



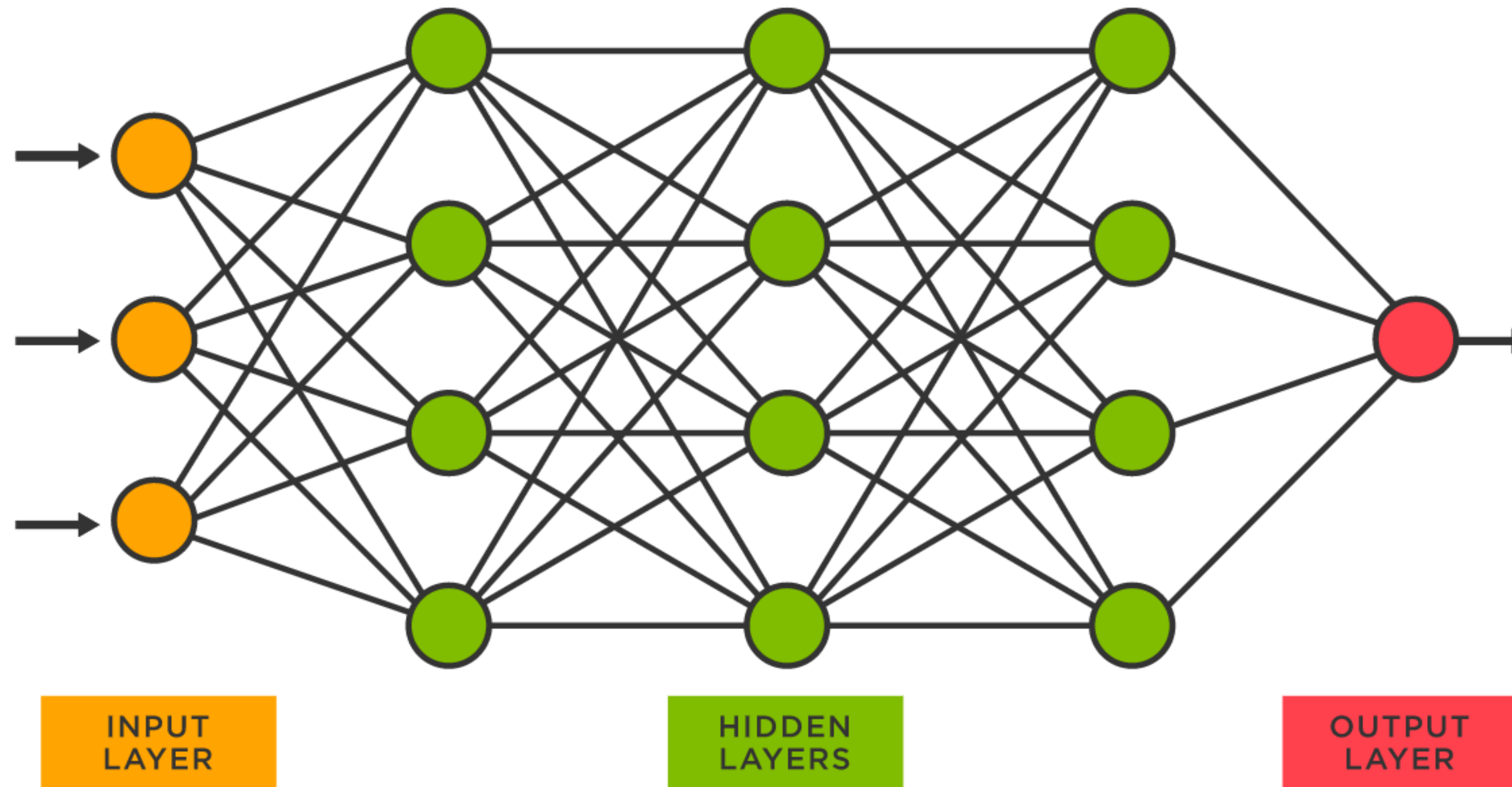
香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 19

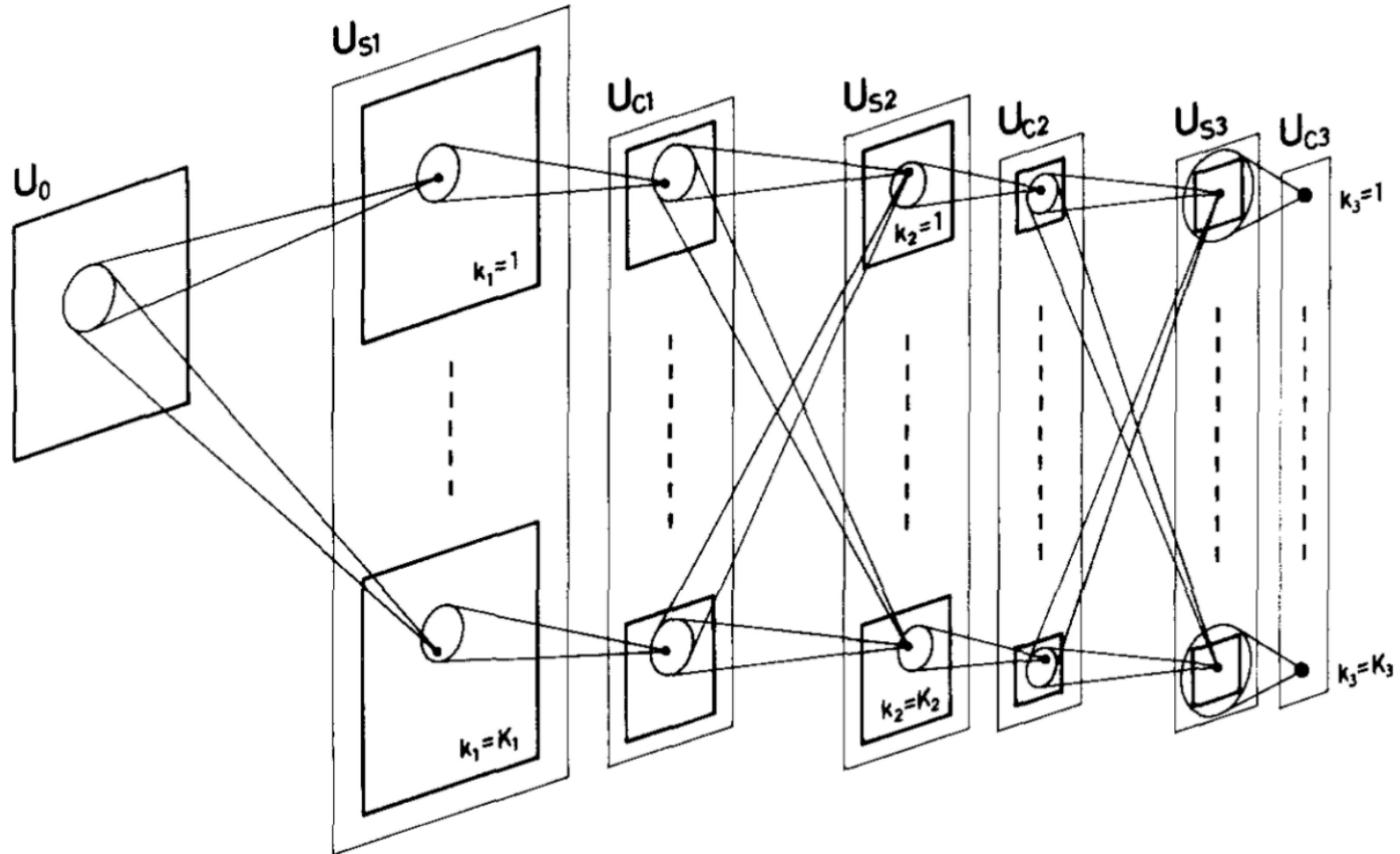
# Neural Networks, Architectures

Junxian He  
Apr 17, 2024

# Recap: Multilayer Perceptron Neural Networks (MLP)

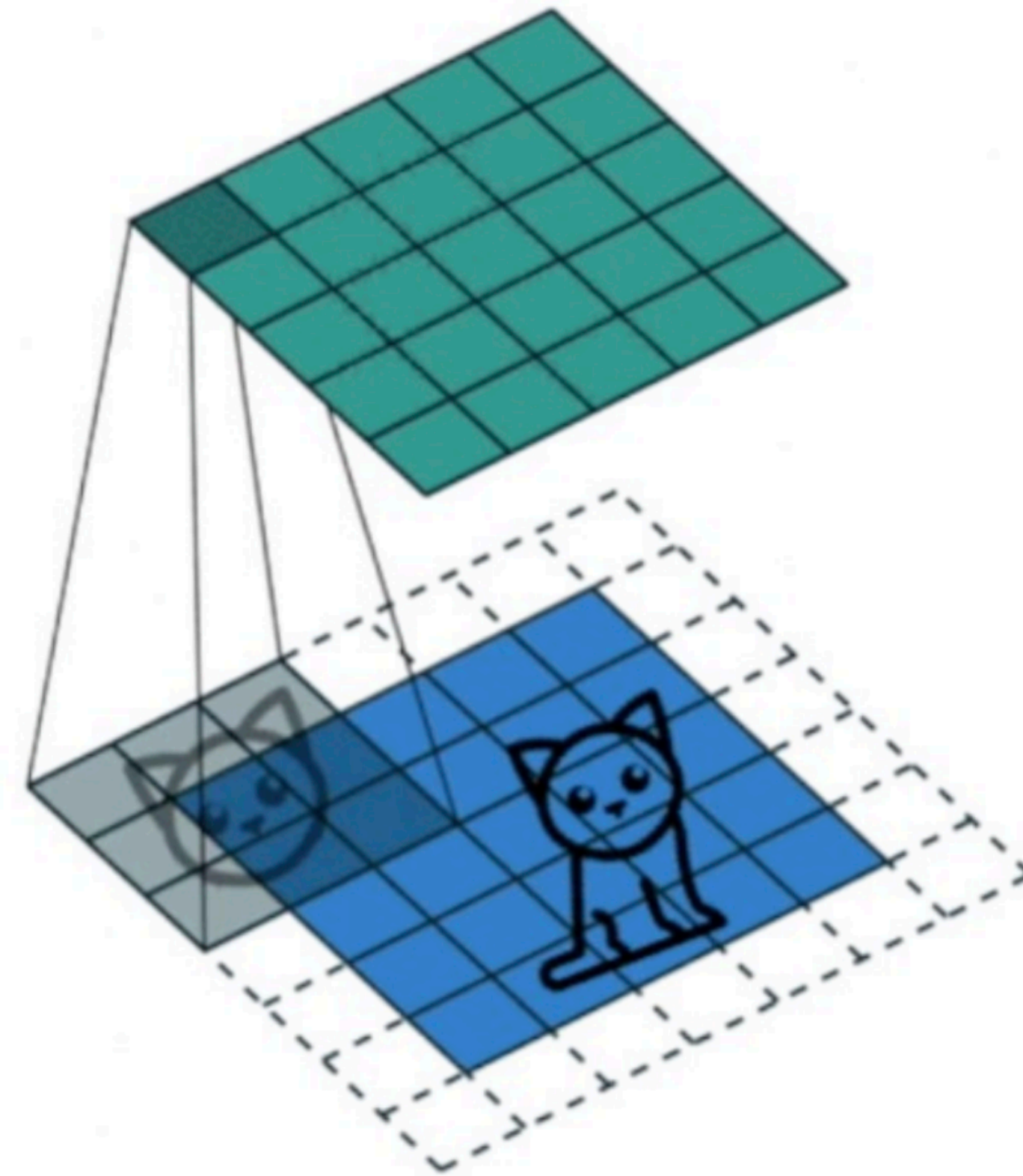


# Convolutional Neural Networks



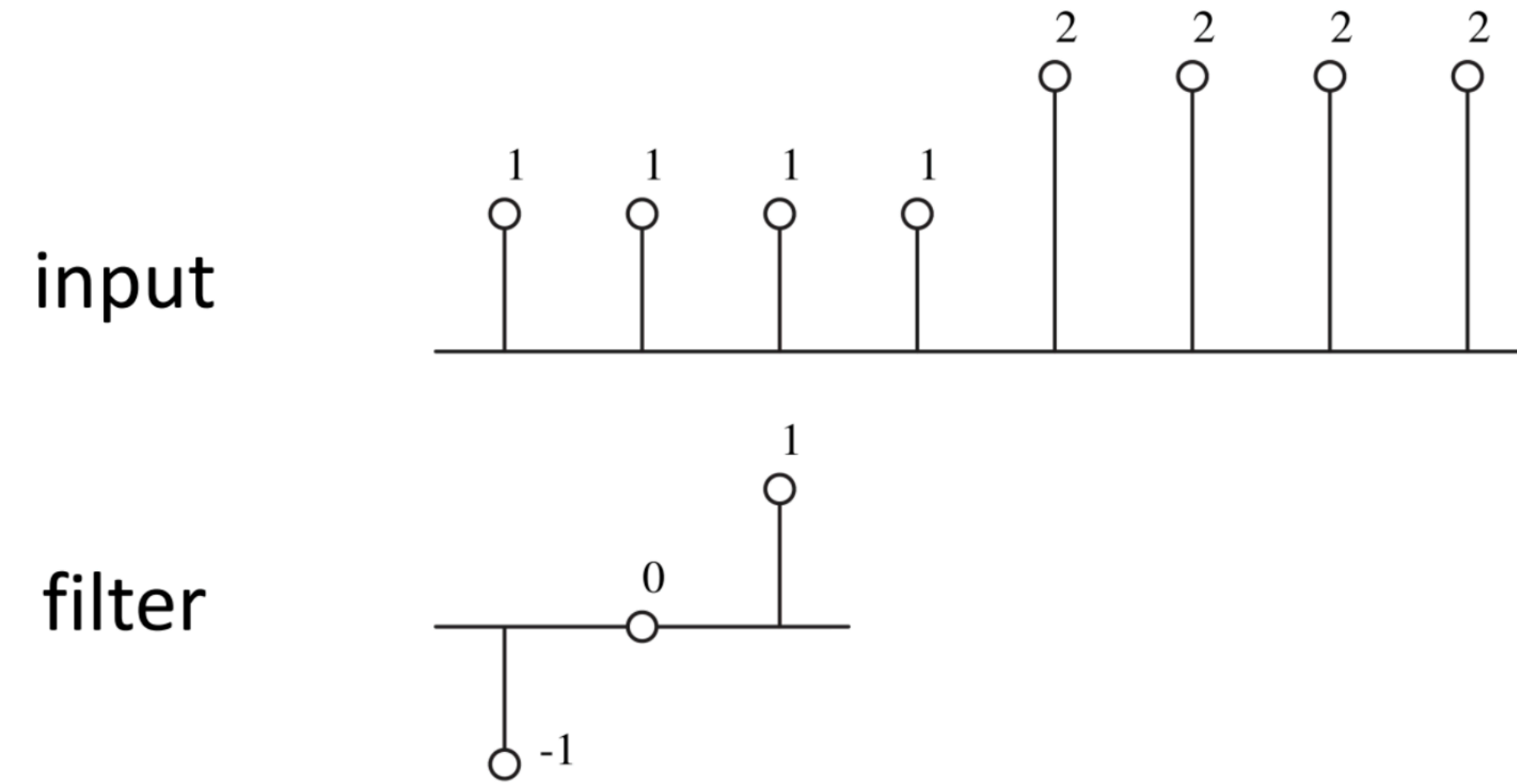
# Convolution is template matching

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





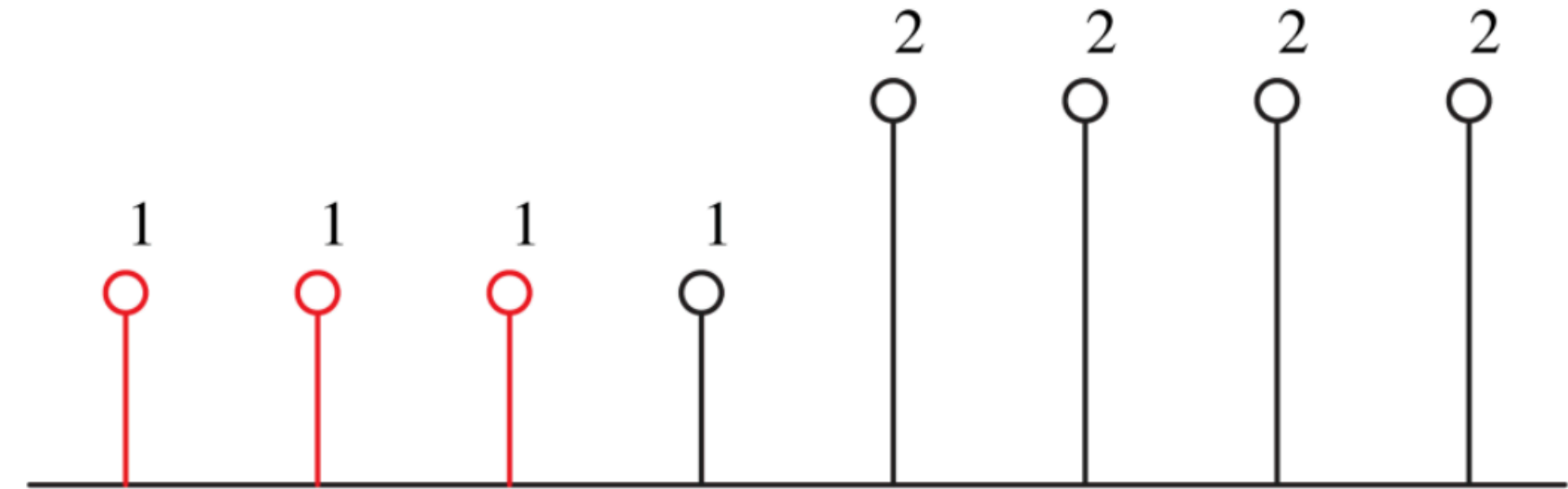
# Convolution: a 1-D example



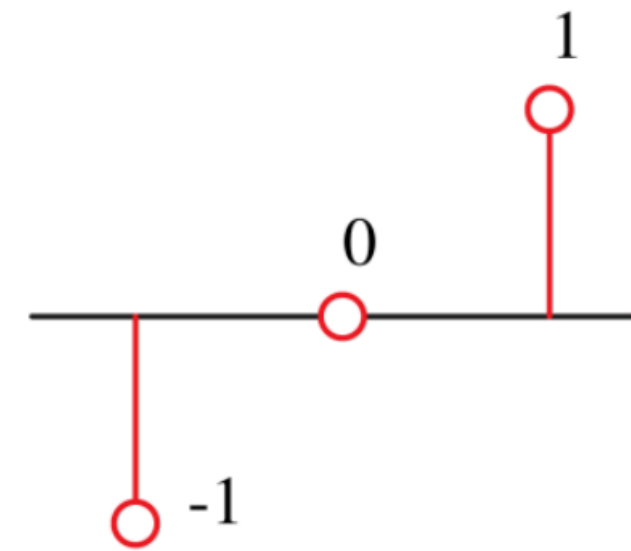
# Convolution: a 1-D example

- sliding window
- dot product

input



filter



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

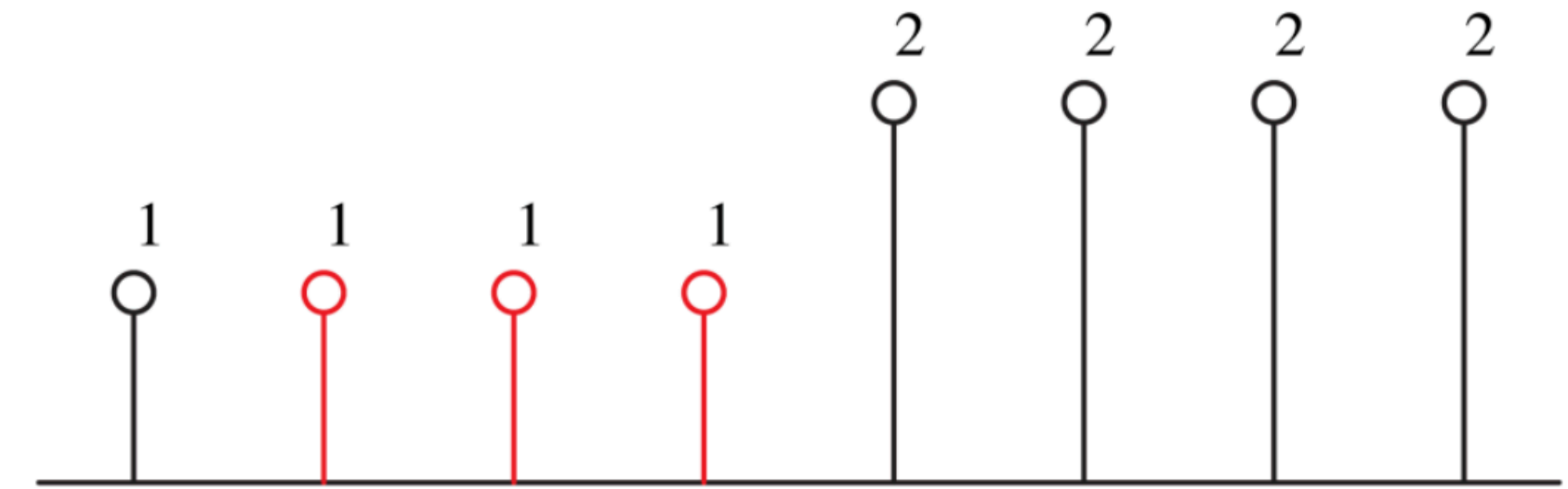
output



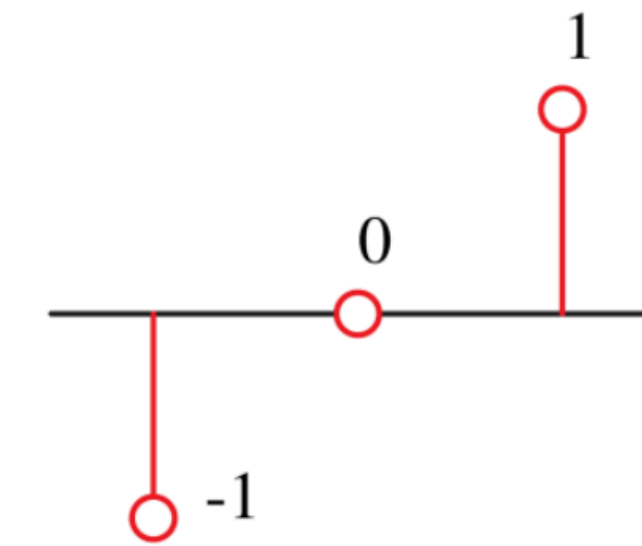
# Convolution: a 1-D example

- sliding window
- dot product

input



filter



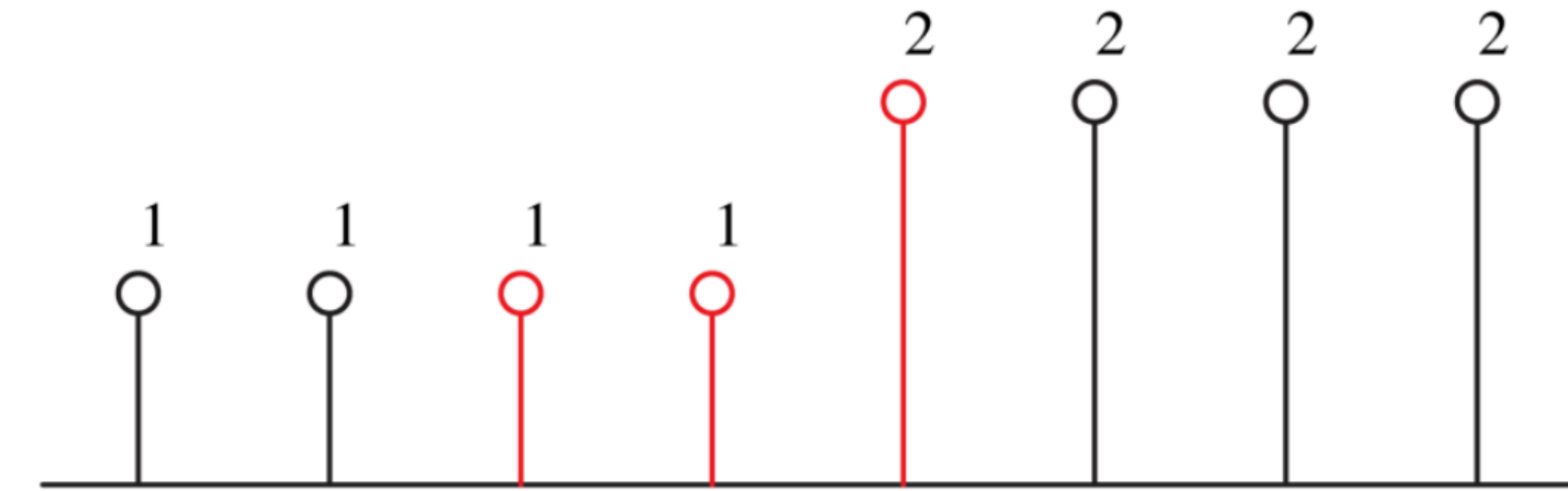
output



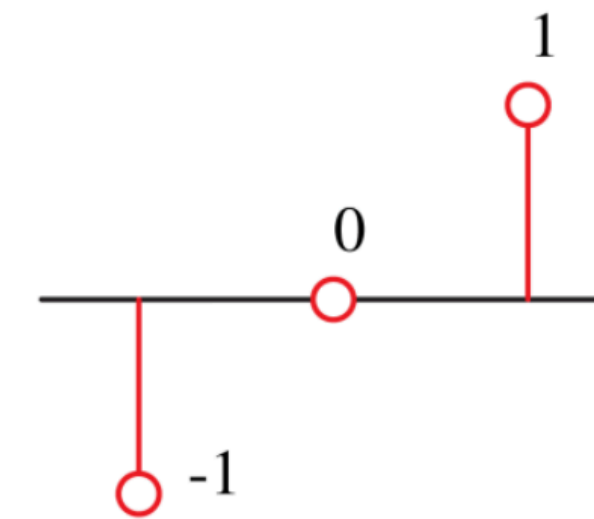
# Convolution: a 1-D example

- sliding window
- dot product

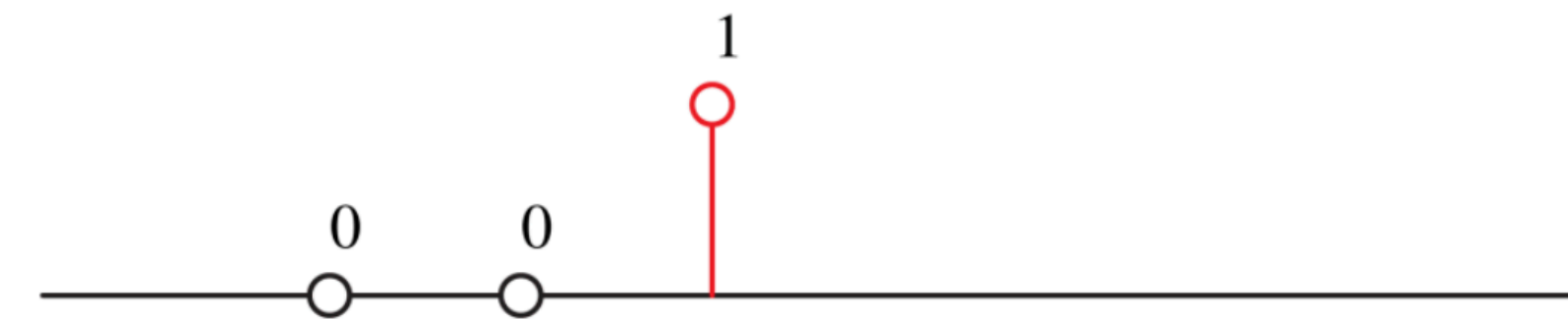
input



filter



output

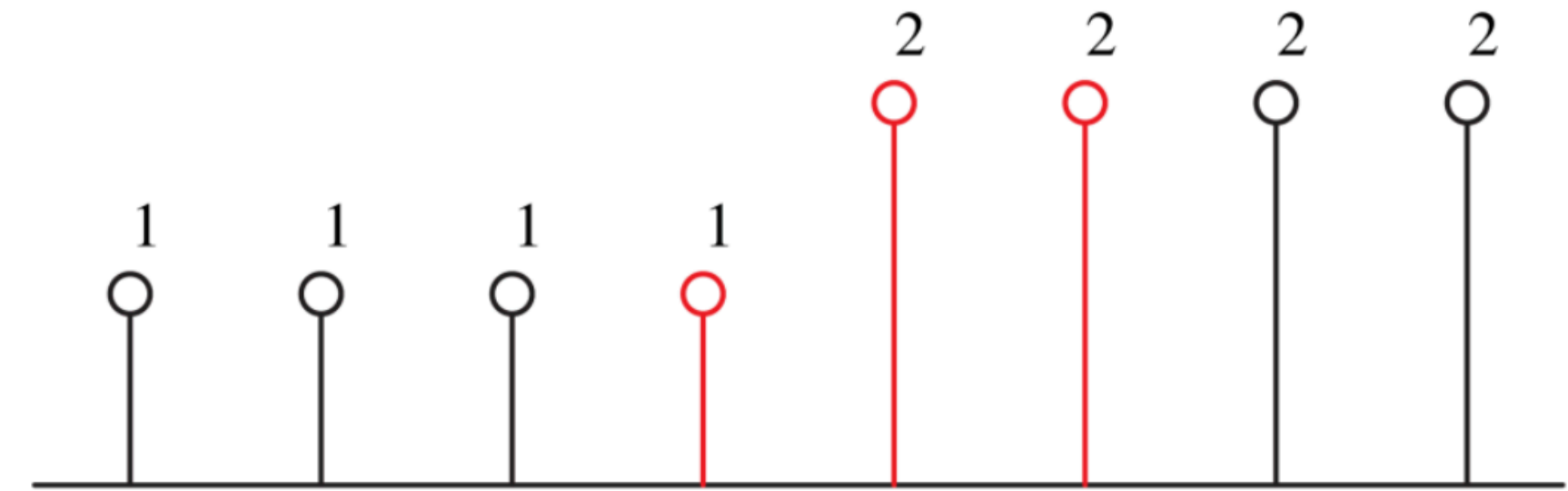




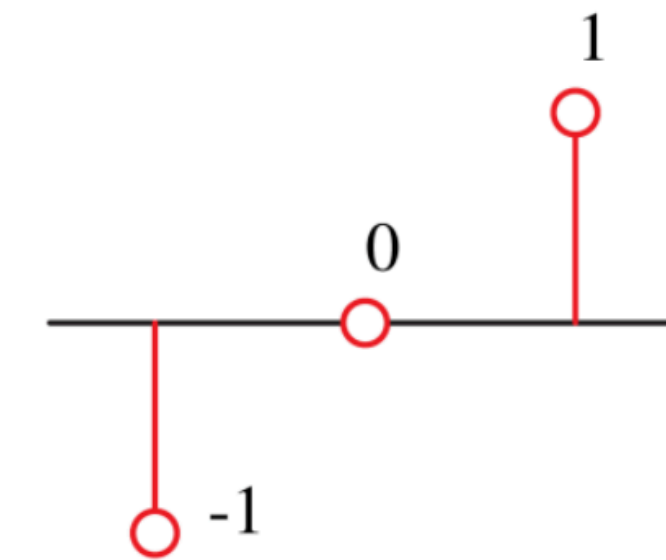
# Convolution: a 1-D example

- sliding window
- dot product

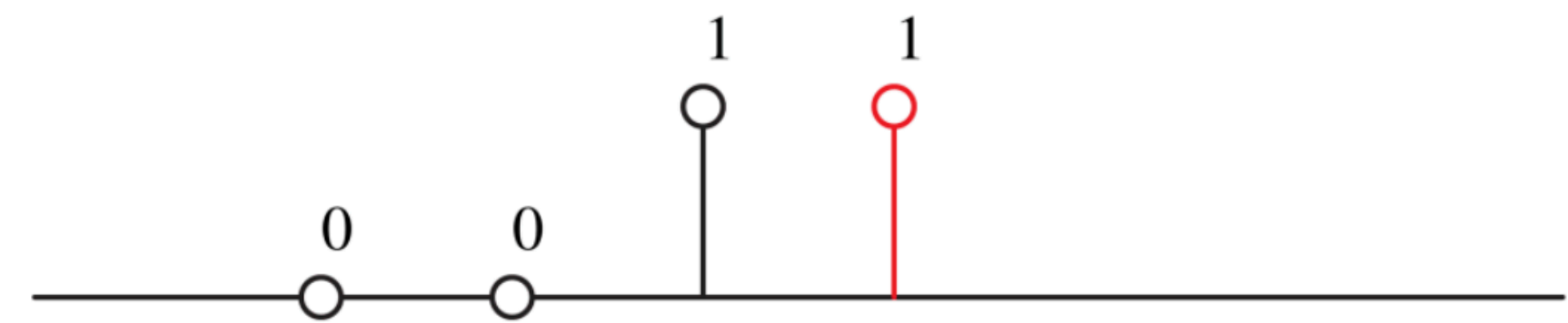
input



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

filter

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

output


# Convolution: a 2-D example

input

<sup>0</sup> 1	<sup>0</sup> 2	<sup>0</sup> 1	0	0	0	0	0
<sup>0</sup> 0	<sup>0</sup> 0	<sup>0</sup> 0	0	0	1	1	0
<sup>0</sup> -1	<sup>1</sup> -2	<sup>1</sup> -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

filter

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

output

-3					

- sliding window
- dot product

# Convolution: a 2-D example

input

0	<sup>0</sup> 1	<sup>0</sup> 2	<sup>0</sup> 1	0	0	0	0
0	<sup>0</sup> 0	<sup>0</sup> 0	<sup>0</sup> 0	0	1	1	0
0	<sup>1</sup> -1	<sup>1</sup> -2	<sup>1</sup> -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

filter

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

output

-3	-4				

- sliding window
- dot product



# Convolution: a 2-D example

input

0	0	<sup>0</sup> 1	<sup>0</sup> 2	<sup>0</sup> 1	0	0	0
0	0	<sup>0</sup> 0	<sup>0</sup> 0	<sup>0</sup> 0	1	1	0
0	1	<sup>1</sup> -1	<sup>1</sup> -2	<sup>1</sup> -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			

- sliding window
- dot product

# 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	<sup>0</sup> 1	<sup>0</sup> 2	<sup>0</sup> 1
0	0	1	1	1	<sup>0</sup> 0	<sup>0</sup> 0	<sup>0</sup> 0
0	0	0	0	0	<sup>0</sup> -1	<sup>0</sup> -2	<sup>0</sup> -1

filter

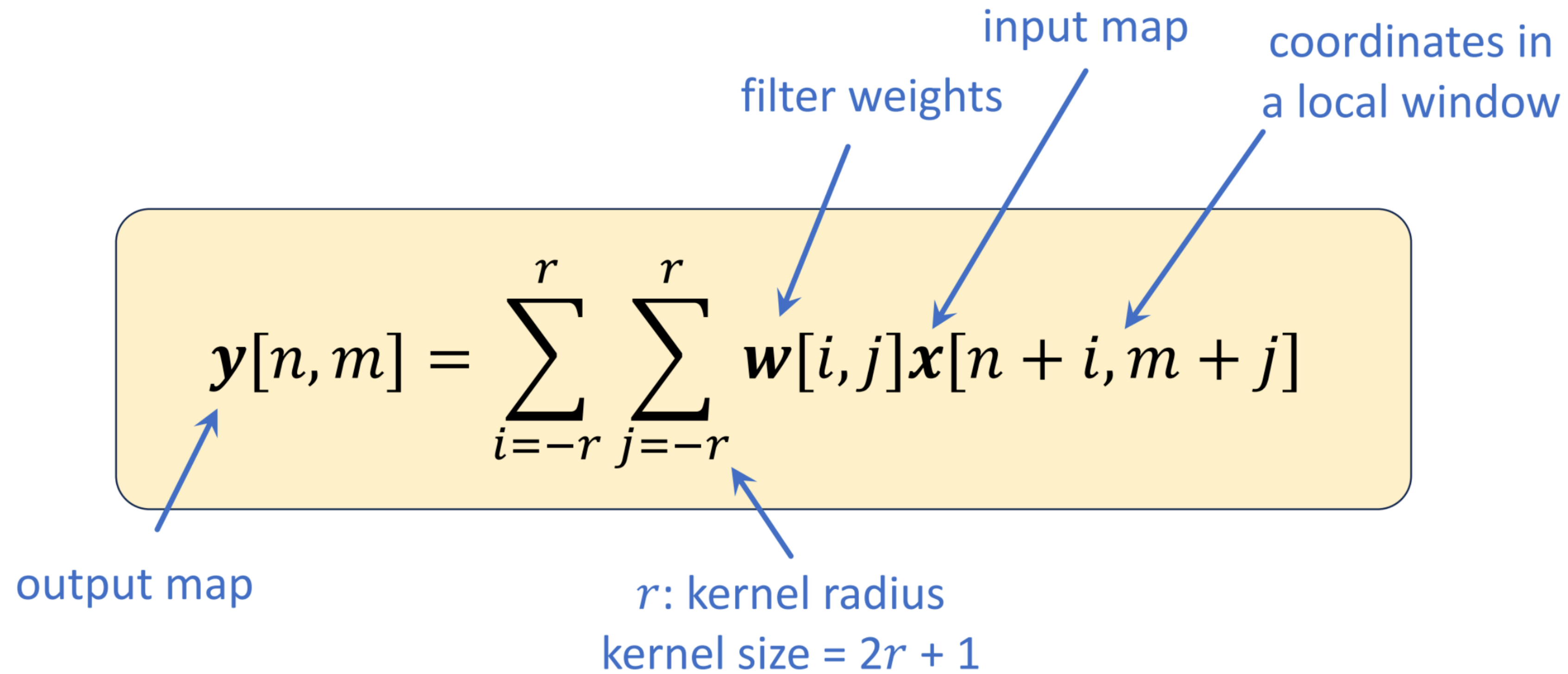
1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

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: a 2-D example

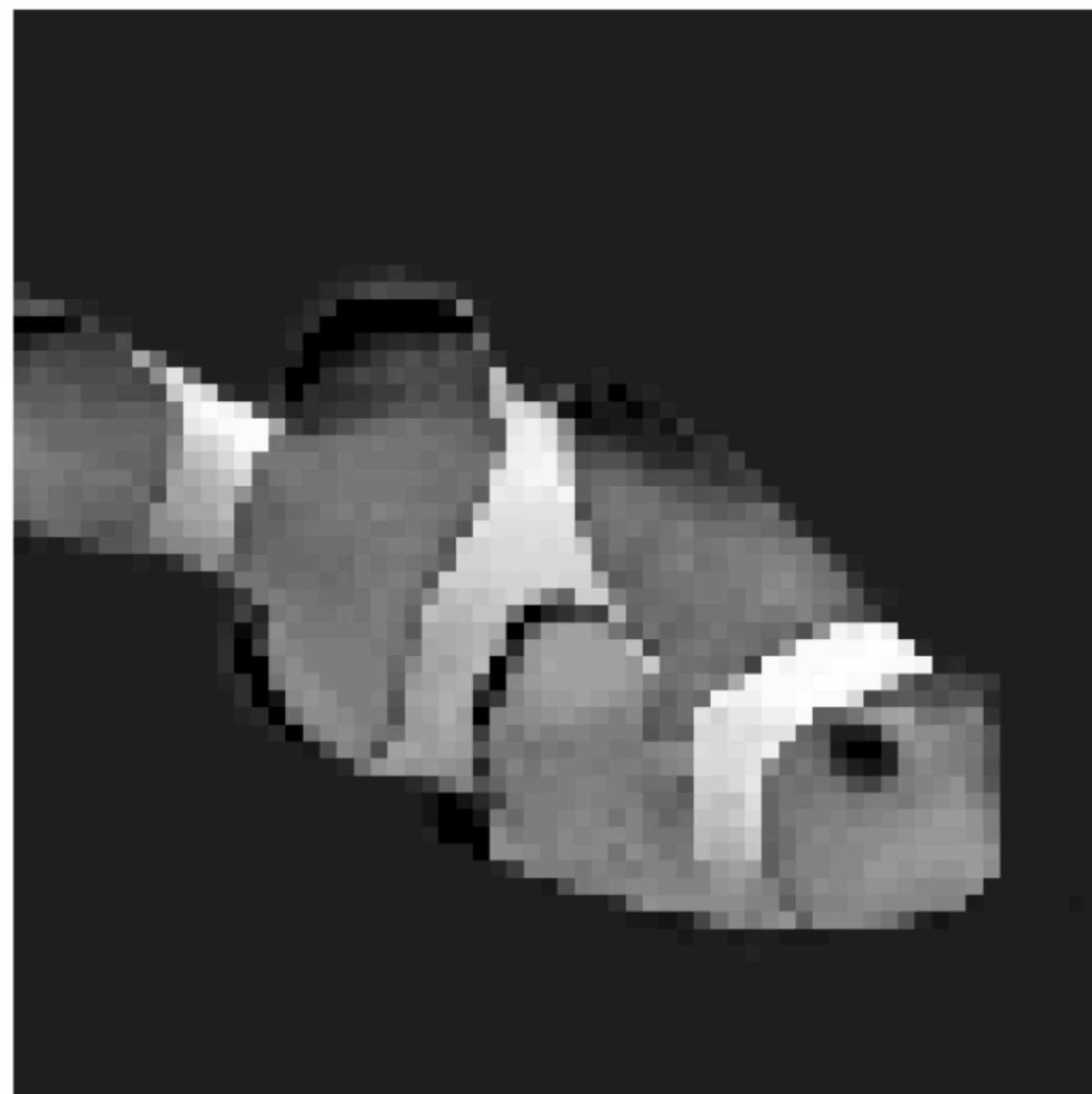


The diagram illustrates the 2-D convolution equation with several labels and arrows pointing to specific parts of the formula:

- output map**: Points to the variable  $y[n, m]$ .
- filter weights**: Points to the variable  $w[i, j]$ .
- input map**: Points to the variable  $x[n + i, m + j]$ .
- coordinates in a local window**: Points to the indices  $i$  and  $j$  in the input map term.
- $r$ : kernel radius** and **kernel size =  $2r + 1$** : Points to the summation limits  $i = -r$  and  $j = -r$ .

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

# Convolution: 2-D



filter

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

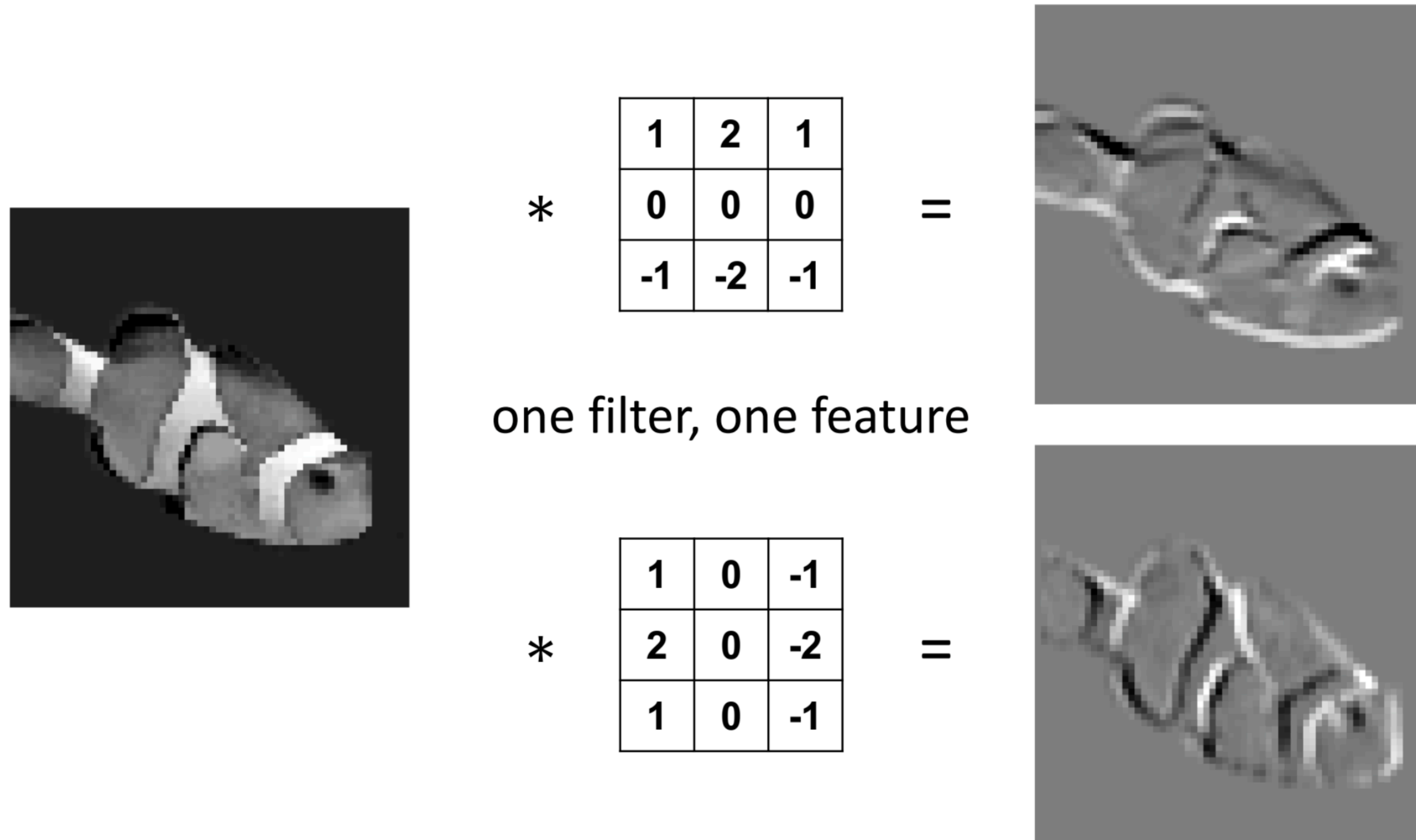
\*

=



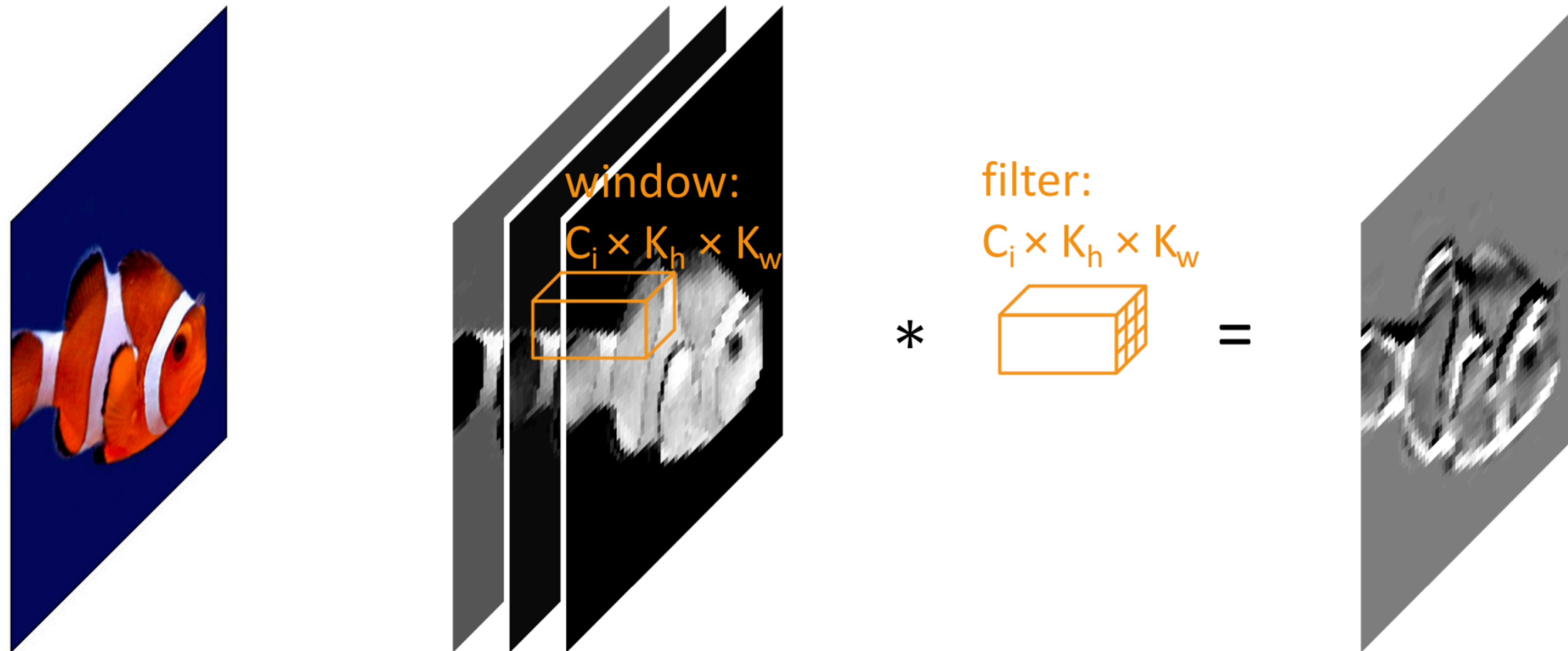


# Convolution: Multi-channel outputs



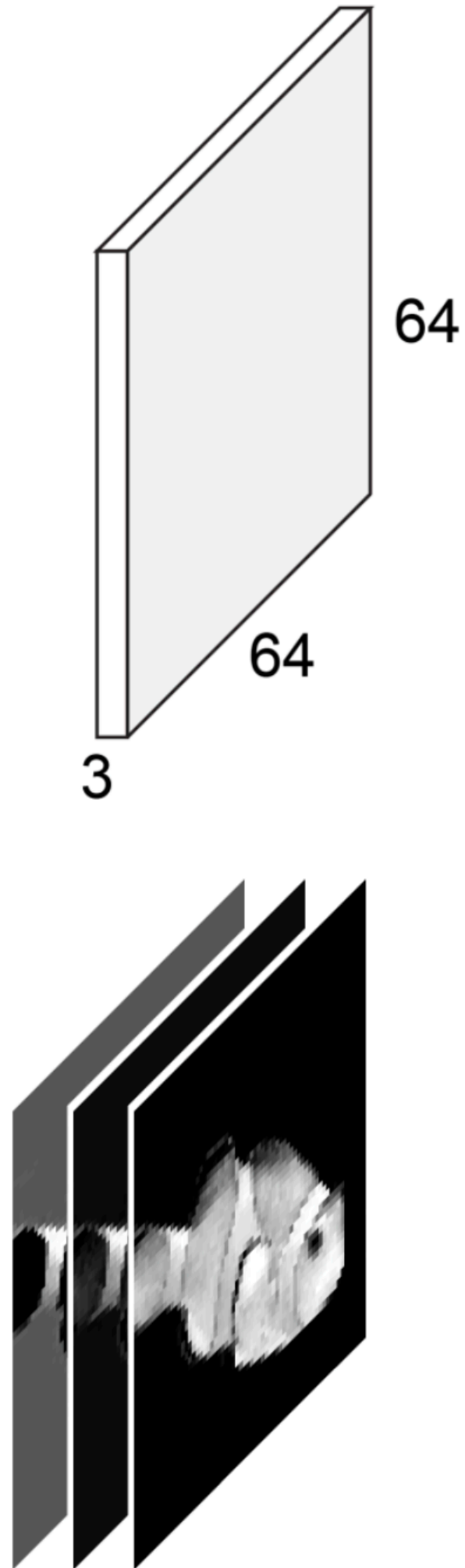
Understanding the filter/kernel as feature extractors

# Convolution: Multi-channel inputs



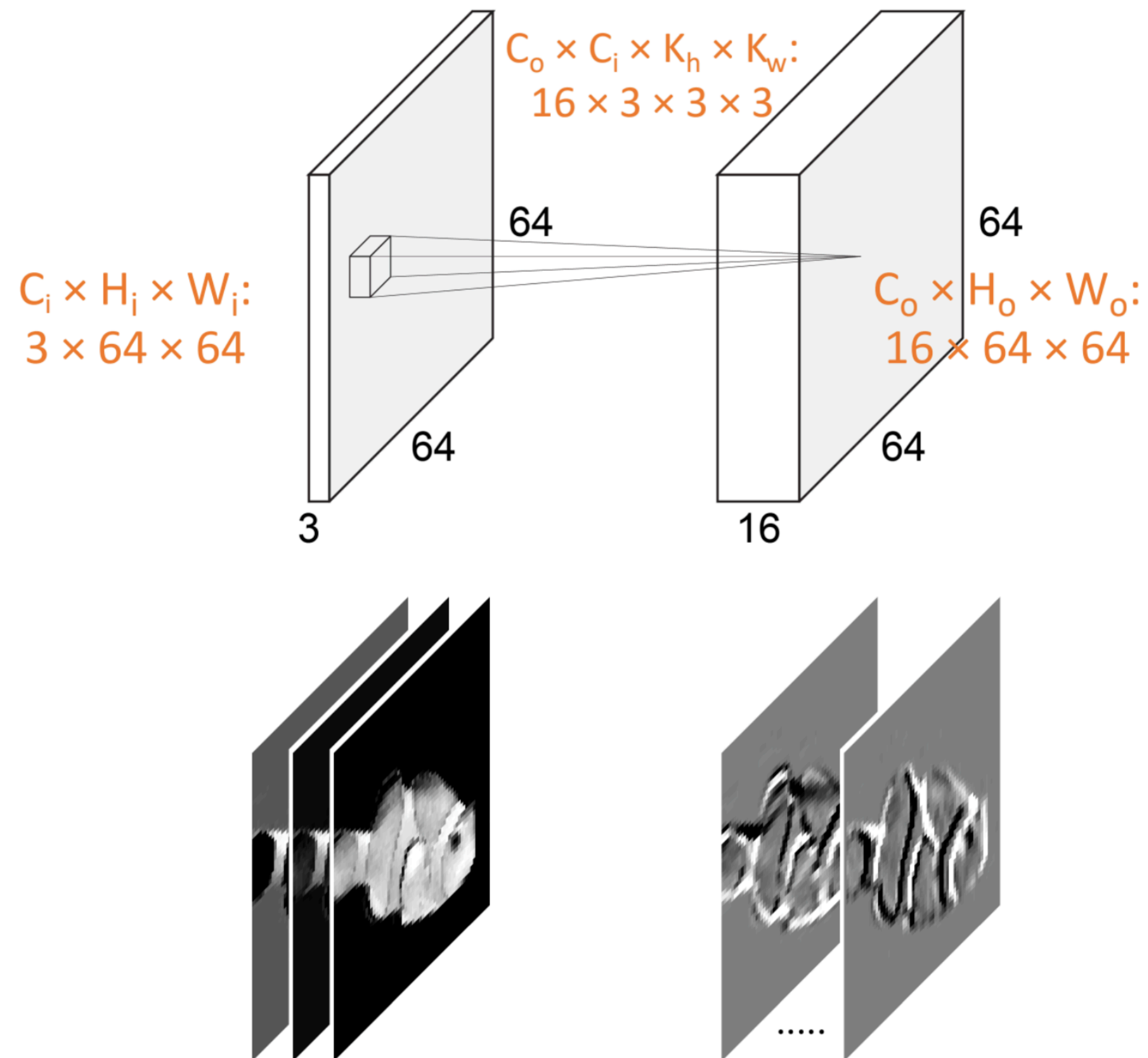
Like (R, G, B) color notations have three features

# 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

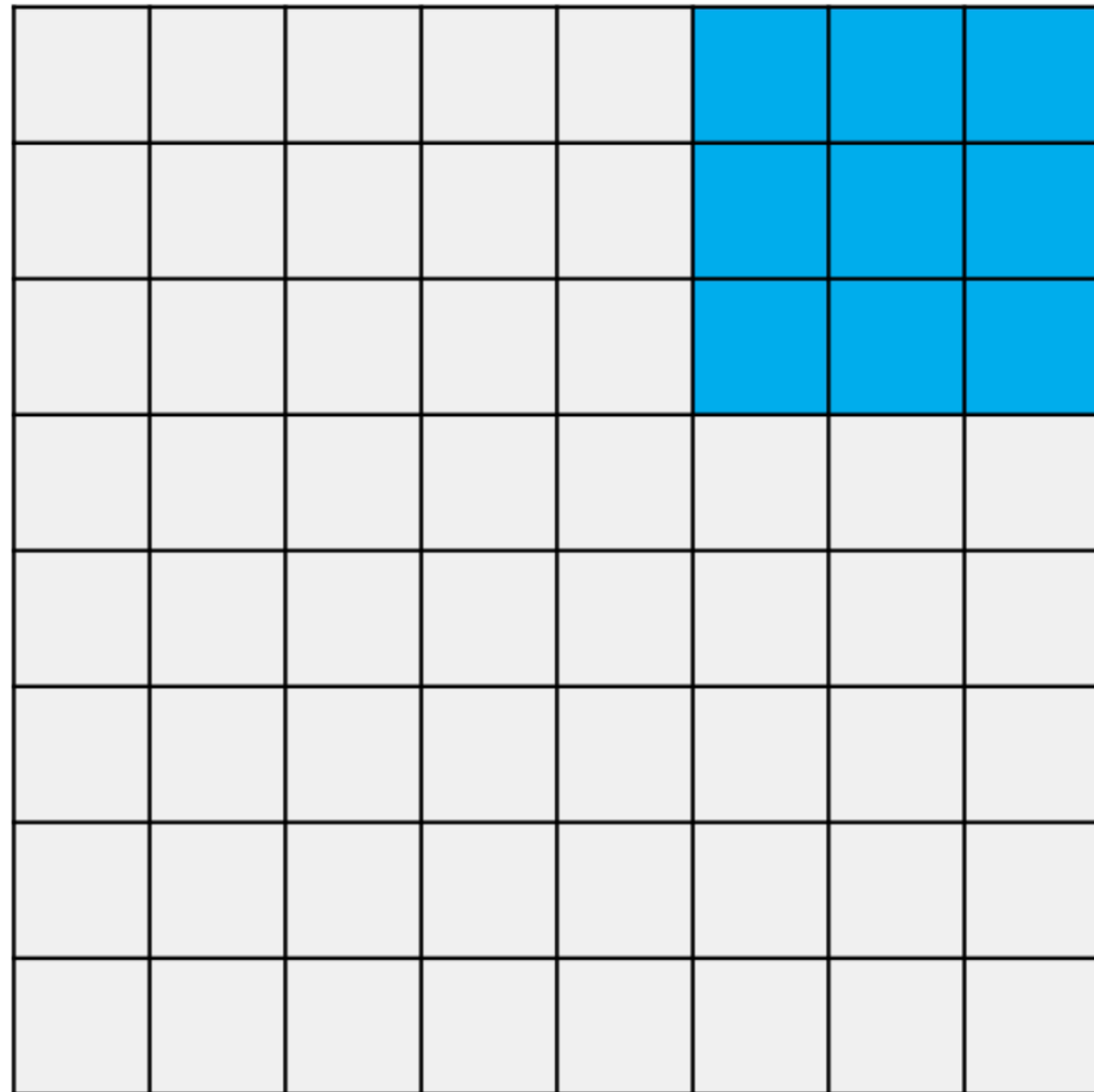


# Convolution: # parameters and # operations

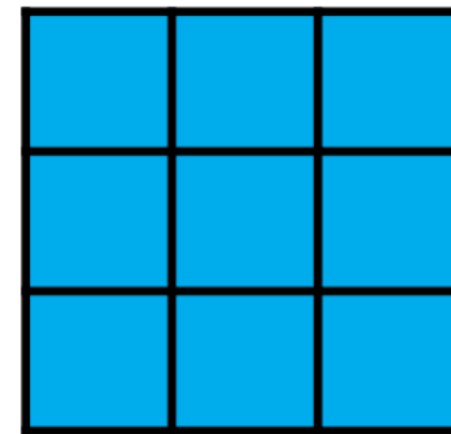
- # parameters
  - weights:  $C_o \times C_i \times K_h \times K_w$
  - bias:  $C_o$
- # floating-point operations (FLOPs)
  - # params  $\times H_o \times W_o$

# Convolution: padding

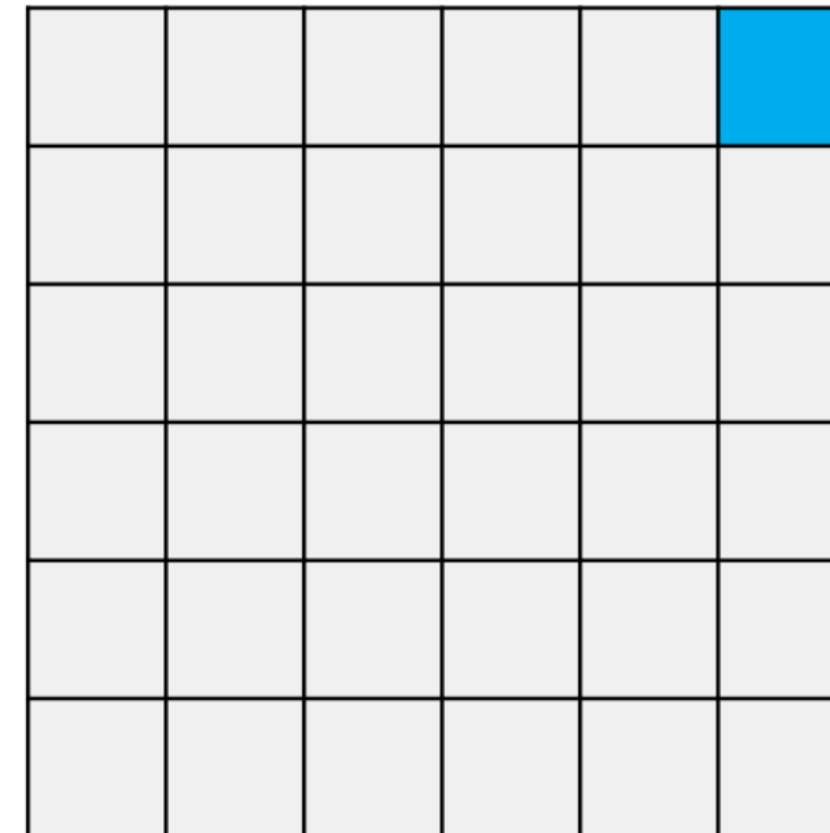
input:  $H \times W = 8 \times 8$



filter



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



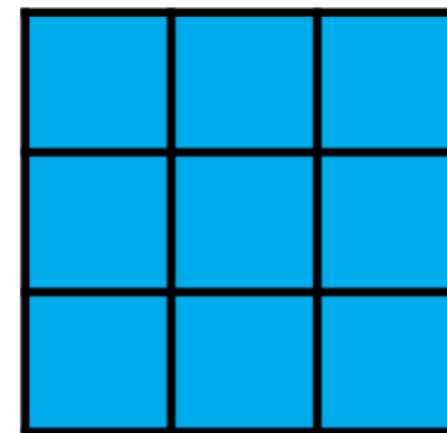
$$H_{\text{out}} = H_{\text{in}} - K_h + 1$$

# 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	0	0	0	0

filter

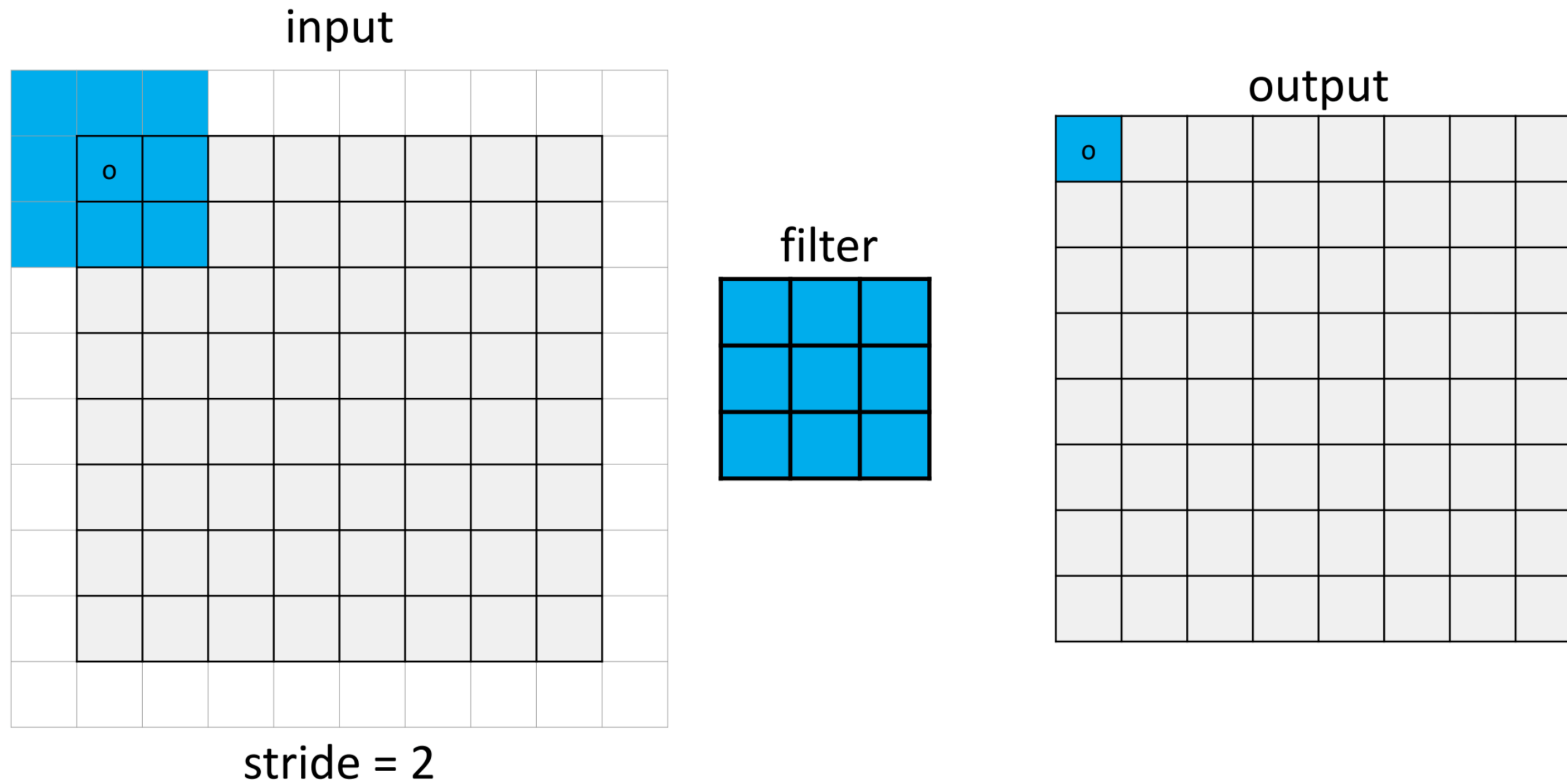


output: H × W = 8 × 8

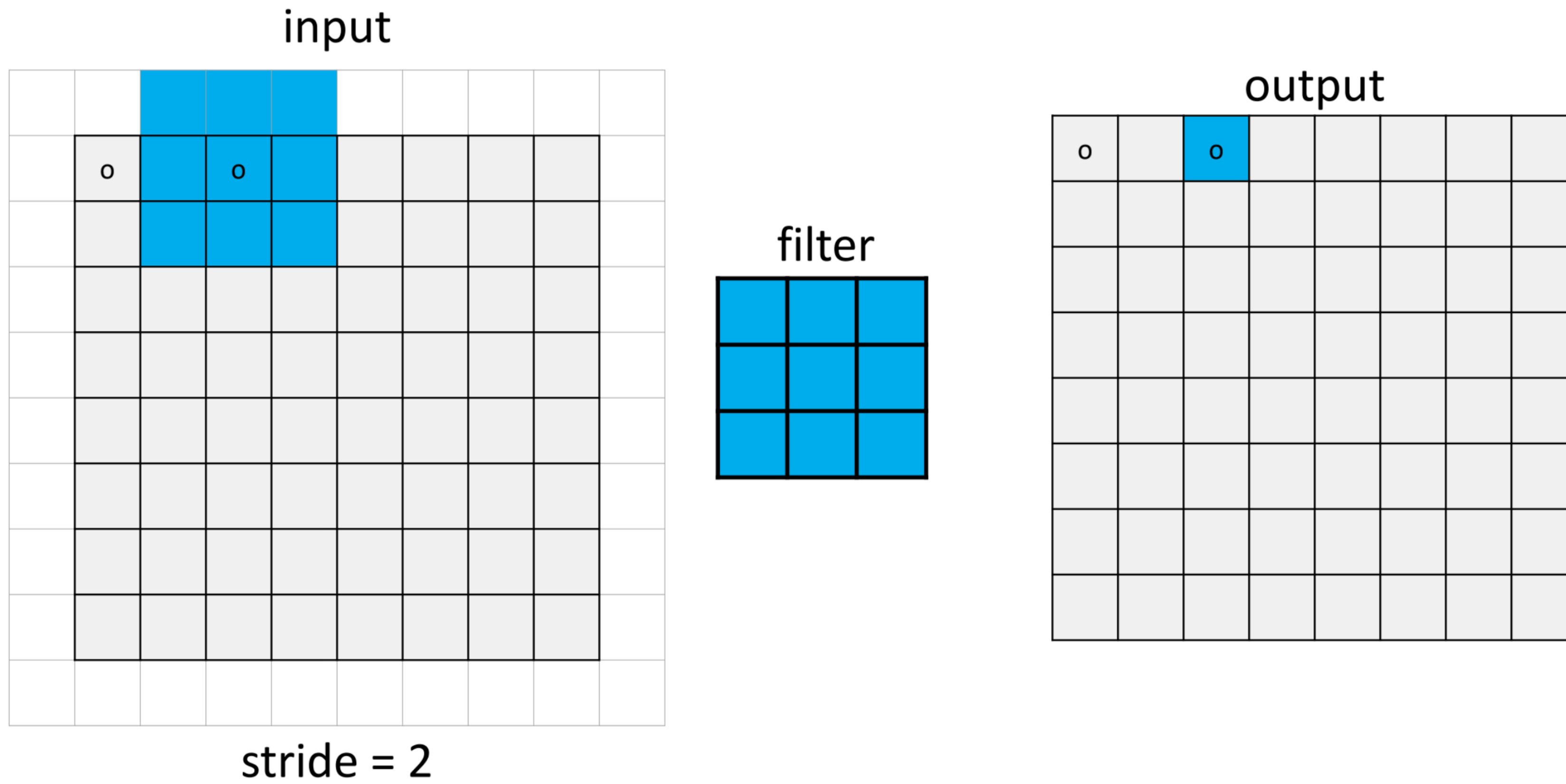

- $\text{pad} = \lfloor \text{kernel\_size} / 2 \rfloor$
- maintains feature map size

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

# Convolution: stride

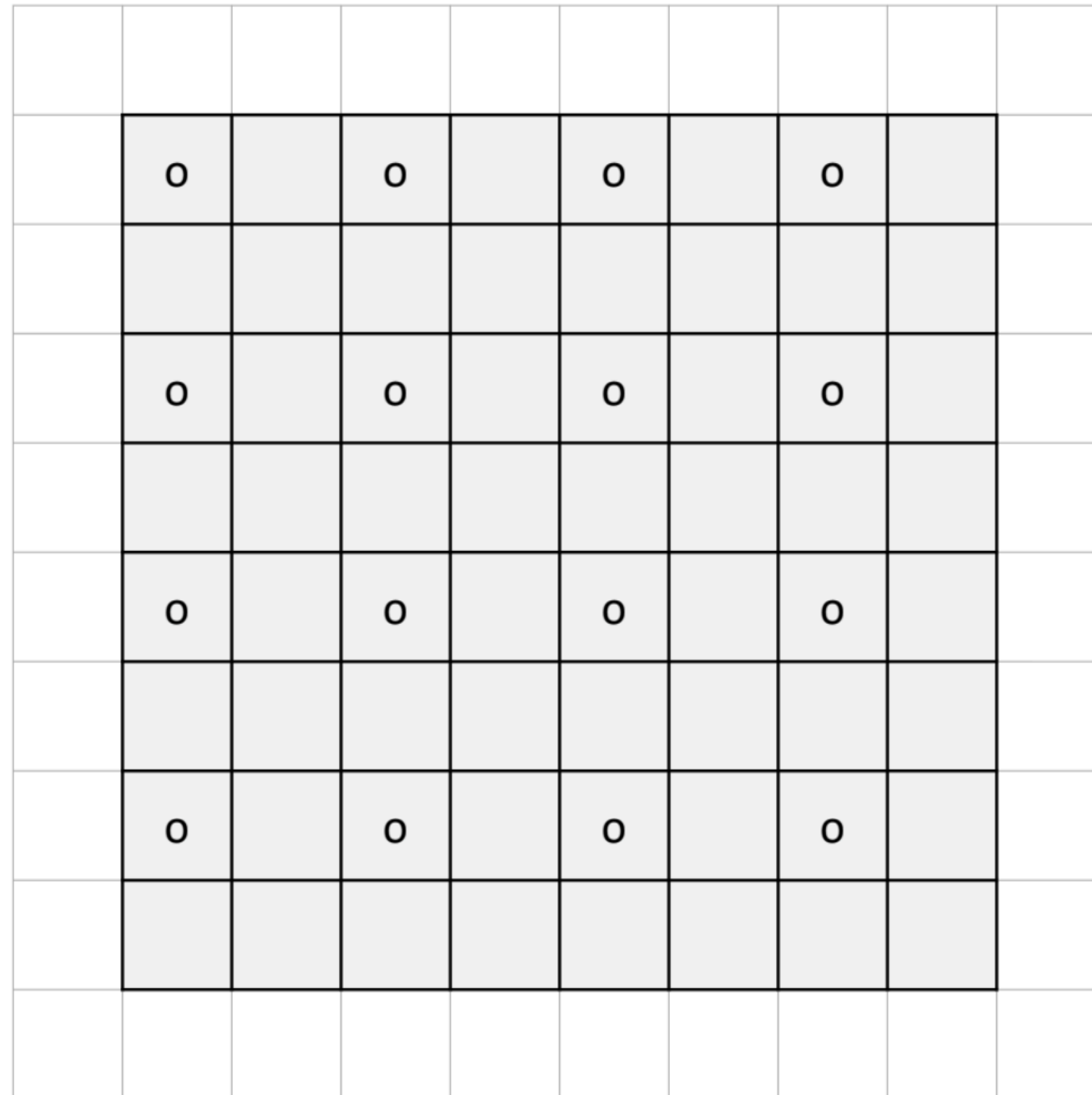


# Convolution: stride



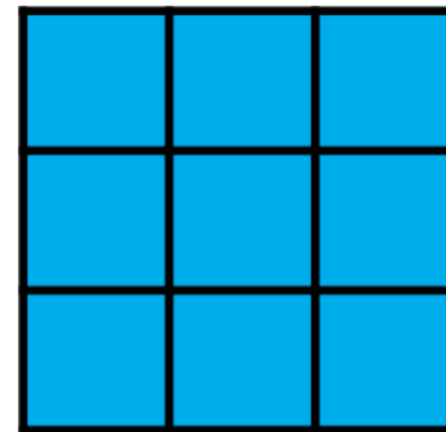
# Convolution: stride

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

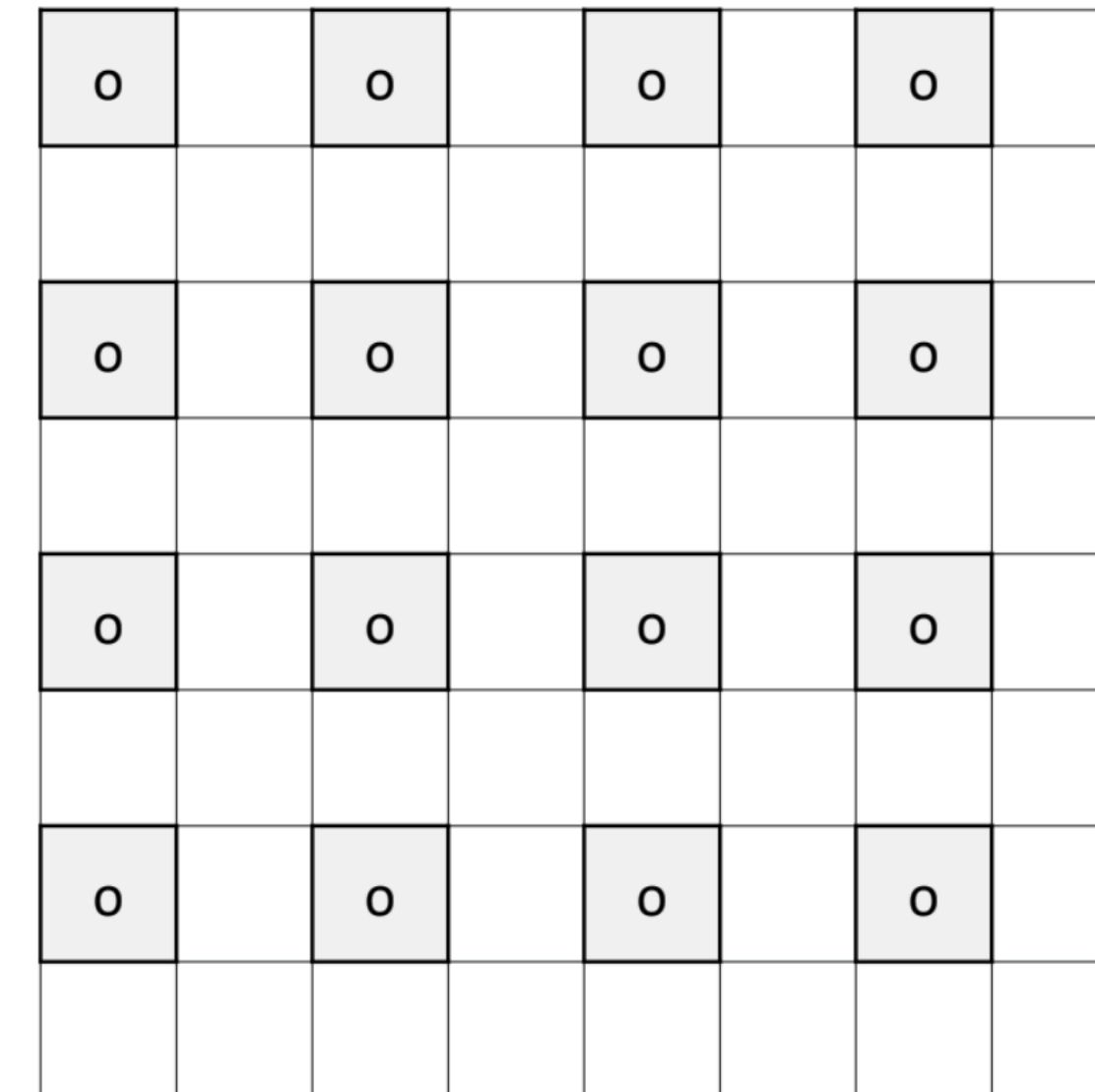


stride = 2

filter



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

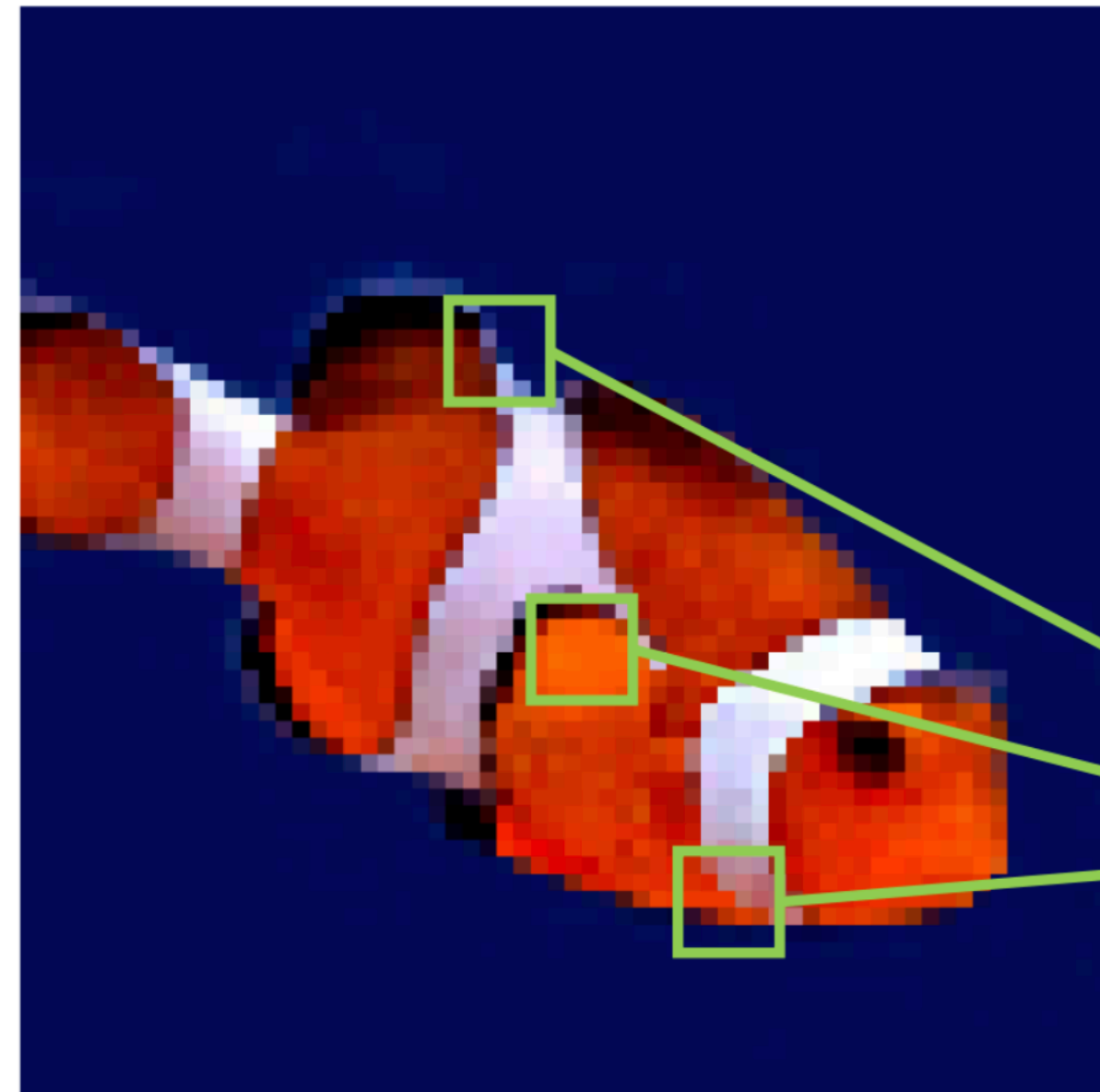


- reduces feature map size
- compress and abstract

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

# Convolution: translation-invariance

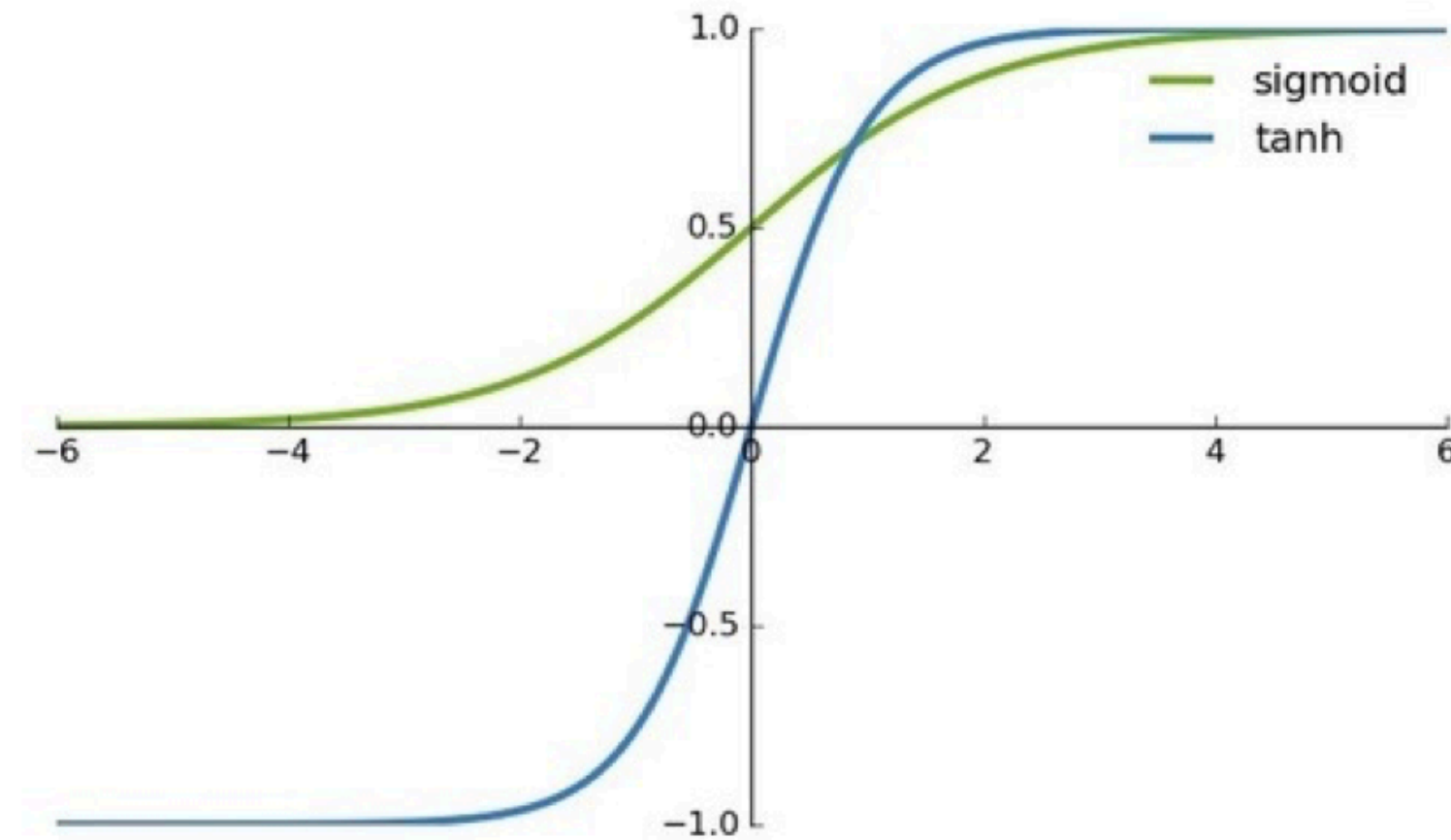
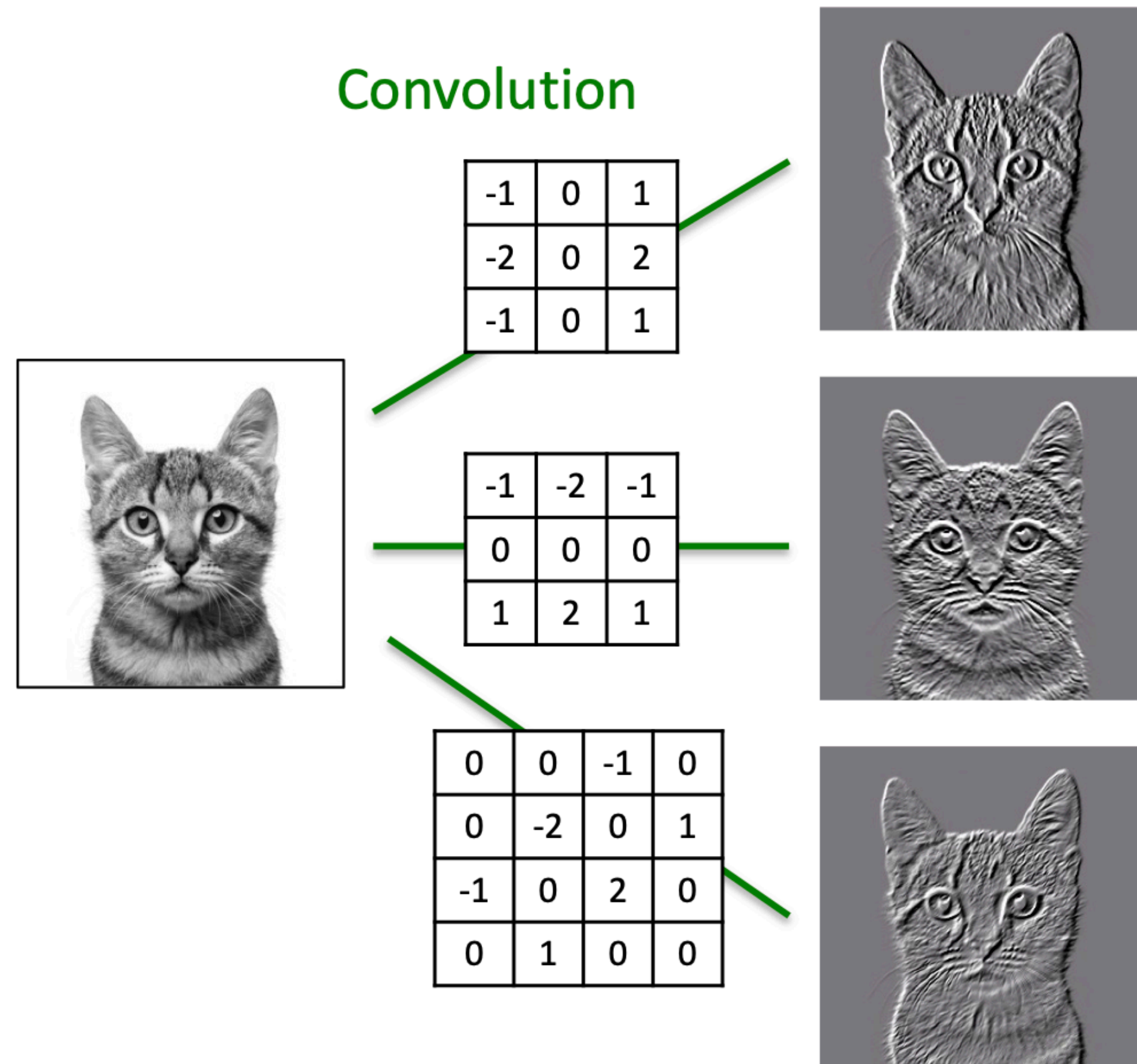
- Process each window in the same way



apply the same weights regardless of the window location

# Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling



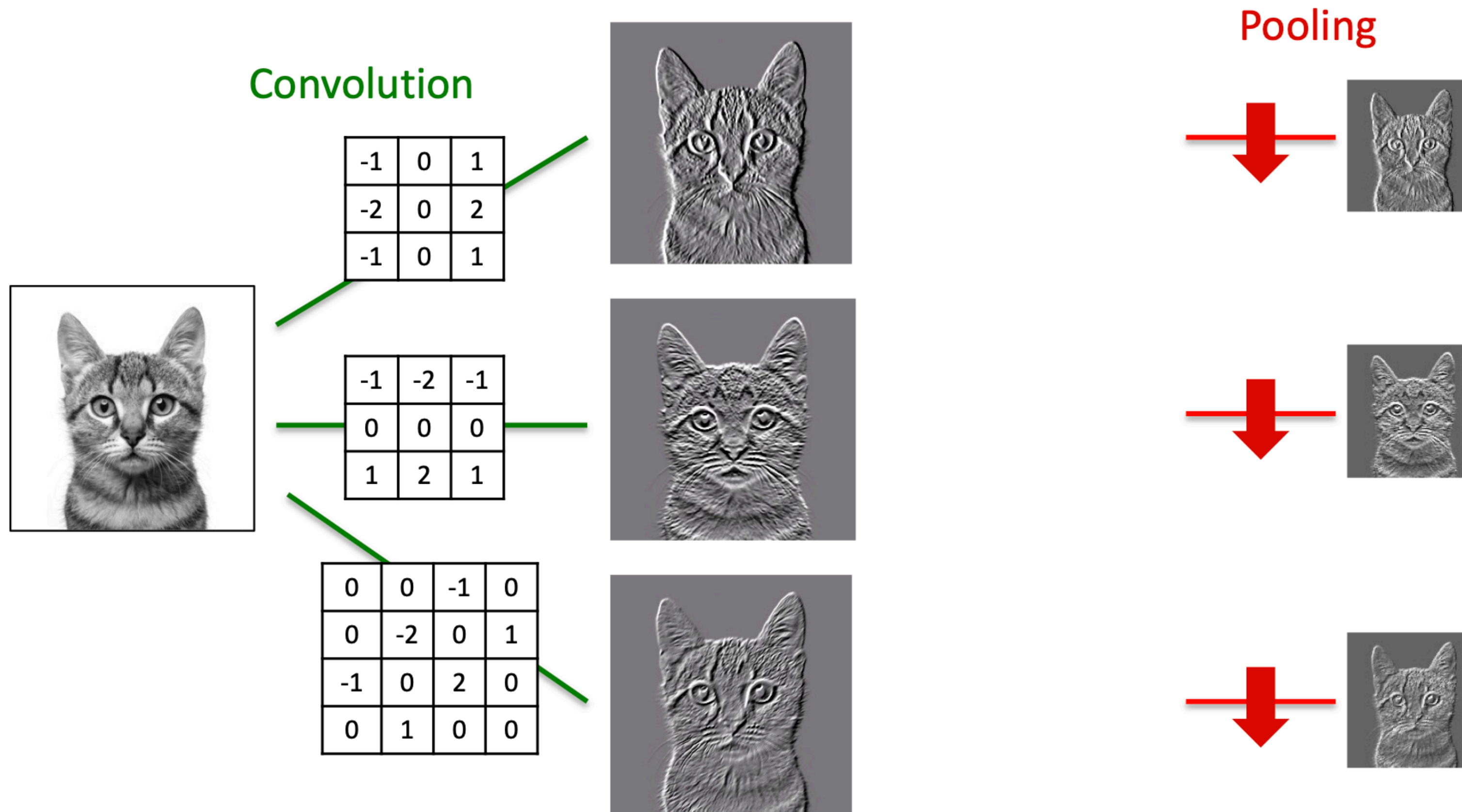
LeNet – tanh activation

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



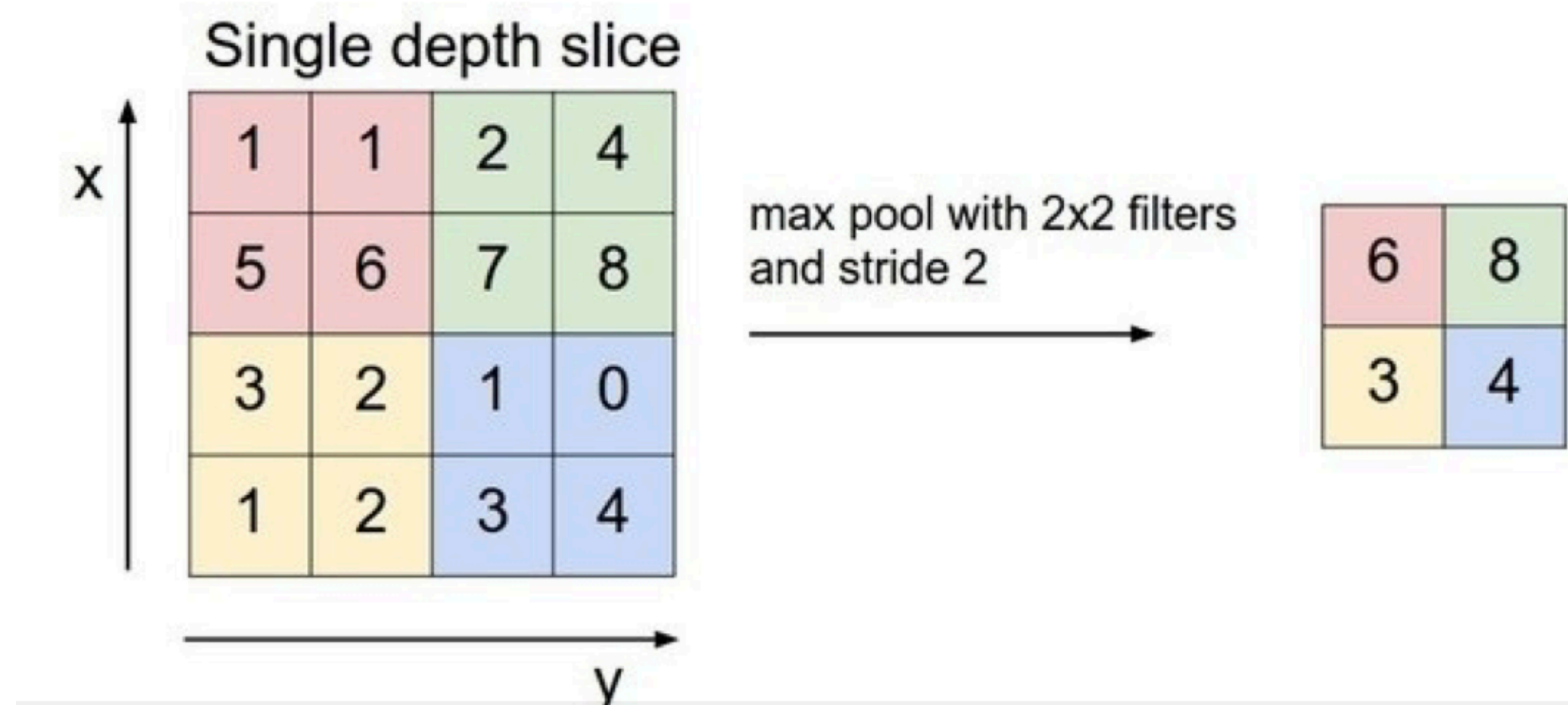
# Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling

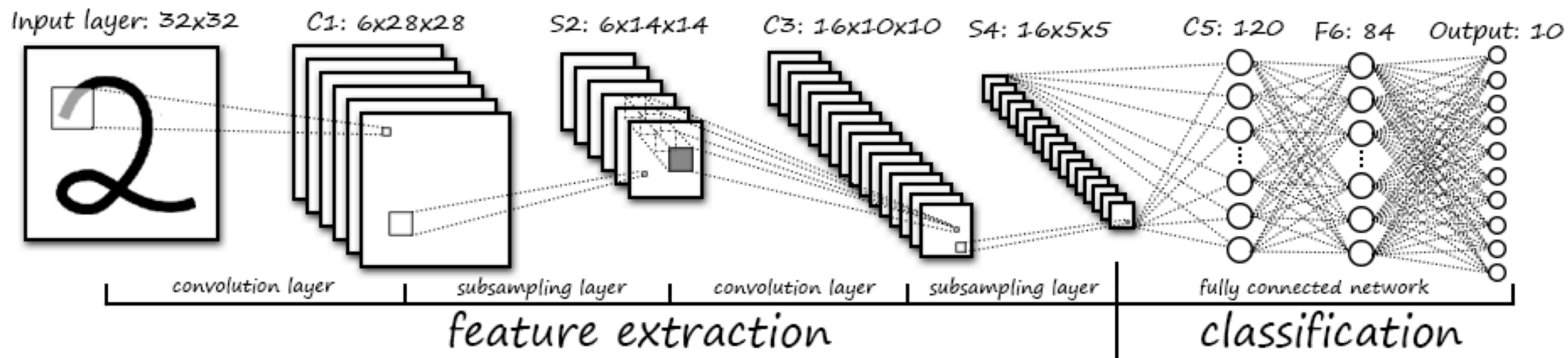


# Pooling = Down-sampling

## 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  
2 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 7 6 9 8 6 1

60,000 original dataset

Test error: 0.95%

MNIST

[1] LeNet 5, LeCun et al. 1998

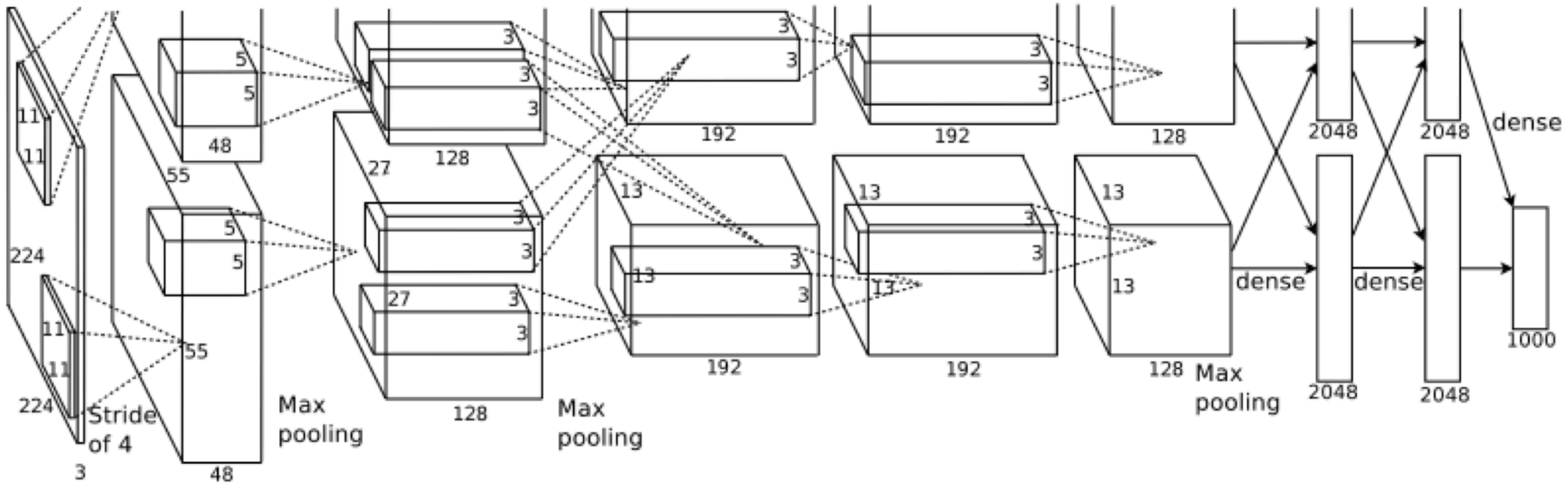


# Misclassified examples on MNIST

True label -> Predicted label

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

# Alex Net



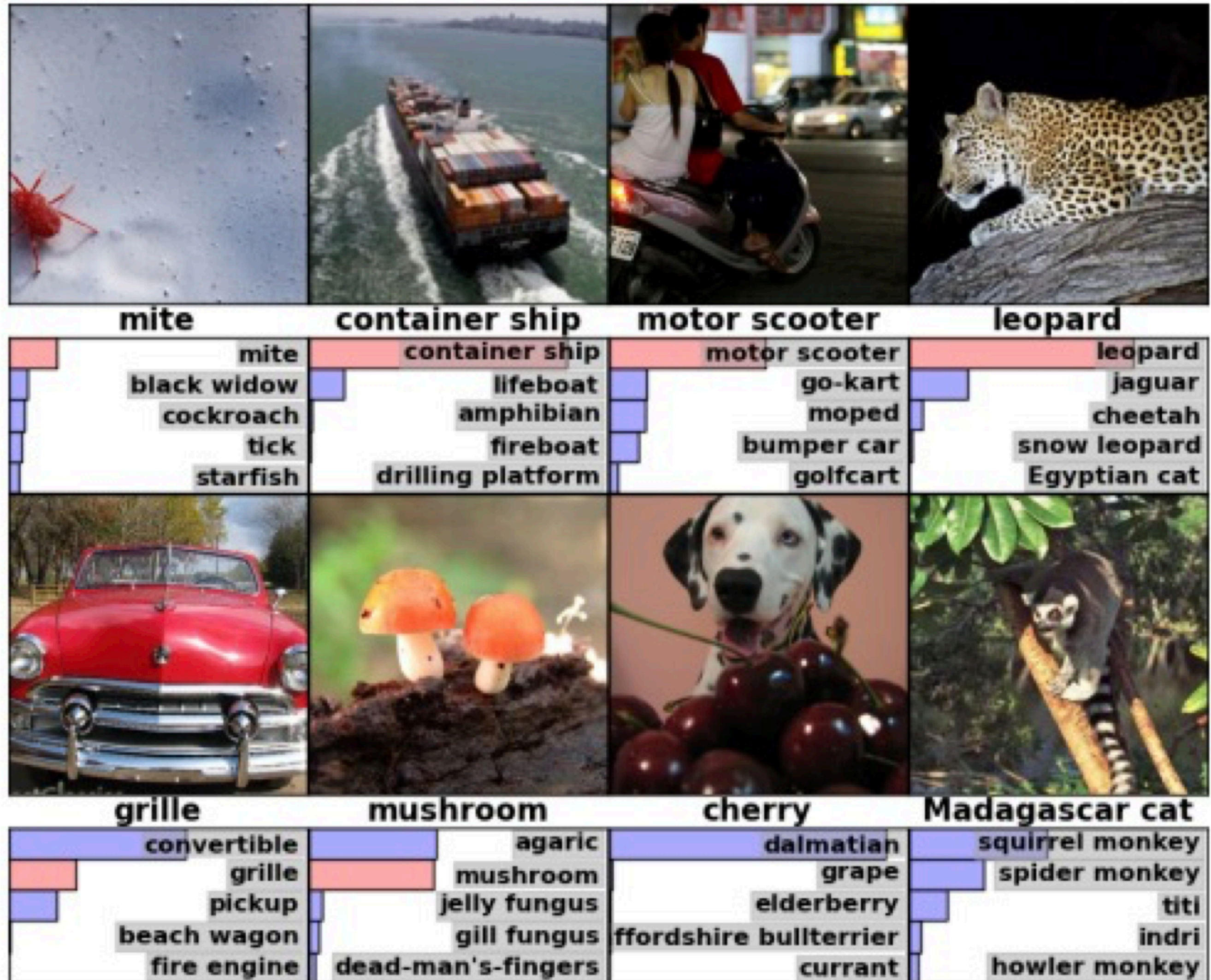
[1] Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NeurIPS 2012.

# 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

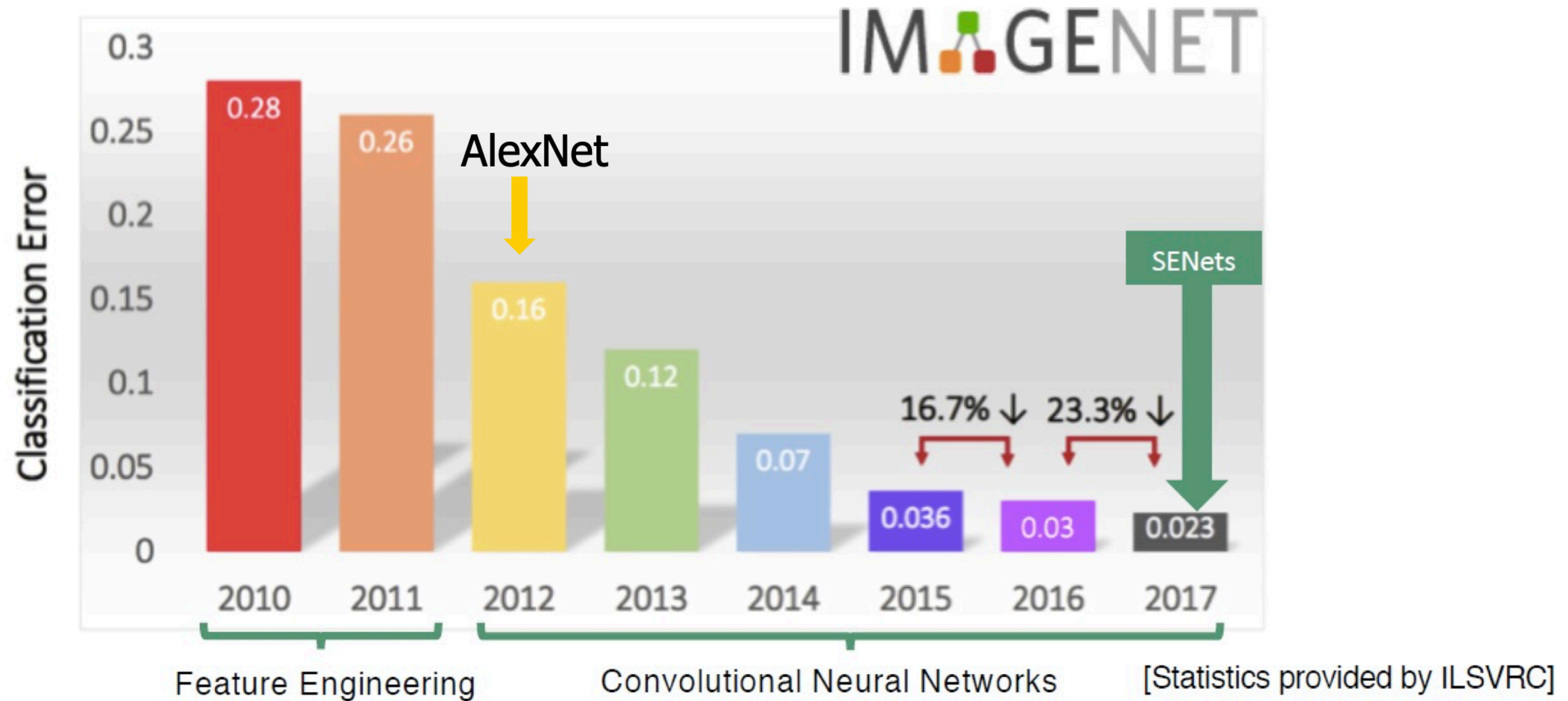


# ImageNet Results

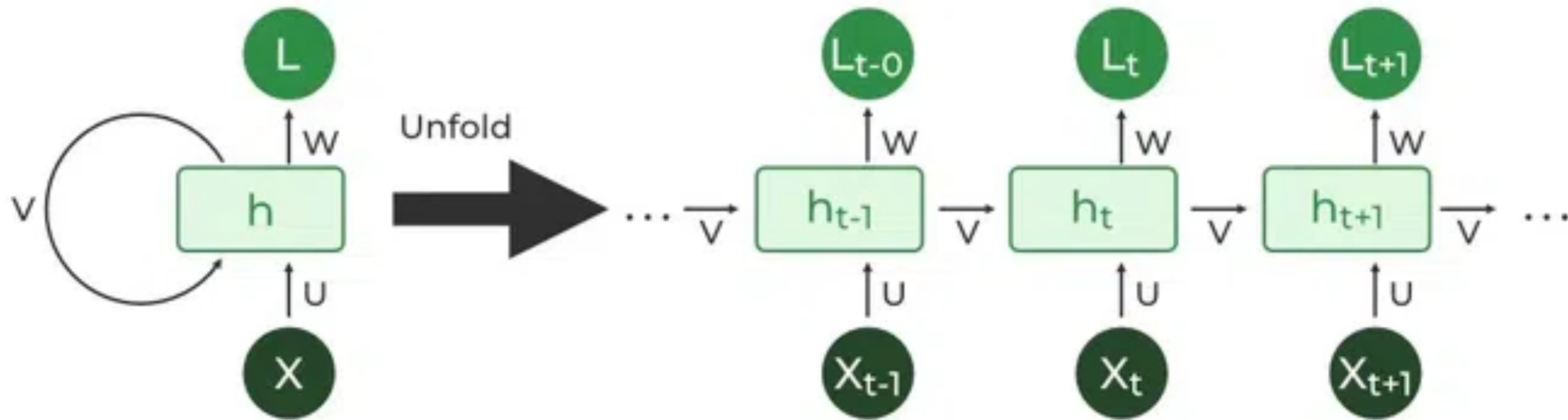




# ImageNet Results



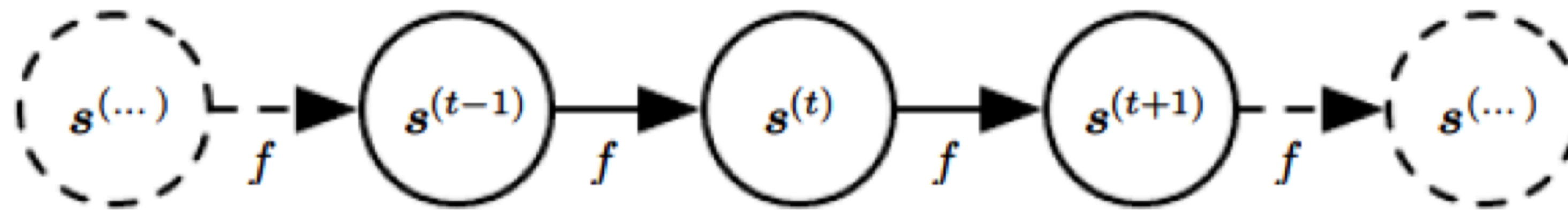
# Recurrent Neural Networks (RNNs)



# 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

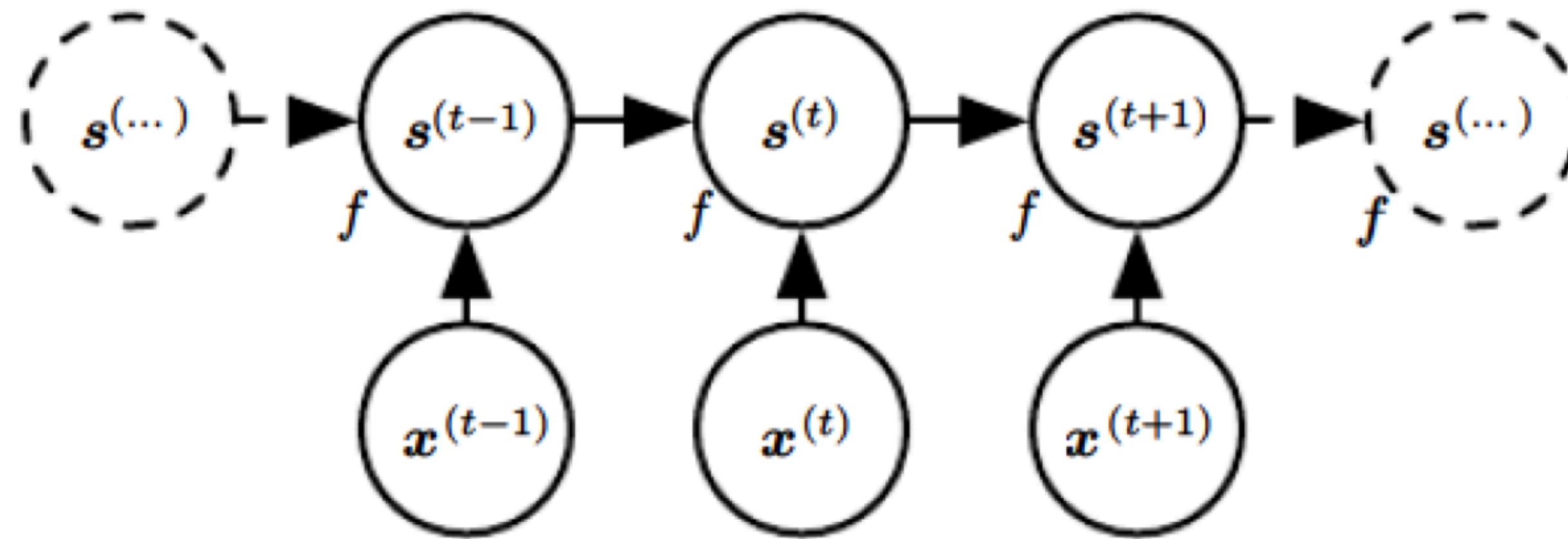


$$s^{(t+1)} = f(s^{(t)}; \theta)$$

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

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

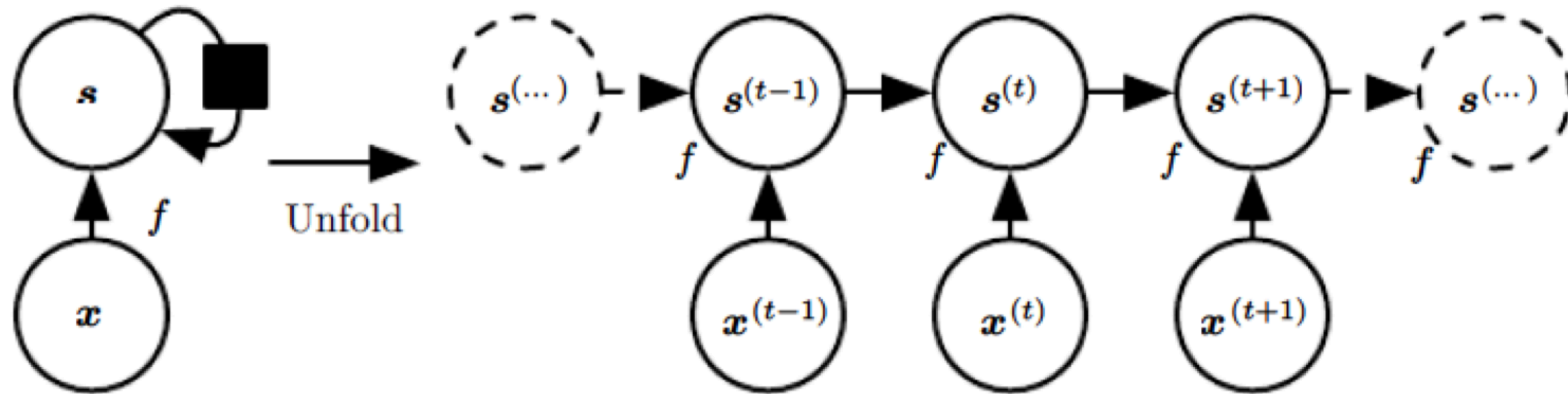
# Computation Graph



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



# Compact view



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

Key: the same  $f$  and  $\theta$   
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

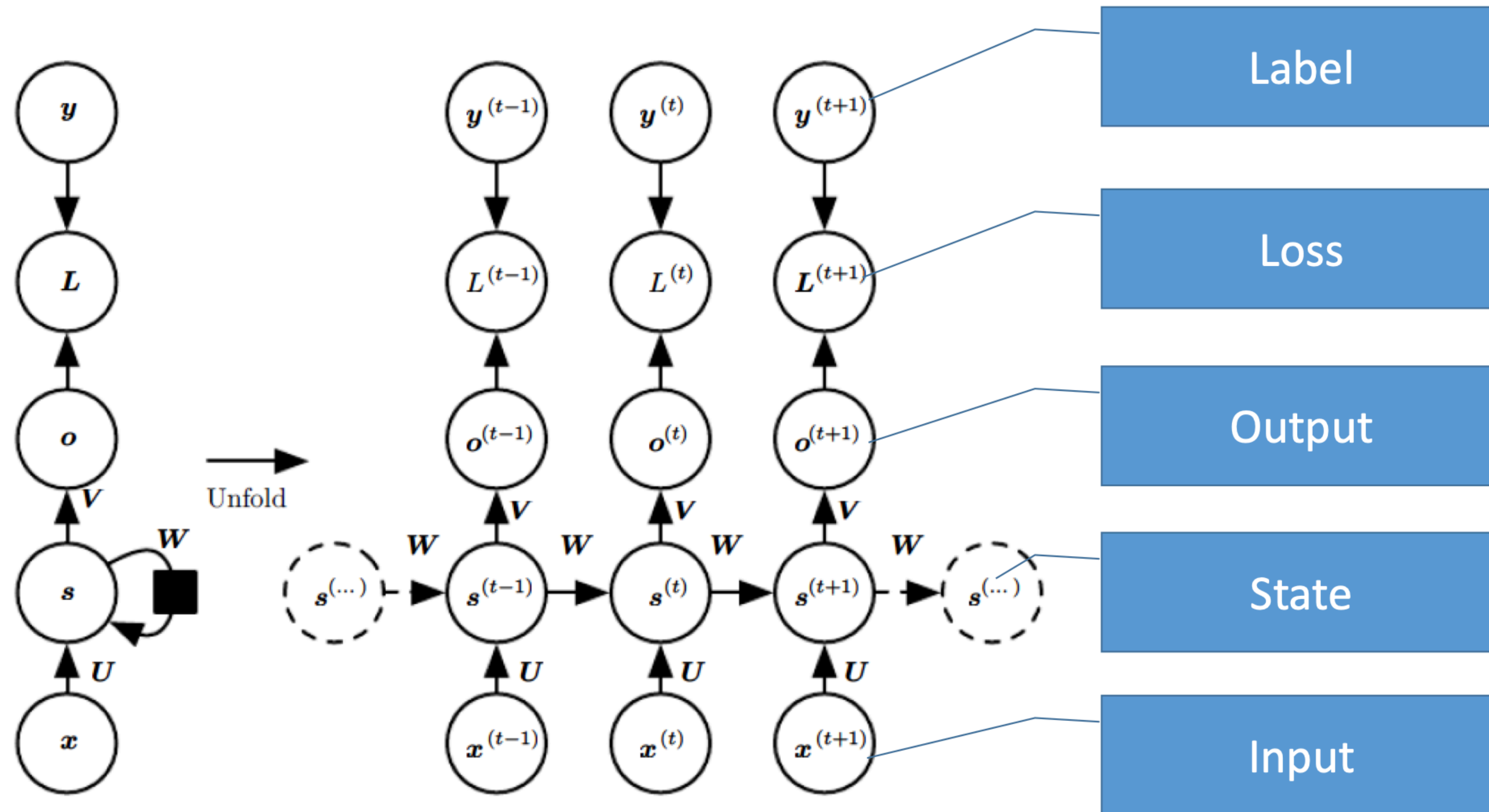


Figure from *Deep Learning*, by Goodfellow, Bengio and Courville

# Recurrent Neural Networks

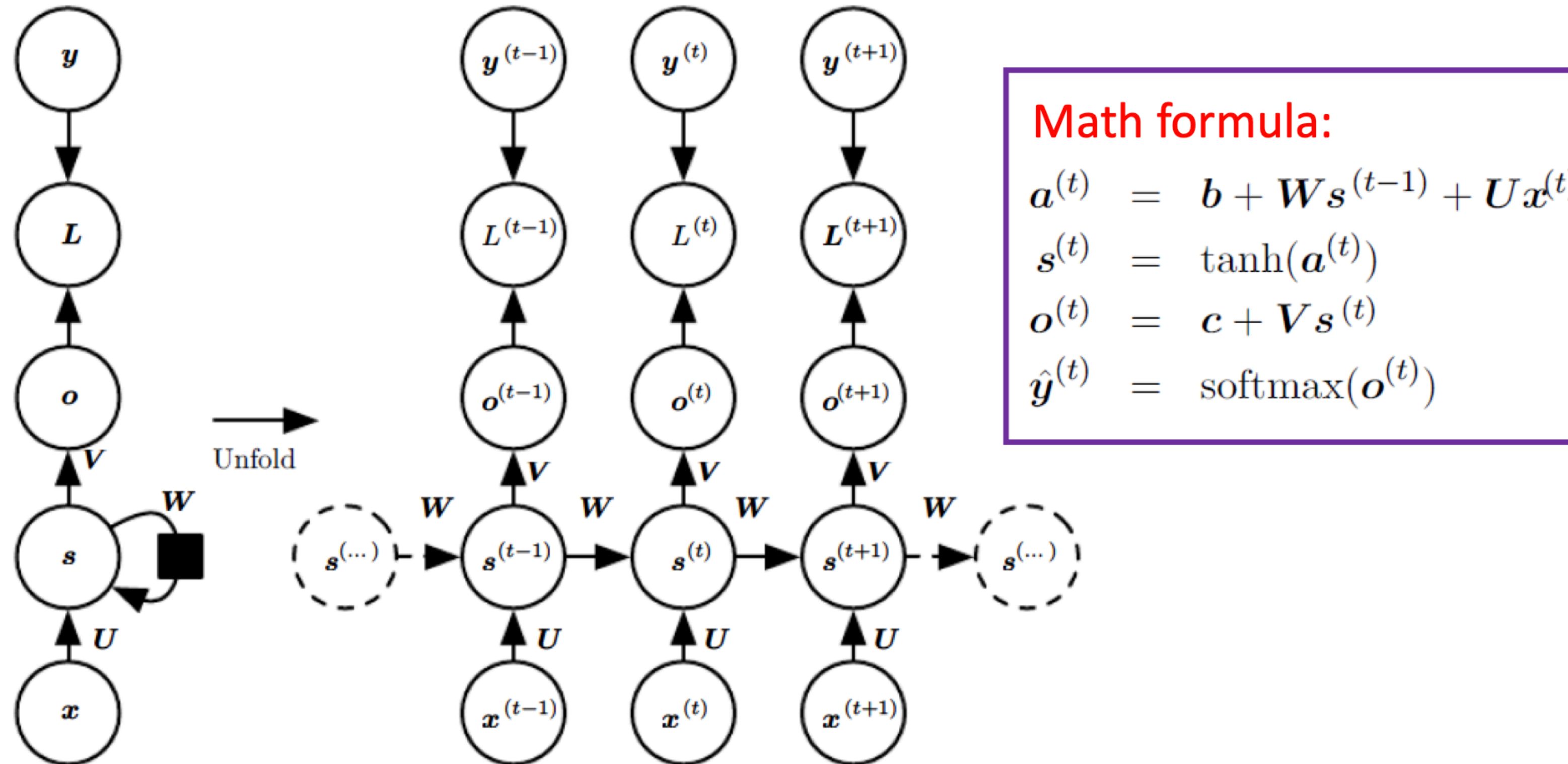


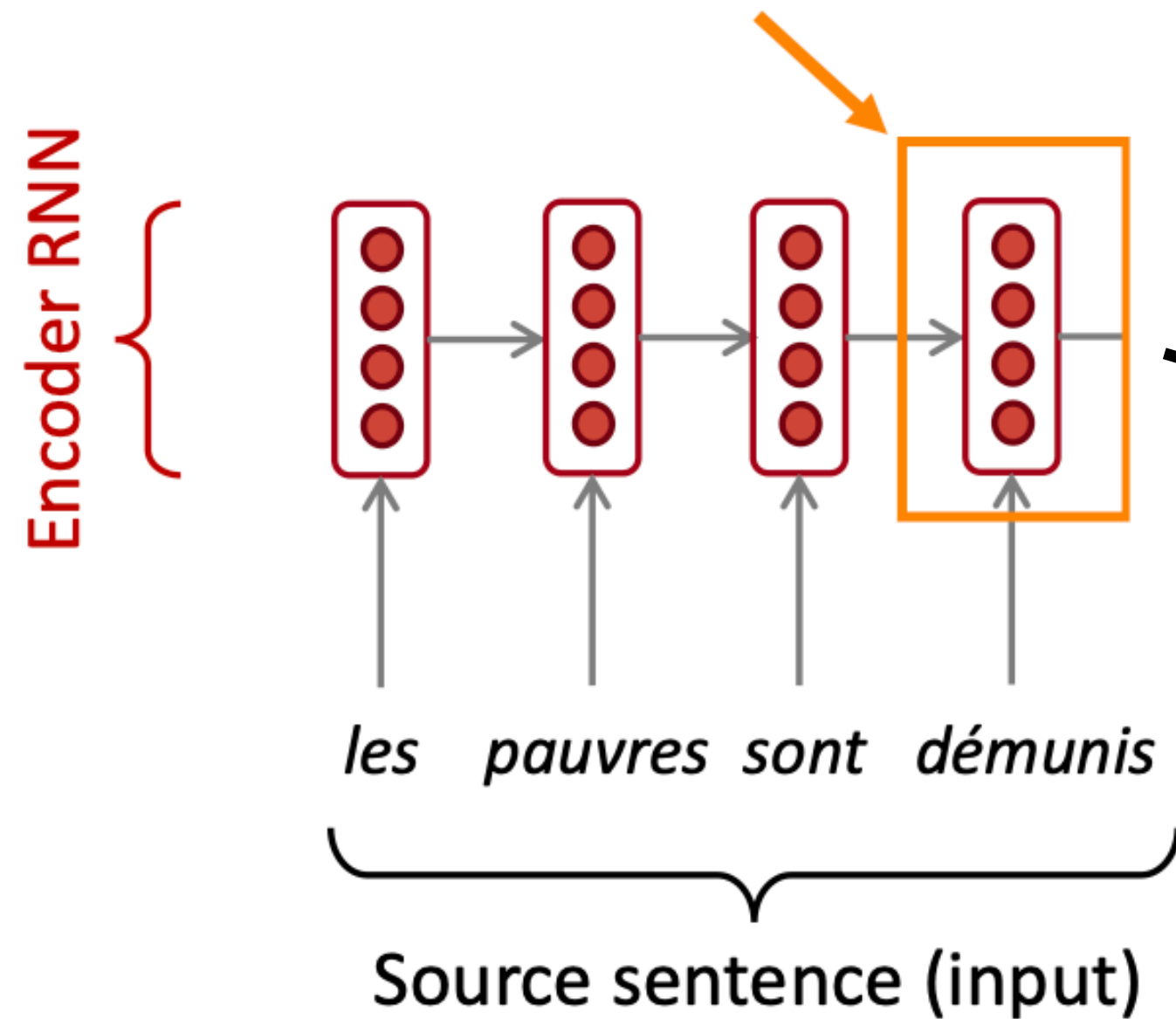
Figure from *Deep Learning*,  
Goodfellow, Bengio and Courville

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

# Sequence-to-Sequence Learning

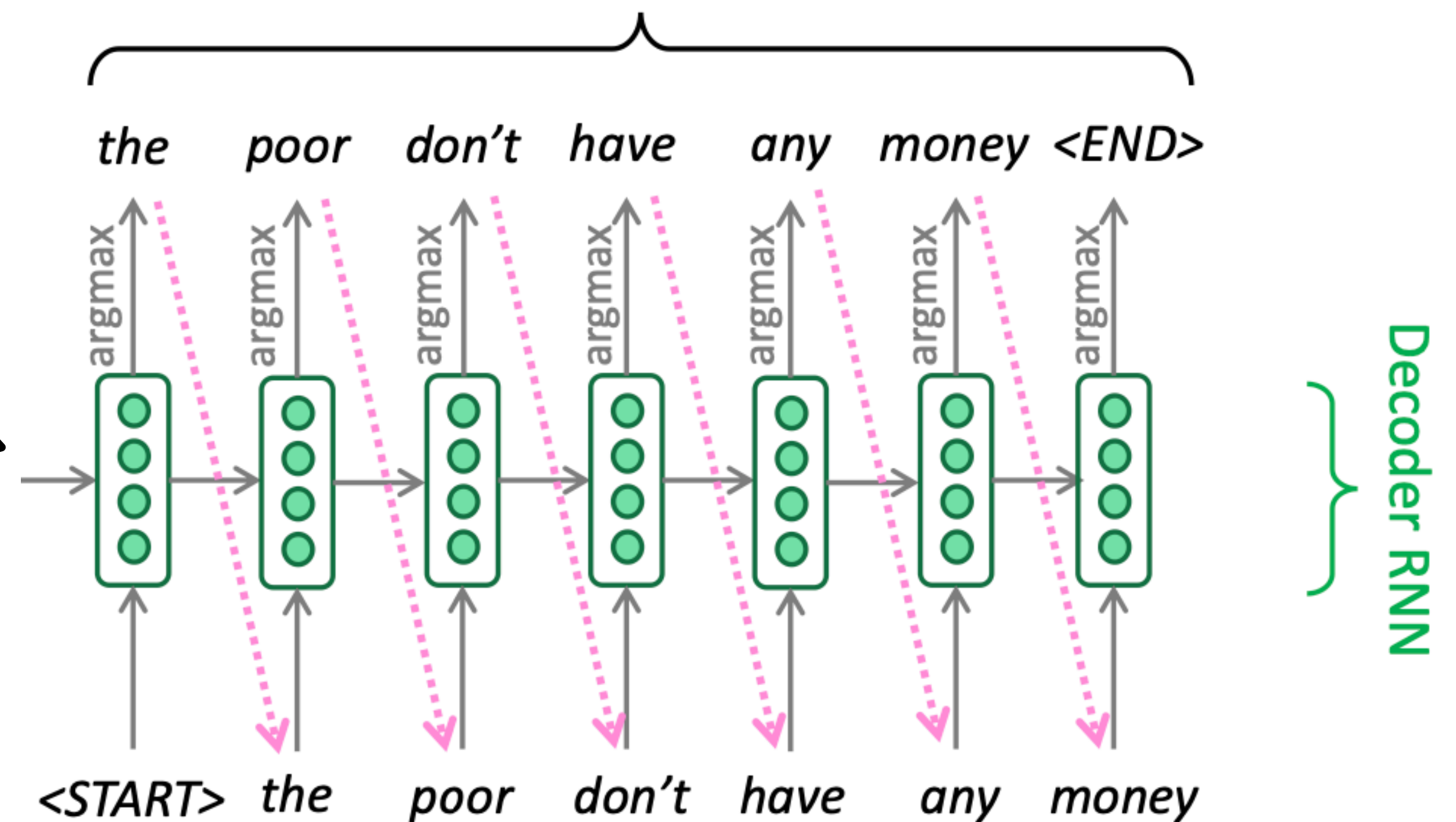
## Example of Neural Machine Translation

Encoding of the source sentence.  
Provides initial hidden state  
for Decoder RNN.



**Encoder RNN** produces an **encoding** of the source sentence.

Target sentence (output)



**Decoder RNN** is a Language Model that generates target sentence conditioned on **encoding**.



# 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

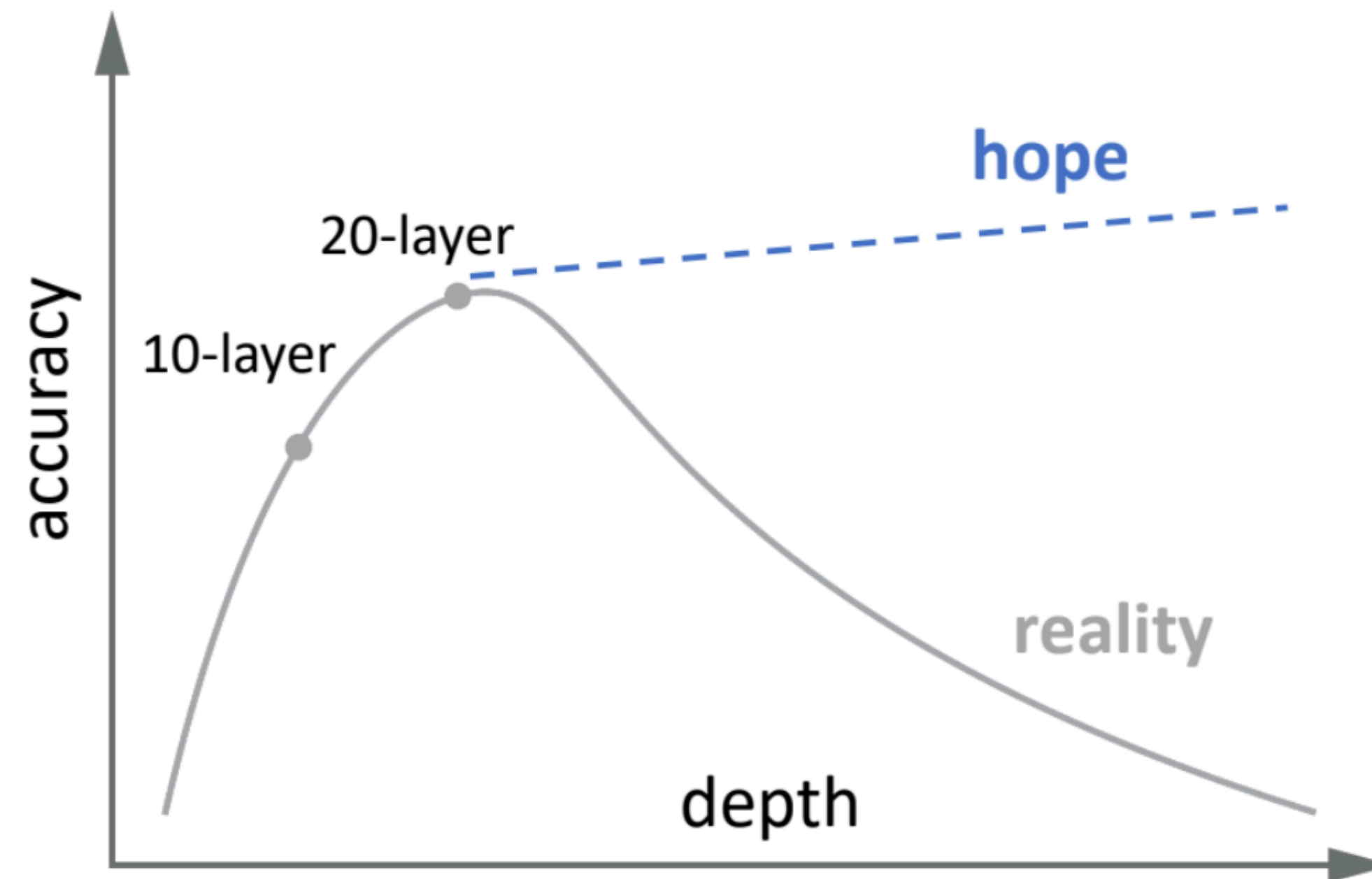
Normalization modules

Deep residual learning



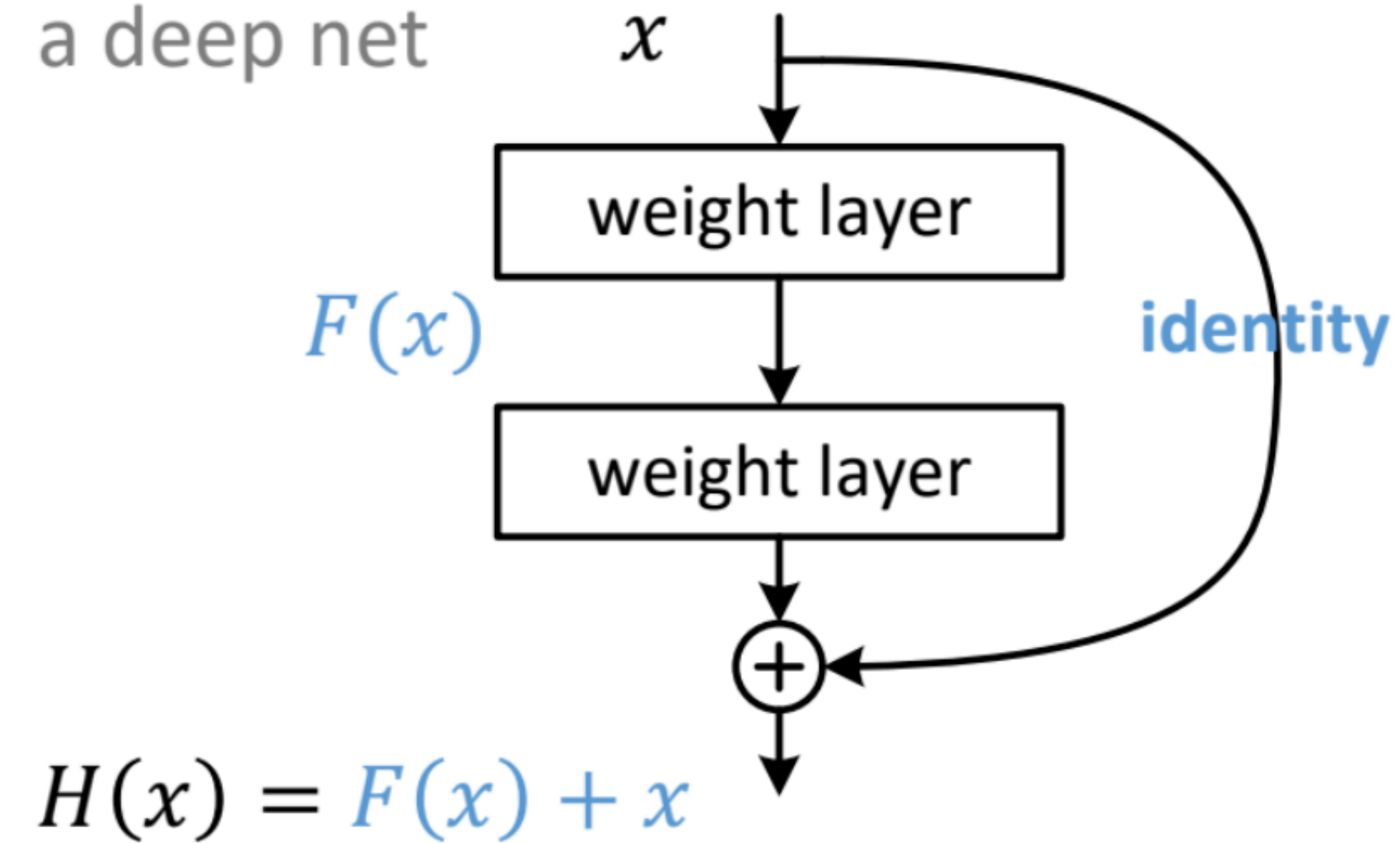
# The Degradation Problem

- Good init + norm enable training deeper models
- Simply stacking more layers?
- Degrade after ~20 layers
- Not overfitting
- Difficult to train



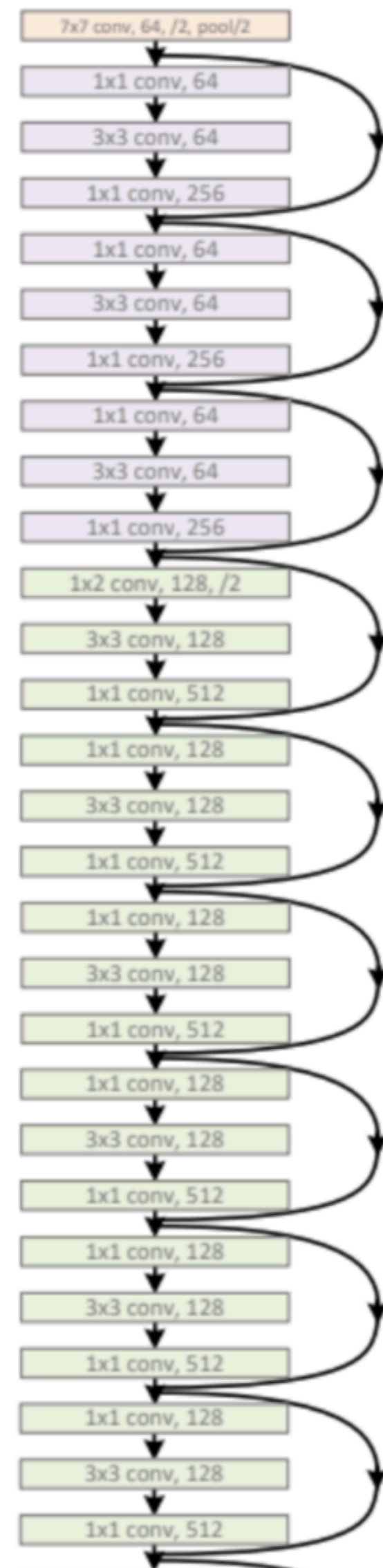
# Deep Residual Learning

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

