

PAT 463/563 (Fall 2025)

Music & AI

Lecture 9: Convolutional Neural Networks

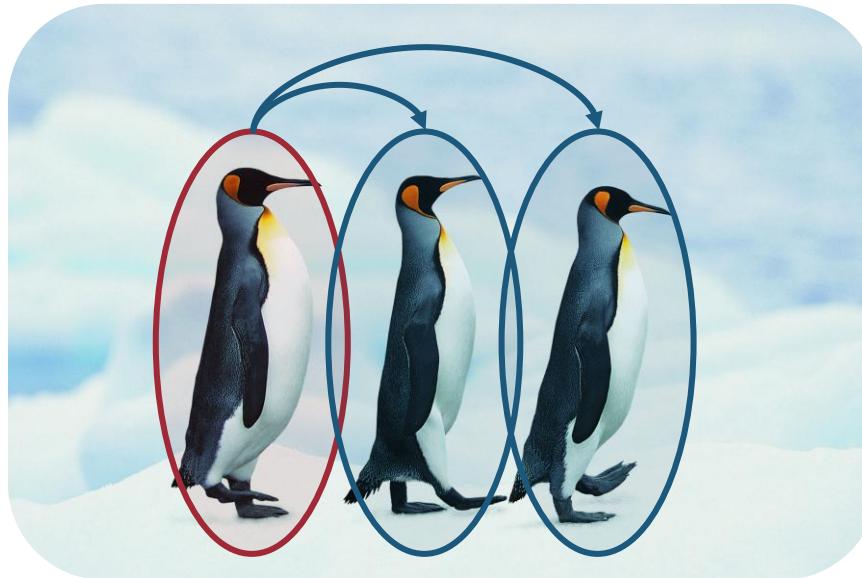
Instructor: Hao-Wen Dong

Convolutional Neural Networks (CNNs)

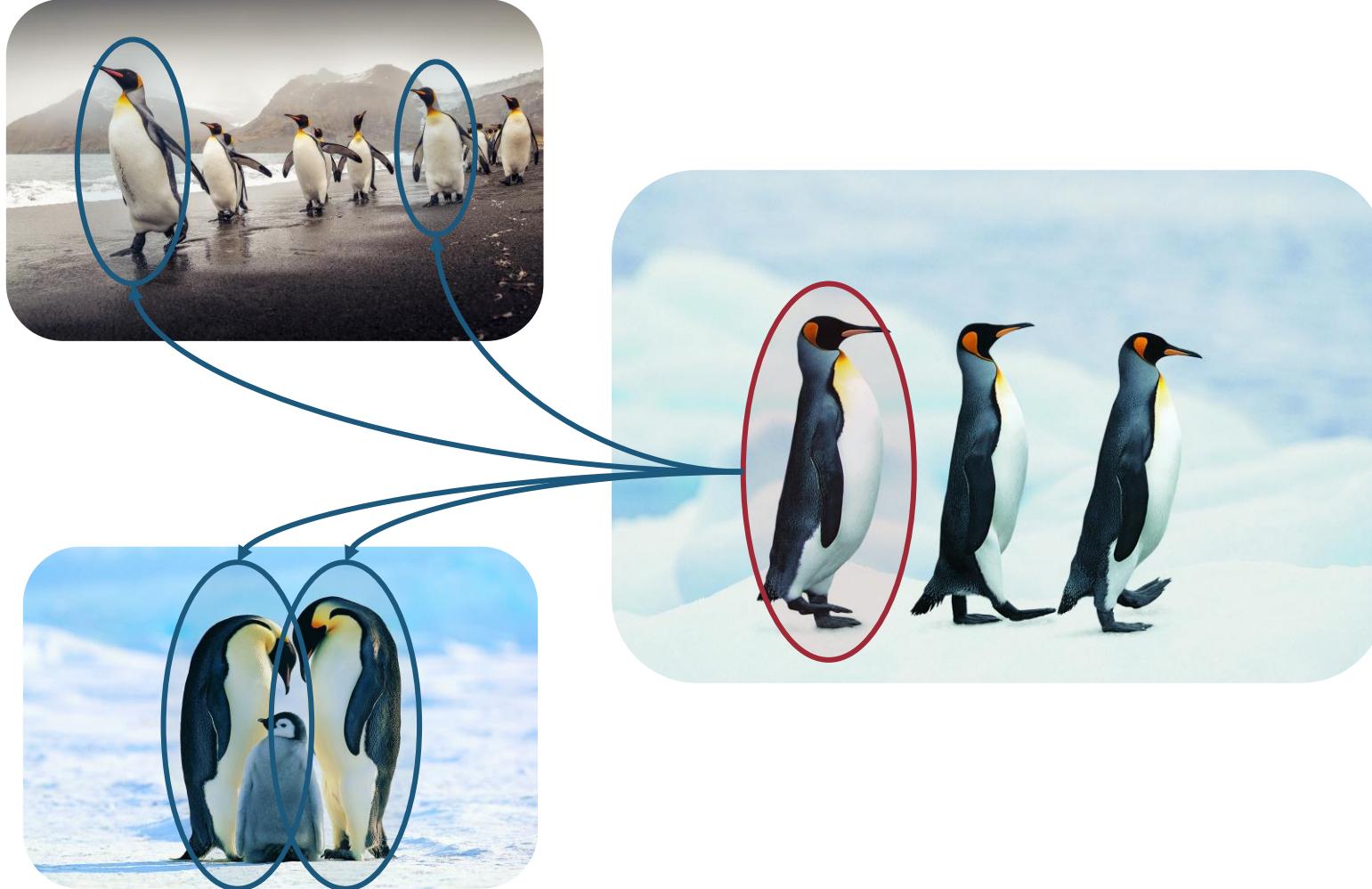
| Convolutional Neural Networks (CNNs)

- Intuition: Learn **reusable local pattern detector**
- Widely used in **computer vision**
- Also used for music and audio
 - Representing music as **piano rolls**
 - Representing audio as **spectrograms**

| Reusable Pattern Detectors



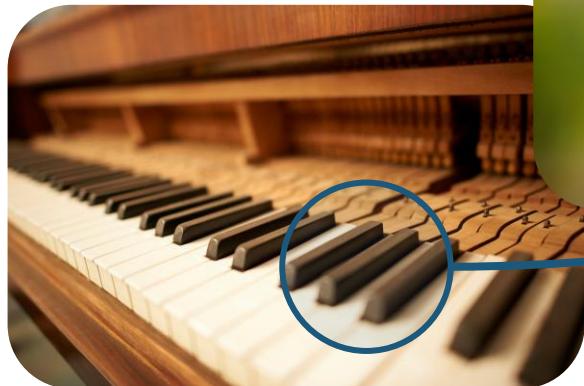
| Reusable Pattern Detectors



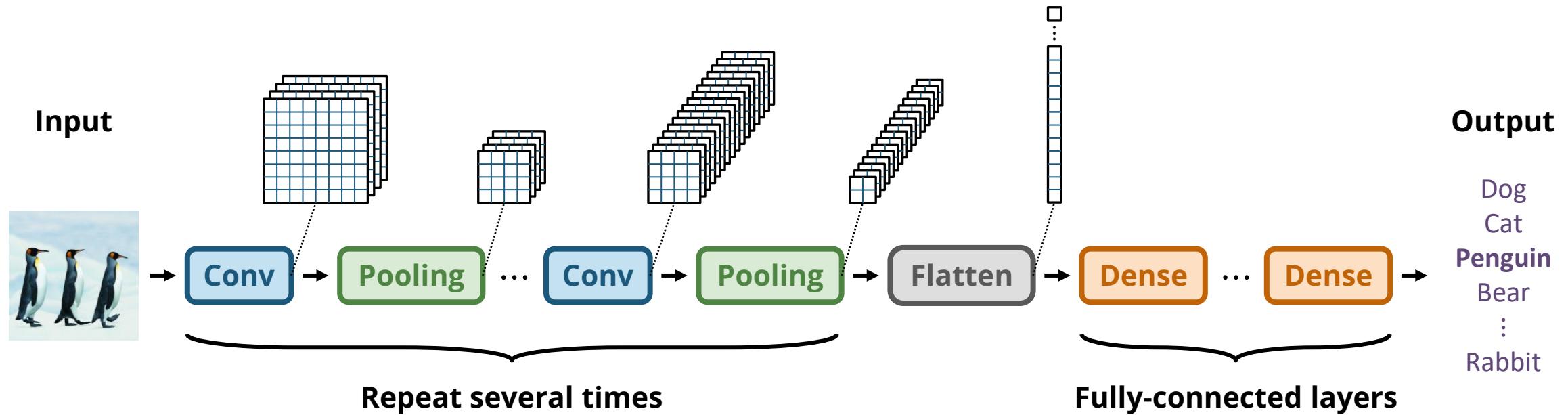
| Reusable Pattern Detectors



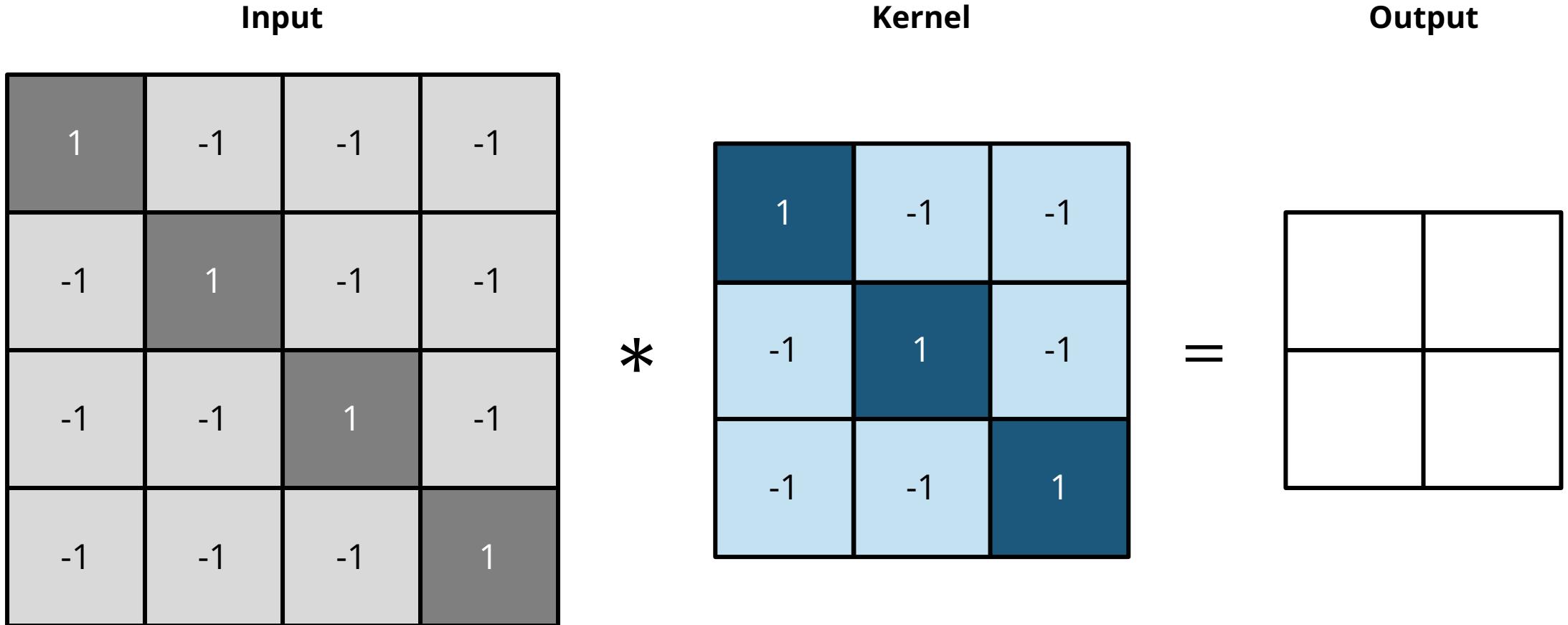
| Reusable Pattern Detectors



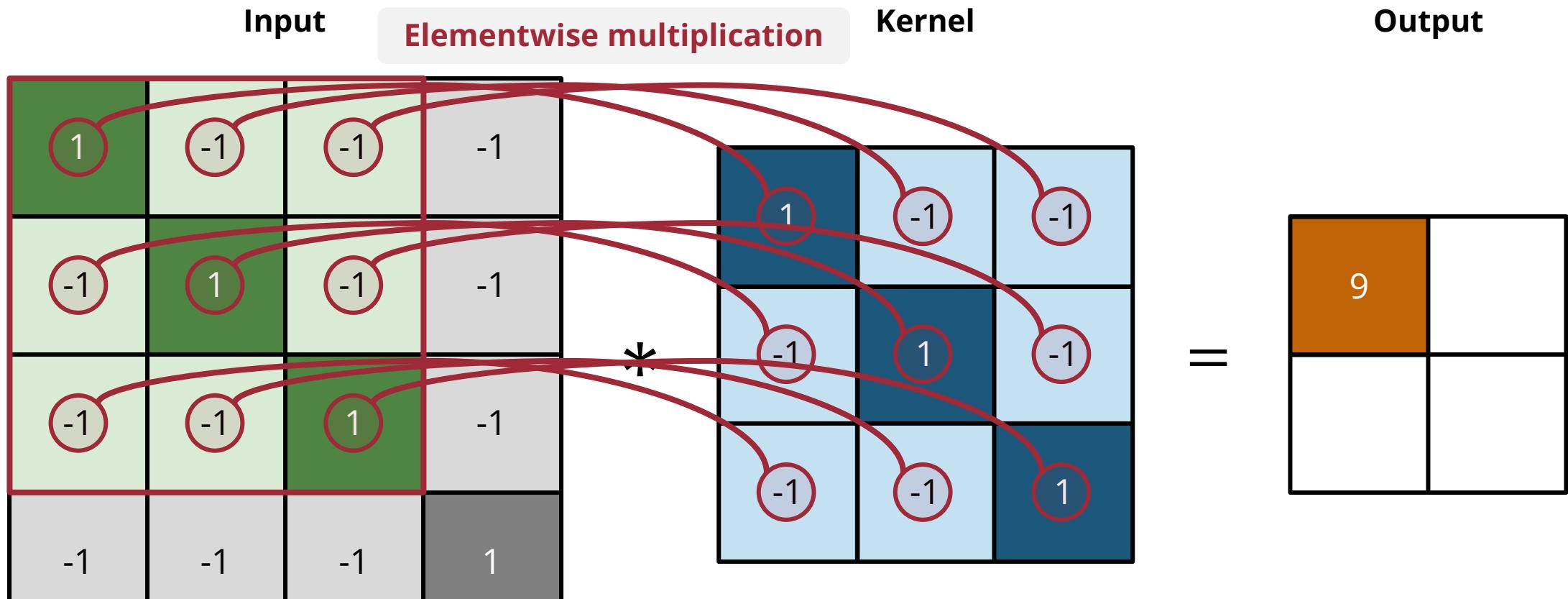
Convolutional Neural Network (CNNs)



2D Convolution

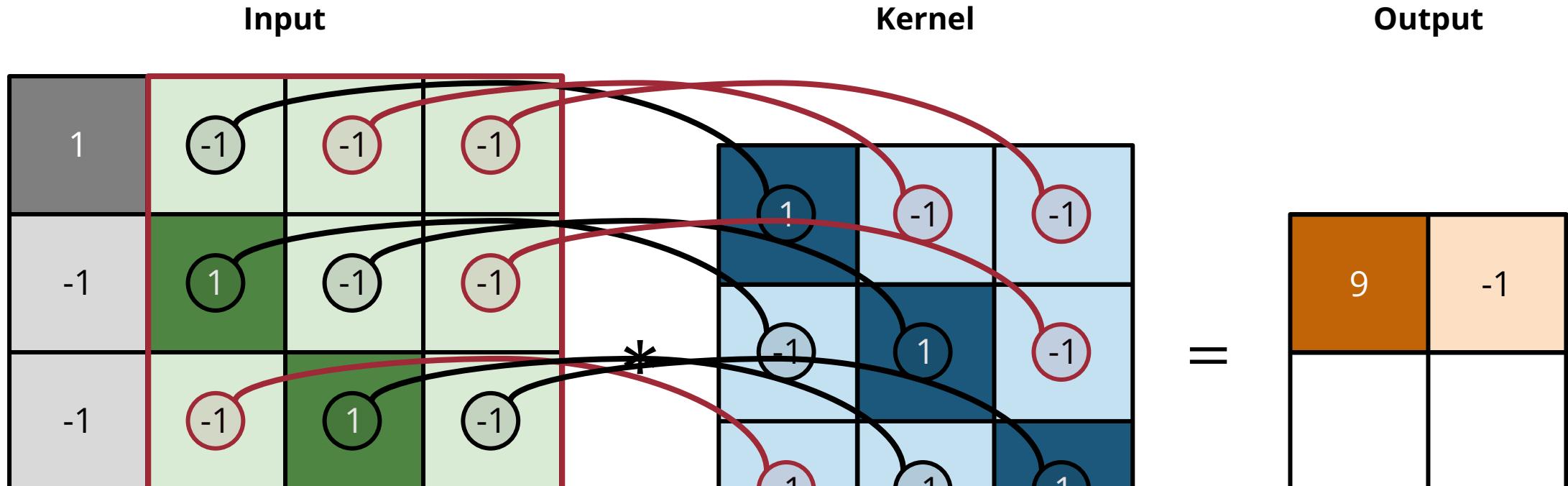


2D Convolution



$$\begin{aligned} (1 \times 1) &+ (-1 \times -1) + (-1 \times -1) \\ + (-1 \times -1) &+ (1 \times 1) + (-1 \times -1) = 9 \\ + (-1 \times -1) &+ (-1 \times -1) + (1 \times 1) \end{aligned}$$

2D Convolution



$$\begin{aligned} & (-1 \times 1) + (-1 \times -1) + (-1 \times -1) \\ & + (-1 \times 1) + (-1 \times 1) + (-1 \times -1) = -1 \\ & + (-1 \times -1) + (1 \times -1) + (-1 \times 1) \end{aligned}$$

2D Convolution

Input			
1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

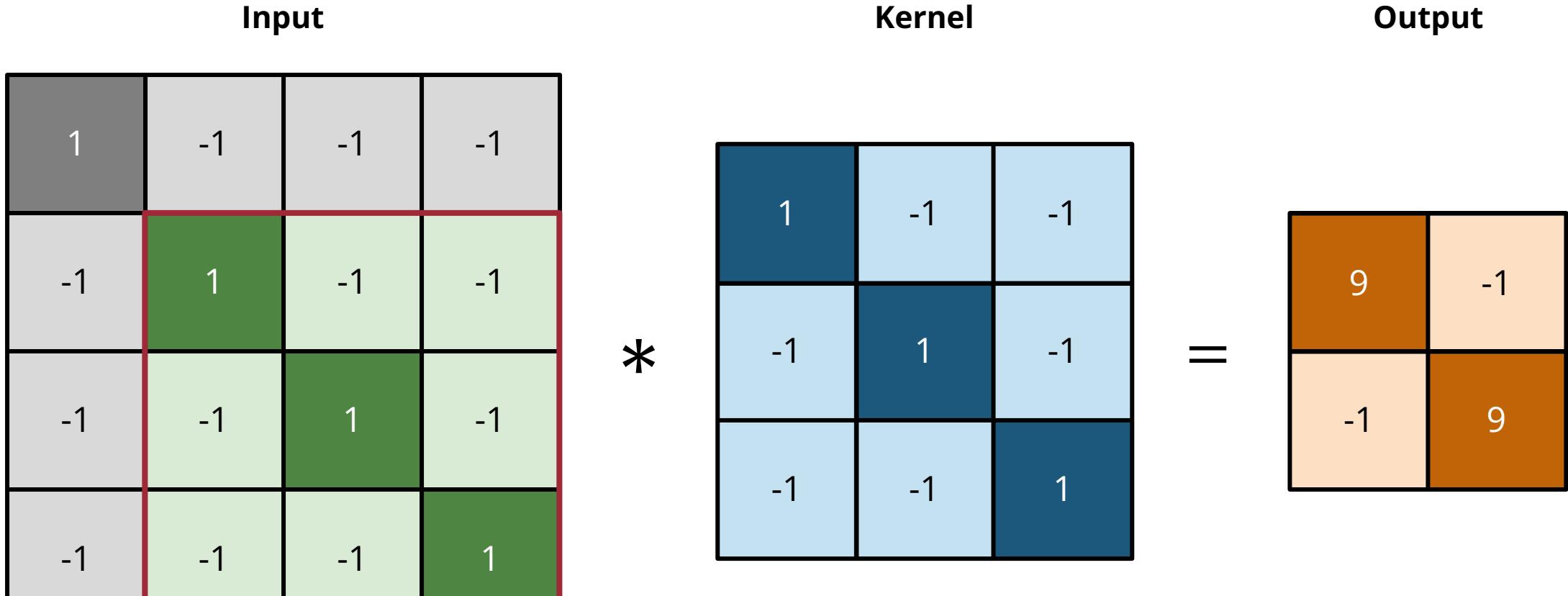
*

1	-1	-1
-1	1	-1
-1	-1	1

Output	
9	-1
-1	

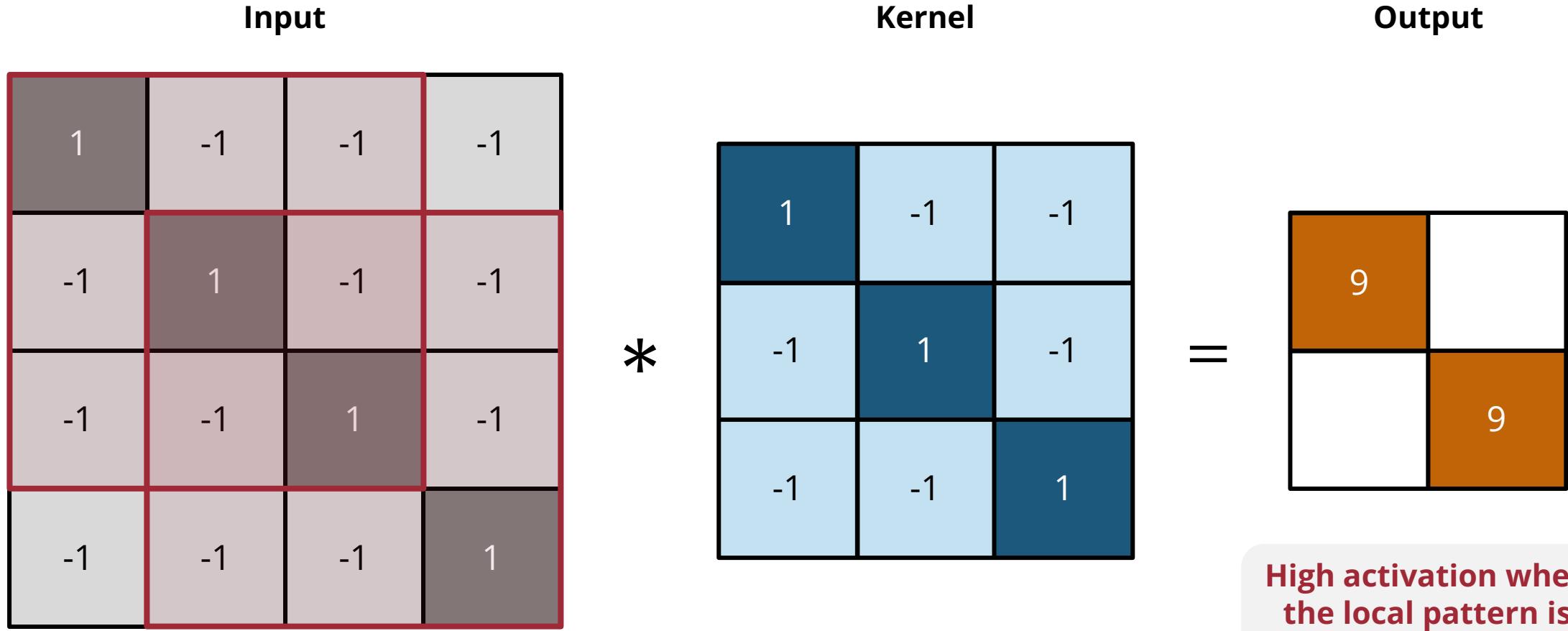
$$\begin{aligned} (-1 \times 1) + (1 \times -1) + (-1 \times -1) \\ + (-1 \times -1) + (-1 \times 1) + (1 \times -1) \\ + (-1 \times -1) + (-1 \times -1) + (-1 \times 1) \end{aligned} = -1$$

2D Convolution



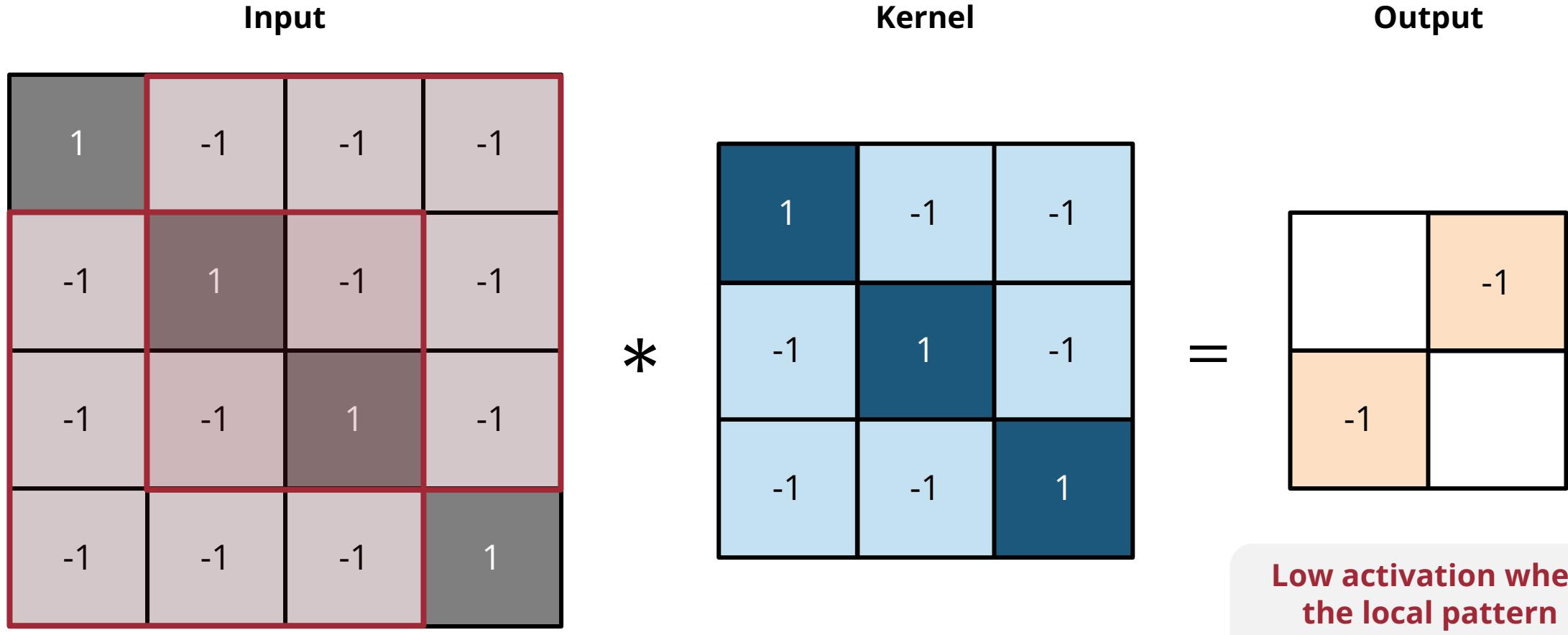
$$\begin{aligned} & (1 \times 1) + (-1 \times -1) + (-1 \times -1) \\ & + (-1 \times -1) + (1 \times 1) + (-1 \times -1) = 9 \\ & + (-1 \times -1) + (-1 \times -1) + (1 \times 1) \end{aligned}$$

2D Convolution



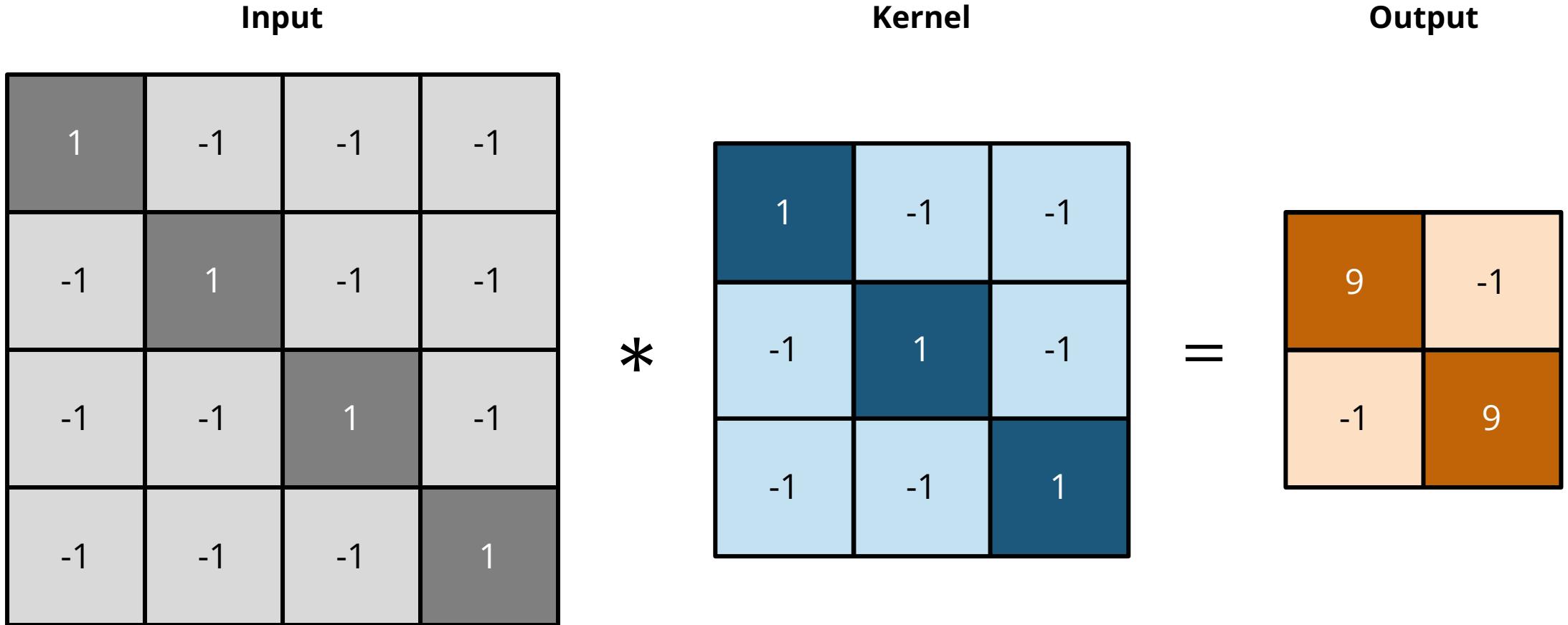
High activation when
the local pattern is
close to the kernel

2D Convolution

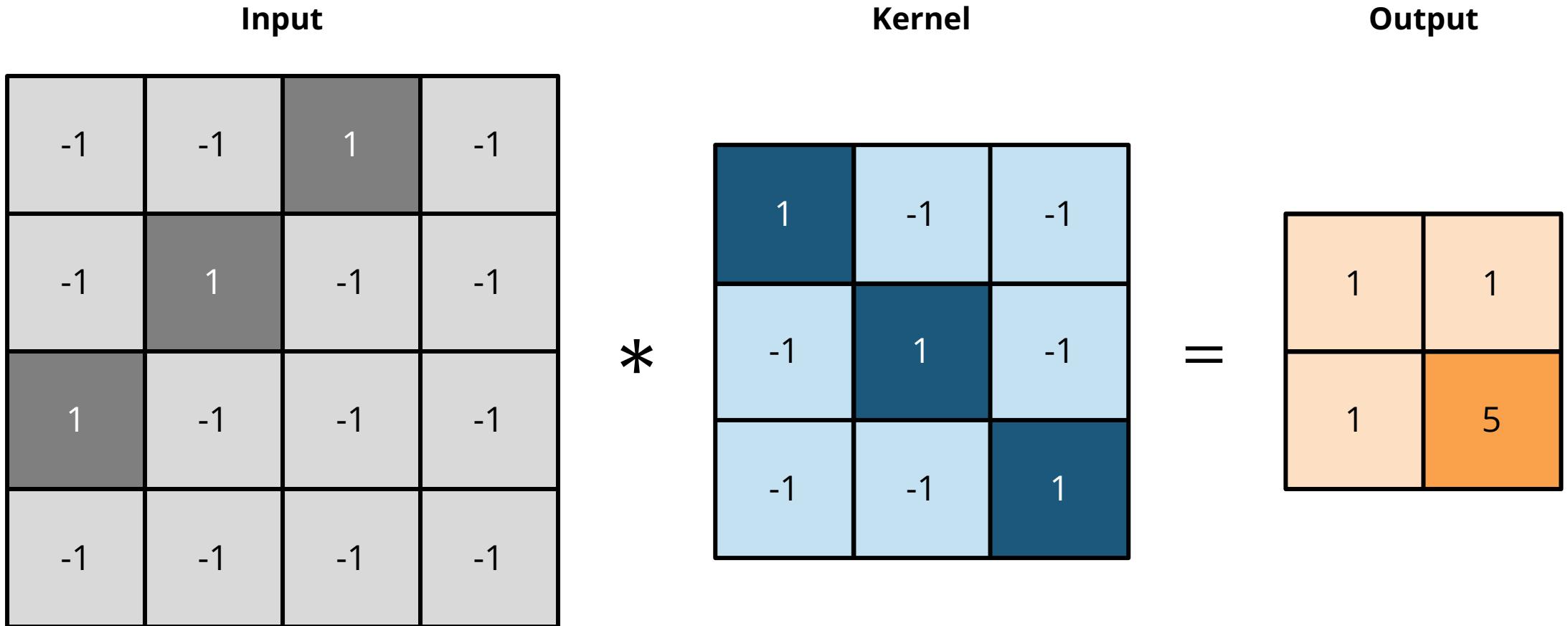


Low activation when
the local pattern
differs from the kernel

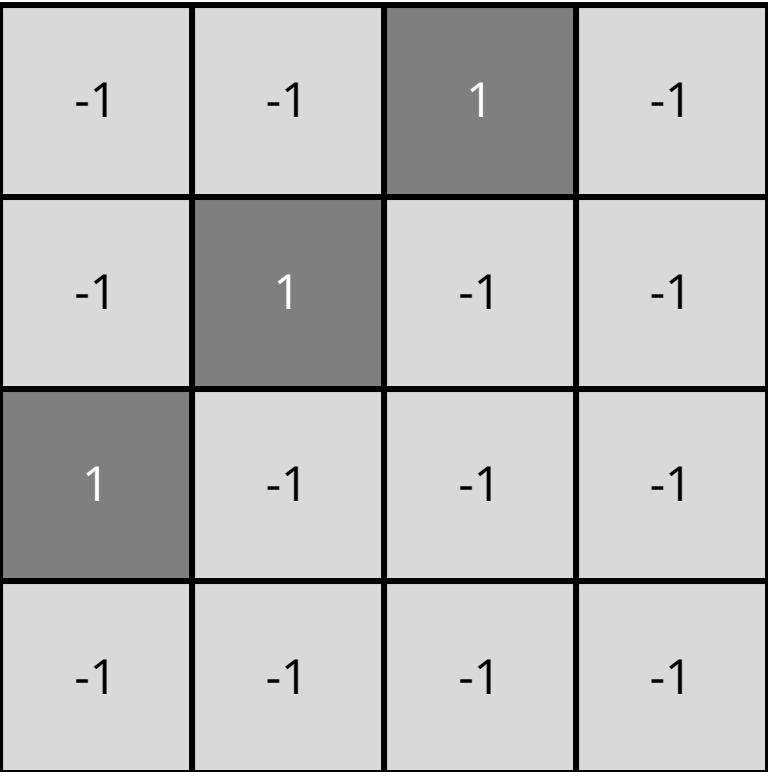
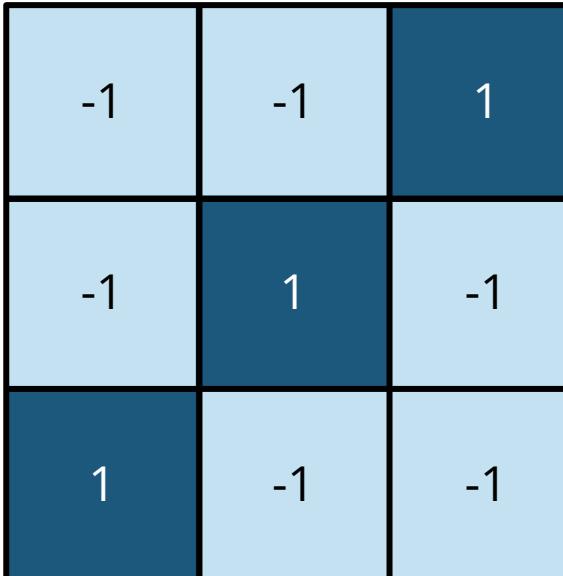
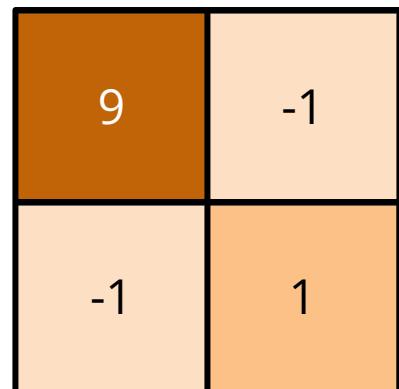
2D Convolution



2D Convolution

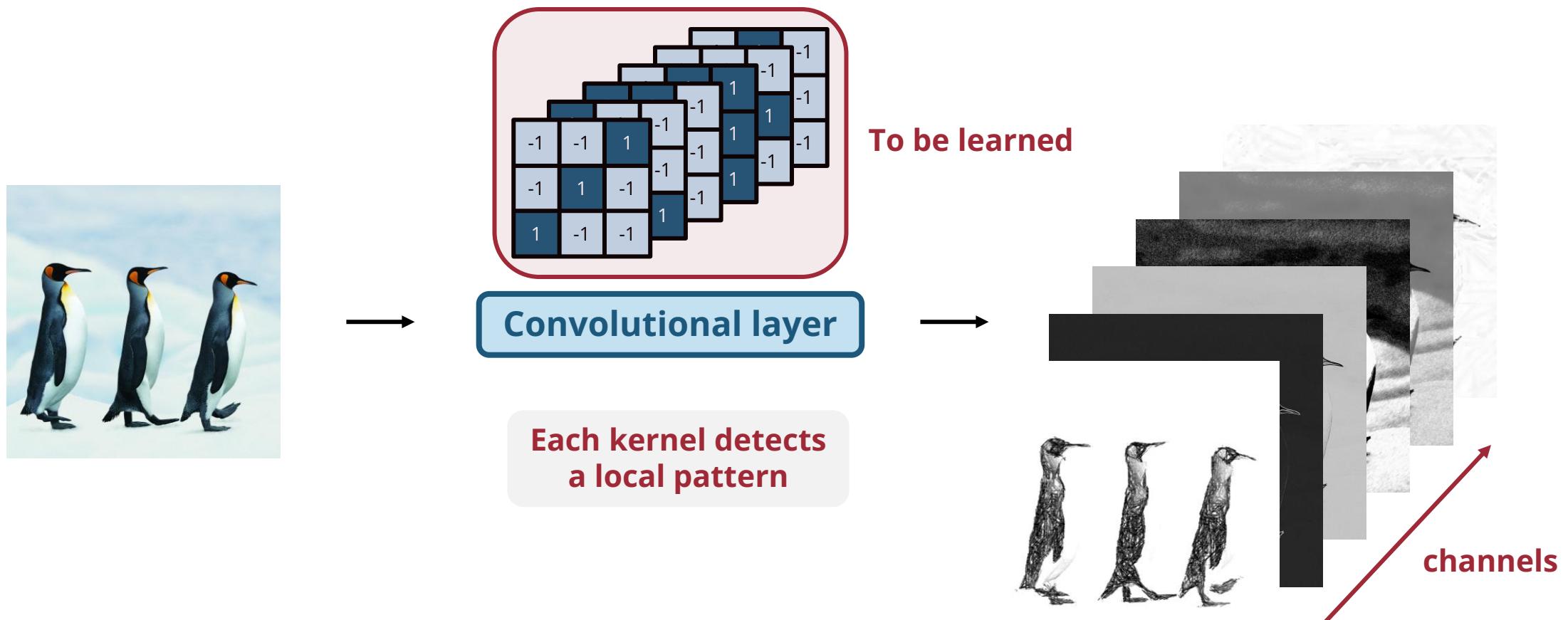


2D Convolution

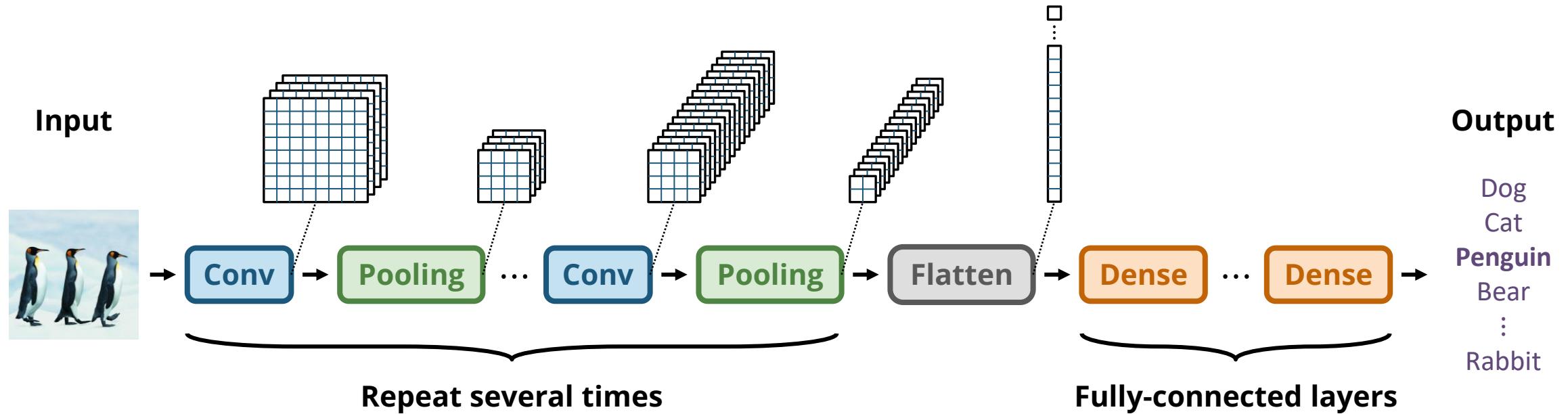
Input		Kernel
 $\begin{bmatrix} -1 & -1 & 1 & -1 \\ -1 & 1 & -1 & -1 \\ 1 & -1 & -1 & -1 \\ -1 & -1 & -1 & -1 \end{bmatrix}$	*	 $\begin{bmatrix} -1 & -1 & 1 \\ -1 & 1 & -1 \\ 1 & -1 & -1 \end{bmatrix}$
		=
		 $\begin{bmatrix} 9 & -1 \\ -1 & 1 \end{bmatrix}$

Convolutional Layer

- A convolutional layer consists of many **learnable kernels** (channels)



Convolutional Neural Network (CNNs)



Padding

padding="valid"

1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

*

1	-1	-1
-1	1	-1
-1	-1	1

=

9	-1
-1	9

padding="same"

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

*

1	-1	-1
-1	1	-1
-1	-1	1

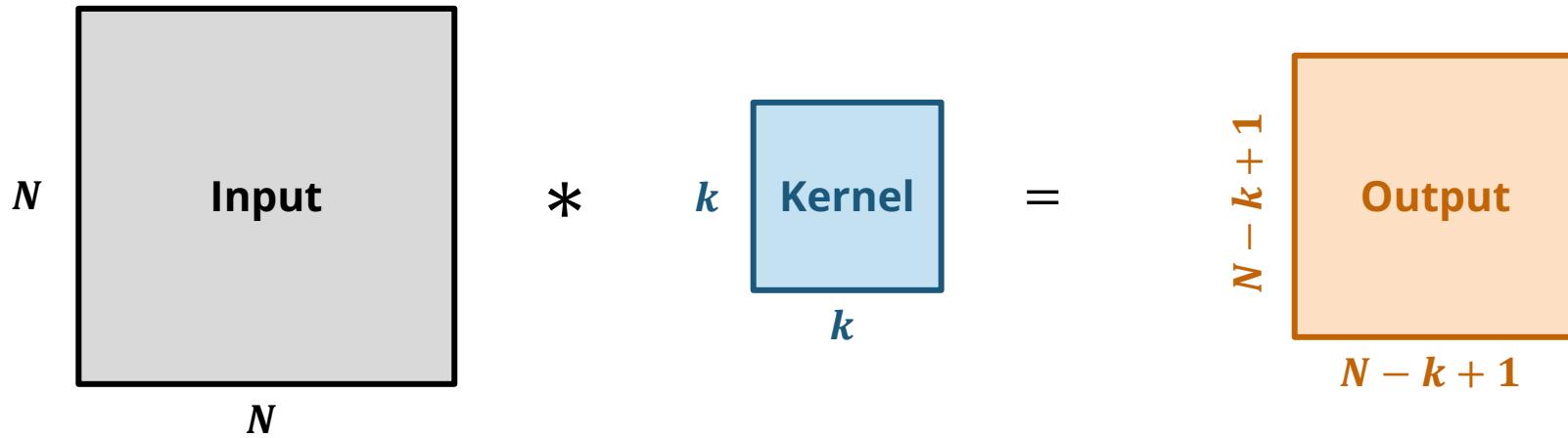
=

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

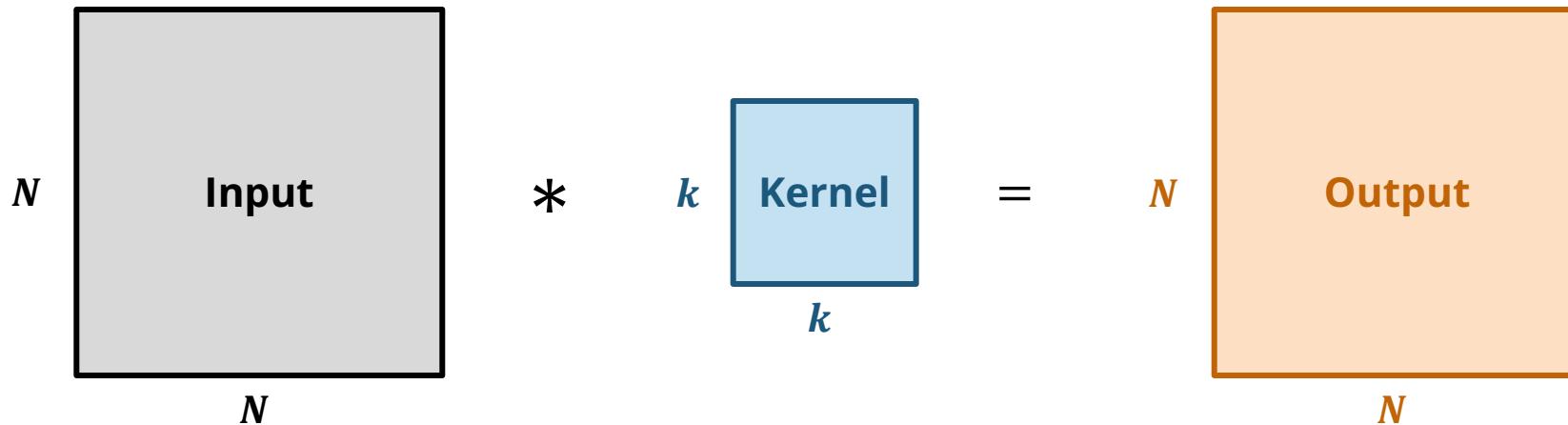
Keep the output of the same size as the input

Shapes

`padding="valid"`



`padding="same"`



| Striding

stride=2

The diagram illustrates the computation of a 2x2 convolution step on a 5x6 input matrix with stride 2, using a 3x3 kernel. The input matrix has values: row 1 [0, 0, 0, 0, 0, 0], row 2 [0, 1, -1, -1, -1, 0], row 3 [0, -1, 1, -1, -1, 0], row 4 [0, -1, -1, 1, -1, 0], row 5 [0, -1, -1, -1, 1, 0]. The kernel is: [1, -1, -1; -1, 1, -1; -1, -1, 1]. The result is a 2x2 matrix: [4, 0; 0, 9]. A green arrow highlights the receptive field of the top-left output unit.

$$\begin{matrix} & \text{stride=2} \\ \begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & -1 & -1 & 0 \\ 0 & -1 & 1 & -1 & -1 & 0 \\ 0 & -1 & -1 & 1 & -1 & 0 \\ 0 & -1 & -1 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} & * \begin{matrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{matrix} = \begin{matrix} 4 & 0 \\ 0 & 9 \end{matrix} \end{matrix}$$

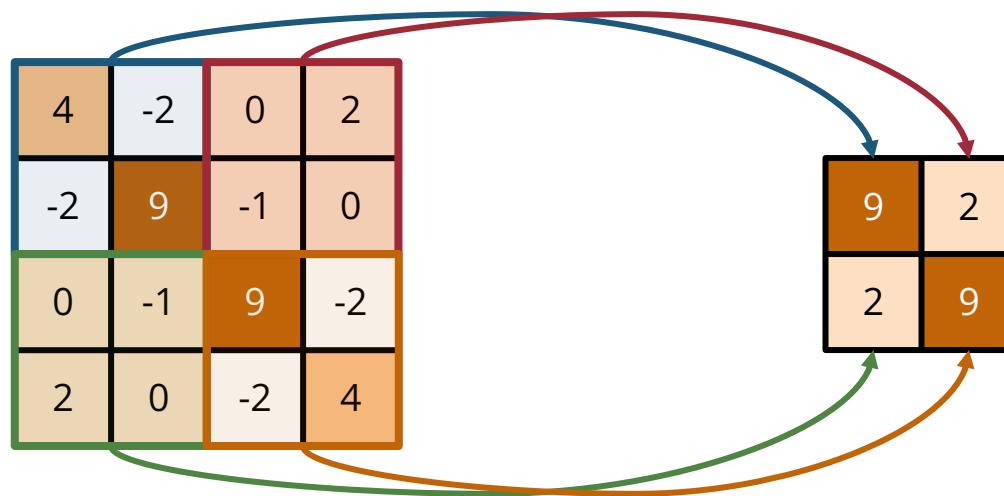
| Striding

stride=3

The diagram illustrates the convolution operation with a stride of 3. It shows two input matrices and one filter kernel. The first input matrix has values: 0, 0, 0; 0, 1, -1; 0, -1, 1; 0, -1, -1; 0, -1, -1; 0, 0, 0. The second input matrix has values: 0, 0, 0; 0, 0, 0; 0, 0, 0; 0, 0, 0; 0, 0, 0; 0, 0, 0. The filter kernel has values: 1, -1, -1; -1, 1, -1; -1, -1, 1. The result of the multiplication is a 2x2 matrix with values 4, 2; 2, 4.

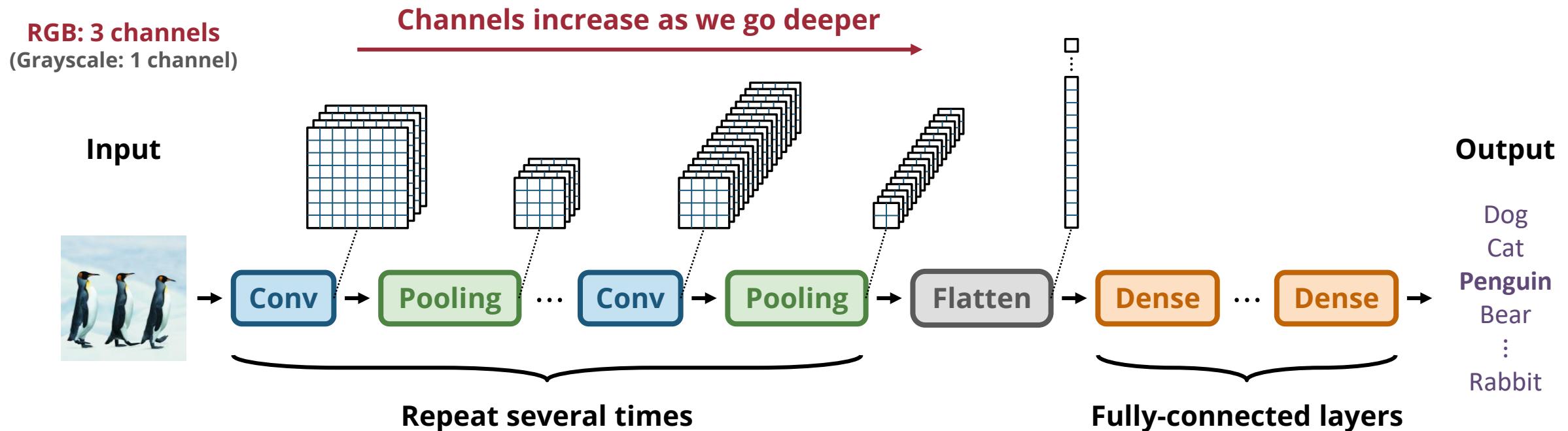
$$\begin{matrix} & \text{stride=3} \\ \begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & -1 & -1 & 0 \\ 0 & -1 & 1 & 1 & -1 & 0 \\ 0 & -1 & -1 & 1 & -1 & 0 \\ 0 & -1 & -1 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} & \begin{matrix} * & \begin{matrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{matrix} & = & \begin{matrix} 4 & 2 \\ 2 & 4 \end{matrix} \end{matrix} \end{matrix}$$

Max Pooling Layer

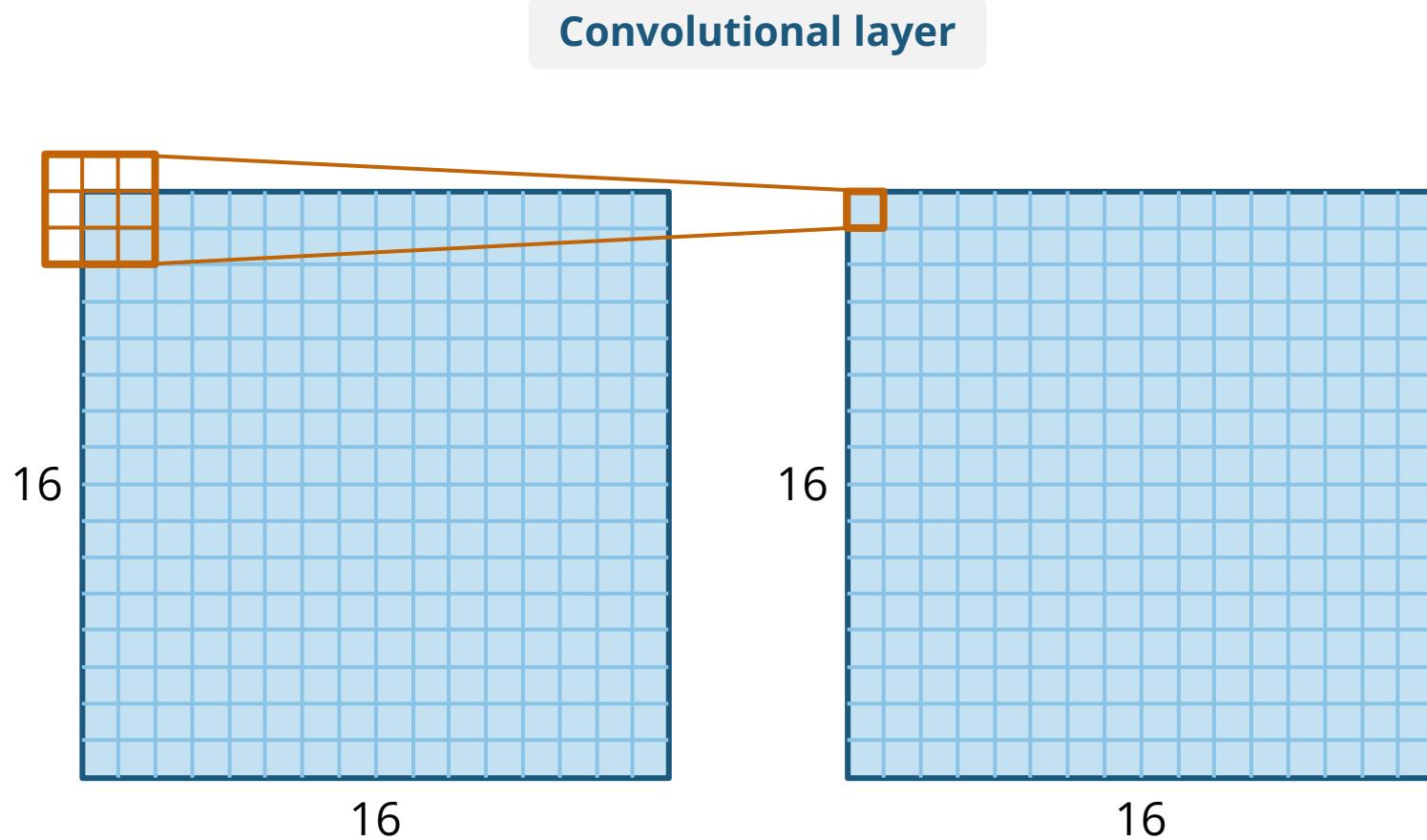


Downsample and keep the
strongest activation in each block

Convolutional Neural Network (CNNs)

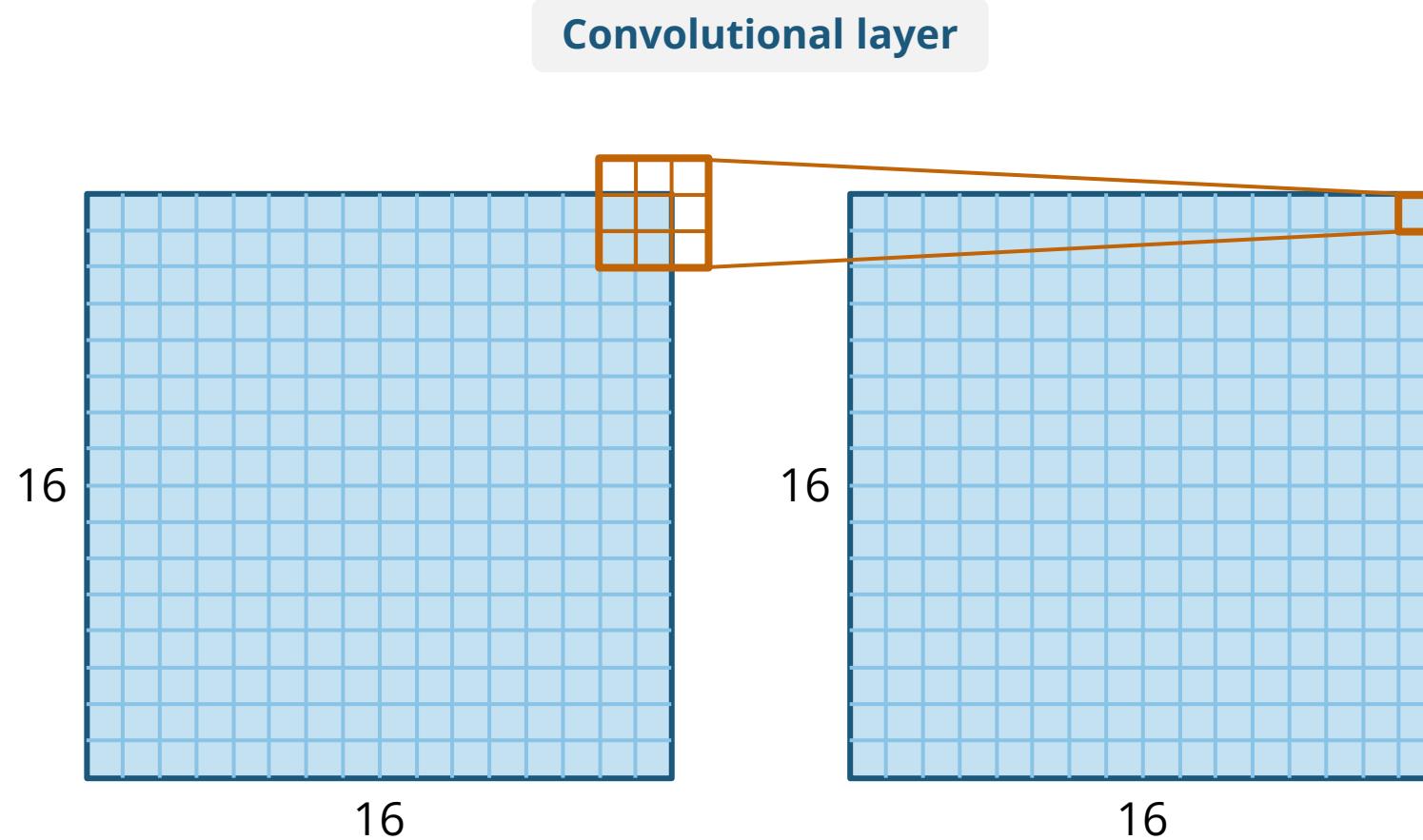


A Toy Example



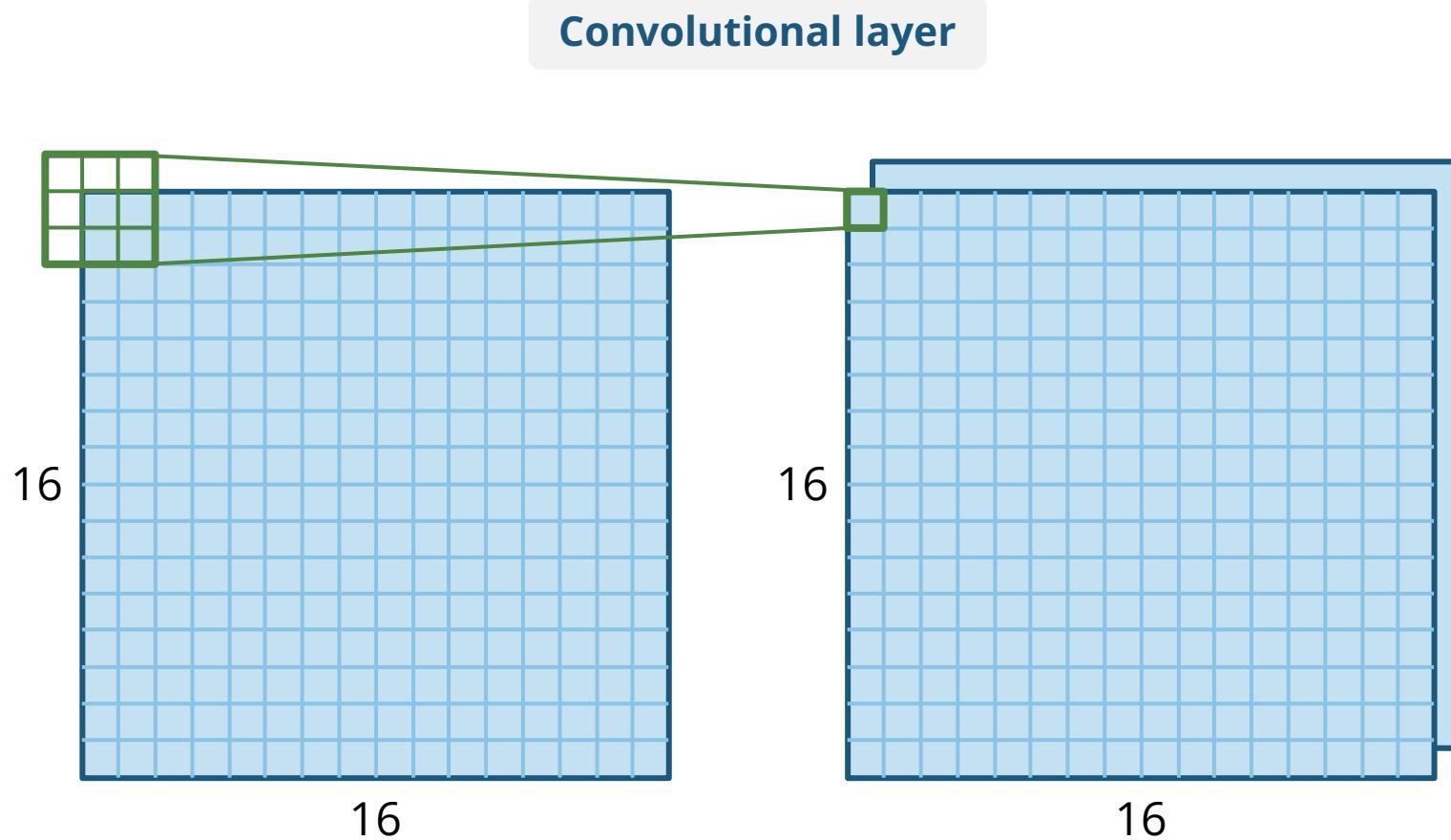
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example



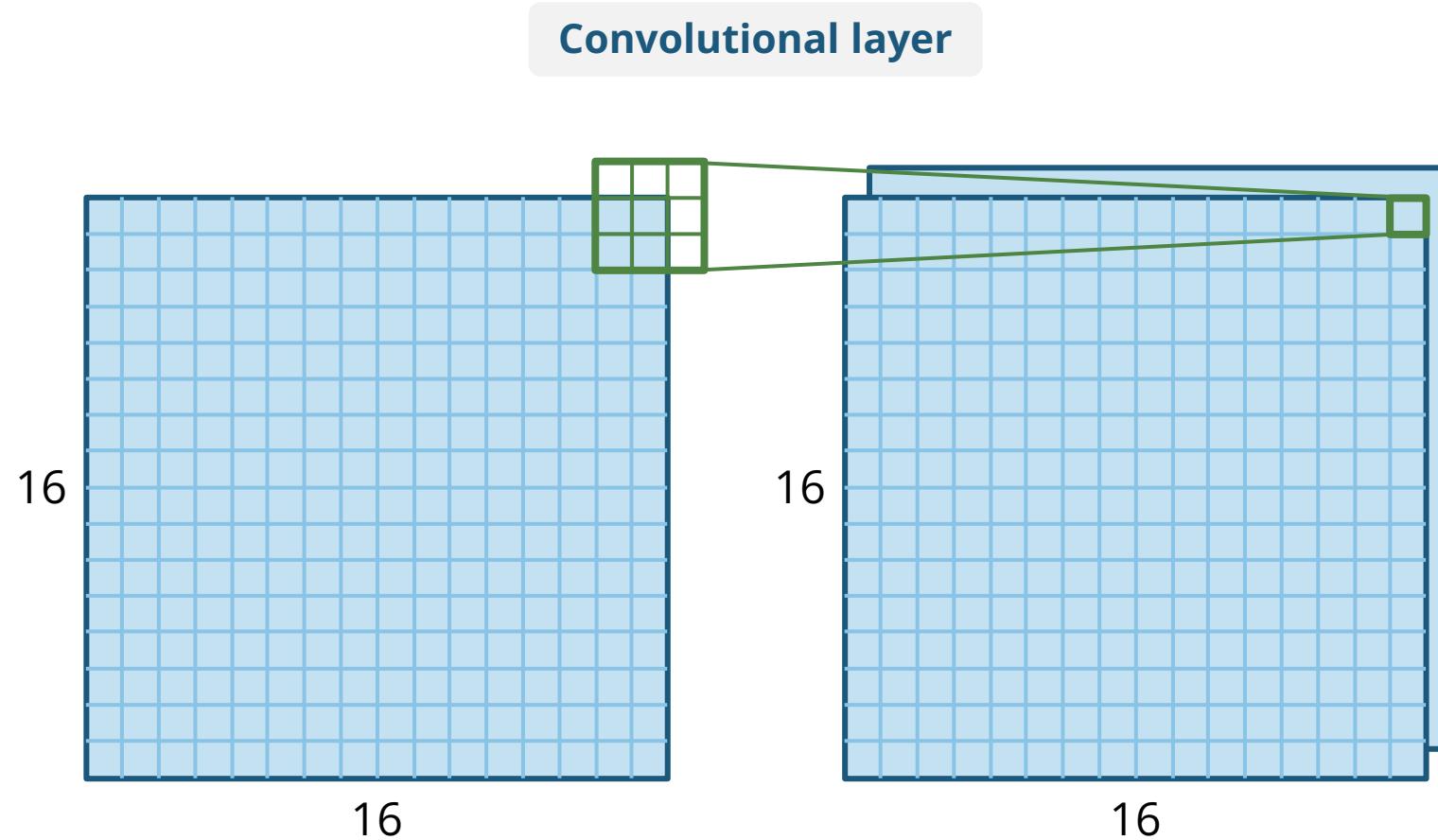
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example



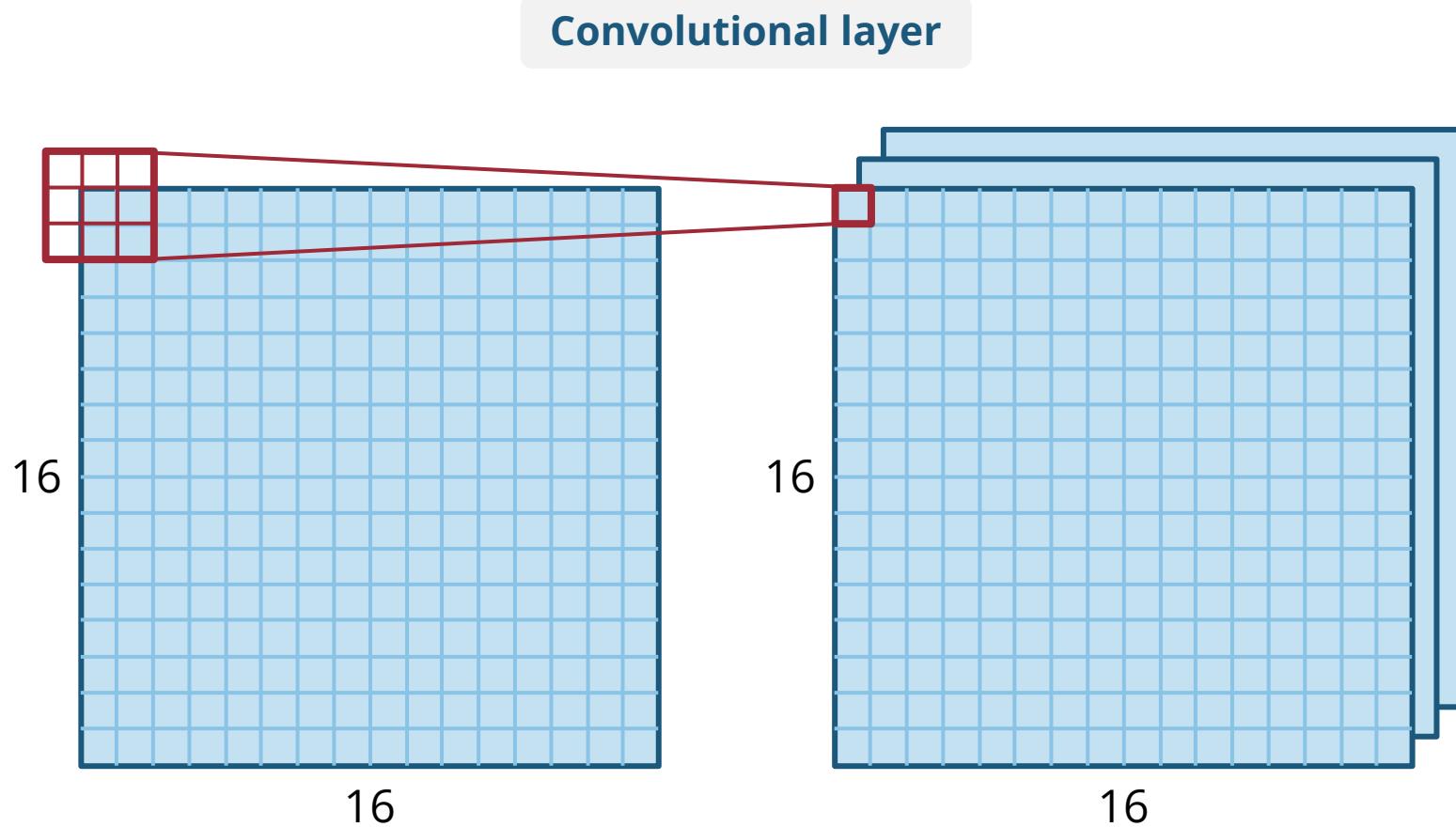
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example



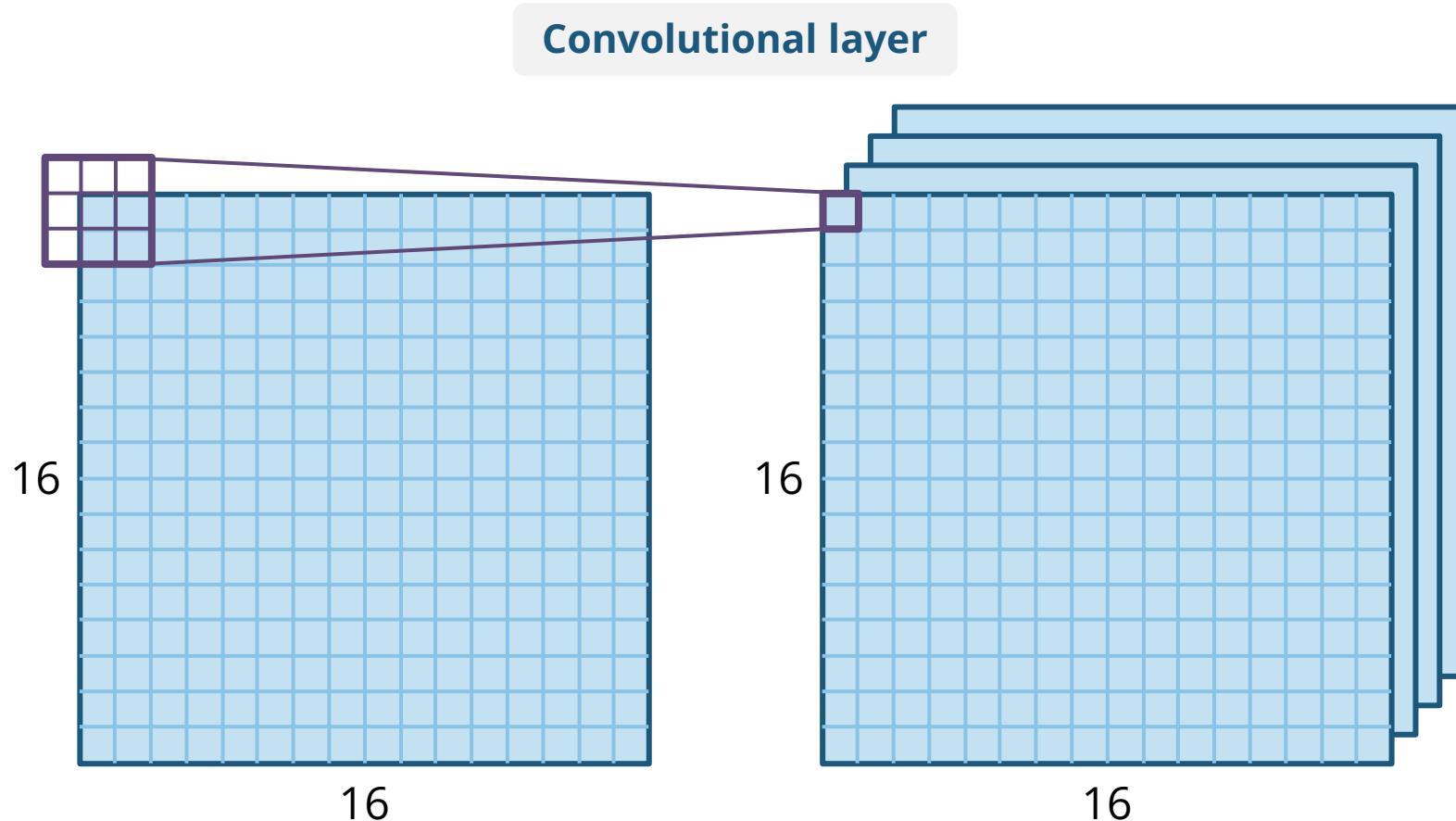
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example



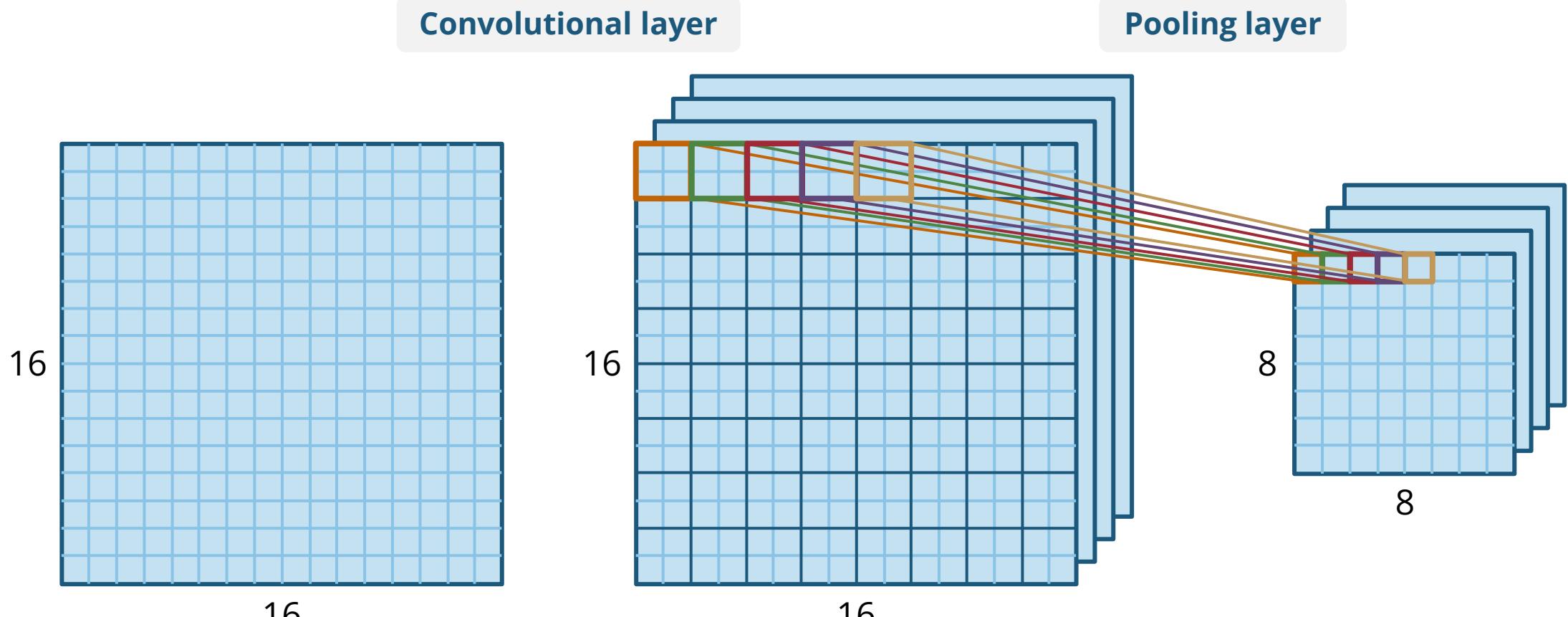
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example



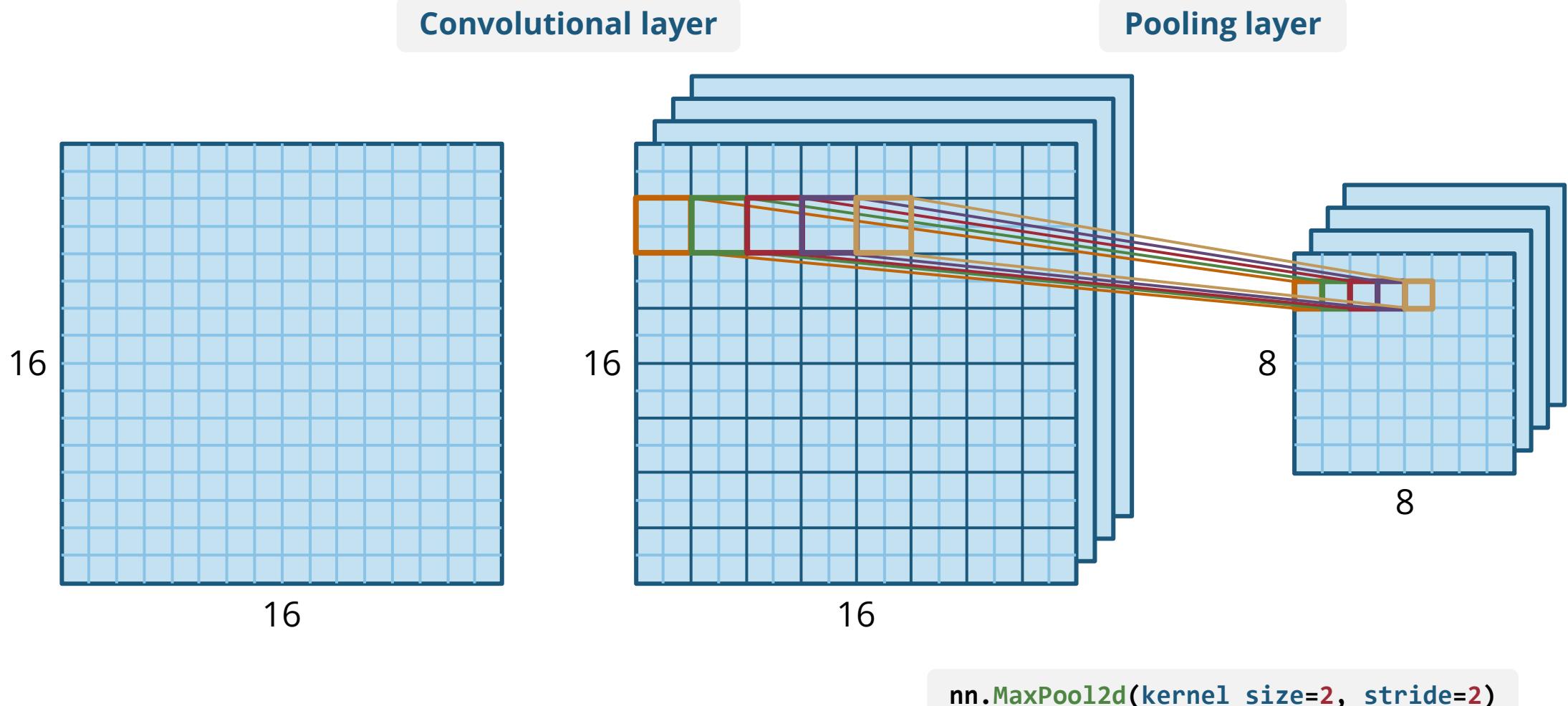
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

A Toy Example

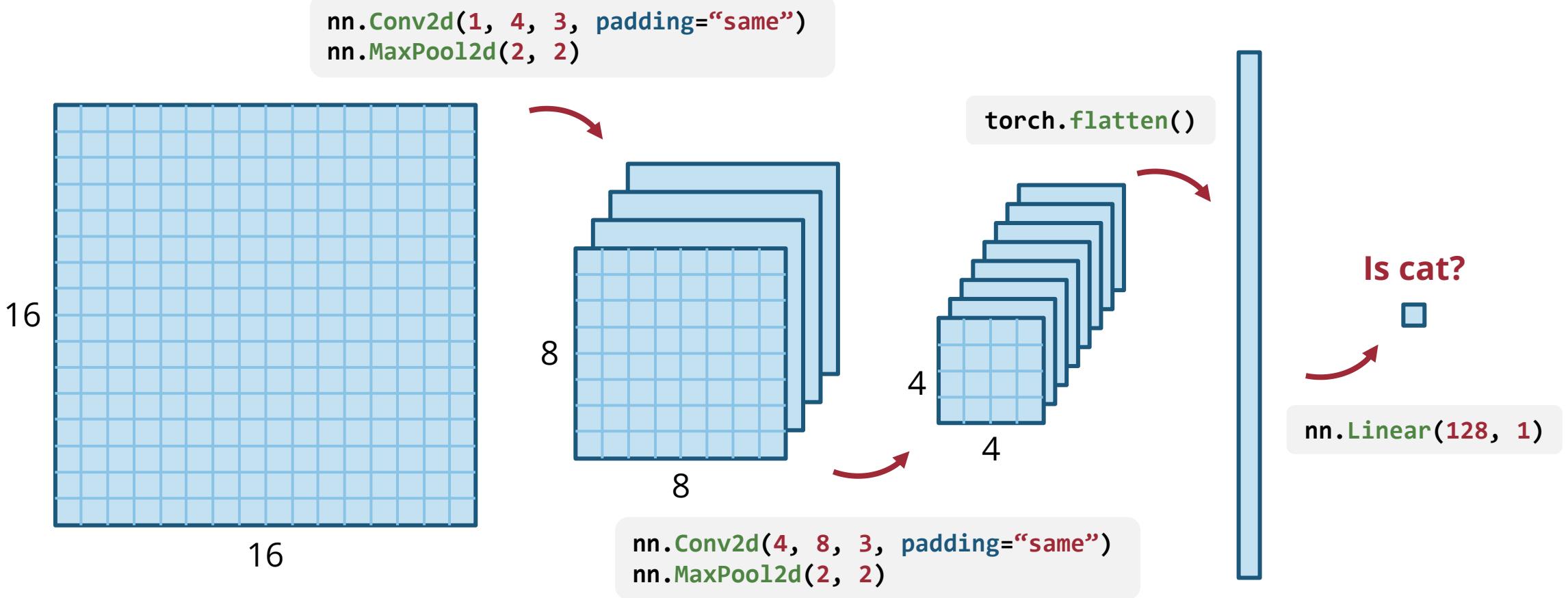


```
nn.MaxPool2d(kernel_size=2, stride=2)
```

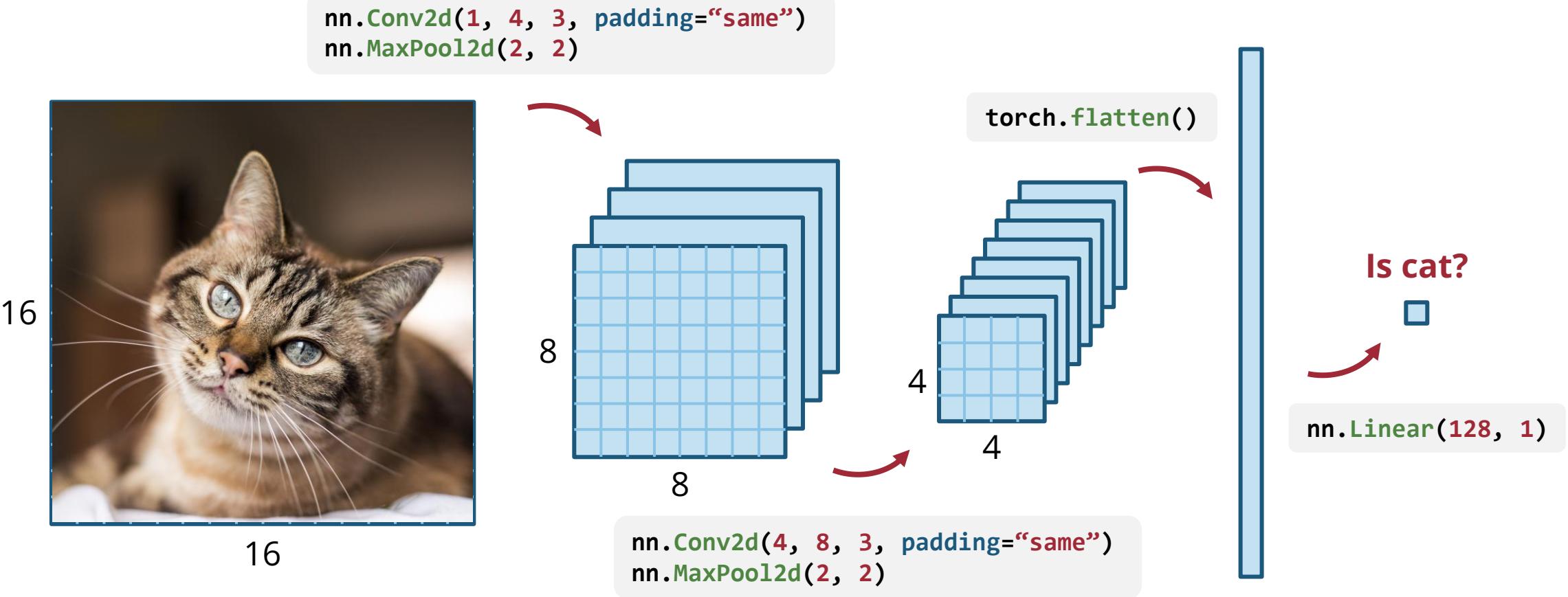
A Toy Example



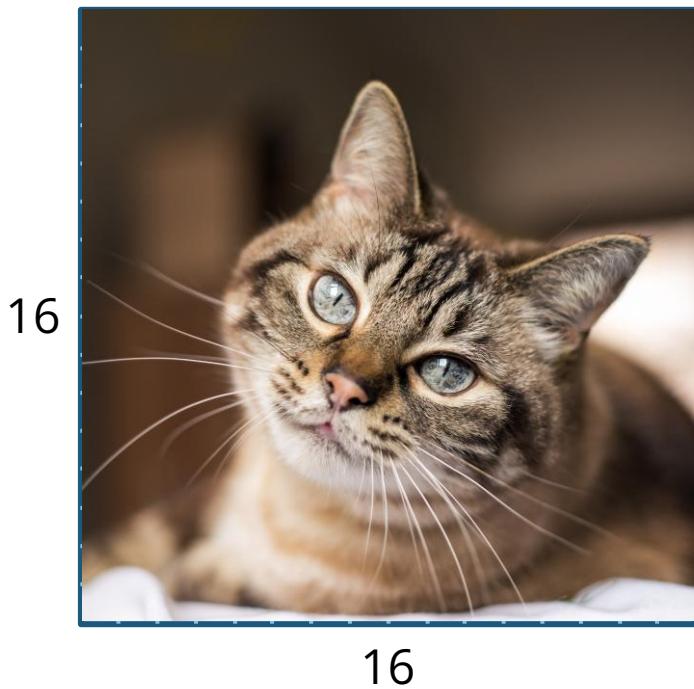
A Toy Example



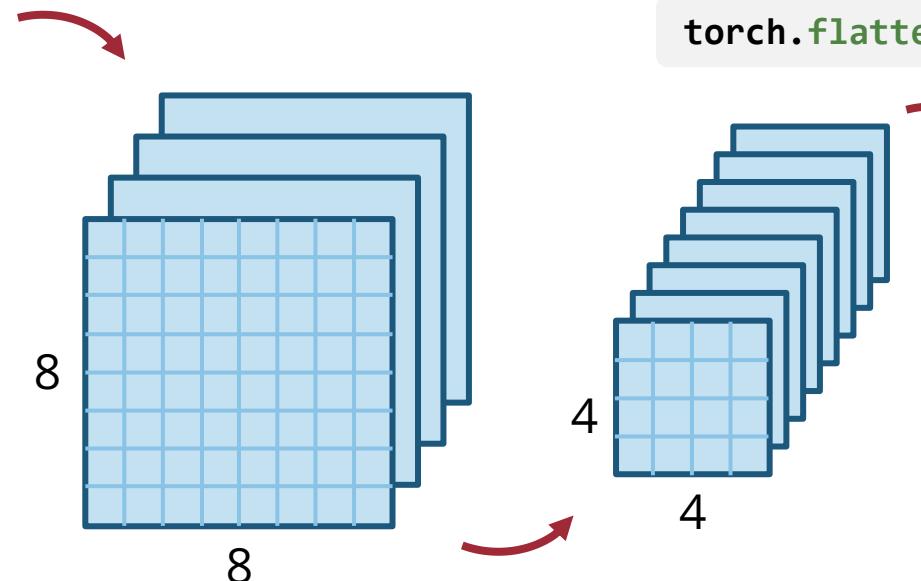
A Toy Example



A Toy Example

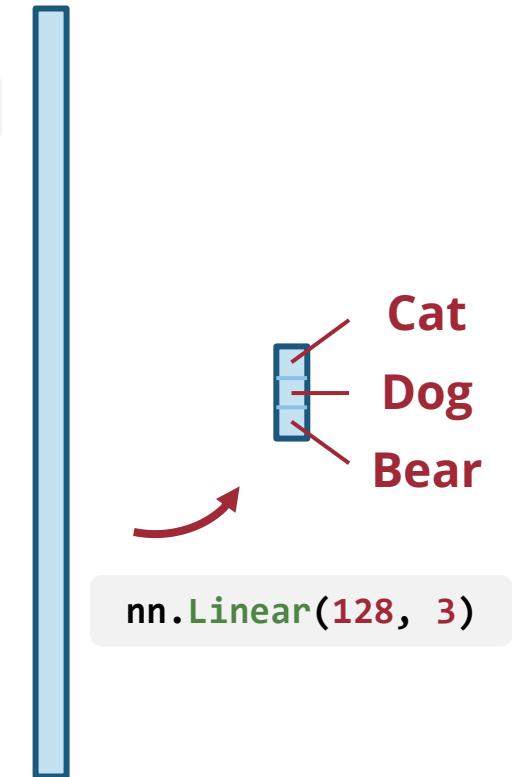


```
nn.Conv2d(1, 4, 3, padding="same")  
nn.MaxPool2d(2, 2)
```



```
nn.Conv2d(4, 8, 3, padding="same")  
nn.MaxPool2d(2, 2)
```

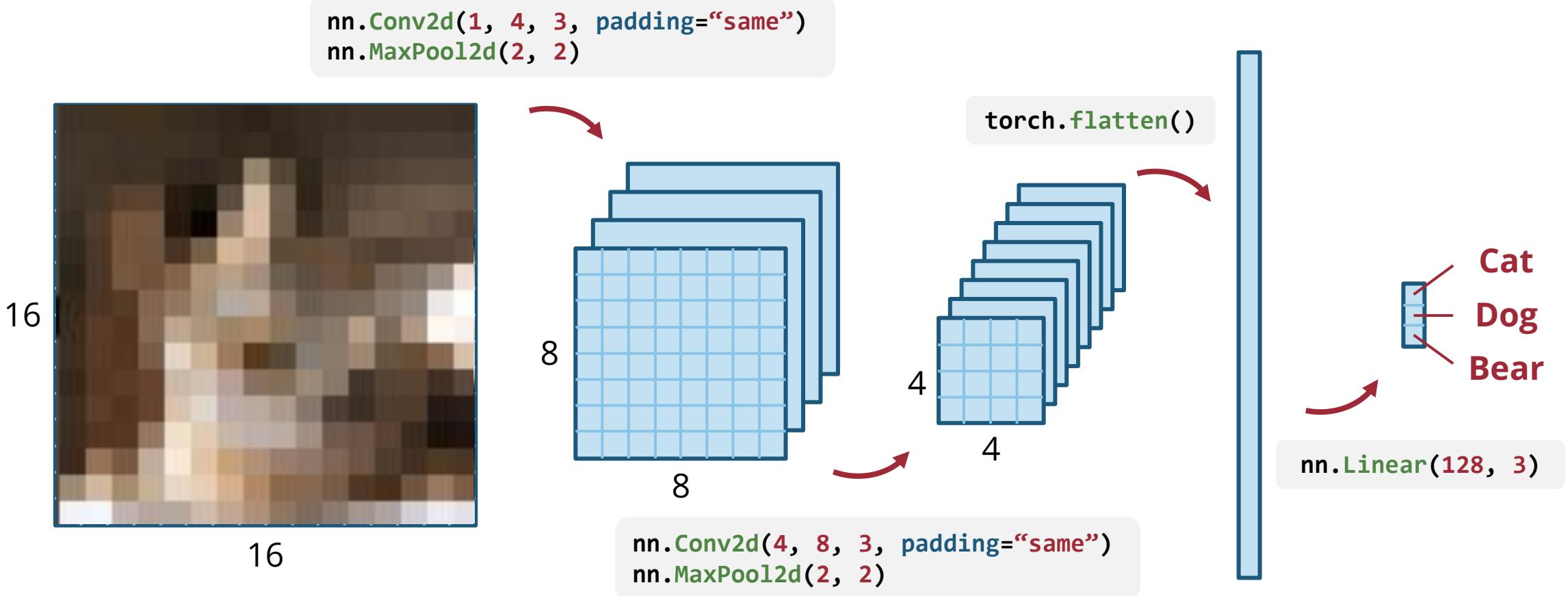
```
torch.flatten()
```



```
nn.Linear(128, 3)
```

Cat
Dog
Bear

A Toy Example



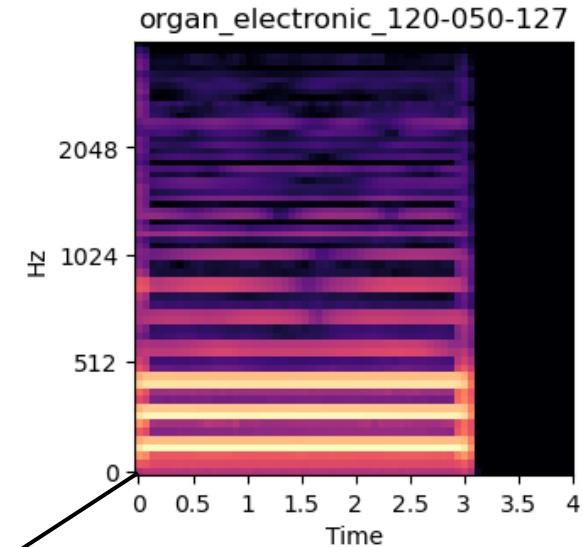
A Real Example

```
class CNN(nn.Module):    Input      Output  
    channels    channels  
    """A basic convolutional neural network."""  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(1, 16, 3, padding="same")  
        self.conv2 = nn.Conv2d(16, 32, 3, padding="same")  
        self.conv3 = nn.Conv2d(32, 64, 3, padding="same")  
        self.conv4 = nn.Conv2d(64, 128, 3, padding="same")  
        self.pool = nn.MaxPool2d(2, 2)  
        self.fc = nn.Linear(128 * 4 * 4, n_classes)  
  
    def forward(self, x):  
        x = self.pool(F.relu(self.conv1(x)))  
        x = self.pool(F.relu(self.conv2(x)))  
        x = self.pool(F.relu(self.conv3(x)))  
        x = self.pool(F.relu(self.conv4(x)))  
        x = torch.flatten(x, 1)  
        x = self.fc(x)  
        return x
```

Kernel size

Input:	1 x 64 x 64	= 16384
→	16 x 32 x 32	= 8192
→	32 x 16 x 16	= 4096
→	64 x 8 x 8	= 2048
→	128 x 4 x 4	

More channels,
lower resolution
as we go deeper



Total number of
features decrease
as we go deeper

A Real Example

Input channels	Output channels	Kernel size
1	16	3

class CNN(nn.Module):

```
    """A basic convolutional neural network."""
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(1, 16, 3, padding="same")
        self.conv2 = nn.Conv2d(16, 32, 3, padding="same")
        self.conv3 = nn.Conv2d(32, 64, 3, padding="same")
        self.conv4 = nn.Conv2d(64, 128, 3, padding="same")
        self.pool = nn.MaxPool2d(2, 2)
        self.fc = nn.Linear(128 * 4 * 4, n_classes)
```

def forward(self, x):

```
    x = self.pool(F.relu(self.conv1(x))) →  $(3 \times 3 \times 1 + 1) \times 16 = 160$ 
    x = self.pool(F.relu(self.conv2(x))) →  $(3 \times 3 \times 16 + 1) \times 32 = 4640$ 
    x = self.pool(F.relu(self.conv3(x))) →  $(3 \times 3 \times 32 + 1) \times 64 = 18496$ 
    x = self.pool(F.relu(self.conv4(x))) →  $(3 \times 3 \times 64 + 1) \times 128 = 73856$ 
    x = torch.flatten(x, 1) →  $(2048 + 1) \times 11 = 22539$ 
    return x
```

How many parameters do we have in each layer?

| Benefits of CNNs

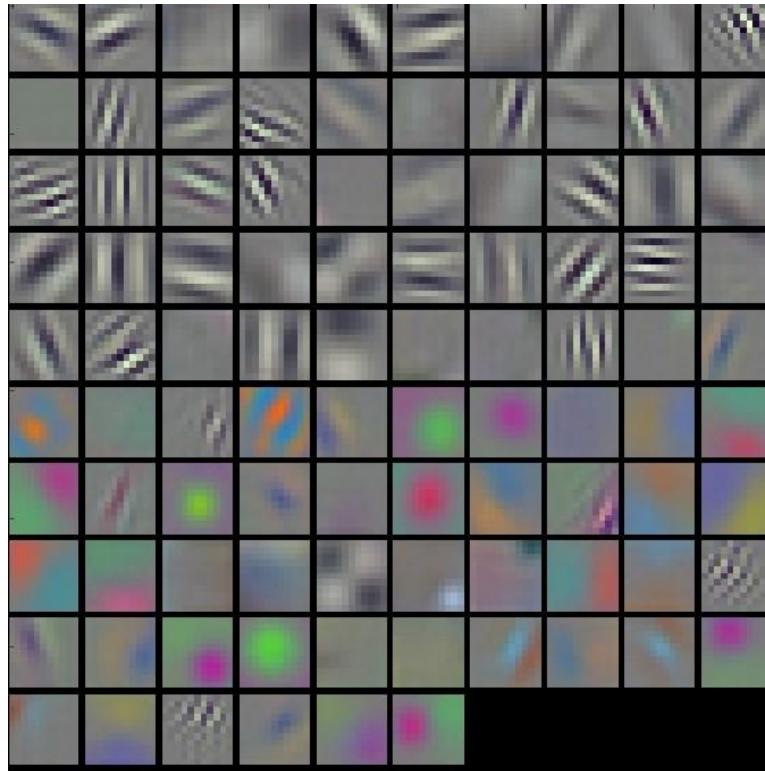
- Learn **local patterns**
- **Invariant to shifts**
 - Also called *translational invariance*
- **Reuse the learned filters** across
 - Different parts of the image
 - Across different images
- **Higher parameter-efficiency** against fully-connected neural network

What does a CNN Learn?

Learned CNN Kernels in a Trained AlexNet

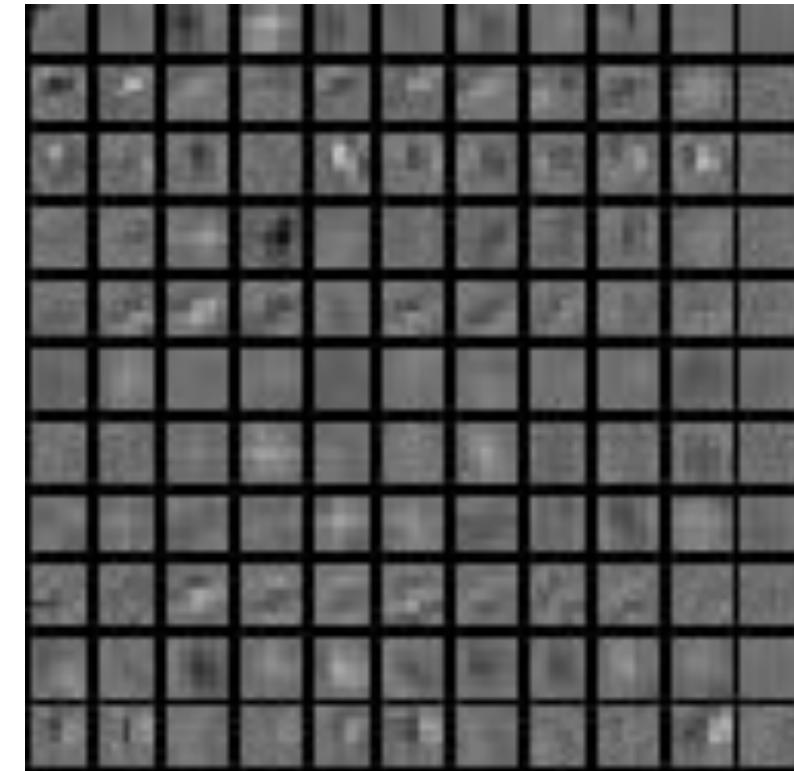
1st convolutional layer

11x11
kernels

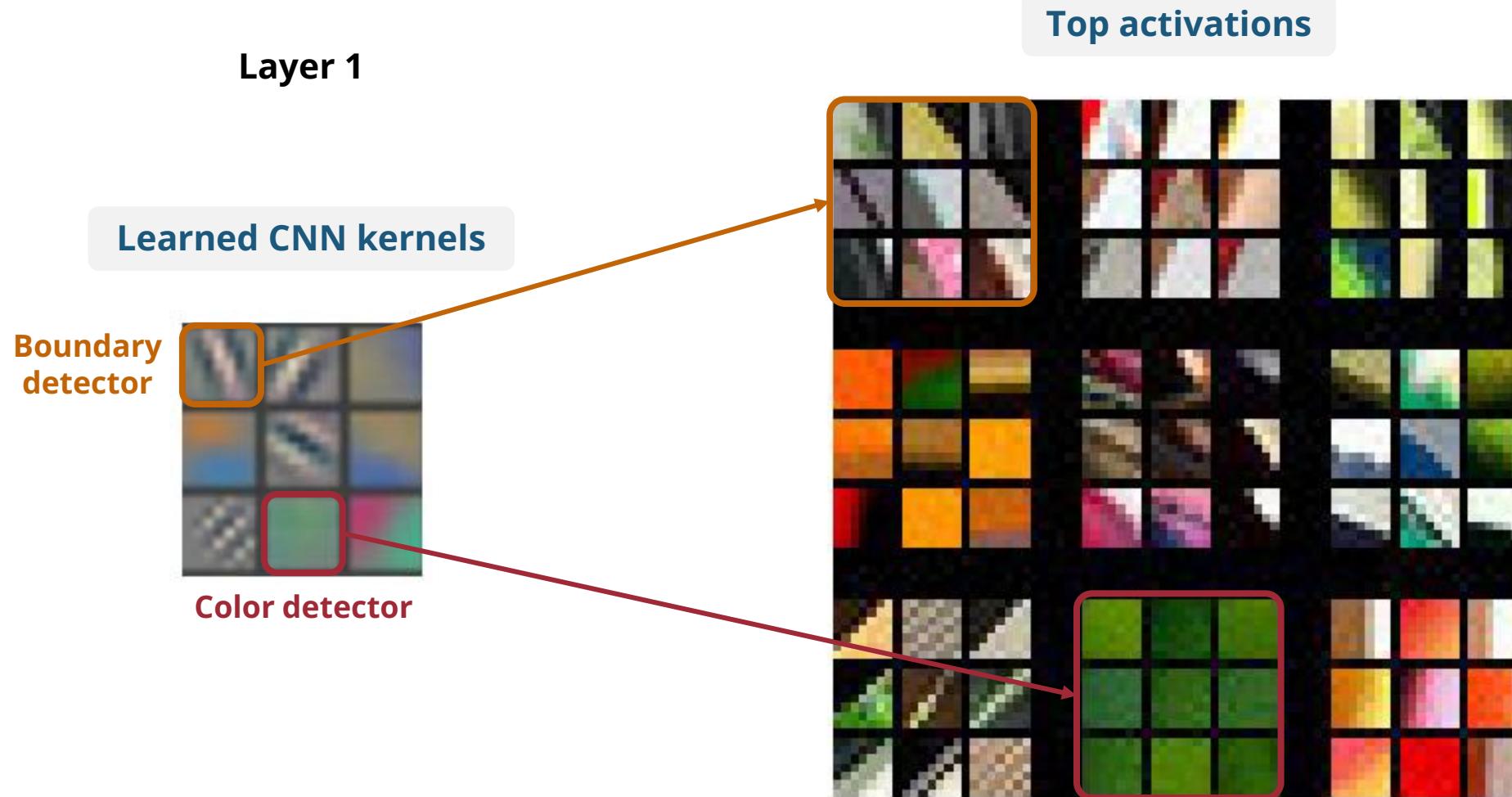


2nd convolutional layer

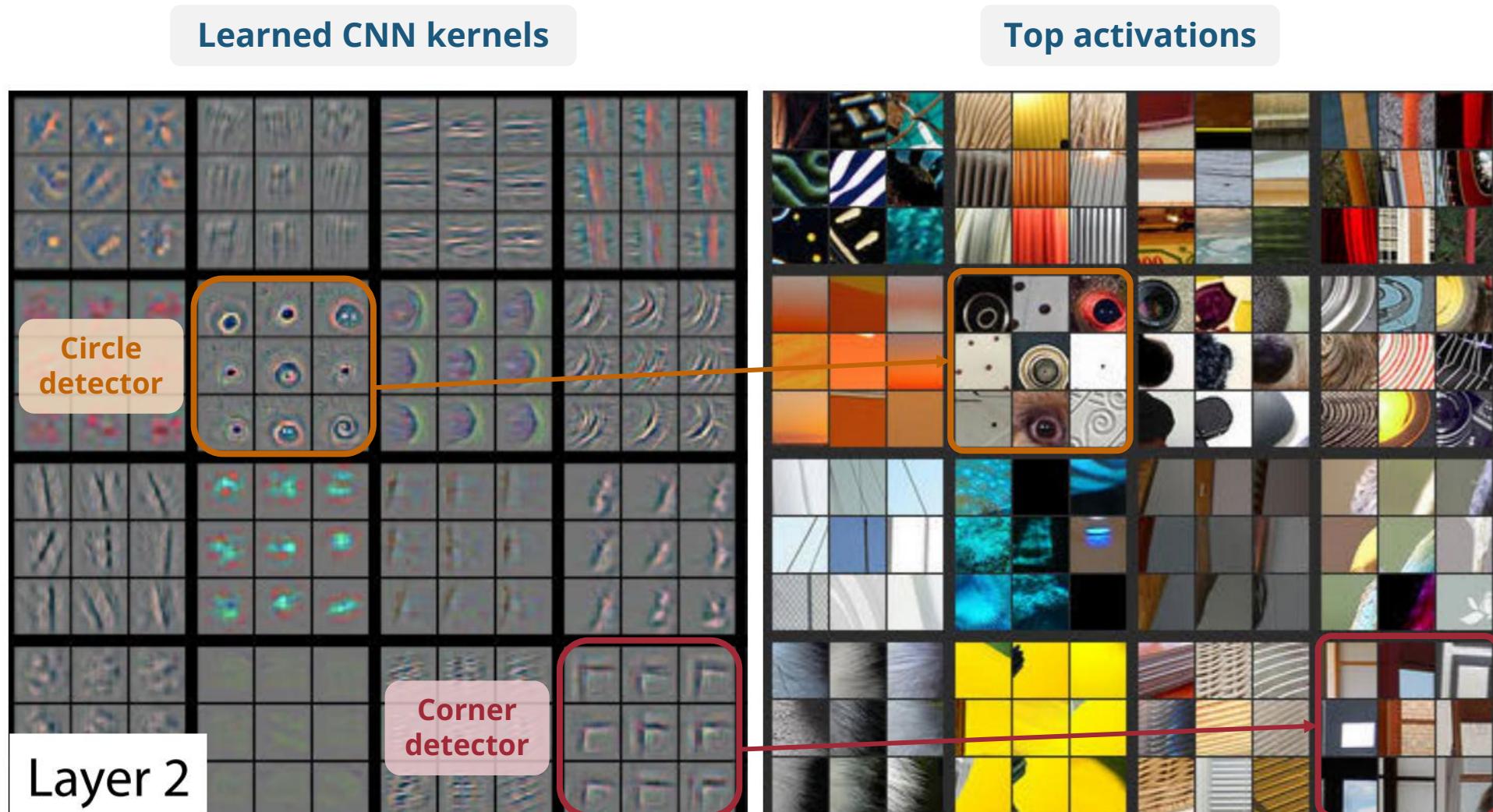
5x5
kernels



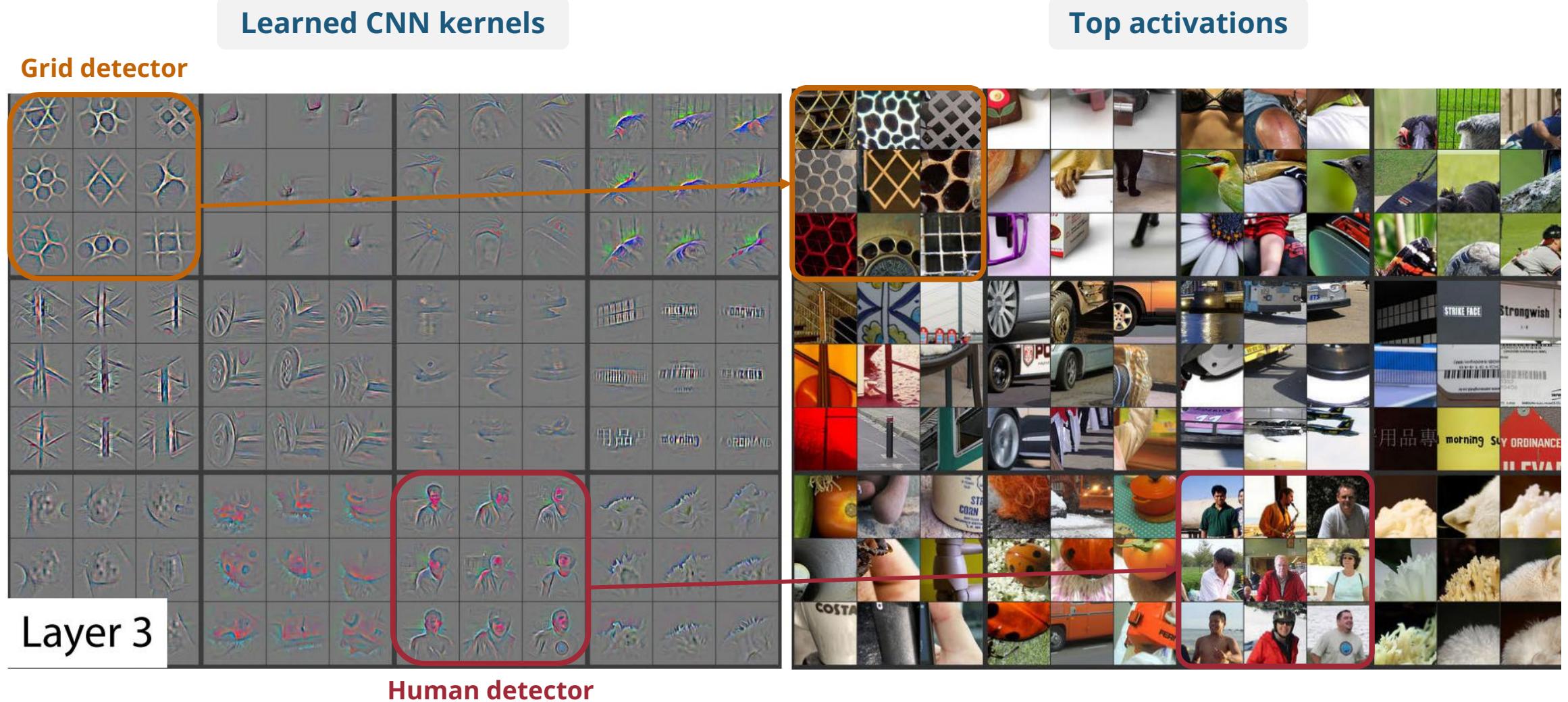
| Learned CNN Kernels in a Trained AlexNet



Learned CNN Kernels in a Trained AlexNet

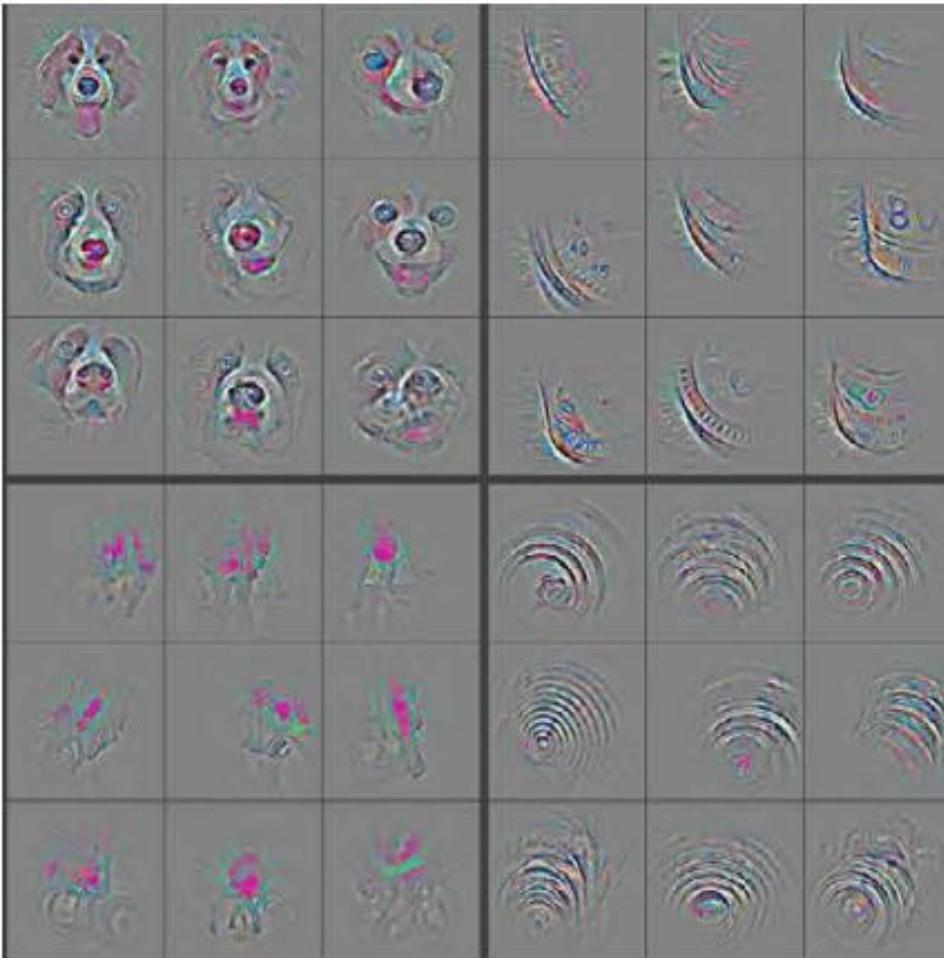


Learned CNN Kernels in a Trained AlexNet

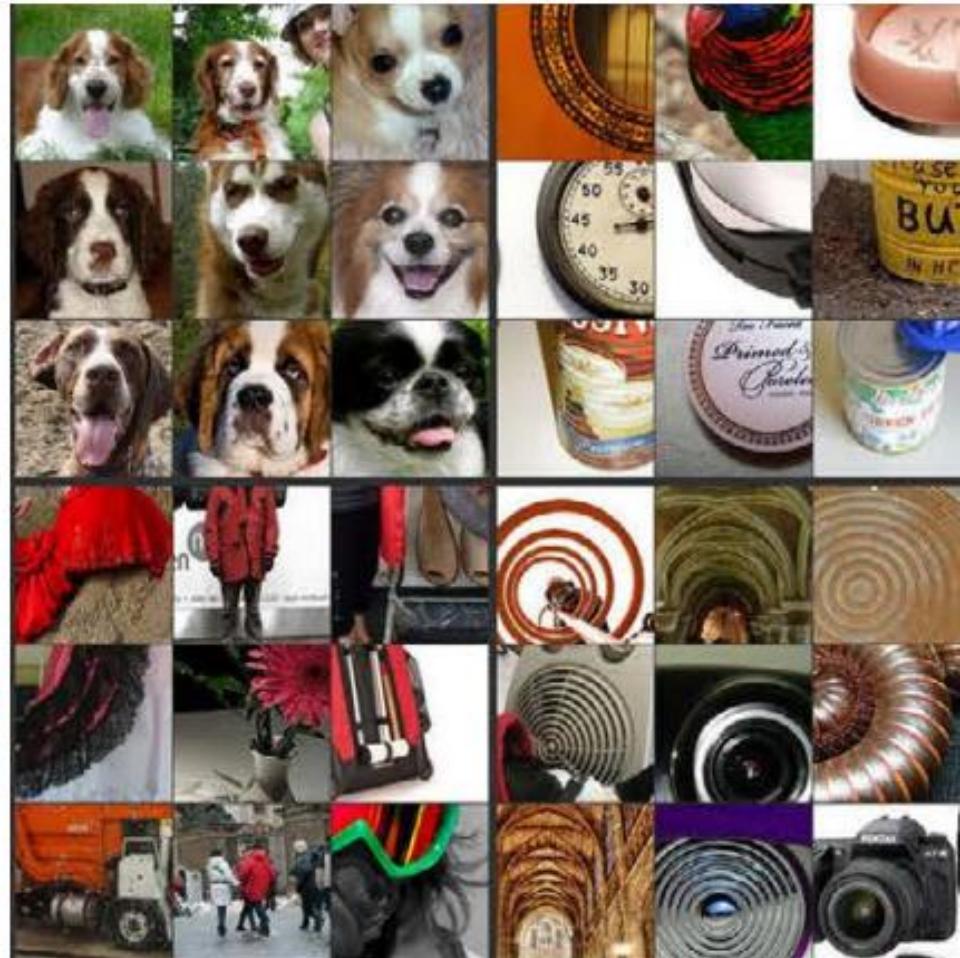


Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

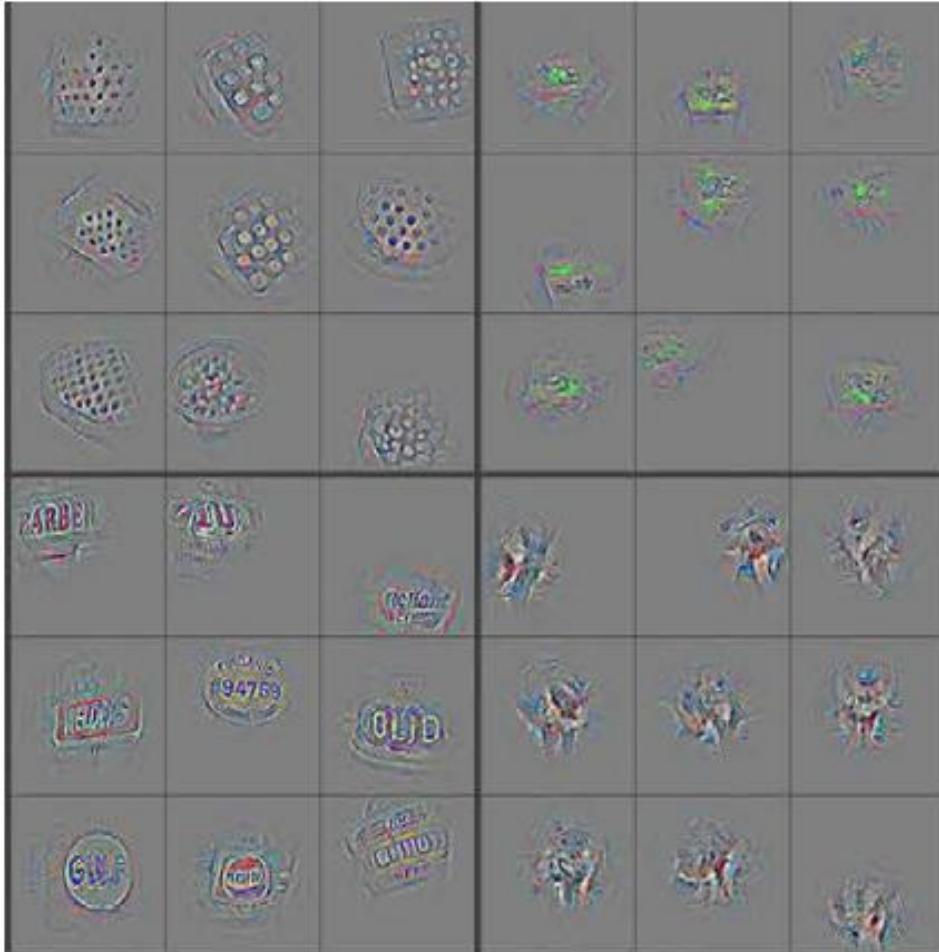


Top activations



Learned CNN Kernels in a Trained AlexNet

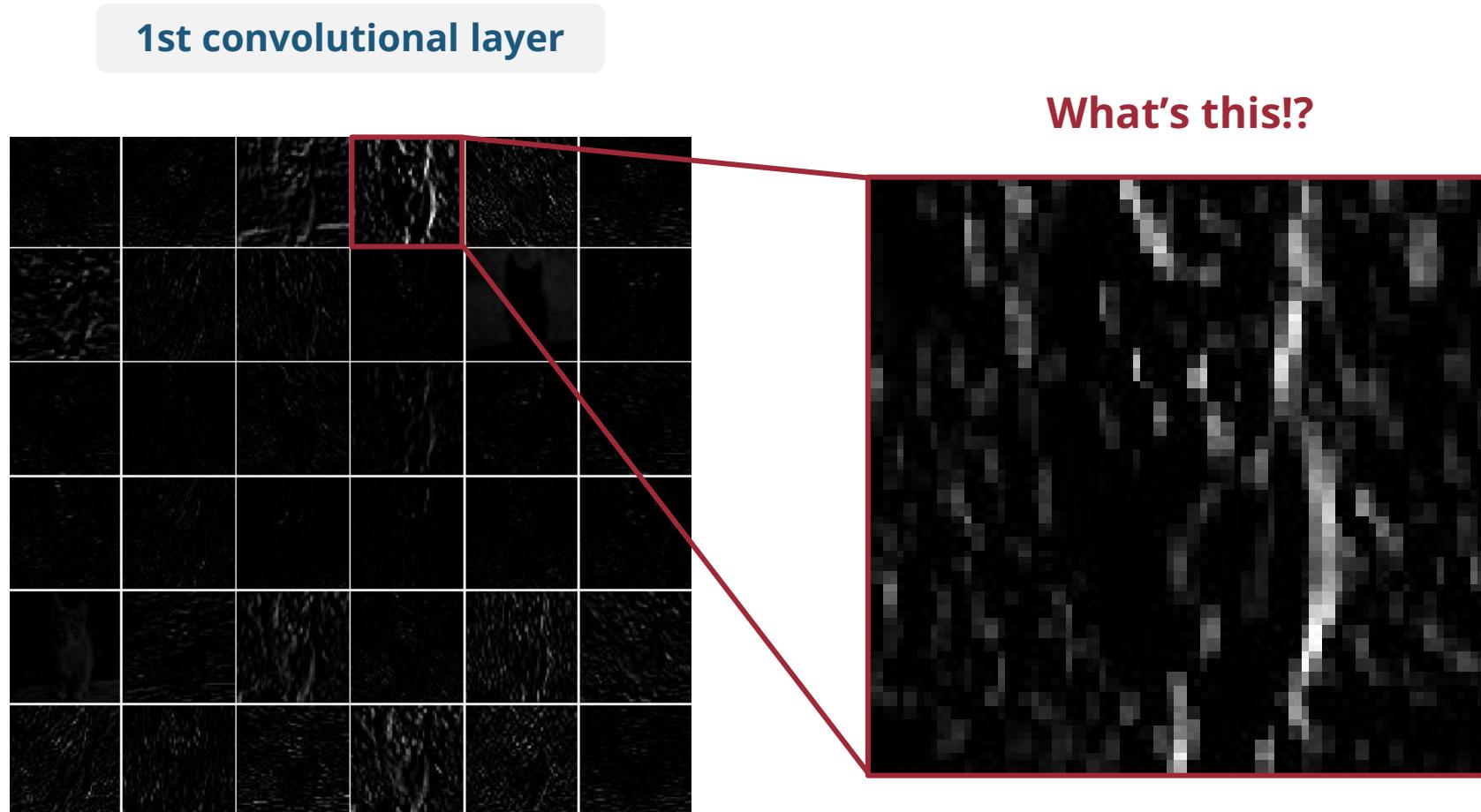
Learned CNN kernels



Top activations

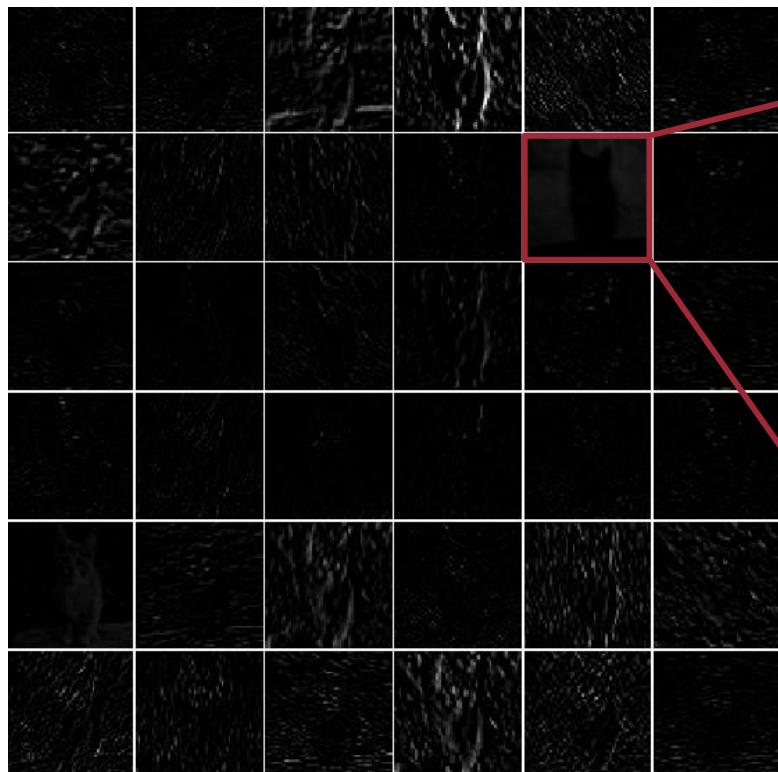


Activations in a Trained AlexNet

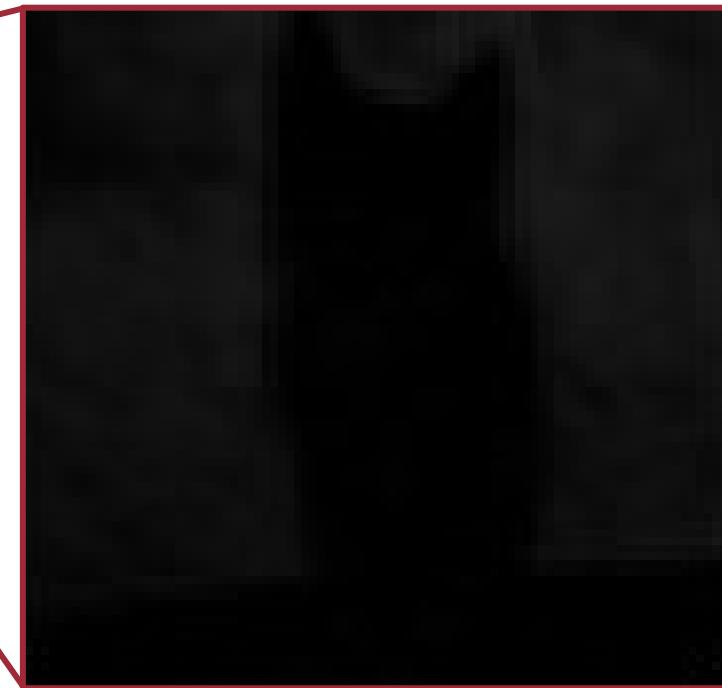


Activations in a Trained AlexNet

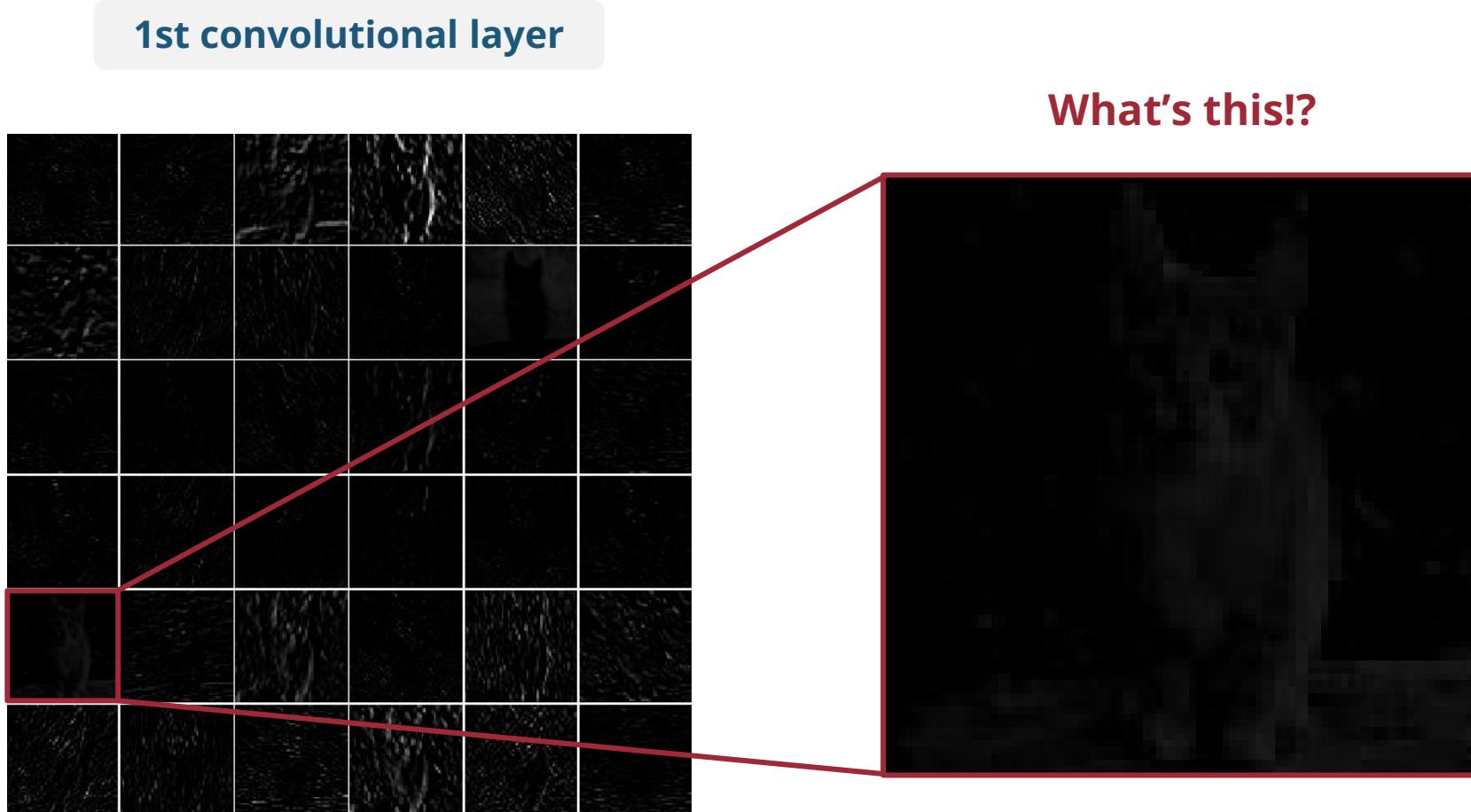
1st convolutional layer



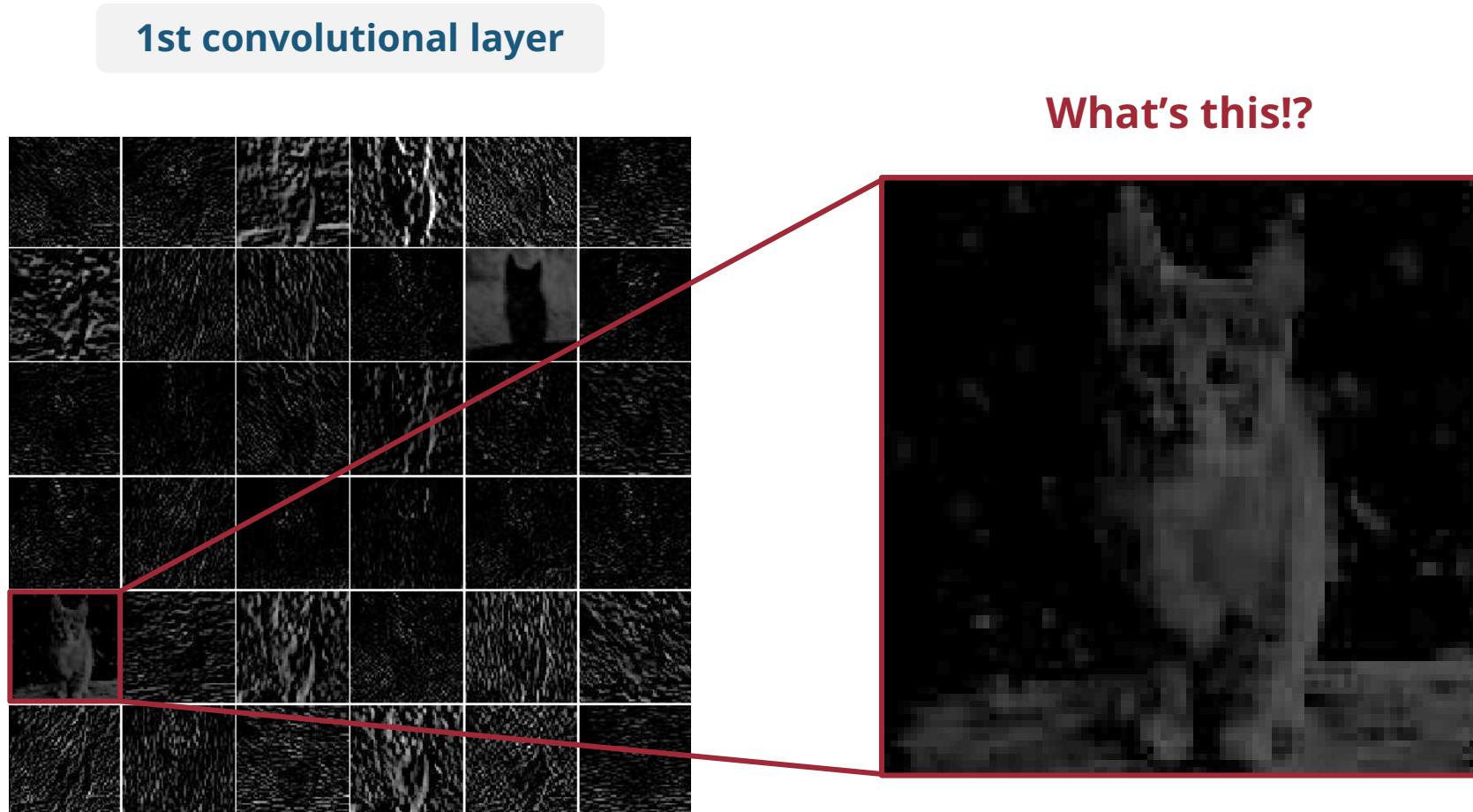
What's this!?



Activations in a Trained AlexNet

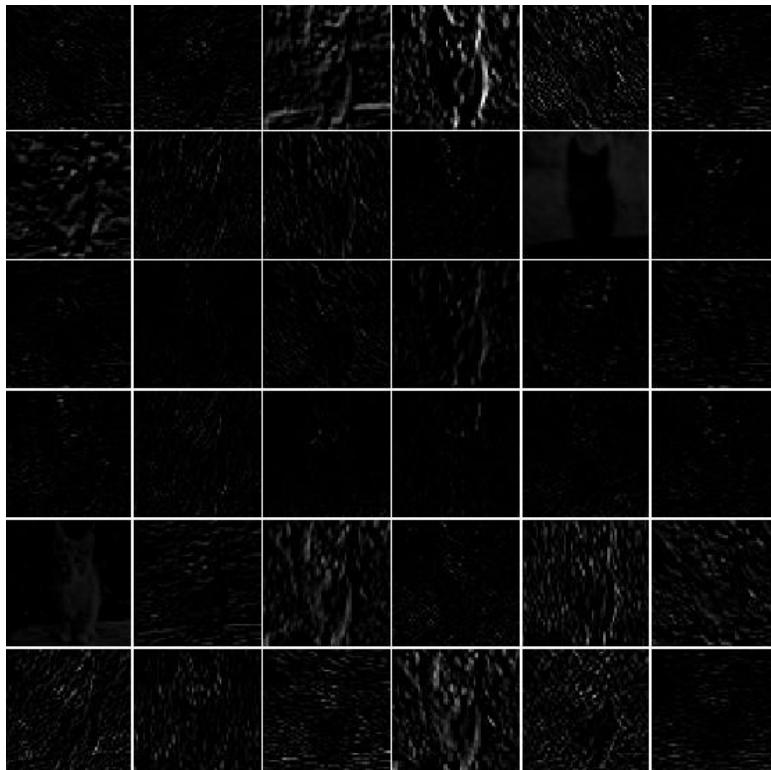


Activations in a Trained AlexNet

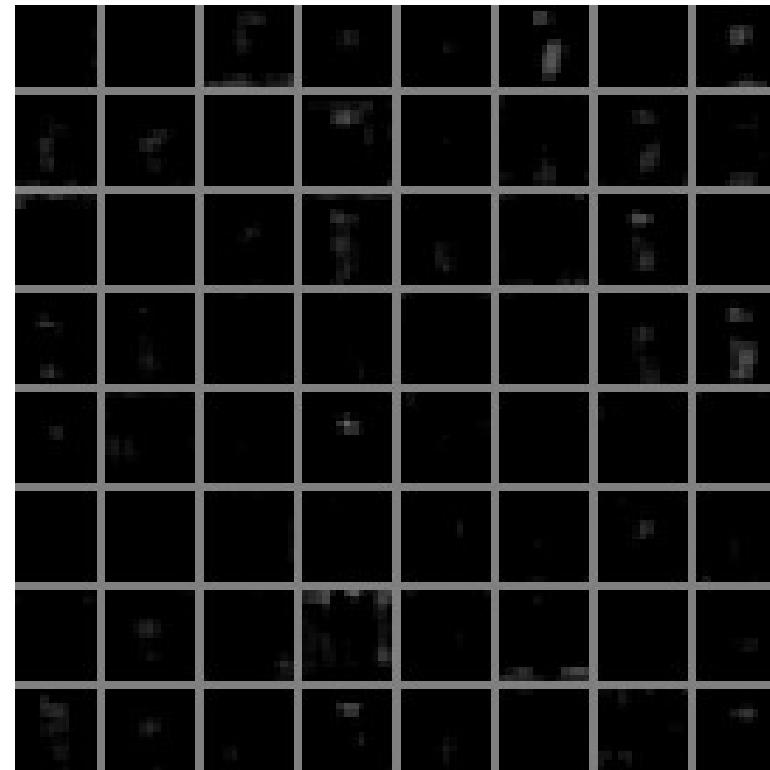


Activations in a Trained AlexNet

1st convolutional layer



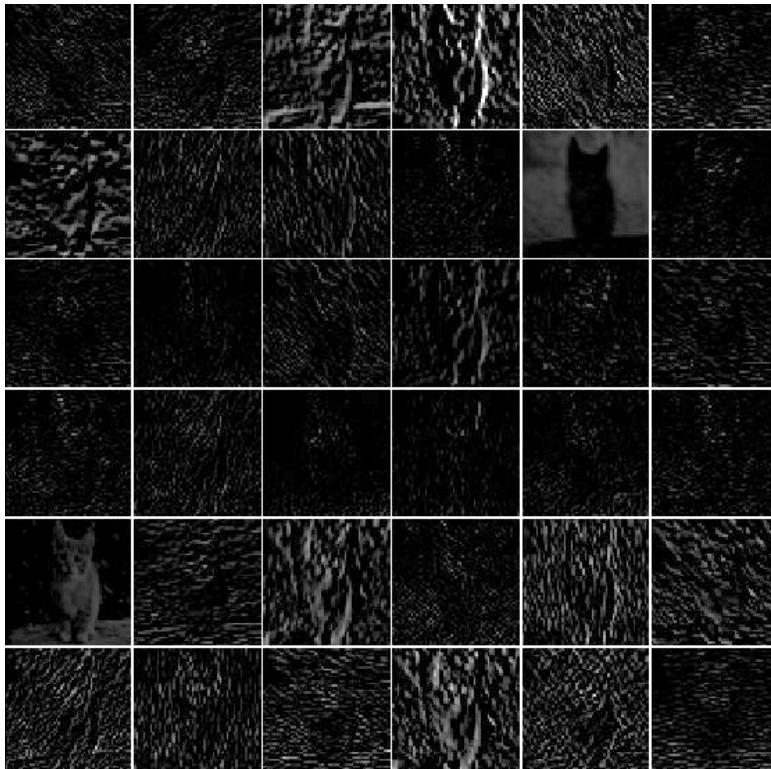
5th convolutional layer



Activations in a Trained AlexNet

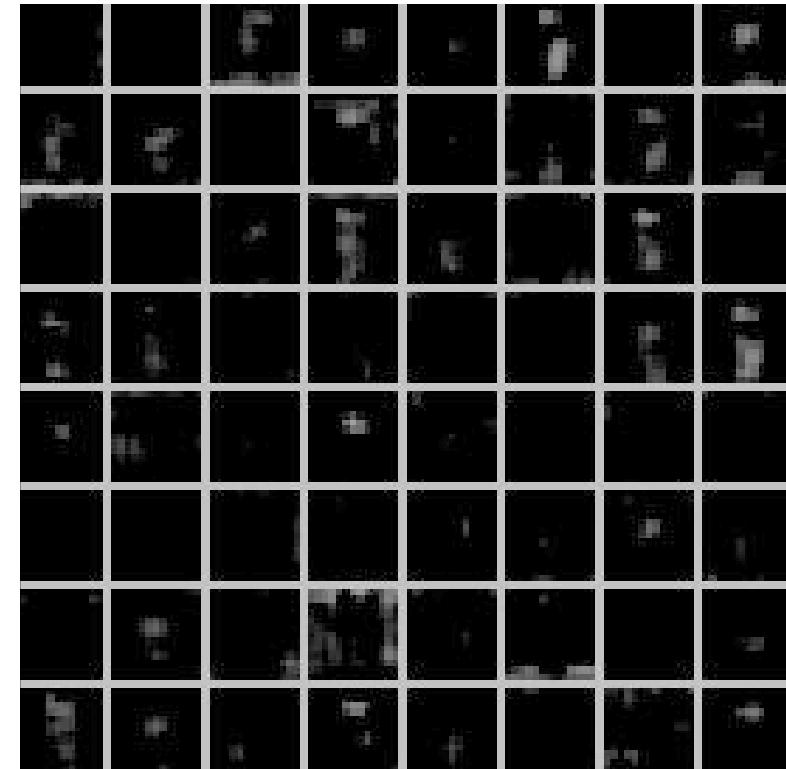
1st convolutional layer

Brightness+
Contrast+

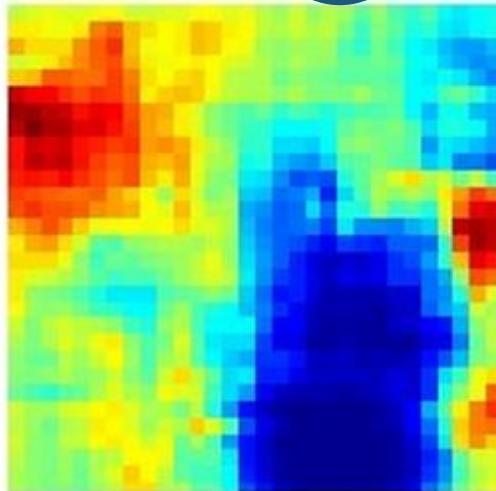
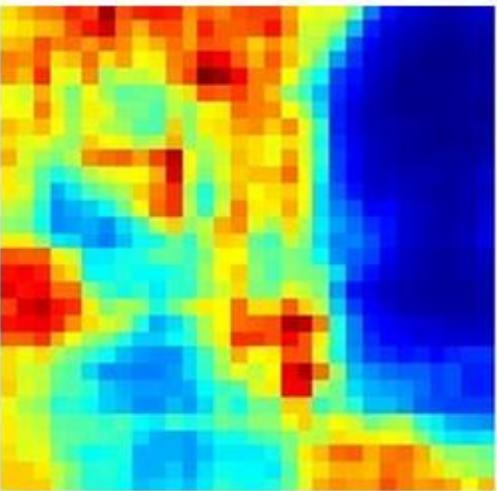
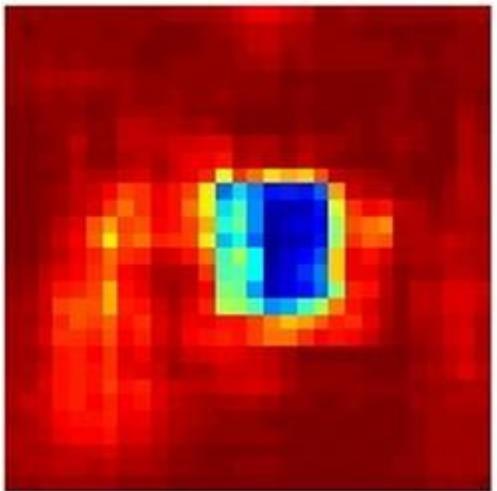
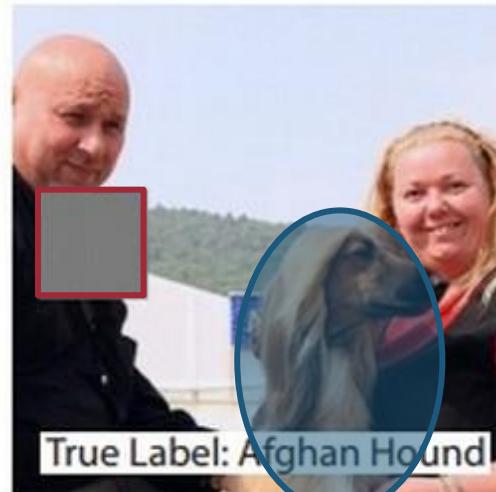


5th convolutional layer

Brightness+
Contrast+

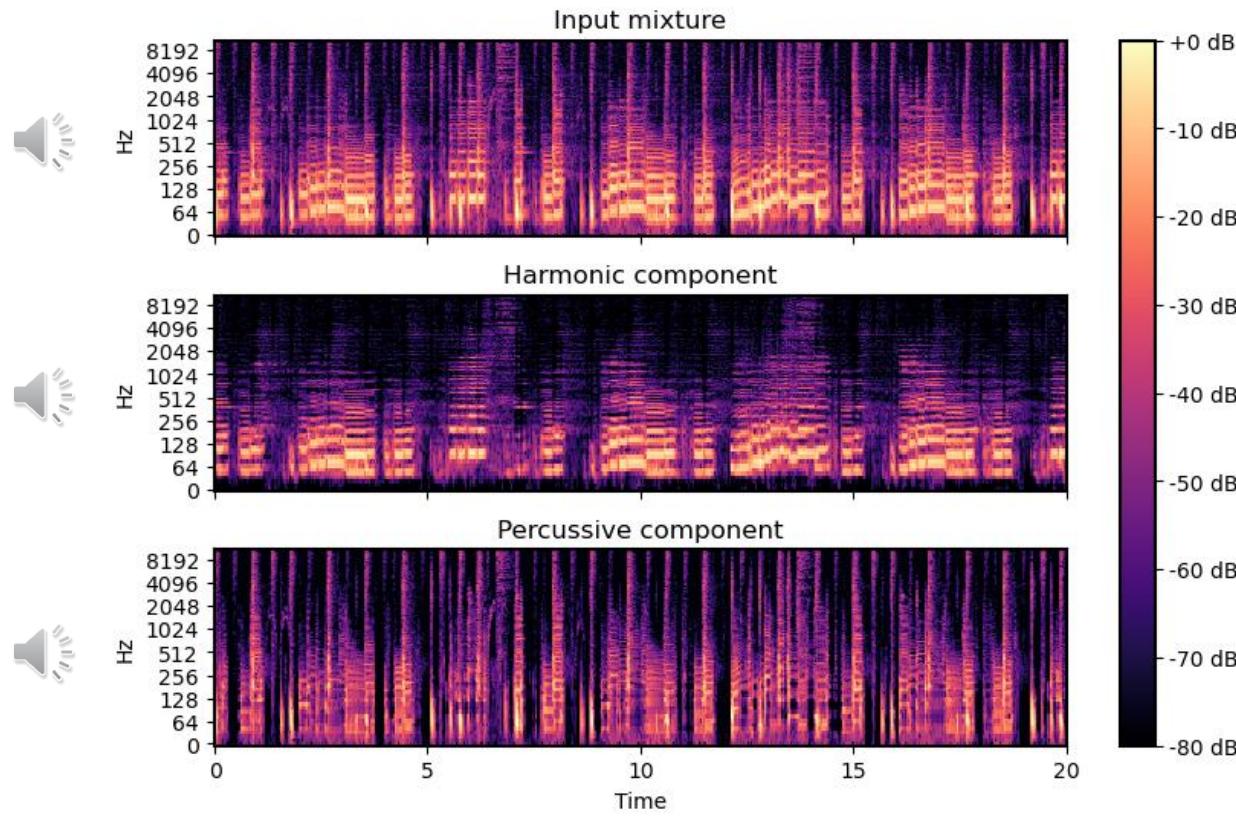


| What does a CNN Learn?



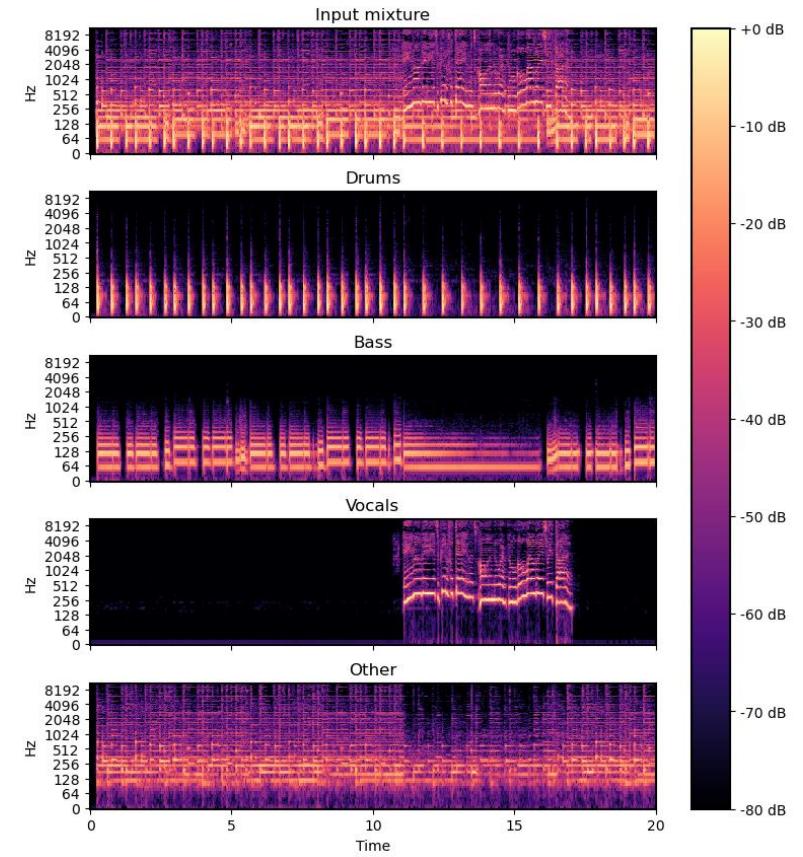
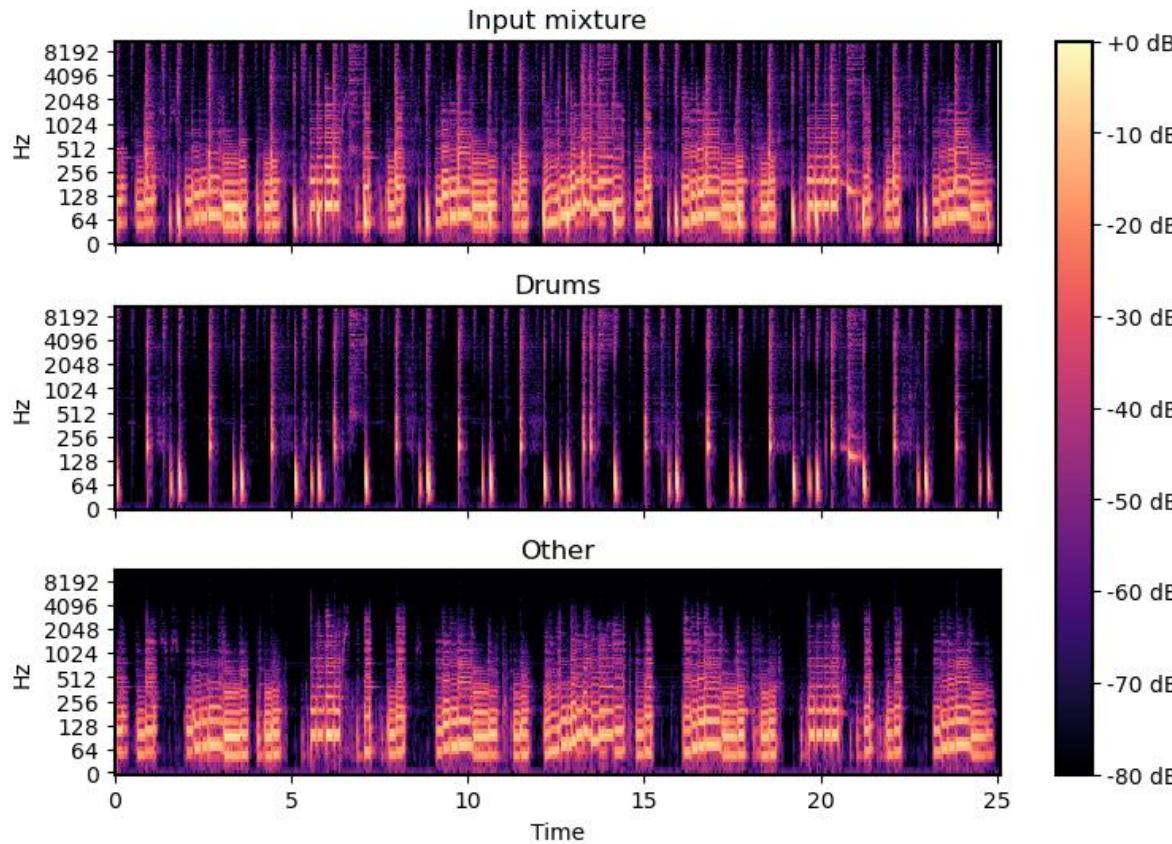
🔥 PA 3: Source Separation

- Instructions will be sent by **emails** and released on the **course website**
- **Part 1:** Harmonic-Percussive Source Separation (HPSS) using **librosa**



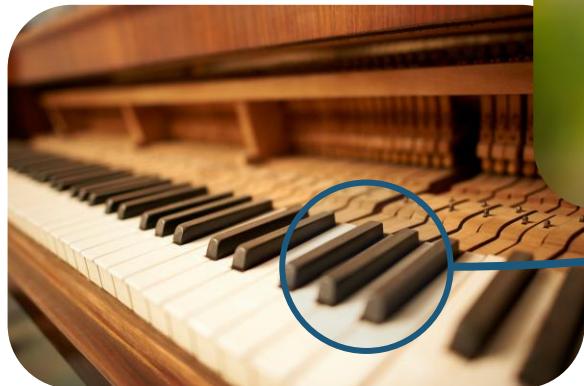
🔥 PA 3: Source Separation

- Part 2: Music Source Separation using **Demucs**

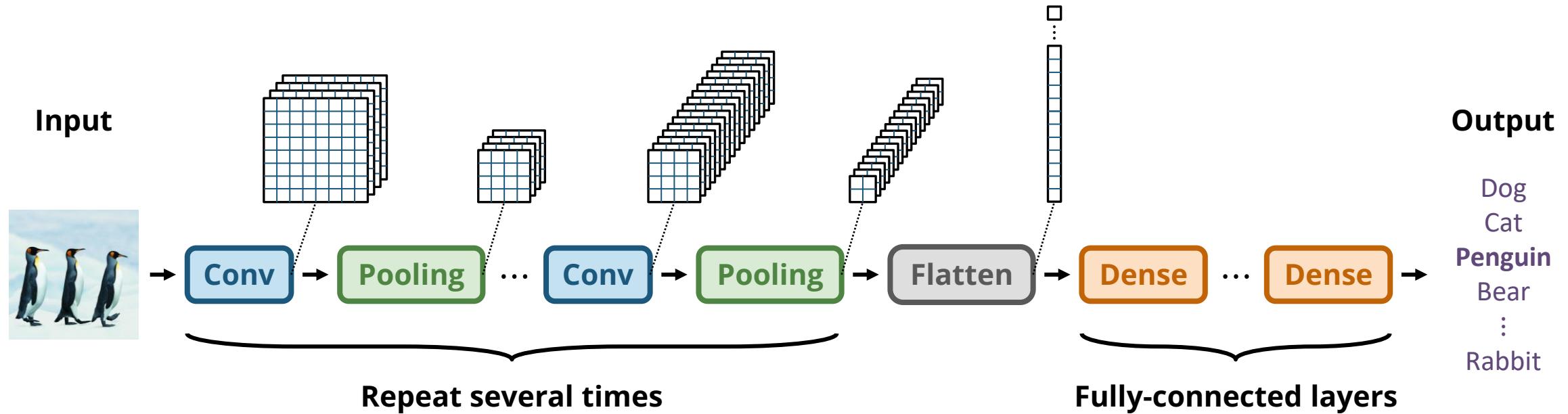


Recap

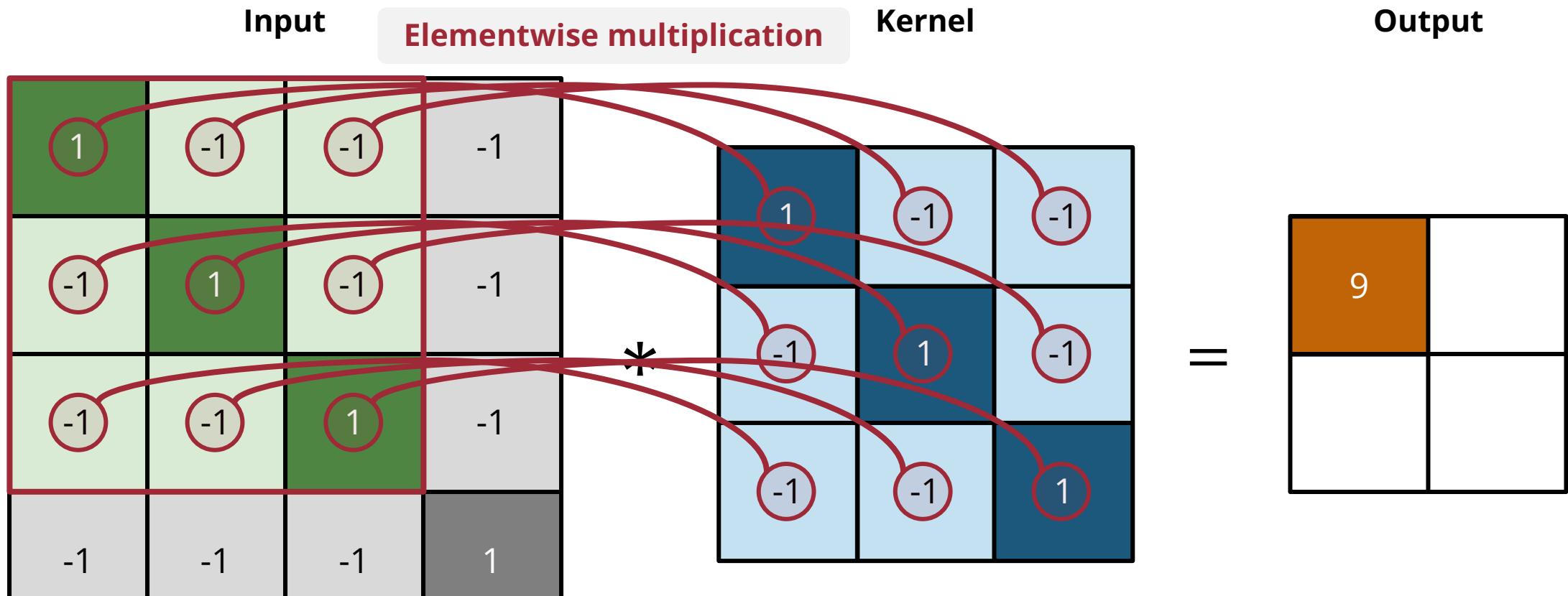
| Reusable Pattern Detectors



Convolutional Neural Network (CNNs)



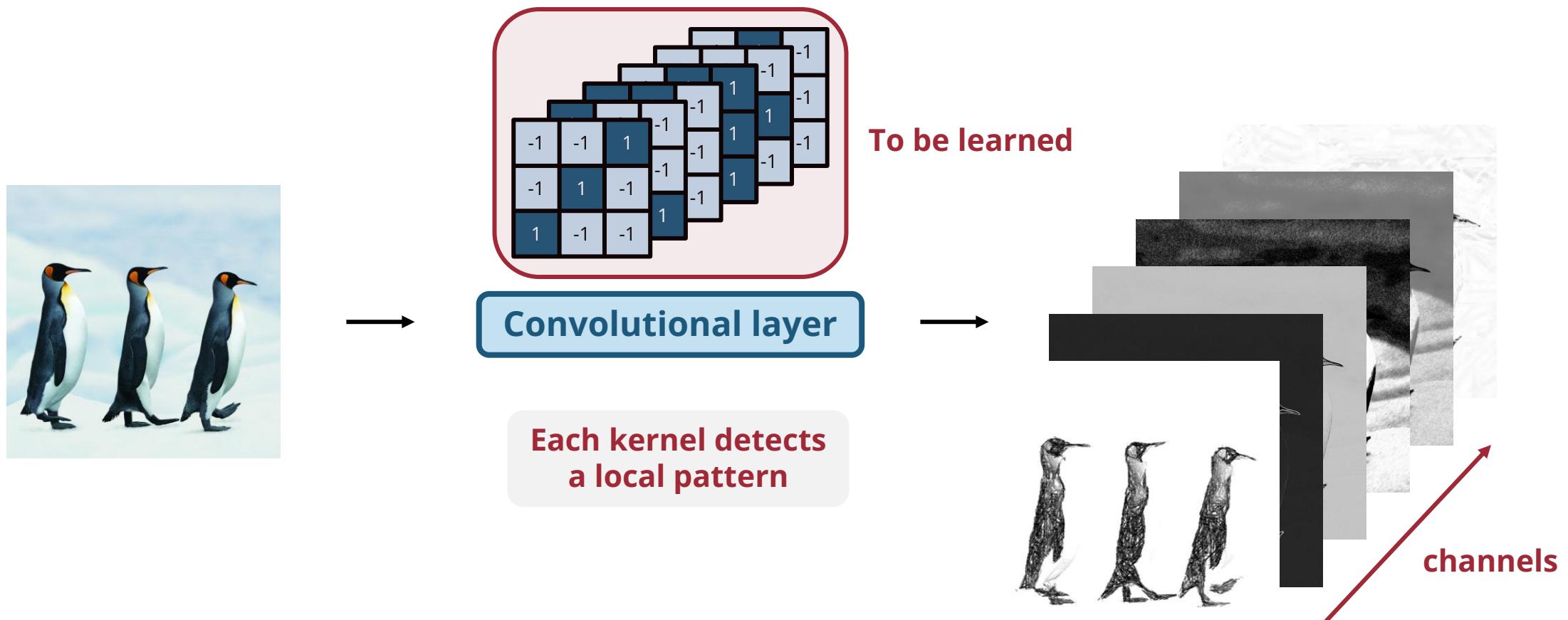
2D Convolution



$$\begin{aligned} (1 \times 1) &+ (-1 \times -1) + (-1 \times -1) \\ + (-1 \times -1) &+ (1 \times 1) + (-1 \times -1) = 9 \\ + (-1 \times -1) &+ (-1 \times -1) + (1 \times 1) \end{aligned}$$

Convolutional Layer

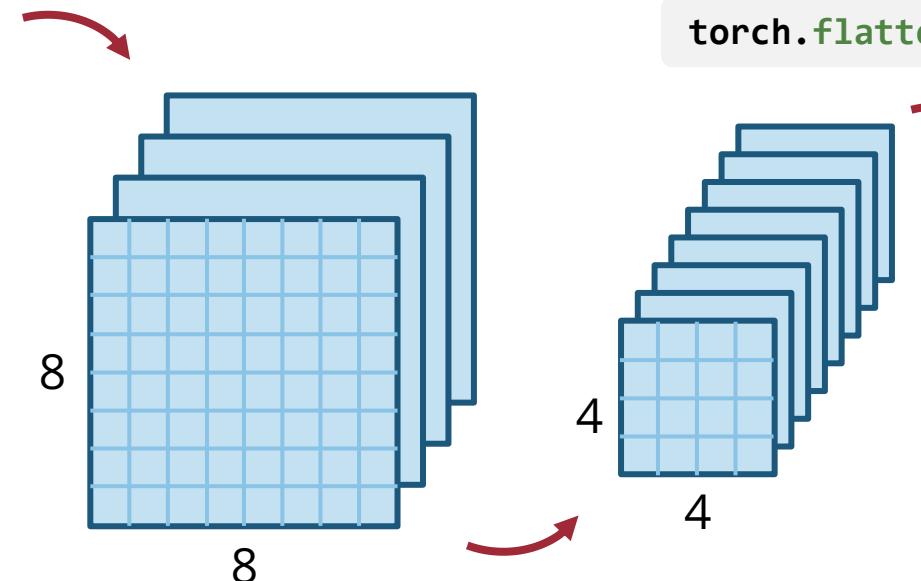
- A convolutional layer consists of many **learnable kernels** (channels)



A Toy Example

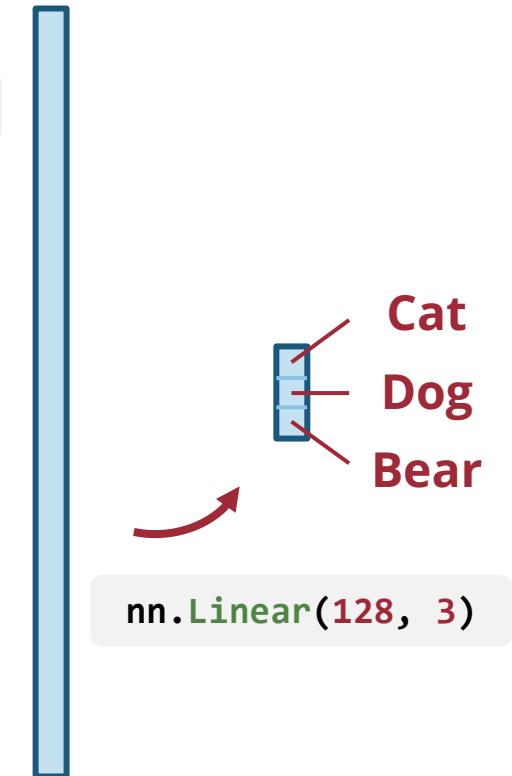


```
nn.Conv2d(1, 4, 3, padding="same")  
nn.MaxPool2d(2, 2)
```



```
nn.Conv2d(4, 8, 3, padding="same")  
nn.MaxPool2d(2, 2)
```

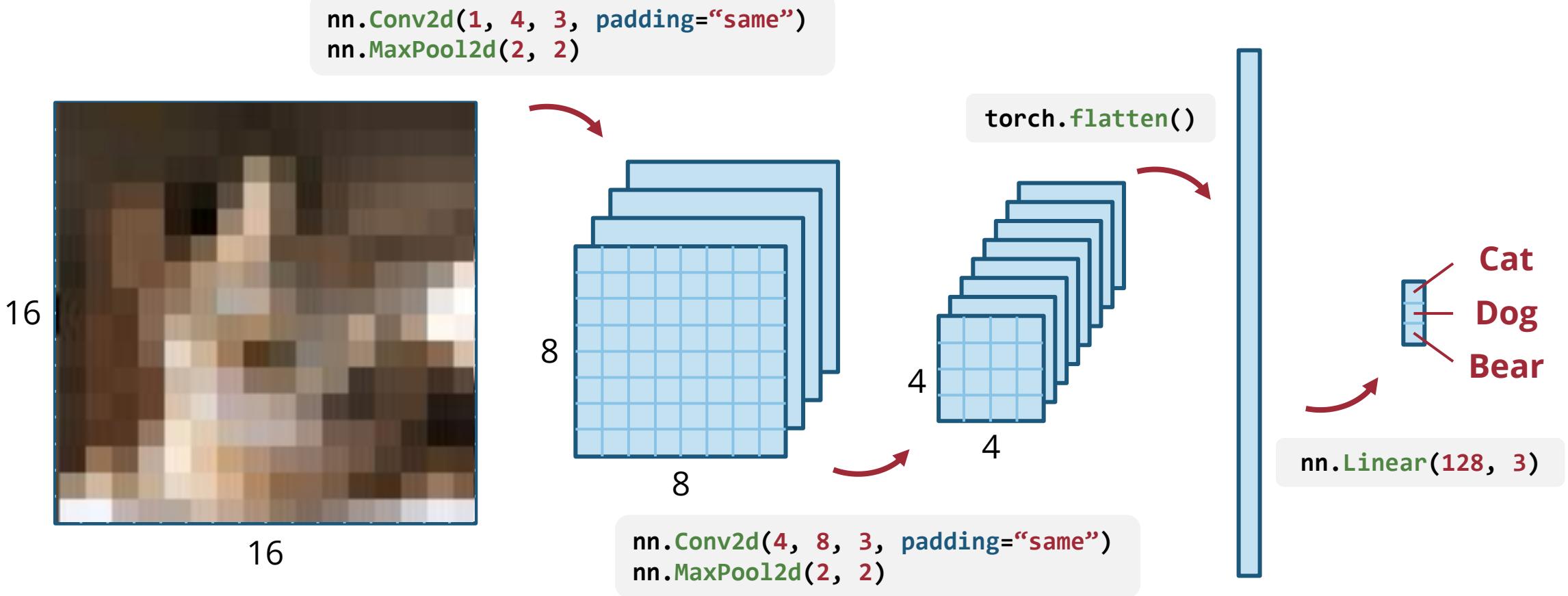
```
torch.flatten()
```



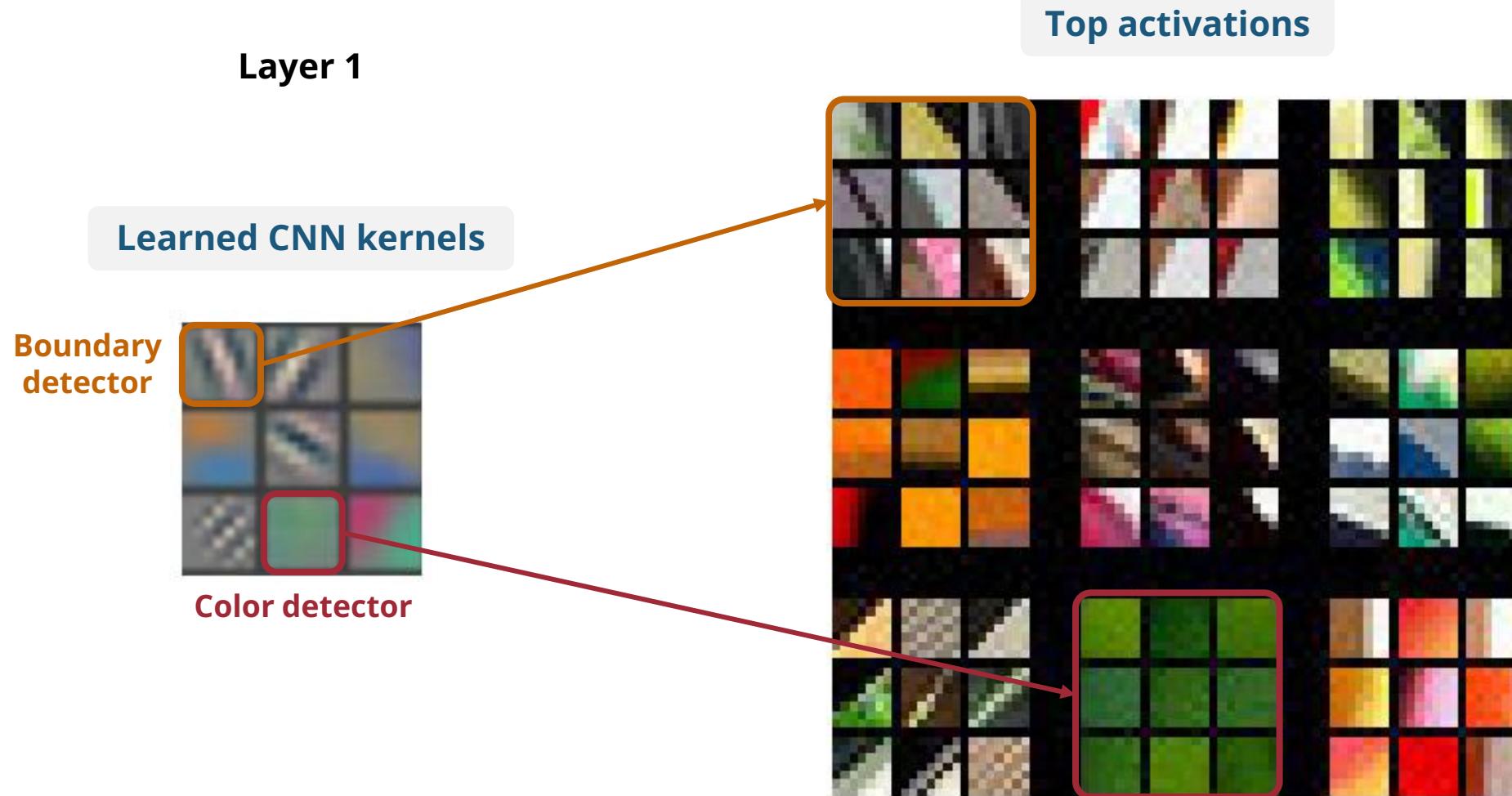
```
nn.Linear(128, 3)
```

Cat
Dog
Bear

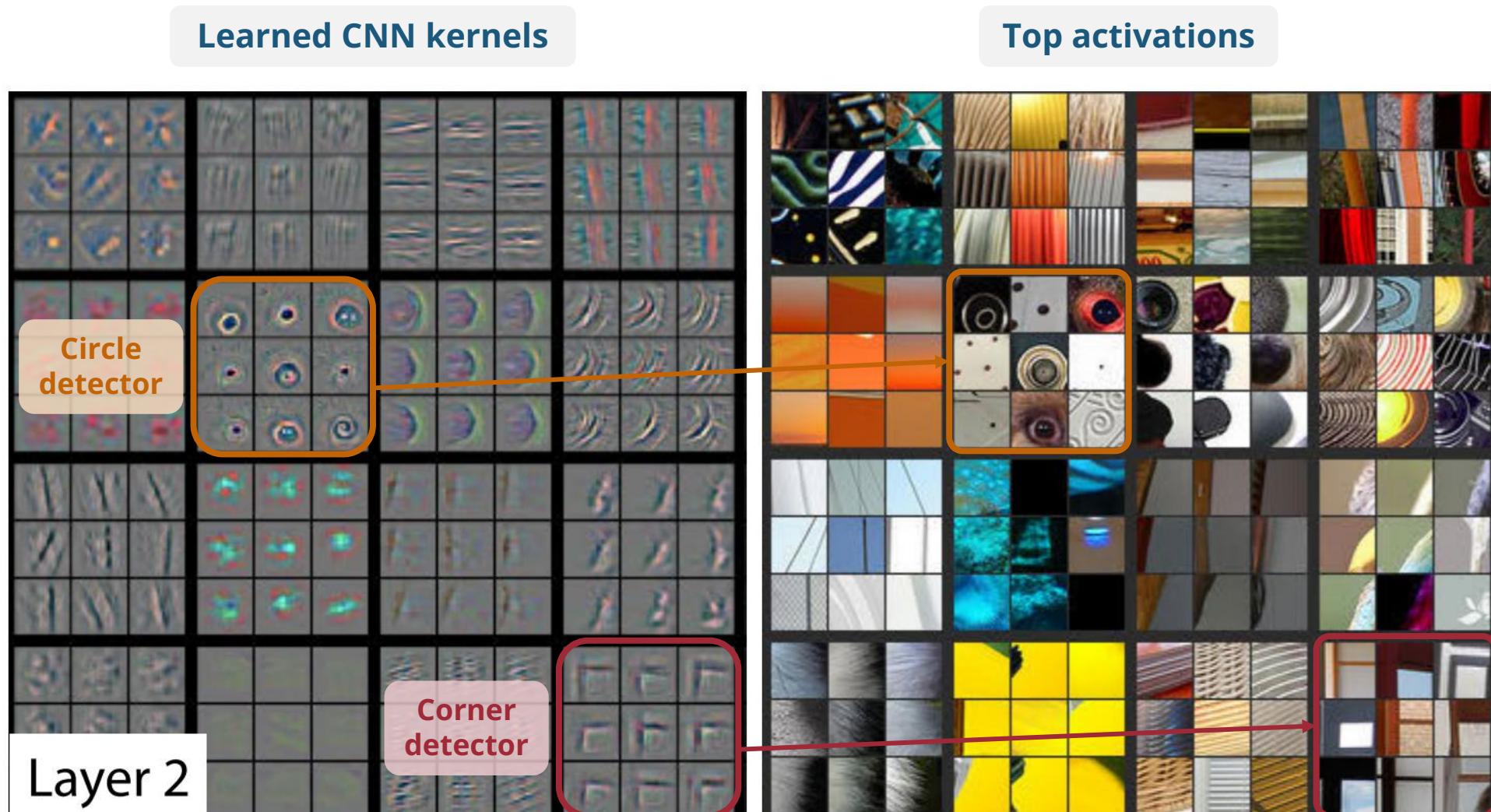
A Toy Example



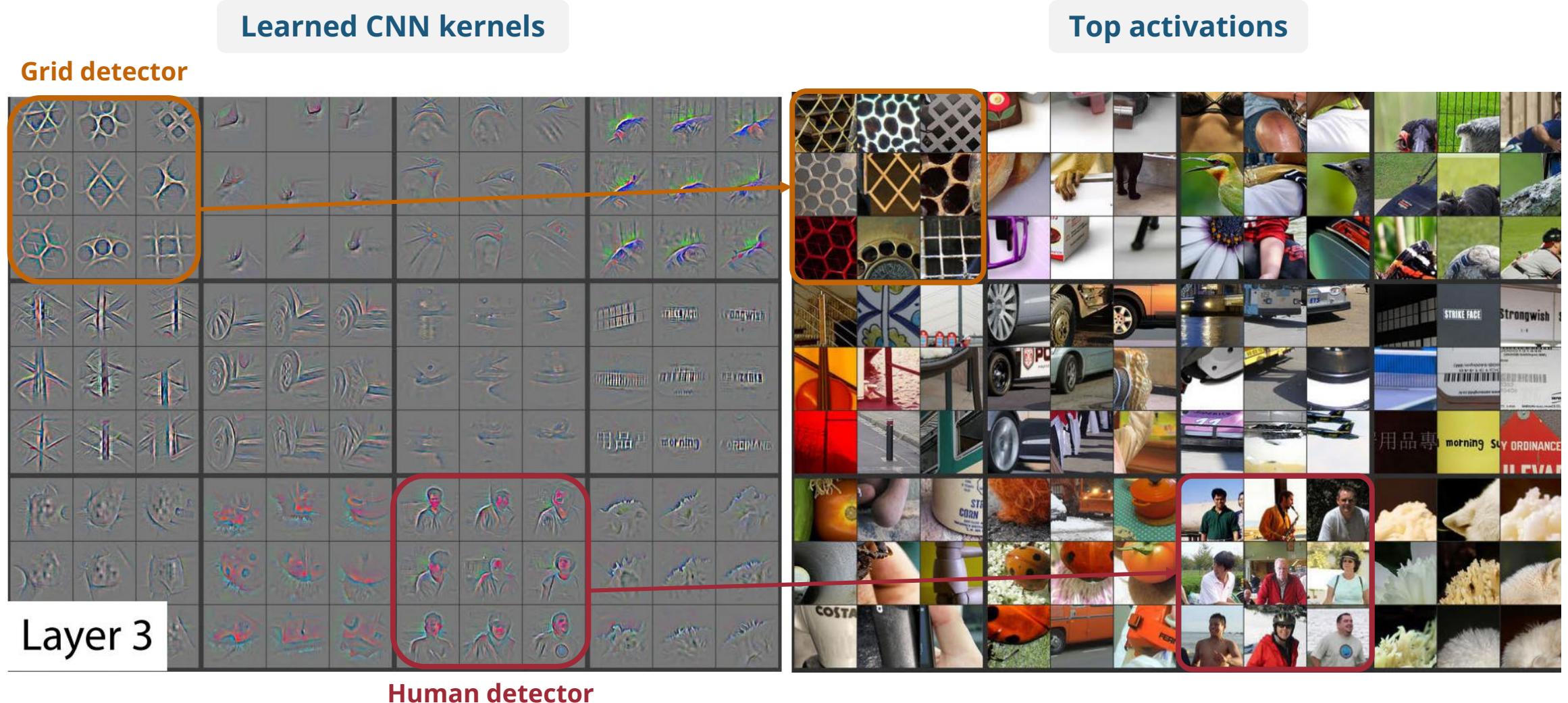
| Learned CNN Kernels in a Trained AlexNet



Learned CNN Kernels in a Trained AlexNet

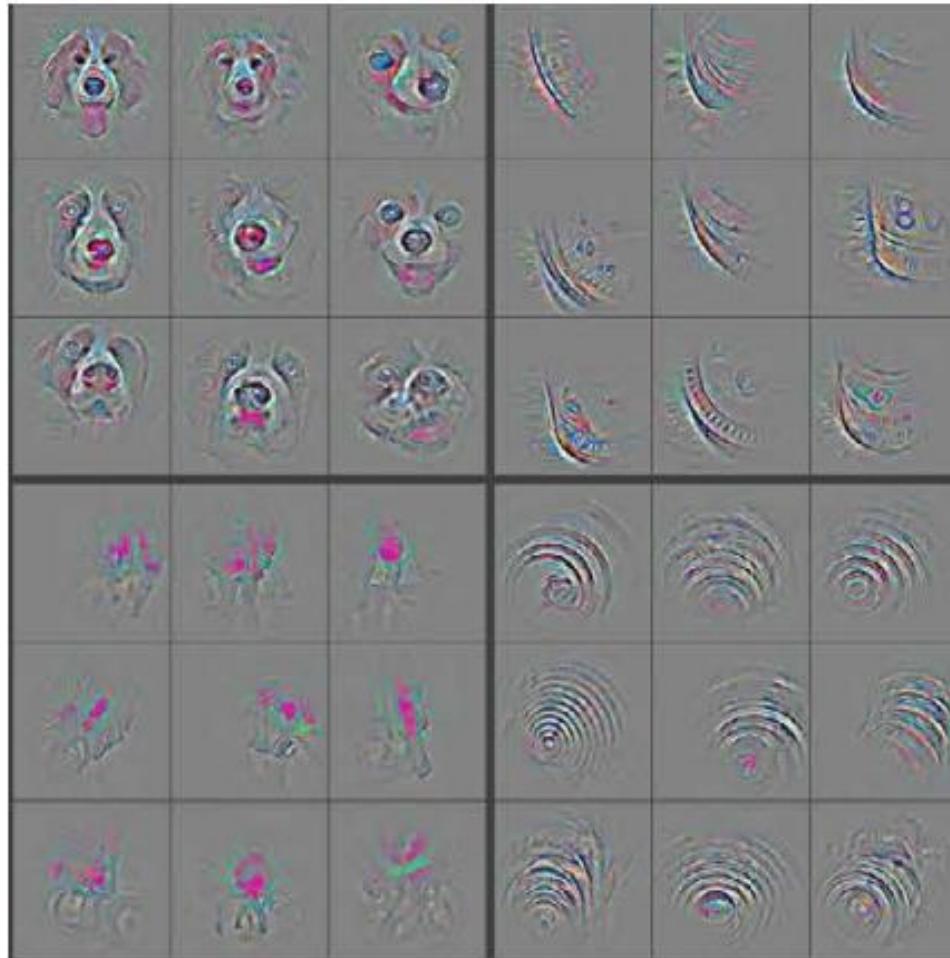


Learned CNN Kernels in a Trained AlexNet

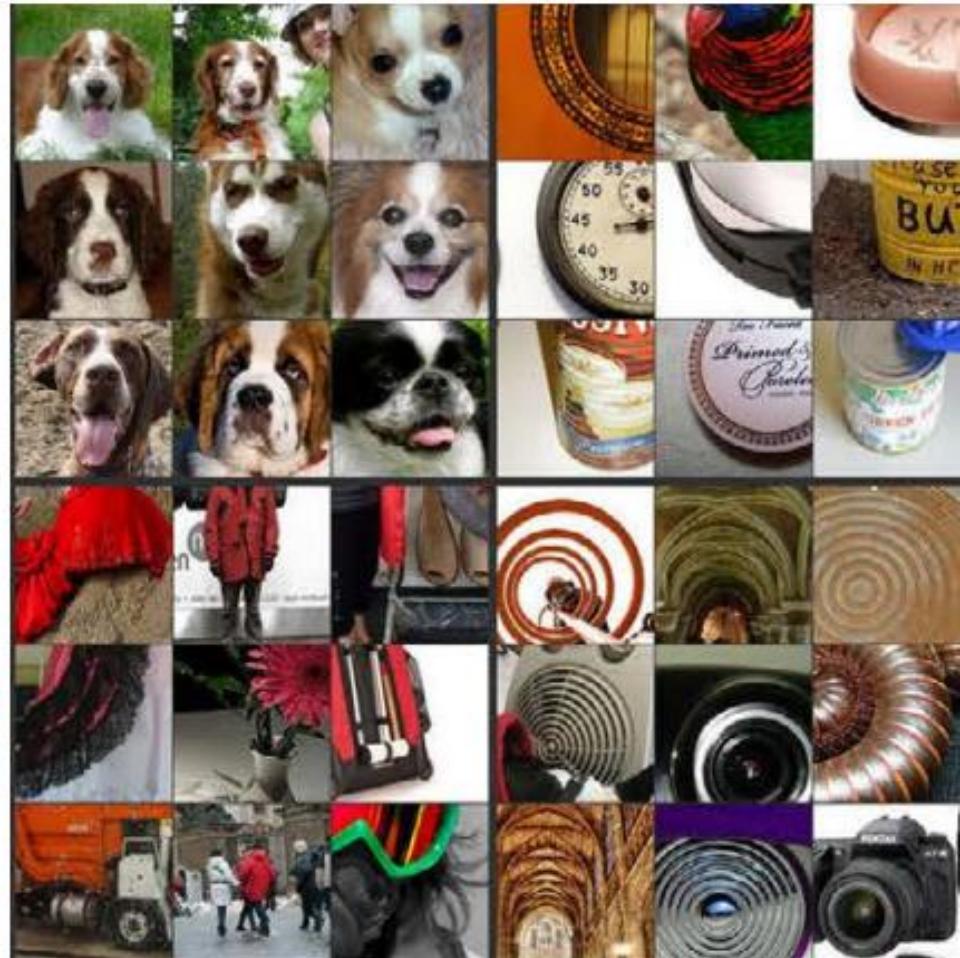


Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

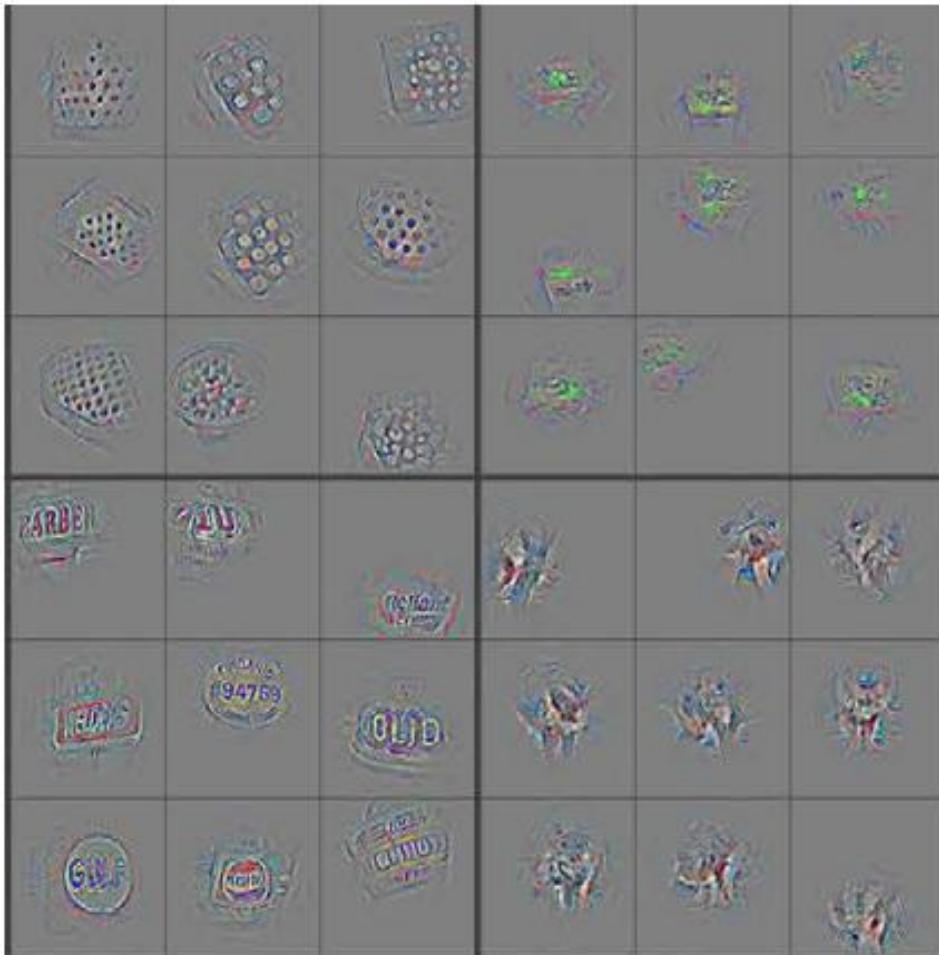


Top activations



Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

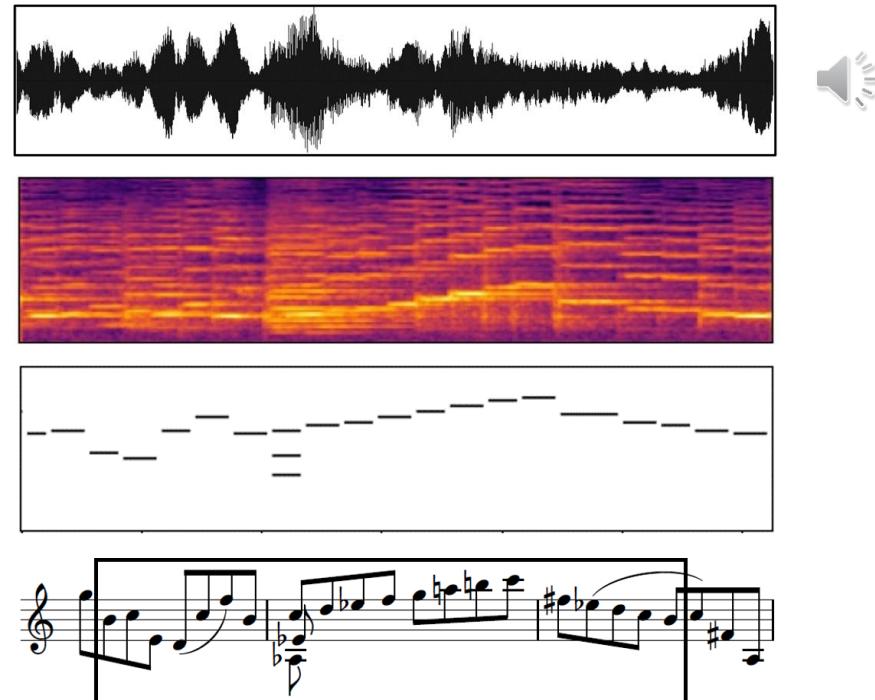


Top activations



Next Lecture

Music Analysis



(Source: Dong et al., 2022)