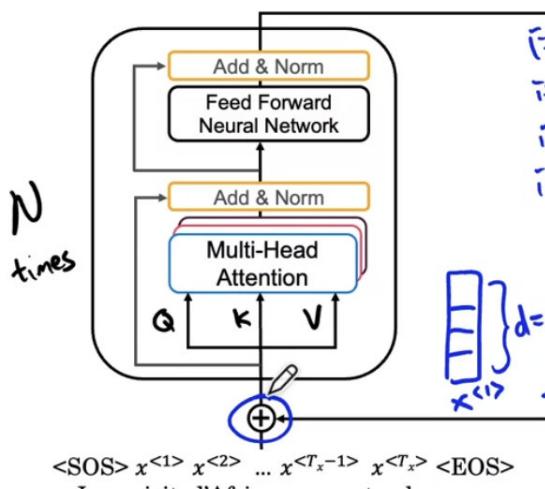
In addition to adding this positional encodings to the embeddings, or the attention representations, we would also pass them through the netword with **residual connections**. The add & norm layer is just like the batchnorm layer to speed up learning.



And these batchnorm-like layer and residual connections are repeated throughout the architecture. Finally, for the output of the decoder block, there are a linear and a softmax layer to predict the next word one at a time.

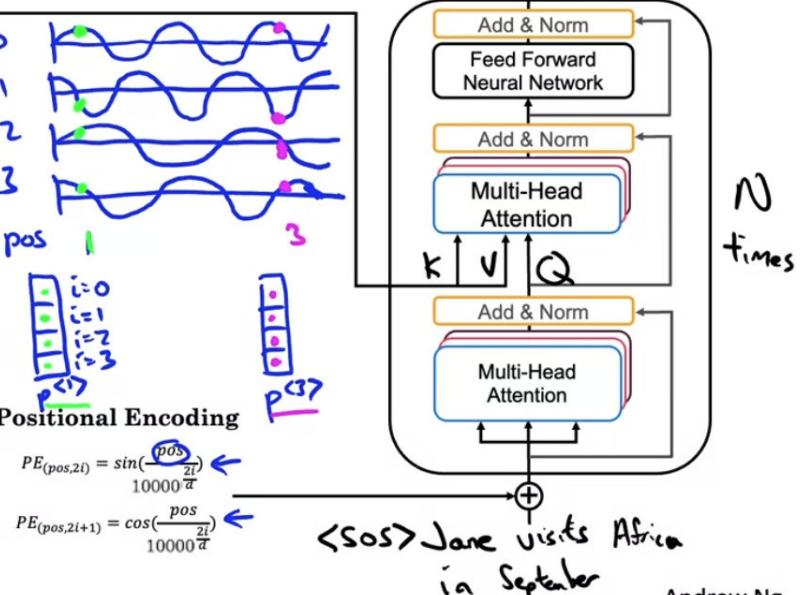
## Transformer Details

### Encoder



<SOS> Jane visits Africa in September <EOS>

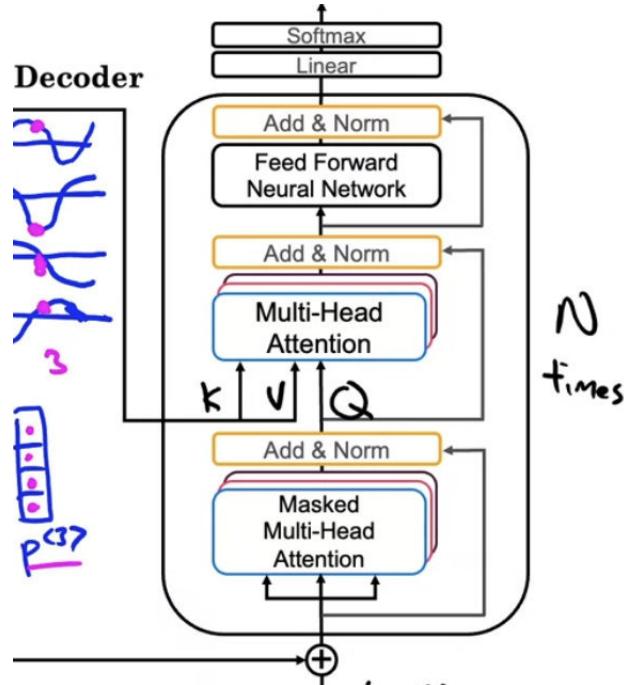
### Decoder



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

Andrew Ng

In literature on transformer, the first multi-head attention block is actually **masked multi-head attention**. This is important only for training process.

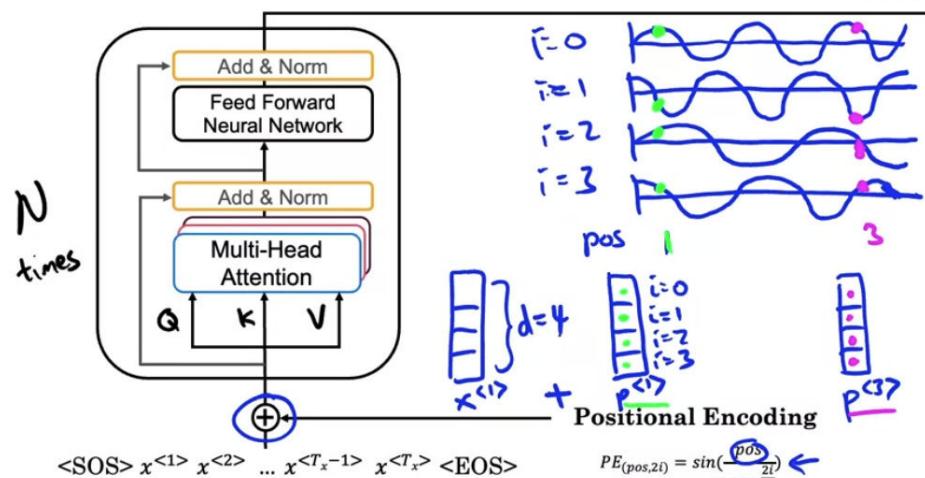


When training the network, you have the access to the entire correct outputs. So we do not need to generate the words one at a time during training. Instead, what masking does is it **blocks out** the last part of sentence to mimic what the network will need to do at test time while doing prediction. In other words, all the masked multi-head attention do is to repeatedly pretend that the network had perfectly translated first few words, and hide the remaining words to see if the network can predict the next word accurately when giving perfect inputs.

Finaly, **THIS IS THE TRANSFORMER NETWORK.**

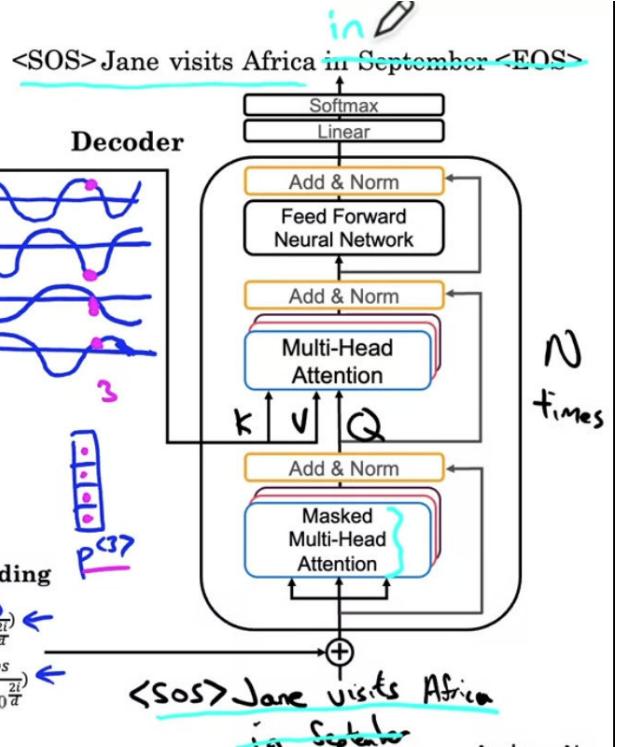
## Transformer Details

### Encoder



<SOS> Jane visits Africa in September <EOS>

### Decoder



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

Andrew Ng

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$