



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 ✓ Dec 8
3. HW3 is out
4. HW4 will be easy (multi-choice QA only)

public leaderboard 50%

private leaderboard

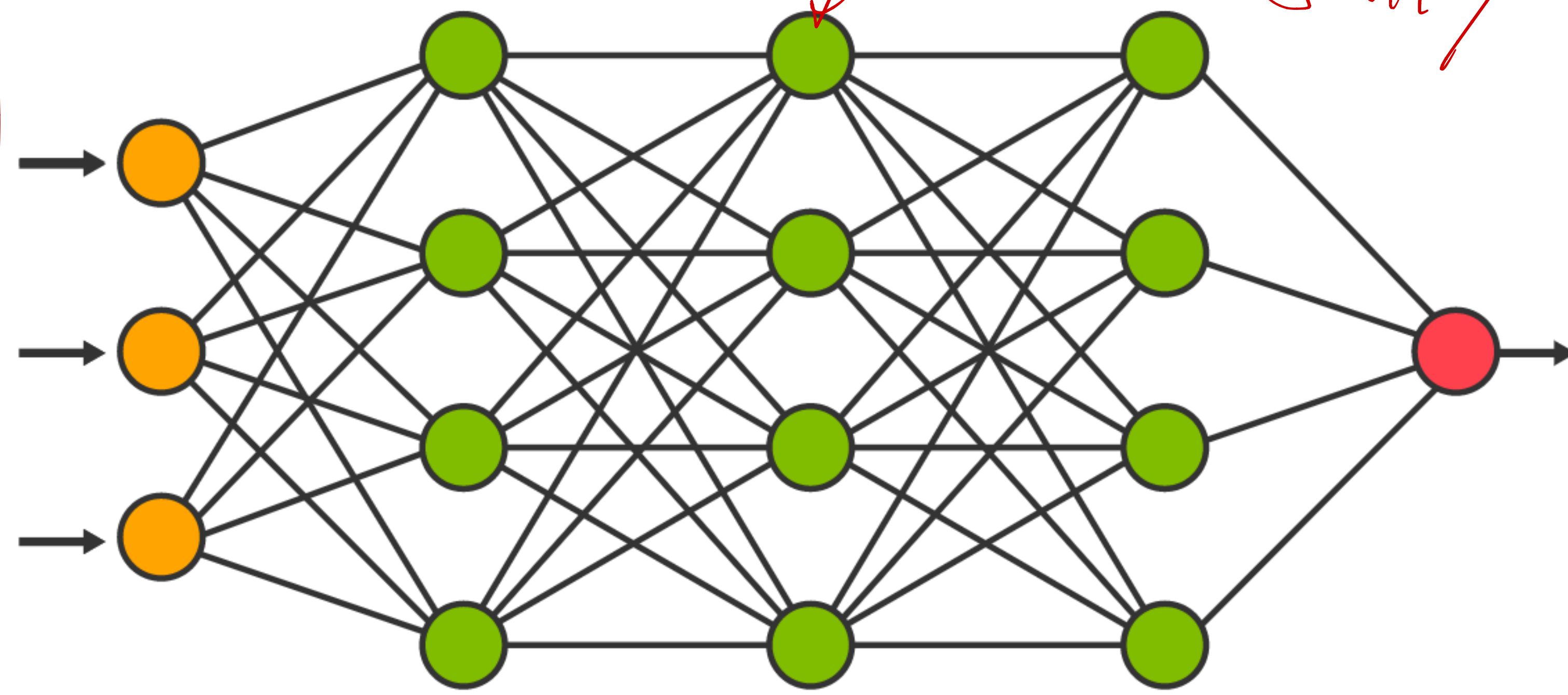


Recap: Multilayer Perceptron Neural Networks (MLP)

*Relu
tanh*

non linear activation

fully-connected



INPUT LAYER

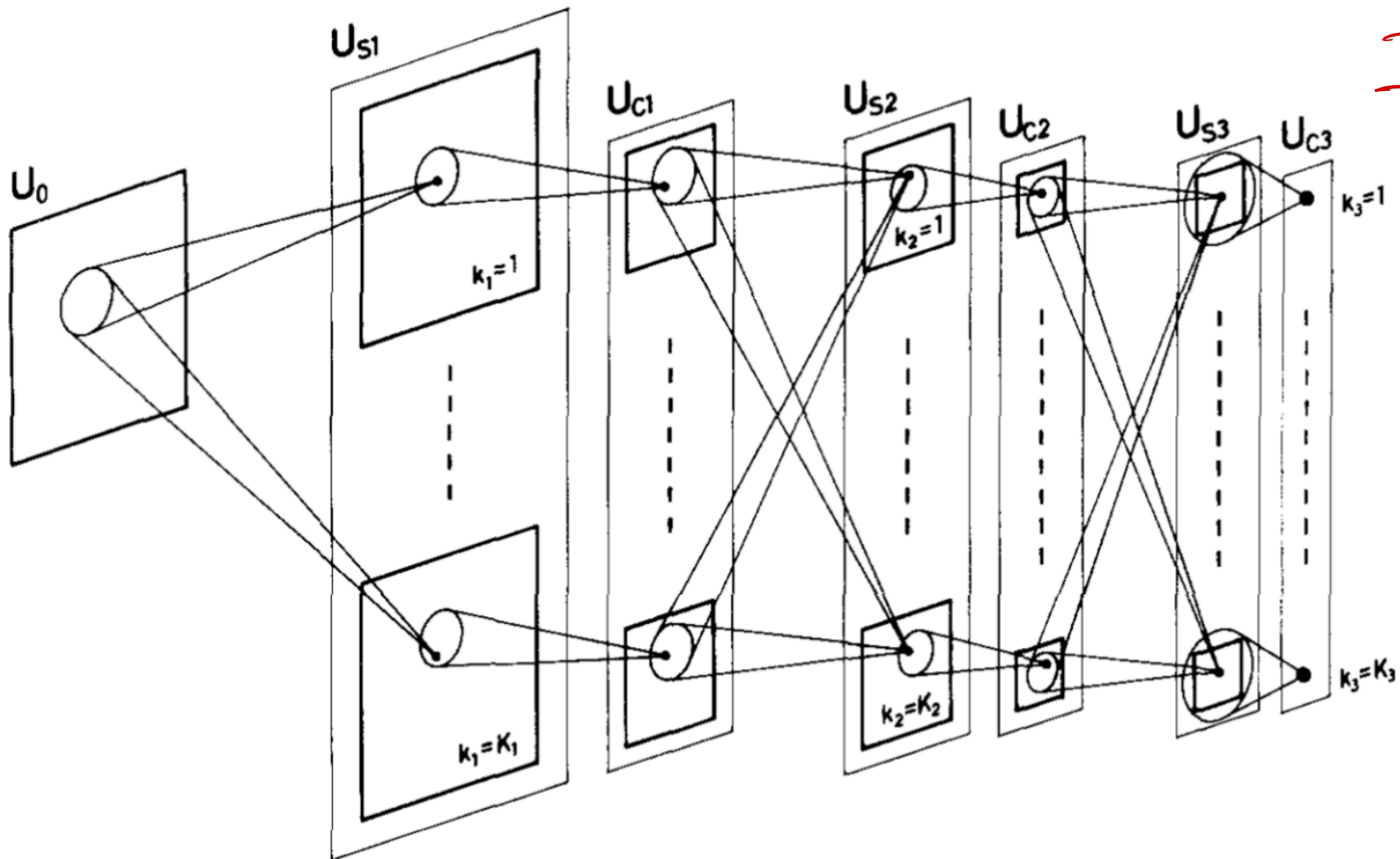
HIDDEN LAYERS

OUTPUT LAYER

3

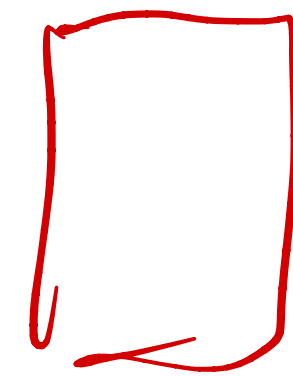
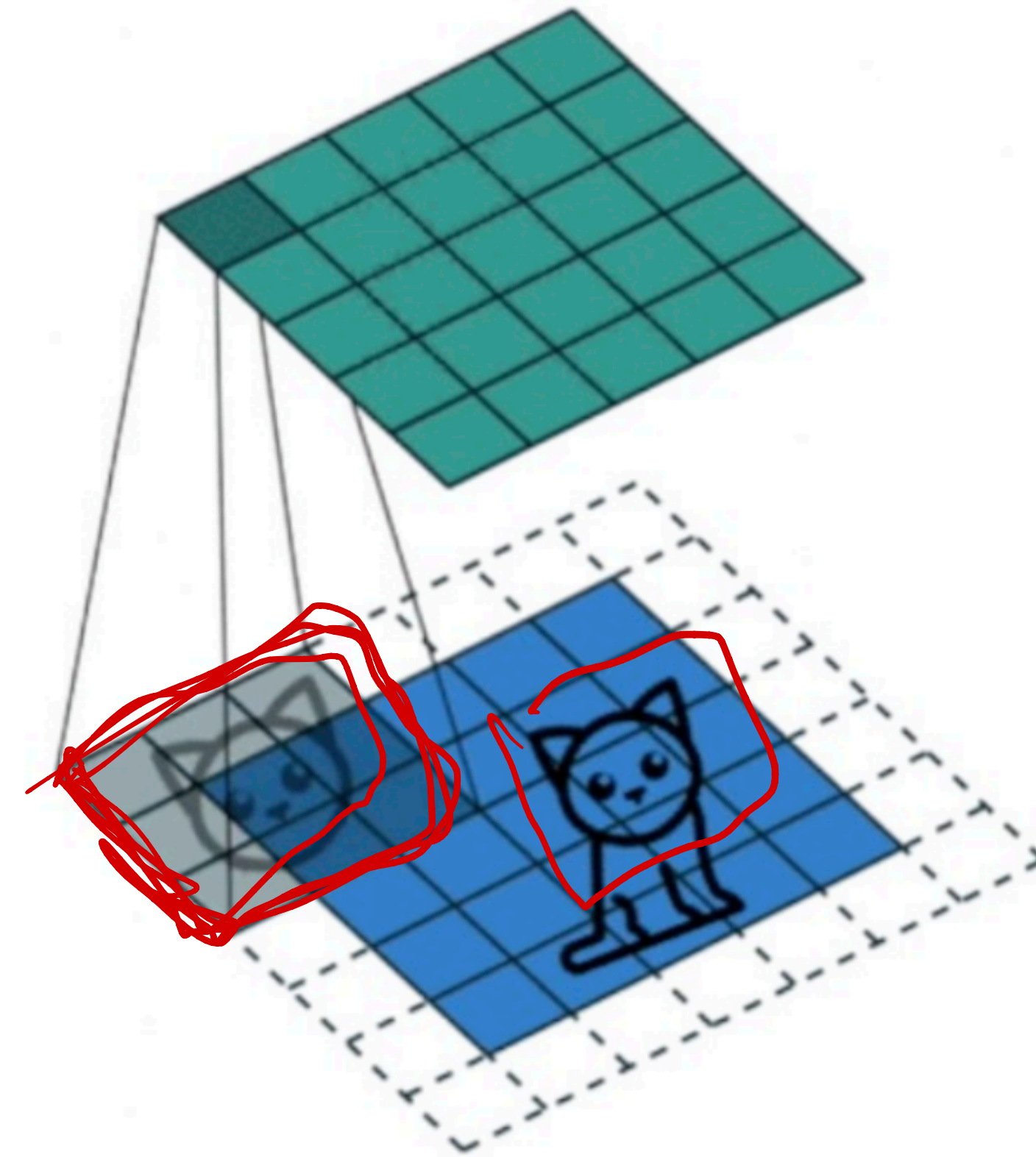
Convolutional Neural Networks

CNN



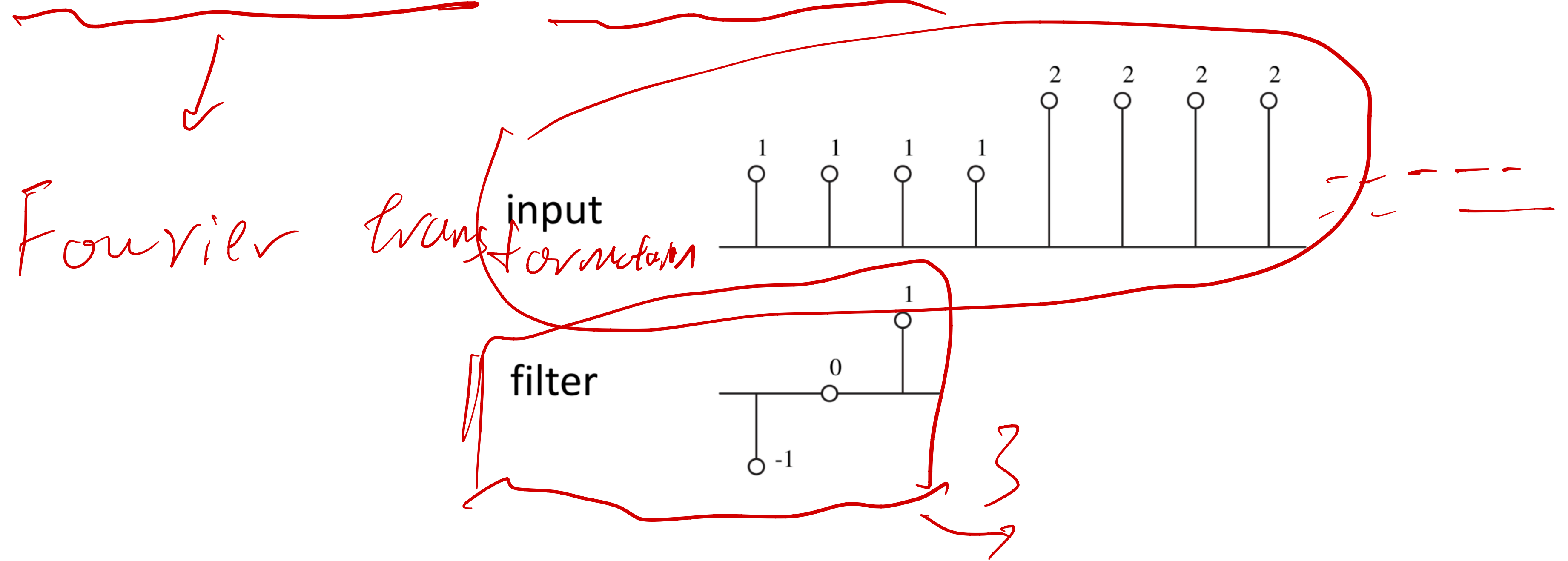
Convolution is template matching

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



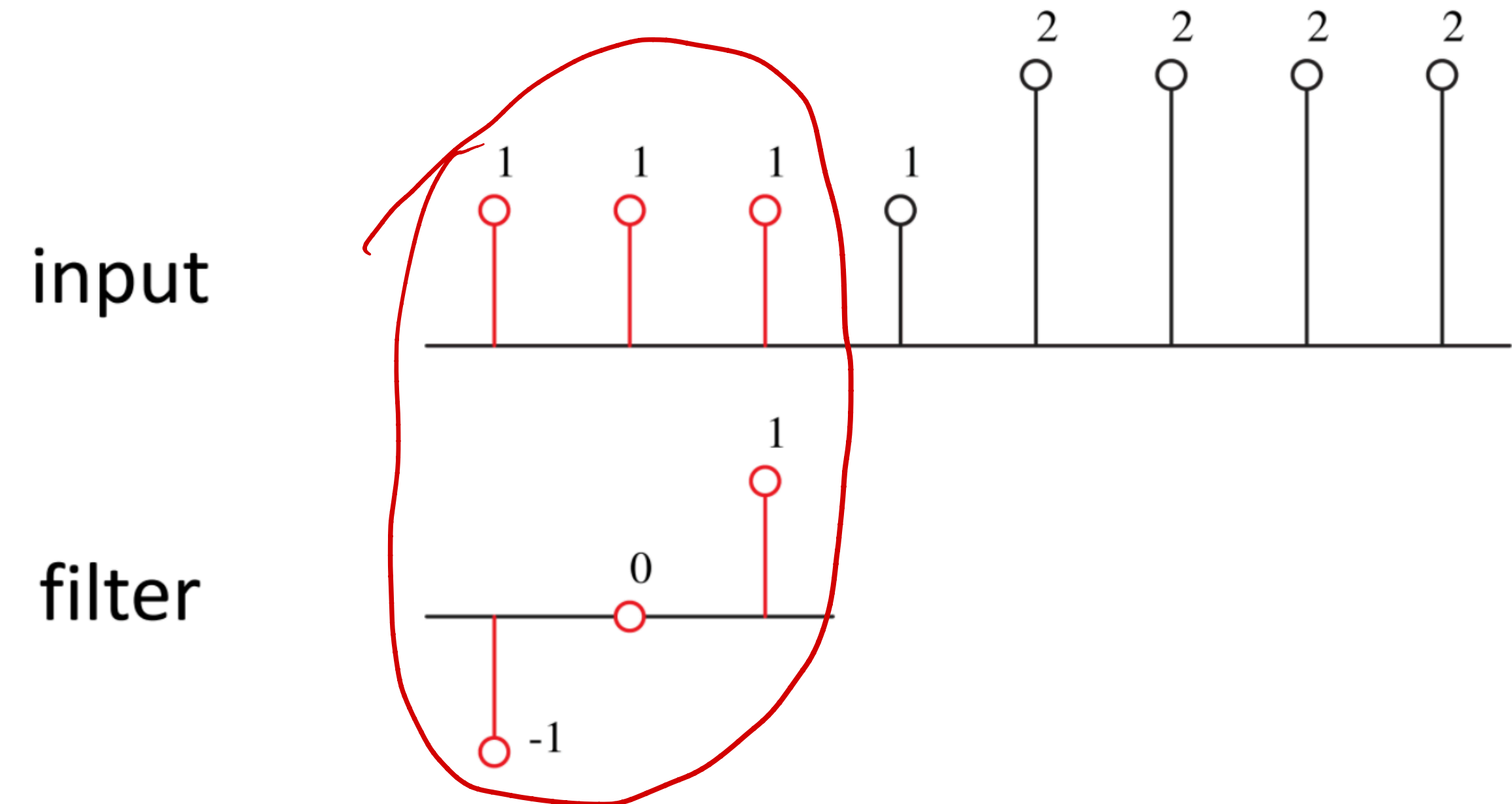
dot prod

Convolution: a 1-D example

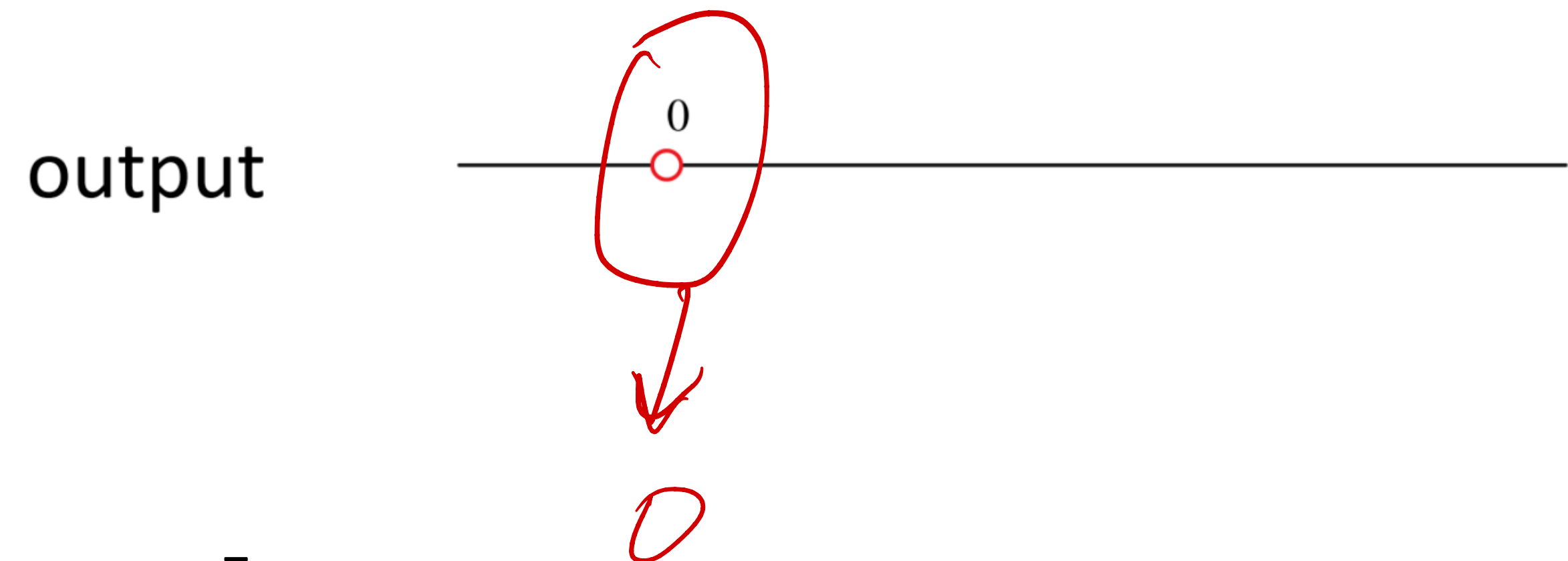


Convolution: a 1-D example

- sliding window
- dot product



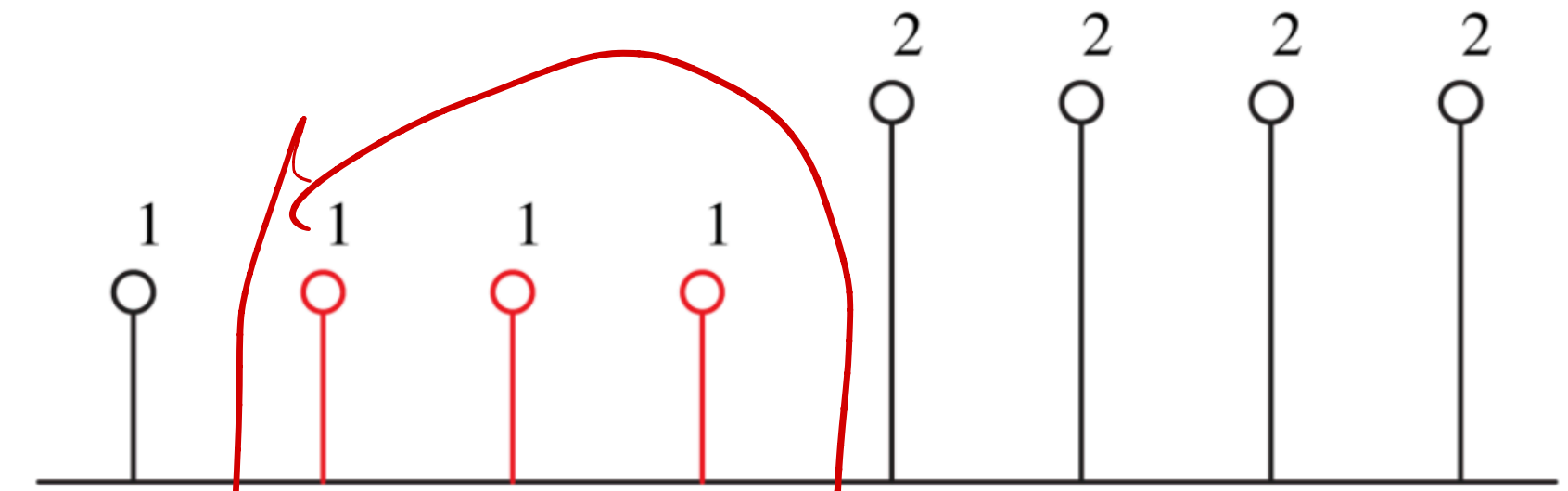
$$-1 \times 1 + 0 \times 1 + 1 \times 1 = 0$$



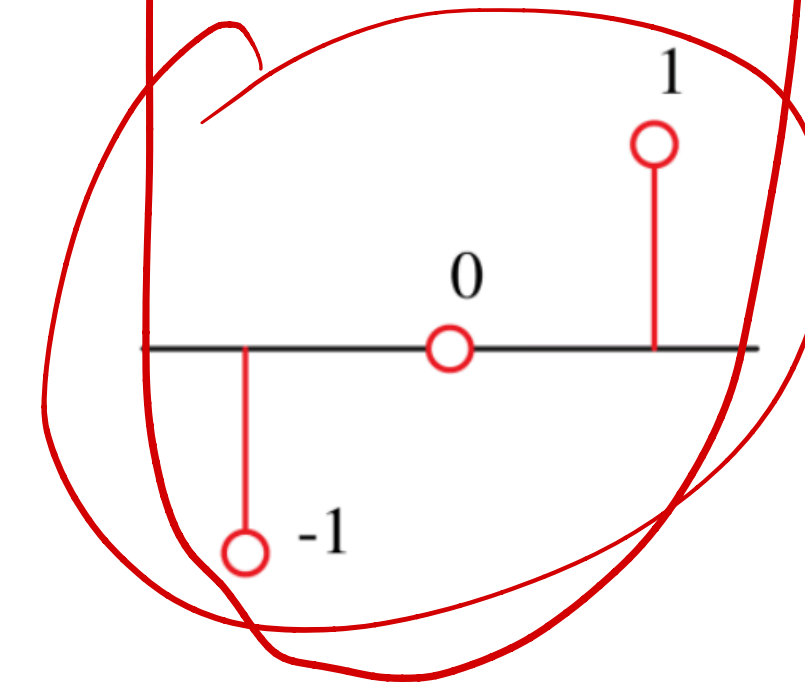
Convolution: a 1-D example

- sliding window
- dot product

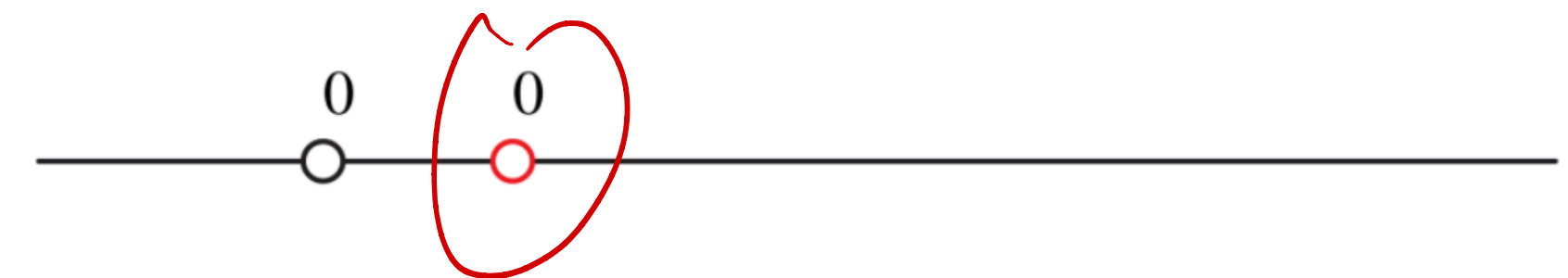
input



filter



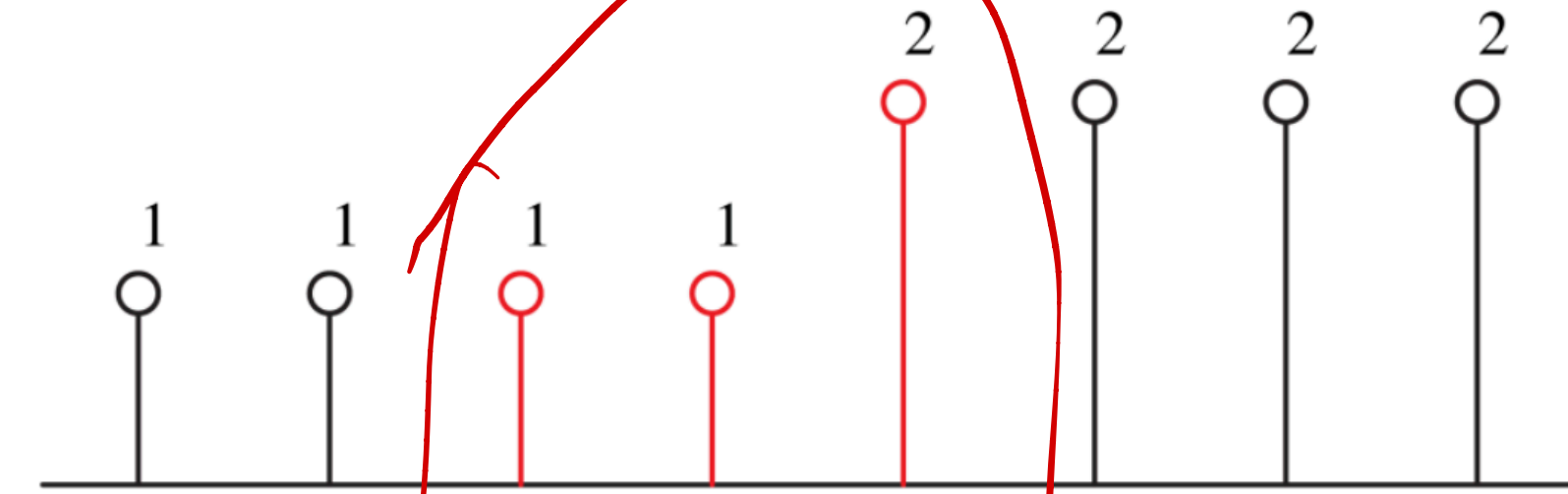
output



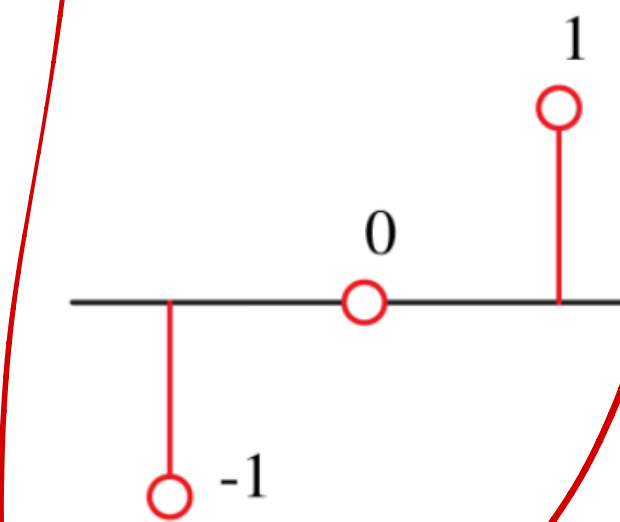
Convolution: a 1-D example

- sliding window
- dot product

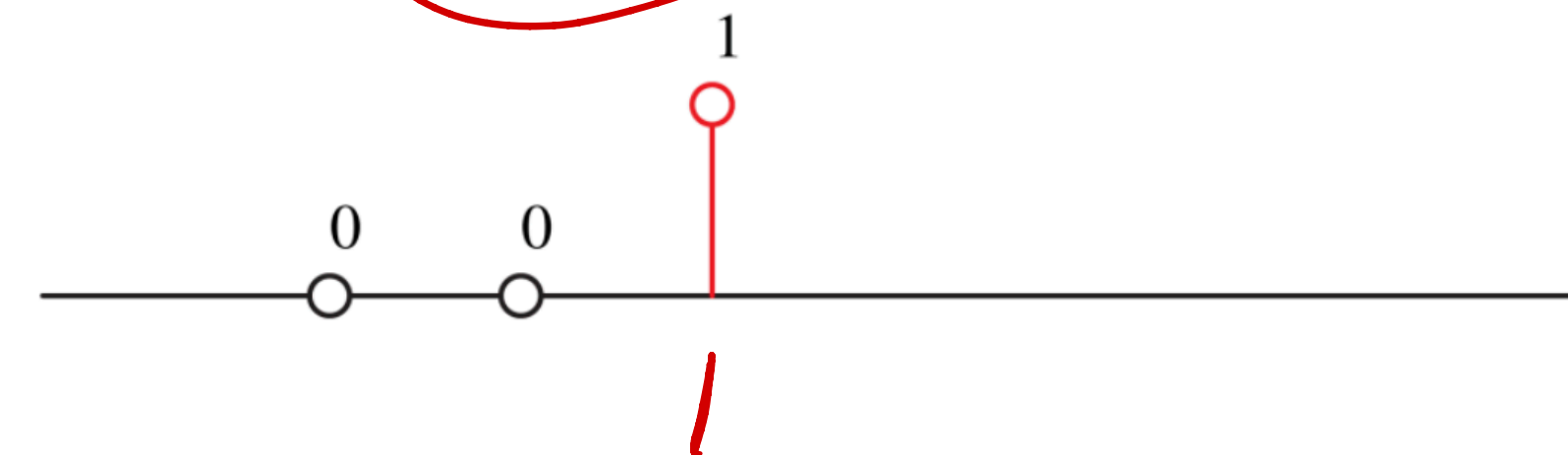
input



filter



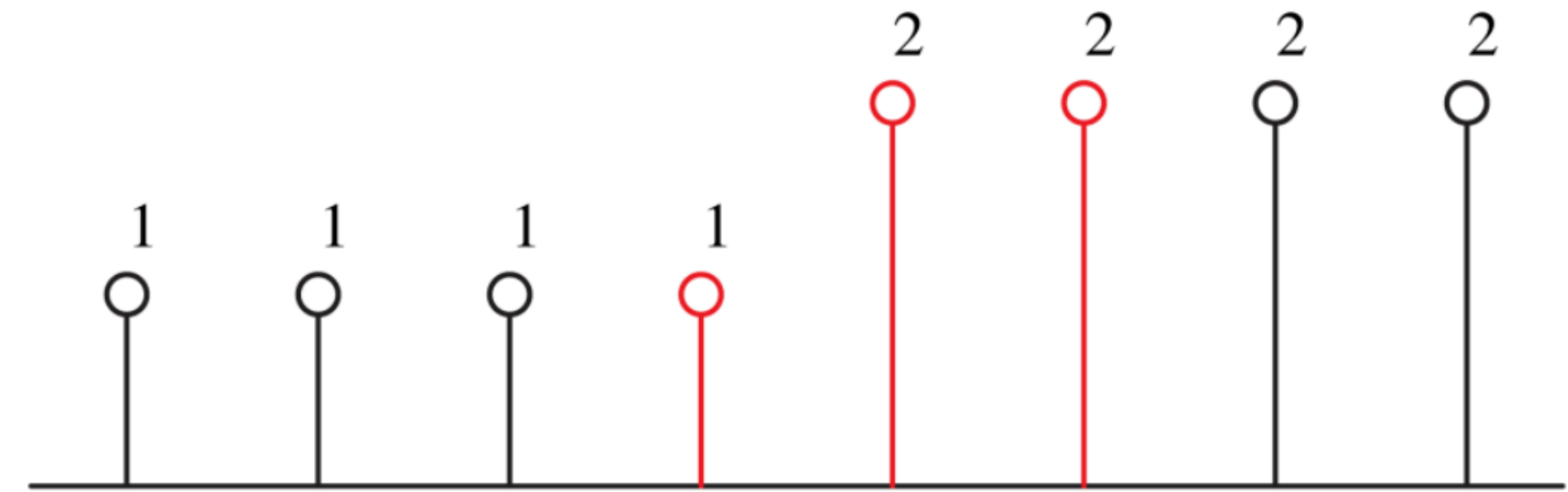
output



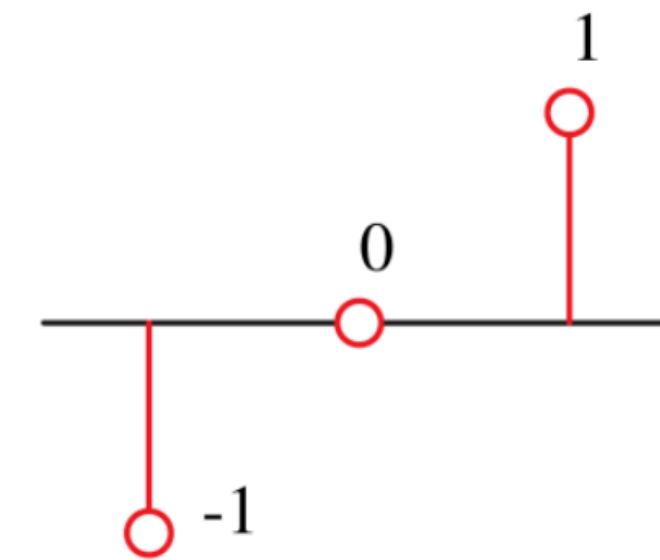
Convolution: a 1-D example

- sliding window
- dot product

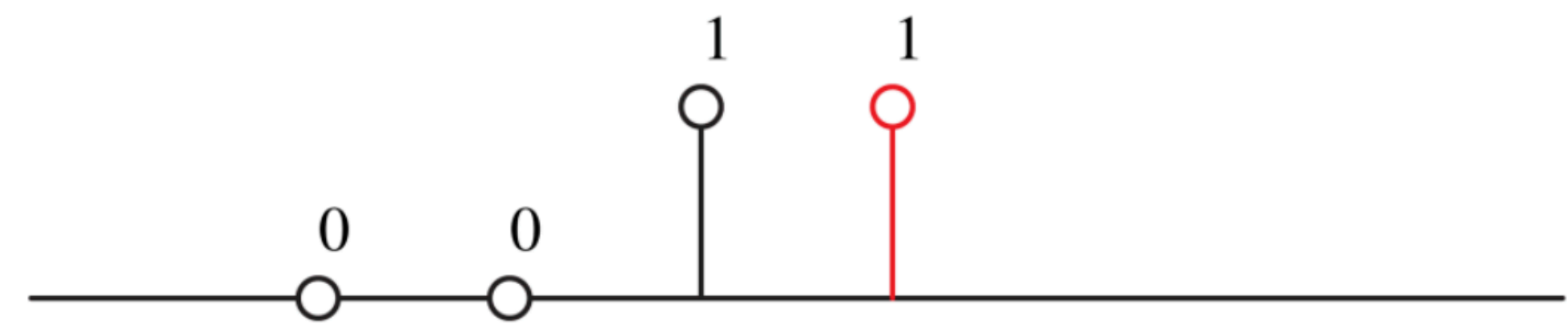
input



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

0

$$1 \times (-2) + 1 \times (-1) = -3$$

input

0 ¹	0 ²	0 ¹	0	0	0	0	0
0 ⁰	0 ⁰	0 ⁰	0	0	1	1	0
0 ⁻¹	1 ⁻²	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

-3					

- 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	1	1	0
0	1 ⁻¹	1 ⁻²	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	-4				

- sliding window
- dot product

Convolution: a 2-D example

input

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

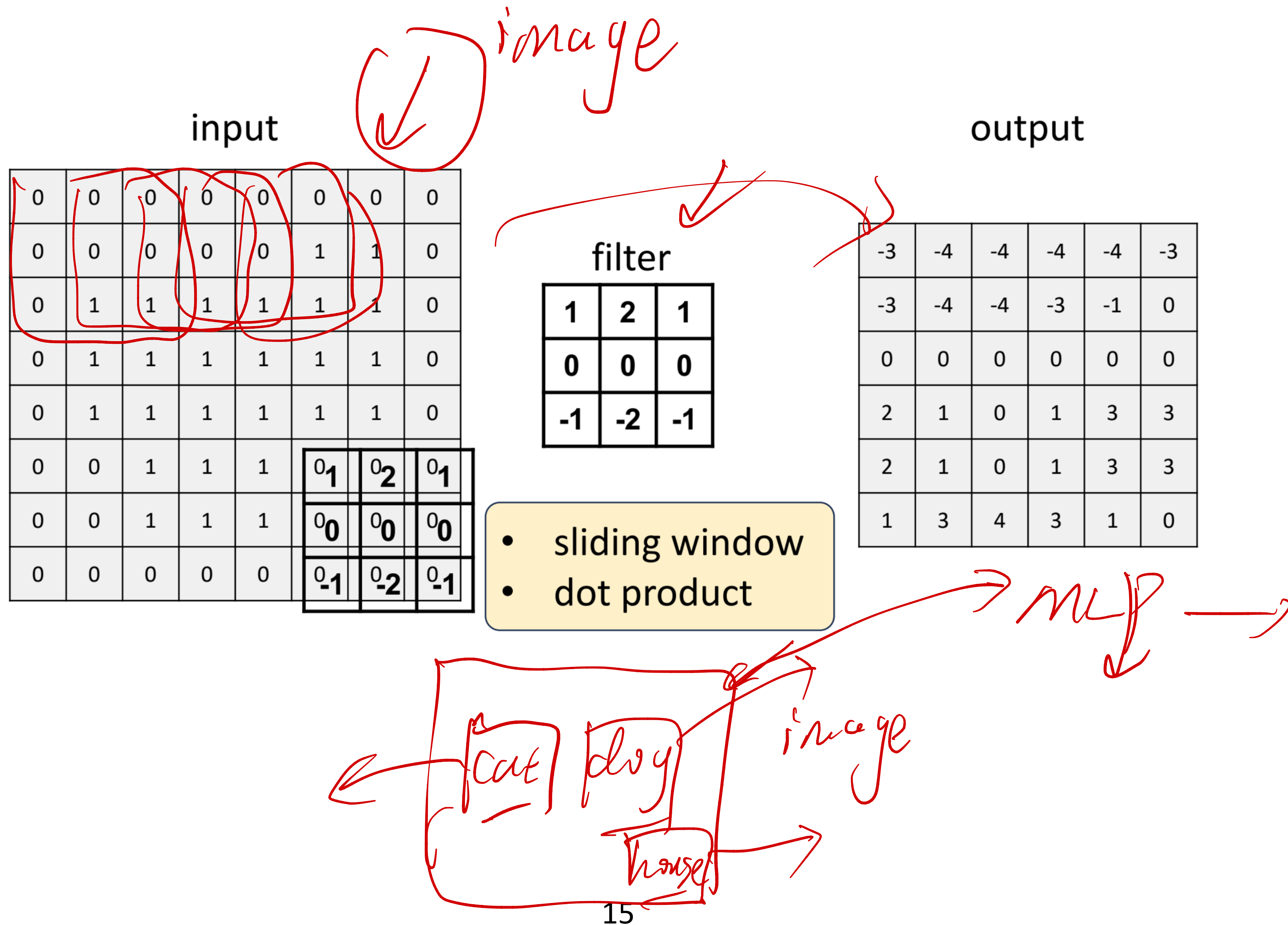
1	2	1
0	0	0
-1	-2	-1

output

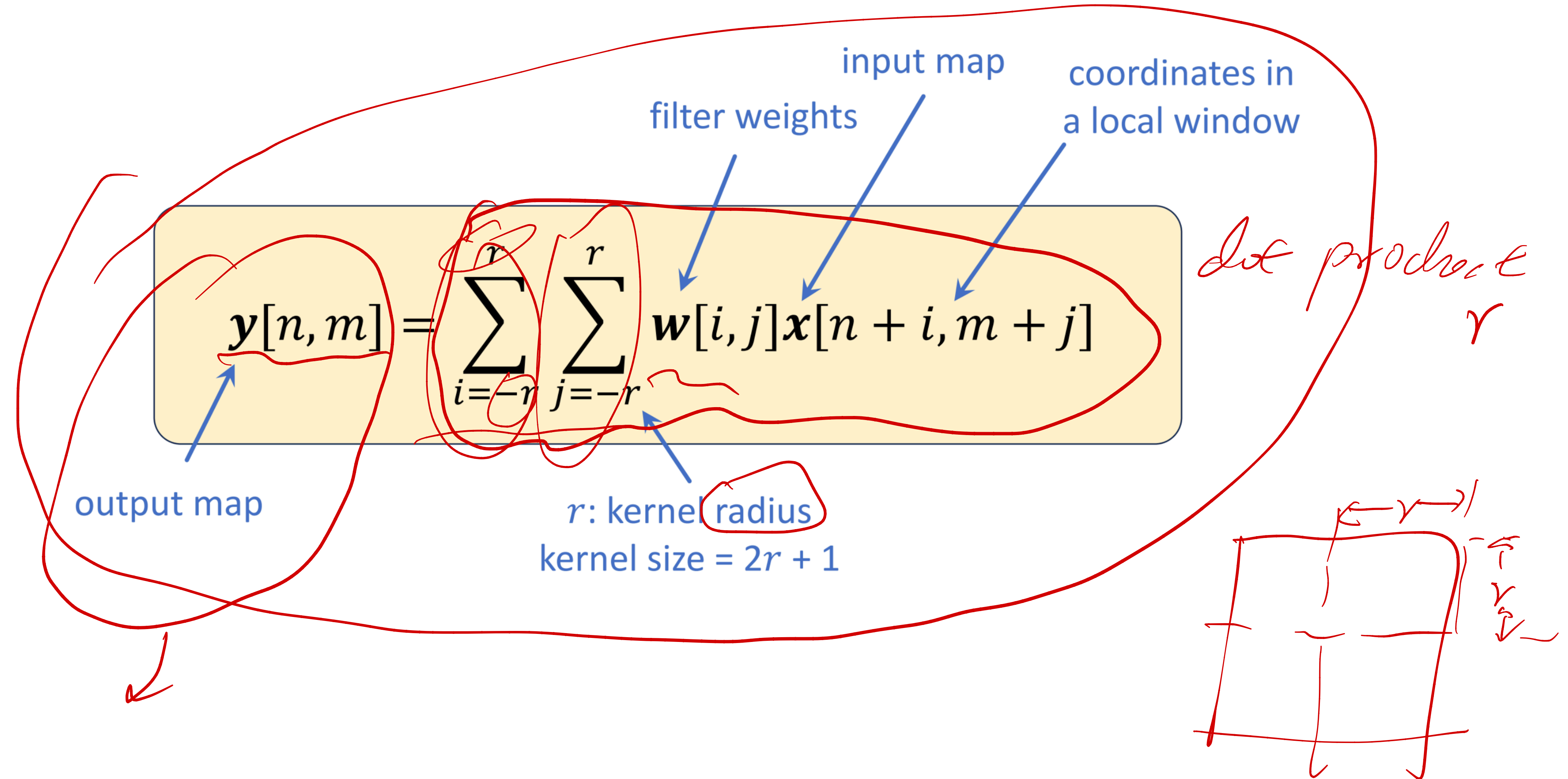
-3	-4	-4			

- sliding window
- dot product

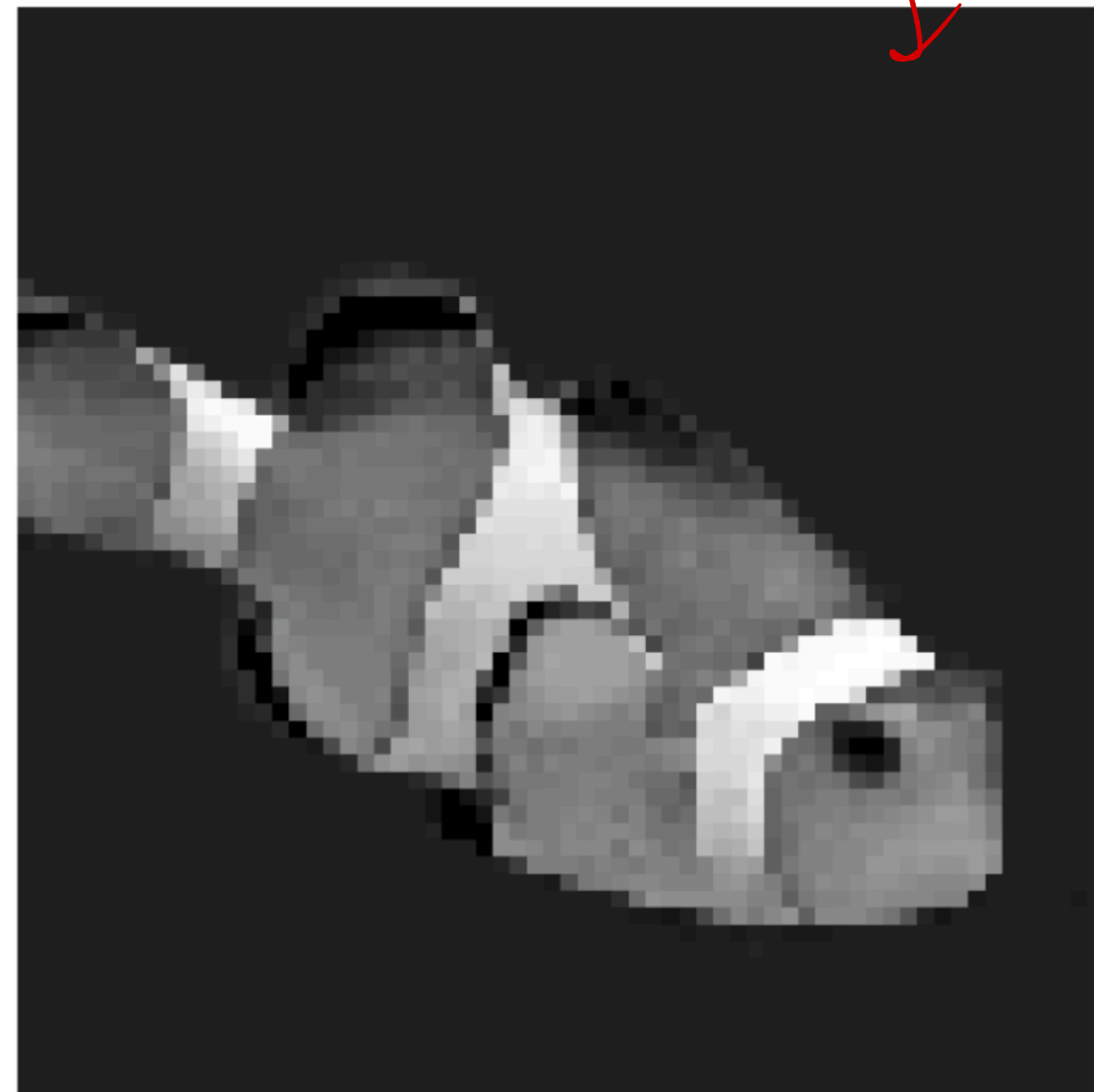
Convolution: a 2-D example



Convolution: a 2-D example



Convolution: 2-D



filter

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

* =

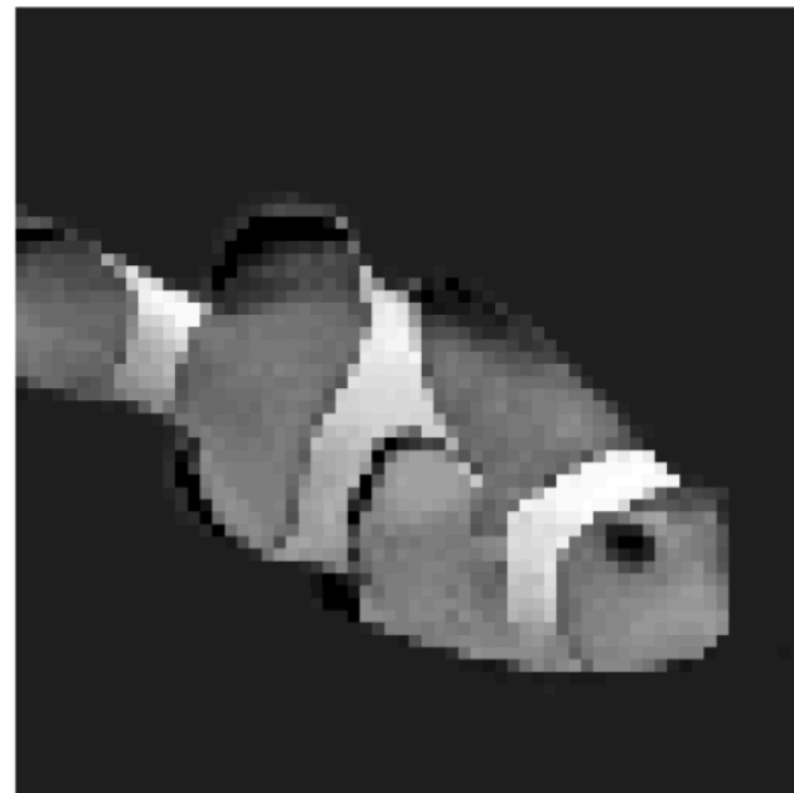


convolution

Convolution: Multi-channel outputs

representation learning

feature extraction



*

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

=



one filter, one feature

*

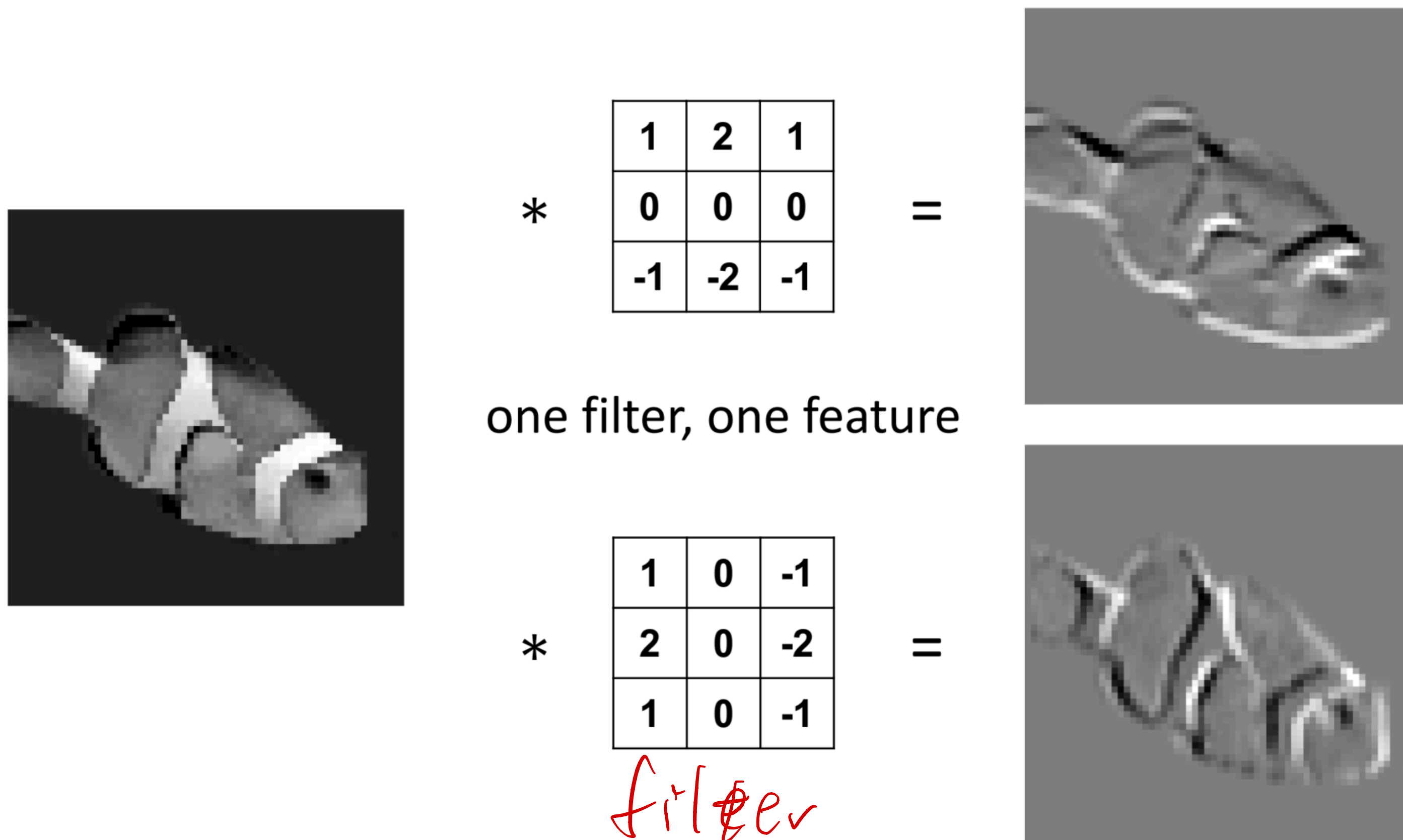
1	0	-1
2	0	-2
1	0	-1

=



color change

Convolution: Multi-channel outputs

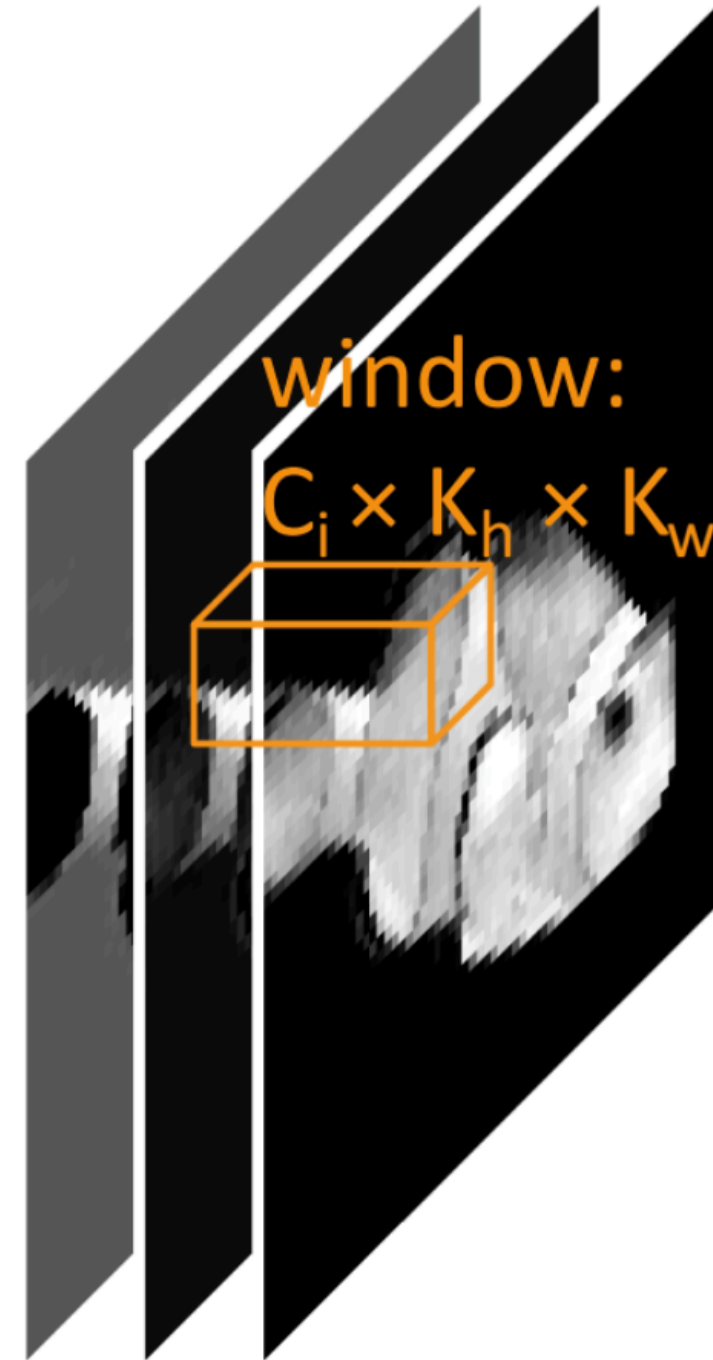


Understanding the filter/kernel as feature extractors

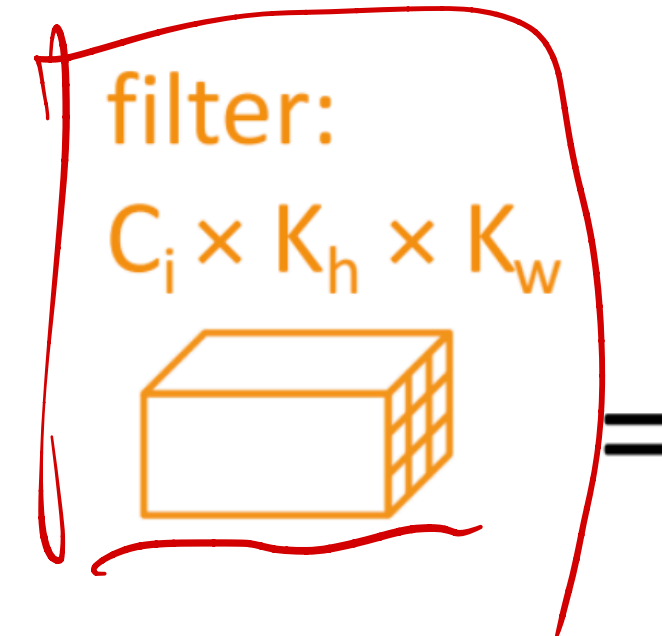
Convolution: Multi-channel inputs

input channel size = 3

$$C_i \times K_h \times K_w$$



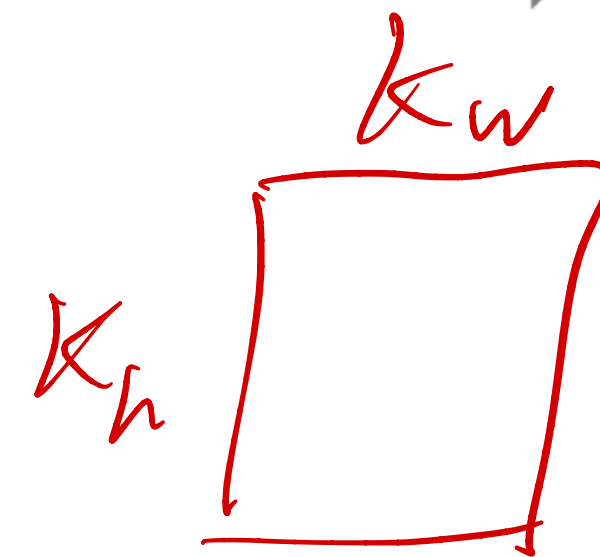
*



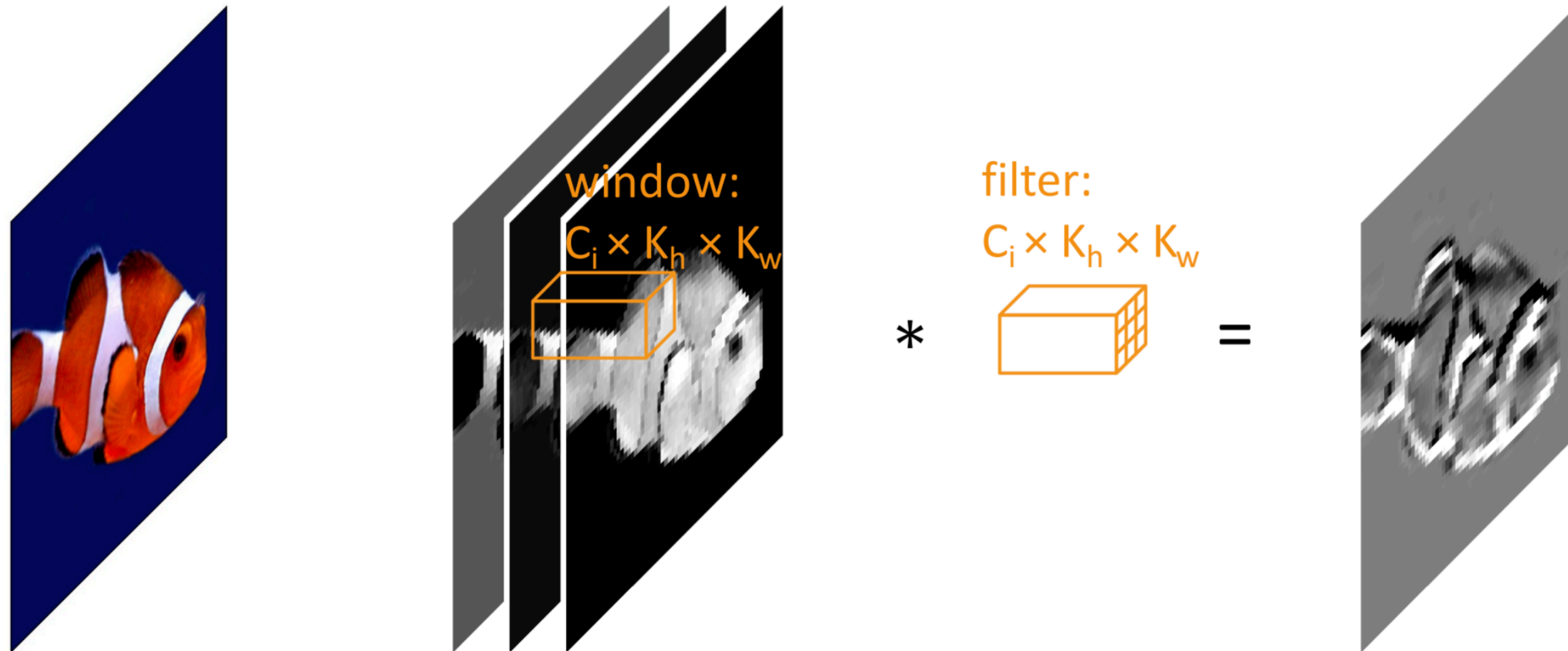
=



(R, G, B)
↓ ↓ ↓

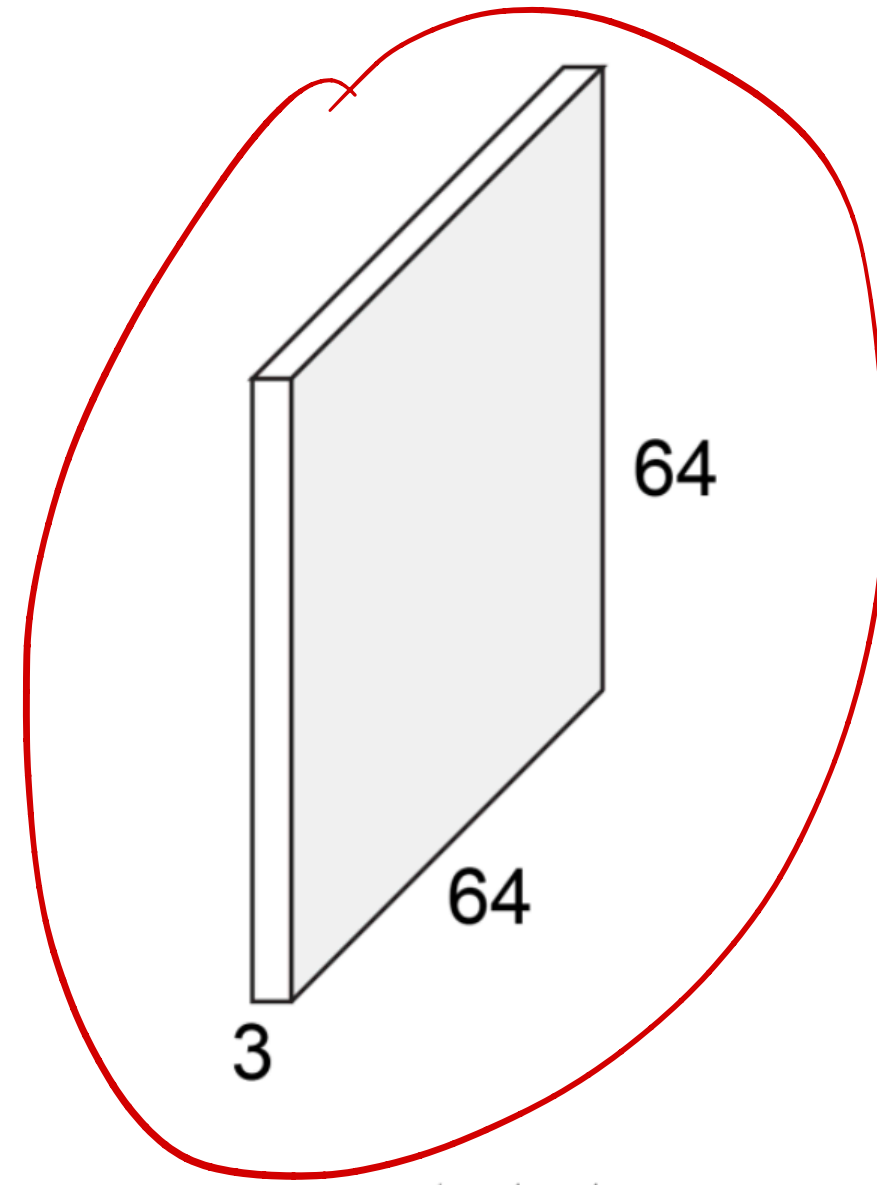


Convolution: Multi-channel inputs



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

Convolution: tensor views

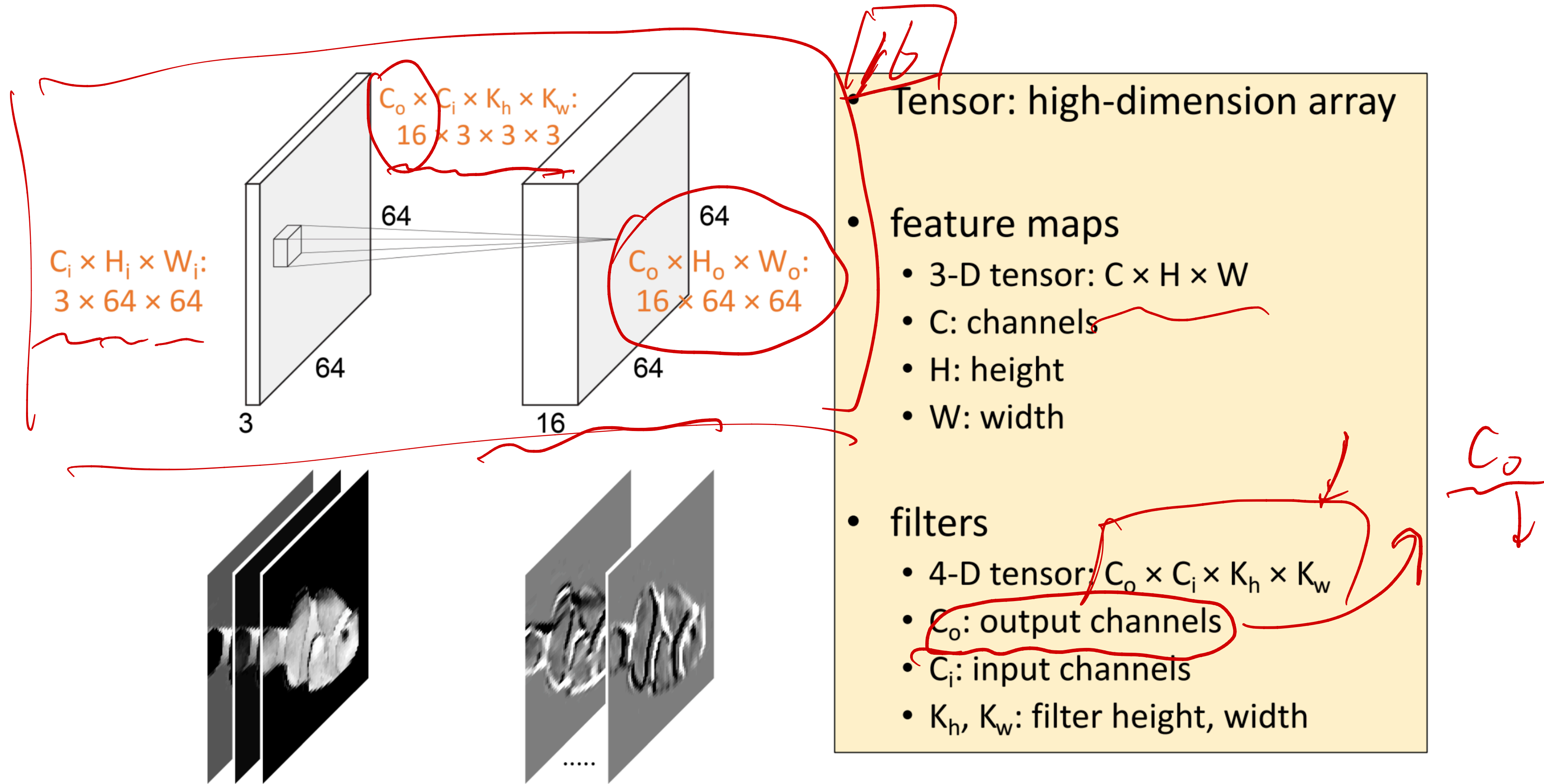


$3 \times 64 \times 64$
↓
RGB

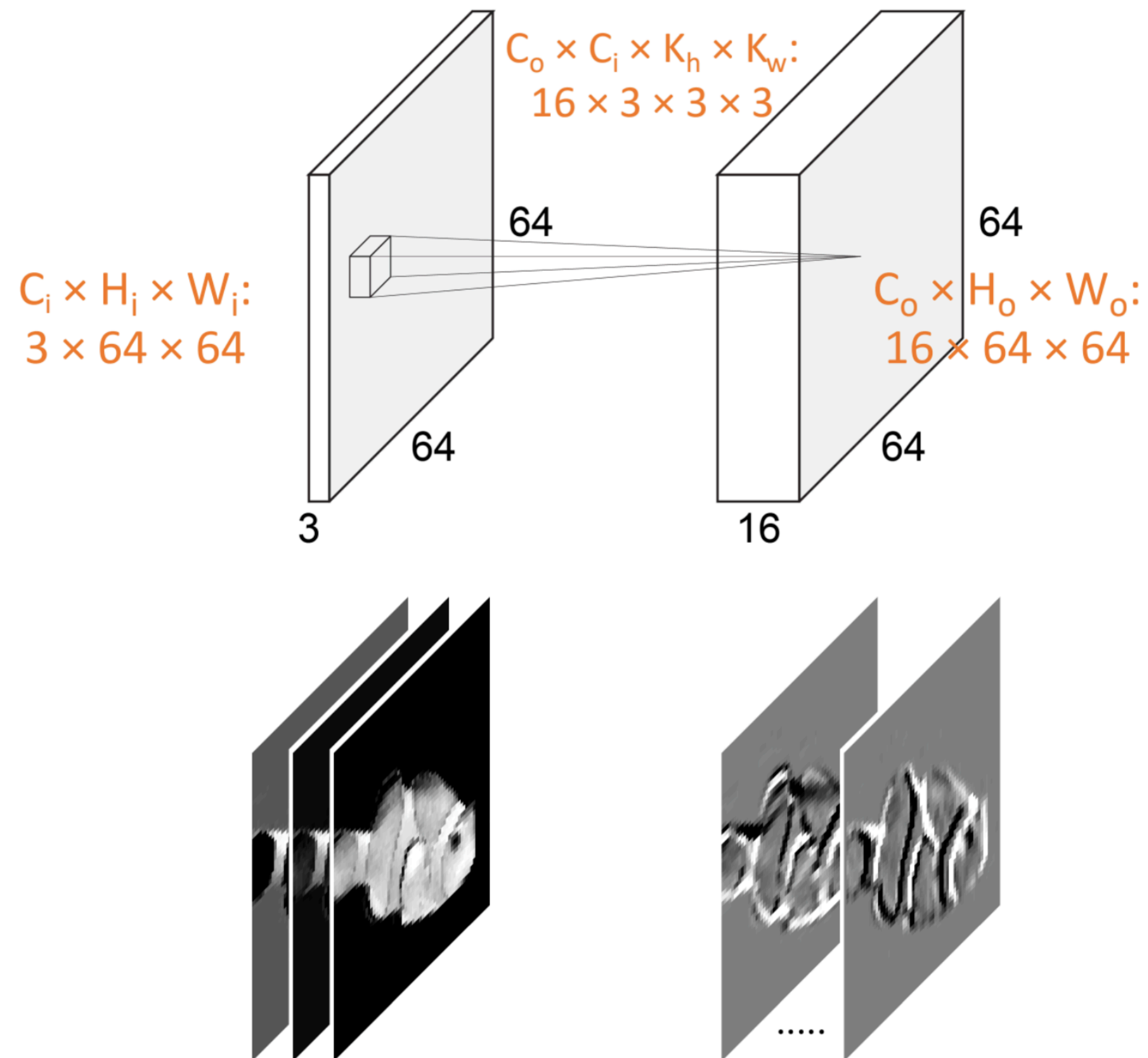


- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C: channels ↙
 - H: height
 - W: width

Convolution: tensor view



Convolution: tensor view



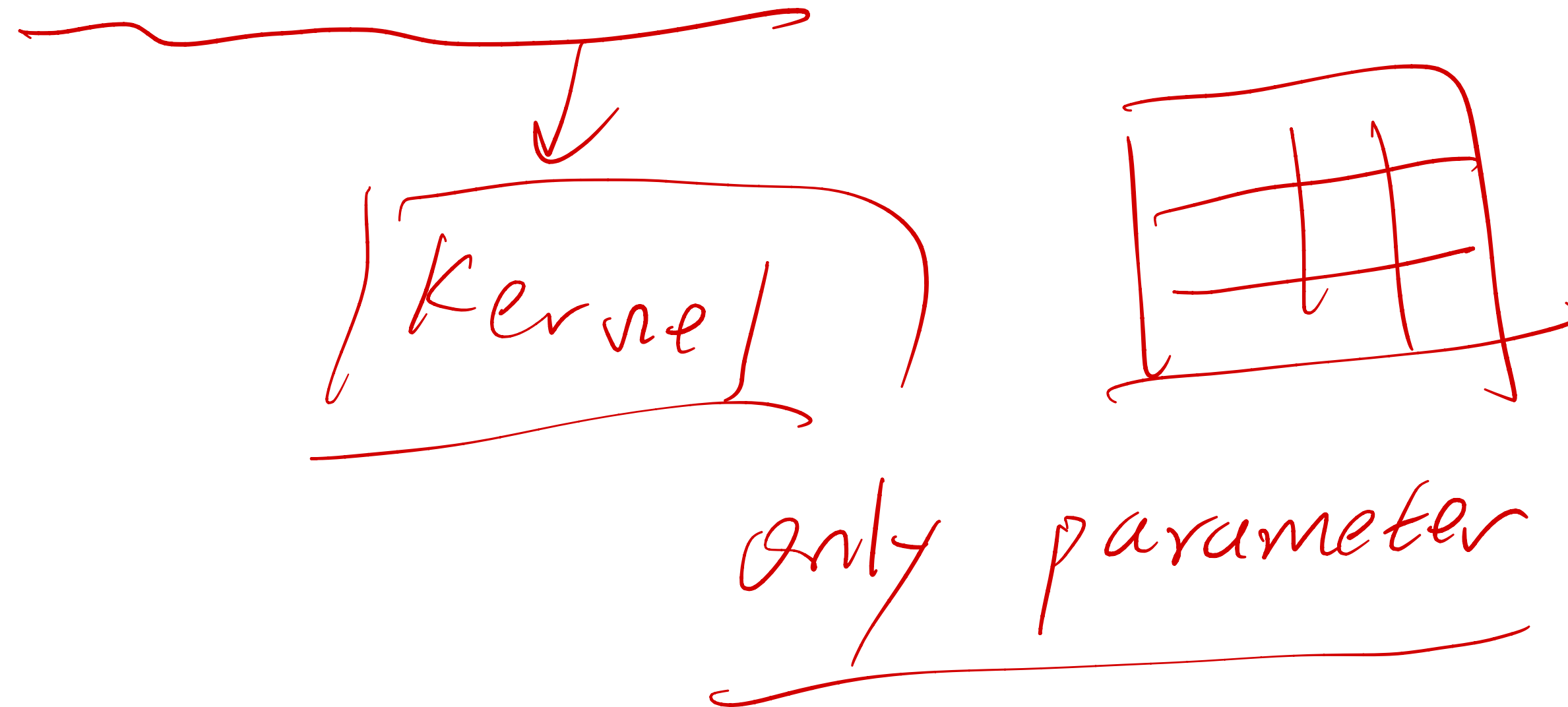
- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C: channels
 - H: height
 - W: width
- filters
 - 4-D tensor: $C_o \times C_i \times K_h \times K_w$
 - C_o : output channels
 - C_i : input channels
 - K_h, K_w : filter height, width

HMM

Stationary

The same filter tensor applies to different locations

Convolution: # parameters and # operations



Convolution: # parameters and # operations

- # parameters
 - weights: $C_o \times C_i \times K_h \times K_w$
 - bias: C_o

wx + (b) bias

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 $\times H_o \times W_o$

$$O(C_o \times \underbrace{H_o \times W_o}_{\substack{\downarrow \\ \uparrow}}) \times C_o \times C_i \times (K_h \times K_w)$$

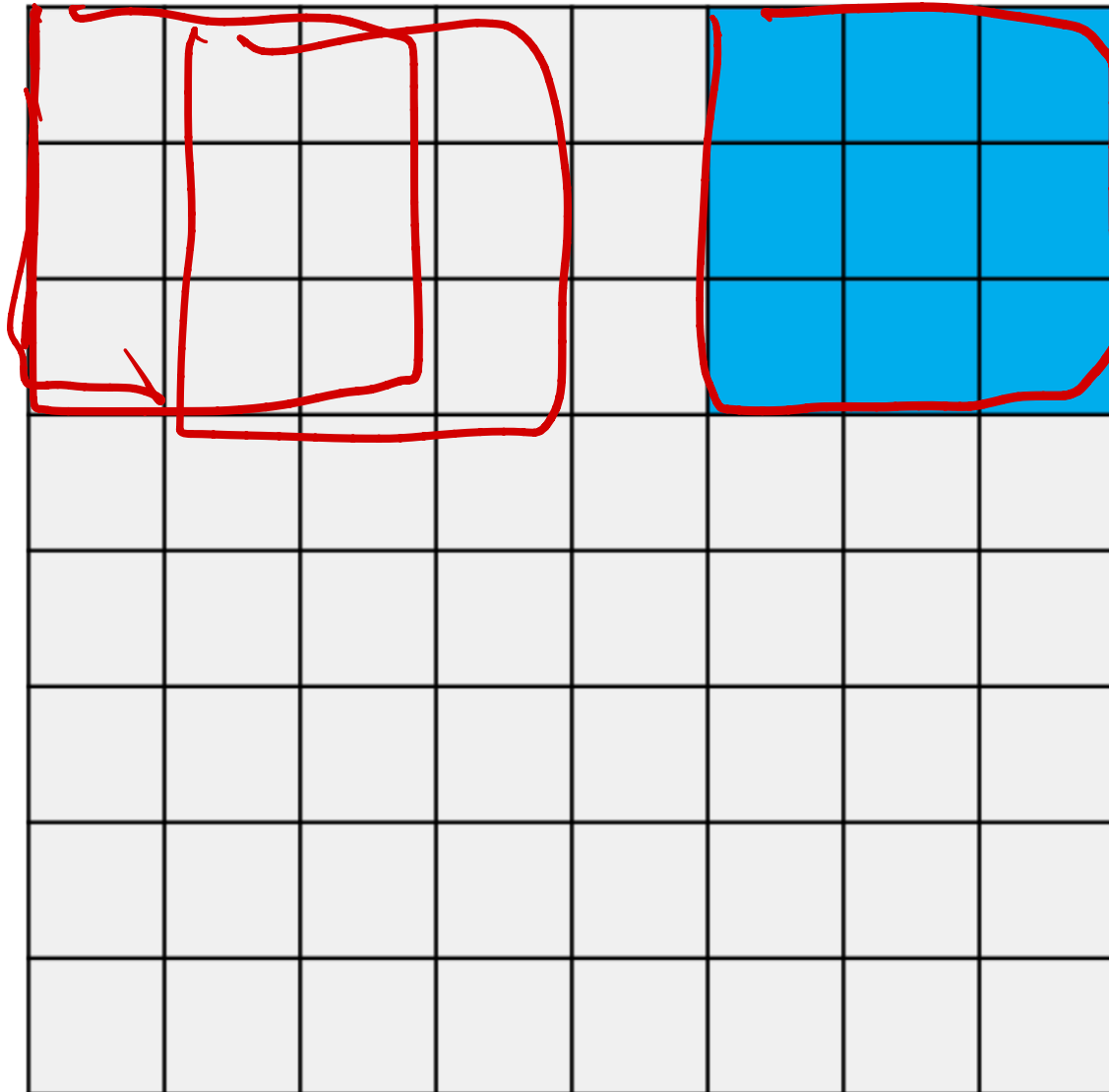
dot product = $O(C_o \times C_i \times K_h \times K_w)$

dot product?

Convolution: padding

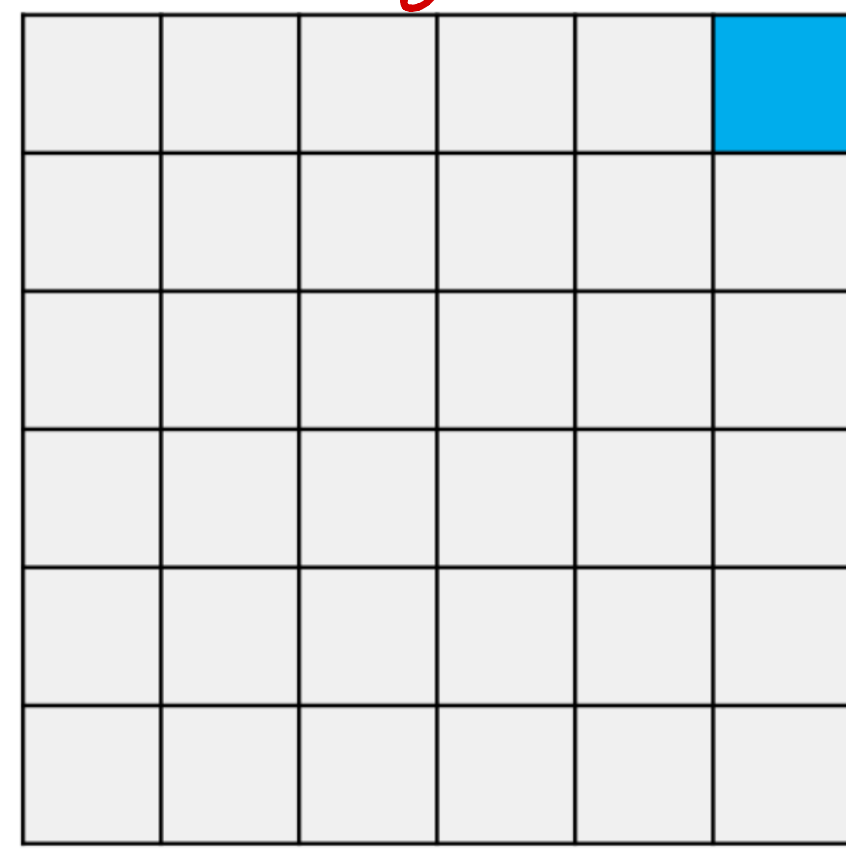
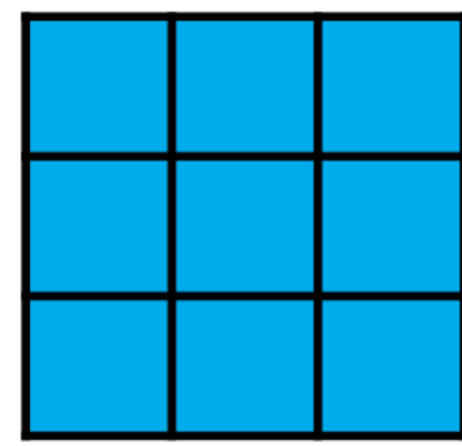
8x8

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



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

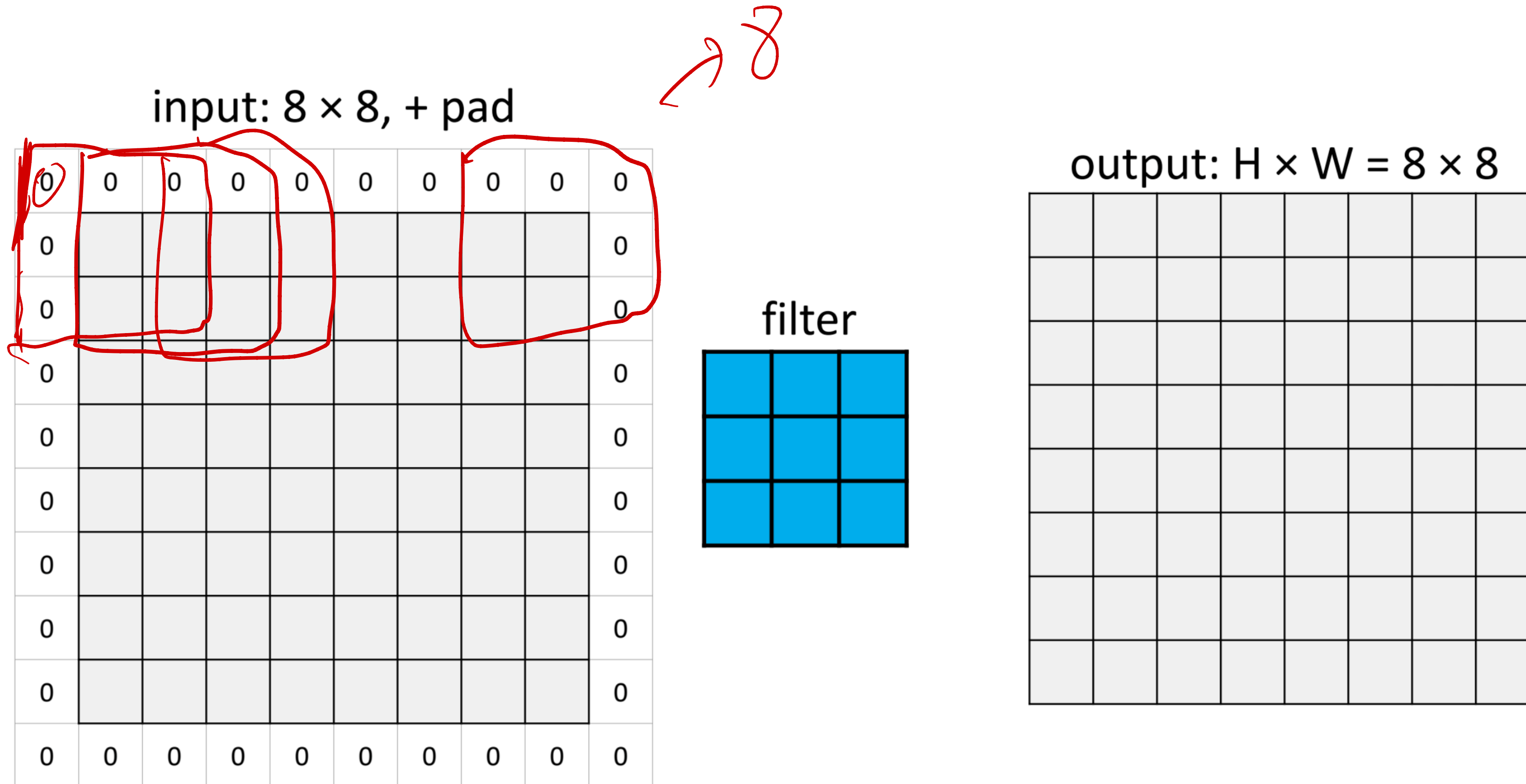
filter



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

$$8 - 3 + 1 = 6$$

Convolution: padding

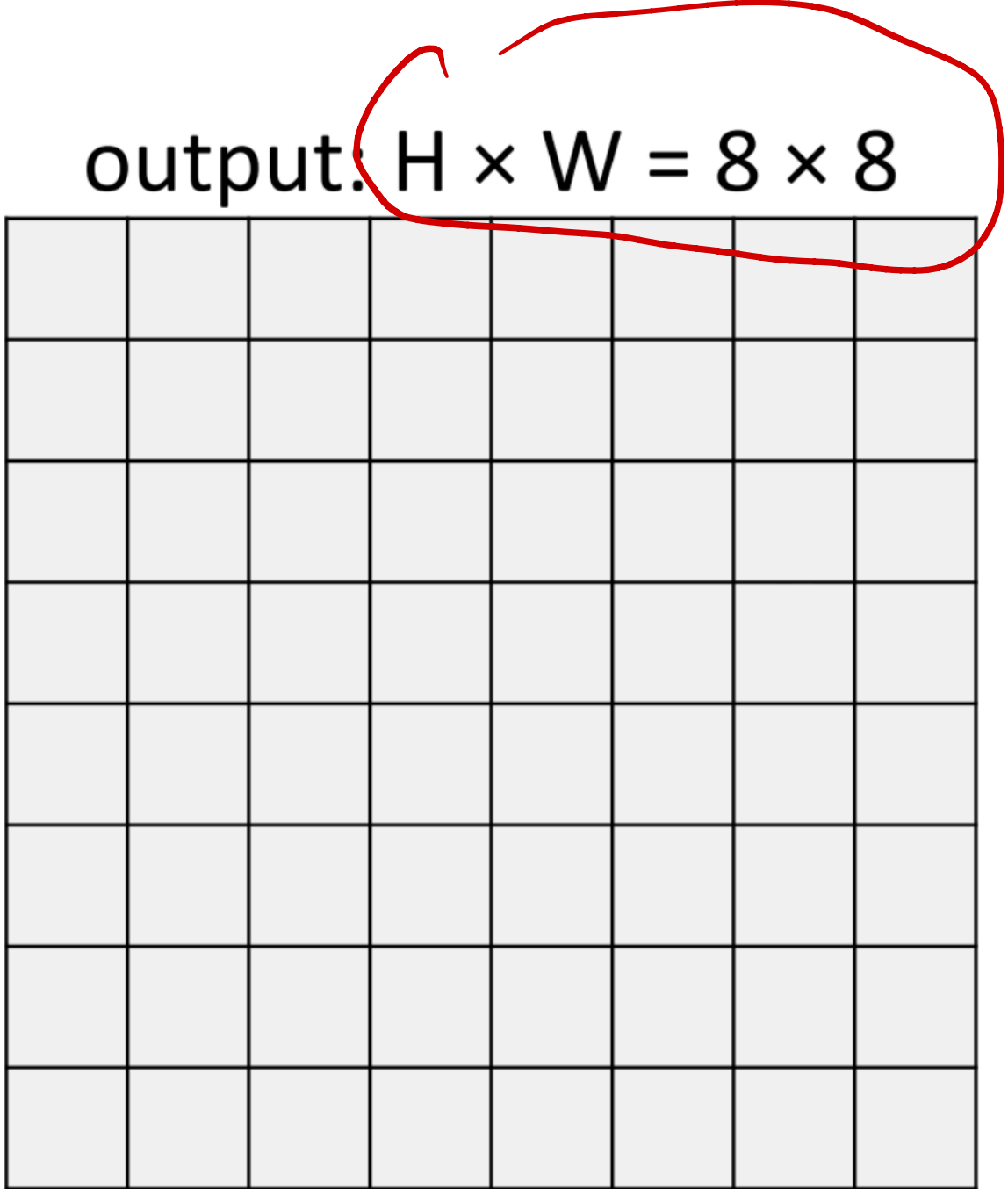
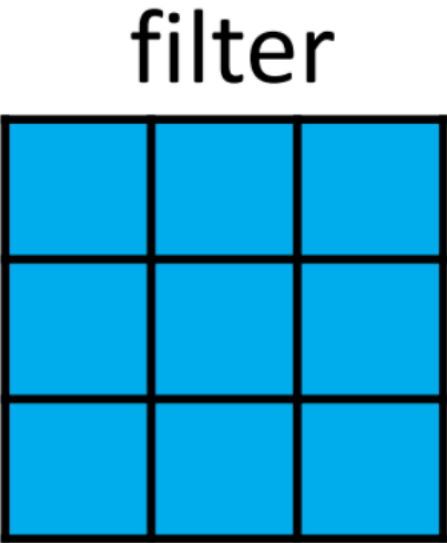


Convolution: padding

pad size

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



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

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

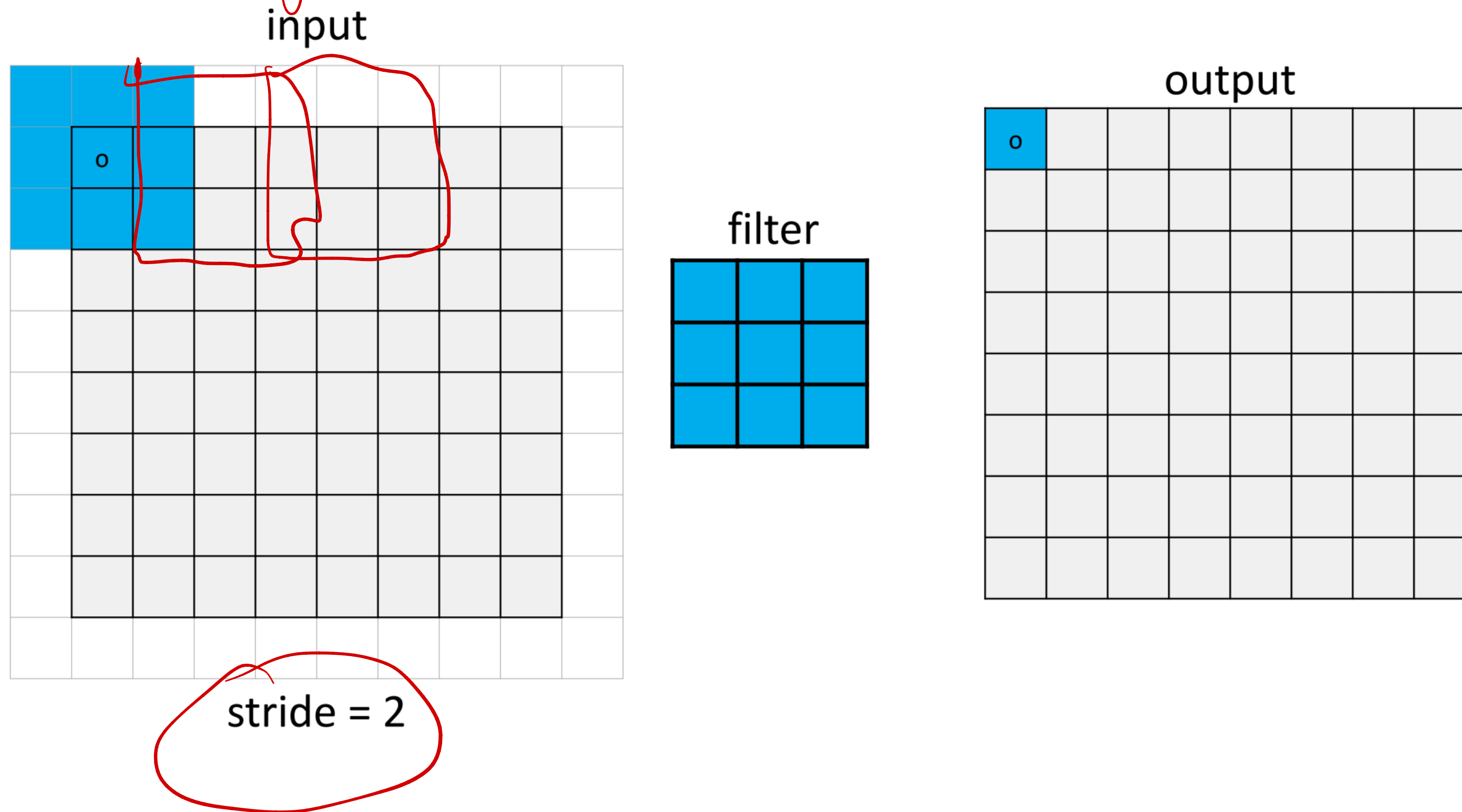
output: $H \times W = 8 \times 8$

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

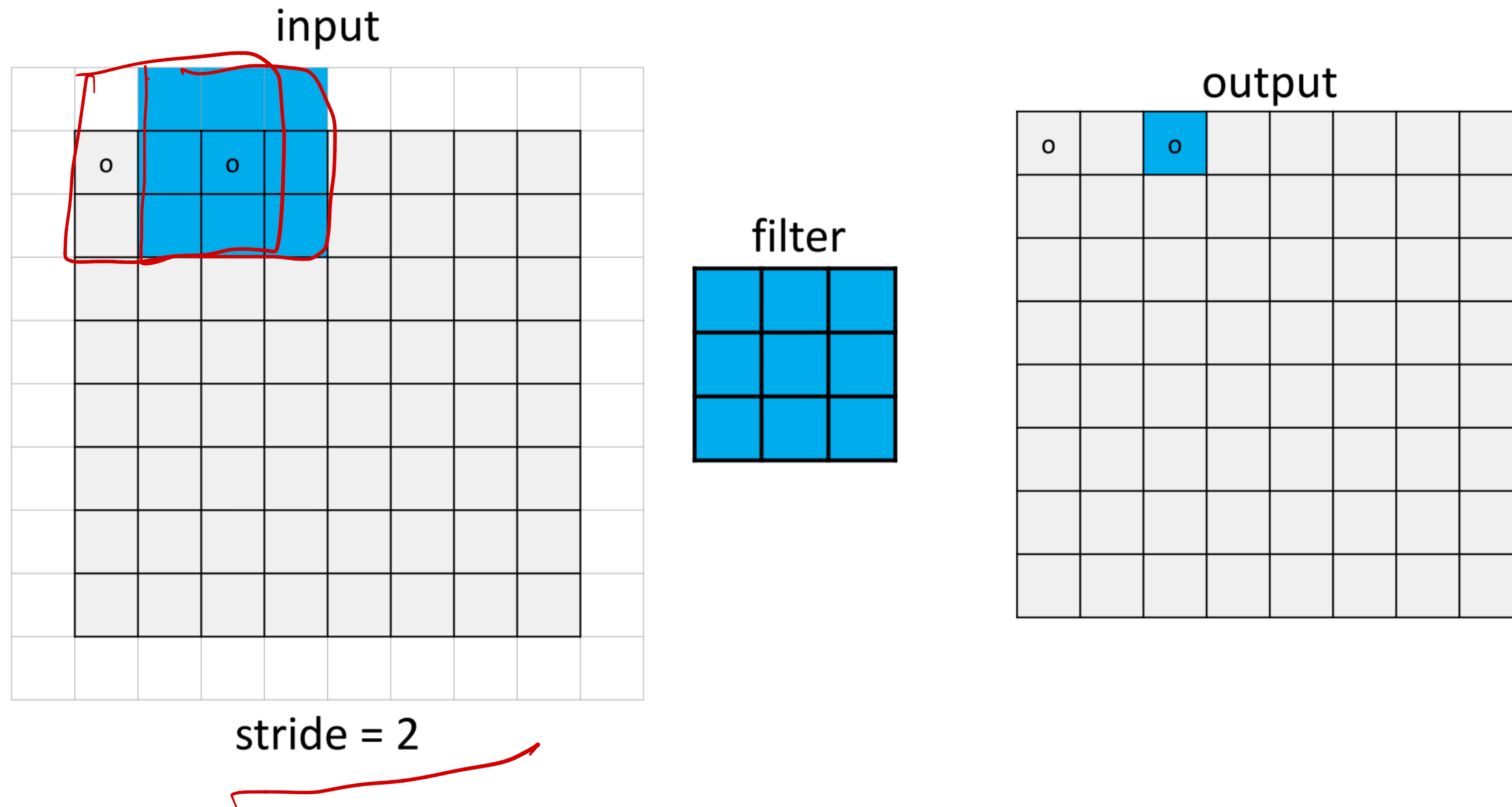
$$H_{\text{out}} = H_{\text{in}} + 2\text{pad}_h - K_h + 1$$

Convolution: stride

sliding step size



Convolution: stride



Convolution: stride

4x4

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

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

filter

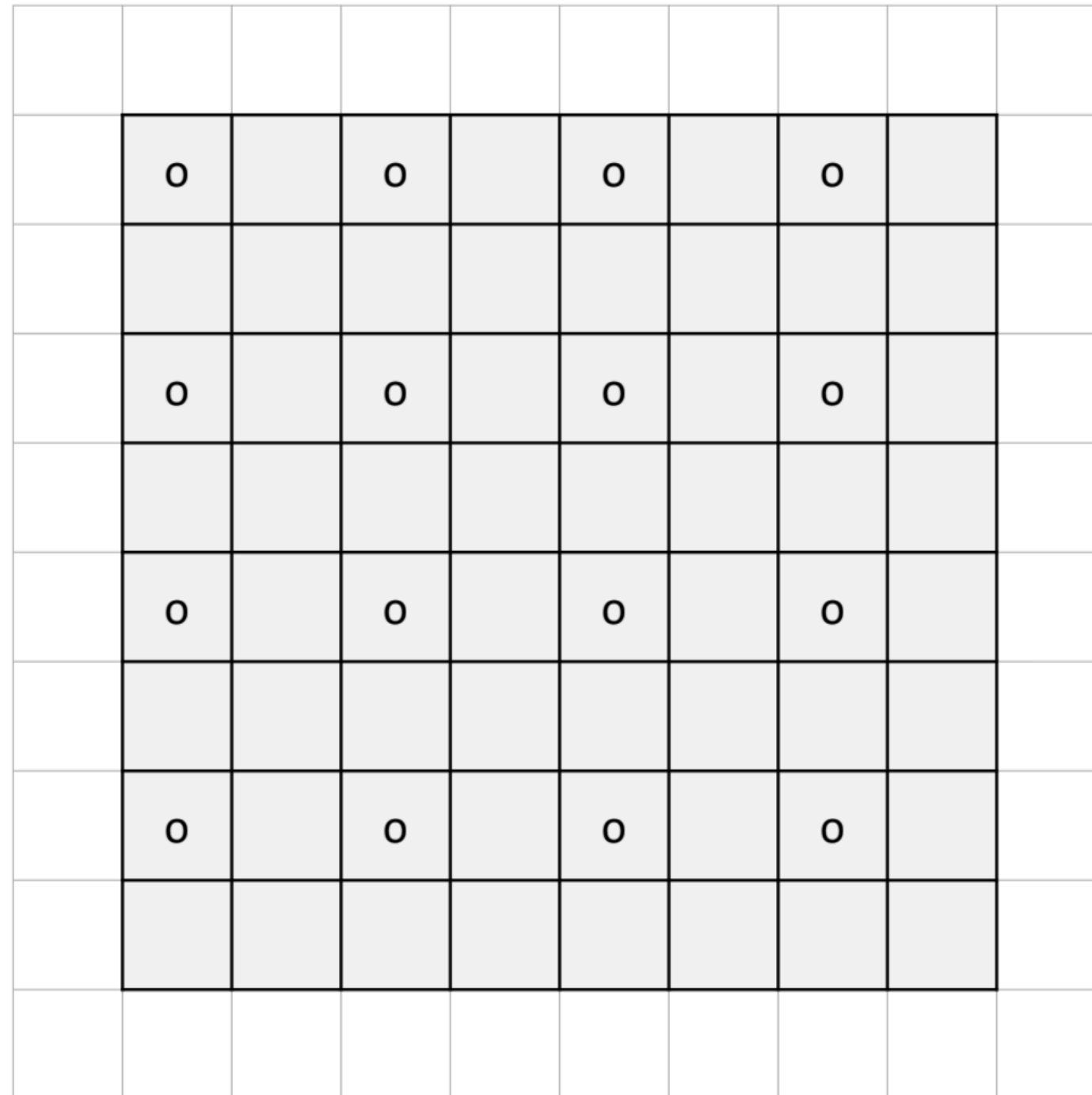
output: $H \times W = 4 \times 4$

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

stride = 2

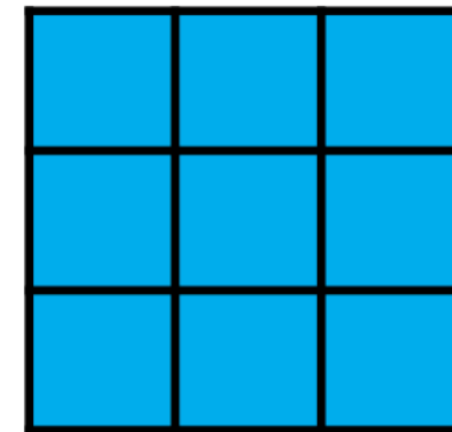
Convolution: stride

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

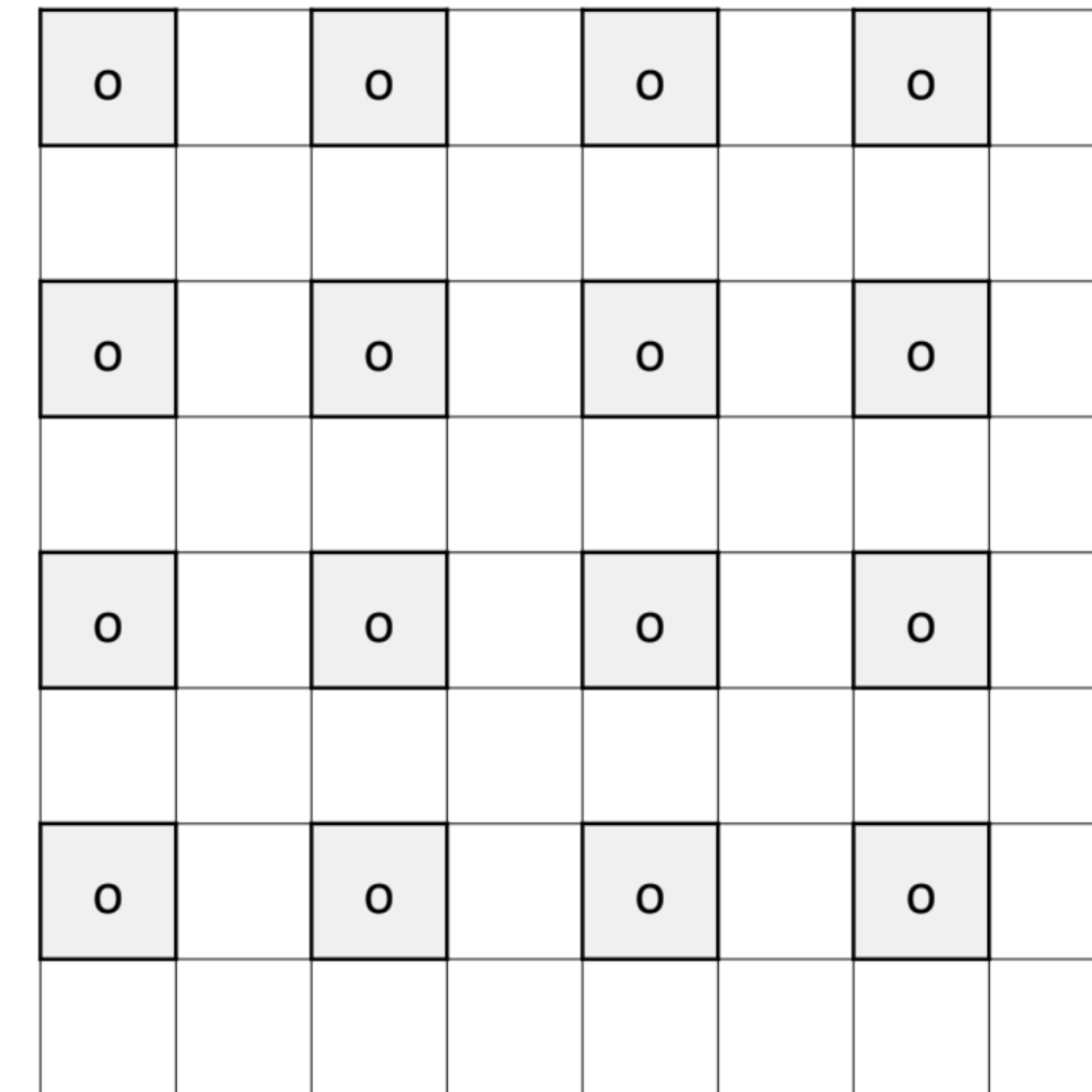


stride = 2

filter



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

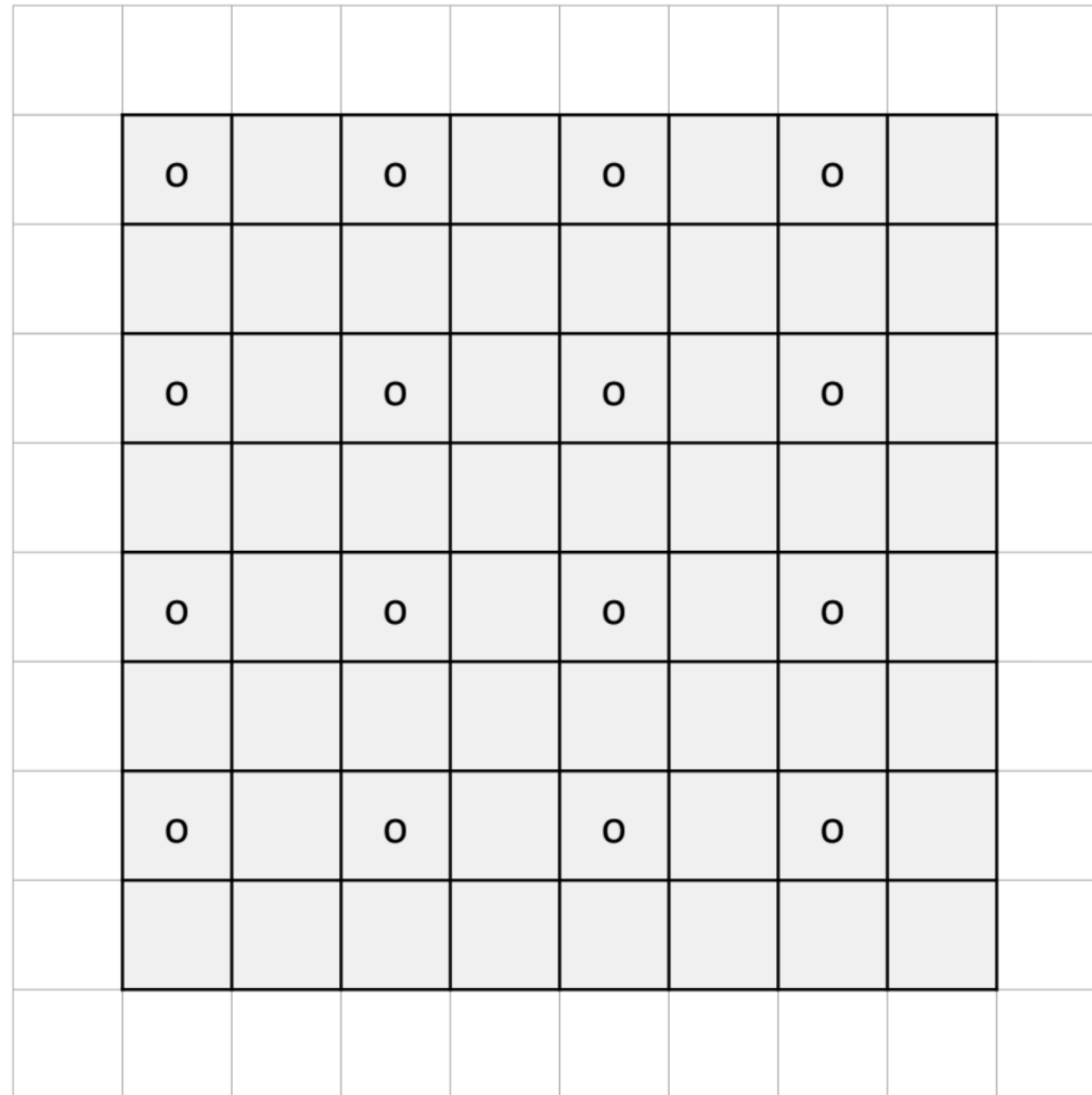


- reduces feature map size
- compress and abstract

input: 2048×2048
↓
 512×512

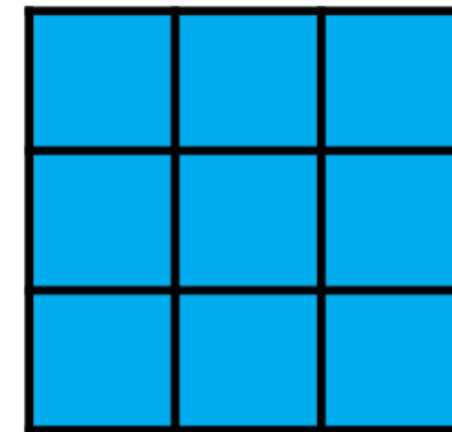
Convolution: stride

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

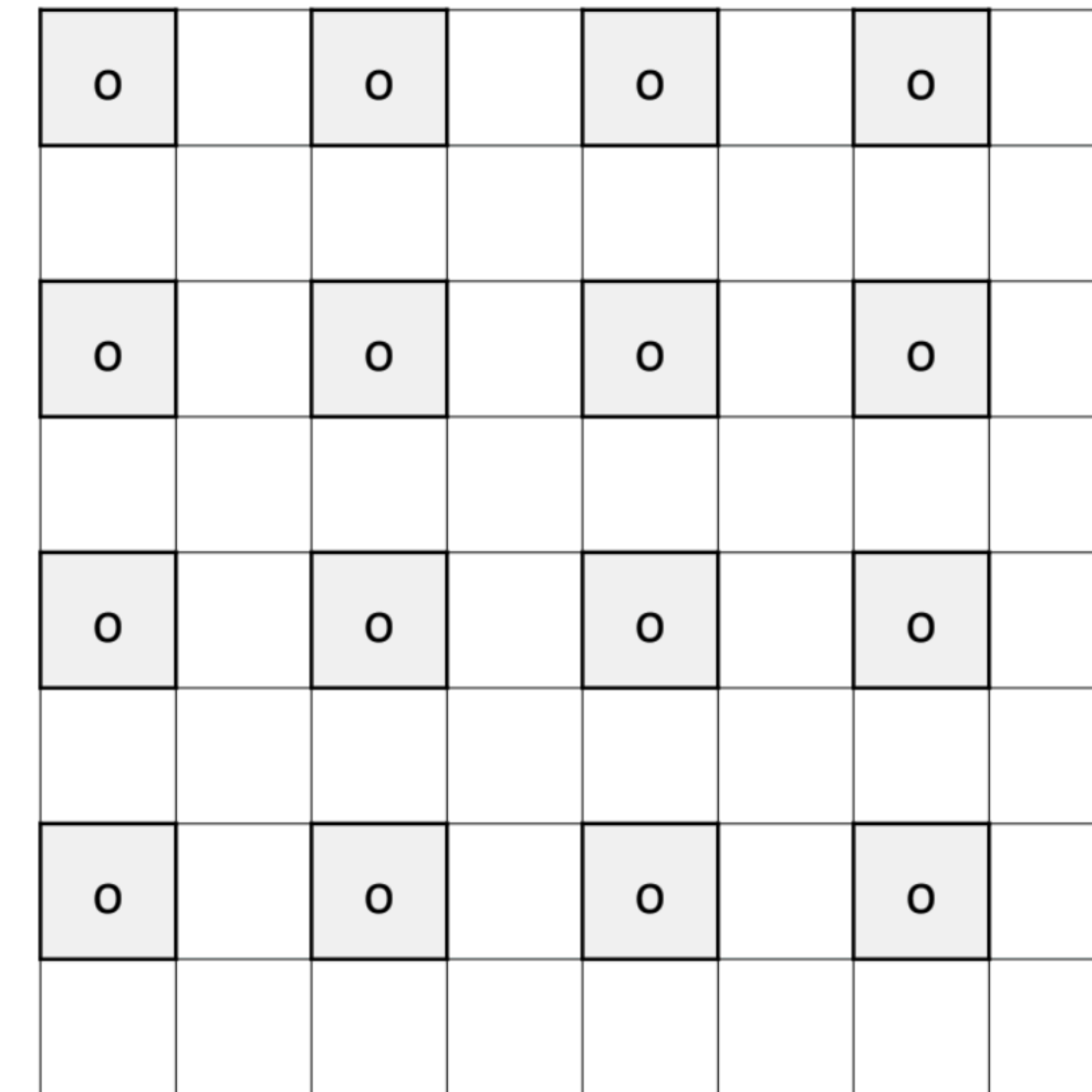


stride = 2

filter



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

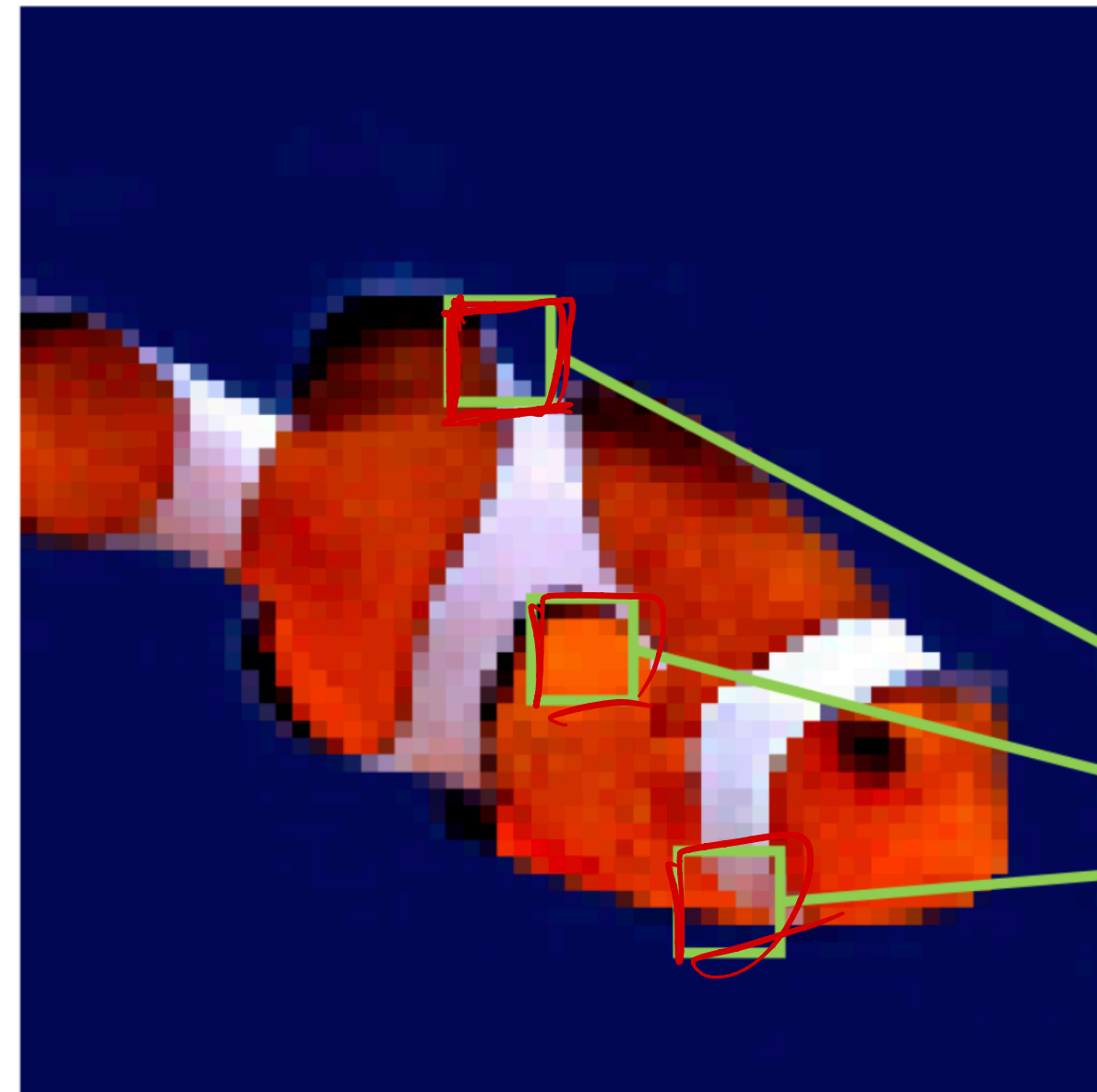


- reduces feature map size
- compress and abstract

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

Convolution: translation-invariance

- Process each window in the same way



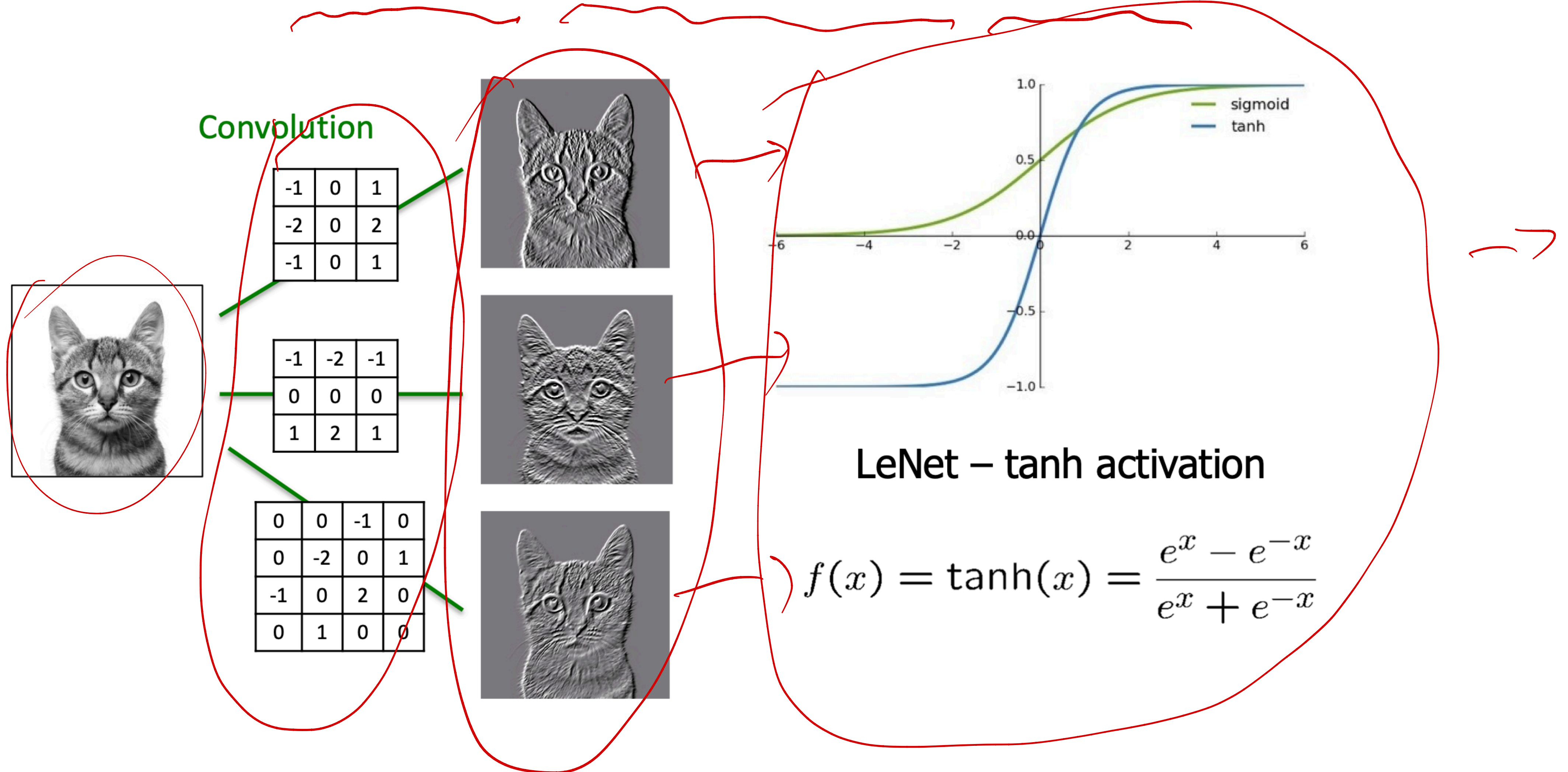
MLP

$w \times t_b$
↓

apply the same weights regardless of the window location

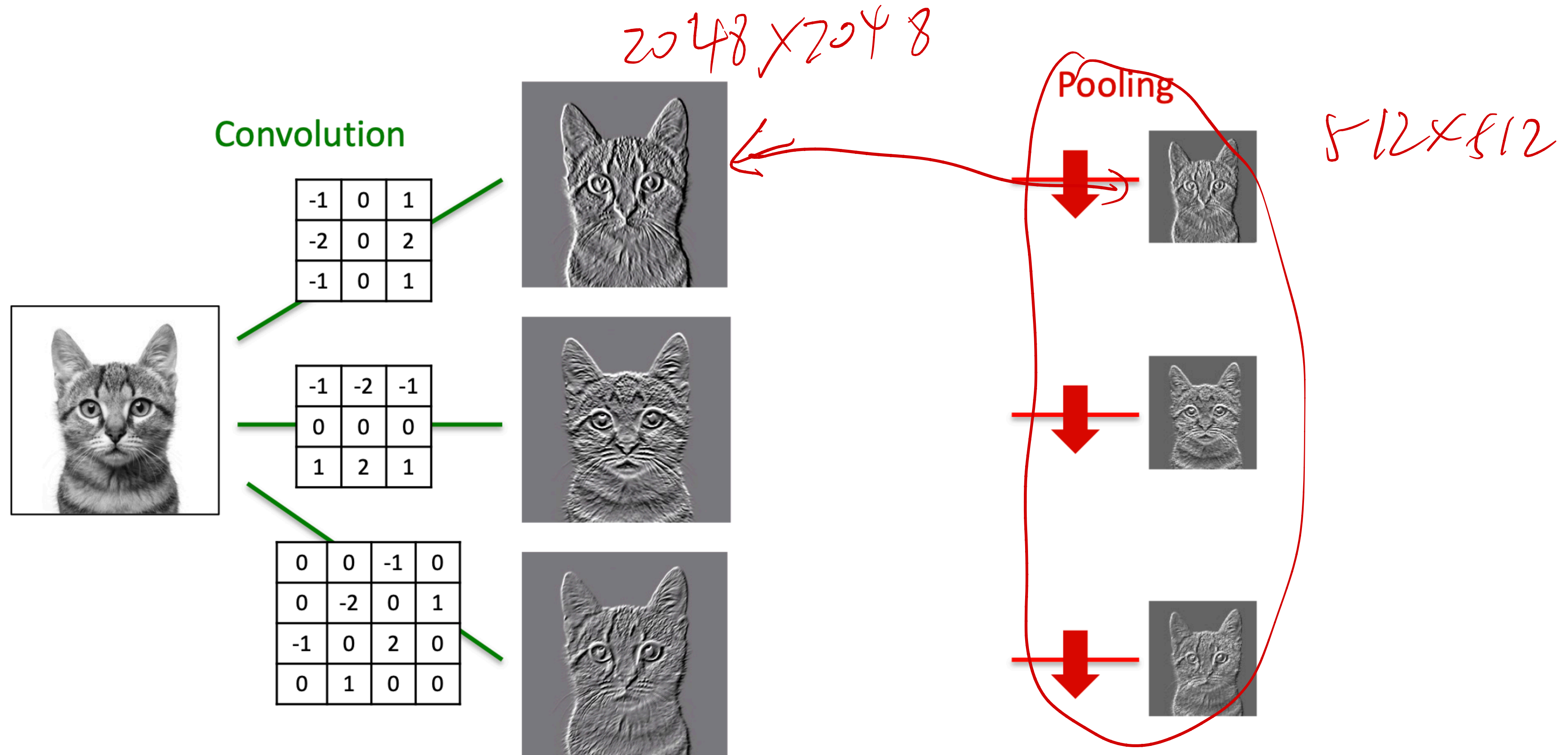
Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling



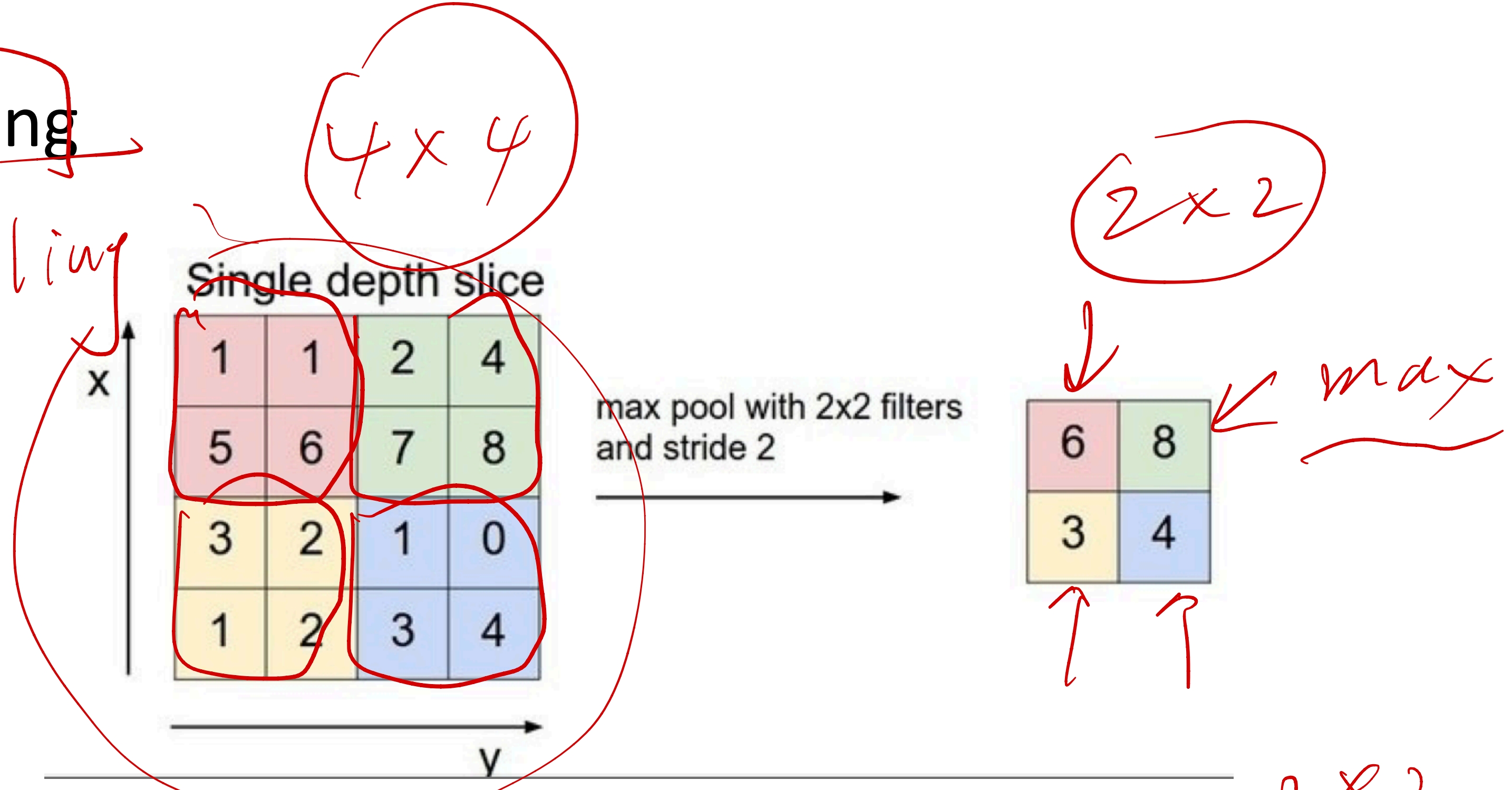
Deep Convolutional Networks

[Convolution + Nonlinear activation] + Pooling



Pooling = Down-sampling ^{do not} like cat

Max pooling
mean pooling



1024 x 1024
 512 x 512

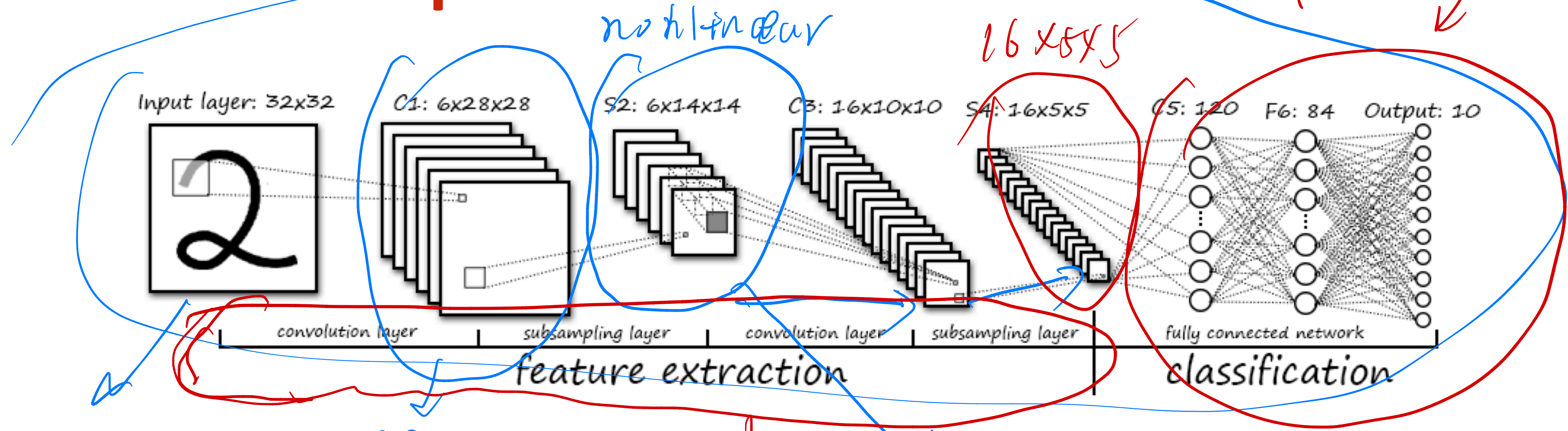
mean pooling

6.5	

2x2

Deep Convolutional Networks

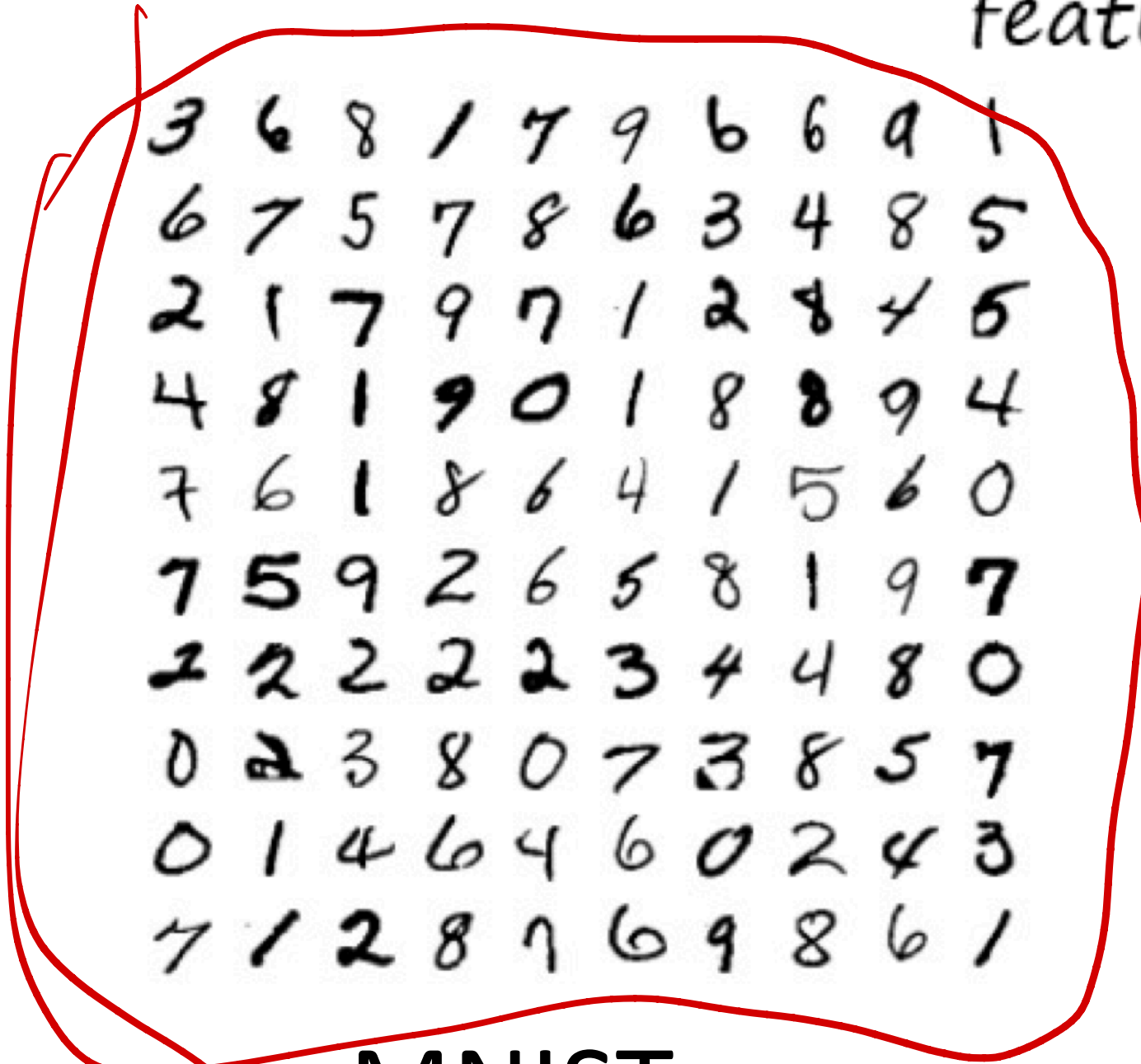
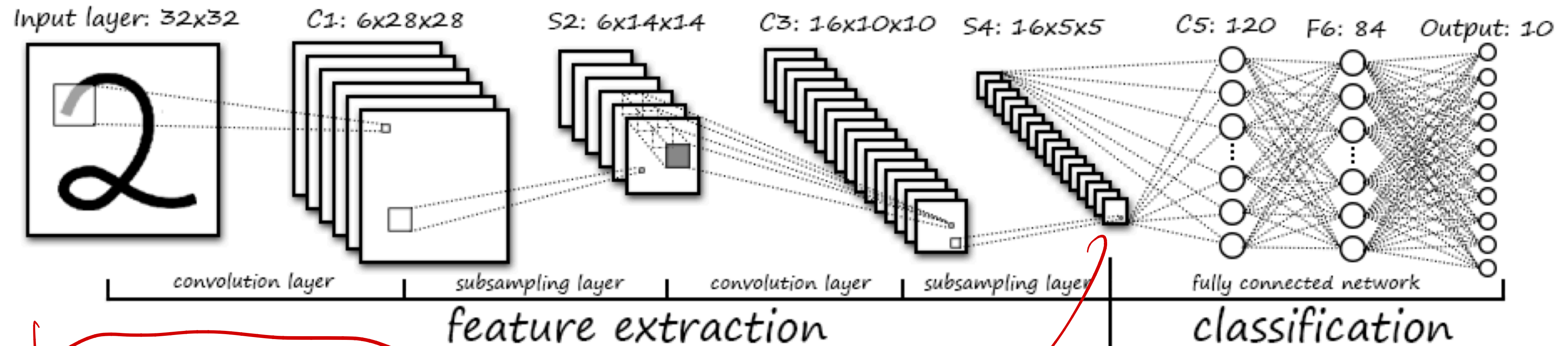
2-layer MLP



32x32
6x28x28
6x14x14
output channel size

[1] LeNet 5, LeCun et al. 1998

Deep Convolutional Networks



MNIST

60,000 original dataset
Test error: 0.95%

[1] LeNet 5, LeCun et al. 1998

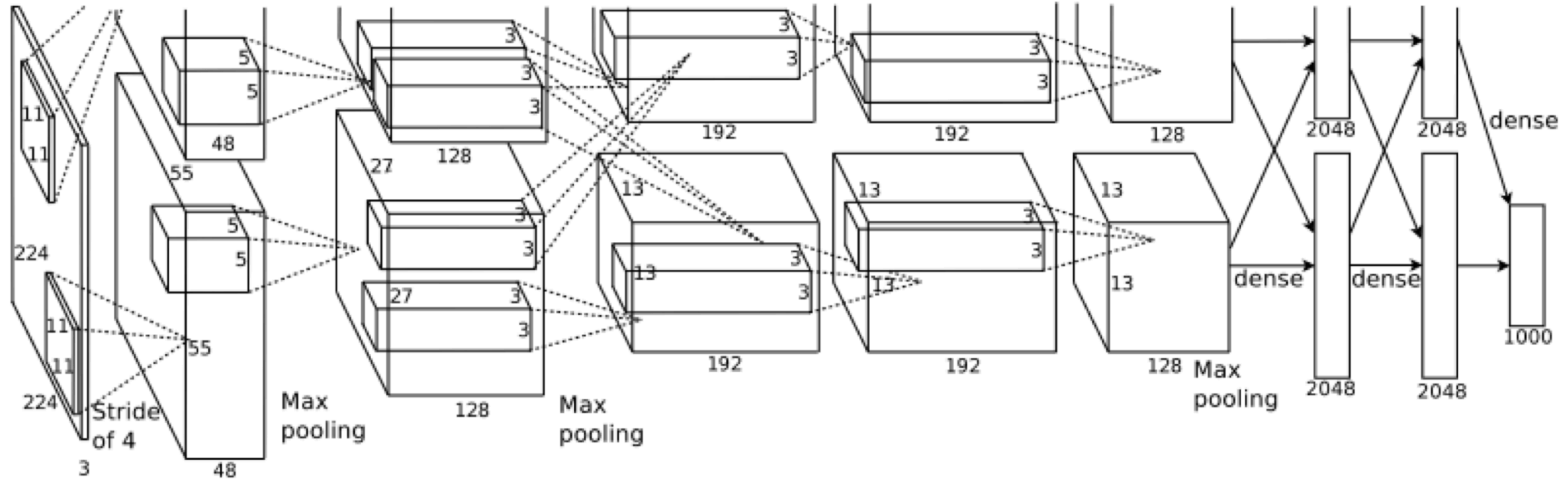
Misclassified examples on MNIST

True label -> Predicted label



Alex Net

milestone



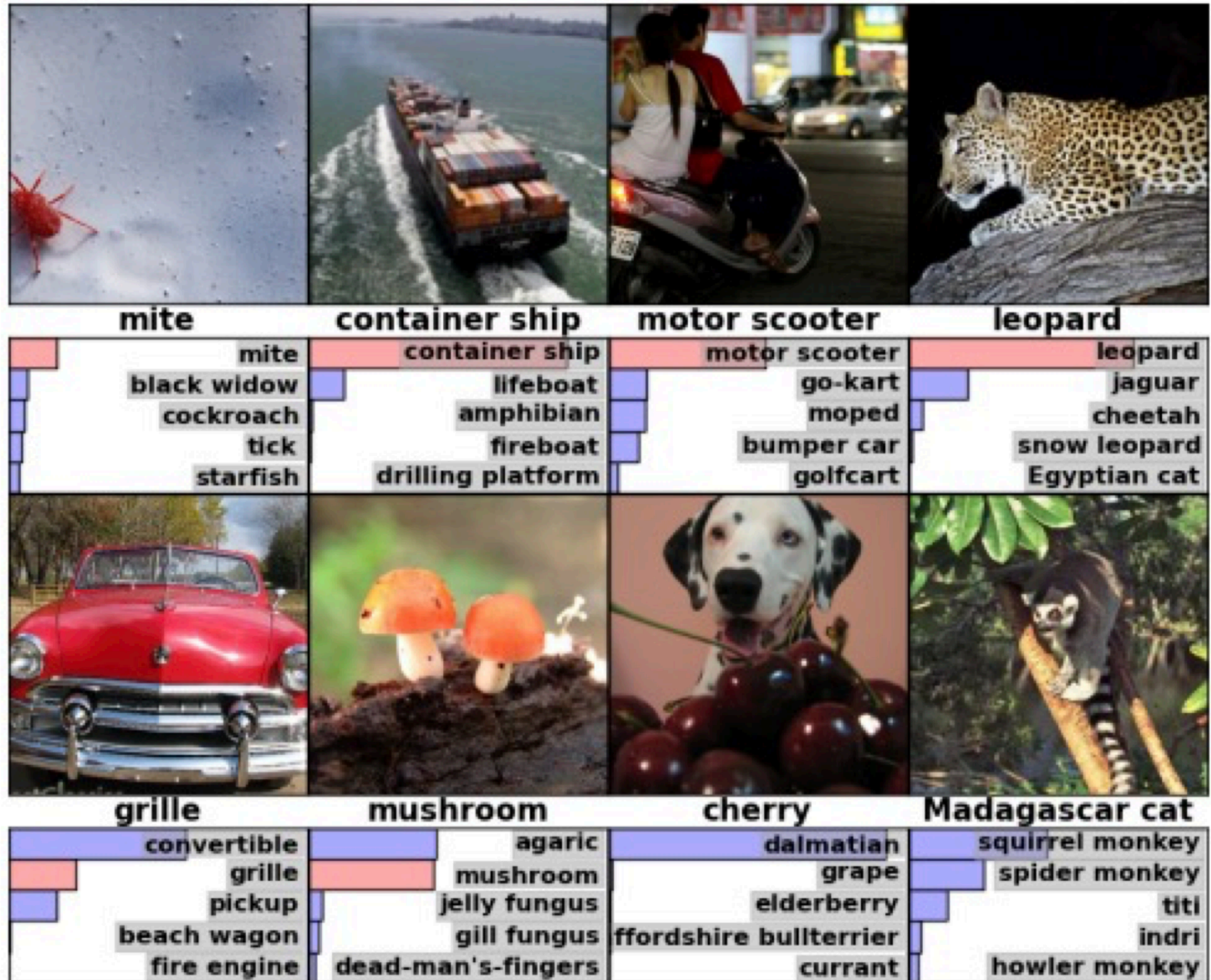
[1] Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton. ImageNet Classification with Deep Convolutional Neural Networks. NeurIPS 2012.

scale is large

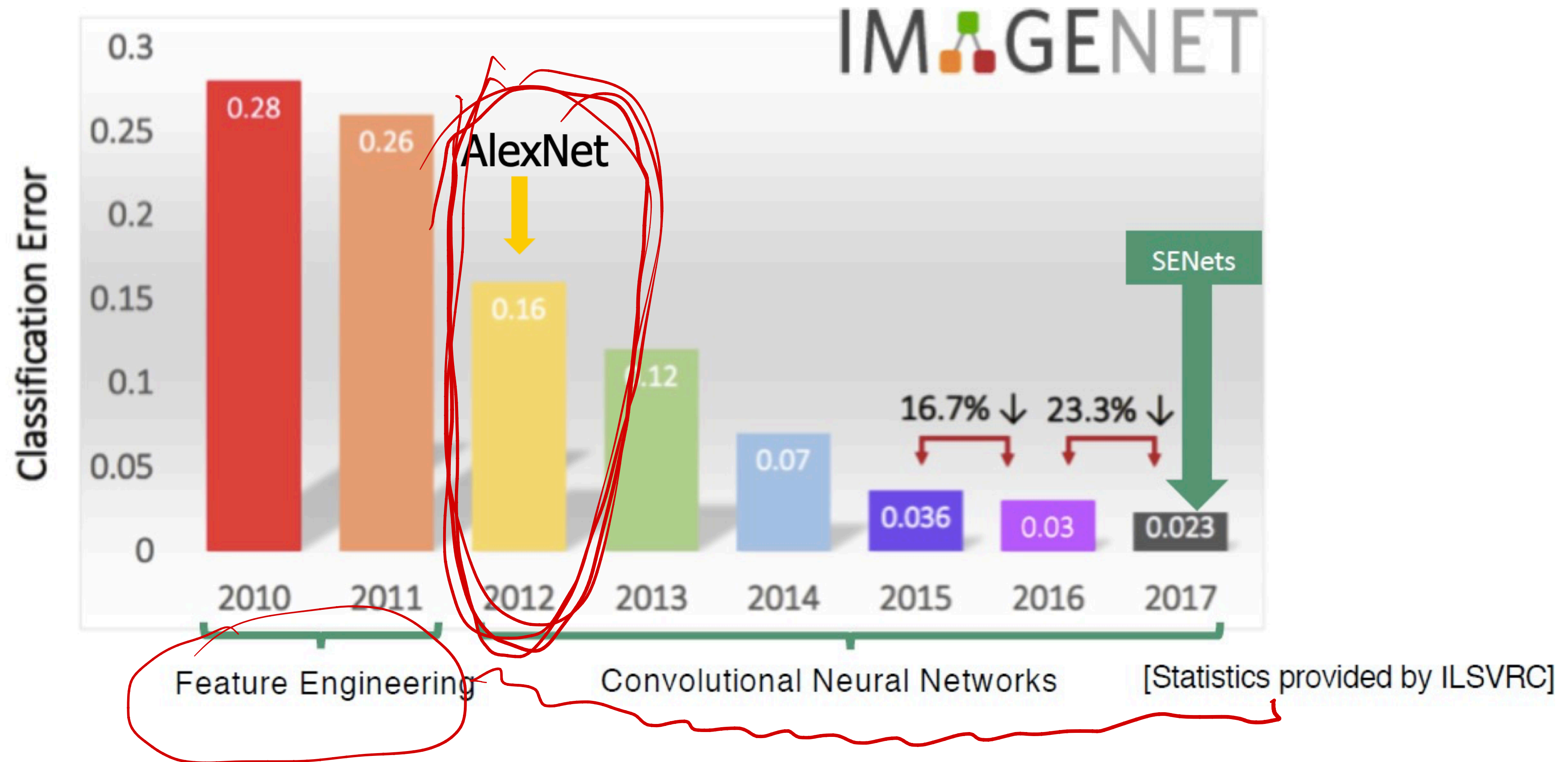
ImageNet

- ❑ 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

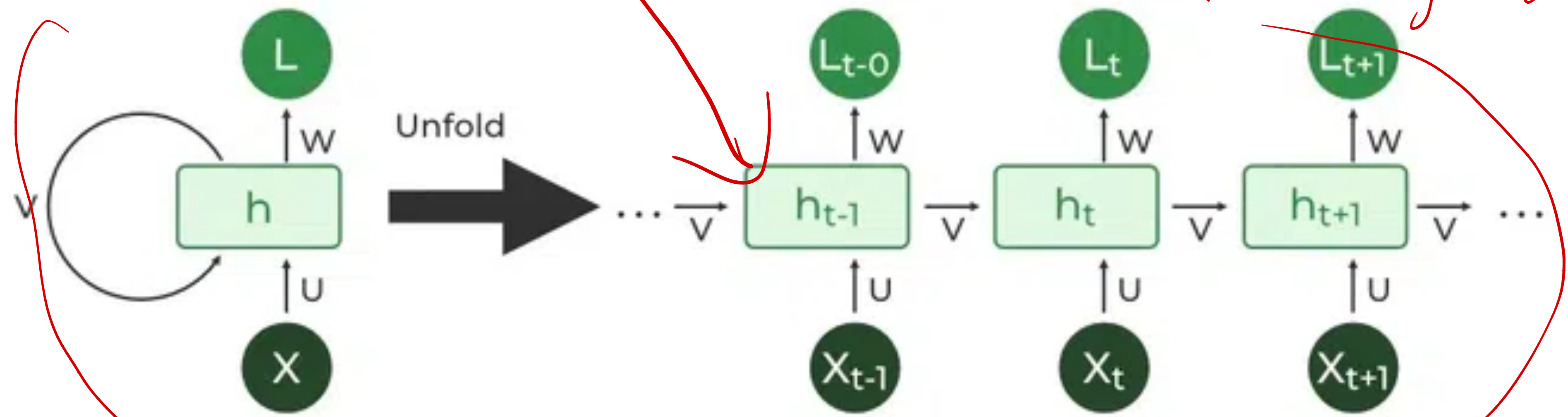


ImageNet Results



Recurrent Neural Networks (RNNs)

CNN for image
RNN for language

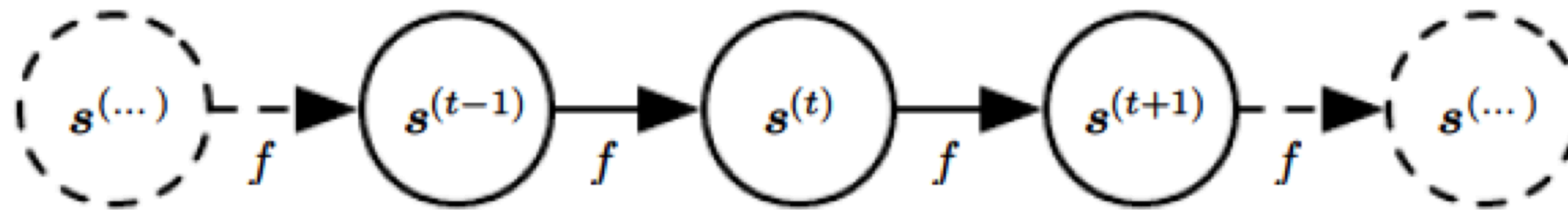


sequence

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



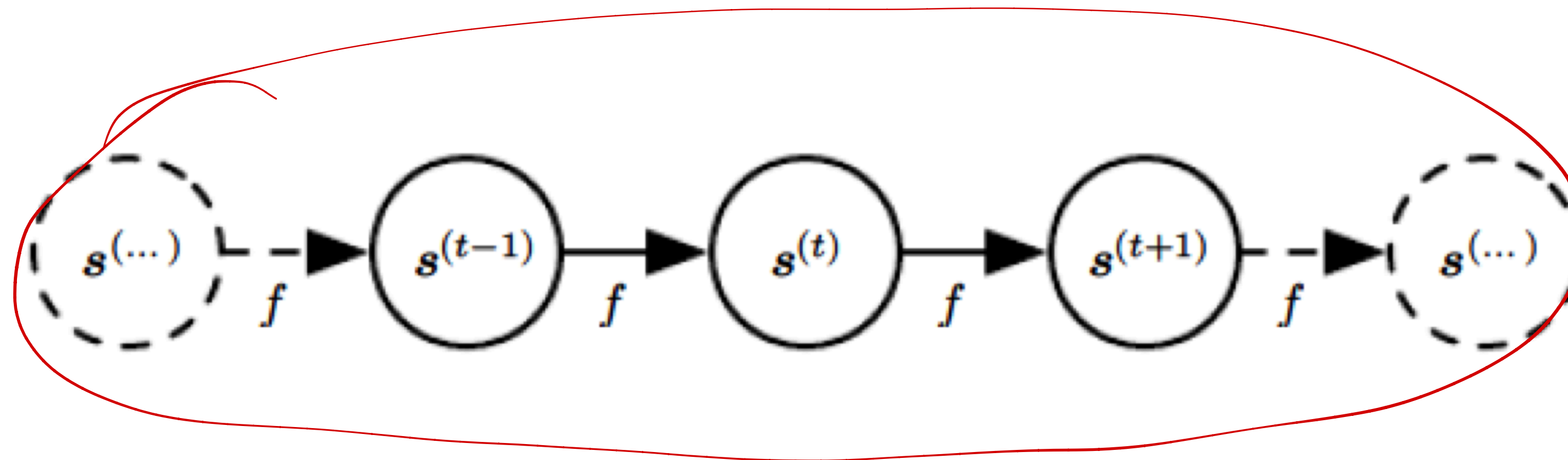
$$s^{(t+1)} = f(s^{(t)}; \theta)$$

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

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

↓ parameters

Computation Graph

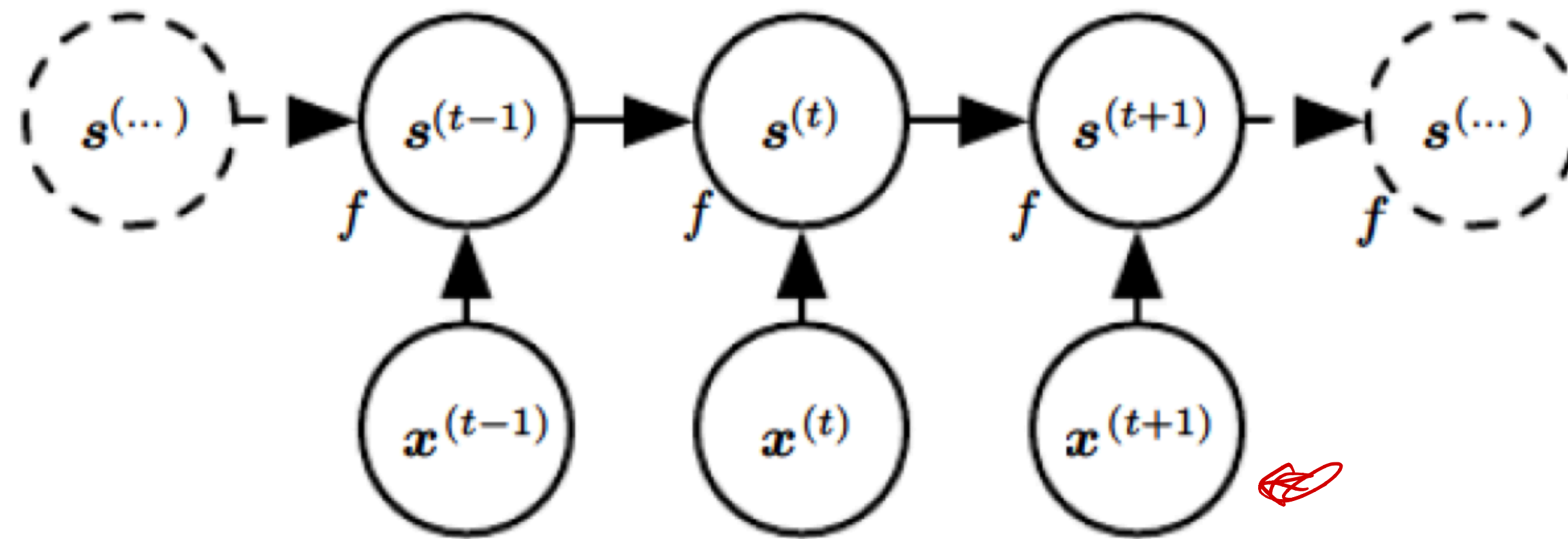


$$s^{(t+1)} = f(s^{(t)}; \theta)$$

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

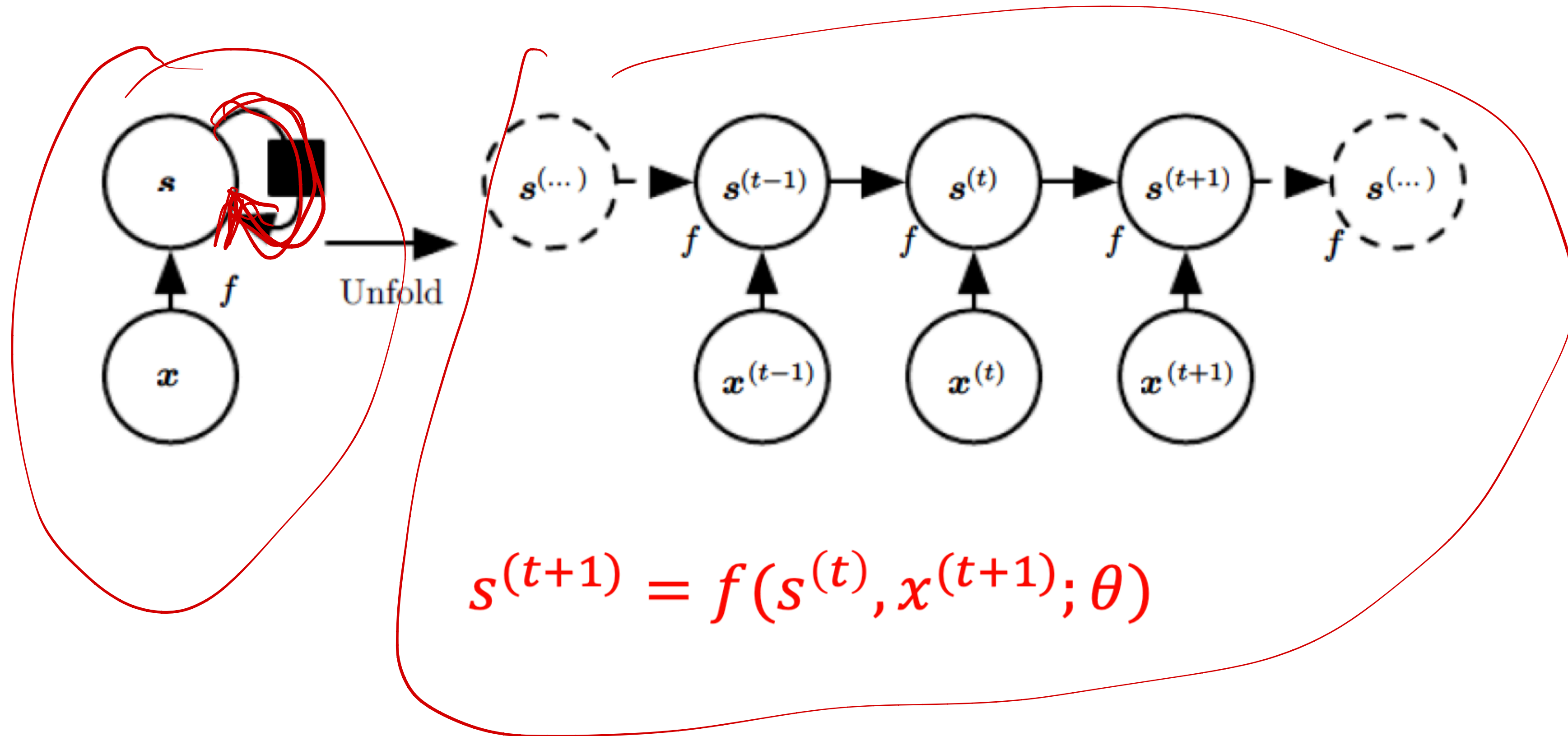
Like Markov model, but here $s^{(t+1)}$ is deterministic given $s^{(t)}$

Computation Graph

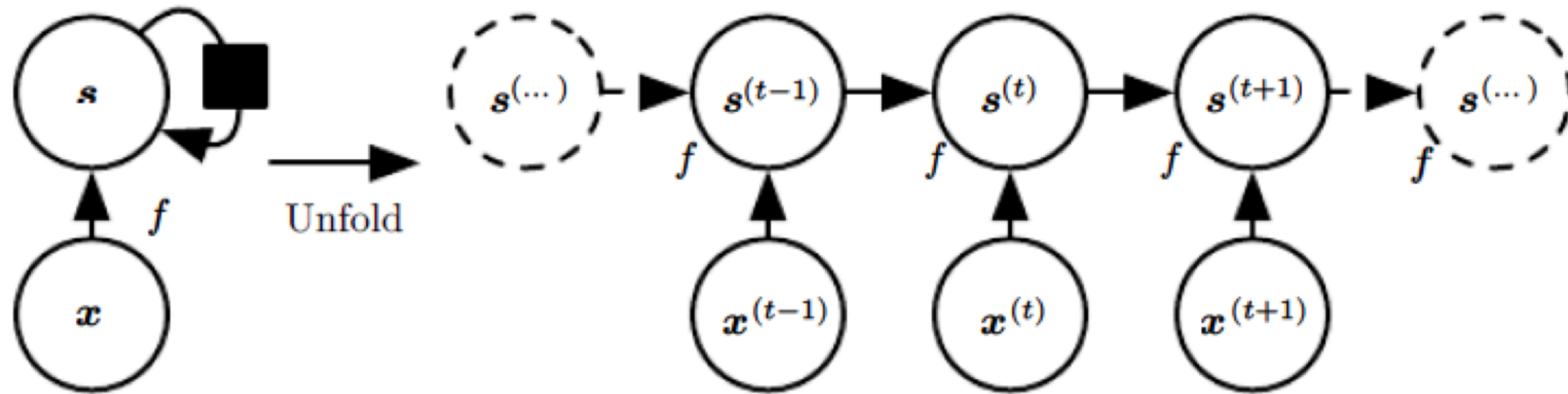


$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Compact view



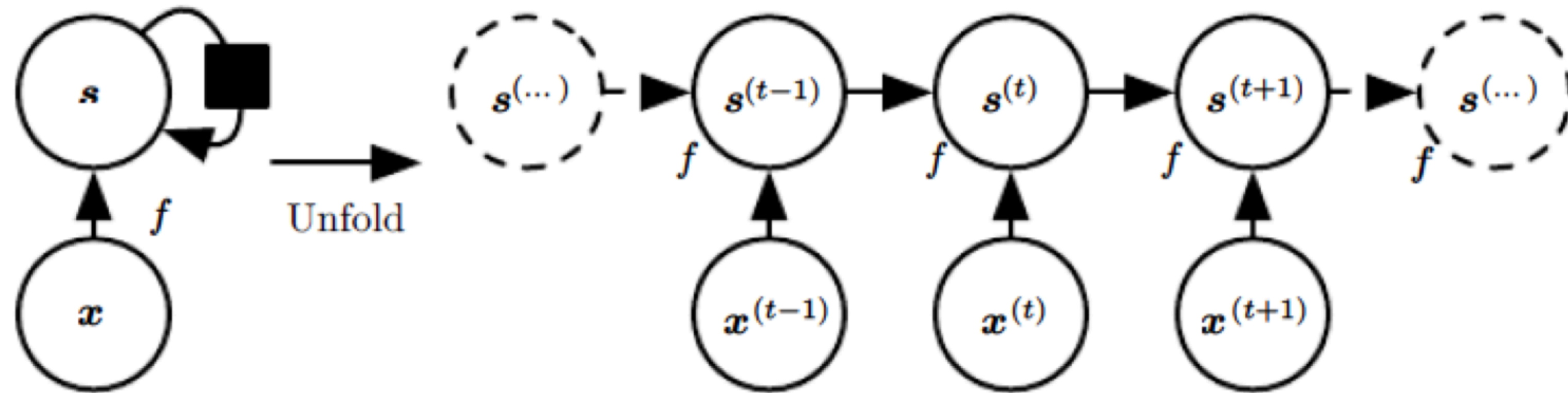
Compact view



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

Key: the same f and θ
for all time steps

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

ω

θ f

Recurrent Neural Networks

Recurrent Neural Networks

- Use **the same** computational function and parameters across different time steps of the sequence

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

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

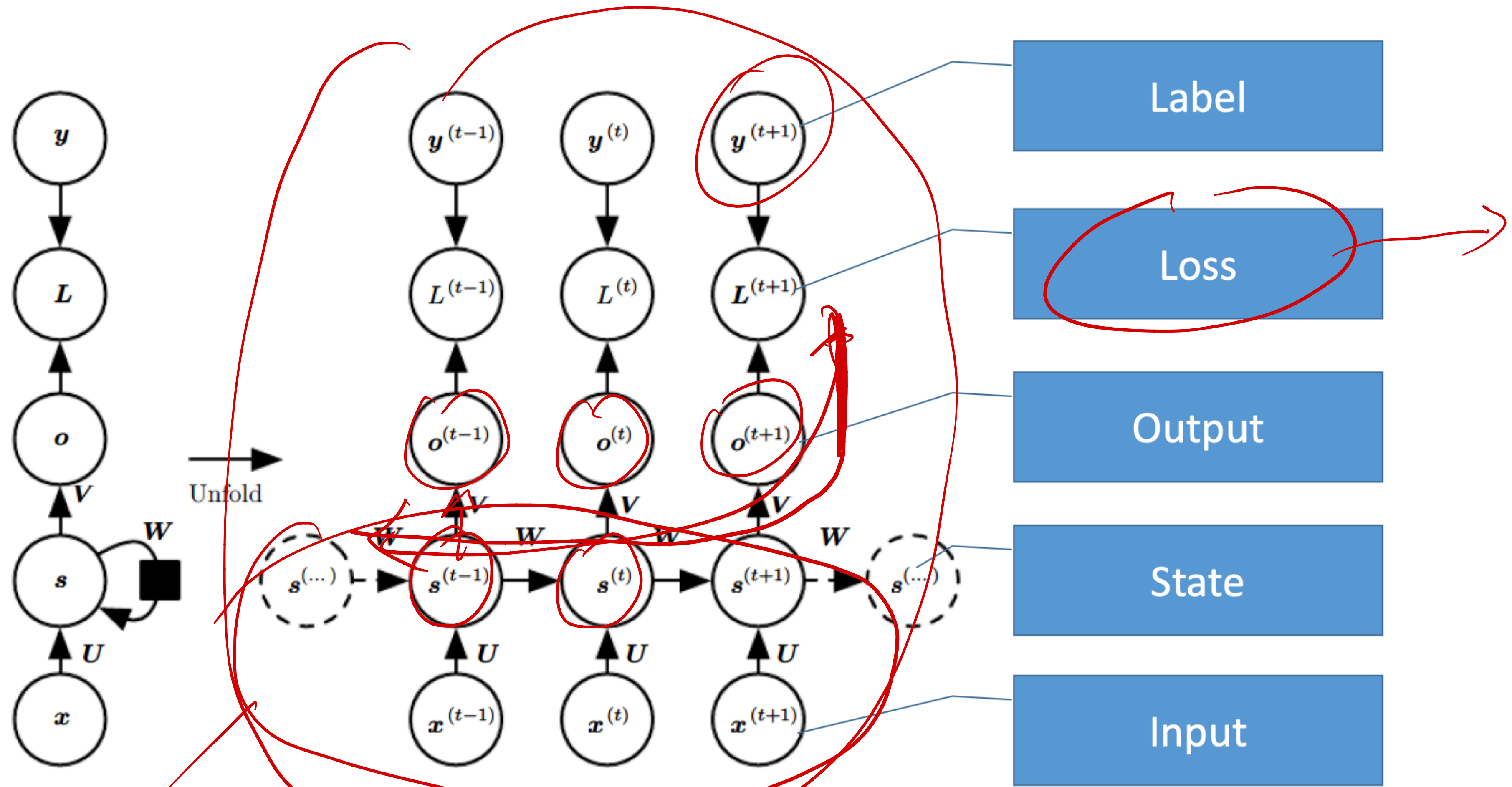
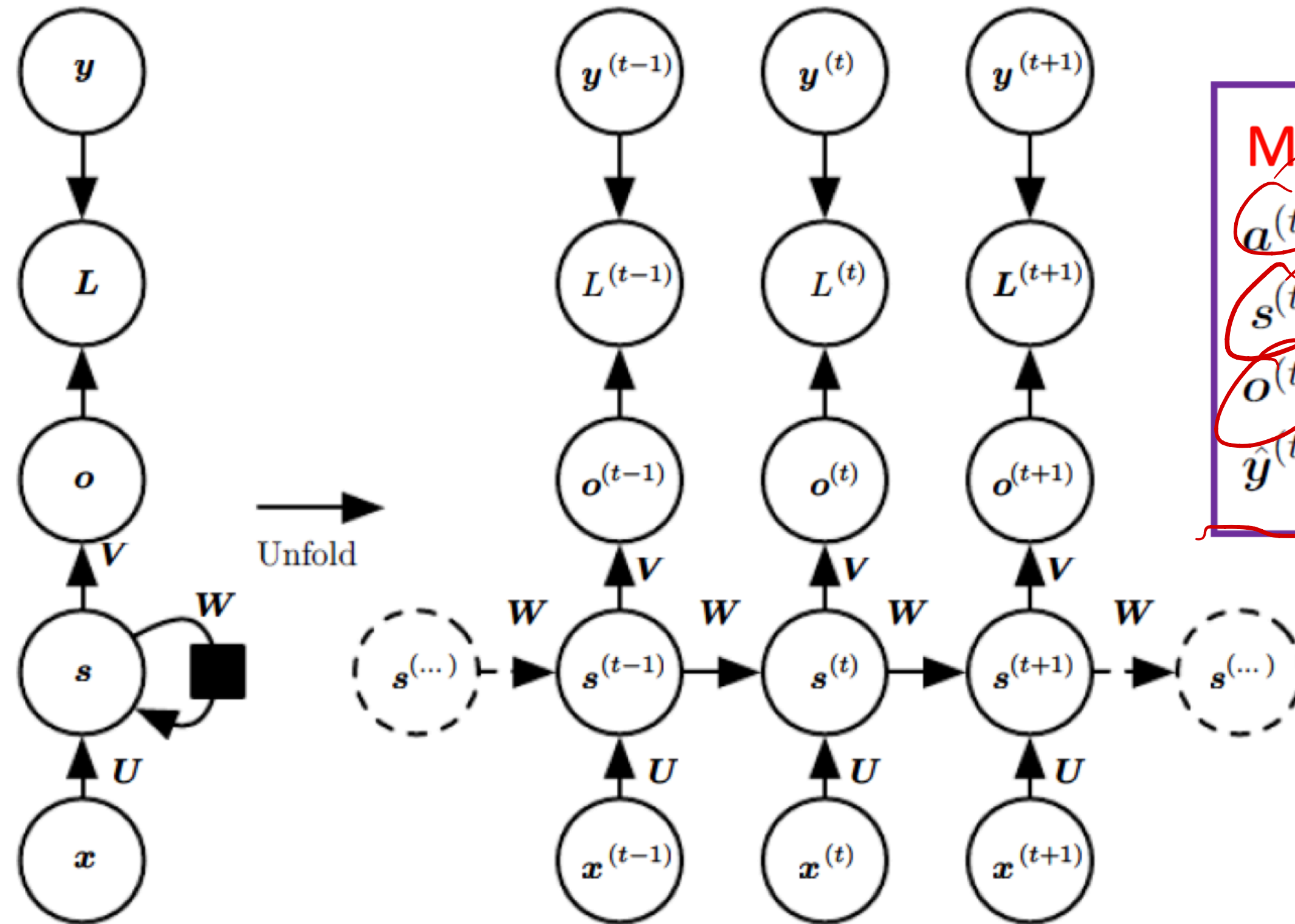


Figure from *Deep Learning*, by Goodfellow, Bengio and Courville

RNN

Recurrent Neural Networks



Math formula:

$$a^{(t)} = b + Ws^{(t-1)} + Ux^{(t)}$$

$$s^{(t)} = \tanh(a^{(t)}) \rightarrow \text{non-linear}$$

$$o^{(t)} = c + Vs^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

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

Sequence-to-Sequence Learning

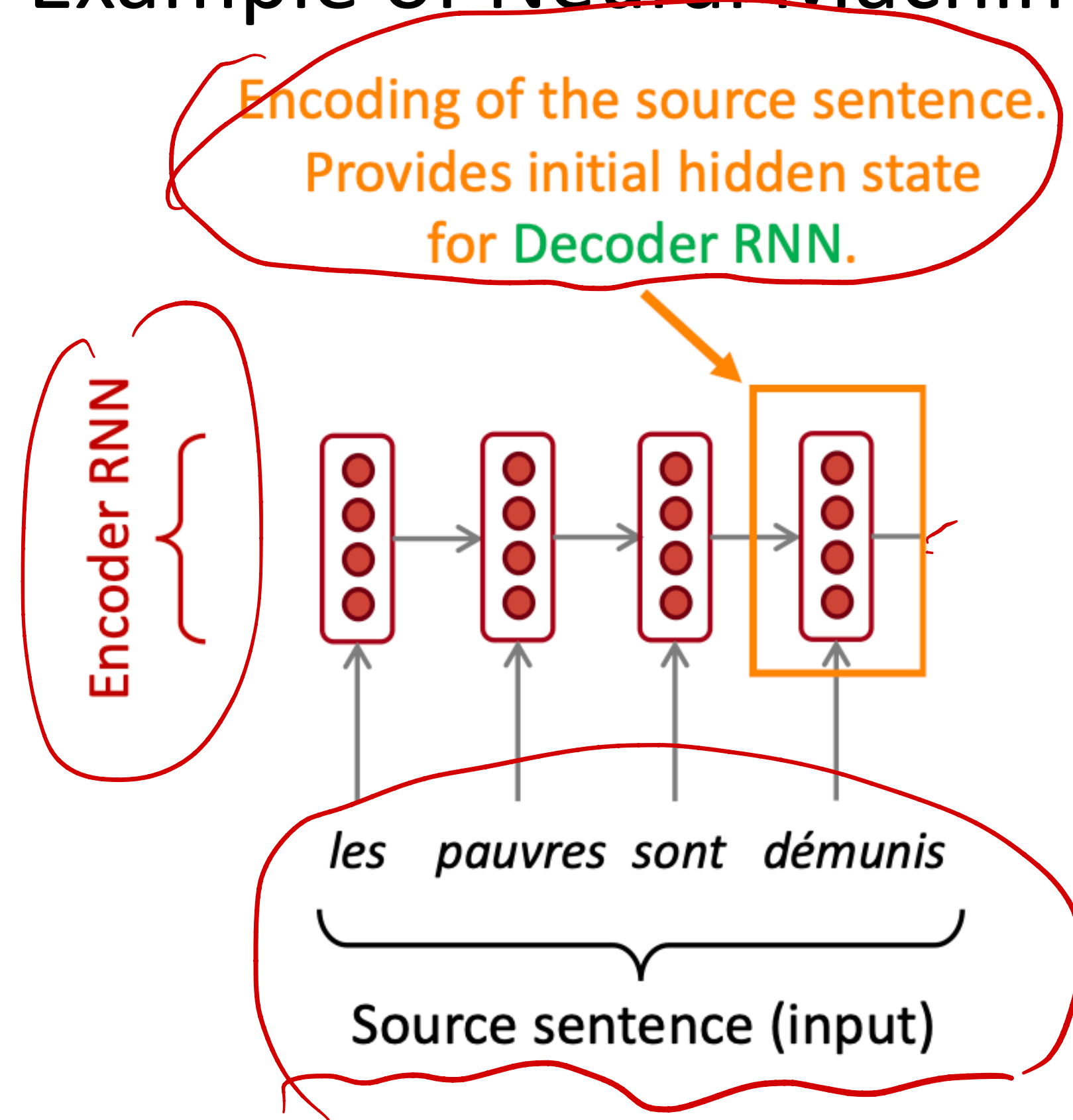


Example of Neural Machine Translation



Sequence-to-Sequence Learning

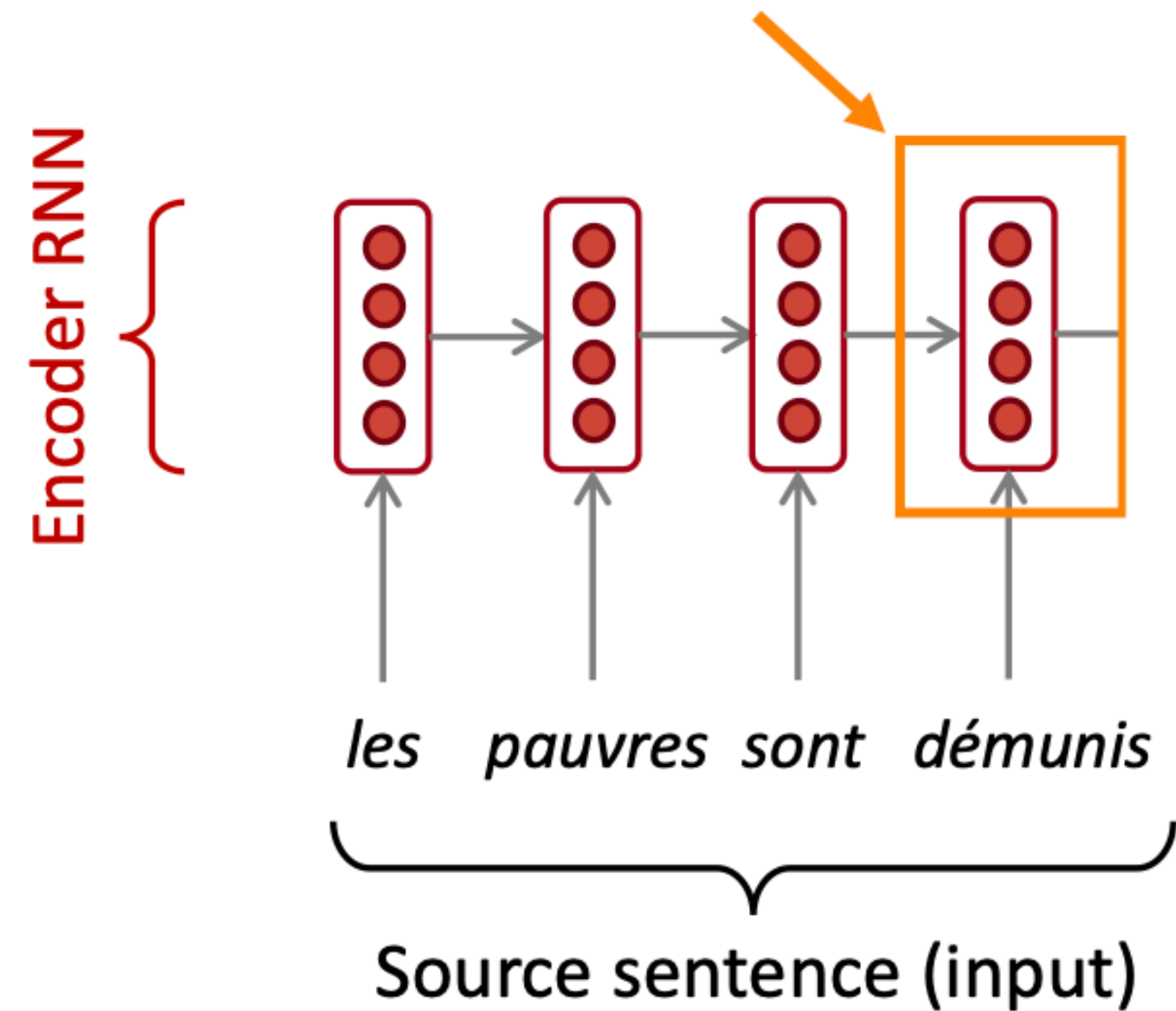
Example of Neural Machine Translation



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.

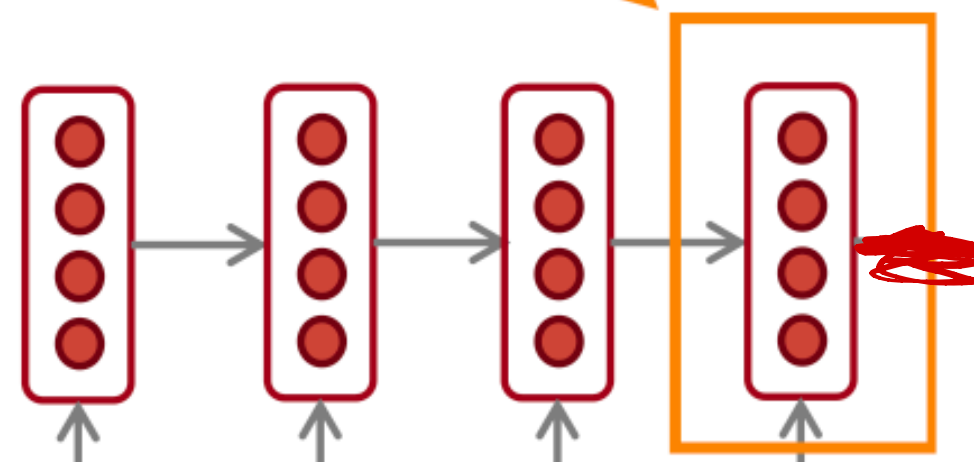
Sequence-to-Sequence Learning

Example of Neural Machine Translation

$$s^t = f(s^{t-1}, x)$$

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

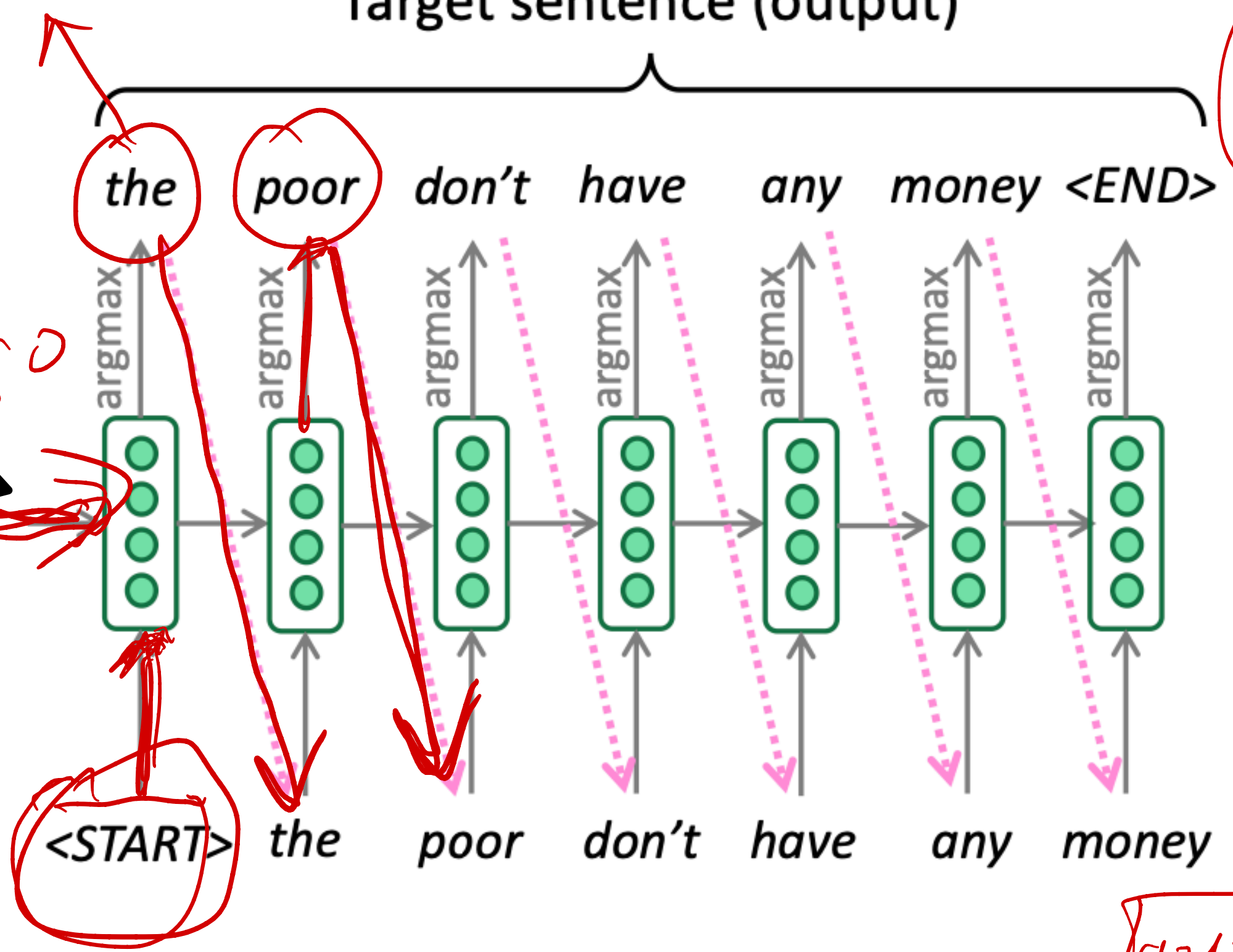
Encoder RNN



les pauvres sont démunis
Source sentence (input)

Encoder RNN produces an encoding of the source sentence.

Target sentence (output)



s^0 ?

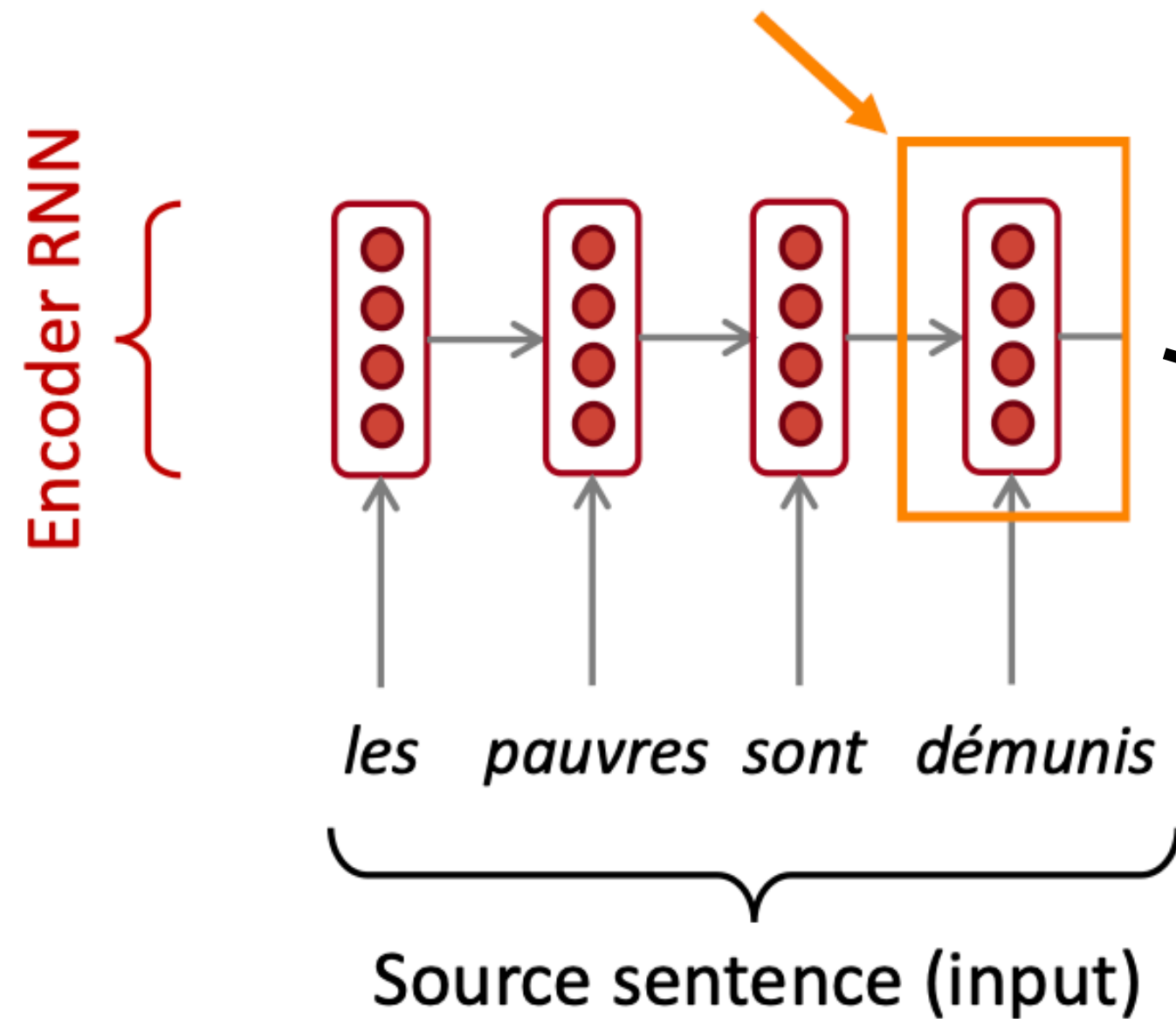
Decoder RNN

autoregressive

Sequence-to-Sequence Learning

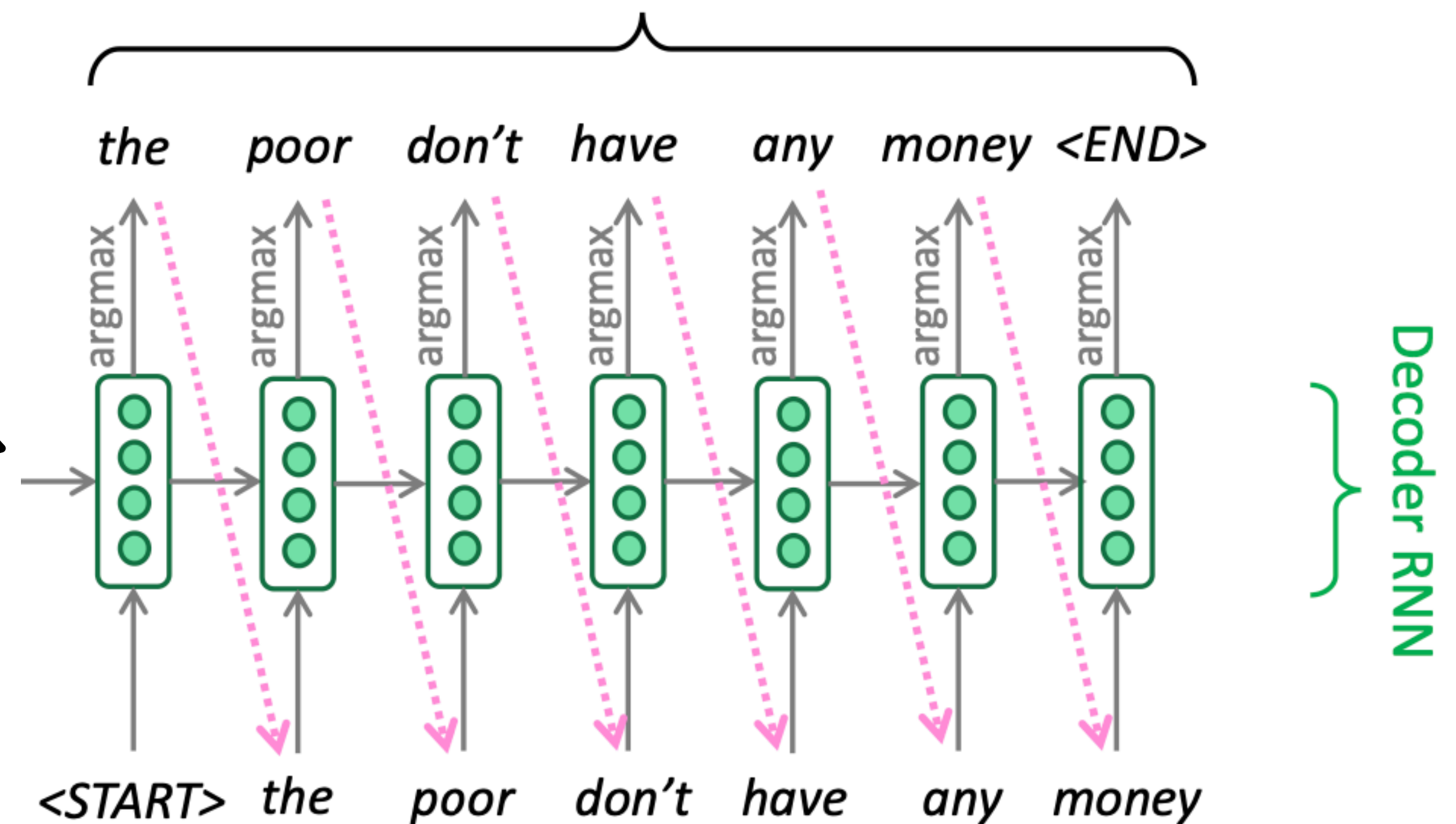
Example of Neural Machine Translation

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



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

large

gradient

explodes
vanish



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

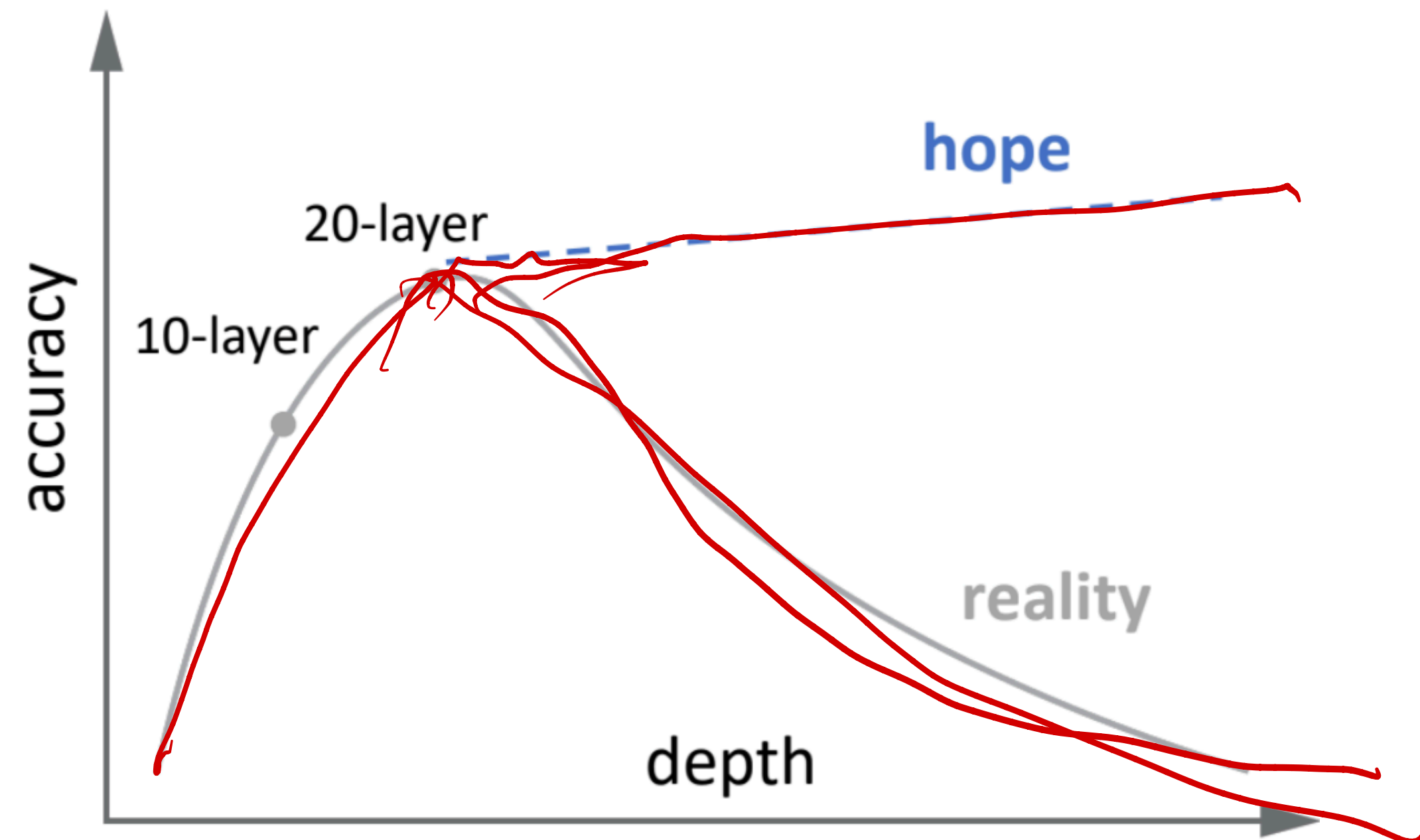
Normalization modules

Deep residual learning



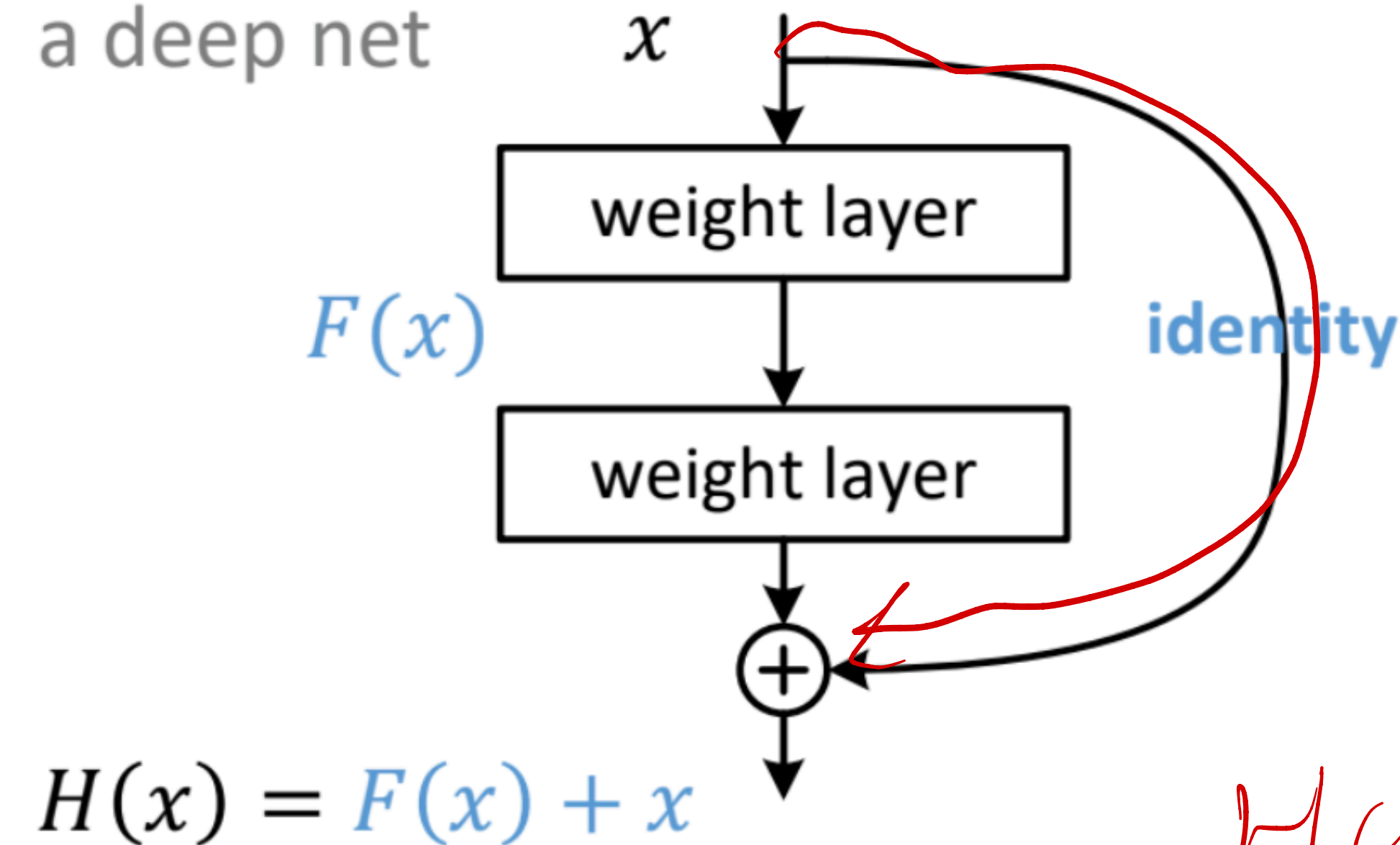
The Degradation Problem

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



Deep Residual Learning

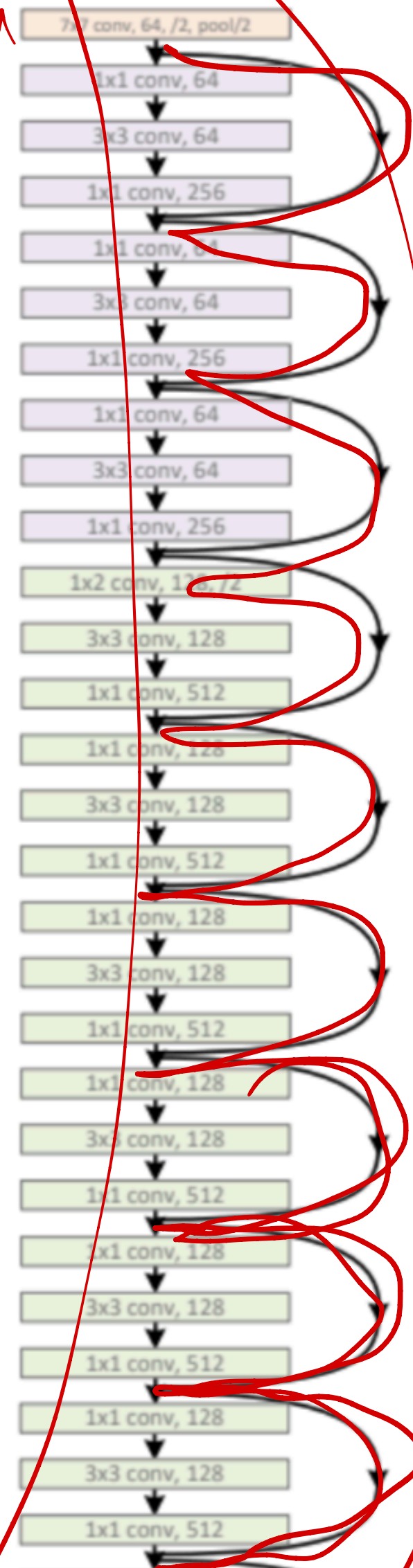
a subnet in
a deep net



$$H(x) = F(x) + x$$

Deep Residual Networks (ResNet)

ResNet-152



deep

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

Transformers

