

<https://introml.mit.edu/>

6.390 Intro to Machine Learning

Lecture 7: Convolutional Neural Networks

Shen Shen

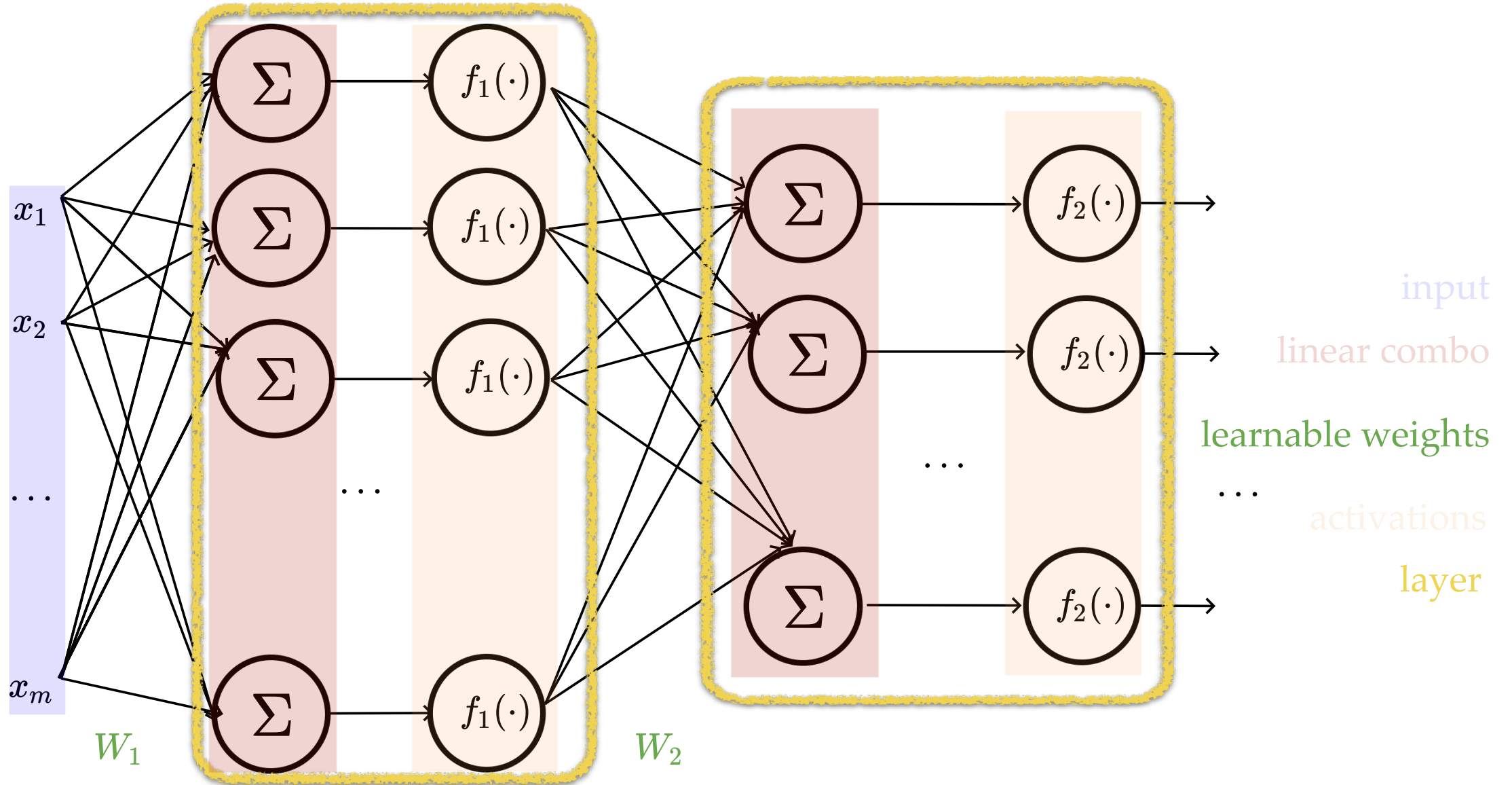
March 22, 2024

(videos edited from [3b1b](#); some slides adapted from [Phillip Isola](#) and [Kaiming He](#))

Outline

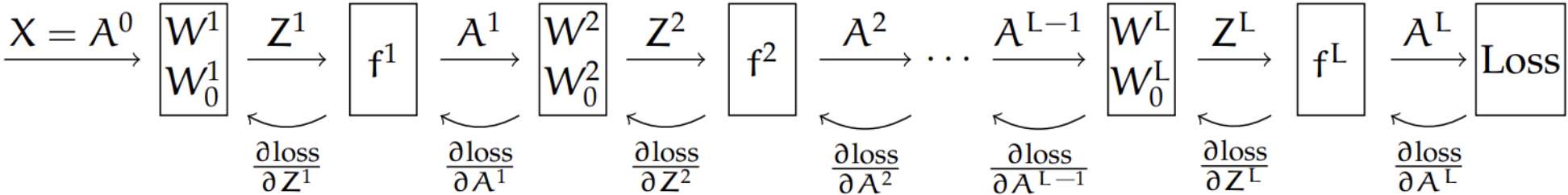
- Recap (fully-connected net)
- Motivation and big picture ideas of CNN
- Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing
 - 3d tensors
- Max pooling
 - Larger window
- Typical architecture and summary

A (feed-forward) neural network is



Recap: Backpropogation

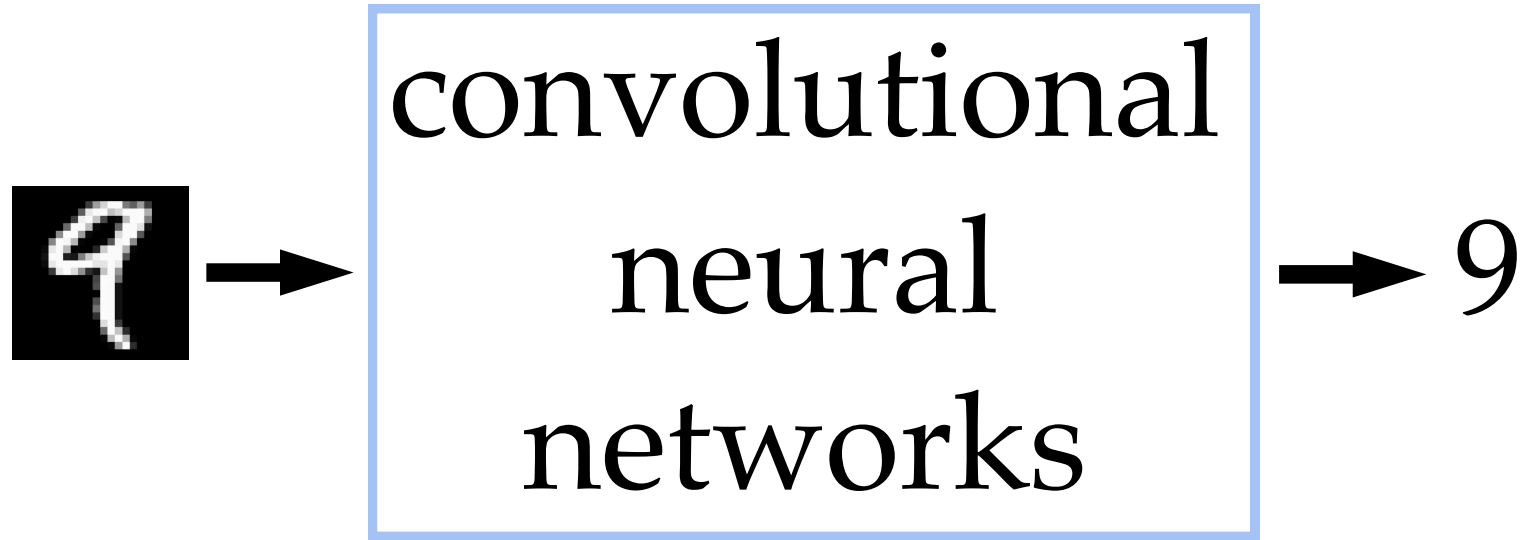
Forward propagation to obtain the output (model's guess)



Backpropagation to obtain gradients with respect to the loss

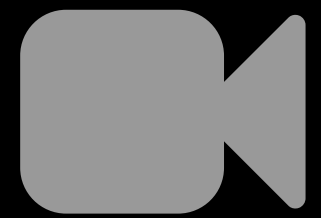
Outline

- Recap (fully-connected net)
- Motivation and big picture ideas of CNN
- Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing
 - 3d tensors
- Max pooling
 - Larger window
- Typical architecture and summary

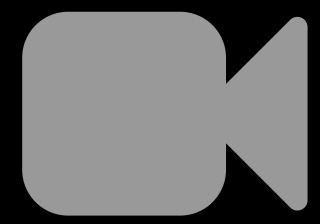
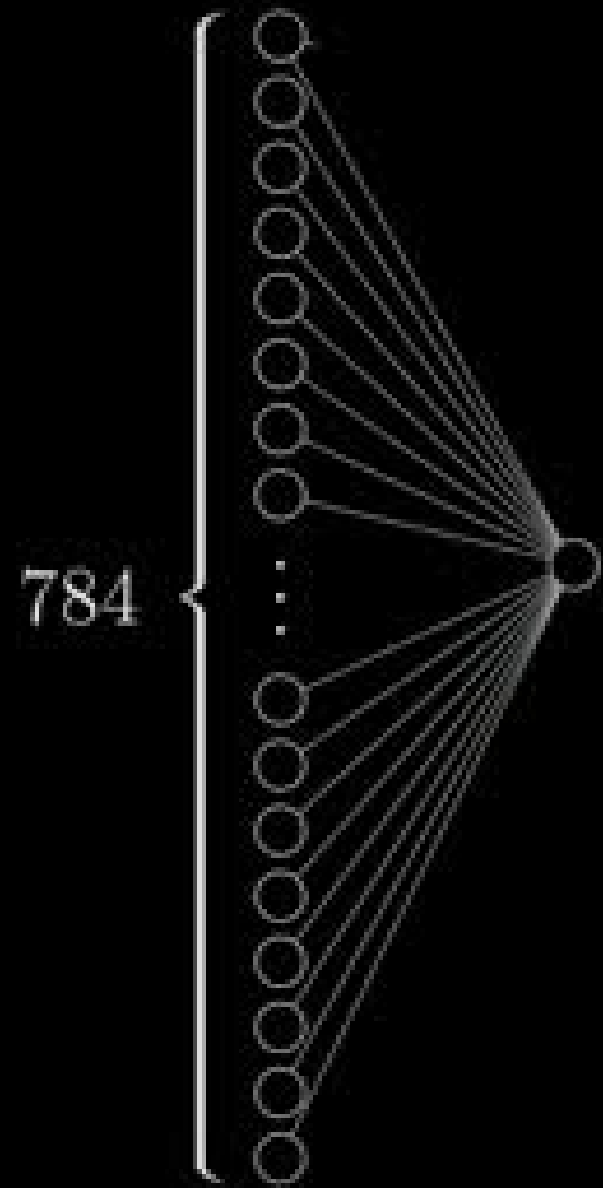


1. Why do we need a special network for images?
2. Why is CNN (the) special network for images?

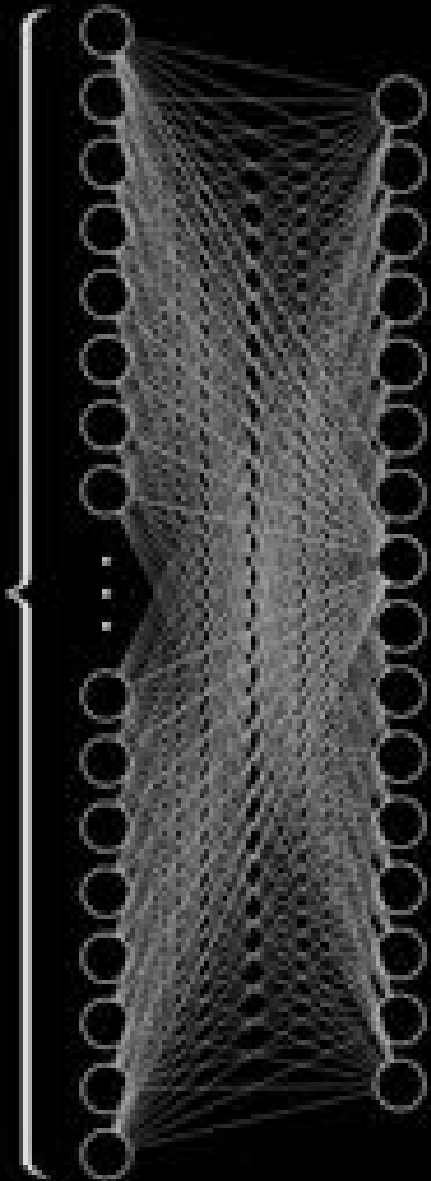
Why do we
need a special net for images?



784 weights per neuron

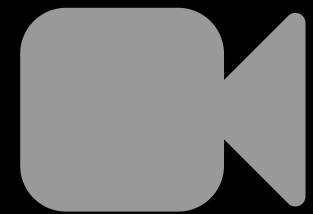


784



784×16 weights

16 biases





👉 426-by-426 grayscale image

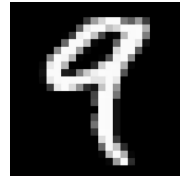
Use the same small 2-layer network?
need to learn $\sim 3\text{M}$ parameters

Imagine even higher-resolution images, or more complex tasks...

Q: Why do we need a specialized network?

A: fully-connected nets don't scale well to (interesting) images

Why do we think



is

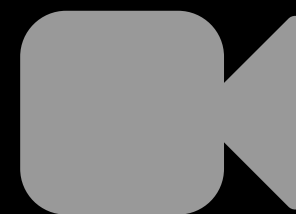
9?

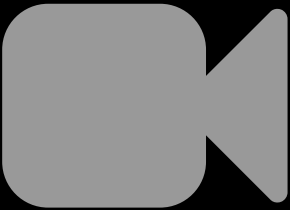
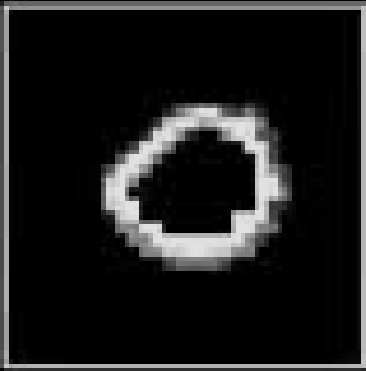
Why do we think any of

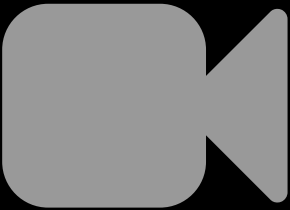
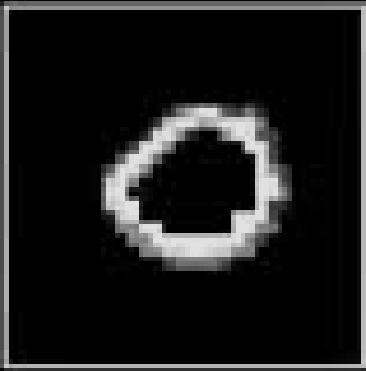


is a

9?



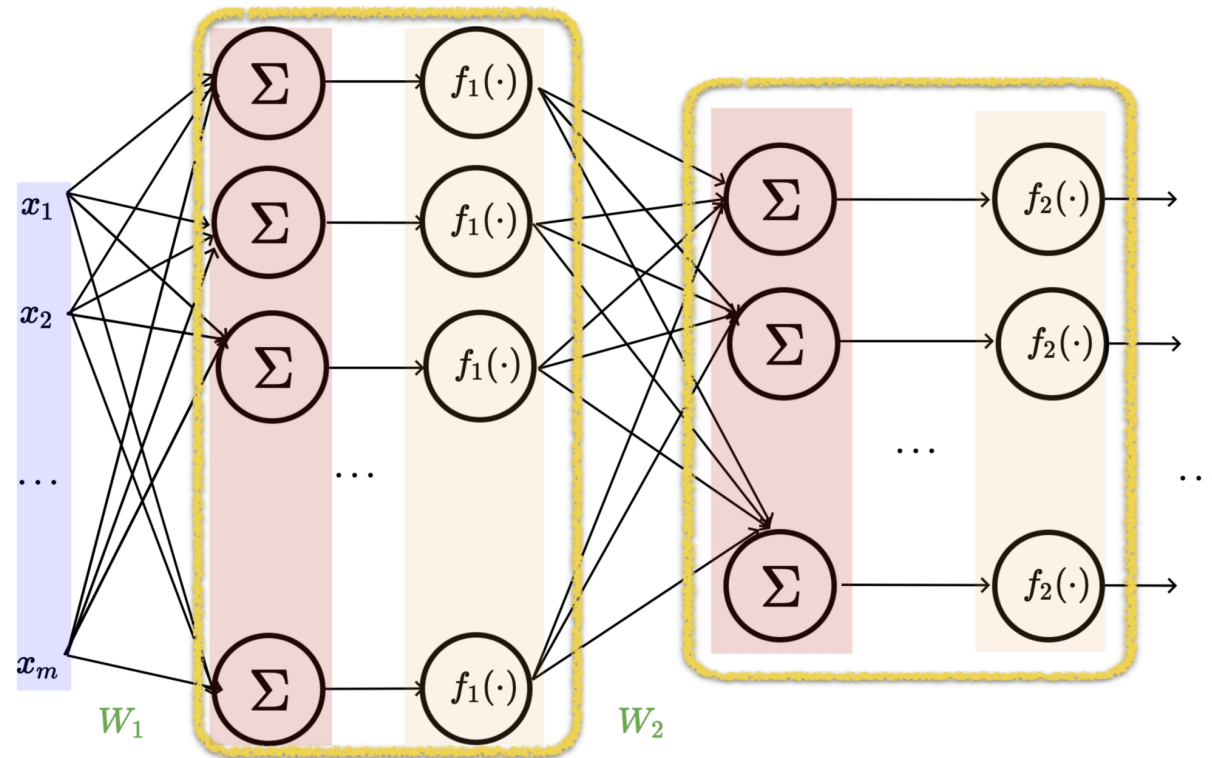




- Visual hierarchy



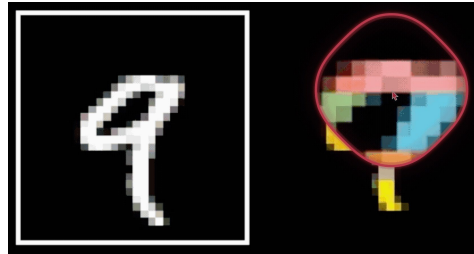
layering would help
take care of that



- Visual hierarchy



- Spatial locality



- Translational invariance



CNN cleverly exploits

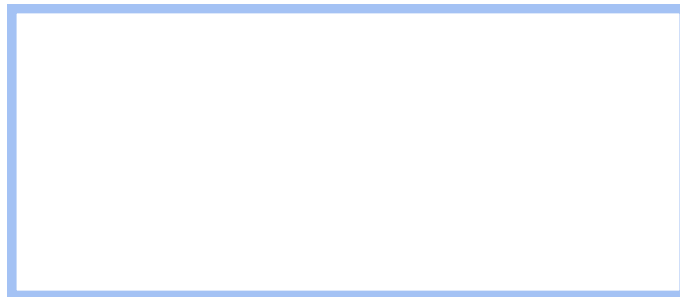
- Visual hierarchy
- Spatial locality
- Translational invariance

via

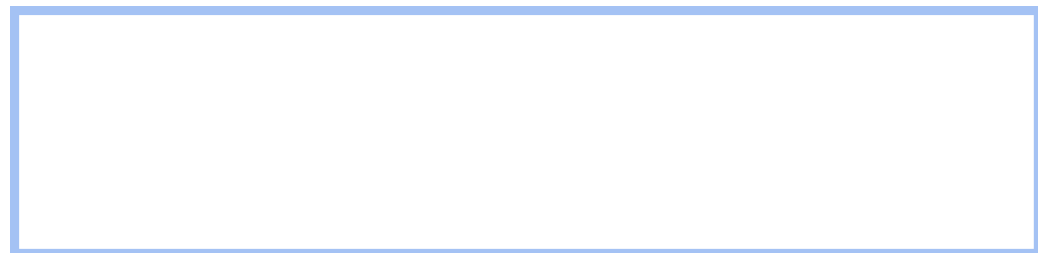
- layering (with nonlinear activations)
- convolution
- pooling

to handle images efficiently

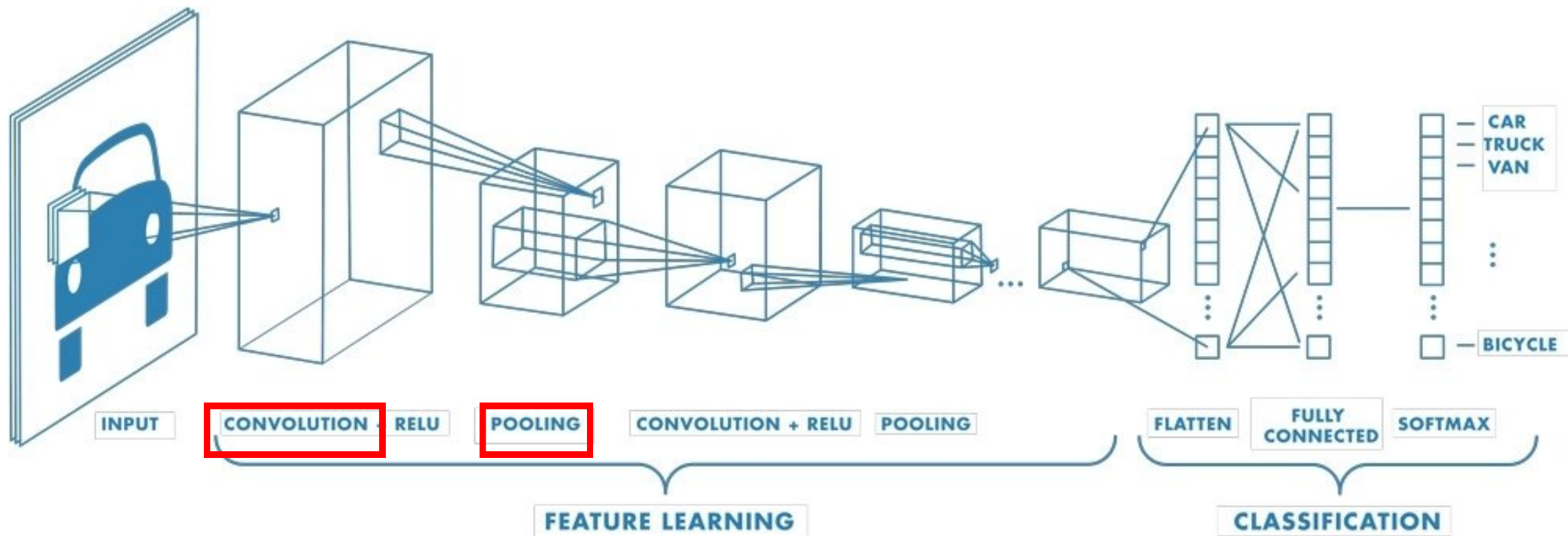
cleverly exploits



via



to handle efficiently



Outline

- Recap (fully-connected net)
- Motivation and big picture ideas of CNN

- Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing

- 3d tensors

- Max pooling
 - Larger window
- Typical architecture and summary

Convolutional layer might sound foreign, but...

Layer	Forward-pass {do}	Back-prop {learn}
Fully-connected	Dot-product	Neurons
Convolutional	Convolution	Filters (kernels)

input image

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

filter

-1	1
----	---

$$(0 * -1) + (1 * 1) = 1$$

output image

1			
---	--	--	--

input image

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

filter

-1	1
----	---

$$(1 * -1) + (0 * 1) = -1$$

output image

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

input image

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

filter

-1	1
----	---

$$(0 * -1) + (1 * 1) = 1$$

output image

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

input image

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

filter

-1	1
----	---

$$(1 * -1) + (1 * 1) = 0$$

output image

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

input image

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

filter

-1	1
----	---

$$(0 * -1) + (1 * 1) = 1$$

$$(-1 * -1) + (1 * 1) = 2$$

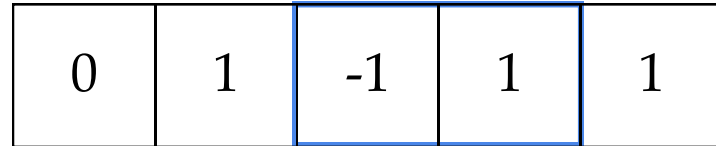
output image

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

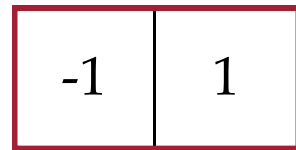
- 'look' locally
- parameter sharing
- "template" matching

- 'look' locally

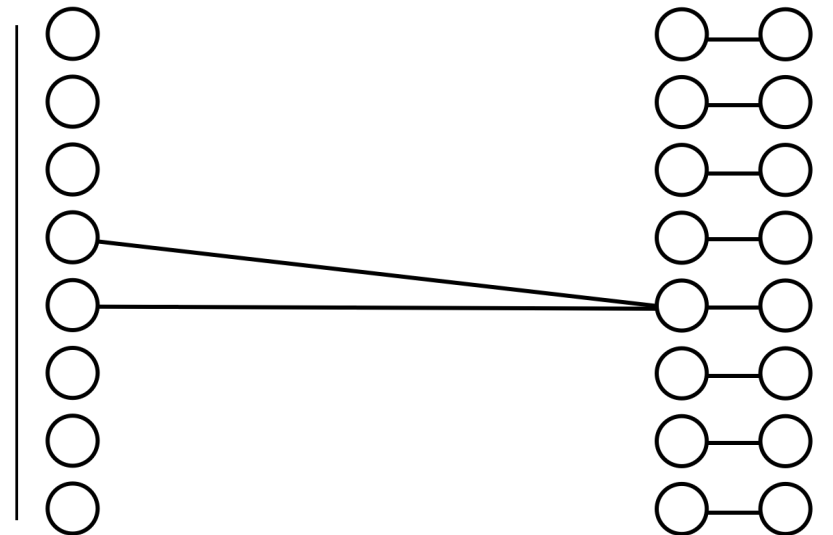
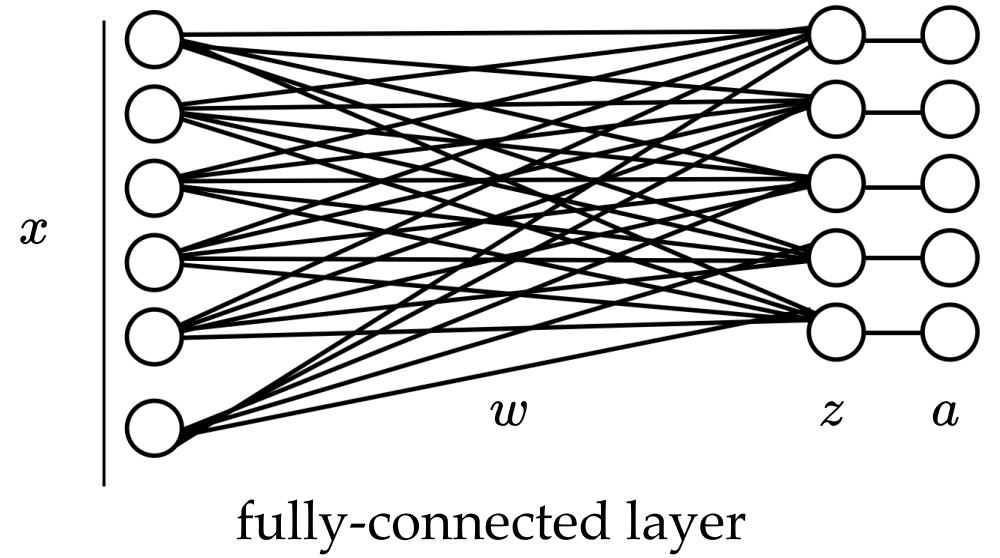
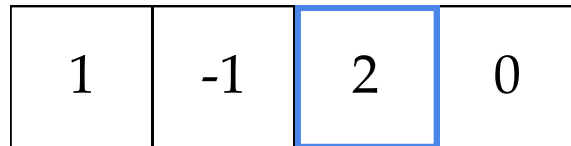
input image



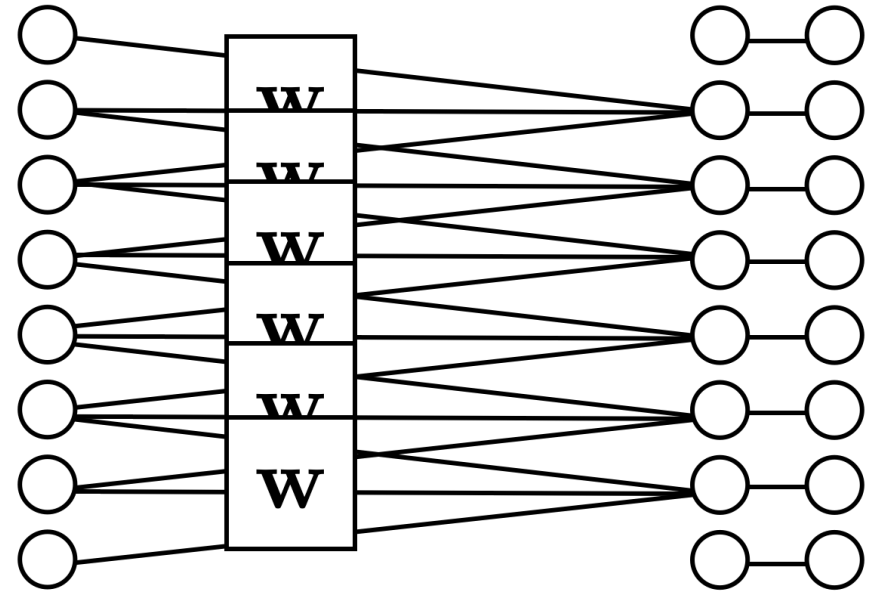
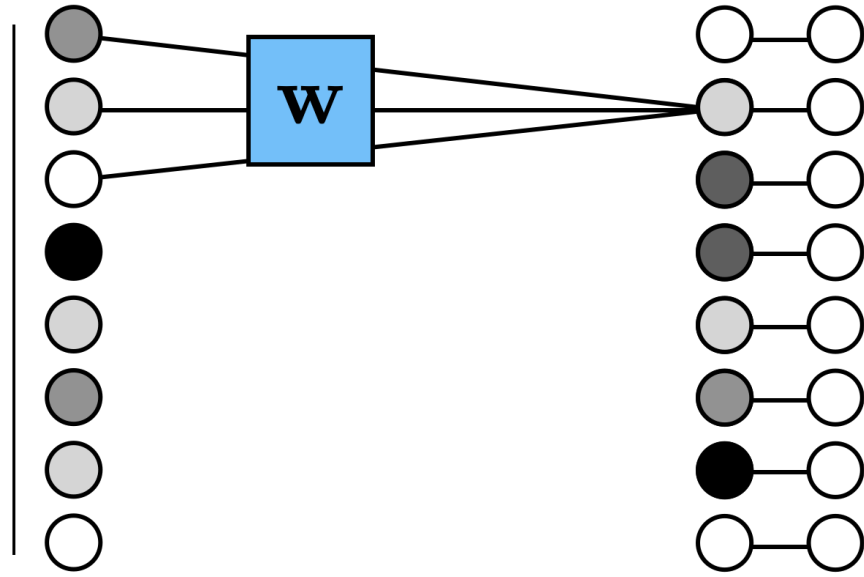
filter



output image



- parameter sharing



- parameter sharing

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

convolve
with

-1	1
----	---

or dot
with

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

=

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



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

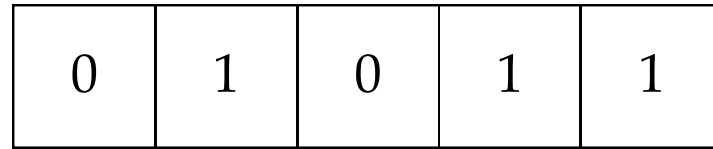
convolve with ?

=

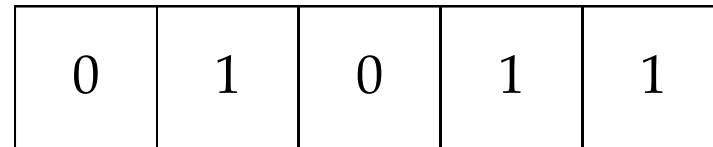
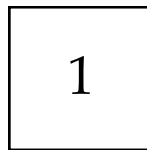
dot-product with ?

=

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



convolve with



dot-product with

$I_{5 \times 5}$



Convolution: a 2-D example

input

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

filter

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

output

-3	-4	-4	-4		

Convolution: a 2-D example

input

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

filter

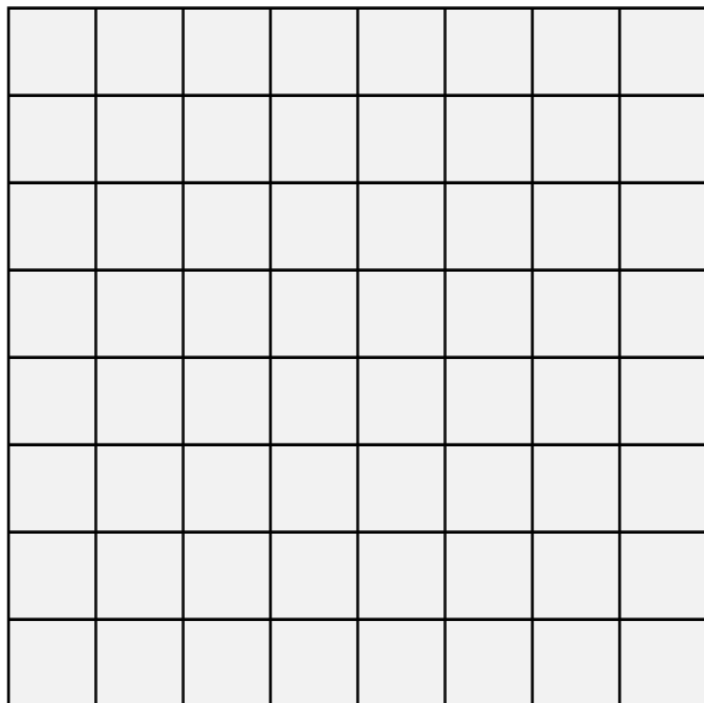
1	2	1
0	0	0
-1	-2	-1

output

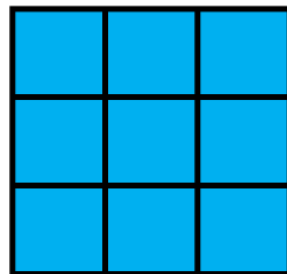
-3	-4	-4	-4	-4	-3
-3	-4	-4	-3	-1	0
0	0	0	0	0	0
2	1	0	1	3	3
2	1	0	1	3	3
1	3	4	3	1	0

Convolution: padding

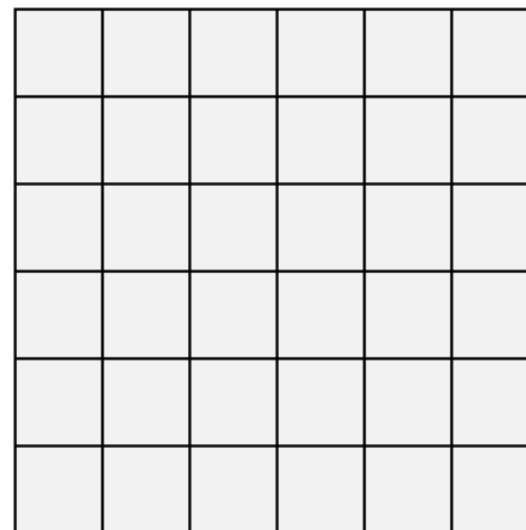
input: $H \times W = 8 \times 8$



filter

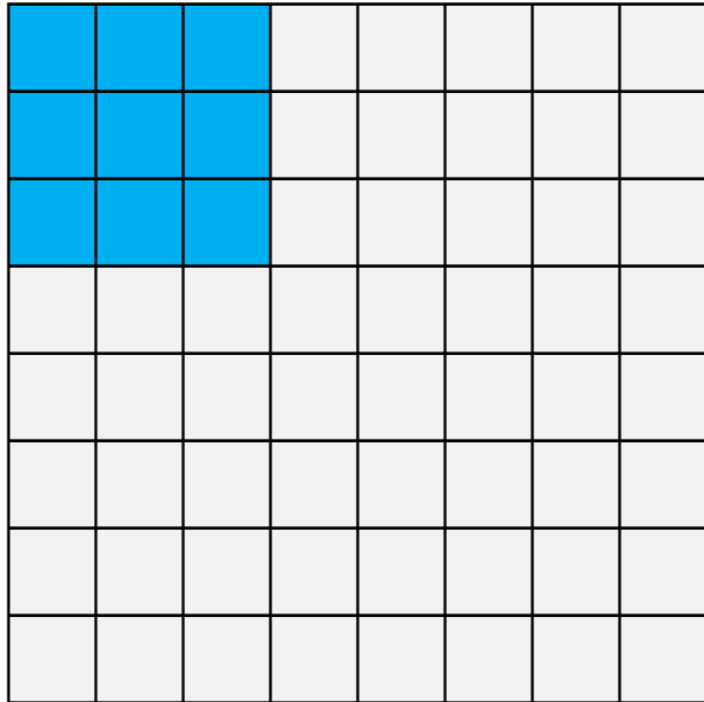


output: $H \times W = 6 \times 6$

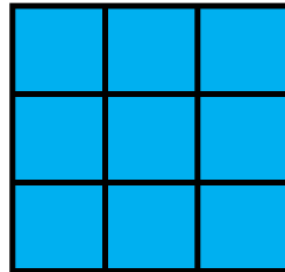


Convolution: padding

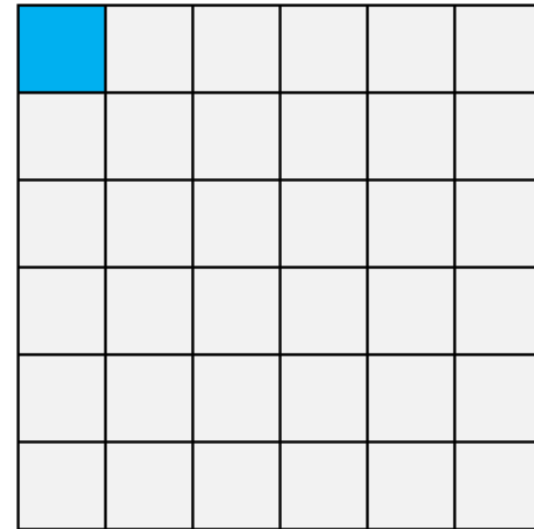
input: $H \times W = 8 \times 8$



filter

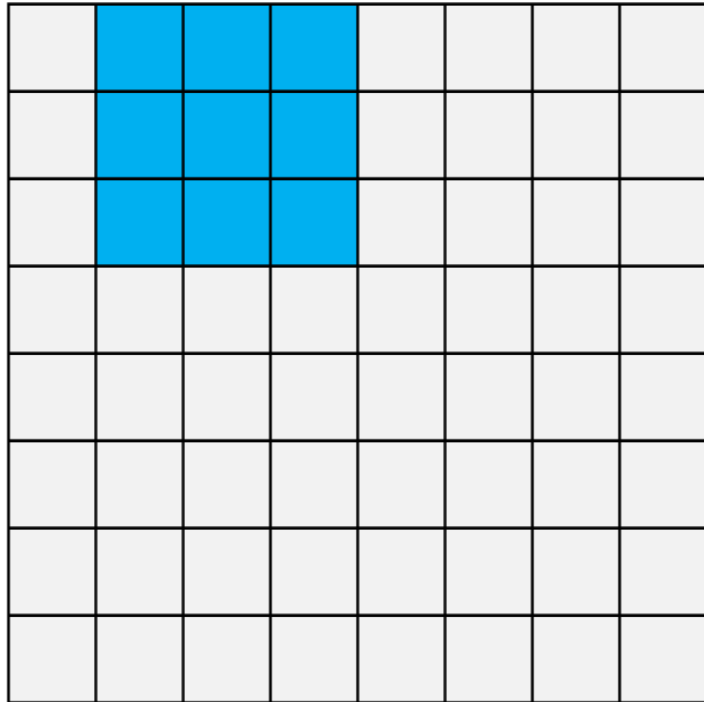


output: $H \times W = 6 \times 6$

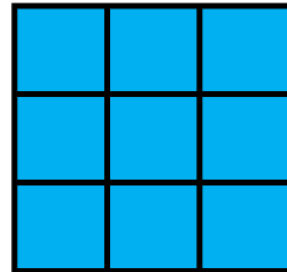


Convolution: padding

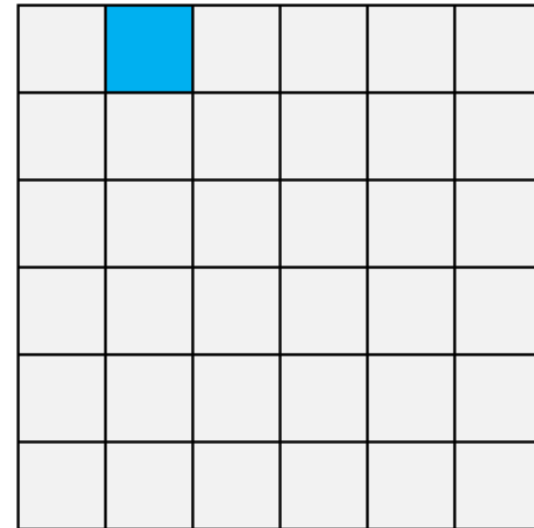
input: $H \times W = 8 \times 8$



filter

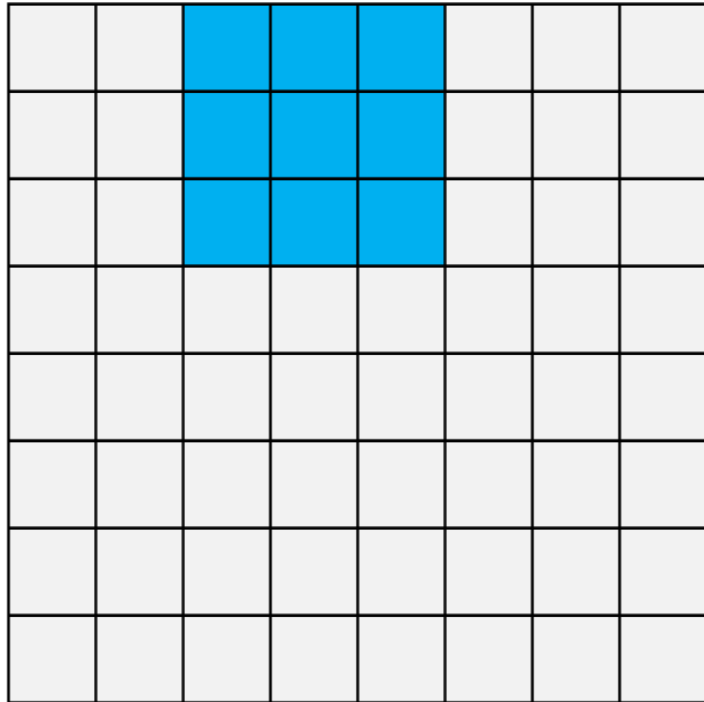


output: $H \times W = 6 \times 6$

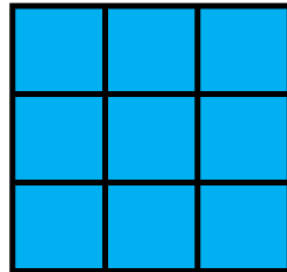


Convolution: padding

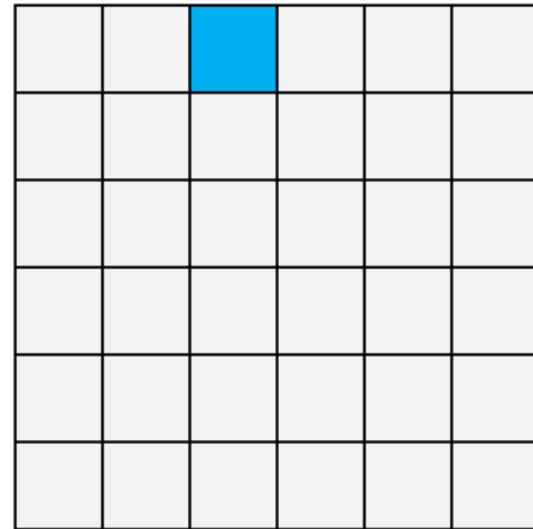
input: $H \times W = 8 \times 8$



filter

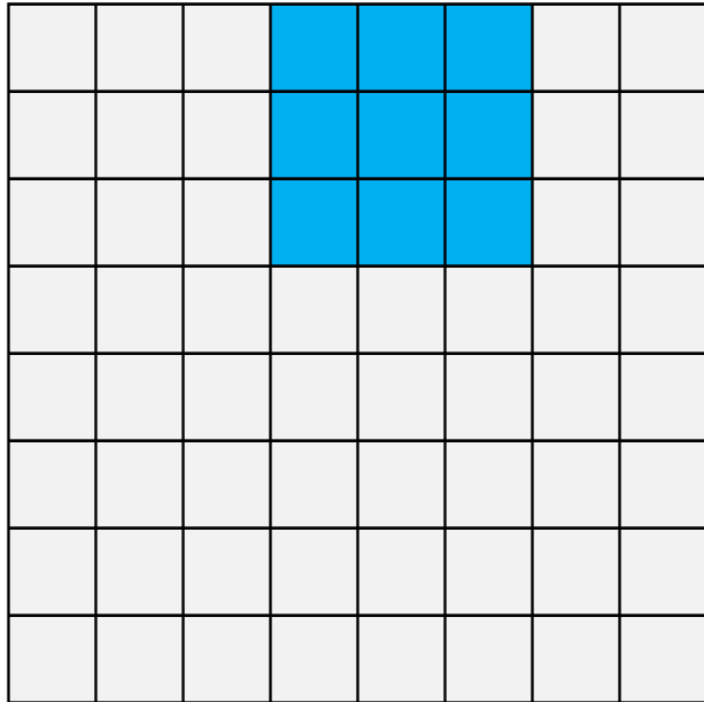


output: $H \times W = 6 \times 6$

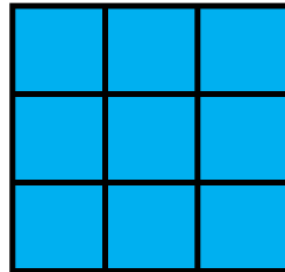


Convolution: padding

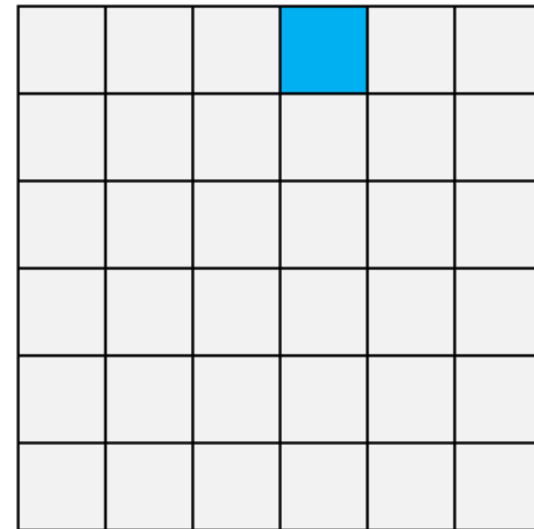
input: $H \times W = 8 \times 8$



filter

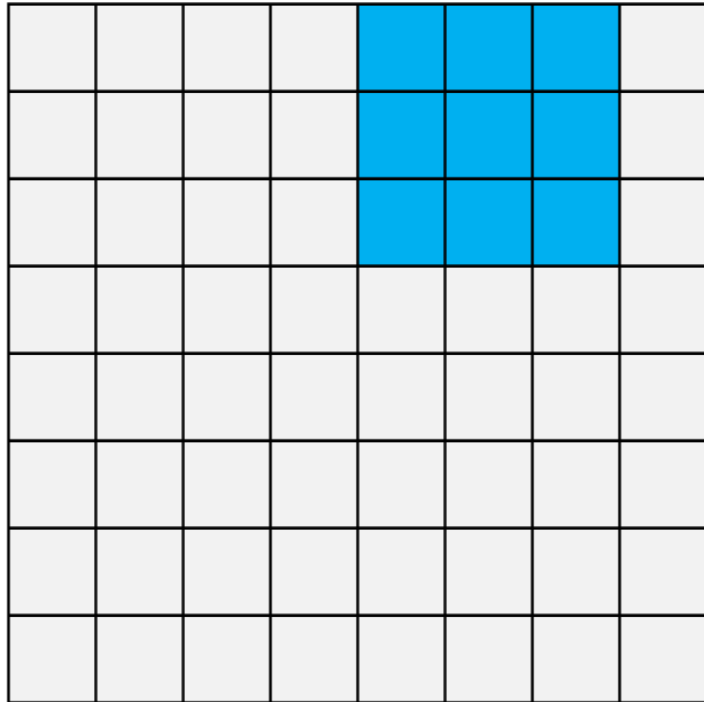


output: $H \times W = 6 \times 6$

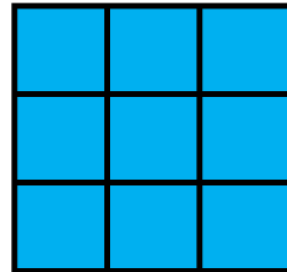


Convolution: padding

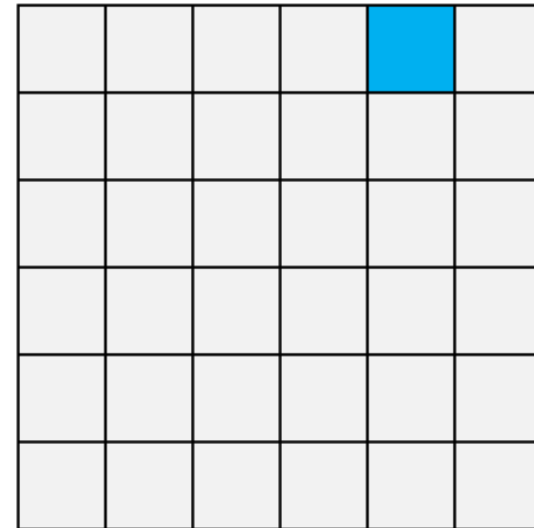
input: $H \times W = 8 \times 8$



filter

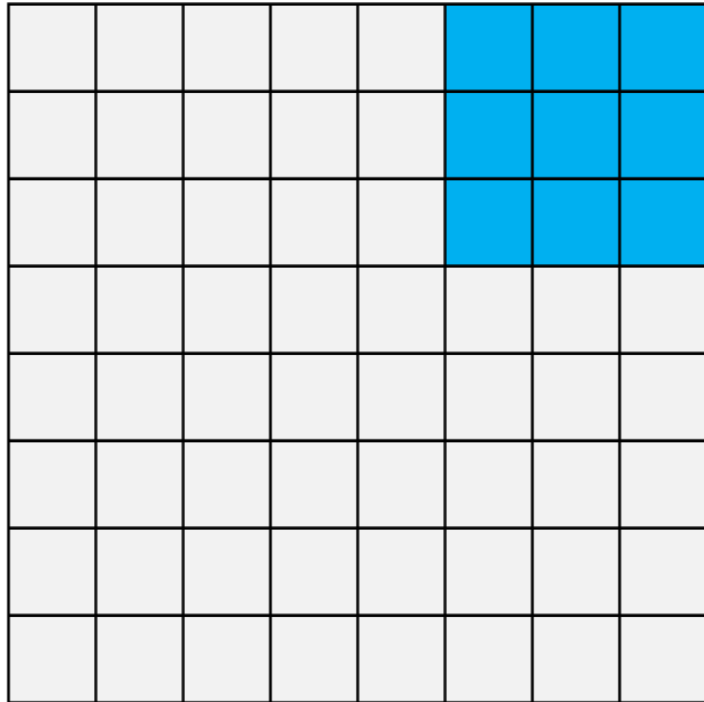


output: $H \times W = 6 \times 6$

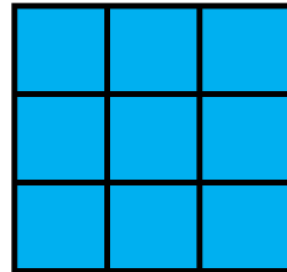


Convolution: padding

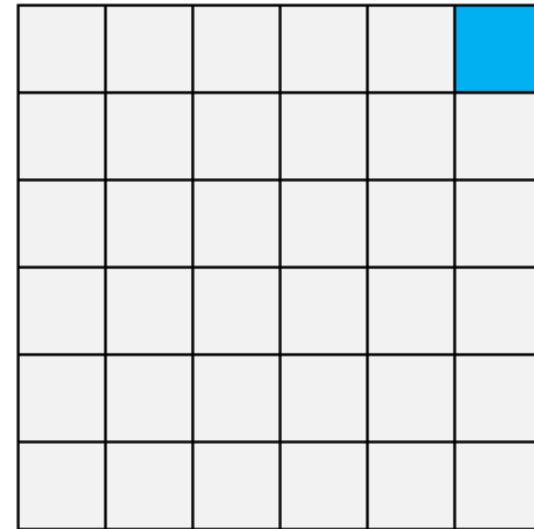
input: $H \times W = 8 \times 8$



filter

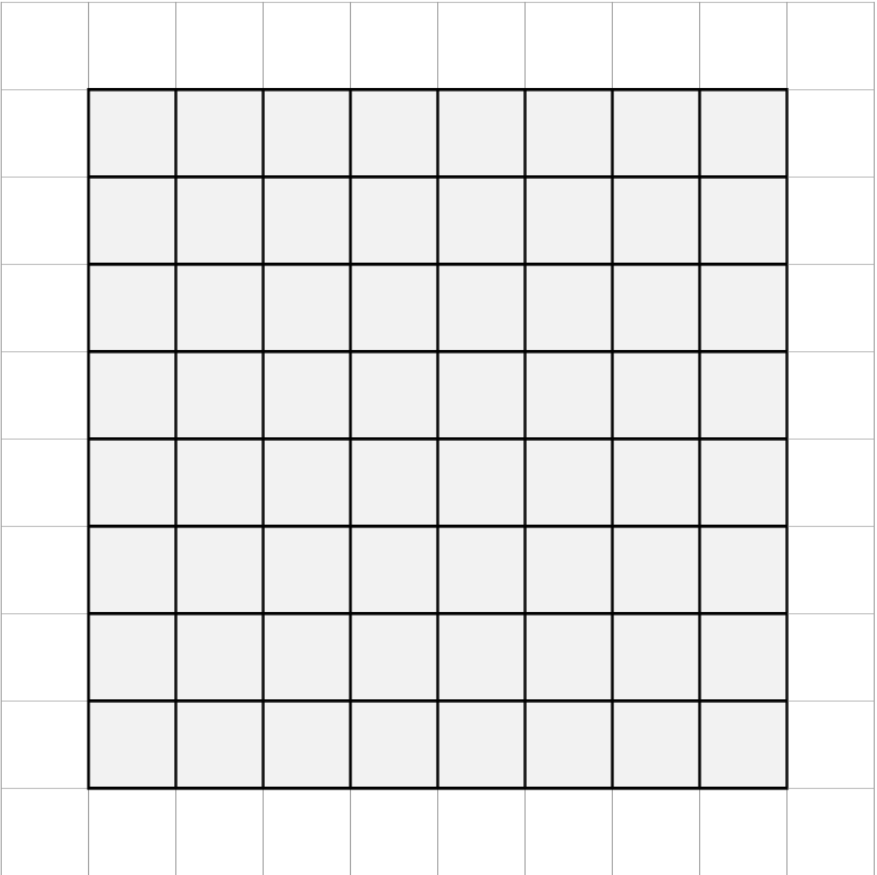


output: $H \times W = 6 \times 6$



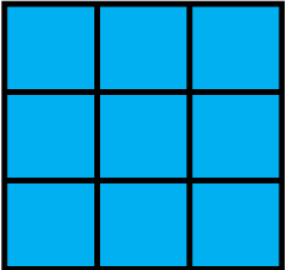
Convolution: stride

input

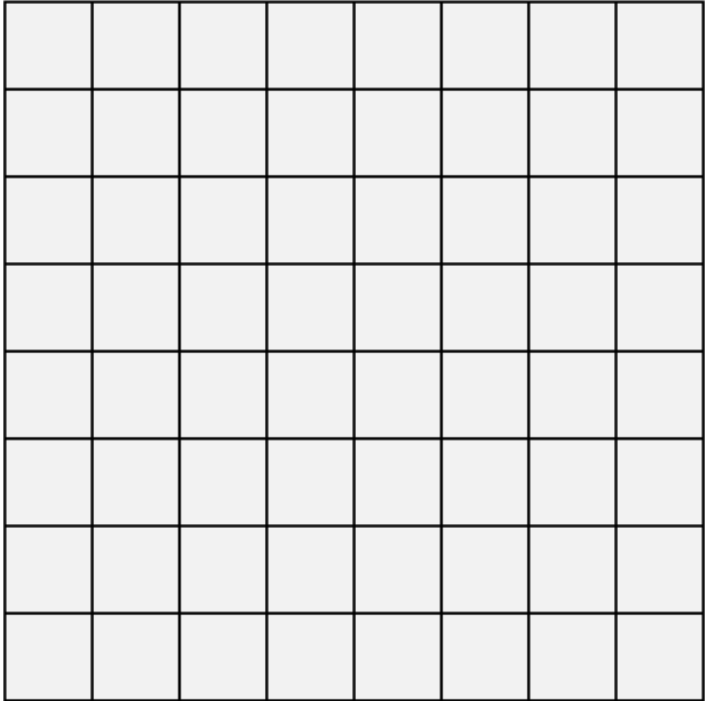


stride = 2

filter

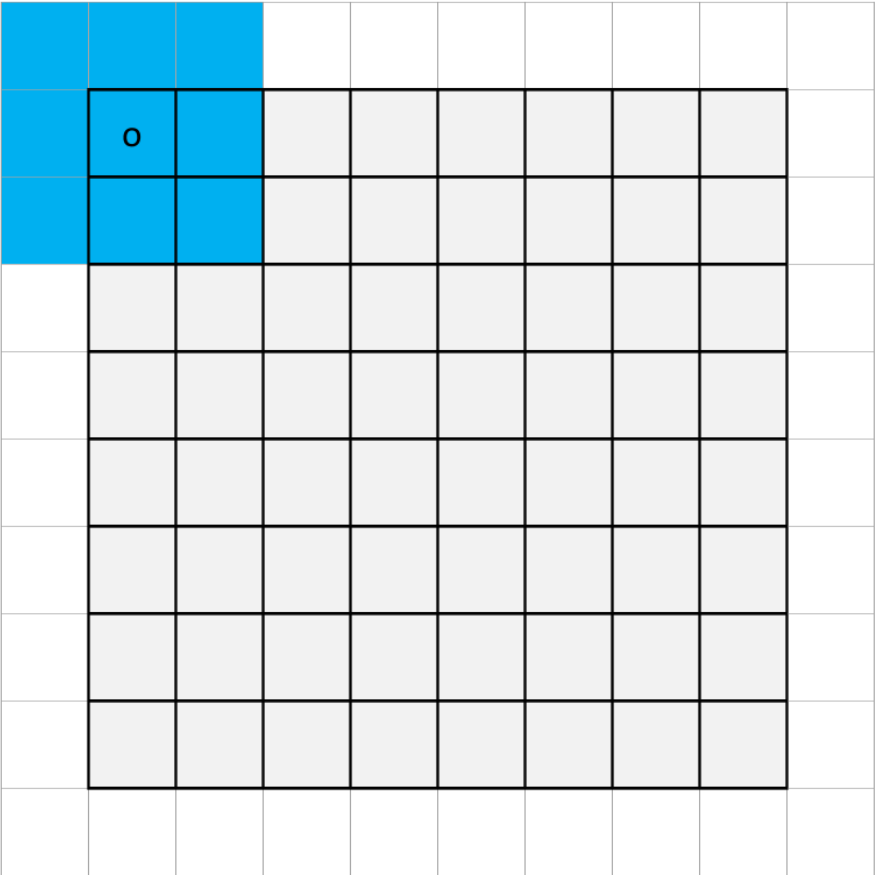


output



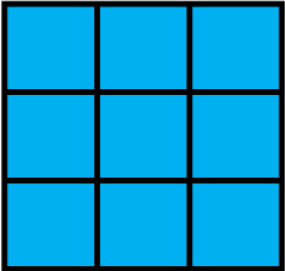
Convolution: stride

input

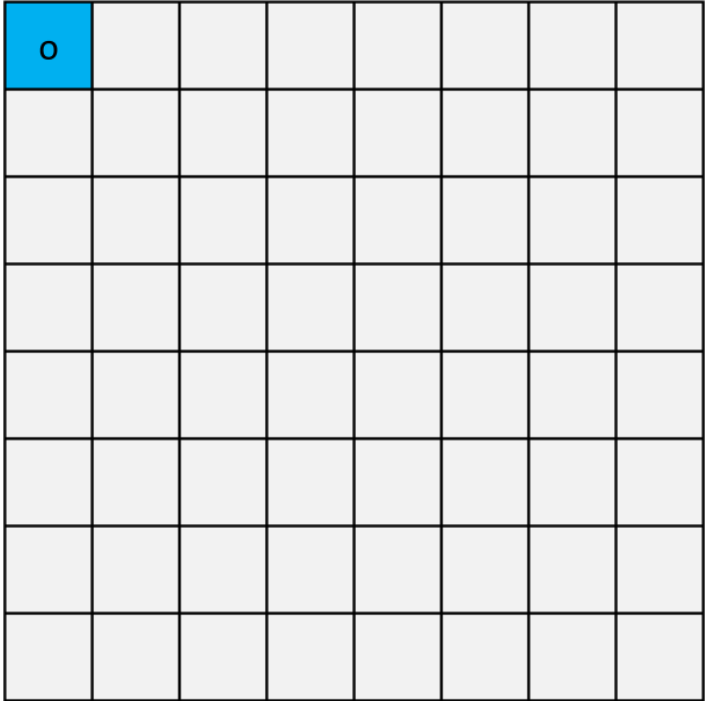


stride = 2

filter

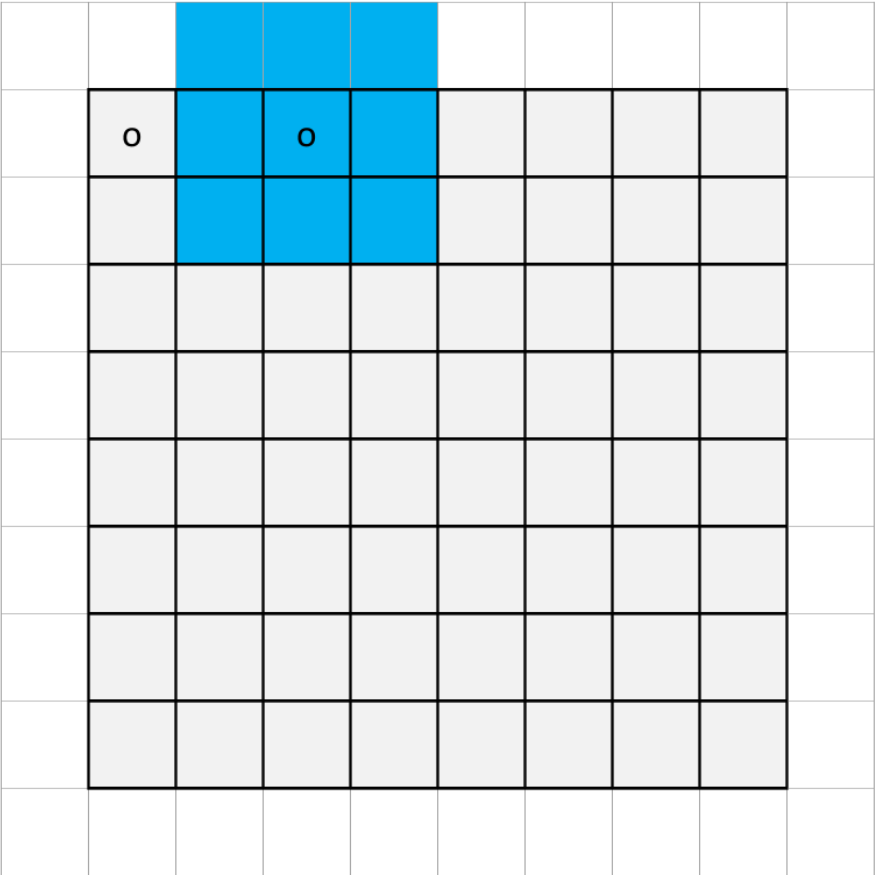


output



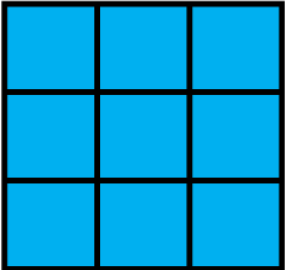
Convolution: stride

input

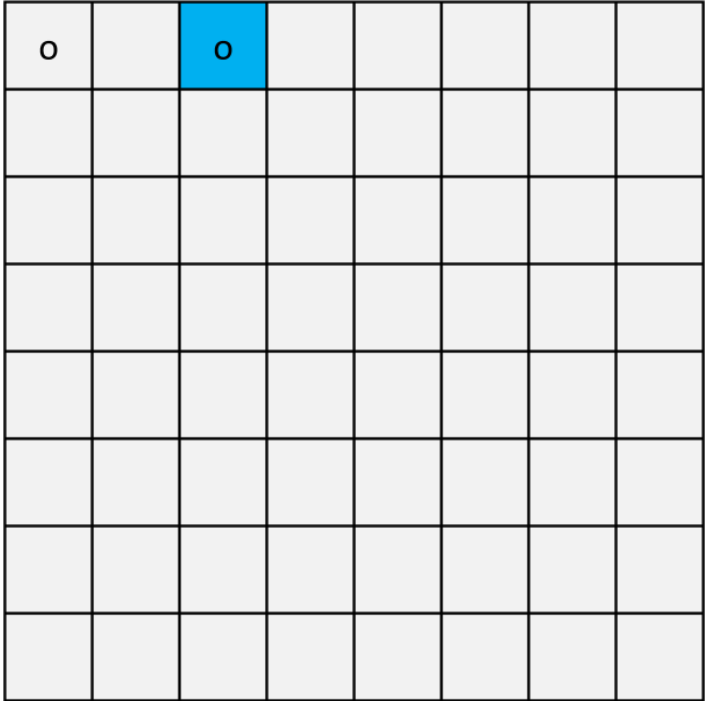


stride = 2

filter

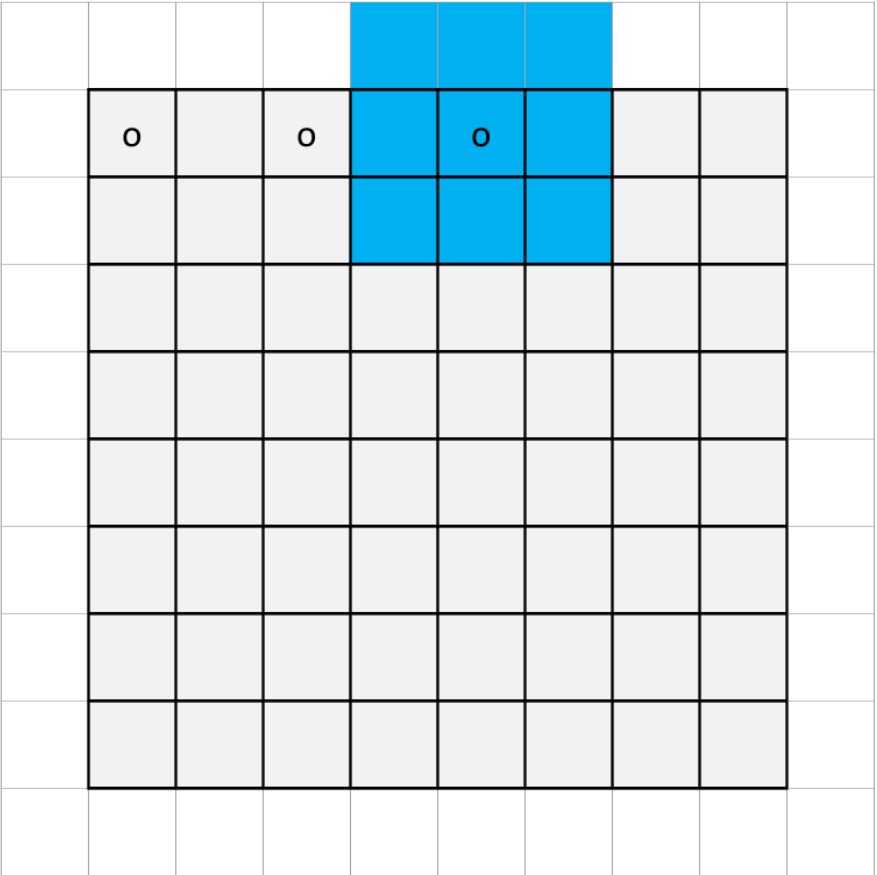


output



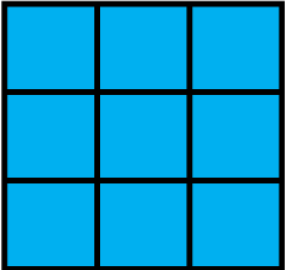
Convolution: stride

input

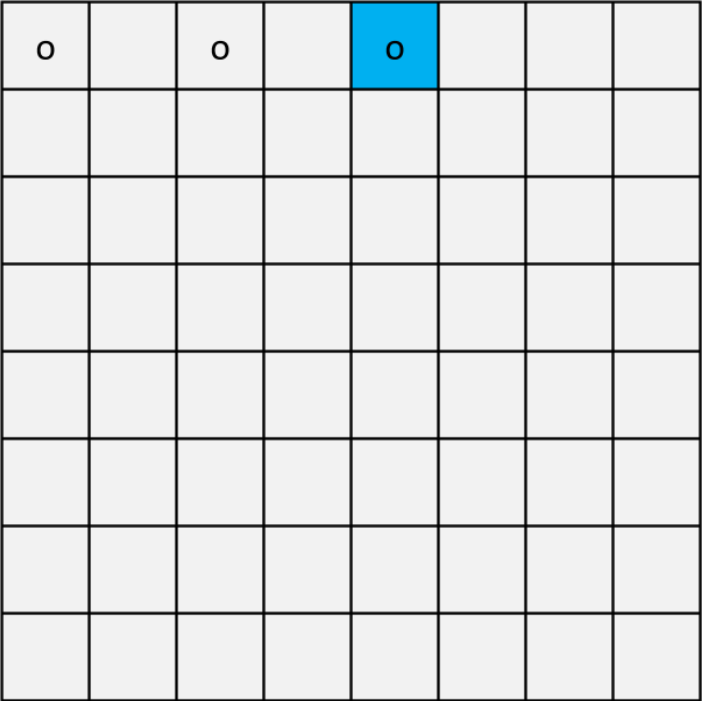


stride = 2

filter



output



Convolution: stride

input

	o		o		o		o		
	o		o		o		o		
	o		o		o		o		
	o		o		o		o		

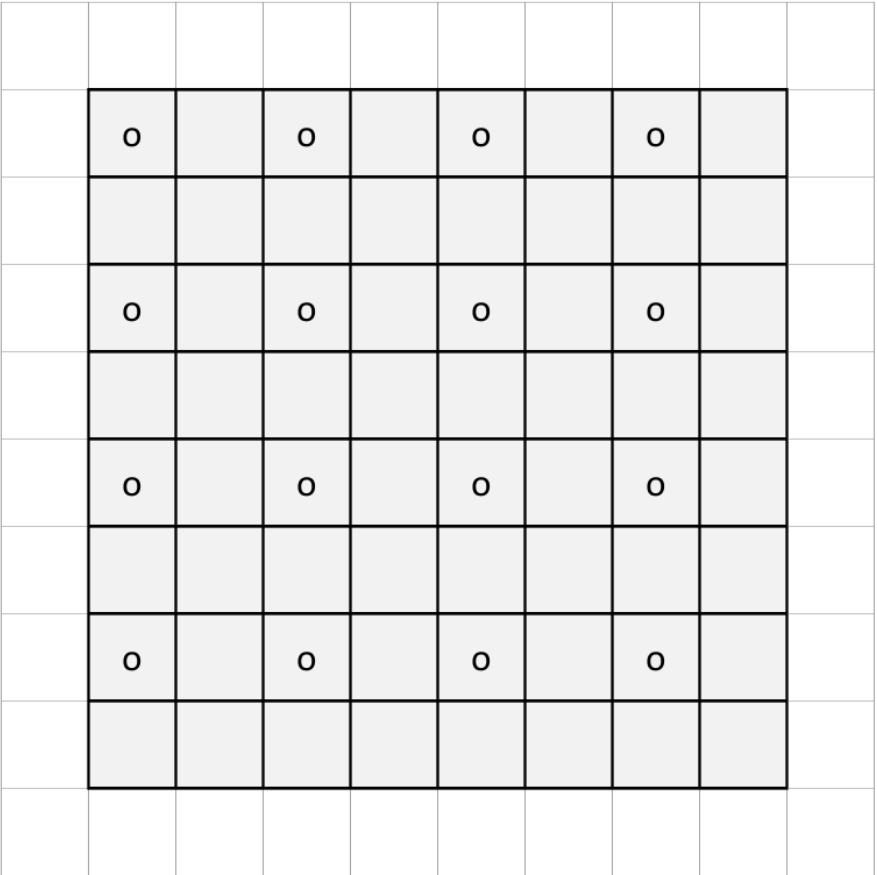
filter

output

o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o		o	

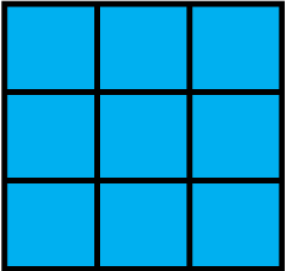
Convolution: stride

input: $H \times W = 8 \times 8$

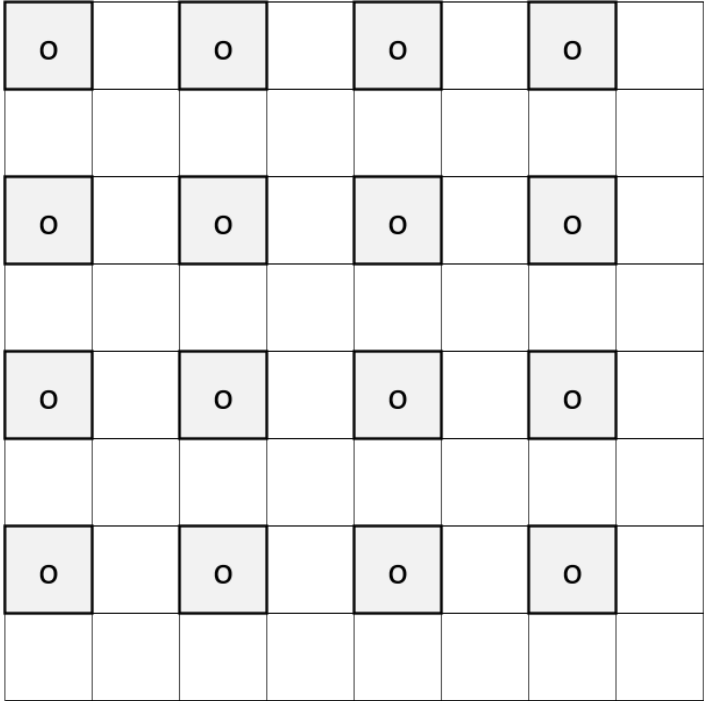


stride = 2

filter

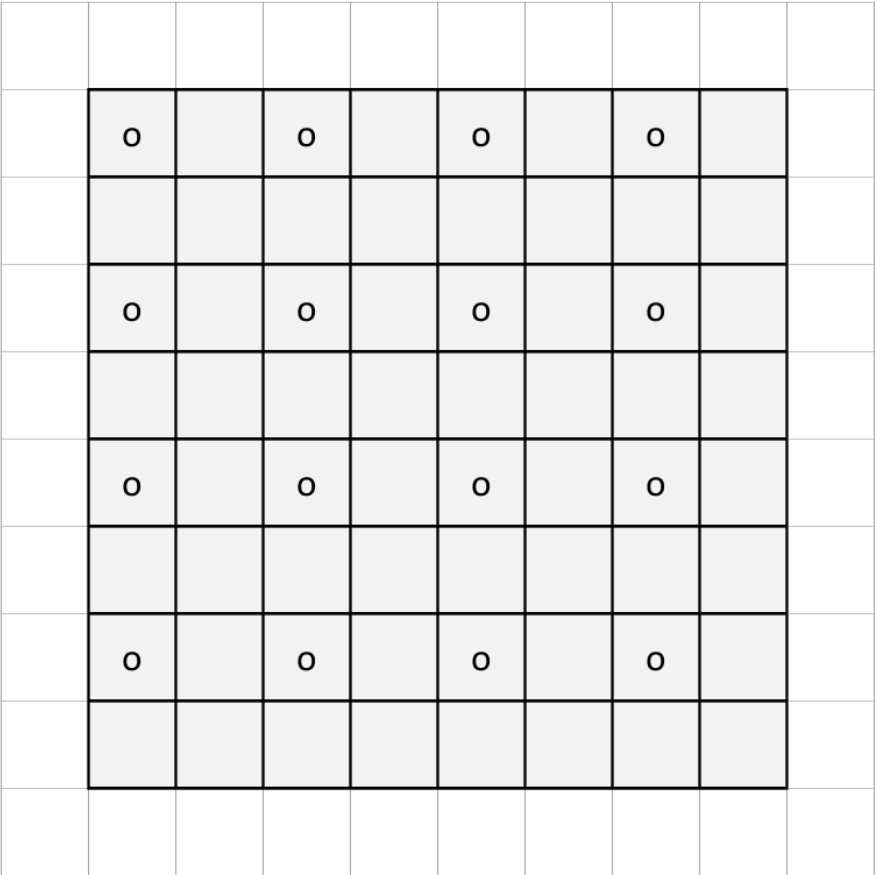


output: $H \times W = 4 \times 4$



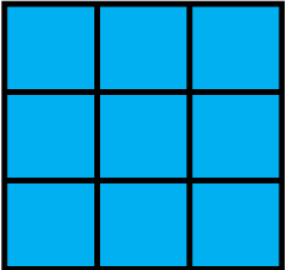
Convolution: stride

input: $H \times W = 8 \times 8$

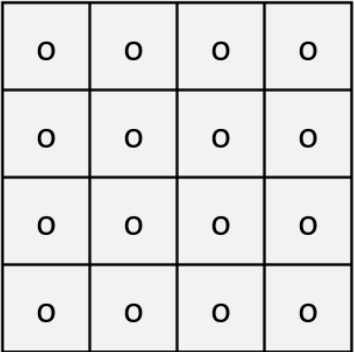


stride = 2

filter



output: $H \times W = 4 \times 4$

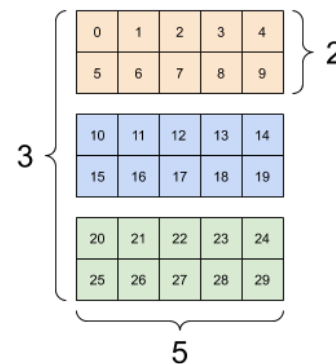
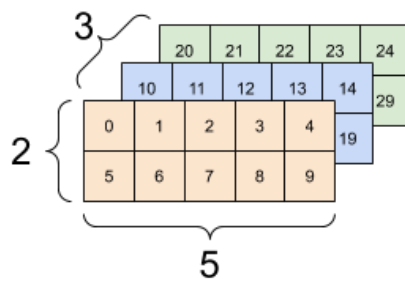
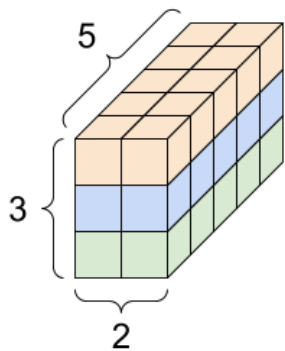
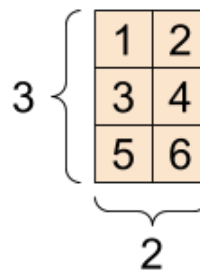
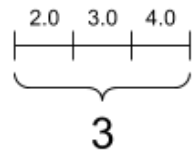


Outline

- Recap (fully-connected net)
 - Motivation and big picture ideas of CNN
 - Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing
- 3d tensors
- Max pooling
 - Larger window
 - Typical architecture and summary

A tender intro to tensor:

4





blue



green



red

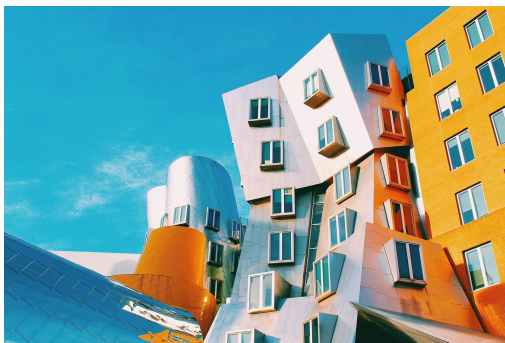
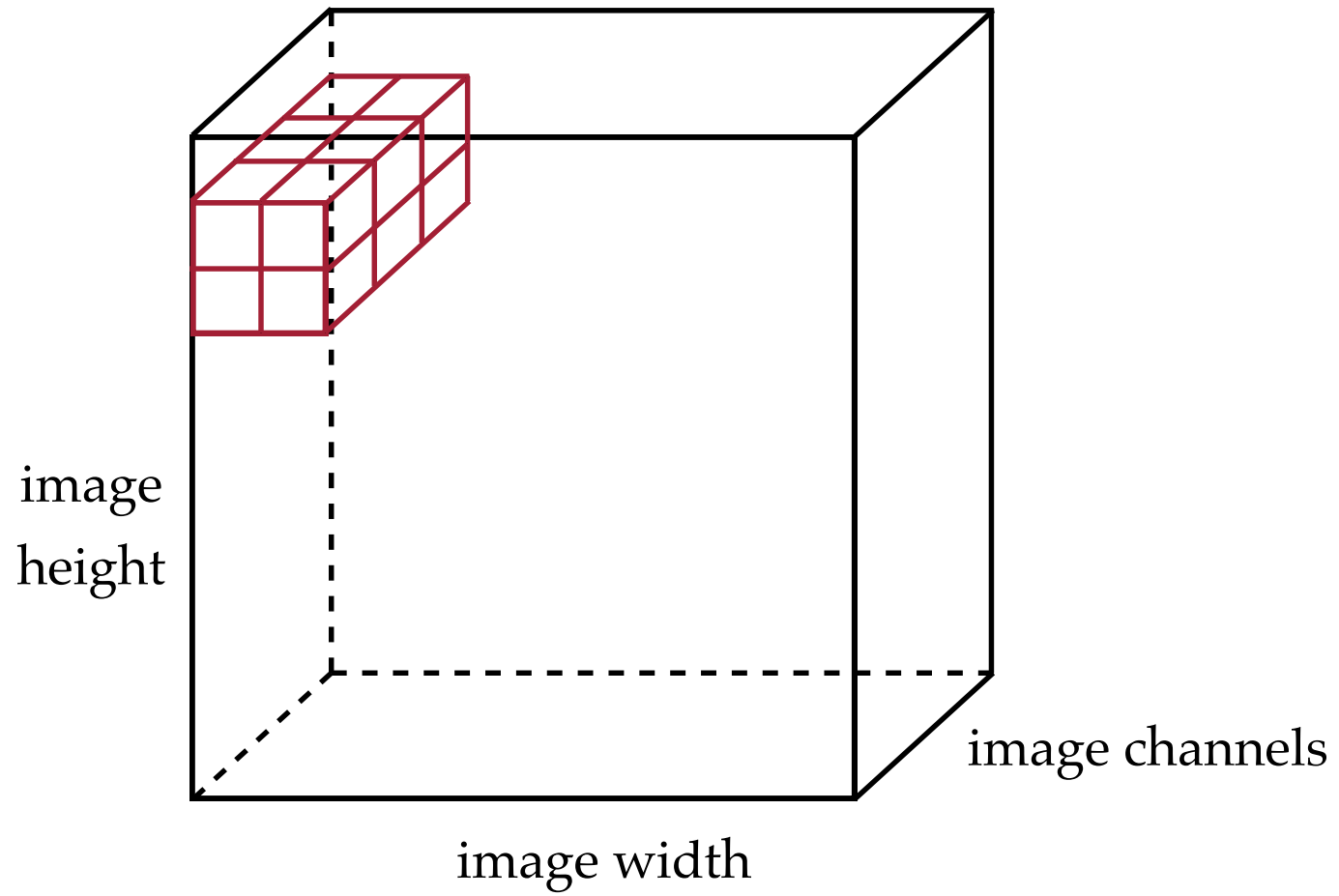


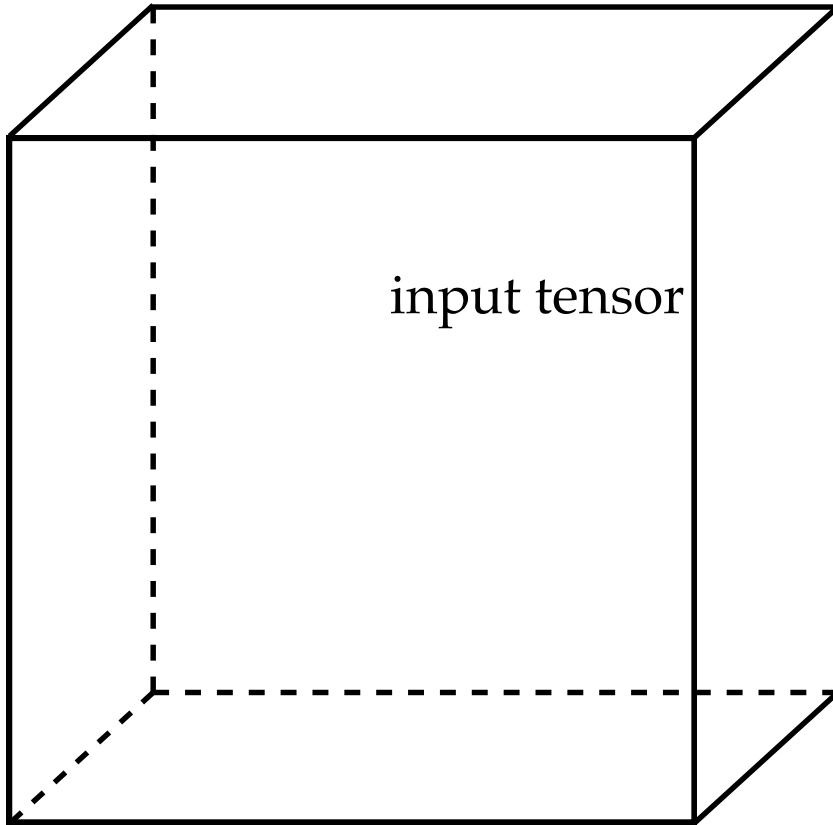
image
height



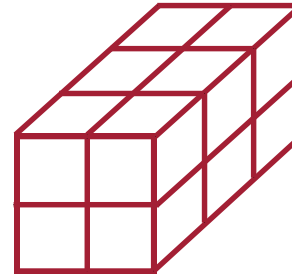
image channels

image width

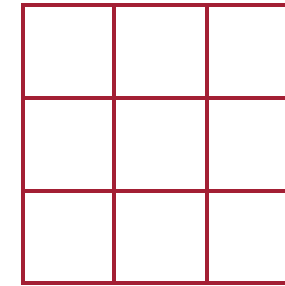




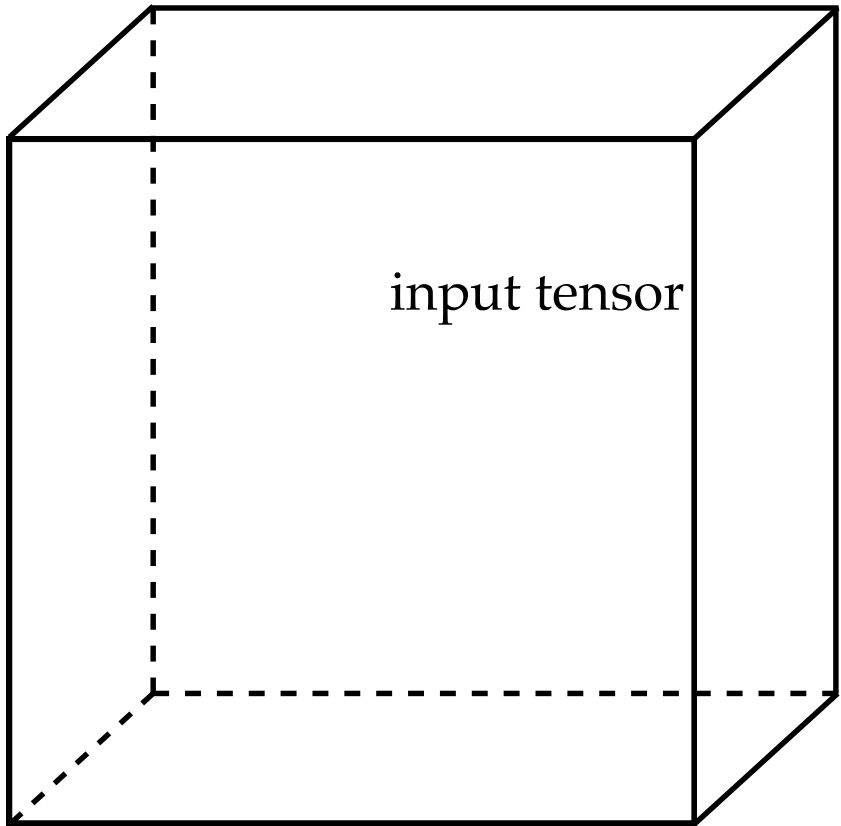
filter



output

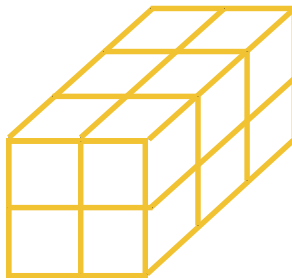
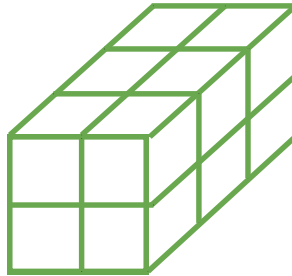
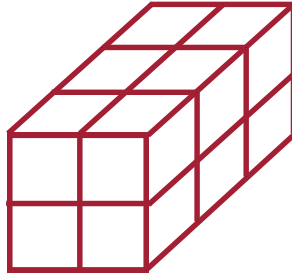


- 3d tensor input, depth d
- 3d tensor filter, depth d
- 2d tensor (matrix) output

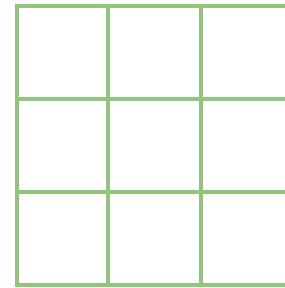
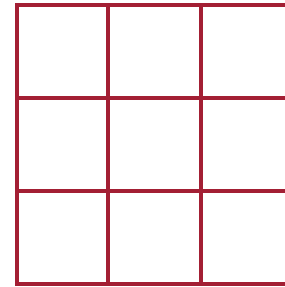


input tensor

filters

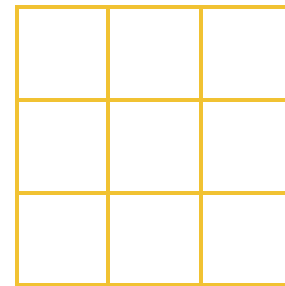


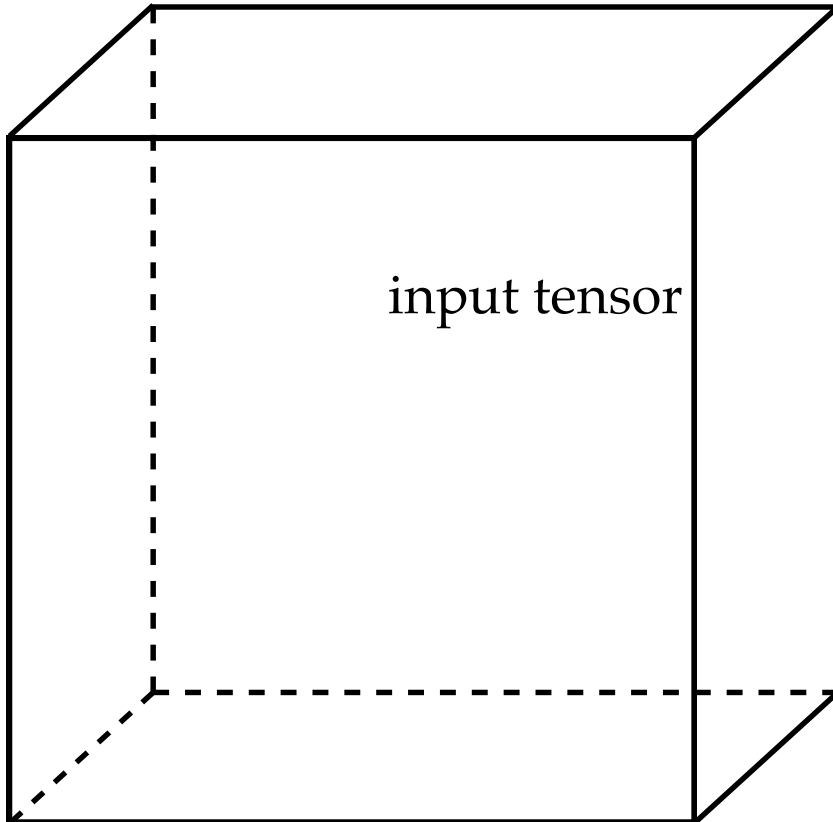
outputs



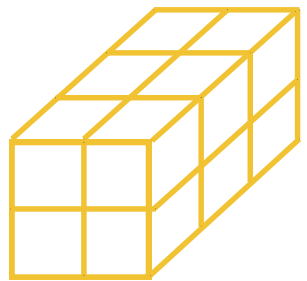
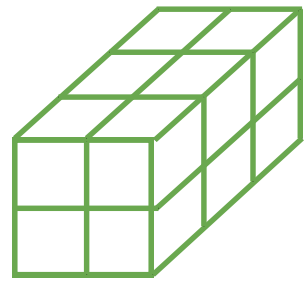
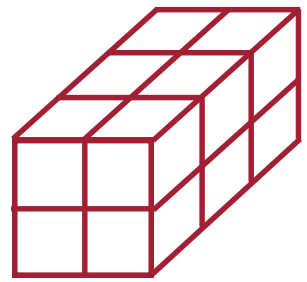
...

...



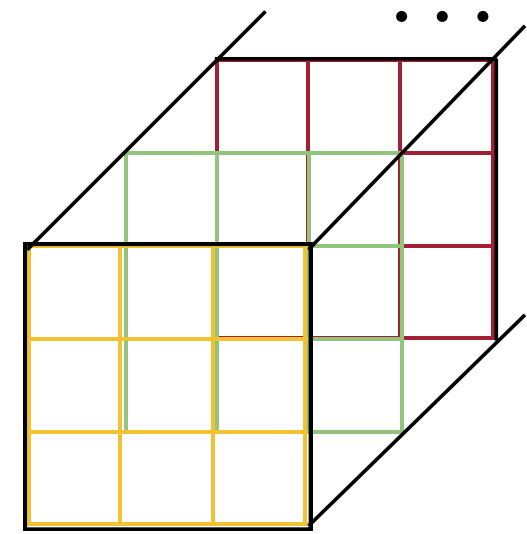


filters

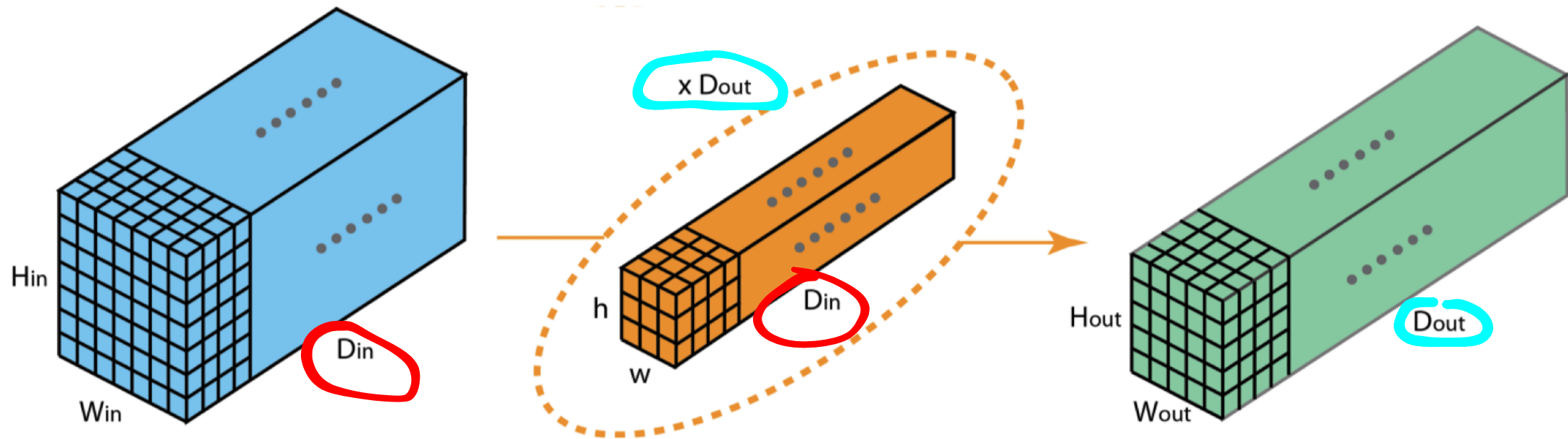


...

output tensor



- 3d tensor input, depth d
- k 3d filters:
 - each filter of depth d
 - each filter makes a 2d tensor (matrix) output
- total output 3d tensor, depth k



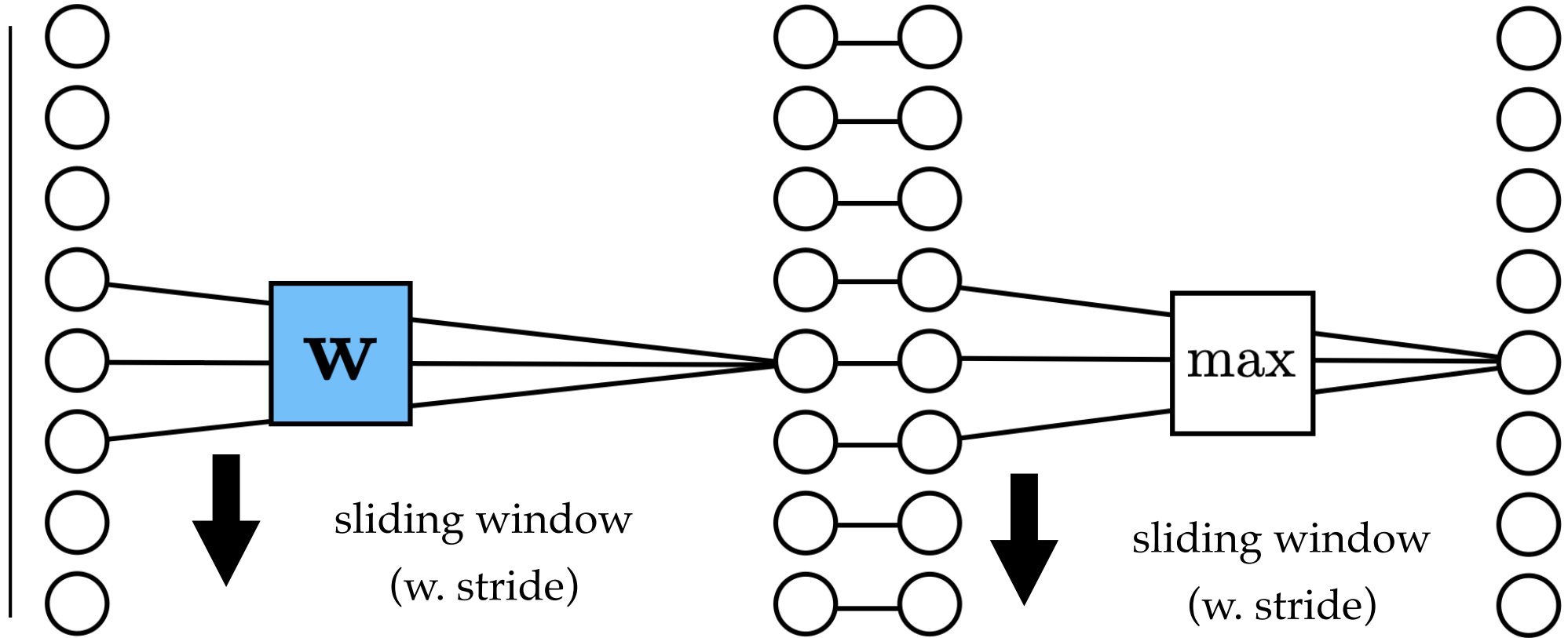
[image credit: medium]

Outline

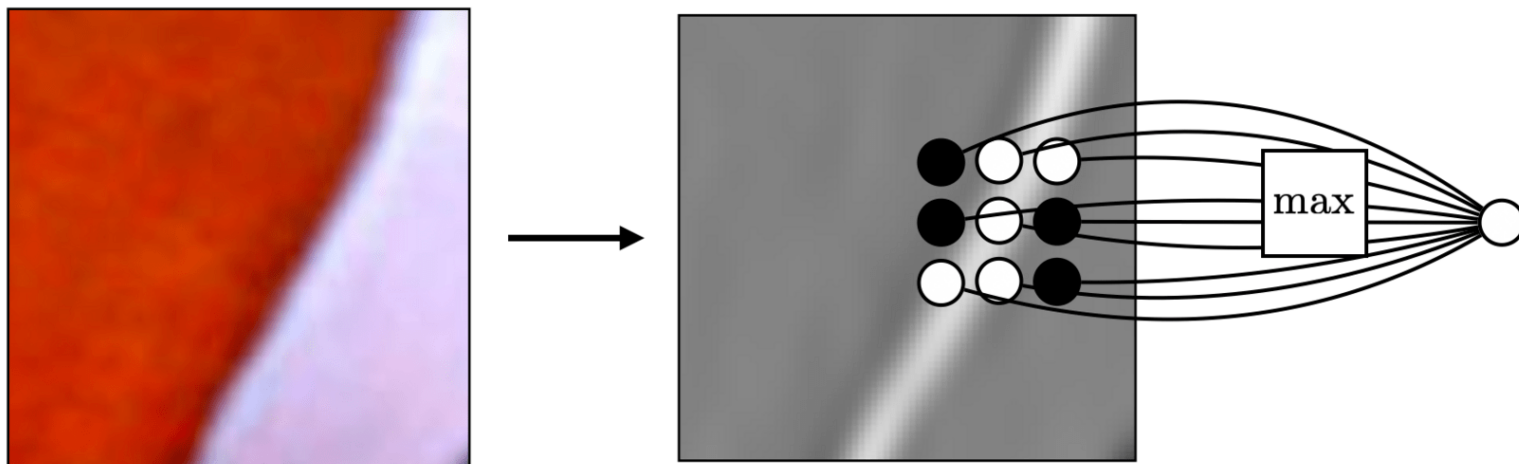
- Recap (fully-connected net)
- Motivation and big picture ideas of CNN
- Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing
 - 3d tensors
- Max pooling
 - Larger window
- Typical architecture and summary

convolution

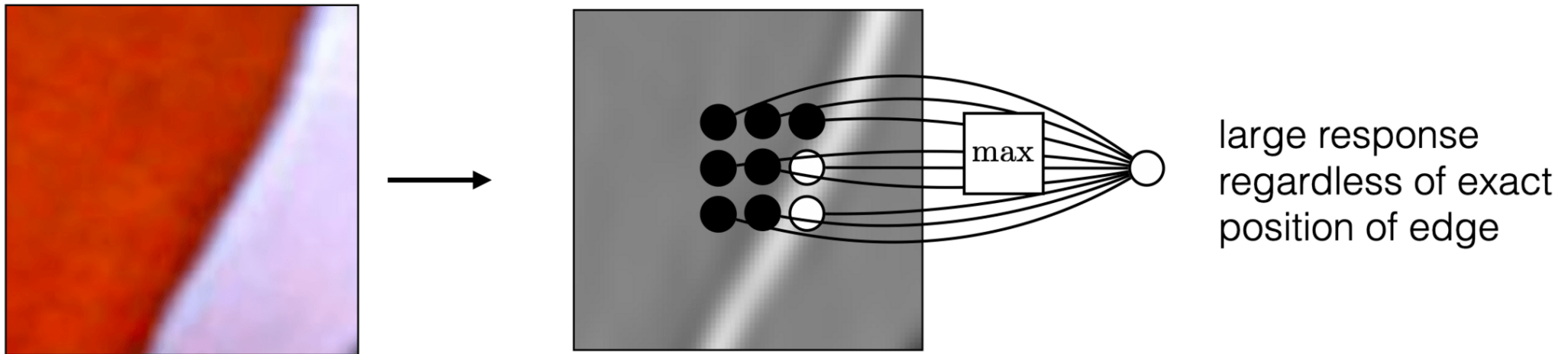
max pooling



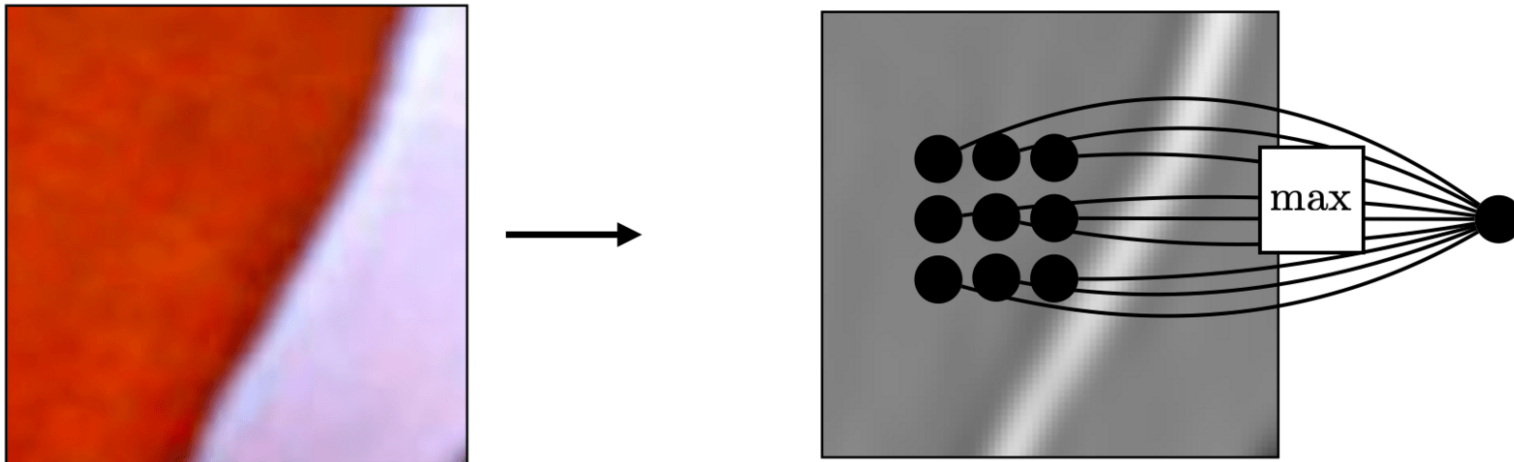
Pooling across spatial locations achieves stability w.r.t. small translations:



Pooling across spatial locations achieves stability w.r.t. small translations:

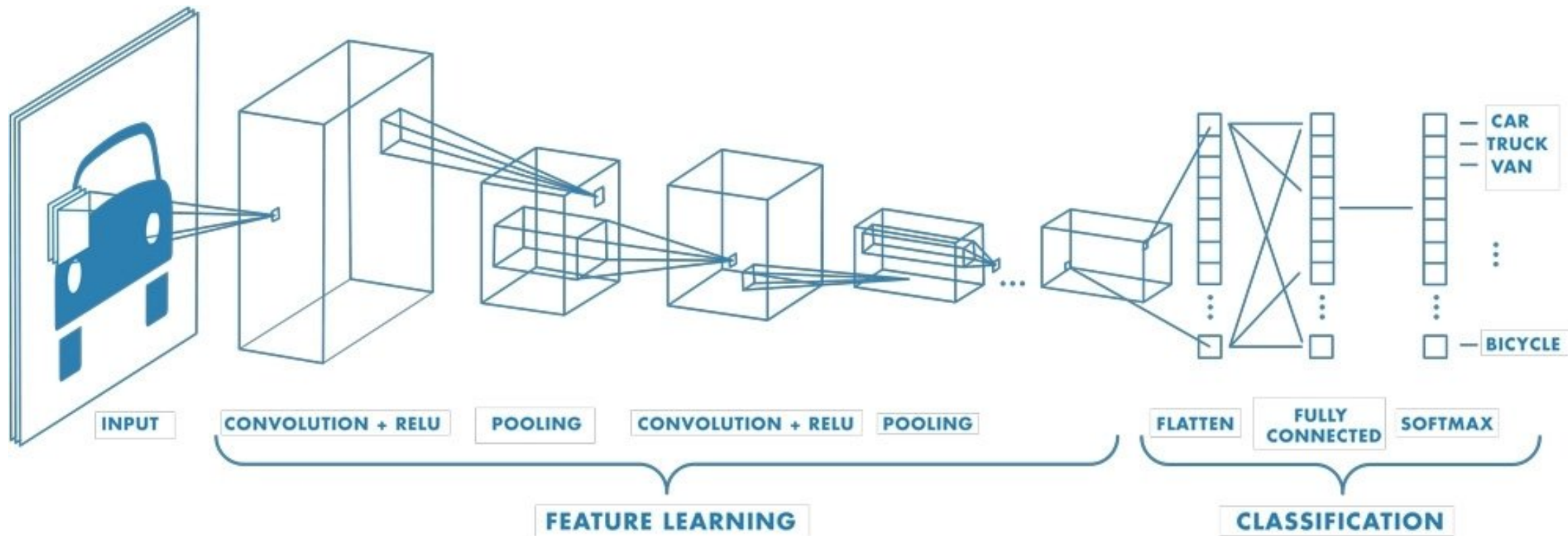


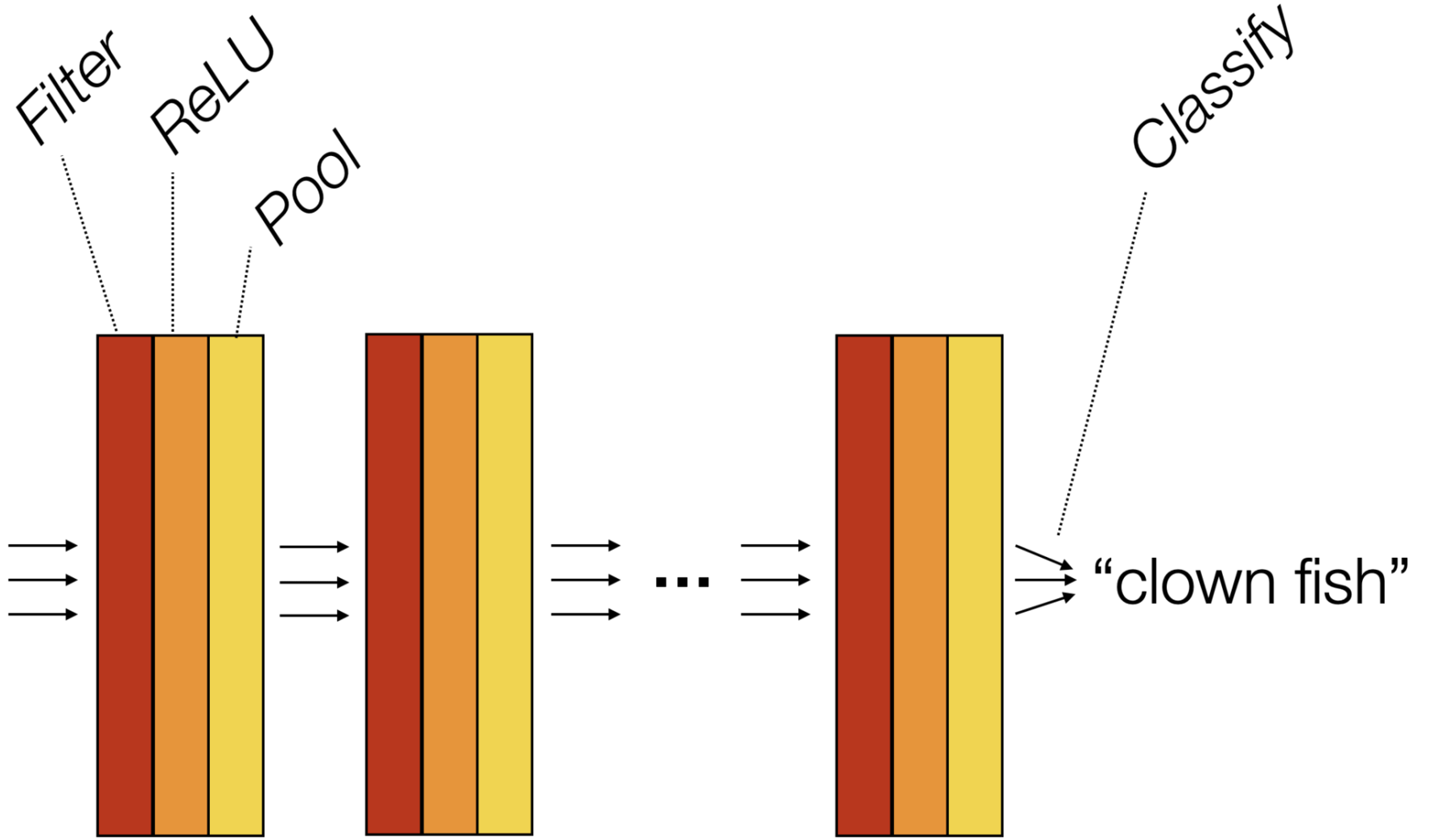
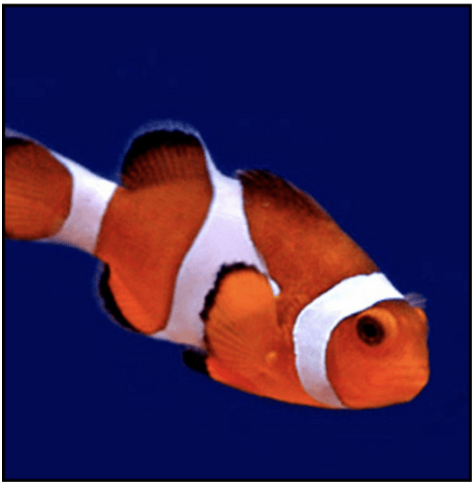
Pooling across spatial locations achieves stability w.r.t. small translations:



Outline

- Recap (fully-connected net)
- Motivation and big picture ideas of CNN
- Convolution operation
 - 1d and 2d convolution mechanics
 - interpretation:
 - local connectivity
 - weight sharing
 - 3d tensors
- Max pooling
 - Larger window
- Typical architecture and summary





We'd love it for you to share some lecture [feedback](#).

Thanks!