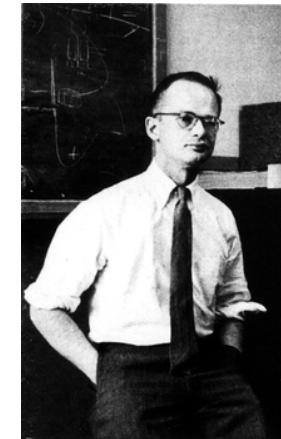
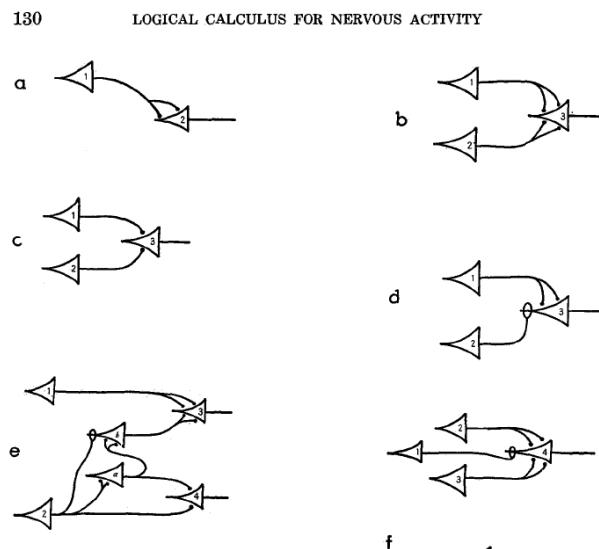


Neural Network



Warren
McCulloch



Walter Pitts

A LOGICAL CALCULUS OF THE IDEAS IMMANENT IN NERVOUS ACTIVITY,
Bulletin of Mathematical Biophysics, Vol. 5, pp. 115-133 (1943).

Neural Network

Here x_1 and x_2 are normalized attribute value of data.

y is the output of the neuron , i.e the class label.

x_1 and x_2 values multiplied by weight values w_1 and w_2 are input to the neuron x .

Value of x_1 is multiplied by a weight w_1 and values of x_2 is multiplied by a weight w_2 .

Given that

- $w_1 = 0.5$ and $w_2 = 0.5$
- Say value of x_1 is 0.3 and value of x_2 is 0.8,

- So, weighted sum is :
- $\text{sum} = w_1 \times x_1 + w_2 \times x_2 = 0.5 \times 0.3 + 0.5 \times 0.8 = 0.55$

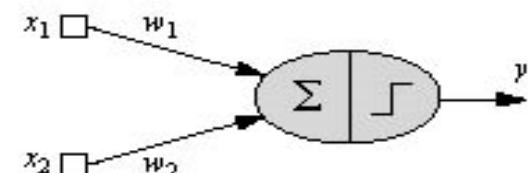
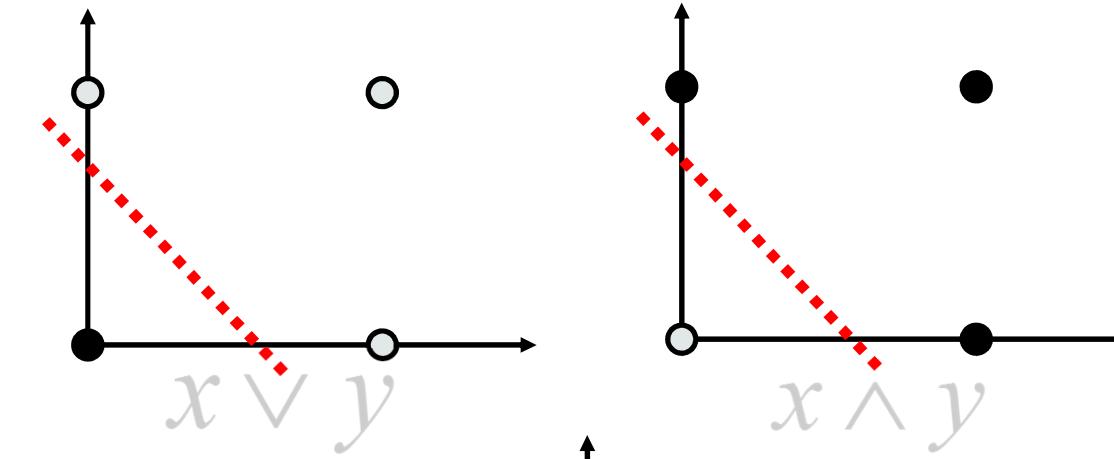


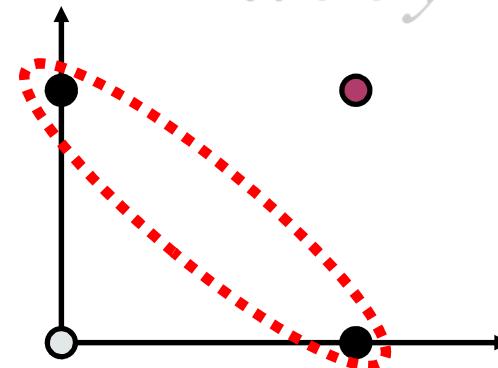
Fig1: an artificial neuron

Why We Need Multi Layer ?

Linear Separable:



Linear inseparable:



Edge Detection



Vertical edges

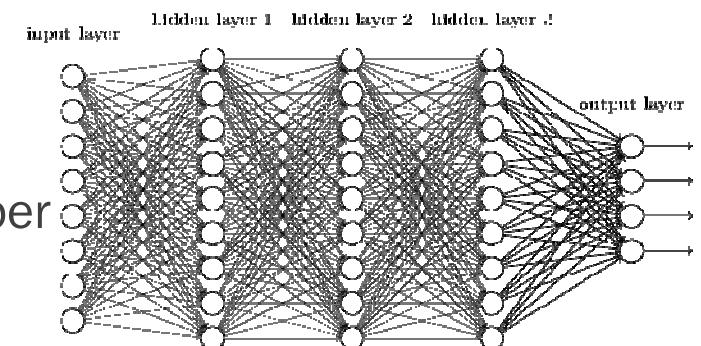


Horizontal edges

How do we detect
these edges

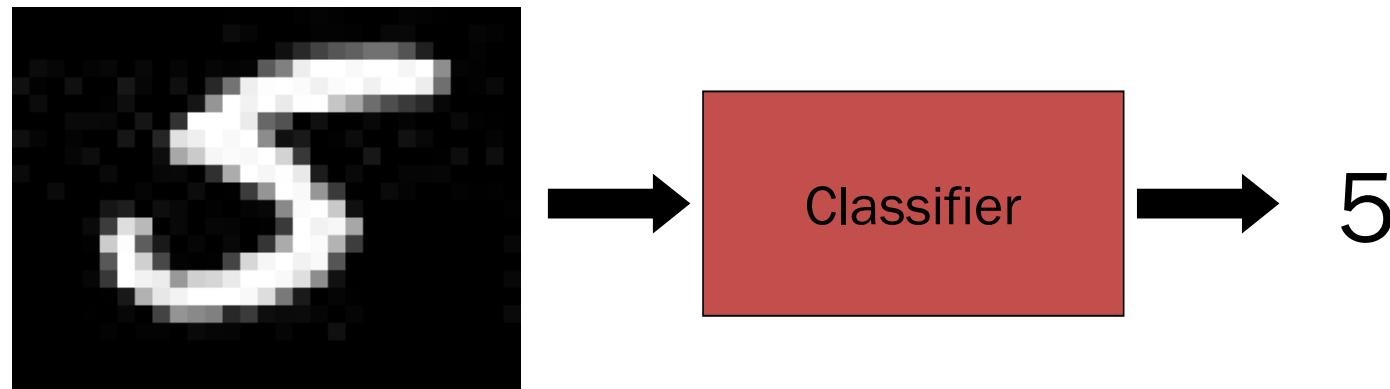
Neural Network?

- ❑ Suppose an image is of the size $68 \times 68 \times 3$
 - Input feature dimension then becomes 12,288
- ❑ If Image size is of $720 \times 720 \times 3$
 - Input feature dimension becomes 1,555,200
- ❑ Number of parameters will swell up to a HUGE number
- ❑ Result in more computational and memory requirements



Another Application

Digit Recognition



$X_1, \dots, X_n \in \{0,1\}$ (Black vs. White pixels)

$Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

Model Parameters

For the Bayes classifier, we need to “learn” two functions, the likelihood and the prior

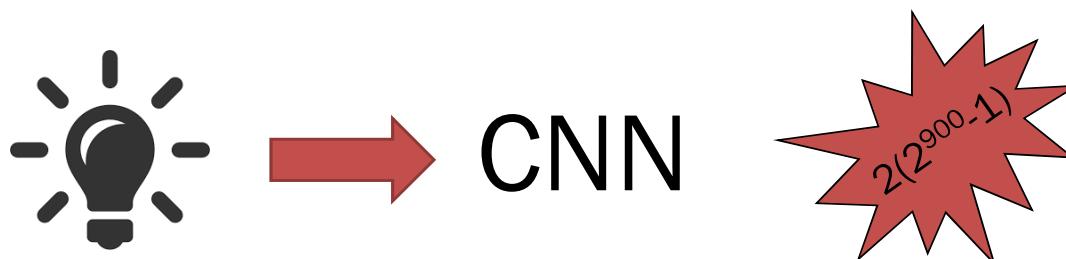
How many parameters are required to specify the prior for our digit recognition example?



Model Parameters

How many parameters are required to specify the likelihood?

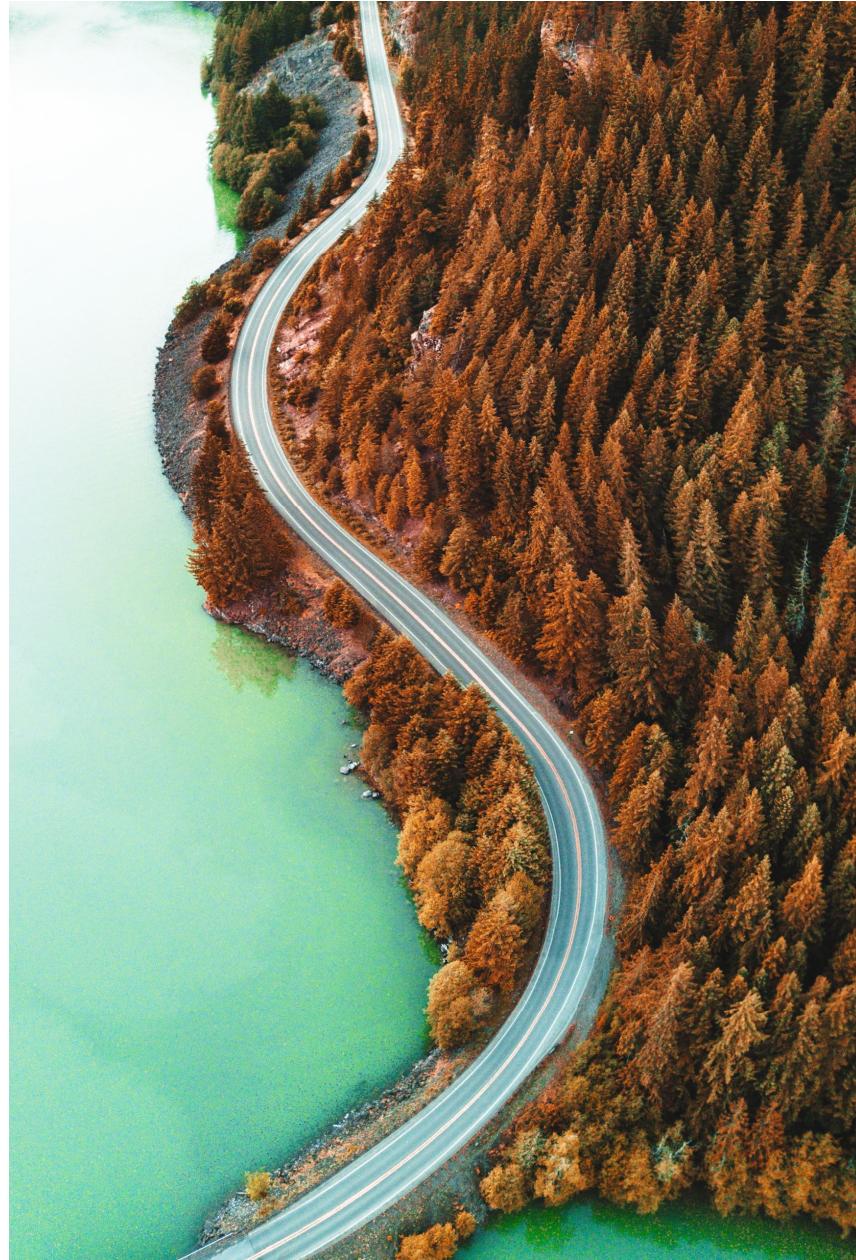
- (Supposing that each image is 30x30 pixels)



Drive into CNN

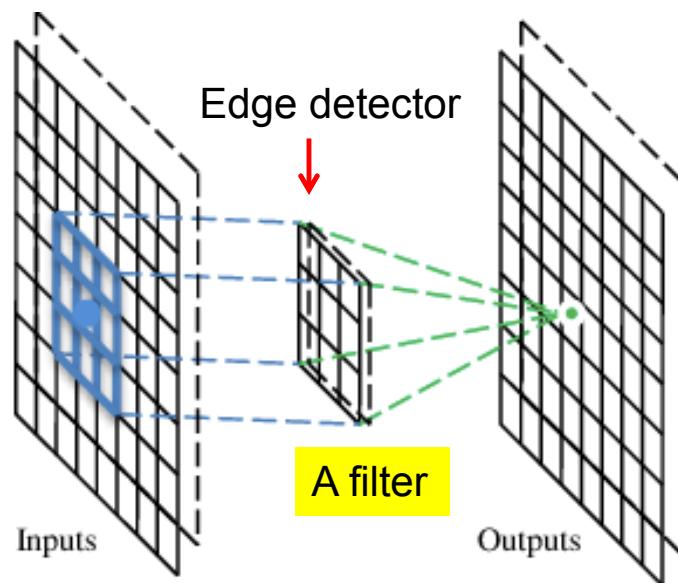
In a convolutional network (ConvNet), there are basically three types of layers:

1. Convolution layer
2. Pooling layer
3. Fully connected layer



A convolutional layer

A CNN is a neural network with some convolutional layers (and some other layers). A convolutional layer has a number of filters that does convolutional operation.



Convolution

These are the network parameters to be learned.

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

: :

Each filter detects a small pattern (3 x 3).

Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Dot
product



6 x 6 image

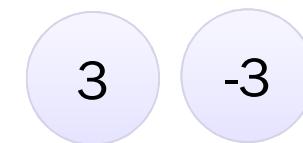
Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

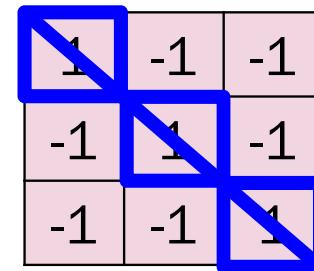
If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0



6 x 6 image

Convolution

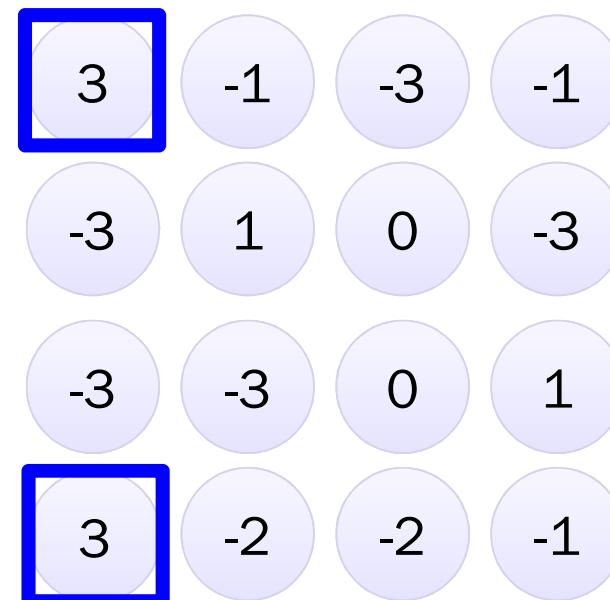


Filter 1

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image



Convolution

-1	1	-1
-1	1	-1
-1	1	-1

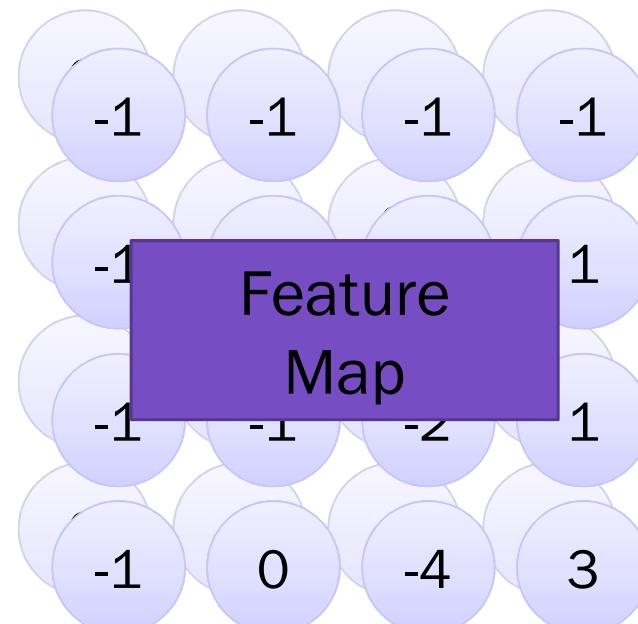
Filter 2

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

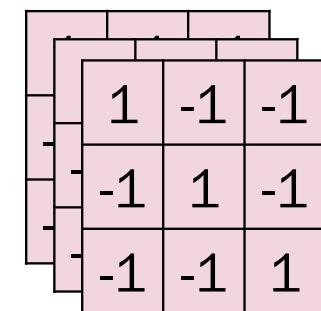
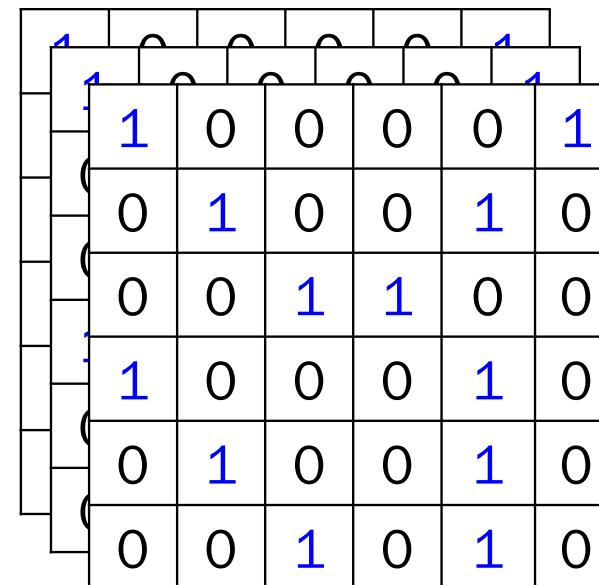
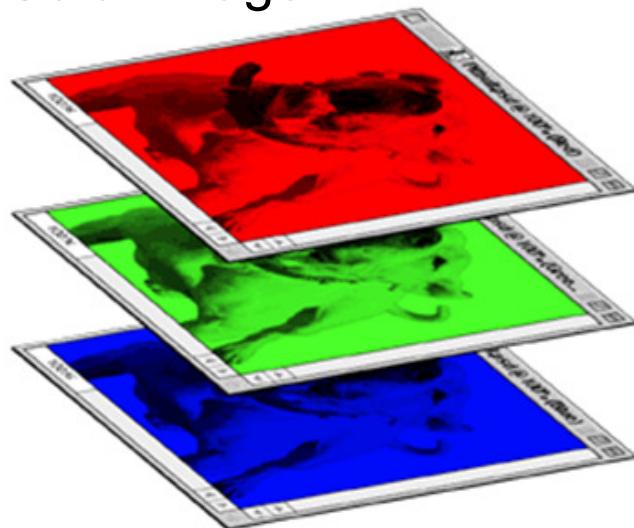
Repeat this for each filter



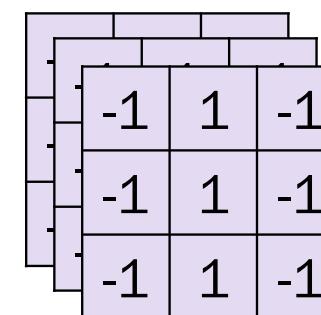
Two 4 x 4 images
Forming 4 x 4 x 2
matrix

Convolution over Volume

Color image

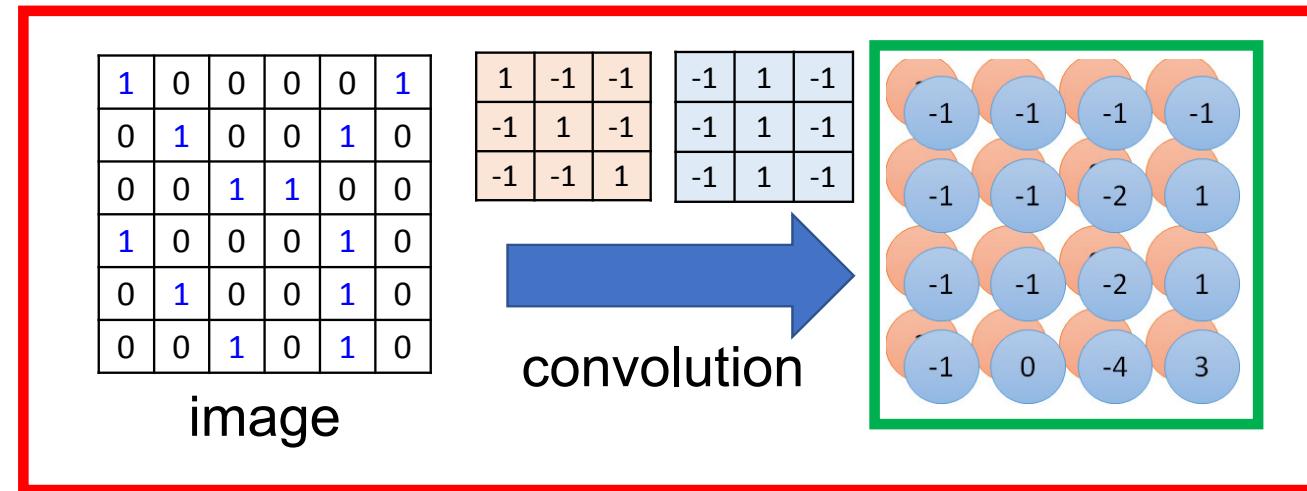


Filter 1



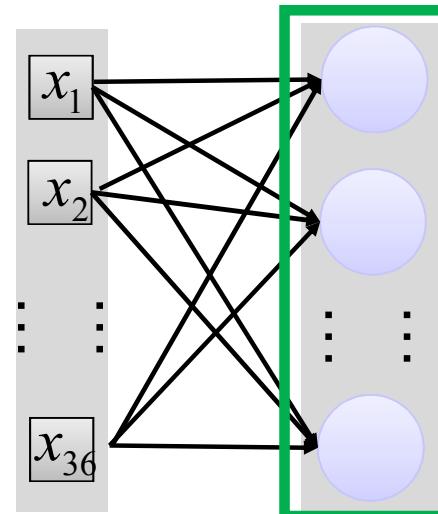
Filter 2

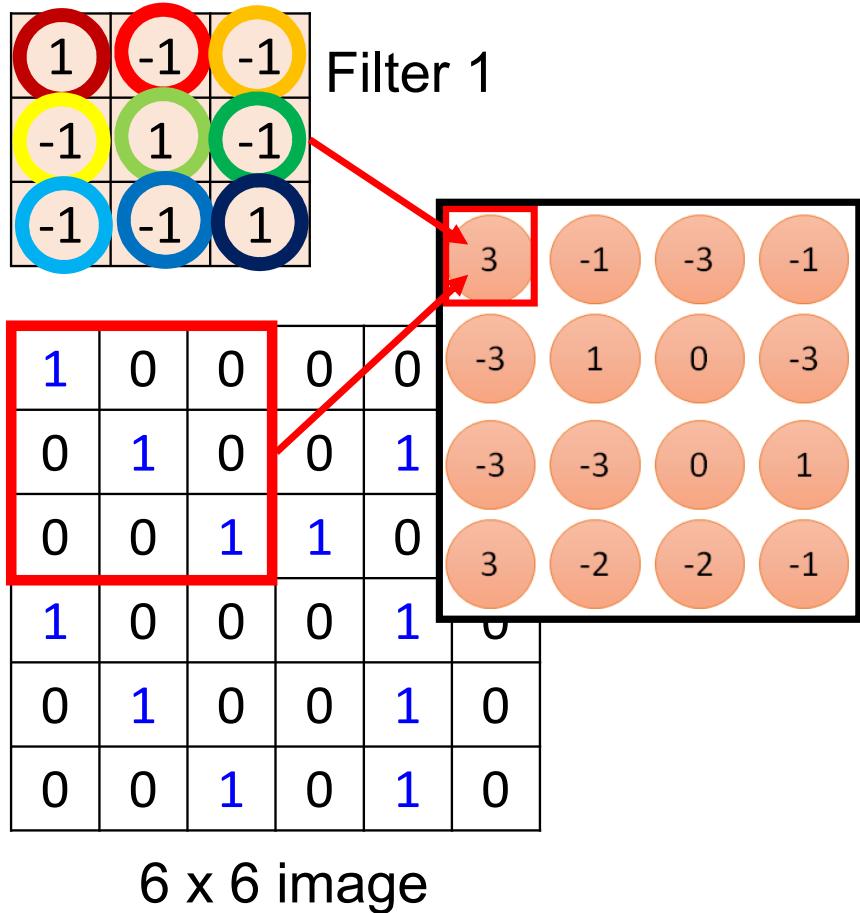
Convolution v.s. Fully Connected



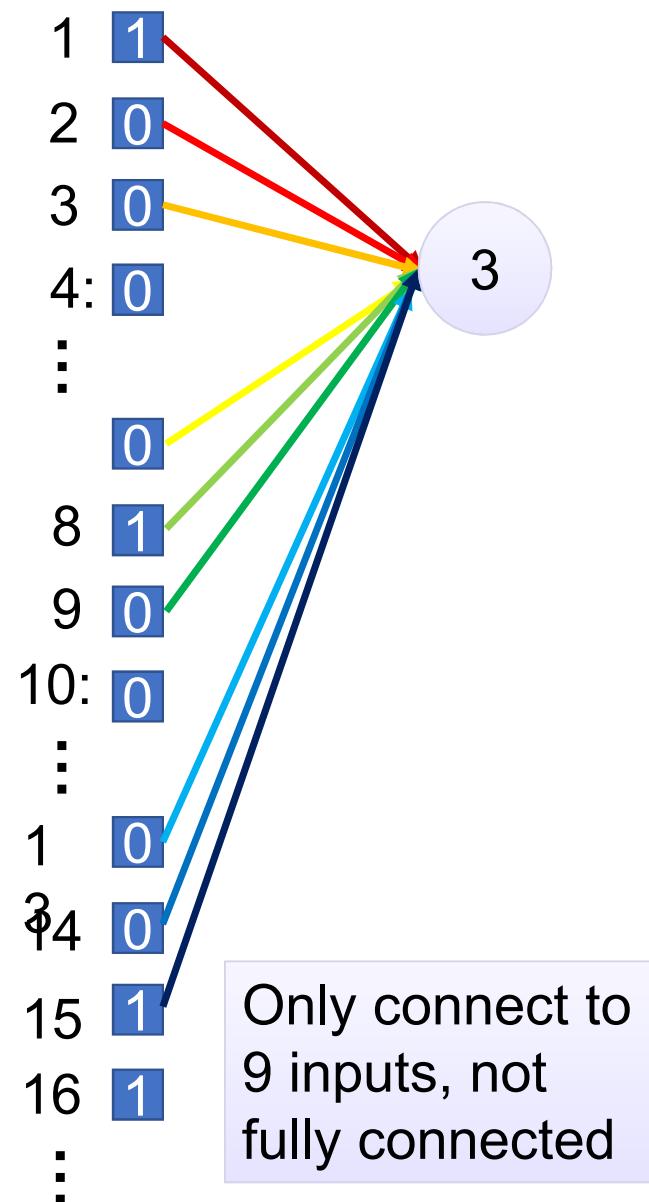
Fully-
connected

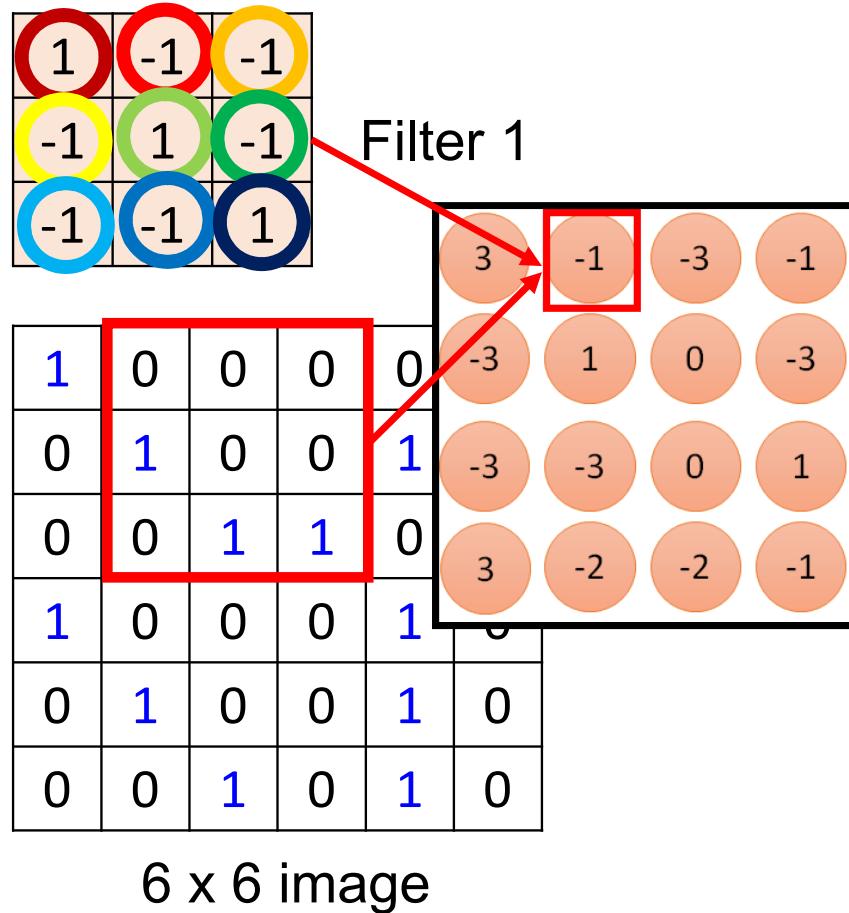
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0





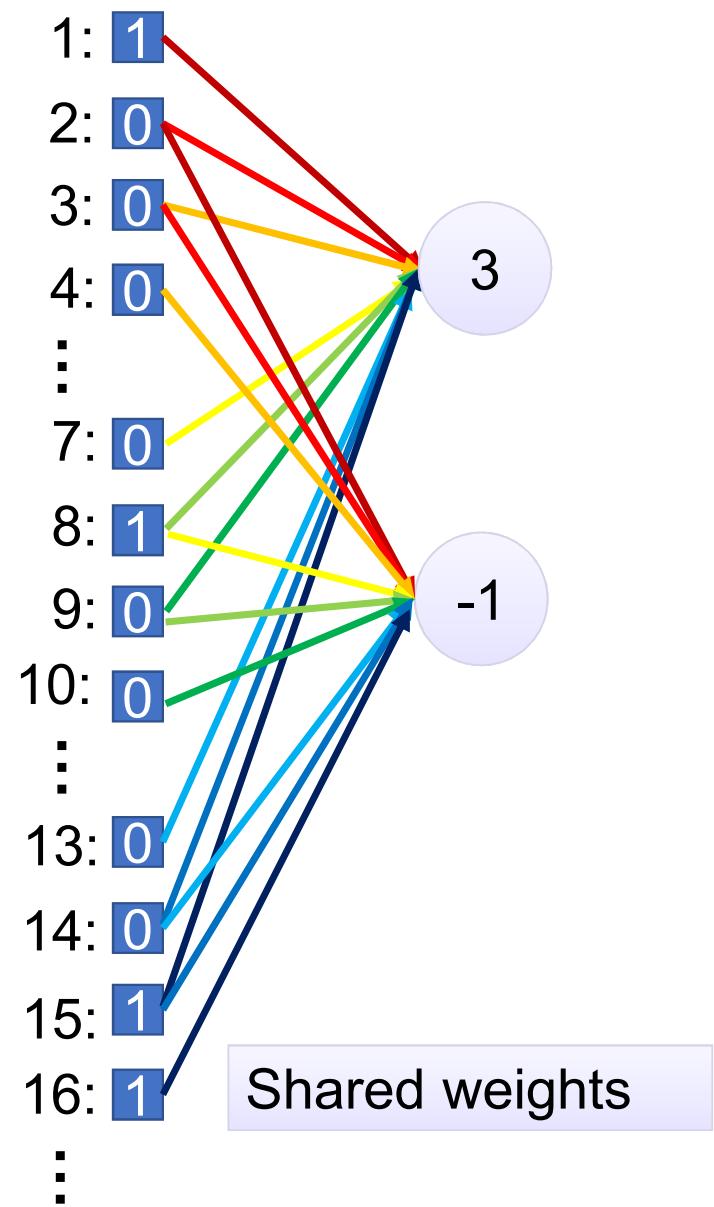
fewer parameters!



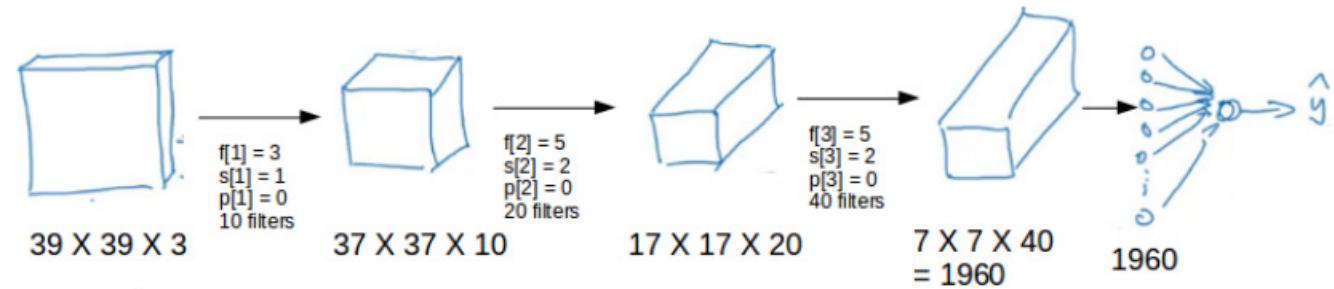


Fewer parameters

Even fewer parameters

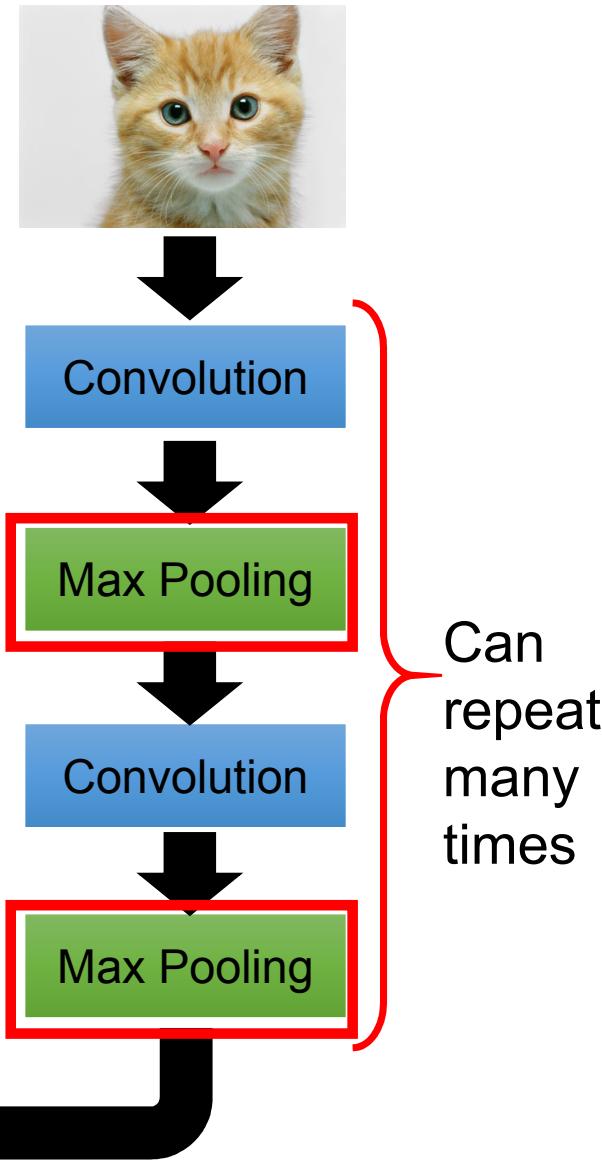
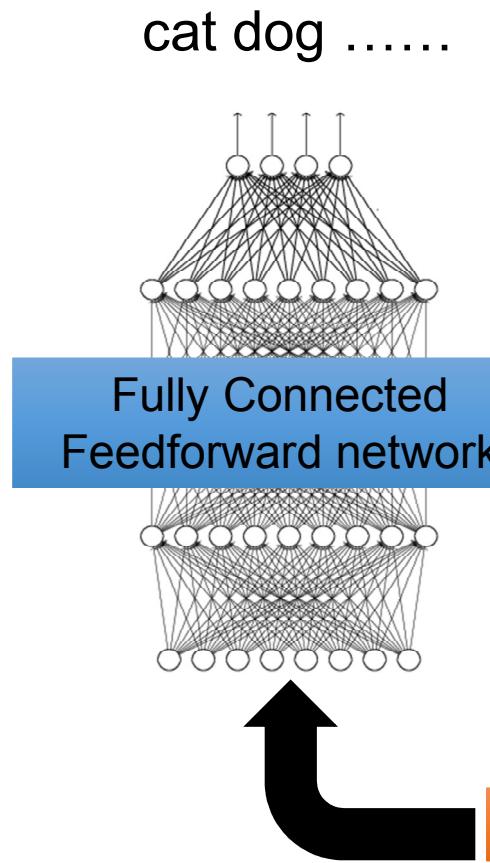


Simple Convolutional Neural Network



- **Size of feature vector :** $(n+2p-f)/s + 1$
- **n :** dimension of matrix
- **p :** size of padding
- **f :** size of filter
- **s :** size of stride

The whole CNN



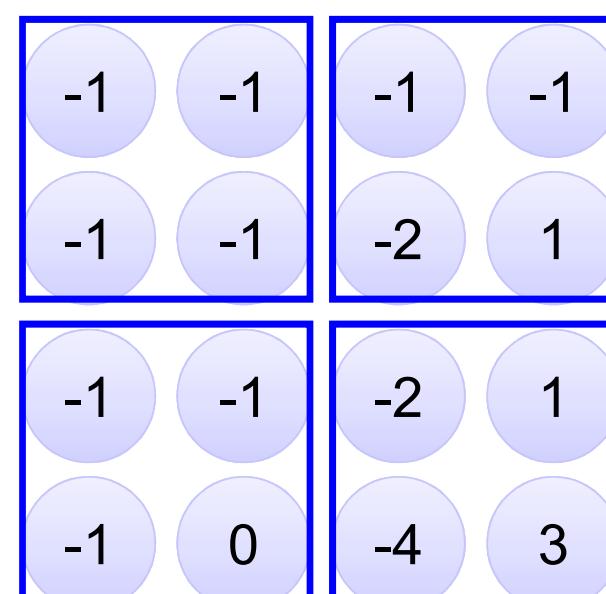
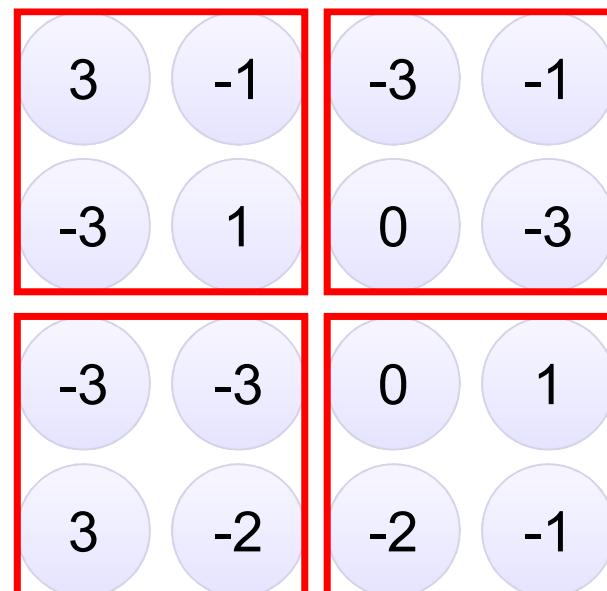
Max Pooling

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2



Why Pooling

- Subsampling pixels will not change the object

bird



Subsampling

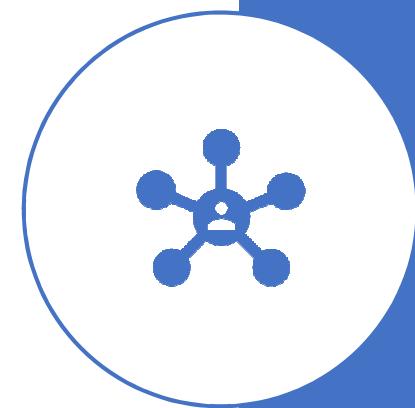
bird



We can subsample the pixels to make image
smaller
fewer parameters to characterize the image

A CNN compresses a fully connected network in two ways:

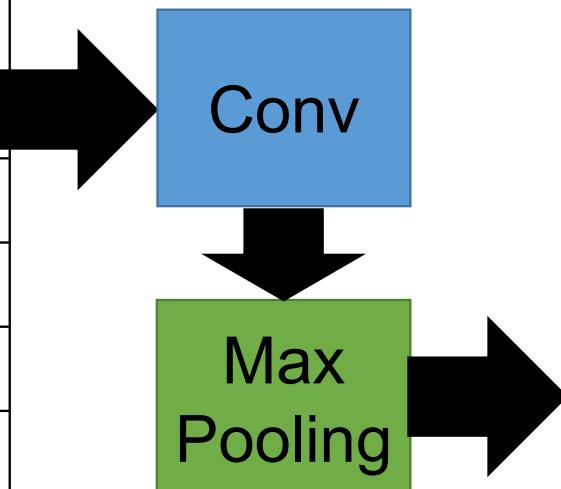
- Reducing number of connections
- Shared weights on the edges
- Max pooling further reduces the complexity



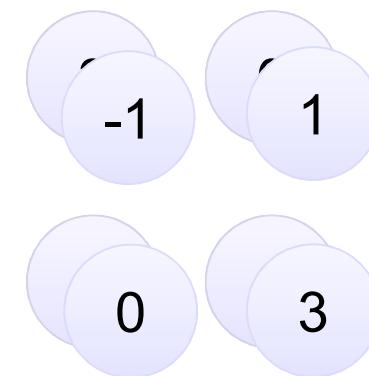
Max Pooling

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image



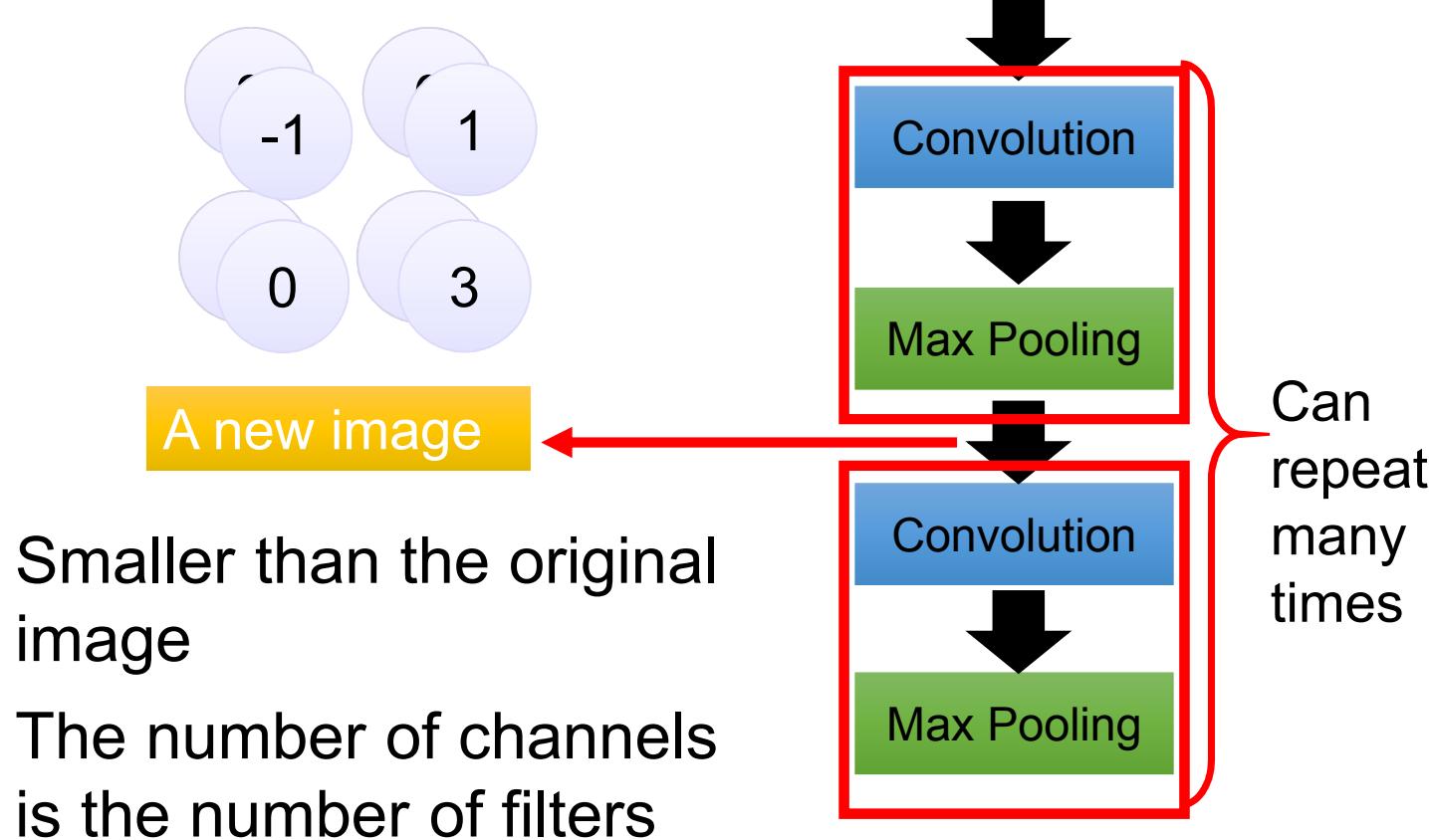
New image
but smaller



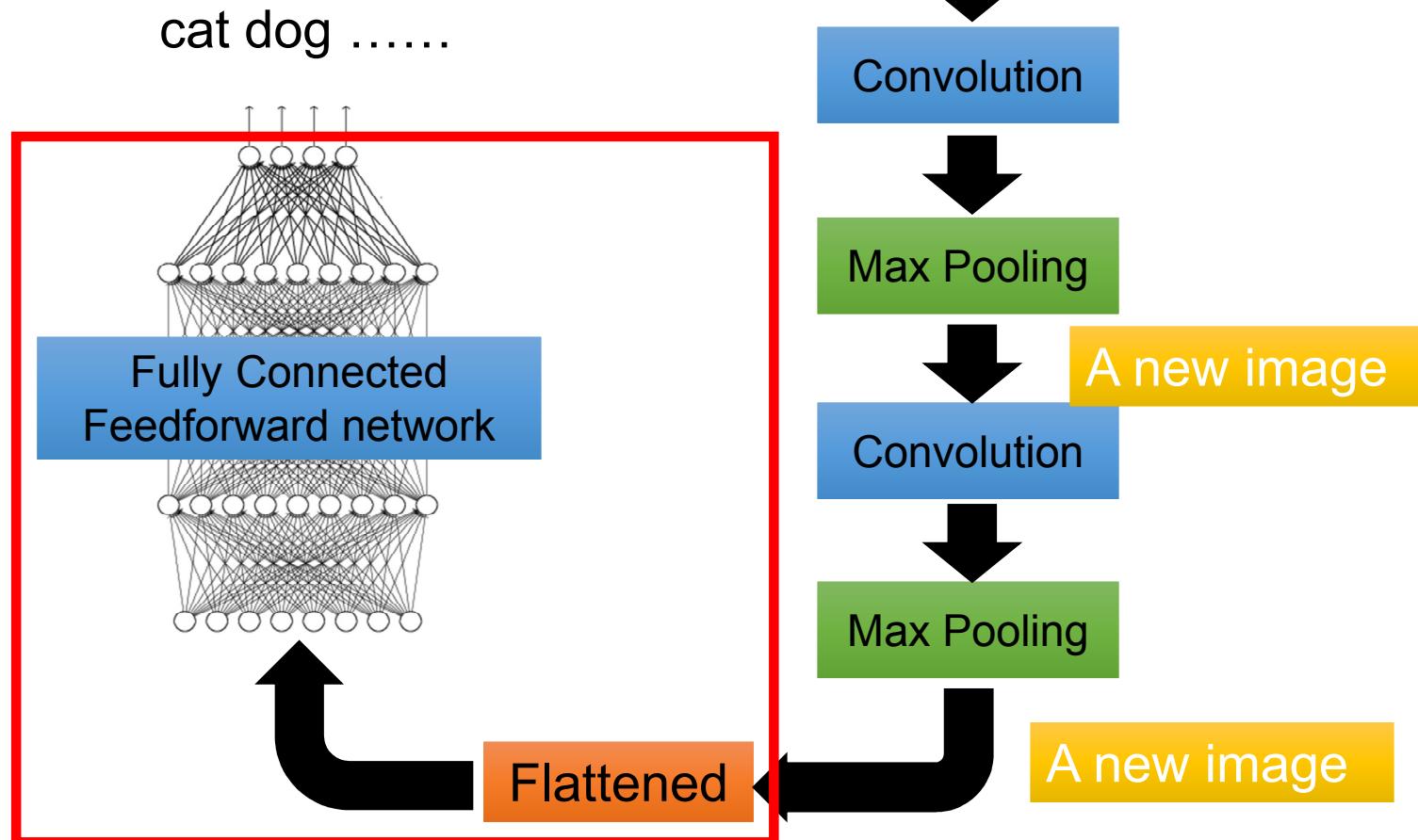
2 x 2 image

Each filter
is a channel

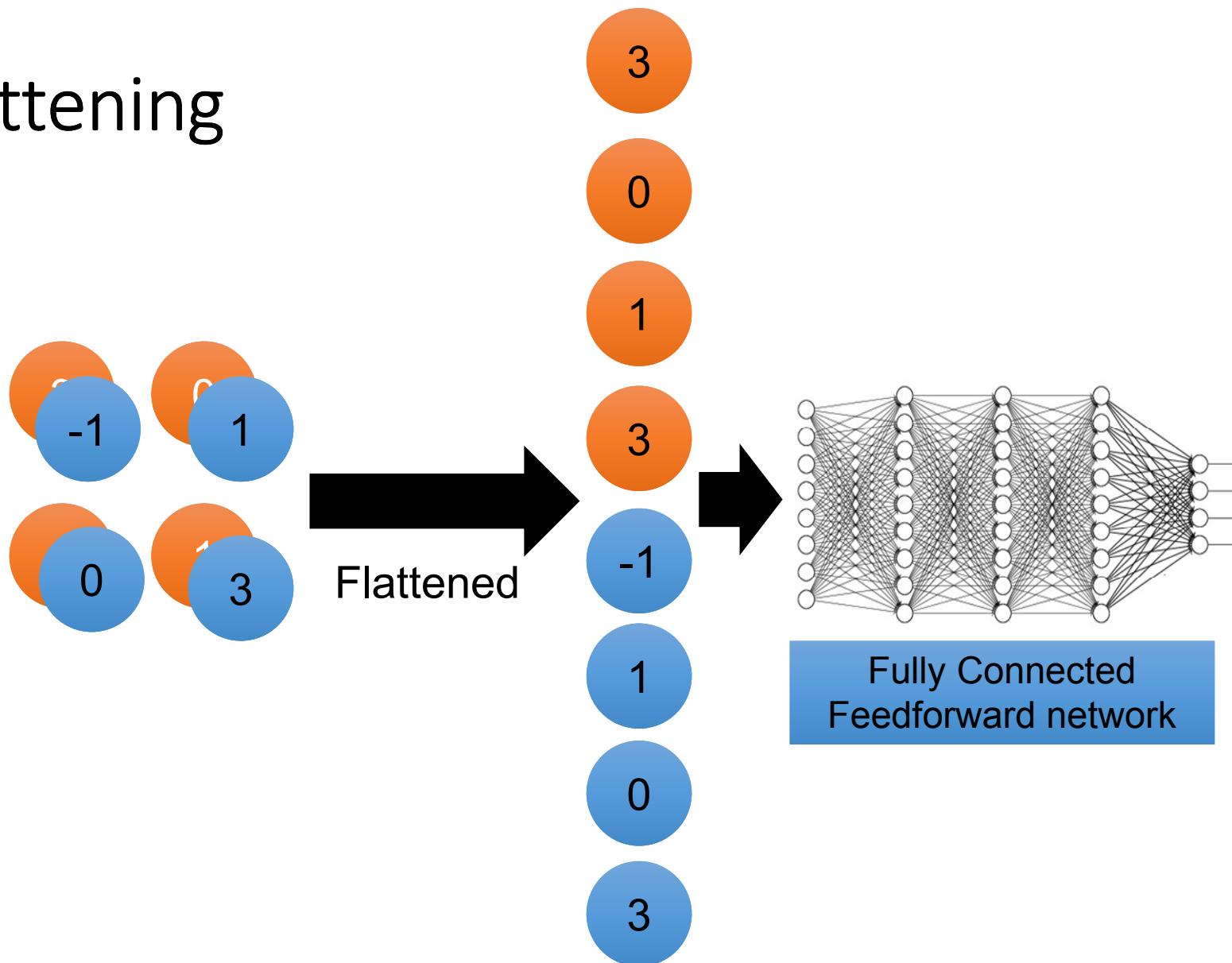
The whole CNN



The whole CNN



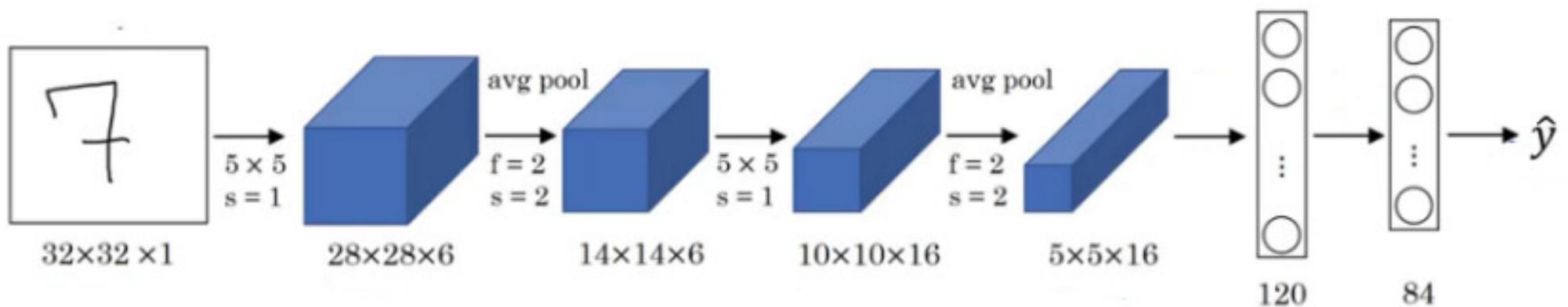
Flattening



Classic Networks

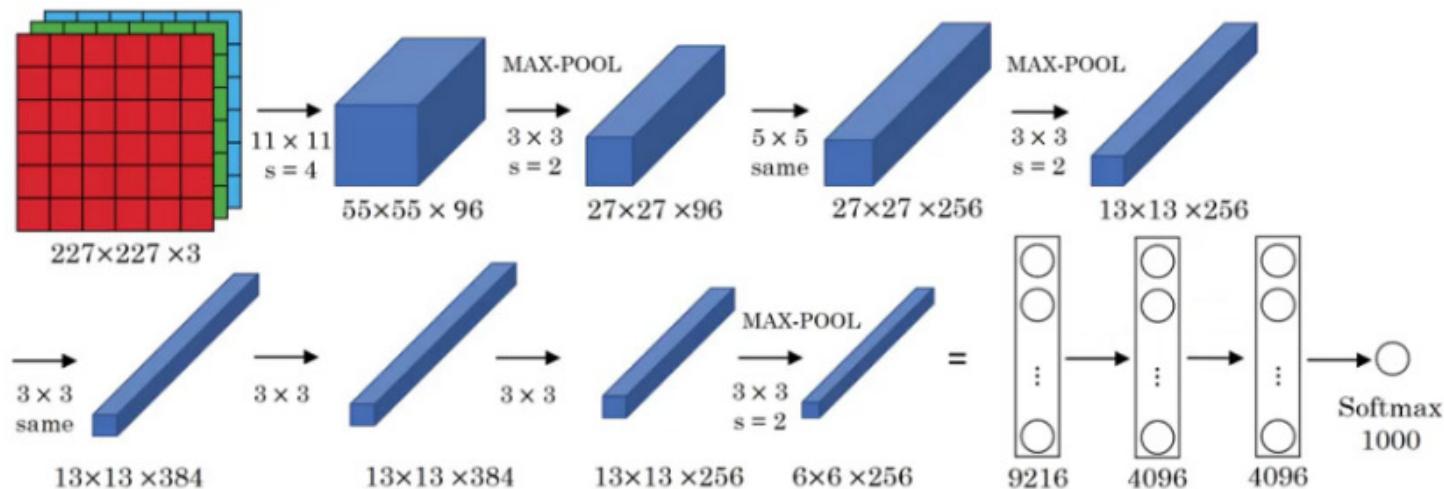
- 1.LeNet-5
- 2.AlexNet
- 3.VGG

LeNet-5



- **Parameters:** 60k
- **Layers flow:** Conv → Pool → Conv → Pool → FC → FC → Output
- **Activation functions:** Sigmoid/tanh and ReLU

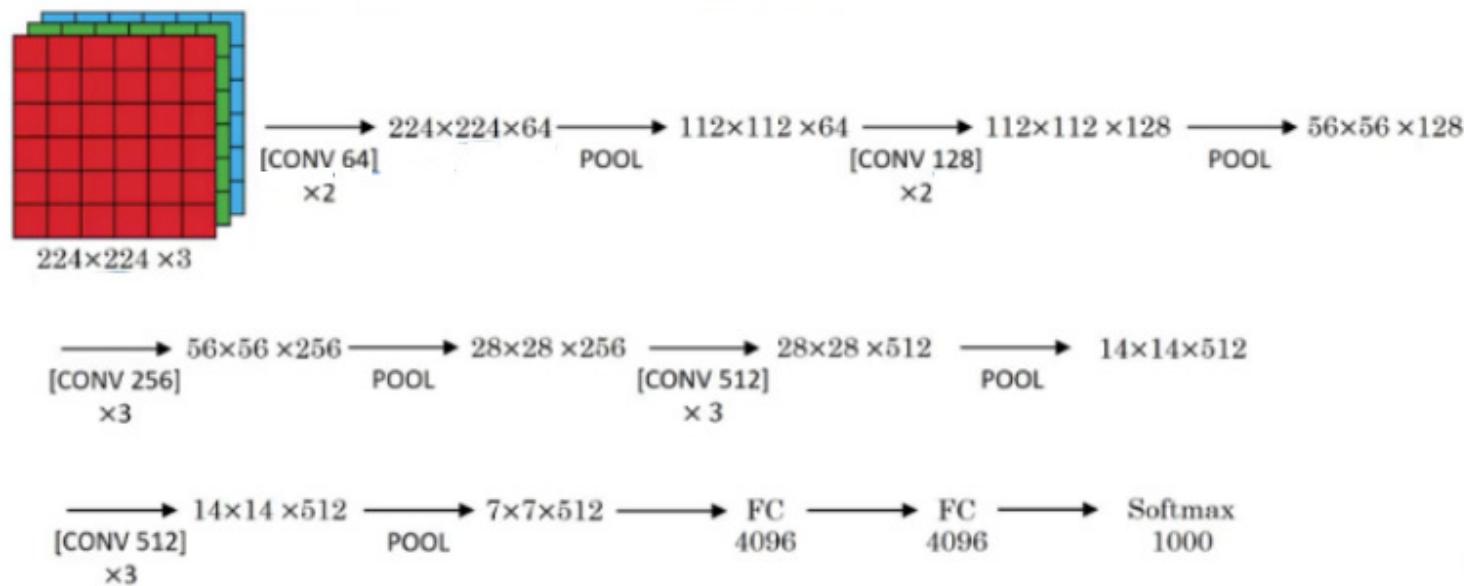
AlexNet



- **Parameters:** 60 million

- **Activation functions:** ReLU

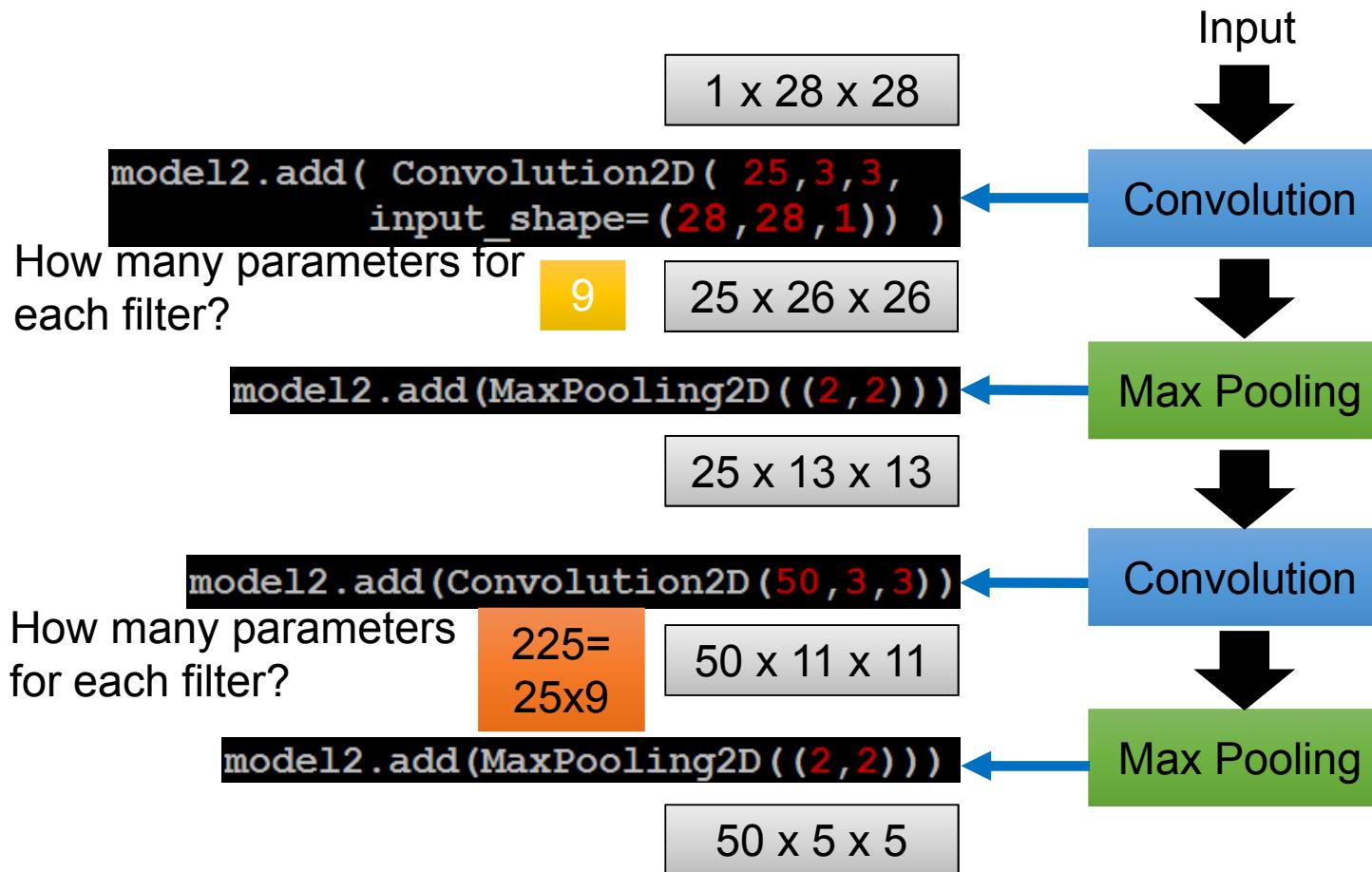
VGG-16



- **Parameters:** 138 million
- **Pool:** MAX with stride 2
- **CONV layer:** stride 1

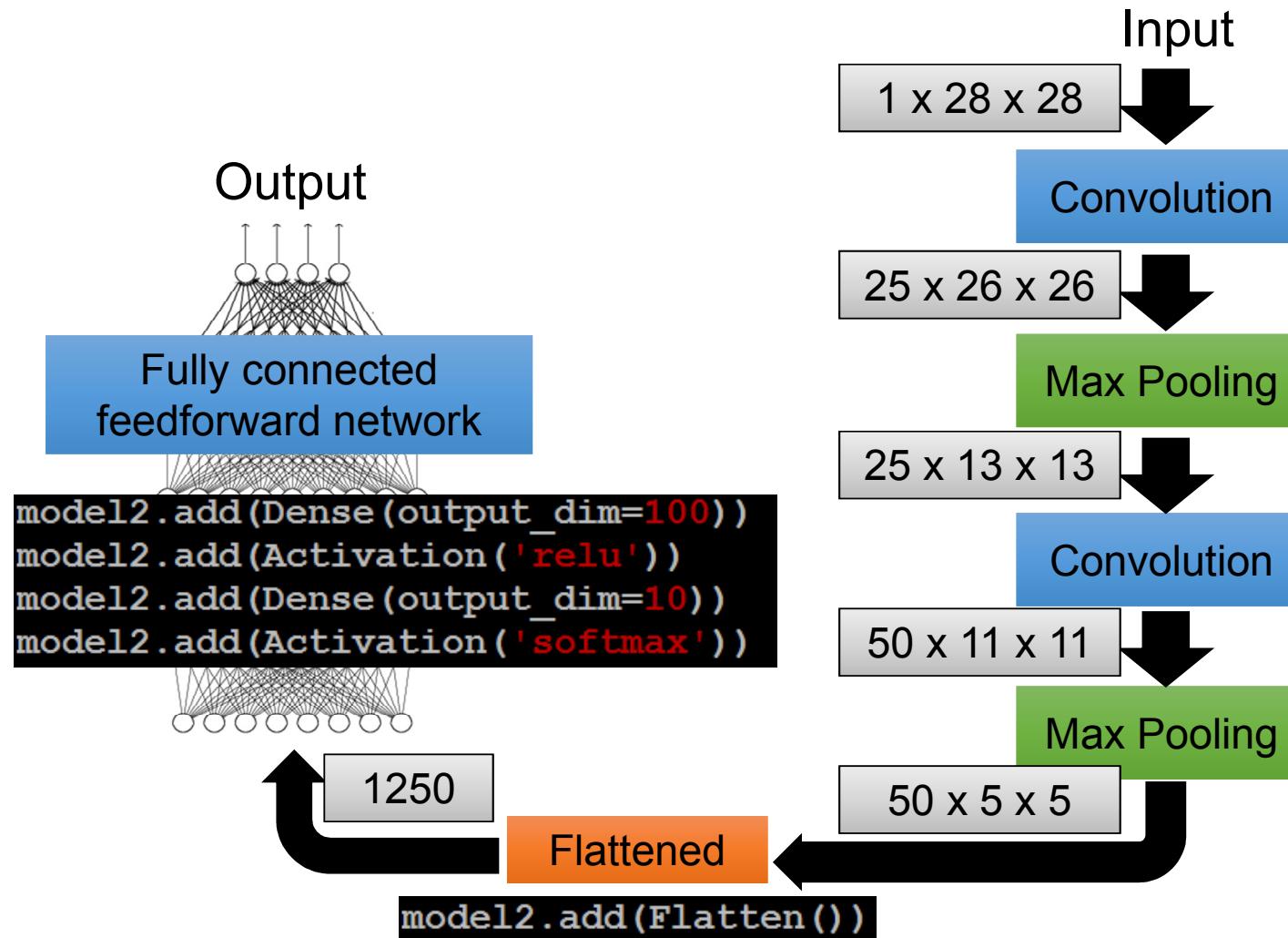
CNN in Keras

Only modified the ***network structure*** and
input format (vector -> 3-D array)



CNN in Keras

Only modified the ***network structure*** and
input format (vector -> 3-D array)



Object Detection using CNN

Classification



CAT

Classification + Localization = Detection



CAT

Object Detection is modeled as a classification problem

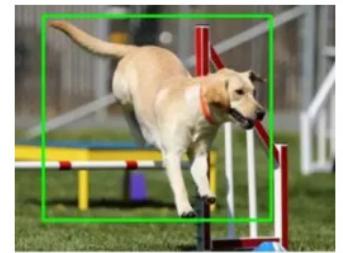
- We take windows of fixed sizes
- Run over input image at all the possible locations
- Feed these patches to an image classifier.
- It predicts the class of the object in the window(or background if none is present)

Problem → Solution

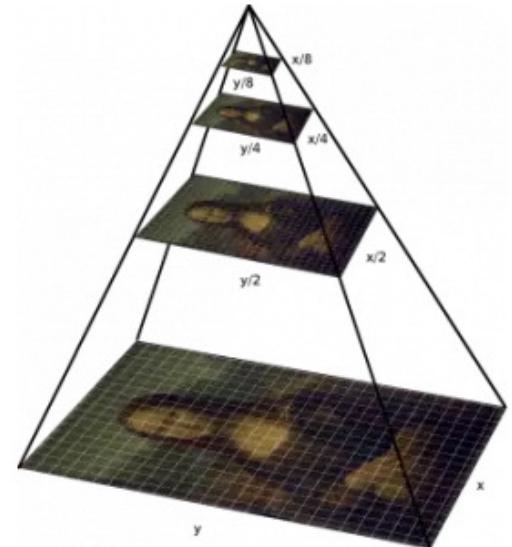
- Resize the image at multiple scales
- Most commonly, the image is downsampled(size is reduced)
- On each of these images, a fixed size window detector is run.
- Now, all these windows are fed to a classifier to detect the object of interest

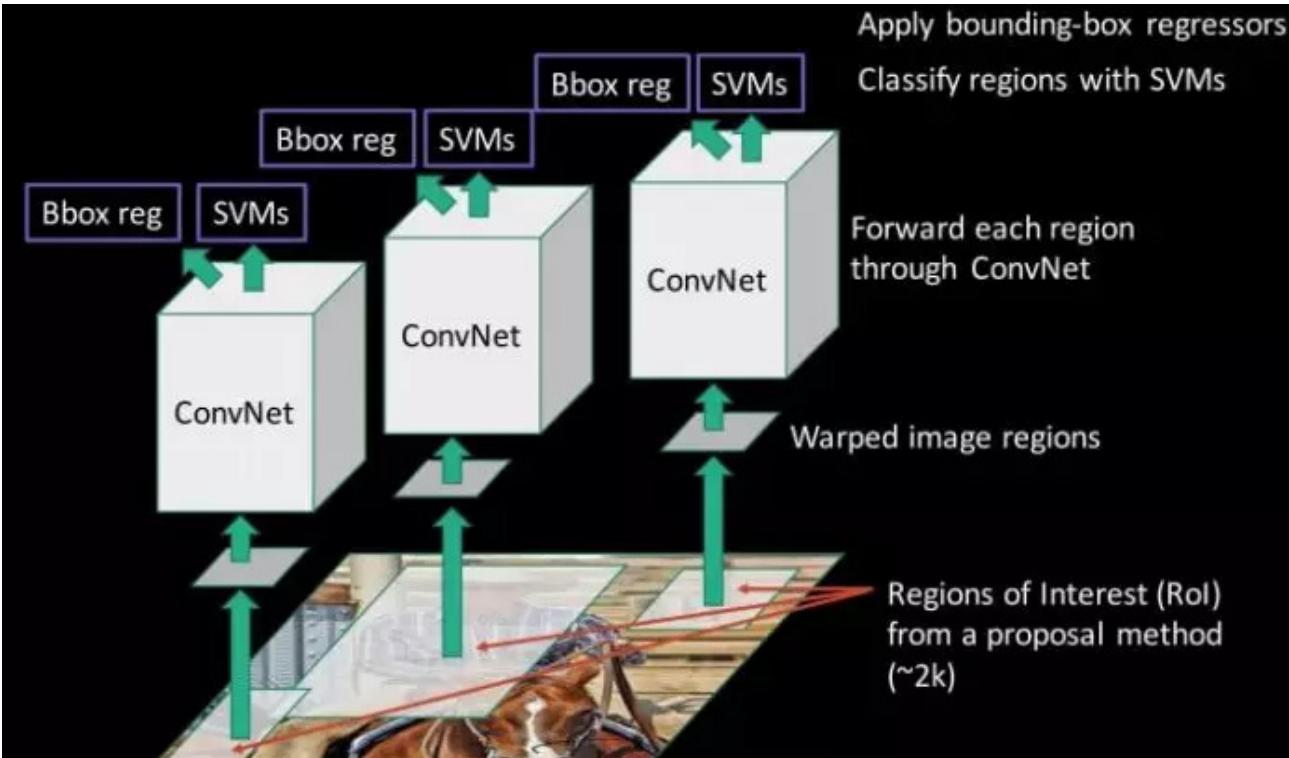


Small sized object



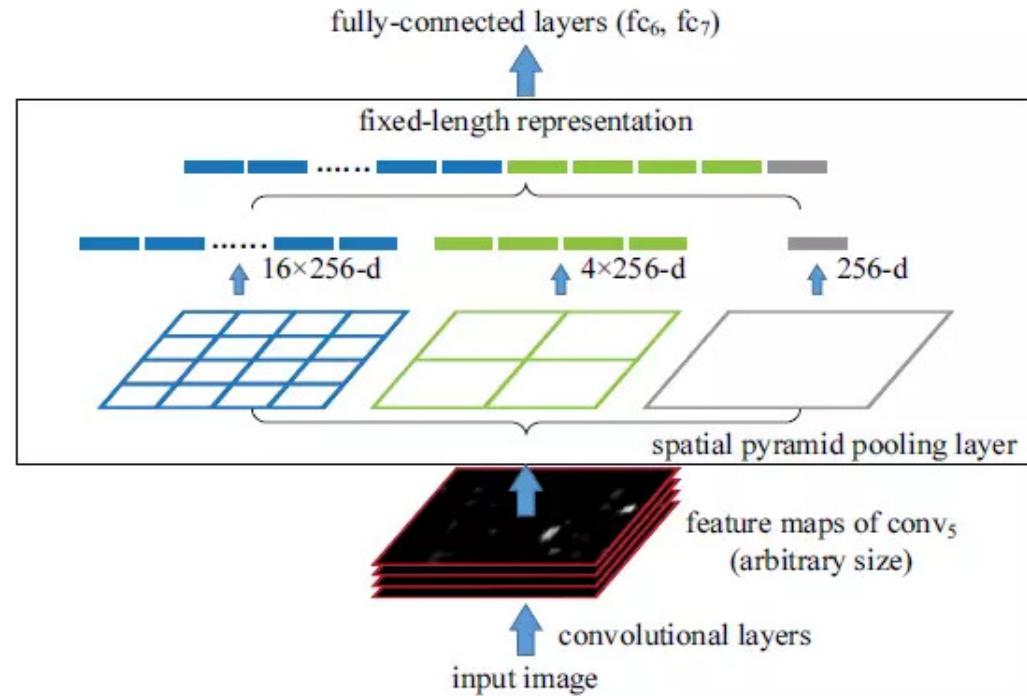
Big sized object. What size do you choose for your sliding window detector?





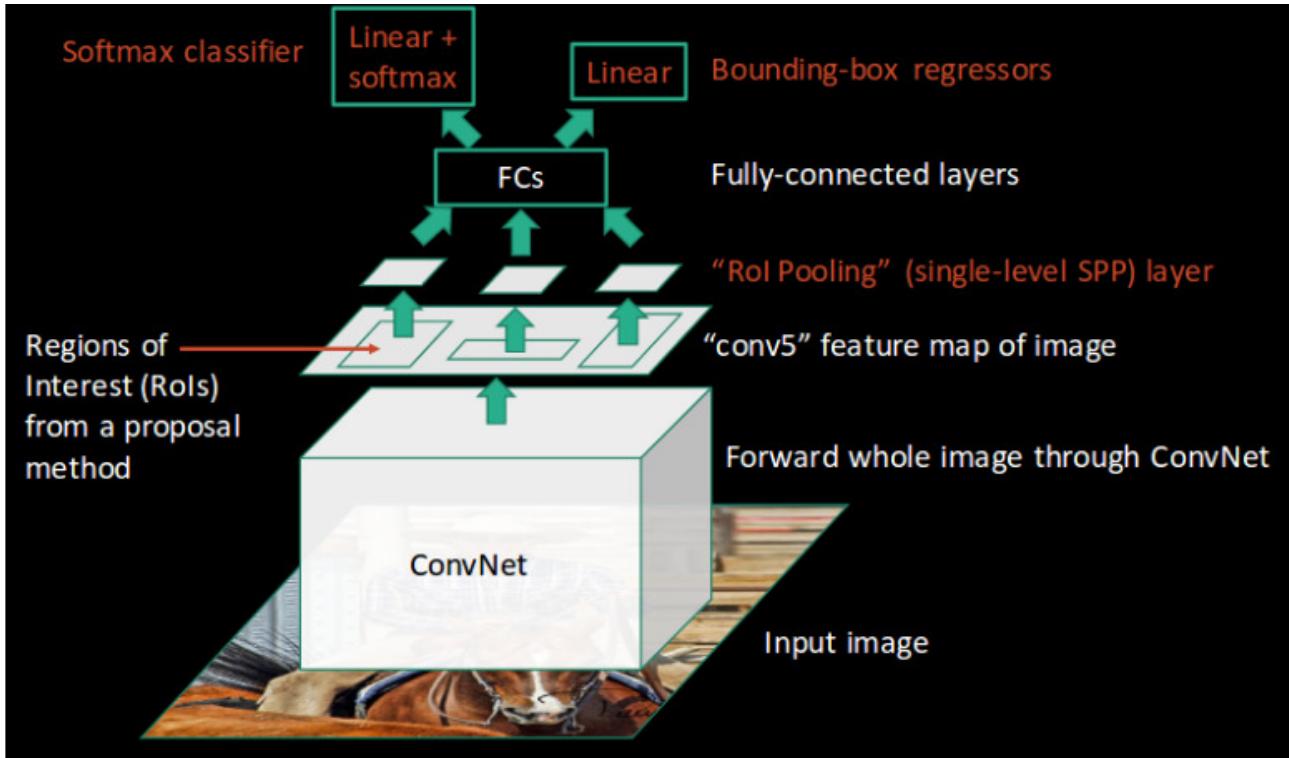
Region-based Convolutional Neural Networks(R-CNN)

- Run Selective Search to generate probable objects (~2k regions)
- Feed these patches to CNN, followed by SVM to predict the class of each patch.
- Optimize patches by training bounding box regression separately.



- Calculate the CNN representation for entire image only once
- **It uses spatial pooling after the last convolutional layer**
- SPP layer divides a region of any arbitrary size into a constant number of bins and max pool is performed on each of the bins
- Since the number of bins remains the same, a constant size vector is produced

Spatial Pyramid Pooling(SPP-net)

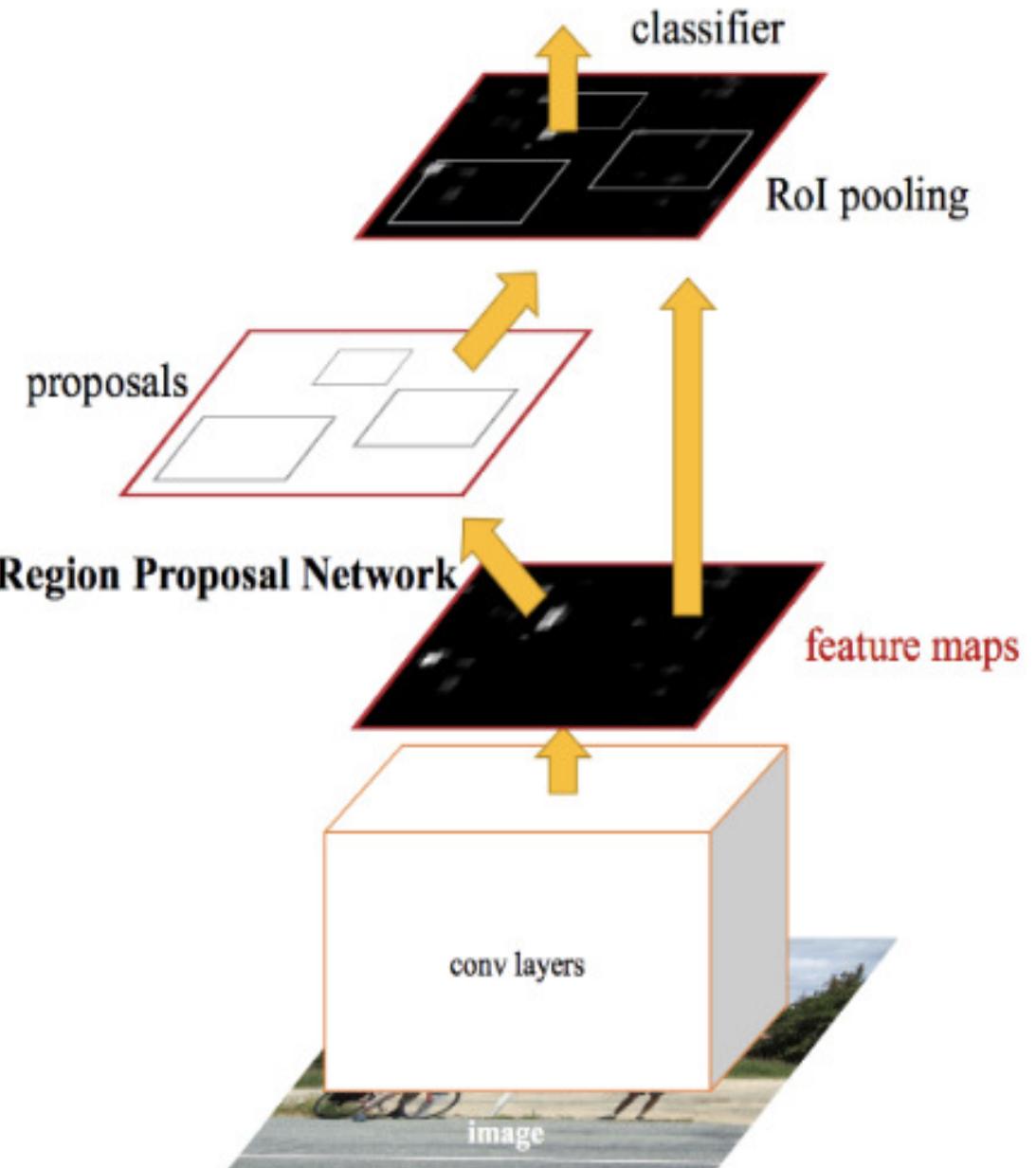


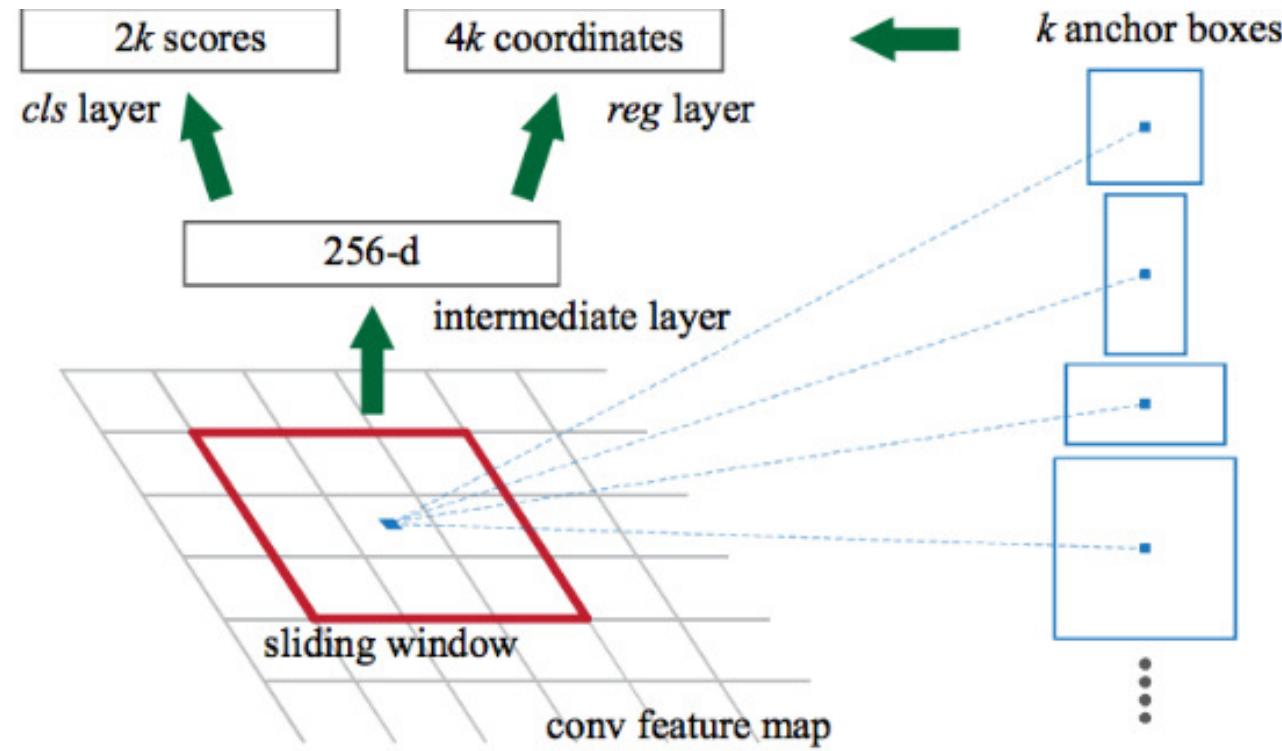
Fast R-CNN

- Fast RCNN uses the ideas from SPP-net and RCNN
- Apply the RoI pooling layer on the extracted regions of interest to make sure all the regions are of the same size
- These regions are passed on to a fully connected network which classifies them, as well as returns the bounding boxes using softmax and linear regression layers simultaneously

Faster R-CNN

- We take an image as input and pass it to the ConvNet which returns the feature map for that image
- **Region Proposal Network (lightweight CNN)** is applied on these feature maps. This returns the object proposals along with their objectness score
- A RoI pooling layer is applied on these proposals to bring down all the proposals to the same size
- Finally, the proposals are passed to a fully connected layer which has a softmax layer and a linear regression layer at its top, to classify and output the bounding boxes for objects.

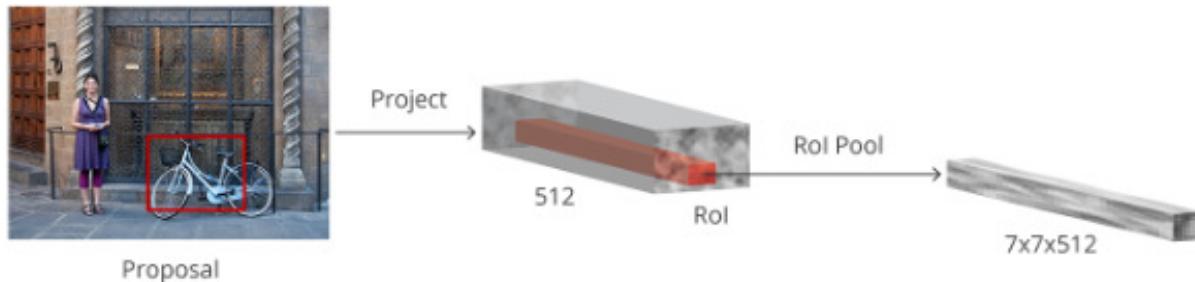




Region Proposal Network (RPN)

- RPN uses a sliding window over the feature maps
- At each window, it generates k Anchor boxes of different shapes and sizes
- For each anchor, RPN predicts two things:
 - first is the probability that an anchor is an object
 - Second is the bounding box regressor for adjusting the anchors to better fit the object

Region Proposal Network (RPN)



- We now have bounding boxes of different shapes and sizes which are passed on to the Roi pooling layer
- It extracts fixed sized feature maps for each anchor
- These feature maps are passed to a fully connected layer
- It has a softmax and a linear regression layer
 - Classifies the object
 - predicts the bounding boxes for the identified objects

Summary of the object detection models

Algorithm	Features	Prediction time / image	Limitations
CNN	Divides the image into multiple regions and then classify each region into various classes.	–	Needs a lot of regions to predict accurately and hence high computation time.
RCNN	Uses selective search to generate regions. Extracts around 2000 regions from each image.	40-50 seconds	High computation time as each region is passed to the CNN separately also it uses three different model for making predictions.
Fast RCNN	Each image is passed only once to the CNN and feature maps are extracted. Selective search is used on these maps to generate predictions. Combines all the three models used in RCNN together.	2 seconds	Selective search is slow and hence computation time is still high.
Faster RCNN	Replaces the selective search method with region proposal network which made the algorithm much faster.	0.2 seconds	Object proposal takes time and as there are different systems working one after the other, the performance of systems depends on how the previous system has performed.

Two stages and Single stage Object Detectors

Two stage Detectors

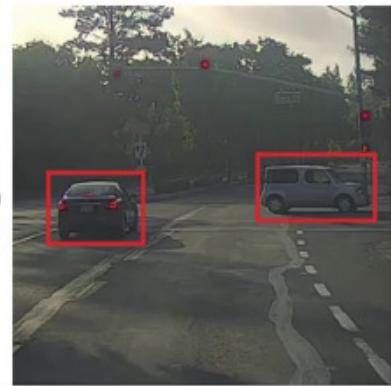
- first generates so-called region proposals — areas of the image that potentially contain an object
- Then it makes a separate prediction for each of these regions
- Examples : R-CNN, Fast R-CNN, Faster R-CNN, Mask R-CNN

One stage Detectors

- These models skip the explicit region proposal stage but apply the detection directly on dense sampled areas
- Examples: Single Shot Detector (SSD), YOLO family

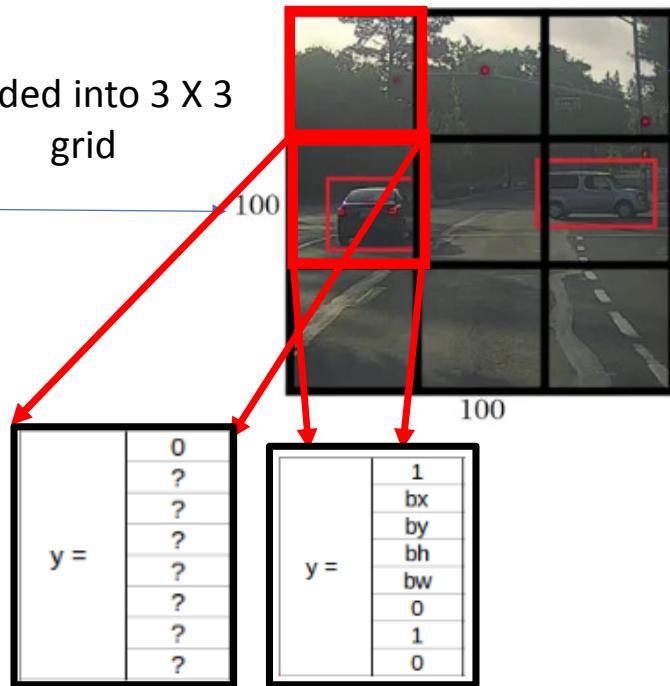
Yogi

How does YOLO Framework Function



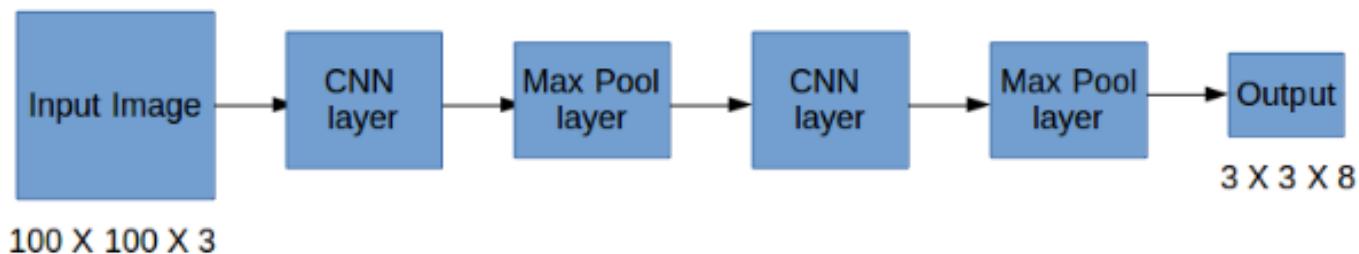
Input Image

Divided into 3 X 3 grid



- **Image classification and localization** are applied on each grid
- Suppose We have 3 classes. Let's say the classes are Pedestrian, Car, and Motorcycle, respectively. So, for each grid cell, the label y will be an eight-dimensional vector

$y =$	pc
	bx
	by
	bh
	bw
	0
	1
	0



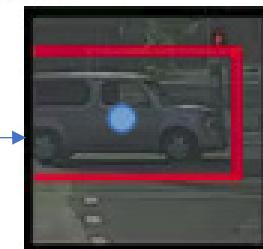
Bounding box in details



Grid contain
bounding box

YOLO assign coordinates to all the grids

(0,0)



(1,1)

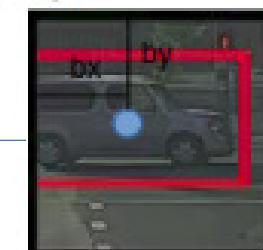
b_x , b_y are the x
and y coordinates
of the midpoint
of the object with
respect to this
grid

$y =$	1
	0.4
	0.3
	0.9
	0.5
	0
	1
	0

b_h : height of the bounding box / height of the grid

b_w : width of the bounding box / width of the grid

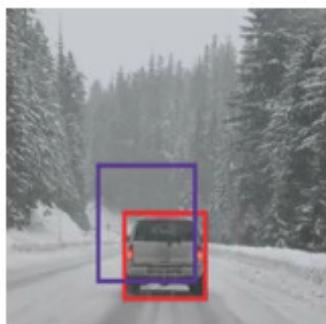
(0,0)



(1,1)

Intersection over Union and Non-Max Suppression

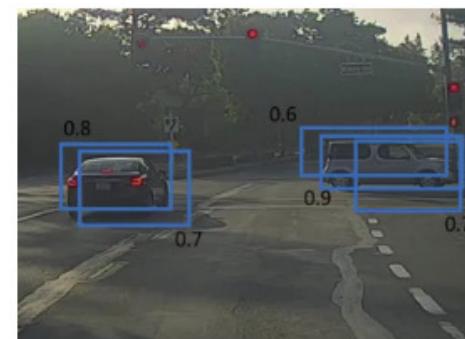
How can we decide whether the predicted bounding box is giving us a good outcome ?



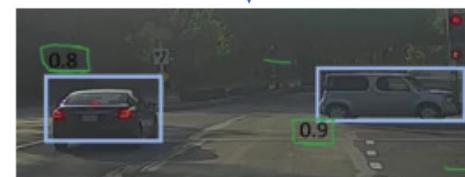
IoU = Area of the intersection / Area of the union

If $\text{IoU} > 0.5$, we accept predicted bounding box

Rather than detecting an object just once, they might detect it multiple times

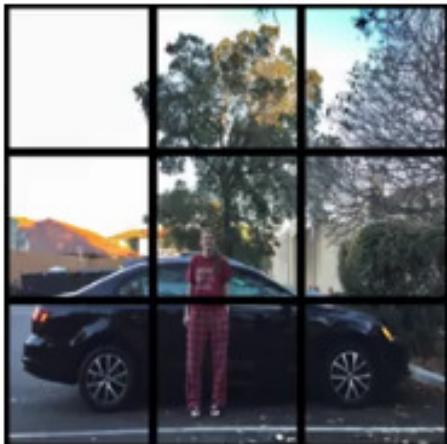


Non-Max Suppression

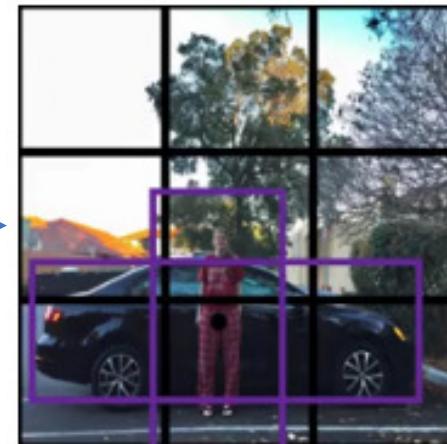


Anchor Box

what if there are multiple objects in a single grid?



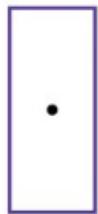
midpoint of both the objects
lies in the same grid



Anchor Box

what if there are multiple objects in a single grid?

Anchor box 1:



Anchor box 2:



- Since the shape of anchor box 1 is similar to the bounding box for the person, the latter will be assigned to anchor box 1 and the car will be assigned to anchor box 2
- The output in this case, instead of $3 \times 3 \times 8$ (using a 3×3 grid and 3 classes), will be $3 \times 3 \times 16$ (since we are using 2 anchors)

y =	pc bx by bh bw c1 c2 c3 pc bx by bh bw c1 c2 c3
-----	----------------------------------------------------------------------------------------------

Anchor box 1

Anchor box 2

You Only Look Once

- Training
 - 3×3 grid with two anchors per grid
 - 3 different object classes
 - **y labels will have a shape of $3 \times 3 \times 16$**
 - suppose if we use 5 anchor boxes per grid
 - number of classes has been increased to 5
 - target will be $3 \times 3 \times 10 \times 5 = 3 \times 3 \times 50$