

DSCI 565: CONVOLUTIONAL NEURAL NETWORKS

*This content is protected and may not
be shared, uploaded, or distributed.*

Ke-Thia Yao

Lecture 8: 2025-09-22

Convolutional Neural Networks

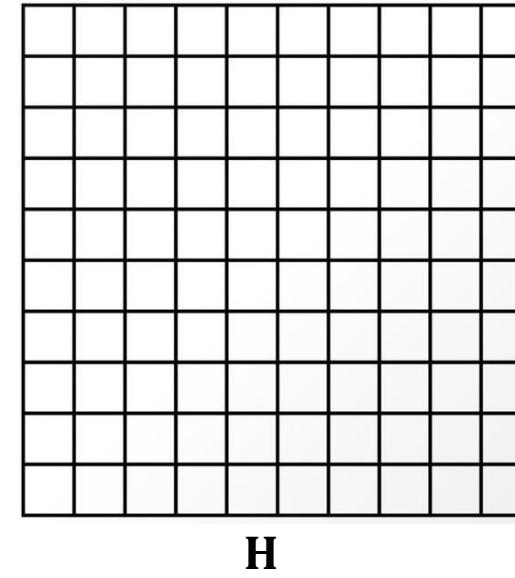
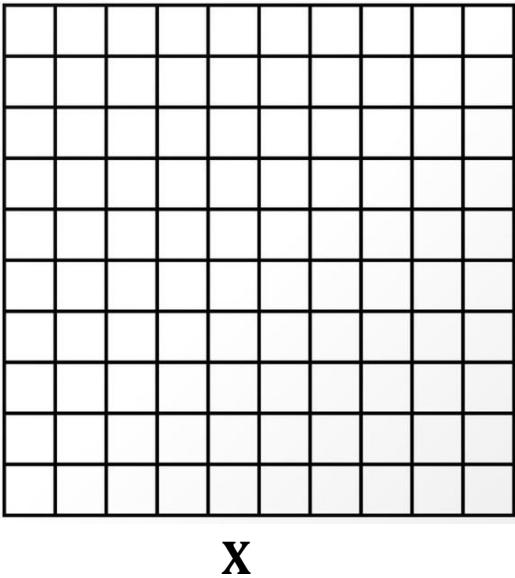
2

- Image data has a grid structure of pixels
- The neural networks that we have seen so far requires flattening of images in into vectors, which destroys the grid structure
- Can we somehow preserve this grid structure?

Initial Attempt

3

- Give an input image X , to preserve the grid structure we require the hidden layer H to have same shape as X
- Let $[X]_{i,j}$ and $[H]_{i,j}$ denote the pixel at (i, j)
- $[H]_{i,j}$ captures some feature related to (i, j)



Initial Attempt

4

- Then, if we use a fully connect architecture as before we need
 - A second-order tensor for the bias $[U]_{i,j}$
 - A fourth-order tensor for the weights $[W]_{i,j,kl}$
- Then, to compute $[H]_{i,j}$

$$[H]_{i,j} = [U]_{i,j} + \sum_k \sum_l [W]_{i,j,k,l} [X]_{k,l}$$

$(10^3)^4 = 10^{12}$

- For a megapixel image (1000x1000), the weight matrix would have 10^{12} (a trillion parameters)
- The largest Large Language Models (LLMs) is about 200 billion parameters
- What else is wrong with this initial attempt?

Translational Invariance

5

- Detecting Waldo at one part of the image should be no different than any other part

$$[\mathbf{H}]_{i,j} = [\mathbf{U}]_{i,j} + \sum_k \sum_l [\mathbf{W}]_{\text{no},k,l} [\mathbf{X}]_{k,l}$$



Locality

6

- To capture the feature at location (i, j) we should not have to look at the entire image
- Only look at nearby pixels $(i + a, j + b)$ and $a, b \in [-\Delta, +\Delta]$

$$\begin{aligned} [\mathbf{H}]_{i,j} &= [\mathbf{U}]_{i,j} + \sum_k \sum_l [\mathbf{W}]_{i,j,k,l} [\mathbf{X}]_{k,l} \\ &= [\mathbf{U}]_{i,j} + \sum_a \sum_b [\mathbf{V}]_{i,j,a,b} [\mathbf{X}]_{i+a,j+b} \end{aligned}$$

Translational
Invariance

$$[\mathbf{H}]_{i,j} = u + \sum_a \sum_b [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a,j+b}$$

Locality

$$[\mathbf{H}]_{i,j} = u + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a,j+b}$$

Channels

7

- So far, we have ignored the fact that some images have color channels
- With color
 - input becomes a third-order tensor $[\mathbf{X}]_{i,j,c}$ and
 - weight because a third-order tensor as well $[\mathbf{V}]_{a,b,c}$
- Also, instead of just one hidden feature map we want multiple feature maps $[\mathbf{H}]_{i,j,d}$, and weight becomes $[\mathbf{V}]_{a,b,c,d}$

$$[\mathbf{H}]_{i,j} = u + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a,j+b}$$

$$[\mathbf{H}]_{i,j,d} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} \sum_c [\mathbf{V}]_{a,b,c,d} [\mathbf{X}]_{i+a,j+b,c}$$

Convolution vs Cross-Correlation

8

- In machine learning, we call the previous formulas convolution
- But, in mathematics it is called cross-correlation
- In mathematics, convolution direction is “flipped”

$$(f * g)(x) = \int f(z)g(x - z)dz.$$

Cross-Correlation Operation

9

- Example of a cross-correlation operation using a 2×2 kernel

Input	Kernel	Output	
$\begin{array}{ c c c } \hline 0 & 1 & 2 \\ \hline 3 & 4 & 5 \\ \hline 6 & 7 & 8 \\ \hline \end{array}$	$* \quad \begin{array}{ c c } \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array}$	$= \quad \begin{array}{ c c } \hline 19 & 25 \\ \hline 37 & 43 \\ \hline \end{array}$	$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19,$ $1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25,$ $3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37,$ $4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$

- The resulting output shape reduce by 1 in each dimension
- In general, if the input size is $n_h \times n_w$ and kernel size is $k_h \times k_w$ output shape is

$$(n_h - k_h + 1, n_w - k_w + 1)$$

Inspiration for Mammal Visual Cortex

10

- Cats (and monkeys) have many neurons in the visual cortex of the brain have small local receptive fields
- They only activate when a limited region of the visual field is stimulated
- Moreover, some of these neurons respond to lines in certain directions

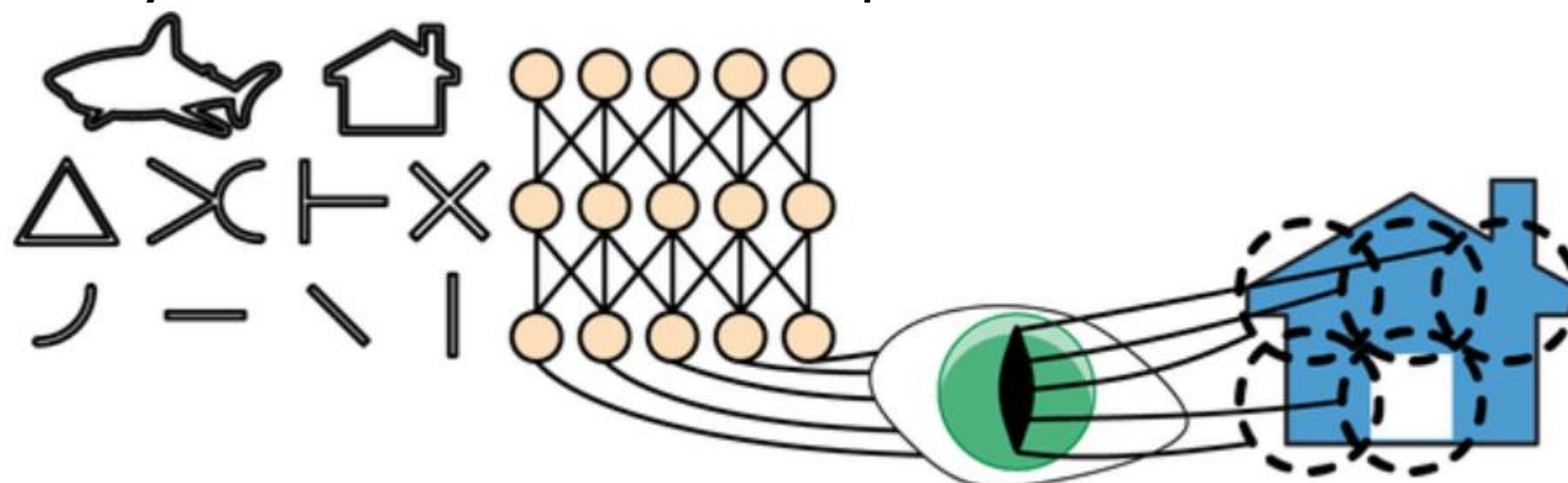
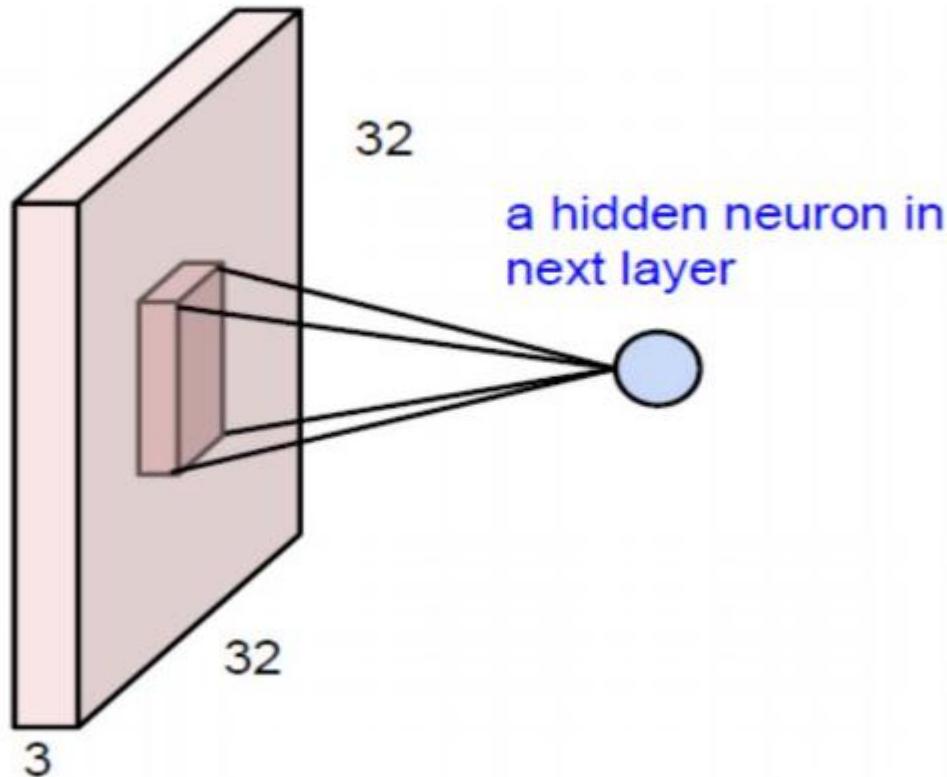


Figure 14-1. Biological neurons in the visual cortex respond to specific patterns in small regions of the visual field called receptive fields; as the visual signal makes its way through consecutive brain modules, neurons respond to more complex patterns in larger receptive fields

Convolution Filter

11

Use convolution filter to mimic neurons with limited receptive fields



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

Notebook

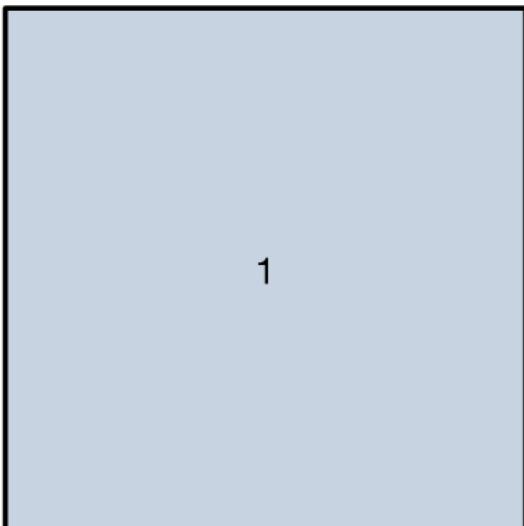
12

□ chapter_convolutional-neural-networks/conv-layer.ipynb

Padding

13

- For convolutions
 - Pixels on edges of images are used less often
 - Size of the image shrinks $\rightarrow n \rightarrow \lceil \frac{n-k+1}{k} \rceil$
- Shrinking becomes problematic if we have multiple convolution layers



kernel size

1×1

1	2	1
2	4	2
1	2	1

2×2

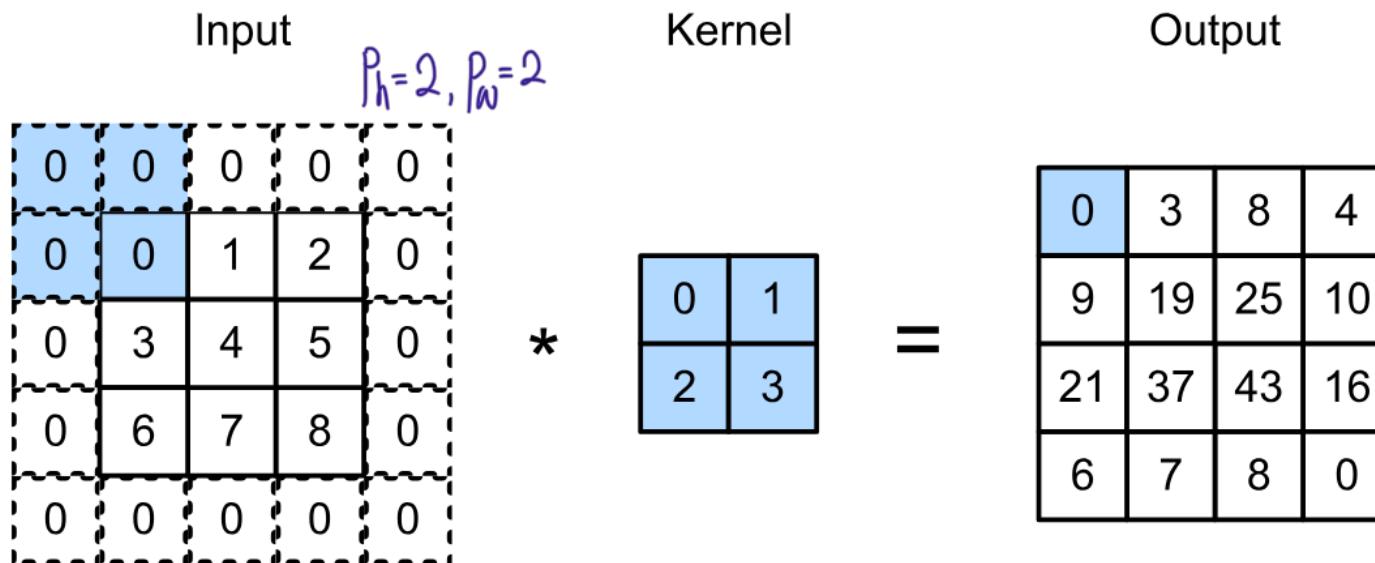
1	2	3	2	1
2	4	6	4	2
3	6	9	6	3
2	4	6	4	2
1	2	3	2	1

3×3

Padding

14

- Padding adds extra zero pixels around the border



- With padding p_h and p_w the output shape is

$$(n_h - k_h + p_h + 1, n_w - k_w + p_w + 1)$$

$$= (3 - 2 + 2 + 1, 3 - 2 + 2 + 1) = (4, 4)$$

Padding and Kernel Size

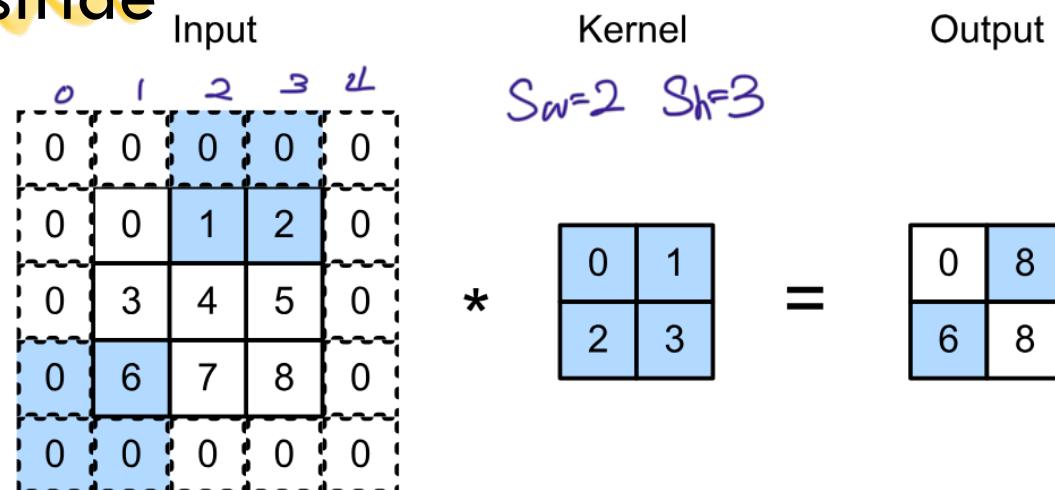
15

- To have the same number of pixels of each side of the border:
 p_h and p_w must be even
- To keep the same image size: $n - k + 1 + p = n$
 $p_h = k_h - 1$ and $p_w = k_w - 1 \Rightarrow p = k - 1$
- Then kernel size must be odd
- Having the same image size makes the interpretation of the hidden layer simpler:
 $[H]_{i,j}$ is calculated using pixels centered around $[X]_{i,j}$
- Note: PyTorch padding Conv2d(..., padding=p) is $p_h = p_w = 2p$

Stride

16

- By default, the convolution is slides one pixel at a time (stride of 1)
- To decrease computational complexity and/or to downsample, we can increase the stride



- With stride s_h and s_w , the output shape is

$$\lfloor (n_h - k_h + p_h + s_h) / s_h \rfloor \times \lfloor (n_w - k_w + p_w + s_w) / s_w \rfloor$$

Notebook

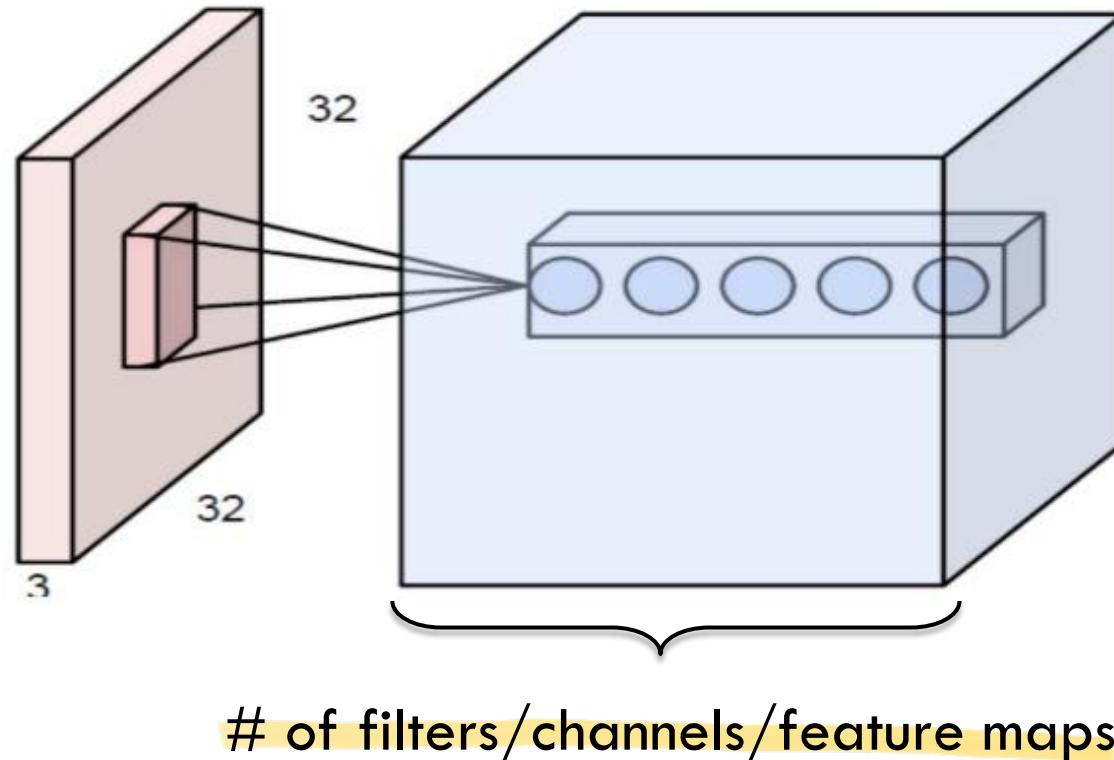
17

- chapter_convolutional-neural-networks/padding-and-strides.ipynb

Multiple Input and Multiple Output Channels

18

- Three input channels (RGB colors) with five output channels



- With input channel size c_i , the kernel shape becomes $c_i \times k_h \times k_w$

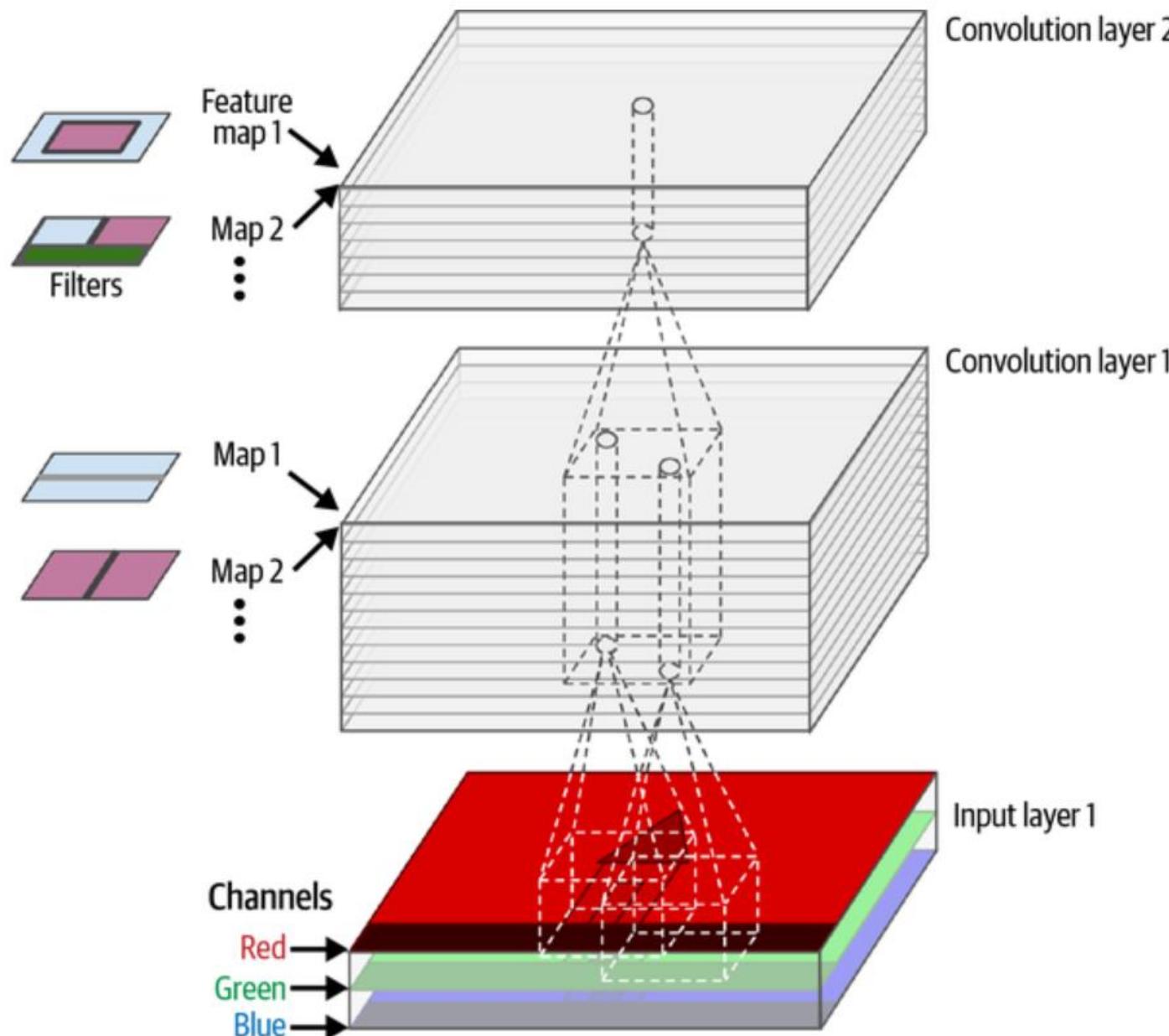
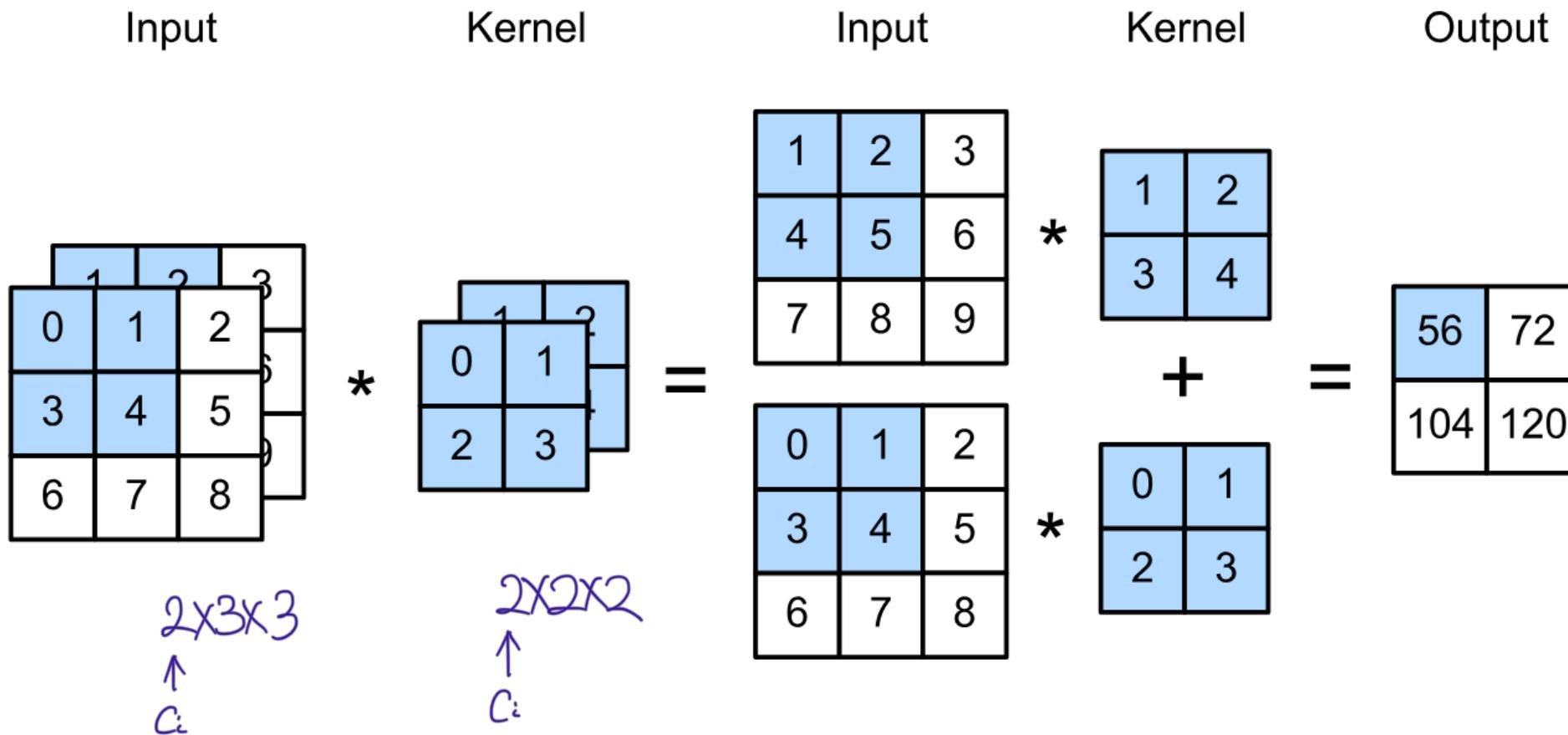


Figure 14-6. Two convolutional layers with multiple filters each (kernels), processing a color image with three color channels; each convolutional layer outputs one feature map per filter

Two Input Channel Examples

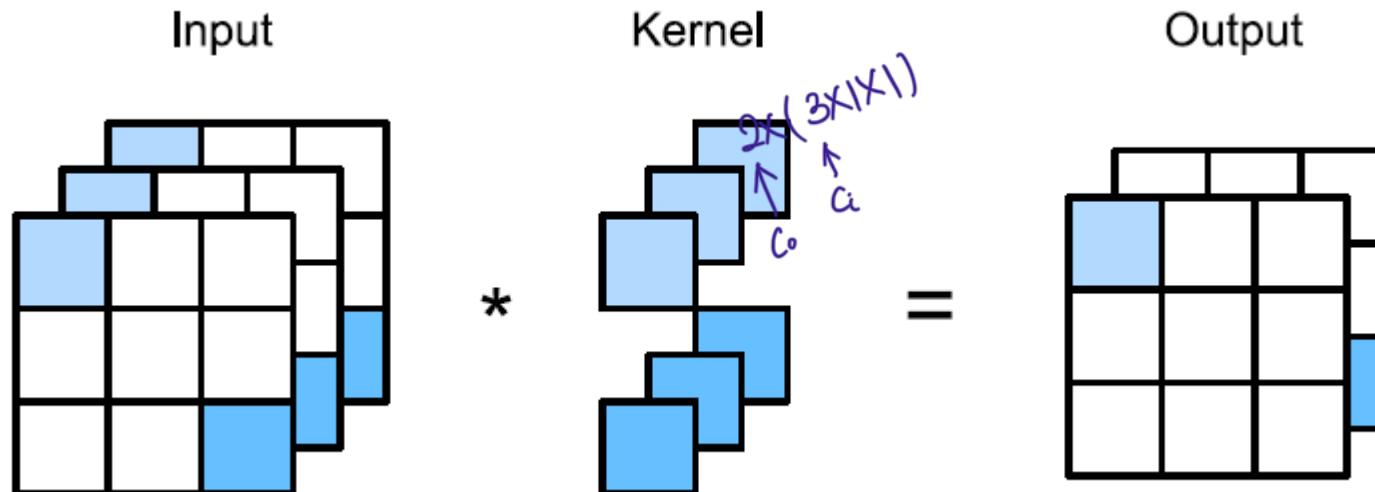
20



1×1 Convolutional Layer \Rightarrow shrink channel size

21

- With a 1×1 convolution, the convolution is only accessing one pixel, and no adjacent pixels
- But this one pixel consists of multiple channels
- \Rightarrow like performing a dot-product of channel vector with kernel's vector



Notebook

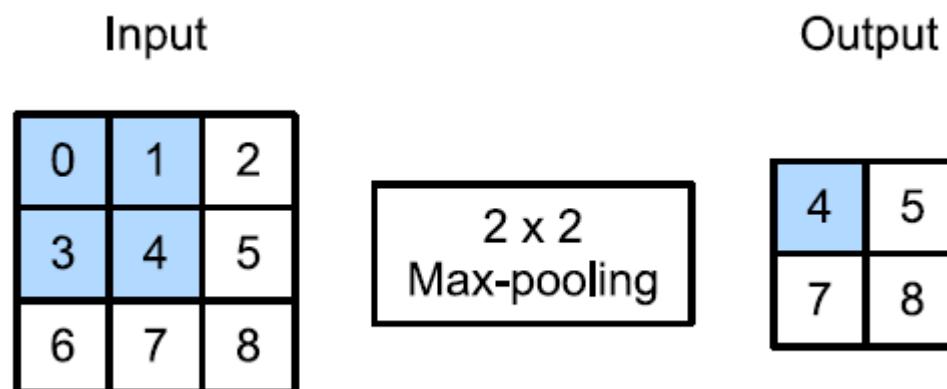
22

- chapter_convolutional-neural-networks/channels.ipynb

Pooling

23

- Pooling combines multiple pixels into one pixel, while keeping the number channels the same size
- Rational for pooling
 - Down sample to reduce size
 - Mitigating sensitivity of convolution to specific location
- Operations to combine pixels: max, average, min



Pooling and Stride

24

- The default stride for a pooling is the same shape as the pooling window shape
- E.g., if the pooling window shape is 2×2 , then the stride is 2×2

See chapter_convolutional-neural-networks/pooling.ipynb

LeNet

25

- Developed by Yann LeCun, LeNet is one of the first CNNs
- Designed for digit recognition, it reached error rate of 1%
- Error rate comparable to support vector machines
- LeNet is in use in ATM machine

