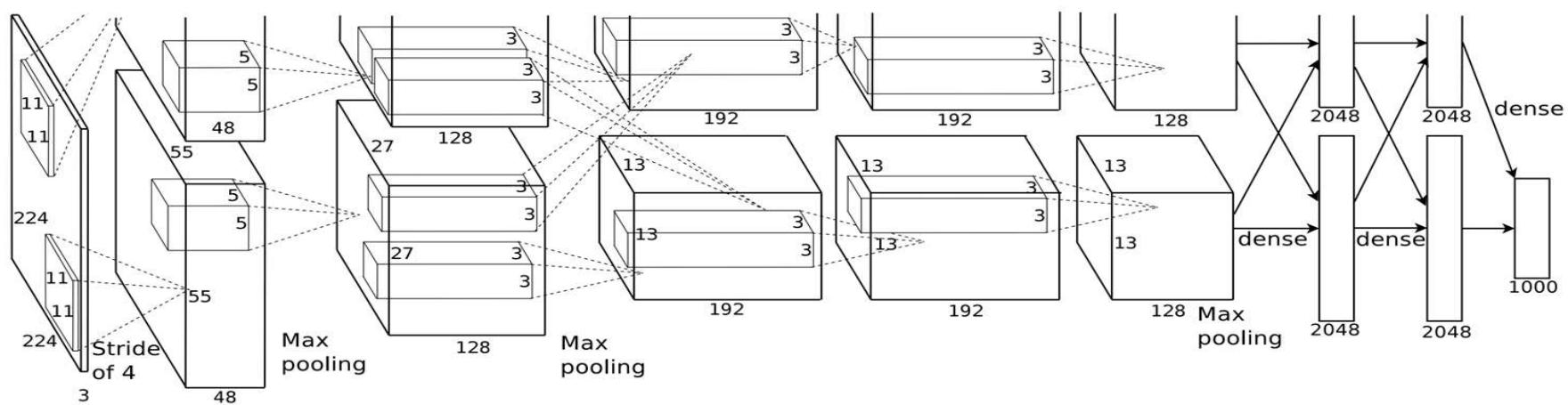


Deep NN



2018 Turing Award

— by ACM @ 2019/3/27



AWARDS & RECOGNITION

ACM Announces 2018 Turing Award Recipients ↗

ACM has named [Yoshua Bengio](#) ↗ of the University of Montreal, [Geoffrey Hinton](#) ↗ of Google, and [Yann LeCun](#) ↗ of New York University recipients of the 2018 ACM A.M. Turing Award for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing. Working independently and together, Hinton, LeCun and Bengio developed conceptual foundations for the field, identified surprising phenomena through experiments, and contributed engineering advances that demonstrated the practical advantages of deep neural networks.

Outline

- ❖ Convolutional Neural Network (CNN)
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
- ❖ Deep Learning Frameworks

Outline

- ❖ Convolutional Neural Network (CNN)
 - ◆ Convolution, Padding, Stride, Pooling
 - ◆ CNN Progress
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
- ❖ Deep Learning Frameworks

2-D Cross Correlation

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>	0	1	2	3
0	1	2													
3	4	5													
6	7	8													
0	1														
2	3														
	=	<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>	19	25	37	43									
19	25														
37	43														

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

$$1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25,$$

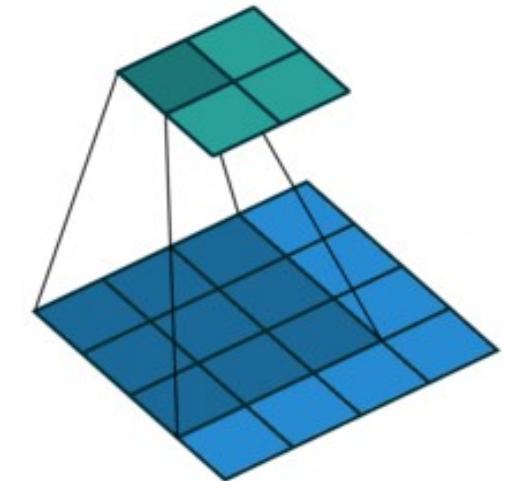
$$3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37,$$

$$4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$$

2-D Convolution Layer

$$\mathbf{Y} = \mathbf{X} \star \mathbf{W} + b$$

\mathbf{W} and b are learnable parameters, b : scalar bias



(vdumoulin@ Github)

$\mathbf{X}: n_h \times n_w$
input matrix

$\mathbf{W}: k_h \times k_w$
kernel matrix

$\mathbf{Y}: (n_h - k_h + 1) \times (n_w - k_w + 1)$
output matrix **feature map**

0	1	2
3	4	5
6	7	8

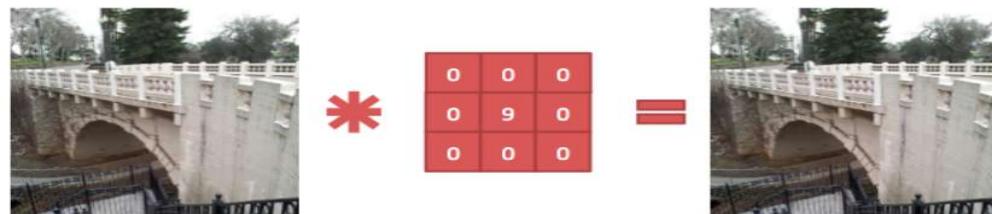
*

0	1
2	3

=

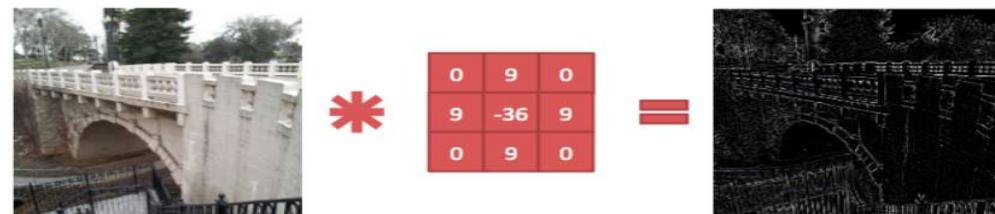
19	25
37	43

$$0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19, \\ 1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25, \\ 3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37, \\ 4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43.$$



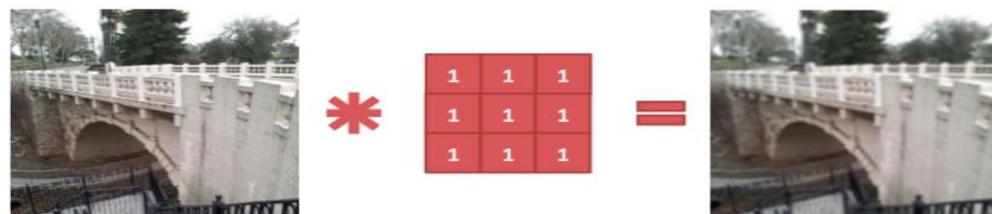
$$\begin{matrix} * & \begin{matrix} 0 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 0 \end{matrix} & = & \text{Input image} \end{matrix}$$

(a) Identity kernel



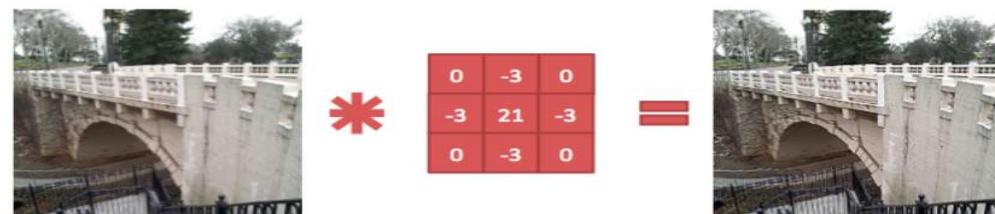
$$\begin{matrix} * & \begin{matrix} 0 & 9 & 0 \\ 9 & -36 & 9 \\ 0 & 9 & 0 \end{matrix} & = & \text{Output image} \end{matrix}$$

(b) Edge detection kernel



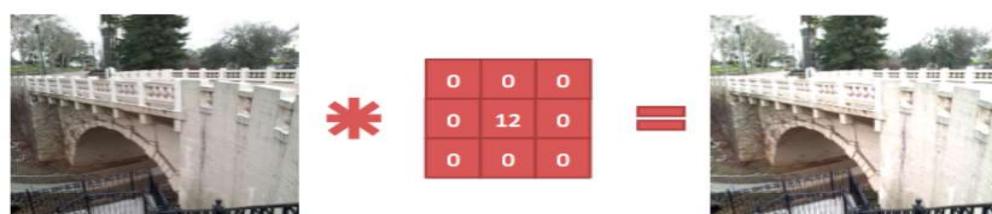
$$\begin{matrix} * & \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} & = & \text{Input image} \end{matrix}$$

(c) Blur kernel



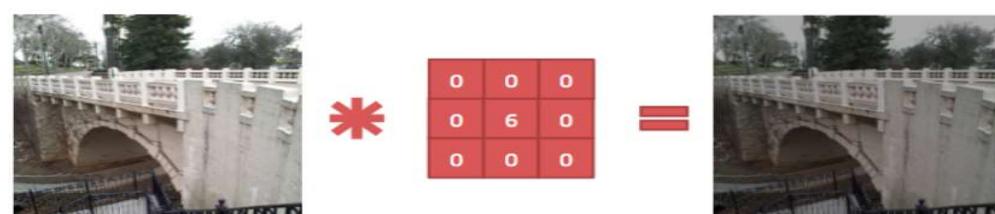
$$\begin{matrix} * & \begin{matrix} 0 & -3 & 0 \\ -3 & 21 & -3 \\ 0 & -3 & 0 \end{matrix} & = & \text{Output image} \end{matrix}$$

(d) Sharpen kernel



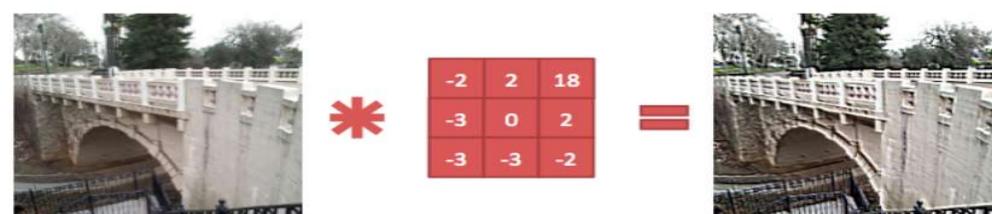
$$\begin{matrix} * & \begin{matrix} 0 & 0 & 0 \\ 0 & 12 & 0 \\ 0 & 0 & 0 \end{matrix} & = & \text{Input image} \end{matrix}$$

(e) Lighten kernel



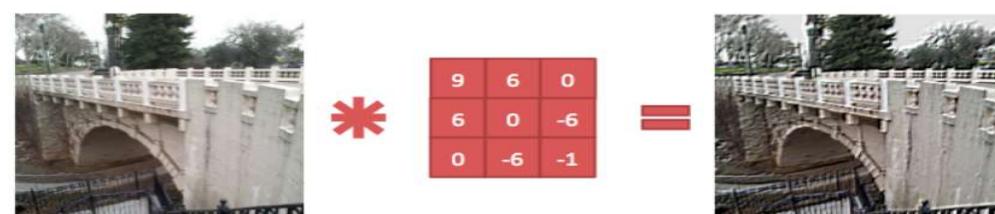
$$\begin{matrix} * & \begin{matrix} 0 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \end{matrix} & = & \text{Output image} \end{matrix}$$

(f) Darken kernel



$$\begin{matrix} * & \begin{matrix} -2 & 2 & 18 \\ -3 & 0 & 2 \\ -3 & -3 & -2 \end{matrix} & = & \text{Input image} \end{matrix}$$

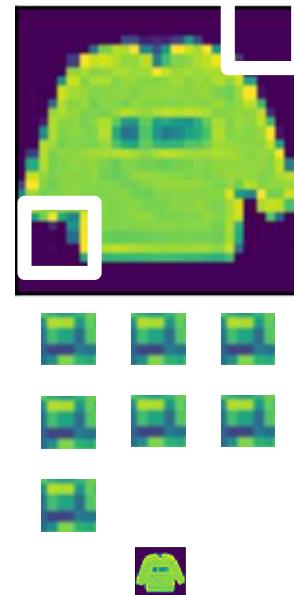
(g) Random kernel 1



$$\begin{matrix} * & \begin{matrix} 9 & 6 & 0 \\ 6 & 0 & -6 \\ 0 & -6 & -1 \end{matrix} & = & \text{Output image} \end{matrix}$$

(h) Random kernel 2

- ❖ Given a 32×32 input image
- ❖ Apply convolutional layer with 5×5 kernel
 - ◆ 28×28 output with 1 layer
 - ◆ 4×4 output with 7 layers
- ❖ **Shape decreases faster with larger kernels**
 - ◆ Shape reduces from $n_h \times n_w$ to
$$(n_h - k_h + 1) \times (n_w - k_w + 1)$$



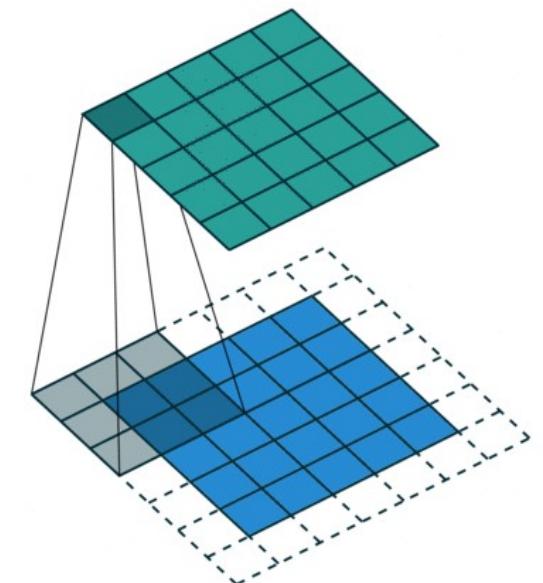
Padding

- Padding adds rows/columns around input
- Padding p_h rows and p_w columns, output shape will be

$$(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1)$$

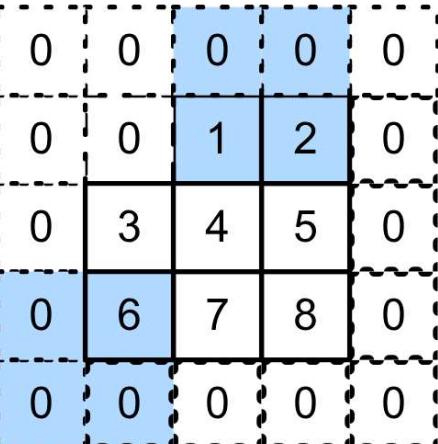
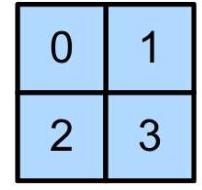
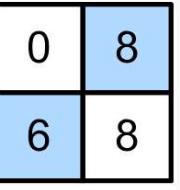
Input					Kernel		Output					
0	0	0	0	0	*	0	1	=	0	3	8	4
0	0	1	2	0		2	3		9	19	25	10
0	3	4	5	0					21	37	43	16
0	6	7	8	0					6	7	8	0
0	0	0	0	0								

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



Stride

- ❖ Stride is the #rows/#columns per slide.
- ❖ Given stride S_h for the height and stride S_w for the width, the output shape is $\lfloor(n_h - k_h + p_h + s_h)/s_h\rfloor \times \lfloor(n_w - k_w + p_w + s_w)/s_w\rfloor$

Input	Kernel	Output	
			<p>Strides of 3 and 2 for height and width</p> <p>$0 \times 0 + 0 \times 1 + 1 \times 2 + 2 \times 3 = 8$ $0 \times 0 + 6 \times 1 + 0 \times 2 + 0 \times 3 = 6$</p>
$*$		$=$	

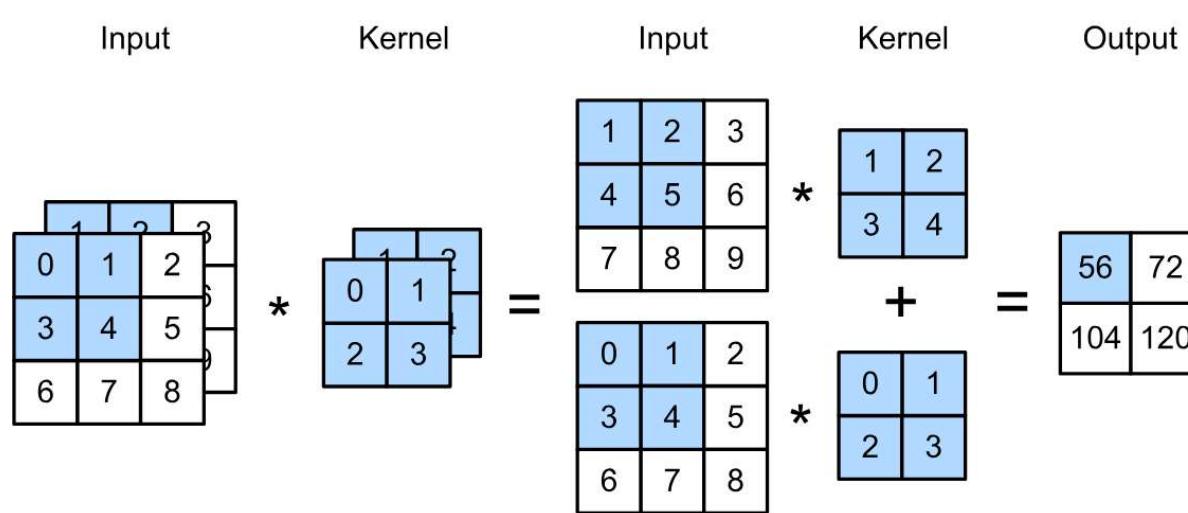
Multiple Input Channels

- ❖ Color image may have three RGB channels
- ❖ Converting to grayscale loses information



Multiple Input Channels

- ❖ Have a kernel for each channel, and then sum results over channels.



Input \mathbf{X} : $c_i \times n_h \times n_w$

Kernel \mathbf{W} : $c_i \times k_h \times k_w$

Output \mathbf{Y} : $m_h \times m_w$

$$\mathbf{Y} = \sum_{i=0}^{c_i} \mathbf{X}_{i,:,:} * \mathbf{W}_{i,:,:}$$

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

* No matter how many inputs channels, we always get single output channel.

Multiple Output Channels

- ❖ We can have multiple 3-D kernels, each one generates a output channel.

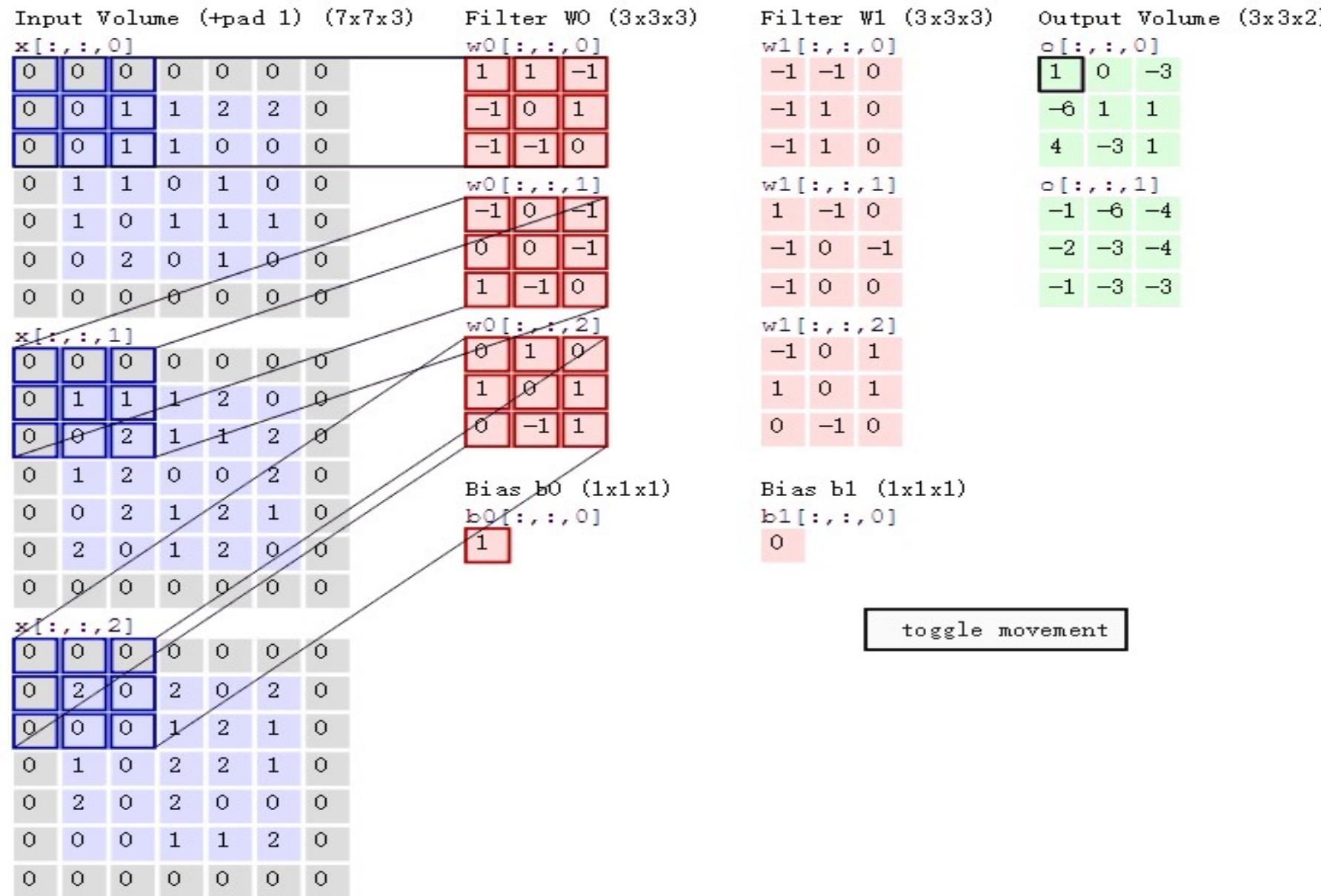
- ◆ Input \mathbf{X} : $c_i \times n_h \times n_w$

$$\mathbf{Y}_{i,:,:} = \mathbf{X} \star \mathbf{W}_{i,:,:}$$

- ◆ Kernel \mathbf{W} : $c_o \times c_i \times k_h \times k_w$

$$for i = 1, \dots, c_o$$

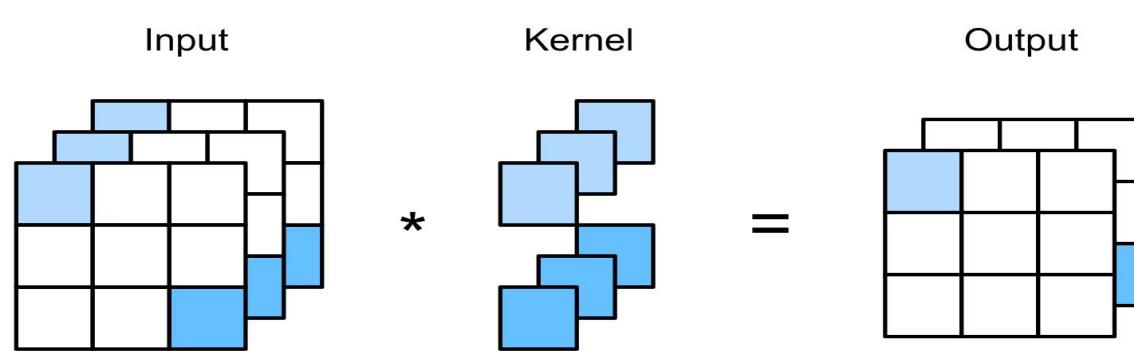
- ◆ Output \mathbf{Y} : $c_o \times m_h \times m_w$



1 x 1 Convolutional Layer

$$k_h = k_w = 1$$

It doesn't recognize spatial patterns, but fuse channels.



Equal to a dense layer with $n_h n_w \times c_i$ input and $c_o \times c_i$ weight.

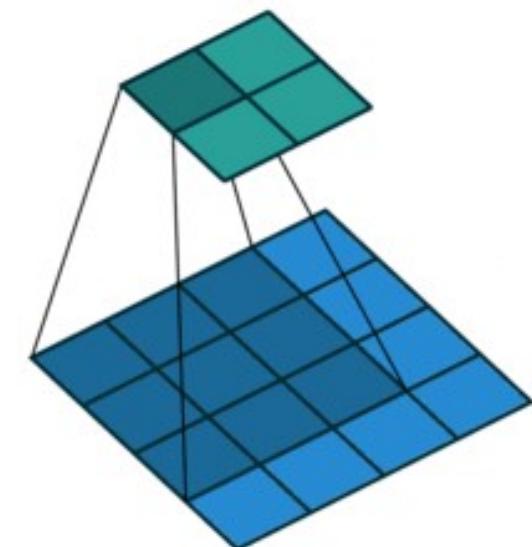
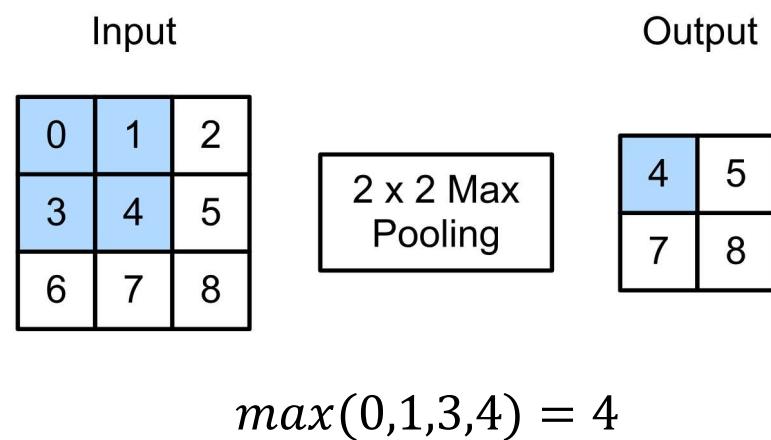
Pooling

- ❖ Why pooling:
 - ◆ reduce the dimensionality when processing data
 - ◆ Alleviate the excessive sensitivity of the convolutional layer to location
 - Lighting, object positions, scales, appearance vary among images
 - Convolution is sensitive to position
 - e.g. Detect vertical edges

```
X [[1. 1. 0. 0. 0.      Y [[ 0.  1.  0.  0.  
     [1. 1. 0. 0. 0.      [ 0.  1.  0.  0.  
     [1. 1. 0. 0. 0.      [ 0.  1.  0.  0.  
     [1. 1. 0. 0. 0.]      [ 0.  1.  0.  0.]
```

2-D Max Pooling

- ❖ Returns the maximal value in the sliding window



Average Pooling

- ❖ Max pooling: the strongest pattern signal in a window
- ❖ Average pooling: replace max with mean in max pooling
 - ◆ The average signal strength in a window



Max pooling

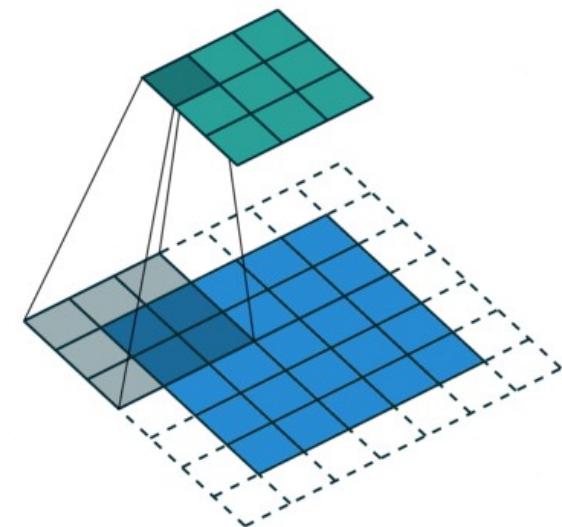


Average pooling

Padding, Stride, and Multiple Channels

- ❖ Pooling layers have similar padding and stride as convolutional layers
- ❖ No learnable parameters
- ❖ Apply pooling for each input channel to obtain the corresponding output channel

#output channels = #input channels



Outline

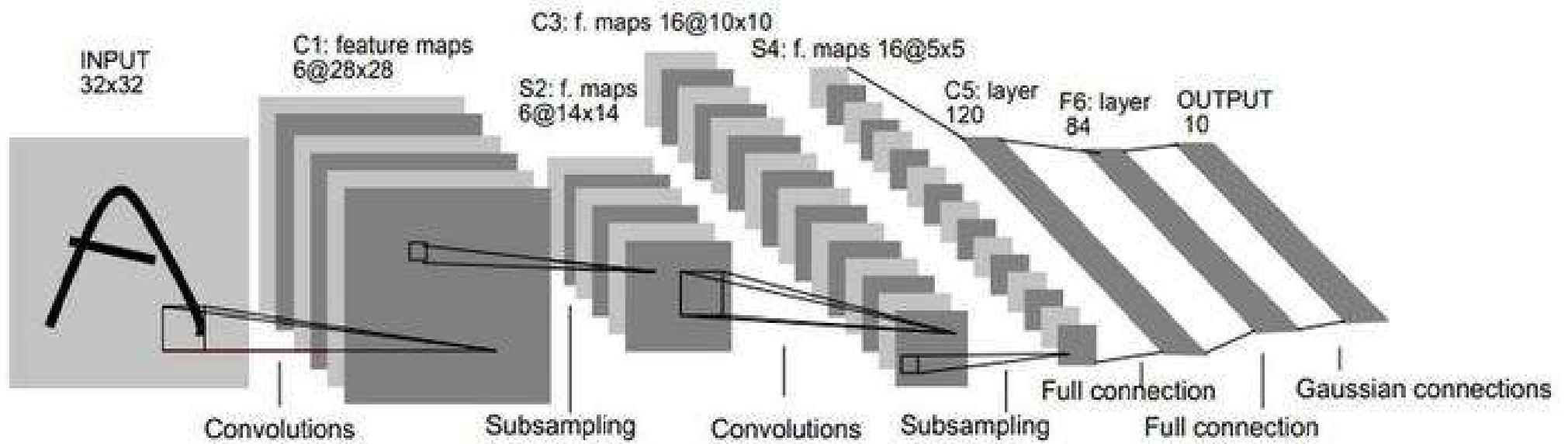
- ❖ Convolutional Neural Network (CNN)
 - ◆ Convolution, Padding, Stride, Pooling
 - ◆ CNN Progress
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
- ❖ Deep Learning Frameworks

CNN Progress

- ❖ LeNet
 - ◆ 2 convolution + pooling layers ; 2 hidden dense hidden layers
- ❖ AlexNet(2012)
 - ◆ Bigger and deeper LeNet; ReLu, Dropout, preprocessing
- ❖ NiN (2013)
 - ◆ 1x1 convolutions + global pooling instead of dense
- ❖ VGG (2014)
- ❖ GoogLeNet(2014)
 - ◆ Inception
- ❖ ResNet (2015)
 - ◆ Residual blocks

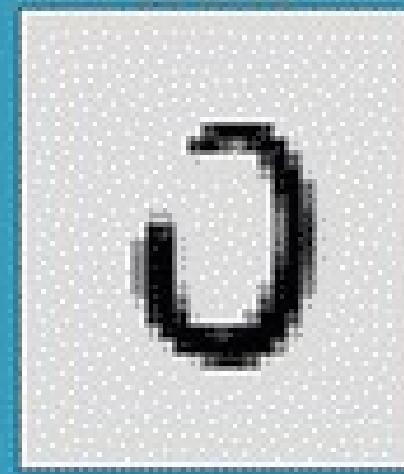
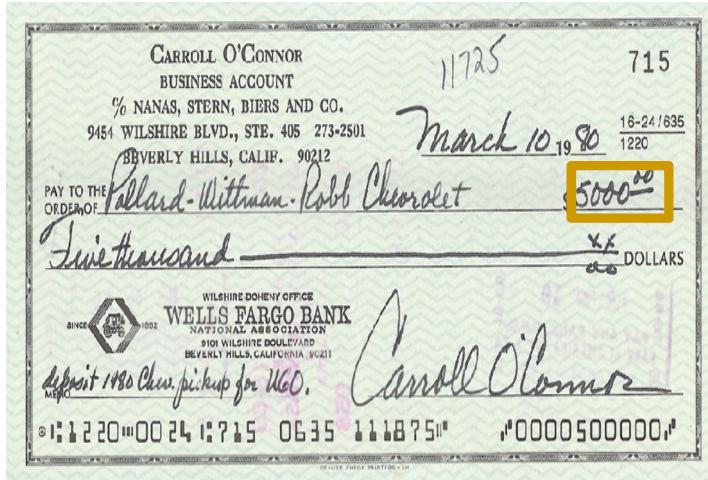
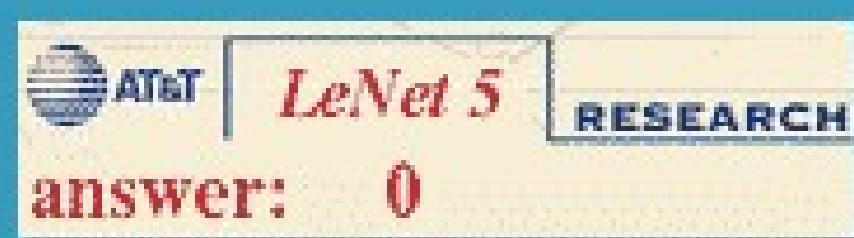
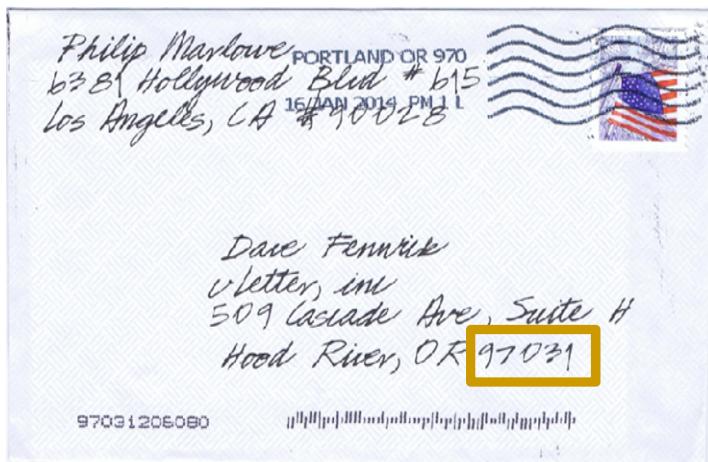
◦ ◦ ◦ ◦ ◦ ◦

LeNet

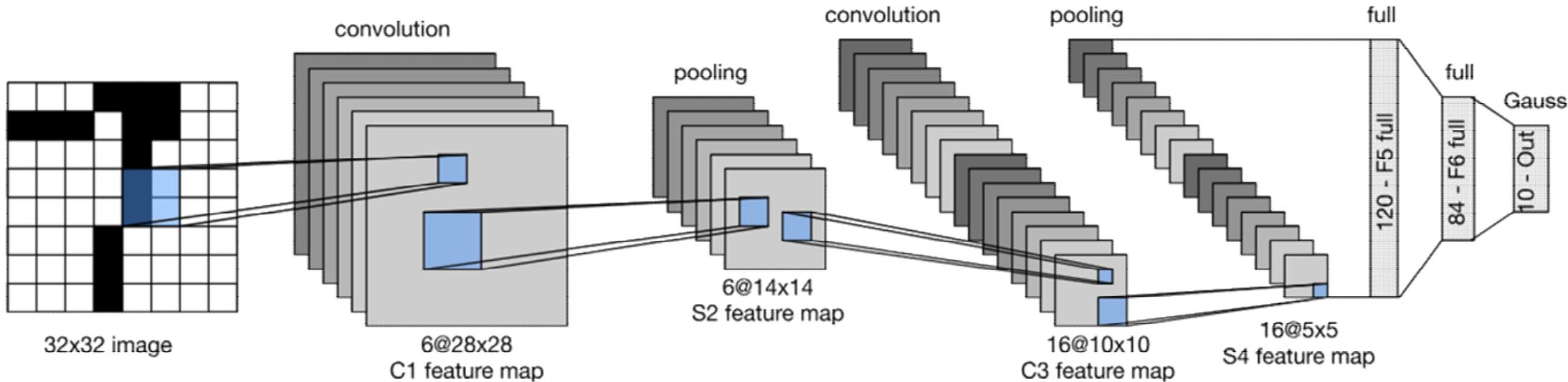


LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.

Handwritten Digit Recognition



LeNet Architecture



- ❖ convolutional layers labeled C_x,
- ❖ subsampling layers labeled S_x
- ❖ fully connected layers labeled F_x

LeNet

❖ convolutional block

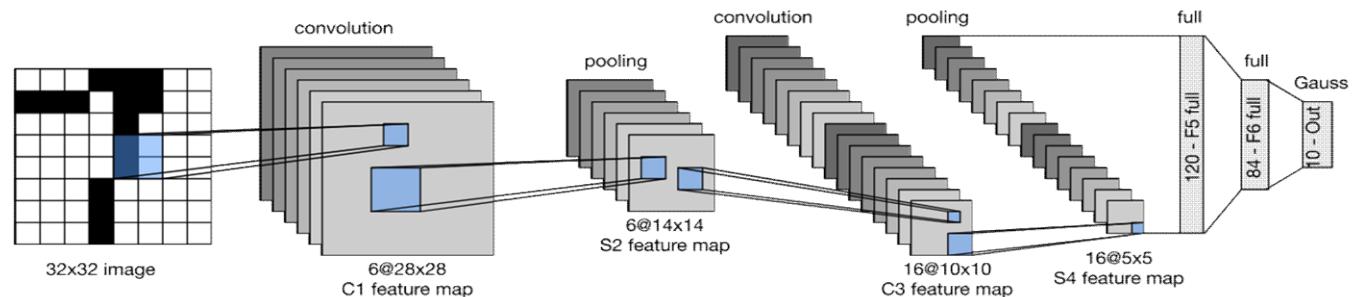
◆ convolutional layer

- ❑ used to recognize the spatial patterns in the image, such as lines and the parts of objects

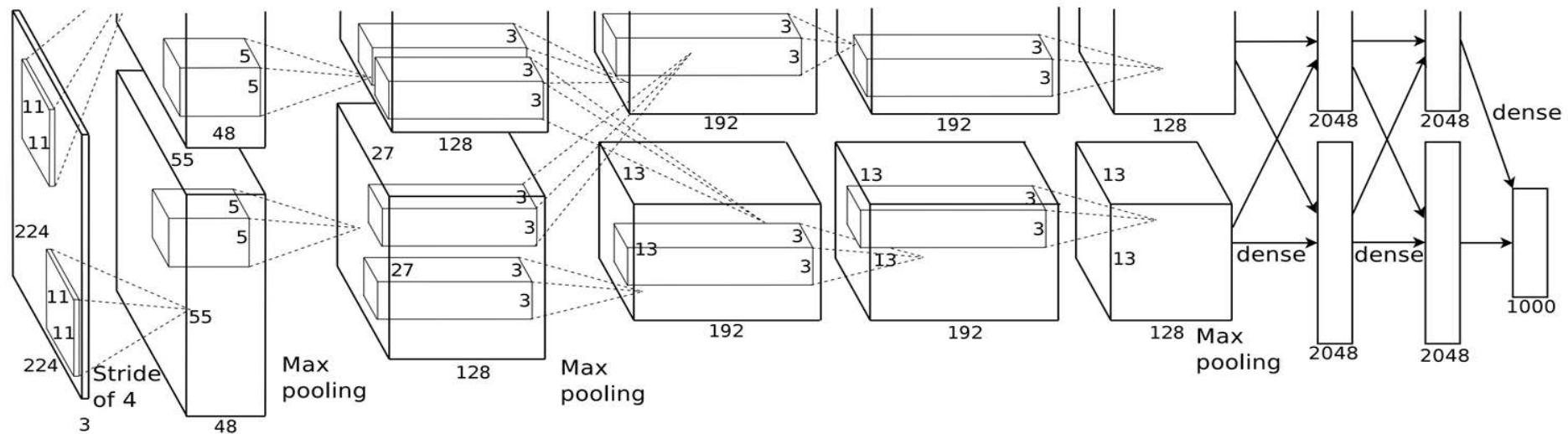
- ❑ each convolutional layer uses a **5*5** window

◆ pooling layer

- ❑ **average pooling** layer is used to reduce the dimensionality
- ❑ the window shape is 2×2 and the stride is 2

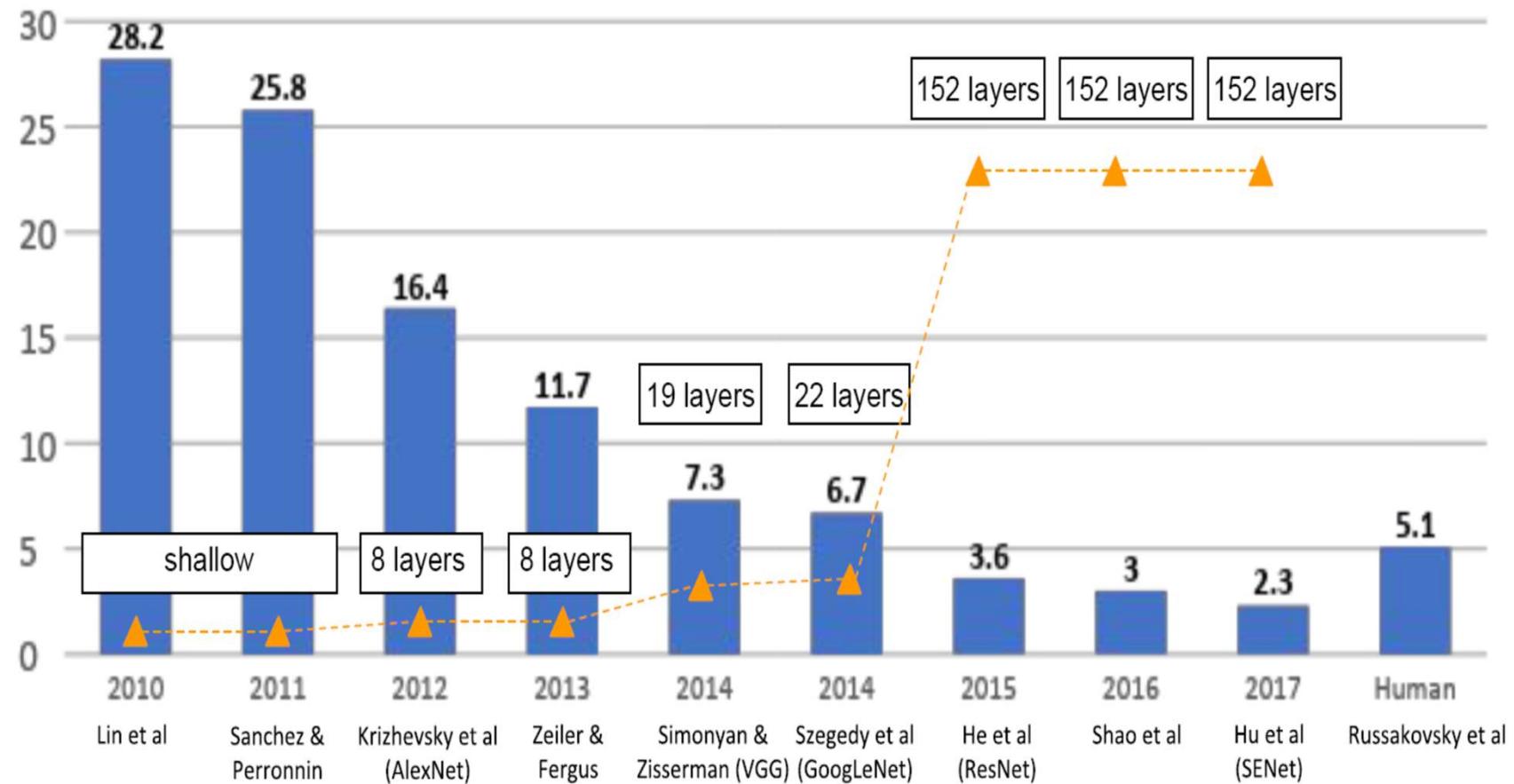


AlexNet



Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems (pp. 1097-1105).

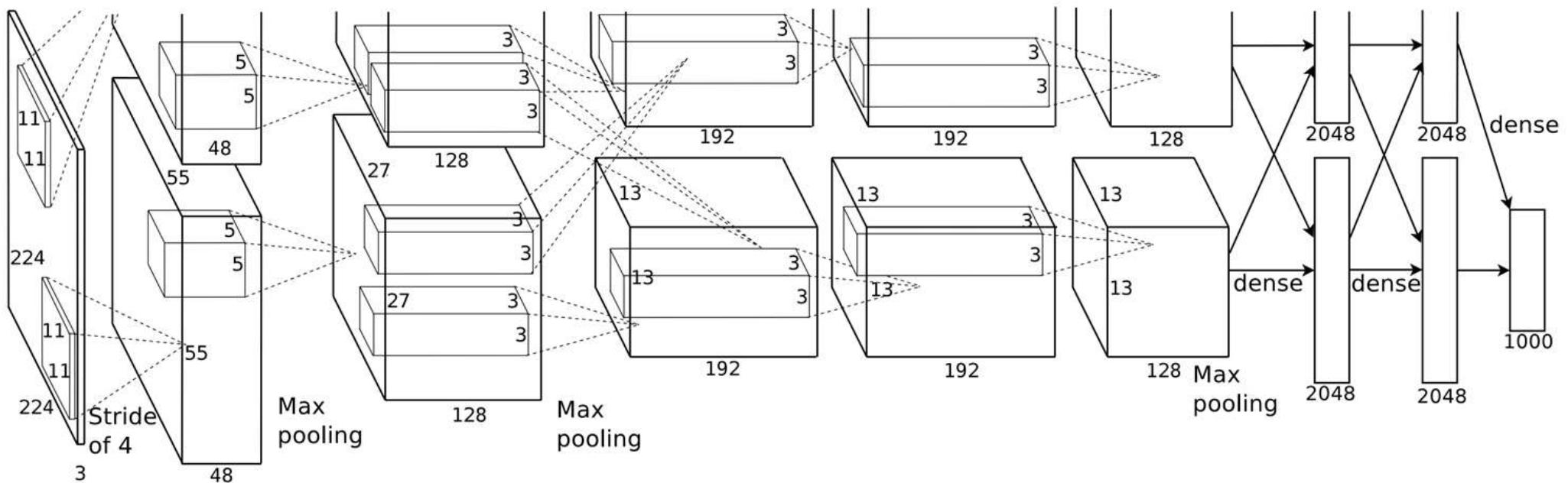
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



AlexNet

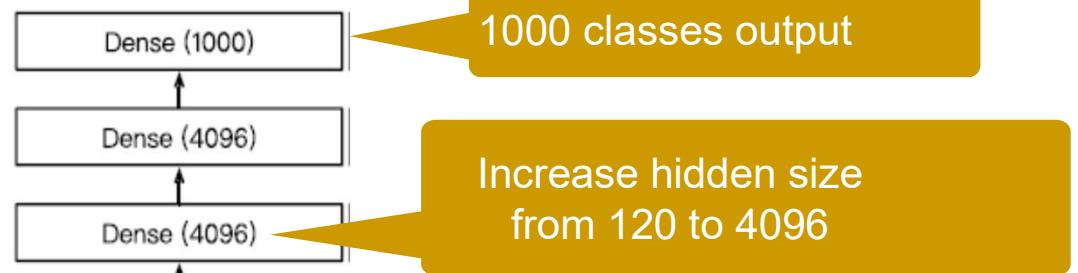
- ❖ A key step from shallow to deep networks !
- ❖ AlexNet won ImageNet competition in 2012
- ❖ Deeper and bigger LeNet
- ❖ Key modifications
 - ◆ Dropout (regularization)
 - ◆ ReLu (training)
 - ◆ Data augmentation、Max Pooling

AlexNet Architecture

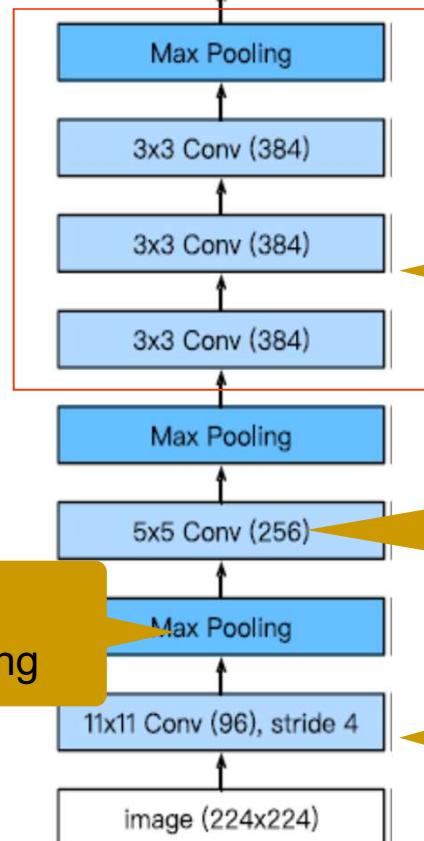
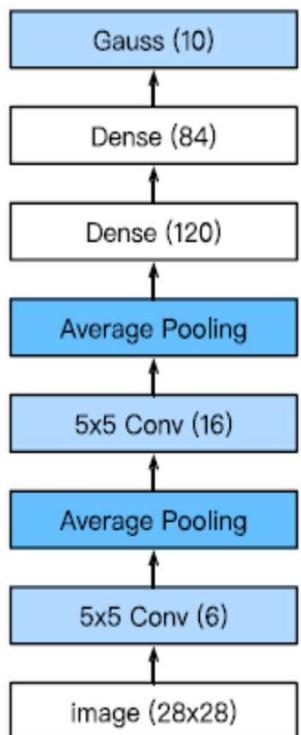


AlexNet Architecture

AlexNet

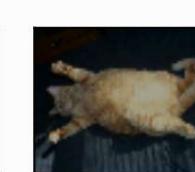


LeNet



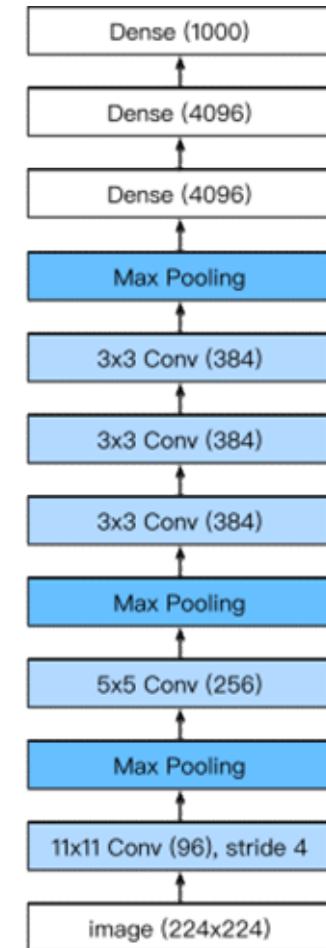
More Tricks

- ❖ Change activation function from sigmoid to **ReLU** (no more vanishing gradient)
- ❖ Add a **dropout** layer after two hidden dense layers (better robustness / regularization)
- ❖ **Data augmentation**

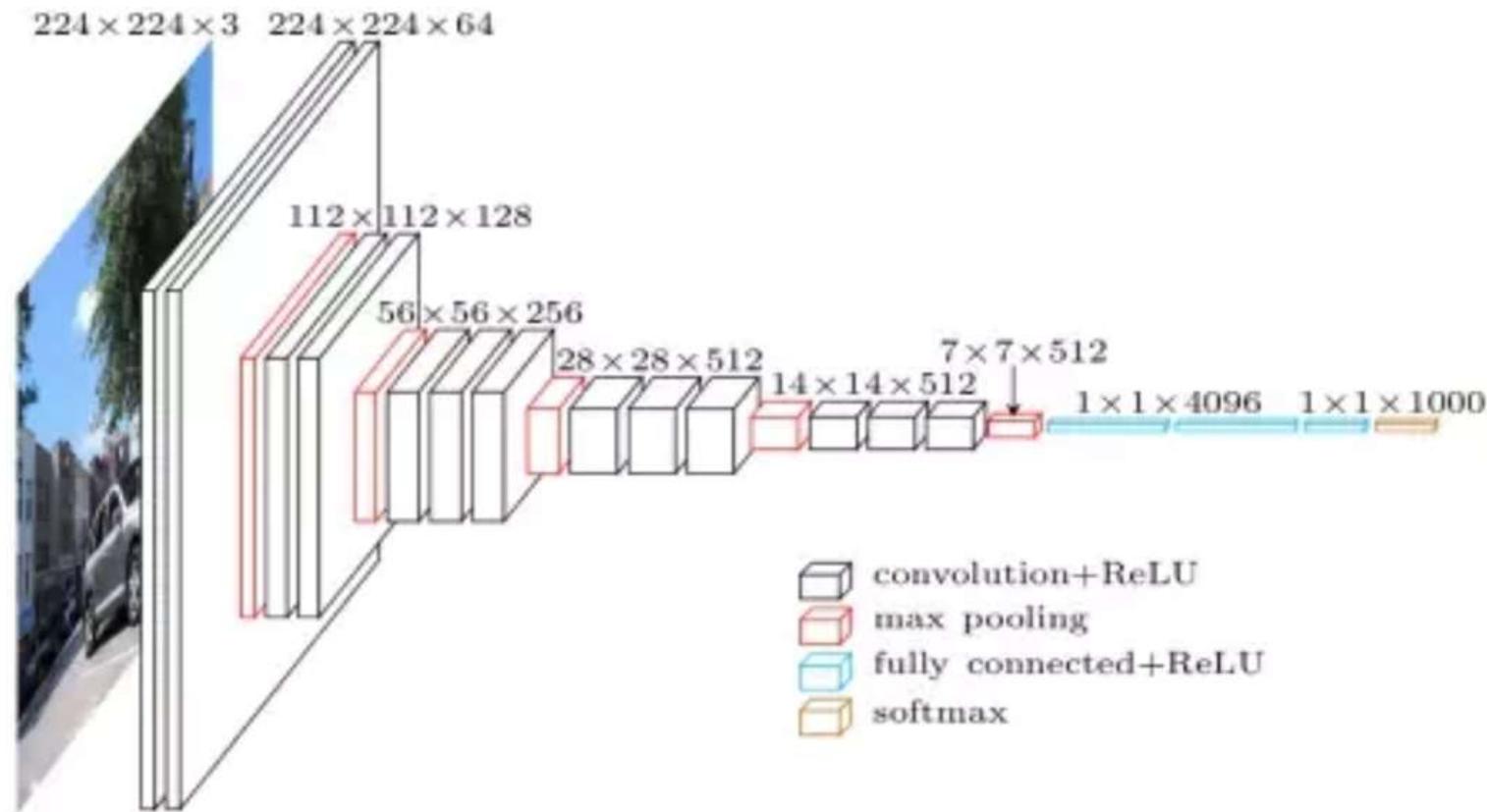


Complexity

	#parameters		FLOP	
	AlexNet	LeNet	AlexNet	LeNet
Conv1	35K	150	101M	1.2M
Conv2	614K	2.4K	415M	2.4M
Conv3-5	3M		445M	
Dense1	26M	0.48M	26M	0.48M
Dense2	16M	0.1M	16M	0.1M
Total	46M	0.6M	1G	4M
Increase	11x	1x	250x	1x

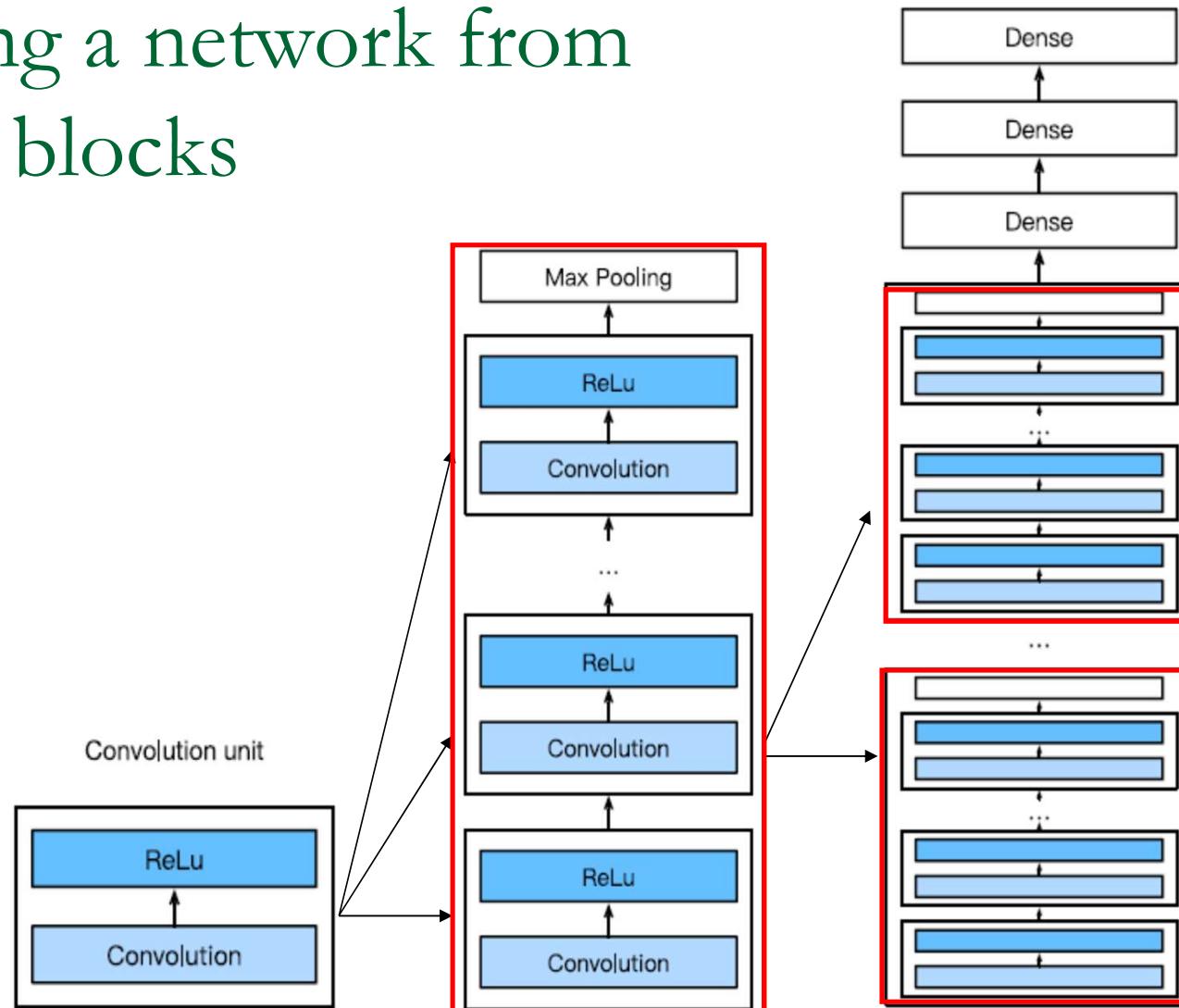


VGG



Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556.

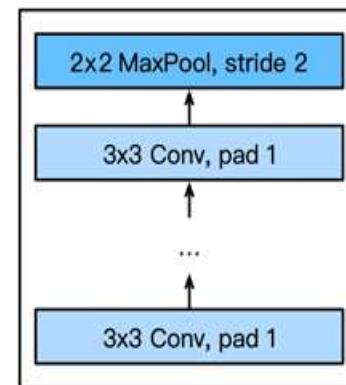
Designing a network from building blocks



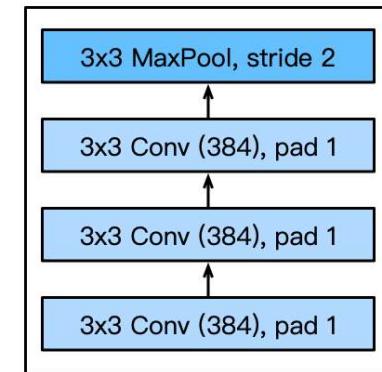
VGG Blocks

- ❖ Deeper vs. wider?
 - ◆ 5x5 convolutions
 - ◆ 3x3 convolutions (more)
 - ◆ **deep & narrow better**
- ❖ VGG block
 - ◆ 3x3 convolutions (pad 1)
(*n* layers, *m* channels)
 - ◆ 2x2 max-pooling
(stride 2)

VGG block

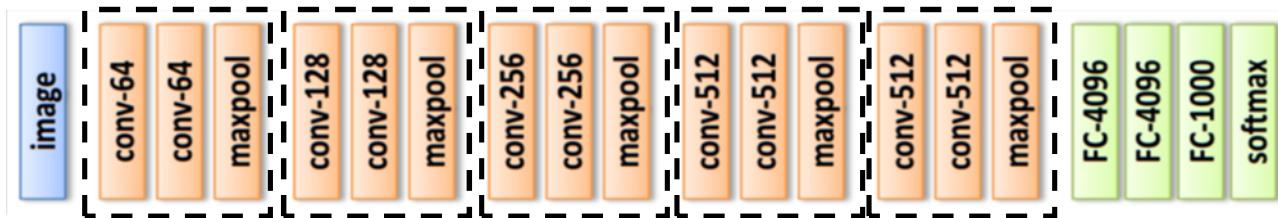


Part of AlexNet



VGG Architecture

- ❖ Multiple VGG blocks followed by dense layers



- ❖ Vary the repeating number to get different architectures, such as VGG-16, VGG-19, ...

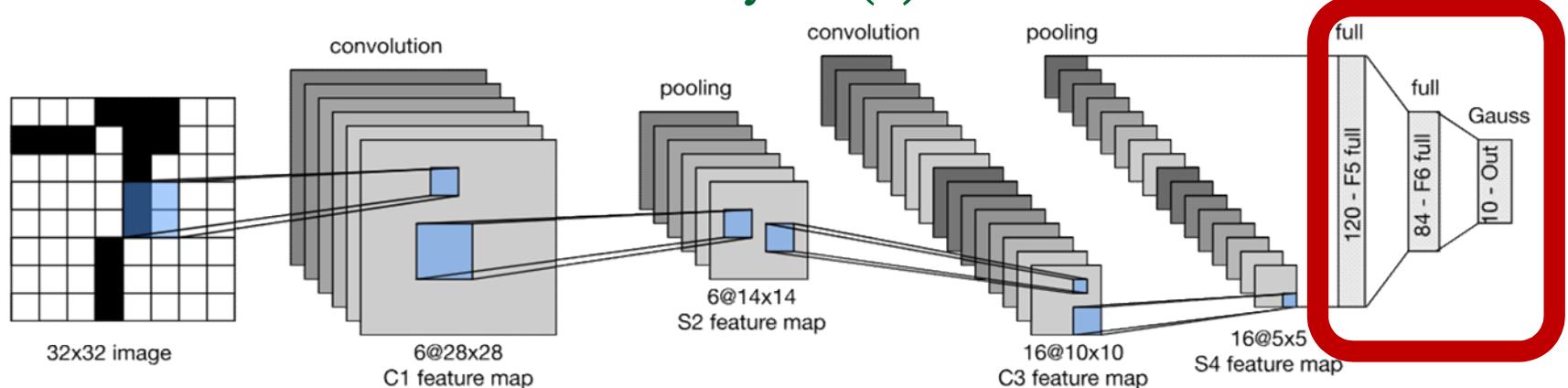
VGG

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Design pattern

- ❖ The design pattern of LeNet, AlexNet, and VGG:
 - ◆ extract the spatial **features** through a sequence of **convolutions** and **pooling layers**
 - ◆ post-process the representations via **fully connected layers**

The Curse of the Last Layer(s)



- ❖ Convolution layers need relatively few parameters

$$c_i \times c_o \times k^2$$

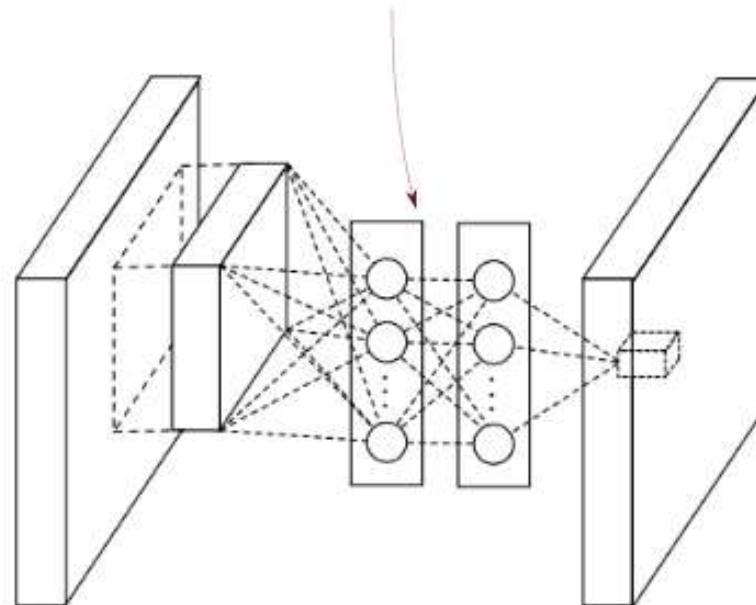
- ❖ Last layer needs many parameters for n classes

$$c \times m_w \times m_h \times n$$

- ❖ LeNet: $16 \times 5 \times 5 \times 120 = 48k$
- ❖ VGG: $512 \times 7 \times 7 \times 4096 = 102M$

Network in Network

Non linear mapping introduced by mlpconv layer consisting of multiple fully connected layers with non linear activation function.



Breaking the Curse of the Last Layer

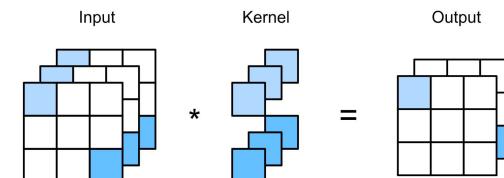
❖ Key Idea

- ❖ **Get rid of the fully connected layer(s)**
- ❖ Convolutions and pooling reduce resolution
(e.g. stride of 2 reduces resolution 4x)

❖ Implementation details

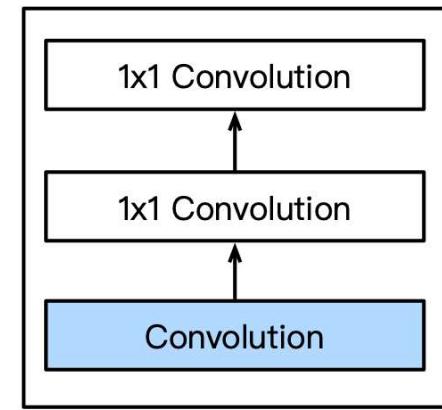
- ❖ Reduce resolution progressively
- ❖ Increase number of channels
- ❖ Use **1x1 convolutions** (they only act per pixel)

❖ Global average pooling in the end



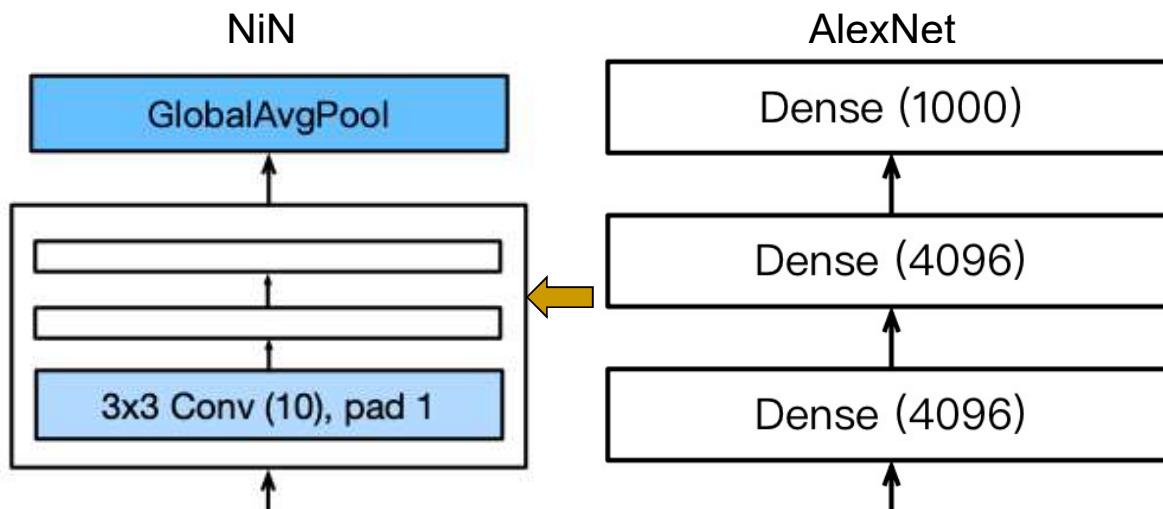
NiN Block

- ❖ A convolutional layer
 - ◆ kernel size, stride, and padding are hyper-parameters
- ❖ Followed by two 1x1 convolutions
 - ◆ 1 stride and no padding, share the same output channels as the first layer
 - ◆ Act as dense layers



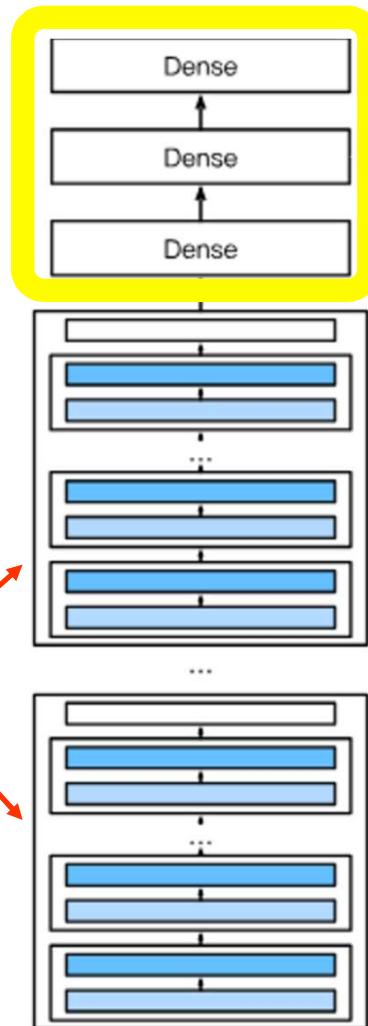
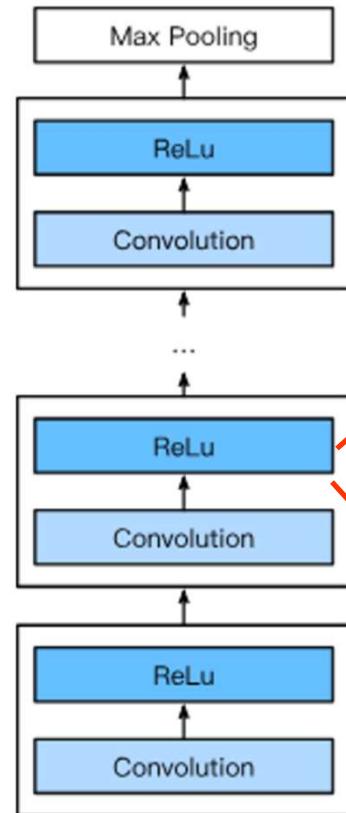
NiN Last Layers

- ❖ Replace AlexNet's dense layers with a NiN block
- ❖ Global average pooling layer to combine outputs



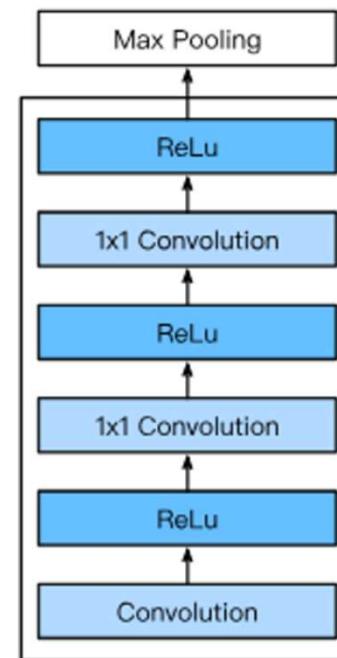
NiN Networks

Convolution block

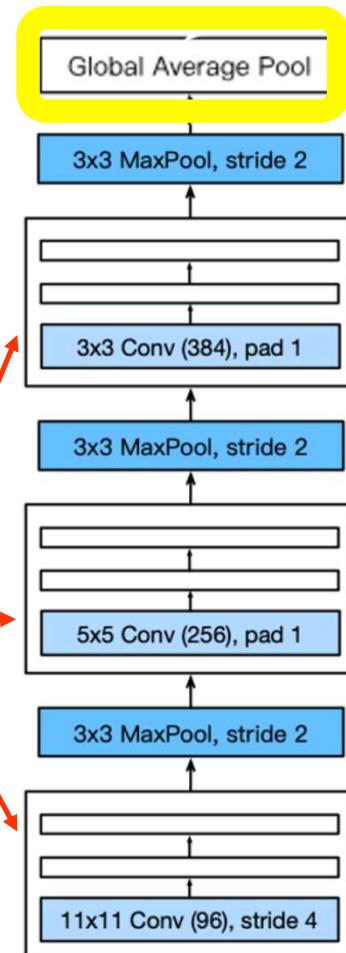


VGG Net

NiN block

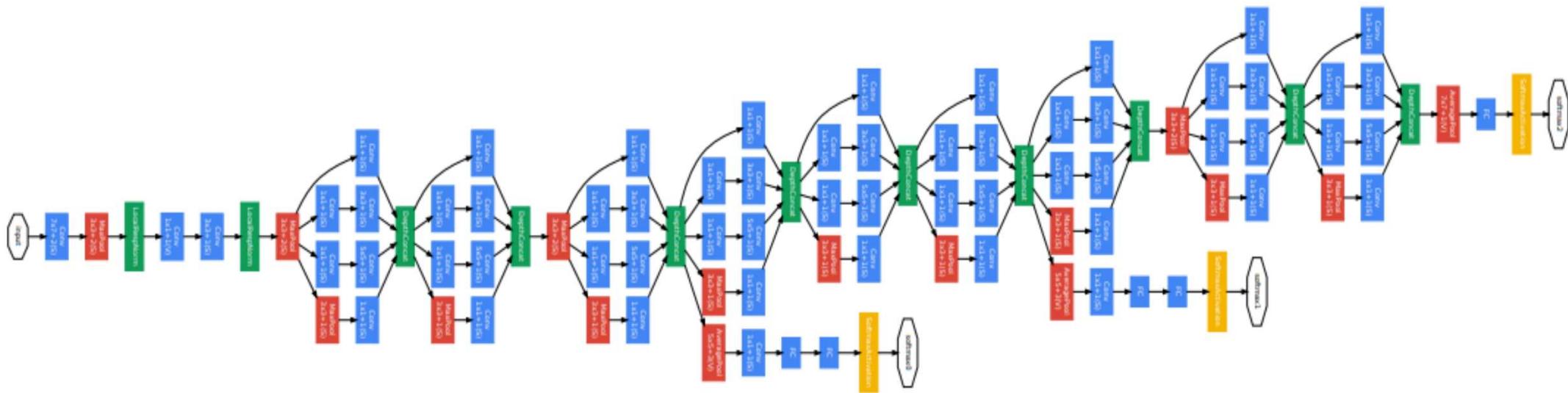


NiN Net



GoogLeNet(2014)

❖ Networks with Parallel Concatenations



Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., & Anguelov, D. & Rabinovich, A. (2015). Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1-9).

Picking the best convolution ...

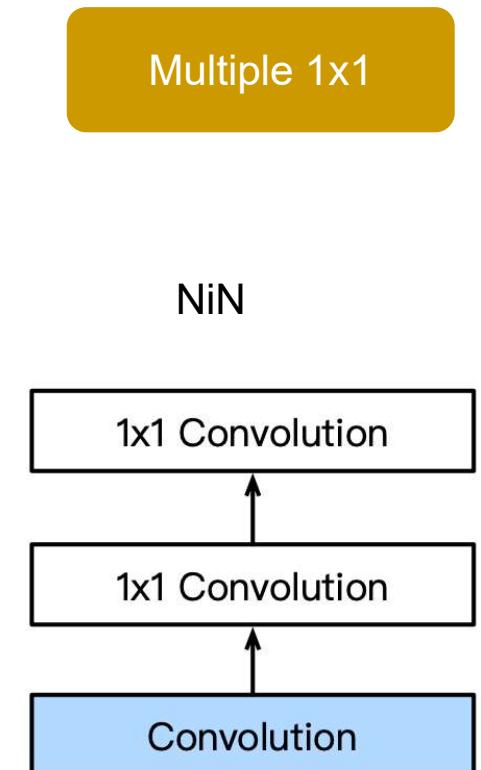
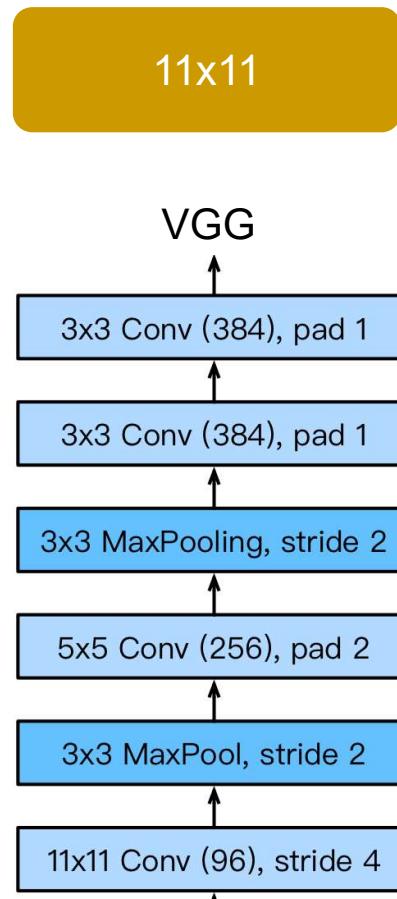
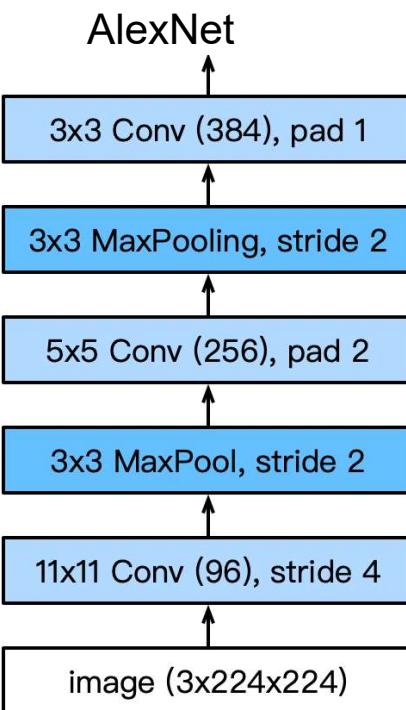
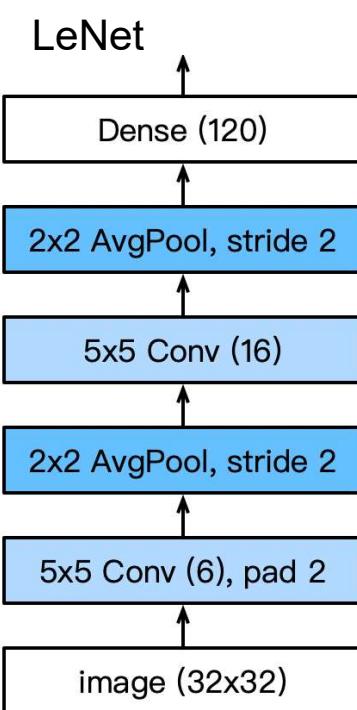
1x1

3x3

5x5

11x11

Multiple 1x1

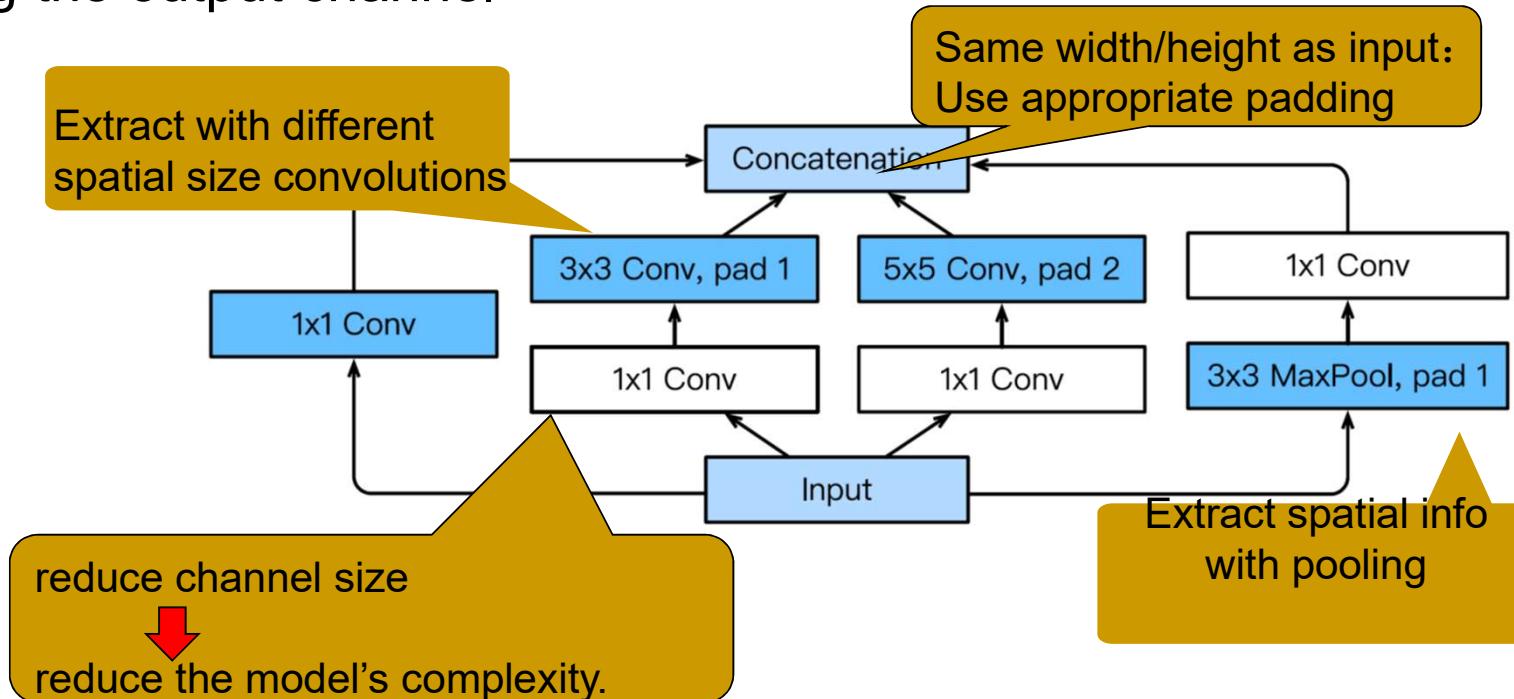


Why choose? Just pick them all !

Inception Blocks

- ✓ The basic convolutional block in GoogLeNet.

4 paths extract information from different aspects, then concatenate along the output channel

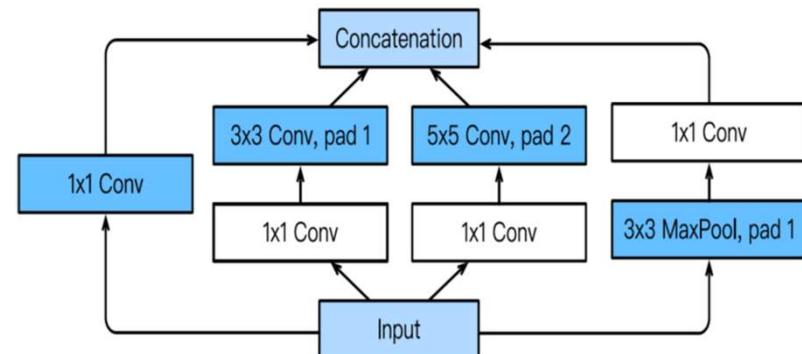


Inception Blocks

Inception blocks have fewer parameters and less computation complexity than a single 3x3 or 5x5 convolutional layer

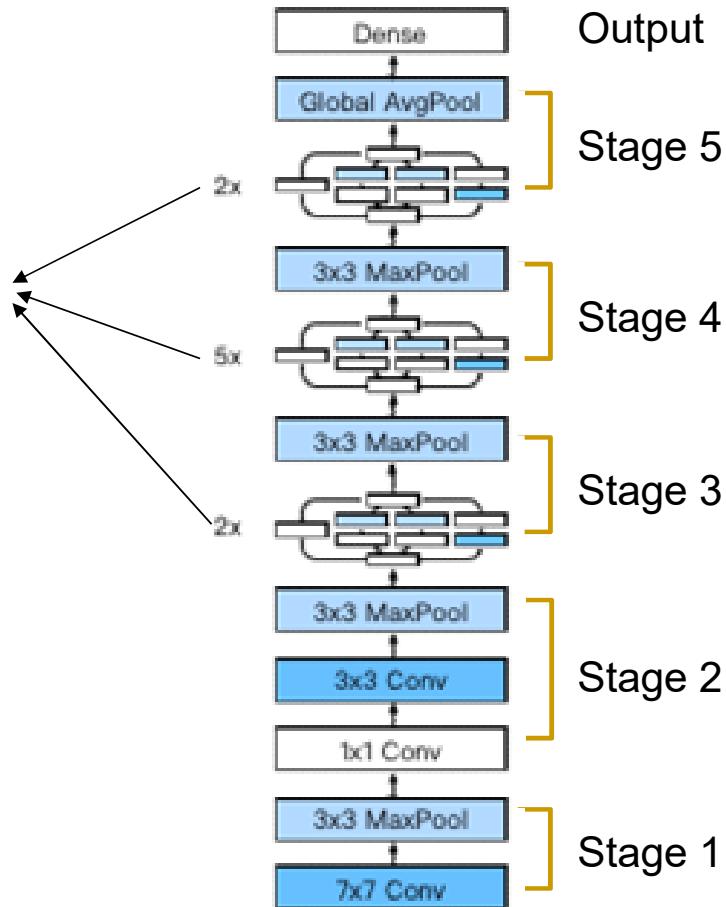
- ❖ Mix of different functions (powerful function class)
- ❖ Memory and compute efficiency (good generalization)

	#parameters	FLOPS
Inception	0.16 M	128 M
3x3 Conv	0.44 M	346 M
5x5 Conv	1.22 M	963 M



GoogLeNet

5 stages with
9 inception blocks

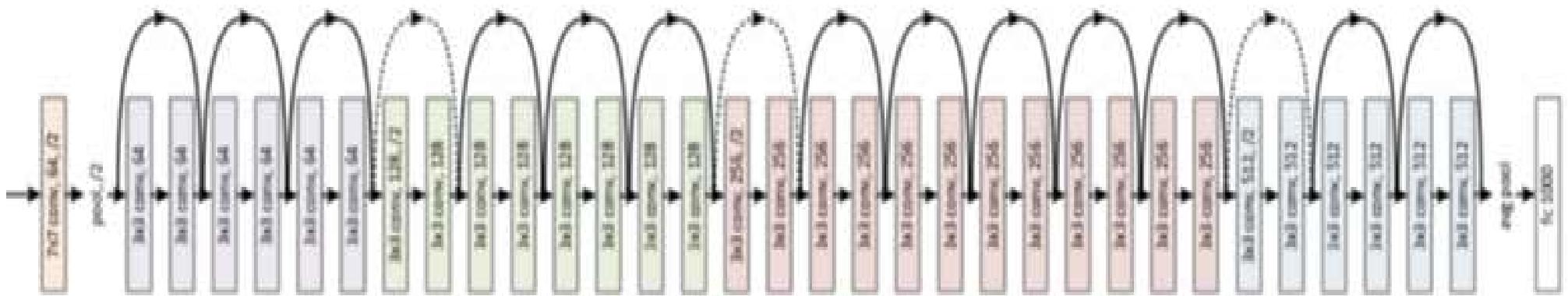


GoogLeNet Incarnation of the Inception Architecture

The many flavors of Inception Networks

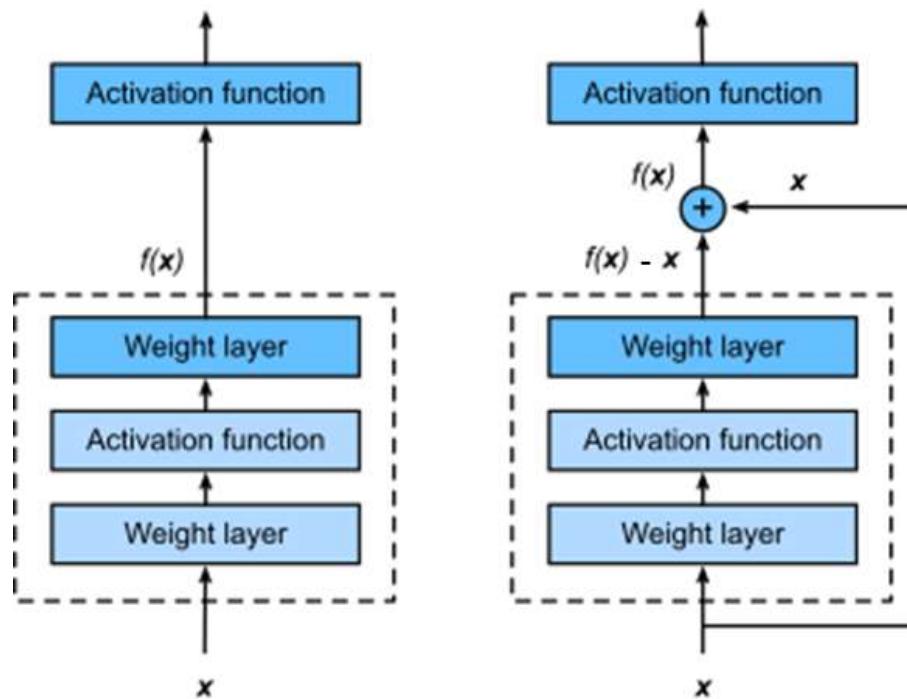
- ❖ Inception-BN (v2) - Add batch normalization
- ❖ Inception-V3 - Modified the inception block
 - ◆ Replace 5x5 by multiple 3x3 convolutions
 - ◆ Replace 5x5 by 1x7 and 7x1 convolutions
 - ◆ Replace 3x3 by 1x3 and 3x1 convolutions
 - ◆ Generally deeper stack
- ❖ Inception-V4 - Add residual connections (more later)

Residual Networks (ResNet)



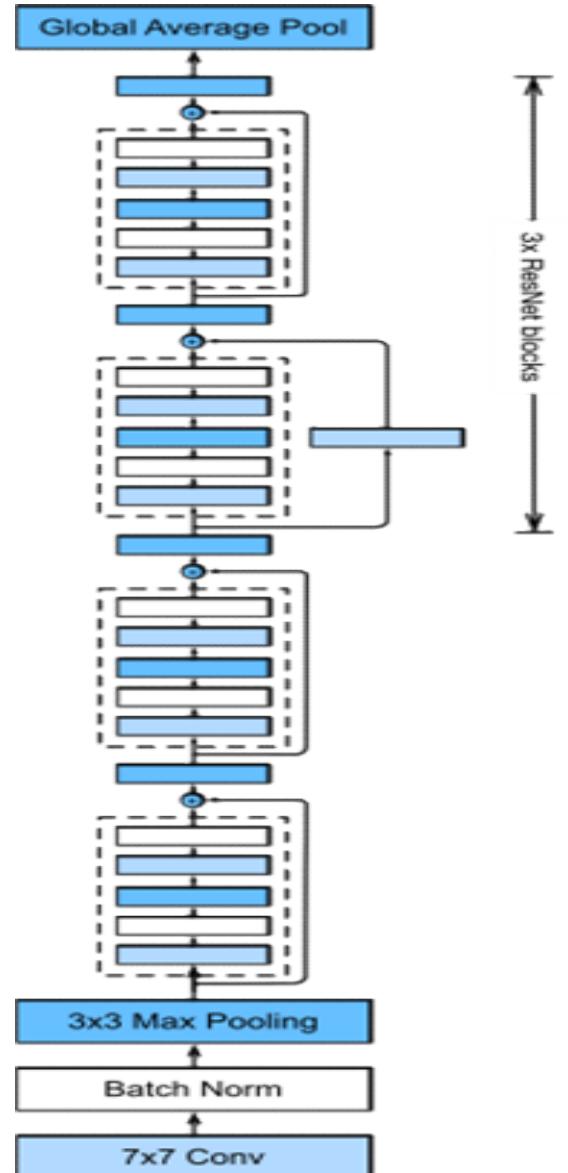
He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. CVPR (pp. 770-778).

Residual Networks (ResNet)

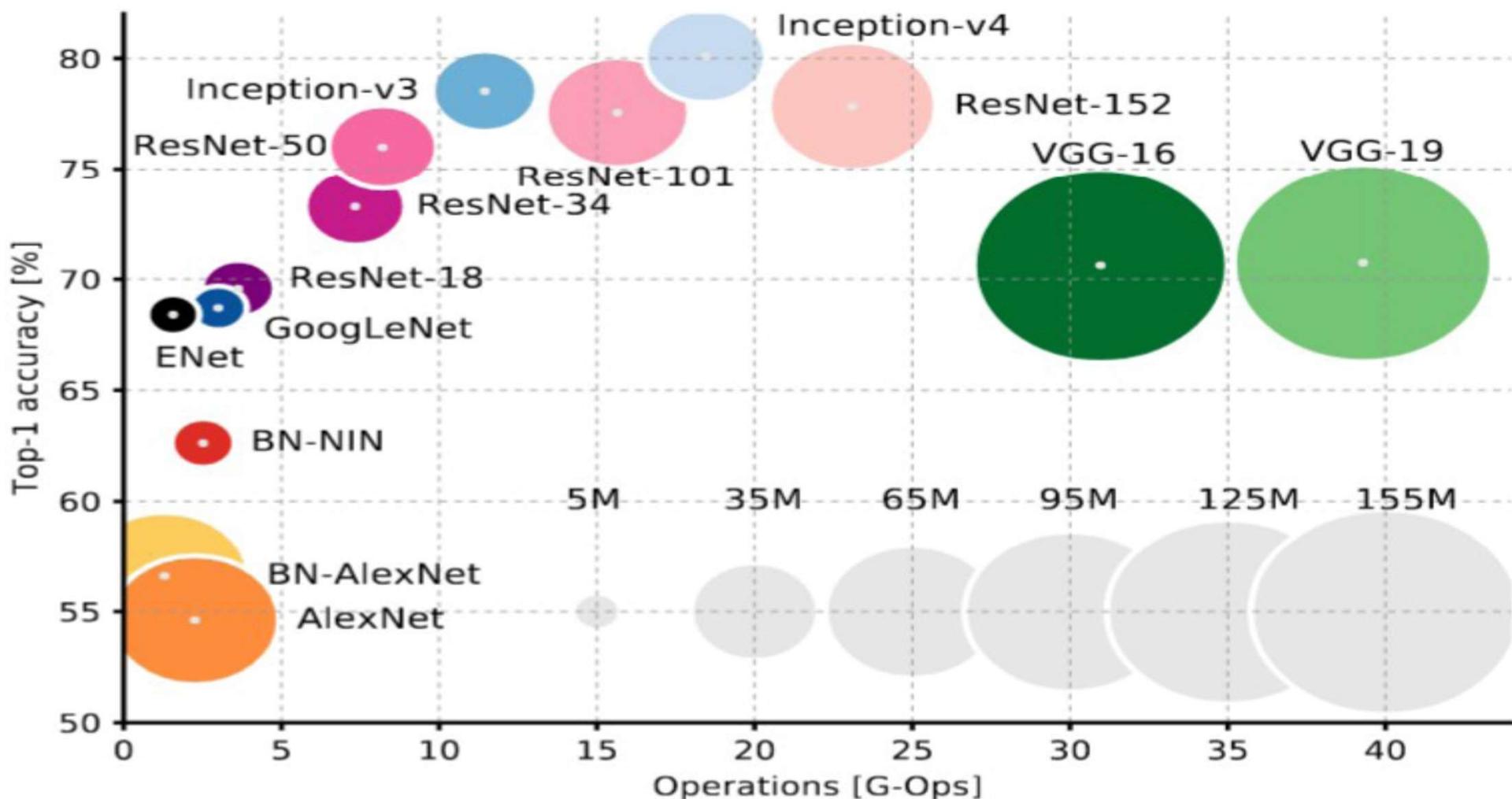


ResNet

- ❖ Same block structure as e.g. VGG or GoogleNet
- ❖ Residual connection to add to the expressiveness
- ❖ Pooling/stride for dimensionality reduction
 - ◆ Down sample per module (stride=2)
- ❖ Batch Normalization for capacity control



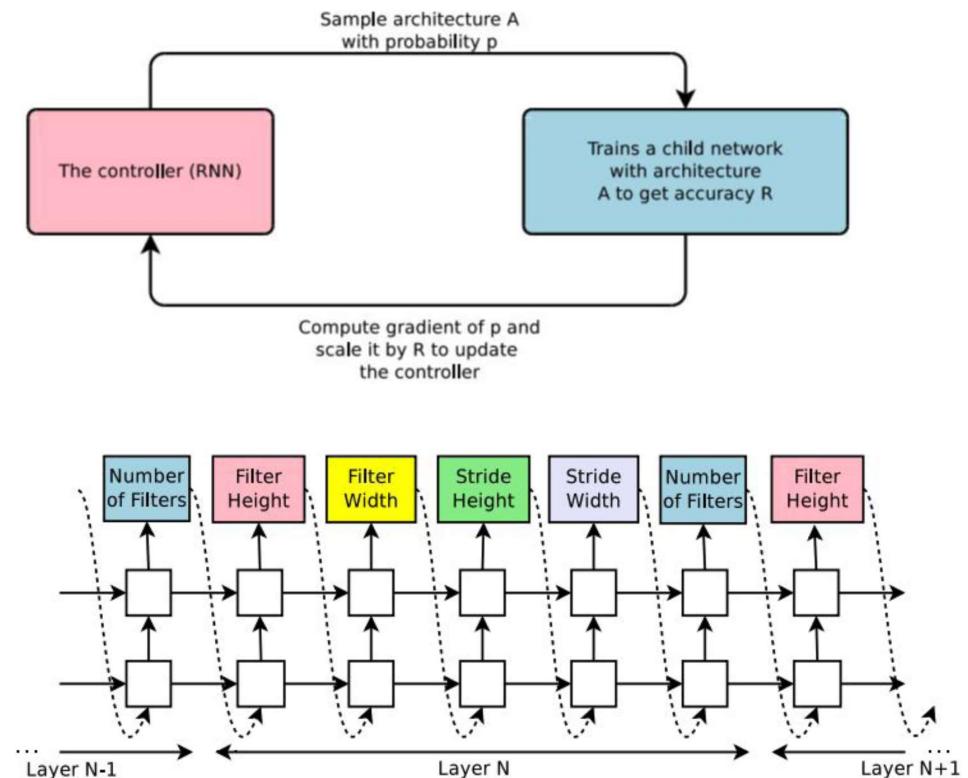
GOPS vs. Accuracy on ImageNet vs. #Parameters



Neural Architecture Search with Reinforcement Learning (NAS)

[Zoph et al. 2016]

- ❖ “Controller” network that learns to design a good network architecture (output a string corresponding to network design)
- ❖ Iterate:
 - ◆ 1) Sample an architecture from search space
 - ◆ 2) Train the architecture to get a “reward” R corresponding to accuracy
 - ◆ 3) Compute gradient of sample probability, and scale by R to perform controller parameter update (i.e. increase likelihood of good architecture being sampled, decrease likelihood of bad architecture)



Outline

- ❖ Convolutional Neural Network (CNN)
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
- ❖ Deep Learning Frameworks

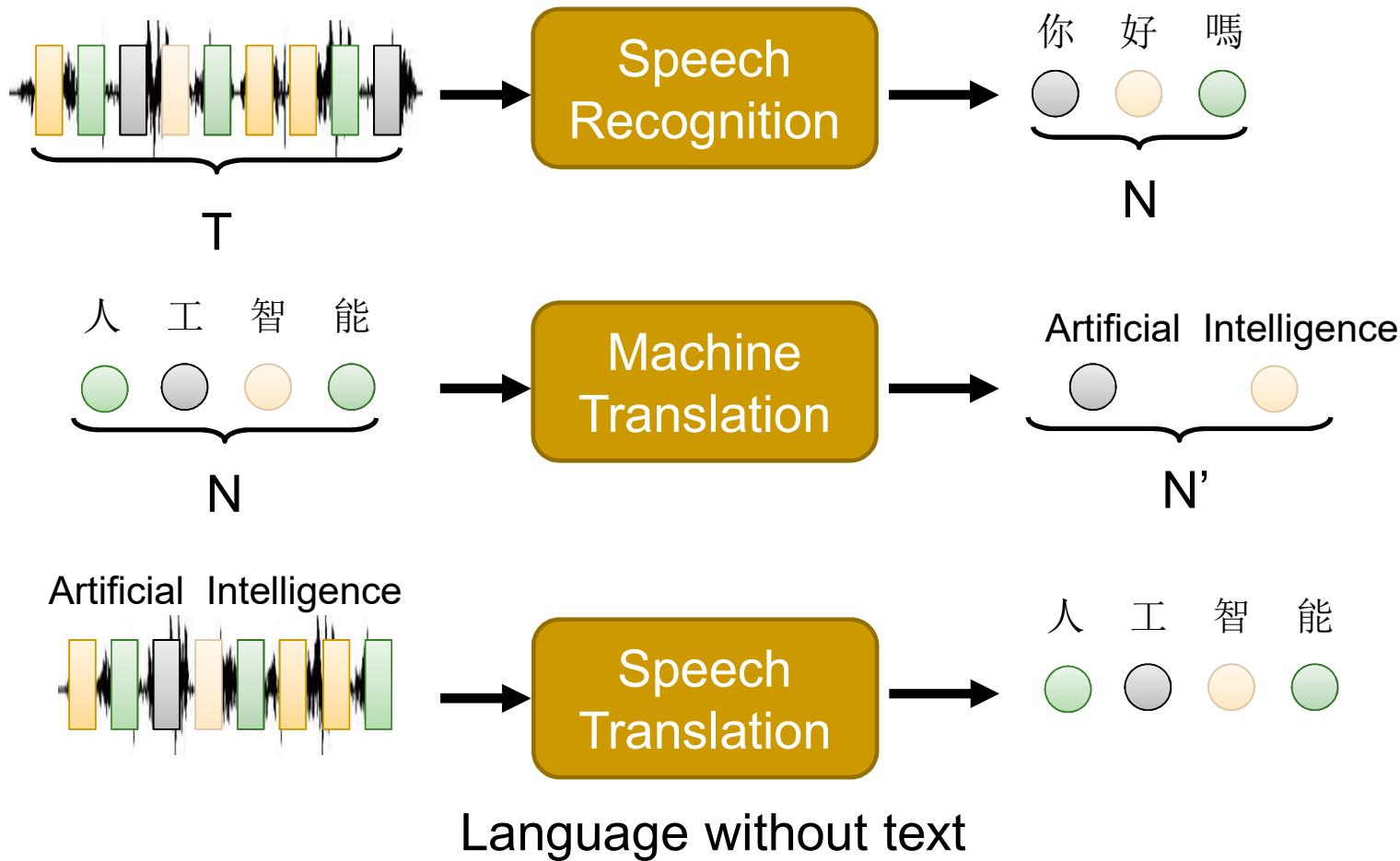
Another Scenario

- ❖ Data is not always generated i.i.d., all drawn from some distribution, but follows sequential order.
 - e.g.
 - ◆ the words in a paragraph are written in sequence
 - ◆ image frames in a video
 - ◆ the audio signal in a conversation
 - ◆ the browsing behavior on a website
- ❖ Not only receive a sequence as an input, but rather might be expected to continue the sequence.

Sequence-to-sequence (Seq2seq)

Input a sequence, output a sequence

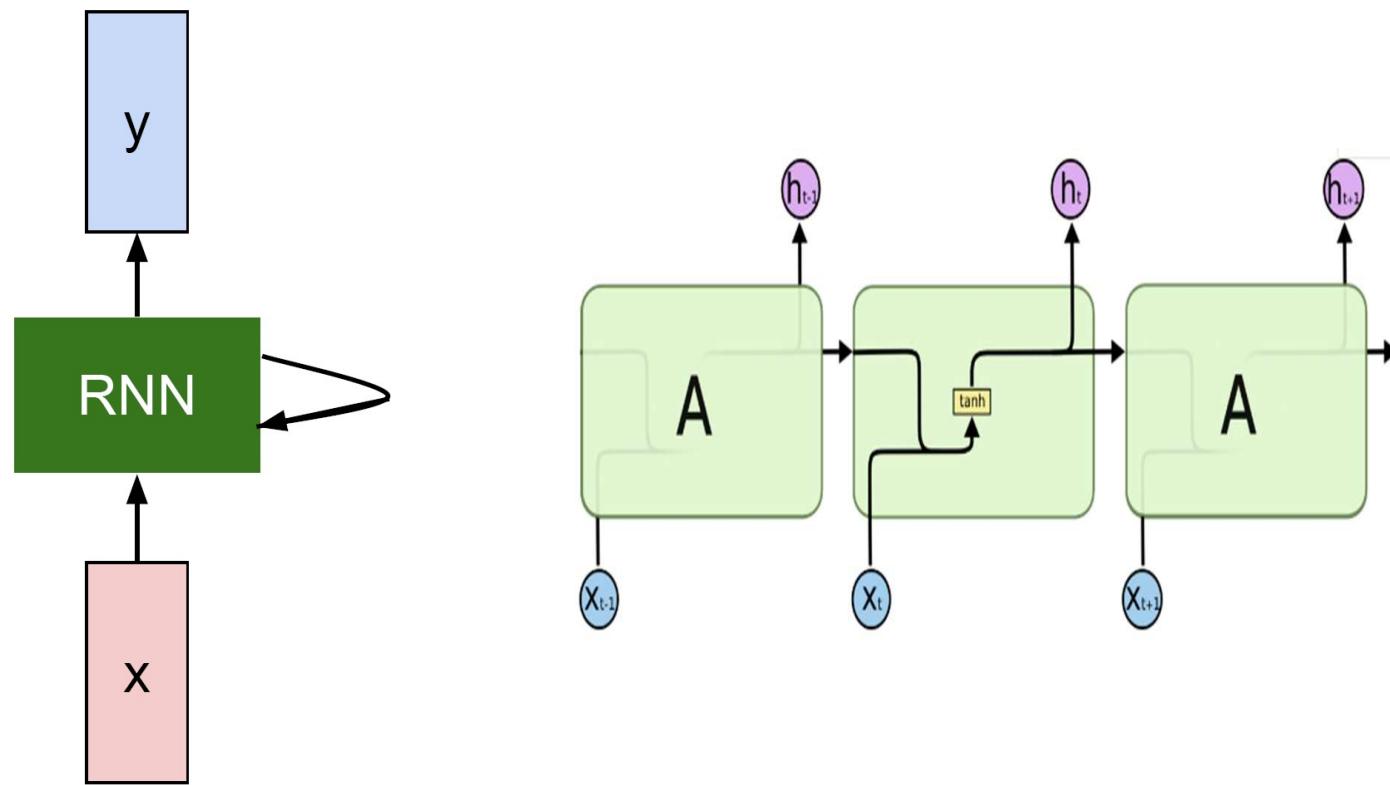
The output length is determined by model.



When the order of data matters.....

Sequence Models !

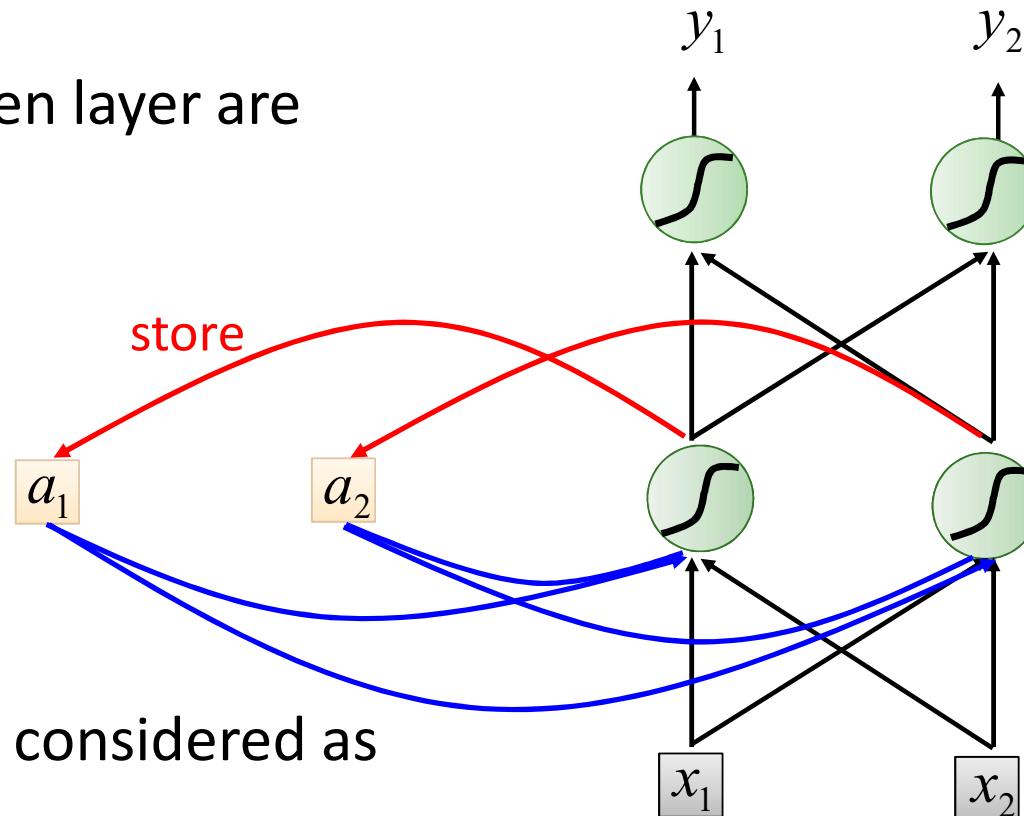
Recurrent Neural Networks



Recurrent Neural Network (RNN)

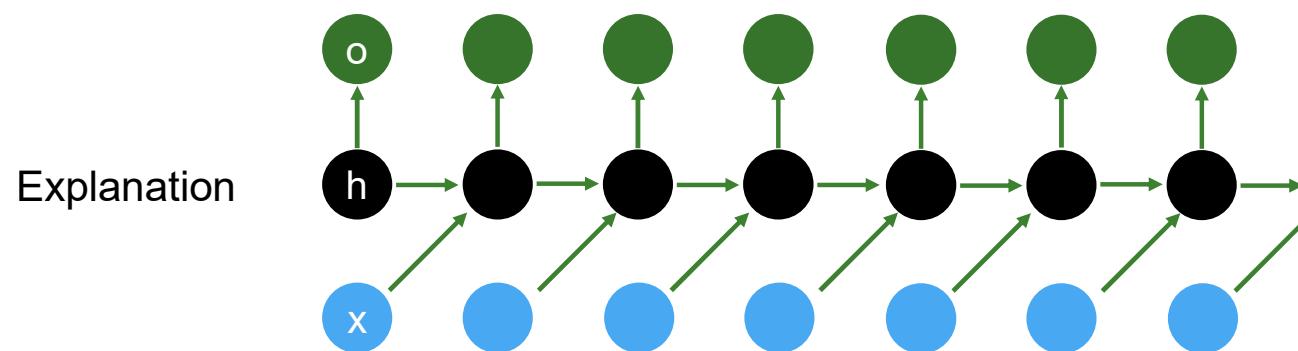
The outputs of the hidden layer are stored in the memory.

Memory can be considered as another input.



Recurrent Neural Networks

(with hidden state)



⑩ Hidden State update

$$\mathbf{h}_t = \phi(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

⑩ Observation update

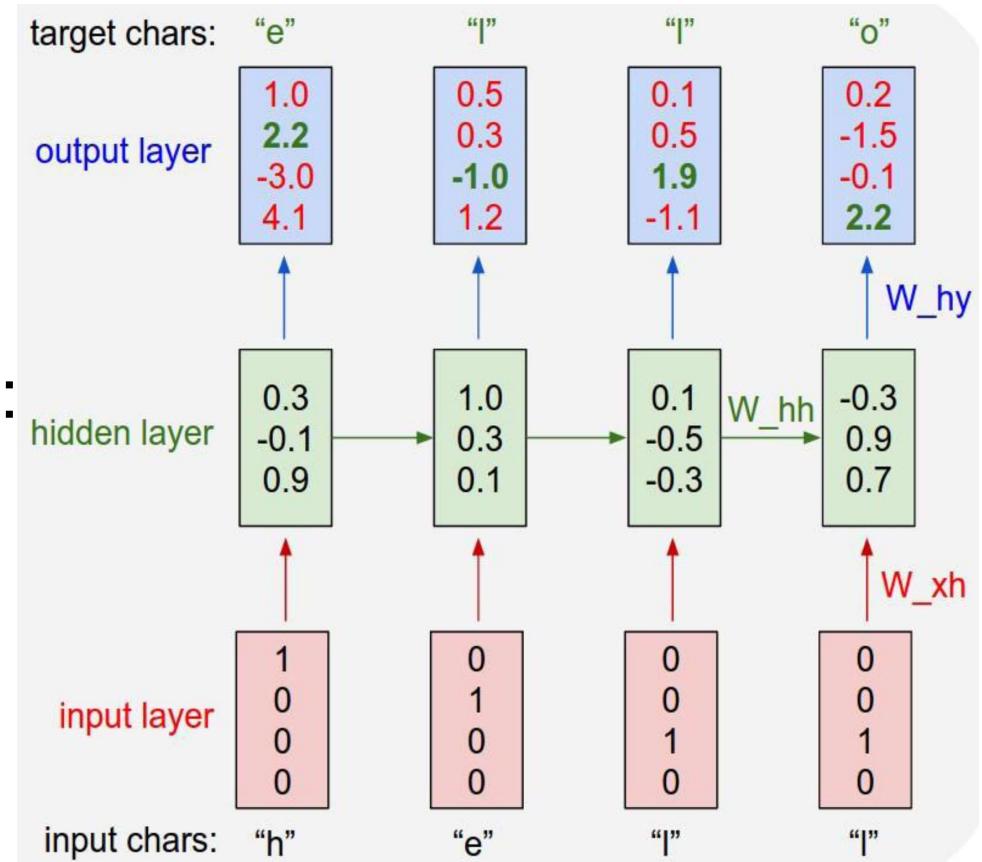
$$\mathbf{o}_t = \phi(\mathbf{W}_{ho}\mathbf{h}_t + \mathbf{b}_o)$$

Example: Character-level Language Model

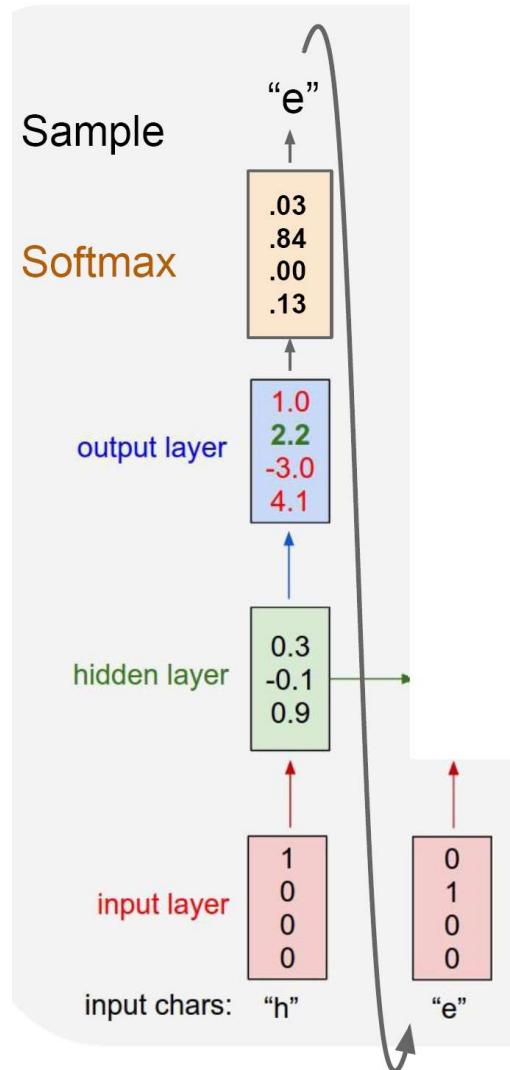
❖ Vocabulary:

[h,e,l,o]

❖ Example **training** sequence:
“hello”



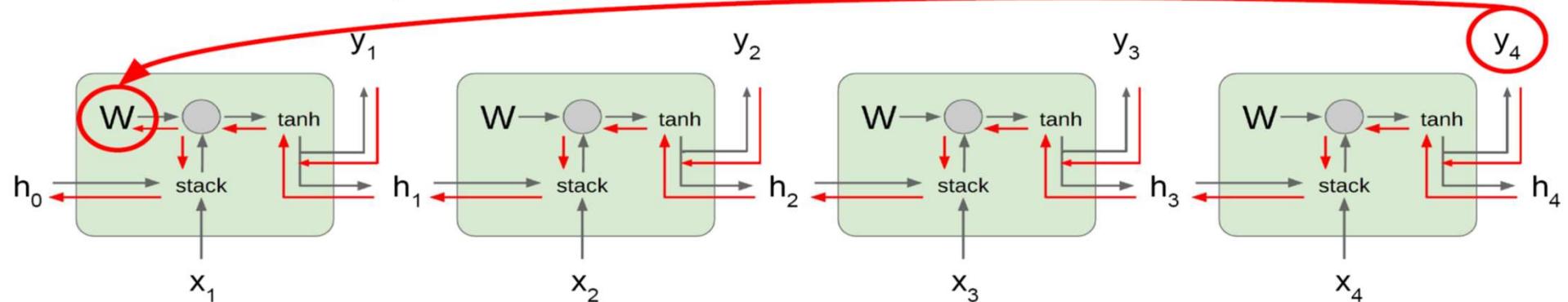
- ❖ At test-time sample characters one at a time, feed back to model



Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
 Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

$$\boxed{\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}}$$

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \frac{\partial h_t}{\partial h_{t-1}} \cdots \frac{\partial h_1}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \boxed{\frac{\partial h_t}{\partial h_{t-1}}} \right) \frac{\partial h_1}{\partial W}$$

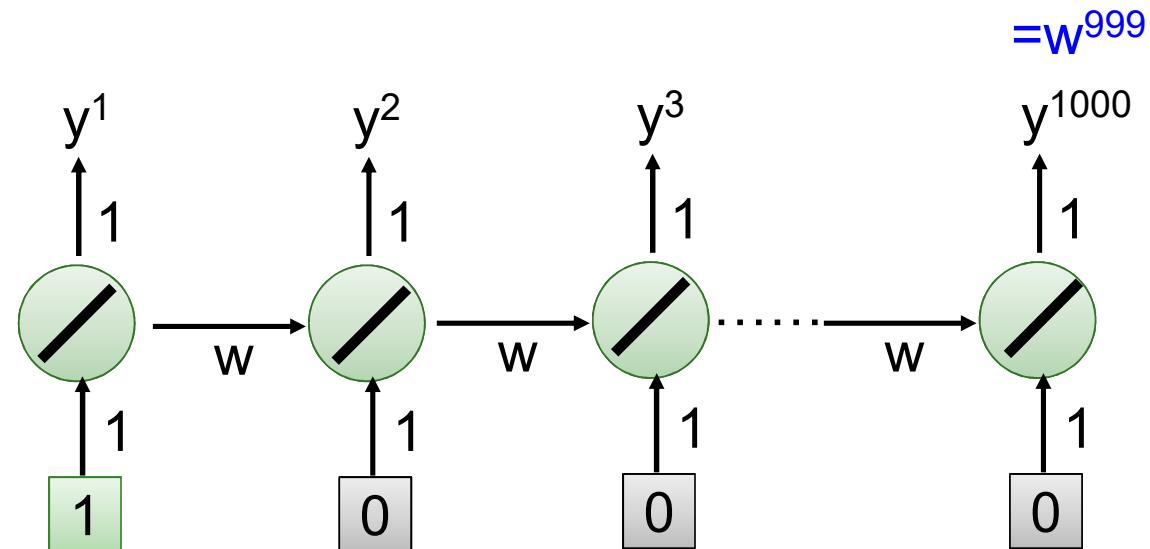
Almost always < 1
 Vanishing gradients

What if we assumed no non-linearity?
 Largest singular value > 1 : Exploding gradients
 Largest singular value < 1 : Vanishing gradients

Gradient Explode/ Vanish

$w = 1$	\rightarrow	$y^{1000} = 1$	Large $\partial L/\partial w$	Small Learning rate?
$w = 1.01$	\rightarrow	$y^{1000} \approx 20000$		
$w = 0.99$	\rightarrow	$y^{1000} \approx 0$	small $\partial L/\partial w$	Large Learning rate?
$w = 0.01$	\rightarrow	$y^{1000} \approx 0$		

Toy Example



Recurrent Neural Networks(RNN)

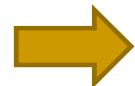
- ✓ Suitable for processing sequences, often applied to the processing of text.
- ❖ Has a problem about gradient vanishing
- ❖ Can not store long-term memory

Modify the structure

❖ Long Short-term Memory (LSTM)

◆ Can deal with gradient vanishing (not gradient explode)

- Memory and input are added
- The influence never disappears unless the forget gate is closed

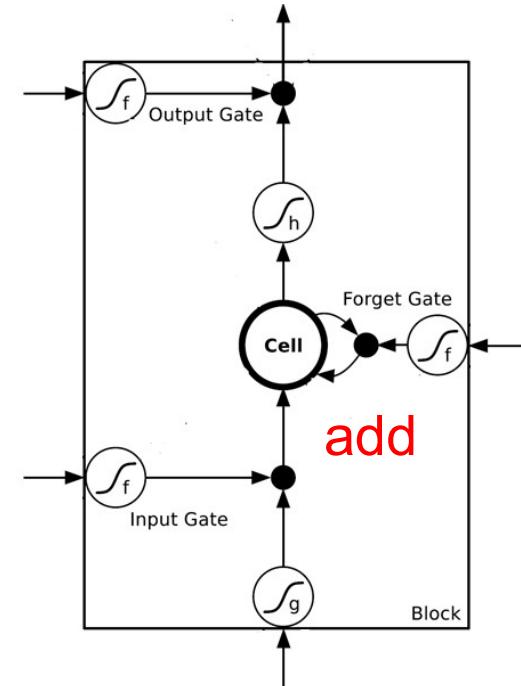


No Gradient vanishing

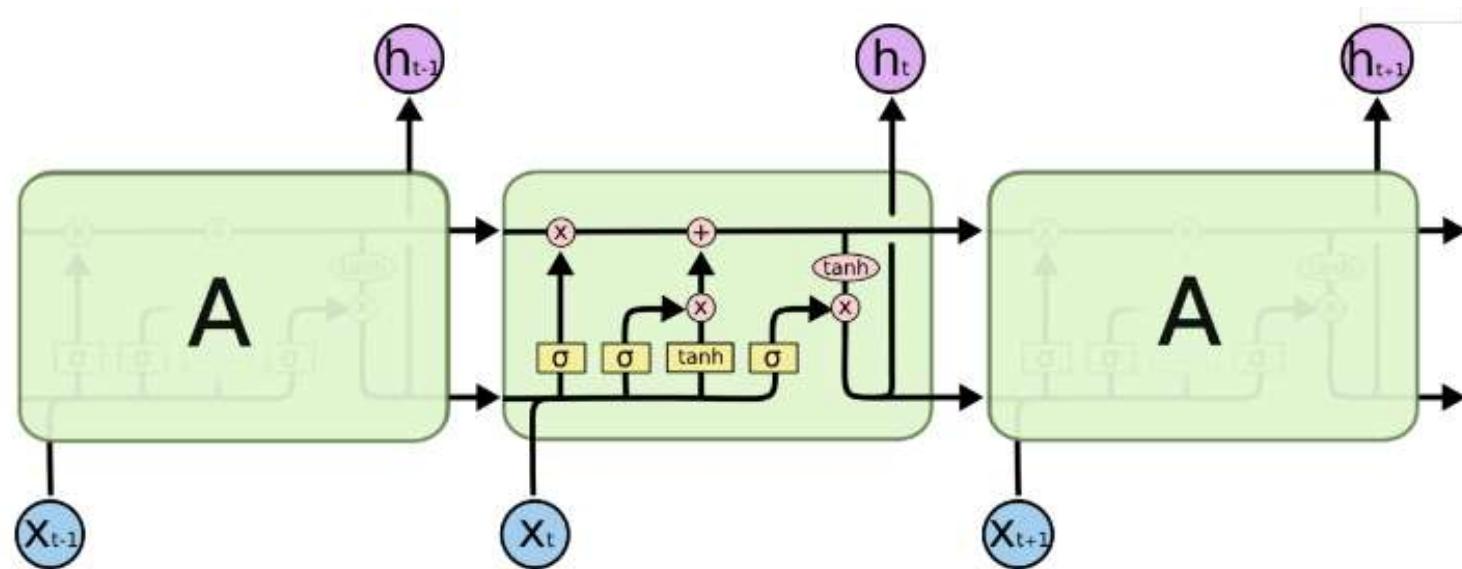
(If the forget gate is opened.)

❖ Gated Recurrent Unit(GRU):

simpler than LSTM



Long Short Term Memory



Hochreiter & Schmidhuber Long Short-term Memory[J]. Neural Computation, 1997,9(8):1735-1780.

Long Short Term Memory

❖ Forget gate

- ◆ Shrink values towards zero, **whether to erase cell**

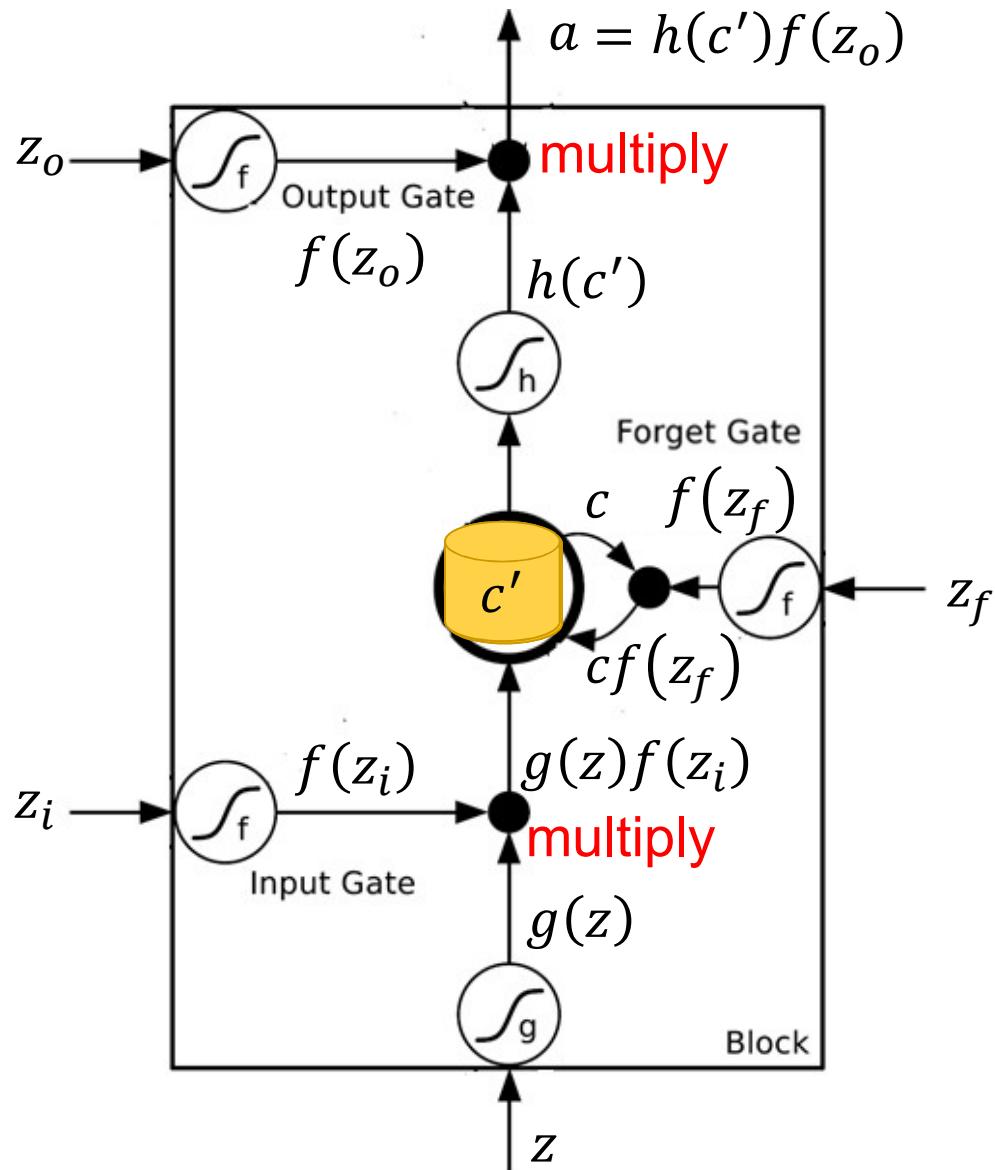
❖ Input gate

- ◆ Decide whether we should ignore the input data

❖ Output gate

- ◆ Decide whether the hidden state is used for the output generated by the LSTM

❖ Hidden state and Memory cell



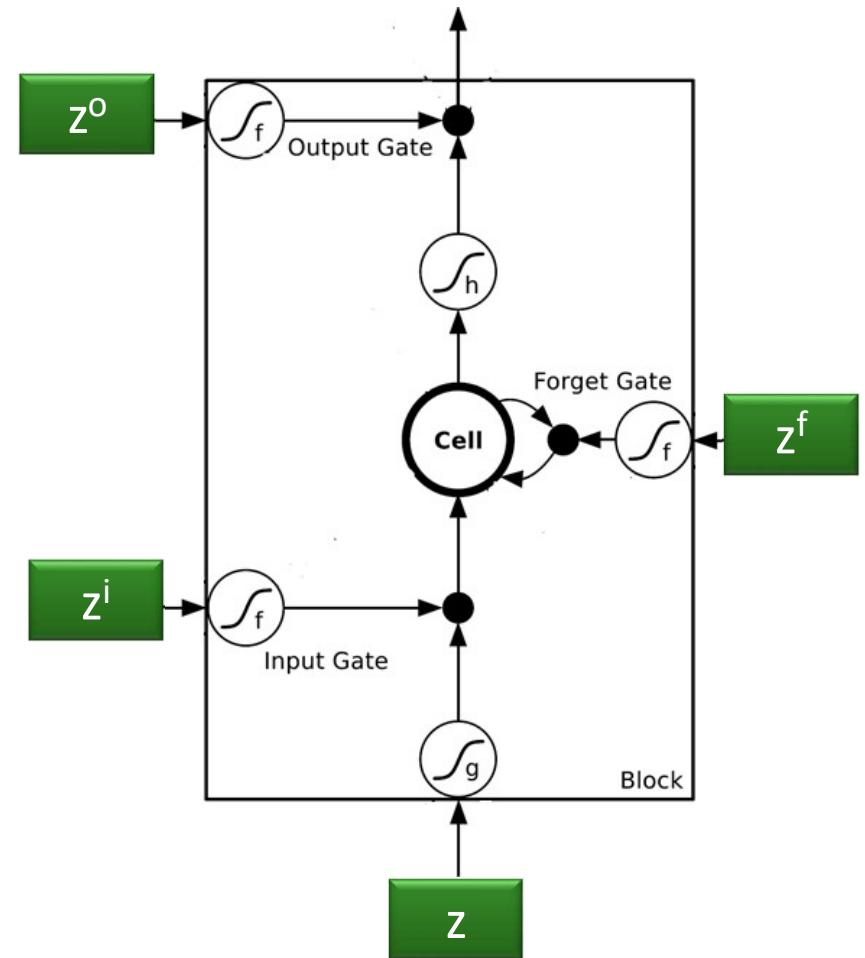
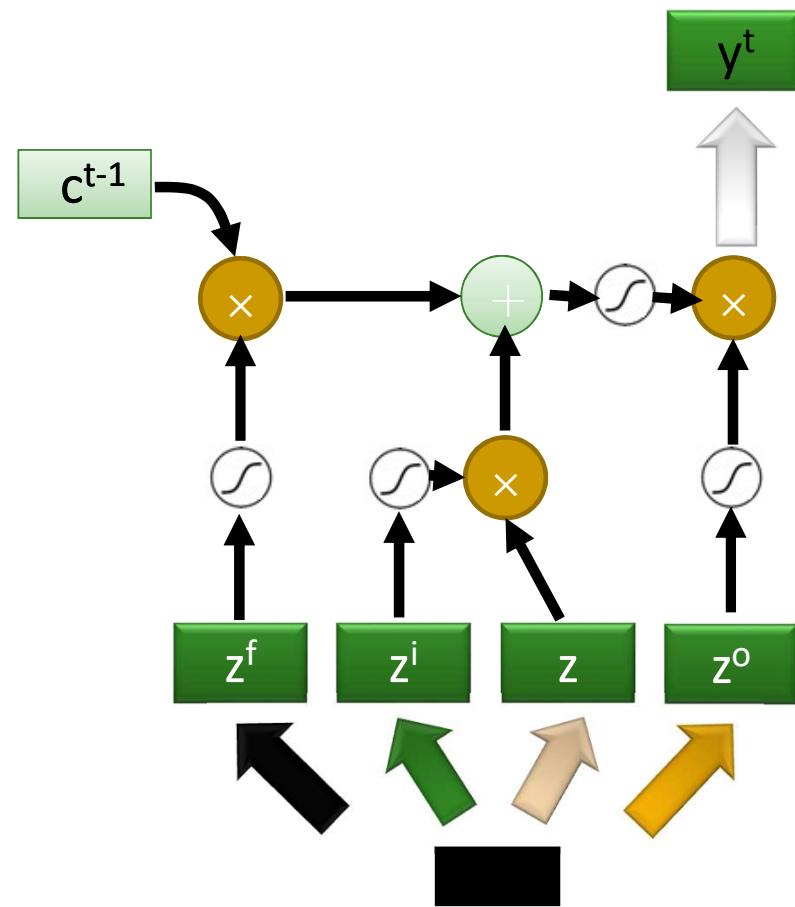
Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

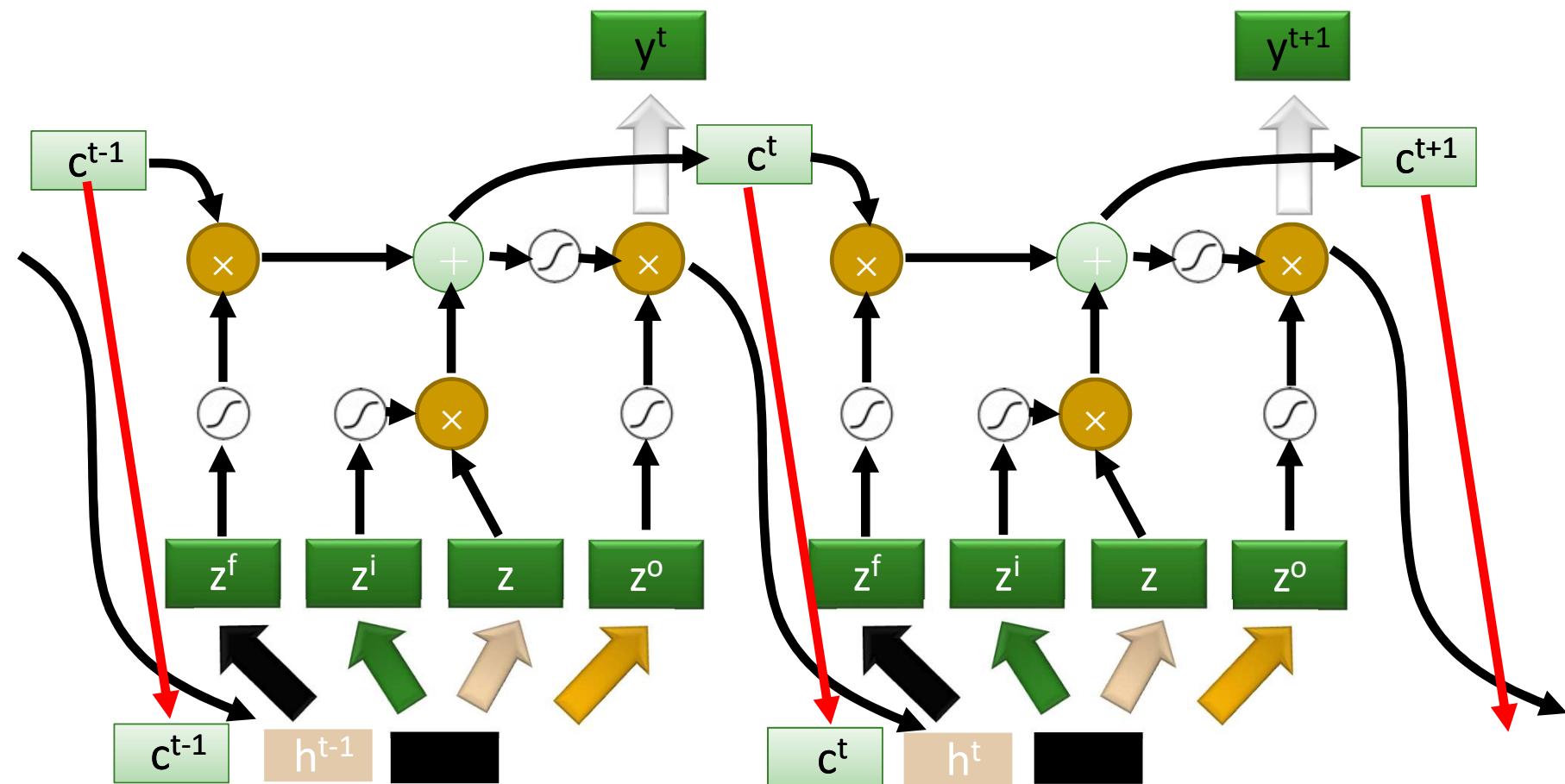
$$c' = g(z)f(z_i) + cf(z_f)$$

LSTM



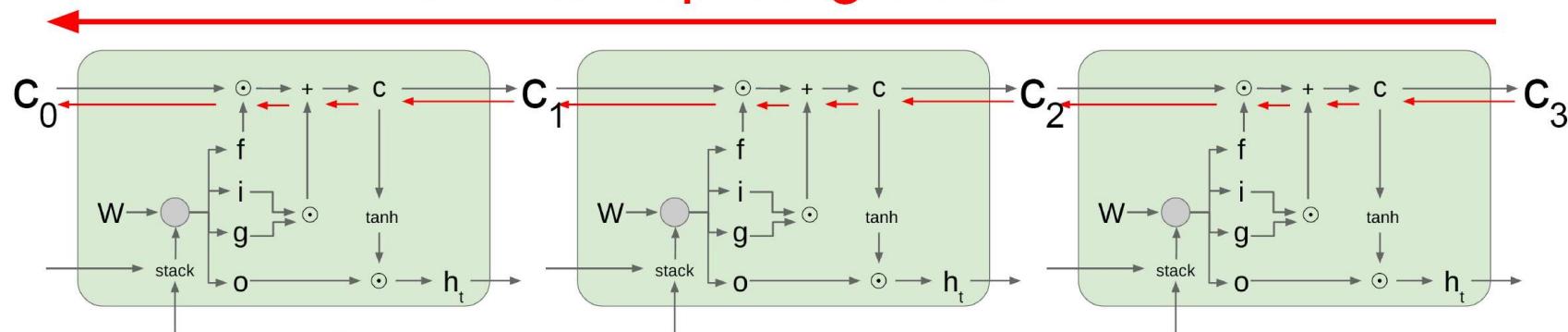
LSTM

Extension: “peephole”

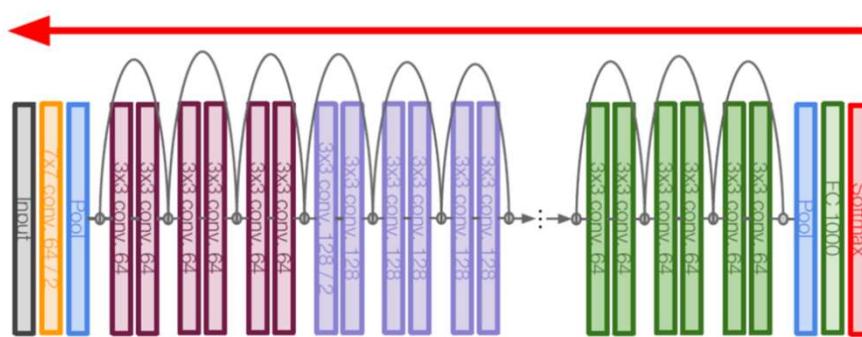


Long Short Term Memory (LSTM): Gradient Flow [Hochreiter et al., 1997]

Uninterrupted gradient flow!



Similar to ResNet!



In between:
Highway Networks

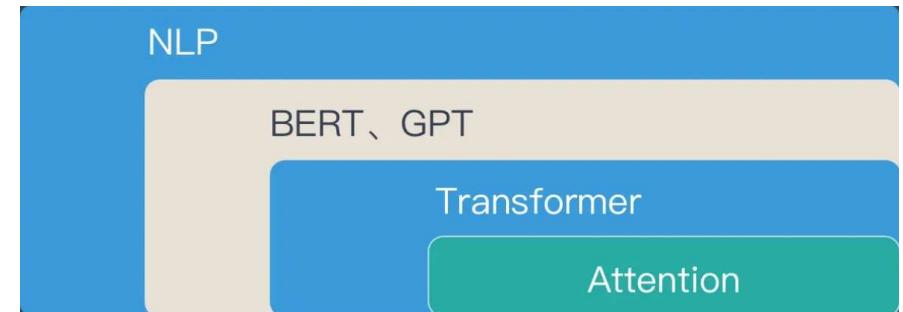
$$g = T(x, W_T)$$

$$y = g \odot H(x, W_H) + (1 - g) \odot x$$

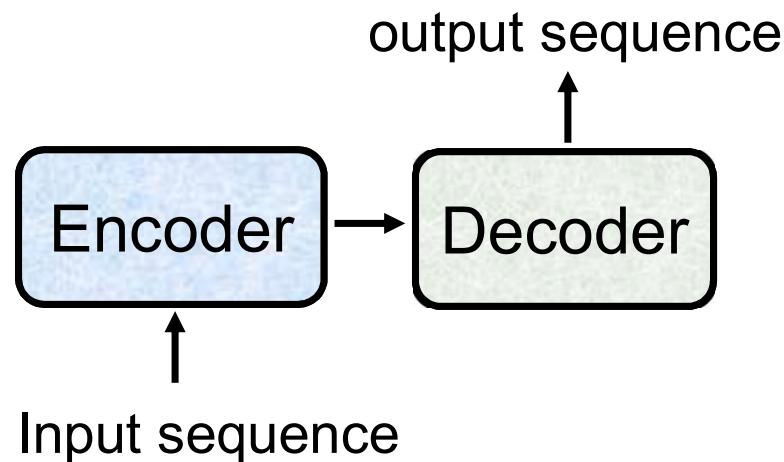
Srivastava et al, "Highway Networks",
ICML DL Workshop 2015

Outline

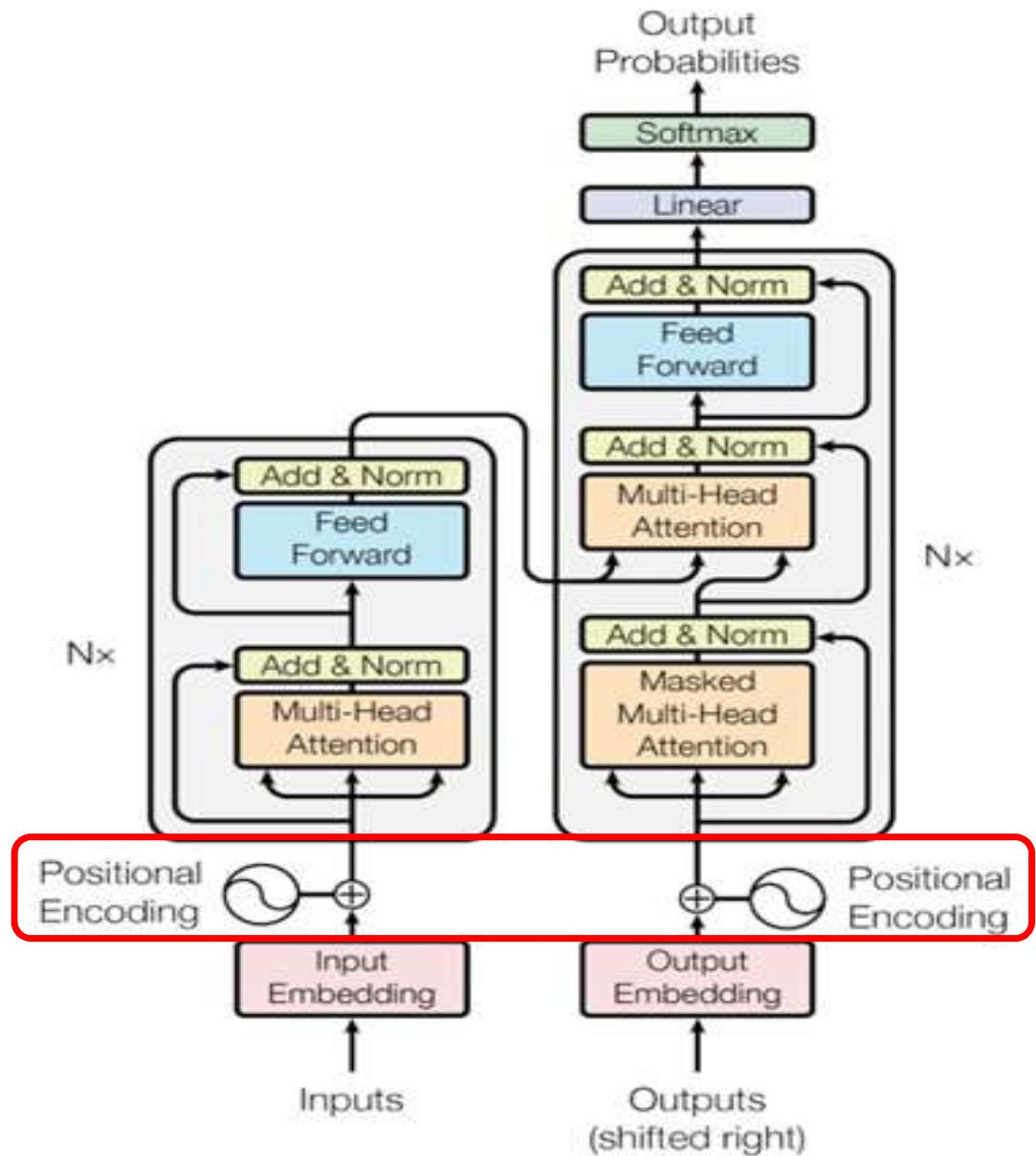
- ❖ Convolutional Neural Network (CNN)
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
 - ◆ Self-attention
 - ◆ Encoder & Decoder
- ❖ Deep Learning Frameworks



Transformer



Attention Is All You Need
<https://arxiv.org/abs/1706.03762>



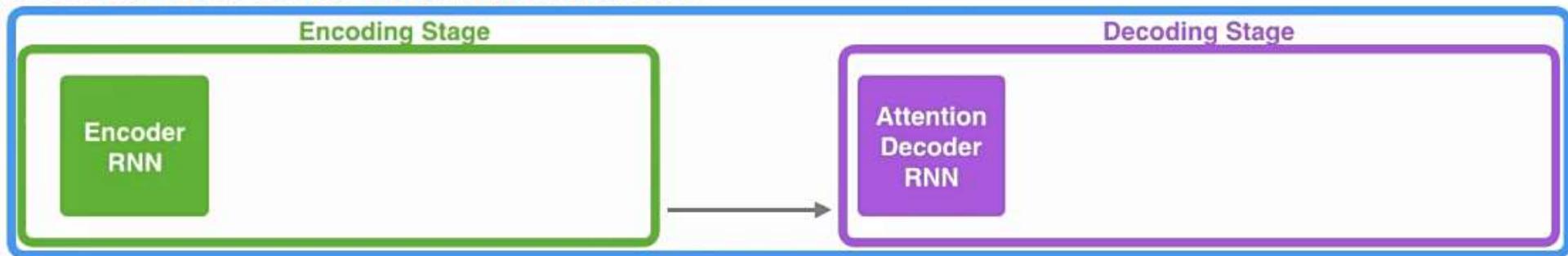
Attention



Sequence to sequence model with attention

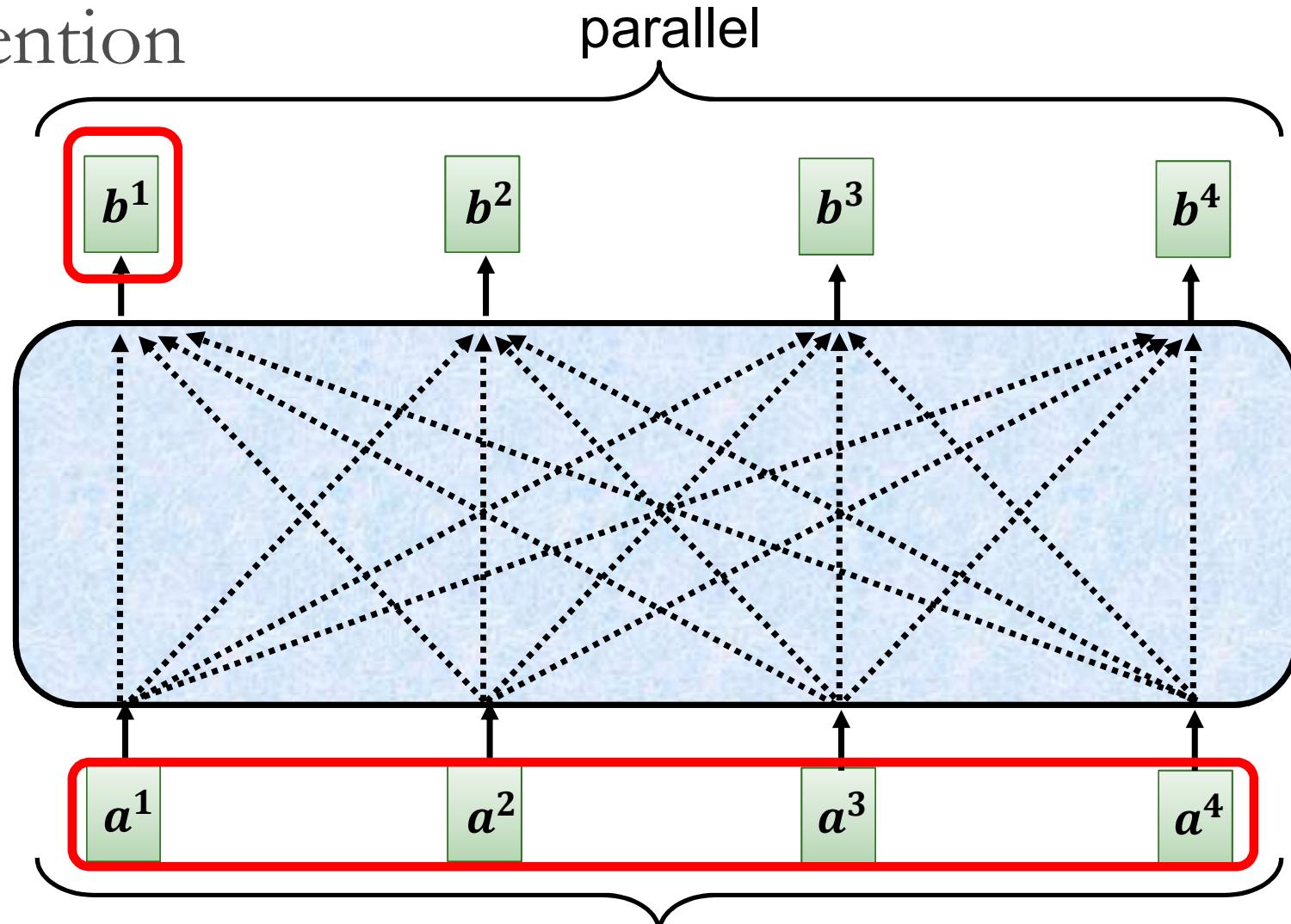
gif5.net

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



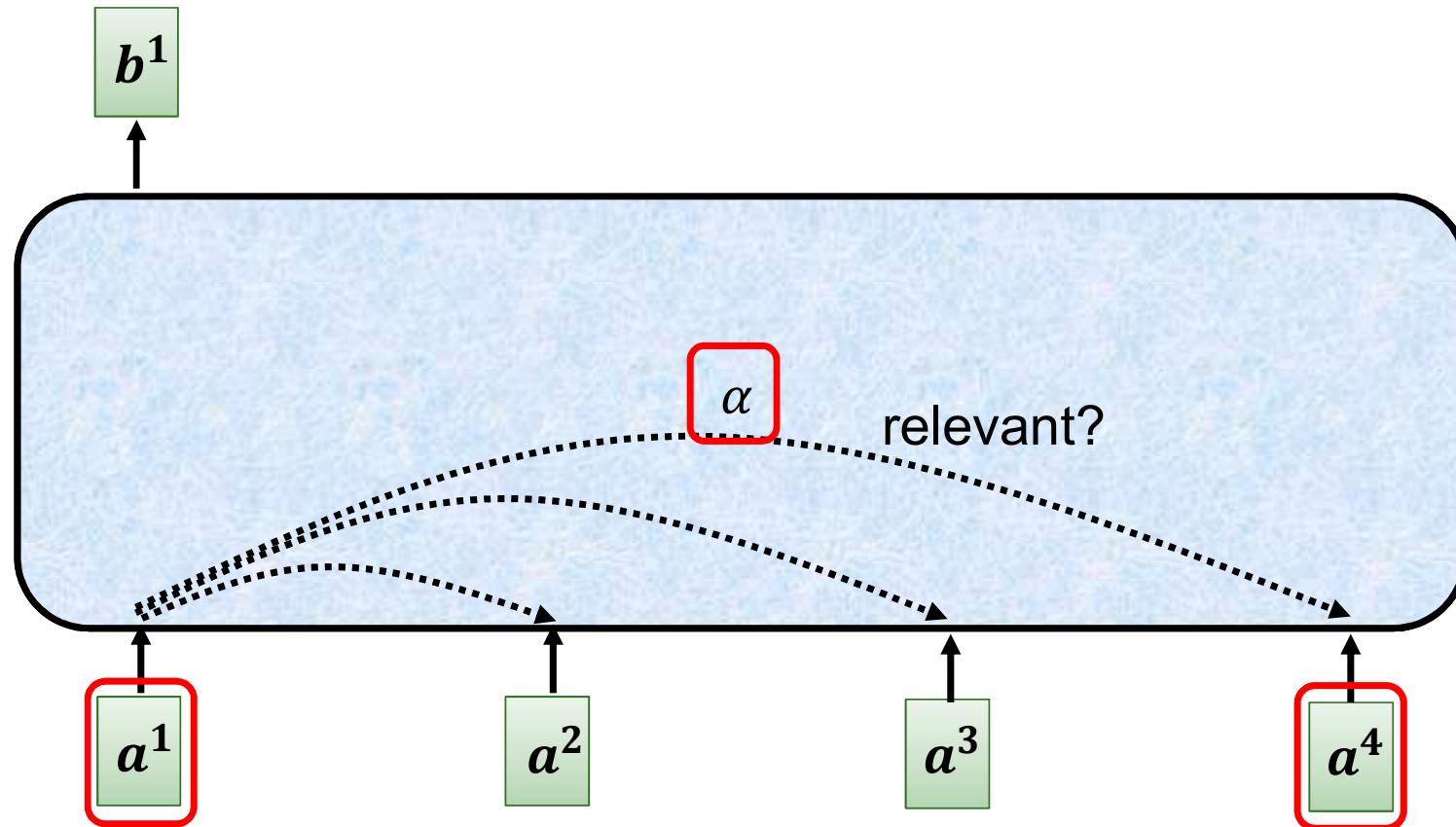
Je suis étudiant

Self-attention



Can be either **input** or the output of a **hidden layer**

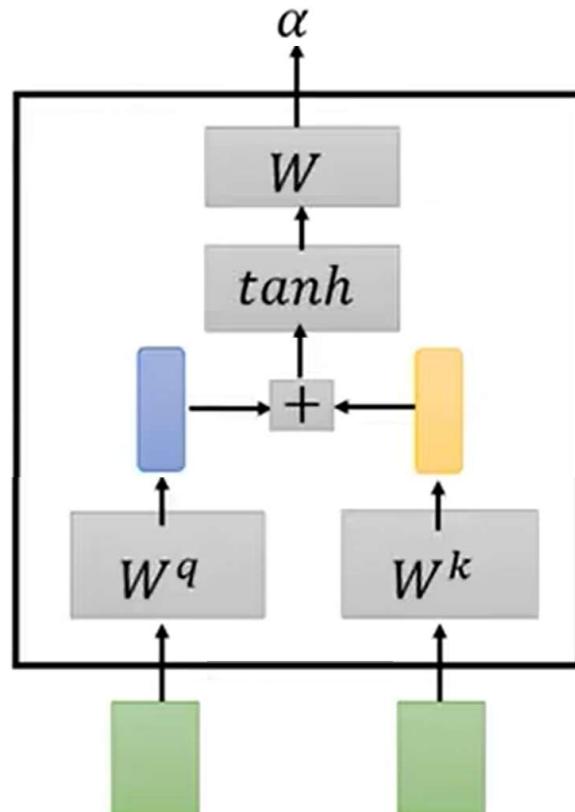
Self-attention



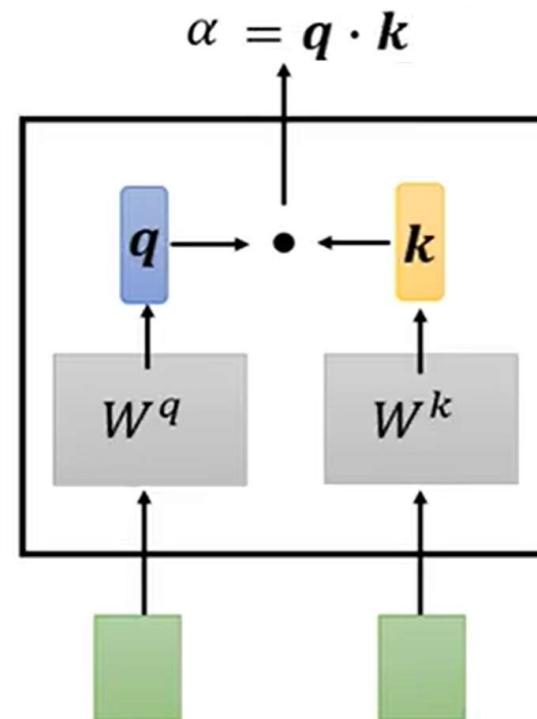
Find the relevant vectors in a sequence

Self-attention

Additive

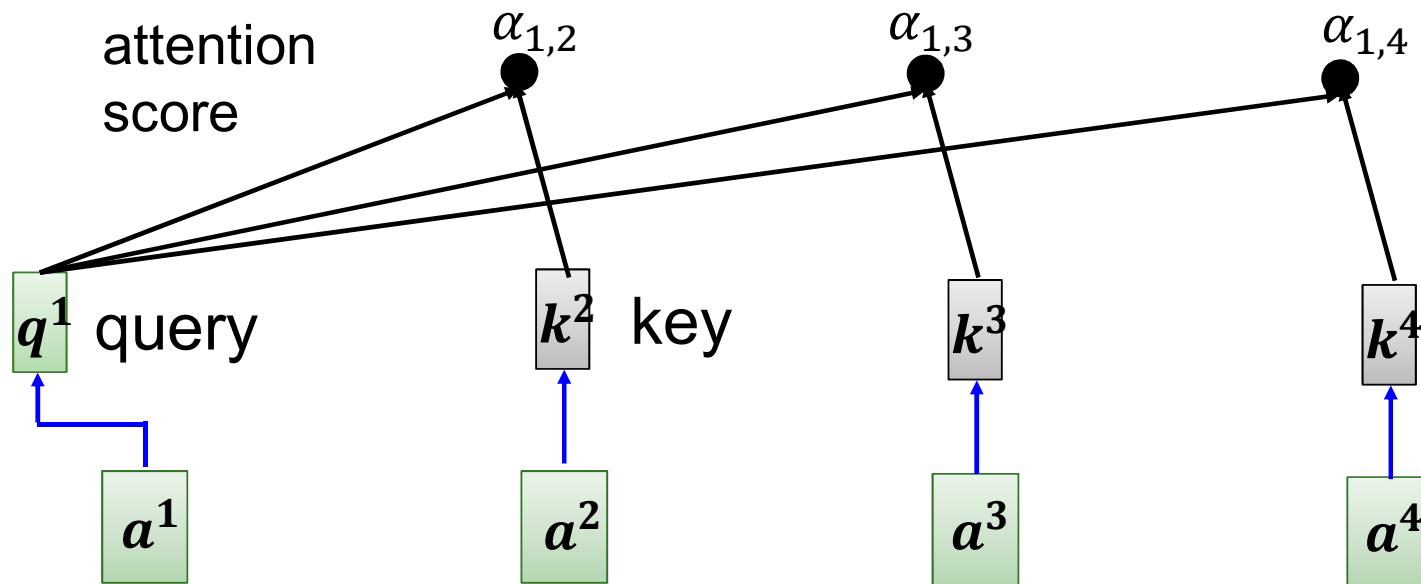


Dot-Product



Self-attention

$$\alpha_{1,2} = q^1 \cdot k^2 \quad \alpha_{1,3} = q^1 \cdot k^3 \quad \alpha_{1,4} = q^1 \cdot k^4$$



$$q^1 = W^q a^1$$

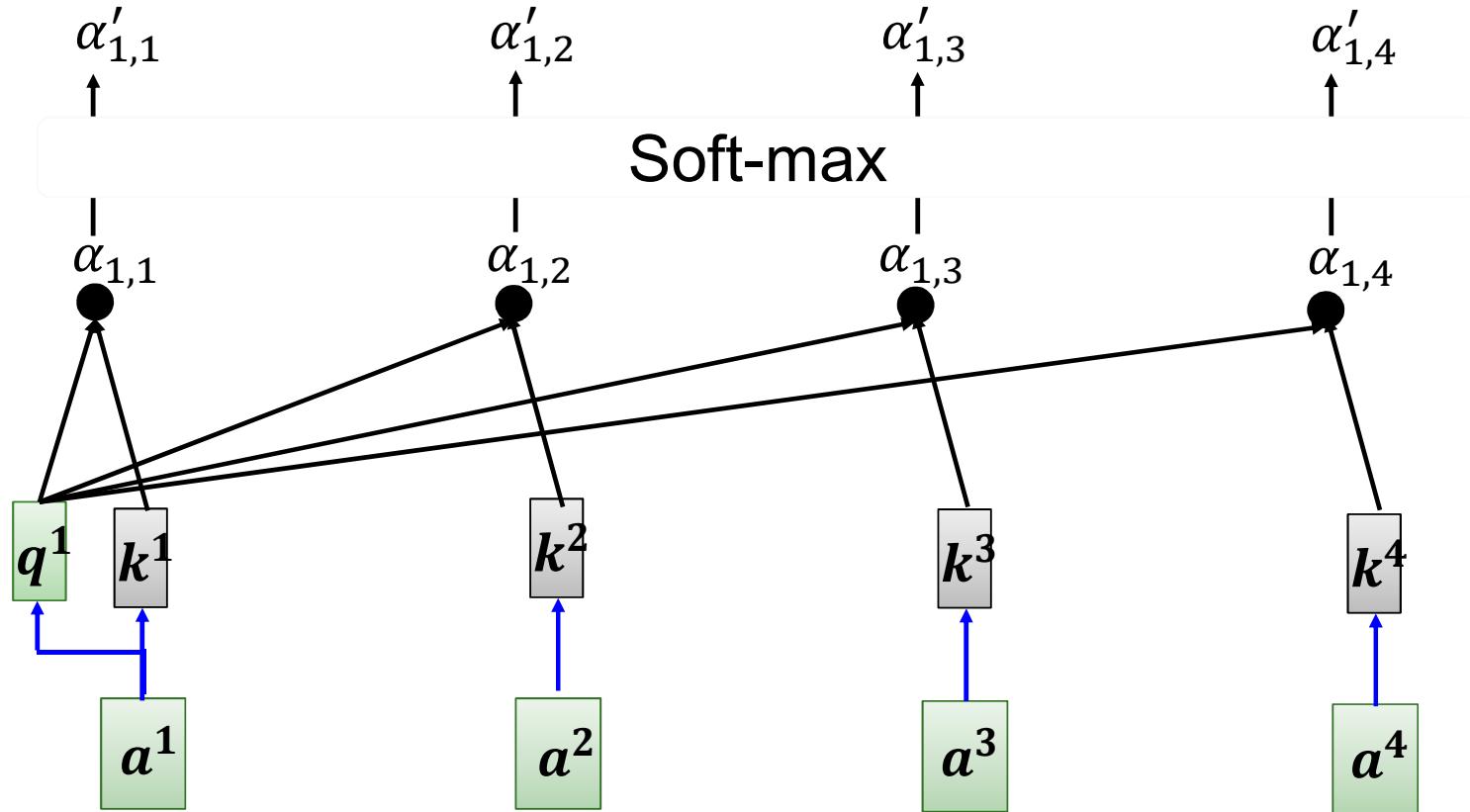
$$k^2 = W^k a^2$$

$$k^3 = W^k a^3$$

$$k^4 = W^k a^4$$

Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



$$q^1 = W^q a^1$$

$$k^2 = W^k a^2$$

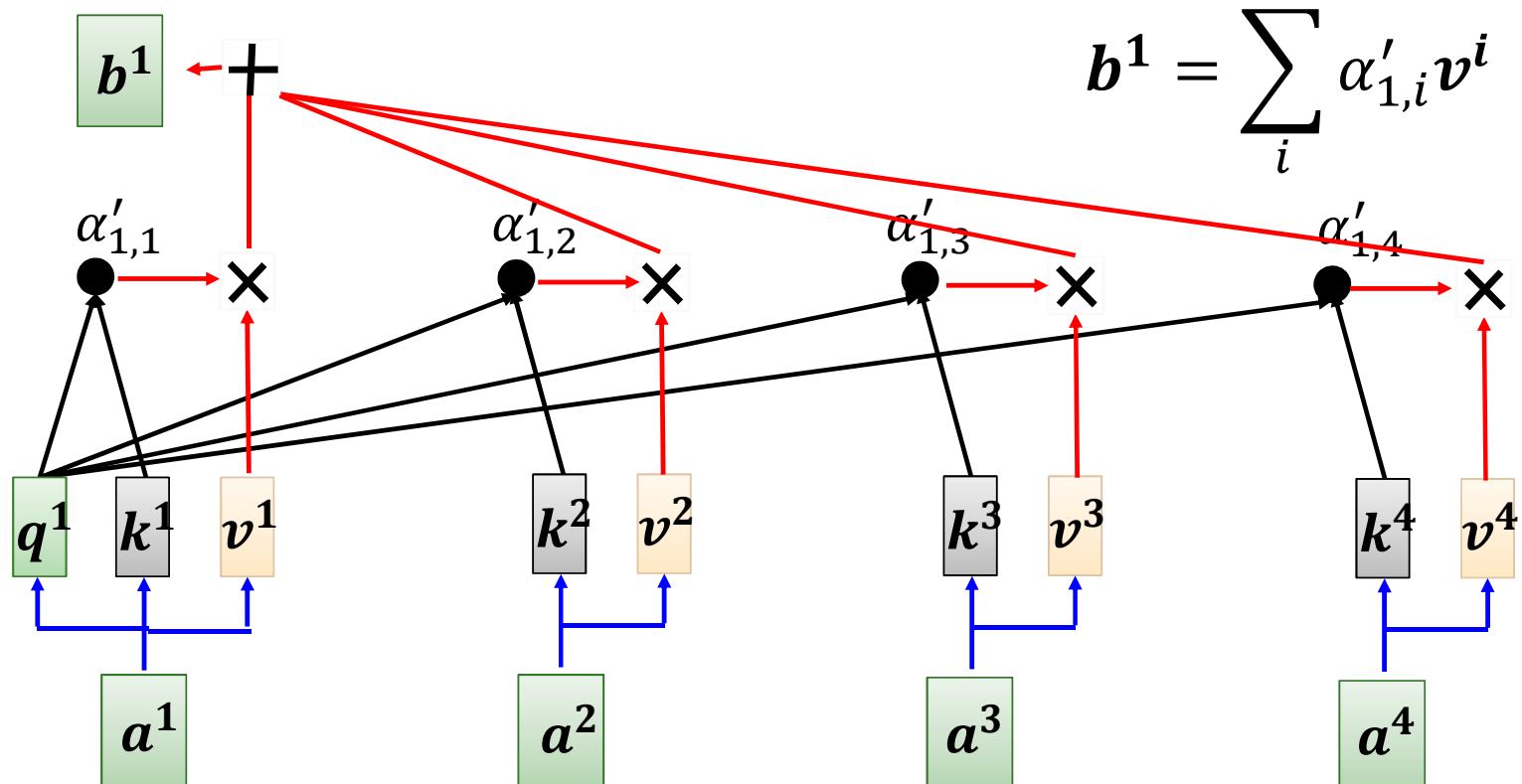
$$k^3 = W^k a^3$$

$$k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

Self-attention

Extract information based on attention scores



$$v^1 = W^v a^1$$

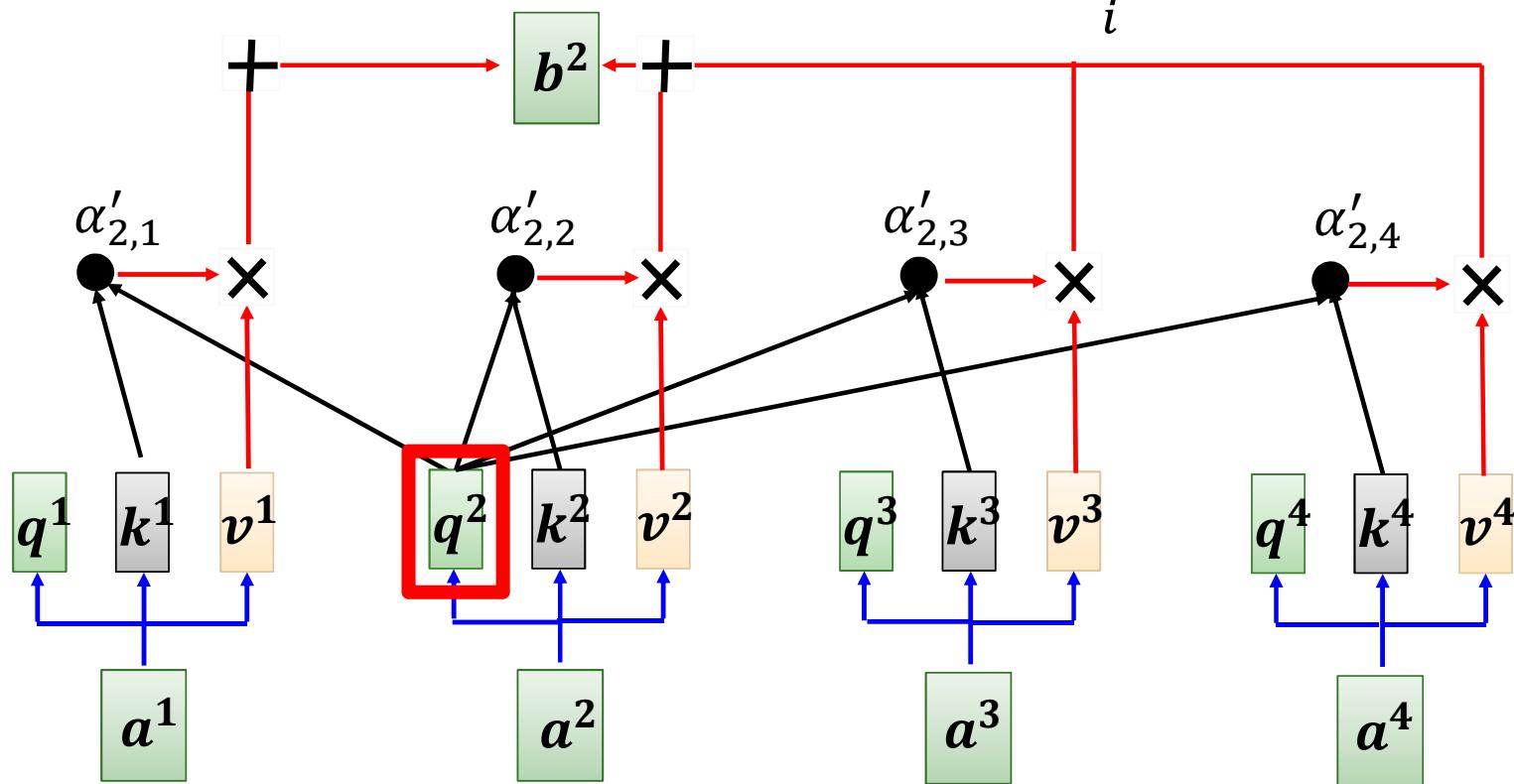
$$v^2 = W^v a^2$$

$$v^3 = W^v a^3$$

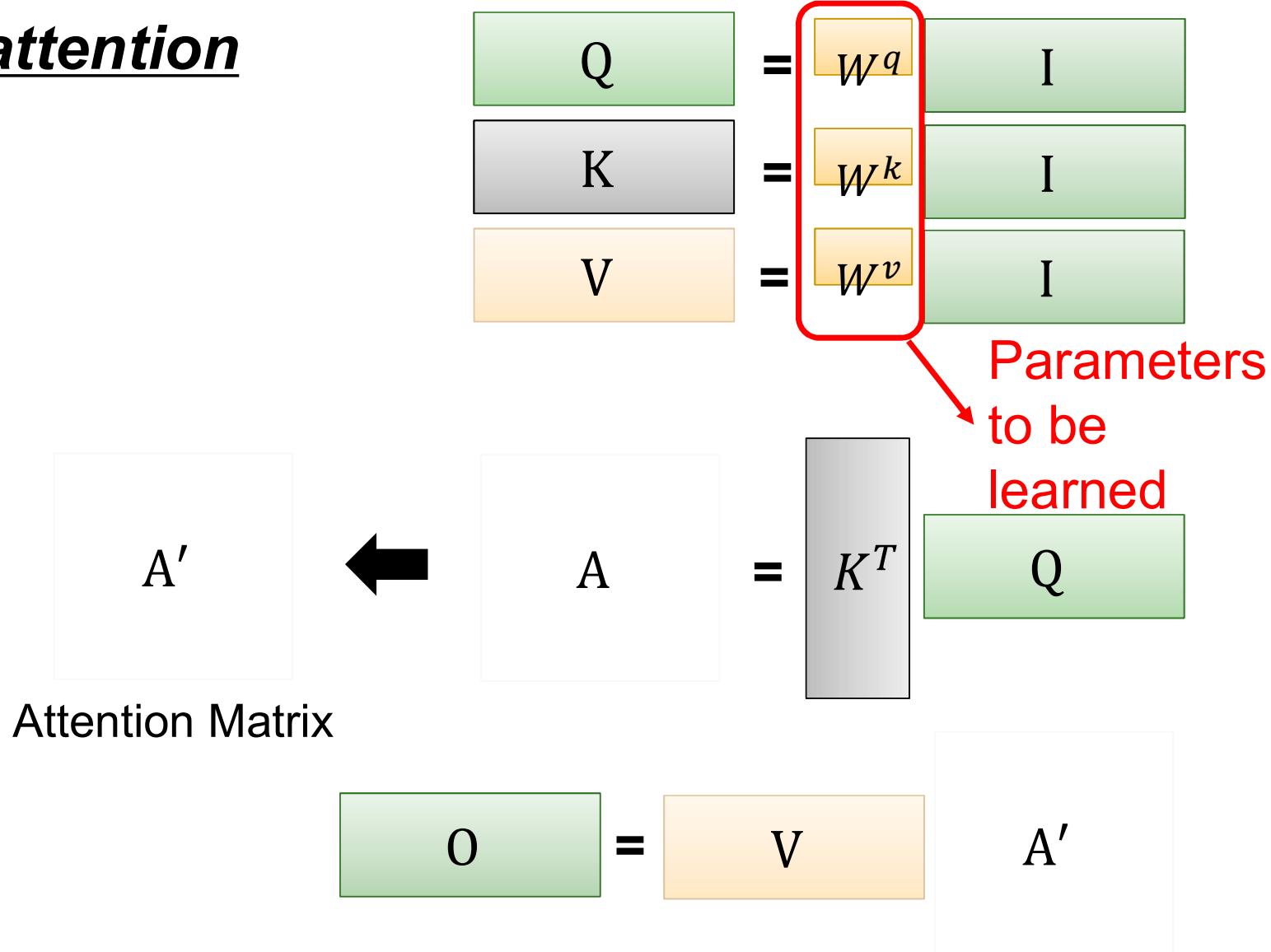
$$v^4 = W^v a^4$$

Self-attention

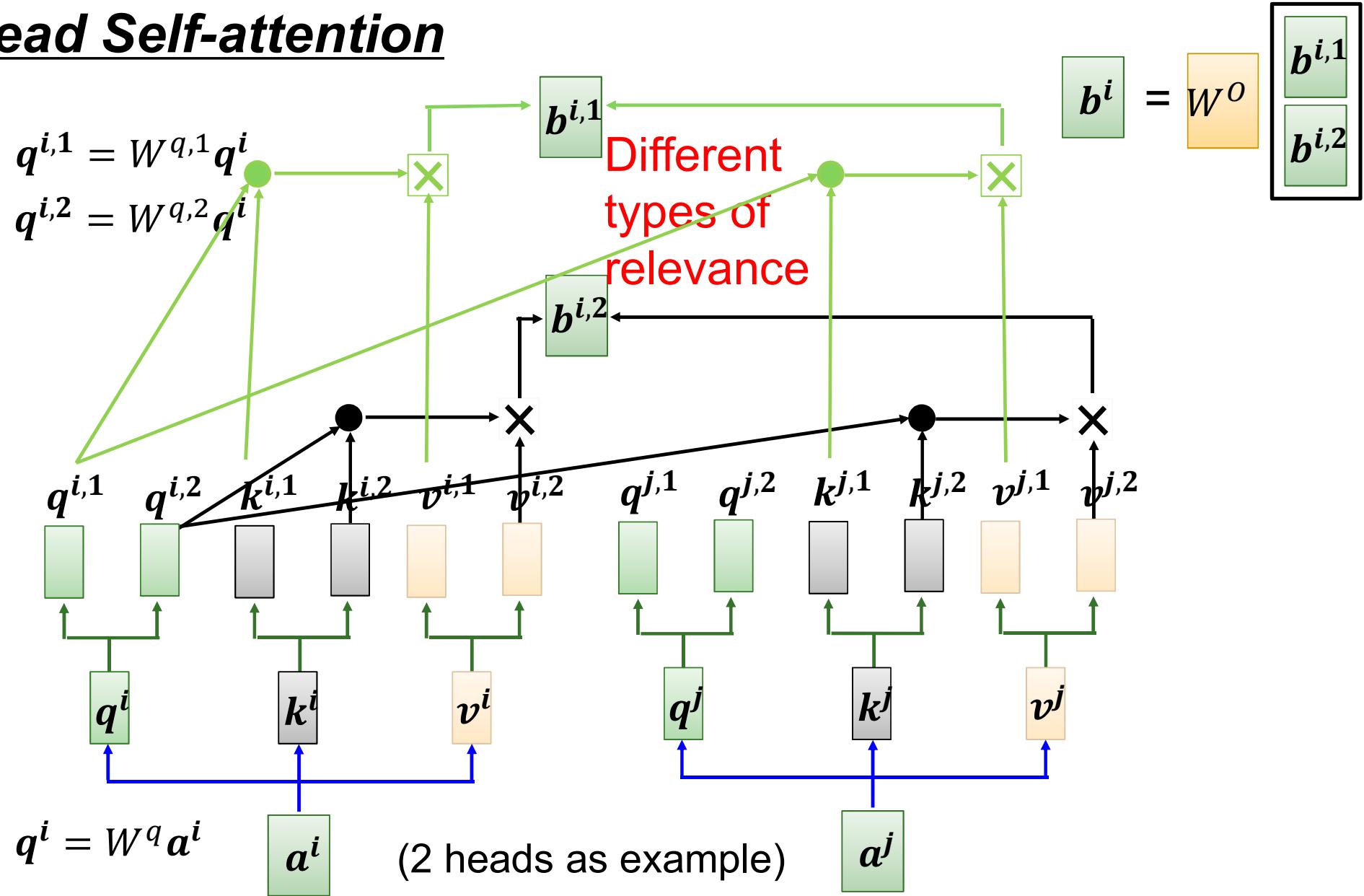
$$b^2 = \sum_i \alpha'_{2,i} v^i$$



Self-attention



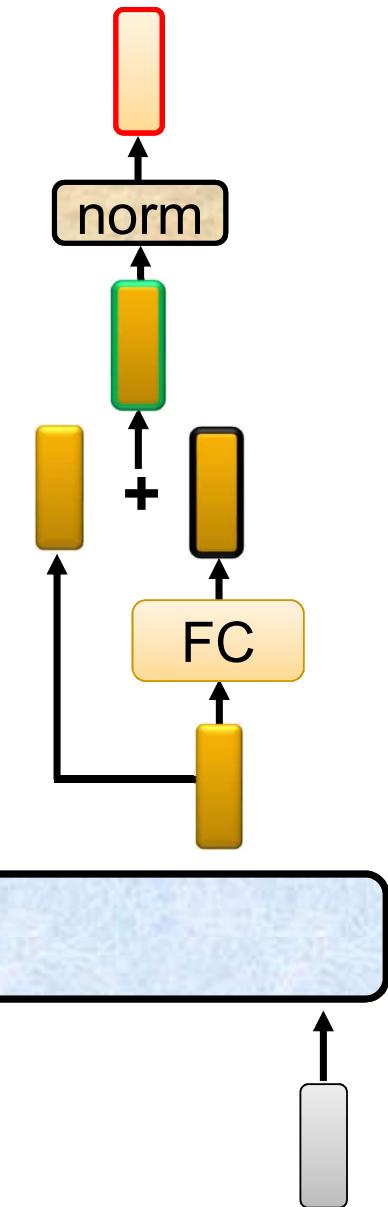
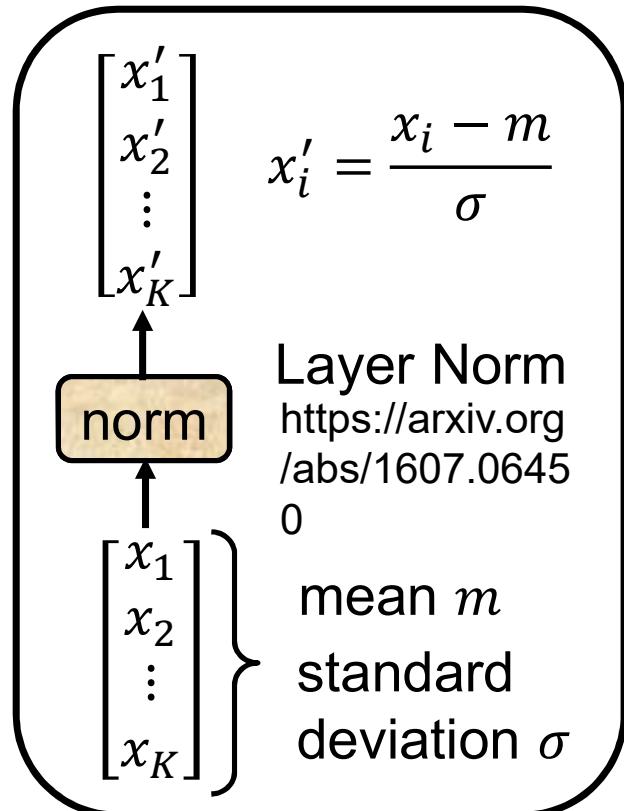
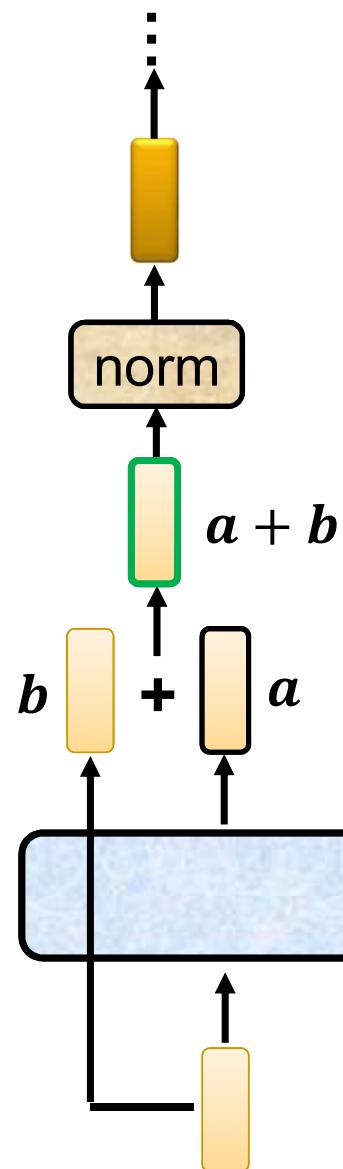
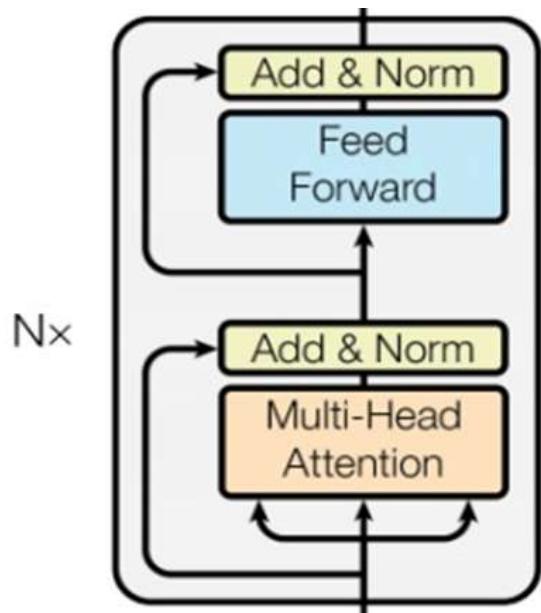
Multi-head Self-attention

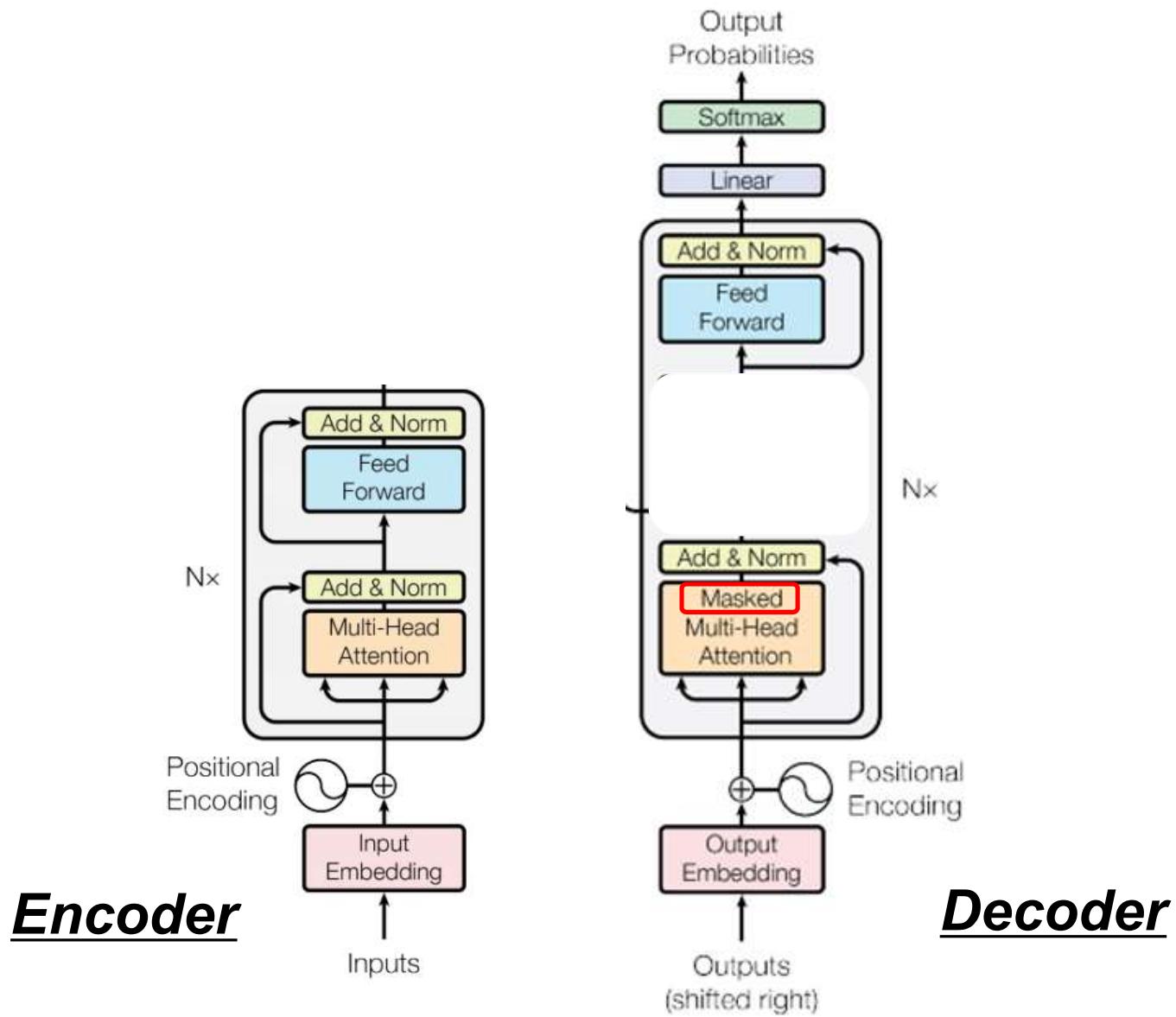


Outline

- ❖ Convolutional Neural Network (CNN)
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
 - ◆ Self-attention
 - ◆ Encoder & Decoder
- ❖ Deep Learning Frameworks

Encoder

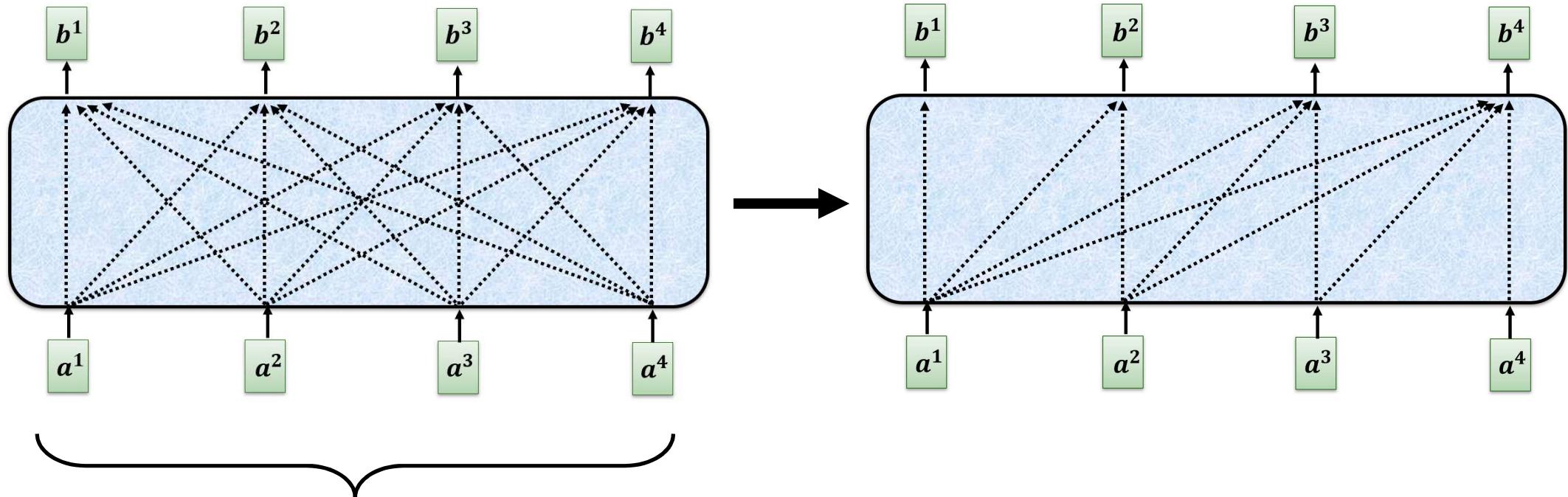




Self-attention



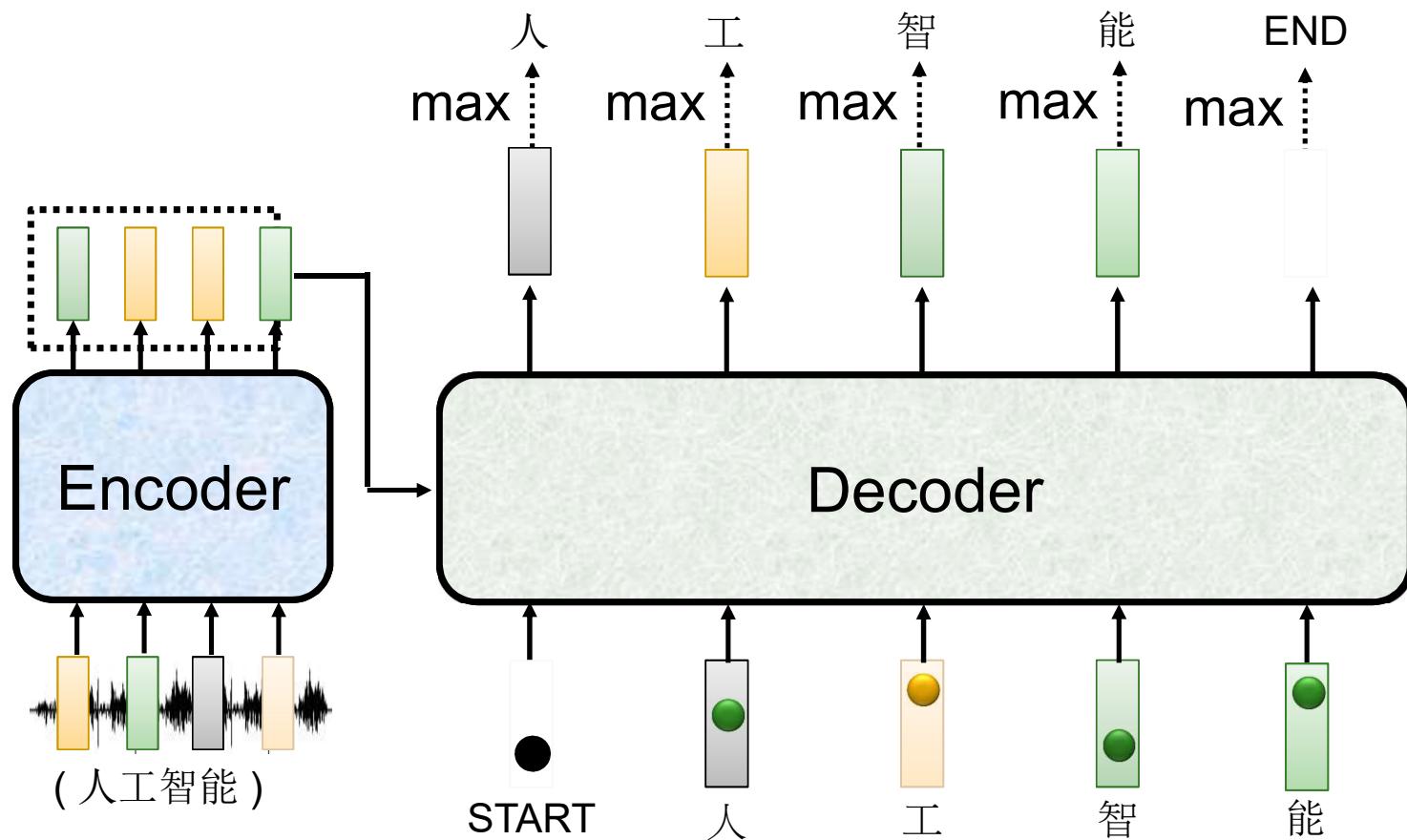
Masked Self-attention



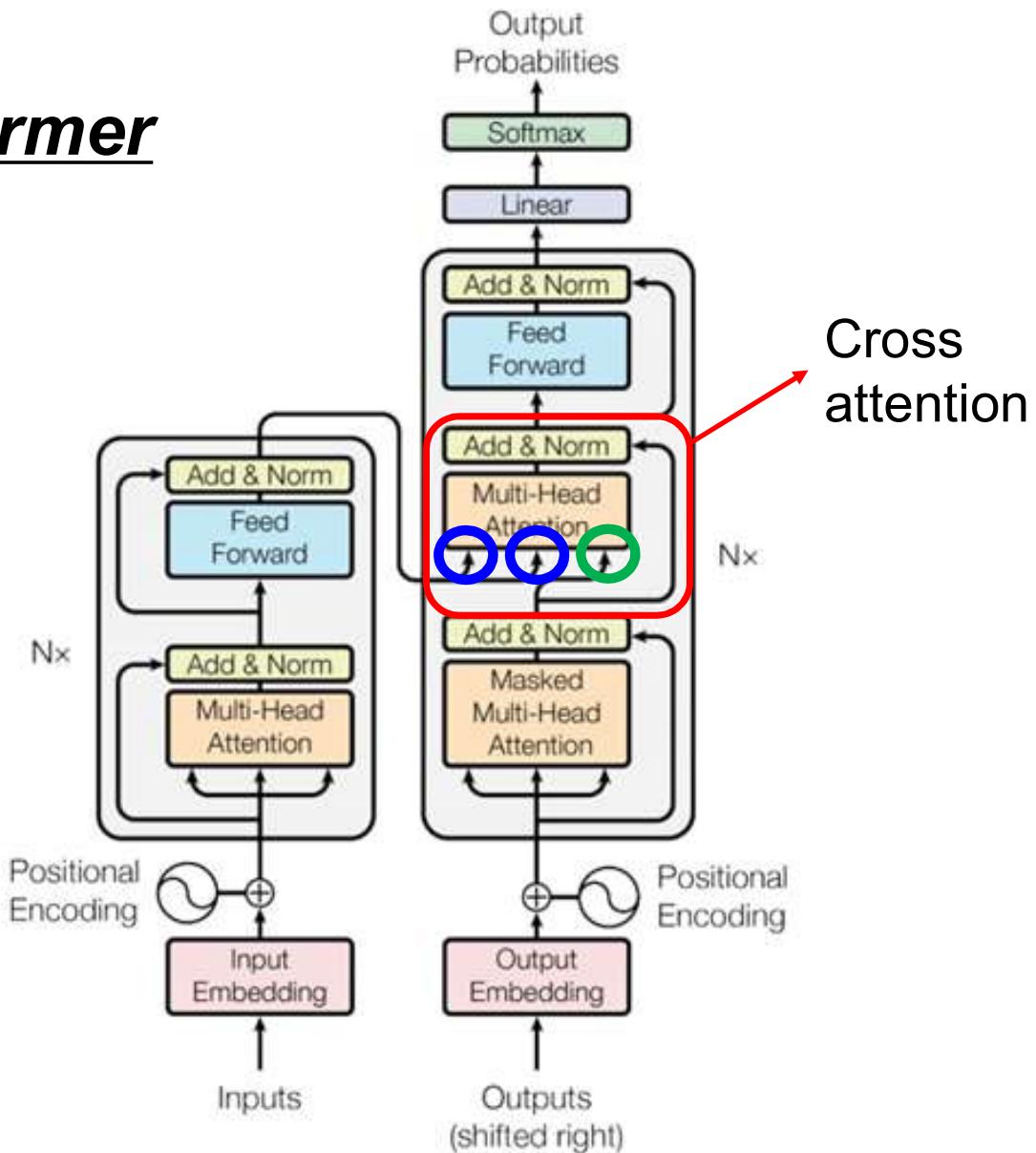
Can be either **input** or **a hidden layer**

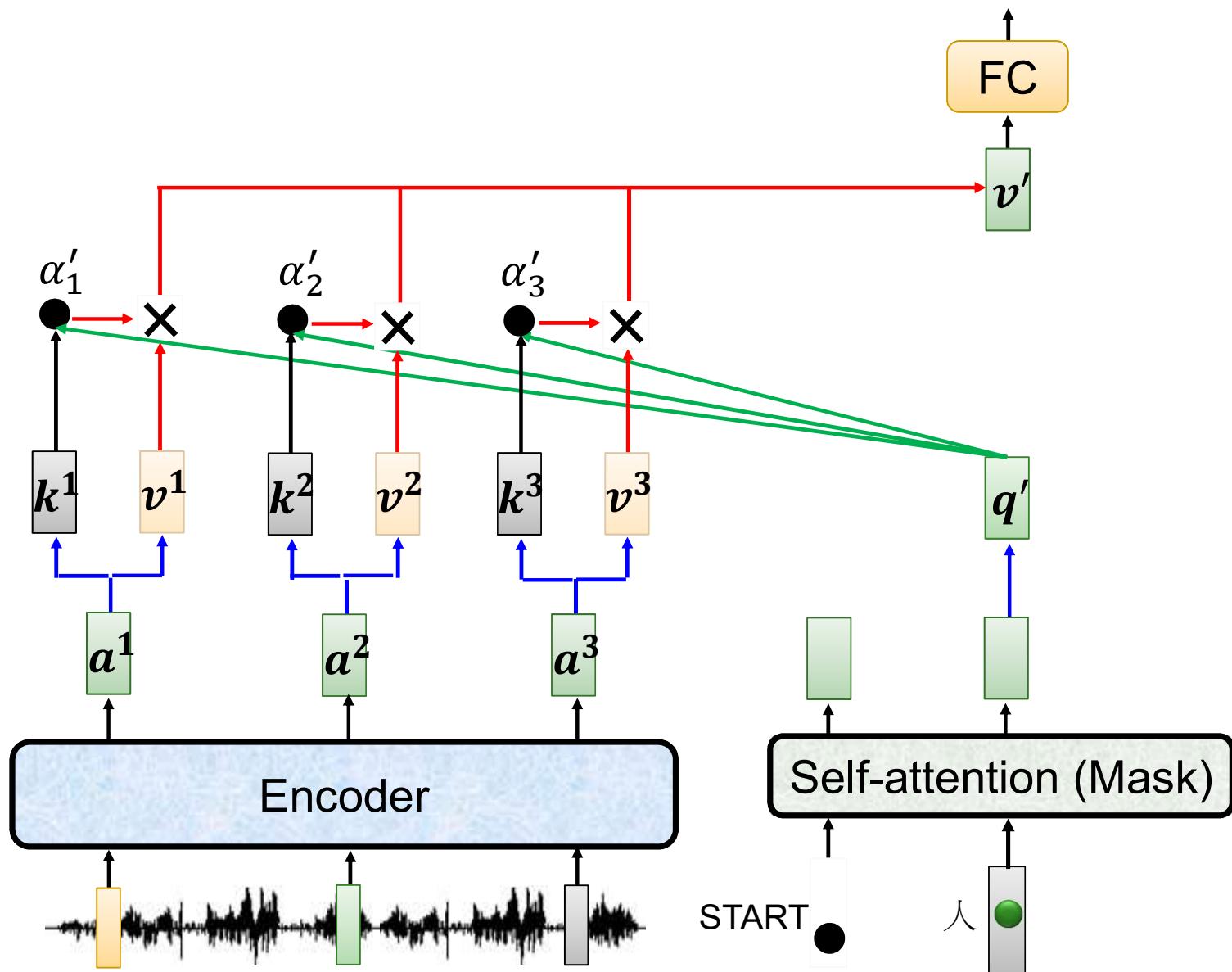
Autoregressive

Stop at here!



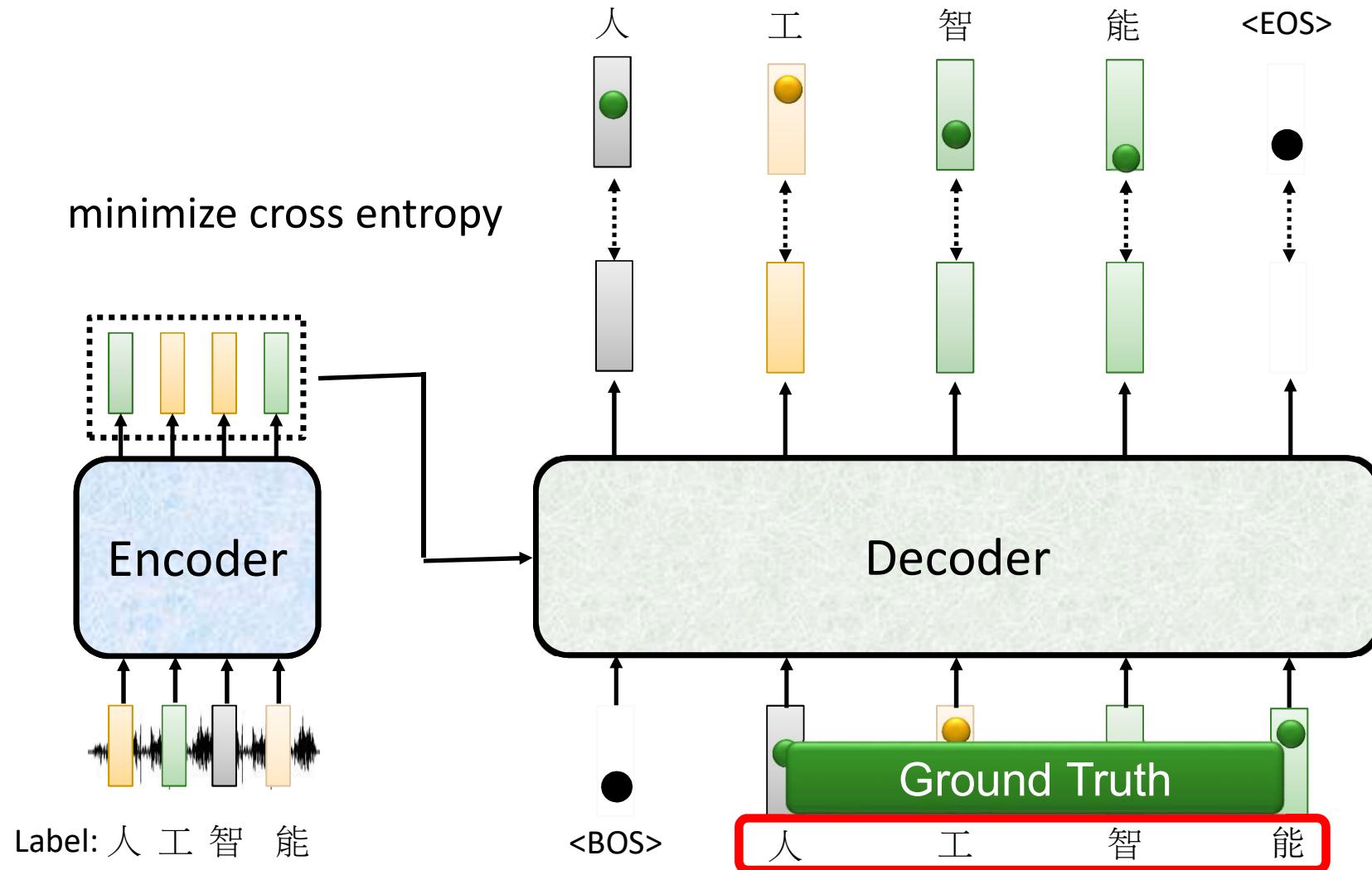
Transformer



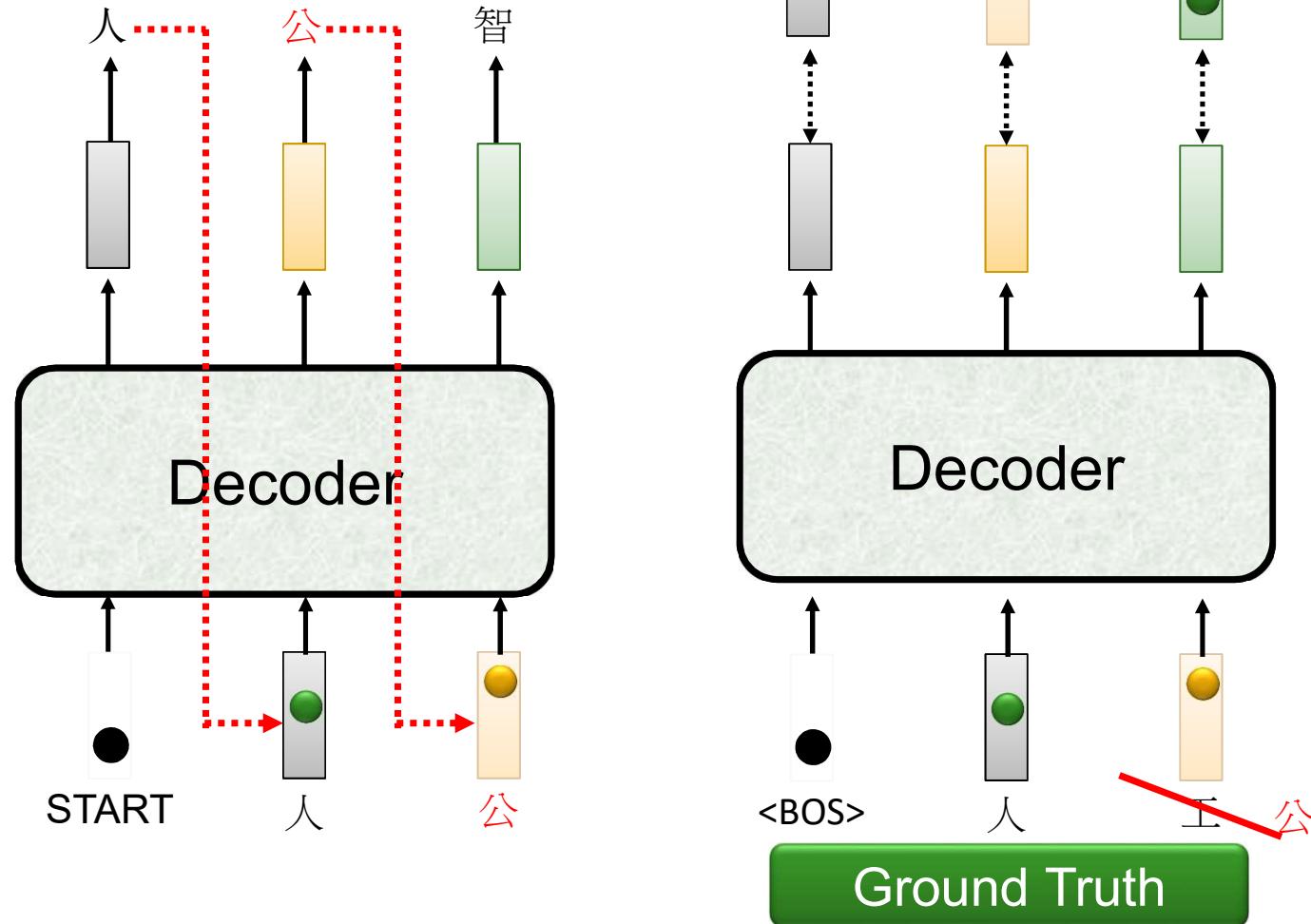


Teacher Forcing: using the ground truth as input.

Training



There is a mismatch! 😞
exposure bias



Outline

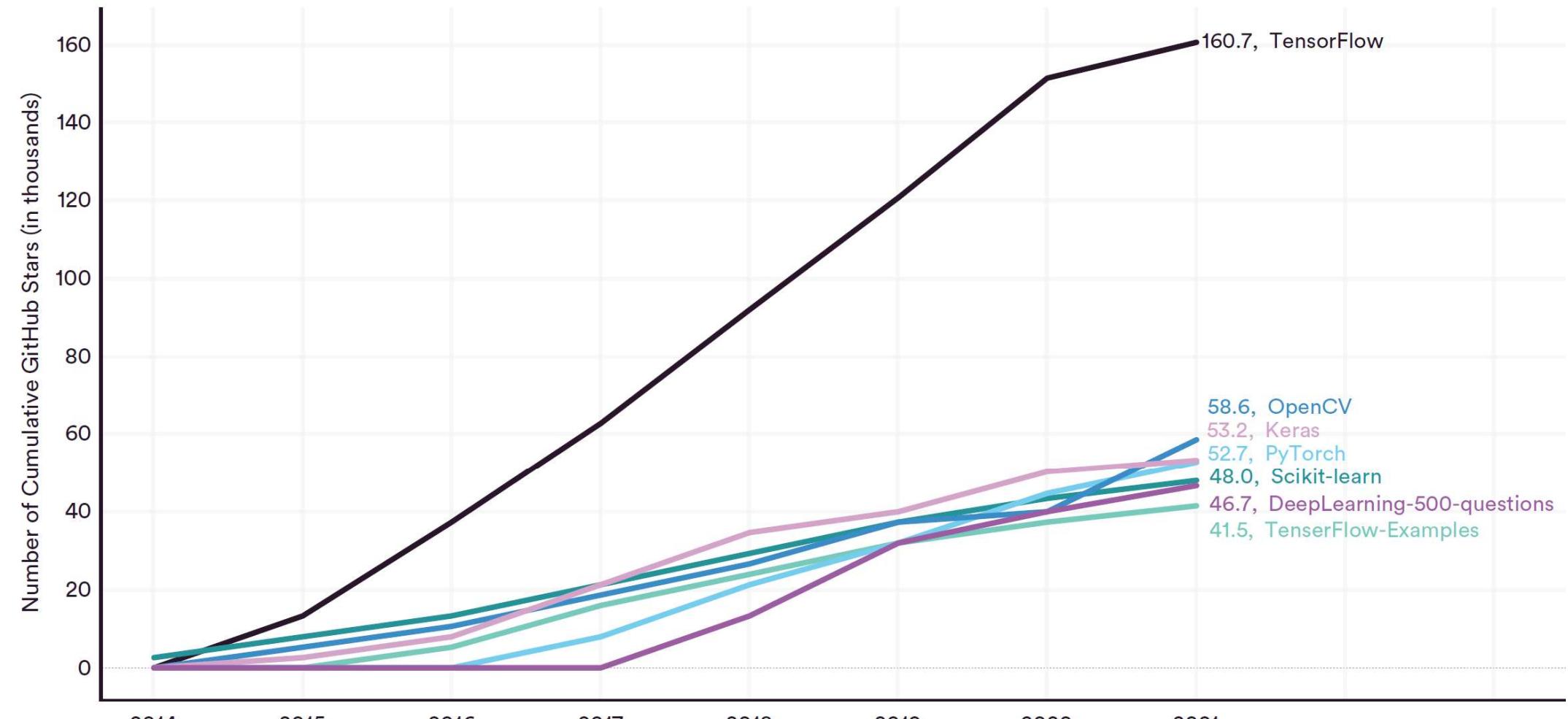
- ❖ Convolutional Neural Network (CNN)
- ❖ Recurrent Neural Network (RNN)
- ❖ Transformer
- ❖ Deep Learning Frameworks

Deep Learning Frameworks



NUMBER of GITHUB STARS by AI LIBRARY (OVER 40K STARS), 2014–21

Source: GitHub, 2021 | Chart: 2022 AI Index Report



from “The AI Index report 2022”

Q & A