



deeplearning.ai

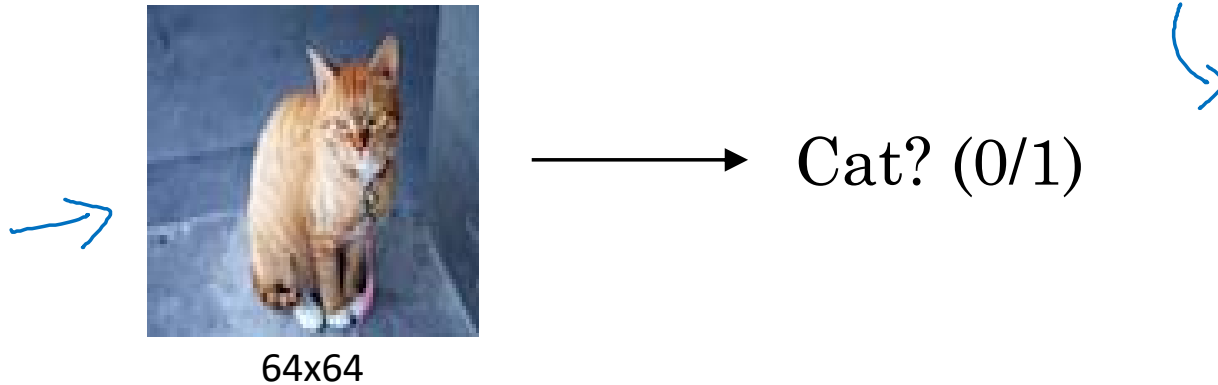
# Convolutional Neural Networks

---

## Computer vision

# Computer Vision Problems

## Image Classification



You have a picture and you want it painted in a different style.

## Neural Style Transfer



## Object detection



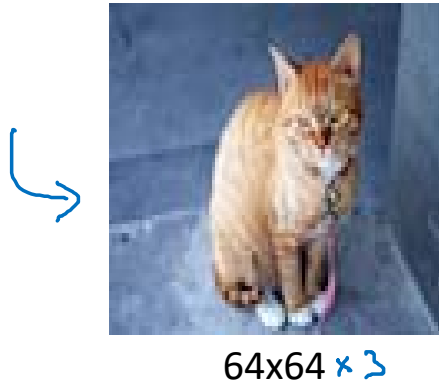
: we not just need to classify them but also need to know where in picture these objects are so that we can draw box around them to detect them.



Andrew Ng

# Deep Learning on large images

one of the prob in comp vision is that img can get really big

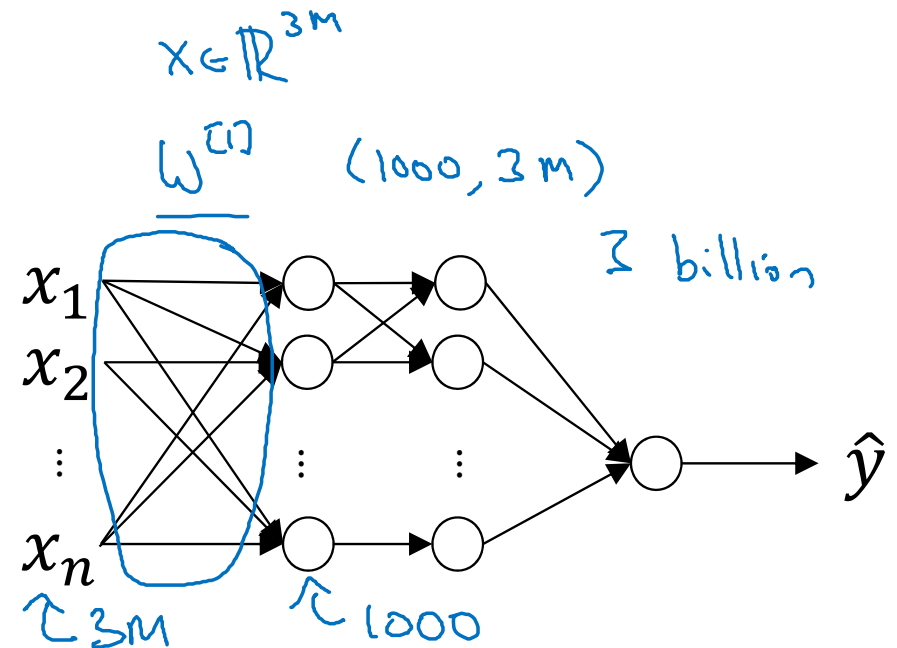


→ Cat? (0/1)

12288



1000x1000 x 3  
= 3 million



With 3 billion data is difficult to prevent a NN from overfitting and also the computational requirements and memory requirements to train a NN with 3 billion is just a bit infeasible. But for computer vision you don't want to be stuck using little images; you want to use large images. To do that you need to better implement the convolution operation, which is one of the fundamental building blocks of convolutional NN.

Andrew Ng



deeplearning.ai

# Convolutional Neural Networks

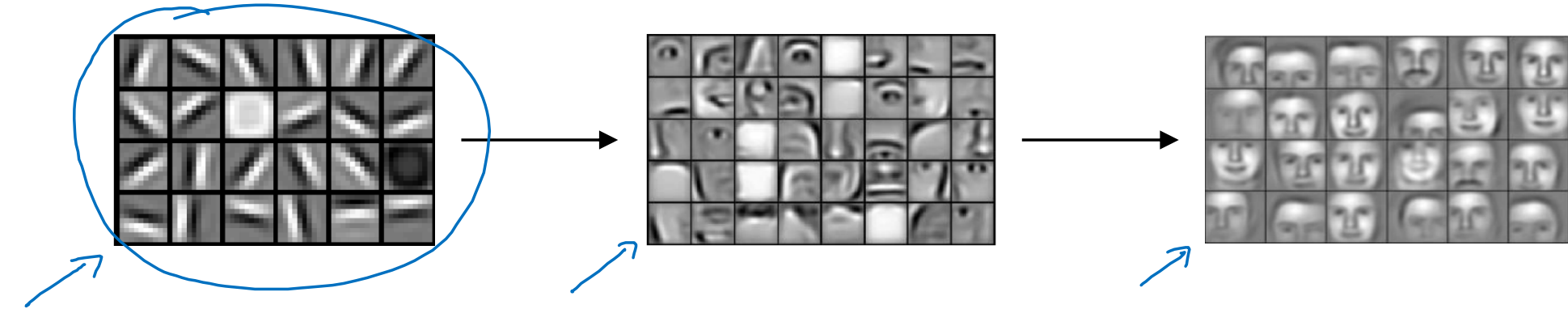
---

## Edge detection example

The convolution operation is one of the fundamental building blocks of a convolutional neural network. Using edge detection as the motivating example in this video, you will see how the convolution operation works.

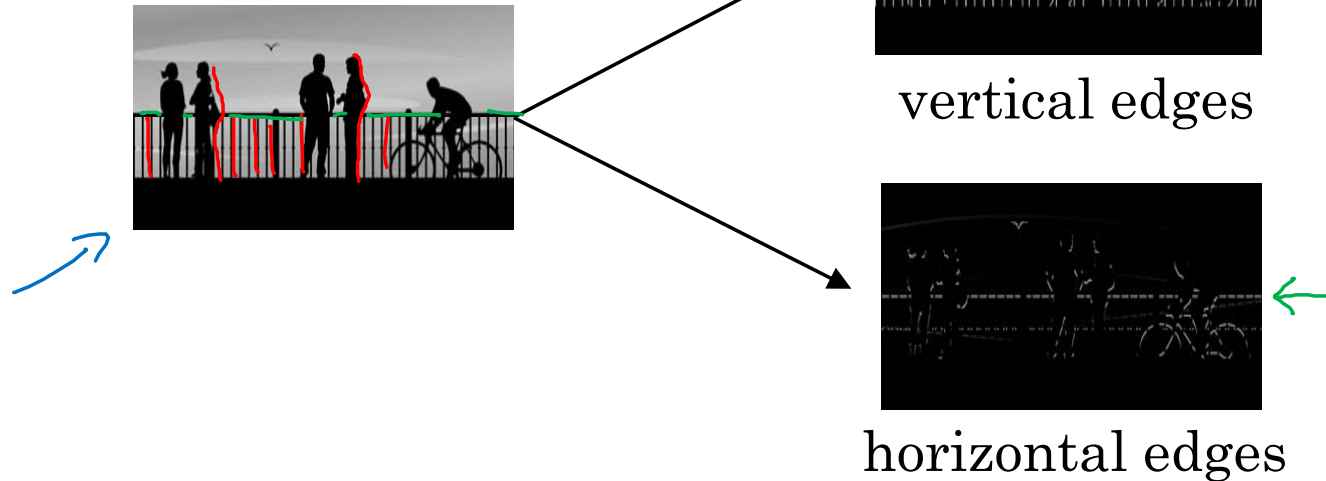
# Computer Vision Problem

In previous videos, I have talked about how the early layers of the neural network might detect edges and then some later layers might detect cause of objects and then even later layers may detect cause of complete objects like people's faces in this case.



In this video we gone see how you can detect edges in an img.

Given a picture like that



Andrew Ng

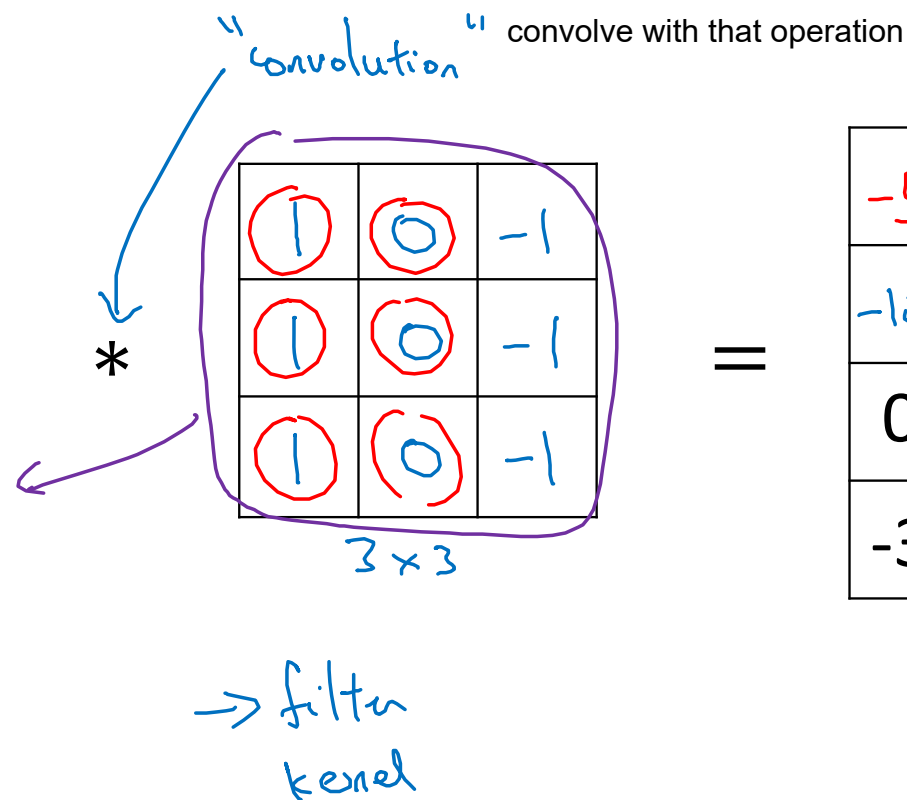
# Vertical edge detection

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

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

6x6

Here is a 6x6 grey scale img



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

4x4

python: conv-forward

tensorflow: tf.nn.conv2d

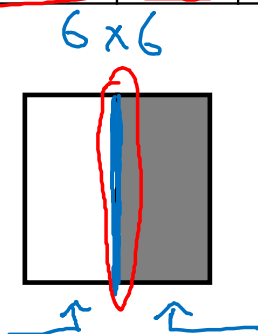
keras: conv2D

Andrew Ng

Why is it doing edge detection.

# Vertical edge detection

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	<u>10</u>	<u>10</u>	<u>0</u>	0	0
10	<u>10</u>	<u>10</u>	<u>0</u>	0	0
10	<u>10</u>	<u>10</u>	<u>0</u>	0	0



\*

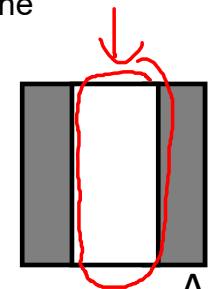
1	0	-1
1	0	-1
1	0	-1

3x3

=

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

4x4  
this white part correspond  
to having detected the  
middle edge



\*



So we did vertical  
edge detection which  
is one of the building blocks of CNN  
U can detect edges using that convolution operation

Andrew Ng



deeplearning.ai

# Convolutional Neural Networks

---

More edge  
detection



# 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

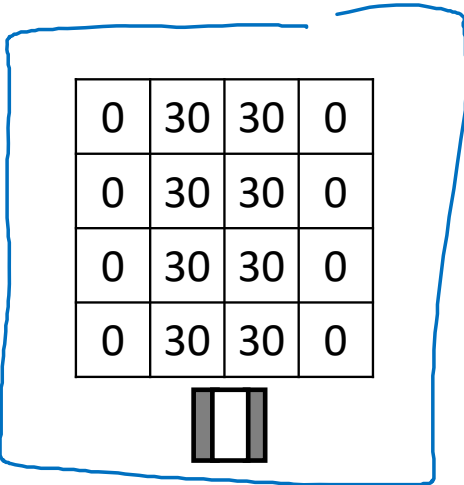


\*


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



=




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

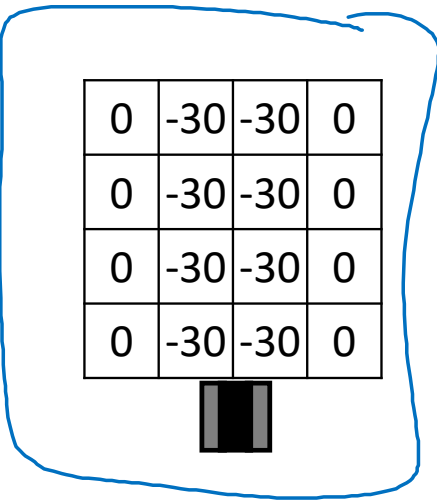


\*


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




=



0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0




# Vertical and Horizontal Edge Detection



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

Vertical



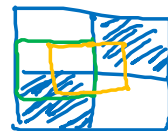
1	1	1
0	0	0
-1	-1	-1

Horizontal

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

6x6

\*



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

=

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

# Learning to detect edges

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



the advantage of this is that it puts some more weight to the center row

→

1	0	-1
2	0	-2
1	0	-1

Sobel filter



convolution



$W_1$	$W_2$	$W_3$
$W_4$	$W_5$	$W_6$
$W_7$	$W_8$	$W_9$

3x3

U can learn these by putting them as parameters and use backprop

=

45°  
70°  
73°

this has other properties

3	0	-3
10	0	-10
3	0	-3

Scharr filter




it can maybe learn to detect edges that are of different degrees

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

The idea that you can put these parameters to be learned is one of the most powerful idea in computer vision later we gone see how u use backprop to learn these 9 parameters

In order to build deep neural networks one modification to the basic convolutional operation that you need to really use is padding



deeplearning.ai

# Convolutional Neural Networks

---

## Padding

These are the 2 drawbacks:

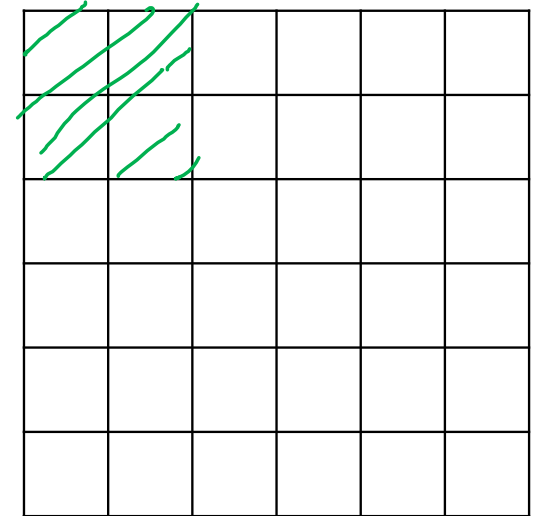
Output size is smaller if you don't apply padding

# Padding

- Shrinky output
- throw away info from edge

What u do for this is pad img with one boarder around the edges so instead of 6x6 u have 8x8 img and with convolut operation u get a 6x6 output

$$\begin{matrix} \begin{matrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{matrix} & = & \begin{matrix} \text{6x6 grid with 3x3 kernel and padding} \end{matrix} \\ 3 \times 3 & & 6 \times 6 \\ f \times f & & \end{matrix}$$



6x6

2 downsides;  
- img shrinks  
- pixels in corners are used much less in output

~~4x4~~

this will be the dimation of output

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

$$6 - 3 + 1 = 4$$

$$\begin{matrix} 6 \times 6 \rightarrow 8 \times 8 \\ n \times n \end{matrix}$$

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

$$\begin{matrix} n + 2p - f + 1 \times n + 2p - f + 1 \\ 6 + 2 - 3 + 1 \times \text{---} = 6 \times 6 \end{matrix}$$

Andrew Ng

by convention u pad with 0 and p is the padding amout and in this case is 1

So u surround that img with these layer of 0.

# Valid and Same convolutions

→ no padding

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

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

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

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

$$3 \times 3$$

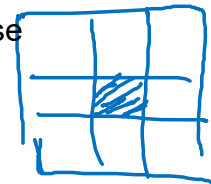
$$p = \frac{3-1}{2} = 1$$

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

$$5 \times 5$$
$$f=5$$

$f$  is usually odd

This is how choose pad in order not change size.



$$1 \times 1$$
$$3 \times 3$$
$$5 \times 5$$
$$7 \times 7$$



deeplearning.ai

# Convolutional Neural Networks

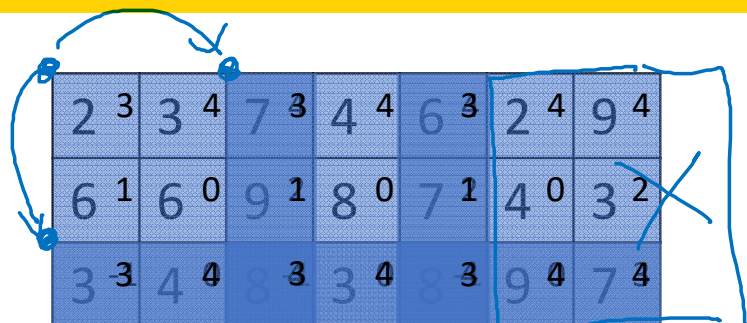
---

Strided  
convolutions

Strided conv is another building block of the convolution as using convolution NN

Say we want to convolve that 7x7 with that 3x3 filter, instead of doing the usual way we gone do it with stride of 2, so we gone move box not just as before by 1 but by 2 this time.

# Strided convolution



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

7x7


the filter must lay entirely within the img or img plus padding  
thats a convention

3	4	4
1	0	2
-1	0	3

3x3

Stride = 2

=



91	100	83
69	91	127
44	72	74

3x3

$\lfloor 2 \rfloor = \text{floor}(2)$

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

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

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

what if the fraction is not an integer??  
in that case we gone round to floor.

Andrew Ng  
this is why we end up  
with this 3x3 matrix



# Summary of convolutions

$n \times n$  image       $f \times f$  filter

padding  $p$       stride  $s$

Output Size:

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

# Technical note on cross-correlation vs. convolution

Just don't confuse the notations of the convolution with the cross-correlation.

Convolution in math textbook:

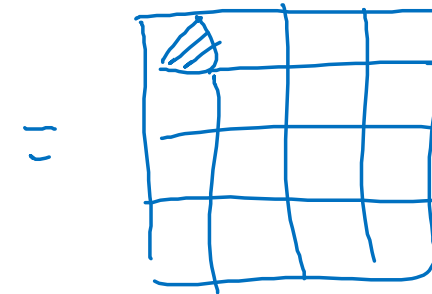
2 <sup>7</sup>	3 <sup>2</sup>	7 <sup>5</sup>	4	6	2
6 <sup>9</sup>	6 <sup>0</sup>	9 <sup>4</sup>	8	7	4
3 <sup>-1</sup>	4 <sup>1</sup>	8 <sup>3</sup>	3	8	9
7	8	3	6	6	3
4	2	1	8	3	4
3	2	4	1	9	8

*	3	4	5
	1	0	2
	-1	9	7

7	2	5
9	0	4
-1	1	3

In math text  
books u flip first the filter  
horizontally and vertically

we do not do it in machine learning



$$(A * B) * C = A * (B * C)$$



deeplearning.ai

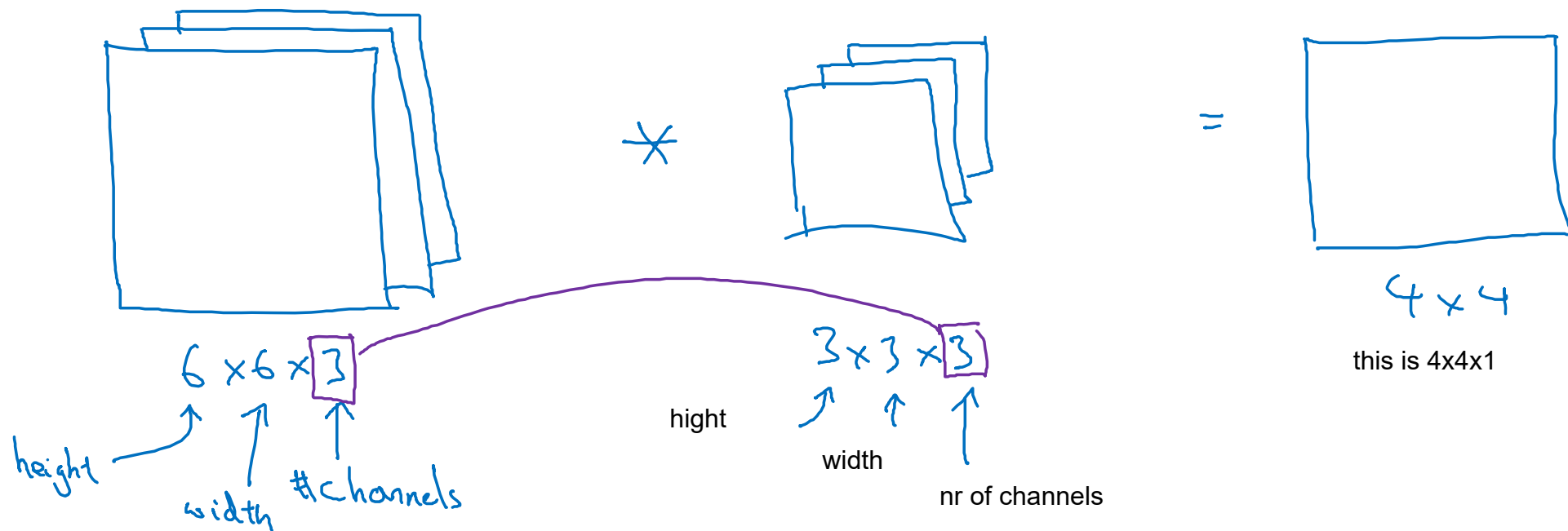
# Convolutional Neural Networks

---

## Convolutions over volumes

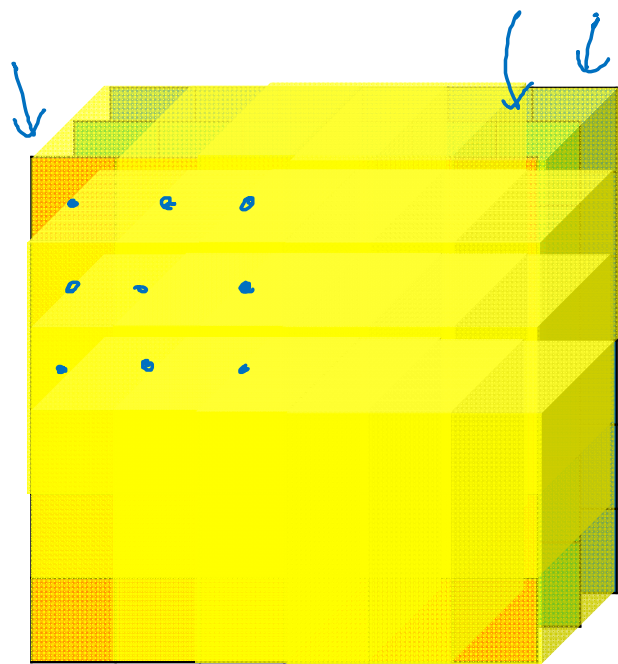
Now we will see how to implement convolutions over 3D img.

# Convolutions on RGB images

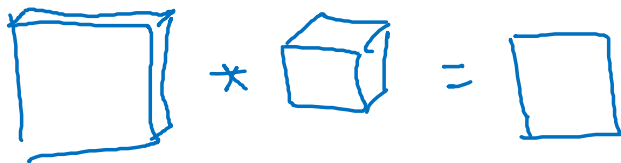


To detect edges in this rgb img u convolve with a 3x3x3 filter

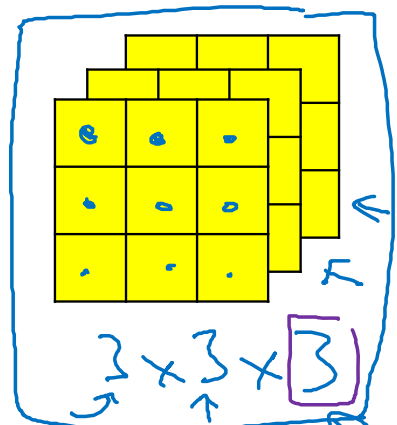
# Convolutions on RGB image



$6 \times 6 \times 3$

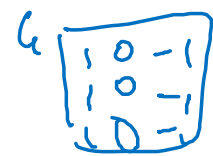
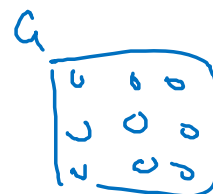
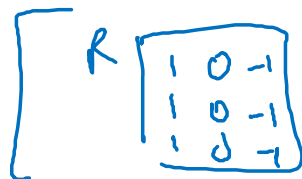
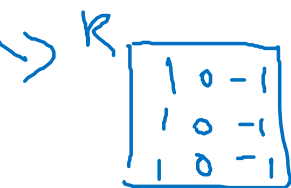


\*

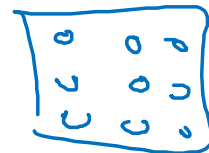


$3 \times 3 \times 3$

27 numbers



3



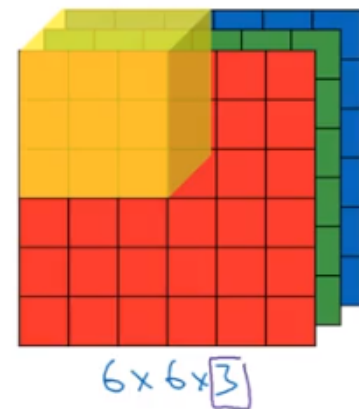
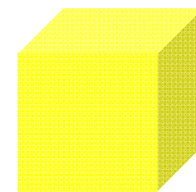
$\rightarrow 3 \times 3 \times 3$   
this detects edges  
only red channel

3



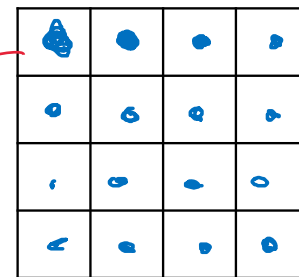
$\rightarrow 3 \times 3 \times 3$

Andrew Ng  
this detects edges in  
r g b channels



$6 \times 6 \times 3$

=



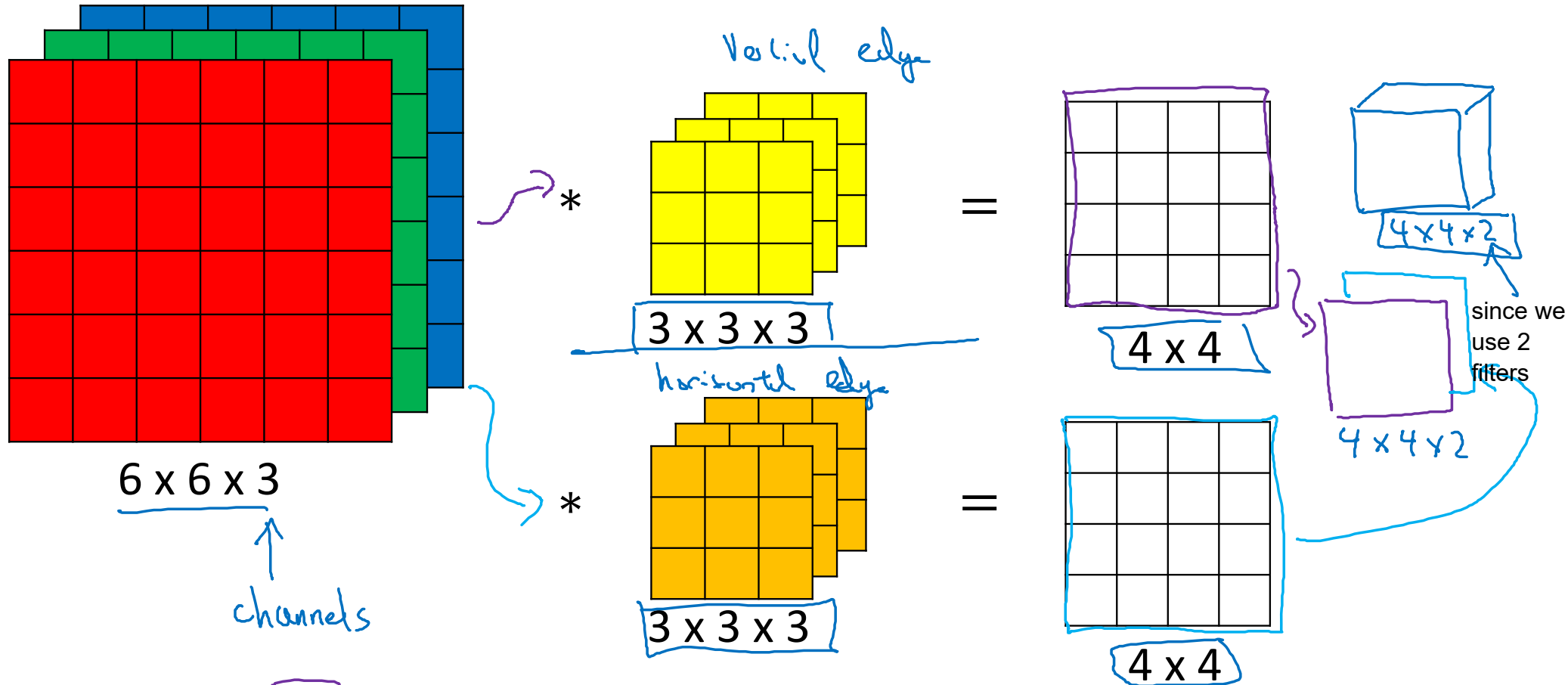
$4 \times 4$

27 multiplications and  
sums to get this number

It is normal sometime to have a filter that  
looks only at green channel or red channel or blue channel.

What if we want to detect vertical and horizontal edges, so not just one at a time

# Multiple filters



Summary:  $n \times n \times n_c$   $\times$   $f \times f \times n_c$   $\rightarrow$   $\frac{n-f+1}{4} \times \frac{n-f+1}{4} \times \frac{n_c'}{2}$

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

We are now ready to see how to build a one layer of a convolution NN. That's through an example



deeplearning.ai

# Convolutional Neural Networks

---

## One layer of a convolutional network

## Example of a layer

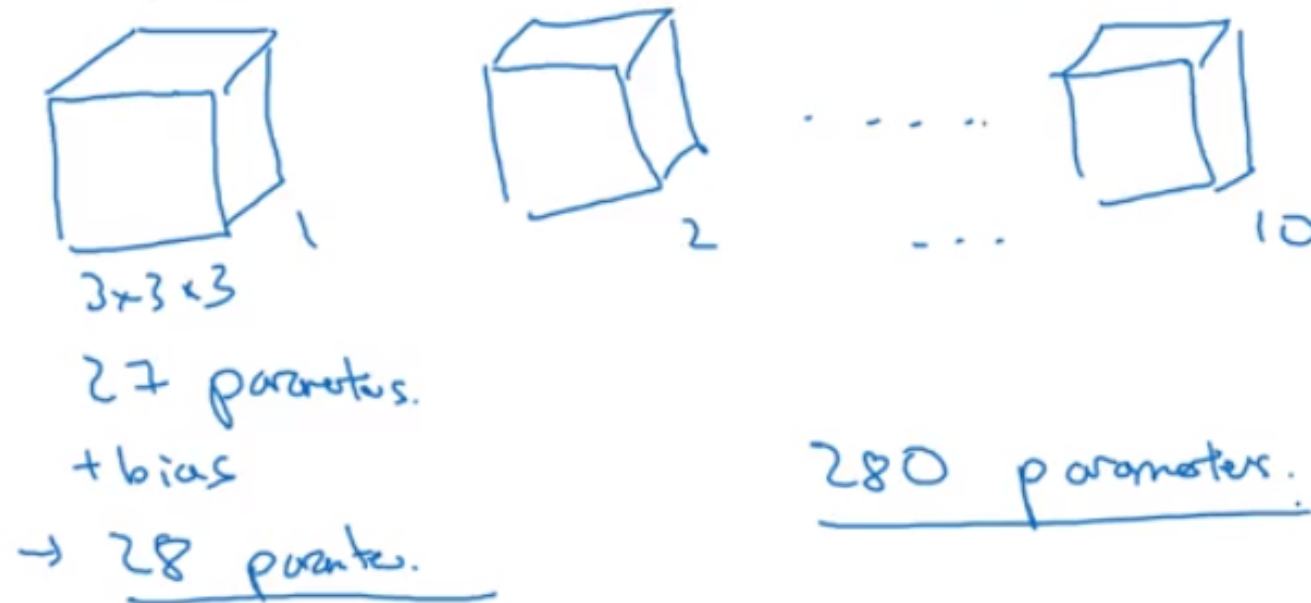


Andrew Ng



# Number of parameters in one layer

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



Andrew Ng

No matter what the size of the img is the size of the parameters remain fixed 280  
and you can use these 10 filters to detect features, vertical edges horizontal ones maybe other features

# Summary of notation

If layer l is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

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

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

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

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

bias:  $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$  ← #filters in layer l.

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

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

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

this means take floor

In the last video, you saw the building blocks of a single layer, of a single convolution layer in the ConvNet. Now let's go through a concrete example of a deep convolutional neural network.



deeplearning.ai

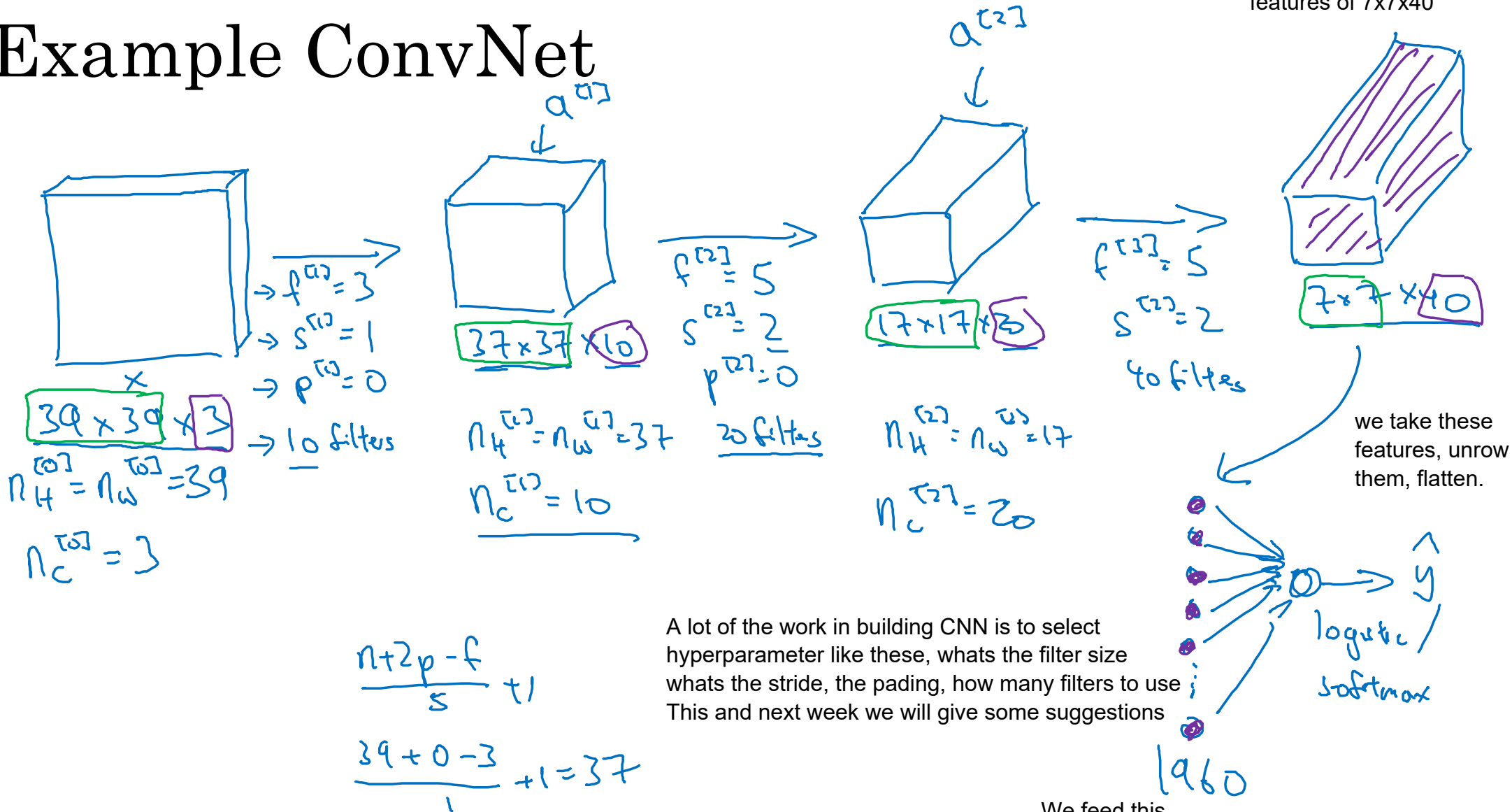
# Convolutional Neural Networks

---

## A simple convolution network example

# Example ConvNet

So what we have done is take that img and transform into this set of features of 7x7x40




A lot of the work in building CNN is to select hyperparameter like these, what's the filter size what's the stride, the padding, how many filters to use This and next week we will give some suggestions

We feed this to a logistic regression or a softmax unit, depending on whether u trying to recognize cat, no cat or any other object

Andrew Ng

# Types of layer in a convolutional network:

- Convolution (conv) ←
  - Pooling (pool) ←
  - Fully connected (FC) ←
- 

We have talked only about the conv layer  
we will see later the pooling layer and the fully  
connected layer

It is possible to build without the pool and fc  
most architecture NN will have a few pool  
layers and FC layers

Other than convolutional layers, ConvNets often also use pooling layers to reduce the size of the representation, to speed the computation, as well as make some of the features that detects a bit more robust.



deeplearning.ai

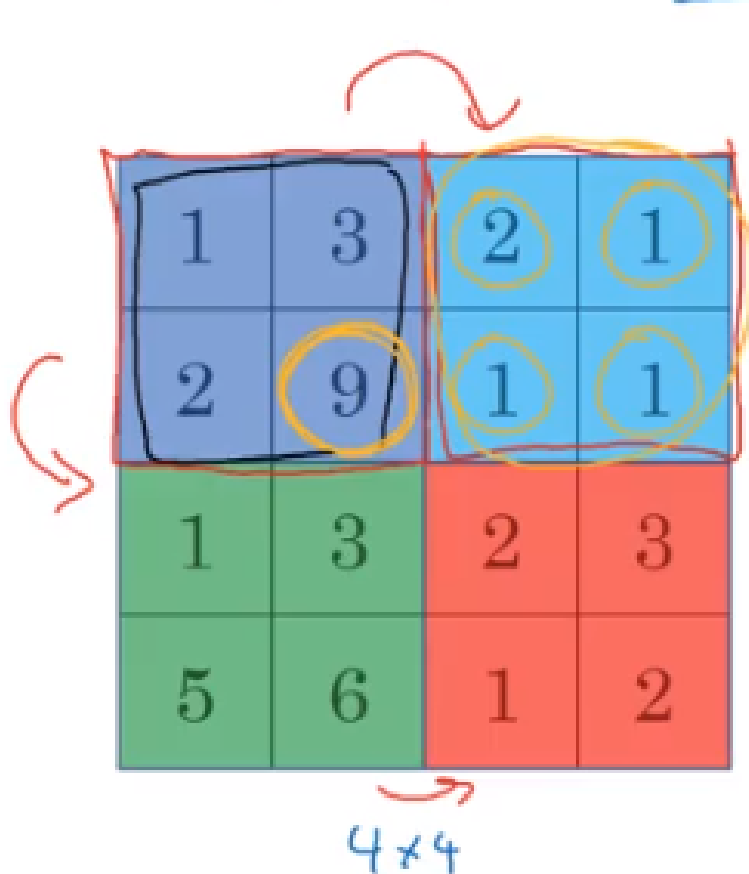
# Convolutional Neural Networks

---

## Pooling layers

# Pooling layer: Max pooling

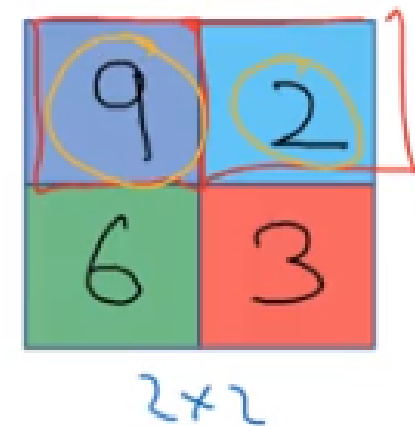
Suppose we want to apply a type of pooling called max pooling



Take input and break into four regions

Intuition is not that great but a lot of people use it as it works well, he does not know very well whether that is the real intuition.

this will be the output



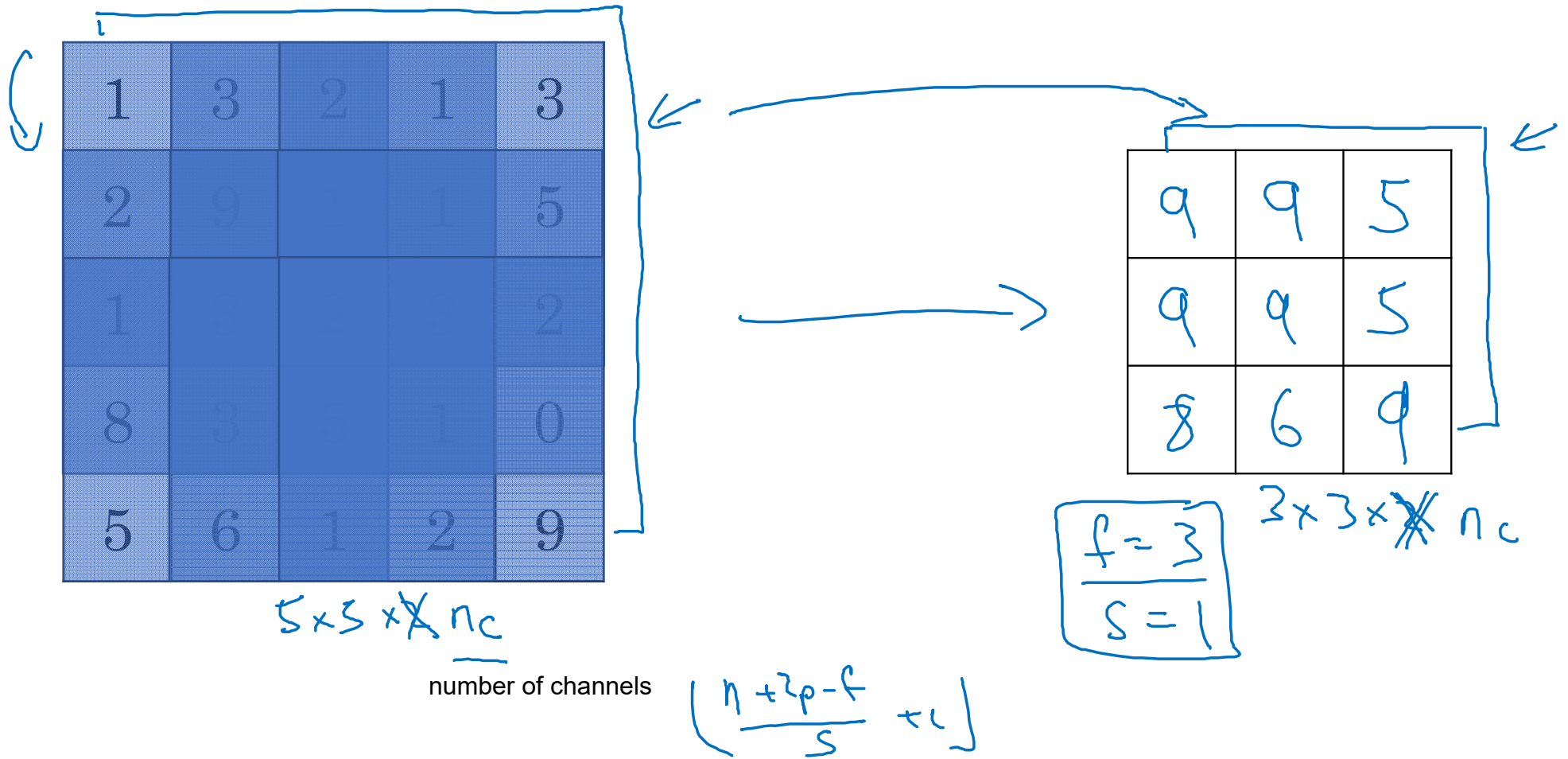
in the output, each will be the max from the four regions in the input

Hypoparameters:  
 $f = 2$   
 $s = 2$

Interesting fact about max pooling is that it has a set of parameters but it has nothing to learn.

So here's the intuition behind what max pooling is doing. If you think of this four by four region as some set of features, the activations in some layer of the neural network, then a large number, it means that it's maybe detected a particular feature. So, the upper left-hand quadrant has this particular feature. It maybe a vertical edge or maybe a higher or whisker if you trying to detect a [inaudible]. Clearly, that feature exists in the upper left-hand quadrant. Whereas this feature, maybe it isn't cat eye detector. Whereas this feature, it doesn't really exist in the upper right-hand quadrant. So what the max operation does is a lots of features detected anywhere, and one of these quadrants, it then remains preserved in the output of max pooling.

# Pooling layer: Max pooling





# Pooling layer: Average pooling

take avg instead of max  
max pooling is used much more

1	3	2	1
2	9	1	1
1	4	2	3
5	6	1	2



3.75	1.25
4	2

$$f=2$$

$$s=2$$

$$\underline{7 \times 7 \times 1000} \rightarrow 1 \times 1 \times 1000$$

# Summary of pooling

Hyperparameters:

f : filter size

s : stride

Max or average pooling

$$f=2, s=2$$

$$f=3, s=2$$

~~$\Rightarrow p$ : padding~~

No parameters to learn!

input

$$n_H \times n_W \times \underline{n_C}$$



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

$$\times \underline{n_C}$$

output



deeplearning.ai

# Convolutional Neural Networks

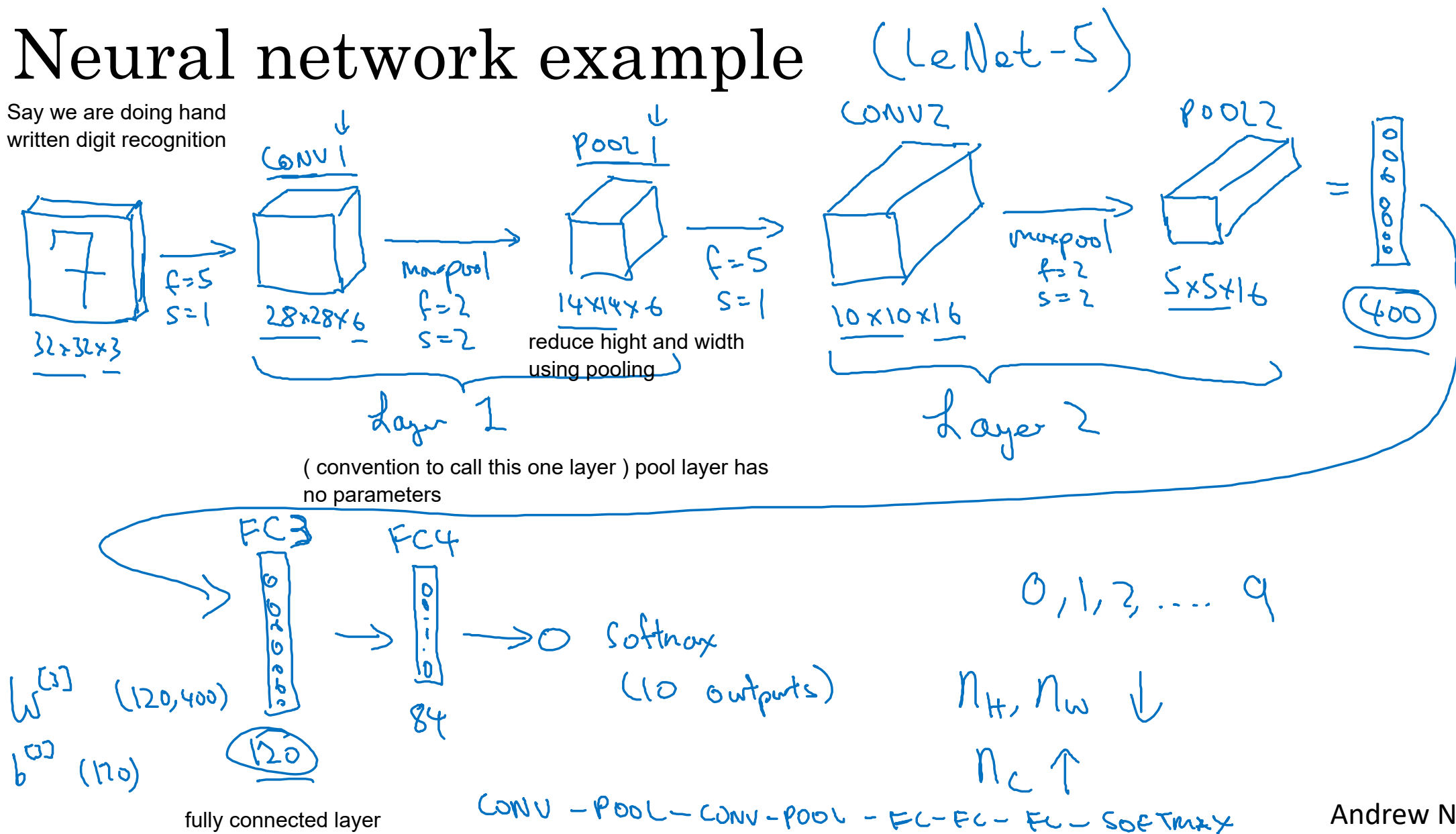
---

## Convolutional neural network example

What we gone do here is inspired by LeNet -5 , very similar to that

# Neural network example

Say we are doing hand written digit recognition



Andrew Ng

# Neural network example

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



deeplearning.ai

# Convolutional Neural Networks

---

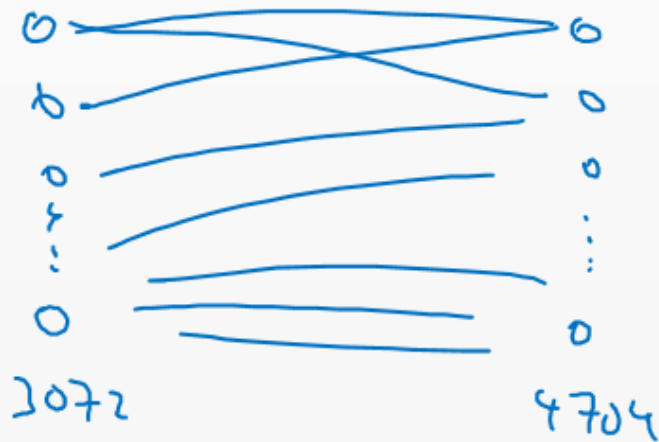
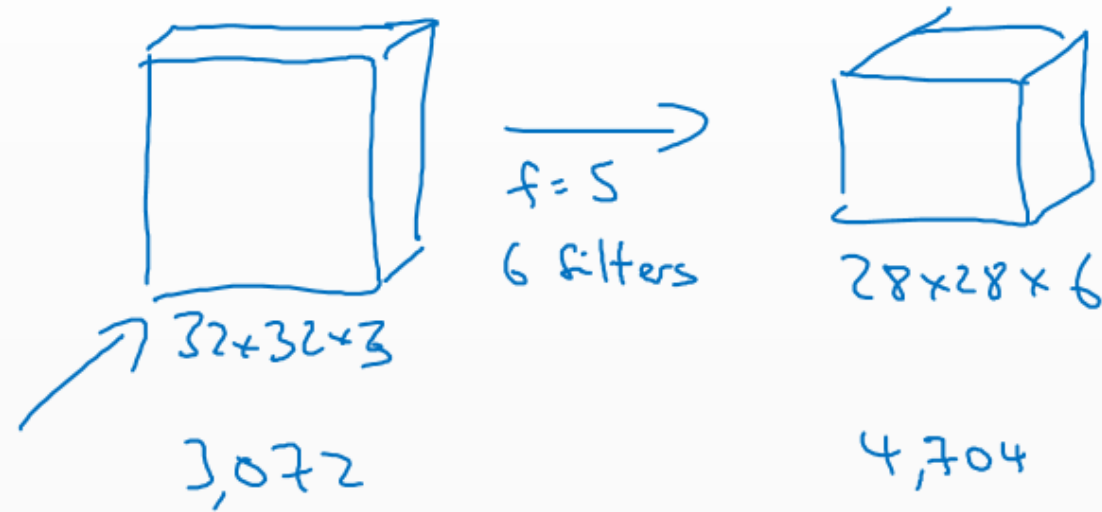
## Why convolutions?

# Why convolutions

I think there are two main advantages of convolutional layers over just using fully connected layers.

And the advantages are:

- parameter sharing and
- sparsity of connections.



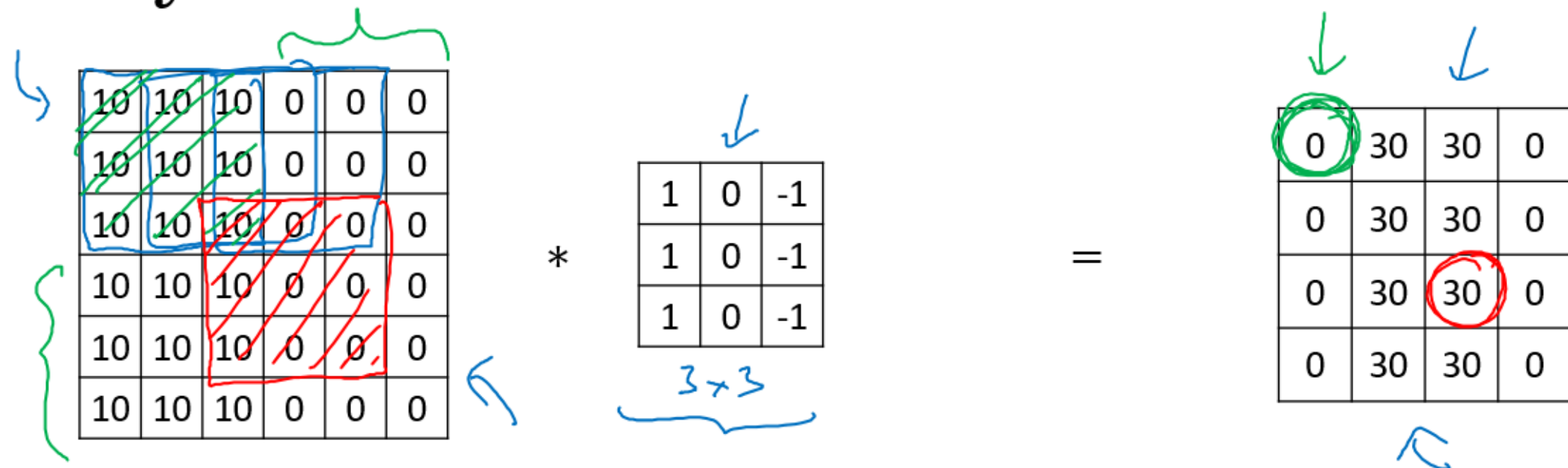
$$5 \times 5 = 25$$

$$26$$

$$6 \times 26 = 156 \text{ parameters}$$

$$3,072 \times 4,704 \approx \underline{14M}$$

# Why convolutions



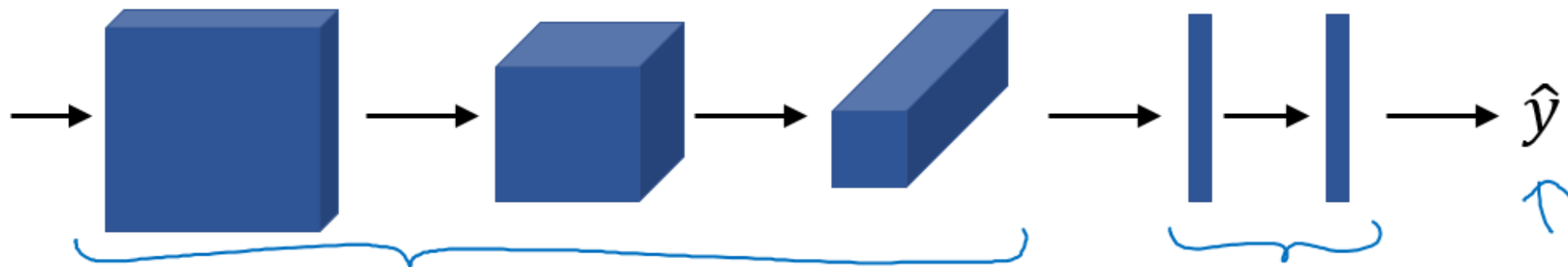
**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

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



# Putting it together

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



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

Use gradient descent to optimize parameters to reduce  $J$