

TRANSFORMER : Attention Is All You Need

엔자이너 연구실

Contents

- ▶ RNN
- ▶ LSTM
- ▶ GRU
- ▶ Transformer : Attention Is All You Need

How to implement GNNs in my study

▶ Attention Is All You Need

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

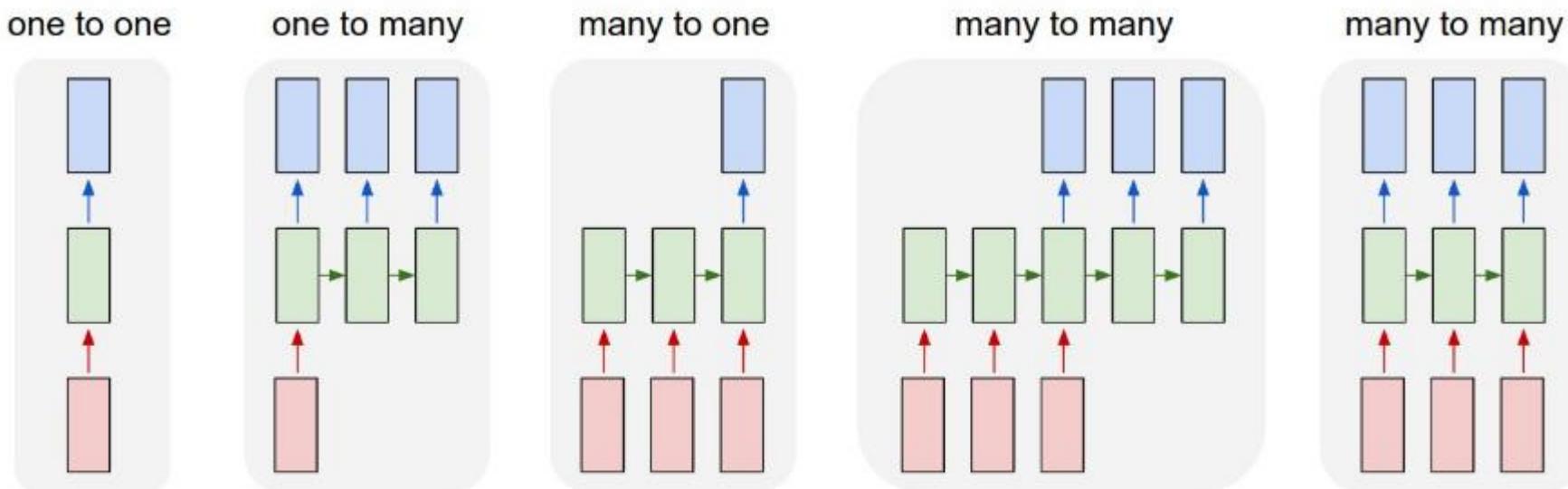
Ilia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

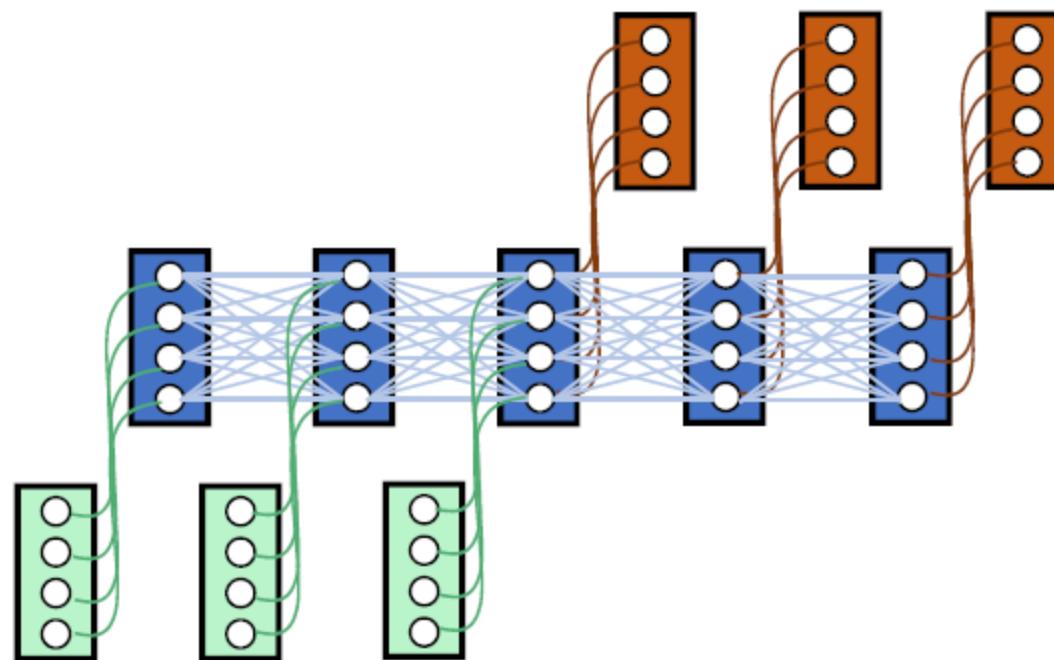
RNN

▶ Many to Many



RNN

▶ Many to Many



RNN

▶ Word Embedding

그림 3-16 말뭉치에서 맥락과 타깃을 만드는 예

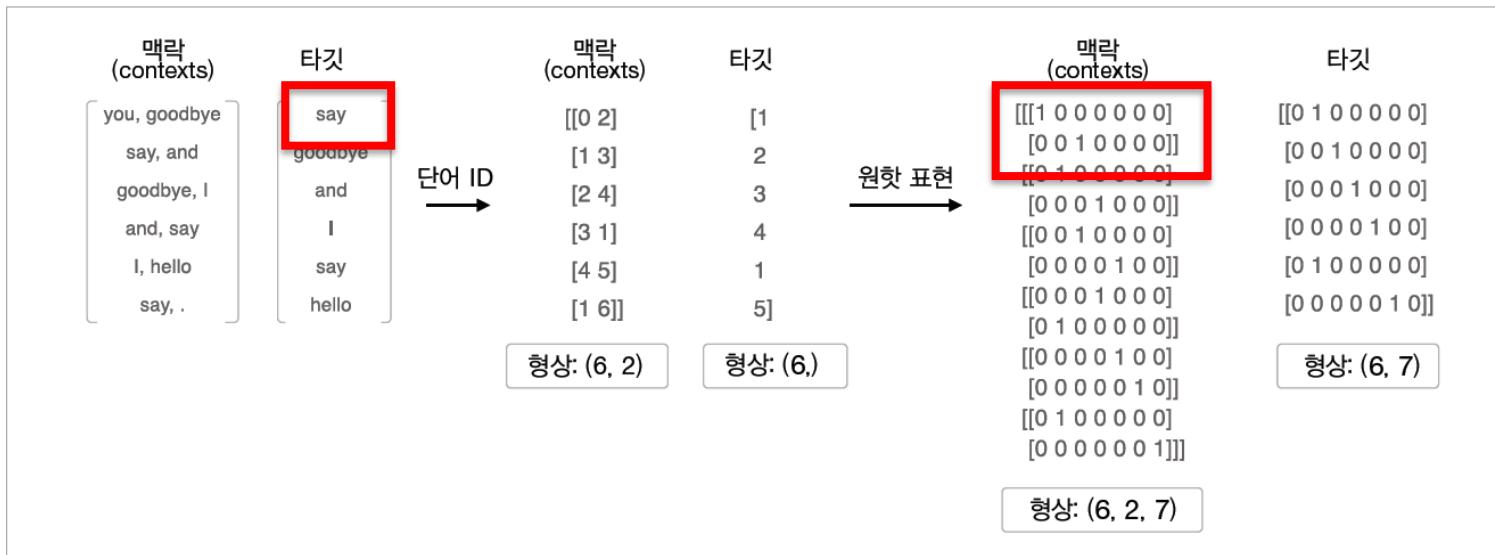
말뭉치	맥락(contexts)	타깃
you <u>say</u> goodbye and I say hello .	you, goodbye	say
you say <u>goodbye</u> and I say hello .	say, and	goodbye
you say goodbye <u>and</u> I say hello .	goodbye, I	and
you say goodbye and <u>I</u> say hello .	and, say	I
you say goodbye and I <u>say</u> hello .	I, hello	say
you say goodbye and I say <u>hello</u> .	say, .	hello

RNN

▶ Word Embedding

▶ You say goodbye and I say hello

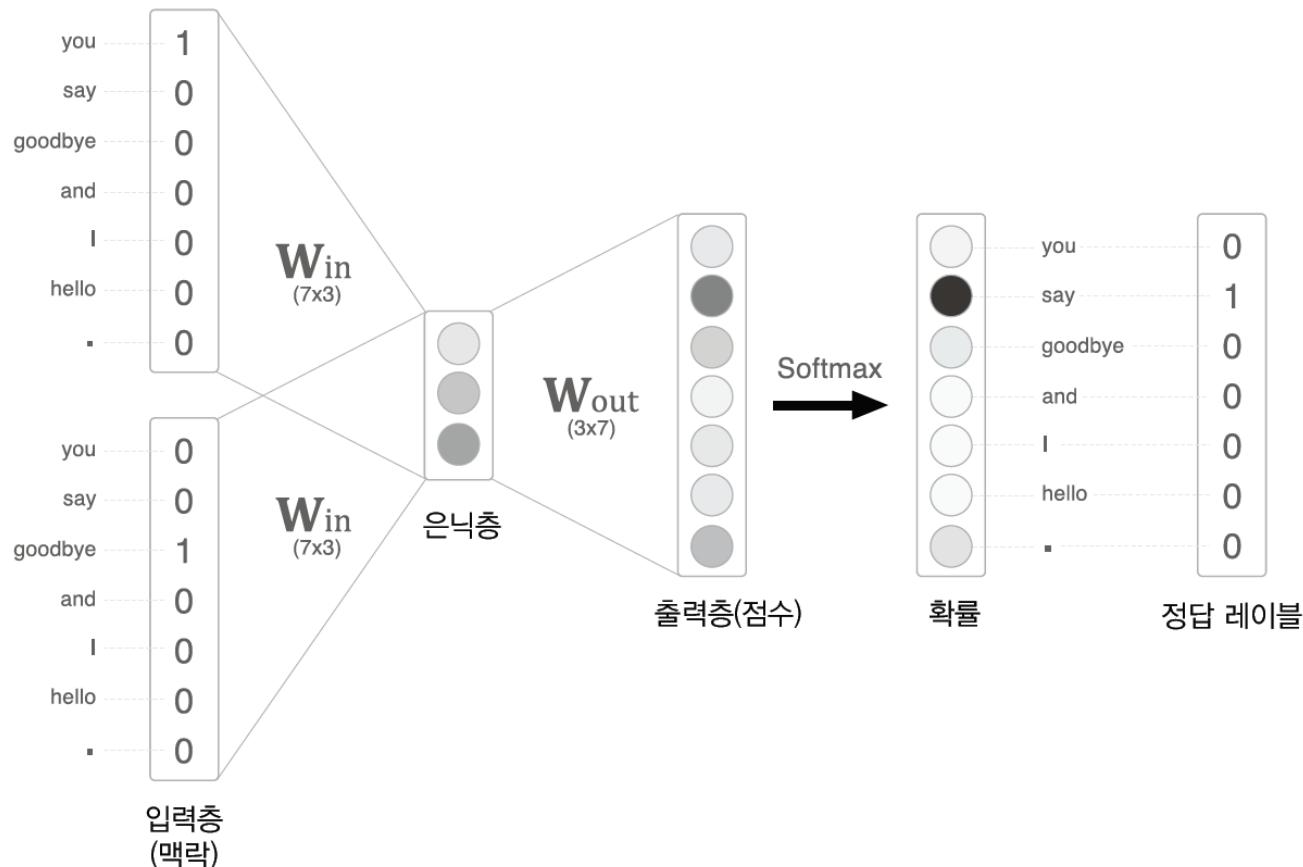
그림 3-18 ‘맥락’과 ‘타깃’을 원핫 표현으로 변환하는 예



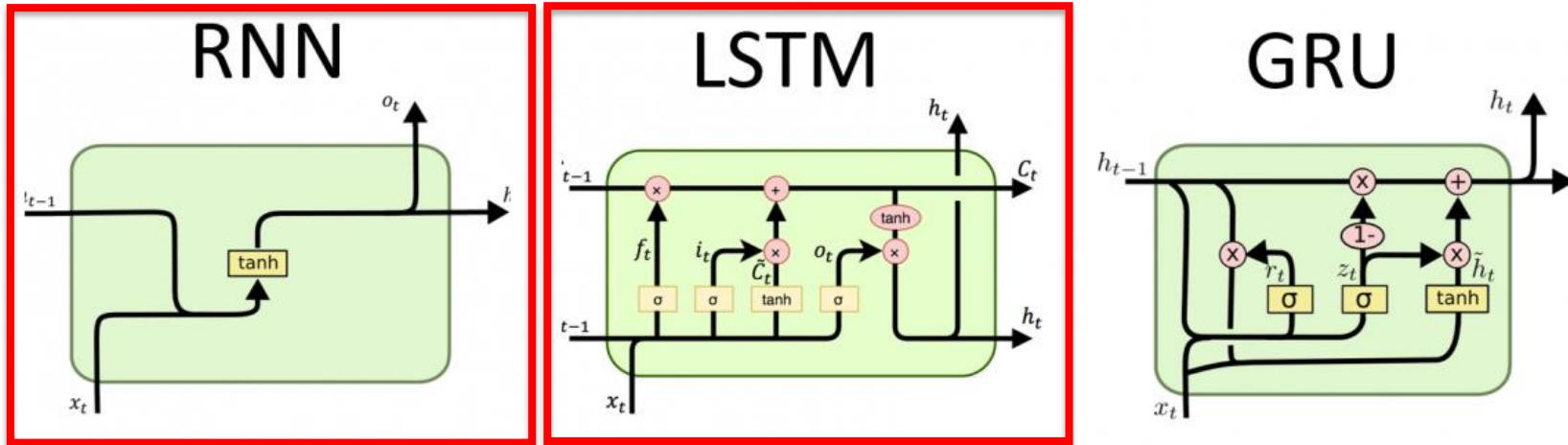
RNN

▶ Word Embedding

그림 3-12 CBOW 모델의 구체적인 예(노드 값의 크기를 흑백의 진하기로 나타냄)

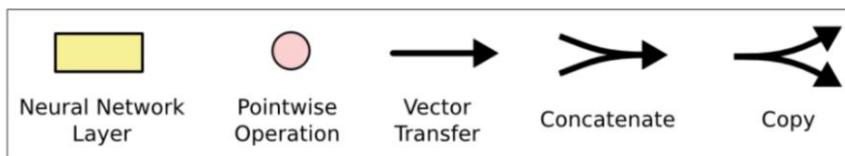
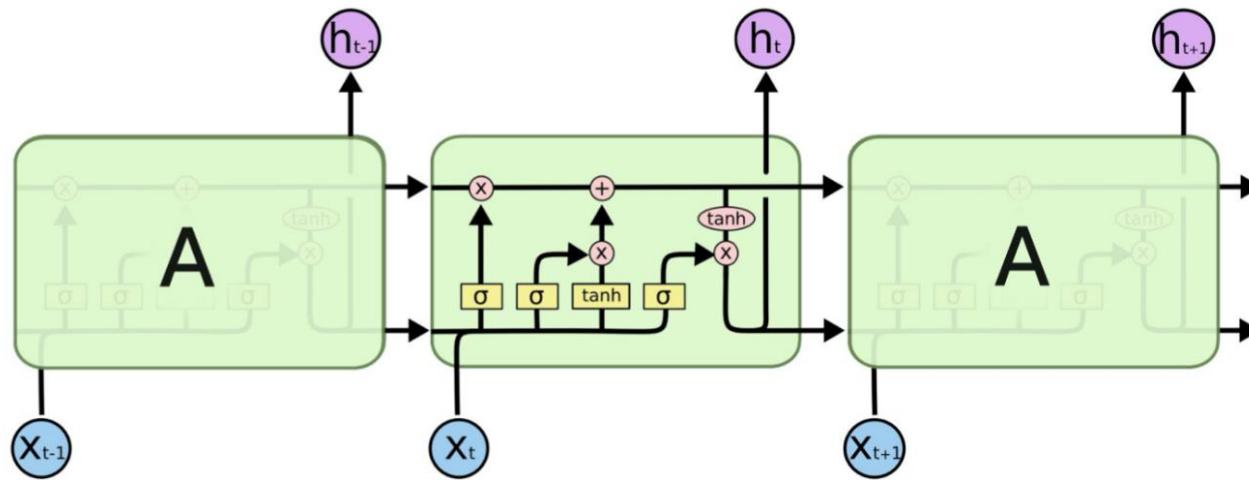


RNN



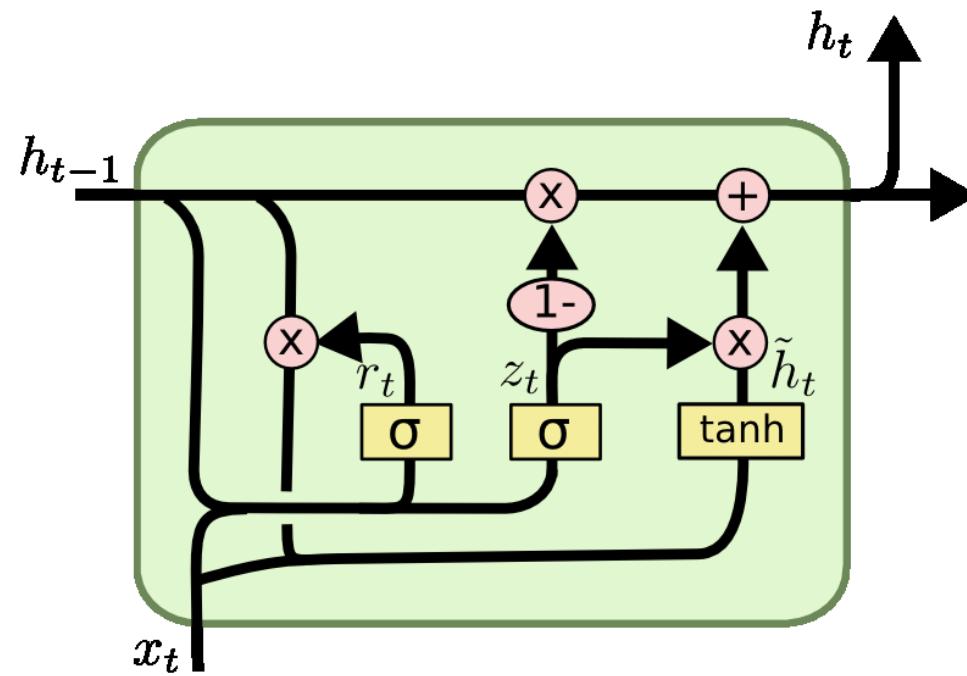
LSTM

▶ Architecture

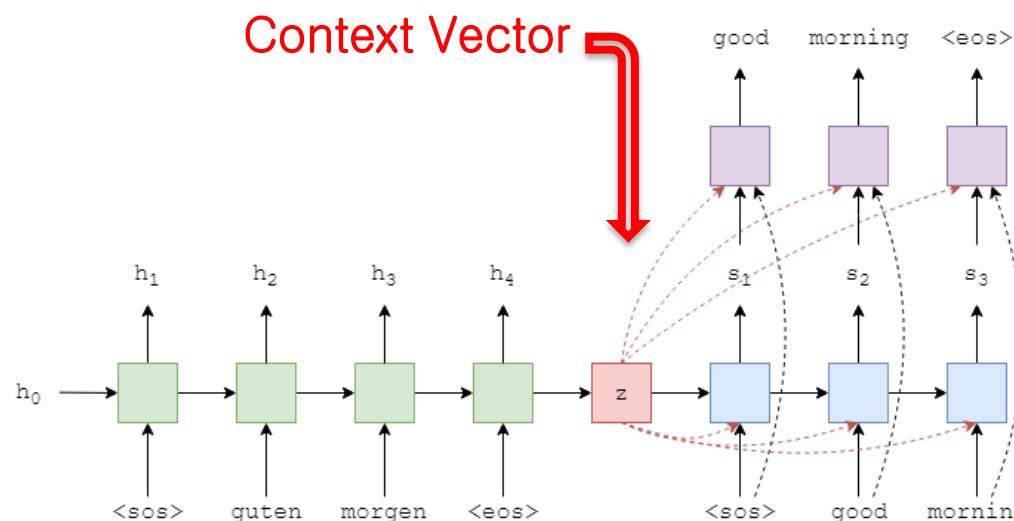
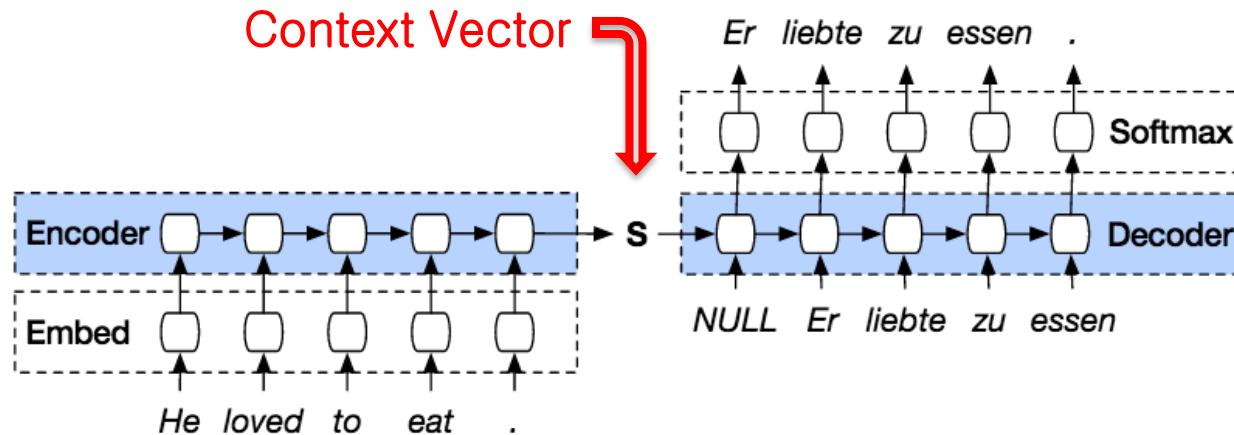


GRU

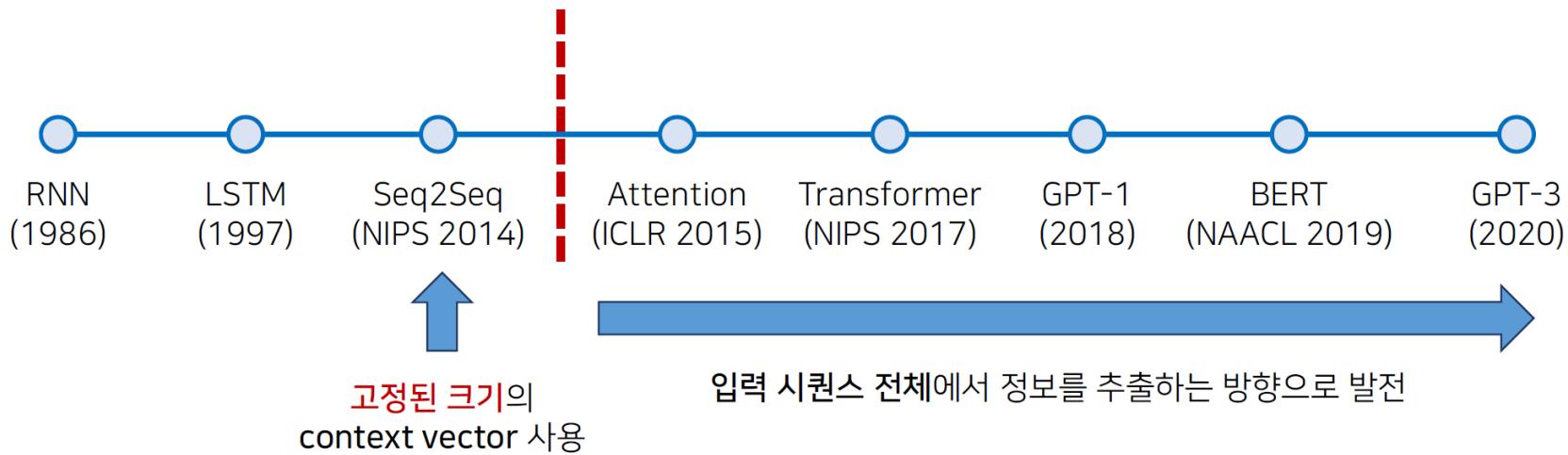
▶ Architecture



Seq2Seq



Seq2Seq



▶ 출처 : 동빈나

Transformer

▶ Attention Is All You Need

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

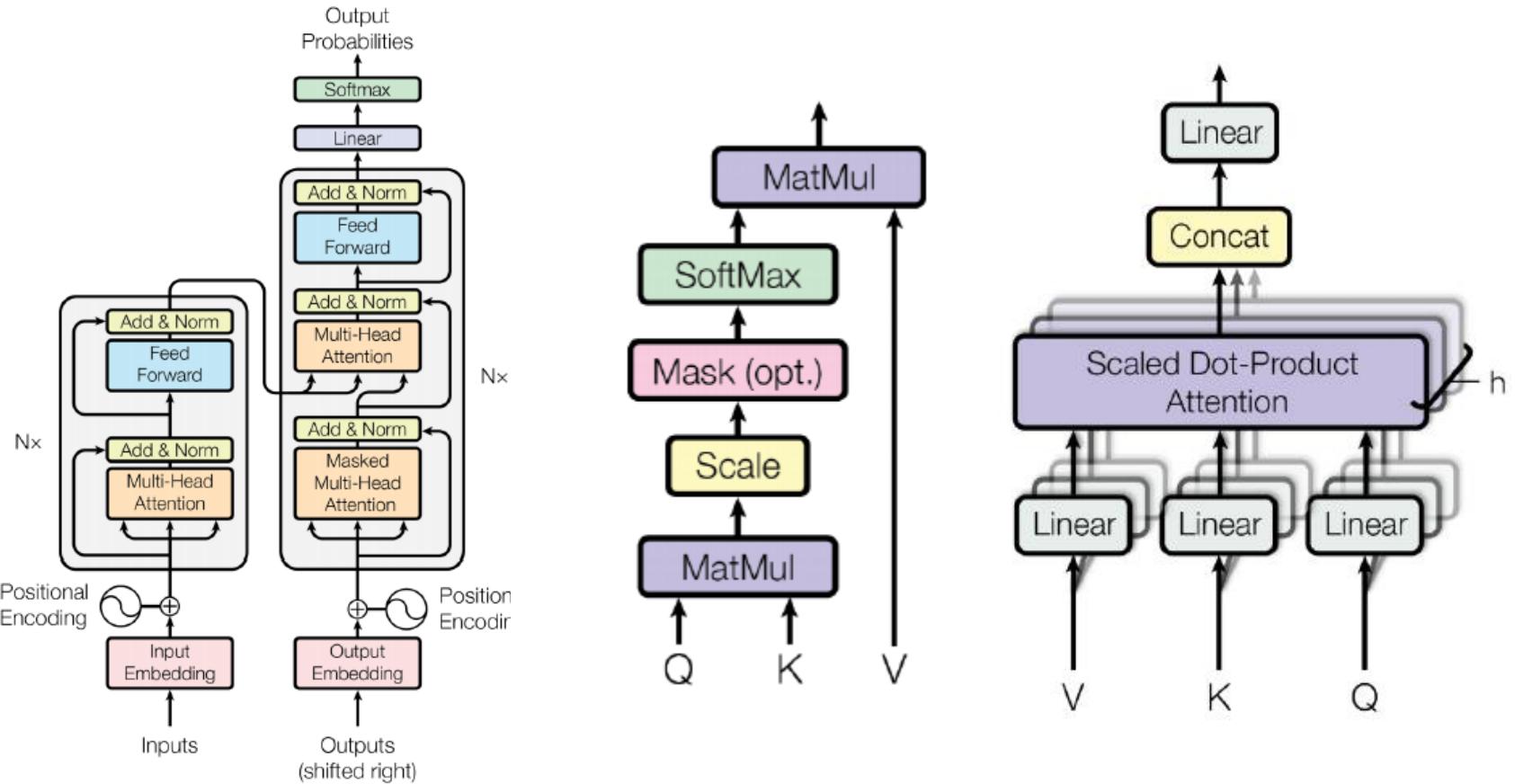
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

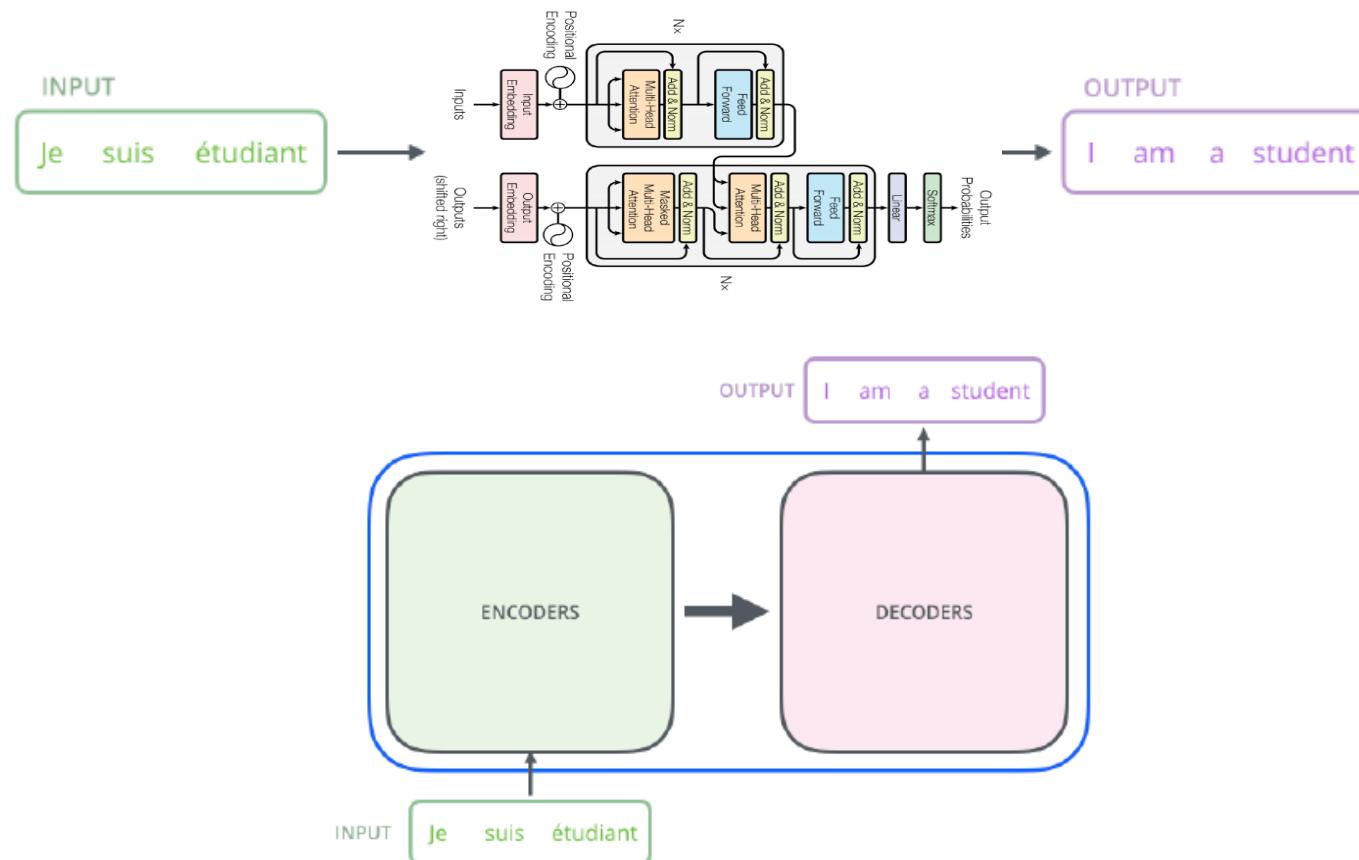
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Transformer

▶ Architecture

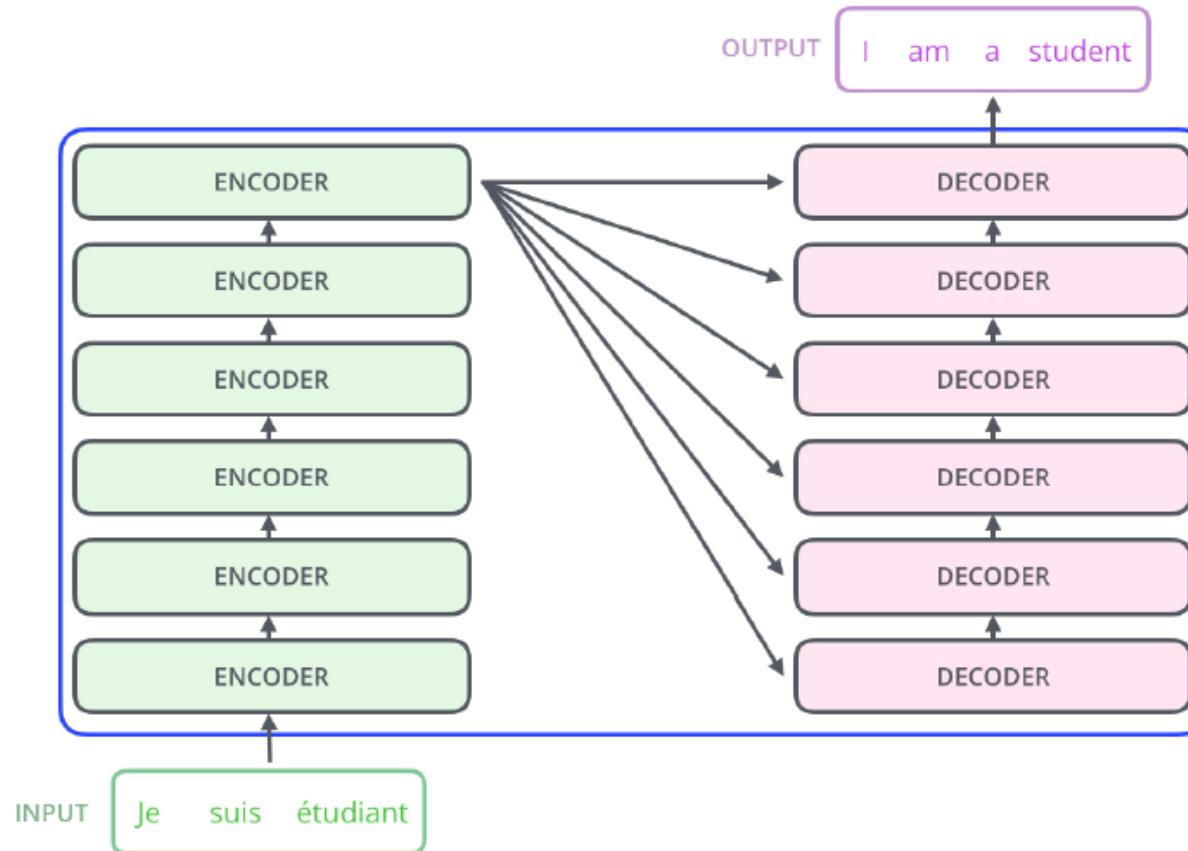


Transformer



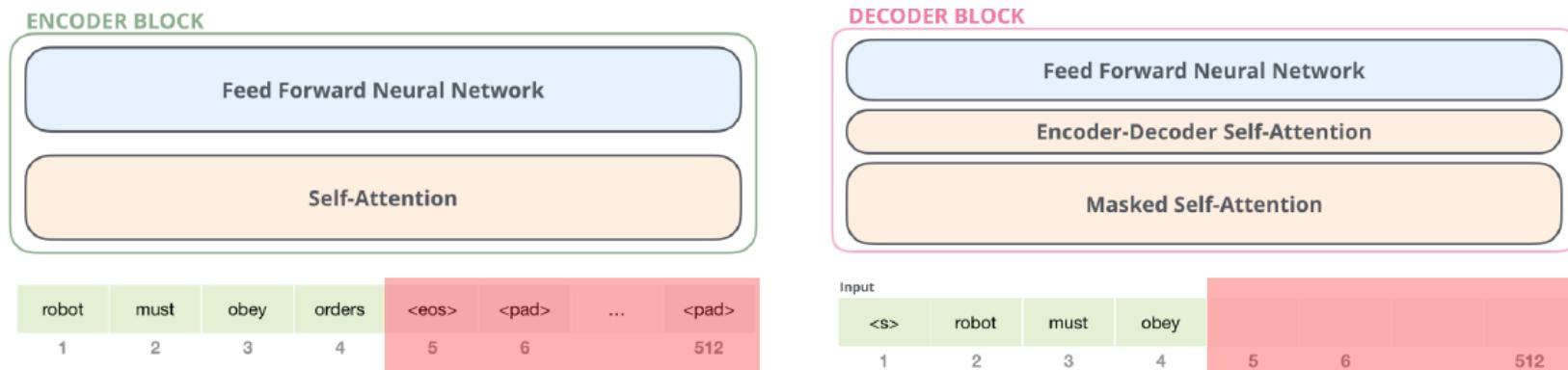
Transformer

▶ Architecture



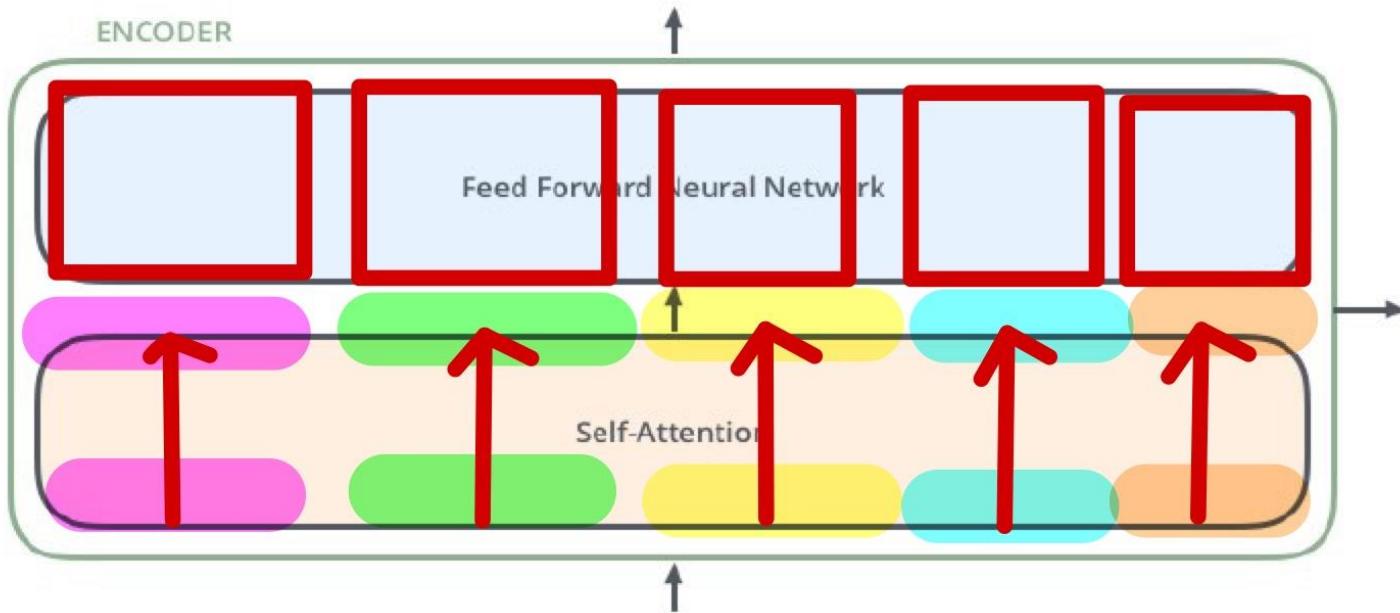
Transformer

▶ Architecture

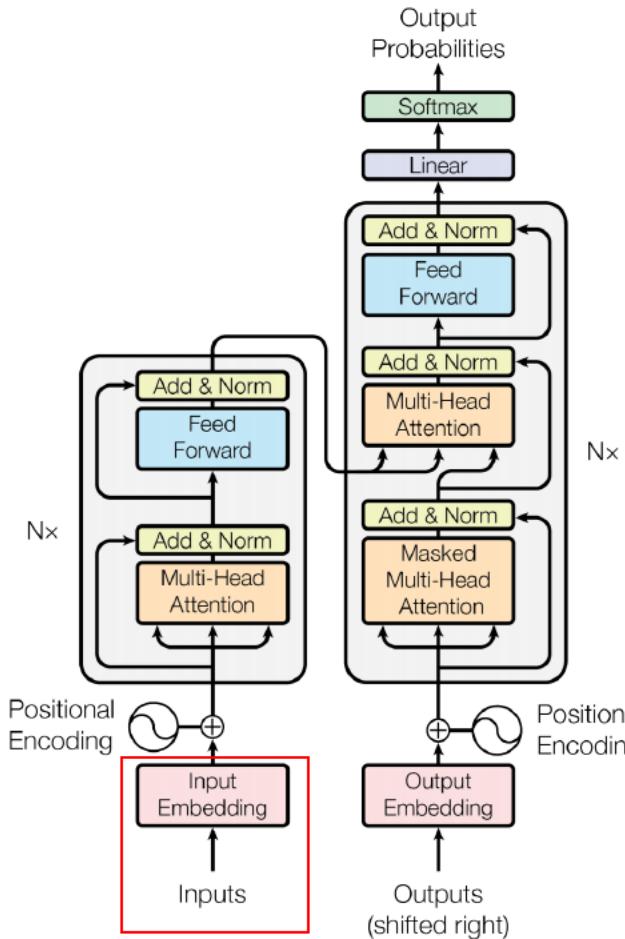


Transformer

▶ Encoder

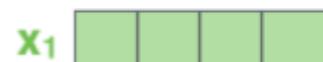


Transformer

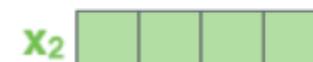


Transformer

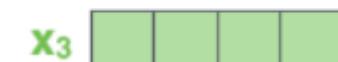
▶ Word Embedding



Je



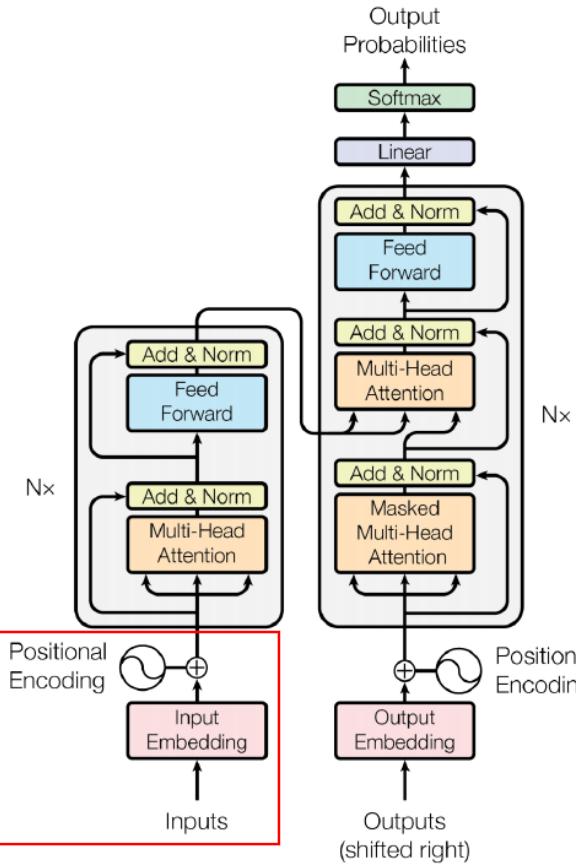
suis



étudiant

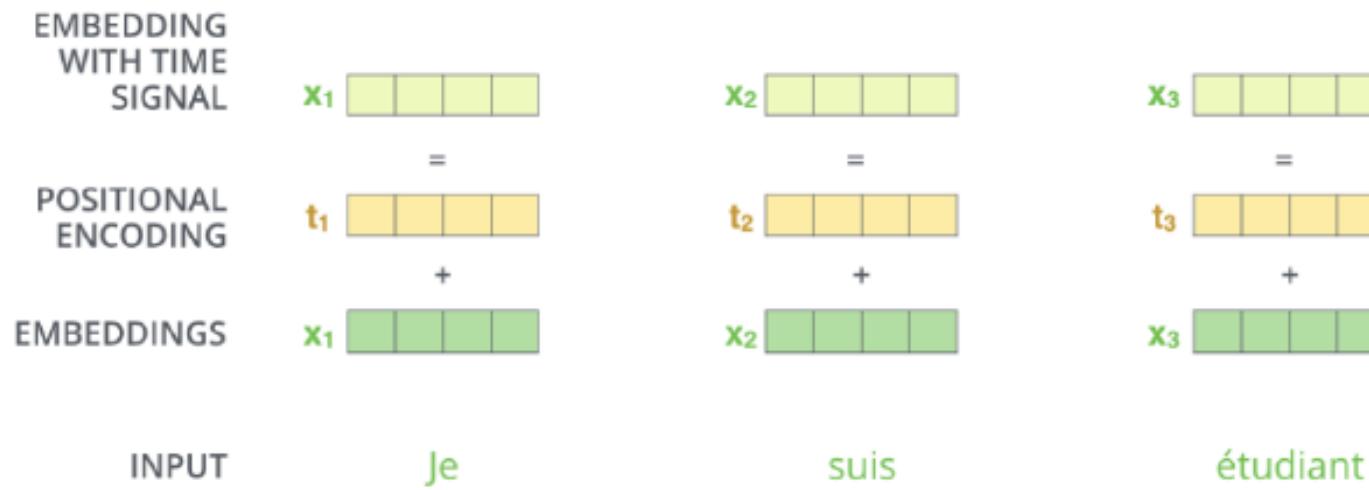
Transformer

▶ Positional Encoding



Transformer

▶ Positional Encoding



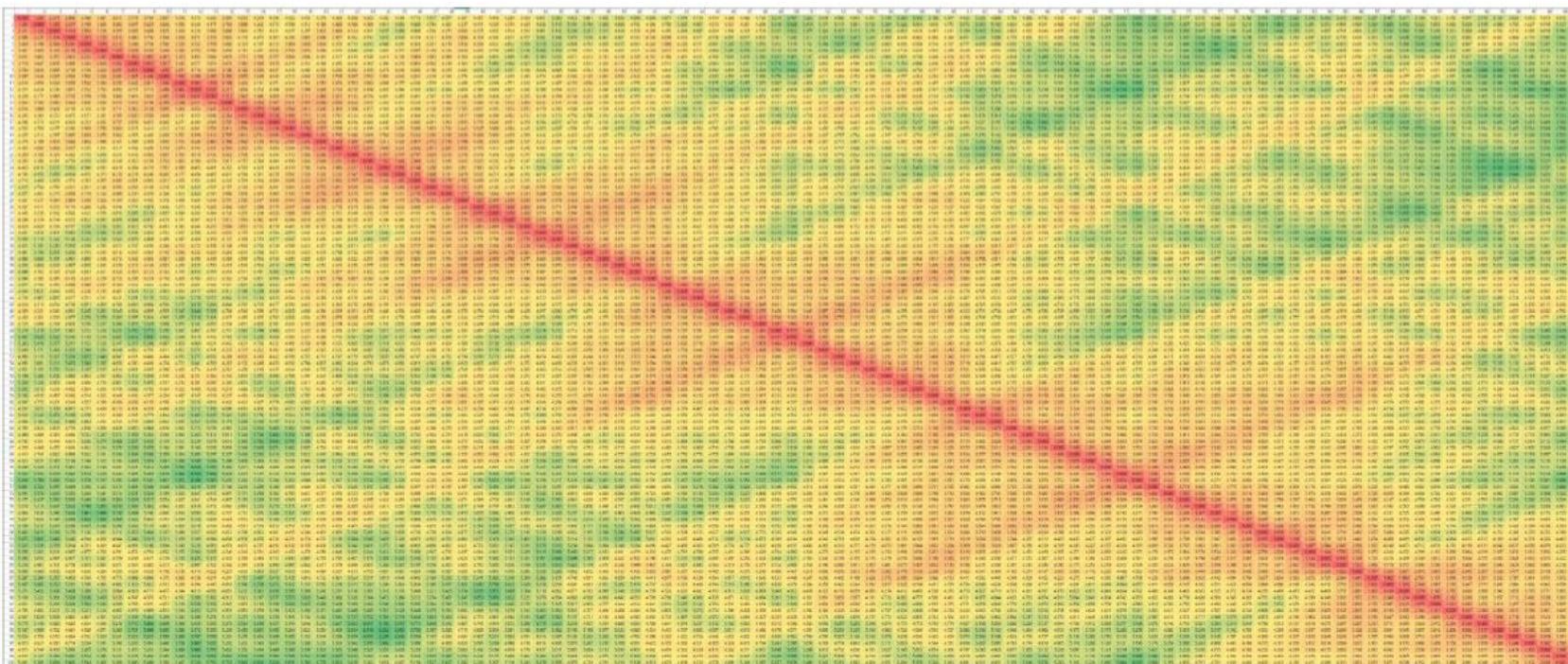
Transformer

▶ Positional Encoding

	X1	X2	X3	X4	X5	X6	X7	X8	X9	X10
X1	0.000	1.275	2.167	2.823	3.361	3.508	3.392	3.440	3.417	3.266
X2	1.275	0.000	1.104	2.195	3.135	3.511	3.452	3.442	3.387	3.308
X3	2.167	1.104	0.000	1.296	2.468	3.067	3.256	3.464	3.498	3.371
X4	2.823	2.195	1.296	0.000	1.275	2.110	2.746	3.399	3.624	3.399
X5	3.361	3.135	2.468	1.275	0.000	1.057	2.176	3.242	3.659	3.434
X6	3.508	3.511	3.067	2.110	1.057	0.000	1.333	2.601	3.169	3.118
X7	3.392	3.452	3.256	2.746	2.176	1.333	0.000	1.338	2.063	2.429
X8	3.440	3.442	3.464	3.399	3.242	2.601	1.338	0.000	0.912	1.891
X9	3.417	3.387	3.498	3.624	3.659	3.169	2.063	0.912	0.000	1.277
X10	3.266	3.308	3.371	3.399	3.434	3.118	2.429	1.891	1.277	0.000

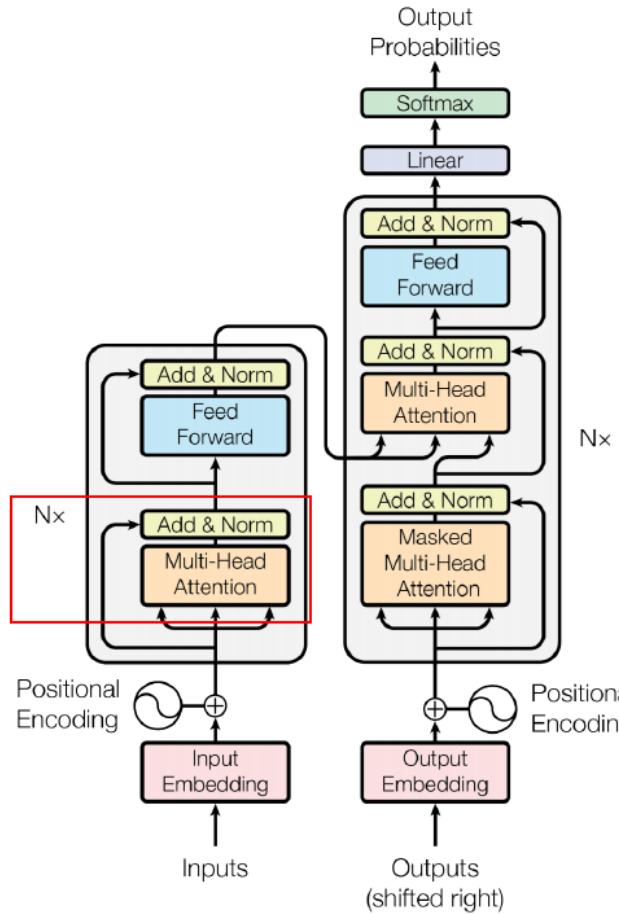
Transformer

▶ Positional Encoding



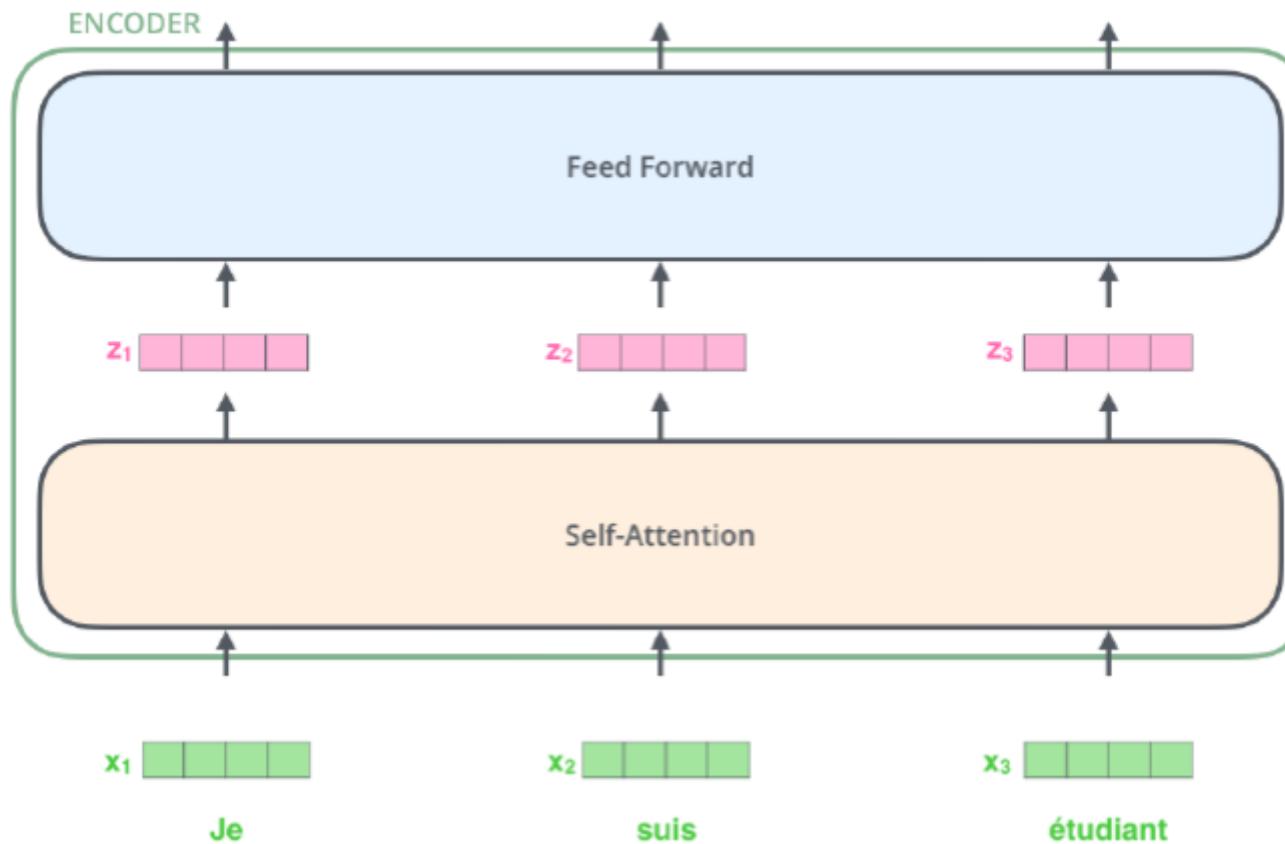
Transformer

▶ Multi-Head Attention



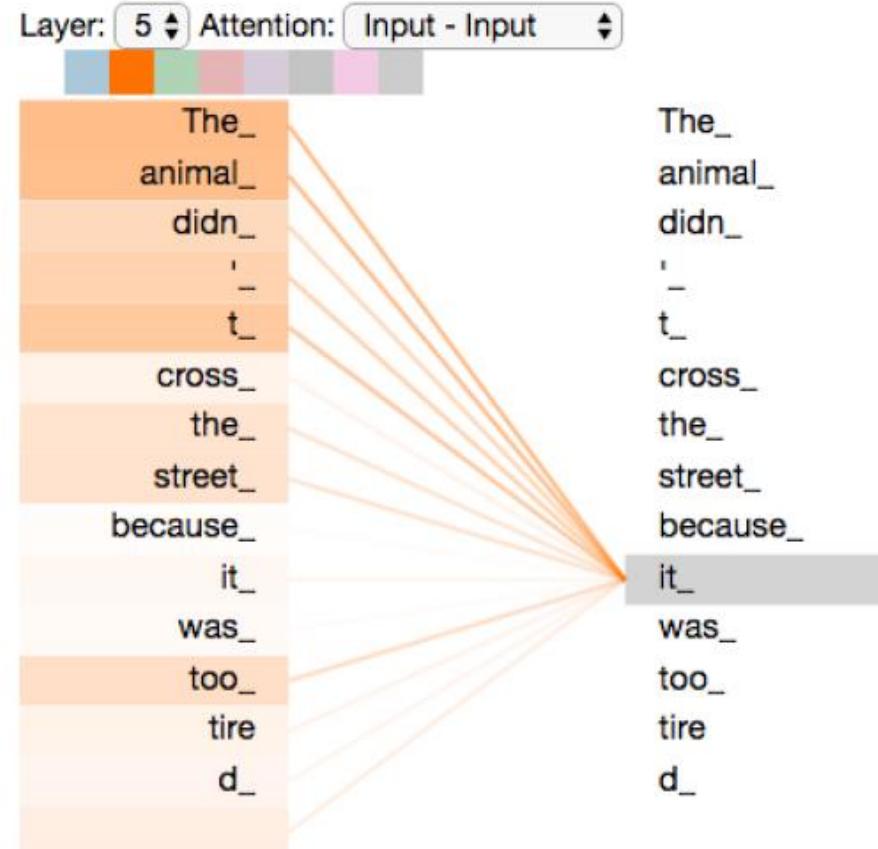
Transformer

▶ Encoder



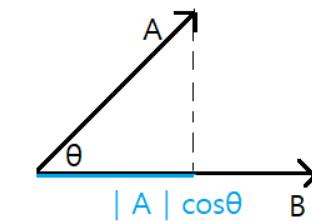
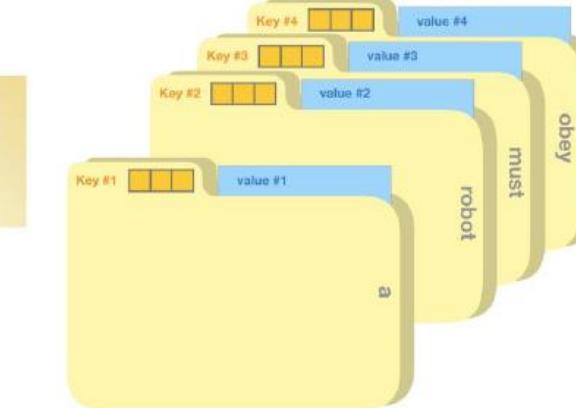
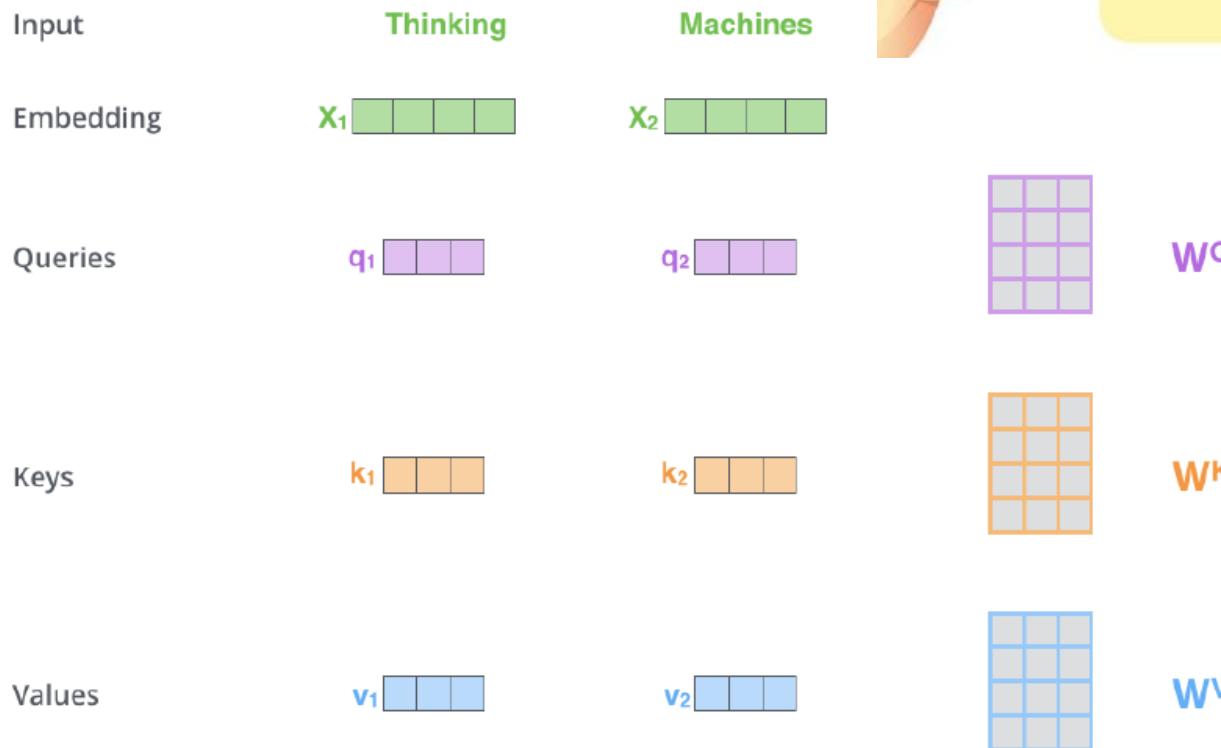
Transformer

▶ Encoder



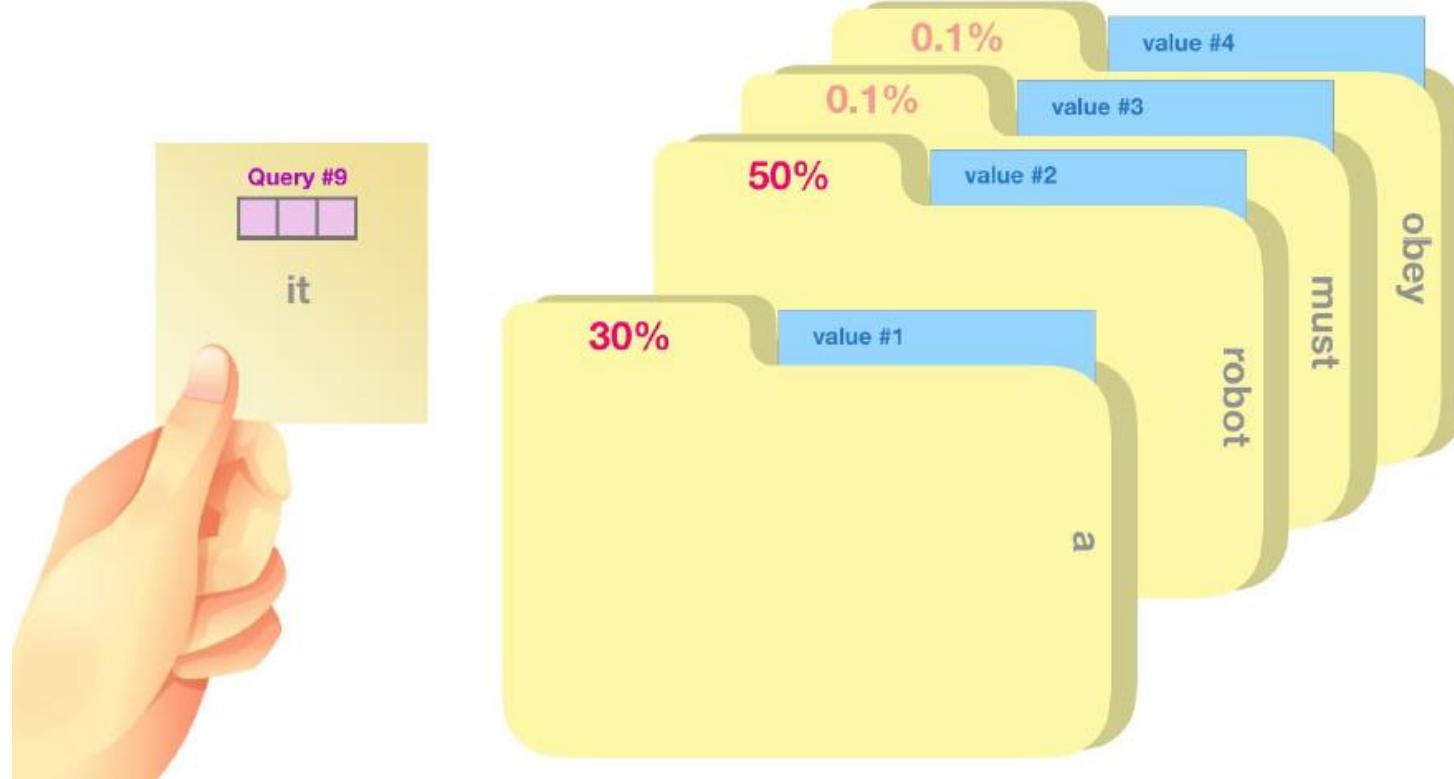
Transformer

▶ Attention



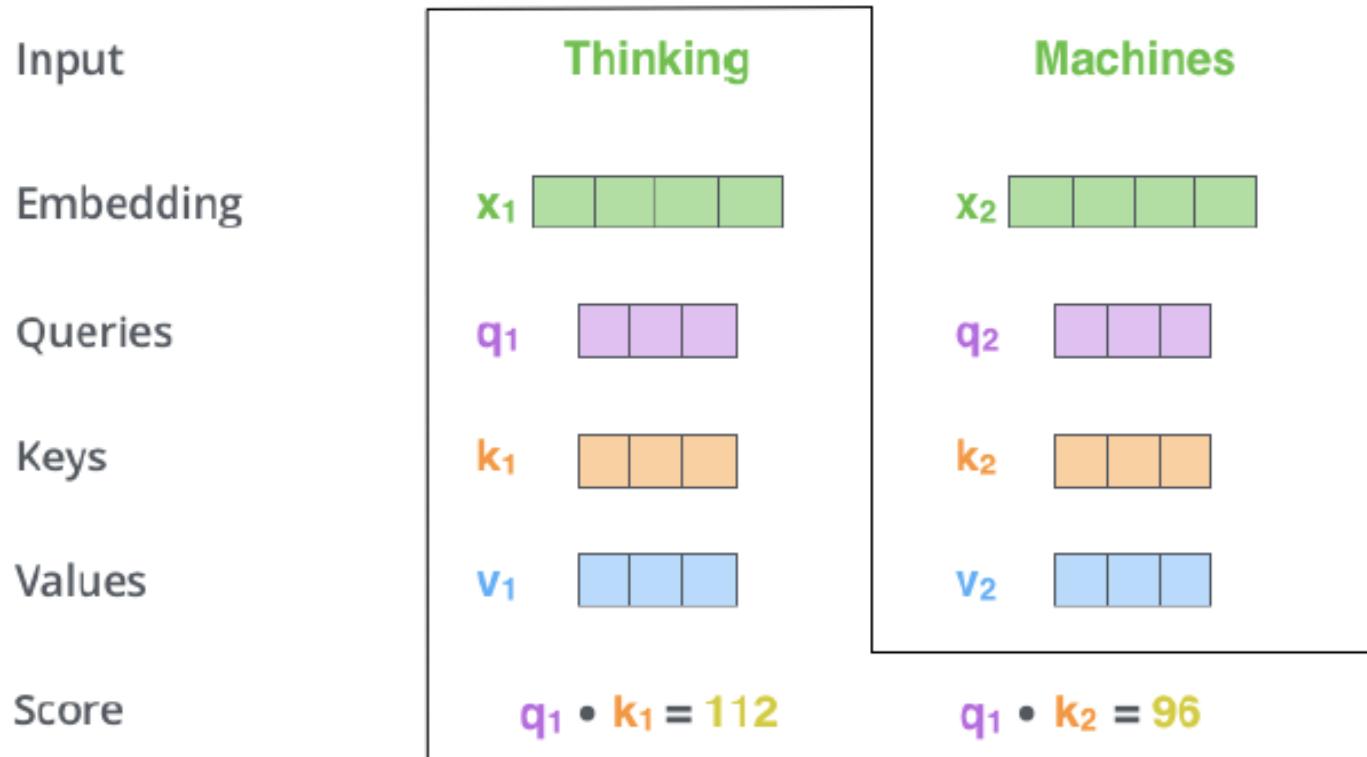
Transformer

▶ Attention



Transformer

▶ Attention



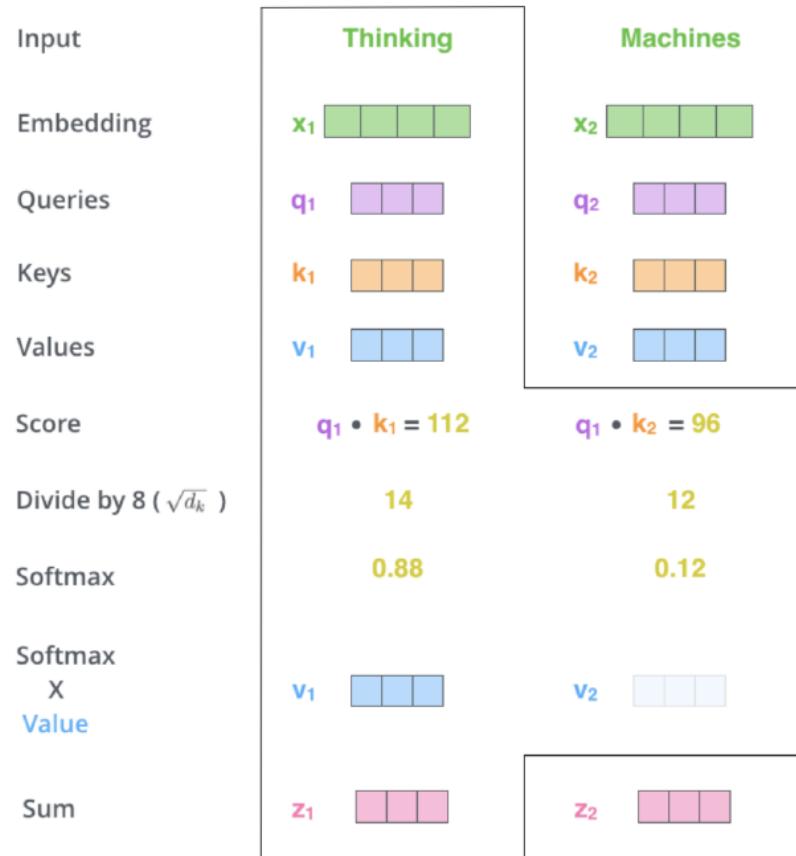
Transformer

▶ Attention

Input		
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

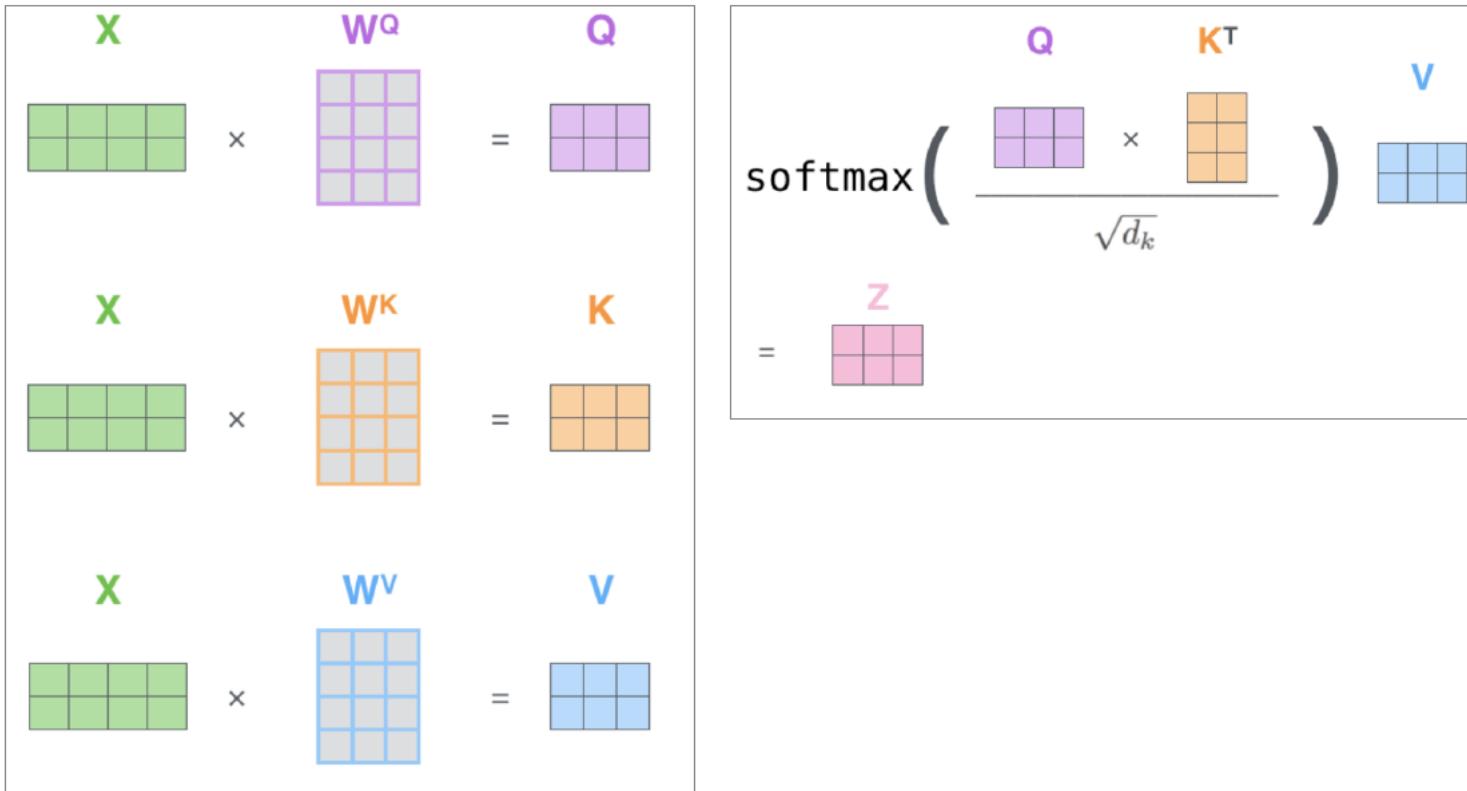
Transformer

▶ Attention



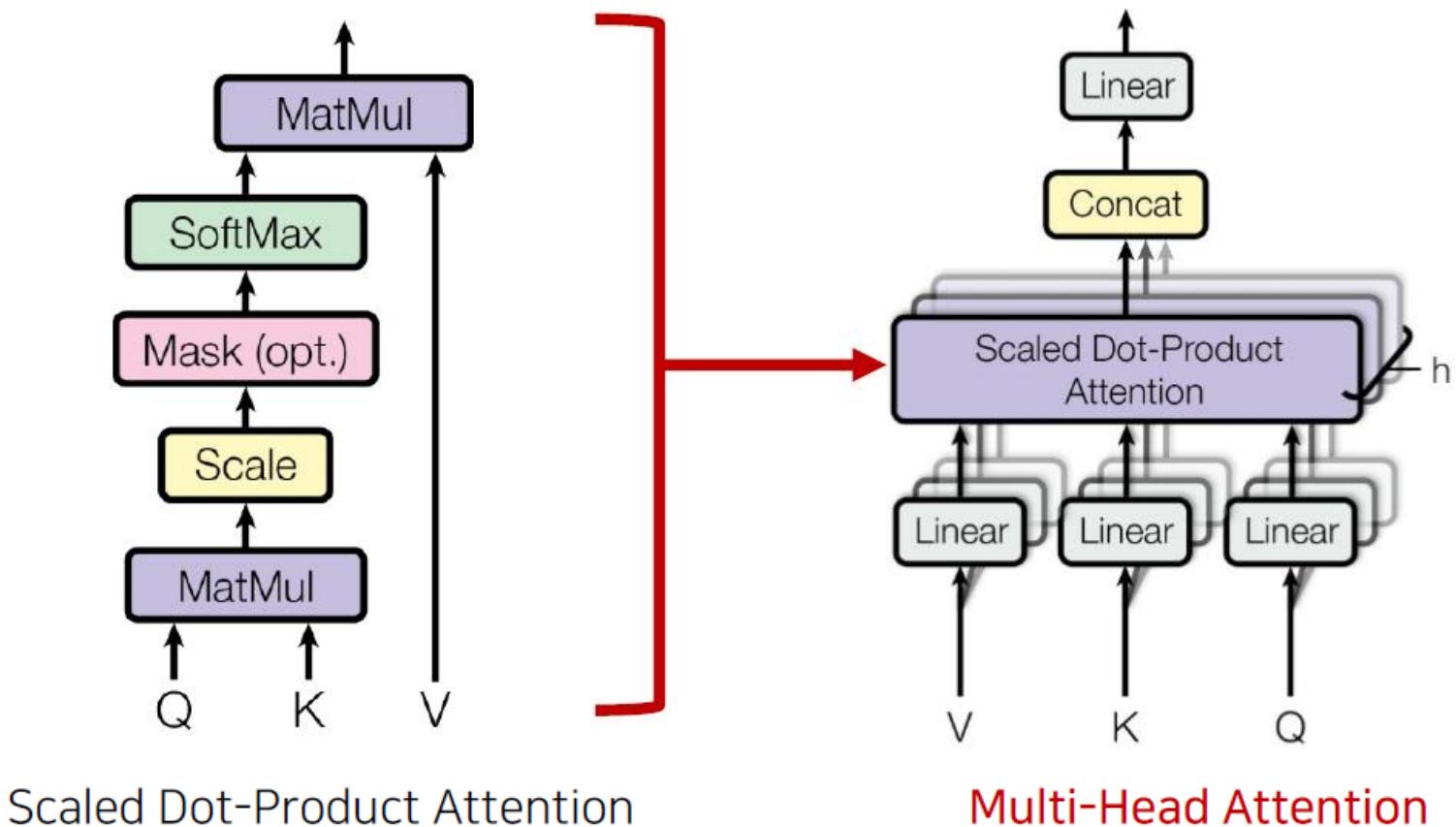
Transformer

▶ Attention



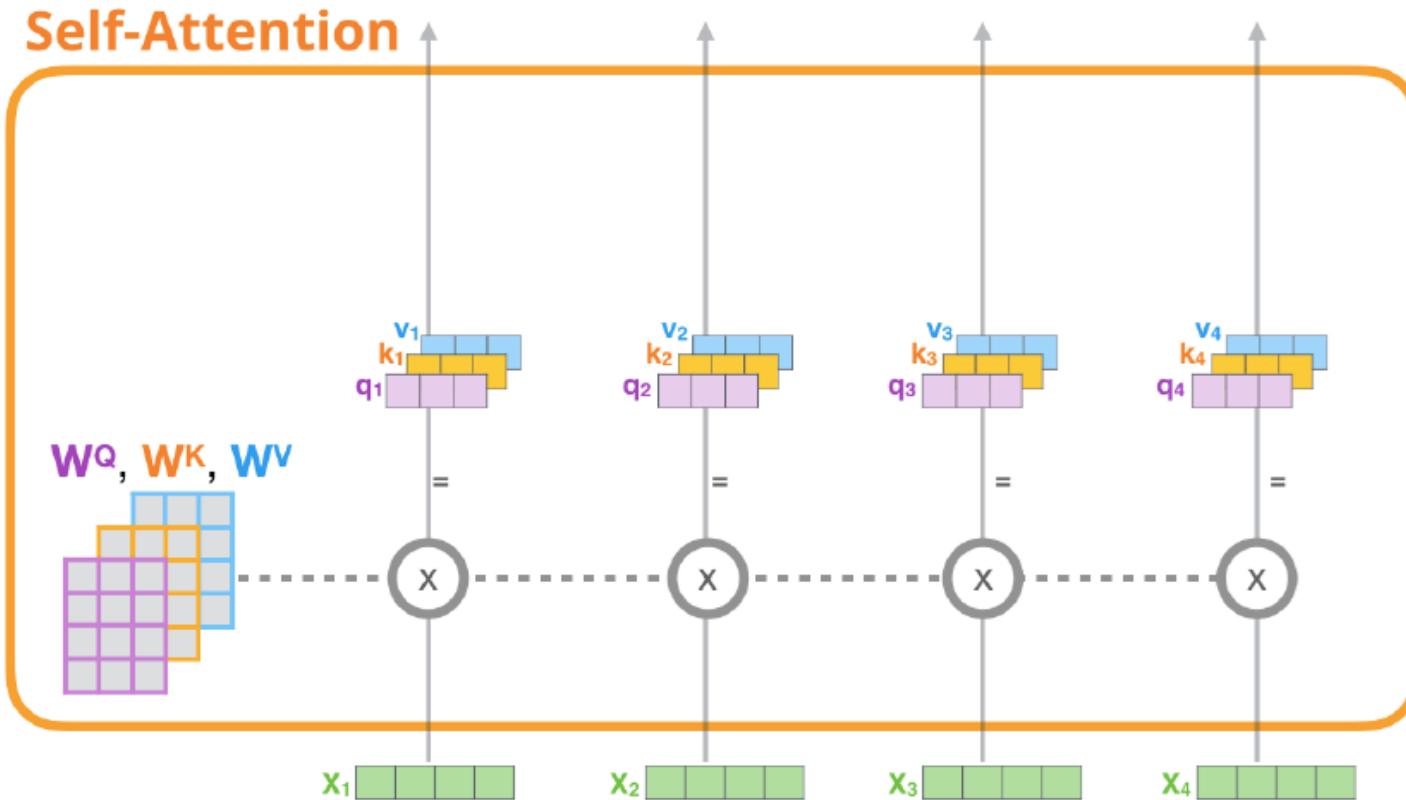
Transformer

▶ Attention



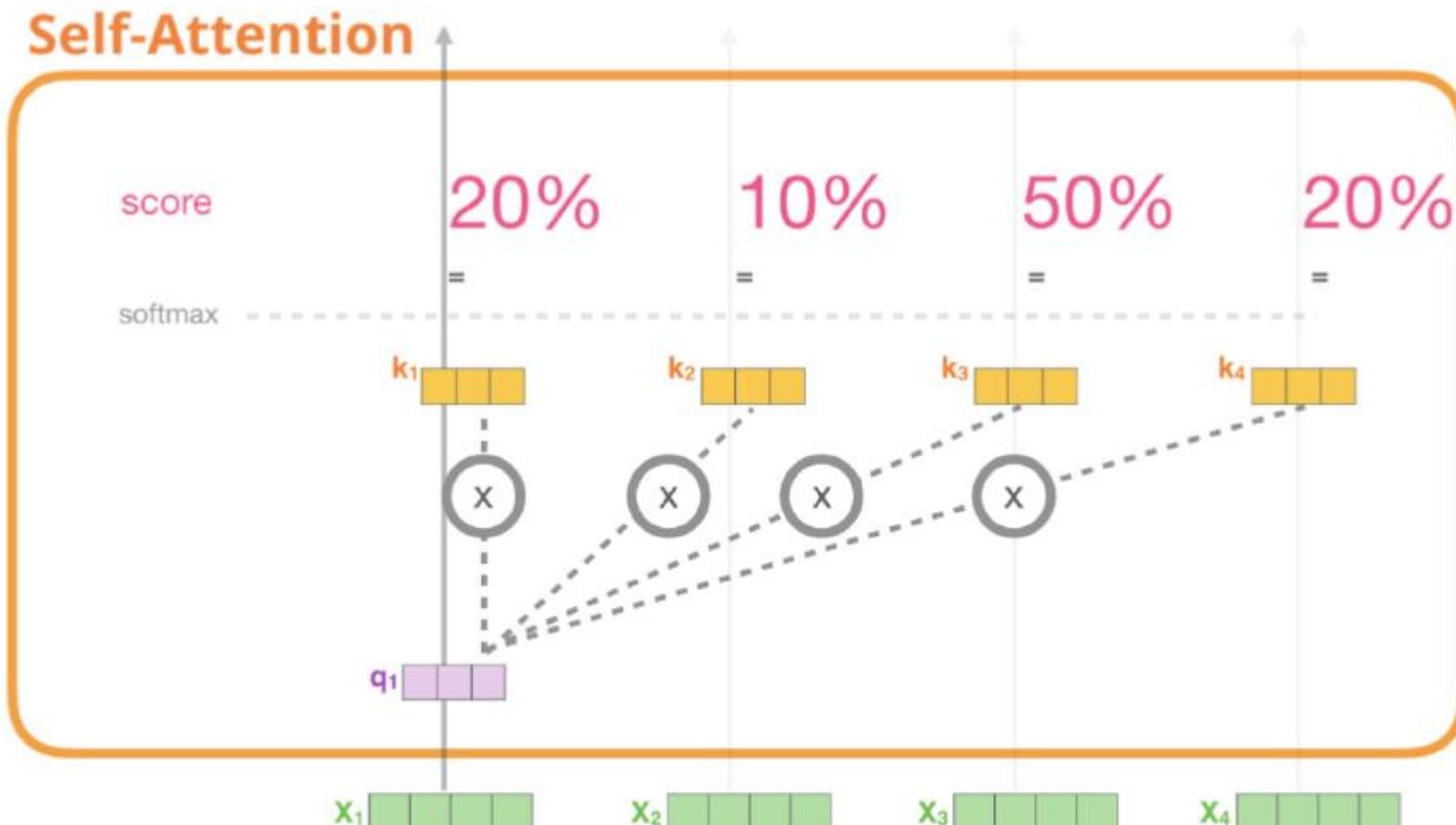
Transformer

▶ Attention



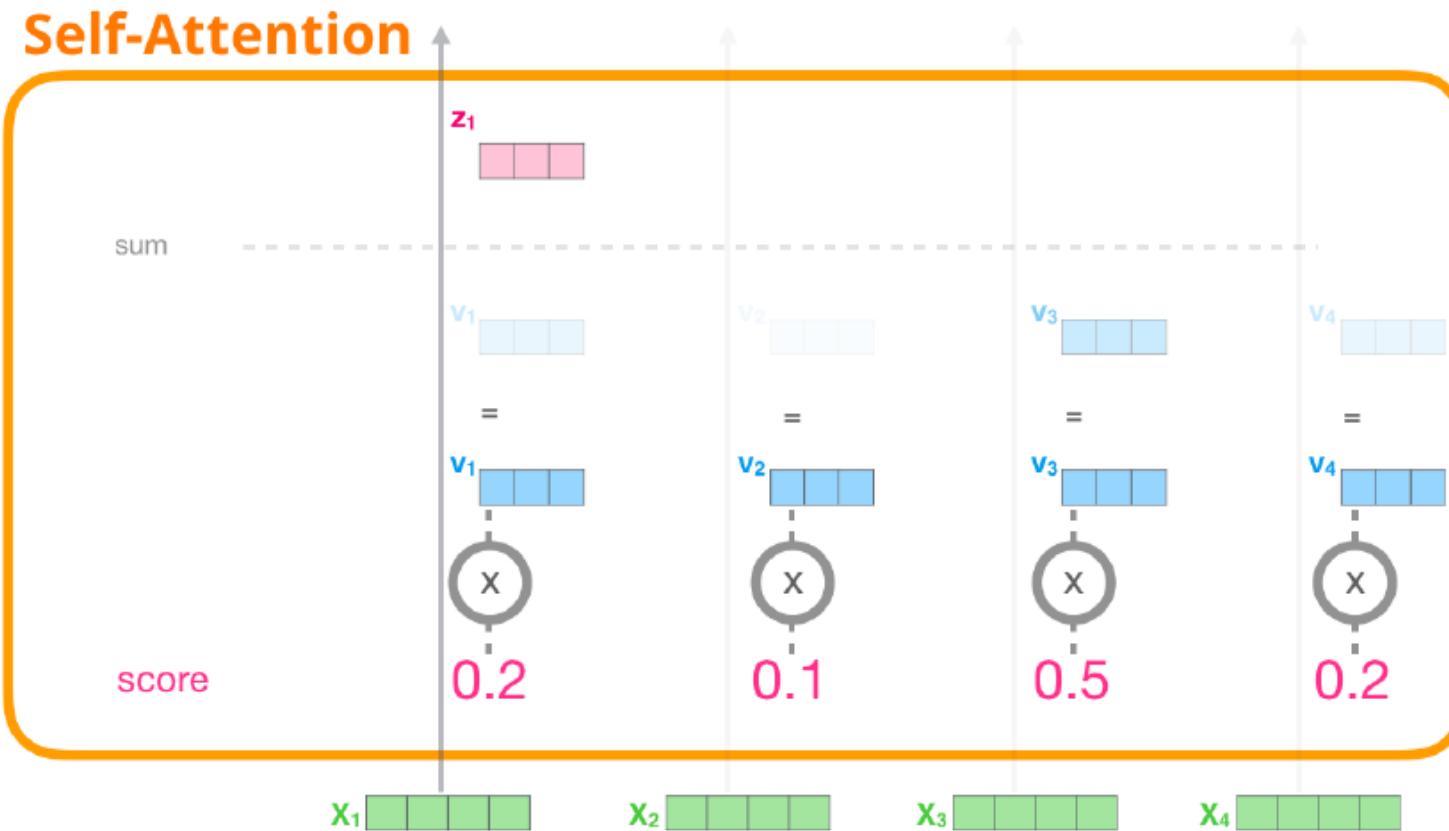
Transformer

▶ Attention



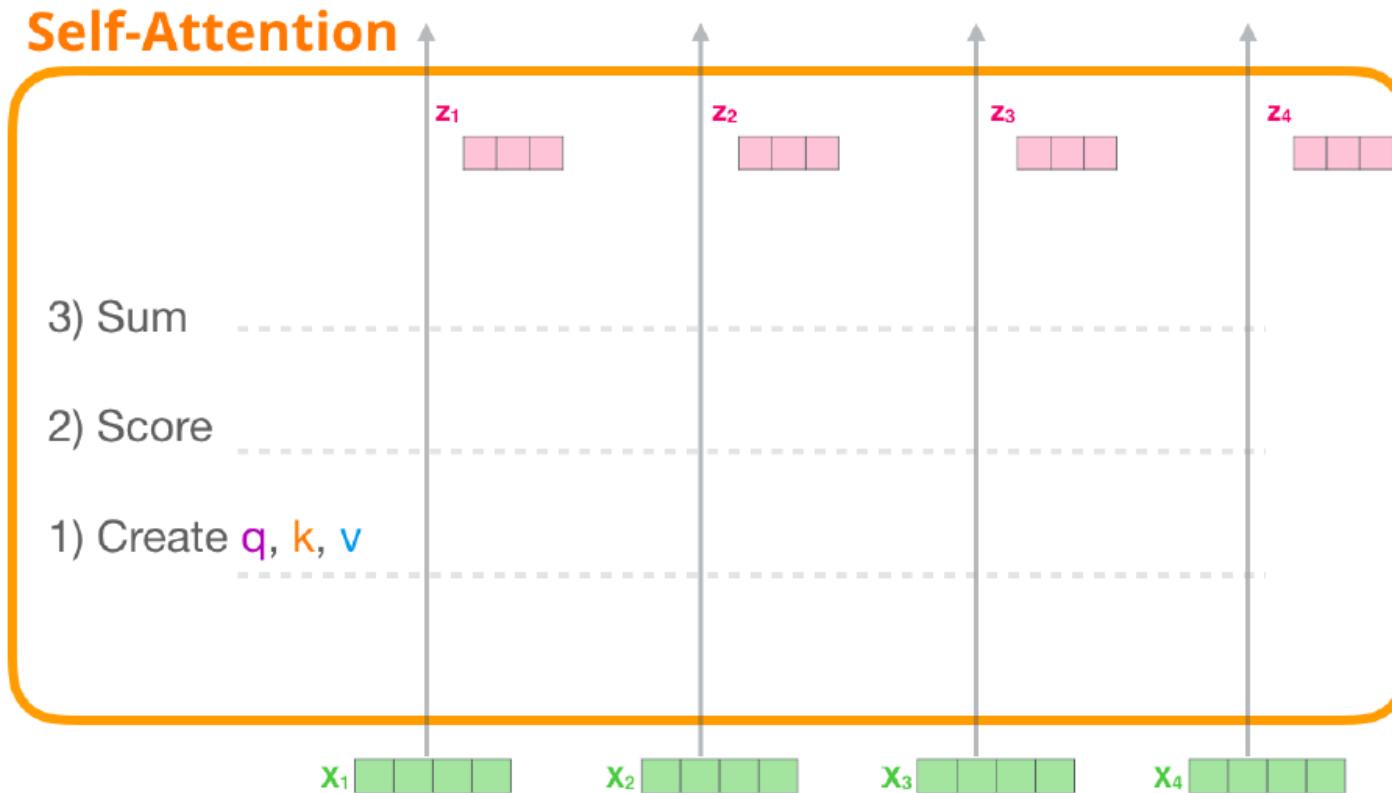
Transformer

▶ Attention



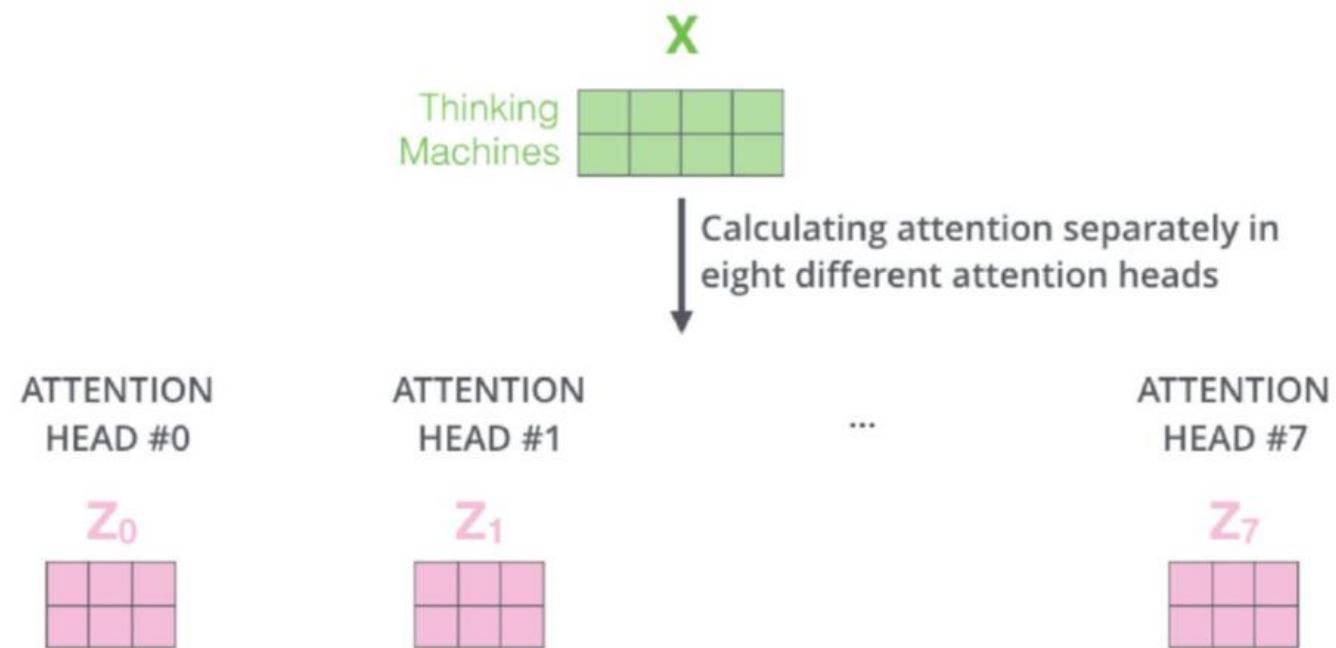
Transformer

▶ Attention



Transformer

▶ Attention



Transformer

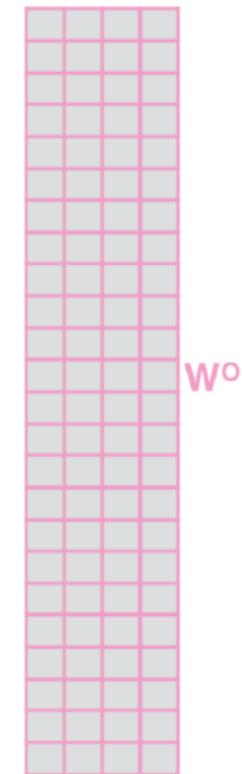
▶ Attention

1) Concatenate all the attention heads



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

X

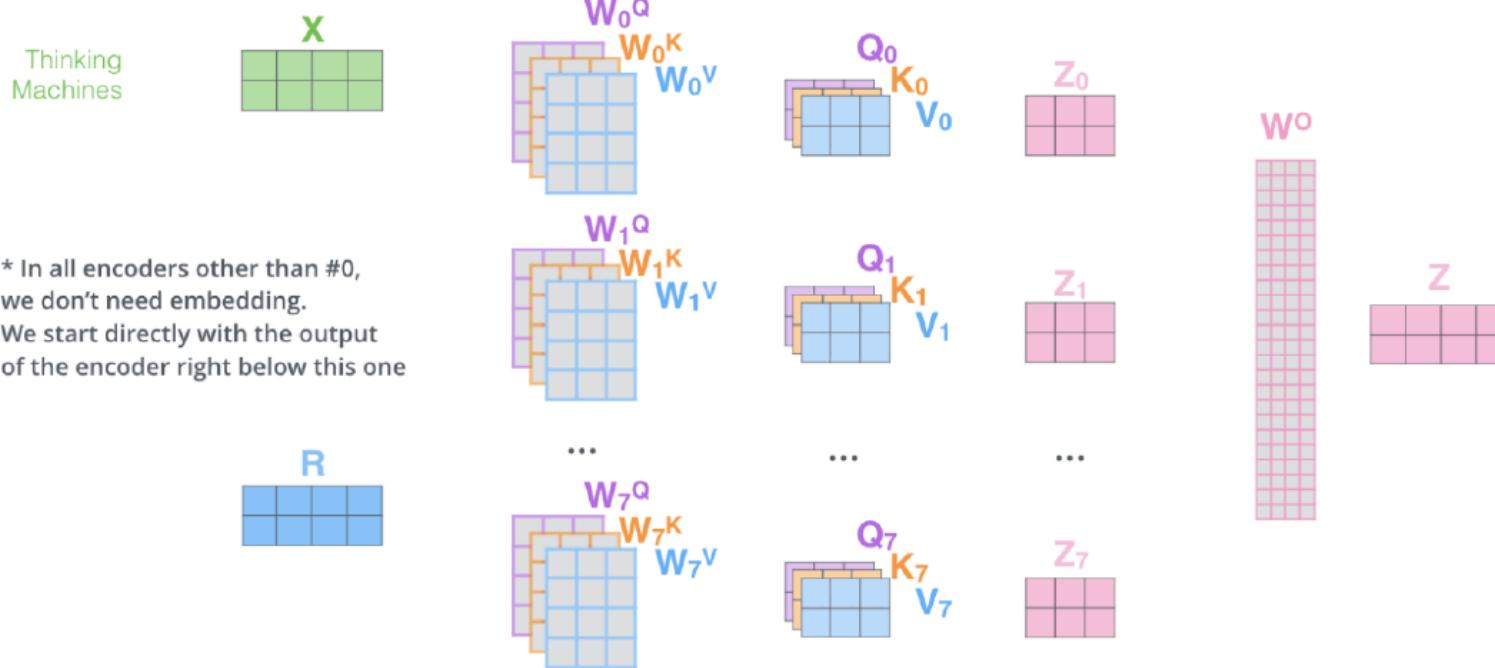


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \begin{matrix} \text{---} & \text{---} & \text{---} & \text{---} \end{matrix} \end{matrix}$$

Transformer

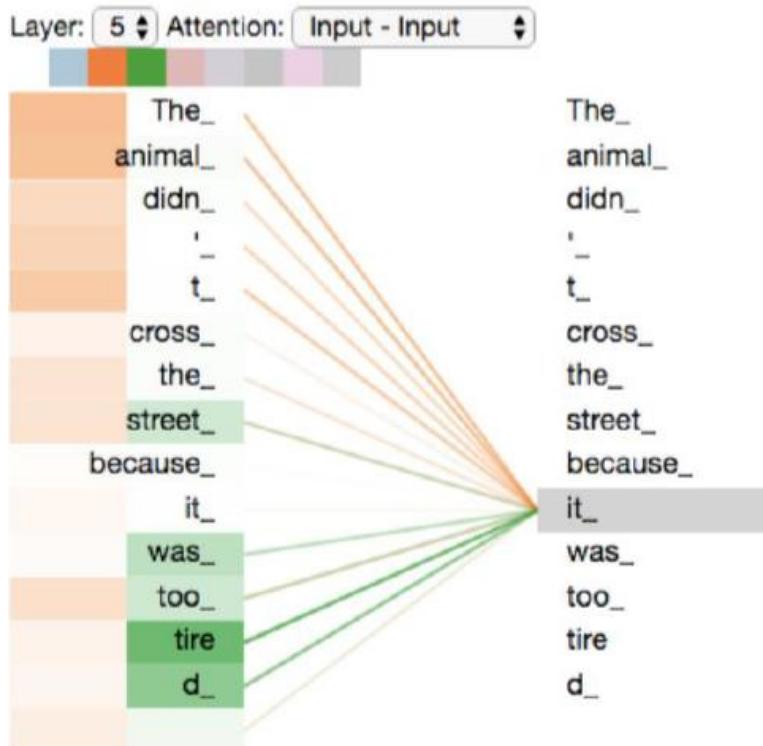
▶ Attention



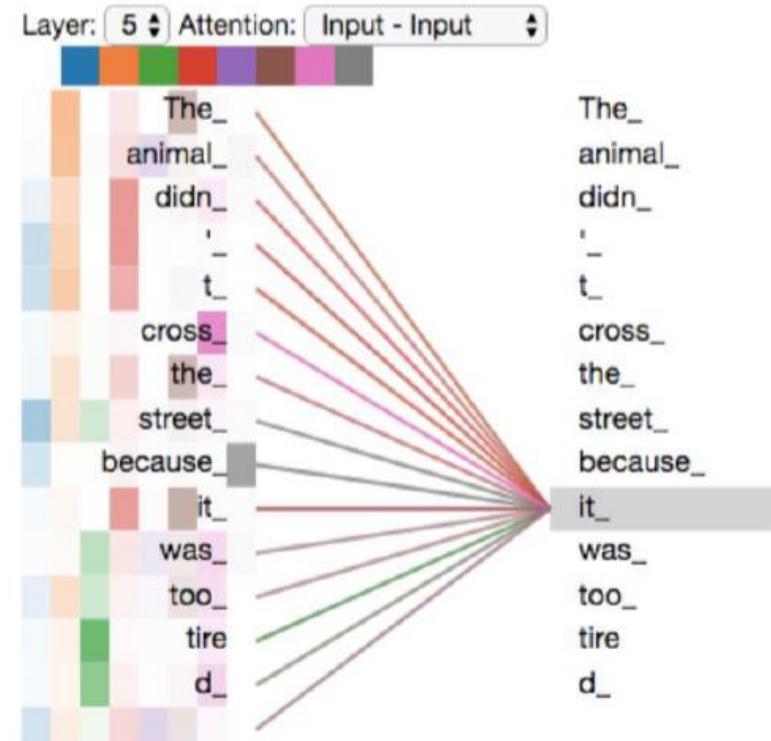
Transformer

▶ Multi Attention

Attention with two heads

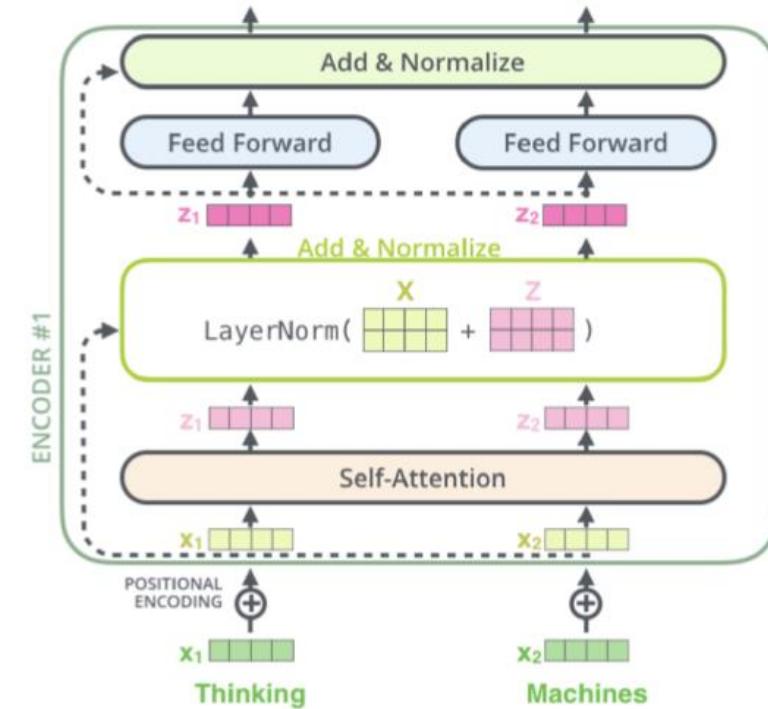
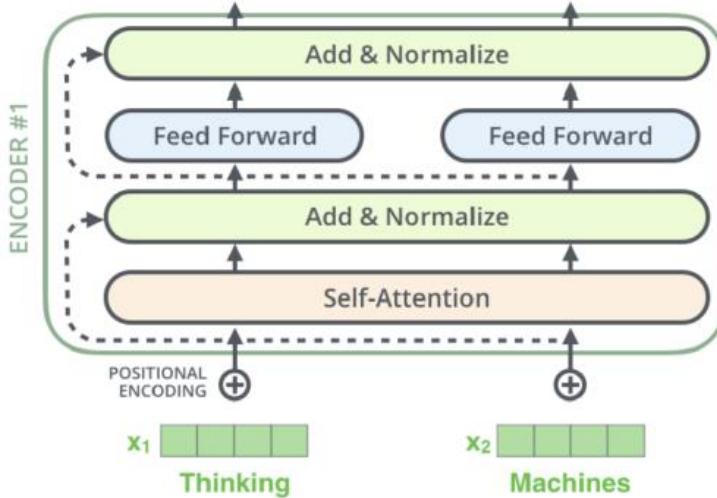


Attention with eight heads



Transformer

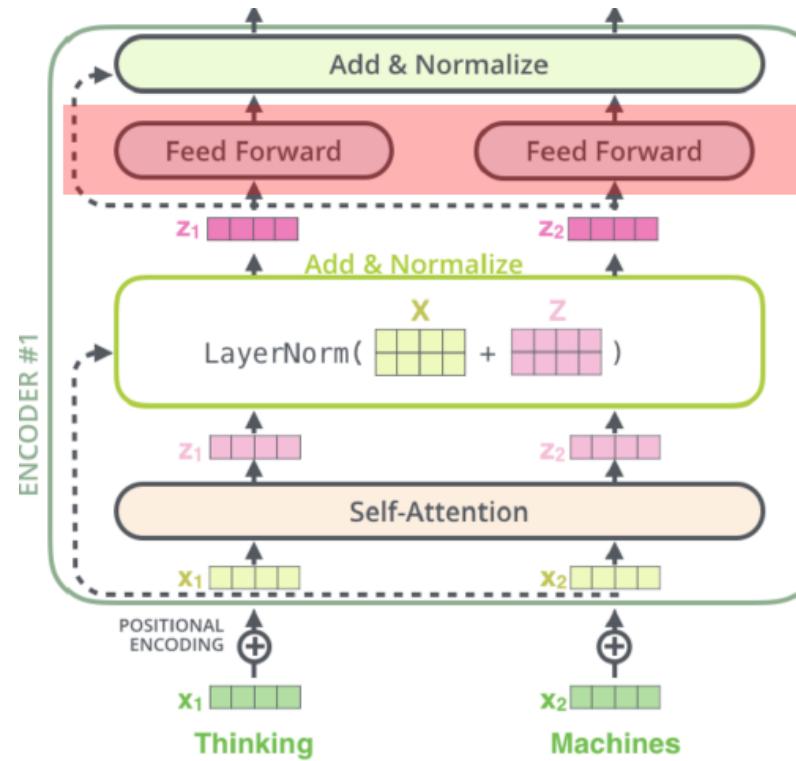
▶ Skip Connection and Normalization



Transformer

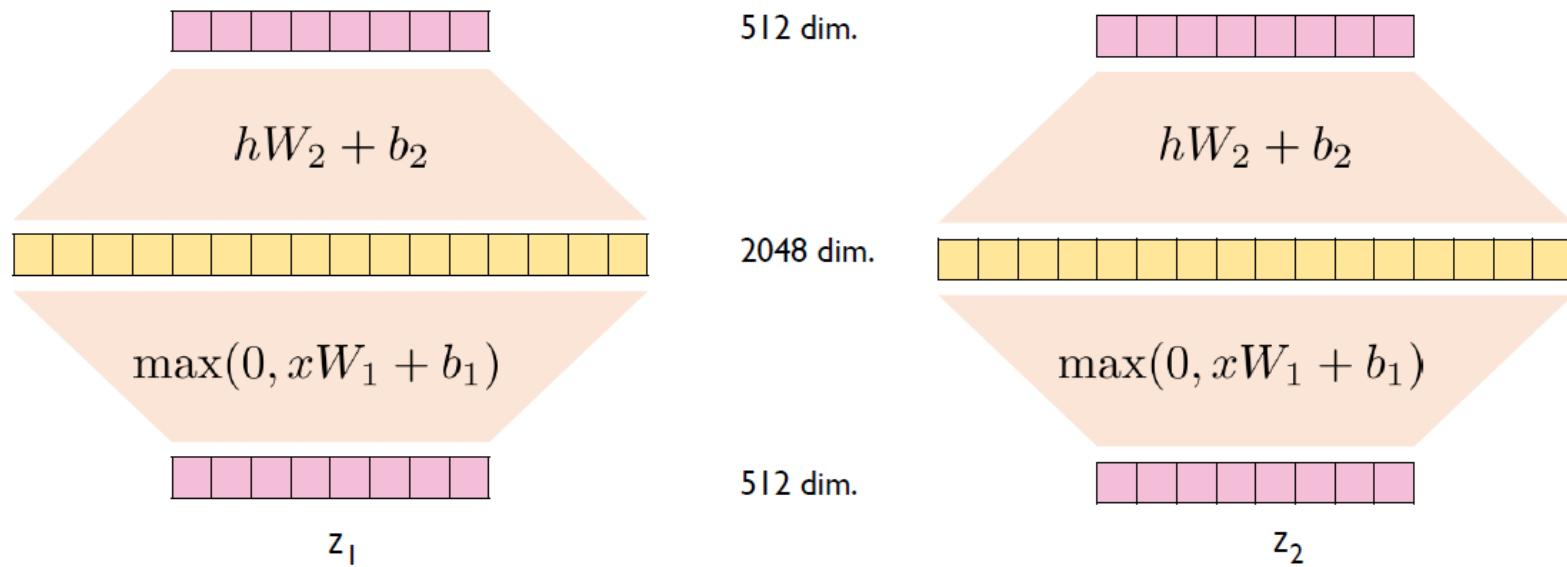
▶ Feed Forward Neural Network

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



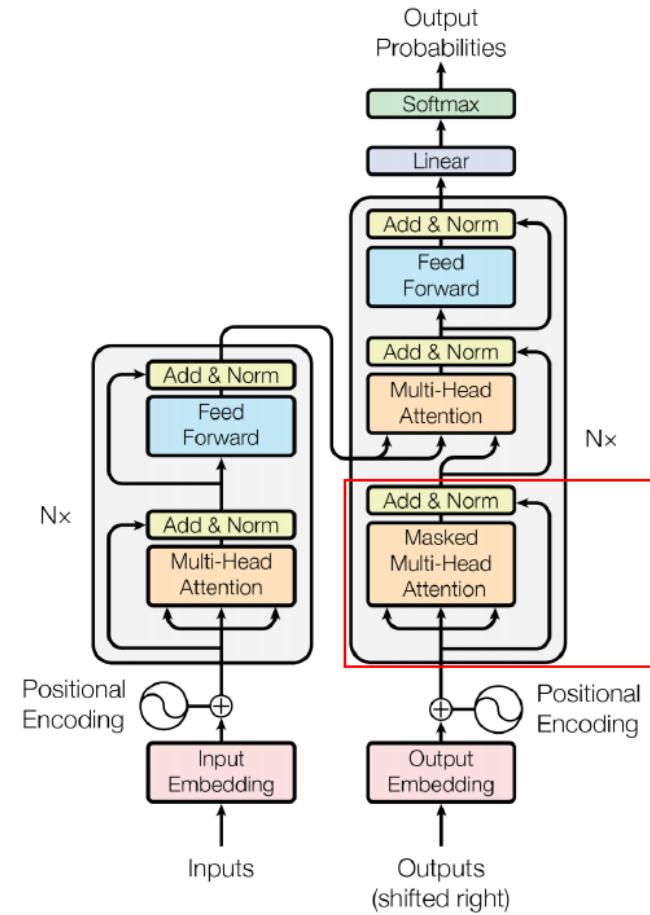
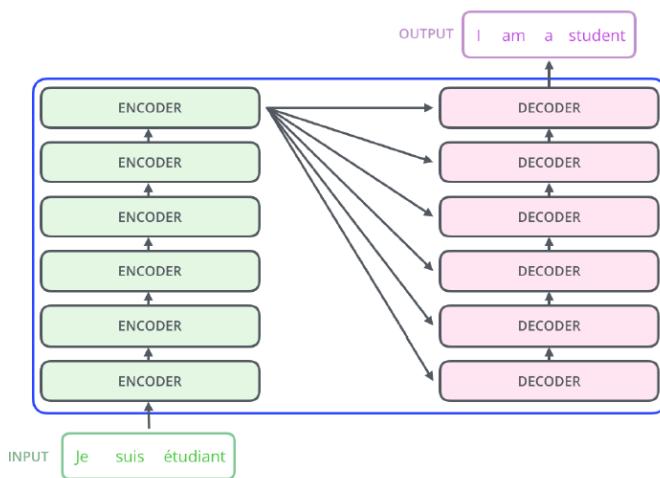
Transformer

▶ Feed Forward Neural Network



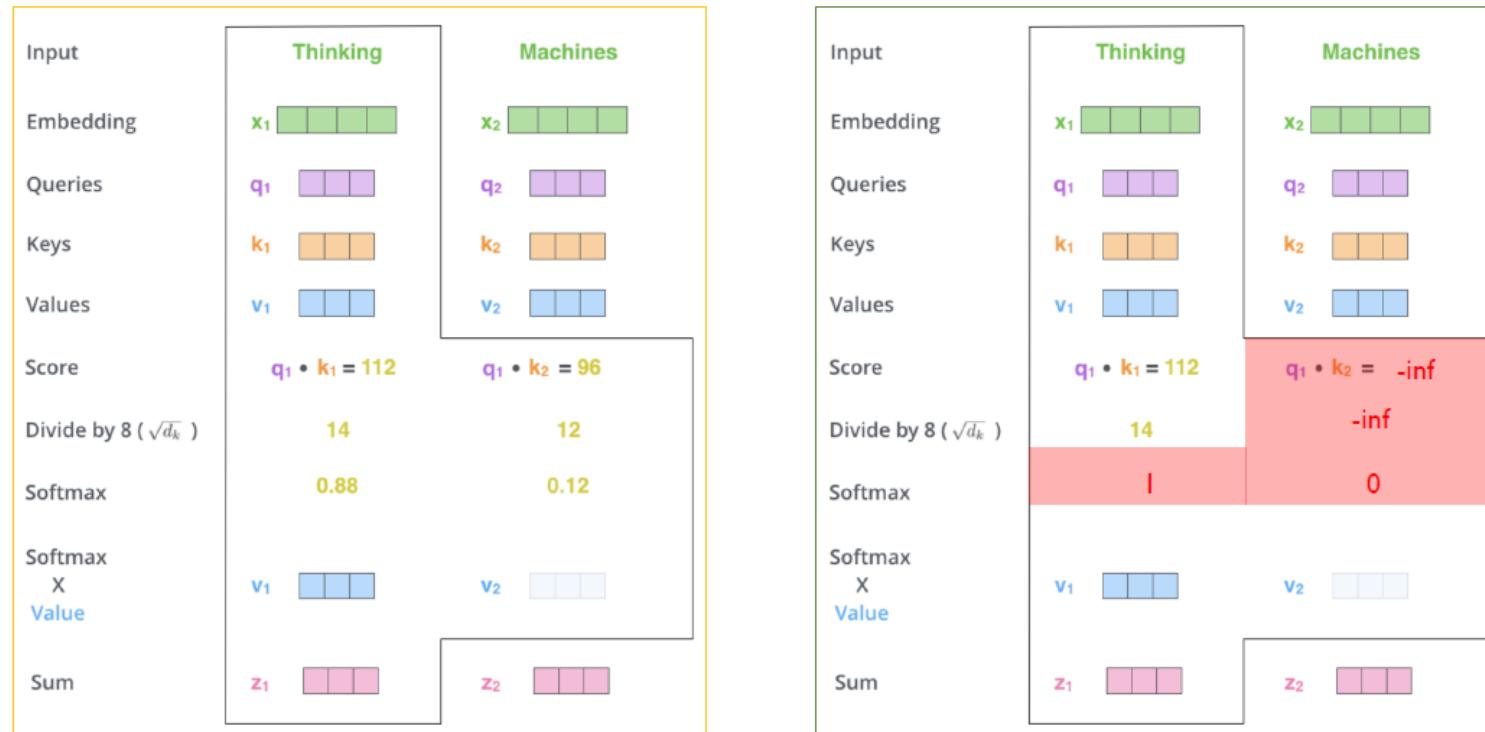
Transformer

▶ Masked Multi Head Attention



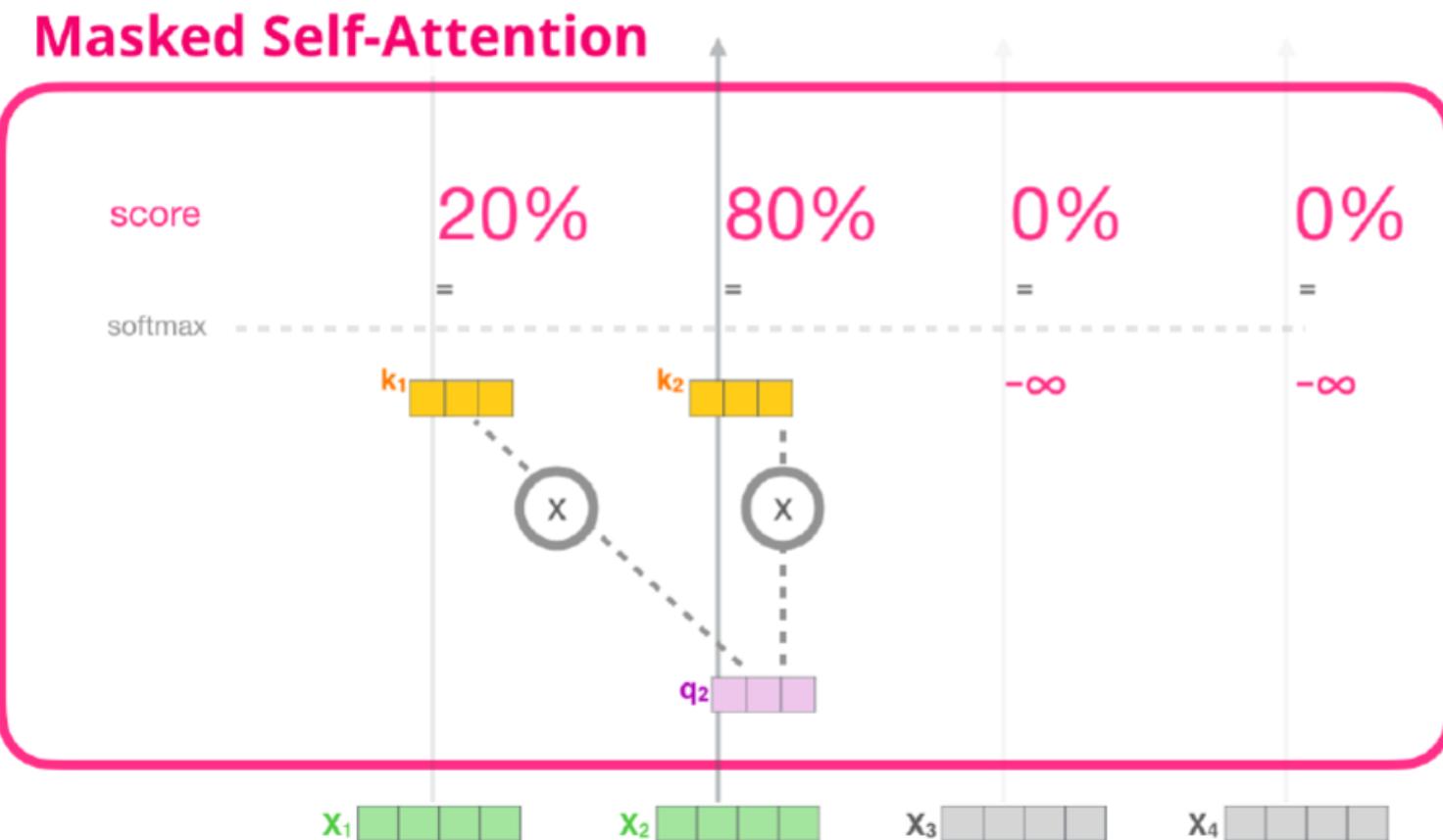
Transformer

▶ Masked Multi Head Attention



Transformer

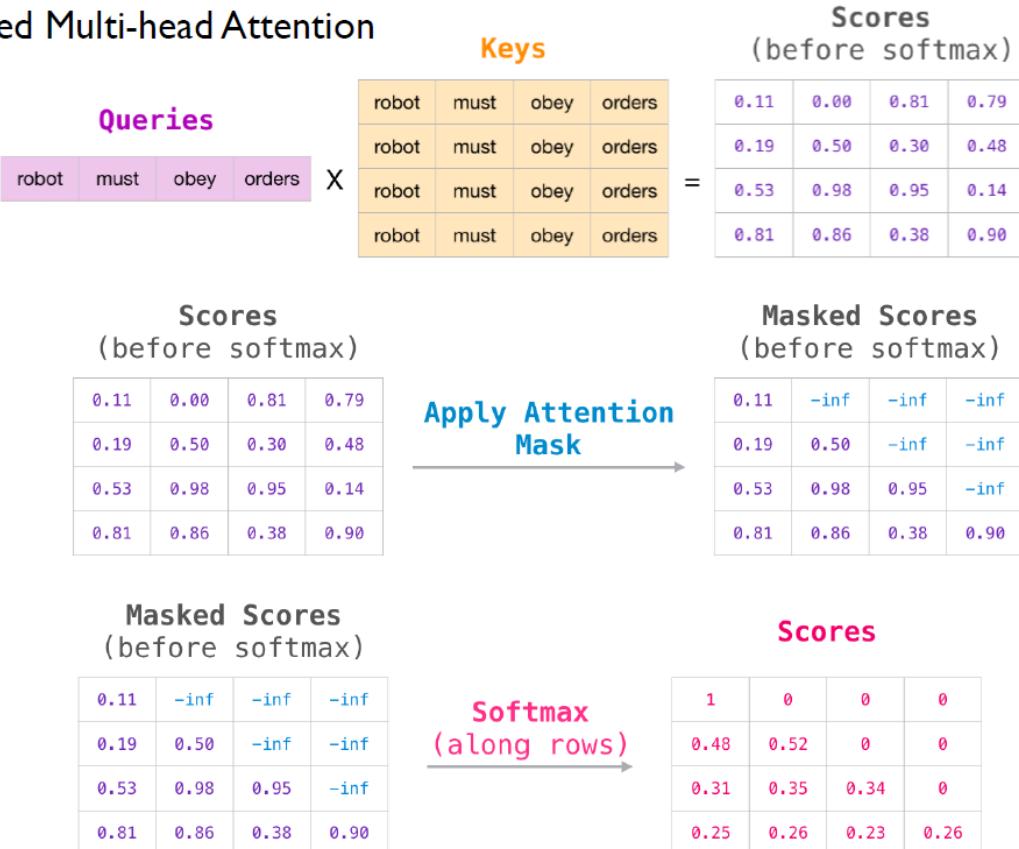
▶ Masked Multi Head Attention



Transformer

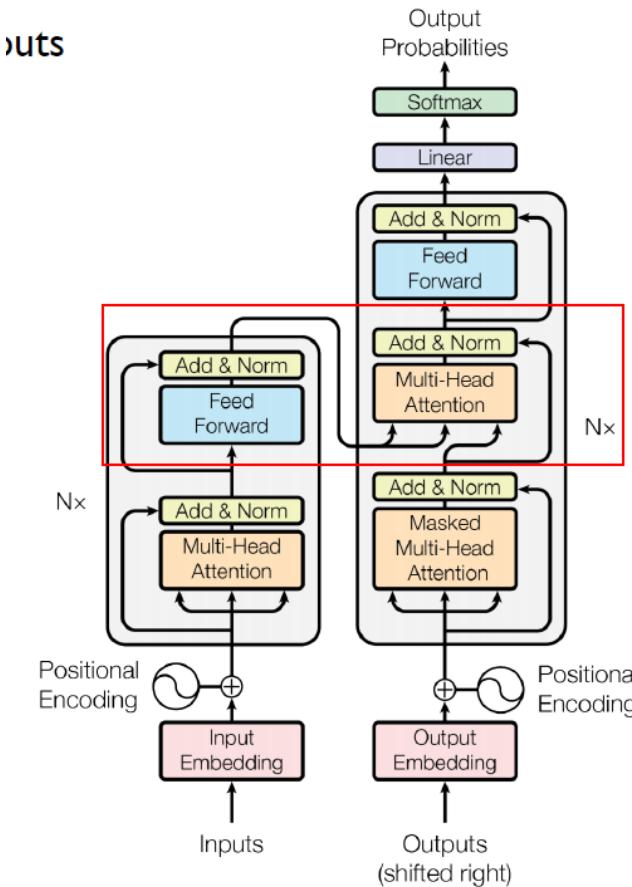
▶ Masked Multi Head Attention

- Masked Multi-head Attention

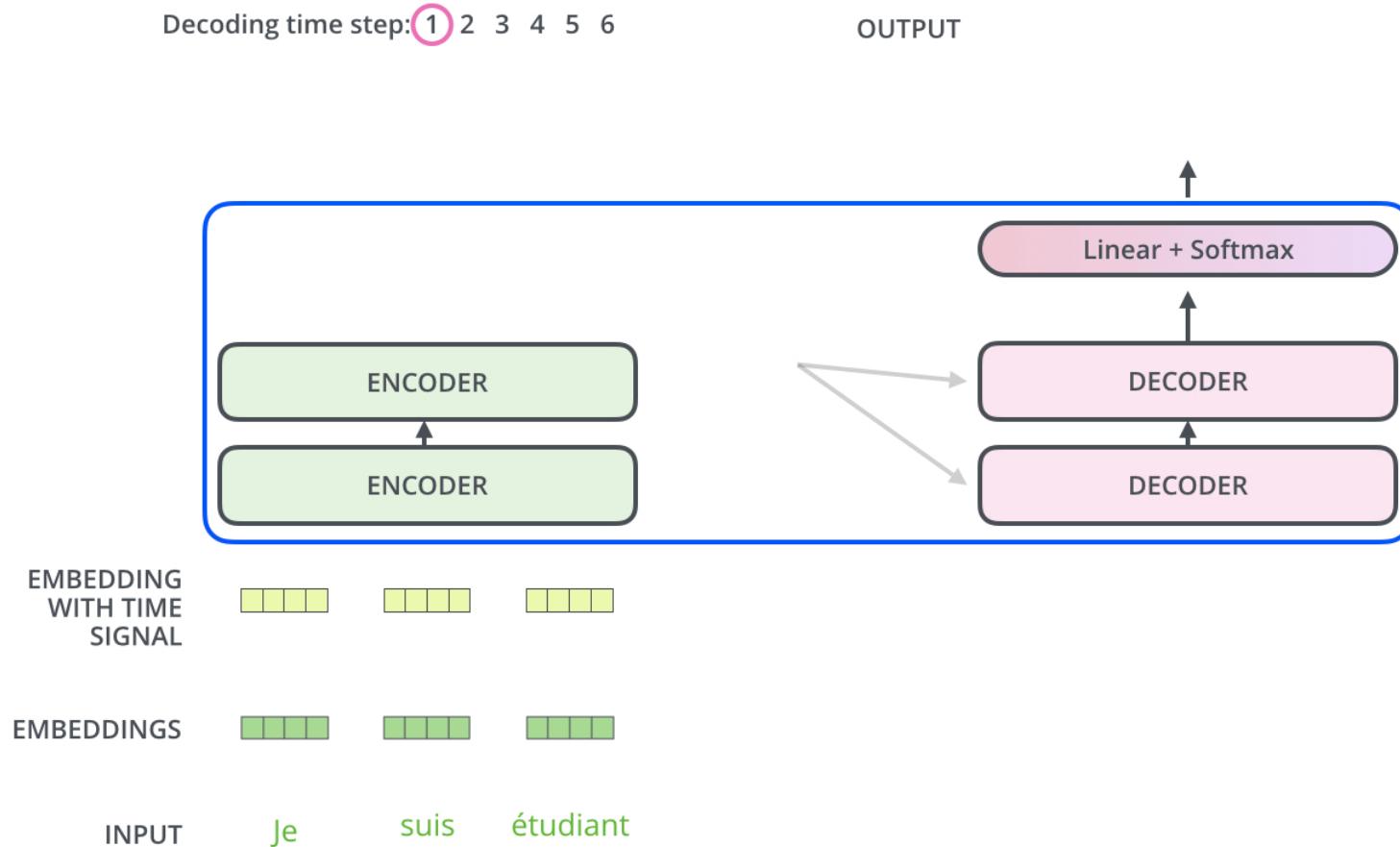


Transformer

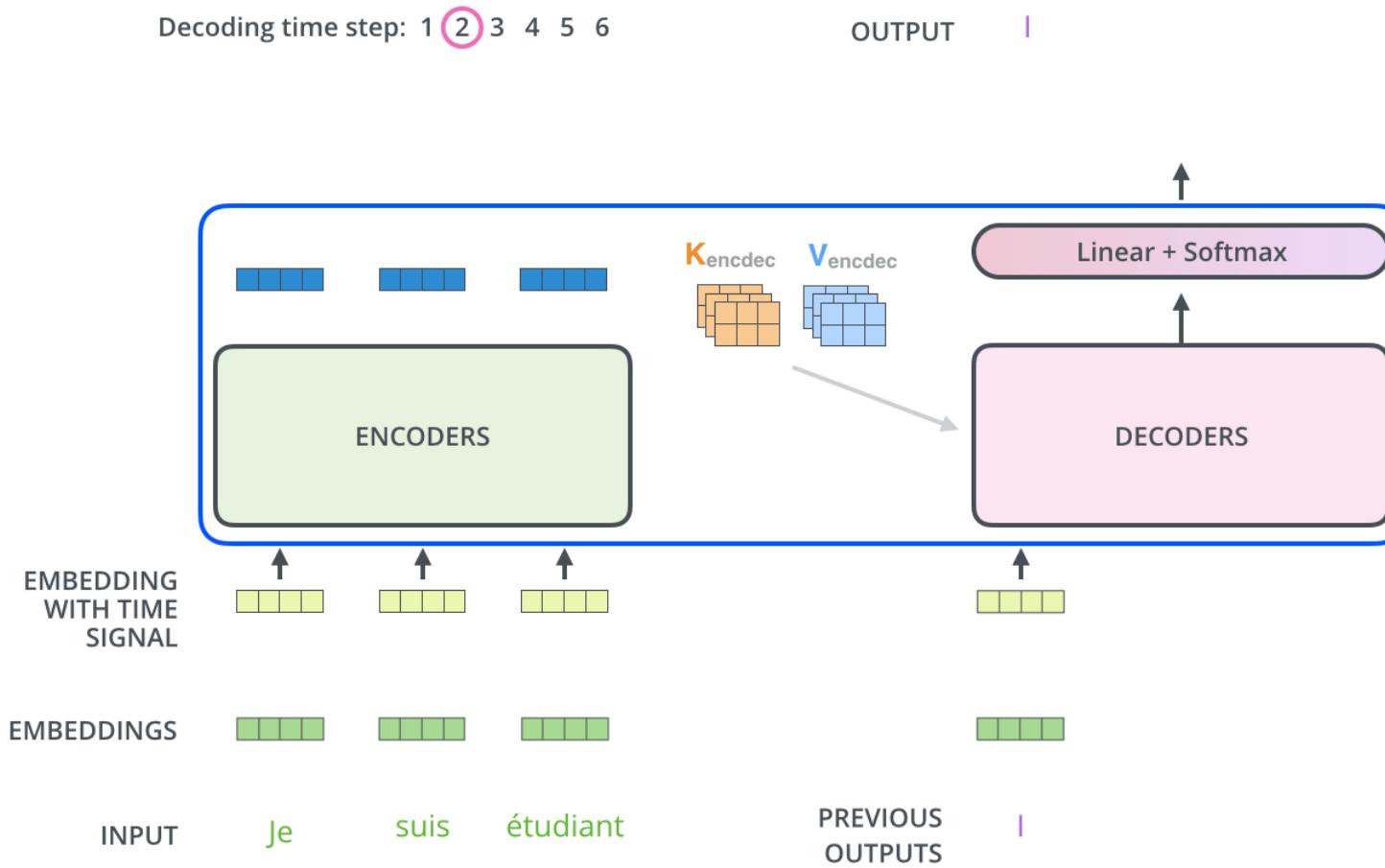
▶ Encoder Decoder Multi Attention



Transformer

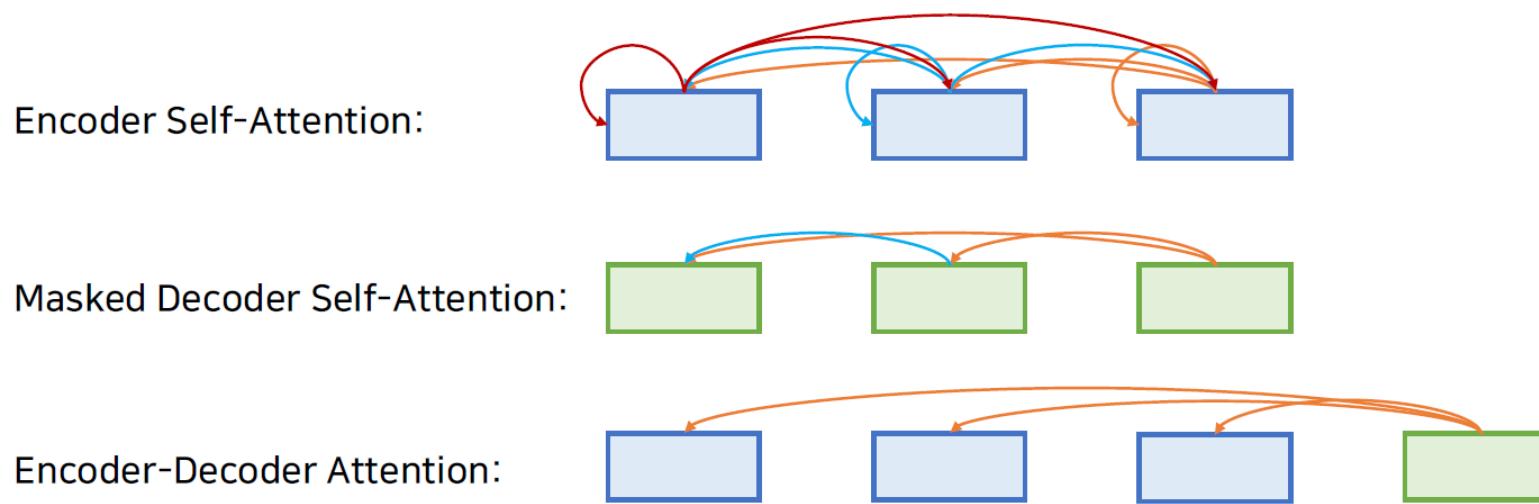


Transformer



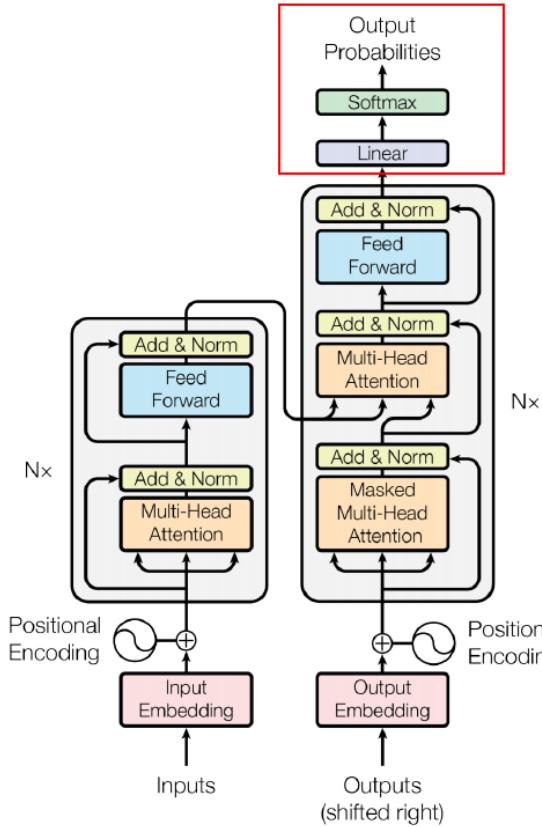
Transformer

▶ 3 Attention



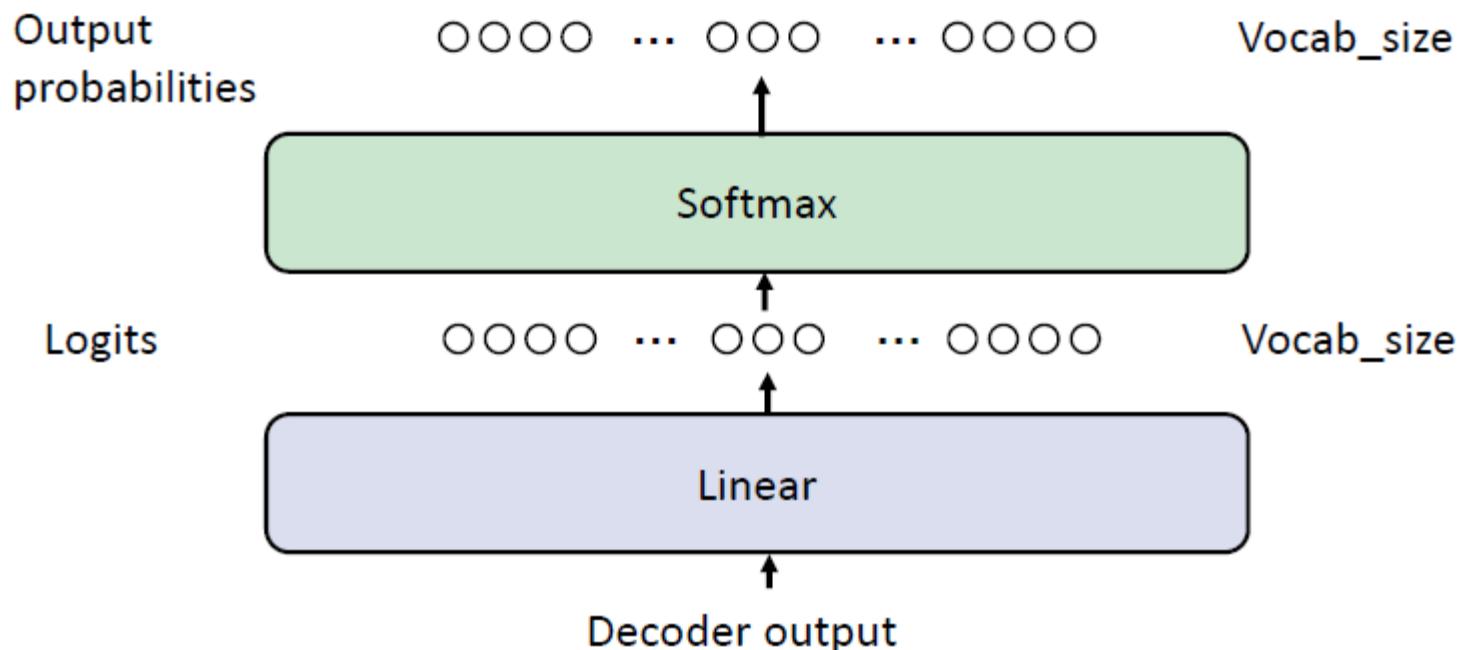
Transformer

▶ Output



Transformer

▶ Output



Transformer

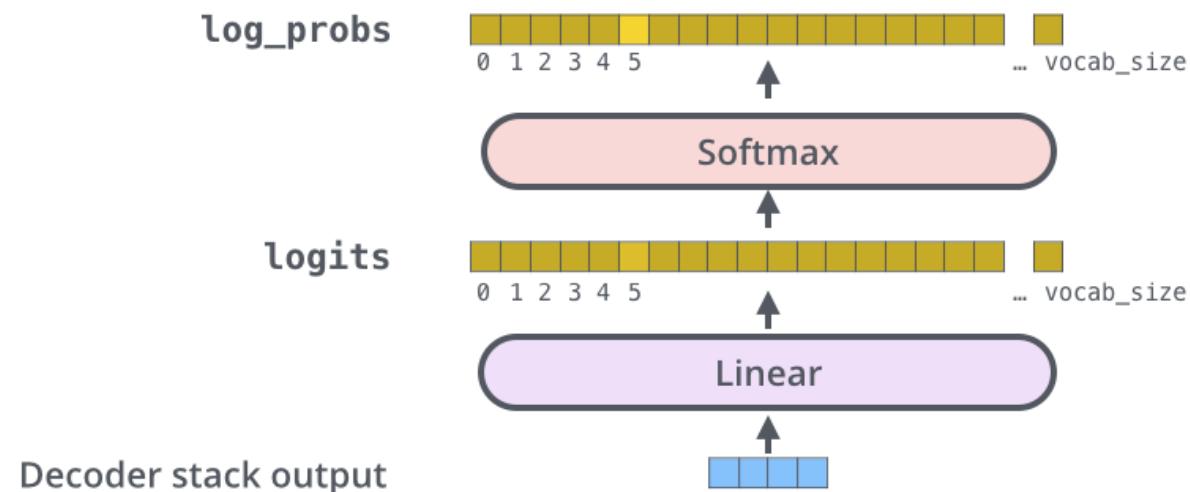
▶ Output

Which word in our vocabulary
is associated with this index?

am

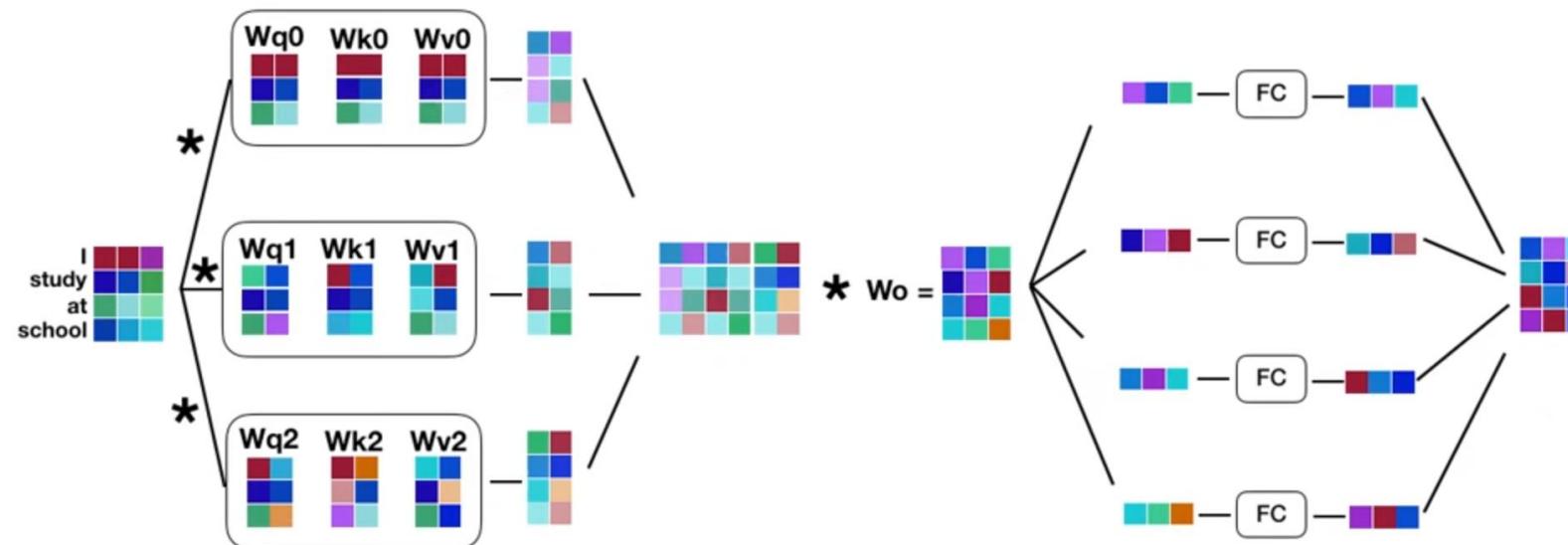
Get the index of the cell
with the highest value
(`argmax`)

5



Summary

Encoder



Summary

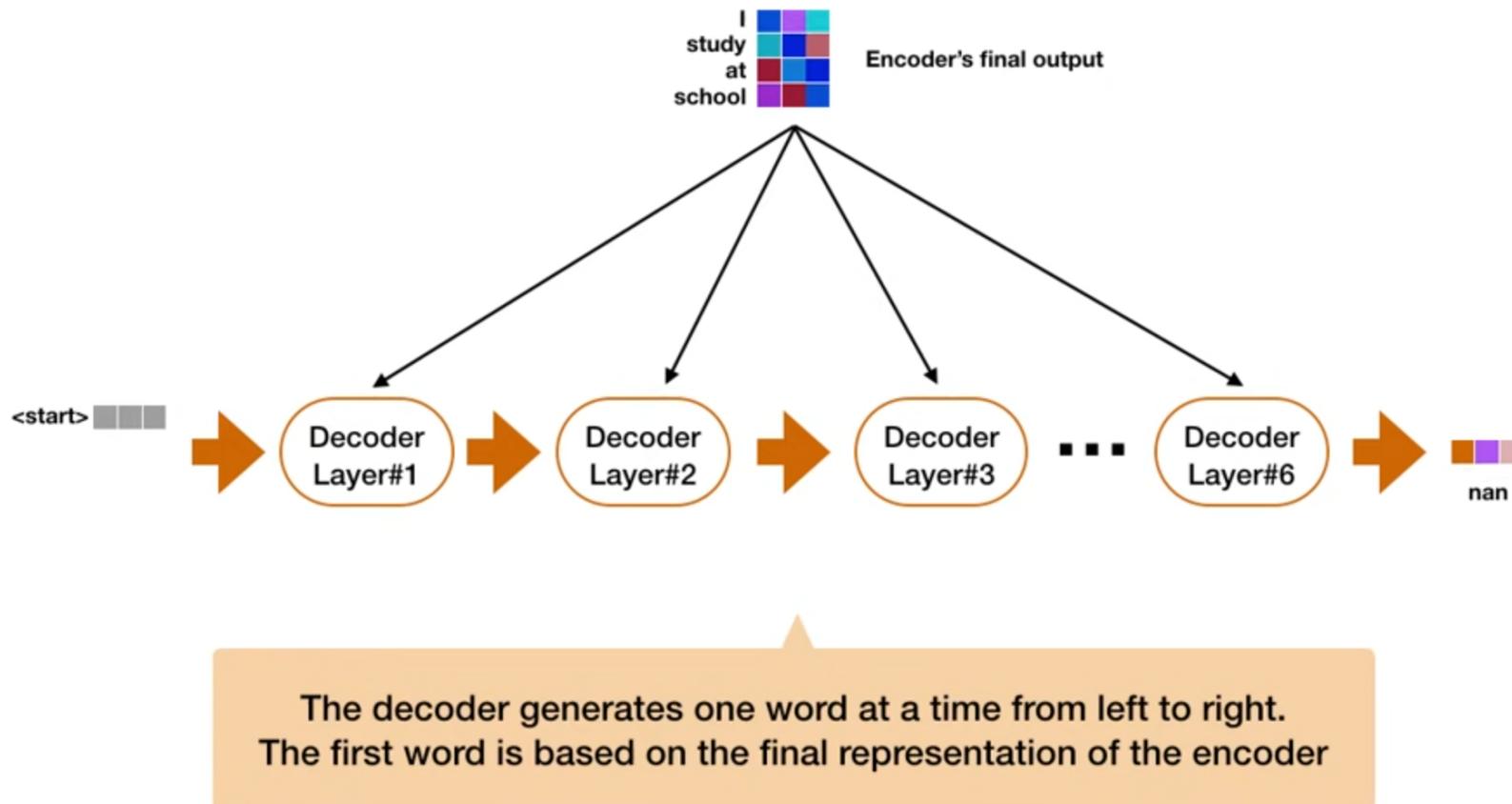
Transformer has 6 encoder layers



Encoder layers are identical but don't share weights

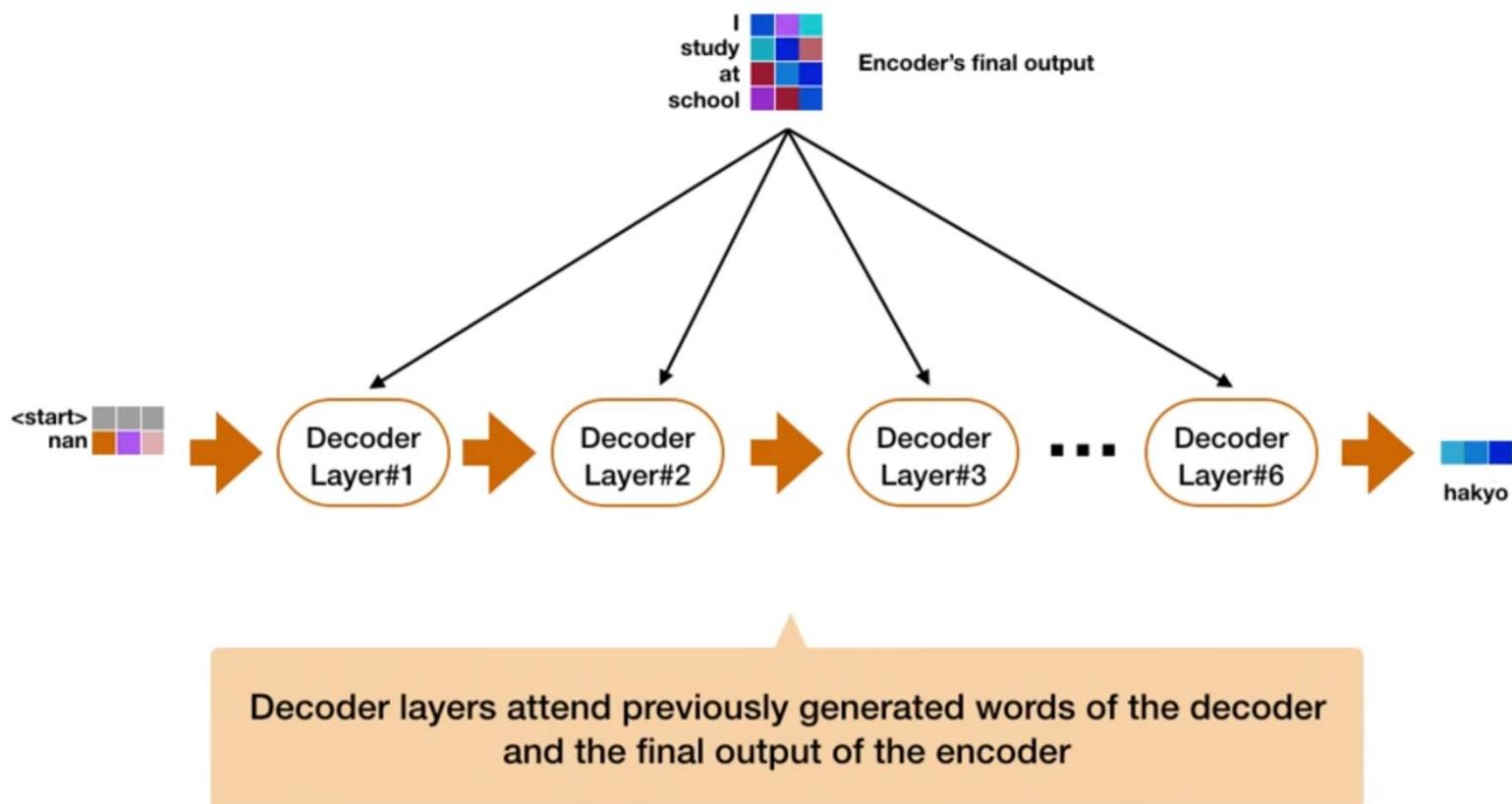
Summary

Transformer has 6 decoder layers



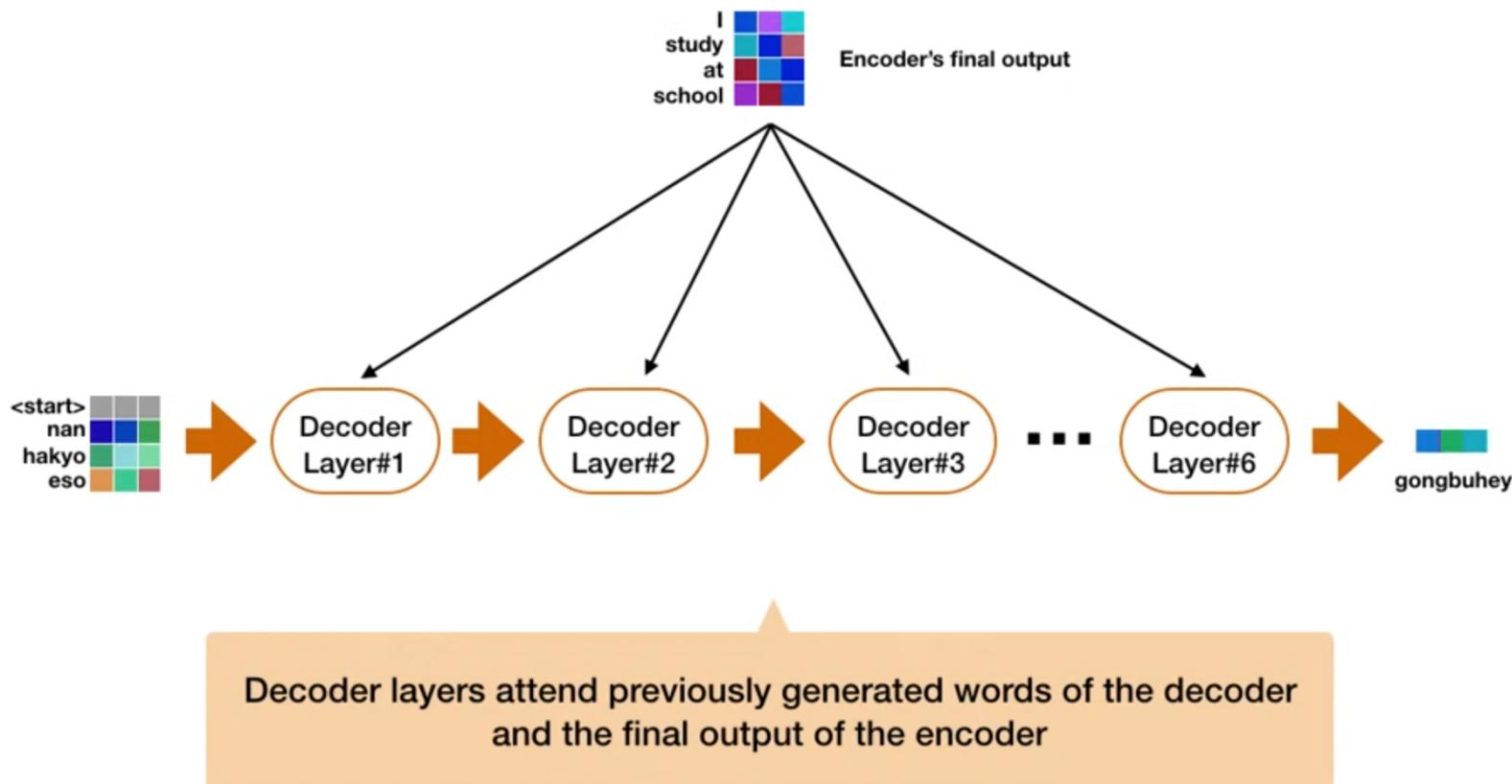
Summary

Transformer has 6 decoder layers



Summary

Transformer has 6 decoder layers



Transformer

