



Generative AI



Generative AI

Generative AI is a branch of artificial intelligence primarily used for creative tasks

- Articles, images, music, and more
 - ChatGPT, Midjourney, Stable Diffusion, Office Copilot, and so on

AR 請幫我寫一封文情並茂、感人肺腑的情書

親愛的 (對方的名字),

我寫下這封情書, 希望用文字的力量將我對你的愛意傳達給你, 就像一顆溫暖的心跳, 傳送到你的內心深處。

每當我想起我們相識的那一刻, 我的心就像一朵綻放的花朵, 充滿了愛的香氣。那時, 我們彼此之間的目光交錯, 彷彿整個世界都靜止了, 只有我倆在那個時刻存在。從那刻起, 你就在我的心中佔據了一個特殊的位置。

每一次和你在一起的時光, 都是我生命中最珍貴的瞬間。當我們一起漫步在星光閃耀的夜晚, 我感受到了無盡的安慰和幸福。你的微笑是我每天的陽光, 溫暖著我的靈魂, 使我感到無比幸運。

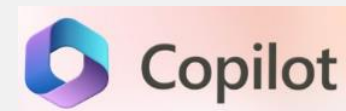
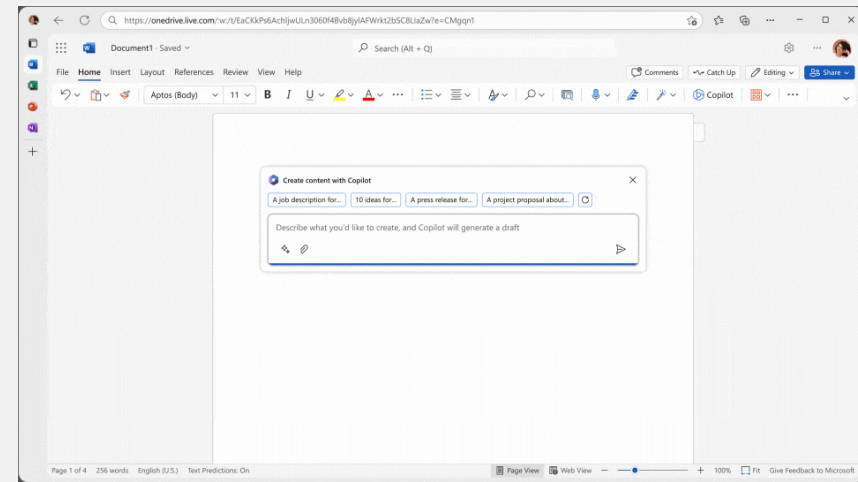
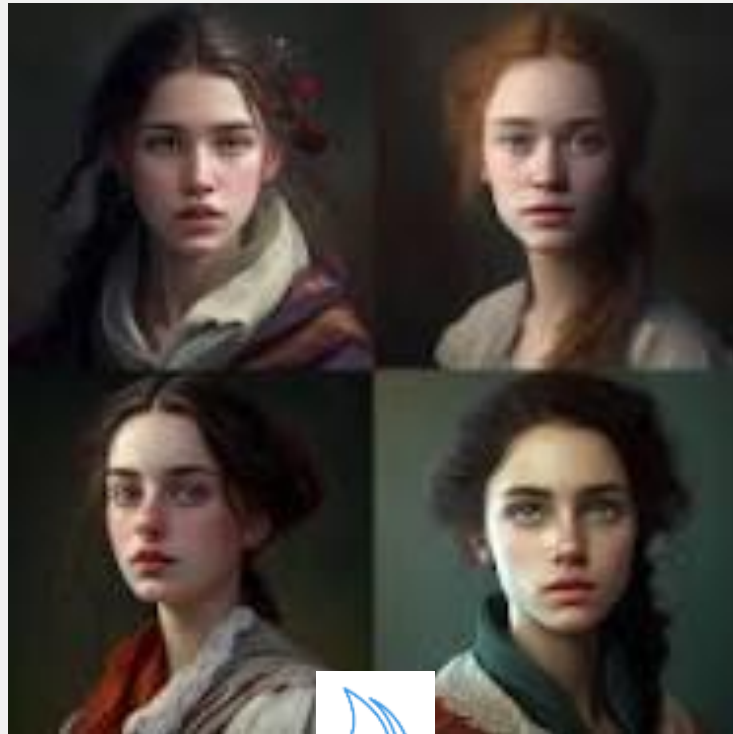
你的存在給了我力量和勇氣去面對生活中的挑戰。在我遇到困難和挫折時, 你總是在我身邊, 給我支持和鼓勵。你的鼓舞讓我相信, 無論遇到什麼困難, 我們都能一起克服它們。

每一次與你分享我的夢想和目標, 你總是無私地支持我, 並為我加油打氣。你的存在讓我相信, 我們可以攜手共同建立一個美好的未來, 實現彼此的夢想。

在我們相處的日子裡, 我深深地愛上了你的一切, 無論是你溫柔的眼神、甜蜜的笑聲還是你的善良和善解人意。你的單純和真誠感染了我, 使我願意把我的一切都奉獻給你。

親愛的, 我希望你能感受到我情書中真摯的情感, 這是我心中對你的深深愛意。無論未來會帶給我們什麼, 我都願意與你攜手共進, 直到永遠。

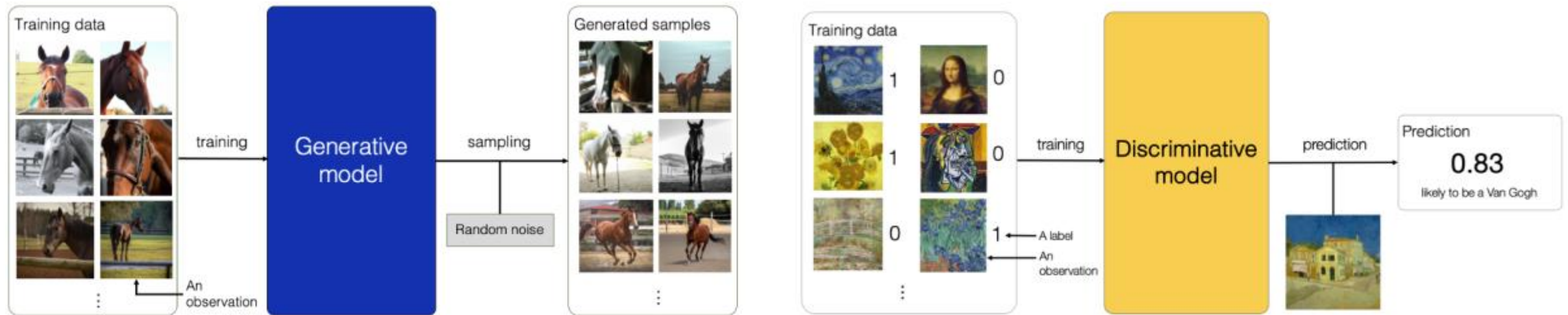
愛你的 (你的名字)



Generative AI

Differences between generative AI and Discriminative AI

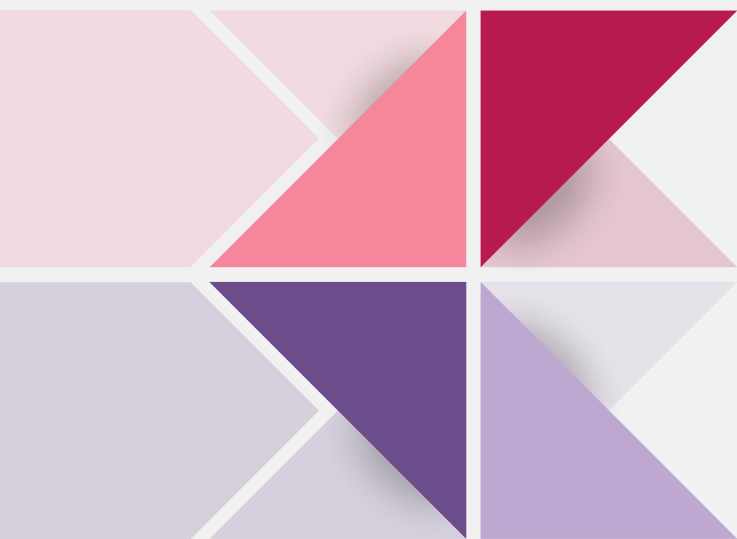
- Discriminative AI - **predict labels or classifications** based on training sets
 - **Predict labels or classifications** based on training sets
- Generative AI - **generate new data samples** based on training sets
 - **Generate new data samples** based on training sets





Generative AI in Language





01

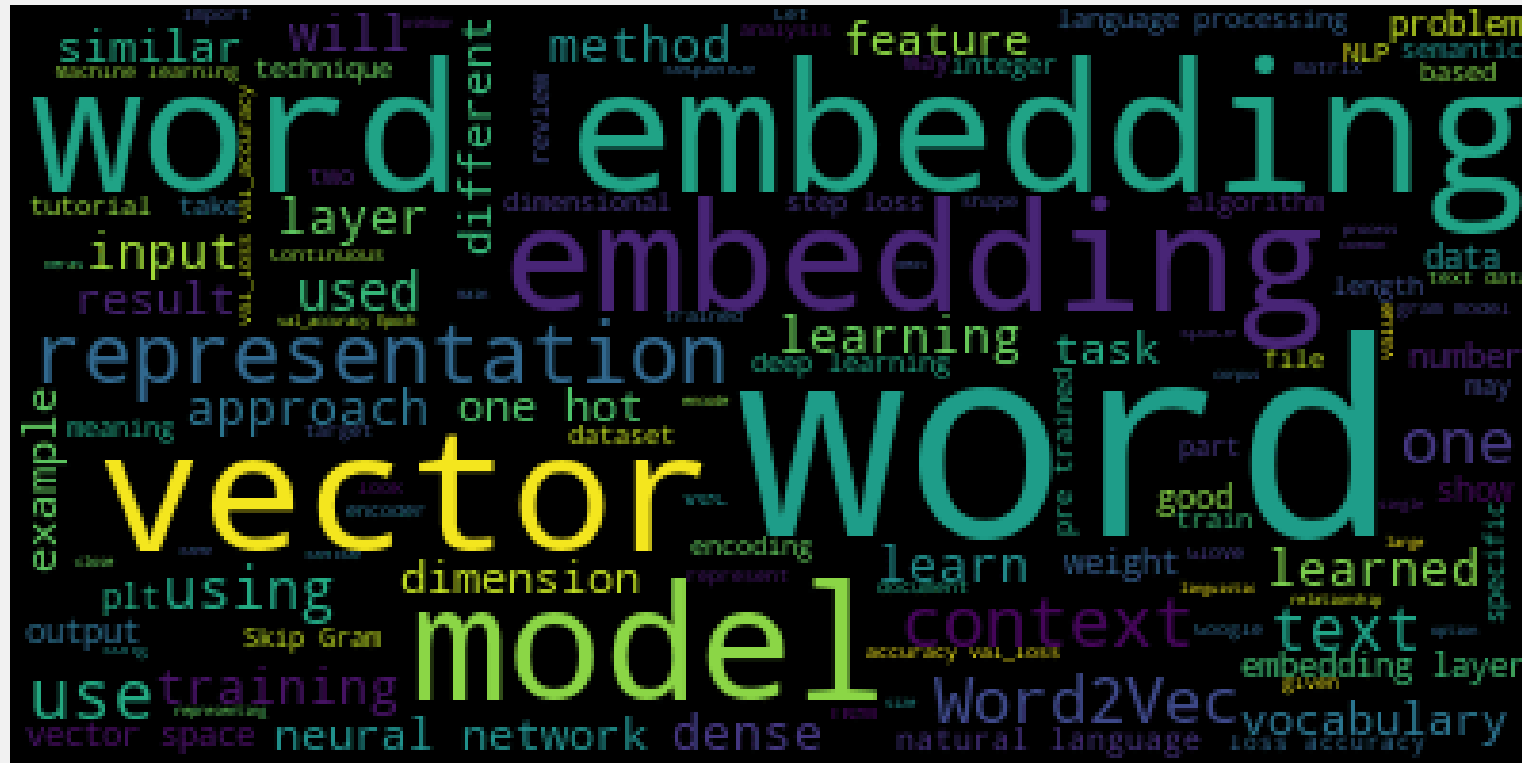
Word Embedding

Word Embedding

It is a method to embed "word" into "value"

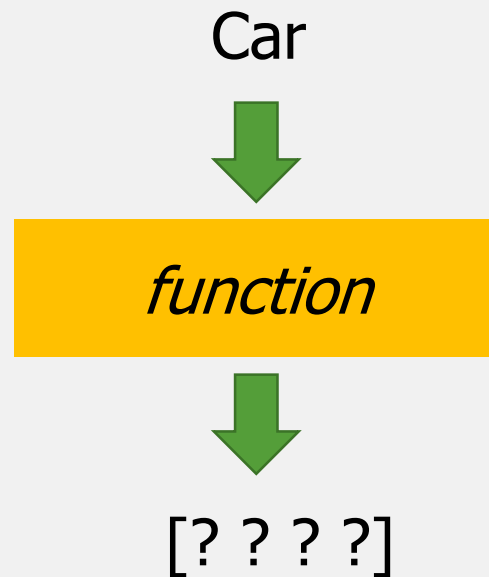
It is a dimension reduction technology

The value is "vector"



Word Embedding

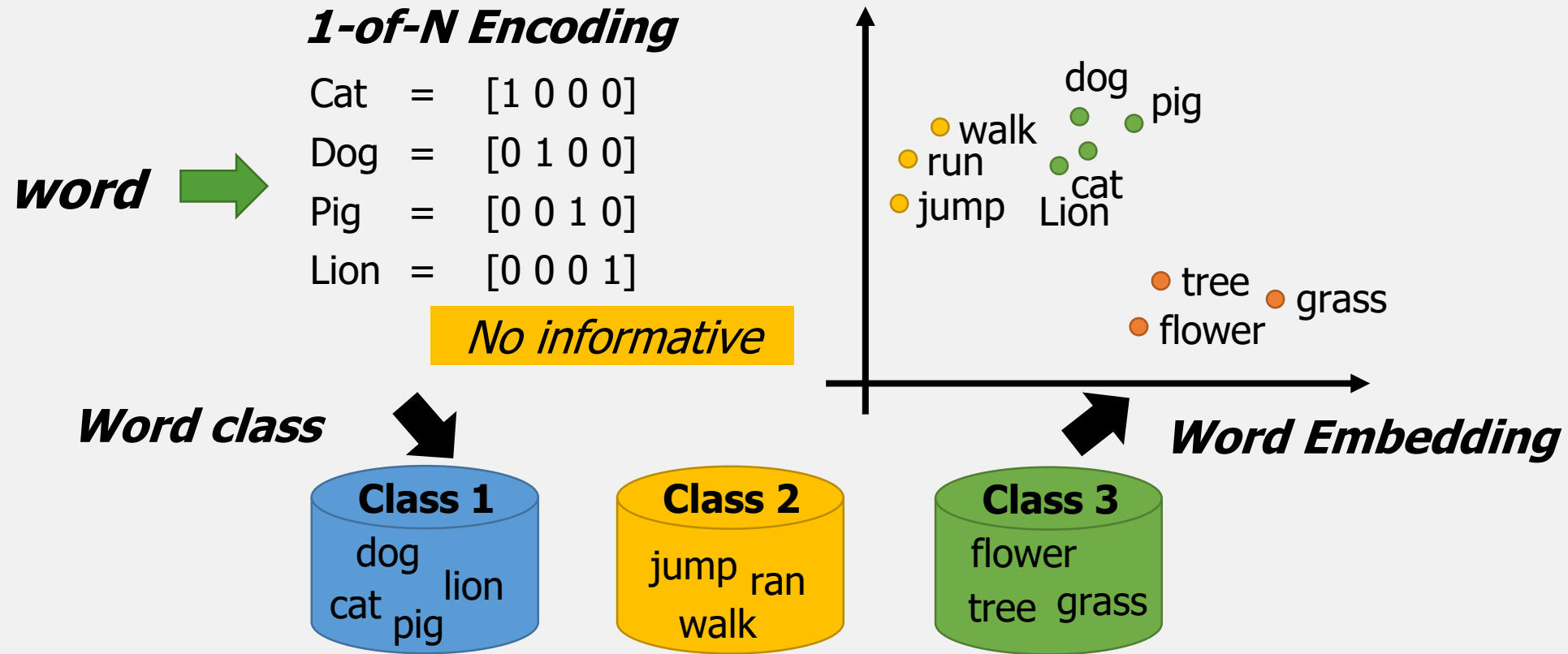
Word vector is generated with unsupervised learning



Word Embedding

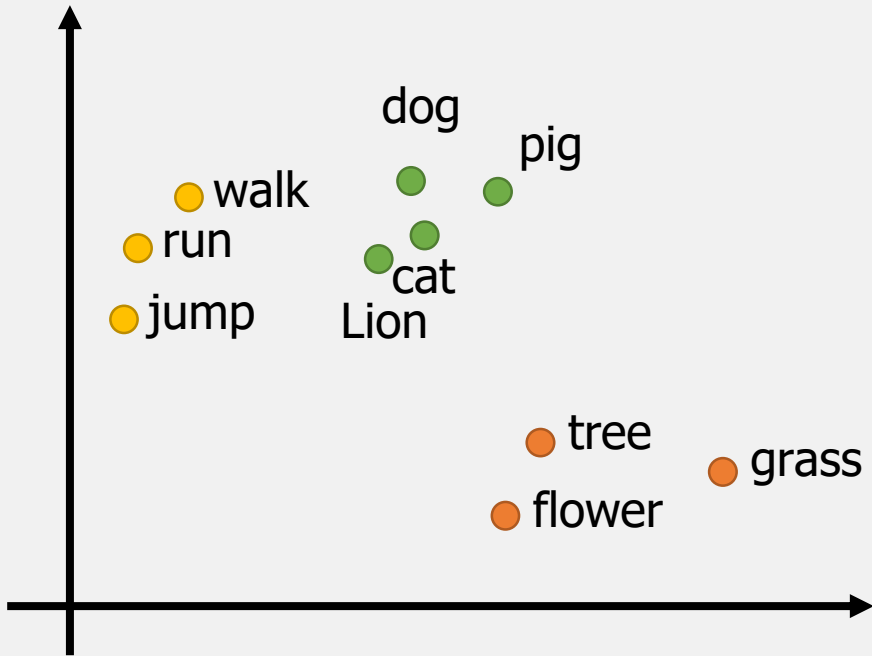
One-hot Encoding

Word vector is generated with unsupervised learning

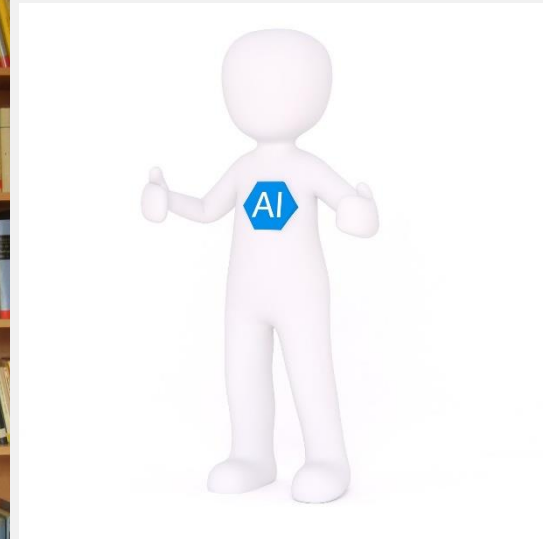


Word Embedding

Word Embedding



Documents



Word Embedding

How to do it?

➤ Counting-based

- if two words w_i and w_j often co-occur, the embedding vector v_i and v_j would be close to each other
- i.e. the larger inner product, the closer to each other



- The popular word embedding technology: Glove vector

Word Embedding

How to do it?

➤ Glove vector

- It uses the overall statistics and context to complete the word embedding

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

Word Embedding

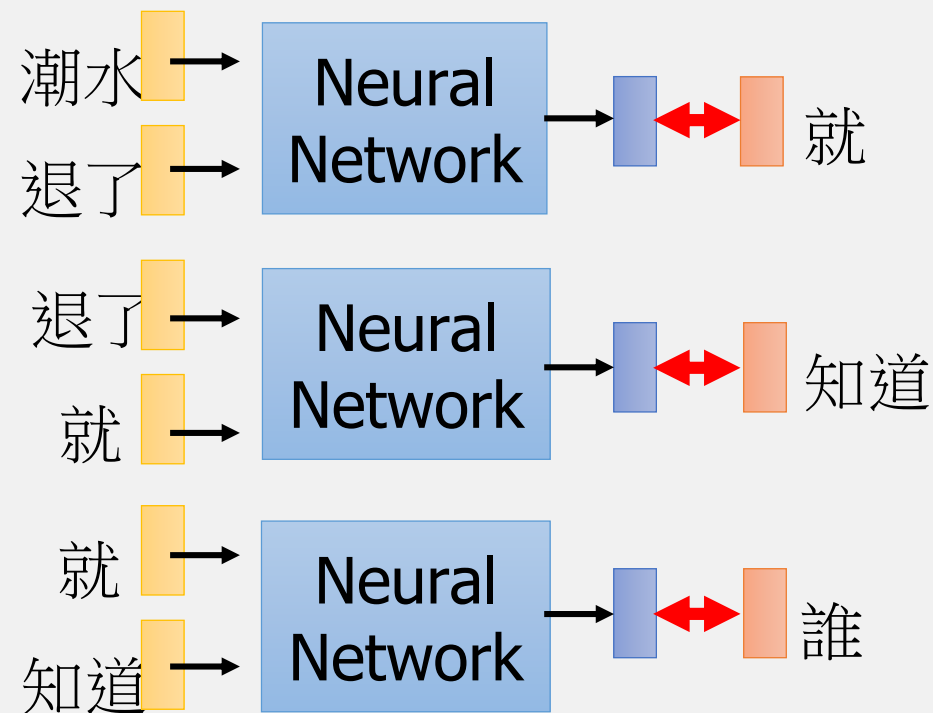
How to do it?

➤ Prediction-based

Pttr data:

道歉 時 需要 ...
本來 是 想說 點 什麼的 ...
不然 你要 投 ...
潮水 退了 就 知道 誰 ...
不爽 不要 買 ...
公道價 八萬 一 ...
.....

**Minimizing
cross entropy**

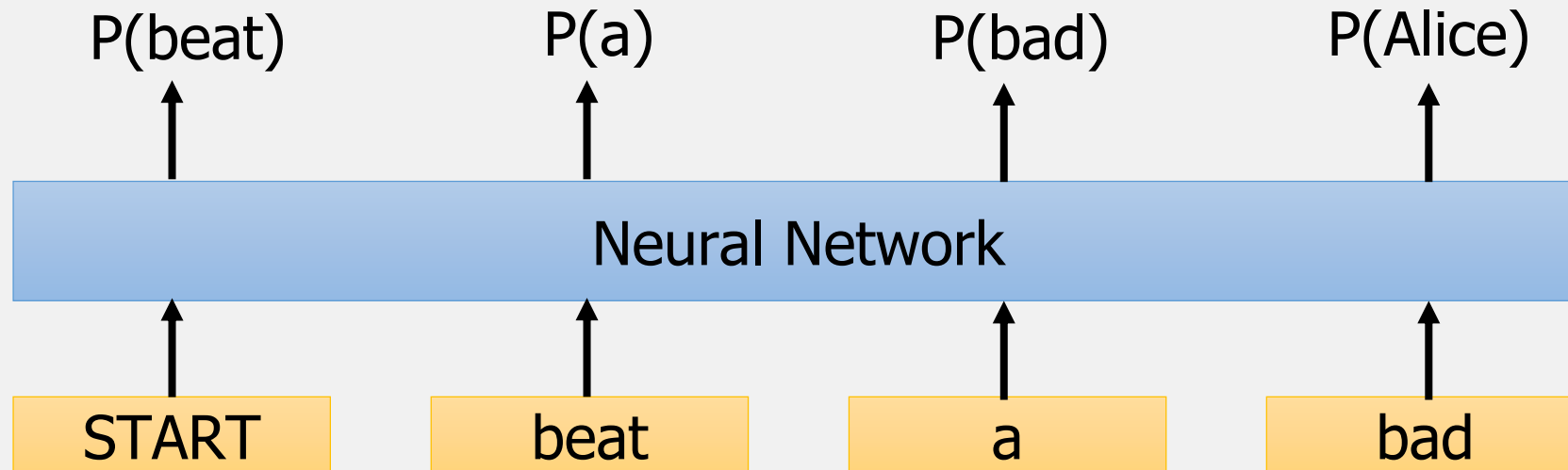


Word Embedding

Prediction-based - using a word x_i to predict the next word x_{i+1}

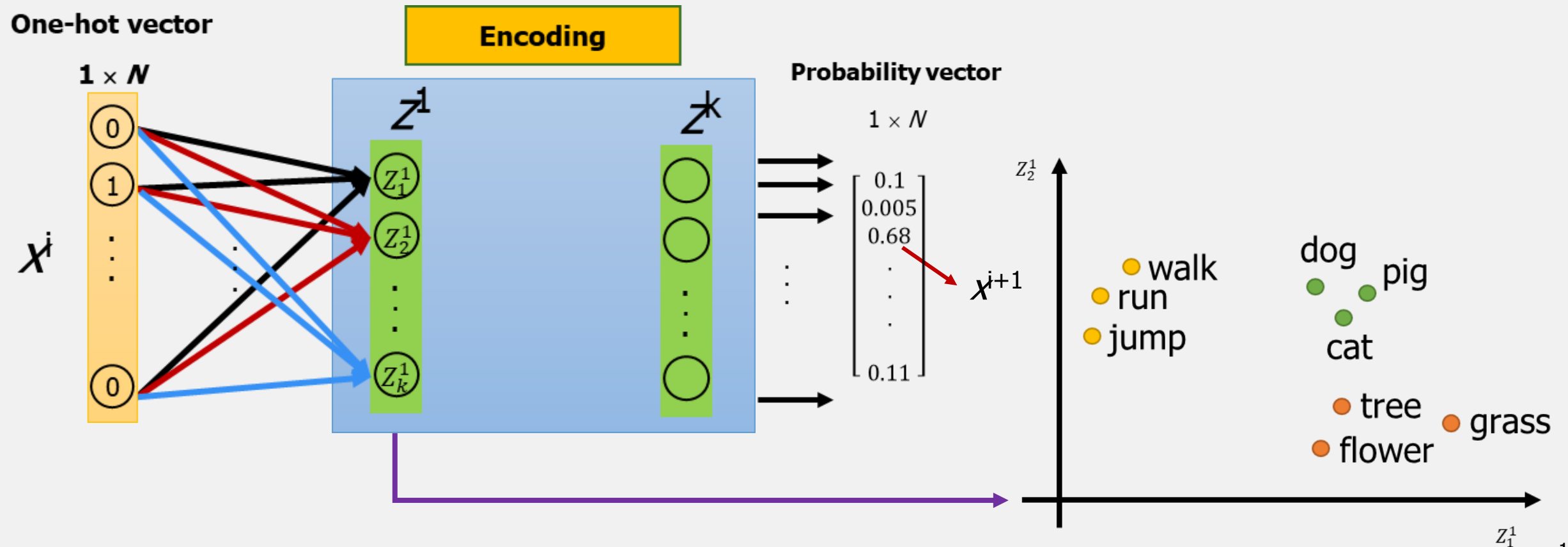
➤ $P(\text{"beat a bad Alice"}) = P(\text{beat}|\text{START})P(\text{a}|\text{beat})P(\text{bad}|\text{a})P(\text{Alice}|\text{bad})$

$P(b|a)$: the probability of model to predict the next word



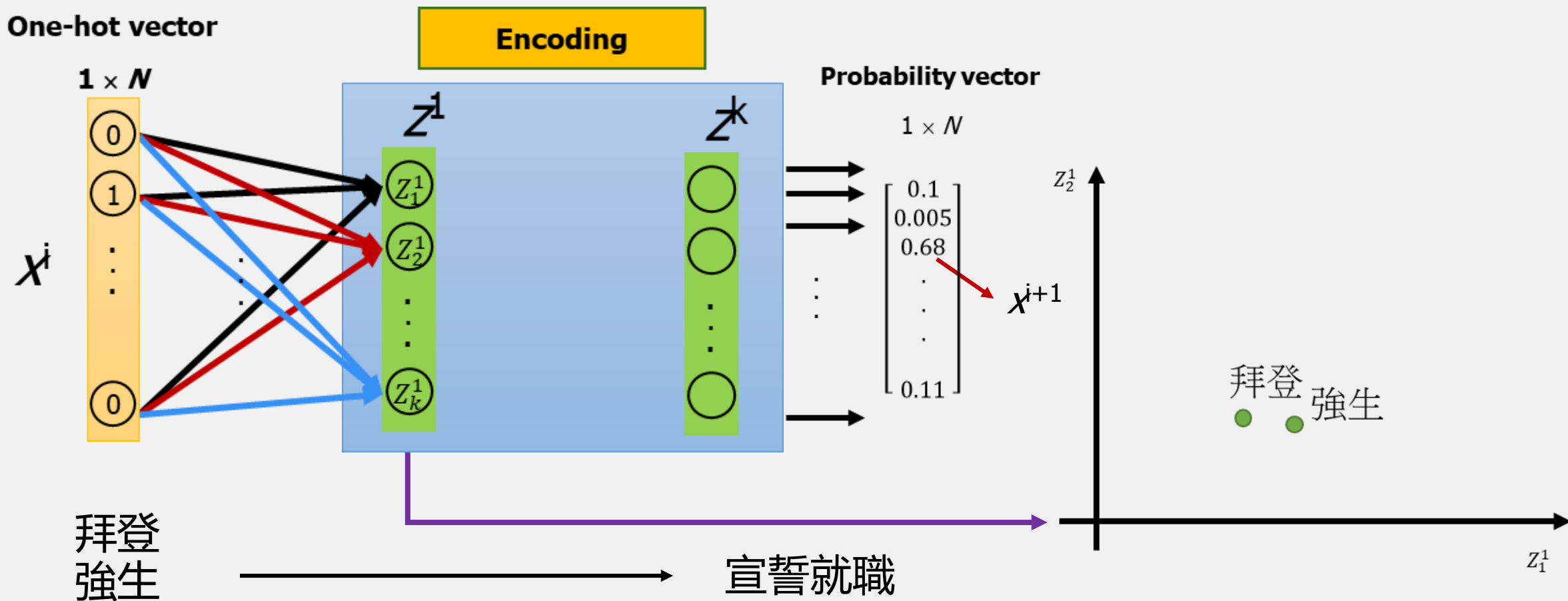
Word Embedding

Prediction-based



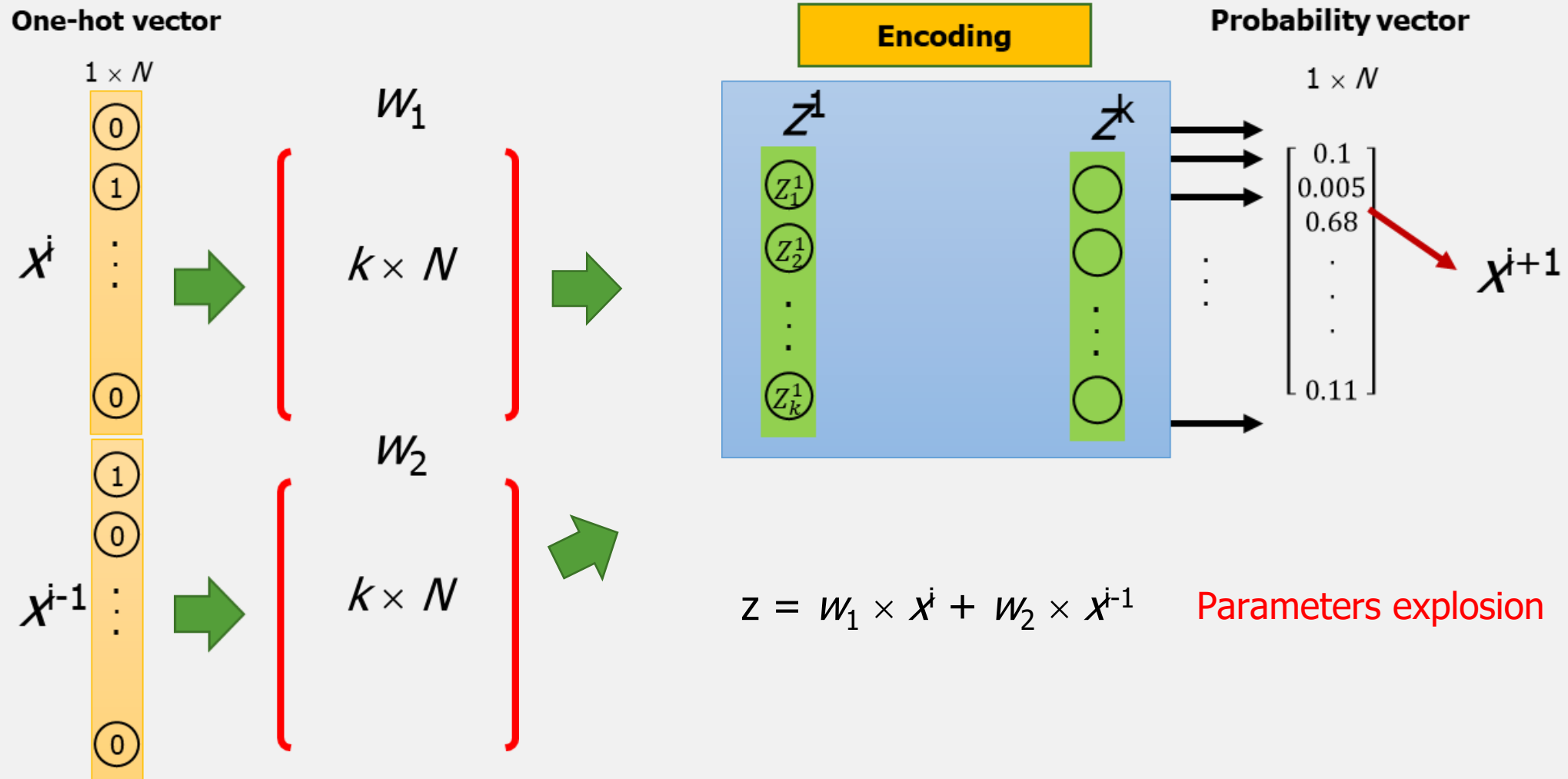
Word Embedding

Prediction-based



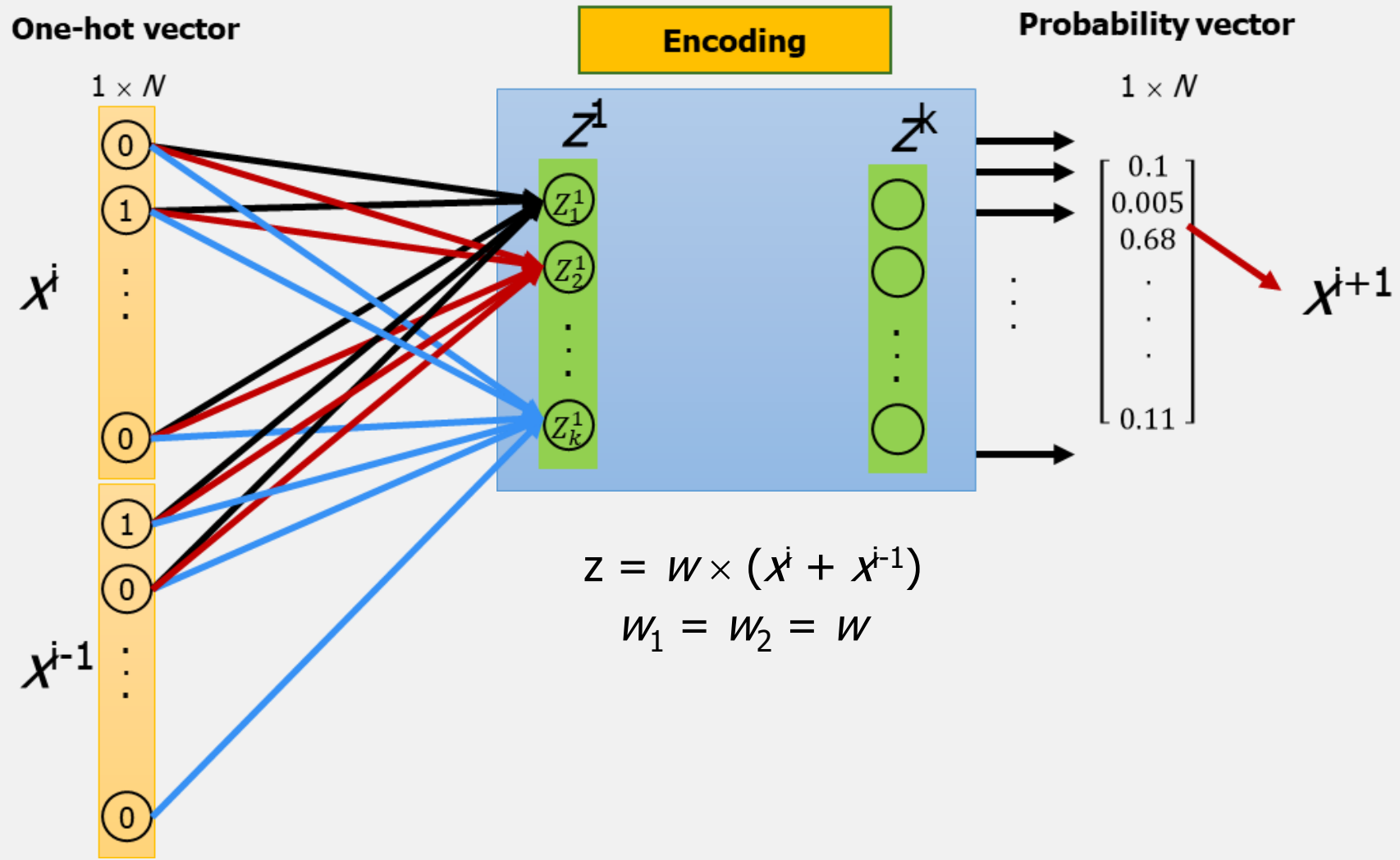
Word Embedding

Prediction-based - More and more words



Word Embedding

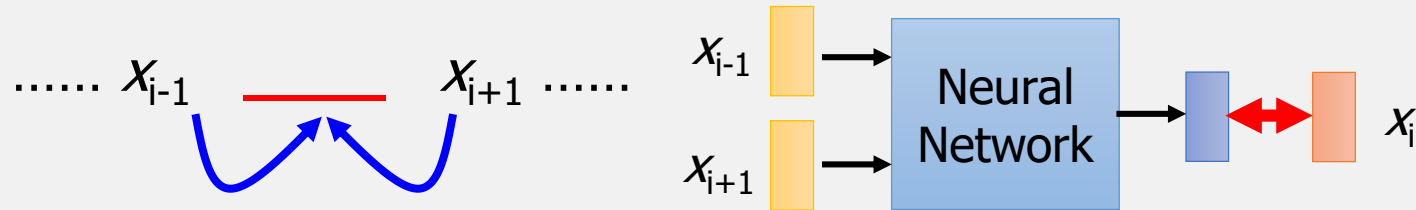
Prediction-based - More and more words



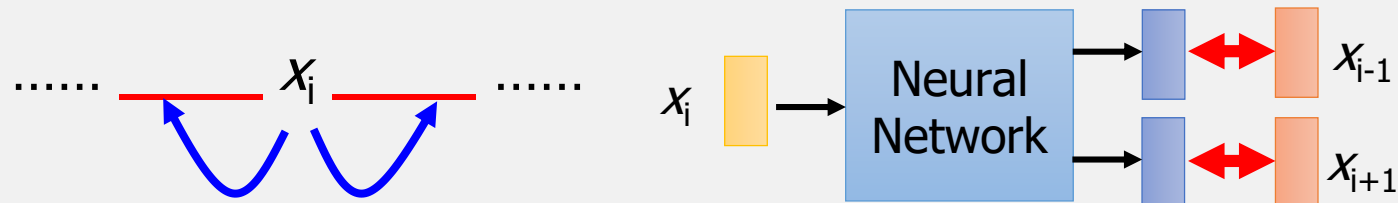
Word Embedding

Other methods

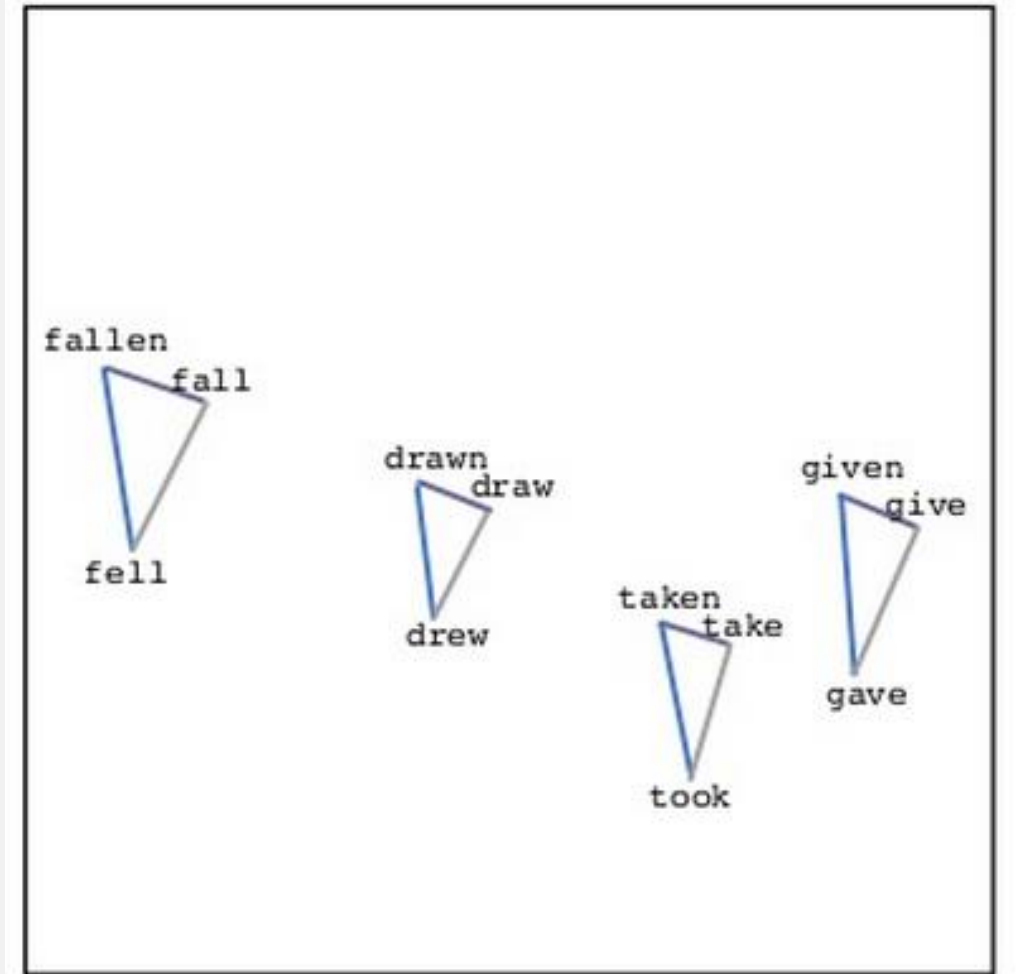
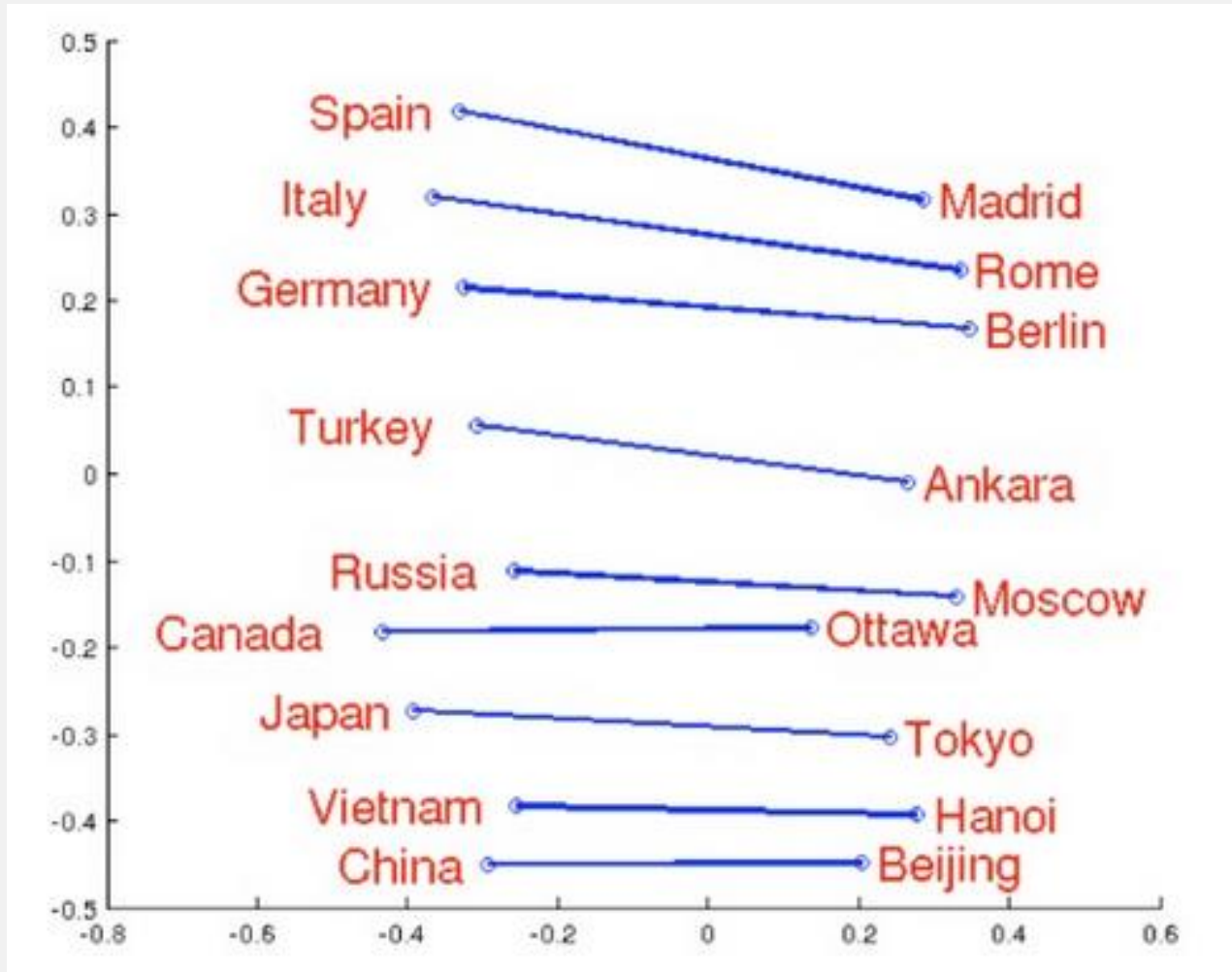
- Continuous bag of word (CBOW)
 - predict a word from given context



- Skip-gram
 - predict a context from given a word



Word Embedding



[Reference](#)

Document Embedding

Could the document embedding be achieved too?

- word -> word sequences with different lengths -> the vector with the same length
- The vector represents the meaning of the word sequence
- A word sequence can be a document or a paragraph



Document Embedding

Bag-of-word (BOW)

- Process a sentence or a document as a bag of words
- Each document is converted to a <word, count> map
- Document similarity
 - Euclidean distance
 - Cosine
 - Dot-product
 - ...

Bag of words (BoW)

Very good drama although it appeared to have a few blank areas leaving the viewers to fill in the action for themselves. I can imagine life being this way for someone who can neither read nor write. This film simply smacked of the real world: the wife who is suddenly the sole supporter, the live-in relatives and their quarrels, the troubled child who gets knocked up and then, typically, drops out of school, a jackass husband who takes the nest egg and buys beer with it. 2 thumbs up... very very very good movie.

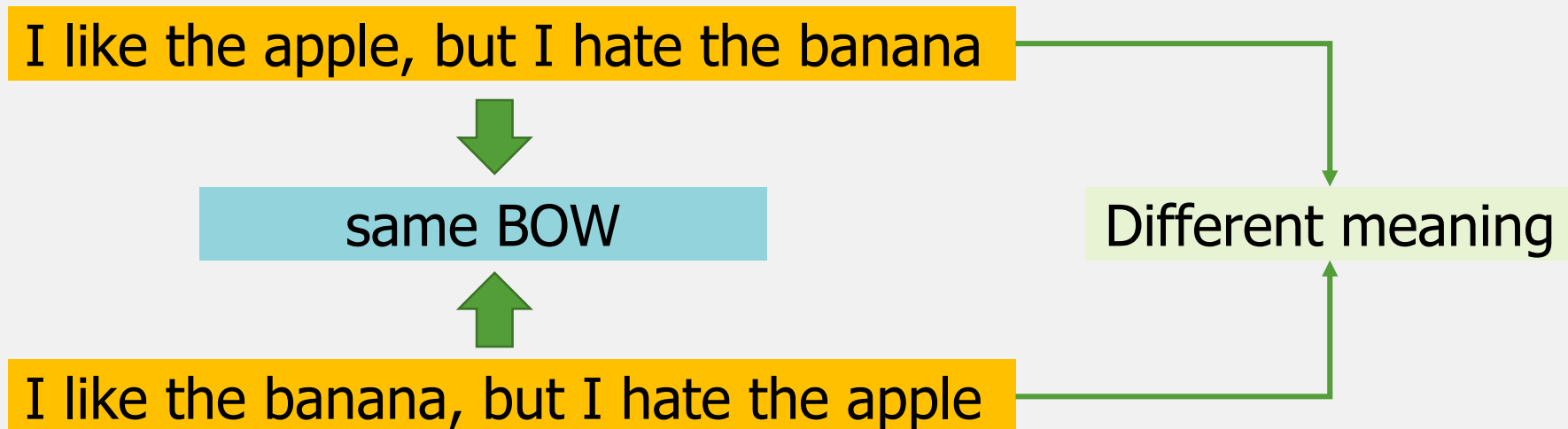


('the', 8),
(',', 5),
('very', 4),
('.', 4),
('who', 4),
('and', 3),
('good', 2),
('it', 2),
('to', 2),
('a', 2),
('for', 2),
('can', 2),
('this', 2),
('of', 2),
('drama', 1),
('although', 1),
('appeared', 1),
('have', 1),
('few', 1),
('blank', 1)
.....

Document Embedding

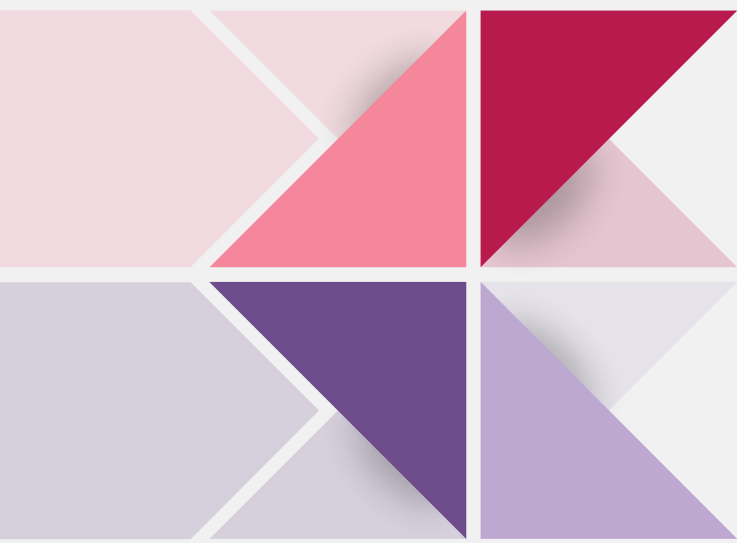
The problem of BOW

- the information of order of the words is ignored



Word Embedding

Example

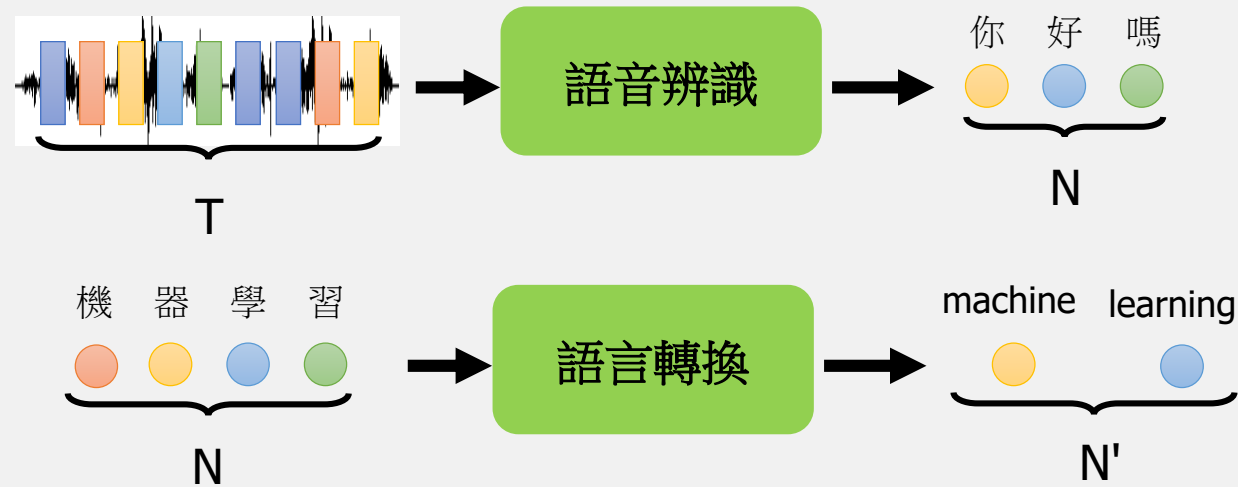


02

Seq2Seq

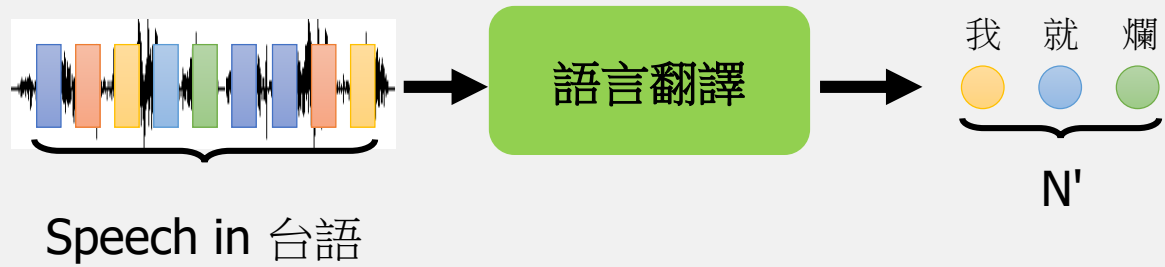
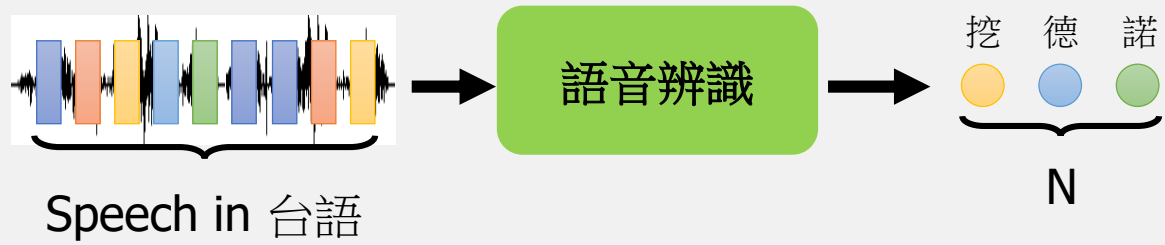
Seq2Seq

Input a sequence, output a sequence whose length is determined by model



Seq2Seq

In 台語



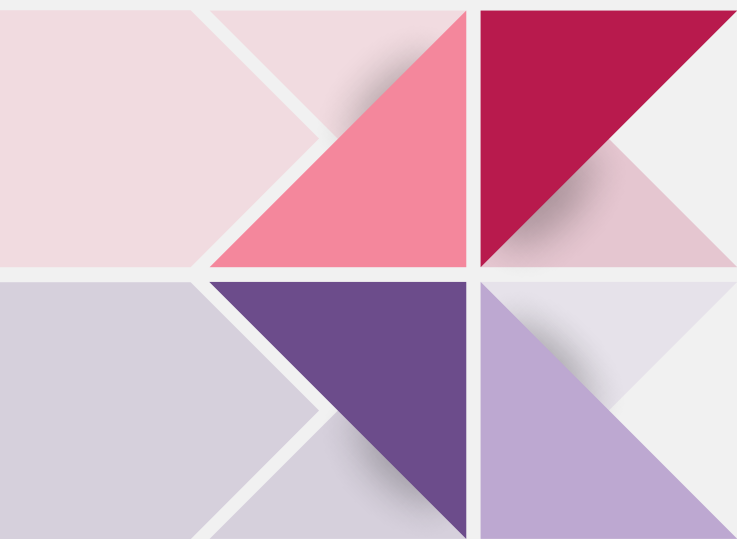
Seq2Seq

Chatbot



Training data:

```
[PERSON 1:] Hi  
[PERSON 2:] Hello ! How are you today ?  
[PERSON 1:] I am good thank you , how are you.  
[PERSON 2:] Great, thanks ! My children and I were just about to watch Game of Thrones.  
[PERSON 1:] Nice ! How old are your children?  
[PERSON 2:] I have four that range in age from 10 to 21. You?  
[PERSON 1:] I do not have children at the moment.  
[PERSON 2:] That just means you get to keep all the popcorn for yourself.  
[PERSON 1:] And Cheetos at the moment!  
[PERSON 2:] Good choice. Do you watch Game of Thrones?  
[PERSON 1:] No, I do not have much time for TV.  
[PERSON 2:] I usually spend my time painting: but, I love the show.
```

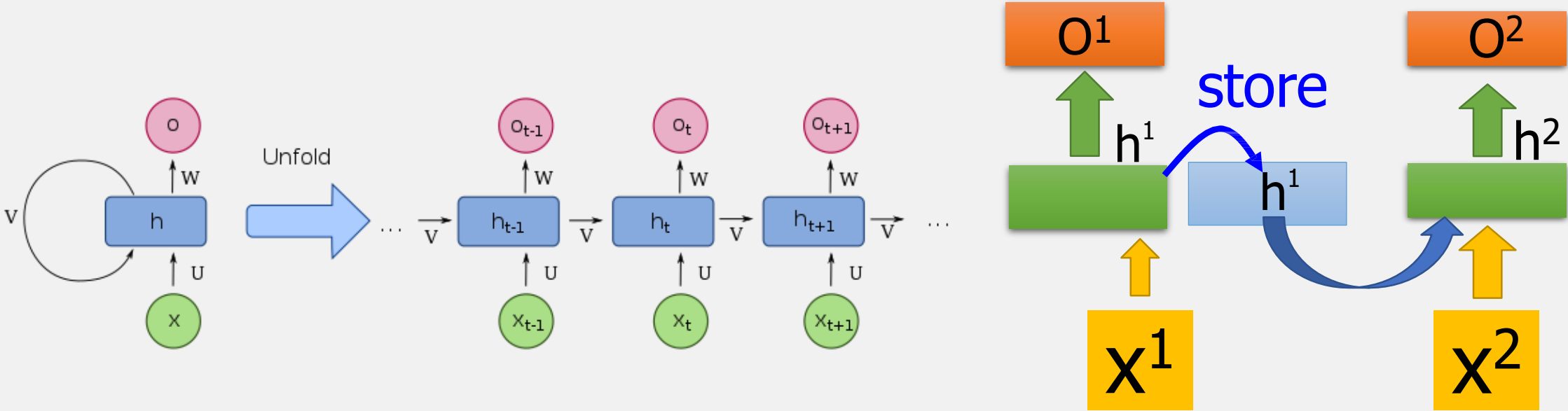


A

RNN and LSTM

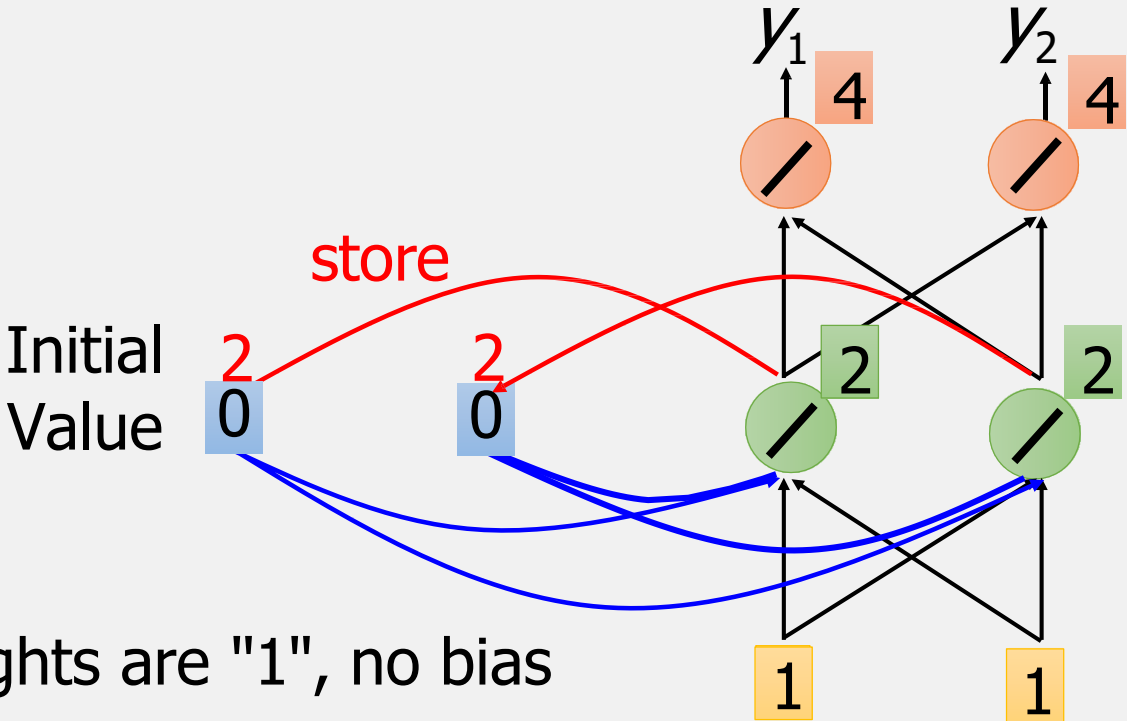
RNN and LSTM

RNN is designed for the time-series problem



RNN and LSTM

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots$
 output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \dots$

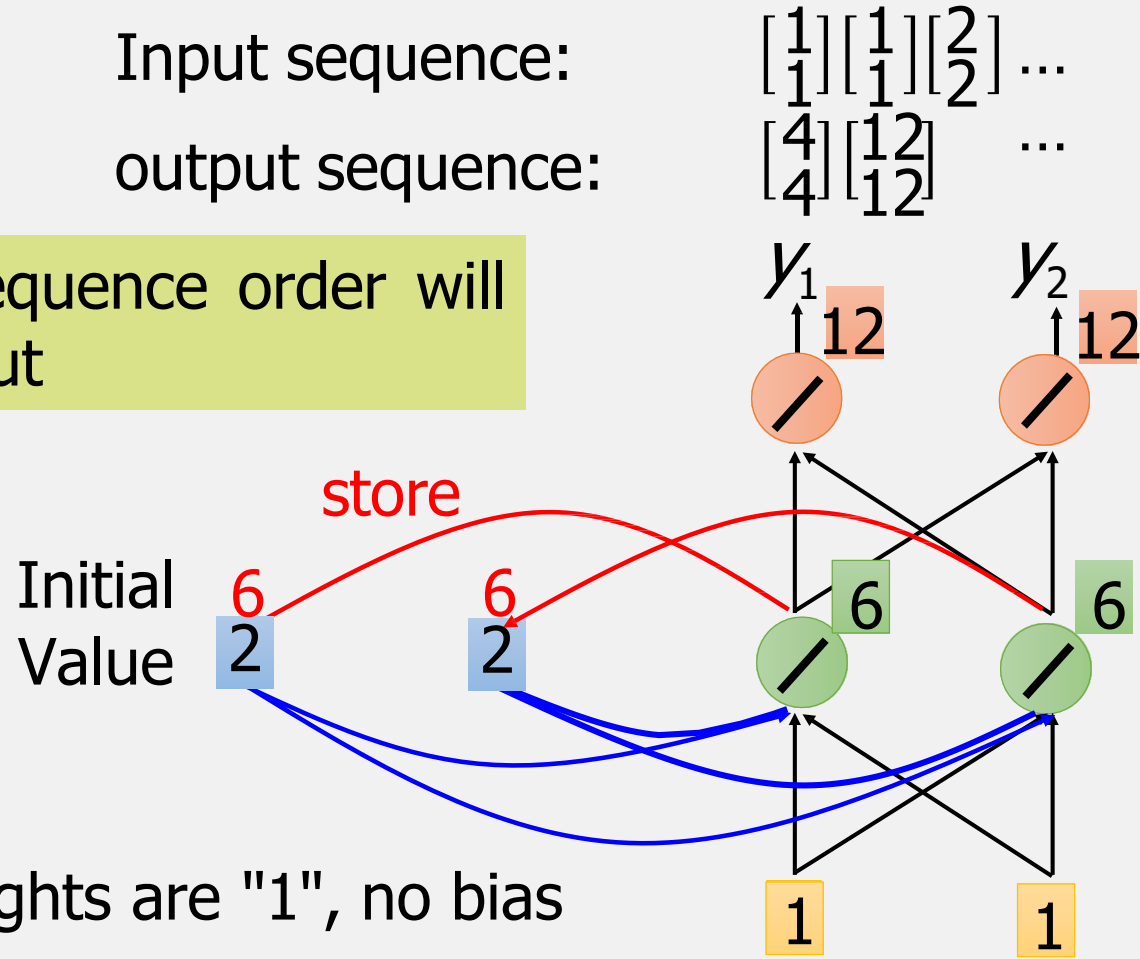


All the weights are "1", no bias
 All activation functions are linear

RNN and LSTM

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots$
output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix} \dots$

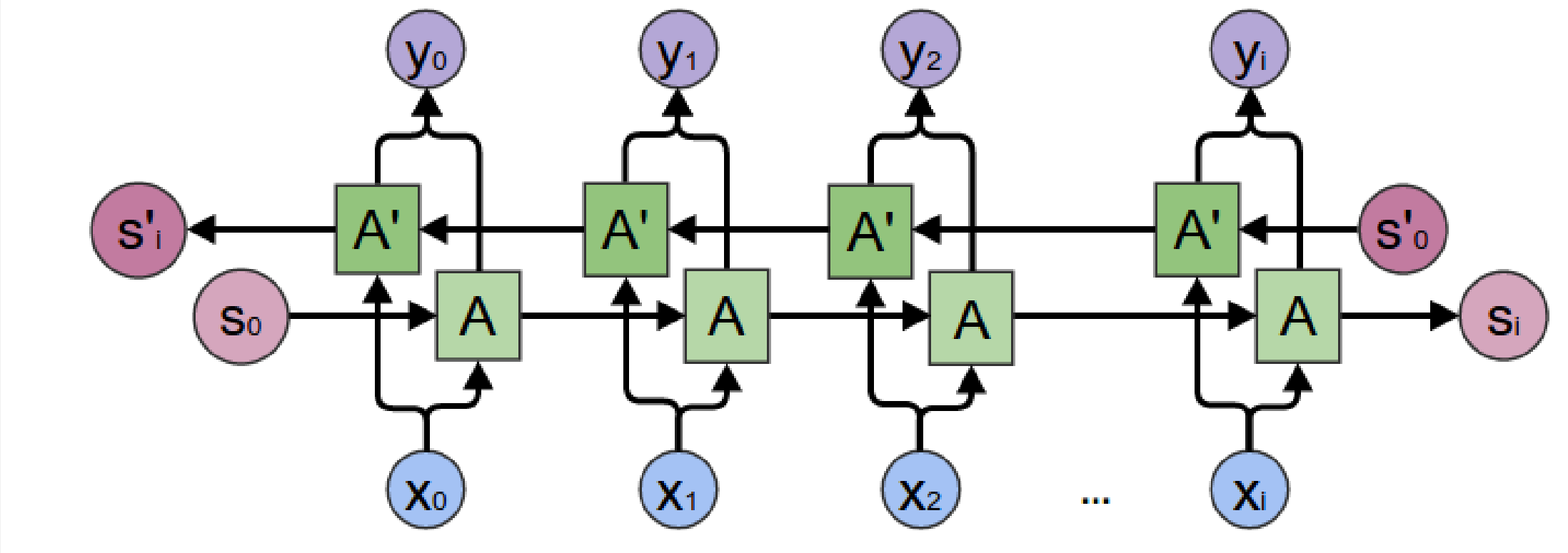
Changing the sequence order will change the output



All the weights are "1", no bias
All activation functions are linear

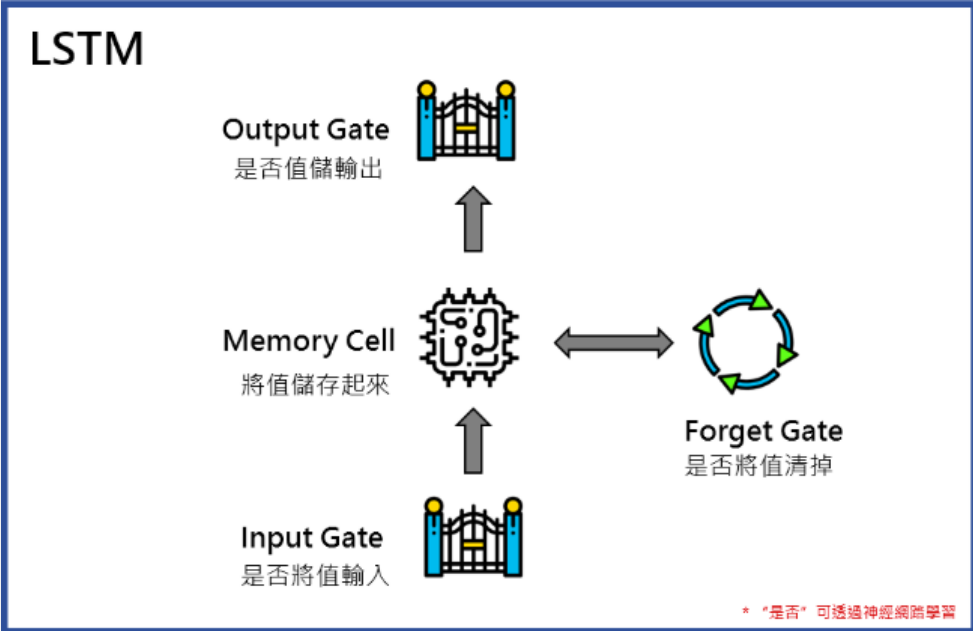
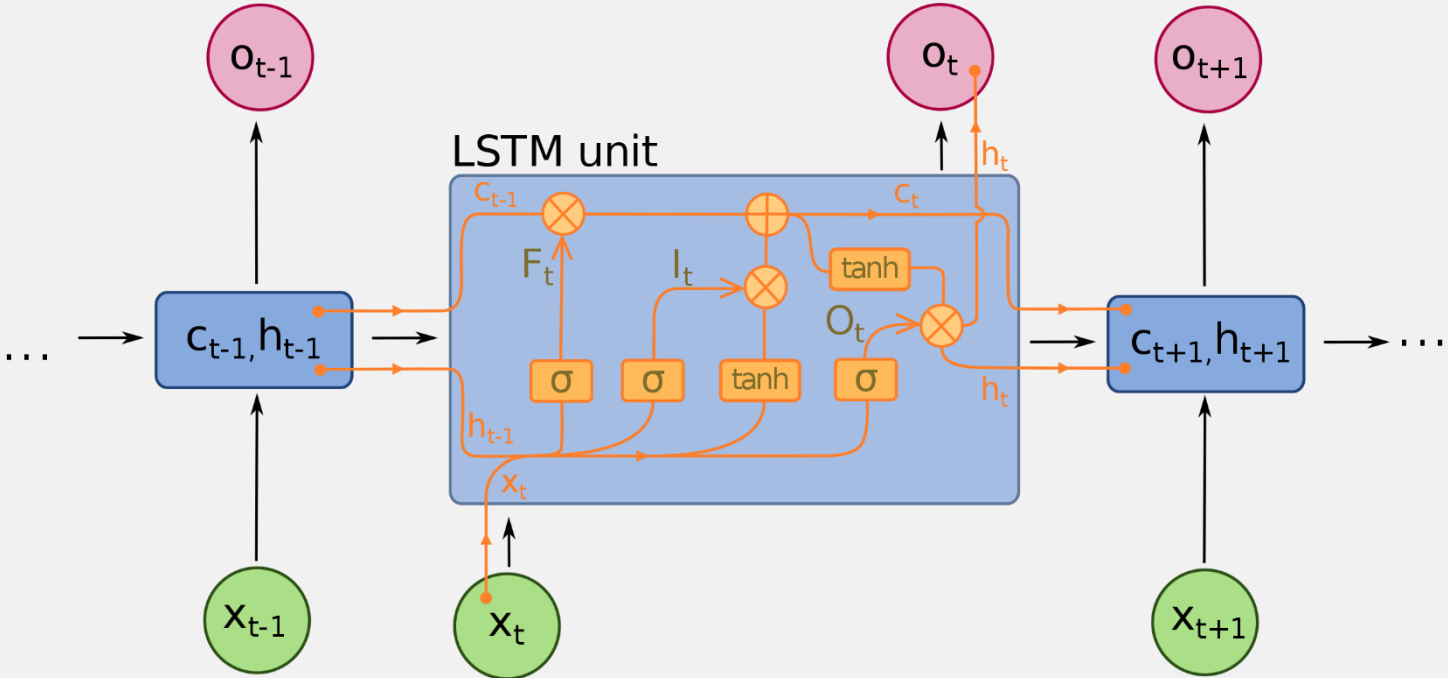
RNN and LSTM

Bi-directional RNN

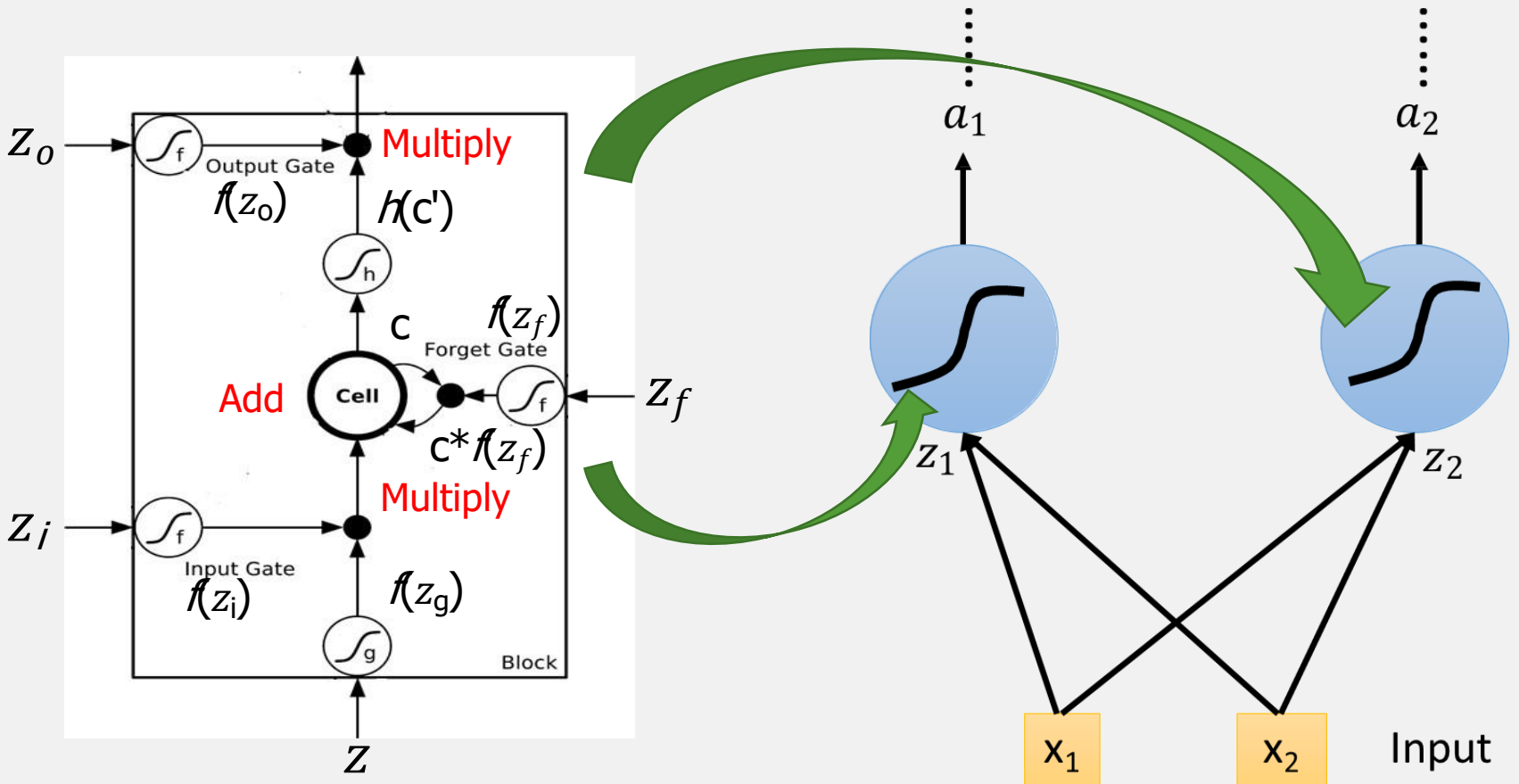


RNN and LSTM

LSTM is the best famous model in RNN



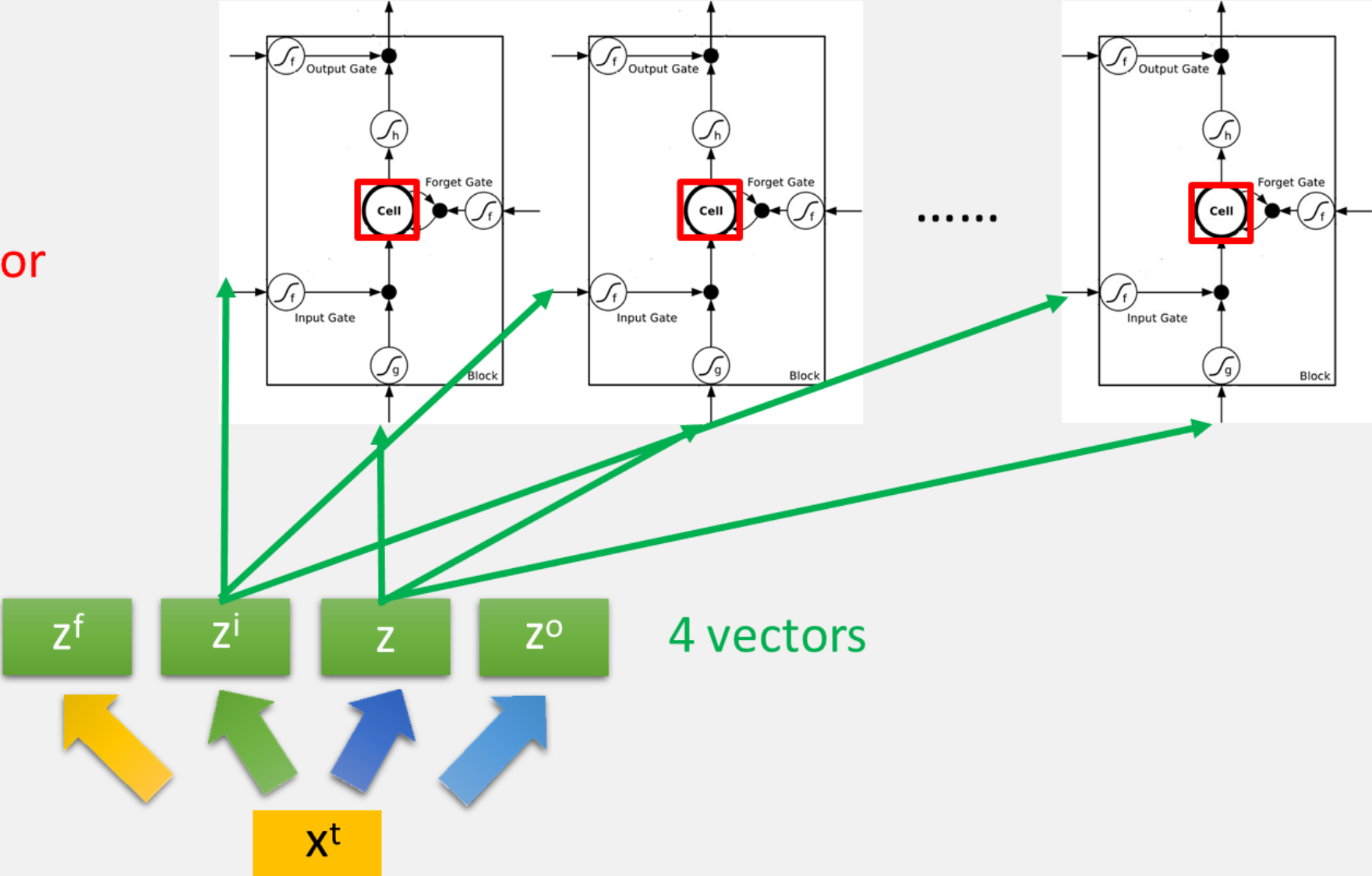
RNN and LSTM



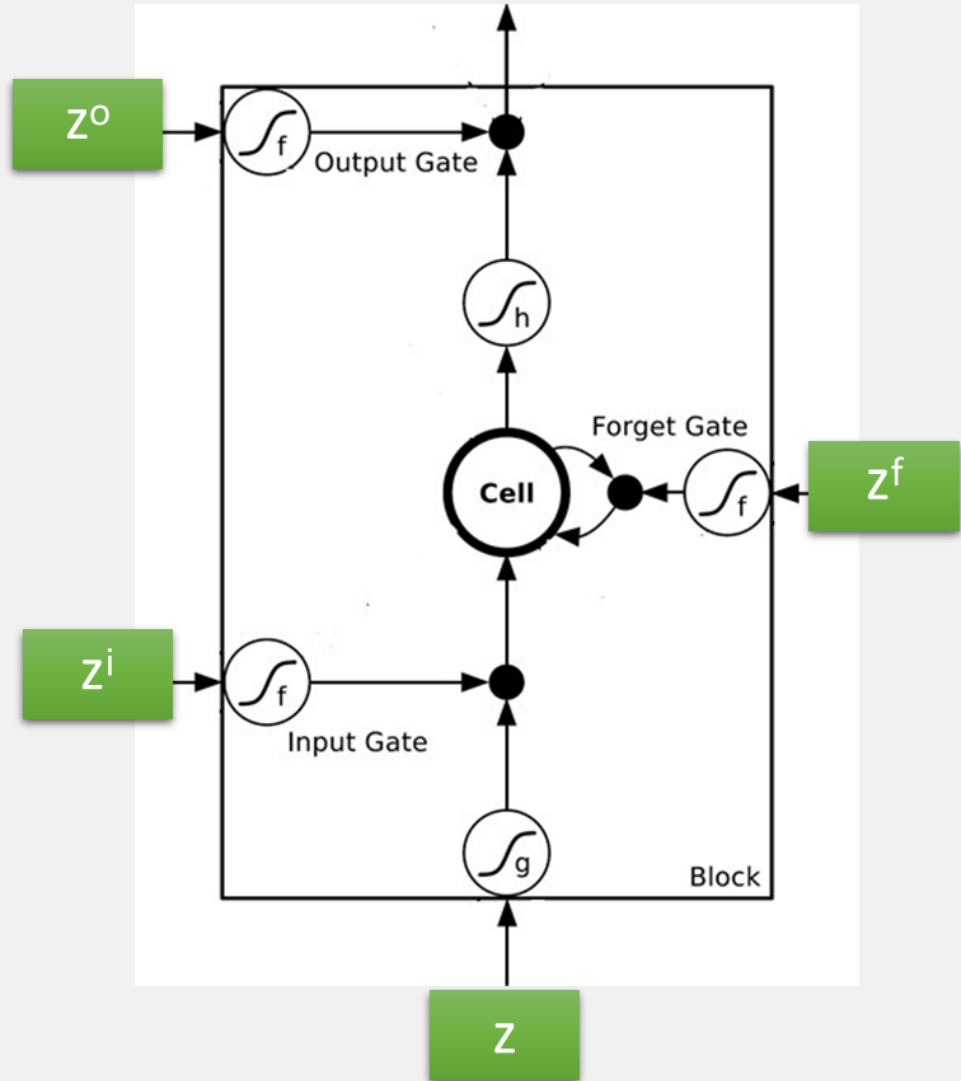
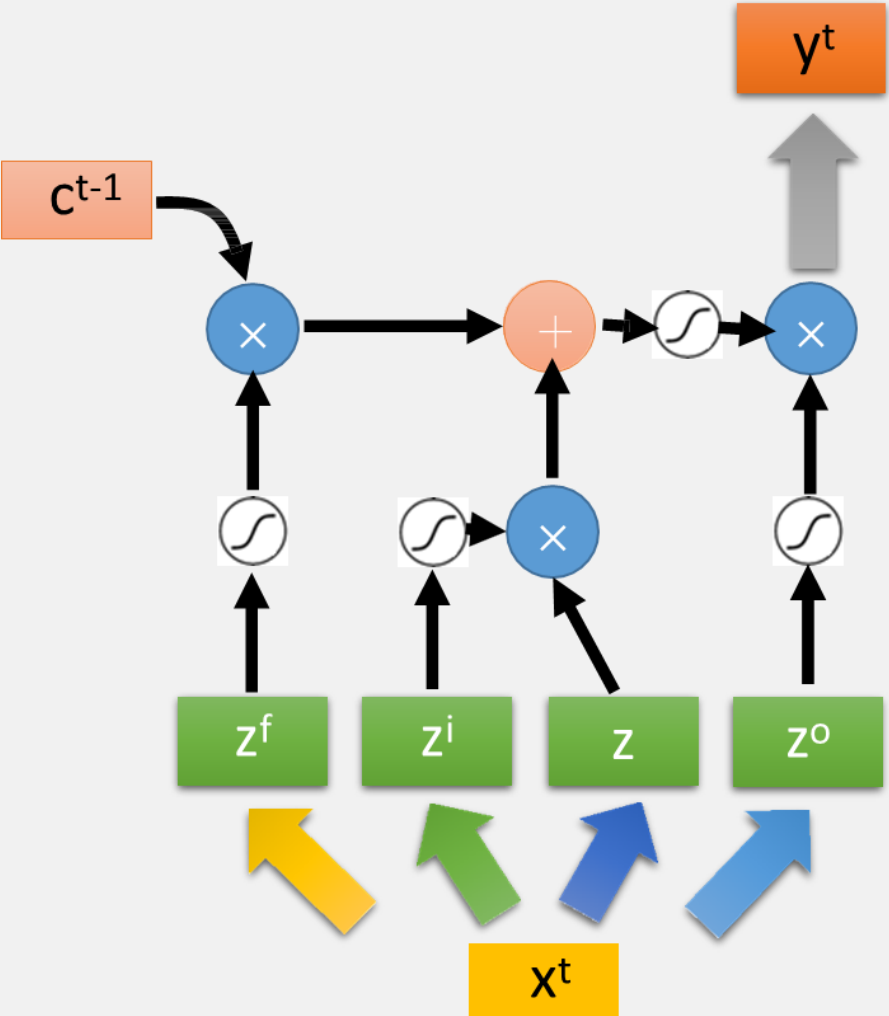
RNN and LSTM

c^{t-1}

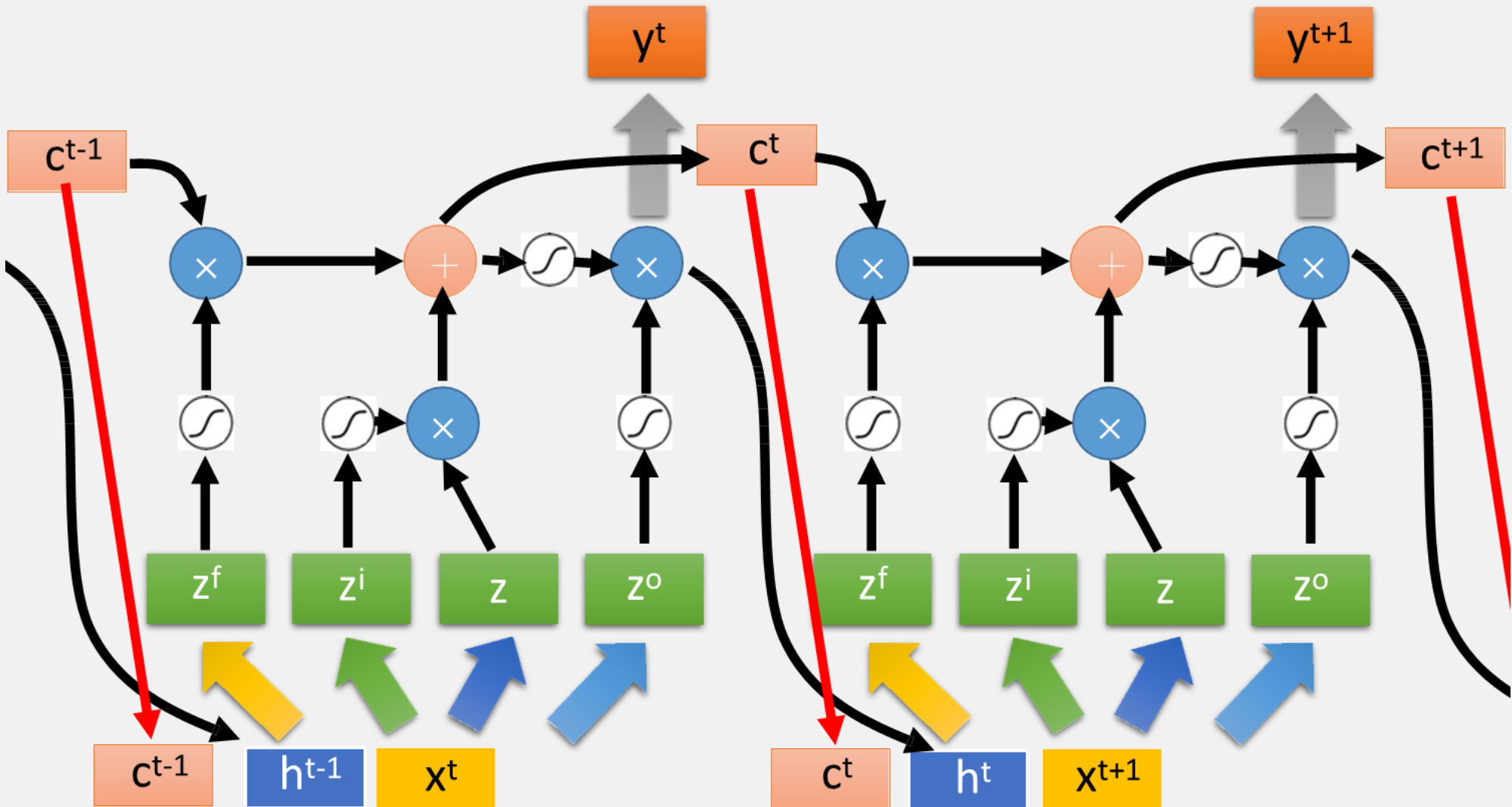
vector

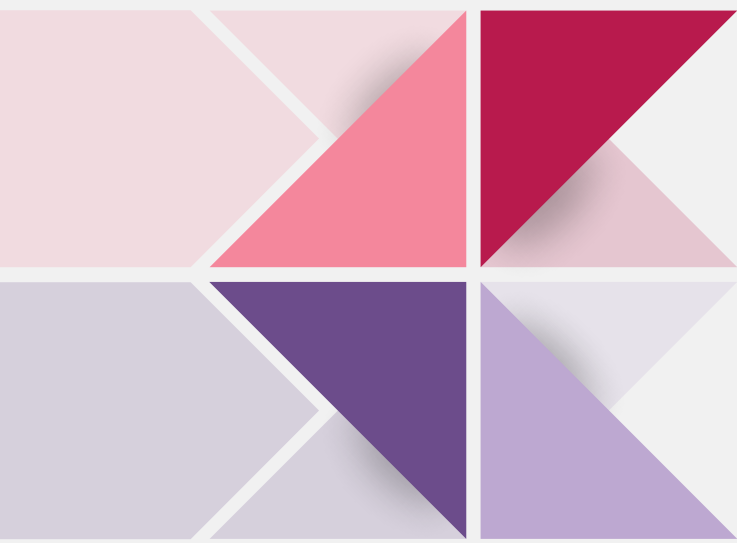


RNN and LSTM



RNN and LSTM





B

Seq2Seq

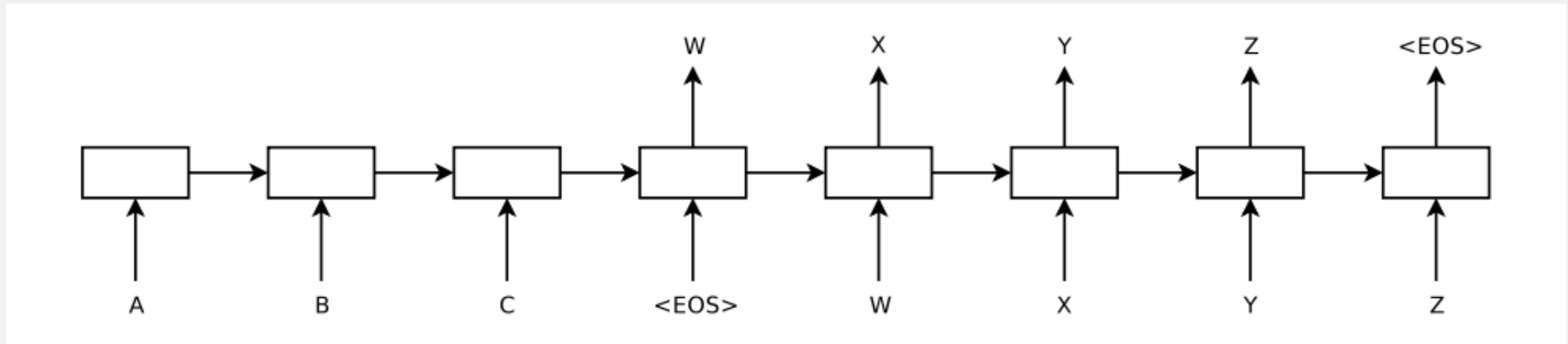
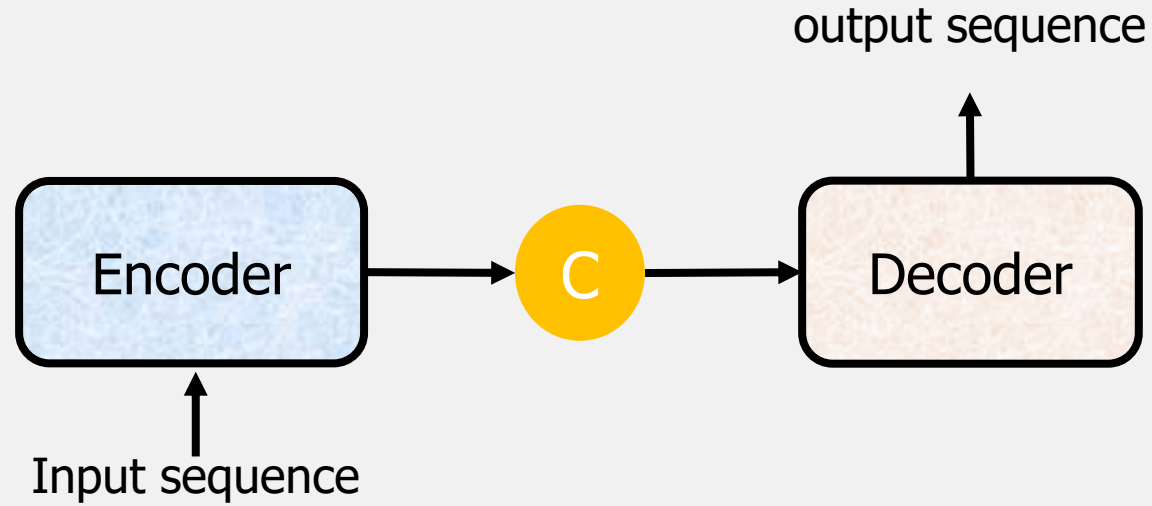
Seq2Seq

Seq2Seq is built by two RNN or LSTM models

The first model, encoder, is to convert the sequence with length M into a vector with a fixed length

The second model, decoder, is to reconstruct the vector back to a sequence with length N

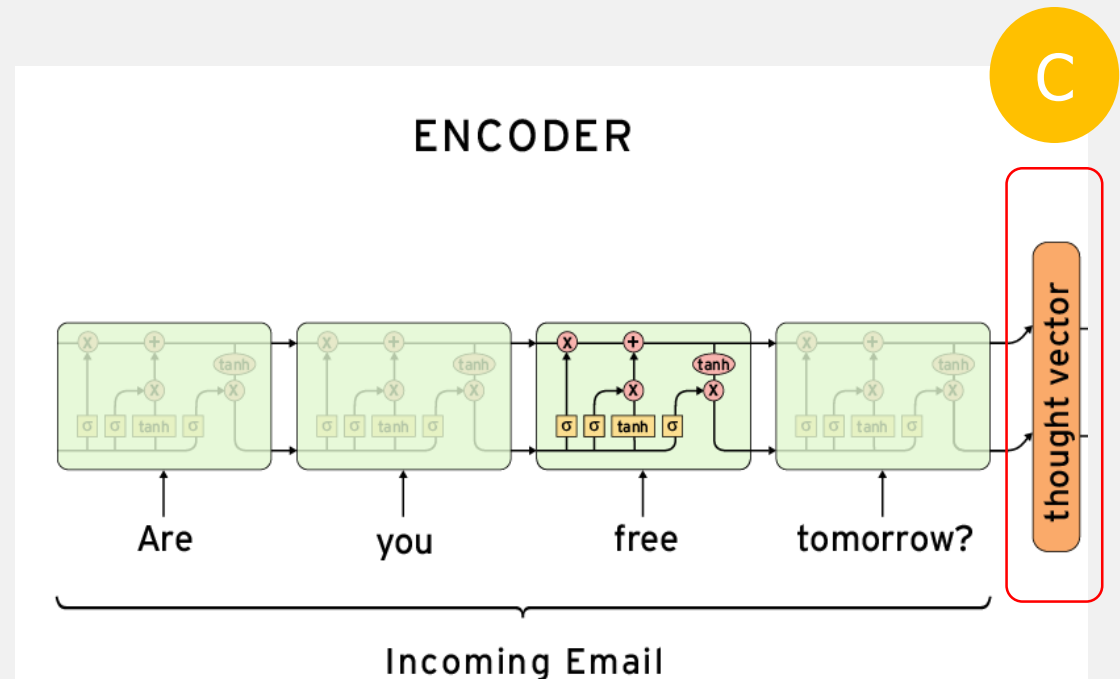
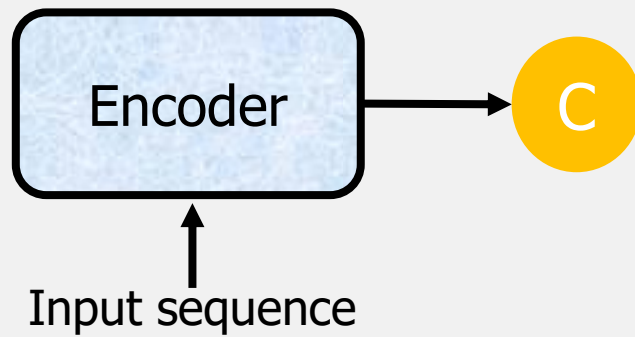
Seq2Seq



Seq2Seq

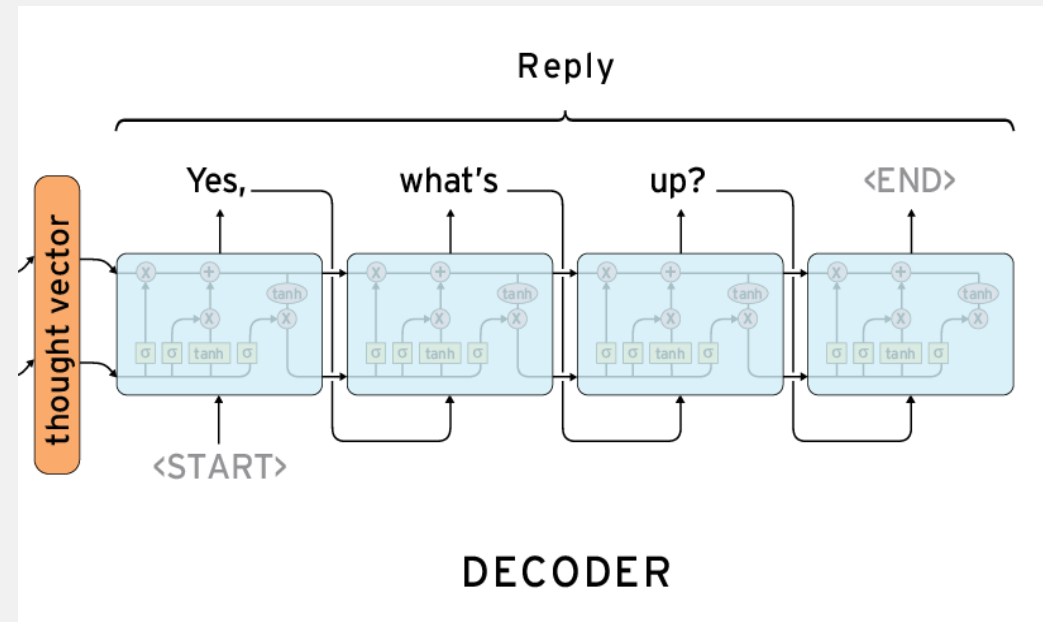
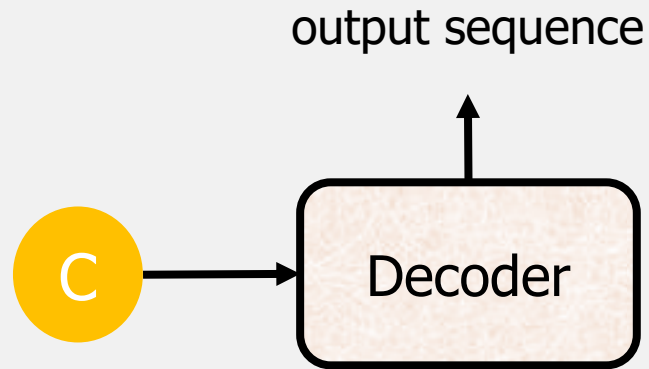
Encoder is responsible for conversion of input sequence into a vector
The vector is context vector composed of important information of input sequence

It can be done by RNN or LSTM

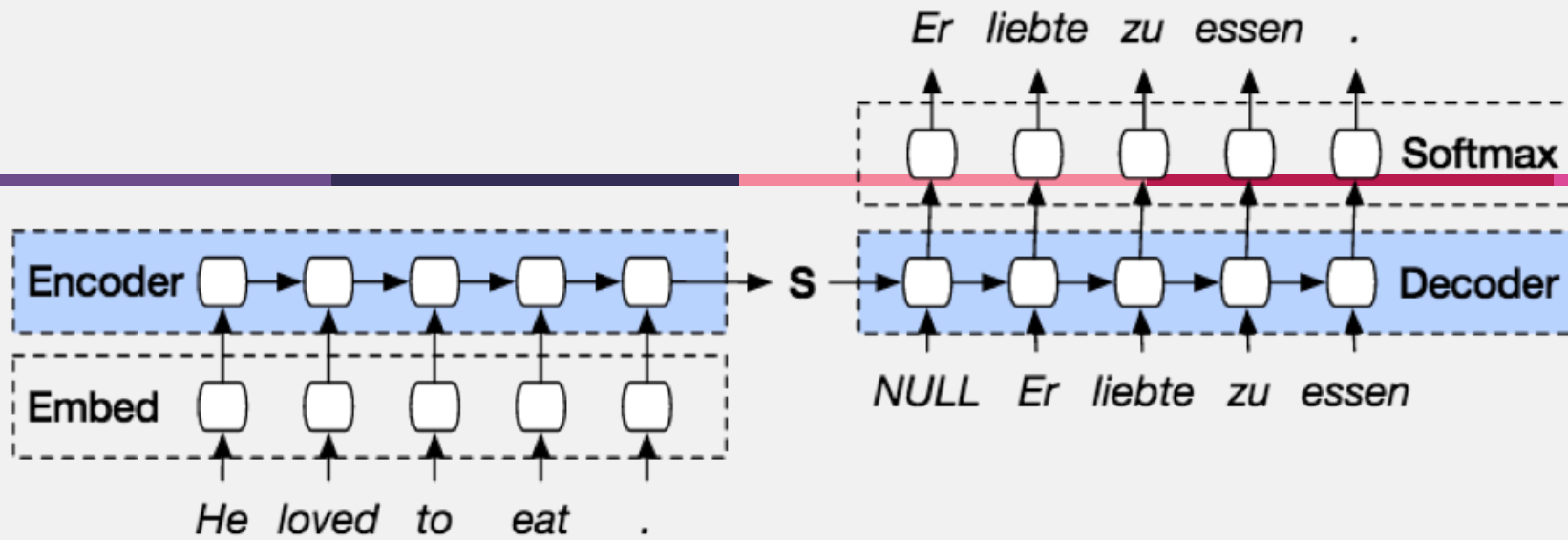


Seq2Seq

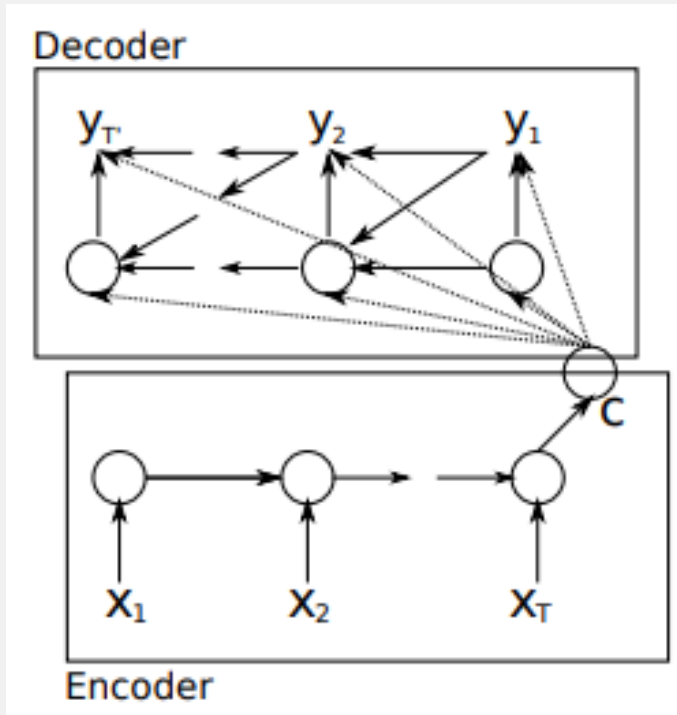
Decoder is responsible for word generation from context vector
It also can be done by RNN or LSTM



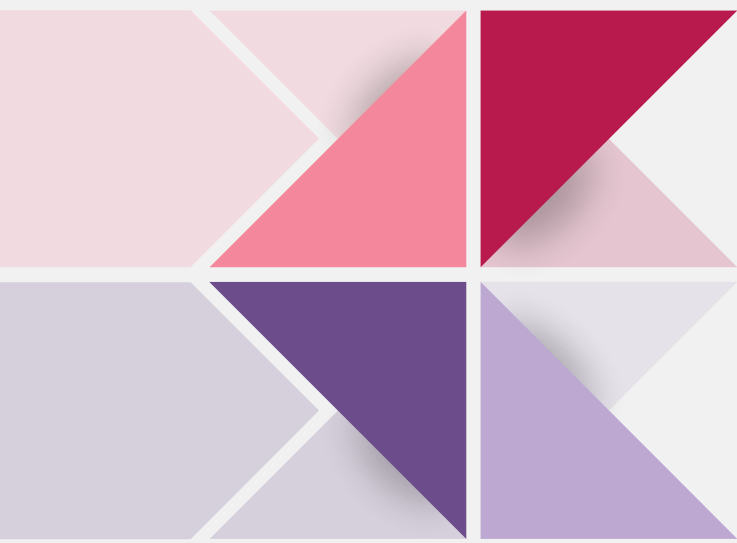
Seq2Seq



Example



$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$



C

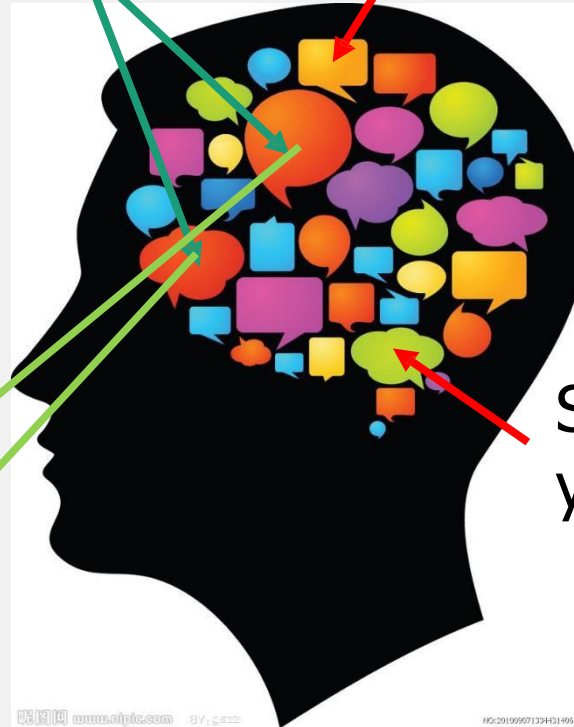
Attention

Attention

What you learned
in these lectures

Breakfast today

What is deep
learning?

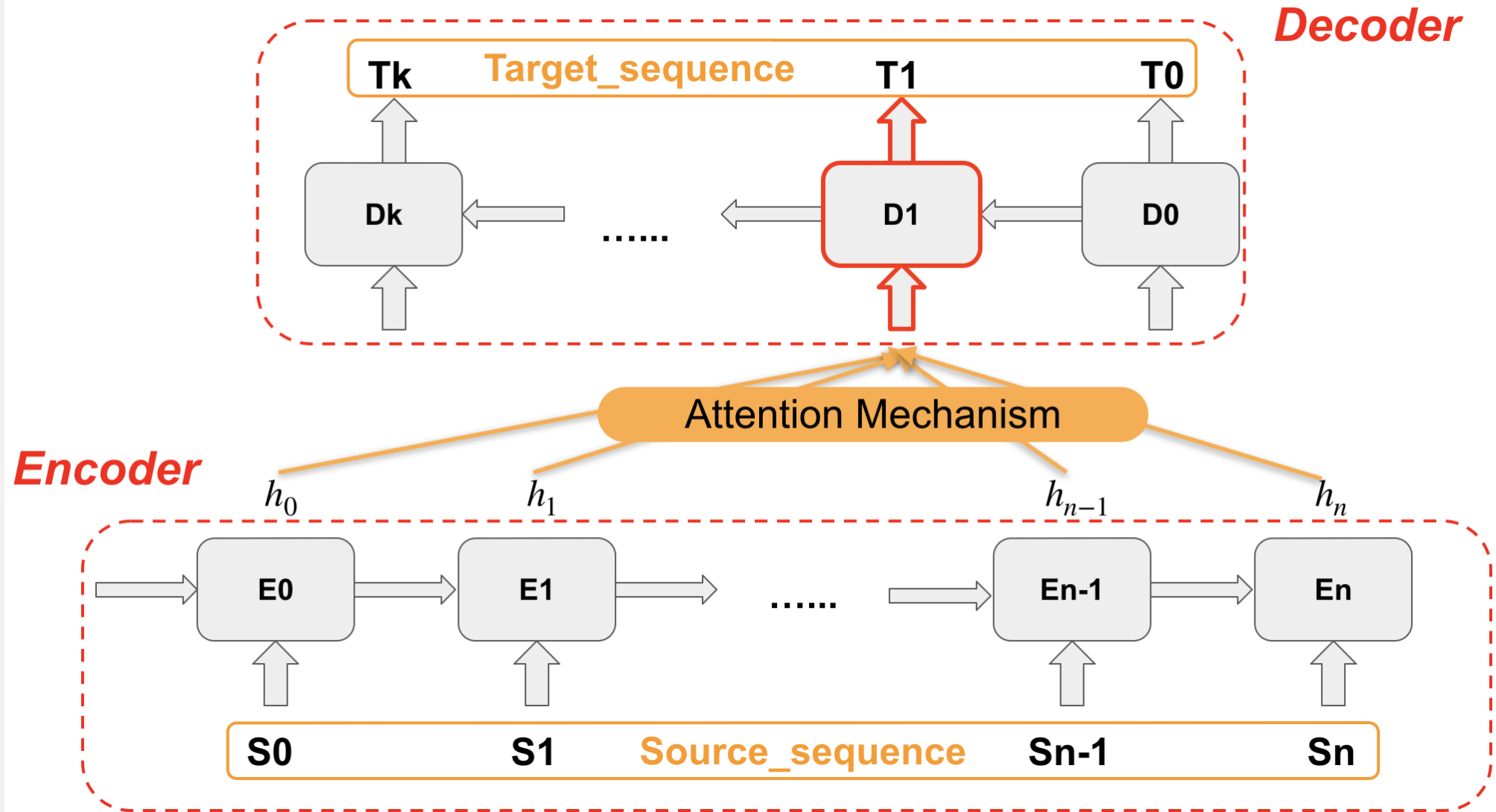


Summer vacation 10
years ago

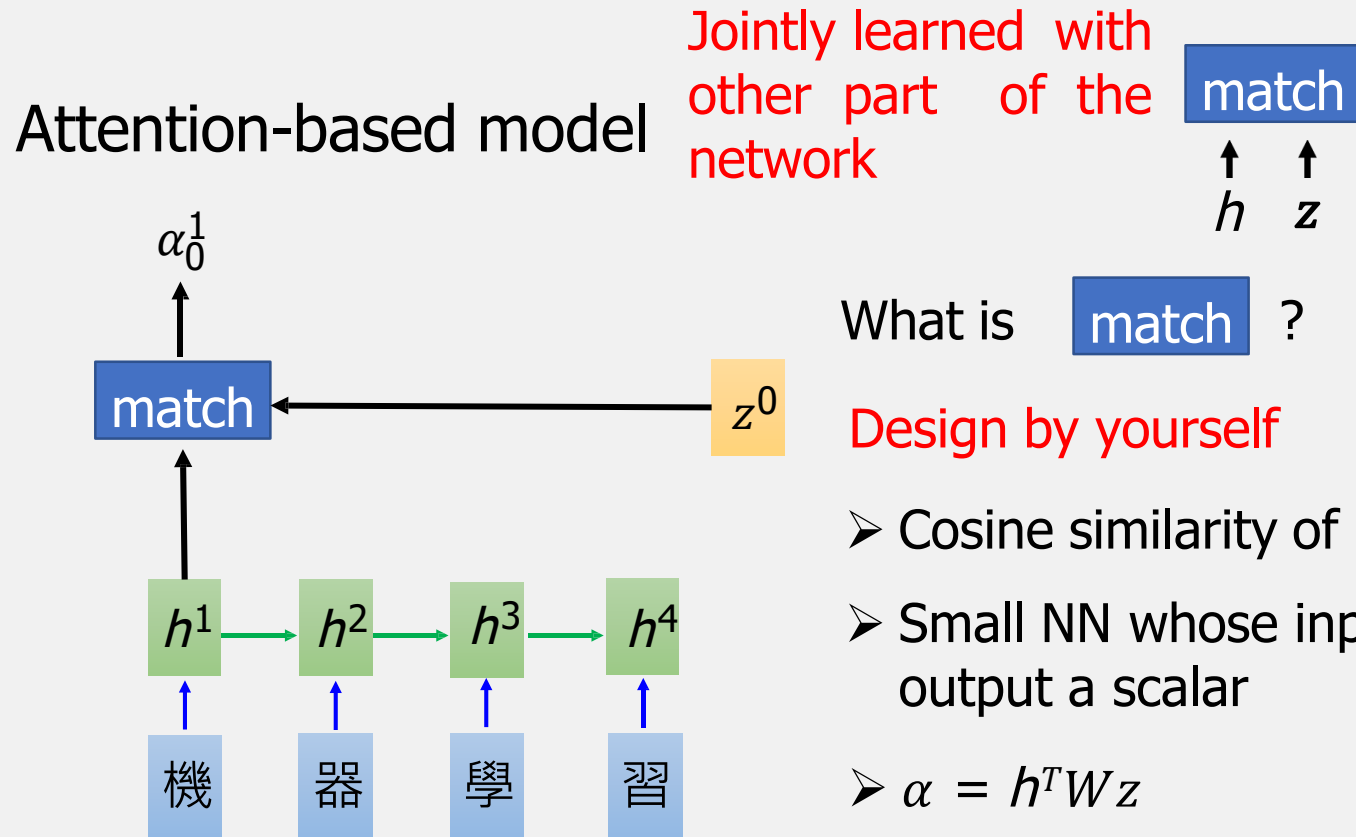
Answer ← **Organize**

[Reference](#)

Attention

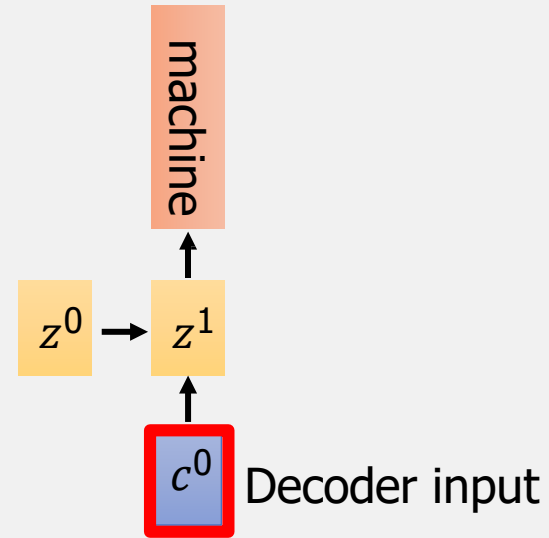
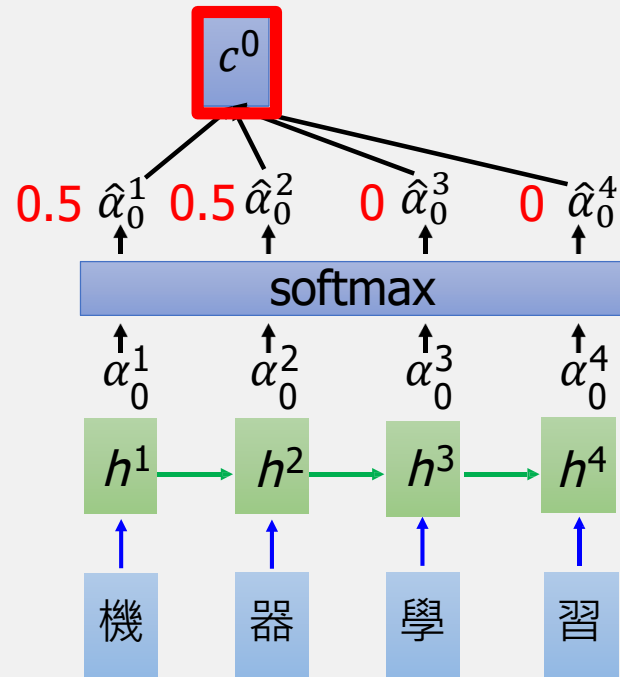


Attention



Attention

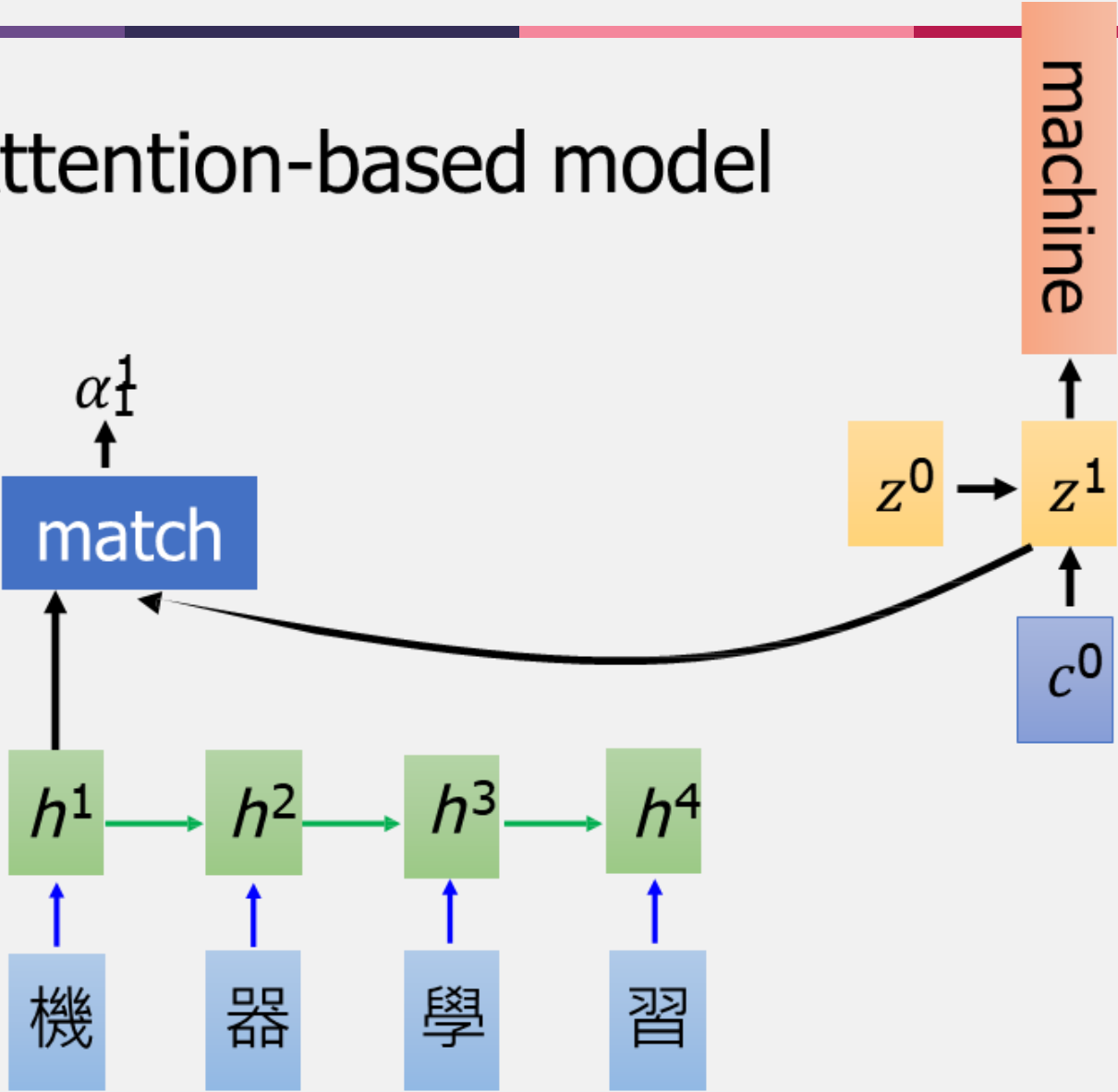
Attention-based model



$$c^0 = \sum \hat{\alpha}_0^i h^i$$
$$= 0.5h^1 + 0.5h^2$$

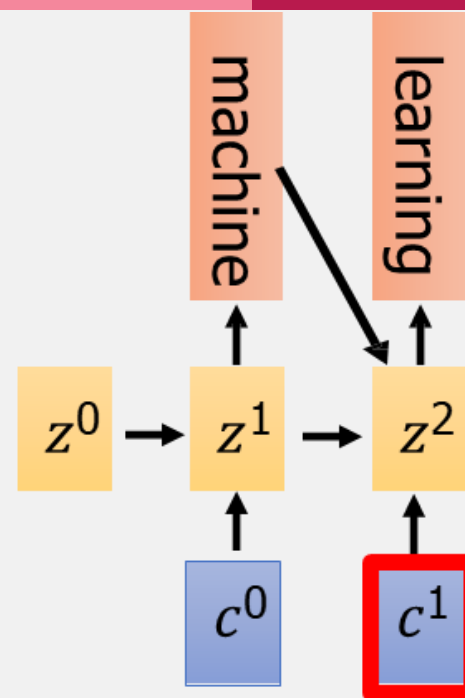
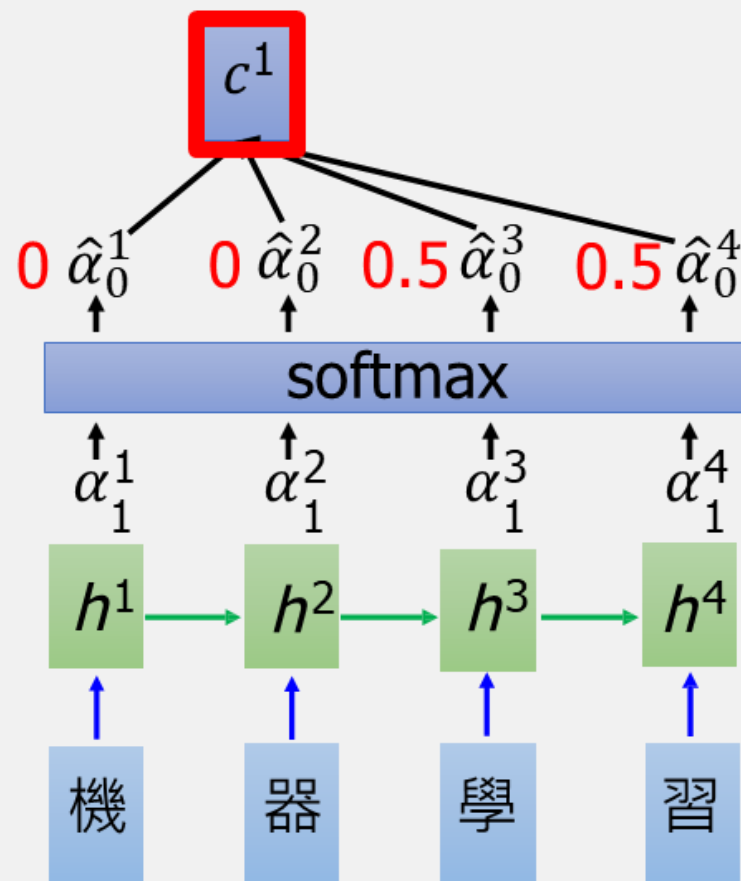
Attention

Attention-based model



Attention

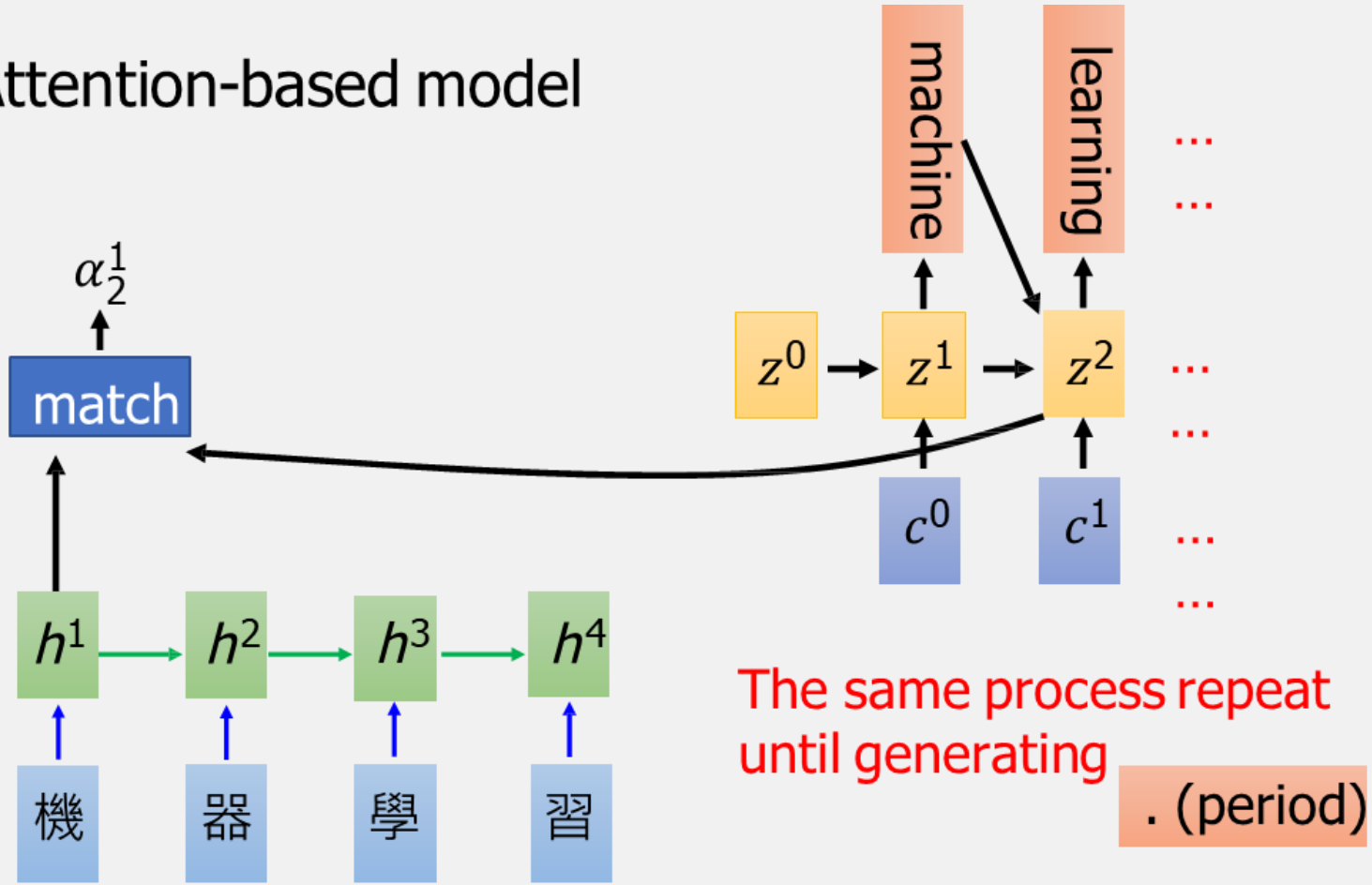
Attention-based model

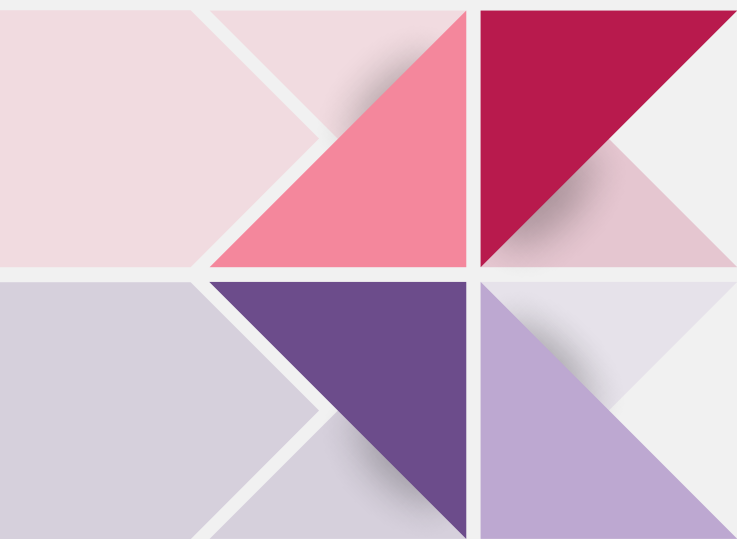


$$c^1 = \sum \hat{\alpha}_0^i h^i$$
$$= 0.5h^3 + 0.5h^4$$

Attention

Attention-based model





D

Self-attention

Attention

gif5.net

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

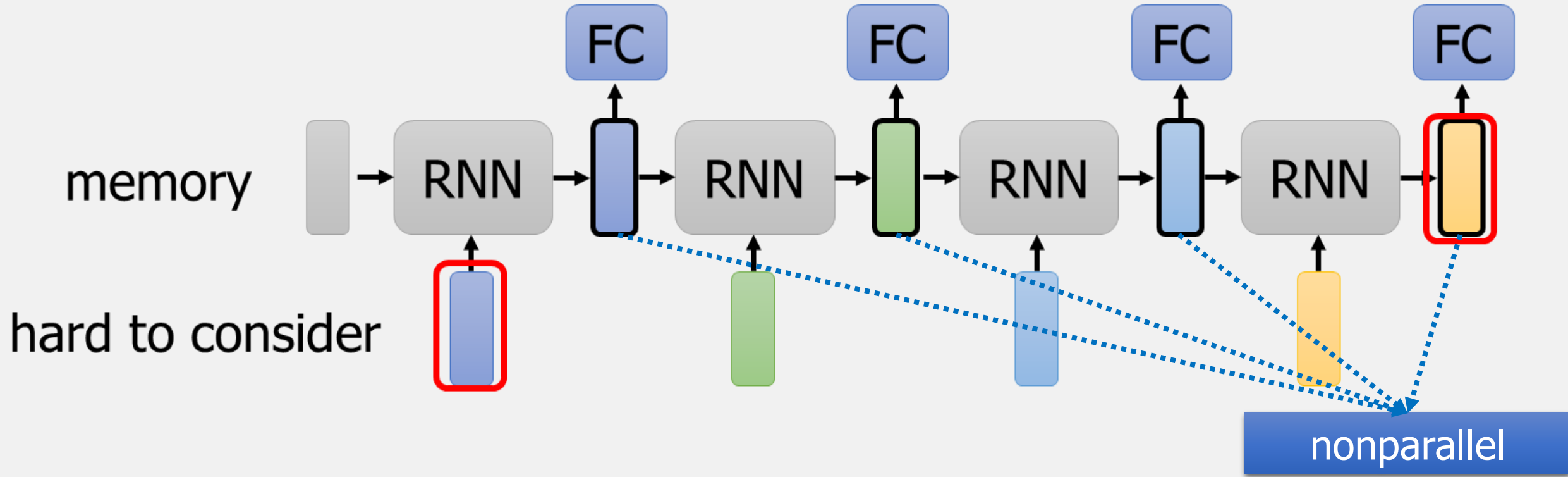


Je

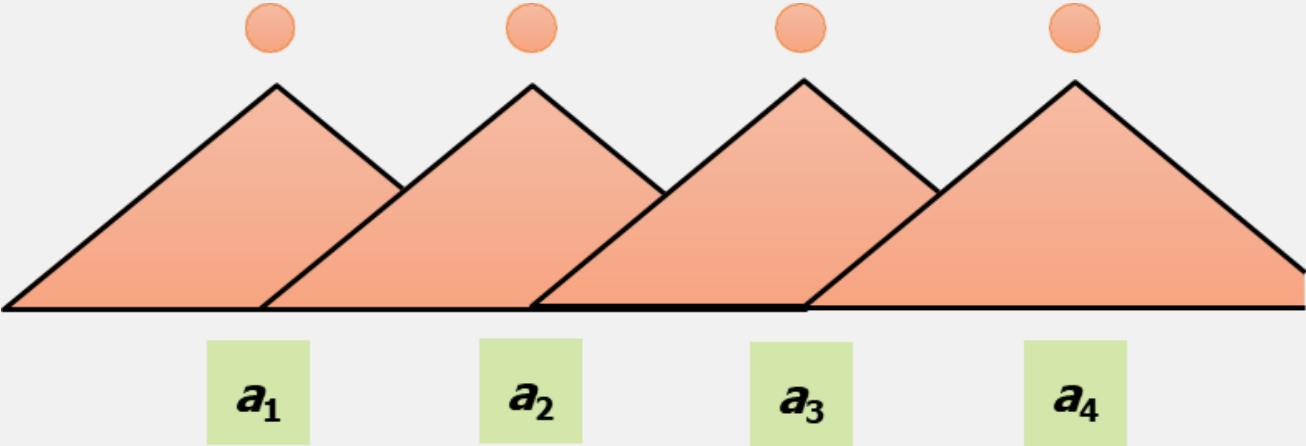
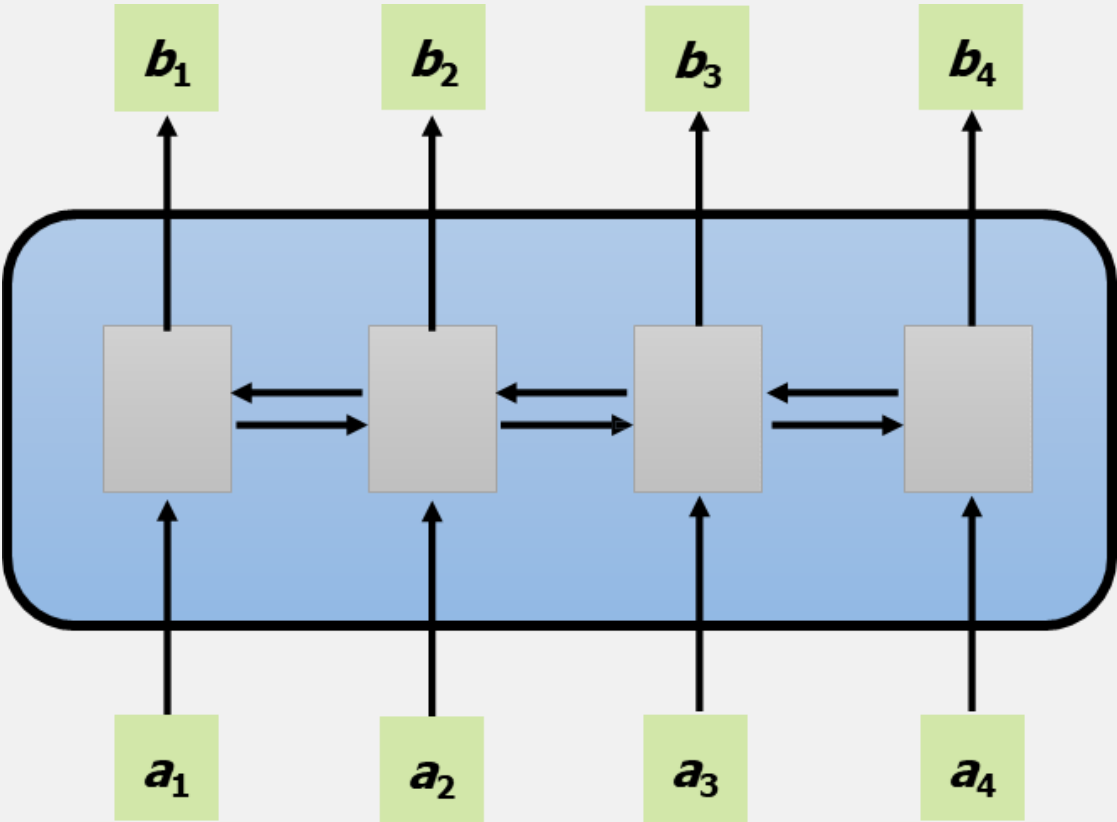
suis

étudiant

Attention



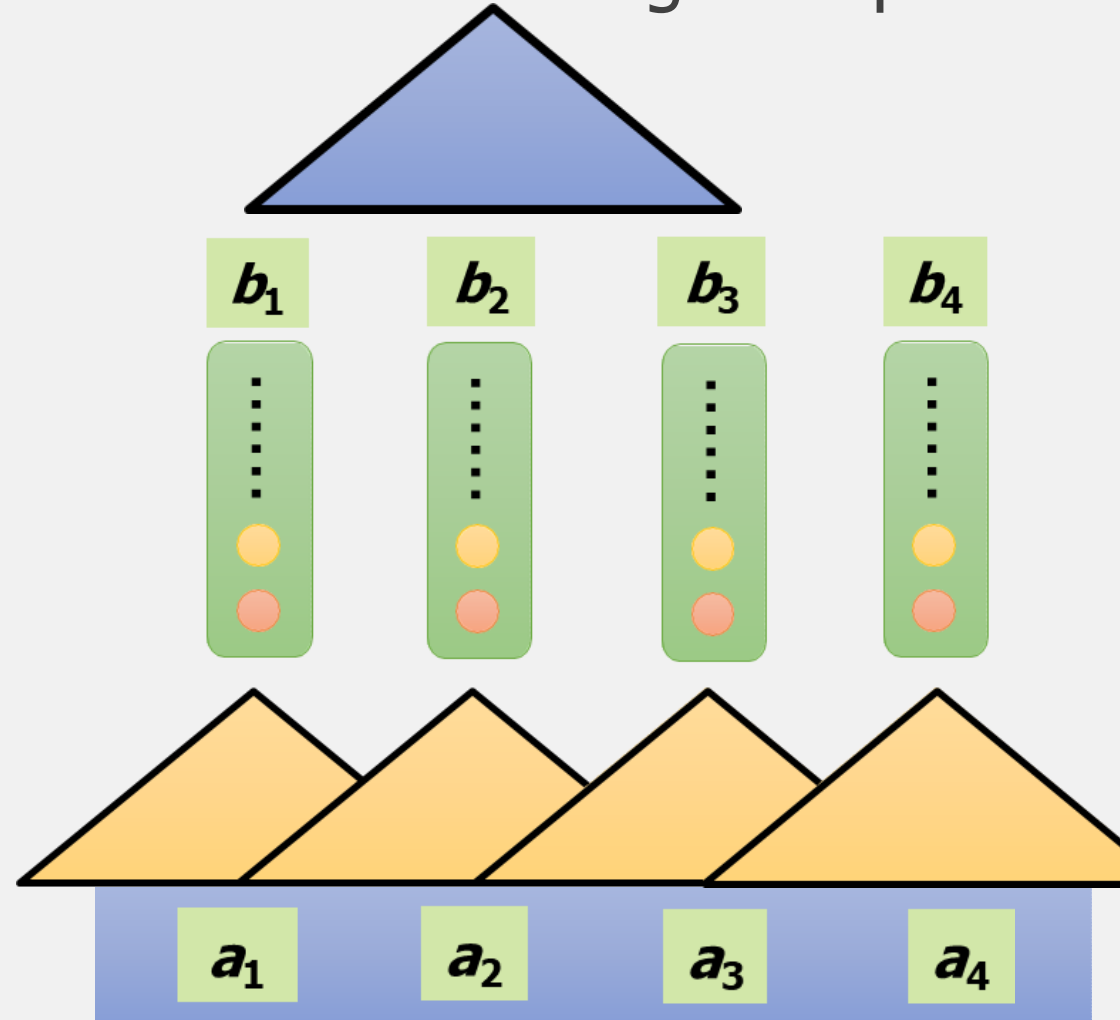
Attention



Using CNN to replace RNN

Attention

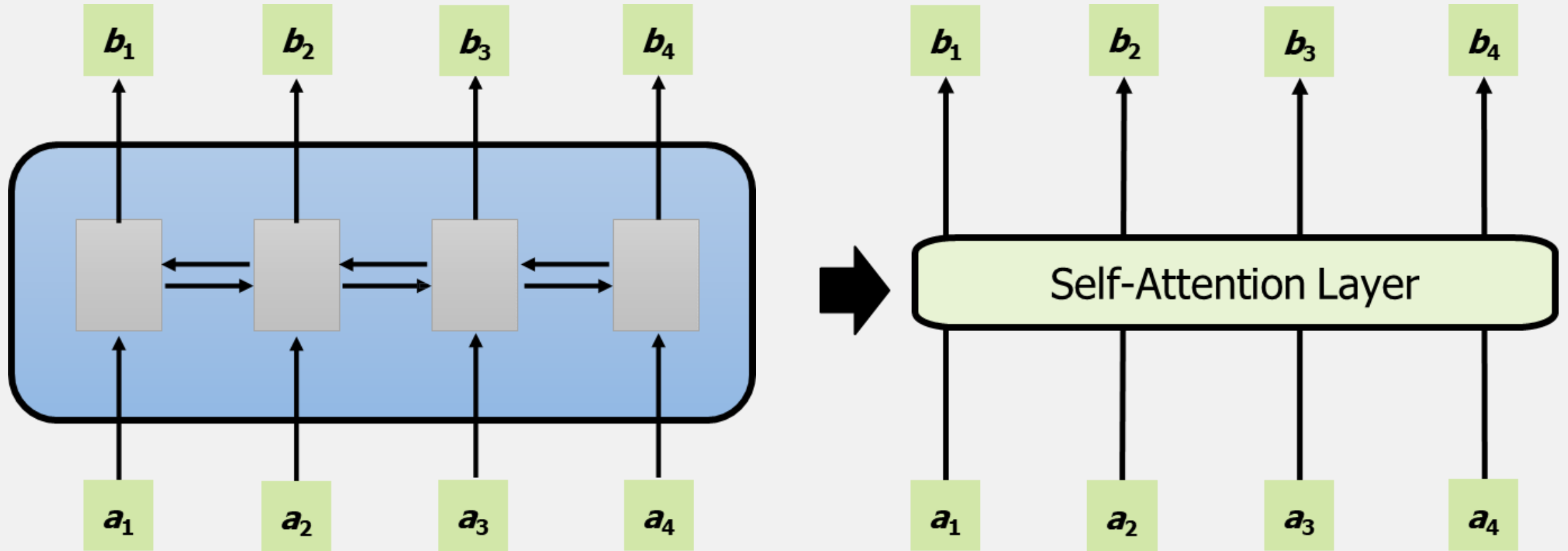
Filters in higher layer can consider longer sequence



Attention

Self-attention

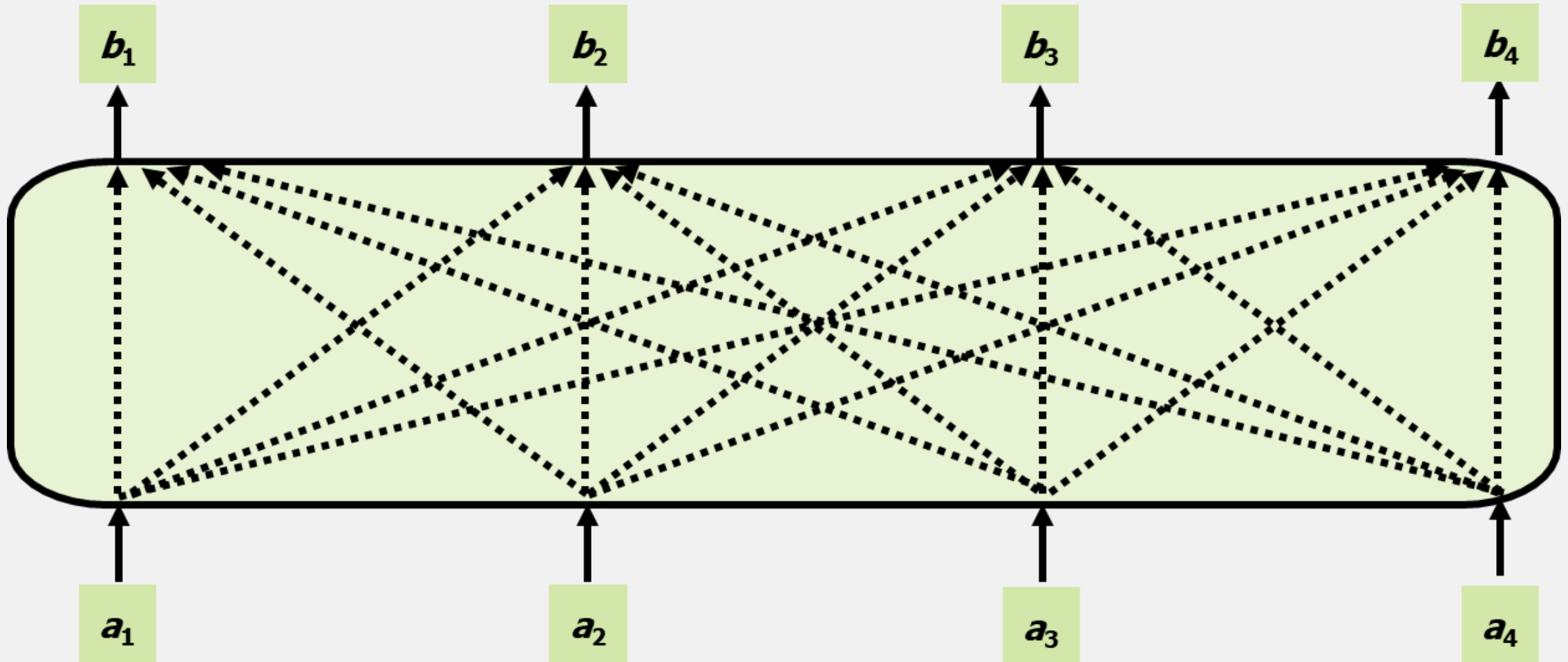
b_i is obtained based on the whole input sequence
 b_1, b_2, b_3, b_4 can be parallelly computed



You can try to replace any thing that has been done by RNN with self-attention

Attention

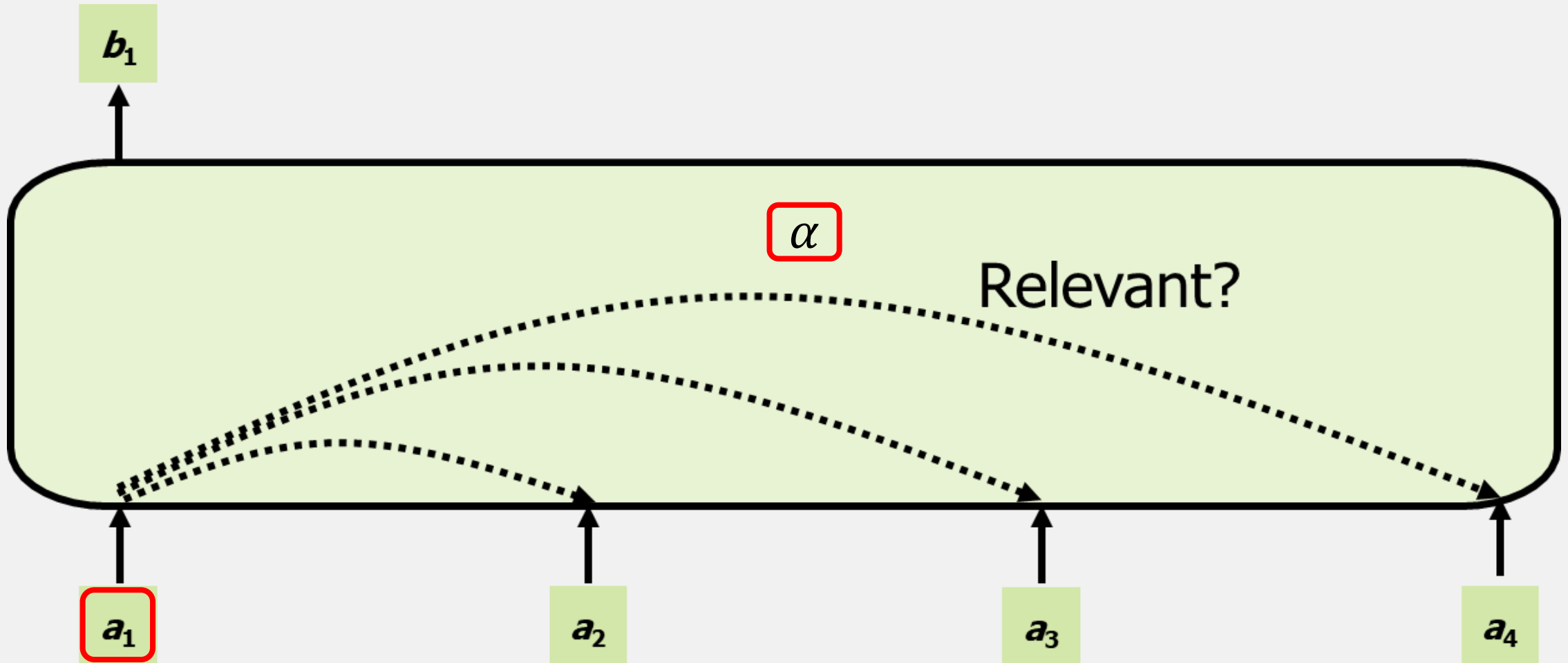
Self-attention



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

Attention

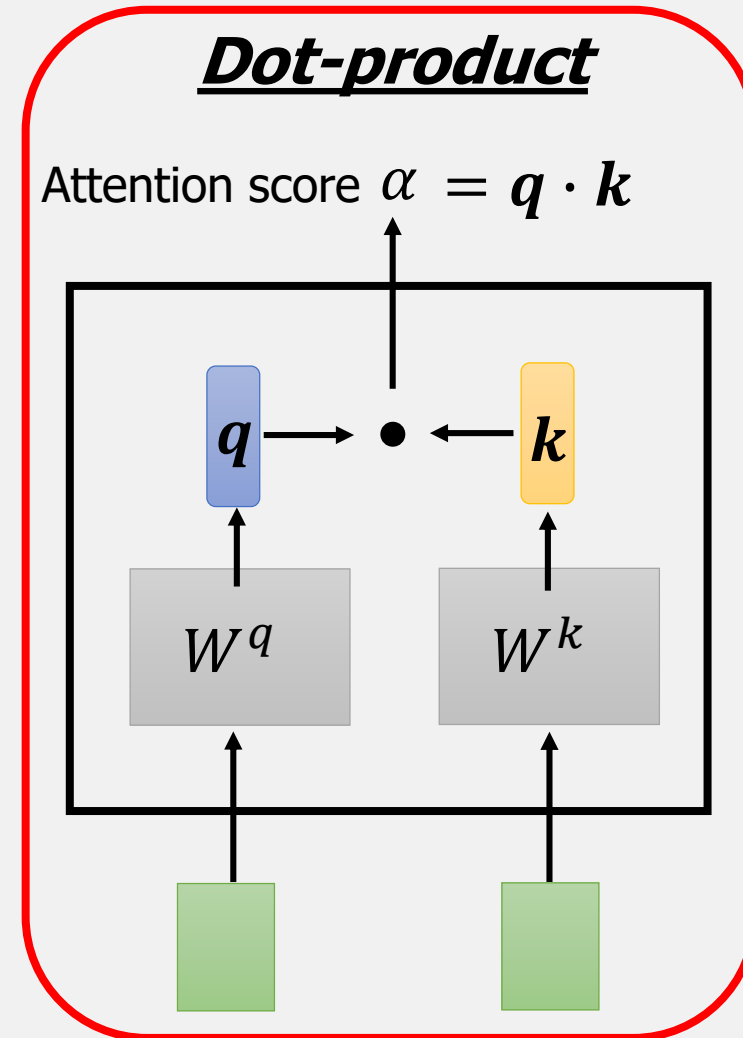
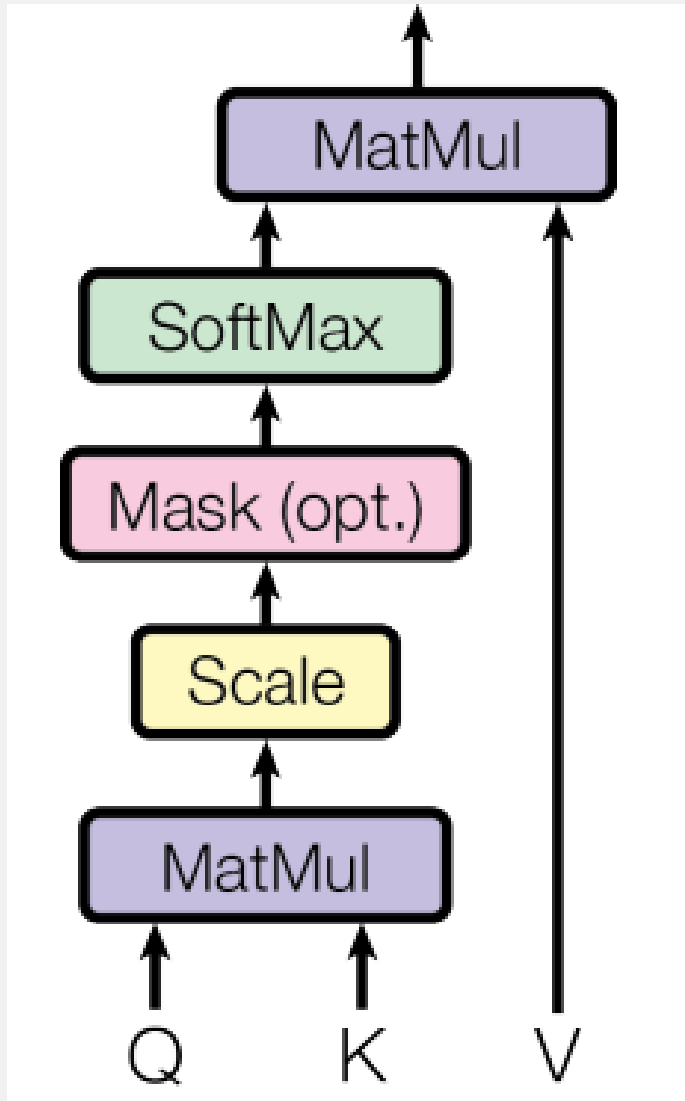
Self-attention



Find the relevant vectors in a sequence

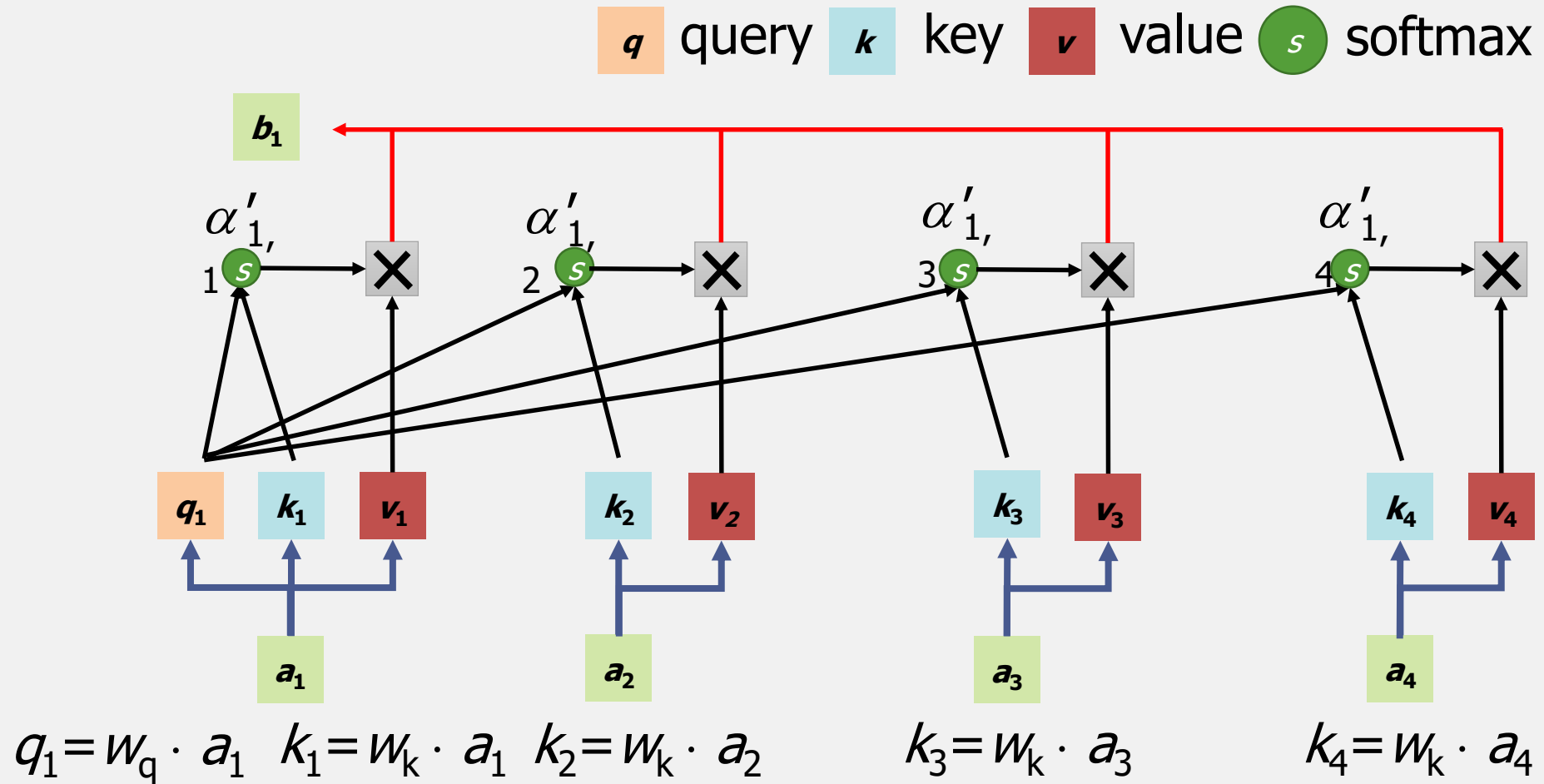
Attention

Self-attention



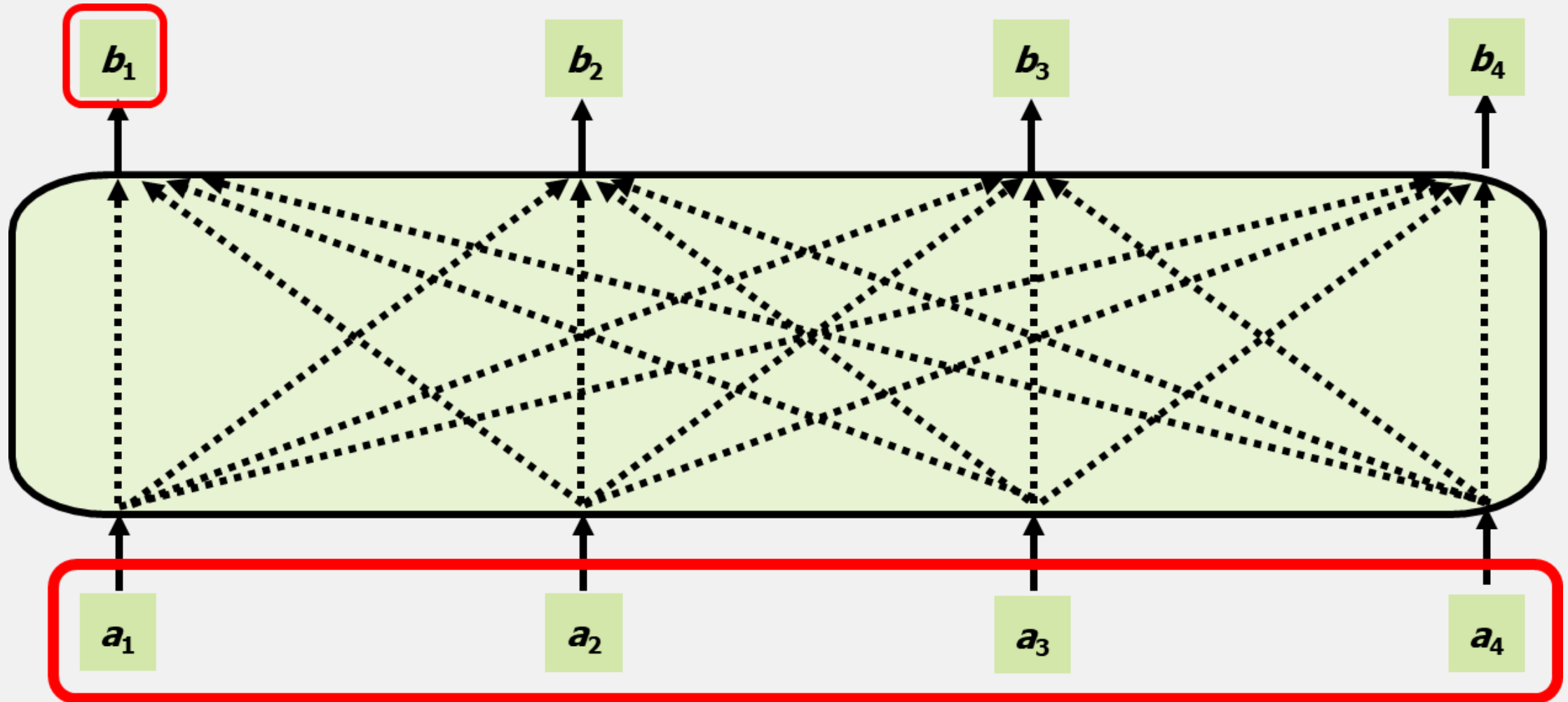
Attention

Self-attention



Attention

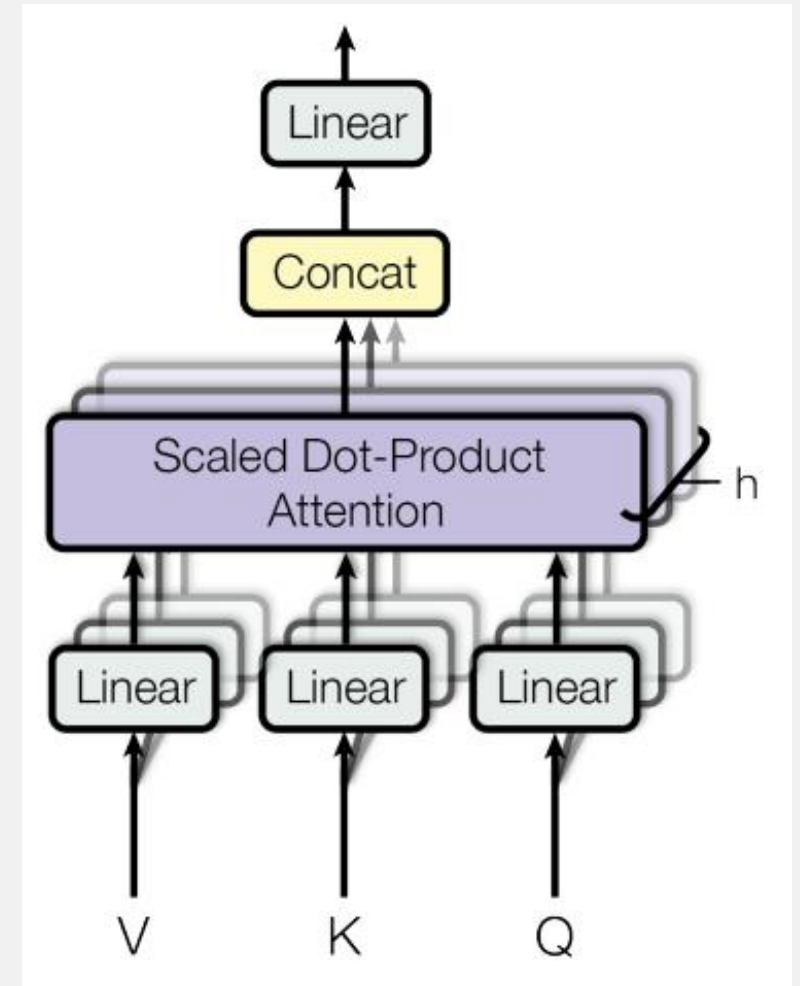
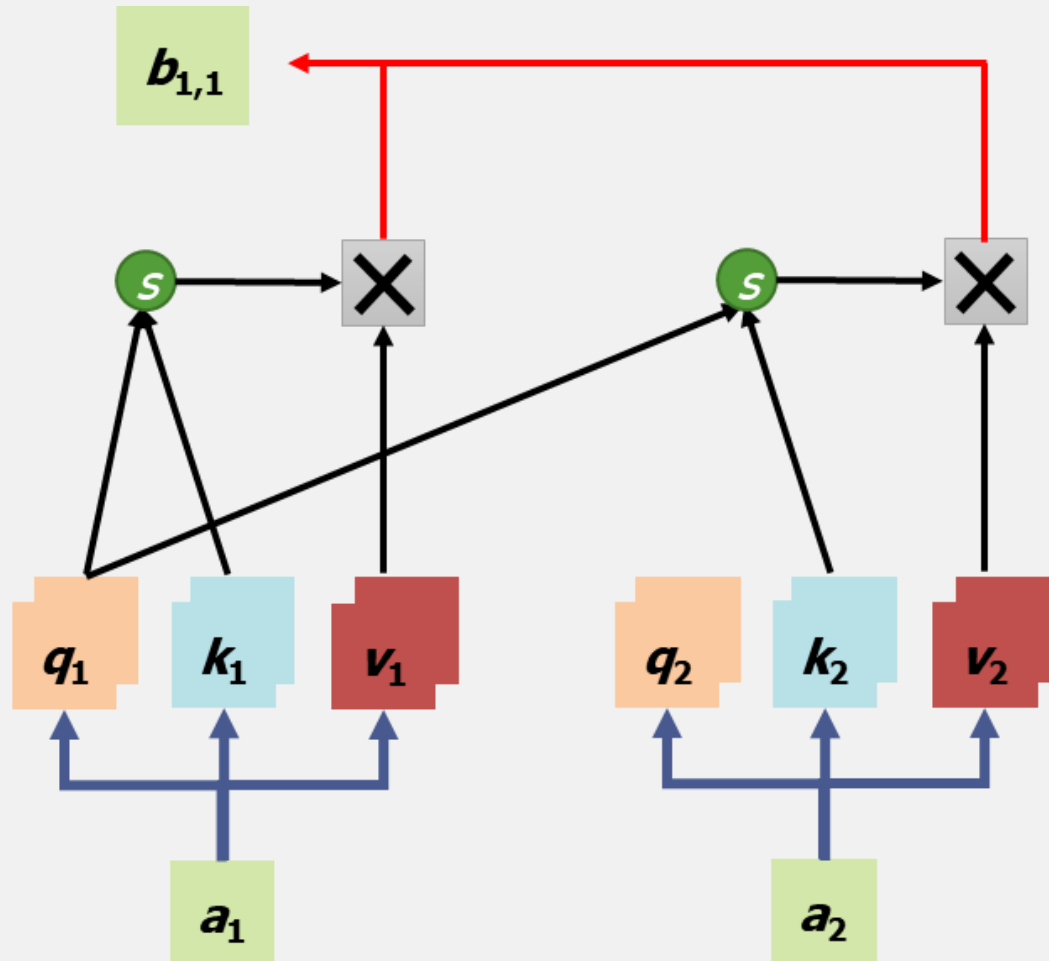
Self-attention



Attention

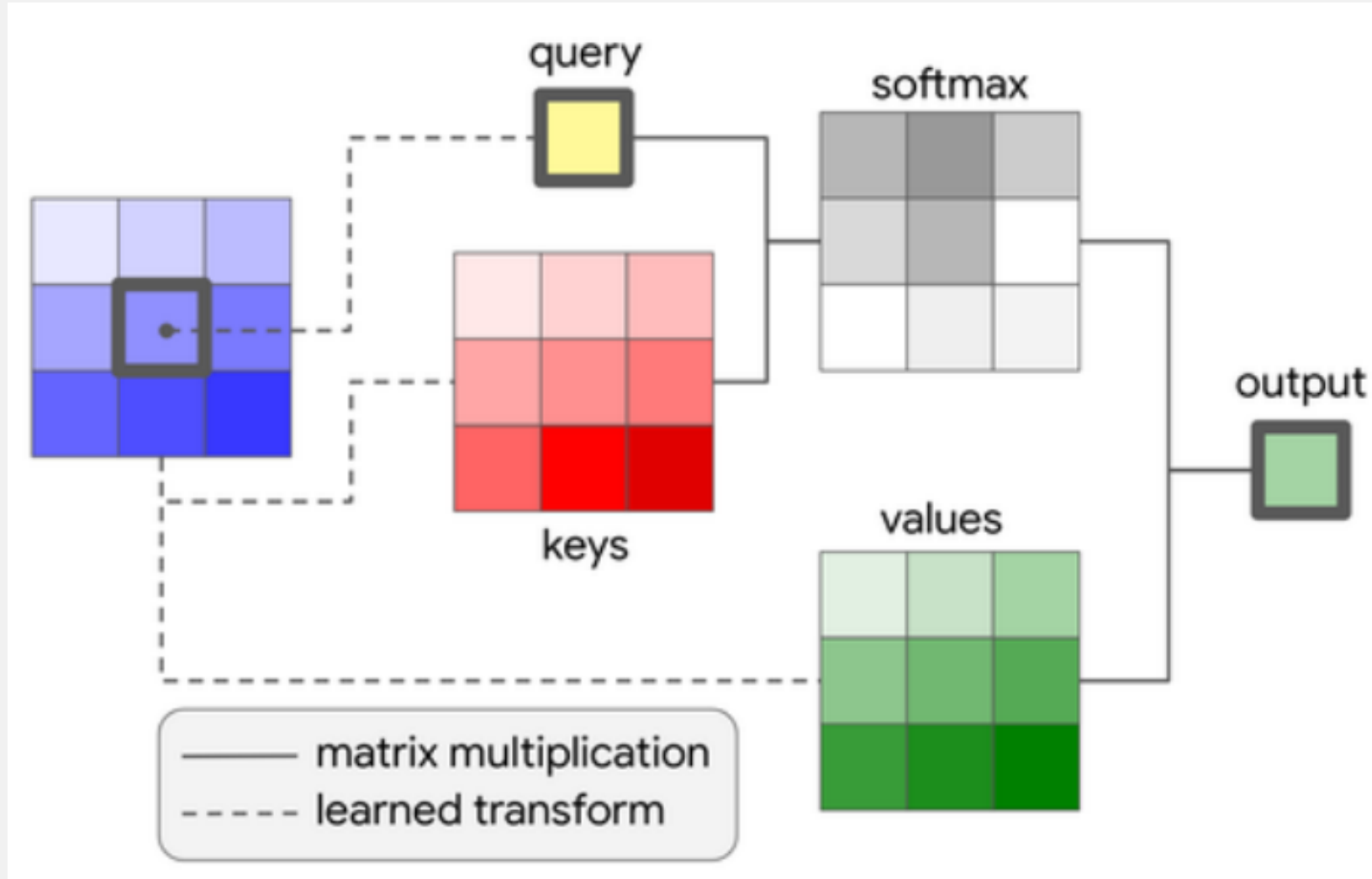
Self-attention

q query k key v value s softmax



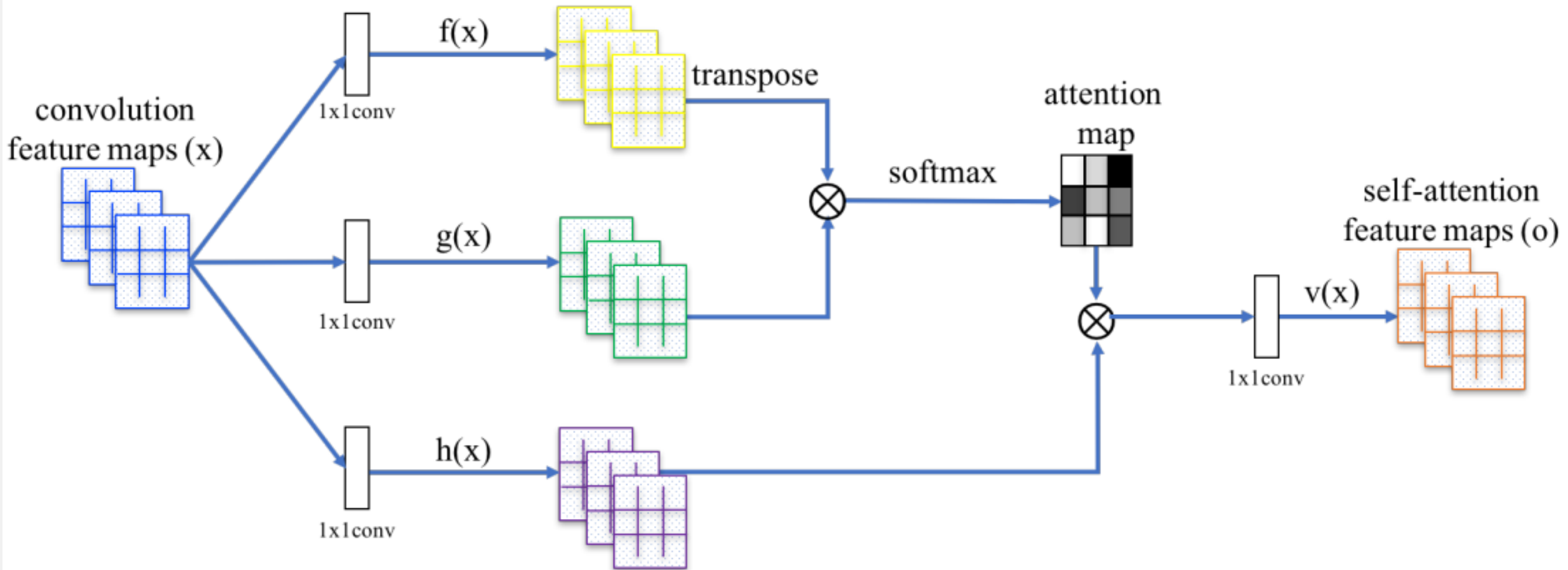
Attention

Self-attention in Vision model



Attention

Self-attention GAN

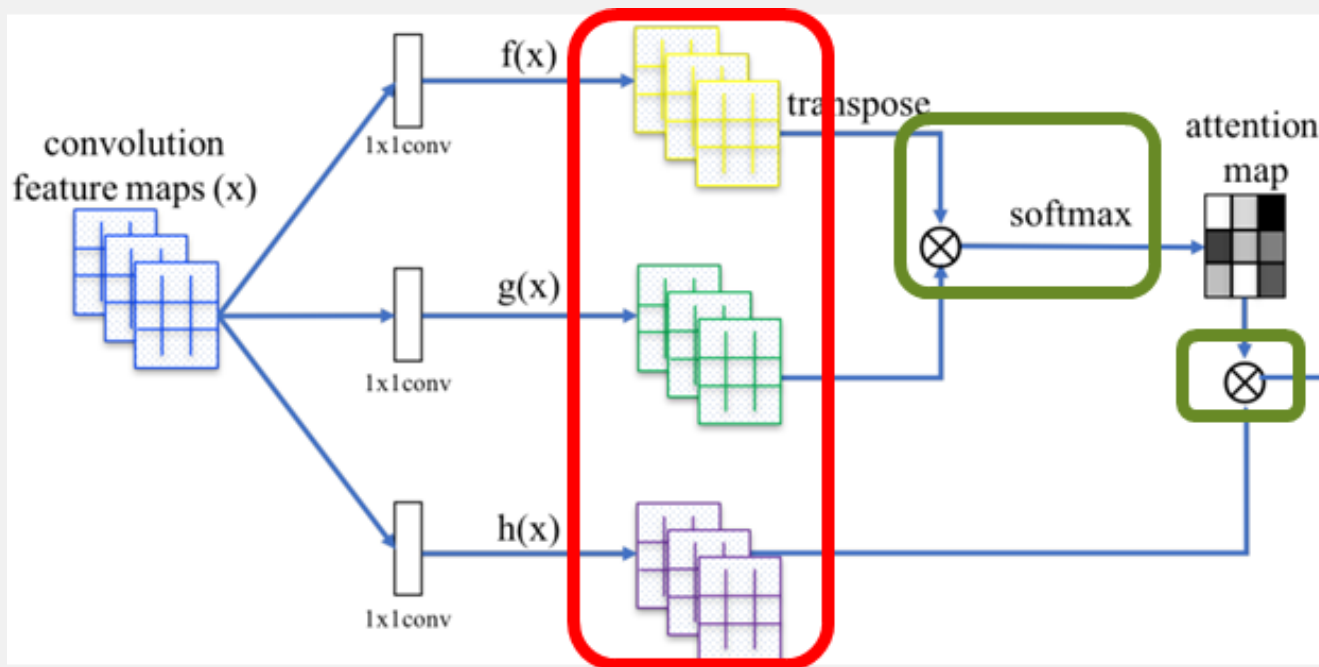


Attention

Self-attention GAN

```
def __init__(self, in_channels):  
    super(Self_Attn, self).__init__()  
    self.in_channels = in_channels  
    self.snconv1x1_theta = snconv2d(in_channels=in_channels, out_channels=in_channels, kernel_size=1, stride=1, padding=0)  
    self.snconv1x1_phi = snconv2d(in_channels=in_channels, out_channels=in_channels, kernel_size=1, stride=1, padding=0)  
    self.snconv1x1_g = snconv2d(in_channels=in_channels, out_channels=in_channels, kernel_size=1, stride=1, padding=0)  
    self.snconv1x1_attn = snconv2d(in_channels=in_channels//2, out_channels=in_channels//2, kernel_size=1, stride=1, padding=0)
```

```
# Theta path  
theta = self.snconv1x1_theta(x)  
theta = theta.view(-1, ch//8, h*w)  
# Phi path  
phi = self.snconv1x1_phi(x)  
phi = self.maxpool(phi)  
phi = phi.view(-1, ch//8, h*w//4)
```



```
# Attn map  
attn = torch.bmm(theta.permute(0, 2, 1), phi)  
attn = self.softmax(attn)
```

```
# g path  
g = self.snconv1x1_g(x)  
g = self.maxpool(g)  
g = g.view(-1, ch//2, h*w//4)
```

```
# Attn_g  
attn_g = torch.bmm(g, attn.permute(0, 2, 1))
```

```
attn_g = attn_g.view(-1, ch//2, h, w)  
attn_g = self.snconv1x1_attn(attn_g)
```

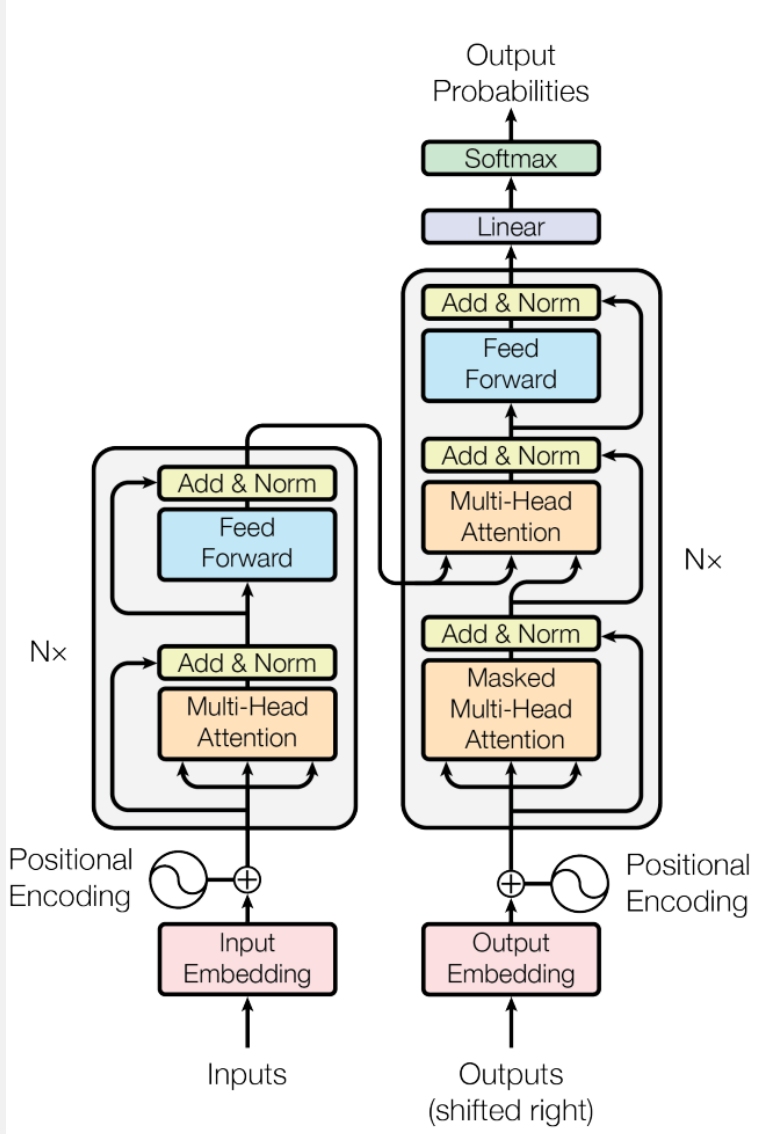
```
# Out  
out = x + self.sigma*attn_g  
return out
```



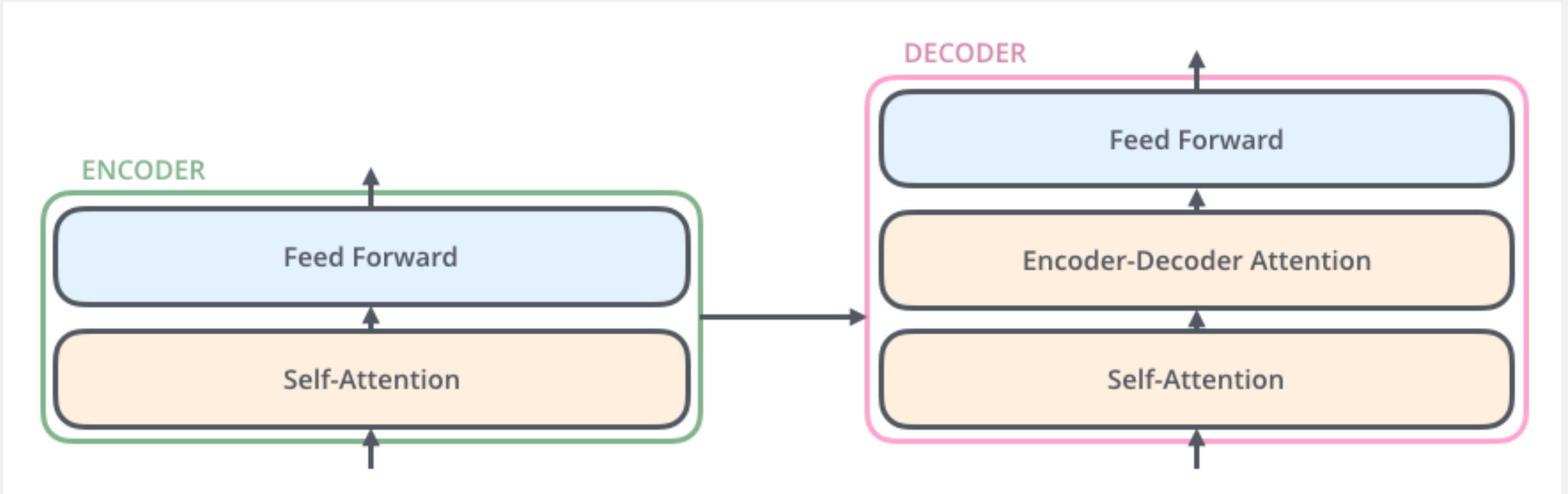
02

Transformer

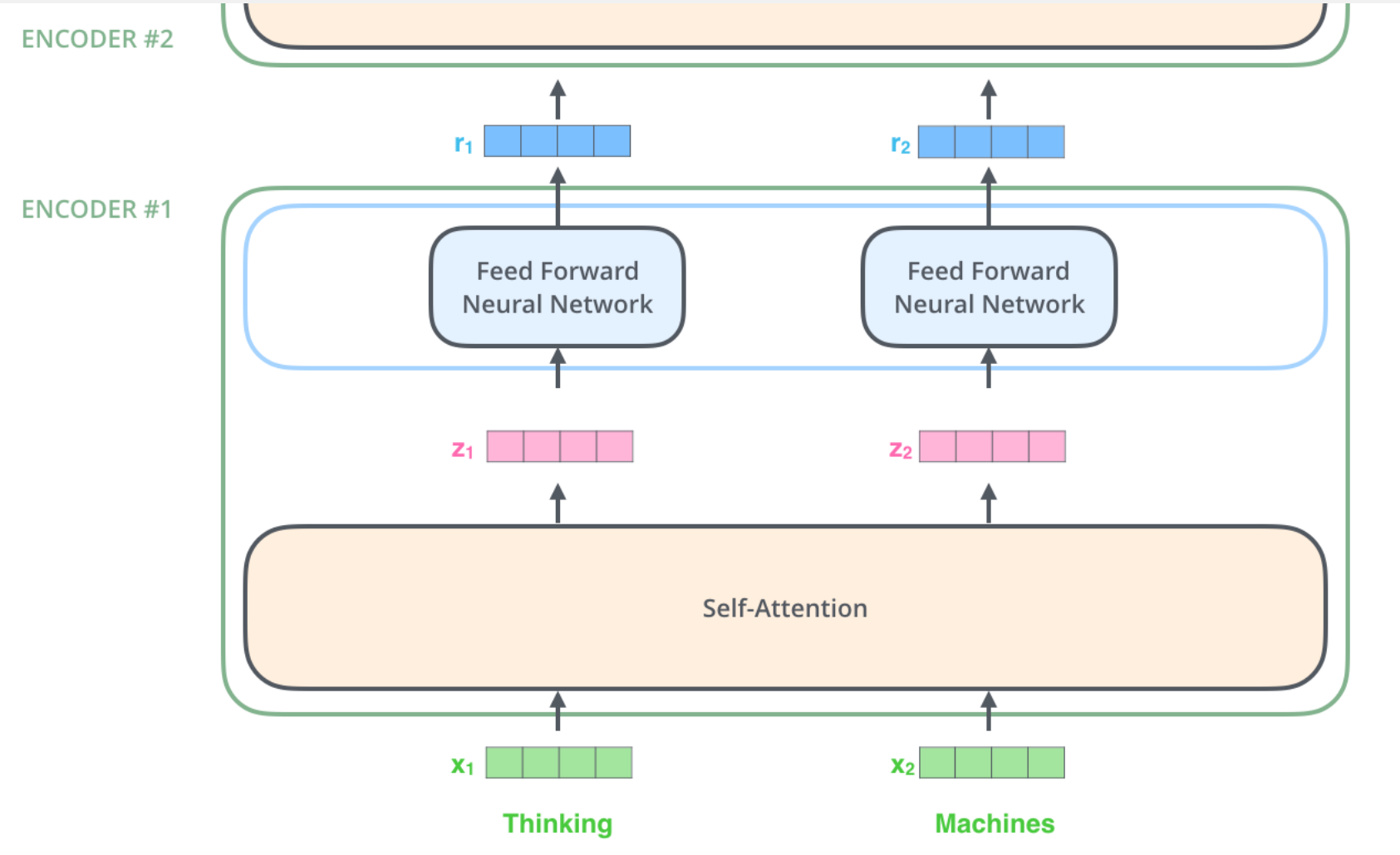
Transformer



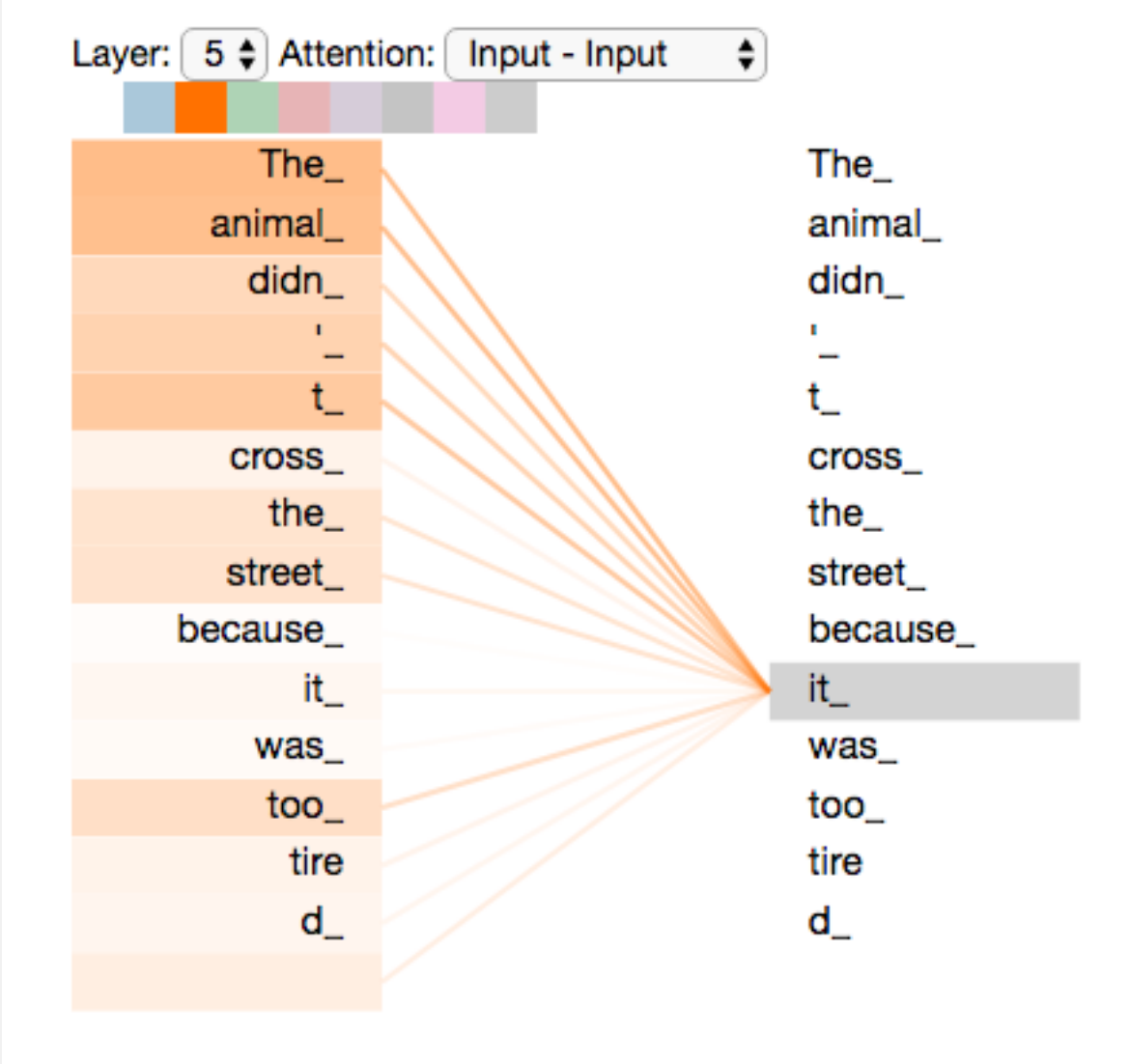
Transformer



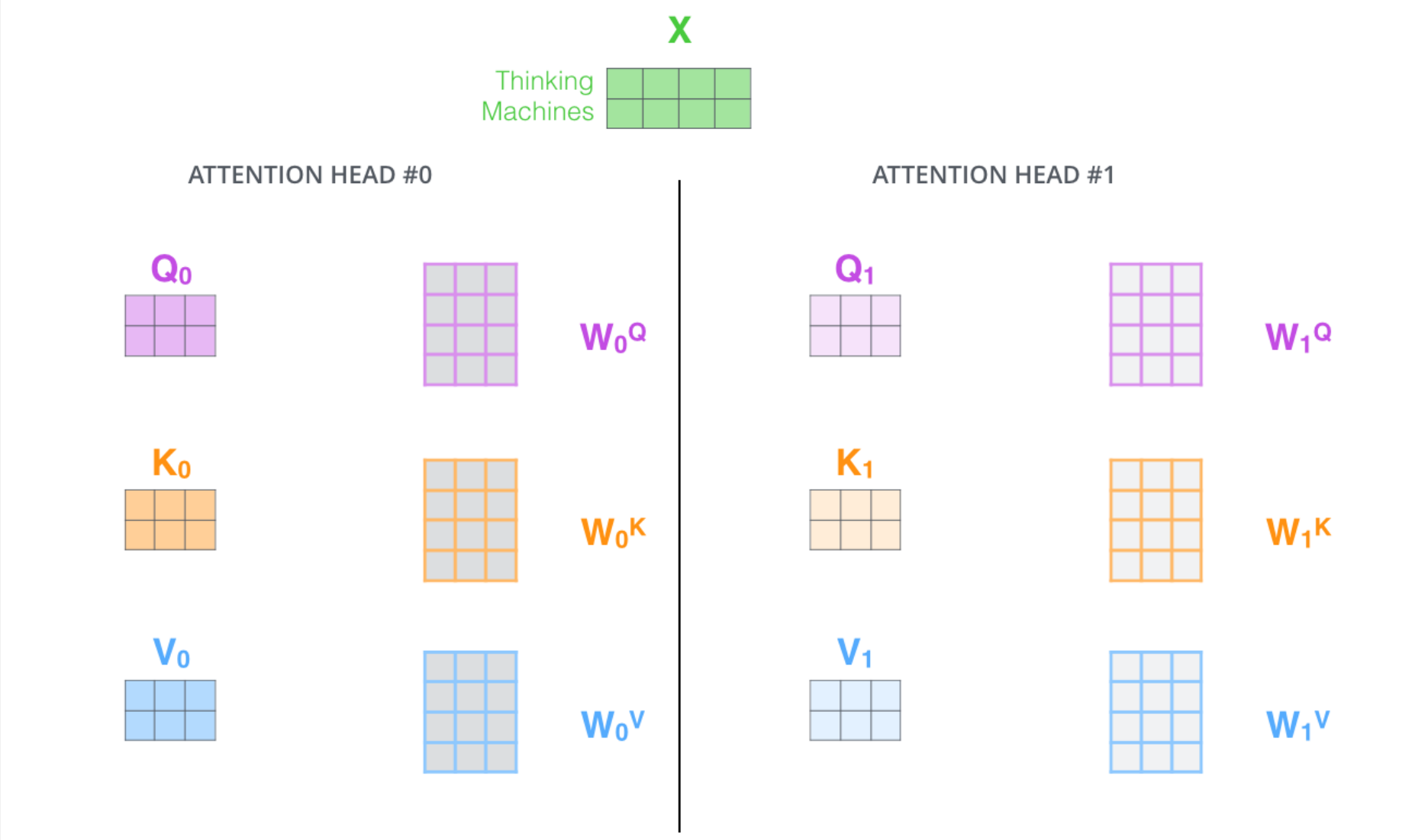
Transformer



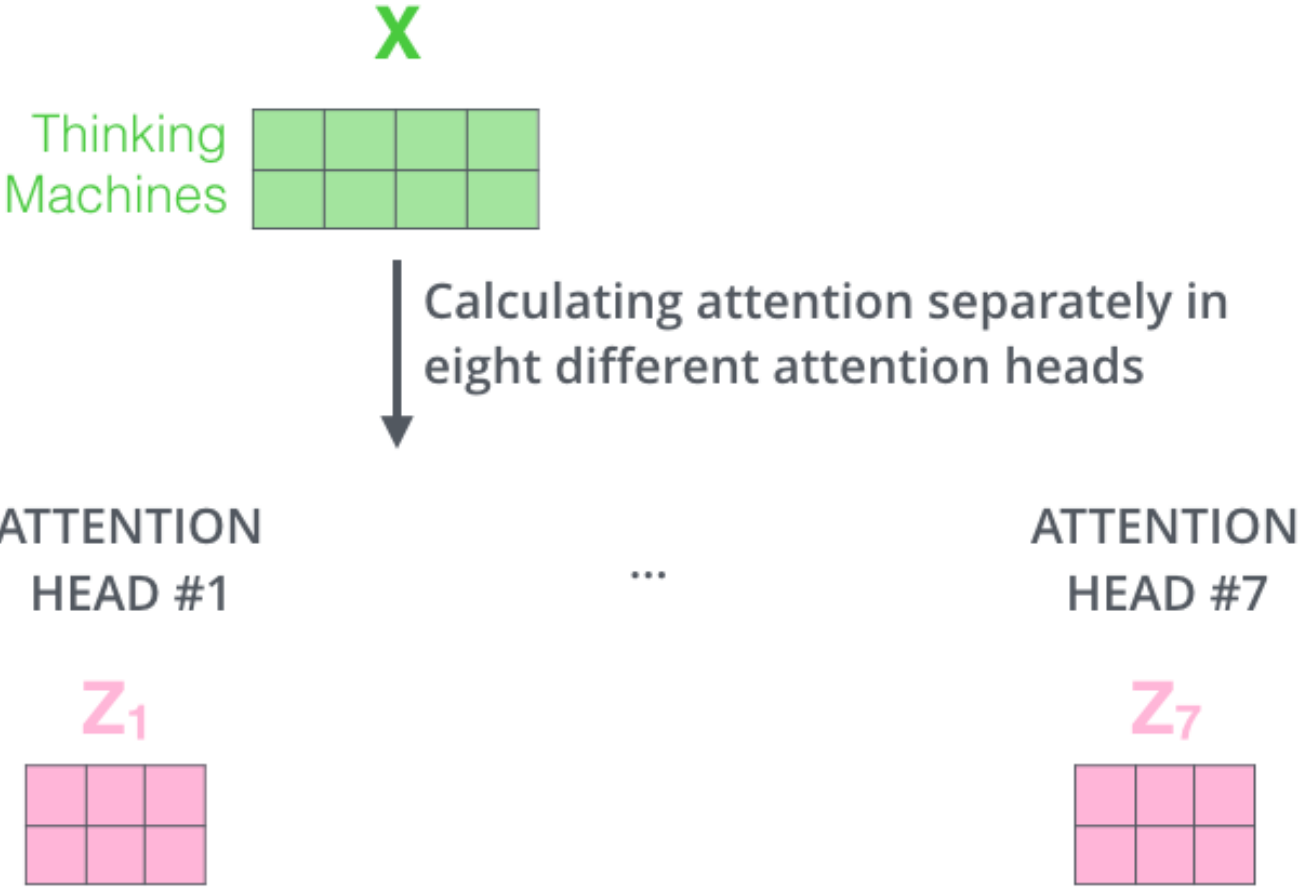
Transformer



Transformer



Transformer



Transformer

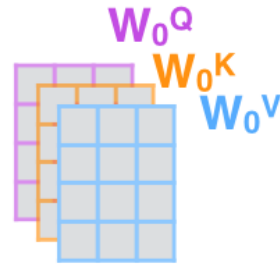
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



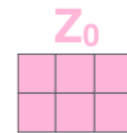
3) Split into 8 heads. We multiply X or R with weight matrices



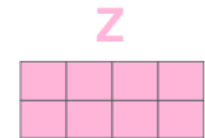
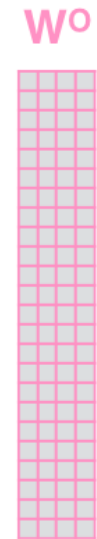
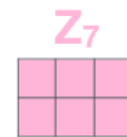
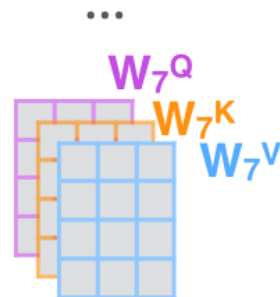
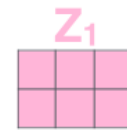
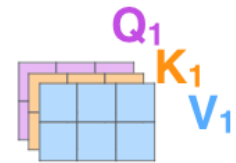
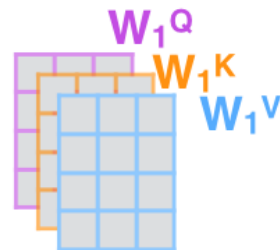
4) Calculate attention using the resulting $Q/K/V$ matrices



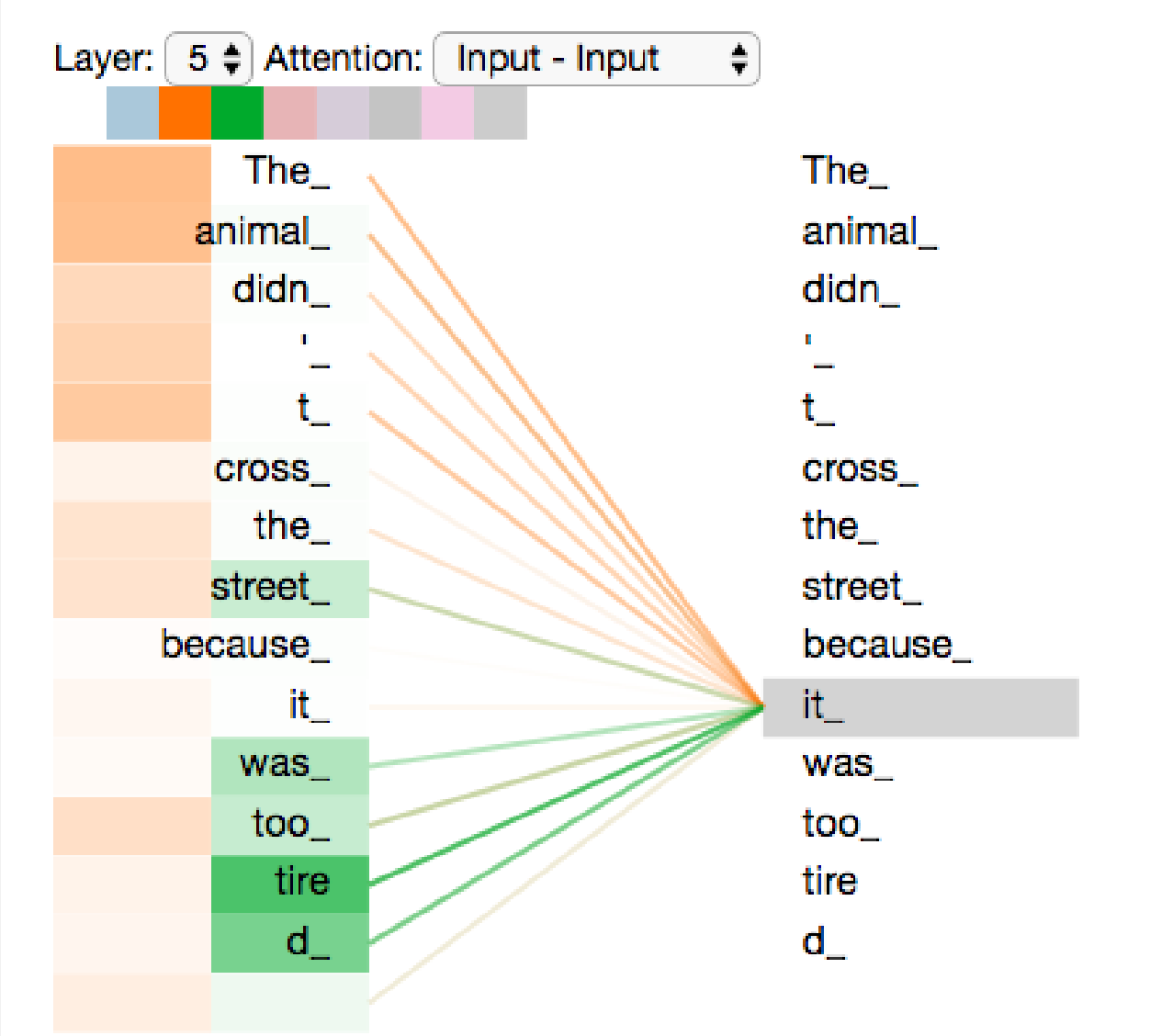
5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



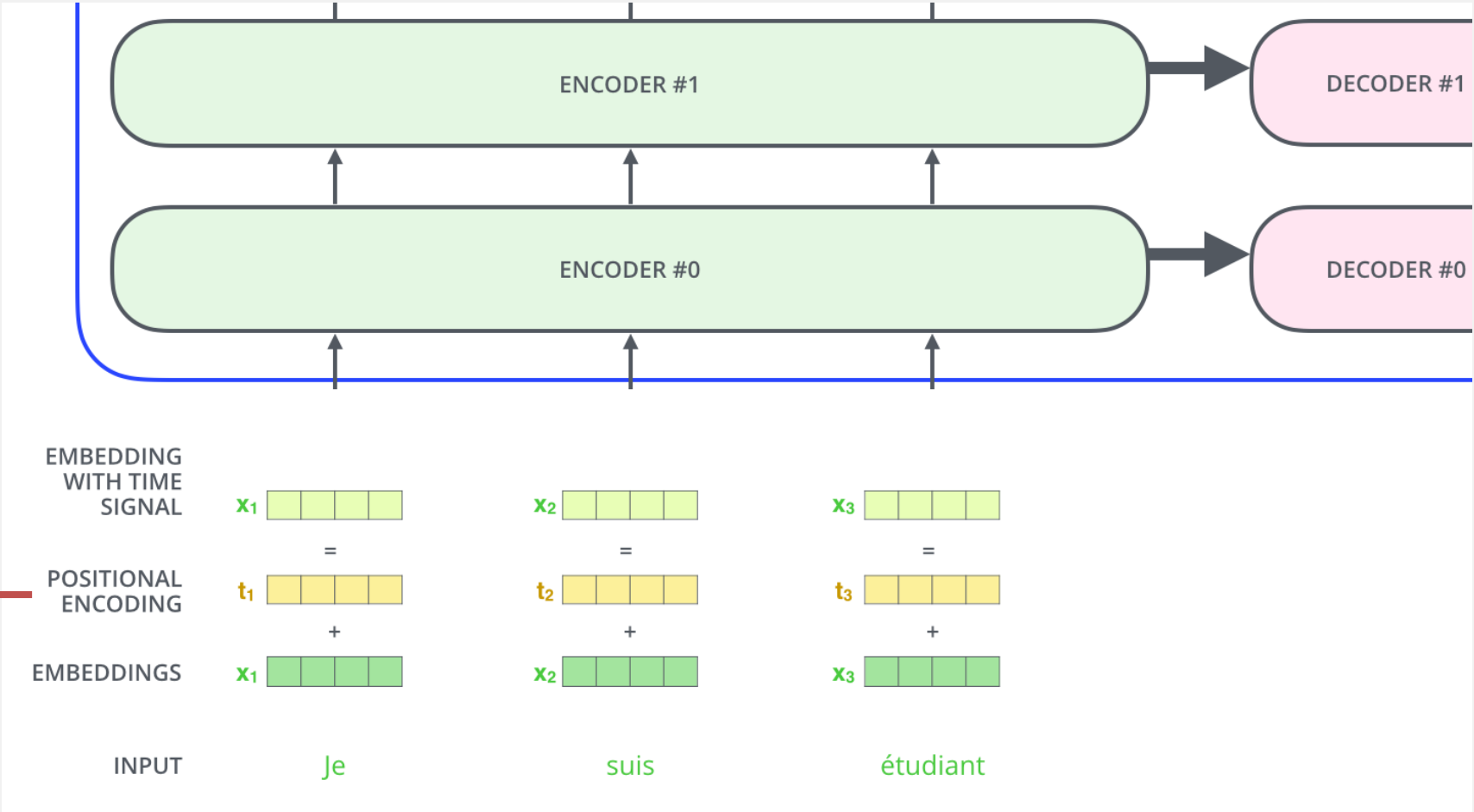
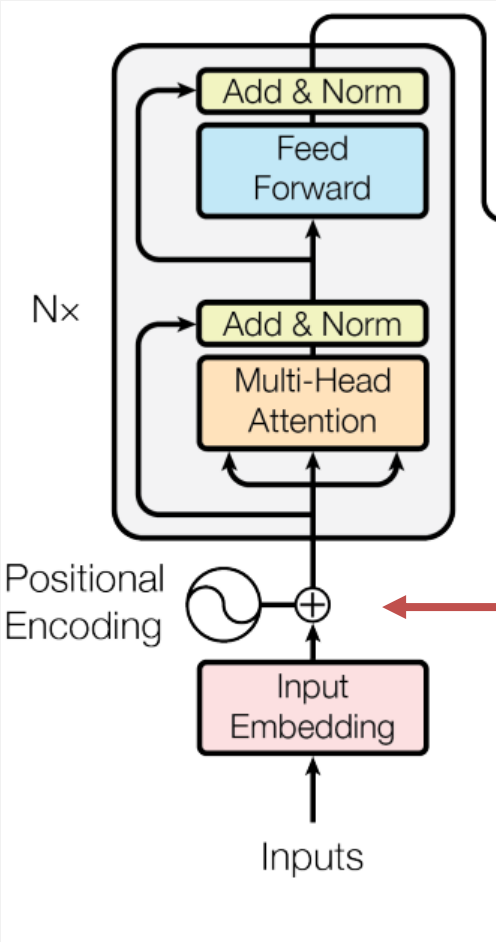
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



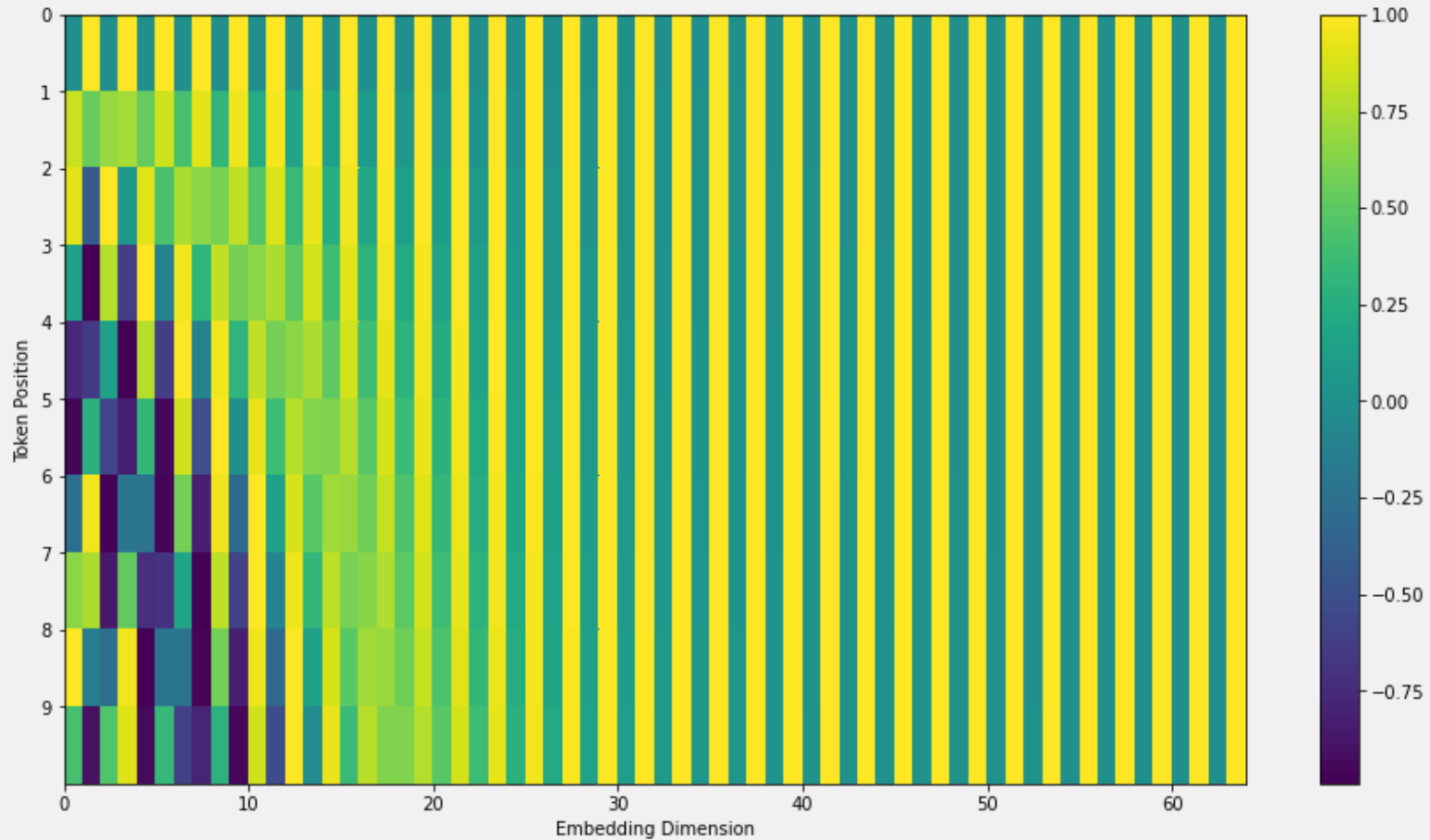
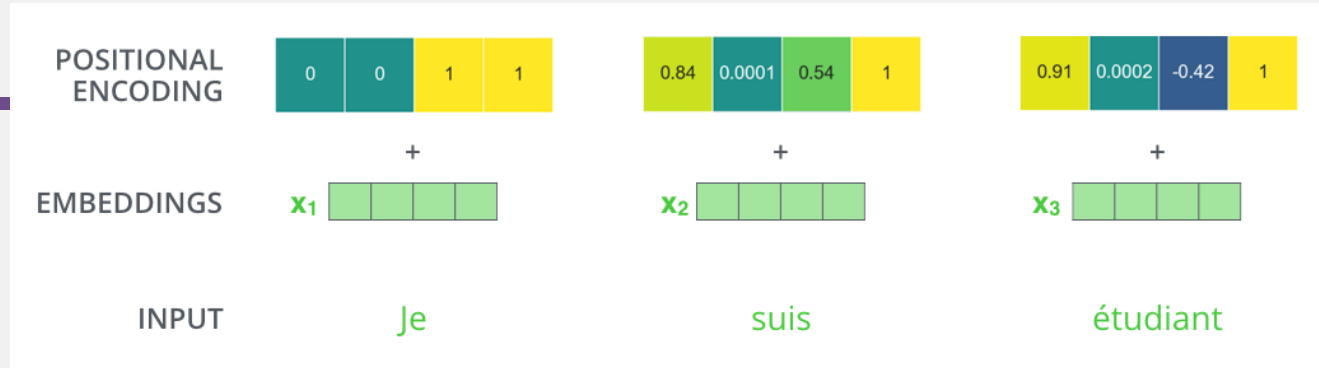
Transformer



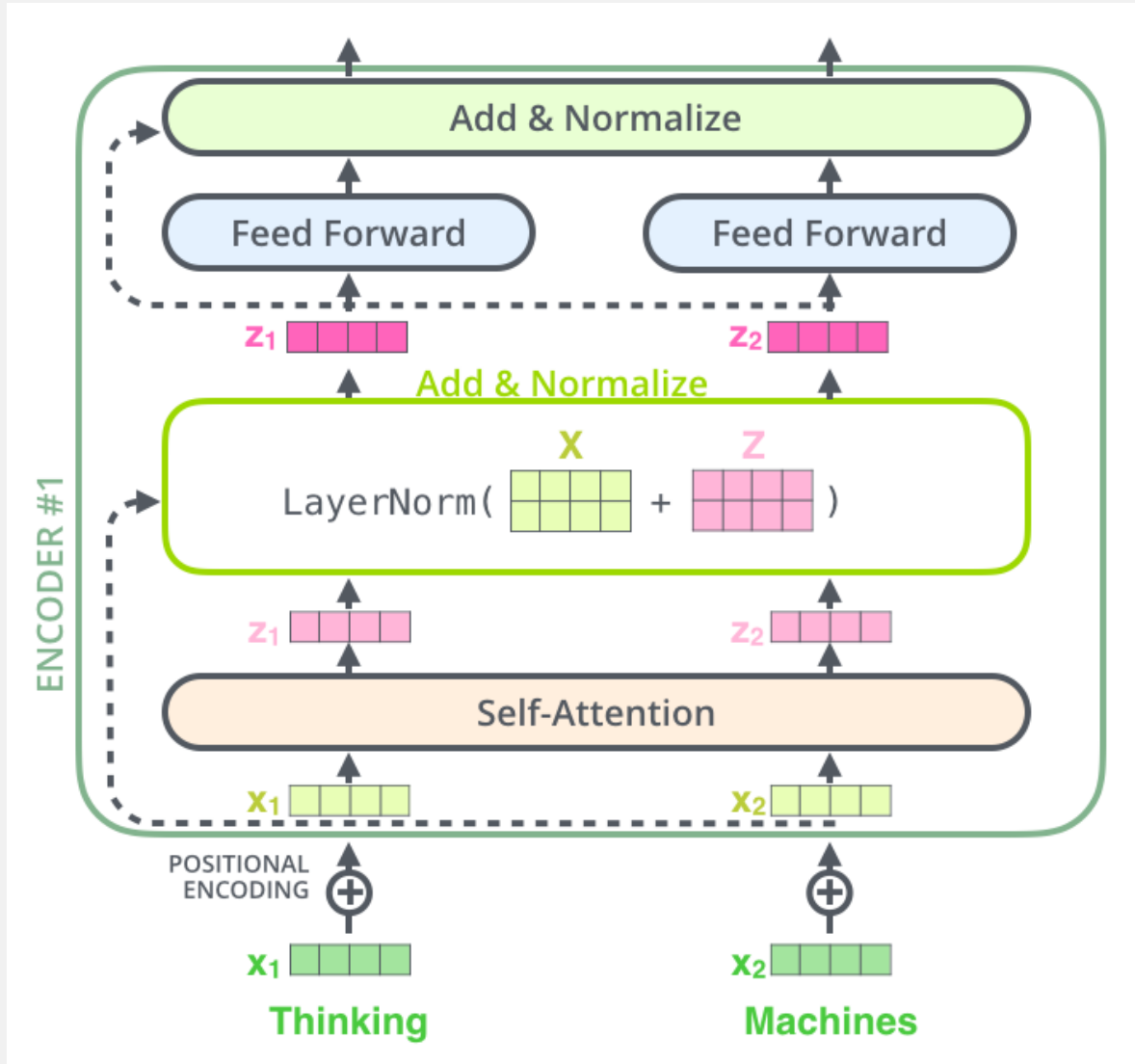
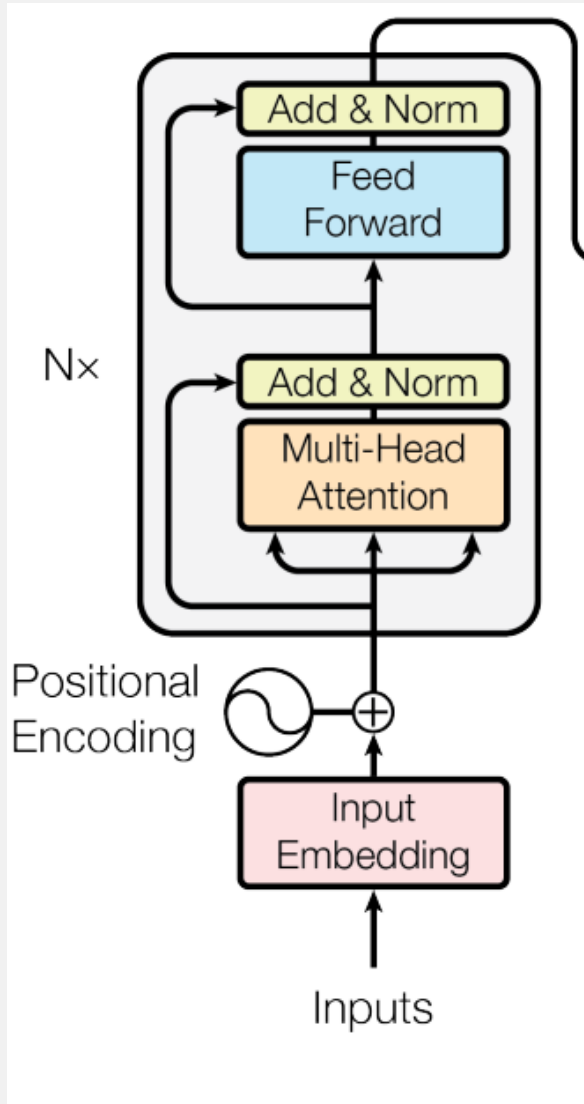
Transformer



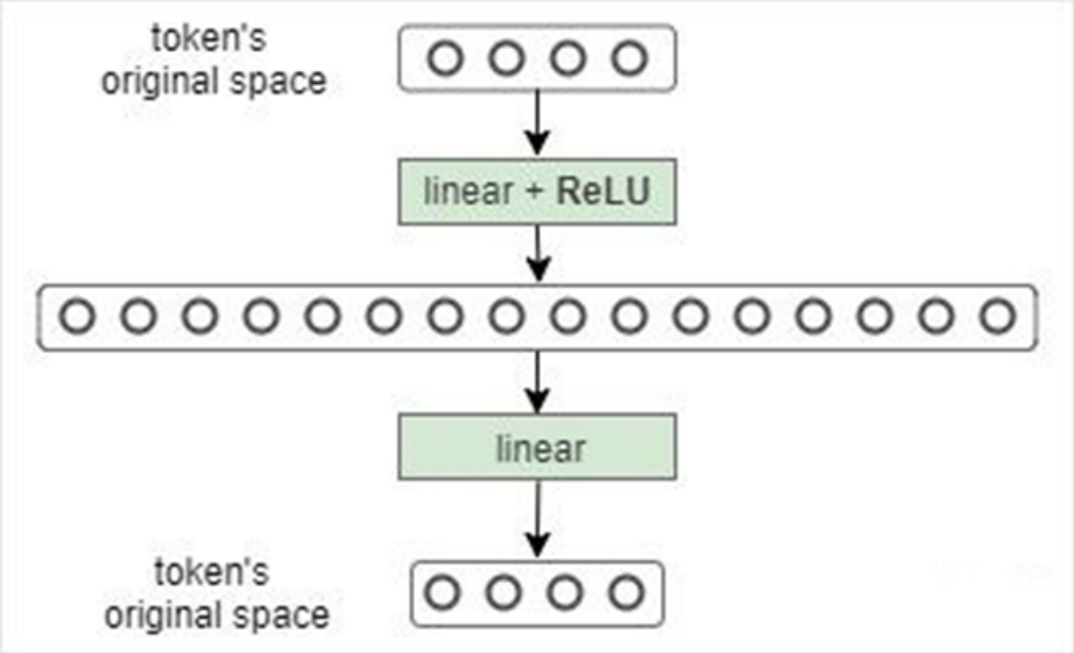
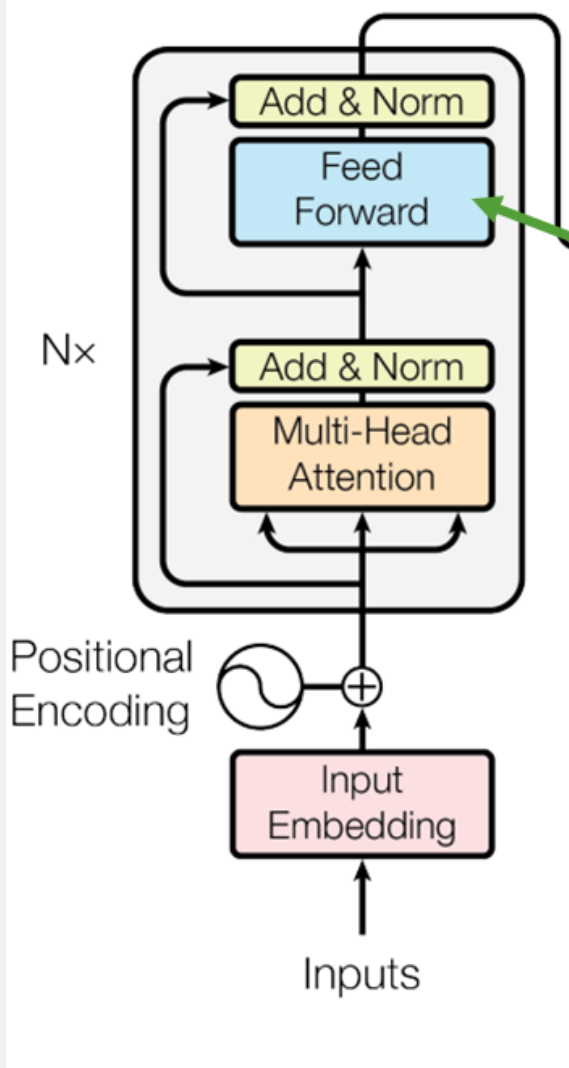
Transformer



Transformer



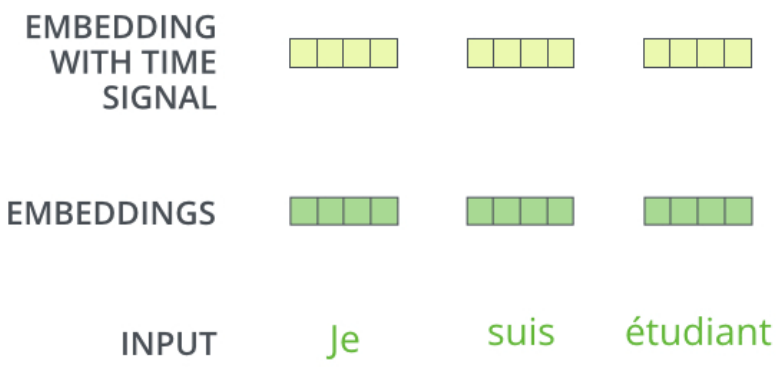
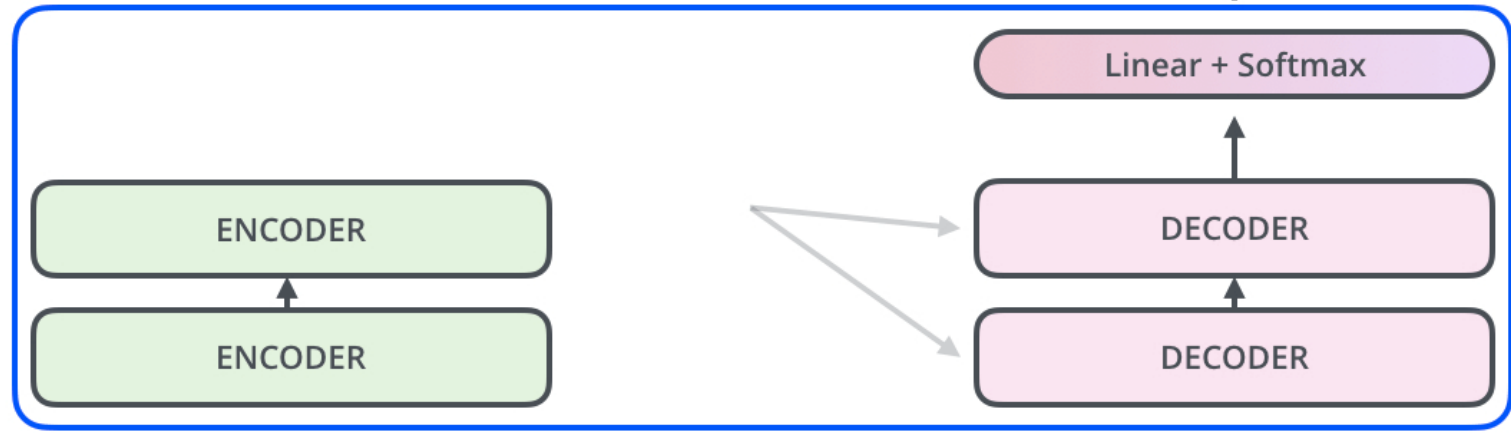
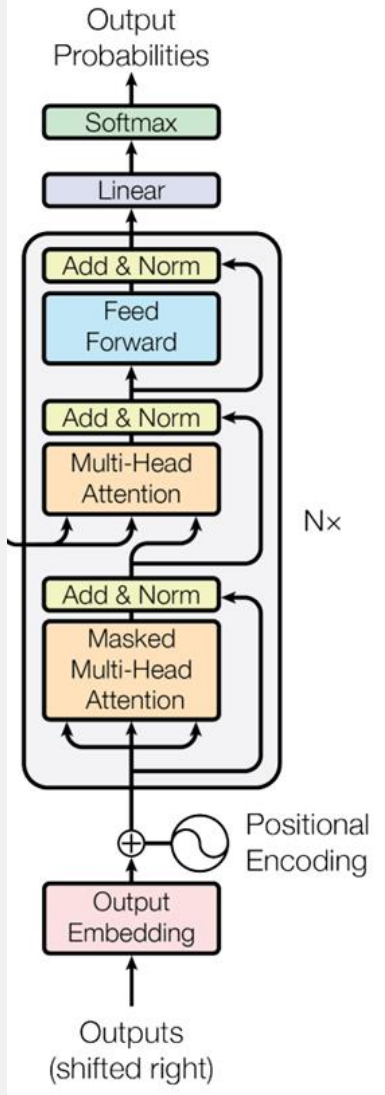
Transformer



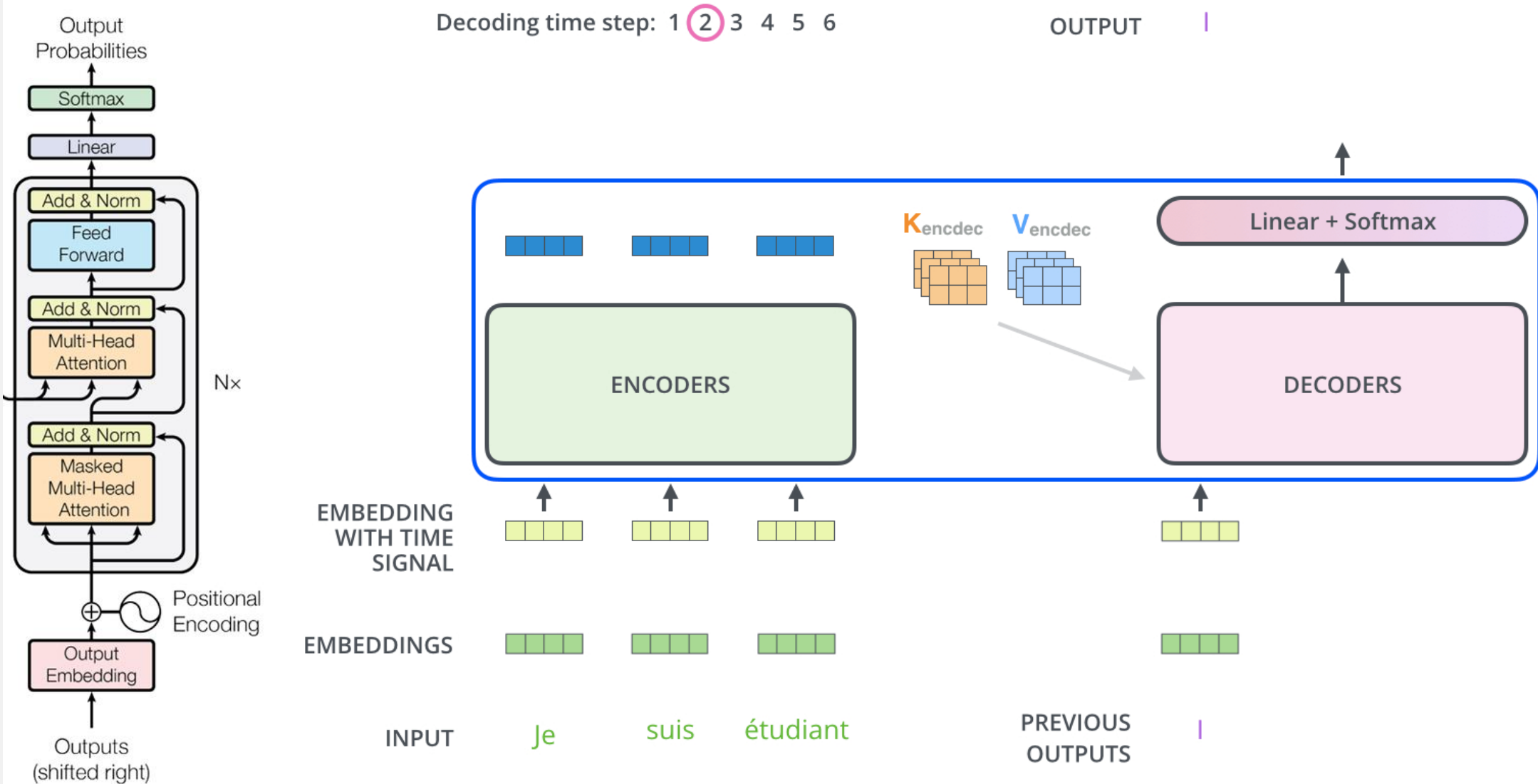
Transformer

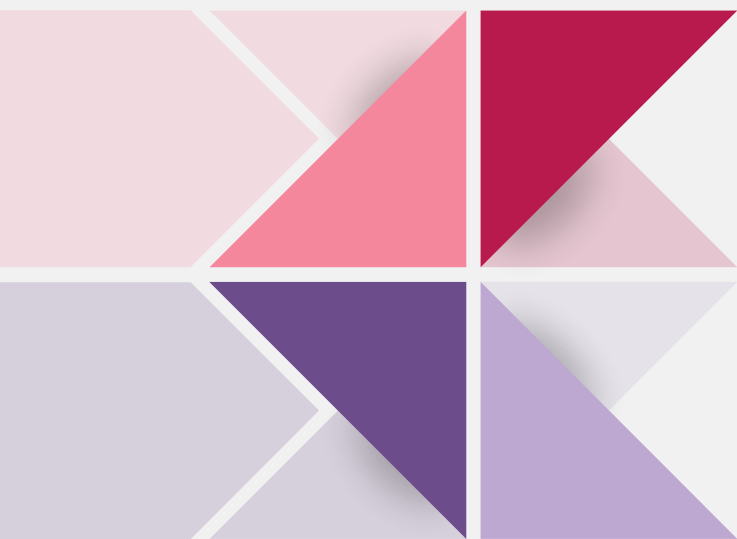
Decoding time step: ① 2 3 4 5 6

OUTPUT



Transformer



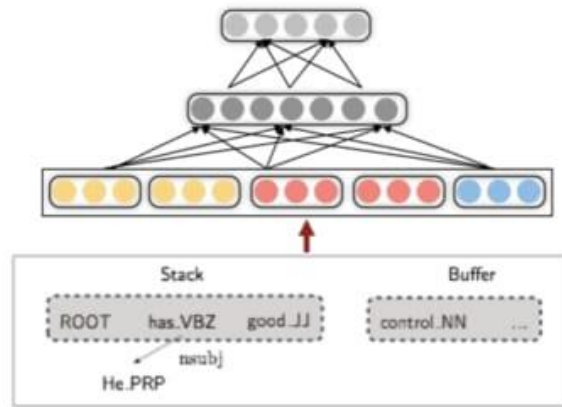


03

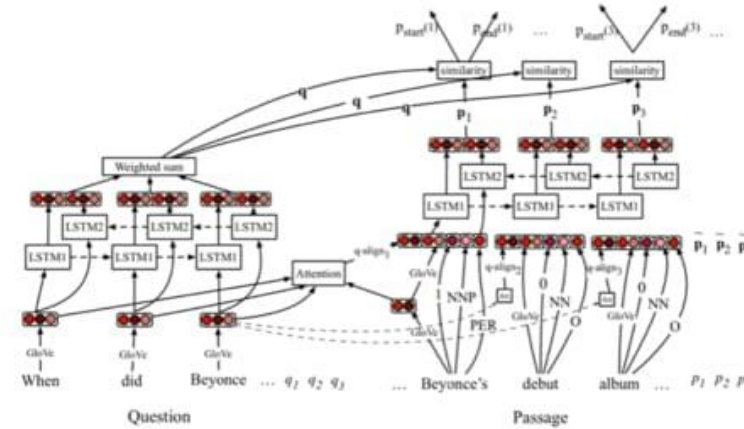
Two Models

BERT

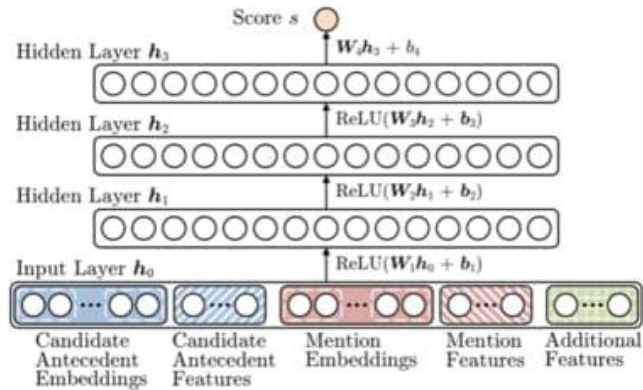
Various Model Architectures for Different NLP Tasks



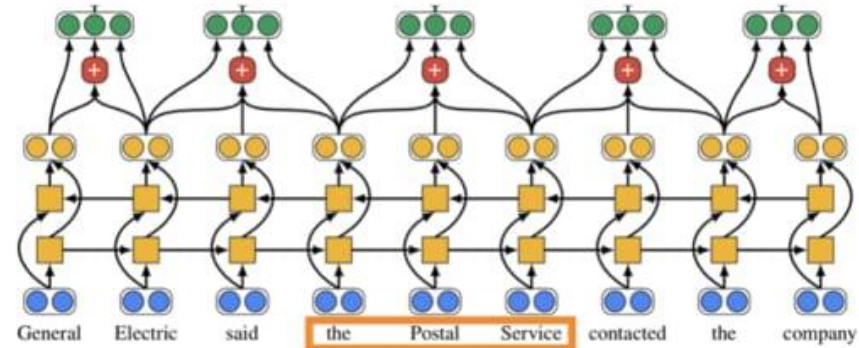
Dependency Parsing



Question Answering



Coreference example 1



Coreference example 2

BERT

Bidirectional Encoder Representation from Transformers

➤ Architecture

- Built by multiple Transformers' encoder

➤ Features

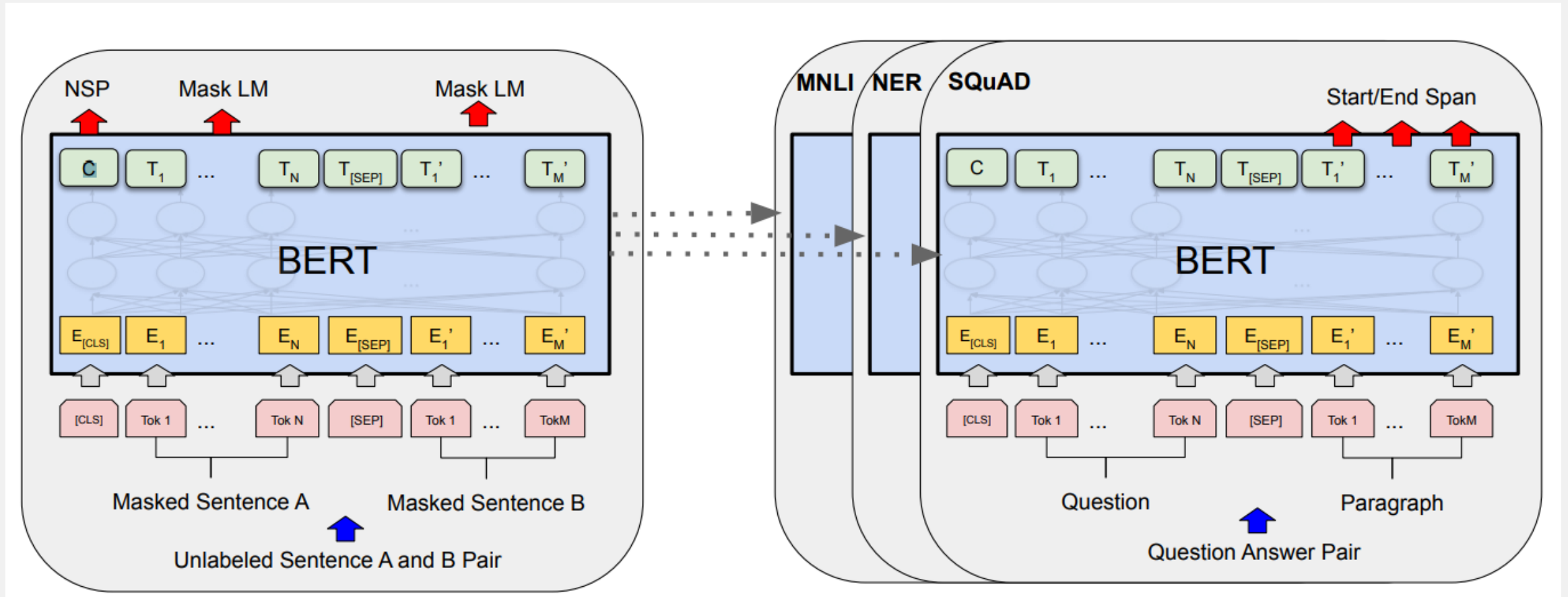
- Bidirectionality
- Pre-training
- Semantic representation

➤ Objective

- Provide a universal language representation model to be fine-tuned for various natural language processing tasks

BERT

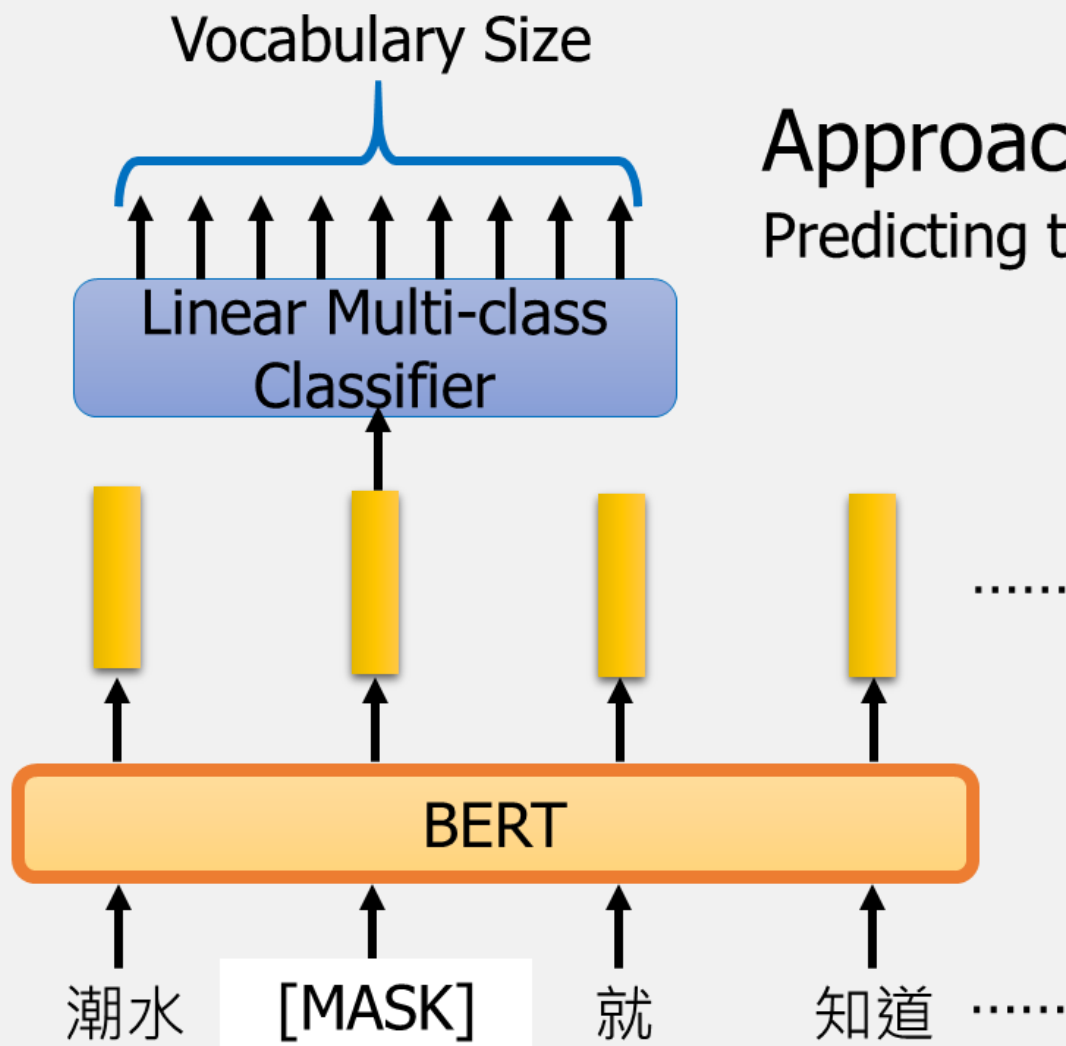
A model for all NLP task



Pre-training

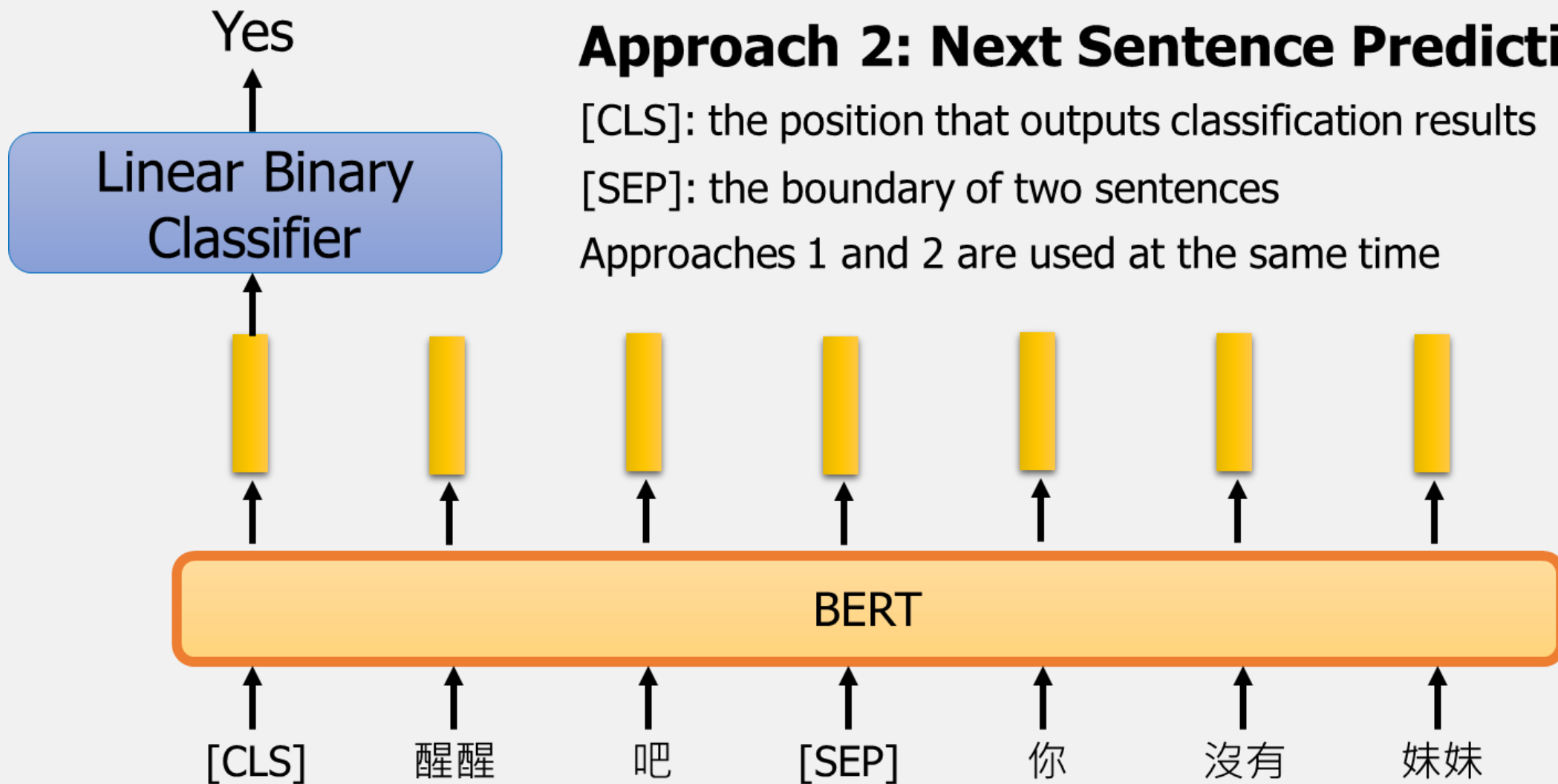
Fine-tuning

BERT



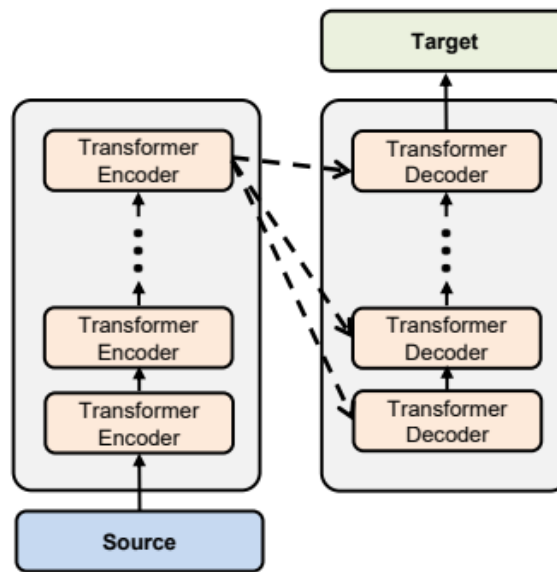
Approach 1: Masked LM
Predicting the masked word

BERT

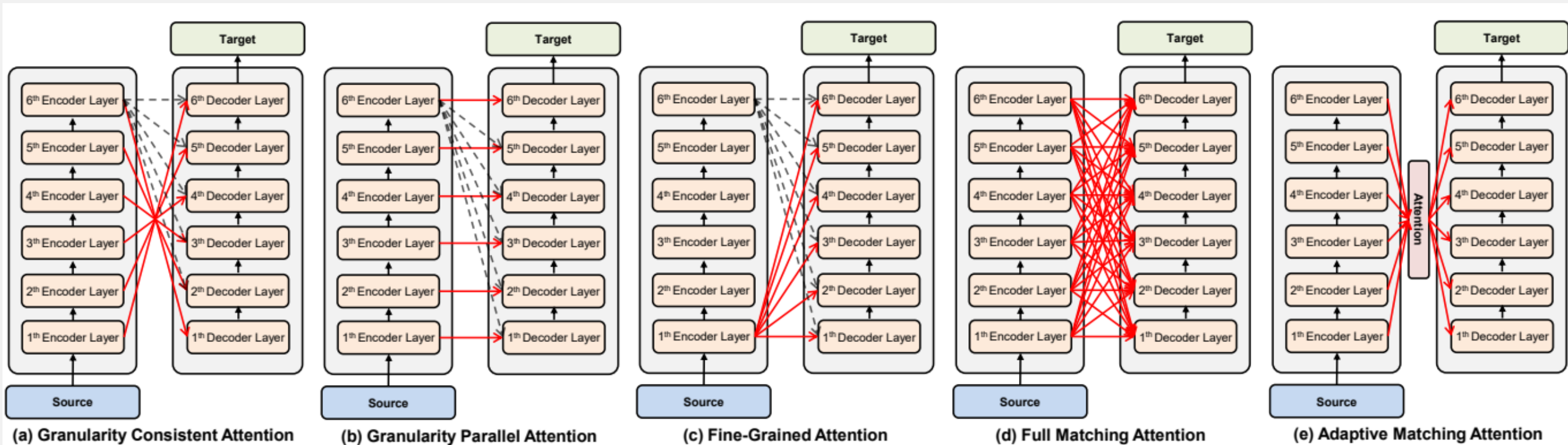


BERT

[Reference](#)



Conventional Transformer



Various Transformer

GPT

Generative Pre-trained Transformer

➤ Architecture

- Also Built by multiple Transformers' encoder

➤ Features

- Generative model
- Large-scale pre-training
- Transfer learning

➤ Objective

- Provide a universal language model to understand and generate for various text generation and dialogue capabilities

GPT vs BERT

Comparisons

➤ Pre-training objectives

- GPT - using "predicting the next word" to learn the language structure and contextual relationships within sentences
- BERT - using "masked language model" and "next sentence prediction" to predict masked words and to judge whether two sentences are consecutive

➤ Input

- GPT - use a unidirectional Transformer architecture and take the text sequence as input
- BERT - employ a bidirectional Transformer architecture and divide the sentence into left and right parts as input

➤ Fine-tuning

- Both models use a similar method of adjusting the last layer of the pre-trained model to adapt to specific tasks