




Deep Learning

6주차

12기 이두형
12기 임효진



What is the best way for translation?

I love you = Nan nul saranghey
난 널 사랑해

Word to Word translation?

Input

Prediction

I

=>

nan (난)

love

=>

saranghey (사랑해)

you

=>

nul (널)

I love you

=>

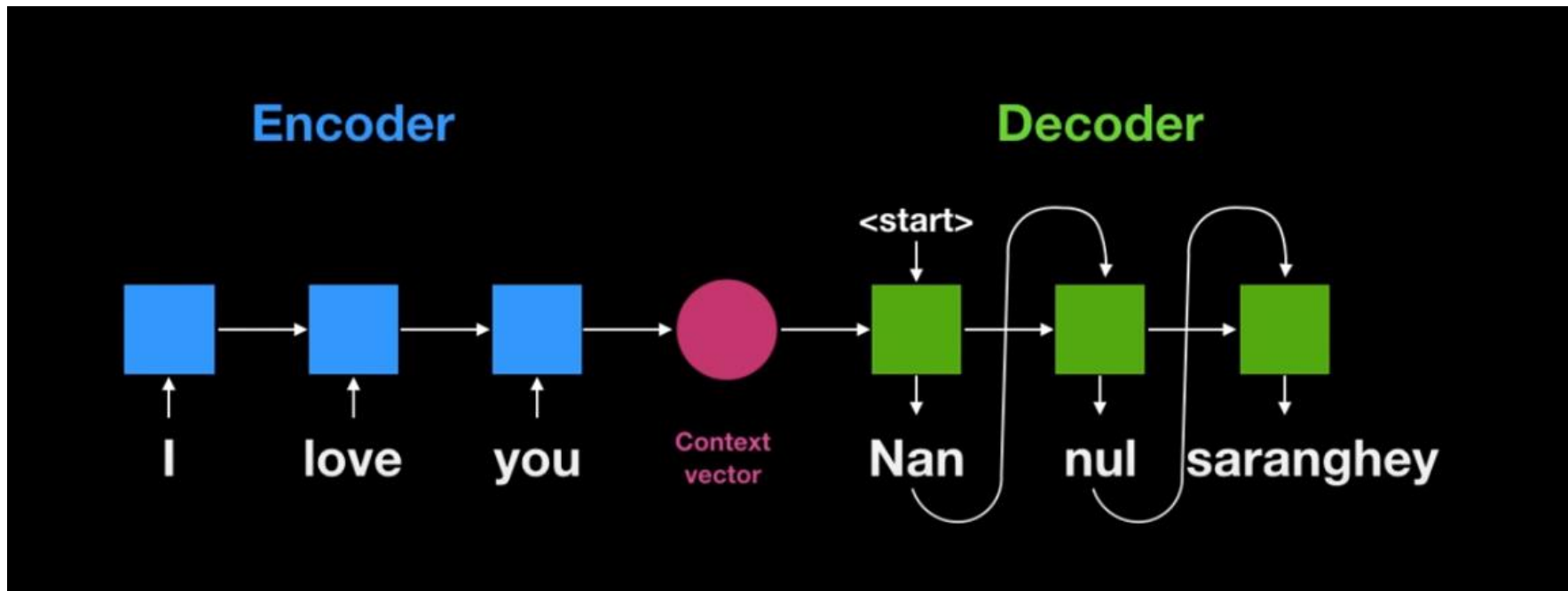
nan saranghey nul

How are you ? = Jal jiney ?

잘 지내

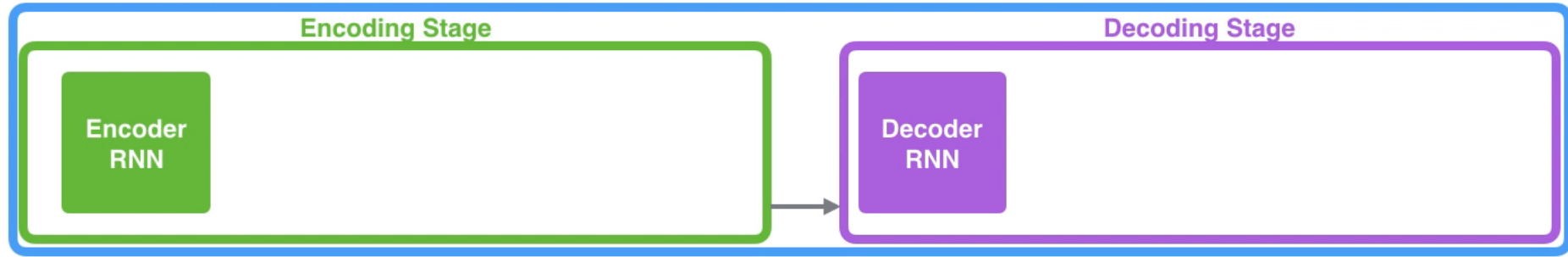
3 words

2 words



Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL



Je

suis

étudiant

Input

Je
suis
étudiant

0.901	-0.651	-0.194	-0.822
-------	--------	--------	--------



-0.351	0.123	0.435	-0.200
--------	-------	-------	--------



0.081	0.458	-0.400	0.480
-------	-------	--------	-------



CONTEXT

0.11
0.03
0.81
-0.62

0.11
0.03
0.81
-0.62

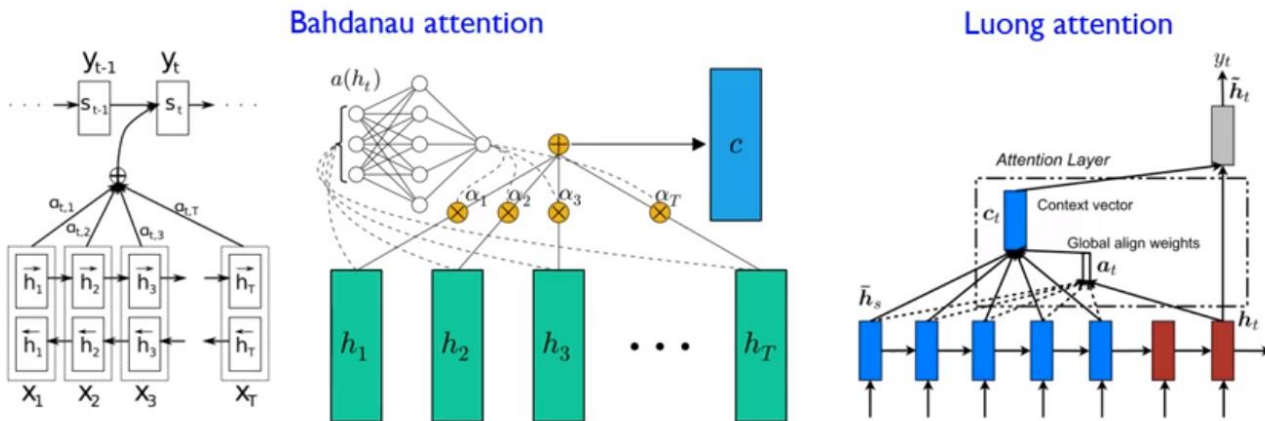
- Attention

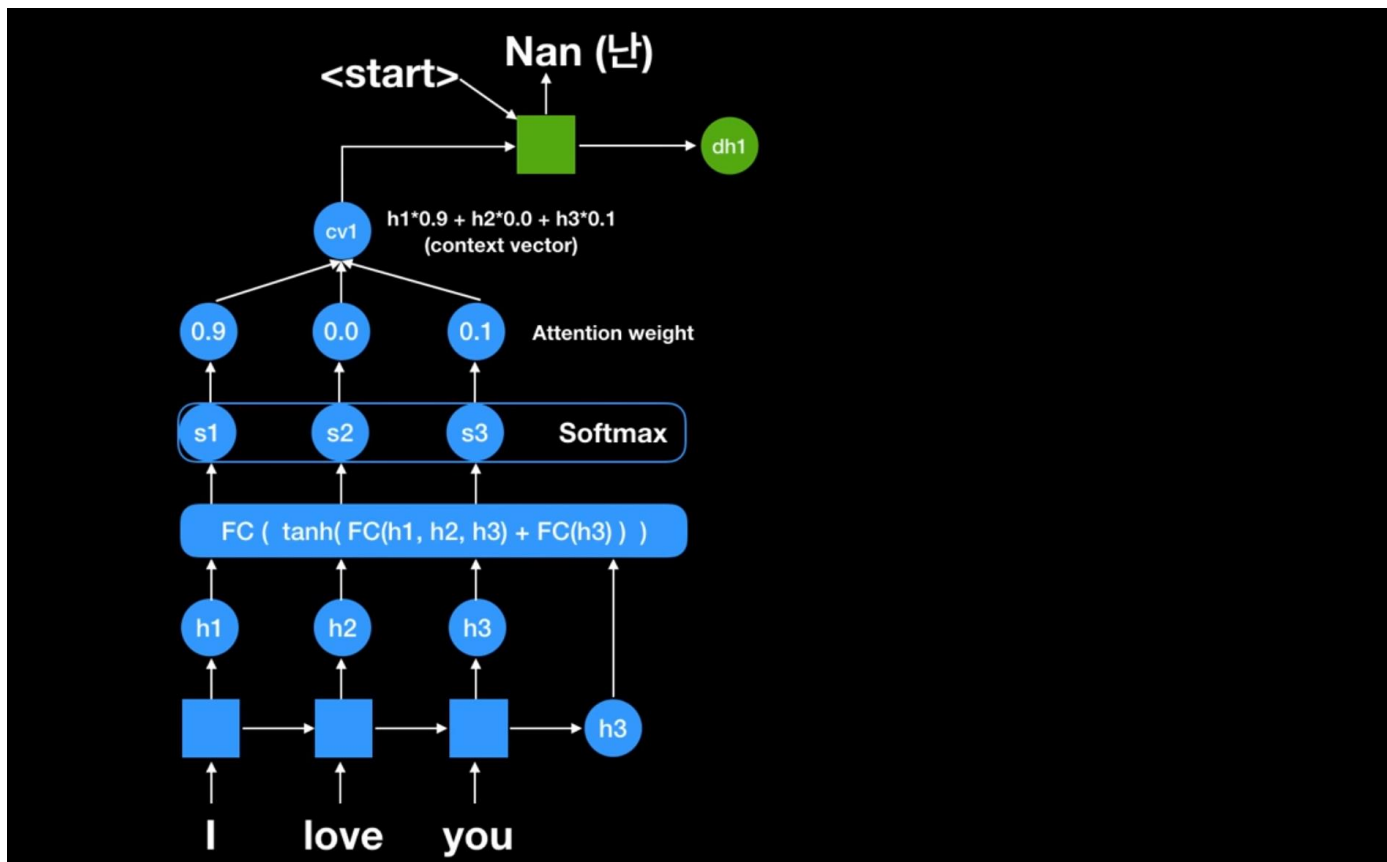
- ✓ Bahdanau attention (Bahdanau et al., 2015)

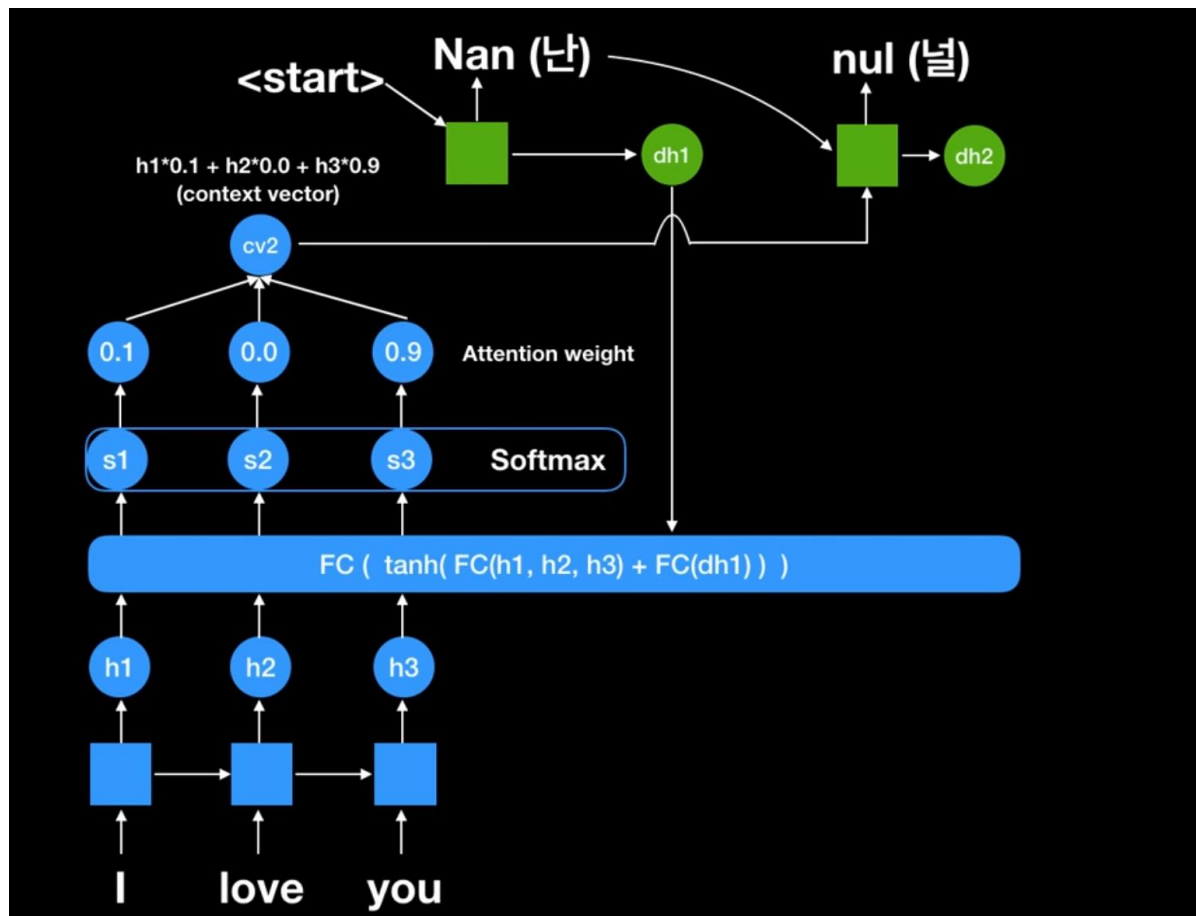
- Attention scores are separated trained, the current hidden state is a function of the context vector and the previous hidden state

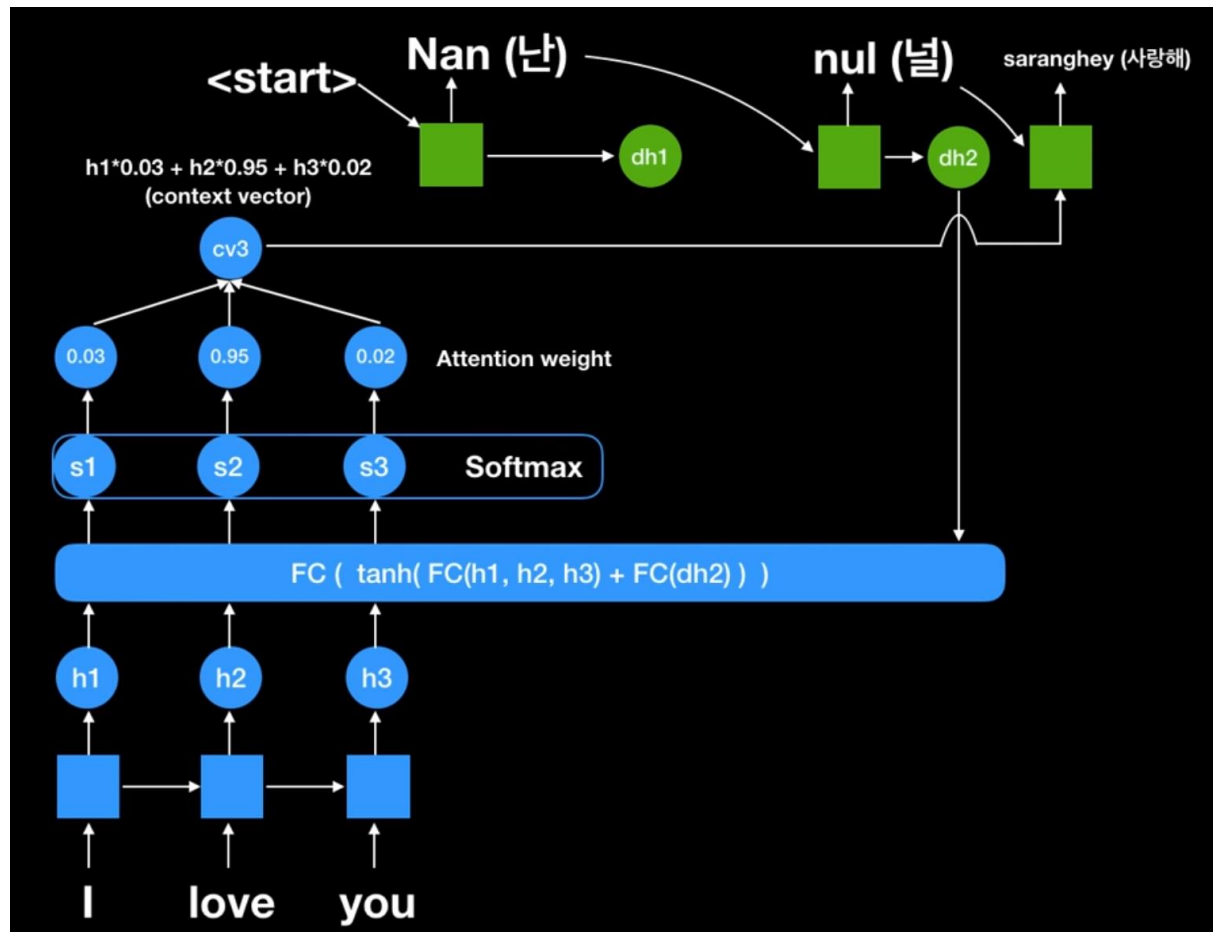
- ✓ Luong attention (Luong et al., 2015)

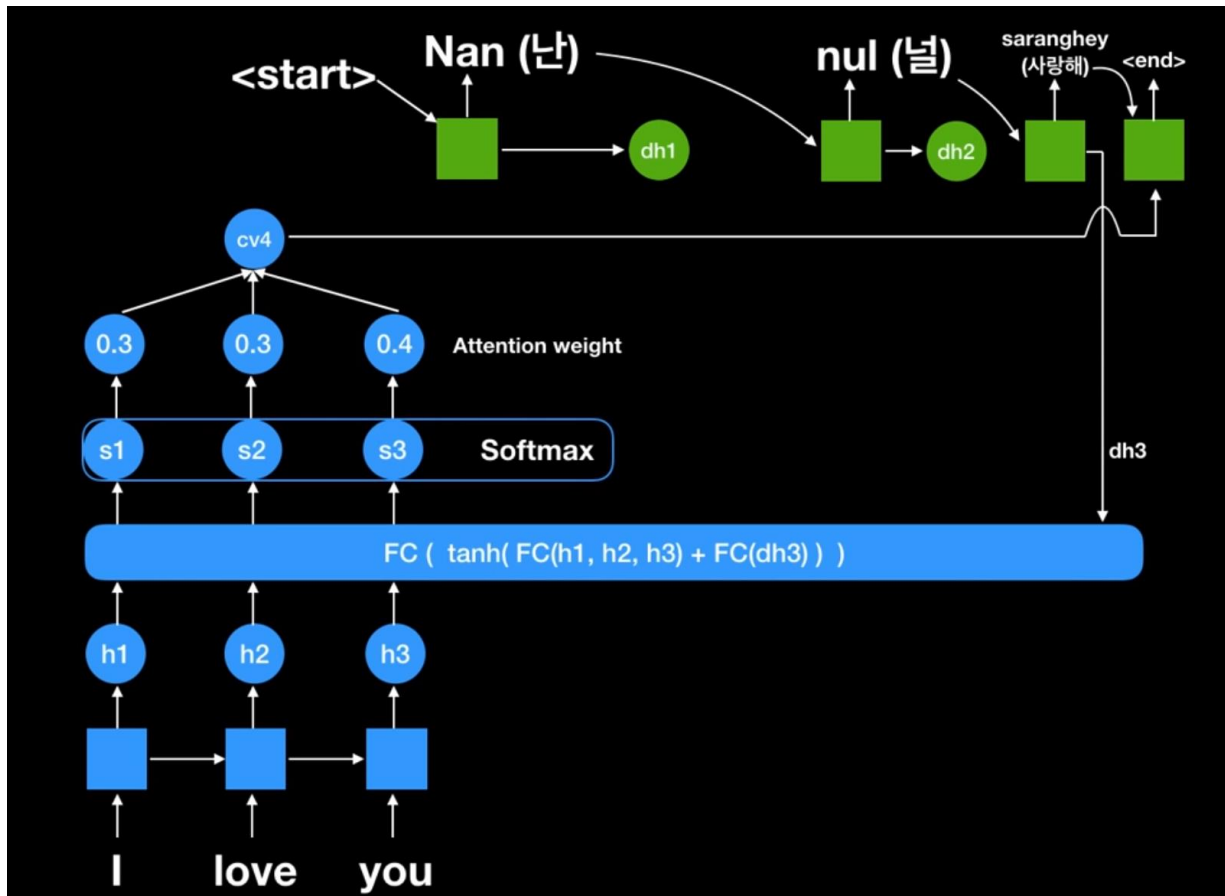
- Attention scores are not trained, the new current hidden state is the simple tanh of the weighed concatenation of the context vector and the current hidden state of the decoder











Attention at time step 4



Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

Encoding Stage



$h_1 h_2 h_3$



Attention Decoding Stage

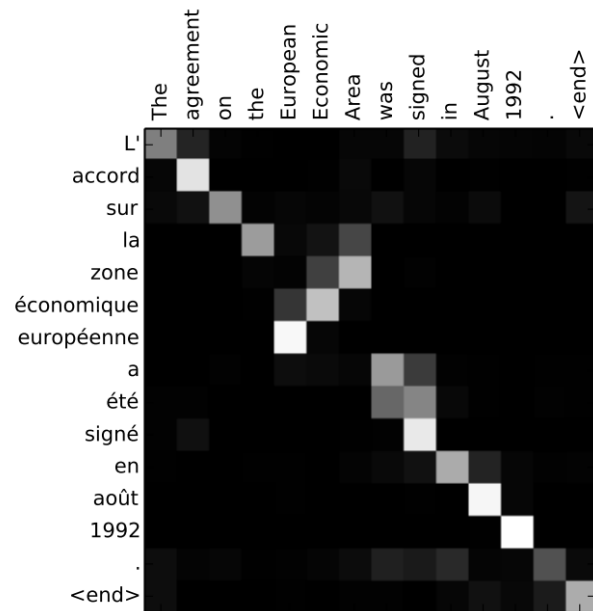
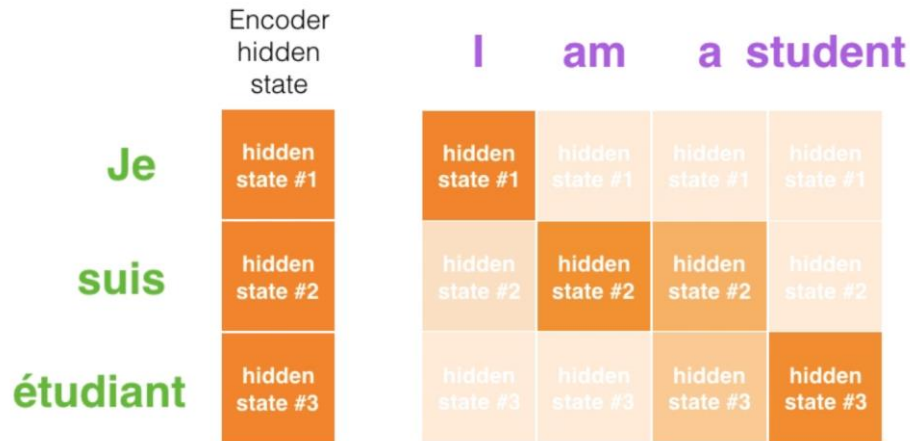


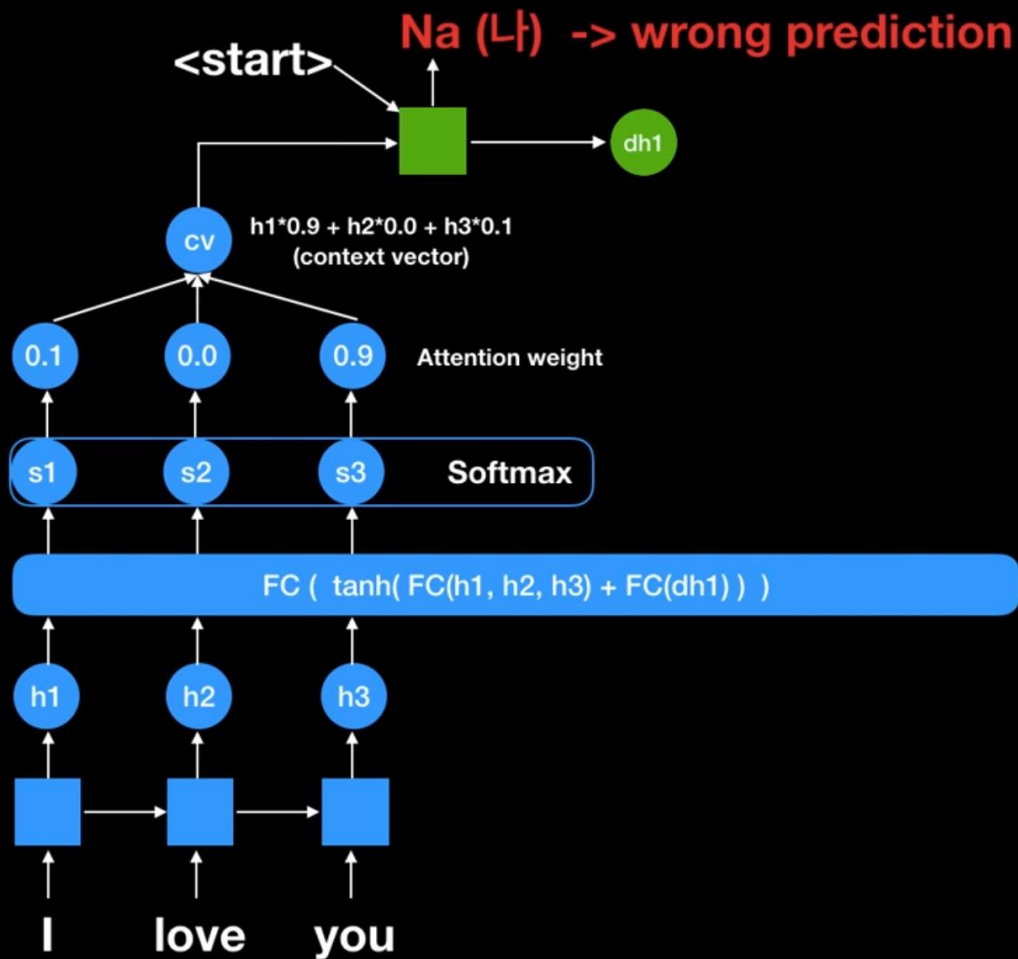
h_{init}

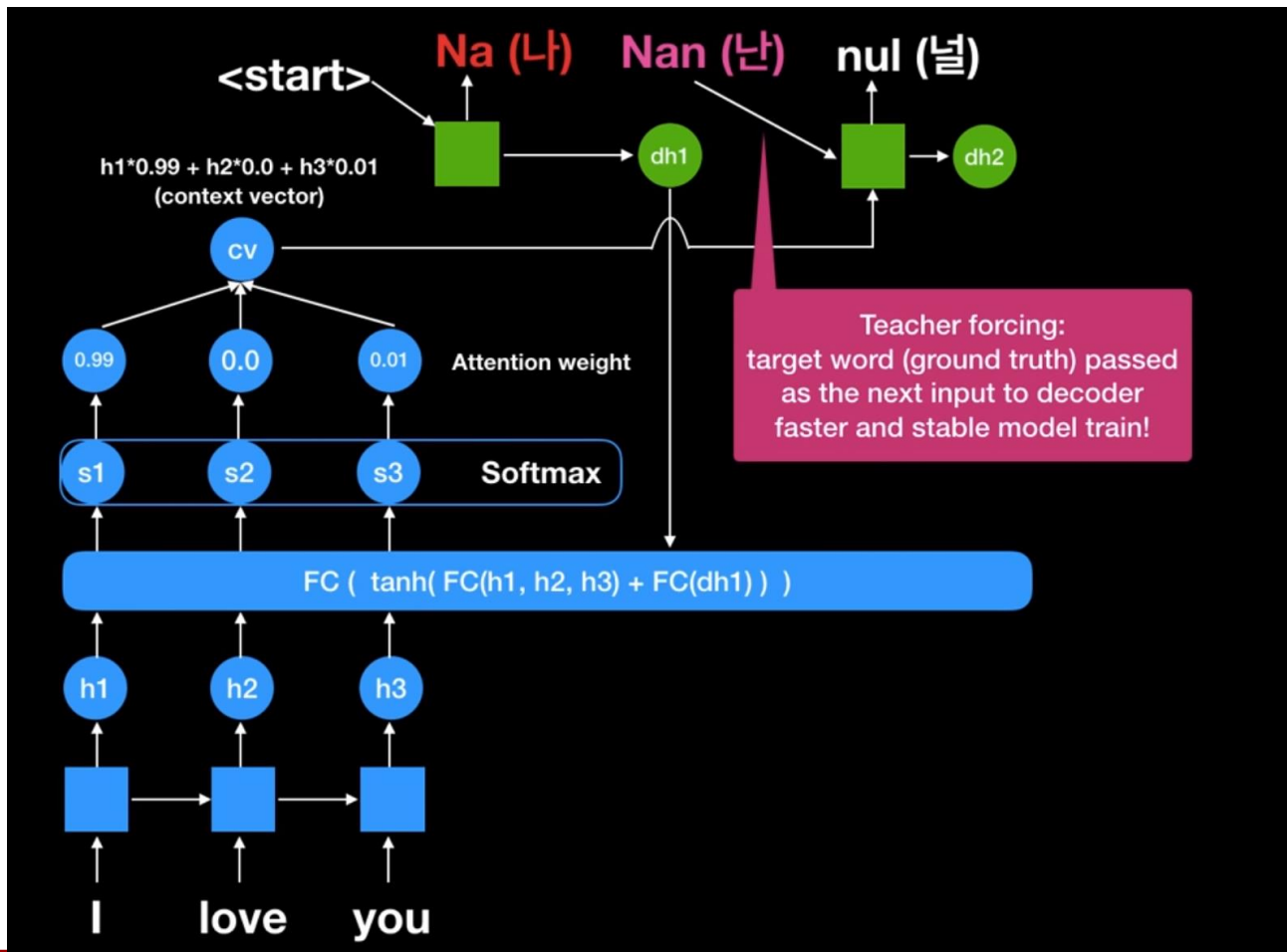


$\langle \text{END} \rangle$

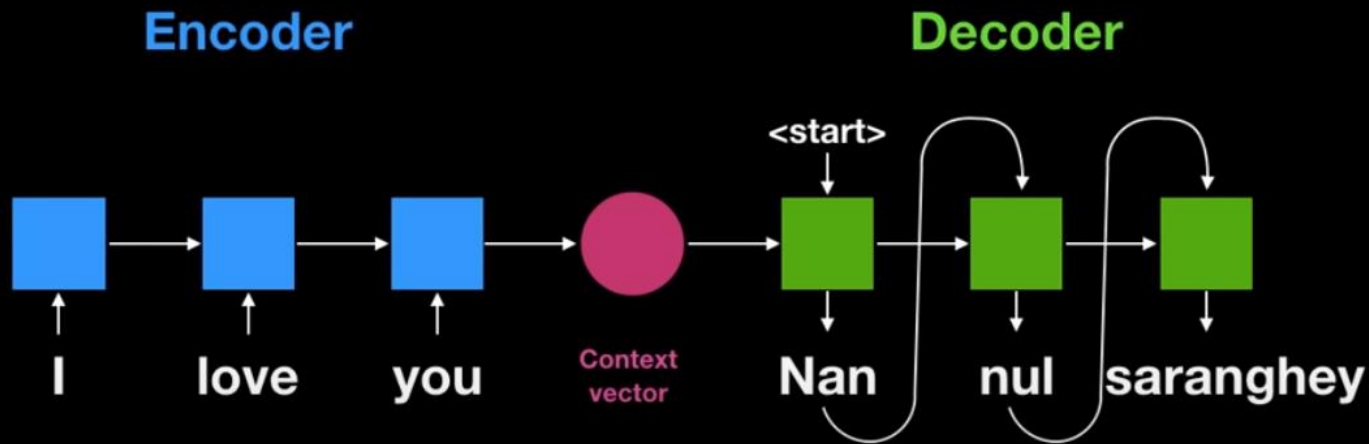
4



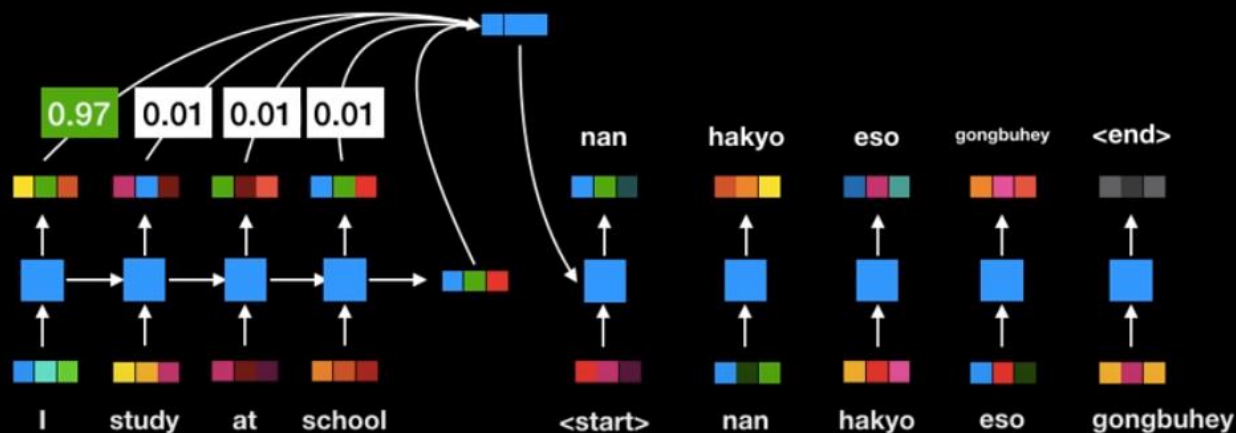




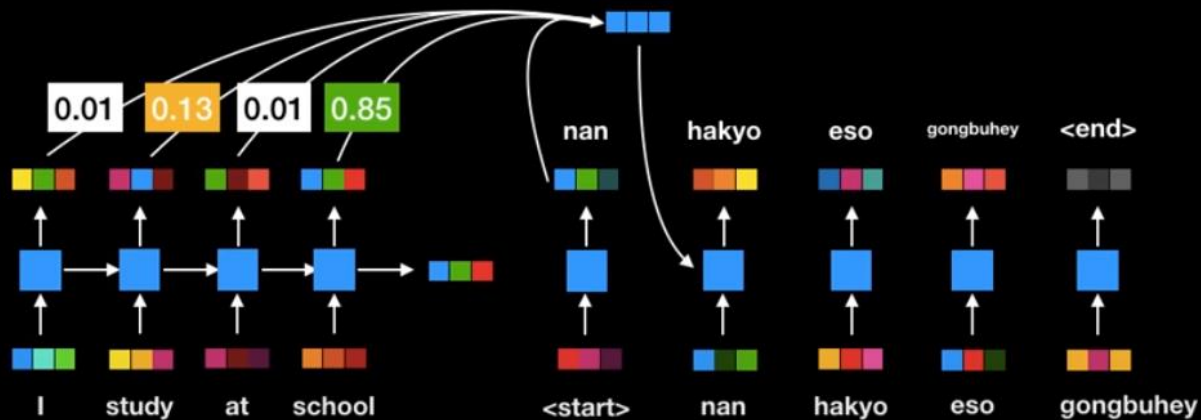
RNN based encoder decoder



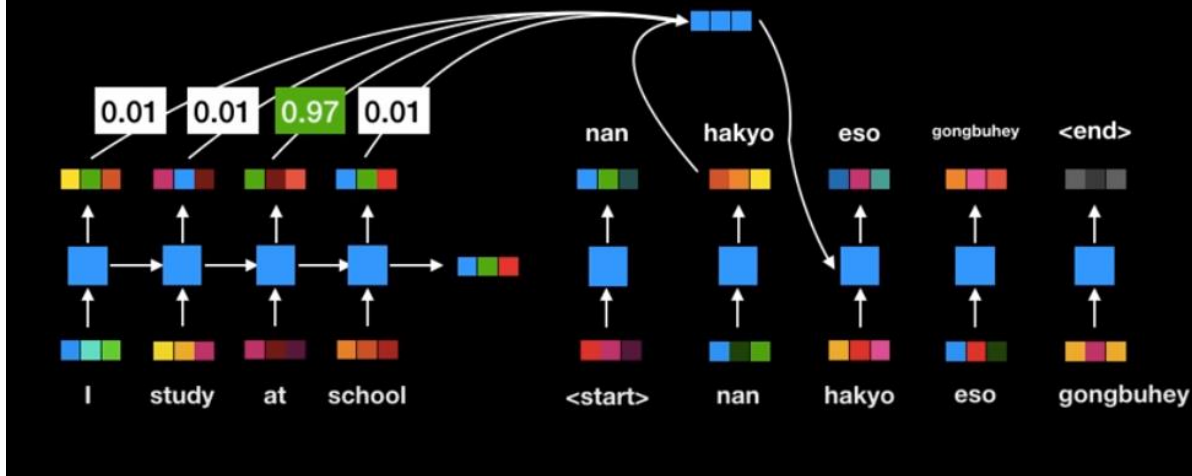
RNN based encoder decoder with attention



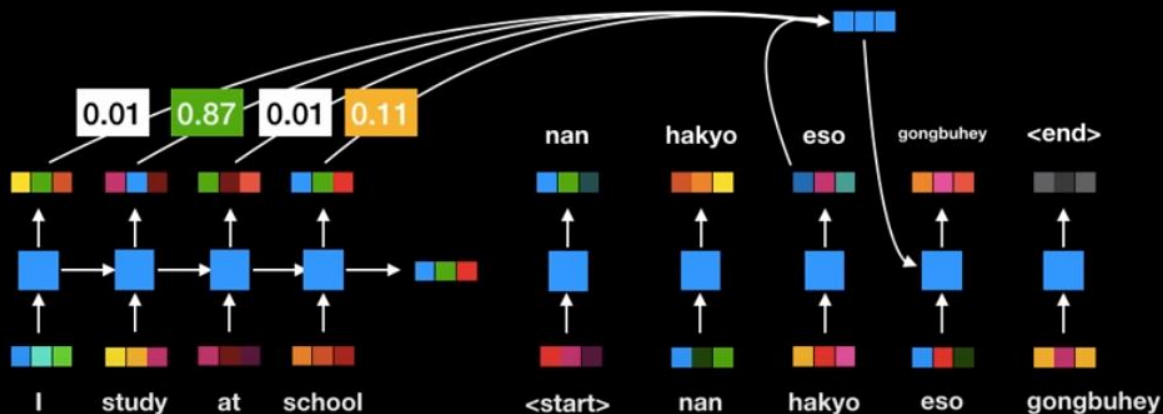
RNN based encoder decoder with attention



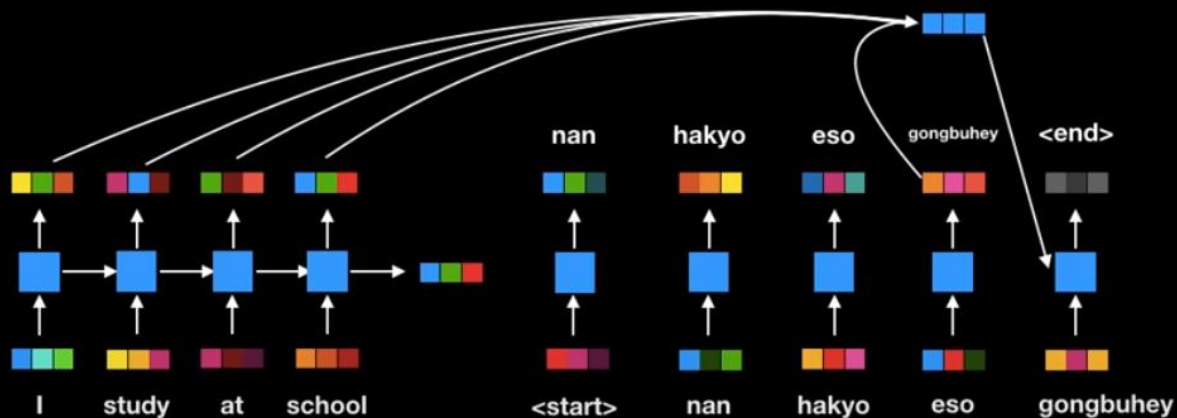
RNN based encoder decoder with attention



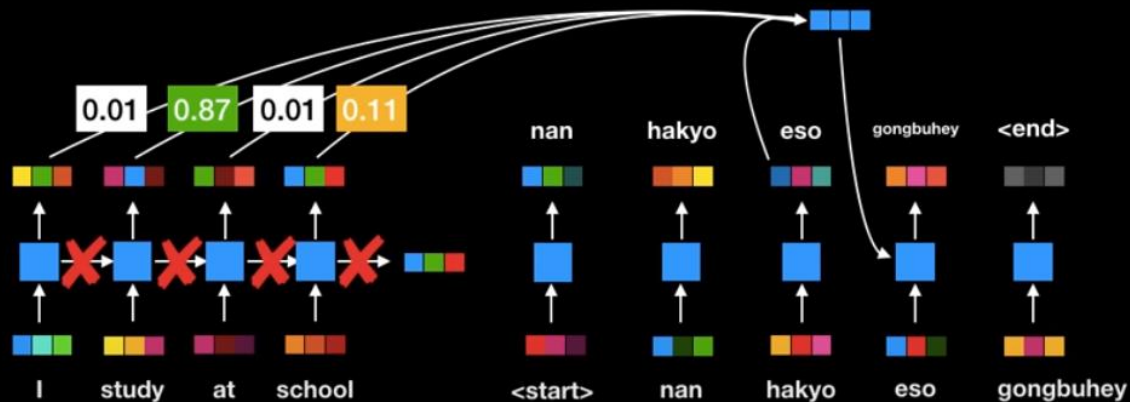
RNN based encoder decoder with attention



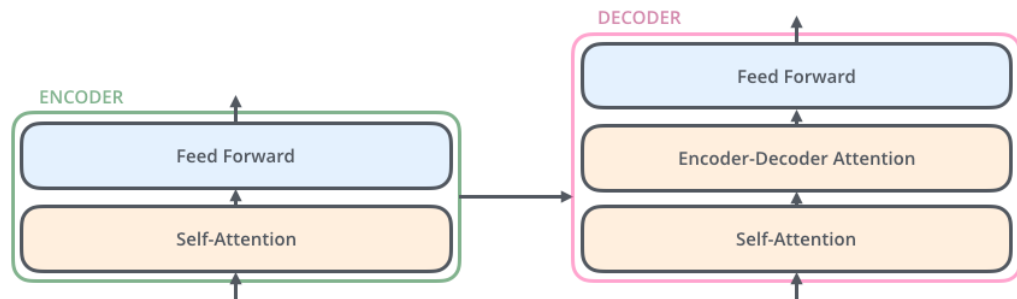
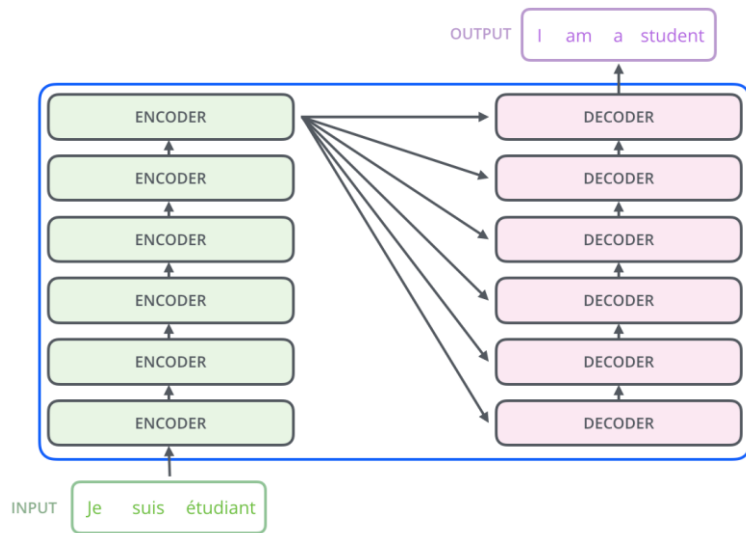
RNN based encoder decoder with attention

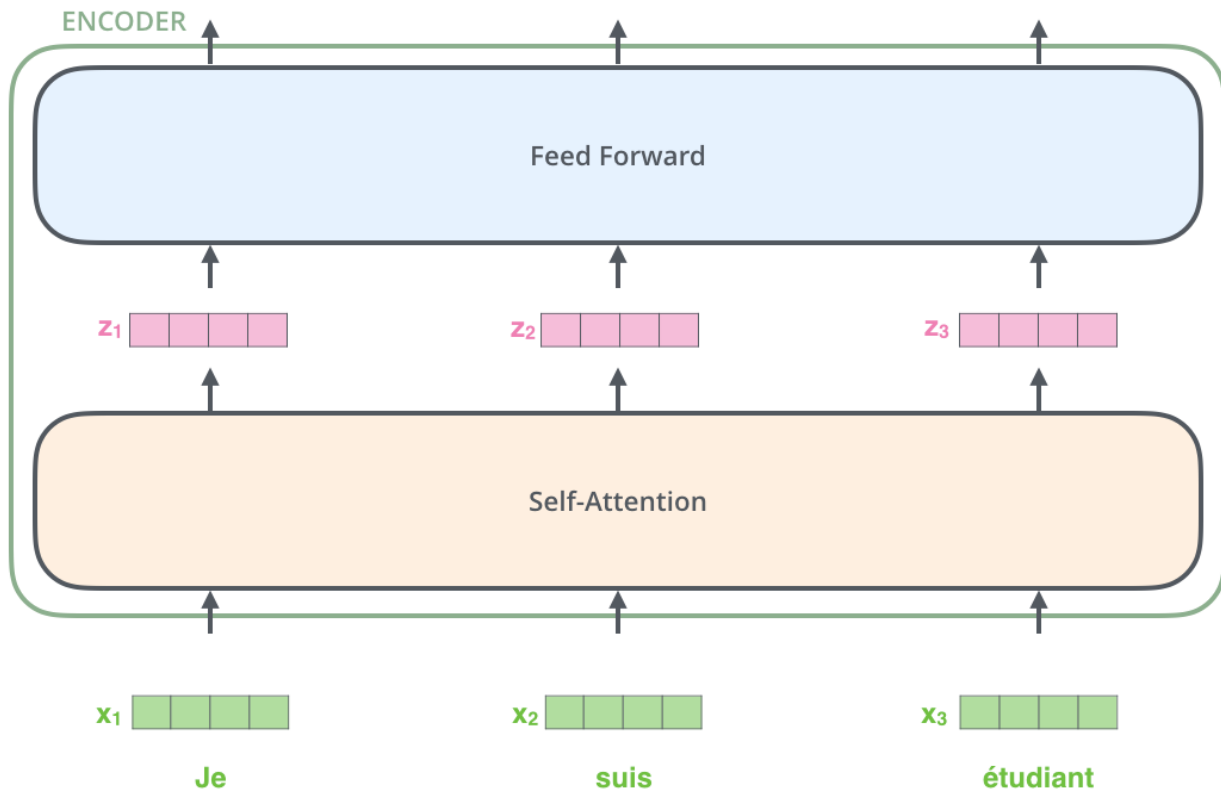


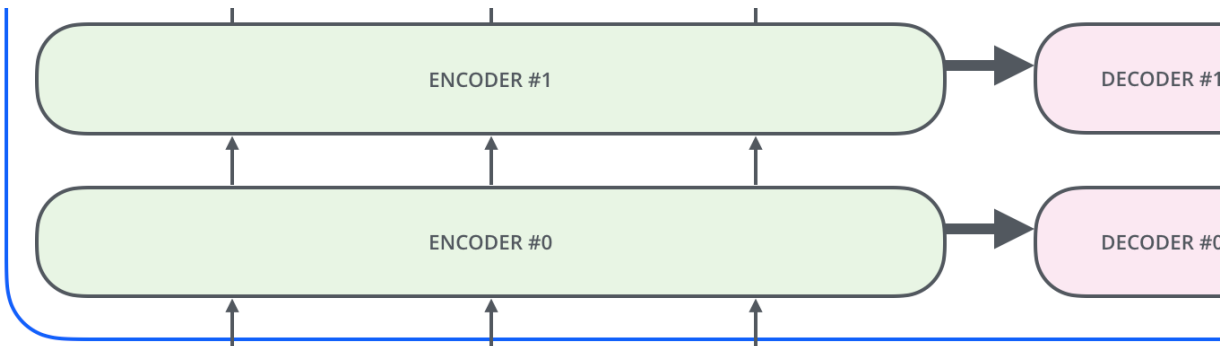
Can we just use attention?



Transformer







EMBEDDING
WITH TIME
SIGNAL

x_1

x_2

x_3

=

=

=

POSITIONAL
ENCODING

t_1

t_2

t_3

+

+

+

EMBEDDINGS

x_1

x_2

x_3

INPUT

Je

suis

étudiant

POSITIONAL
ENCODING

+

+

+

EMBEDDINGS

x_1

x_2

x_3

INPUT

Je

suis

étudiant

- Two properties that a good positional encoding scheme should have

✓ The norm of encoding vector is the same for all positions

✓ The further the two positions, the larger the distance

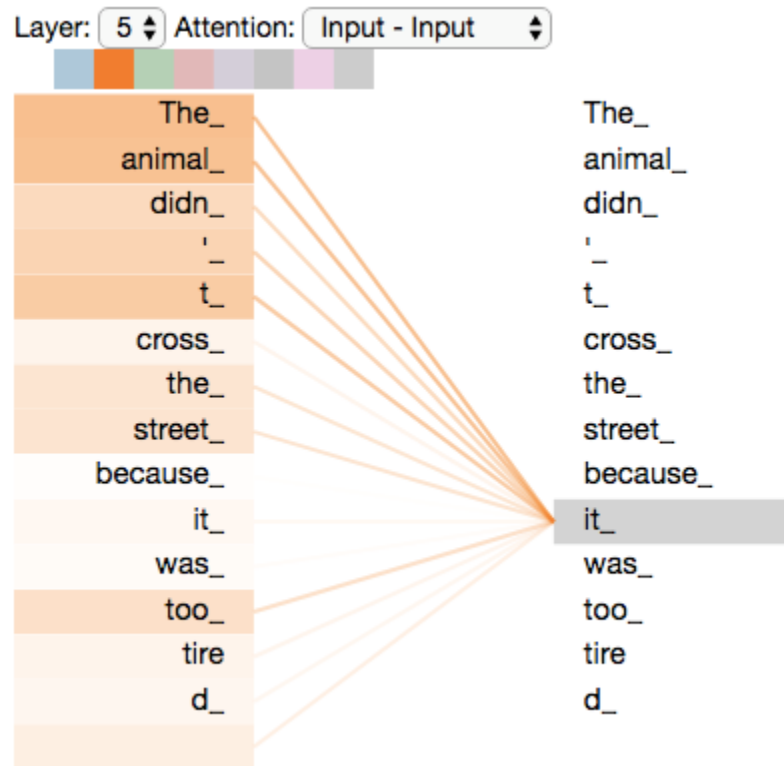
- A Simple Example ($n = 10$, $\text{dim} = 10$)

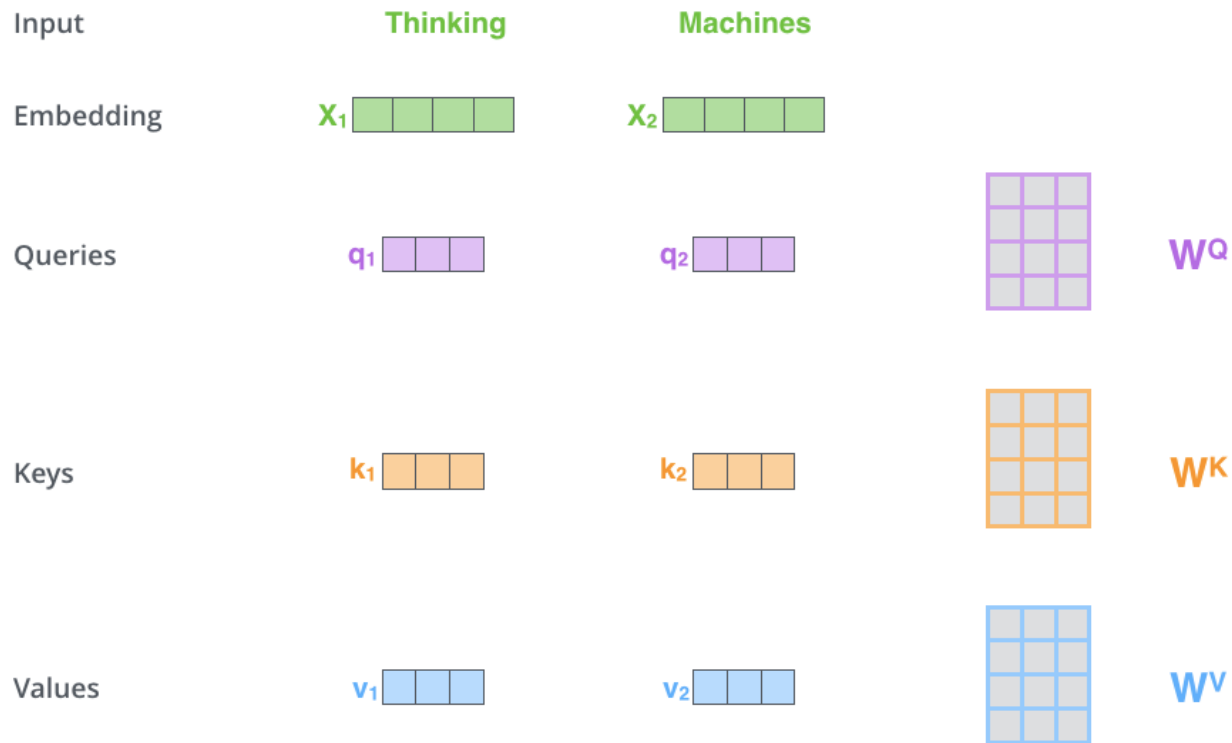
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10
X1	0.000	1.000	0.000	1.000	0.000	1.000	0.000	1.000	0.000	1.000
X2	0.841	0.674	0.638	0.839	0.461	0.922	0.325	0.962	0.227	0.982
X3	0.909	-0.093	0.983	0.408	0.818	0.699	0.615	0.852	0.442	0.928
X4	0.141	-0.798	0.875	-0.155	0.991	0.368	0.838	0.678	0.634	0.841
X5	-0.757	-0.983	0.366	-0.668	0.942	-0.022	0.970	0.452	0.793	0.723
X6	-0.959	-0.526	-0.312	-0.965	0.680	-0.408	0.996	0.192	0.911	0.579
X7	-0.279	0.275	-0.847	-0.952	0.267	-0.730	0.915	-0.082	0.981	0.415
X8	0.657	0.896	-0.992	-0.632	-0.207	-0.938	0.734	-0.350	0.999	0.235
X9	0.989	0.932	-0.681	-0.109	-0.635	-0.999	0.473	-0.591	0.966	0.046
X10	0.412	0.360	-0.057	0.450	-0.919	-0.904	0.161	-0.788	0.882	-0.144

"The animal didn't cross the street because it was too tired"



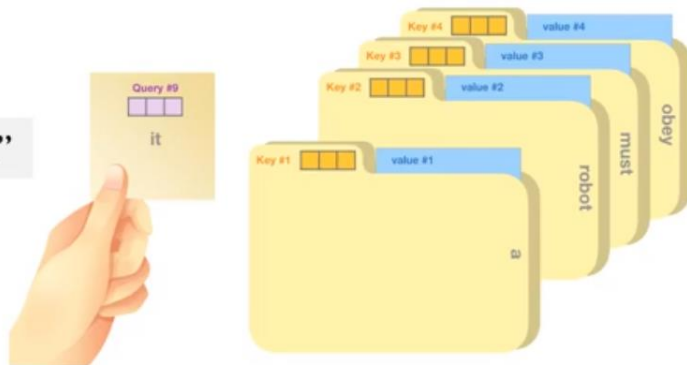


- Self-Attention in Detail

✓ Step 1: Create three vectors from each of the encoder's input vectors

- **Query:** The query is a representation of the current word used to score against all the other words (using their keys). We only care about the query of the token we're currently processing.
- **Key:** Key vectors are like labels for all the words in the segment. They're what we match against in our search for relevant words.
- **Value:** Value vectors are actual word representations, once we've scored how relevant each word is, these are the values we add up to represent the current word.

“A robot must obey the orders given it”



Input

Embedding

Queries

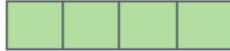
Keys

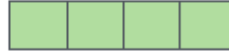
Values

Score

Thinking

Machines

x_1 

x_2 

q_1 

q_2 

k_1 

k_2 

v_1 

v_2 

$$q_1 \cdot k_1 = 112$$

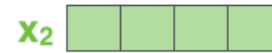
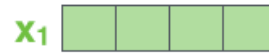
$$q_1 \cdot k_2 = 96$$

Input

Thinking

Machines

Embedding



Queries



Keys



Values



Score

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

Divide by 8 ($\sqrt{d_k}$)

14

12

Softmax

0.88

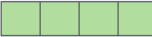
0.12

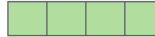
Input

Thinking

Machines

Embedding

x_1 

x_2 

Queries

q_1 

q_2 

Keys

k_1 

k_2 

Values

v_1 

v_2 

Score

$q_1 \cdot k_1 = 112$

$q_1 \cdot k_2 = 96$

Divide by 8 ($\sqrt{d_k}$)

14

12

Softmax

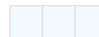
0.88

0.12

Softmax

X
Value

v_1 

v_2 

Sum

z_1 















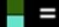





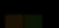















































z_2 

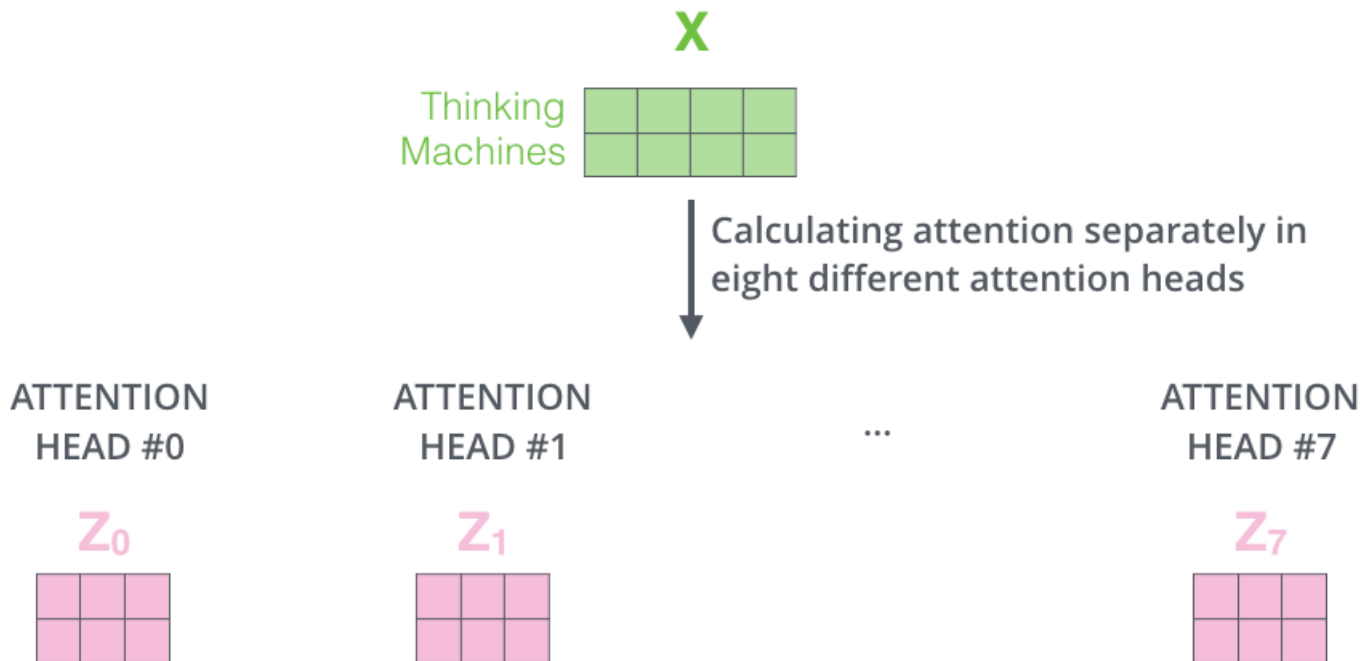
$$\begin{array}{c} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline \end{array} \end{array} \times \begin{array}{c} W^Q \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{array} = \begin{array}{c} Q \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array}$$

$$\begin{array}{c} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline \end{array} \end{array} \times \begin{array}{c} W^K \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{array} = \begin{array}{c} K \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array}$$

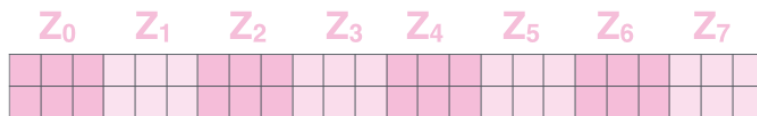
$$\begin{array}{c} \text{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline \end{array} \end{array} \times \begin{array}{c} W^V \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{array} = \begin{array}{c} V \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array}$$

$$\begin{array}{c} \text{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array} \times \begin{array}{c} K^T \\ \begin{array}{|c|c|} \hline & \\ \hline & \\ \hline & \\ \hline \end{array} \end{array} \\ \hline \sqrt{d_k} \\ \text{softmax} \left(\frac{\quad}{\quad} \right) \begin{array}{c} V \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array} \\ = \begin{array}{c} Z \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{array}$$

	Query * Key ^T	Score	Softmax	Value	Softmax * Value	Σ Softmax * Value (Attention layer output)
I	I * I  *  = 130	0.92		I		
	I * study  *  = 50	0.05		study		
	I * at  *  = 20	0.02		at		
	I * school  *  = 10	0.01		school		
study	study * I  *  = 30	0.02		I		
	study * study  *  = 110	0.70		study		
	study * at  *  = 20	0.03		at		
	study * school  *  = 70	0.25		school		
at	at * I  *  = 30	0.03		I		
	at * study  *  = 50	0.10		study		
	at * at  *  = 90	0.80		at		
	at * school  *  = 40	0.07		school		
school	school * I  *  = 30	0.01		I		
	school * study  *  = 80	0.27		study		
	school * at  *  = 23	0.02		at		
	school * school  *  = 160	0.70		school		

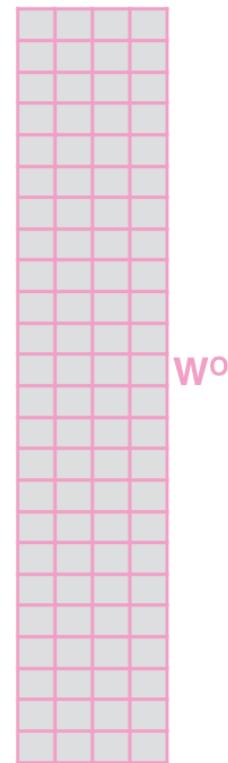


1) Concatenate all the attention heads

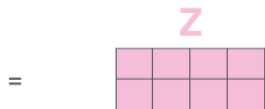


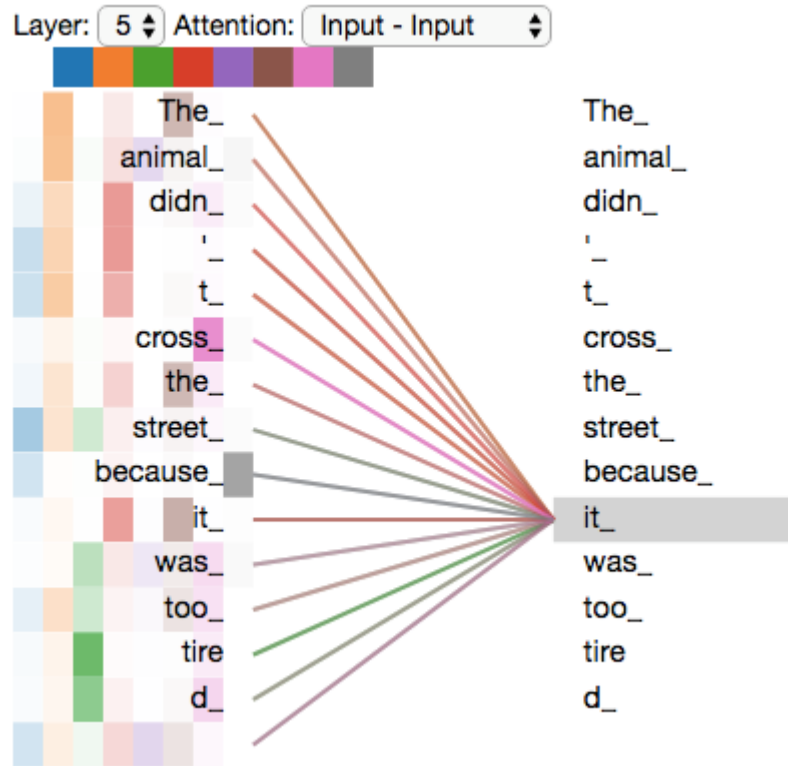
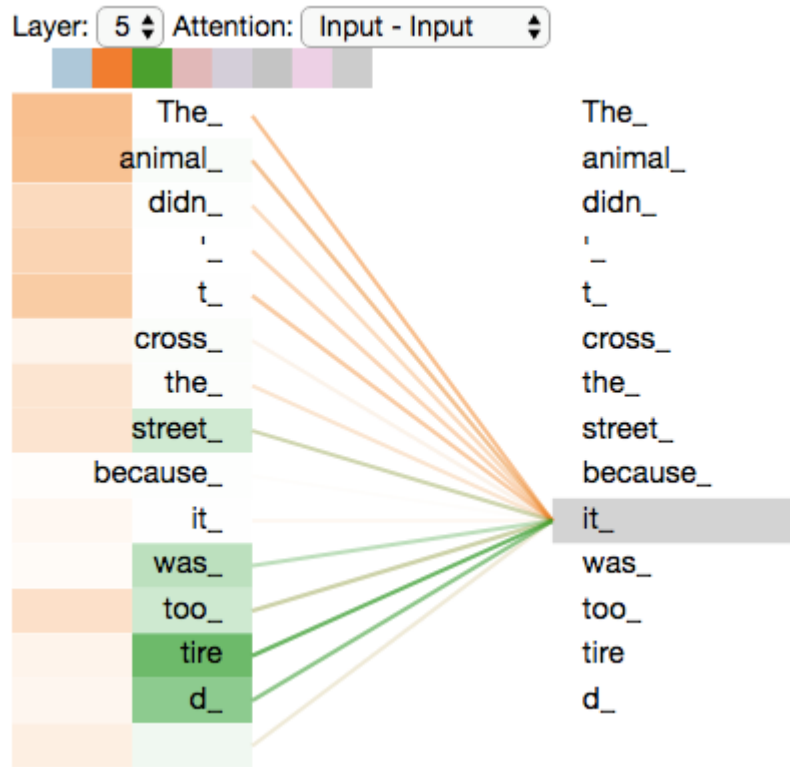
2) Multiply with a weight matrix W^O that was trained jointly with the model

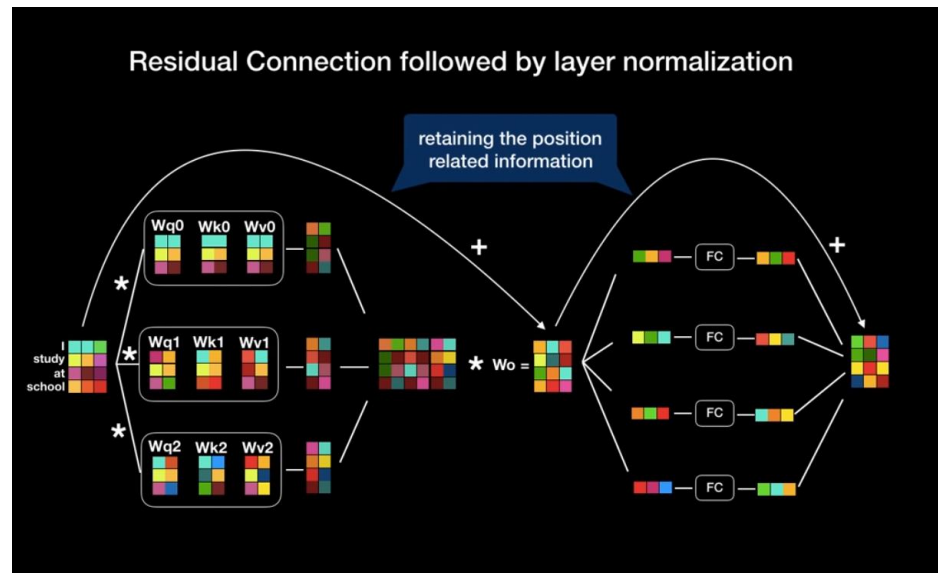
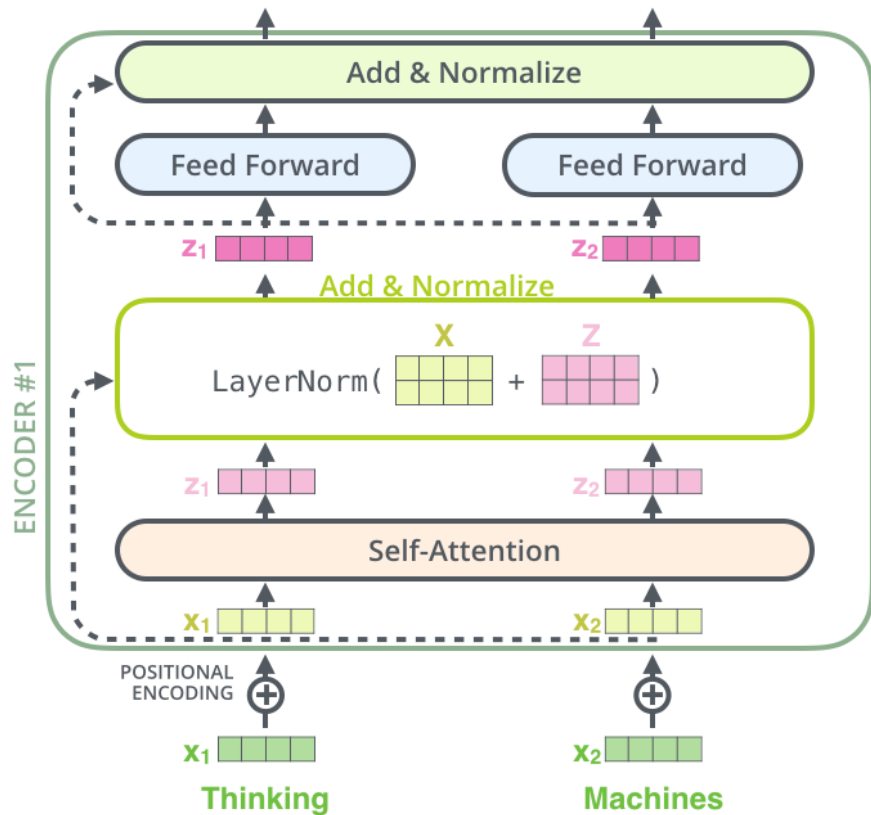
\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN





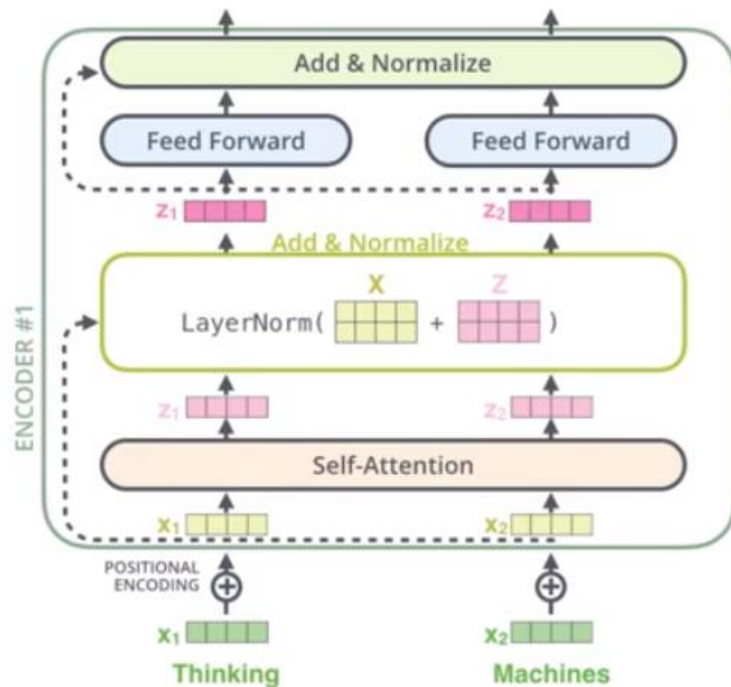


- Position-wise Feed-Forward Networks

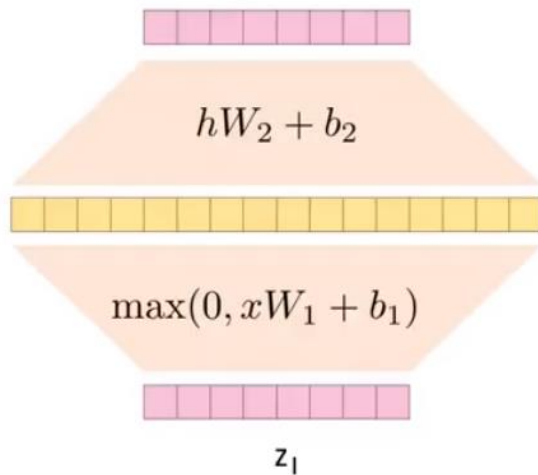
- ✓ Fully connected feed-forward network
- ✓ Applied to each position separately and identically

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

- ✓ The linear transformations are the same across different positions
- ✓ They use different parameters from layer to layer



- Position-wise Feed-Forward Networks

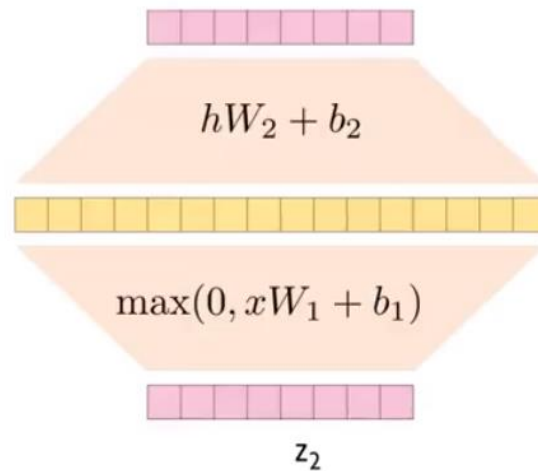


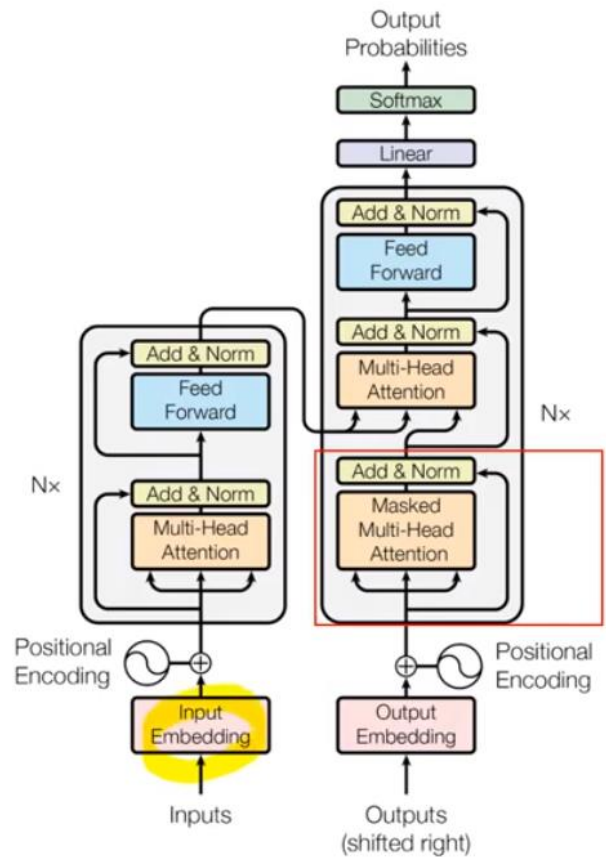
512 dim.

2048 dim.

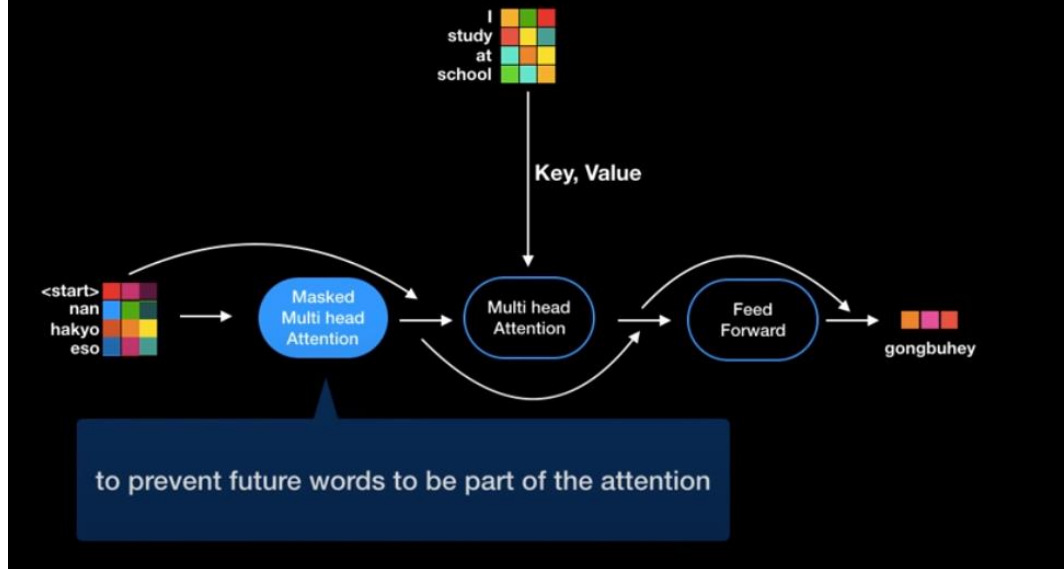
512 dim.

z_1



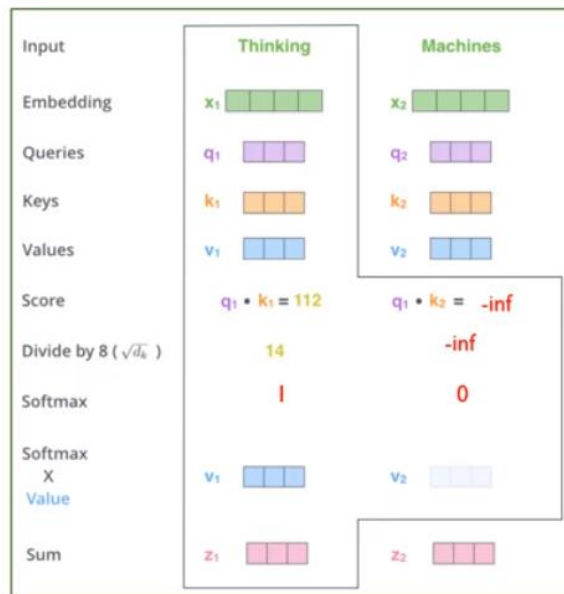
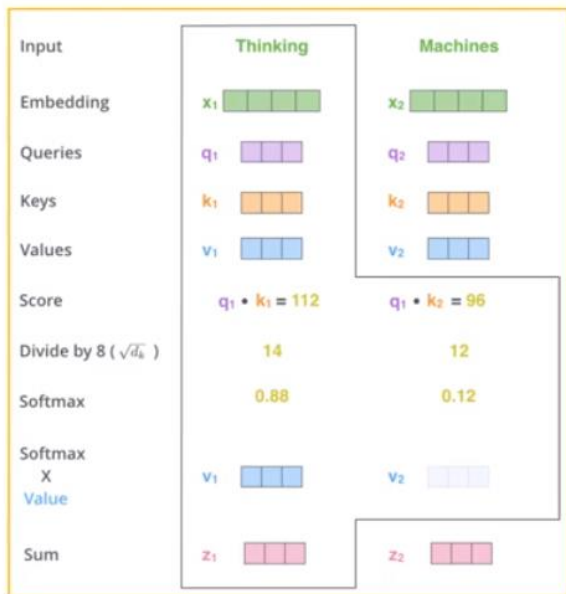


Decoder Layer



- Masked Multi-head Attention

- ✓ Self attention layers in the decoder is only allowed to attend to **earlier positions** in the output sequence, which is done by masking future positions (setting them to $-\text{inf}$) before the softmax step in the self attention calculation.



- Masked Multi-head Attention

Queries				X	Keys				=	Scores (before softmax)			
robot	must	obey	orders		robot	must	obey	orders		0.11	0.00	0.81	0.79
robot	must	obey	orders		robot	must	obey	orders		0.19	0.50	0.30	0.48
robot	must	obey	orders		robot	must	obey	orders		0.53	0.98	0.95	0.14
robot	must	obey	orders		robot	must	obey	orders		0.81	0.86	0.38	0.90

Scores (before softmax)				Apply Attention Mask →	Masked Scores (before softmax)			
0.11	0.00	0.81	0.79		0.11	-inf	-inf	-inf
0.19	0.50	0.30	0.48		0.19	0.50	-inf	-inf
0.53	0.98	0.95	0.14		0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90		0.81	0.86	0.38	0.90

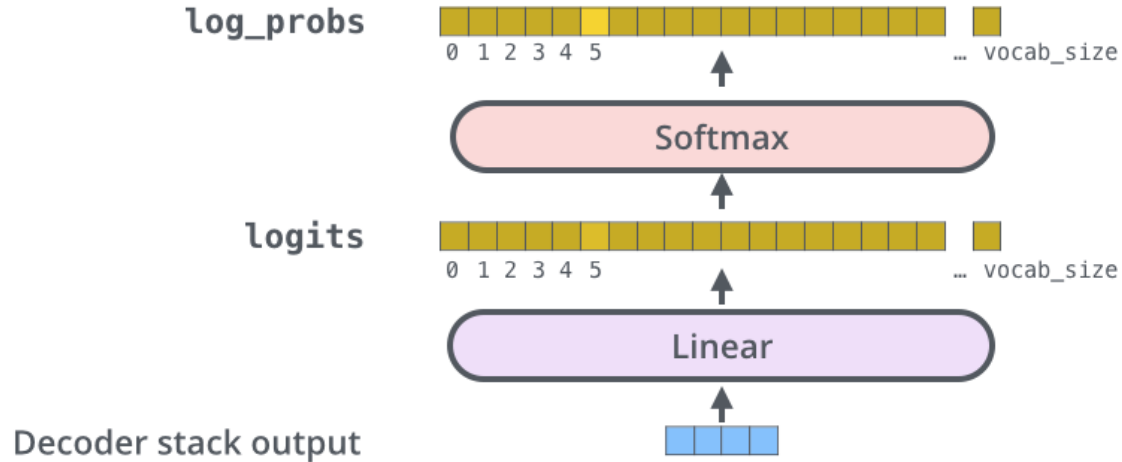
Masked Scores (before softmax)				Softmax (along rows)	Scores			
0.11	-inf	-inf	-inf		1	0	0	0
0.19	0.50	-inf	-inf		0.48	0.52	0	0
0.53	0.98	0.95	-inf		0.31	0.35	0.34	0
0.81	0.86	0.38	0.90		0.25	0.26	0.23	0.26

Which word in our vocabulary
is associated with this index?

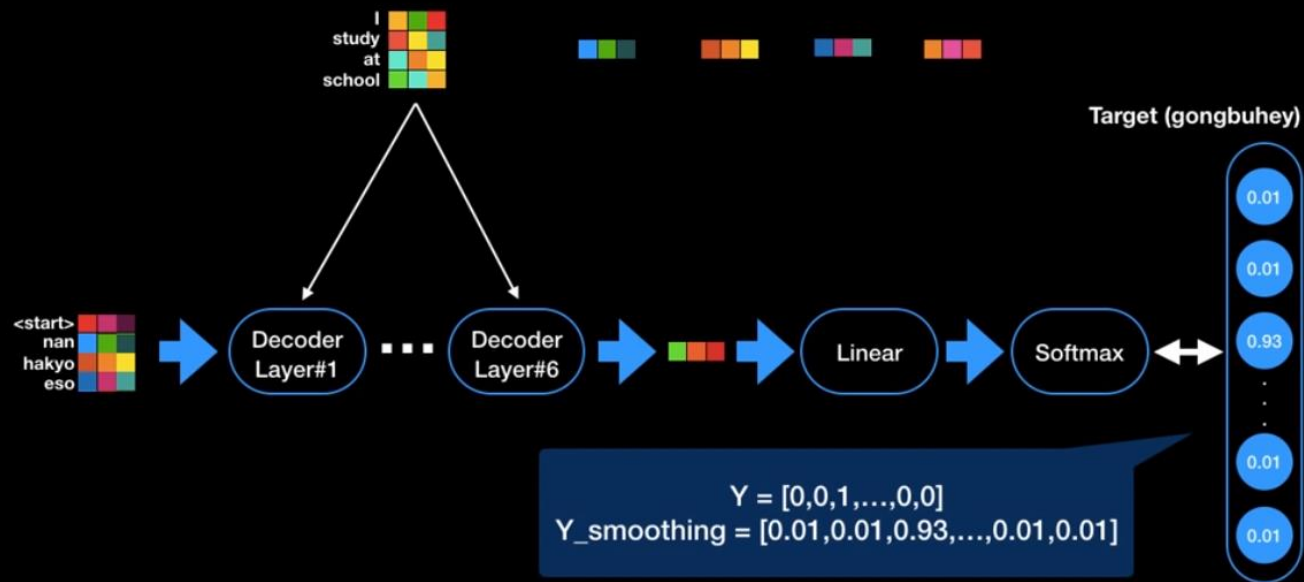
am

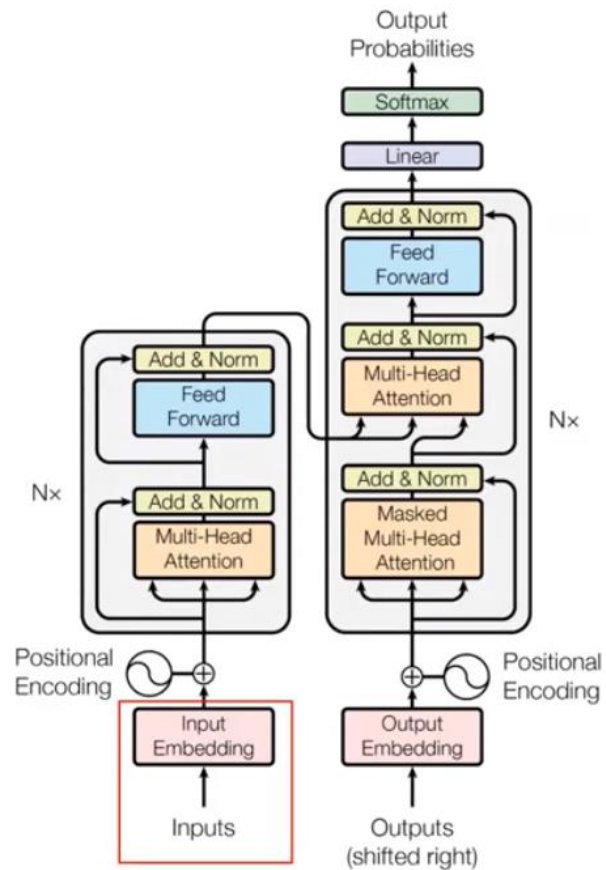
Get the index of the cell
with the highest value
(**argmax**)

5



Label Smoothing



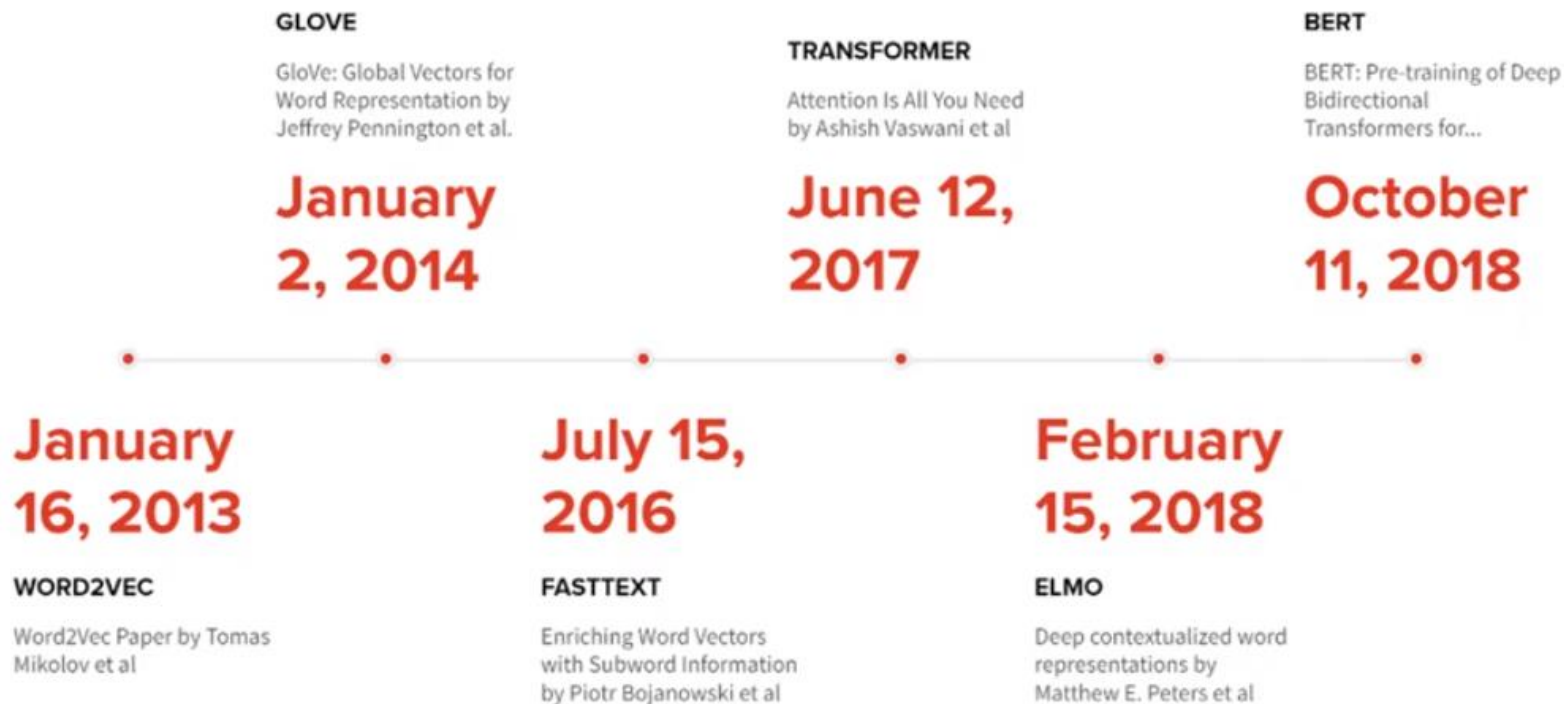


- Performance (in terms of BLEU score)

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

- Transformer variations

	N	d_{model}	d_{ff}	h	d_k	d_v	P_{drop}	ϵ_{ls}	train steps	PPL (dev)	BLEU (dev)	params $\times 10^6$		
base	6	512	2048	8	64	64	0.1	0.1	100K	4.92	25.8	65		
(A)					1	512	512				5.29	24.9		
					4	128	128				5.00	25.5		
					16	32	32				4.91	25.8		
					32	16	16				5.01	25.4		
(B)					16					5.16	25.1	58		
					32					5.01	25.4	60		
(C)	2									6.11	23.7	36		
	4									5.19	25.3	50		
	8									4.88	25.5	80		
		256					32	32				5.75	24.5	28
		1024					128	128				4.66	26.0	168
			1024									5.12	25.4	53
			4096									4.75	26.2	90
(D)							0.0				5.77	24.6		
							0.2				4.95	25.5		
								0.0				4.67	25.3	
								0.2				5.47	25.7	
(E)	positional embedding instead of sinusoids									4.92	25.7			
big	6	1024	4096	16					0.3	300K	4.33	26.4	213	



reference

- <http://jalammar.github.io/illustrated-transformer/>
- Attention Is All You Need
<https://arxiv.org/abs/1706.03762>