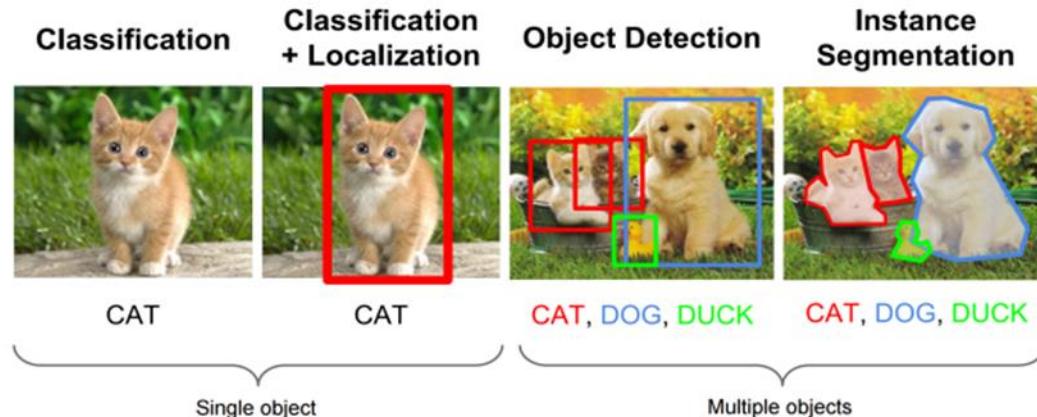
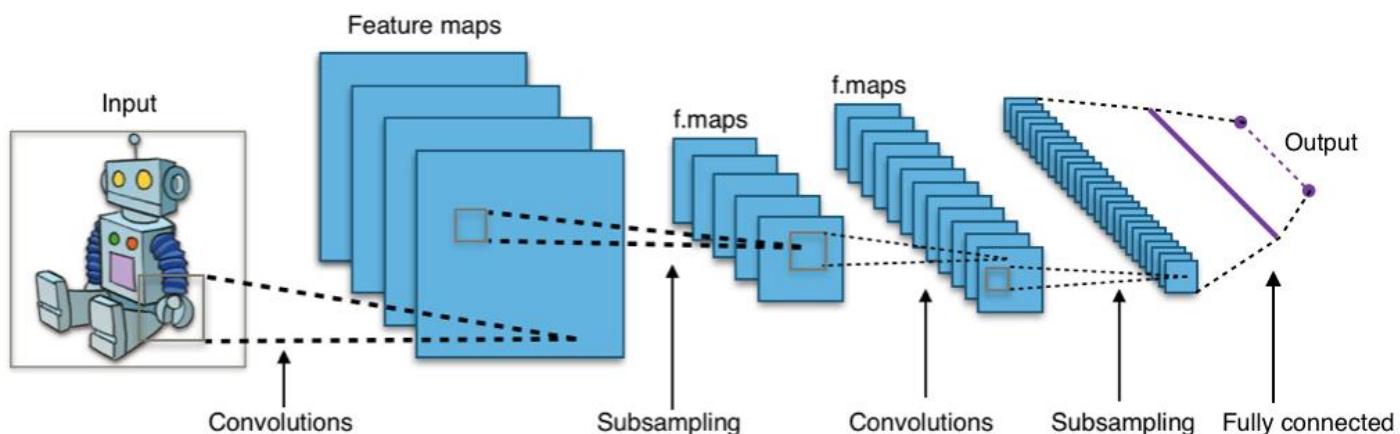


# Convolutional Neural Network



<https://www.kdnuggets.com/2018/09/object-detection-image-classification-yolo.html>

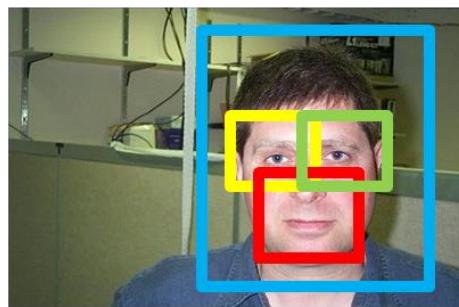


# Learning a Hierarchy of Feature Extractors

Goal: **Learn useful higher-level features** from images

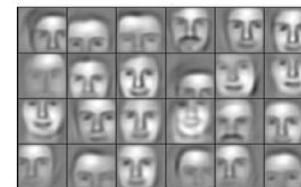
- Each layer of hierarchy extracts features from output of previous layer

Input data



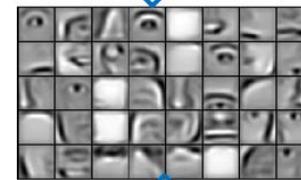
Lee et al., ICML 2009;  
CACM 2011

Feature representation



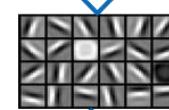
3rd layer  
“Objects”

High-level  
features



2nd layer  
“Object parts”

Mid-level  
features



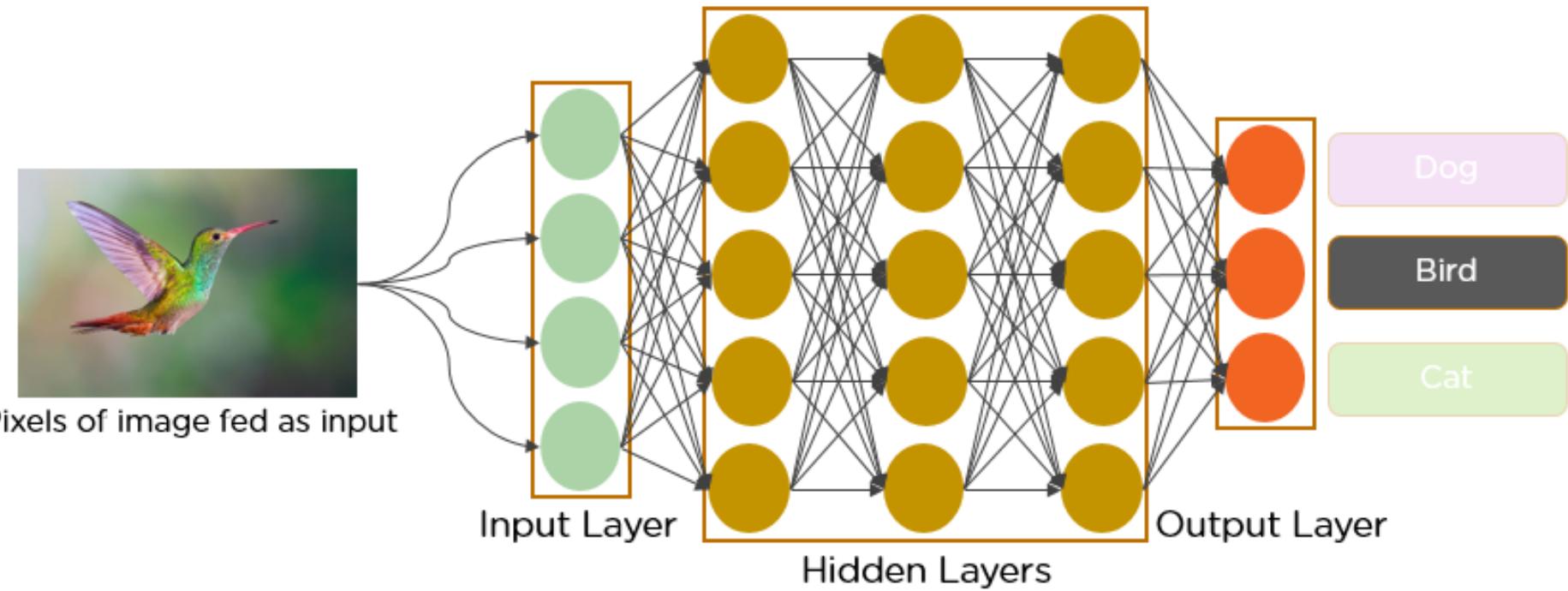
1st layer  
“Edges”

Low-level  
features



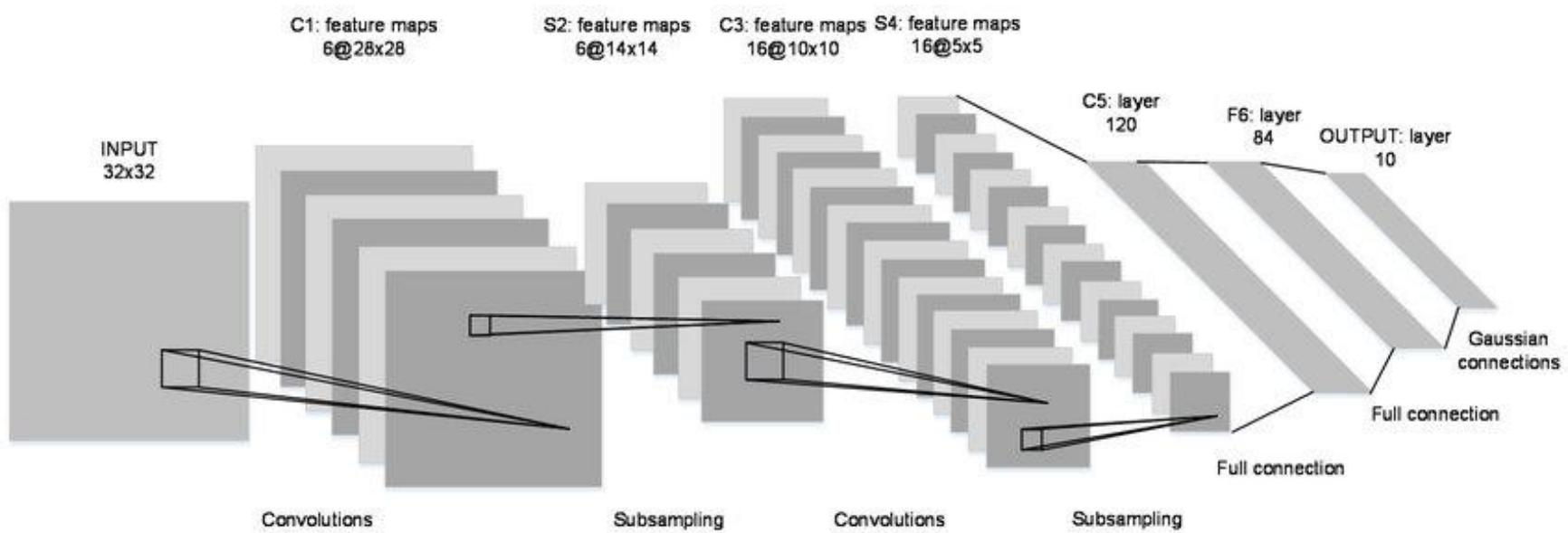
Pixels

# Convolutional Neural Network (CNN)



<https://www.analyticsvidhya.com/blog/2021/05/convolutional-neural-networks-cnn/>

# Background of CNNs

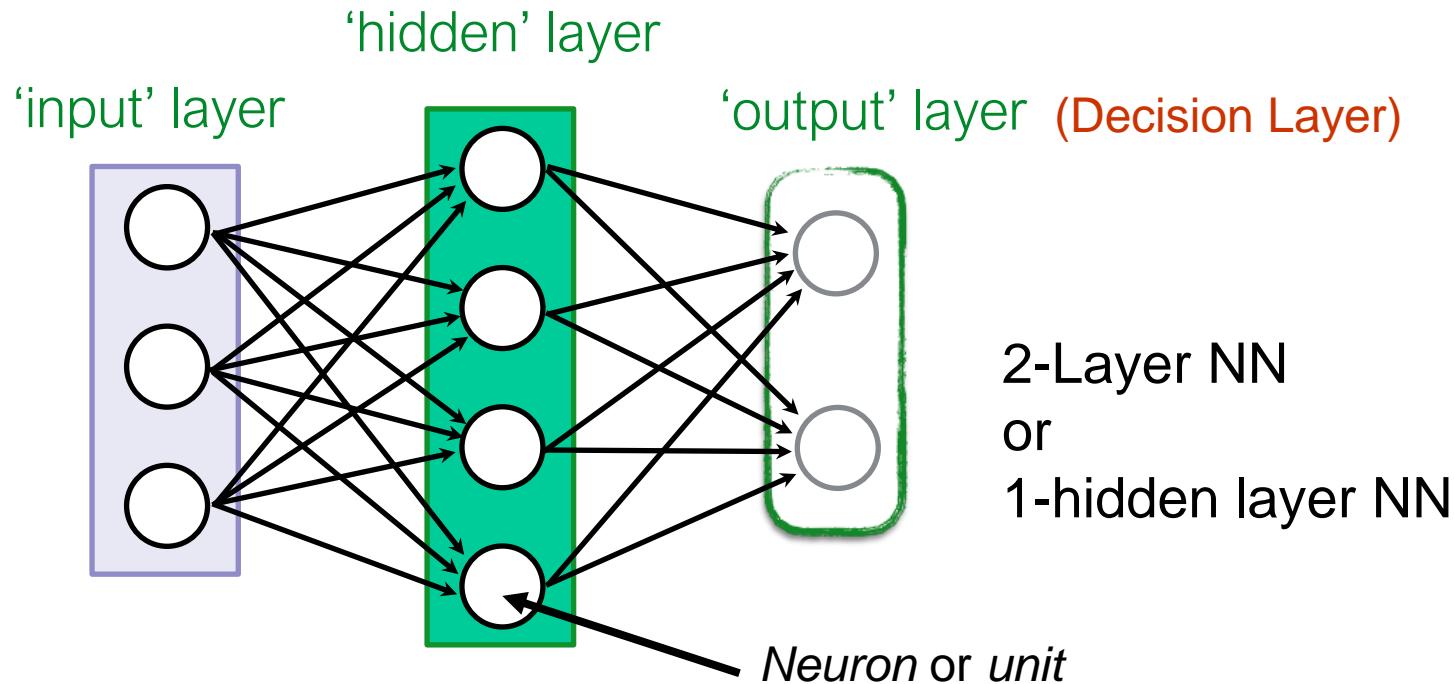


[https://www.researchgate.net/figure/The-LeNet-5-Architecture-a-convolutional-neural-network\\_fig4\\_321586653](https://www.researchgate.net/figure/The-LeNet-5-Architecture-a-convolutional-neural-network_fig4_321586653)

Connect a bunch of perceptrons together ...

# Neural Network

a collection of connected perceptrons



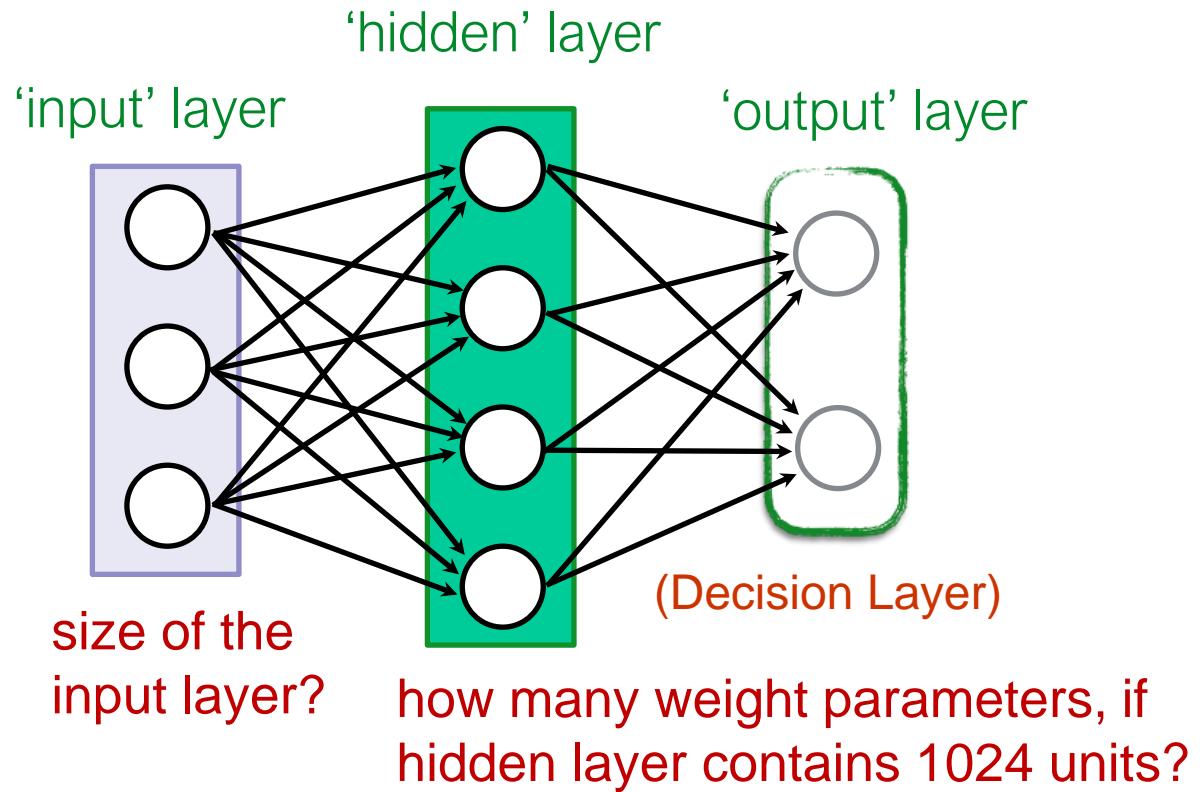
...also called a Multi-layer Perceptron (MLP)

# Neural Network for Image Data

- Suppose the input vector is an image.
- You have to classify if the picture contains a bird or not.



$256 \times 256$



# Image Data

High dimensionality

Patterns

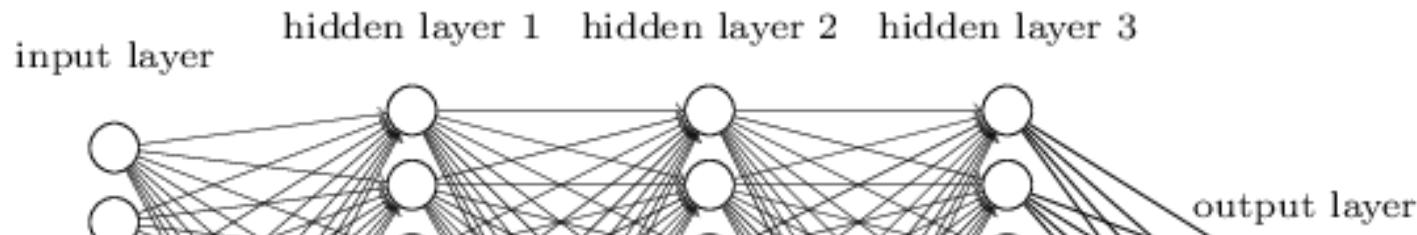
Has redundancies

112	112	113	115	124	137	147	151	155	155	155	155	156	155	150	146	136
112	114	116	120	130	143	151	152	154	155	156	158	157	154	151	150	146
112	116	119	124	134	147	153	152	153	156	163	168	168	162	158	157	154
116	119	120	120	130	142	149	150	154	168	175	166	165	171	168	156	164
120	125	129	130	135	144	152	155	162	171	178	166	165	173	175	170	171
122	129	136	137	140	146	154	158	167	168	166	162	163	170	177	180	171

Pixel values of a block of an image

# Smaller Network: CNN

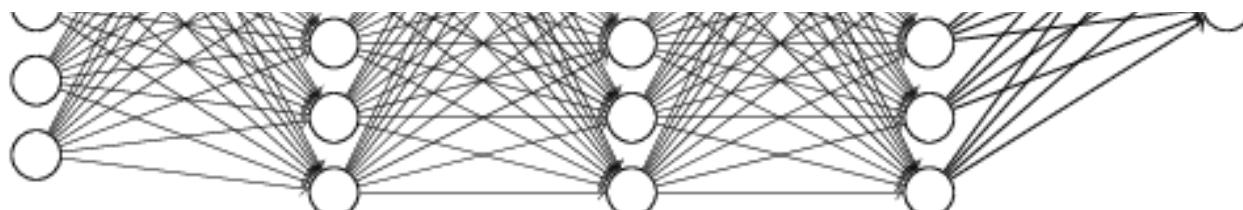
- Main questions:
- Do we need all the edges in this fully connected model?
- Can some of the weights be shared?



CNN are just neural networks

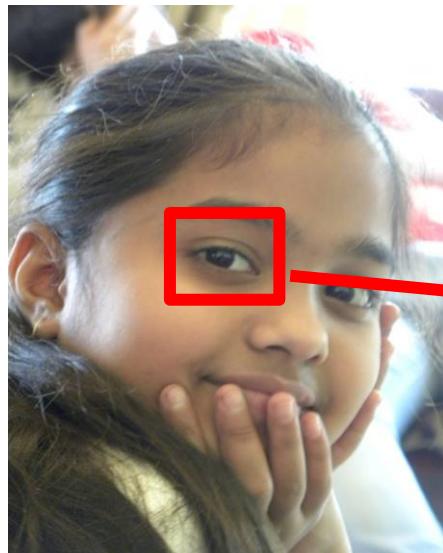
with

3D activations and weight sharing



# Image Understanding

- Patterns to learn can be very small compared to the image size.
- Example: Eye detector filter in images.



A small pattern detector can be used in place of the entire weight parameter matrix.



“eye” detector

# Pattern Learning

- Same pattern can be at different positions in images.
- For examples – eyes.



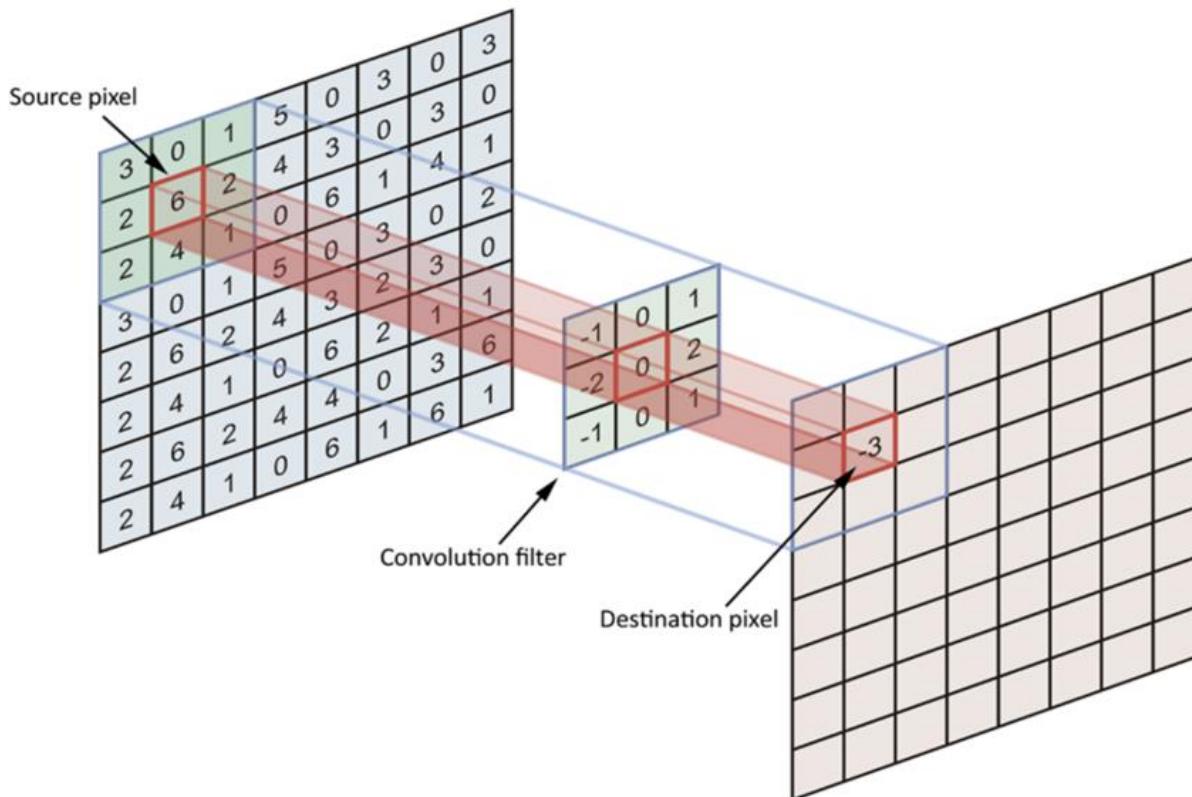
Same pattern detector can move around and detect eyes at different positions !

# Convolution Neural Networks (CNN)

- Use parameter sharing.
- Small pattern detectors called **filters** are used to convolve over the entire image.
- These filters are learned through NN training in the same way as in fully connected networks.
- Just like a hidden layer in a fully connected layer, convolution layers are used in CNNs.
- To handle large size of image data, pooling layers are introduced.

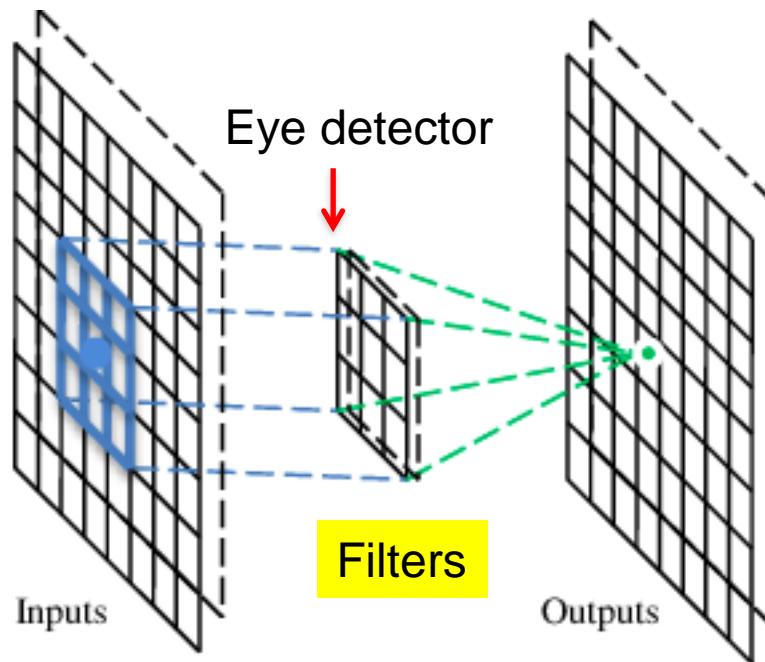
# Convolution Operation

- Convolution with an image block outputs a real number.
- When we slide this filter over the entire image, shifting block by block, we get an array of real numbers.



# Convolutional Layer

A convolutional layer has a number of filters that perform convolutional operation.



# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

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

# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

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

# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

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

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

# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

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

$$x = \begin{bmatrix} 5 \\ 1 \\ 7 \\ 7 \\ 2 \\ 5 \\ 5 \\ 3 \\ 5 \\ 3 \\ 3 \\ 4 \\ 7 \\ 8 \\ 1 \\ 5 \\ 7 \\ 8 \\ 2 \\ 3 \\ 1 \end{bmatrix} \quad W = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 1 \\ 1 \\ 0 \\ -1 \\ 1 \\ 0 \\ -1 \end{bmatrix}$$

- Convert the image block and filter to single column vectors and calculate their dot product.

$$\begin{aligned} W^T x &= 5 \times 1 + 1 \times 0 + 7 \times (-1) + 7 \times 2 + 2 \times 0 + 5 \times (-2) \\ &\quad + 5 \times 1 + 3 \times 0 + 5 \times (-1) = 2 \end{aligned}$$

# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

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

Slide the filter by one pixel  
and repeat the process.

$$\begin{aligned}1 \times 1 + 7 \times 0 + 7 \times (-1) + 2 \times 2 + 5 \times 0 + 7 \times (-2) \\+ 3 \times 1 + 5 \times 0 + 4 \times (-1) = -17\end{aligned}$$

# Convolution Layer Example

5	1	7	7	8	1
7	2	5	7	8	1
5	3	5	4	4	2
4	4	3	3	2	1
3	3	4	7	8	1
5	7	8	2	3	1

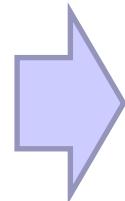
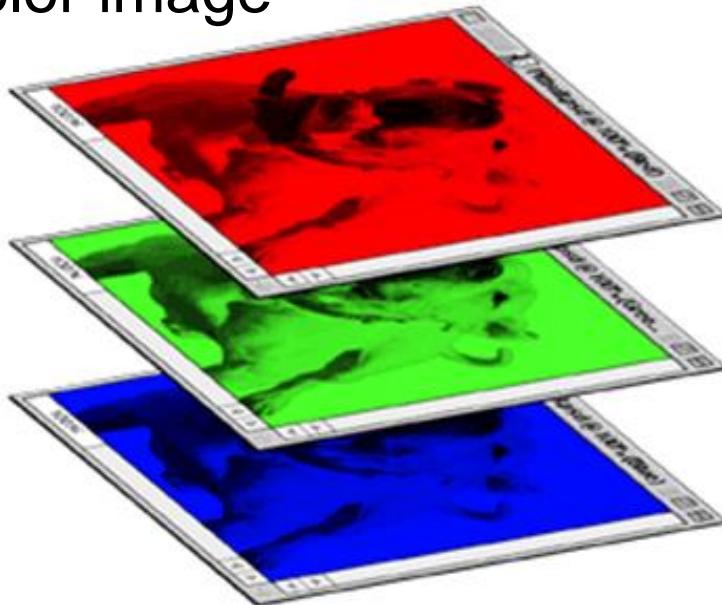
1	0	-1
2	0	-2
1	0	-1

and the process  
continues in this way...

$$7 \times 1 + 7 \times 0 + 8 \times (-1) + 5 \times 2 + 7 \times 0 + 8 \times (-2) + \\ 5 \times 1 + 4 \times 0 + 4 \times (-1) = -6$$

# Convolution on RGB Images

Color image



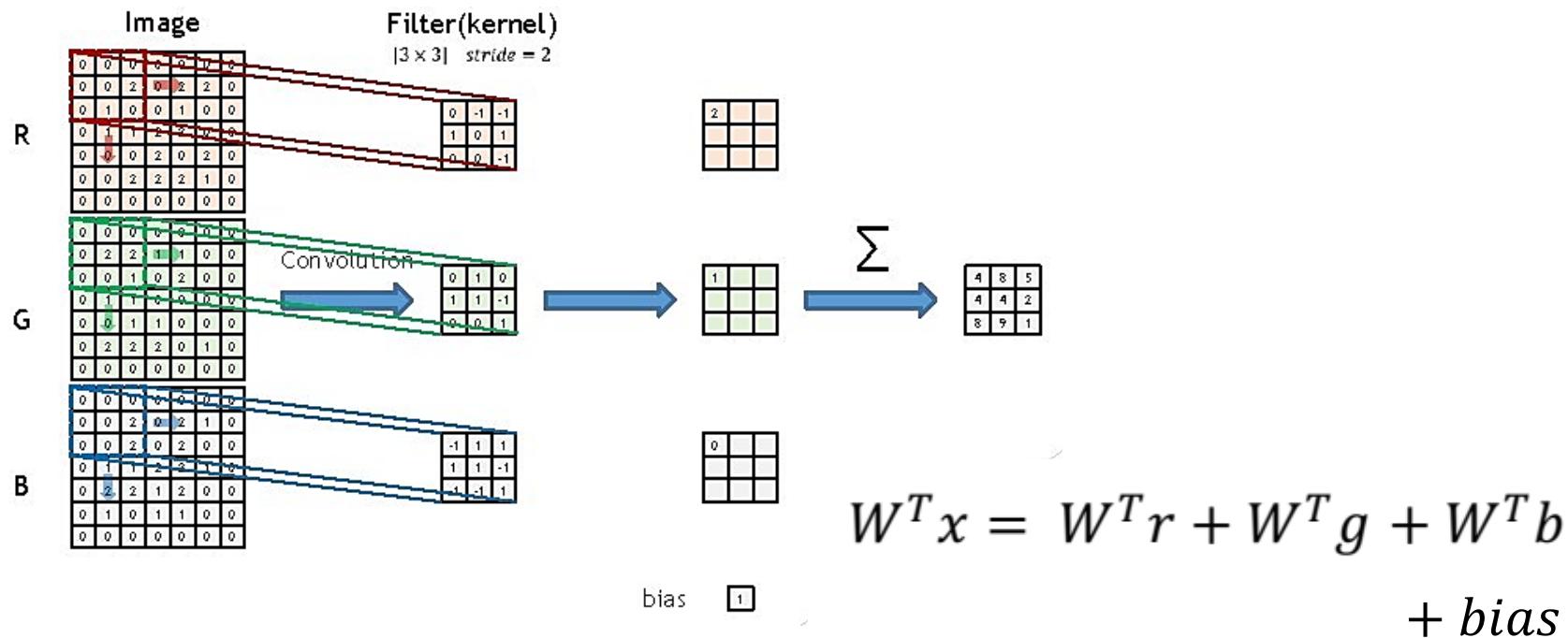
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Filter

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

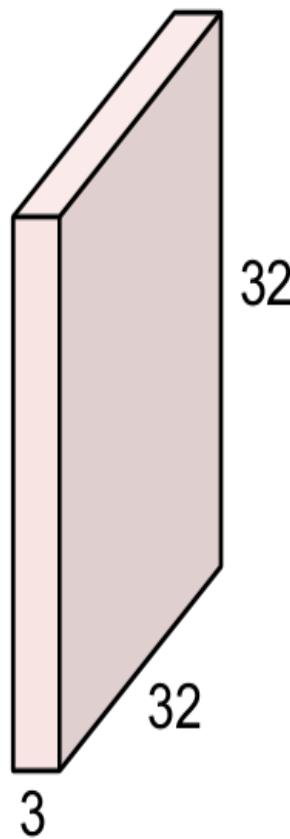
# Convolution on RGB Images

- Make three copies of the filter and apply convolution operation on each channel.
- You will get three values, one for each channel.
- Add all three values of filter mapping to get a single output.



# A Convolution Layer

32x32x3 image



5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

Figure source: Fei Fei Li's Lecture slides, Stanford University

# A Convolution Layer

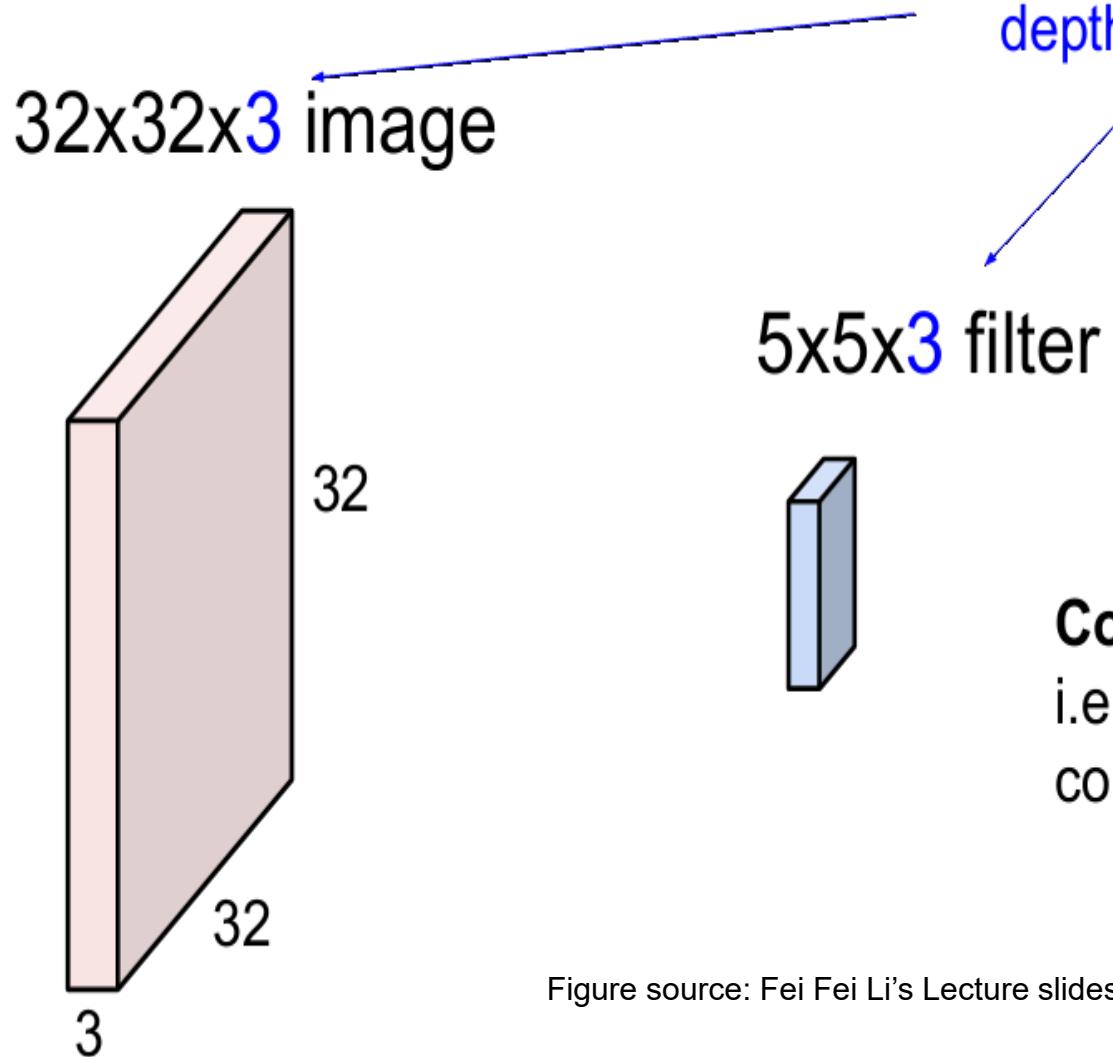


Figure source: Fei Fei Li's Lecture slides, Stanford University

# A Convolution Layer

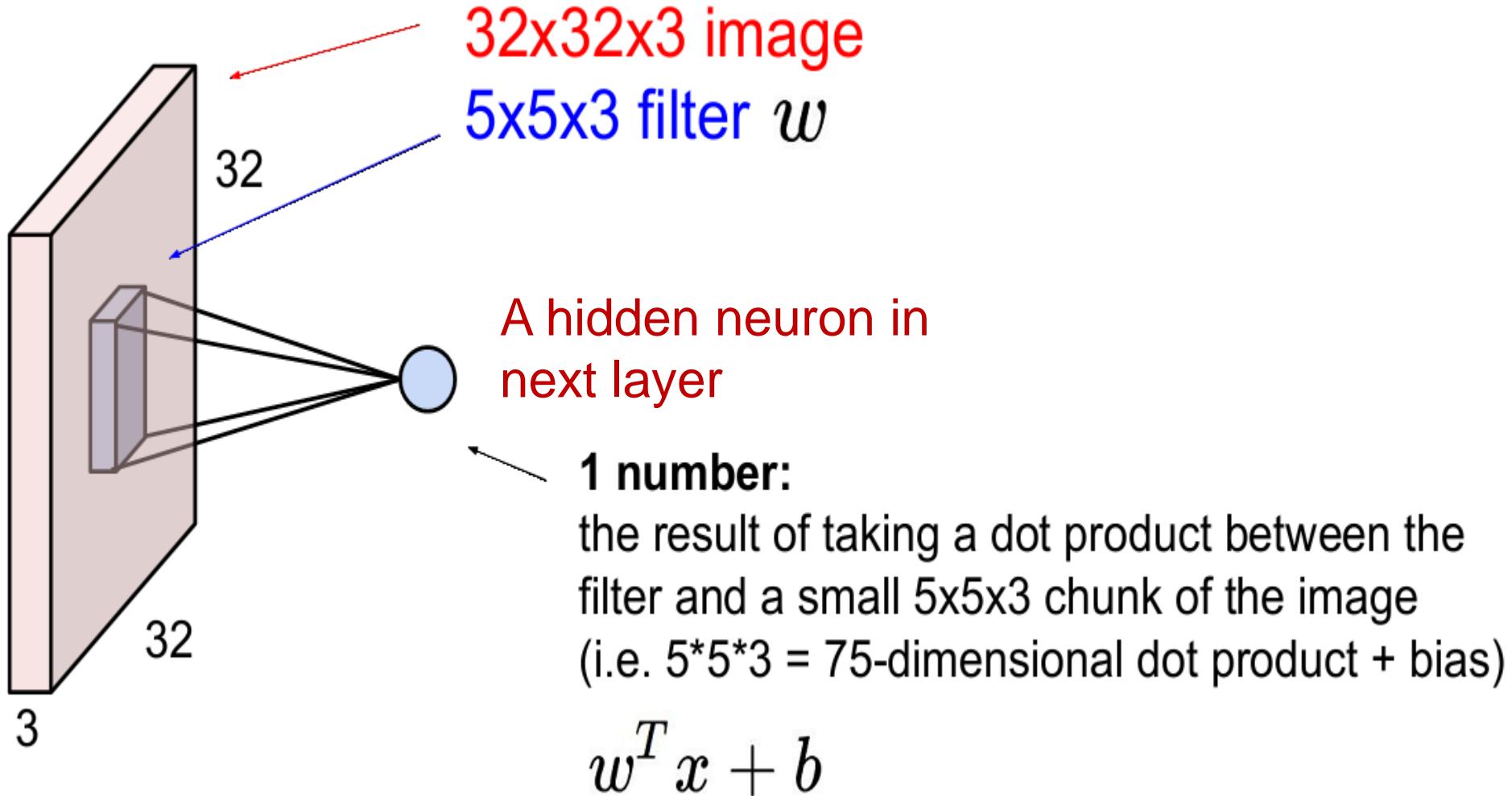
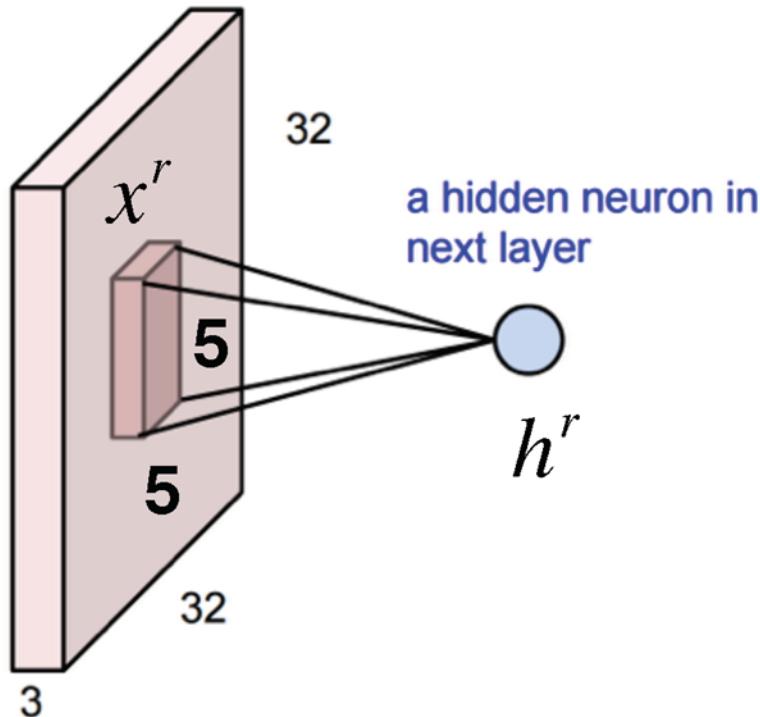


Figure source: Fei Fei Li's Lecture slides, Stanford University

# A Convolution Layer 3D Activations



Example: consider the region of the input " $x^r$ "

With output neuron  $h^r$

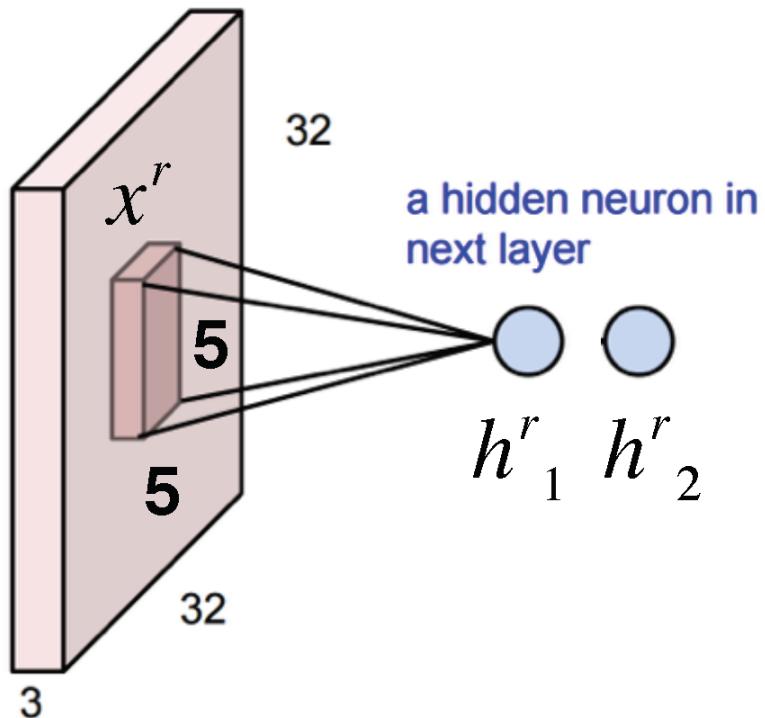
Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$



Sum over 3 axes

# A Convolution Layer 3D Activations

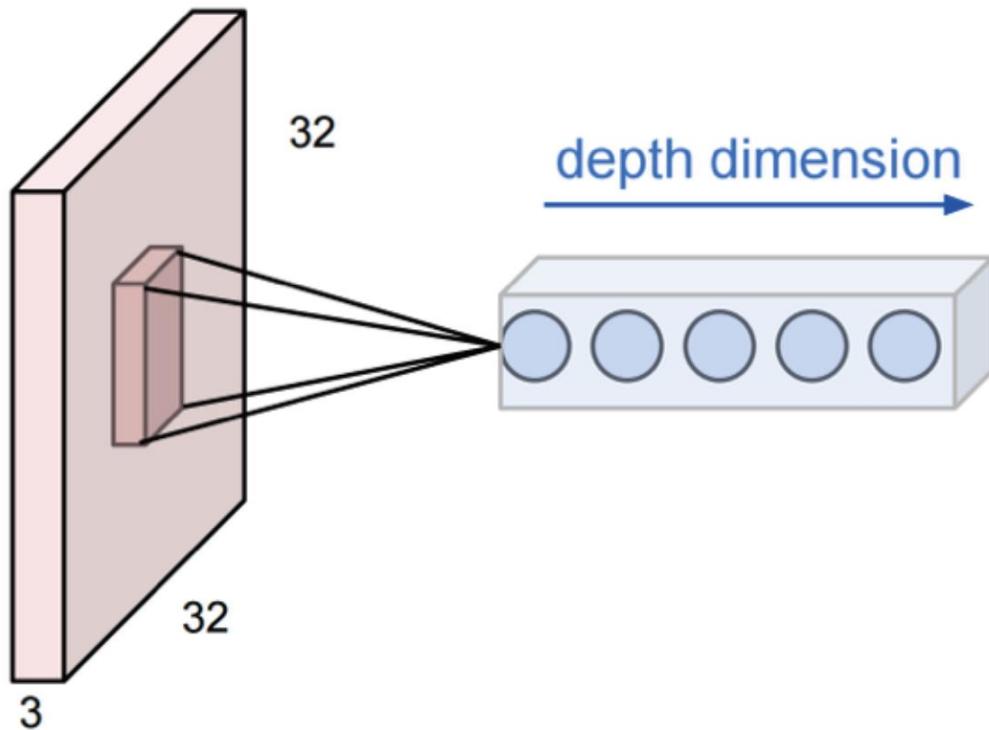


With **2** output neurons

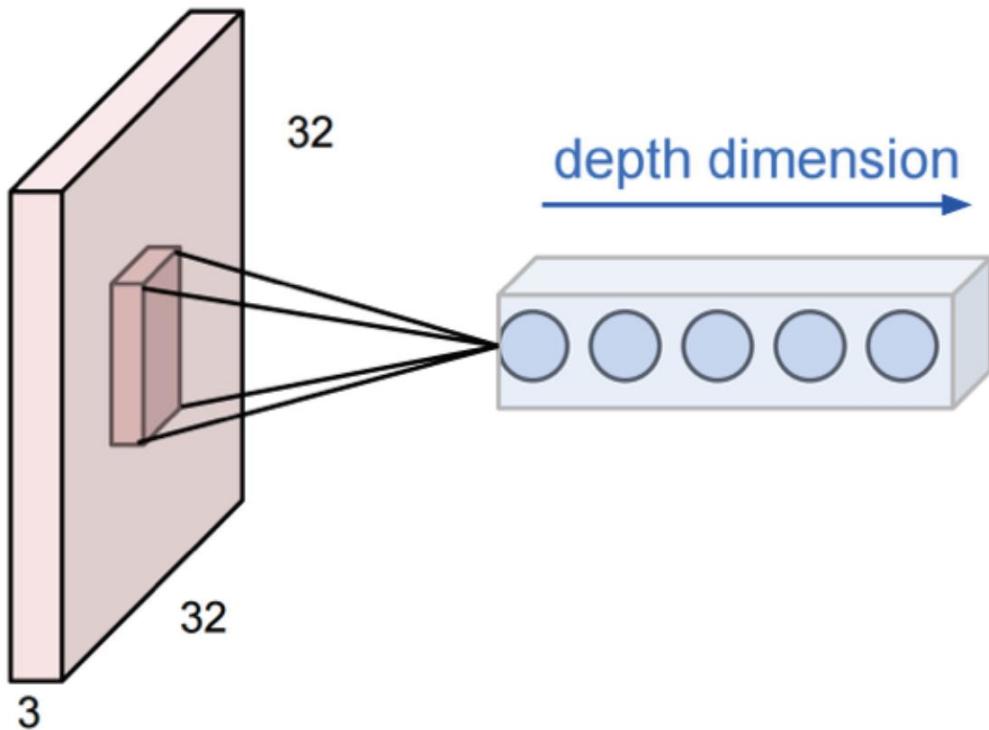
$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_1$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$

# A Convolution Layer 3D Activations



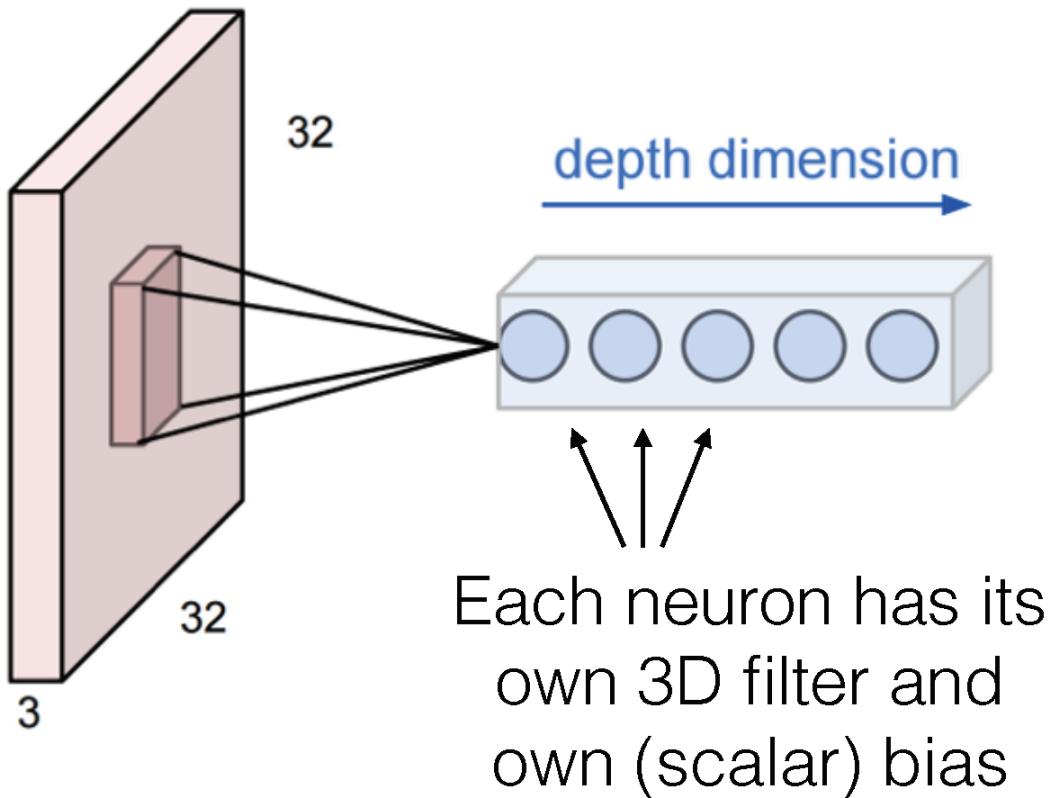
# A Convolution Layer 3D Activations



We can keep adding more outputs

These form a column in the output volume:  
[depth x 1 x 1]

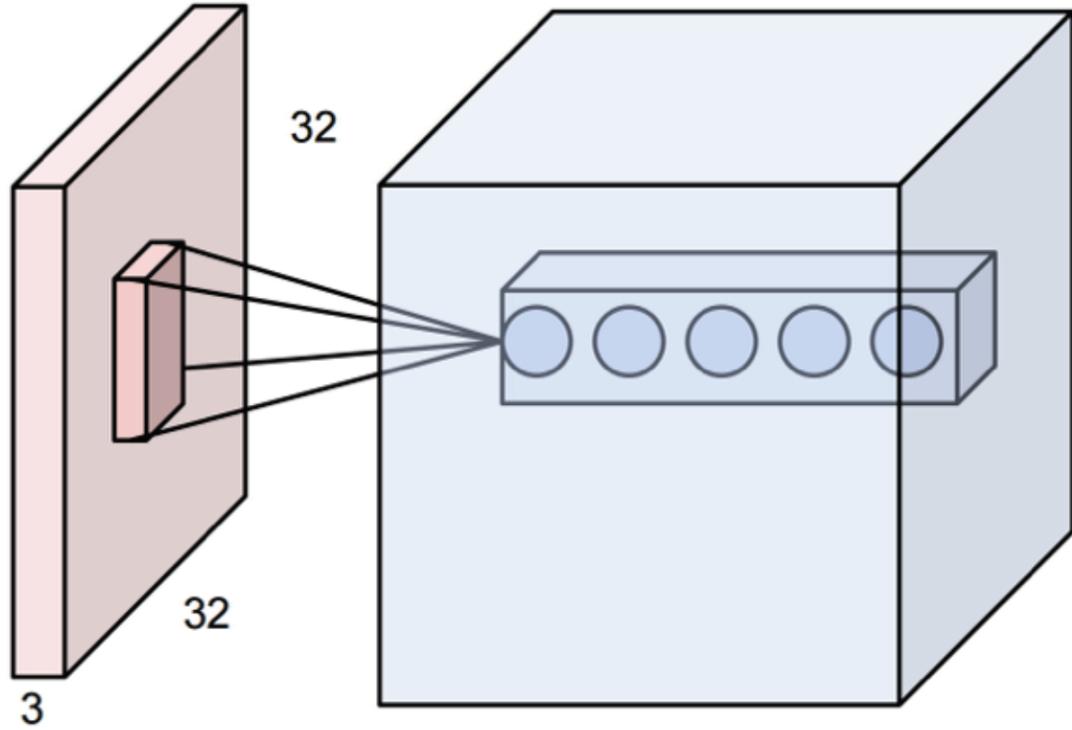
# A Convolution Layer 3D Activations



We can keep adding more outputs

These form a column in the output volume:  
[depth x 1 x 1]

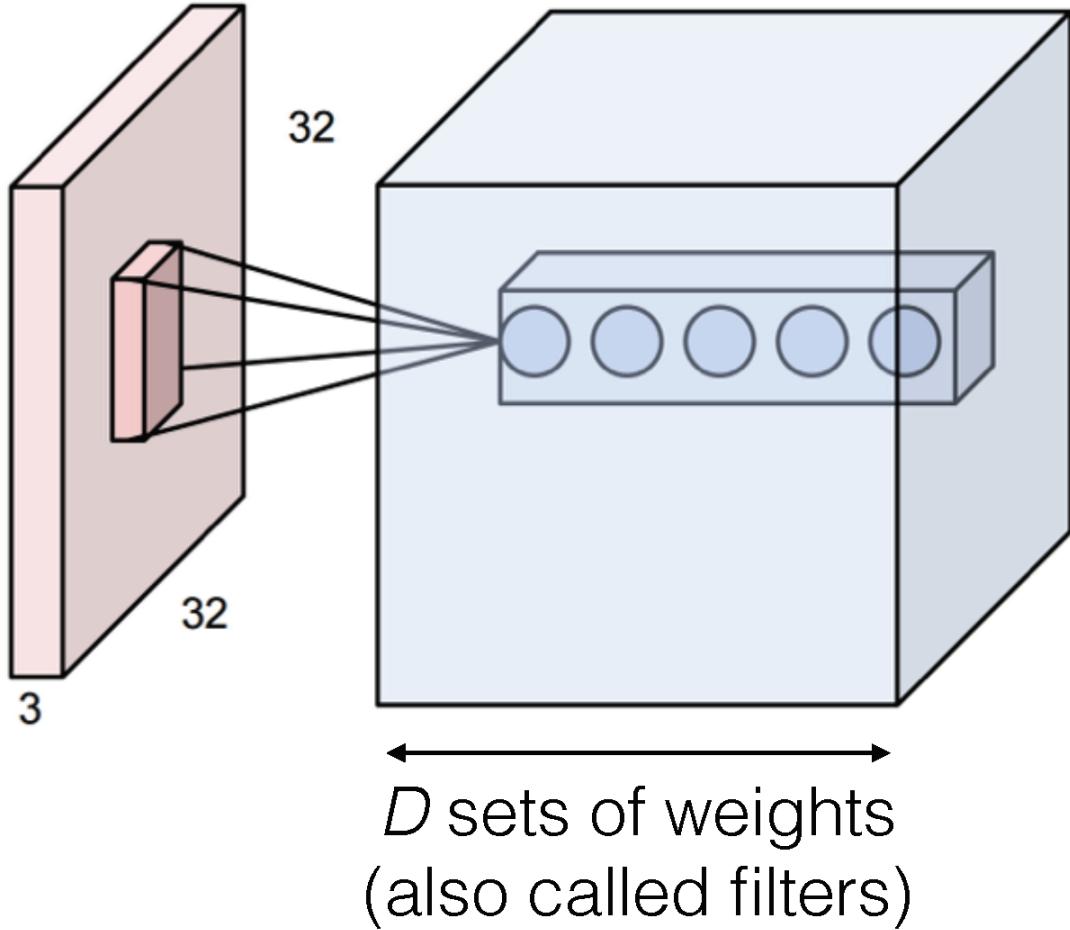
# A Convolution Layer 3D Activations



Now repeat this across the input

$D$  sets of weights  
(also called filters)

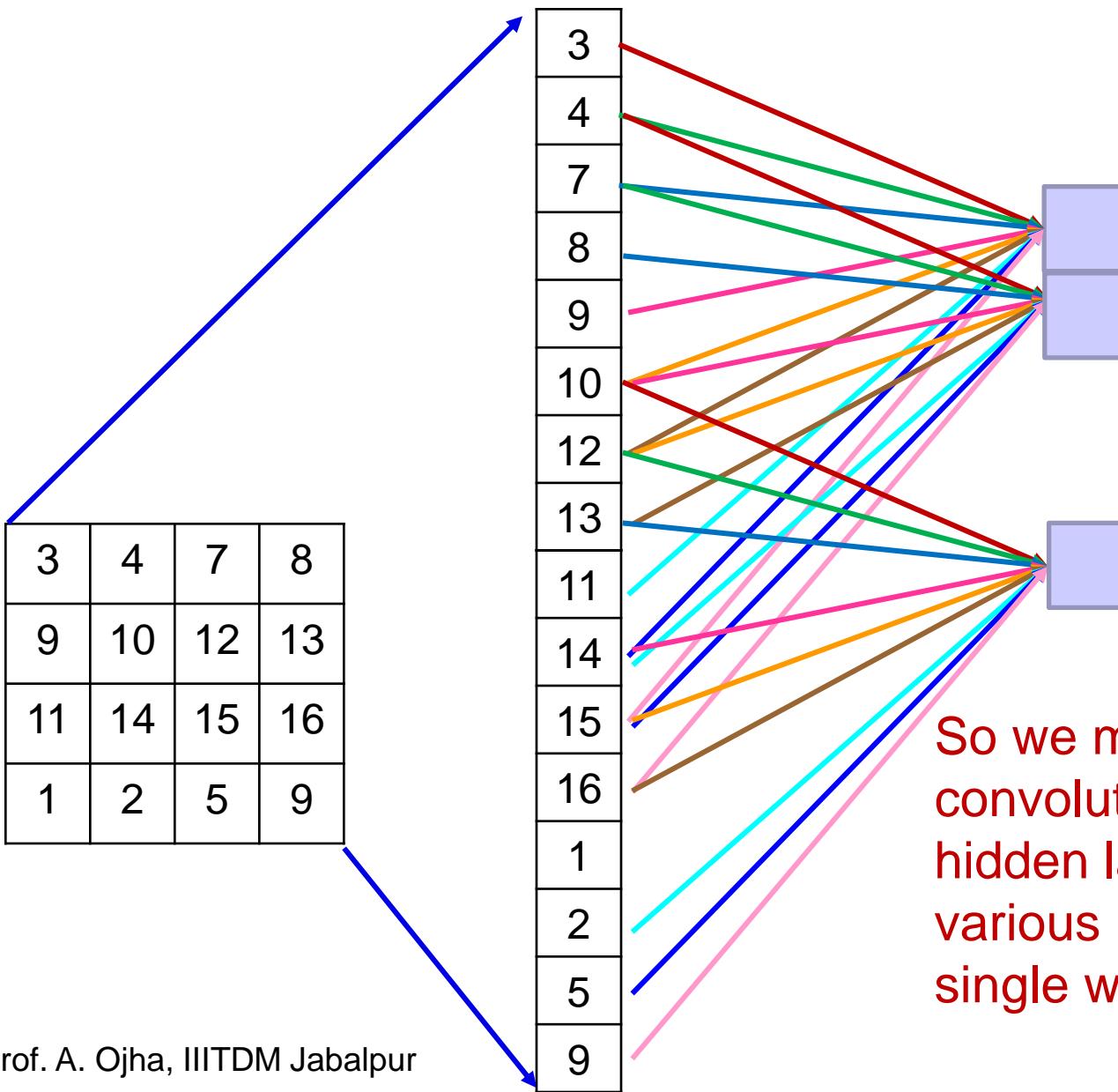
# A Convolution Layer 3D Activations



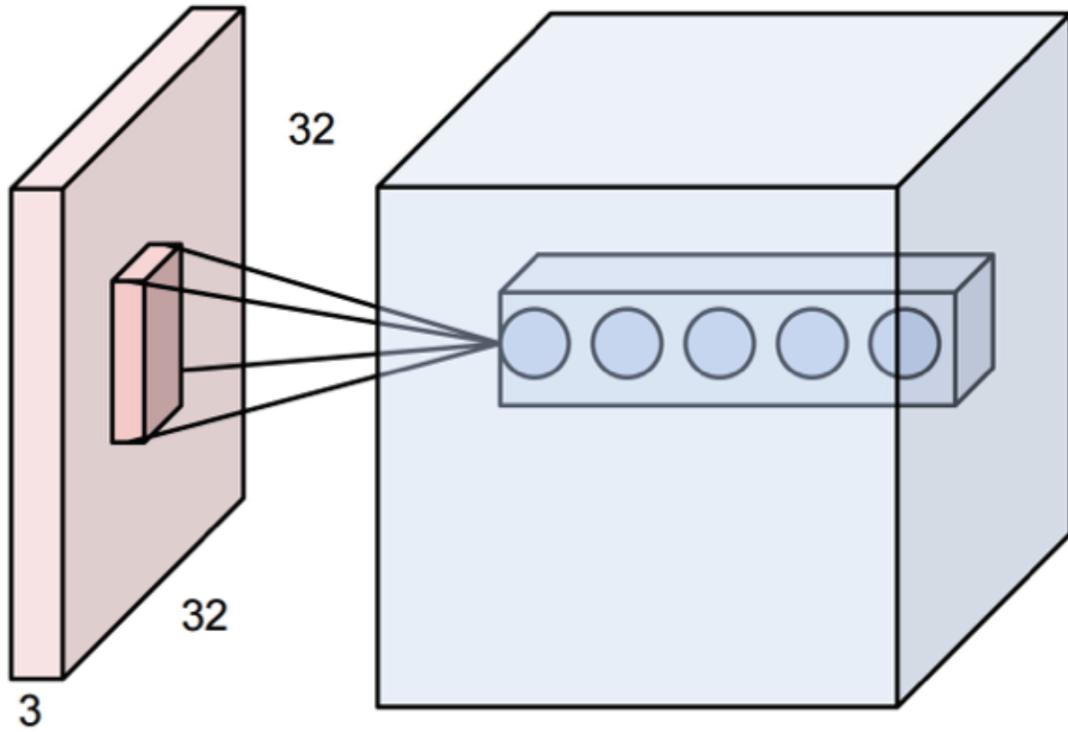
Now repeat this across the input

**Weight sharing:**  
Each filter shares  
the same weights  
(but each depth  
index has its own  
set of weights)

# Weight Sharing



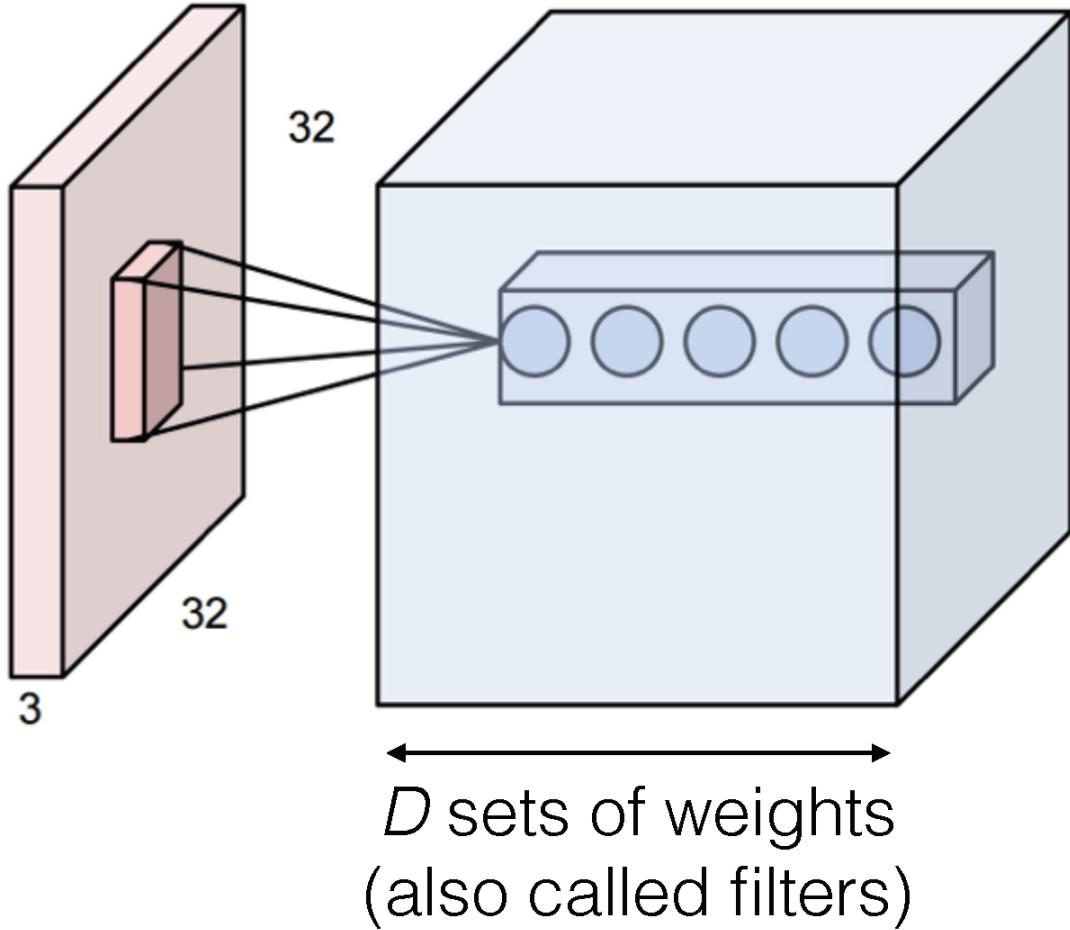
# A Convolution Layer 3D Activations



With weight sharing,  
this is called  
**convolution**

$D$  sets of weights  
(also called filters)

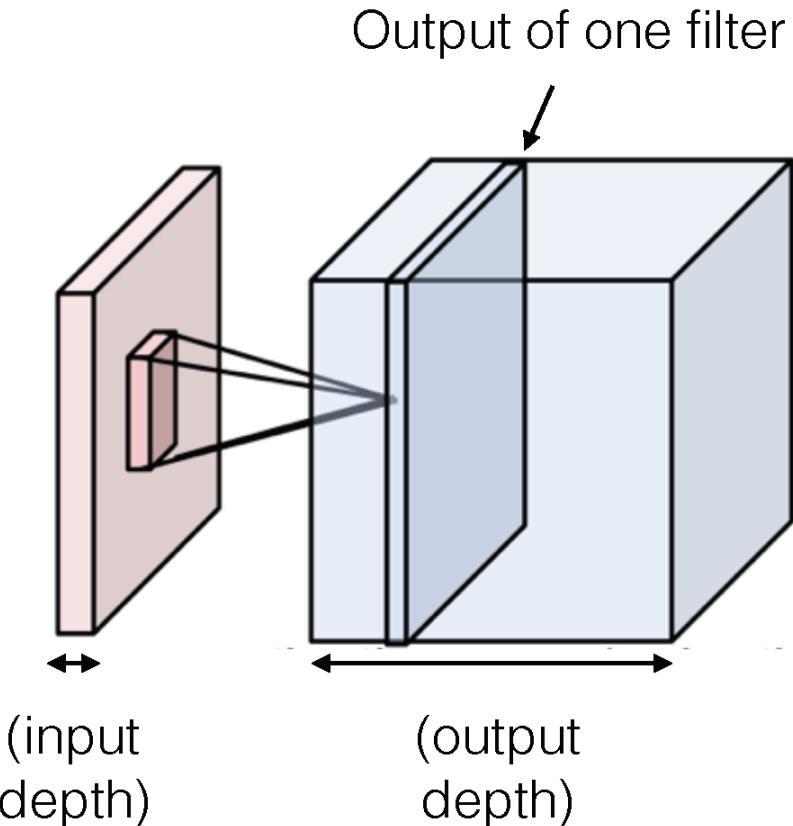
# A Convolution Layer 3D Activations



With weight sharing,  
this is called  
**convolution**

Without weight sharing,  
this is called a  
**locally connected layer**

# A Convolution Layer 3D Activations



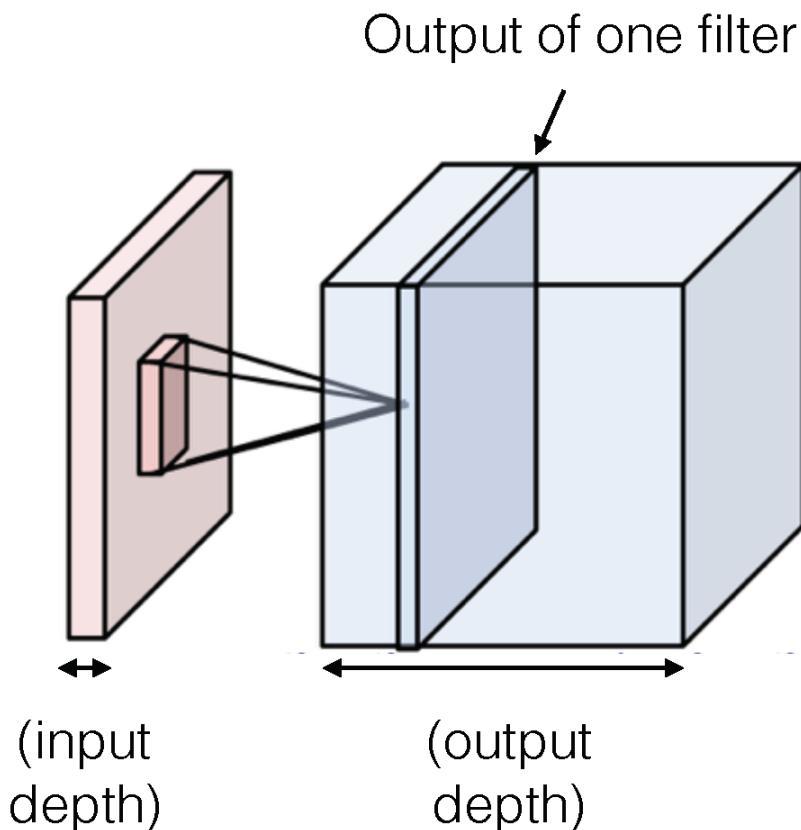
One set of weights gives one slice in the output

To get a 3D output of depth  $D$ , use  $D$  different filters

In practice, ConvNets use many filters ( $\sim 64$  to 1024)

# A Convolution Layer

## 3D Activations



One set of weights gives  
one slice in the output

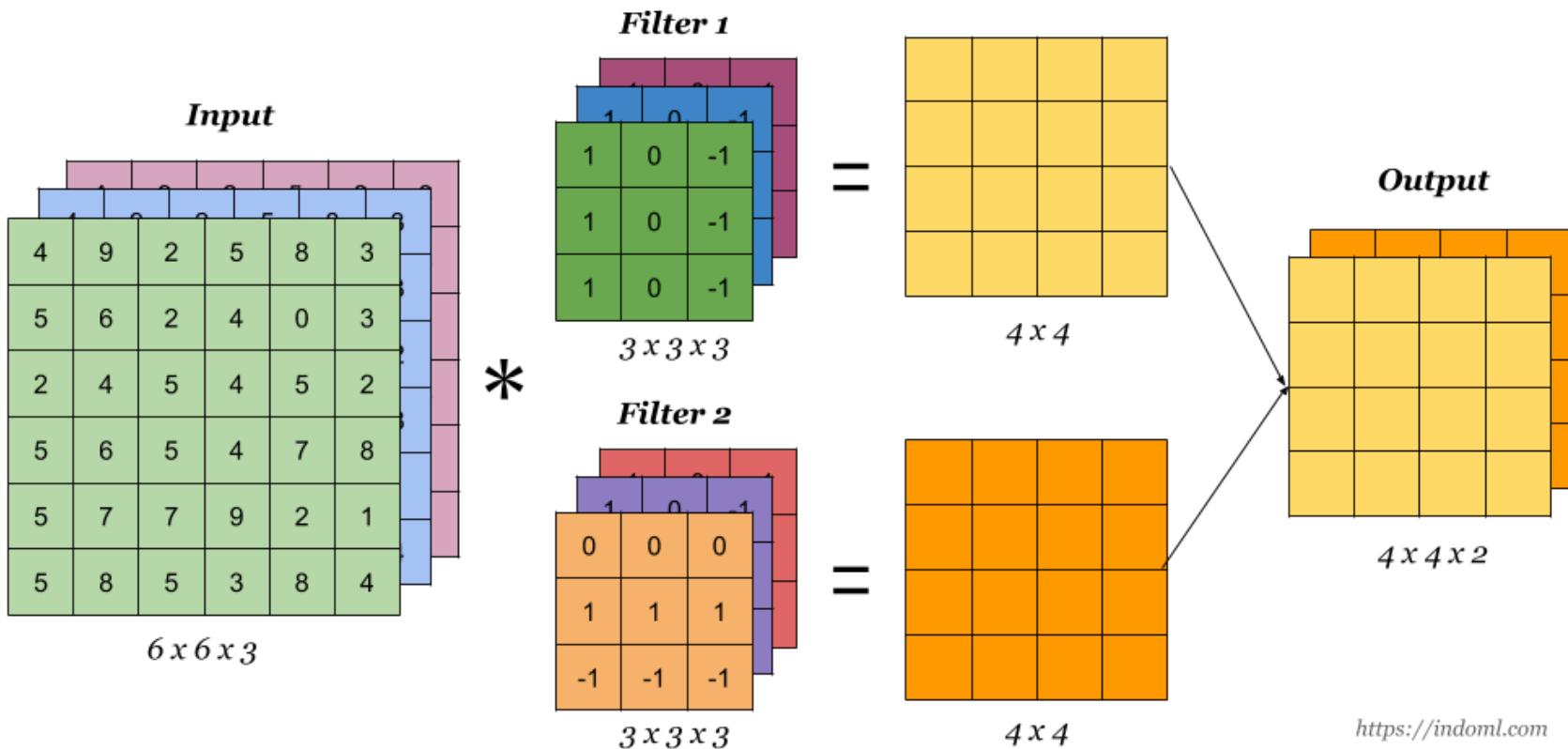
To get a 3D output of depth  $D$ ,  
use  $D$  different filters

In practice, ConvNets use  
many filters (~64 to 1024)

All together, the weights are **4** dimensional:  
(output depth, input depth, kernel height, kernel width)

# Feature Maps

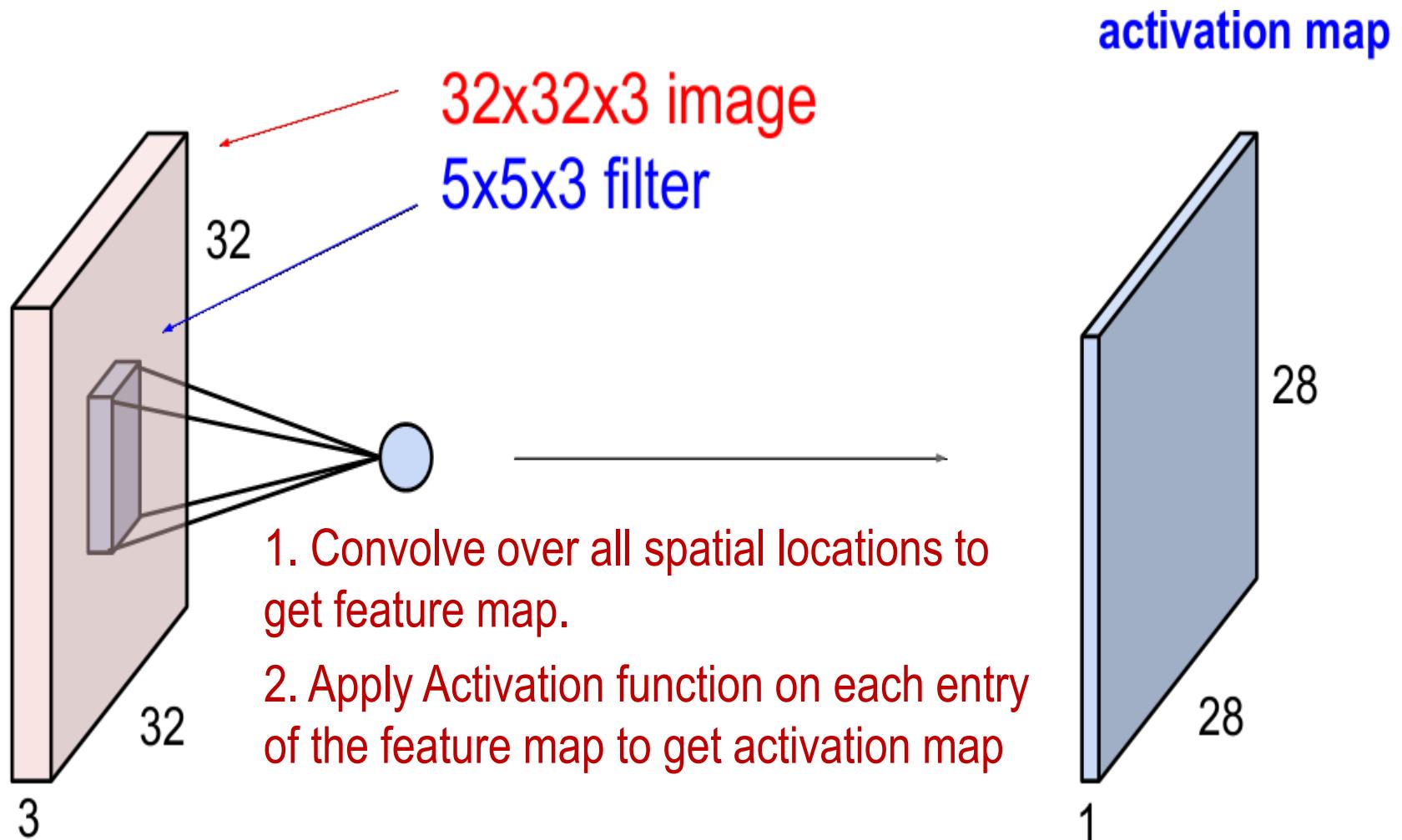
- With 2 filters, output will contain 2 channels, known as feature maps, one for each filter.
- Each feature map applies its own bias.



# Activation Maps

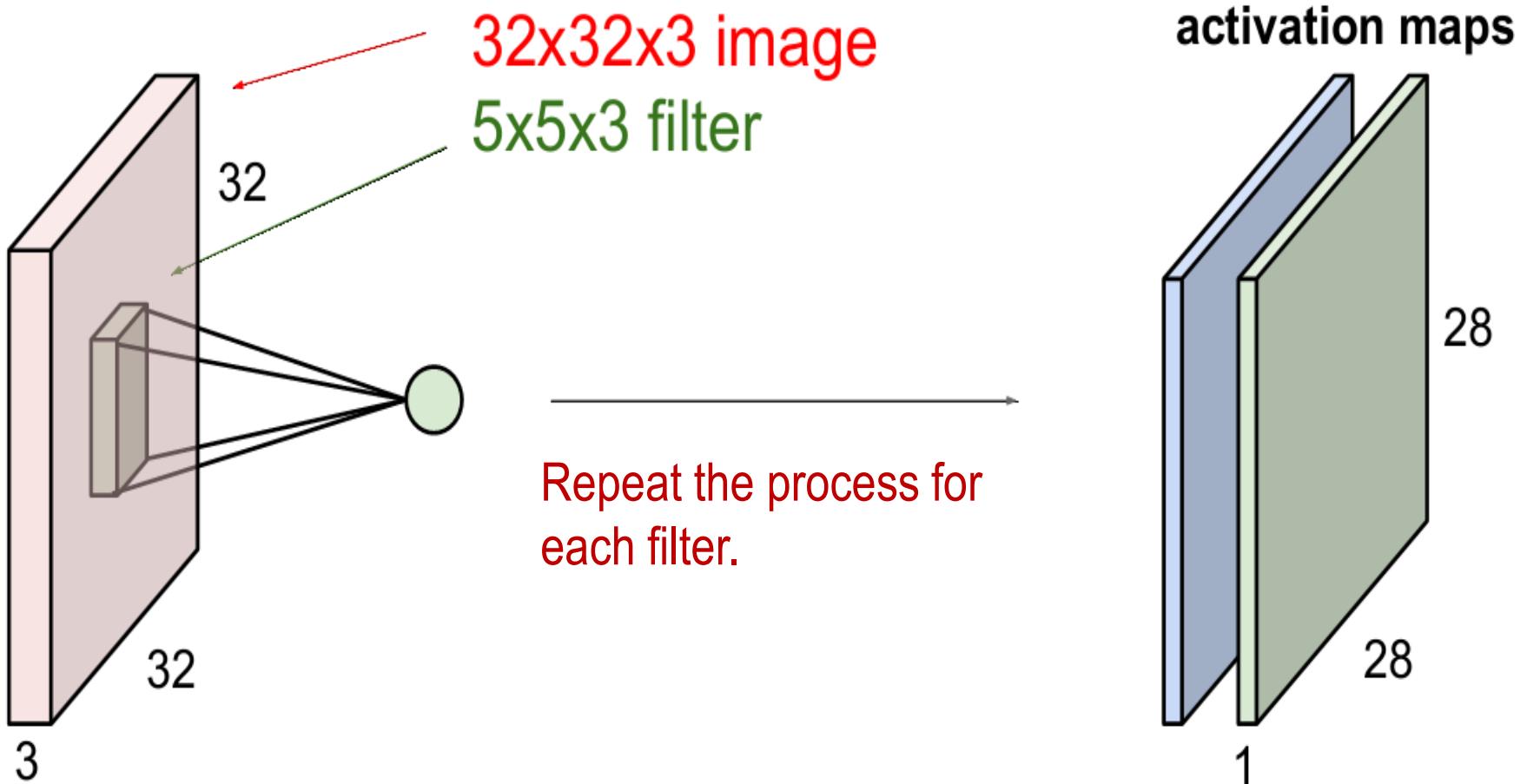
- An activation function is applied on all the entries of each feature map to get the activation maps.

# Filters as activation maps



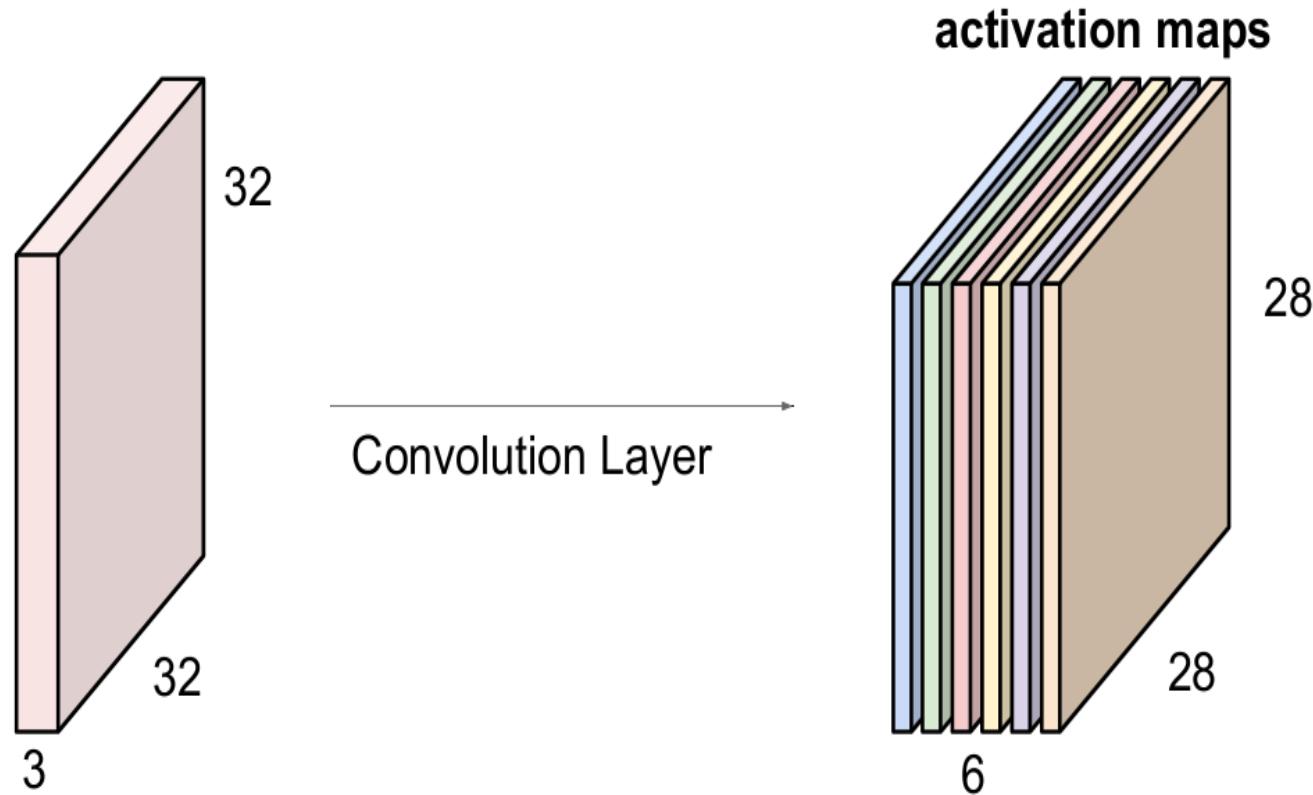
One set of weights gives one slice in the output.

# Filters as activation maps



# Filters as activation maps

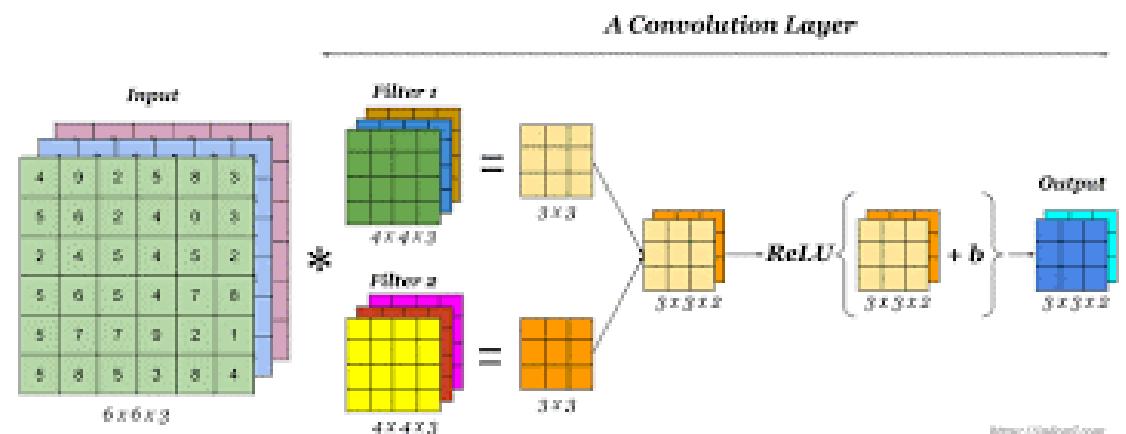
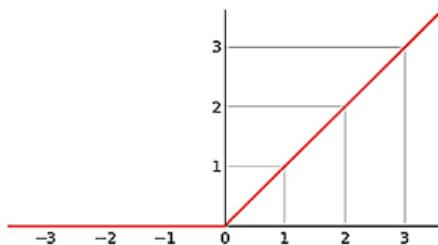
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

# Activation Maps

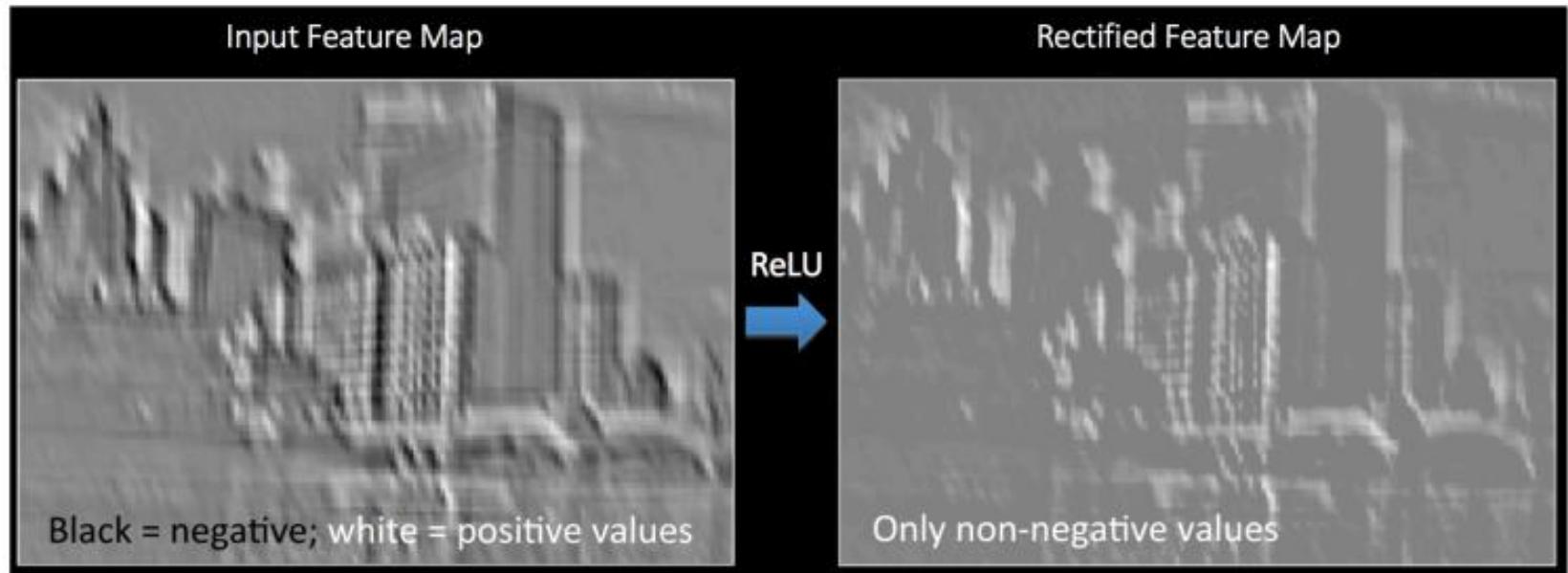
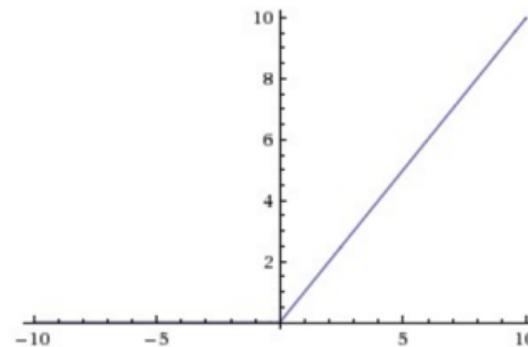
- An activation function is applied on all the entries of each feature map to get the activation maps.



$$a = f(Z) = f(W^T X + b)$$

# Activation Maps - Example

Output = Max(zero, Input)



# Filters and Activation Maps

- A CNN is a sequence of convolutional layers, interspersed with activation functions (and possibly other layer types)

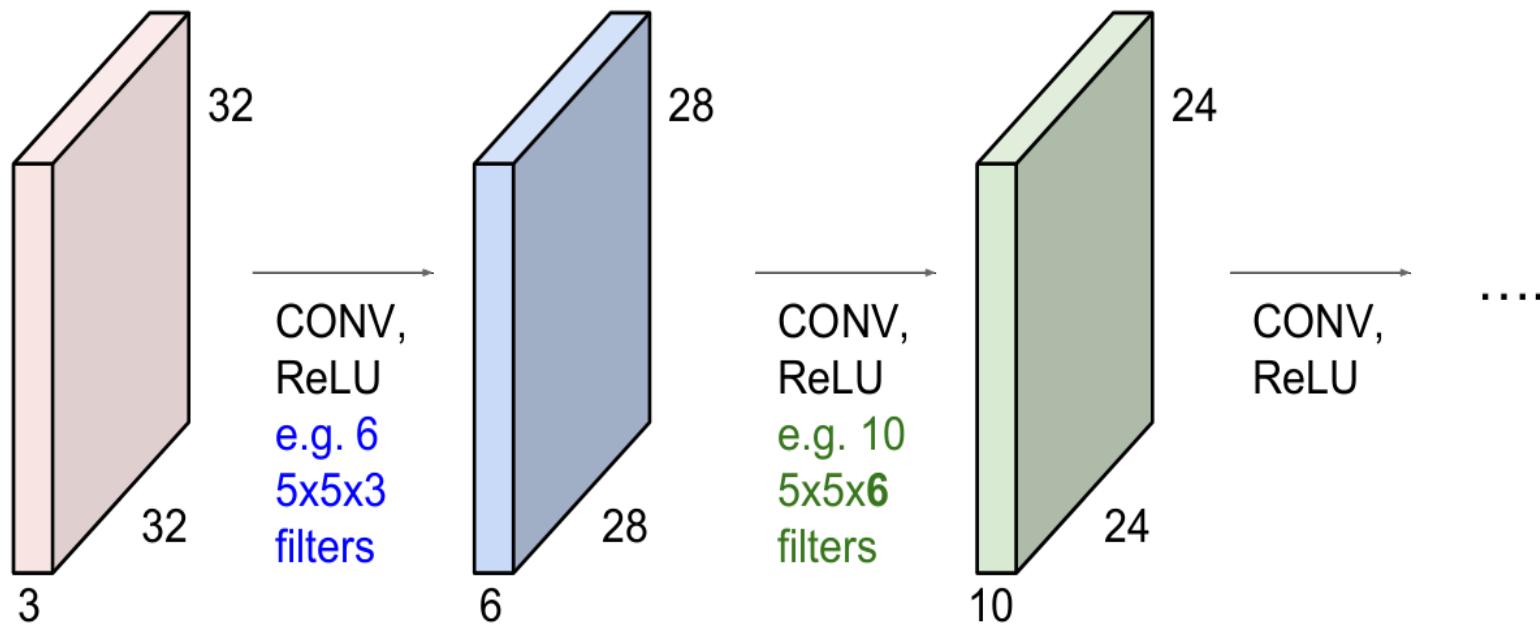
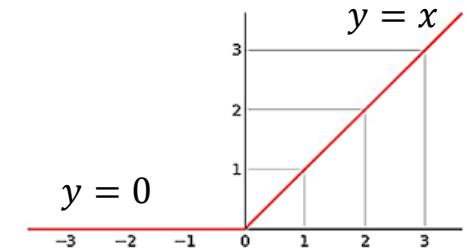


Figure source: Fei Fei Li's Lecture slides, Stanford University

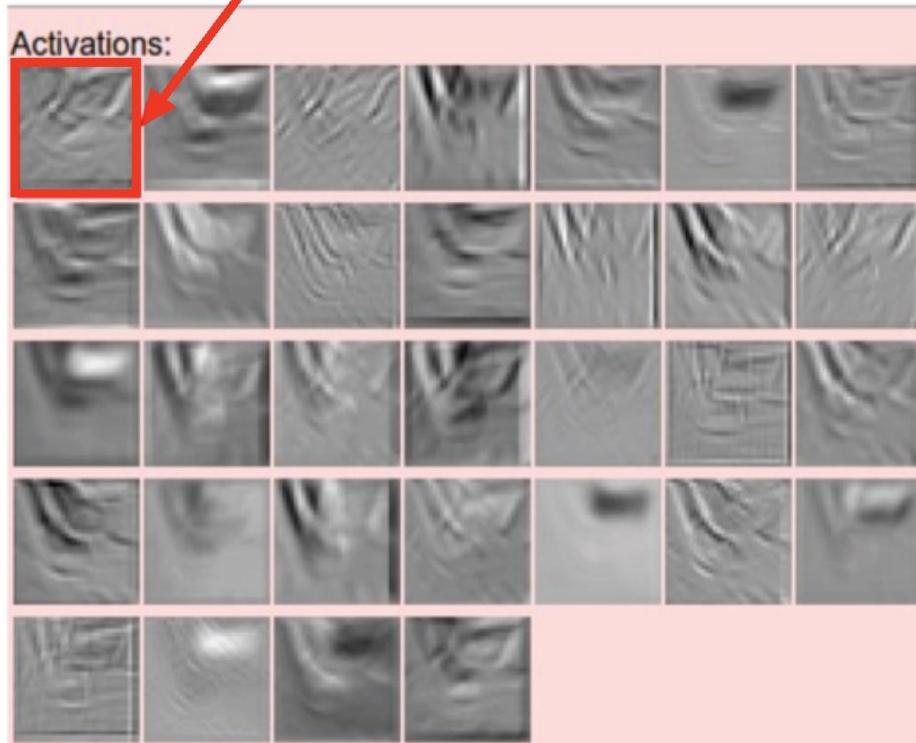
# Filters and Activation Maps

We can unravel the 3D cube and show each layer separately:  
(Input)



one filter = one depth slice (or activation map)

(32 filters, each 3x5x5)



We call the layer convolutional because it is related to convolution of two signals:

$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x - n_1, y - n_2]$$

elementwise multiplication and sum of a filter and the signal (image)

# Filters and Activation Maps

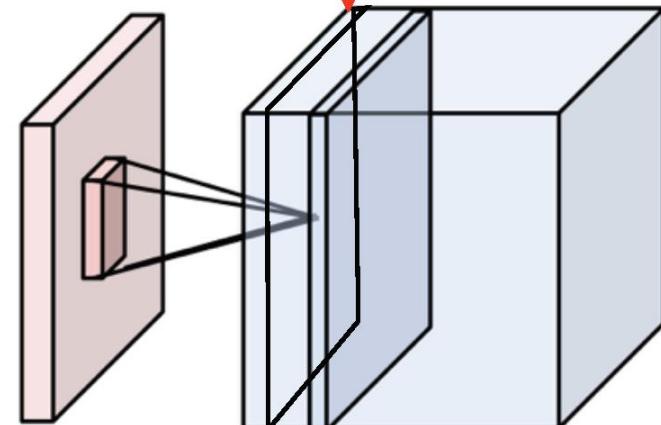
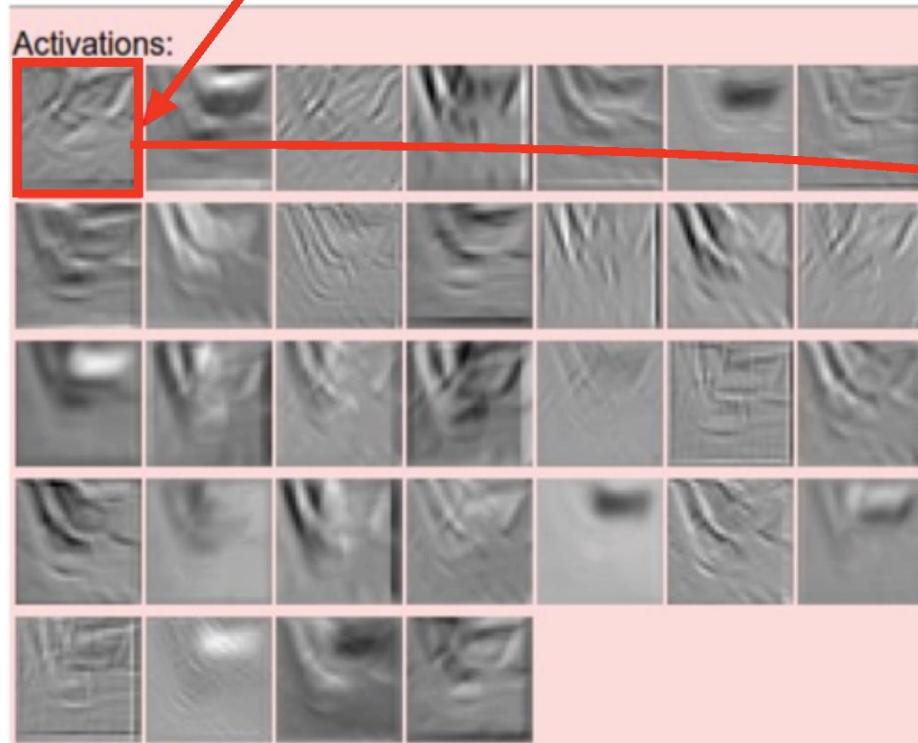
We can unravel the 3D cube and show each layer separately:

(Input)



one filter = one depth slice (or activation map)

(32 filters, each 3x5x5)



# Filters and Activation Maps

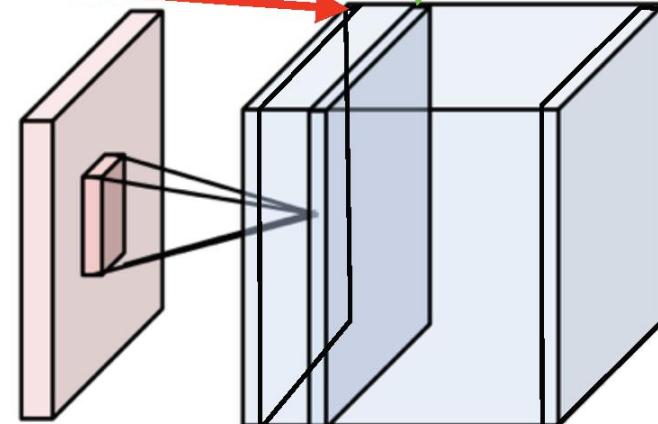
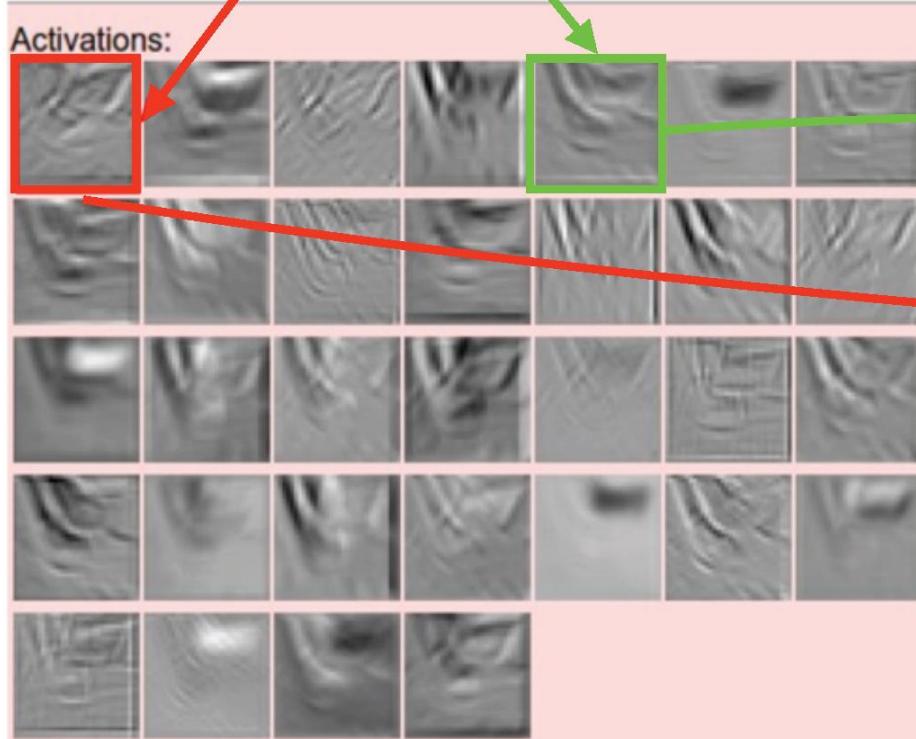
We can unravel the 3D cube and show each layer separately:

(Input)



one filter = one depth slice (or activation map)

(32 filters, each 3x5x5)



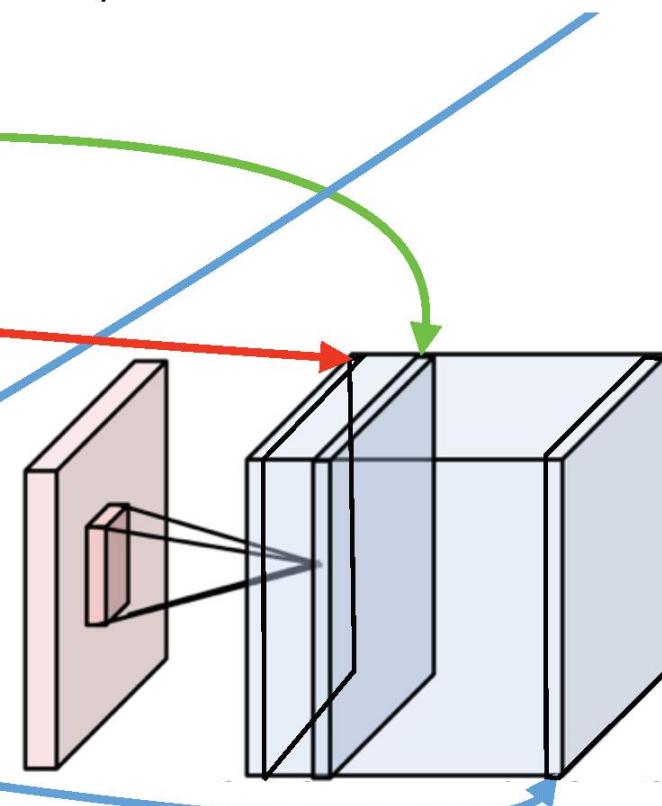
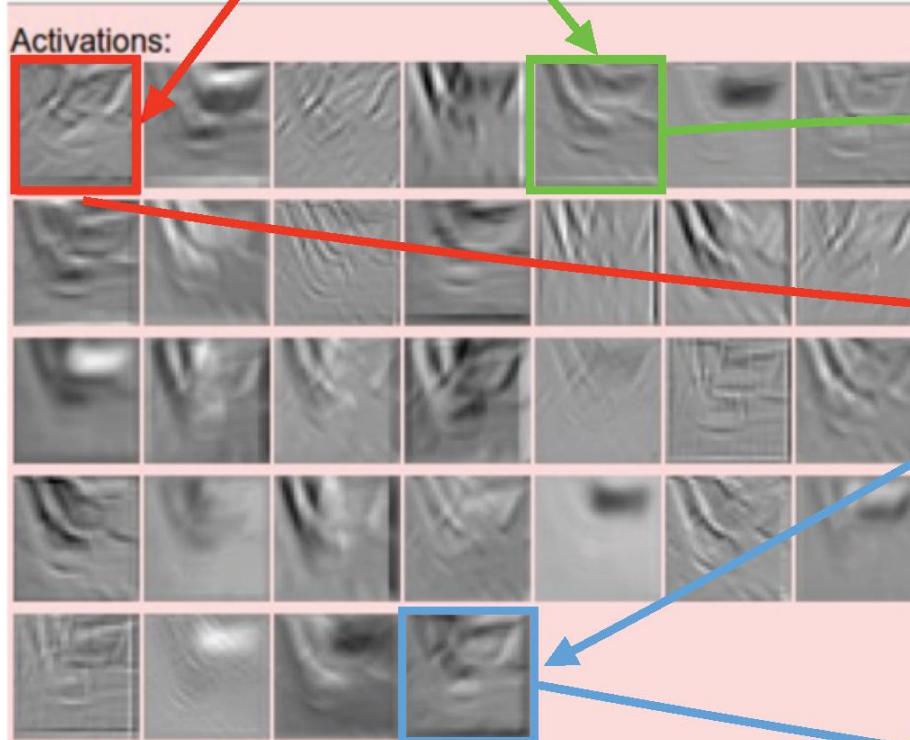
# Filters and Activation Maps

We can unravel the 3D cube and show each layer separately:

(Input)

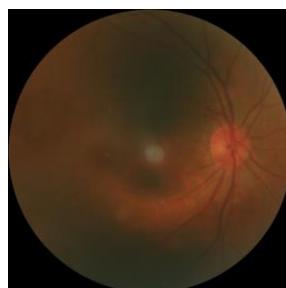


one filter = one depth slice (or activation map)  
(32 filters, each 3x5x5)



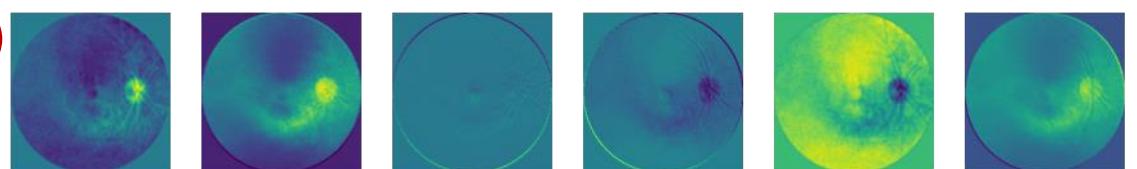
# Activation Maps: Example

Peking university international competition on ocular disease intelligent recognition (ODIR-2019)

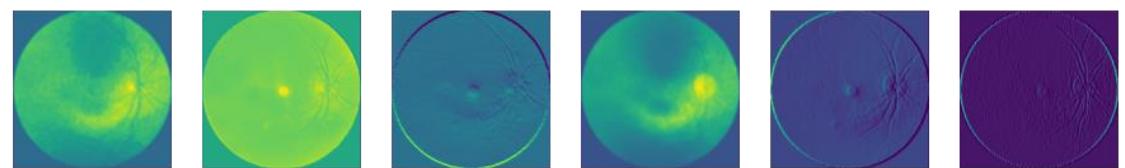


Actual fundus image  
with cataract in left eye

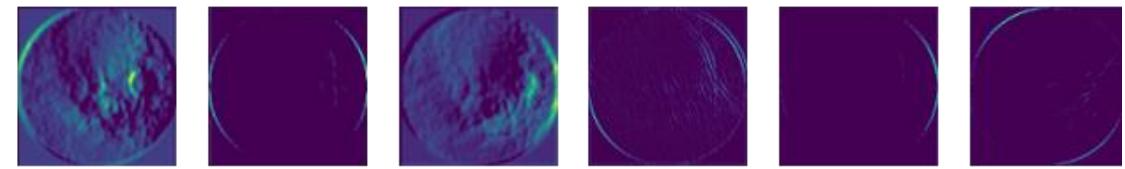
Level of abstraction  
increases at the high  
level layers



1<sup>st</sup> Convolutional Layer Activation Maps



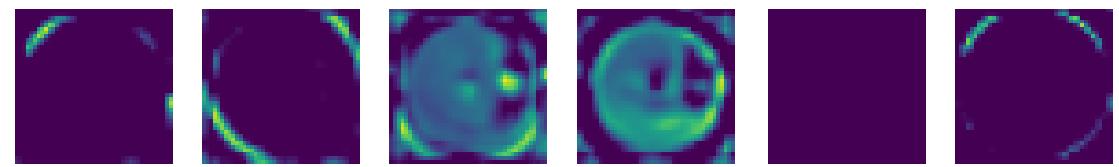
2<sup>nd</sup> Convolutional Layer Activation Maps



4<sup>th</sup> Convolutional Layer Activation Maps



7<sup>th</sup> Convolutional Layer Activation Maps

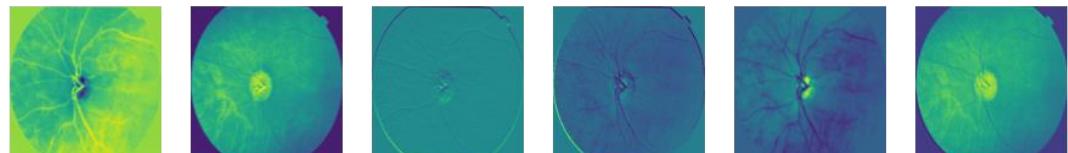


10<sup>th</sup> Convolutional Layer Activation Maps

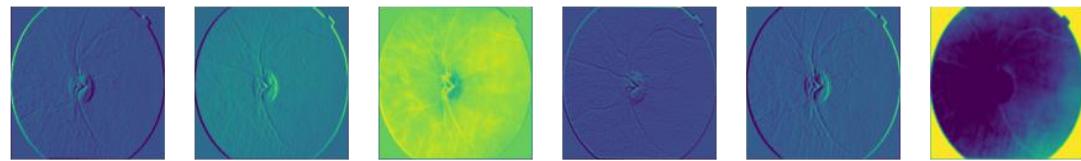
# Activation Maps: Example



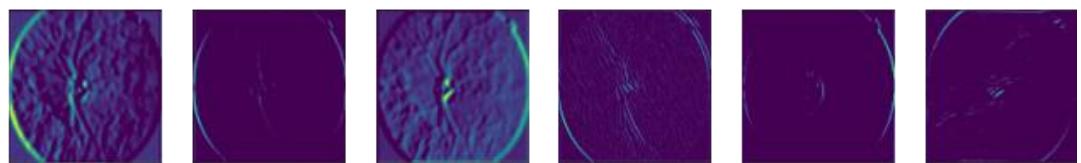
Actual image of  
glaucoma eye disease  
from Drishti-GS dataset



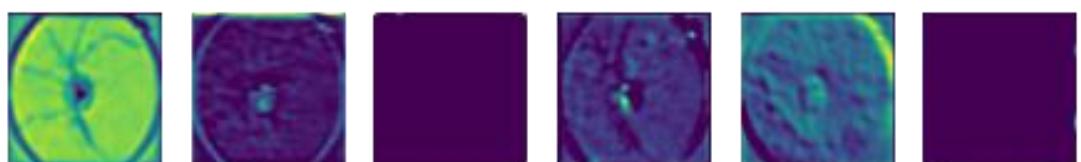
1<sup>st</sup> Convolutional Layer Activation Maps



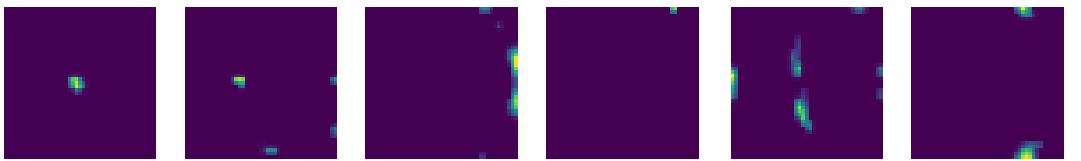
2<sup>nd</sup> Convolutional Layer Activation Maps



4<sup>th</sup> Convolutional Layer Activation Maps



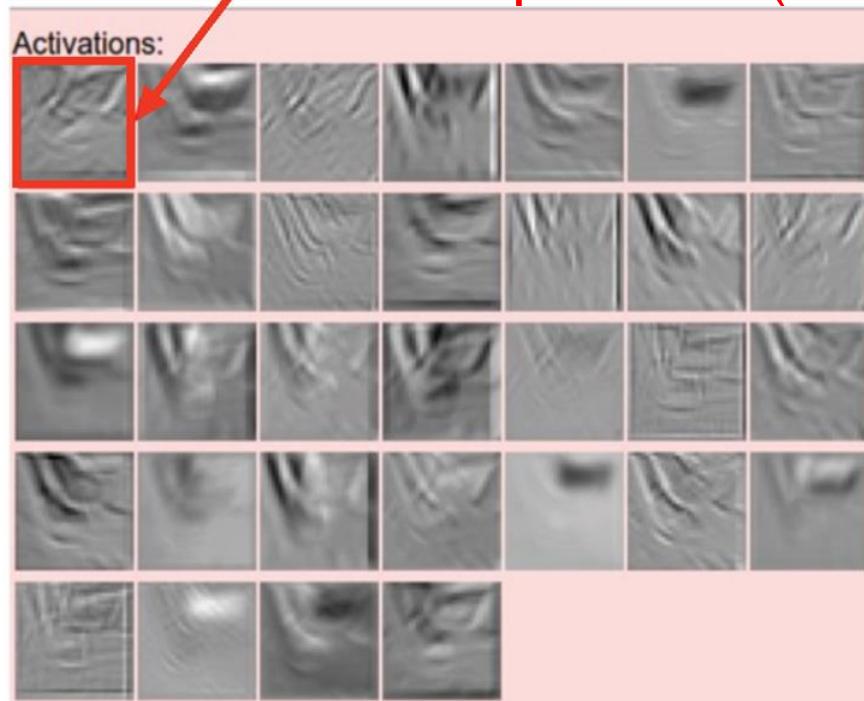
7<sup>th</sup> Convolutional Layer Activation Maps



10<sup>th</sup> Convolutional Layer Activation Maps

# Filters and Redundancies

One filter = one depth slice (or activation map)



- Observe that some filters detect almost same patterns.

# Convolution Layer

- Number of filters in a single Convolution layer are chosen by the user of the CNN (Depends on application and data input type).
- Typically, a sizable number of filters are chosen to extract the maximum number of features.
- So, what do you learn here ?

# Demos

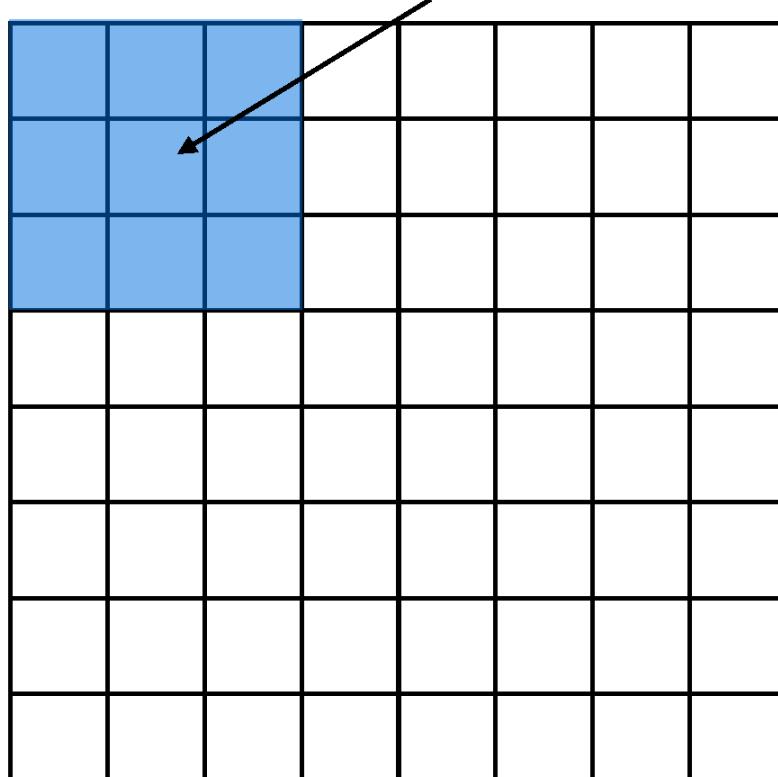
- <http://cs231n.stanford.edu/>
- <http://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

# Convolution layer and Stride

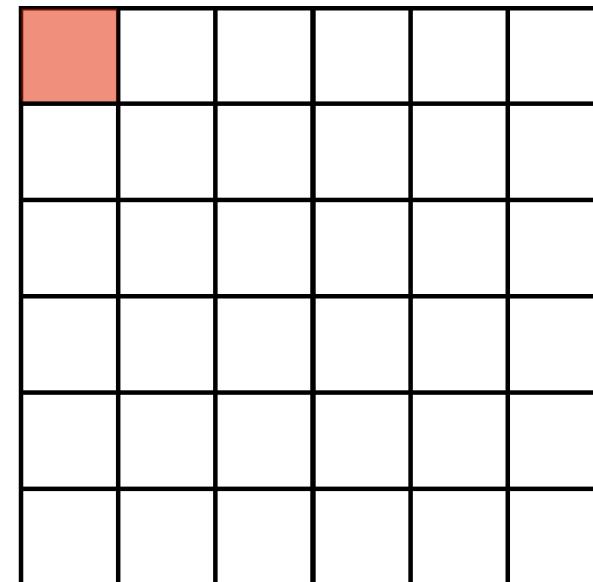
During convolution, the weights “slide” along the input to generate each output

*how many cells filter moves next  
in one step*

**Weights**



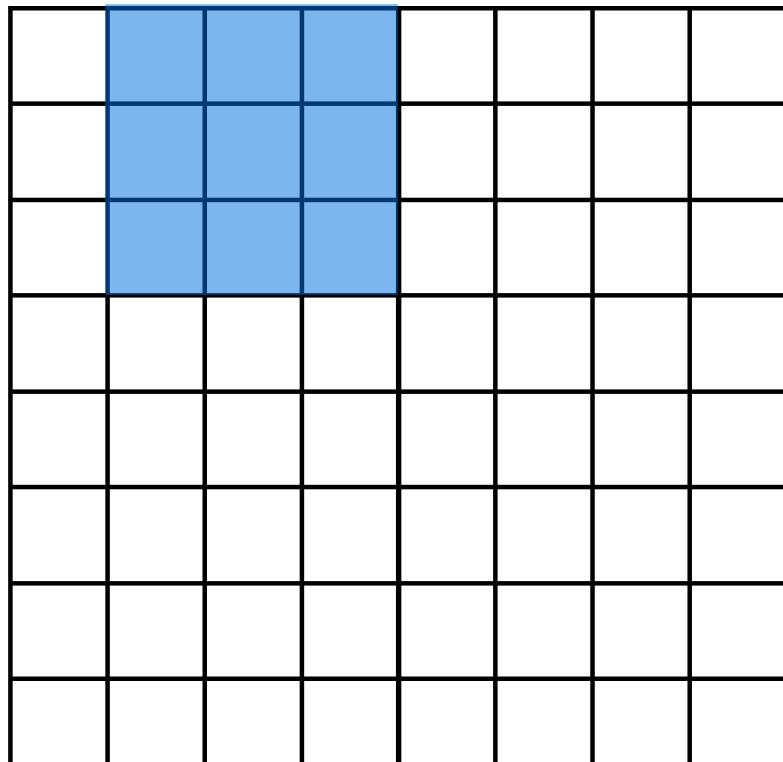
**Input**



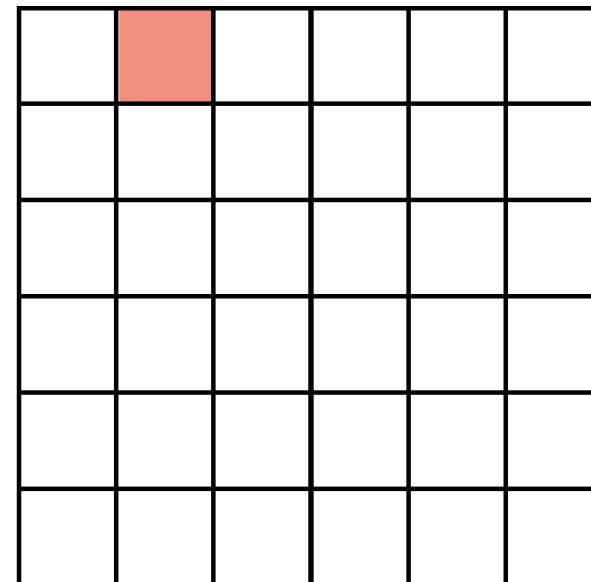
**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



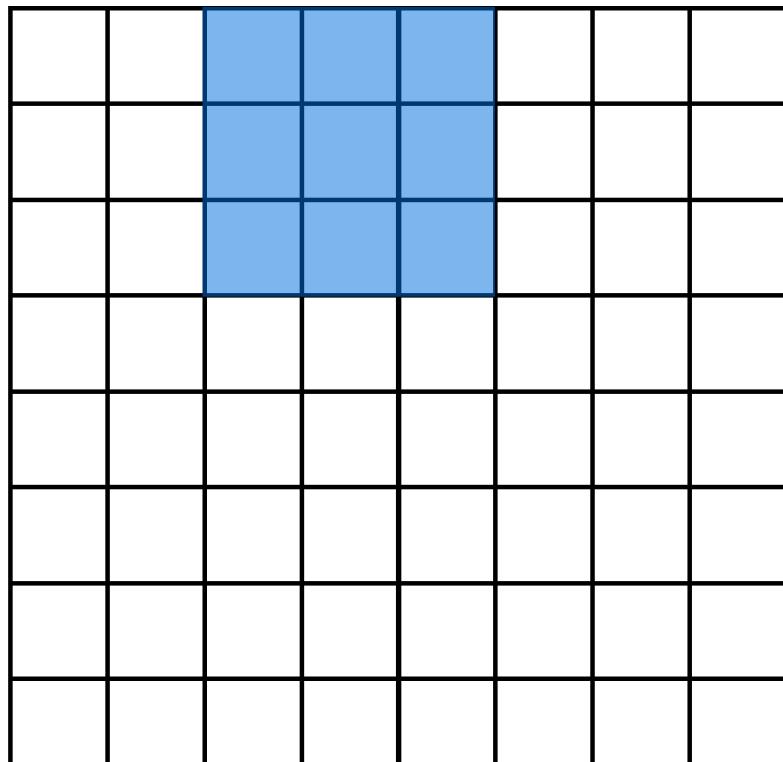
**Input**



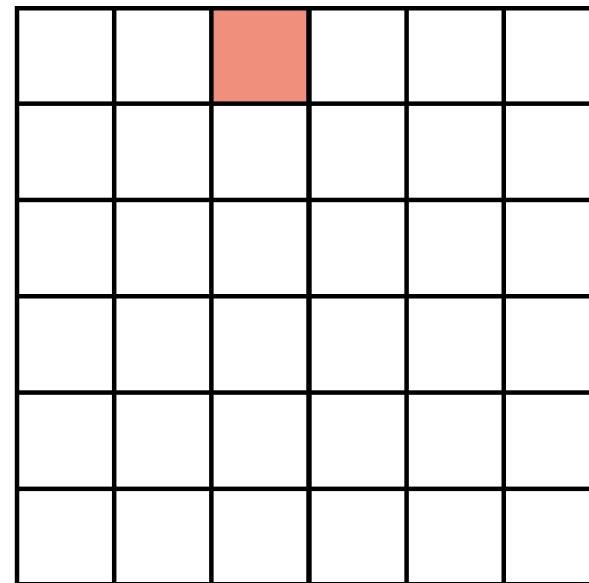
**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



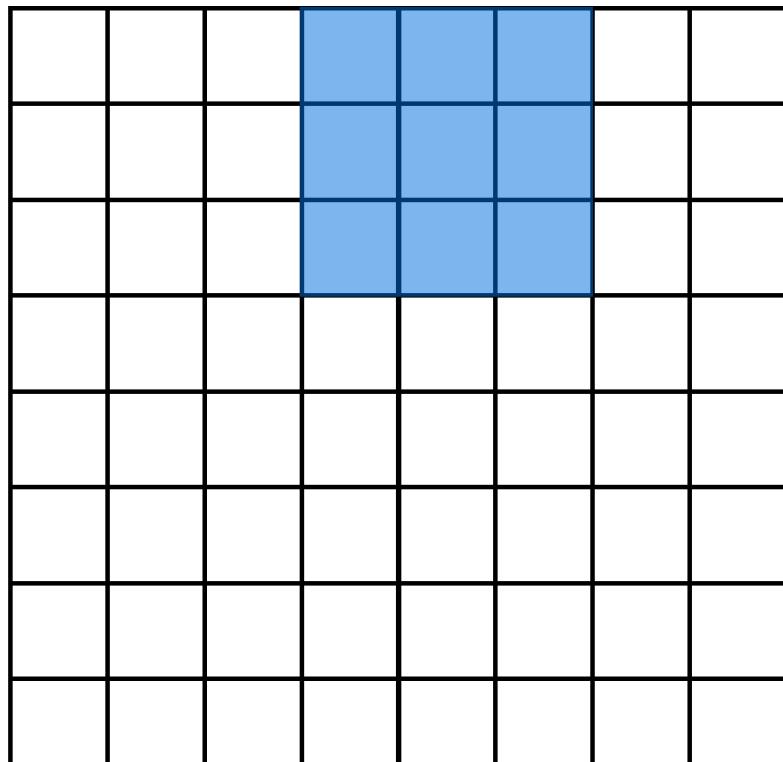
**Input**



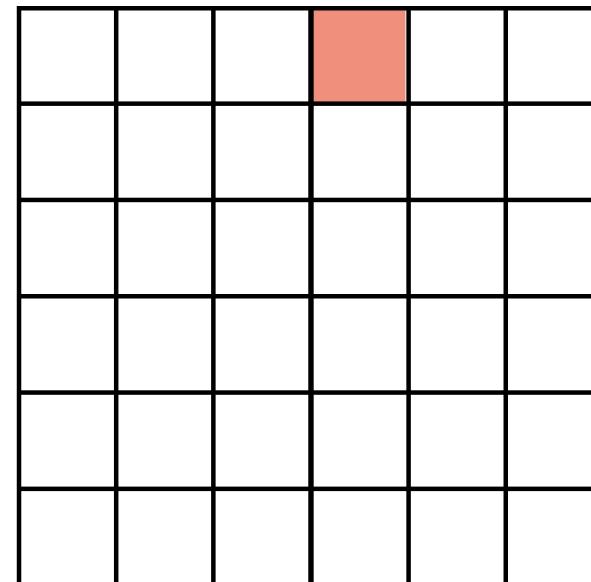
**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



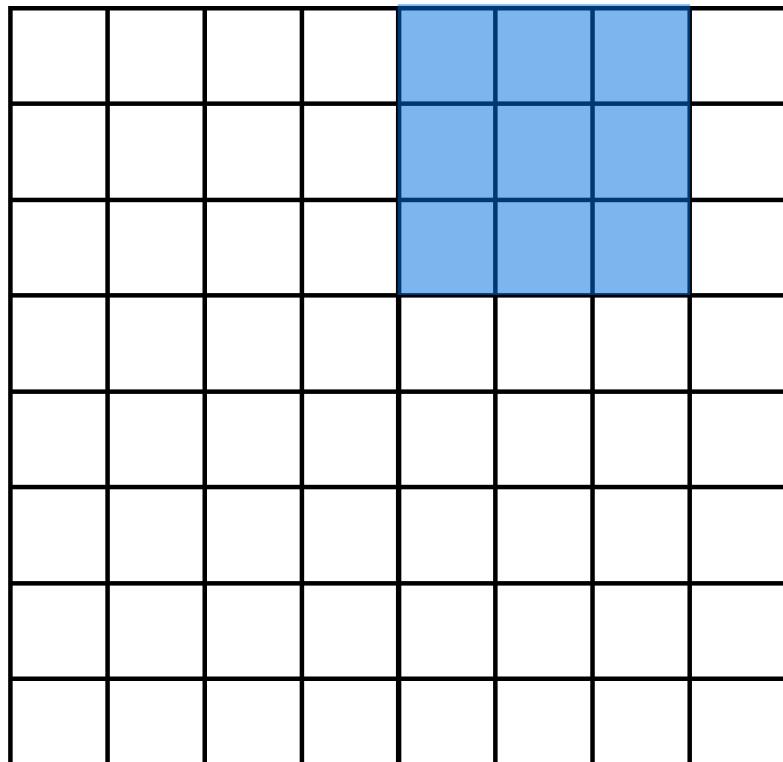
**Input**



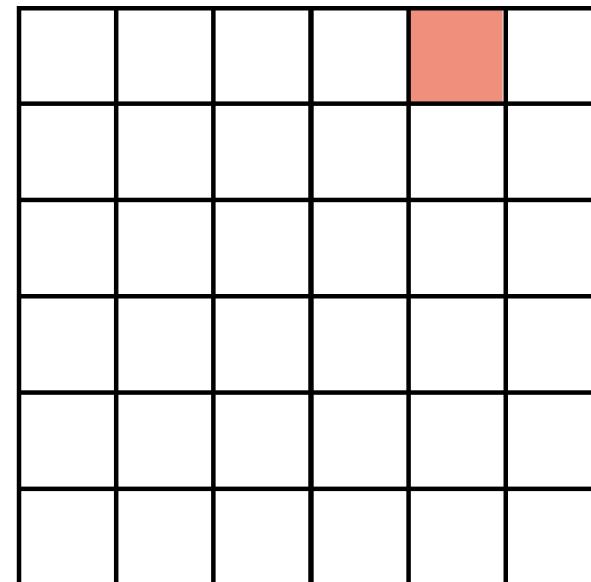
**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



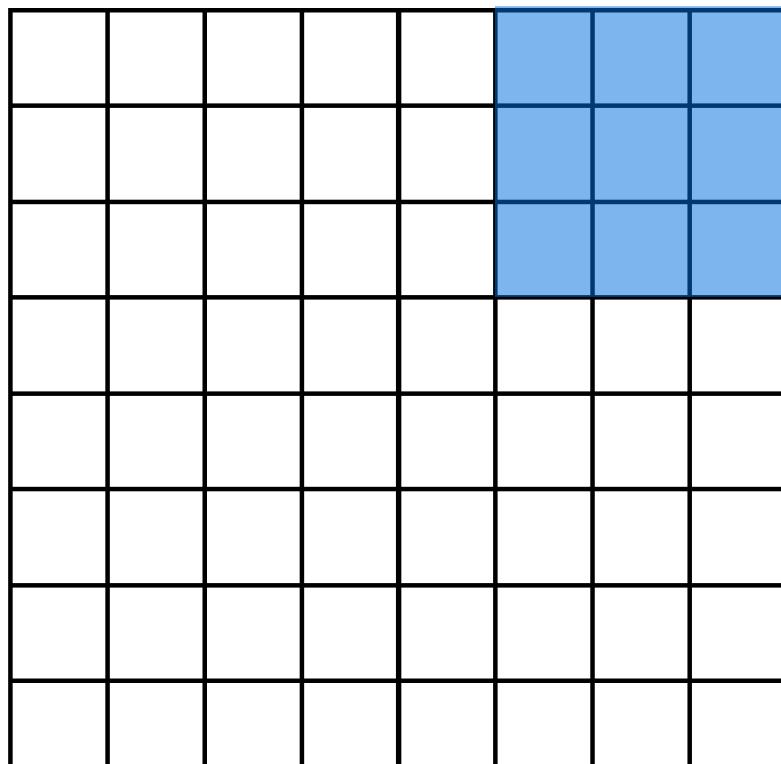
**Input**



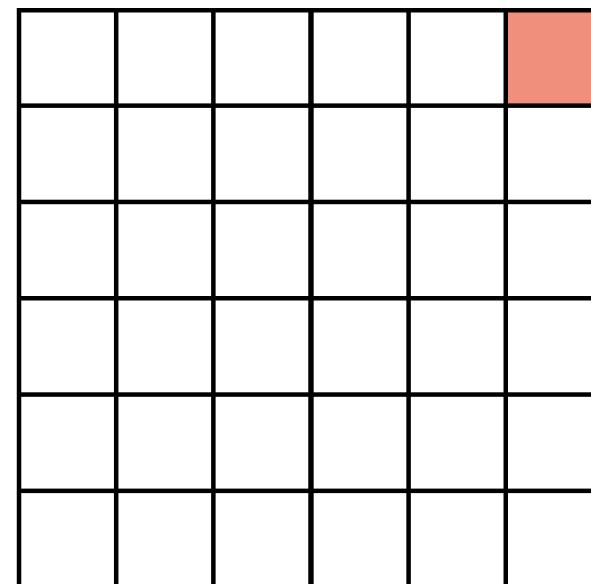
**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



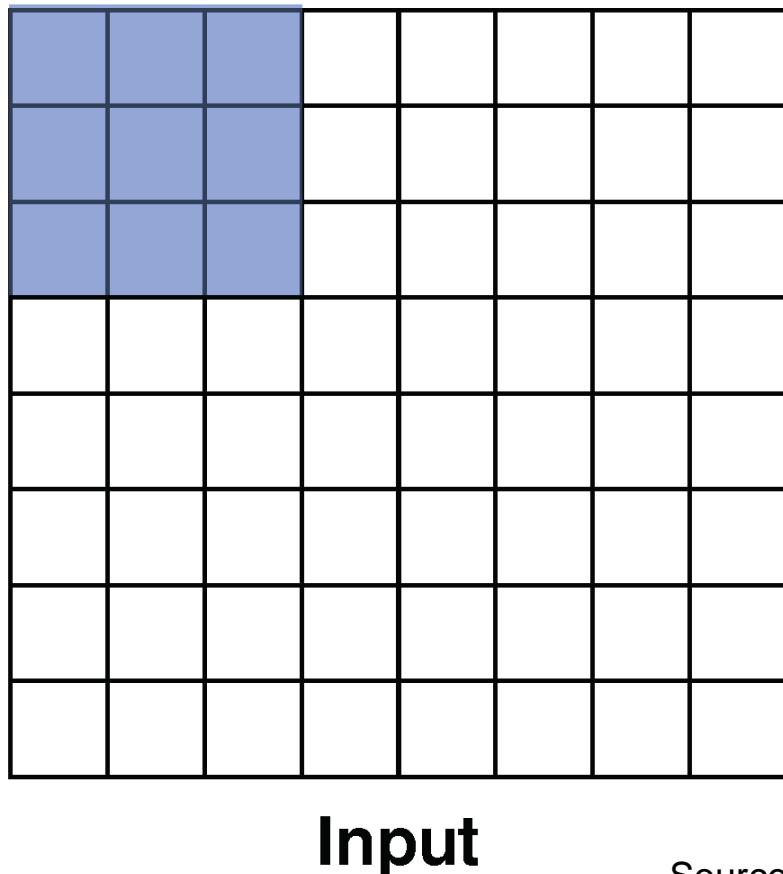
**Input**



**Output**

# Convolution layer and Stride

During convolution, the weights “slide” along the input to generate each output



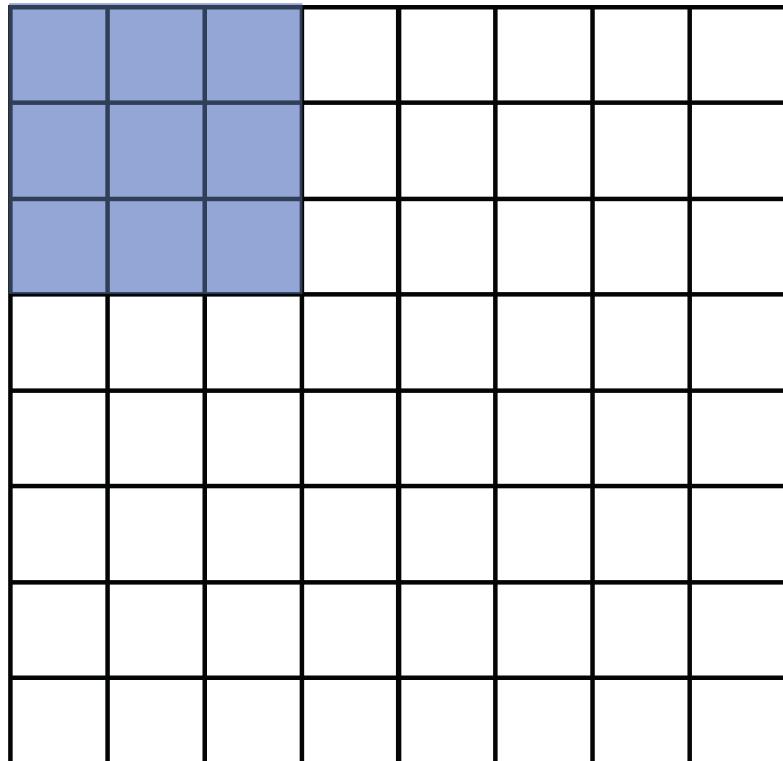
Recall that at each position,  
we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

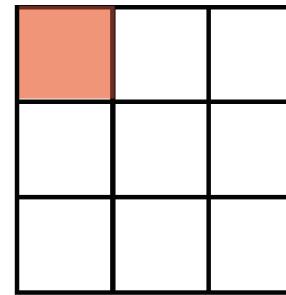
*(channel, row, column)*

# Convolution layer and Stride

But we can also convolve with a **stride**, e.g. stride = 2



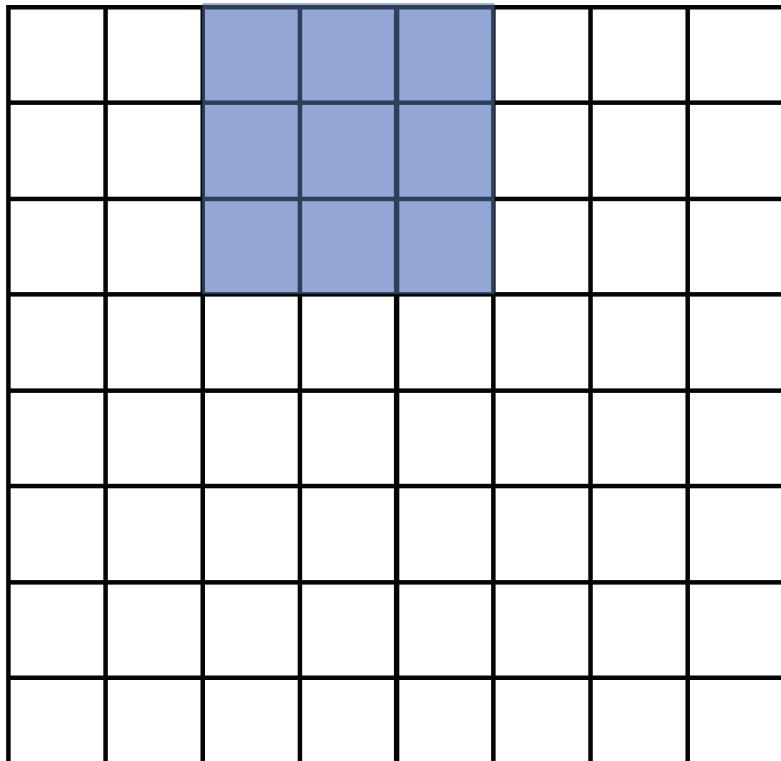
**Input**



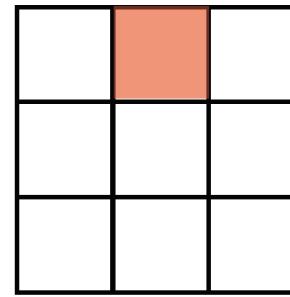
**Output**

# Convolution layer and Stride

But we can also convolve with a **stride**, e.g. stride = 2



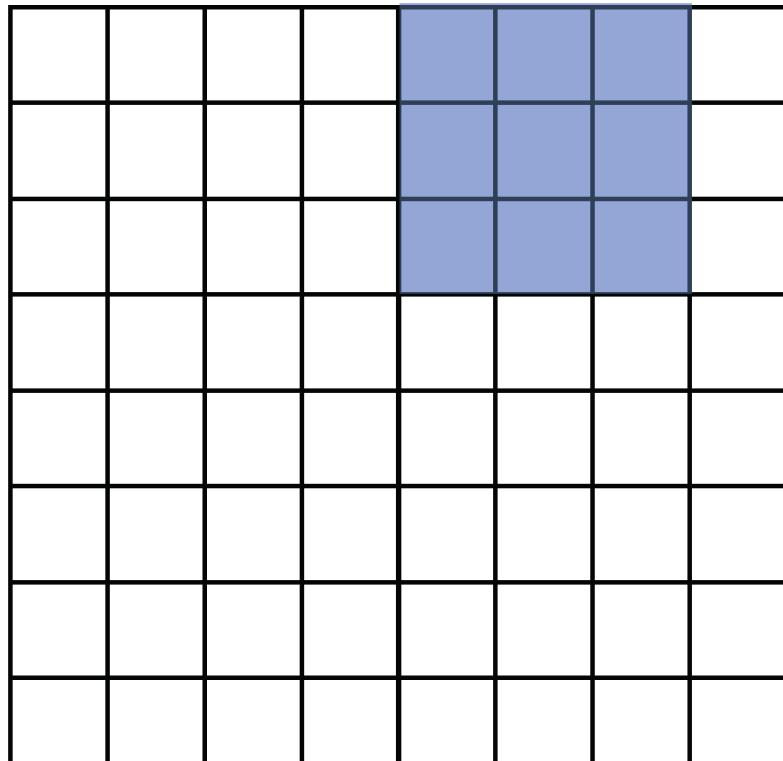
**Input**



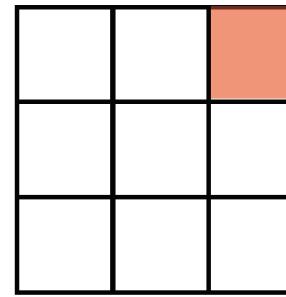
**Output**

# Convolution layer and Stride

But we can also convolve with a **stride**, e.g. stride = 2



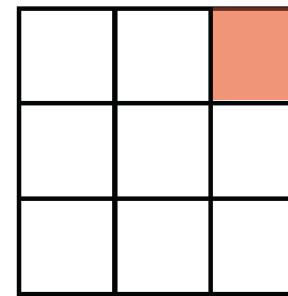
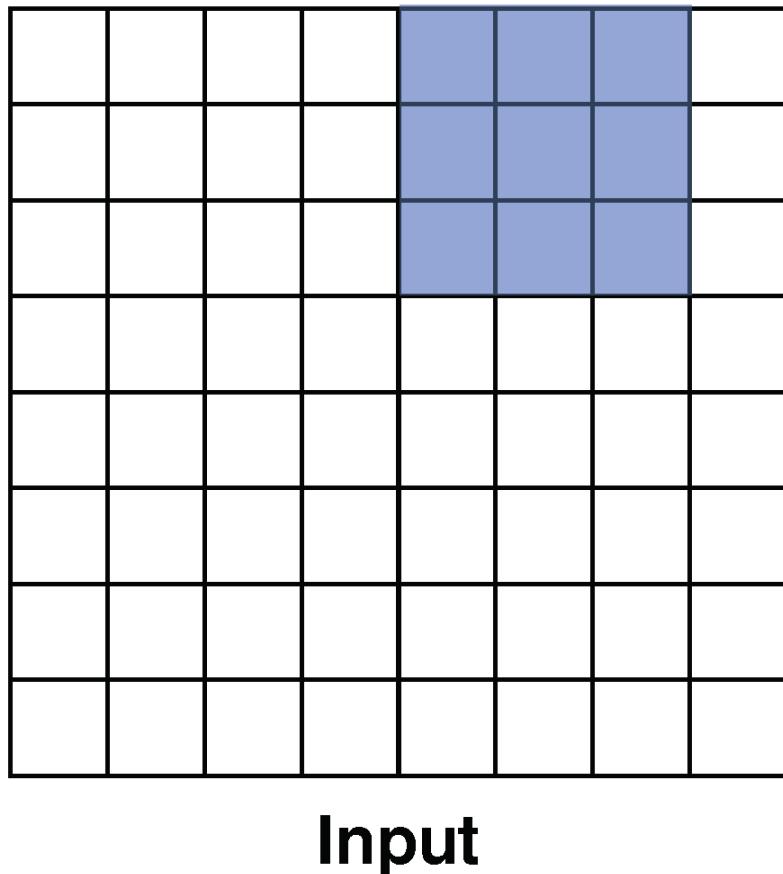
**Input**



**Output**

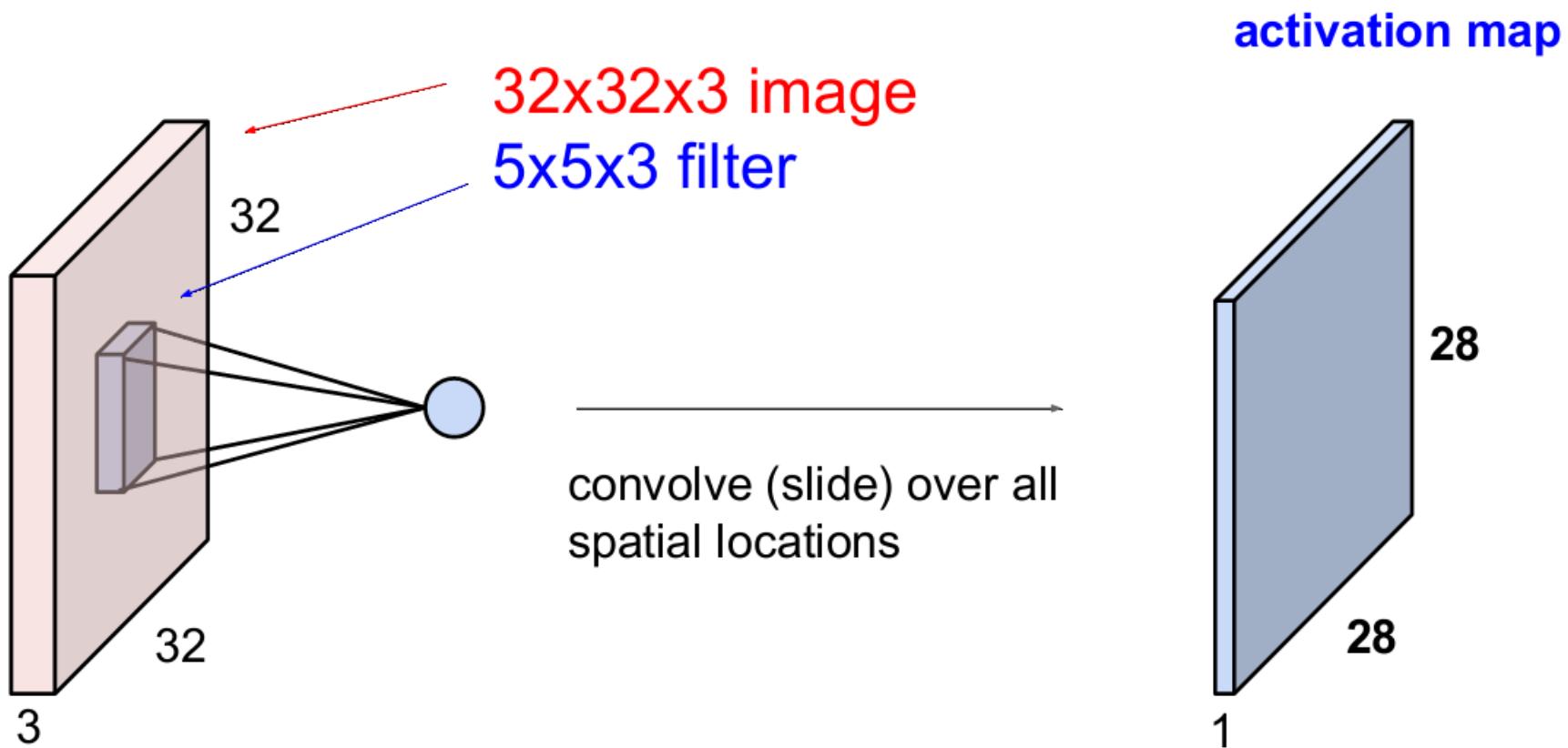
# Convolution layer and Stride

But we can also convolve with a **stride**, e.g. stride = 2



- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

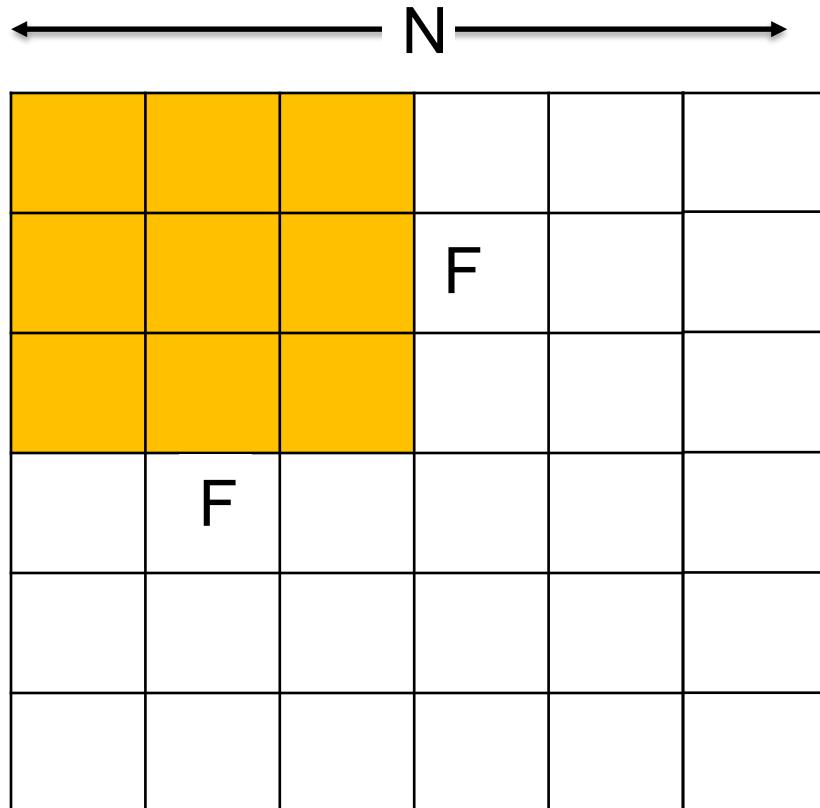
# Convolution layers and image dimension



# Output Volume

- Operation with filters produces an output volume.
- Filter values are initially chosen randomly.
- Gradually, values of filters are updated in each iteration. Each filter learns to detect a small pattern.
- Suppose the image is  $256 \times 256 \times 3$  and there are 8 filters, each of size  $5 \times 5$ , what will be the output volume ?
- Each filter  $5 \times 5$  will be converted to a  $5 \times 5 \times 3$  filter volume.
- Then it will convolve with each block of  $5 \times 5 \times 3$  in the image to give just one output value.
- Whole convolution process with one filter will give a layer of size  $252 \times 252$ .
- For 8 filters, it will be a  $252 \times 252 \times 8$  volume.

# Convolution Layer and Output Dimension



Output size:

$$\frac{N - F}{stride} + 1$$

$Stride = 1$

$N = 6$  and  $F = 3$ ,  $Stride S = 1$ .  
What will be the output volume?

# Convolution Layers and Output Dimensions

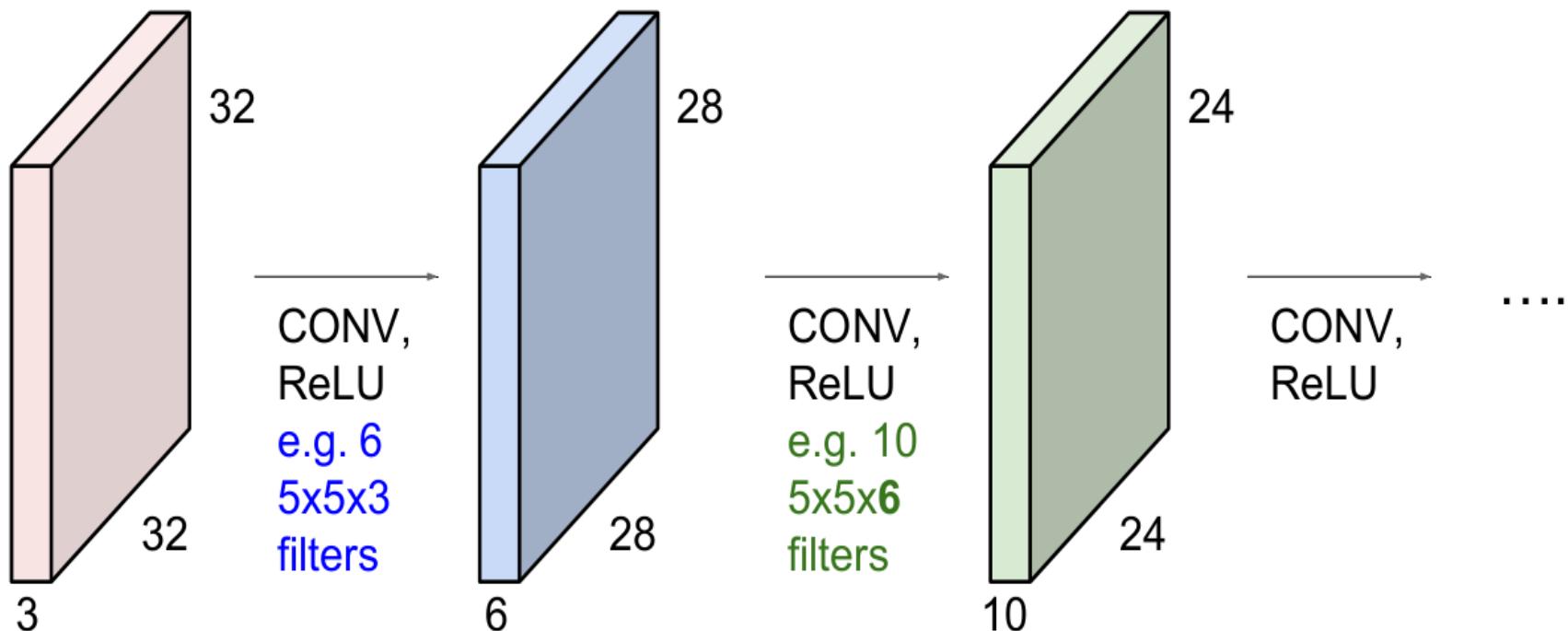


Figure source: Fei Fei Li's Lecture slides, Stanford University

# Zero Padding

- In practice: Common to zero pad the border
- Two advantages
  - Dimension reduction is controlled.
  - Corner pixels have better contribution.

0	0	0	0	0	0	0	0
0							0
0							0
0							0
0							0
0							0
0							0
0	0	0	0	0	0	0	0

# Zero Padding

- Example. Input image of size  $7 \times 7$
- $3 \times 3$  filter, applied with stride 1
- Pad with 1 pixel border. what will be the output?
  - $7 \times 7$
- CONV layers mostly use stride  $S = 1$ .
- So, what should be the size of zero padding to maintain the image's height and width ?
- $(F-1)/2$
- Example
  - $F = 3 \Rightarrow$  zero pad with 1
  - $F = 5 \Rightarrow$  zero pad with 2

# Input and Output Volumes

- Suppose the input volume is  $64 \times 64 \times 3$
- Apply 10 filters each of size  $5 \times 5$
- Padding  $P=2$ , stride  $S=1$
- What will be the output volume?
- Number of parameters in a layer?
- How many parameters to learn?

(Hint: After convolution apply bias)

# Input and Output Volumes

- Now device a formula if the input volume is of the size :  $WXHX3$
- Width :  $\left\lfloor \frac{W-F+2P}{S} \right\rfloor + 1$
- Height :  $\left\lfloor \frac{H-F+2P}{S} \right\rfloor + 1$
- Depth = number of filters.

# Convolution vs Fully Connected Layer

- Fully connected layer

- The input image size is  $32 \times 32 \times 3$
- Flatten it first (convert to 1-d vector)
- Dimension of input vector:  $32 \times 32 \times 3 = 3072 \times 1$
- If there are 10 units (neuron) in the first hidden layer of a FC network, then each unit looks at the full input vector

$$a = f(Z) = f(W^T X + b).$$

Trainable parameters

Weights:  $3072 \times 10$

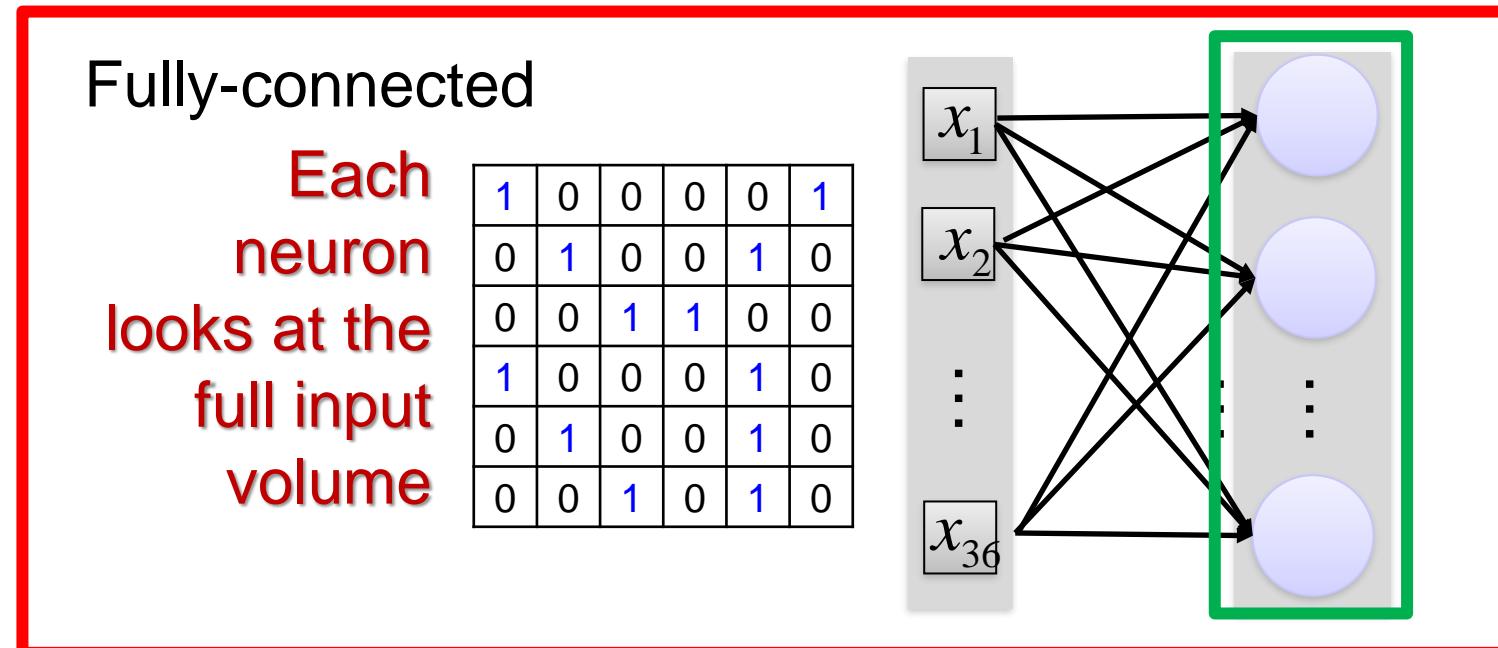
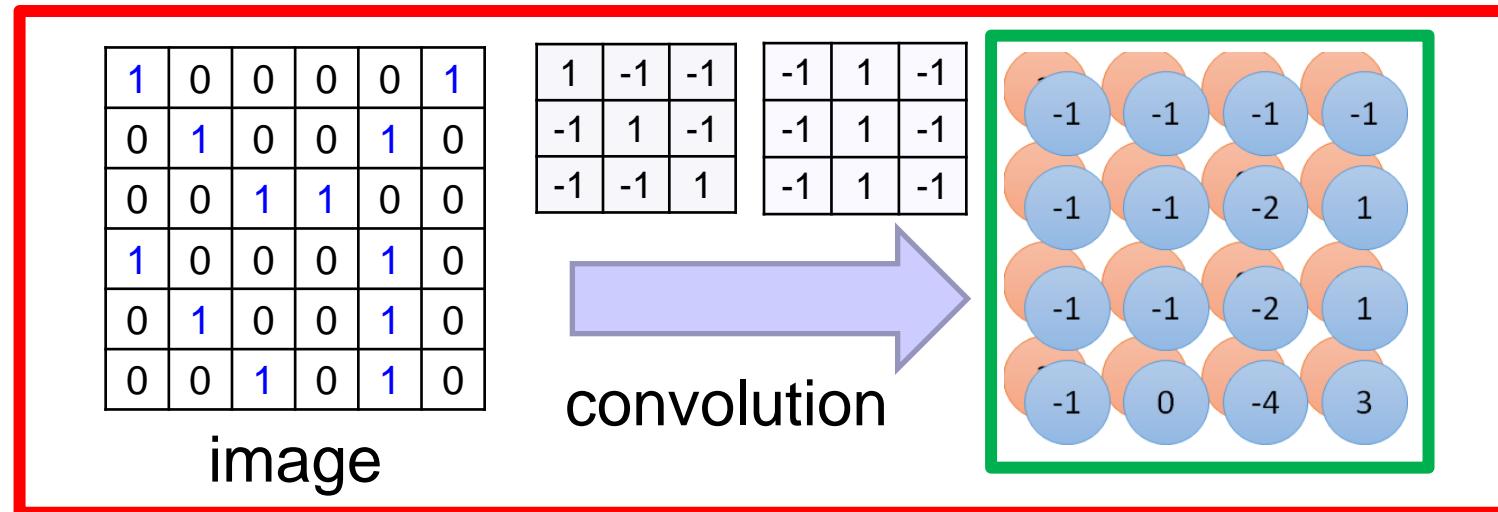
Bias: 10

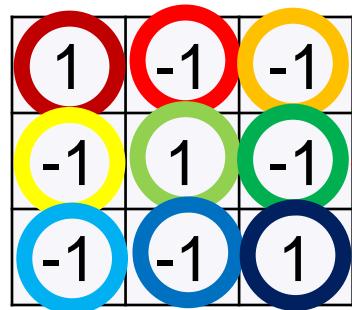
Total trainable parameters :  
30730

# Convolution vs Fully Connected Layer

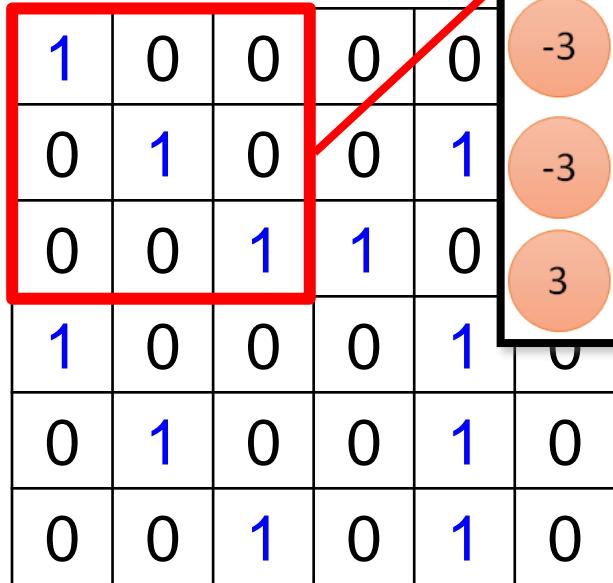
- Convolution layer
  - The input image size is  $32 \times 32 \times 3$
  - If there are 10 filters, each of size  $3 \times 3 \times 3$  in the first Conv layer of the network.
  - Total trainable parameters ?
  - 280

# Convolution vs Fully Connected layer



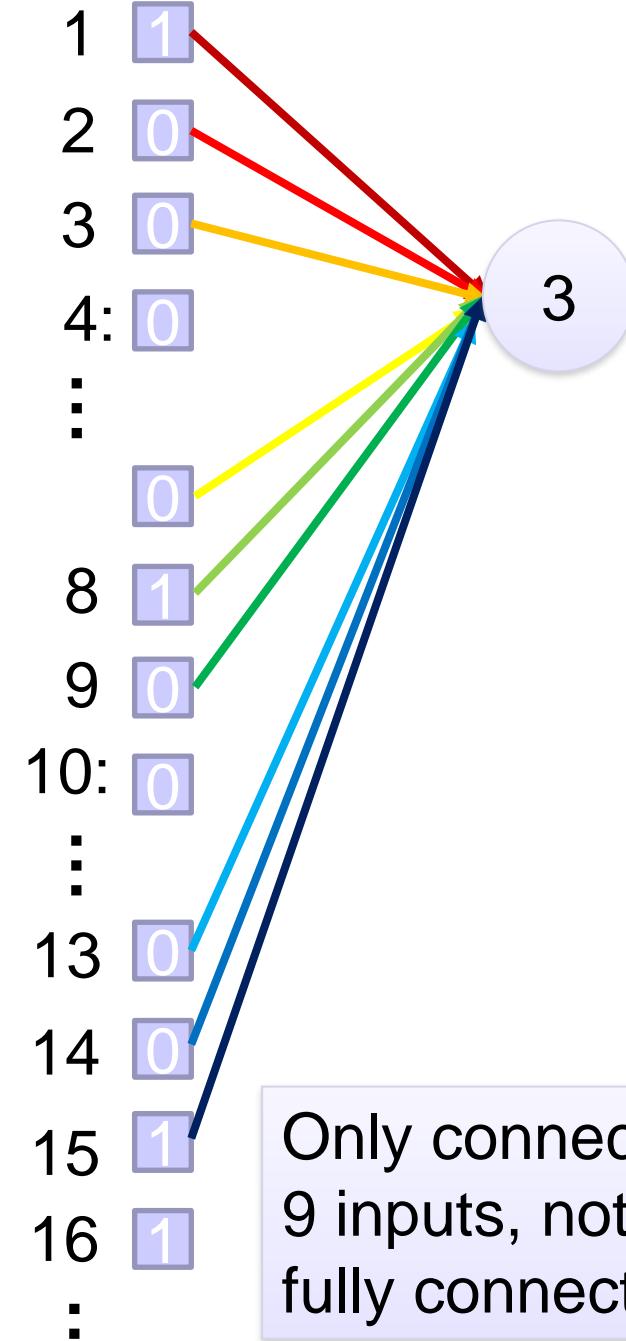
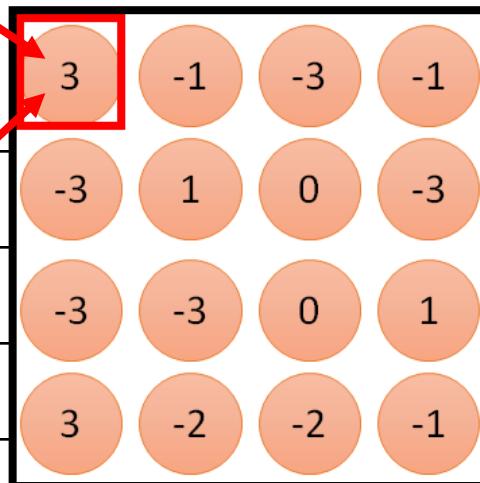


Filter 1

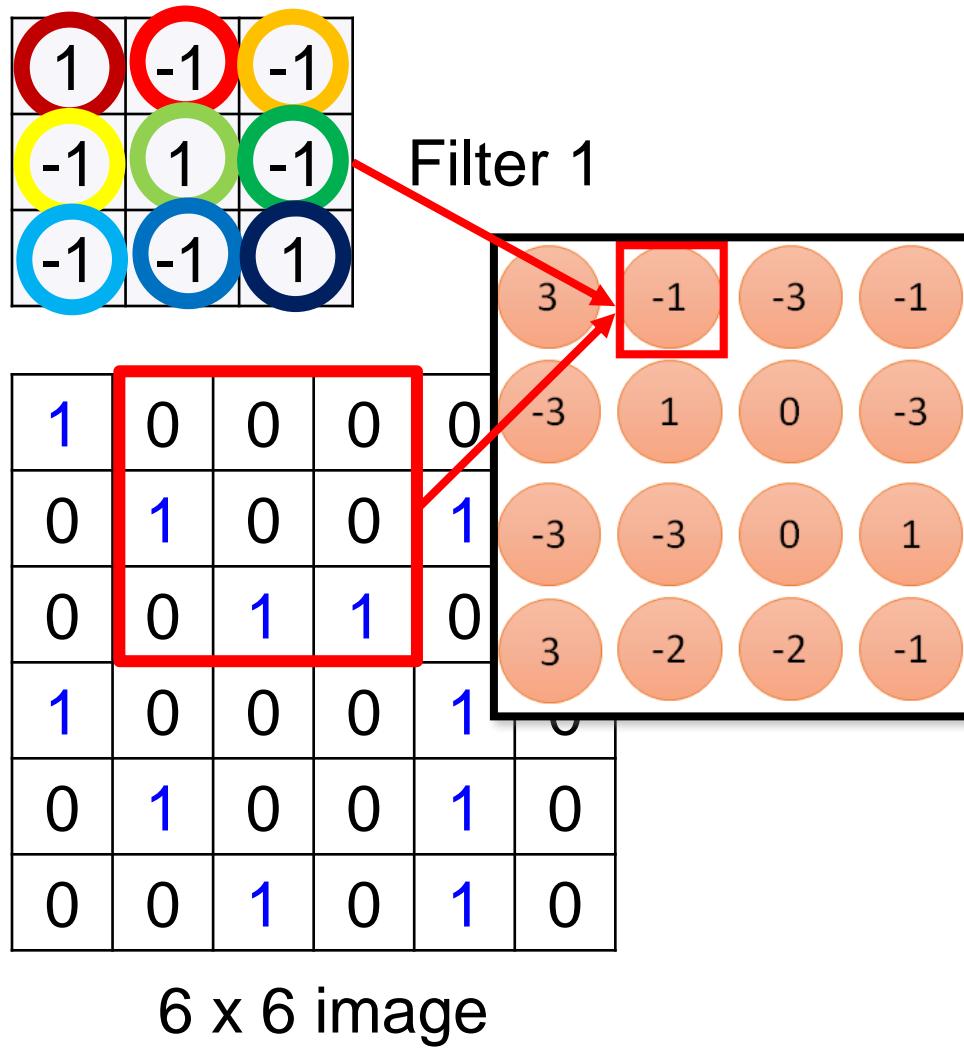


$6 \times 6$  image

fewer parameters!

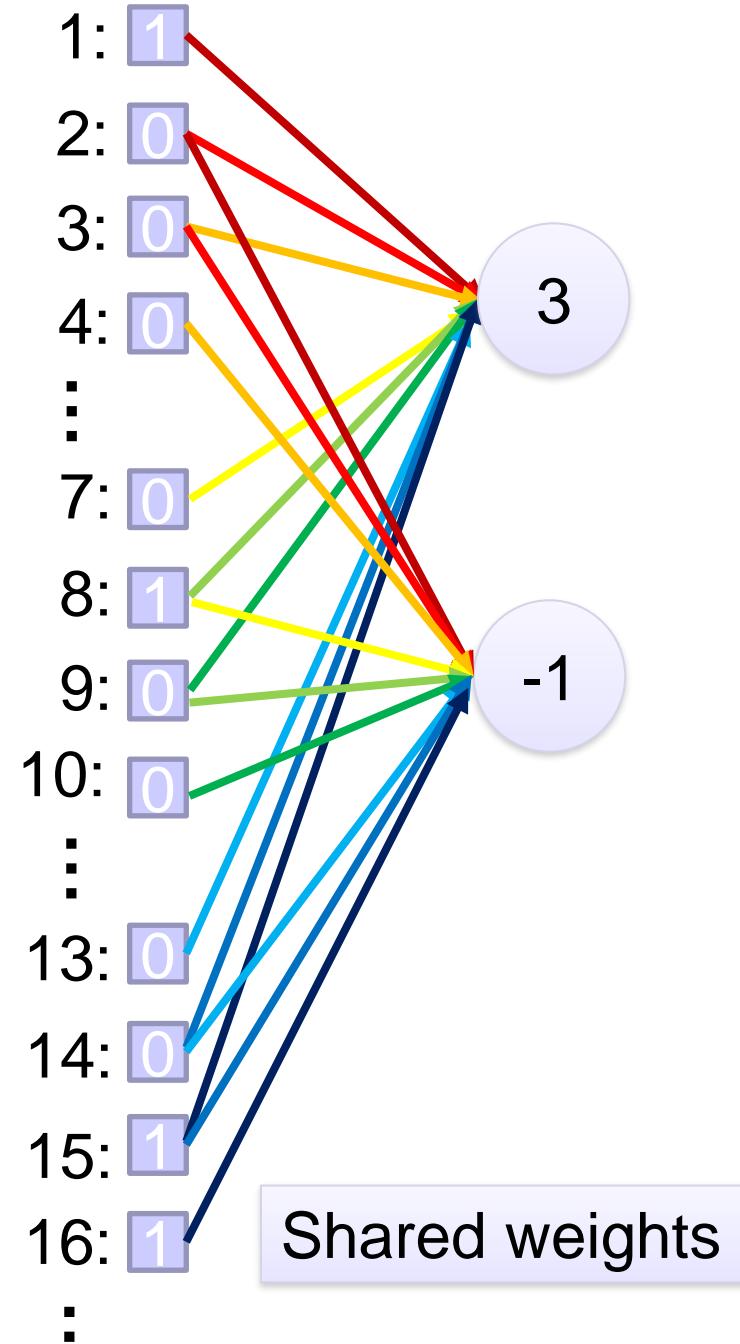


Only connect to 9 inputs, not fully connected

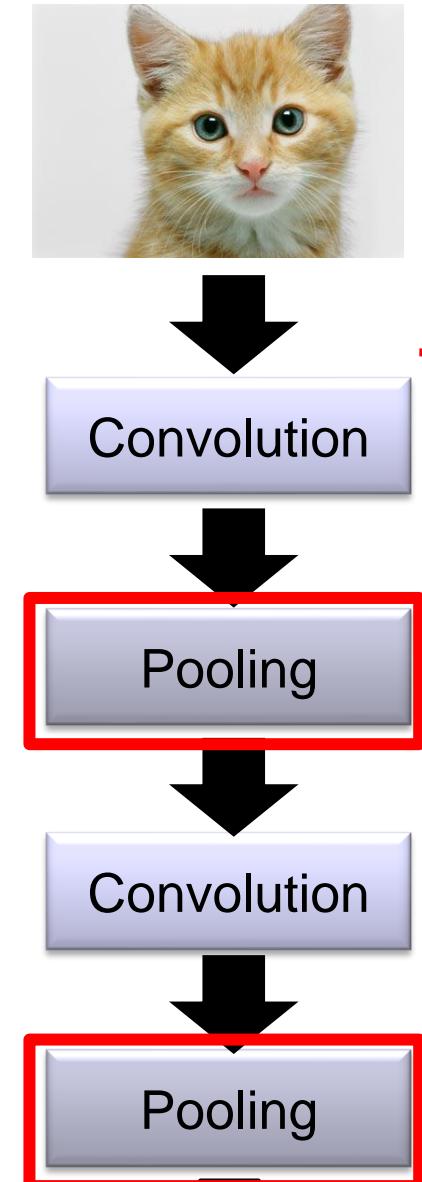
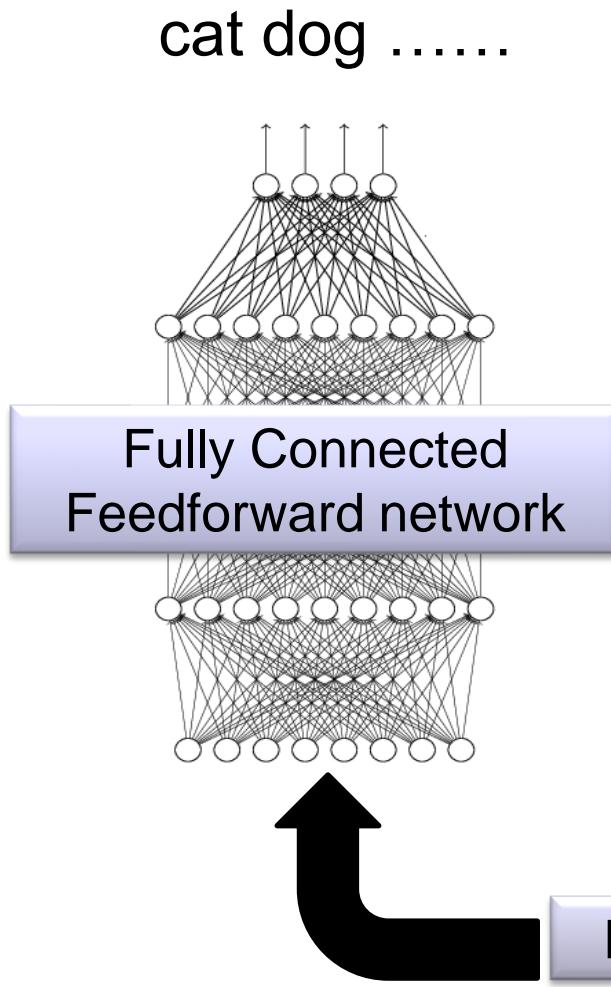


Fewer parameters

Even fewer parameters



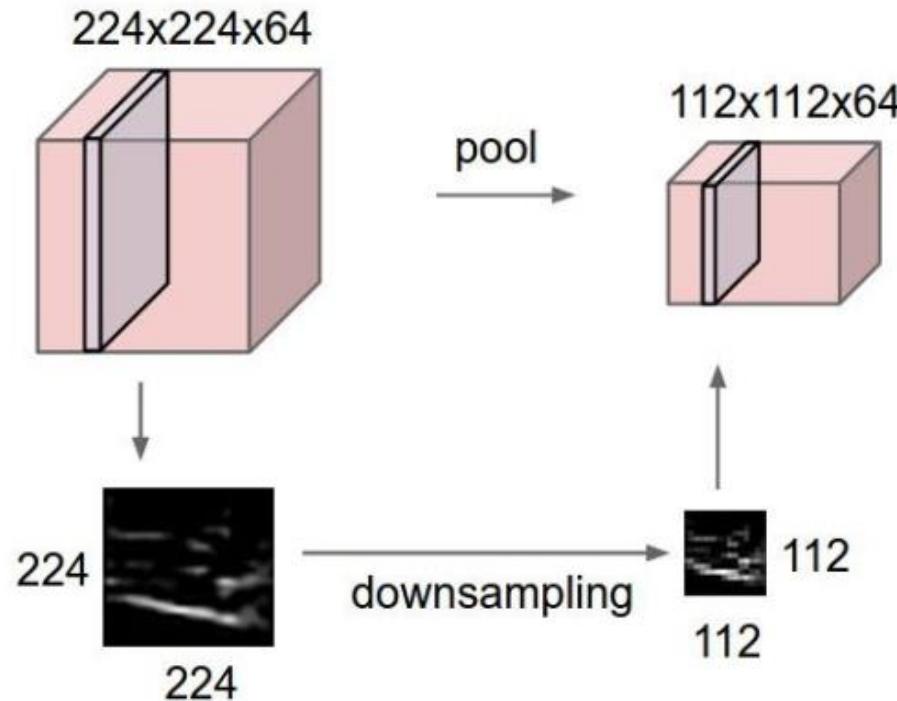
# The whole CNN Architecture



# Pooling layer

- When the image is large, one needs to compress the same for managing the input and output volumes between different convolution layers.
- Pooling layers make the representations smaller.
- Output volumes are more manageable.
- Operates over each activation map independently.

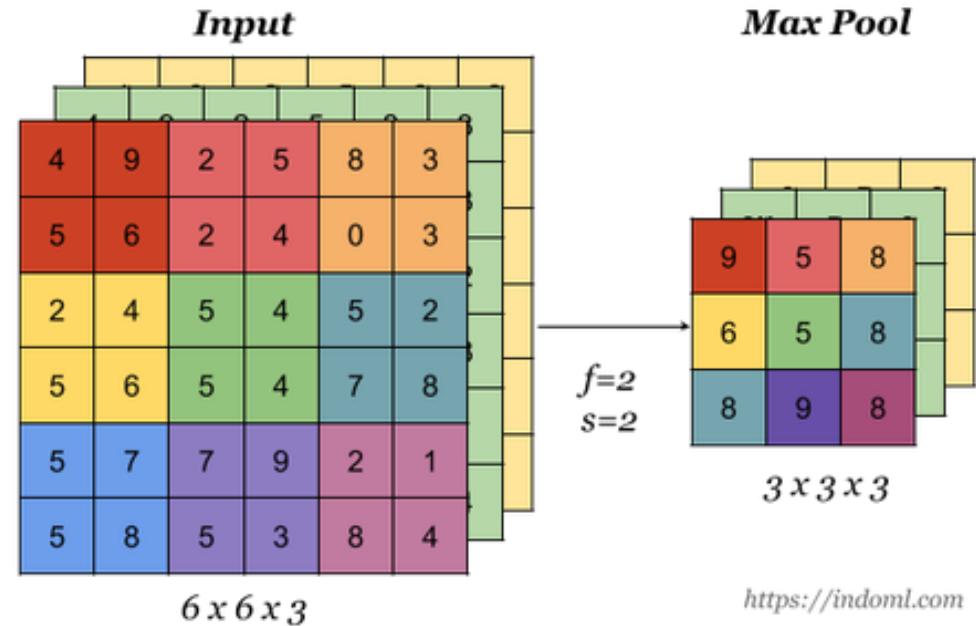
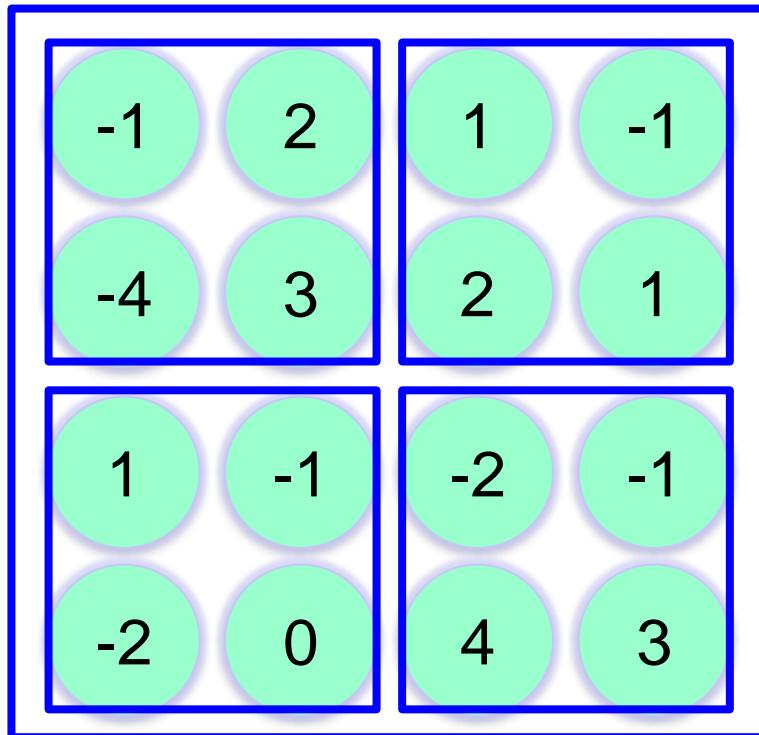
# Pooling layer



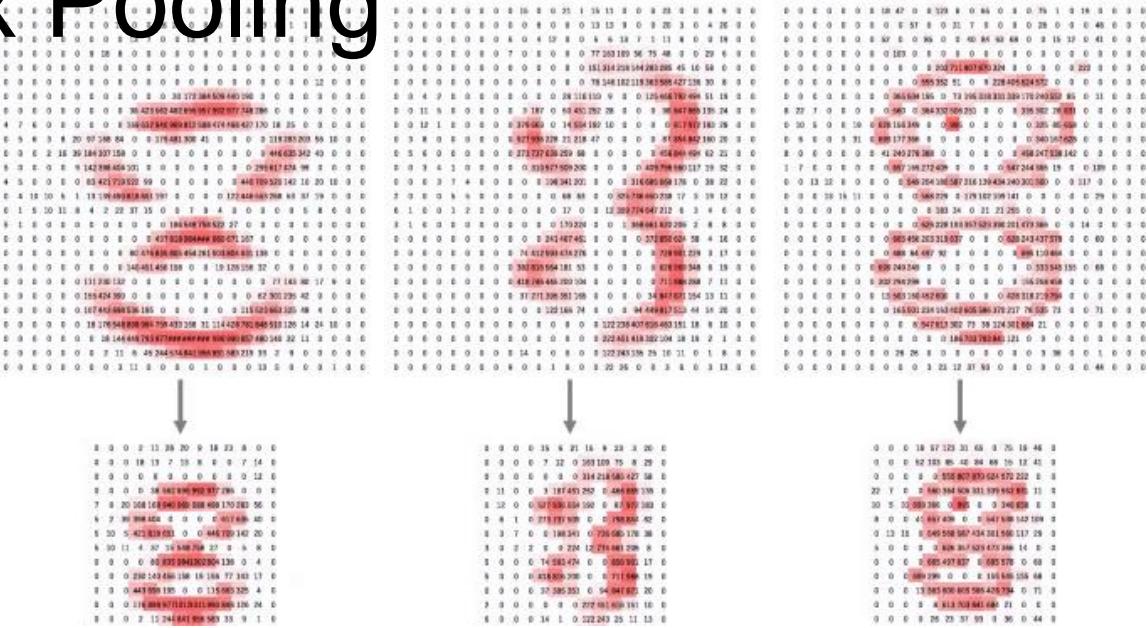
- Max pooling
- Average pooling
- Others: Min pooling.

# Max Pooling

- With  $2 \times 2$  filters and stride 2

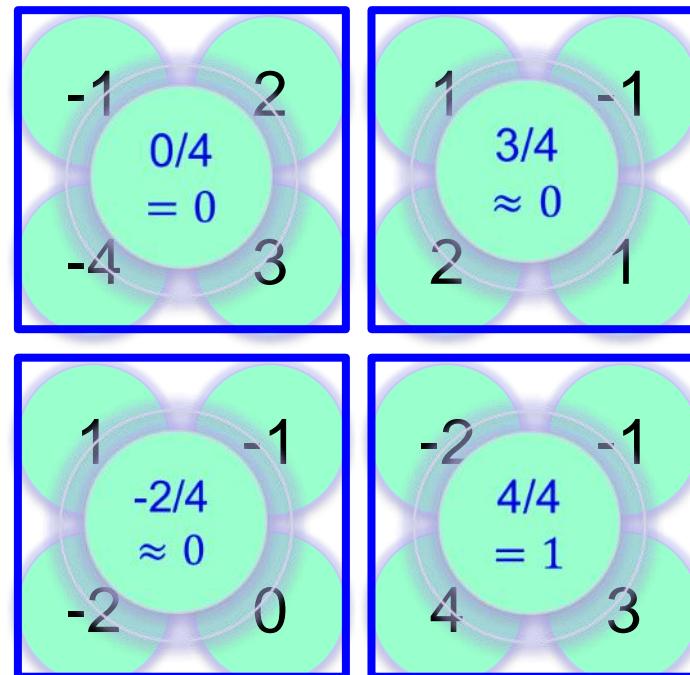


# Max Pooling



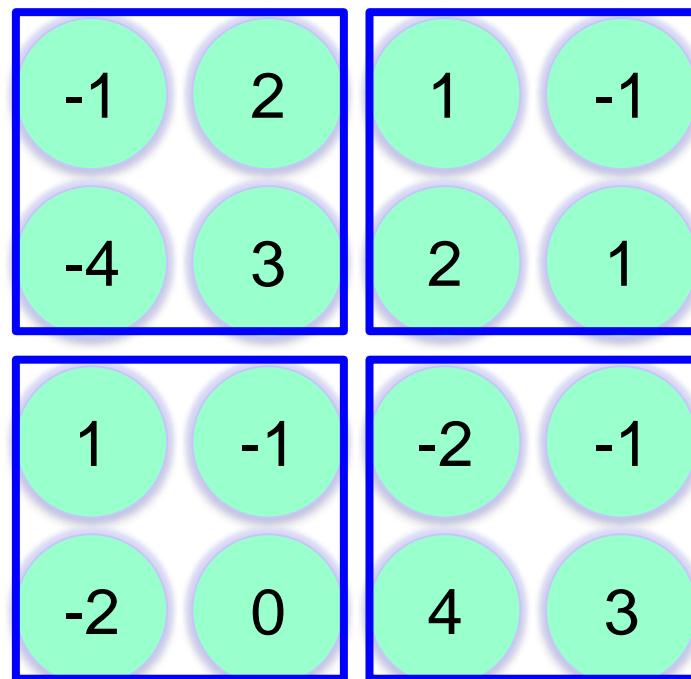
# Average Pooling

- With  $2 \times 2$  filters and stride 2



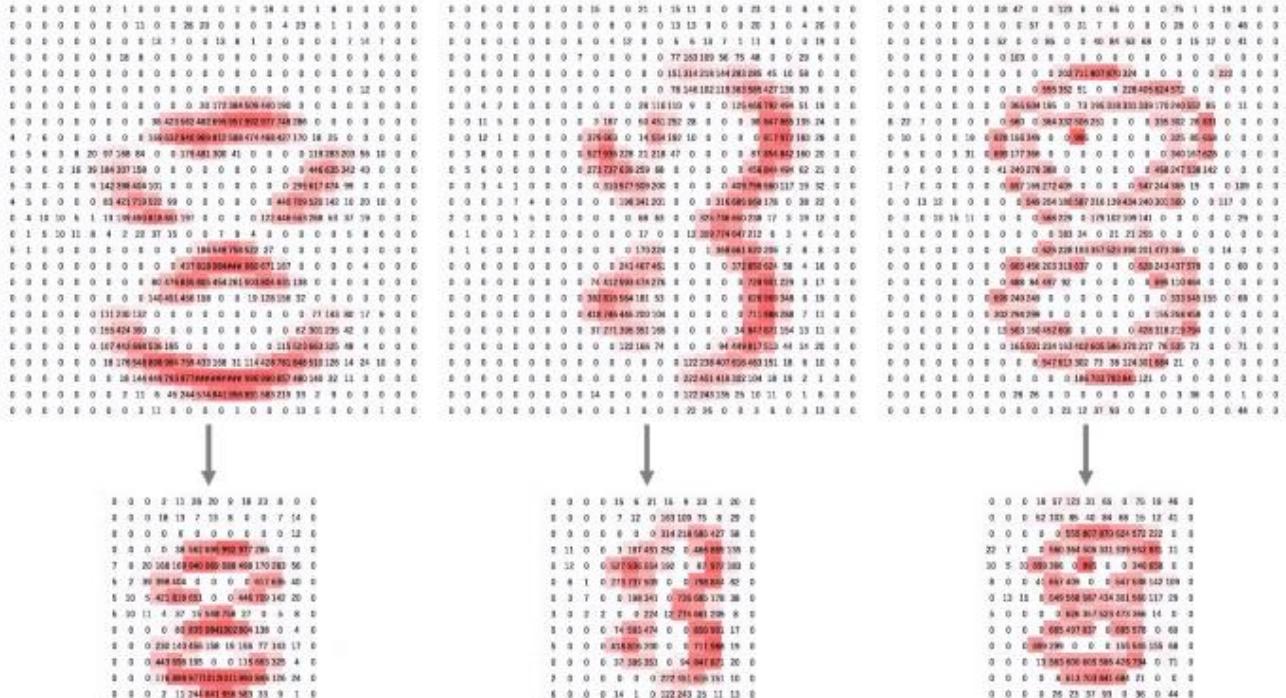
# Min Pooling

- With  $2 \times 2$  filters and stride 2



# Max Pooling – Another Example

- With  $2 \times 2$  filters and stride 2
- Note that images in the upper row are downsampled by a factor of 2, from  $26 \times 26$  to  $13 \times 13$



# Pooling Advantage

- Pooling is a kind of subsampling.
- Subsampling the image does not lose image information.

A bowl of food on a plate

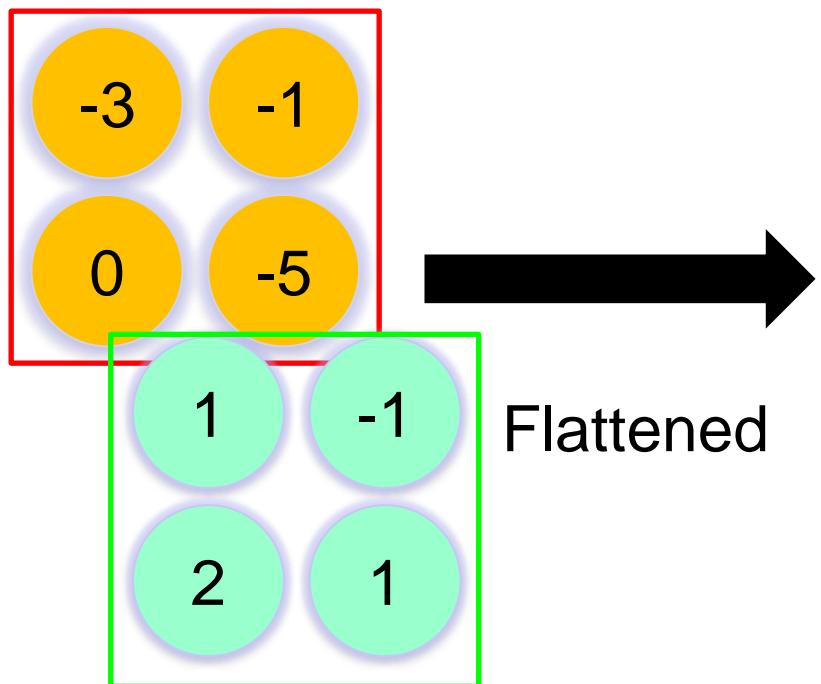


A bowl of food on a plate

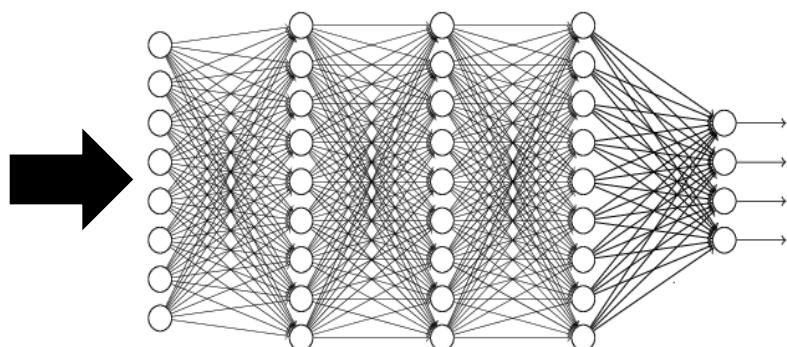


Fewer parameters needed to extract features for pattern recognition and characterizing the image.

# Flattening

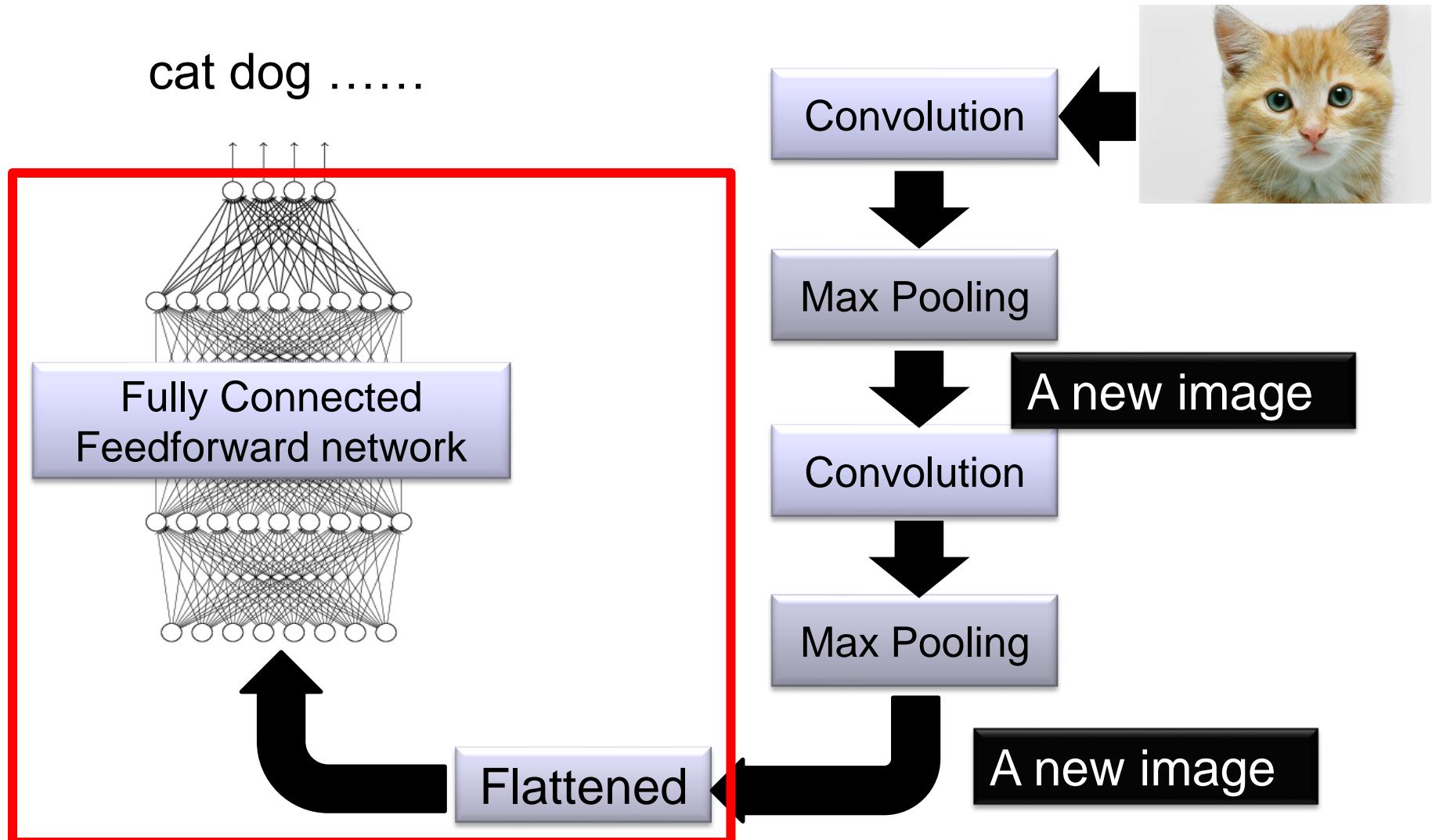


what's left:  
Fully Connected  
Layers

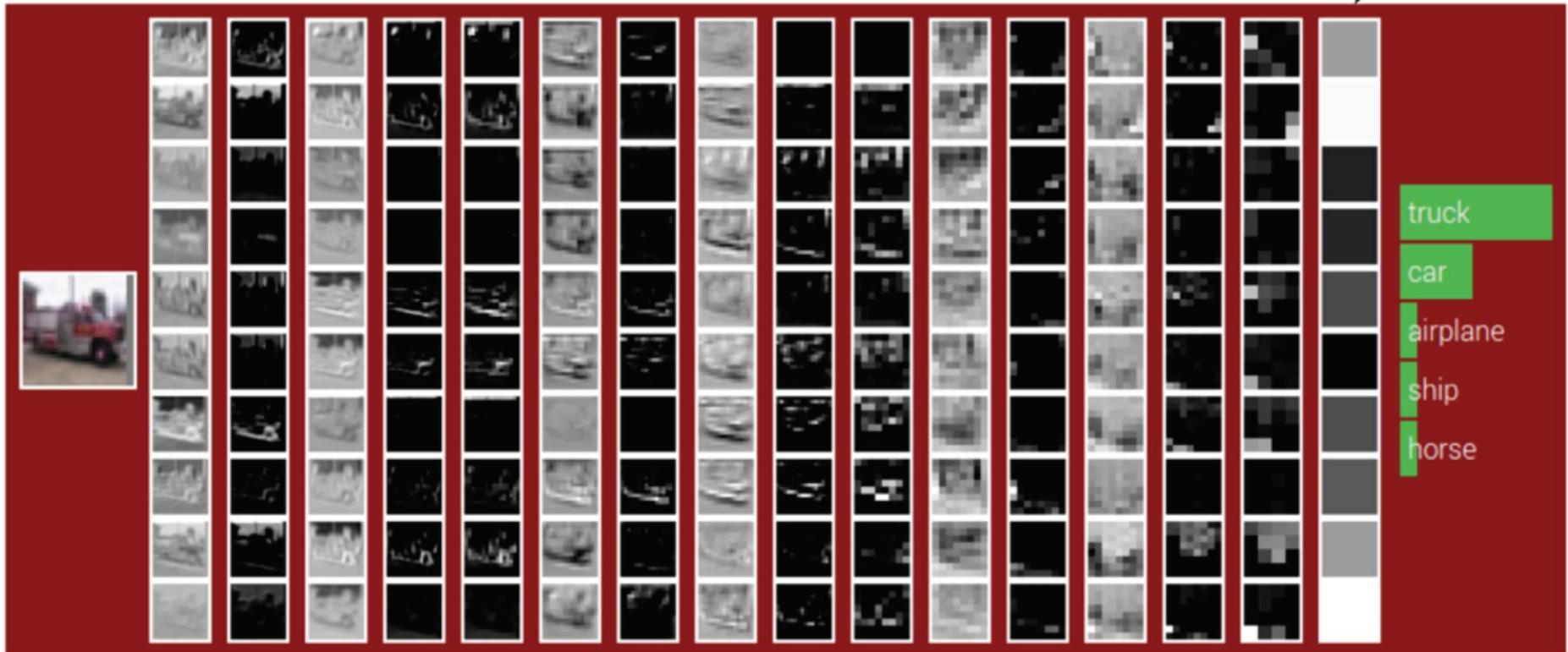
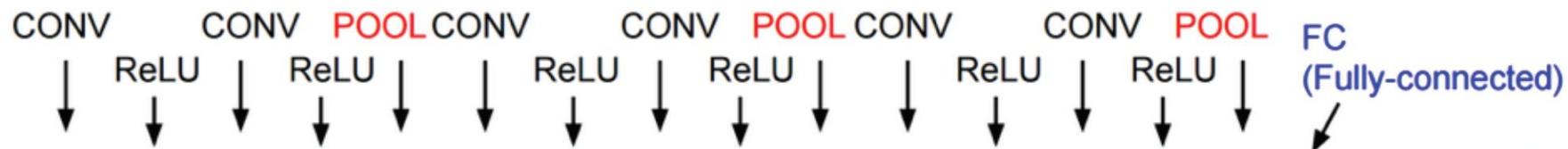


Fully Connected  
Feedforward network

# A Generic CNN Architecture



# Another Example ConvNet



10x3x3 conv filters, stride 1, pad 1

2x2 pool filters, stride 2

Figure: Andrej Karpathy

# CNN Advantage in Dealing High Dimension Data

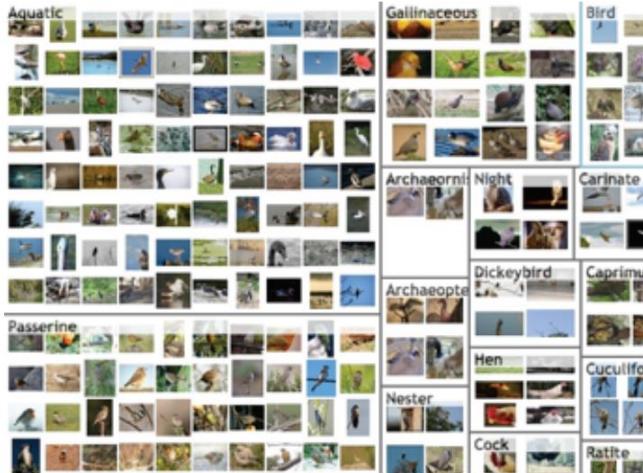
- Number of connections are reduced (Sparse interactions).
- Parameter sharing.
- Equivariant representations.
- Ability to work with inputs of variable size.
- Pooling layers compress the output volumes.
- As a result, number of parameters to learn are much less.
- The number of training parameters depends on how many filters you choose for training.

# Training a CNN

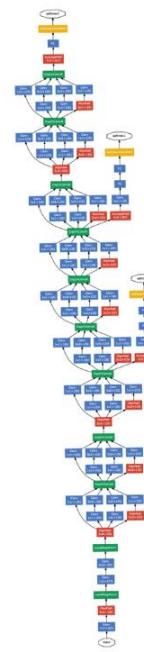
**How do you actually train these things?**

**Roughly Speaking:**

Gather  
labeled data



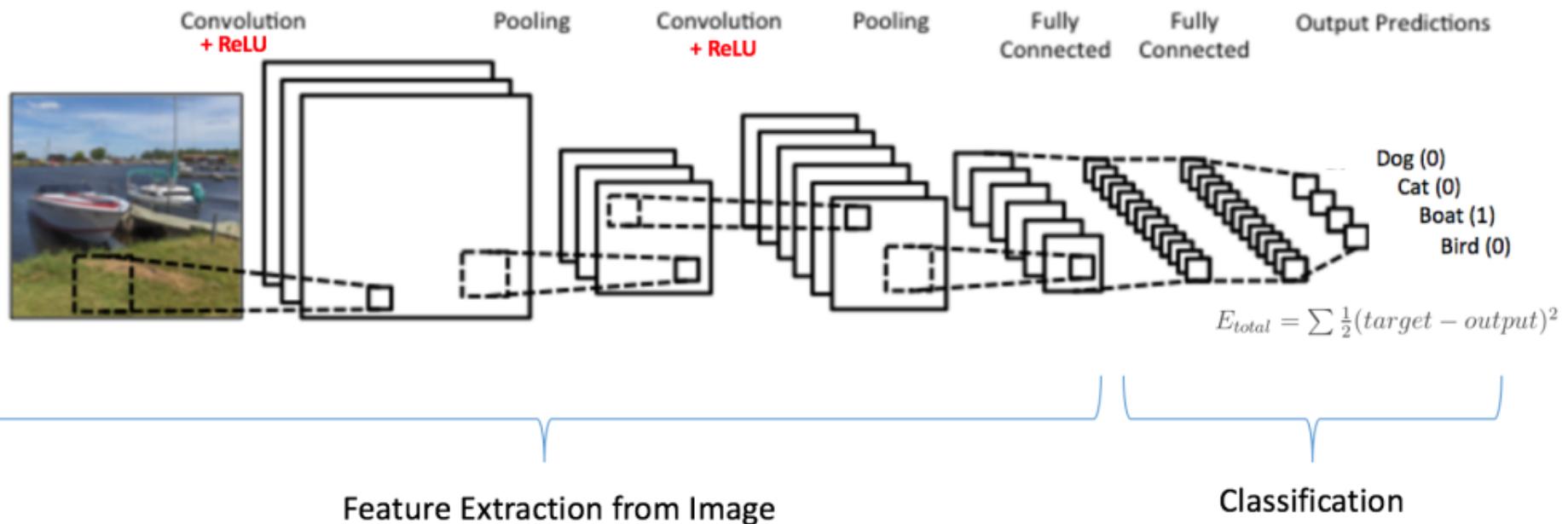
Find a ConvNet  
architecture



Minimize  
the loss



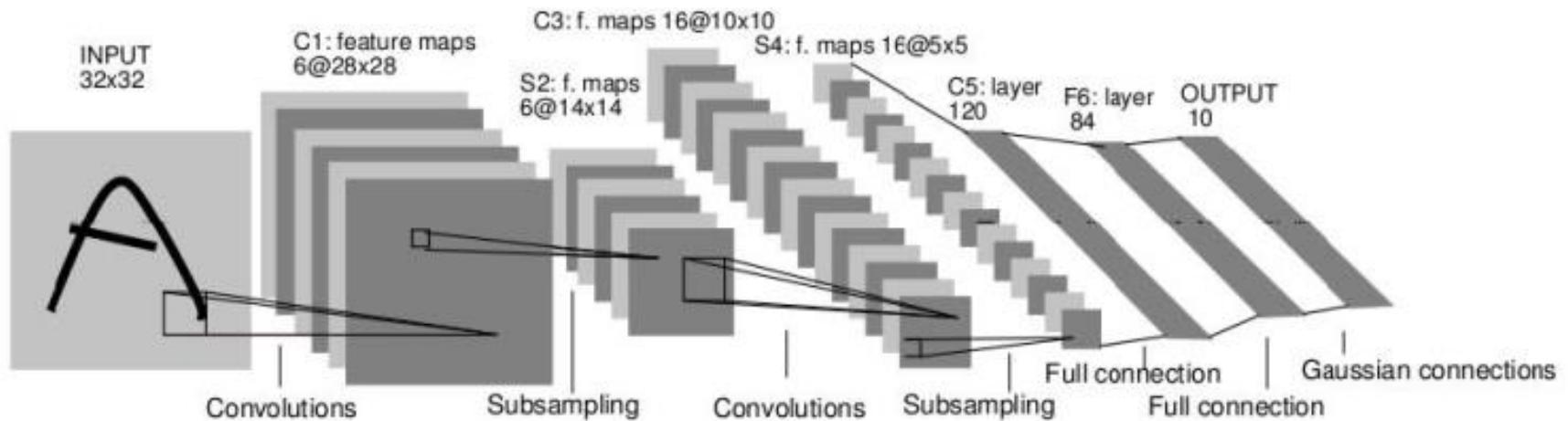
# Training a CNN - Example



# CNN Architectures

- Case Studies
  - AlexNet - Krizhevsky, Sutskever, Hinton, 2012
    - 8 layers
  - VGG - Simonyan and Zisserman, 2014
    - 16-19 layers
  - GoogLeNet - Szegedy et al., 2014
    - 22 layers
  - ResNet – He et al., 2015
    - 152-layers
- Also....
  - NiN (Network in Network) -dense network
  - Wide ResNet -fractal network
  - ResNeXT -squeezeNet
  - Stochastic Depth

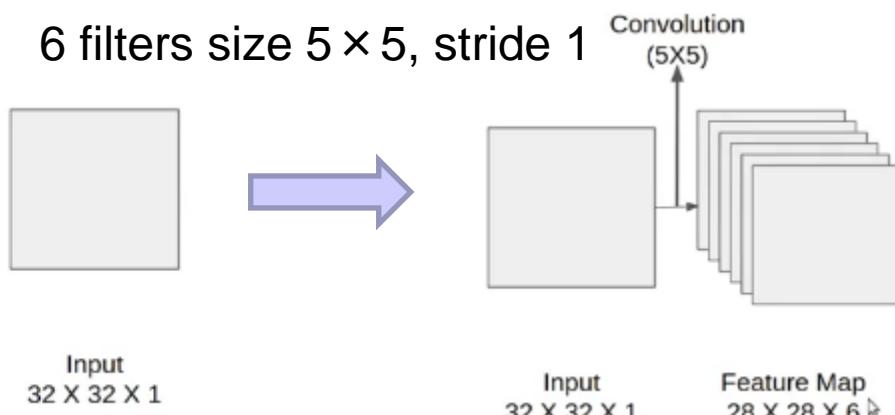
# LeNet-5 (LeCun 1998)



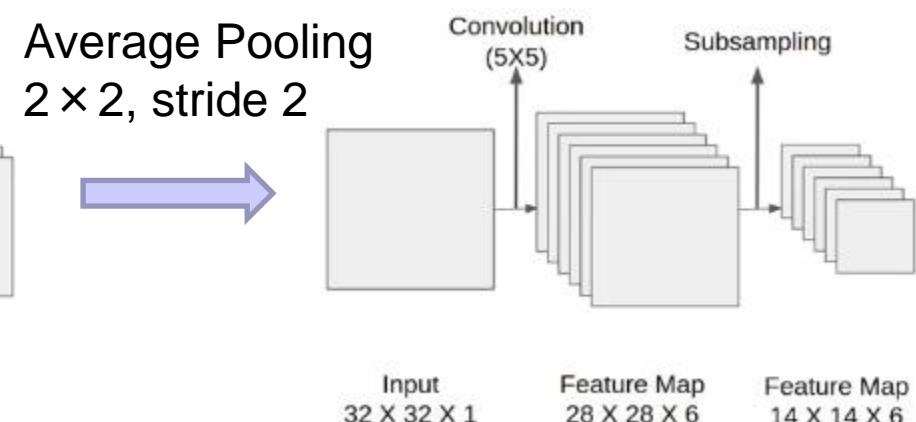
- 5 layers with learnable parameters
- How many parameters?

# LeNet-5 Architecture

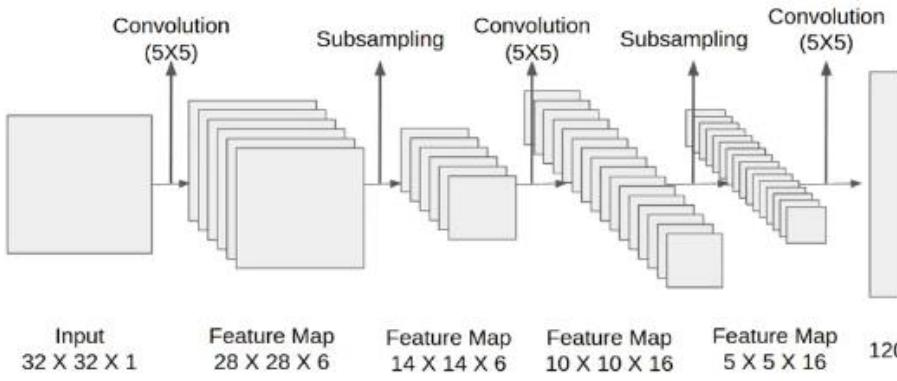
6 filters size  $5 \times 5$ , stride 1



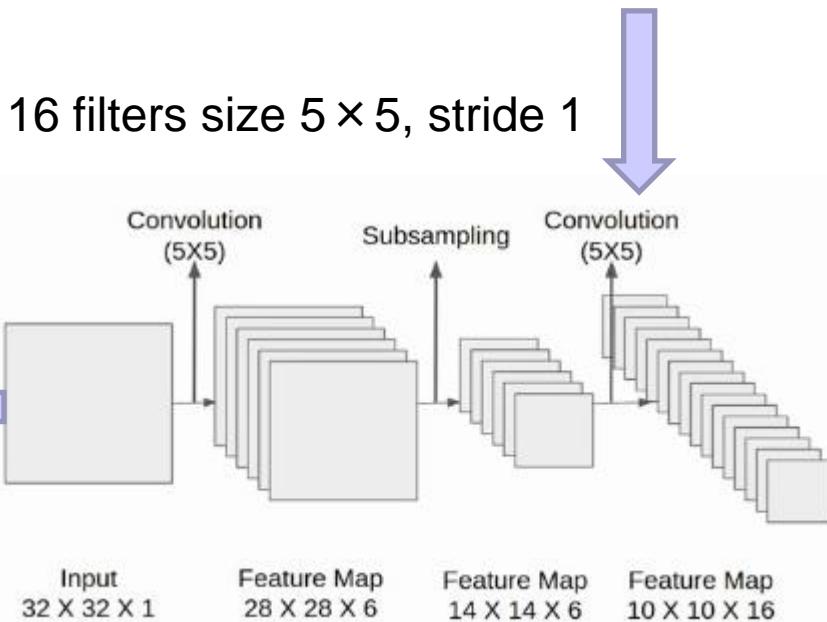
Average Pooling  
 $2 \times 2$ , stride 2



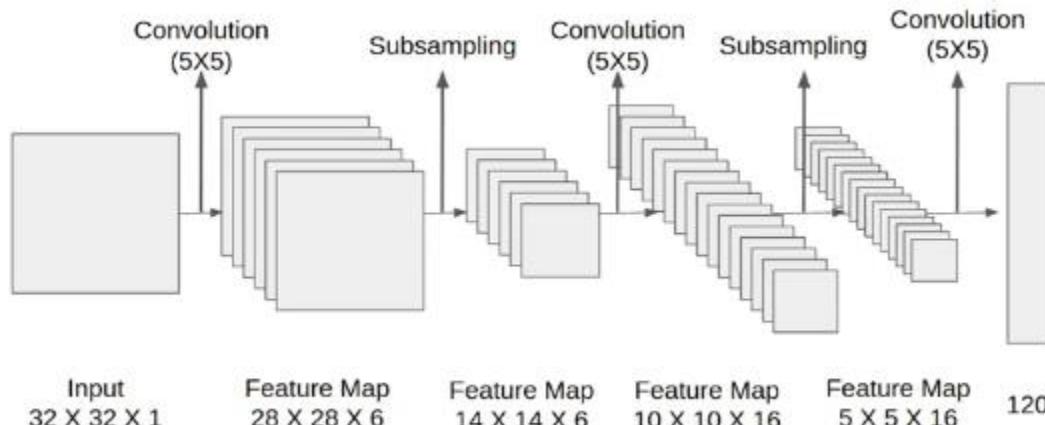
120 filters size  $5 \times 5$ , stride 1



16 filters size  $5 \times 5$ , stride 1

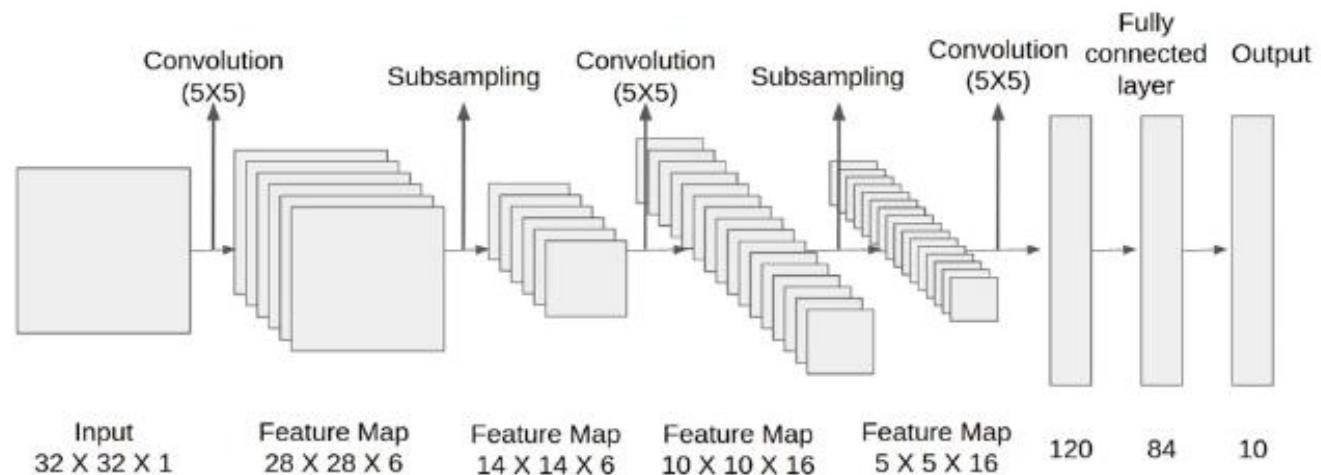
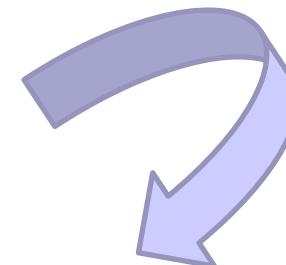


# LeNet-5 Architecture



Fully connected layer with 84 neurons

Last layer is the output layer with 10 neurons and Softmax function

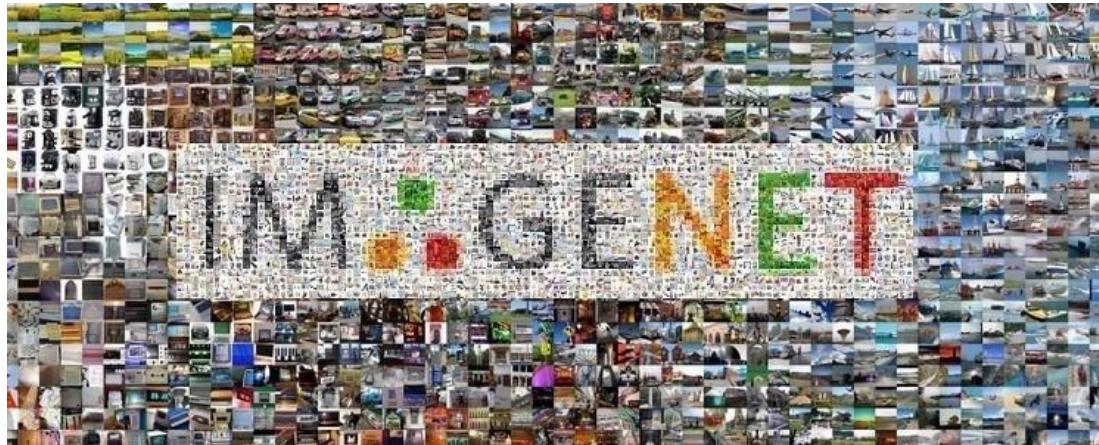


# LeNet-5 Architecture

Layer	# filters / neurons	Filter size	Stride	Size of feature map	Activation function
Input	-	-	-	32 X 32 X 1	
Conv 1	6	5 * 5	1	28 X 28 X 6	tanh
Avg. pooling 1		2 * 2	2	14 X 14 X 6	
Conv 2	16	5 * 5	1	10 X 10 X 16	tanh
Avg. pooling 2		2 * 2	2	5 X 5 X 16	
Conv 3	120	5 * 5	1	120	tanh
Fully Connected 1	-	-	-	84	tanh
Fully Connected 2	-	-	-	10	Softmax

- How many parameters?
- Around 60,000

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



- The idea conceived by Prof. Fei-Fei Li
- Launched in 2009.
- Visual Dataset contains more than 15 million of labeled high-resolution images covering almost 22,000 categories.
- High quality dataset, highly popular among researchers to test their image classification model on this dataset.

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

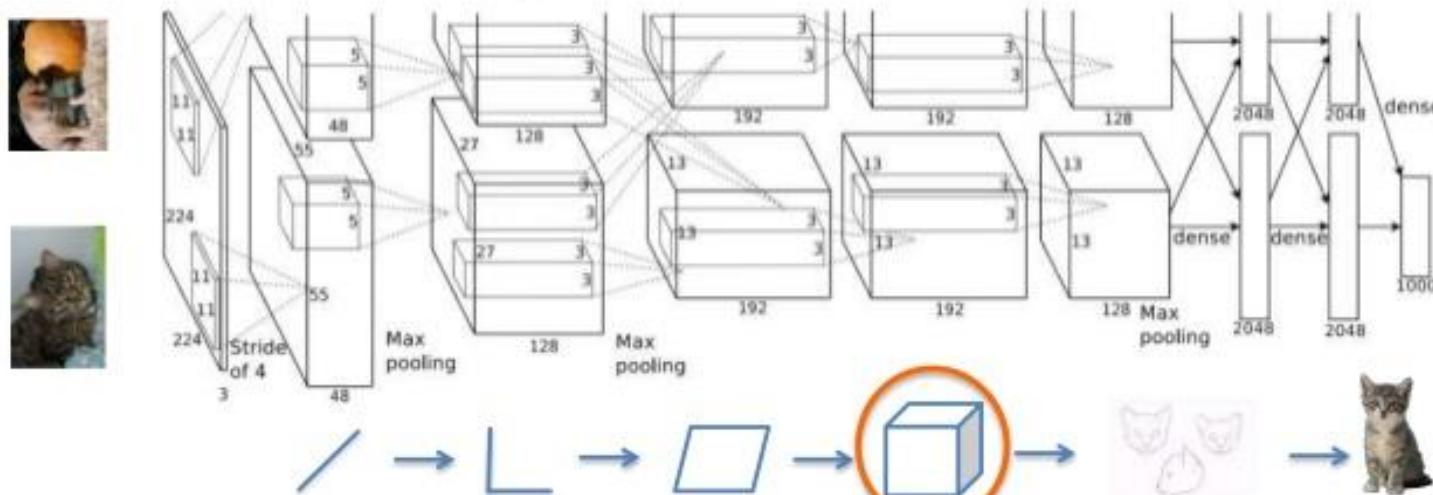
- Launched in 2010; held yearly till 2017.
- Uses the smaller portion of the ImageNet consisting of only 1000 categories.
- The total count of training images is 1.3 million, accompanied by 50,000 validation images, and 1,00,000 testing images.
- Tasks: image classification, object detection, object localization
- Minimal top-1 and top-5 error rates achievers - winner.
- Top-5 error rate - the percent of images where the correct label is not one of the model's five most likely labels

# AlexNet Architecture

- ImageNet Classification with Deep Convolutional Neural Networks

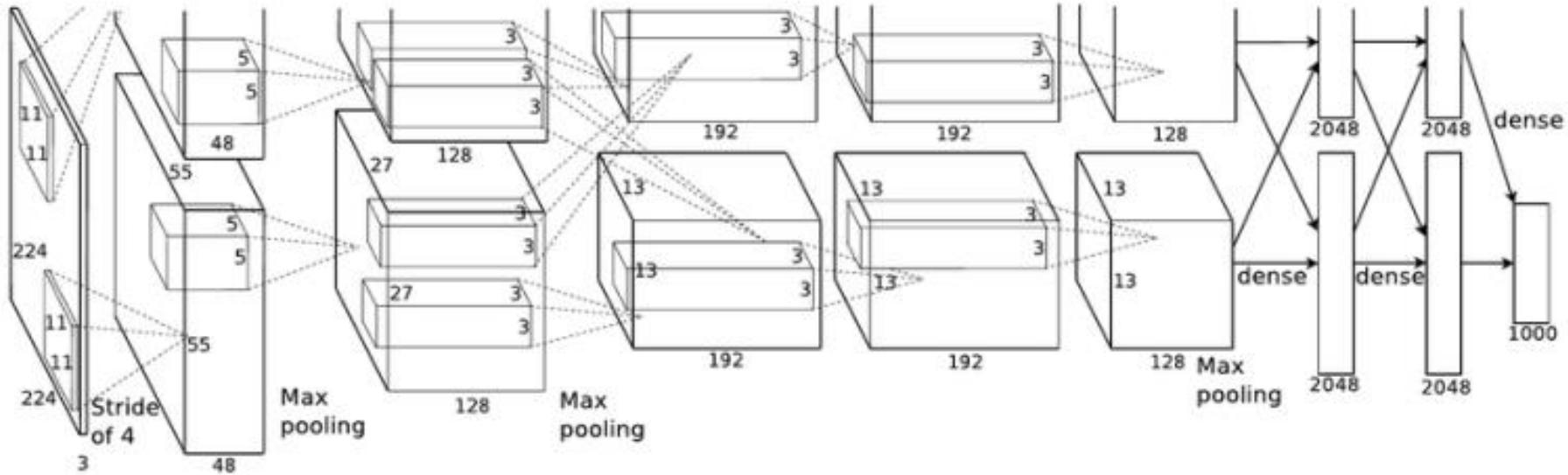
## AlexNet (Krizhevsky et al. 2012)

*The class with the highest likelihood is the one the DNN selects*



When AlexNet is processing an image, this is what is happening at each layer.

# AlexNet Architecture



- 8 layers: 5 convolutional layers + 3 fully-connected layers
- ReLU activation function to add nonlinearity & improve convergence rate
- Multiple GPUs for faster training
- 60 million parameters - susceptible to overfitting
- Data augmentation and dropouts to avoid overfitting
- Overlapping pooling to reduce the error.

# AlexNet Architecture

## Convolution and Maxpooling Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
Input	-	-	-	-	227 x 227 x 3	-
Conv 1	96	11 x 11	4	-	55 x 55 x 96	ReLU
Max Pool 1	-	3 x 3	2	-	27 x 27 x 96	-
Conv 2	256	5 x 5	1	2	27 x 27 x 256	ReLU
Max Pool 2	-	3 x 3	2	-	13 x 13 x 256	-
Conv 3	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 4	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 5	256	3 x 3	1	1	13 x 13 x 256	ReLU
Max Pool 3	-	3 x 3	2	-	6 x 6 x 256	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-

# AlexNet Architecture

## Convolution and Maxpooling Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
Input	-	-	-	-	227 x 227 x 3	-
Conv 1	96	11 x 11	4	-	55 x 55 x 96	ReLU
Max Pool 1	-	3 x 3	2	-	27 x 27 x 96	-
Conv 2	256	5 x 5	1	2	27 x 27 x 256	ReLU
Max Pool 2	-	3 x 3	2	-	13 x 13 x 256	-
Conv 3	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 4	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 5	256	3 x 3	1	1	13 x 13 x 256	ReLU
Max Pool 3	-	3 x 3	2	-	6 x 6 x 256	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-

- Width and height of output =  $(227-11)/4 + 1 = 55$
- Number of parameters in the first layer
  - $11 \times 11 \times 3 \times 96 = 34848$

# AlexNet Architecture

## Convolution and Maxpooling Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
Input	-	-	-	-	227 x 227 x 3	-
Conv 1	96	11 x 11	4	-	55 x 55 x 96	ReLU
Max Pool 1	-	3 x 3	2	-	27 x 27 x 96	-
Conv 2	256	5 x 5	1	2	27 x 27 x 256	ReLU
Max Pool 2	-	3 x 3	2	-	13 x 13 x 256	-
Conv 3	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 4	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 5	256	3 x 3	1	1	13 x 13 x 256	ReLU
Max Pool 3	-	3 x 3	2	-	6 x 6 x 256	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-

- Next layer: Overlapping Max Pooling

# AlexNet Architecture

## Convolution and Maxpooling Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
Input	-	-	-	-	227 x 227 x 3	-
Conv 1	96	11 x 11	4	-	55 x 55 x 96	ReLU
Max Pool 1	-	3 x 3	2	-	27 x 27 x 96	-
Conv 2	256	5 x 5	1	2	27 x 27 x 256	ReLU
Max Pool 2	-	3 x 3	2	-	13 x 13 x 256	-
Conv 3	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 4	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 5	256	3 x 3	1	1	13 x 13 x 256	ReLU
Max Pool 3	-	3 x 3	2	-	6 x 6 x 256	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-

- 2<sup>nd</sup> convolution operation
- Max-pooling layer

# AlexNet Architecture

## Convolution and Maxpooling Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
Input	-	-	-	-	227 x 227 x 3	-
Conv 1	96	11 x 11	4	-	55 x 55 x 96	ReLU
Max Pool 1	-	3 x 3	2	-	27 x 27 x 96	-
Conv 2	256	5 x 5	1	2	27 x 27 x 256	ReLU
Max Pool 2	-	3 x 3	2	-	13 x 13 x 256	-
Conv 3	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 4	384	3 x 3	1	1	13 x 13 x 384	ReLU
Conv 5	256	3 x 3	1	1	13 x 13 x 256	ReLU
Max Pool 3	-	3 x 3	2	-	6 x 6 x 256	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-

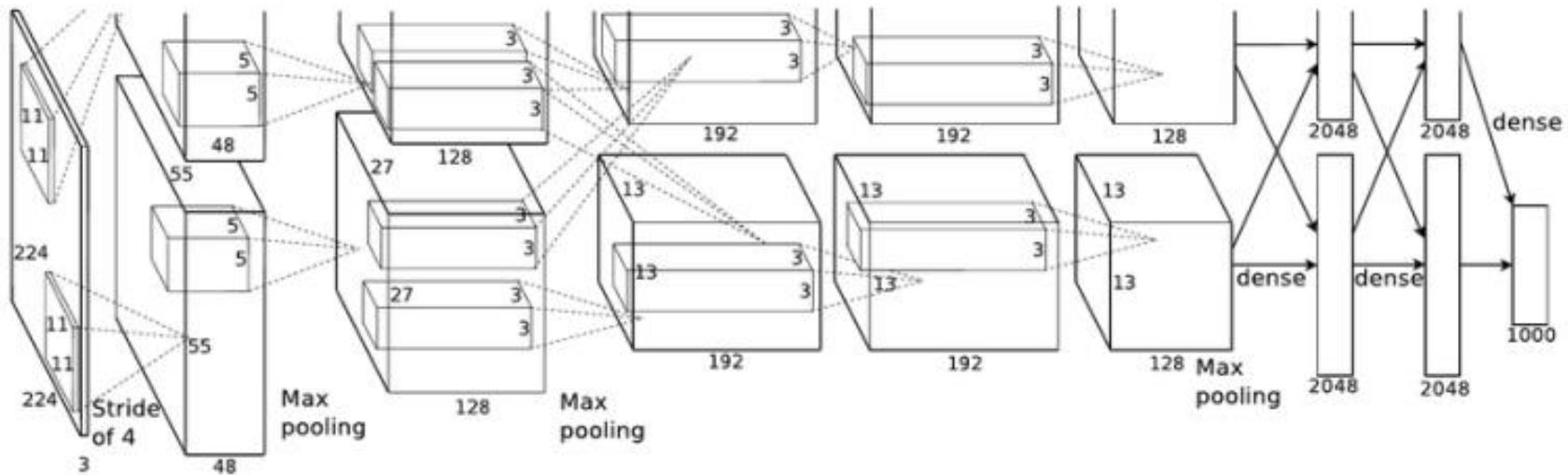
# AlexNet Architecture

## Fully Connected and Dropout Layers

Layer	# filters / neurons	Filter size	Stride	Padding	Size of feature map	Activation function
-	-	-	-	-	-	-
-	-	-	-	-	-	-
-	-	-	-	-	-	-
Dropout 1	rate = 0.5	-	-	-	6 x 6 x 256	-
Fully Connected 1	-	-	-	-	4096	ReLU
Dropout 2	rate = 0.5	-	-	-	4096	-
Fully Connected 2	-	-	-	-	4096	ReLU
Fully Connected 3	-	-	-	-	1000	Softmax

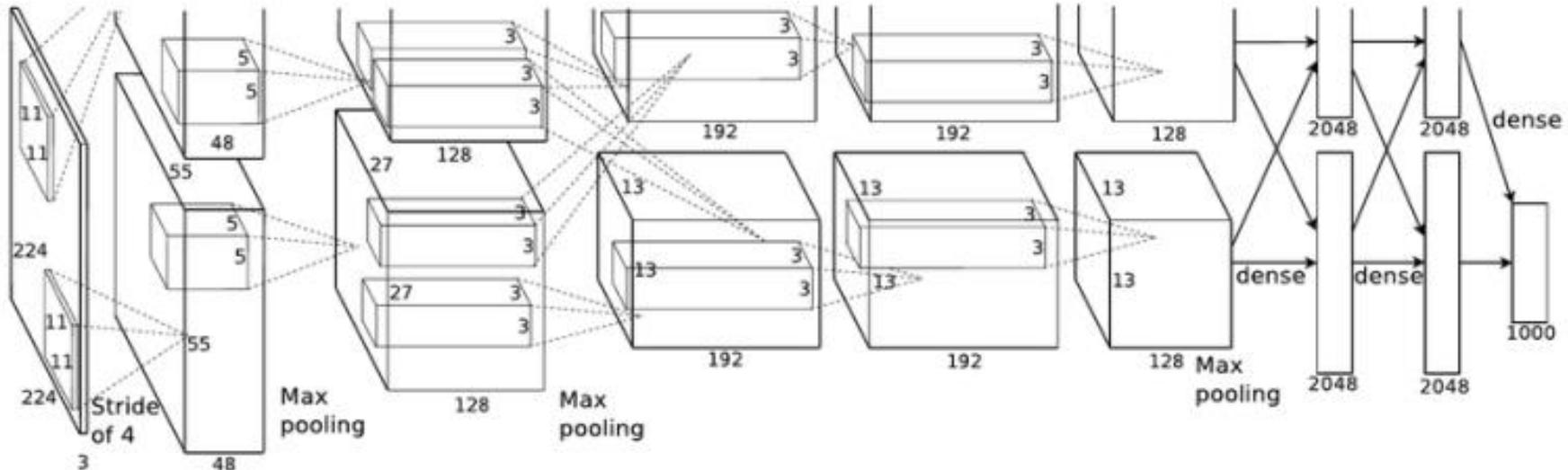
- Last fully connected layer or output layer with 1000 neurons as we have 1000 classes in the data set.
- Parameters: Dropout rate 0.5, Batch size = 128, Weight decay term: 0.0005, Momentum term  $\alpha = 0.9$ , learning rate  $\eta = 0.01$ , manually reduced by factor of ten on monitoring validation loss.

# AlexNet Architecture



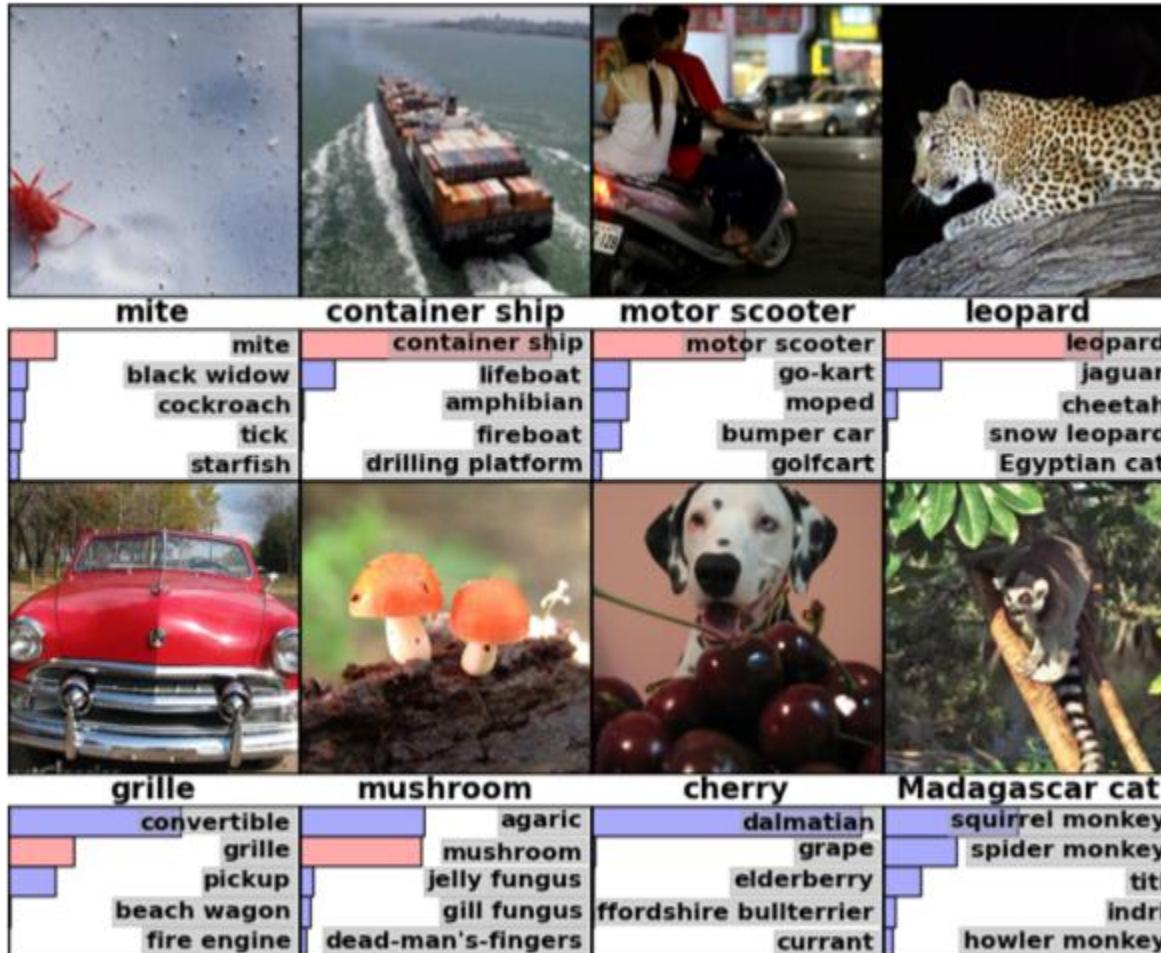
- Popularized the use of ReLUs
- Used heavy data augmentation (flipped images, random crops of size  $227 \times 227$ , color normalization, etc. )
- 62.3 million learnable parameters

# AlexNet Architecture



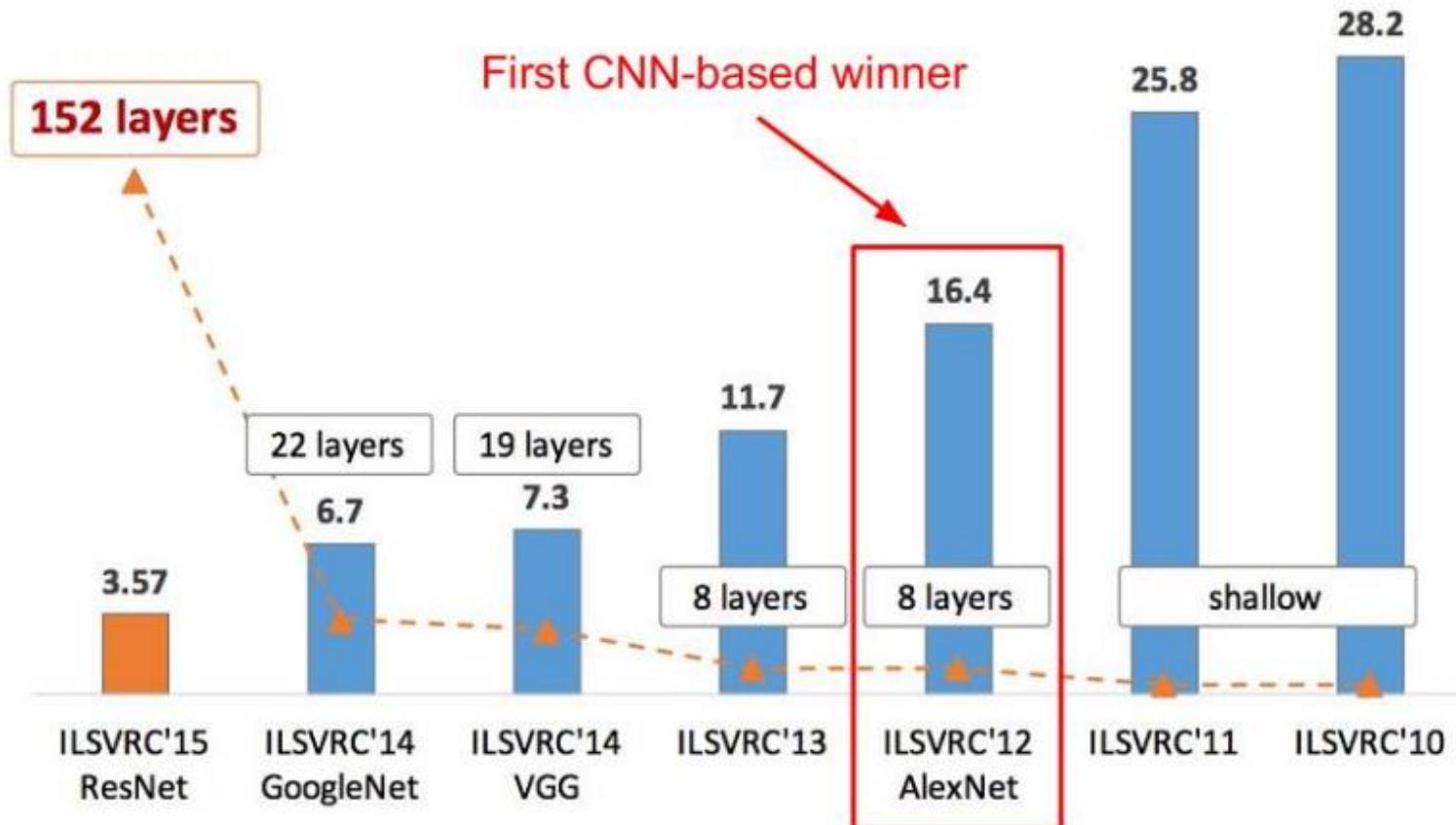
- AlexNet was trained on a GTX 580 GPU with only 3 GB of memory which couldn't fit the entire network.
- The network was split across 2 GPUs, with half of the neurons (feature maps) on each GPU.
- See a split in the architecture diagram.

# AlexNet Results



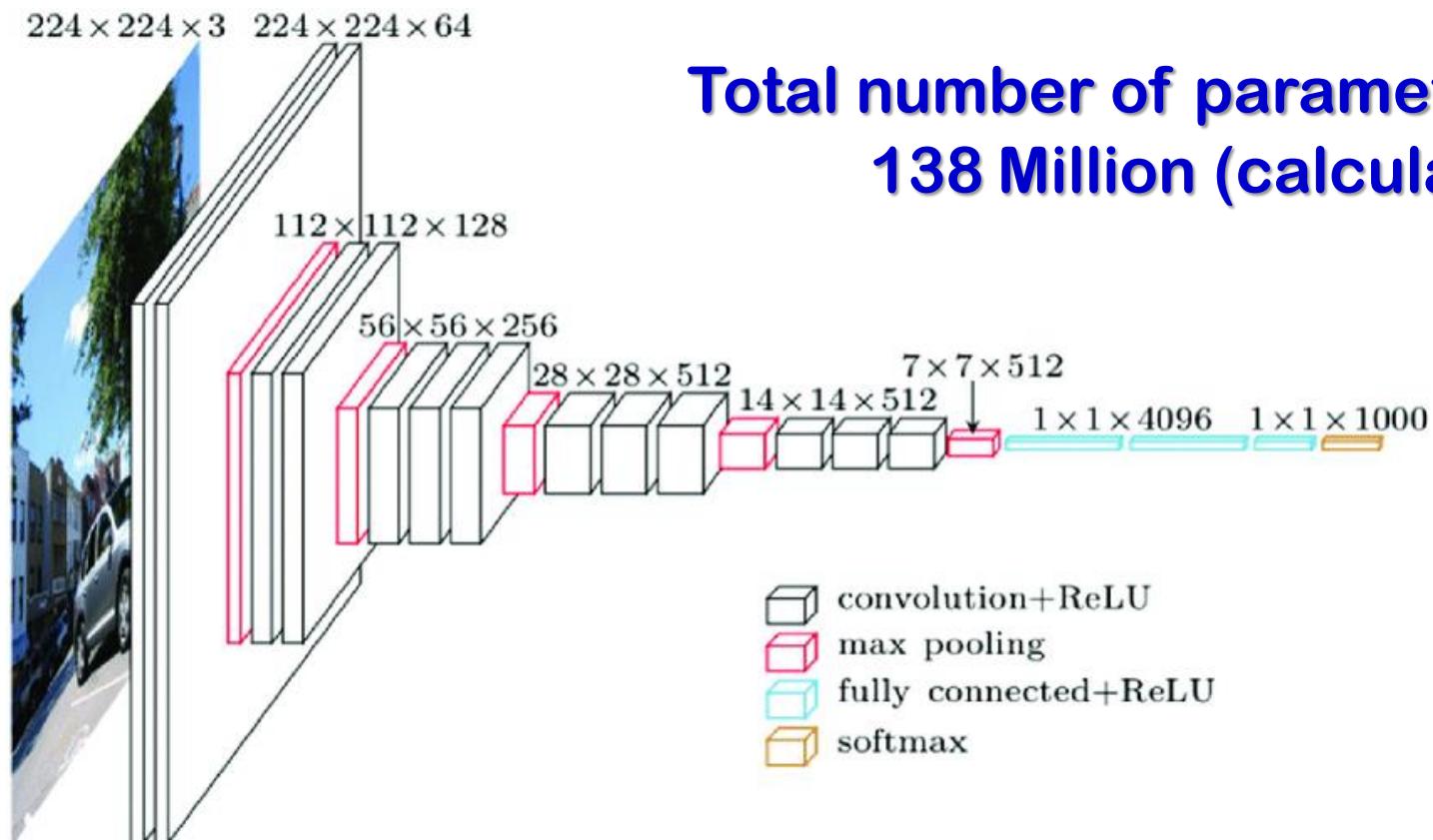
# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

Classification: Top-5 error



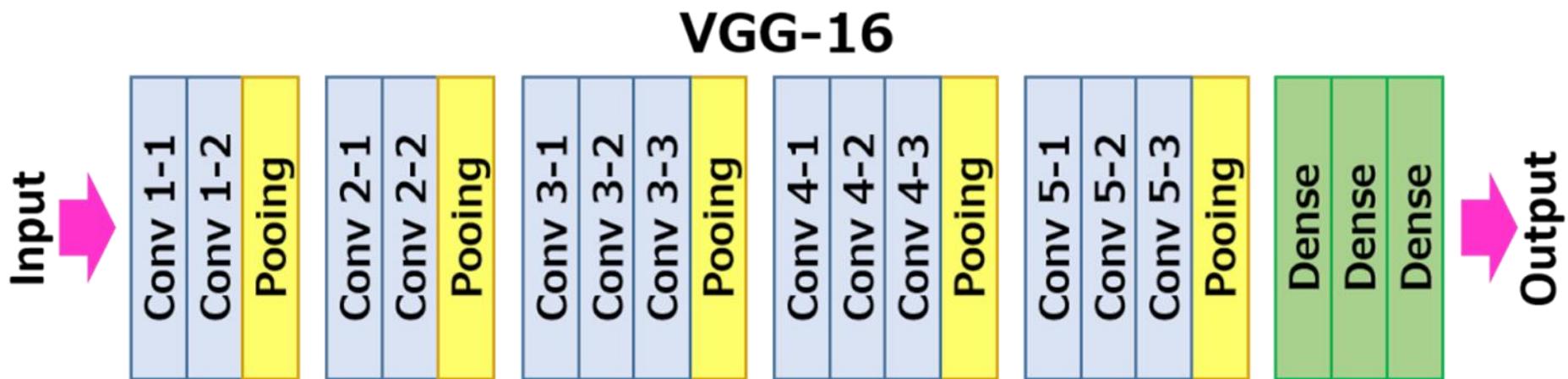
# VGG

- K. Simonyan and A. Zisserman (2014) “Very Deep Convolutional Networks for Large-Scale Image Recognition”



# VGG

- K. Simonyan and A. Zisserman (2014) “Very Deep Convolutional Networks for Large-Scale Image Recognition”



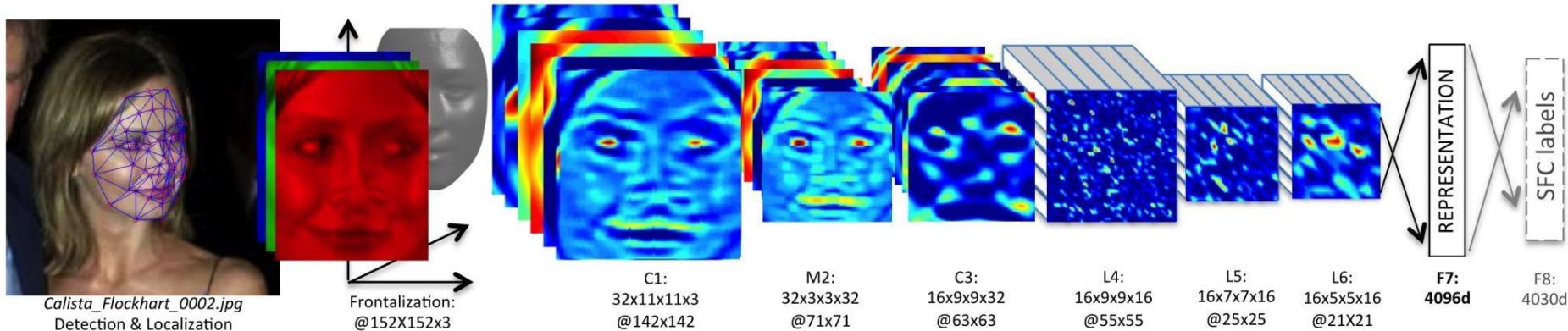
**Reduce the size of the convolution kernel and increase the number of convolution layers**

# VGG Configurations

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 <b>conv3-256</b> <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 <b>conv3-512</b> <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 <b>conv3-512</b> <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

# Industry Deployment

- Used in Facebook, Google, Microsoft
- Image Recognition, Speech Recognition, ....
- Fast at test time



Taigman et al. DeepFace: Closing the Gap to Human-Level Performance in Face Verification, CVPR'14

# References

- Ian Goodfellow, Deep Learning, MIT Press, 2016
- Andrew Ng, Coursera Course on CNN
- Lecture slides of Prof. Aparajita Ojha, IIITDM Jabalpur (CS617 Basics of Deep Learning)
- Lecture slides of Prof. Ming Li Lecture slides, University of Waterloo, Canada (CS 898: Deep Learning and Its Applications, Spring 2017).
- Lecture slides of Prof. Fei-Fei Li, Justin Johnson, Serena Yeung, Stanford University, (CS231n Convolution Neural Networks for Visual recognition).
- Lecture slides of Prof. Ioannis Gkioulekas, Carnegie Mellon University (16-385 Computer Vision, Spring 2019).