

# Convolutional Neural Networks

Bùi Tiến Lên

2023



KHOA CÔNG NGHỆ THÔNG TIN  
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

# Contents

---



- 1. Architecture**
- 2. Convolutional layer**
- 3. Pooling layer**
- 4. Fully connected layer**
- 5. Visualizing and Understanding**
- 6. Network Design**
- 7. Miscellaneous convolutions**
- 8. Applications**



# Notation

Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design  
Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications  
Datasets  
Applications

symbol	meaning	operator	meaning
$a, b, c, N \dots$	scalar number	$w^T$	transpose
$\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$	column vector	$\mathbf{X} \mathbf{Y}$	matrix
$\mathbf{X}, \mathbf{Y} \dots$	matrix	$\mathbf{X}^{-1}$	multiplication
$\mathbb{R}$	set of real numbers	$\mathbf{X} \odot \mathbf{Y}$	inverse
$\mathbb{Z}$	set of integer numbers		an element-wise
$\mathbb{N}$	set of natural numbers		matrix-vector
$\mathbb{R}^D$	set of vectors		multiplication
$\mathcal{X}, \mathcal{Y}, \dots$	set		
$\mathcal{A}$	algorithm		



# Architecture



Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# Convolutional neural network

---

**Convolutional neural networks (CNNs) use three basic ideas:**

- *local receptive fields*
- *shared weights*
- *pooling*



Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# What are learnt in CNNs?

---

CNNs has two interesting properties

- The patterns they learn are translation invariant
- They can learn spatial hierarchies of patterns



Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

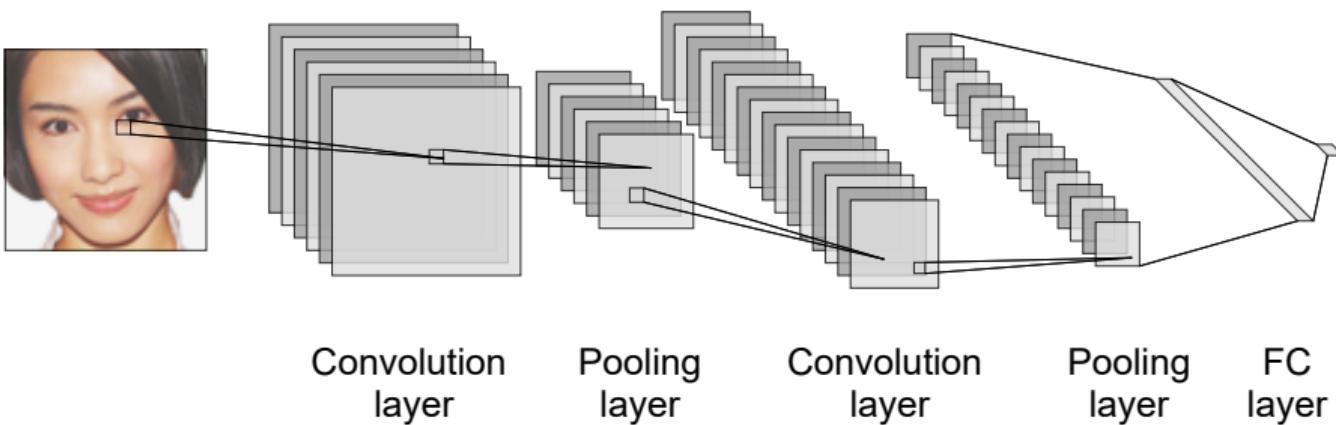
Datasets

Applications

# A basic architecture of CNN

- A common pattern

Convolution + ReLU → Pool → Fully connected → Softmax loss





## Convolutional layer

- Single Input, Single Output
- Multiple Input, Multiple Output
- Tensor Operation



Architecture

## Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

## Pooling layer

Fully connected layer

Visualizing and Understanding

## Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

## Miscellaneous convolutions

Transposed convolution

Dilated convolutions

## Applications

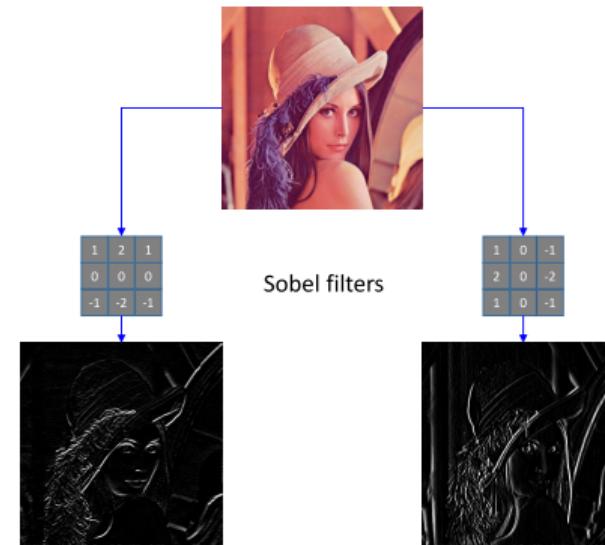
Datasets

Applications

# Convolutional layer

The fundamental difference between a densely connected layer and a convolution layer is this:

- Dense layers learn **global patterns** in their input feature space
- Convolutional layers learn **local patterns** in their input feature space





Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# Convolutional types

---

There are three common kinds of convolution

- 1D convolution
- **2D convolution**
- 3D convolution



# Image processing

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

Applications

In image processing, an input image is convolved with a kernel to create an output image

- **Input image** (*input feature map*)  $In$
- **Output image** (*output feature map*)  $Out$
- **Filter** (*kernel*) requires three parameters
  - The *spatial extent* (kernel size)  $F$
  - The *stride*  $S$  (distance between two consecutive positions of the kernel)
  - The amount of *zero padding*  $P$  (number of zeros concatenated at the beginning and at the end of an axis)



# Convolution vs. Cross-correlation

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

Miscellaneous convolutions

Transposed convolution  
Dilated convolutions

Applications

Datasets  
Applications

- Convolution operation (used in digital processing)

$$Out(x, y) = (K * In)(x, y) = \sum_m \sum_n In(x - m, y - n) K(m, n) \quad (1)$$

- Cross-correlation operation (used in CNN)

$$Out(x, y) = (K * In)(x, y) = \sum_m \sum_n In(x + m, y + n) K(m, n) \quad (2)$$



# Convolution vs. Cross-correlation (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

Miscellaneous convolutions

Transposed convolution  
Dilated convolutions

Applications  
Datasets  
Applications

The following figure illustrates 2D cross-correlation operation.

- Input size is  $3 \times 3$
- Output size is  $2 \times 2$
- Kernel has  $F = 2, S = 1, P = 0$

Input	Kernel	Output																	
<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td><td>2</td></tr><tr><td>3</td><td>4</td><td>5</td></tr><tr><td>6</td><td>7</td><td>8</td></tr></table>	0	1	2	3	4	5	6	7	8	$*$	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>0</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table> $=$ <table border="1" style="display: inline-table; vertical-align: middle;"><tr><td>19</td><td>25</td></tr><tr><td>37</td><td>43</td></tr></table>	0	1	2	3	19	25	37	43
0	1	2																	
3	4	5																	
6	7	8																	
0	1																		
2	3																		
19	25																		
37	43																		

- The shaded portions are the first output element as well as the input and kernel tensor elements used for the output computation

$$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19$$



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

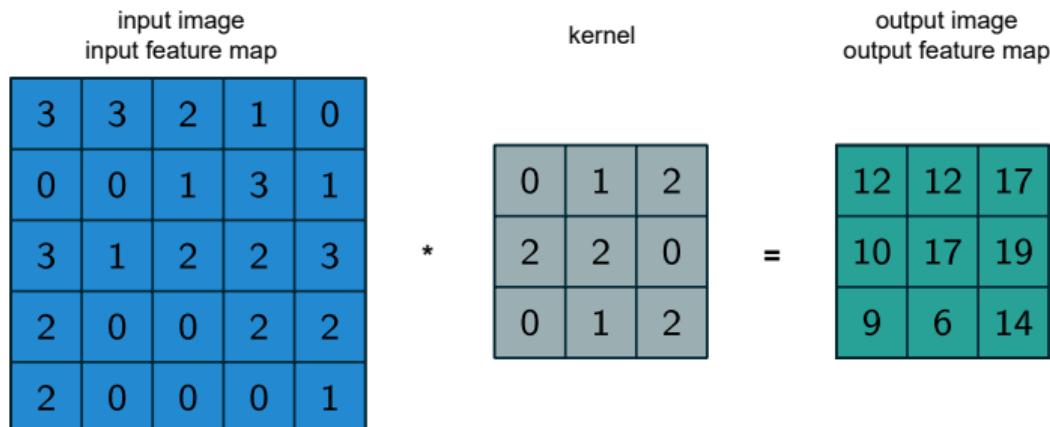
Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

# Example 1

- Input image size is  $5 \times 5$
- Output image size is  $3 \times 3$
- Kernel has  $F = 3, S = 1, P = 0$





Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# Example 1 (cont.)

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 <sub>0</sub>	2 <sub>1</sub>	1 <sub>2</sub>	0
0	0 <sub>2</sub>	1 <sub>2</sub>	3 <sub>0</sub>	1
3	1 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2 <sub>0</sub>	1 <sub>1</sub>	0 <sub>2</sub>
0	0	1 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>
3	1	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>2</sub>	3	1
3 <sub>2</sub>	1 <sub>2</sub>	2 <sub>0</sub>	2	3
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0 <sub>0</sub>	1 <sub>1</sub>	3 <sub>2</sub>	1
3	1 <sub>2</sub>	2 <sub>2</sub>	2 <sub>0</sub>	3
2	0 <sub>0</sub>	0 <sub>1</sub>	2 <sub>2</sub>	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1 <sub>0</sub>	3 <sub>1</sub>	1 <sub>2</sub>
3	1	2 <sub>2</sub>	2 <sub>2</sub>	3 <sub>0</sub>
2	0	0 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2 <sub>2</sub>	0 <sub>2</sub>	0 <sub>0</sub>	2	2
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 <sub>1</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3	0
2 <sub>0</sub>	0 <sub>2</sub>	0 <sub>0</sub>	2	0
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 <sub>1</sub>	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>	0 <sub>2</sub>
2	0	0	2	2 <sub>0</sub>
2	0	0	0	1 <sub>2</sub>

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

## Example 2

- Input image size is  $5 \times 5$
- Output image size is  $3 \times 3$
- Kernel has  $F = 3, S = 2, P = 1$

input image input feature map					kernel	output image output feature map		
3	3	2	1	0		*	0	1
0	0	1	3	1		2	2	0
3	1	2	2	3		0	1	2
2	0	0	2	2				
2	0	0	0	1				

=

output image output feature map		
6	17	3
8	17	13
6	4	4



Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# Example 2 (cont.)

0_0	0_1	0_2	0	0	0	0	0
0_2	3_2	3_0	2	1	0	0	0
0_0	0_1	0_2	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	1

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0_0	0_1	0_2	0	0	0
0	3	3_2	2_2	1_0	0	0	0
0	0	0_0	1_1	3_2	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0_0	0_1	0_2
0	3	3	2	1_2	0_2	0_0	0_0
0	0	0	1	3_0	1_1	0_2	0_0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	1	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	1	0	2	2	0	0
0	2	2	0	0	0	1	0
0	0	1	0	2	2	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	2	0	0	1	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	2	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0



# Multiple Input Channels

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

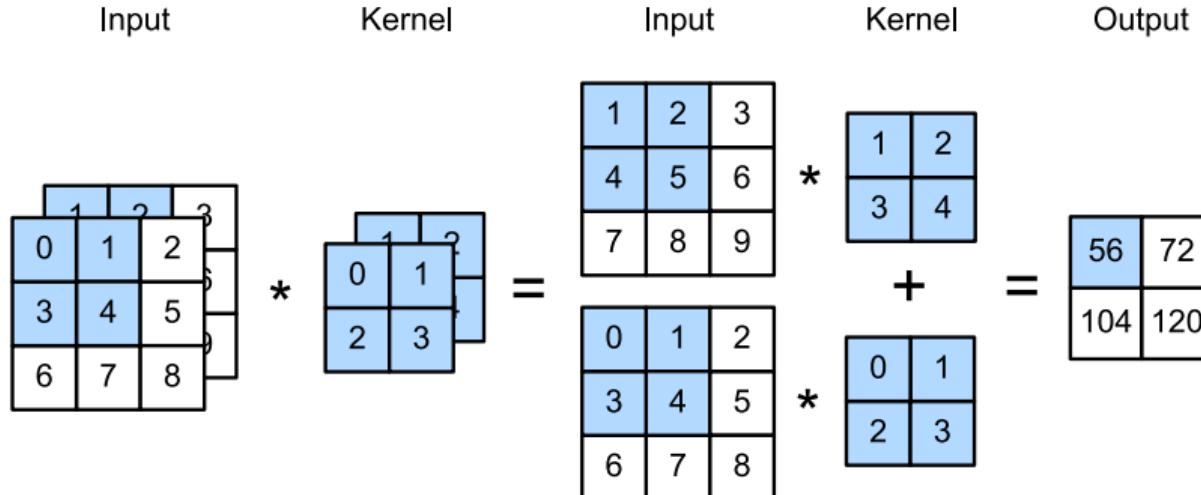
Dilated convolutions

Applications

Datasets

Applications

- Convolution computation with 2 input channels.



- The shaded portions are the first output element as well as the input and kernel tensor elements used for the output computation

$$(1 \times 1 + 2 \times 2 + 4 \times 3 + 5 \times 4) + (0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3) = 56$$



# Multiple Input, Output Channels

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

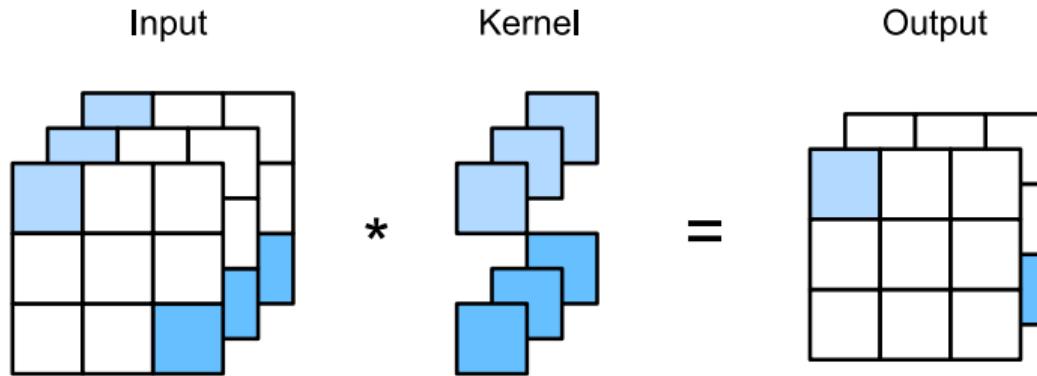
Miscellaneous convolutions

Transposed convolution  
Dilated convolutions

Applications

Datasets  
Applications

- Convolution computation for 3 input channels and 2 output channels.





Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

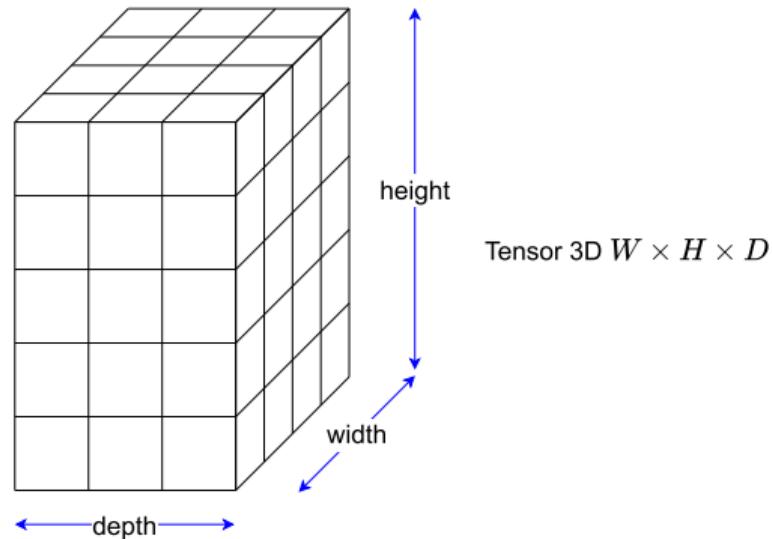
Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

# 3D tensor

- Convolutions operate over 3D tensors, called *feature maps*, with two spatial axes (*width W* and *height H*) as well as a *depth axis D* (also called the *channels axis*)





# Convolution operator

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

Applications

- Accepts a volume of size  $W_{in} \times H_{in} \times D_{in}$  and a group of filters with four hyperparameters:
  - Number of filters  $K$
  - Their spatial extent  $F$
  - The stride  $S$
  - The amount of zero padding  $P$



# Convolution operator (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

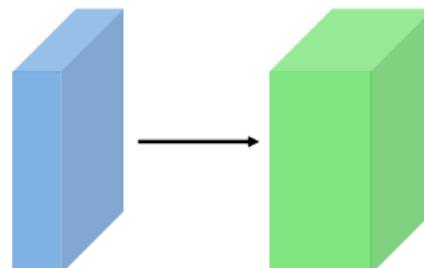
Applications

- Produces a volume of size  $W_{out} \times H_{out} \times D_{out}$  where:

$$W_{out} = \frac{W_{in} - F + 2P}{S} + 1$$

$$H_{out} = \frac{H_{in} - F + 2P}{S} + 1$$

$$D_{out} = K \quad (3)$$



$W_{in} \times H_{in} \times D_{in}$

$W_{out} \times H_{out} \times D_{out}$



# Convolution operator (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

Applications

In summary,

- With parameter sharing, it introduces  $F \times F \times D_{in}$  weights per filter, for a total of  $F \times F \times D_{in} \times K$  **weights** and  $K$  **biases**.
- The filters depth is always equal to the input volume depth.
- The output volume depth is equal to the number of filters
- In the output volume, the  $d$ -th depth slice (of size  $W_{out} \times H_{out}$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.
- Normally, in the same conv layer, all filters have the same dimensions so that special optimized routines can be invoked.



# Explanation

---

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

Applications

- We can imagine a filter as a neuron.
- In a convolution layer we have filters.
  - Each filter is a weight matrix which is convoluted over the input volume.
  - Each time the filter is applied it outputs a single value.
  - The result of convolving the filter over all the input volume (*input feature maps*) is another volume called **activation map** (*output feature map*).
  - A filter's depth is always equal to the input volume depth.
- Instead of having only one filter, we have a stack of filter which produce a stack of activation maps (one for each filter)



# Explanation (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

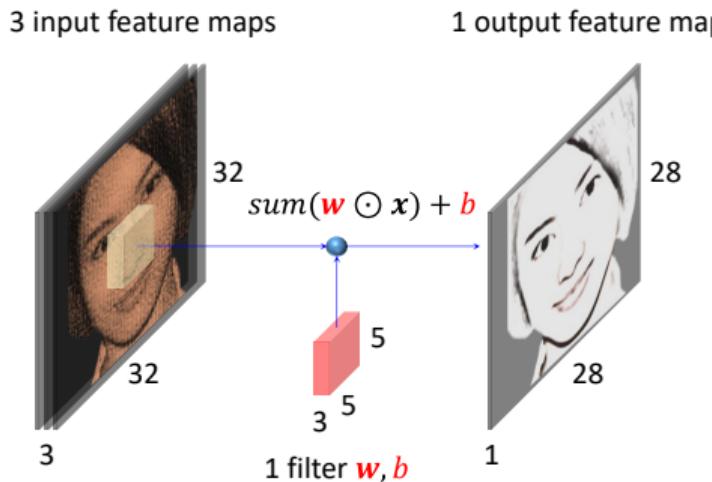
Miscellaneous convolutions

Transposed convolution  
Dilated convolutions

Applications

Datasets  
Applications

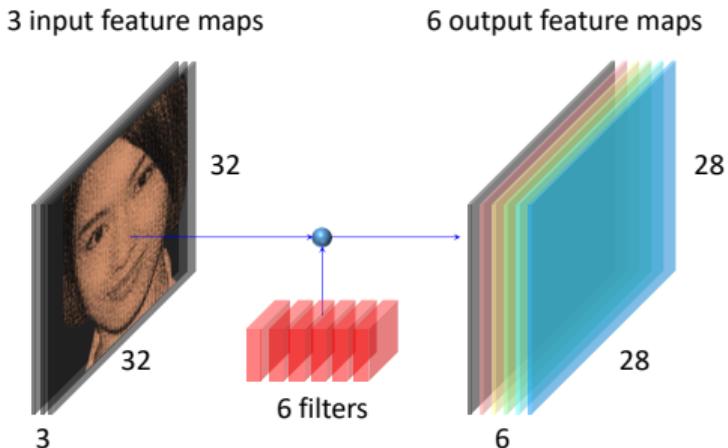
- The input volume  $32 \times 32 \times 3$  is convolved with the filter  $5 \times 5 \times 3$  ( $S = 1$  and  $P = 0$ ) and the result is the activation map  $28 \times 28 \times 1$





# Explanation (cont.)

- The input volume  $32 \times 32 \times 3$  is convolved with the stack of 6 filter  $5 \times 5 \times 3$  and the result is the stack of 6 activation map  $28 \times 28 \times 1$  ( $28 \times 28 \times 6$ )



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications



# Pooling layer



Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design  
Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications  
Datasets  
Applications

# Pooling layer

- The reason to use pooling layer (downsampling) is to reduce the number of feature-map coefficients to process
- Pooling layers also provide invariance to small translations of the input
- The most common kind of pooling is *max pooling*
- Another kind of pooling is average pooling; however, it's more informative to look at the *maximal presence* of different features than at their *average presence*



# Pooling operator

- Accepts a volume of size  $W_{in} \times H_{in} \times D_{in}$  and one filter with two hyperparameters:
  - Their spatial extent  $F$
  - The stride  $S$
- Produces a volume of size  $W_{out} \times H_{out} \times D_{out}$  where:

$$W_{out} = \frac{W_{in} - F}{S} + 1$$

$$H_{out} = \frac{H_{in} - F}{S} + 1$$

$$D_{out} = D_{in}$$

(4)

- The most common form is the filter of size  $2 \times 2$  applied with a stride of 2.



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

## Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

## Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

## Miscellaneous convolutions

Transposed  
convolution

Dilated convolutions

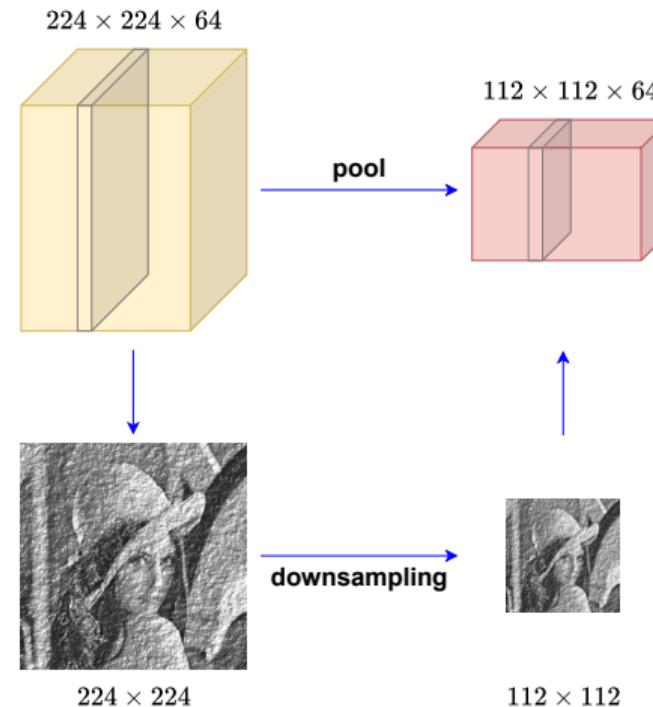
## Applications

Datasets

Applications

# Example

- The input volume of size  $[224 \times 224 \times 64]$  is pooled with filter size 2, stride 2 into output volume of size  $[112 \times 112 \times 64]$





Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

## Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design  
Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

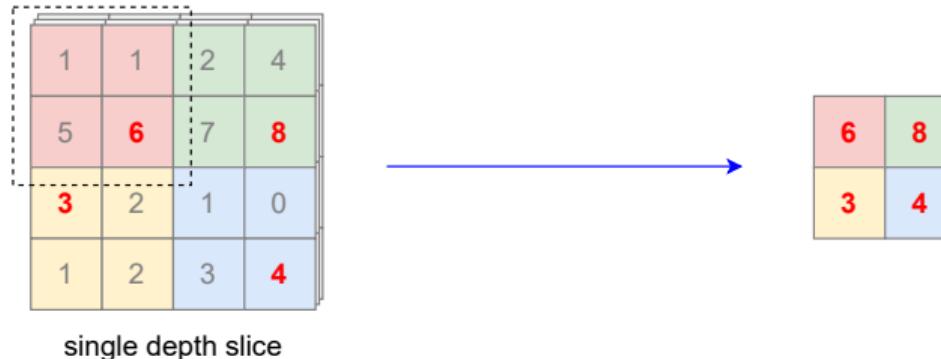
Miscellaneous  
convolutions  
Transposed  
convolution  
Dilated convolutions

Applications  
Datasets  
Applications

# Example (cont.)

- Max pooling with filter size 2, stride 2

maxpool with  $2 \times 2$  filters  
and stride 2





# Fully connected layer



# Fully connected layer

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

- Fully connected layers (FC layers), also called *affine layers*, produce the *high-level reasoning* in the DNN.
- Neurons in a fully connected layer have full connections to all activations in the previous layer, as in regular Neural Networks.
- Their activations can hence be computed with a matrix multiplication followed by a bias offset.



Architecture

Convolutional  
layerSingle Input, Single  
OutputMultiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layerVisualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutionsTransposed  
convolution

Dilated convolutions

Applications

Datasets

Applications

# VGGNet

Layer name	Variables		Parameters	
	Layout	Size	Weights	Biases
INPUT	[224x224x3]	224*224*3=150K		
CONV3-64	[224x224x64]	224*224*64=3.2M	(3*3*3)*64=1,728	64
CONV3-64	[224x224x64]	224*224*64=3.2M	(3*3*64)*64=36,864	64
POOL2	[112x112x64]	112*112*64=800K		
CONV3-128	[112x112x128]	112*112*128=1.6M	(3*3*64)*128=73,728	128
CONV3-128	[112x112x128]	112*112*128=1.6M	(3*3*128)*128=147,456	128
POOL2	[56x56x128]	56*56*128=400K		
CONV3-256	[56x56x256]	56*56*256=800K	(3*3*128)*256=294,912	256
CONV3-256	[56x56x256]	56*56*256=800K	(3*3*256)*256=589,824	256
CONV3-256	[56x56x256]	56*56*256=800K	(3*3*256)*256=589,824	256
POOL2	[28x28x256]	28*28*256=200K		
CONV3-512	[28x28x512]	28*28*512=400K	(3*3*256)*512=1,179,648	512
CONV3-512	[28x28x512]	28*28*512=400K	(3*3*512)*512=2,359,296	512
CONV3-512	[28x28x512]	28*28*512=400K	(3*3*512)*512=2,359,296	512
POOL2	[14x14x512]	14*14*512=100K		
CONV3-512	[14x14x512]	14*14*512=100K	(3*3*512)*512=2,359,296	512
CONV3-512	[14x14x512]	14*14*512=100K	(3*3*512)*512=2,359,296	512
CONV3-512	[14x14x512]	14*14*512=100K	(3*3*512)*512=2,359,296	512
POOL2	[7x7x512]	7*7*512=25K		
FC	[1x1x4096]	4096	7*7*512*4096=102,760,448	4096
FC	[1x1x4096]	4096	4096*4096=16,777,216	4096
FC	[1x1x1000]	1000	4096*1000=4,096,000	1000

- Total variables (memory):  $24M \times 4 \text{ bytes} \simeq 93\text{MB}/\text{image}$  (only forward!  $\sim \times 2$  for backward)
- Total parameters: 138M parameters
- The most of the network parameters are in the FC layer and most of the variables (memory) required by the network is used in the first 2 Conv Layers.



# Visualizing and Understanding



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

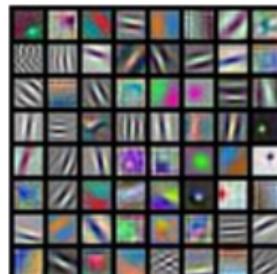
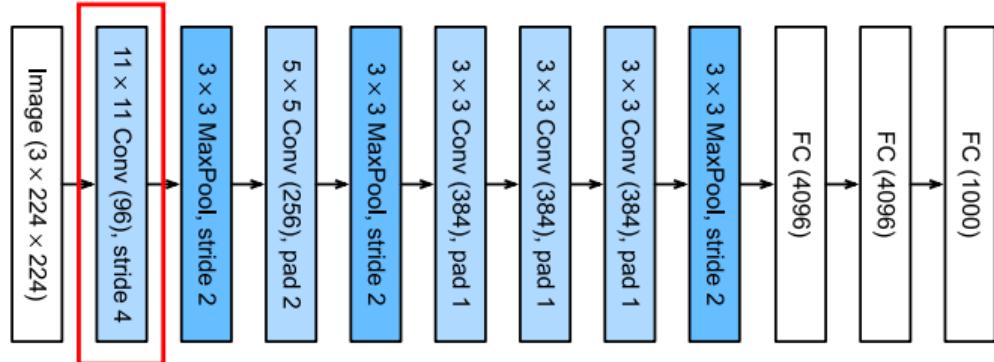
Applications

Datasets

Applications

# First Layer

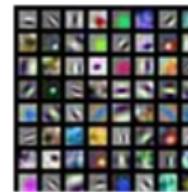
- Visualize first-layer filters/kernels (raw weights) directly



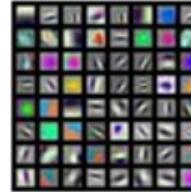
AlexNet:  
 $64 \times 3 \times 11 \times 11$



ResNet-18:  
 $64 \times 3 \times 7 \times 7$



ResNet-101:  
 $64 \times 3 \times 7 \times 7$



DenseNet-121:  
 $64 \times 3 \times 7 \times 7$



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications  
Datasets  
Applications

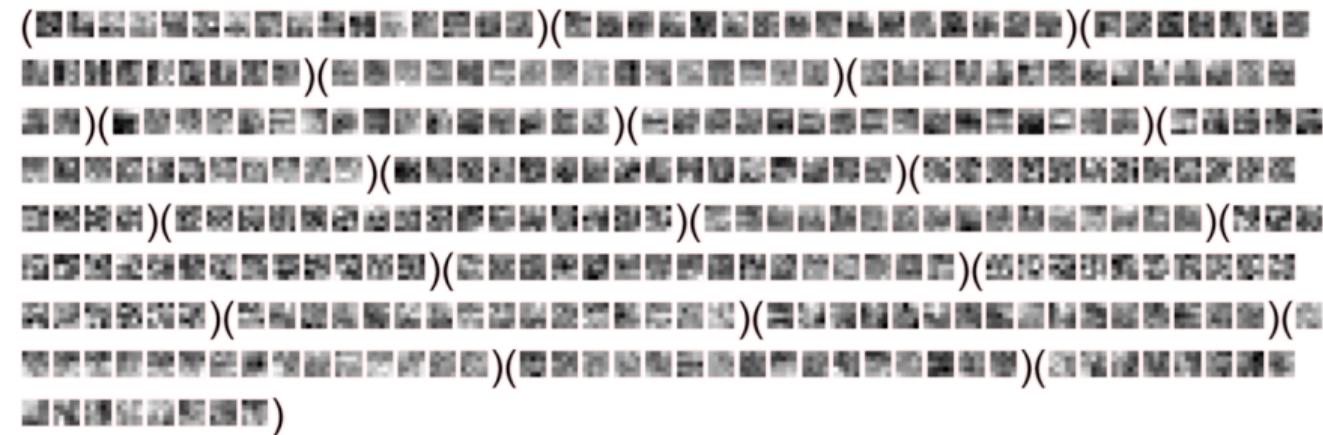
# First Layer (cont.)

- We can visualize filters at higher layers, but not that interesting (these pictures are taken from ConvNetJS CIFAR-10 demo)

Layer 1 filters



Layer 2 filters





Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

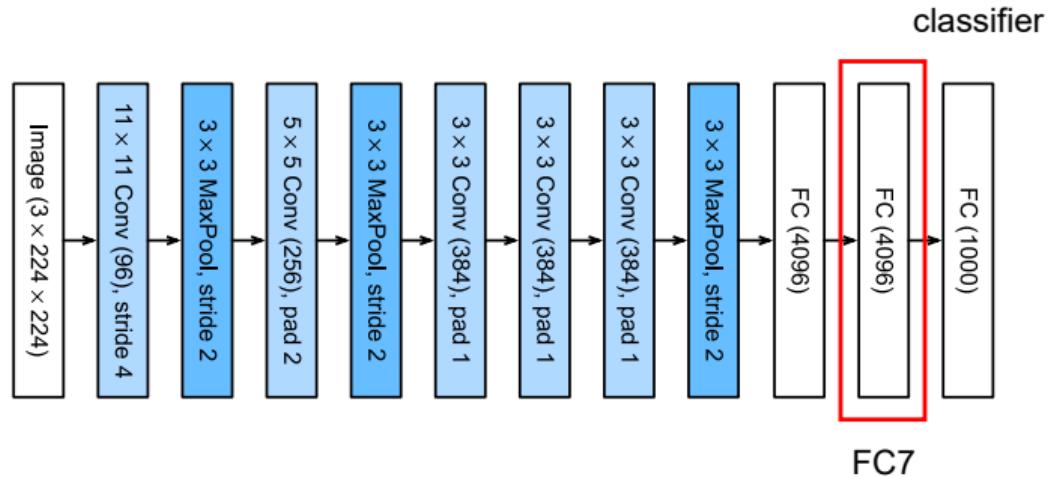
Dilated convolutions

Applications

Datasets

Applications

# Last Layer



- FC7: 4096-dimensional feature vector for an image (layer immediately before the classifier)
- Run the network on many images, collect the feature vectors



# Last Layer: Nearest Neighbors

Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

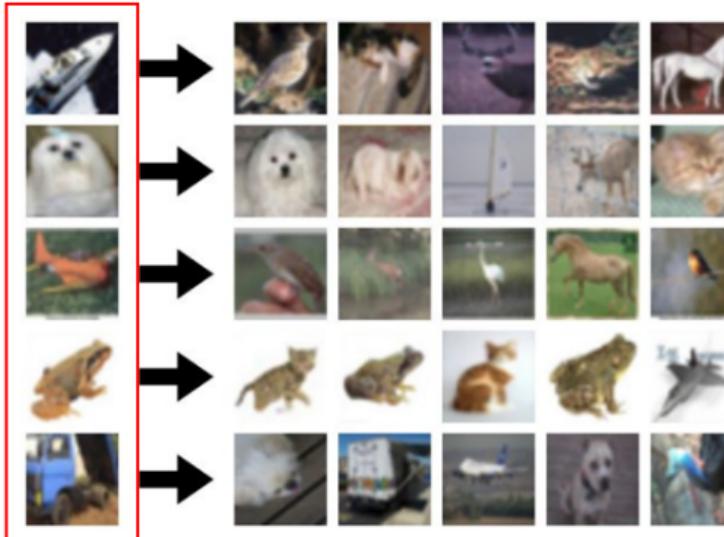
Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

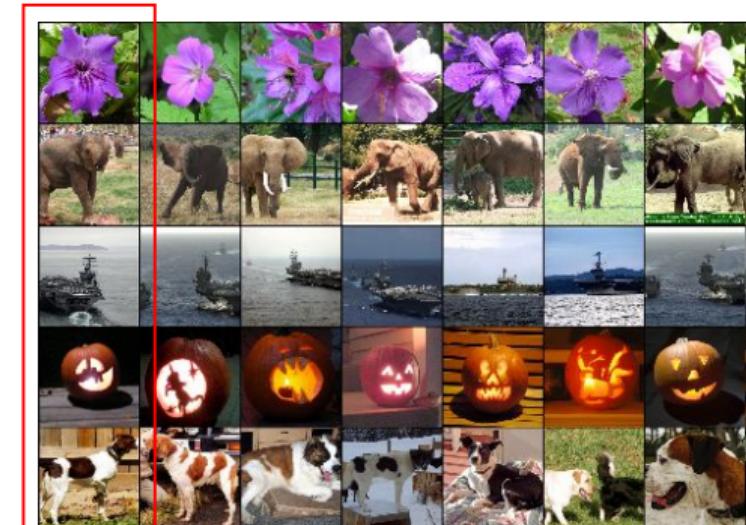
Applications  
Datasets  
Applications

Nearest neighbors in **pixel space**



test images

Nearest neighbors in **feature space**



test images



# Last Layer: Dimensionality Reduction

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

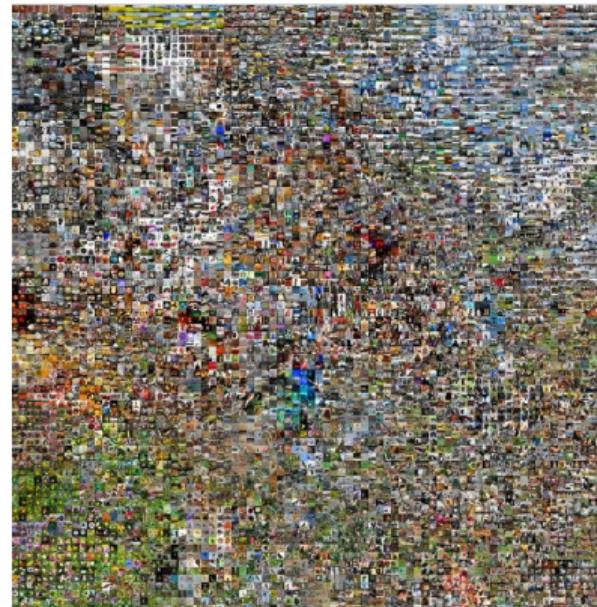
Dilated convolutions

Applications

Datasets

Applications

- Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions
  - Simple algorithm: Principal Component Analysis (PCA)
  - More complex: t-SNE





Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

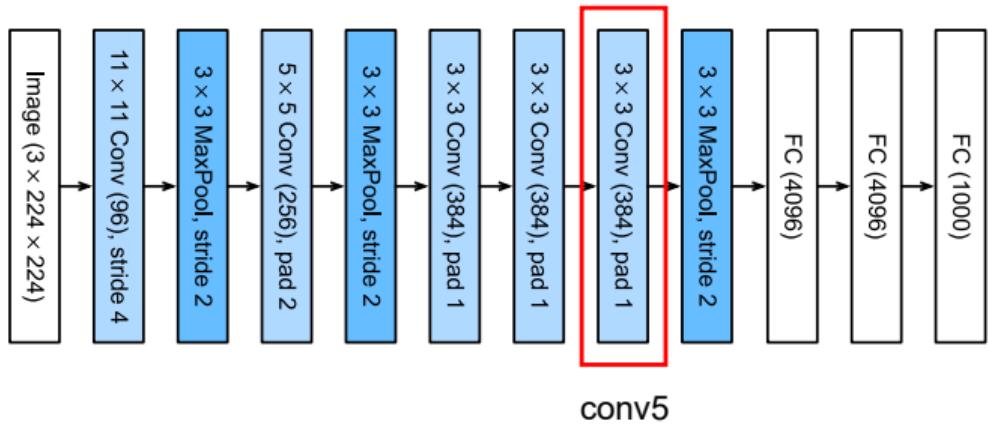
Transposed  
convolution  
Dilated convolutions

Applications

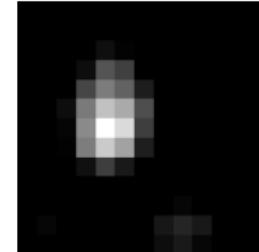
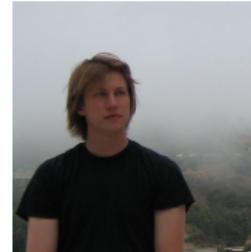
Datasets  
Applications

# Visualizing Activations

- A view of the  $13 \times 13$  activation of the 151<sup>st</sup> channel on the conv5 layer



Input image



Activation map



# Mapping activations back to pixels

Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

## (Guided) backpropagation method

- Forward an image through the network
- Choose a feature map and an activation
- Zero out all values except for the one of interest
- Propagate that value back to the image



# Which pixels matter: Saliency

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

Miscellaneous convolutions

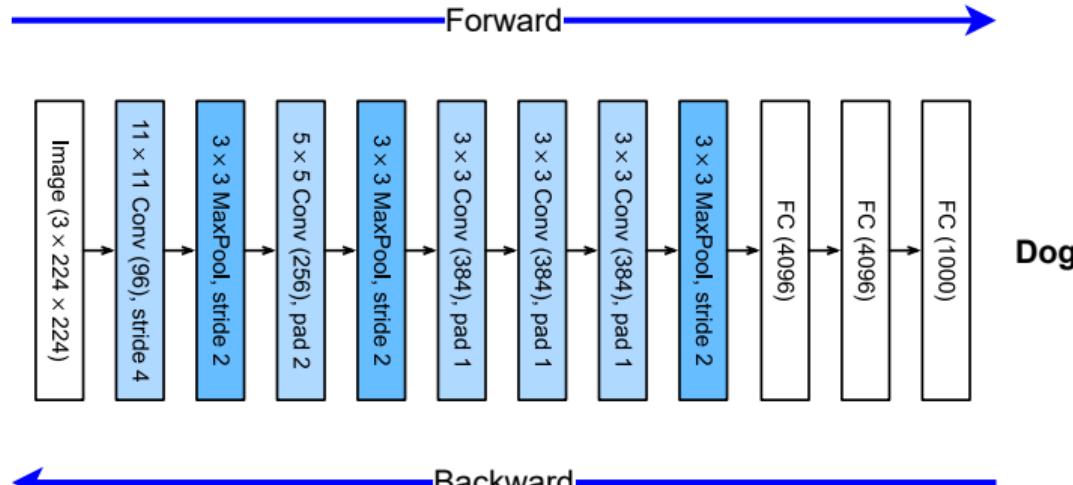
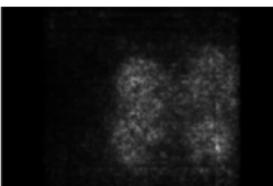
Transposed convolution

Dilated convolutions

Applications

Datasets  
Applications

- Compute gradient of (unnormalized) class score with respect to image pixels,
- Take absolute value and max over RGB channels





# Which pixels matter: Saliency (cont.)

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

- Superimposing the class activation heatmap on the original picture





# Generate Images: Gradient Ascent

## (Guided) backpropagation

**Find** the part of an image that a neuron responds to

## Gradient ascent

**Generate** a synthetic image that maximally activates a neuron

$$I_{syn} = \arg \max_I S_c(I) - \lambda \|I\|_2^2 \quad (5)$$

where  $S_c$  score for class  $c$  (before Softmax)

1. Initialize image  $I$  to zeros
2. Repeat:
  - 2.1 Forward image to compute current scores
  - 2.2 Backprop to get gradient of neuron value with respect to image pixels
  - 2.3 Make a small update to the image



# Generate Images: Gradient Ascent

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output

Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

Applications



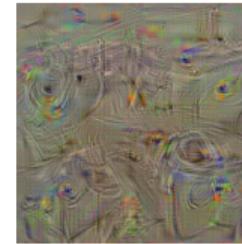
dumbbell



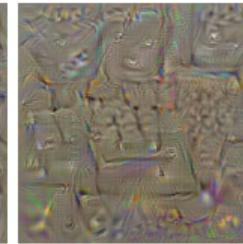
cup



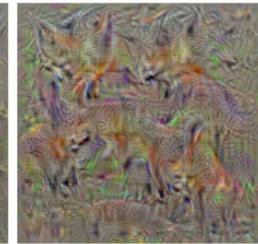
dalmatian



washing machine



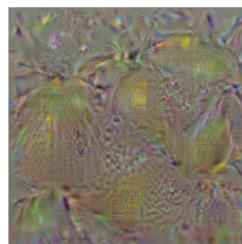
computer keyboard



kit fox



bell pepper



lemon



husky



goose



ostrich



limousine



# Network Design

- Networks Using Blocks
- Networks with Parallel Concatenations
- Residual Networks



# Blocks

---

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

Applications

Datasets

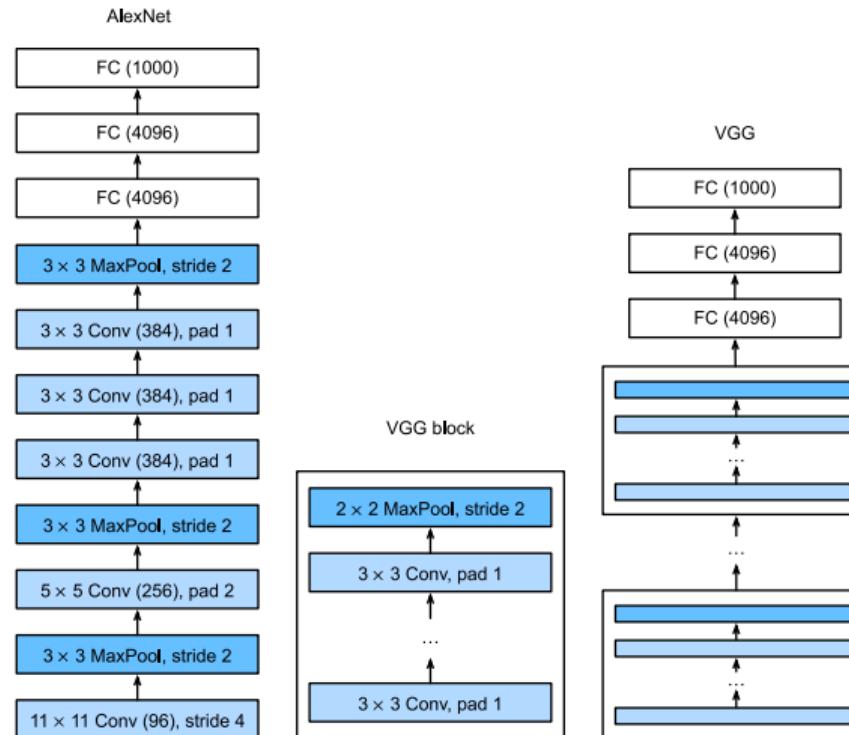
Applications

- The design of neural network architectures had grown progressively more abstract, moving from thinking in terms of individual neurons to whole layers, and now to blocks, repeating patterns of layers.
- The basic building block of classic CNNs is a sequence of the following:
  - a convolutional layer with padding to maintain the resolution
  - a nonlinearity such as a ReLU
  - a pooling layer such as a max pooling layer



# VGG Network

- From AlexNet to VGG that is designed from building blocks.





# Inception Blocks

Architecture  
Convolutional layer

Single Input, Single Output  
Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

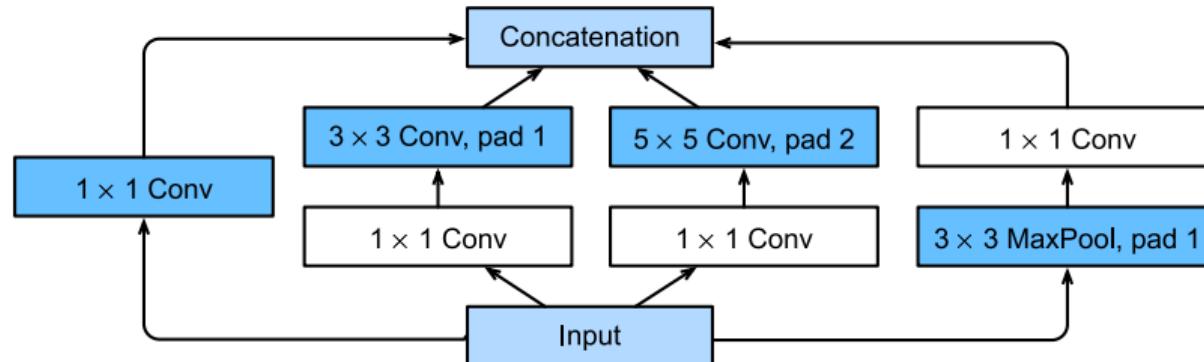
Dilated convolutions

Applications

Datasets

Applications

- Inception block employs a combination of variously-sized kernels





Architecture

Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

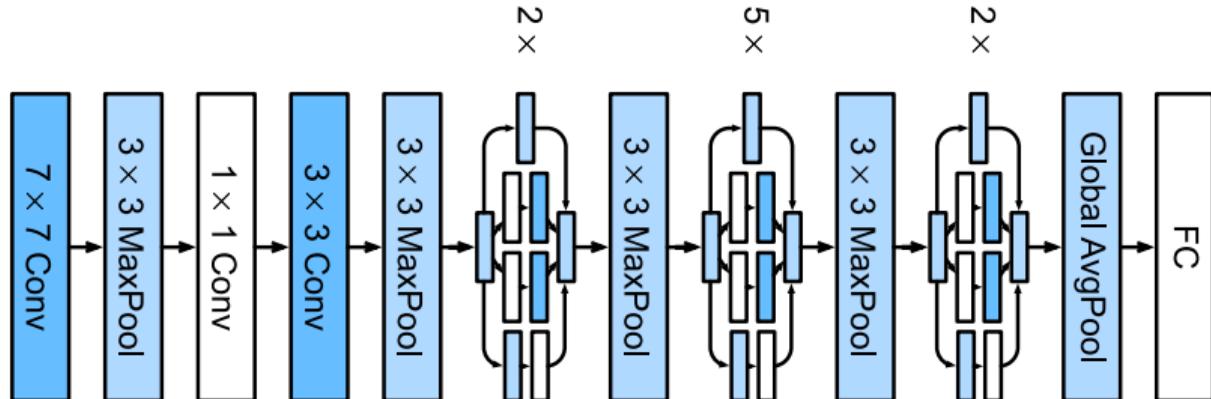
Applications

Datasets

Applications

# GoogleNet

- GoogLeNet uses a stack of a total of 9 inception blocks and global average pooling to generate its estimates.





# Learning Model

Architecture  
Convolutional layer

Single Input, Single Output  
Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design  
Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

Miscellaneous convolutions

Transposed convolution  
Dilated convolutions

Applications

Datasets  
Applications

- Consider  $\mathcal{F}$ , **the class of functions** (hypothesis set) that a specific network architecture can reach.
- Given a data set  $\mathcal{D}$ , algorithm  $A$  finds the best estimated  $\hat{f}$  of the “truth” function  $f^*$

$$\hat{f} = \arg \min_f L(f \mid \mathcal{D}) \text{ subject to } f \in \mathcal{F}. \quad (6)$$



Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks

Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

Dilated convolutions

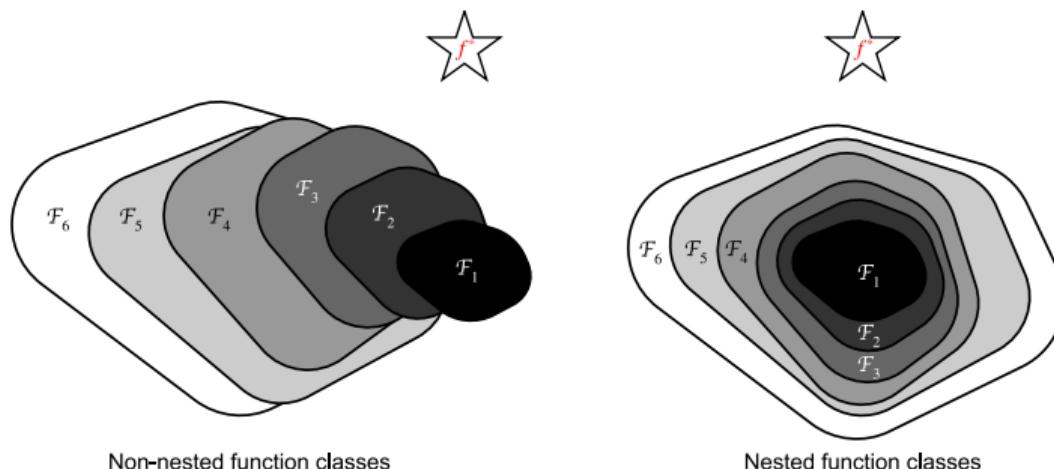
Applications

Datasets

Applications

# Learning Model (cont.)

- For non-nested function classes, a larger (indicated by area) function class does not guarantee to get closer to the “truth” function  $f^*$ . This does not happen in nested function classes.





Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design  
Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

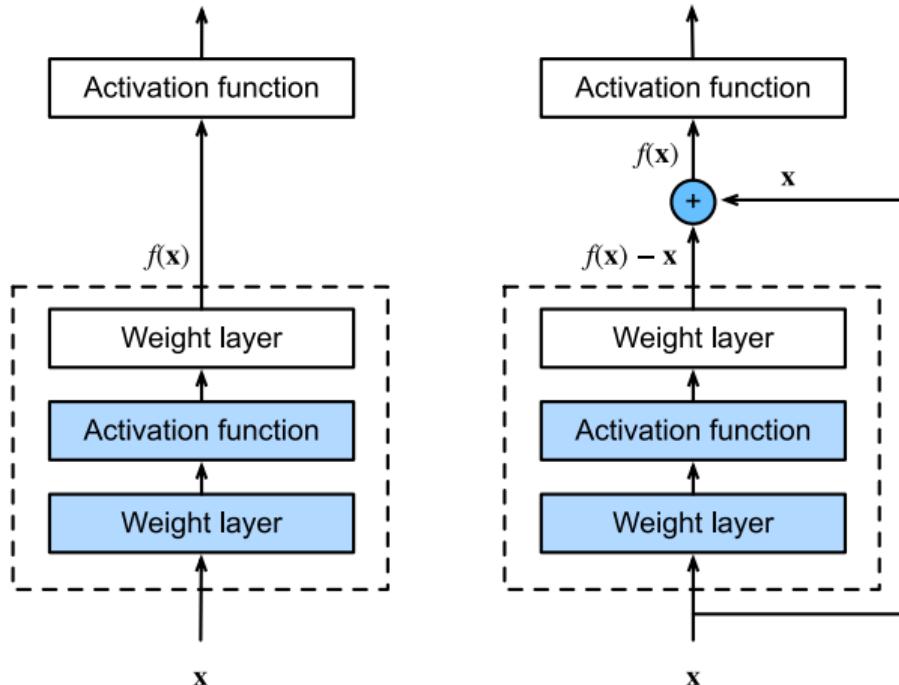
Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications  
Datasets  
Applications

# Residual Blocks

- A regular block (left) and a residual block (right).





Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

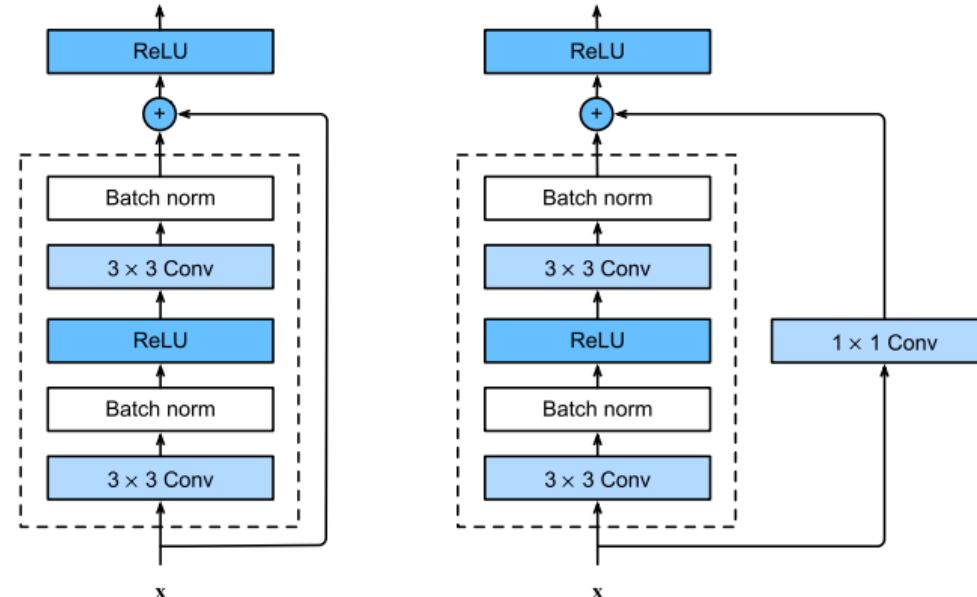
Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

# Residual Blocks (cont.)

- ResNet block with and without  $1 \times 1$  convolution.





# ResNet Model

Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations

Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution

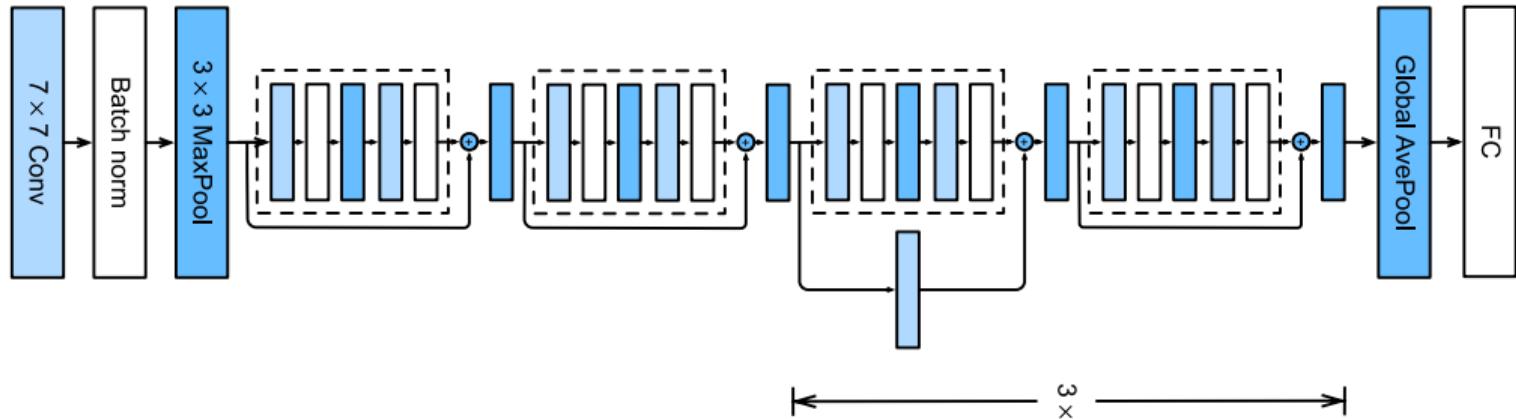
Dilated convolutions

Applications

Datasets

Applications

- The ResNet-18 architecture





## Miscellaneous convolutions

- Transposed convolution
- Dilated convolutions



# Transposed convolution

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets

Applications

kernel  $w \rightarrow$  matrix  $C$

- Convolution can be interpreted as a matrix multiplication

convolution  $C \longleftrightarrow$  transposed convolution  $C^T$



# Transposed convolution (cont.)

Architecture  
Convolutional layer

Single Input, Single Output  
Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks  
Networks with Parallel Concatenations  
Residual Networks

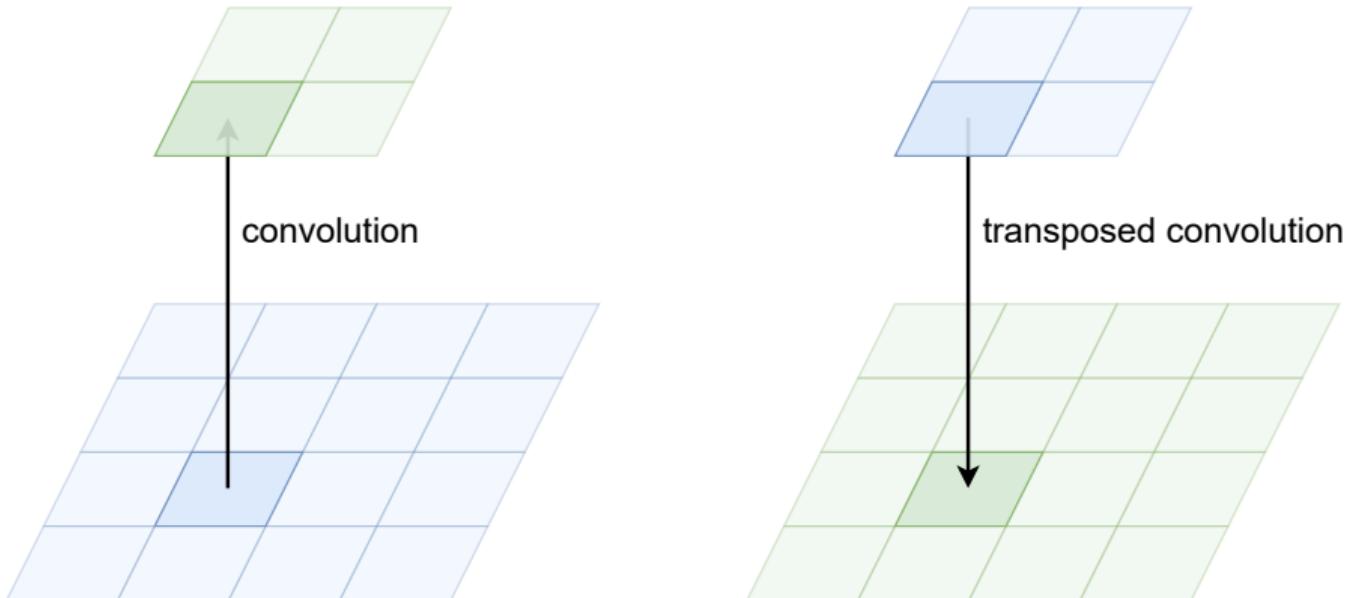
Miscellaneous convolutions

Transposed convolution

Dilated convolutions

Applications

Datasets  
Applications





# Transposed Convolution (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

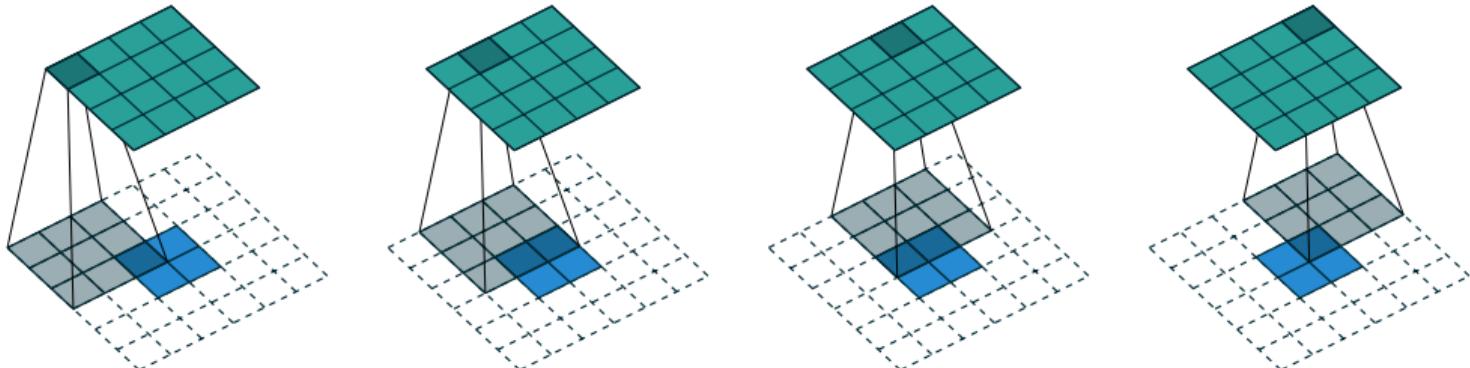
Applications

Datasets

Applications

- No zero padding, unit strides, transposed

	convolution	transposed convolution
<b>kernel size</b>	$k$	$k' = k$
<b>stride</b>	$s = 1$	$s' = 1$
<b>padding</b>	$p = 0$	$p' = k - 1$
<b>output size</b>		$o' = i' + (k - 1)$





# Transposed Convolution (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

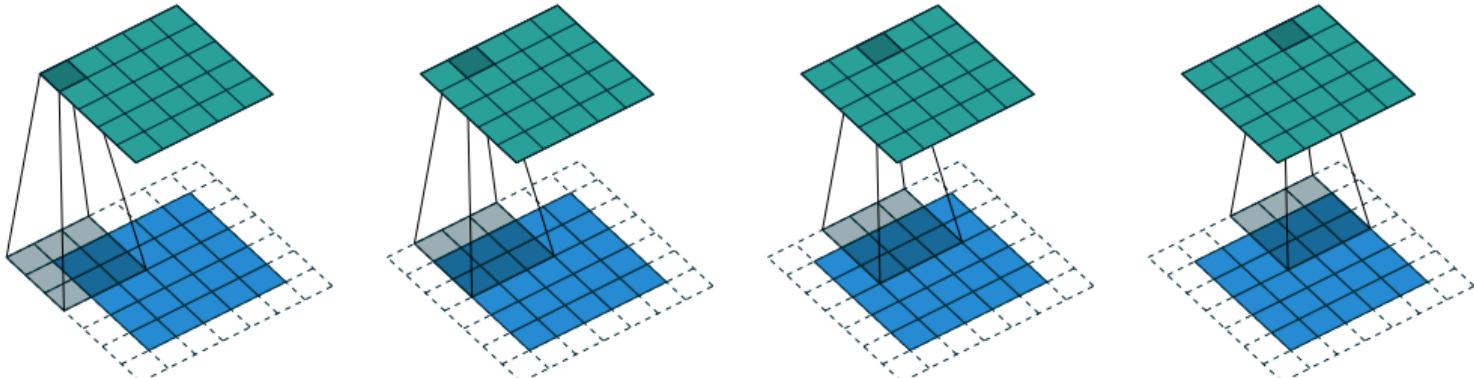
Applications

Datasets

Applications

- Zero padding, unit strides, transposed

	convolution	transposed convolution
<b>kernel size</b>	$k$	$k' = k$
<b>stride</b>	$s = 1$	$s' = 1$
<b>padding</b>	$p$	$p' = k - p - 1$
<b>output size</b>		$o' = i' + (k - 1) - 2p$





# Transposed Convolution (cont.)

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output

Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

Dilated convolutions

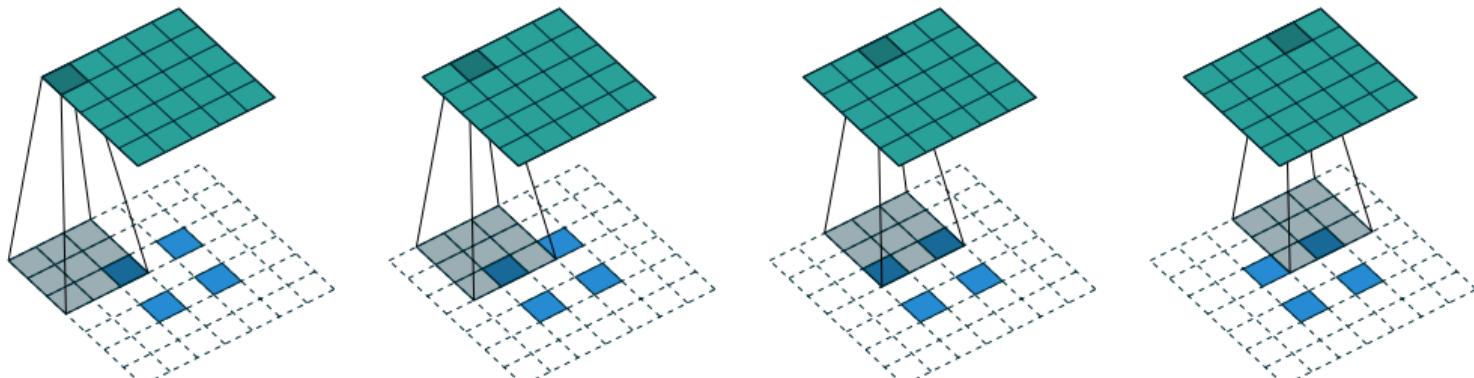
Applications

Datasets

Applications

- No zero padding, non-unit strides, transposed:
  - adding  $s - 1$  zeros between each input unit

	<b>convolution</b>	<b>transposed convolution</b>
<b>kernel size</b>	$k$	$k' = k$
<b>stride</b>	$s$	$s' = 1$
<b>padding</b>	$p = 0$	$p' = k - 1$
<b>output size</b>		$o' = s(i' - 1) + k$



[Architecture](#)[Convolutional  
layer](#)[Single Input, Single  
Output](#)[Multiple Input,  
Multiple Output](#)[Tensor Operation](#)[Pooling layer](#)[Fully connected  
layer](#)[Visualizing and  
Understanding](#)[Network Design](#)[Networks Using Blocks](#)[Networks with Parallel  
Concatenations](#)[Residual Networks](#)[Miscellaneous  
convolutions](#)[Transposed  
convolution](#)[Dilated convolutions](#)[Applications](#)[Datasets](#)[Applications](#)

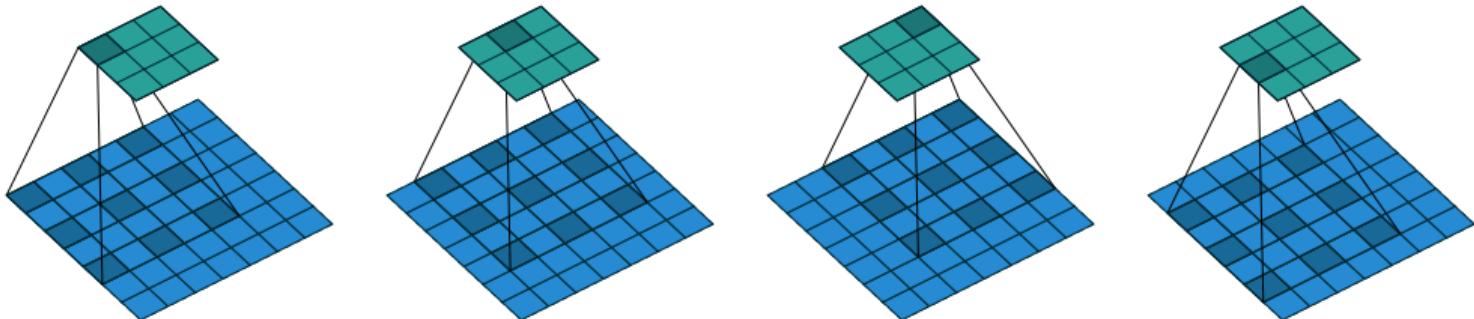
# Dilated convolutions

## Concept 1

**Dilated convolutions** “inflate” the kernel by inserting spaces between the kernel elements. The dilation “rate” is controlled by an additional hyperparameter  $d$ .

- For any  $i, k, s$  and  $p$ , and for a dilation rate  $d$ ,

$$o = \left\lfloor \frac{i + 2p - k - (k - 1)(d - 1)}{s} \right\rfloor + 1. \quad (7)$$





# Applications

- Datasets
- Applications



# MNIST

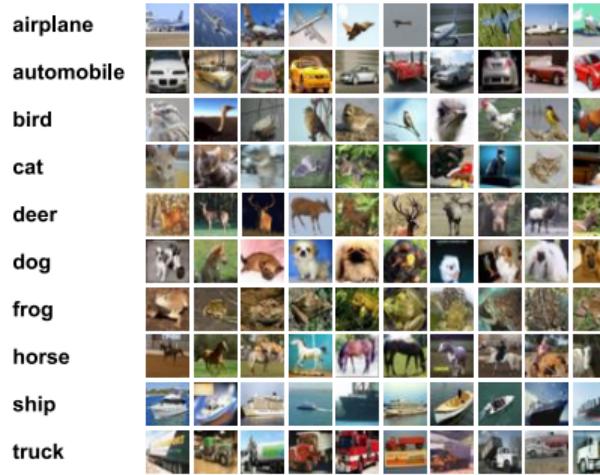
- The MNIST database is a large database of handwritten digits that is commonly used for training various image processing systems
- The MNIST database contains 60,000 training images and 10,000 testing images

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



# CIFAR

- The CIFAR-10 and CIFAR-100 are labeled subsets of the 80 million tiny images dataset. They were collected by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton.
- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.





Architecture  
Convolutional  
layer

Single Input, Single  
Output

Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications

Datasets  
Applications

# ImageNet

- The ImageNet project is a large visual database designed for use in visual object recognition software research





# Image classification

Architecture  
Convolutional layer

Single Input, Single Output

Multiple Input, Multiple Output  
Tensor Operation

Pooling layer

Fully connected layer

Visualizing and Understanding

Network Design

Networks Using Blocks

Networks with Parallel Concatenations

Residual Networks

Miscellaneous convolutions

Transposed convolution

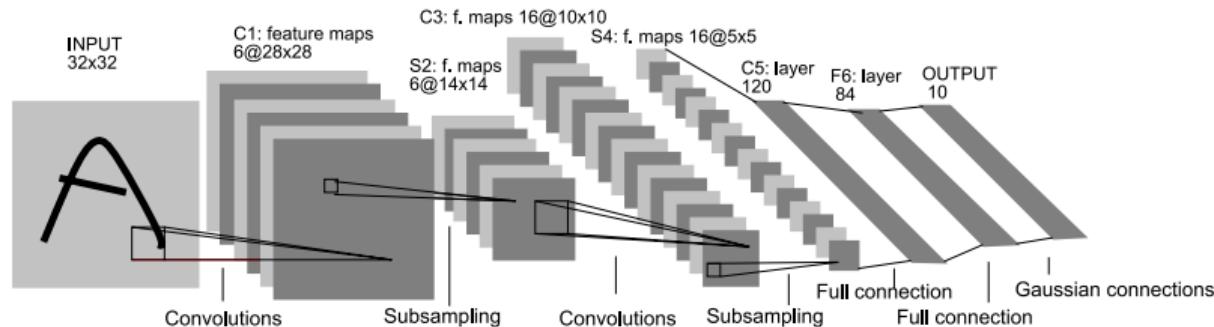
Dilated convolutions

Applications

Datasets

Applications

- Architecture of LeNet-5 a Convolutional Neural Network, here for digit recognition, but is extended to character recognition.





Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

Transposed  
convolution  
Dilated convolutions

Applications

Datasets

Applications

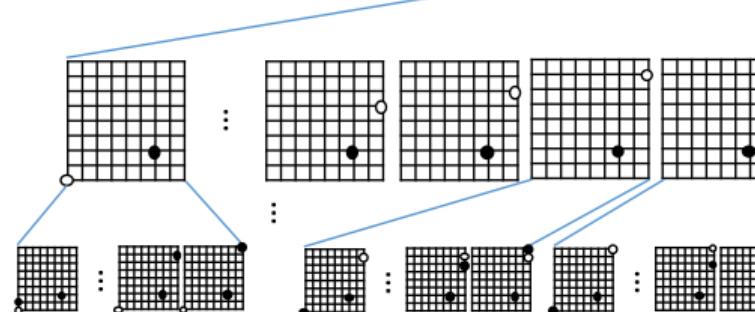
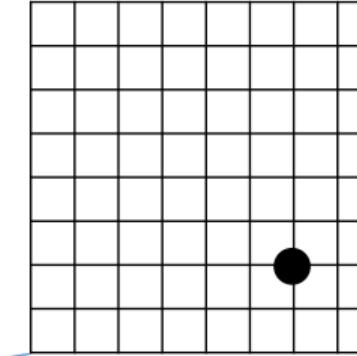
# Games

- Go game

Game tree

Depth

$d = 0$



$d = 1$

$d = 2$



Architecture  
Convolutional  
layer

Single Input, Single  
Output  
Multiple Input,  
Multiple Output  
Tensor Operation

Pooling layer

Fully connected  
layer

Visualizing and  
Understanding

Network Design

Networks Using Blocks  
Networks with Parallel  
Concatenations  
Residual Networks

Miscellaneous  
convolutions

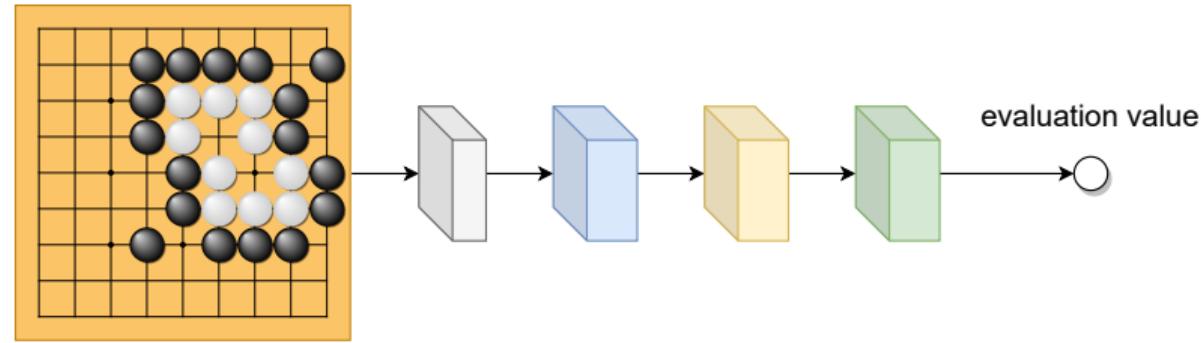
Transposed  
convolution  
Dilated convolutions

Applications

Datasets

Applications

# Games (cont.)





## References

---

- Goodfellow, I., Bengio, Y., and Courville, A. (2016).  
*Deep learning.*  
MIT press.
- Lê, B. and Tô, V. (2014).  
*Cở sở trí tuệ nhân tạo.*  
Nhà xuất bản Khoa học và Kỹ thuật.
- Russell, S. and Norvig, P. (2021).  
*Artificial intelligence: a modern approach.*  
Pearson Education Limited.