

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# Convolutional Neural Networks

---

## Computer vision

# Computer Vision Problems

## Image Classification

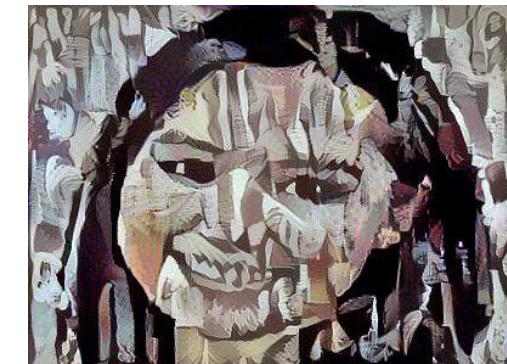
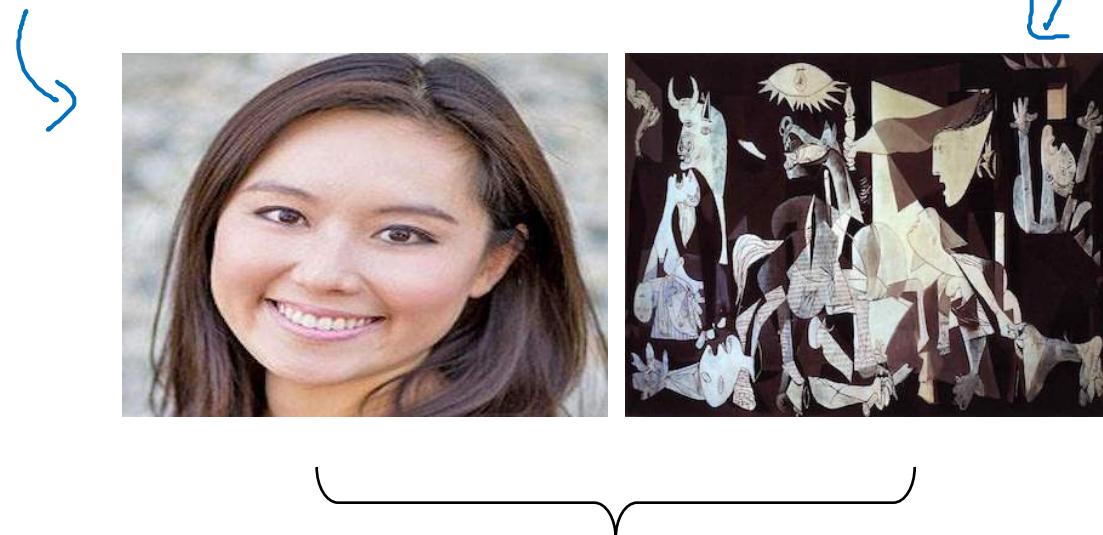


Cat? (0/1)

## Object detection



## Neural Style Transfer



Andrew Ng

# Deep Learning on large images



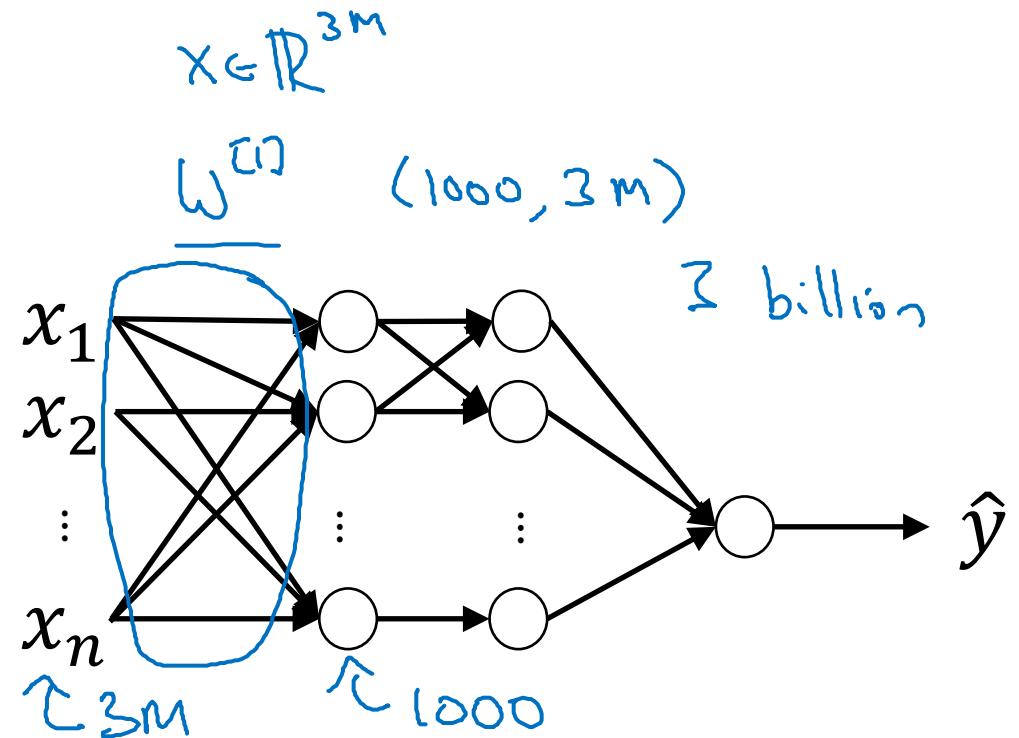
$64 \times 64 \times 3$

→ Cat? (0/1)

12288



$1000 \times 1000 \times 3$   
= 3 million





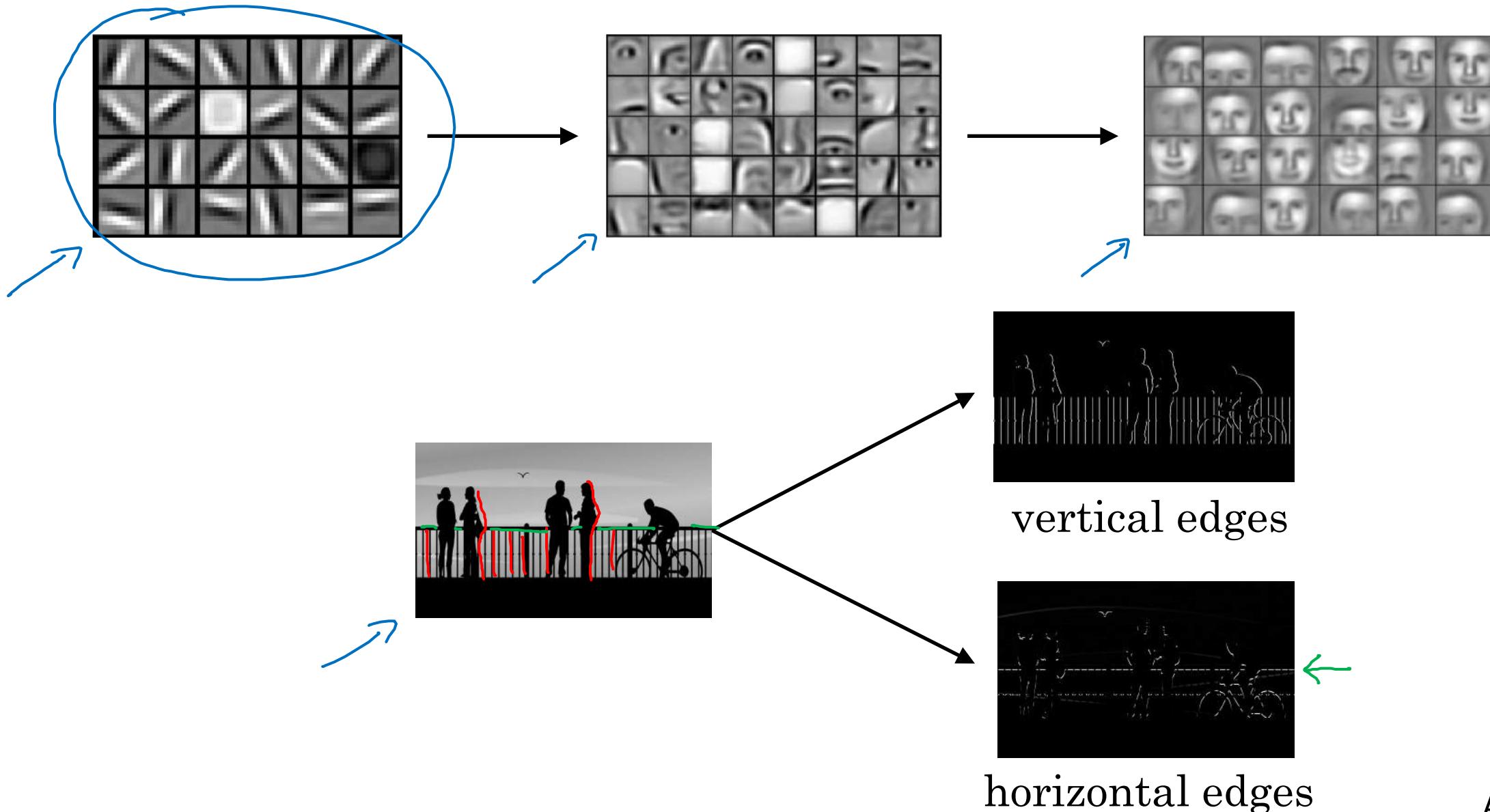
deeplearning.ai

# Convolutional Neural Networks

---

## Edge detection example

# Computer Vision Problem

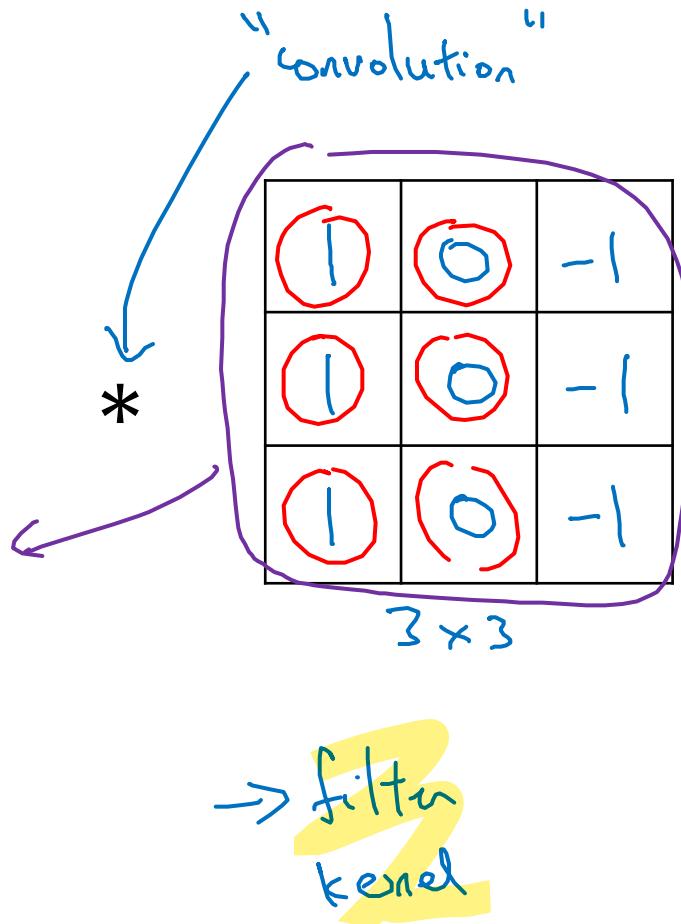


# Vertical edge detection

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

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

$6 \times 6$



=

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

$4 \times 4$

# Vertical edge detection

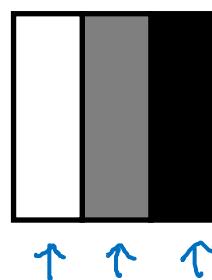
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0



$3 \times 3$

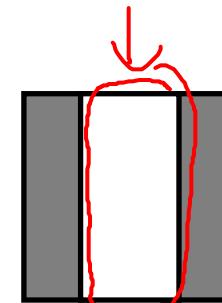
$$\begin{matrix} * & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} \\ & 3 \times 3 \end{matrix}$$

\*



$\uparrow \uparrow \uparrow$

$$= \begin{matrix} \downarrow & \begin{matrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{matrix} \\ & 4 \times 4 \end{matrix}$$



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

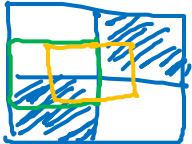
Vertical

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

$6 \times 6$



\*

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

→

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

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

↑

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

Sobel filter

convolution

\*

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

$\{ \}$   
 $3 \times 3$

these parameters can also be learnt vs. hard coding, allowing NN to learn various edges (low level features)

=  
 $45^\circ$   
 $70^\circ$   
 $73^\circ$

↑




deeplearning.ai

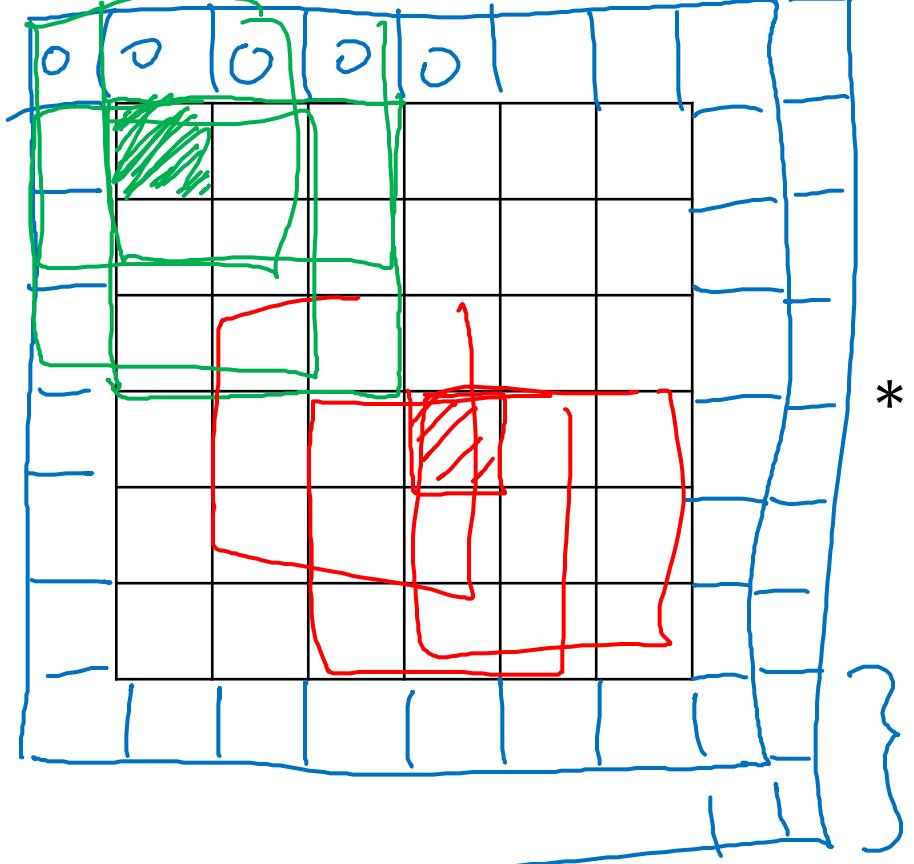
# Convolutional Neural Networks

---

## Padding

# Padding

- shrinky output
- throw away info from edge



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

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

$$\begin{matrix} n - (f-1) \\ n-f+1 \times n-f+1 \\ 6-3+1=4 \end{matrix}$$

$P = \text{padding} = 1$

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

$\frac{6 \times 6}{4 \times 4}$

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

# Valid and Same convolutions

→ n → padding

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

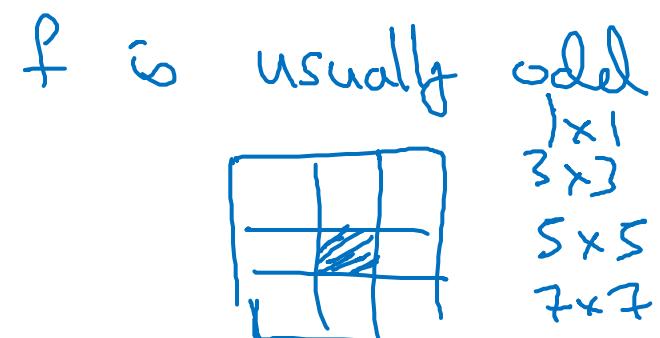
$$6 \times 6 \quad * \quad 3 \times 3 \quad \rightarrow \quad 4 \times 4$$

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

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

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

$$3 \times 3 \quad p = \frac{3-1}{2} = 1 \quad \left| \begin{array}{c} 5 \times 5 \\ f=5 \\ p=2 \end{array} \right.$$





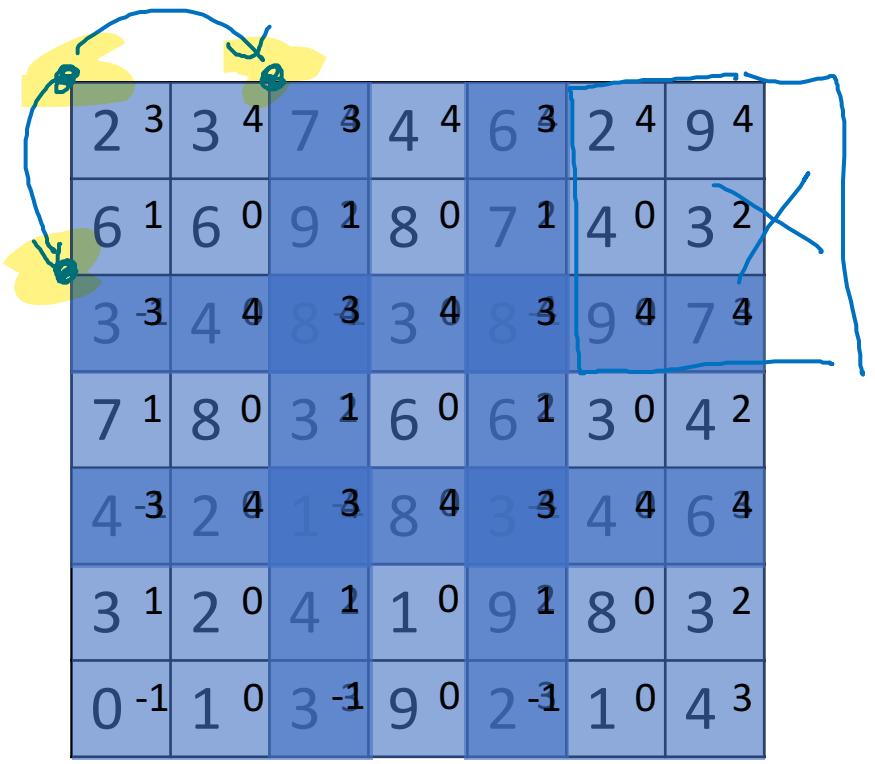
deeplearning.ai

# Convolutional Neural Networks

---

## Strided convolutions

# Strided convolution



7x7

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

$$\begin{array}{c}
 \begin{array}{|c|c|c|} \hline 3 & 4 & 4 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} & * & \begin{array}{|c|c|c|} \hline 91 & 100 & 83 \\ \hline 69 & 91 & 127 \\ \hline 44 & 72 & 74 \\ \hline \end{array} \\
 = & & \begin{array}{c} \text{Stride } = 2 \\ \hline 3 \times 3 \end{array} \\
 & & \lfloor \frac{z}{2} \rfloor = \text{floor}(z) \quad \downarrow
 \end{array}$$

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

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

# Summary of convolutions

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

padding  $p$       stride  $s$

Output size:

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

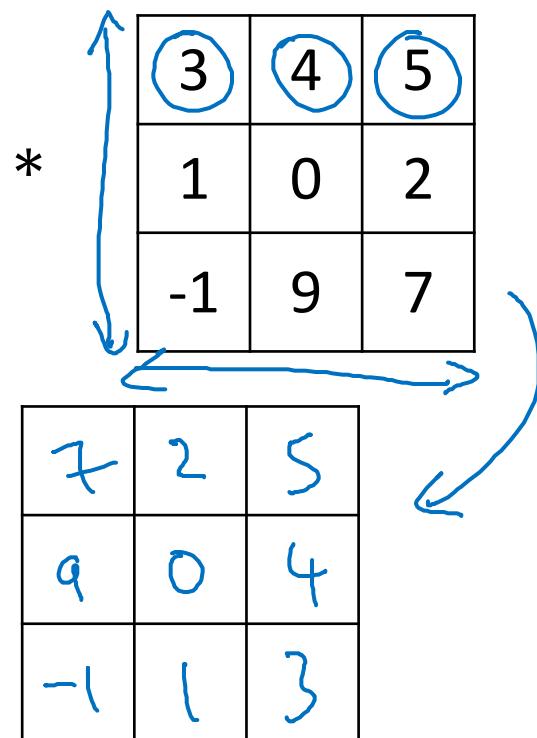

# Technical note on cross-correlation vs. convolution

## Convolution in math textbook:

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

flipped across horizontal & vertical should be:

7 9 -1  
2 0 1  
5 4 3



In ML - we skip the flipping step, but still call it convolution (and not cross-correlation)

#tradition

10	20	30
20	40	50
30	50	60

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



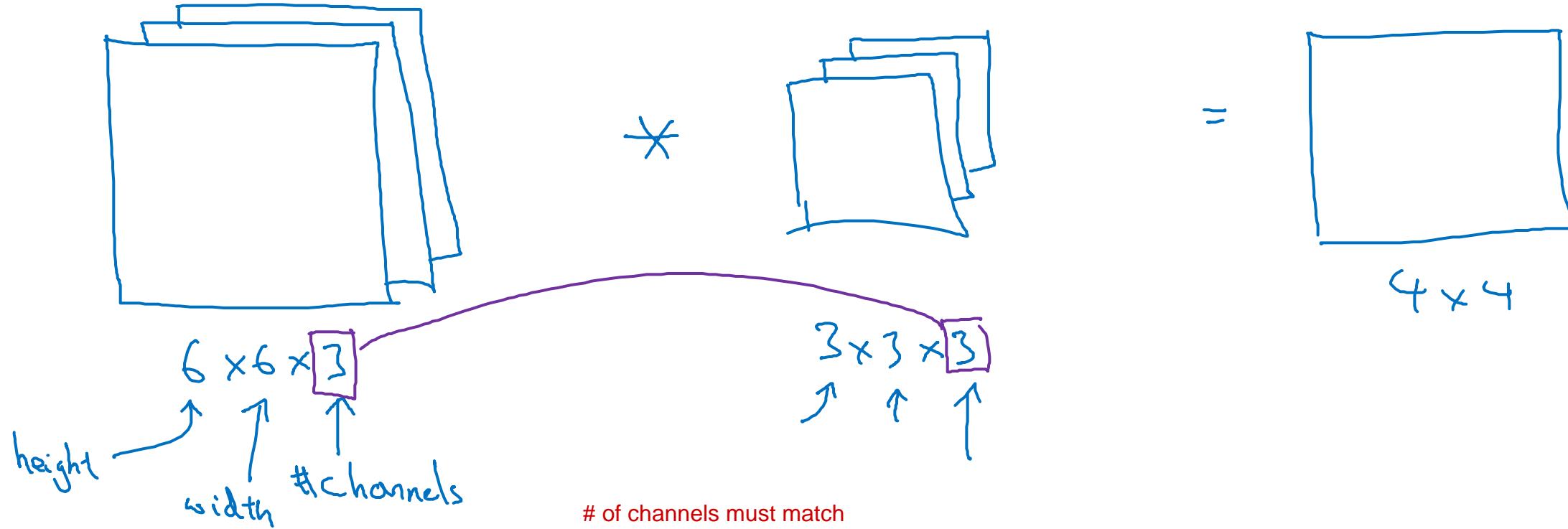
deeplearning.ai

# Convolutional Neural Networks

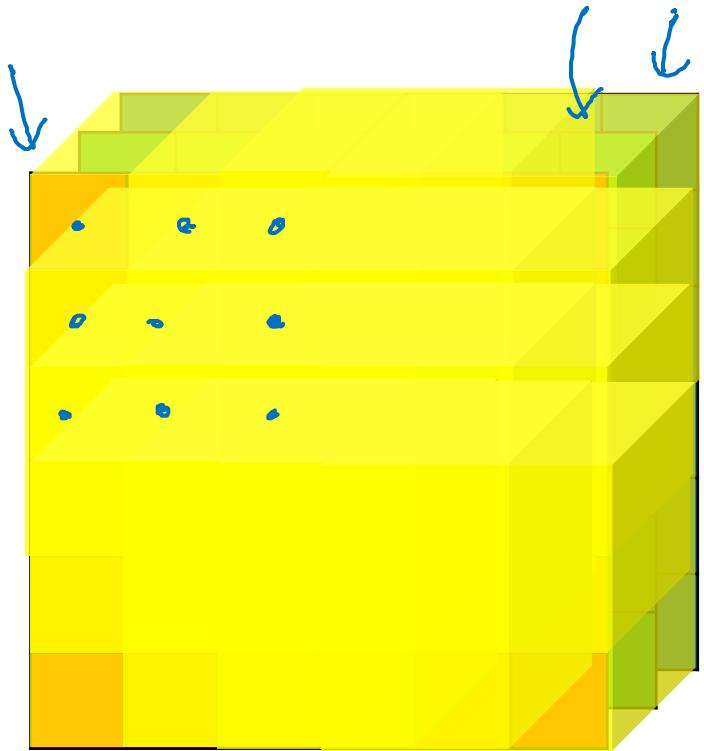
---

## Convolutions over volumes

# Convolutions on RGB images

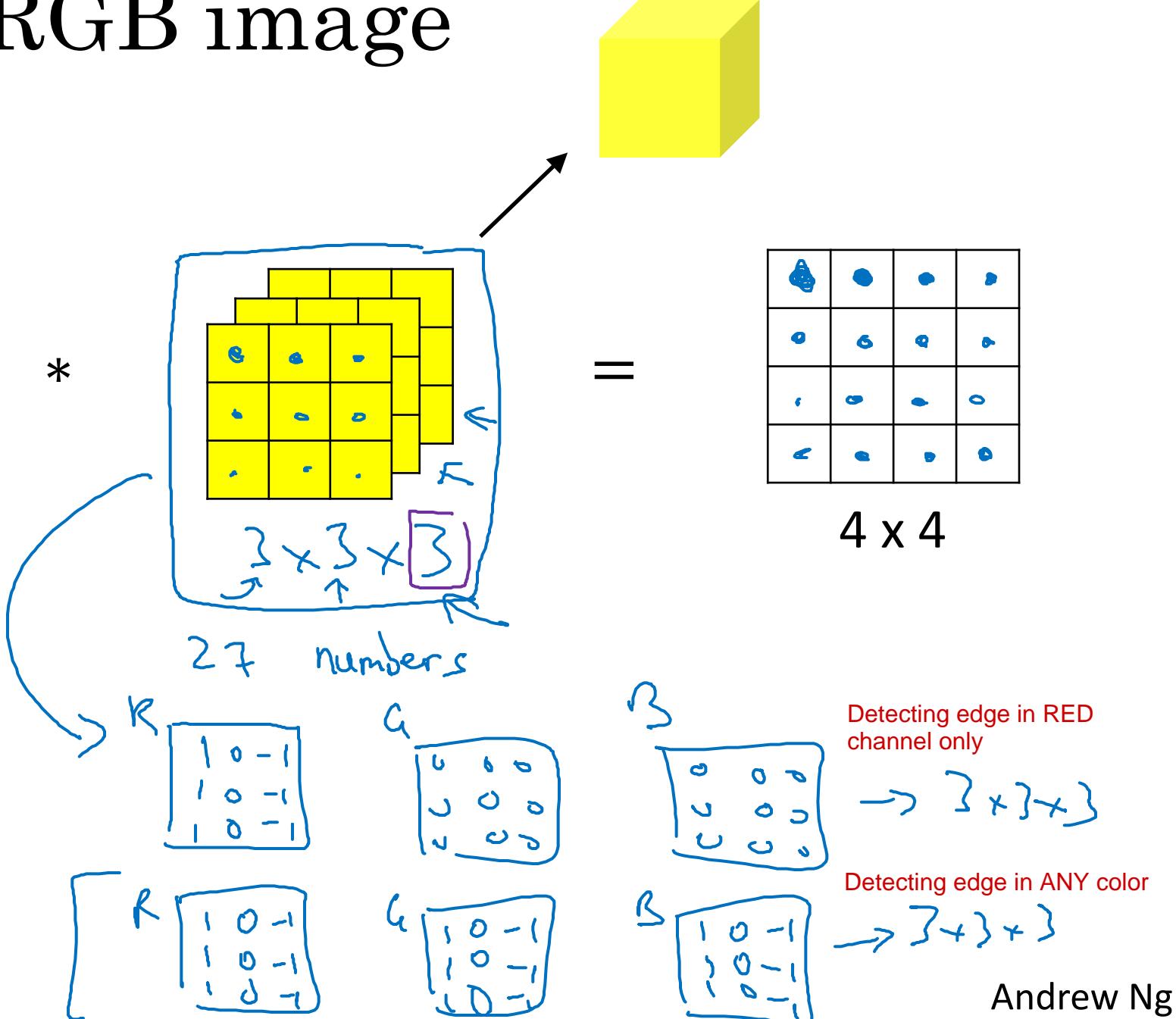


# Convolutions on RGB image



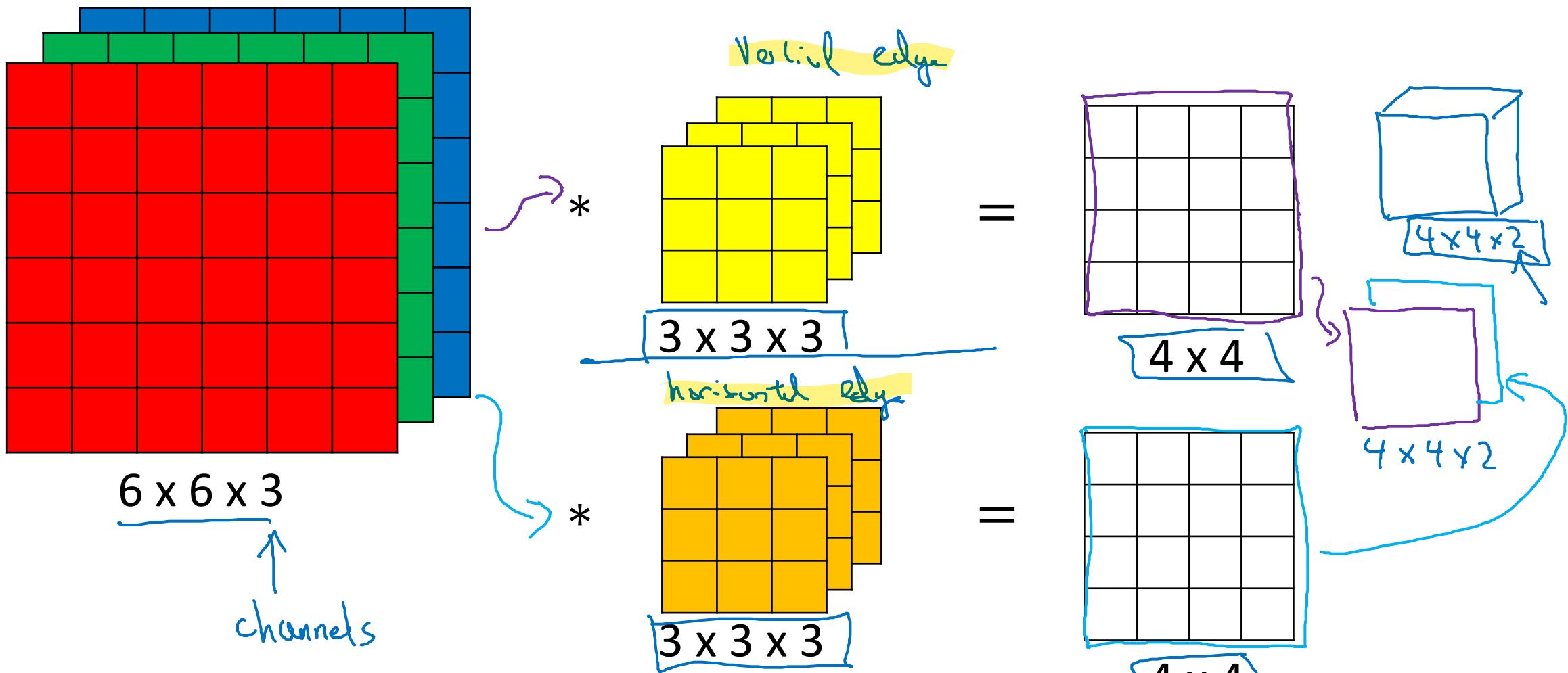
$$\text{Input} \times \text{Kernel} = \text{Output}$$

$6 \times 6 \times 3$



Andrew Ng

# Multiple filters



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

$6 \times 6 \times 3 \quad 3 \times 3 \times 3 \quad \frac{4-3+1}{4} \times \frac{4-3+1}{4} \times 2$



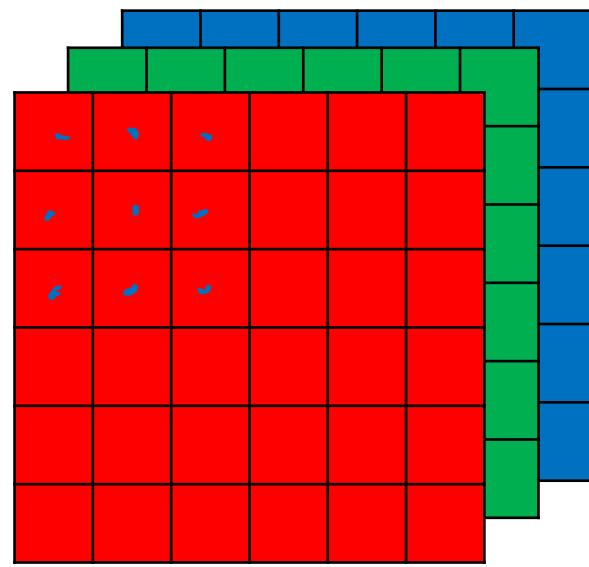
deeplearning.ai

# Convolutional Neural Networks

---

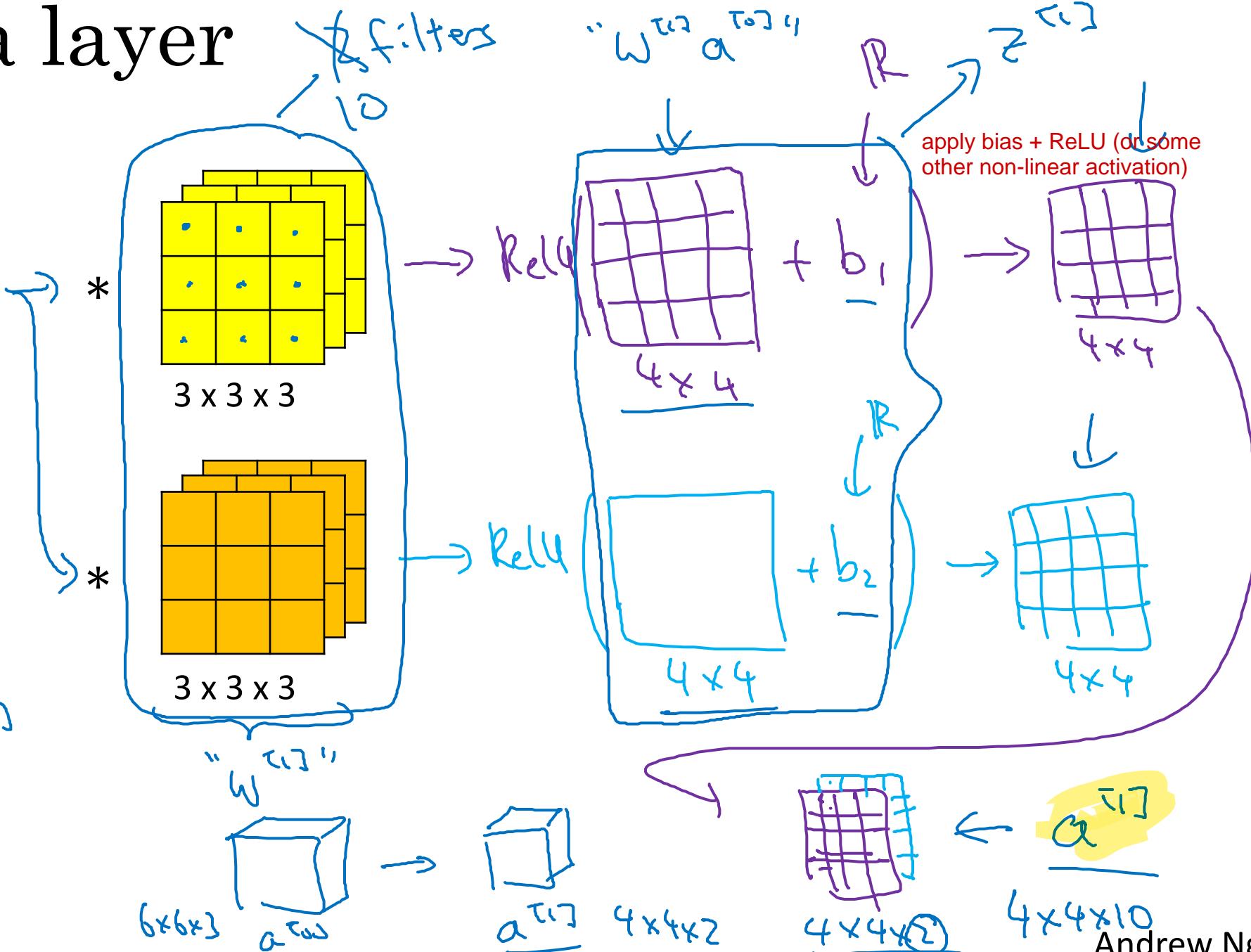
One layer of a  
convolutional  
network

# Example of a layer



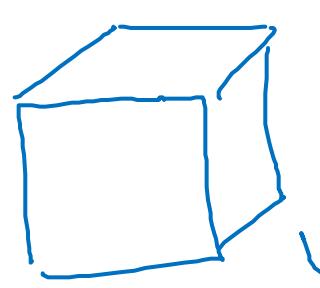
$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$$

$$a^{(1)} = g(z^{(1)})$$



# 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?

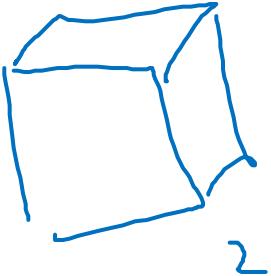


$3 \times 3 \times 3$

27 parameters.

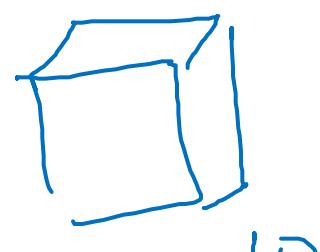
+ bias

→ 28 parameters.



2

...  
...



10

280 parameters.

# Summary of notation

If layer l is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

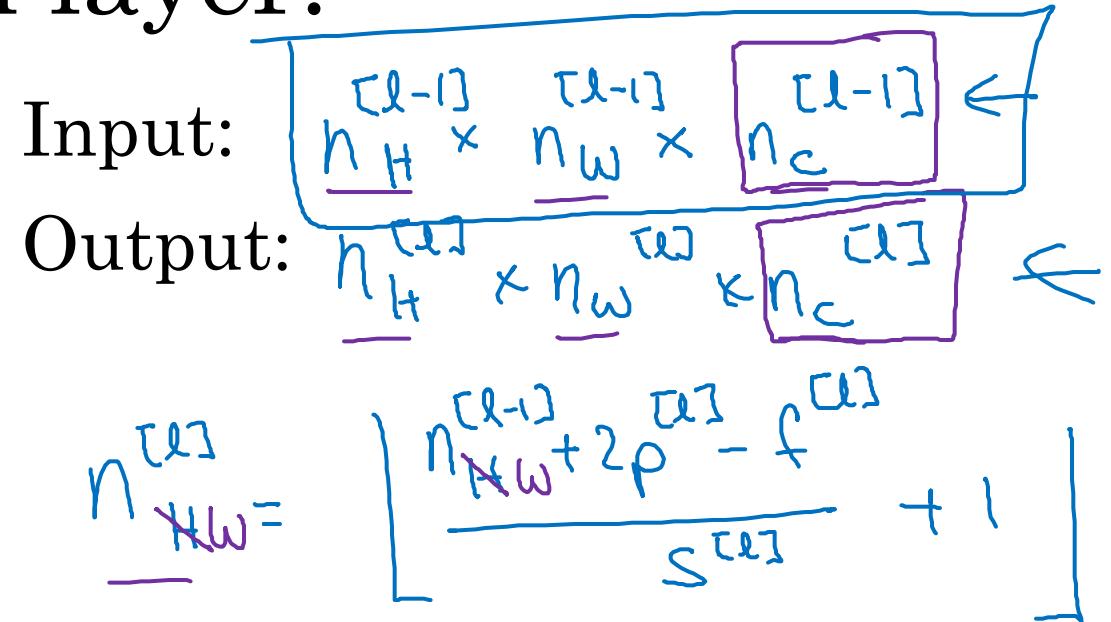
$n_c^{[l]}$  = number of filters

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

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

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

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



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

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



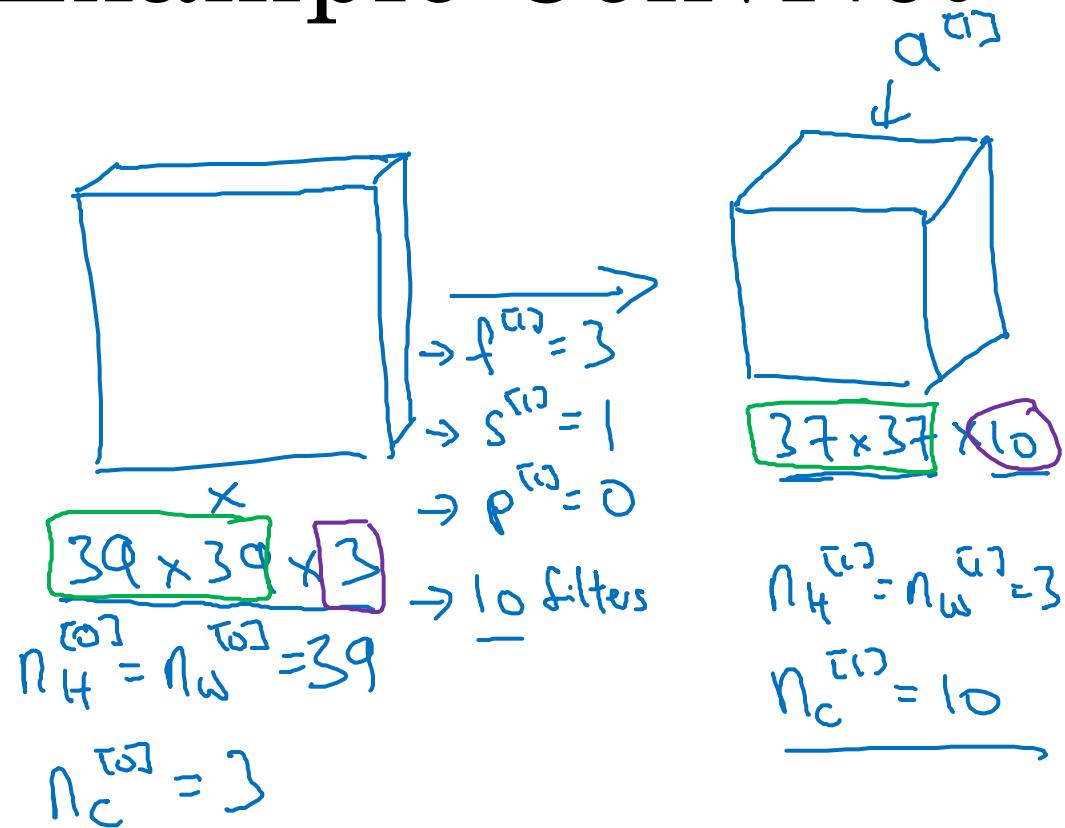
deeplearning.ai

# Convolutional Neural Networks

---

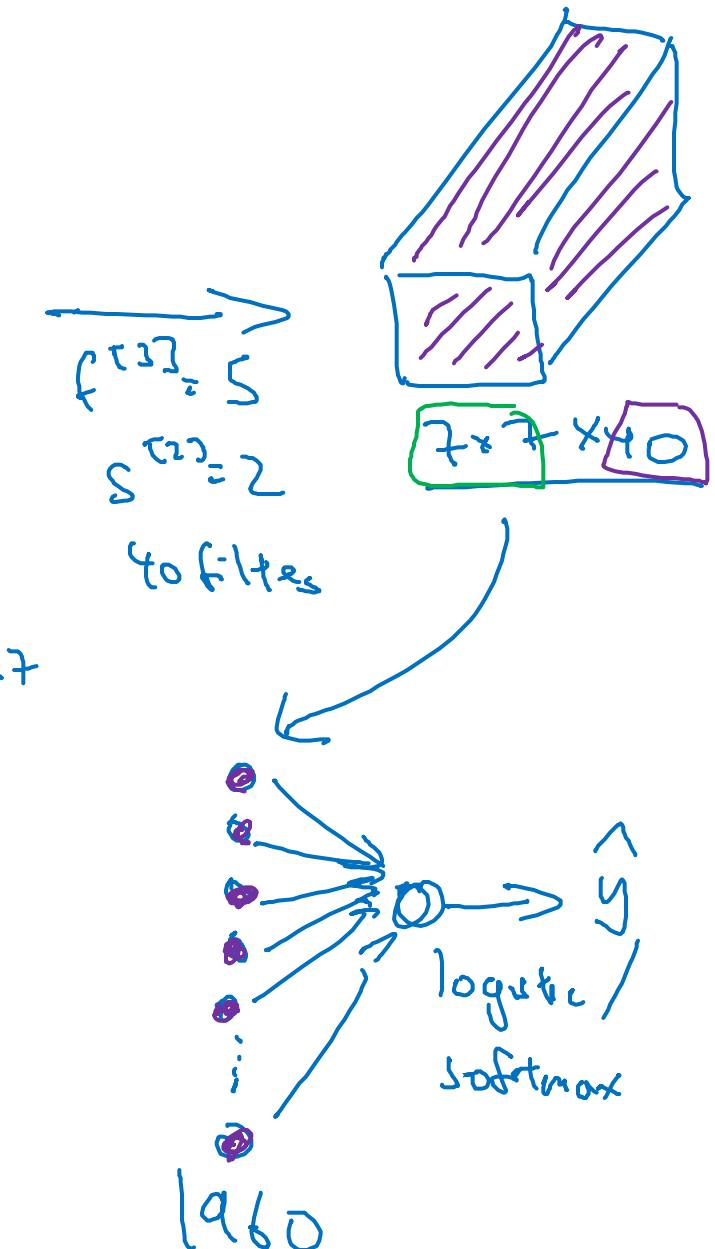
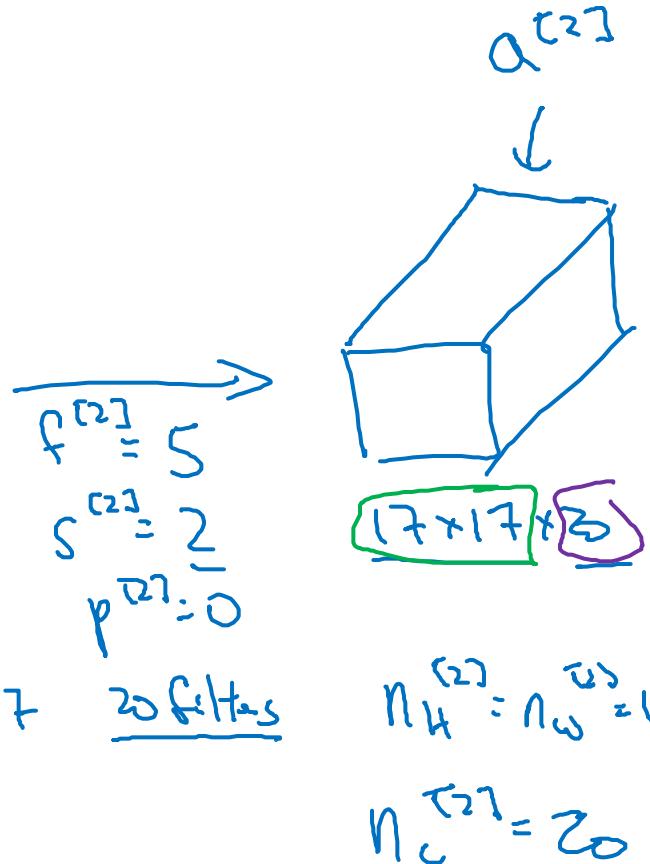
## A simple convolution network example

# Example ConvNet



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

$$\frac{39+0-3}{1} + 1 = 37$$



# Types of layer in a convolutional network:

- Convolution (Conv) ←
- Pooling (pool) ←
- Fully connected (Fc) ←



deeplearning.ai

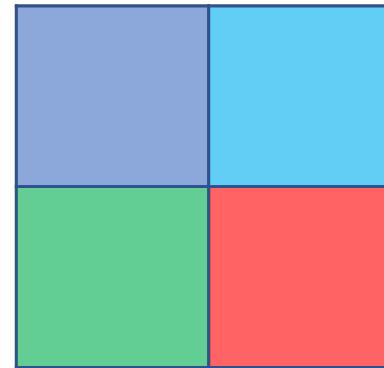
# Convolutional Neural Networks

---

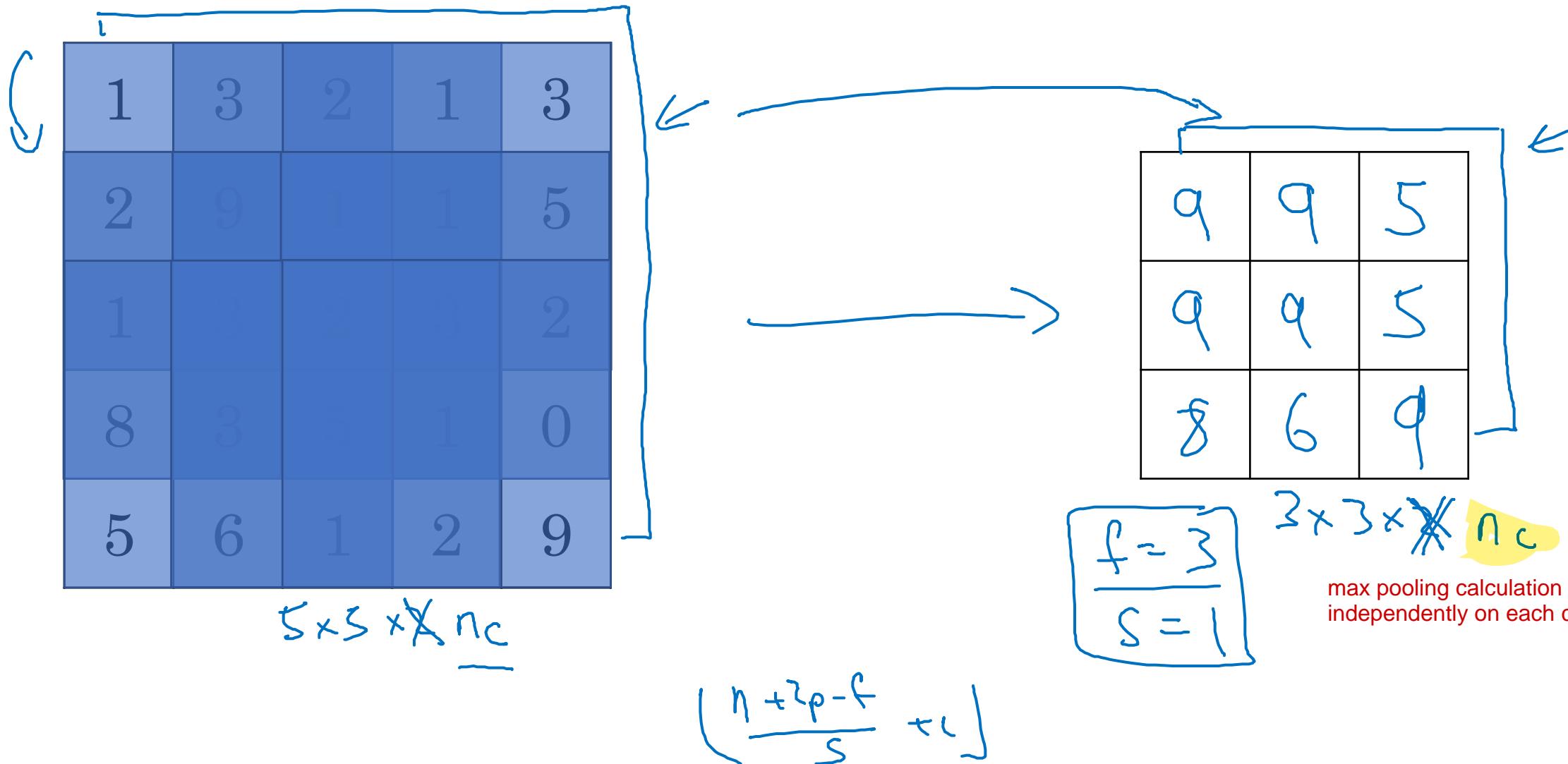
## Pooling layers

# Pooling layer: Max pooling

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



# Pooling layer: Max pooling



# Pooling layer: Average pooling

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



3.75	1.25
4	2

$$f=2$$

$$s=2$$

$$\underbrace{7 \times 7 \times 1000}_{\rightarrow} \rightarrow 1 \times 1 \times 1000$$

# Summary of pooling

Hyperparameters:

$f$  : filter size

$$f=2, s=2$$

$s$  : stride

$$f=3, s=2$$

Max or average pooling

$\rightarrow p$ : padding.

No parameters to learn!

$$n_H \times n_w \times n_c$$



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

$$\times n_c$$



deeplearning.ai

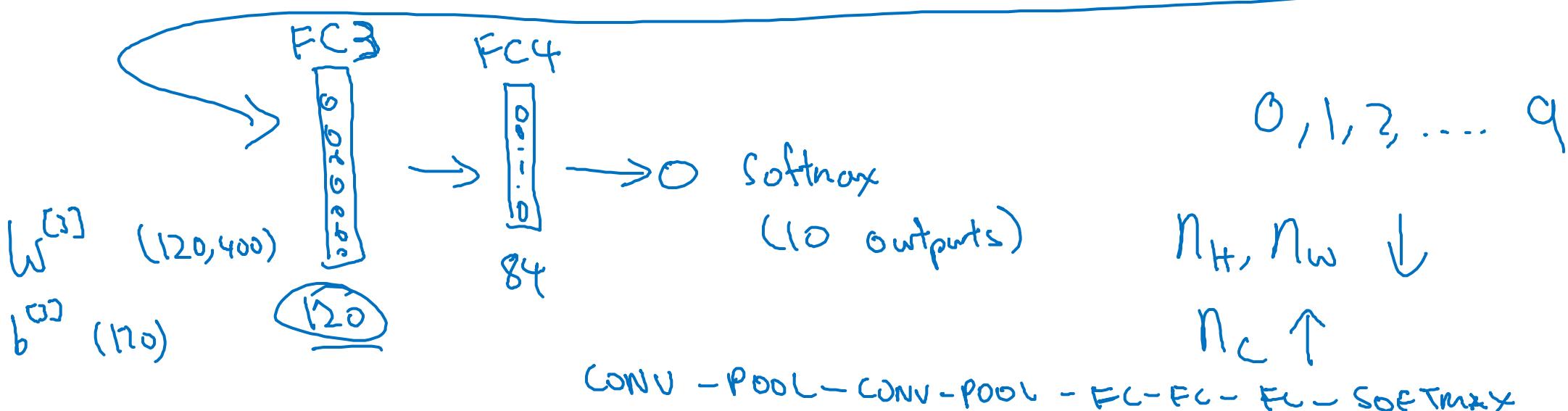
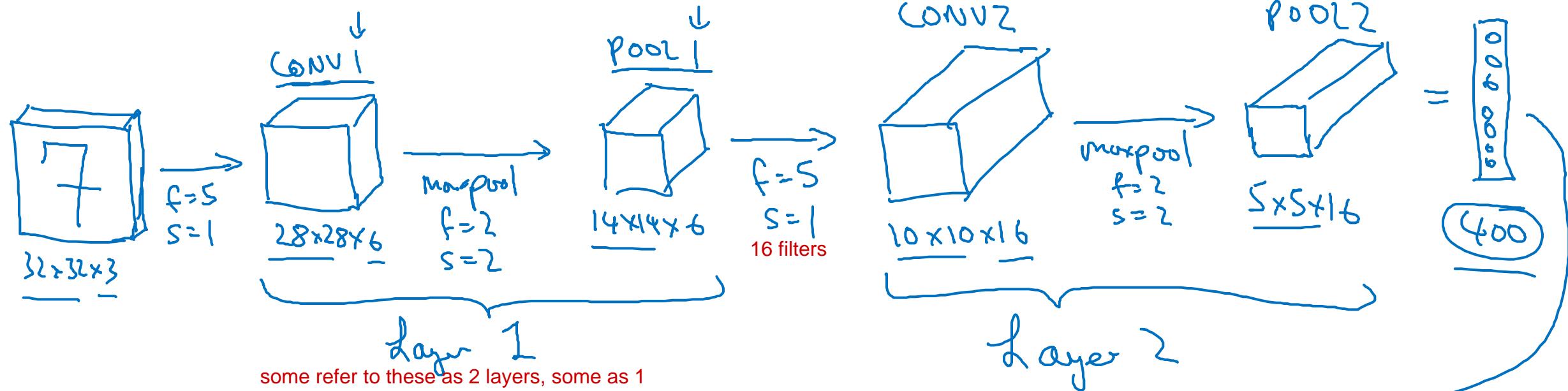
# Convolutional Neural Networks

---

## Convolutional neural network example

# Neural network example

(LeNet-5)



# Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072 $a^{[3]}$	0
CONV1 (f=5, s=1)	(28,28,8)	6,272	$(5*5*3+1)*8$ 608 ←
POOL1	(14,14,8)	1,568	0 ←
CONV2 (f=5, s=1)	(10,10,16)	1,600	$(5*5*8+1)*16$ 3216 ←
POOL2	(5,5,16)	400	0 ←
FC3	(120,1)	120	48120 {
FC4	(84,1)	84	10164 }
Softmax	(10,1)	10	850



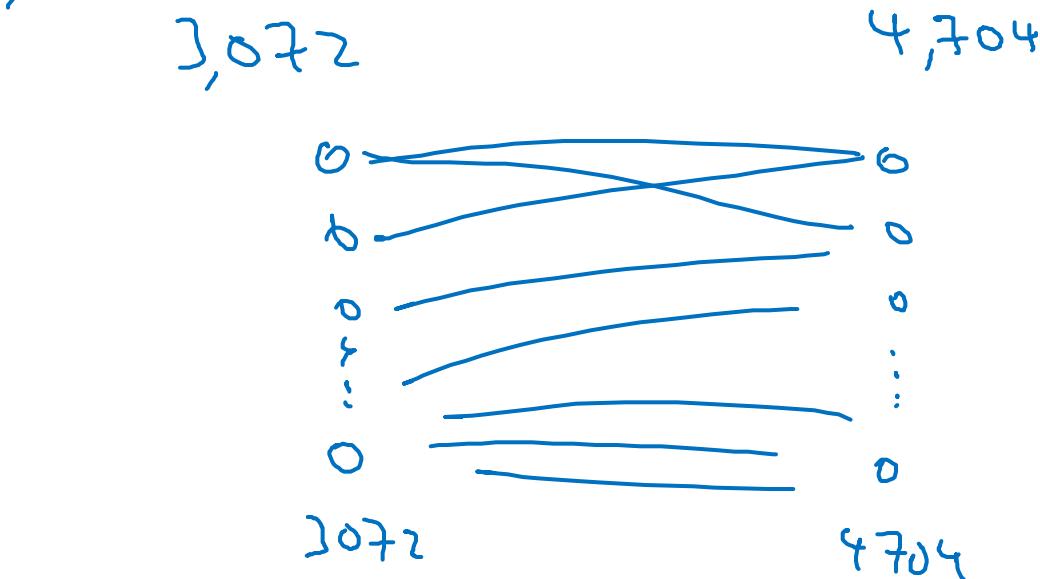
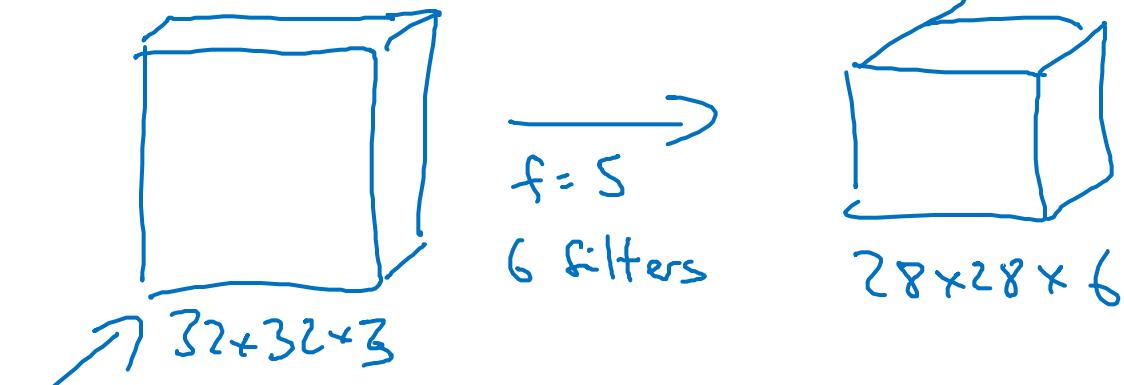
deeplearning.ai

# Convolutional Neural Networks

---

## Why convolutions?

# Why convolutions



$$5 \times 5 - 25$$

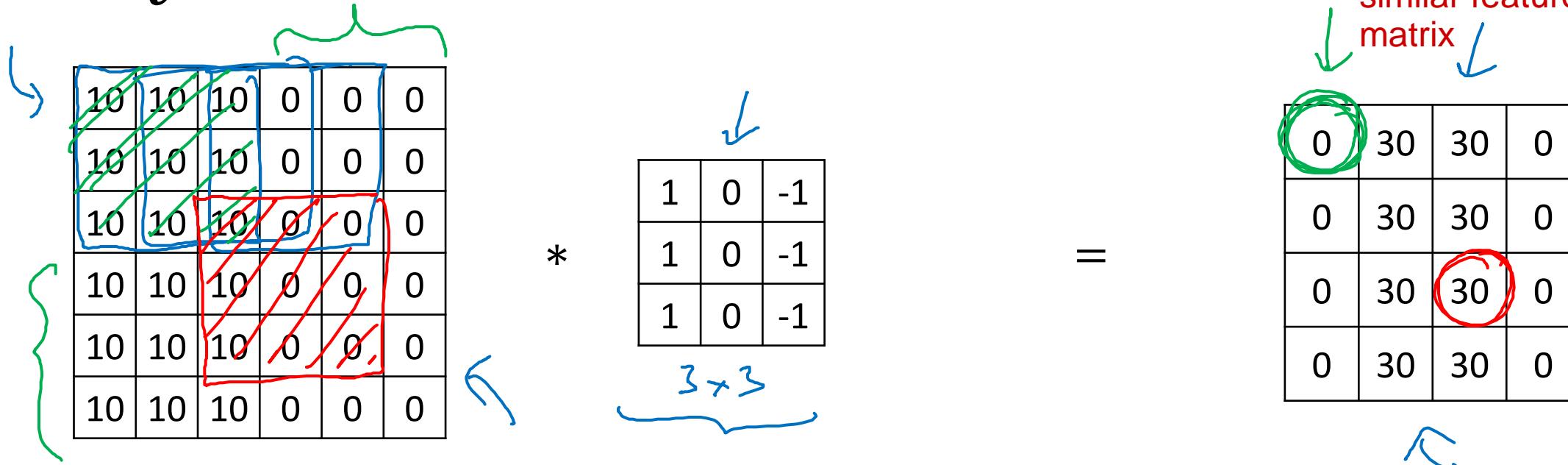
$$26$$

$$\cancel{6 \times 25} = \cancel{156} \quad \text{Parameters}$$
$$6 * (5 \times 5 \times 3 + 1) = 456$$

+ Parameter sharing  
+ Sparsity of connections

$$3,072 \times 4,704 \approx 14M$$

# Why convolutions



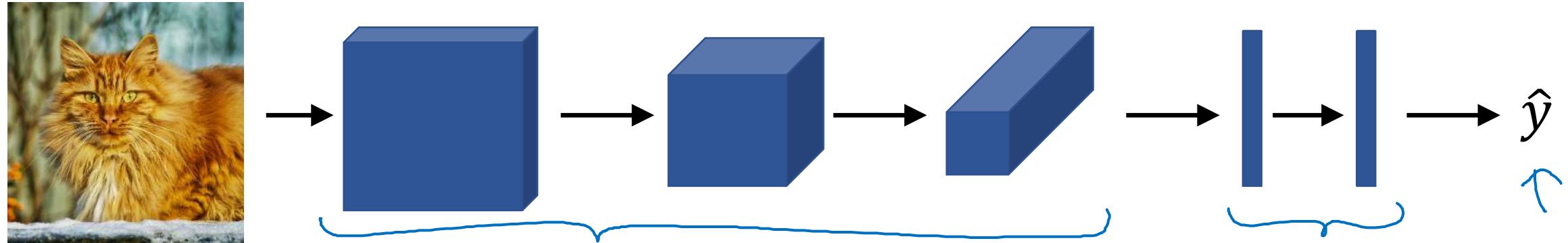
**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

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

Image shifted a few pixels over should still result in similar features and output matrix

# 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  $J$

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# Case Studies

---

Why look at  
case studies?

# Outline

Classic networks:

- LeNet-5 ←
- AlexNet ←
- VGG ←

ResNet (152)

Inception



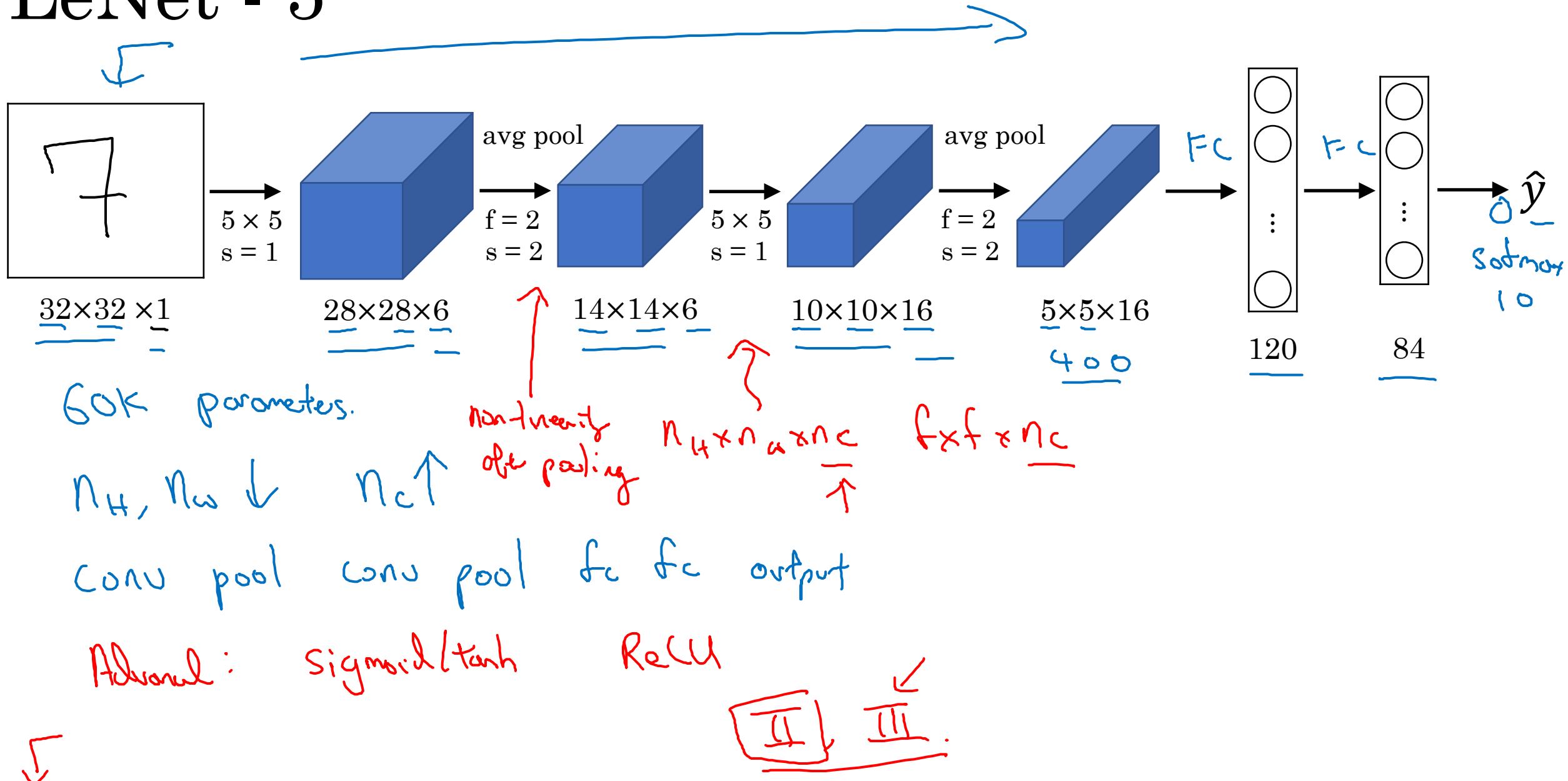
deeplearning.ai

## Case Studies

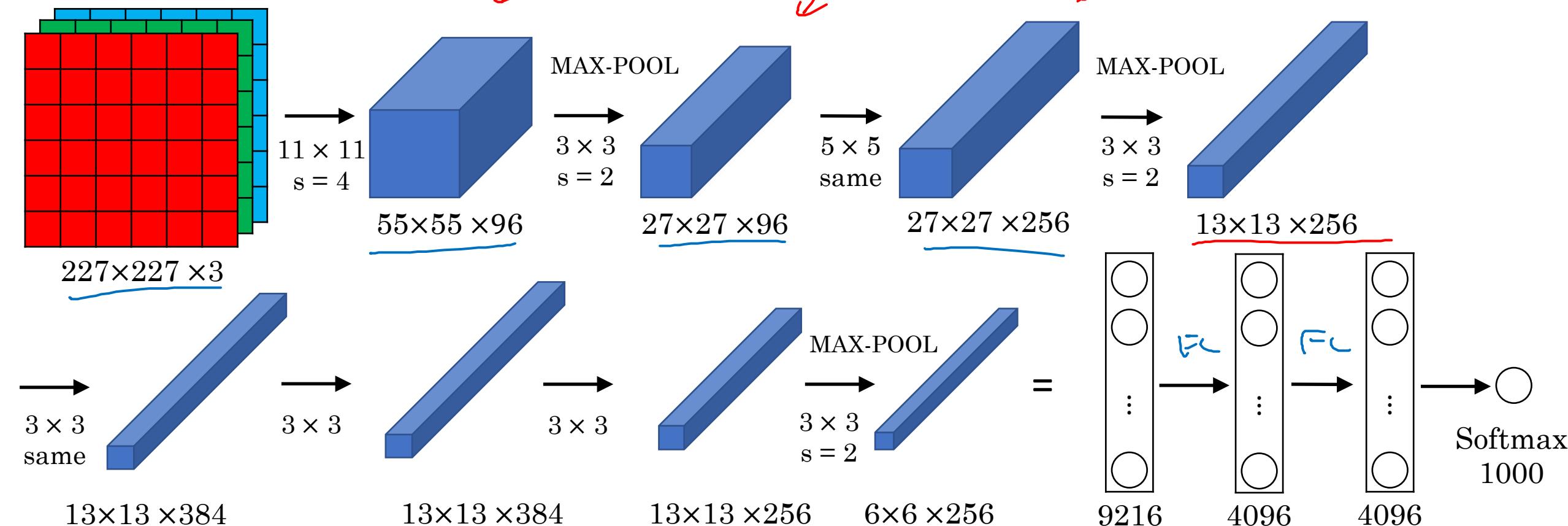
---

## Classic networks

# LeNet - 5

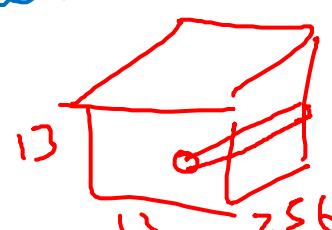


# AlexNet



- Similar to LeNet, but much bigger.
- ReLU

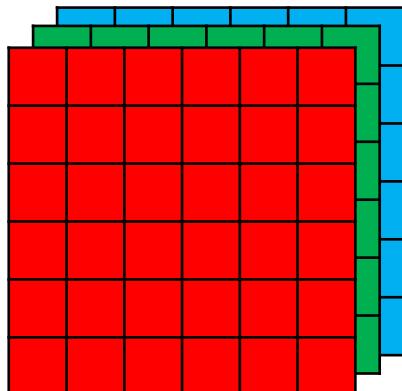
- Multiple GPUs.  
- Local Response Normalization (LRN)



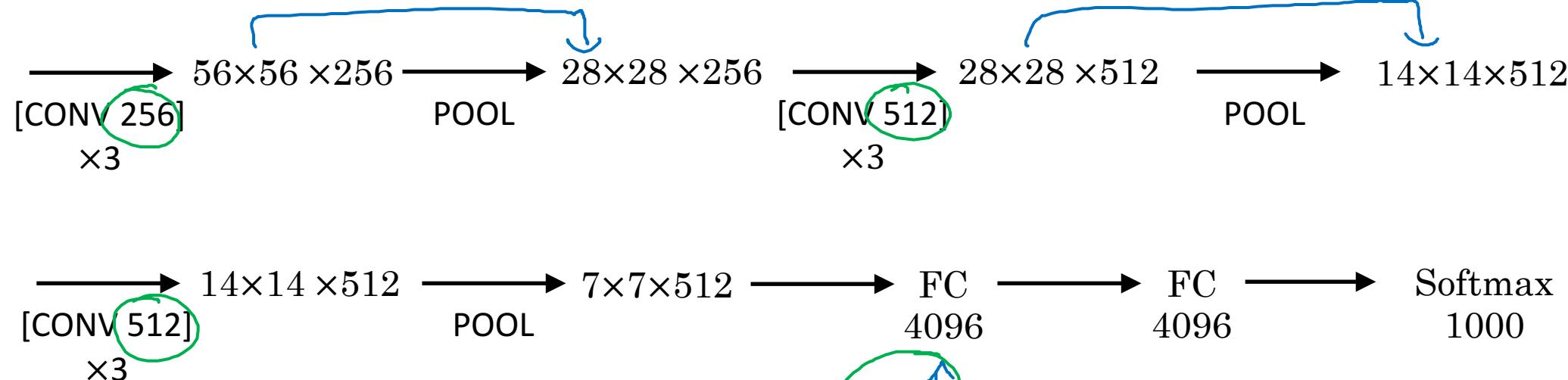
$\sim 60M$  Parameters

# VGG - 16

Use same filter config  
 $\text{CONV} = 3 \times 3 \text{ filter}, s = 1, \text{ same}$



$224 \times 224 \times 3$



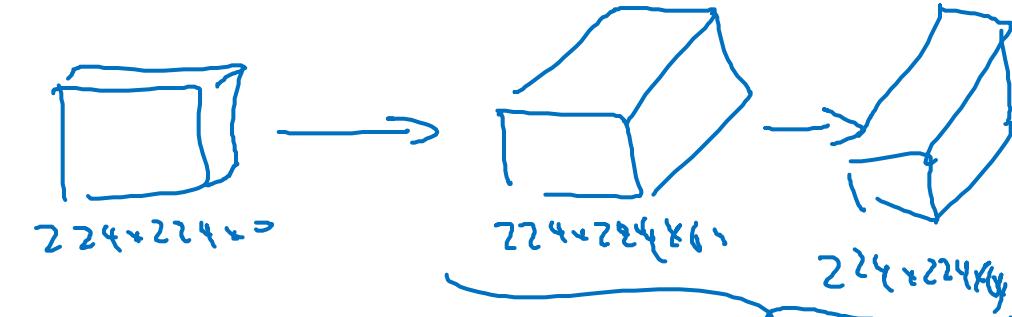
$n_H, n_W \downarrow$

$n_C \uparrow$

$\sim 38M$

# VGG-19

MAX-POOL =  $2 \times 2$ ,  $s = 2$



$56 \times 56 \times 128$  → POOL →  $28 \times 28 \times 512$

$28 \times 28 \times 512$  → POOL →  $14 \times 14 \times 512$

$14 \times 14 \times 512$  → POOL →  $7 \times 7 \times 512$  → FC 4096 → FC 4096 → Softmax 1000



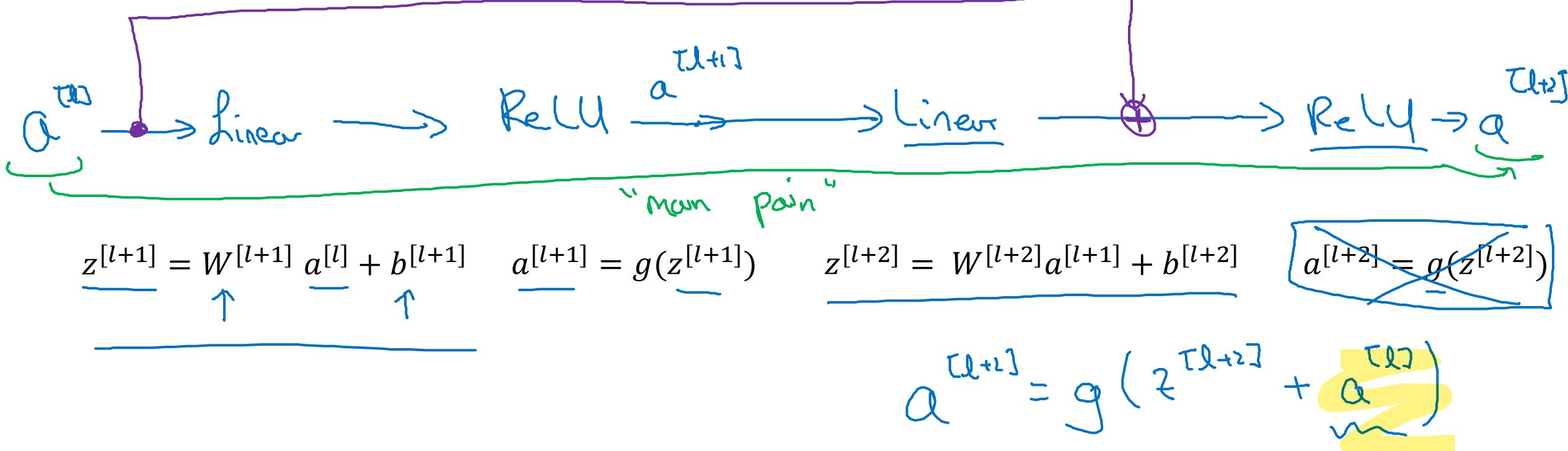
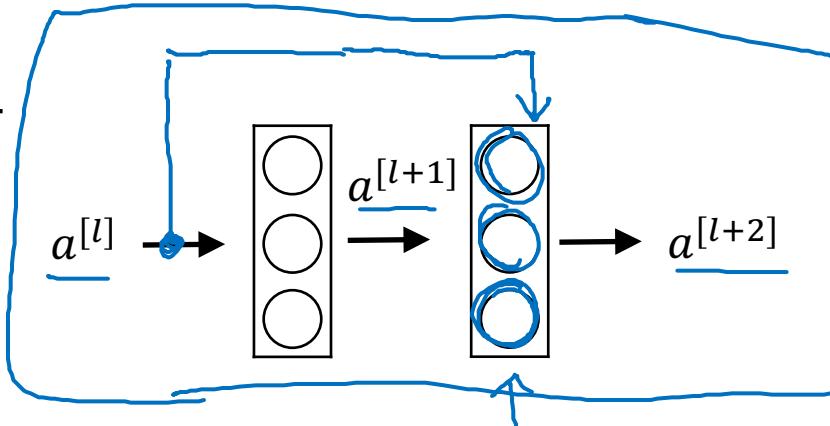
deeplearning.ai

## Case Studies

---

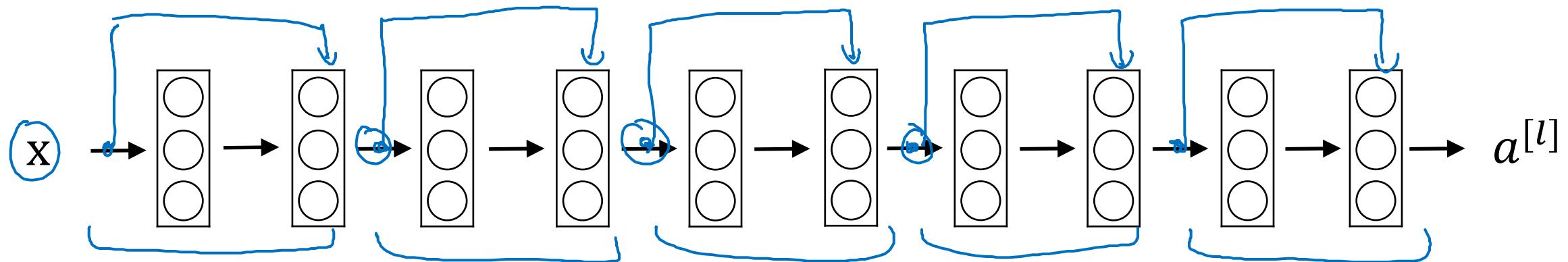
# Residual Networks (ResNets)

# Residual block



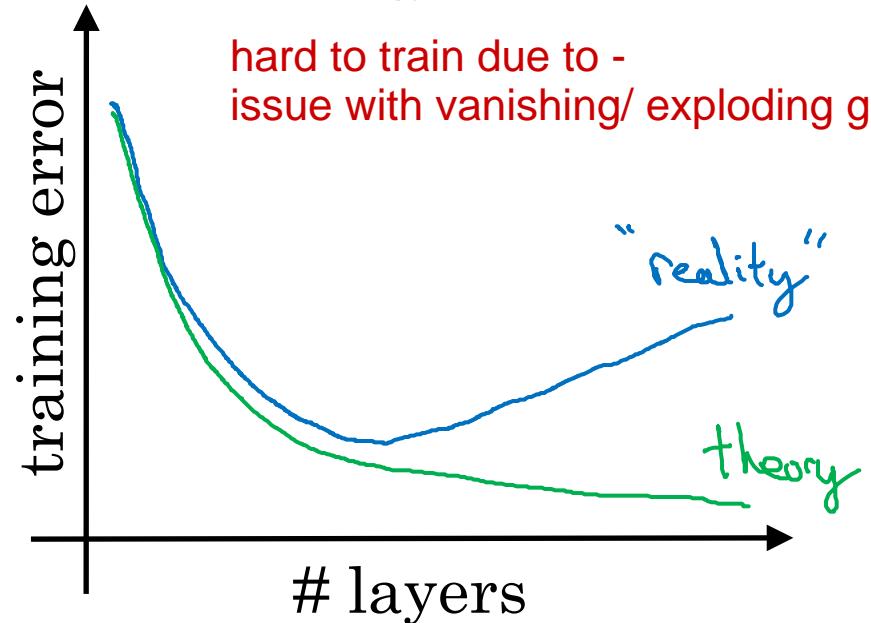
'a' added to 'z'  
dimensions are kept same

# Residual Network

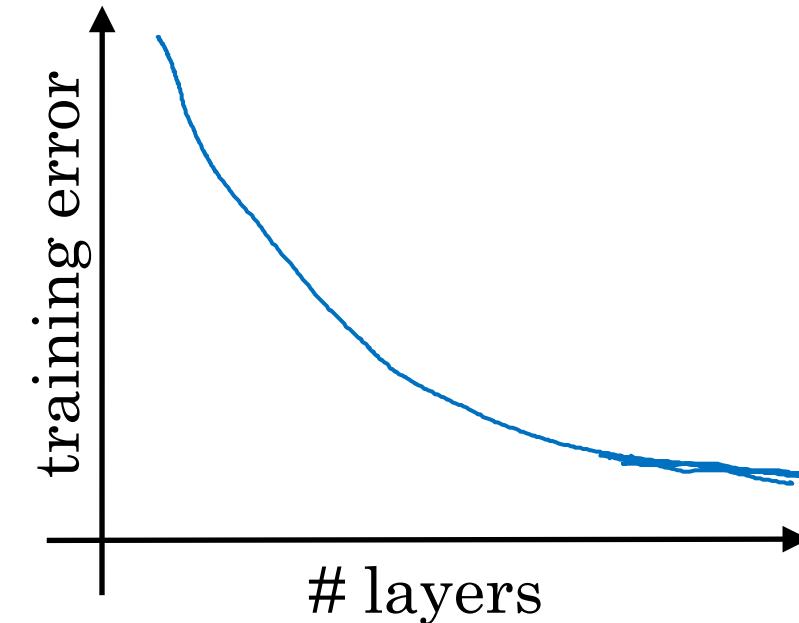


Plain

hard to train due to -  
issue with vanishing/ exploding gradients



ResNet





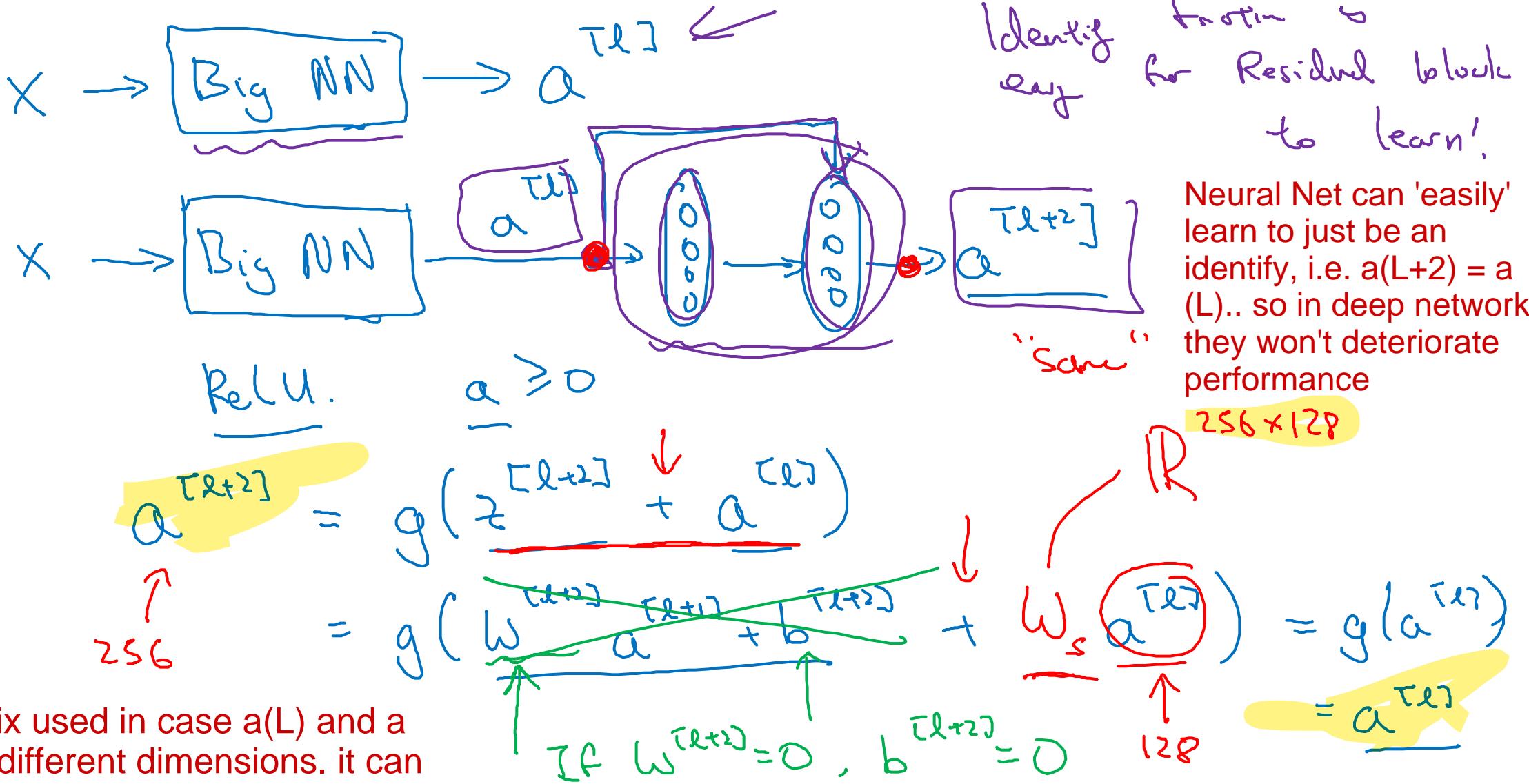
deeplearning.ai

## Case Studies

---

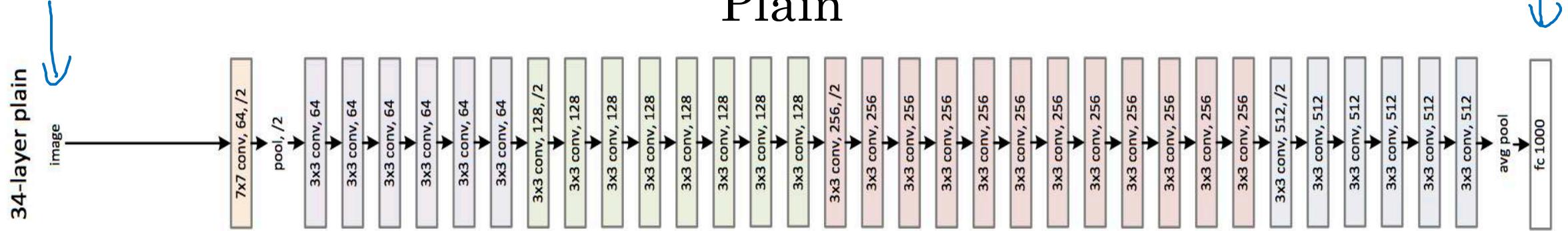
# Why ResNets work

# Why do residual networks work?

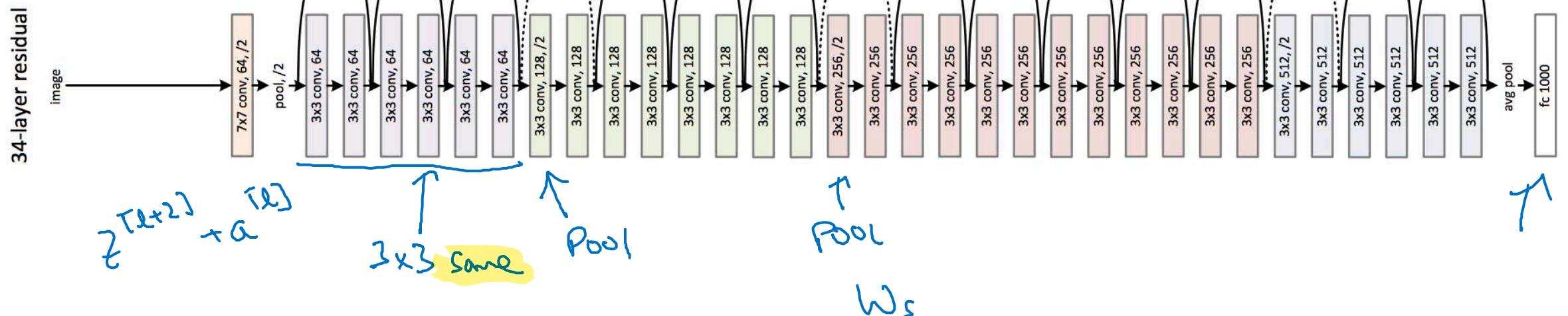


# ResNet

Plain



ResNet





deeplearning.ai

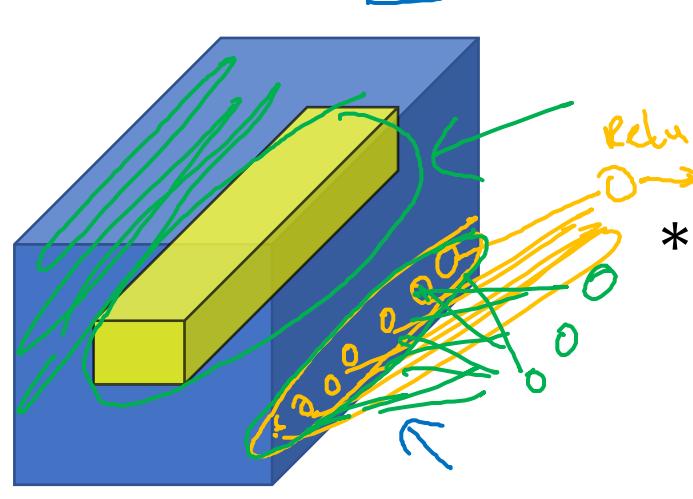
## Case Studies

---

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

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

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



\*

2



=

32  $\rightarrow$  # filters.

$n_c^{[l+1]}$

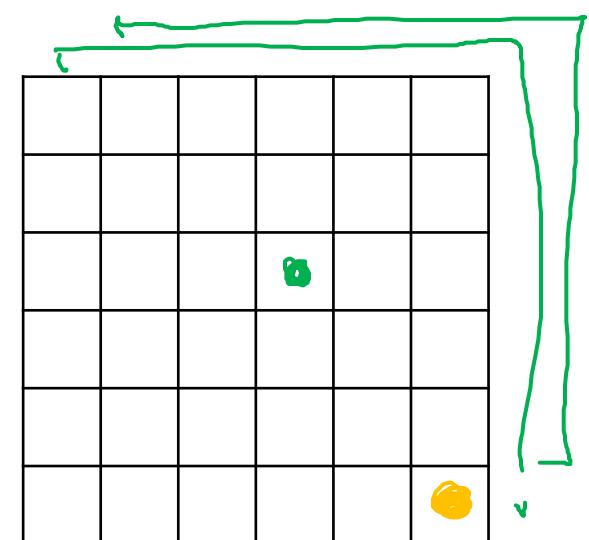
=

ReLU

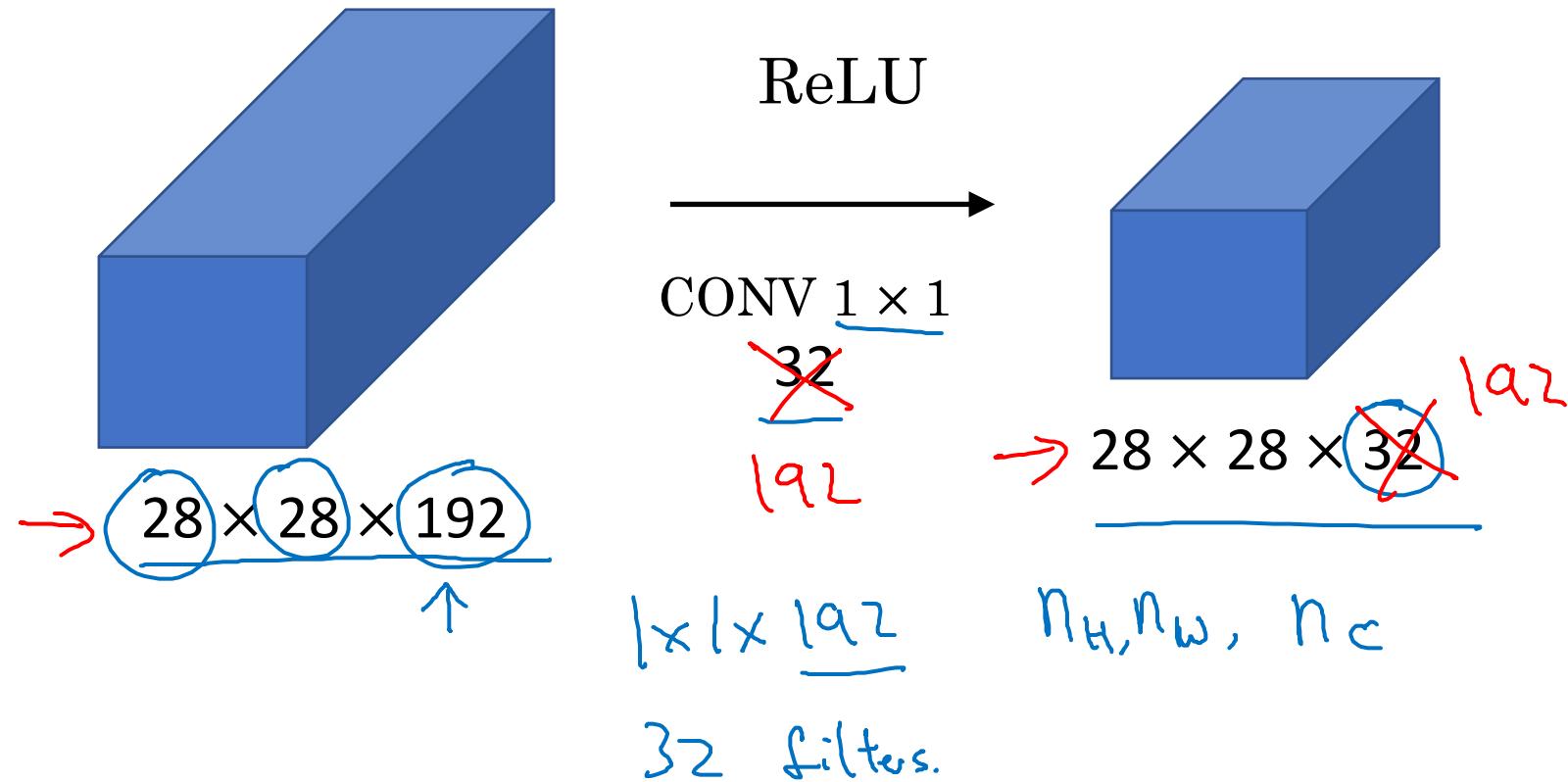
Network  $\hookrightarrow$  Network

1  $\times$  1  $\times$  32

2	4	6	...



# Using $1 \times 1$ convolutions





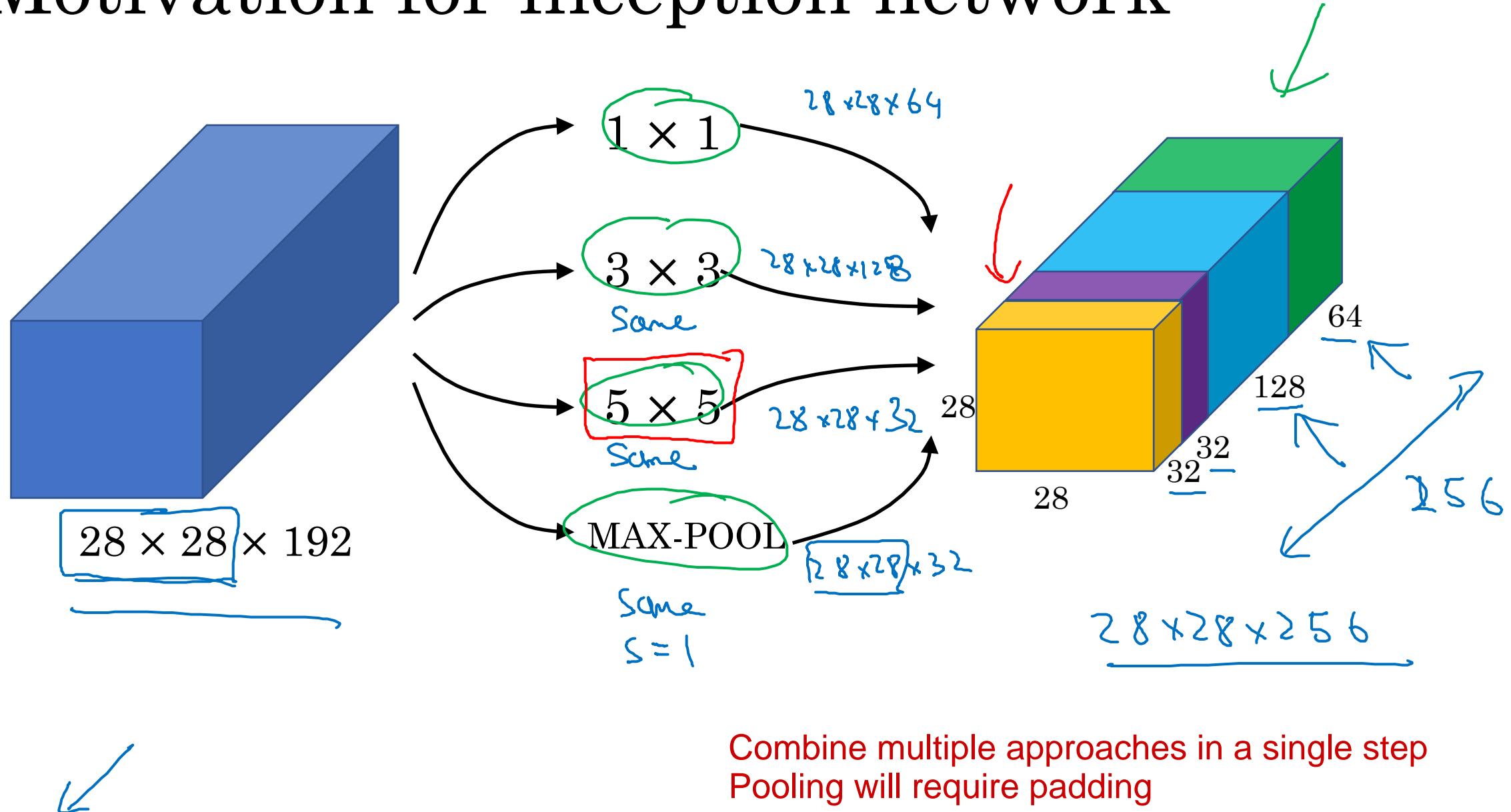
deeplearning.ai

## Case Studies

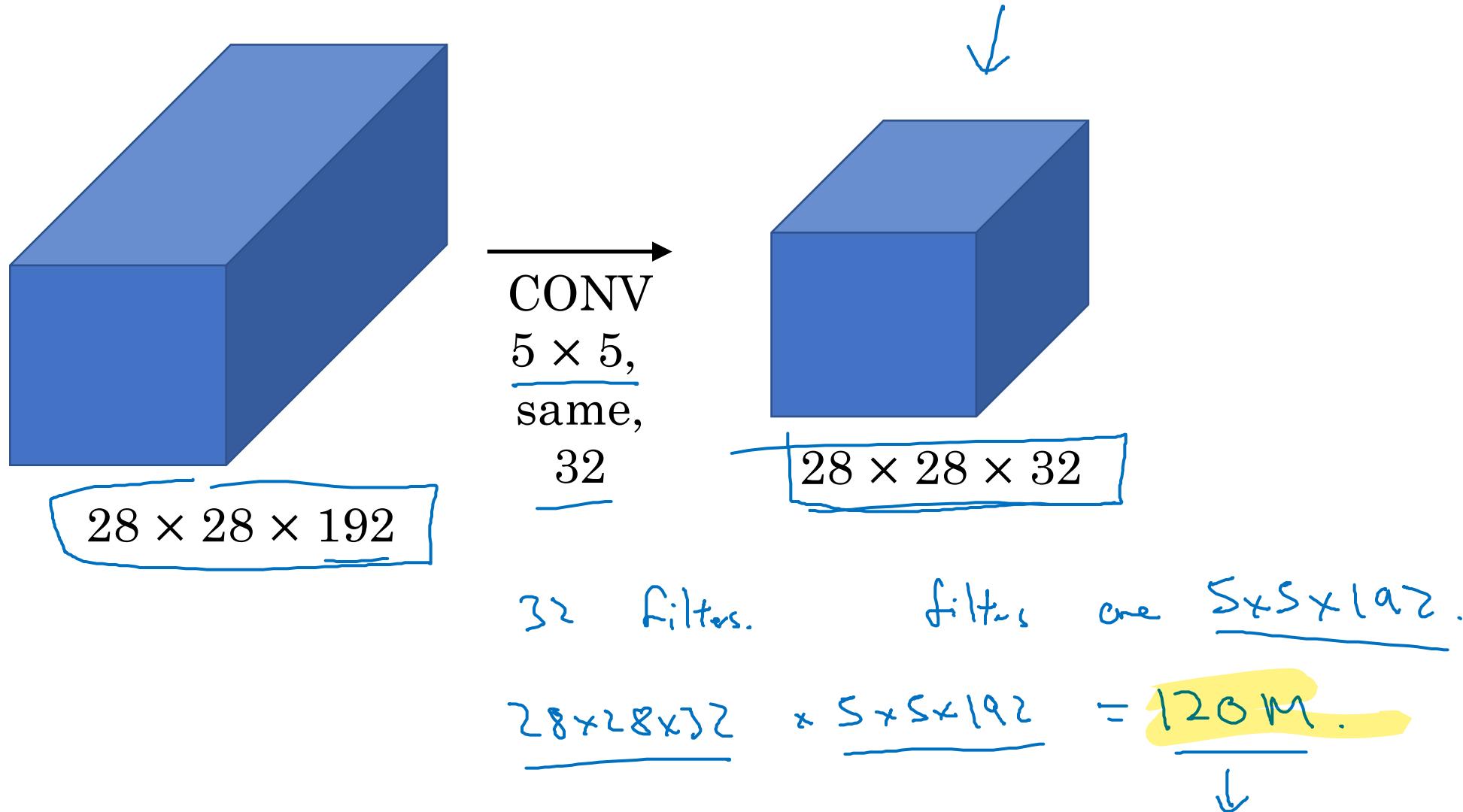
---

Inception network  
motivation

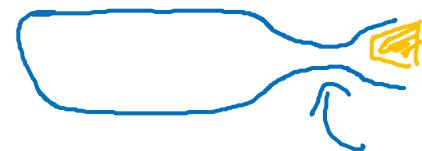
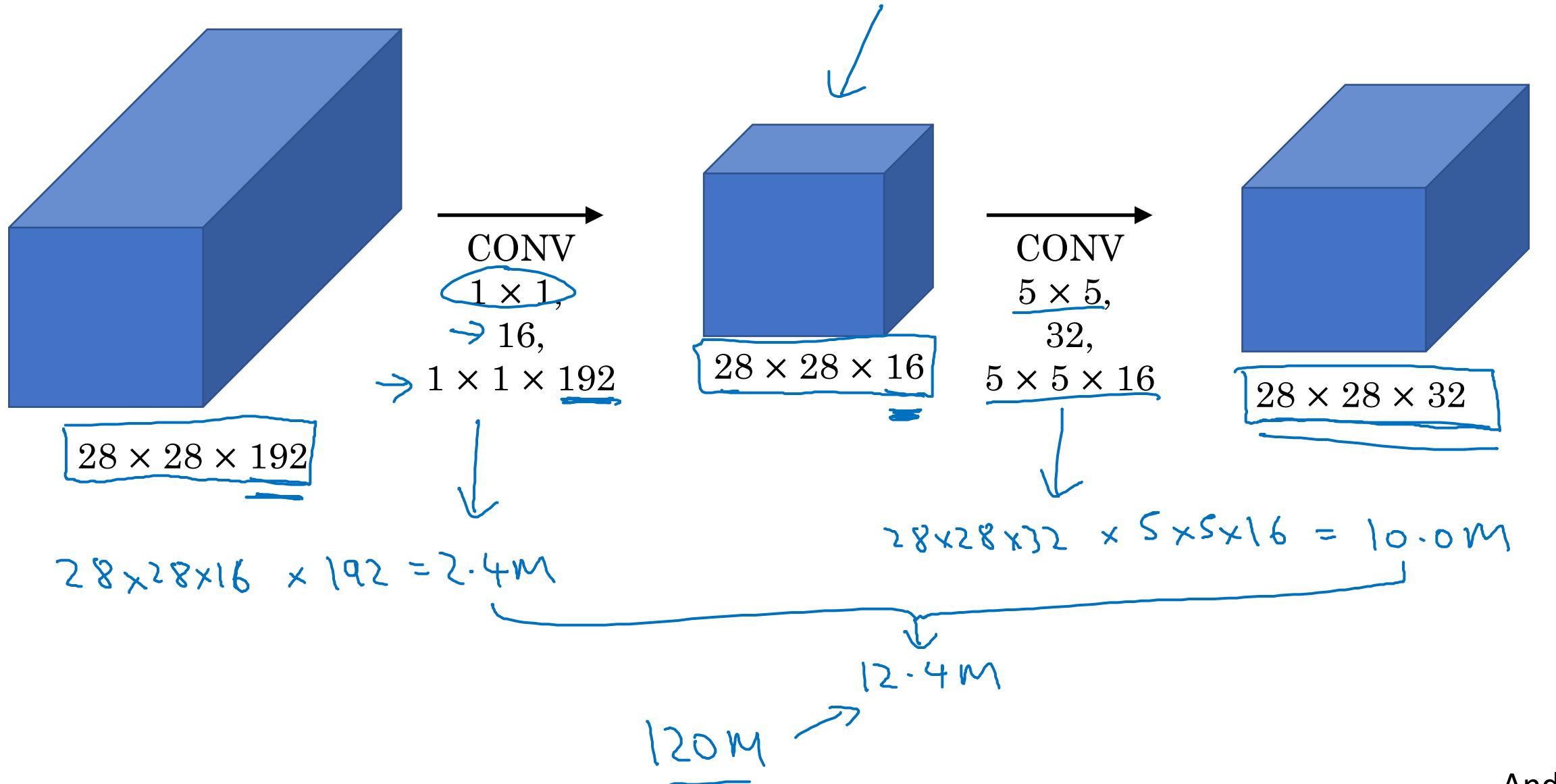
# Motivation for inception network



# The problem of computational cost



# Using $1 \times 1$ convolution





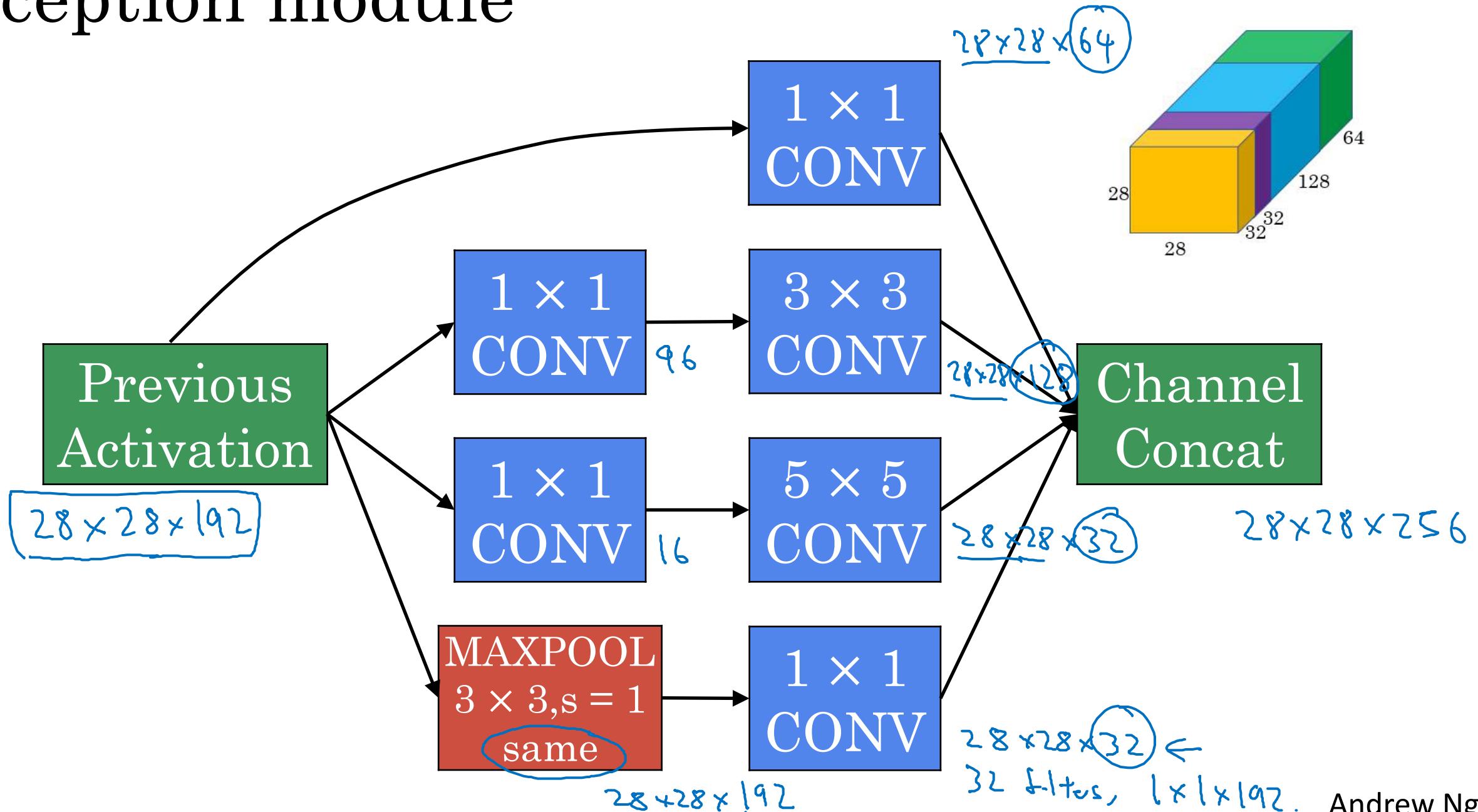
deeplearning.ai

## Case Studies

---

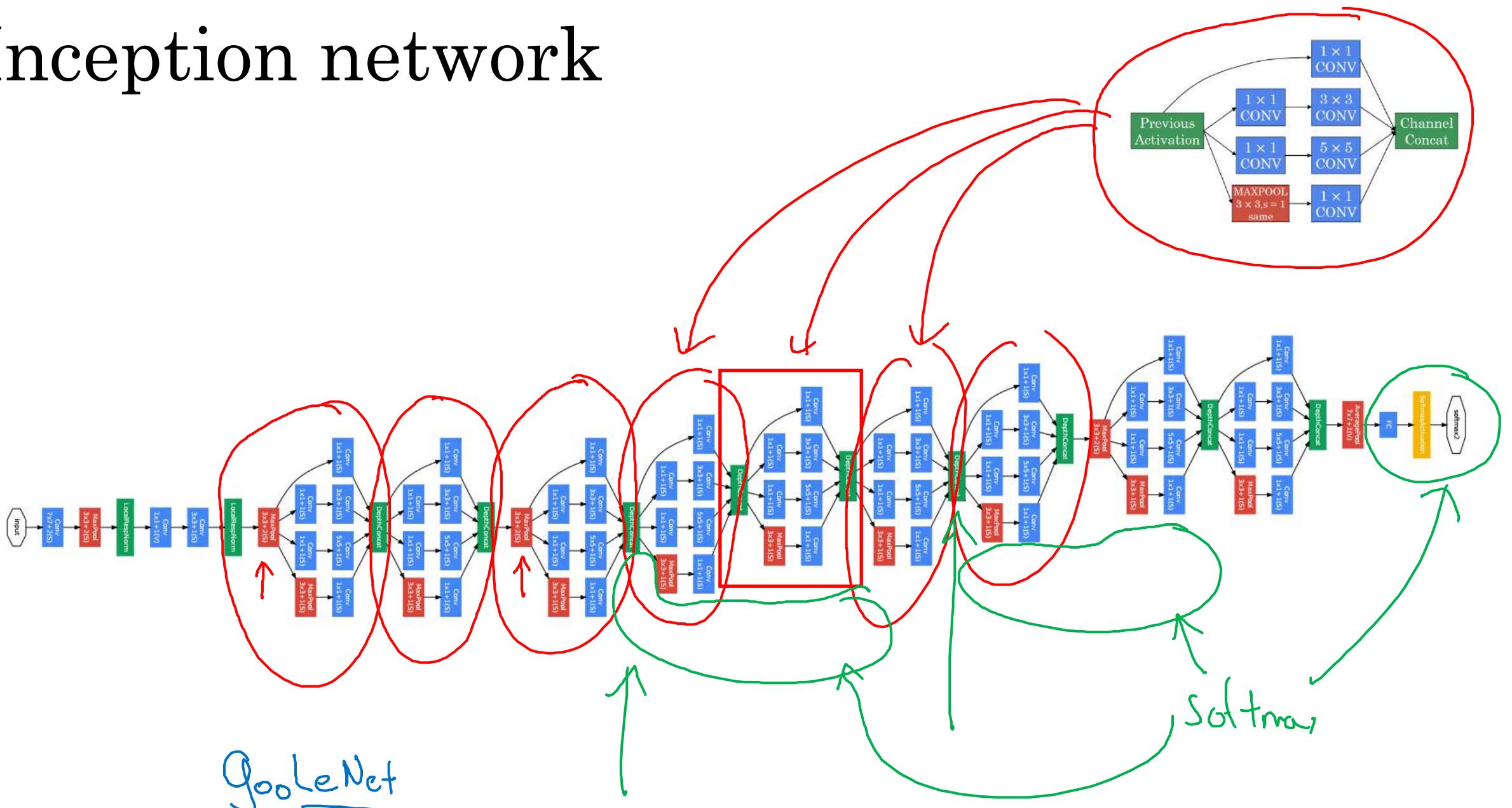
### Inception network

# Inception module



$28 \times 28 \times 32 \leftarrow$   
32 filters,  $1 \times 1 \times 192$ . Andrew Ng

# Inception network







deeplearning.ai

# Convolutional Neural Networks

---

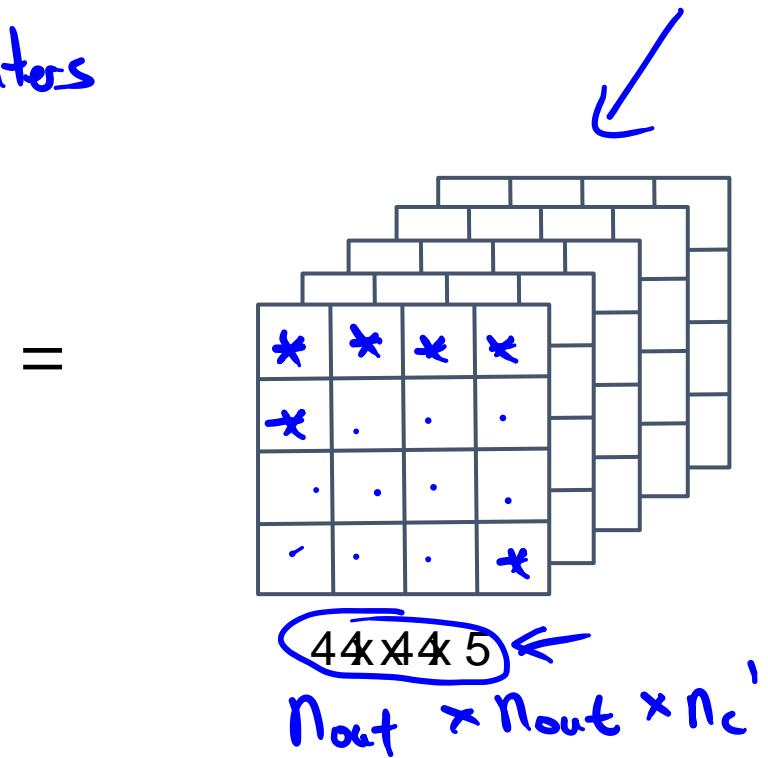
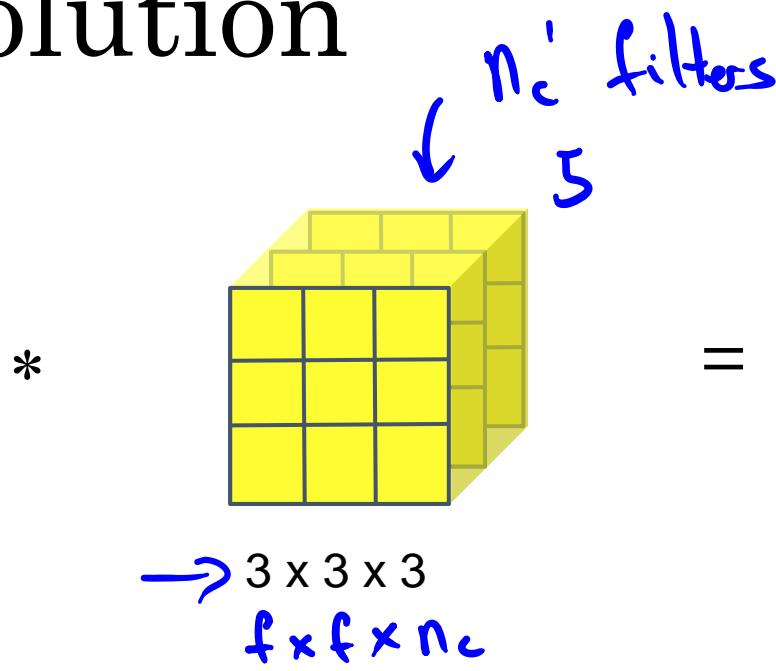
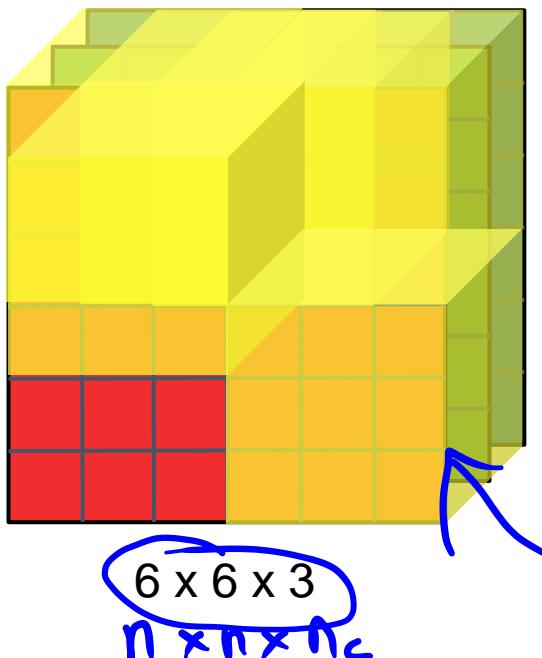
## MobileNet

# Motivation for MobileNets

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



# Normal Convolution



$$\text{Computational cost} = \frac{\# \text{filter params}}{3 \times 3 \times 3} \times \frac{\# \text{filter positions}}{4 \times 4}$$

$$\rightarrow \underline{2160} = \underline{\uparrow} \quad \underline{\uparrow}$$

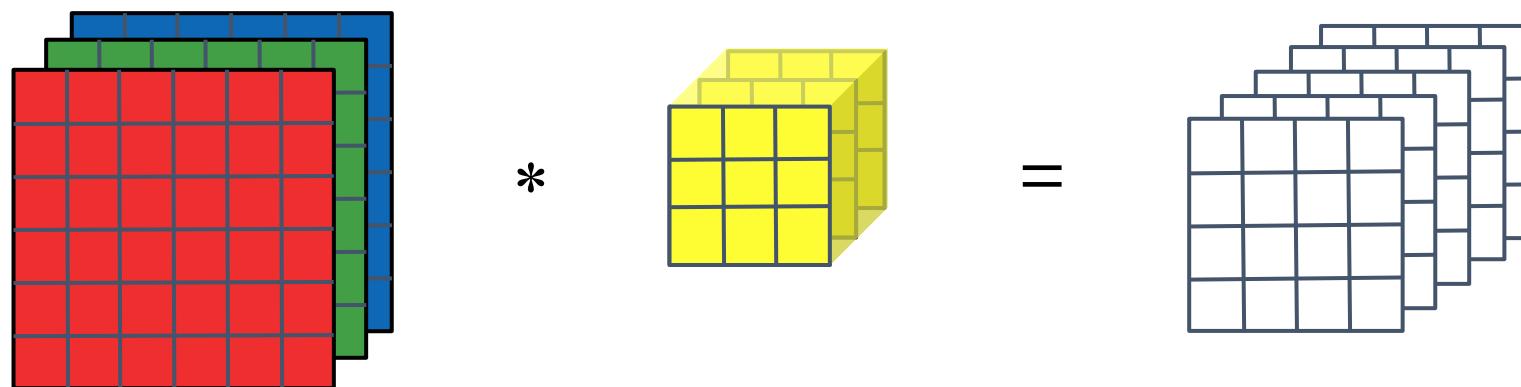
$$x \quad \# \text{of filters}$$

$$\times \quad 5$$

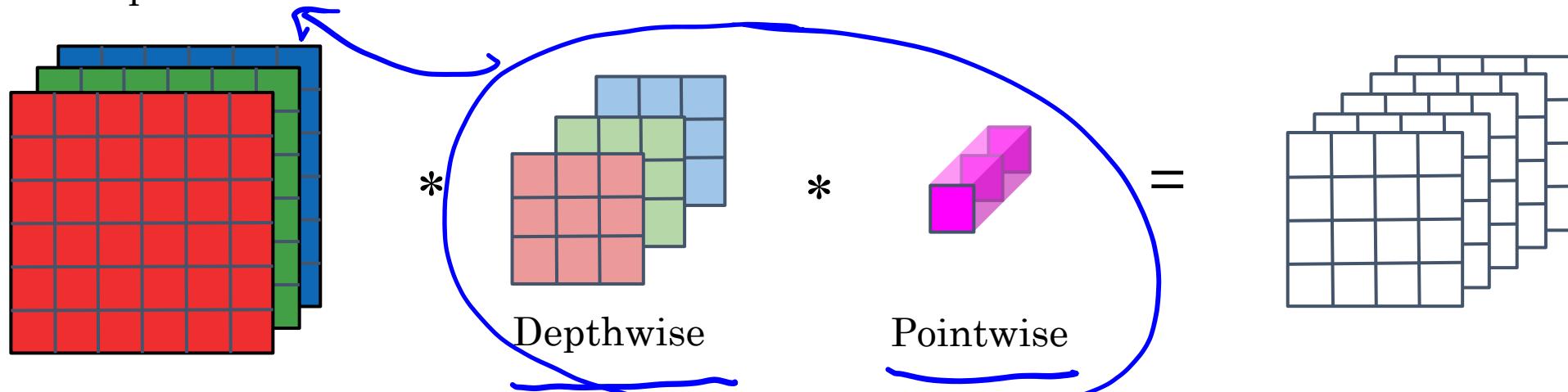
$$\uparrow$$

# Depthwise Separable Convolution

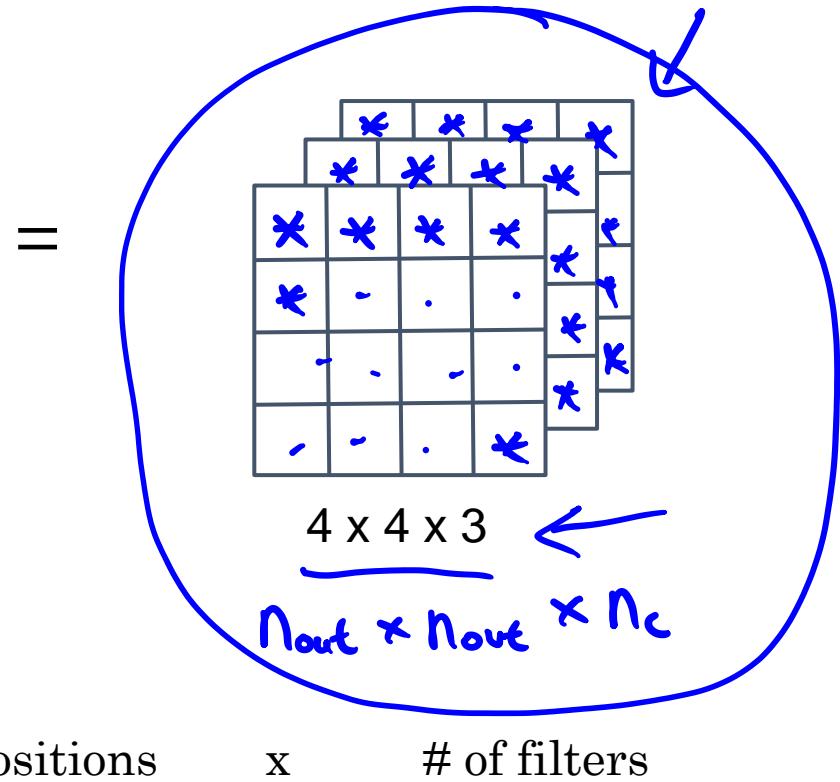
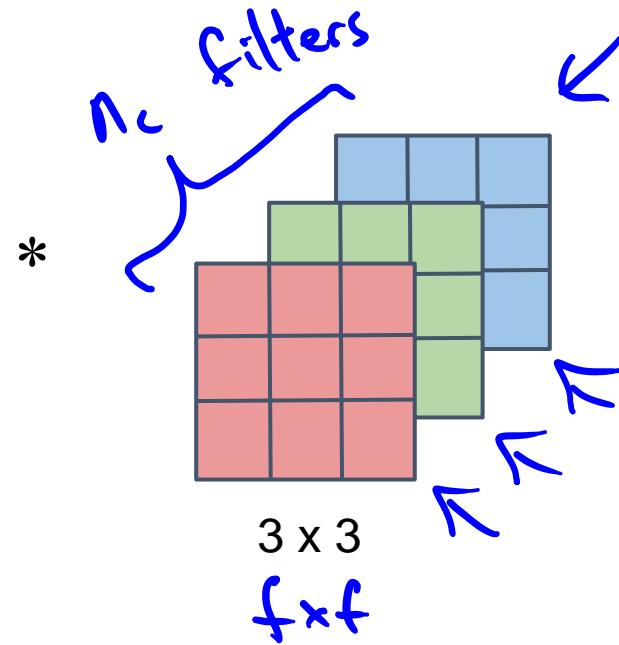
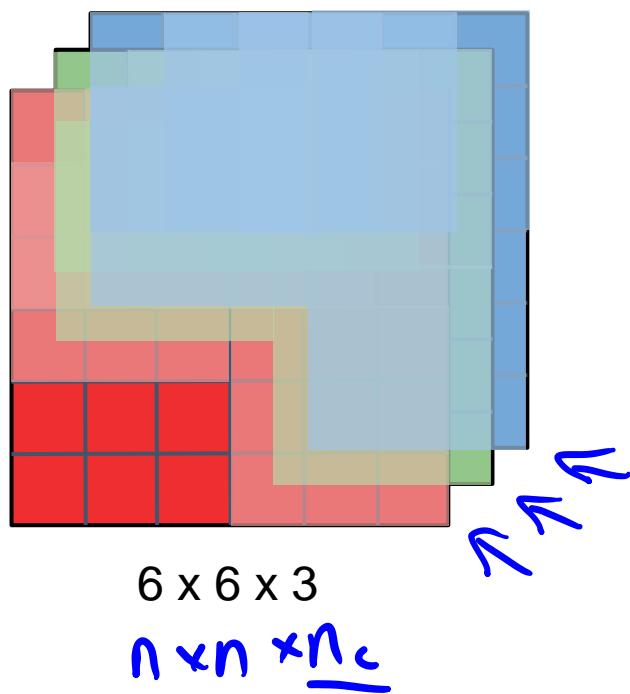
Normal Convolution



Depthwise Separable Convolution



# Depthwise Convolution

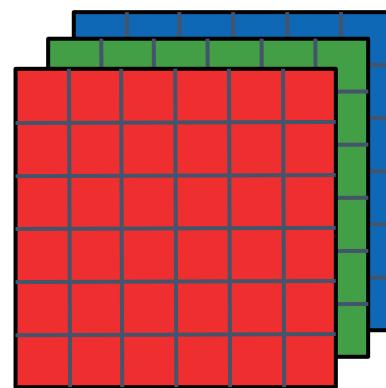


Computational cost = #filter params  $\times$  # filter positions  $\times$  # of filters

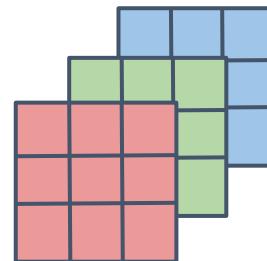
$$432 = \underbrace{3 \times 3}_{\# \text{filter params}} \times \underbrace{4 \times 4}_{\# \text{filter positions}} \times \underbrace{3}_{\# \text{of filters}}$$

# Depthwise Separable Convolution

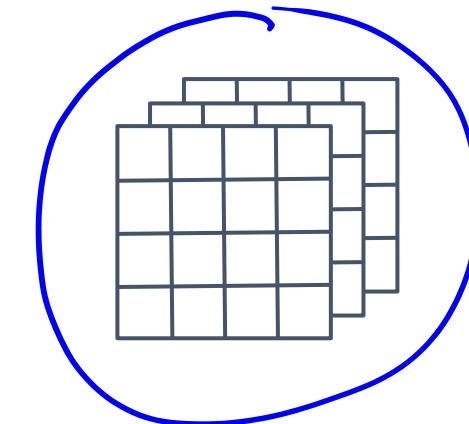
Depthwise Convolution



\*

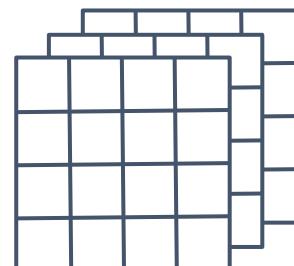


=



432

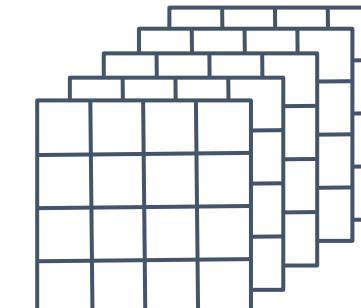
Pointwise Convolution



\*

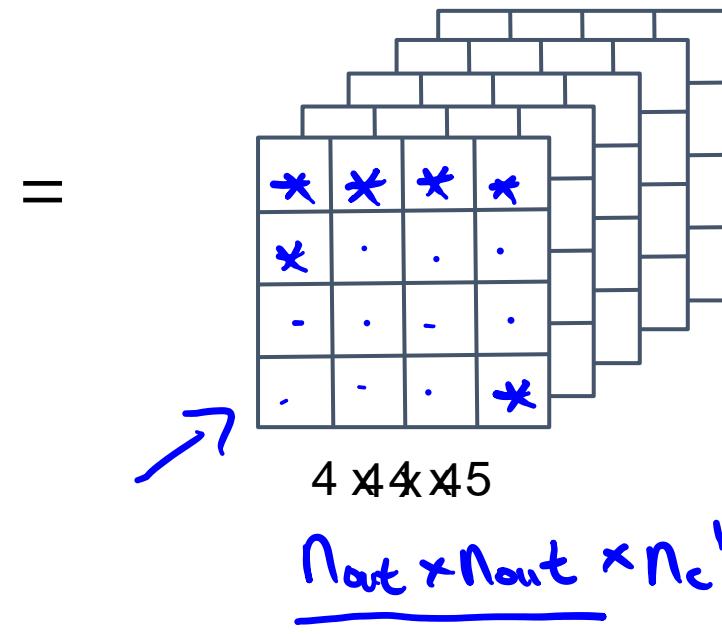
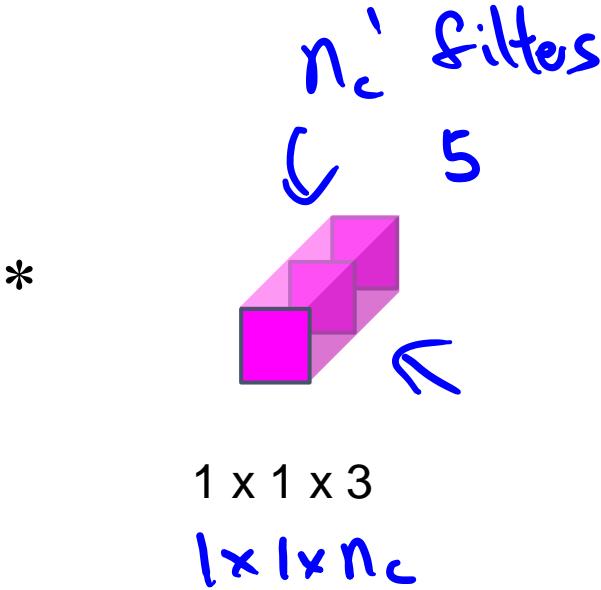
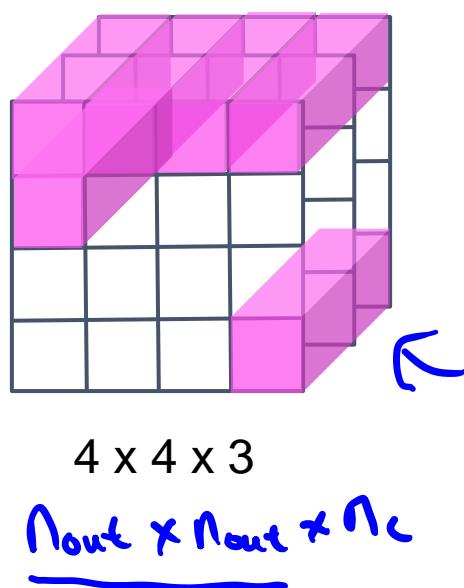


=



4x4x5

# Pointwise Convolution

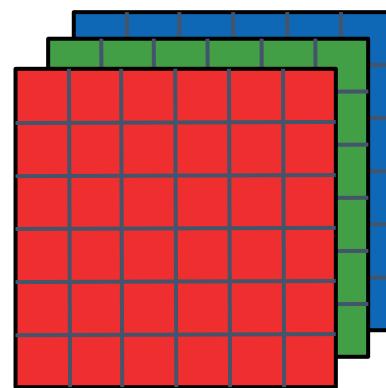


Computational cost = #filter params  $\times$  # filter positions  $\times$  # of filters

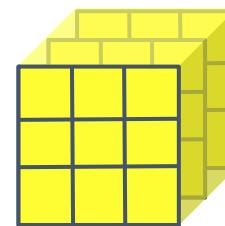
$$240 = 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

# Depthwise Separable Convolution

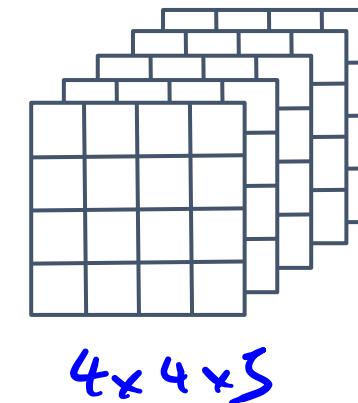
Normal Convolution



\*



=

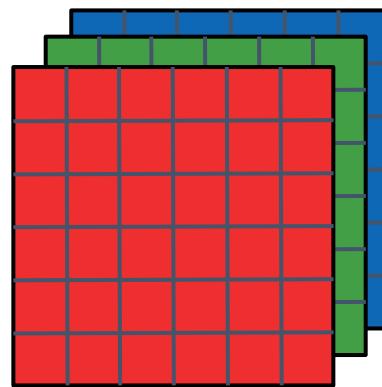


$$4 \times 4 \times 5 \times 3 \times 3 \times 3 \Rightarrow 4 \times 4 \times 3 \times (5 \times 9) = 2160$$

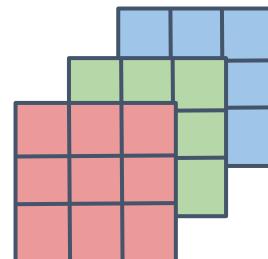
changed from  
x to +

$$\begin{aligned} \text{ratio} &= 1/5 + 1/9 \\ &= 1/n_{c'} + 1/f^{**2} \end{aligned}$$

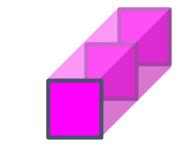
Depthwise Separable Convolution



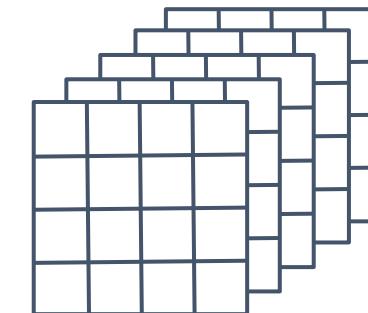
\*



\*



=



Depthwise

Pointwise

$$4 \times 4 \times 3 \times 3 + 4 \times 4 \times 5 \times 3 \Rightarrow 4 \times 4 \times 3 \times (5 + 9) = 672$$

# Cost Summary

Cost of normal convolution

2160

Cost of depthwise separable convolution

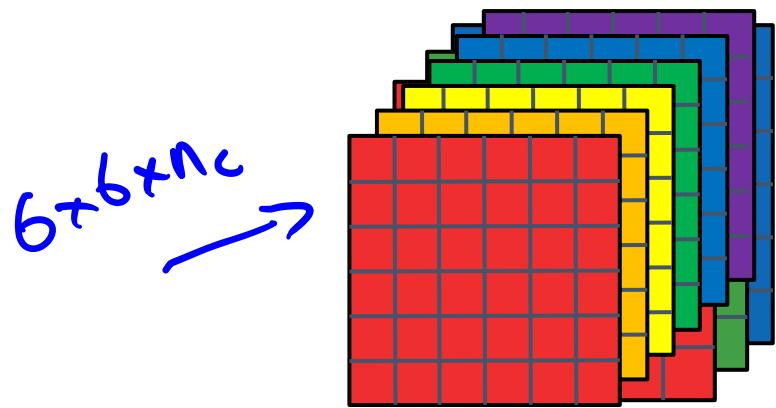
$$\begin{array}{l} \text{depthwise} + \text{pointwise} \\ 432 + 240 = 672 \end{array}$$

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

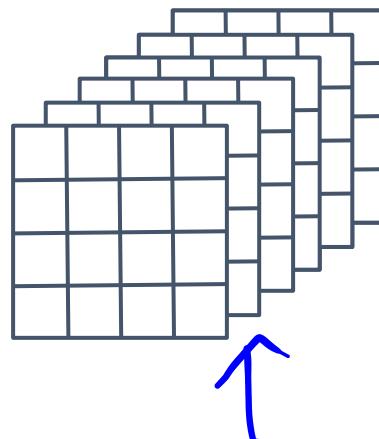
$$\begin{aligned} &= \frac{1}{n_c} + \frac{1}{f^2} \\ &\quad \frac{1}{s} + \frac{1}{q} \\ &= \frac{1}{512} + \frac{1}{3^2} \\ &\quad \nwarrow \quad \swarrow \\ &\text{n10 times cheaper} \end{aligned}$$

# Depthwise Separable Convolution

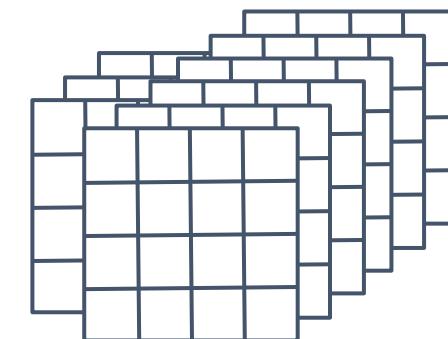
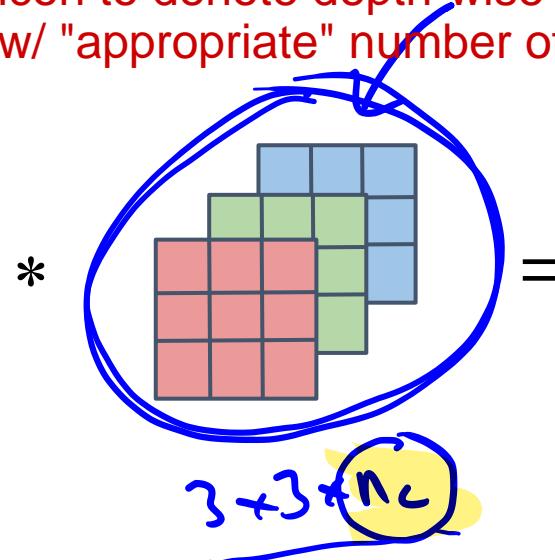
Depthwise Convolution



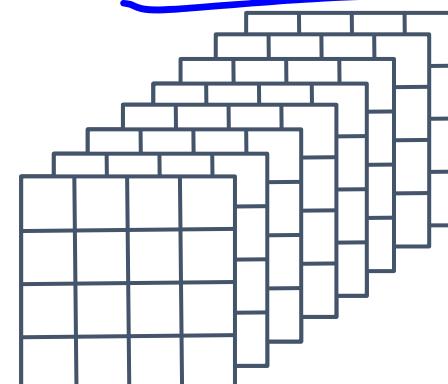
Pointwise Convolution



icon to denote depth-wise convolution  
w/ "appropriate" number of layers



$4 \times 4 \times n_c$



$4 \times 4 \times 8$ ,  
 $n_c$



deeplearning.ai

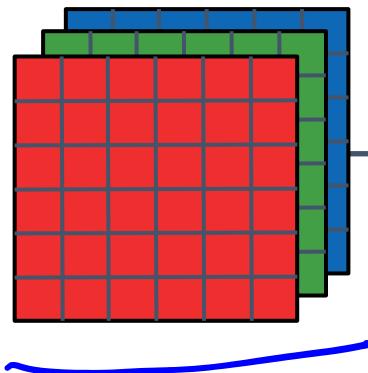
# Convolutional Neural Networks

---

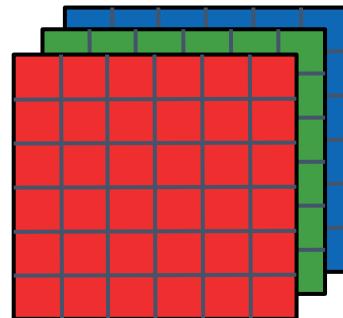
## MobileNet Architecture

# MobileNet

MobileNet v1



MobileNet v2



13 times

POOL, FC, SOFTMAX

17 times

Residual Connection

POOL, FC,  
SOFTMAX

Expansion

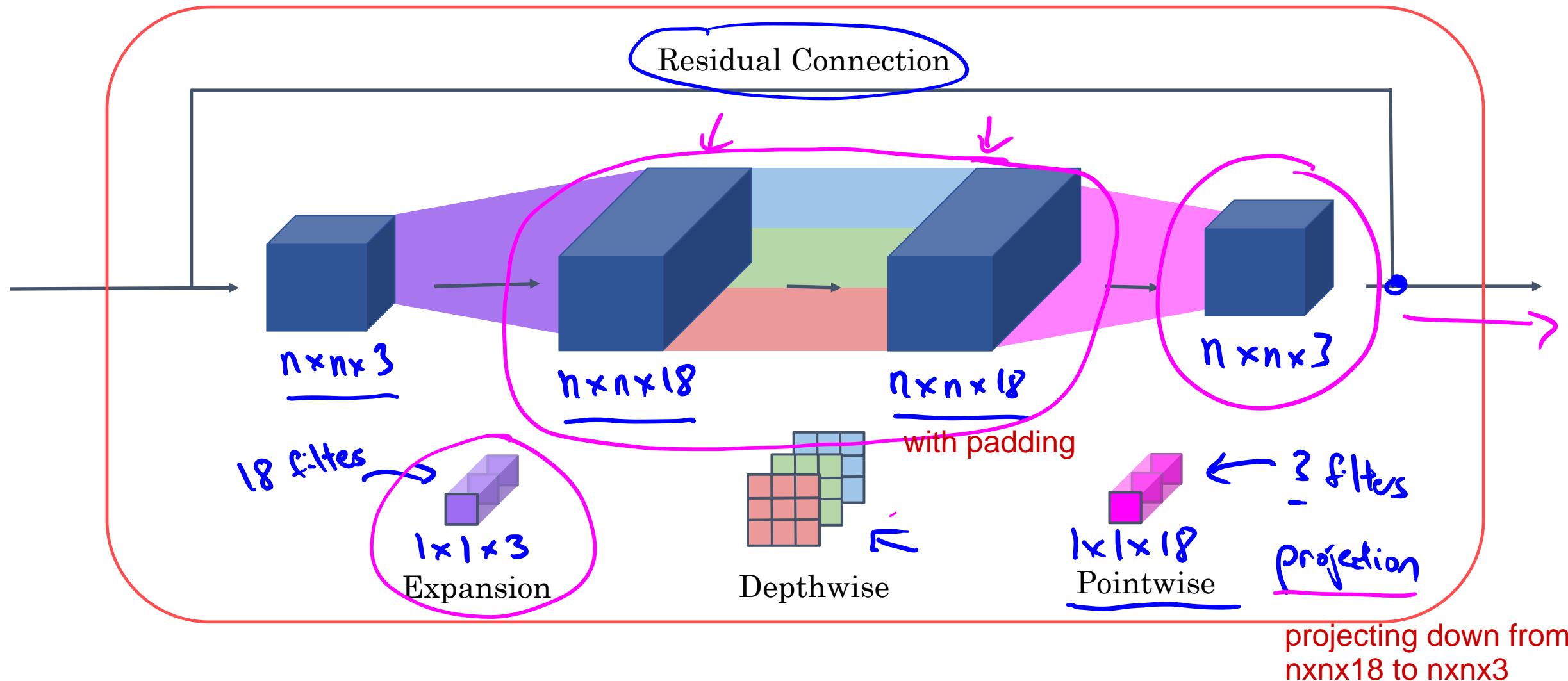
Depthwise

Projection

Bottleneck

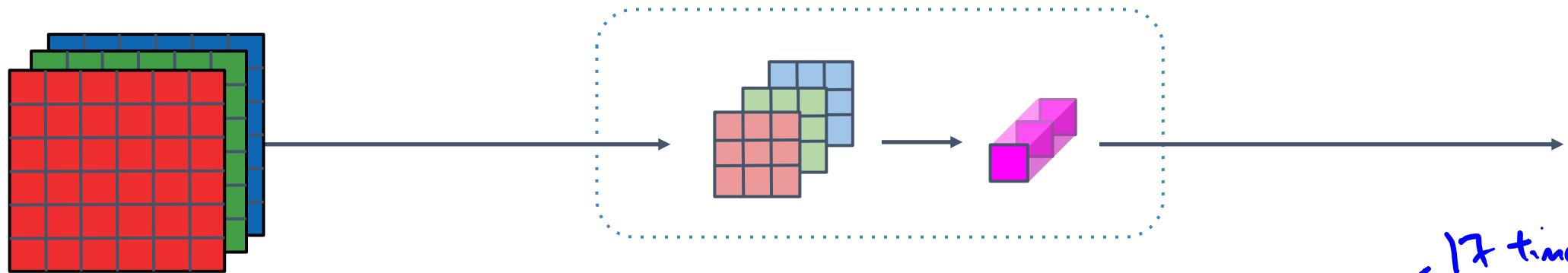
# MobileNet v2 Bottleneck

clever idea here is - keep memory low (in term of activations moved from L to L+1), by projecting down... but it allows NN to learn a richer / complex function within the layer

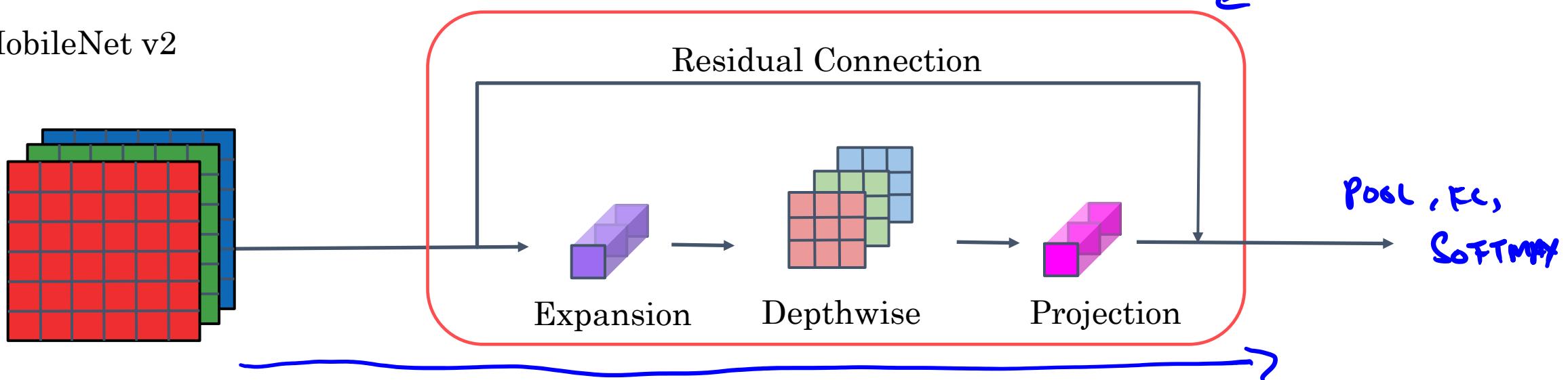


# MobileNet

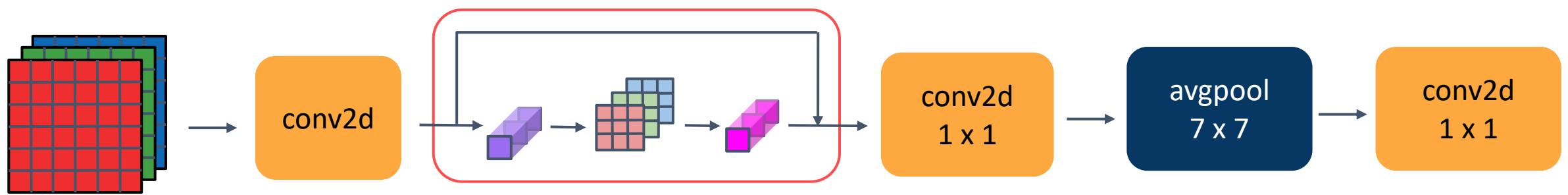
MobileNet v1



MobileNet v2



# MobileNet v2 Full Architecture





deeplearning.ai

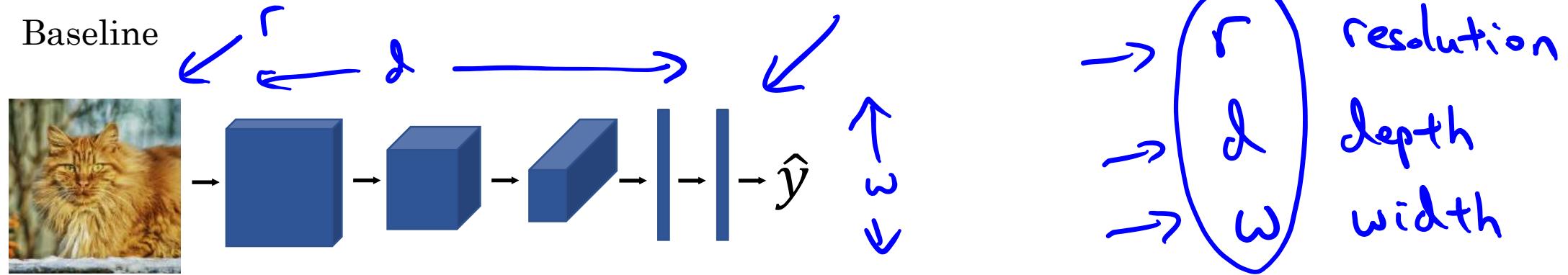
# Convolutional Neural Networks

---

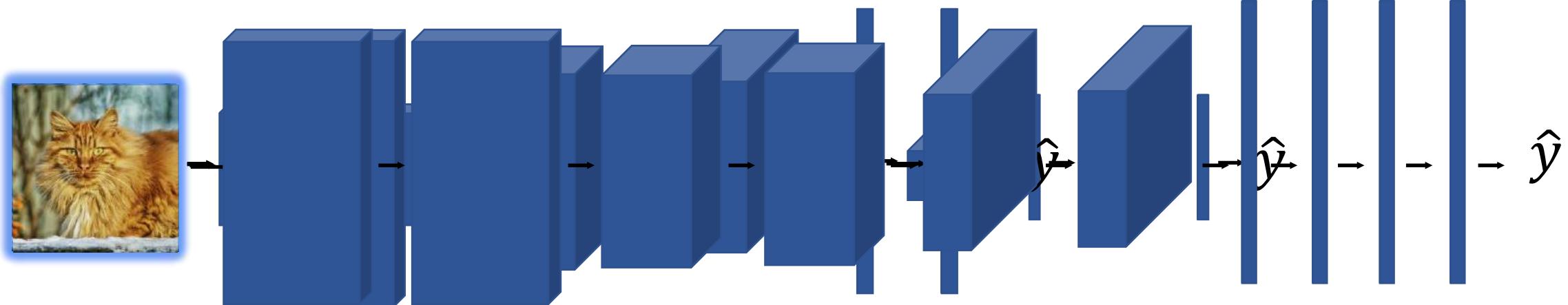
## EfficientNet

How to scale the network up or down based on the compute available on a device - think in terms of r/ d/ w and look for an open source implementation of efficient net to calibrate

# EfficientNet



Width Rescaling





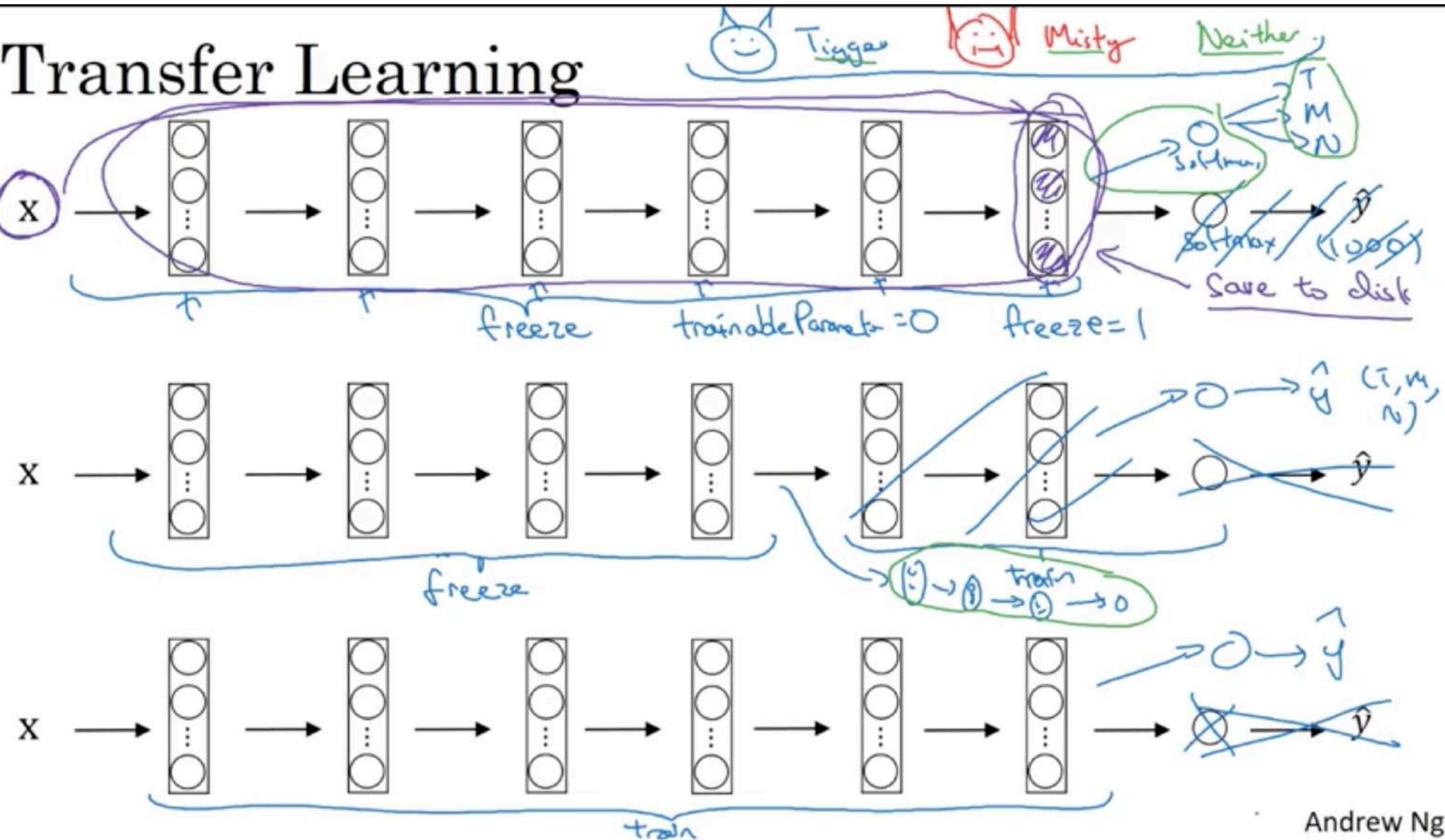
deeplearning.ai

Practical advice for  
using ConvNets

---

Transfer Learning

# Transfer Learning



1st approach: few images - take imagenet (1000classes), replace the last layer with 3 unit softmax. Convert  $X$  to  $X'$  feature set by running it through the network upto the replaced layer and save to disk. Then it's like training a single layer shallow NN using  $X'$ .  
 2nd approach - have more pics to train on, then retrain >1 layers at the end. Extreme case 3rd approach - use trained weights as initialization points, and train entire network.



deeplearning.ai

Practical advice for  
using ConvNets

---

Data augmentation

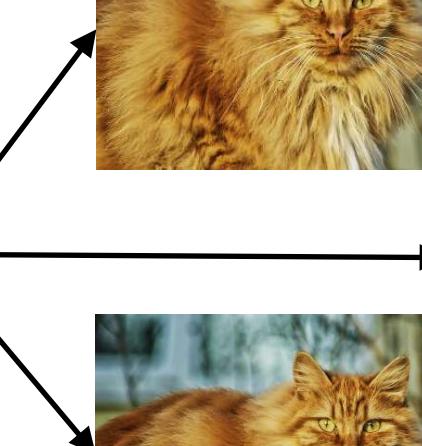
# Common augmentation method

Mirroring



yc

Random Cropping

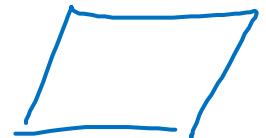


Rotation

Shearing

Local warping

...



# Color shifting



R G B  
↓ ↓ ↓  
+20,-20,+20



-20,+20,+20



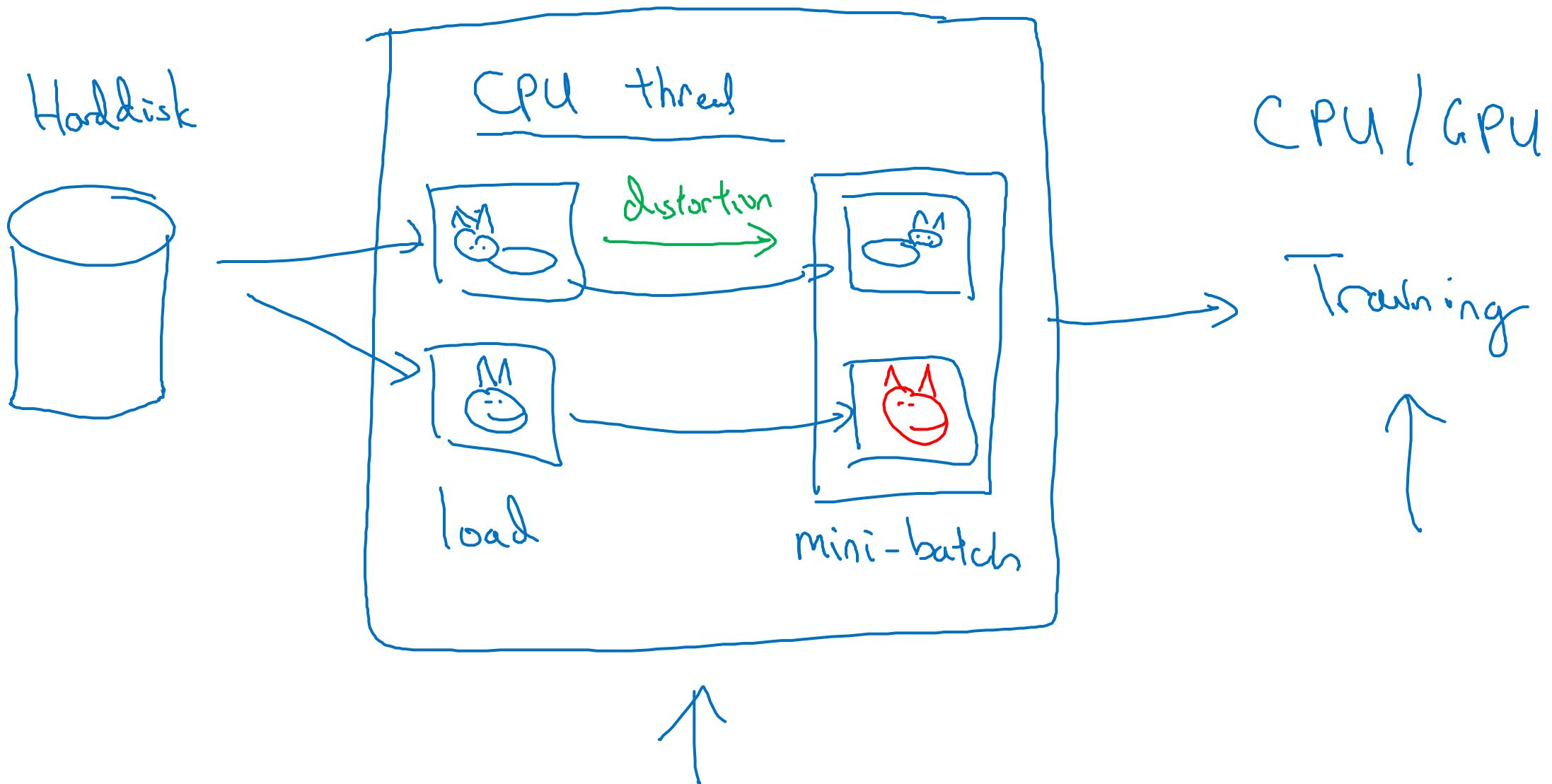
+5,0,+50



y

Advanced:  
PCA  
[ml-class.org](http://ml-class.org)  
[ AlexNet paper  
[ "PCA color augmentation."  
R B      G

# Implementing distortions during training





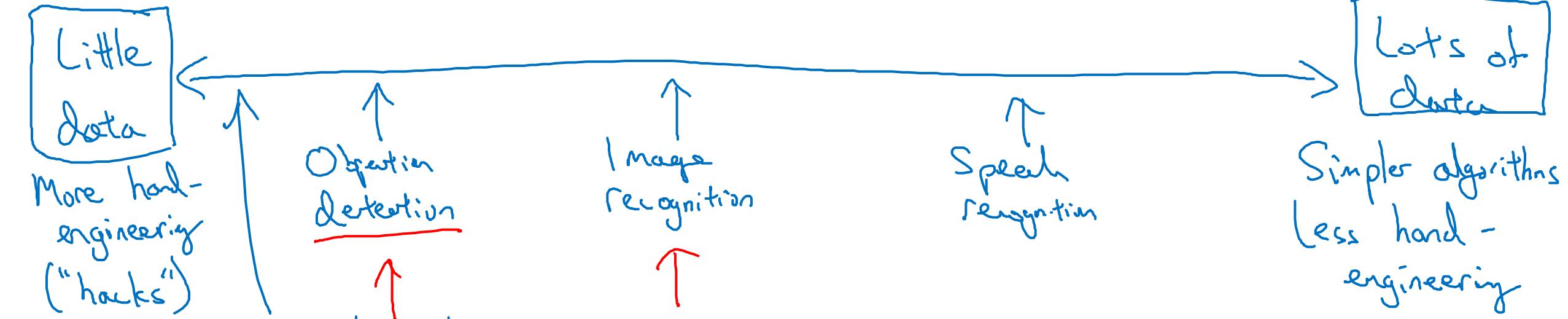
deeplearning.ai

Practical advice for  
using ConvNets

---

The state of  
computer vision

# Data vs. hand-engineering



Two sources of knowledge

- • Labeled data  $(x, y)$
- • Hand engineered features/network architecture/other components



# Tips for doing well on benchmarks/winning competitions

## Ensembling

3 - 15 networks

$\rightarrow \hat{y}$

- Train several networks independently and average their outputs

Good for competition, not so much for production system, because of high compute requirement + necessity to store data.

## Multi-crop at test time

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

10-crop

Take image - create 1 flipped copy. Create 5 crops for each image = "10-crop"



1

+

4

+

1

+

4



# Use open source code

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

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# Object Detection

---

## Object localization

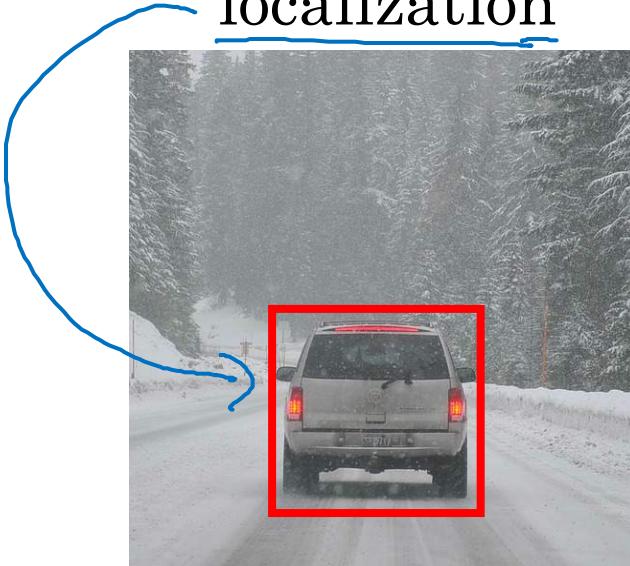
# What are localization and detection?

Image classification



"Car"

Classification with  
localization



"Car"

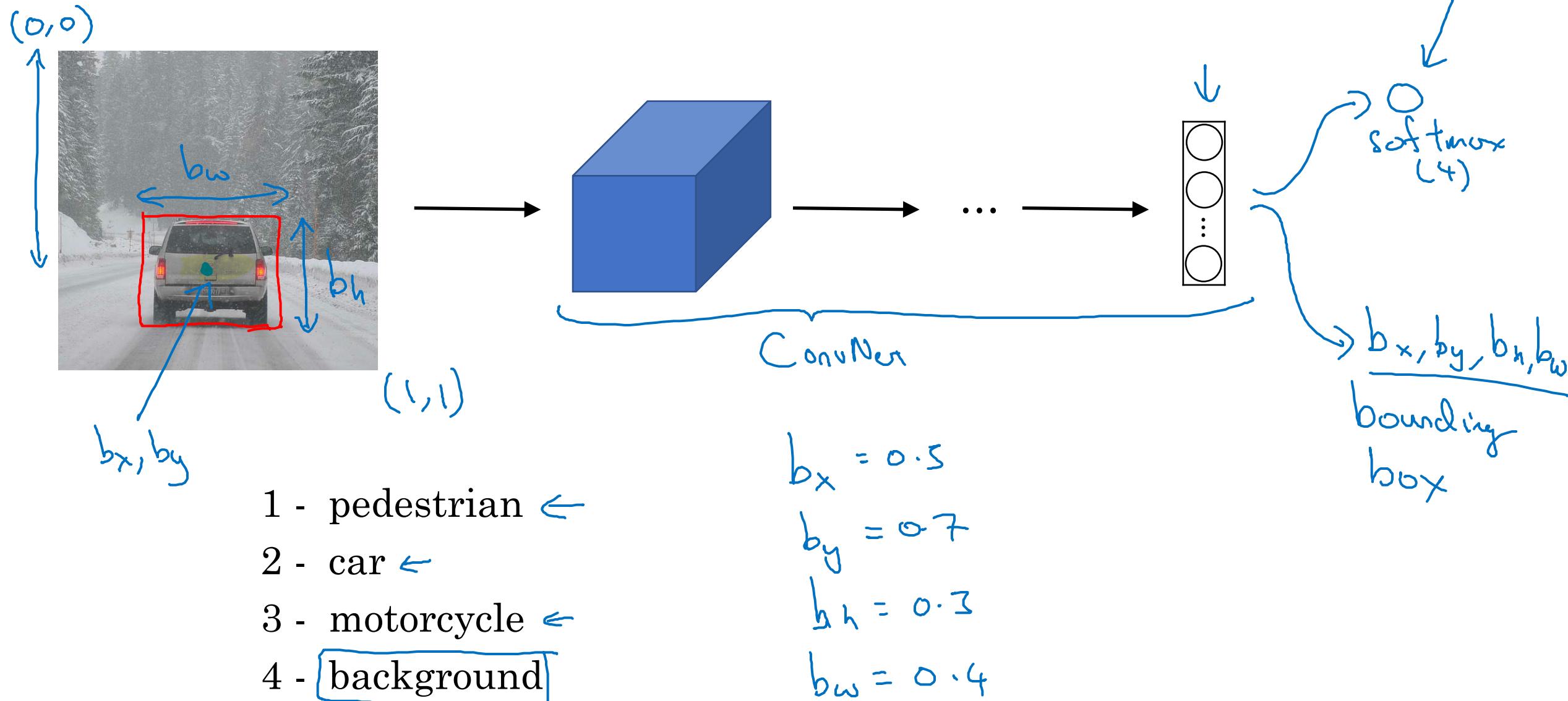
1 object

Detection



multiple  
objects

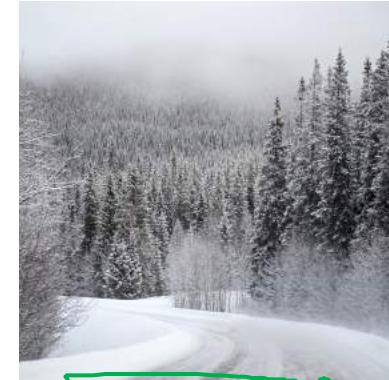
# Classification with localization



# Defining the target label $y$

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

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

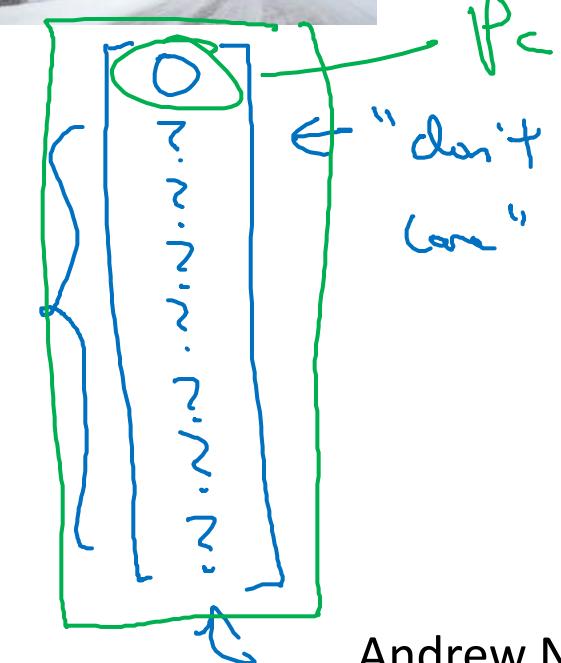
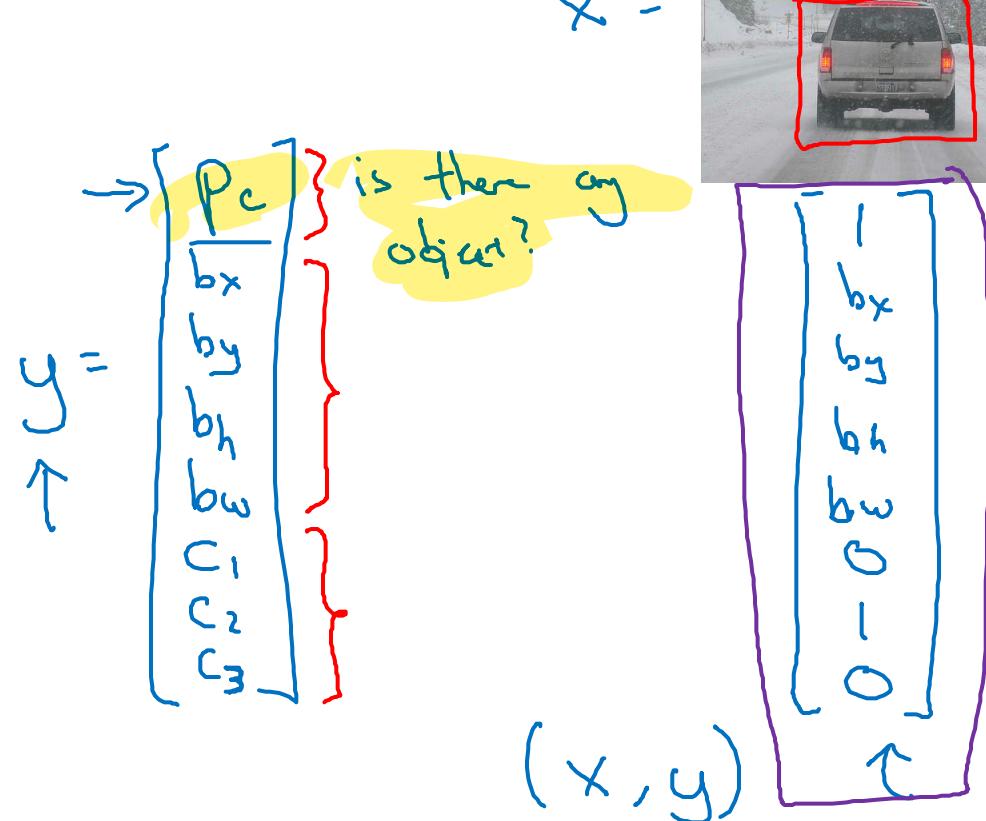


$$L(\hat{y}, y) = \text{RMSE} \text{ (ideally use RMSE + LogLoss)}$$

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

$$(\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0$$

in case of  $P_c = 0$  (i.e. background) we only care about  $P_c$ , not other values





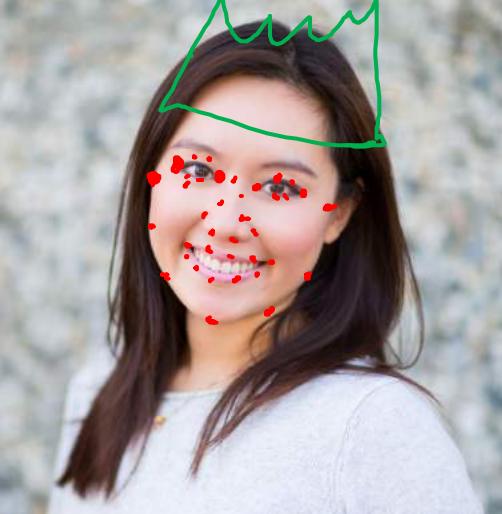
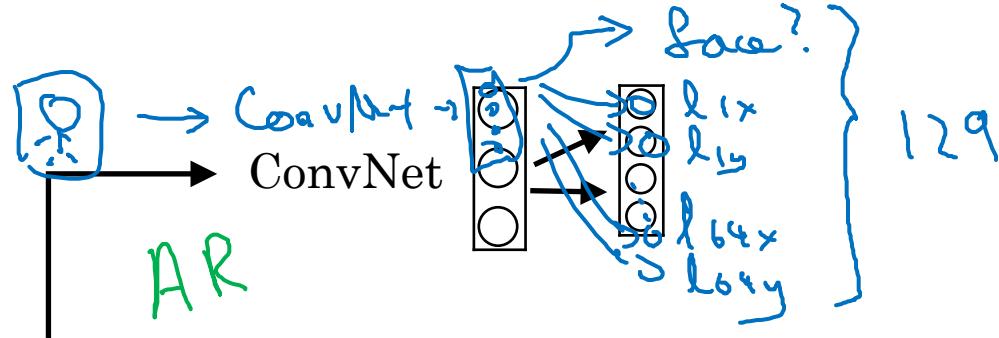
deeplearning.ai

# Object Detection

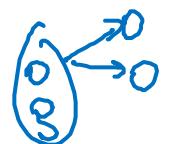
---

## Landmark detection

# Landmark detection



$b_x, b_y, b_h, b_w$



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

X, Y

$l_{1x}, l_{1y},$   
:  
 $l_{32x}, l_{32y}$



deeplearning.ai

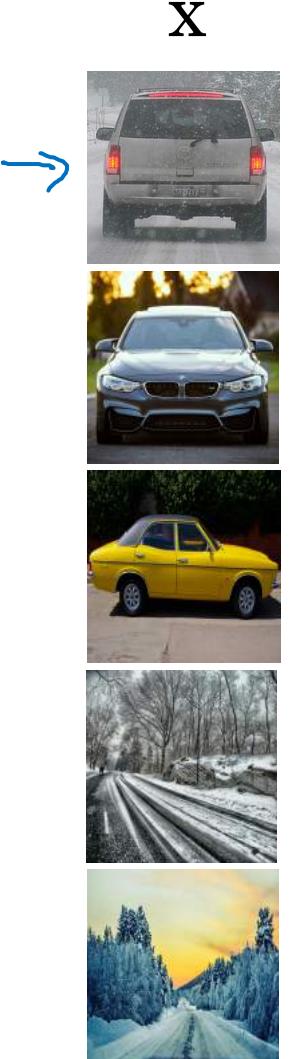
# Object Detection

---

## Object detection

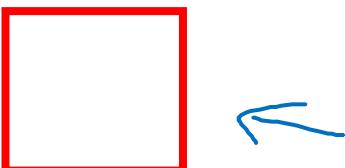
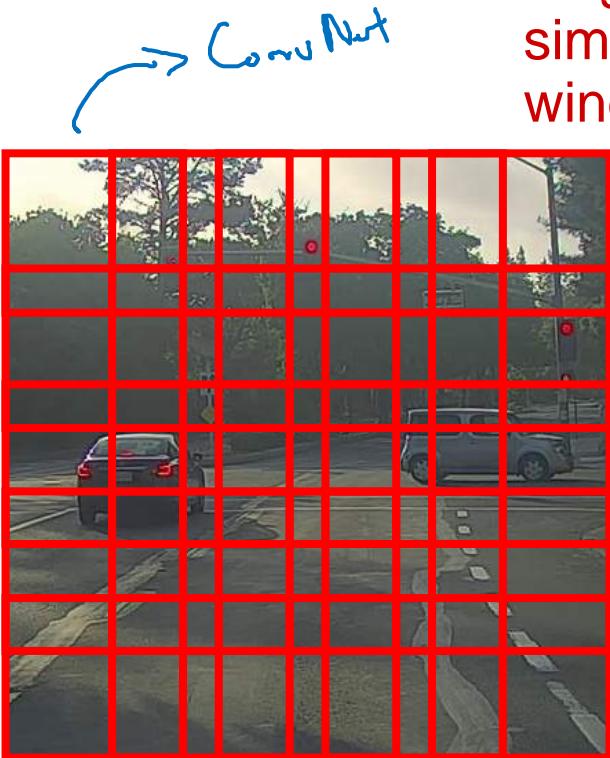
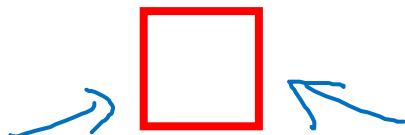
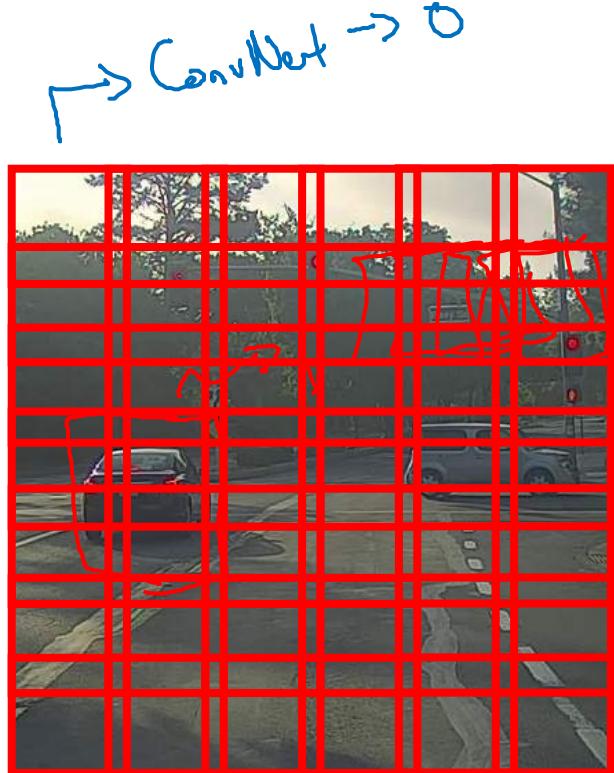
# Car detection example

Training set:



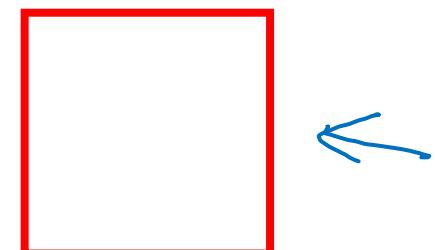
$\xrightarrow{\hspace{1cm}}$  ConvNet  $\rightarrow y$

# Sliding windows detection



Computation cost

Disadvantage - computational cost of running multiple sliding windows.  
Originally (pre-CNN), classifiers were simple linear classifiers, so sliding window was feasible. Not so with CNNs.





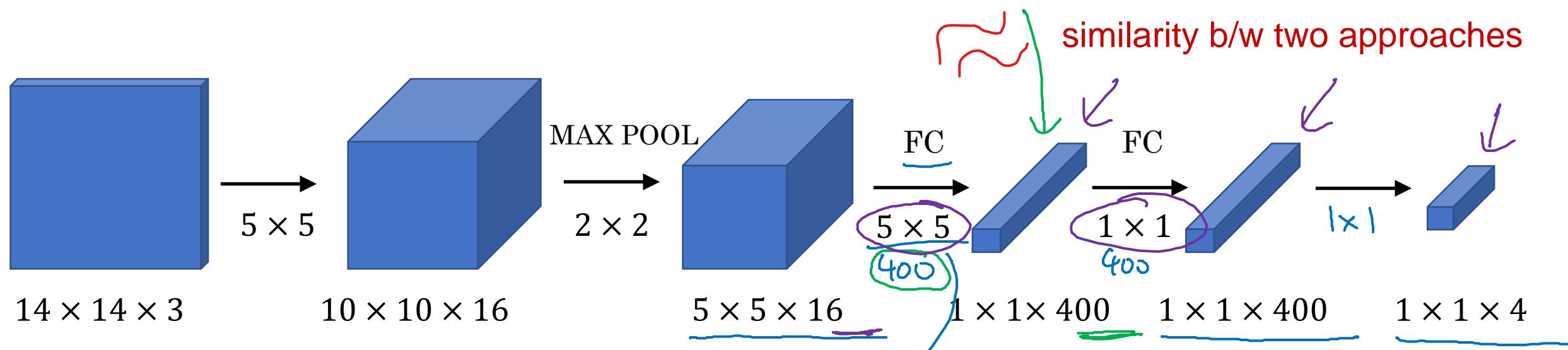
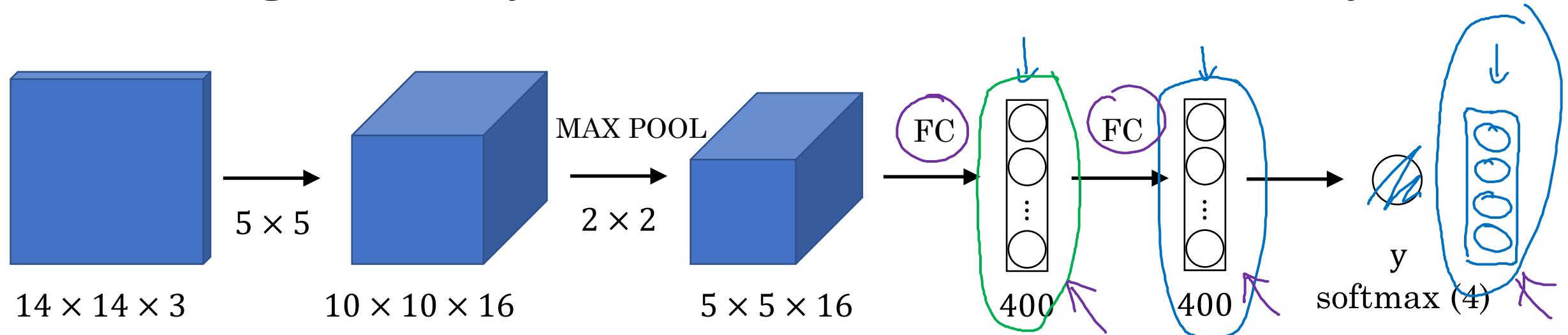
deeplearning.ai

# Object Detection

---

## Convolutional implementation of sliding windows

# Turning FC layer into convolutional layers

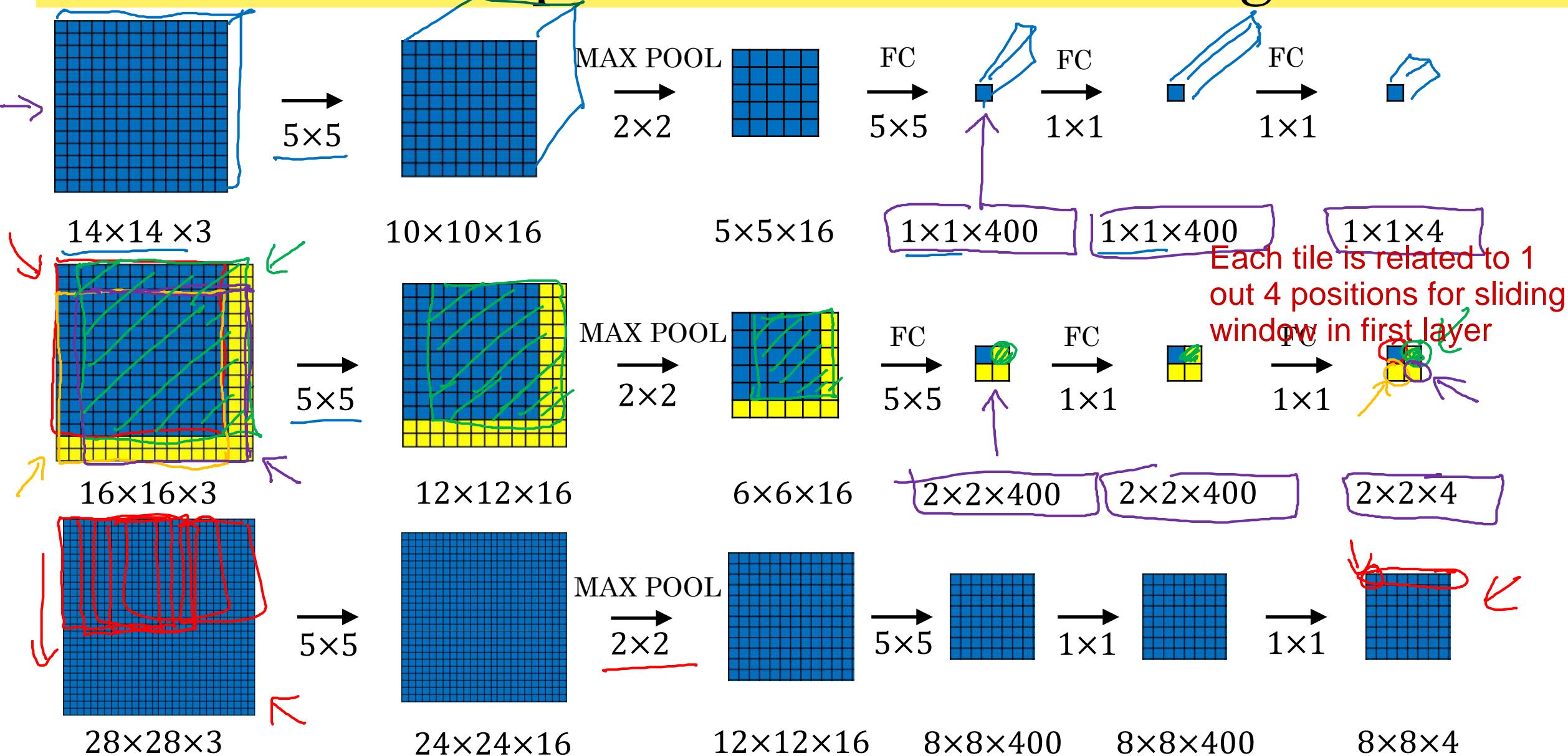


can think of as: 400 equations & 400 variables

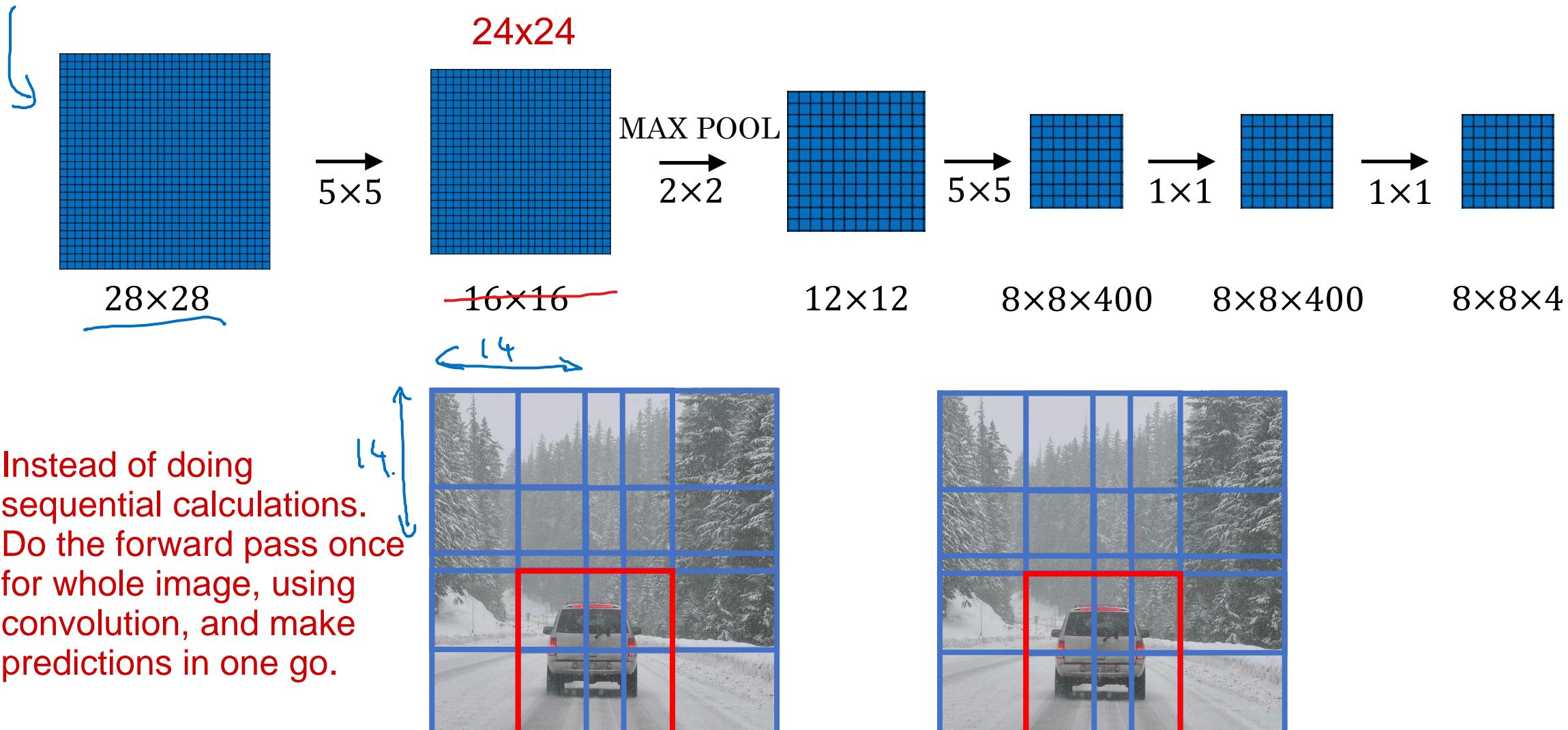
$5 \times 5 \times 16$

Andrew Ng

# Convolution implementation of sliding windows



# Convolution implementation of sliding windows



Negative - position of bounding box is still not going to be very accurate



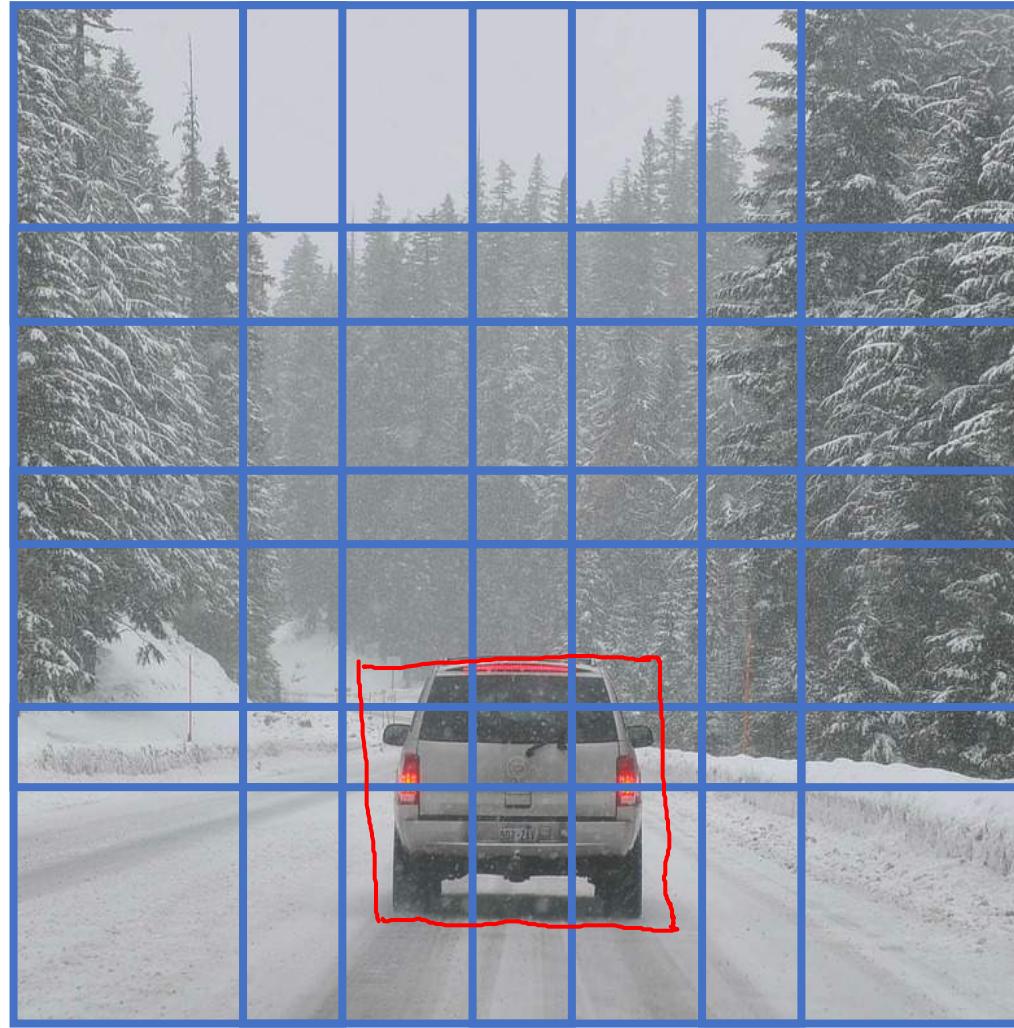
deeplearning.ai

# Object Detection

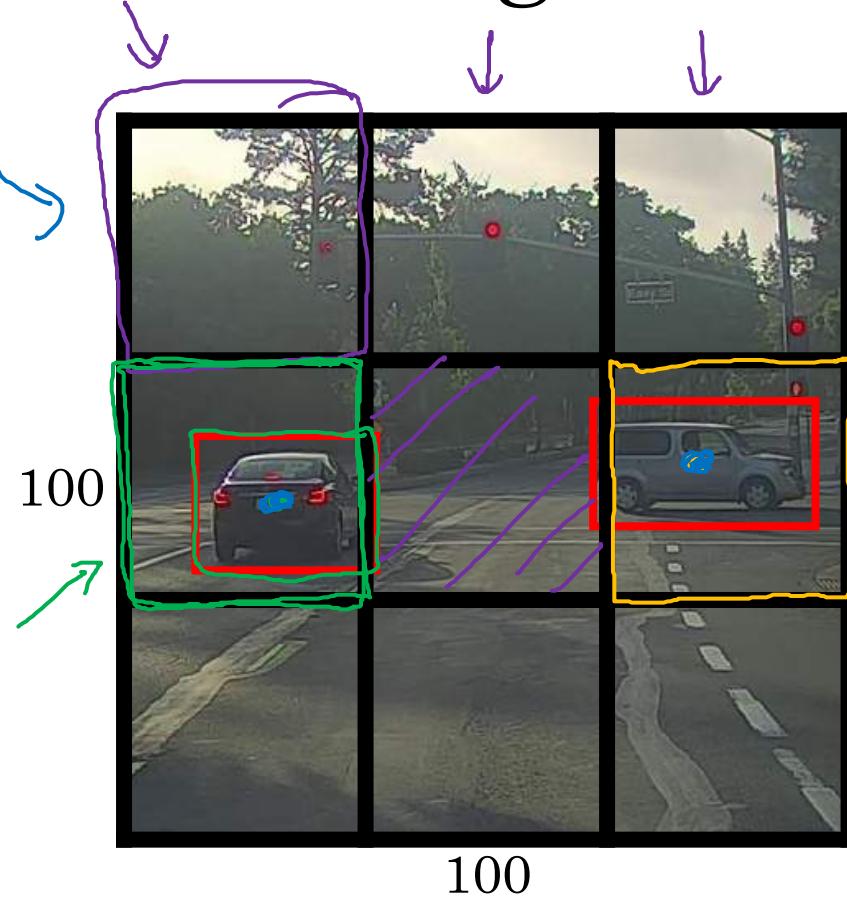
---

Bounding box  
predictions

# Output accurate bounding boxes



# YOLO algorithm



You Only Look Once

Apply fine grid to the image (19x19) - to each grid apply image classification & localization algorithm.

Object is assigned to the grid where its 'center' lies.

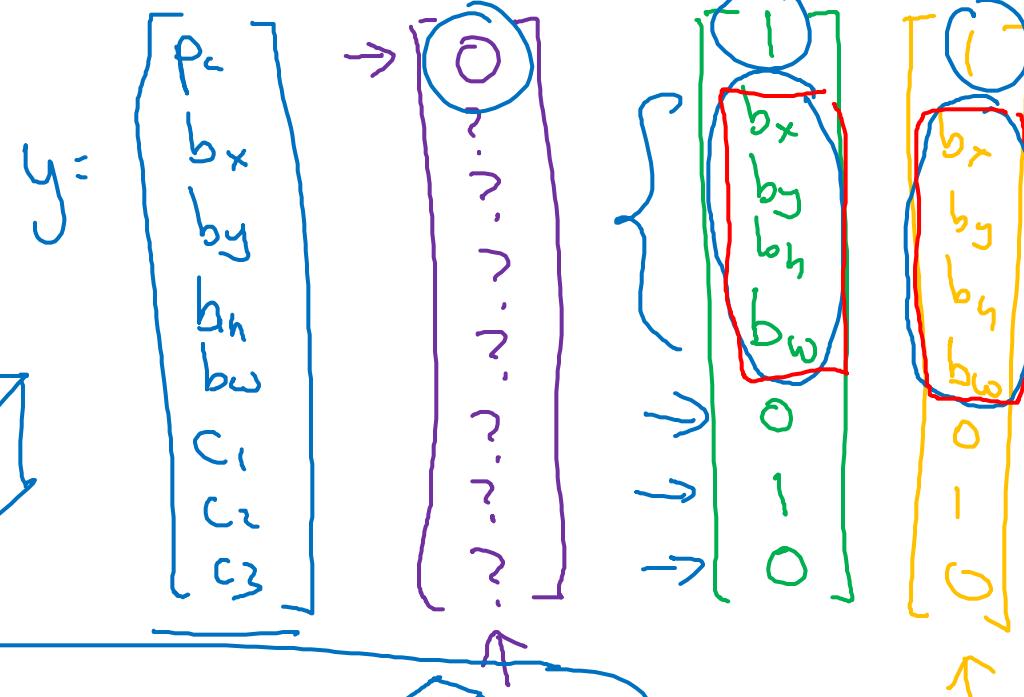
Labels for training

For each grid cell:

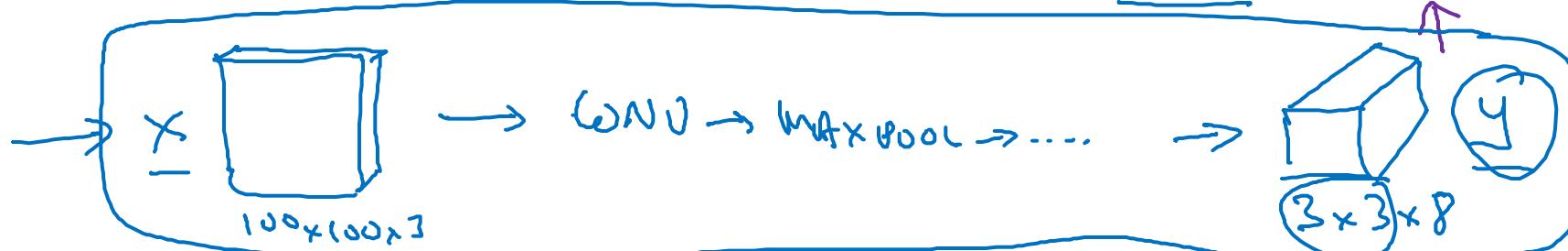
we want 8 o/p values  
for each cell on grid,  
for total of  $3 \times 3 \times 8$

Target output:

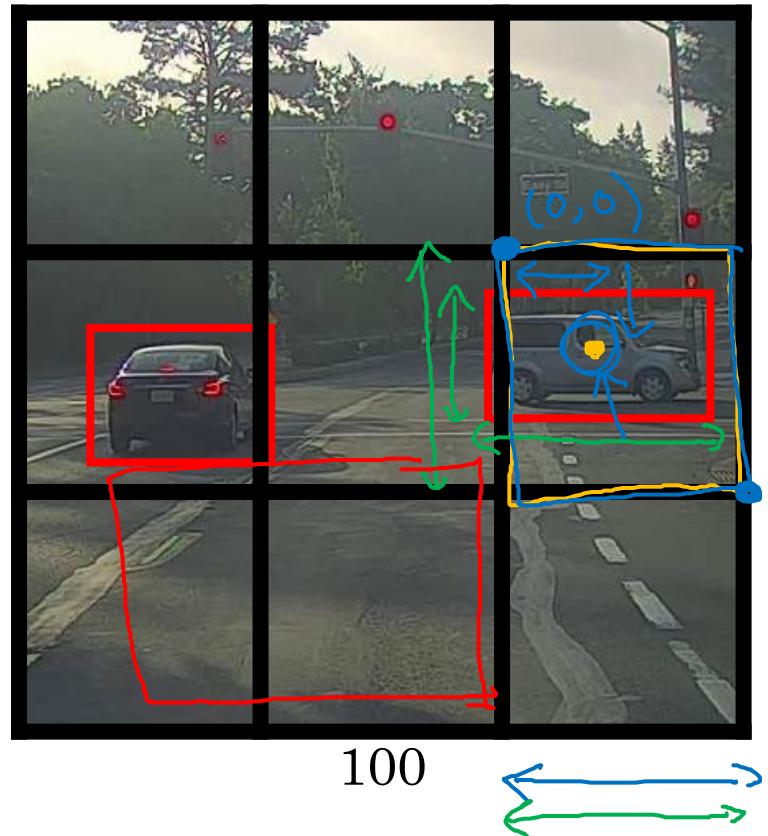
$$3 \times 3 \times 8$$



361



# Specify the bounding boxes



$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ o \\ o \end{bmatrix}$$

0.4 } between 0 and 1  
0.3 }  
0.9 }  
0.5 }  
  
bottom two values can be >1

Could be >1



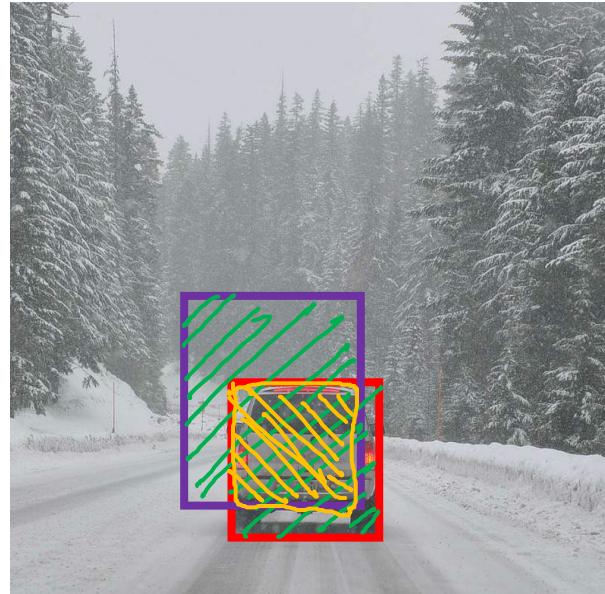
deeplearning.ai

# Object Detection

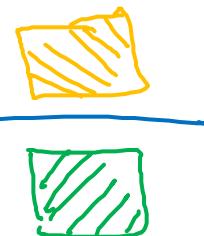
---

Intersection  
over union

# Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of intersection}}{\text{Size of union}}$$


“Correct” if  $\text{IoU} \geq 0.5$

by convention  $>0.5$  is considered good enough

0.6

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



deeplearning.ai

# Object Detection

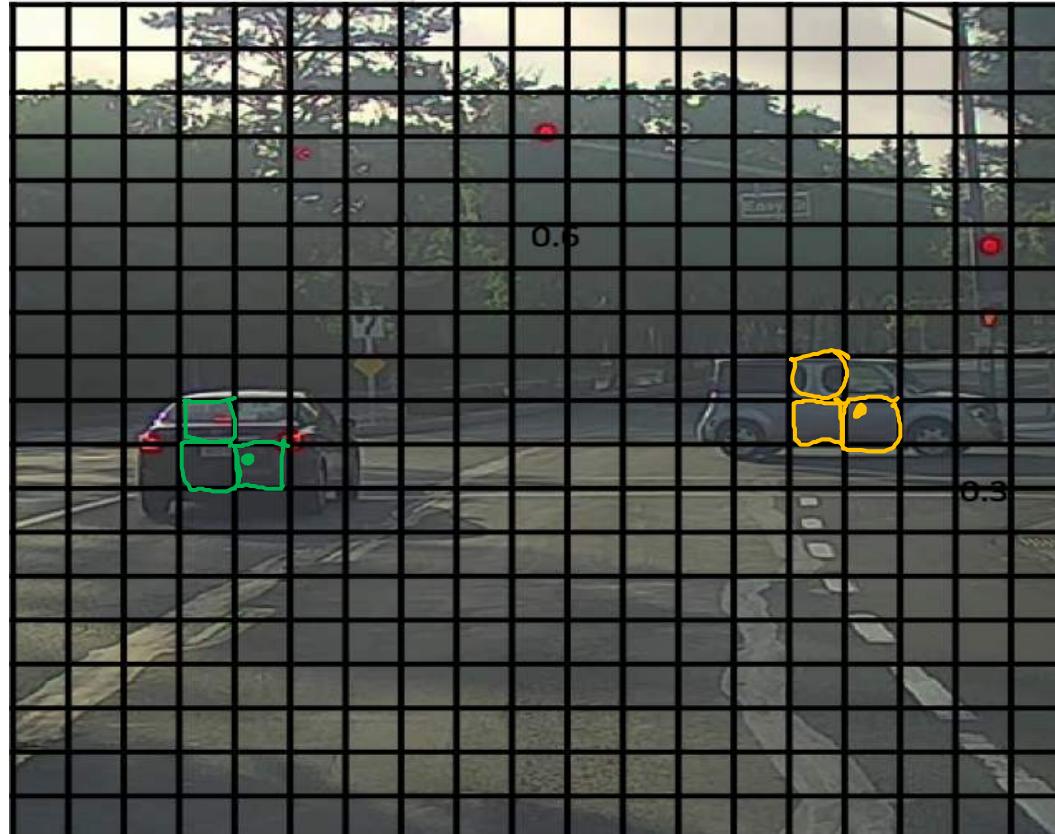
---

## Non-max suppression

# Non-max suppression example

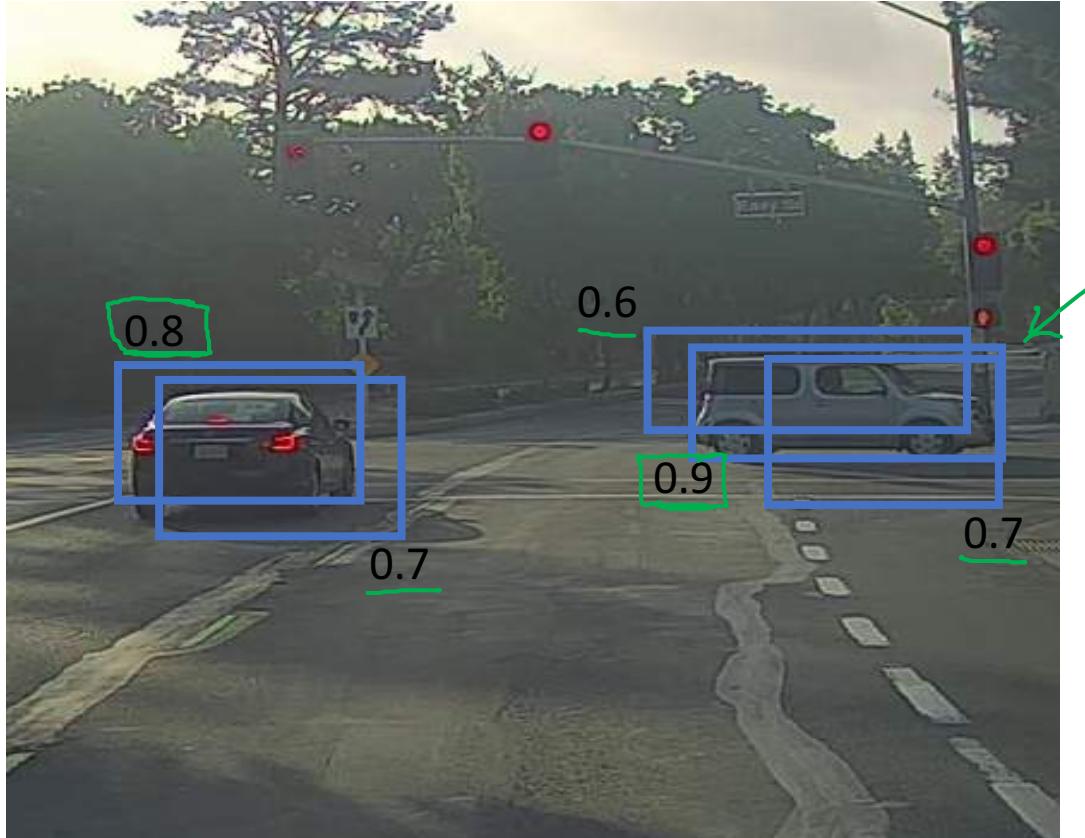


# Non-max suppression example



19x19

# Non-max suppression example



P<sub>c</sub>

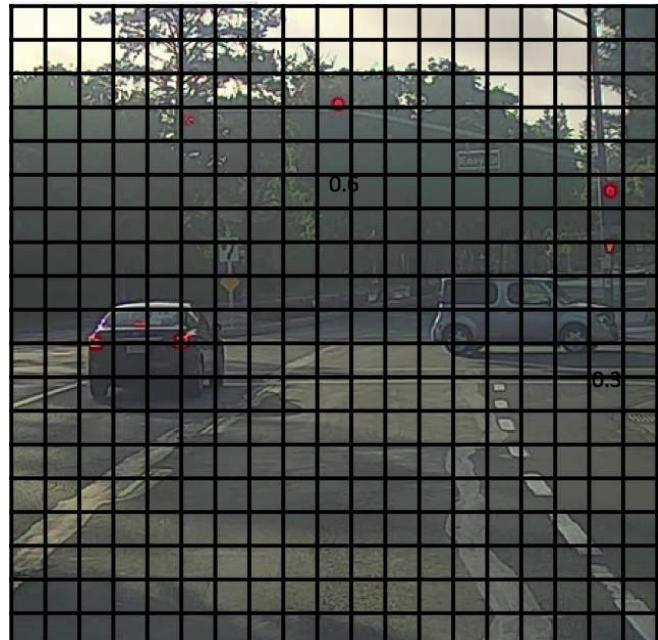
Issue - many boxes will detect any single car and generate  $P_c > 0$

- pick the one with the max value (most confident detection)

- suppress other boxes with significant overlap with the max-box (high IoU)

in the image on left - only two boxes 0.8 and 0.9 will be highlighted

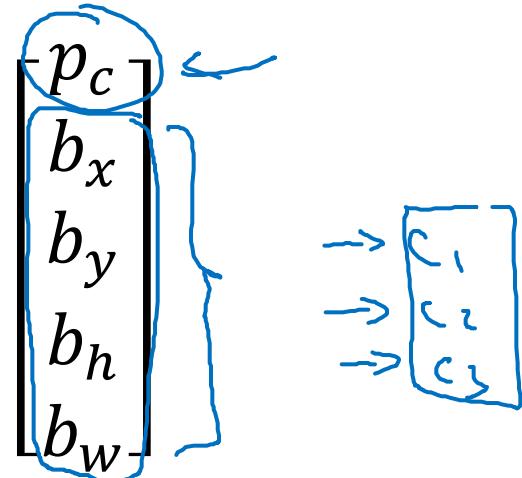
# Non-max suppression algorithm



19×19

Each output prediction is:

361 outputs



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



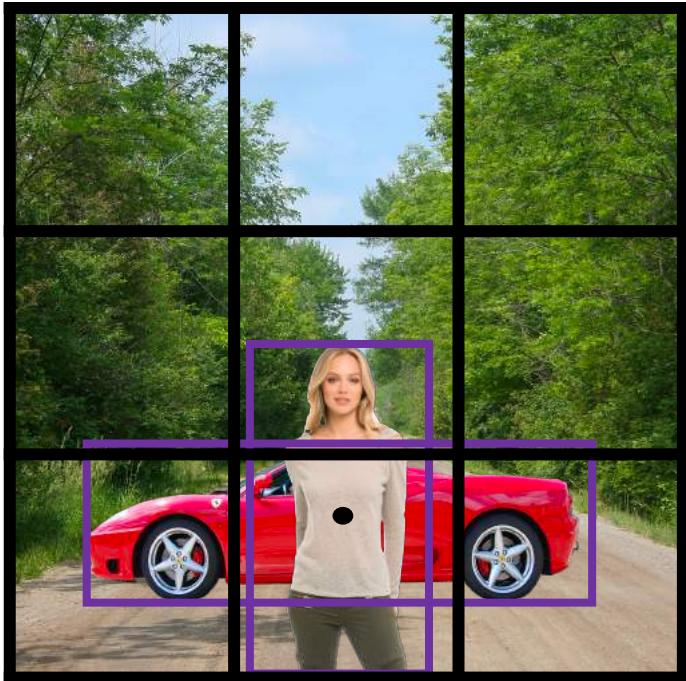
deeplearning.ai

# Object Detection

---

## Anchor boxes

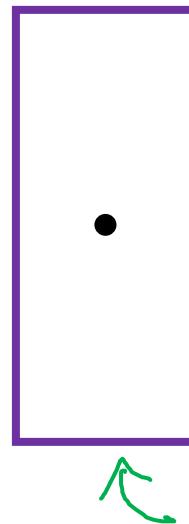
# Overlapping objects:



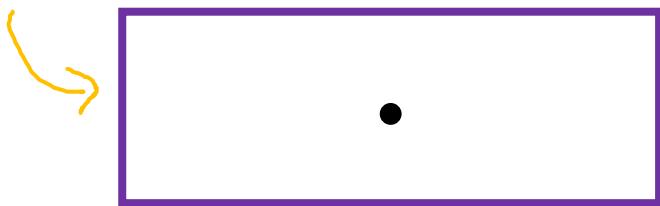
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

ISSUE:  
mid-point of both objects fall in  
same grid - implying we'll know  
there is a landmark object here,  
but only be able to pick one

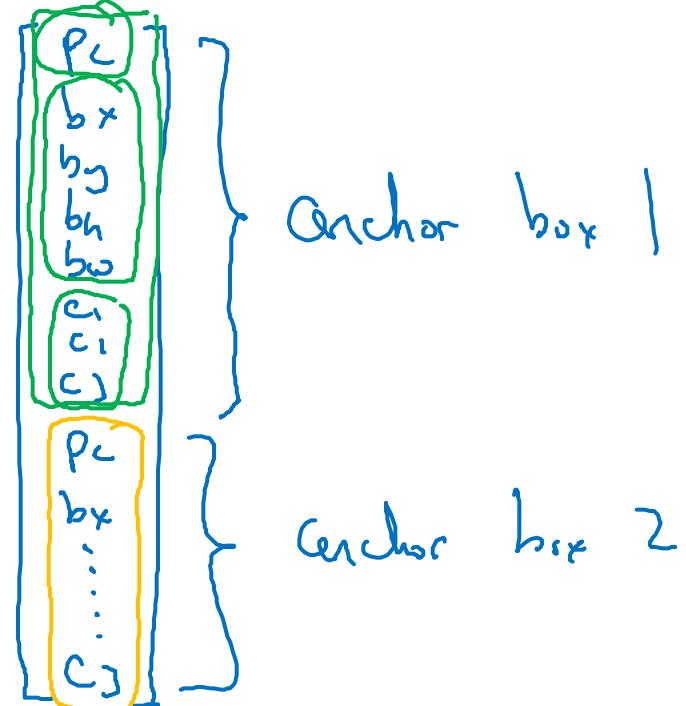
Anchor box 1:



Anchor box 2:



$$y = \begin{array}{l} 8 \times 2 \text{ o/p vector} \\ \vdots \end{array}$$



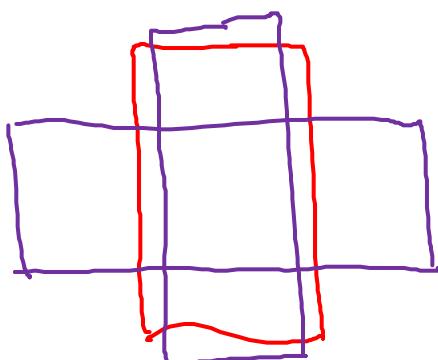
# Anchor box algorithm

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.

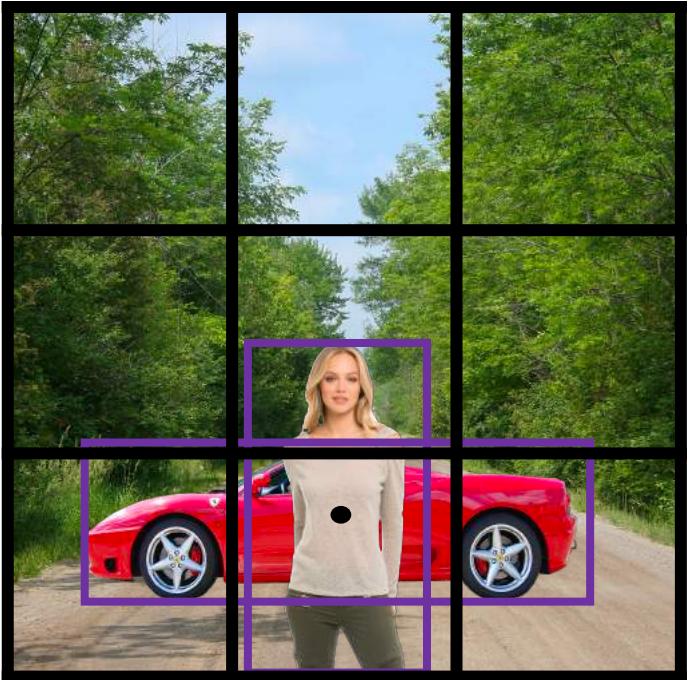
(grid cell, anchor box)

Output y:

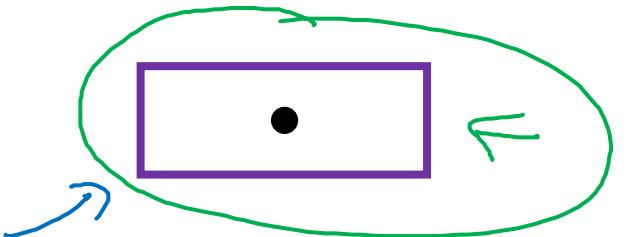
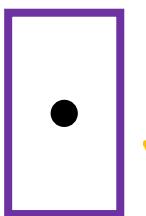
$3 \times 3 \times 16$

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

# Anchor box example



Anchor box 1:      Anchor box 2:



What it won't capture:  
(a) 3 objects in same grid  
(b) 2 objects with same anchor box

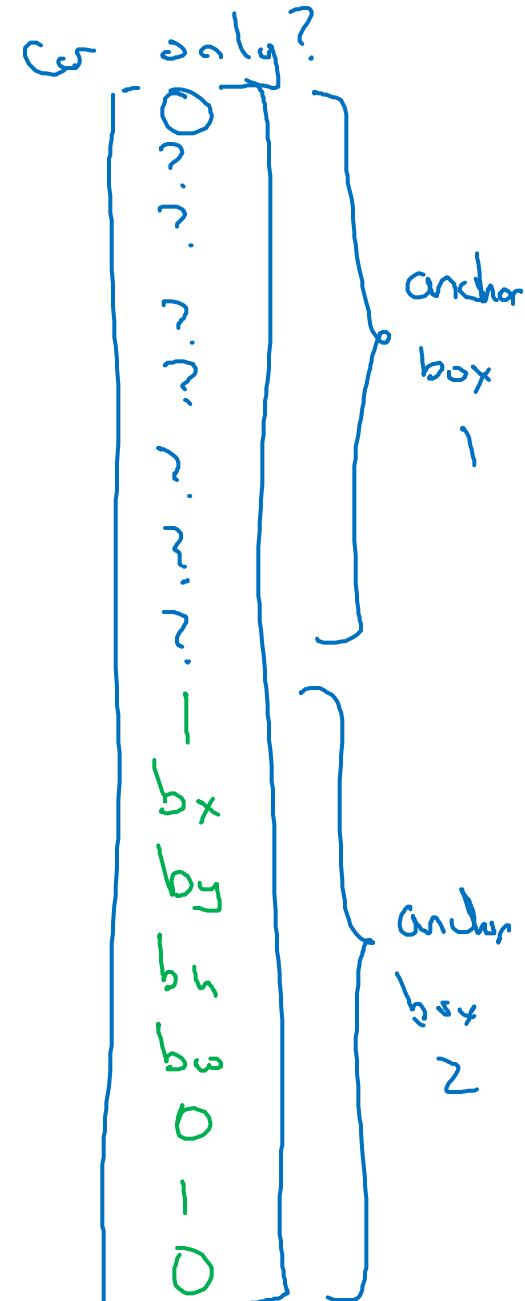
Chances of such collision  
are reduced as we make  
the grid finer

$$y =$$

In advanced  
research - people use  
more than 2 anchor  
boxes, or use KNN  
algorithm

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

$$\begin{bmatrix} 1 \\ bx \\ by \\ bh \\ bw \\ 1 \\ 0 \\ 0 \\ 1 \\ bx \\ by \\ bh \\ bw \\ 0 \\ 1 \\ 0 \end{bmatrix}$$





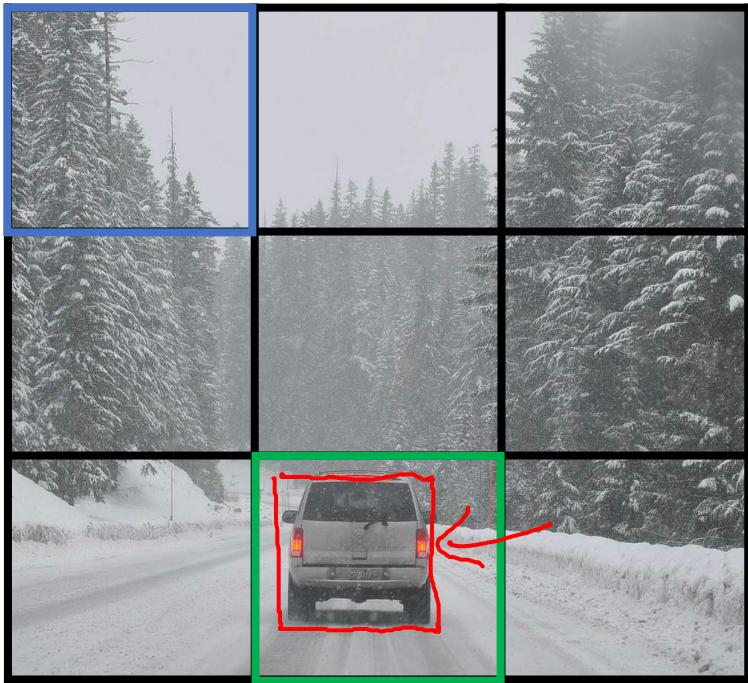
deeplearning.ai

# Object Detection

---

Putting it together:  
YOLO algorithm

# Training



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

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

$3 \times 3 \times 16$

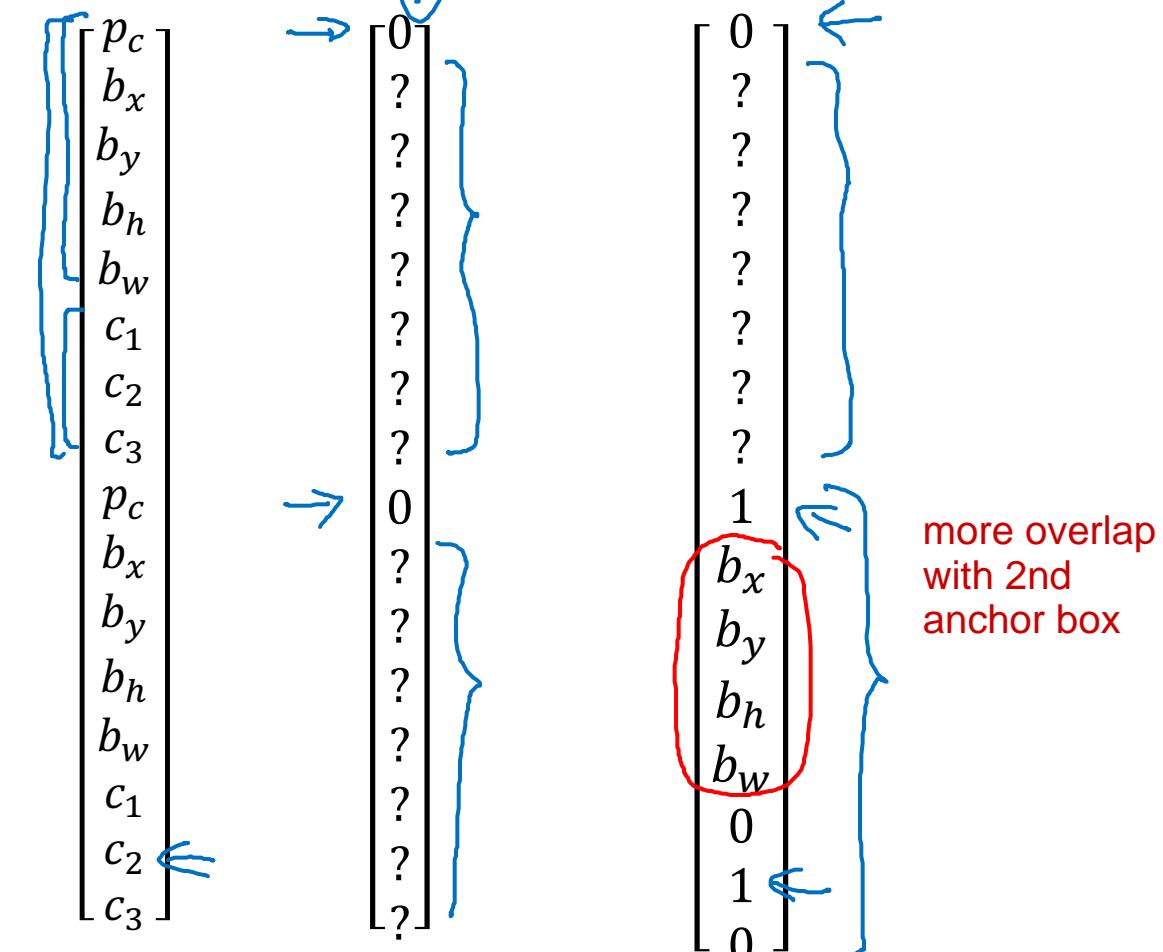
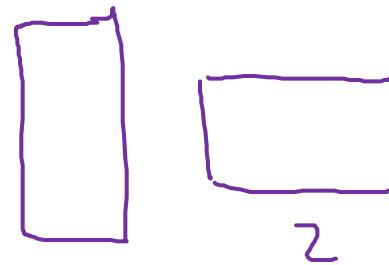
$\uparrow$   
#anchors

$\uparrow$   
5 + #classes

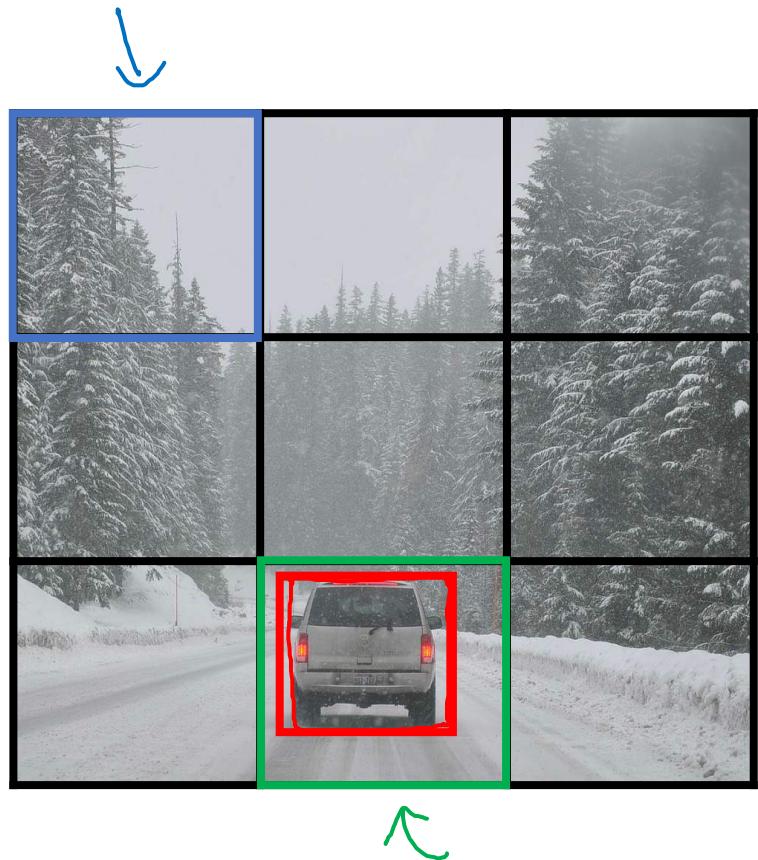
training set for each image

- 1 - pedestrian
- 2 - car
- 3 - motorcycle

$y =$



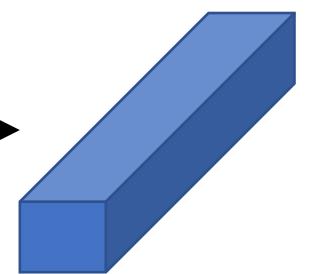
# Making predictions



→

...

→



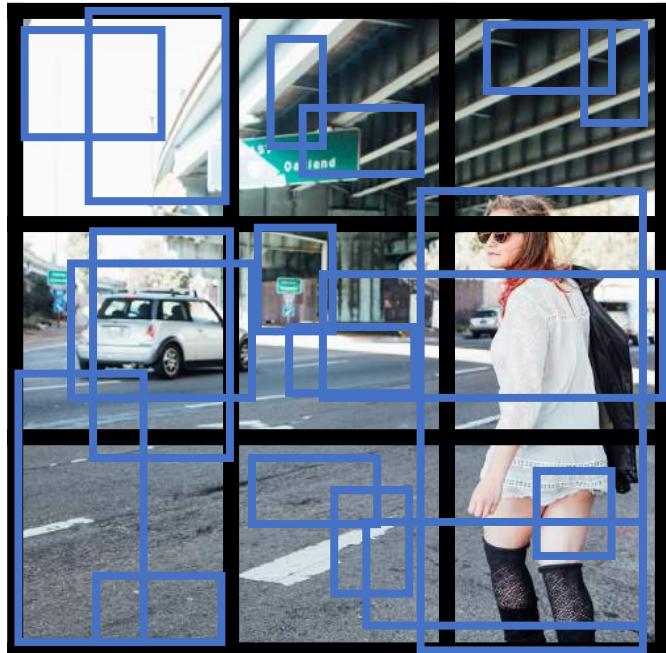
$y =$

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

Annotations for the output vector:

- Blue arrows point to the first four elements:  $p_c$ ,  $b_x$ ,  $b_y$ , and  $b_h$ .
- Red annotations point to the last four elements:  $b_w$ ,  $c_1$ ,  $c_2$ , and  $c_3$ .
- A blue arrow points to the bottom element:  $c_3$ .

# Outputting the non-max suppressed outputs



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



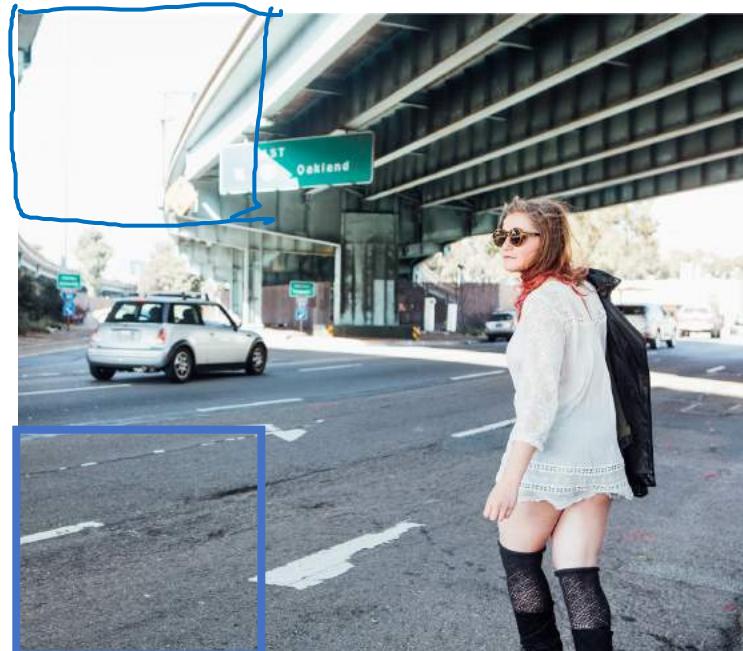
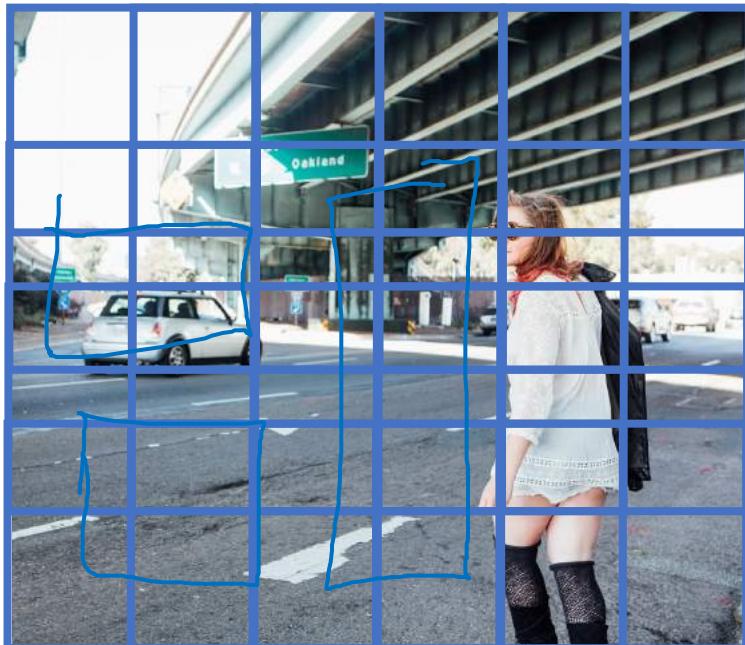
deeplearning.ai

# Object Detection

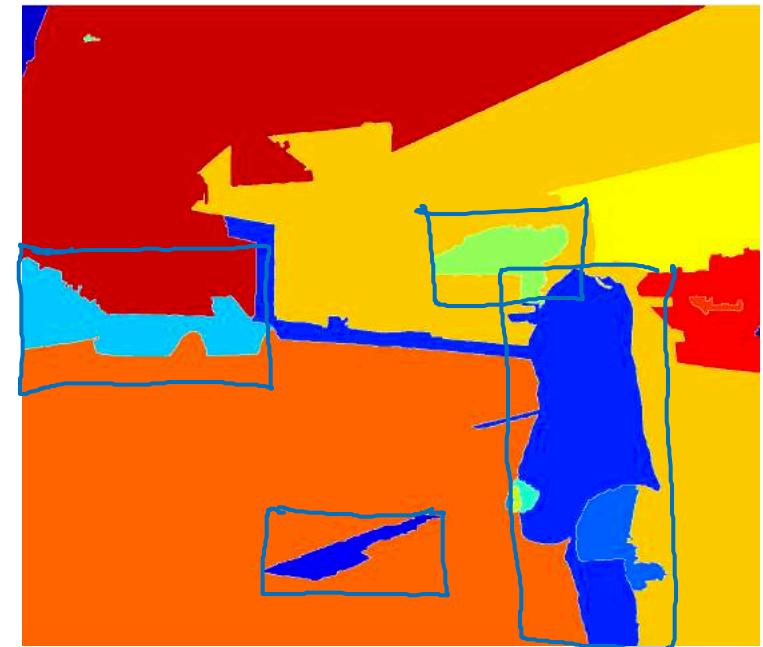
---

Region proposals  
(Optional)

# Region proposal: R-CNN



↖



Segmentation algorithm

~ 2,000

Instead of running convolution on full image, run a segmentation algorithm to identify interesting regions in the image, and then run convolution on boxes around those regions (~ find maybe 2000 blobs and run convolution on those)

# Faster algorithms

Generally, these approaches are still slower than YOLO

→ R-CNN:

this approach is still slow

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

it predicts the bounding box, and doesn't just trust the bounding box it was given from segmentation

Fast R-CNN:

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

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

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

[Girshik, 2015. Fast R-CNN]

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

Andrew Ng



deeplearning.ai

# Convolutional Neural Networks

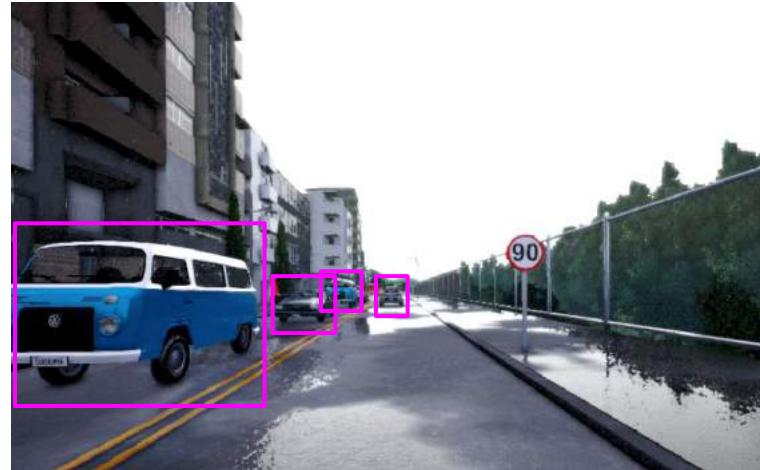
---

## Semantic segmentation with U-Net

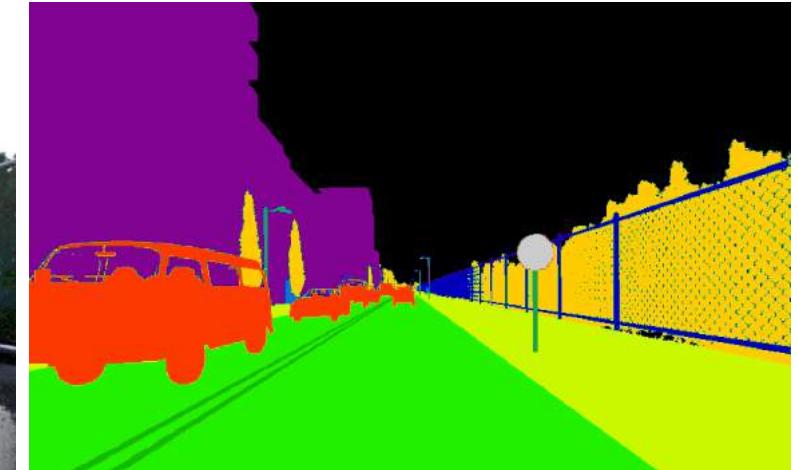
# Object Detection vs. Semantic Segmentation



Input image



Object Detection

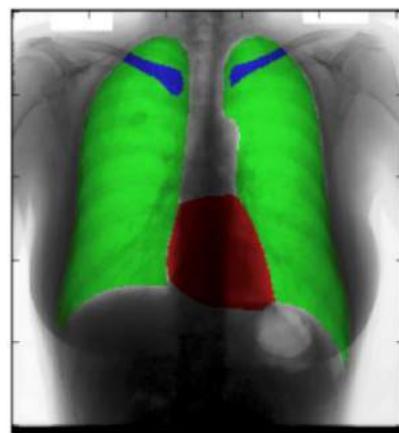


Semantic Segmentation

it attempts to segment each pixel - for example, what pixels represent drivable road

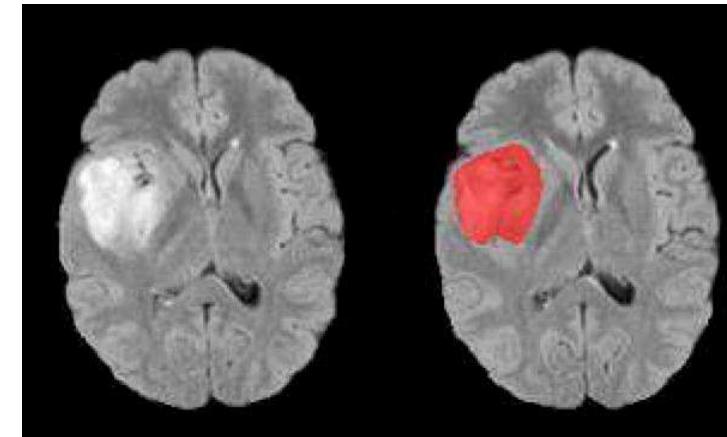
# Motivation for U-Net

algorithm highlights different parts - heart, lungs, clavicle



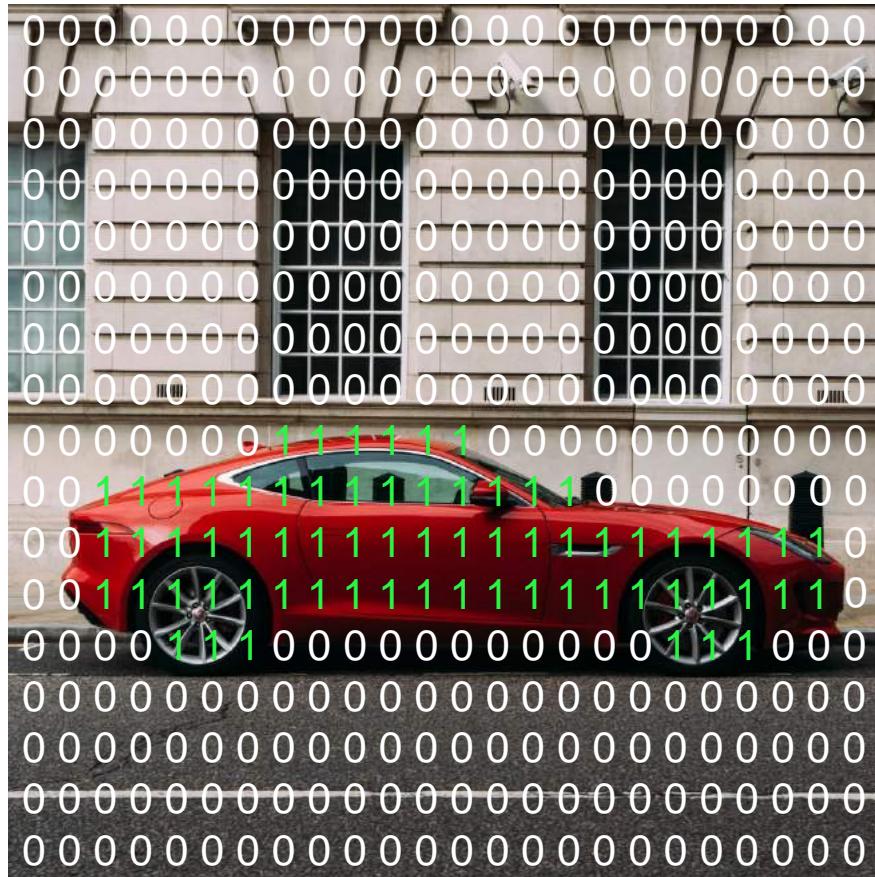
Chest X-Ray

the algorithm highlights the area of tumor - which is useful for surgery planning



Brain MRI

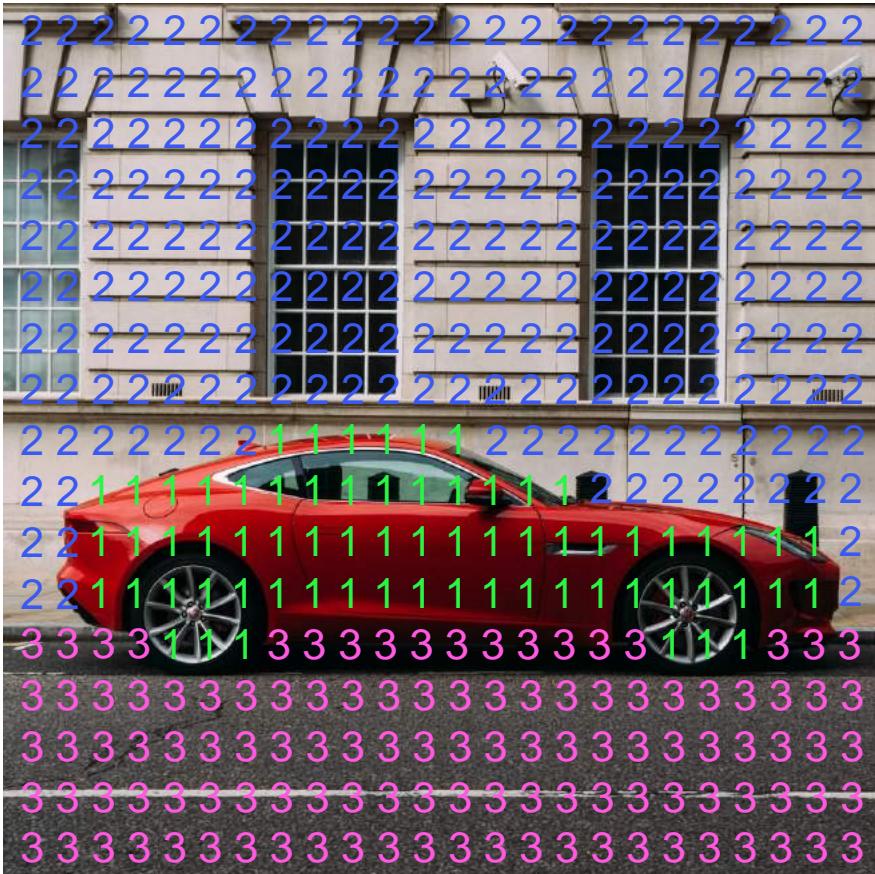
# Per-pixel class labels



1. Car  
0. Not Car

# Per-pixel class labels

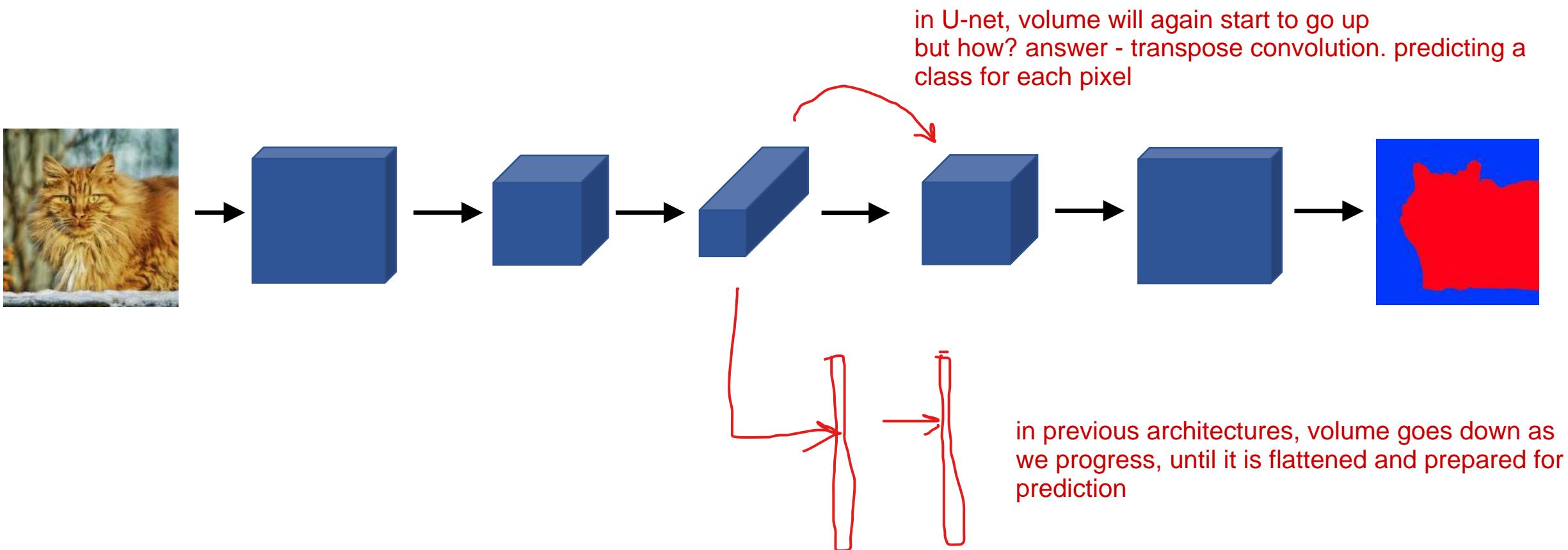
Much more complex output in this case  
- not just giving object classification (0/1)  
- or, just classification & bounding box  
but a whole matrix of output



1. Car
  2. Building
  3. Road

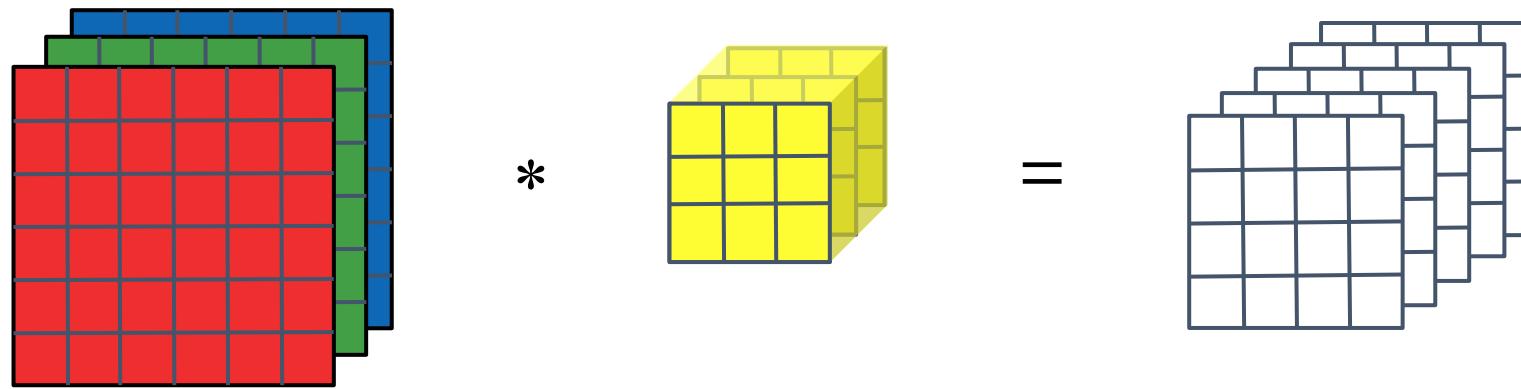
## Segmentation Map

# Deep Learning for Semantic Segmentation

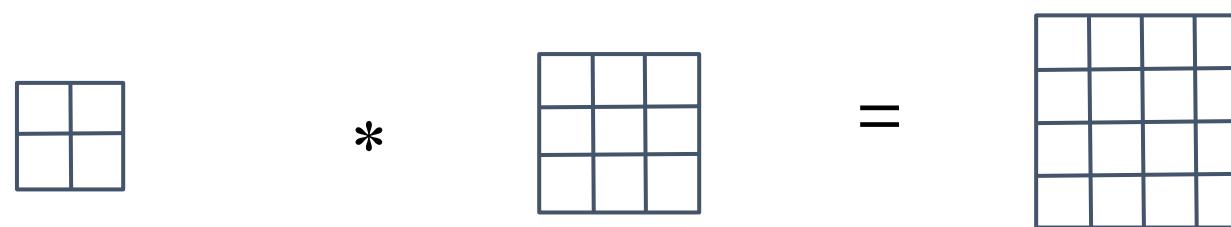


# Transpose Convolution

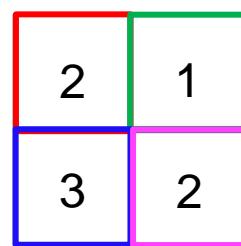
Normal Convolution



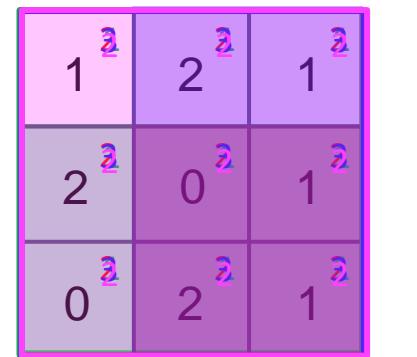
Transpose Convolution



# Transpose Convolution



$2 \times 2$



weight filter

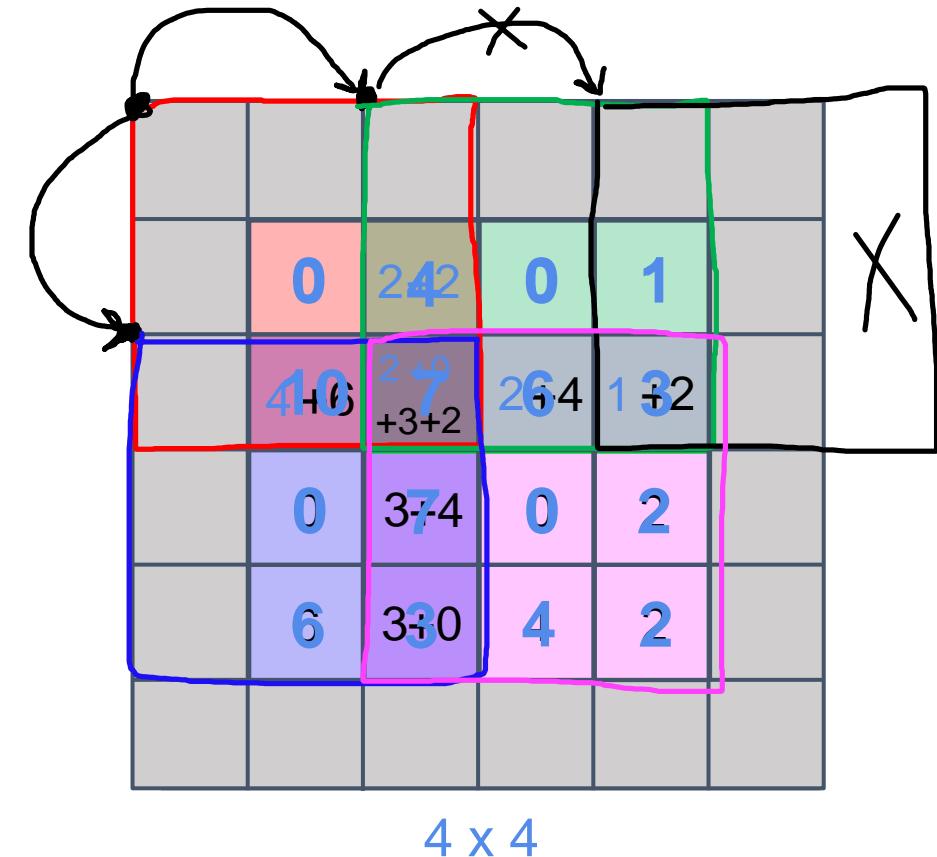
output:  
0 4 0 1  
10 7 6 3  
0 7 0 2  
6 3 4 2

filter  $f \times f = 3 \times 3$

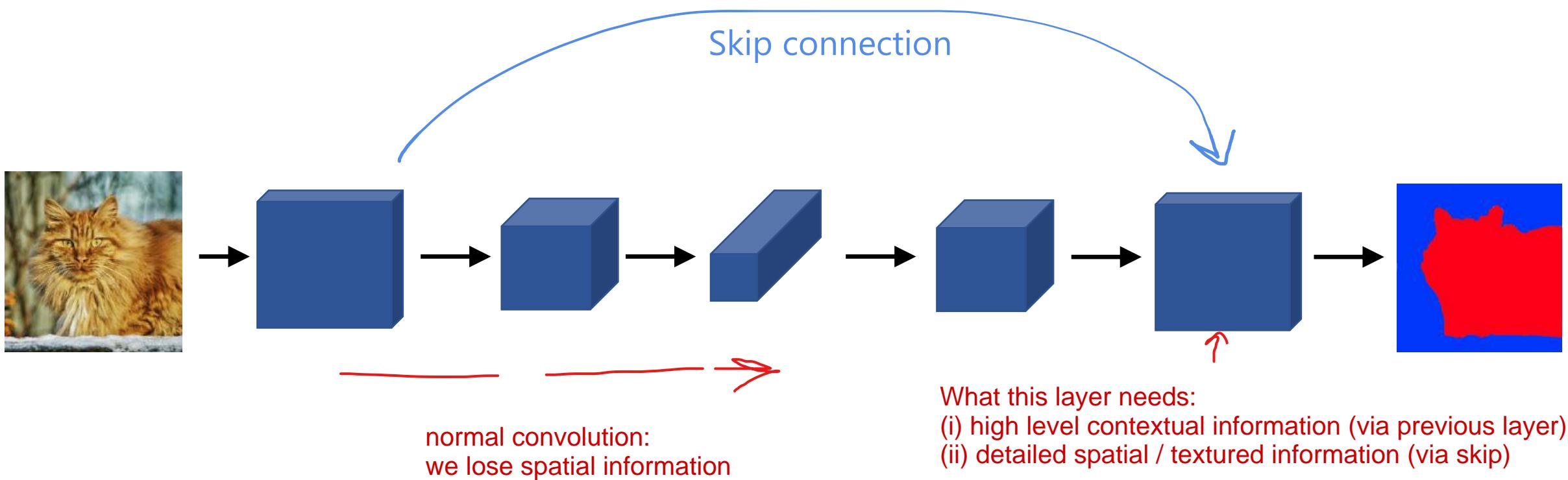
padding  $p = 1$

stride  $s = 2$

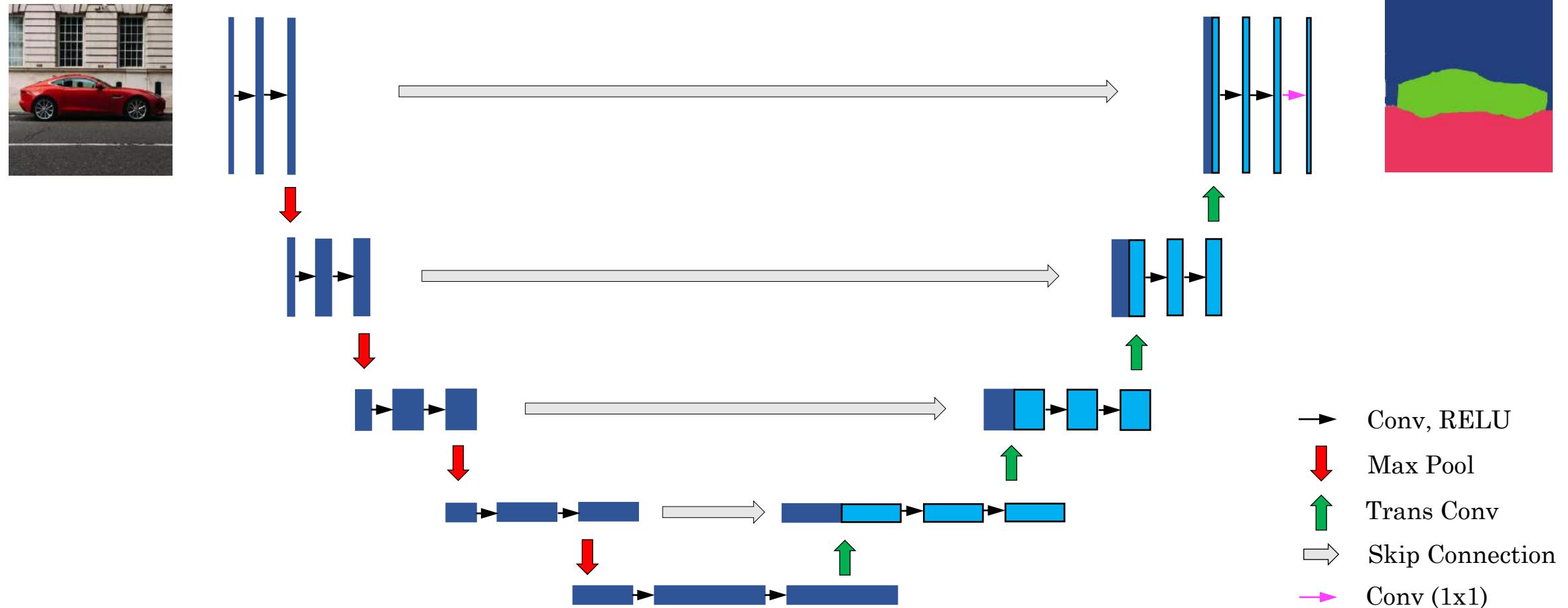
Multiple first value in A (=2) with Filter (B) and fill C (non-padded region)  
- then repeat with A (=1), filling C after stride 2. Add values in overlapping boxes



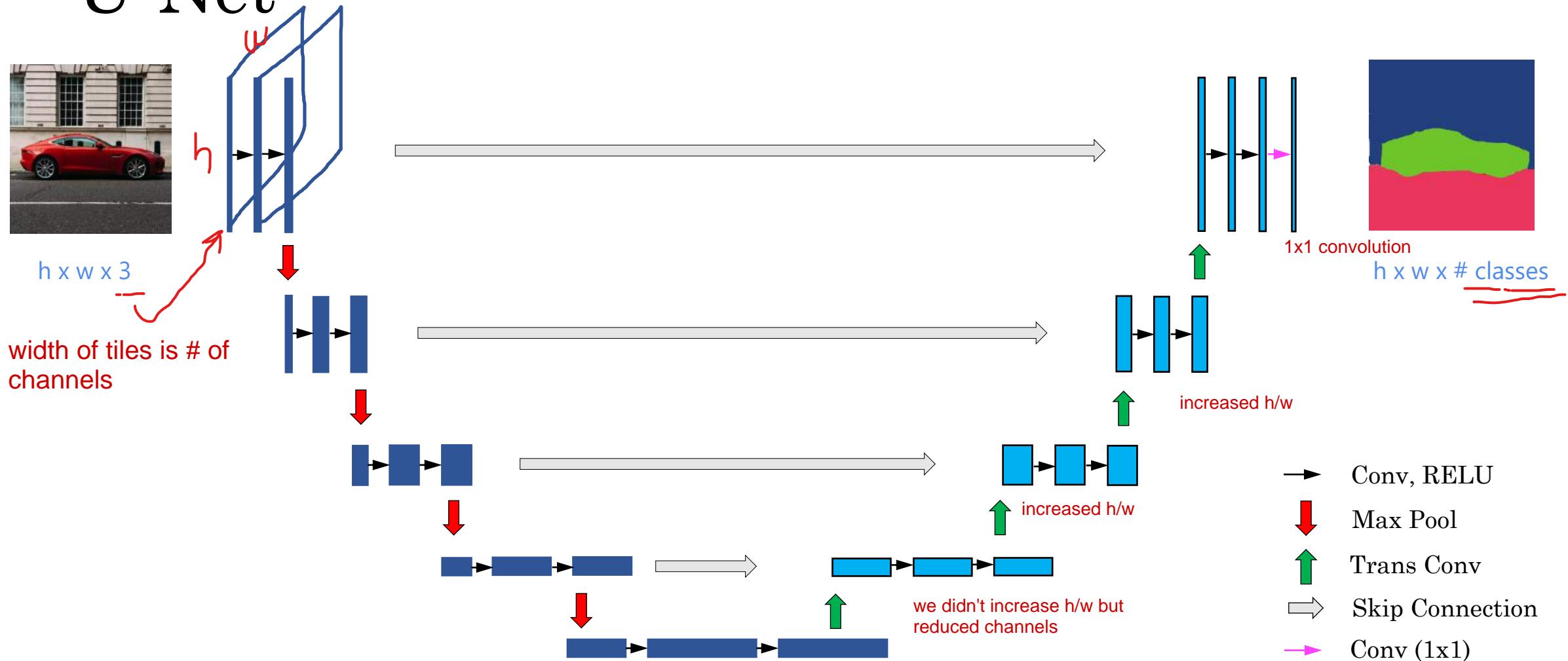
# Deep Learning for Semantic Segmentation



# U-Net



# U-Net



# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



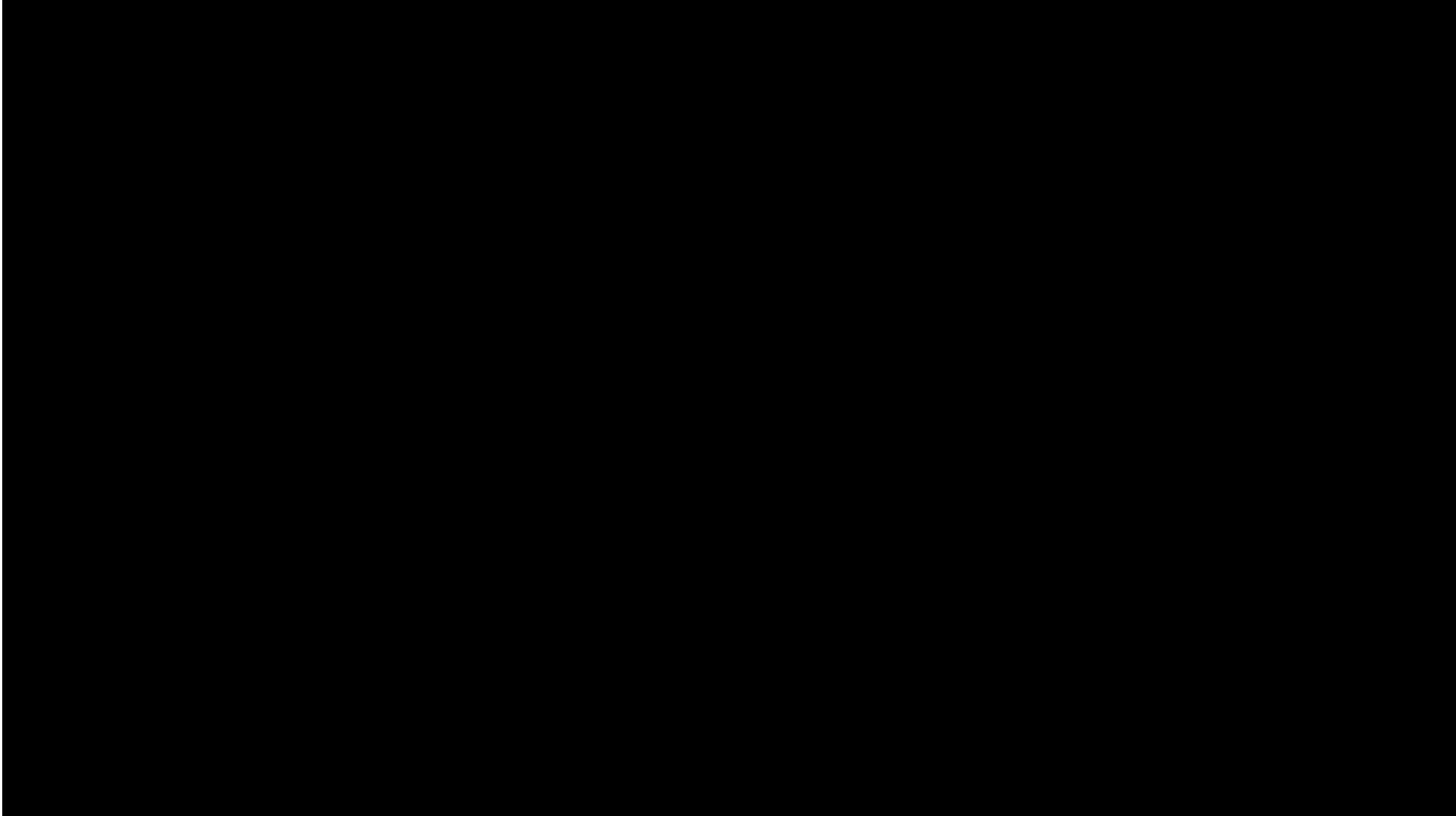
deeplearning.ai

# Face recognition

---

What is face  
recognition?

# 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 ←



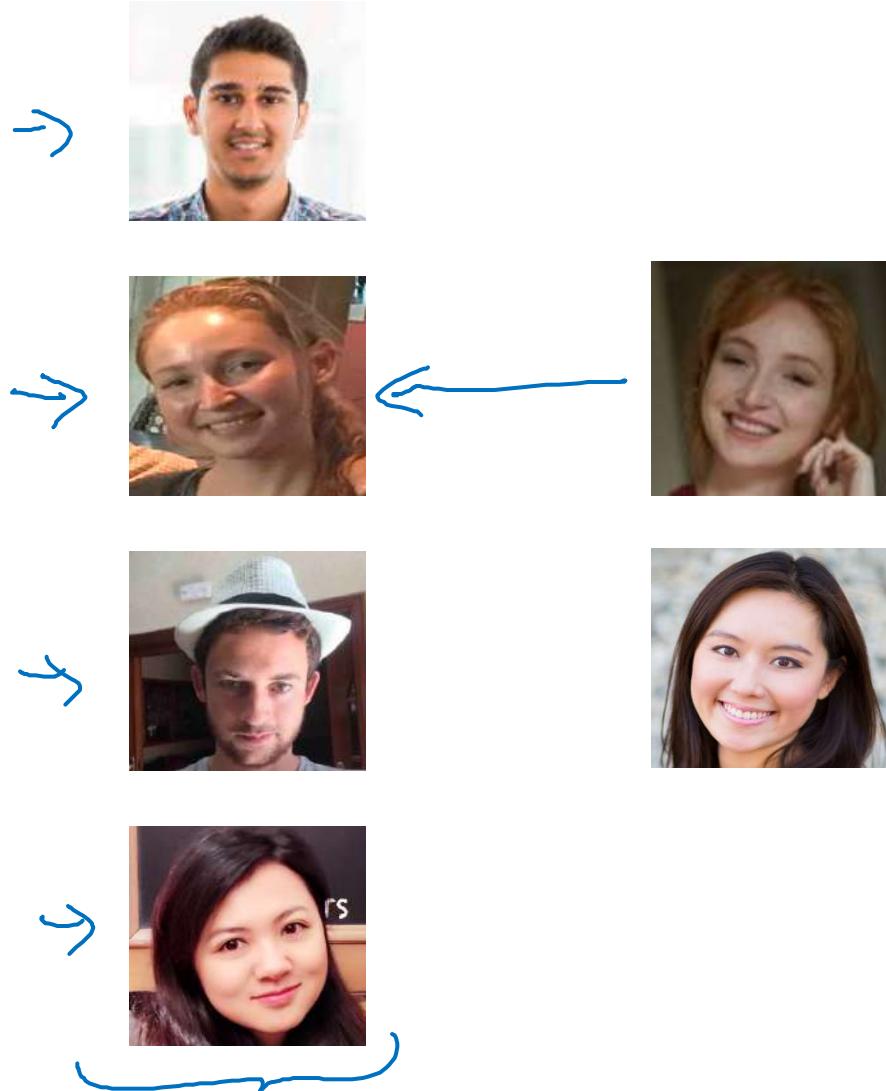
deeplearning.ai

# Face recognition

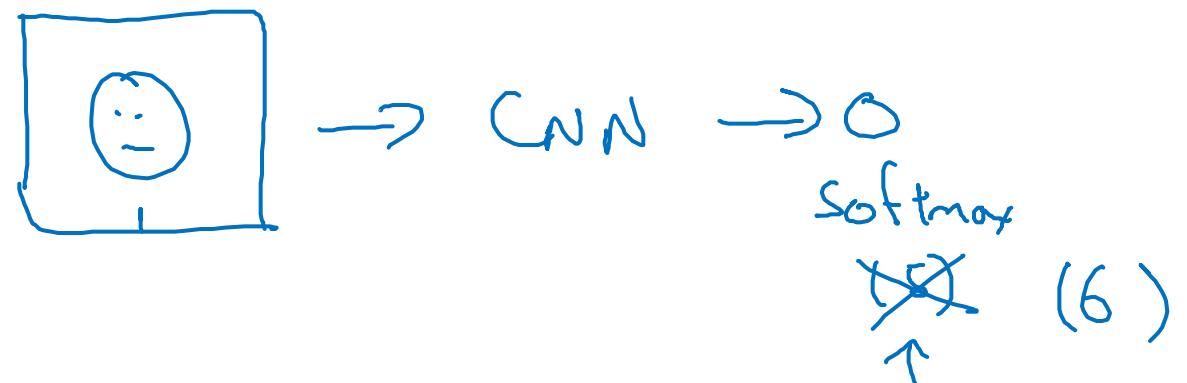
---

# One-shot learning

# One-shot learning



Learning from one example to recognize the person again



not a good approach - will need to expand last layer and retrain network every time someone new joins the team

# Learning a “similarity” function

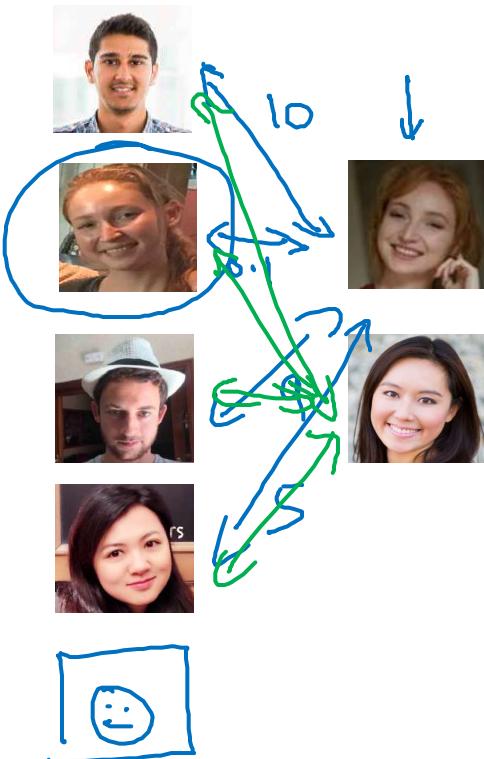
→  $d(\underline{\text{img1}}, \underline{\text{img2}})$  = degree of difference between images

If  $d(\text{img1}, \text{img2}) \leq \tau$

$> \tau$

“some”  
“different”

Verification.



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

Key idea - learn a "similarity function"

How to train a NN to learn this function "d" - use a Siamese network  
Training images for Siamese network can be a different set of images (with >1 for each person)



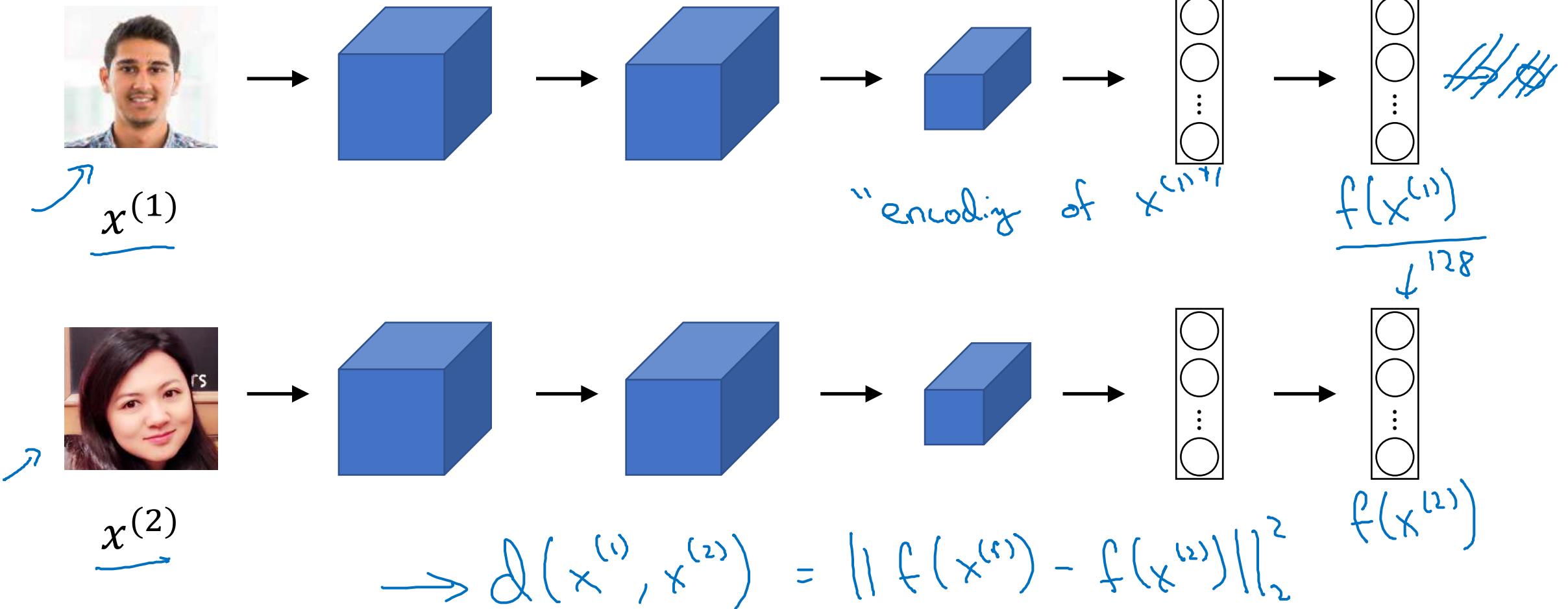
deeplearning.ai

# Face recognition

---

## Siamese network

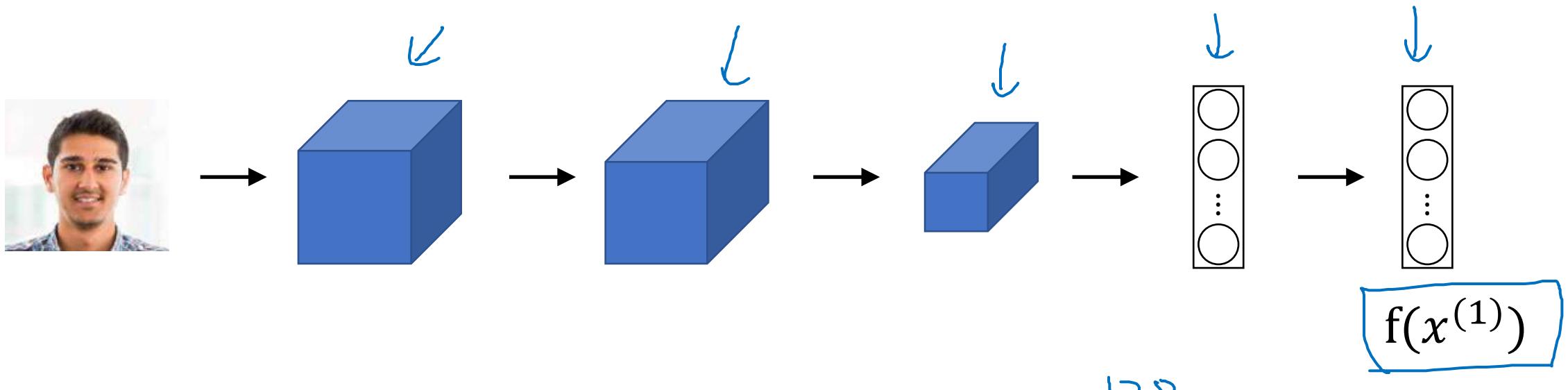
# Siamese network



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

Andrew Ng

# 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.



deeplearning.ai

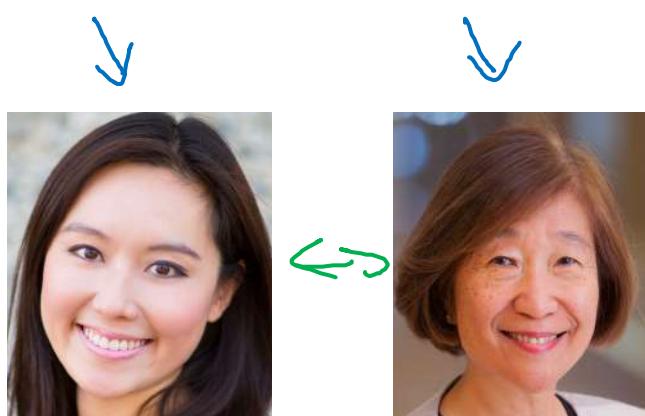
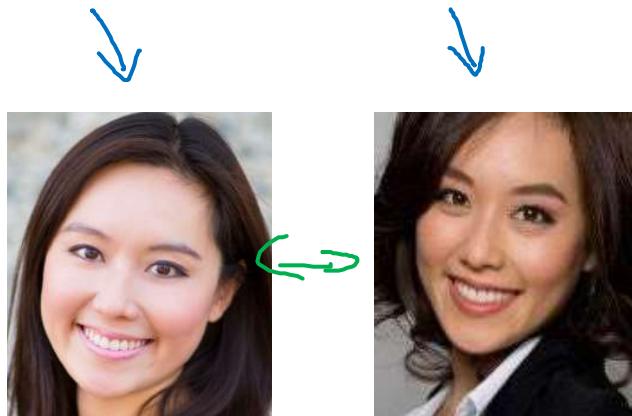
# Face recognition

---

## Triplet loss

# Learning Objective

Triplet loss - you are looking at 3 images at each alpha (margin) is added to ensure 'trivial' solution of 0 is eliminated



Anchor	Positive		Anchor	Negative
<u>A</u>	<u>P</u>	$\frac{d(A, P)}{d(A, P)} = 0.5$	<u>A</u>	<u>N</u>
$\Rightarrow 0.2$			$\frac{d(A, N)}{d(A, N)} = \cancel{0.5} \quad 0.7$	
Want: $\frac{\ f(A) - f(P)\ ^2}{d(A, P)} + \alpha \leq \frac{\ f(A) - f(N)\ ^2}{d(A, N)}$				
$\frac{\ f(A) - f(P)\ ^2}{0} - \frac{\ f(A) - f(N)\ ^2}{0} + \alpha \leq 0 \quad \text{Margin}$				

# Loss function

- Loss is either 'positive' or 'zero'
- if it is negative, we've reached our objective and no longer want to optimize

Given 3 images

$A, P, N$ :

$$\underline{L(A, P, N)} = \max \left( \boxed{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \delta} \right), 0 \right)$$

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

$A, P$   
 $T$

Training set:  $\underbrace{10k}_{\infty}$  pictures of  $\frac{1k}{\infty}$  persons

you need  $>1$  picture for each person for training

# 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.

$$\underbrace{\|f(A) - f(P)\|^2}_{\text{Distance between } A \text{ and } P} + \alpha \leq \underbrace{\|f(A) - f(N)\|^2}_{\text{Distance between } A \text{ and } N}$$

Choose triplets that're “hard” to train on.

$$\begin{aligned} \cancel{d(A, P)} + \alpha &\leq d(A, N) \\ \frac{d(A, P)}{\downarrow} &\approx \frac{d(A, N)}{\uparrow} \end{aligned}$$

Face Net  
Deep Face



# Training set using triplet loss

Anchor



Positive



Negative



:

:

:



J

$$d(x^{(i)}, x^{(j)})$$



deeplearning.ai

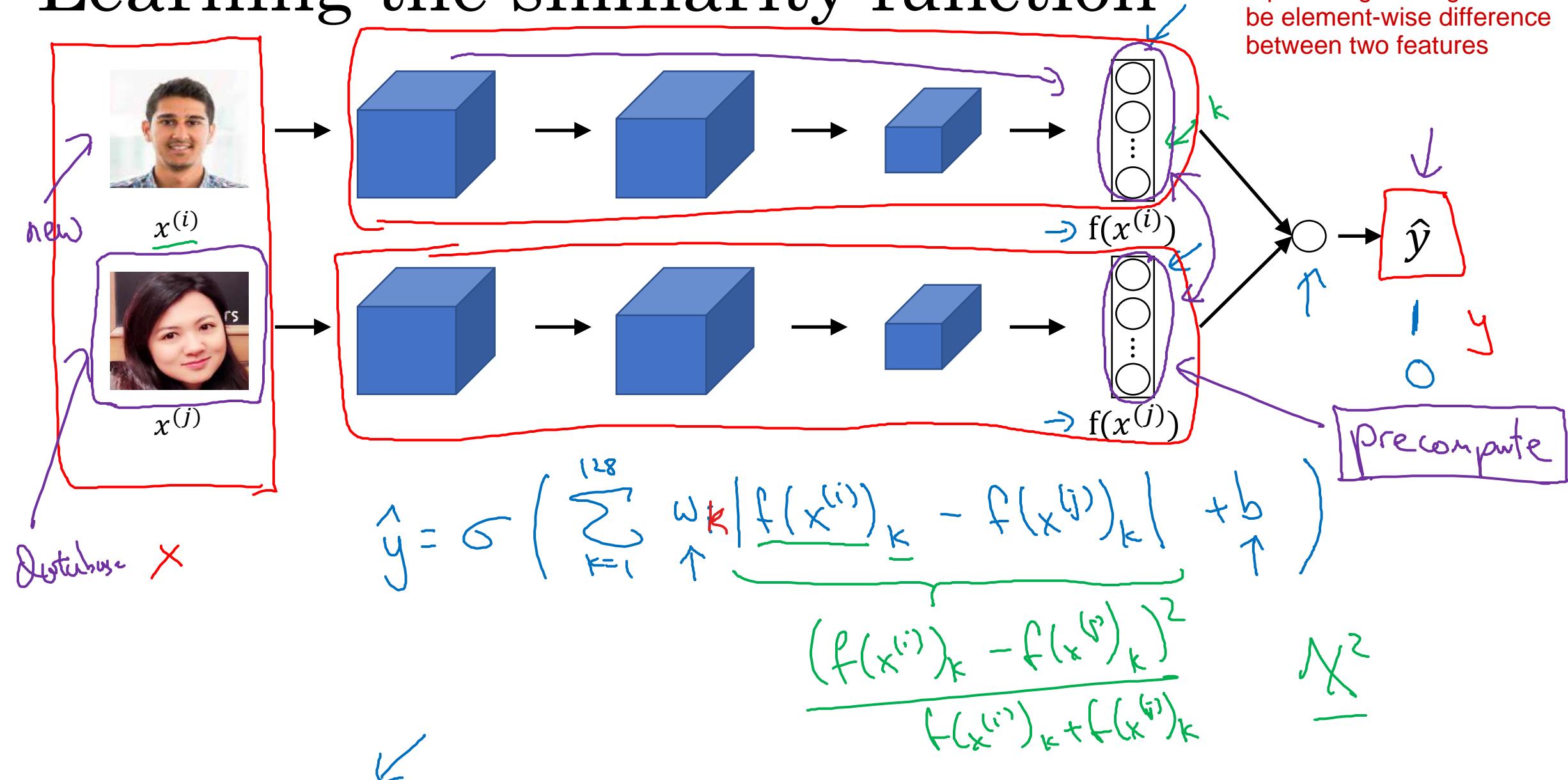
# Face recognition

---

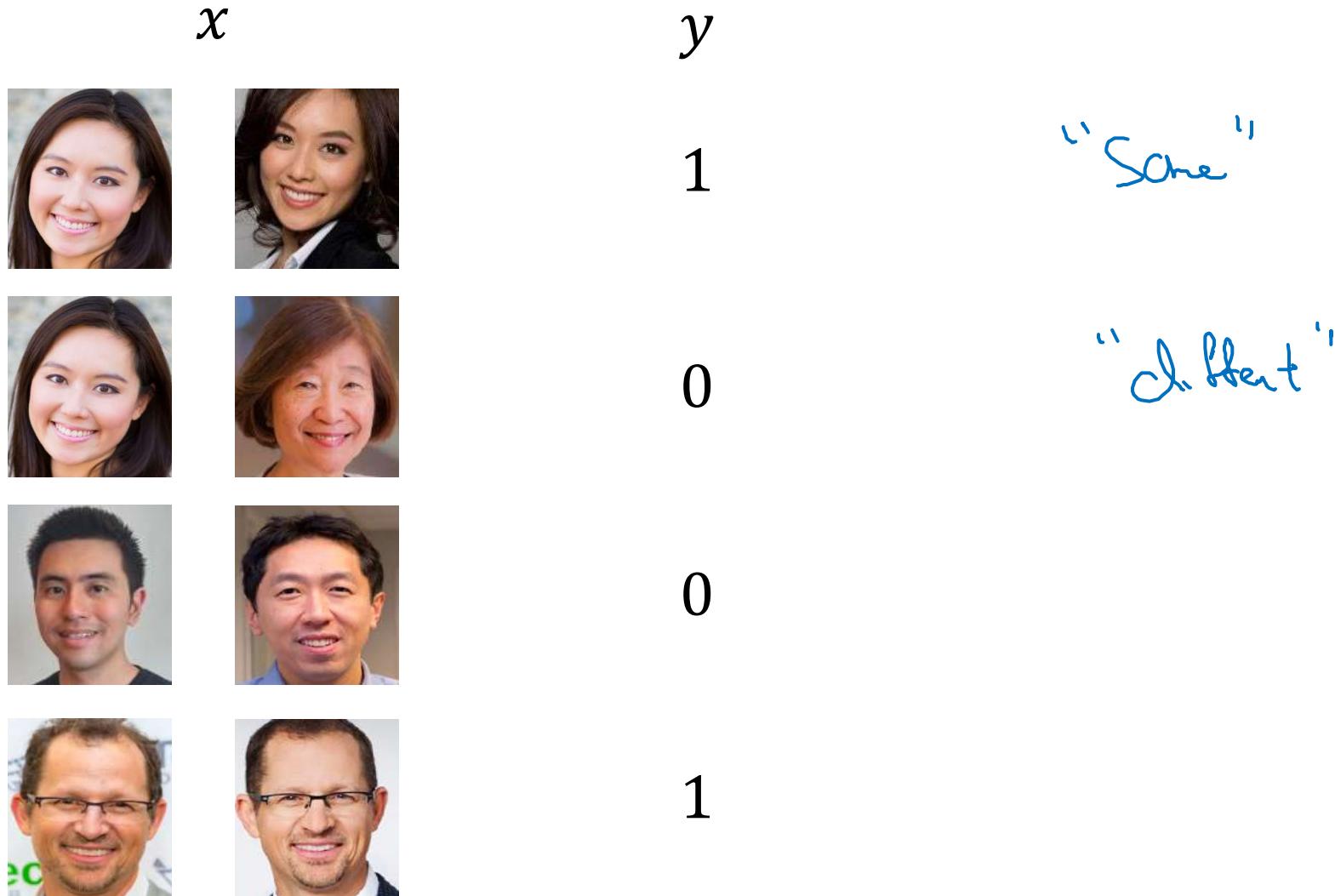
# Face verification and binary classification

# Learning the similarity function

Input to logistic regression can be element-wise difference between two features



# Face verification supervised learning





deeplearning.ai

# Neural Style Transfer

---

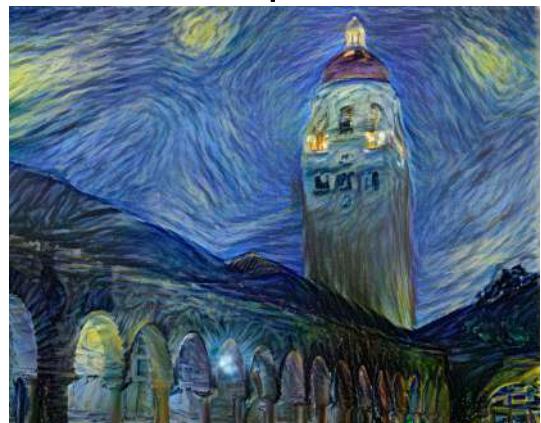
What is neural style  
transfer?

# Neural style transfer



Content ( $c$ )

Style ( $s$ )



Generated image ( $g$ )



Content ( $c$ )

Style ( $s$ )



Generated image ( $g$ )



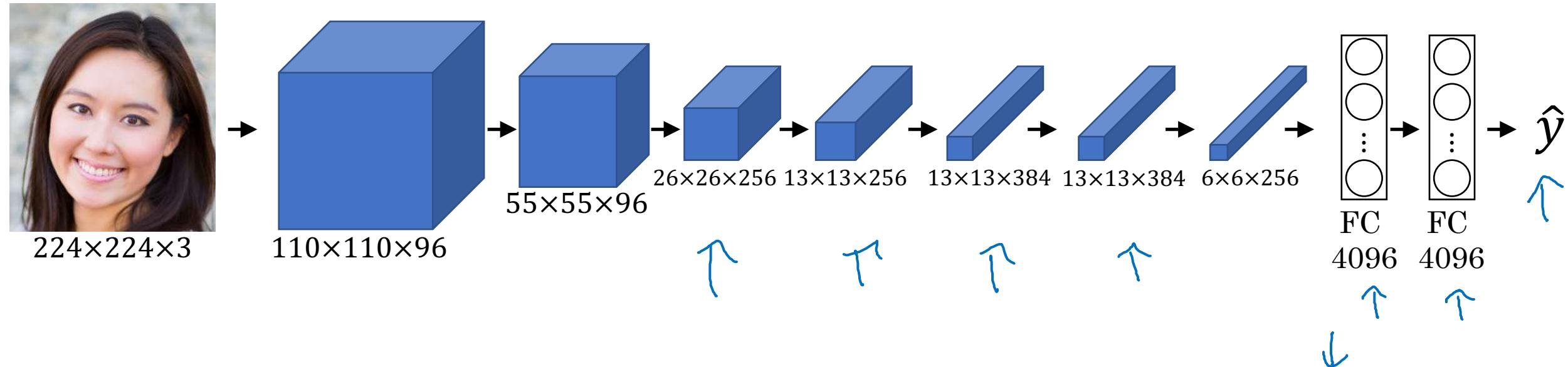
deeplearning.ai

# Neural Style Transfer

---

What are deep  
ConvNets learning?

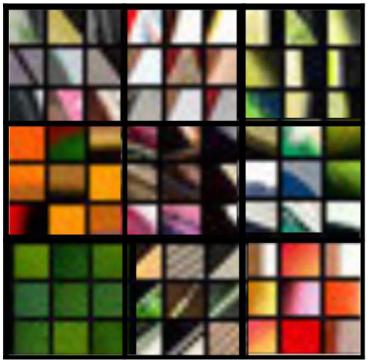
# Visualizing what a deep network is learning



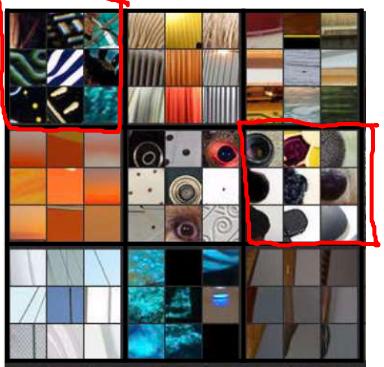
Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

Repeat for other units.

# Visualizing deep layers



Layer 1



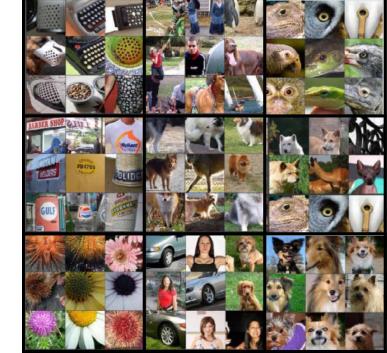
Layer 2



Layer 3

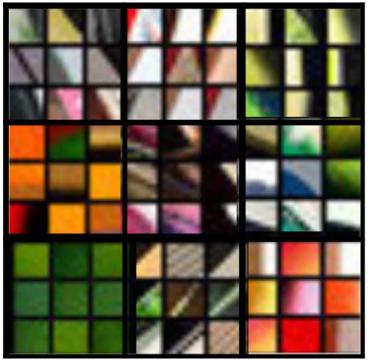


Layer 4



Layer 5

# Visualizing deep layers: Layer 1



Layer 1



Layer 2



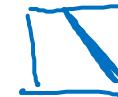
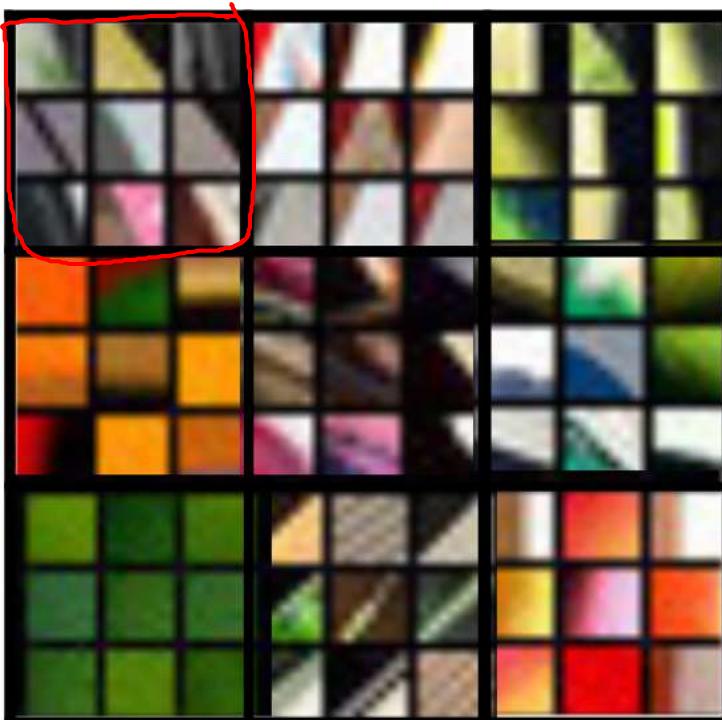
Layer 3



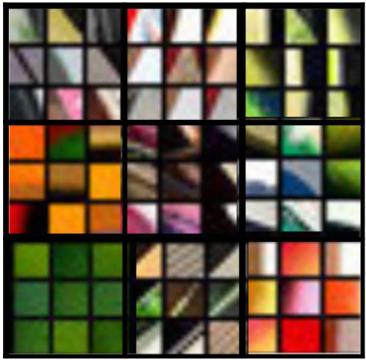
Layer 4



Layer 5



# Visualizing deep layers: Layer 2



Layer 1



Layer 2



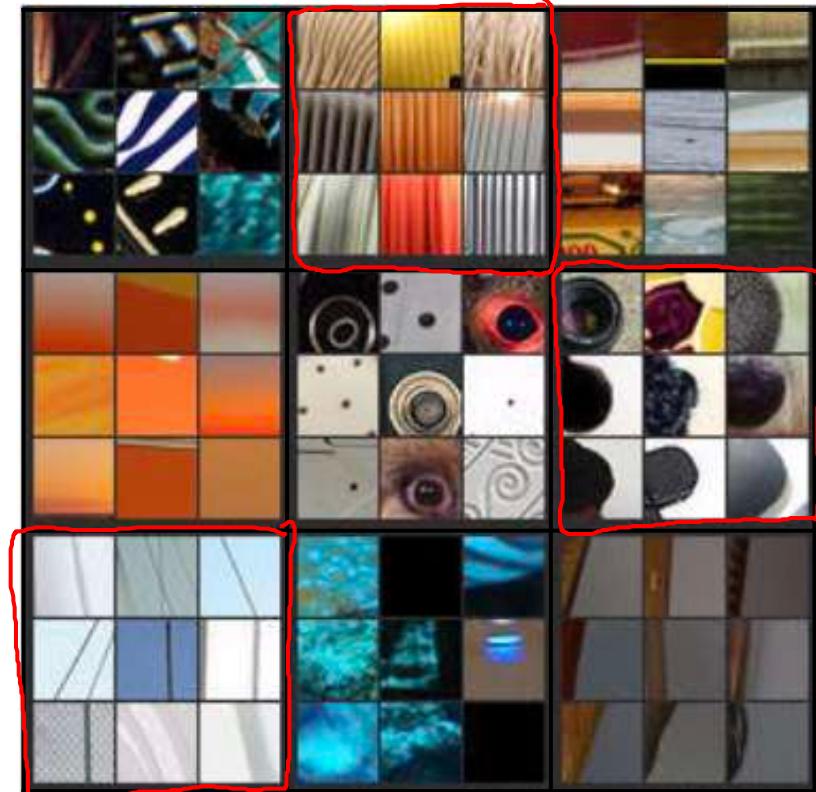
Layer 3



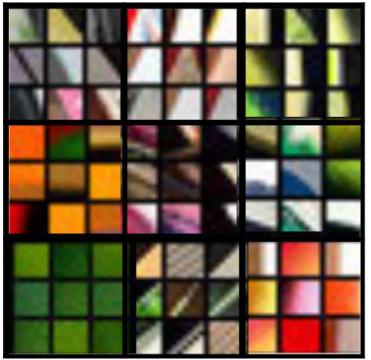
Layer 4



Layer 5



# Visualizing deep layers: Layer 3



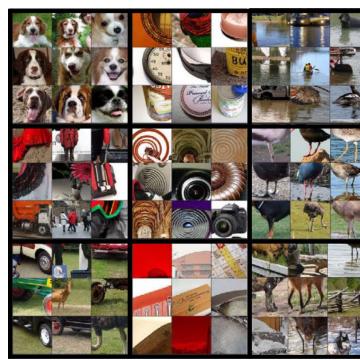
Layer 1



Layer 2



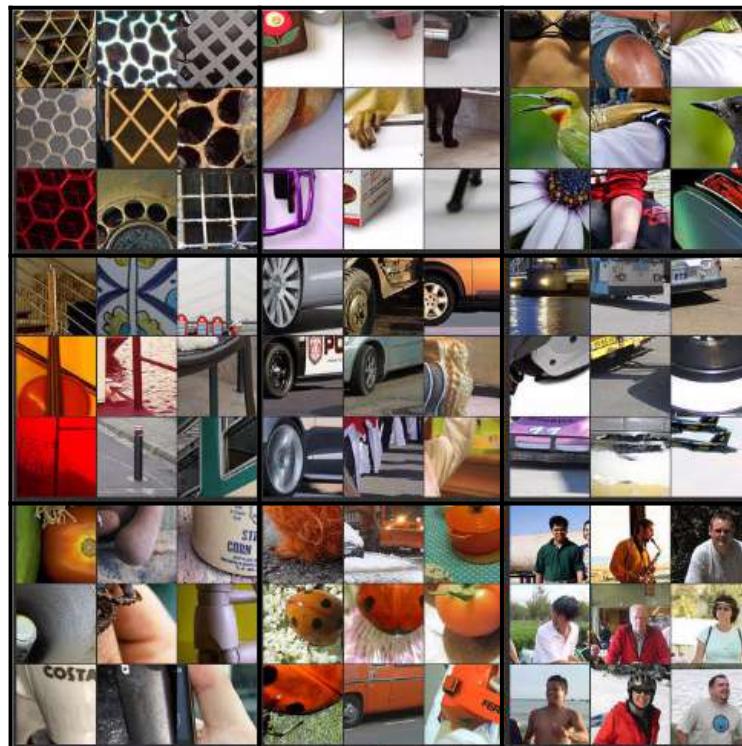
Layer 3



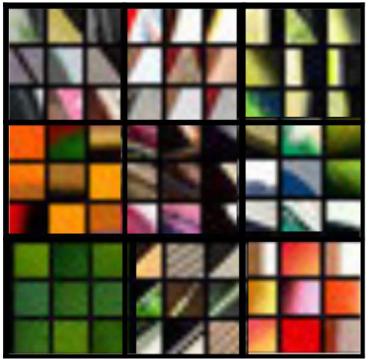
Layer 4



Layer 5



# Visualizing deep layers: Layer 3

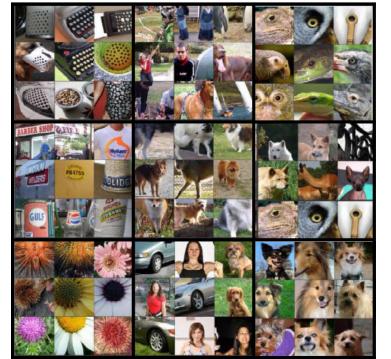


Layer 1



L

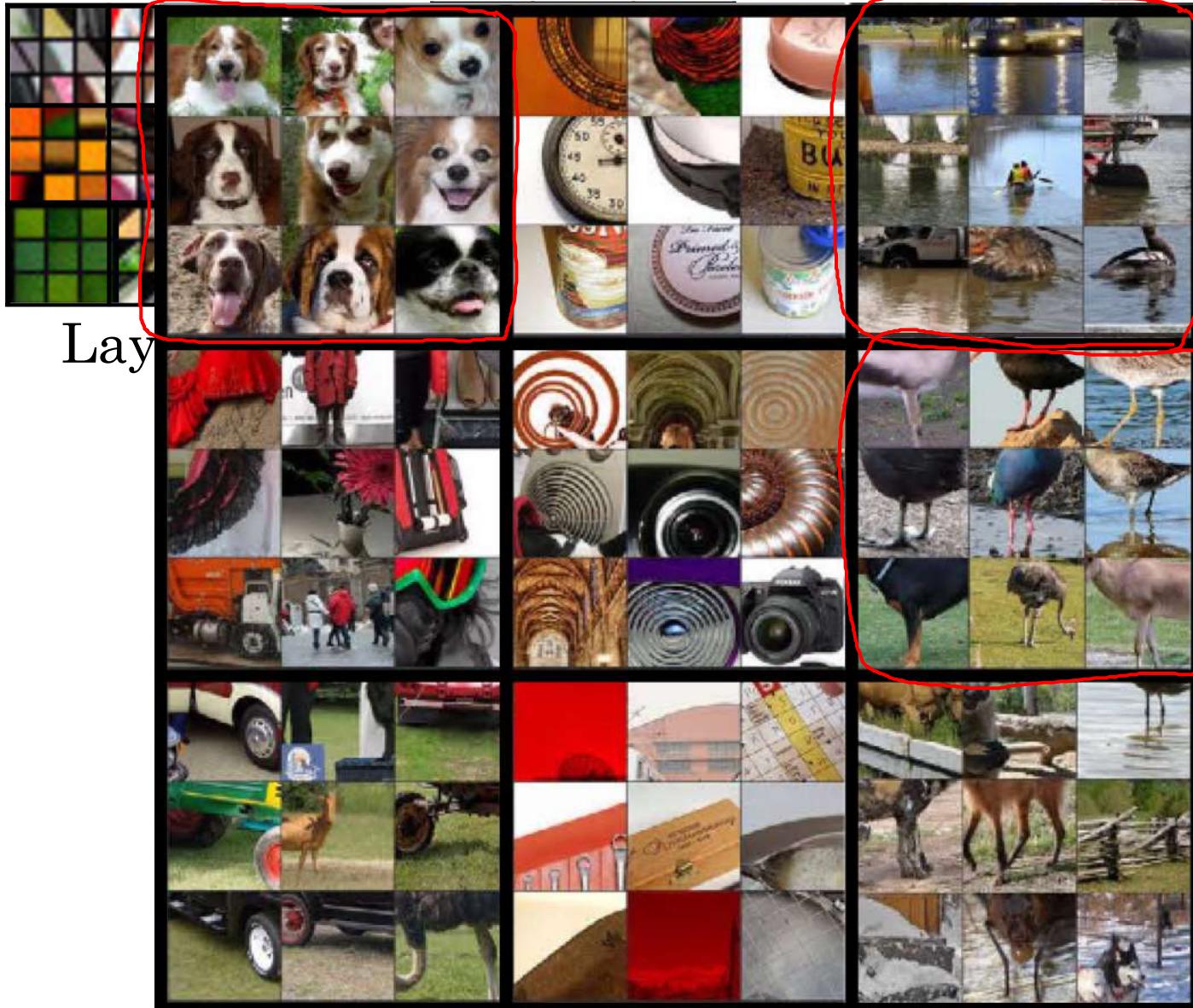
4



Layer 5



# Visualizing deep layers: Layer 4



Layer 4

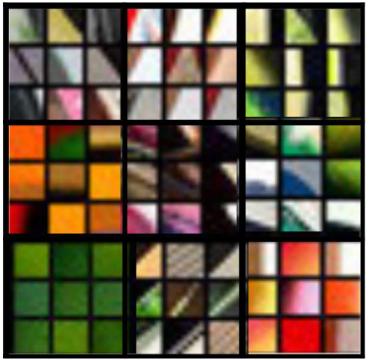


Layer 4



Layer 5

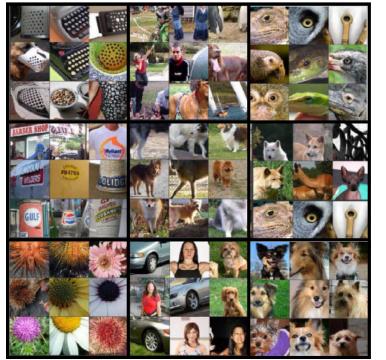
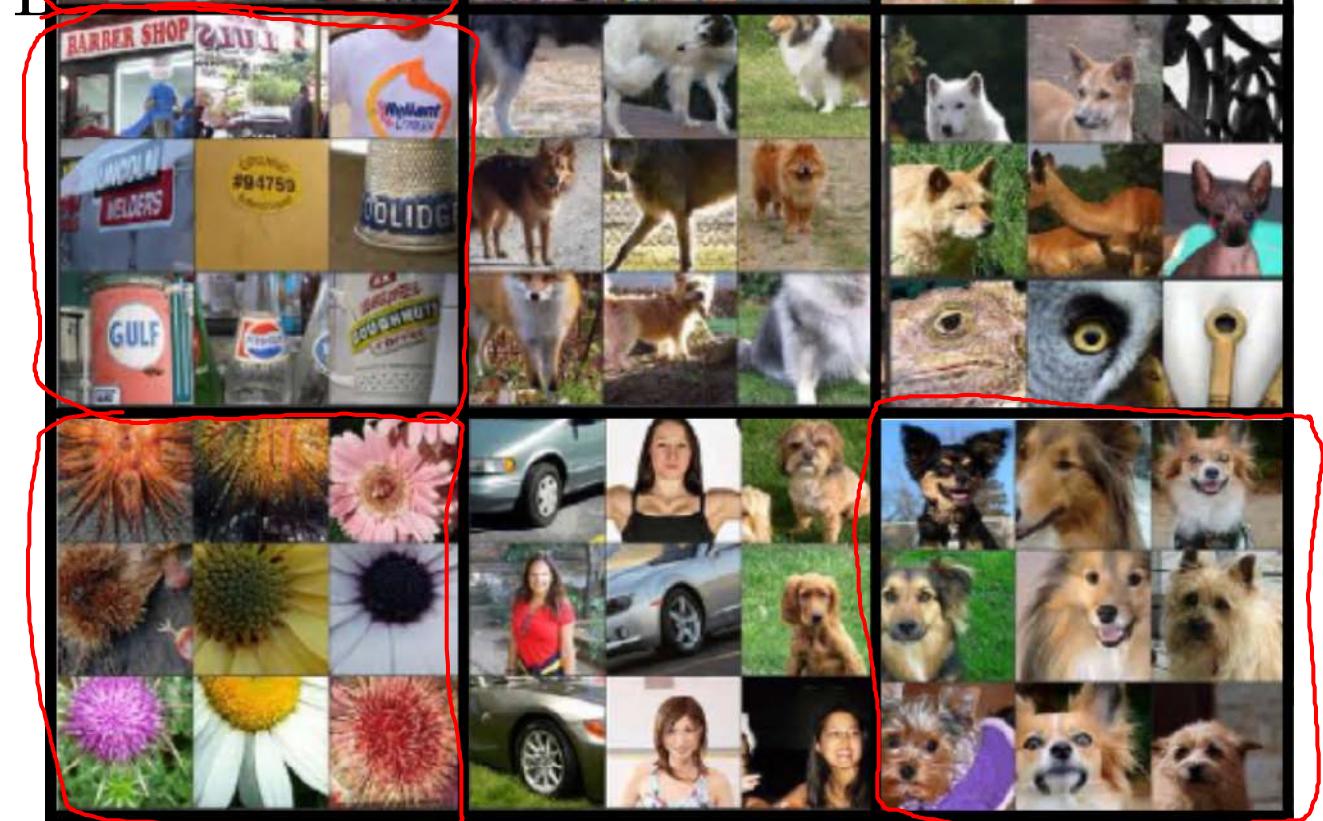
# Visualizing deep layers: Layer 5



Layer 1



L2



Layer 5



deeplearning.ai

# Neural Style Transfer

---

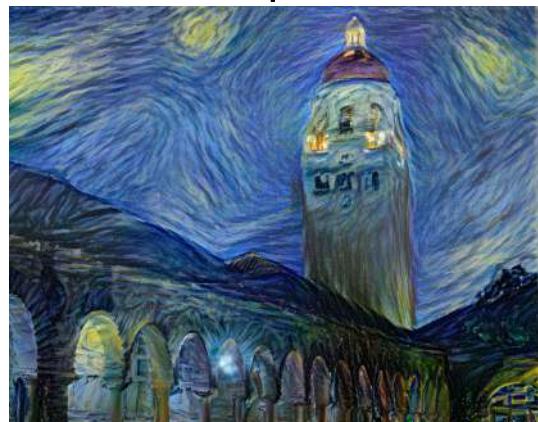
## Cost function

# Neural style transfer cost function



Content C

Style S



Generated image G

We define two parts to the cost function: Content & Style

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

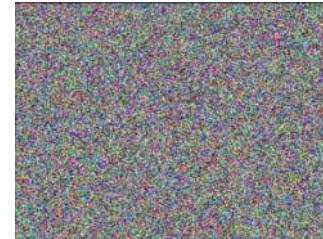
# Find the generated image G

1. Initiate G randomly

$G: 100 \times 100 \times 3$

$\xrightarrow{\text{RGB}}$

start with white-noise image



2. Use gradient descent to minimize  $\underline{J(G)}$

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





deeplearning.ai

# Neural Style Transfer

---

## Content cost function

# Content cost function

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

'l' layer is neither too shallow nor too deep

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

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



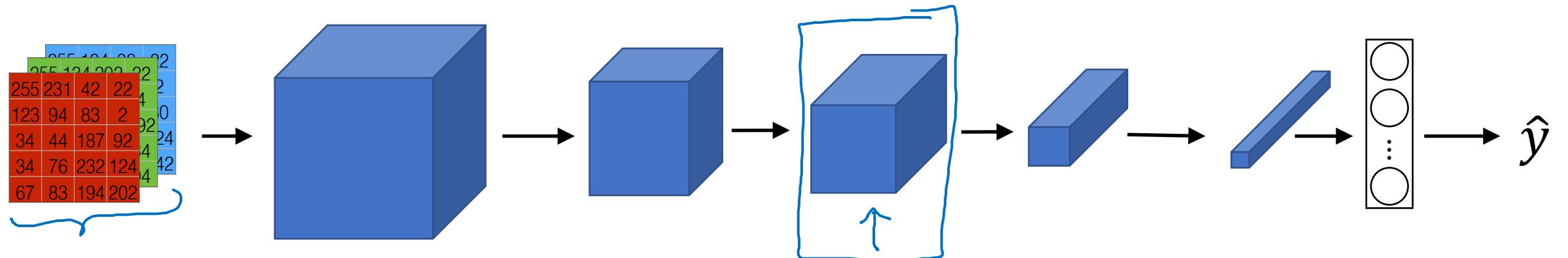
deeplearning.ai

# Neural Style Transfer

---

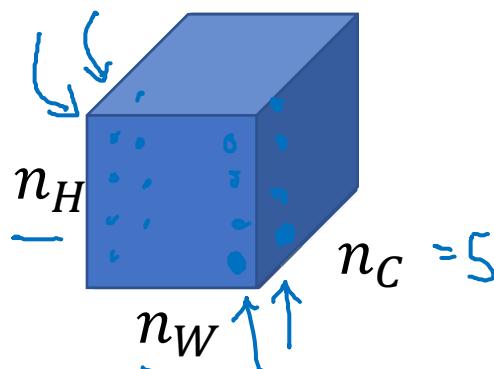
## Style cost function

# Meaning of the “style” of an image



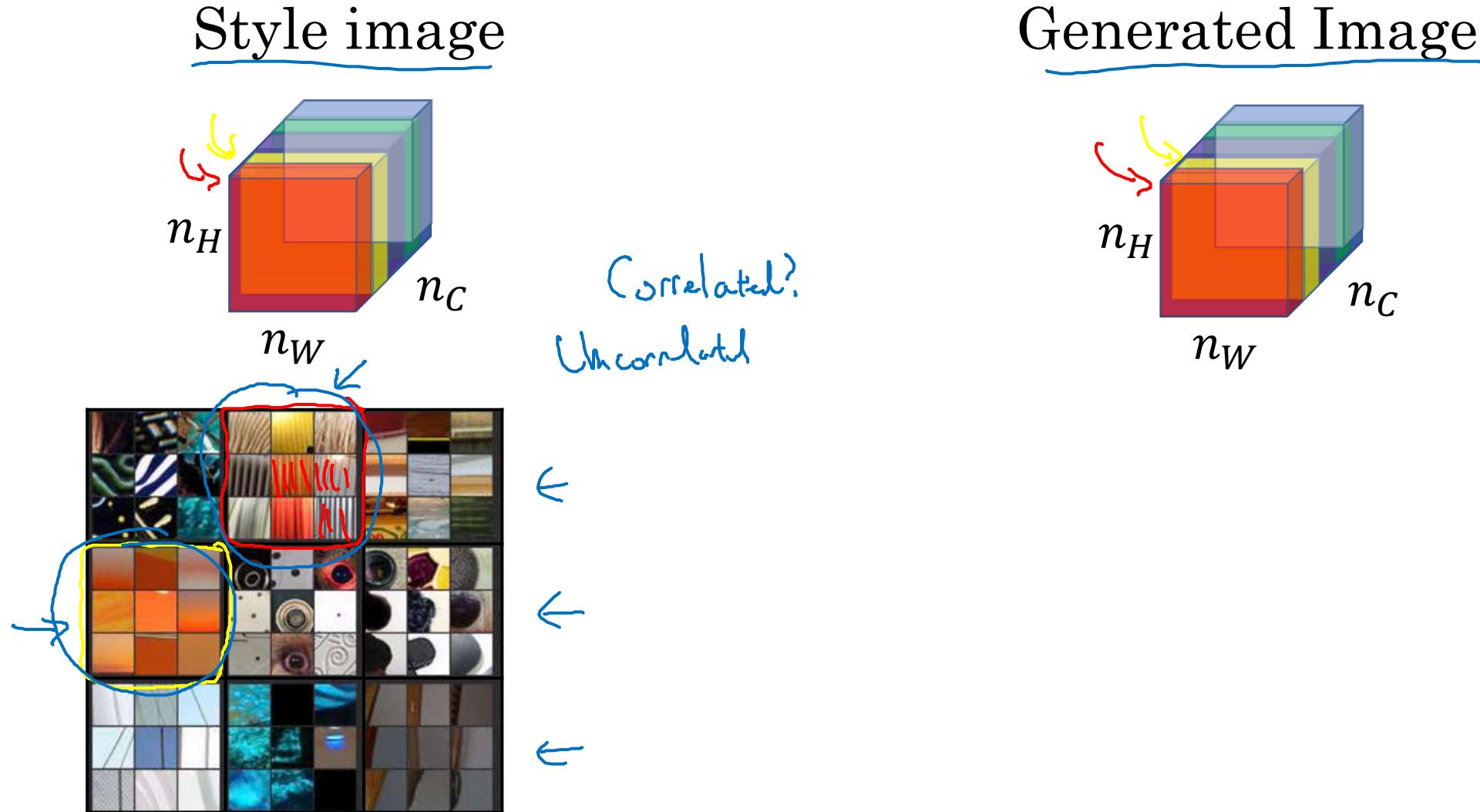
Say you are using layer  $l$ 's activation to measure “style.”

Define style as correlation between activations across channels.



How correlated are the activations  
across different channels?

# Intuition about style of an image



# Style matrix

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

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

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

H W C  
↓ ↓ ↓

$n_c$

$$G_{kk'}^{[l]} \quad \text{for } k, k' = 1, \dots, n_c^{[l]}$$

$$k = 1, \dots, n_c^{[l]}$$

"Gram matrix"

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

# Style cost function

$$\left\| G^{[l](s)} - G^{[l](G)} \right\|_F^2$$

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

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

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



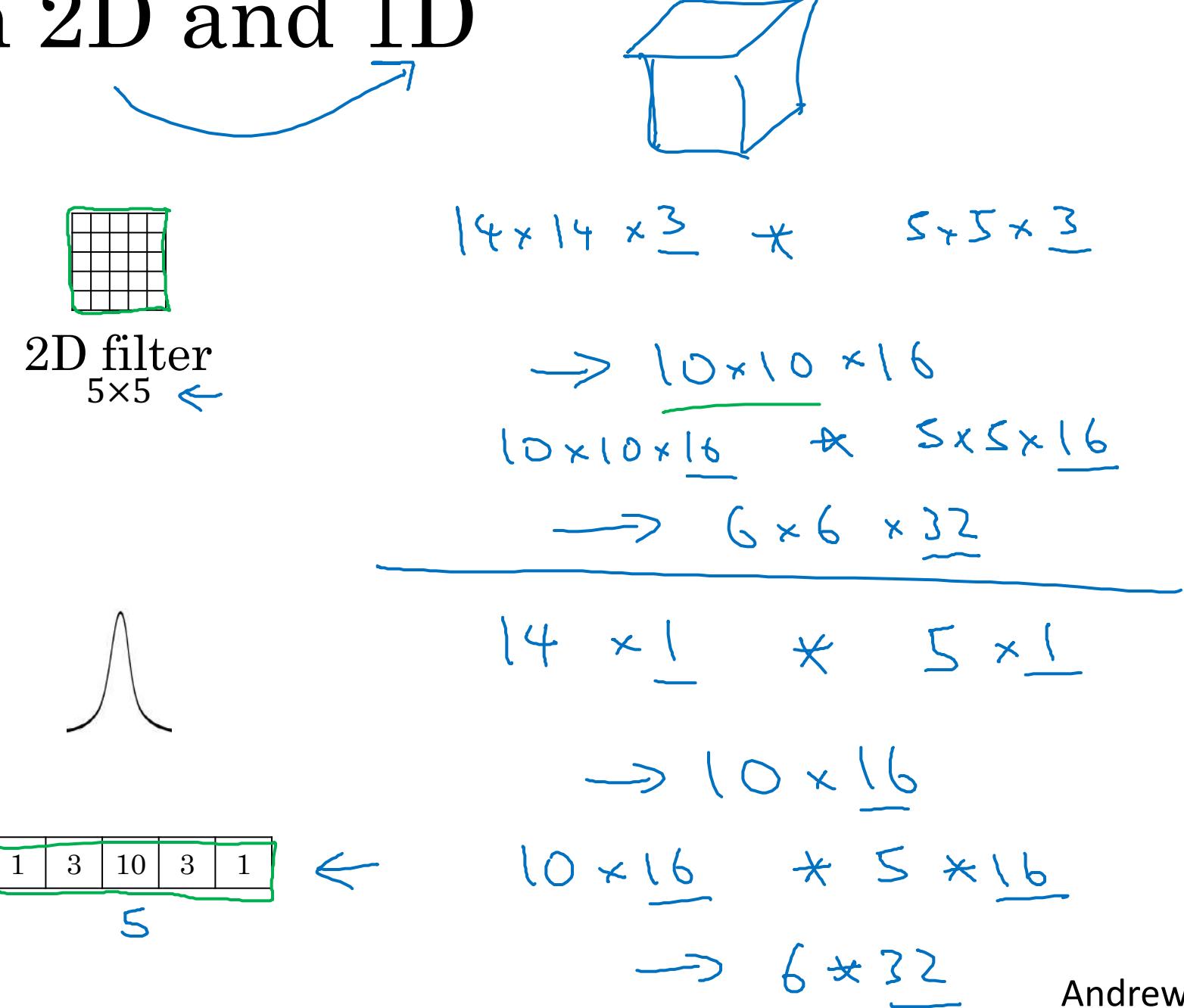
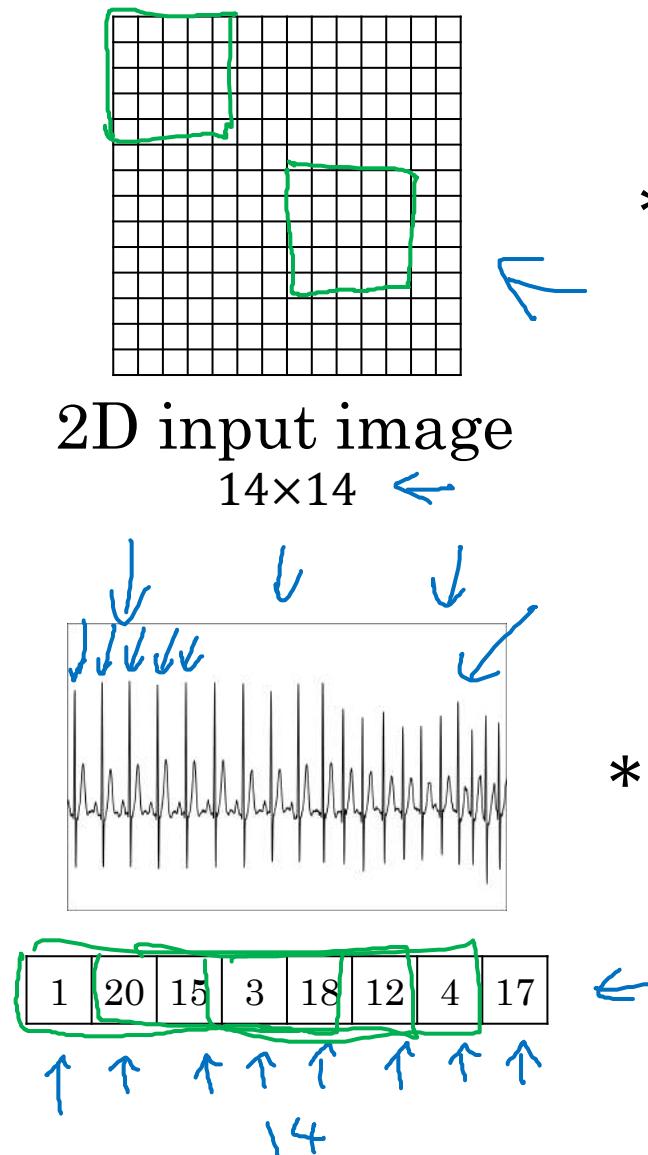
deeplearning.ai

# Convolutional Networks in 1D or 3D

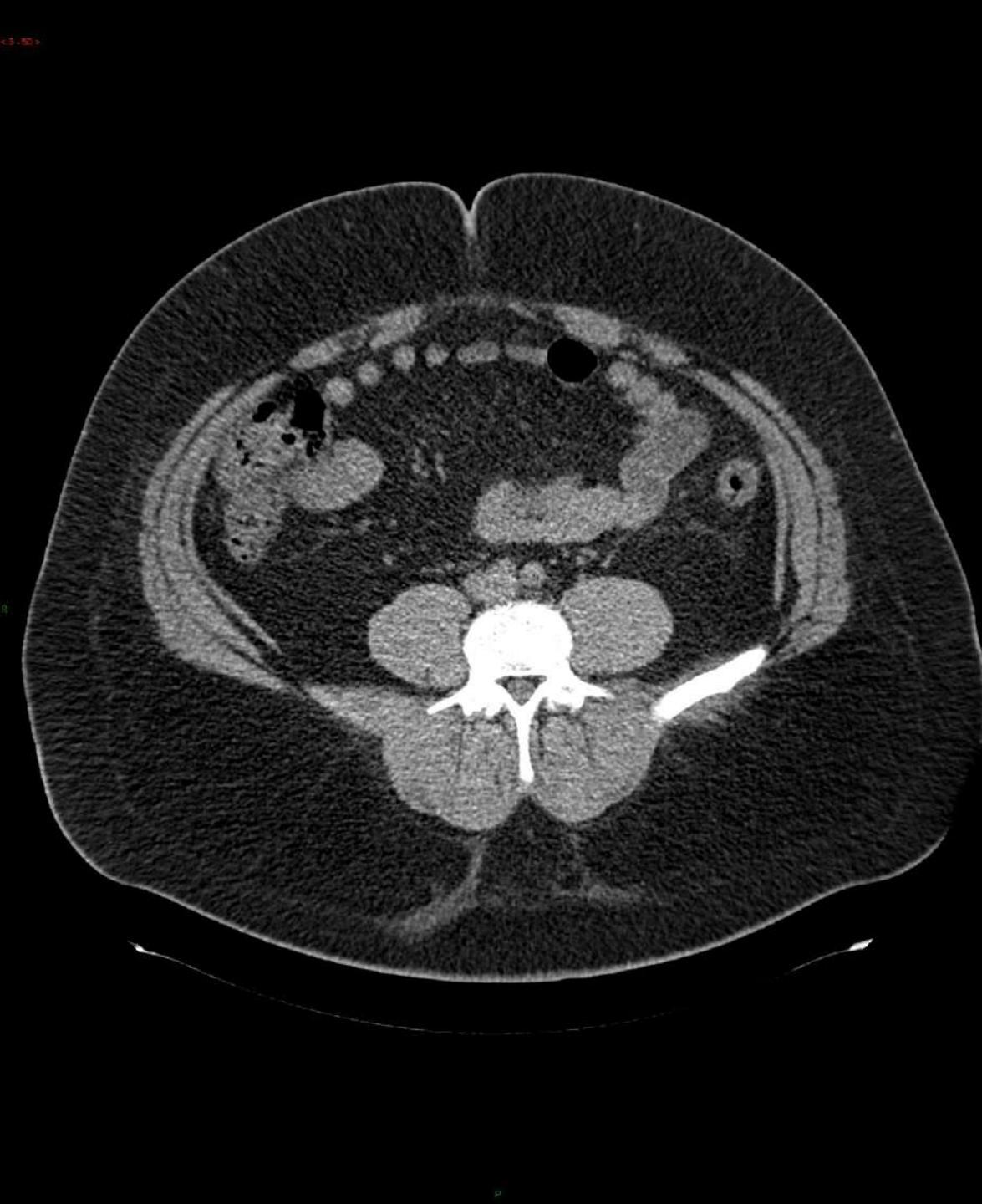
---

1D and 3D  
generalizations of  
models

# Convolutions in 2D and 1D



# 3D data



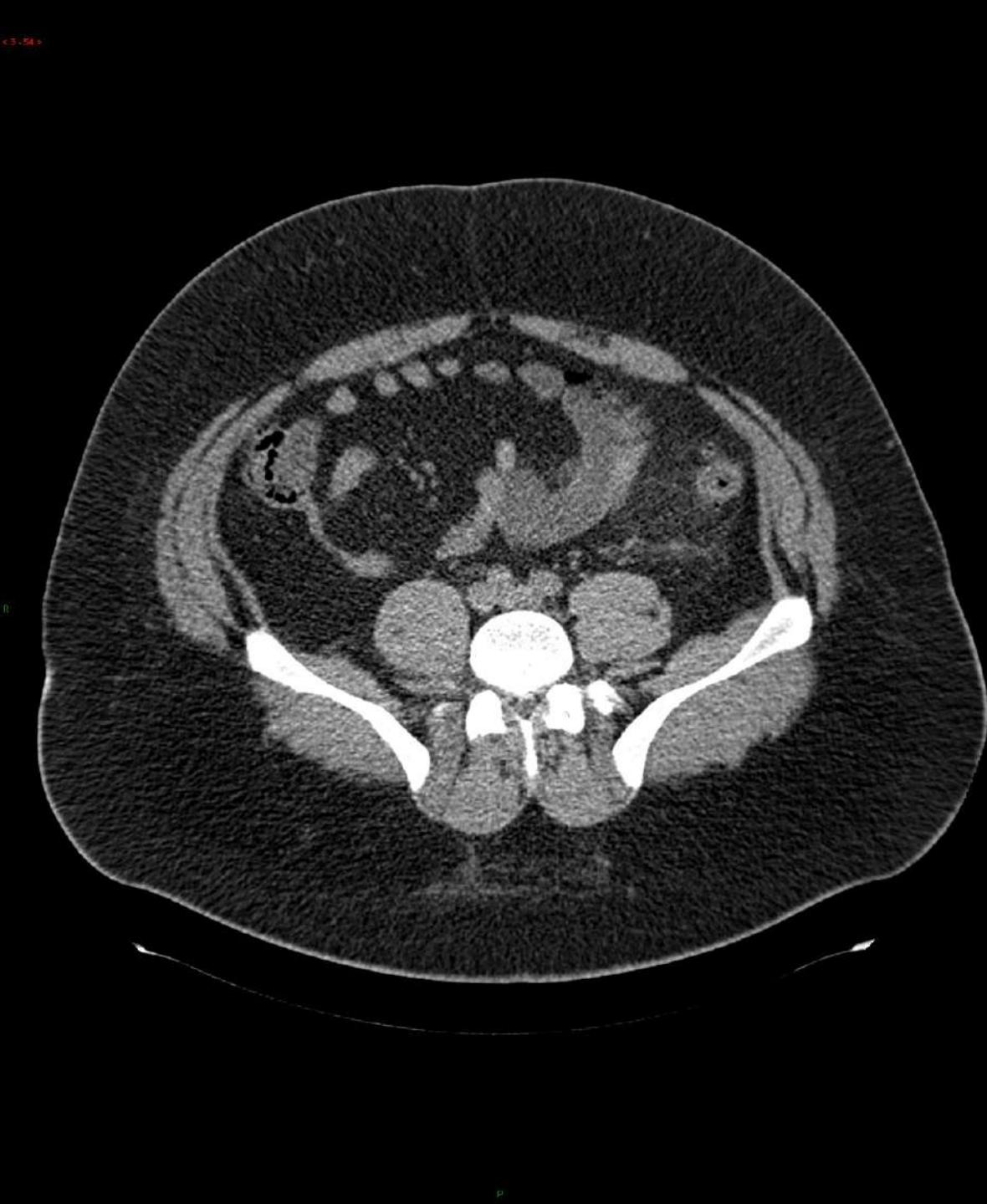
Andrew Ng

3D data



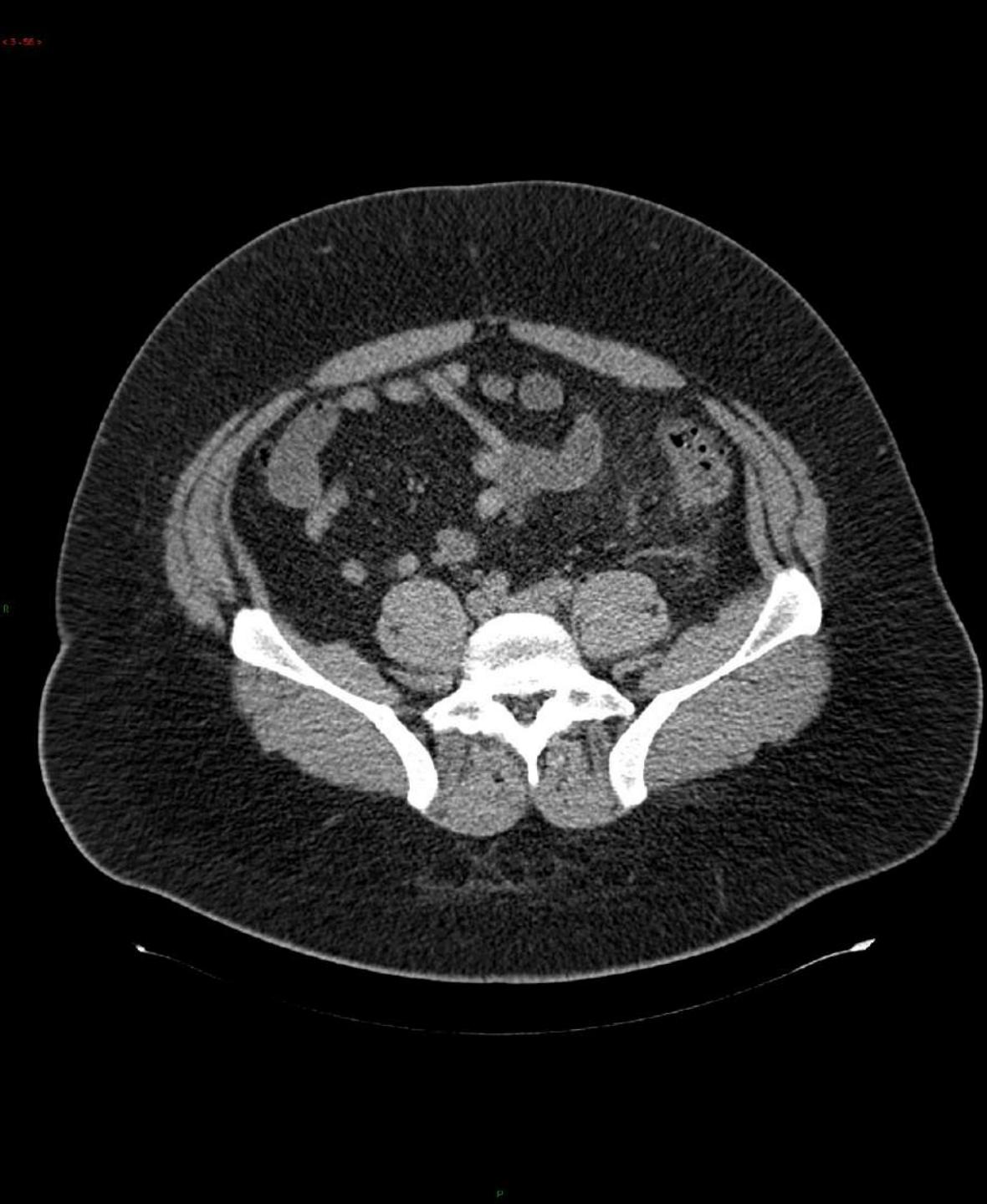
Andrew Ng

# 3D data

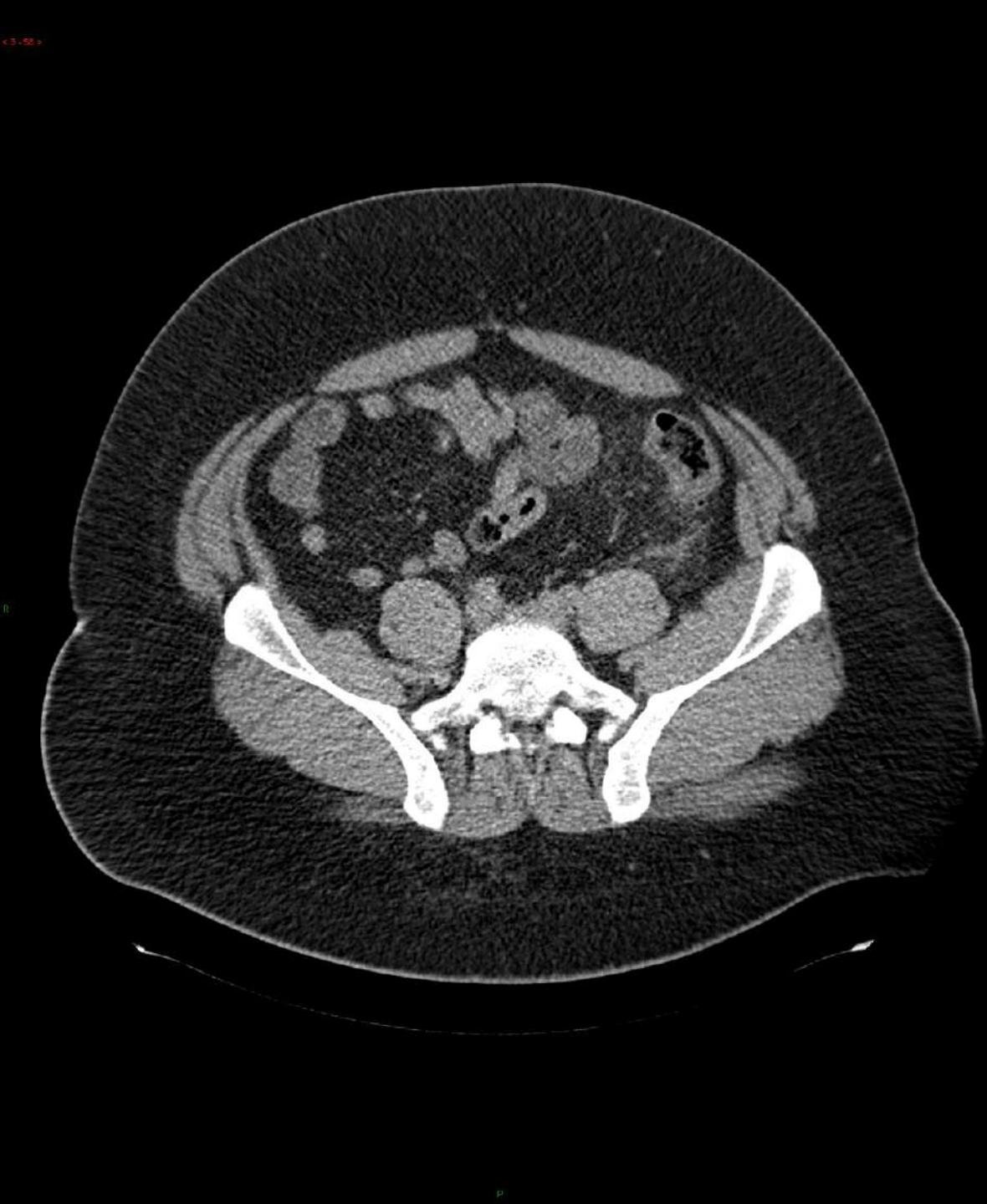


Andrew Ng

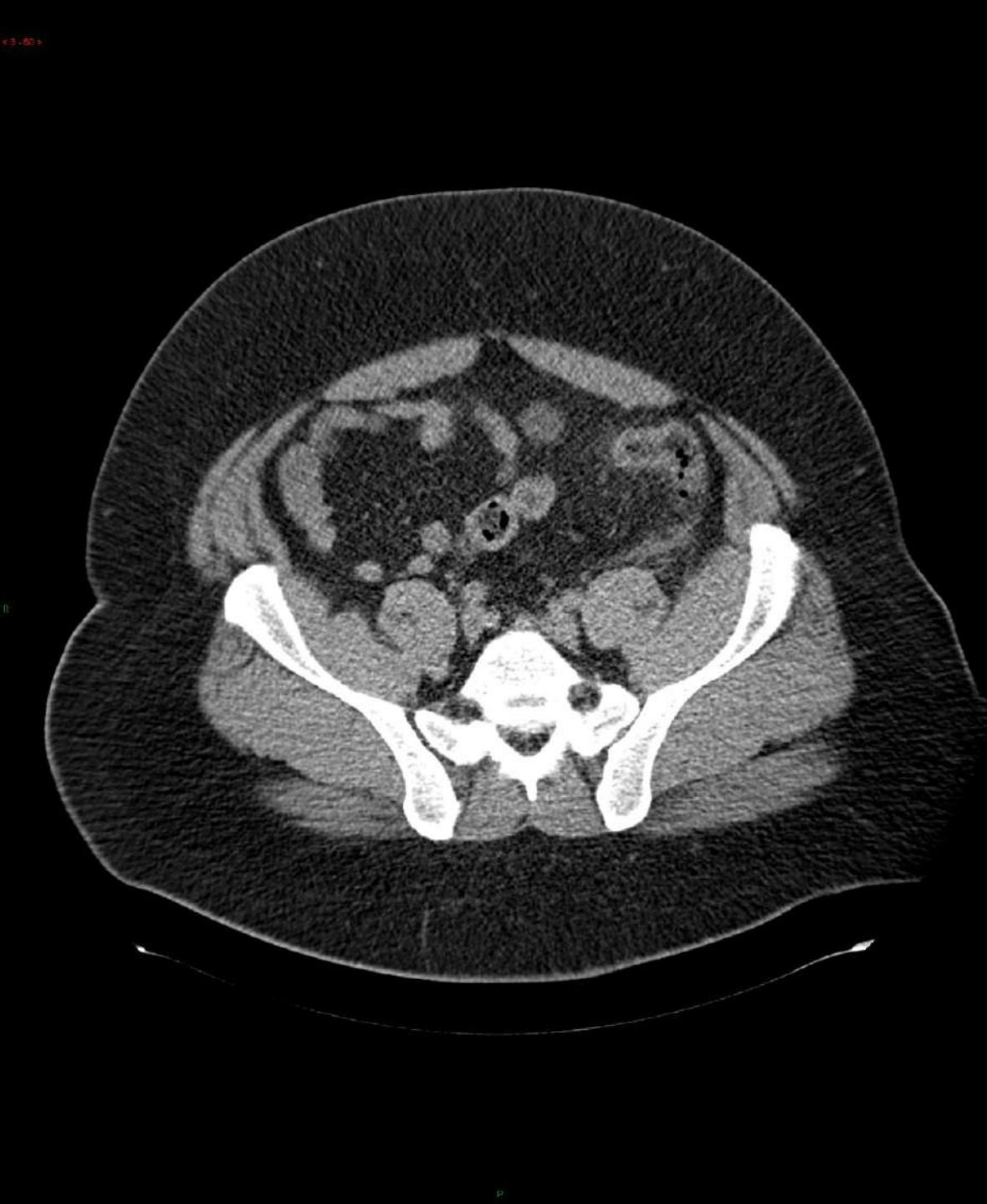
# 3D data



# 3D data



# 3D data



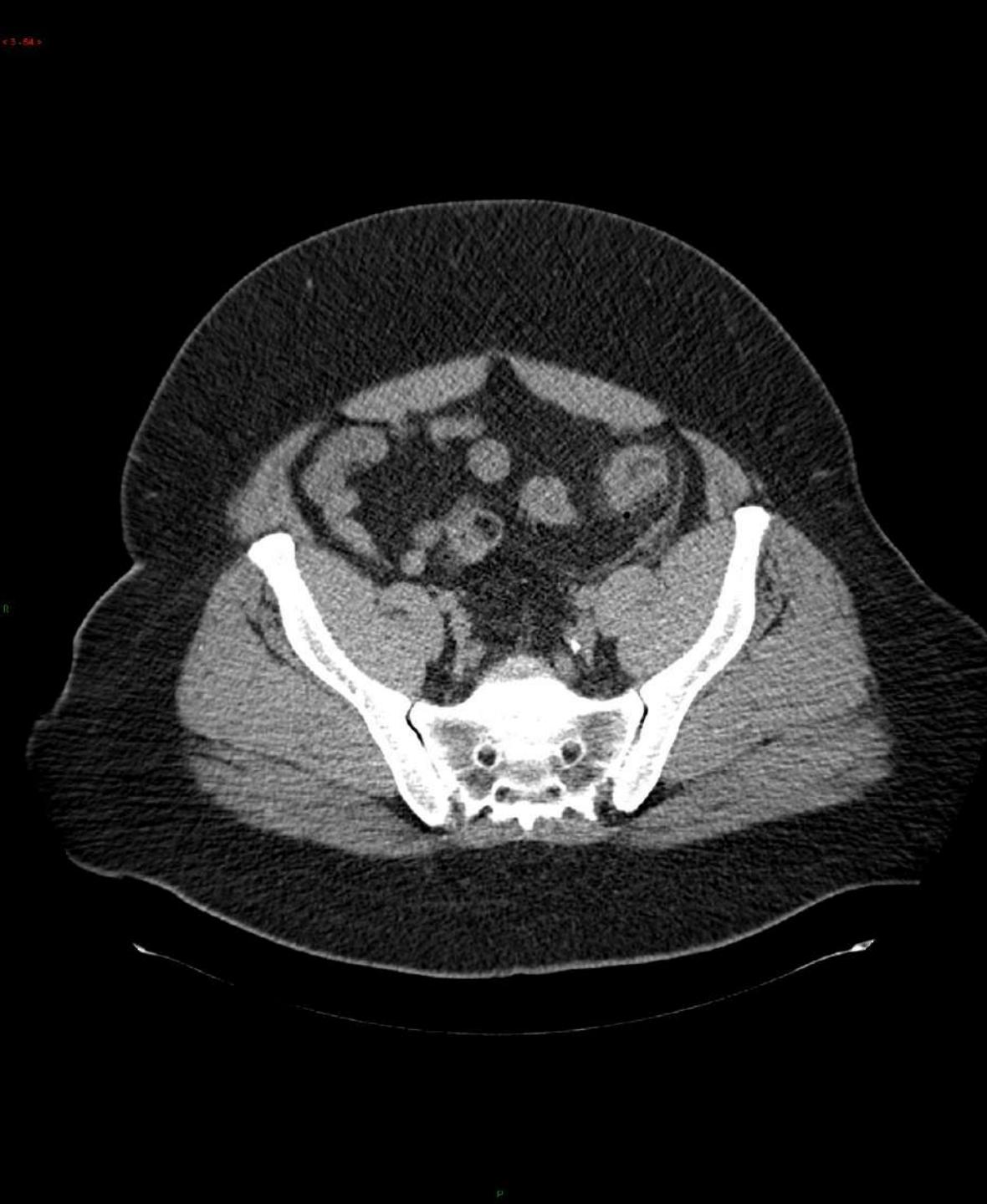
Andrew Ng

# 3D data



Andrew Ng

# 3D data



Andrew Ng

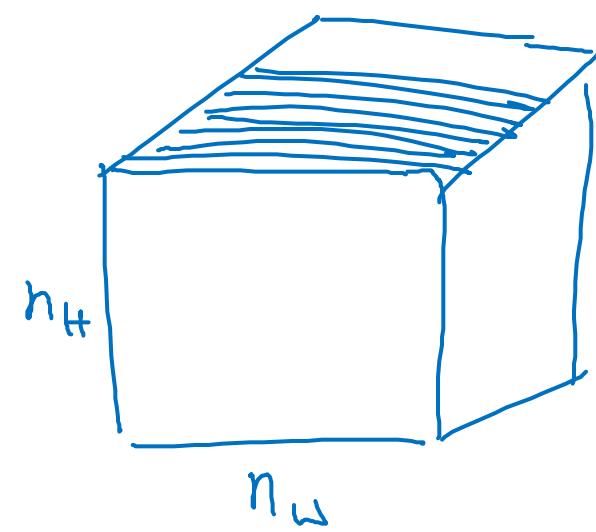
3D data



# 3D data



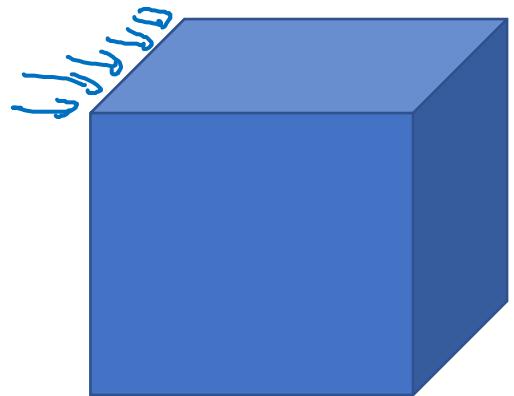
# 3D data



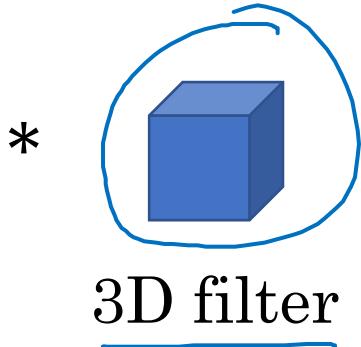
Andrew Ng

# 3D convolution

# of channels can be considered 4th dimension



3D volume



3D filter, with 3rd dimension not same as 3rd dimension of input volume

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

4x14x14  $\times 1$

$\times \quad \underline{5 \times 5 \times 5} \times 1$

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

$\times \underline{5 \times 5 \times 5 \times 16}$

32 filters

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