

DEEP LEARNING

Lecture 5: Basics of Convolutional Neural Networks

Dr. Yang Lu

Department of Computer Science and Technology

luyang@xmu.edu.cn



CNN Applications

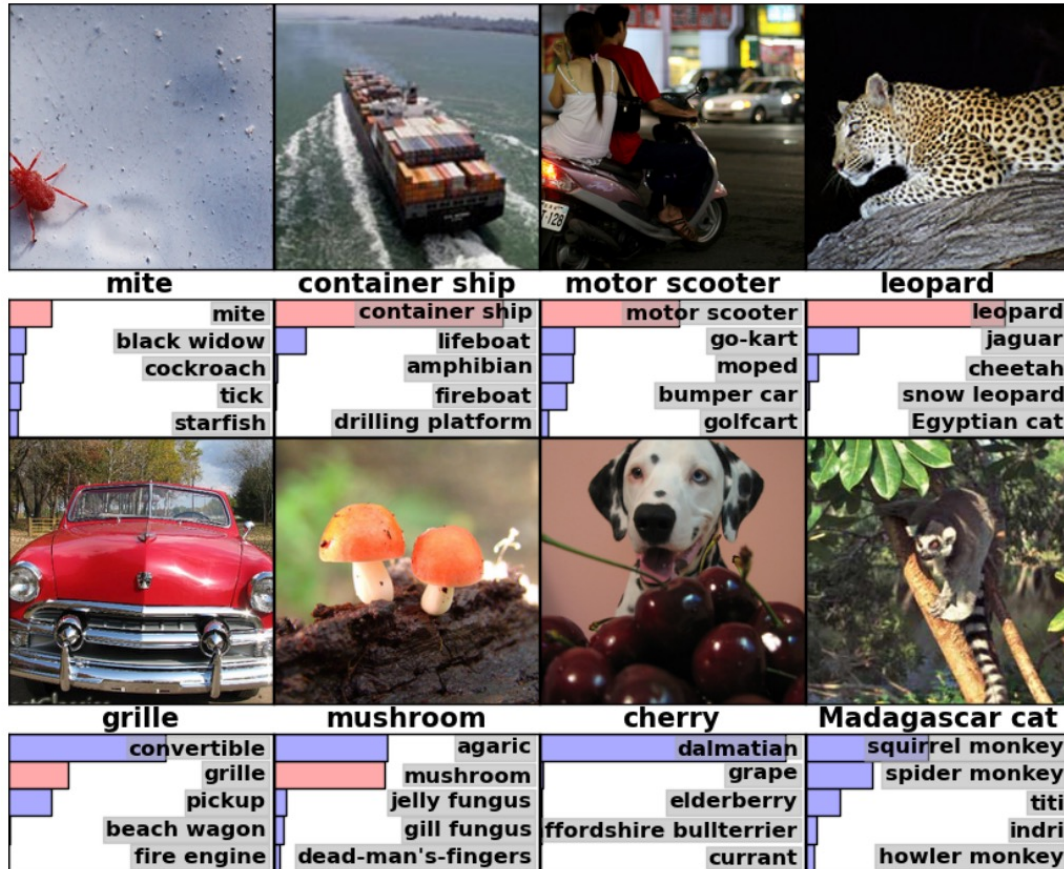


Image classification

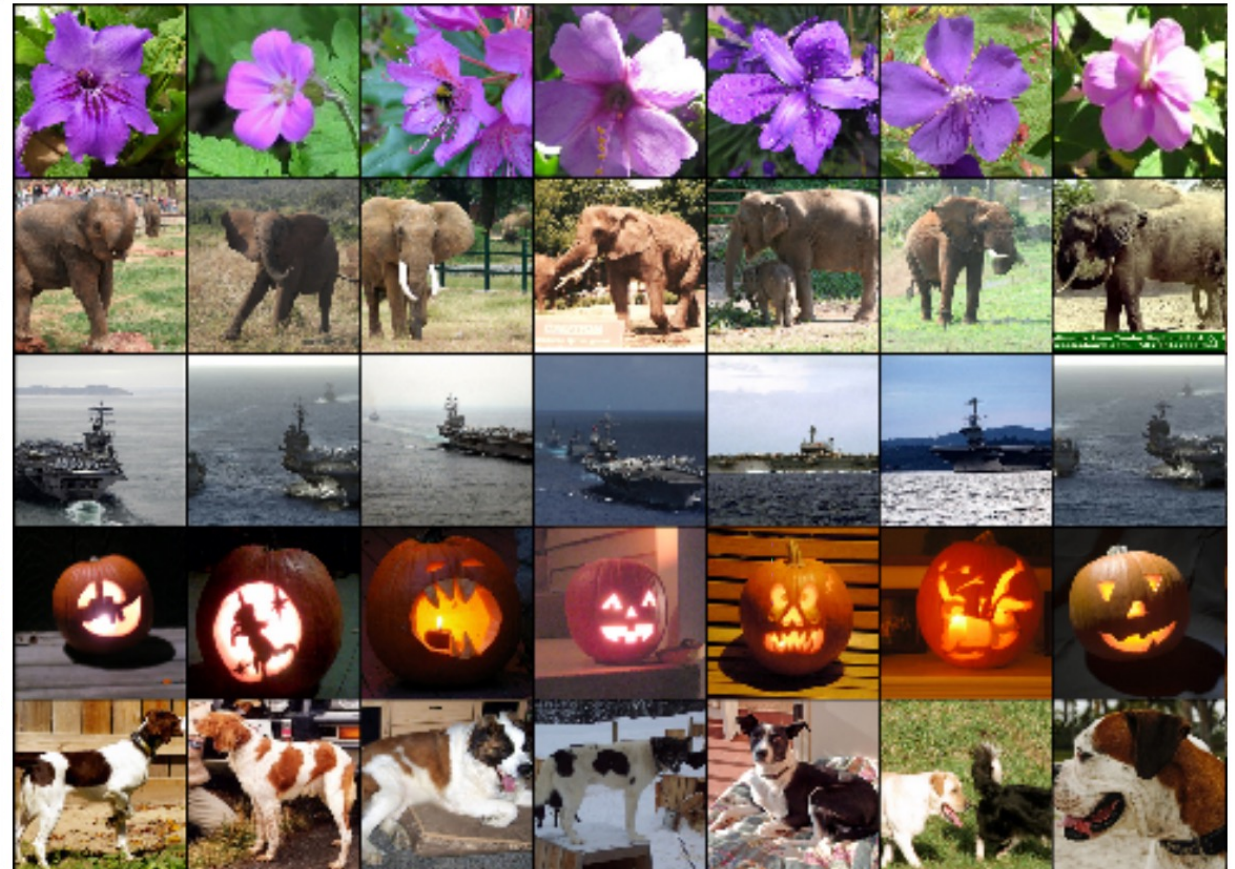


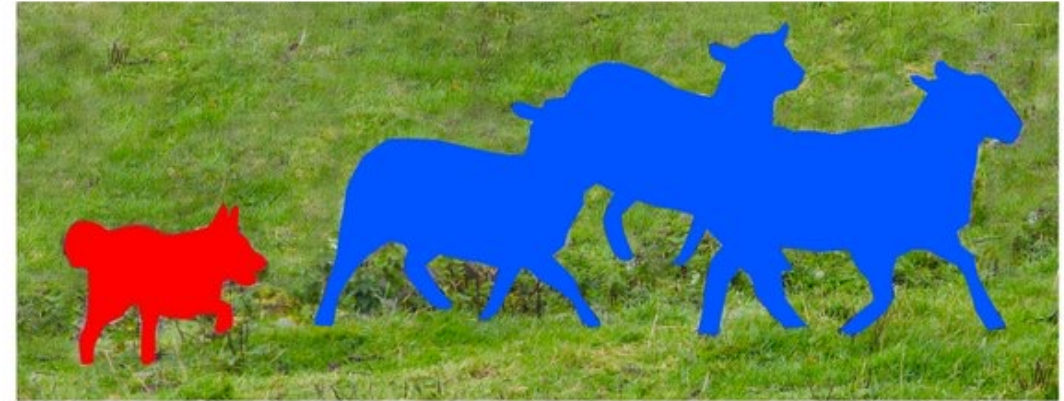
Image retrieval



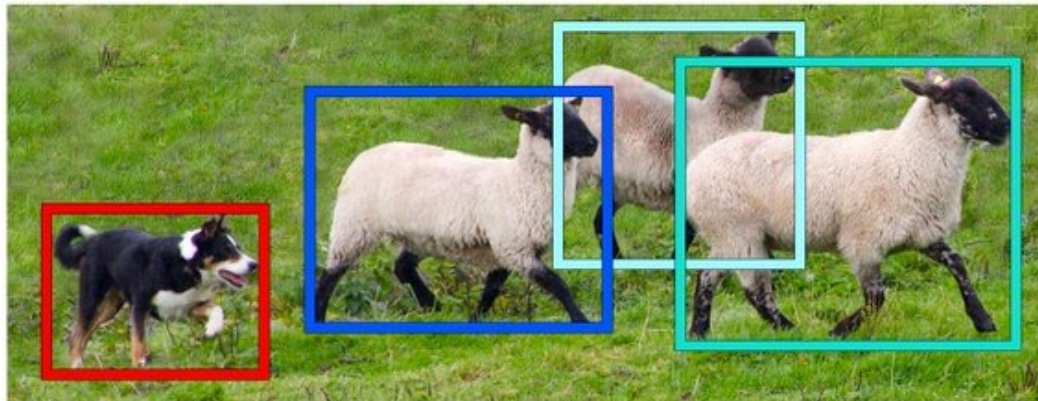
CNN Applications



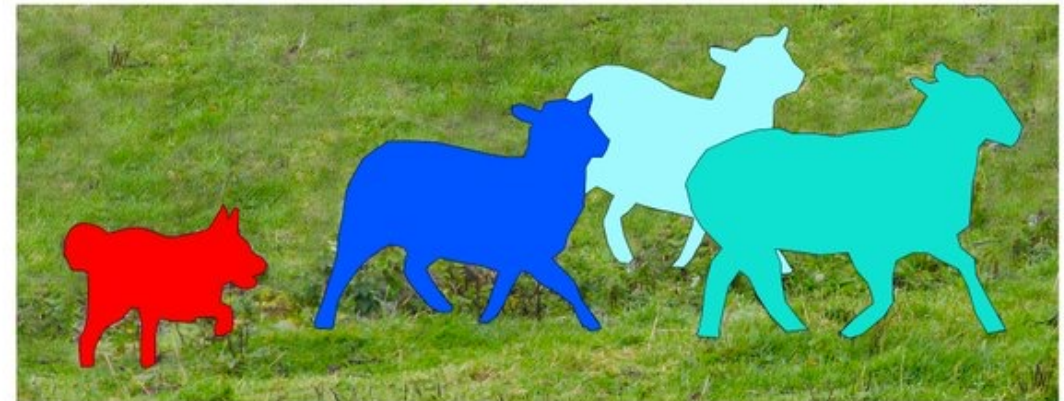
Image Recognition



Semantic Segmentation



Object Detection



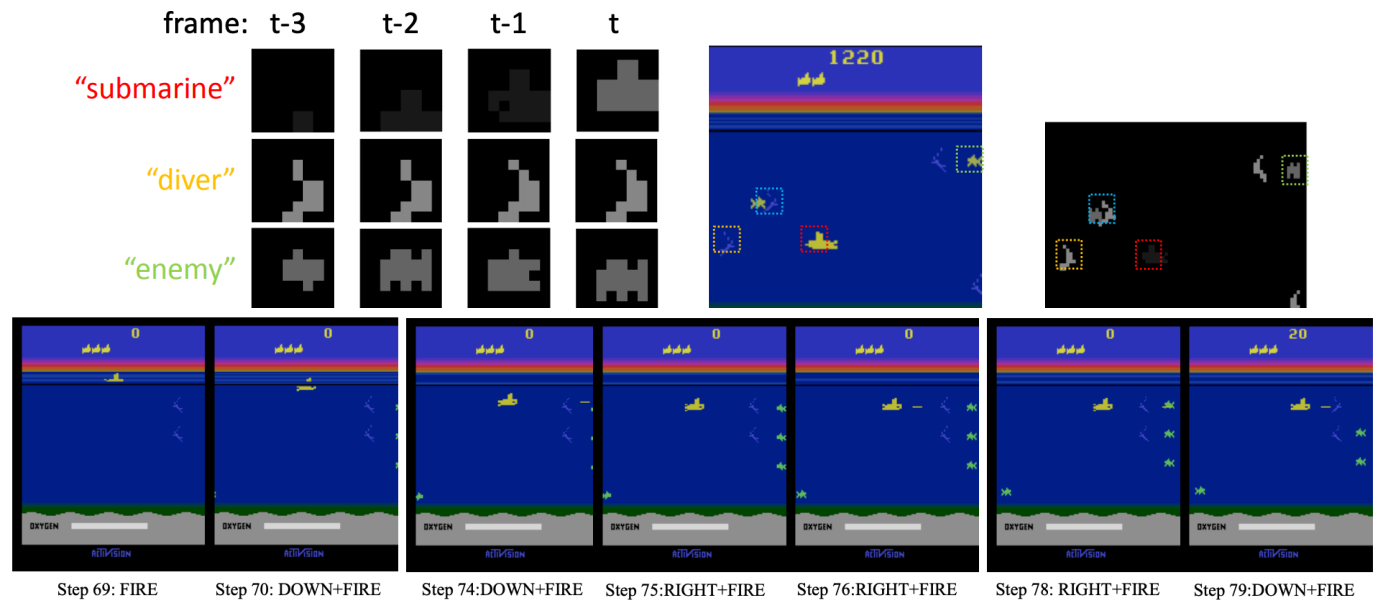
Instance Segmentation



CNN Applications



Pose estimation



Real-time Atari game play



CNN Applications













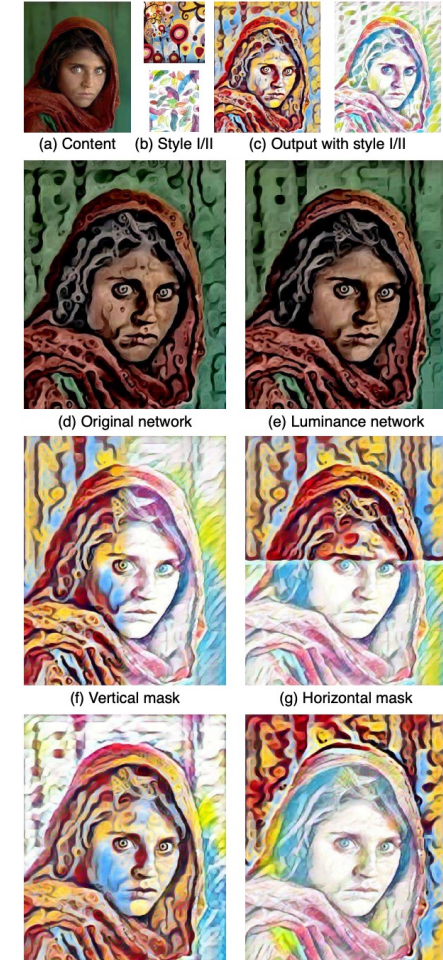
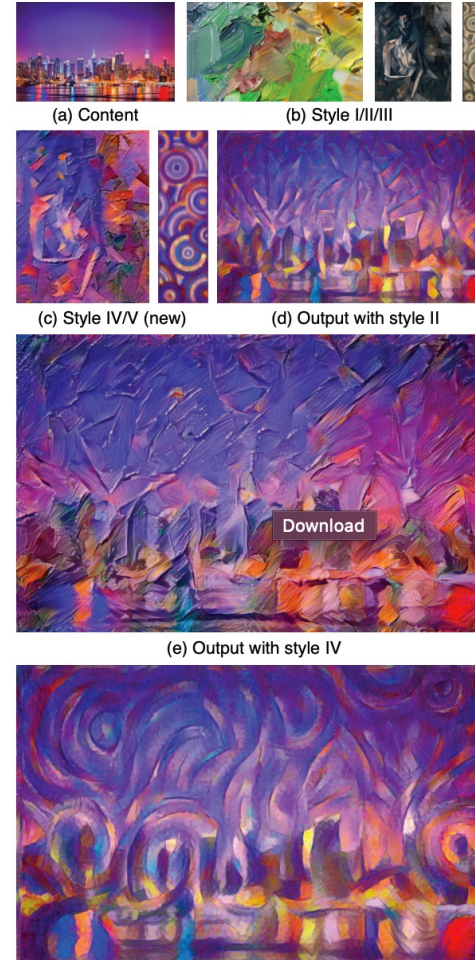
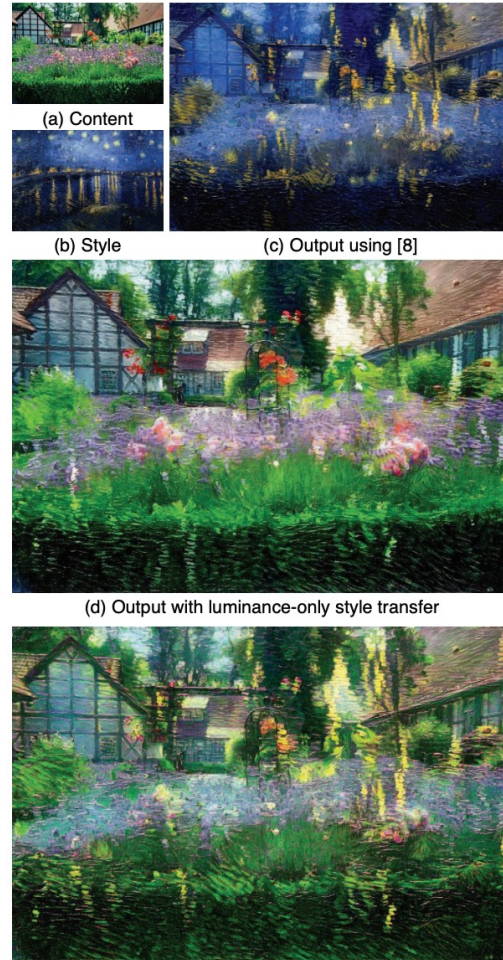
<p>A person riding a motorcycle on a dirt road.</p> 	<p>Two dogs play in the grass.</p> 	<p>A skateboarder does a trick on a ramp.</p> 	<p>A dog is jumping to catch a frisbee.</p> 
<p>A group of young people playing a game of frisbee.</p> 	<p>Two hockey players are fighting over the puck.</p> 	<p>A little girl in a pink hat is blowing bubbles.</p> 	<p>A refrigerator filled with lots of food and drinks.</p> 
<p>A herd of elephants walking across a dry grass field.</p> 	<p>A close up of a cat laying on a couch.</p> 	<p>A red motorcycle parked on the side of the road.</p> 	<p>A yellow school bus parked in a parking lot.</p> 
Describes without errors	Describes with minor errors	Somewhat related to the image	Unrelated to the image

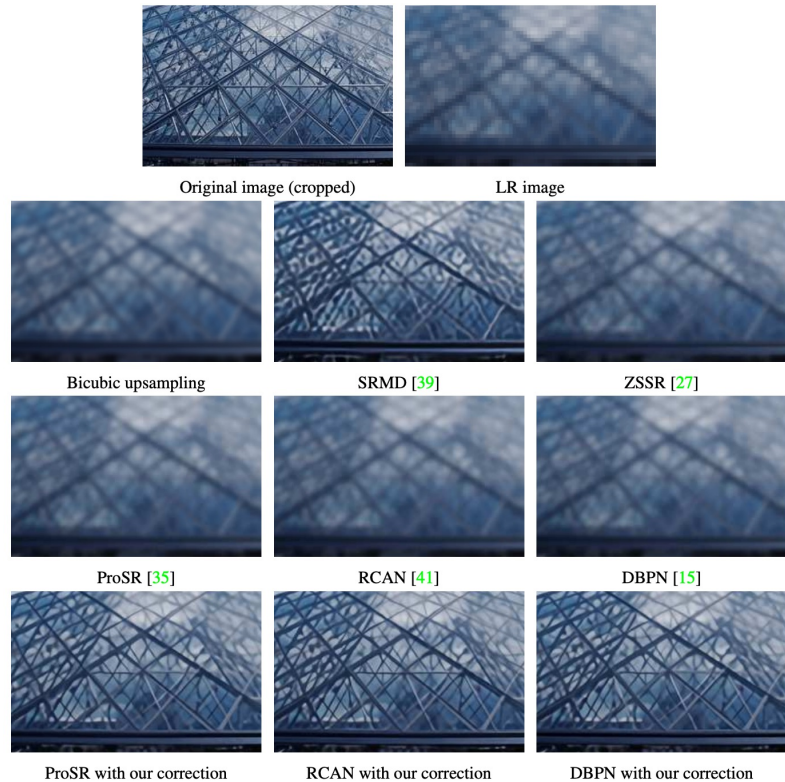
Image captioning

CNN Applications

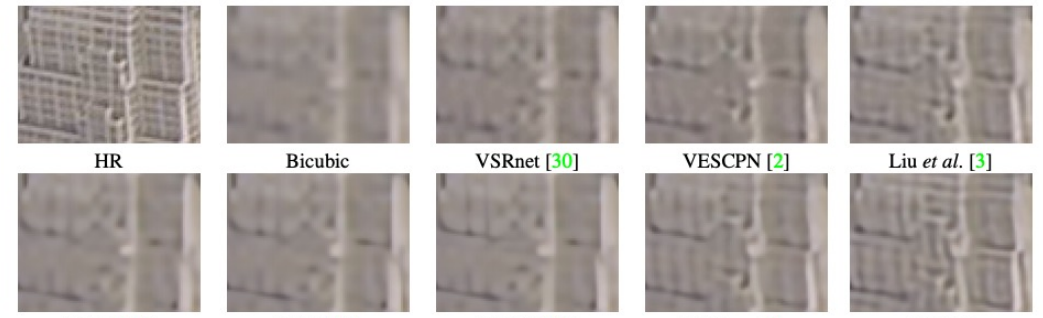


Style transfer

CNN Applications



City/BI



DBPN [22]

RDN [21]

RCAN [23]

TOFlow [5]

TDAN



Walk/BD



HR

Bicubic

DRVSR [4]

DUF [1]

TDAN

Image super-resolution

Video super-resolution



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



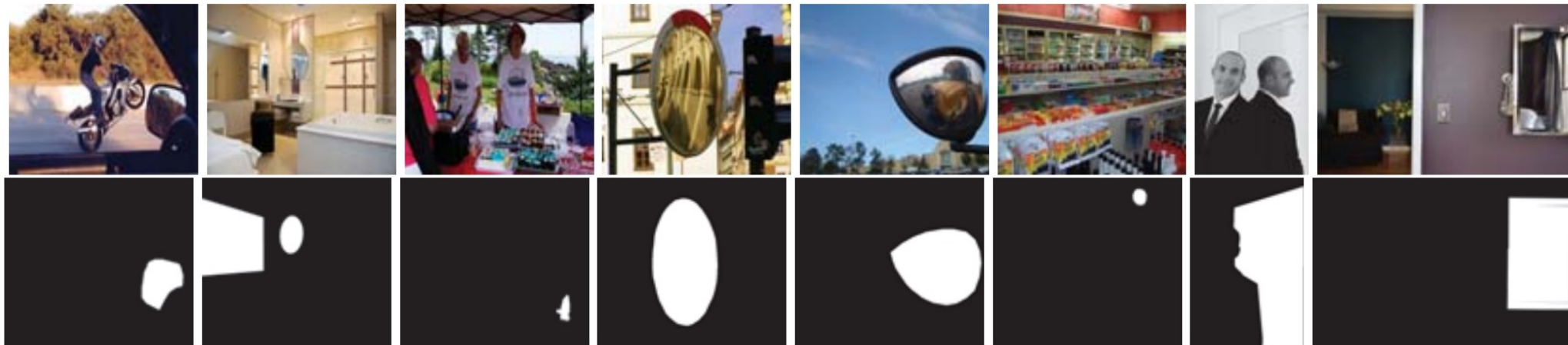
廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

CNN Applications



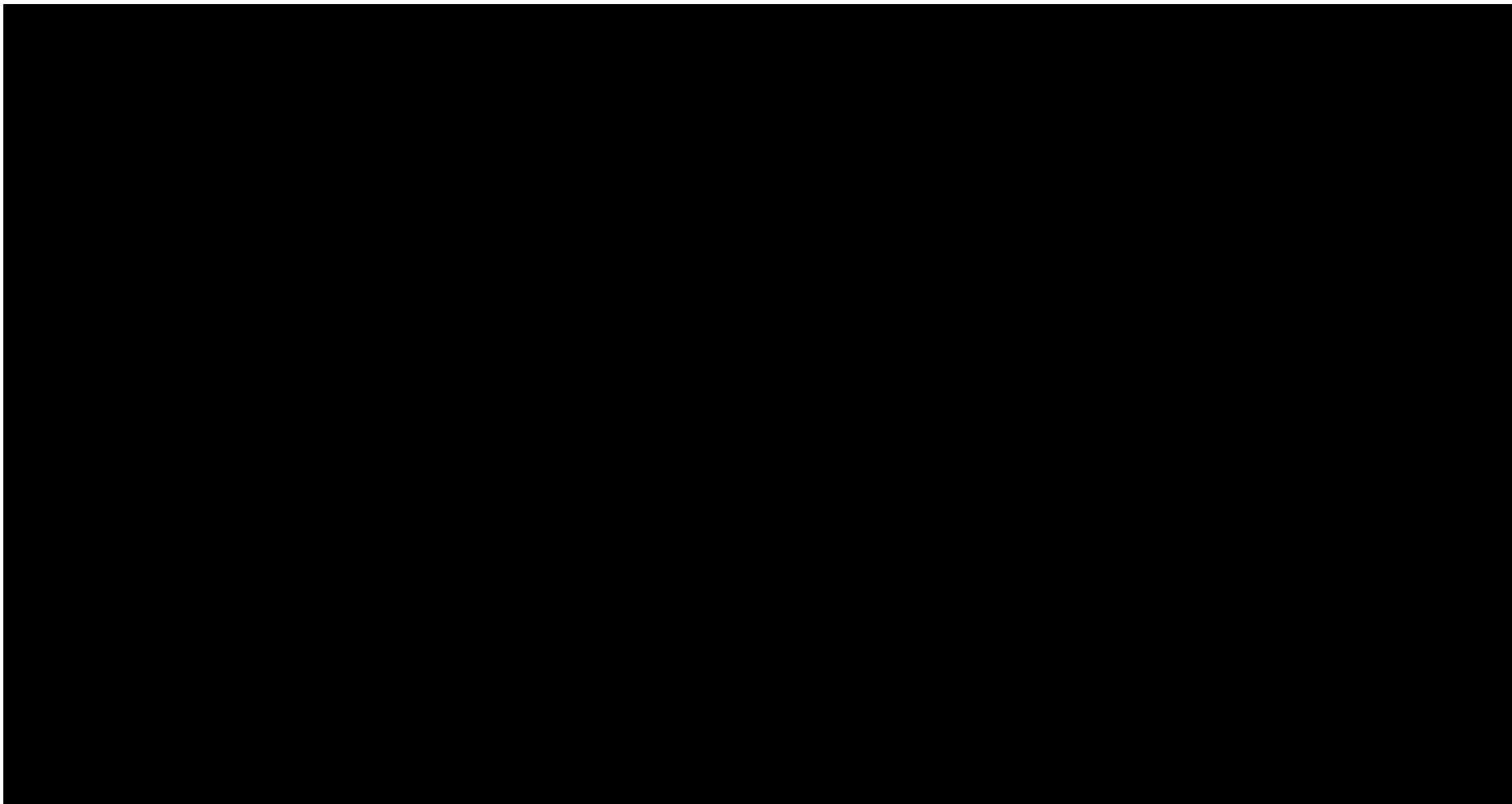
Rain and fog removal



Mirror detection



CNN Applications



Video frame interpolation



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)

Video source: <https://www.youtube.com/watch?v=5qAiffYFJhc>

Paper source: <https://arxiv.org/abs/2103.16206>



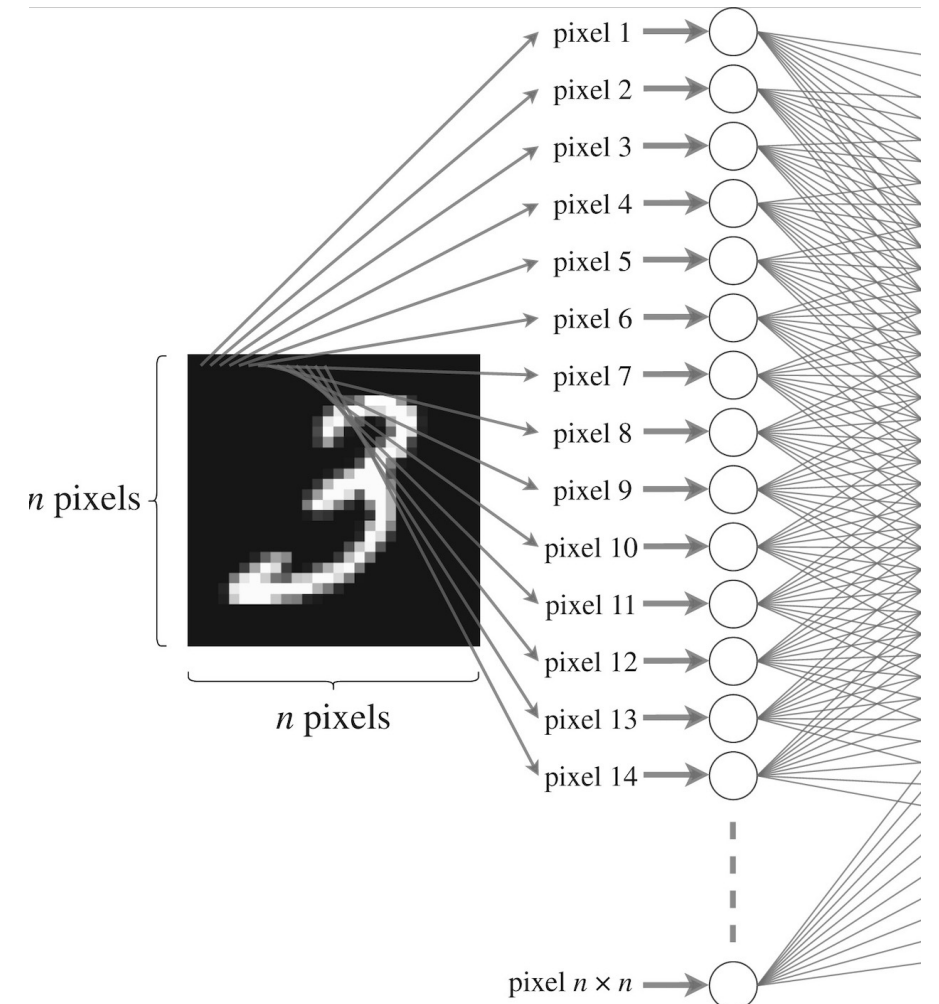
廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

Convolutional Neural Networks

- Recall in Lecture 2, we **vectorize** an image as the input of a neural network.
- What is the problem here?

Can maintain 2D structure through the whole network?



Convolutional Neural Networks

- Convolutional neural networks (CNNs) are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers.
- Use the non-vectorized image as input with a 2D weight, which are called a filter or a kernel.
- We call the hidden outputs in CNNs feature map.

Convolutional Neural Networks

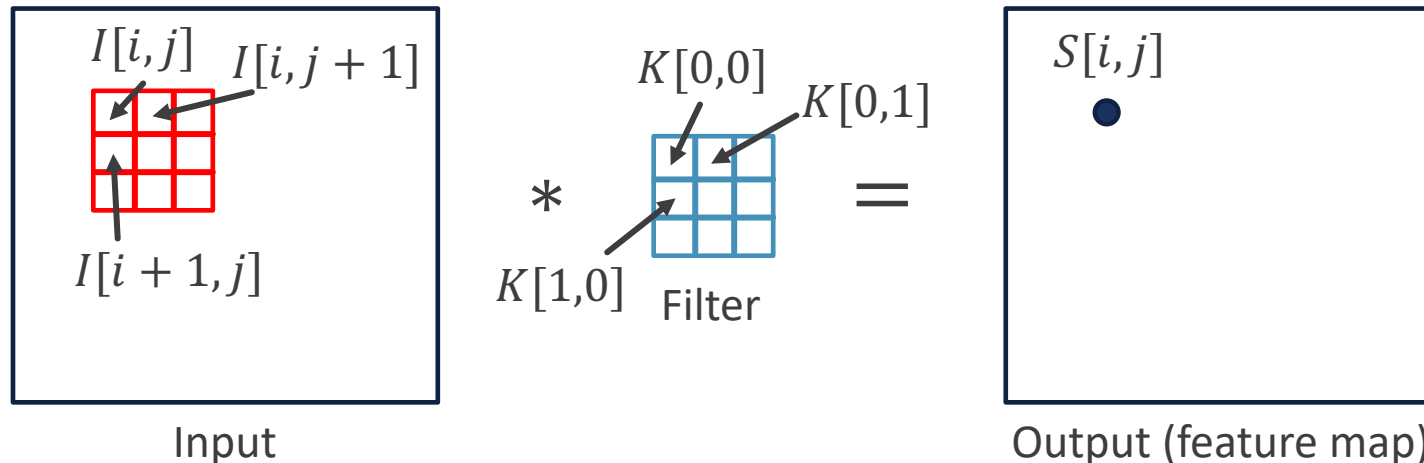
- The operation $*$ between the input image I and filter K to produce a new image S is called **convolution**, which is defined as:

$$S[i, j] = (I * K)[i, j] = \sum_m \sum_n I[i + m, j + n] K[m, n].$$

pixel position

offset

- MLP: input x , weight W , output h .
- CNN: input I , weight K , output S .



Convolution

Input: 3×3

1	2	3
4	5	6
7	8	9

×

Filter: 2×2

10	20
30	40

=

Output: 2×2

370	

$$1 \times 10 + 2 \times 20 + 4 \times 30 + 5 \times 40 = 370$$



Convolution

Input: 3×3

1	2	3
4	5	6
7	8	9

×

Filter: 2×2

10	20
30	40

=

Output: 2×2

370	470

$$2 \times 10 + 3 \times 20 + 5 \times 30 + 6 \times 40 = 470$$



Convolution

Input: 3×3

1	2	3
4	5	6
7	8	9



Filter: 2×2

10	20
30	40



Output: 2×2

370	470
670	

$$4 \times 10 + 5 \times 20 + 7 \times 30 + 8 \times 40 = 670$$



Convolution

Input: 3×3

1	2	3
4	5	6
7	8	9

×

Filter: 2×2

10	20
30	40

=

Output: 2×2

370	470
670	770

Now the problem: the size of the new image after convolution is shrunk.

$$5 \times 10 + 6 \times 20 + 8 \times 30 + 9 \times 40 = 770$$

Padding

- In order to keep the dimension of input and output matrix the same, we add padding.

Input: 3×3 + 1×1 padding

0	0	0	0	0
0	1	2	3	0
0	4	5	6	0
0	7	8	9	0
0	0	0	0	0

×

Filter: 3×3

10	20	30
40	50	60
30	20	10

=

Output: 3×3

300		

$$\begin{aligned} &0 \times 10 + 0 \times 20 + 0 \times 30 + \\ &0 \times 40 + 1 \times 50 + 2 \times 60 + \\ &0 \times 30 + 4 \times 20 + 5 \times 10 = \\ &300 \end{aligned}$$

Padding

- In order to keep the dimension of input and output matrix the same, we add padding.

Input: 3×3 + 1×1 padding

0	0	0	0	0
0	1	2	3	0
0	4	5	6	0
0	7	8	9	0
0	0	0	0	0

×

Filter: 3×3

10	20	30
40	50	60
30	20	10

=

Output: 3×3

300	600	

$$\begin{aligned} &0 \times 10 + 0 \times 20 + 0 \times 30 + \\ &1 \times 40 + 2 \times 50 + 3 \times 60 + \\ &4 \times 30 + 5 \times 20 + 6 \times 10 = \\ &600 \end{aligned}$$

Padding

- In order to keep the dimension of input and output matrix the same, we add padding.

Input: 3×3 + 1×1 padding

0	0	0	0	0
0	1	2	3	0
0	4	5	6	0
0	7	8	9	0
0	0	0	0	0

×

Filter: 3×3

10	20	30
40	50	60
30	20	10

=

Output: 3×3

300	600	500

$$\begin{aligned} &0 \times 10 + 0 \times 20 + 0 \times 30 + \\ &2 \times 40 + 3 \times 50 + 0 \times 60 + \\ &5 \times 30 + 6 \times 20 + 0 \times 10 = \\ &500 \end{aligned}$$

Stride

- Stride: skip a location of image.

Input: 3×3 + 1×1 padding

0	0	0	0	0
0	1	2	3	0
0	4	5	6	0
0	7	8	9	0
0	0	0	0	0

×

Filter: 3×3

10	20	30
40	50	60
30	20	10

Stride: 2×2

=

Output: 2×2

300	

$$\begin{aligned} &0 \times 10 + 0 \times 20 + 0 \times 30 + \\ &0 \times 40 + 1 \times 50 + 2 \times 60 + \\ &0 \times 30 + 4 \times 20 + 5 \times 10 = \\ &300 \end{aligned}$$



Stride

- Stride: skip a location of image.

Input: 3×3 + 1×1 padding

0	0	0	0	0
0	1	2	3	0
0	4	5	6	0
0	7	8	9	0
0	0	0	0	0

×

Filter: 3×3

10	20	30
40	50	60
30	20	10

Stride: 2×2

=

Output: 2×2

300	500

$$0 \times 10 + 0 \times 20 + 0 \times 30 + 2 \times 40 + 3 \times 50 + 0 \times 60 + 5 \times 30 + 6 \times 20 + 0 \times 10 = 500$$



Output Size

- Input size: $n_h \times n_w$; filter size: $k_h \times k_w$; padding size: $p_h \times p_w$, stride size: $s_h \times s_w$.

- Output size:

$$\left\lfloor \frac{n_h + 2p_h - k_h}{s_h} + 1 \right\rfloor \times \left\lfloor \frac{n_w + 2p_w - k_w}{s_w} + 1 \right\rfloor$$

- For example:

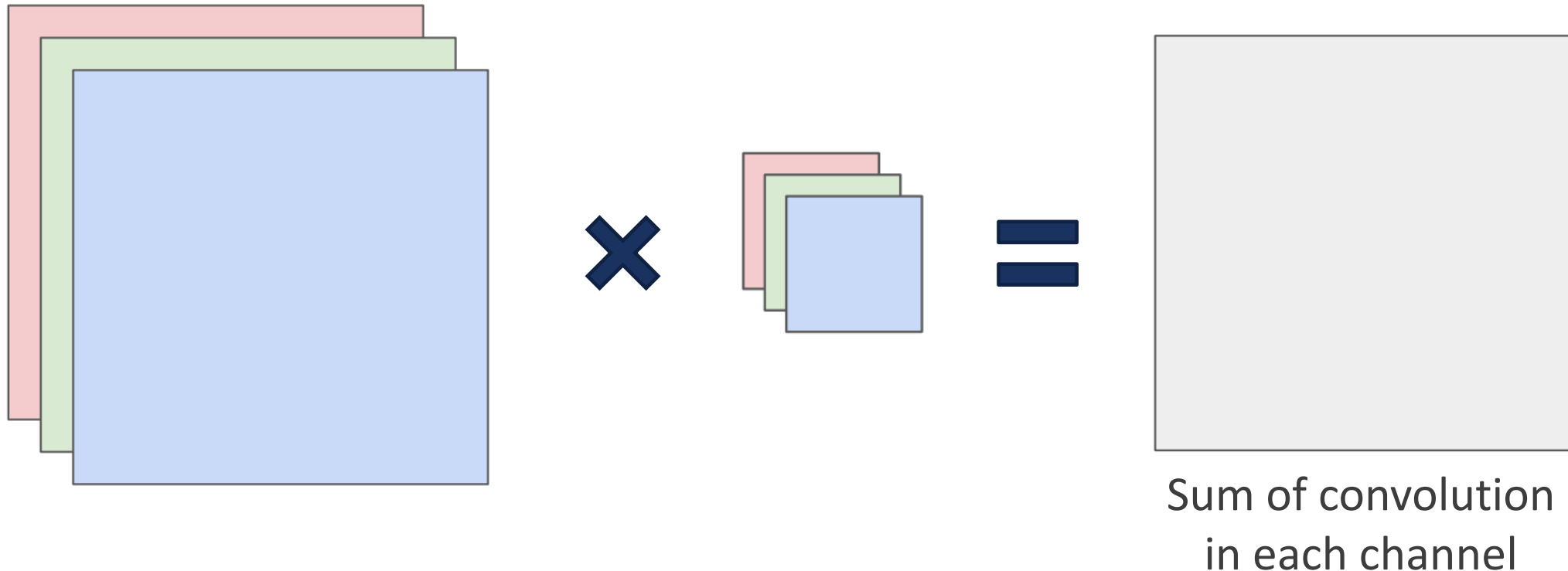
- Input size 3×3 , filter size 3×3 , padding size 1×1 , stride size 2×2 .

- Output size $\left\lfloor \frac{3+2-3}{2} + 1 \right\rfloor \times \left\lfloor \frac{3+2-3}{2} + 1 \right\rfloor = 2 \times 2$.



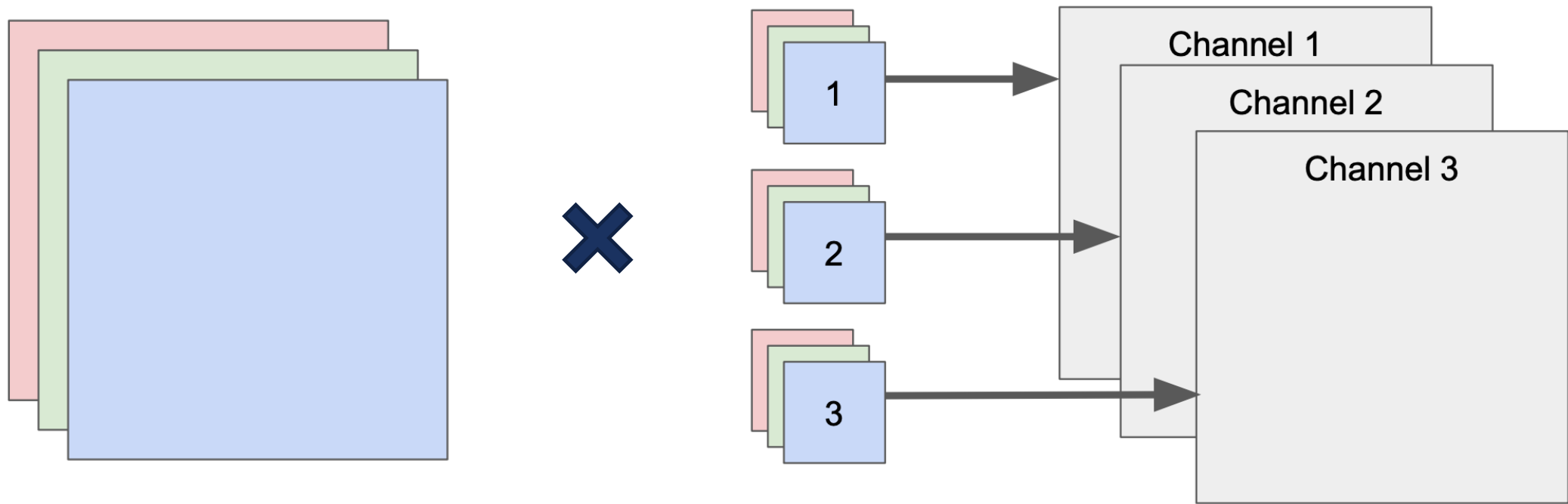
Channel and Depth

- The depth of the filter is same as the **channel** of the input image.
 - For an RGB image, we have three channels: red, green and blue.



Channel and Depth

- The depth of the feature map is a hyperparameter.
- It corresponds to the number of filters we would like to use, each learning to **look for something different** in the input.



Output Size with Depth

- Input size: $n_h \times n_w \times c_{in}$; filter size: $k_h \times k_w \times c_{in}$; filter number: c_{out} , padding size: $p_h \times p_w$, stride size: $s_h \times s_w$.

- Output size:

$$\left\lfloor \frac{n_h + 2p_h - k_h}{s_h} + 1 \right\rfloor \times \left\lfloor \frac{n_w + 2p_w - k_w}{s_w} + 1 \right\rfloor \times c_{out}$$

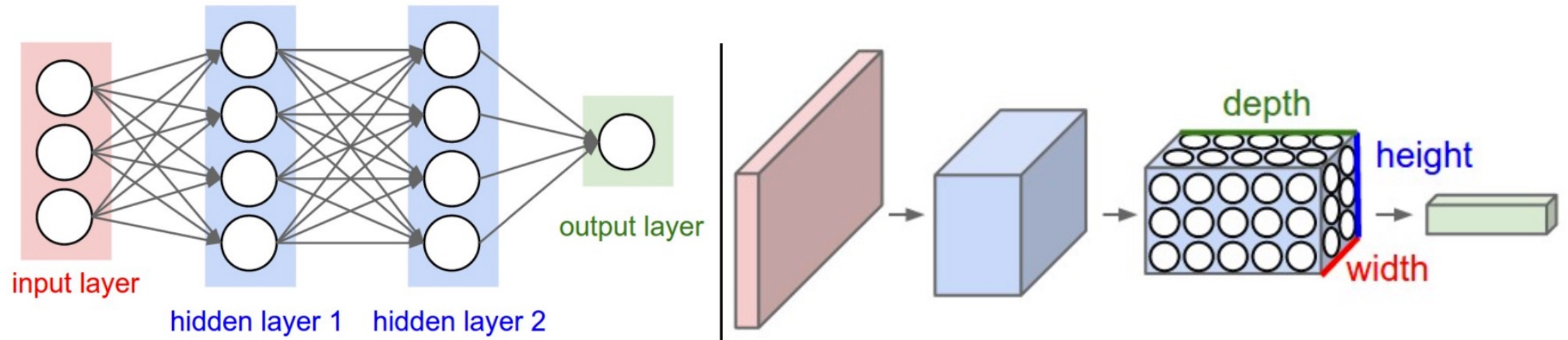
- For example:

- Input size $5 \times 5 \times 3$, filters size $3 \times 3 \times 3$, filter number 5, padding size 1×1 , stride size 2×2 .

- Output size $\left\lfloor \frac{5+2-3}{2} + 1 \right\rfloor \times \left\lfloor \frac{5+2-3}{2} + 1 \right\rfloor \times 5 = 3 \times 3 \times 5$.

Channel and Depth

- Every layer of a CNN transforms the 3D input volume to a 3D output volume of neuron activations.

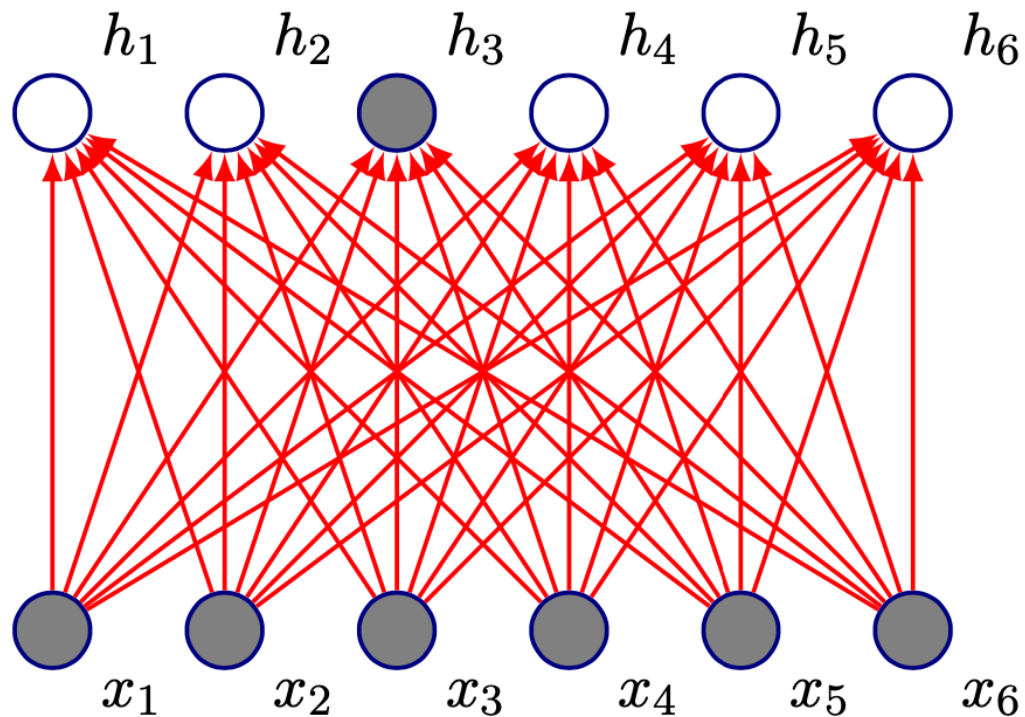


A regular 3-layer Neural Network.

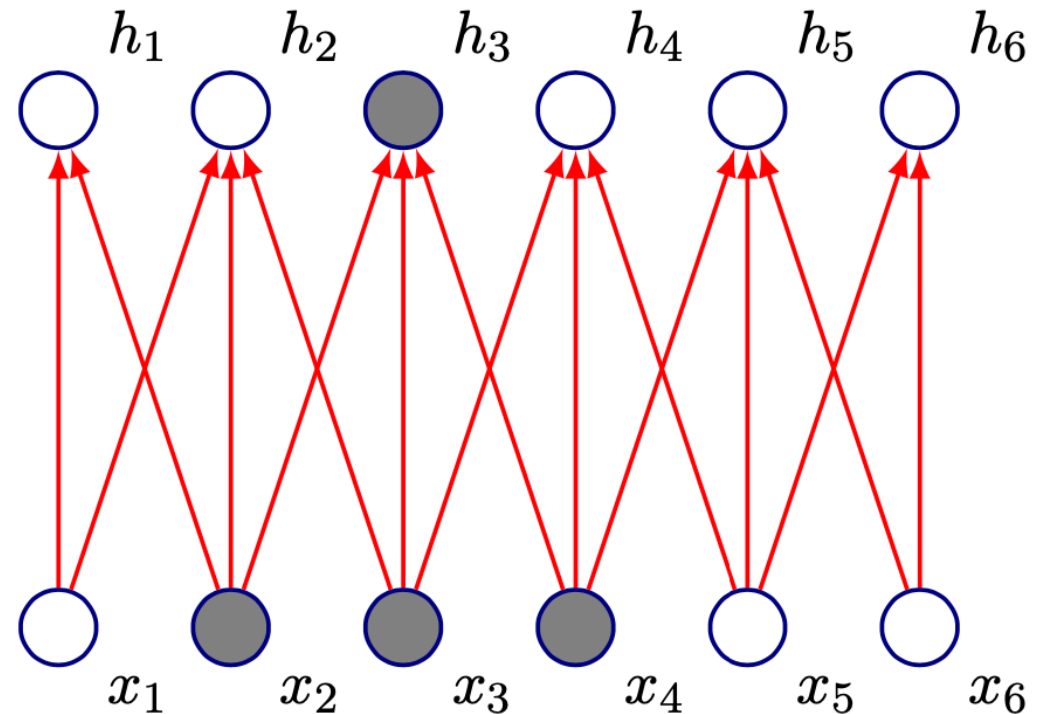
A CNN arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers



Sparse Connectivity



Fully connected network: h_3 is computed by full matrix multiplication with no sparse connectivity.



Kernel of size 3, moved with stride of 1. h_3 only depends on x_2, x_3, x_4 .



Sparse Connectivity

- Input: $55 \times 55 \times 3$, output: $55 \times 55 \times 96$.
- If we adopt a fully connected layer, the number of parameters for one single layer is:

$$(55 \times 55 \times 3 + 1) \times 55 \times 55 \times 96 = 2,635,670,400$$

- Now, if we use 96 11×11 filters with 5×5 padding and 1×1 stride.
- We can reduce the number of parameters to

$$(11 \times 11 \times 3 + 1) \times 55 \times 55 \times 96 = 105,705,600$$

- It is still unacceptable.

We use different filters for each pixel



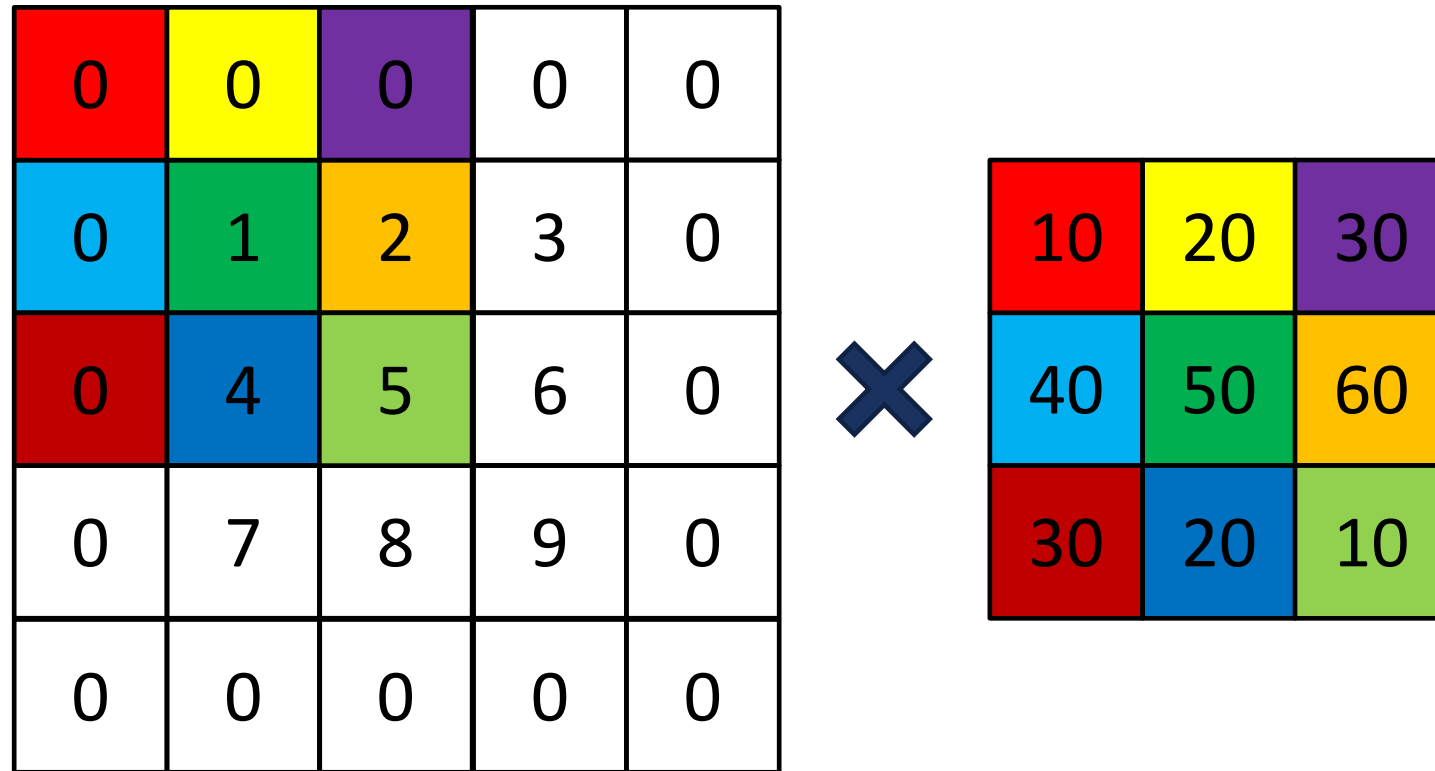
Parameter Sharing

- We can dramatically reduce the number of parameters by making one reasonable assumption:

If one filter is useful to compute at some spatial position (x_1, y_1) , it should also be useful to compute at a different position (x_2, y_2) .

- We are going to constrain the neurons in each channel to use the same weights and bias.

Parameter Sharing



A filter is fixed for all pixel positions in a channel

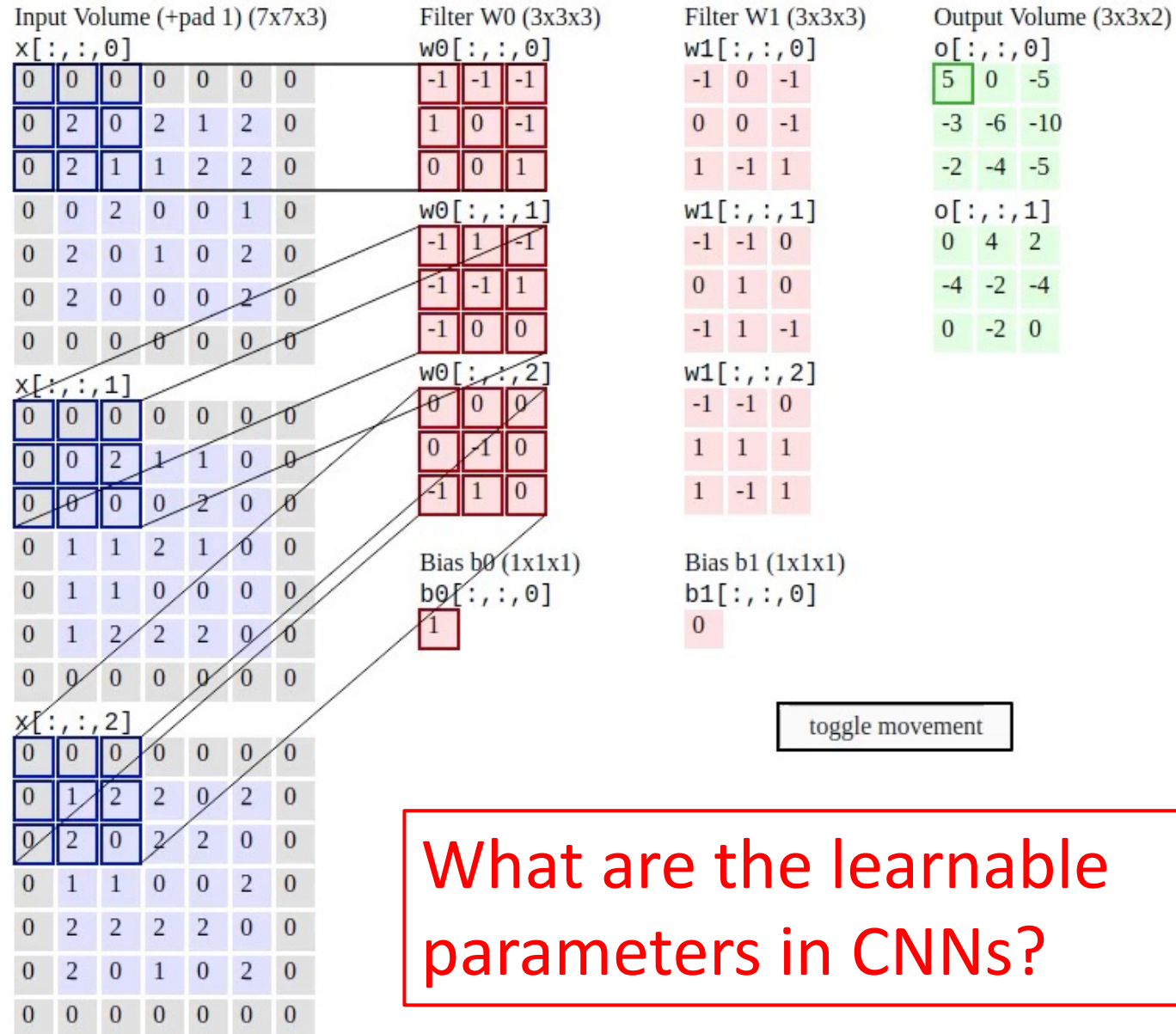
Parameter Sharing

- Rather than learning a separate set of filter parameters for every location during convolution, we learn only one set.
- This does not affect the runtime of forward propagation, but it does further reduce the storage requirements.
- In the previous example, the number of parameters are reduced to

$$96 \times (3 \times 11 \times 11 + 1) = 34,944$$

- 3000 times smaller than the non-sharing one.
- 75,400 times smaller than the fully connected one.





Parameter Sharing

- Each of the 96 learned filters is of size $11 \times 11 \times 3$.
- If detecting a horizontal edge is important at some location in the image, it should intuitively be **useful at some other location**.
- **No need to relearn** to detect a horizontal edge at every one of the 55×55 distinct locations in the Conv layer output volume.



Compare with the famous Sobel filter for edge detection:

X – Direction Kernel

-1	0	1
-2	0	2
-1	0	1

Y – Direction Kernel

-1	-2	-1
0	0	0
1	2	1

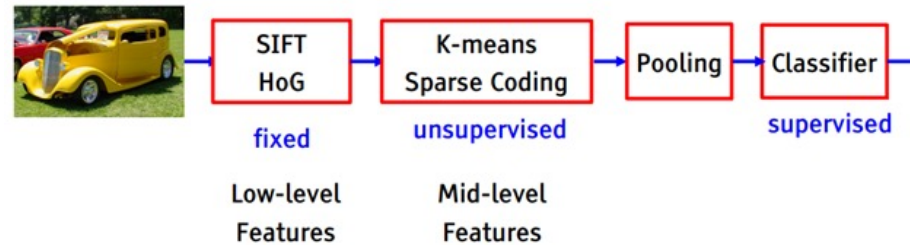


Parameter Sharing

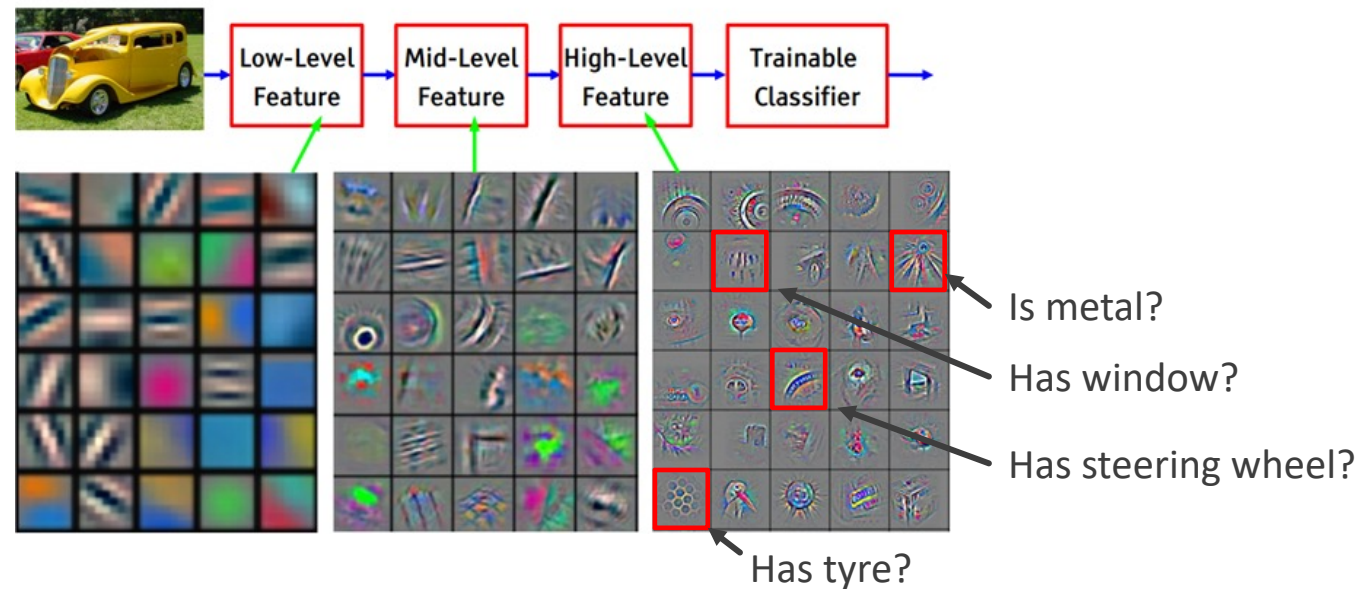
- Parameter sharing is not only for reducing the number of parameters.
- It can also be treated as a regularized method for preventing overfitting.
- It forces the filters to learn **some common patterns** over the whole image, rather than some patterns specific at some positions.

CNNs vs. Traditional Methods

Object recognition 2006-2012



State of the art object recognition using CNNs

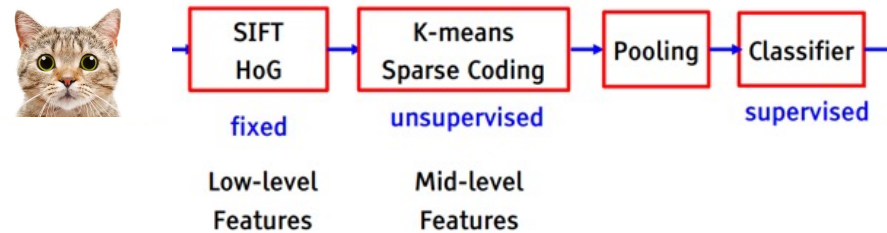


High-level feature contains **semantic information**, indicating if this image has certain components or not.



CNNs vs. Traditional Methods

Object recognition 2006-2012

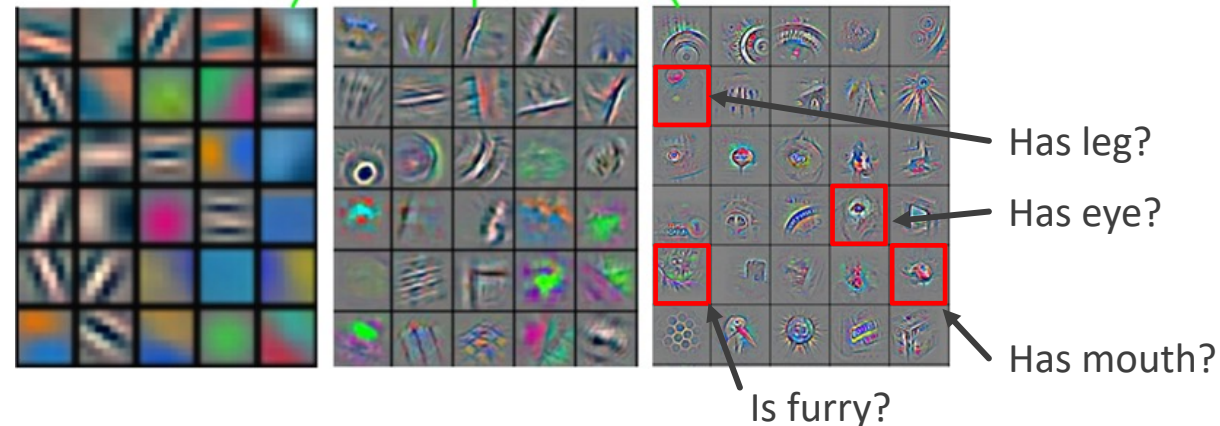


State of the art object recognition using CNNs



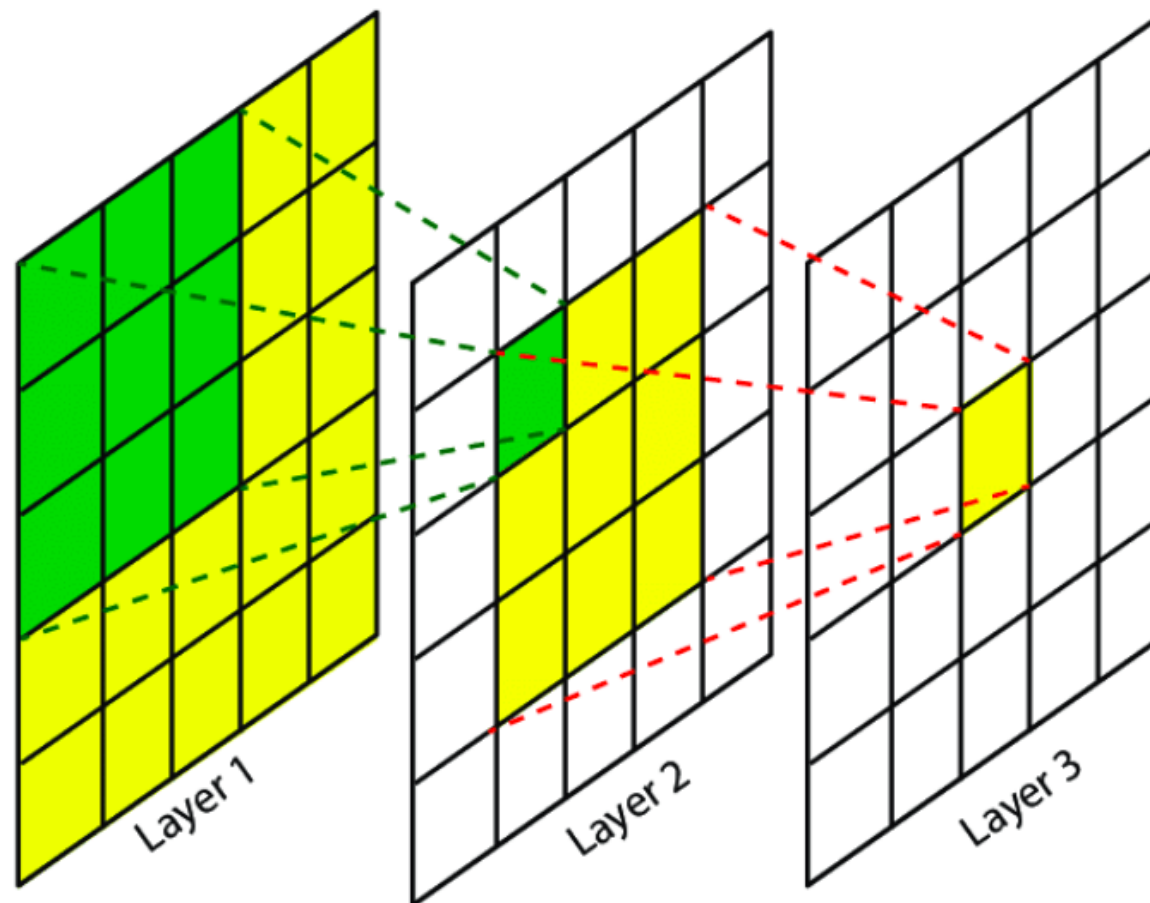
High-level features are also called representations.

High-level feature contains semantic information, indicating if this image has certain components or not.



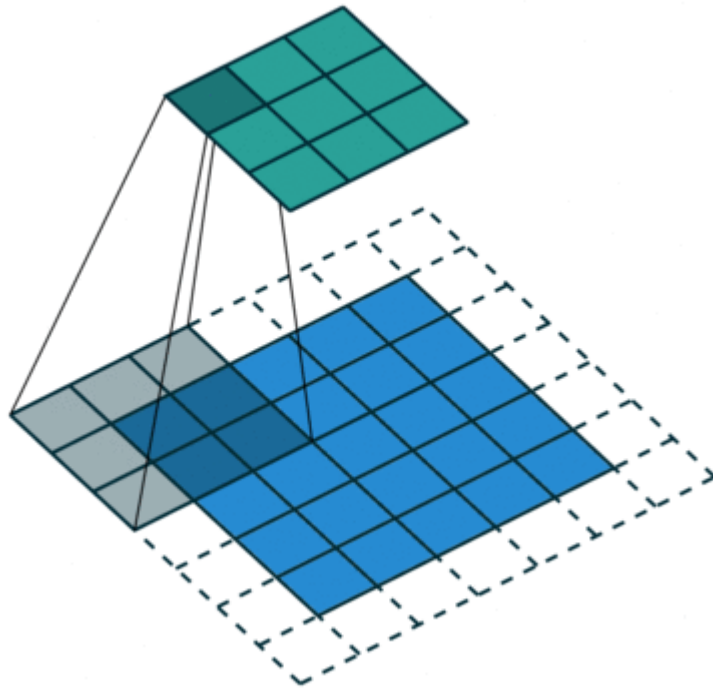
Receptive Field

- The **receptive field** is defined as the size of the region in the input that produces the feature.
- It is a measure of association of an output feature to the input region.
- Increasing model depth is a straightforward way to increase receptive field.
- Is there any other way to increase the receptive field?

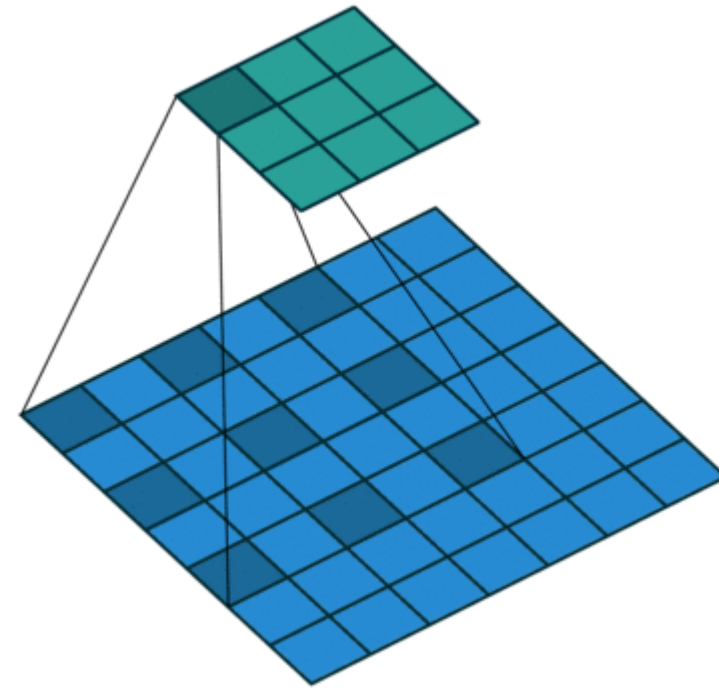


Dilated Convolution

- As an alternative, **dilated convolution** allows to merge spatial information across the inputs much more aggressively with fewer layers.



Standard Convolution ($d = 1$)



Dilated Convolution ($d = 2$)



Output Size with Dilated Convolution

- Input size: $n_h \times n_w \times c_{in}$; filter size: $k_h \times k_w \times c_{in}$; filter number: c_{out} ; padding size: $p_h \times p_w$, stride size: $s_h \times s_w$; dilation size: $d_h \times d_w$.

- Output size:

$$\left\lfloor \frac{n_h + 2p_h - (d_h \times k_h - 1) - 1}{s_h} + 1 \right\rfloor \times \left\lfloor \frac{n_w + 2p_w - (d_w \times k_w - 1) - 1}{s_w} + 1 \right\rfloor \\ \times c_{out}$$

- For example:

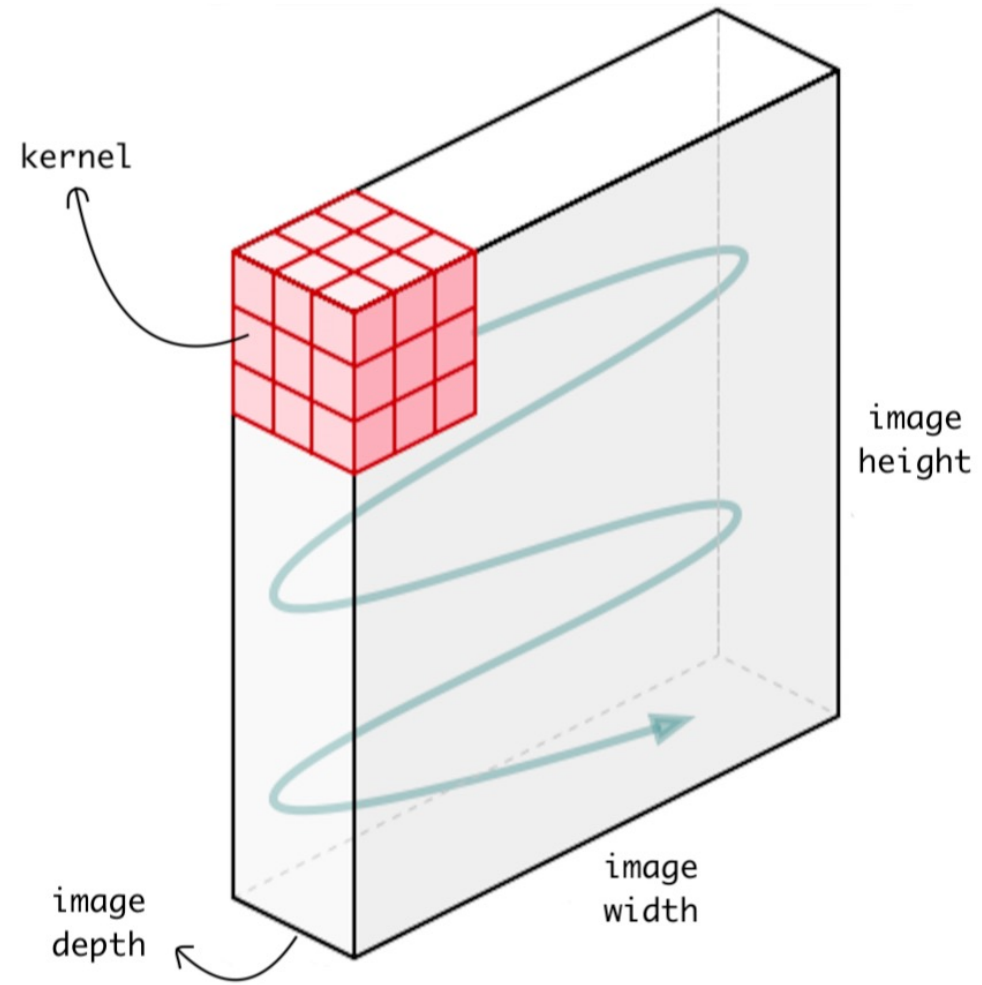
- Input size $11 \times 11 \times 3$, filter size $3 \times 3 \times 3$, padding size 1×1 , stride size 2×2 , dilation size 2×2 .

- Output size $\left\lfloor \frac{11+2-(2 \times 3-1)-1}{2} + 1 \right\rfloor \times \left\lfloor \frac{11+2-(2 \times 3-1)-1}{2} + 1 \right\rfloor \times 3 = 4 \times 4 \times 3$.



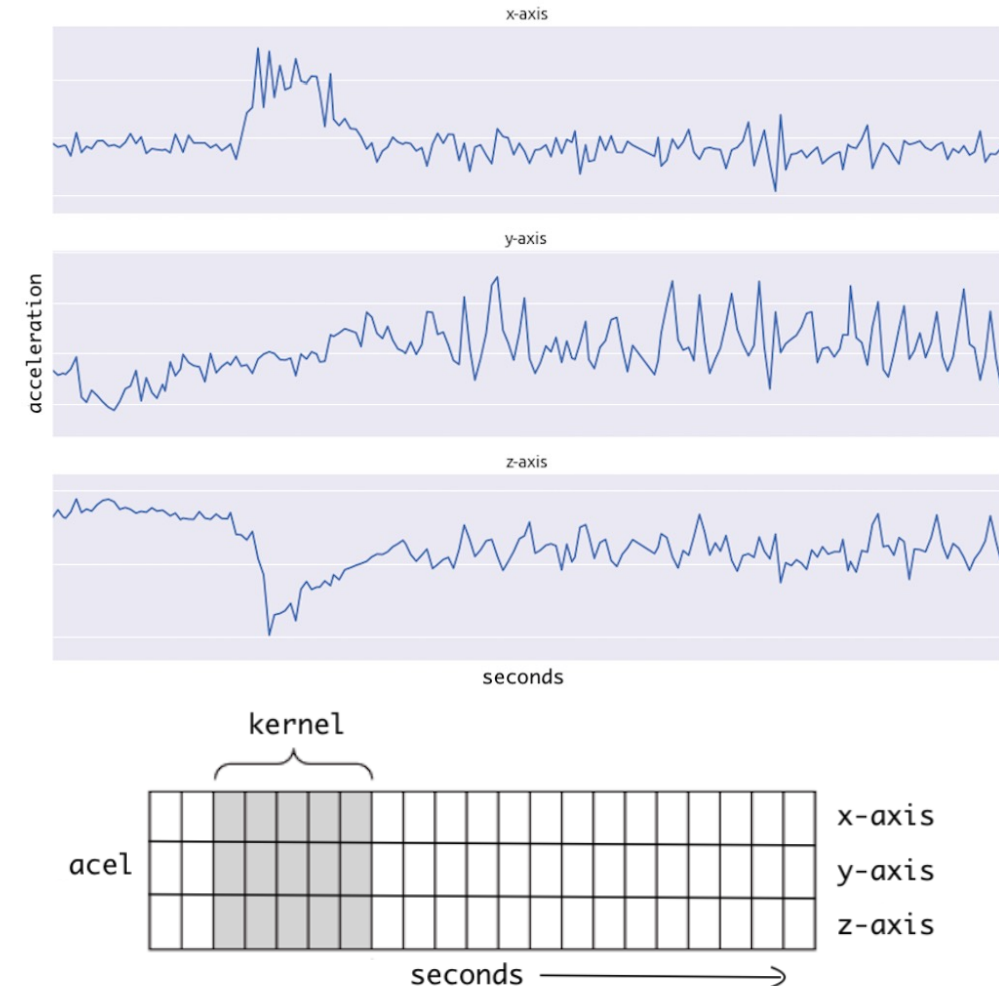
1D and 3D CNNs

- For image data, we are actually using 2D CNN.
- It means the input data and the filter are both 2-dimensional.



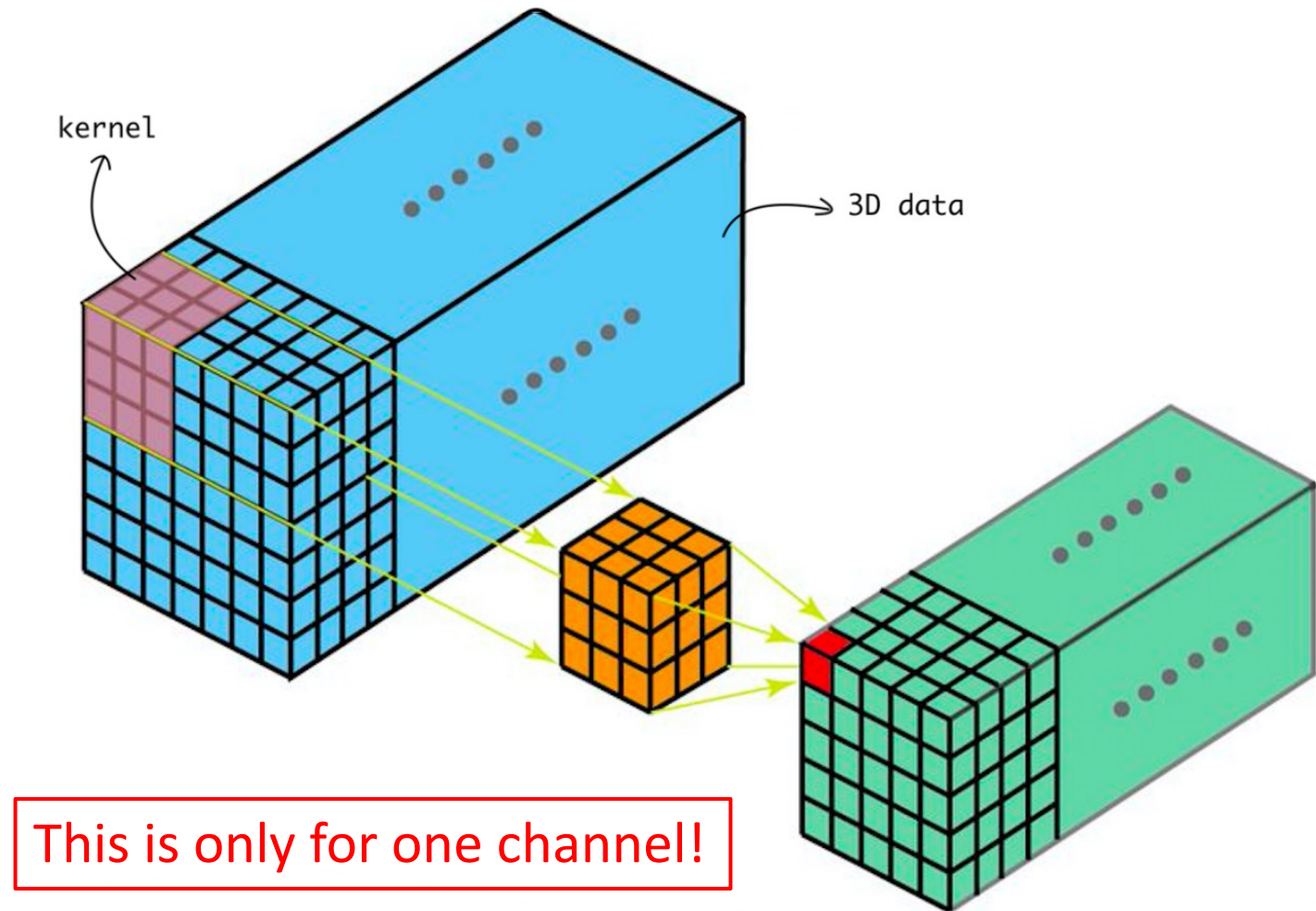
1D and 3D CNNs

- For 1D data like signal or time-series data, we can adopt 1D CNN with 1D filter.
- It is just a sparse and parameter-shared version of MLP.



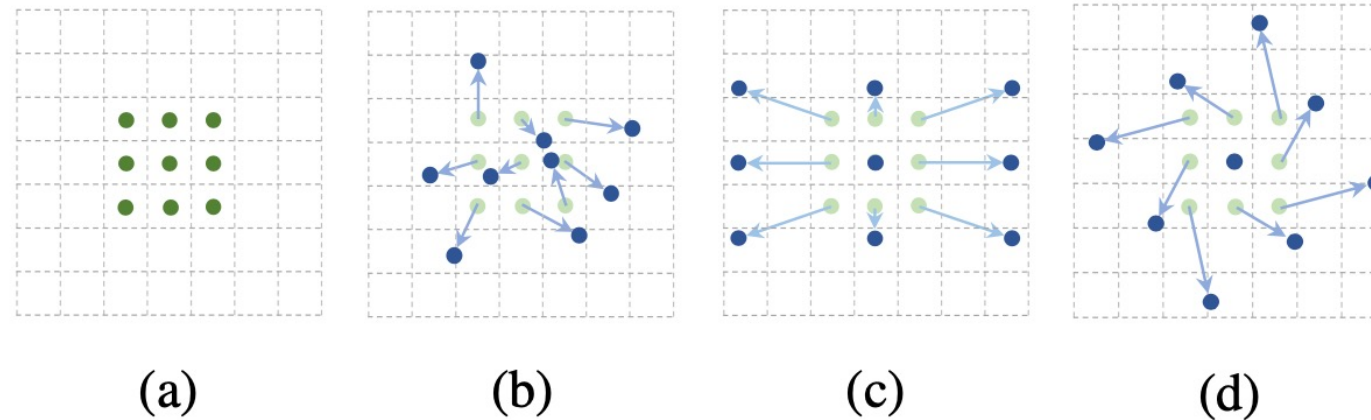
1D and 3D CNNs

- For 3D data like hyperspectral images, medical images, or videos, we can generalize 2D CNNs to 3D.



Deformable Convolutional Networks

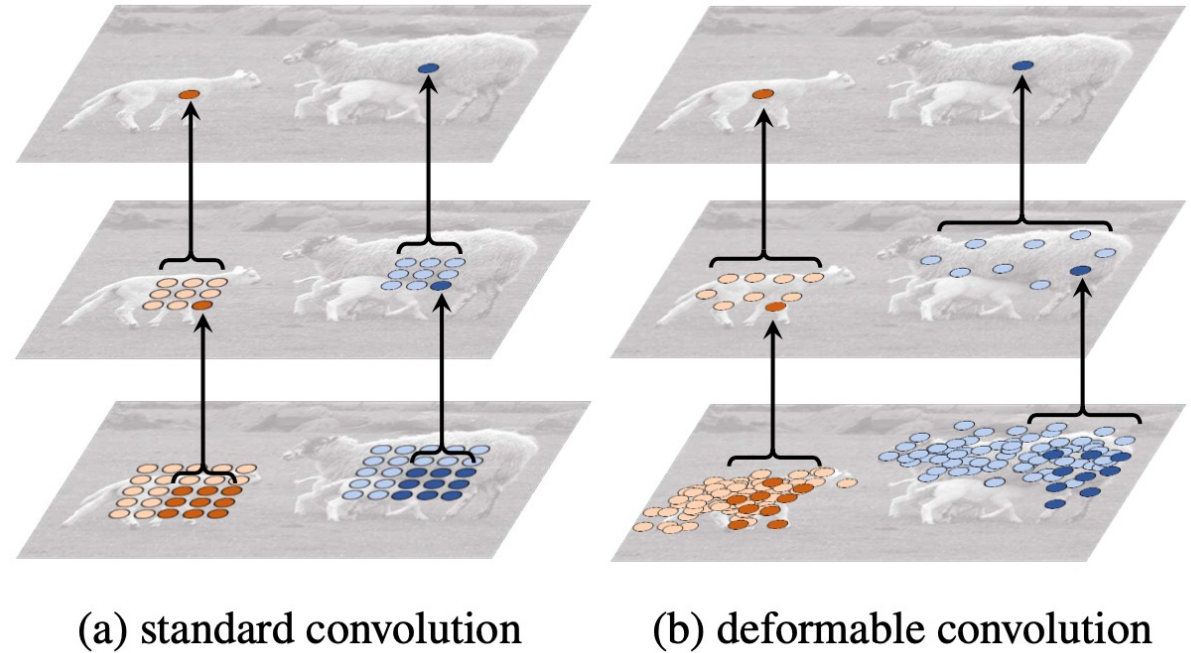
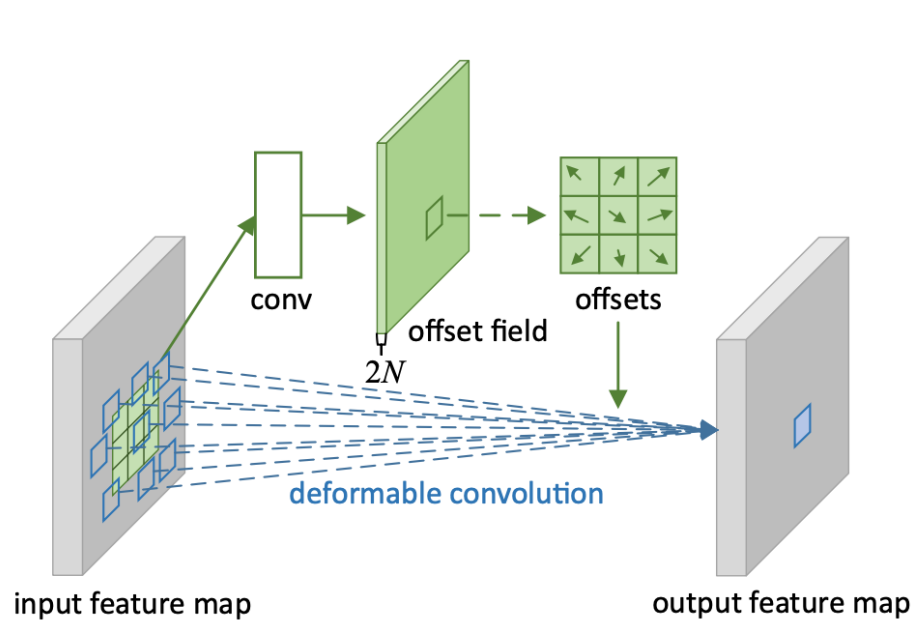
- One problem of CNNs: the shape of filters is fixed (e.g. square-like), but the shape of objects is variant.



- (a) regular sampling grid (green points) of standard convolution.
- (b) deformed sampling locations (dark blue points) with augmented offsets (light blue arrows) in deformable convolution.
- (c)(d) are special cases of (b).



Deformable Convolutional Networks



$$\mathbf{y}(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{x}(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n).$$

Filter weights and position offsets are learnable parameters



Convolutional Layer in PyTorch

Docs > torch.nn > Conv2d



CONV2D

CLASS `torch.nn.Conv2d(in_channels: int, out_channels: int, kernel_size: Union[T, Tuple[T, T]], stride: Union[T, Tuple[T, T]] = 1, padding: Union[T, Tuple[T, T]] = 0, dilation: Union[T, Tuple[T, T]] = 1, groups: int = 1, bias: bool = True, padding_mode: str = 'zeros')` [\[SOURCE\]](#)

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

$$\text{out}(N_i, C_{out_j}) = \text{bias}(C_{out_j}) + \sum_{k=0}^{C_{in}-1} \text{weight}(C_{out_j}, k) \star \text{input}(N_i, k)$$

where \star is the valid 2D **cross-correlation** operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

Parameters

- **in_channels** (*int*) – Number of channels in the input image
- **out_channels** (*int*) – Number of channels produced by the convolution
- **kernel_size** (*int or tuple*) – Size of the convolving kernel
- **stride** (*int or tuple, optional*) – Stride of the convolution. Default: 1
- **padding** (*int or tuple, optional*) – Zero-padding added to both sides of the input. Default: 0
- **padding_mode** (*string, optional*) – 'zeros', 'reflect', 'replicate' or 'circular'. Default: 'zeros'
- **dilation** (*int or tuple, optional*) – Spacing between kernel elements. Default: 1
- **groups** (*int, optional*) – Number of blocked connections from input channels to output channels. Default: 1
- **bias** (*bool, optional*) – If `True`, adds a learnable bias to the output. Default: `True`



```
import torch
rand_input = torch.randn(20, 16, 50, 100)
```

```
# With square kernels and equal stride
m = torch.nn.Conv2d(16, 33, 3, stride=2)
output = m(rand_input)
print(output.shape)
```

```
torch.Size([20, 33, 24, 49])
```

```
# non-square kernels and unequal stride and with padding
m = torch.nn.Conv2d(16, 33, (3, 5), stride=(2, 1), padding=(4, 2))
output = m(rand_input)
print(output.shape)
```

```
torch.Size([20, 33, 28, 100])
```

```
# non-square kernels and unequal stride and with padding and dilation
m = torch.nn.Conv2d(16, 33, (3, 5), stride=(2, 1), padding=(4, 2), dilation=(3, 1))
output = m(rand_input)
print(output.shape)
```

```
torch.Size([20, 33, 26, 100])
```

Input shape in PyTorch is:
 $[N, C_{in}, H, W]$



Convolutional Layer in TensorFlow

```
tf.keras.layers.Conv2D(  
    filters, kernel_size, strides=(1, 1), padding='valid', data_format=None,  
    dilation_rate=(1, 1), groups=1, activation=None, use_bias=True,  
    kernel_initializer='glorot_uniform', bias_initializer='zeros',  
    kernel_regularizer=None, bias_regularizer=None, activity_regularizer=None,  
    kernel_constraint=None, bias_constraint=None, **kwargs  
)
```

Arguments	
filters	Integer, the dimensionality of the output space (i.e. the number of output filters in the convolution).
kernel_size	An integer or tuple/list of 2 integers, specifying the height and width of the 2D convolution window. Can be a single integer to specify the same value for all spatial dimensions.
strides	An integer or tuple/list of 2 integers, specifying the strides of the convolution along the height and width. Can be a single integer to specify the same value for all spatial dimensions. Specifying any stride value != 1 is incompatible with specifying any dilation_rate value != 1.
padding	one of "valid" or "same" (case-insensitive).



```
import tensorflow as tf
rand_input = tf.random.normal([20, 50, 100, 16])
```

```
# With square kernels and equal stride
m = tf.keras.layers.Conv2D(33, 3, strides=2)
output = m(rand_input)
print(output.shape)
```

```
(20, 24, 49, 33)
```

```
# non-square kernels and unequal stride and with padding
p = tf.keras.layers.ZeroPadding2D((4, 2))
m = tf.keras.layers.Conv2D(33, (3, 5), strides=(2, 1))
output = m(p(rand_input))
print(output.shape)
```

```
(20, 28, 100, 33)
```

```
# non-square kernels and equal stride and with padding and dilation
p = tf.keras.layers.ZeroPadding2D((4, 2))
m = tf.keras.layers.Conv2D(33, (3, 5), strides=(1, 1), dilation_rate=(3, 1))
output = m(p(rand_input))
print(output.shape)
```

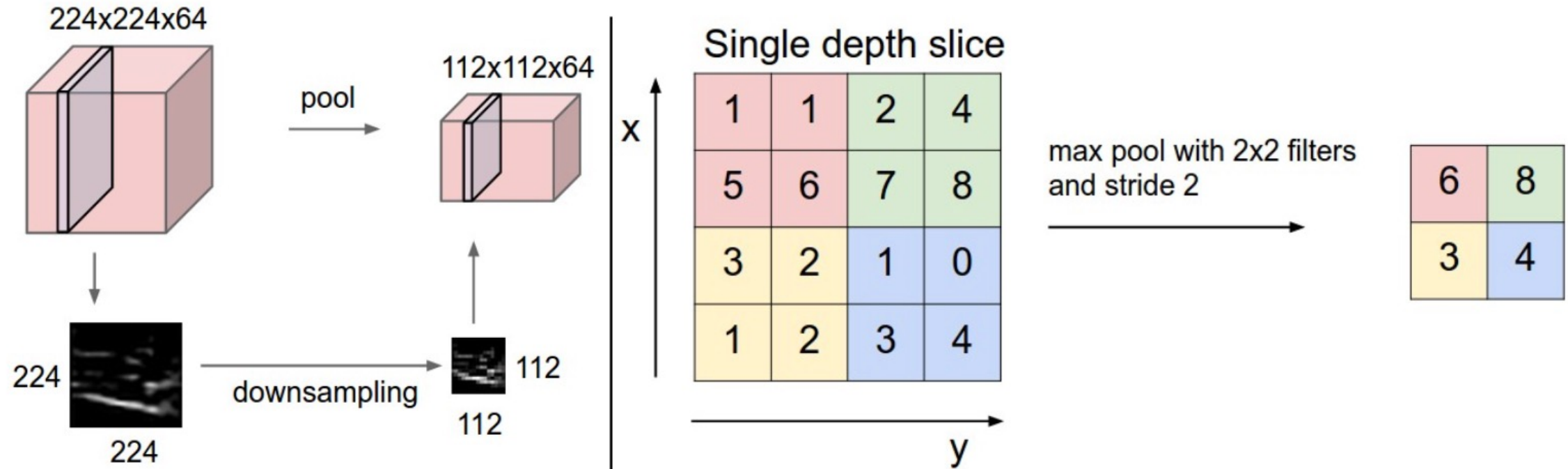
```
(20, 52, 100, 33)
```

Input shape in TF is:
 $[N, H, W, C_{in}]$

Special padding size can't
be assigned in conv2d

stride>1 is incompatible
with dilation_rate >1

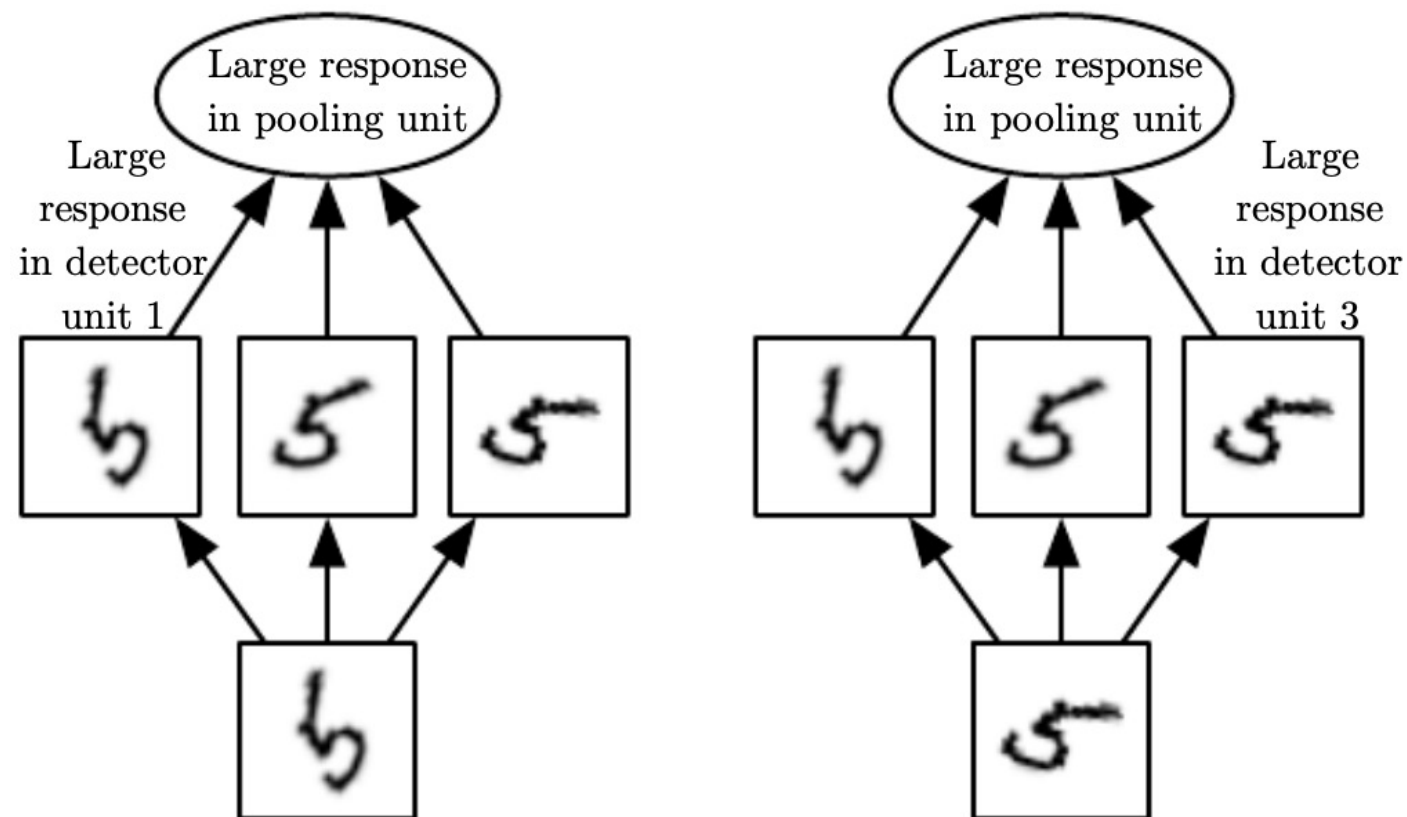
Pooling



- Pooling layer downsamples the volume spatially, independently in each channel of the input volume.

Pooling

- On one hand, pooling increases the receptive field.
- On the other hand, pooling introduces invariance.



The max pooling unit then has a large activation regardless of which detector unit was activated.



Pooling Layer in PyTorch

Docs > torch.nn > MaxPool2d



MAXPOOL2D

CLASS `torch.nn.MaxPool2d(kernel_size: Union[T, Tuple[T, ...]], stride: Optional[Union[T, Tuple[T, ...]]] = None, padding: Union[T, Tuple[T, ...]] = 0, dilation: Union[T, Tuple[T, ...]] = 1, return_indices: bool = False, ceil_mode: bool = False)` [\[SOURCE\]](#)

Applies a 2D max pooling over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C, H, W) , output (N, C, H_{out}, W_{out}) and `kernel_size` (kH, kW) can be precisely described as:

$$\text{out}(N_i, C_j, h, w) = \max_{m=0, \dots, kH-1} \max_{n=0, \dots, kW-1} \text{input}(N_i, C_j, \text{stride}[0] \times h + m, \text{stride}[1] \times w + n)$$

Parameters

- **kernel_size** – the size of the window to take a max over
- **stride** – the stride of the window. Default value is `kernel_size`
- **padding** – implicit zero padding to be added on both sides
- **dilation** – a parameter that controls the stride of elements in the window
- **return_indices** – if `True`, will return the max indices along with the outputs. Useful for `torch.nn.MaxUnpool2d` later
- **ceil_mode** – when `True`, will use `ceil` instead of `floor` to compute the output shape



Pooling Layer in PyTorch

```
# pool of square window of size=3, stride=2  
m = torch.nn.MaxPool2d(3, stride=2)  
output = m(rand_input)  
print(output.shape)
```

```
torch.Size([20, 16, 24, 49])
```

```
# pool of non-square window  
m = torch.nn.MaxPool2d((3, 2), stride=(2, 1))  
output = m(rand_input)  
print(output.shape)
```

```
torch.Size([20, 16, 24, 99])
```



Pooling Layer in TensorFlow

```
tf.keras.layers.MaxPool2D(  
    pool_size=(2, 2), strides=None, padding='valid', data_format=None, **kwargs  
)
```

Arguments

pool_size	integer or tuple of 2 integers, window size over which to take the maximum. (2, 2) will take the max value over a 2x2 pooling window. If only one integer is specified, the same window length will be used for both dimensions.
strides	Integer, tuple of 2 integers, or None. Strides values. Specifies how far the pooling window moves for each pooling step. If None, it will default to pool_size .
padding	One of "valid" or "same" (case-insensitive). "valid" adds no zero padding. "same" adds padding such that if the stride is 1, the output shape is the same as input shape.
data_format	A string, one of channels_last (default) or channels_first . The ordering of the dimensions in the inputs. channels_last corresponds to inputs with shape (batch, height, width, channels) while channels_first corresponds to inputs with shape (batch, channels, height, width). It defaults to the image_data_format value found in your Keras config file at <code>~/.keras/keras.json</code> . If you never set it, then it will be "channels_last".



Pooling Layer in TensorFlow

```
# pool of square window of size=3, strides=2  
m = tf.keras.layers.MaxPooling2D(3, strides=2)  
output = m(rand_input)  
print(output.shape)
```

```
(20, 24, 49, 16)
```

```
# pool of non-square window  
m = tf.keras.layers.MaxPooling2D((3, 2), strides=(2, 1))  
output = m(rand_input)  
print(output.shape)
```

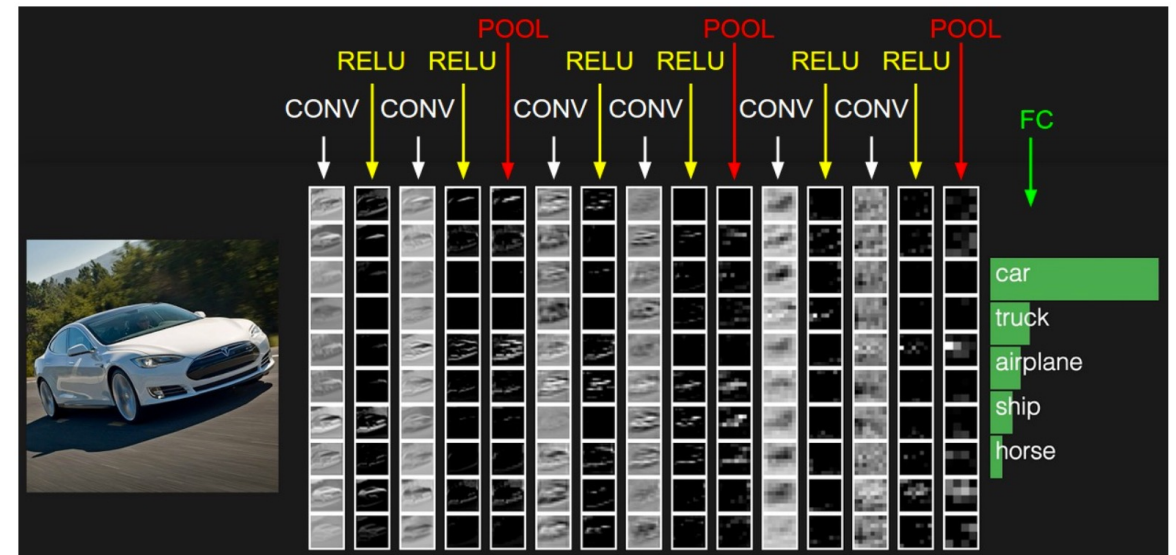
```
(20, 24, 99, 16)
```



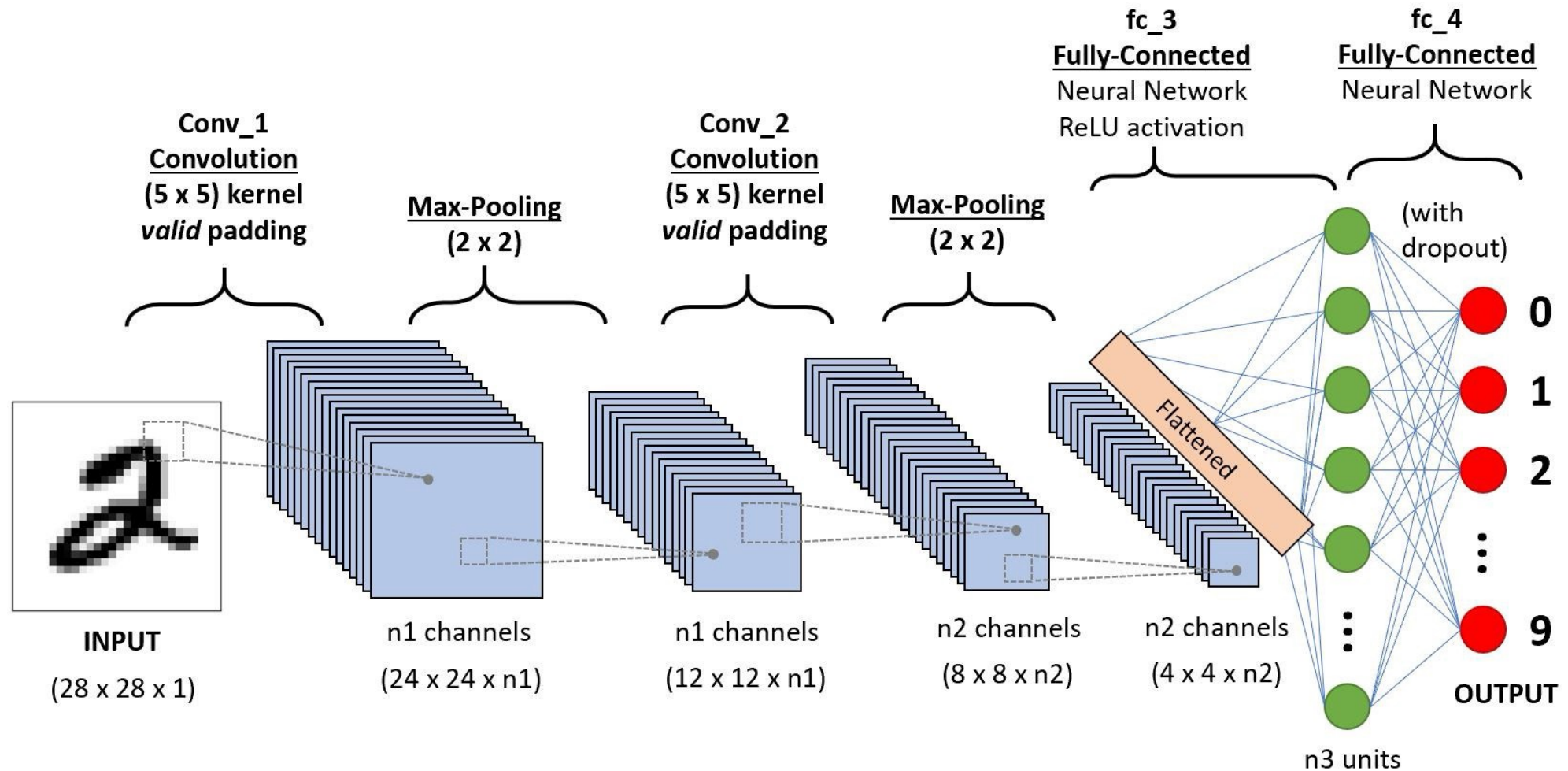
CNN Architecture

A typical CNN consists of four basic modules:

- **CONV layer** will compute the output of neurons that are connected to local regions in the input.
- **RELU layer** will apply an elementwise activation function.
- **FC (i.e. fully-connected) layer** will compute the class scores.
- **POOL layer** will perform a downsampling operation along the spatial dimensions (width, height).

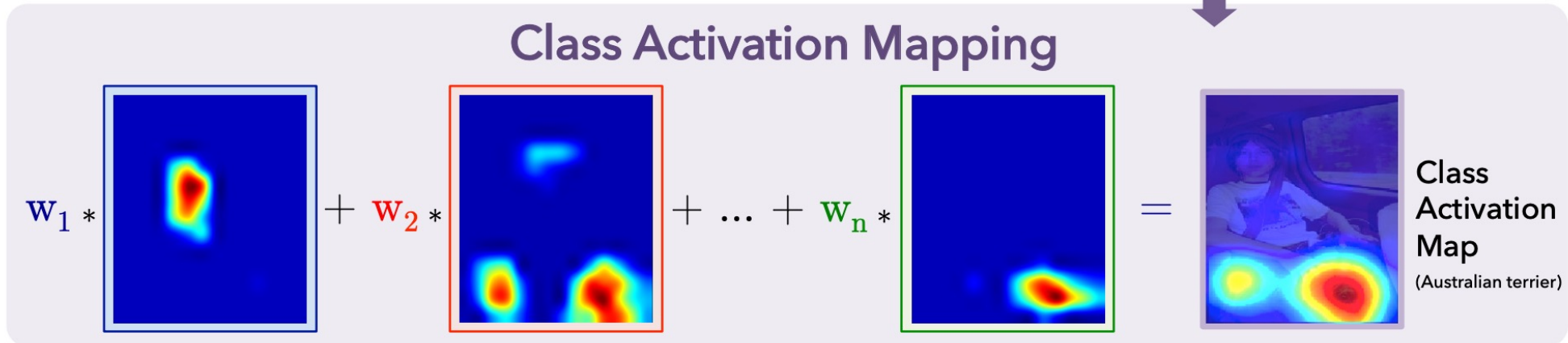
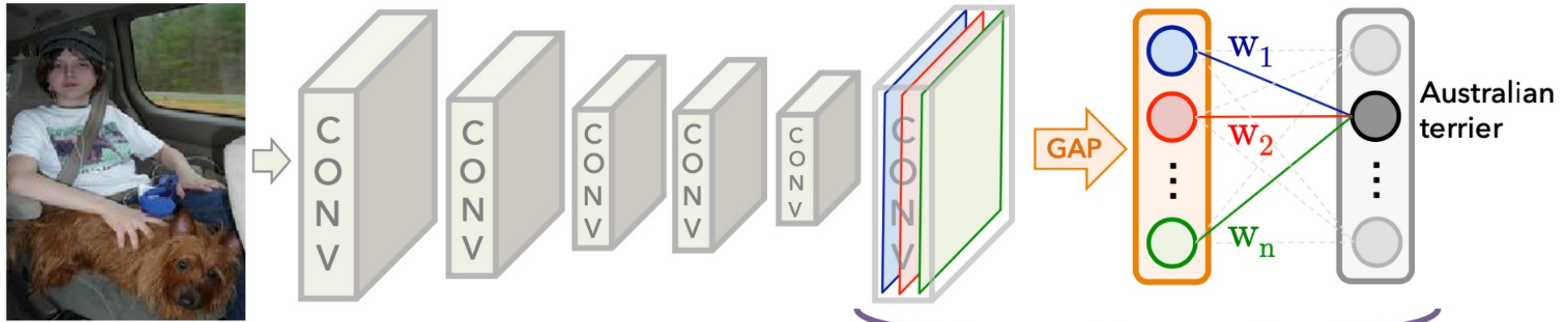


CNN Architecture



CNN Architecture

- We can also use **global average pooling (GAP)** to replace flattening.



CNN Architecture

- The most common form of a CNN architecture:
 - stacks a few CONV-RELU layers;
 - follows them with POOL layers;
 - repeats this pattern until the image has been merged spatially to a small size;
 - transits to fully-connected layers to produce output (e.g. class scores).
- The most common CNN architecture follows the pattern:
INPUT -> [[CONV -> RELU]*N -> POOL?]*M -> [FC -> RELU]*K -> FC
 - The * indicates repetition, and the indicates an optional pooling layer.
 - $N \geq 0$ (and usually $N \leq 3$), , $K \geq 0$ (and usually $K < 3$).

CNN Architecture

- INPUT -> FC.
 - A simple linear classifier. Here $N = M = K = 0$.
- INPUT -> CONV -> RELU -> FC.
 - Only CONV layer and RELU layer are used.
- INPUT -> [CONV -> RELU -> POOL]*2 -> FC -> RELU -> FC.
 - There is a single CONV layer between every POOL layer.
- INPUT -> [CONV -> RELU -> CONV -> RELU -> POOL]*3 -> [FC -> RELU]*2 -> FC.
 - Two CONV layers stacked before every POOL layer.
 - This is generally a good idea for larger and deeper networks, because multiple stacked CONV layers can **develop more complex features** of the input volume **before the destructive pooling** operation.

Filter Size

Is a stack of three 3×3 CONV layers equivalent to a single 7×7 CONV layer?

- No. There are several disadvantages for using filters with large size:
 - **Less powerful**: the neurons would be computing a linear function over the input, while the three stacks of CONV layers contain non-linearities that make their features more expressive.
 - **More parameters**: if both the input and output of a layer have depth C , 7×7 CONV layer would contain $C \times (7 \times 7 \times C) = 49C^2$, while the three 3×3 CONV layers only contains $3 \times C \times (3 \times 3 \times C) = 27C^2$.
- Intuitively, stacking CONV layers with tiny filters as opposed to having one CONV layer with big filters allows us to **express more powerful features of the input, and with fewer parameters**.

Layer Sizing Patterns

The common rules of thumb for sizing the architectures:

- The **INPUT layer** (that contains the image) should be divisible by 2 many times.
 - E.g. 32 (e.g. CIFAR-10), 64, 96 (e.g. STL-10), or 224 (e.g. common ImageNet ConvNets), 384, and 512.
- The **CONV layers** should be
 - using small filters (e.g. 3x3 or at most 5x5),
 - using a stride of 1x1,
 - padding the input volume with zeros in such way that the conv layer **does not alter the spatial dimensions of the input.**



Layer Sizing Patterns

- The **pool layers** are in charge of downsampling the spatial dimensions of the input.
 - The most common setting is to use max-pooling with 2×2 receptive fields, and with a stride of 2×2 .
 - Note that this discards exactly 75% of the activations in an input volume (due to downsampling by 2 in both width and height).
 - Another slightly less common setting is to use 3×3 receptive fields with a stride of 2.
 - It is very uncommon to see receptive field sizes for max pooling that are larger than 3 because the **pooling is then too lossy and aggressive**. This usually leads to worse performance.

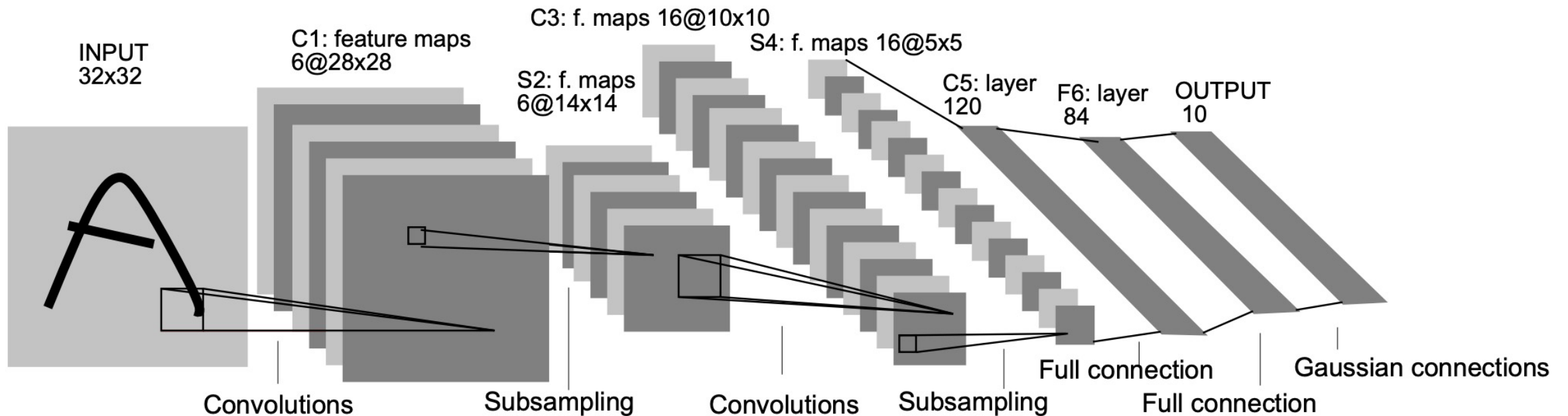
Classical CNN Architectures

Gradient-based learning applied to document recognition

[Y LeCun, L Bottou, Y Bengio...](#) - Proceedings of the ..., 1998 - [ieeexplore.ieee.org](#)

Multilayer neural networks trained with the back-propagation algorithm constitute the best example of a successful gradient based learning technique. Given an appropriate network architecture, gradient-based learning algorithms can be used to synthesize a complex decision surface that can classify high-dimensional patterns, such as handwritten characters, with minimal preprocessing. This paper reviews various methods applied to handwritten character recognition and compares them on a standard handwritten digit ...

☆ 57 [Cited by 30256](#) [Related articles](#) [All 38 versions](#)



Architecture of LeNet-5

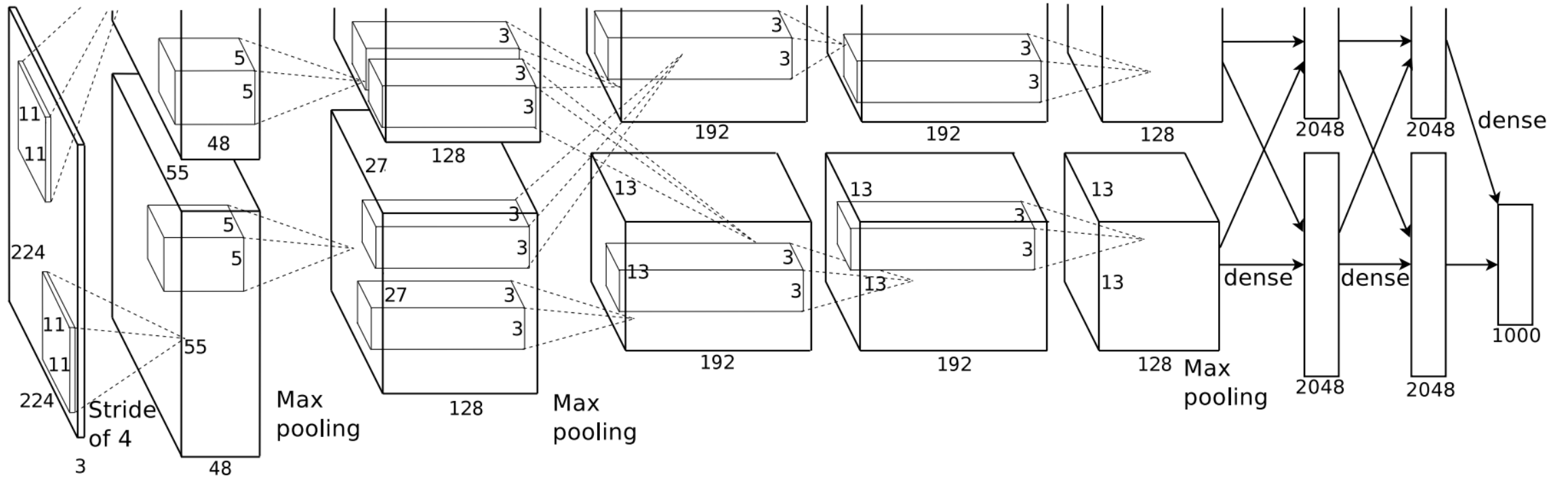


Classical CNN Architectures

Imagenet classification with deep convolutional neural networks

[A Krizhevsky, I Sutskever, GE Hinton - Advances in neural ..., 2012 - papers.nips.cc](#)
We trained a large, deep convolutional neural network to classify the 1.3 million high-resolution images in the LSVRC-2010 ImageNet training set into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 39.7% and 18.9% which is ...

☆ 77 **Cited by 70691** Related articles All 133 versions ⇄



Architecture of AlexNet

Conclusion

After this lecture, you should know:

- What is convolution and filter.
- What are the commonly used layers in CNN.
- How to calculate the output size of after a convolutional layer.
- Why do we need pooling.
- What are the typical CNN architectures.

Suggested Reading

- cs231n CNN tutorial
- Conv Nets: A Modular Perspective
- Understanding Convolutions



Thank you!

- Any question?
- Don't hesitate to send email to me for asking questions and discussion. 😊