

CNN

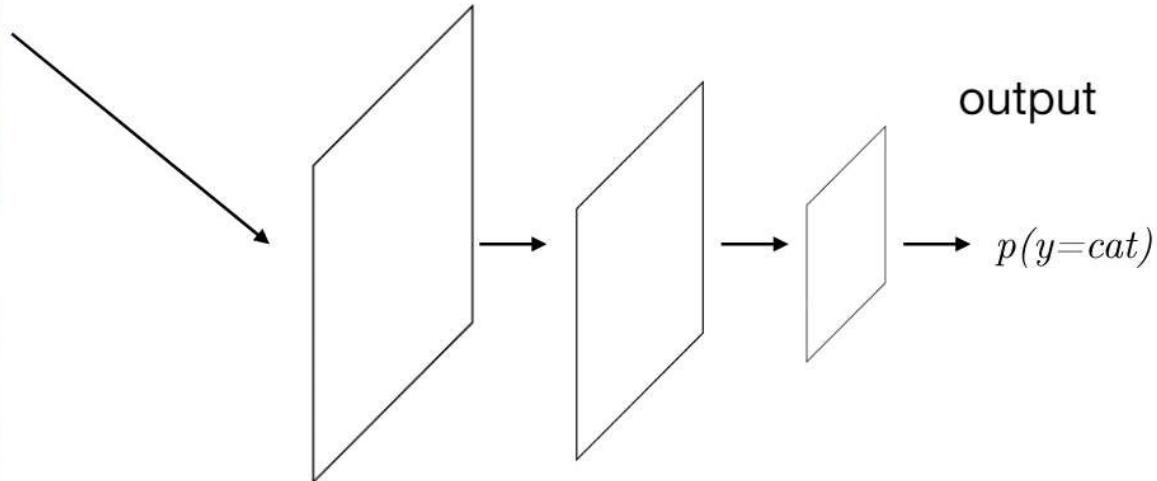
CNN for image classification



Image Source:
twitter.com%2Fcats&psig=AOvVaw30_o-PCM-K21DiMAJQimQ4&ust=1553887775741551



Image Source: <https://www.pinterest.com/pin/244742560974520446>



Object detection



Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 779-788).

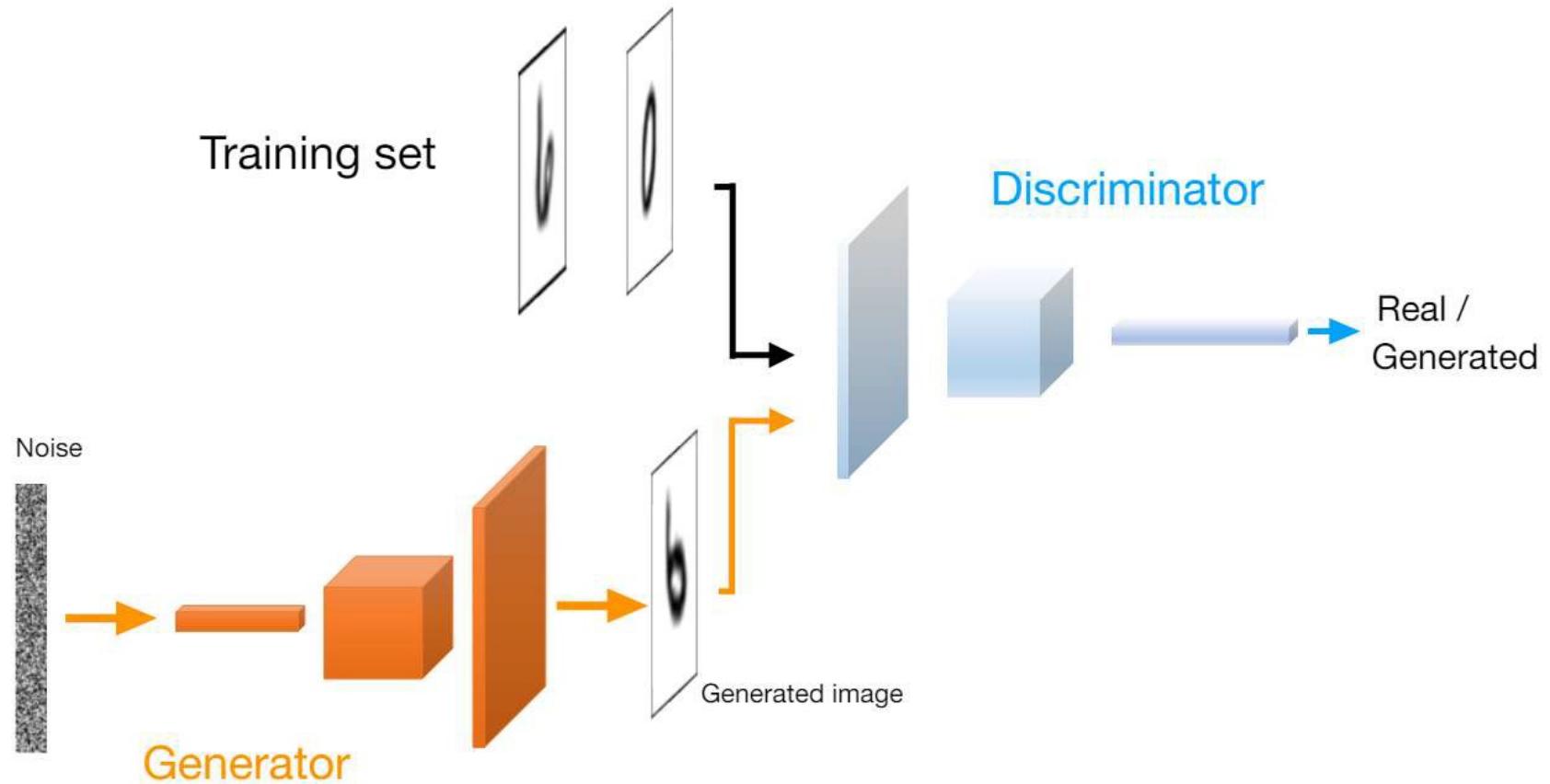
Image segmentation



Figure 2. **Mask R-CNN** results on the COCO test set. These results are based on ResNet-101 [15], achieving a *mask AP* of 35.7 and running at 5 fps. Masks are shown in color, and bounding box, category, and confidences are also shown.

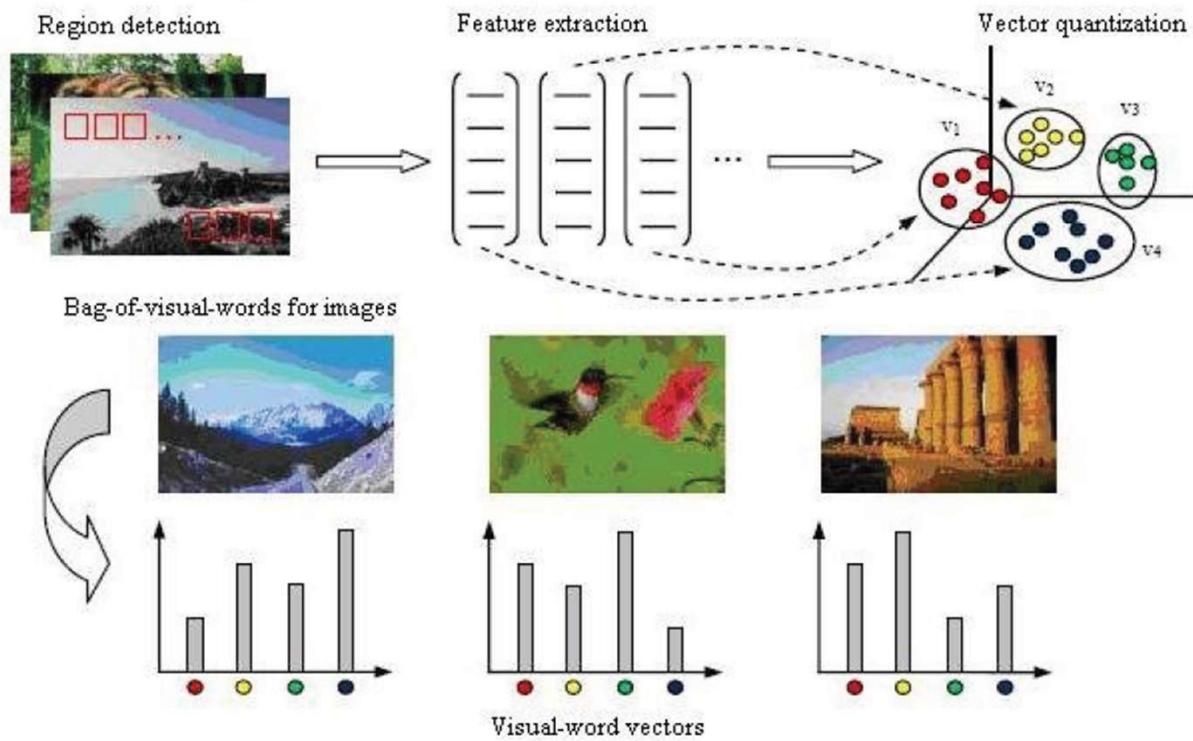
He, Kaiming, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask R-CNN." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961-2969. 2017.

Image generation



History of Computer Vision

Before Deep Learning: Hand-Crafted Features



History of Computer Vision

Effect of Deep Learning: Comparison between Deep Learning and Traditional Models

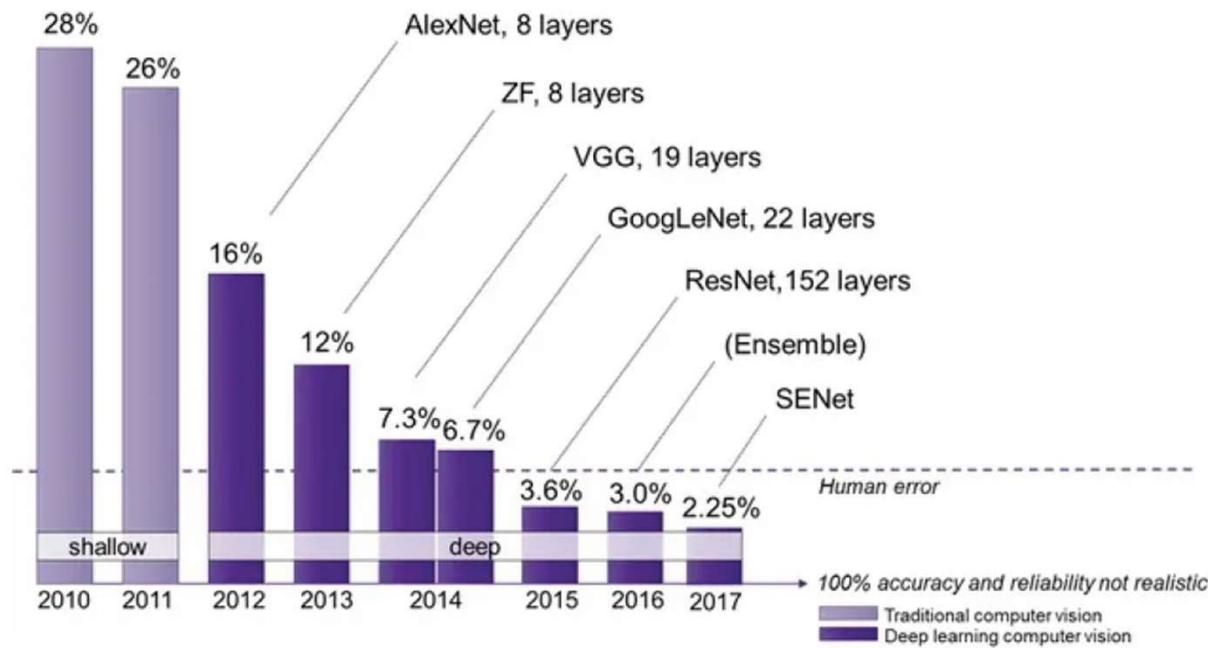
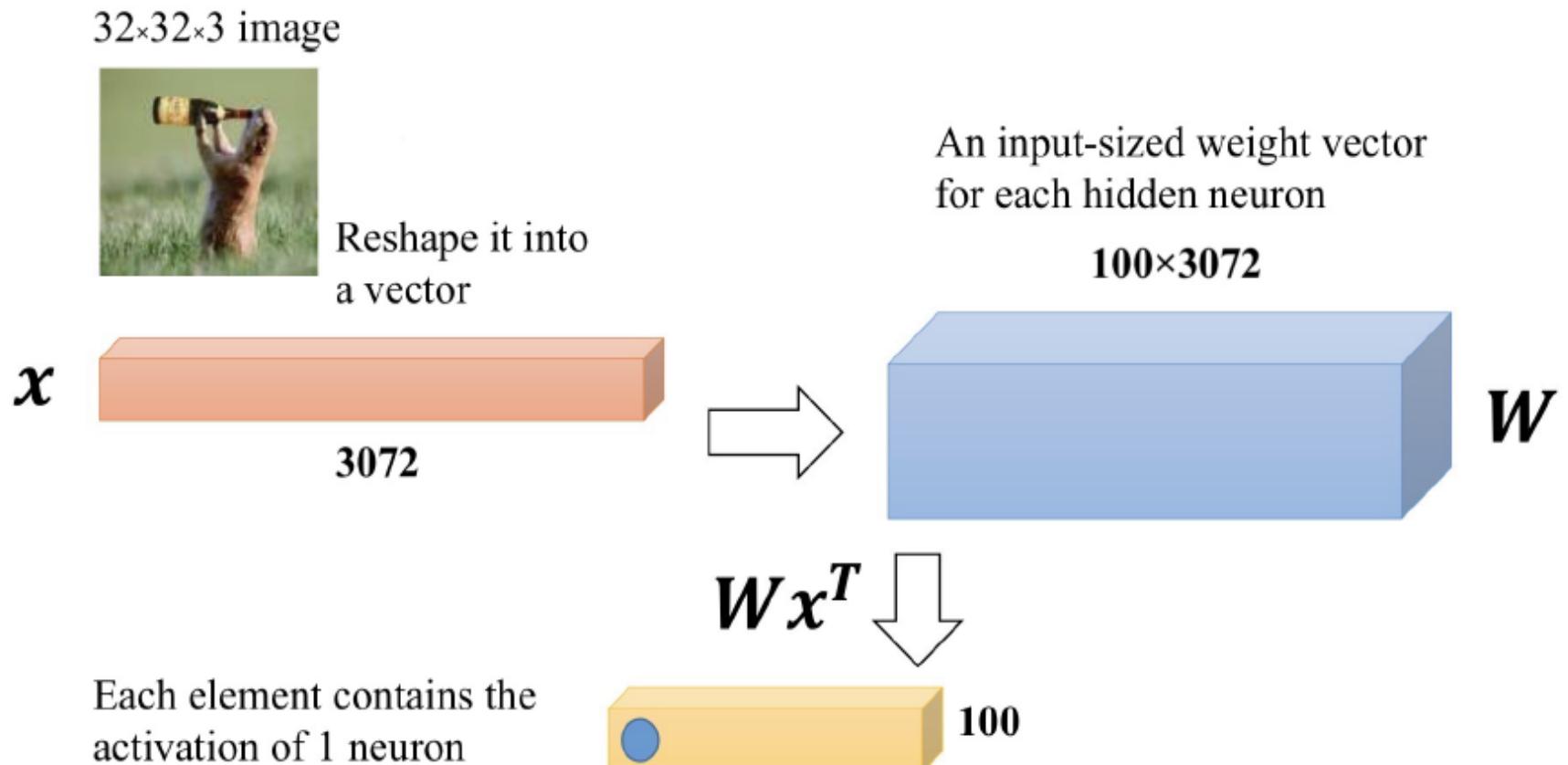


Figure: ImageNet Competition Results Over Time, [Source](#)

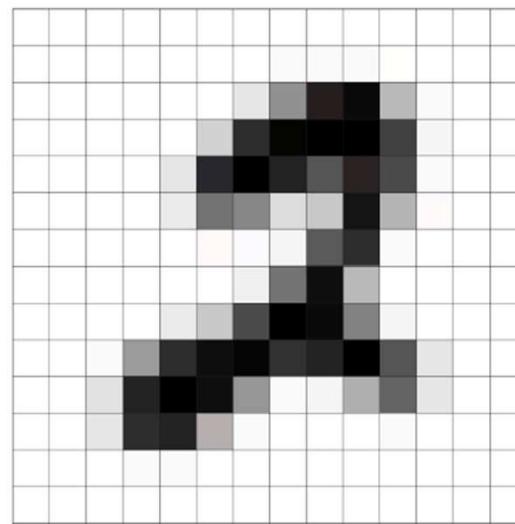
Using Fully Connected Artichecture

- What we've been using?
 - Dense Vector Multiplication



- What is the Problem with FC ?
- Neurons in a single layer operate independently and **do not share connections**, even for inputs.
- To be shift-invariant, **many samples** (in various time or locations) must be shown to them.
- **Regular Neural Nets don't scale well to full images.**
 - Parameters would **add up** quickly!
 - Full connectivity is wasteful and the huge number of parameters would quickly **lead to overfitting**.

Multi layer perceptron on image data

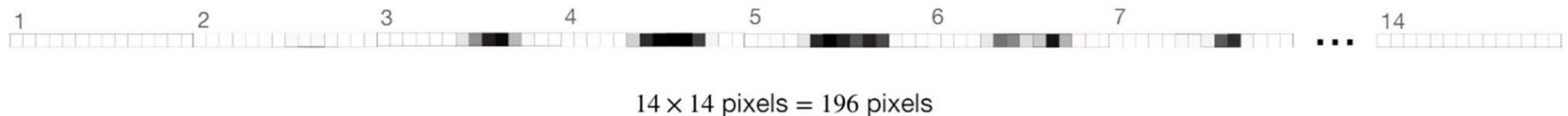


14x14-dimensional
image

Concatenate the rows



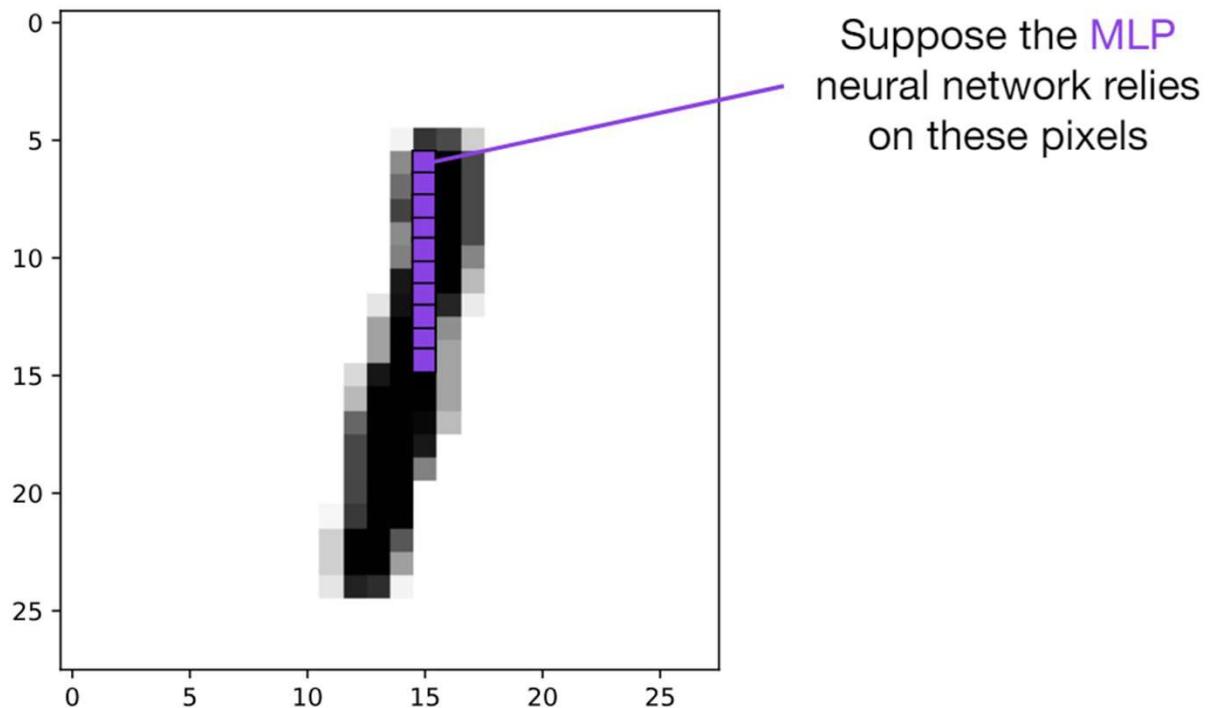
get a 196-dimensional row vector



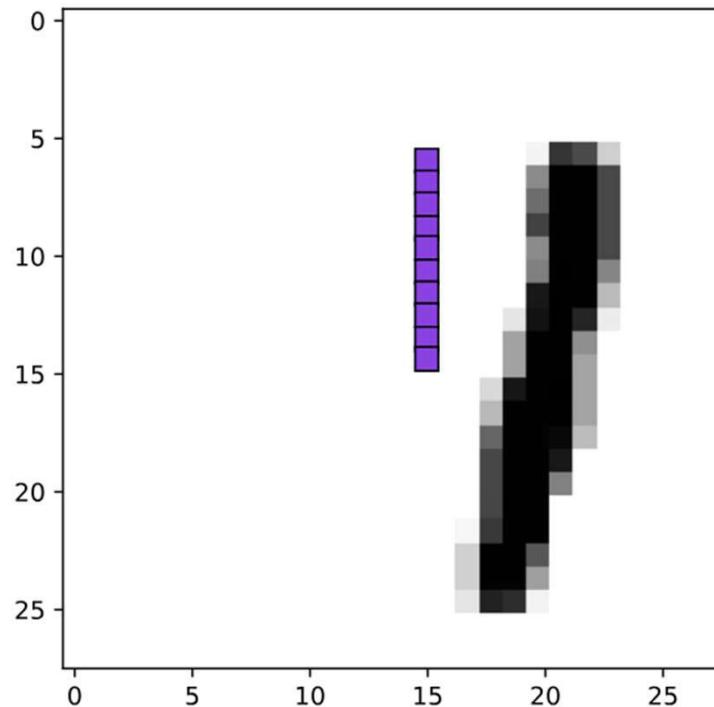
Training examples



Problem with MLP for images



Problem with MLP for images



An **MLP** will not
recognize the digit
if it is in slightly
different position

Why Image Classification is Hard

Different lighting, contrast, viewpoints, etc.



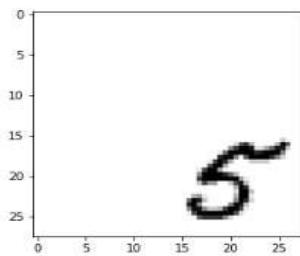
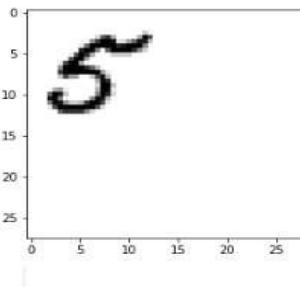
Image Source:
twitter.com%2Fcats&psig=AOvVaw30_o-PCM-K21DiMAJQimQ4&ust=1553887775741551



Image Source: https://www.123rf.com/photo_76714328_side-view-of-tabby-cat-face-over-white.html



Or even simple translation



This is hard for traditional methods like multi-layer perceptrons, because the prediction is basically based on a sum of pixel intensities

Main concept behind convolutional neural networks

- Sparse-connectivity: A single element in the feature map is connected to only a small patch of pixels. (This is very different from connecting to the whole input image, in the case of multilayer perceptrons.)
- Parameter-sharing: The same weights are used for different patches of the input image.
- Many layers: Combining extracted local patterns to global patterns

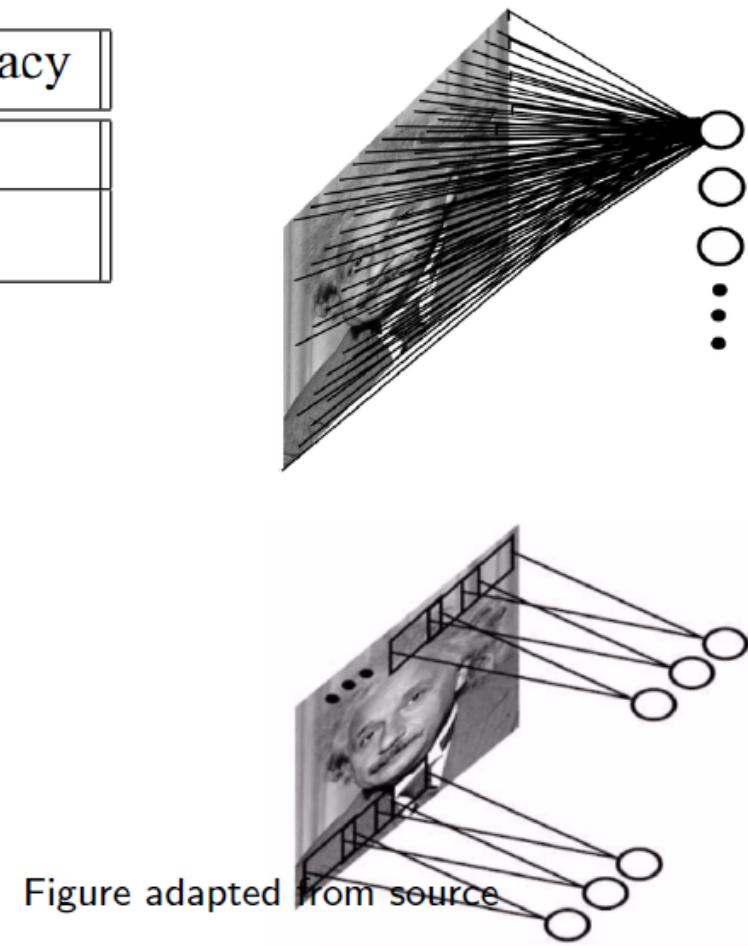
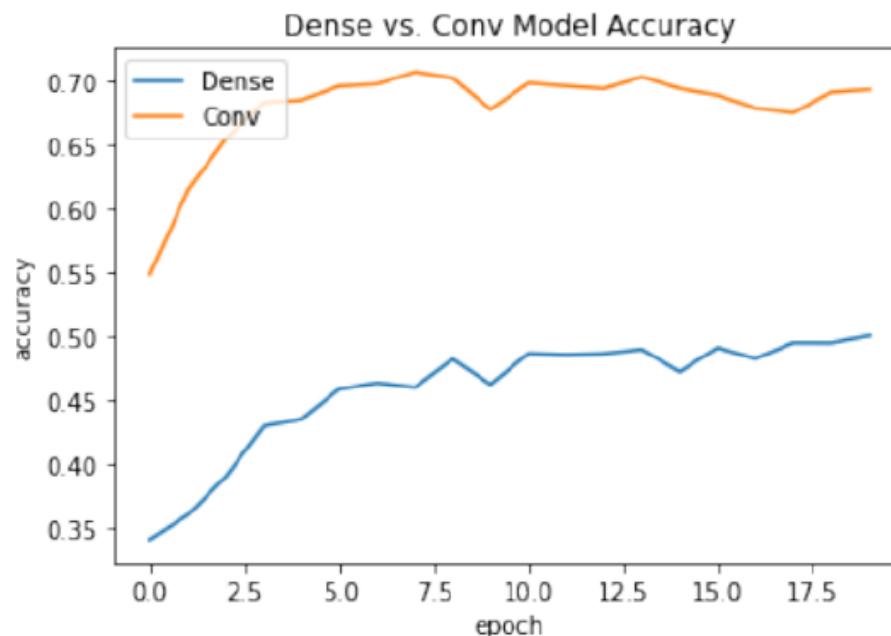
Fully connected Vs Convolution

Why Use Convolution:

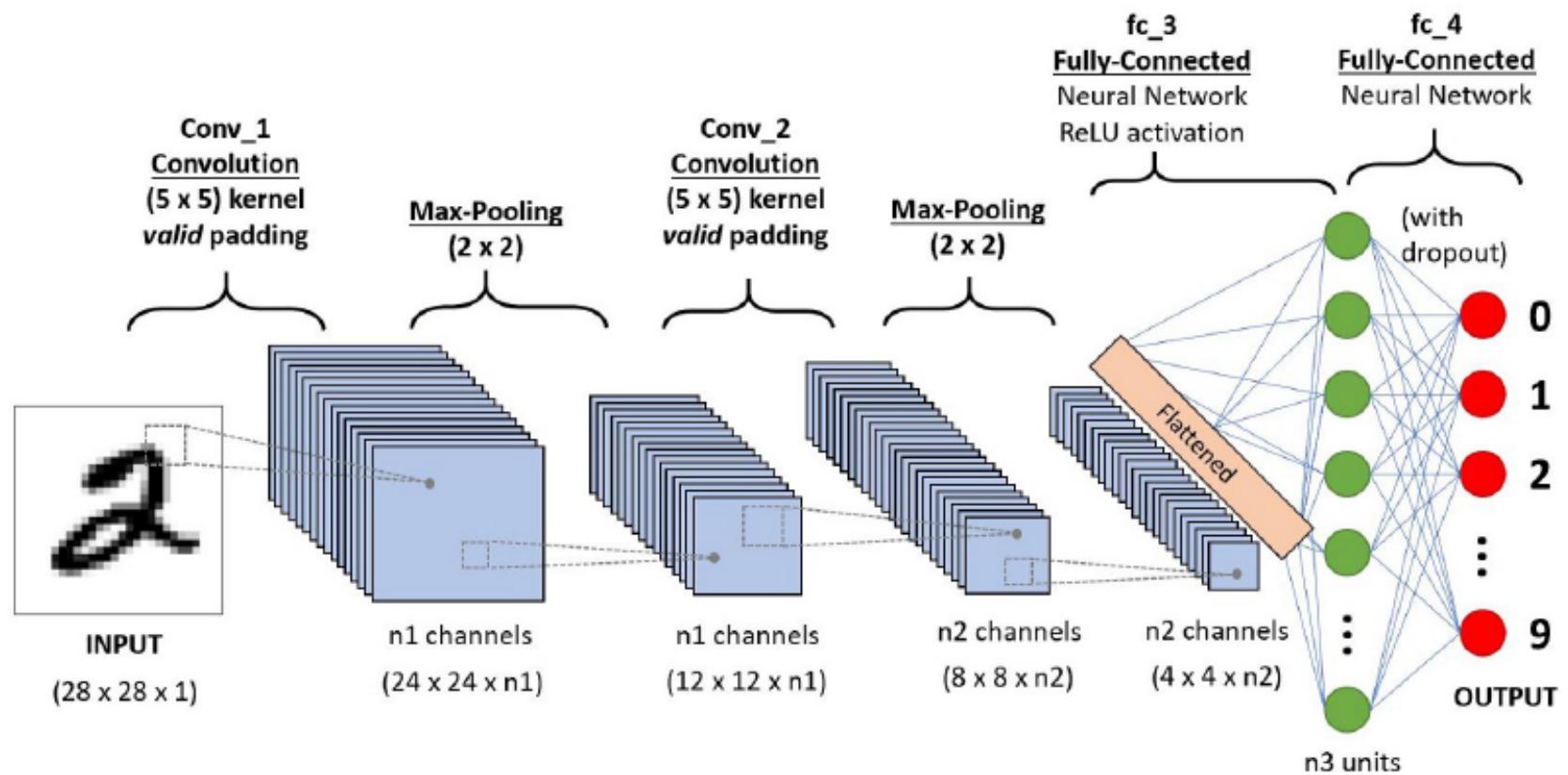
- **Exploit spatial structure:** Convolution sees local patterns by using filters over small patches of the image.
- **Reduce parameters:** By sharing weights across the image, convolution reduce the number of parameters.
- **Build hierarchical features:** Convolution starts with small patterns, then combines them to recognize more complex patters.

FC vs CNN

Model Architecture	# Params	Test Accuracy
CNN	2M	67.8%
FCN	9M	50.6%



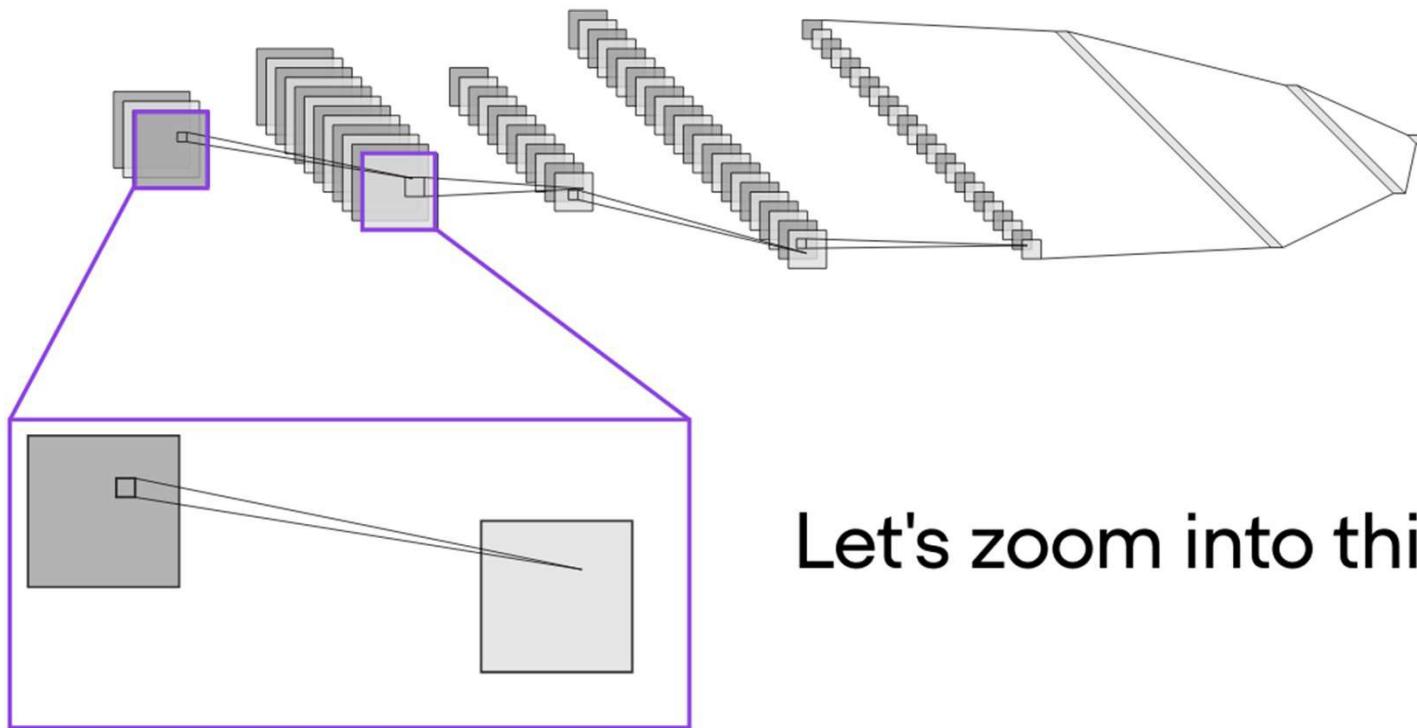
Looking at convolutional layers in more detail



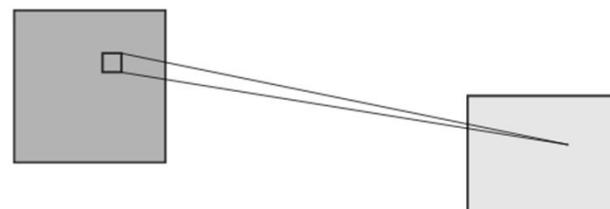
CNN basic structure

Three Main Types of Layers

- **Convolutional Layer**
 - Output of neurons are connected to **local regions** in the input.
 - Applying the **same filter** across the entire image.
 - CONV layer's parameters consist of a set of **learnable filters**.
- **Pooling Layer**
 - Performs a **downsampling** operation along the spatial dimensions.
- **Fully-Connected Layer**
 - Typically used in the final stages of the network to combine high-level features and **make predictions**.

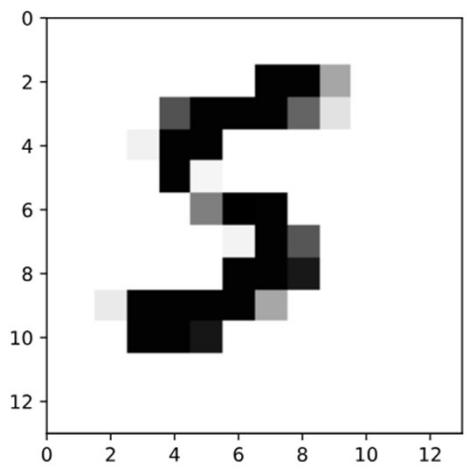


Let's zoom into this part

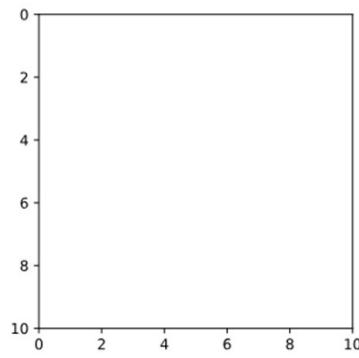


Input (image)

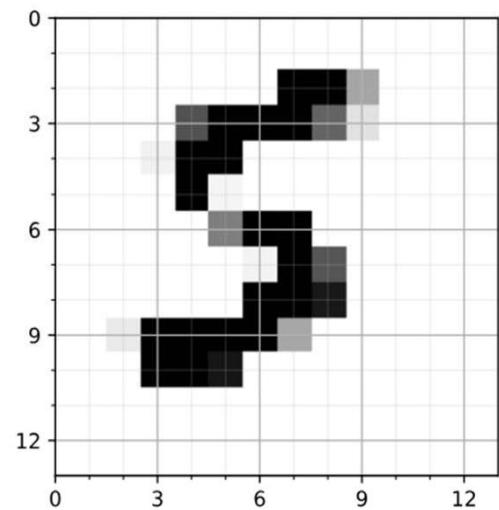
Feature map



Input (image)



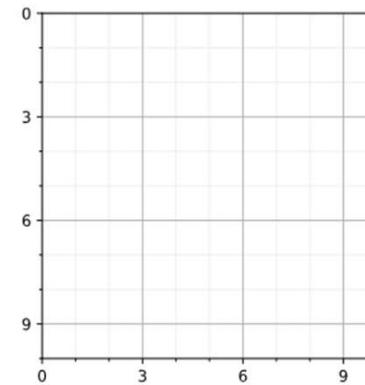
Feature map



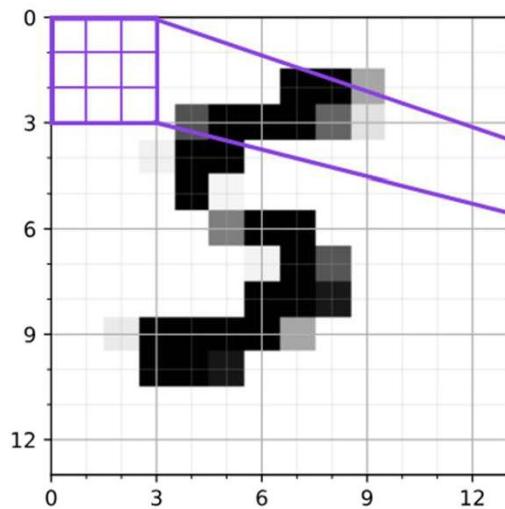
Input (image)

3x3 feature detector (kernel, filter)

A 3x3 grid of squares, each outlined in purple. This represents a 3x3 feature detector or kernel used in convolutional neural networks.

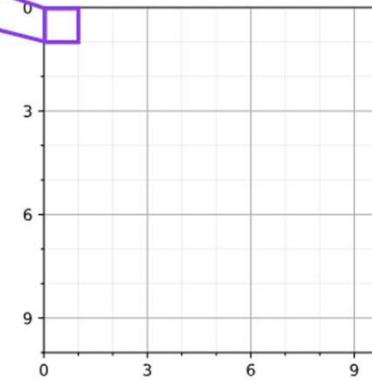


Feature map

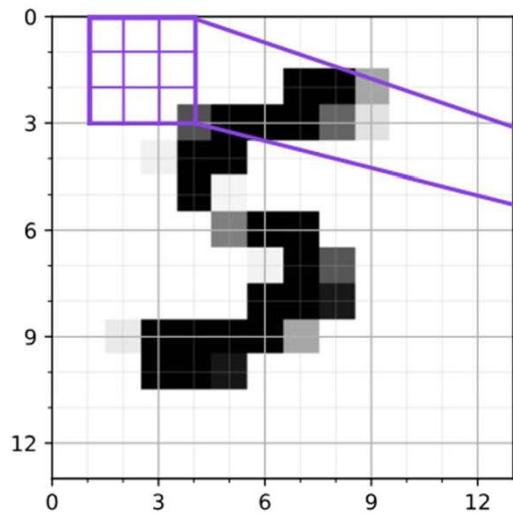


Input (image)

Slide feature detector over inputs

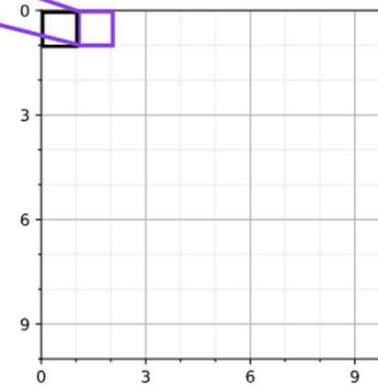


Feature map

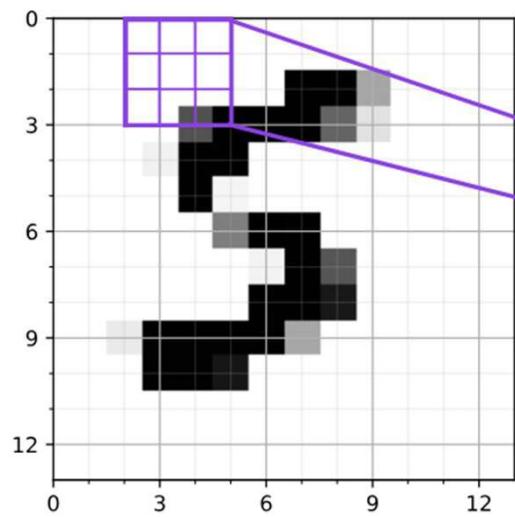


Input (image)

Slide feature detector over inputs

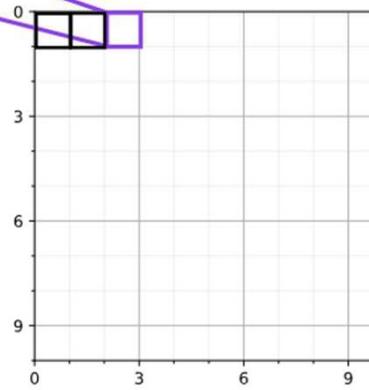


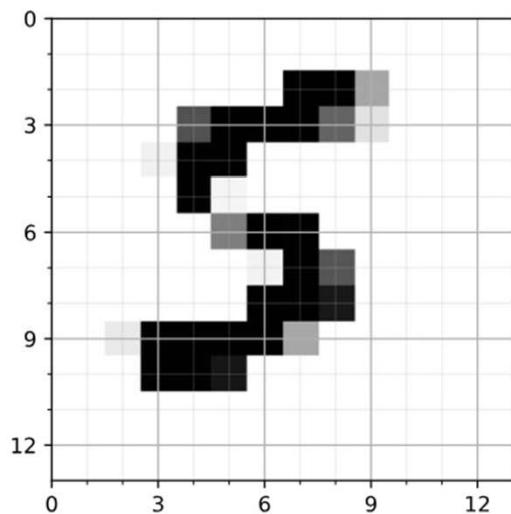
Feature map



Input (image)

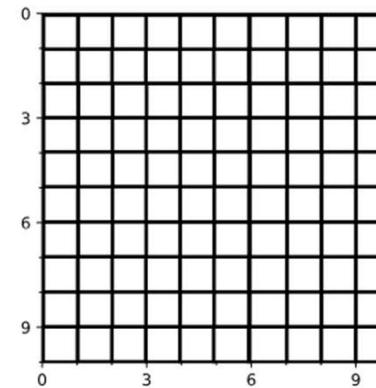
Slide feature detector over inputs



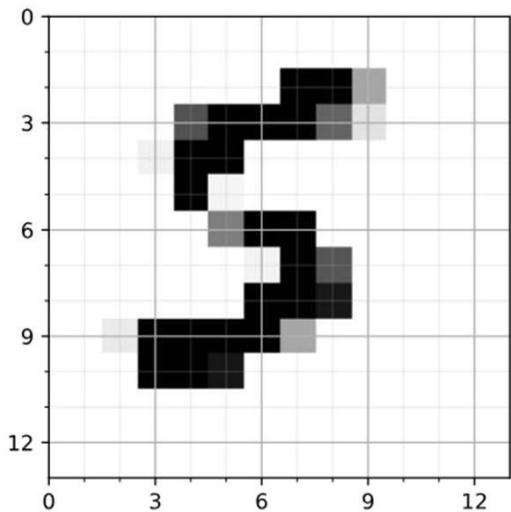


Input (image)

Slide feature detector over inputs

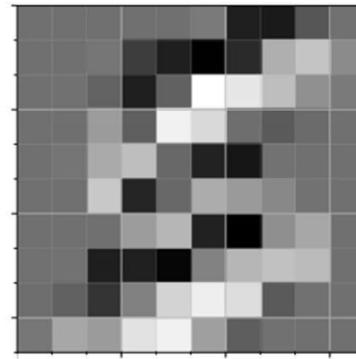


Feature map



Input (image)

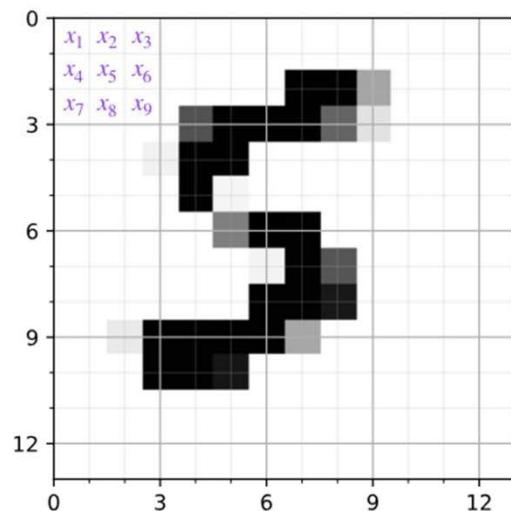
Slide feature detector over inputs



Feature map

What is happening during this operation?

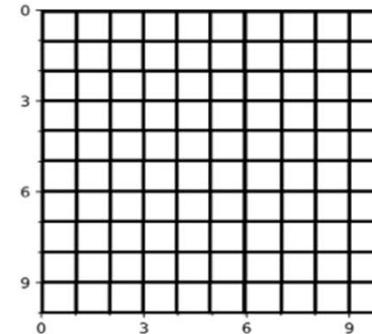
$$\begin{array}{|c|c|c|} \hline w_1 & w_2 & w_3 \\ \hline w_4 & w_5 & w_6 \\ \hline w_7 & w_8 & w_9 \\ \hline \end{array}$$



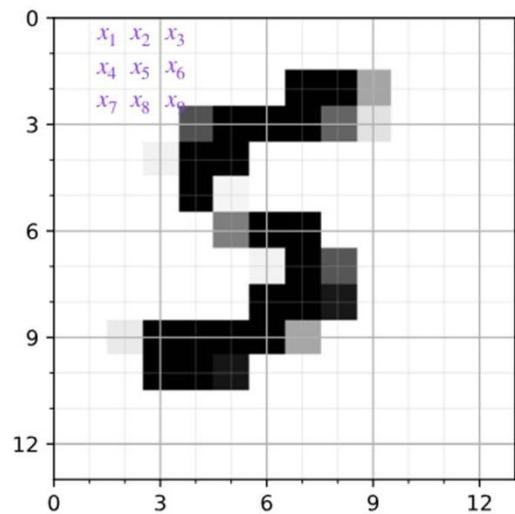
The inputs (x 's) differ as we slide over the image.

The weights (w 's) do not differ. → weight sharing

$$z = b + \sum_j w_j x_j$$



w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



weight sharing

- Rationale: A feature detector that works well in one region may also work well in another region
- A reduction in parameters to fit (compared to MLP)

Convolution

- This is how we calculate the convolutional layer's output:

$$\text{ConvolvedFeature}(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i+a, j+b)K(a, b)$$

I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel

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

4		

Convolved
Feature

- This is how we calculate the convolutional layer's output:

$$ConvolvedFeature(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i + a, j + b)K(a, b)$$

I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved
Feature

- This is how we calculate the convolutional layer's output:

$$ConvolvedFeature(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i + a, j + b)K(a, b)$$

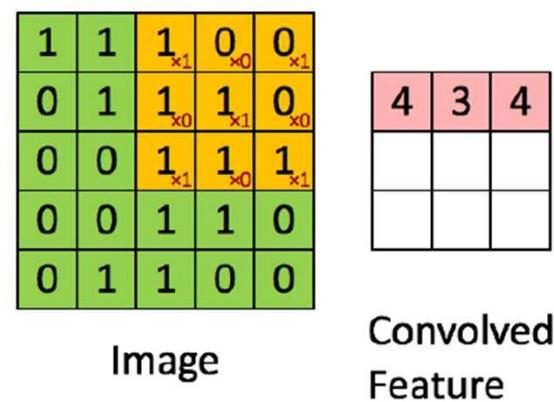
I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel



- This is how we calculate the convolutional layer's output:

$$ConvolvedFeature(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i + a, j + b)K(a, b)$$

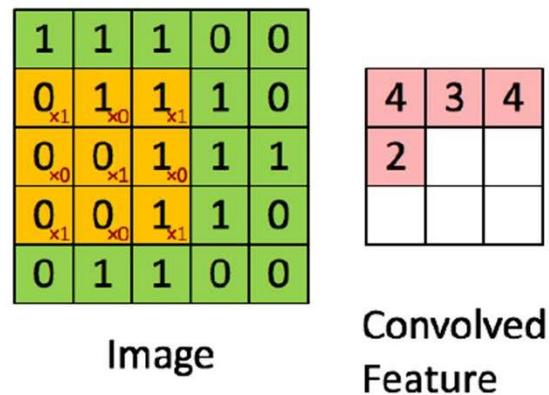
I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel



- This is how we calculate the convolutional layer's output:

$$ConvolvedFeature(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i + a, j + b)K(a, b)$$

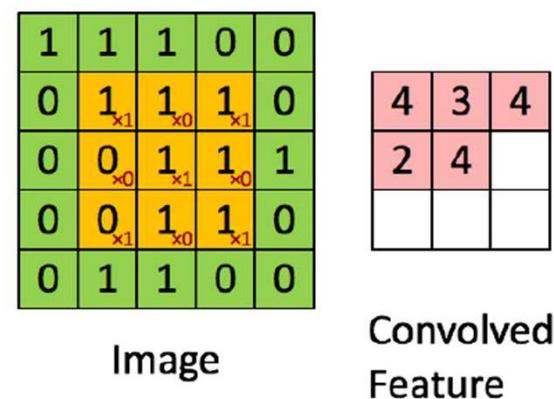
I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel



- This is how we calculate the convolutional layer's output:

$$ConvolvedFeature(i, j) = (I * K)(i, j) = \sum_{a=0}^{k_h-1} \sum_{b=0}^{k_w-1} I(i + a, j + b)K(a, b)$$

I: Input Image

K: Our Kernel

k_h and k_w : The height and width of the Kernel

1	0	1
0	1	0
1	0	1

Figure:
Convolv-
ing Kernel

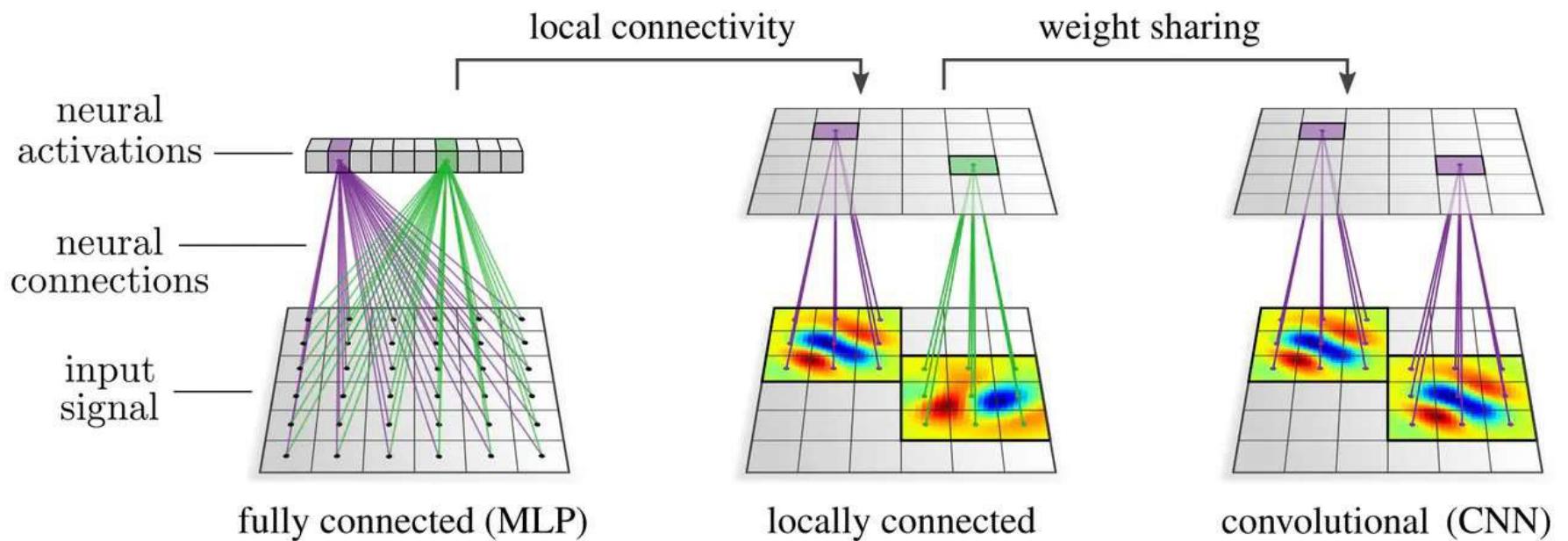
1	1	1	0	0
0	1	1	1	0
0	0	1	<small>\times_1</small>	<small>\times_0</small>
0	0	1	<small>\times_0</small>	<small>\times_1</small>
0	1	1	<small>\times_1</small>	<small>\times_0</small>

Image

4	3	4
2	4	3
2	3	4

Convolved
Feature

Locality and Weight sharing



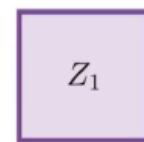
Backpropagation in CNNs

- Same overall concept as before: Multivariable chain rule, but now with an additional weight sharing constraint

Convolution

a_1	a_2	a_3	a_4	a_5
a_6	a_7	a_8	a_9	a_{10}
a_{11}	a_{12}	a_{13}	a_{14}	a_{15}
a_{16}	a_{17}	a_{18}	a_{19}	a_{20}
a_{21}	a_{22}	a_{23}	a_{24}	a_{25}

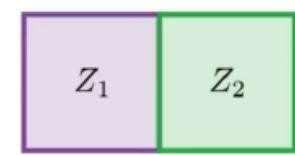
w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



Convolution

a_1	a_2	a_3	a_4	a_5
a_6	a_7	a_8	a_9	a_{10}
a_{11}	a_{12}	a_{13}	a_{14}	a_{15}
a_{16}	a_{17}	a_{18}	a_{19}	a_{20}
a_{21}	a_{22}	a_{23}	a_{24}	a_{25}

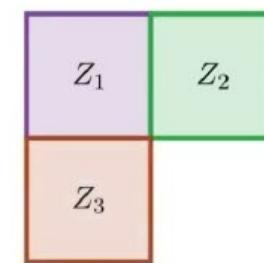
w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



Convolution

a_1	a_2	a_3	a_4	a_5
a_6	a_7	a_8	a_9	a_{10}
a_{11}	a_{12}	a_{13}	a_{14}	a_{15}
a_{16}	a_{17}	a_{18}	a_{19}	a_{20}
a_{21}	a_{22}	a_{23}	a_{24}	a_{25}

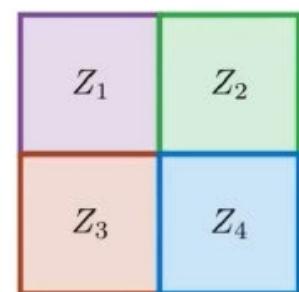
w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



Convolution

a_1	a_2	a_3	a_4	a_5
a_6	a_7	a_8	a_9	a_{10}
a_{11}	a_{12}	a_{13}	a_{14}	a_{15}
a_{16}	a_{17}	a_{18}	a_{19}	a_{20}
a_{21}	a_{22}	a_{23}	a_{24}	a_{25}

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



$$z_1 = w_1 \times a_1 + w_2 \times a_2 + w_3 \times a_3 + w_4 \times a_4 + \dots + w_9 \times a_{13}$$

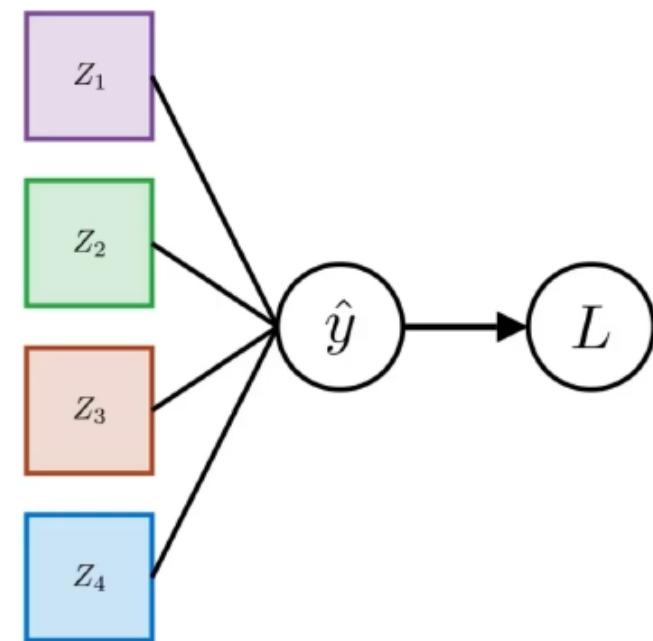
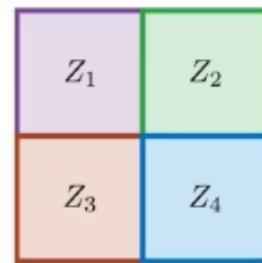
$$z_2 = w_1 \times a_3 + w_2 \times a_4 + w_3 \times a_5 + w_4 \times a_8 + \dots + w_9 \times a_{15}$$

$$z_3 = w_1 \times a_{11} + w_2 \times a_{12} + w_3 \times a_{13} + w_4 \times a_{16} + \dots + w_9 \times a_{23}$$

$$z_4 = w_1 \times a_{13} + w_2 \times a_{14} + w_3 \times a_{15} + w_4 \times a_{18} + \dots + w_9 \times a_{25}$$

- For easier computation, we will assume one layer of convolution and a perceptron as our whole network.
- **Recall:** $w_i^* = w_i - \alpha \times \frac{\partial L}{\partial w_i}$

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9



- We can easily calculate the gradients:

$$\begin{aligned}
 \frac{\partial L}{\partial w_1} &= \frac{\partial z_1}{\partial w_1} \frac{\partial L}{\partial z_1} + \frac{\partial z_2}{\partial w_1} \frac{\partial L}{\partial z_2} + \frac{\partial z_3}{\partial w_1} \frac{\partial L}{\partial z_3} + \frac{\partial z_4}{\partial w_1} \frac{\partial L}{\partial z_4} \\
 &= a_1 \frac{\partial L}{\partial z_1} + a_3 \frac{\partial L}{\partial z_2} + a_{11} \frac{\partial L}{\partial z_3} + a_{13} \frac{\partial L}{\partial z_4}
 \end{aligned}$$

a_1	a_2	a_3	a_4	a_5
a_6	a_7	a_8	a_9	a_{10}
a_{11}	a_{12}	a_{13}	a_{14}	a_{15}
a_{16}	a_{17}	a_{18}	a_{19}	a_{20}
a_{21}	a_{22}	a_{23}	a_{24}	a_{25}

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

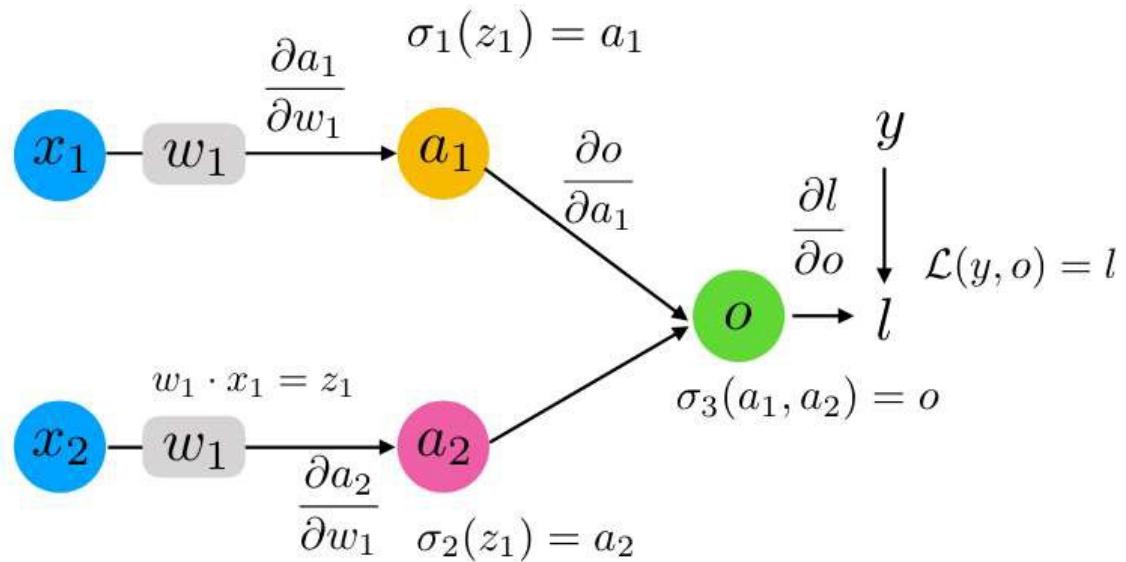
$\frac{\partial L}{\partial z_1}$	$\frac{\partial L}{\partial z_2}$
$\frac{\partial L}{\partial z_3}$	$\frac{\partial L}{\partial z_4}$

$$\begin{aligned}
 \frac{\partial L}{\partial w_1} &= a_1 \frac{\partial L}{\partial z_1} + a_3 \frac{\partial L}{\partial z_2} + a_{11} \frac{\partial L}{\partial z_3} + a_{13} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_2} &= a_2 \frac{\partial L}{\partial z_1} + a_4 \frac{\partial L}{\partial z_2} + a_{12} \frac{\partial L}{\partial z_3} + a_{14} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_3} &= a_3 \frac{\partial L}{\partial z_1} + a_5 \frac{\partial L}{\partial z_2} + a_{13} \frac{\partial L}{\partial z_3} + a_{15} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_4} &= a_6 \frac{\partial L}{\partial z_1} + a_8 \frac{\partial L}{\partial z_2} + a_{16} \frac{\partial L}{\partial z_3} + a_{18} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_5} &= a_7 \frac{\partial L}{\partial z_1} + a_9 \frac{\partial L}{\partial z_2} + a_{17} \frac{\partial L}{\partial z_3} + a_{19} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_6} &= a_8 \frac{\partial L}{\partial z_1} + a_{10} \frac{\partial L}{\partial z_2} + a_{18} \frac{\partial L}{\partial z_3} + a_{20} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_7} &= a_{11} \frac{\partial L}{\partial z_1} + a_{13} \frac{\partial L}{\partial z_2} + a_{21} \frac{\partial L}{\partial z_3} + a_{23} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_8} &= a_{12} \frac{\partial L}{\partial z_1} + a_{14} \frac{\partial L}{\partial z_2} + a_{22} \frac{\partial L}{\partial z_3} + a_{24} \frac{\partial L}{\partial z_4} \\
 \frac{\partial L}{\partial w_9} &= a_{13} \frac{\partial L}{\partial z_1} + a_{15} \frac{\partial L}{\partial z_2} + a_{23} \frac{\partial L}{\partial z_3} + a_{25} \frac{\partial L}{\partial z_4}
 \end{aligned}$$

$$\begin{array}{ccc}
 \begin{array}{|c|c|c|} \hline
 a_1 & a_2 & a_3 \\ \hline
 a_6 & a_7 & a_8 \\ \hline
 a_{11} & a_{12} & a_{13} \\ \hline
 \end{array} & \times \frac{\partial L}{\partial z_1} & + \\
 \begin{array}{|c|c|c|} \hline
 a_3 & a_4 & a_5 \\ \hline
 a_8 & a_9 & a_{10} \\ \hline
 a_{13} & a_{14} & a_{15} \\ \hline
 \end{array} & \times \frac{\partial L}{\partial z_2} & + \\
 \begin{array}{|c|c|c|} \hline
 a_{11} & a_{12} & a_{13} \\ \hline
 a_{16} & a_{17} & a_{18} \\ \hline
 a_{21} & a_{22} & a_{23} \\ \hline
 \end{array} & \times \frac{\partial L}{\partial z_3} & + \\
 \begin{array}{|c|c|c|} \hline
 a_{13} & a_{14} & a_{15} \\ \hline
 a_{18} & a_{19} & a_{20} \\ \hline
 a_{23} & a_{24} & a_{25} \\ \hline
 \end{array} & \times \frac{\partial L}{\partial z_4} & =
 \end{array}$$

$\frac{\partial L}{\partial w_1}$	$\frac{\partial L}{\partial w_2}$	$\frac{\partial L}{\partial w_3}$
$\frac{\partial L}{\partial w_4}$	$\frac{\partial L}{\partial w_5}$	$\frac{\partial L}{\partial w_6}$
$\frac{\partial L}{\partial w_7}$	$\frac{\partial L}{\partial w_8}$	$\frac{\partial L}{\partial w_9}$

Graph with Weight Sharing



Upper path

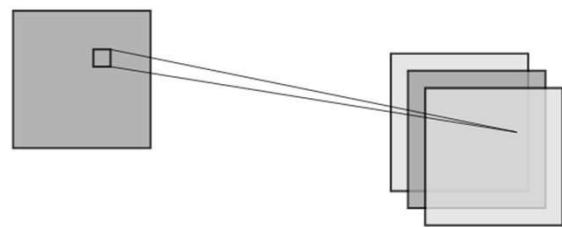
$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} + \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_1} \quad (\text{multivariable chain rule})$$

Lower path

- Which problem is solved with weigh sharing and locality?

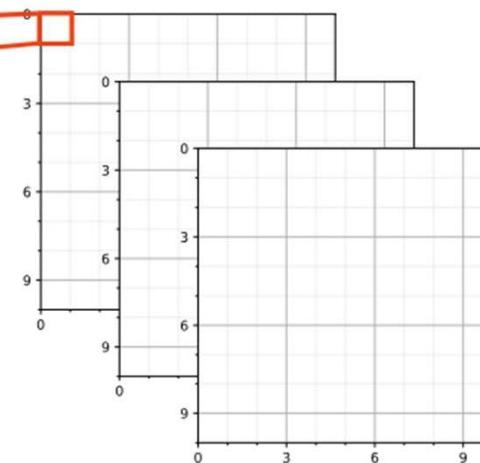
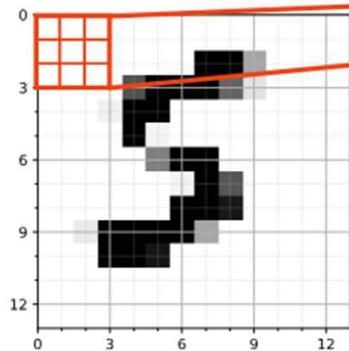
Multiple Output Channels

1 input channel, 3 output channels

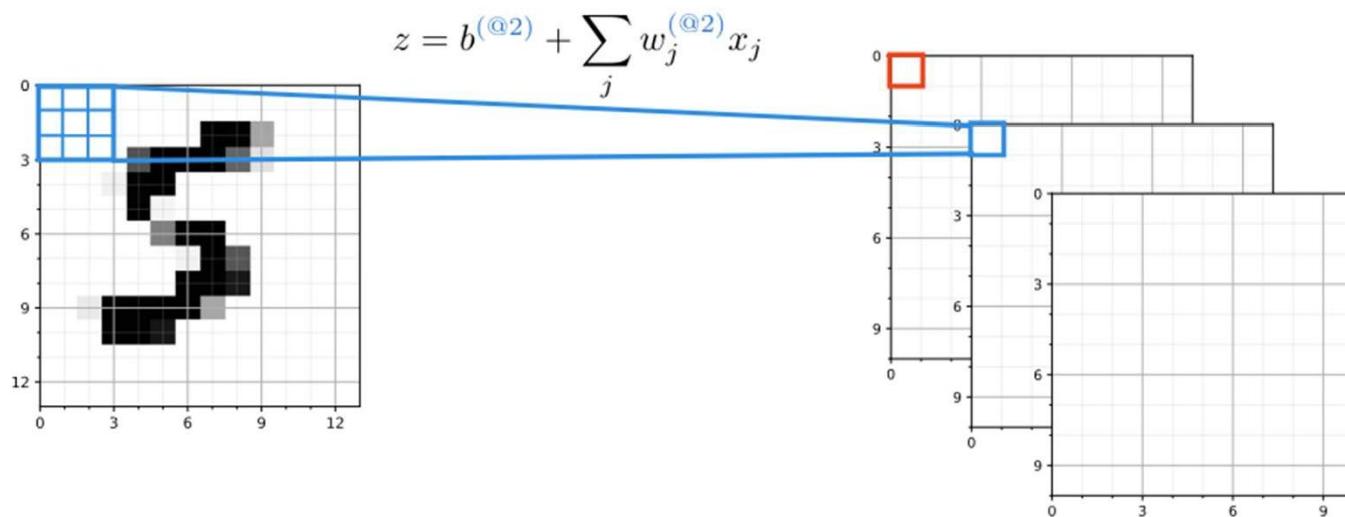


1 input channel 3 output channel

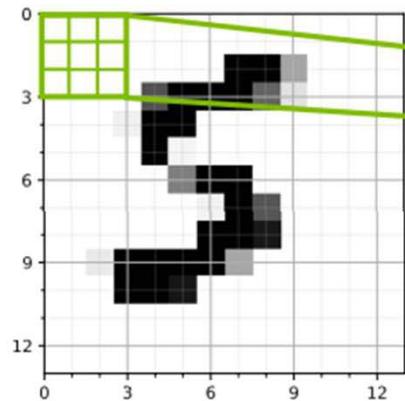
$$z = b^{(\text{@1})} + \sum_j w_j^{(\text{@1})} x_j$$



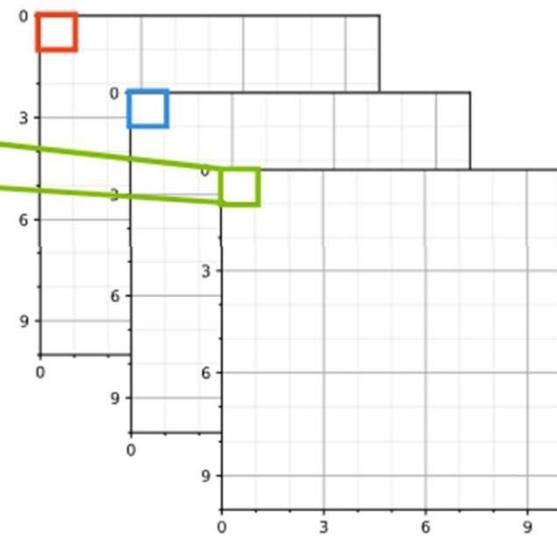
1 input channel 3 output channel



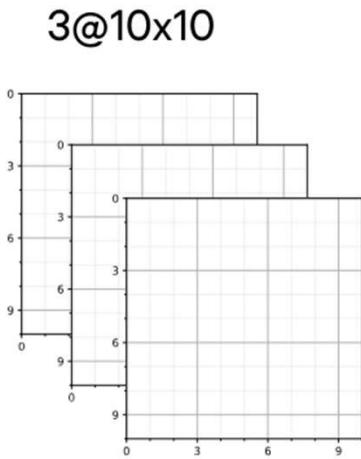
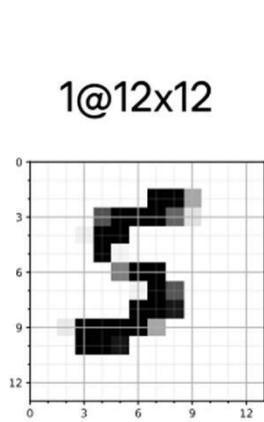
1 input channel 3 output channel



$$z = b^{(\text{@}3)} + \sum_j w_j^{(\text{@}3)} x_j$$



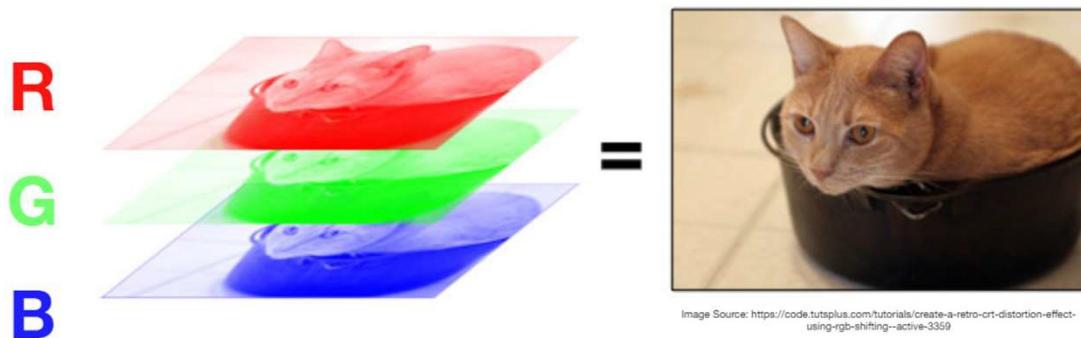
1 input channel 3 output channel



Multiple "feature detectors" (kernels)
are used to create multiple feature maps

Multiple input channel

RGB Image

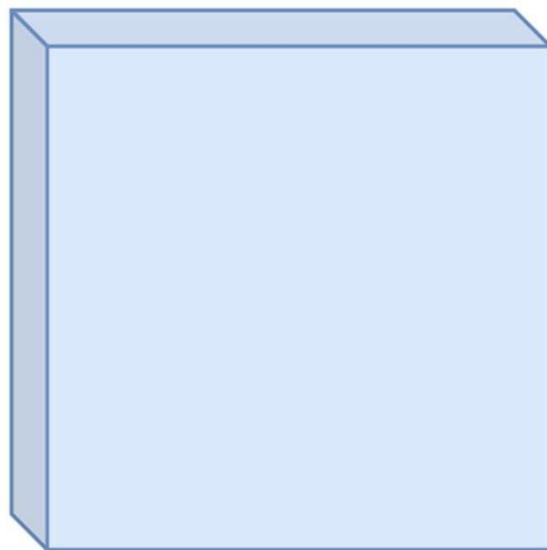


Color image as a stack of matrices

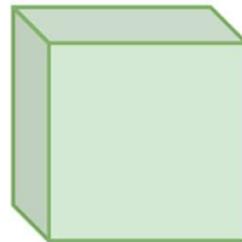
Input channels

Filters always extend the full depth of the input.
(#Input channels == #Filter Channels)

32x32x3 Image

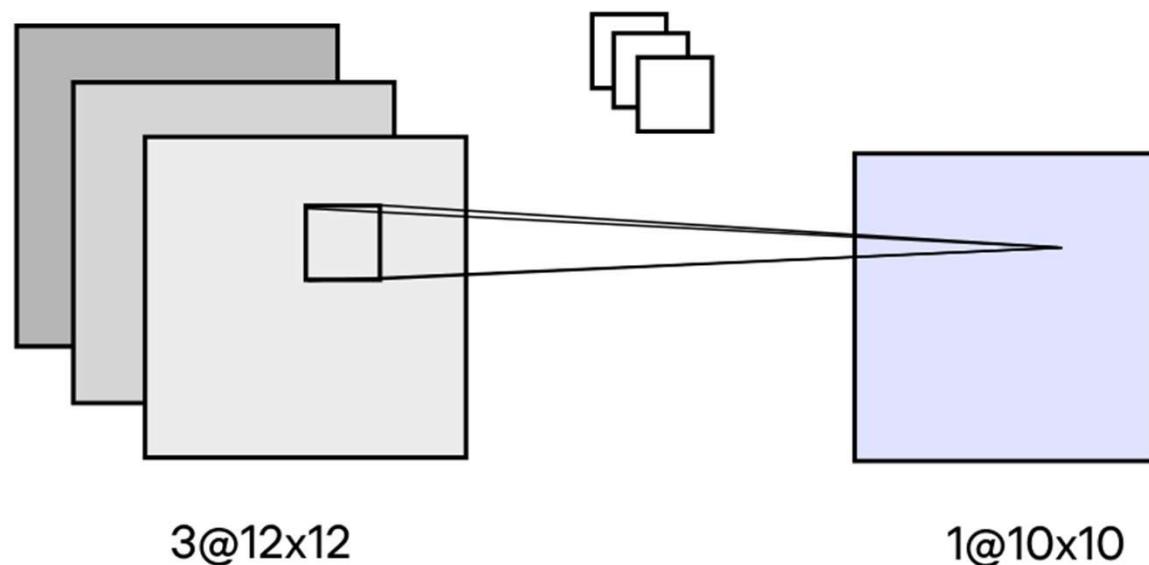


5x5x3 Filter

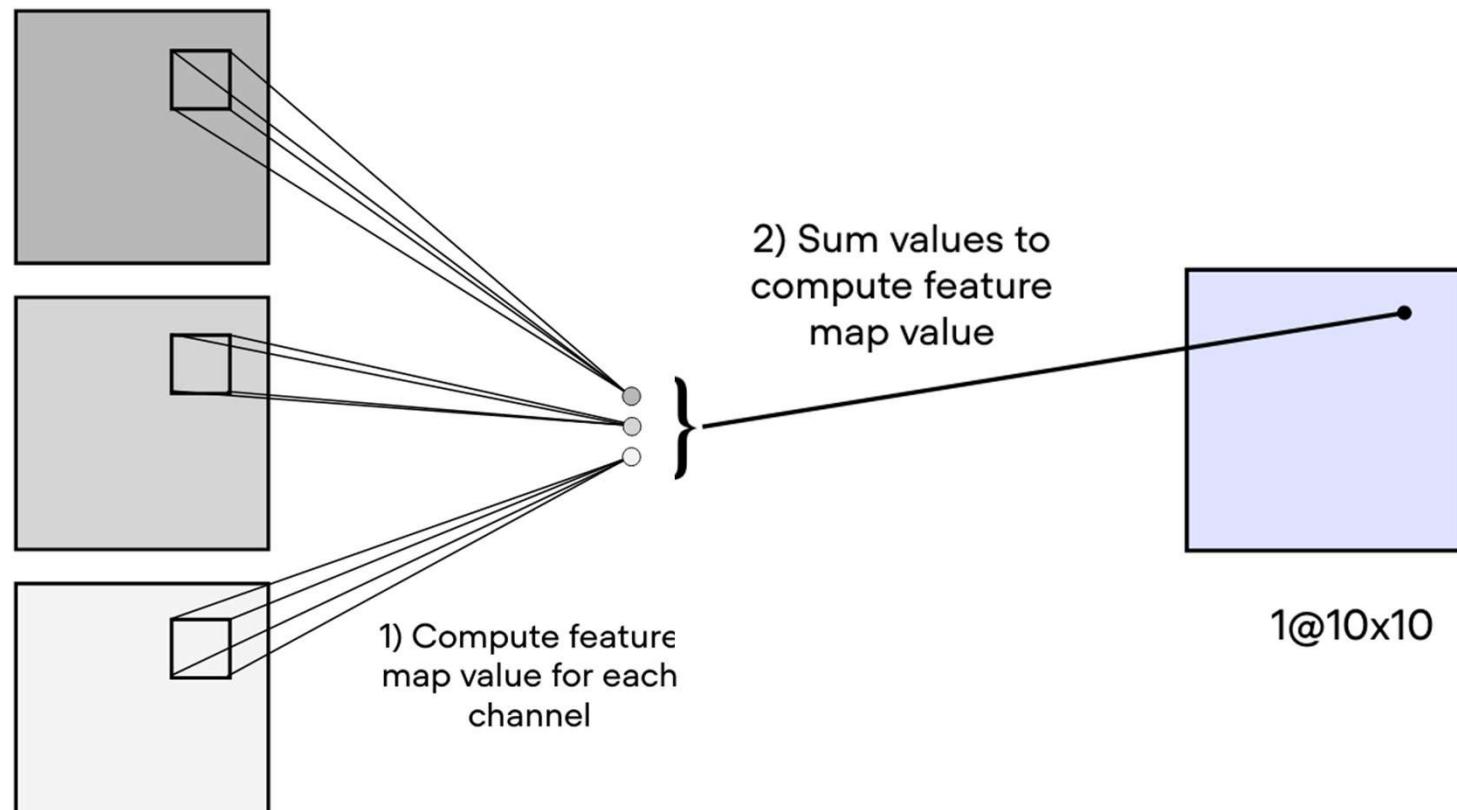


3 input channel, 1 output channel

kernel has 3 channels

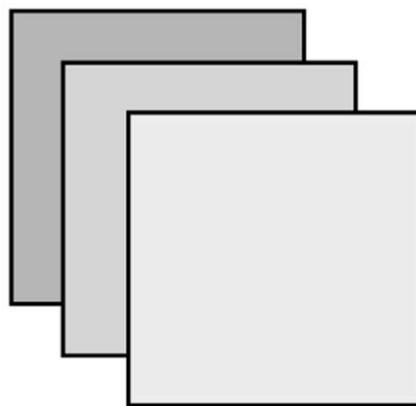


3 input channel, 1 output channel



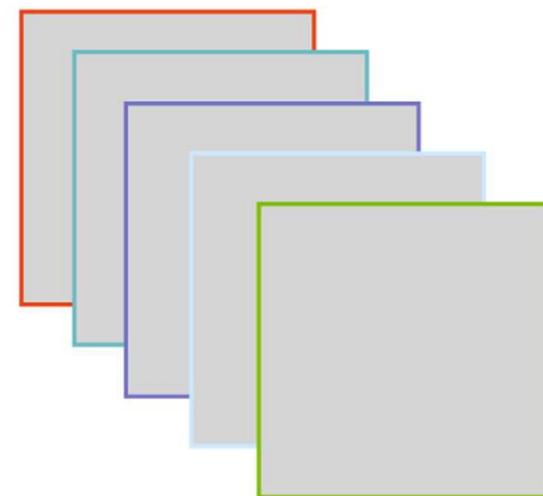
Multiple input and output channels

3 **input** channels



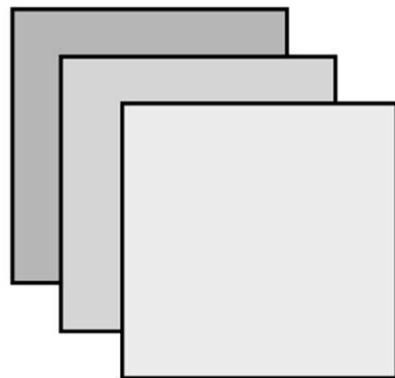
3@64x64

5 **output** channels

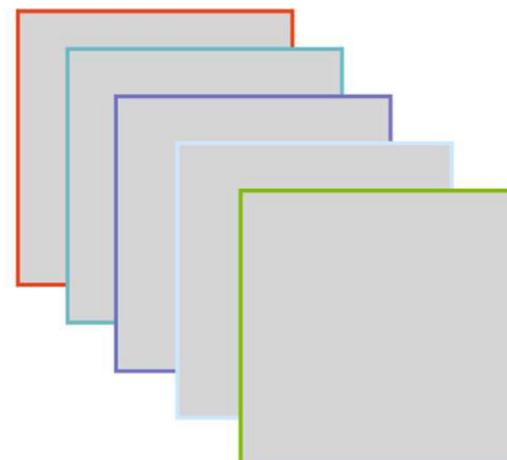


5@64x64

Multiple input and output channels



3@64x64

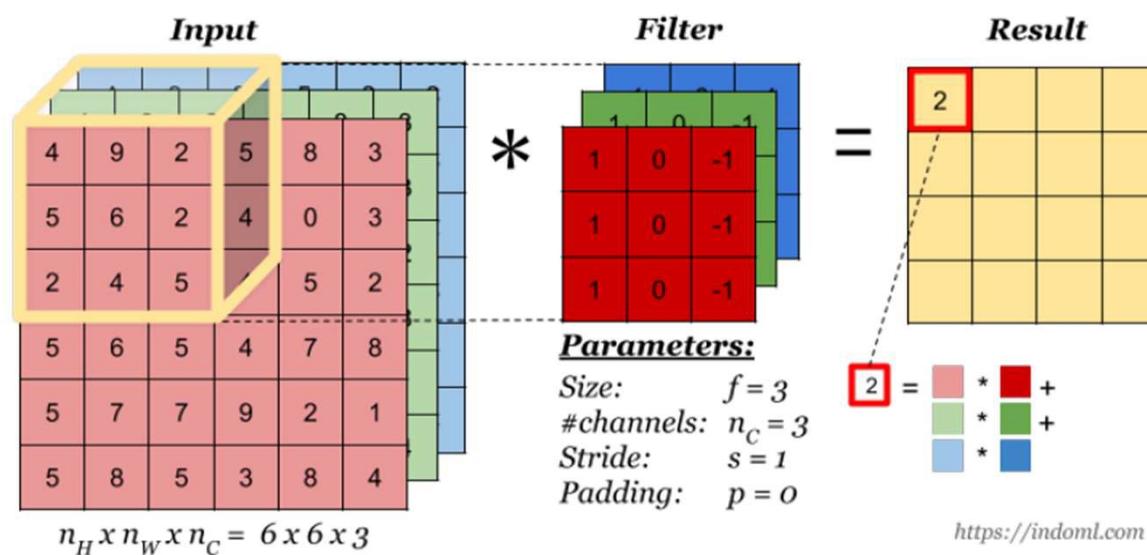


5@64x64

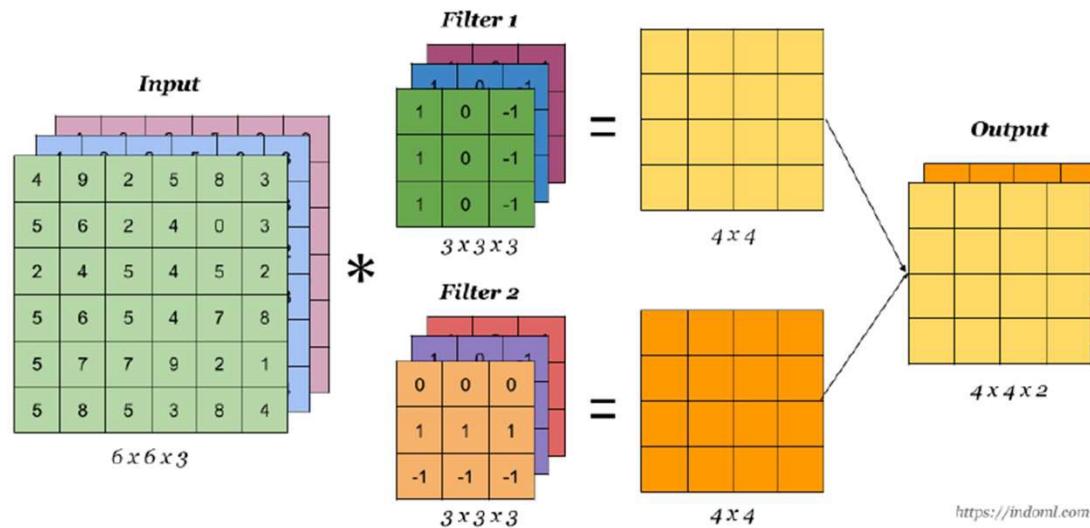
5 kernels with
3 channels each



- As said before: "Filters always extend the full depth of the input.
(#Input channels == #Filter Channels)"

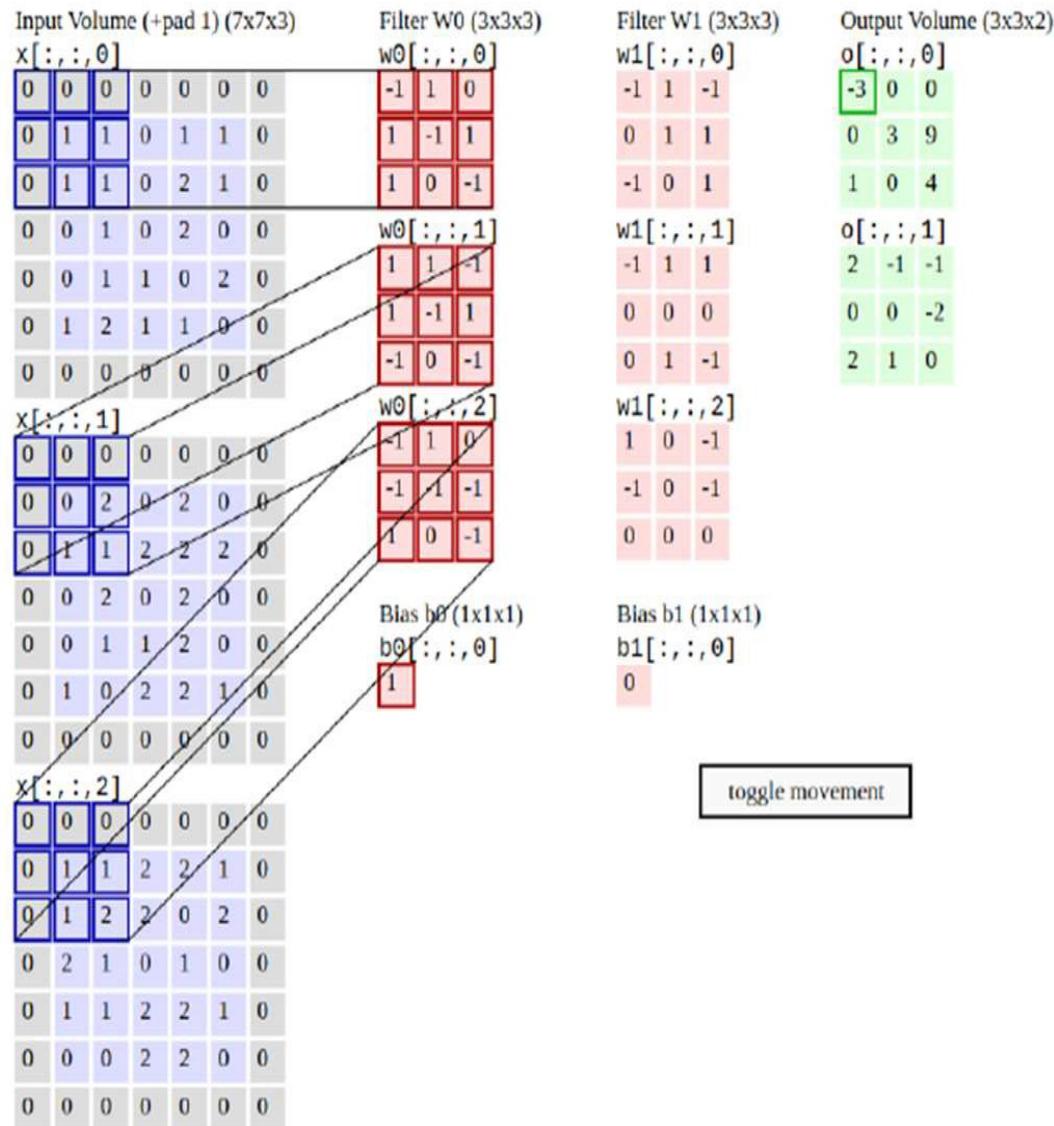


- As said before: "Filters always extend the full depth of the input volume.
(#Input channels == #Filter Channels)"
- Each filter results in one output channel. We can apply multiple different filters to obtain multiple output channels. Each channel can learn something distinct.



<https://indoml.com>

Let's have a closer look to how the calculations with more than one filter work:



Input Volume (+pad 1) (7x7x3)
 $x[:, :, \theta]$

0	0	0	0	0	0	0
0	1	1	0	1	1	0
0	1	1	0	2	1	0
0	0	1	0	2	0	0
0	0	1	1	0	2	0
0	1	2	1	1	0	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)
 $w0[:, :, \theta]$

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

Filter W1 (3x3x3)
 $w1[:, :, \theta]$

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

Output Volume (3x3x2)
 $o[:, :, \theta]$

-3	0	0
0	3	9
1	0	4
2	-1	-1
0	0	-2
2	1	0

$x[:, :, 1]$

0	0	0	0	0	0	0
0	0	2	0	2	0	0
0	1	1	2	2	2	0
0	0	2	0	2	0	0
0	0	1	1	2	0	0
0	1	0	2	2	1	0
0	0	0	0	0	0	0

$w0[:, :, 2]$

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

$w1[:, :, 2]$

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

Bias b0 (1x1x1)
 $b0[:, :, 0]$

1

Bias b1 (1x1x1)
 $b1[:, :, \theta]$

0

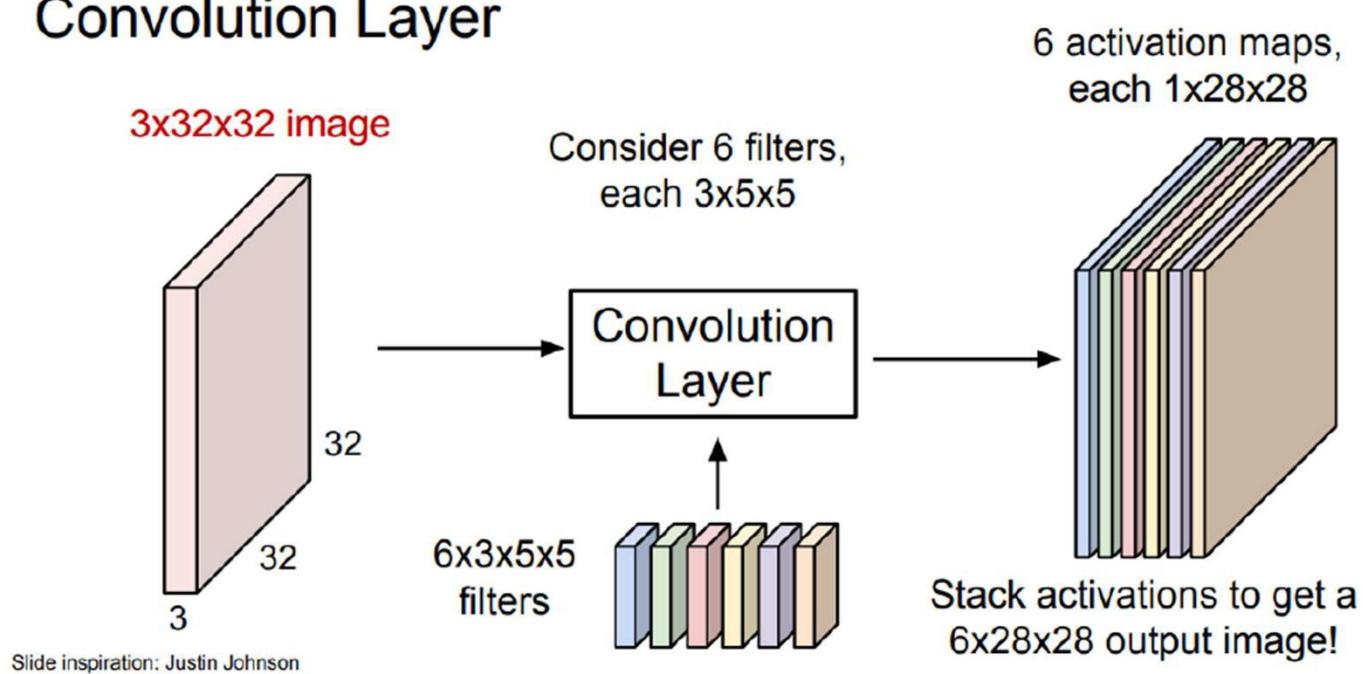
$x[:, :, 2]$

0	0	0	0	0	0	0
0	1	1	2	2	1	0
0	1	2	2	0	2	0
0	2	1	0	1	0	0
0	1	1	2	2	1	0
0	0	0	2	2	0	0
0	0	0	0	0	0	0

toggle movement

- We apply several filters and stack the output together to get a multilayer output, with each layer representing something it learned.

Convolution Layer

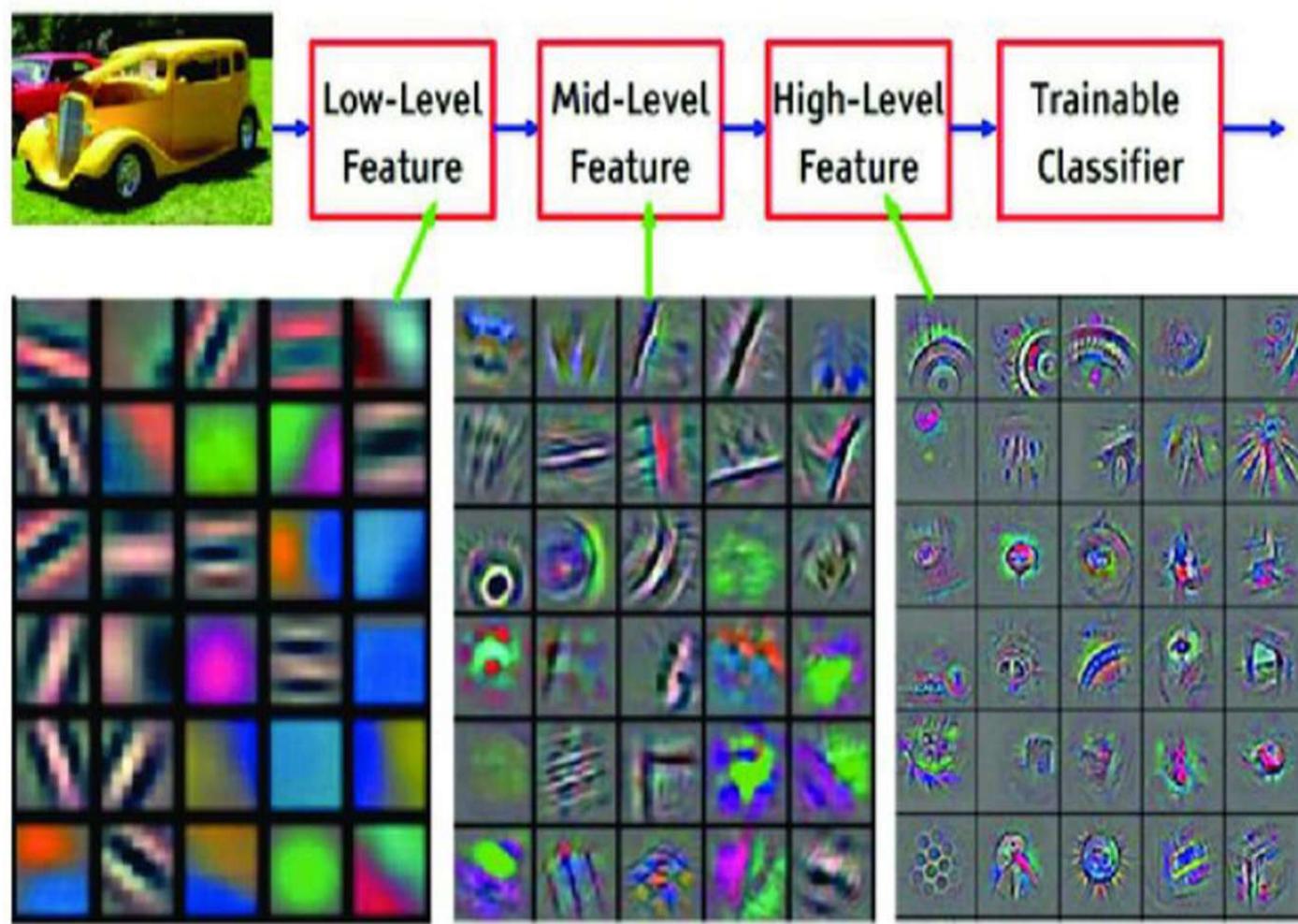


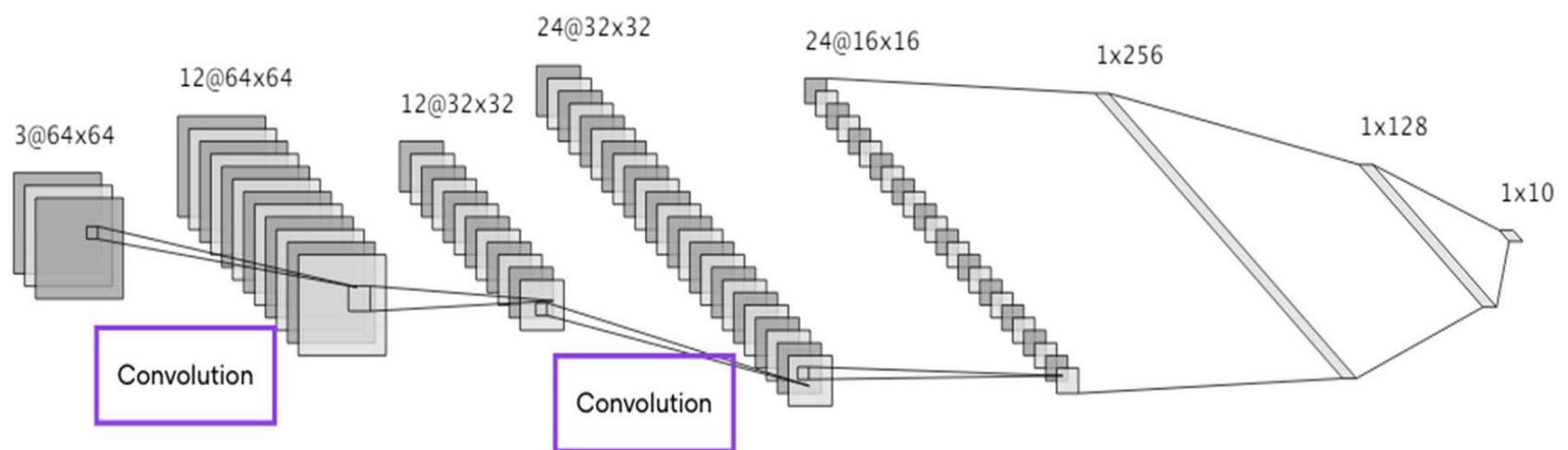
Question: What do convolutional filters at different levels of a ConvNet learn?

Answer:

- Filters in the **early layers** typically detect **simple features** such as edges, textures, and basic shapes.
- As we move **deeper** into the network, the filters learn **more complex and abstract features**, such as specific parts of objects (e.g., a nose, eyes, or other high-level patterns).

- We apply several filters and stack the output together to get a multilayer output, with each layer representing something it learned.
- Some kernels learn to recognize vertical lines, some circles, and some, specific objects:



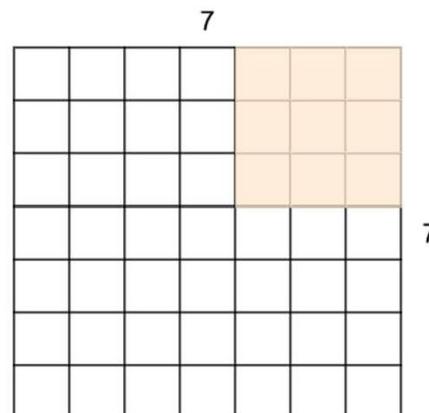


What is a Stride

The amount of movement between applications of the filter to the input image is referred to as the stride, and it is almost always symmetrical in height and width dimensions.

Closer look

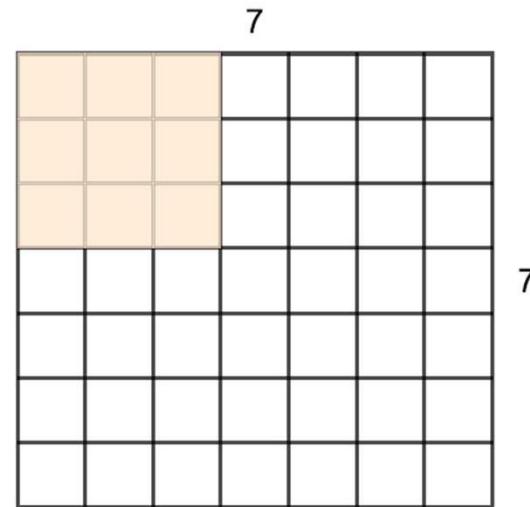
- 7x7 input with 3x3 filter
- This was a Stride 1 filter
- => Outputs 5x5



Strides

Now let's use Stride 2

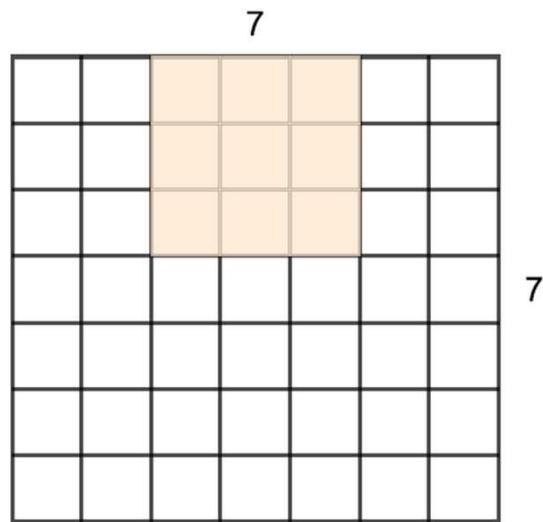
- 7x7 input with 3x3 filter



Strides

Now let's use Stride 2

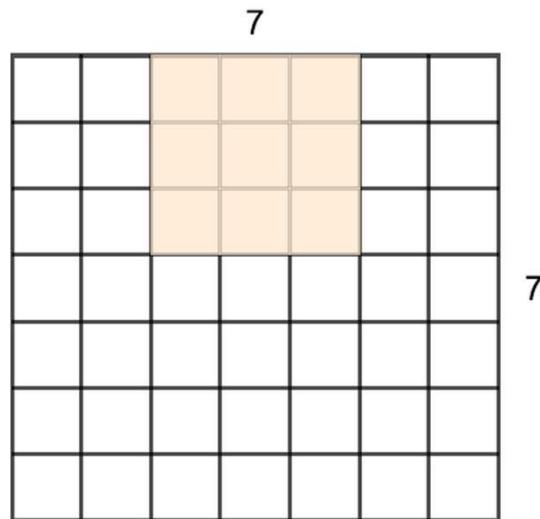
- 7x7 input with 3x3 filter



Strides

Now let's use Stride 2

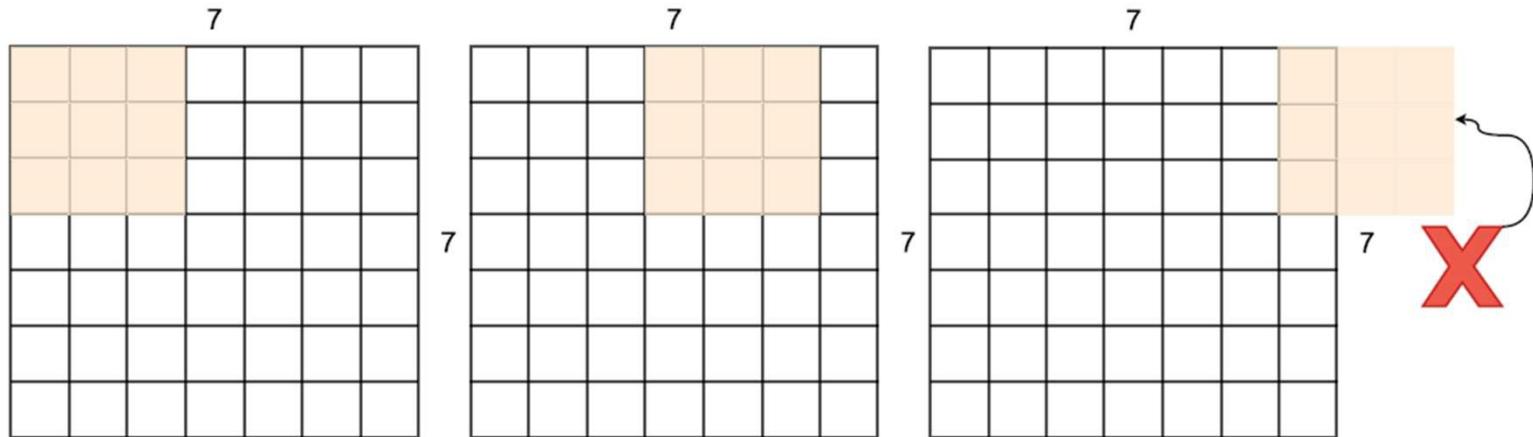
- 7x7 input with 3x3 filter



Strides

Stride 3?

- 7x7 input with 3x3 filter

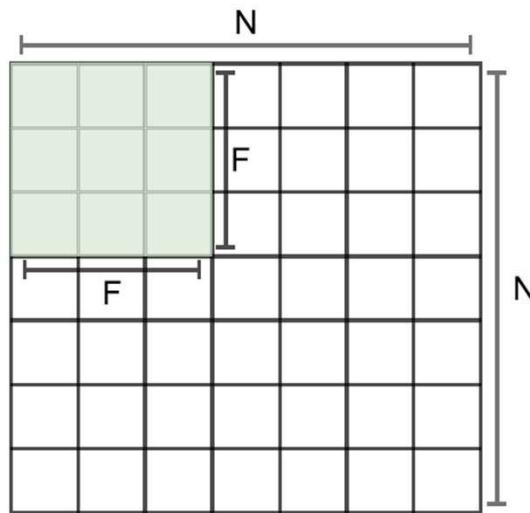


Strides

So 7x7 input with 3x3 filter and stride 3 doesn't work!

Let's do the calculations:

$$\text{OutputSize} = (N - F) / \text{Stride} + 1$$

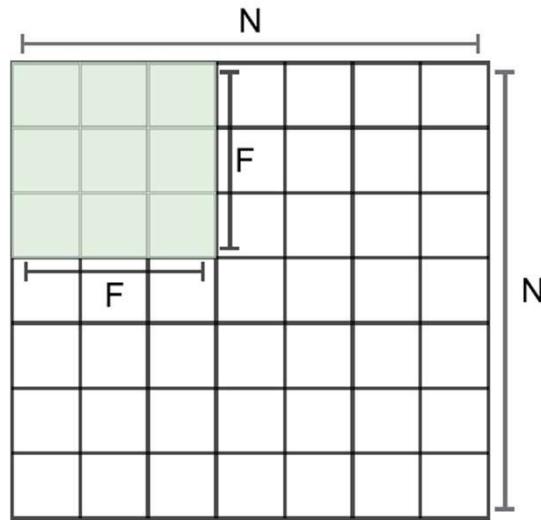


Strides

$$OutputSize = (N - F)/Stride + 1$$

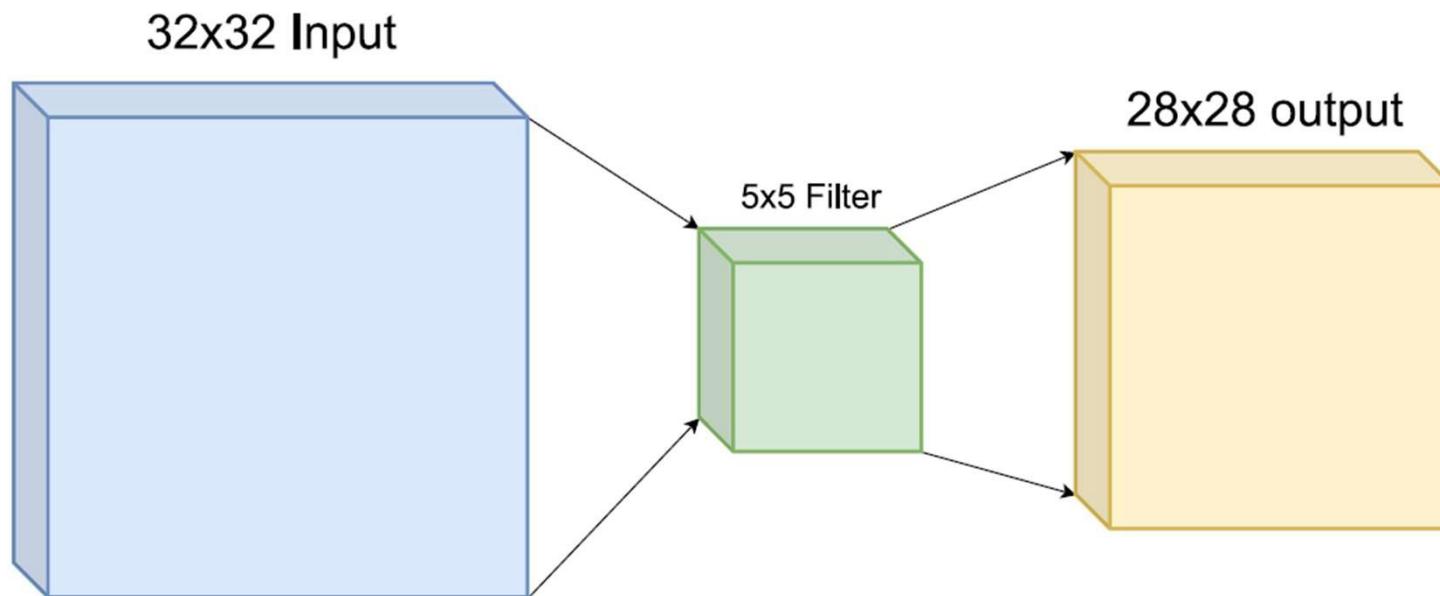
$N = 7, F = 3 \Rightarrow$

- Stride 1 $\Rightarrow (7 - 3)/1 + 1 = 5$
- Stride 2 $\Rightarrow (7 - 3)/2 + 1 = 3$
- Stride 3 $\Rightarrow (7 - 3)/3 + 1 = 2.33 :))$



Padding

1. Borders don't get enough attention.
2. Outputs shrink!



Padding

- We can use Padding to preserve the dimensionality
- It ensures that all pixels are used equally frequently

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 3 & 3 & 4 & 4 & 7 & 0 & 0 \\ \hline 0 & 9 & 7 & 6 & 5 & 8 & 2 & 0 \\ \hline 0 & 6 & 5 & 5 & 6 & 9 & 2 & 0 \\ \hline 0 & 7 & 1 & 3 & 2 & 7 & 8 & 0 \\ \hline 0 & 0 & 3 & 7 & 1 & 8 & 3 & 0 \\ \hline 0 & 4 & 0 & 4 & 3 & 2 & 2 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|c|c|} \hline -10 & -13 & 1 & & & & \\ \hline -9 & 3 & 0 & & & & \\ \hline & & & & & & \\ \hline \end{array} \quad 6 \times 6$$

$6 \times 6 \rightarrow 8 \times 8$

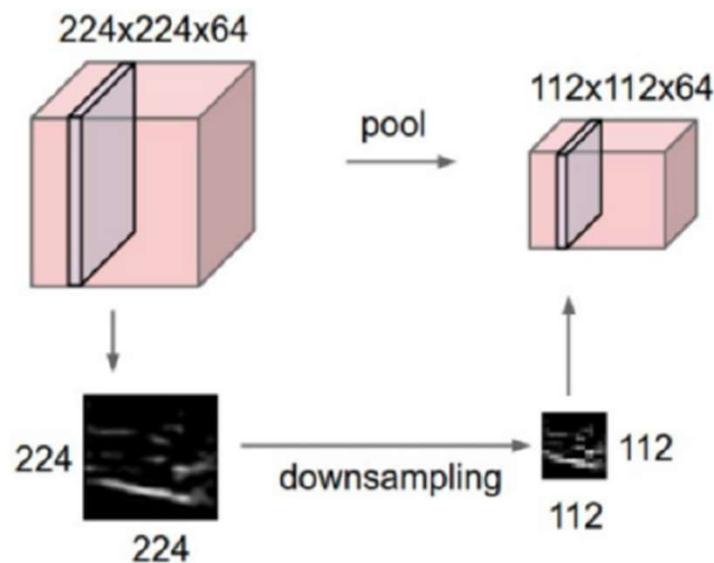
Padding

$$OutputSize = (N + 2P - F)/Stride + 1$$

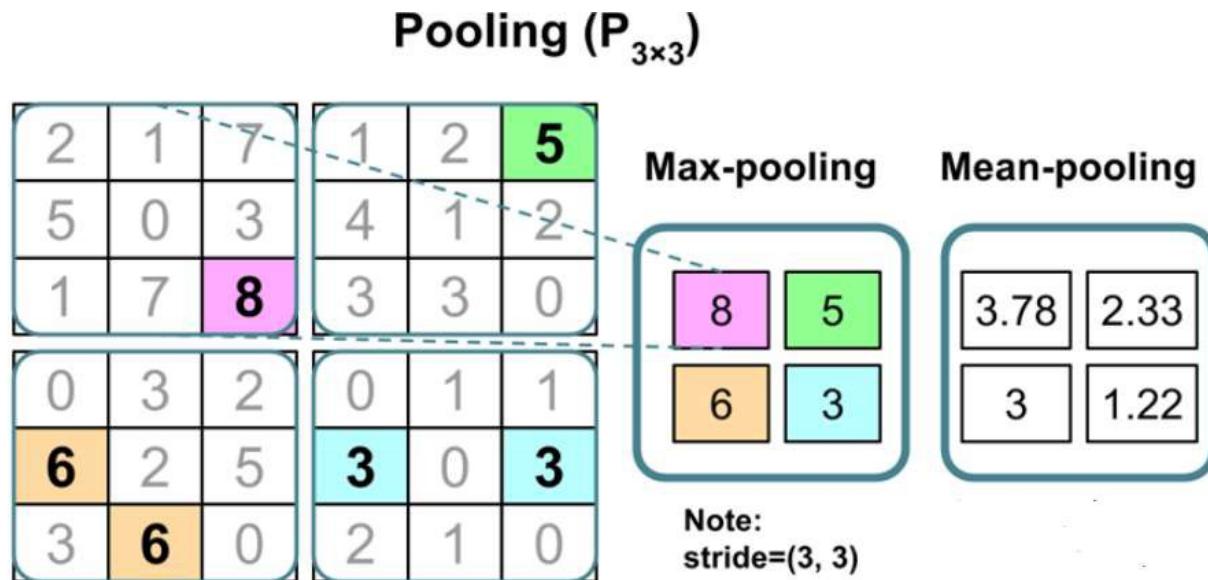
It is common to use $P = (F - 1)/2$ with stride 1 to preserve the input size.

Pooling

Pooling is performed in neural networks to reduce variance and computation complexity



Pooling Layers Can Help With Local Invariance



Note that typical pooling layers do not have any learnable parameters

Downside: Information is lost.

May not matter for classification, but applications where relative position is important (like face recognition)

In practice for CNNs: some image preprocessing still recommended

Pooling benefits

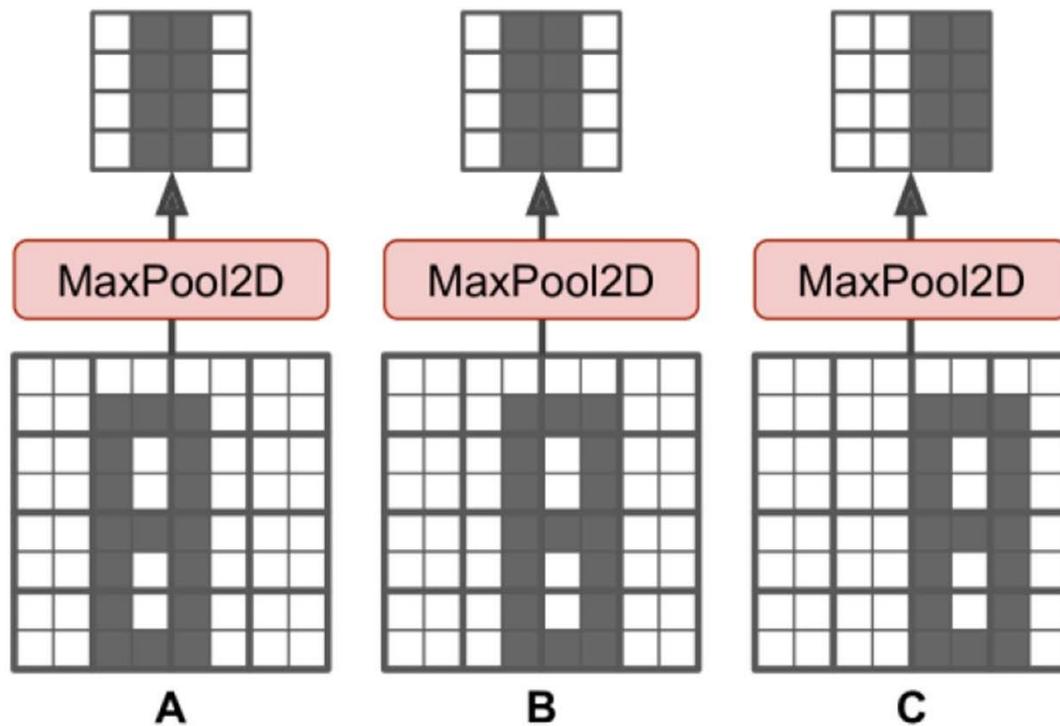
A Pooling layer's goal is to sub-sample the input in order to reduce:

- ▶ The computational load
- ▶ The memory usage
- ▶ The number of parameters
- ▶ Limiting the risk of overfitting

Max Pooling

Most common type of pooling layer

Invariance to small translations



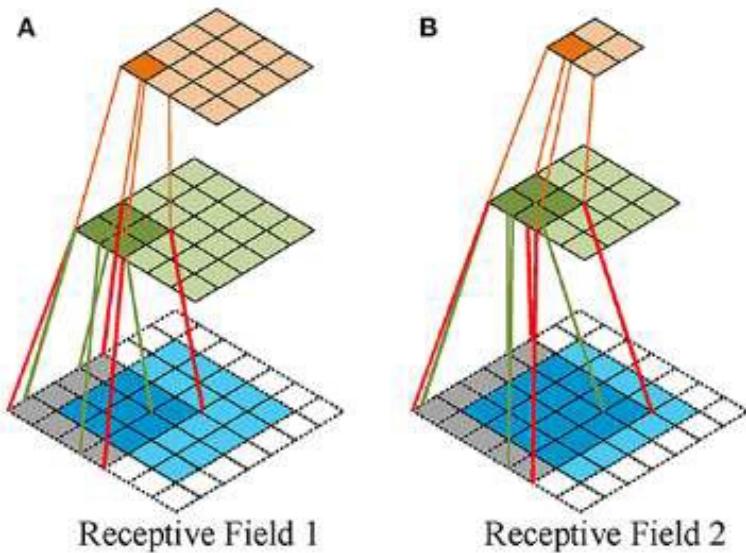
Smaller filter size or larger ?

- Using smaller filter size is better or larger?

Distributing the scan

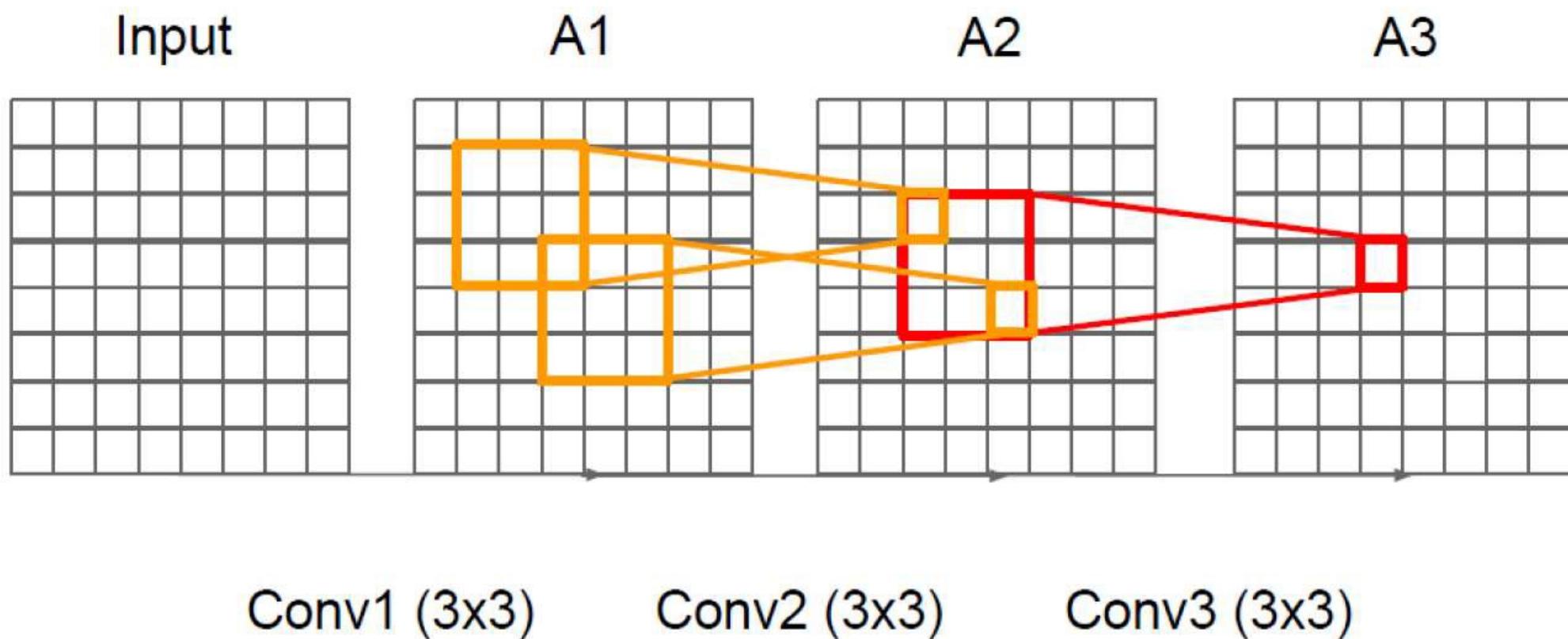
- We can distribute the pattern matching over two layers and still achieve the same block analysis at the second layer
 - The first layer evaluates smaller blocks of pixels
 - The next layer evaluates blocks of outputs from the first layer

Receptive field

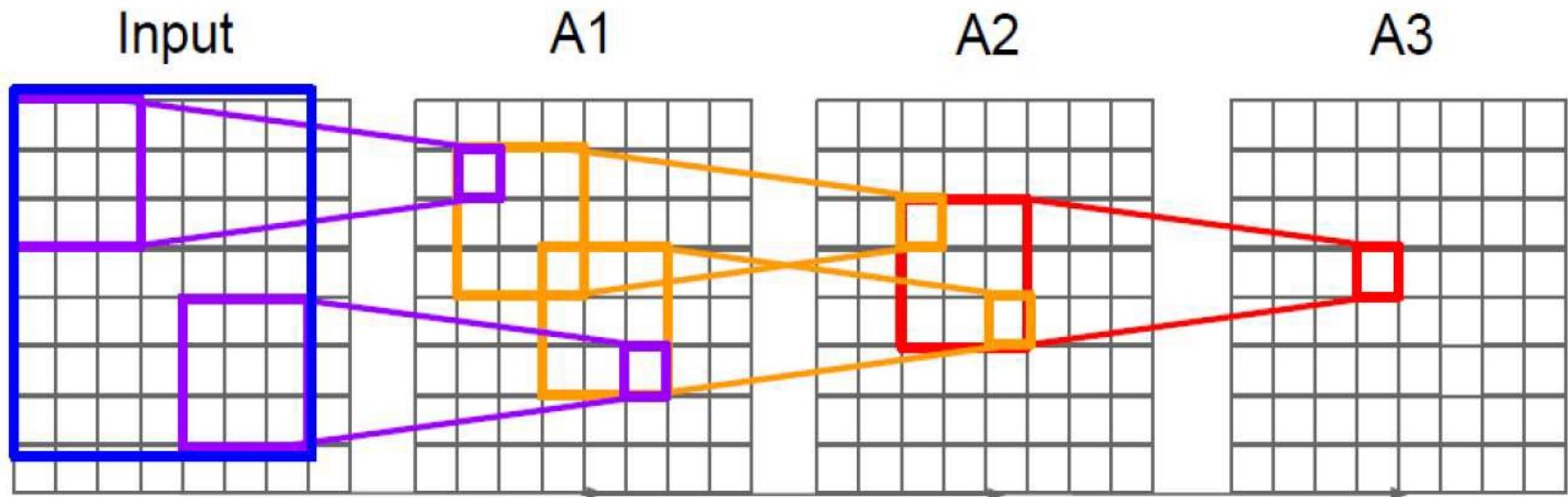


Receptive field a deep CNN. (A) The first convolving: a 3×3 kernel over a 5×5 input padded with a 1×1 border of zeros using unit strides. The second convolving: a 2×2 kernel using unit strides. (B) The first convolving: a 4×4 kernel over a 5×5 input padded with a 1×1 border of zeros using unit strides. The second convolving: a 2×2 kernel using 2×2 strides.

What is the effective receptive field of three 3×3 conv (stride 1) layers?



What is the effective receptive field of three 3×3 conv (stride 1) layers?



Power of small filter

- Suppose input is $H \times W \times C$ and we use convolutions of C filters

Stride 1 and Padding to preserve H , W

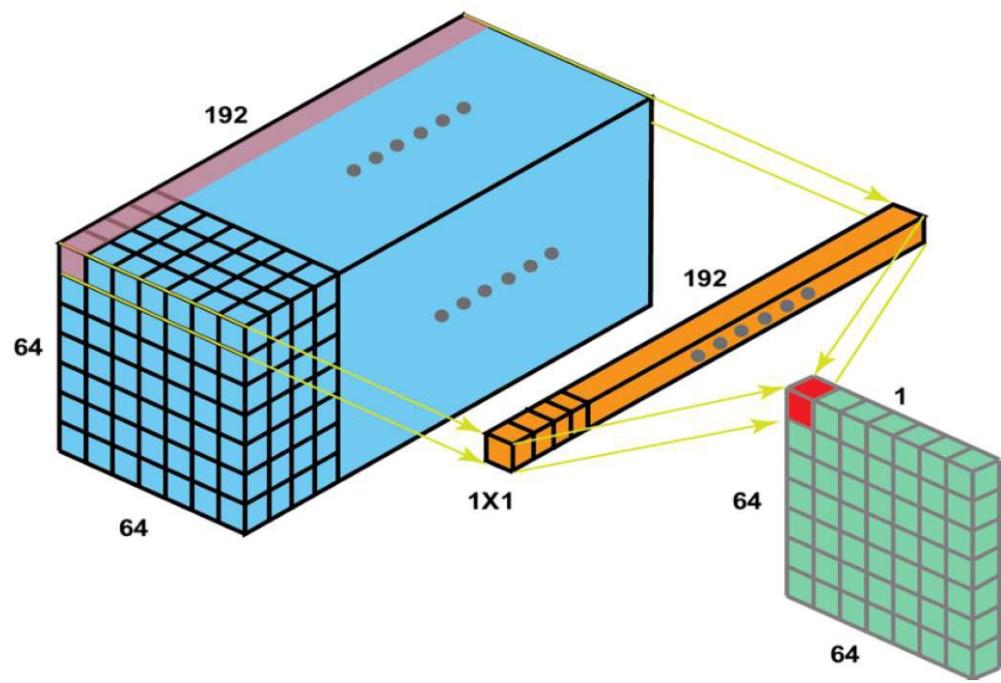
- One convolution with 7×7 filters

- Number of weights : $(7 \times 7 \times C) \times C$ (for the number of filter) = $49C^2$

- 3 convolution with 3×3 filters

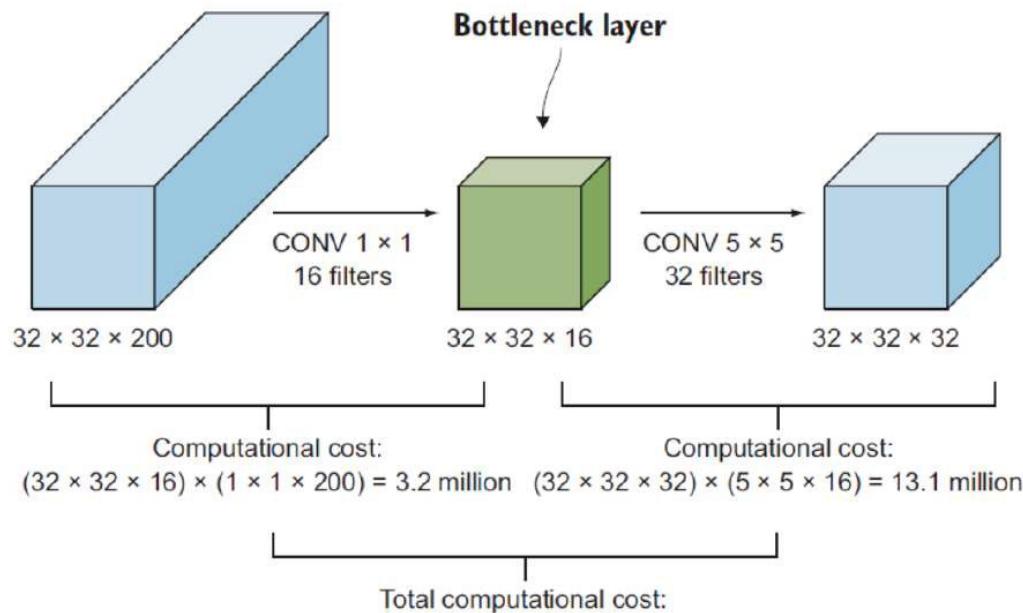
- Number of weights: $(3 \times 3 \times C) \times C + (3 \times 3)C \times C + (3 \times 3 \times C) \times C = 27C^2$

1×1 Convolution



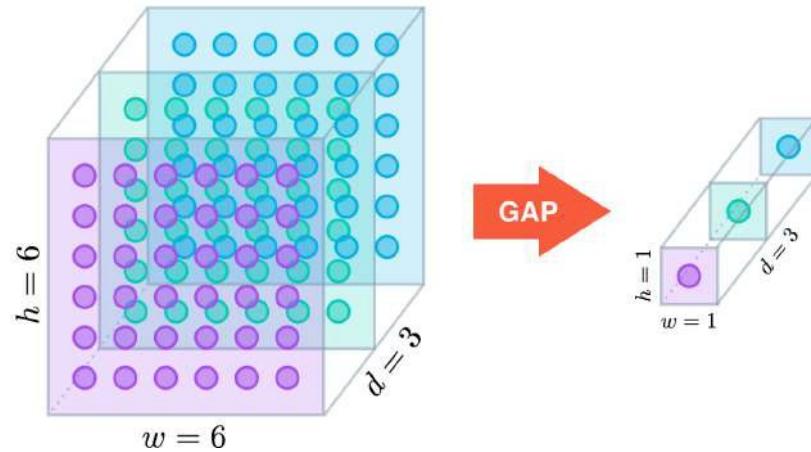
1×1 Convolution use case

- Assume we want to transform a $32 \times 32 \times 200$ tensor to a $32 \times 32 \times 32$ one using 32 5×5 filters. Thus, we need $(32 \times 32 \times 200) \times (5 \times 5 \times 32) \approx 163M$ operations!
- Instead we can use 1×1 convolution:



Global Average Pooling

Similar to max pooling layers, GAP layers are used to reduce the spatial dimensions of a three-dimensional tensor. However, GAP layers perform a more extreme type of dimensionality reduction, where a tensor with dimensions $h \times w \times d$ is reduced in size to have dimensions $1 \times 1 \times d$. GAP layers reduce each $h \times w$ feature map to a single number by simply taking the average of all hw values.



Why Gap

- GAP is used to replace the traditional fully connected layers in CNN.
- There is no parameter to optimize in the GAP thus overfitting is avoided at this layer
- GAP sums out the spatial information, thus it is more robust to spatial translations of the input.

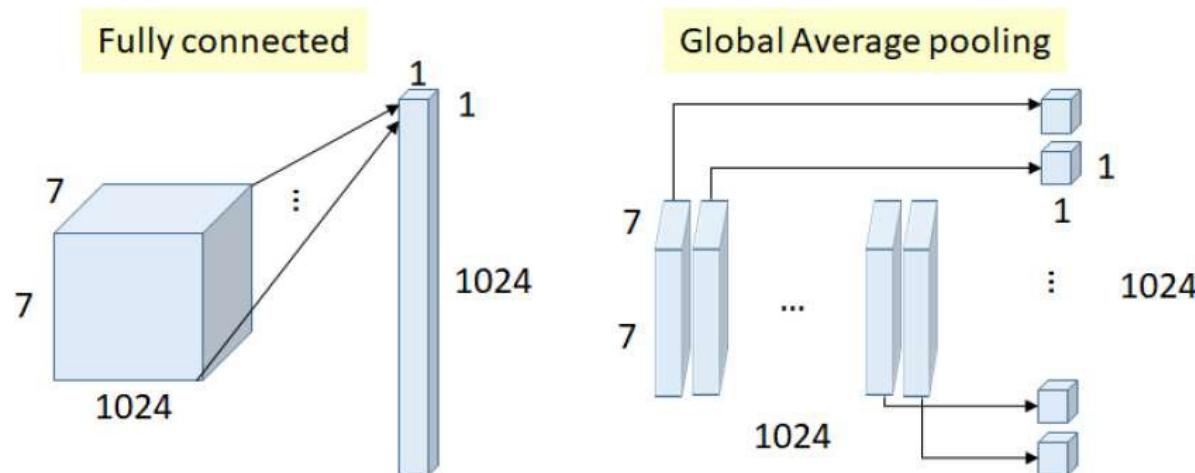


Figure: FC Layer vs GAP Layer