

PAT 498/598 (Winter 2025)

# Music & AI

## Lecture 10: Convolutional Neural Networks

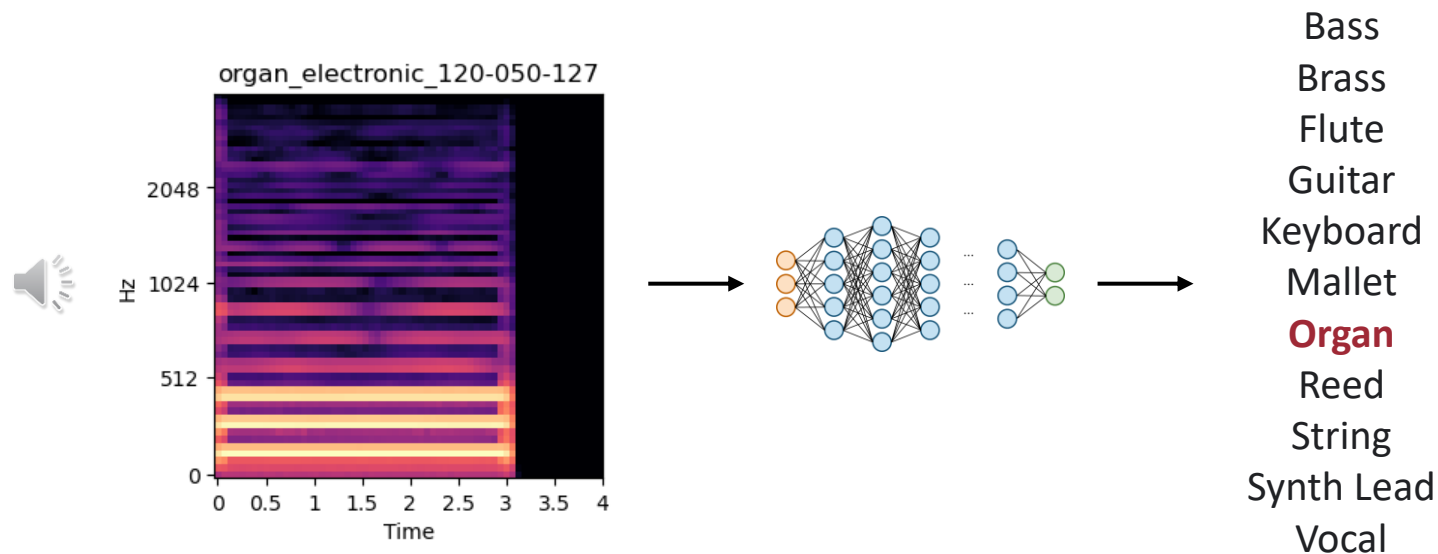
Instructor: Hao-Wen Dong



SCHOOL OF MUSIC, THEATRE & DANCE  
PERFORMING ARTS TECHNOLOGY  
UNIVERSITY OF MICHIGAN

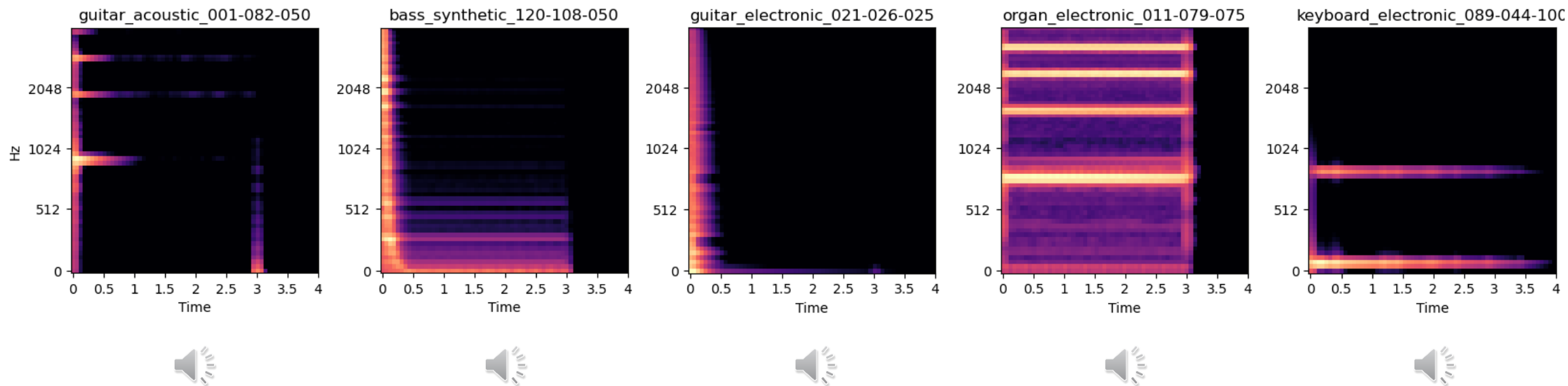
# Homework 3: Musical Note Classification using CNNs

- Train a CNN that can classify audio files into their **instrument families**
  - **Input:** 64x64 mel spectrogram
  - **Output:** 11 instrument classes
  - Using the **NSynth** dataset (Engel et al., 2017)



# NSynth Dataset

- A collection of 305,979 **single-shot musical notes** (Engel et al., 2017)
  - Produced from 1,006 **commercial sample libraries**
  - With different **MIDI pitches** (21–108) and **velocities** (25, 50, 75, 100, 127)



## Homework 3: Musical Note Classification using CNNs

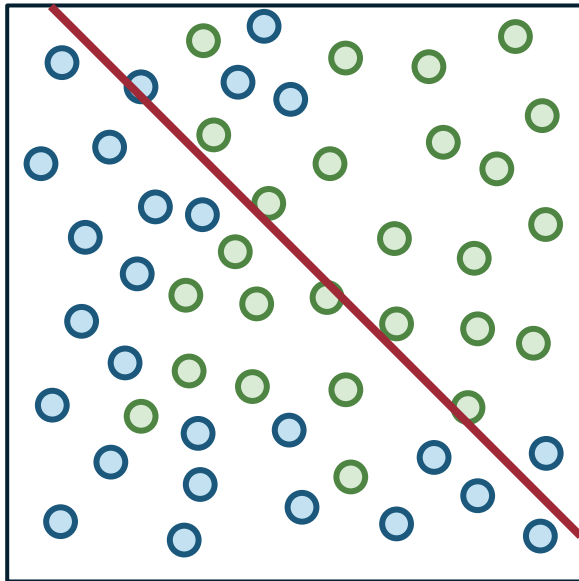
- Instructions will be released on Gradescope
- Due at **11:59pm ET** on **February 17**
- Late submissions: **1 point deducted per day**

## (Recap) In-distribution vs Out-of-distribution

- **Key:** Make the training distribution **closer to** the target distribution
- First, we need to **define our target distribution**
- Then, we can try to
  - Collect a **diverse** dataset covering that covers different parts of the target distribution
  - Apply **data augmentation** to fill the gaps in the distribution

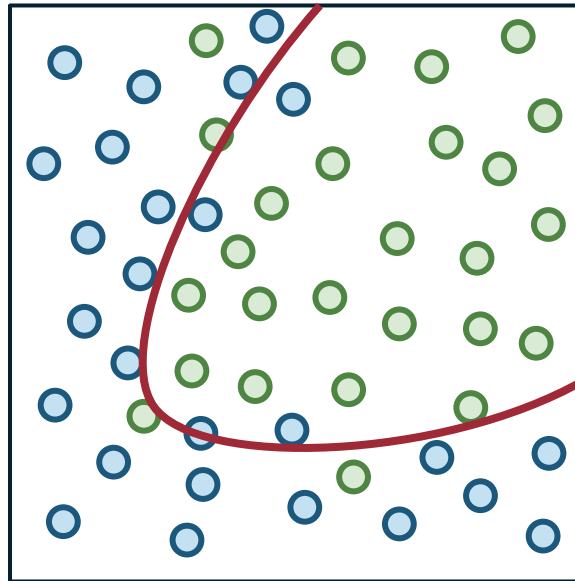
# (Recap) Overfitting & Underfitting

Underfitting

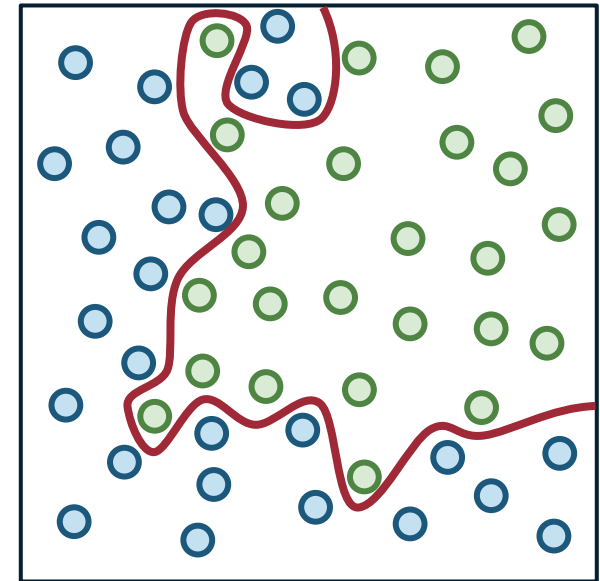


Model too **in**expressive

Good fit!



Overfitting



Model too **ex**pressive

# (Recap) Train-Validation-Test Split



# (Recap) Train-Validation-Test Split

**Training**



**Test**





# (Recap) Train-Validation-Test Split

**Training**



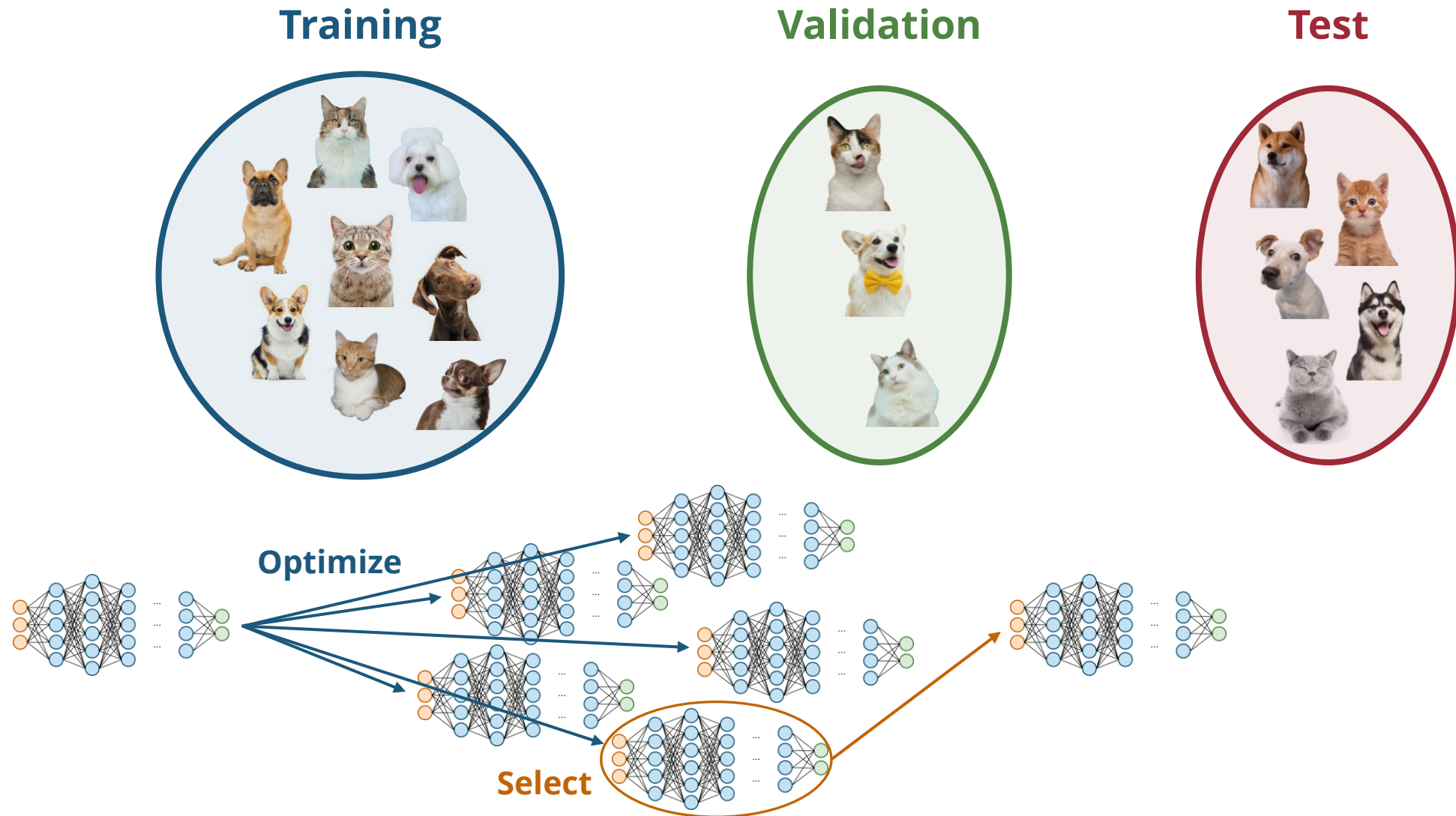
**Validation**



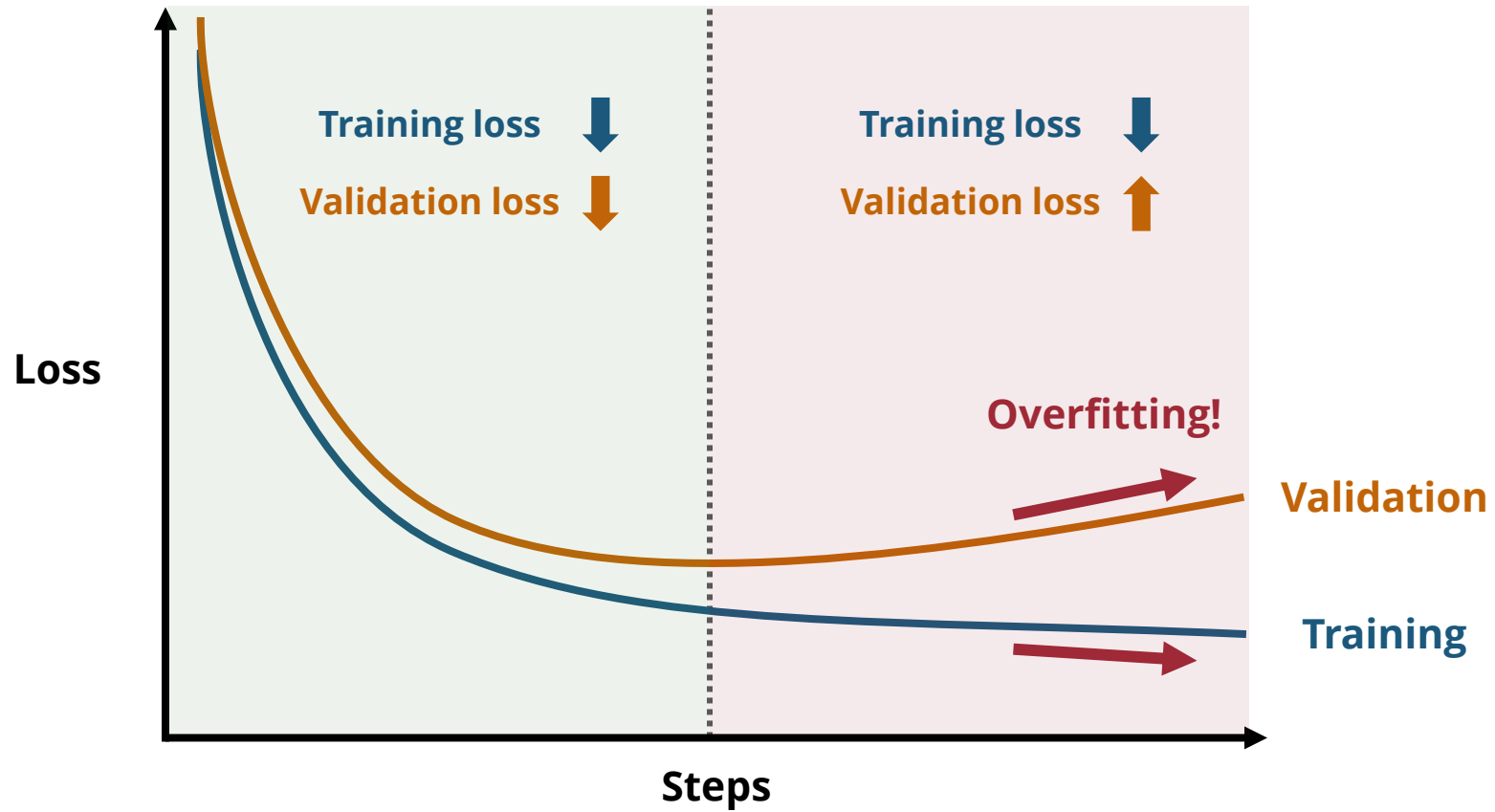
**Test**



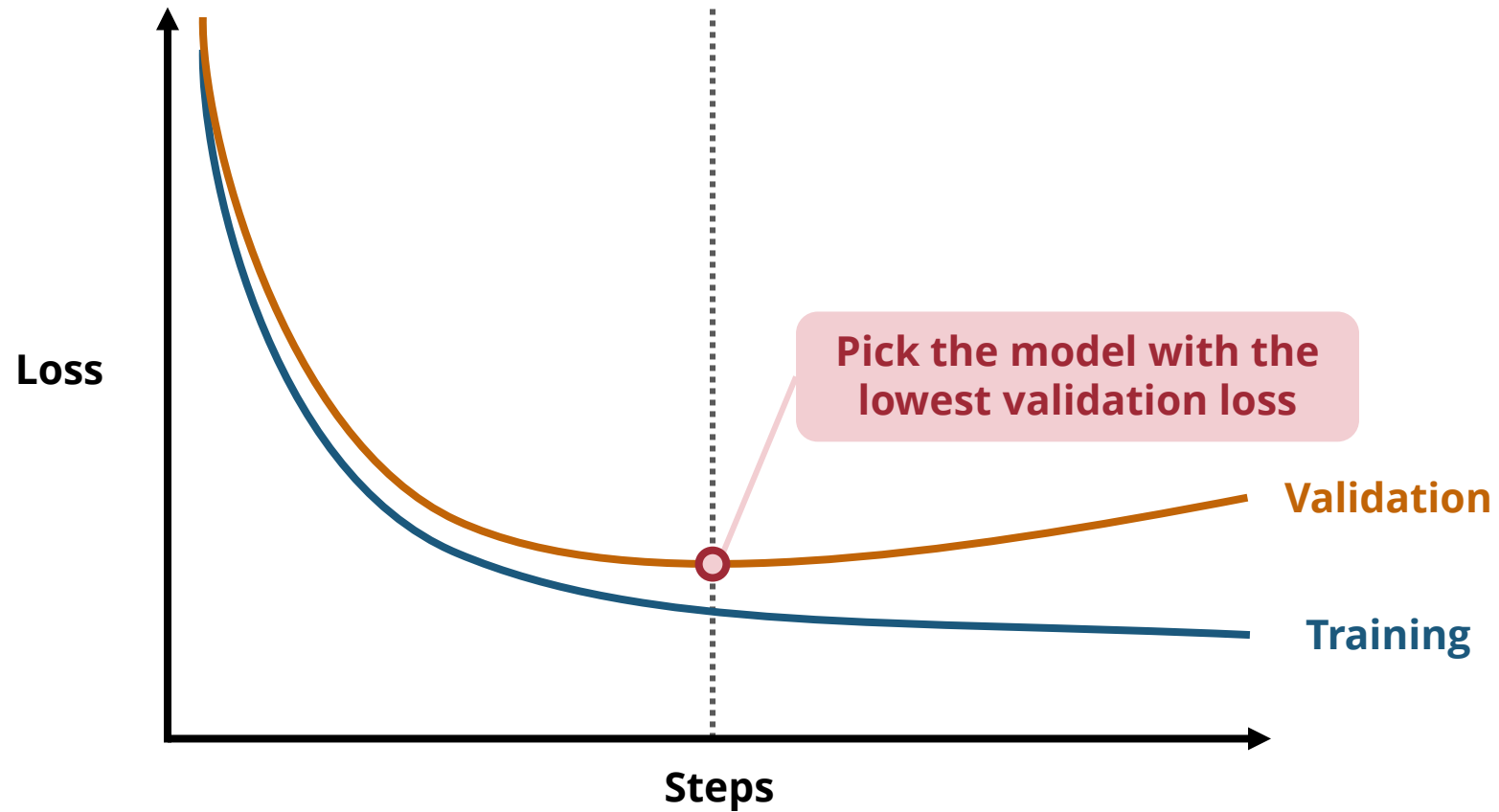
# (Recap) Training-Validation-Test Pipeline



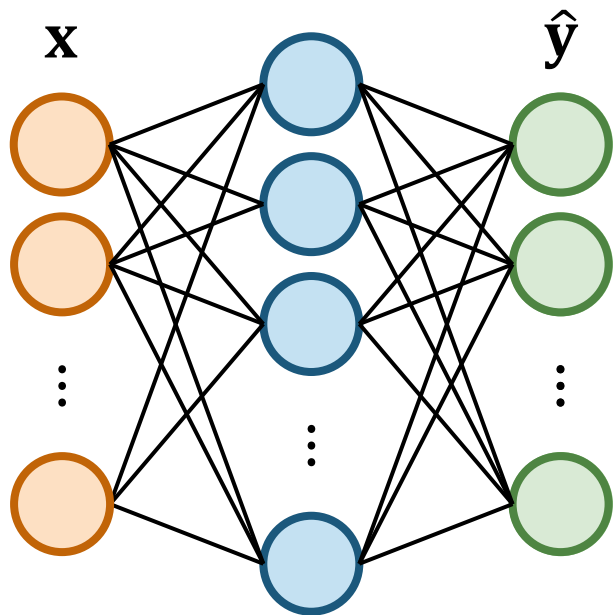
# (Recap) Training vs Validation Losses



# (Recap) Training vs Validation Losses

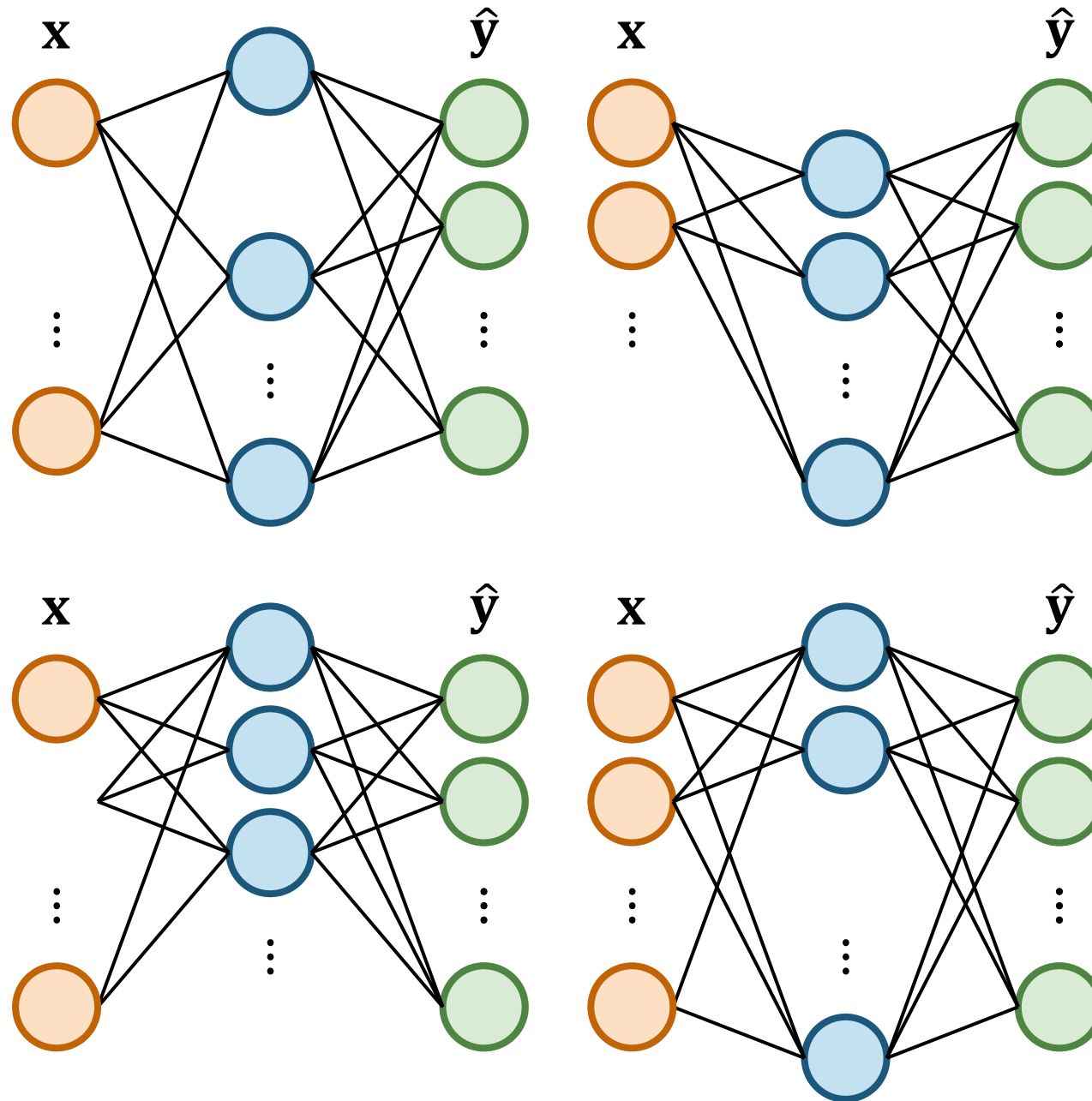


# (Recap) Dropout

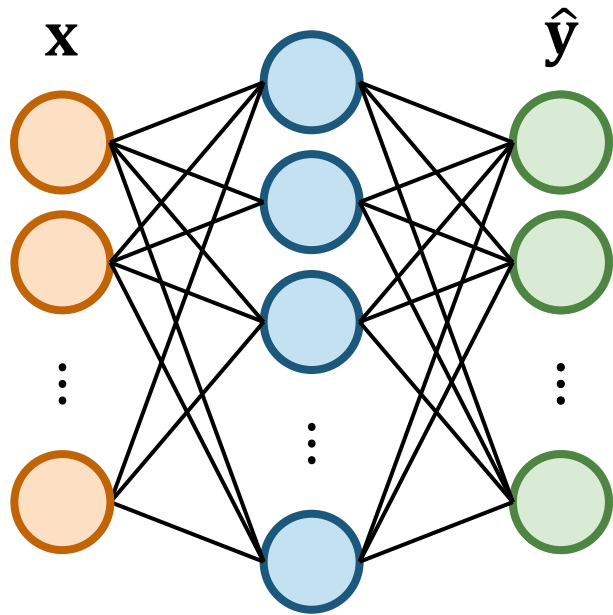


Each neuron may be removed with probability  $p$  during training

Dropout rate

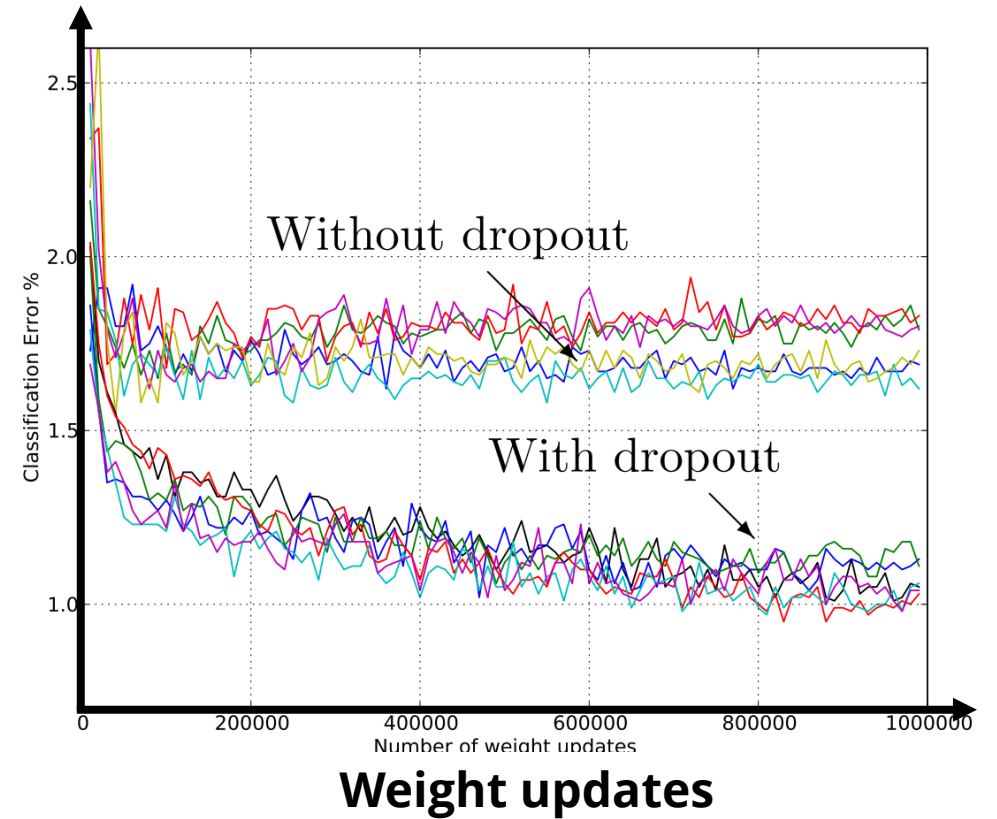


# (Recap) Dropout



Each neuron may be removed with probability  $p$  during training

Test error rate



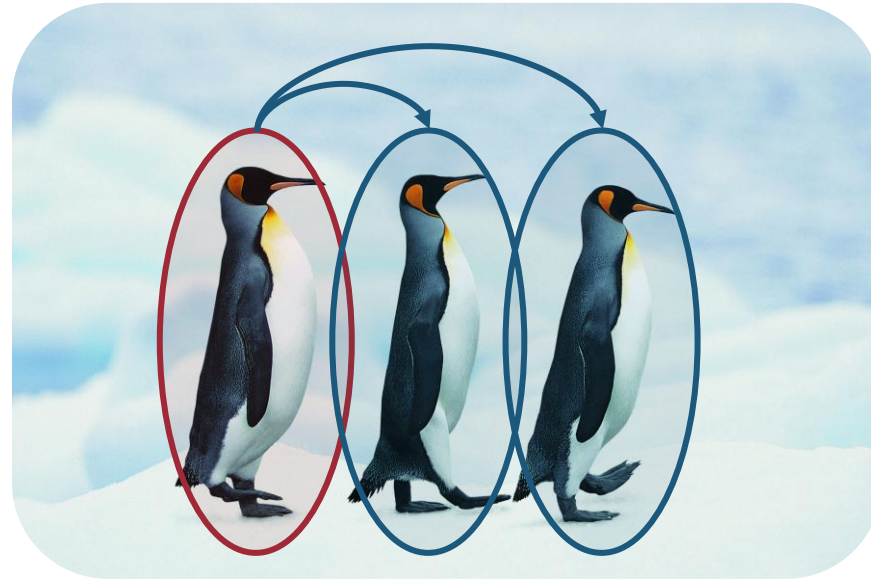
# Convolutional Neural Networks (CNNs)

# Convolutional Neural Networks (CNNs)

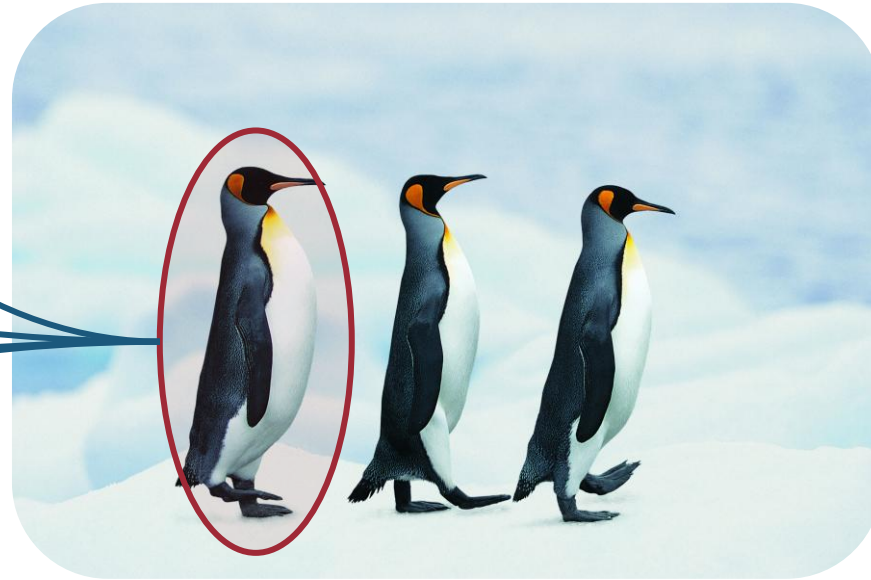
- **Intuition:** Learn **reusable local pattern detector**
- Widely used in **computer vision**
- Also used for music and audio
  - Representing music as **piano rolls**
  - Representing audio as **spectrograms**



# Reusable Pattern Detectors



# Reusable Pattern Detectors



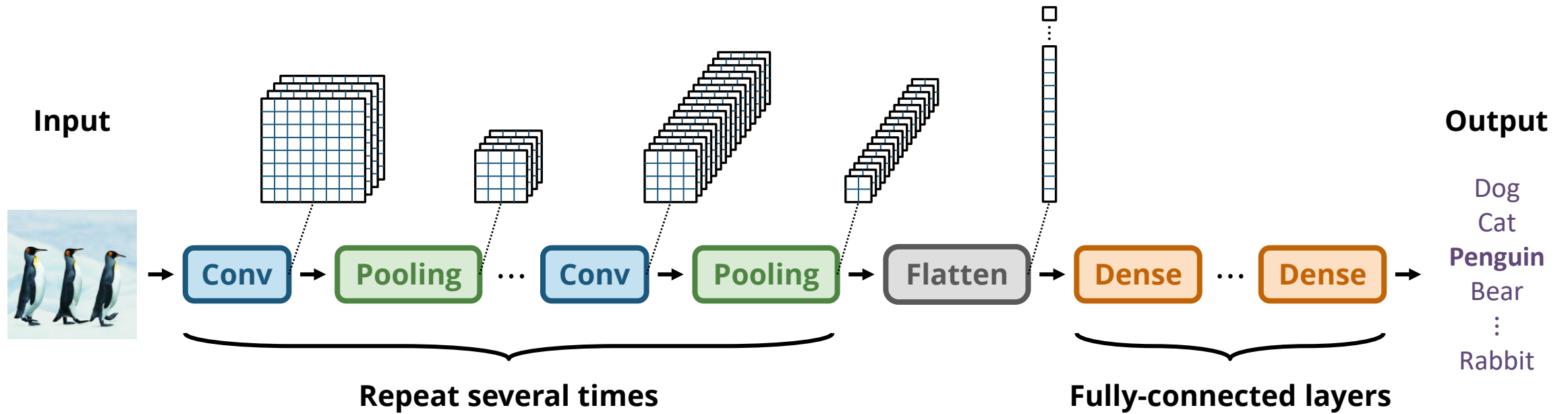
# Reusable Pattern Detectors



# Reusable Pattern Detectors



# Convolutional Neural Network (CNNs)



# 2D Convolution

Input

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

Kernel

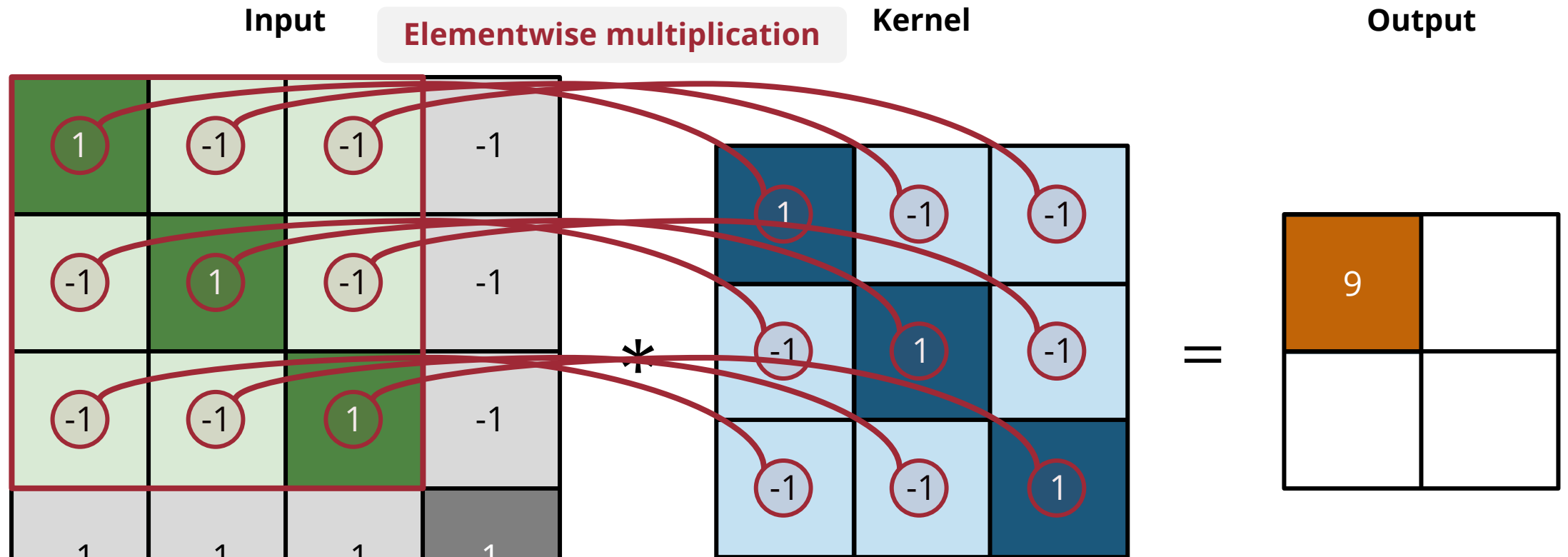
1	-1	-1
-1	1	-1
-1	-1	1

\*

=

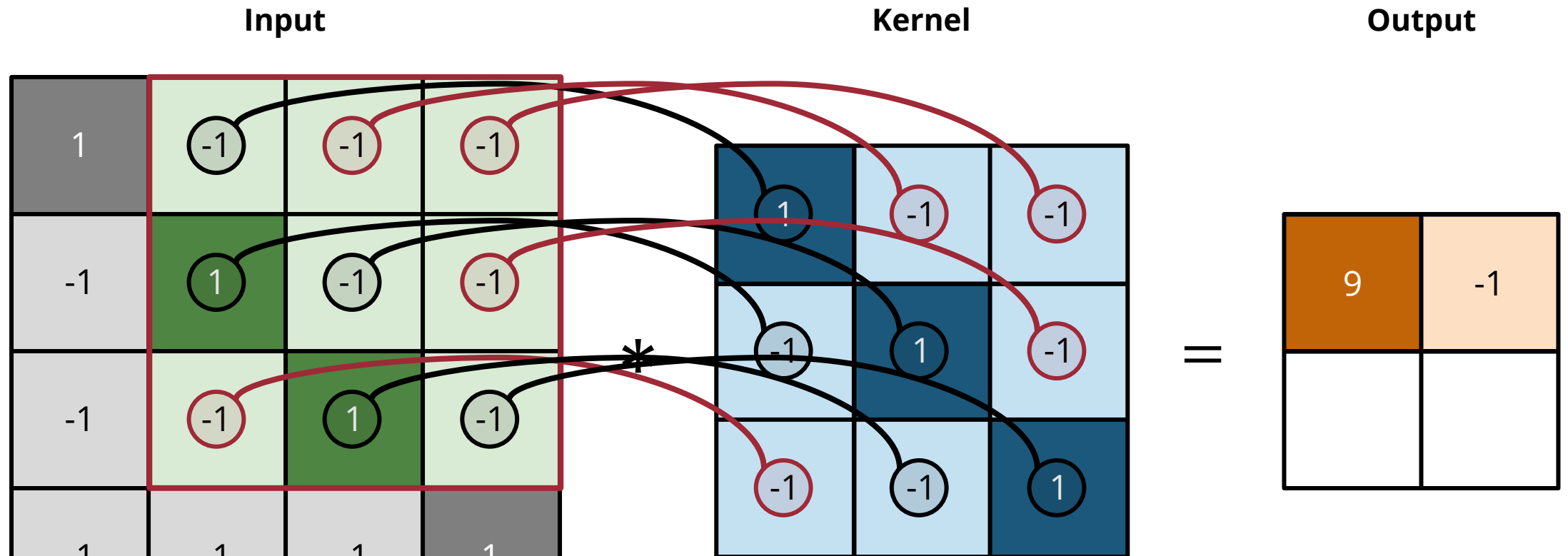
Output


# 2D Convolution



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

# 2D Convolution



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



# 2D Convolution

Input

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

Kernel

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

Input

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

Kernel

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

Output

9	-1
-1	9

\*

=

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

# 2D Convolution

Input

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

\*

Kernel

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

=

Output

9	
	9

High activation when the local pattern is close to the kernel

# 2D Convolution

Input

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

Kernel

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

\*

=

Output

	-1
-1	

**Low activation when  
the local pattern  
differs from the kernel**

# 2D Convolution

Input

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

\*

Kernel

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

=

Output

9	-1
-1	9

# 2D Convolution

Input

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

\*

Kernel

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

=

Output

1	1
1	5

# 2D Convolution

Input

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

Kernel

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

\*

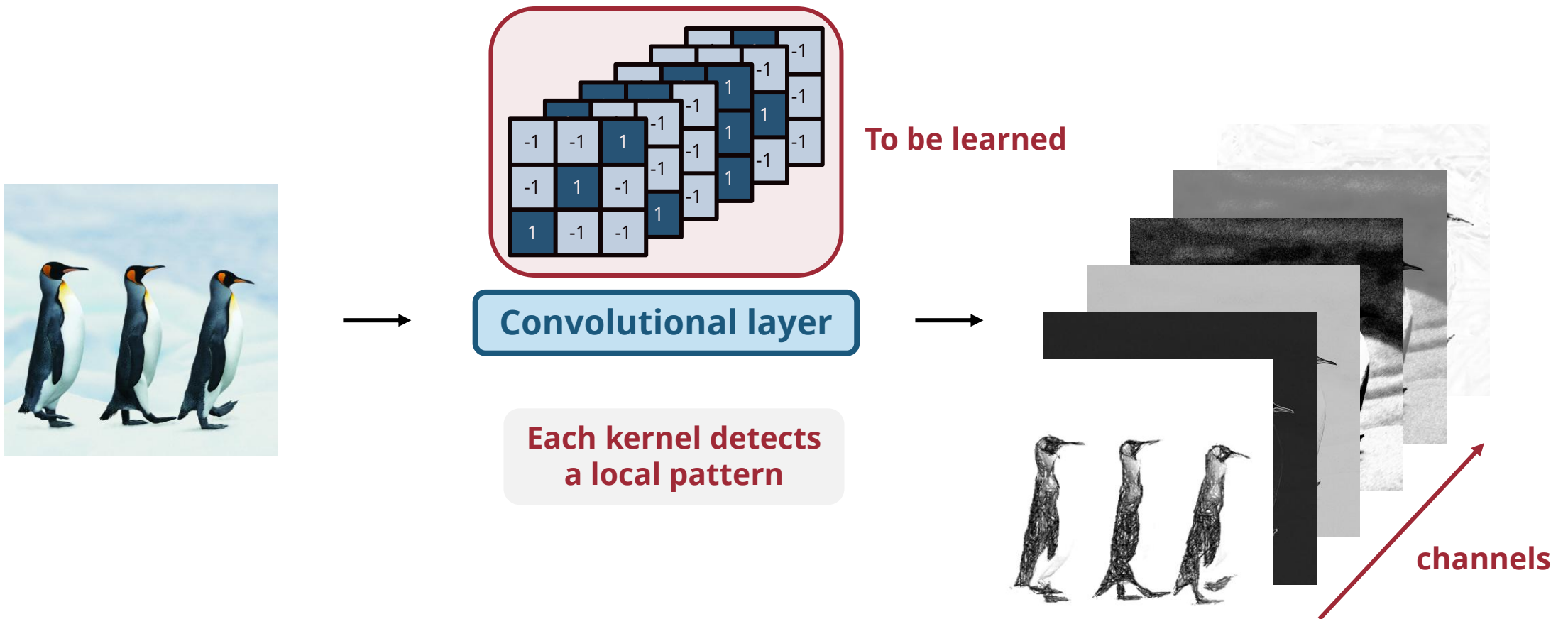
=

Output

9	-1
-1	1

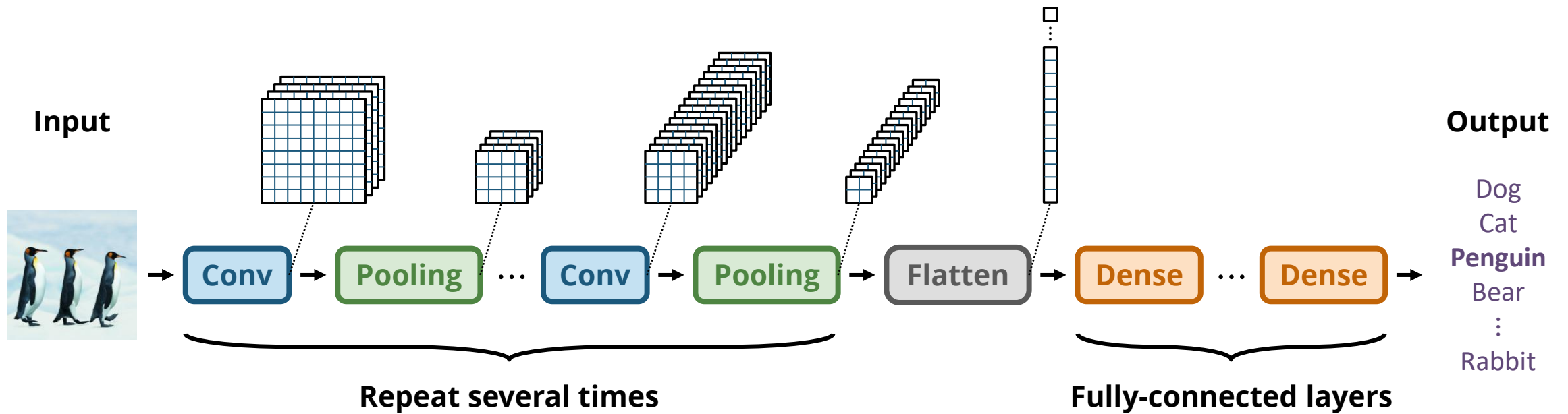
# Convolutional Layer

- 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

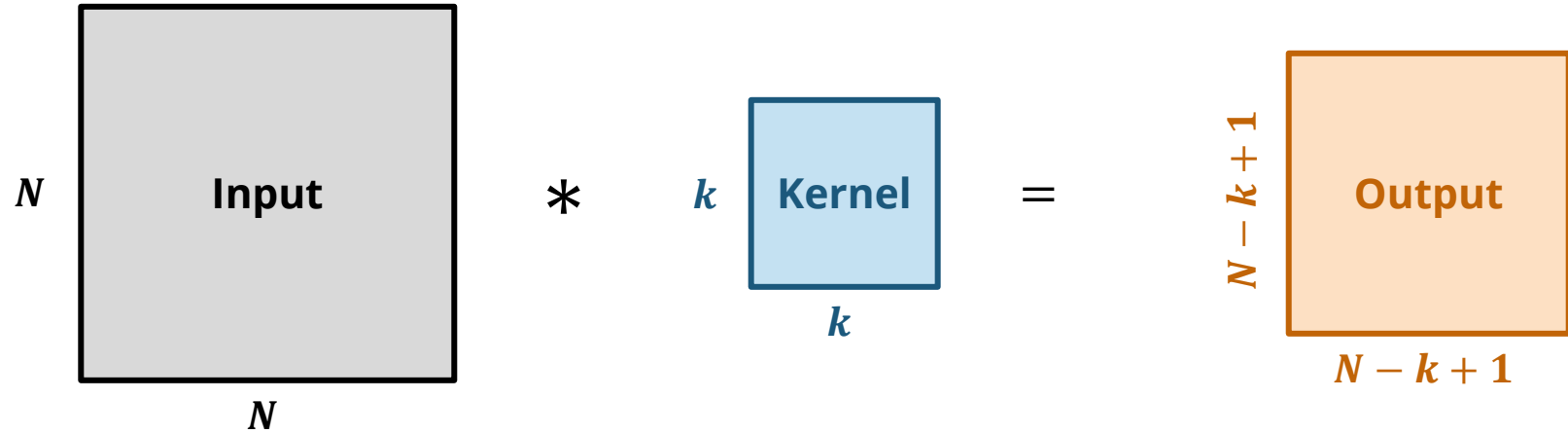
=

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

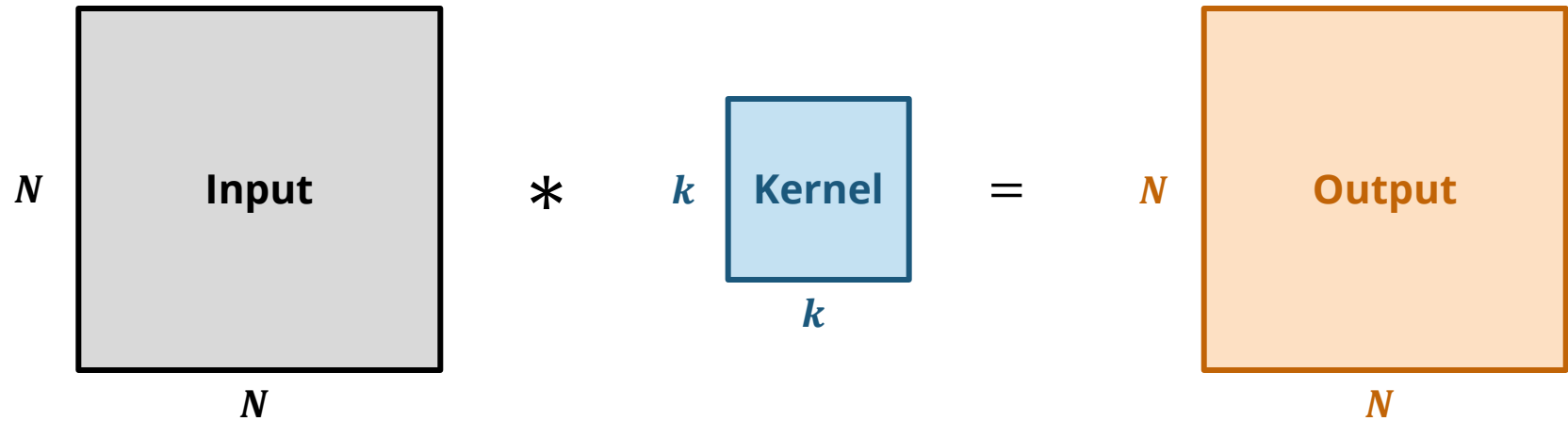
Keep the output of the same size as the input

# Shapes

padding="valid"



padding="same"



# Striding

stride=2

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

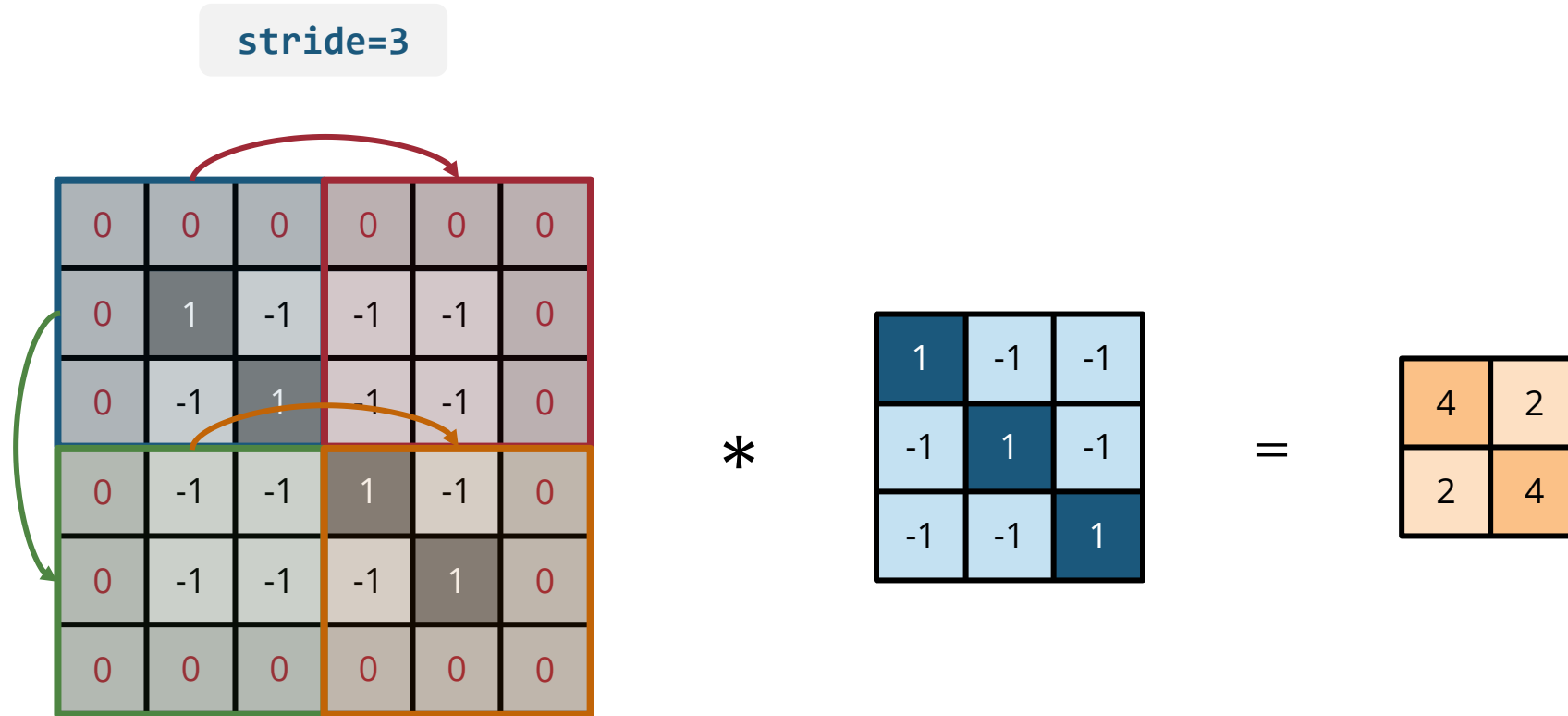
\*

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

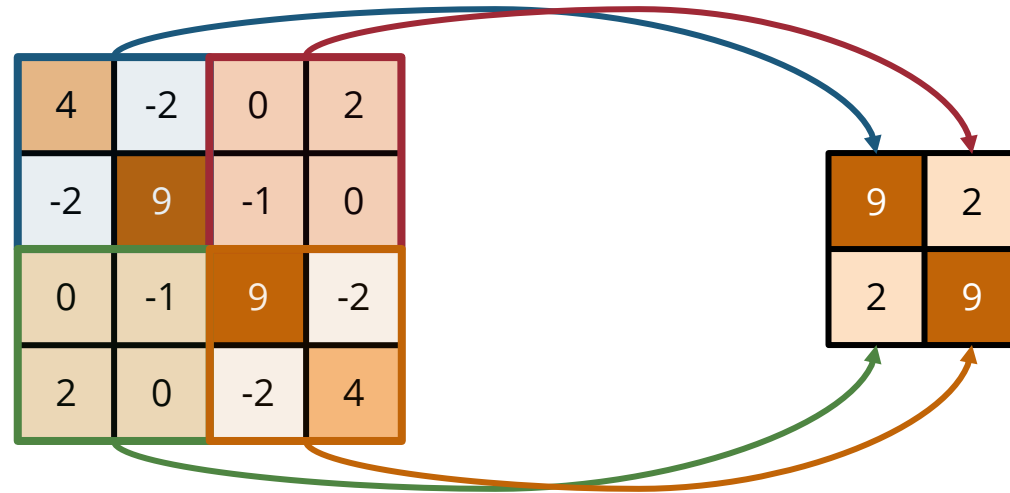
=

4	0
0	9

# Striding

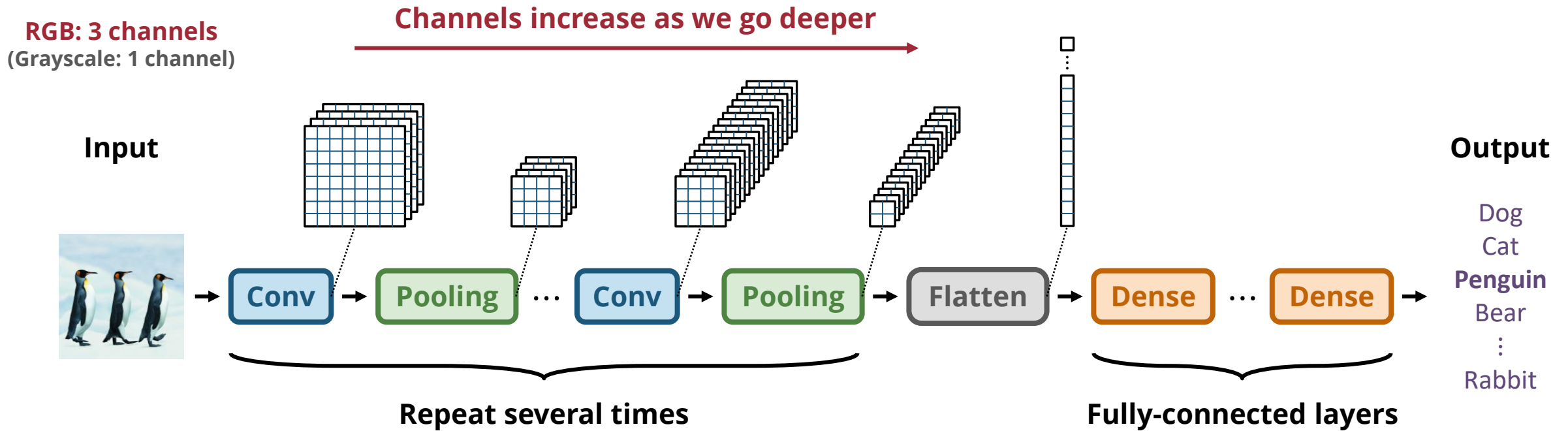


# Max Pooling Layer



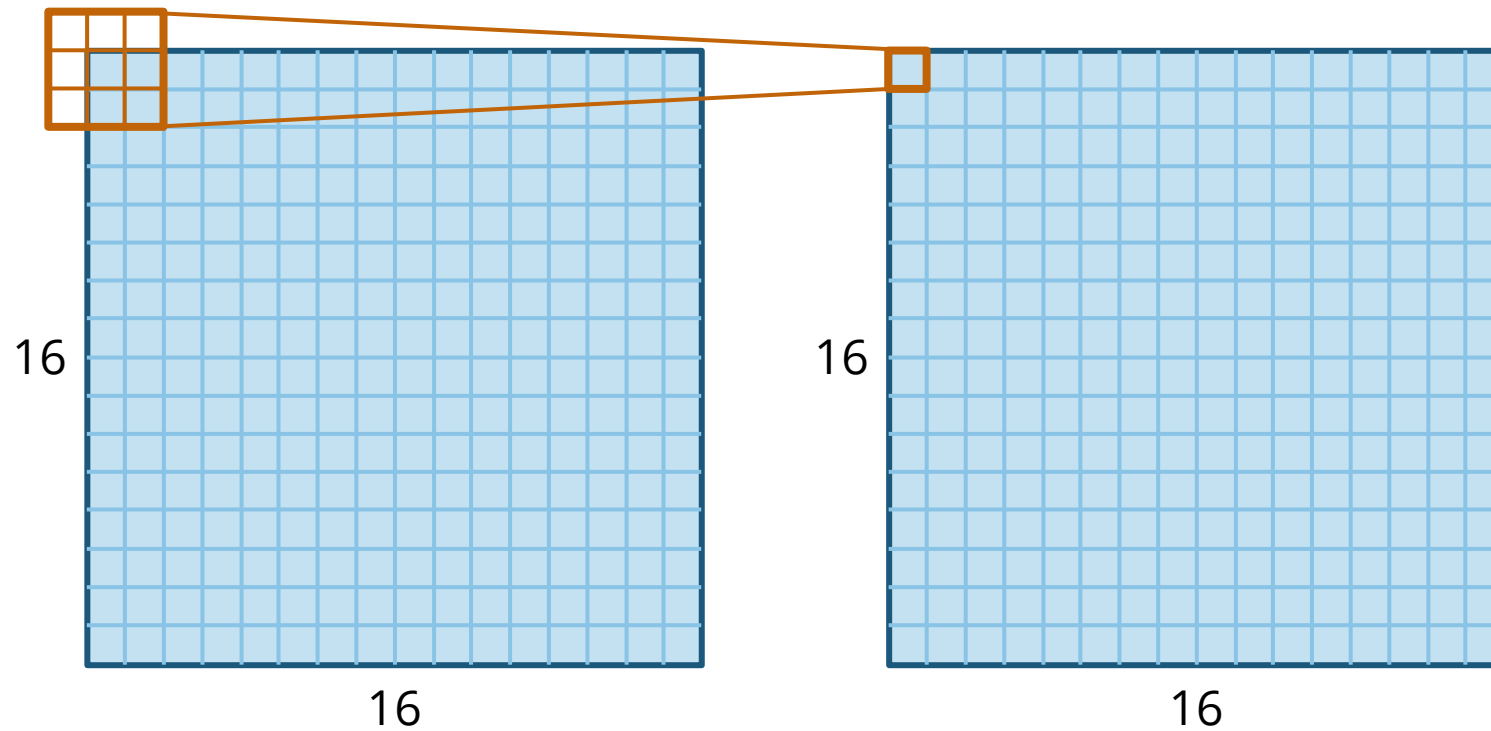
**Downsample and keep the strongest activation in each block**

# Convolutional Neural Network (CNNs)



# A Toy Example

Convolutional layer

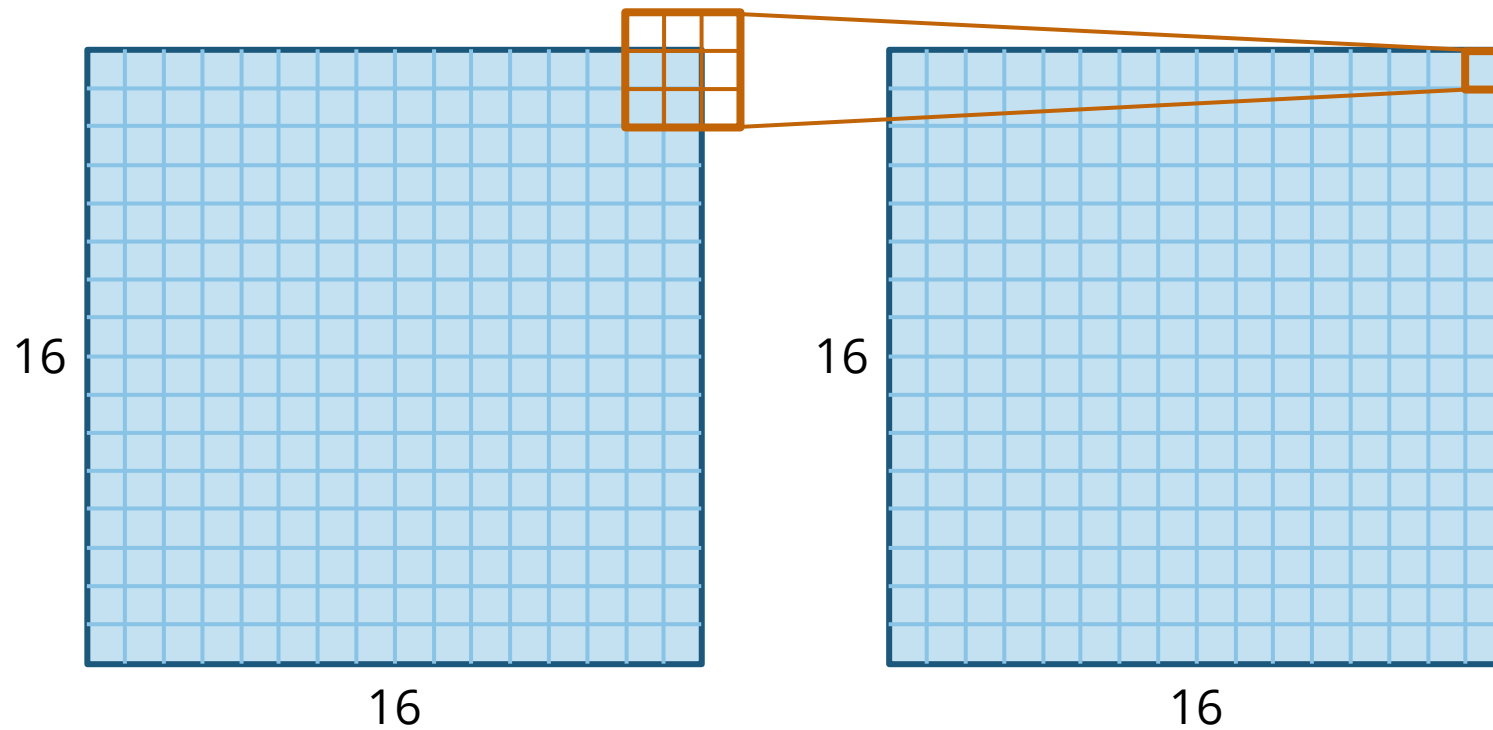


```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```



# A Toy Example

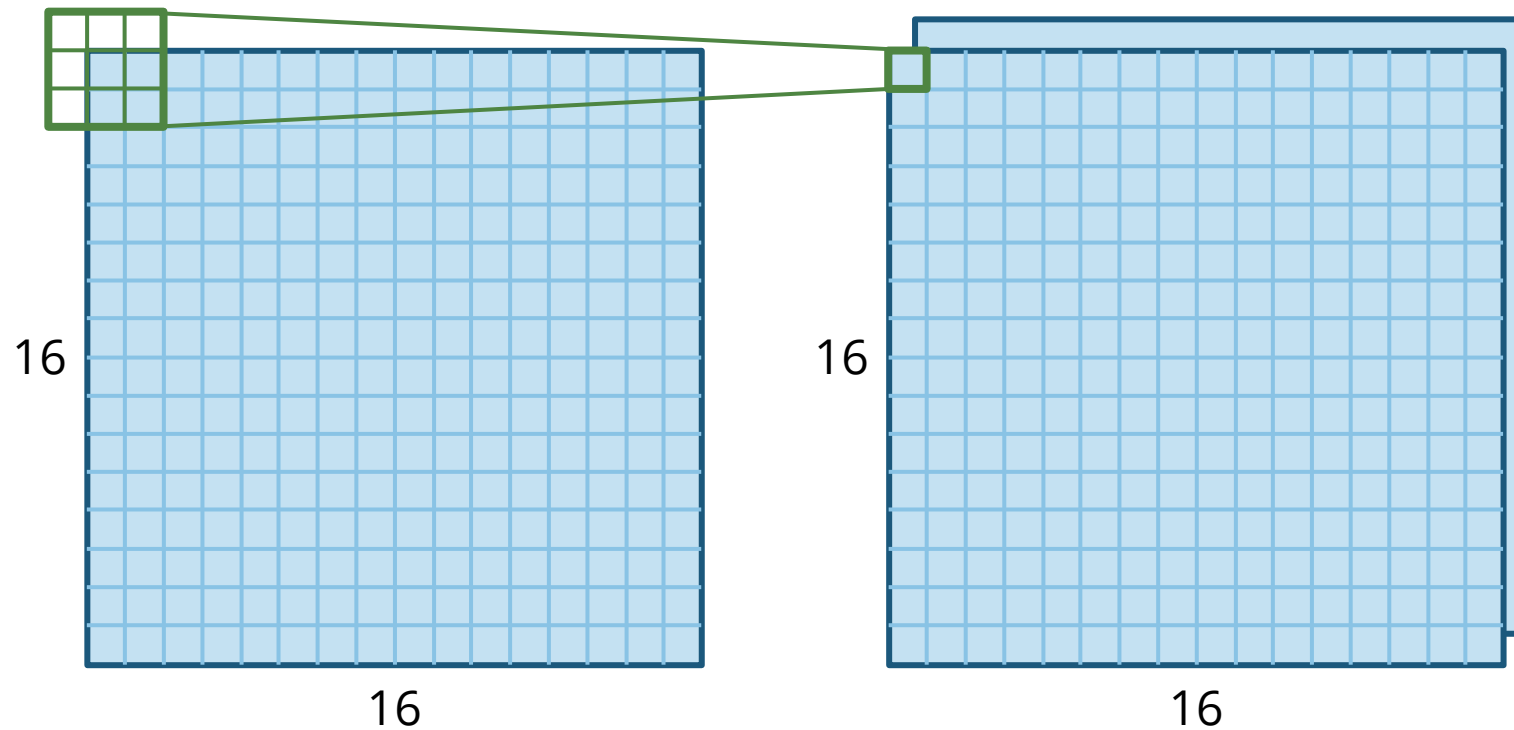
Convolutional layer



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

# A Toy Example

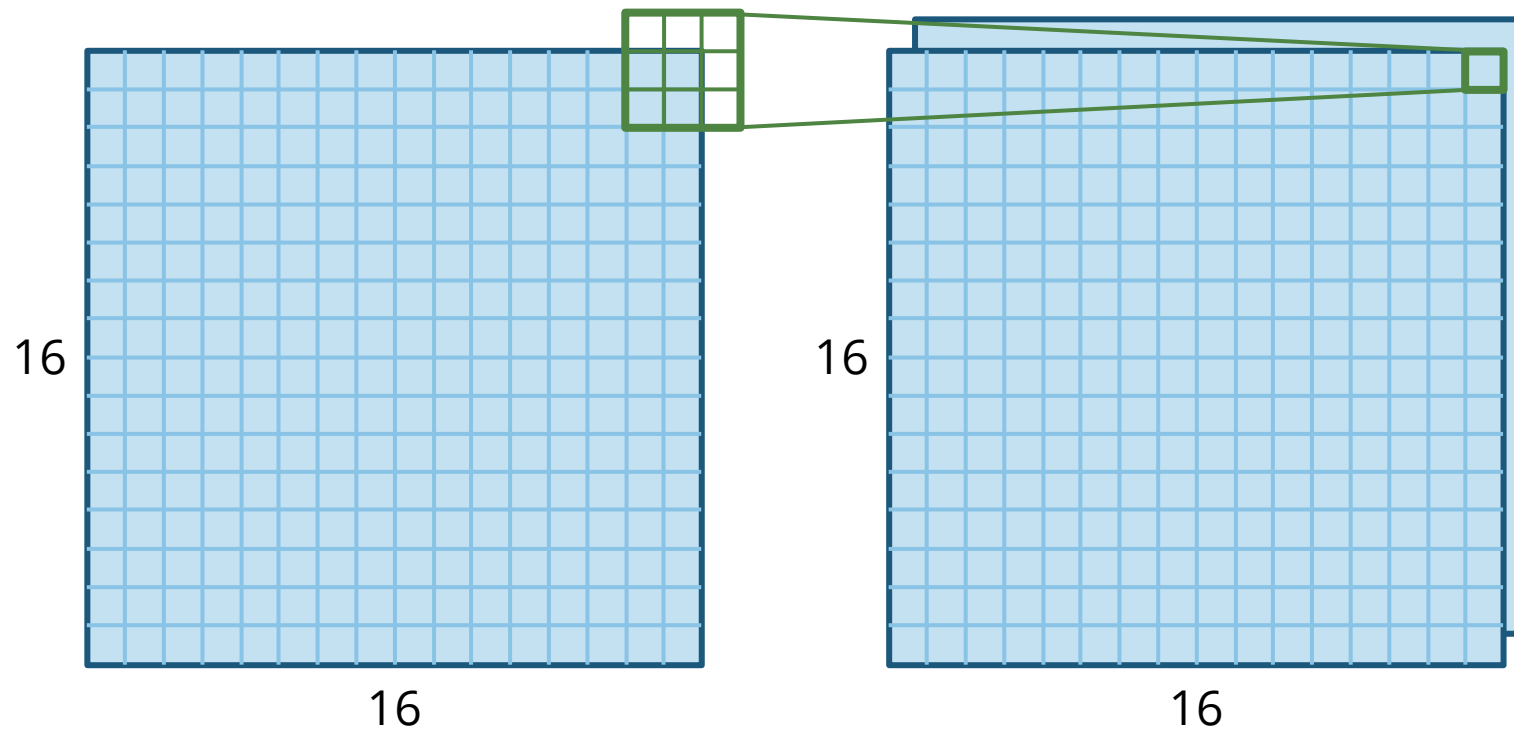
Convolutional layer



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

# A Toy Example

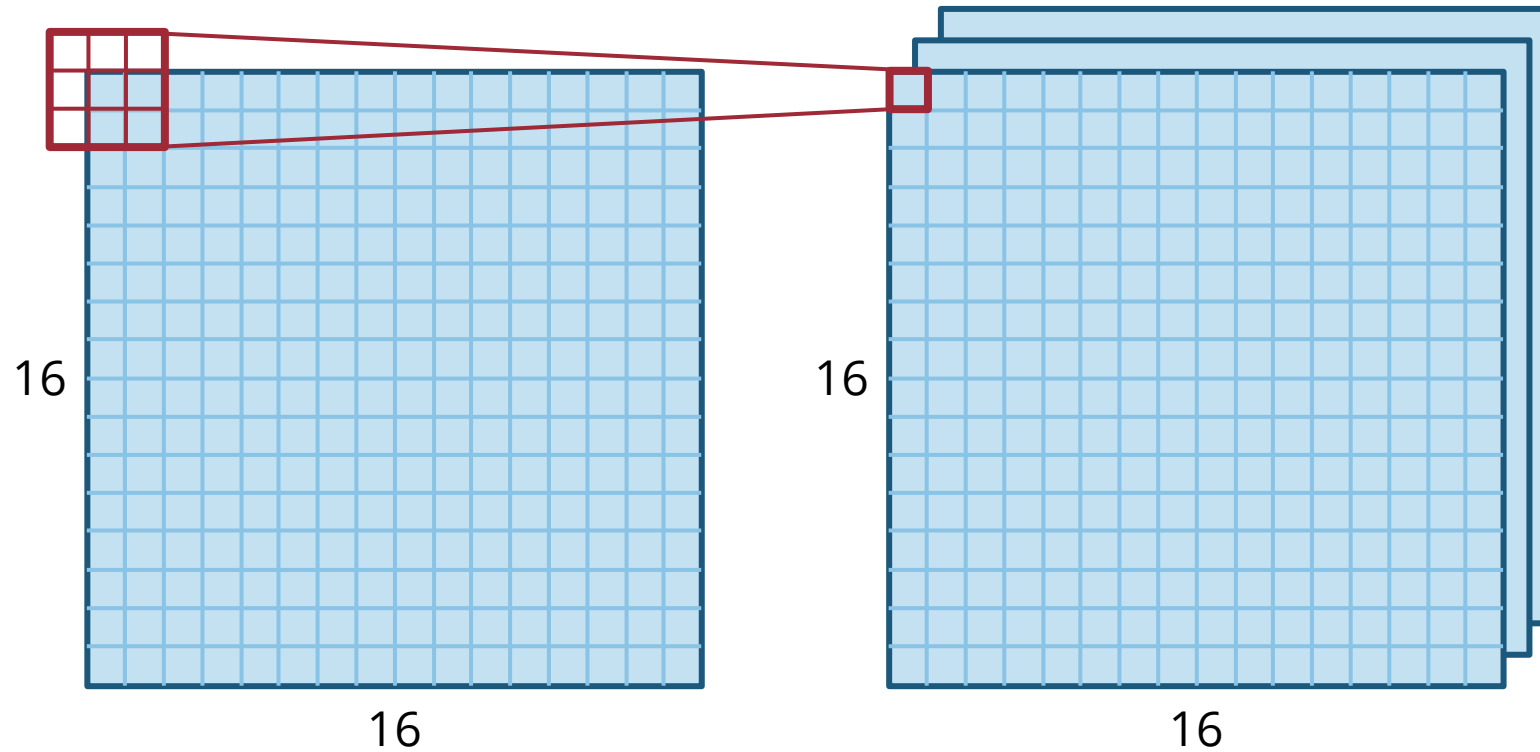
Convolutional layer



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

# A Toy Example

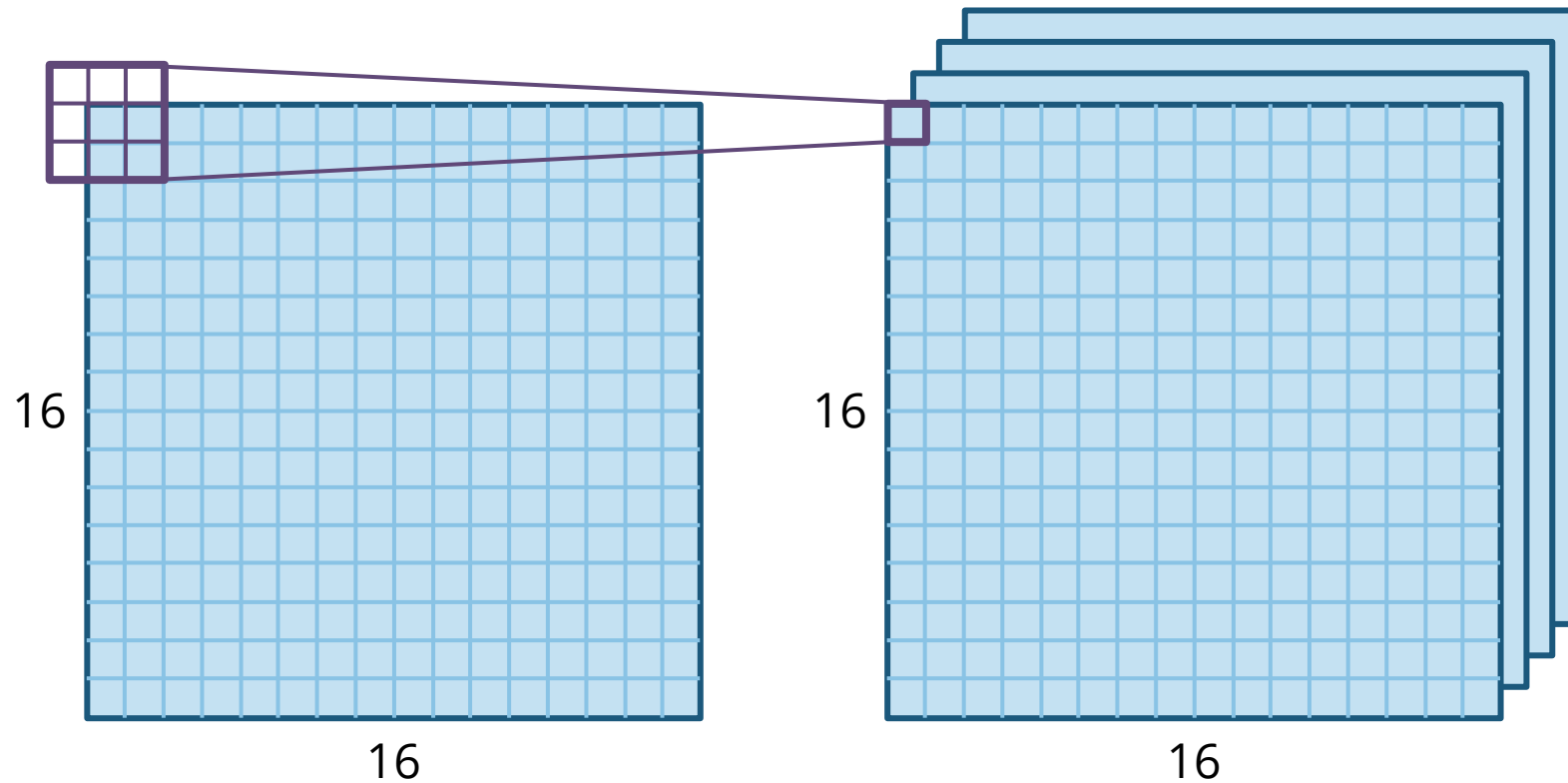
Convolutional layer



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

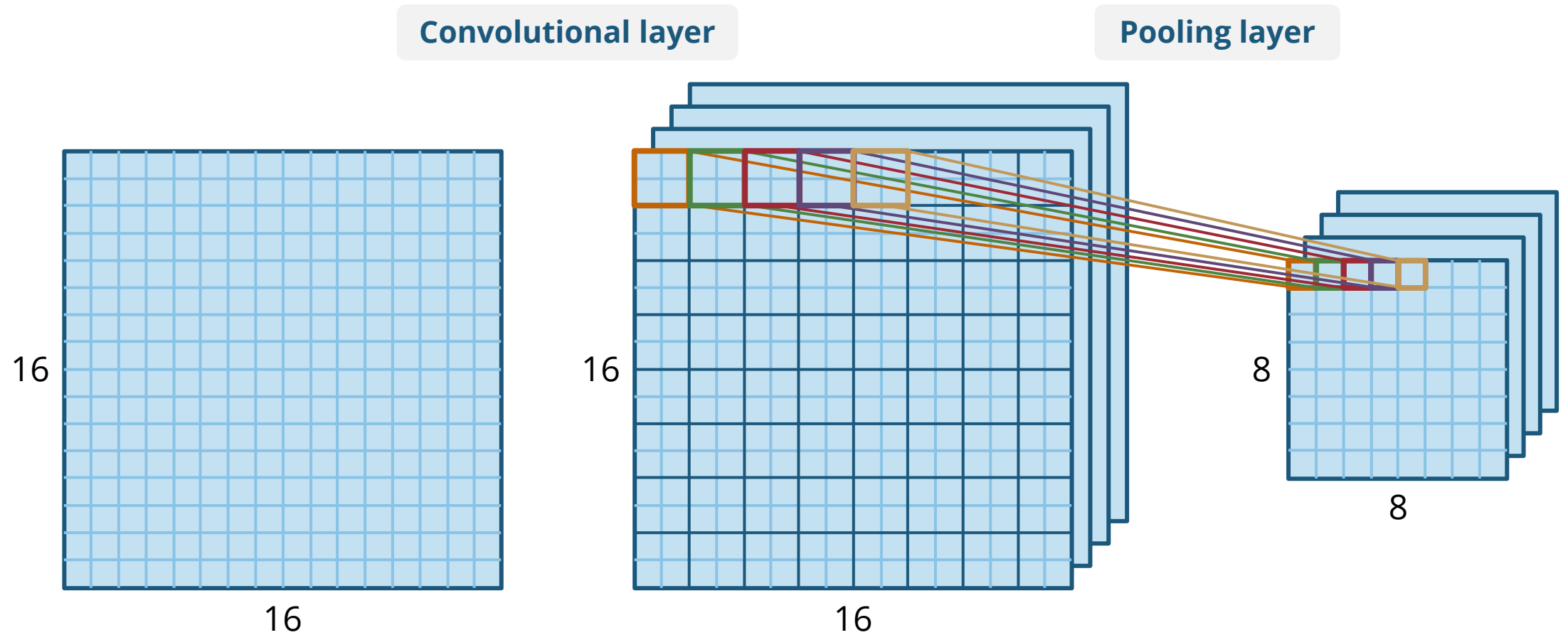
# A Toy Example

Convolutional layer



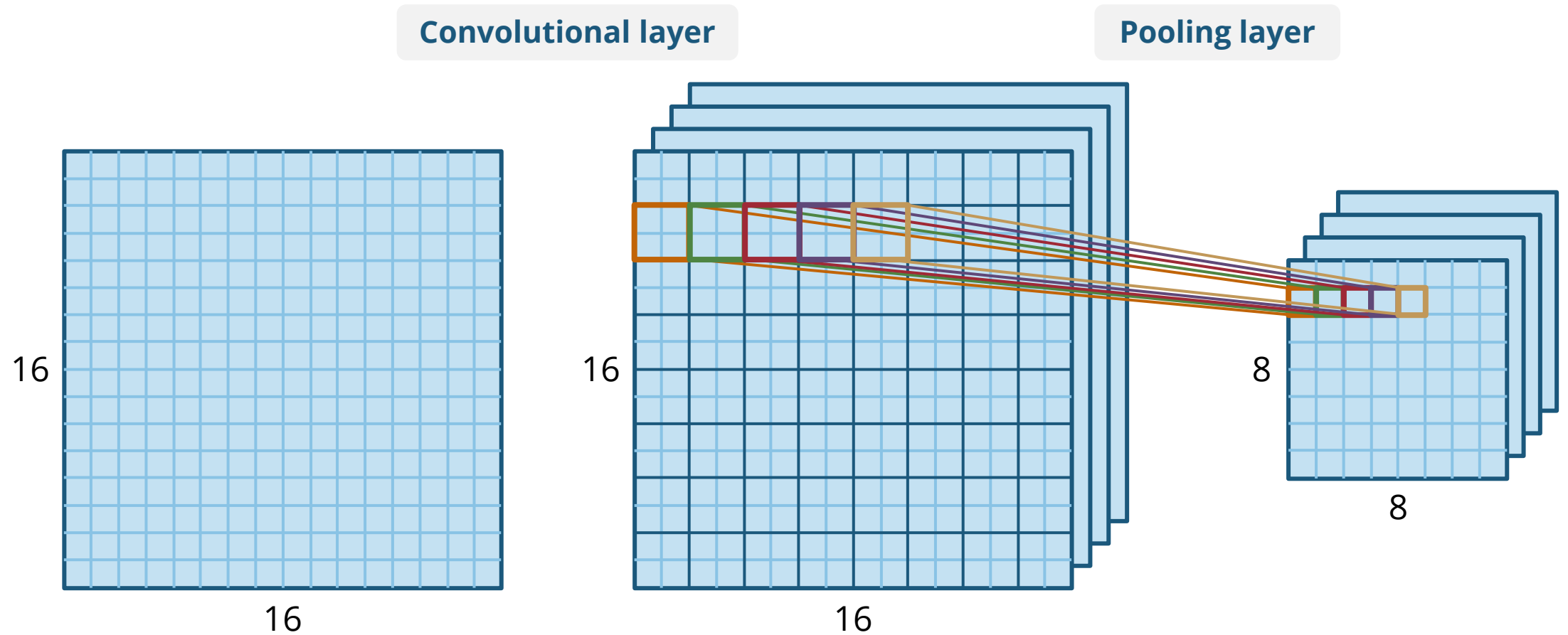
```
nn.Conv2d(in_channels=1, out_channels=4, kernel_size=3, padding="same")
```

# A Toy Example



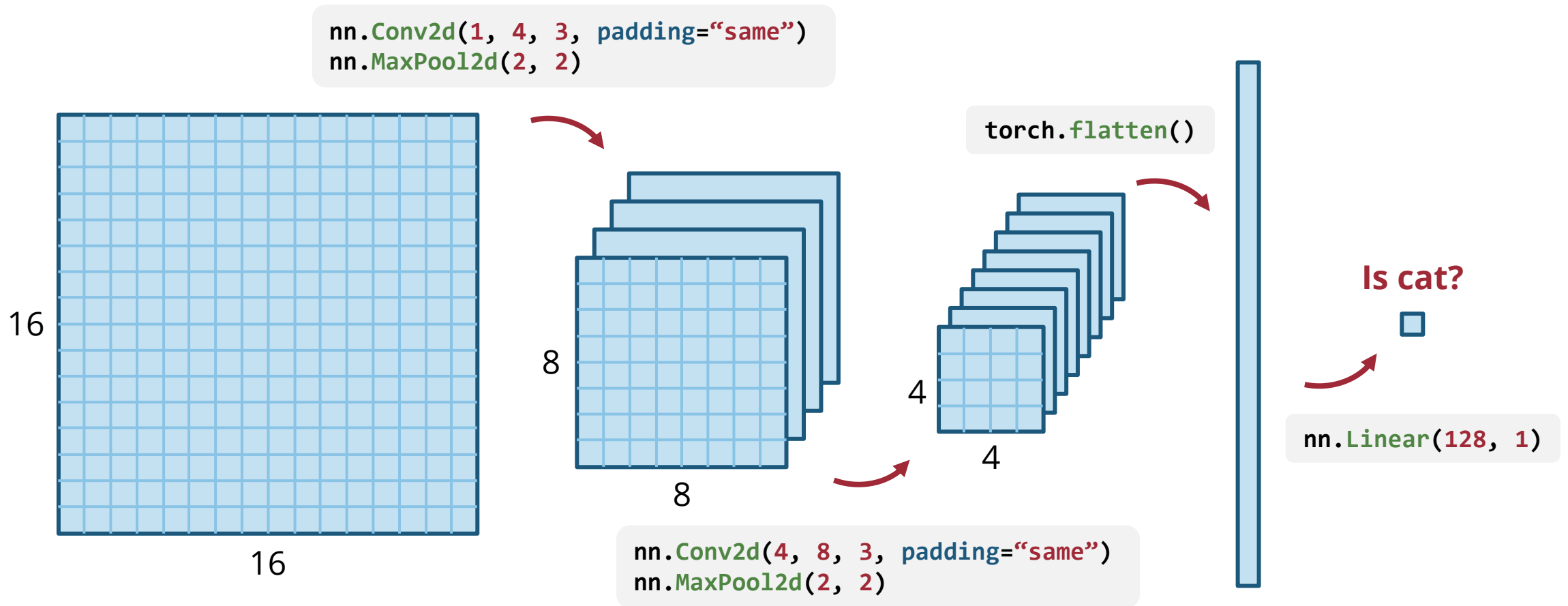
```
nn.MaxPool2d(kernel_size=2, stride=2)
```

# A Toy Example



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

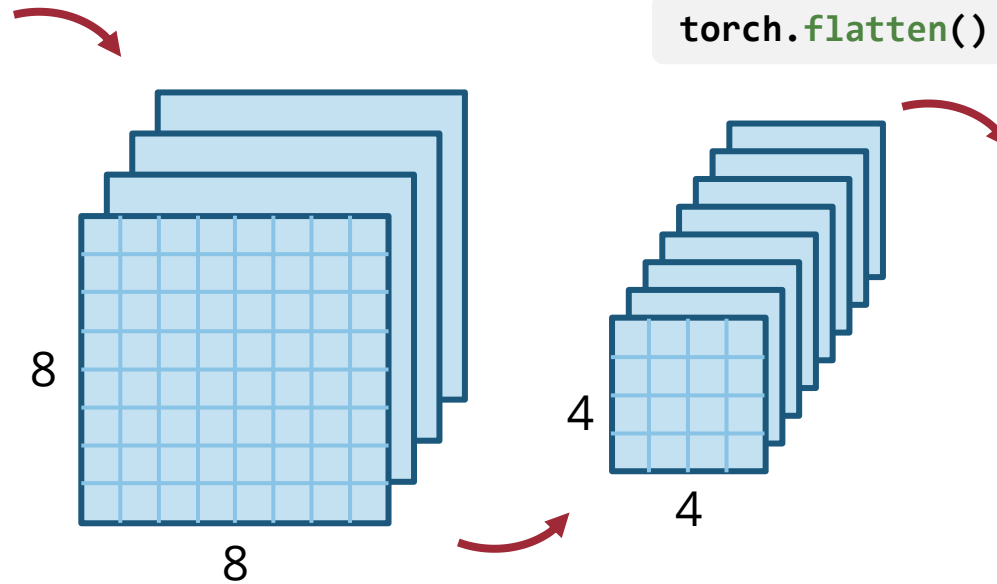
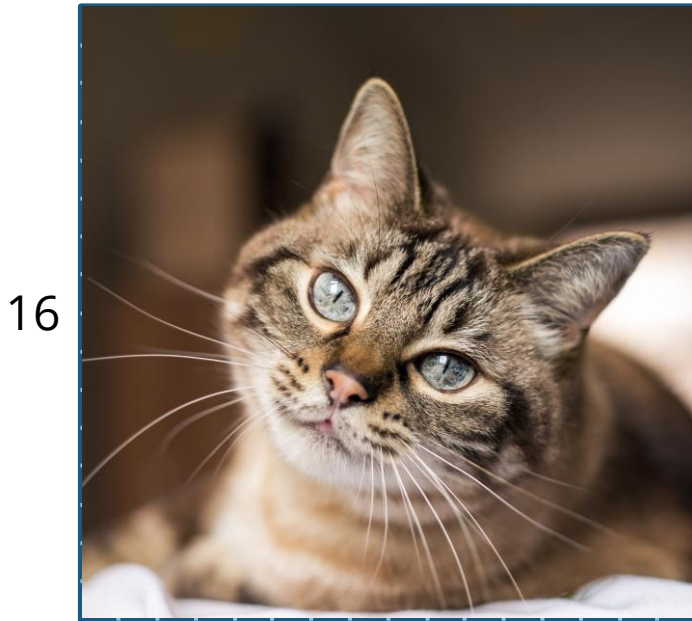
# A Toy Example





# 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)
```



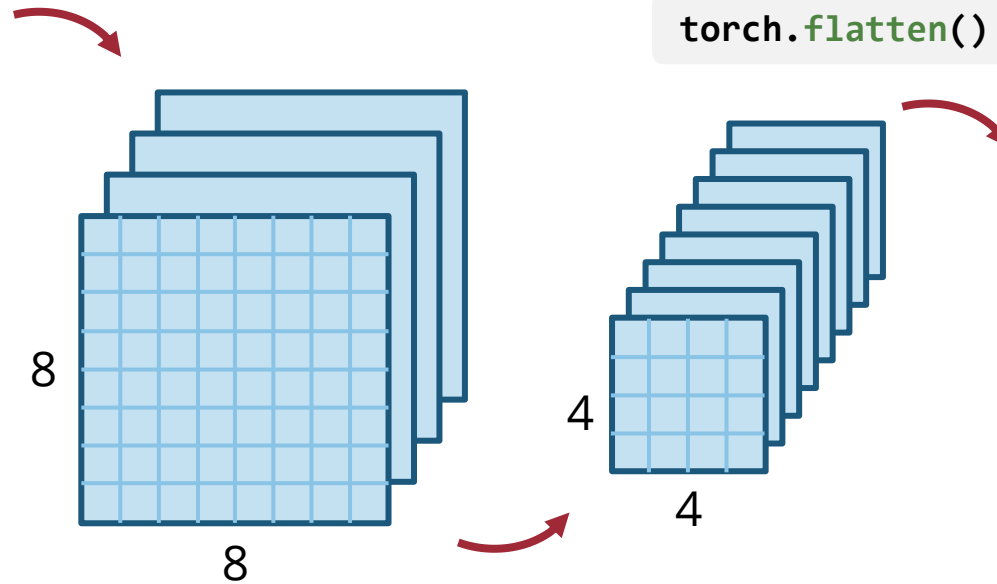
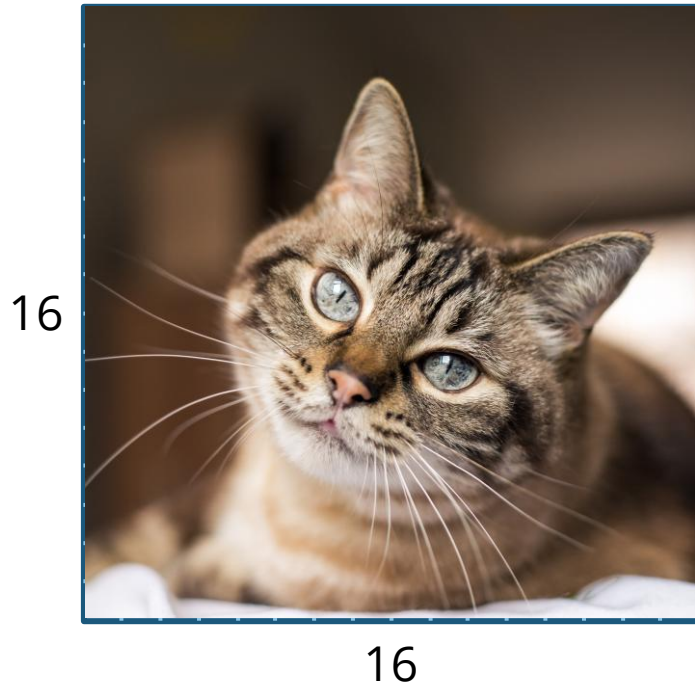
```
nn.Linear(128, 1)
```

Is cat?

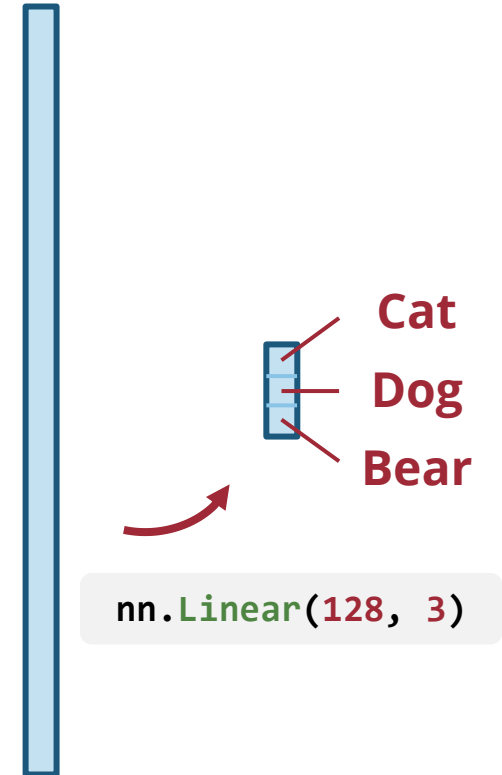


# A Toy Example

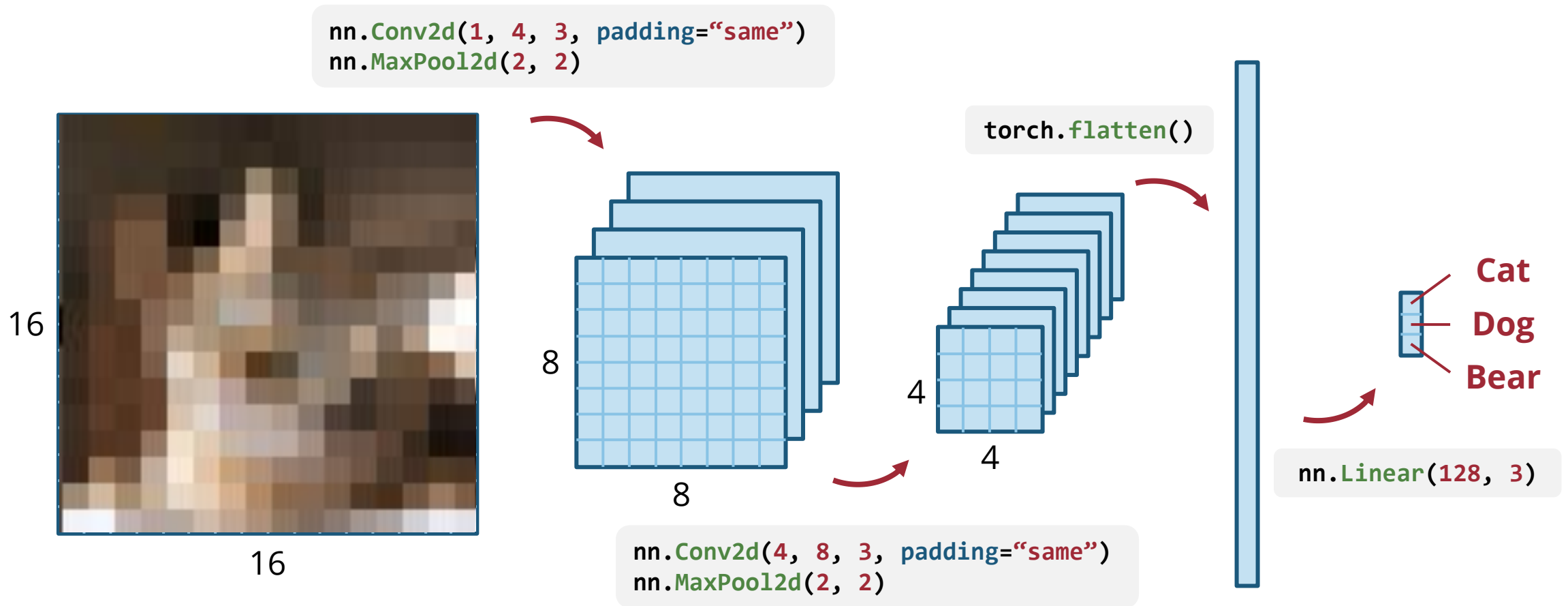
```
nn.Conv2d(1, 4, 3, padding="same")  
nn.MaxPool2d(2, 2)
```



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



# A Toy Example



# A Real Example

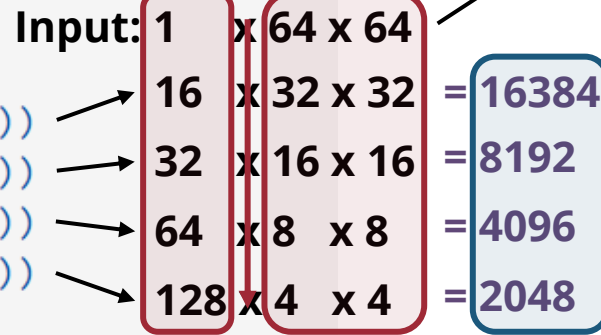
```

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)))
        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
    
```

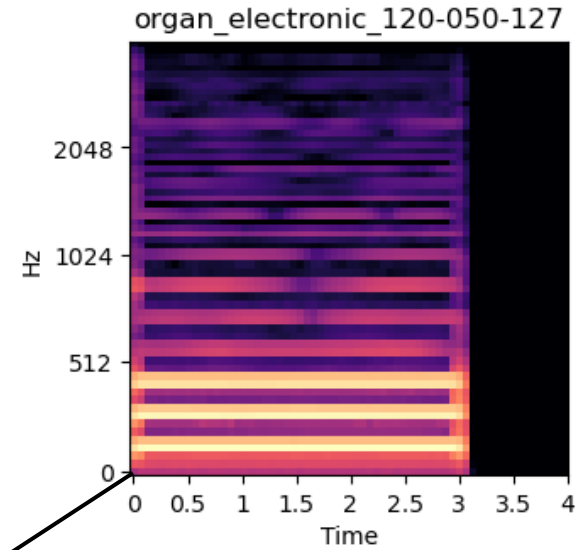
**Input channels**      **Output channels**

**Kernel size**



**More channels,  
lower resolution  
as we go deeper**

**Total number of  
features decrease  
as we go deeper**



# A Real Example

```
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)))
        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
```

**Input channels** (1), **Output channels** (16), **Kernel size** (3)

**How many parameters do we have in each layer?**

$(3 \times 3 \times 1 + 1) \times 16$	= 160
$(3 \times 3 \times 16 + 1) \times 32$	= 4640
$(3 \times 3 \times 32 + 1) \times 64$	= 18496
$(3 \times 3 \times 64 + 1) \times 128$	= 73856
$(2048 + 1) \times 11$	= 22539

# Benefits of CNNs

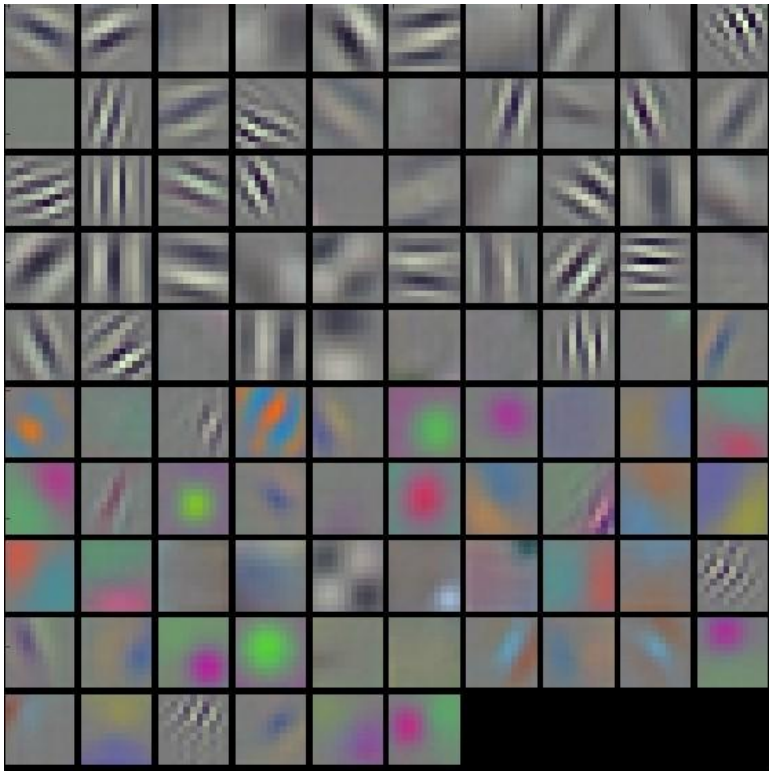
- Learn **local patterns**
- **Invariant to shifts**
  - Also called *translational invariance*
- **Reuse the learned filters** across
  - Different parts of the image
  - Across different images
- **Higher parameter-efficiency** against fully-connected neural network

What does a CNN Learn?

# Learned CNN Kernels in a Trained AlexNet

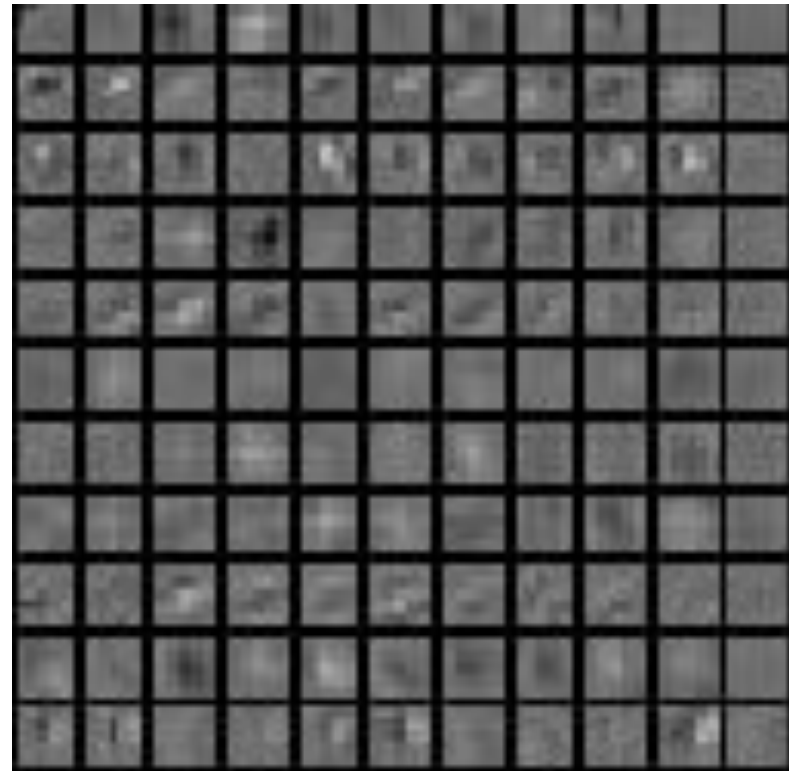
1st convolutional layer

11x11  
kernels



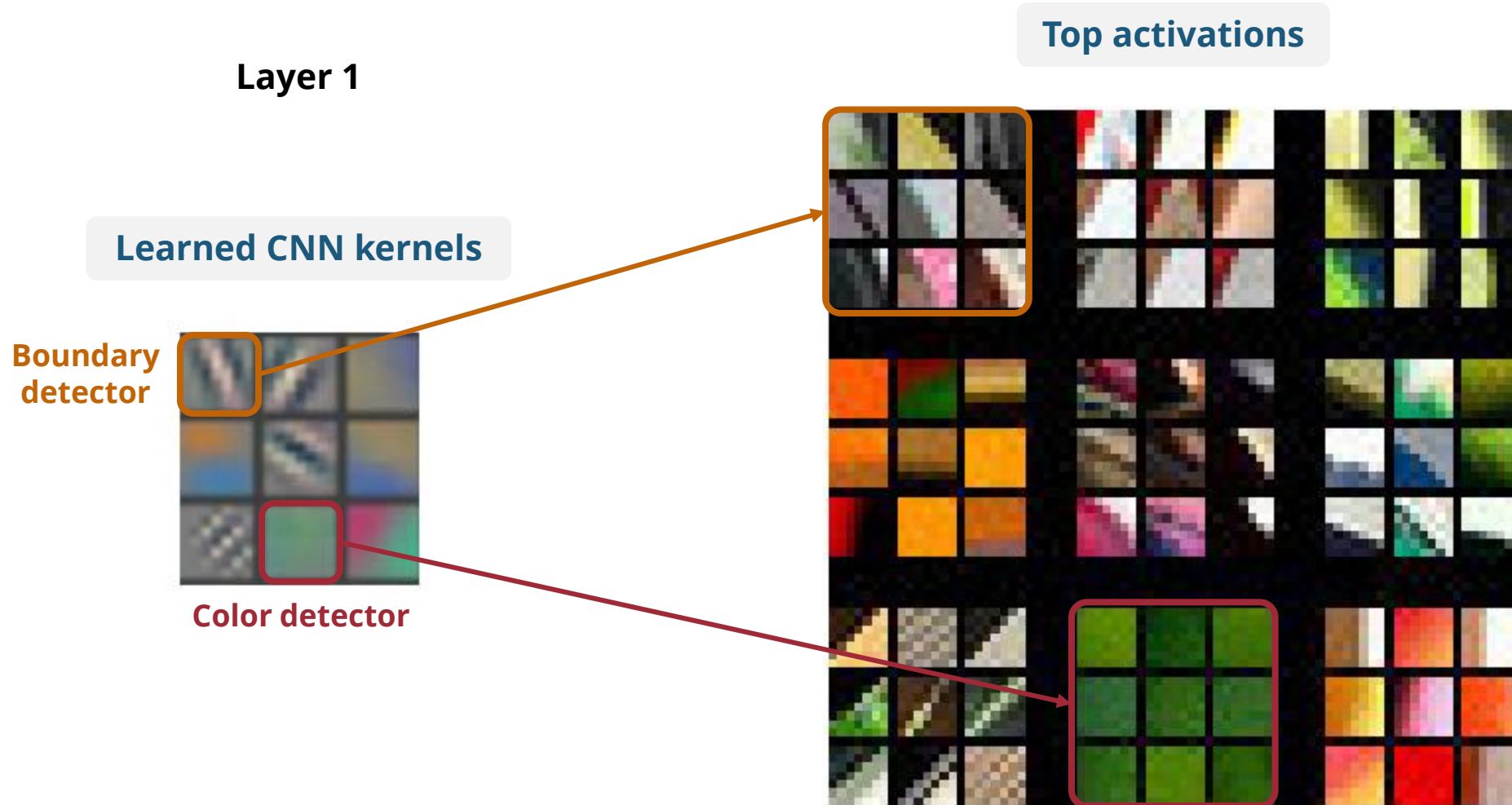
2nd convolutional layer

5x5  
kernels

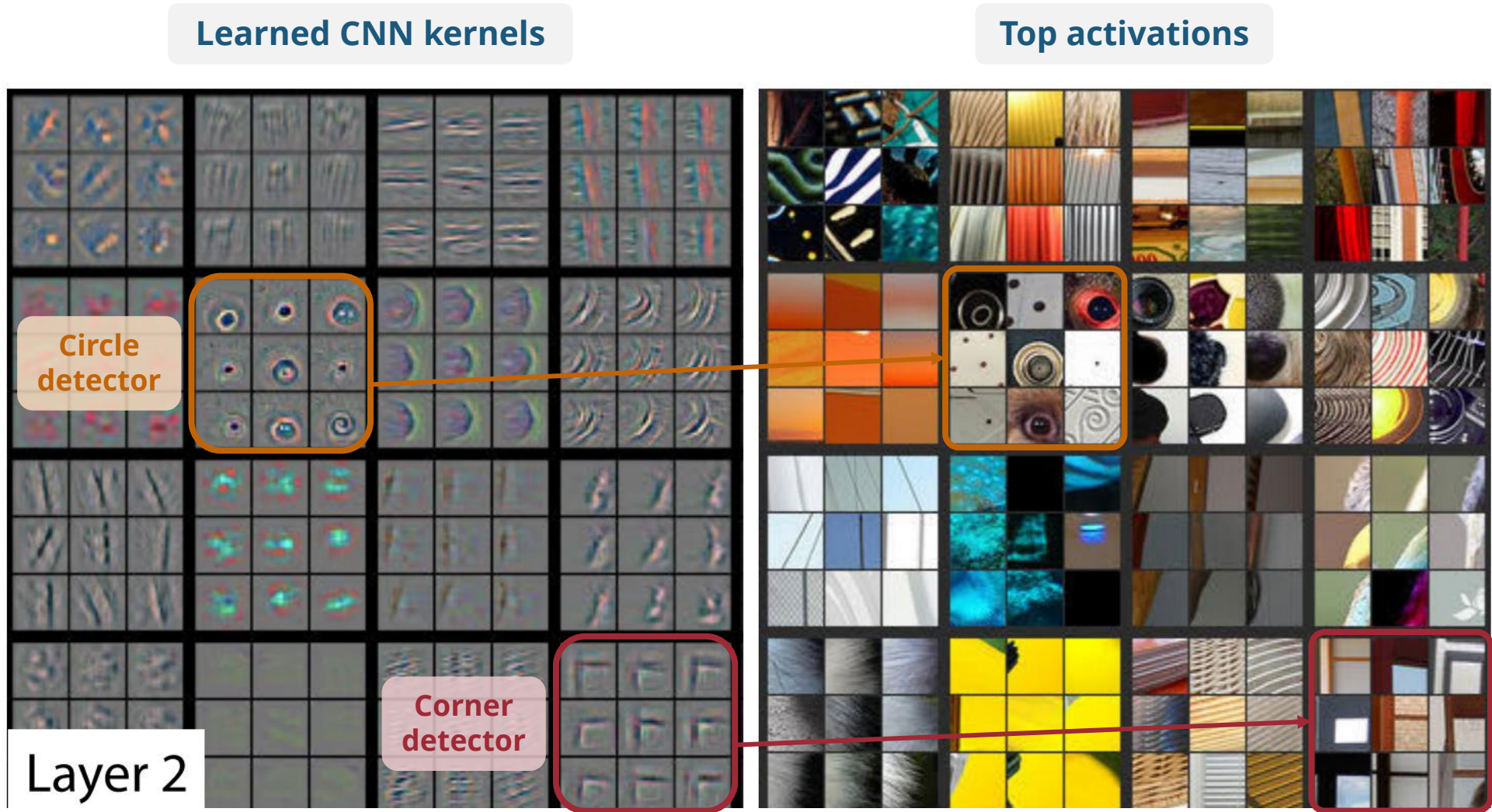




# Learned CNN Kernels in a Trained AlexNet



# Learned CNN Kernels in a Trained AlexNet

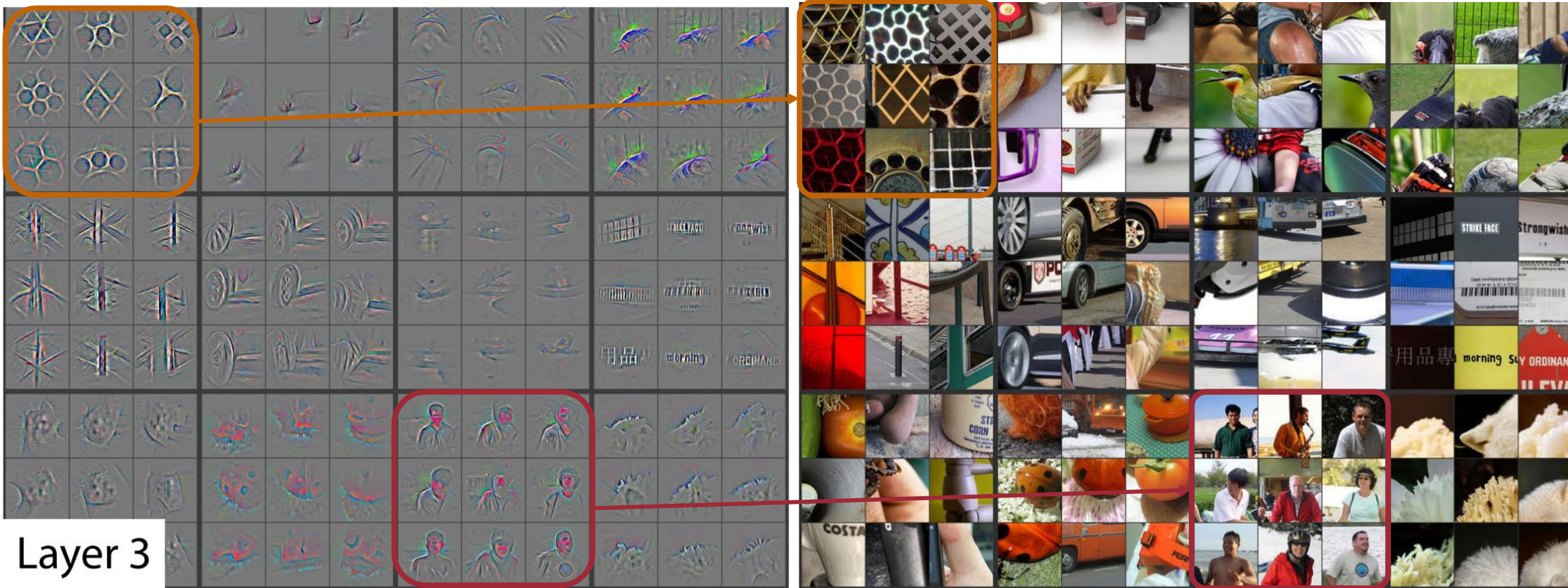


# Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

Top activations

Grid detector

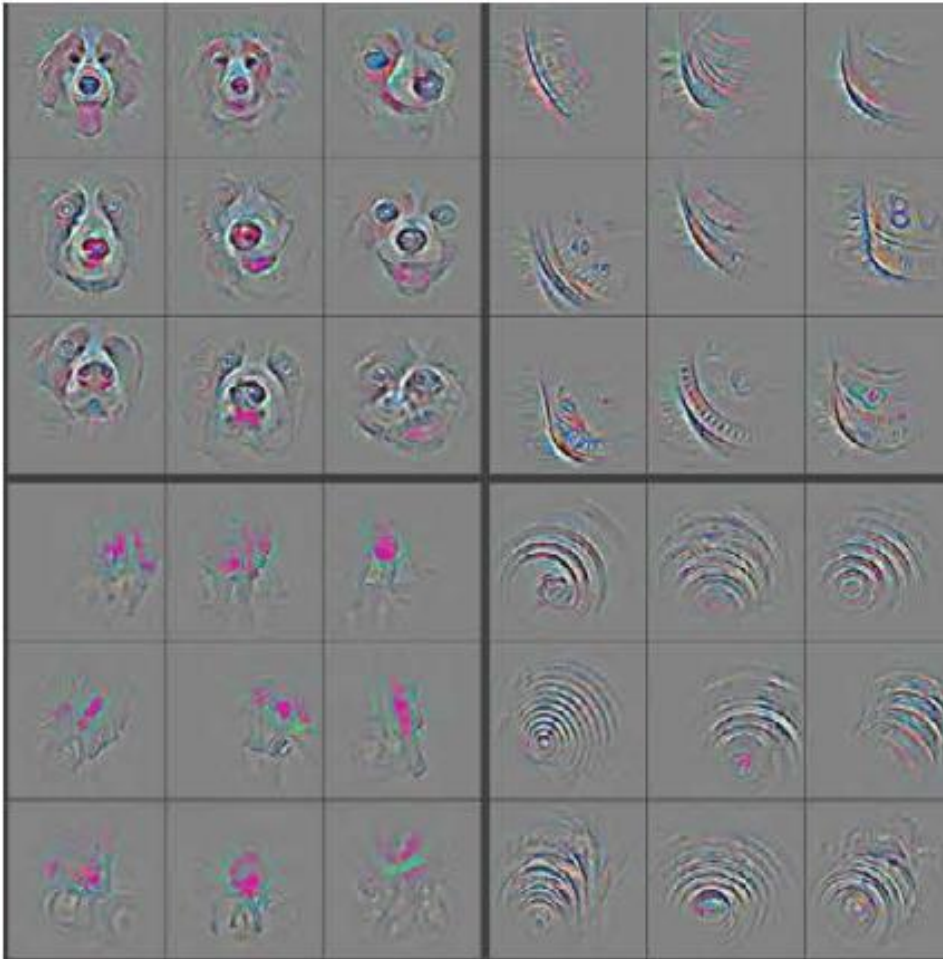


Layer 3

Human detector

# Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

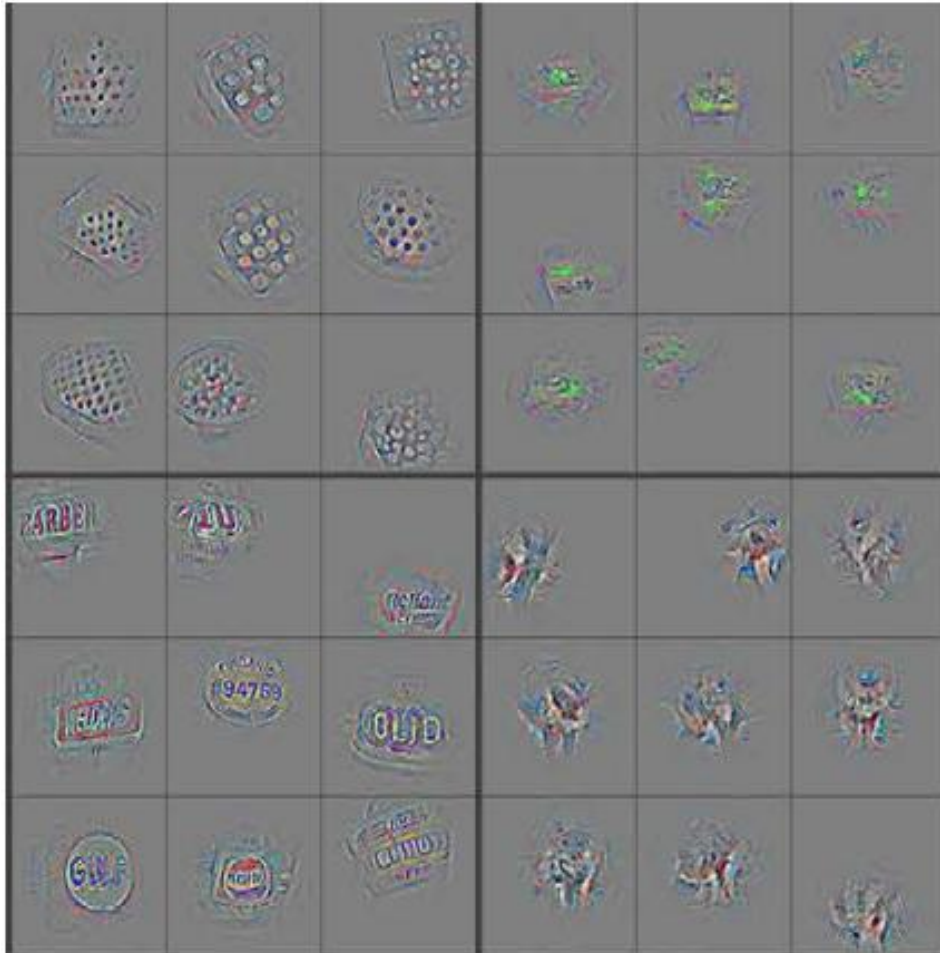


Top activations



# Learned CNN Kernels in a Trained AlexNet

Learned CNN kernels

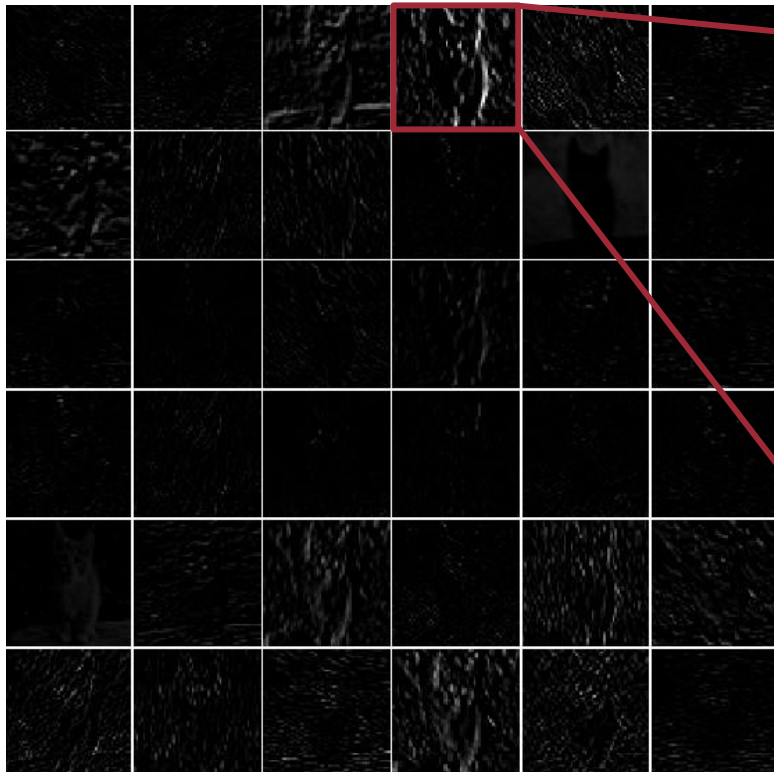


Top activations

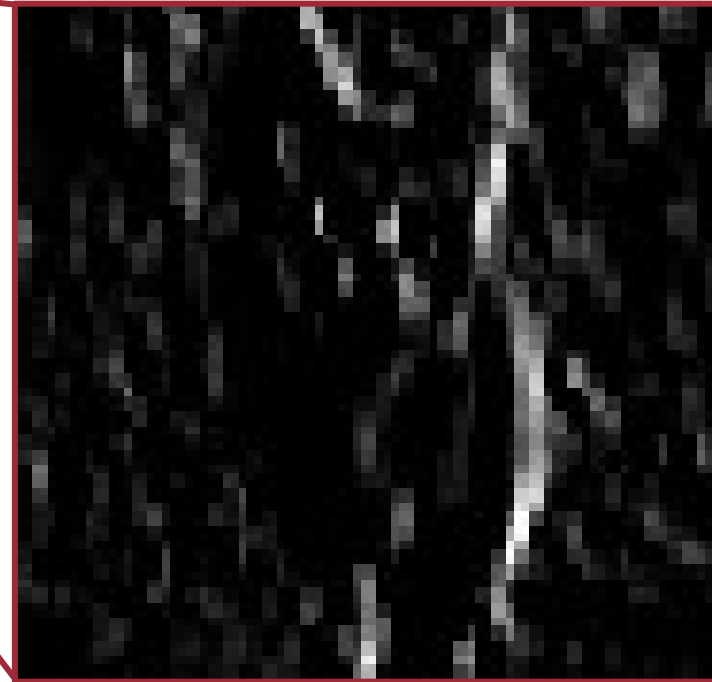


# Activations in a Trained AlexNet

1st convolutional layer

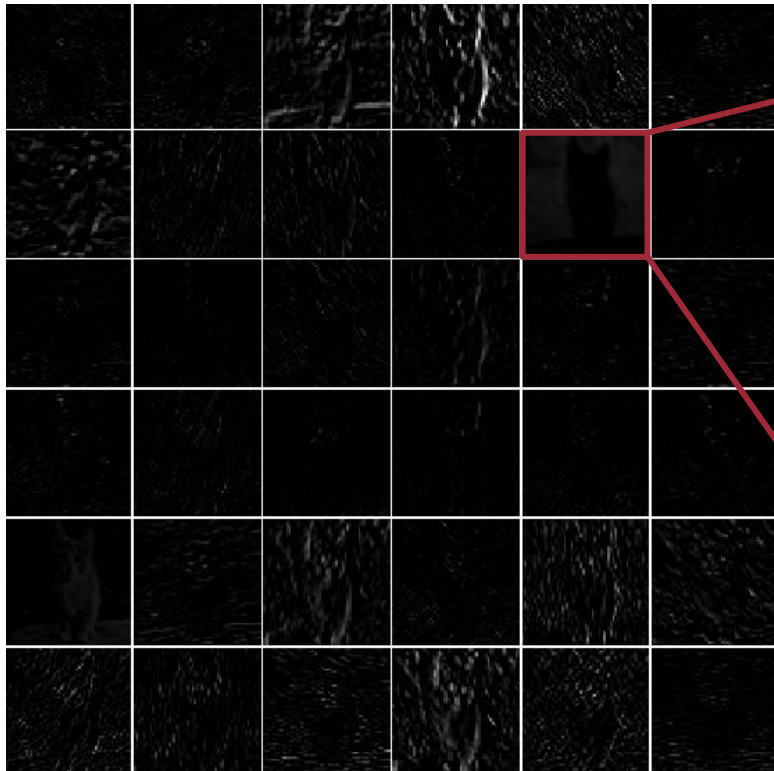


What's this!?

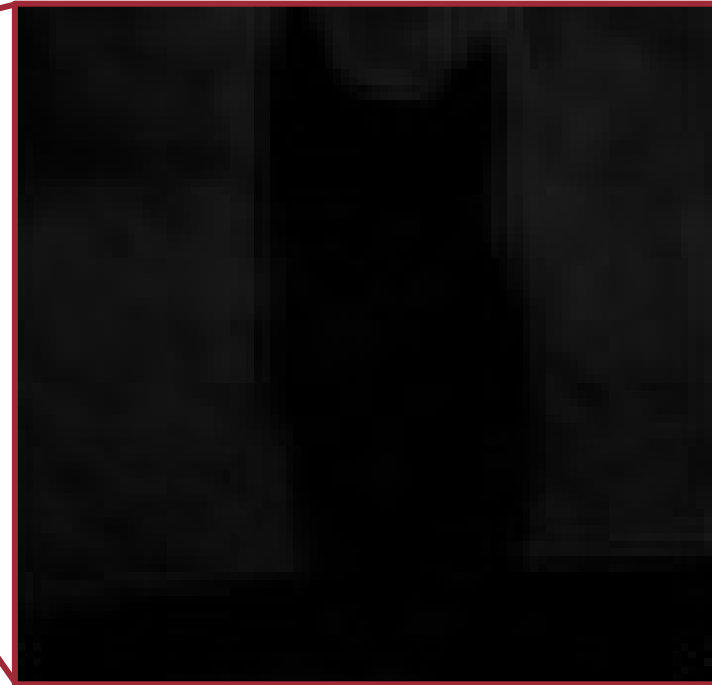


# Activations in a Trained AlexNet

1st convolutional layer

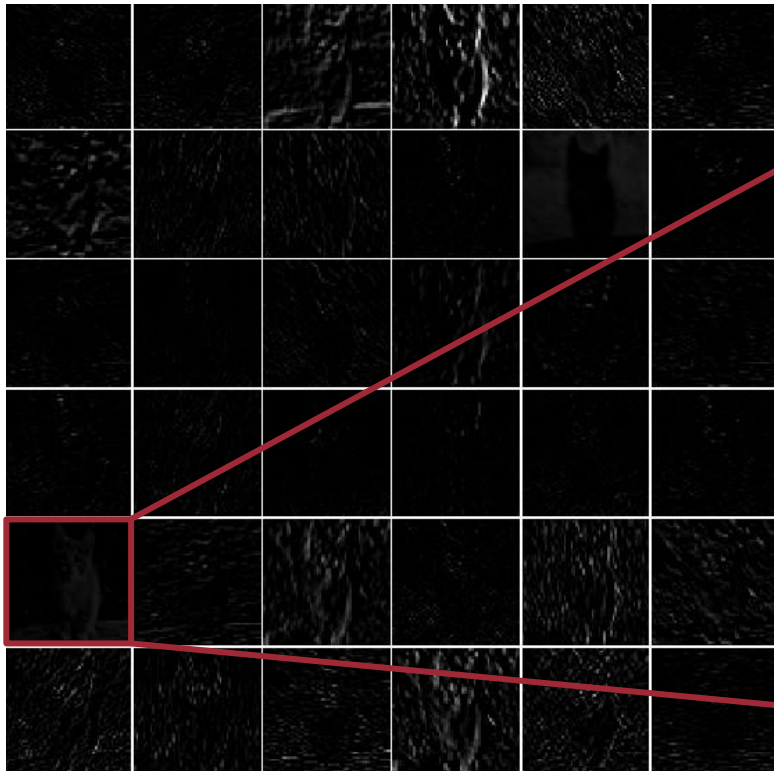


What's this!?



# Activations in a Trained AlexNet

1st convolutional layer



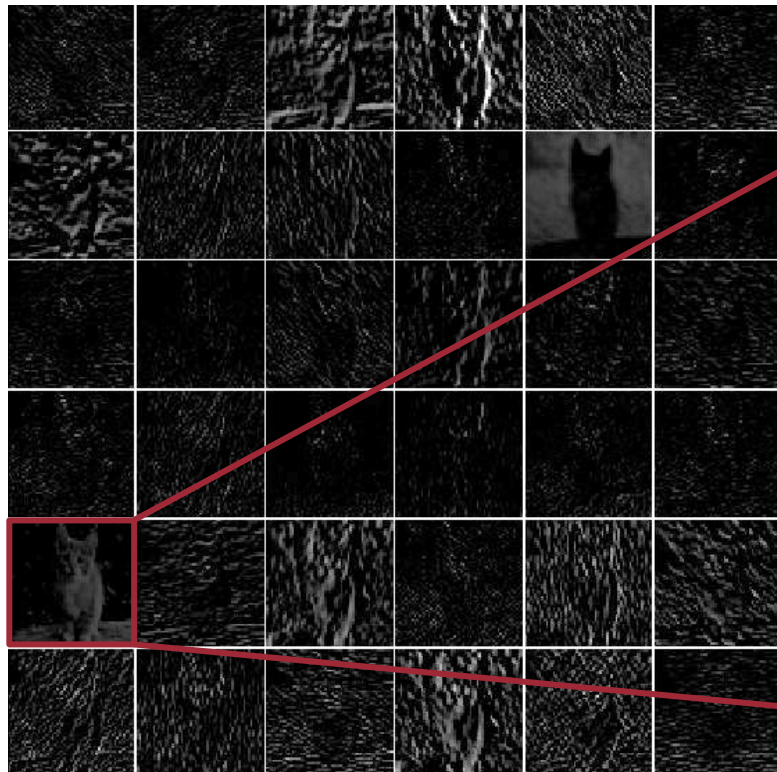
What's this!?





# Activations in a Trained AlexNet

1st convolutional layer



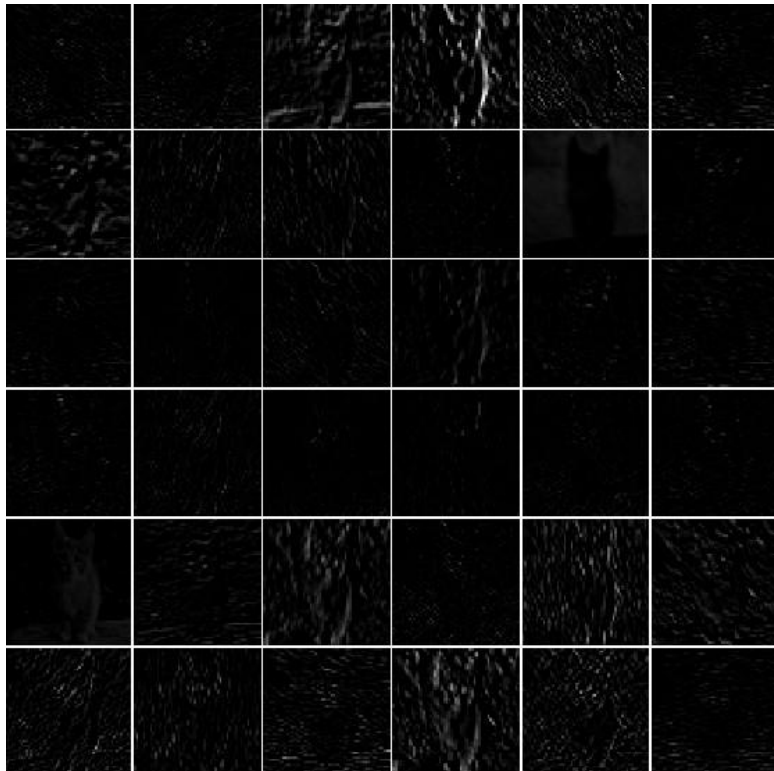
What's this!?



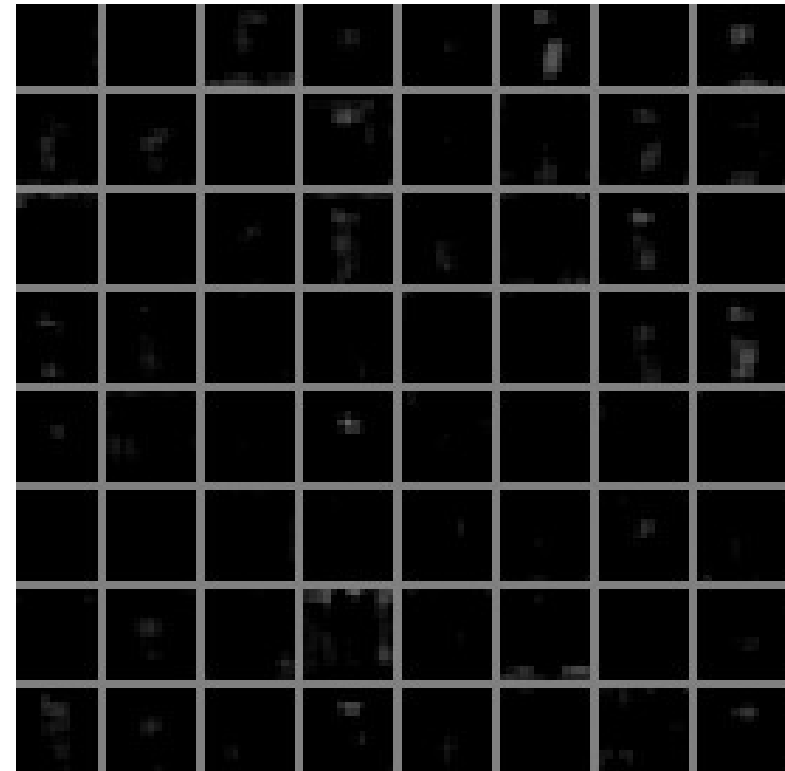
Brightness+  
Contrast+

# Activations in a Trained AlexNet

1st convolutional layer



5th convolutional layer

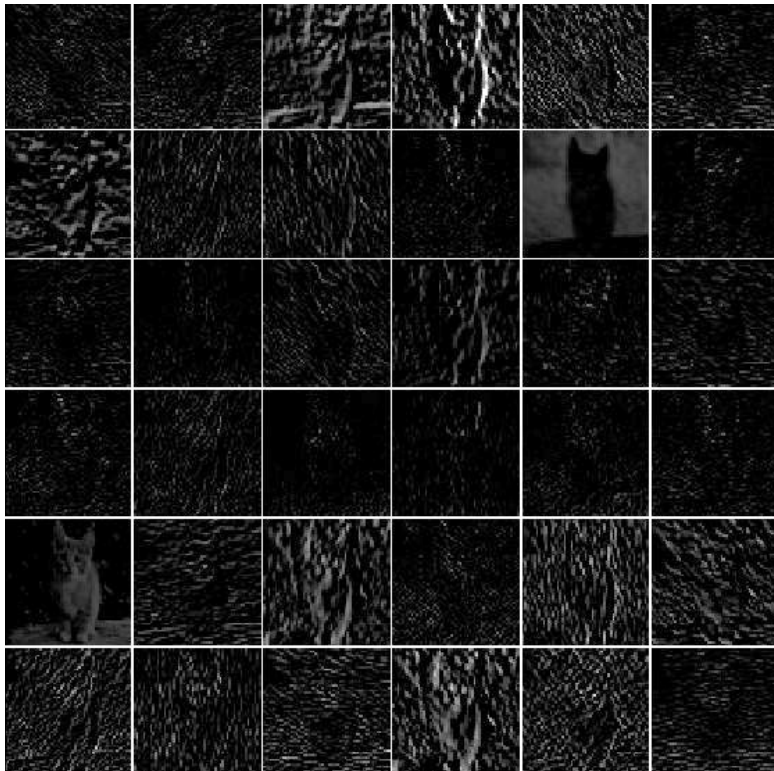


# Activations in a Trained AlexNet

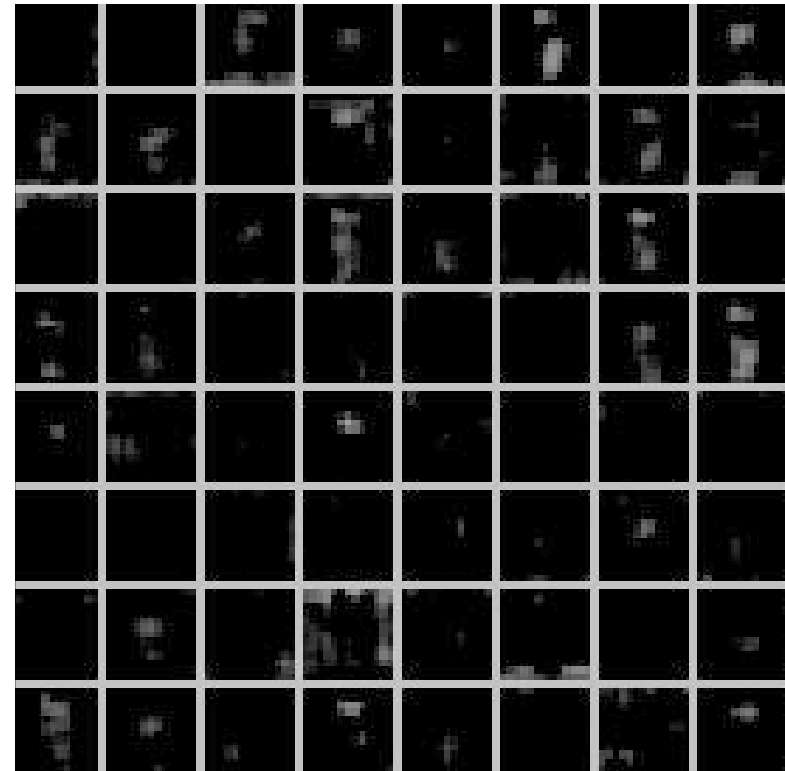
1st convolutional layer

5th convolutional layer

Brightness+  
Contrast+



Brightness+  
Contrast+



# What does a CNN Learn?

