



PyTorch 배우는 딥러닝 입문 (2)

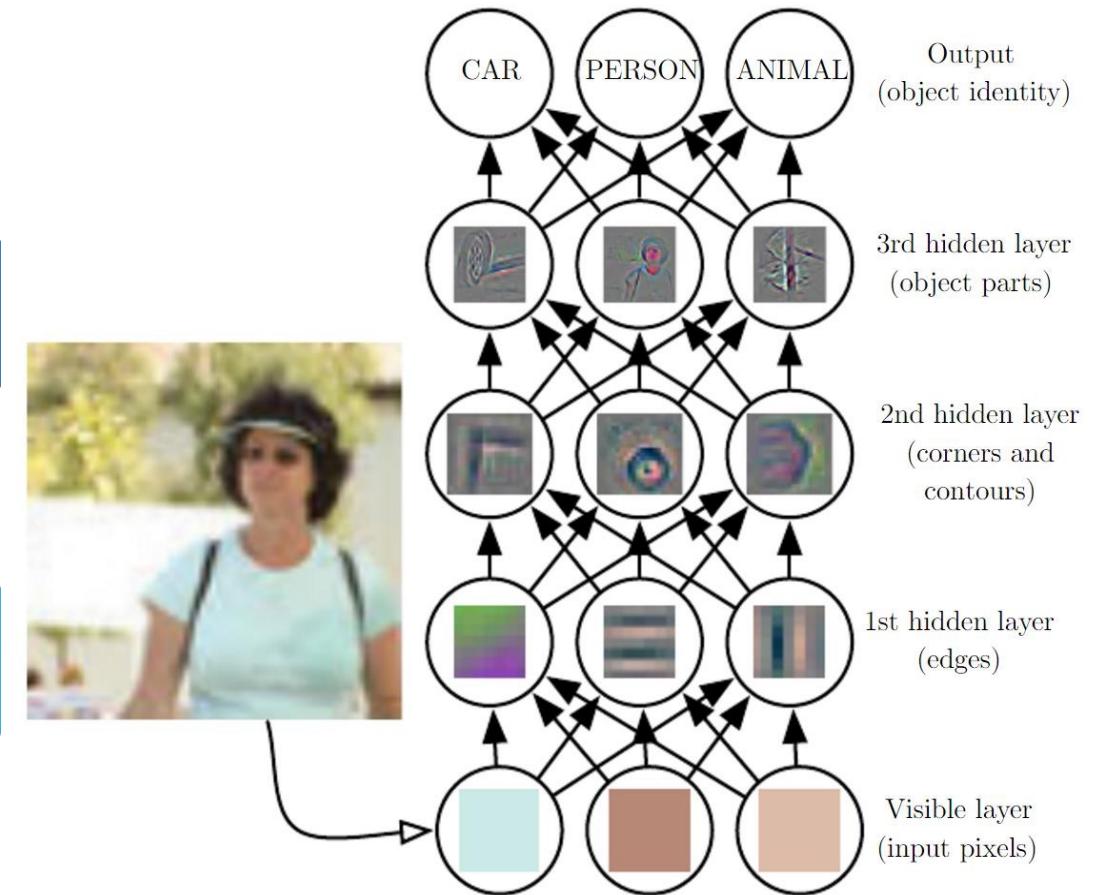
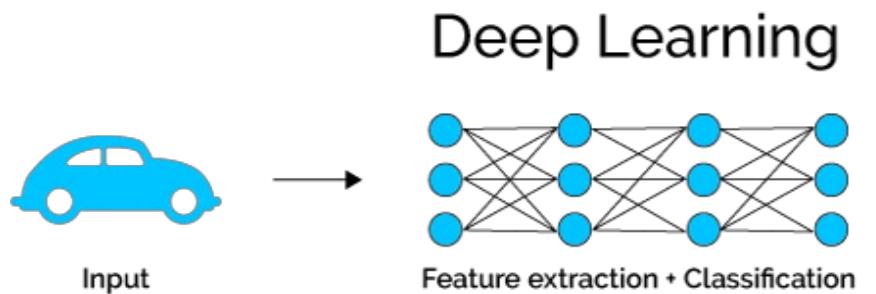
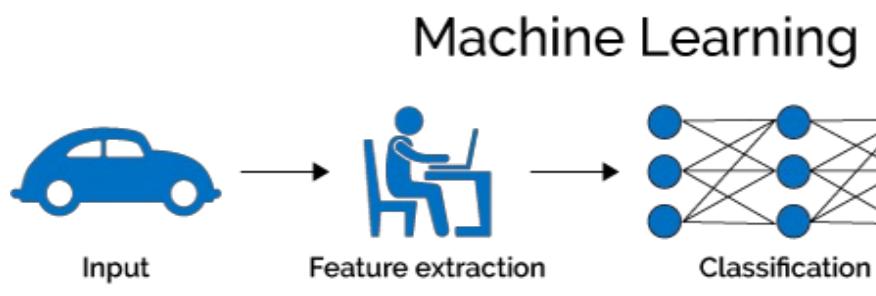
Introduction to Convolutional Neural Network

이영완

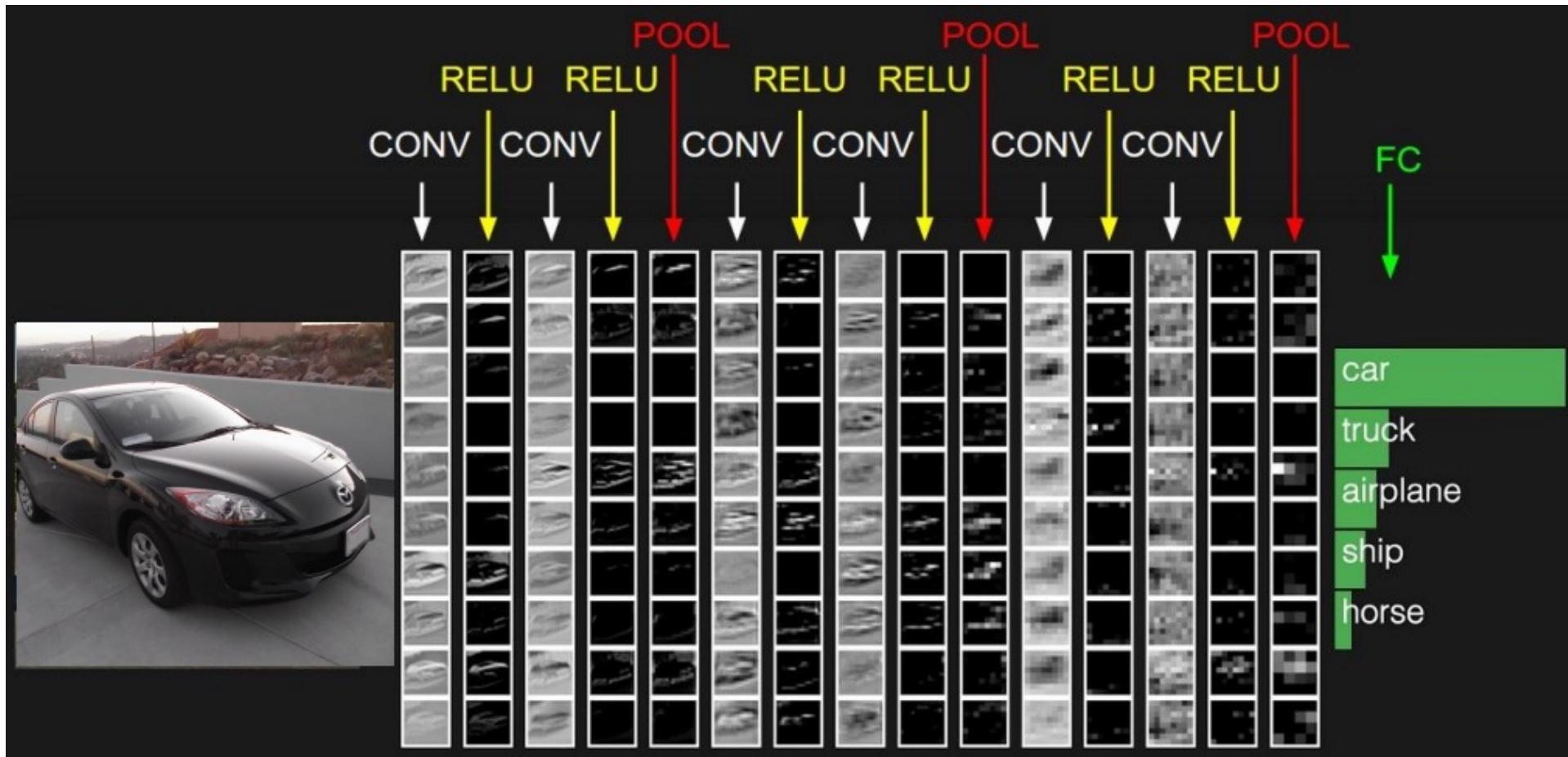
한국전자통신연구원

Deep Learning is Representation Learning

(aka Feature Learning)



Convolution Neural Network (CNN)



CNNs are everywhere

Classification



Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

CNNs are everywhere

Classification



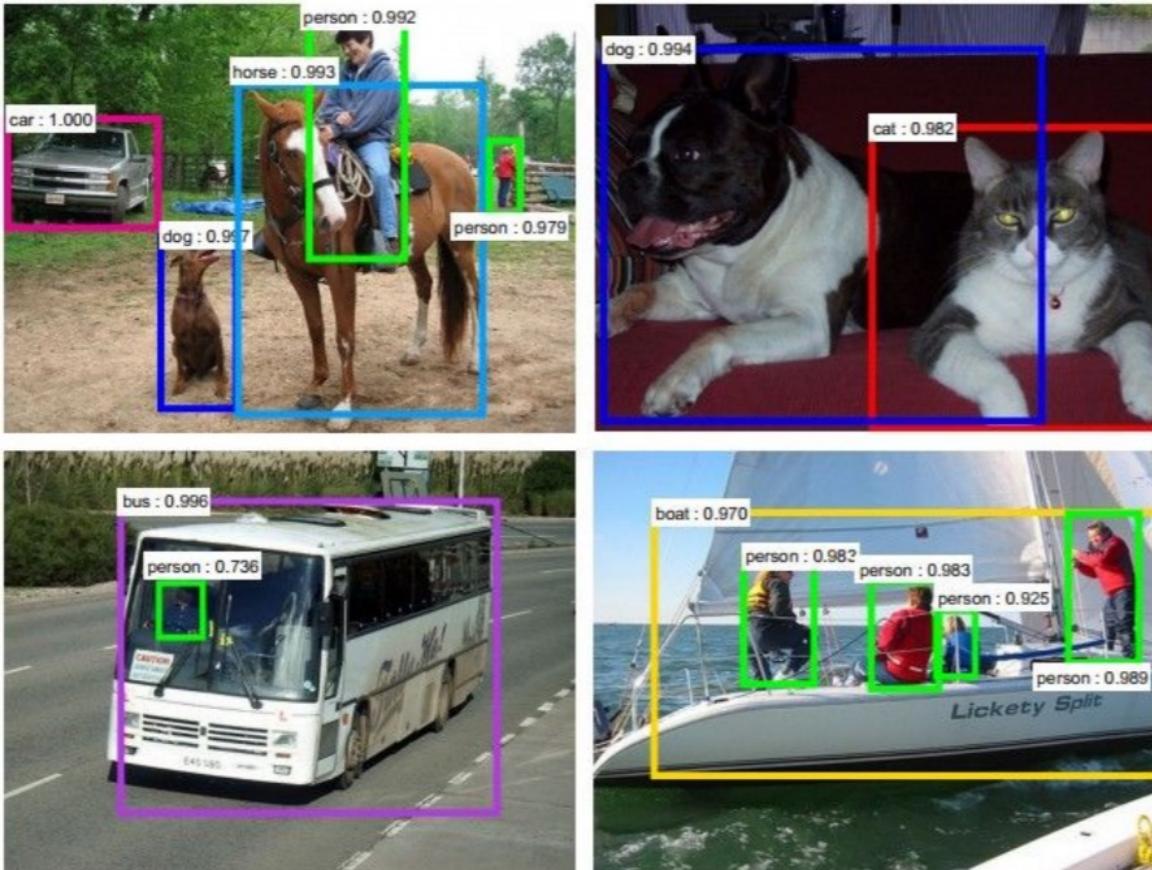
Retrieval



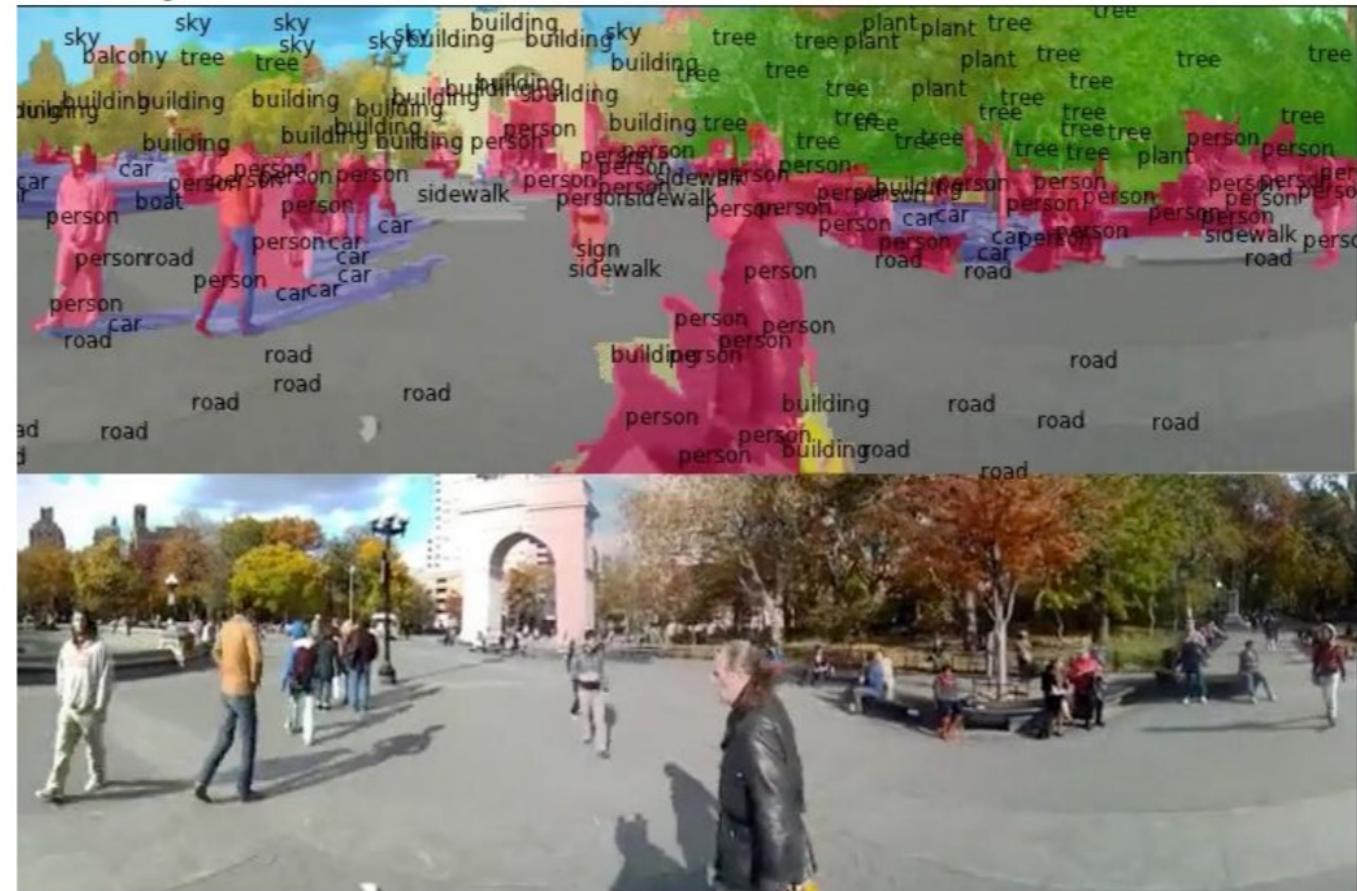
Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

CNNs are everywhere

Detection



Segmentation



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.

[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Figures copyright Clement Farabet, 2012.
Reproduced with permission.

[Farabet et al., 2012]

CNNs are everywhere



self-driving cars



[This image](#) by GBPublic_PR is
licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

(these are the GPUs on rye01.stanford.edu)

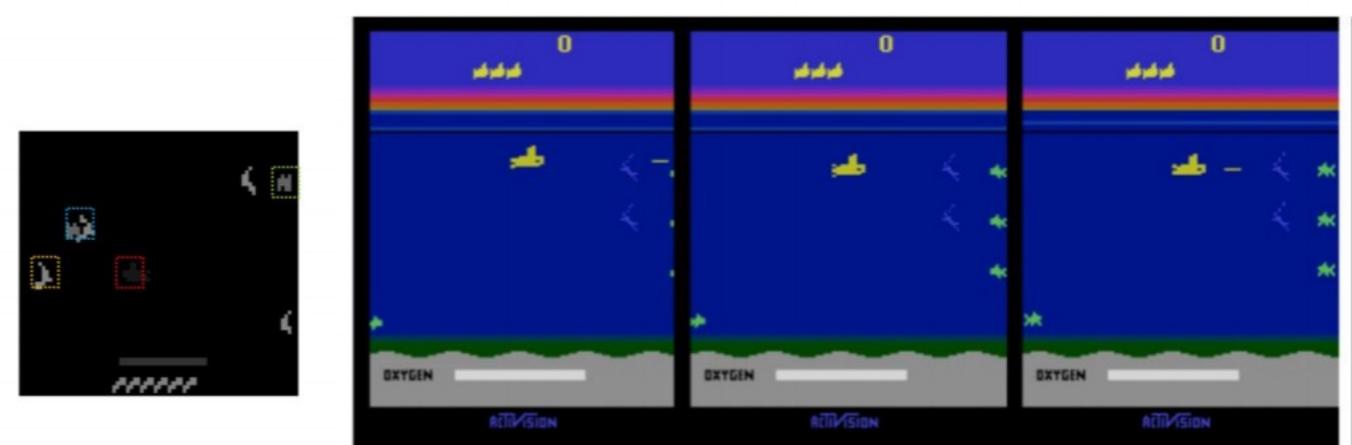
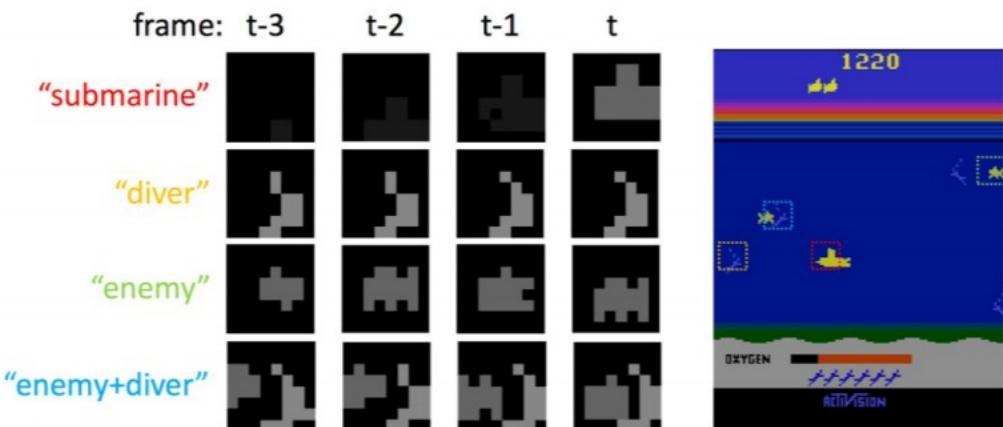
Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

CNNs are everywhere



Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.

[Toshev, Szegedy 2014]



[Guo et al. 2014]

Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reproduced with permission.

CNNs are everywhere

No errors



A white teddy bear sitting in the grass



A man riding a wave on top of a surfboard

Minor errors



A man in a baseball uniform throwing a ball



A cat sitting on a suitcase on the floor

Somewhat related



A woman is holding a cat in her hand



A woman standing on a beach holding a surfboard

Image Captioning

[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]

All images are CC0 Public domain:
<https://pixabay.com/en/luggage-antique-cat-1643010/>
<https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/>
<https://pixabay.com/en/surf-wave-summer-sport-litoral-1668716/>
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>
<https://pixabay.com/en/handstand-lake-meditation-496008/>
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

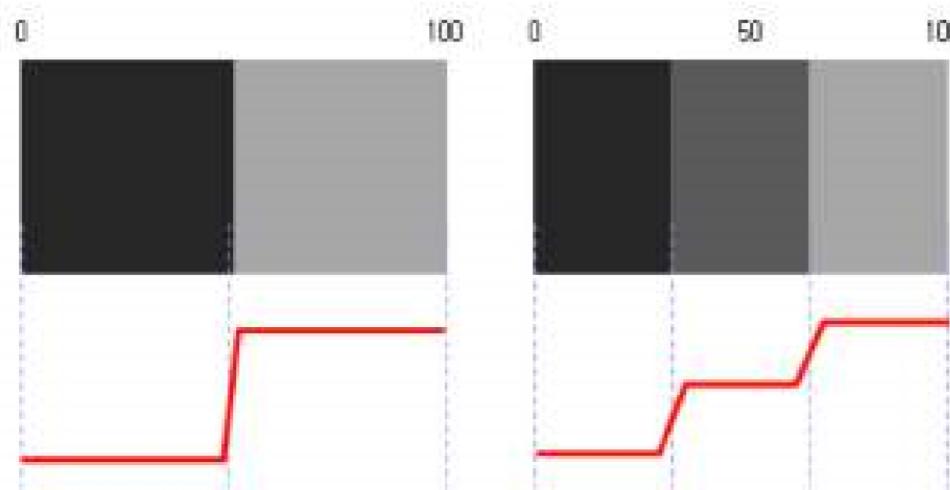
Captions generated by Justin Johnson using [Neuraltalk2](#)

Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.

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

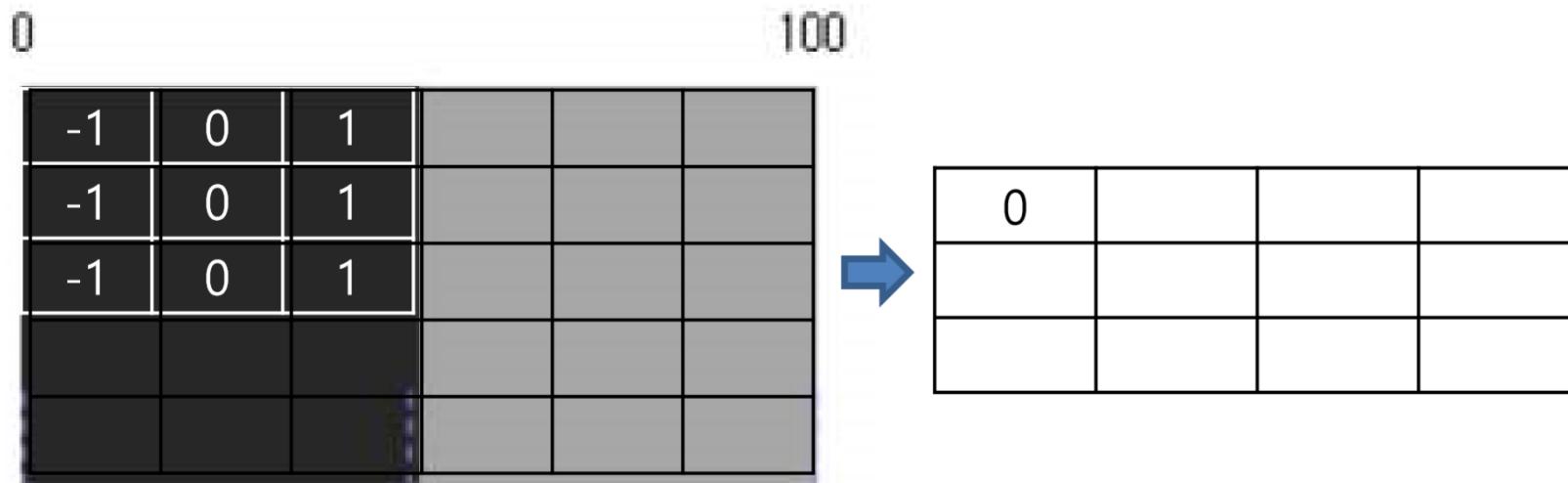
수직 에지를 찾는
컨볼루션 필터



Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.

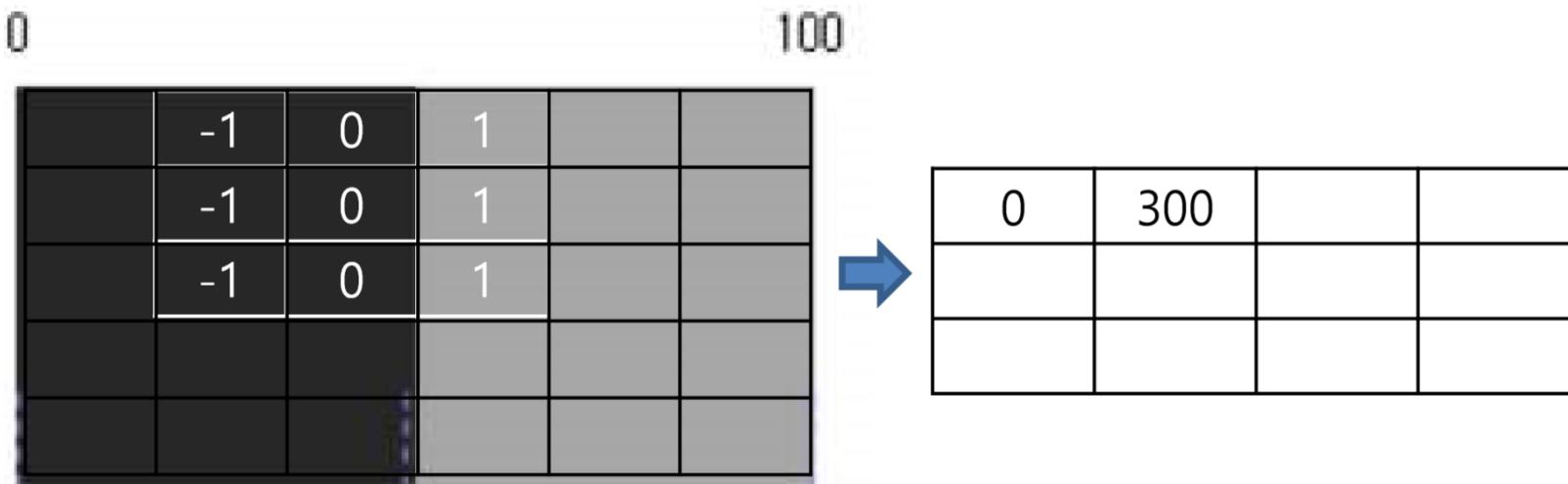
수직 엣지를 찾는
컨볼루션 필터



Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.

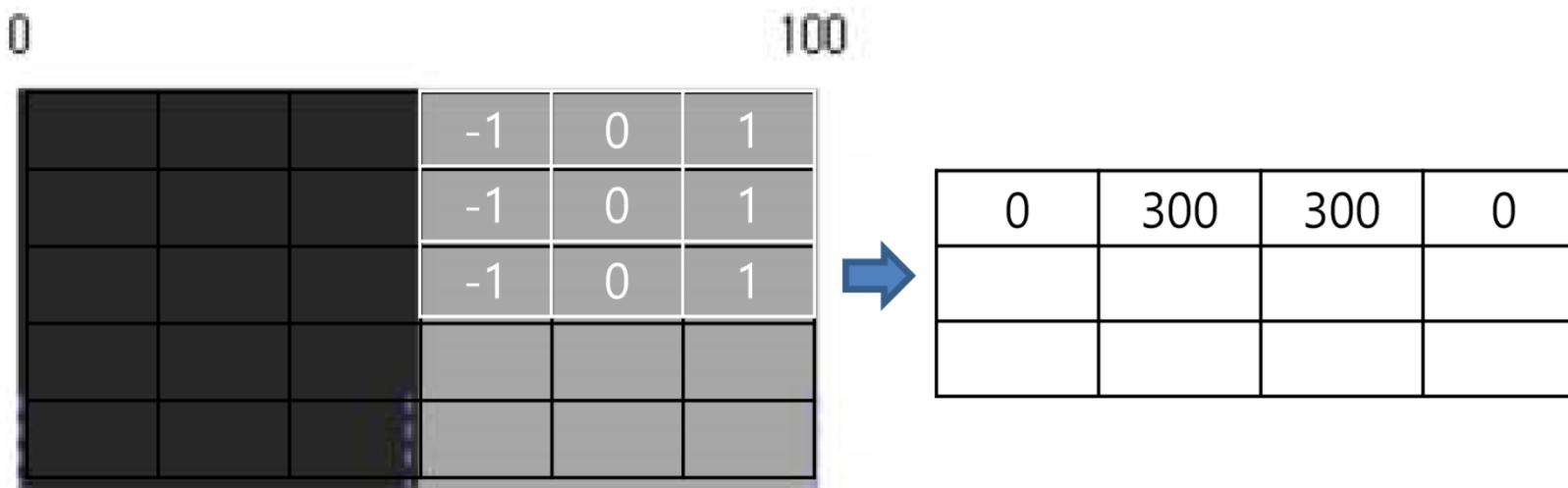
수직 엣지를 찾는
컨볼루션 필터



Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.

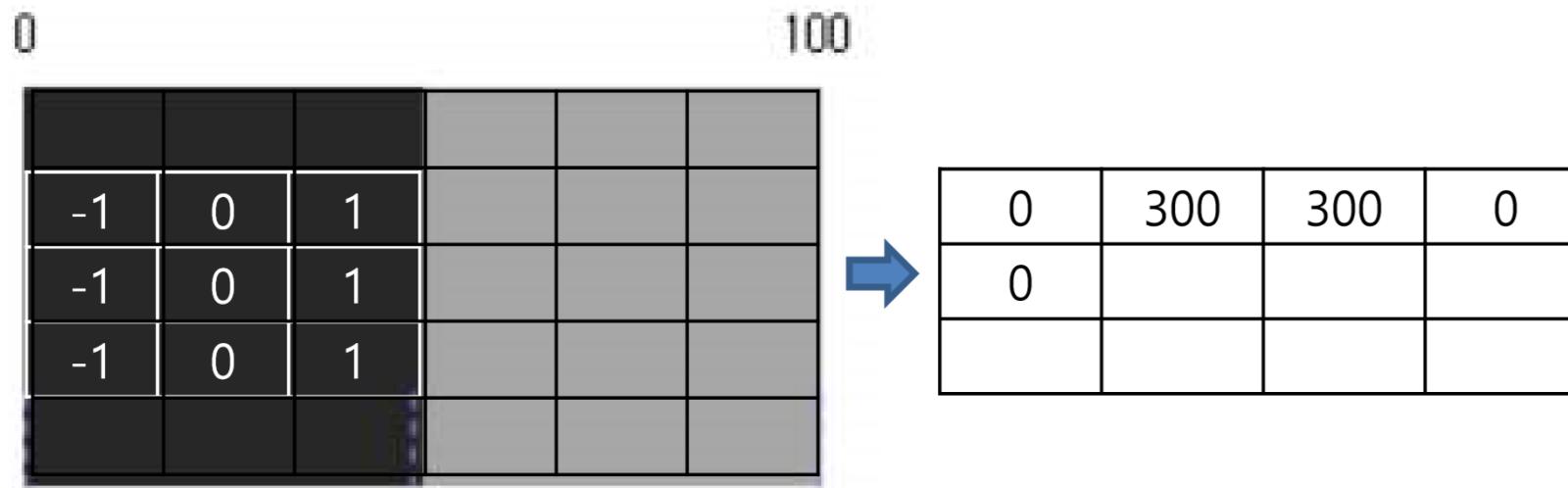
수직 엣지를 찾는
컨볼루션 필터



Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(feature)을 추출할 수 있다.

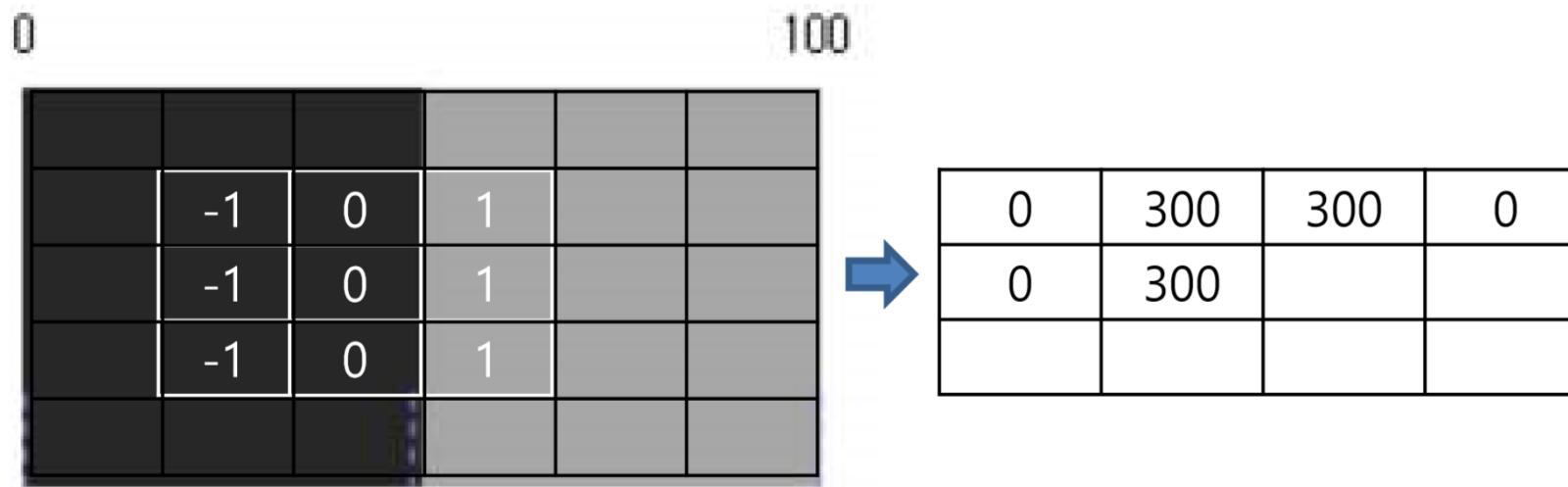
수직 엣지를 찾는 컨볼루션 필터



Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 특징(feature)을 추출할 수 있다.

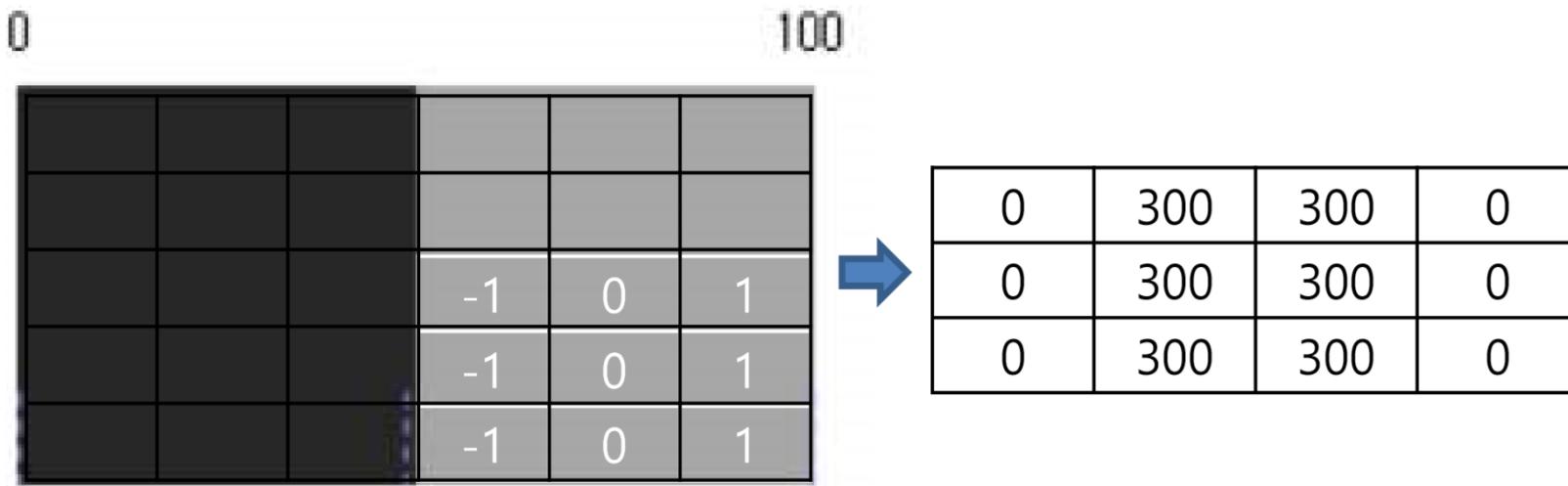
수직 엣지를 찾는 컨볼루션 필터



Convolutional Neural Network(CNN)

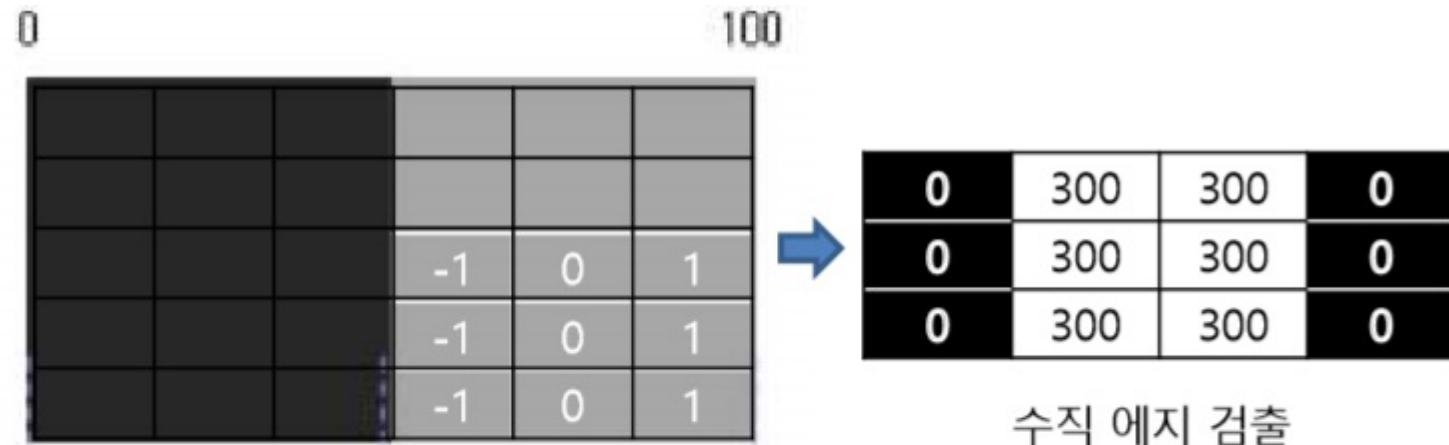
Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.

수직 엣지를 찾는
컨볼루션 필터



Convolutional Neural Network(CNN)

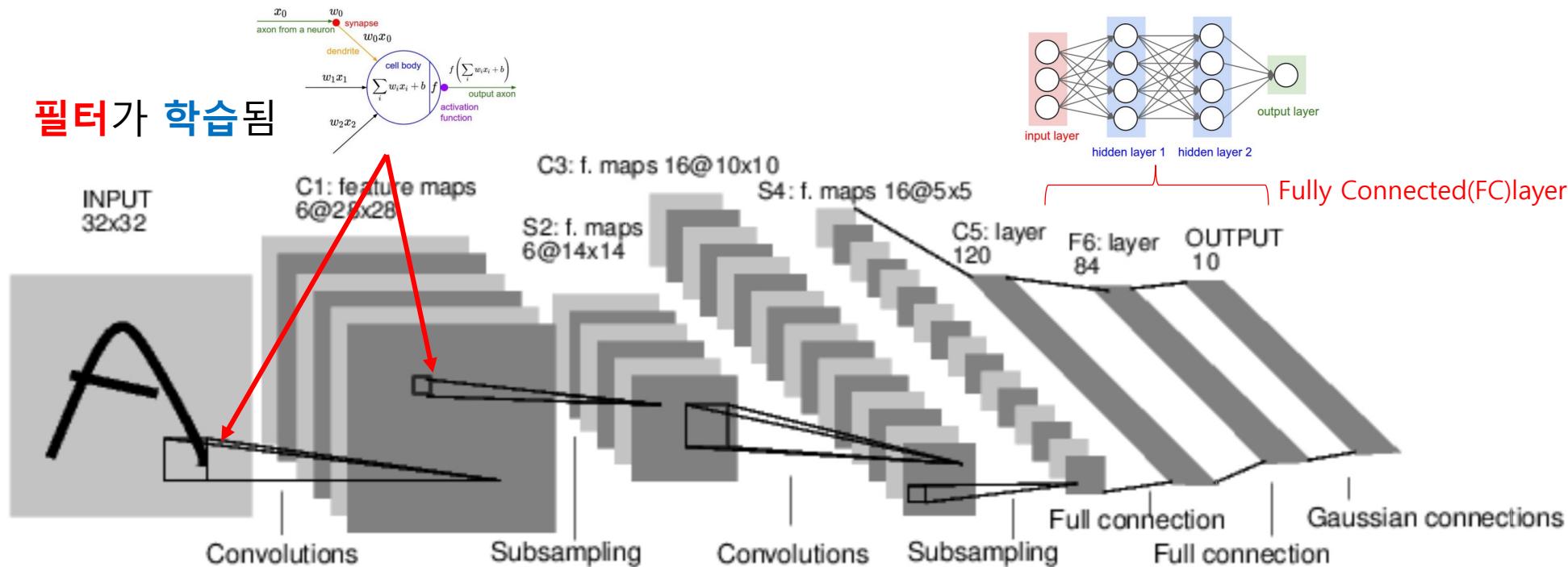
Convolution Filter : 이미지의 특징(**feature**)을 추출할 수 있다.



다양한 필터로 다양한
특징을 추출할 수 있다

Convolutional Neural Network(CNN)

Convolution Filter : 이미지의 **특징(feature)**을 추출할 수 있다.
→ Convolutional Neural Network는 컨볼루션 필터를 **학습**



Visualization

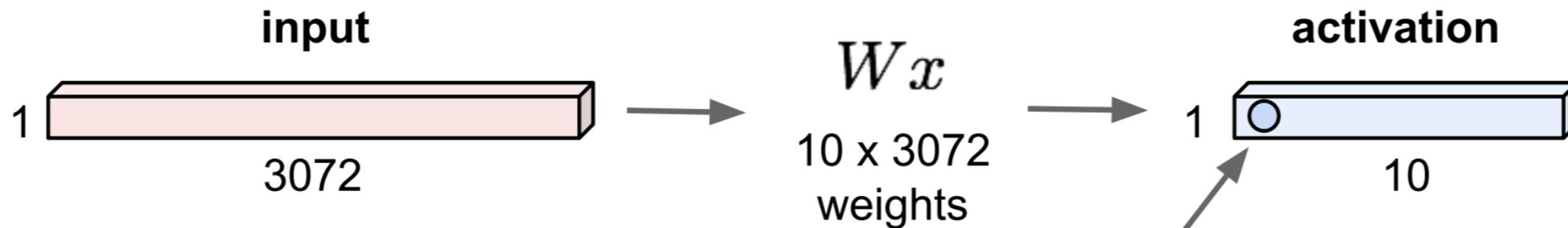
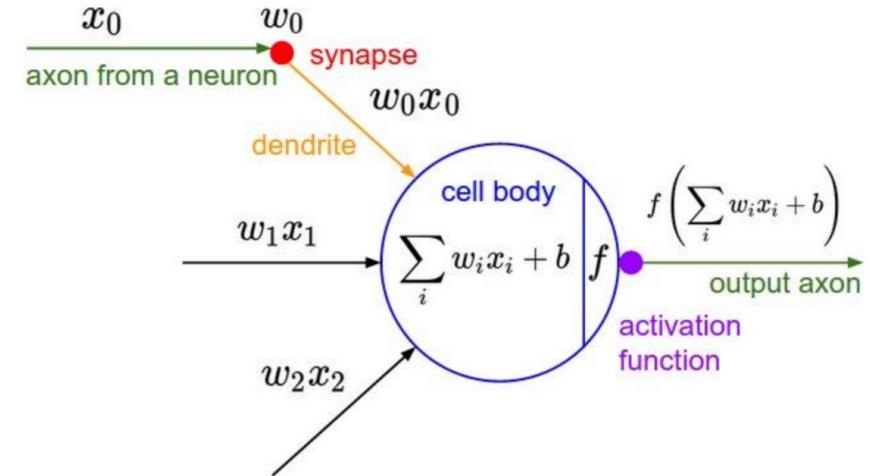
<https://poloclub.github.io/cnn-explainer/>

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Delving deep into CNNs

Fully Connected Layer

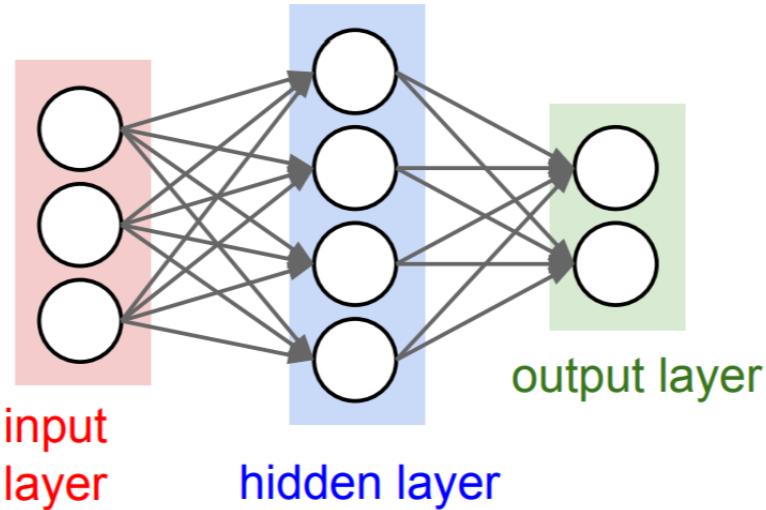
32x32x3 image -> stretch to 3072 x 1



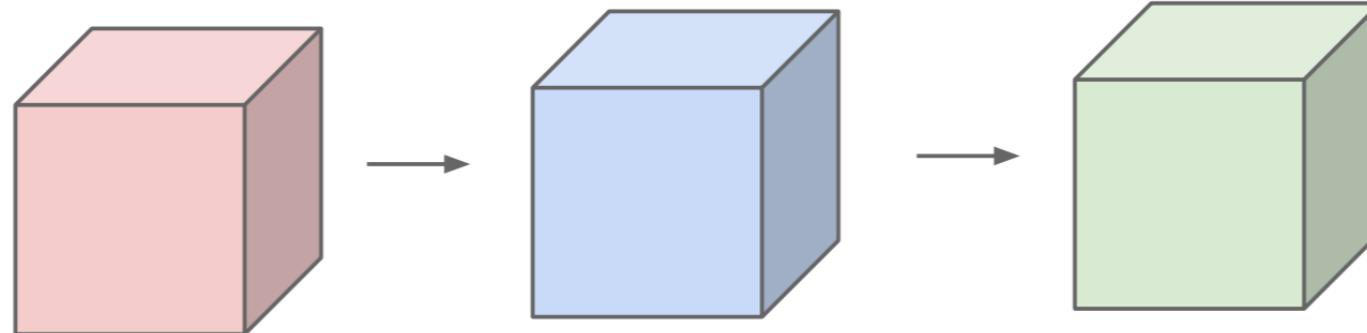
1 number:
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

Fully Connected(FC)

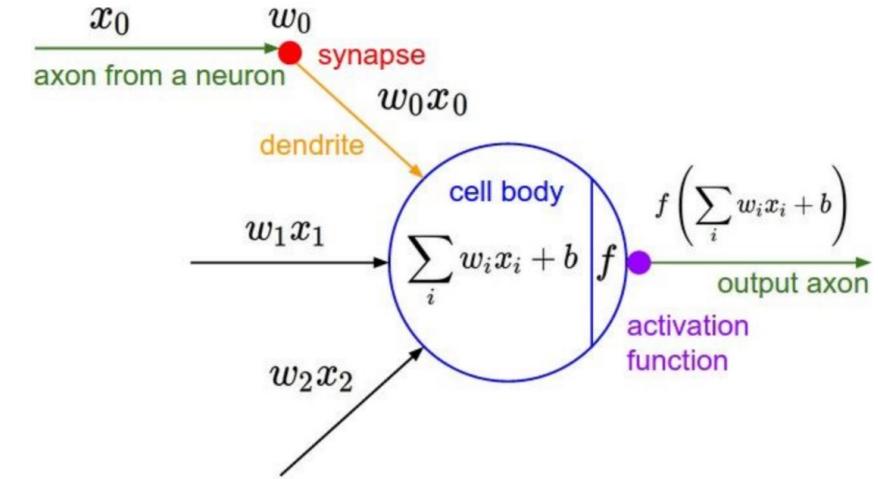
before:



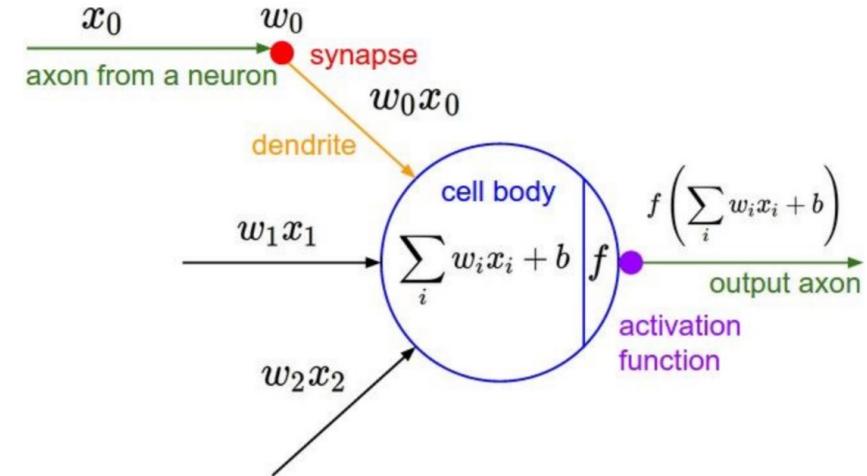
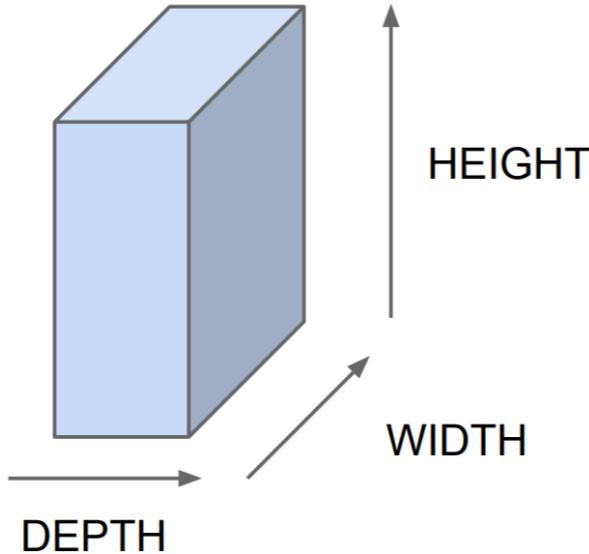
now:



Convolution



All Neural Net activations arranged in 3 dimensions:



For example, a CIFAR-10 image is a 32x32x3 volume
32 width, 32 height, 3 depth (RGB channels)

pytorch Tensor dimension : N x C x H x W
→ N : batch size , C : Channel(depth)

Convolutional Neural Networks are just Neural Networks BUT: **1. Local connectivity**

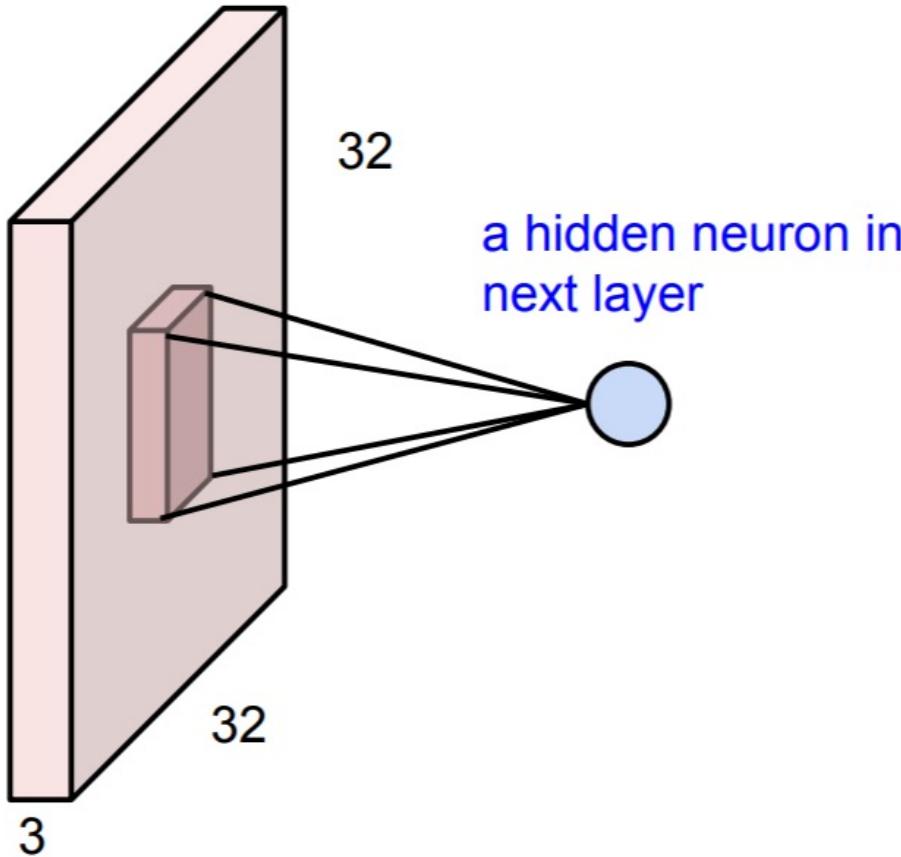


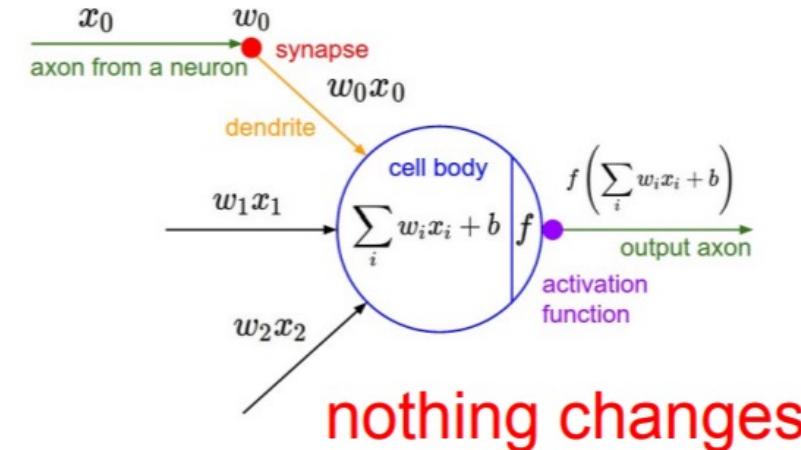
image: 32x32x3 volume

before: full connectivity: 32x32x3 weights

now: one neuron will connect to, e.g. 5x5x3 chunk and only have 5x5x3 weights.

note that connectivity is:

- local in space (5x5 inside 32x32)
- but full in depth (all 3 depth channels)



Convolutional Neural Networks are just Neural Networks BUT: **1. Local connectivity**

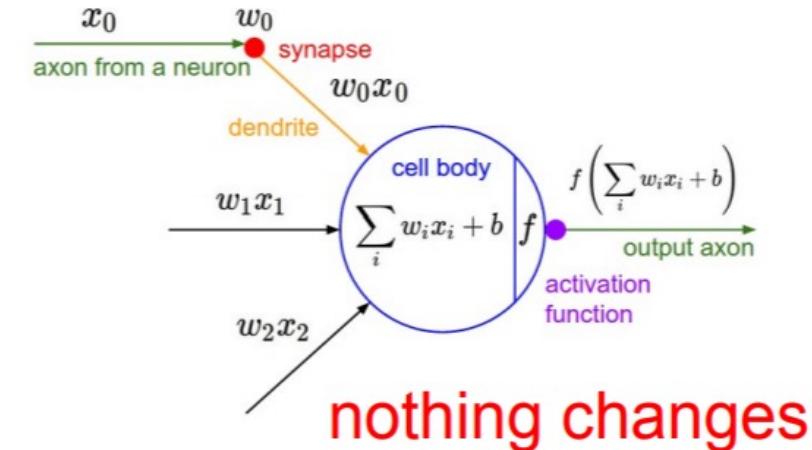
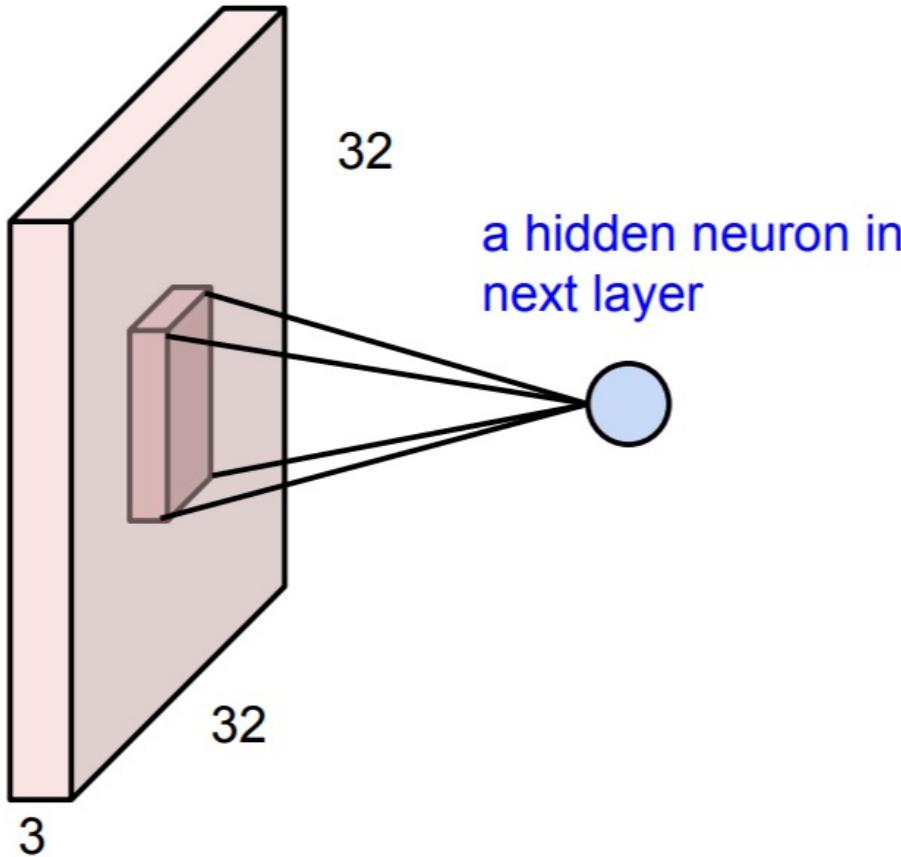


image: 32x32x3 volume

before: full connectivity: 32x32x3 weights

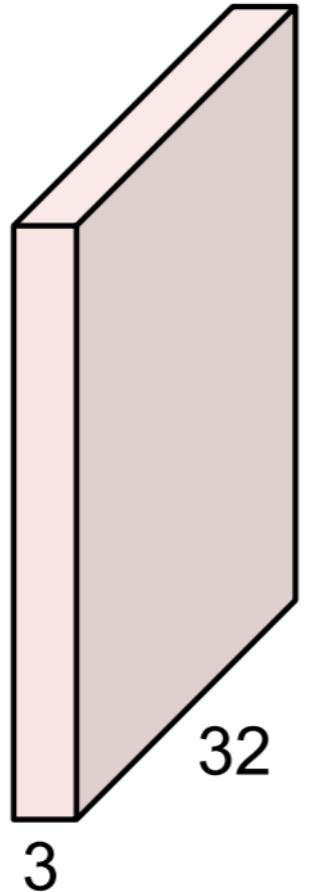
now: one neuron will connect to, e.g. 5x5x3 chunk and only have 5x5x3 weights.

note that connectivity is:

- local in space (5x5 inside 32x32)
- but full in depth (all 3 depth channels)

Convolution Layer

32x32x3 image



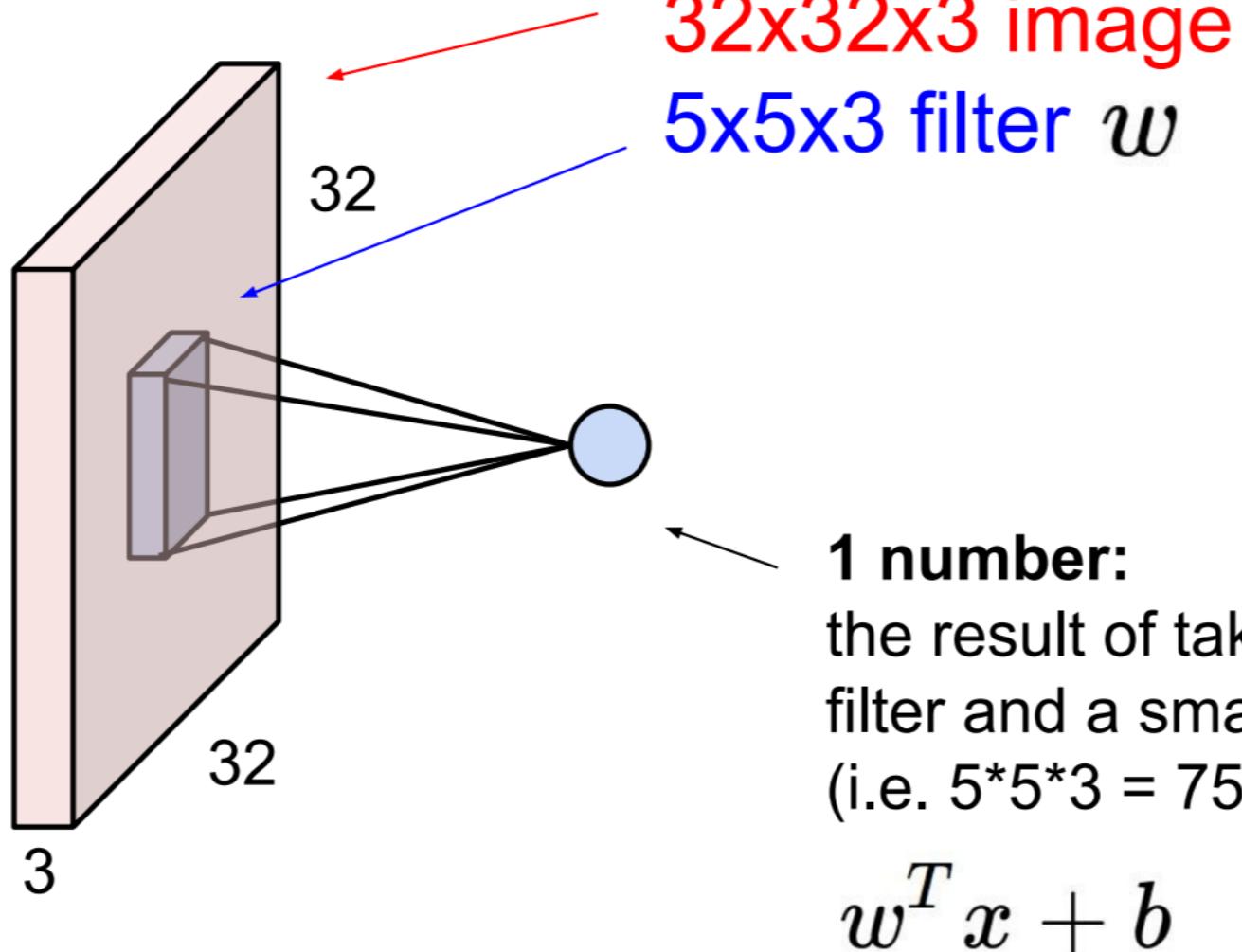
5x5x3 filter



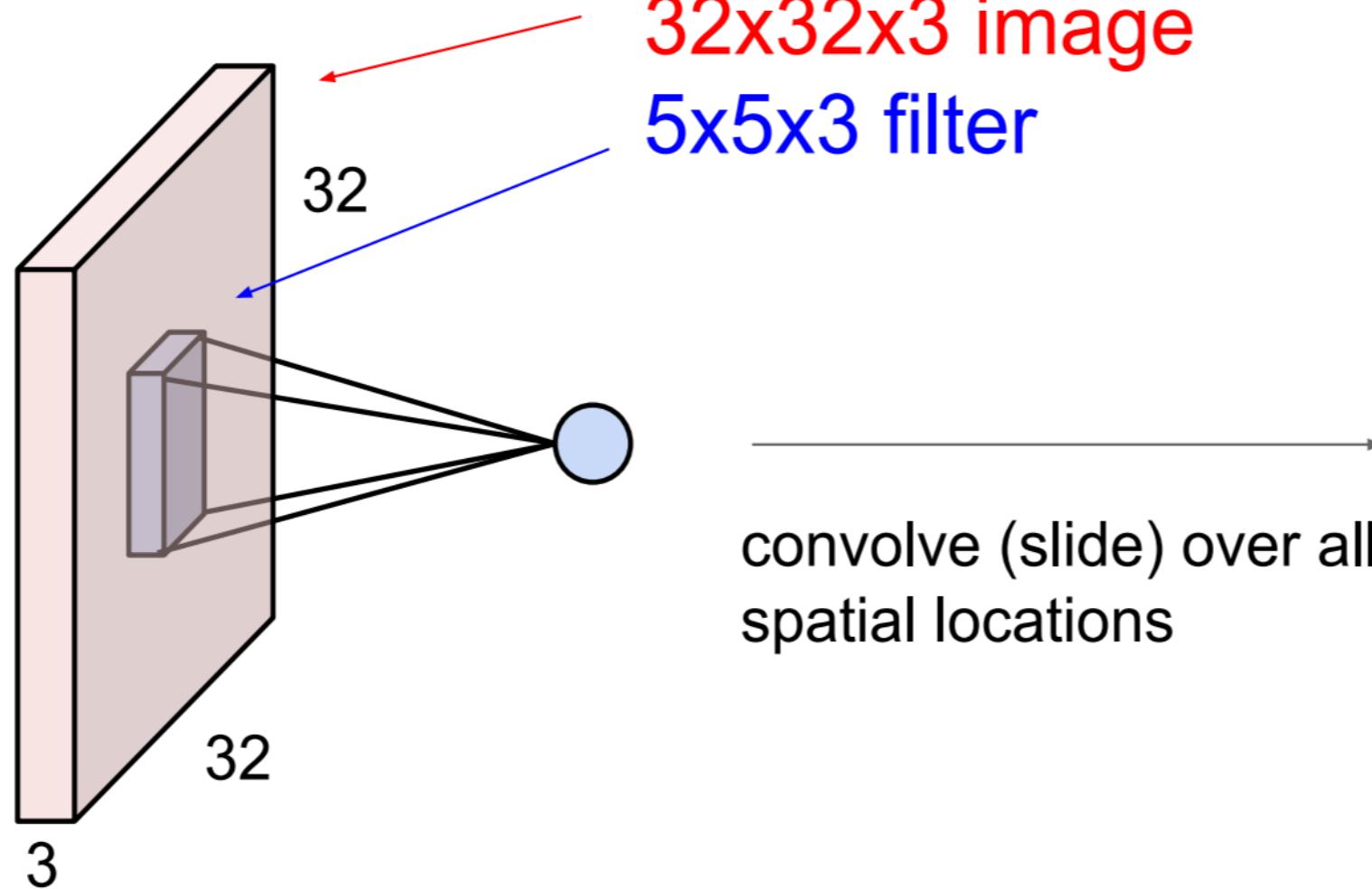
Filters always extend the full depth of the input volume

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

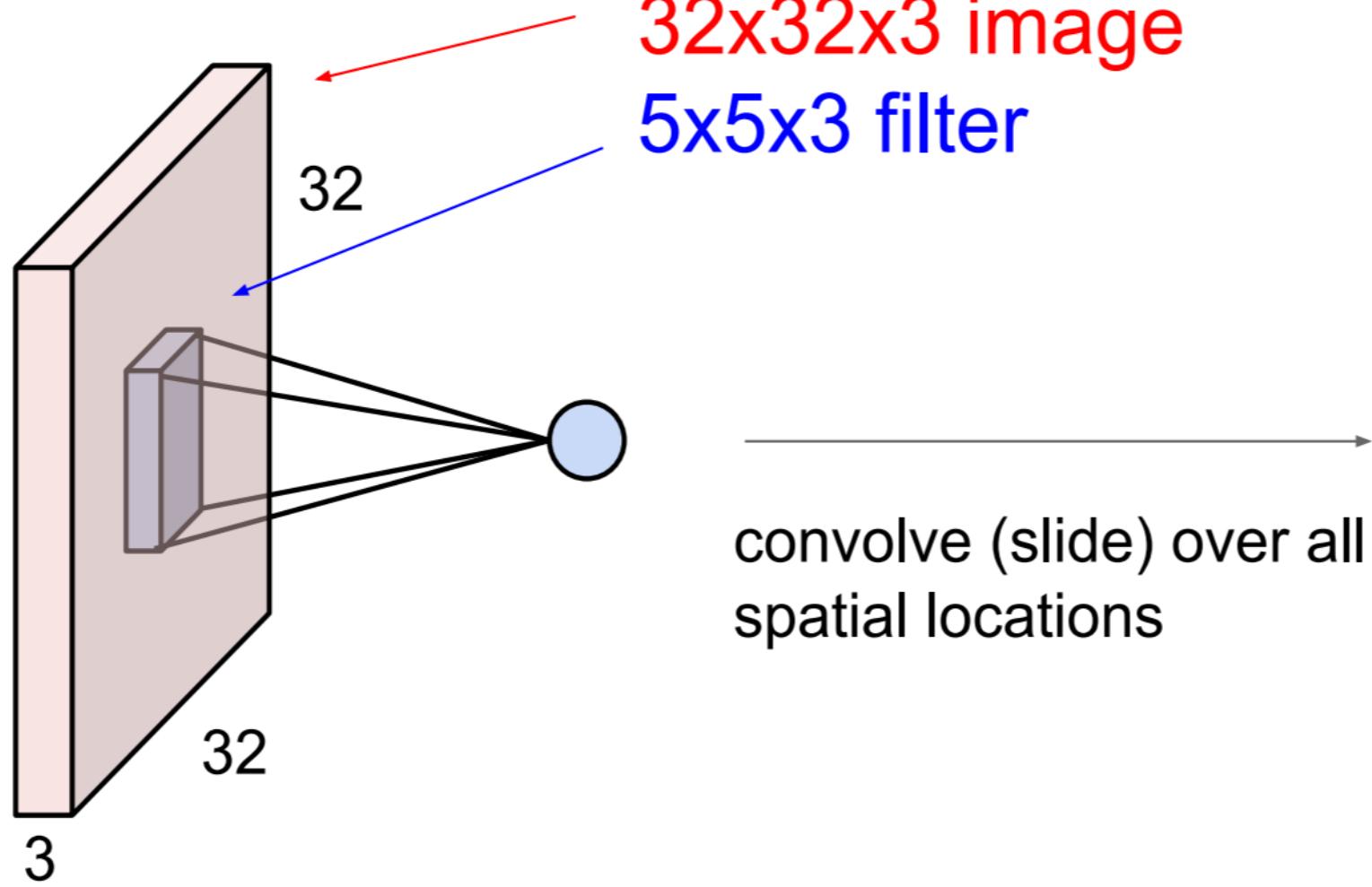
Convolution Layer



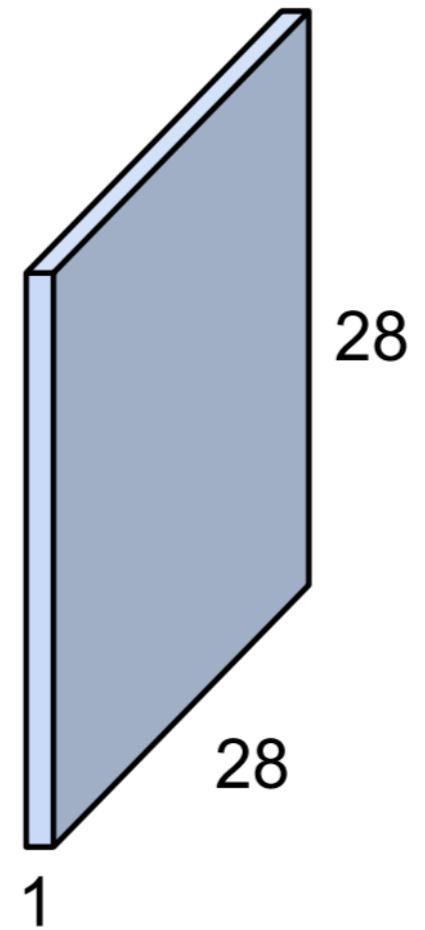
Convolution Layer



Convolution Layer

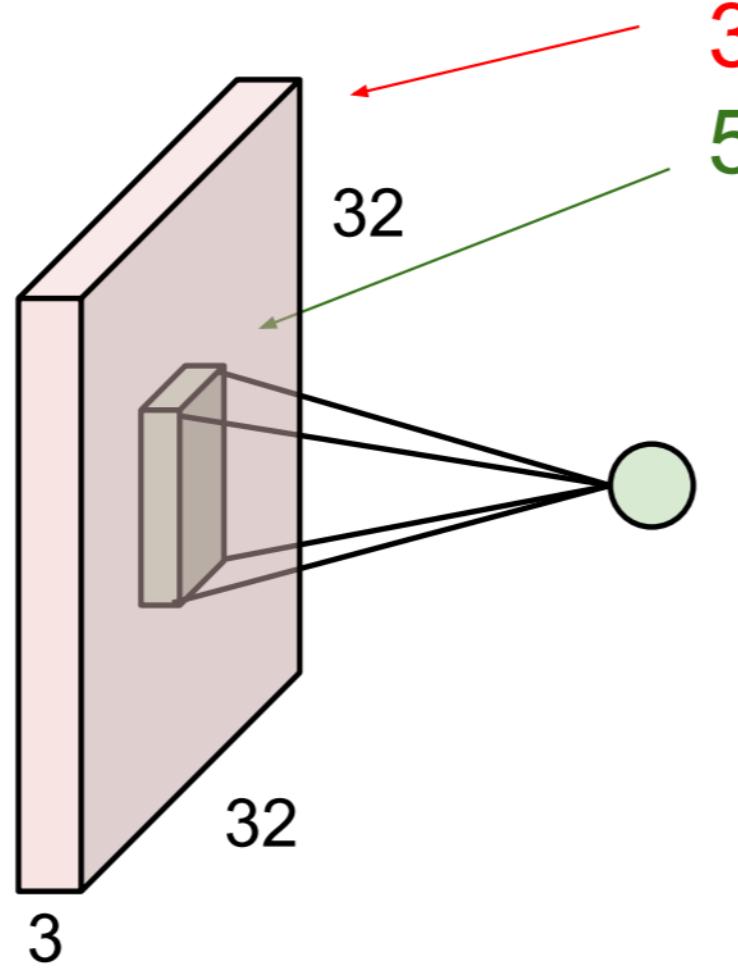


activation map



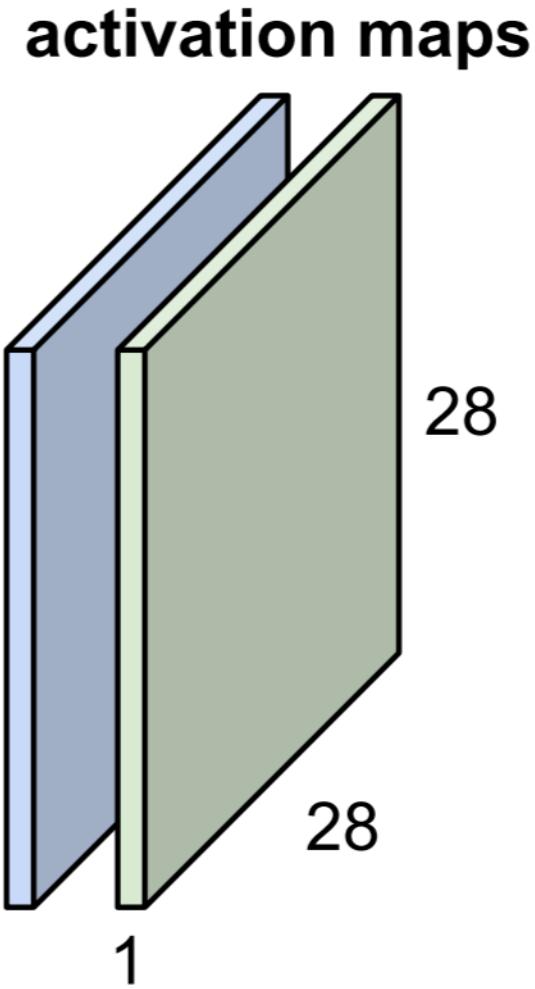
Convolution Layer

consider a second, green filter

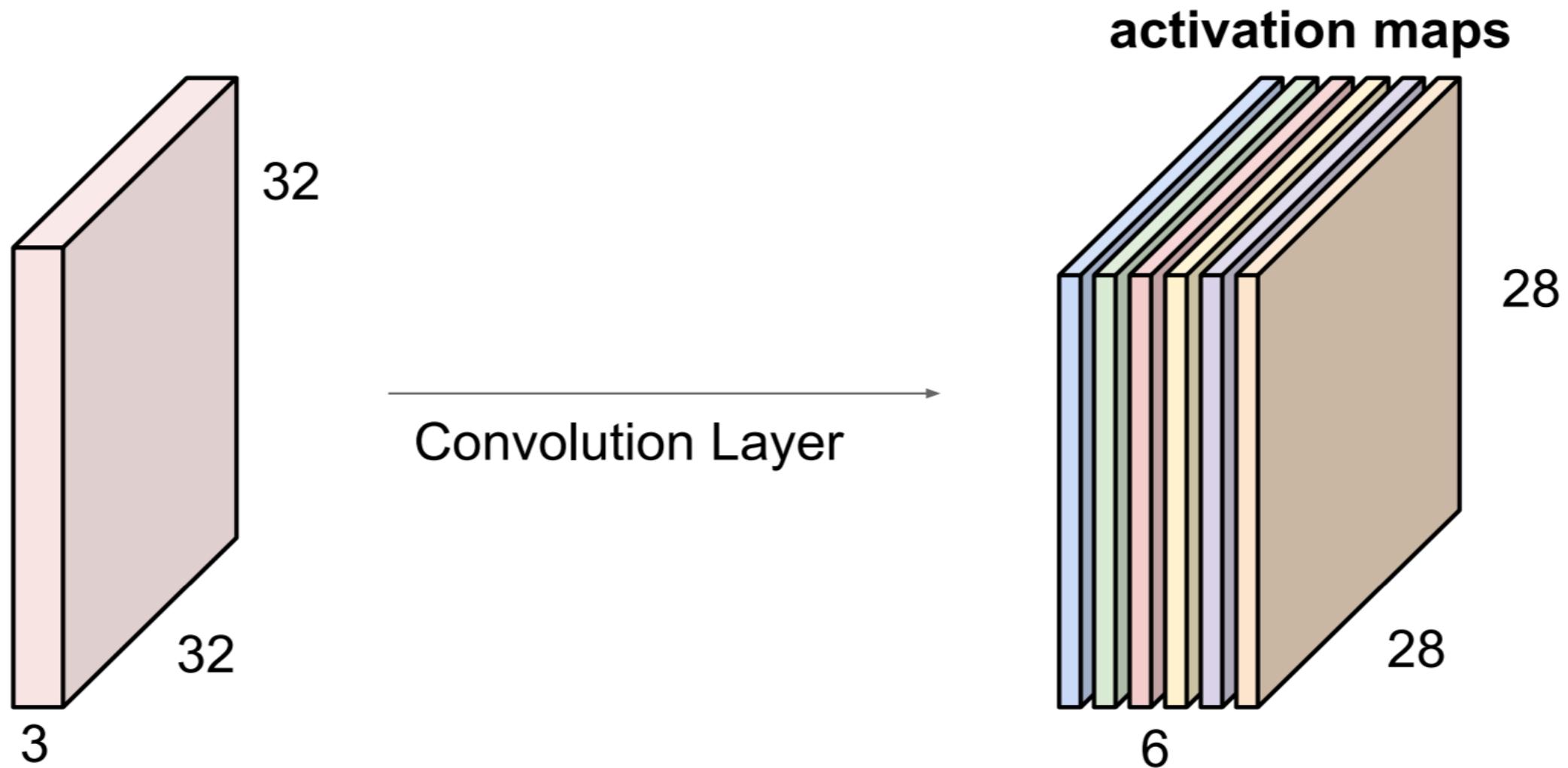


32x32x3 image
5x5x3 filter

convolve (slide) over all
spatial locations

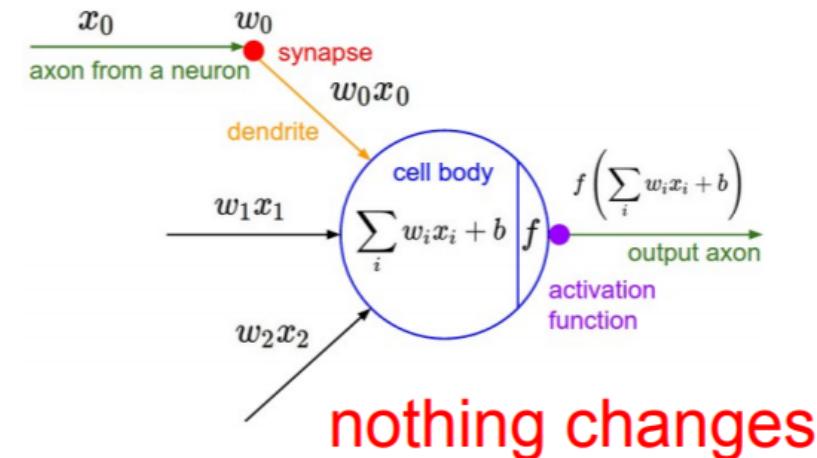
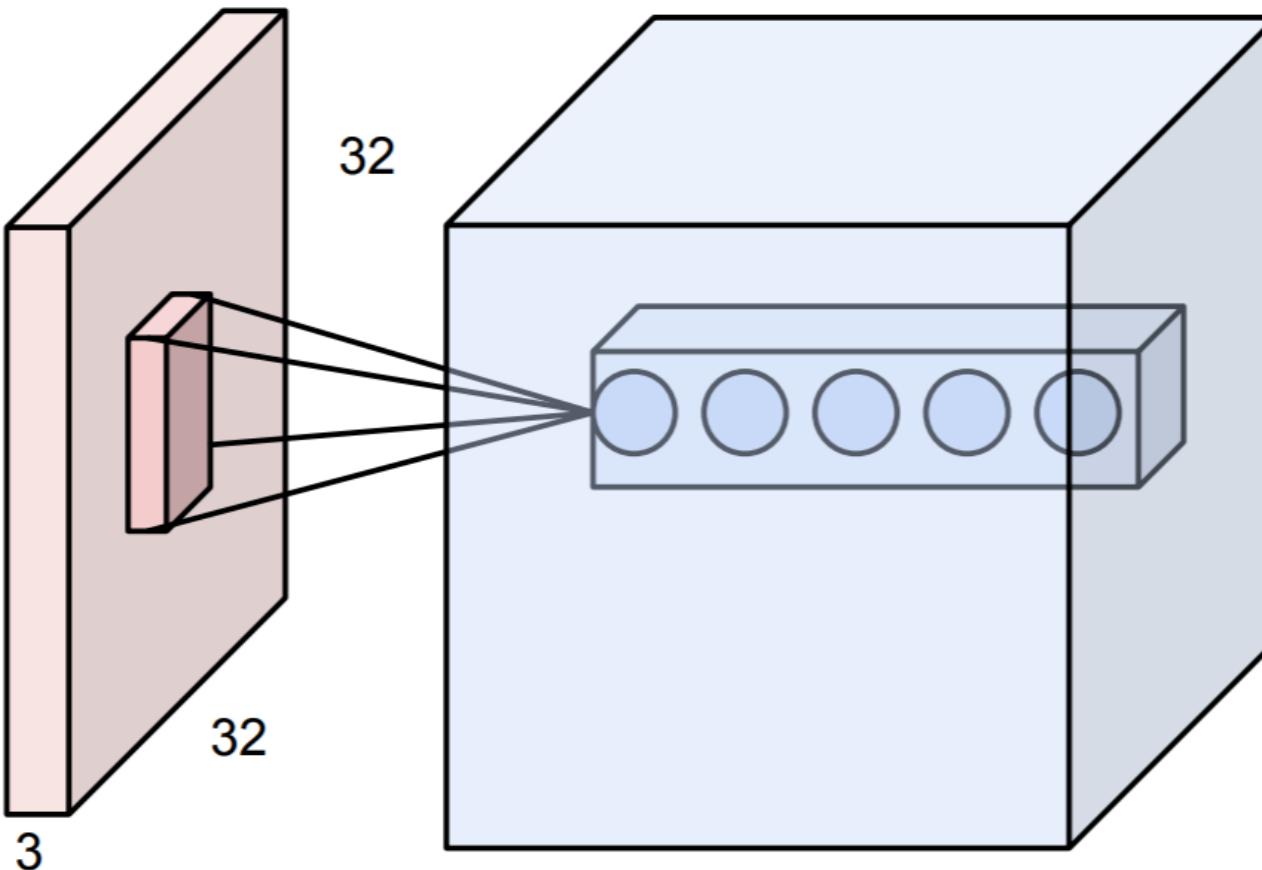


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



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

Convolutional Neural Networks
are just Neural Networks BUT:
1. Local connectivity

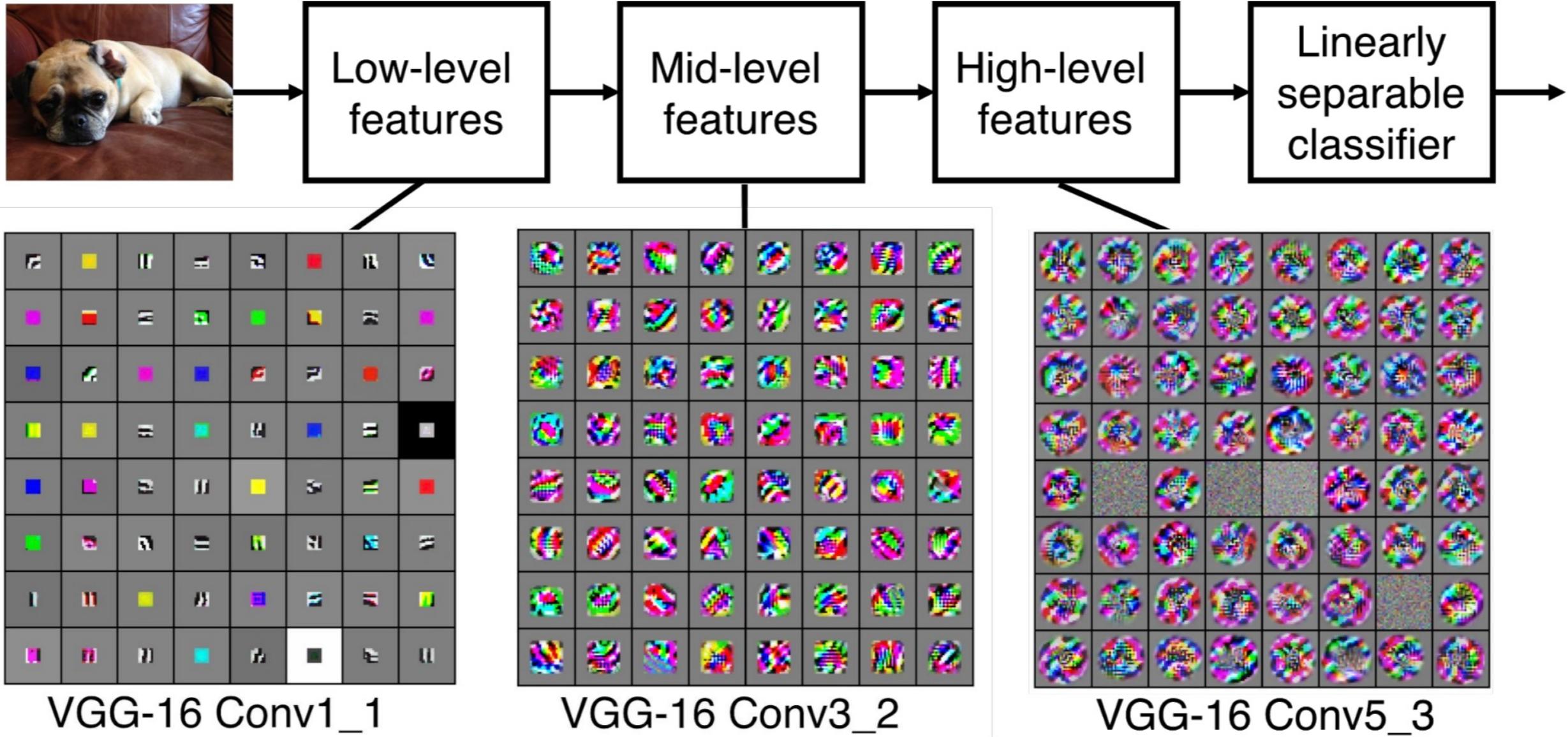


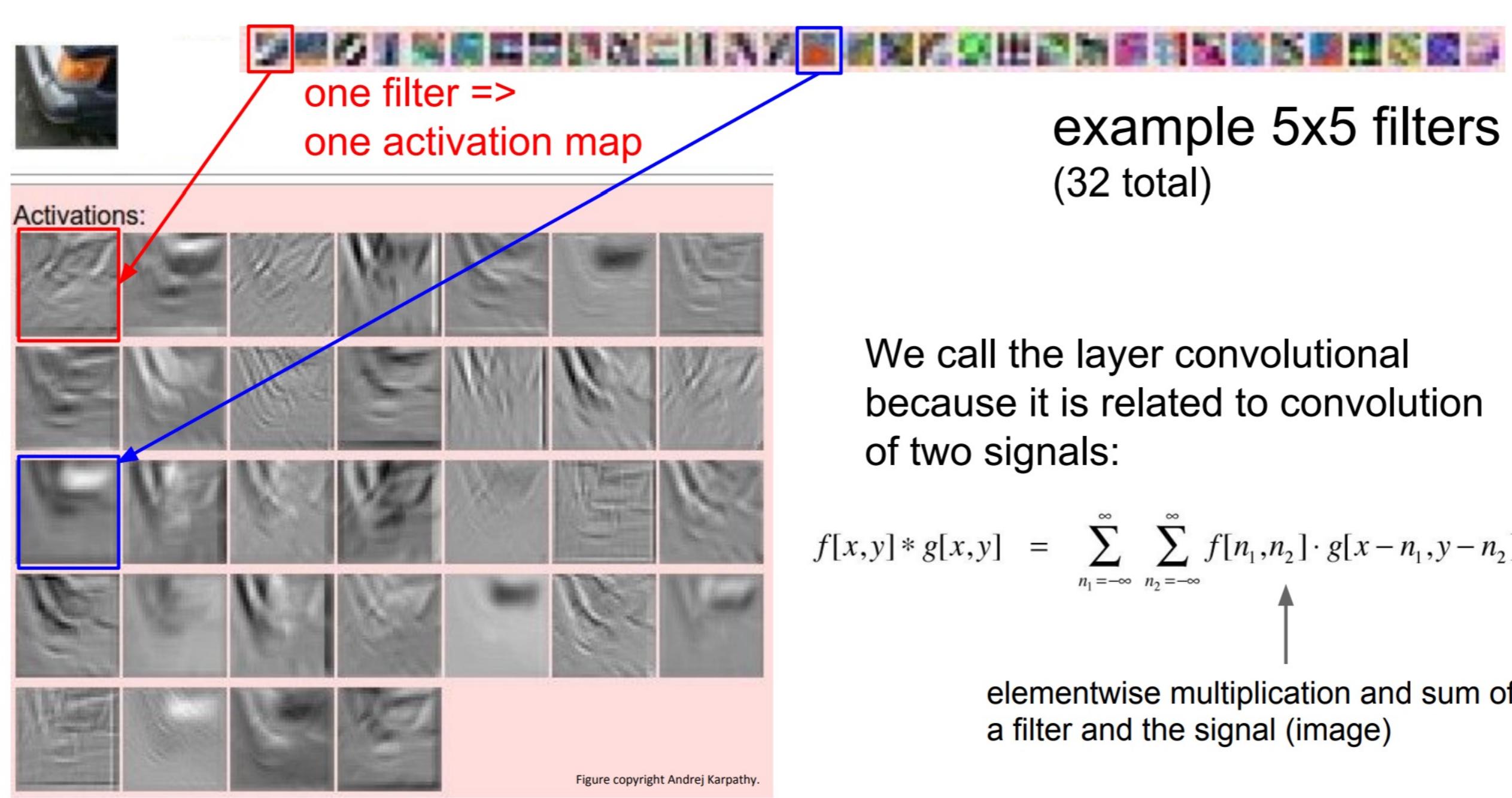
These form a single
[1 x 1 x depth]
“depth column” in the
output volume

Preview

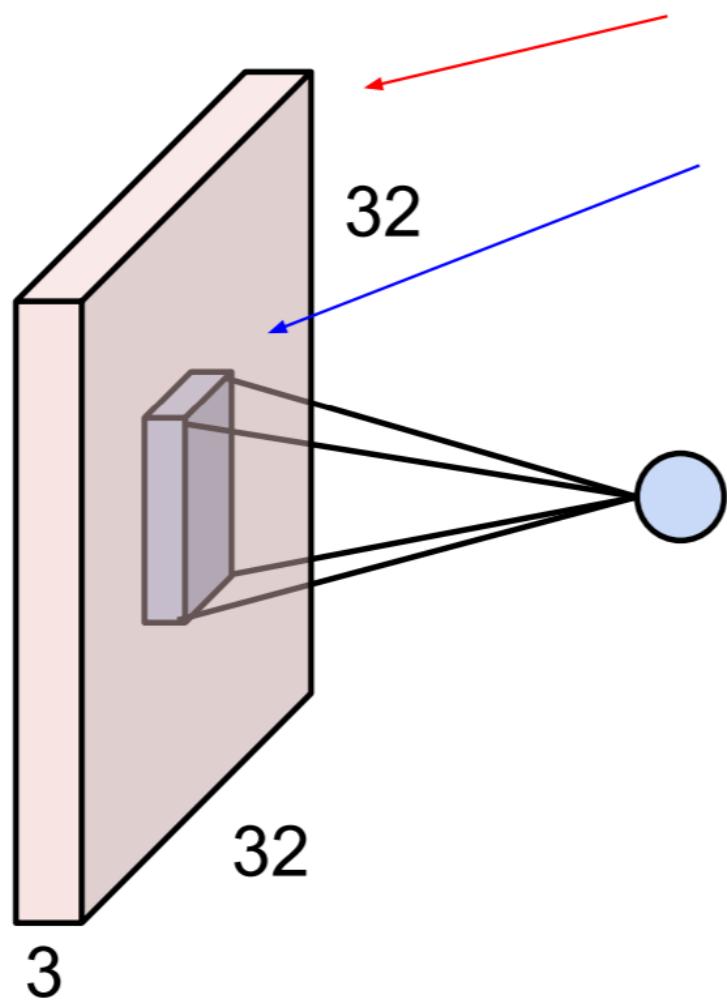
[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].





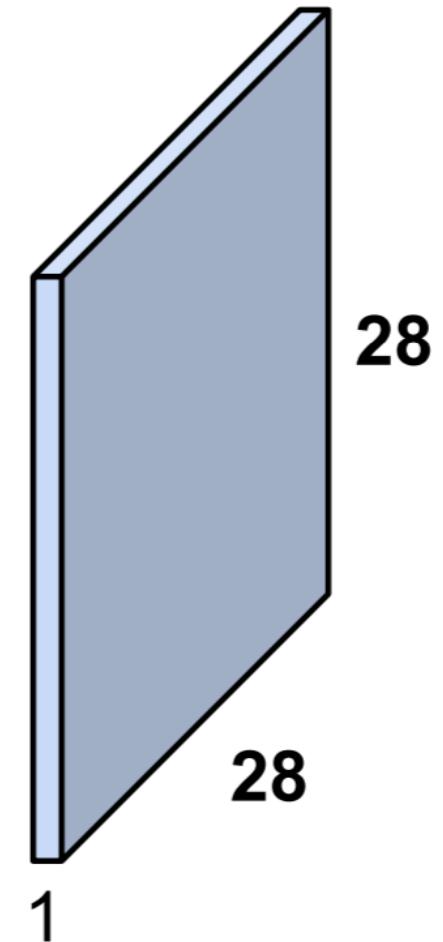
A closer look at spatial dimensions:



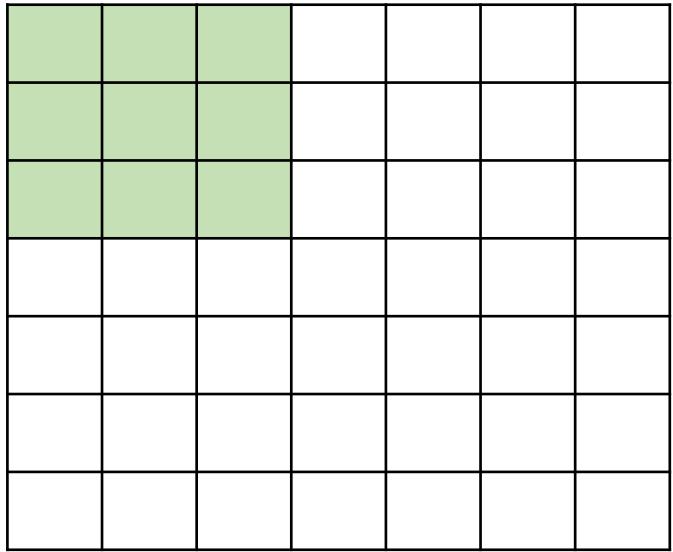
32x32x3 image
5x5x3 filter

convolve (slide) over all
spatial locations

activation map



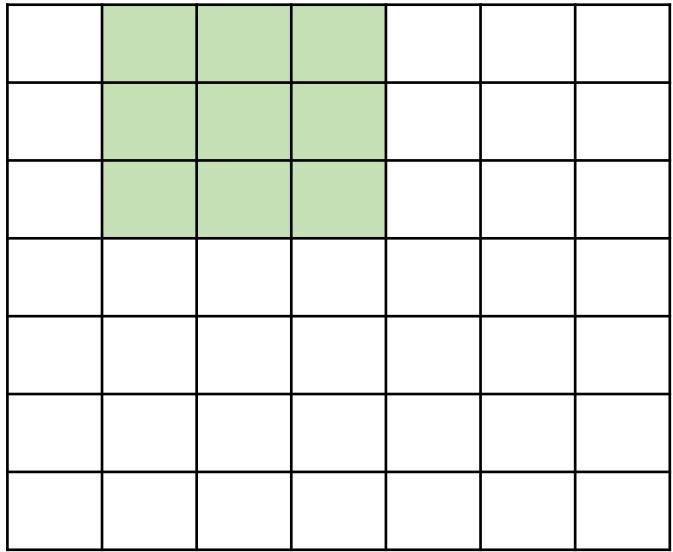
7



7

7x7 input (spatially)
assume 3x3 filter

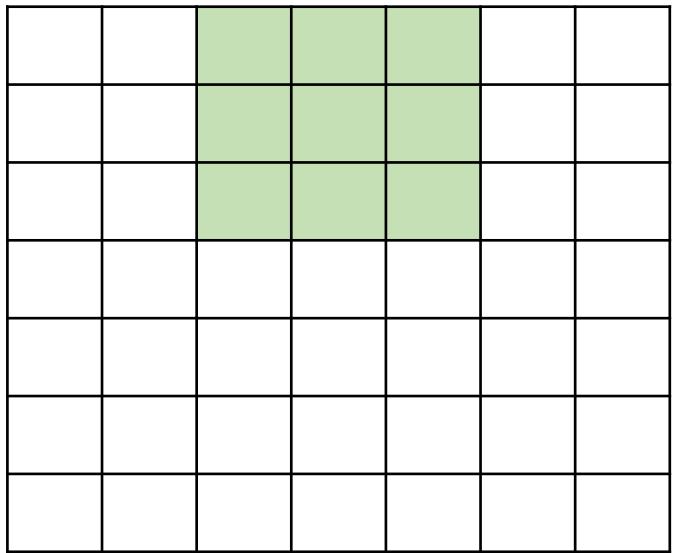
7



7

7x7 input (spatially)
assume 3x3 filter

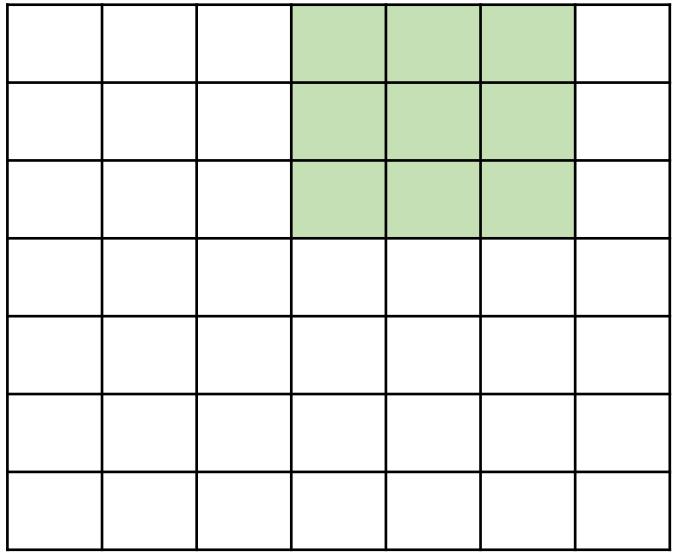
7



7

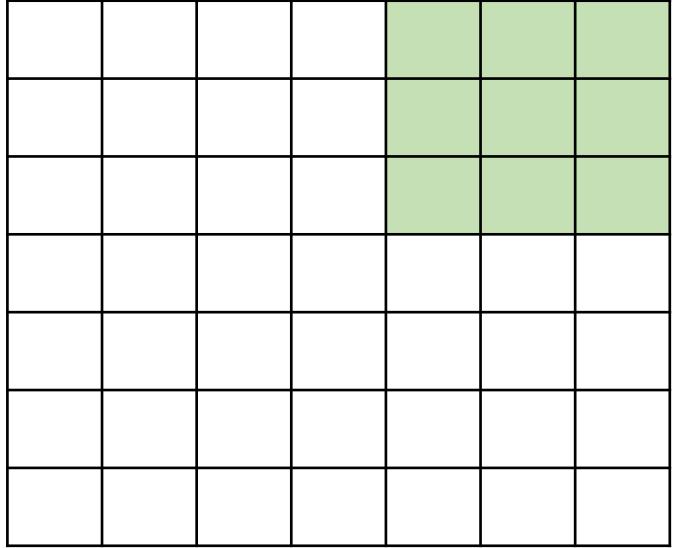
7x7 input (spatially)
assume 3x3 filter

7



7

7x7 input (spatially)
assume 3x3 filter

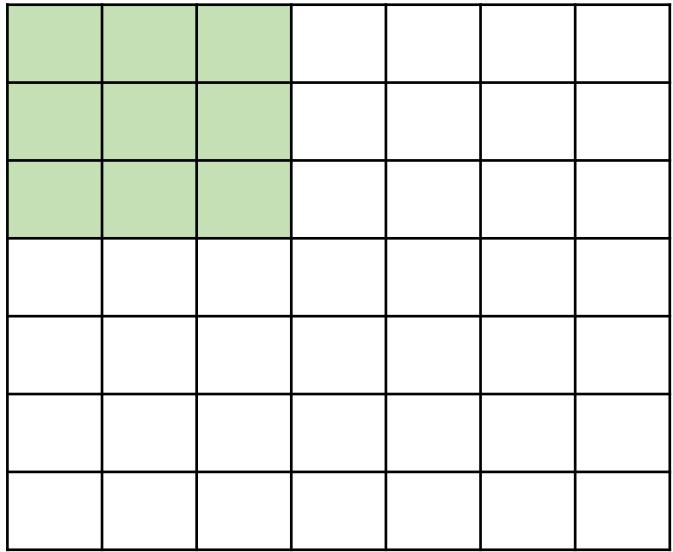


7

7x7 input (spatially)
assume 3x3 filter

→ **5x5 output**

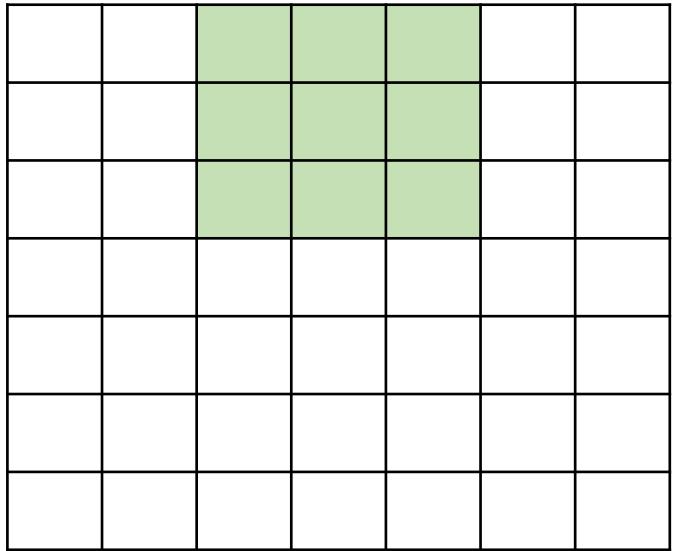
7



7

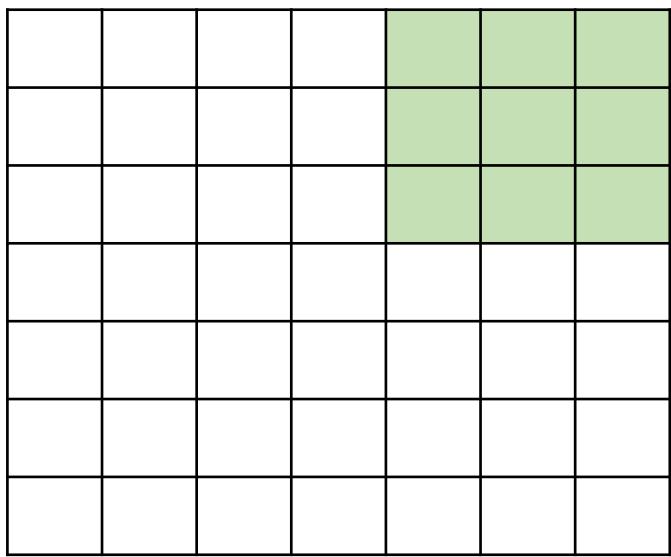
7x7 input (spatially)
assume 3x3 filter
with **stride 2**

7



7

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

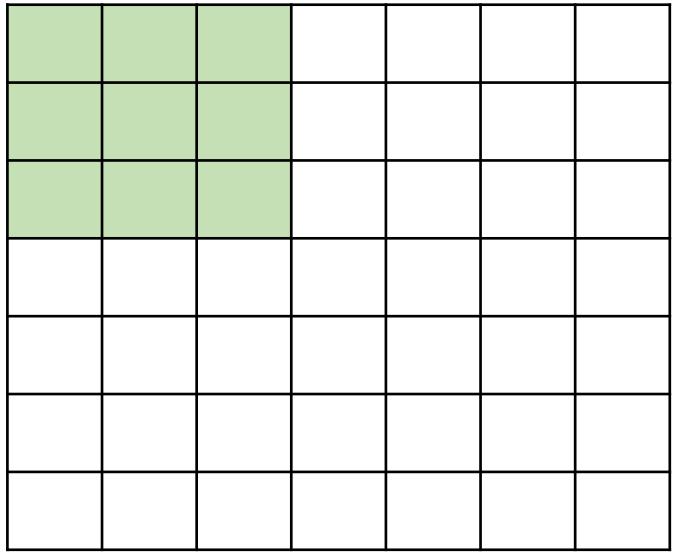


7

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

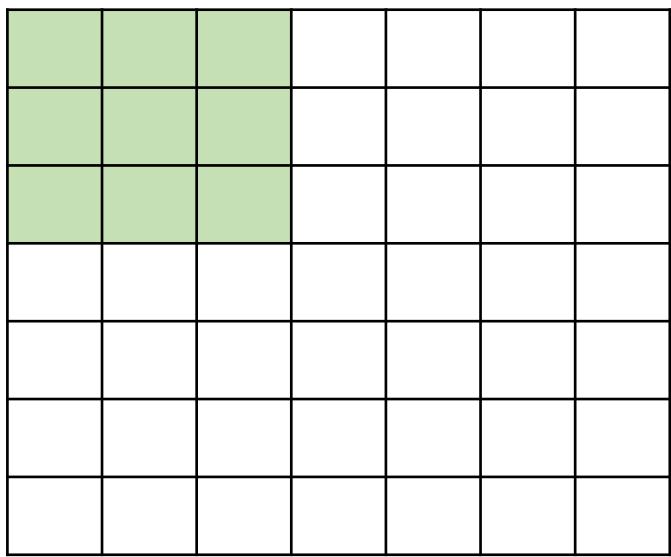
→ 3x3 output

7



7

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

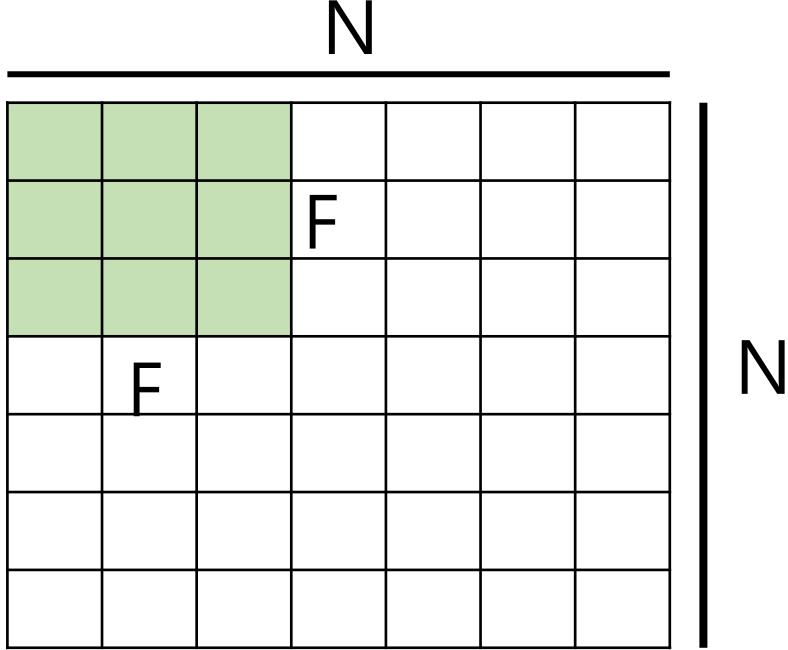


7

7

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

→ Doesn't fit!!



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

e.g. $N = 7, F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 :\backslash$

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	0	0	0	0

7x7 input (spatially)

assume 3x3 filter with **stride 1 & pad 1 pixel**

→ What is the output?

(recall:)

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

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	0	0	0	0

7x7 input (spatially)

assume 3x3 filter with **stride 1 & pad 1 pixel**

→ What is the output? **7x7 output!**

In general, zero padding with $(F - 1)/2$ pixel
preserves spatial size.

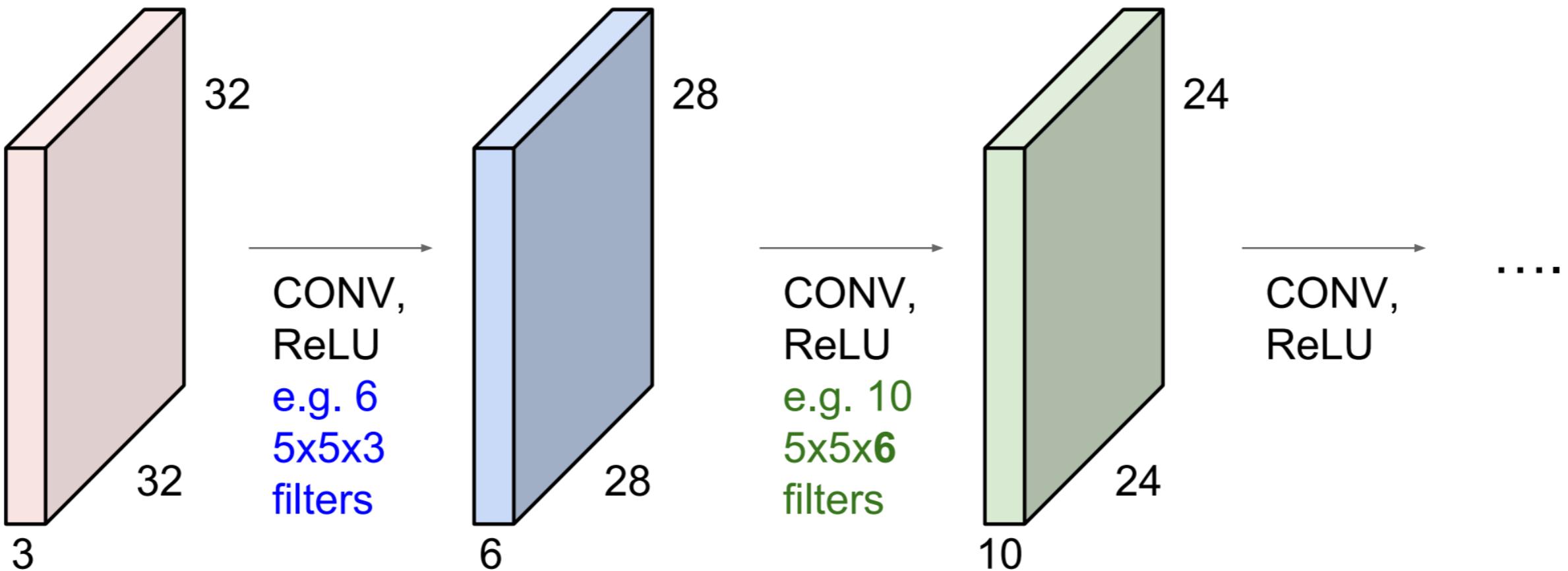
(given filters of size FxF, stride 1)

e.g., F=3 → zero pad with 1

F=5 → zero pad with 2

F=7 → zero pad with 3

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.

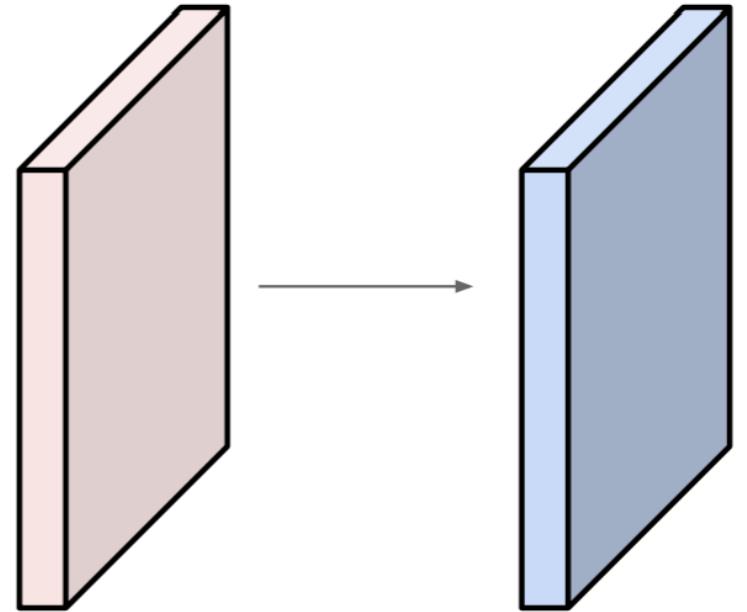


Examples time:

Input volume : **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size ??



Examples time:

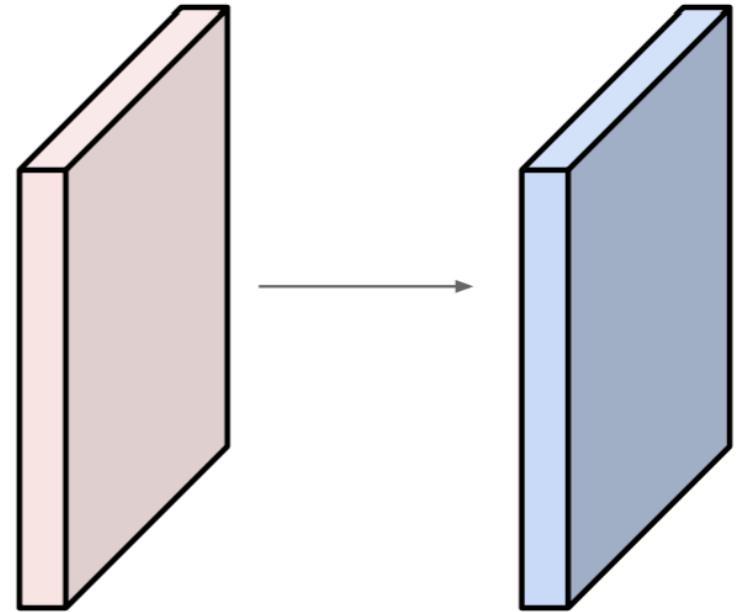
Input volume : **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size :

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

→ **32x32x10**



$$W_{output} = \frac{I - K + 2 * P}{S} + 1$$

I : input

K : filter size

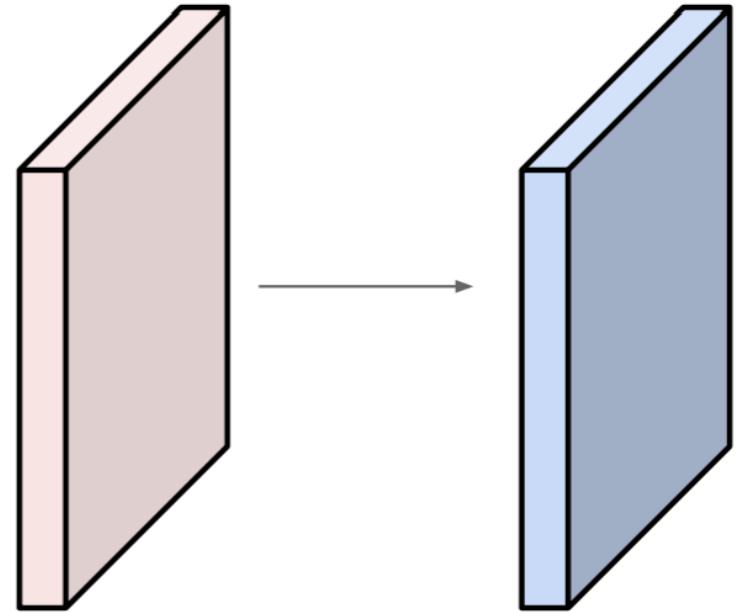
P : pad

Examples time:

Input volume : **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?



Examples time:

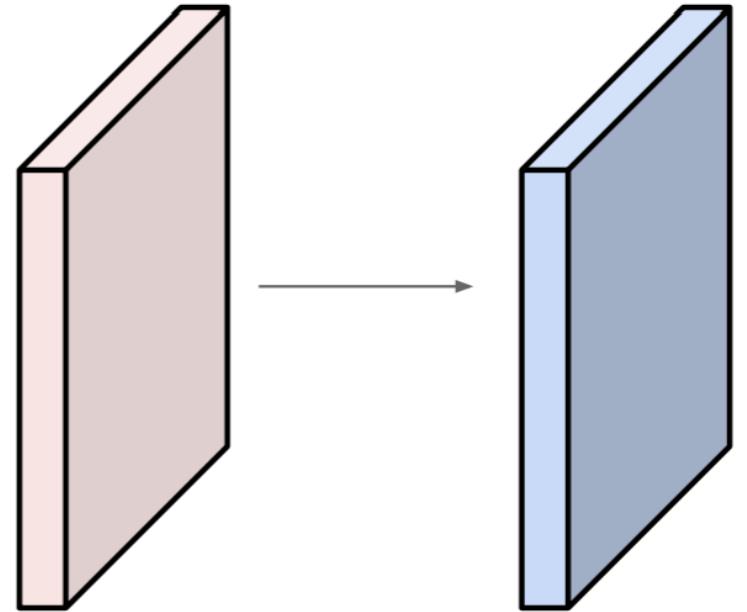
Input volume : **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?

→ each filter has $5 \times 5 \times 3 + 1 = 76$ params (+1 for bias)

→ **76x10 = 760**



Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:

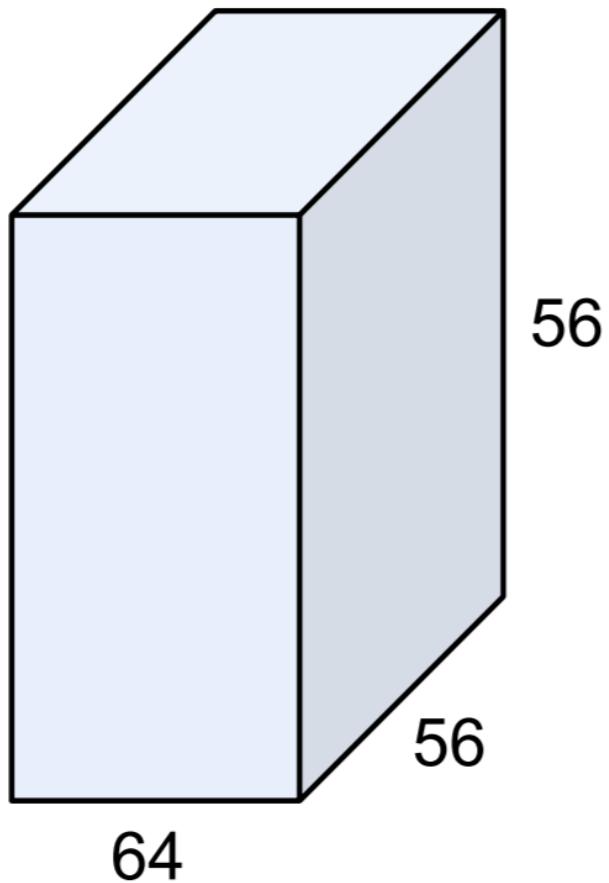
- Number of filters K ,
- their spatial extent F ,
- the stride S ,
- the amount of zero padding P .

- 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:

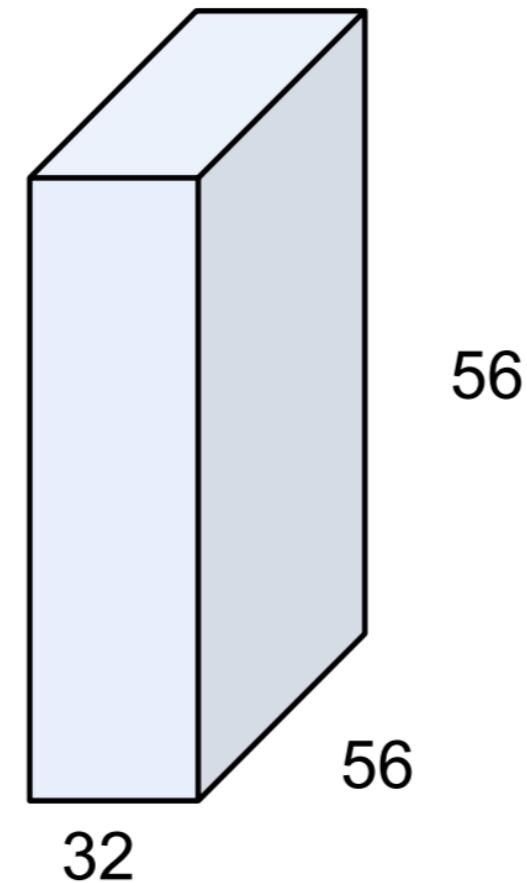
- $K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$
- $F = 3, S = 1, P = 1$
 - $F = 5, S = 1, P = 2$
 - $F = 5, S = 2, P = ? \text{ (whatever fits)}$
 - $F = 1, S = 1, P = 0$

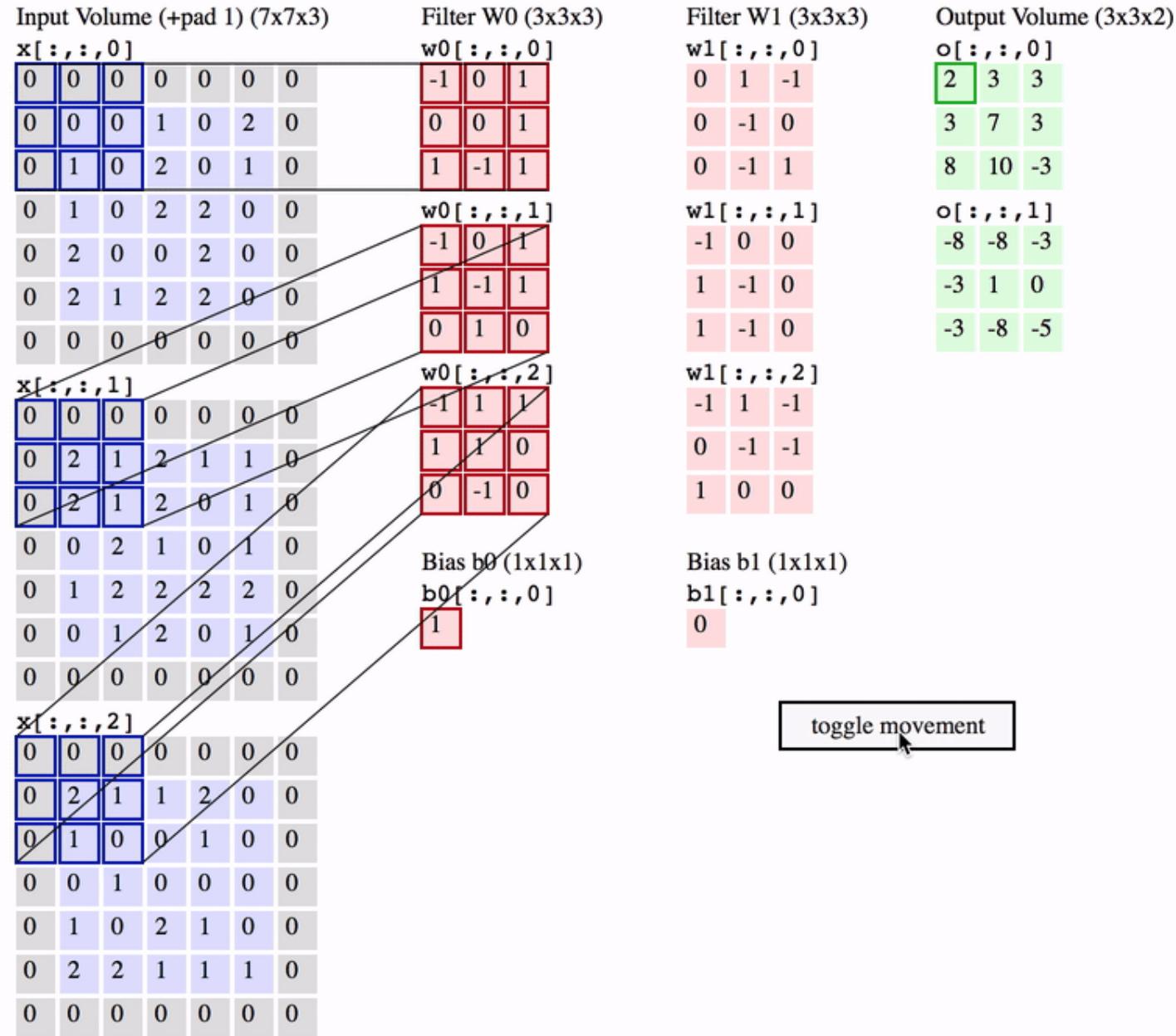
(btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

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





Examples CONV layer in PyTorch

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 $W_2 \times H_2 \times D_2$
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. v)
 - $D_2 = K$

• NOTE

Depending of the size of your kernel, several (of the last) columns of the input might be lost, because it is a valid cross-correlation, and not a full cross-correlation. It is up to the user to add proper padding.

• NOTE

When $groups == in_channels$ and $out_channels == K * in_channels$, where K is a positive integer, this operation is also termed in literature as depthwise convolution.

In other words, for an input of size $(N, C_{in}, H_{in}, W_{in})$, a depthwise convolution with a depthwise multiplier K , can be constructed by arguments ($in_channels = C_{in}$, $out_channels = C_{in} \times K$, ..., $groups = C_{in}$).

Conv2d

CLASS `torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros')`

[SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

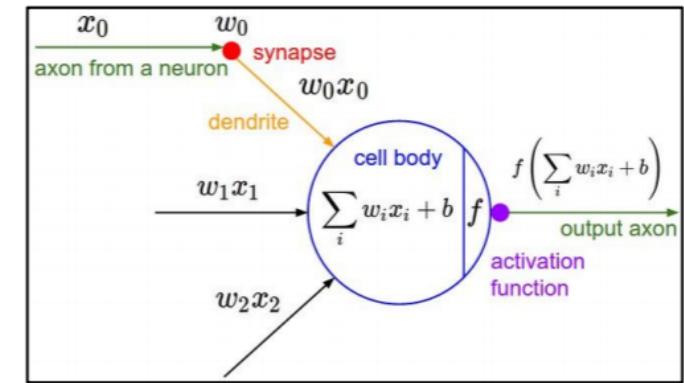
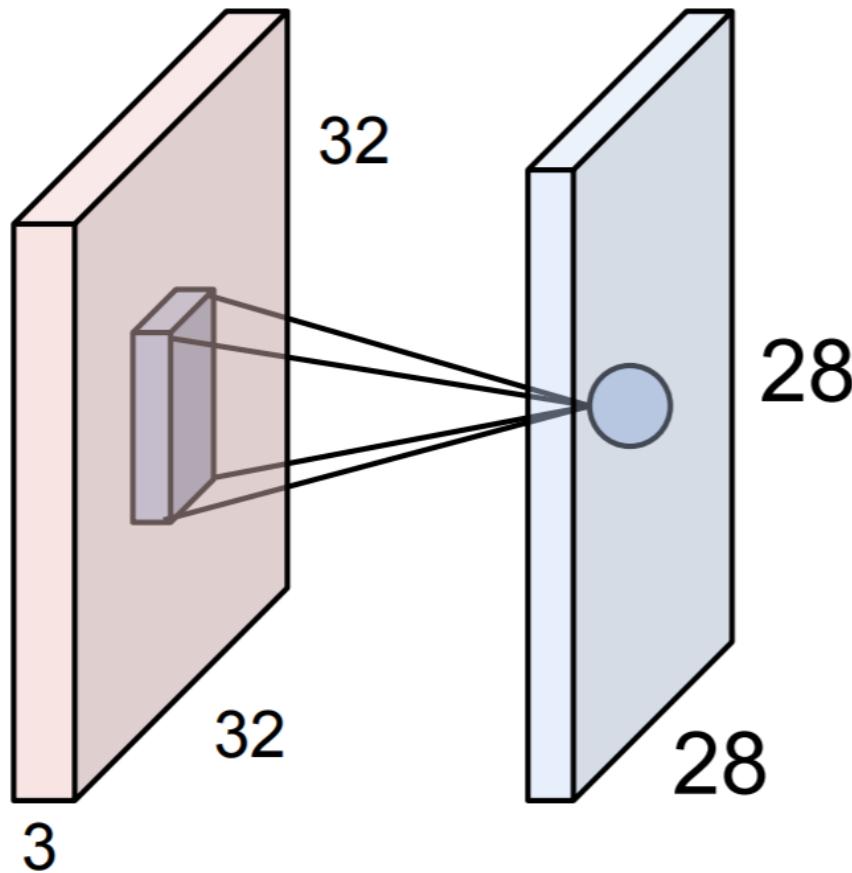
In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

$$\text{out}(N_i, C_{out_j}) = \text{bias}(C_{out_j}) + \sum_{k=0}^{C_{in}-1} \text{weight}(C_{out_j}, k) \star \text{input}(N_i, k)$$

where \star is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

- `stride` controls the stride for the cross-correlation, a single number or a tuple.
- `padding` controls the amount of implicit zero-paddings on both sides for `padding` number of points for each dimension.
- `dilation` controls the spacing between the kernel points; also known as the à trous algorithm. It is harder to describe, but this [link](#) has a nice visualization of what `dilation` does.
- `groups` controls the connections between inputs and outputs. `in_channels` and `out_channels` must both be divisible by `groups`. For example,
 - At `groups=1`, all inputs are convolved to all outputs.
 - At `groups=2`, the operation becomes equivalent to having two conv layers side by side, each seeing half the input channels, and producing half the output channels, and both subsequently concatenated.
 - At `groups= in_channels`, each input channel is convolved with its own set of filters, of size: $\left\lfloor \frac{out_channels}{in_channels} \right\rfloor$.

The brain/neuron view of CONV Layer

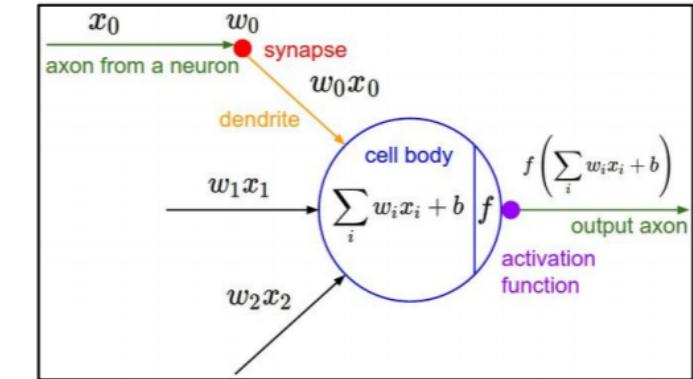
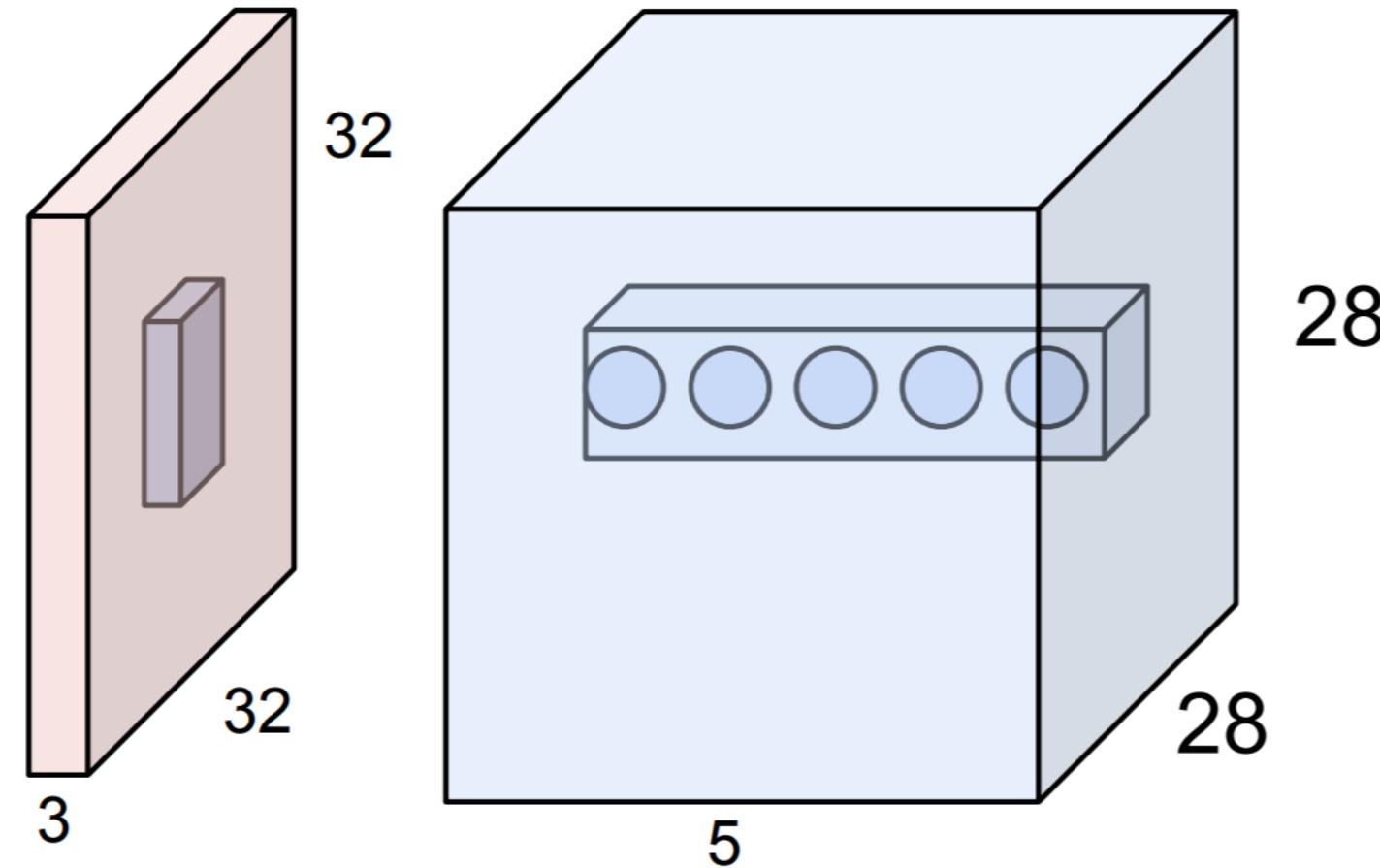


An activation map is a 28x28 sheet of neuron outputs:

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



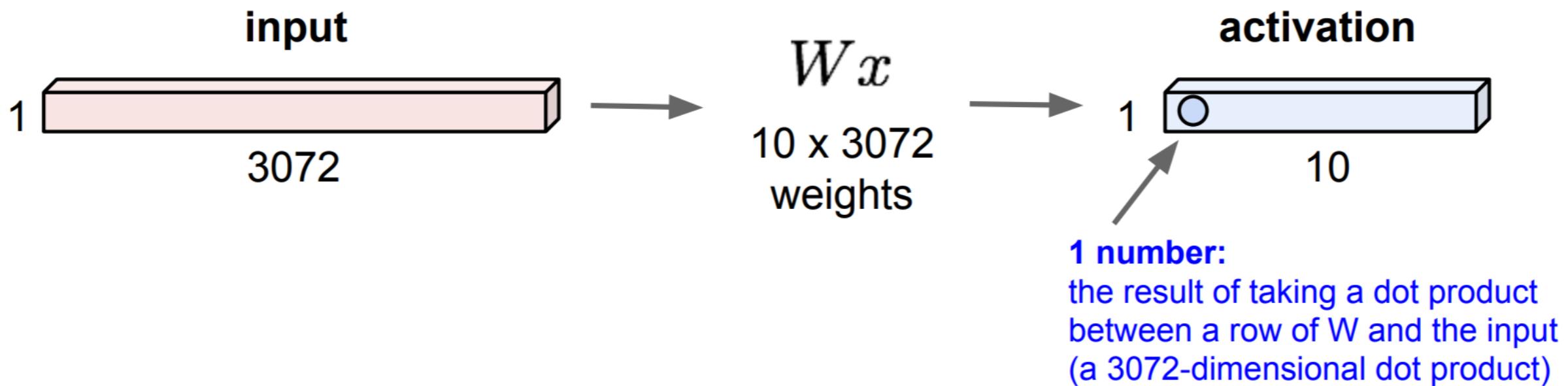
E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
(28x28x5)

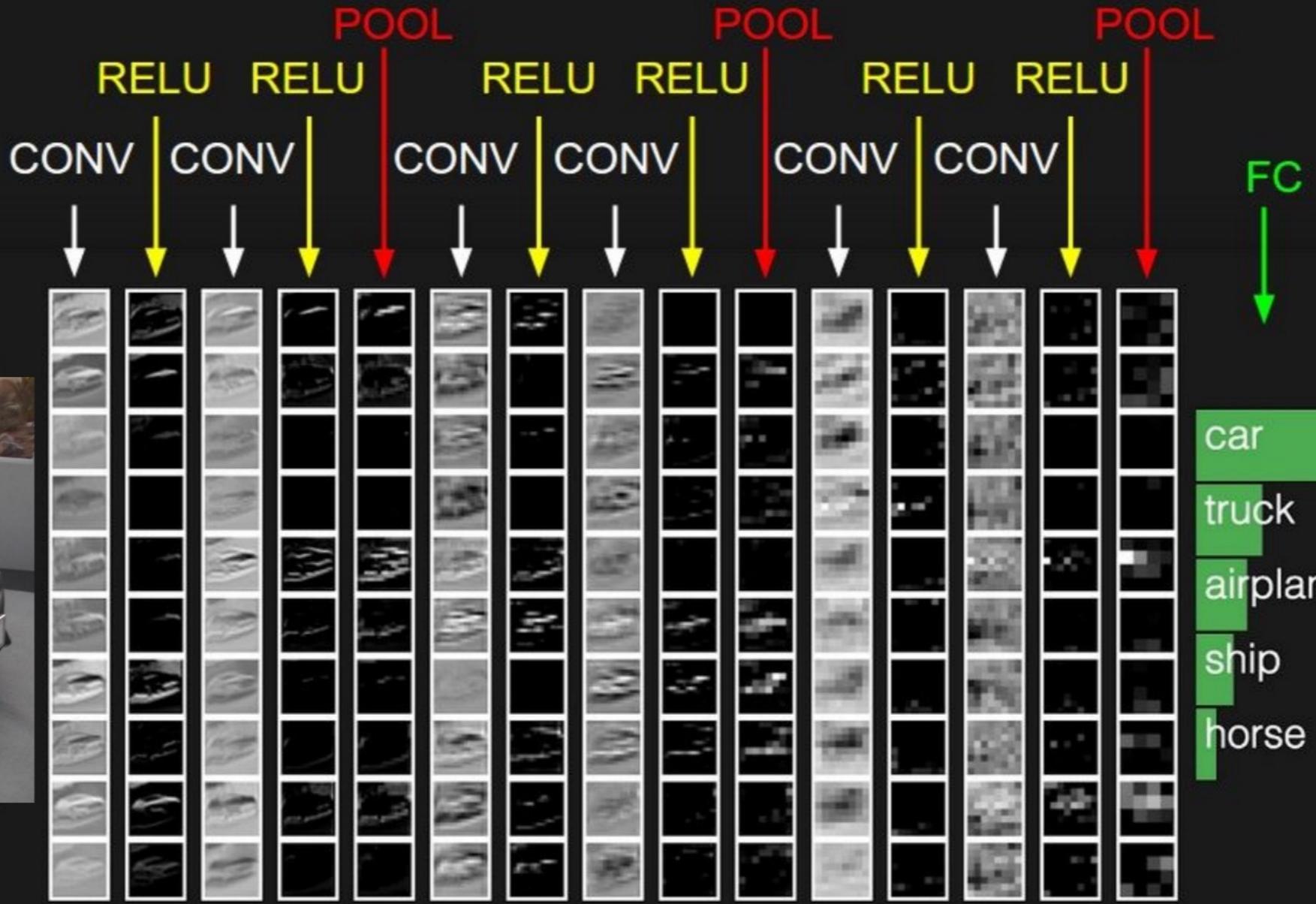
There will be 5 different
neurons all looking at the same
region in the input volume

Reminder: Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

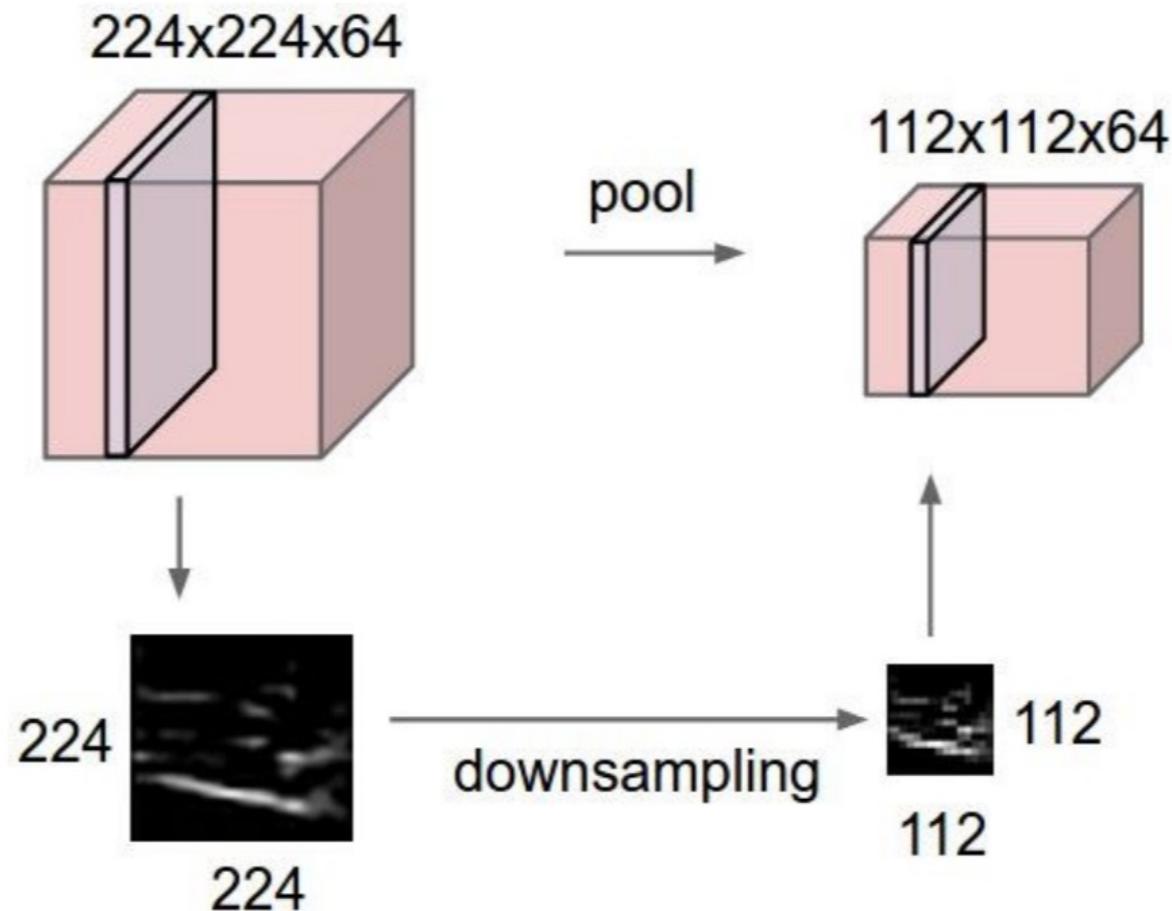
Each neuron
looks at the full
input volume





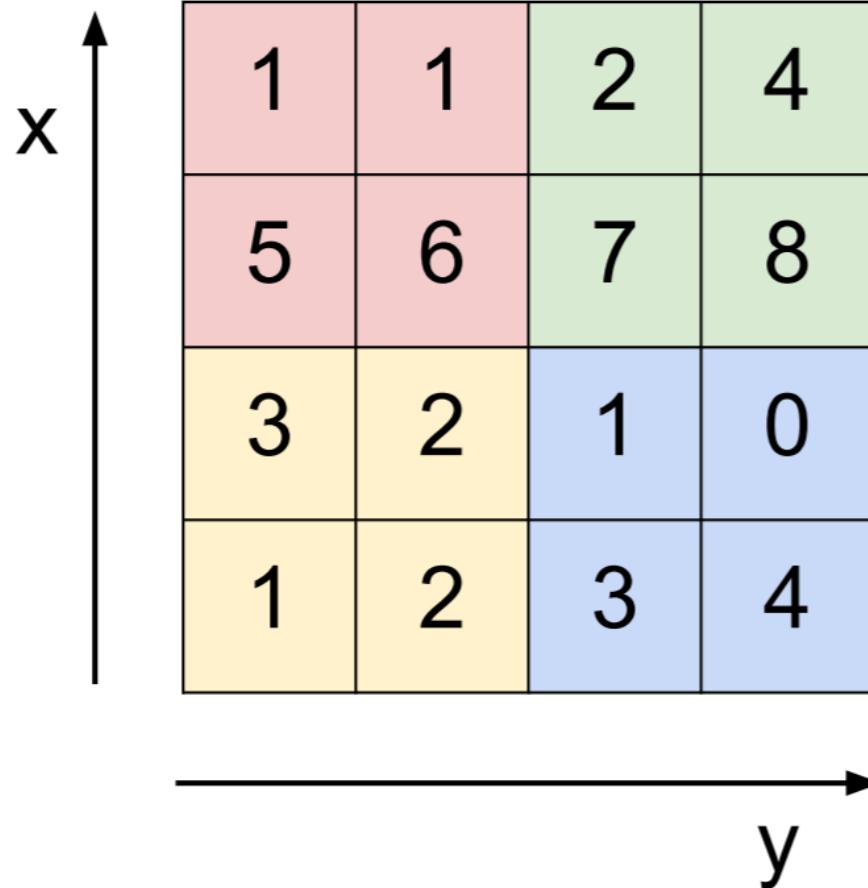
Pooling layer

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



MAX POOLING

Single depth slice



max pool with 2x2 filters
and stride 2

6	8
3	4

Common settings:

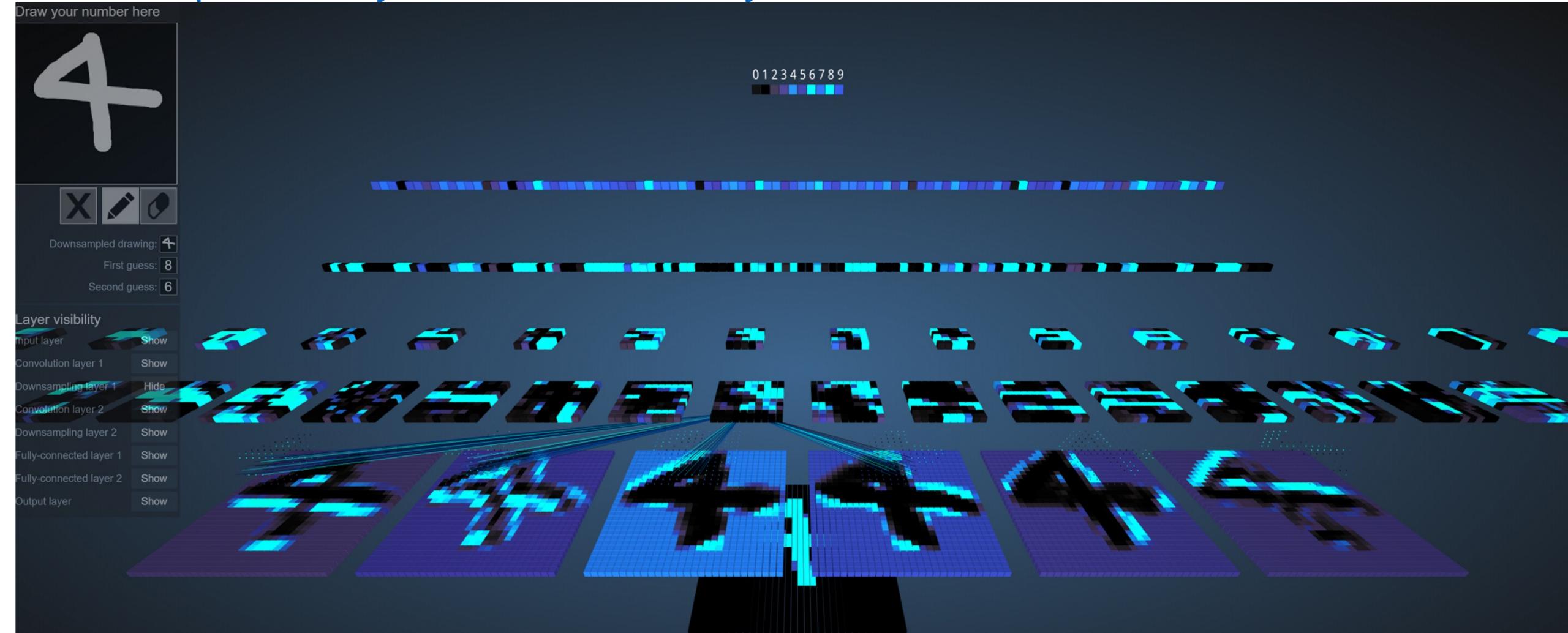
$F = 2, S = 2$

$F = 3, S = 2$

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Visualization

<https://scs.ryerson.ca/~aharley/vis/conv/>



Visualization

<https://poloclub.github.io/cnn-explainer/>



Visualization

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

ConvNetJS CIFAR-10 demo

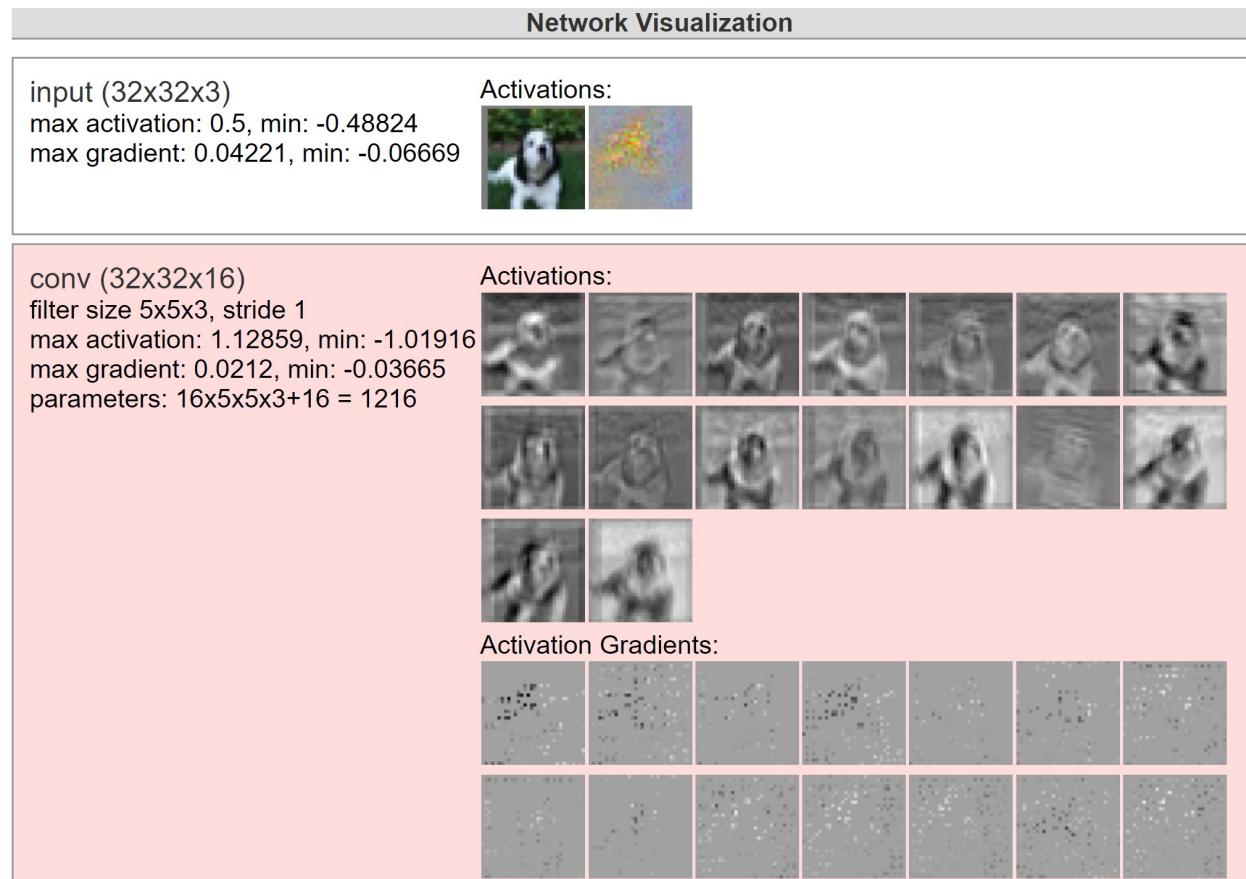
Description

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to parse the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

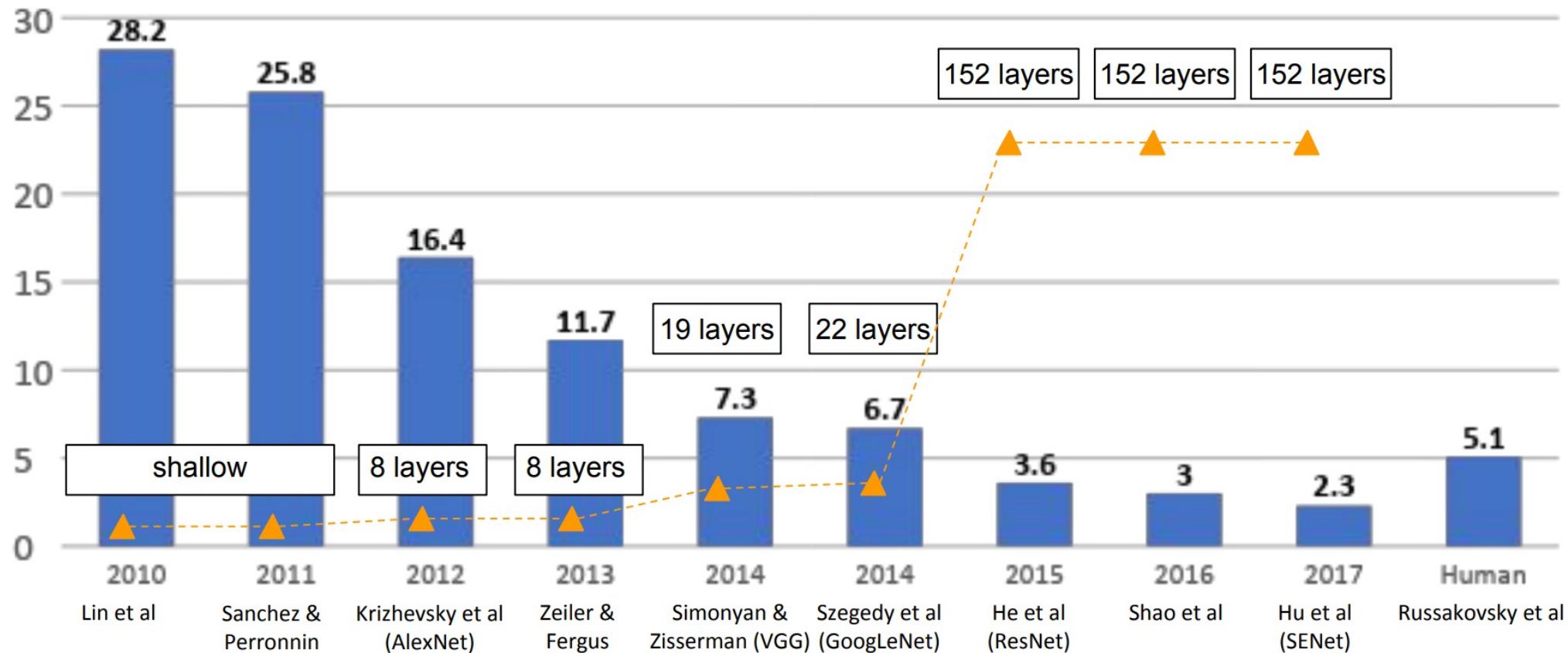
By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).



Takeaway

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



References

- <http://cs231n.github.io/convolutional-networks/>
- <http://cs231n.stanford.edu/>
- https://deeplearningzerotoall.github.io/season2/lec_pytorch.html