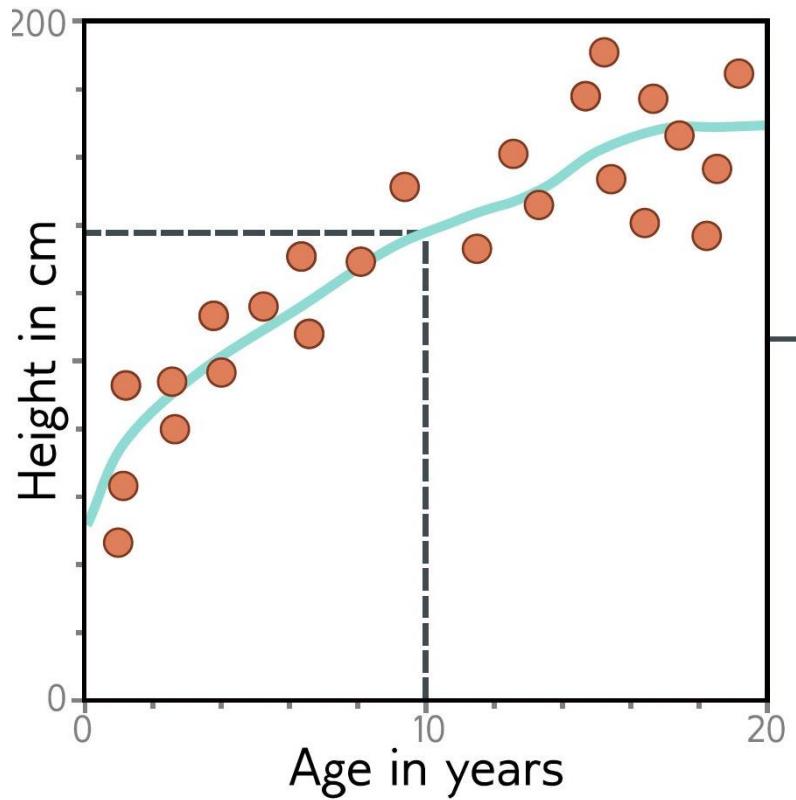


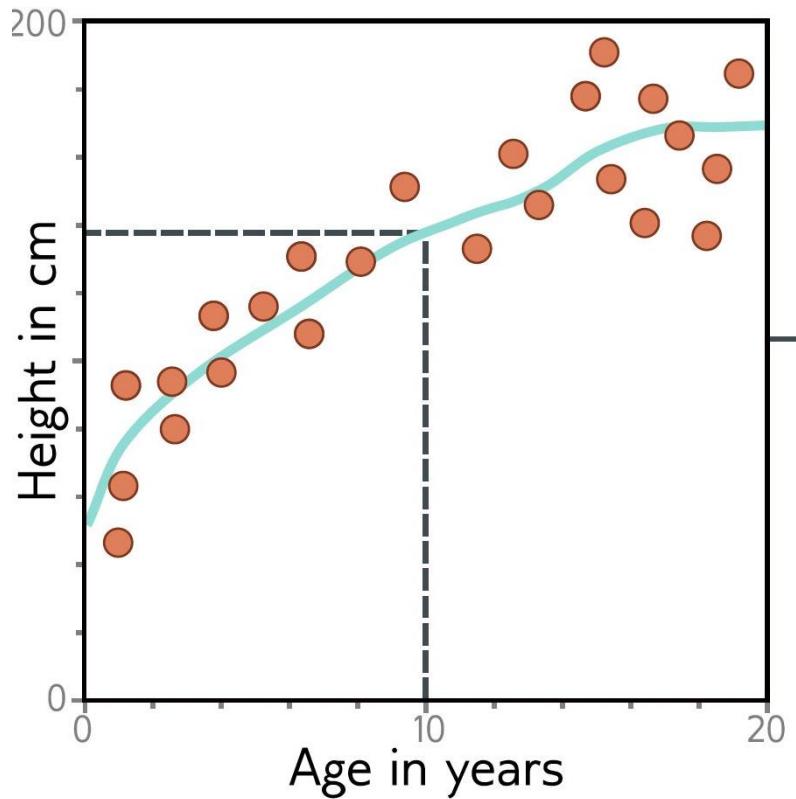
Machine Learning & Neural Networks – 7PAM2021

Supervised Learning in a Nutshell



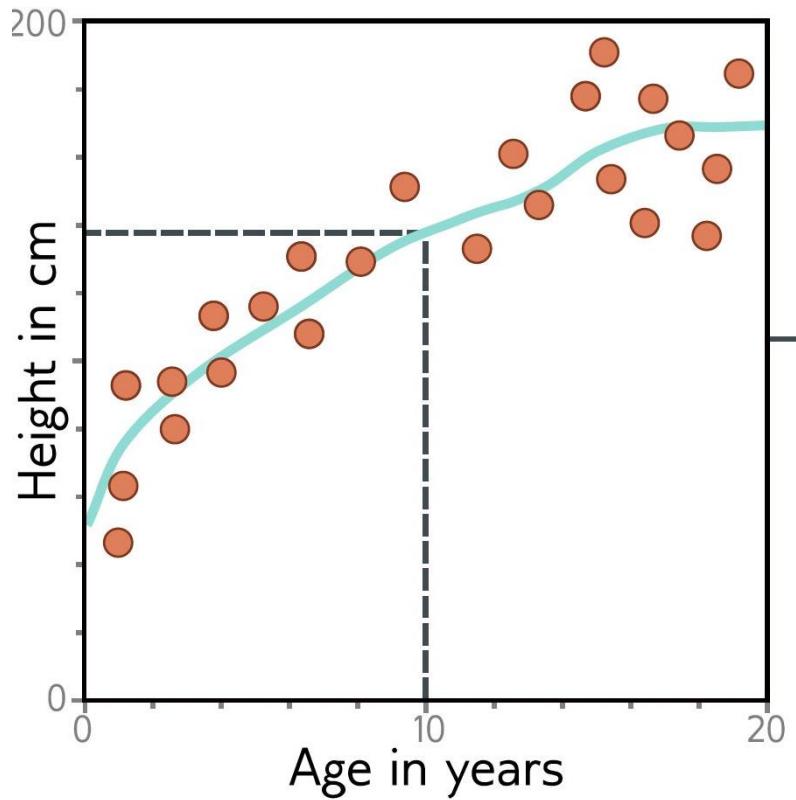
$$y = f[x, \phi]$$

Supervised Learning in a Nutshell



$$\begin{aligned}y &= f[x, \phi] \\&= \phi_0 + \phi_1 x\end{aligned}$$

Supervised Learning in a Nutshell

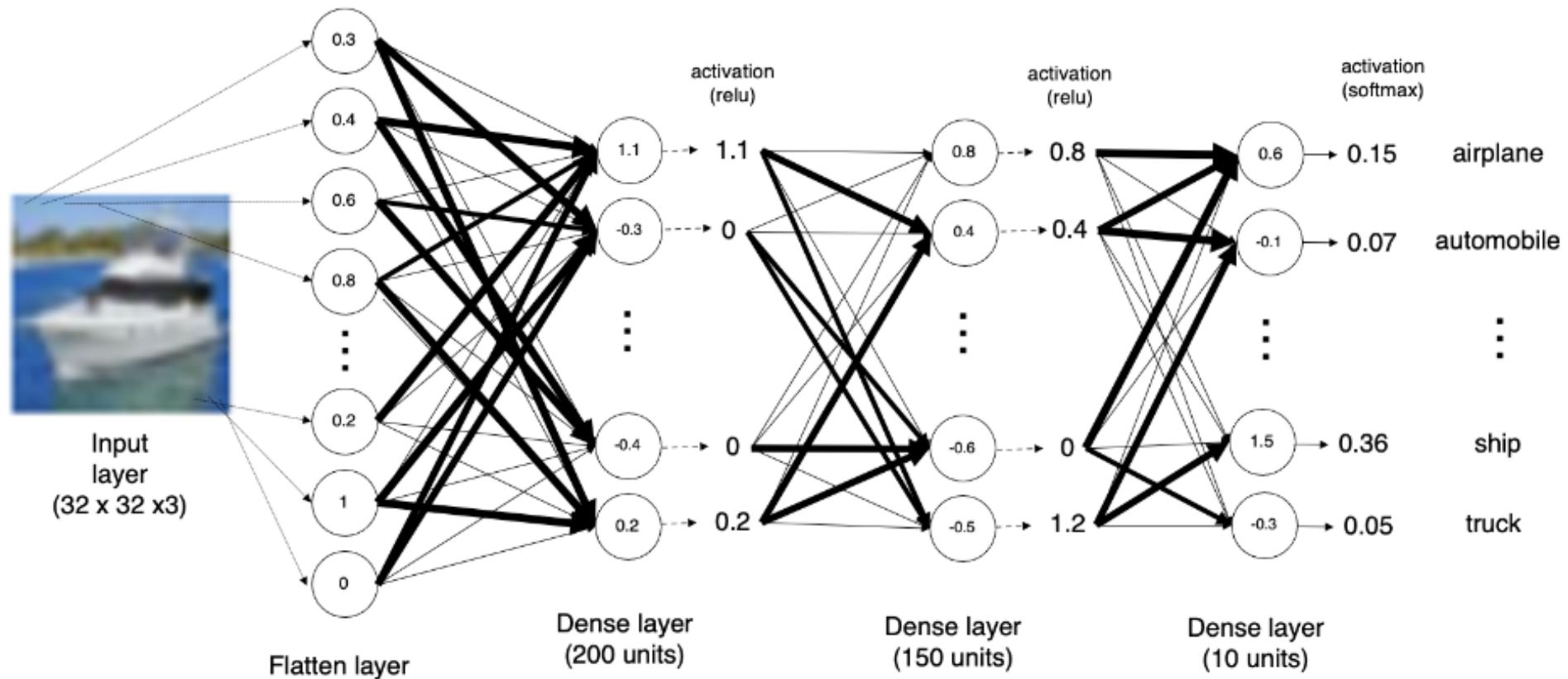


$$y = f[x, \phi]$$

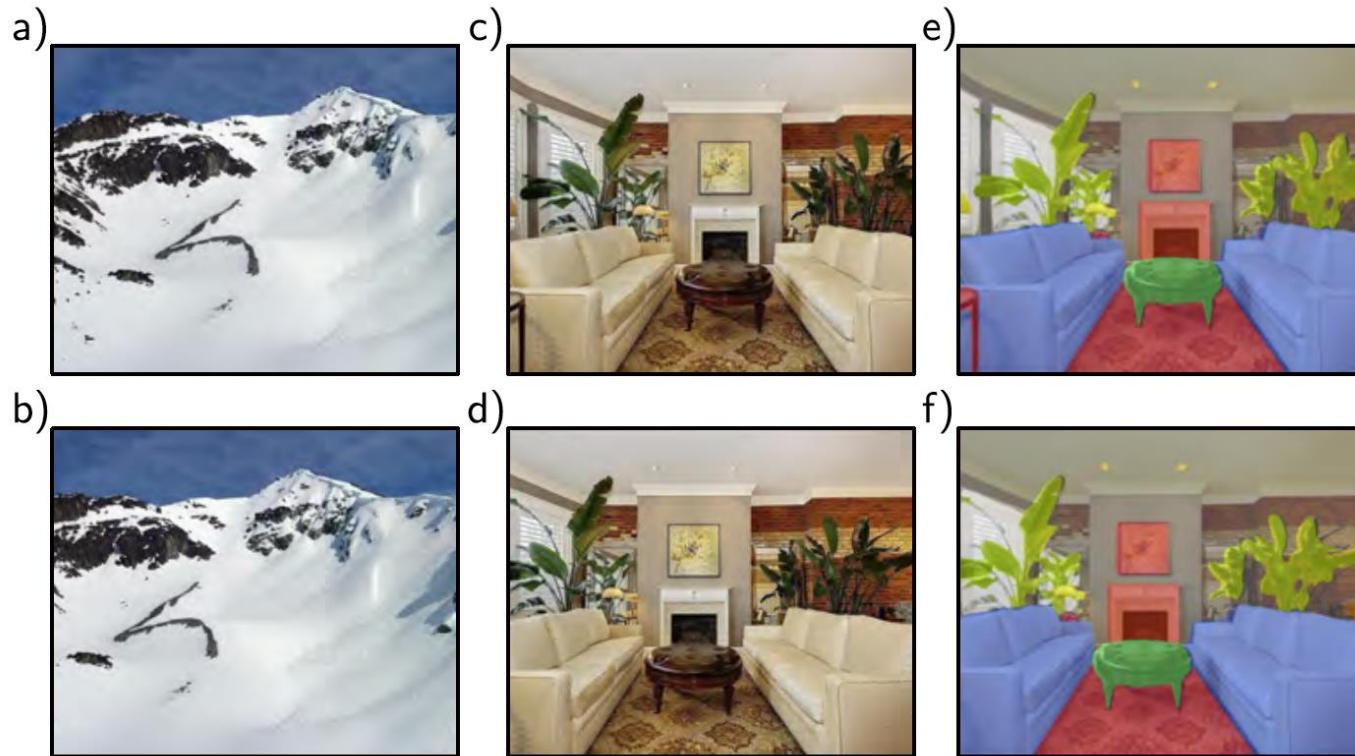
$$= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x].$$

Why do we need CNNs?

The architecture of MLPs is fundamentally limited in processing image data due to their fully connected nature, which does not take advantage of the spatial and hierarchical structure of images.



Invariance and equivariance



$$f[t[x]] = f[x]$$

Image Classification

One of the most common forms of supervised learning is identifying the contents of an image

This is a task that humans are **VERY** good at—in most cases—but is harder than it looks for a computer

Predictions of labels are often expressed as probabilities—95% chance the image is a cat



Cat? Or Dog?

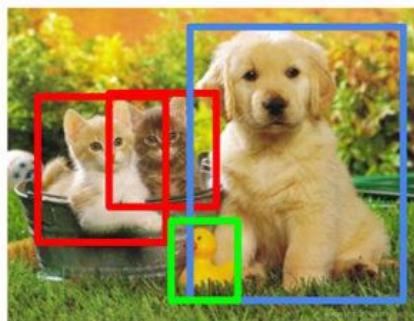
Object Detection

Classification



CAT

Object Detection



CAT, DOG, DUCK

Like image classification, object detection attempts to predict the contents of an image, but now with multiple objects, not just one

The models are trained to predict the “Bounding Boxes” of objects in the image, and apply labels for the contents of each box

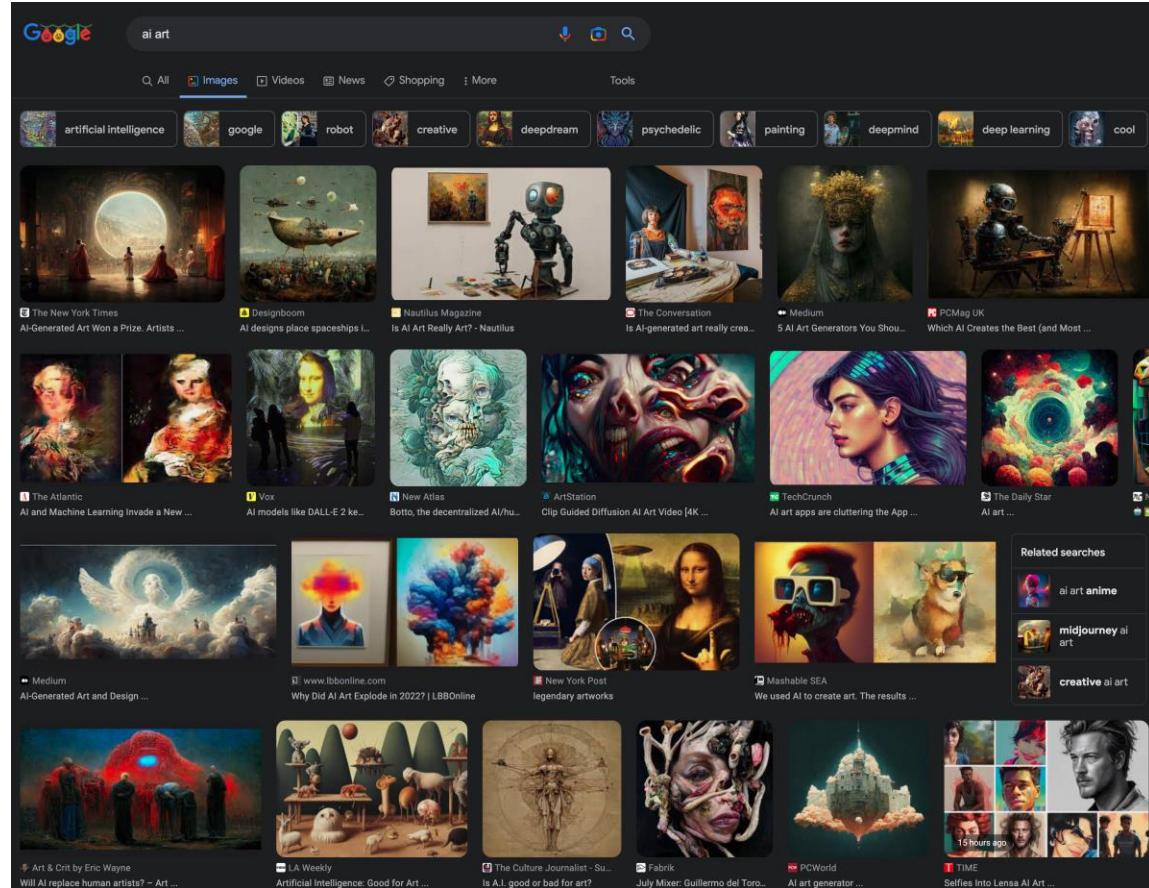


Image Generation

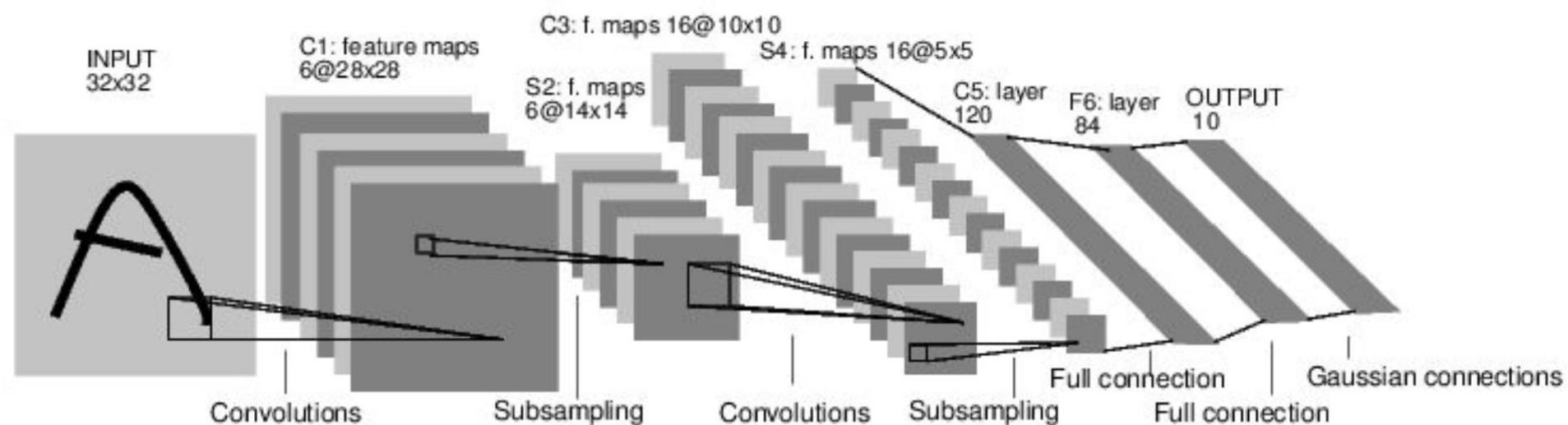
Increasingly popular in the last few years, training models to create images, either from scratch or with a prompt, has many applications.

Raise of Stable Diffusion has led to an explosion of “AI Art”, but other models include VAEs and Generalised Adversarial Nets (GANs)

Can be used for image restoration and scene reconstruction.

Convolutional Layers

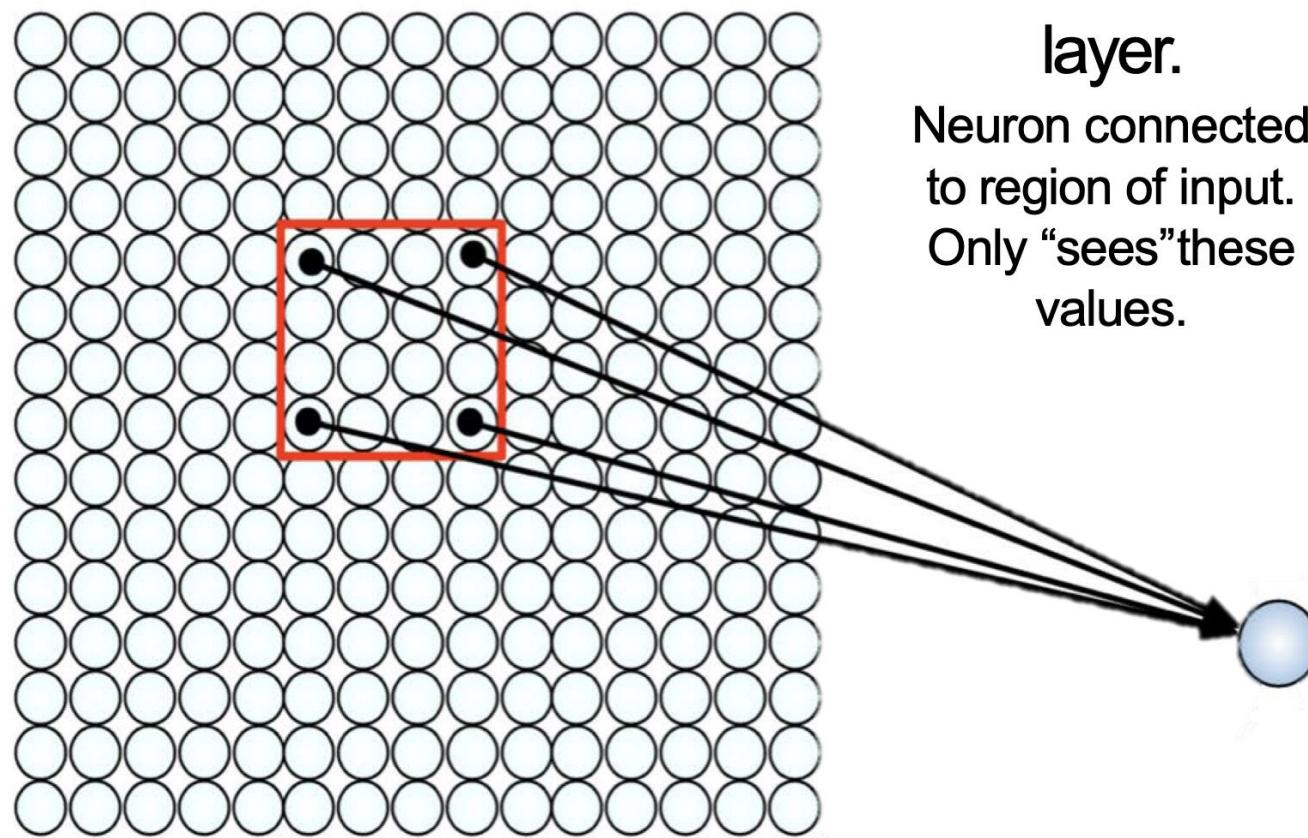
Convolutional Neural Networks



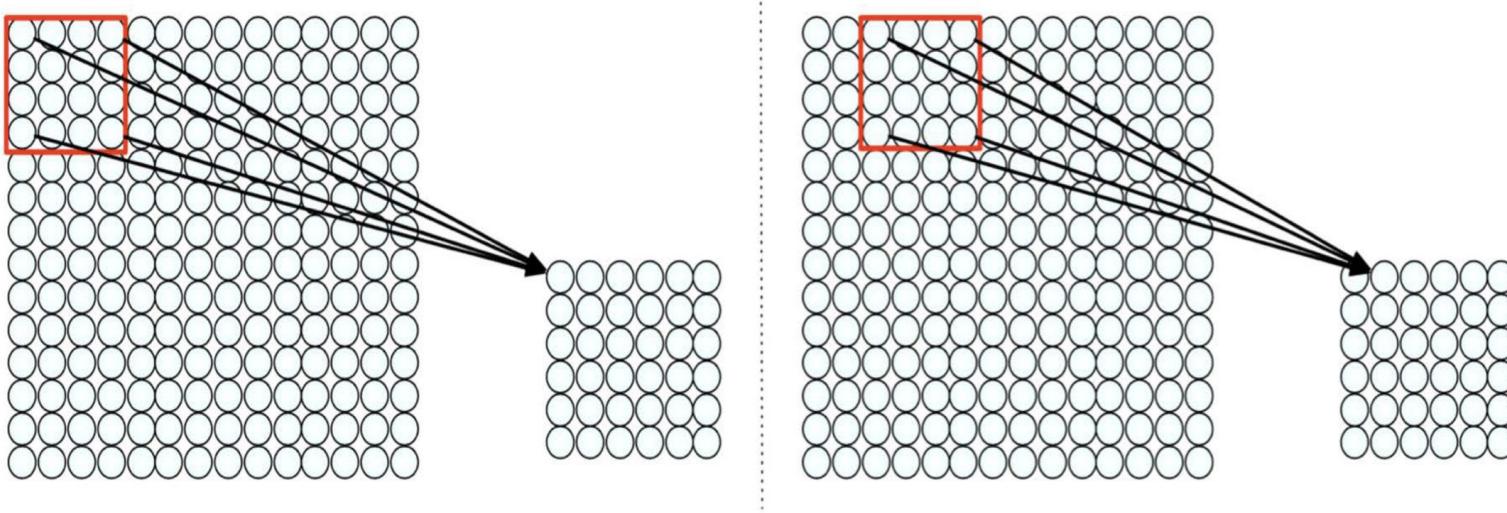
Input: 2D
image.
Array of pixel
values

Idea: connect
patches of input to
neurons in hidden
layer.

Neuron connected
to region of input.
Only “sees”these
values.



Using Spatial Structure

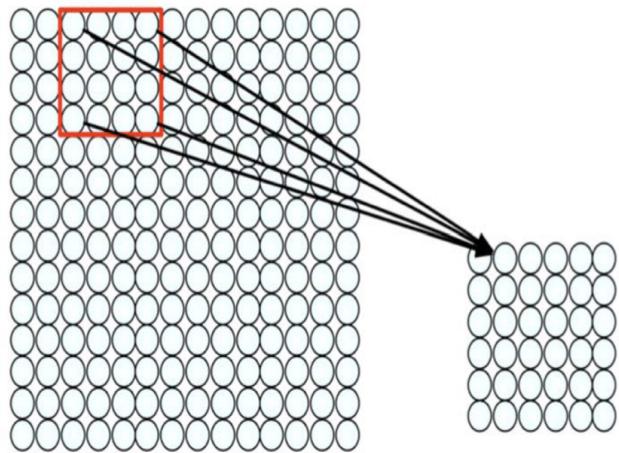


Connect patch in input layer to a single neuron in subsequent layer:

Use a sliding window to define connections.

How can we weight the patch to detect particular features?

Feature Extraction with Convolution



- Filter of size 4×4 : 16 different weights
- Apply this same filter to 4×4 patches in input
- Shift by 2 pixels for next patch

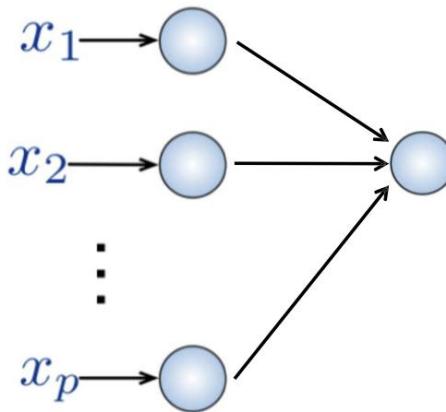
This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Fully Connected Neural Network

Input:

- 2D image
- Vector of pixel values

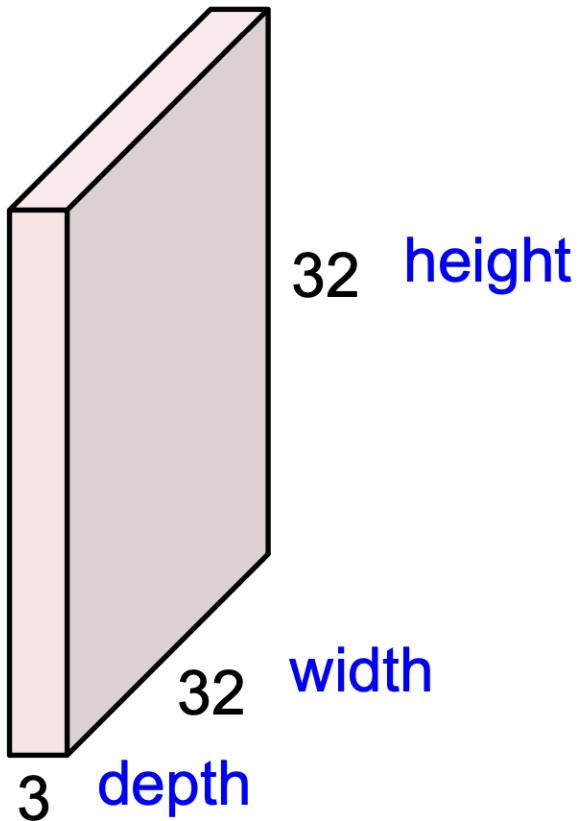
**Fully Connected:**

- Each neuron in hidden layer connected to all neurons in input layer
- No spatial information
- Many, many parameters

Key idea: Use spatial structure in input to inform architecture of the network

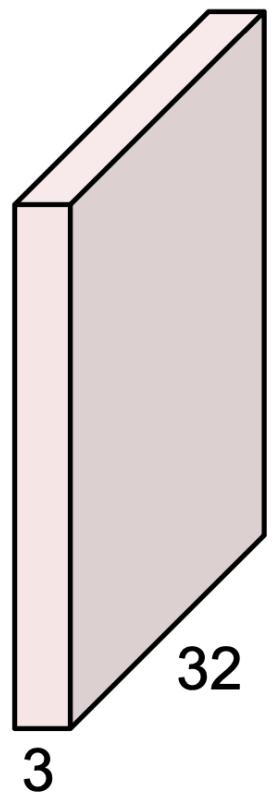
Convolution Layer

32x32x3 image



Convolution Layer

32x32x3 image



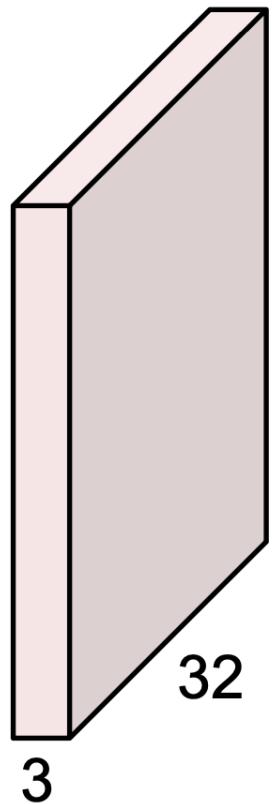
5x5x3 filter



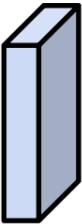
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

32x32x3 image



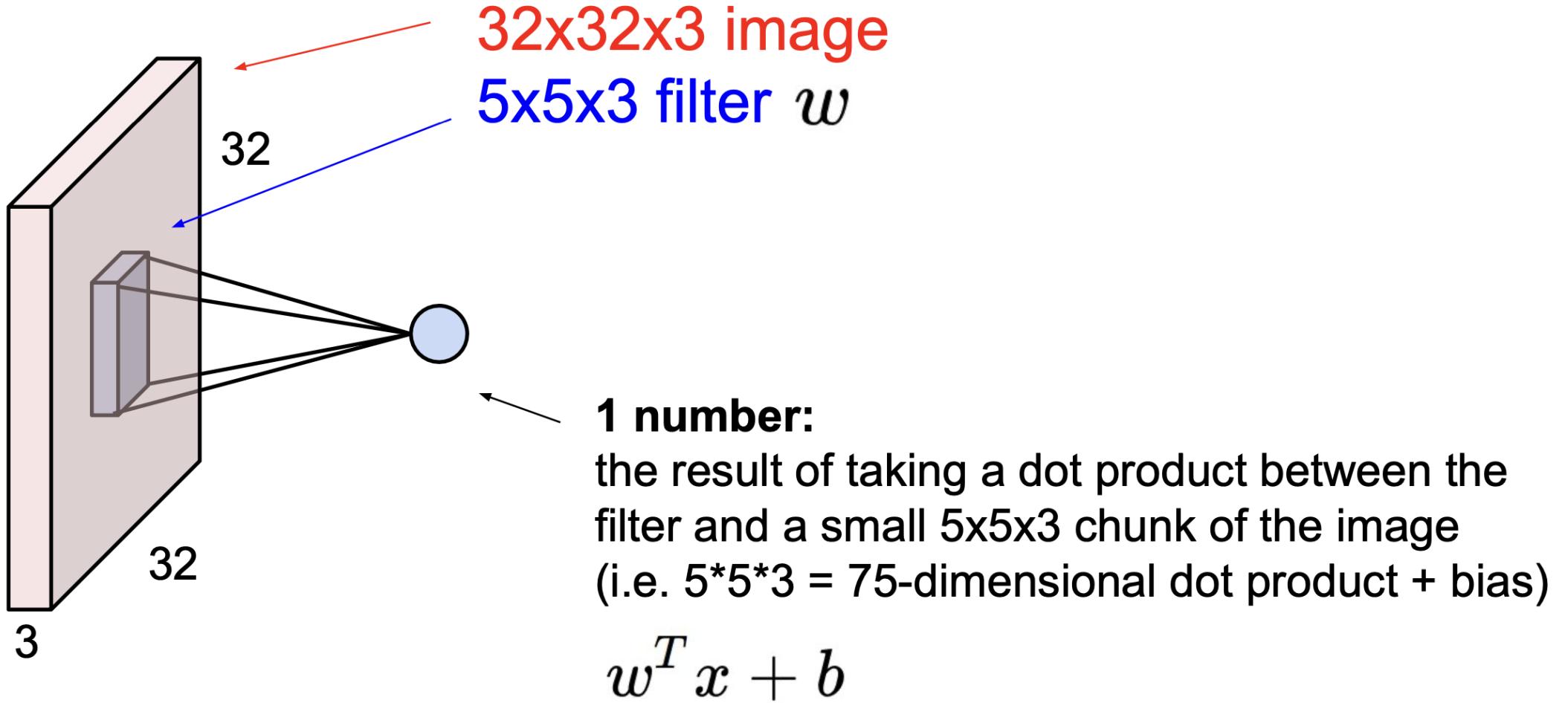
5x5x3 filter



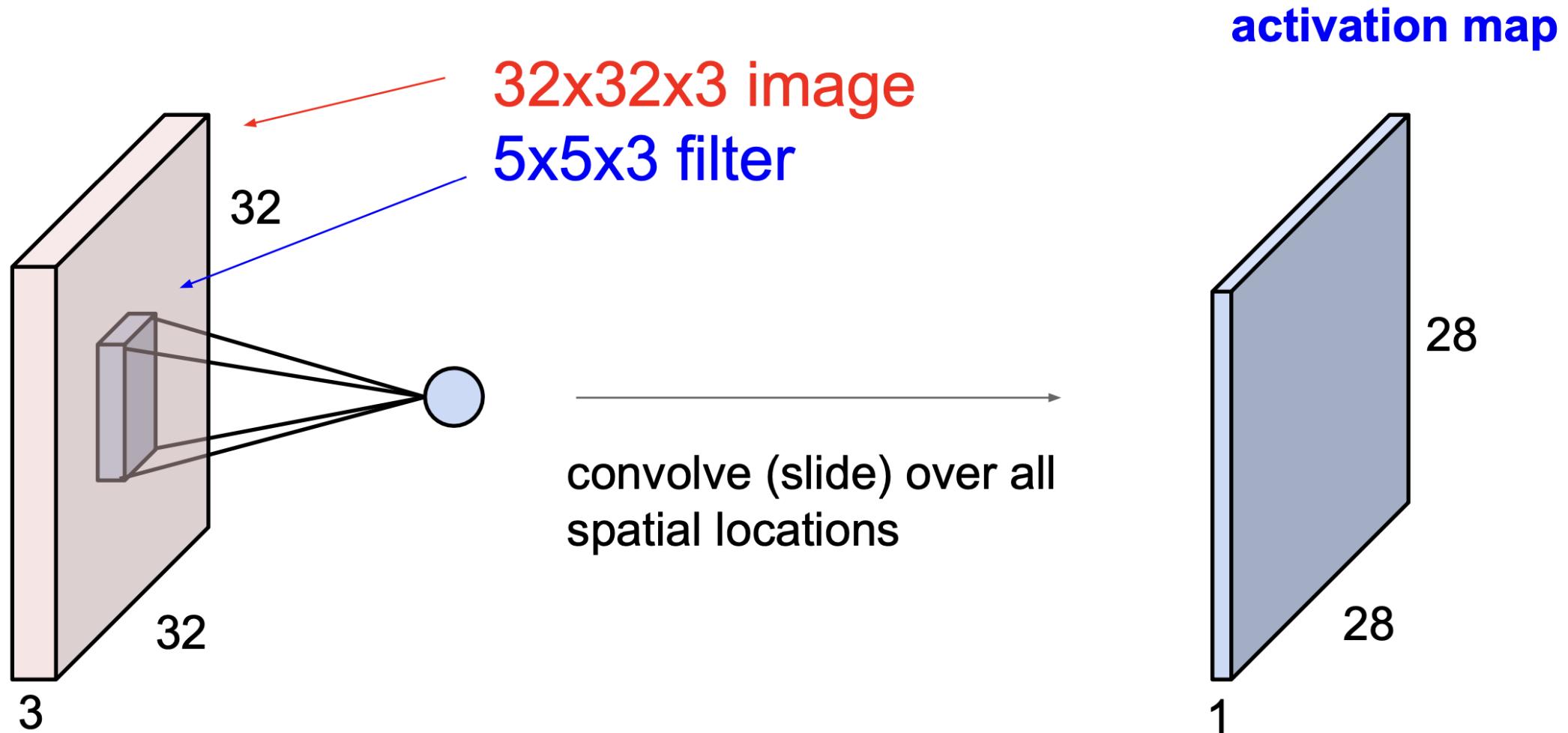
Filters always extend the full depth of the input volume

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

Convolution Layer



Convolution Layer

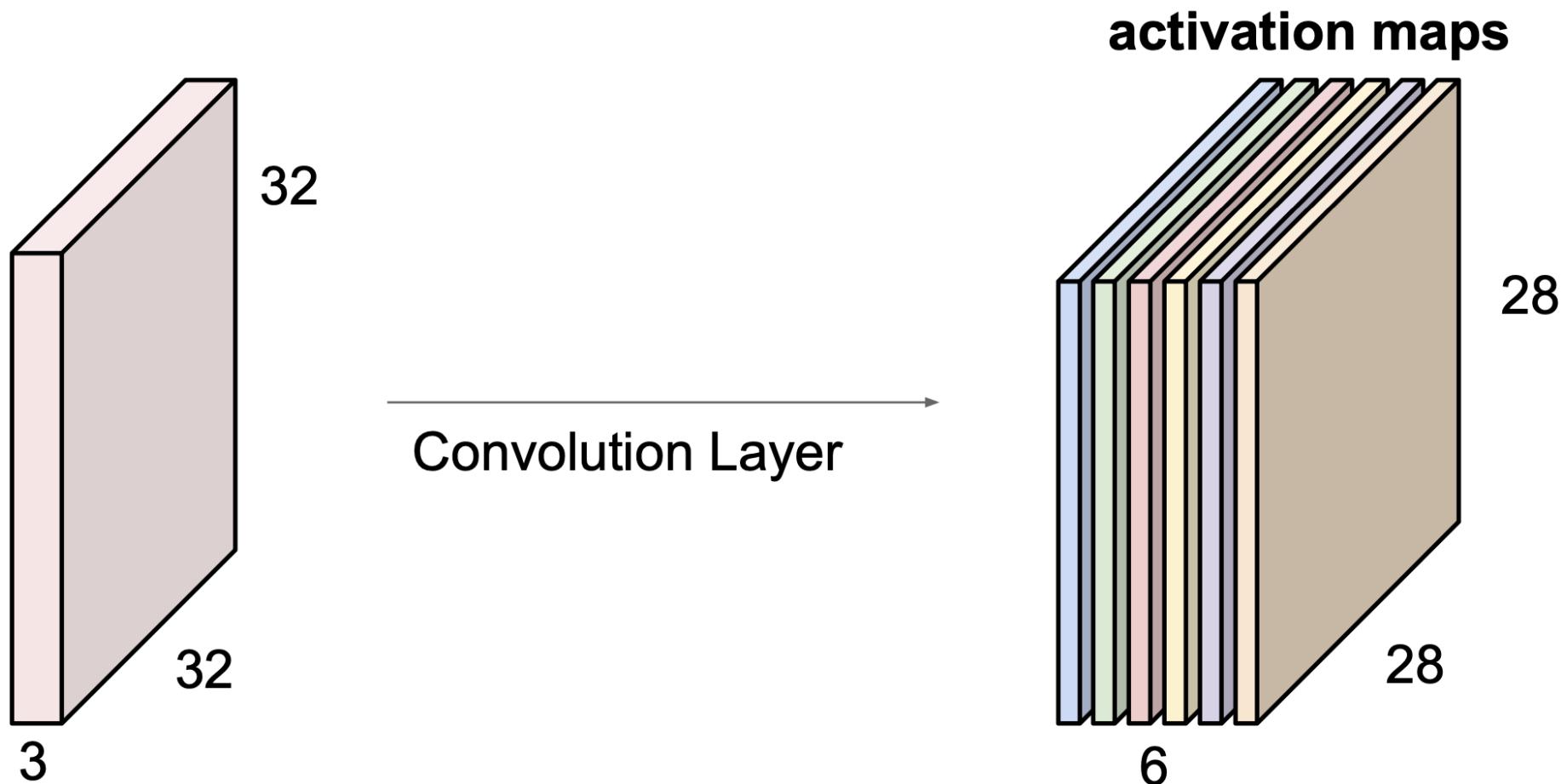


Convolution Layer

consider a second, green filter

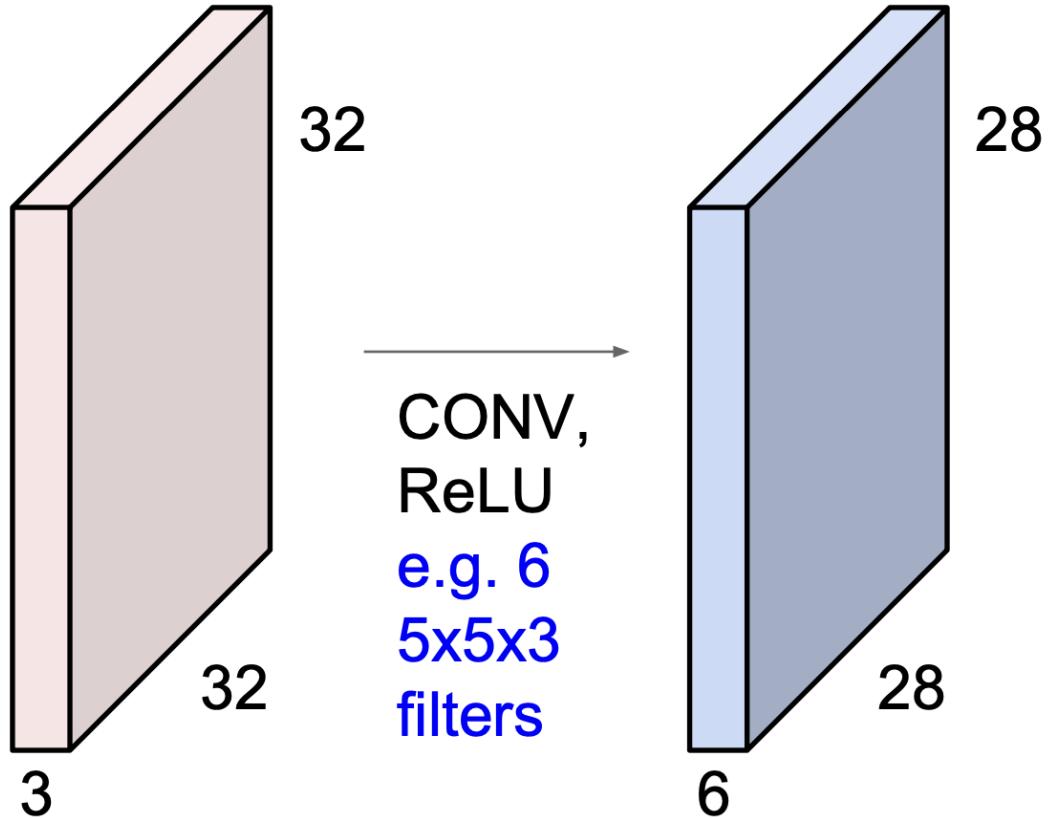


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

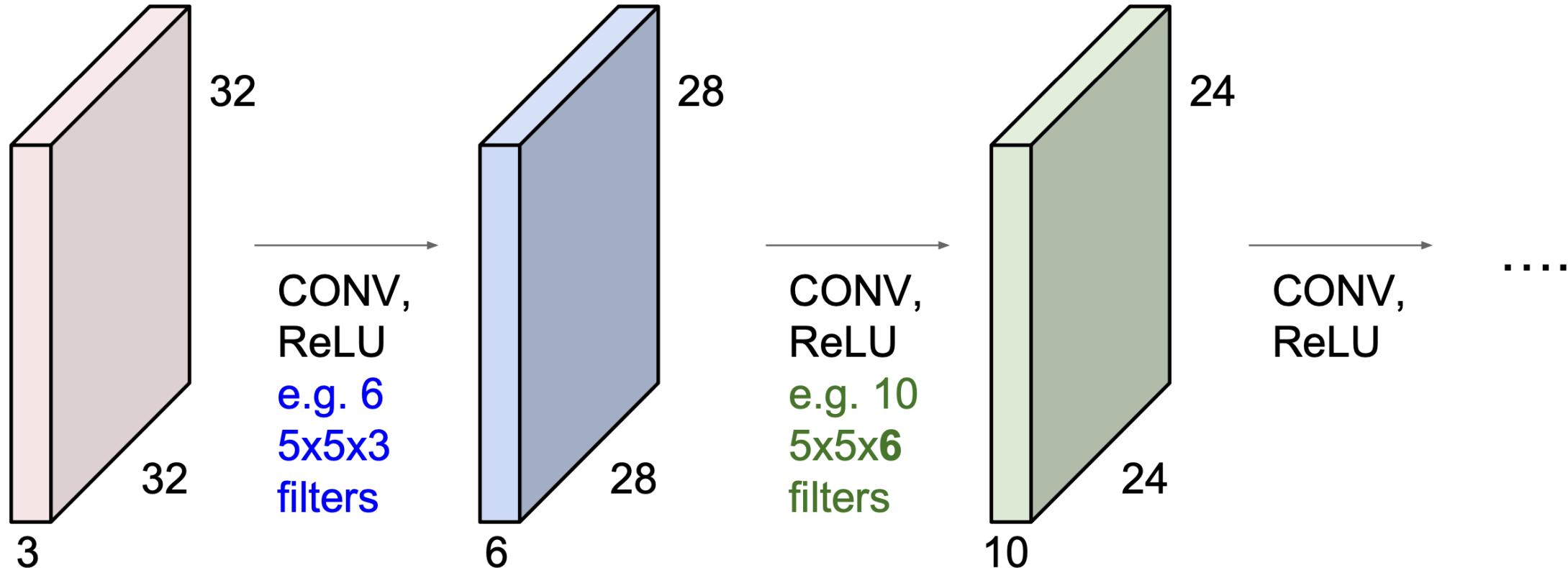


We stack these up to get a “new image” of size 28x28x6!

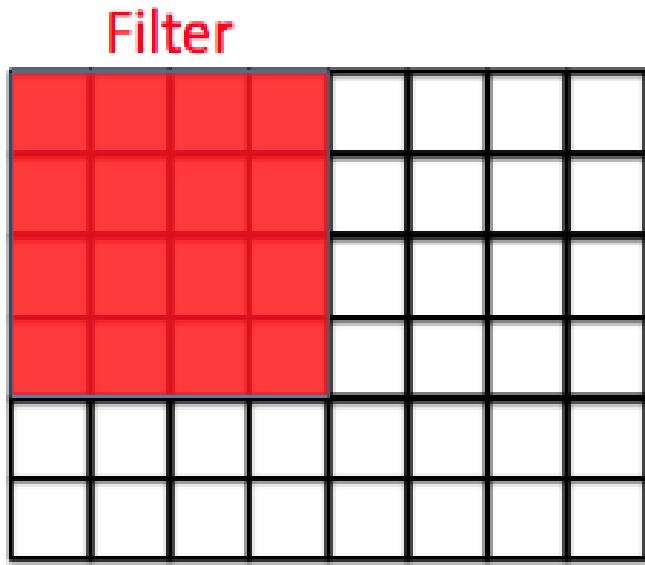
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



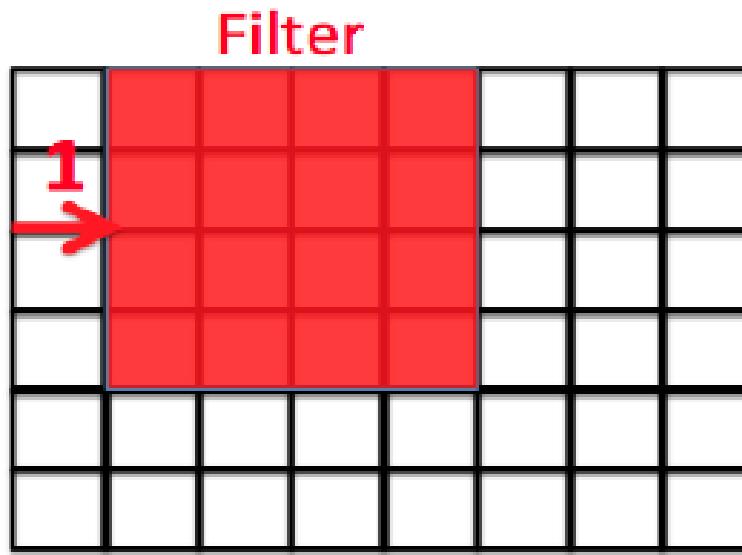
Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



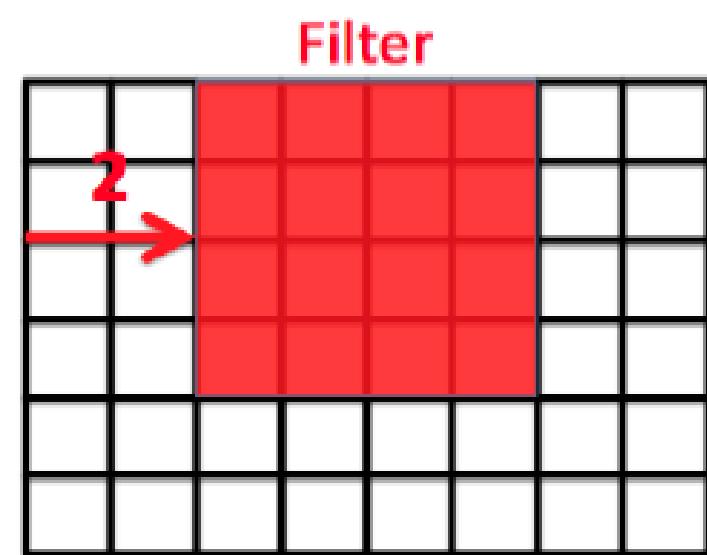
- **Stride** refers to the number of pixels by which the filter/kernel is moved across the input image or matrix when performing the convolution operation. A stride of 1 means the filter moves one pixel at a time. A larger stride results in downsampling the input and leads to a smaller output dimension.
- **Padding** is a technique used to preserve the spatial dimensions of the input volume such that the spatial size of the output volume is the same as the input volume when using a convolutional layer. Padding adds a certain number of rows and columns to the input image/matrix, typically with zeros (zero-padding). This allows the convolution operation to be applied to the bordering elements of the input.
- A **kernel** or **filter** in the context of CNNs is a small matrix used to apply effects such as blurring, sharpening, edge detection, etc. This kernel is convolved over the input matrix or image to produce a feature map indicating the presence of specific features or patterns in the input.
- A **feature map** is the output generated by applying a filter to the input. It represents the responses of the filter throughout the spatial extent of the input. Essentially, it is a mapped output of where certain features are detected in the input.
- **Convolution** is the operation of adding each element of the image to its local neighbors, weighted by the kernel. This operation blends the pixels within a local region of the image to produce a single average output of that region.
- **Pooling** is a technique used to reduce the dimensions of the feature maps. It reduces the computational complexity and the number of parameters. The most common pooling operation is max pooling, which selects the maximum element from the region of the feature map covered by the filter. There is also average pooling, which calculates the average of the elements in that region.



Image



Image



Image

Convolution operation

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

- The kernel, a yellow matrix, slides horizontally and then vertically across the original green matrix (image) to detect features by calculating the dot product sum of image and kernel values, outputting a new matrix.
- This kernel, with initially random values, is a trainable parameter.

Producing Feature Maps



Original



Sharpen



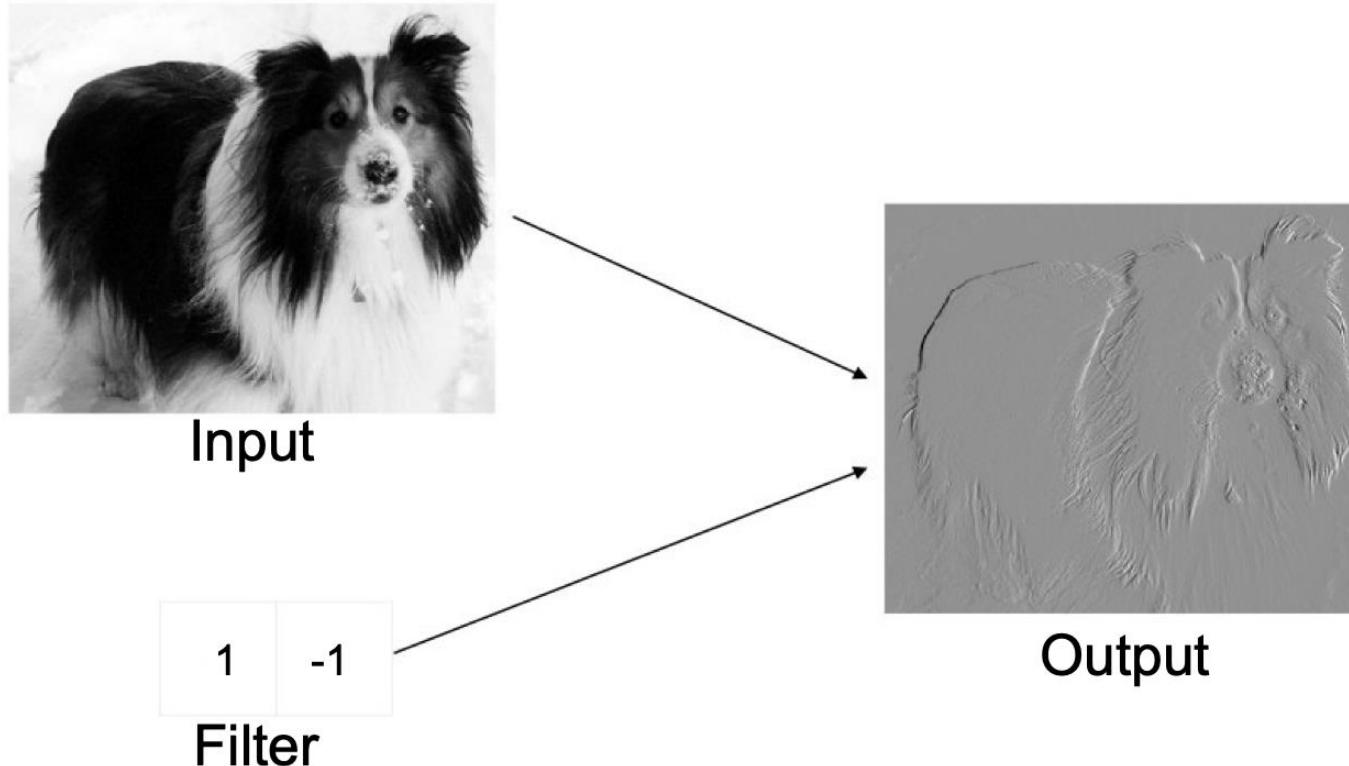
Edge Detect



“Strong” Edge
Detect

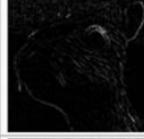
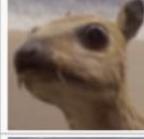
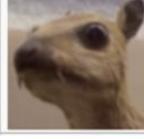
A simple pattern: Edges

How can we detect edges with a kernel?

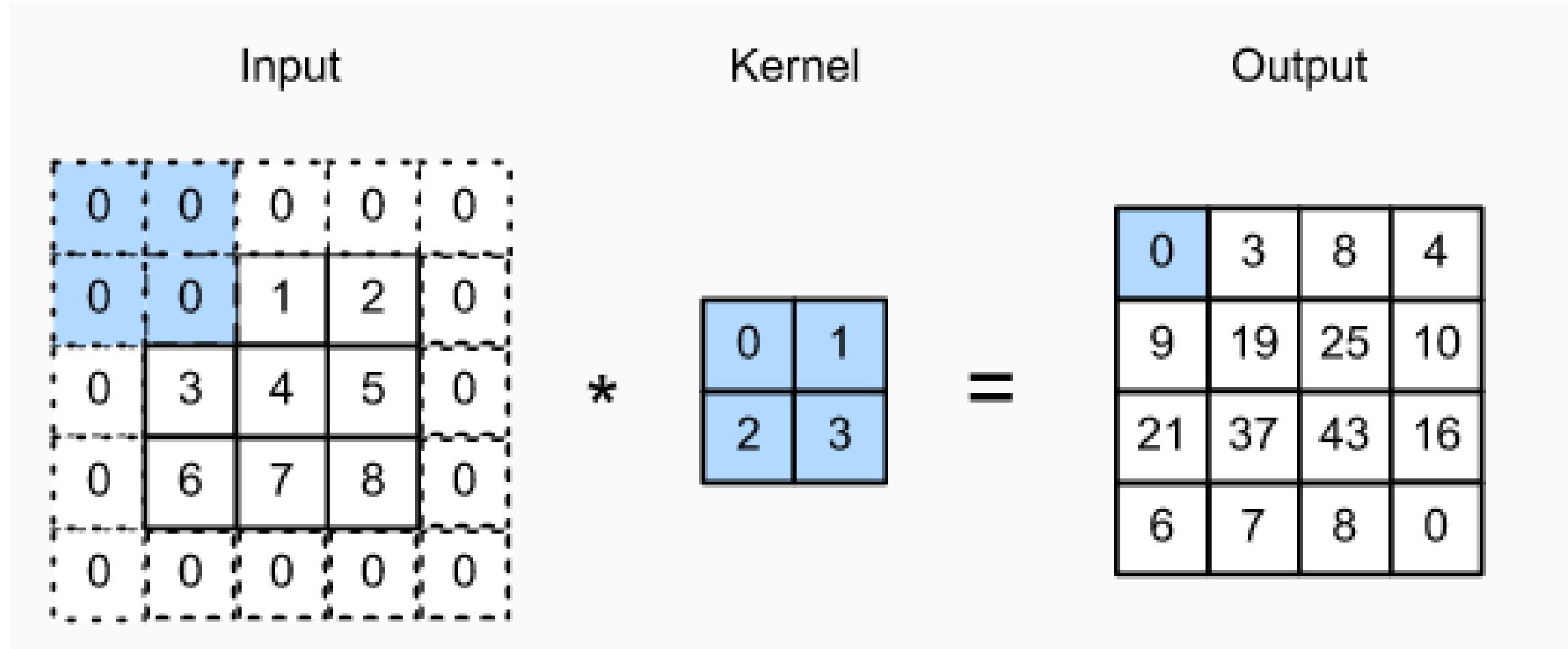


(Goodfellow 2016)

Simple Kernels / Filters

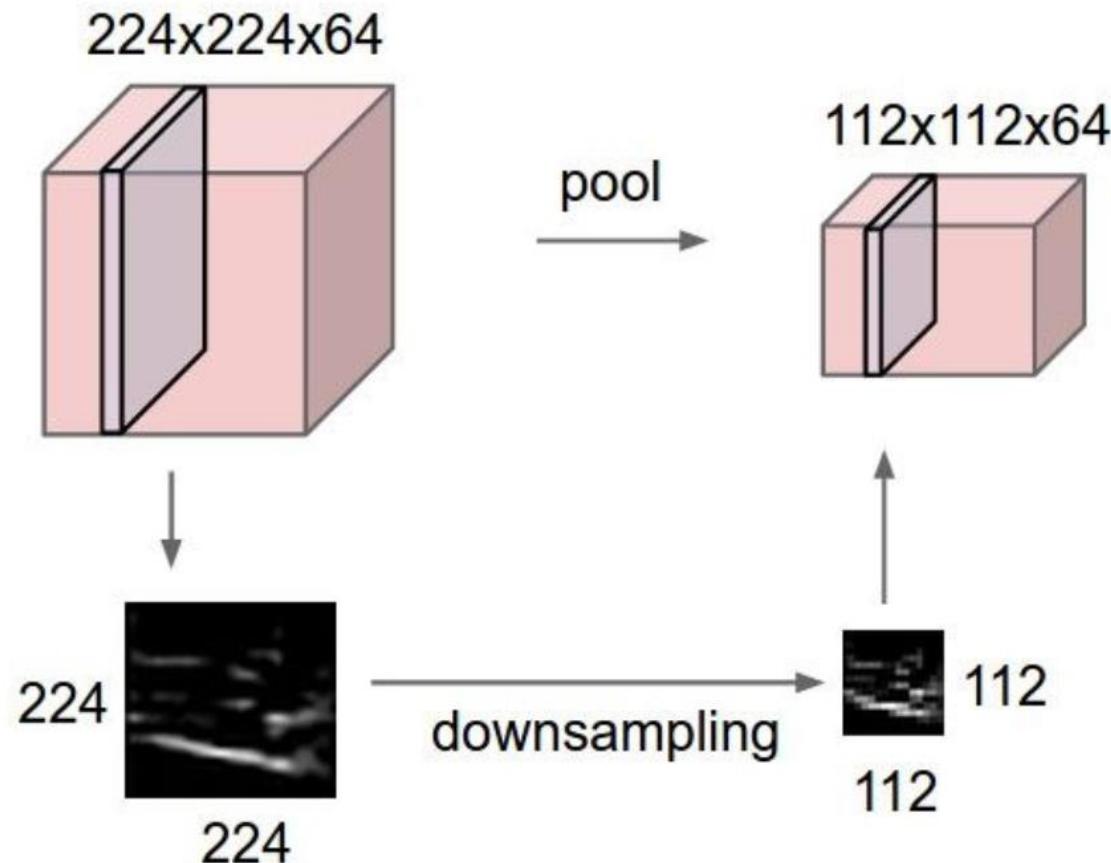
Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

Padding



Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:

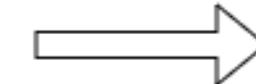


Max Pooling - Selects the brighter pixels from the image.

4	3	1	5
1	3	4	8
4	5	4	3
6	5	9	4

$$\text{Max}([4, 3, 1, 3]) = 4$$

4	3	1	5
1	3	4	8
4	5	4	3
6	5	9	4



4	8
6	9

A pooling layer is another building block of a CNN. Its function is to progressively reduce the spatial size of the representation to reduce the network complexity and computational cost.

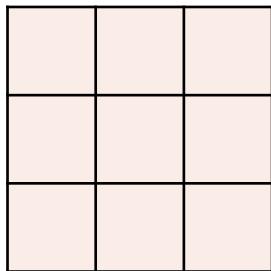
Average Pooling - retains much information about the “less important” elements of a block, or pool.



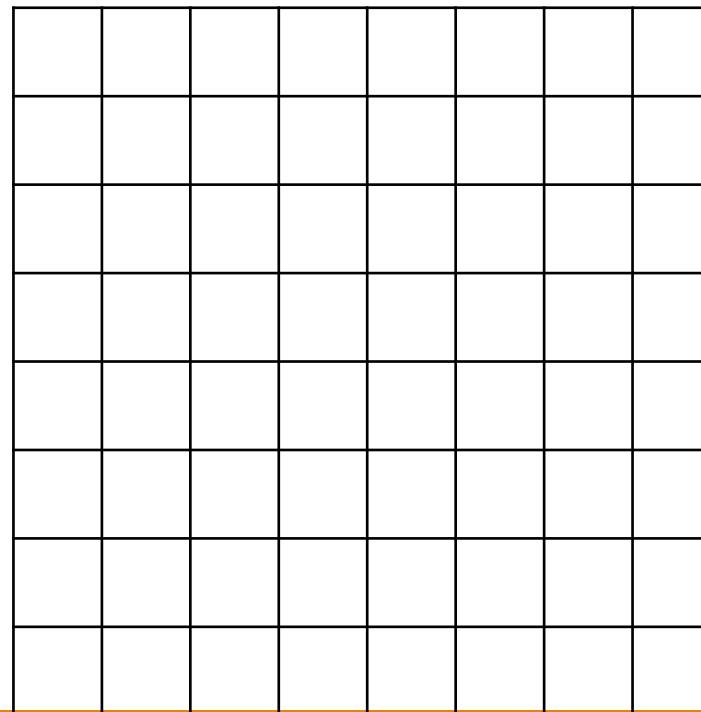
A pooling layer is another building block of a CNN. Its function is to progressively reduce the spatial size of the representation to reduce the network complexity and computational cost.

Convolutional Layers

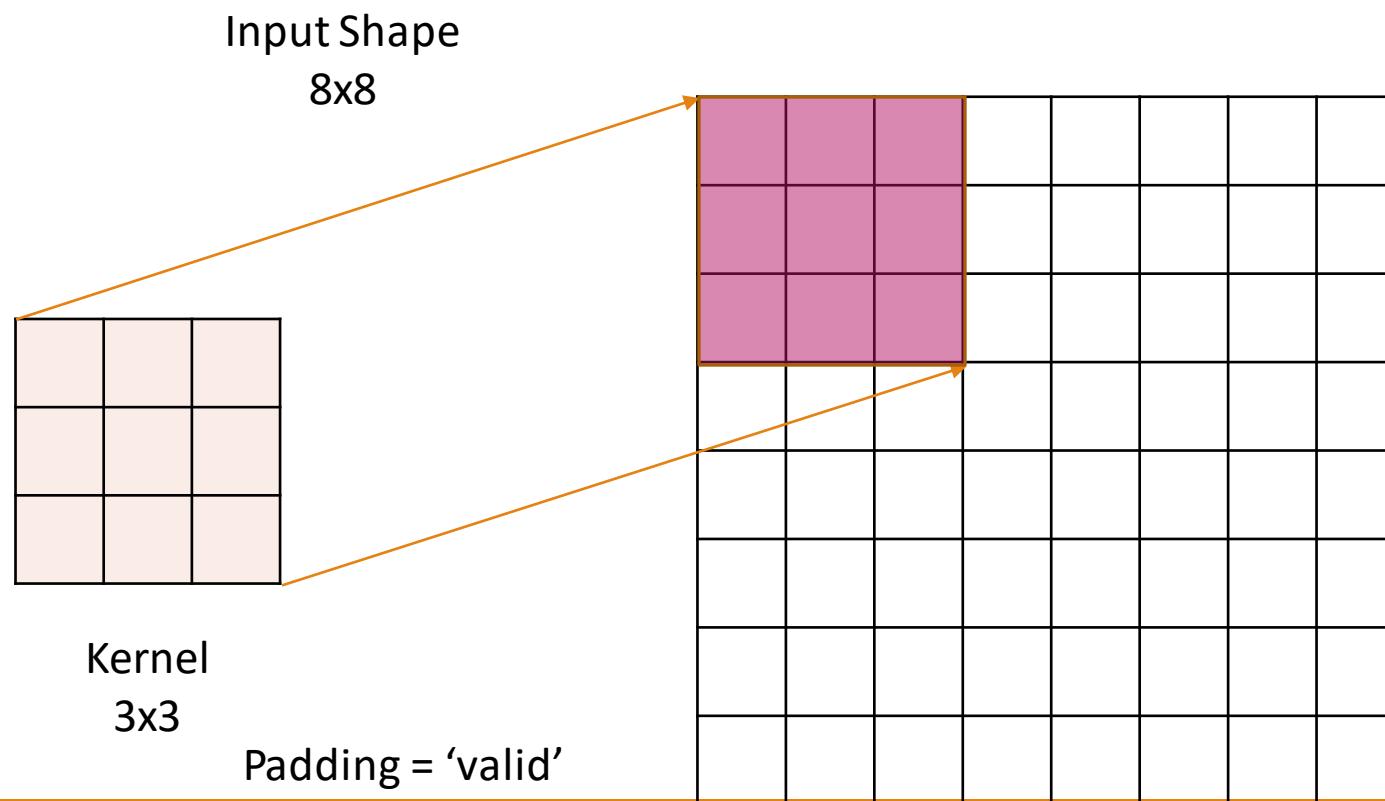
Input Shape
8x8



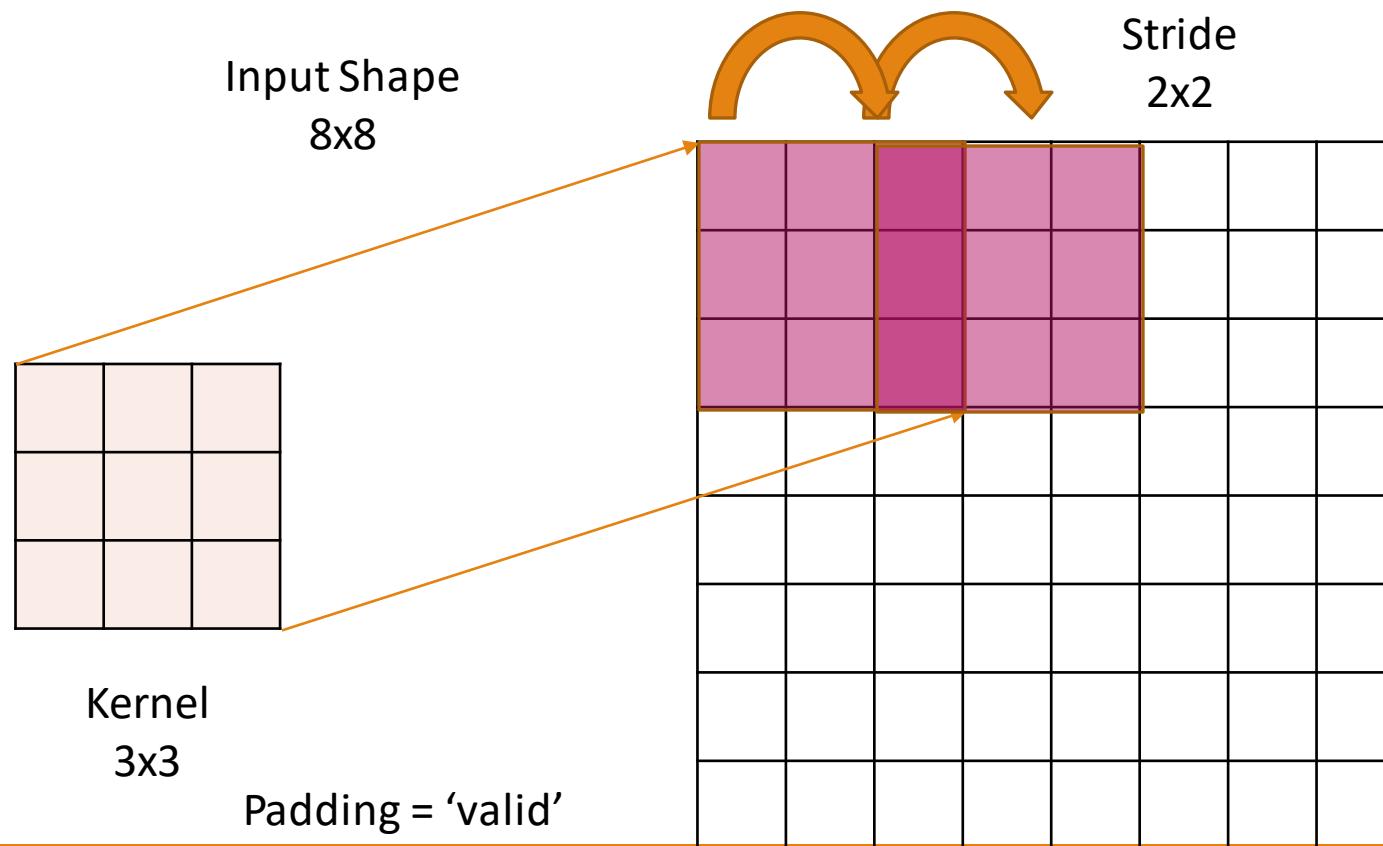
Kernel
3x3
Padding = 'valid'



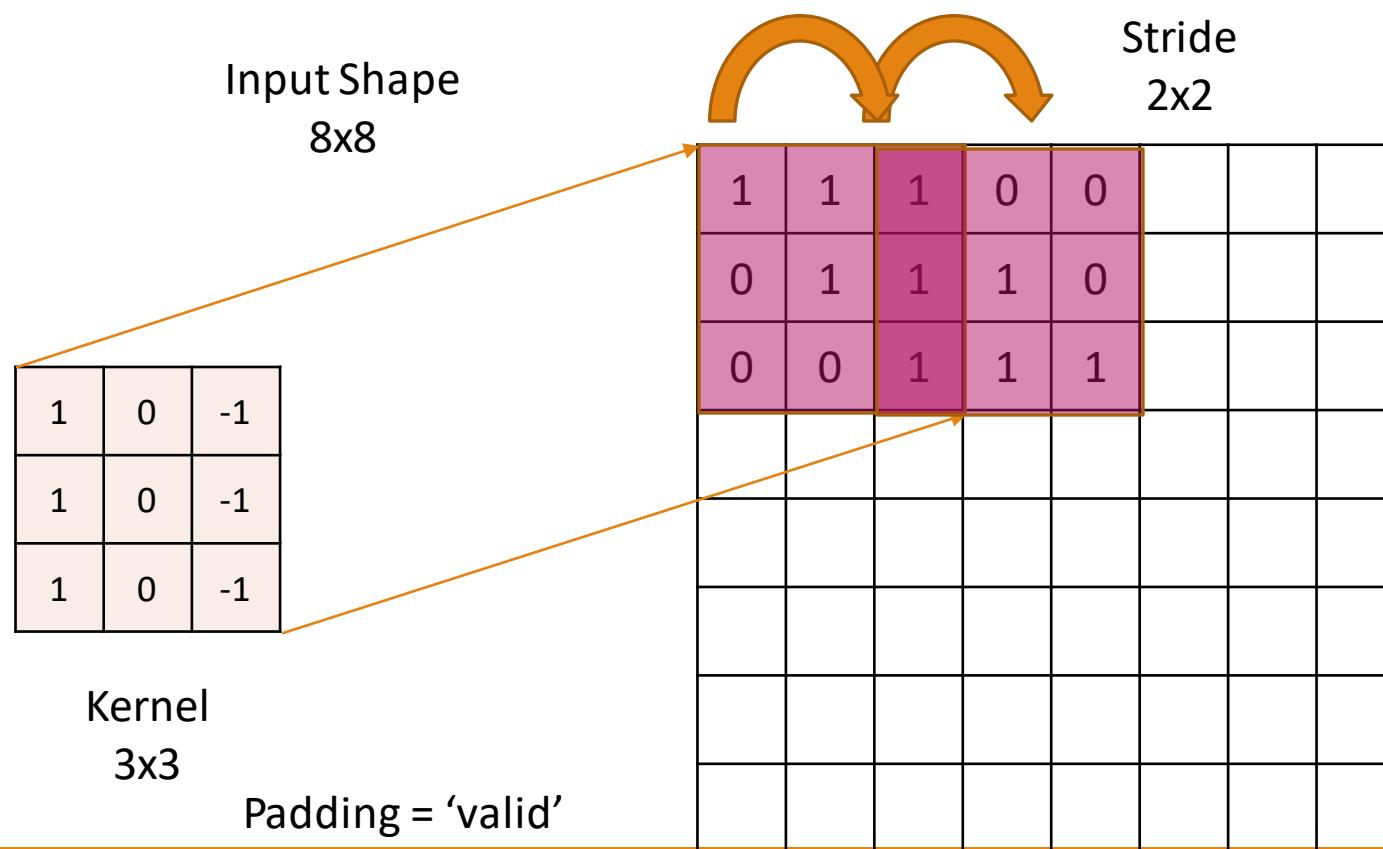
Convolutional Layers



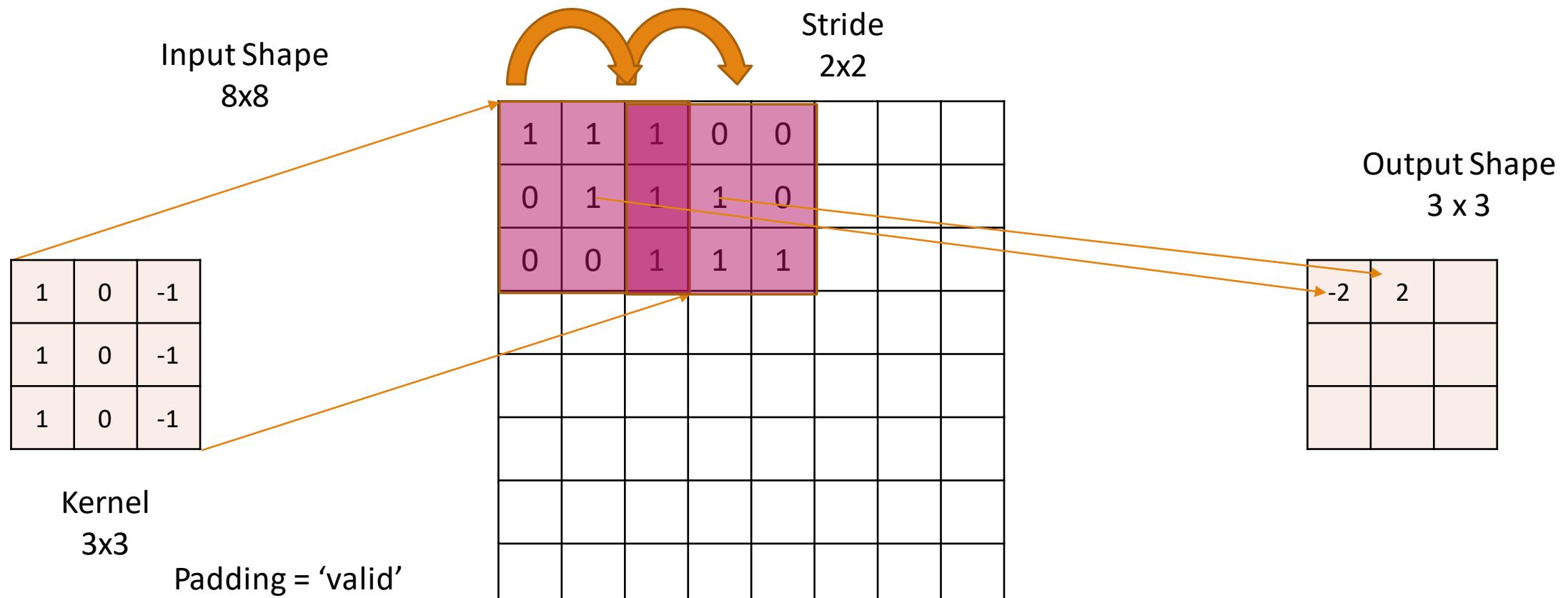
Convolutional Layers



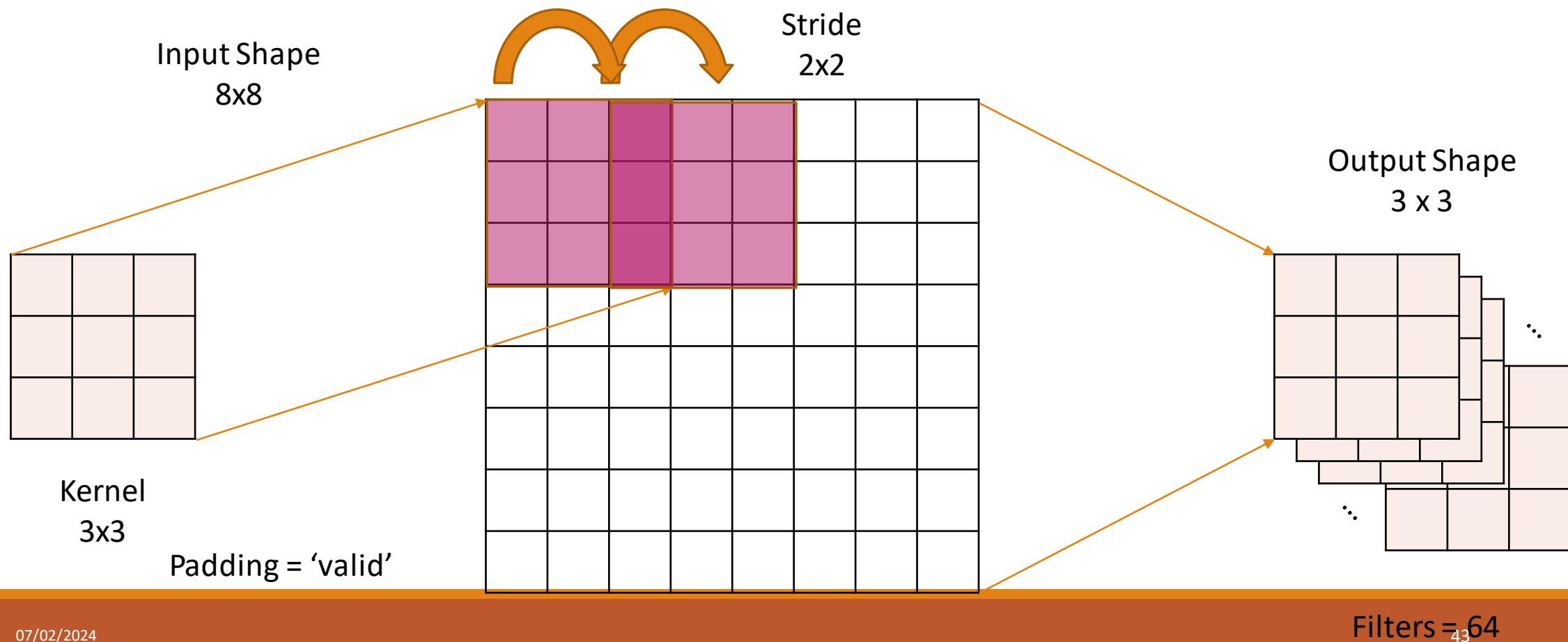
Convolutional Layers

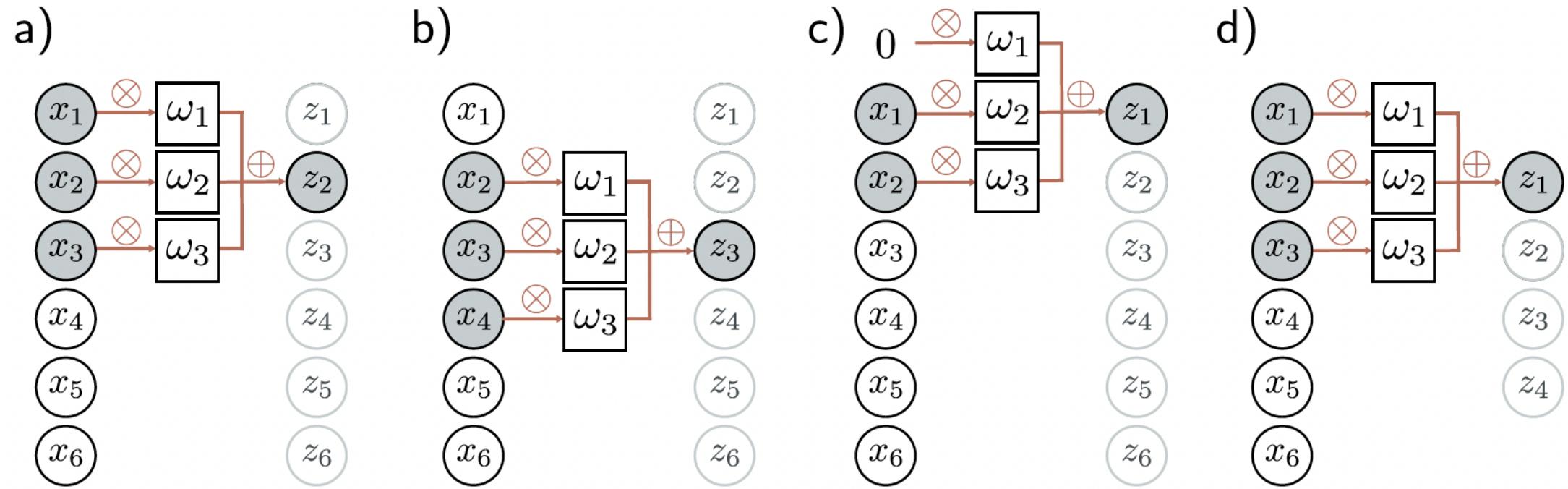


Convolutional Layers

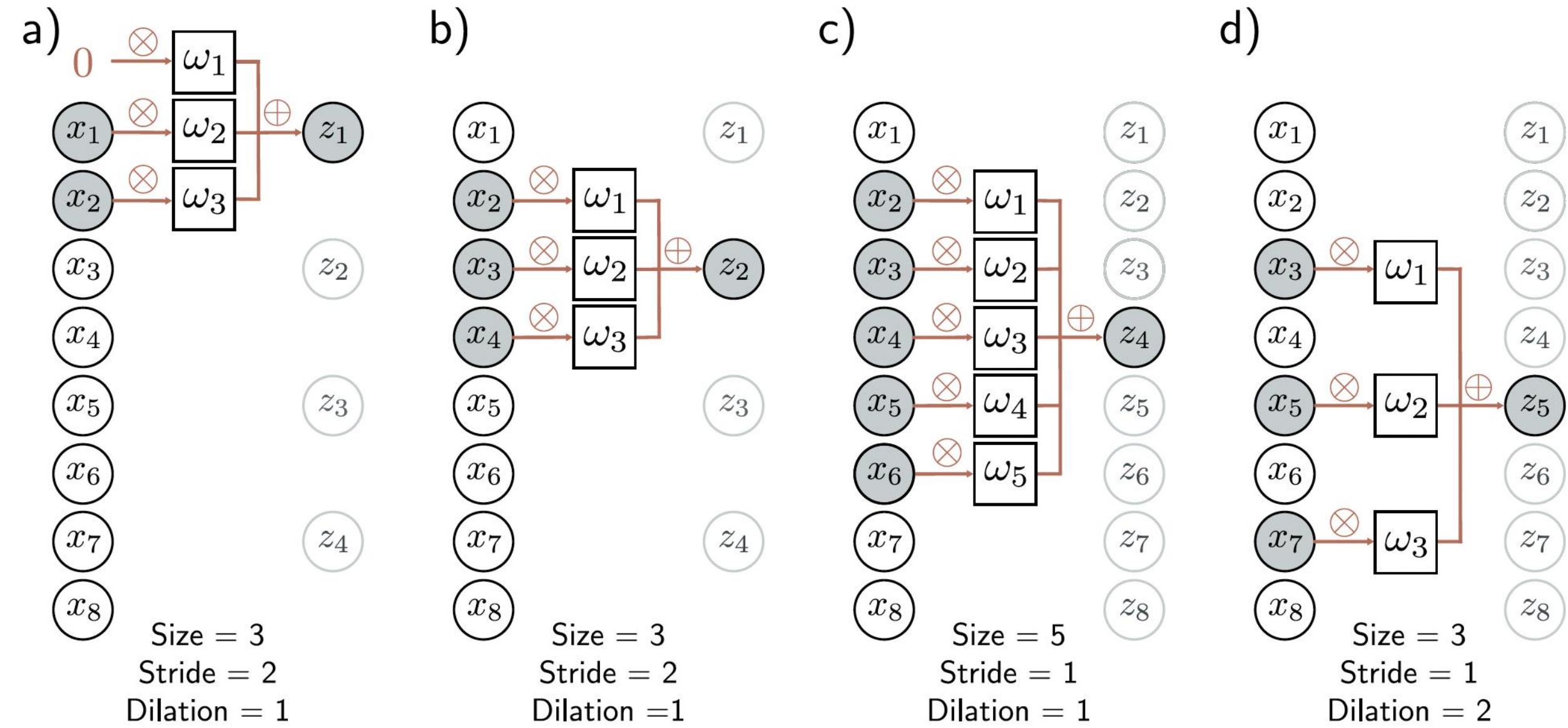


Convolutional Layers





$$z_i = \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1},$$

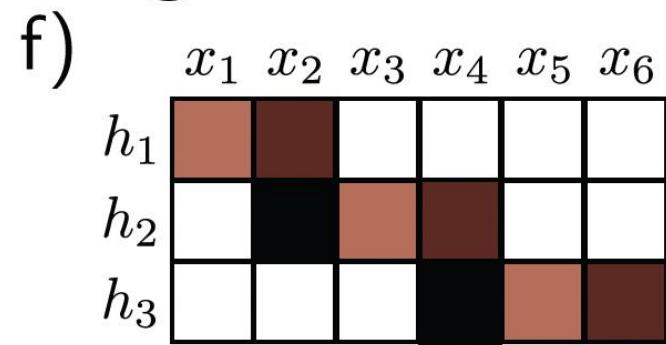
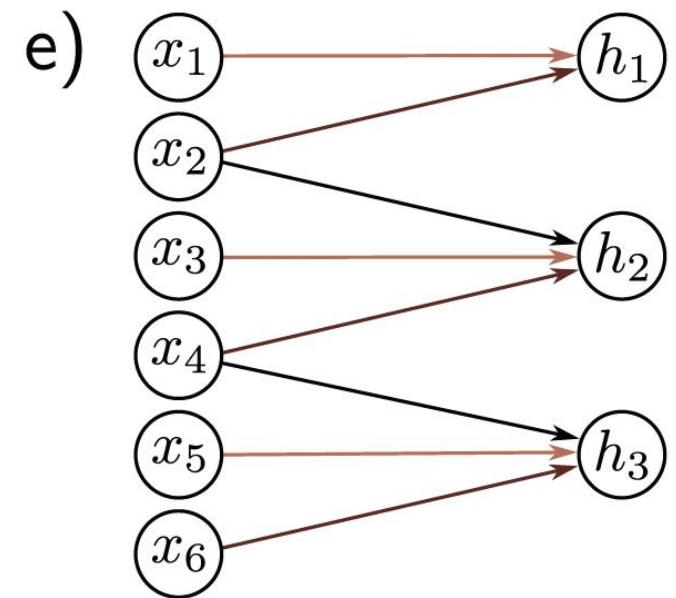
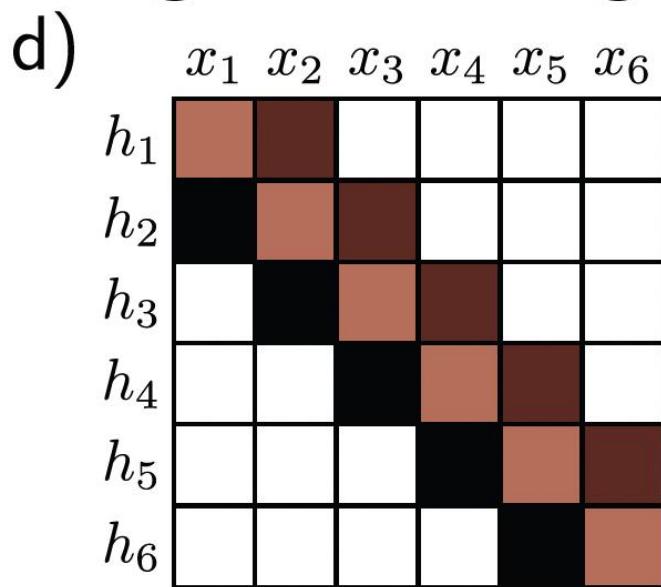
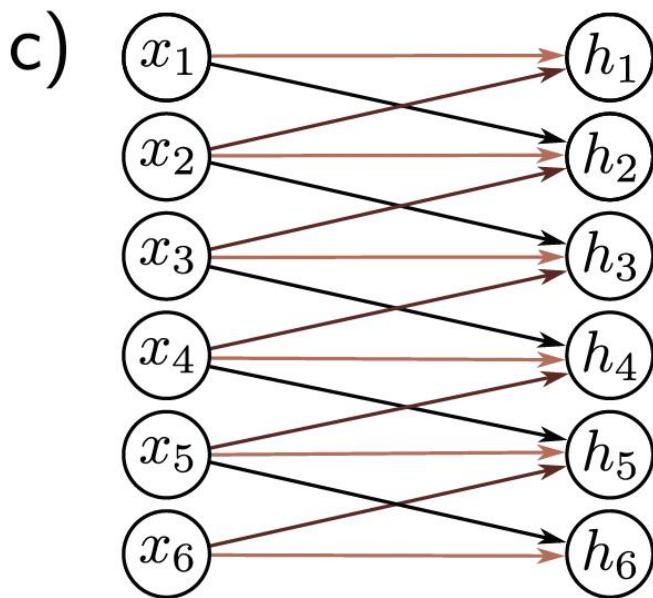
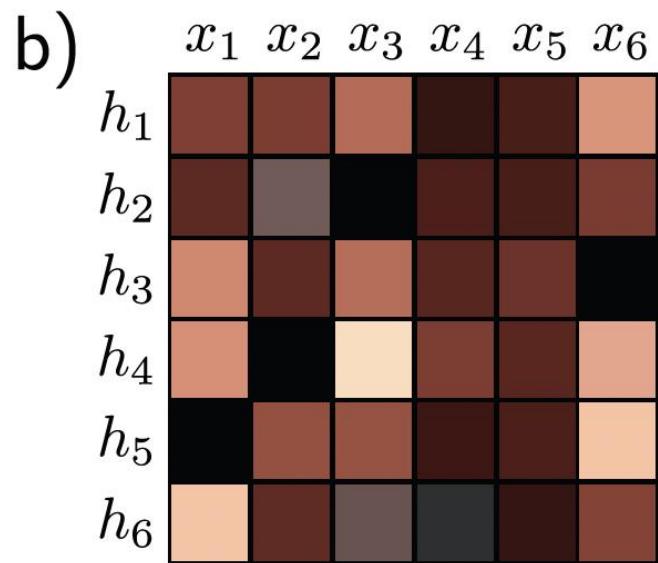
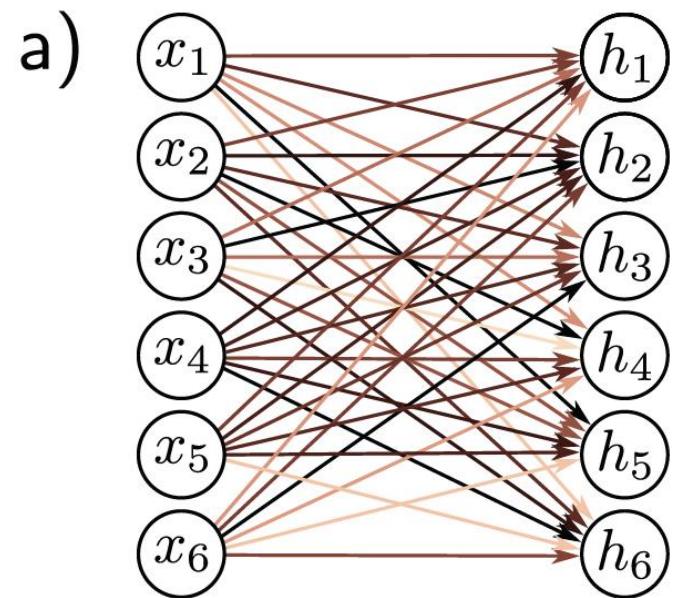


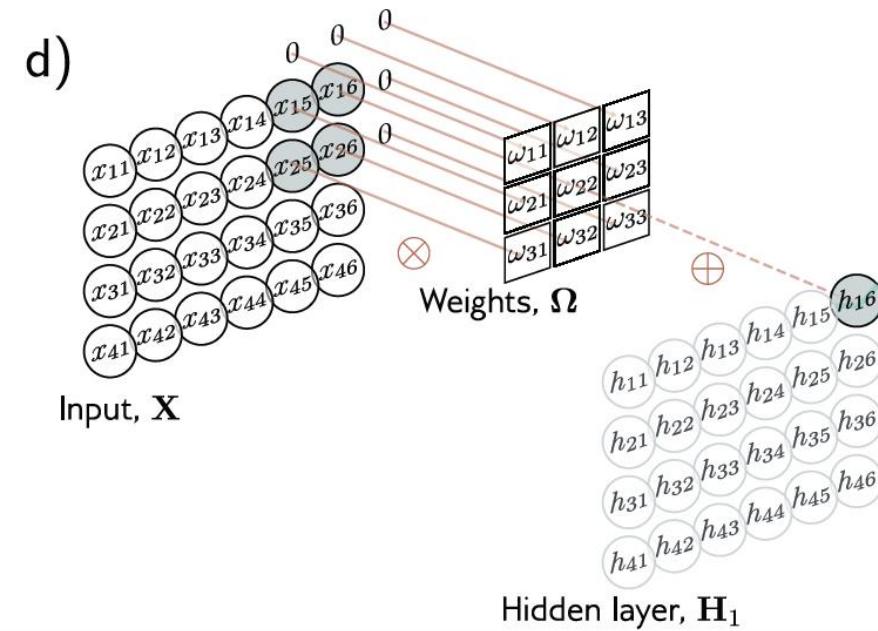
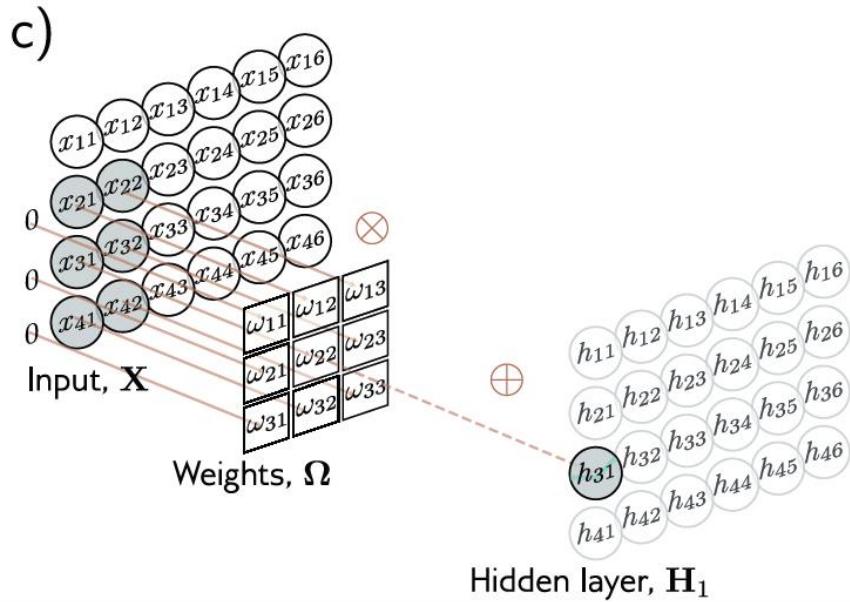
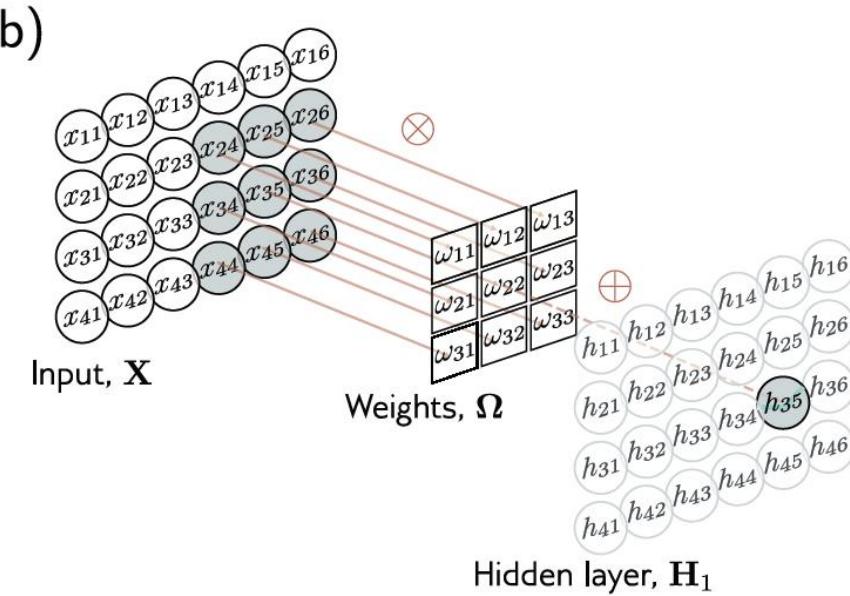
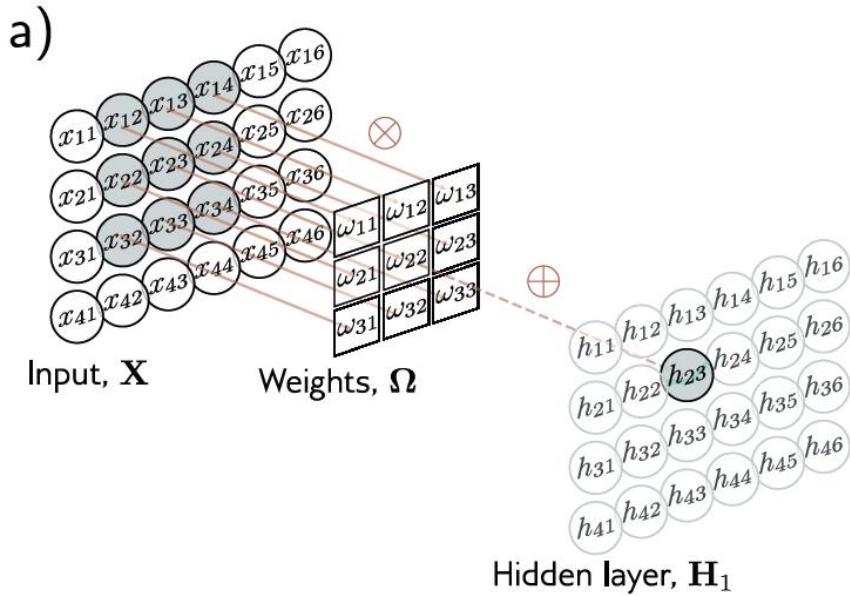
Convolutional layers

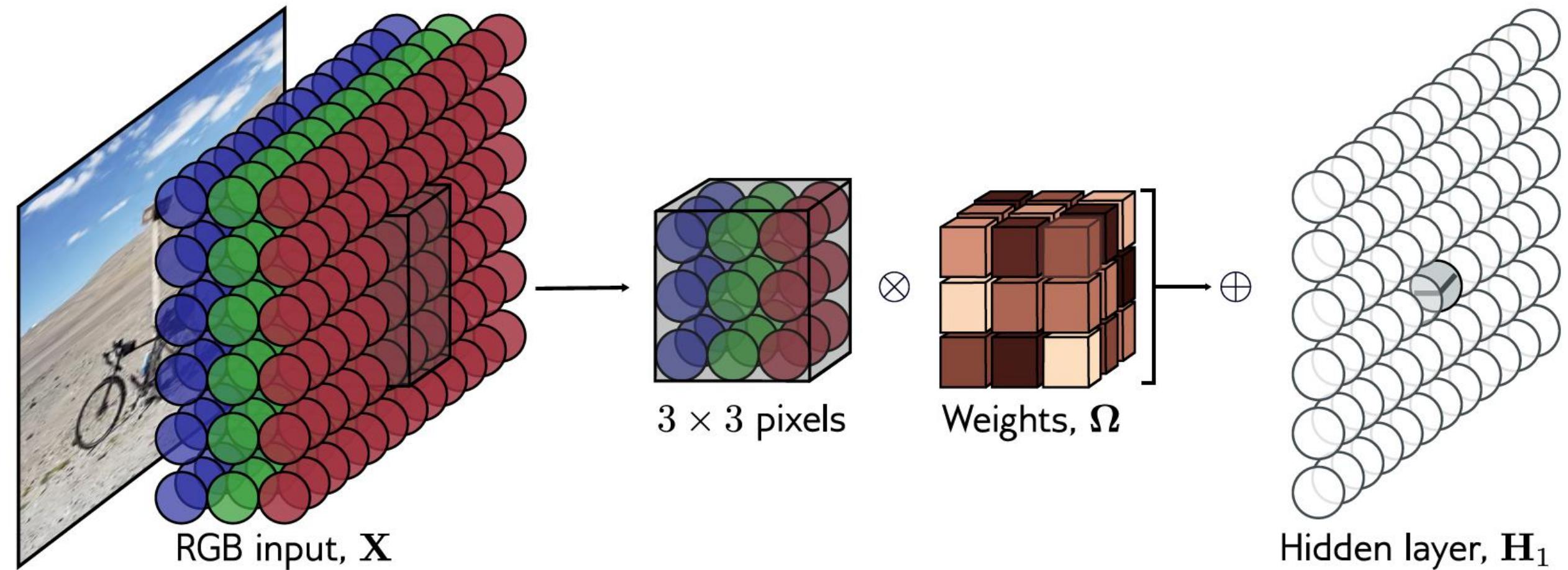
$$\begin{aligned} h_i &= a [\beta + \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}] \\ &= a \left[\beta + \sum_{j=1}^3 \omega_j x_{i+j-2} \right], \end{aligned}$$

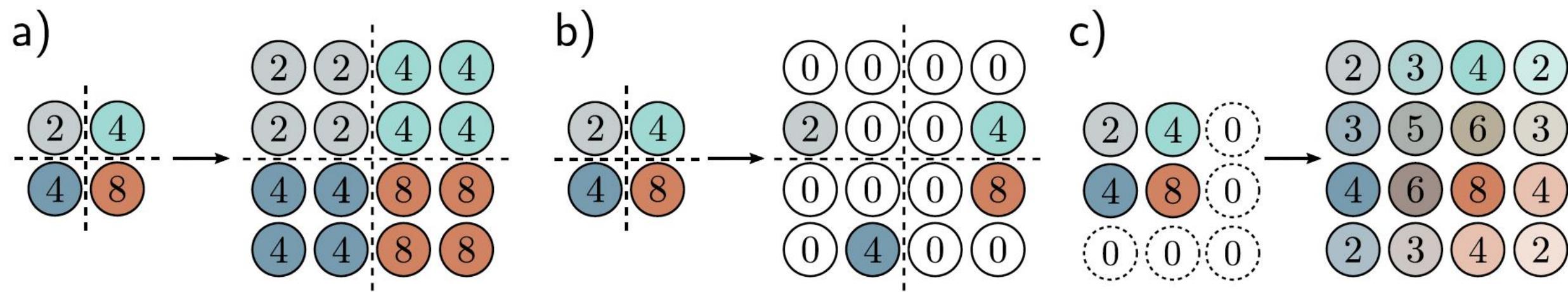
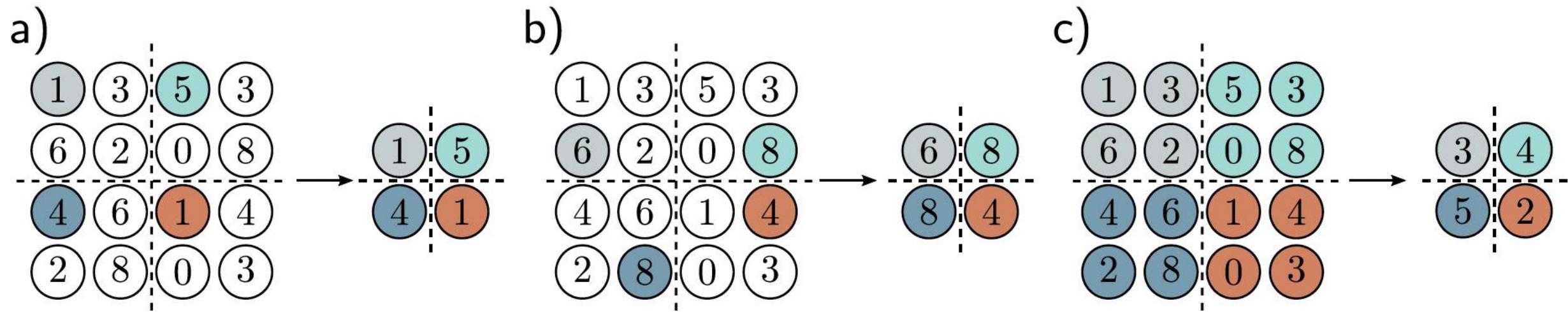
fully connected layer

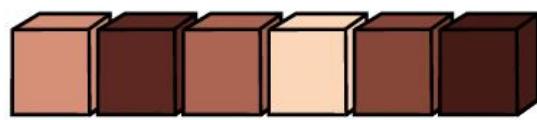
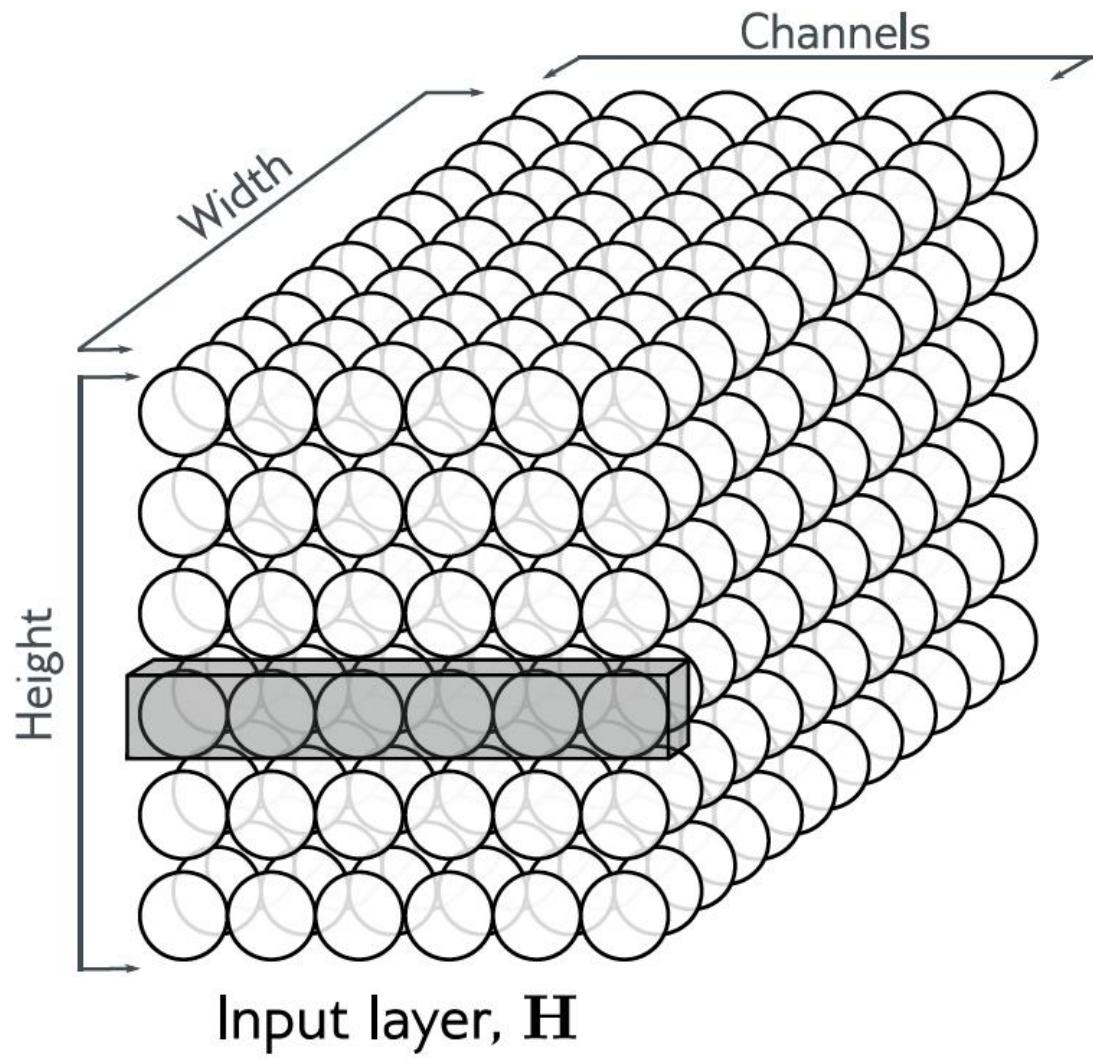
$$h_i = a \left[\beta_i + \sum_{j=1}^D \omega_{ij} x_j \right].$$



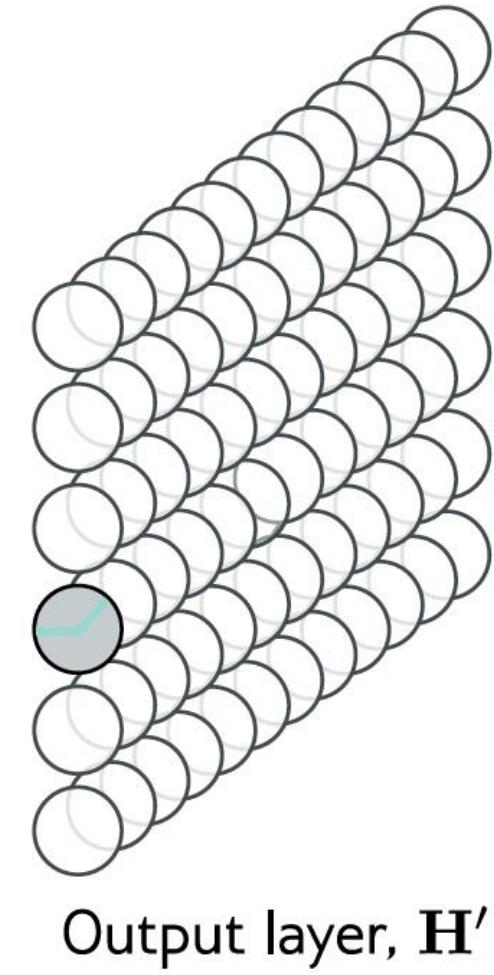


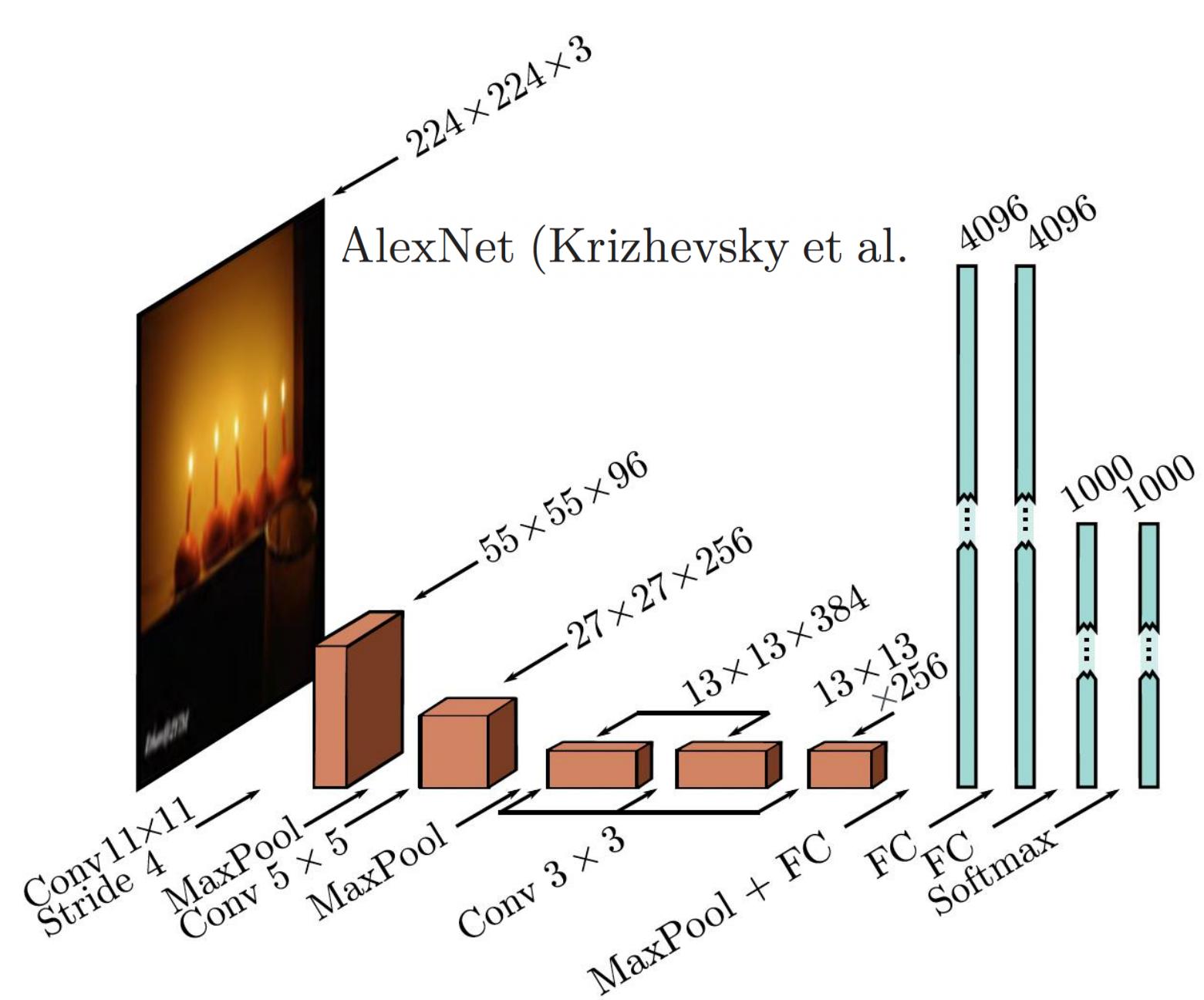




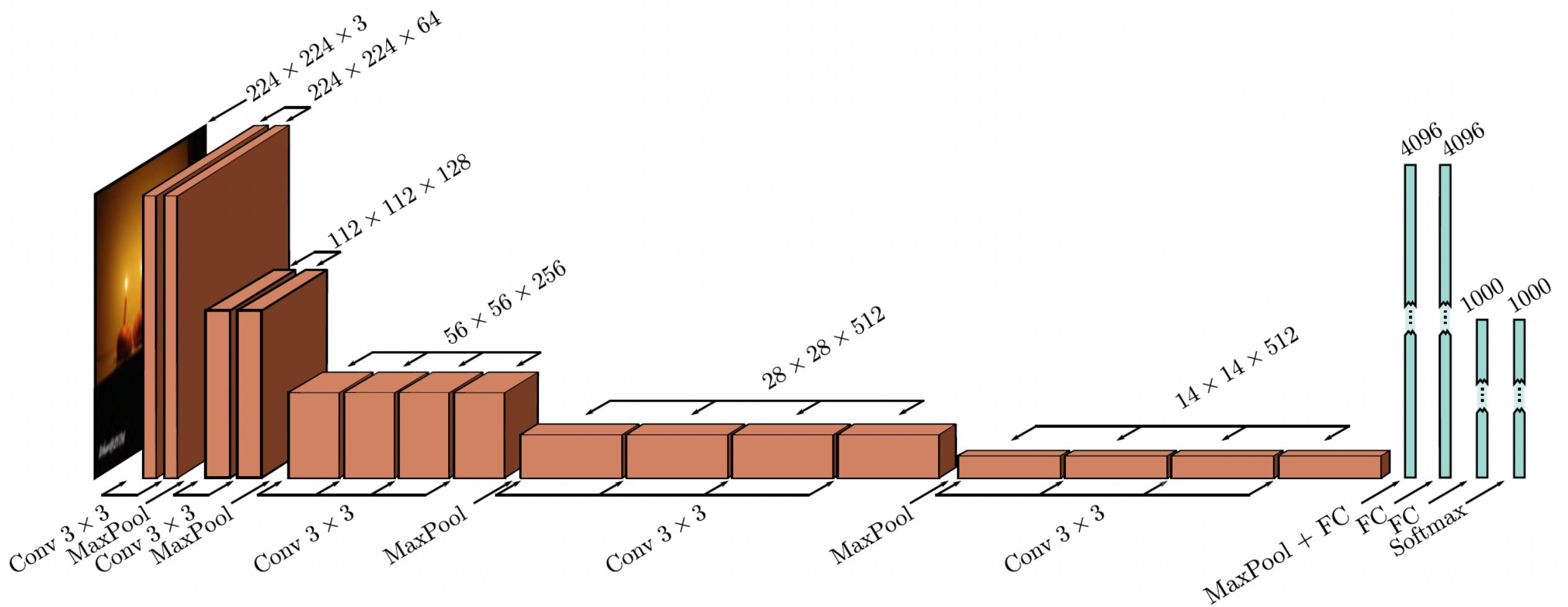


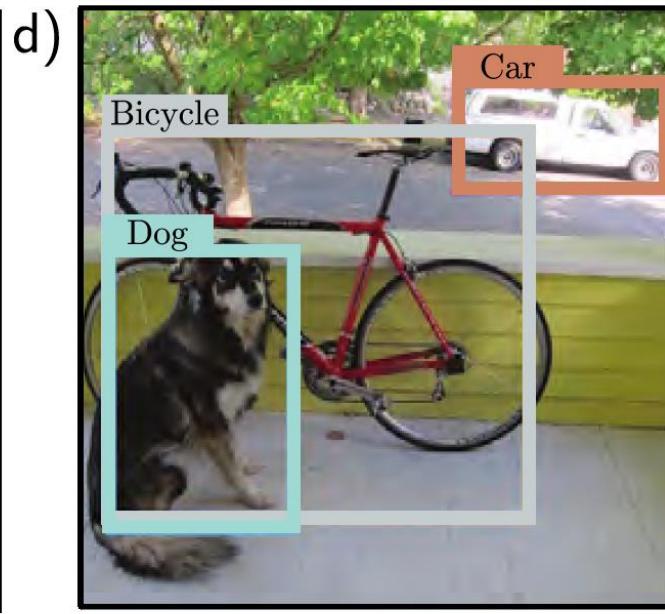
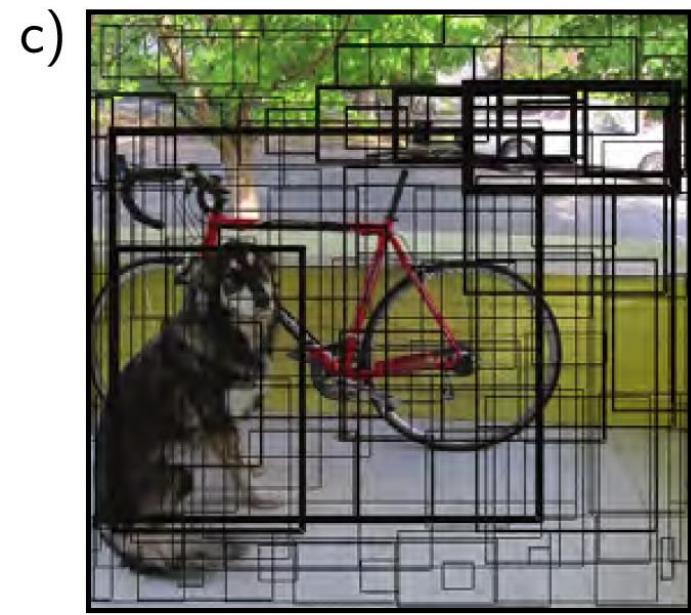
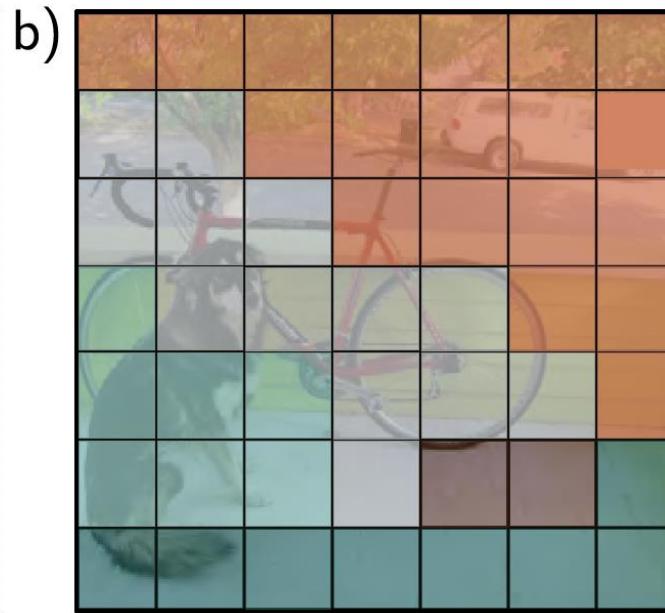
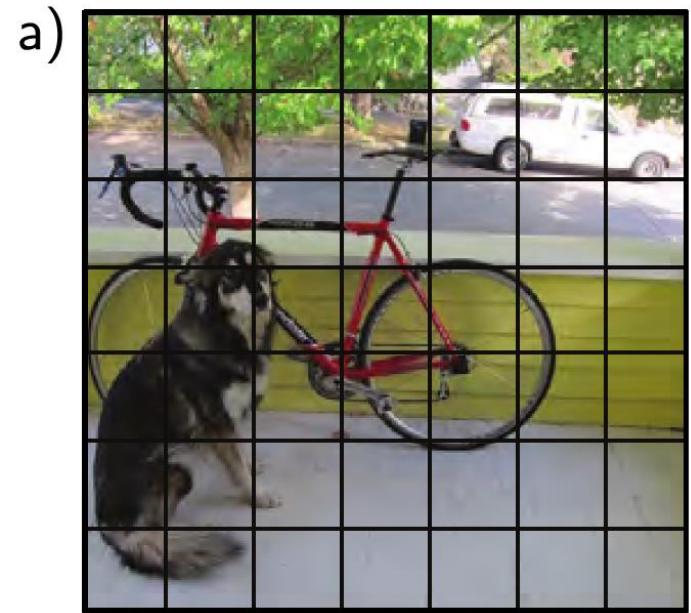
Weights, Ω





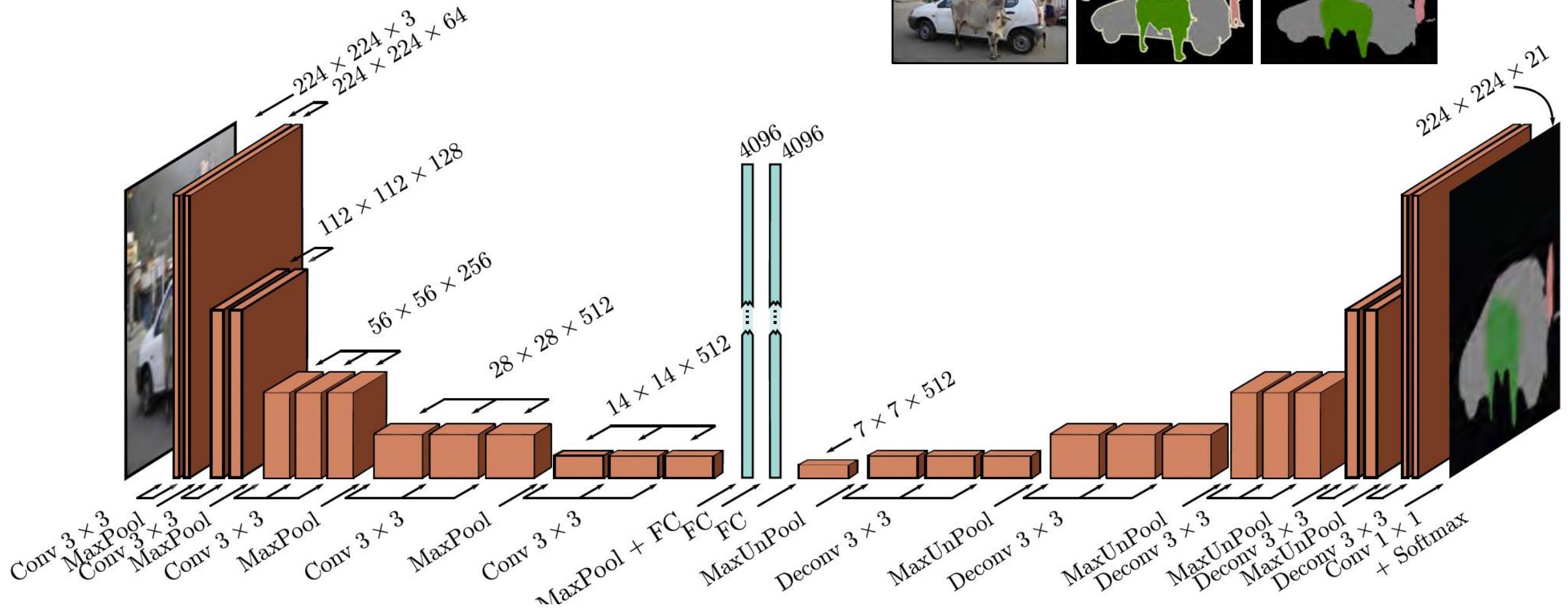
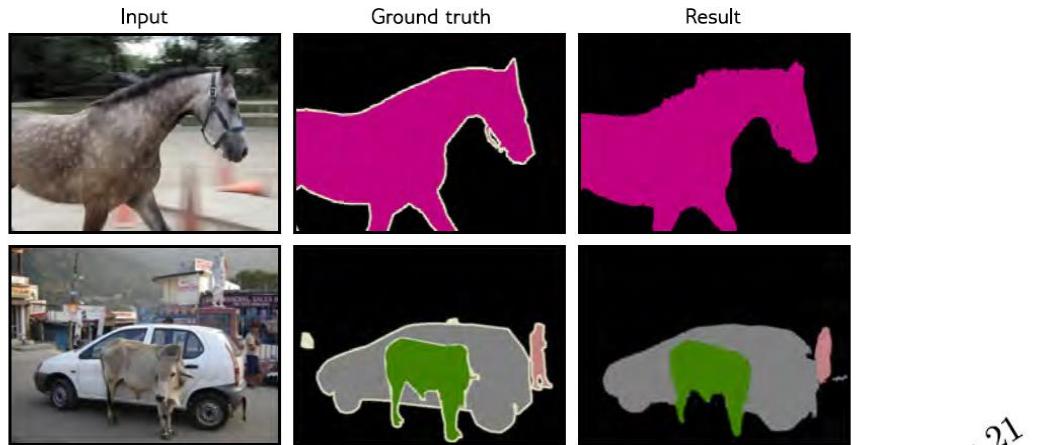
VGG network (Simonyan & Zisserman, 2014)

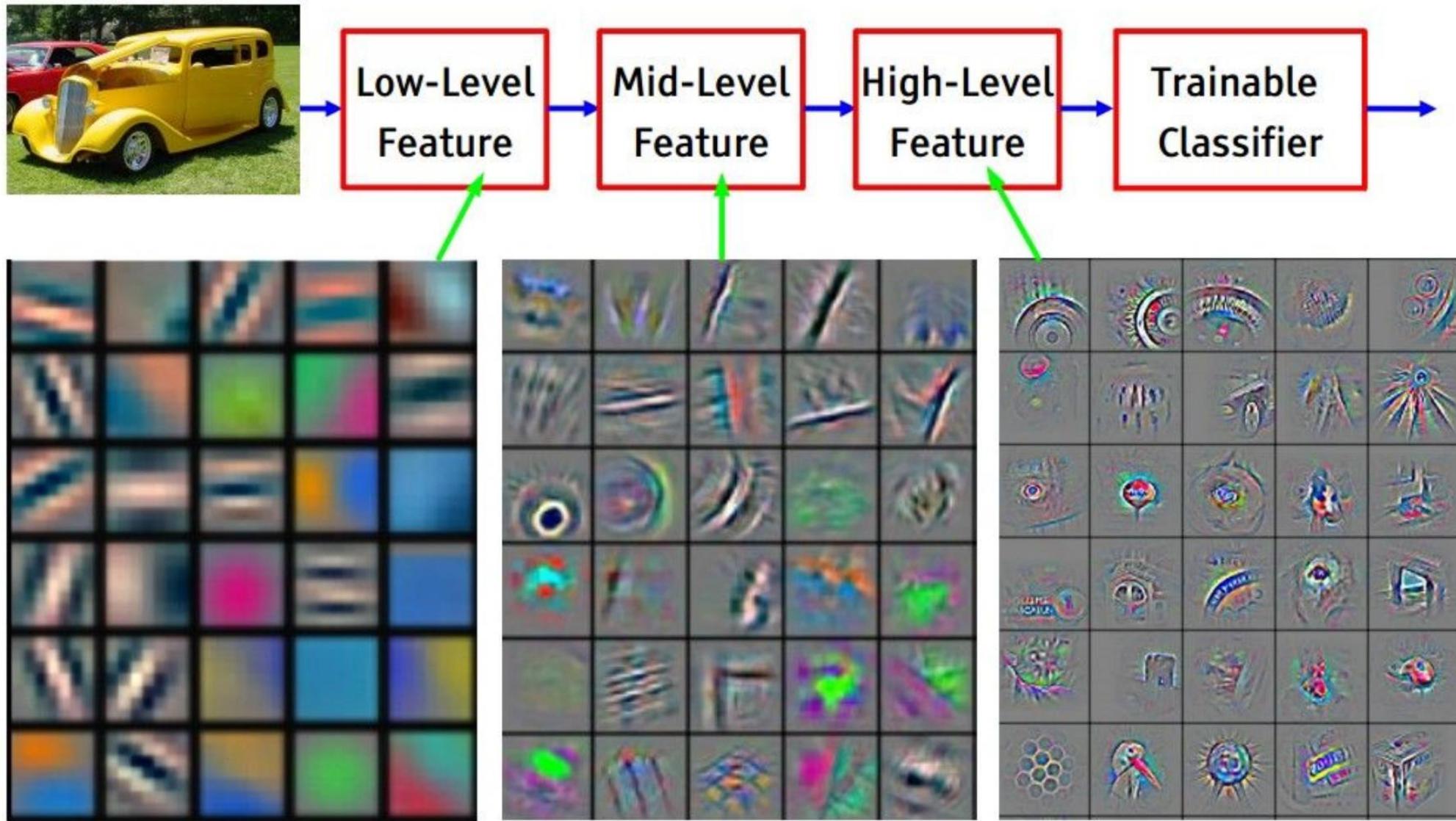




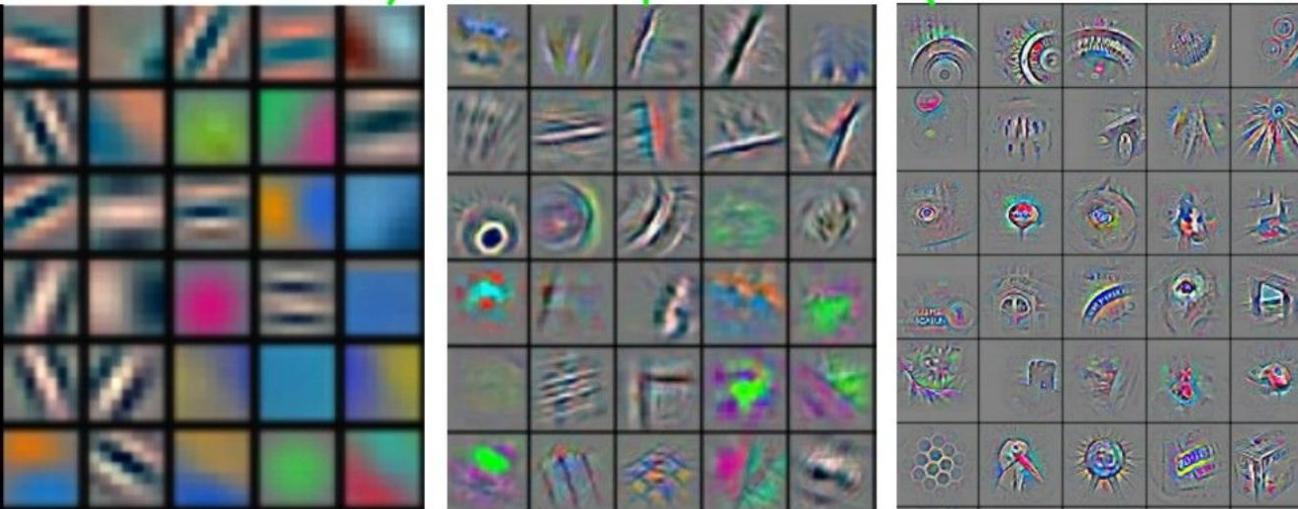
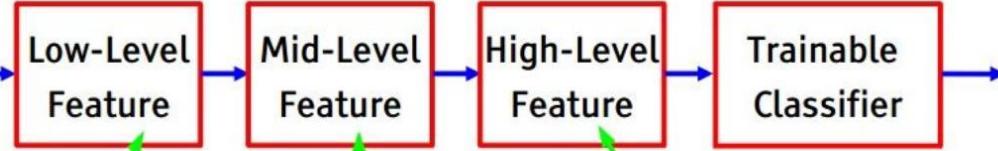
YOLO object detection

Semantic segmentation



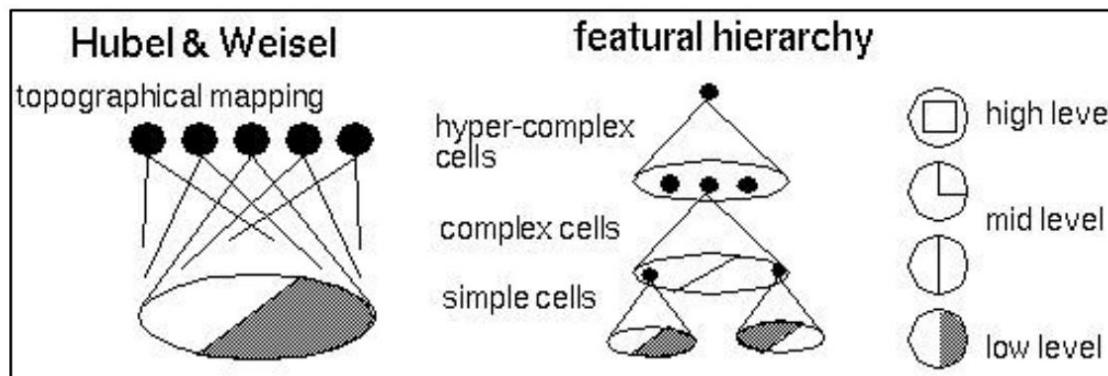


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



[From recent Yann LeCun slides]

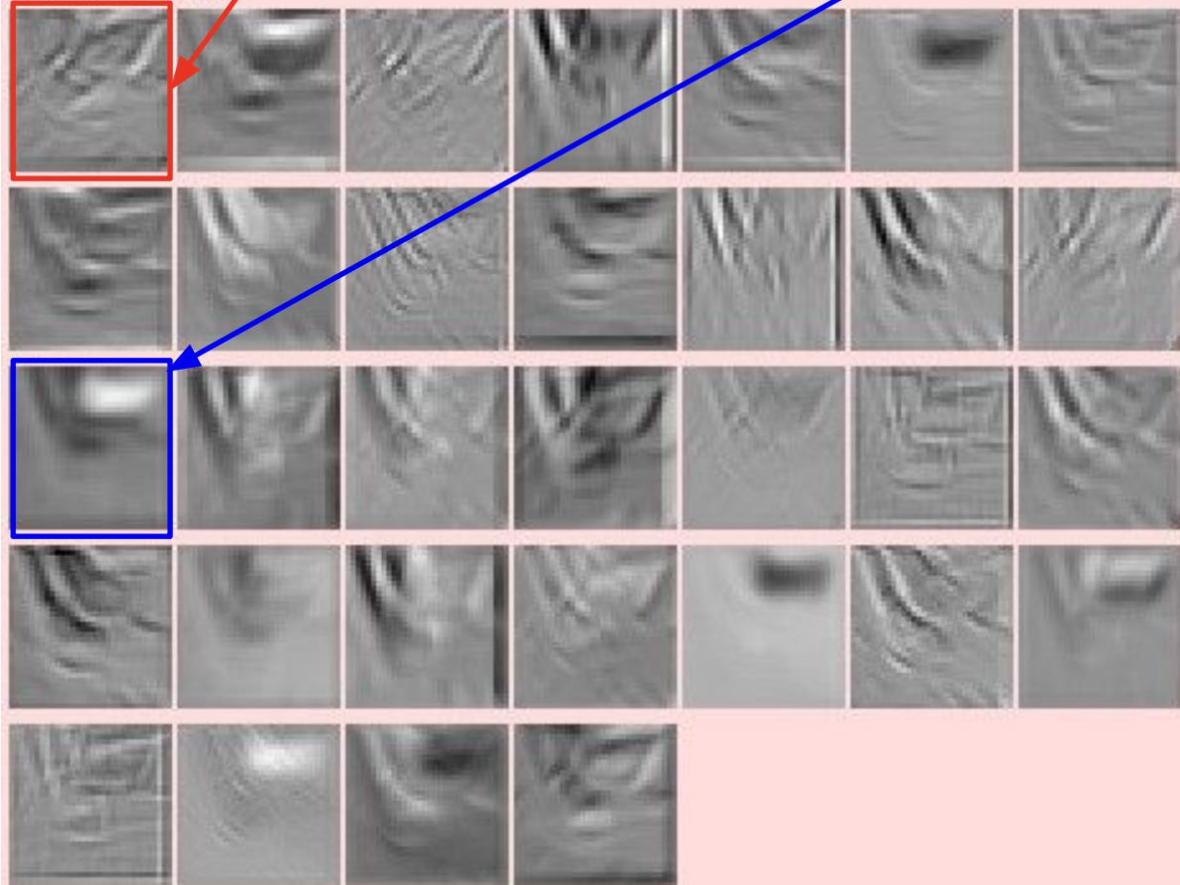
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]





one filter =>
one activation map

Activations:



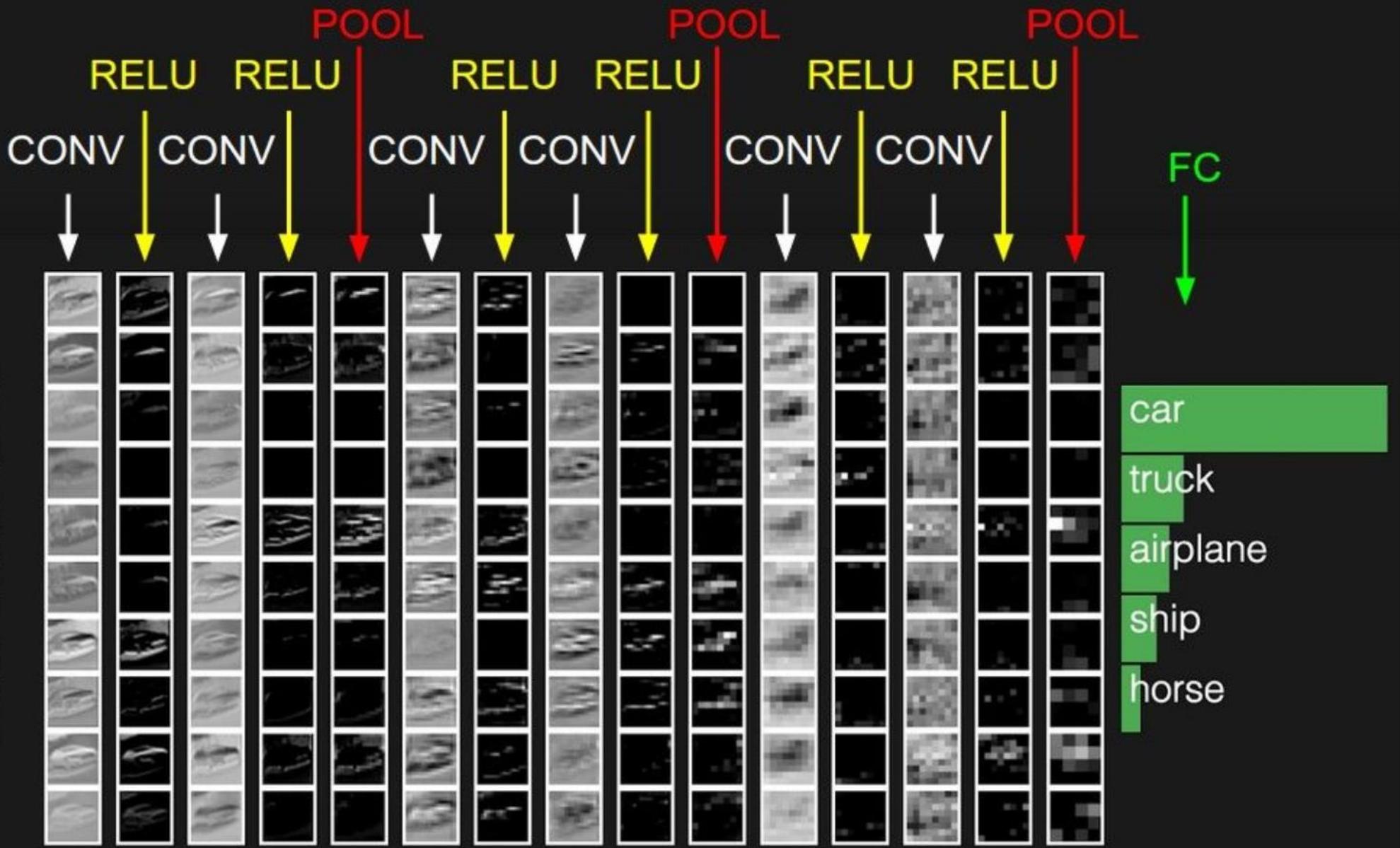
example 5x5 filters
(32 total)

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

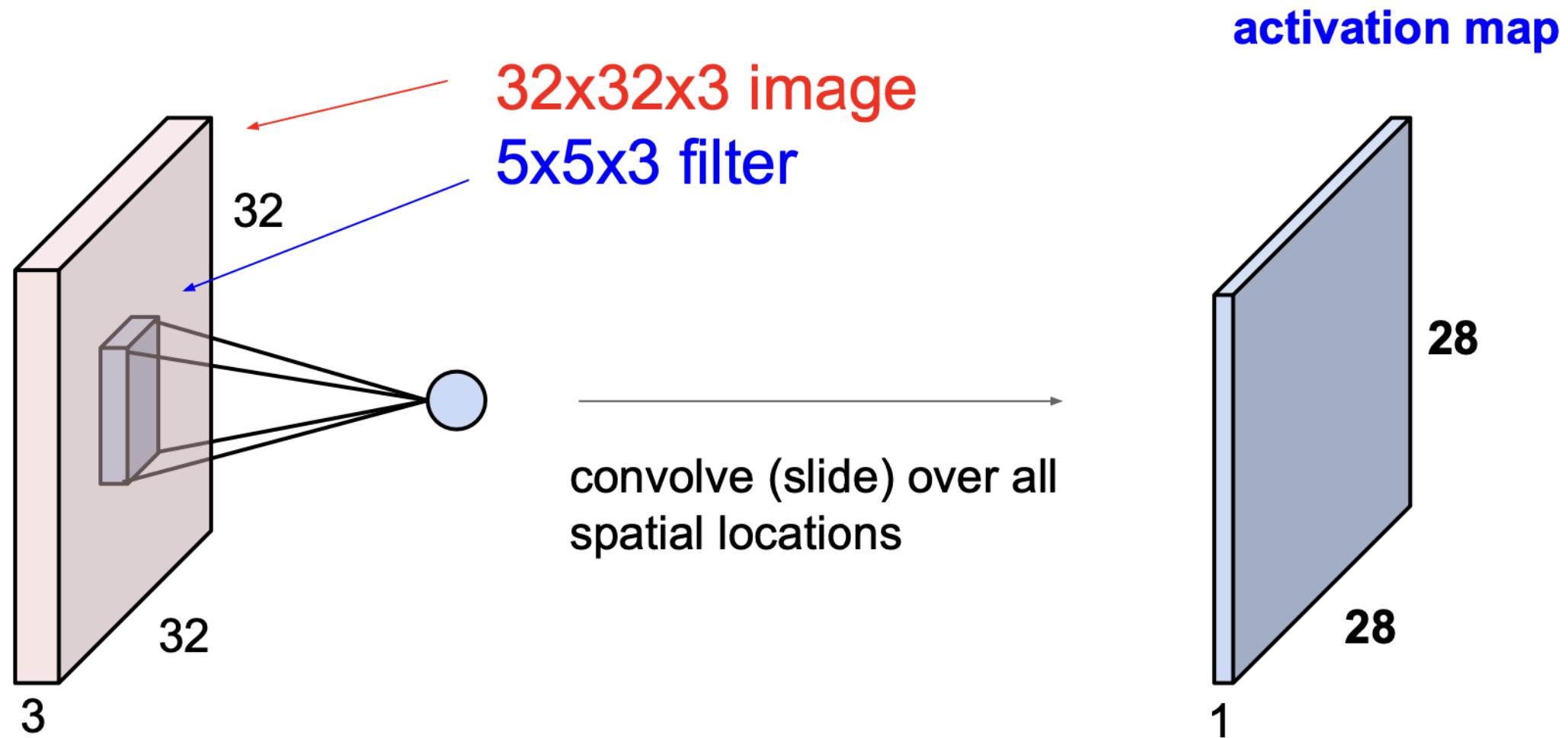
$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x - n_1, y - n_2]$$



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

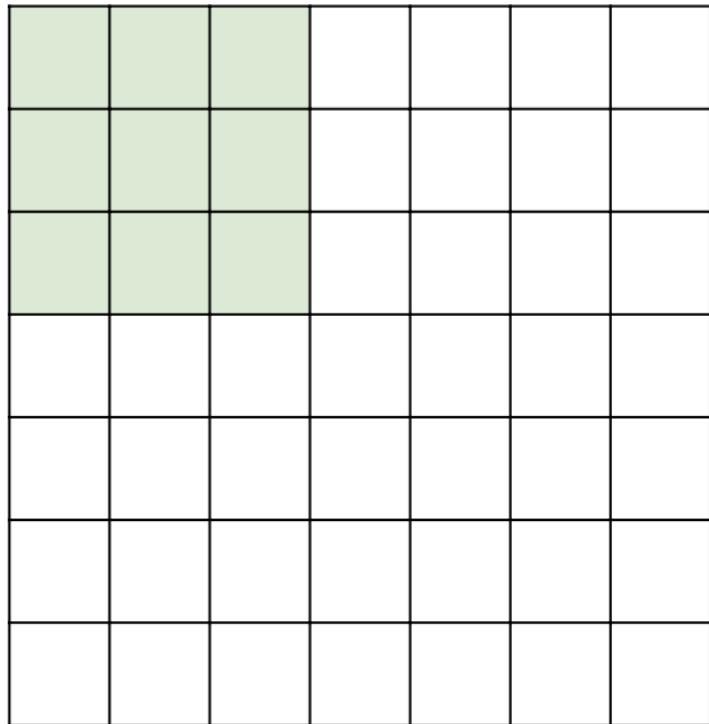


A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

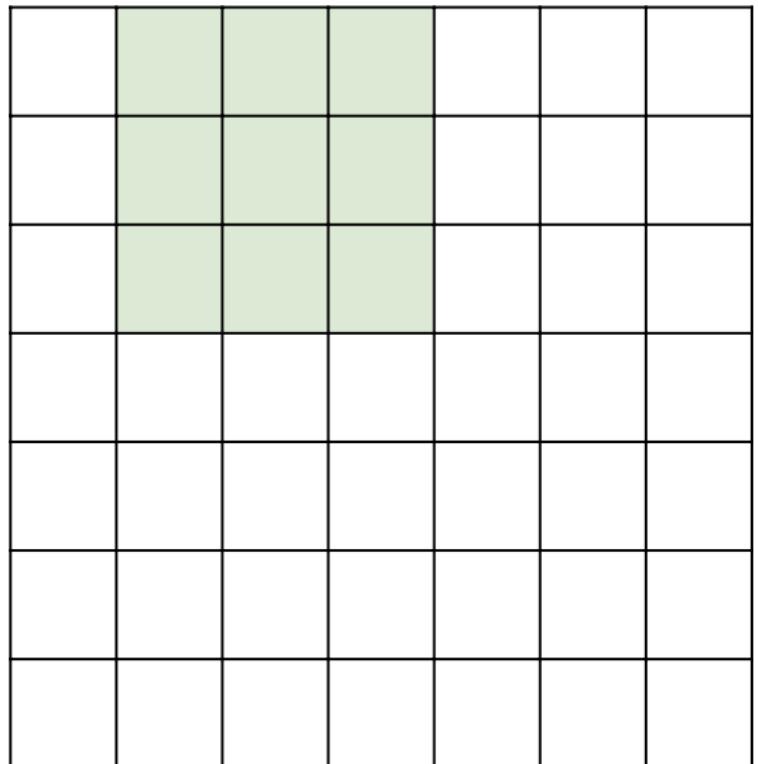


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

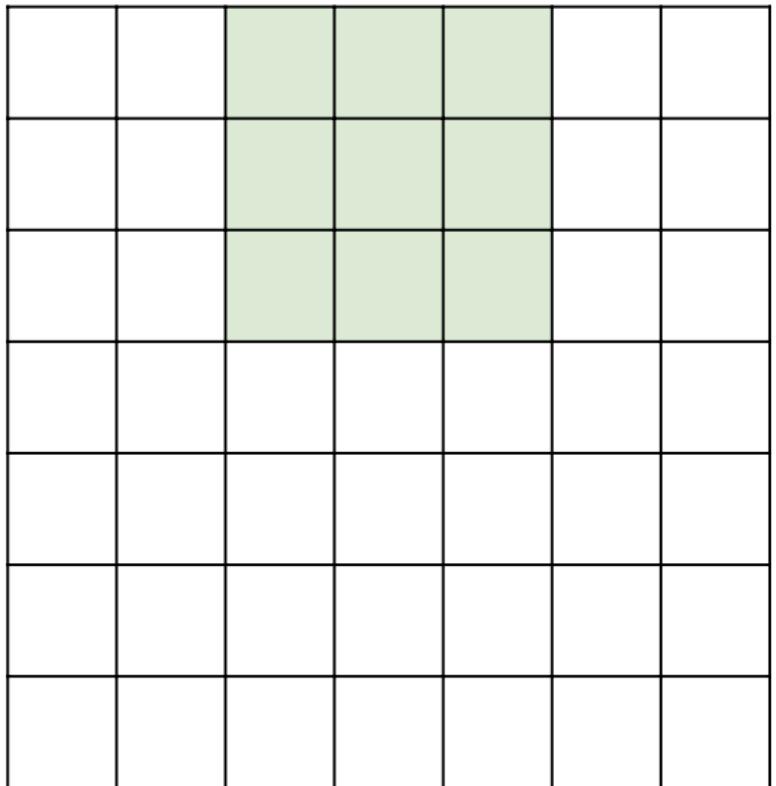


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

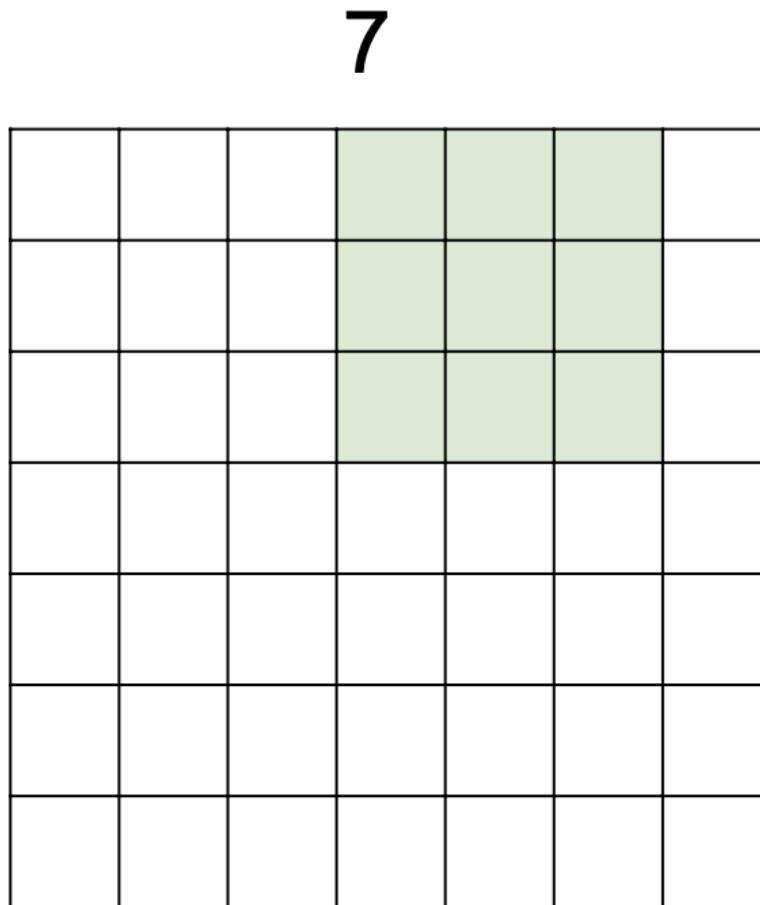
7



7x7 input (spatially)
assume 3x3 filter

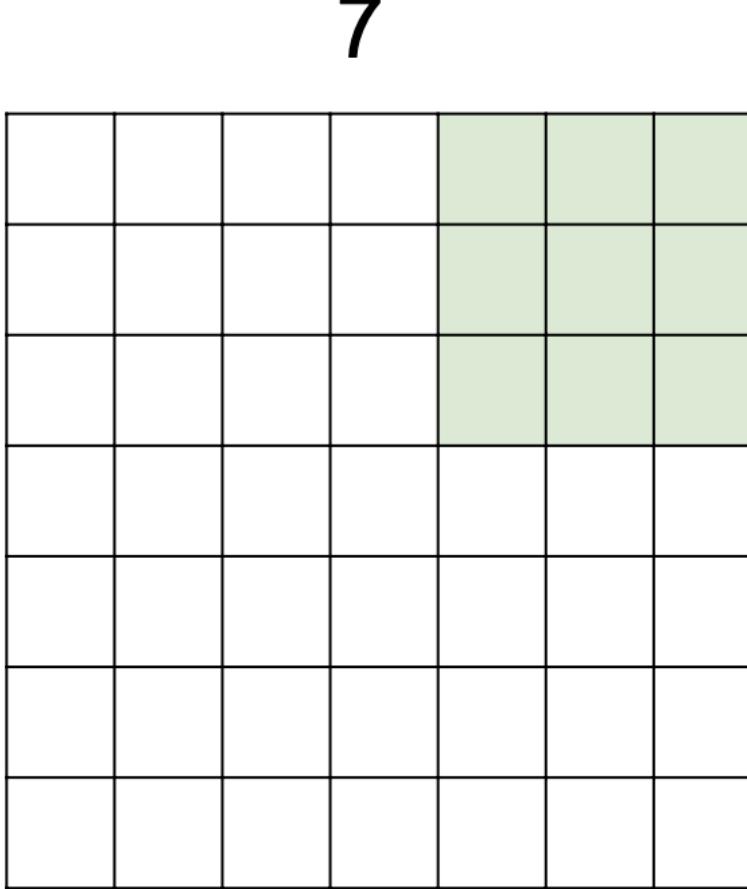
7

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

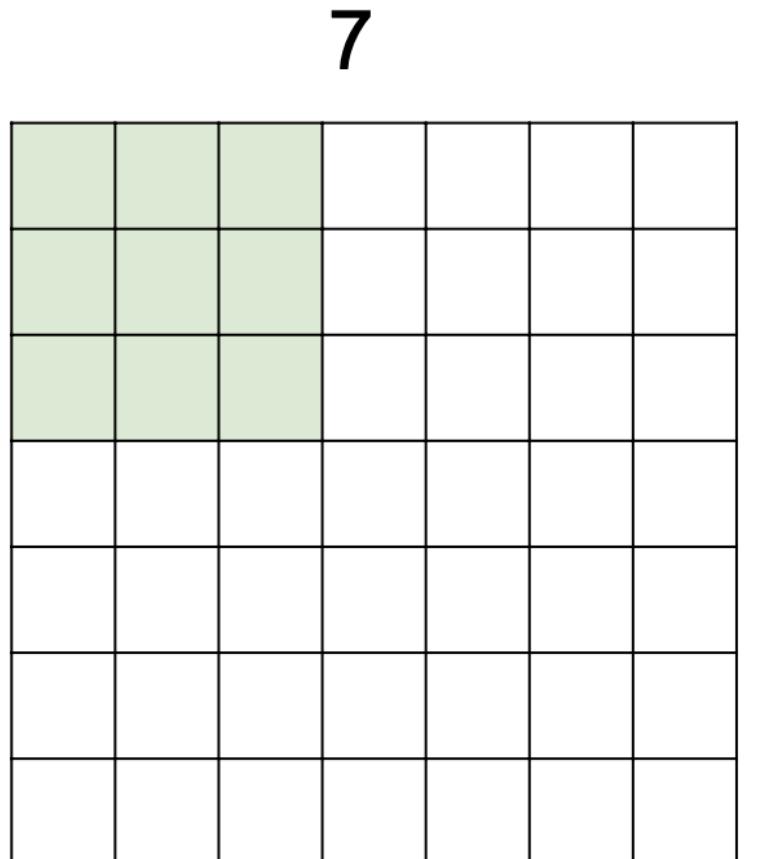
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

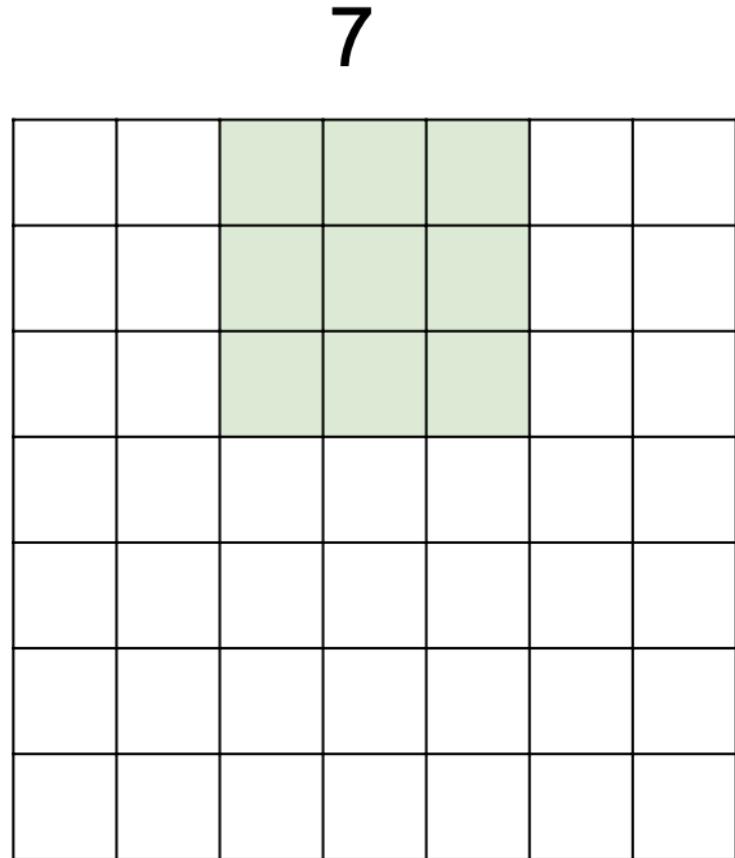
=> 5x5 output

A closer look at spatial dimensions:



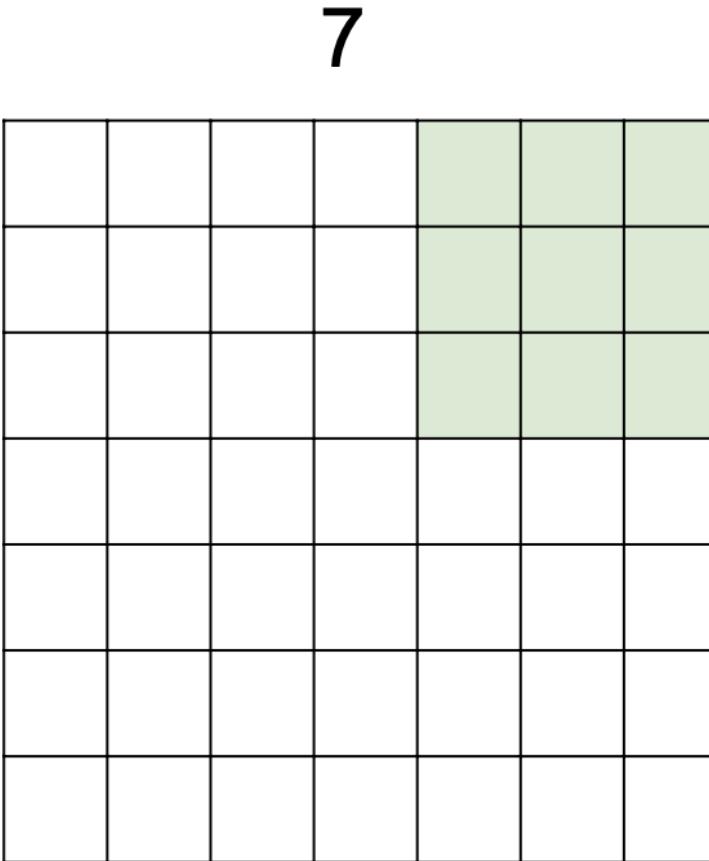
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



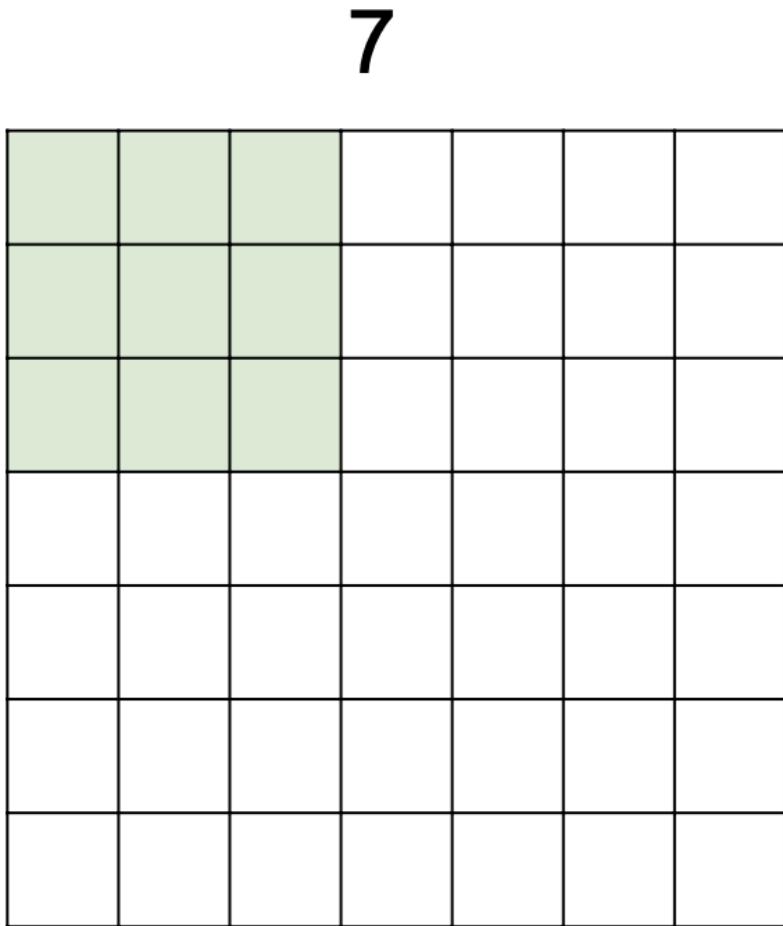
7x7 input (spatially)
assume 3x3 filter
applied with stride 2

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

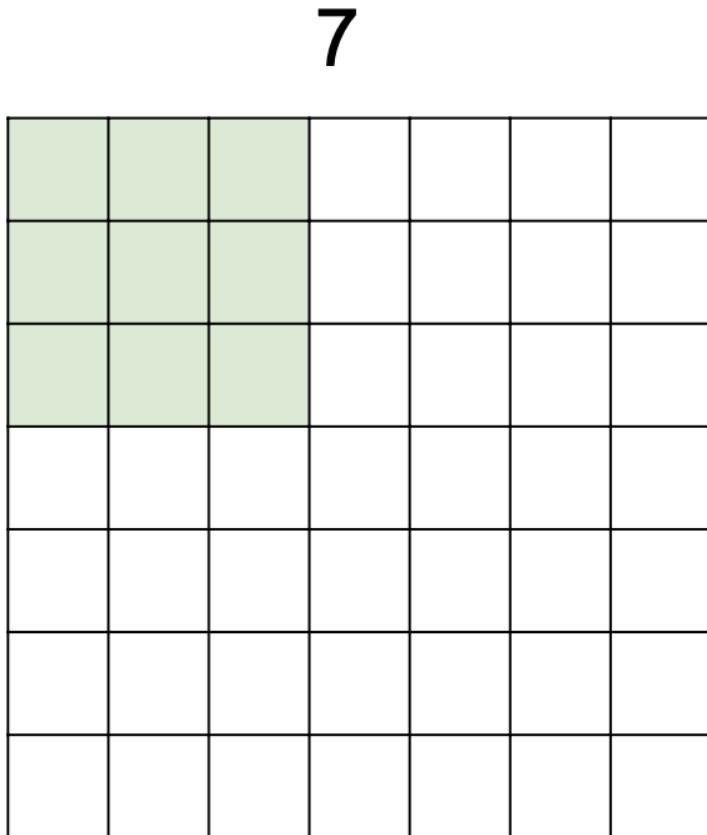
A closer look at spatial dimensions:



7

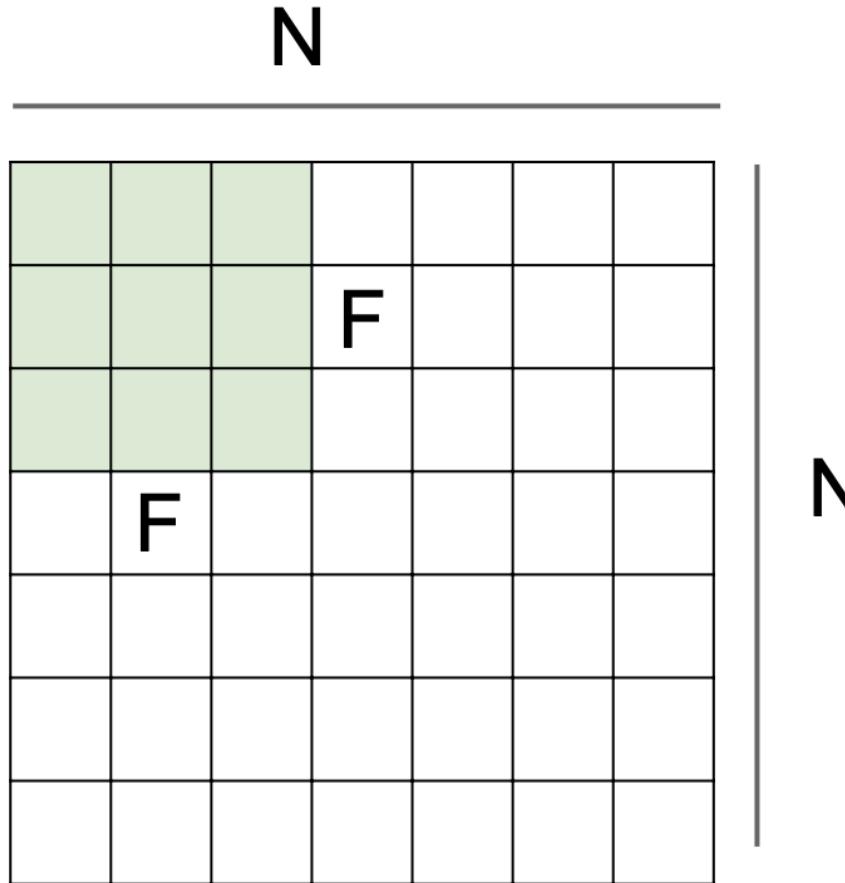
7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.



N

Output size:
(N - F) / stride + 1

e.g. $N = 7, F = 3$:

$$\text{stride } 1 \Rightarrow (7 - 3)/1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3)/2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3)/3 + 1 = 2.33 : \backslash$$

In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

(recall:)

$$(N - F) / \text{stride} + 1$$

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

7x7 output!

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$. (will preserve size spatially)

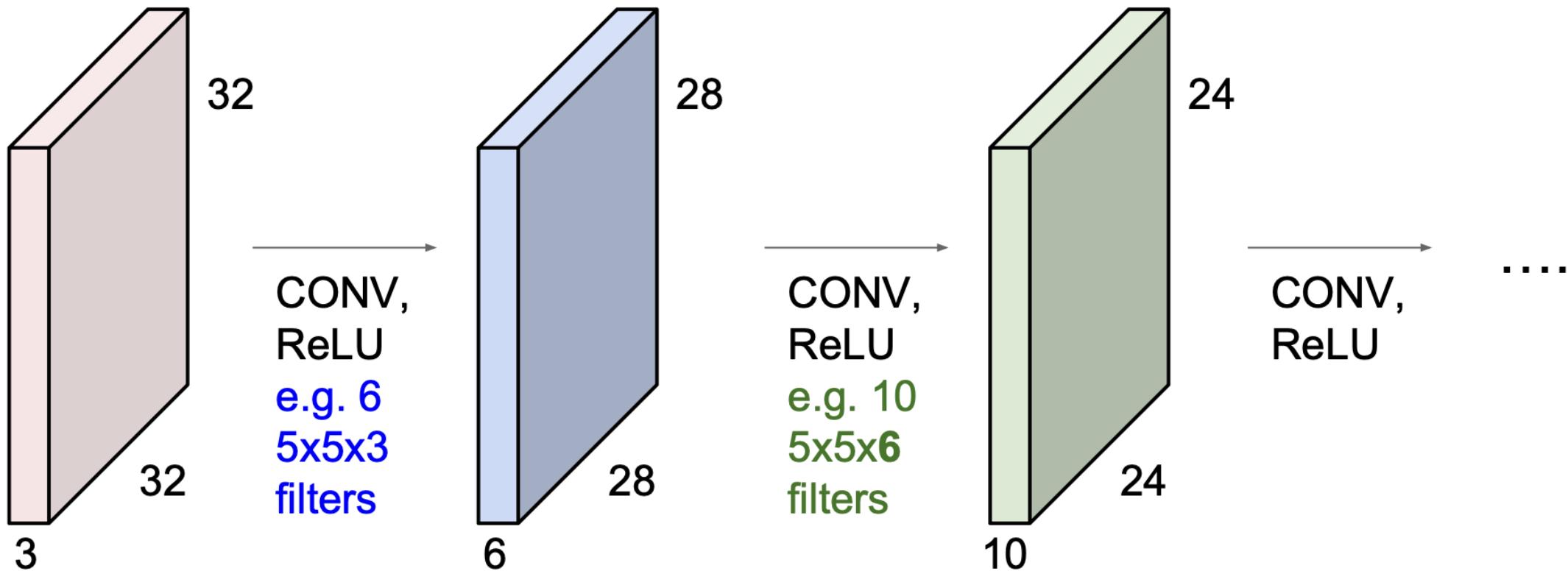
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

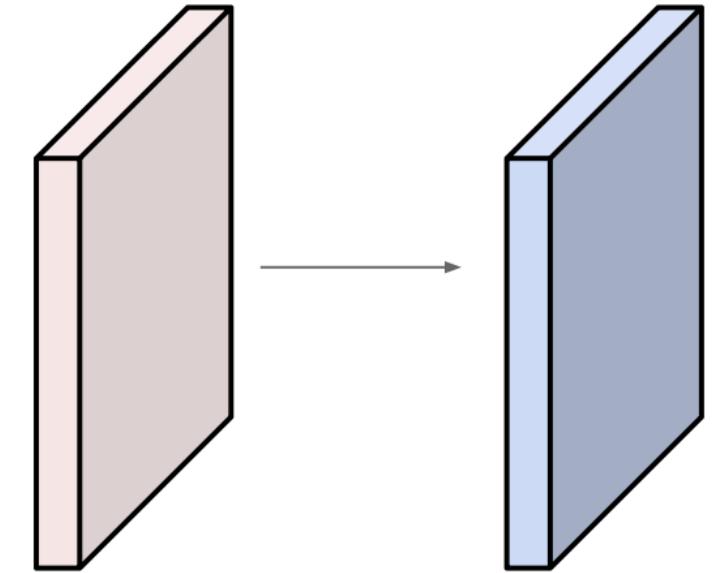
Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 → 28 → 24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

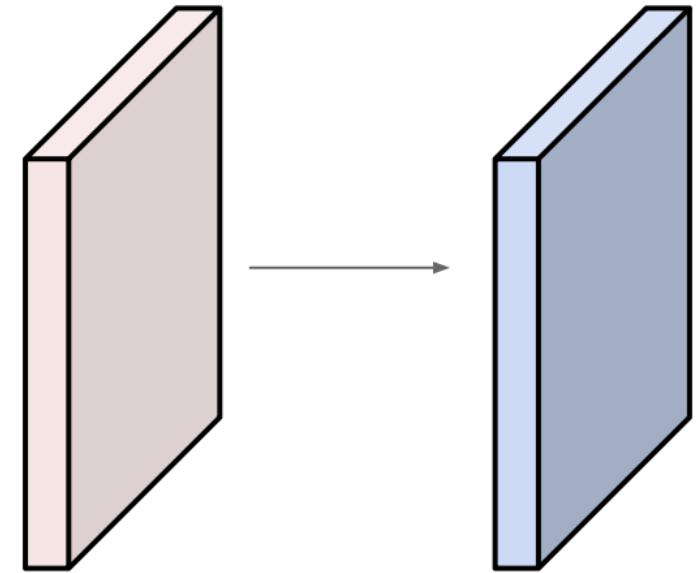


Output volume size: ?

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride **1**, pad **2**



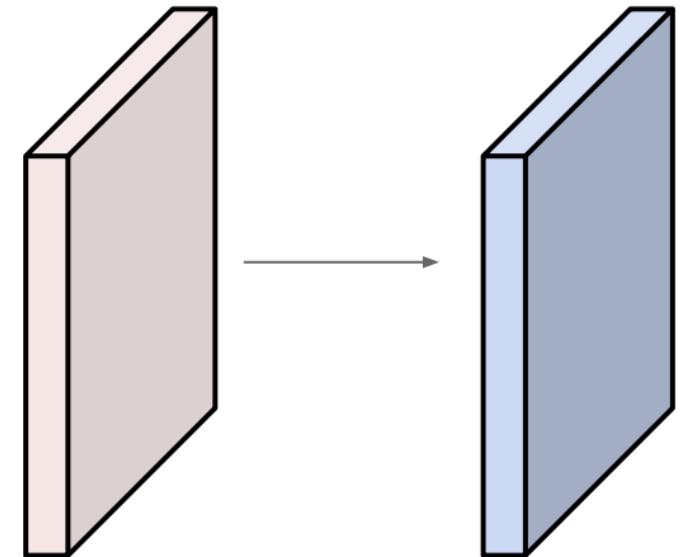
Output volume size:

$(32+2*2-5)/1+1 = 32$ spatially, so

32x32x10

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

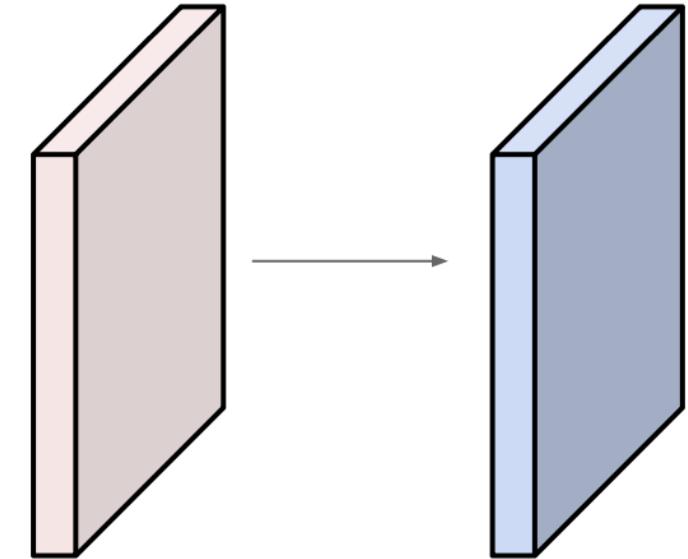


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

10 **5x5** filters with stride 1, pad 2



Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params (+1 for bias)

$$\Rightarrow 76 * 10 = 760$$

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Common settings:

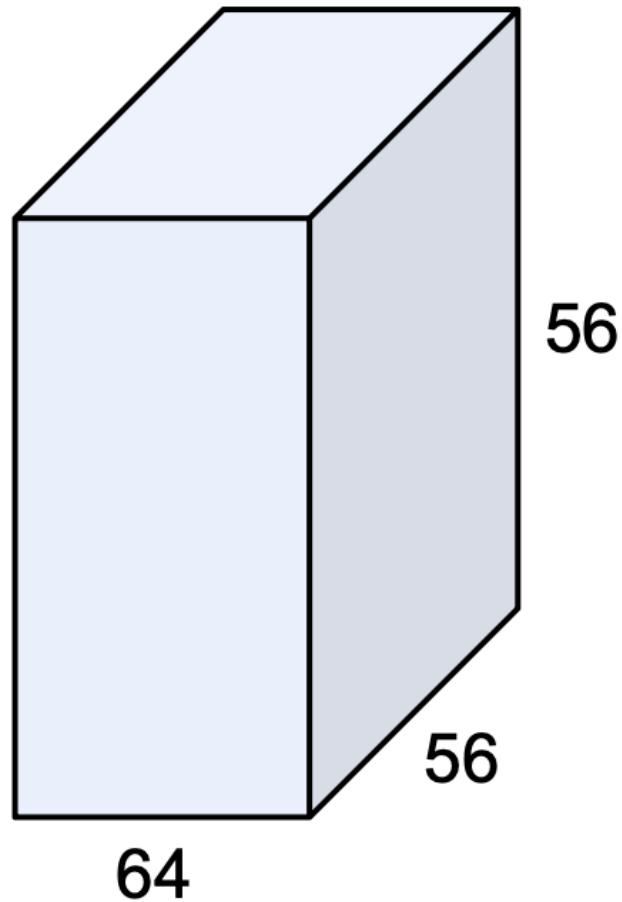
Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

$K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$

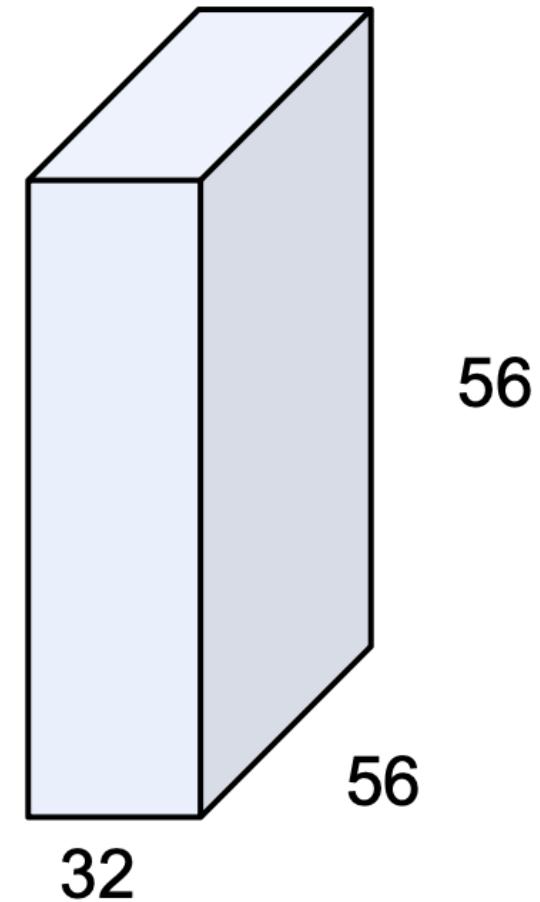
- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ? \text{ (whatever fits)}$
- $F = 1, S = 1, P = 0$

(btw, 1x1 convolution layers make perfect sense)

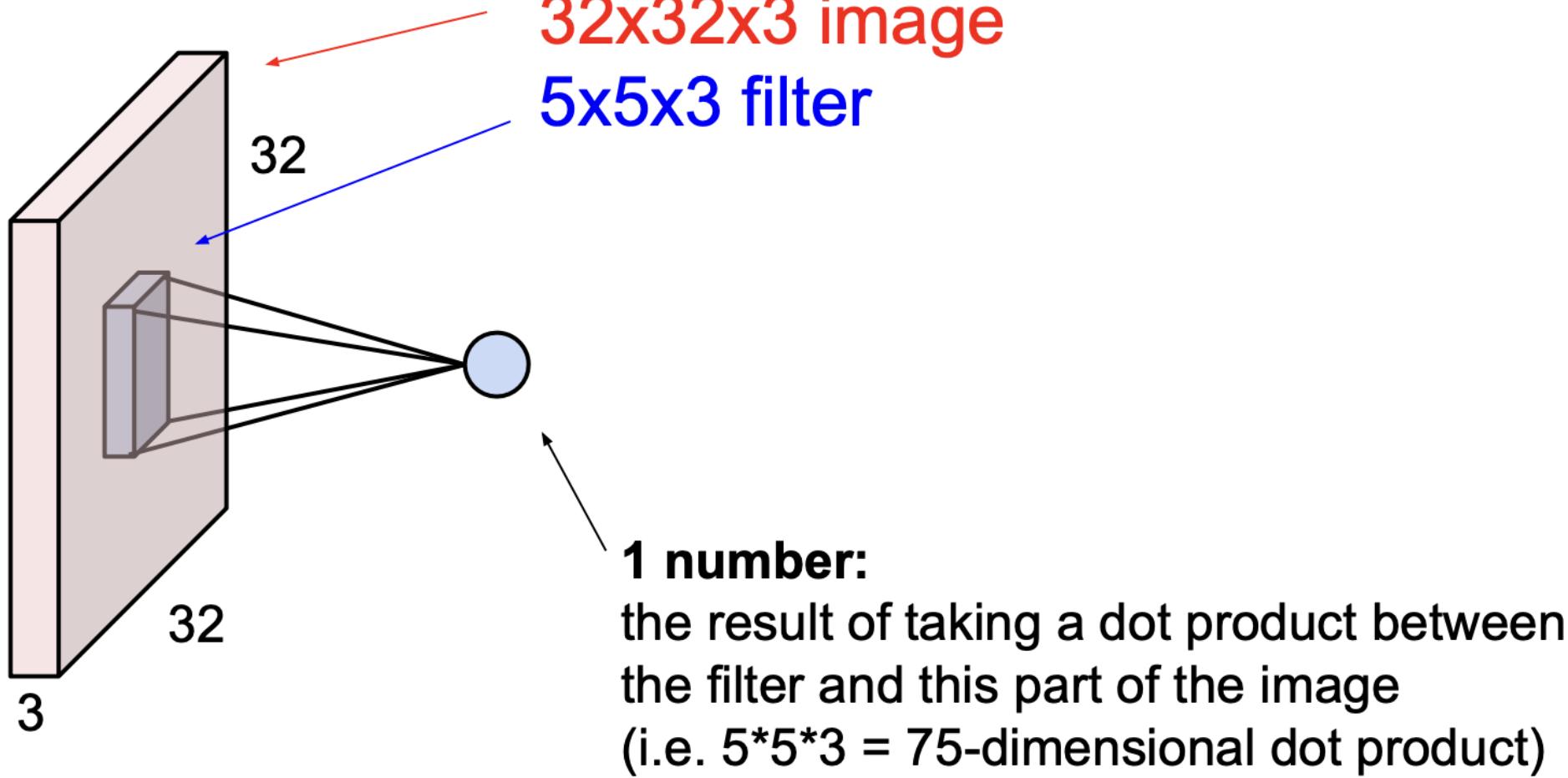


1x1 CONV
with 32 filters

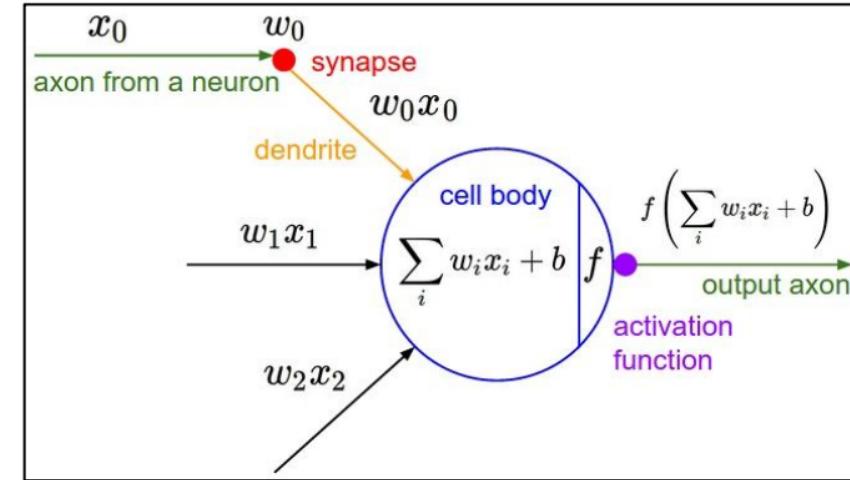
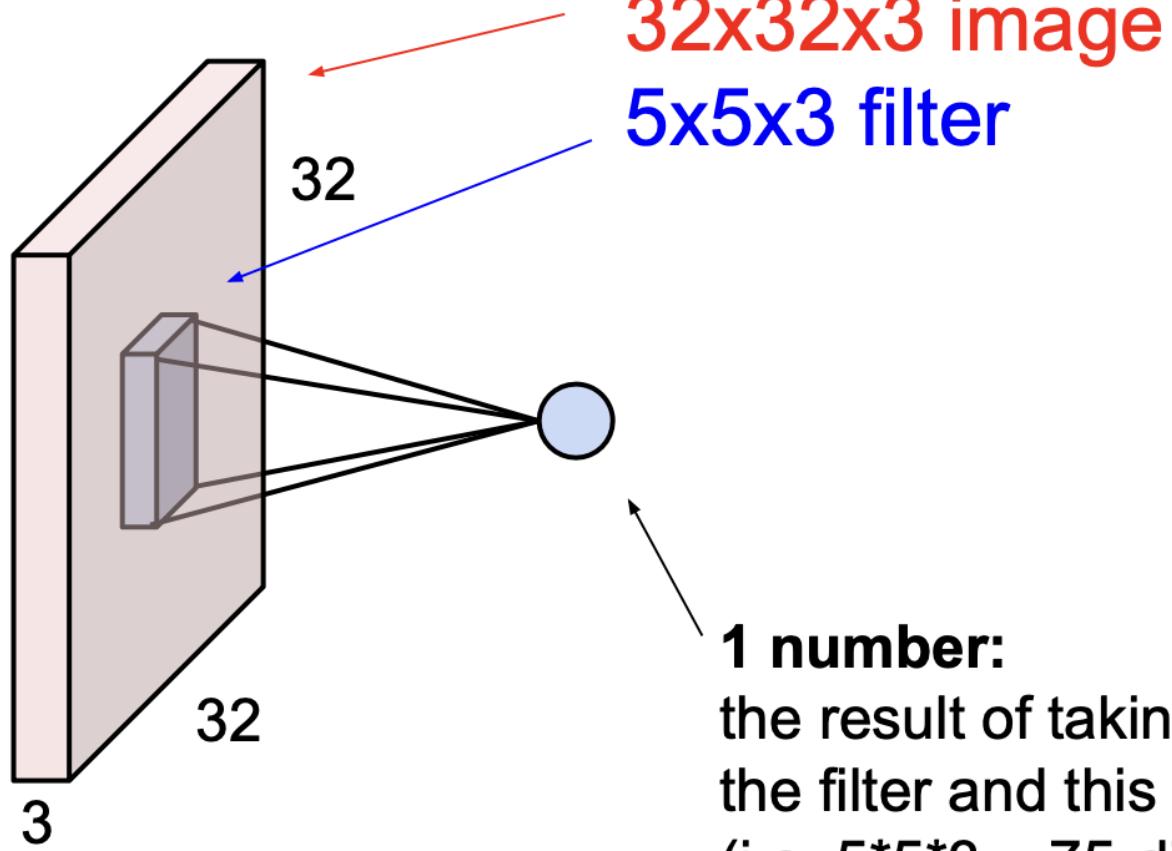
(each filter has size
1x1x64, and performs a
64-dimensional dot
product)



The brain/neuron view of CONV Layer

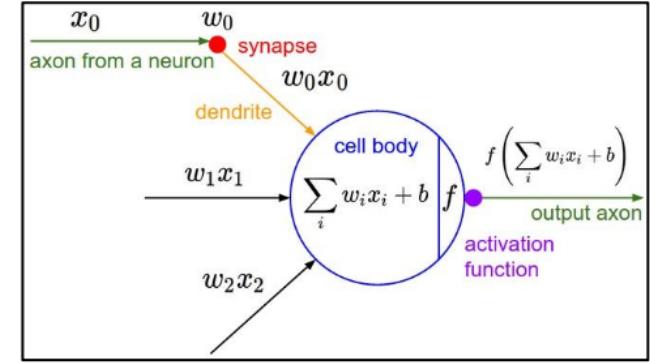
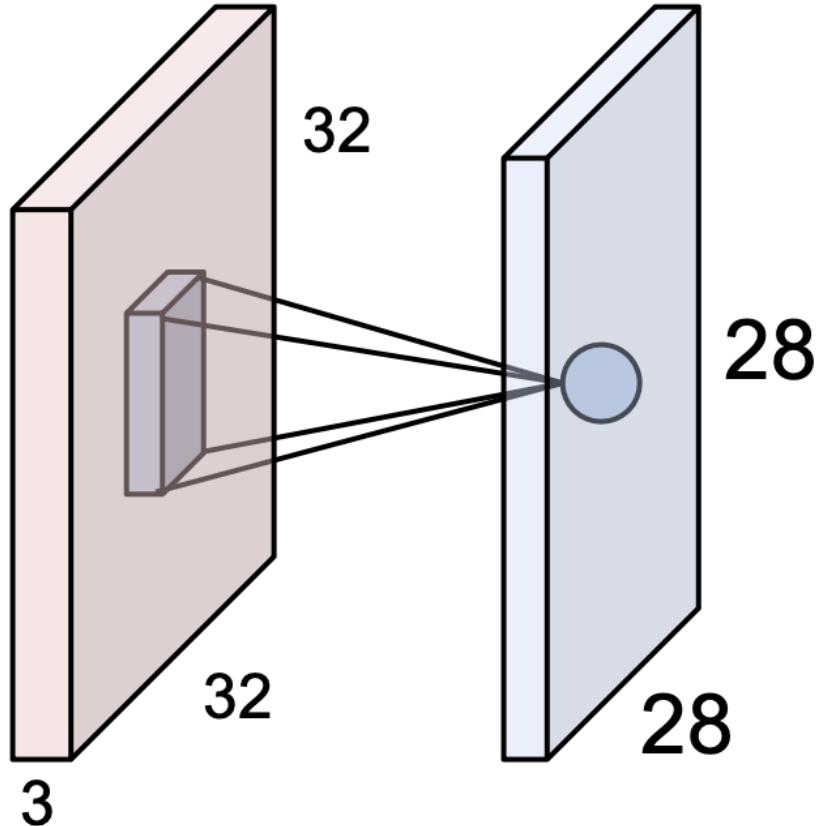


The brain/neuron view of CONV Layer



It's just a neuron with local connectivity...

The brain/neuron view of CONV Layer

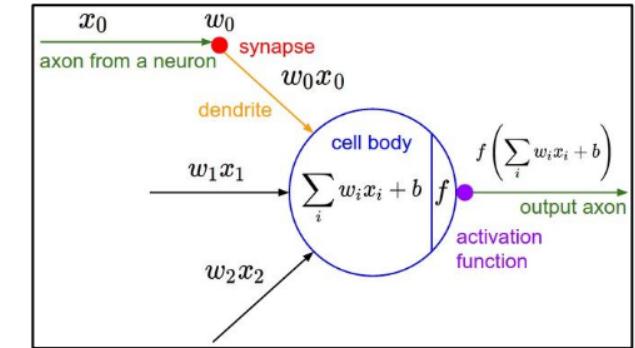
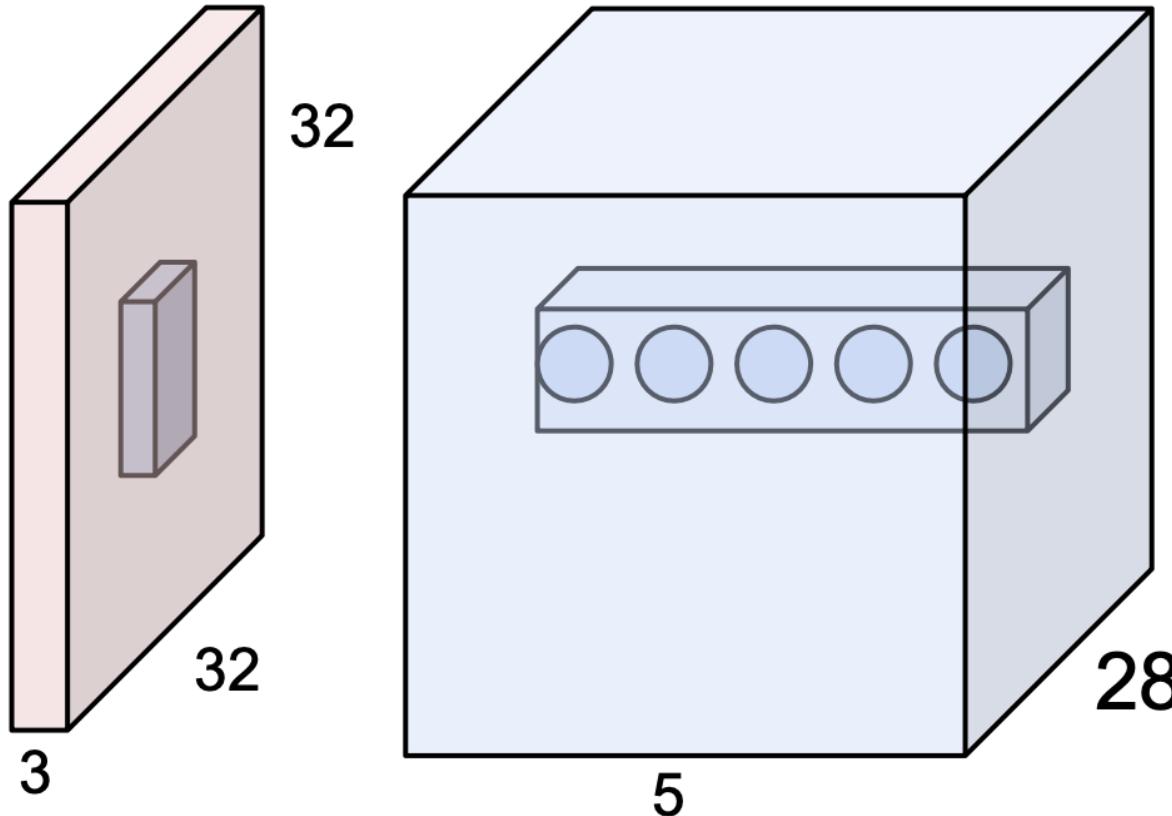


An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

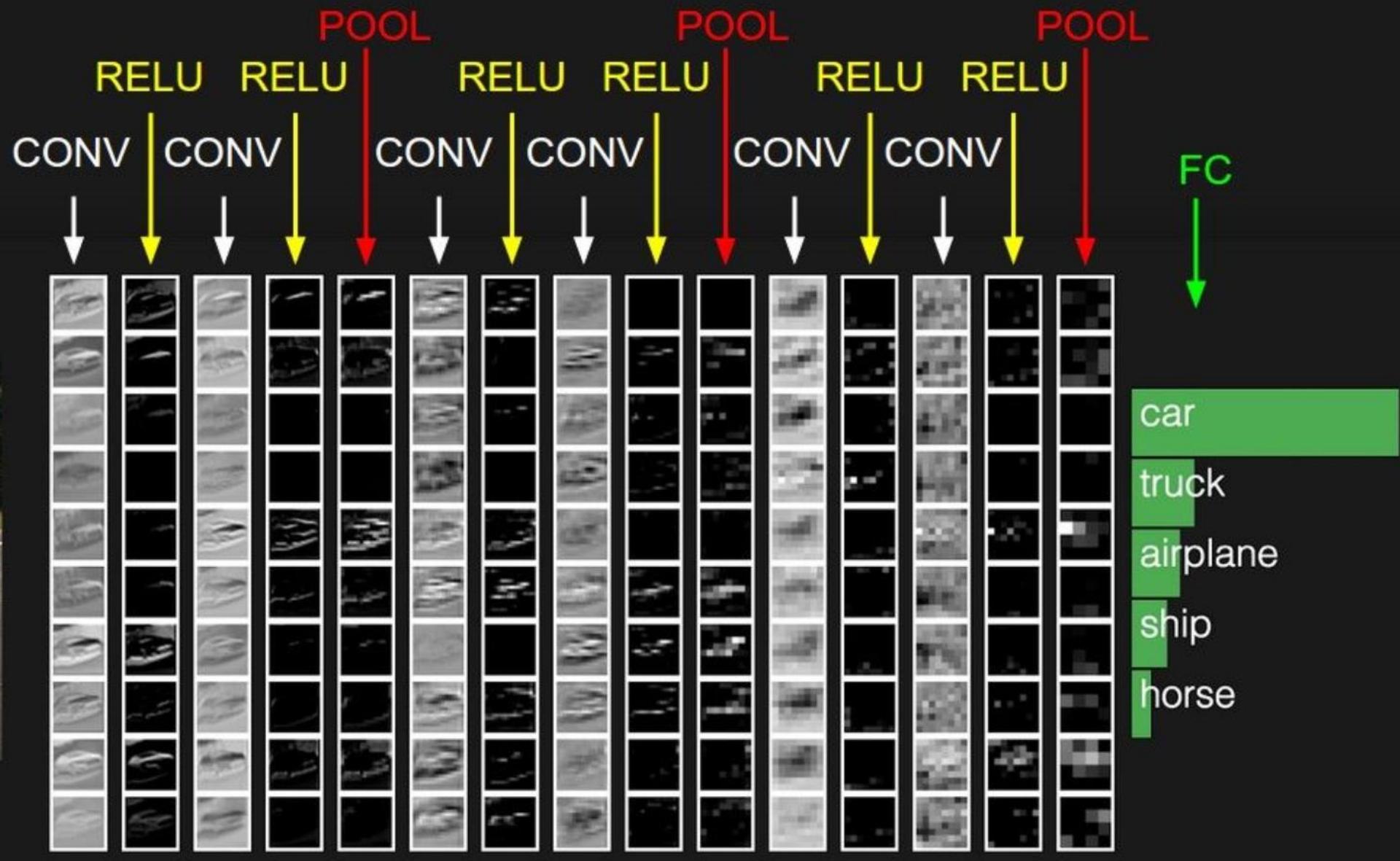
“5x5 filter” -> “5x5 receptive field for each neuron”

The brain/neuron view of CONV Layer



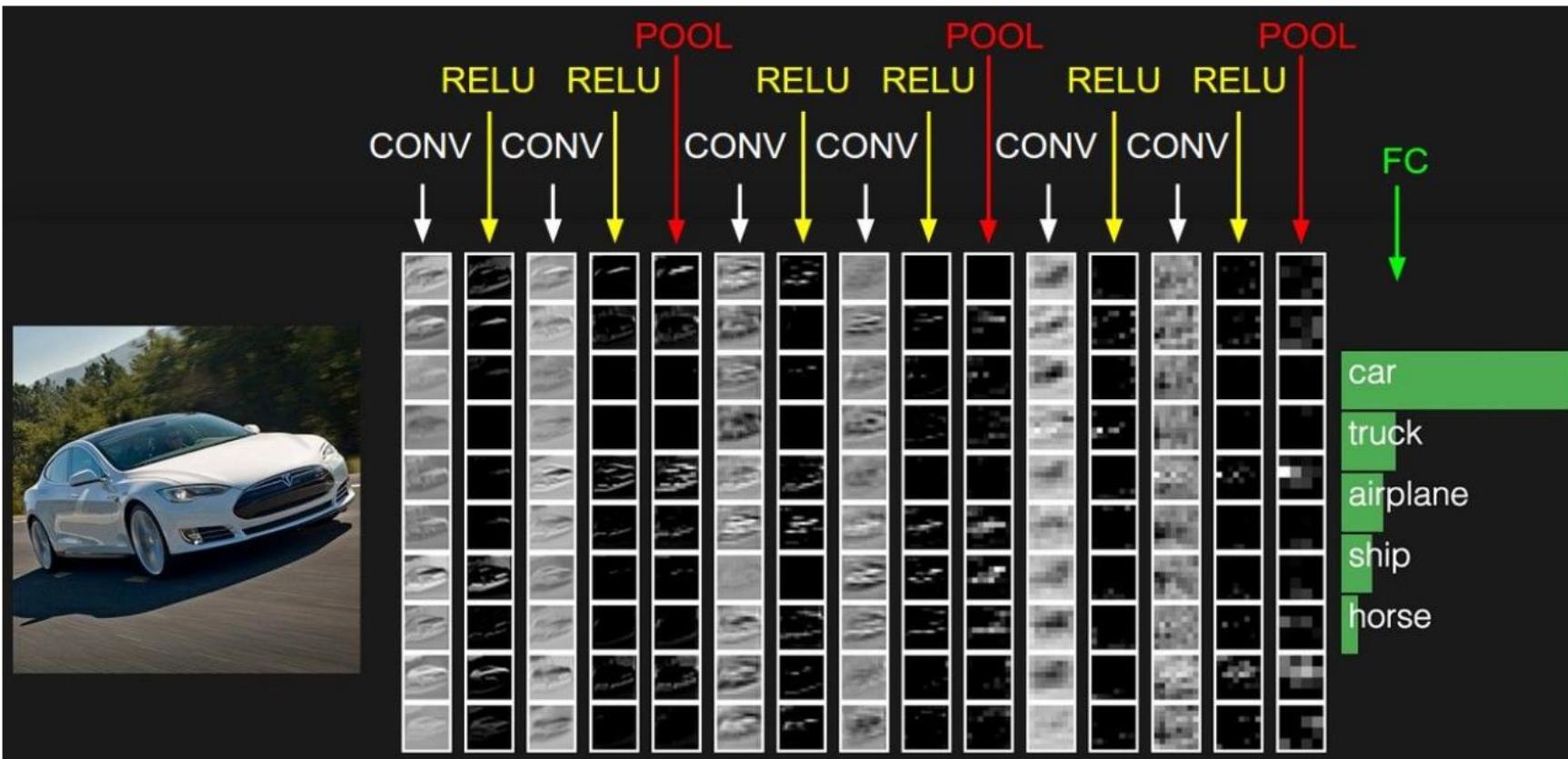
E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
($28 \times 28 \times 5$)

There will be 5 different
neurons all looking at the same
region in the input volume



Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



LeNet-5

- *Gradient Based Learning Applied To Document Recognition - Y. Lecun, L. Bottou, Y. Bengio, P. Haffner; 1998*
- Helped establish how we use CNNs today
- Replaced manual feature extraction

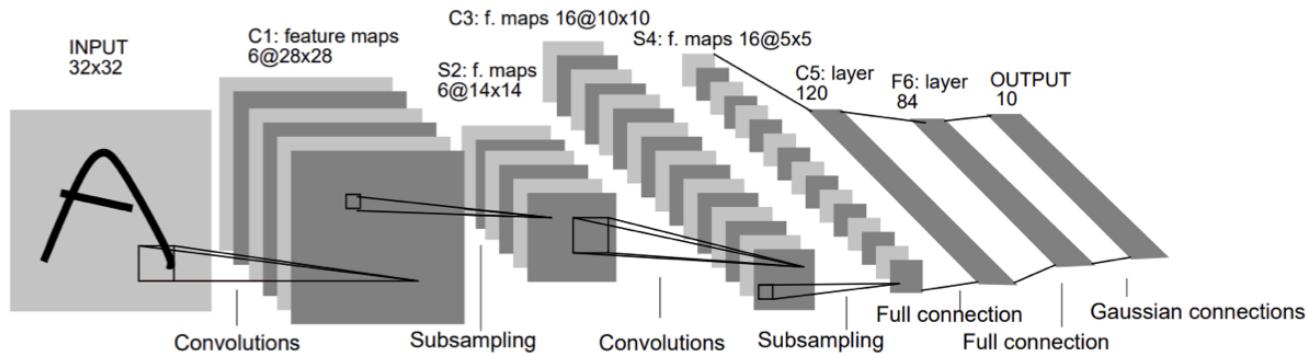
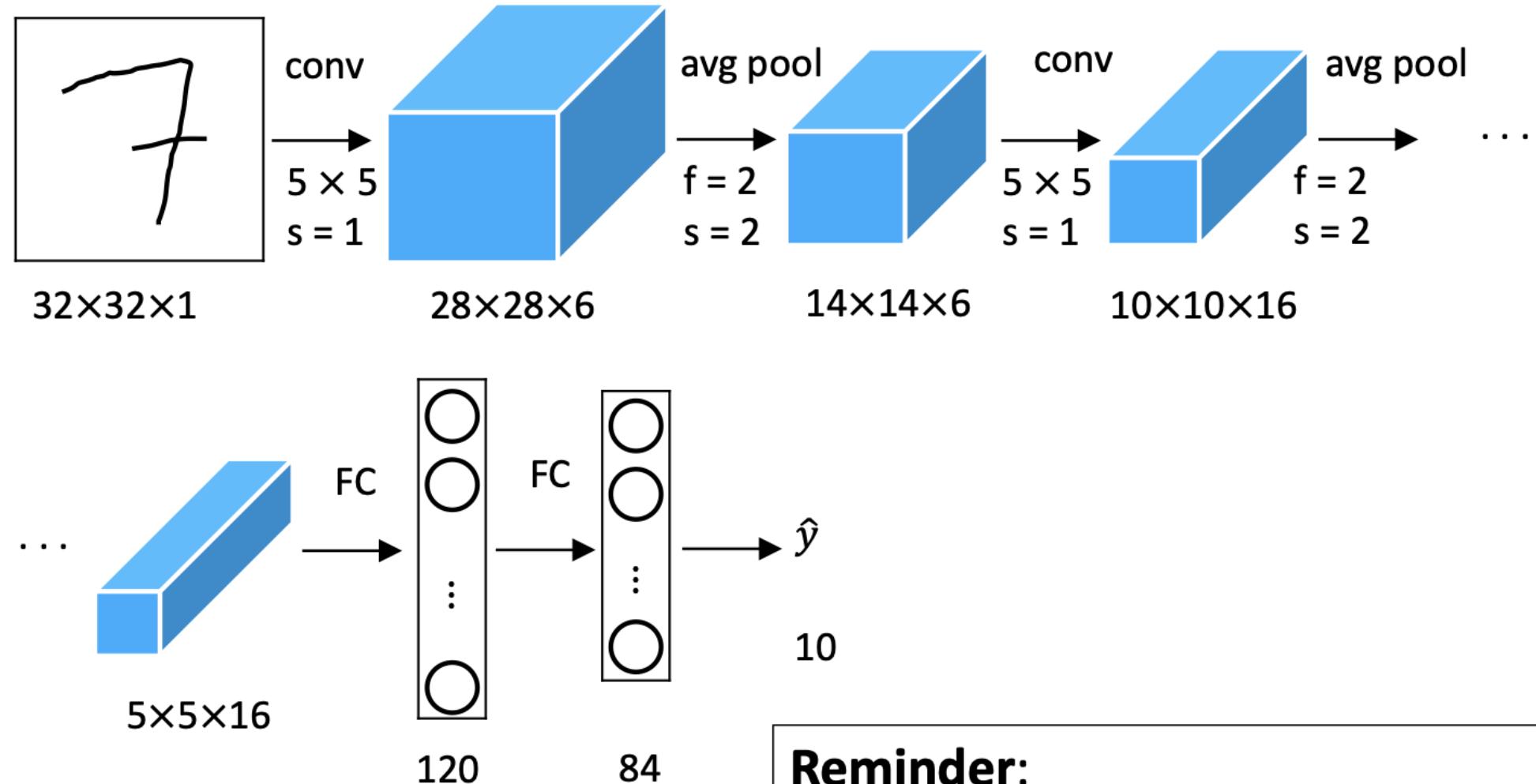


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

[LeCun et al., 1998]

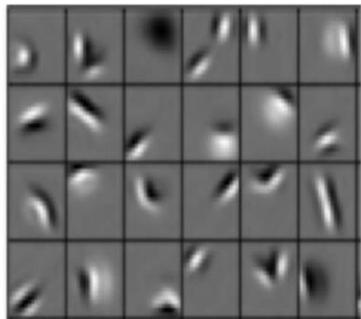
LeNet-5



Reminder:

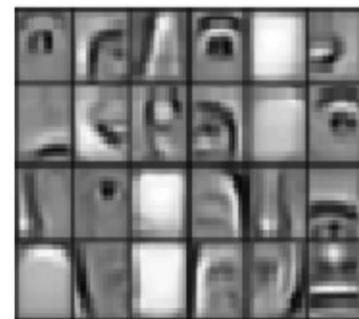
Key idea:
learn hierarchy of features
directly from the data
(rather than hand-engineering them)

Low level features



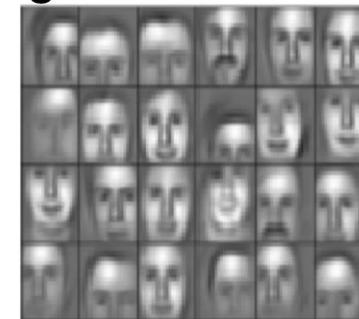
Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



Facial structure

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

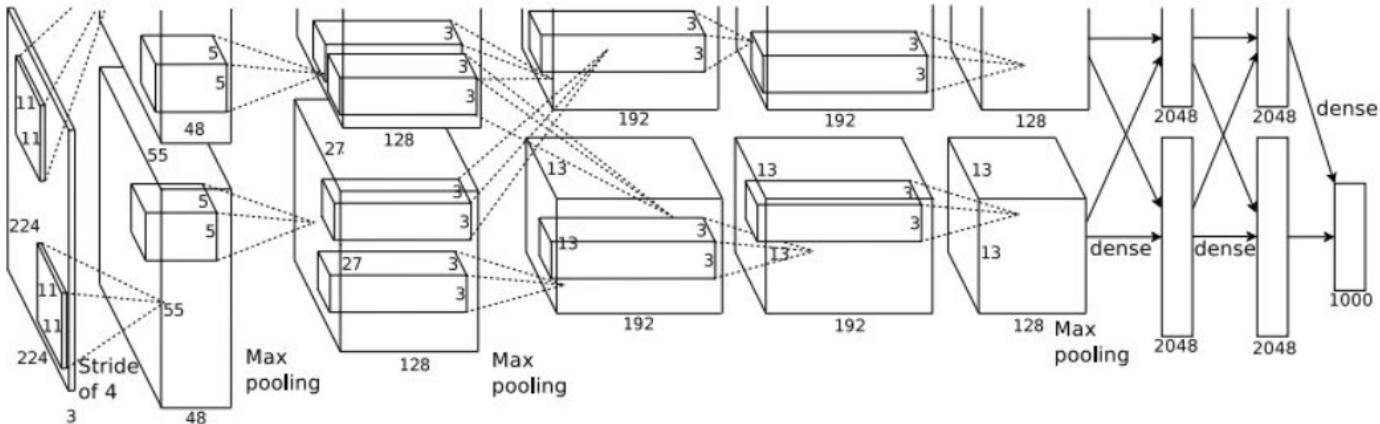
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

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 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	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	conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150K$ params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 3) \times 64 = 1,728$

CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2M$ params: $(3 \times 3 \times 64) \times 64 = 36,864$

POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800K$ params: 0

CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 64) \times 128 = 73,728$

CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6M$ params: $(3 \times 3 \times 128) \times 128 = 147,456$

POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400K$ params: 0

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 128) \times 256 = 294,912$

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$

CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800K$ params: $(3 \times 3 \times 256) \times 256 = 589,824$

POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200K$ params: 0

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: 0

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100K$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$

POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25K$ params: 0

FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$

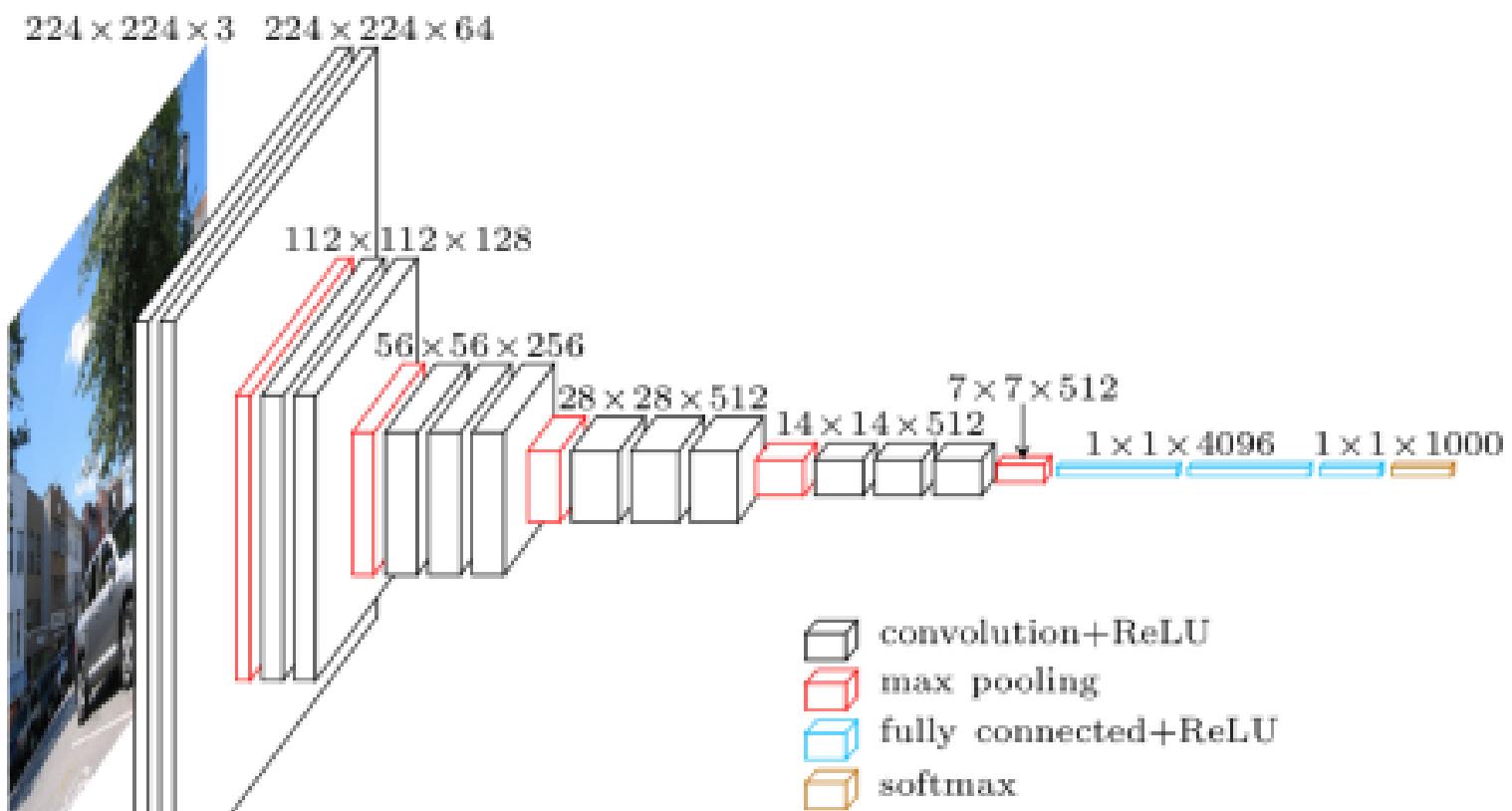
FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	conv3-64
conv3-64	conv3-64	conv3-64	conv3-64
maxpool			
conv3-128	conv3-128	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
conv1-256	conv3-256	conv3-256	conv3-256
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
conv1-512	conv3-512	conv3-512	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
conv1-512	conv3-512	conv3-512	conv3-512
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

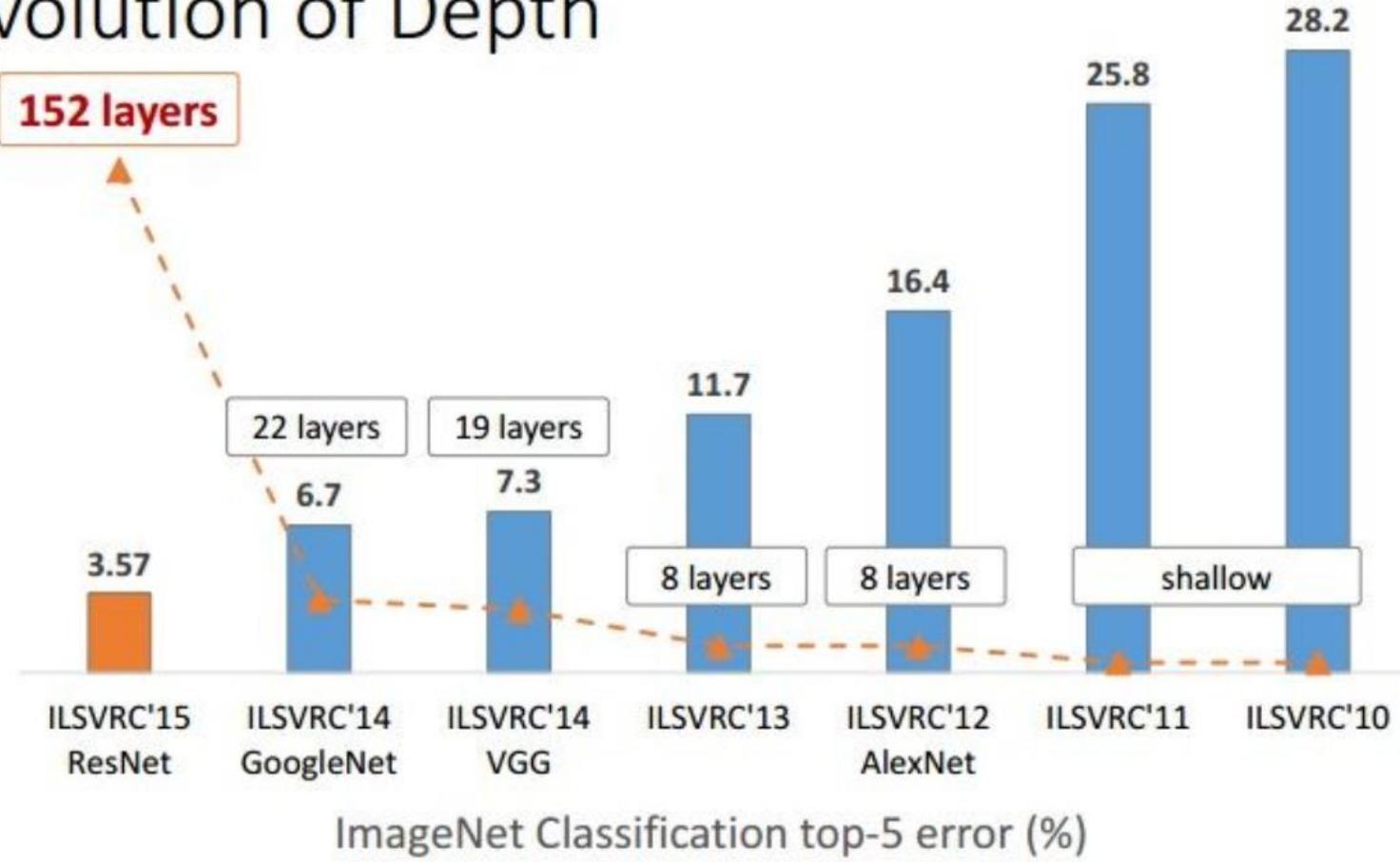
VGGNet16

- Replaces layers with large kernels with stacked (3,3) kernels. This uses fewer parameters but produces the same receptive field.
- Uses pooling layers instead of strides, again reducing the number of parameters.
- Even with these savings, VGGNet16 is a BIG model—almost 600MB.

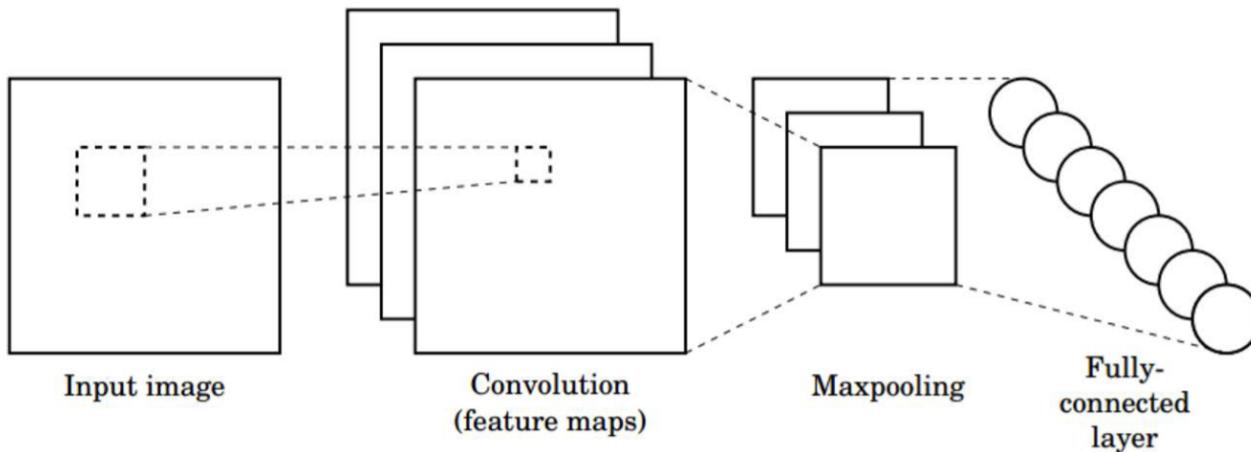


<https://arxiv.org/abs/1409.1556v6>

Revolution of Depth



CNNs for Classification



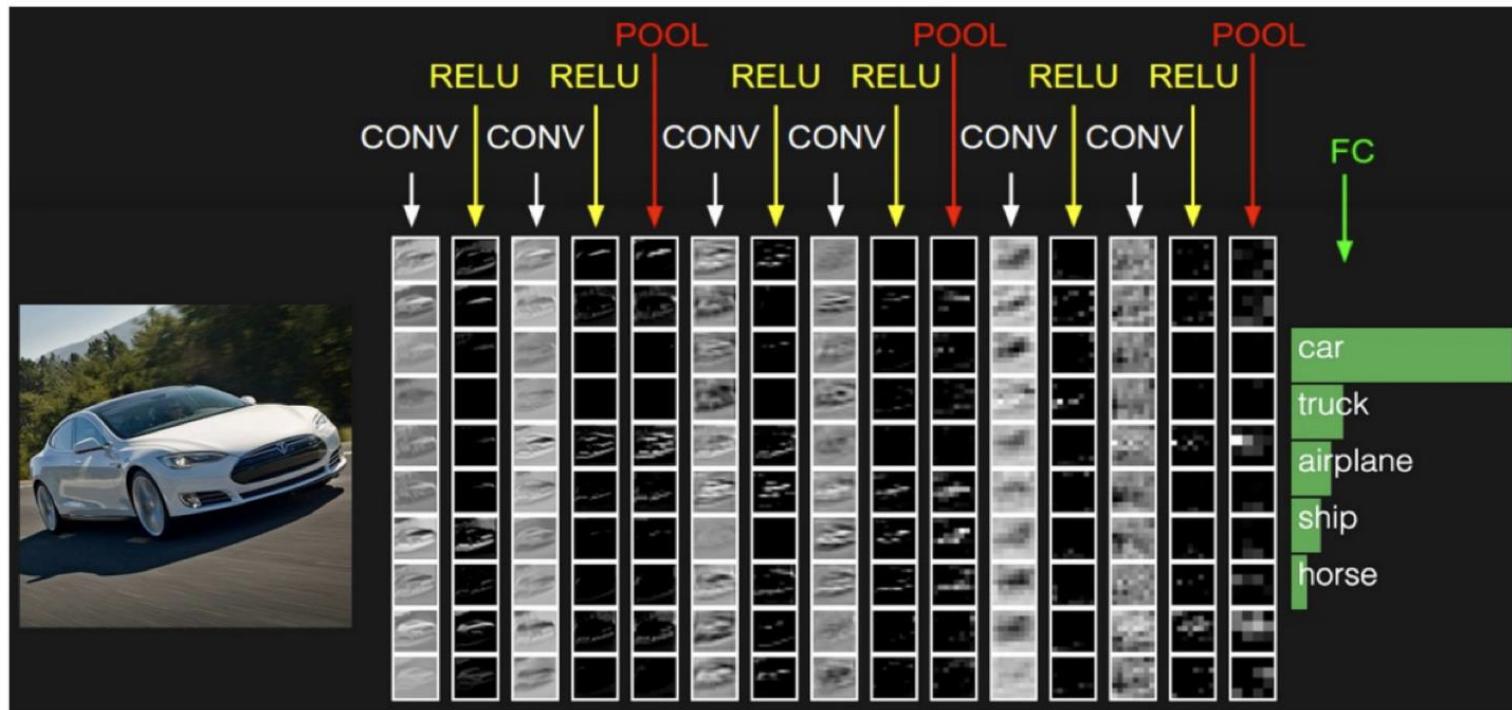
- 1. Convolution:** Apply filters to generate feature maps.
- 2. Non-linearity:** Often ReLU.
- 3. Pooling:** Downsampling operation on each feature map.

Train model with image data.

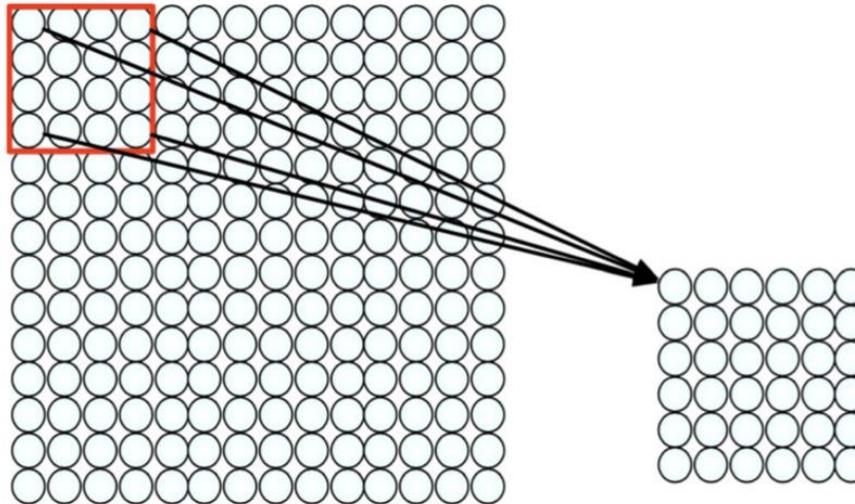
Learn weights of filters in convolutional layers.

```
tf.keras.layers.Conv2D  
tf.keras.activations.*  
tf.keras.layers.MaxPool2D
```

Example – Six convolutional layers



Convolutional Layers: Local Connectivity



`tf.keras.layers.Conv2D`

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

**4x4 filter:
matrix of
weights w_{ij}**

$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

for neuron (p,q) in hidden layer

- 1) applying a window of weights
- 2) computing linear combinations
- 3) activating with non-linear function

```

model = keras.models.Sequential([
    keras.layers.Conv2D(64, 7, activation="relu", padding="same",
                       input_shape=[28, 28, 1]),
    keras.layers.MaxPooling2D(2),
    keras.layers.Conv2D(128, 3, activation="relu", padding="same"),
    keras.layers.Conv2D(128, 3, activation="relu", padding="same"),
    keras.layers.MaxPooling2D(2),
    keras.layers.Conv2D(256, 3, activation="relu", padding="same"),
    keras.layers.Conv2D(256, 3, activation="relu", padding="same"),
    keras.layers.MaxPooling2D(2),
    keras.layers.Flatten(),
    keras.layers.Dense(128, activation="relu"),
    keras.layers.Dropout(0.5),
    keras.layers.Dense(64, activation="relu"),
    keras.layers.Dropout(0.5),
    keras.layers.Dense(10, activation="softmax")
])

```

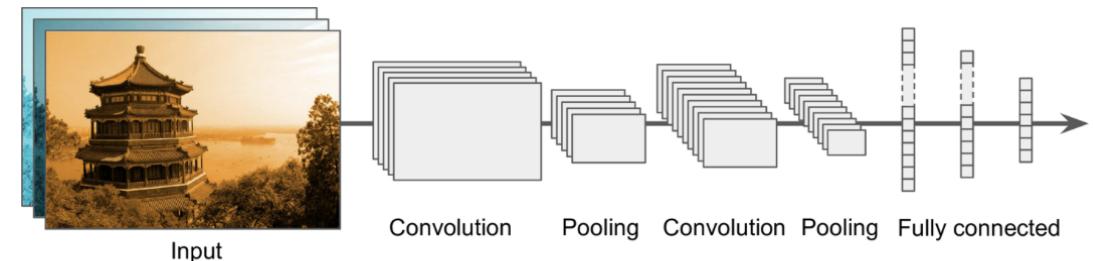
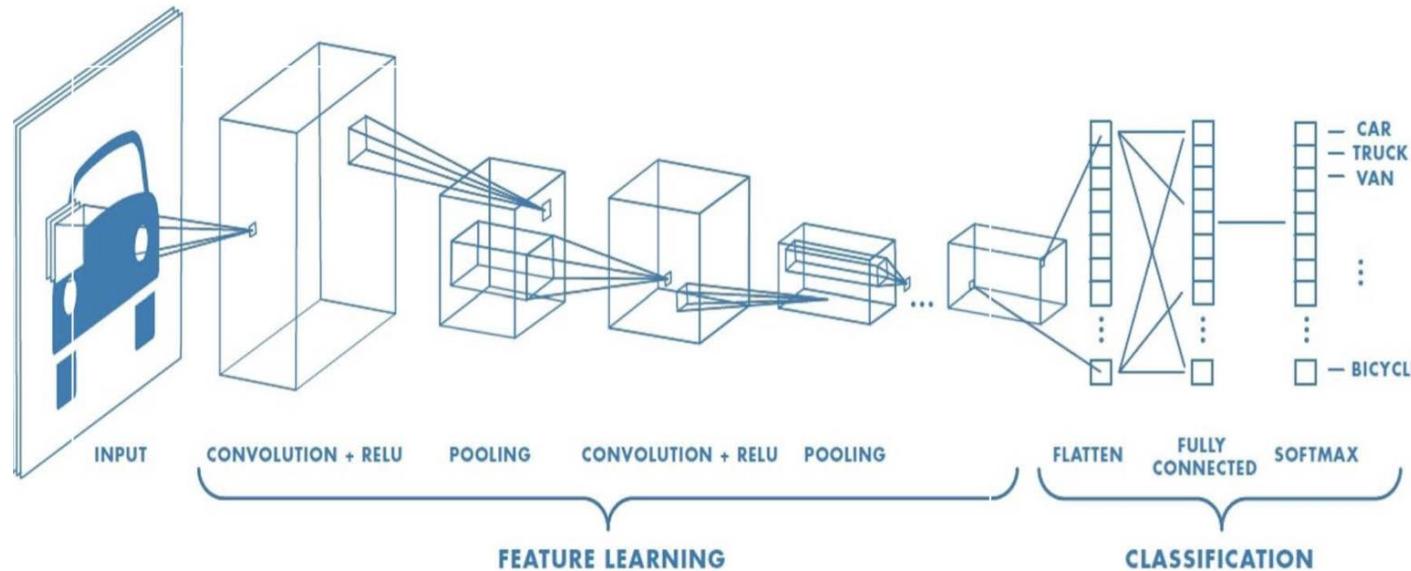


Figure 14-11. Typical CNN architecture

CNNs for Classification: Class Probabilities



$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

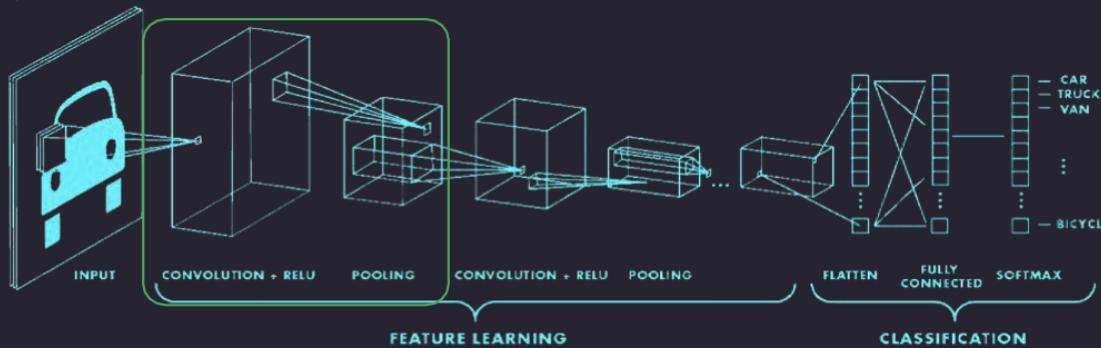
Putting it all together

```
import tensorflow as tf

def generate_model():
    model = tf.keras.Sequential([
        # first convolutional layer
        tf.keras.layers.Conv2D(32, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),  
  

        # second convolutional layer
        tf.keras.layers.Conv2D(64, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),  
  

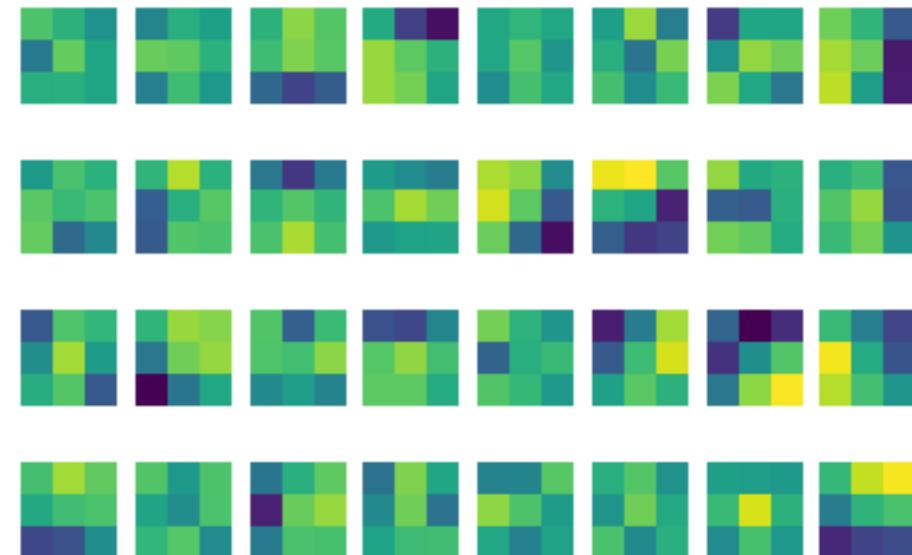
        # fully connected classifier
        tf.keras.layers.Flatten(),
        tf.keras.layers.Dense(1024, activation='relu'),
        tf.keras.layers.Dense(10, activation='softmax')
        # 10 outputs
    ])
    return model
```



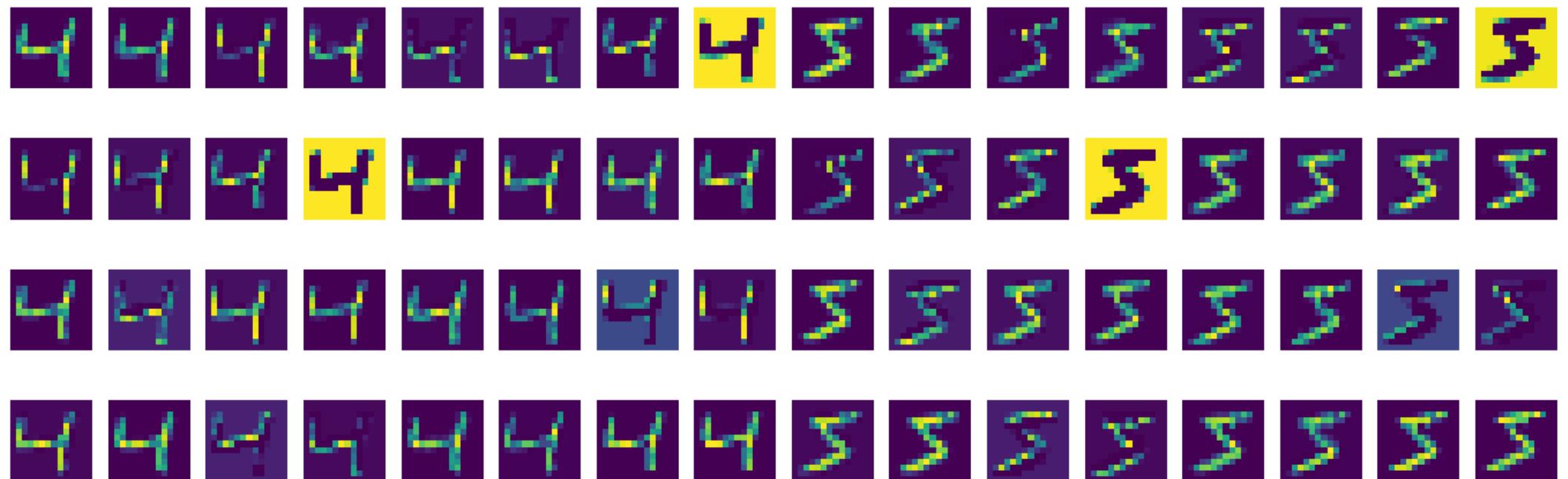
Convolutional Filters

The filters in a convolutional layer are learned from a random initialisation

The filters job is to pick out details in the image, so they learn to do things like edge detection, mapping transitions between colours, identifying areas of high or low contrast



Convolutional Filters



Filter Shapes

The shapes of filters can be somewhat misleading—if we think about a layer with 32 filters and a (3,3) kernel, we might think that it has a shape of (3,3,32), like an image.

But it all depends on the input—each filter has the same number of channels as the input tensor

So the first Conv Layer might have as input a (28,28,1) image, and 32 (3,3) filters. So the shape of the filters is:

$$(3,3,1,32)$$

But the second layer, which has 64 filters, will have shapes like this:

$$(3,3,32,64)$$

Each of the 64 filters, has the shape (3,3,32)—one kernel for each input channel

Kernel sizes

Consider the number of parameters that are needed for a convolutional layer with (k, k) kernels, c_{in} input channels and c_{out} output channels:

$$n = k^2 * c_{in} * c_{out}$$

So a layer with (5,5) kernels has $\frac{25}{9}$ more parameters than a layer with (3,3) kernels

Kernel sizes

Consider the number of parameters that are needed for a convolutional layer with (k, k) kernels, c_{in} input channels and c_{out} output channels:

$$n = k^2 * c_{in} * c_{out}$$

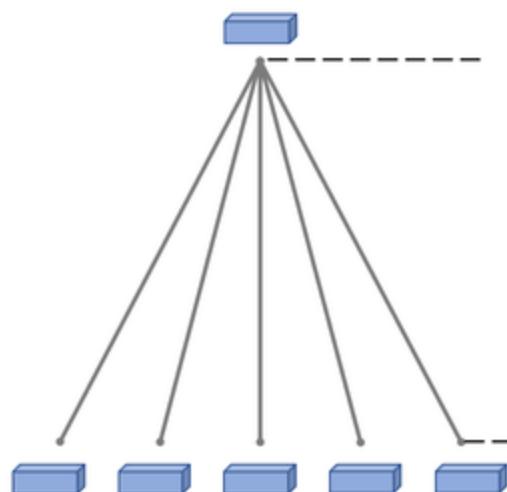
So a layer with $(5,5)$ kernels has $\frac{25}{9}$ more parameters than a layer with $(3,3)$ kernels

If we consider the outputs of these layers, the $(5,5)$ reduces 5 input pixels to 1 output pixel—reducing the size by a factor of 5. The $(3,3)$ kernels reduce by a factor of 3—so we get 2 or 3 output pixels, depending on padding.

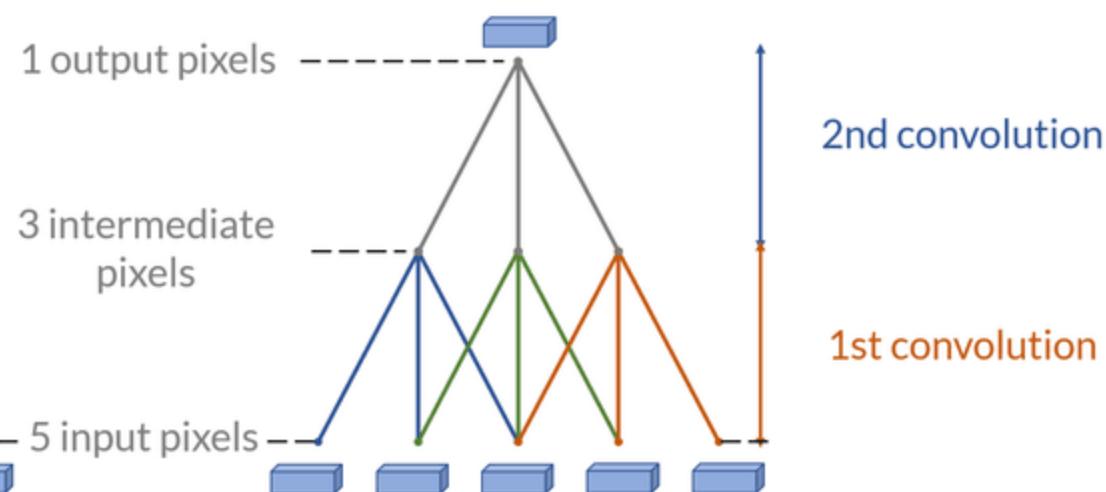
But we can stack the $(3,3)$ layers to get the same effect as a $(5,5)$ layer, with less parameters...

Kernel Sizes

5x5 convolution



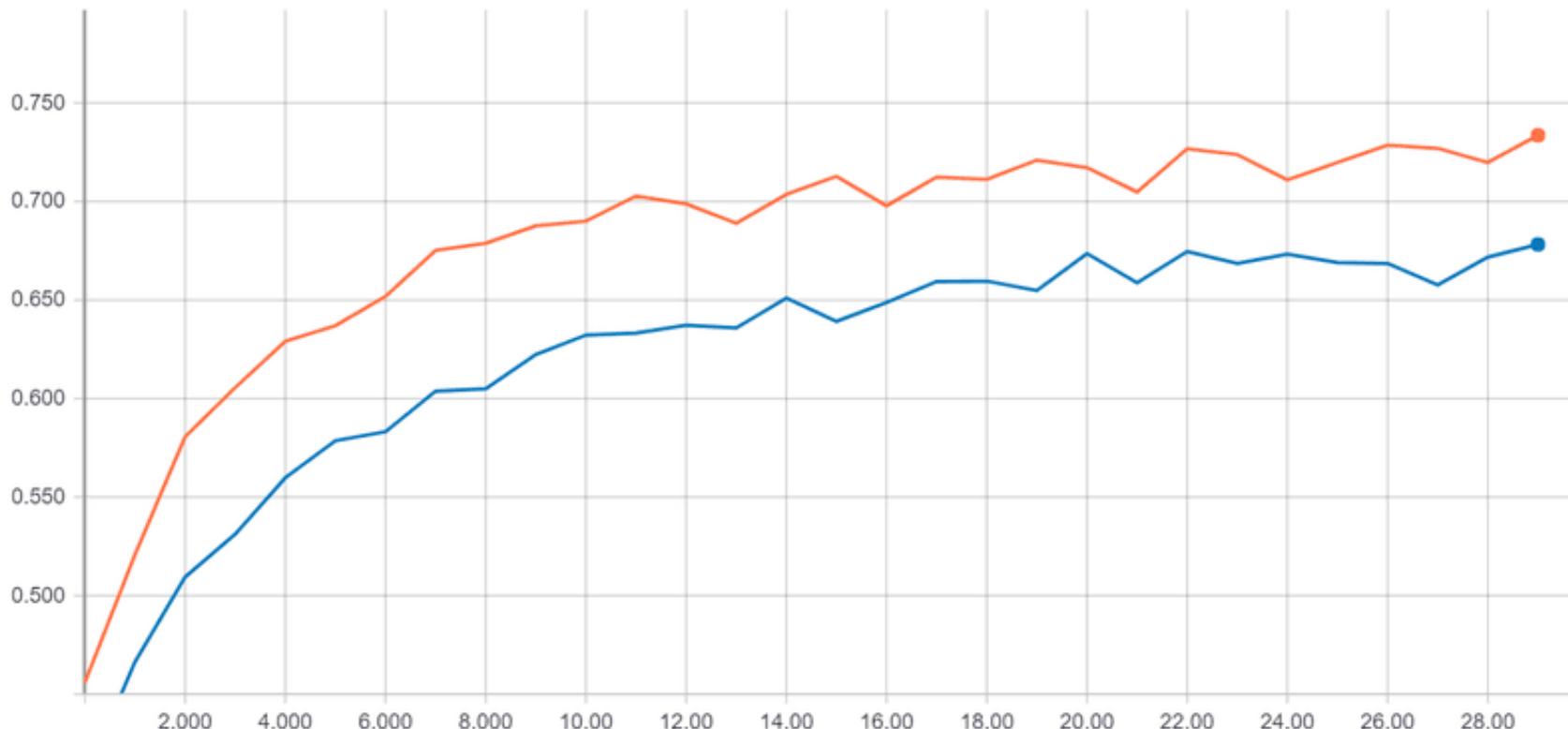
Stacked 3x3 convolutions



$$n = 5^2 * c_{in} * c_{out}$$

$$n = 2 * 3^2 * c_{in} * c_{out}$$

Smaller kernels are lighter... and better



Using stacked layers with smaller kernels uses fewer parameters in total, but often yields better results.

Key development in CNN Architectures...

Validation Accuracy on a 3x3-based Convnet (orange) and the equivalent 5x5-based Convnet (blue)
<https://www.sicara.fr/blog-technique/2019-10-31-convolutional-layer-convolution-kernel>

Strides vs Pooling

STRIDED CONVOLUTIONS

Downsamples the image at the same time as performing the convolution.

Useful when you have large input images and want to downscale them quickly.

Larger Kernels can produce the same receptive fields as multiple layers with small kernels.

POOLING

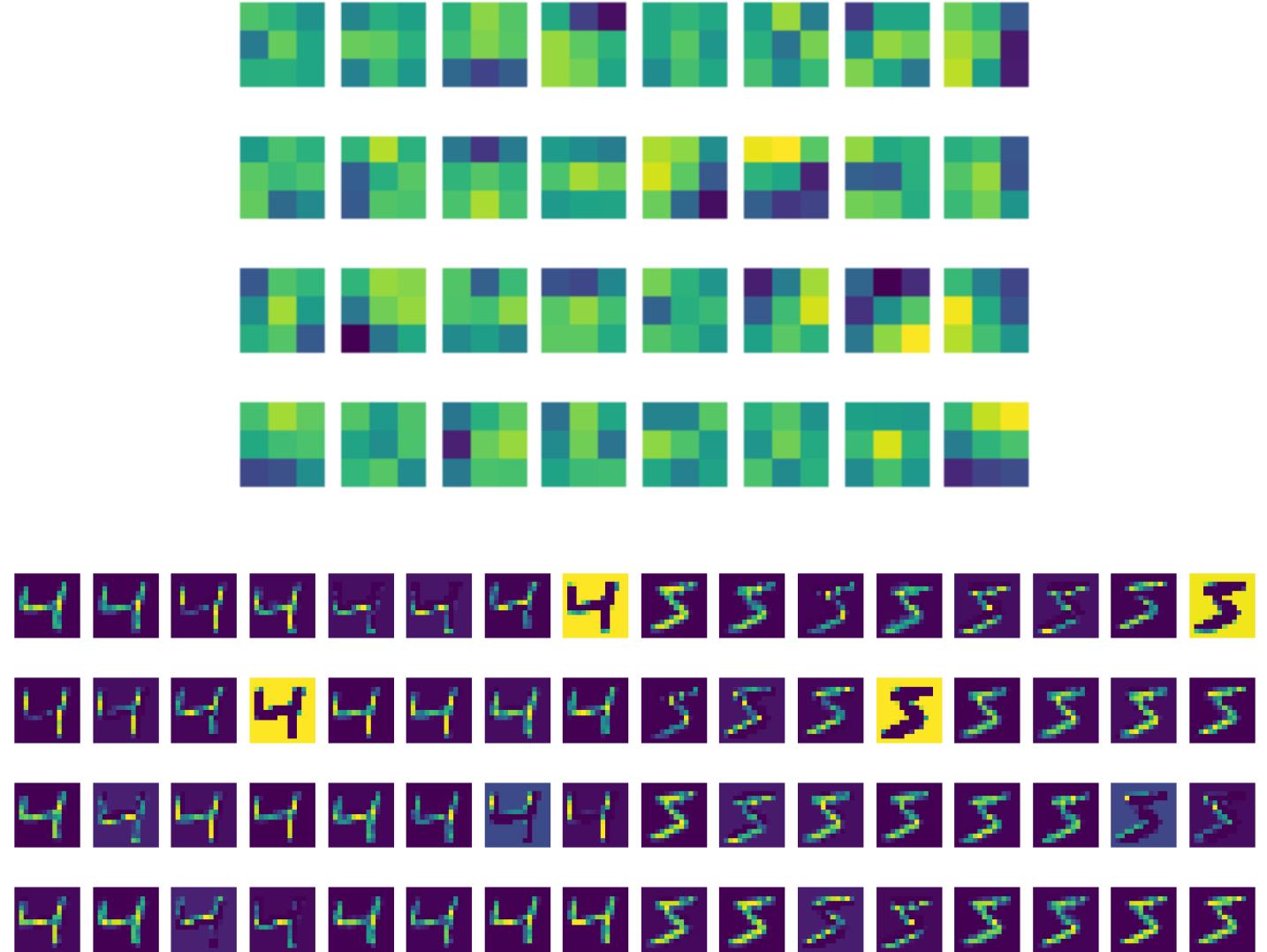
Uses max pooling, average pooling, and other techniques to reduce the size of an input image.

Computationally cheap, both in terms of speed and parameters (none for a pooling layer).

No difficulties with propagating gradients.

Summary: Convolutional Layers

- Convolutional layers are a computers eye. They learn filters which highlight and enhance different features in an image—like edges
- Designing CNNs is tricky business, we have to consider the number of filters, the kernel size, and strides vs pooling



VGGNet

ILSVRC 2014

In 2014, there was a major shift in CNN design, as seen in the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in that year.

Up to this point, the winners of this challenge had large kernels in their initial layers—(11,11) in the case of AlexNet in 2012.

As we have seen, using smaller kernel sizes can lead to lighter models and better results, and this was shown by the 2nd place model in ILSVRC 2014.

VGGNet used only (3,3) kernels, while most other models had (7,7) and (5,5) kernels, at least in their initial layers.

AlexNet: <https://papers.nips.cc/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-Abstract.html>

Depth in VGGNet

There are multiple versions of VGGNet, with a varying number of layers. The most popular have 16 and 19 layers.

Ideally, adding more layers would lead to a better model—because we would have more parameters with which to learn with.

So why not use a VGGNet50 or VGGNet100?

Depth in VGGNet

There are multiple versions of VGGNet, with a varying number of layers. The most popular have 16 and 19 layers.

Ideally, adding more layers would lead to a better model—because we would have more parameters with which to learn with.

So why not use a VGGNet50 or VGGNet100?

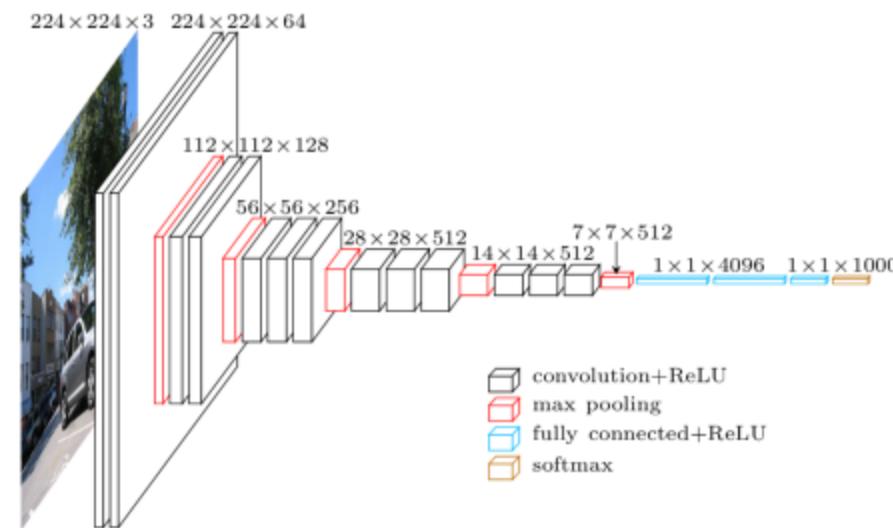
If we keep adding layers to the network, we run into problems with backpropagation. Namely, the Vanishing Gradient Problem.

The ReLU activations in VGGNet can suffer from vanishing gradients, especially in very deep networks.

With too many layers, it becomes impossible to train early convolutions.

Summary: VGGNet

- VGGNet, created by the Visual Geometry Group, consists of blocks of convolutional layers, each with a max pooling layer.
- The network uses small kernel sizes to reduce the number of parameters, while maintaining the receptive field.
- VGGNet suffers in terms of depth, more layers do not markedly improve the network and can make it perform worse.



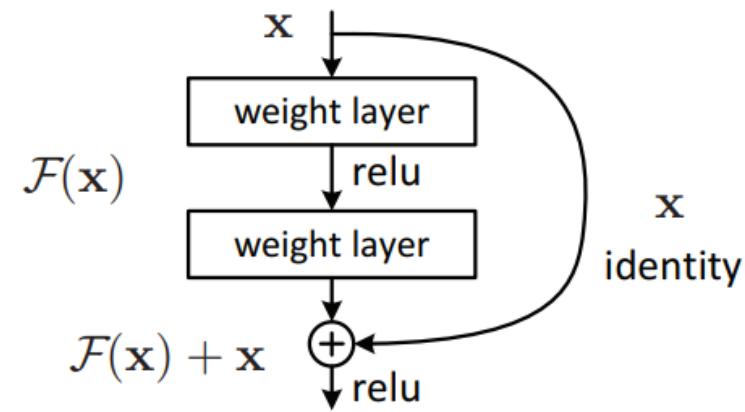
Deep Residual Networks

ResNet – ILSVRC 2015

After VGGNet took the computer vision field by storm, it was quickly unseated as the top dog by ResNet—Deep Residual Networks.

ResNet sets out to tackle the key problem with earlier CNNs, the difficulty in adding deeper layers to the network, without the gradients vanishing.

It does this by learning a Residual Function in its layers.

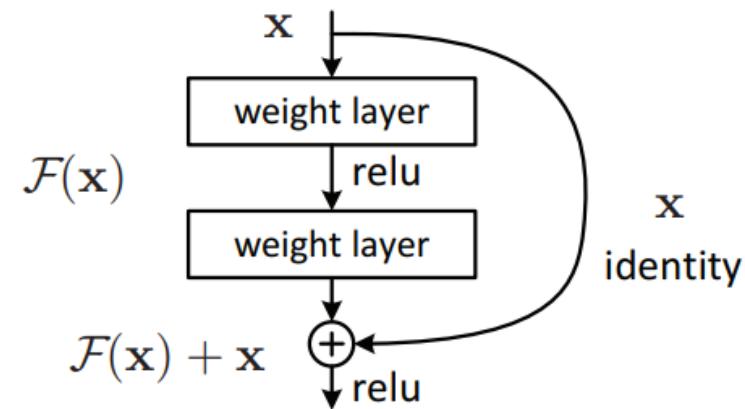


Skip Connections

The key tool in ResNets Architecture is the skip connection.

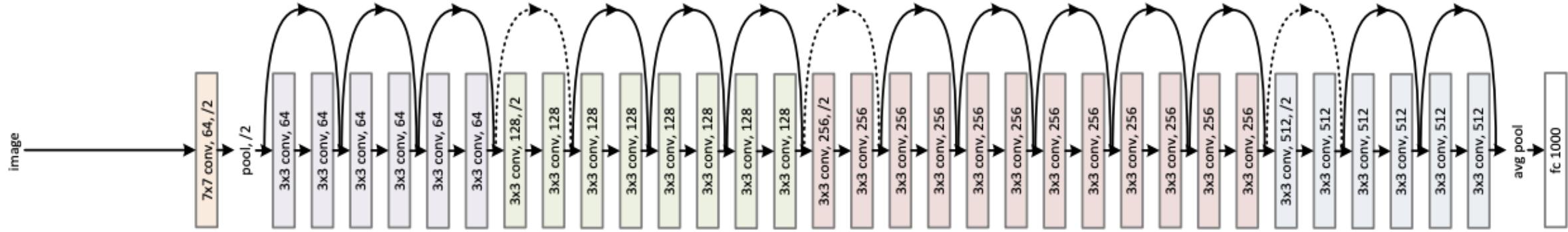
Instead of learning the function $H(x)$ to map the input to the receptive field, it learns the **residual function**, $F(x) = H(x) - x$, with two stacked convolutional layers.

The input, x , is then passed through an “identity” layer, and added to $F(x)$.

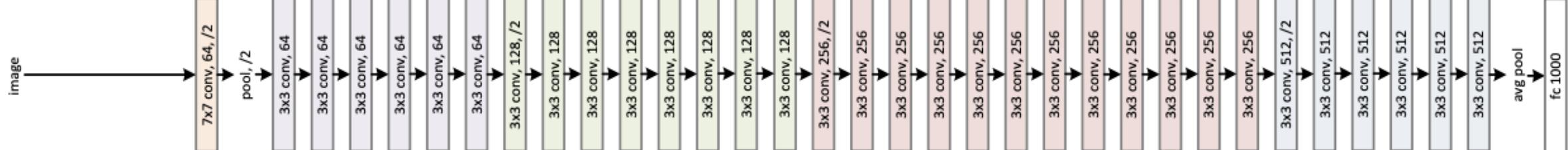


ResNet: <https://arxiv.org/abs/1512.03385>

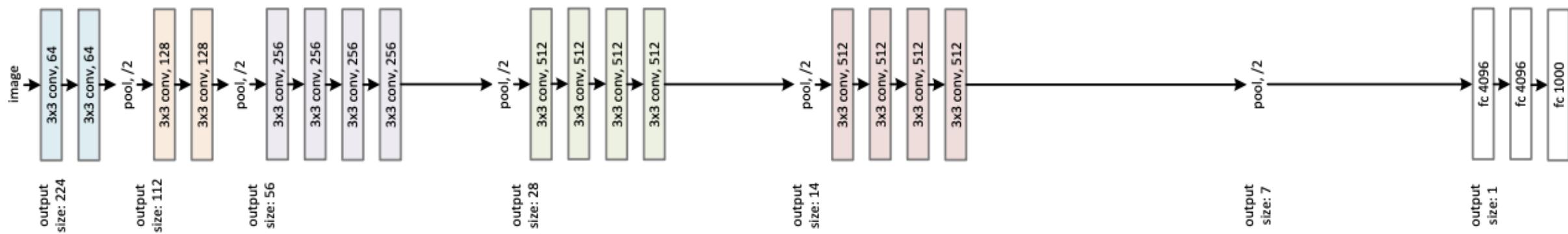
34-layer residual



34-layer plain



VGG-19



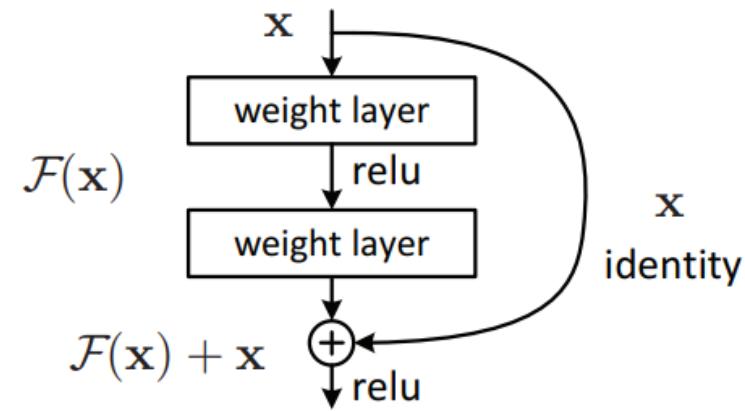
Architecture of ResNet-34 (top), a plain 34-layer network (middle), and VGGNet19 (bottom).
 Solid skip connections are identity mappings, dotted lines change the dimensions.

Back Propagation in ResNet

The skip connections act as a shortcut for the gradients of the loss function.

By skipping over the convolutions in the deeper layers of the network, the early layers become much easier to train, because the effect of the vanishing gradients is greatly reduced.

ResNet also has significantly fewer parameters than VGGNet (~3.6 Billion FLOPs vs. ~19.6 Billion FLOPs).



Summary: Deep Residual Networks

- ResNet utilises skip connections to learn residual mapping between layers of the network.
- The skip connections allow the network to be significantly deeper than previous models, due to the reduced vanishing gradients.
- ResNet gains in performance to much greater depths than previous models, with models from ResNet-34 and ResNet-50, to ResNet-100 and ResNet-152. All with fewer parameters than VGGNet19

model	top-1 err.	top-5 err.
VGG-16 [41]	28.07	9.33
GoogLeNet [44]	-	9.15
PReLU-net [13]	24.27	7.38
plain-34	28.54	10.02
ResNet-34 A	25.03	7.76
ResNet-34 B	24.52	7.46
ResNet-34 C	24.19	7.40
ResNet-50	22.85	6.71
ResNet-101	21.75	6.05
ResNet-152	21.43	5.71

Table 3. Error rates (%), **10-crop** testing) on ImageNet validation. VGG-16 is based on our test. ResNet-50/101/152 are of option B that only uses projections for increasing dimensions.

Questions?
