

ADVANCED DEEP LEARNING MODELS FOR NATURAL LANGUAGE PROCESSING

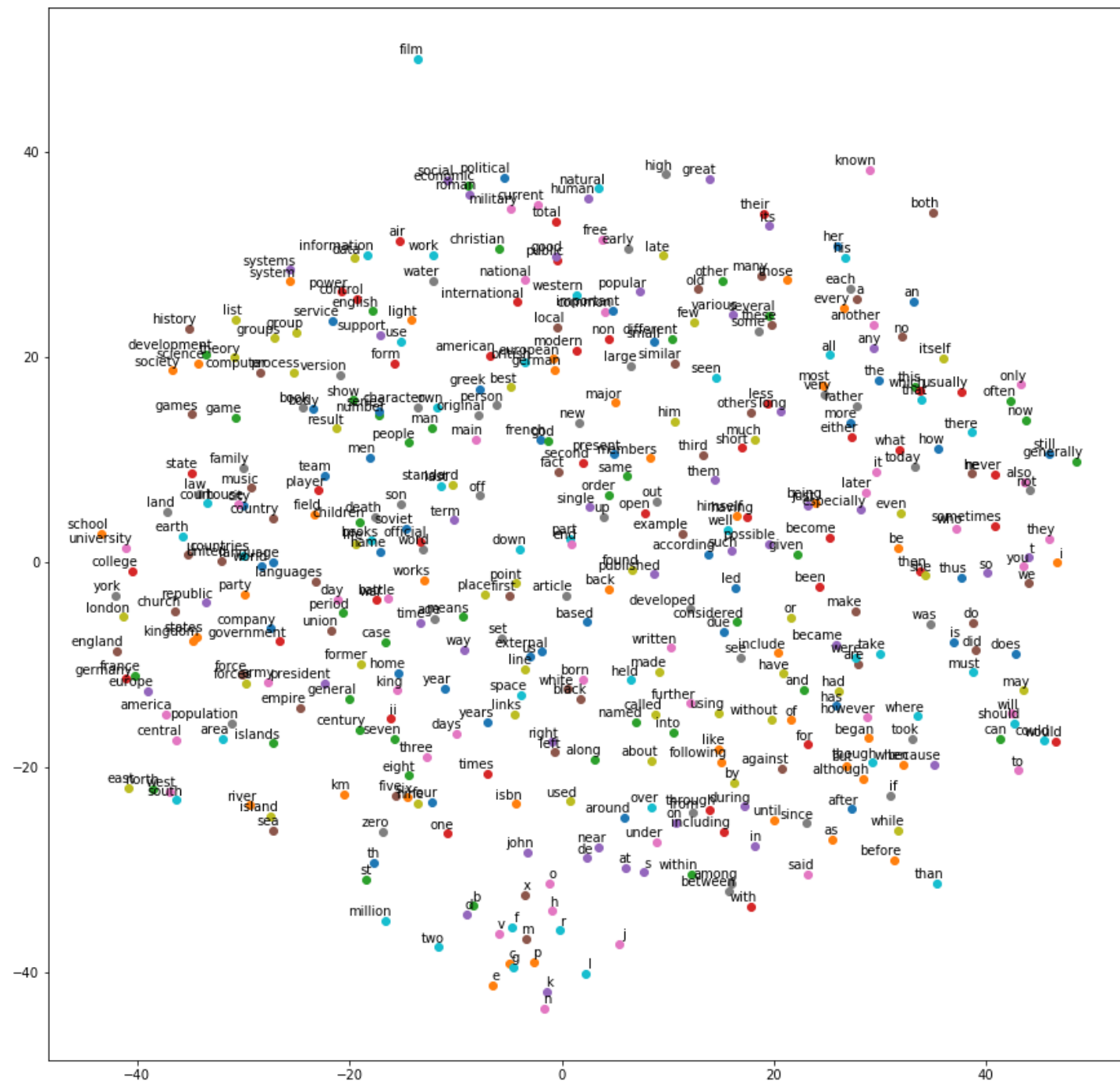
自然語言處理的進階深度模型

張家瑋 博士

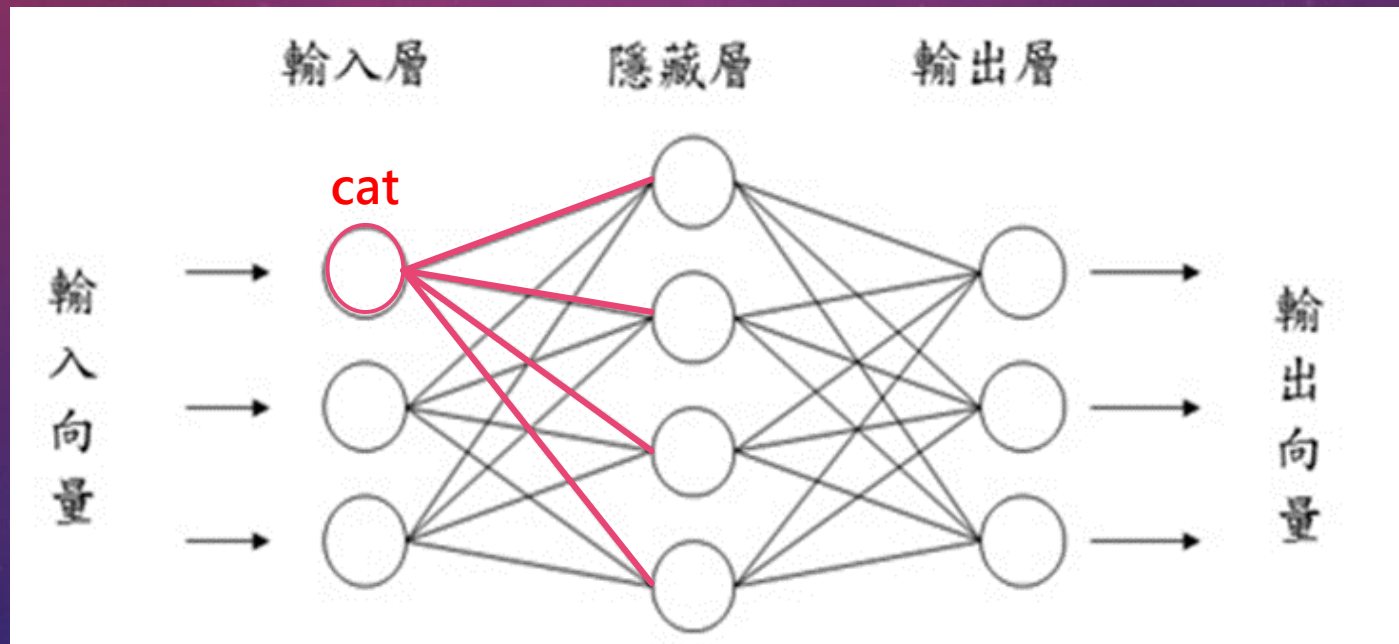
國立臺中科技大學資訊工程系助理教授

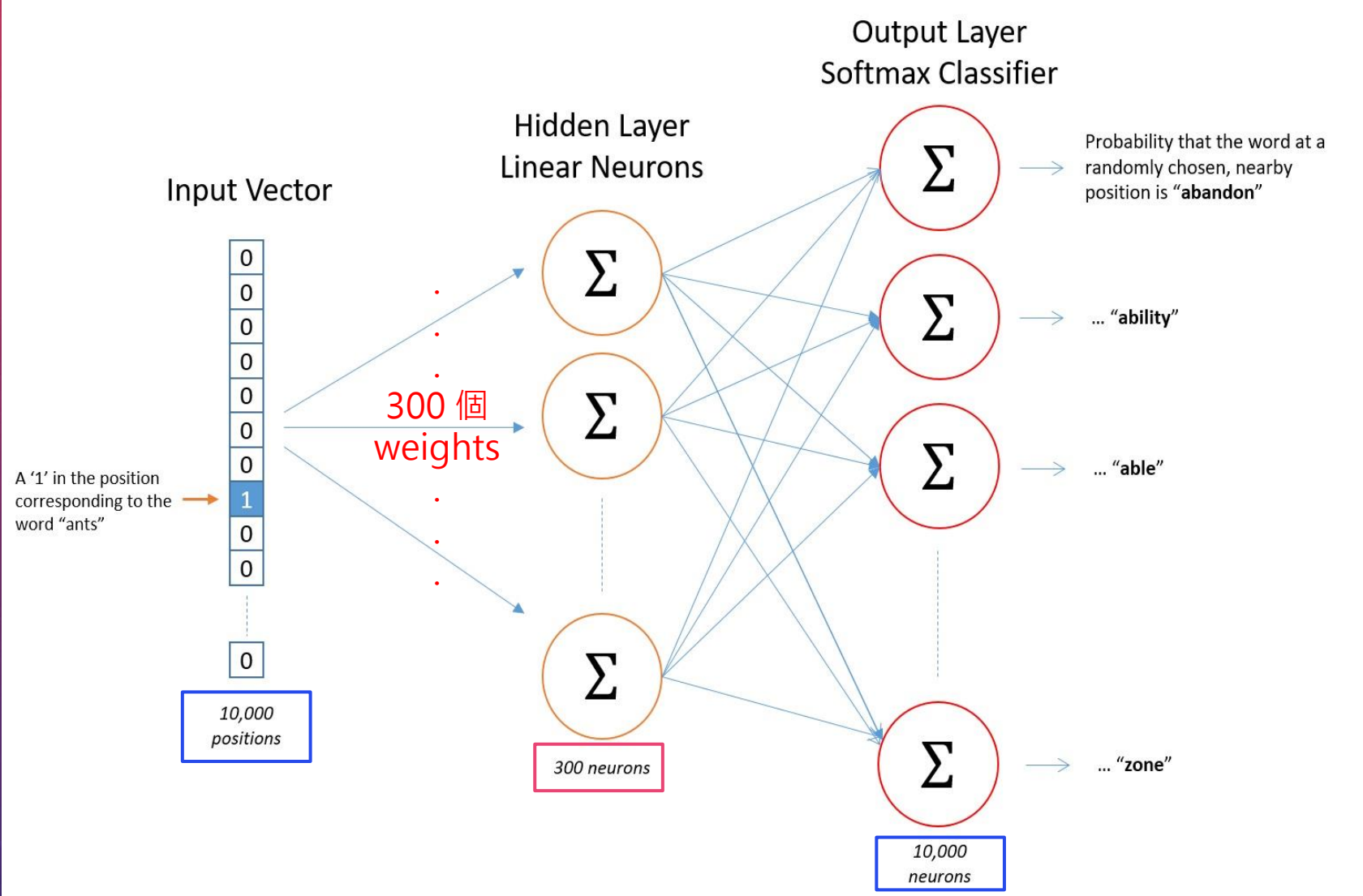


REVIEW



The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
cat -> [0. 0. 1. 0. 0. 0. 0. 0. 0. 0.]
jump -> [0. 0. 0. 1. 0. 0. 0. 0. 0. 0.]
over -> [0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
the -> [0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 0. 1. 0. 0. 0.]
The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 0. 1. 0. 0. 0.]
ate -> [0. 0. 0. 0. 0. 0. 0. 1. 0. 0.]
my -> [0. 0. 0. 0. 0. 0. 0. 0. 1. 0.]
homework -> [0. 0. 0. 0. 0. 0. 0. 0. 0. 1.]

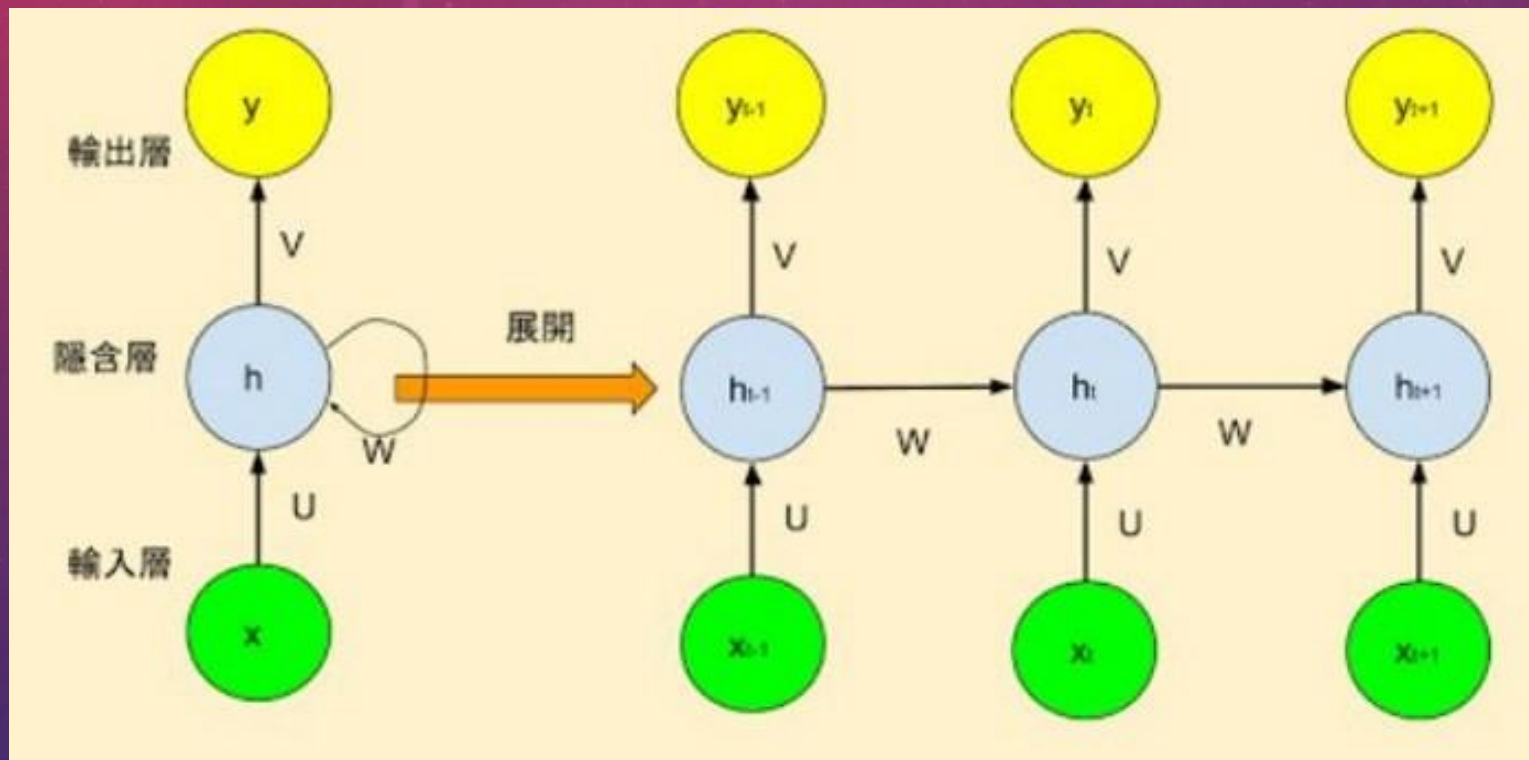




RECURRENT NEURAL NETWORK, RNN



RNN STRUCTURE



*tanh is usually used to be the activation function in RNN

PROBLEMS OF RNN

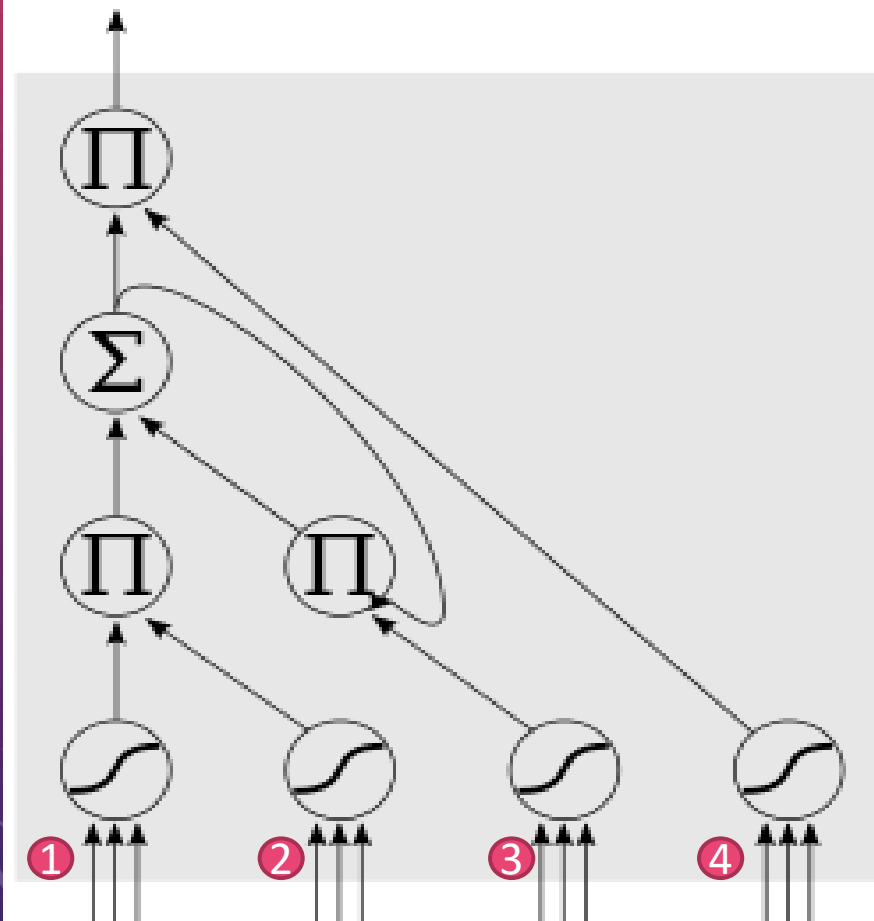
- RNN 對 short-term 敏感
- RNN 對 long-term 容易遺忘
- 越前面的字越會遺忘... 這就是RNN的梯度消失...

$$\begin{aligned}h_t &= W * h_{t-1} + U * x_t + b \\h_t &= (W^2 * h_{t-2} + W * U * x_{t-1} + W * b) + U * x_t + b \\&\dots \\h_t &= (W^t * h_0 + \dots) + U * x_t + b\end{aligned}$$

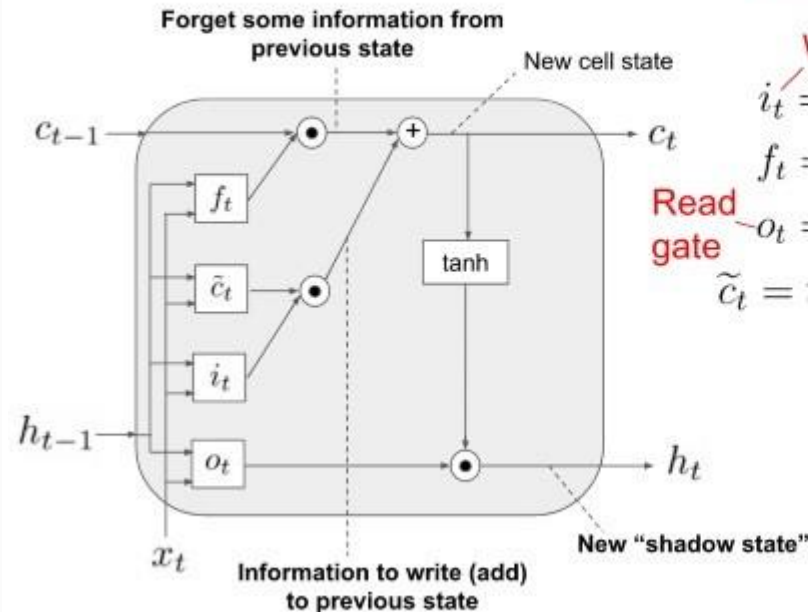
公式 h_0 對 h_t 的影響力為 w 的 t 次方，通常， w 會小於 1，離目標字越遠的字詞經過越多次傳遞(連乘)... 影響力幾乎不見了

LONG SHORT TERM MEMORY, LSTM

LSTM STRUCTURE



LSTM Cell



Conceptually, we lose information

We use shadow state to calculate gates

LSTM equations

Write gate

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i) \quad \text{Input gate}$$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f) \quad \text{Forget gate}$$

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o) \quad \text{Output gate}$$

$$\tilde{c}_t = \tanh(W c h_{t-1} + U c x_t + b) \quad \text{Memory cell candidate}$$

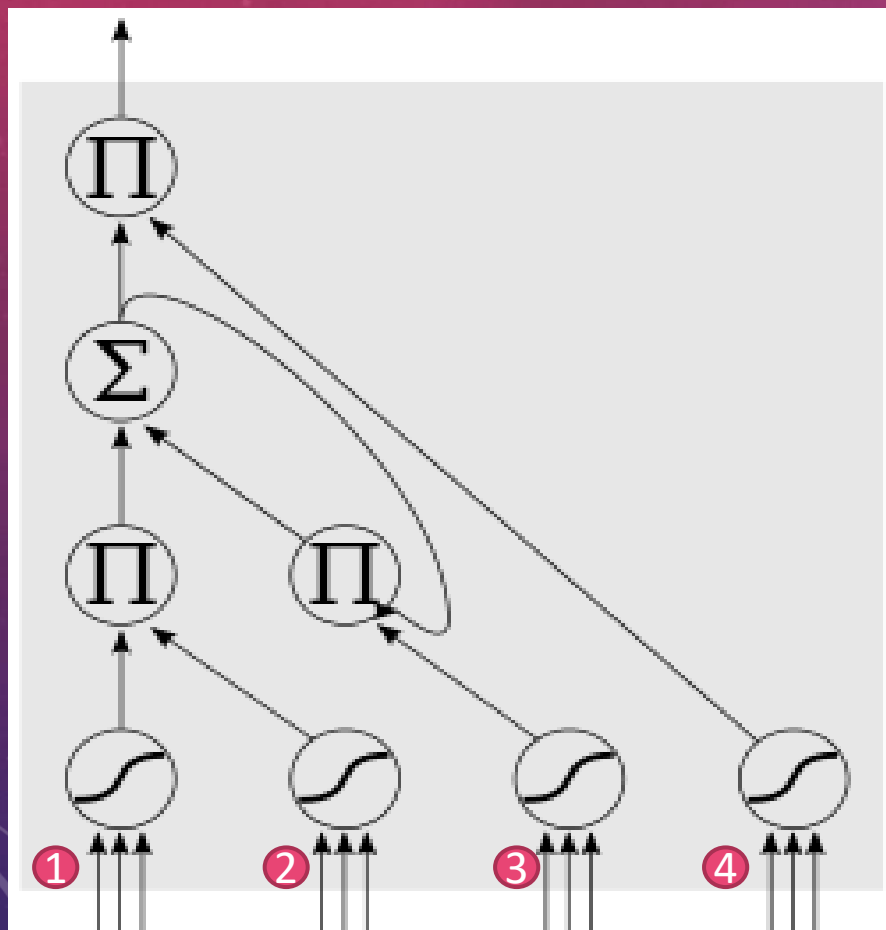
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \quad \text{Memory cell}$$

$$h_t = o_t \circ \tanh(c_t) \quad \text{Shadow state}$$

$$y_t = h_t \quad \text{Cell Output}$$

Read occurs after writing

LSTM STRUCTURE



1. 純的 input，下面三個gate會影響這個input能否被記住或作為輸出

2. Input gate，如果值為0就擋住，不給進下一層 (sigmoid)

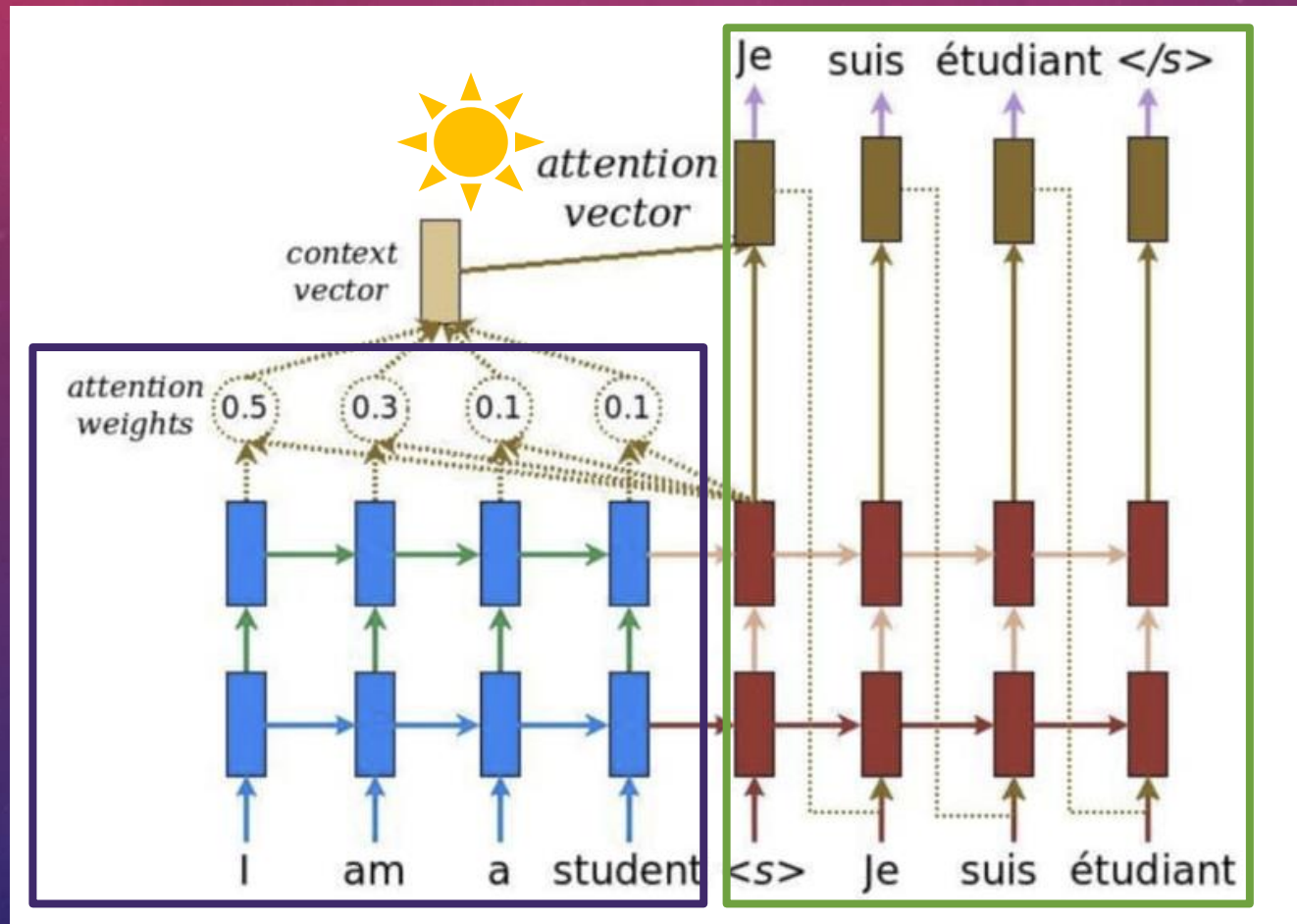
3. Forget gate，如果值近似0，就把區塊裡記住的值忘掉 (sigmoid)

4. Output gate，決定在區塊記憶中的input是否能輸出 (sigmoid)

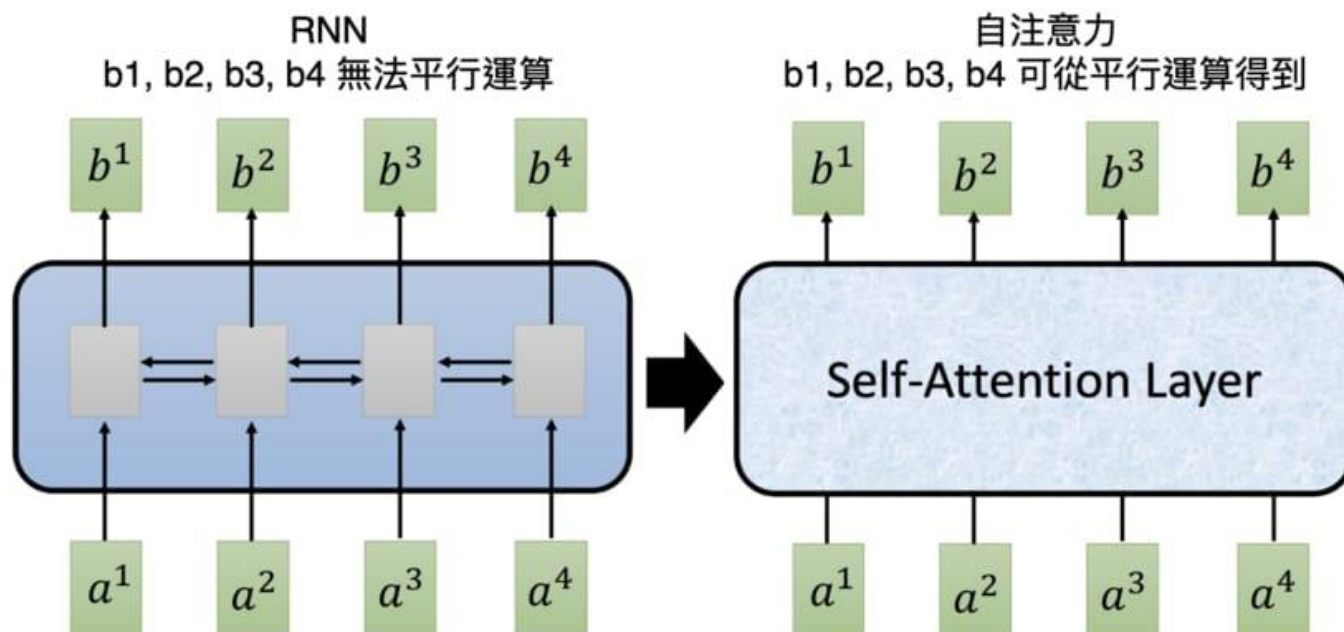
TRANSFORMER



WHAT IS ATTENTION?

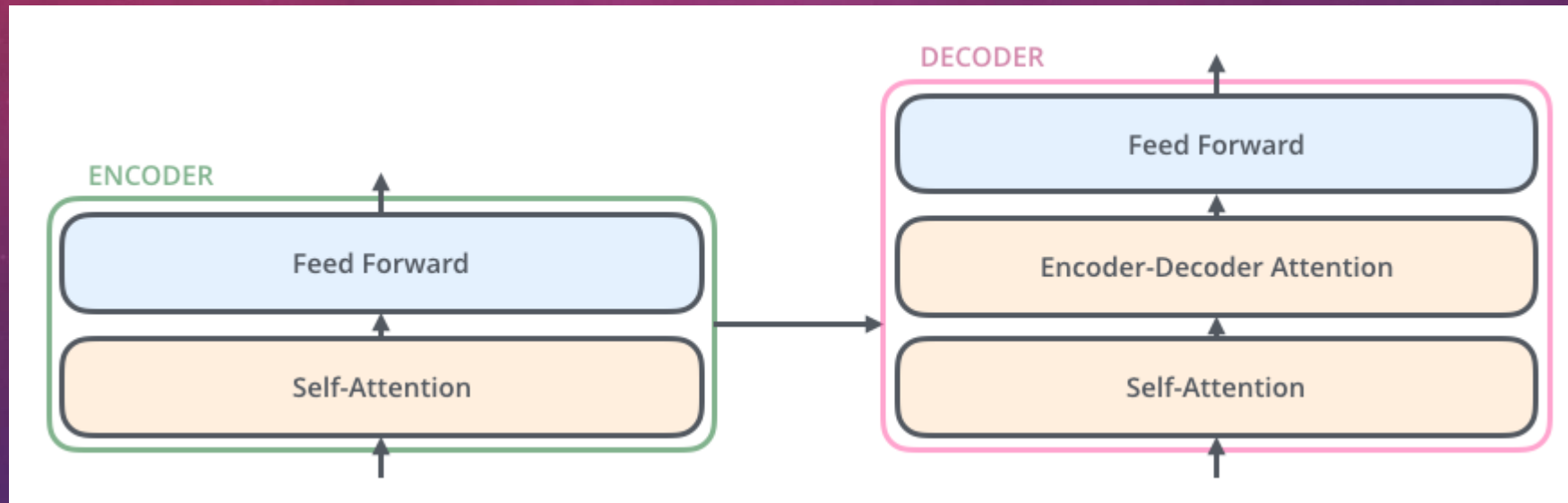


WHY

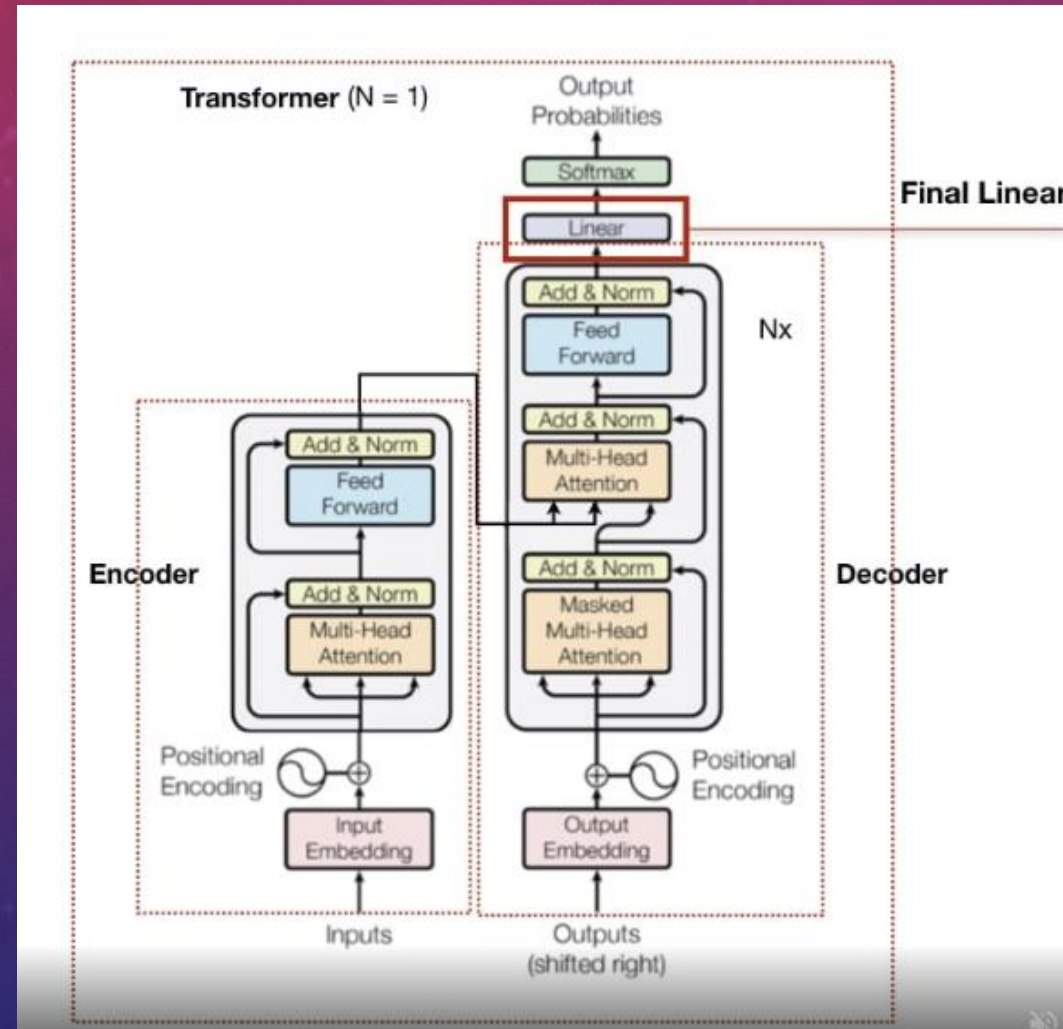


無法有效地平行運算

TRANSFORMER STRUCTURE



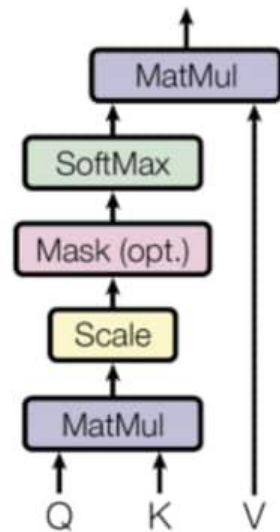
TRANSFORMER STRUCTURE



TRANSFORMER STRUCTURE

論文內示意圖

Scaled Dot-Product Attention

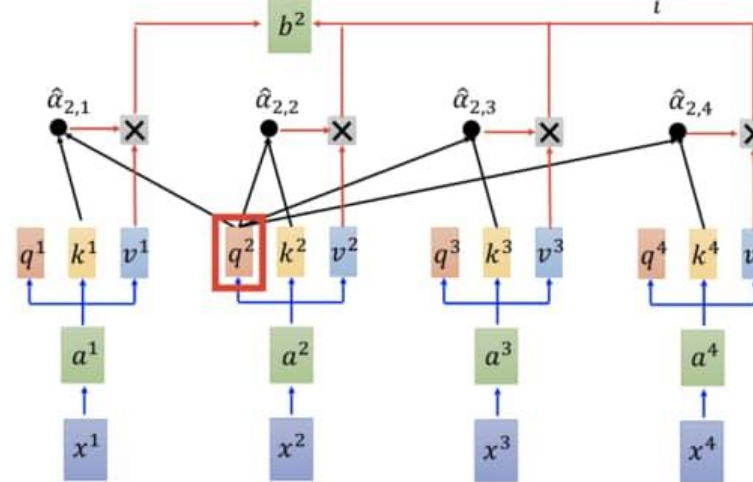


教授課程內示意圖

Self-attention

拿每個 query q 去對每個 key k 做 attention

$$b^2 = \sum_i \hat{a}_{2,i} v^i$$



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

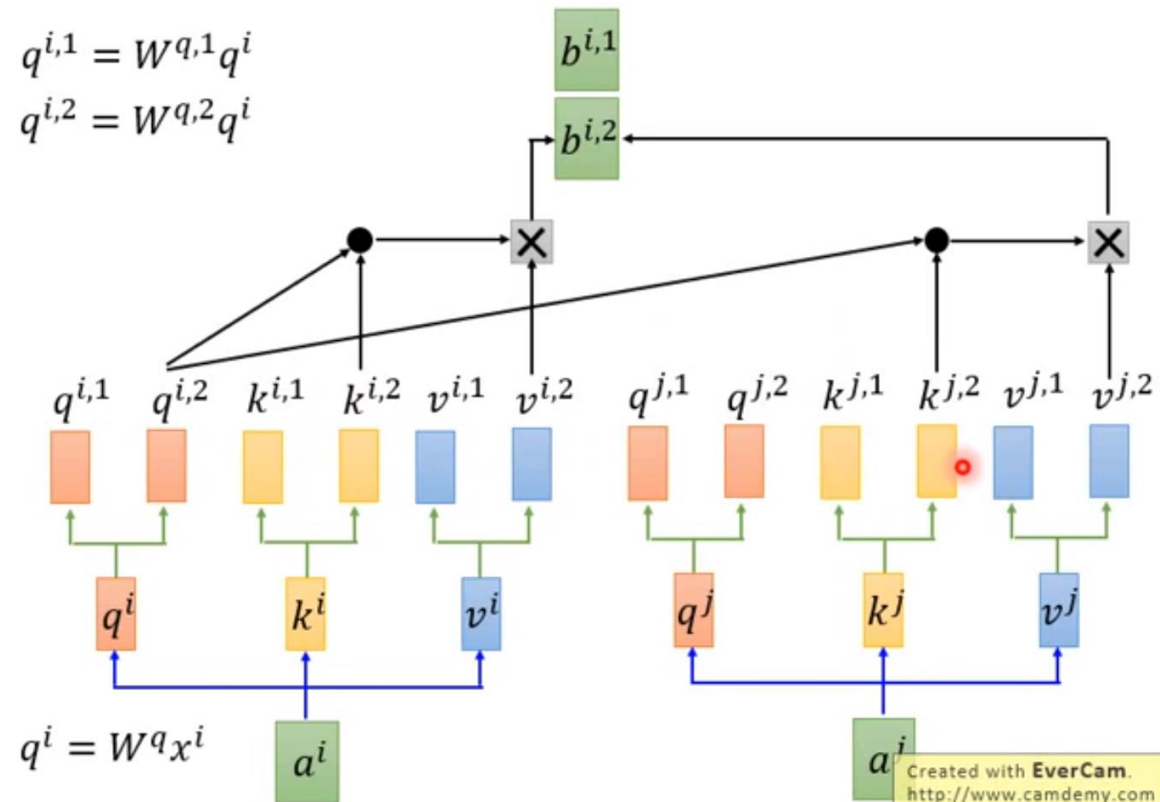
TRANSFORMER STRUCTURE

Multi-head Self-attention

(2 heads as example)

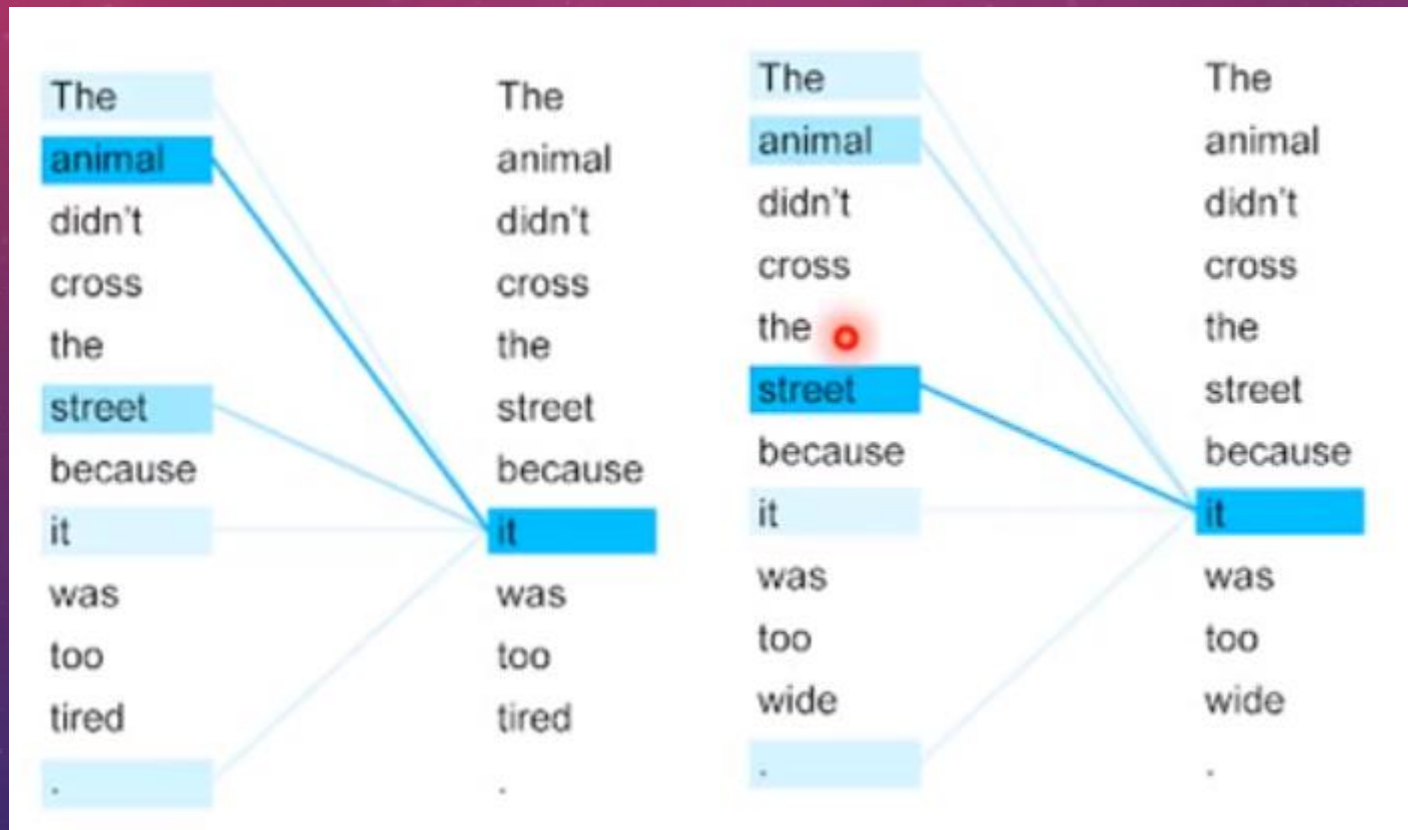
$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



Created with EverCam.
<http://www.camdemy.com>

RESULTS



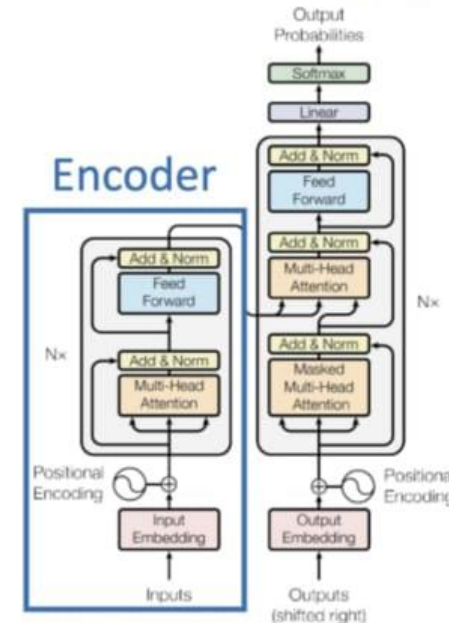
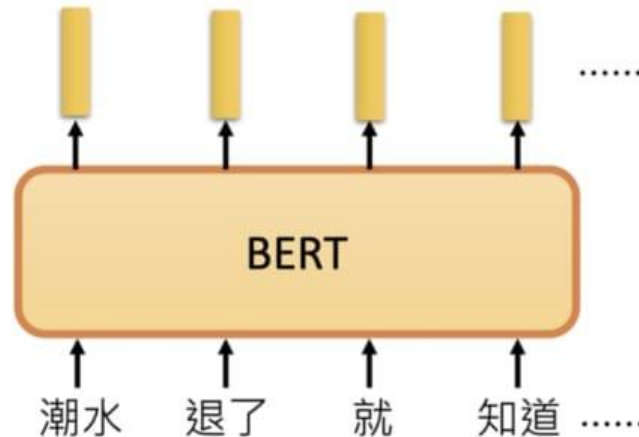


BIDIRECTIONAL ENCODER REPRESENTATIONS FROM TRANSFORMERS, BERT

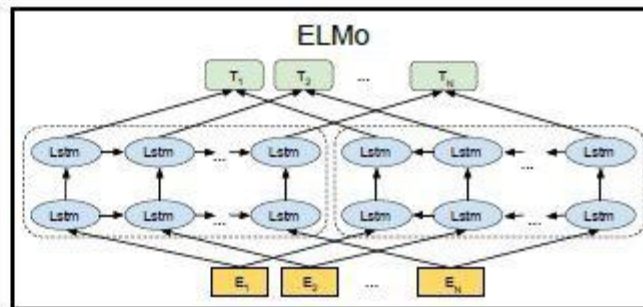
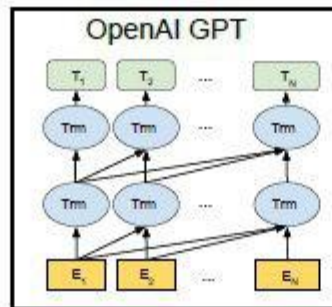
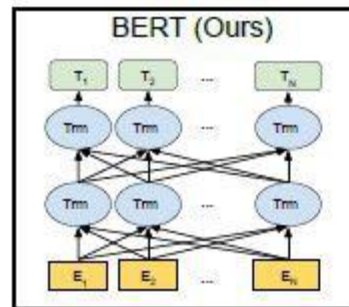
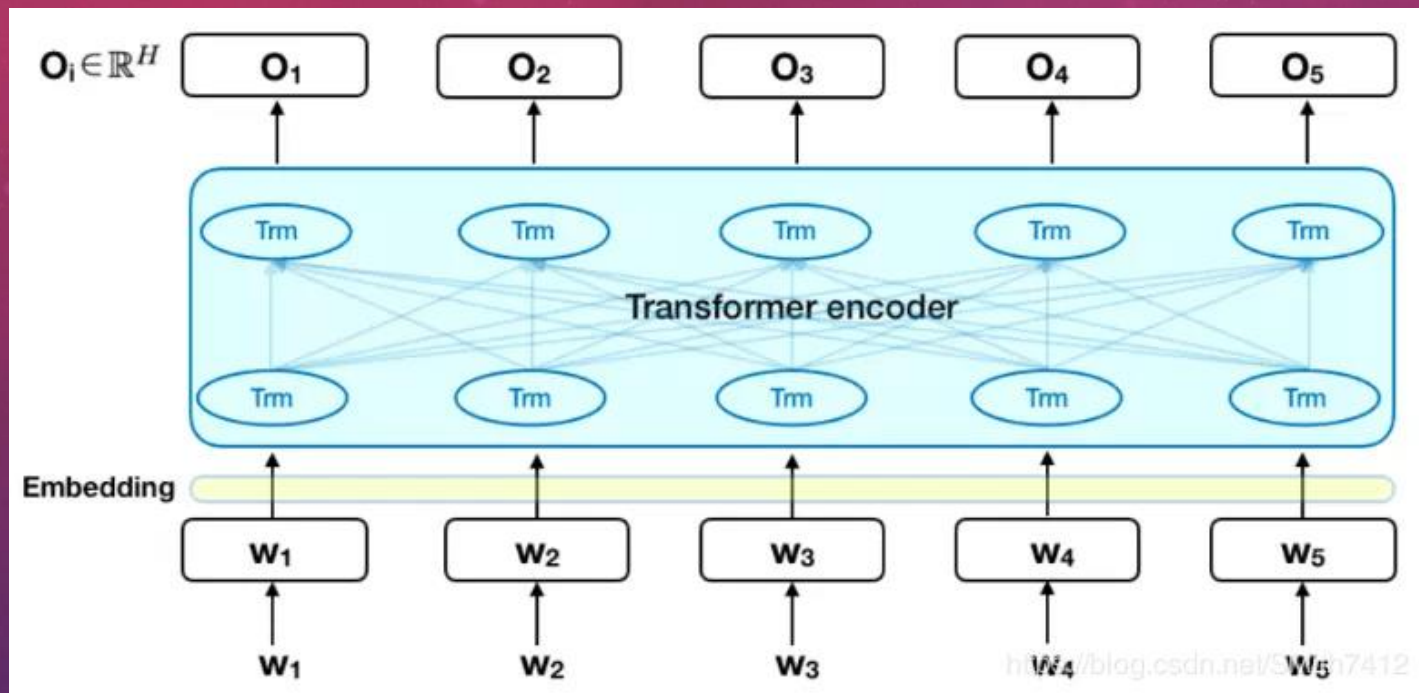
BERT= ENCODER OF TRANSFORMER

Bidirectional Encoder
Representations from Transformers
(BERT)

- BERT = Encoder of Transformer
Learned from a large amount of text
without annotation



BERT STRUCTURE

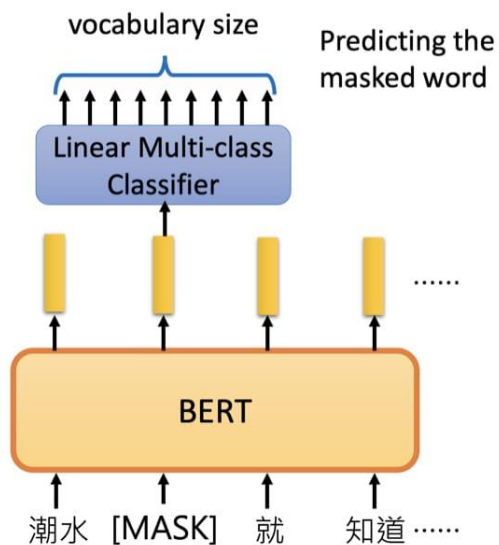


TRAINING OF BERT

預訓練任務 1：克漏字填空

Training of BERT

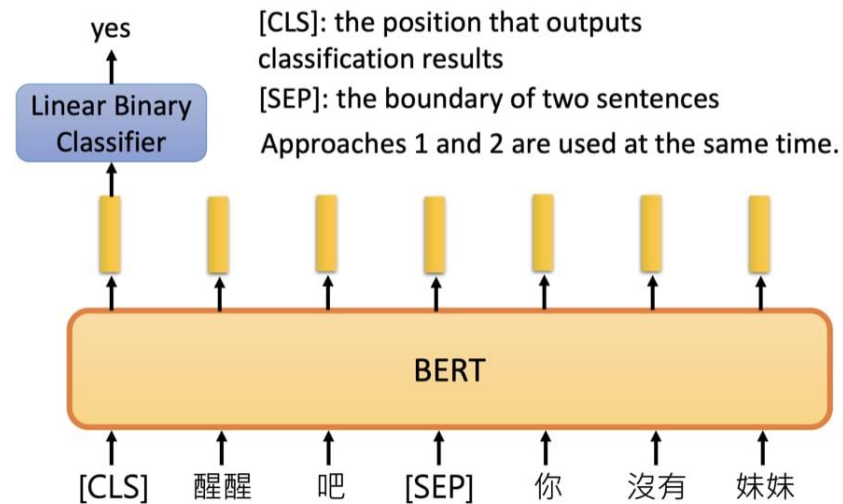
- Approach 1:
Masked LM



預訓練任務 2：下個句子預測

Training of BERT

Approach 2: Next Sentence Prediction



KEY POINTS OF BERT

BERT sentence pair encoding
(with tensors for PyTorch implementation)

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]	[PAD]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\#ing}$	$E_{[SEP]}$	
	+	+	+	+	+	+	+	+	+	+	+	
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B	
	+	+	+	+	+	+	+	+	+	+	+	
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}	
tokens_tensor	5566	1	2	3	4	9527	5	6	7	8	9527	0
segments_tensor	0	0	0	0	0	0	1	1	1	1	1	0
masks_tensor	1	1	1	1	1	1	1	1	1	1	1	0



THANK YOU