

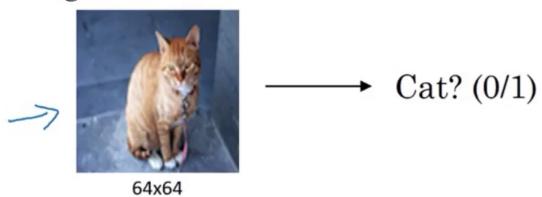
Course 4: Convolutional Neural Network

week 1: Foundations of CNN

Computer Vision

Computer Vision Problems

Image Classification



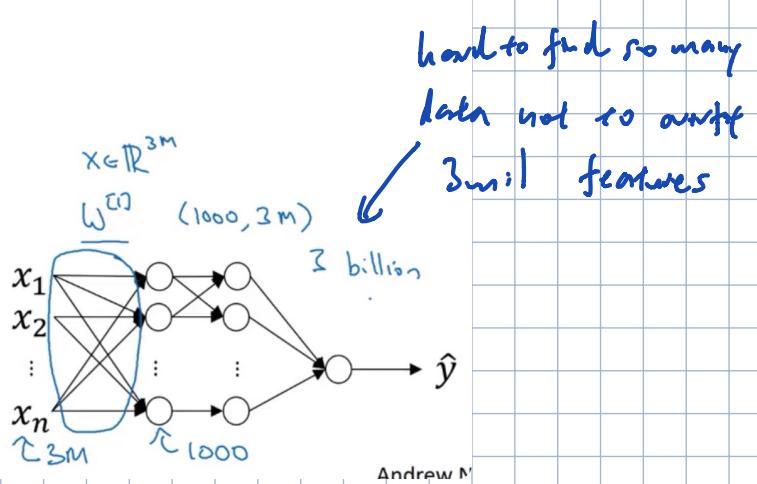
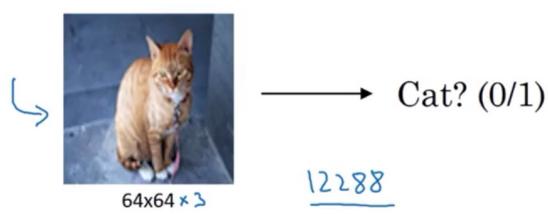
Neural Style Transfer



Object detection



Deep Learning on large images



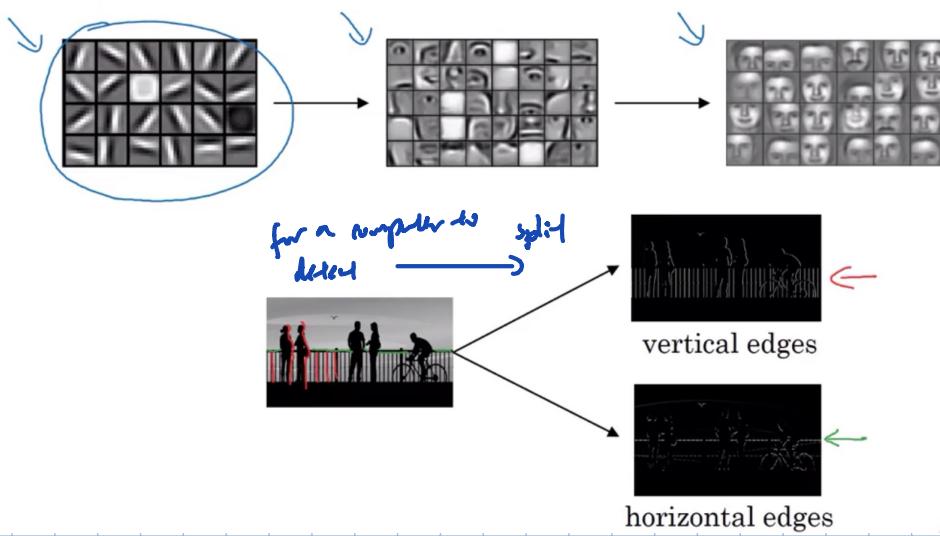
Hence, standard Fully Connected NN would work \rightarrow too many features, and weights have huge number of elements.

\hookrightarrow CNN

Edge detection example

The convolution operation is one of the fundamental building blocks of a CNN. → use edge detection as a motivation example here.

Computer Vision Problem



Vertical edge detection

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

6x6 grayscale

"convolution"

*

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

$\begin{matrix} 3 \times 3 \\ \text{filter or Kernel} \end{matrix}$

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

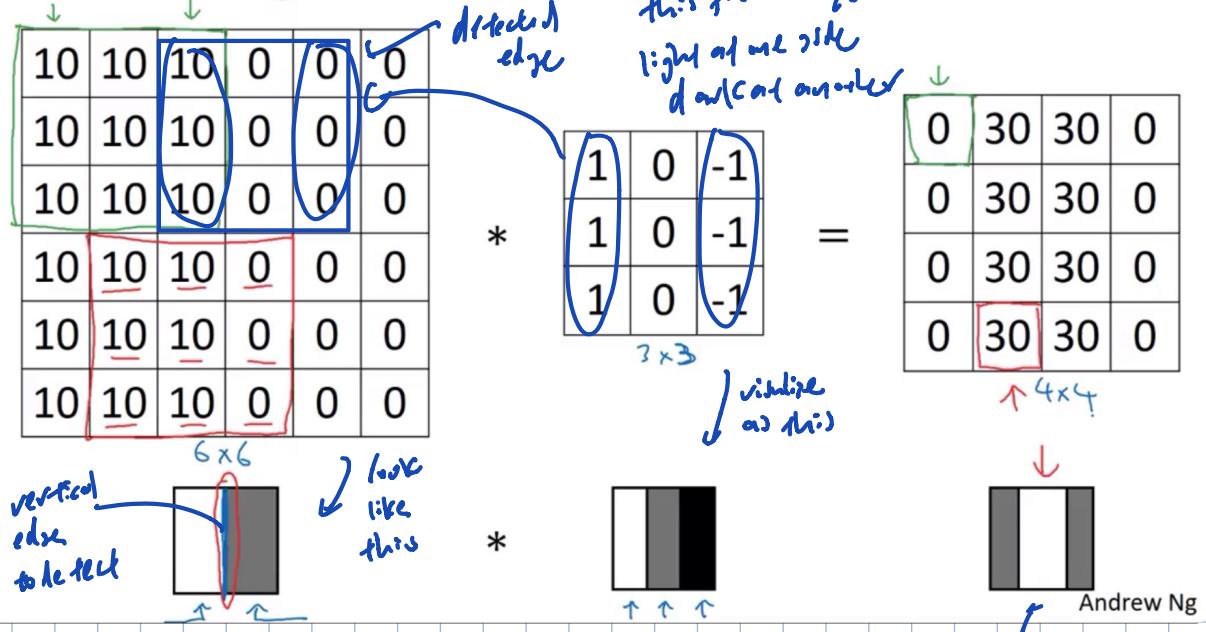
$\begin{matrix} 4 \times 4 \\ \text{output} \end{matrix}$

$$\begin{aligned}
 & 3 \times 1 + 1 \times 1 + 2 \times 1 & 0 \times 1 + 0 \times 5 + 1 \times 7 \\
 & + 0 \times 0 + 0 \times 5 + 0 \times 7 & + 1 \times 0 + 1 \times 0 + 2 \times 0 \\
 & + 1 \times 1 + 1 \times 5 + 1 \times 2 & + 2 \times (-1) + 9 \times (-1) + 5 \times (-1) \\
 \therefore -5 & = -4
 \end{aligned}$$

python: conv2d function (tf.layers)
tensorflow! `tf.nn.conv2d`
`keras Conv2D`

why this filter can act as edge detection?

Vertical edge detection

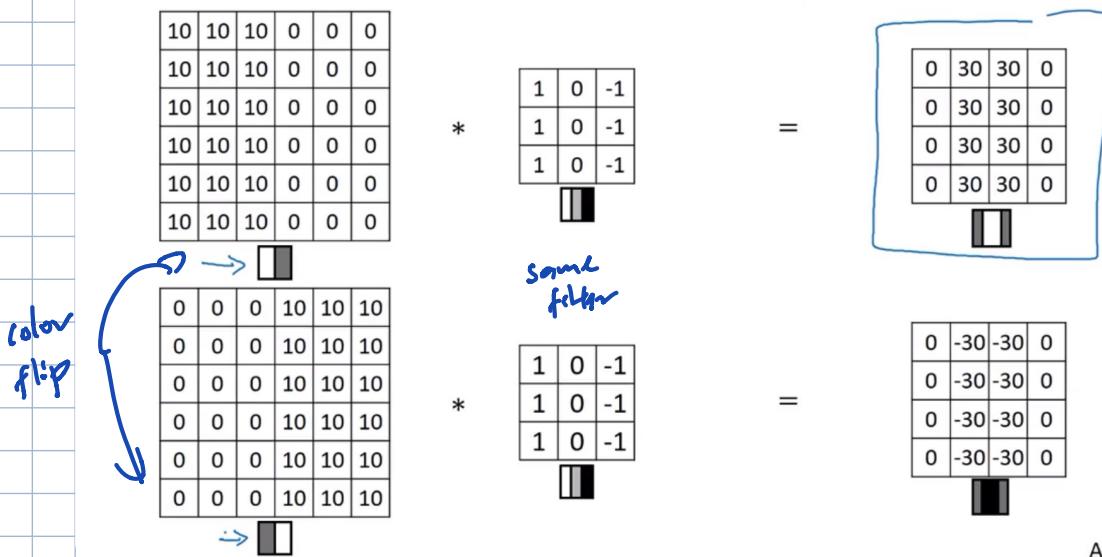


looks thick because we work with
small image, if 100x100 image it will
be quite accurate

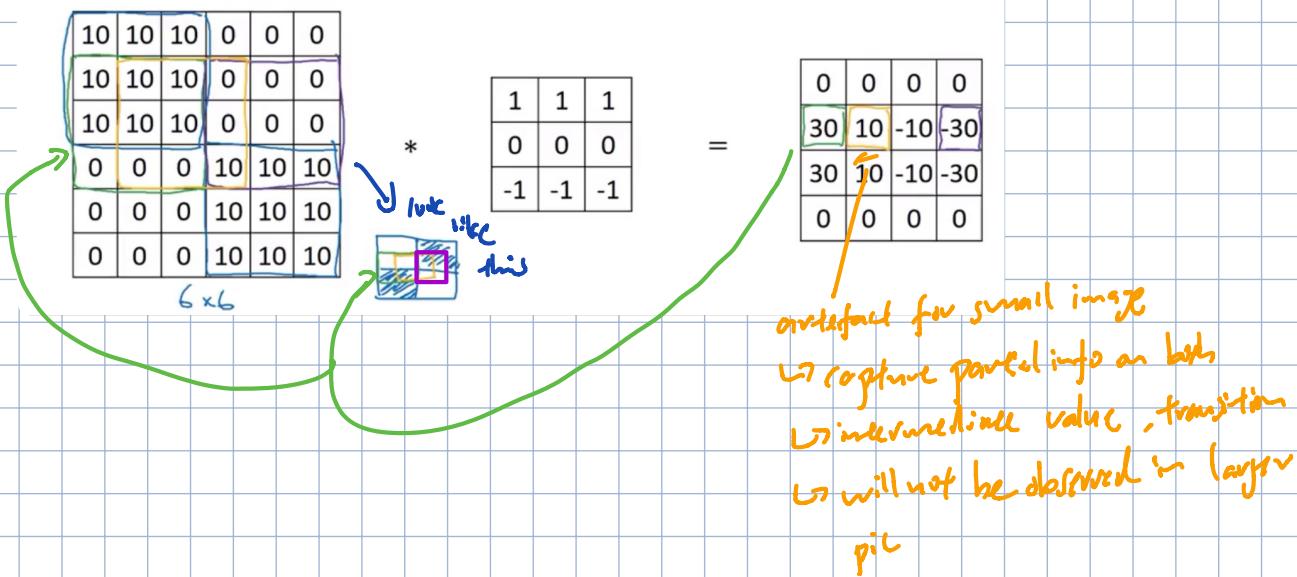
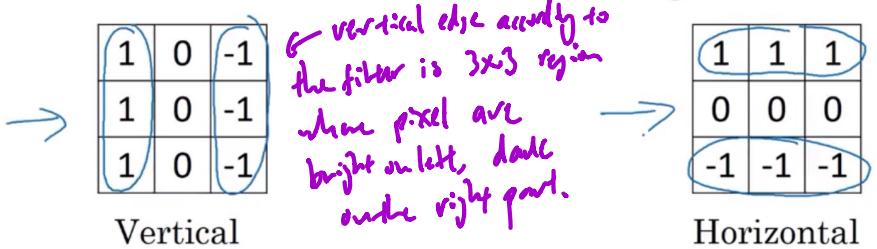
More Edge Detection

Previous example is about detecting vertical edge. Here we investigate positive and negative edges, i.e. difference between light to dark versus dark to light edge transitions.

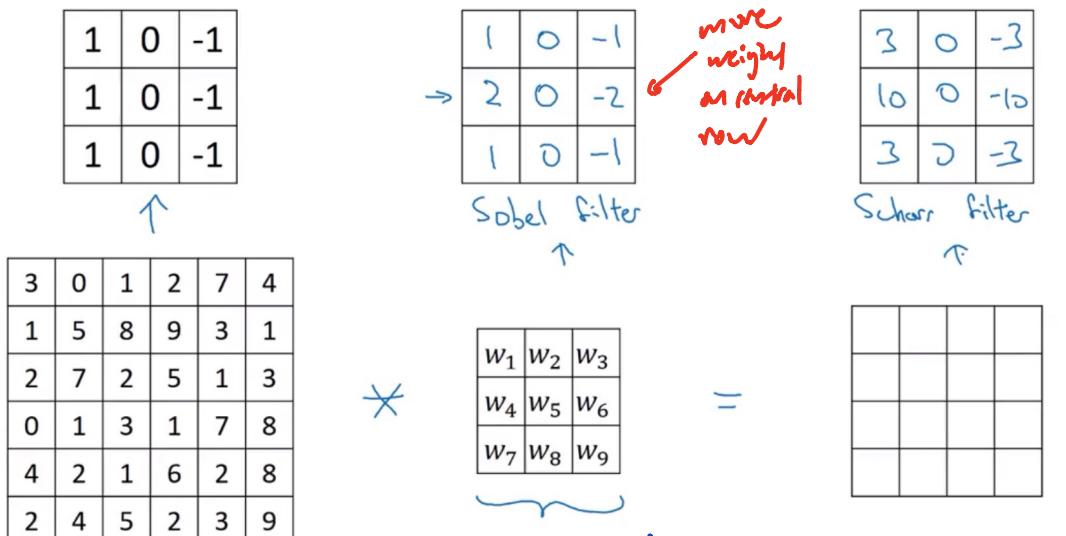
Vertical edge detection examples



Vertical and Horizontal Edge Detection



Learning to detect edges



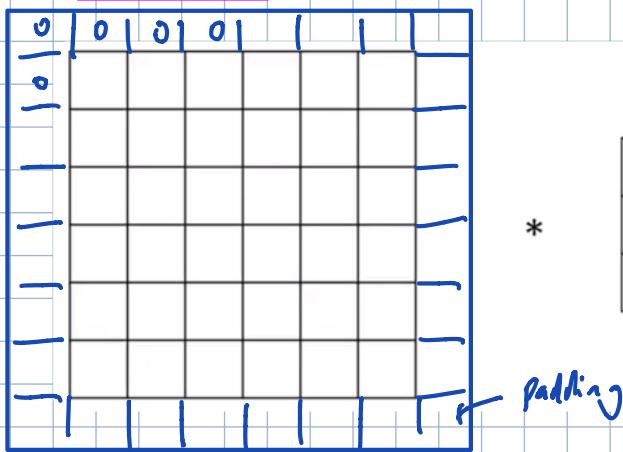
we can learn these instead of hand pick them

→ gives good edge detector

↳ can be chosen to learn

45°, 75° edge

Padding



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

previously: $6 \times 6 \times 3 \times 3 \rightarrow 4 \times 4$; now: $\xrightarrow[\text{padding}]{} 8 \times 8 \times 3 \times 3 \rightarrow 6 \times 6$

$\hookrightarrow n \times n \times f \times f \rightarrow (n-f+1) \times (n-f+1)$
image filter assume stride = 1

Downside: ① every time we apply a convolutional operator, the image shrinks
 \hookrightarrow image will be really small after a few layers

② the image at the corners/sides are only being used once
 \hookrightarrow throwing away info at the pixel at the sides

$p = \text{padding} = 1 \Rightarrow$ padding all around with an extra border of 1 pixels.

$$\begin{aligned} &\hookrightarrow (n + 2p - f + 1) \times (n + 2p - f + 1) \\ &= (6 + 2 - 3 + 1) \times (6 + 2 - 3 + 1) \\ &= 6 \times 6 \end{aligned}$$

Valid and Same convolutions

In terms of how much to pad, there are 2 common choices:

- ① valid convolution
- ② same convolution

"Valid": no padding

$$n \times n \times f \times f \rightarrow (n-f+1) \times (n-f+1)$$

e.g. $6 \times 6 \times 3 \times 3 \rightarrow 4 \times 4$

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

$$(n \times n + p) \times f \times f \rightarrow (n+2p-f+1) \times (n+2p-f+1)$$

e.g. if we want output size same as input size

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

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

By convention in CV, f is always odd number

↪ if f is even, then need asymmetric padding

↪ odd filter has a central position → distinguisher, nice to have

Strided Convolutions

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

$$\begin{array}{ccc} * & \begin{array}{|c|c|c|}\hline 3 & 4 & 4 \\ \hline 1 & 0 & 2 \\ \hline -1 & 0 & 3 \\ \hline \end{array} & = \\ & \begin{array}{|c|c|c|}\hline 91 & 100 & 83 \\ \hline 67 & 91 & 127 \\ \hline 44 & 72 & 74 \\ \hline \end{array} & \\ & 3 \times 3 & \\ & \text{stride} = 2 & \end{array}$$

7×7

In general:

$$\begin{matrix} n \times n & \times & f \times f \end{matrix} \quad \rightarrow \quad \left(\frac{n+2p-f}{s} + 1 \right) \times \left(\frac{n+2p-f}{s} + 1 \right)$$

padding p , stride s

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

→ what if this fraction is not an integer? then we take the floor

↳ round the number down

i.e. the way this is implemented is that, the box multiplication is only performed when the box is fully contained within the image or image + padding

(in short, filter has to be completely within image or image + padding - then perform computation)

Technical note on cross-correlation vs. convolution

Convolution in math textbook:

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

$$* \quad \begin{array}{|c|c|c|} \hline 3 & 4 & 5 \\ \hline 1 & 0 & 2 \\ \hline -1 & 9 & 7 \\ \hline \end{array}$$

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

flip
first
and
axis

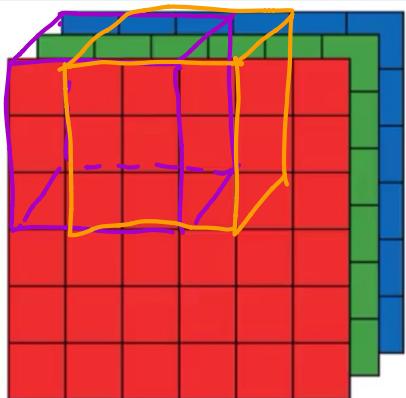
use this
filter

this is proper convolution from mathbook

we skip flipping in DL, so it should be called cross-correlation instead.

but we don't
care, just
call
convolution
anyway

Convolutions over volume



$6 \times 6 \times 3$
height width channel

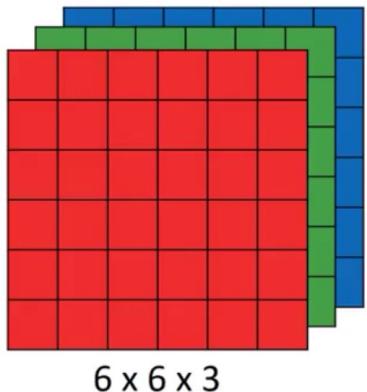
$$* \quad \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix} = \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix}$$

$3 \times 3 \times 3$
height width channel
6728 numbers
columnwise
multipl. then sum

4×4

vertical edge for
red channel

Multiple filters



$6 \times 6 \times 3$

detect vertical and horizontal edge at the same time

$$* \quad \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix} = \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix}$$

vertical edge
 $3 \times 3 \times 3$
horizontal edge

$$* \quad \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix} = \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix}$$

$3 \times 3 \times 3$

4×4

$4 \times 4 \times 2$

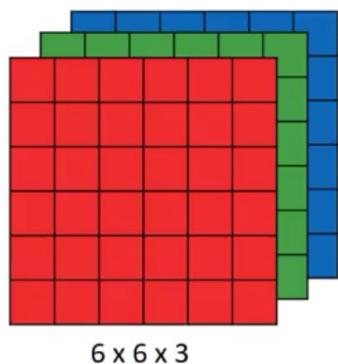
Summary: $n \times n \times n_c \times f \times f \times n_c \rightarrow (n-f+1) \times (n-f+1) \times n_c'$

e.g. $6 \times 6 \times 3 \times 3 \times 3 \times 3 \xrightarrow{\text{same}} (6-3+1) \times (6-3+1) \times 3' = 4 \times 4 \times 3'$

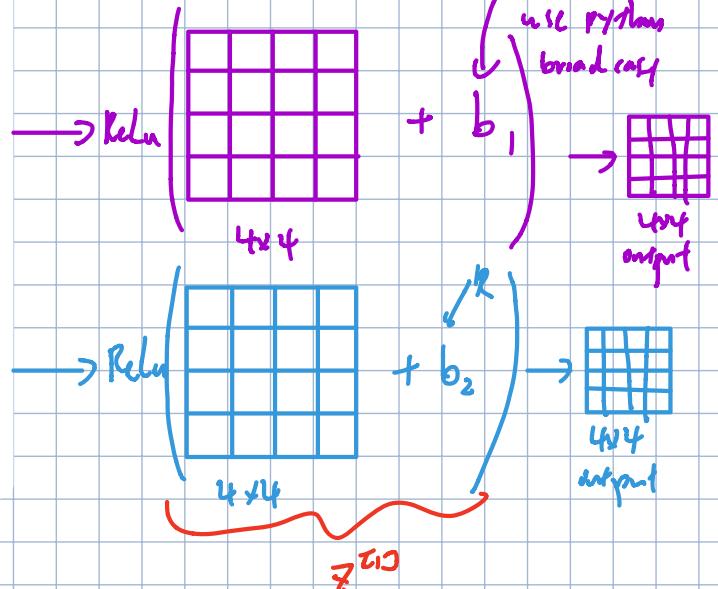
assume
Stride = 1
no. of filters
we use

One layer of a Convolutional Network

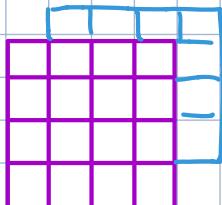
Example of a layer convolving with 2 filters



$$\begin{aligned} * & \quad \text{3x3x3} \\ * & \quad \text{3x3x3} \end{aligned}$$



then we stack up the 2 outputs



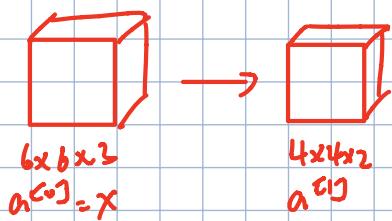
From $6 \times 6 \times 3 \rightarrow 4 \times 4 \times 2$
1 layer of convolutional neural network

Recall that 1 step of forward prop is:

$$Z^{(1)} = W^{(1)} a^{(0)} + b^{(1)}$$

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

Here, $6 \times 6 \times 3$ image $\rightarrow a^{(0)}$ or X
the filters $3 \times 3 \times 3$ is similar to $W^{(1)}$
i.e. here 1 layer



In this example, there is 2 filters (or 2 features)

If 10 filters are used, then we will get $4 \times 4 \times 10$ output volume

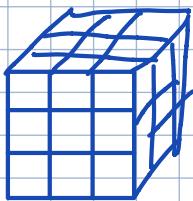
Number of parameters in one layer

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?

each filter is:

$$3 \times 3 \times 3$$

$\hookrightarrow 27$ parameters



$$+ \text{bias} = 28 \text{ parameters}$$

so in total we have 280 parameters

No matter your input image size, the parameters in this layer will still be the same.
so these 10 filters can detect $u, H, 45^\circ M$ for large images as well.

Summary of notation

If layer l is a convolution layer:

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

$$p^{[l]} = \text{padding} \quad \begin{matrix} \text{valid: no padding} \\ \text{same: input size} = \text{output size} \end{matrix}$$

$$s^{[l]} = \text{stride}$$

Input to this layer:

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

height, weight
 of activation
 input from
 previous layer

\downarrow
 layer
 of activation
 input from
 previous layer

Output of this layer

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

\downarrow
 # of filters
 in layer l

$\hookrightarrow n_H^{[l]}:$ $\left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$
 height &
 output width

$$\Rightarrow \text{Each filter is: } f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \quad \begin{matrix} \checkmark \text{ number of neurons from} \\ \text{previous layer} \end{matrix}$$

$$\Rightarrow \text{Activation: } a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

(softmax)
 + nonlinearity
 ReLU

$$\Rightarrow \text{ReLU activation: } A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

⇒ weight:

$$f_{\text{enc}} \times f_{\text{enc}} \times n_c^{ce-1} \times n_c^{ce} \times T_{\text{total}}$$

of filters

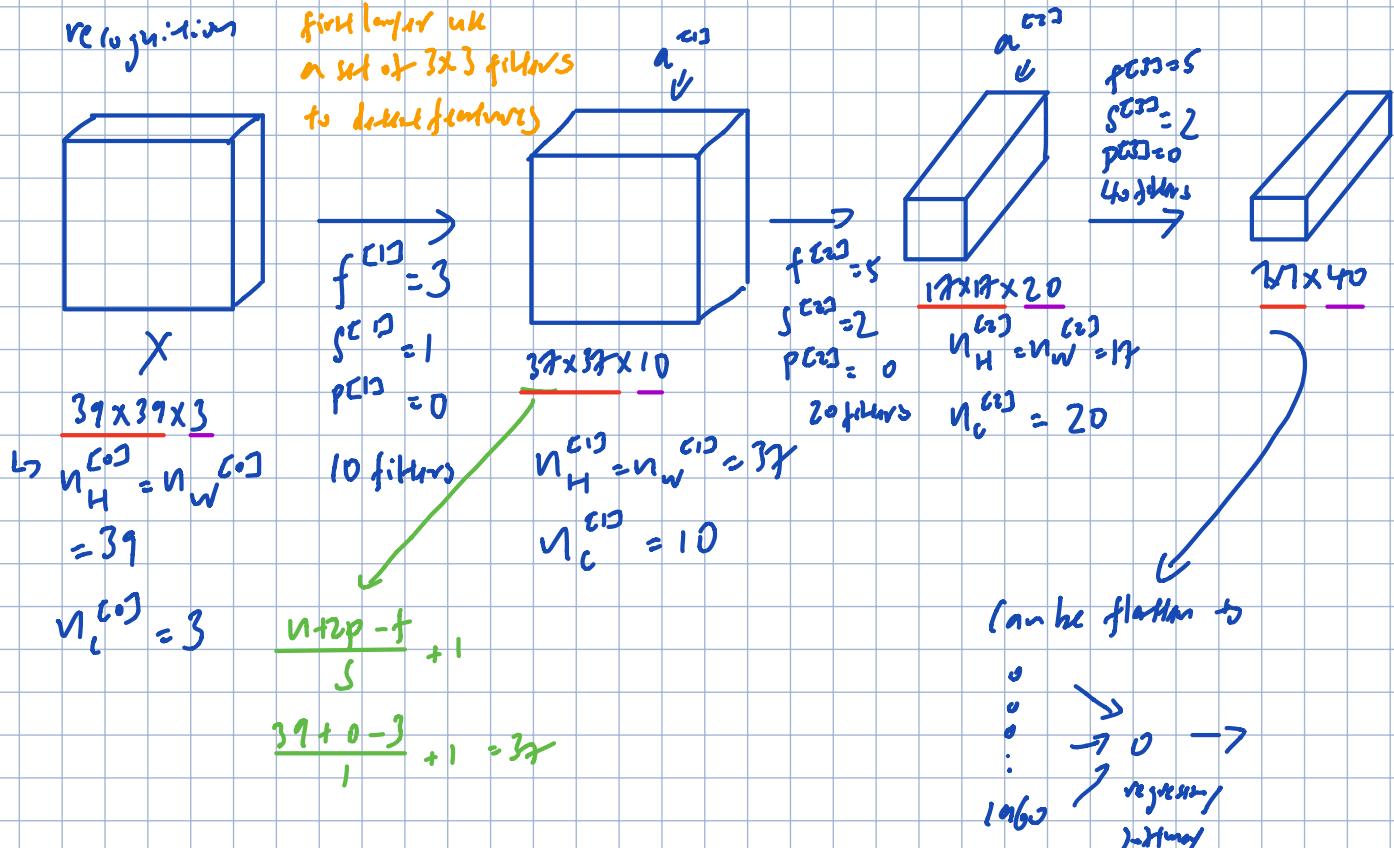
⇒ bias: 1 real number for each filter

$n_c^{(k)}$, a vector of this dimension

Let's do it: $C(1, 1, 1, u_6^{(1)})$ for convenience

Simple Convolutional Network Example

Say we have an image, we want to do image classification or image recognition first layer will consist of 3x3 filters.



usually image size decreases with layers
while channel size increase with layers

Types of layer in a convolutional network

① convolution (conv)

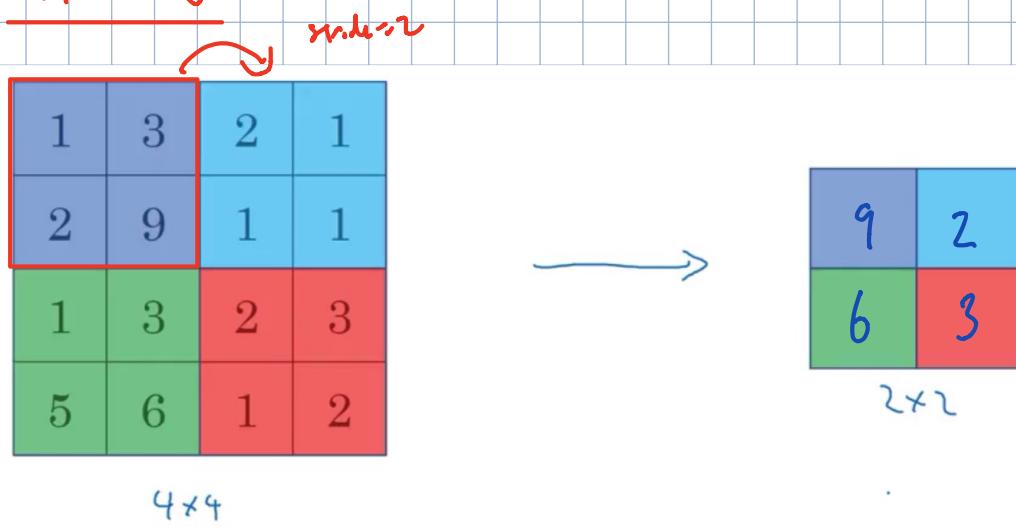
② pooling (Pool)

③ Fully connected (FC)

Pooling Layers

Pooling layers can be used to reduce the size of the representation, hence speed up the computation & also makes some of the features more robust.

Max Pooling



As if we apply a filter size of 2 , stride of 2

↳ hyper parameters of max pooling. ↳ nothing for gradient descent to learn

Intuition: input are the activations of some layer of NN, max pooling detect a particular feature in the selected area

↳ If these features are detected anywhere in this filter, then keep a high number.

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

$5 \times 5 \times n_c$



9	9	5
9	9	5
8	6	1

$$f=3 \quad 3 \times 3 \times n_c$$

$$S=1.$$

formula to determine
output size at layer V
also works for max pooling

If there are multiple channel, do max pooling independently on each channel.

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

$$f=2$$

$$S=2$$

This is used when very deep in the NN, collapse $7 \times 7 \times 1000 \rightarrow 1 \times 1 \times 1000$

Summary of pooling

Hyperparameters:

f : filter size $f=2, S=2$

s : stride $f=3, S=2$

Max or average pooling

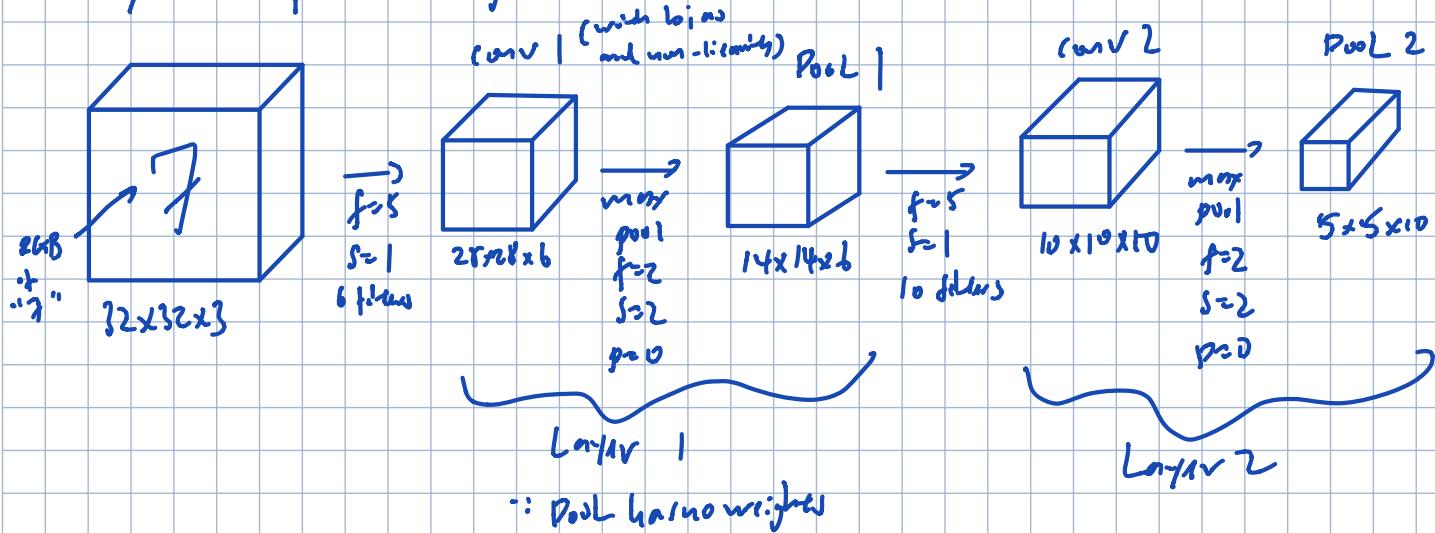
$\rightarrow p$: padding

No parameters to learn!

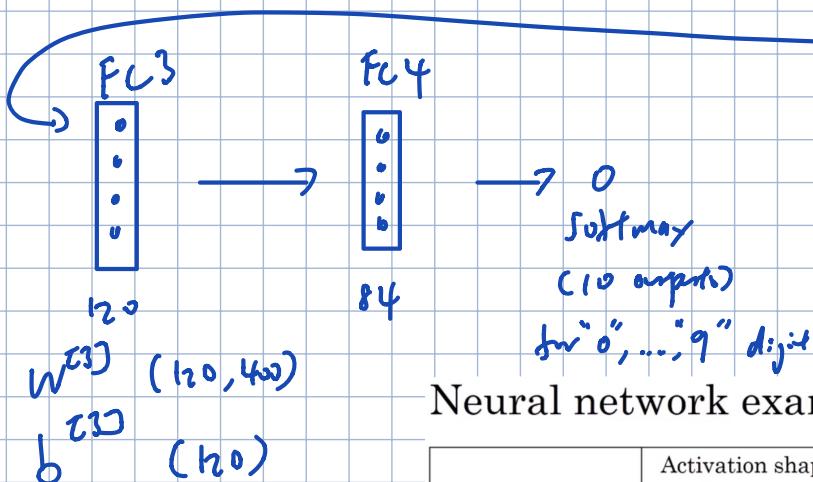
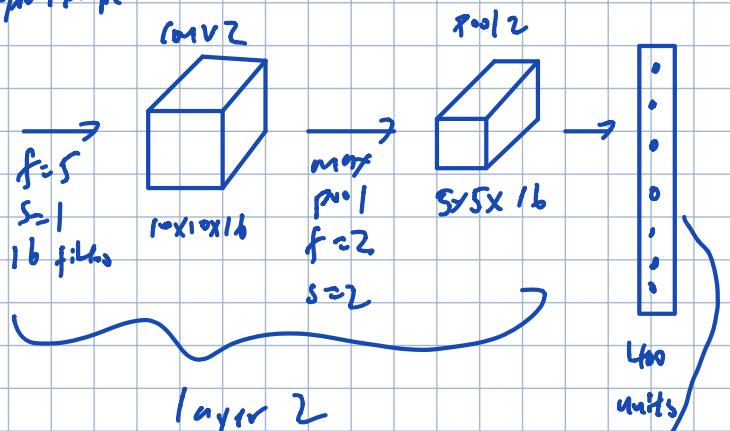
$$\begin{aligned} & N_H \times N_W \times N_C \\ & \downarrow \\ & \left\lfloor \frac{N_H-f+1}{S} \right\rfloor \times \left\lfloor \frac{N_W-f}{S} + 1 \right\rfloor \\ & \times N_C \end{aligned}$$

CNN example

Say we input an image which one $32 \times 32 \times 3$



another layer example



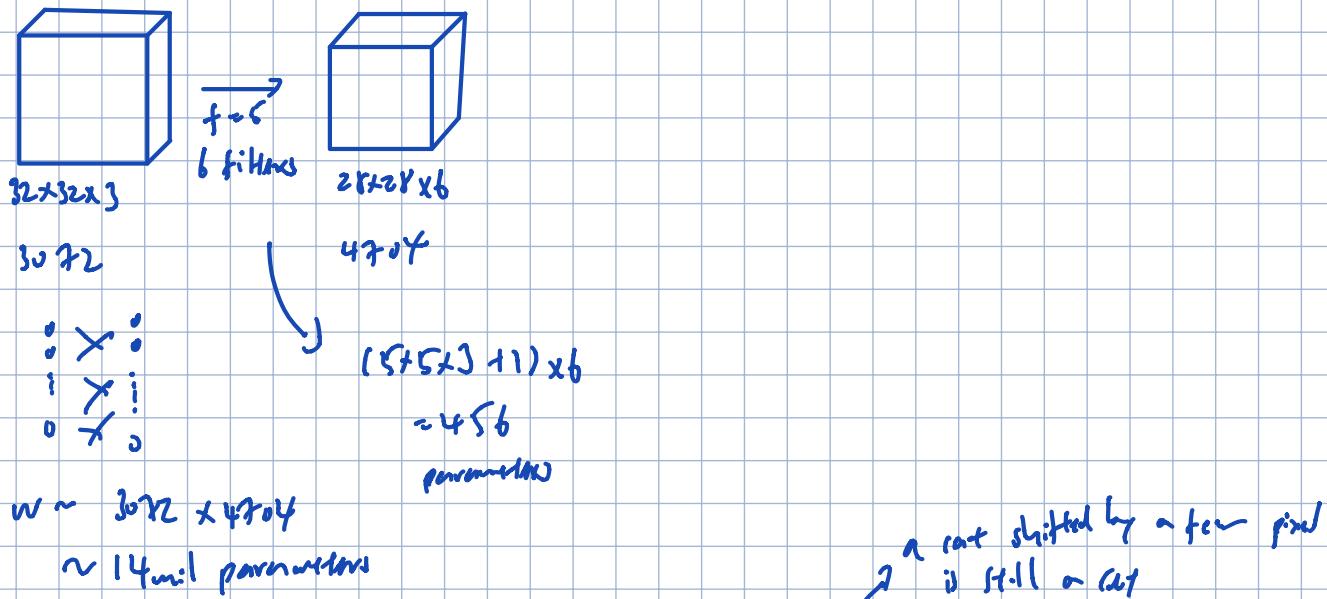
Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32, 32, 3)	3,072	0
CONV1 ($f=5$, $s=1$)	(28, 28, 8)	6,272	208 $\rightarrow (5 \times 5 \times 3 + 1) \times 8 = 608$
POOL1	(14, 14, 8)	1,568	0
CONV2 ($f=5$, $s=1$)	(10, 10, 16)	1,600	416 $\rightarrow (5 \times 5 \times 8 + 1) \times 16 = 3216$
POOL2	(5, 5, 16)	400	0
FC3	(120, 1)	120	48,001 $\rightarrow 120 \times 84 + 1 = 48,120$
FC4	(84, 1)	84	10,081 $\rightarrow 120 \times 10 + 1 = 1,001$
Softmax	(10, 1)	10	841 $\rightarrow 84 \times 10 + 1 = 841$

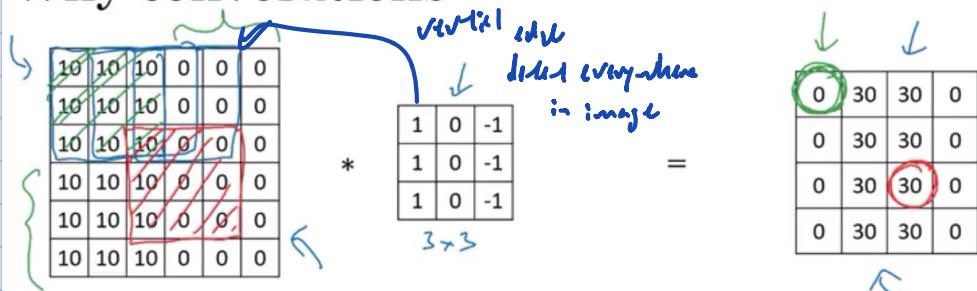
why convolution?

Advantages of convolution over just dense layers:

- ① Parameter sharing
- ② Sparsity of connections



Why convolutions

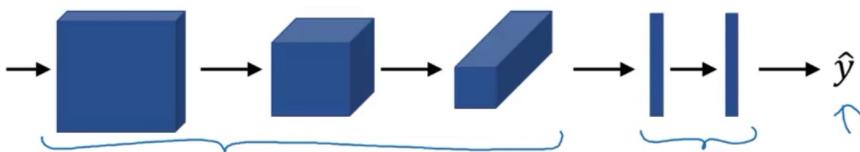


Parameter sharing: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

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

Putting it together

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



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 .

Week 2: Deep Convolutional models: case studies

Some of the networks we will be looking at:

LeNet-5

AlexNet

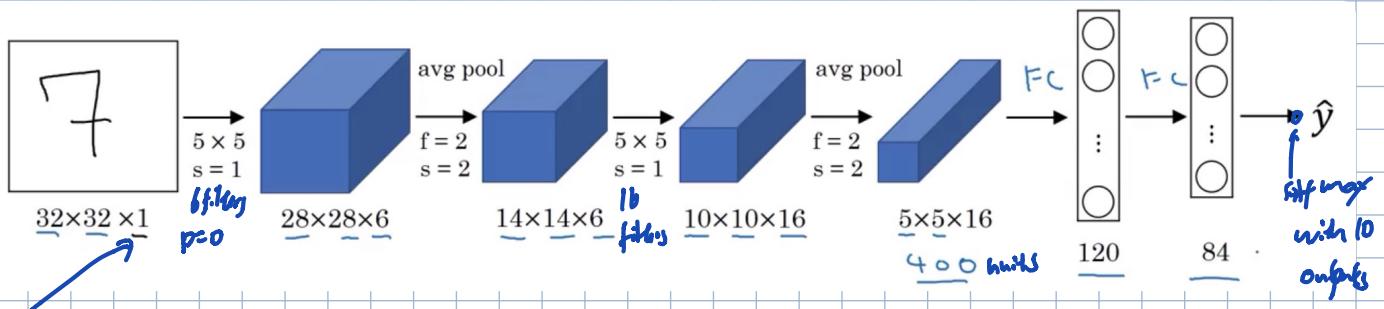
VGG

ResNet

Inception

Classic Networks

LeNet-5



LeNet-5 was trained on grayscale image
has about 60K parameters

back then people use sigmoid / tanh and
not ReLU.

As we progress deeper into NN,

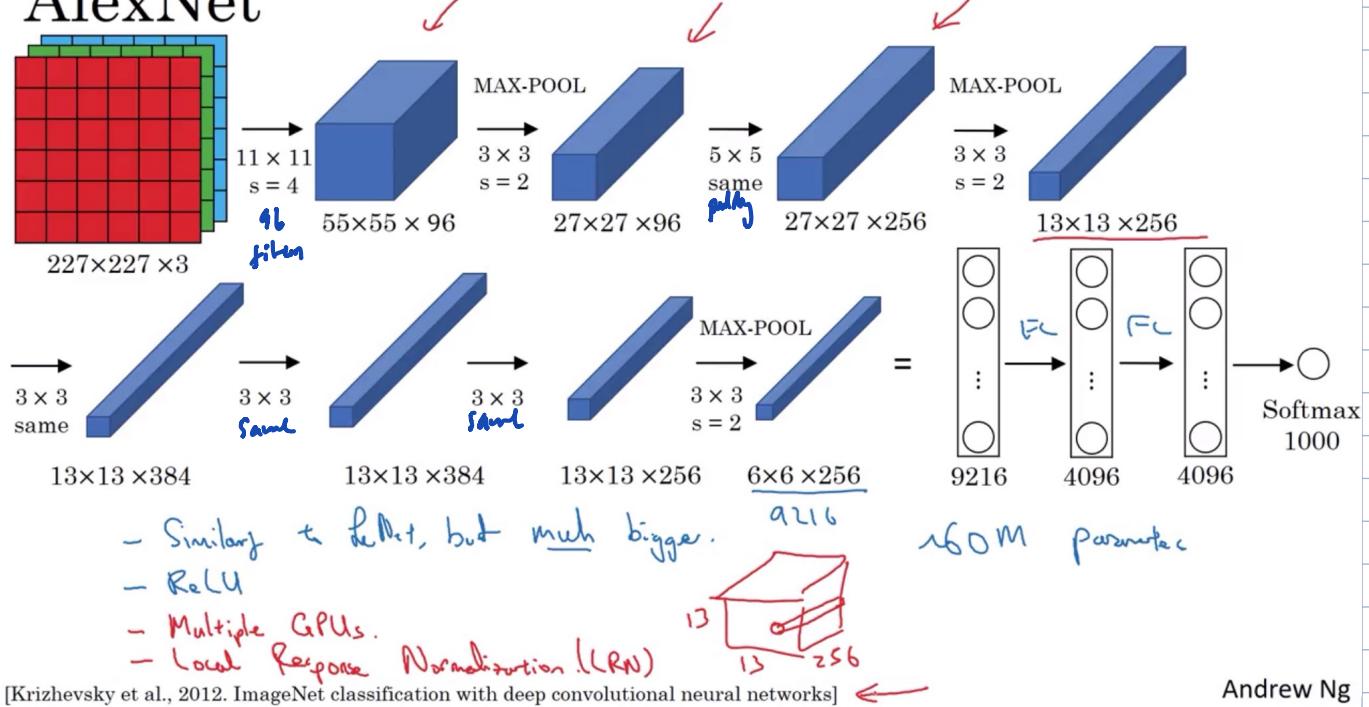
$n_H, n_W \downarrow$ $n_C \uparrow$

a common structure:

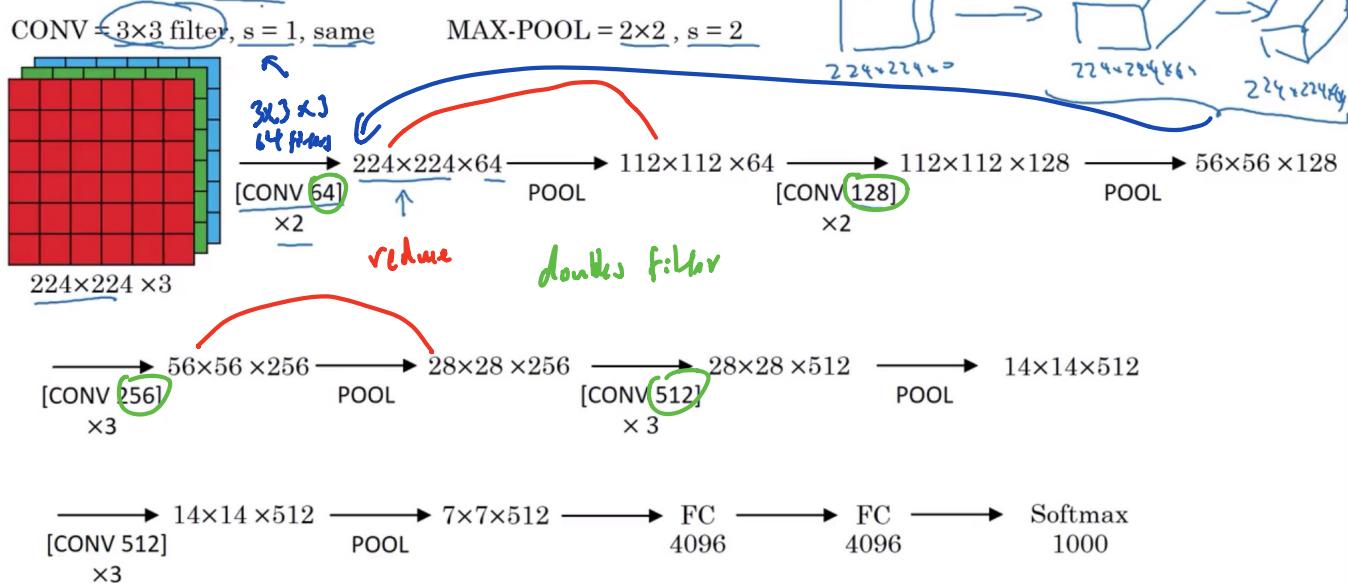
(conv — pool) — conv — pool — fc — fc — output

AlexNet

AlexNet



VGG - 16



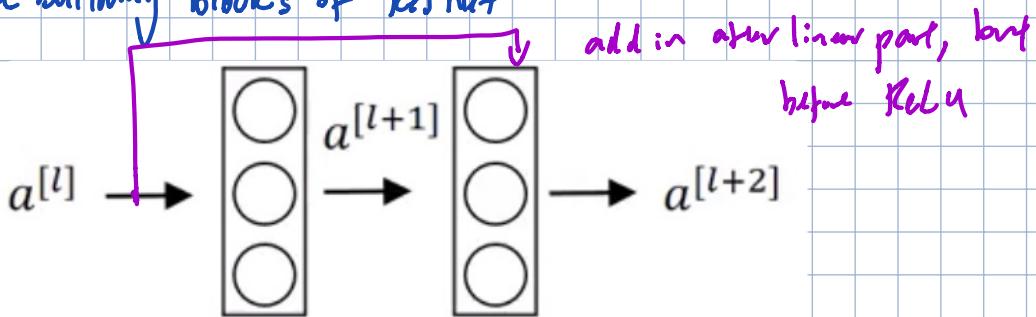
Residual Network (ResNets)

Very deep neural networks are difficult to train due to vanishing and exploding gradients. Here, we learn about skip connections which allows us to take the activations from one layer and feed them to another layer much deeper in NN.

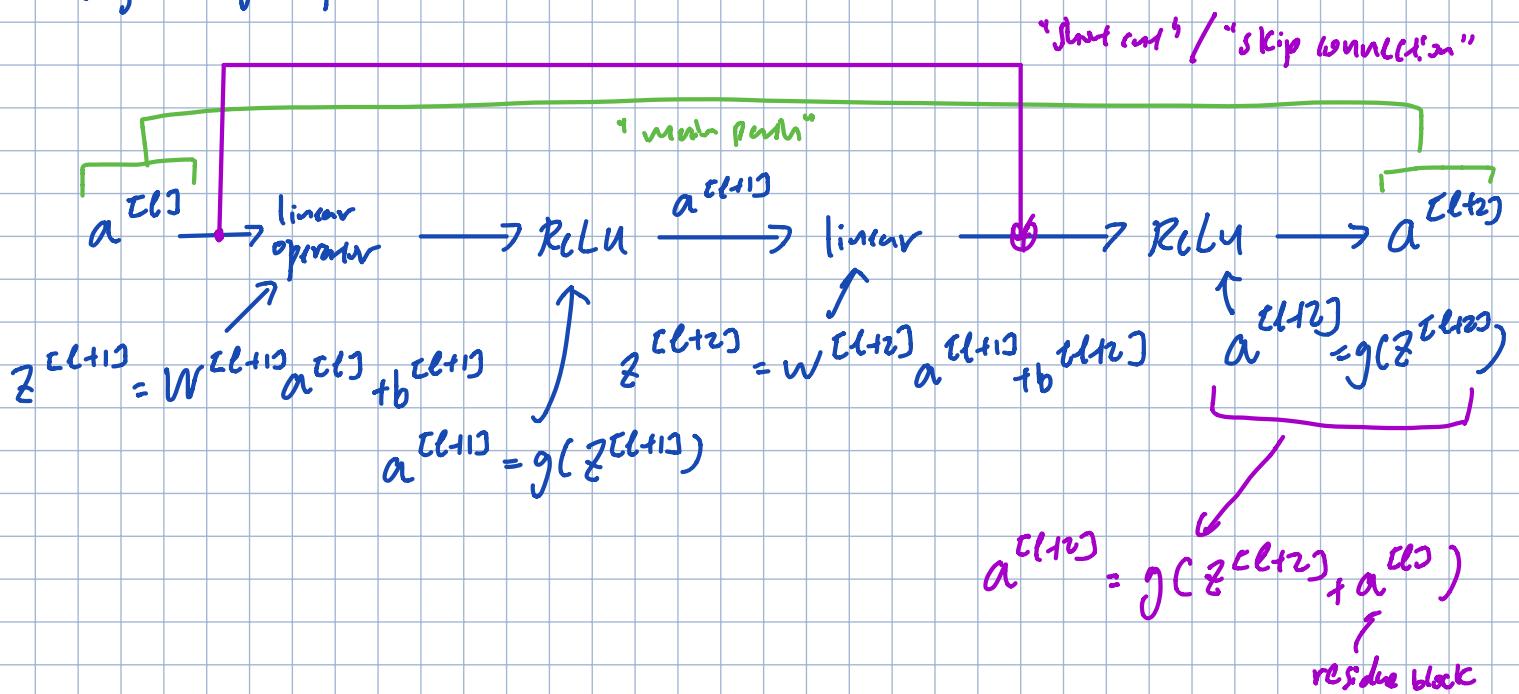
↳ we can build ResNet which we can train very deep networks.

Residual block

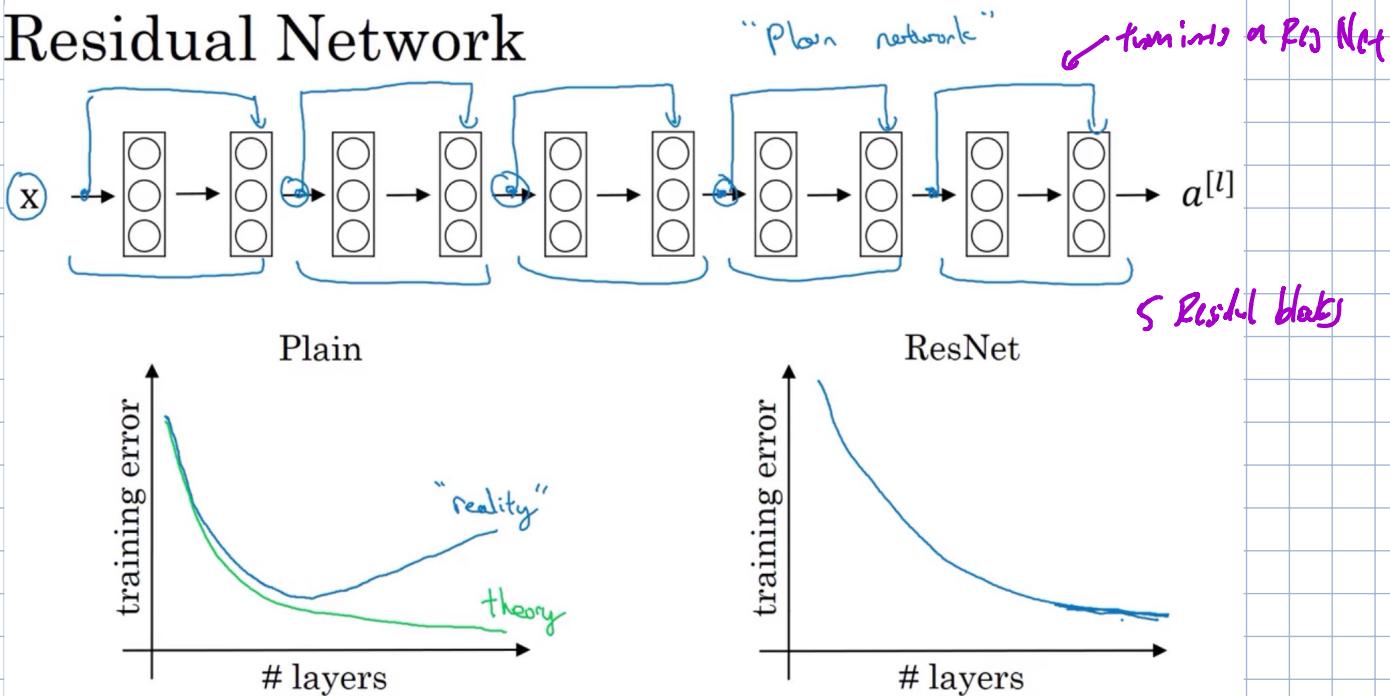
They are the building blocks of ResNet



To go through steps above:

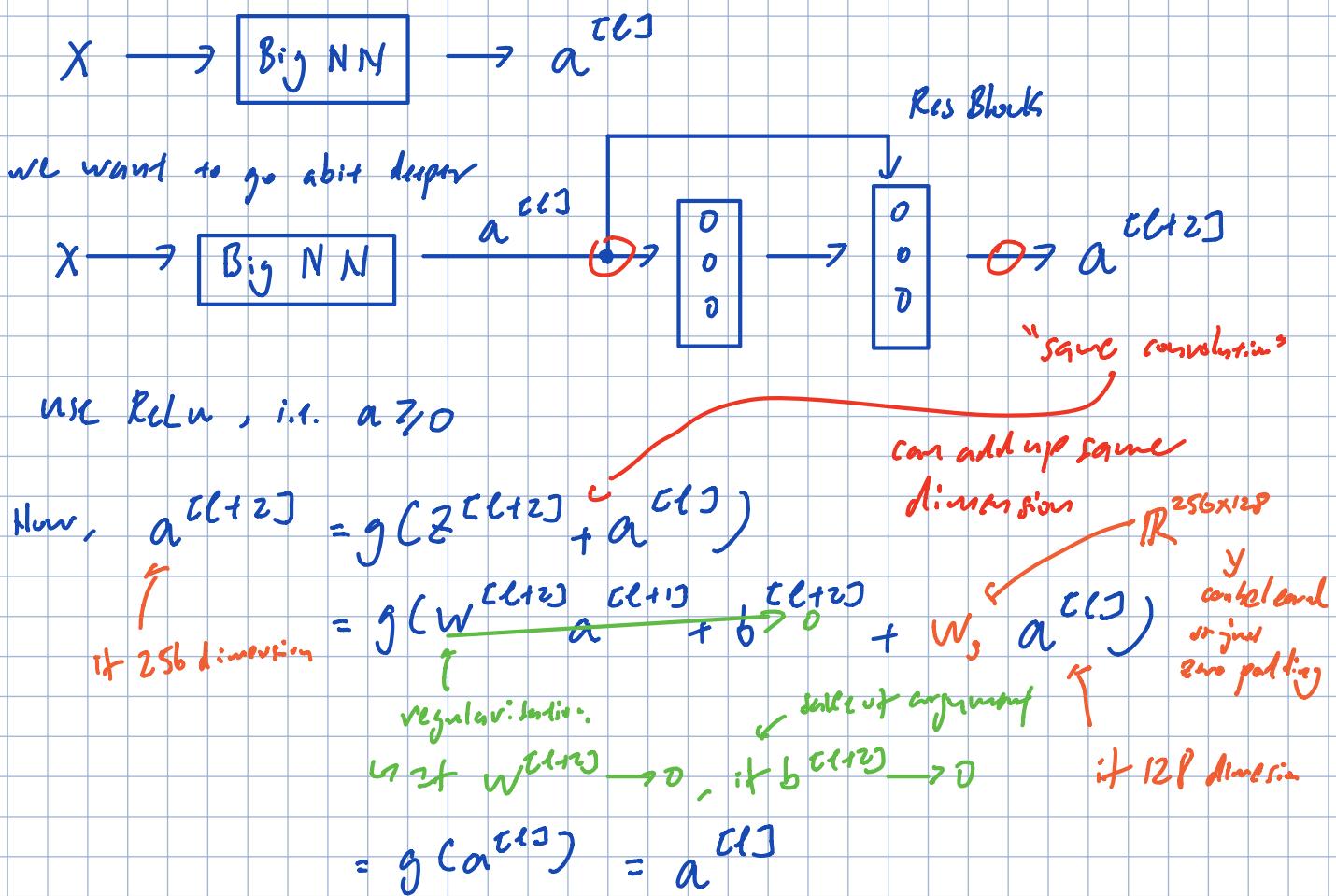


Residual Network



why ResNets work?

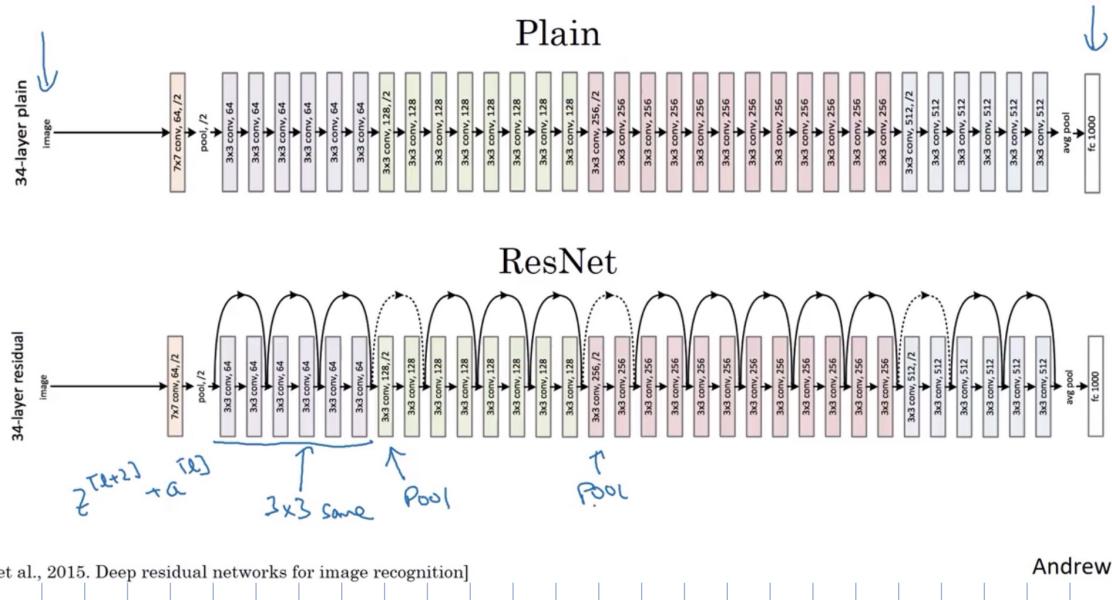
L.9.



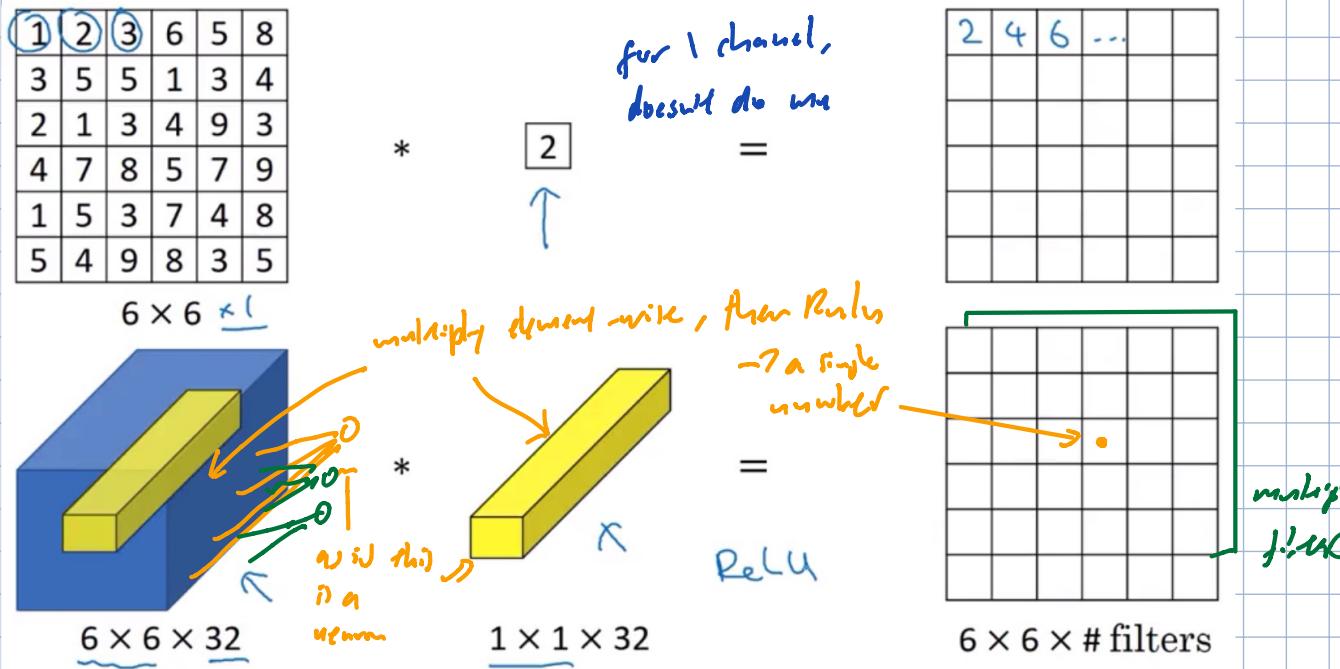
I identity function is easy for Residual block to learn

This means the extra 2 layers does not hurt the ability of a deeper NN.

ResNet



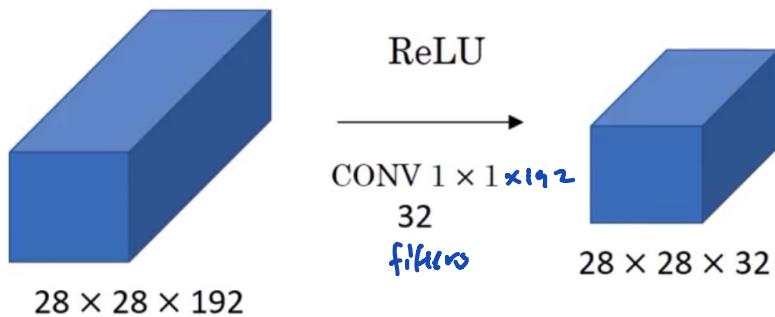
Networks in Networks and 1x1 convolutions



called 1x1 convolution

or Network in Network

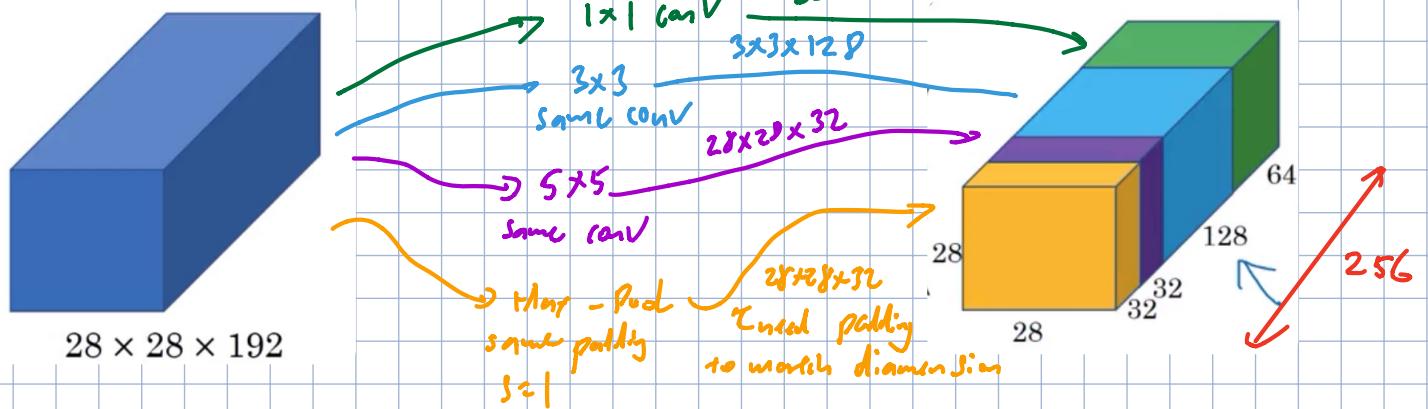
Using 1×1 convolutions



1×1 conv shrink N_C
pooling divide N_H, N_W

say if we want to shrink height, width \rightarrow use pooling layer
if we want to shrink # channels?

Inception Network Motivation



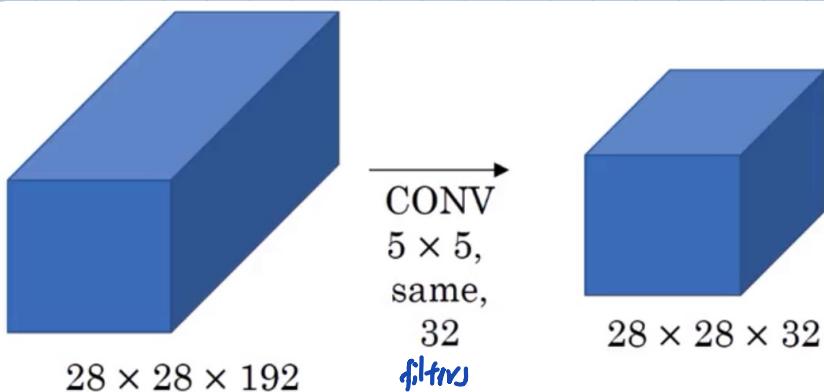
e.g. we input $28 \times 28 \times 192$ volume, the inception layer do is instead of choosing what filter size we want in a conv layer or even having a conv or pooling layer, we can do them all!

We have 1 inception module input $28 \times 28 \times 192$, and output $28 \times 28 \times 256$

↳ Basic idea: instead of pick one of these filter size or pooling, we do them all and concatenate all the outputs. Let the network learn whatever parameters it wants to use, or whatever the combinations of these filter size it wants.

✗ But this description of inference has a problem: computational cost!

e.g. we focus on the 5×5 filter in the example



each filter is $5 \times 5 \times 192$

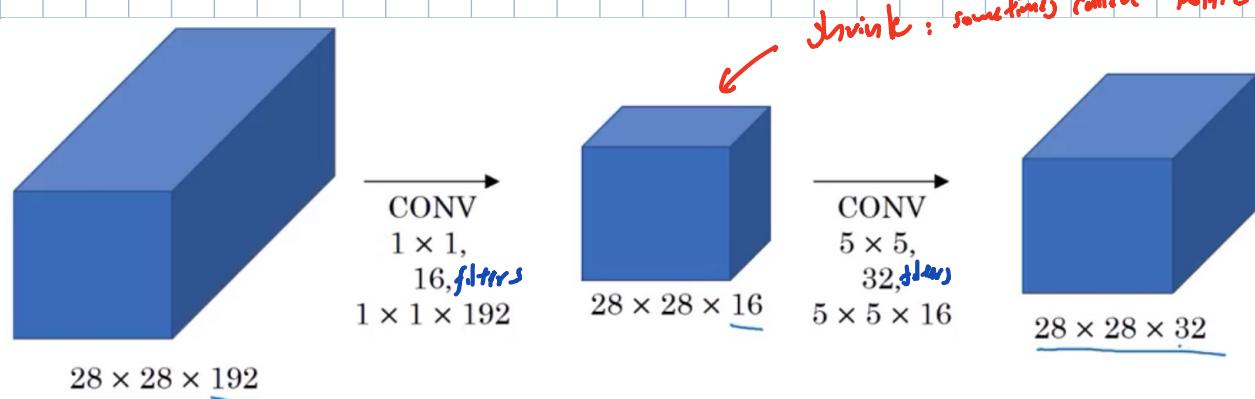
↳ we need to compute

$28 \times 28 \times 32$ numbers

↳ for each number, need
 $5 \times 5 \times 192$ multiplications

↳ total 120 million multiplications

Using 1×1 conv, we can reduce # of multiplications:



Notice the output dimension are the same for the 2nd approach

To apply 1×1 conv, 16 filters, each is $1 \times 1 \times 192$

$$\underbrace{28 \times 28 \times 16}_{\# \text{ of outputs}} \times (192 \times 1 \times 1) \approx 2.4 \text{ mil} \quad \leftarrow \text{first layer}$$

of outputs

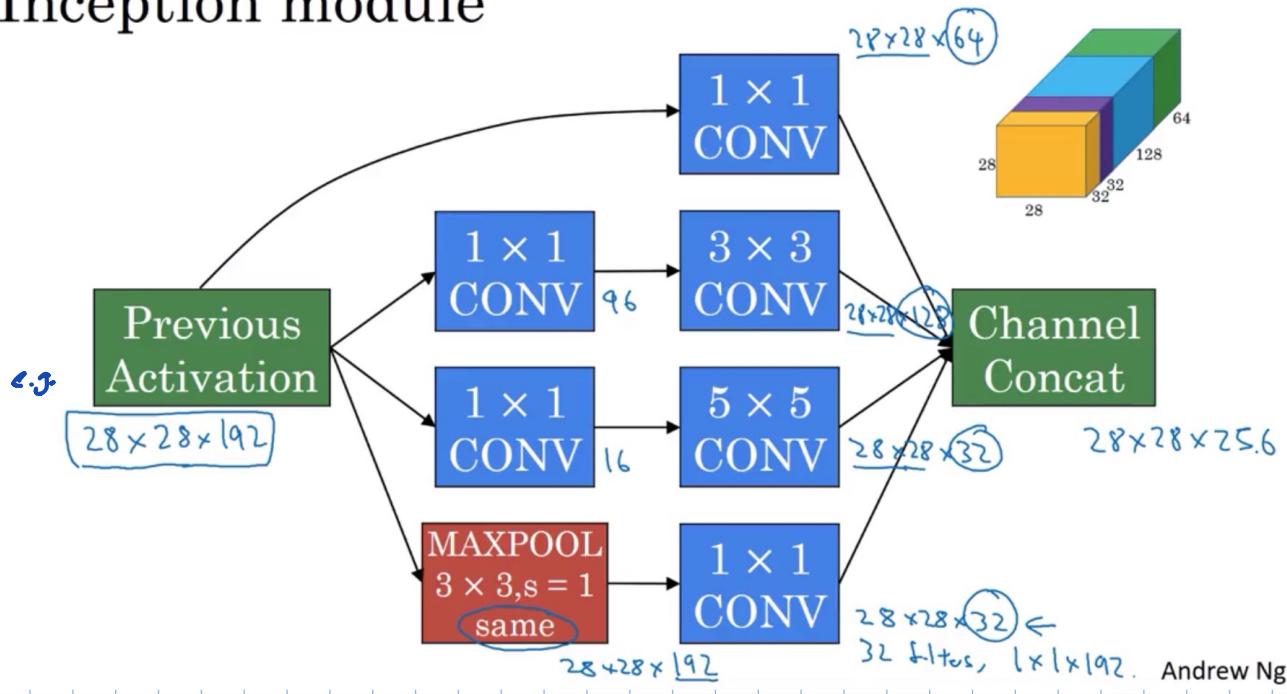
$$\underbrace{28 \times 28 \times 32}_{\# \text{ of outputs}} \times (5 \times 5 \times 16) = 10.0 \text{ mil} \quad \leftarrow \text{second layer}$$

} 12.4 mil computations

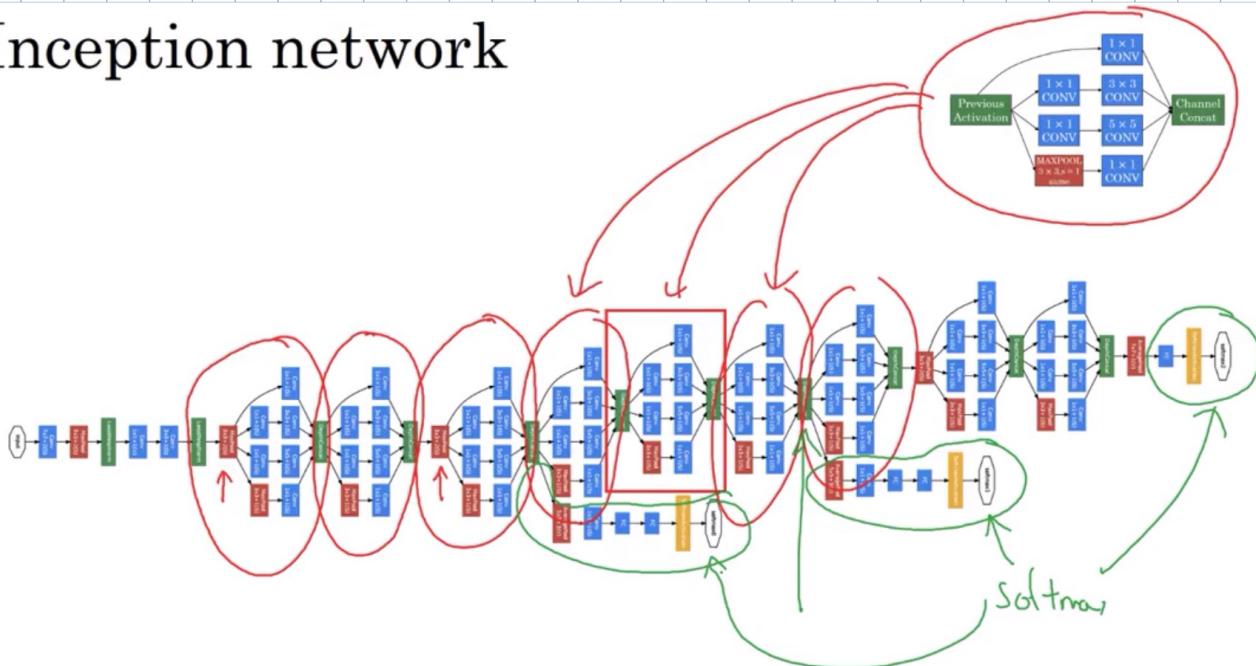
in total

Inception Network

Inception module



Inception network

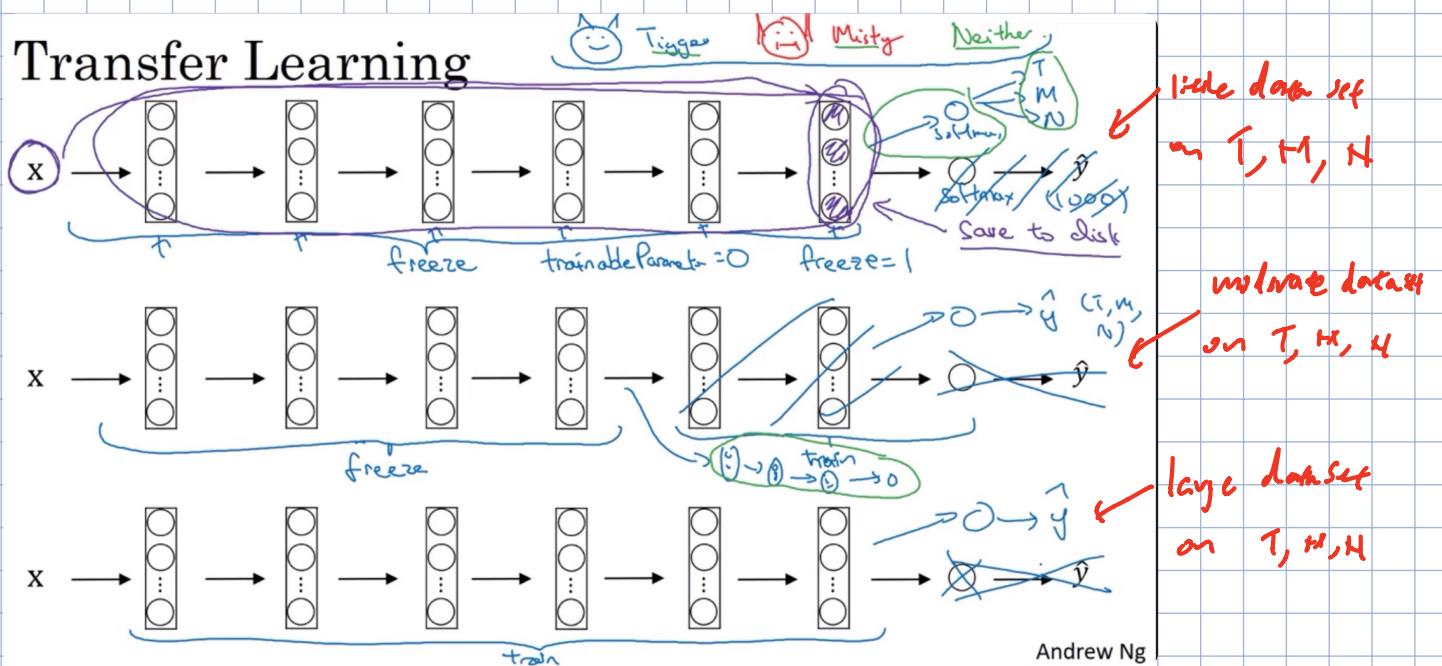


Szegedy et al., 2014, Going Deeper with Convolutions]

Andrew Ng

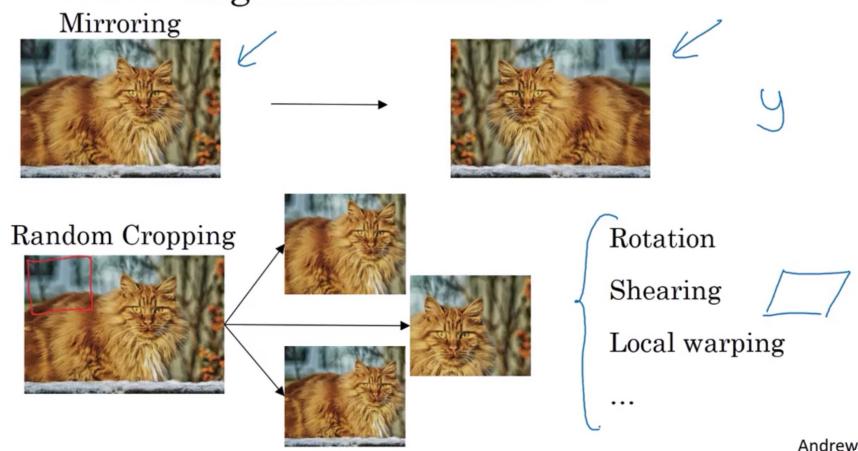
Transfer Learning

Transfer Learning

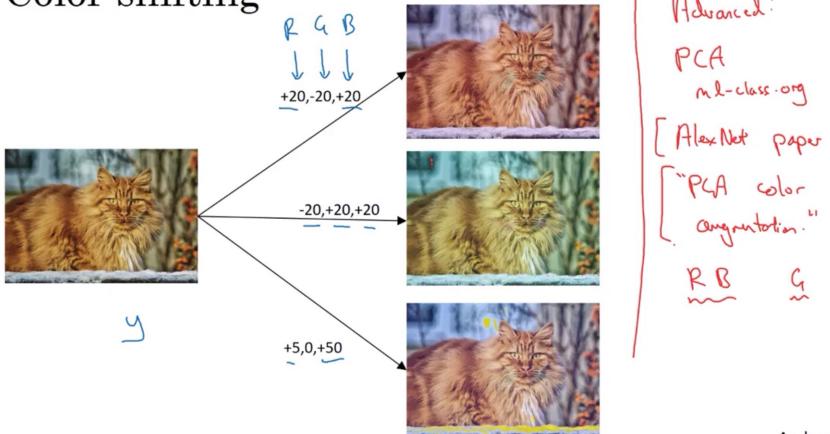


Data augmentation

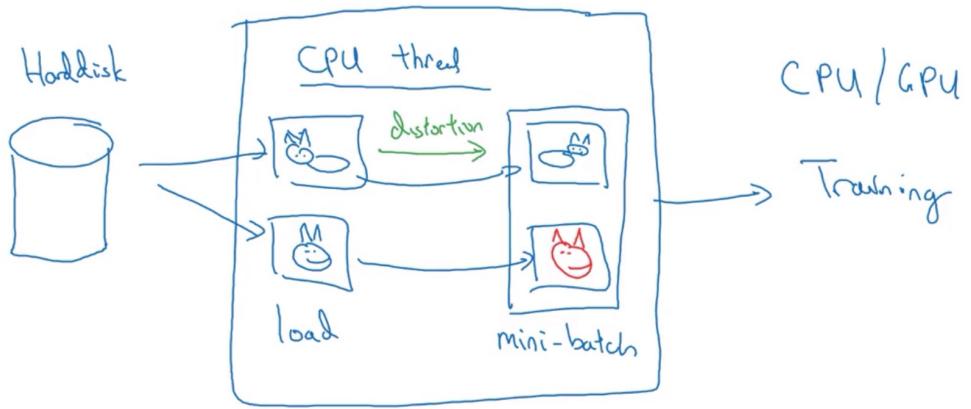
Common augmentation method



Color shifting

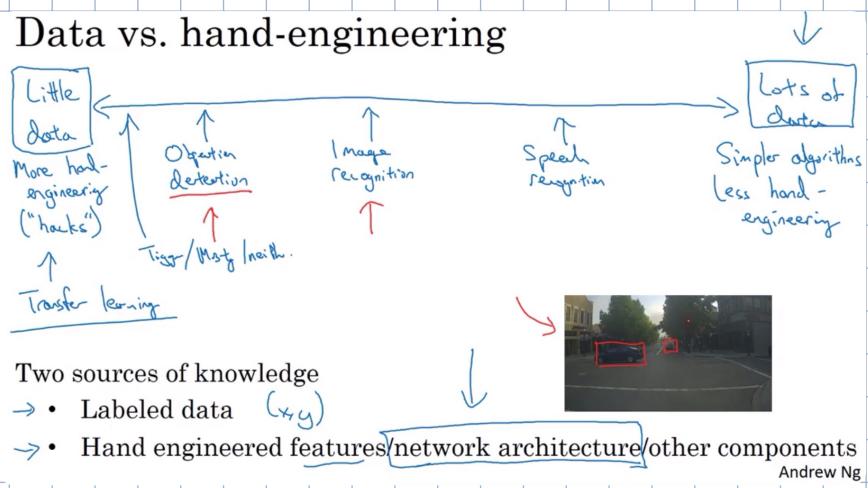


Implementing distortions during training



State of Computer Vision

Data vs. hand-engineering



Tips for doing well on benchmarks/winning competitions

Ensembling
• Train several networks independently and average their outputs

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

Multi-crop at test time

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

10-crop



Use open source code

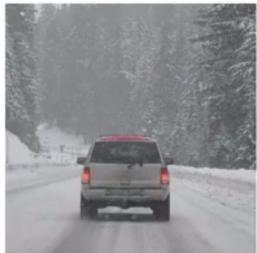
- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset

WEEK 3 : Object Detection

Object Localisation

What are localization and detection?

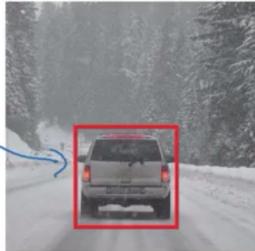
Image classification



"Car"

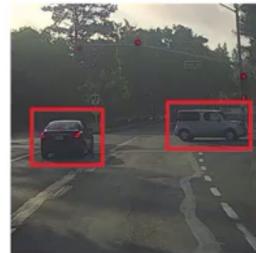
1 object

Classification with localization



"Car"

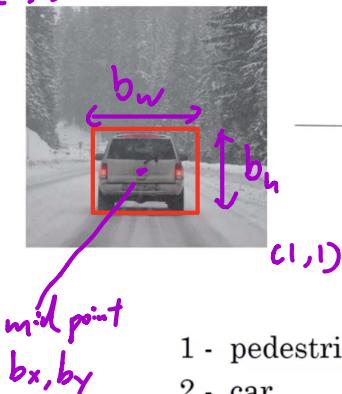
Detection



multiple objects

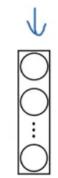
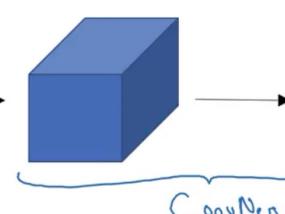
Classification with localization

(c_{obj})



(1,1)

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



softmax (4)

4 objects

e.g. here: $b_x = 0.5$
 $b_y = 0.7$
 $b_h = 0.3$
 $b_w = 0.4$

b_x, b_y, b_h, b_w for
bounding box

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)

logistic regression

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

P_c : is there any object?
 if object is class 1, 2, 3, then $P_c=1$
 if background or none, $P_c=0$
 Probability that there is one object (1 of the 3)

* assume only 1 object

In practice, can use logistic loss

c-1.

$$x =$$



$$y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x =$$



$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

don't care

Loss Function

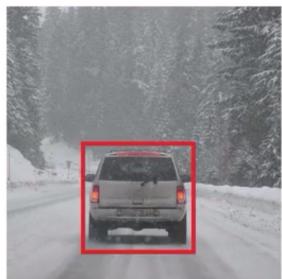
$$\begin{aligned} \mathcal{L}(\hat{y}, y) &= (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_r - y_r)^2 \quad \text{if } y_i = P_i = 1 \\ &= (\hat{y}_1 - y_1)^2 \quad \text{if } y_i = P_i = 0 \end{aligned}$$

(!! y has 8 components here)

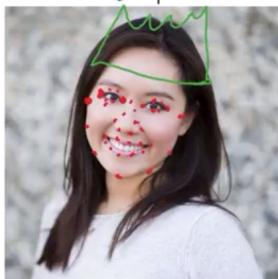
Landmark Detection

In general, we can have a NN just output x and y coordinates of important points on image (sometimes called landmark), that we want the NN to recognise.

Landmark detection



b_x, b_y, b_h, b_w



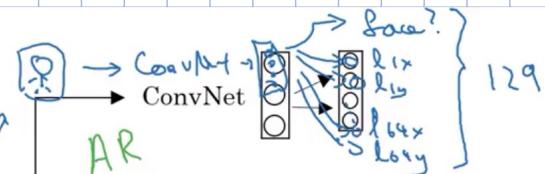
$l_{1x}, l_{1y},$
 $l_{2x}, l_{2y},$
 $l_{3x}, l_{3y},$
 $l_{4x}, l_{4y},$
 \vdots
 l_{64x}, l_{64y}

} coordinates
 X, Y



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

Andre



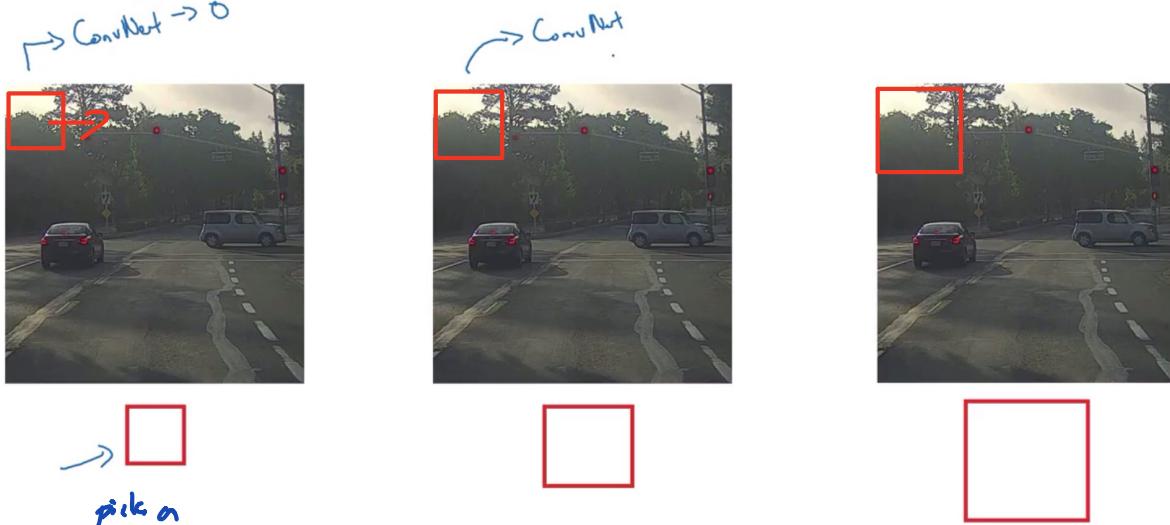
Landmarks has to be consistently defined across all images during labelling.

Object Detection

Car detection example



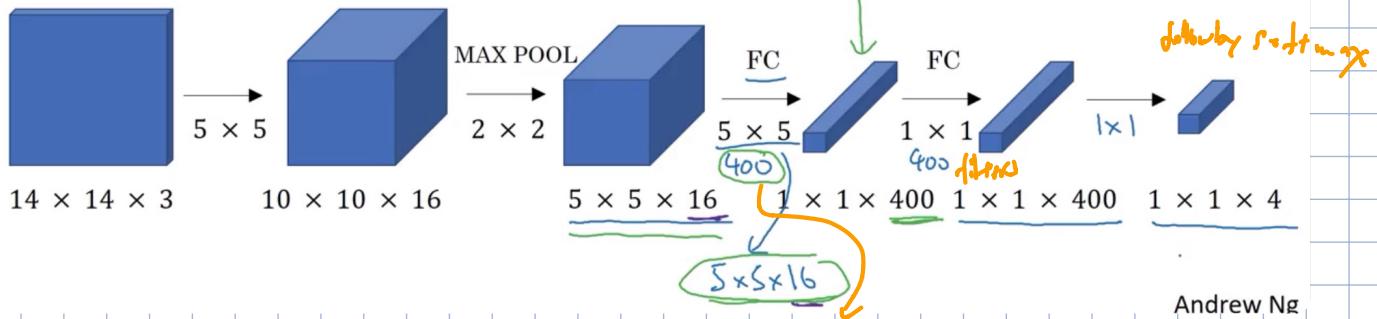
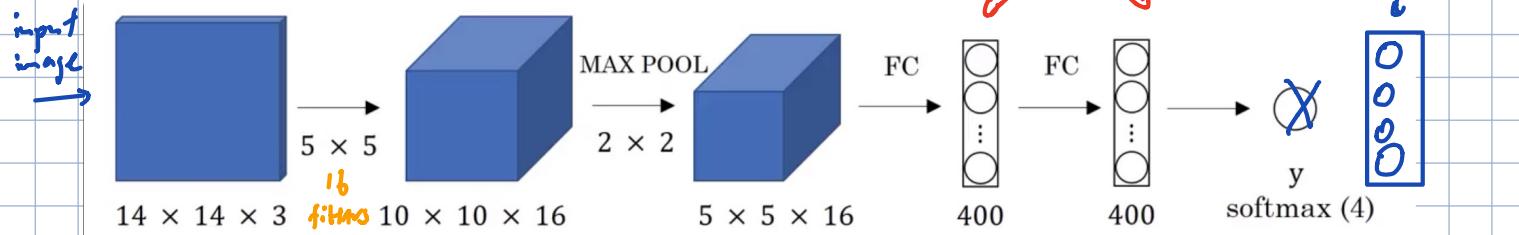
Sliding windows detection



Disadvantage: computational cost

Convolutional Implementation of Sliding Windows

Turn FC layer → convolutional layers

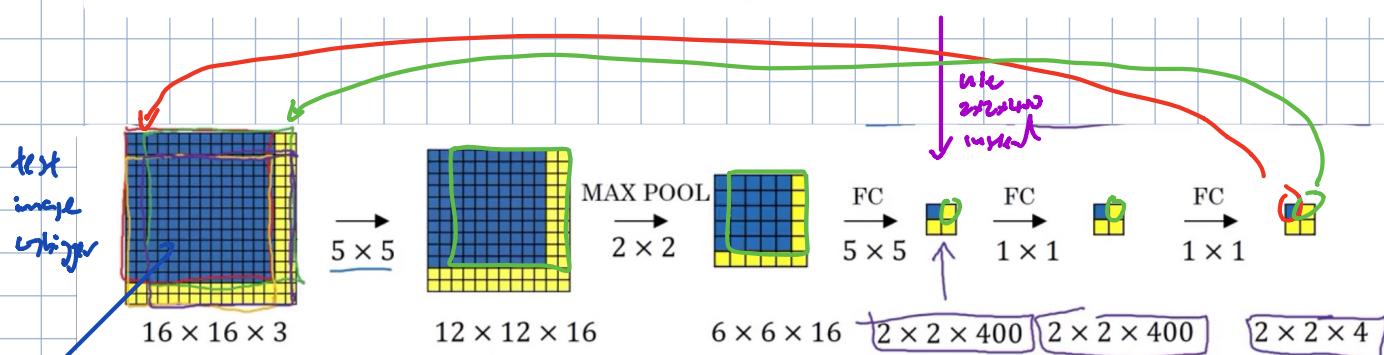
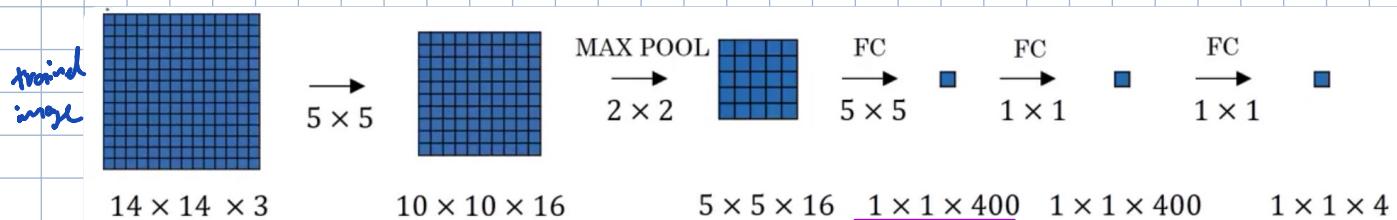


Andrew Ng

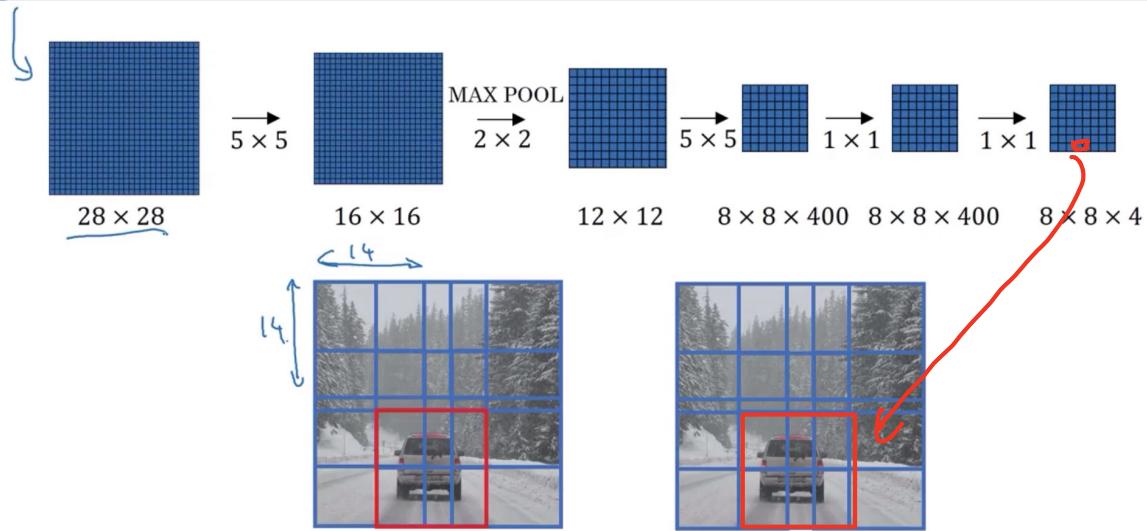
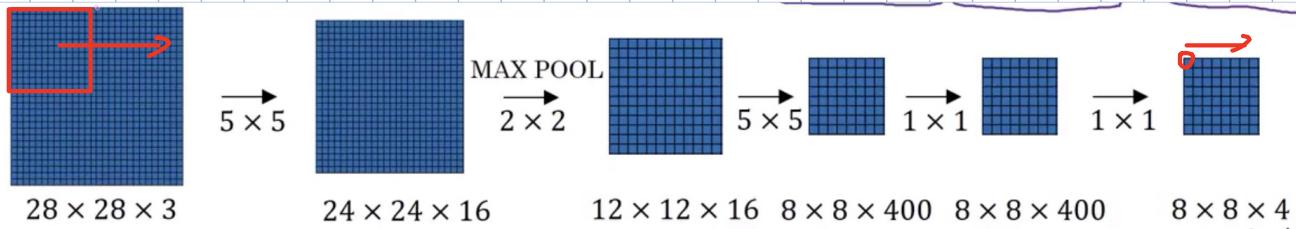
400 filters
5x5x16 (note of all 16 channels)

Convolution implementation of sliding windows

only front is drawn, they are volumes



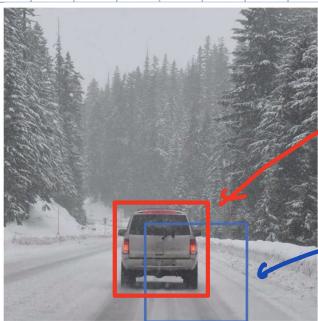
bigger test e.g.



Instead of doing sequentially, to slide the window,
we can implement entire image - and make all predictions at
the same time , i.e. 1 formal pass to predict (after training)
instead of each sliding window 1 formal pass to predict

However, the position of the bounding boxes are not too accurate

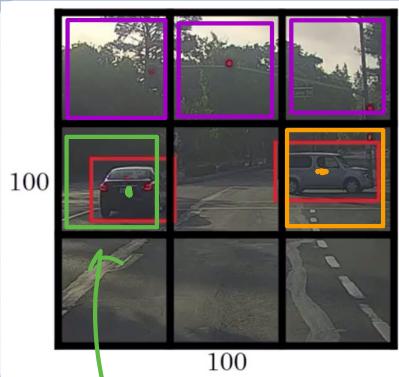
Bounding Box Prediction



ideally, should be this box, also may not be a square, best box might be a rectangle.

using sliding window, this might be the best box

YOLO algorithm



$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{pos of box} \quad \text{car}$$

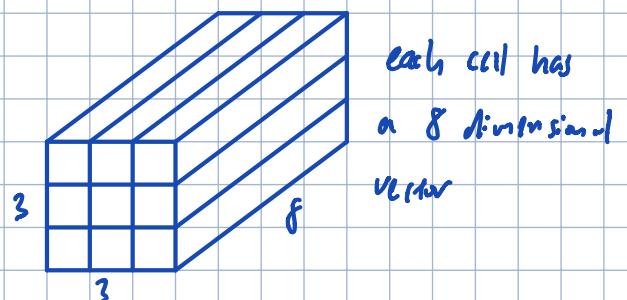
For each grid cell, we specify label y :

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix} \xrightarrow{p_c = 0 \text{ or } 1}$$

3 classes
(cat, car, background)

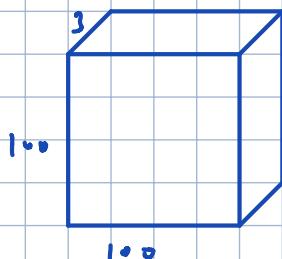
i.e. here we have a vector like this for each grid cell

Total output volume: $3 \times 3 \times 8$ (target output)

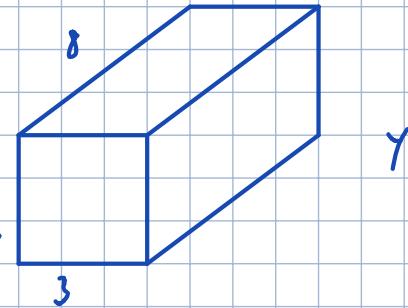


1 object is assigned to 1 grid cell (based on the location of the midpoint)

To train our NN, input image $100 \times 100 \times 3$

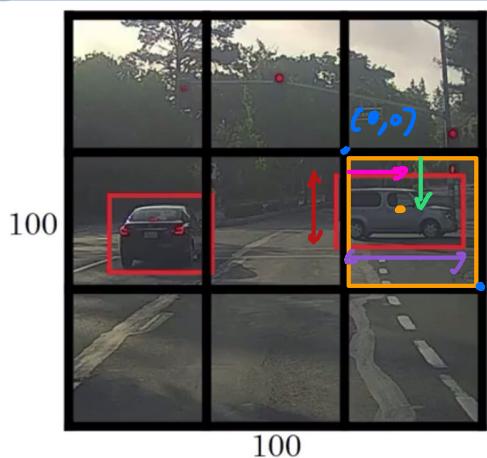


\rightarrow conv \rightarrow maxpool $\rightarrow \dots \rightarrow$



match target value

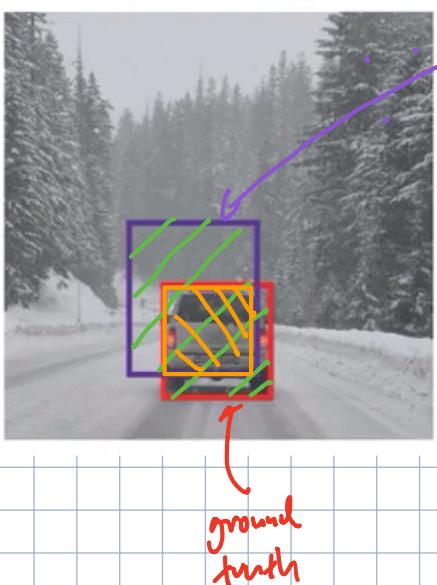
How do we encode bounding box (bx, by, bh, bw)



$$y = \begin{bmatrix} 1 \\ bx \\ by \\ bh \\ bw \\ 0 \\ 1 \\ 0 \end{bmatrix}_{(1,1)}$$

between 0 and 1
specified relative to the grid cell
 y can be larger than 1

Intersection over Union



predicted
bounding
box

ground
truth

Intersection over Union (IoU)

$$\approx \frac{\text{size of } \cap}{\text{size of } \cup}$$

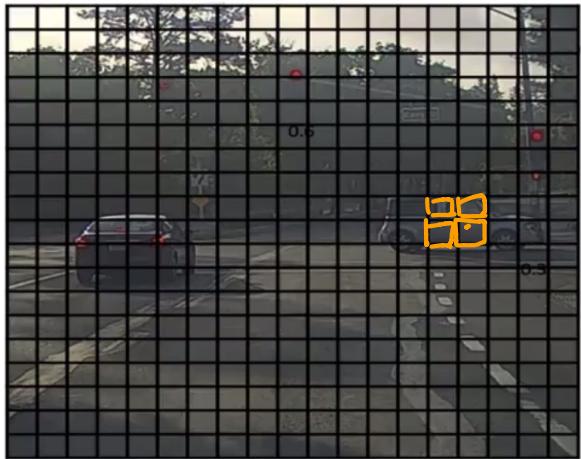
\hookrightarrow "convert" if $IoU \geq 0.5$

\hookrightarrow used to judge if the bounding box is correct or not

More generally, IoU is a measure of the overlap between 2 bounding boxes.

Non-max suppression

One of the problems of object detection is that algorithm may find multiple detections of the same objects.



multiple boxes think they got the center point.



multiple grid cells say my P_C is large
non-max suppression looks at the probability first

↪ take the largest one

↪ looks at remaining box and all those with high IoU

↪ suppress those box

Non-max suppression algorithm



19

e.g. $19 \times 19 \times 5$ output volume

for simplification

each output prediction is:

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

① Discard all boxes with $P_c \leq 0.6$, i.e. throw all low probability box

② while there are any remaining boxes:

 ↳ pick the box with largest P_c

 ↳ Output that as prediction

 ↳ discard any remaining box with $I_{IoU} \geq 0.5$ with the
 box output (in the previous step)

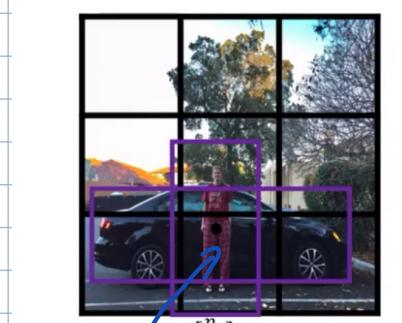
→ boxes are either high P_c , or discarded

Anchor Boxes

what if a grid cell want to detect multiple objects?

↳ use anchor boxes

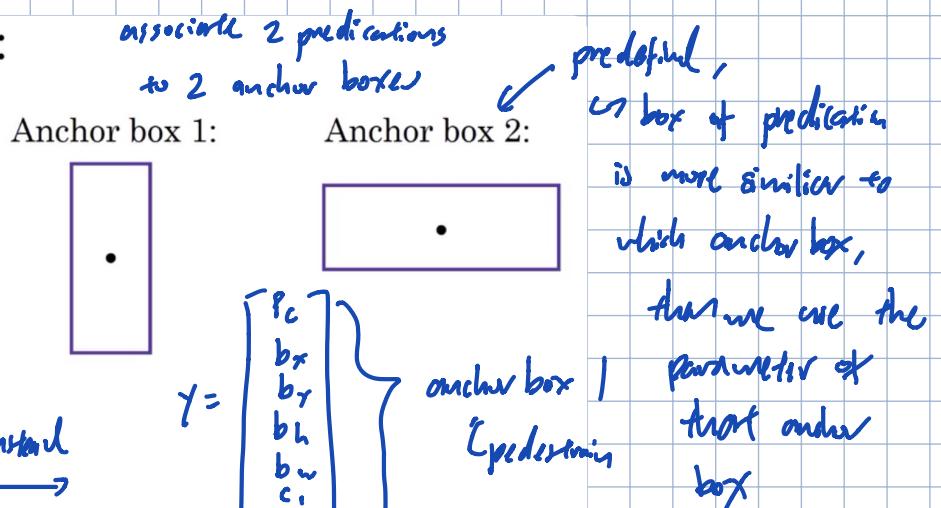
Overlapping objects:



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

for this grid cell
only output more than 1 class
3 classes

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



notice both midpoints of car and person are in the same grid cell.

Anchor box algorithm

Previously:

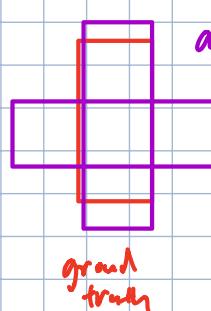
Each object in training image is assigned to grid cell that contains that object's midpoint.

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.



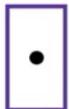
→ object assigned to
(grid cell, anchor box)
pair

output y :
 $3 \times 3 \times 16$
 $\rightarrow 2 \times 8$
2 bytes

Anchor box example



Anchor box 1: Anchor box 2:



pedestrian more similar
to an anchor 1

car more similar to,
another box 2

$$y =$$

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

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

car only?

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

pedestrian
car

This algorithm don't handle well:

- ① 2 anchor box, 3 object in the same grid cell } need tiebreaker
- ② 2 object, same anchor box shape in same grid cell }

↳ quite rare if use finer grid cell.

↳ anchor box determinization

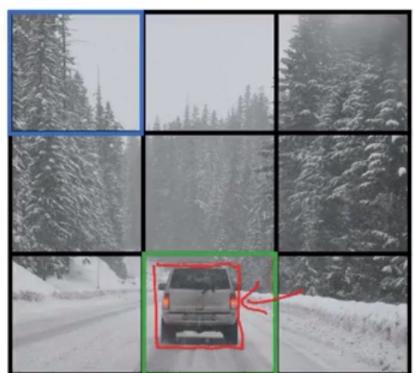
↳ used to be by hand

↳ now use k-means

YOLO Algorithm

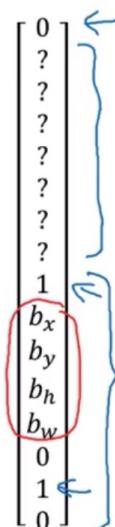
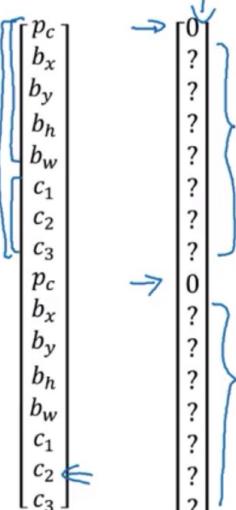
want to train to detect
assume 2 another box

Training



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

$$y =$$



more similar to
audax Lox 2

y is $\underbrace{3 \times 3}_{\text{ }} \times 2 \times 8$

$19 \times 19 \times 16$

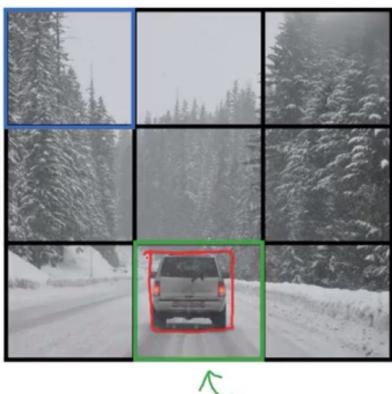
$$\uparrow \quad \overline{A} \quad S + \# \text{classes}$$

#anchors

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

För
gräsl

Making predictions

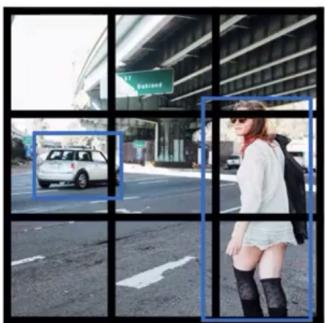


$$y =$$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \xleftarrow{\text{O}} \quad \begin{bmatrix} 0 \\ \dots \\ 0 \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ 0 \end{bmatrix} \xleftarrow{\text{O}}$$

Andrew Ng

Outputting the non-max suppressed outputs



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



Region Proposals

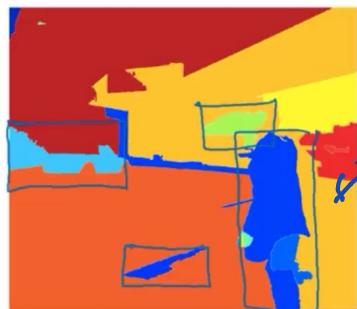
Region proposal: R-CNN



sliding window



CNN
most box
has no object



Segmentation algorithm
~2,000 blobs
R-CNN

Faster algorithms

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

Fast R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions. ↪

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

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

[Girshik, 2015. Fast R-CNN]

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

WEEK 4: Face Recognition & Neural Style Transfer

What is face recognition?

Face verification vs. face recognition

→ Verification

- Input image, name/ID
- Output whether the input image is that of the claimed person

$$\begin{array}{r} 1:1 \\ \hline 99\% \end{array}$$

99.0

→ 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")

$$\begin{array}{r} 1:K \\ \hline K=100 \end{array}$$

One shot Learning

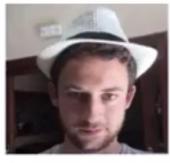
For most face recognition applications, we need to recognise a person given just 1 single image → 1 shot learning



model
only seen 1 image



same person?



Learning from one example to recognise
the person again

we can try:



→ CNN → O
softmax

say 4 employee in DB

↑
people want
to pass gate

↑
training image size
(only 1) too small
(5 output)

Instead, we are going to let the model learn a "similarity" function

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

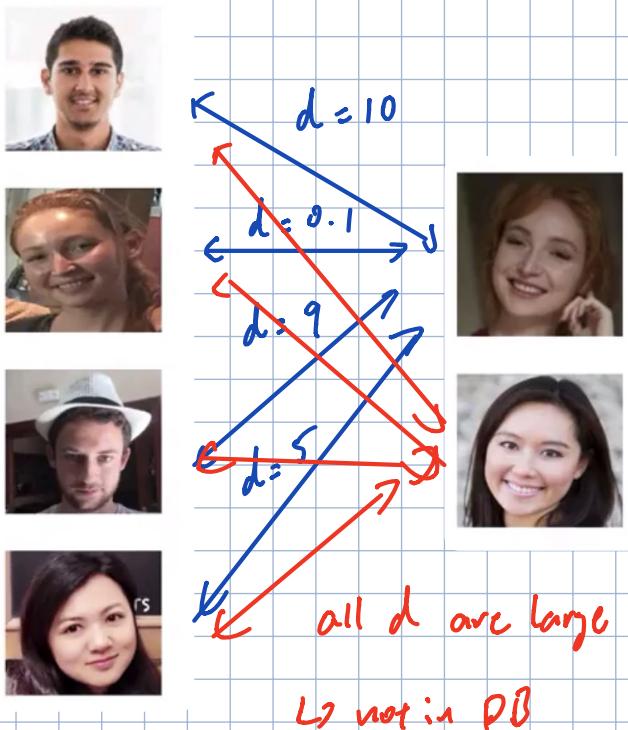
Small number: same person

Big number: different person

During recognition time, it

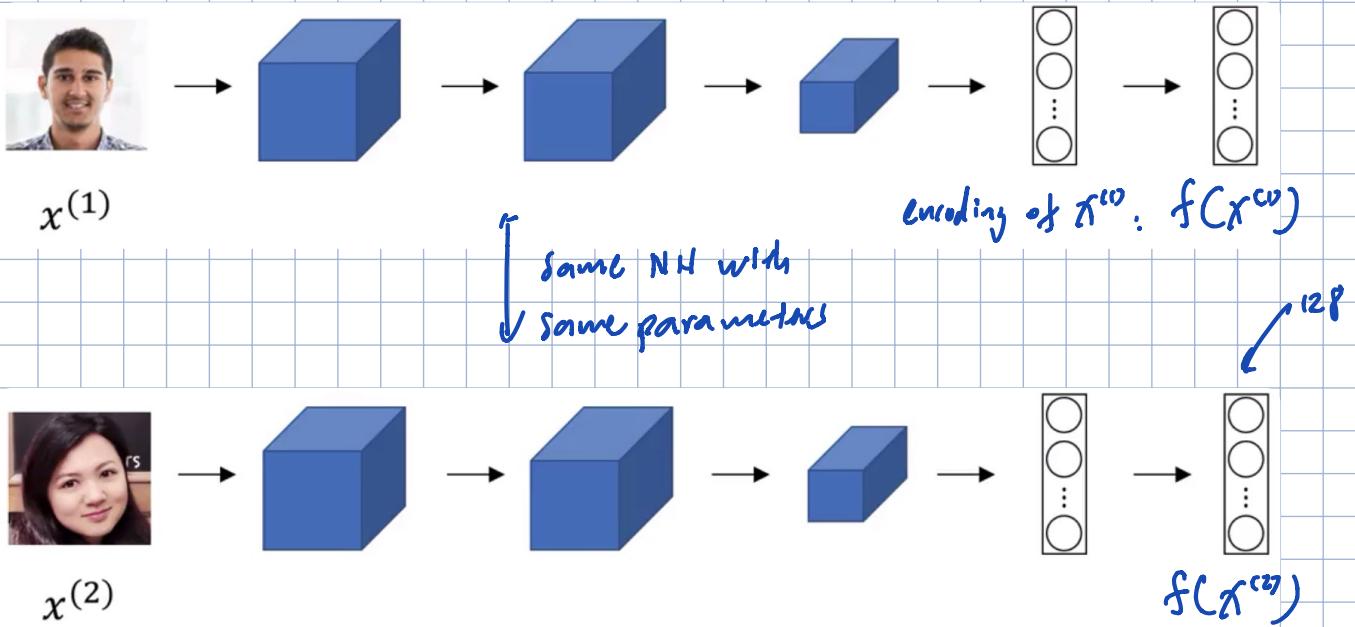
$$d(\text{img1}, \text{img2}) \leq T \quad \begin{array}{l} \rightarrow \text{same person} \\ \rightarrow \text{different person} \end{array} \quad \left. \begin{array}{l} \text{face} \\ \text{verification} \\ \text{task} \end{array} \right\}$$

For face recognition task:



we have to learn
 $d(\text{img1}, \text{img2})$

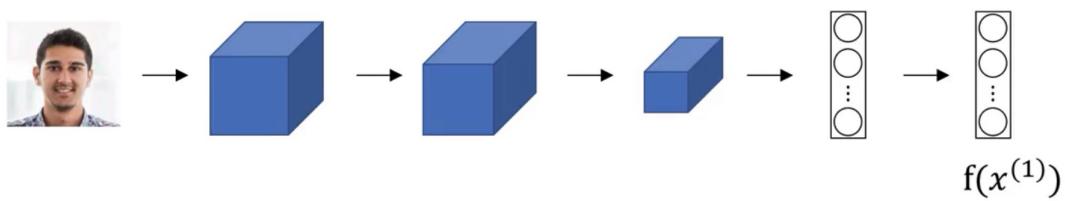
Siamese Network



we have to make sure these 128 units are good representation of the pictures

$$\text{then : } d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

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.

use back prop, get parameters to make these are small
 \therefore different parameters
 → different encoding

Triplet Loss

To apply the triplet loss, we need to compare pairs of images



Anchor **A** Positive **P**



Anchor **A** Negative **N**

we want encoding to
be same

T /
we want something to
be different

We always look at 3 images at a time:

↳ anchor, positive, negative image

$$\left\| f(A) - f(P) \right\|^2 - \left\| f(A) - f(N) \right\|^2 \leq 0 \Leftrightarrow P = N$$

✓ NH can actually make the 2 terms all 0
 $\Rightarrow f(\text{img}) = \vec{0}$) trivially satisfy this equation

↳ To prevent this from happening, we need to modify.

Home :

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \delta \leq 0$$

\uparrow
margin

Triplet Loss Function

This loss function is defined on triplets of images:

Given 3 images: A, P, N

single triplet example

$$L(A, P, N) = \max \left(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \Delta, 0 \right)$$

so

so long we achieve the objective of making that ≤ 0 , then loss = 0
if that thing > 0 , we will have a positive loss

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

e.g. training set: 10K pairs of 1K persons

take 10K pairs to generate / select the triplets, then train NN

↳ we need multiple pic of the same person to train this model
∴ need (A, P) pair

How do we choose these triplets to form the training set?

During training, if A, P, N are chosen randomly,

$d(A, P) + \Delta \leq d(A, N)$ is easily satisfied.

↑

$$\|f(A) - f(P)\|^2 + \Delta \leq \|f(A) - f(N)\|^2 \quad \leftarrow \begin{array}{l} \text{dissimilarity } A, N \text{ are much} \\ \text{different} \end{array}$$

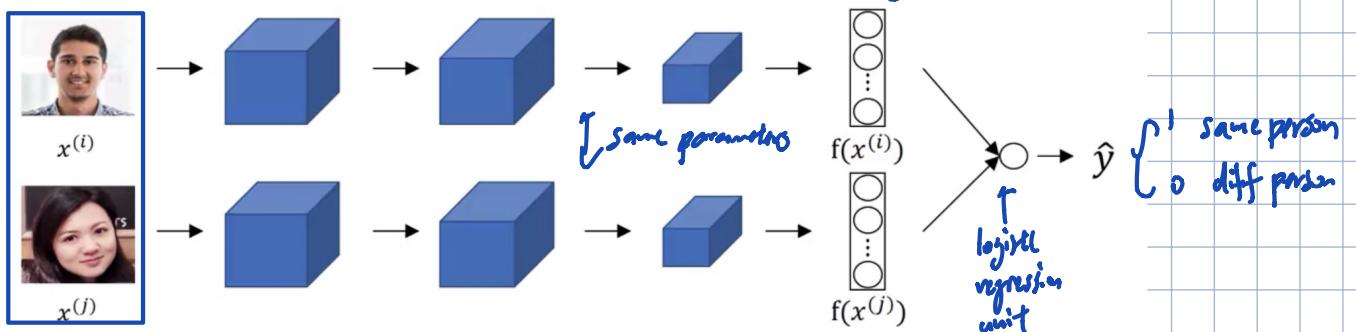
↳ choose triplets that are "hard" to train on.

$$d(A, P) + \Delta \leq d(A, N)$$

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

Face verification and Binary classification

Learning the similarity function



$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

element-wise

$$\text{or}$$

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$$

$x^{(i)}$

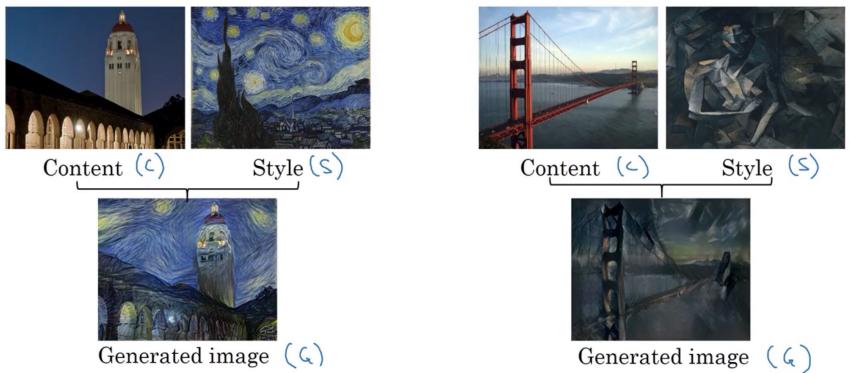
Face verification supervised learning

x	y
	1 "Same"
	0 "Different"
	0
	1

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

Neural Style Transfer

Neural style transfer

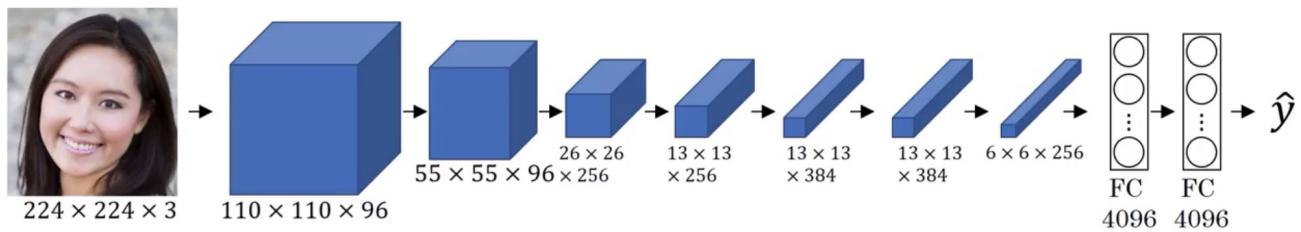


[Images generated by Justin Johnson]

We need to look at the features extracted by ConvNet at various layers.

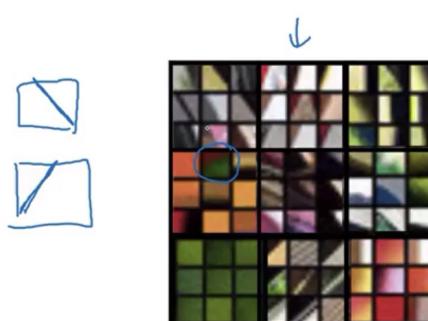
What are deep ConvNets learning?

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.



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

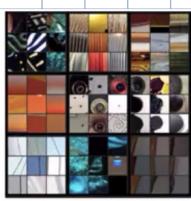
Andrew Ng

Later layers see larger patches

1st layer sees small patches
highly activated nodes



Layer 1



Layer 2



Layer 3

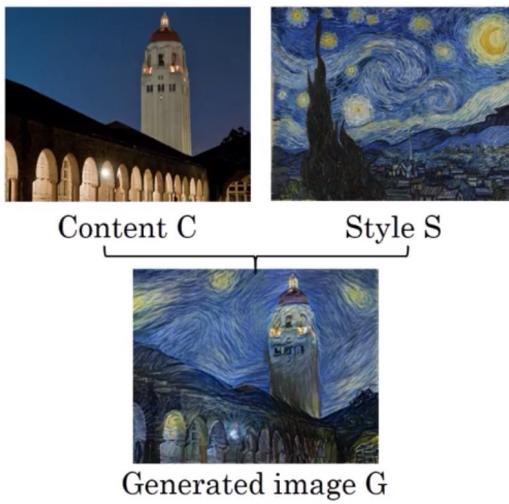


Layer 4



Layer 5

Neural Style Transfer: Cost function



$J(G) \rightarrow$ a measure of how good is the image.

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

$\uparrow \quad \uparrow$
how similar in content

$\uparrow \quad \uparrow$
how similar in style

Find the generated image G

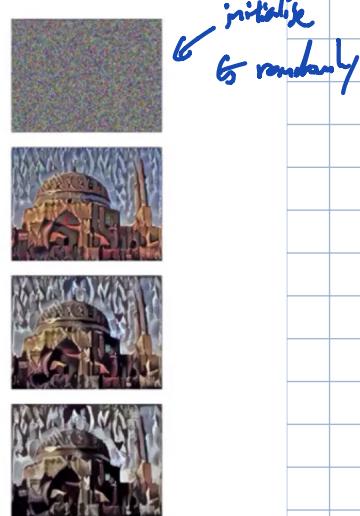
1. Initiate G randomly

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

\uparrow
RGB

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

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



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

Andrew Ng

Content Cost Function

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

- Say you use hidden layer l to compute content cost.
- Use pre-trained ConvNet. (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

usually ℓ is not too small or too large

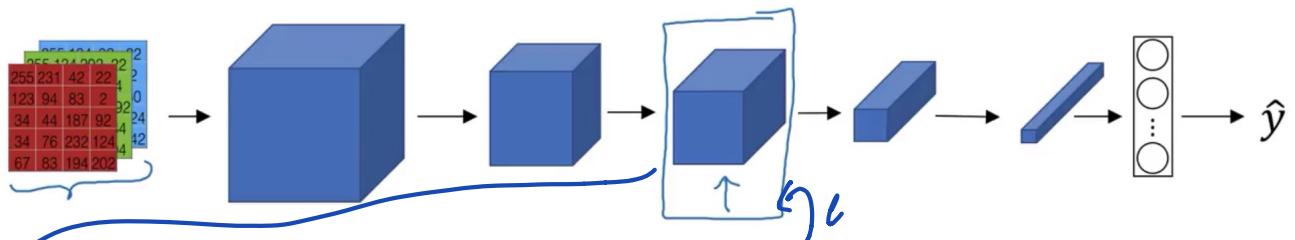
$$J_{content}(C, G)$$

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

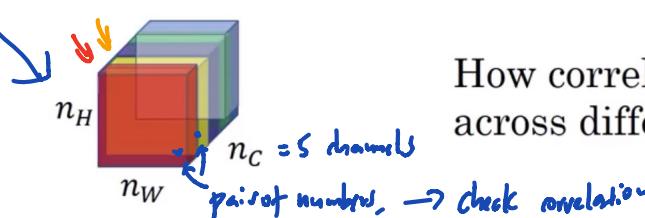
element-wise
assume unrolled into a vector

Style Cost Function

Meaning of the “style” of an image



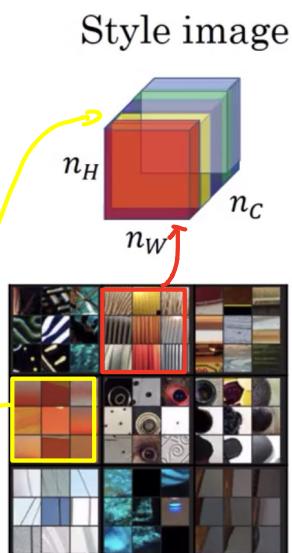
Say you are using layer l 's activation to measure “style.”
Define style as correlation between activations across channels.



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

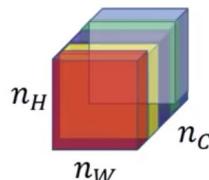
Andrew Ng

Intuition about style of an image



Correlation:
image wide
vertical
also how
orange art
to features
occurs
together

Generated Image



So correlation gives us info on how similar style between style image and generated image

Style matrix ↪ measures all the correlations

Let $a_{i,j,k}^{(l)}$ = activation at (i,j,k)

$\begin{matrix} u \\ v \\ w \end{matrix}$; c

$G^{(l)}$ is $n_c^{(l)} \times n_c^{(l)}$ ↪ style matrix

for 1 element in G

$$G_{Kk'}^{(l) (s)} = \sum_{i=1}^{n_u^{(l)}} \sum_{j=1}^{n_w^{(l)}} a_{ijk}^{(l) (s)} a_{ijk'}^{(l) (s)}$$

style

for general image

↳ $G_{kk'}^{(l)}$ measures the correlation between channel K and K' (where $k=1, \dots, n_c^{(l)}$)

"gram matrix" in linear algebra

$$G_{Kk'}^{(l) (s)} = \sum_{i=1}^{n_u^{(l)}} \sum_{j=1}^{n_w^{(l)}} a_{ijk}^{(l) (s)} a_{ijk'}^{(l) (s)}$$

$$J_{style}^{(l)}(S, G) = \|G_{Kk'}^{(l) (s)} - G_{Kk'}^{(l) (G)}\|_F^2$$

$$\text{normalize} = \frac{1}{(2n_u^{(l)} n_w^{(l)} n_c^{(l)})^2} \sum_{K} \sum_{K'} (G_{Kk'}^{(l) (s)} - G_{Kk'}^{(l) (G)})^2$$

overall style loss function

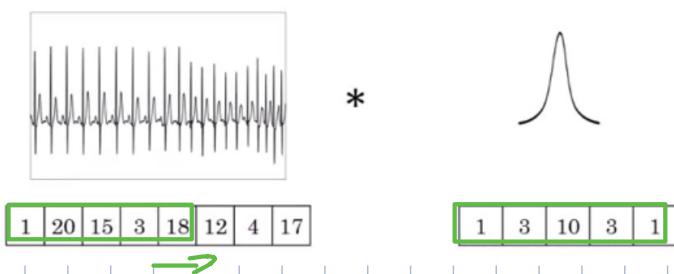
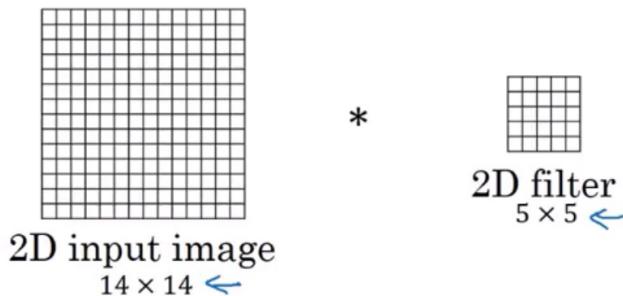
$$J_{style}(S, G) = \sum_l \sum_t \lambda^{(l)} J_{style}^{(l)}(S, G)$$

weights: hyperparameters

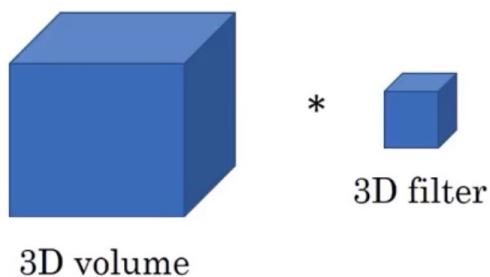
recap:

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

1D and 3D Generalisations



3D convolution



input $14 \times 14 \times 14 \times 1$ channel
 $\times 5 \times 5 \times 5 \times 1$ filter
 $\rightarrow 10 \times 10 \times 10 \times 1$ if there are
 1 filters