



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 5212
Machine Learning
Lecture 19

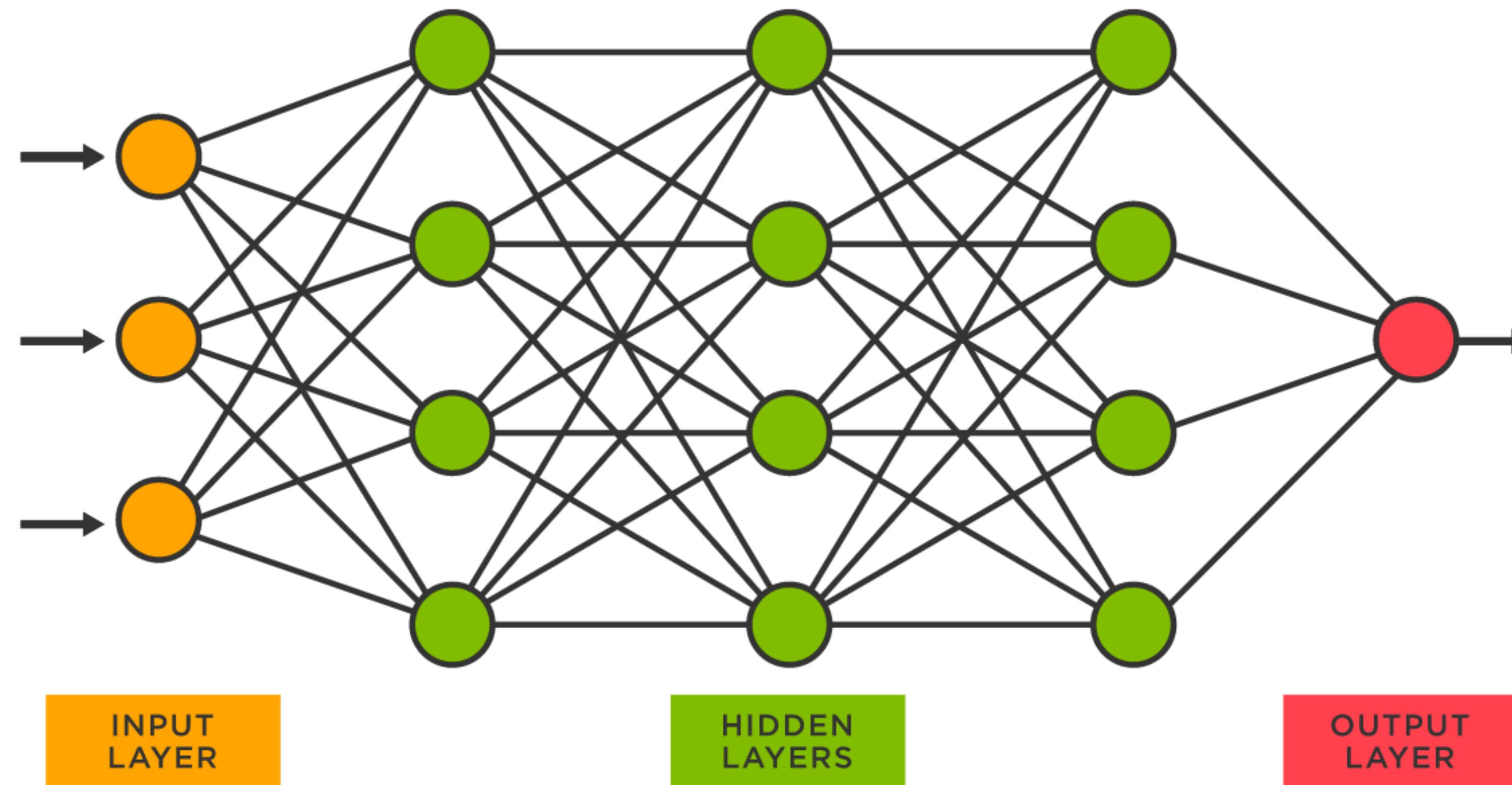
Neural Networks, Architectures

Junxian He
Nov 12, 2024

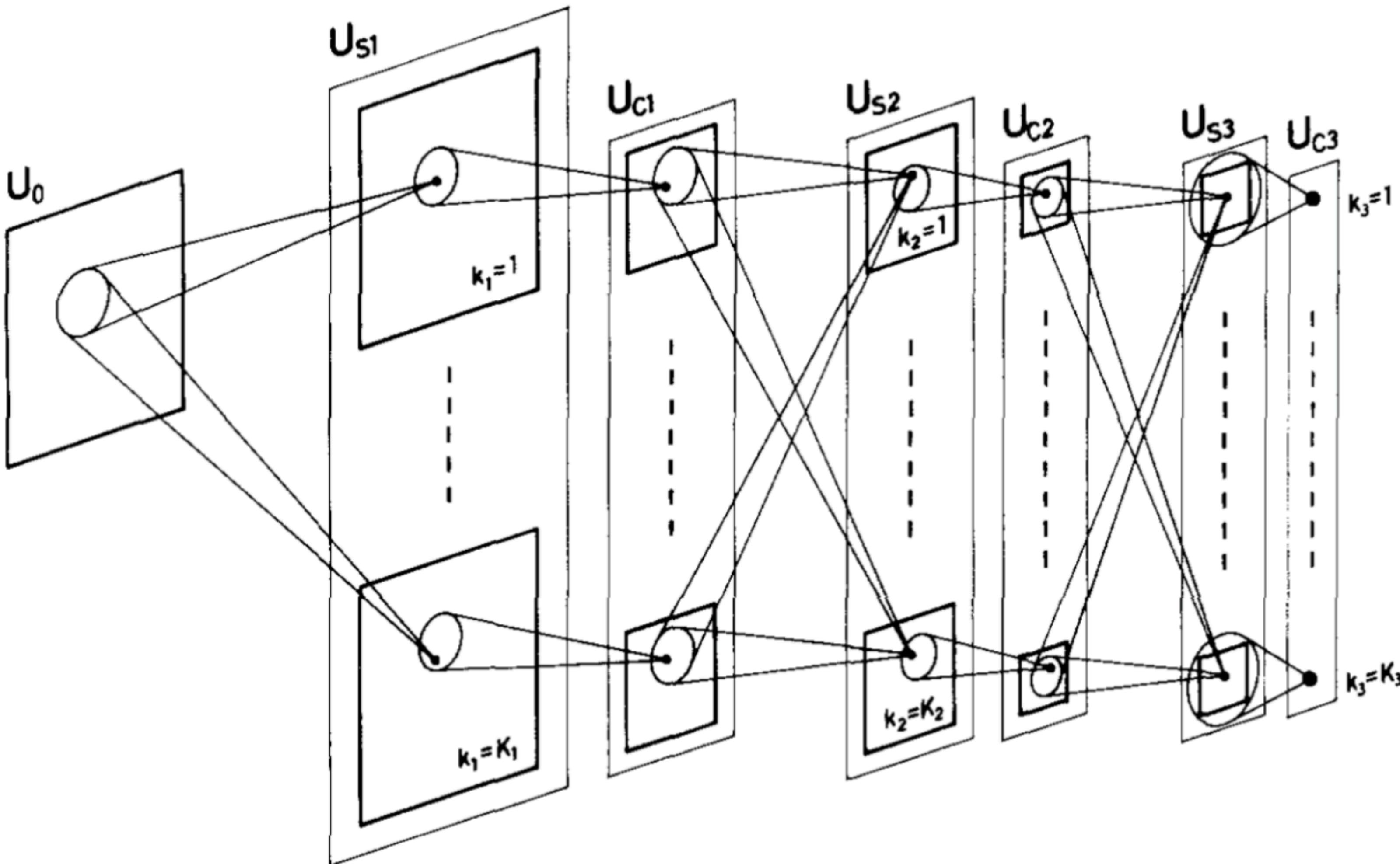
Announcements

1. Midterm Exam paper check session is at this evening (7-8pm), 3520
2. Programming Assignment is out, please start early
3. HW3 is out
4. HW4 will be easy (multi-choice QA only)

Recap: Multilayer Perceptron Neural Networks (MLP)

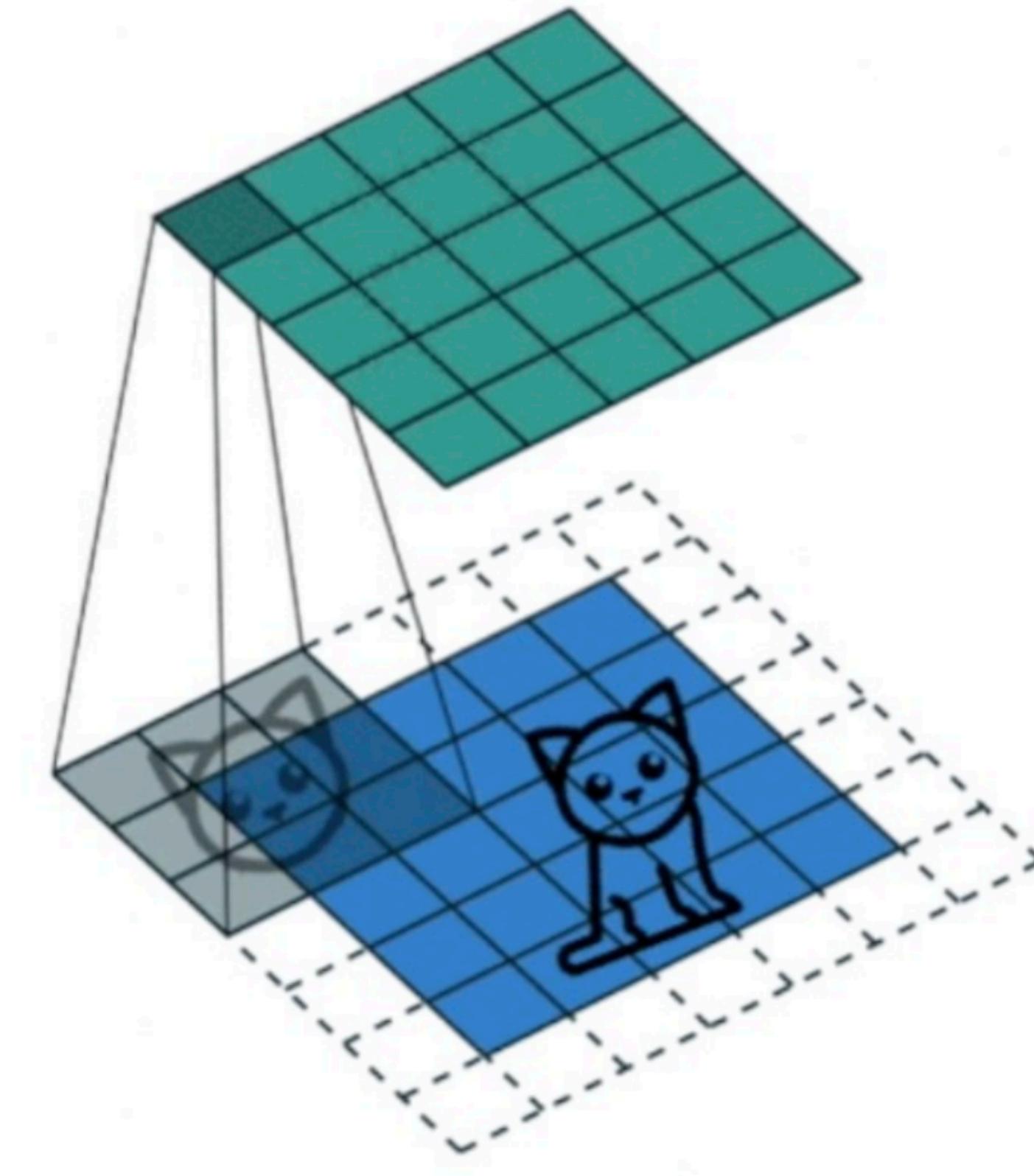


Convolutional Neural Networks



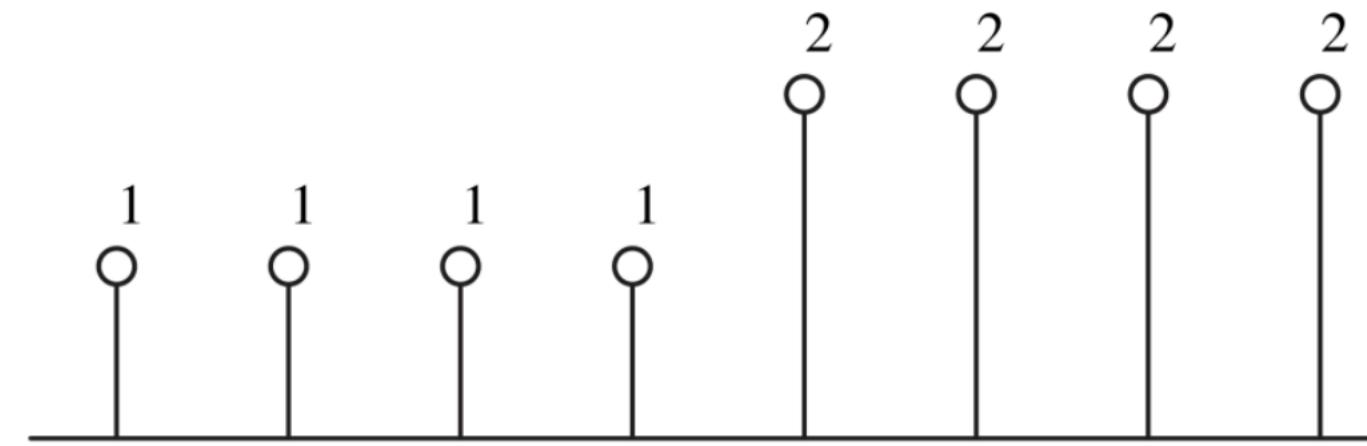
Convolution is template matching

- with a sliding window
- abstract templates
- similarity measured by dot product
- stronger activation, better matching

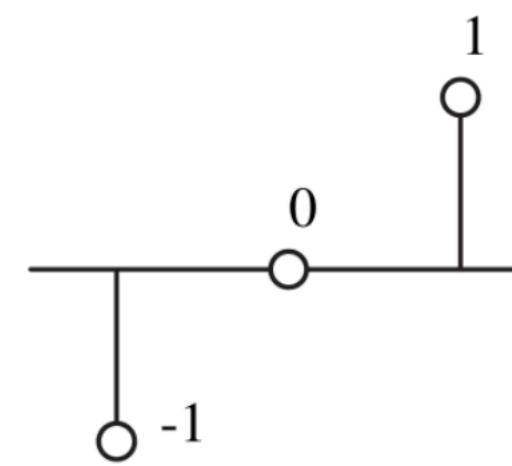


Convolution: a 1-D example

input



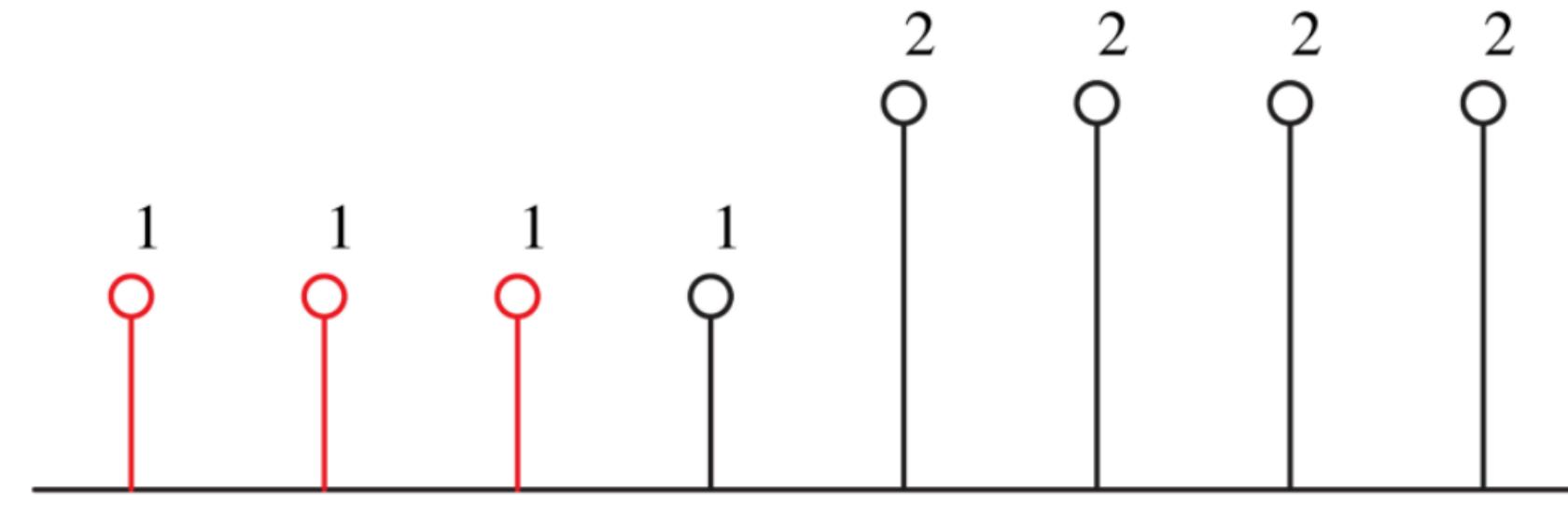
filter



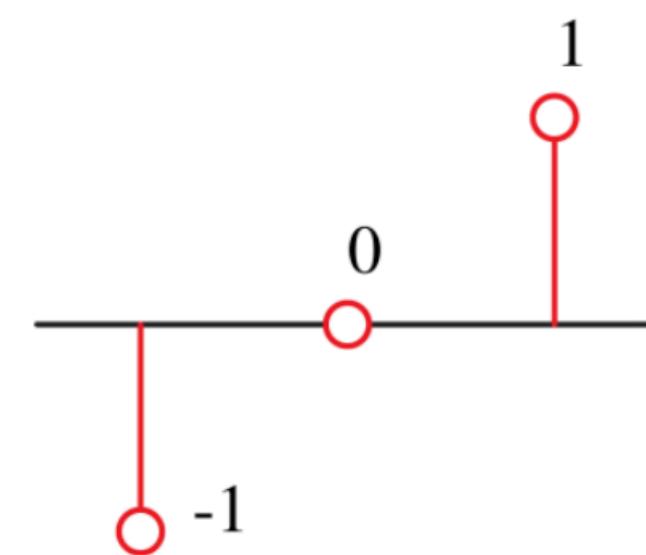
Convolution: a 1-D example

- sliding window
- dot product

input



filter



output

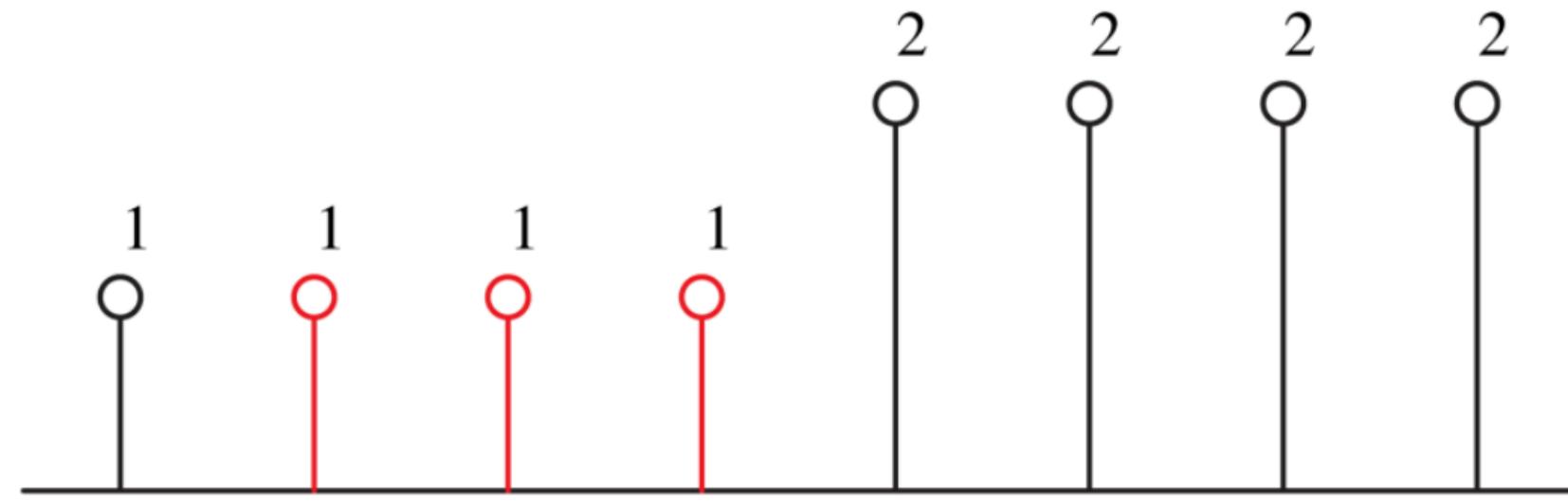
$$-1 \times 1 + 0 \times 1 + 1 \times 1$$



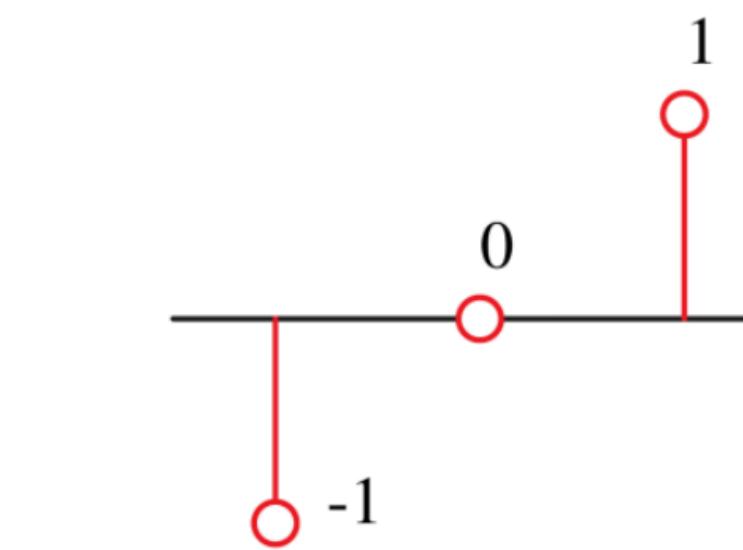
Convolution: a 1-D example

- sliding window
- dot product

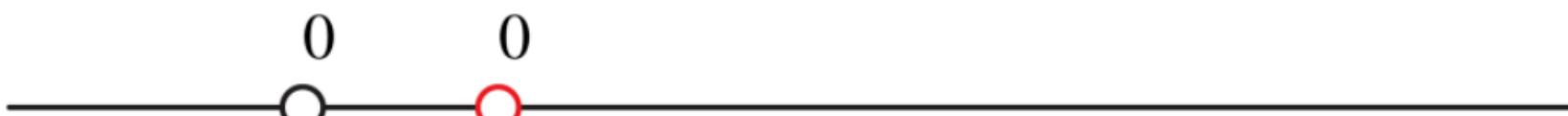
input



filter



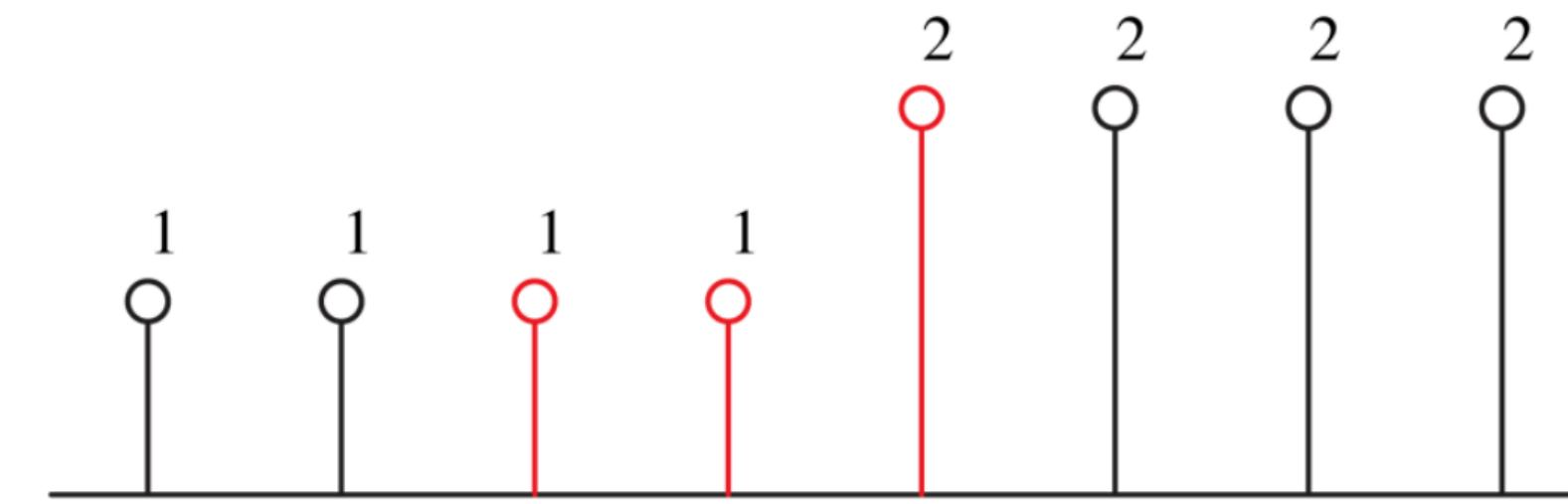
output



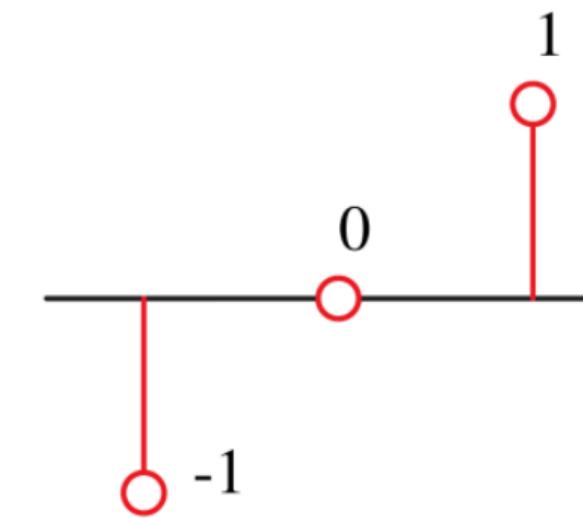
Convolution: a 1-D example

- sliding window
- dot product

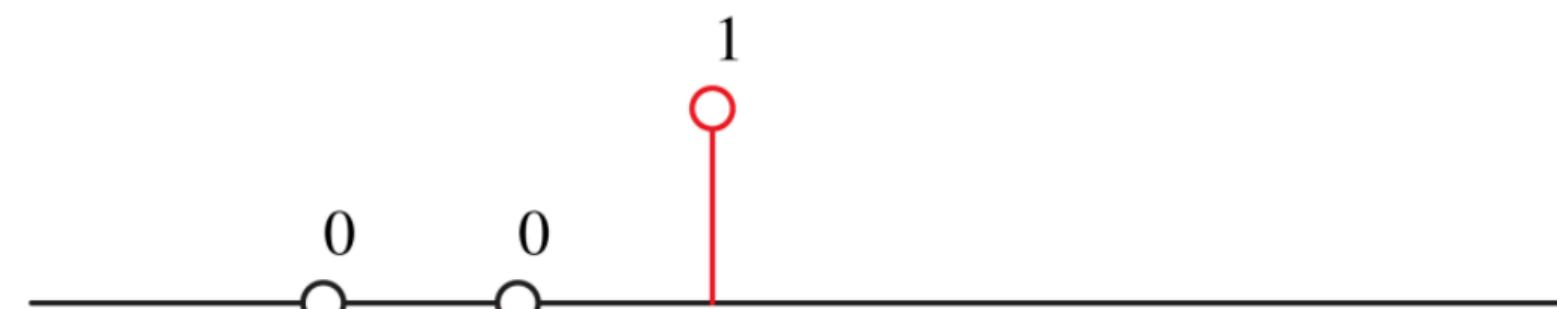
input



filter



output



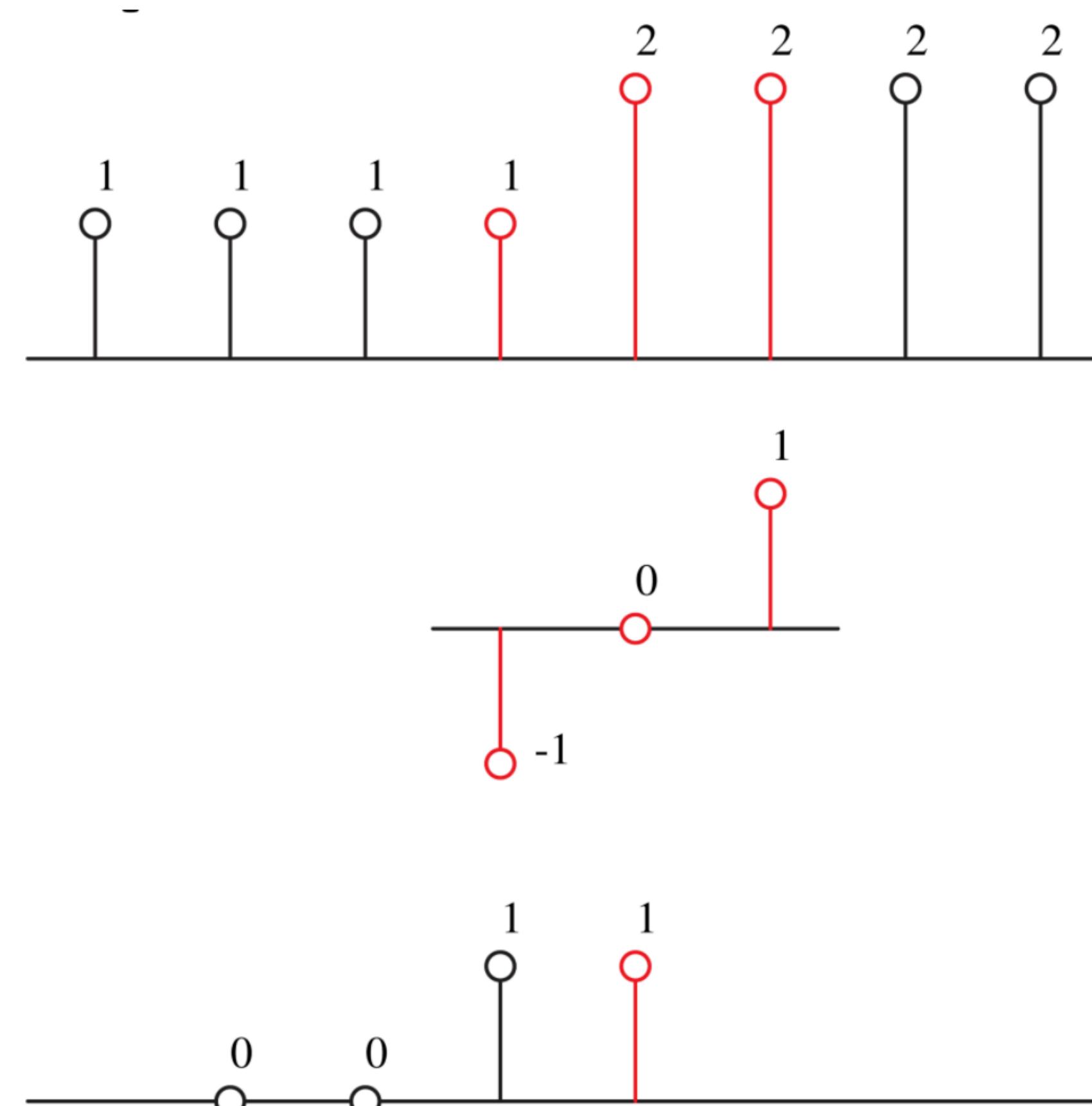
Convolution: a 1-D example

- sliding window
- dot product

input

filter

output



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	0	0
0	0	1	1	1	0	0	0
0	0	0	0	0	0	0	0

output

filter

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

Convolution: a 2-D example

input

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

output

-3					

filter

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

- sliding window
- dot product

Convolution: a 2-D example

input

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

output

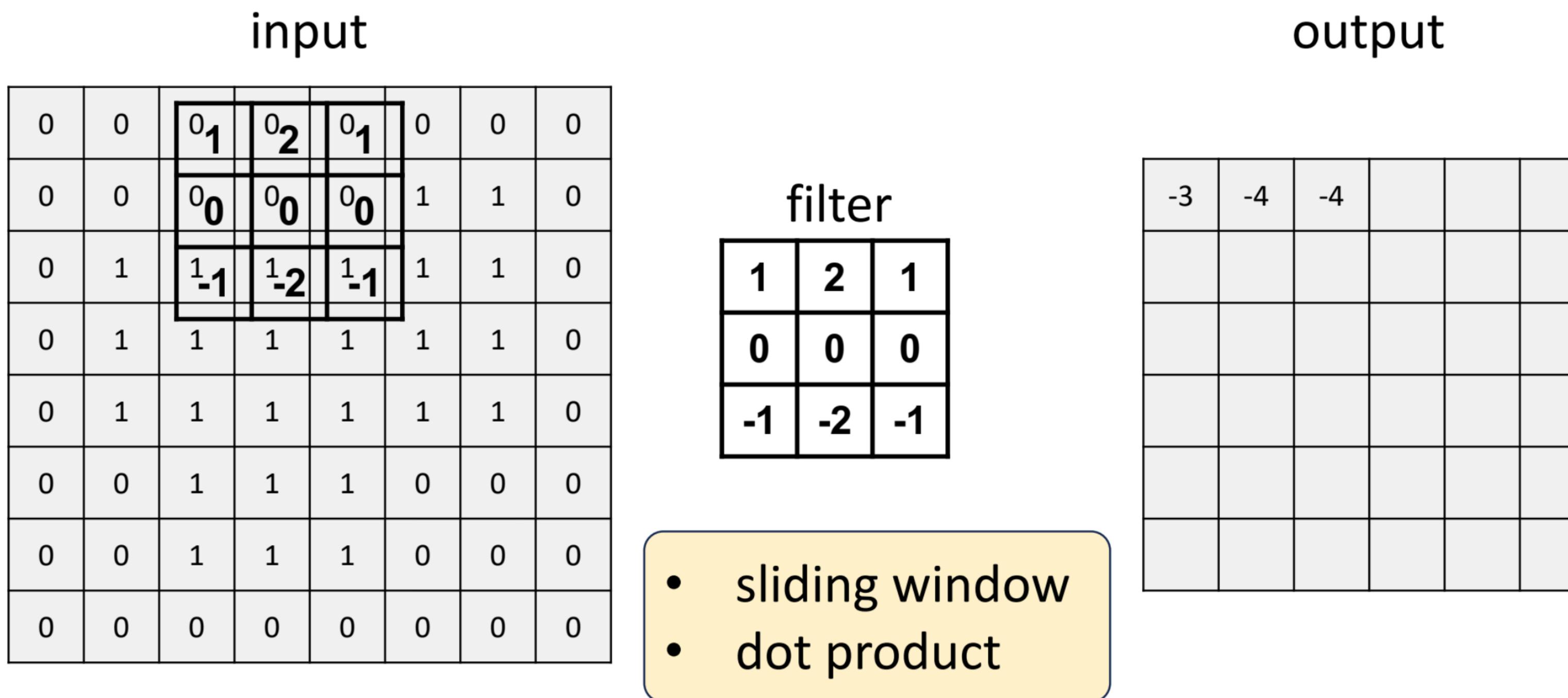
-3	-4					

filter

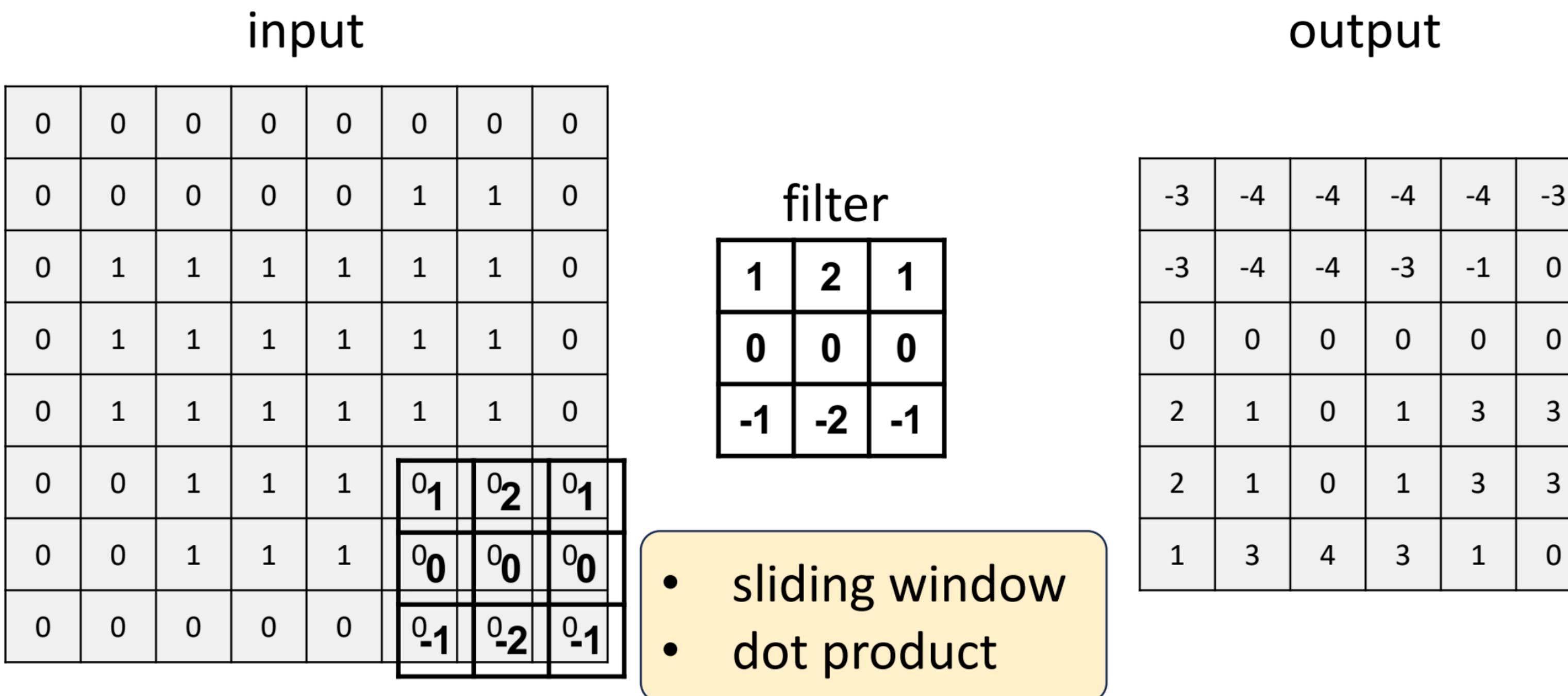
1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

Convolution: a 2-D example



Convolution: a 2-D example



Convolution: a 2-D example

$$y[n, m] = \sum_{i=-r}^r \sum_{j=-r}^r w[i, j]x[n + i, m + j]$$

output map

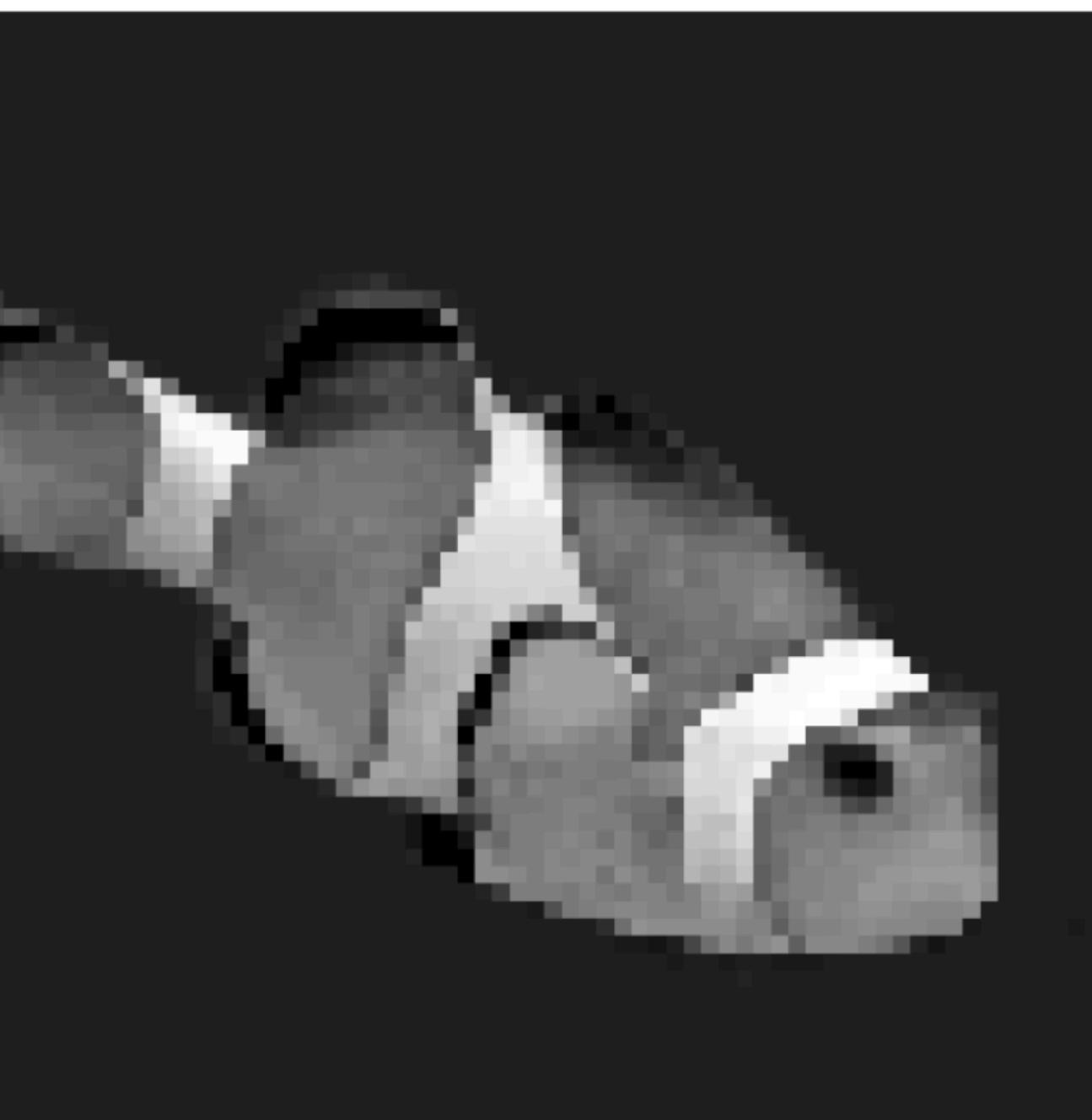
input map

filter weights

coordinates in a local window

r : kernel radius
kernel size = $2r + 1$

Convolution: 2-D



filter

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -2 & -1 \\ \hline \end{array} =$$



Convolution: Multi-channel outputs



$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -2 & -1 \\ \hline \end{array} & = \end{matrix}$$

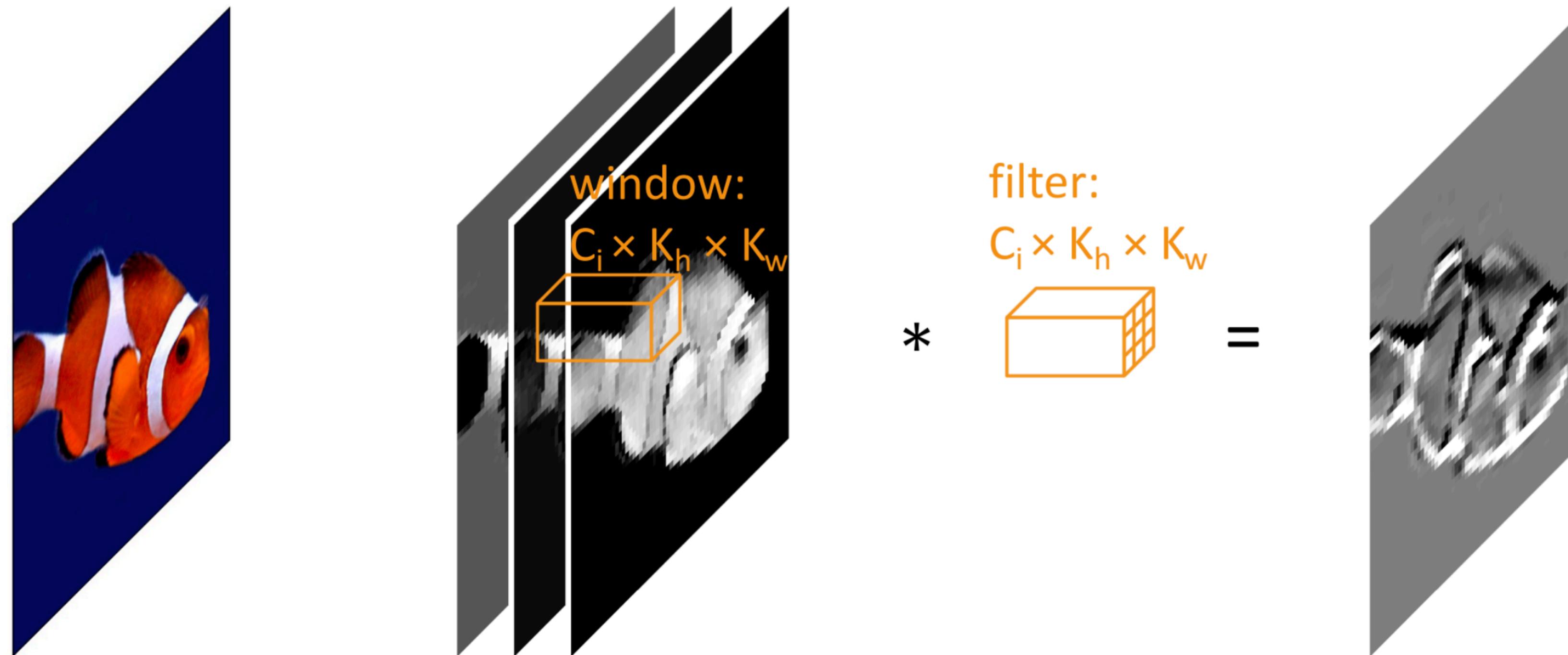
one filter, one feature

$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} & = \end{matrix}$$



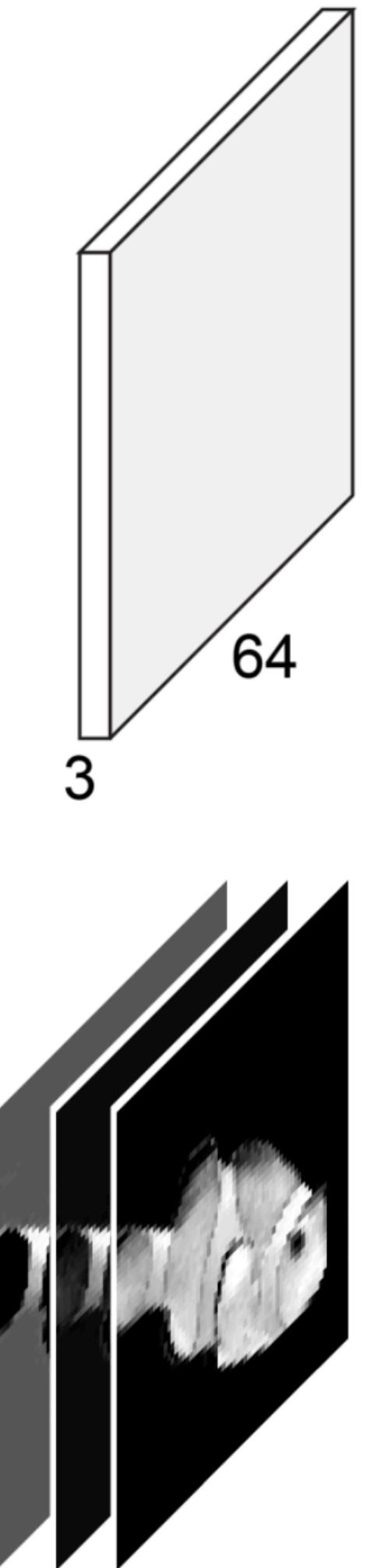
Understanding the filter/kernel as feature extractors

Convolution: Multi-channel inputs



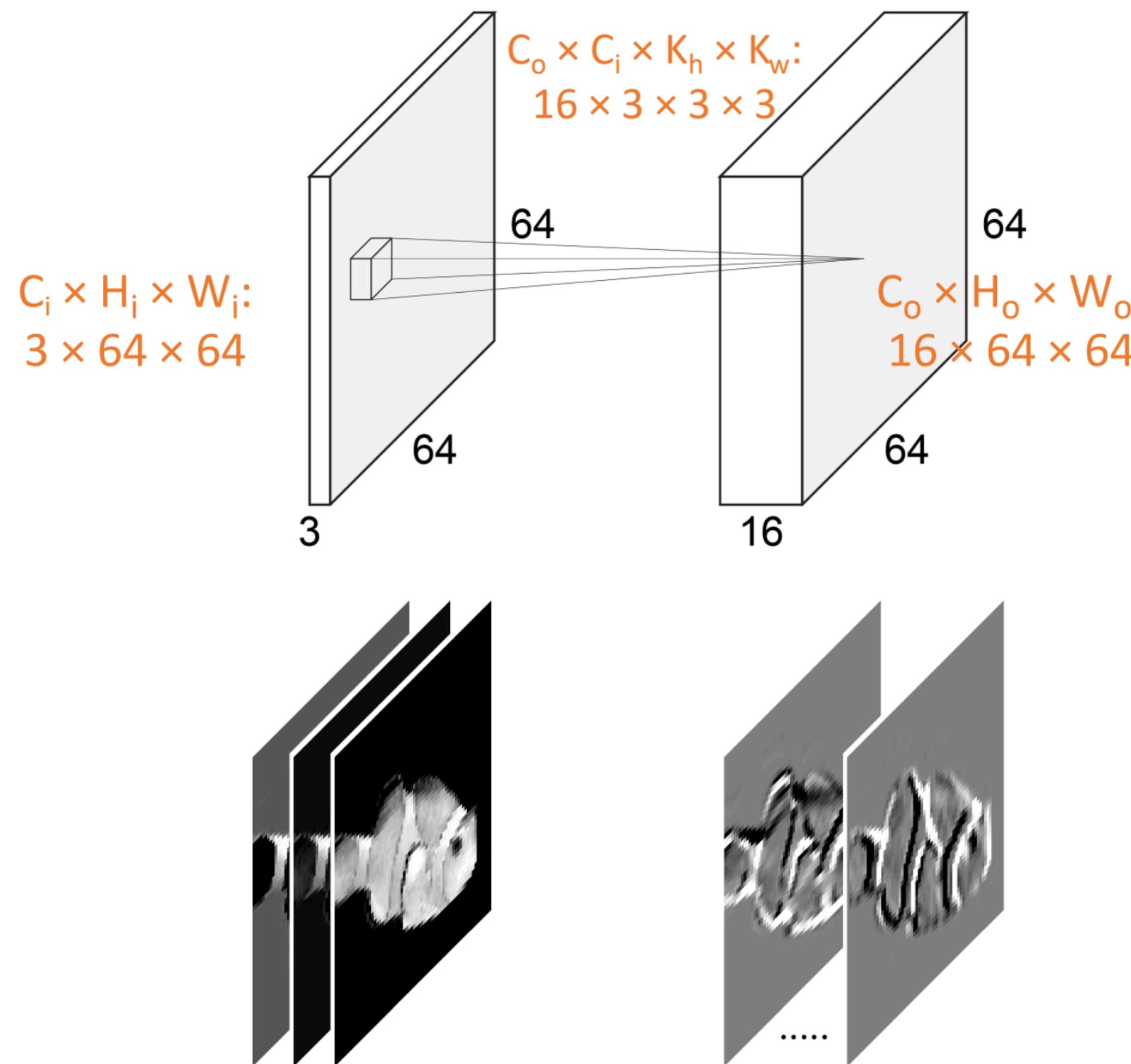
Like (R, G, B) color notations have three features

Convolution: tensor views



- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C: channels
 - H: height
 - W: width

Convolution: tensor view



- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C : channels
 - H : height
 - W : width
- filters
 - 4-D tensor: $C_o \times C_i \times K_h \times K_w$
 - C_o : output channels
 - C_i : input channels
 - K_h, K_w : filter height, width

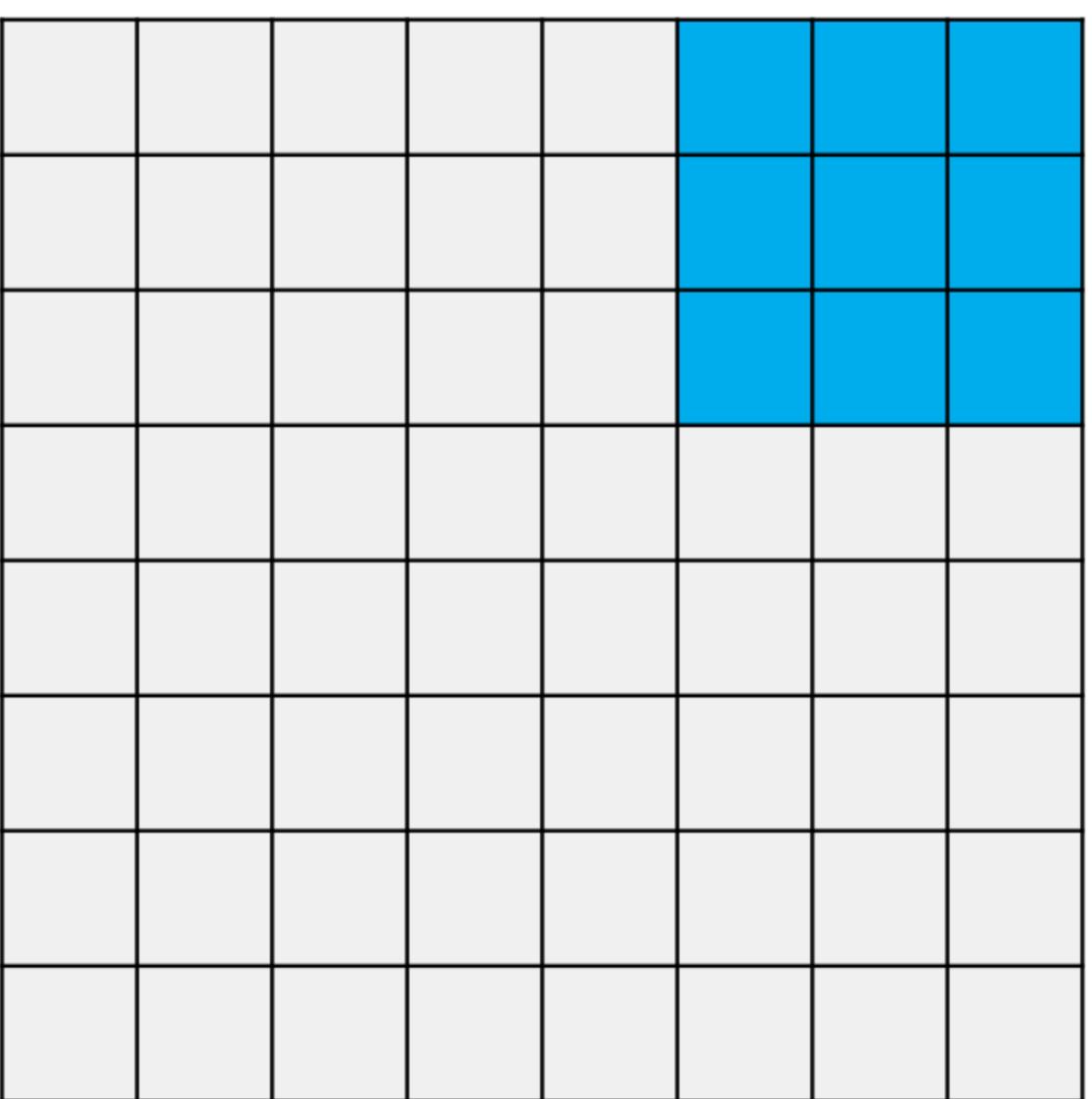
The same filter tensor applies to different locations

Convolution: # parameters and # operations

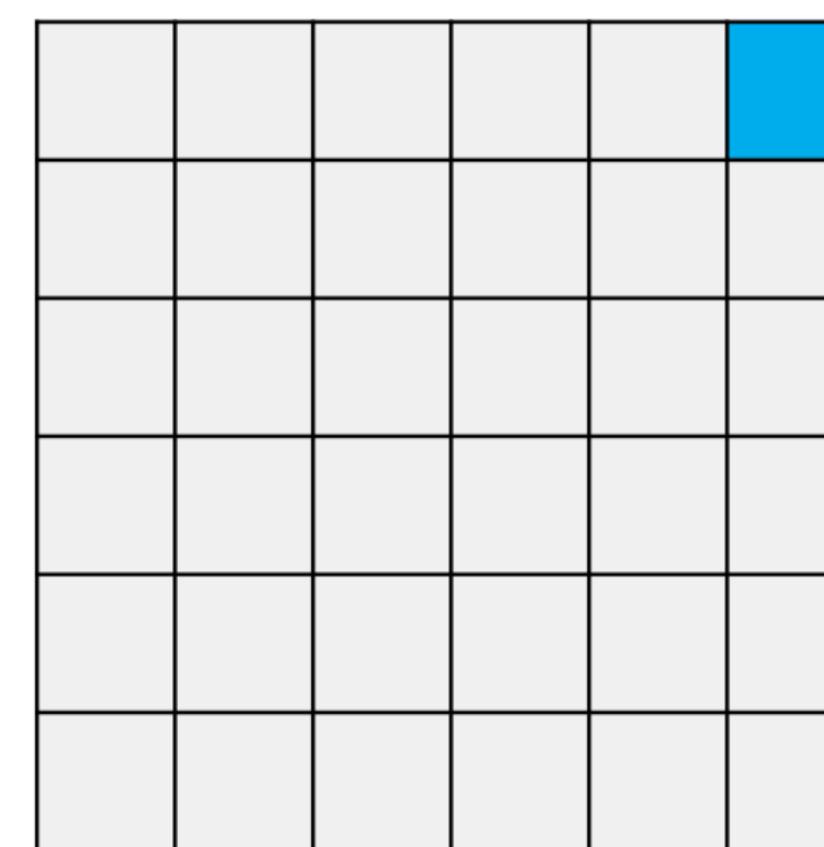
- # parameters
 - weights: $C_o \times C_i \times K_h \times K_w$
 - bias: C_o
- # floating-point operations (FLOPs)
 - # params $\times H_o \times W_o$

Convolution: padding

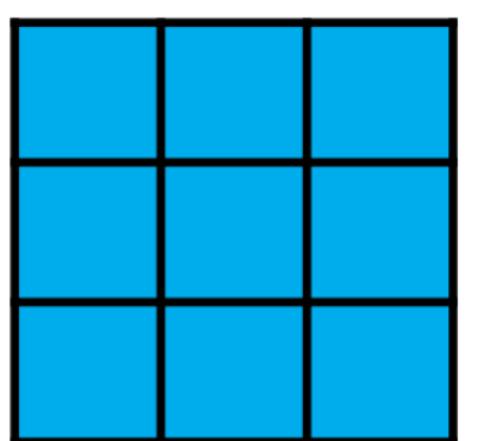
input: $H \times W = 8 \times 8$



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



filter



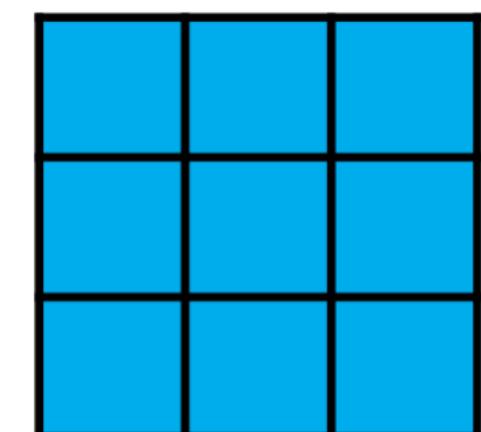
$$H_{\text{out}} = H_{\text{in}} - K_h + 1$$

Convolution: padding

input: 8×8 , + pad

0	0	0	0	0	0	0	0	0	0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0	0	0	0	0	0	0	0	0	0

filter

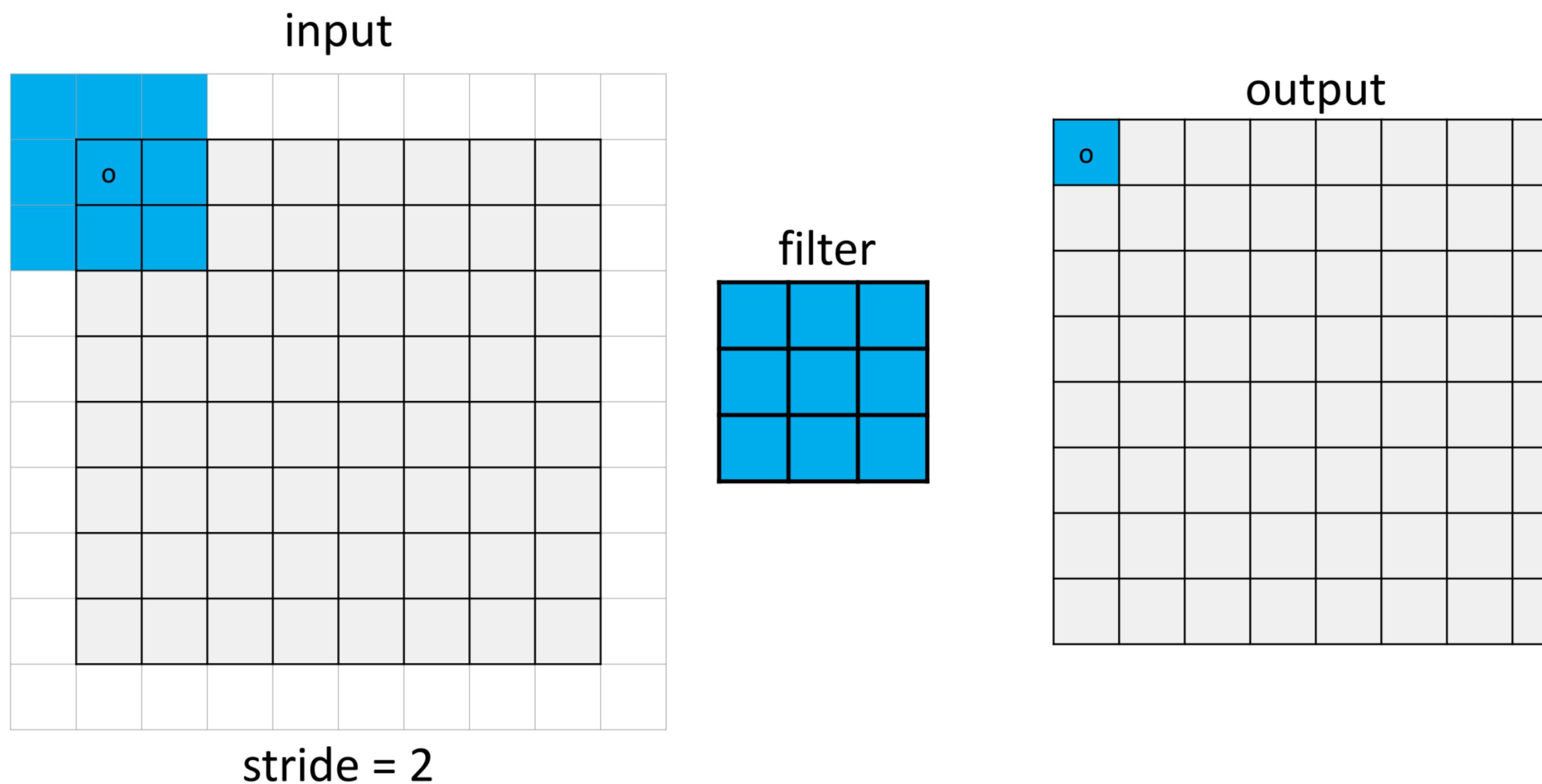


output: $H \times W = 8 \times 8$

- $\text{pad} = [\text{kernel_size} / 2]$
- maintains feature map size

$$H_{\text{out}} = H_{\text{in}} + 2\text{pad}_h - K_h + 1$$

Convolution: stride

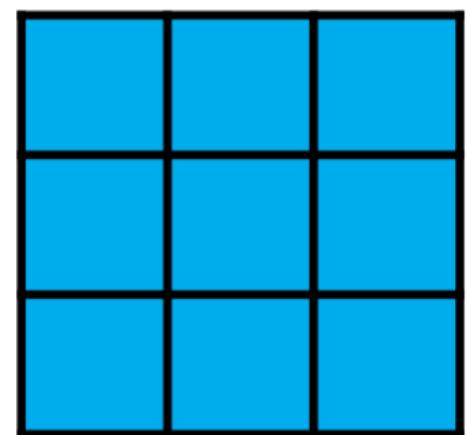


Convolution: stride

input

stride = 2

filter



output

Convolution: stride

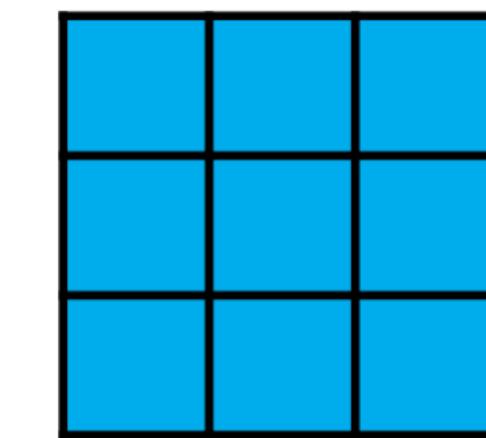
input: $H \times W = 8 \times 8$

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

stride = 2

- reduces feature map size
- compress and abstract

filter



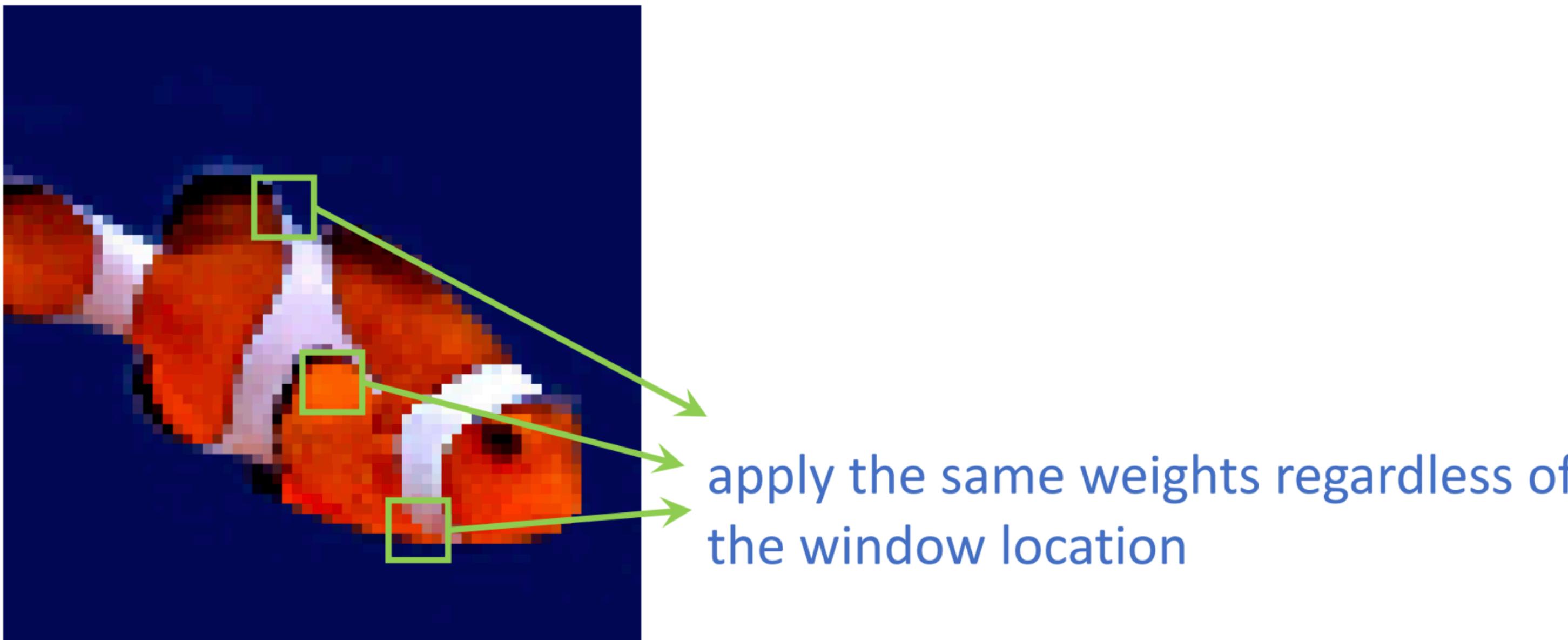
output: $H \times W = 4 \times 4$

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

$$H_{out} = \lfloor (H_{in} + 2\text{pad}_h - K_h) / \text{str} \rfloor + 1$$

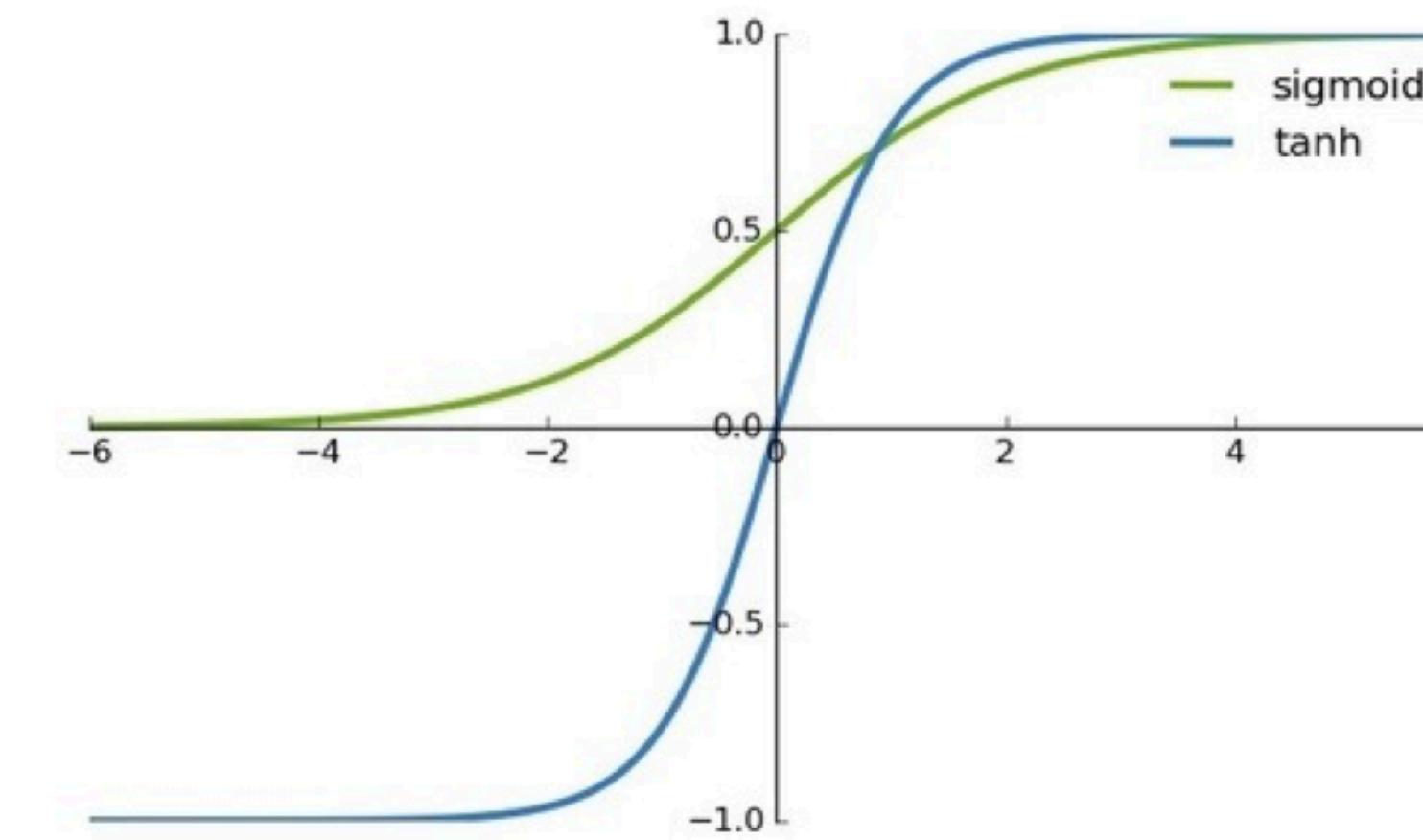
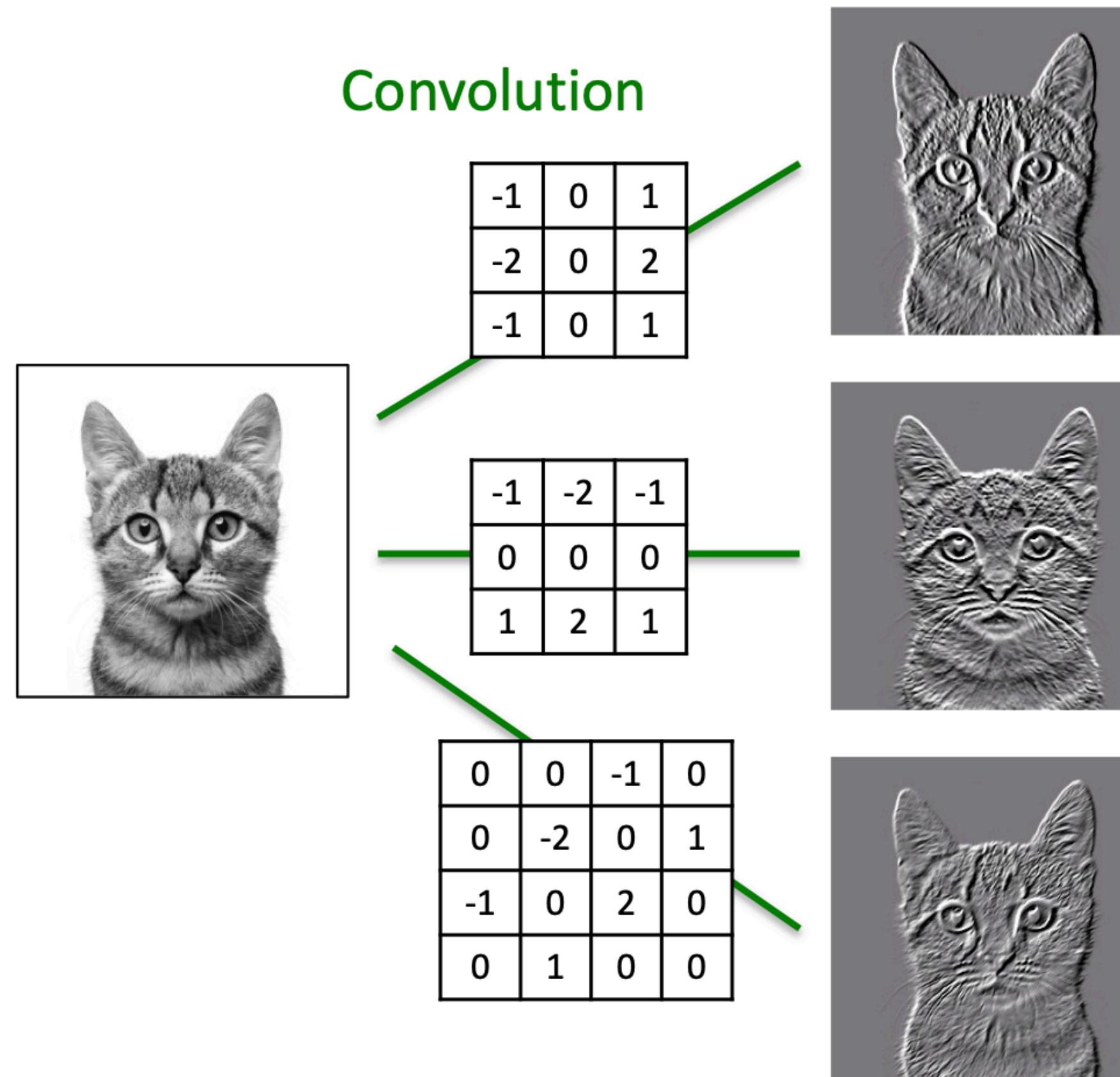
Convolution: translation-invariance

- Process each window in the same way



Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling

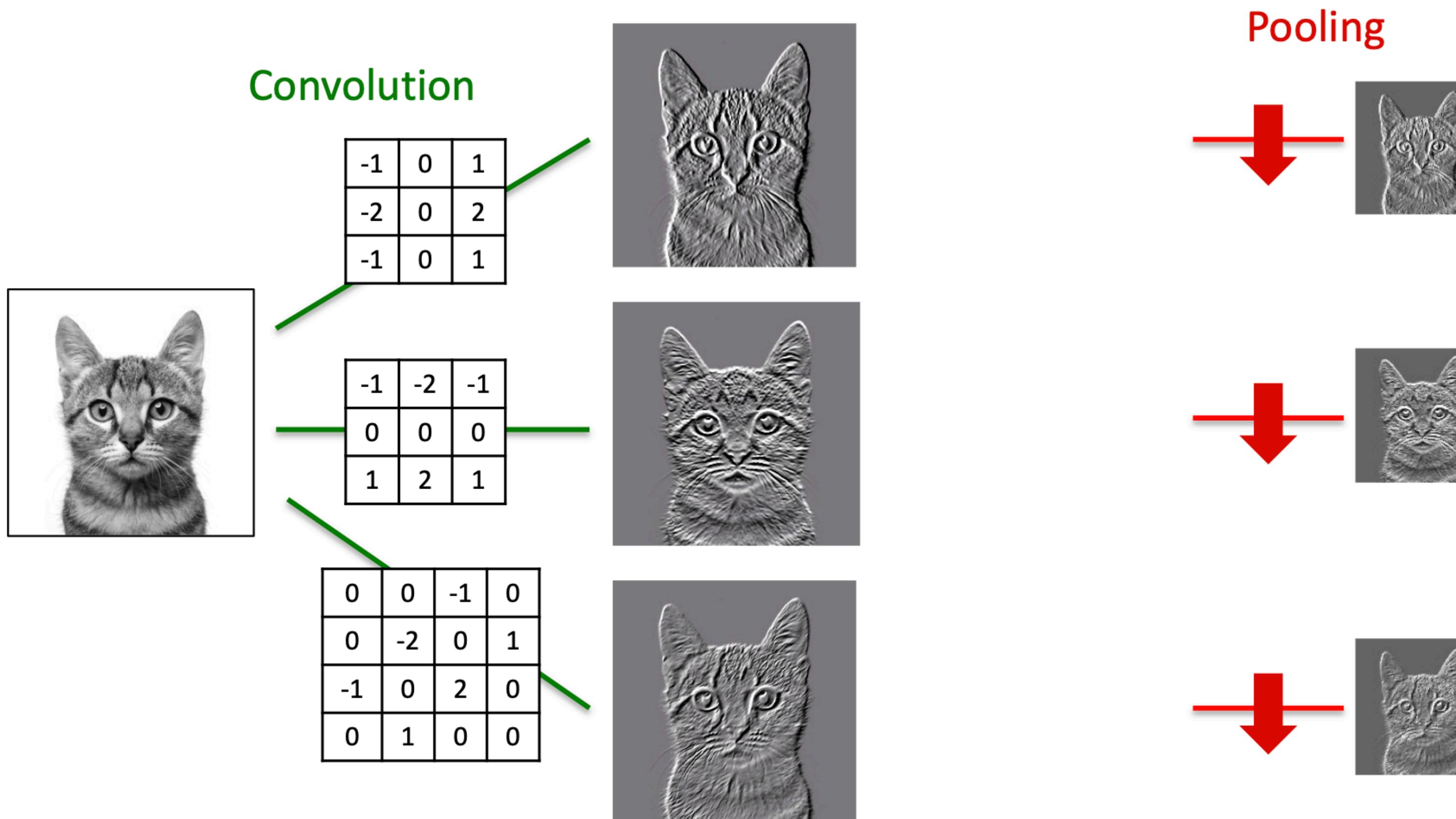


LeNet – tanh activation

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

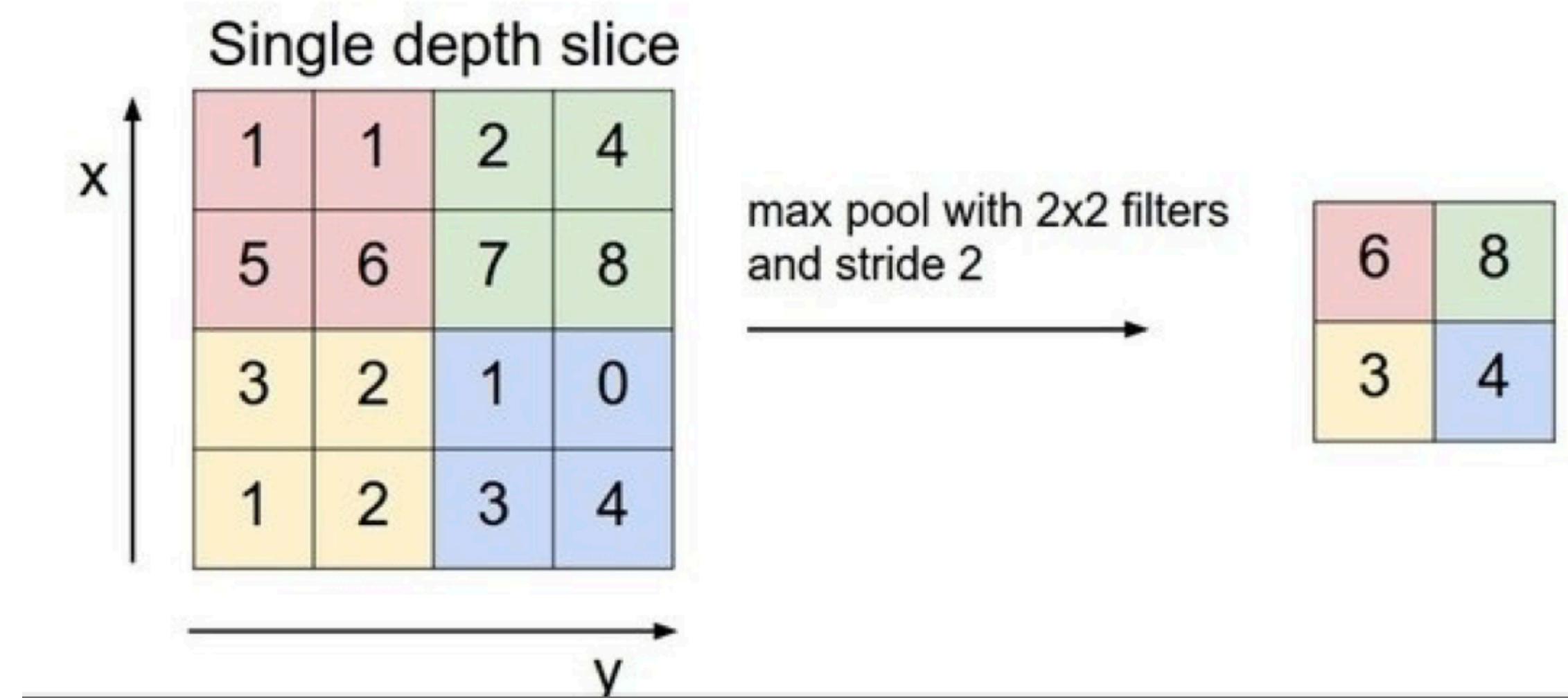
Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling

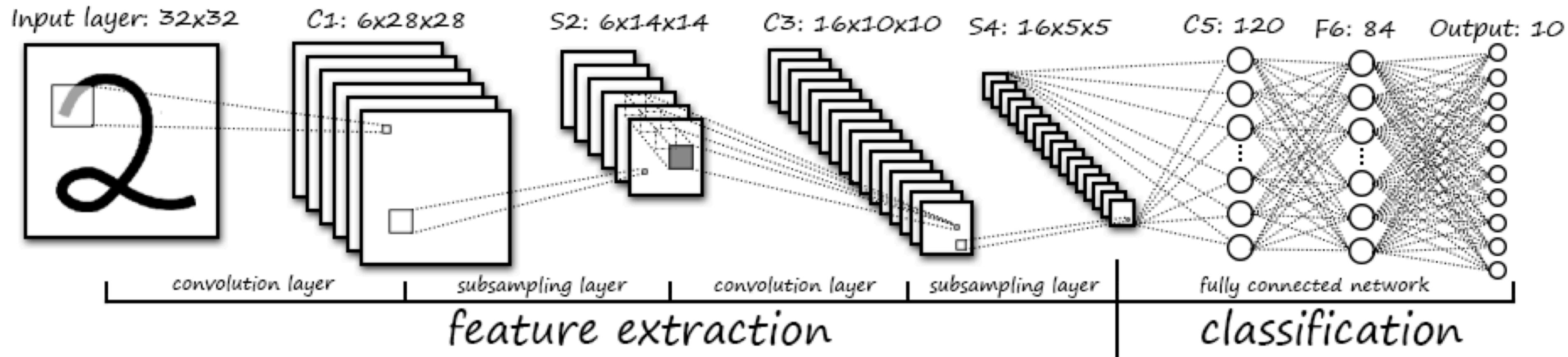


Pooling = Down-sampling

Max pooling



Deep Convolutional Networks



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

60,000 original dataset

Test error: 0.95%

MNIST

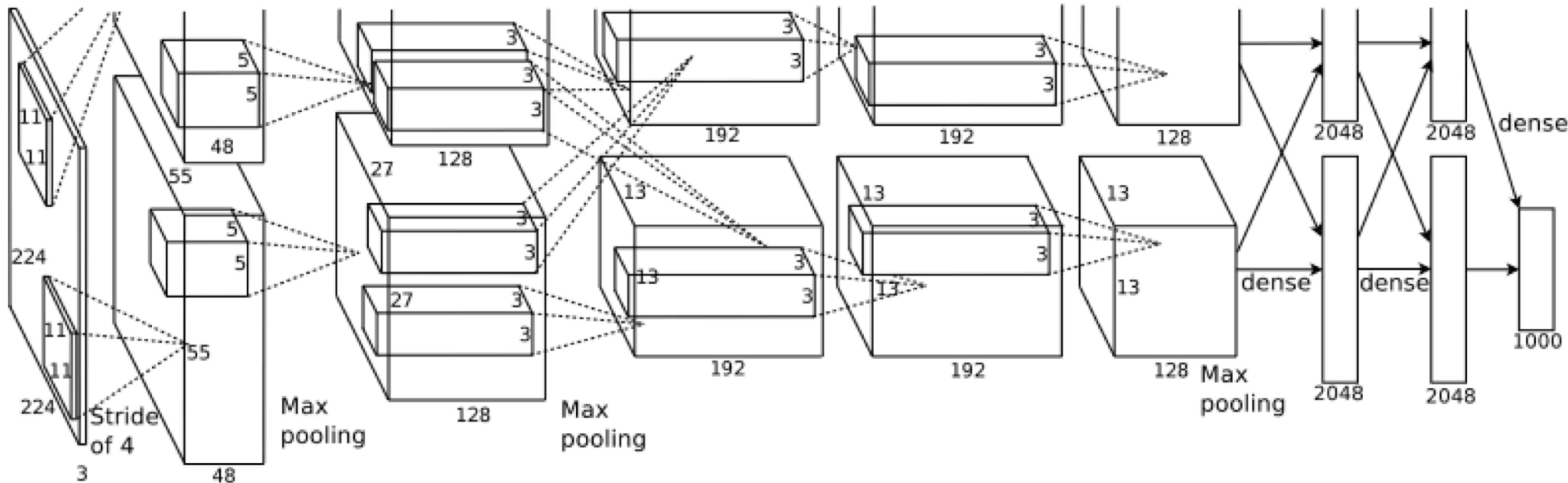
[1] LeNet 5, LeCun et al. 1998

Misclassified examples on MNIST

True label -> Predicted label

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

Alex Net



[1] Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NeurIPS 2012.

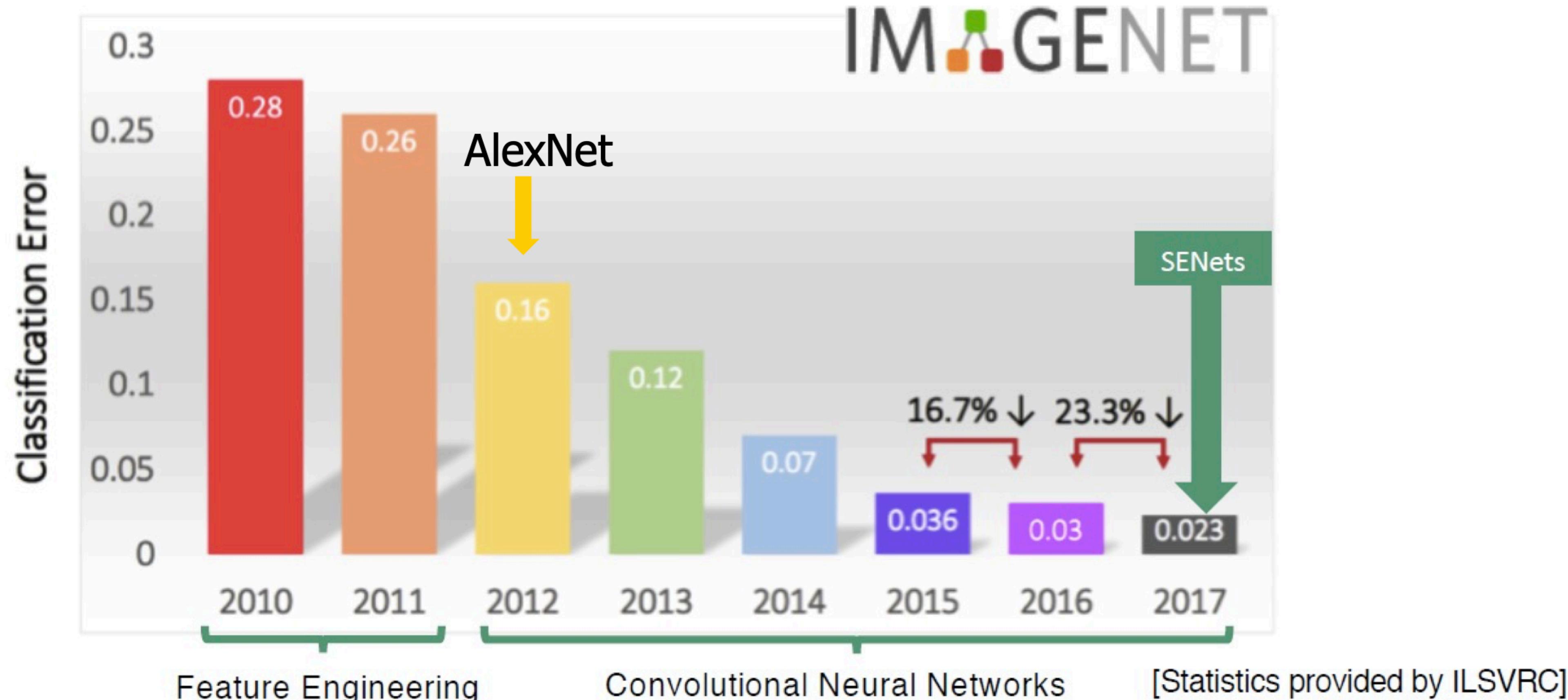
ImageNet

- ❑ 15M images
- ❑ 22K categories
- ❑ Images collected from Web
- ❑ Human labelers (Amazon's Mechanical Turk crowd-sourcing)
- ❑ ImageNet Large Scale Visual Recognition Challenge (ILSVRC-2010)
 - 1K categories
 - 1.2M training images (~1000 per category)
 - 50,000 validation images
 - 150,000 testing images
- ❑ RGB images
- ❑ Variable-resolution, but this architecture scales them to 256x256 size

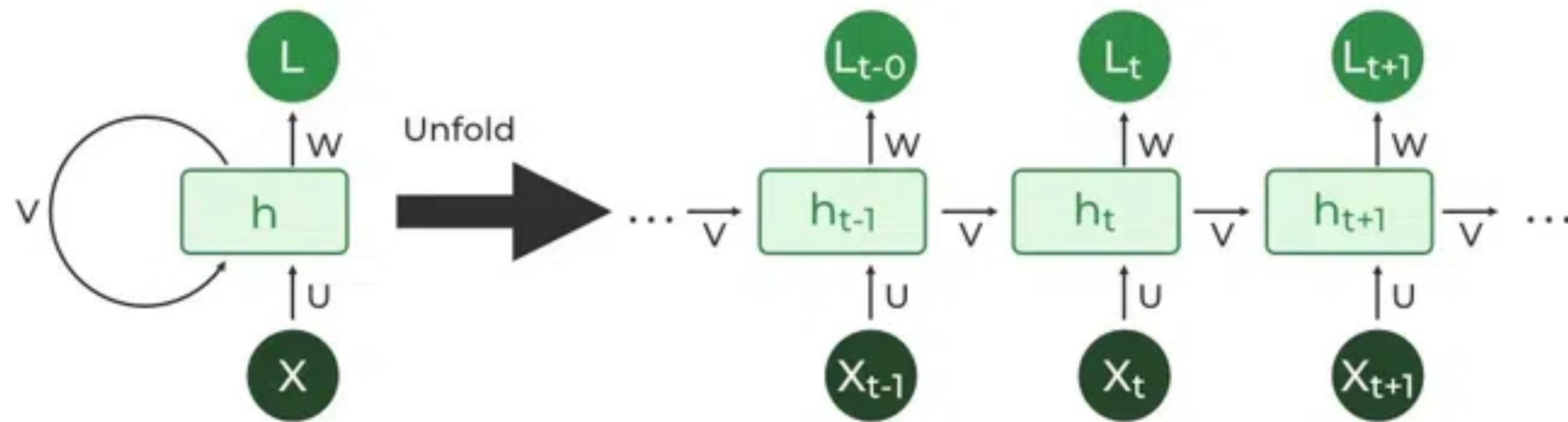
ImageNet Results



ImageNet Results



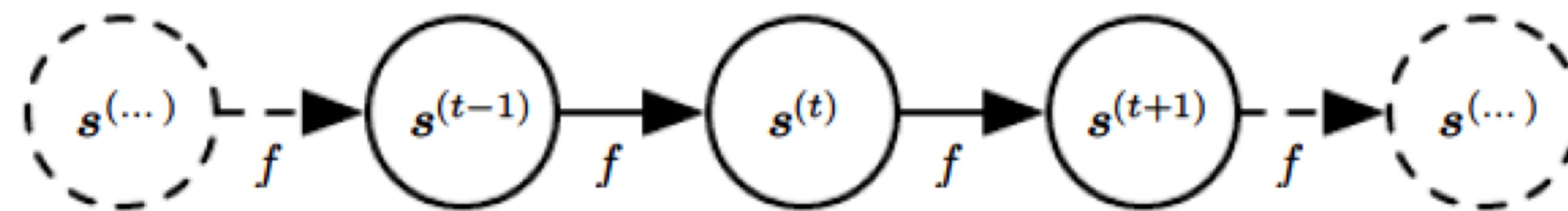
Recurrent Neural Networks (RNNs)



Recurrent Neural Networks

- Dates back to (Rumelhart *et al.*, 1986)
- A family of neural networks for handling sequential data, which involves variable length inputs or outputs
- Especially, for natural language processing (NLP)

Computation Graph

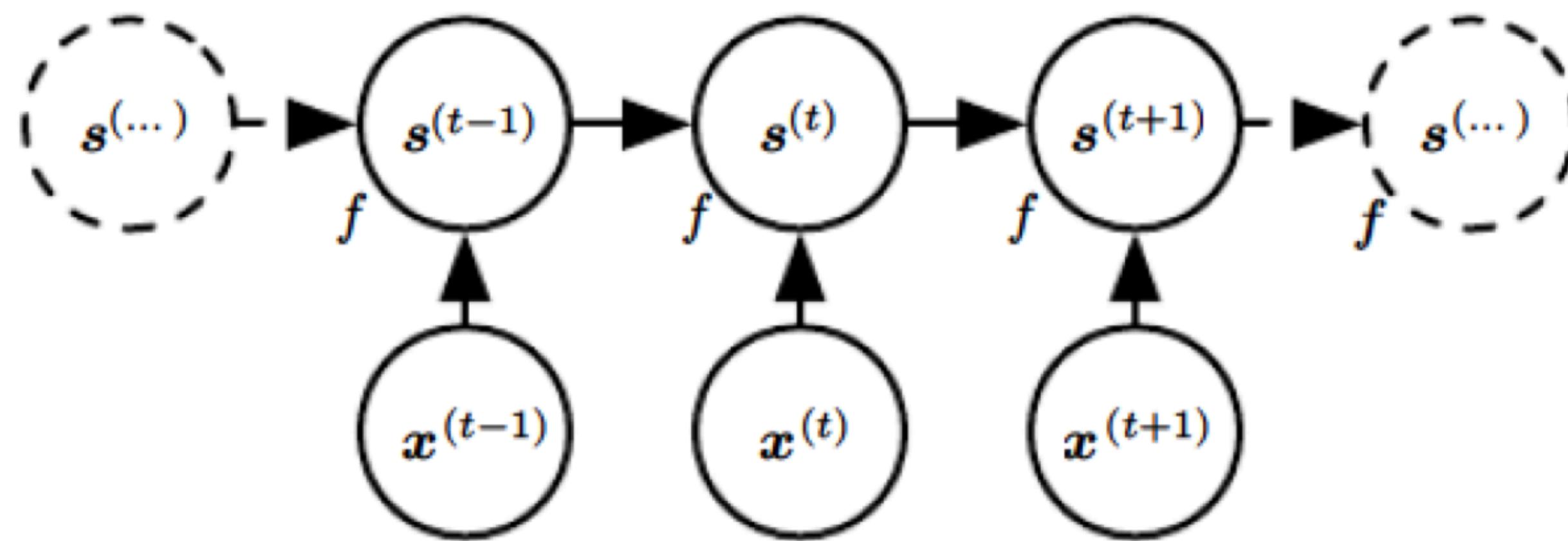


$$s^{(t+1)} = f(s^{(t)}; \theta)$$

Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

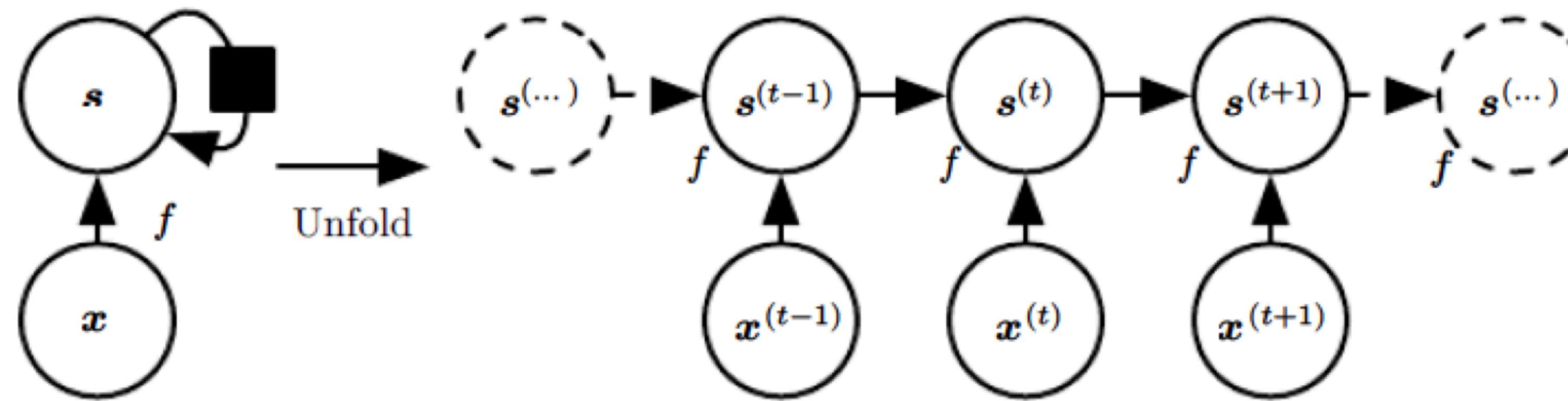
Like Markov model, but here $s^{(t+1)}$ is deterministic given $s^{(t)}$

Computation Graph



$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Compact view



$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Key: the same f and θ
for all time steps

Like stationary HMM

Recurrent Neural Networks

- Use **the same** computational function and parameters across different time steps of the sequence
- Each time step: takes the input entry **and the previous hidden state** to compute the output entry
- Loss: typically computed every time step

Recurrent Neural Networks

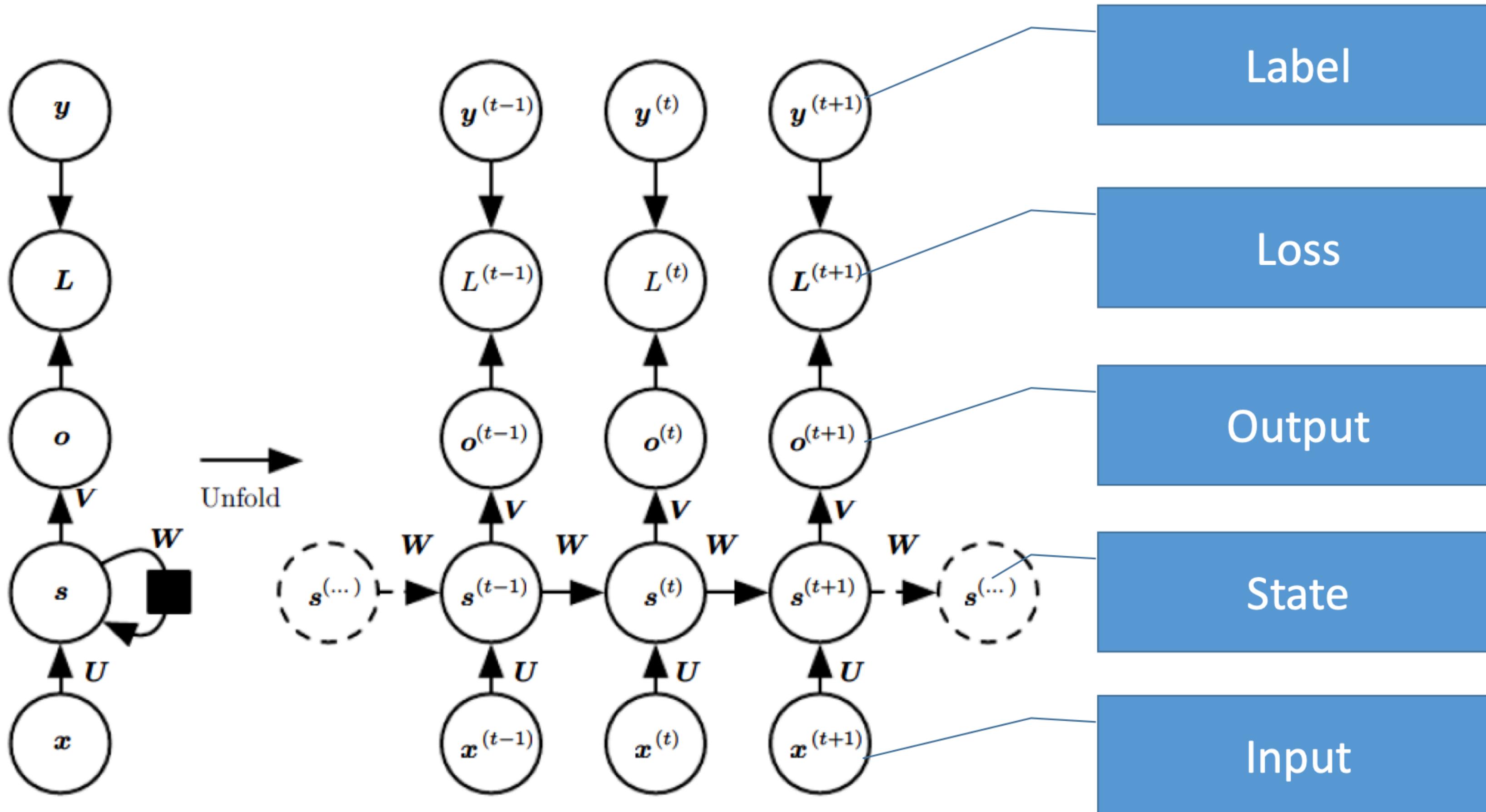


Figure from *Deep Learning*, by Goodfellow, Bengio and Courville

Recurrent Neural Networks

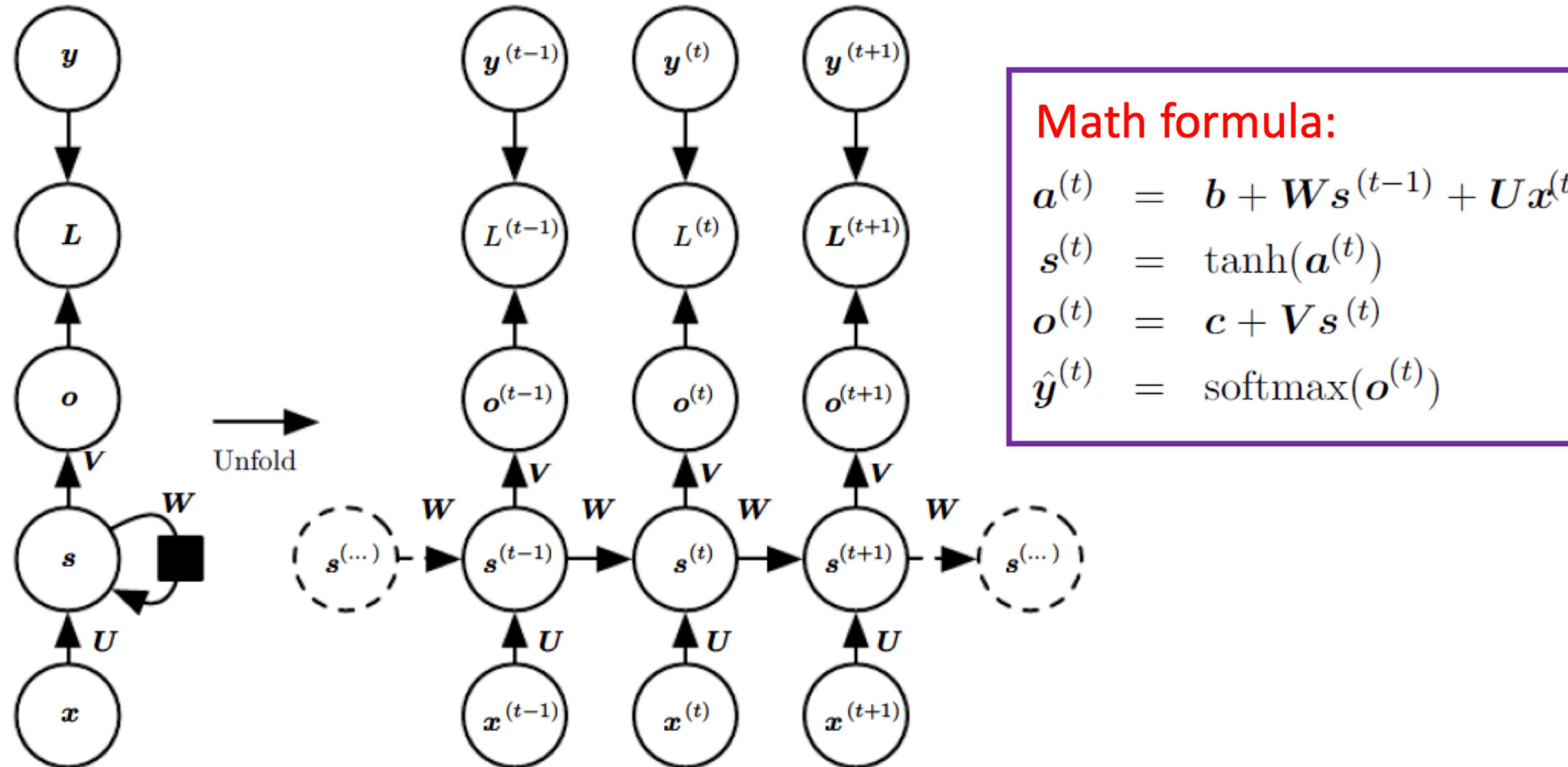
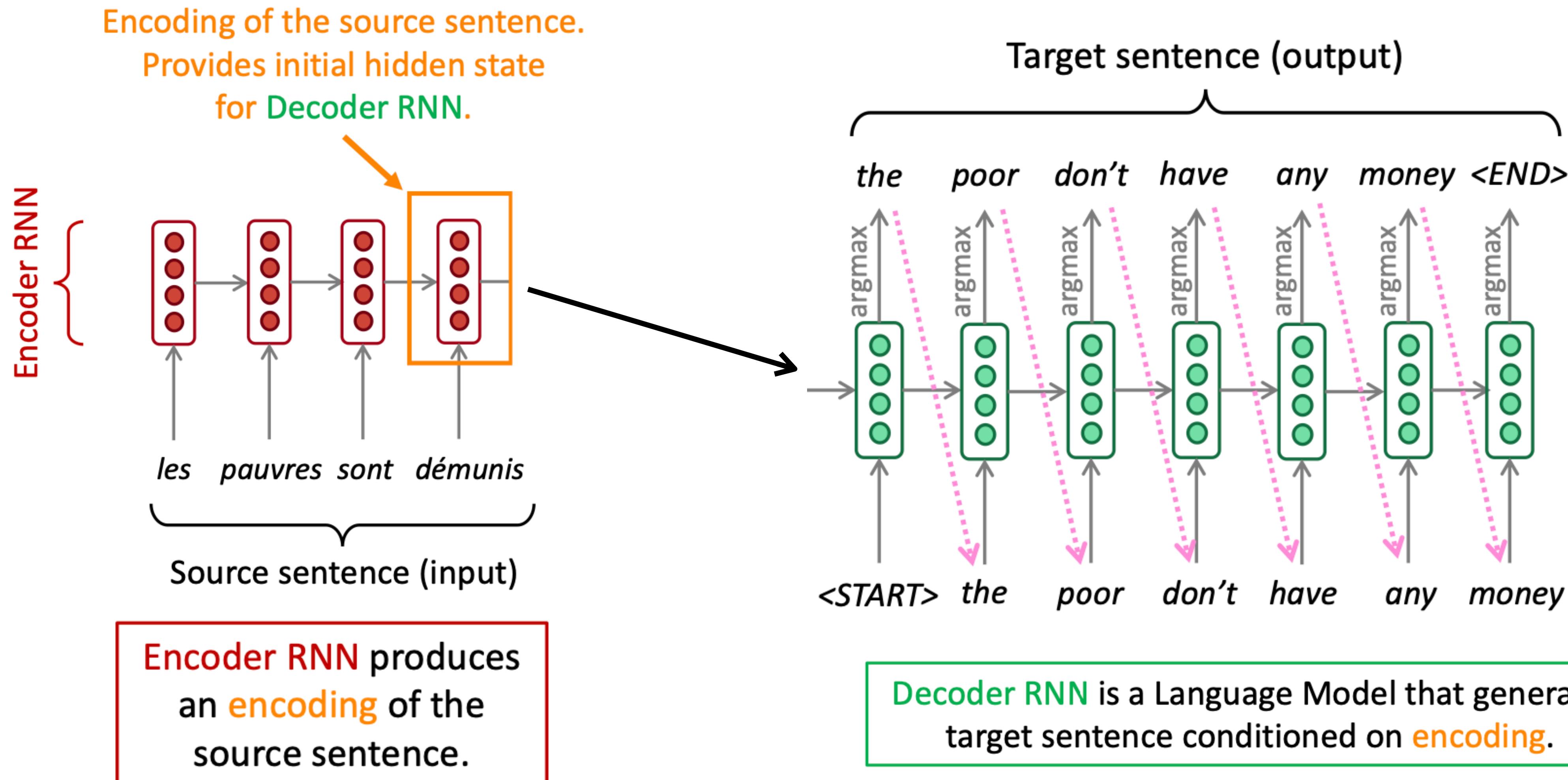


Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

There are many variants of RNNs since the functional form to compute $s^{(t)}$ can vary, e.g., LSTM

Sequence-to-Sequence Learning

Example of Neural Machine Translation



Residual Connection

We want deeper and deeper NNs, but going deep is difficult

- Troubles accumulate w/ more layers
- Signals get distorted when propagated
- in forward and backward passes

Commonly used techniques to train “Deep” NNs:

Weight initialization

Normalization modules

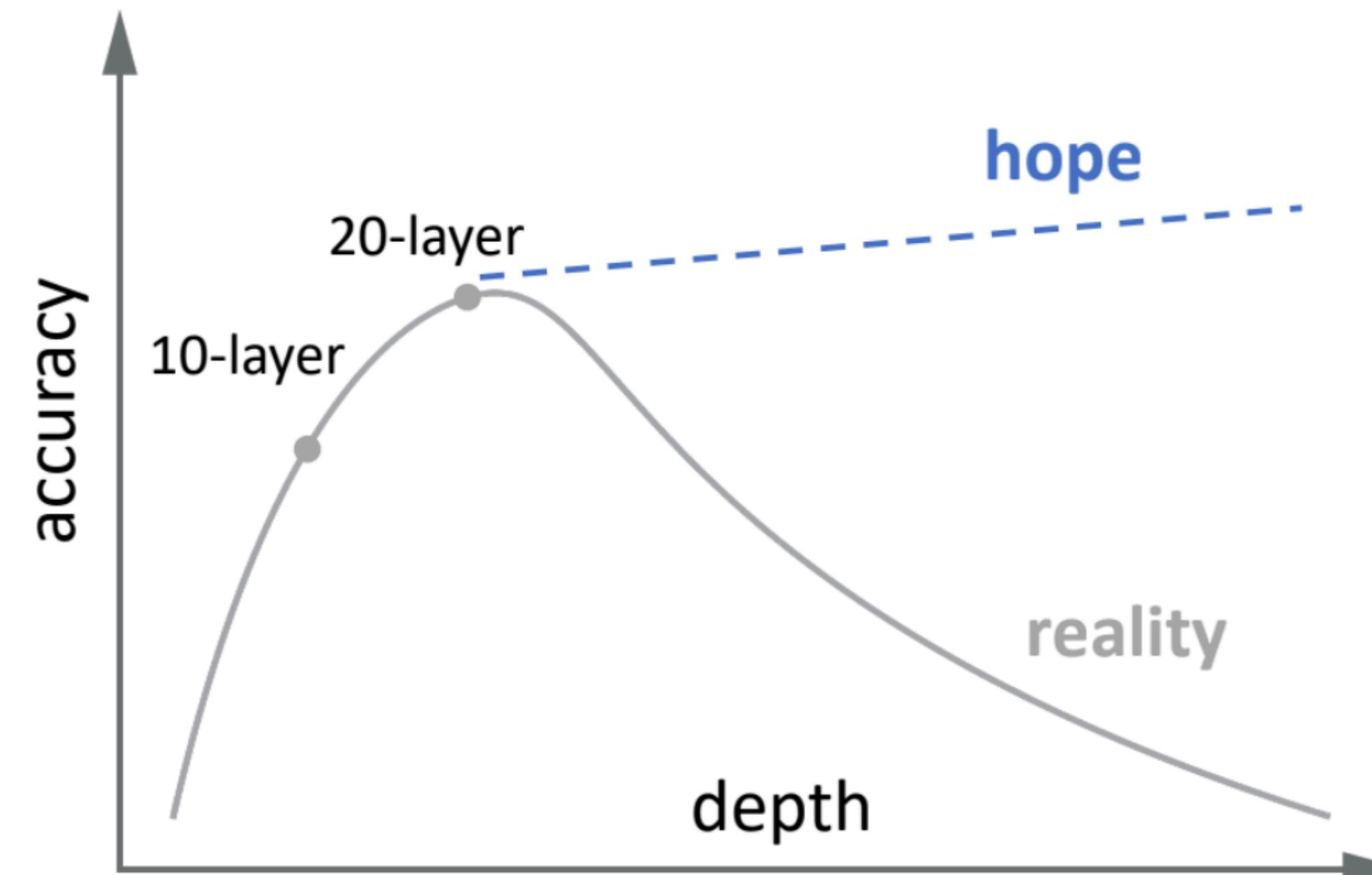
Deep residual learning



The Degradation Problem

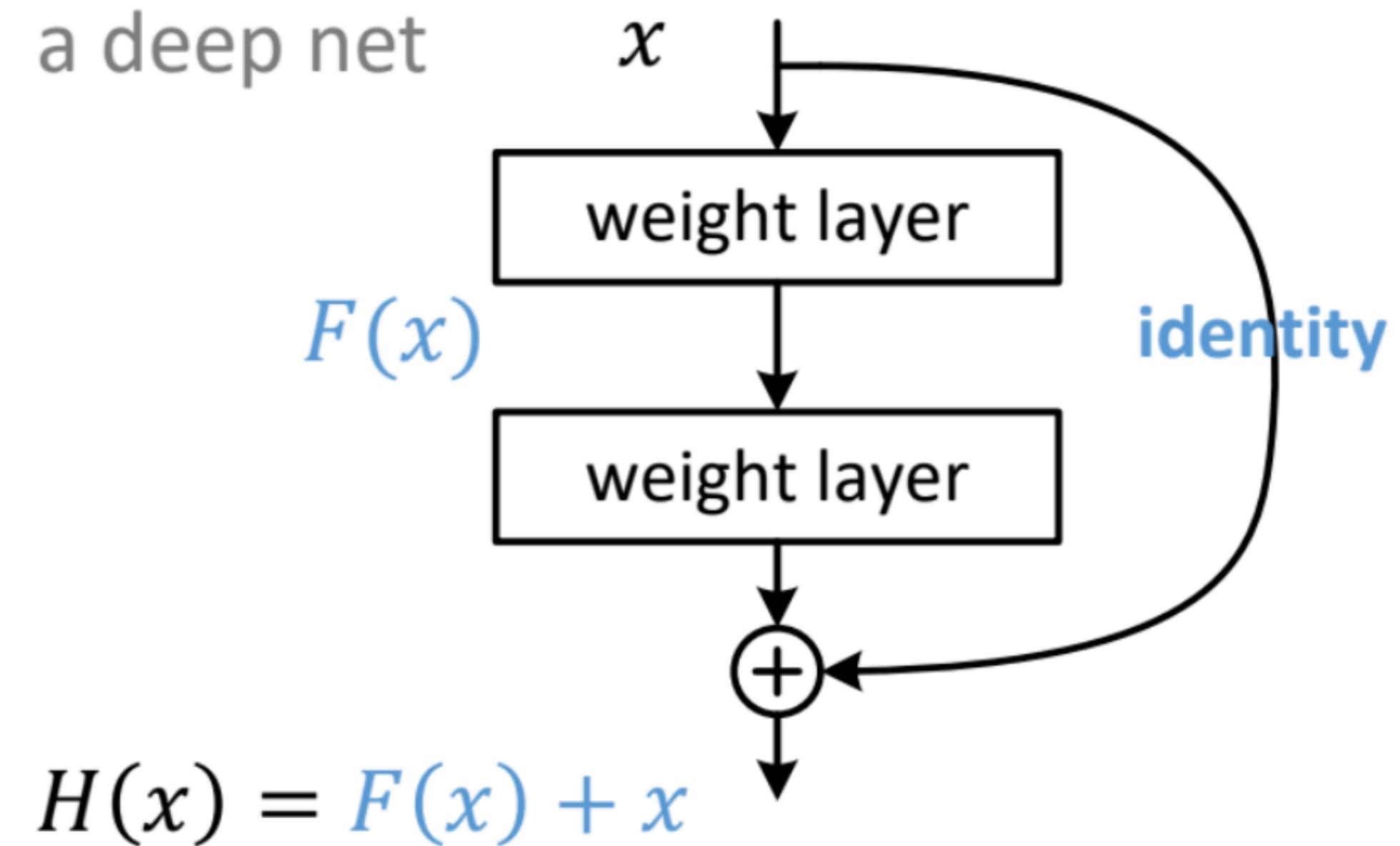
- Good init + norm enable training deeper models
- Simply stacking more layers?

- Degrade after ~ 20 layers
- Not overfitting
- Difficult to train



Deep Residual Learning

a subnet in
a deep net



Deep Residual Networks (ResNet)

ResNet-152



Kaiming He et al. Deep Residual Learning for Image Recognition. CVPR 2016.

Transformers

