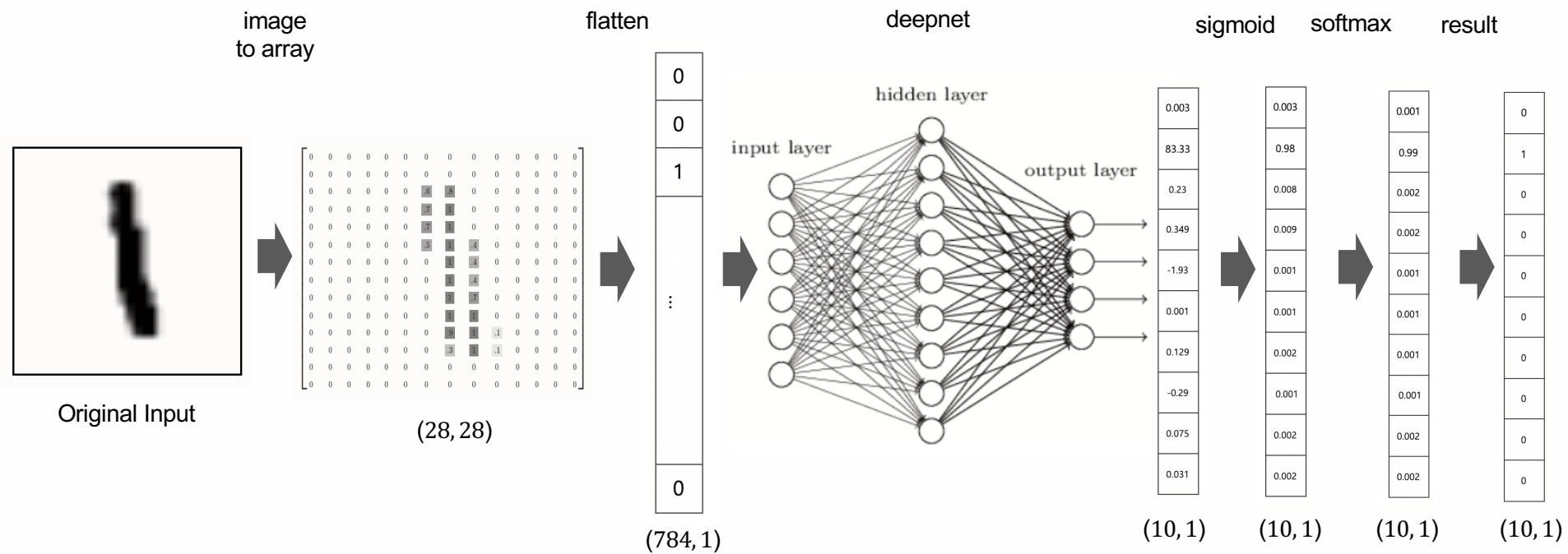


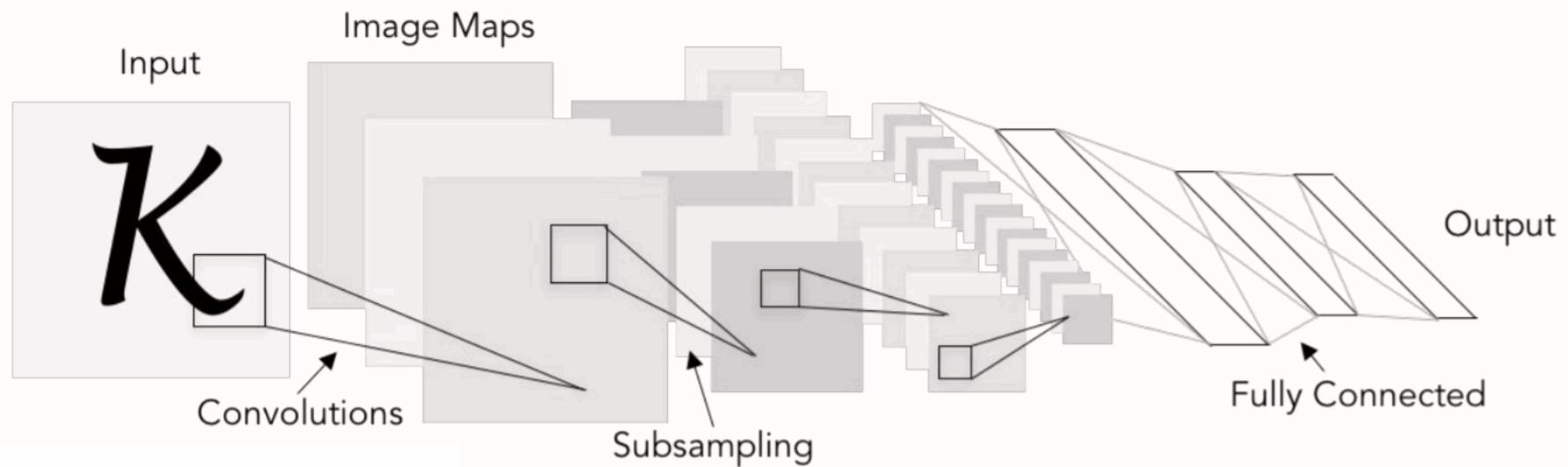
lecture 6

Dong-Geol Choi
Hanbat Nat'l Univ.

Remind



Convolutional Neural Networks

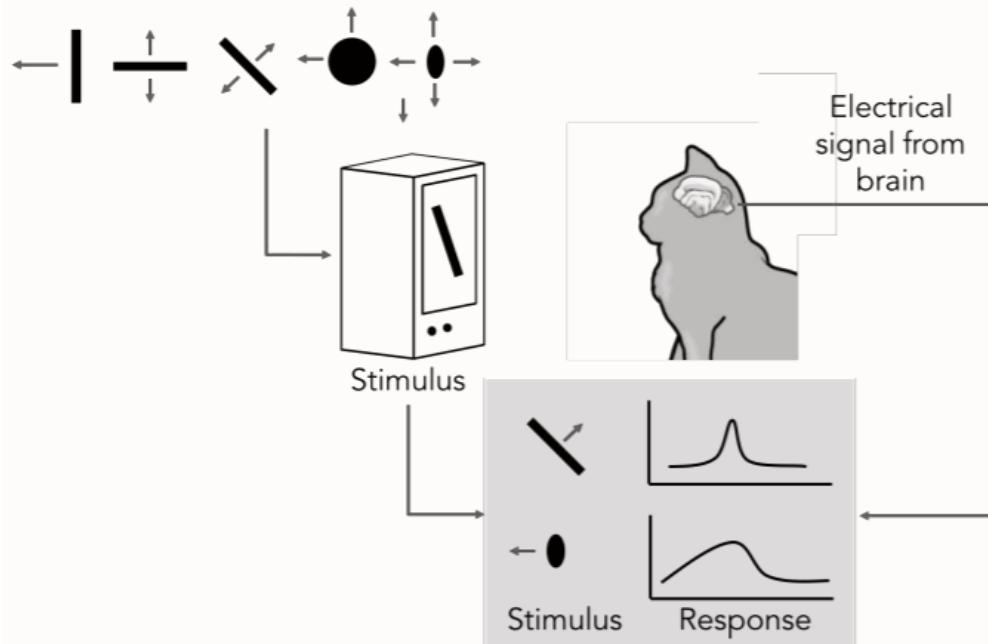


LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.

Hubel and Wiesel Cat Experiment



Hubel and Wiesel Cat Experiment : Results

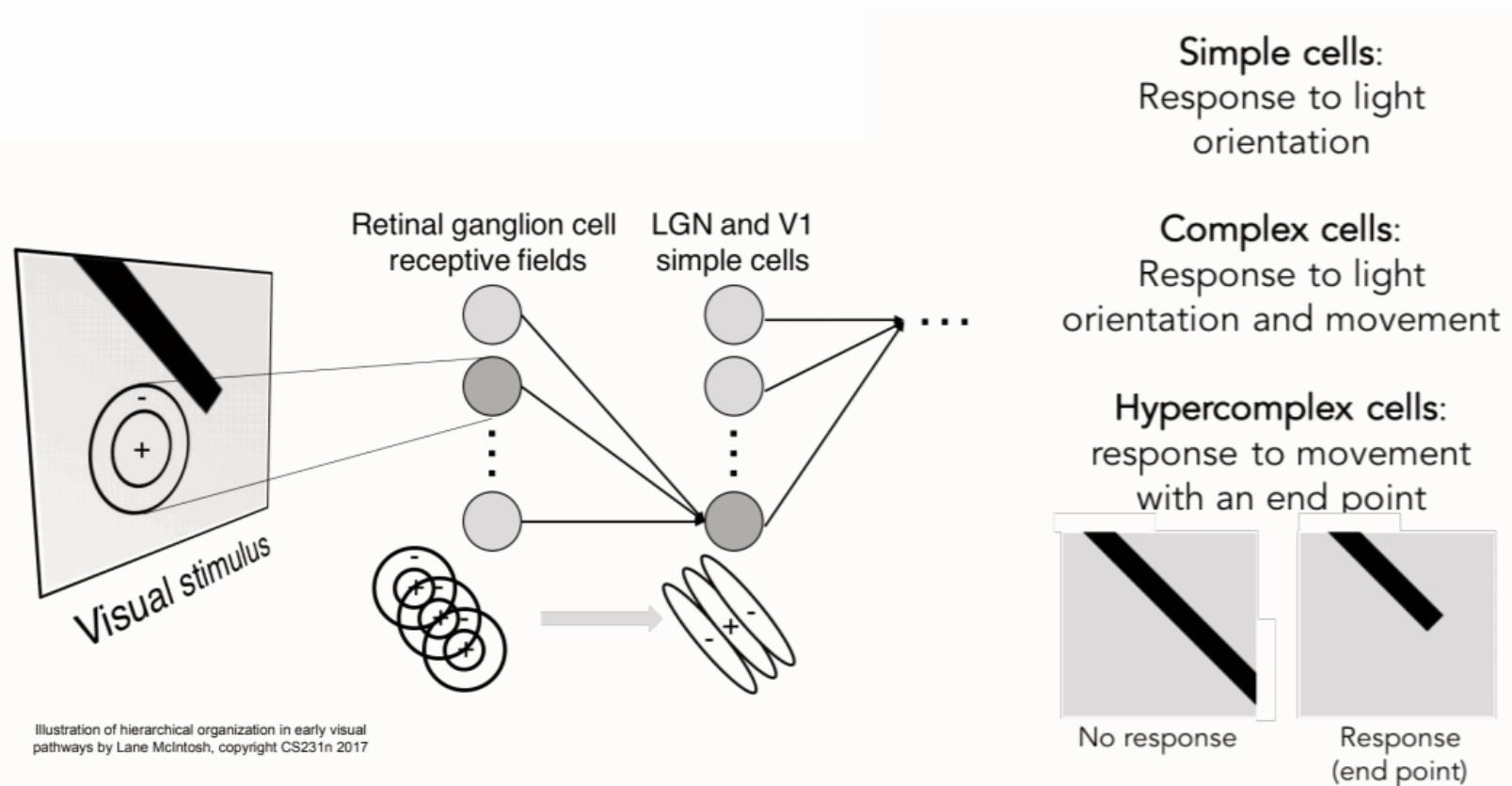


This experiment found that neurons responded to things like oriented **edges** and **shapes**.

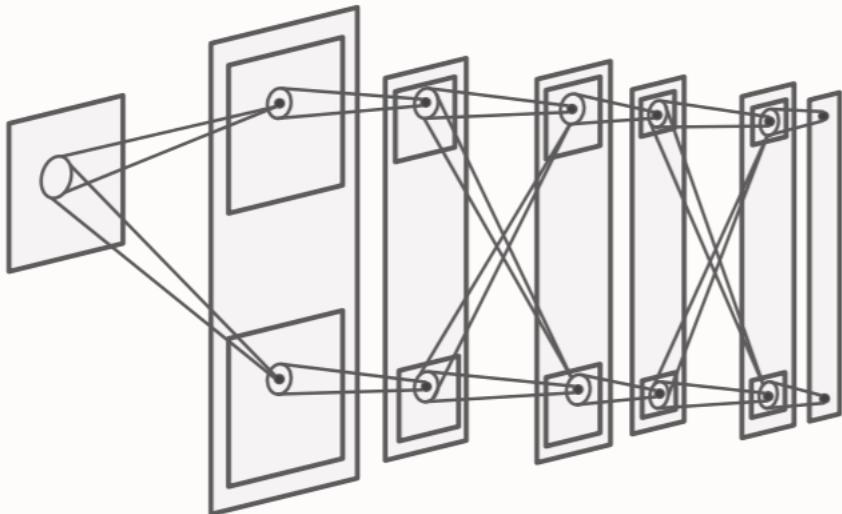
Experiments have found that neurons have a **hierarchy**.

[Cat image](#) by CNX OpenStax is licensed under CC BY 4.0; changes made

Hubel and Wiesel Cat Experiment : Results



Neocognitron [Fukushima 1980]



The first NN to use the idea of simple / complex cells.

“sandwich” architecture (SCSCSC···)
simple cells: modifiable parameters
complex cells: perform pooling

Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]

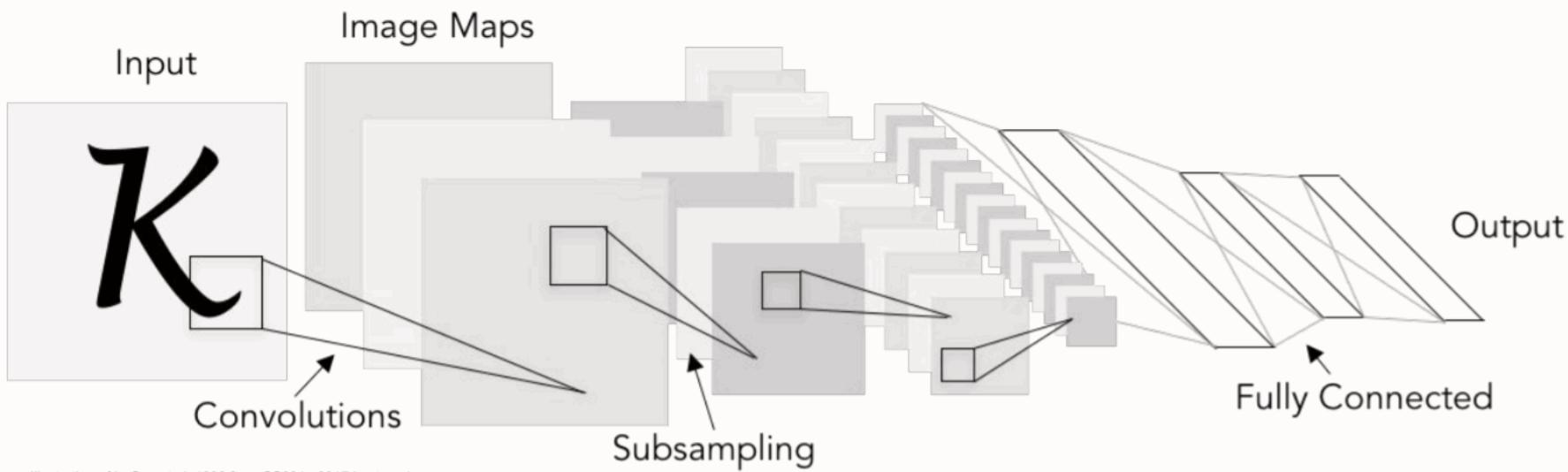


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

In 1998, Yann LeCun first trained NN Apply backprob and gradient-based learning.

It works well for recognizing postal code numbers and is widely used for real postal code recognition in postal services.

ImageNet Classification with Deep Convolutional Neural Networks [Krizhevsky, Sutskever, Hinton, 2012]

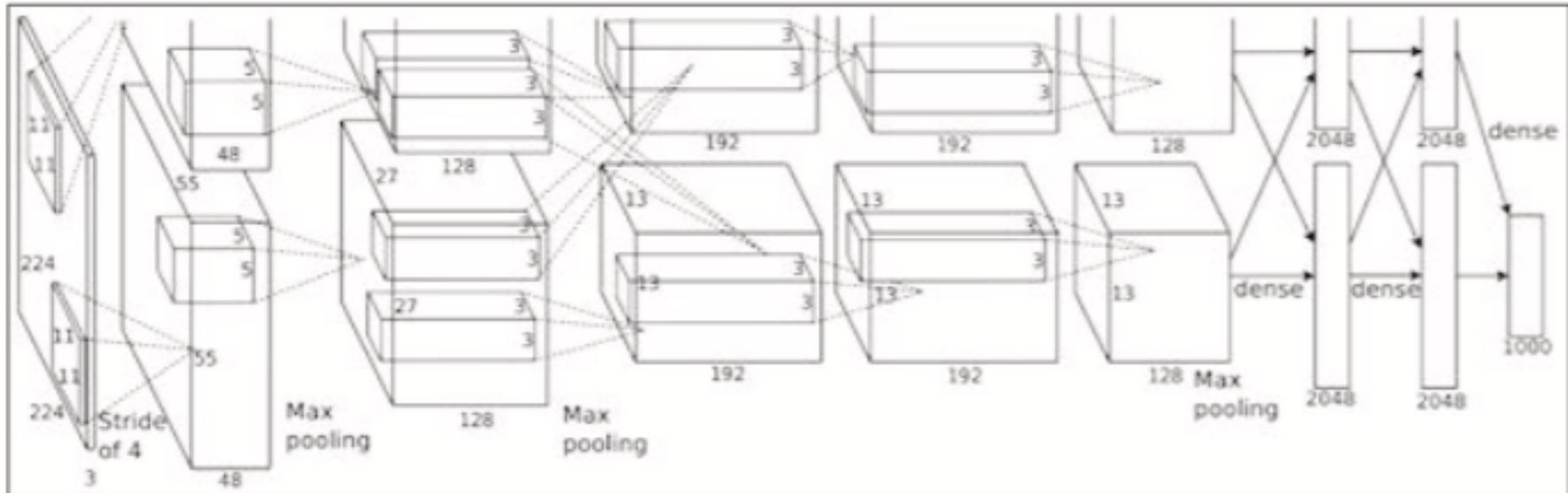


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

The most important point is that now we have a **big data** is there and we can take advantage of the **GPU**.

ConvNets are everywhere : Classification, Retrieval

Classification



Retrieval



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

ConvNets are everywhere : Detection, Segmentation

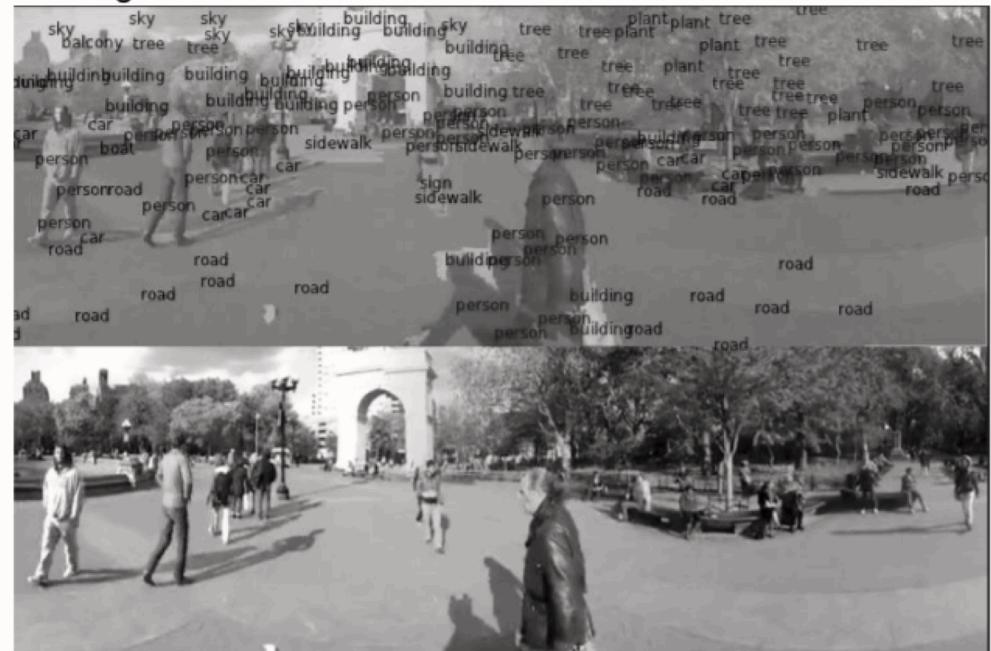
Detection



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

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

Segmentation



Figures copyright Clement Farabet, 2012.
Reproduced with permission.

[*Farabet et al., 2012*]

ConvNets are everywhere : Video

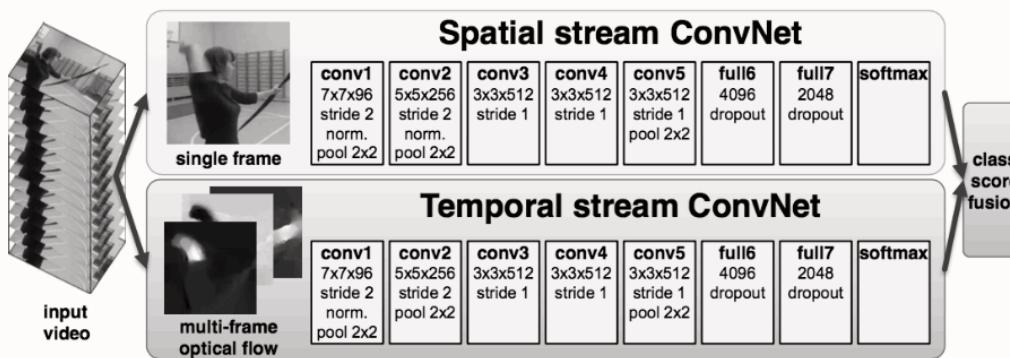


Figure 1: Two-stream architecture for video classification.

ConvNets are everywhere : self-driving cars



self-driving cars

Photo by Lane McIntosh. Copyright CS231n 2017.



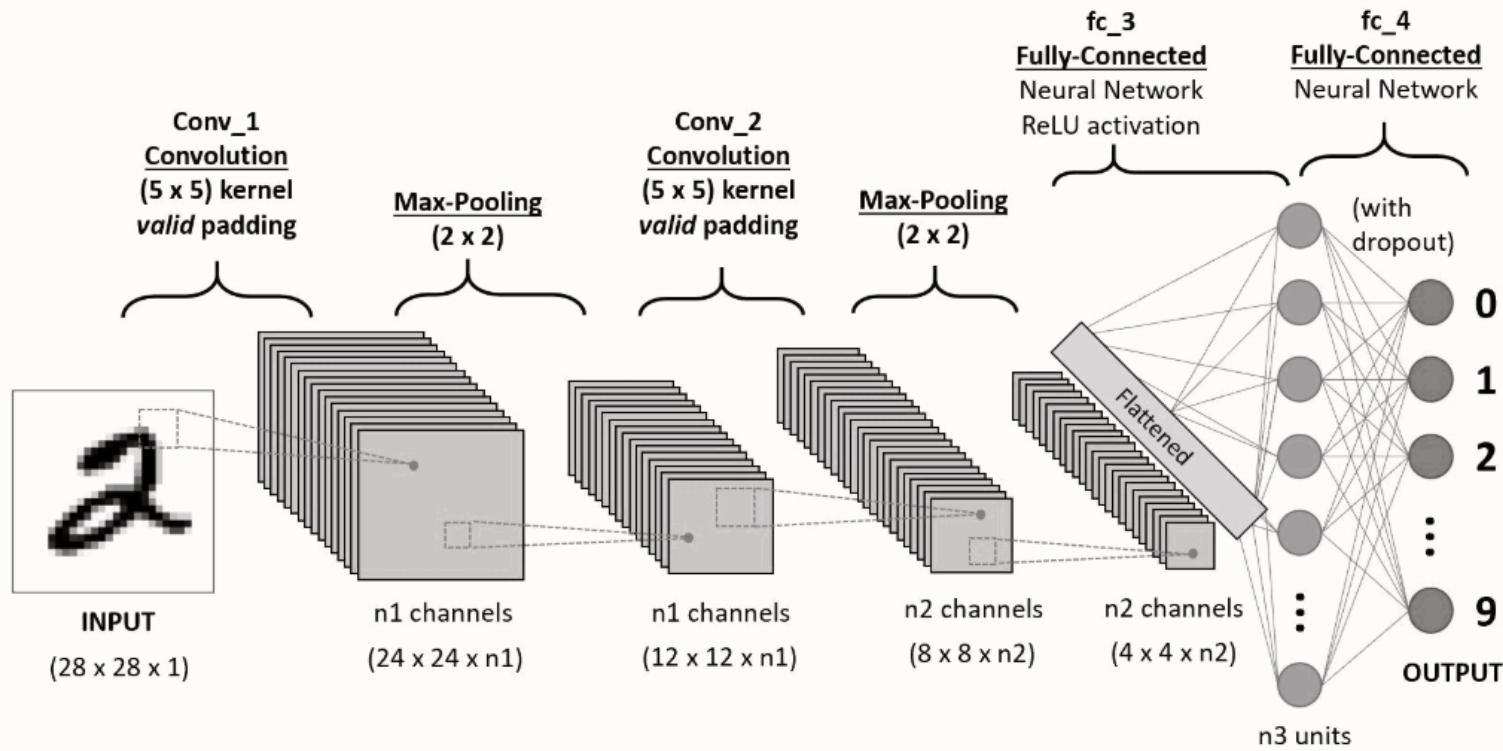
[This image](#) by GBPublic_PR is
licensed under [CC-BY 2.0](#)

NVIDIA Tesla line

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

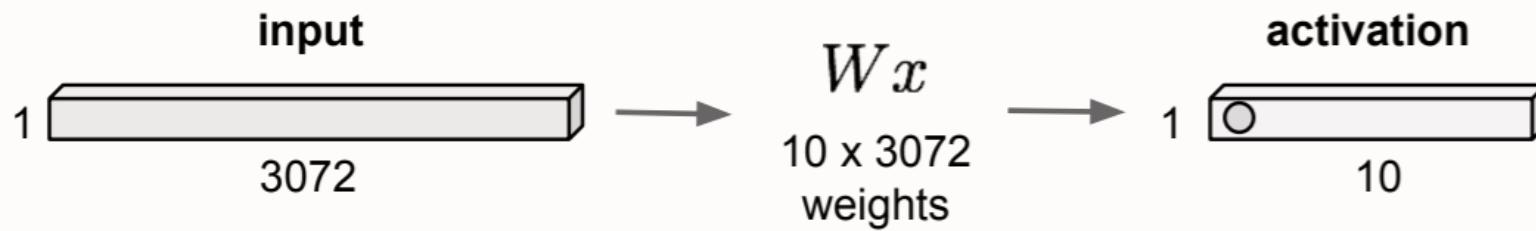
Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

Convolutional Neural Networks



Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



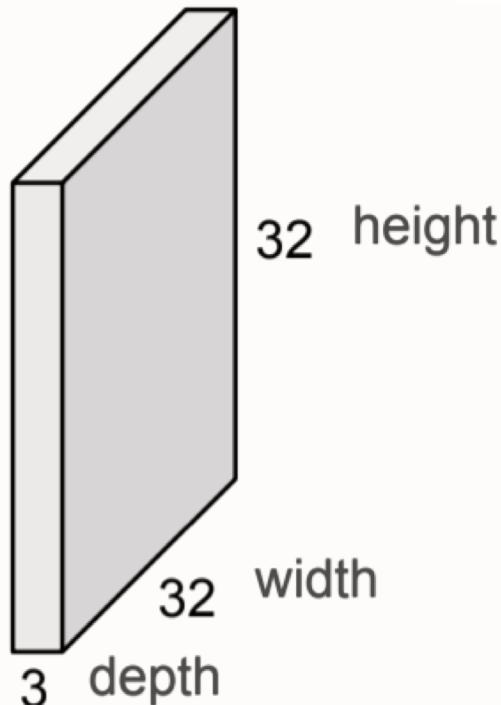
Fully Connected Layer **loses spatial information**.

The Fully Connected layer does **not work if the image size is different**.

Since the Fully Connected layer has **too many parameters**, it is difficult to stack several.

Convolution Layer

32x32x3 image



Keep spatial information.

It works **regardless of the size of the input image**.

There are **fewer parameters** than FC layer.

What is Depth?

RED Channel



Green Channel



Blue Channel

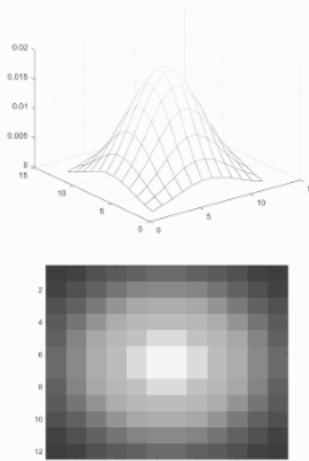


$$\text{Depth} = \text{RED Channel} + \text{Green Channel} + \text{Blue Channel}$$
$$3 = 1 + 1 + 1$$

Filter? Filtering?



Image



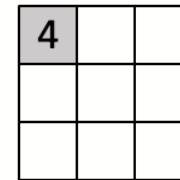
Kernel



Output

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

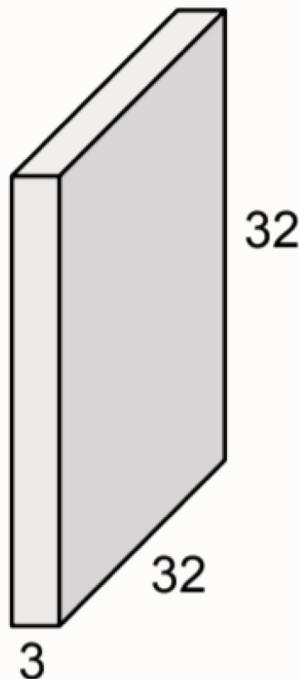


Convolved Feature

Filtering in an image is a process of **extracting a desired component**.

Convolution Layer (1)

32x32x3 image



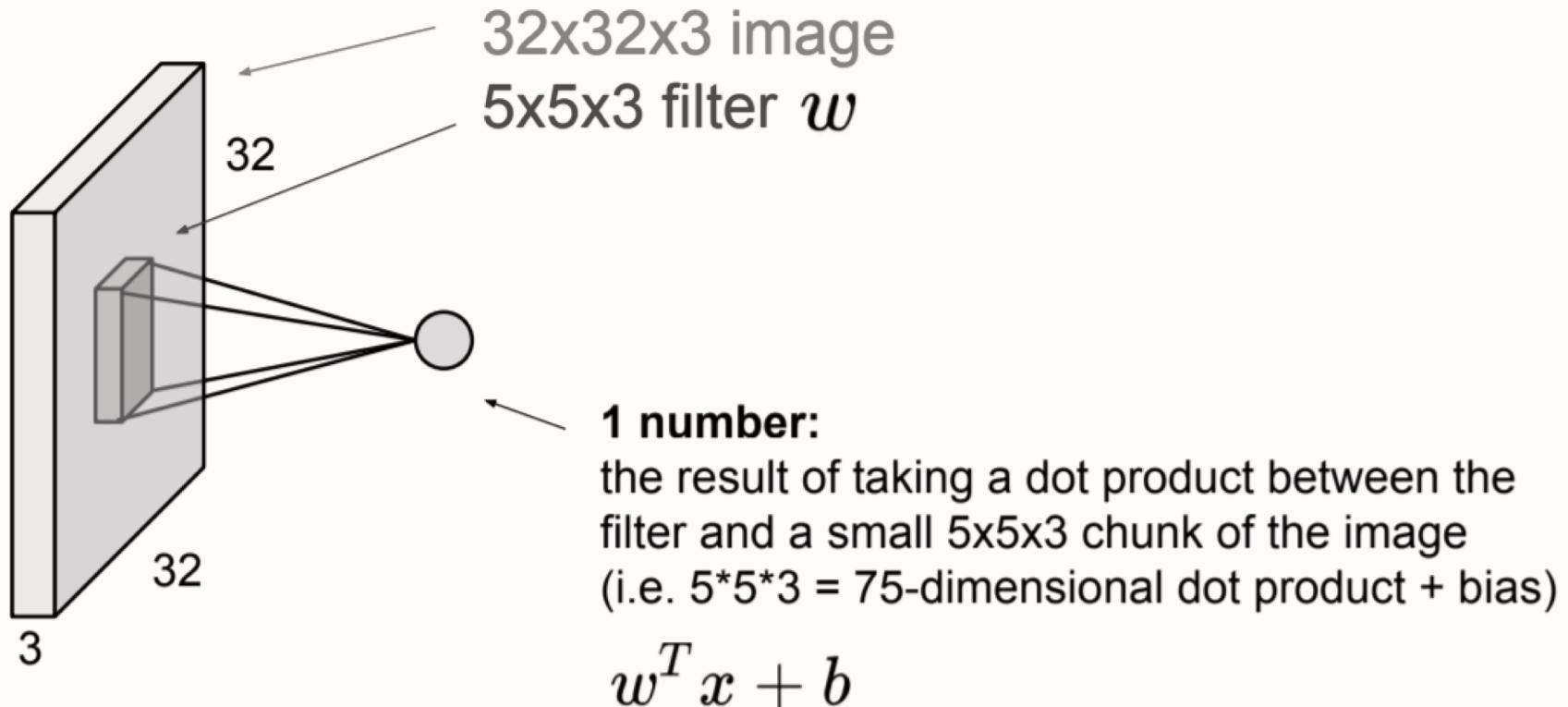
5x5x3 filter



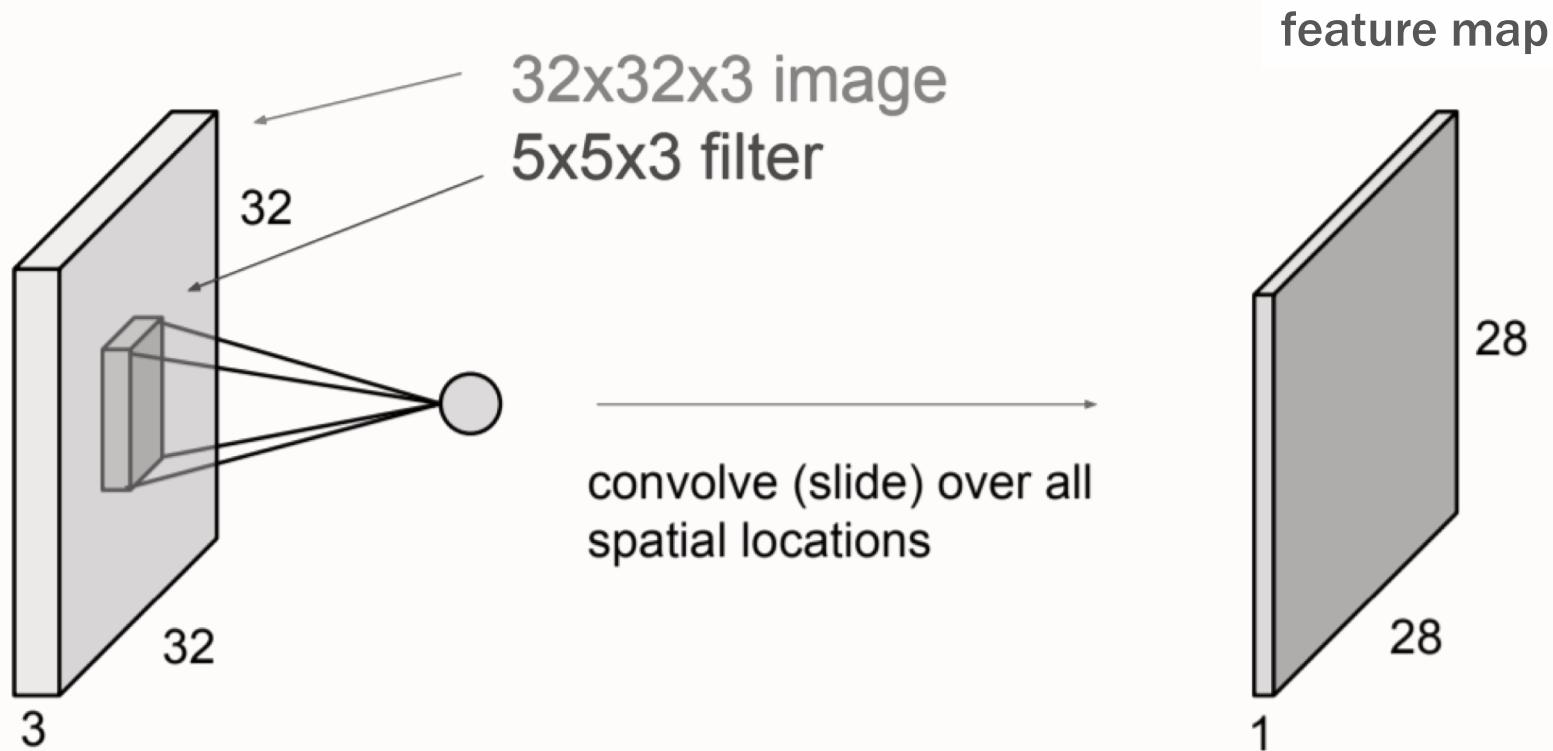
Filters always extend the full depth of the input volume

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

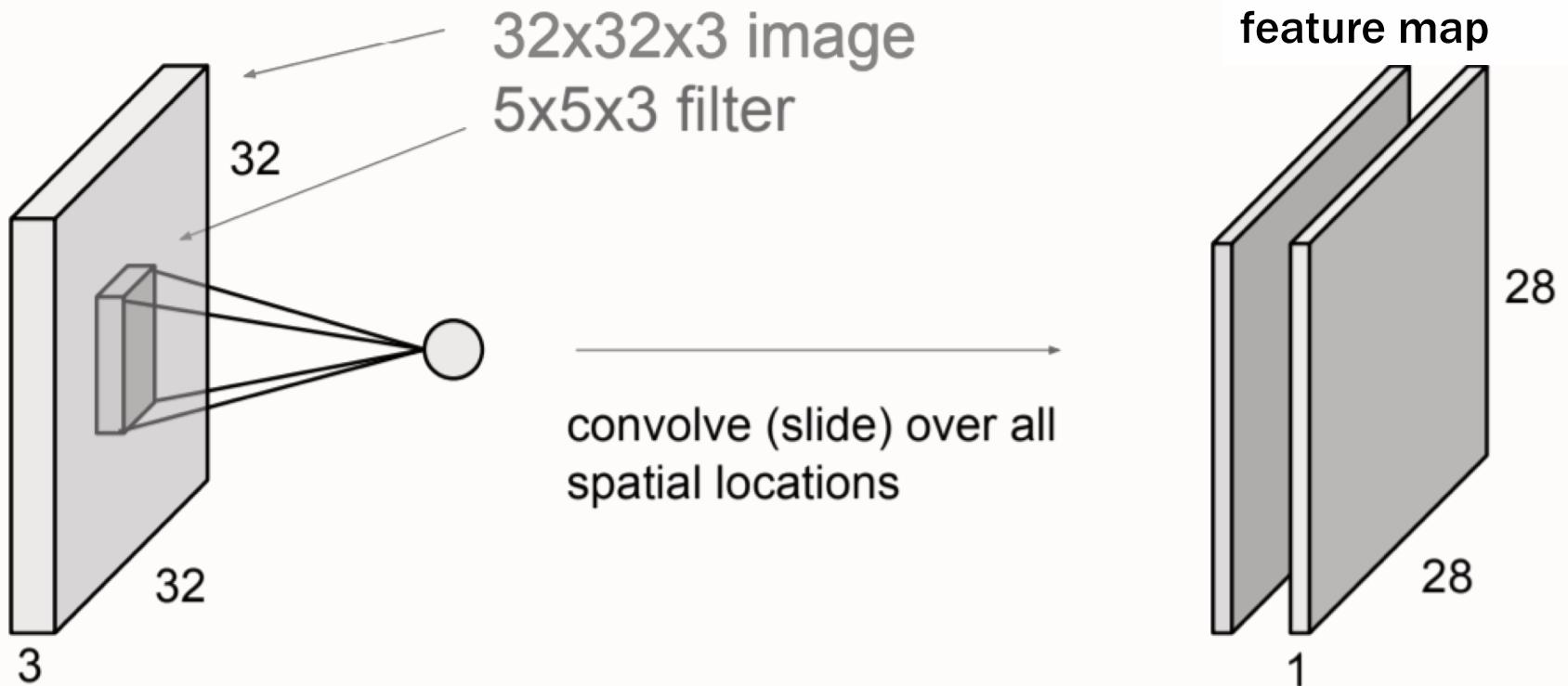
Convolution Layer (2)



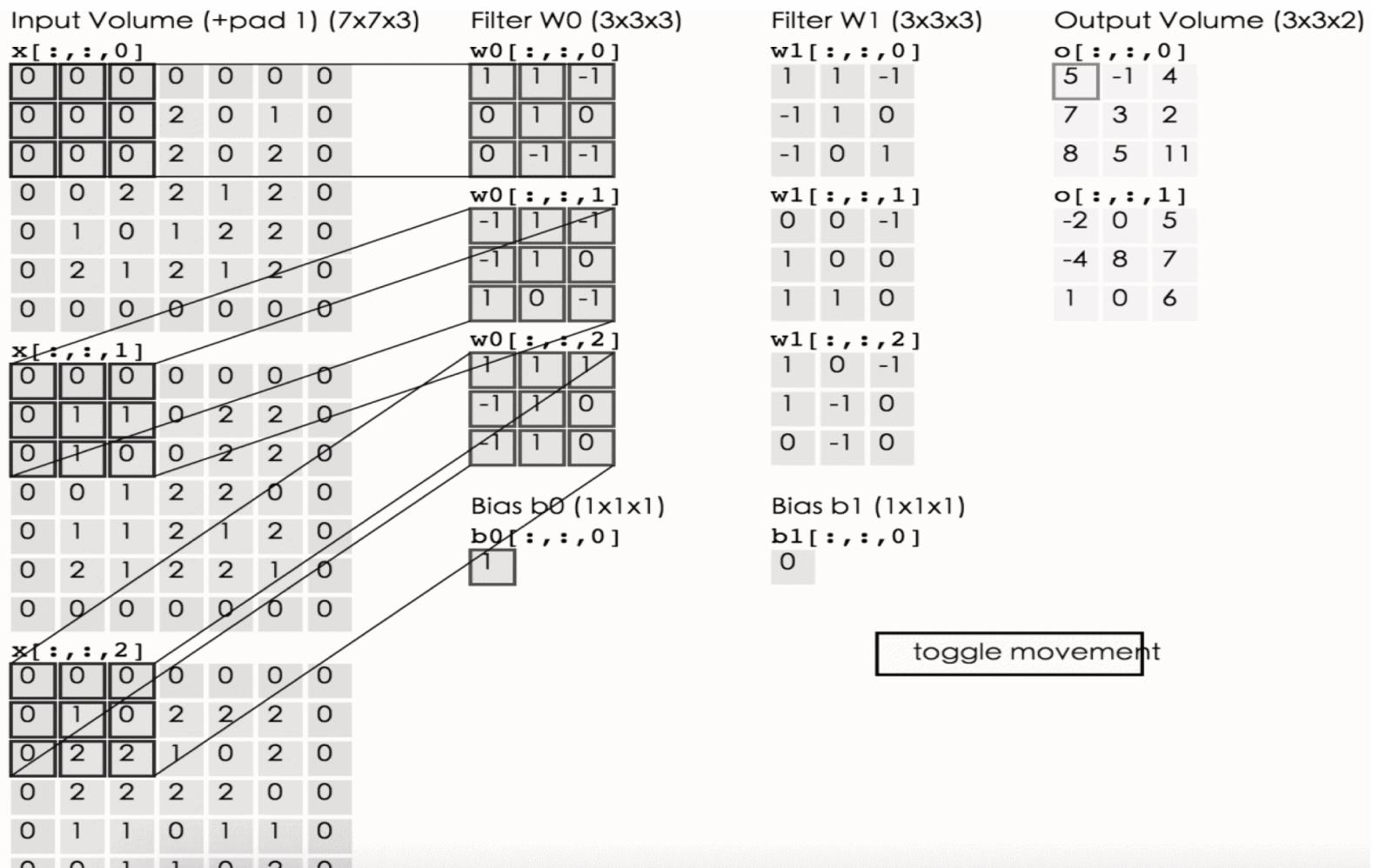
Convolution Layer (3)



Convolution Layer (4)

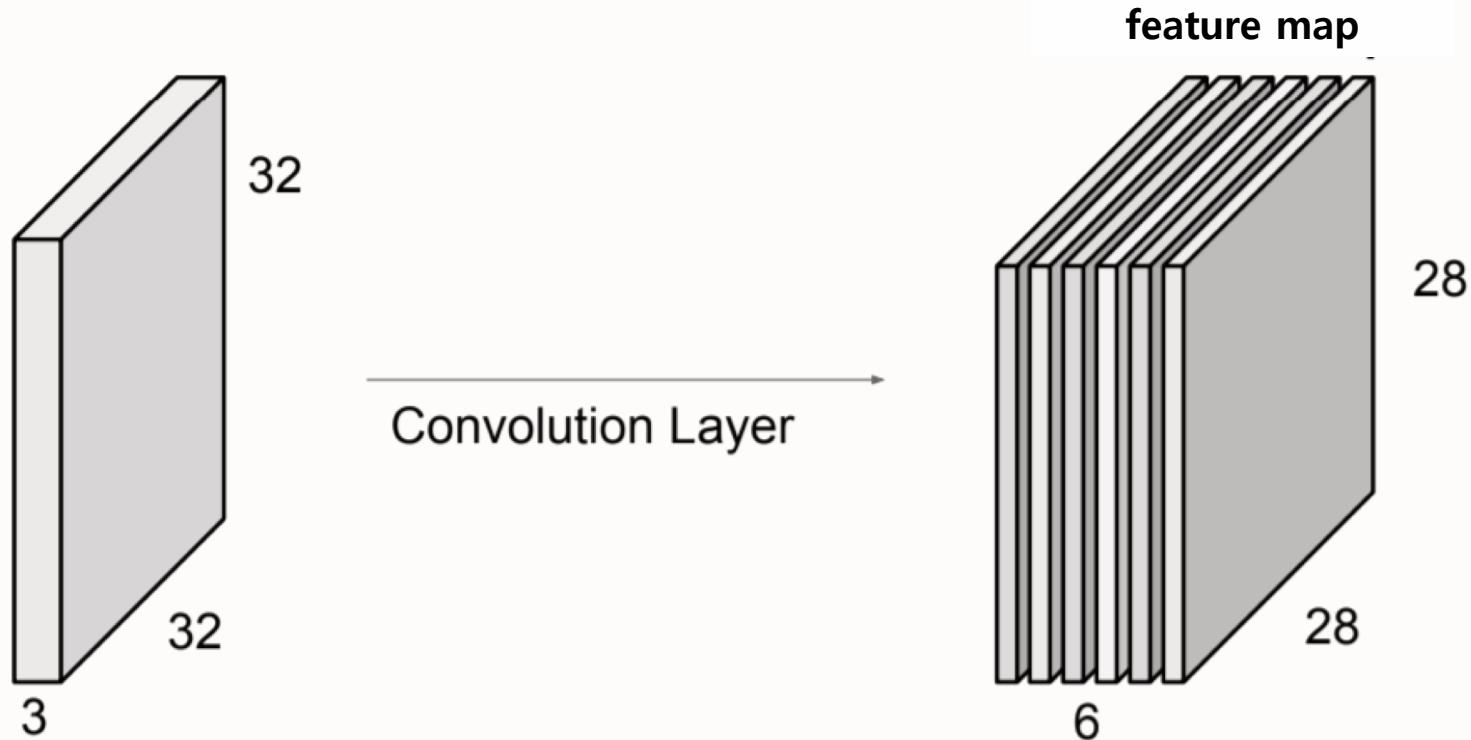


Convolution Layer (5)



Convolution Layer (6)

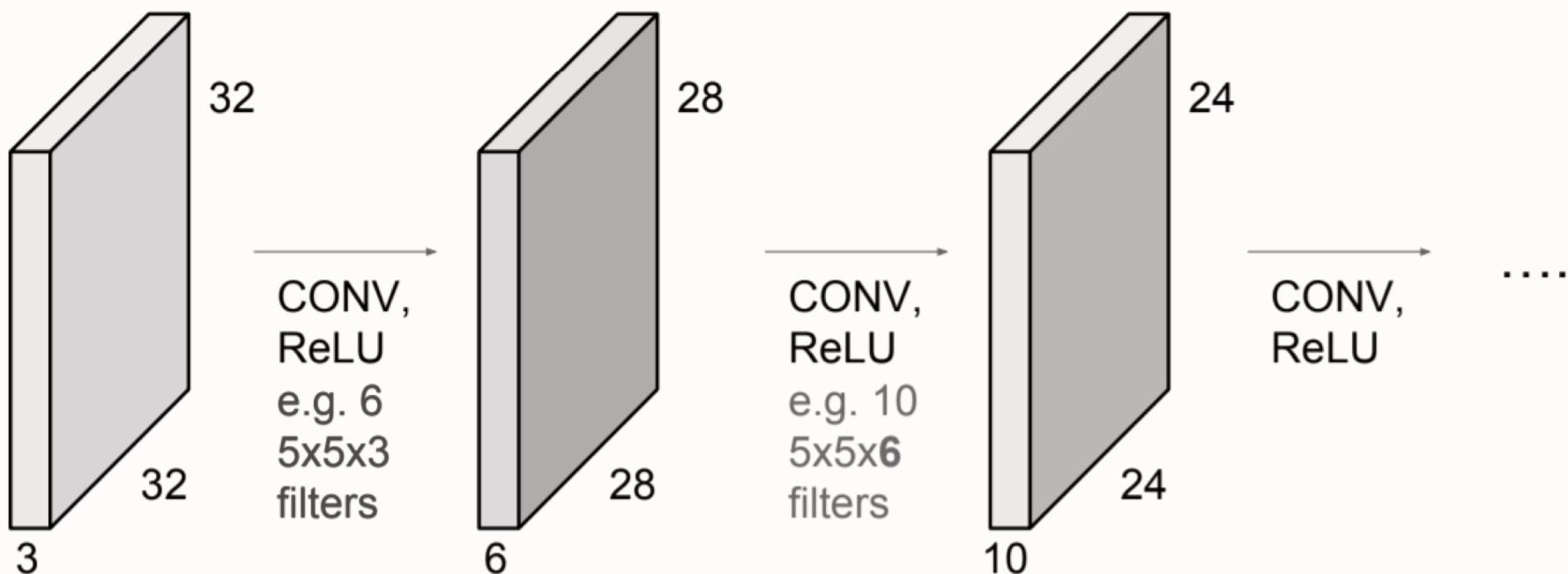
For example, if we had 6 5x5 filters, we'll get 6 separate feature map



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

Convolution Layer (7)

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

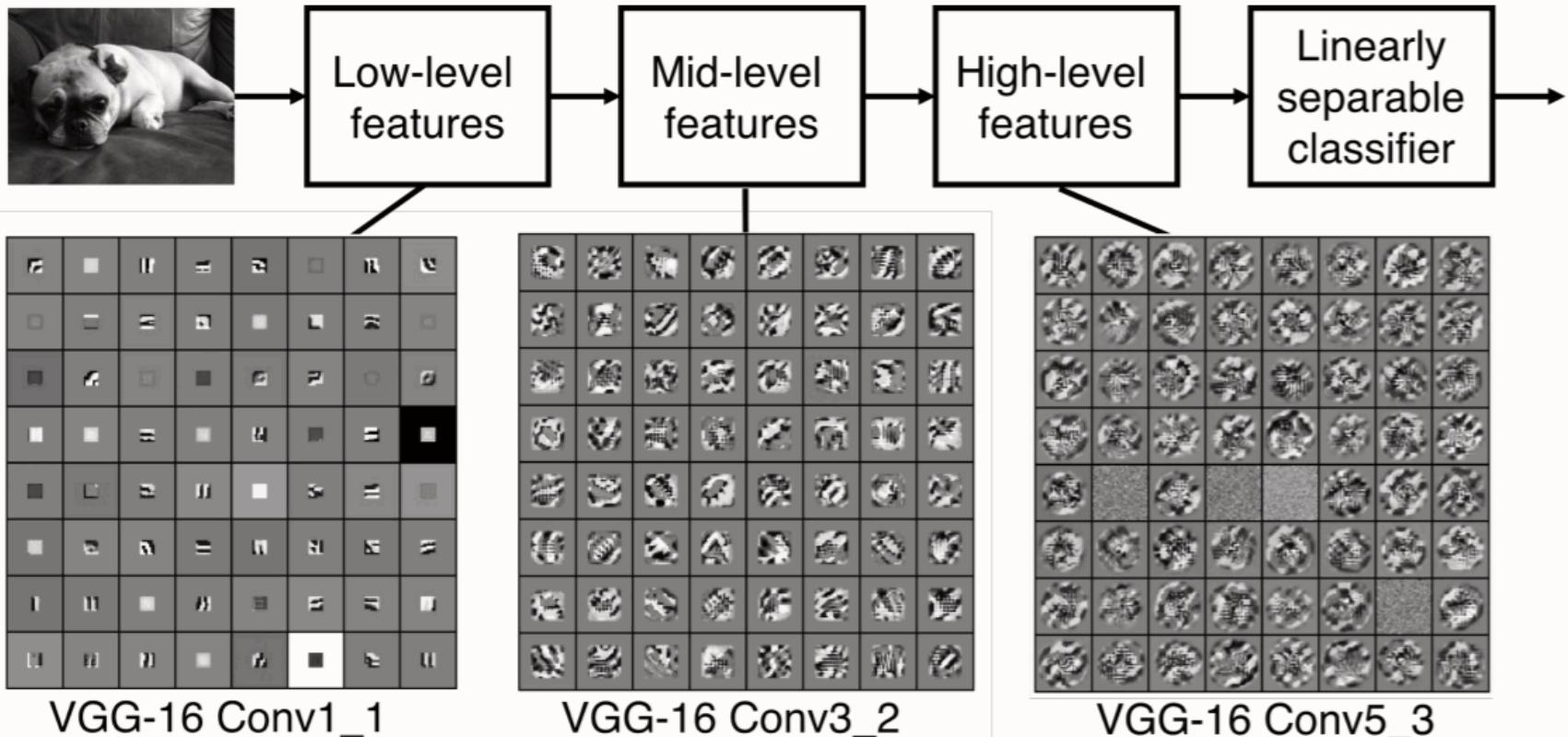


Convolution Layer (8)

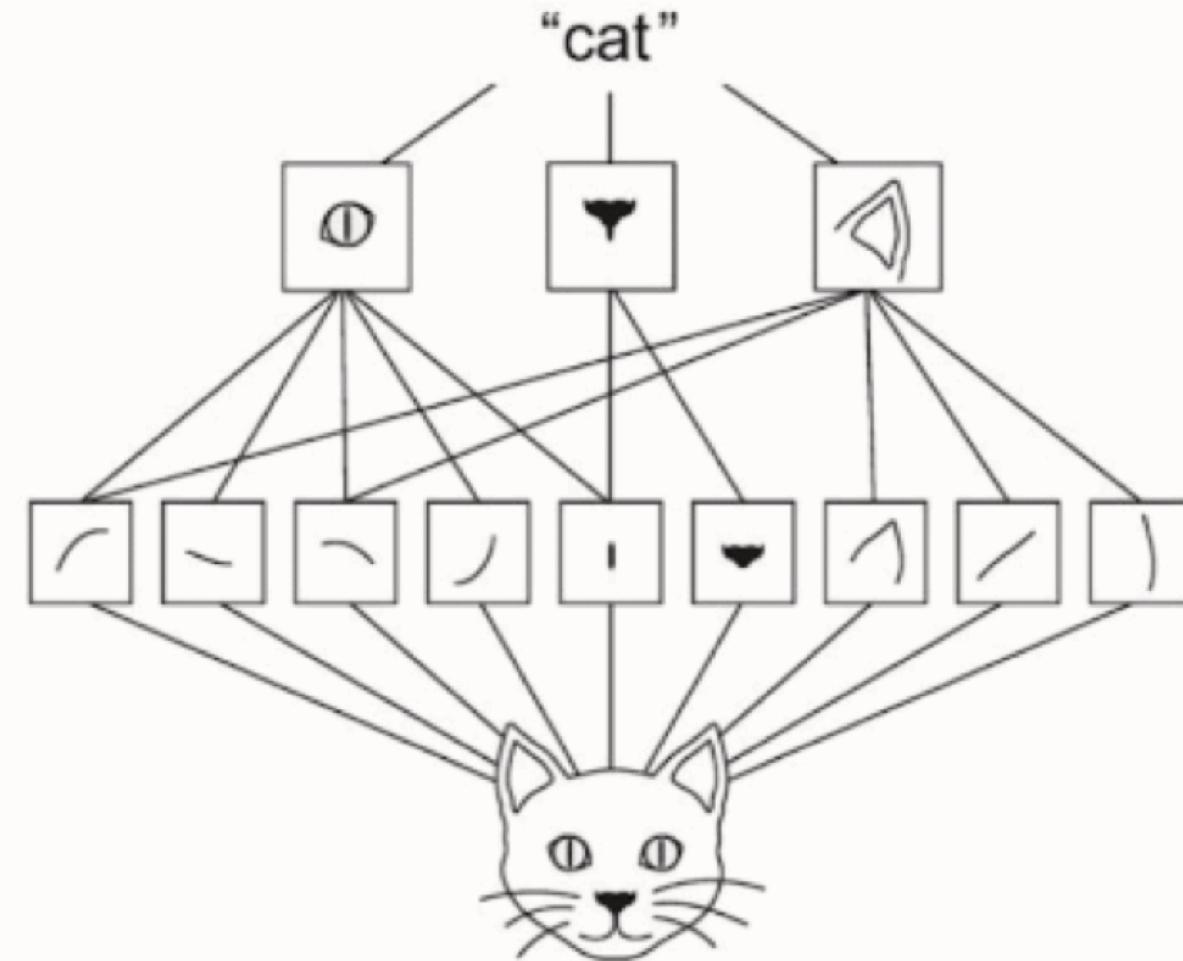
Preview

[Zeiler and Fergus 2013]

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

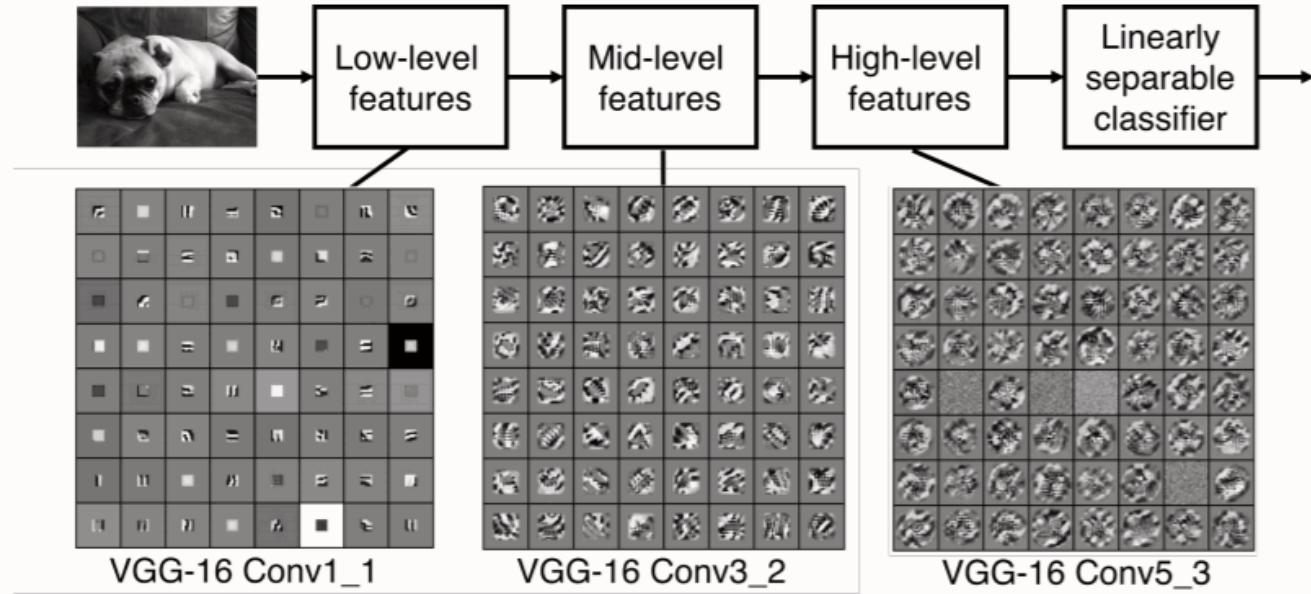


Convolution Layer (9)



Convolution Layer (10)

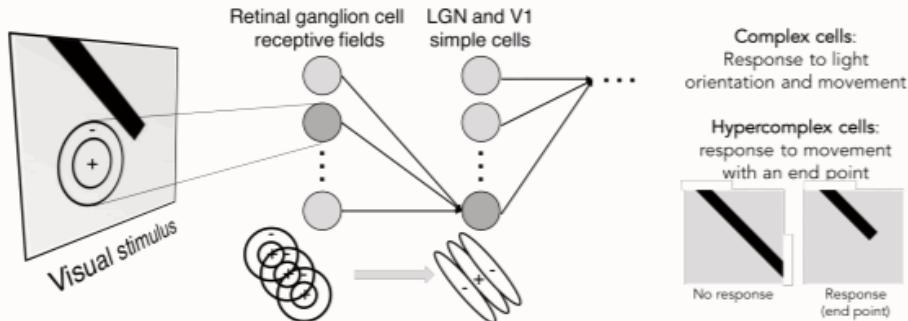
Preview



VGG-16 Conv1_1

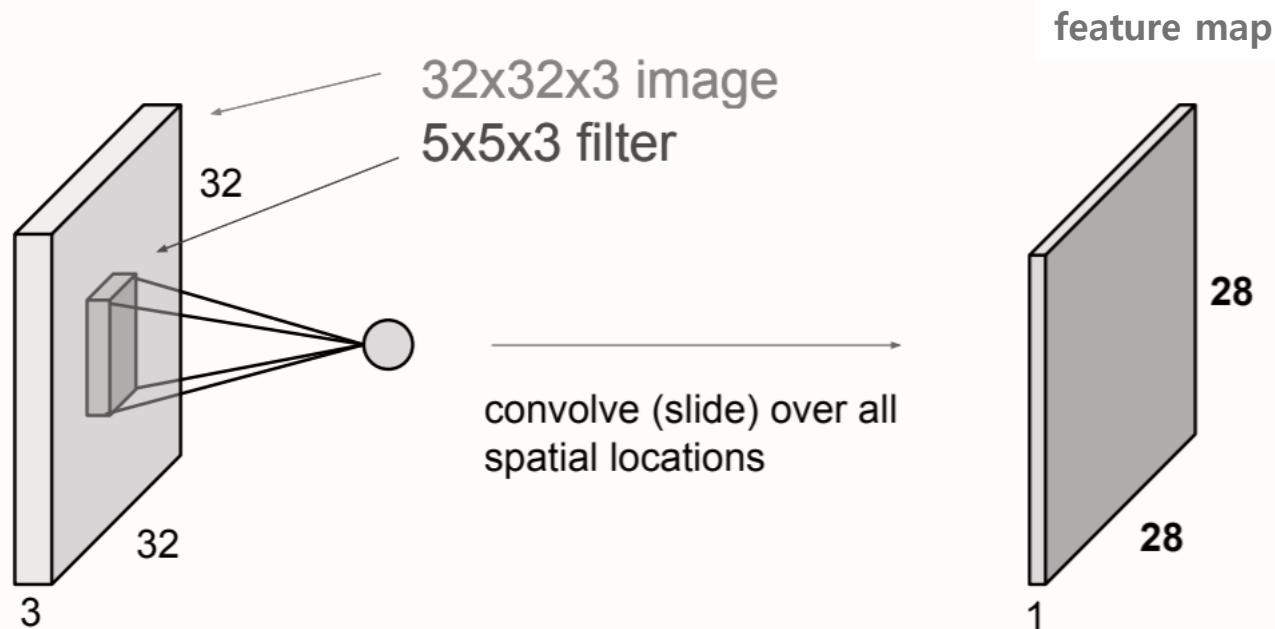
VGG-16 Conv3_2

VGG-16 Conv5_3



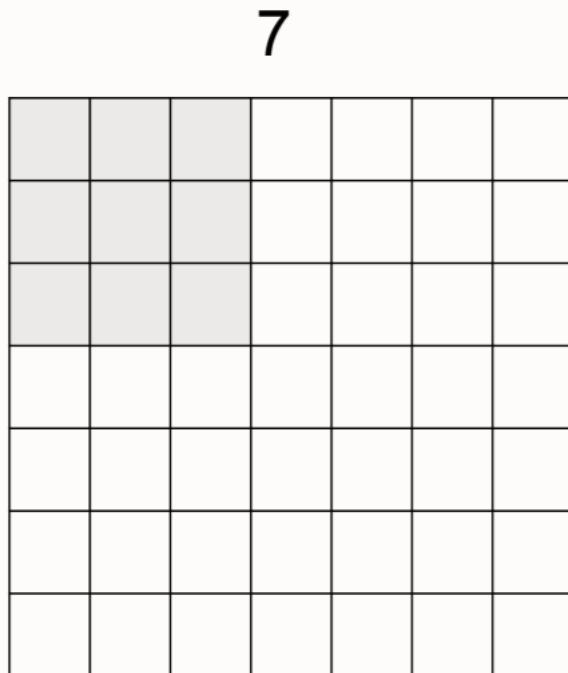
Convolution Layer (11)

A closer look at spatial dimensions:



Convolution Layer (12)

A closer look at spatial dimensions:

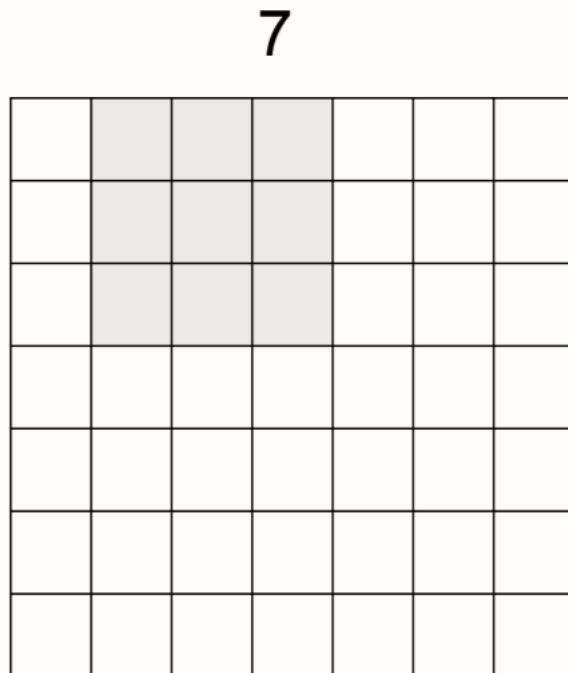


7x7 input (spatially)
assume 3x3 filter

7

Convolution Layer (13)

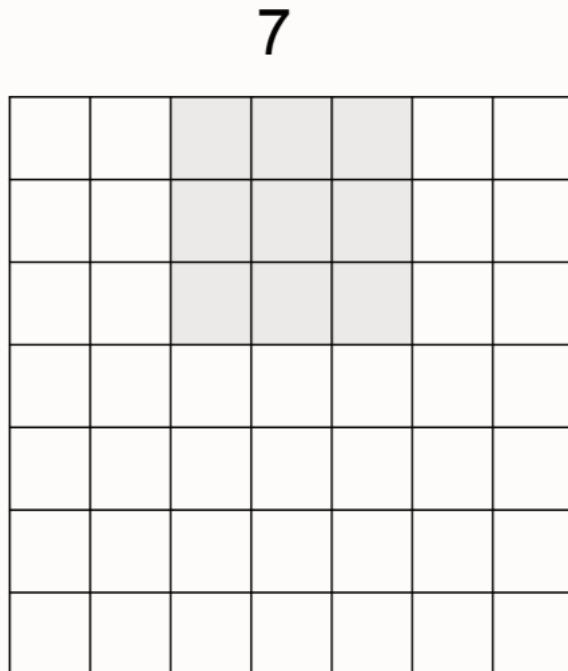
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

Convolution Layer (14)

A closer look at spatial dimensions:

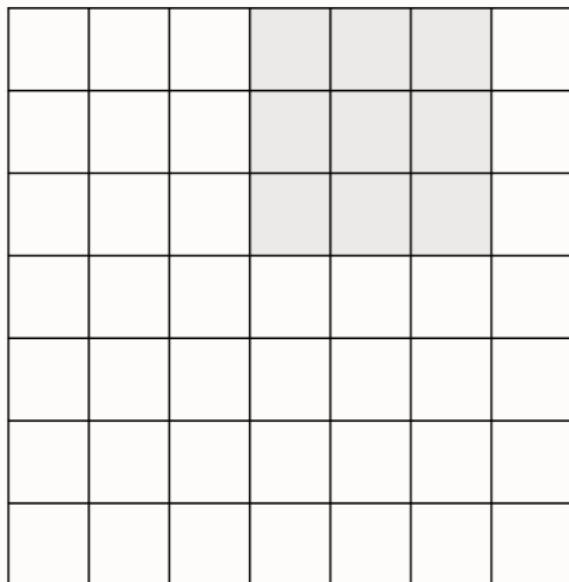


7x7 input (spatially)
assume 3x3 filter

Convolution Layer (15)

A closer look at spatial dimensions:

7

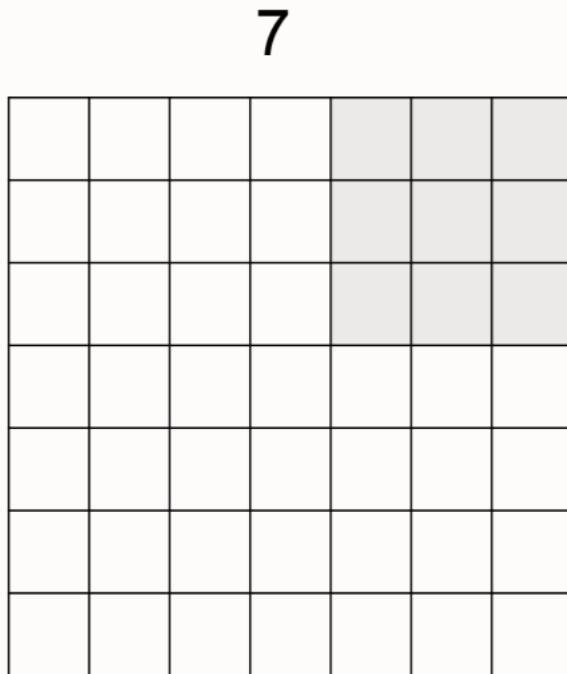


7x7 input (spatially)
assume 3x3 filter

7

Convolution Layer (16)

A closer look at spatial dimensions:

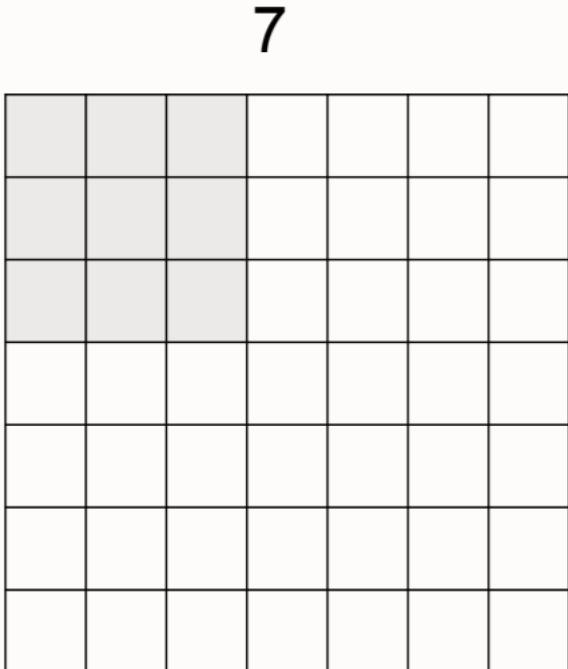


7x7 input (spatially)
assume 3x3 filter

=> 5x5 output

Convolution Layer (17)

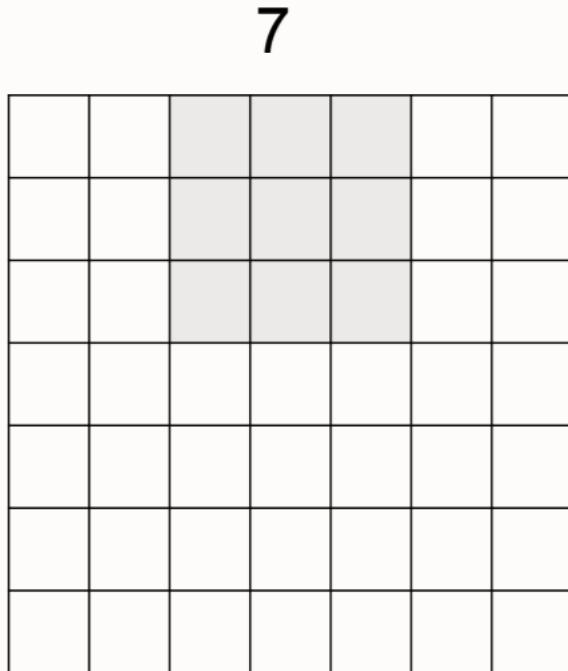
A closer look at spatial dimensions:



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

Convolution Layer (18)

A closer look at spatial dimensions:

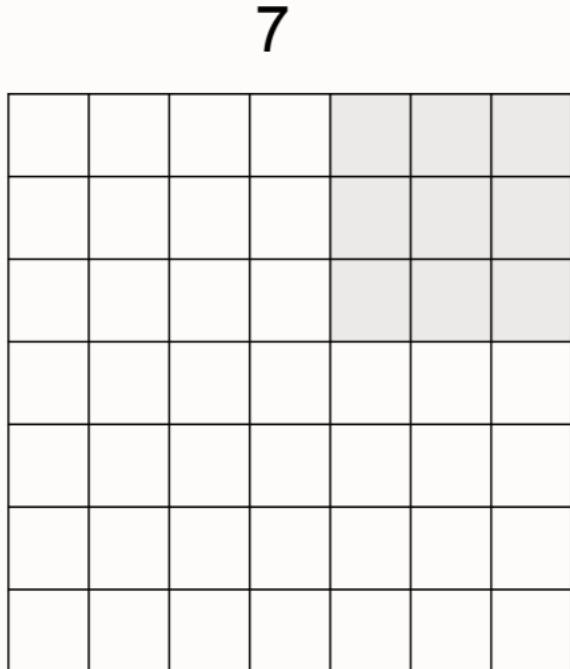


7

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

Convolution Layer (19)

A closer look at spatial dimensions:

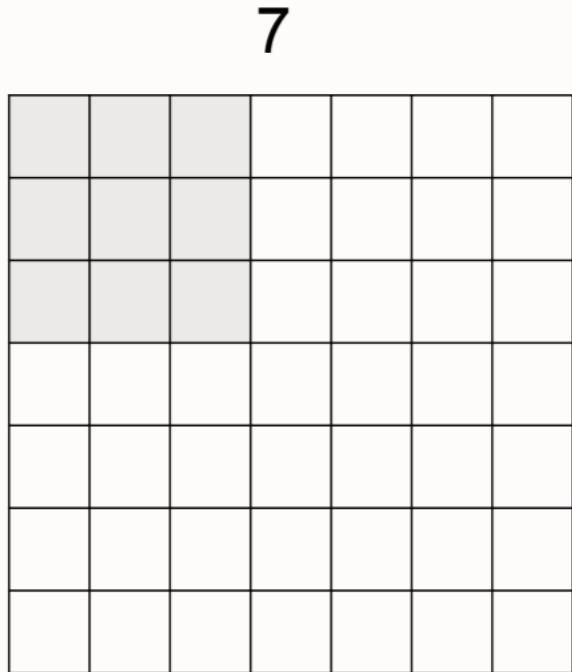


7

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

Convolution Layer (20)

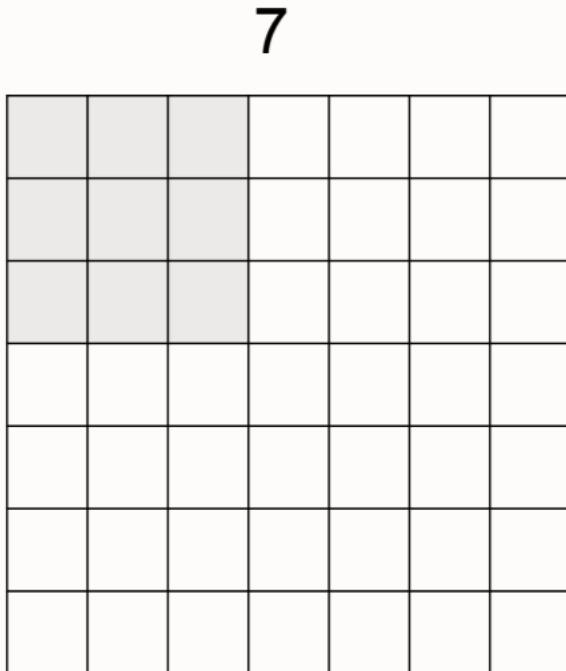
A closer look at spatial dimensions:



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

Convolution Layer (21)

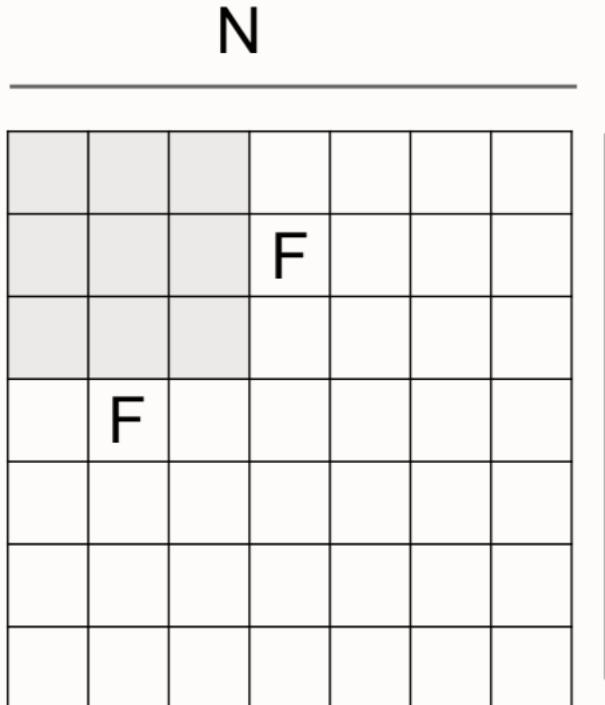
A closer look at spatial dimensions:



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

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

Convolution Layer (22)



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

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

Convolution Layer (23)

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:)

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

Convolution Layer (24)

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!

Convolution Layer (25)

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

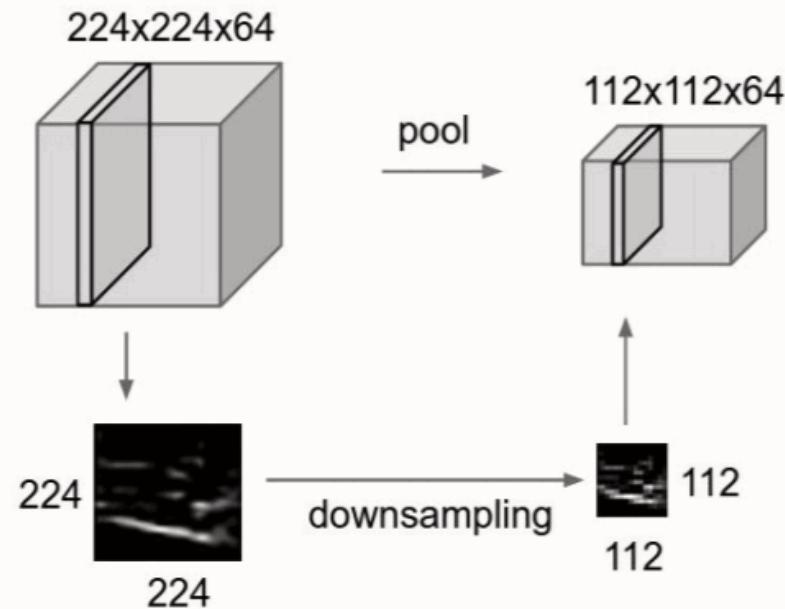
$F = 7 \Rightarrow$ zero pad with 3

By using zero padding, you can **maintain the size of the image**.

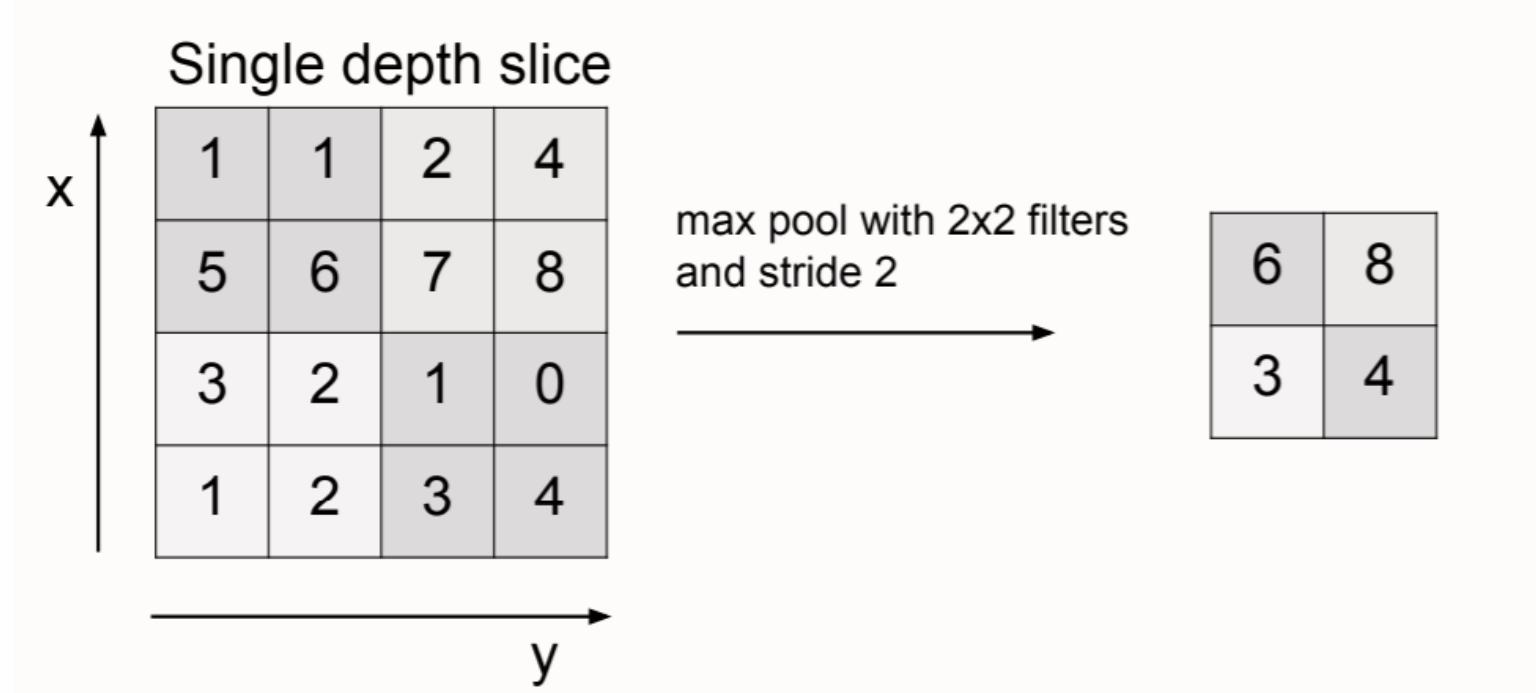
The filter allows you to see the feature map in **more detail**.

Pooling Layer

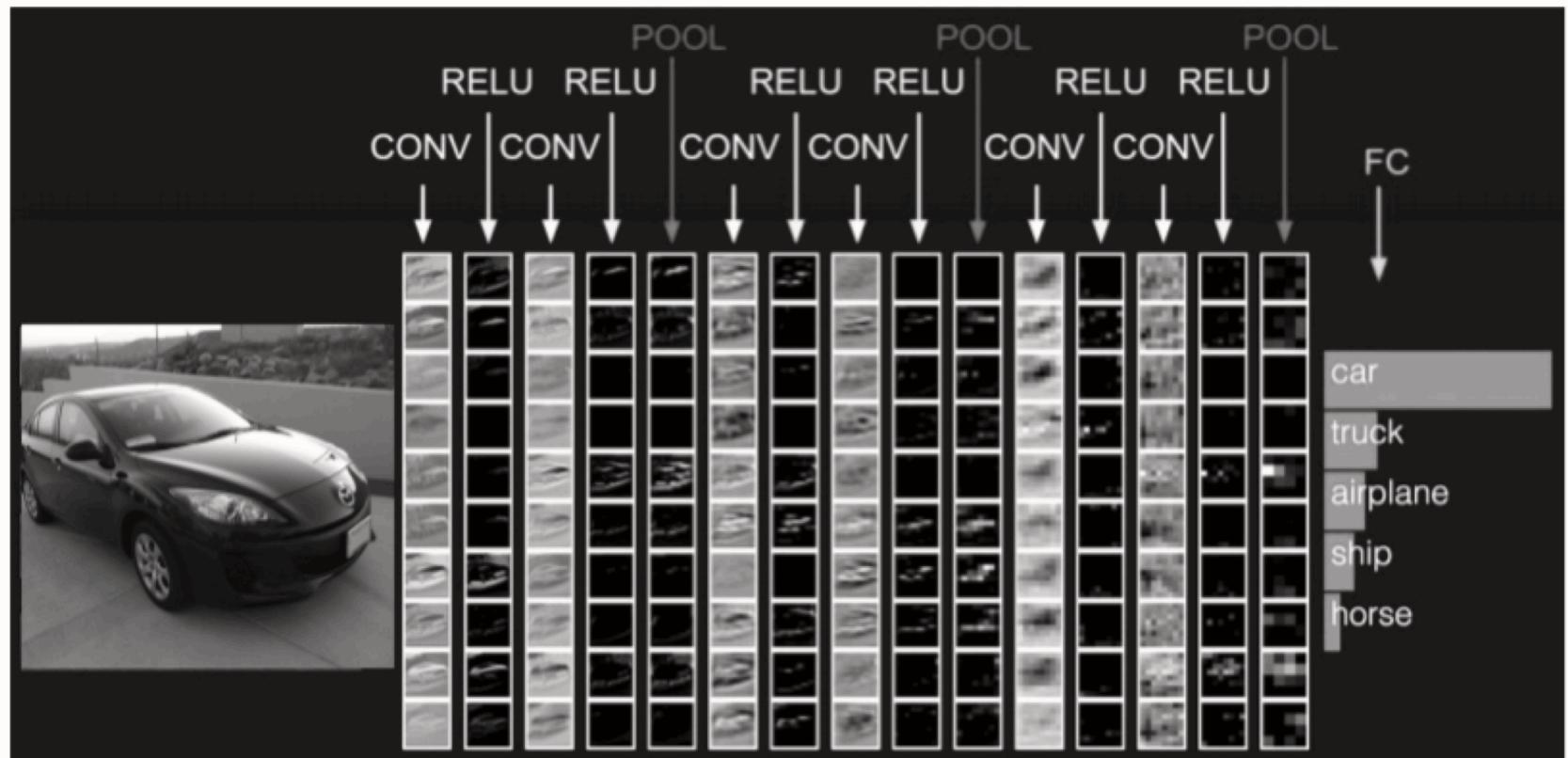
- makes the representations smaller and more manageable
- operates over each activation map independently:



Pooling Layer



Summary



Today's generic NN architecture.