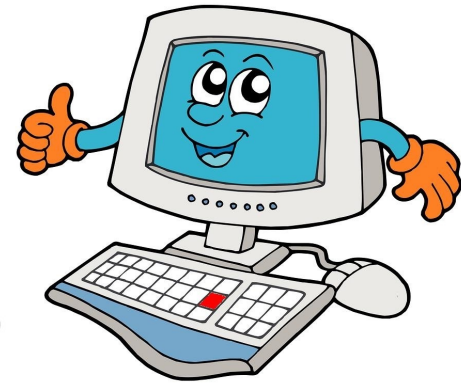
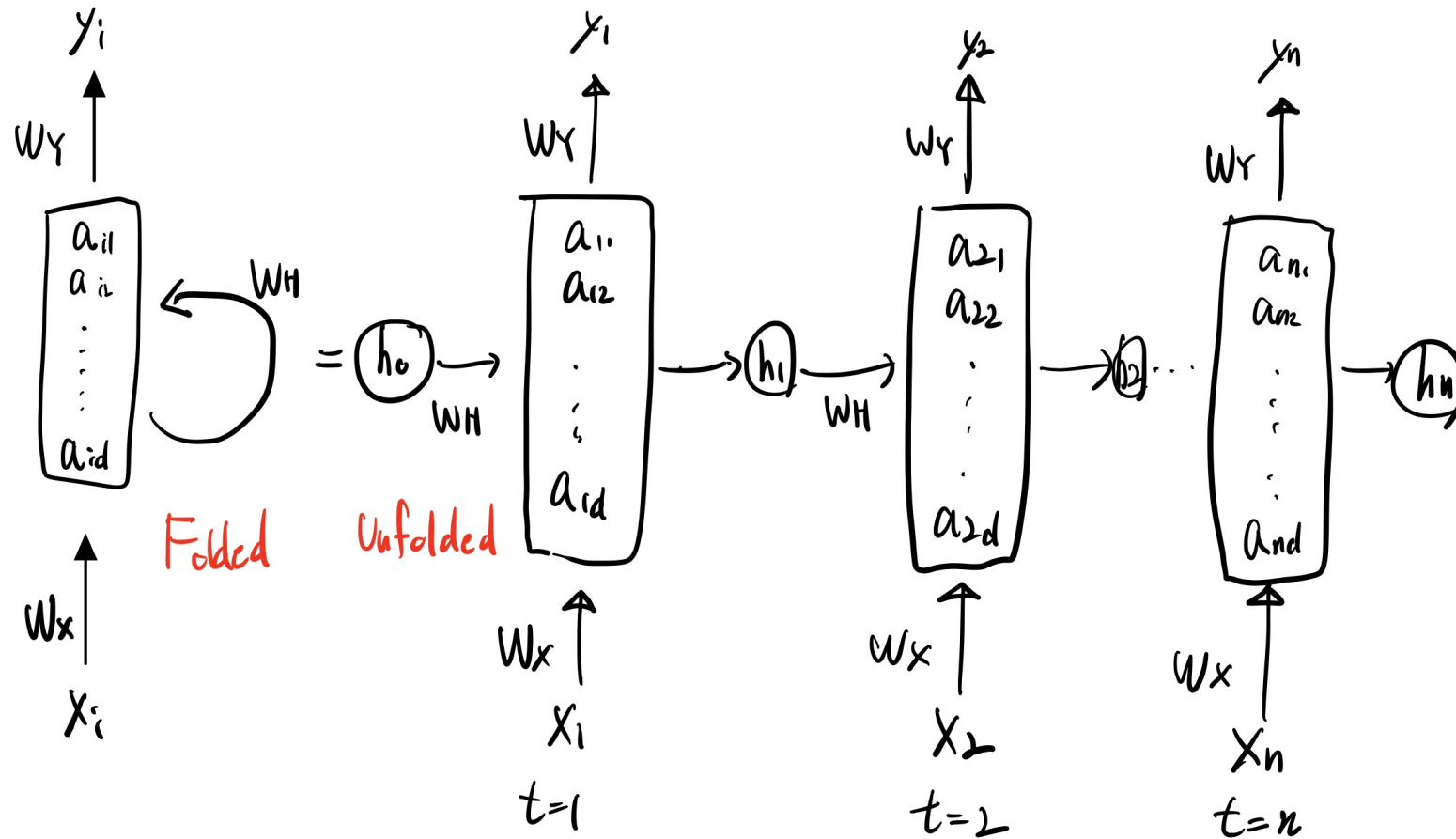


# Recurrent Neural Network

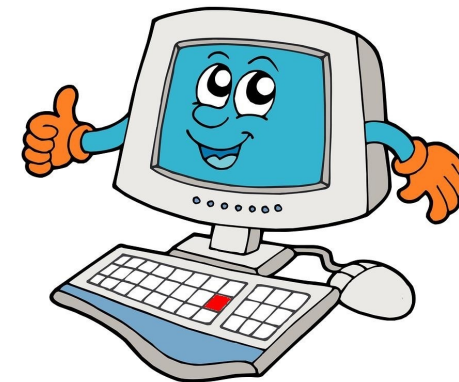
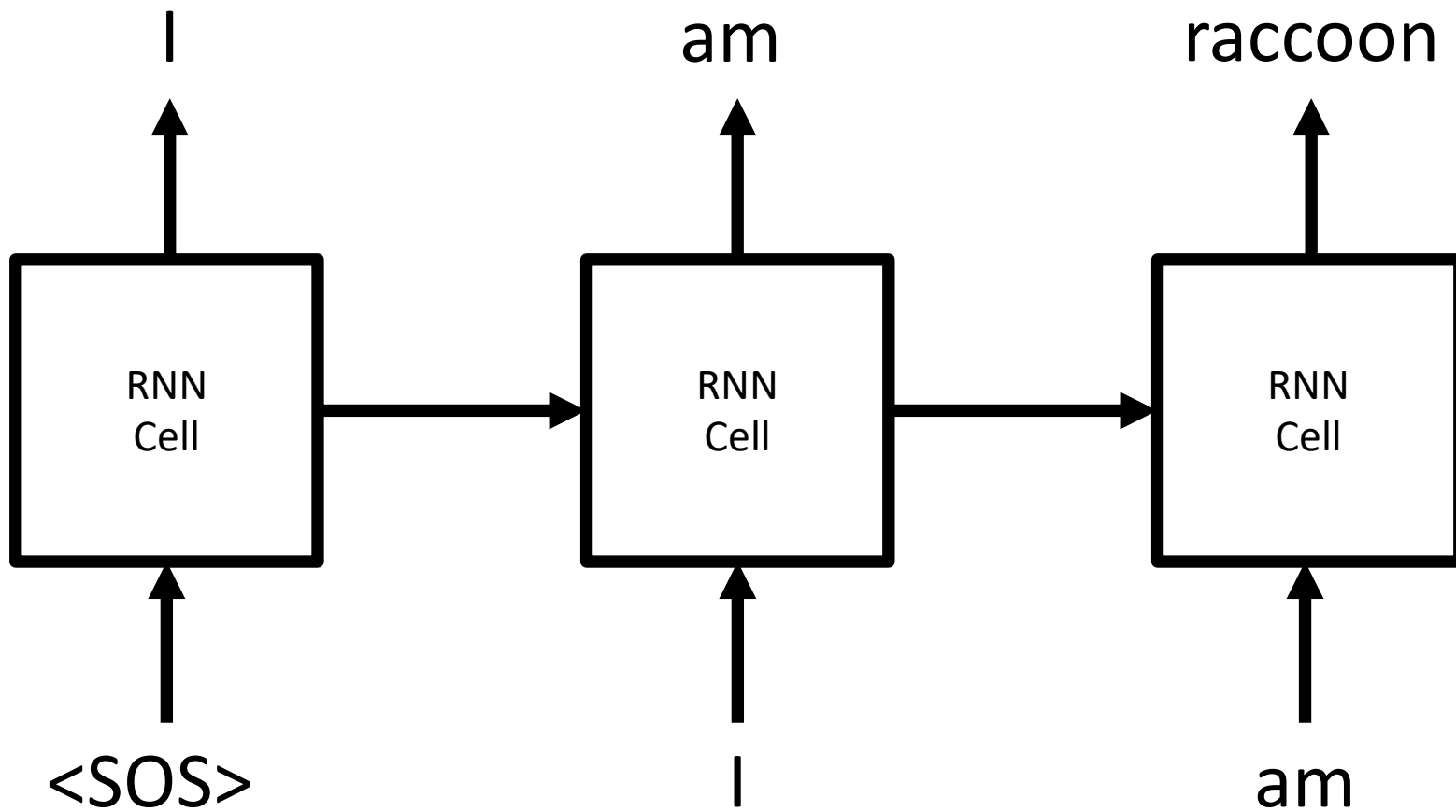
## : Processing Sequence Data



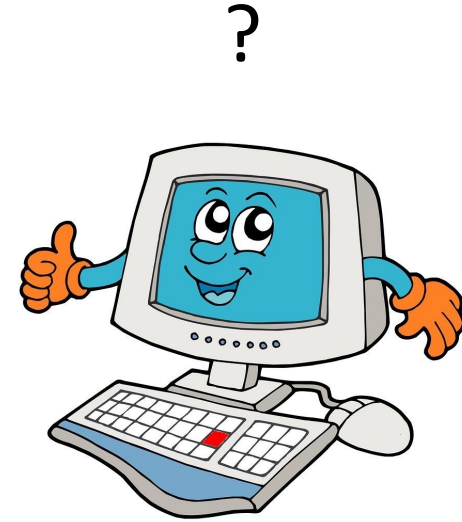
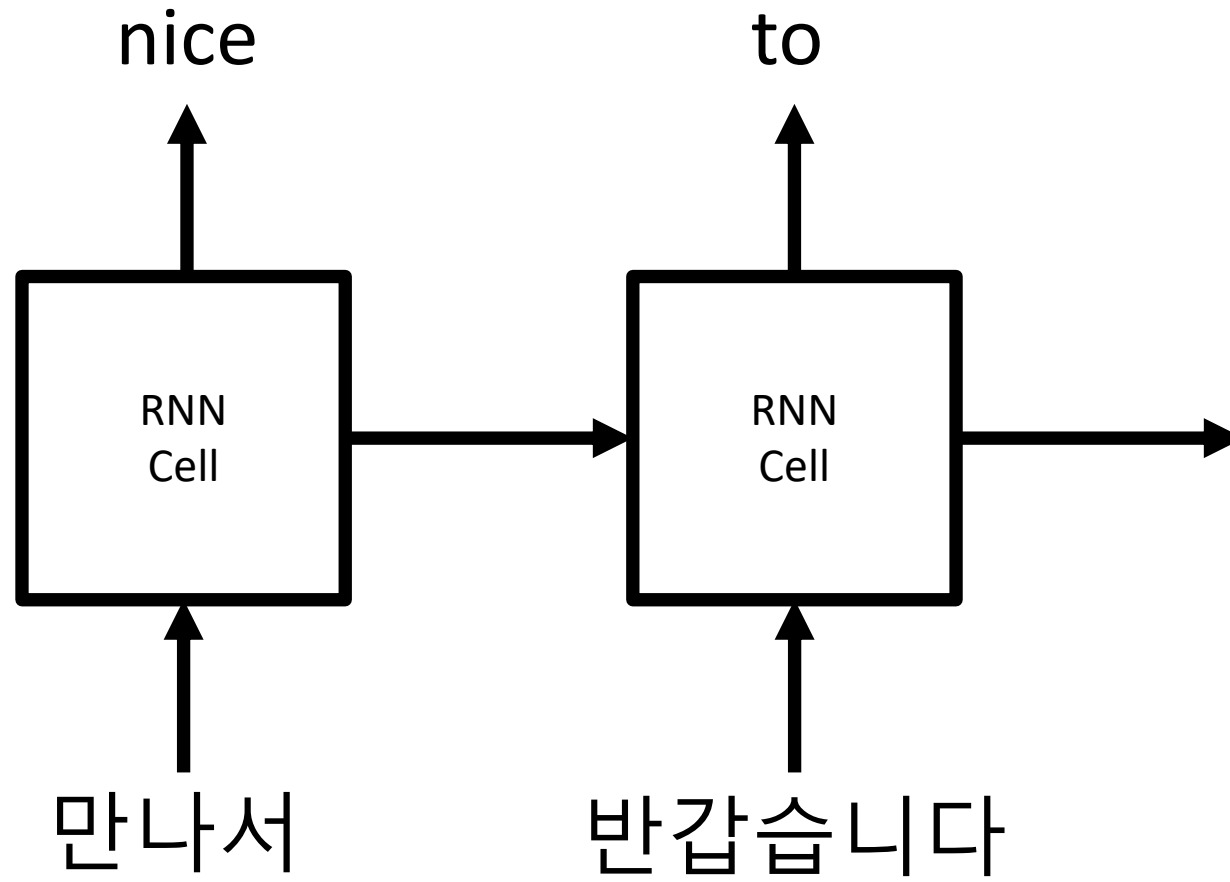
## Reminder



## Text Generation



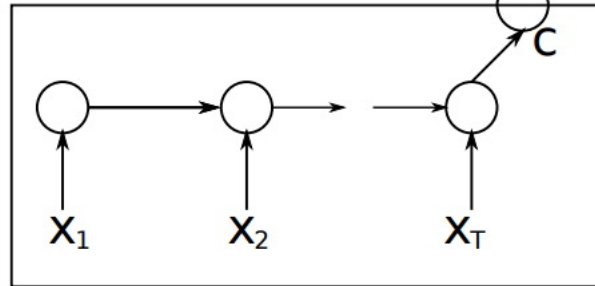
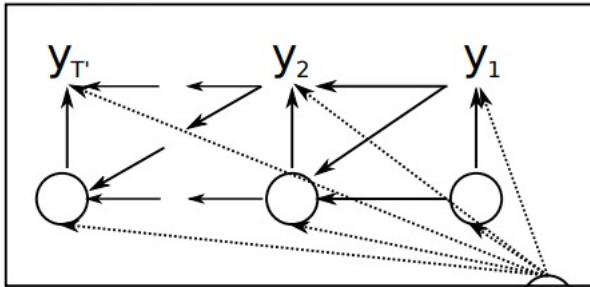
## Machine Translation Example



RNN Encoder-Decoder Model:

**Learning Phrase Representations using RNN Encoder–Decoder for  
Statistical Machine Translation (2014, Cho et al.)**

Decoder



Encoder

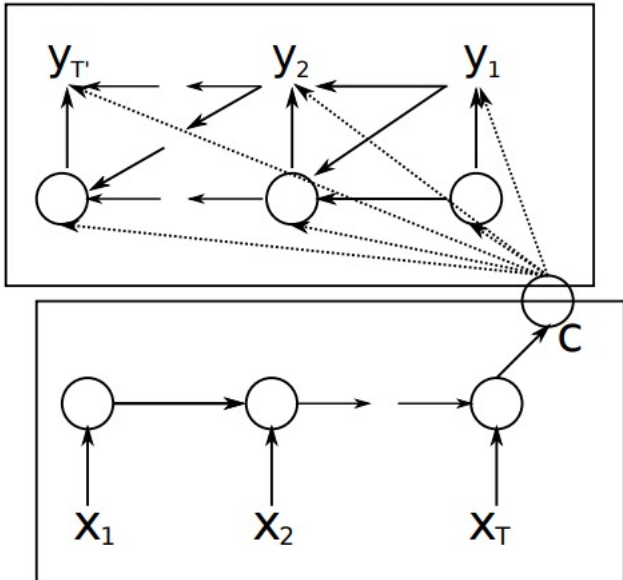
Main objective of Encoder-decoder model:

*Text Generation*

## RNN Encoder-Decoder Model:

### Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation (2014, Cho et al.)

Decoder



Encoder

$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T)$$

given input sequence  $x$ , model computes **output sequence  $y$ 's conditional distribution**.  
 $T \neq T'$

$$h_{<t>} = f(h_{<t-1>}, y_{t-1}, c)$$

$h_{<t>}$  = decoder hidden state at time  $t$   
 $c$  = encoders' summary vector

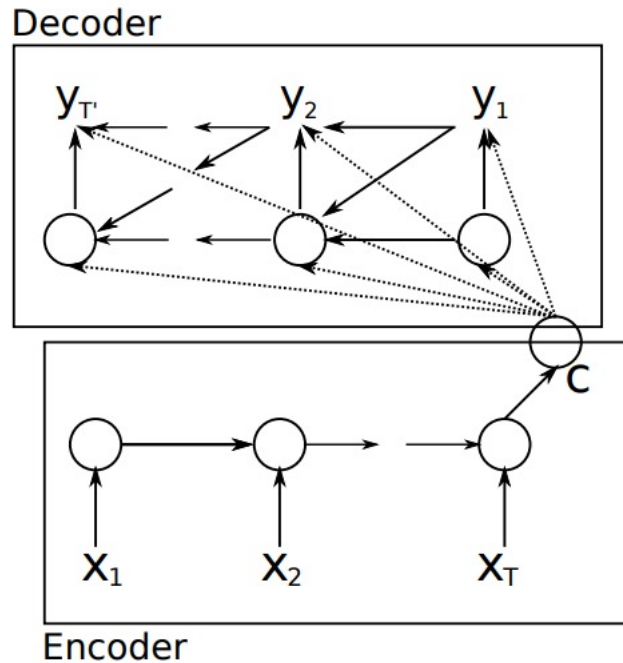
$$P(y_t | y_{t-1}, y_{t-2}, \dots, y_1, c) = g(h_{<t>}, y_{t-1}, c)$$

$f$  and  $g$  is activation function.  
 $g$  must produce probabilities.



RNN Encoder-Decoder Model:

## Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation (2014, Cho et al.)



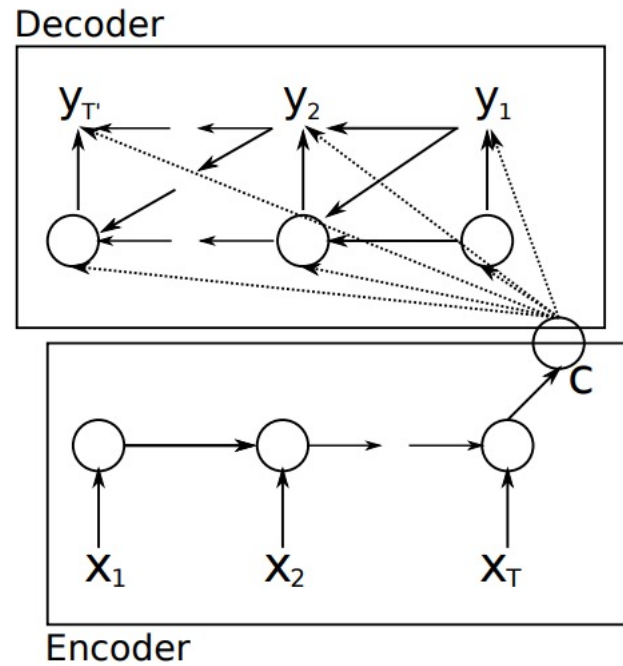
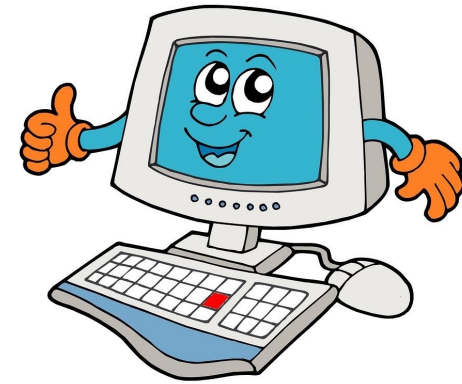
$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(y_n | x_n)$$

two components of the proposed model are jointly trained to **maximize the conditional log-likelihood**

## RNN Encoder-Decoder

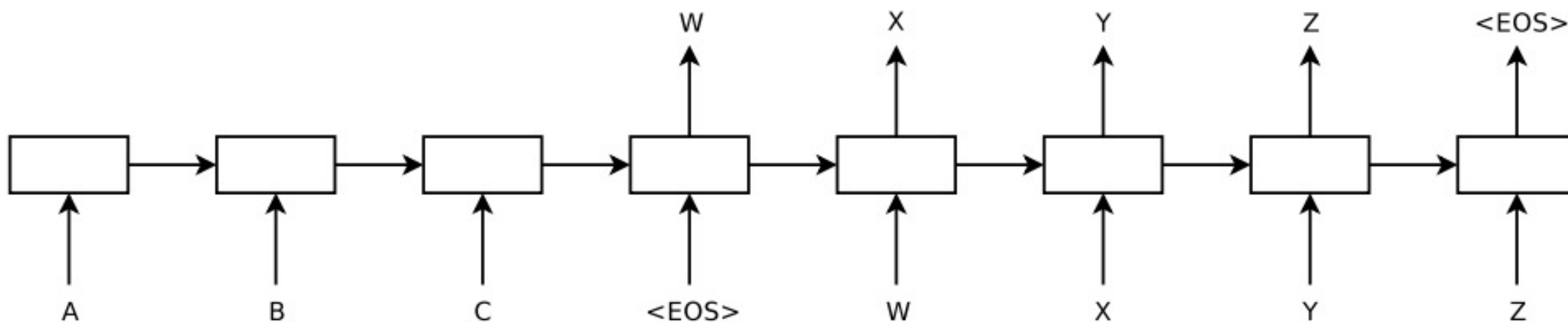
nice to meet you

만나서 반갑습니다





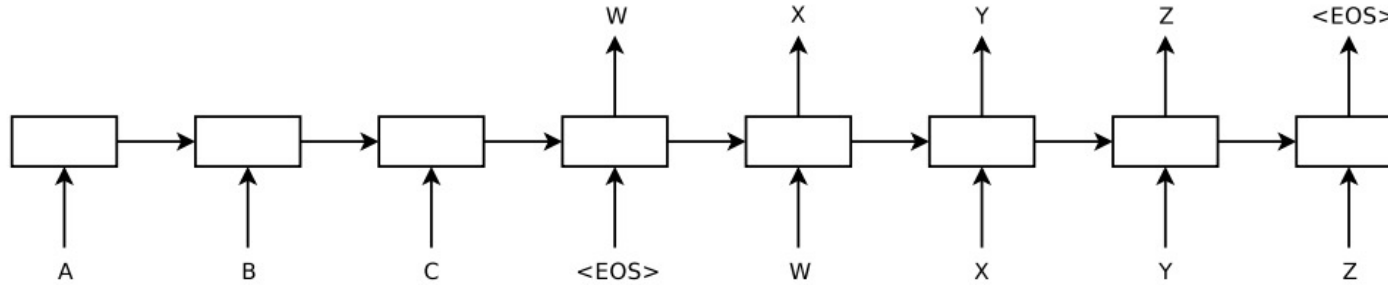
## Sequence to Sequence Learning with Neural Networks (2014, Sutskever et al.)



input sequence: ABC -> encoder

(objective) output sequence: WXYZ -> decoder

## Sequence to Sequence Learning with Neural Networks (2014, Sutskever et al.)

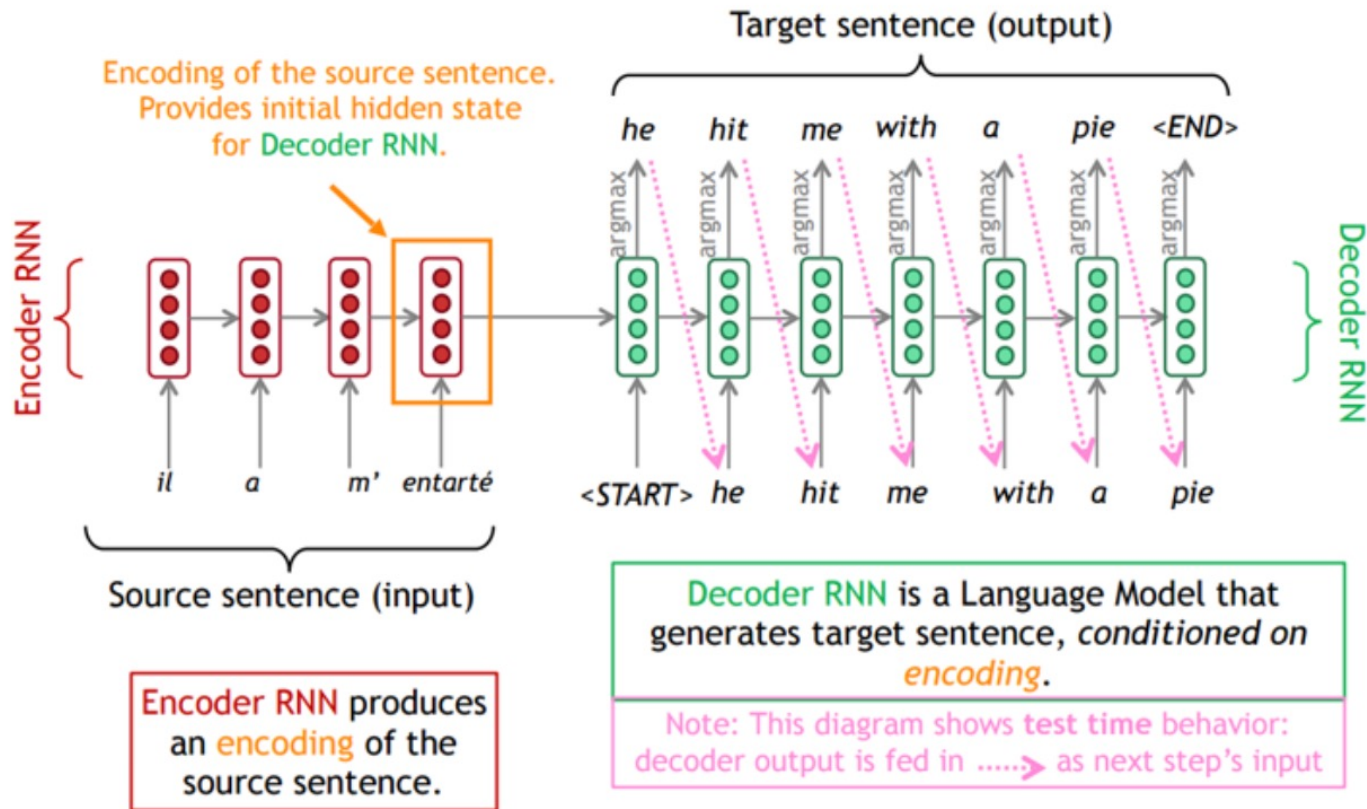


difficult to train the RNN (due to long term dependency) → **LSTM**

**<EOS> symbol** enables the model to define a distribution over *sequences of all possible lengths*

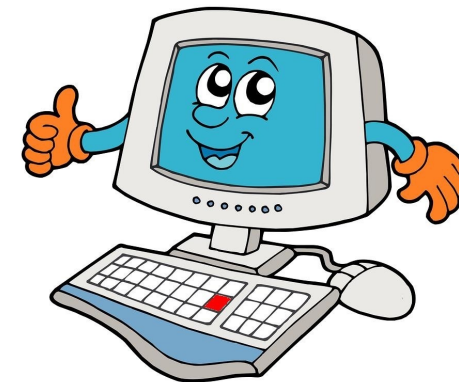
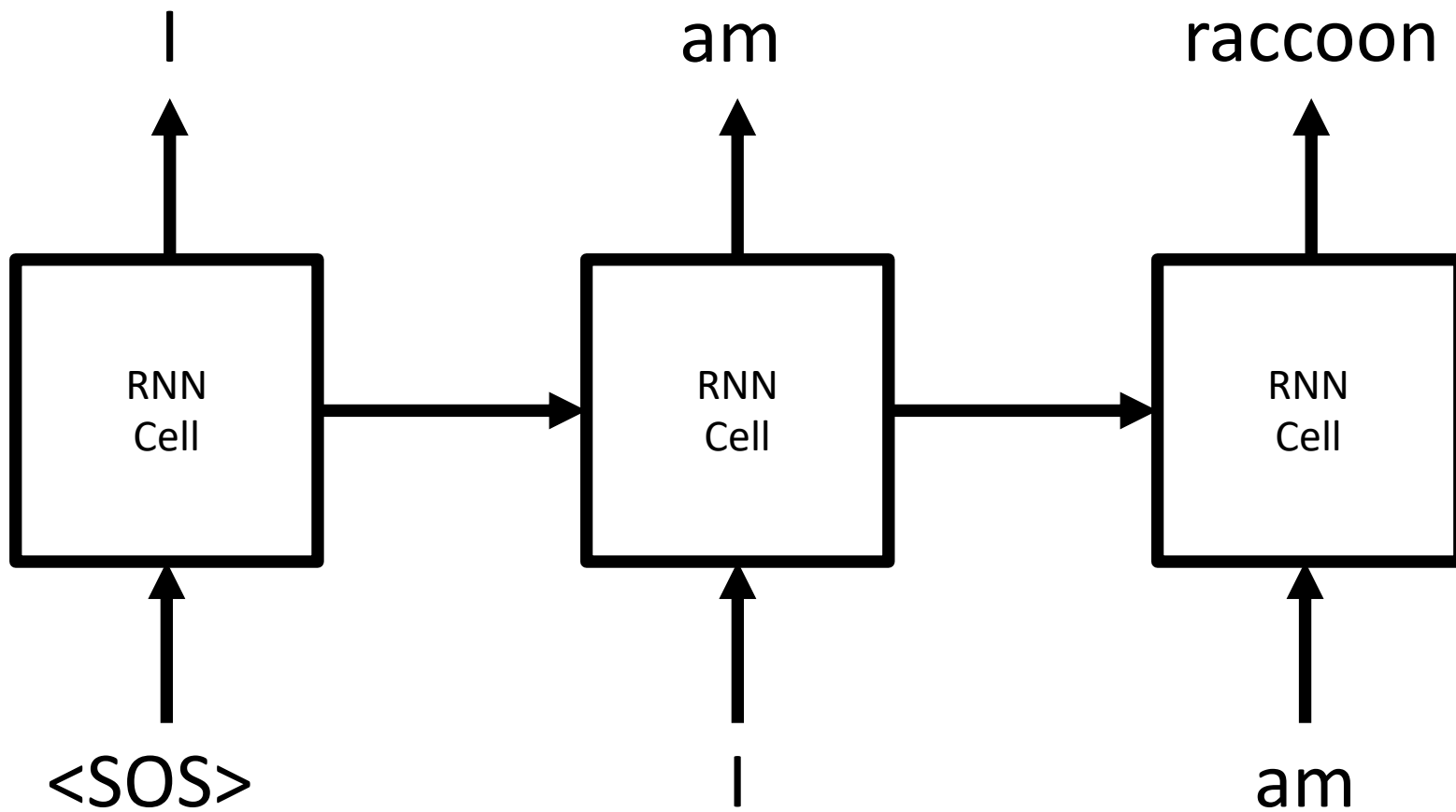
$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

# seq2seq

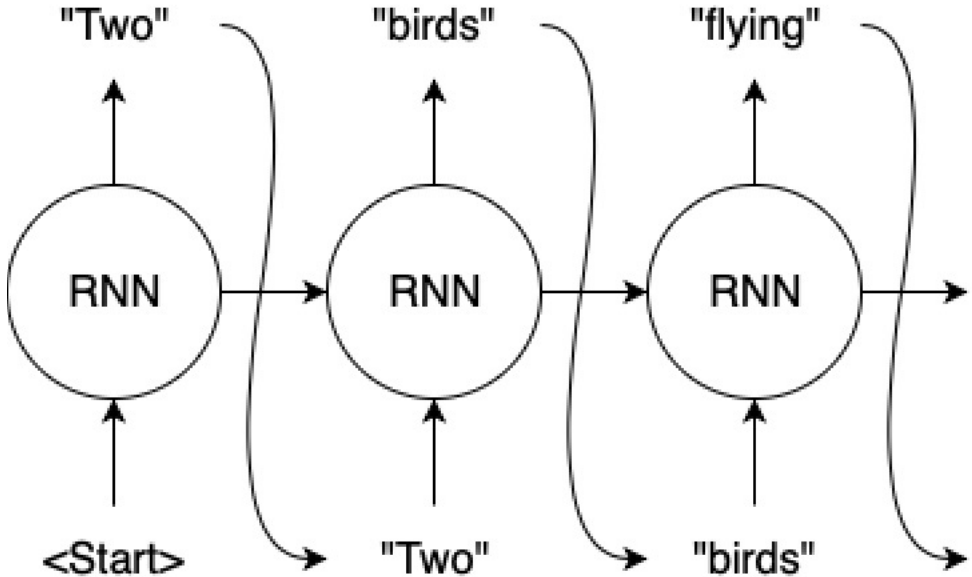


we want **conditional probability** of output sequence by given input sequence

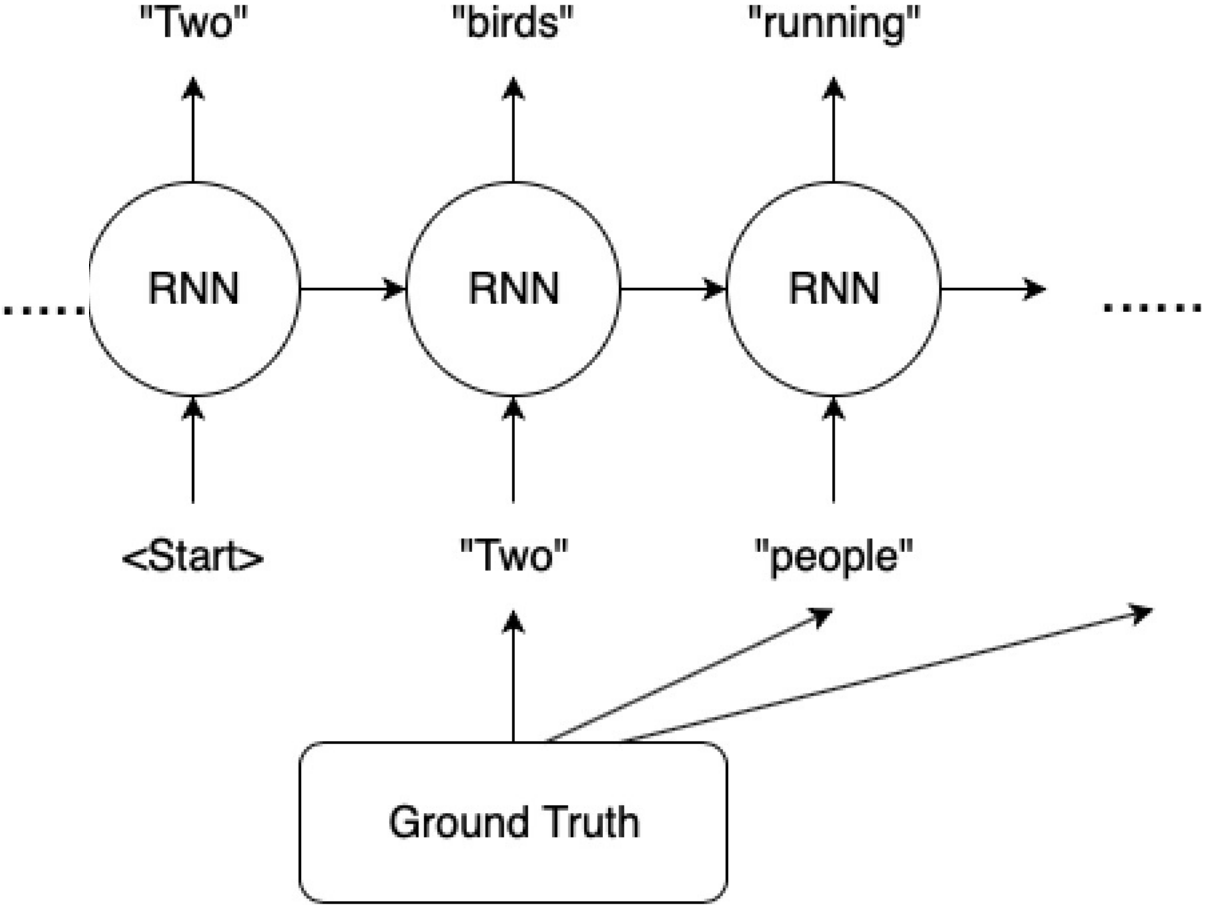
## Text Generation



# Teacher Forcing



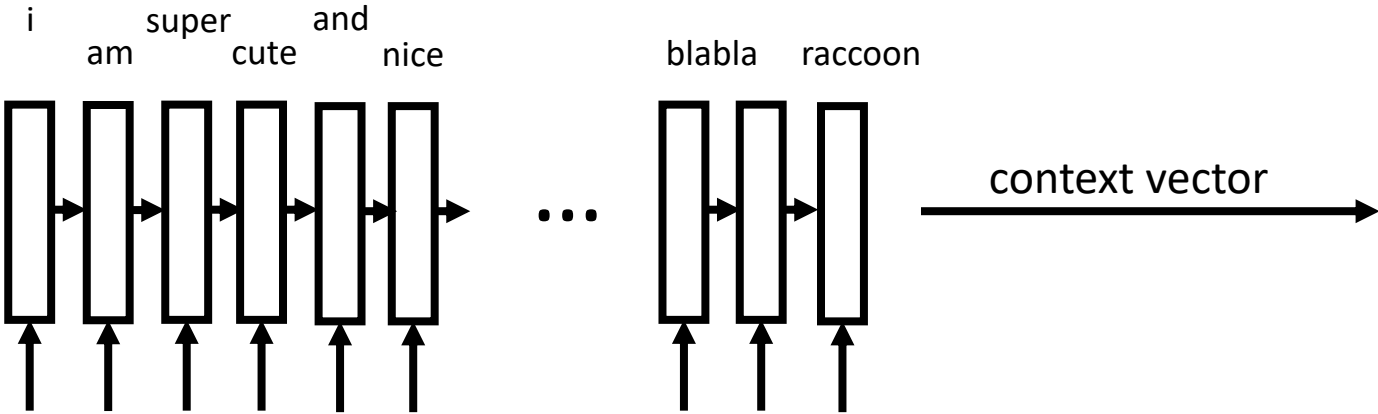
Without Teacher Forcing



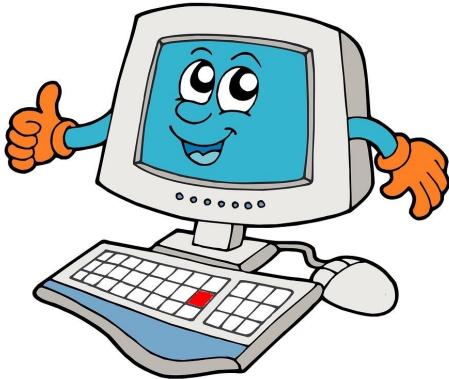
With Teacher Forcing

# Code Review

seq2seq

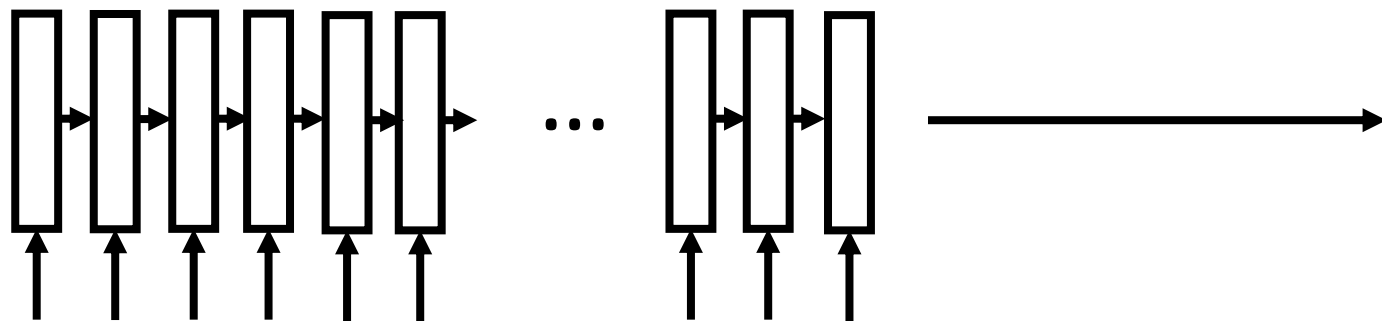


?



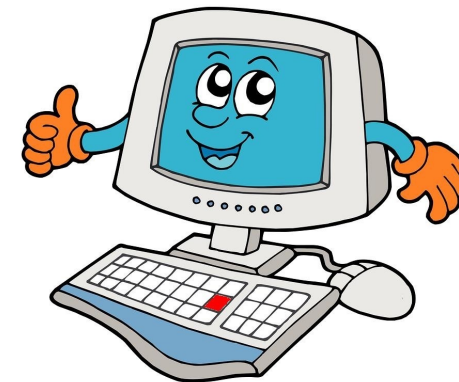
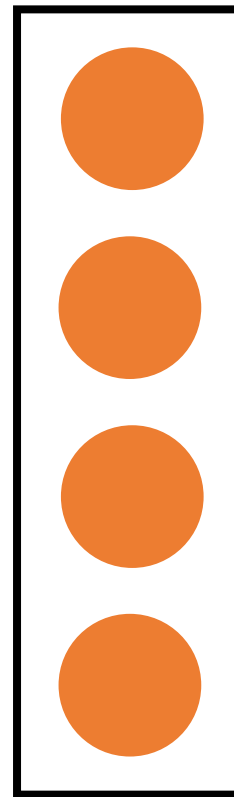


seq2seq



저는/너구리가/아니라/라쿤입니다/어쩌구저쩌구... tmi

context vector  
(4x1 vector)



## Attention Visualizations

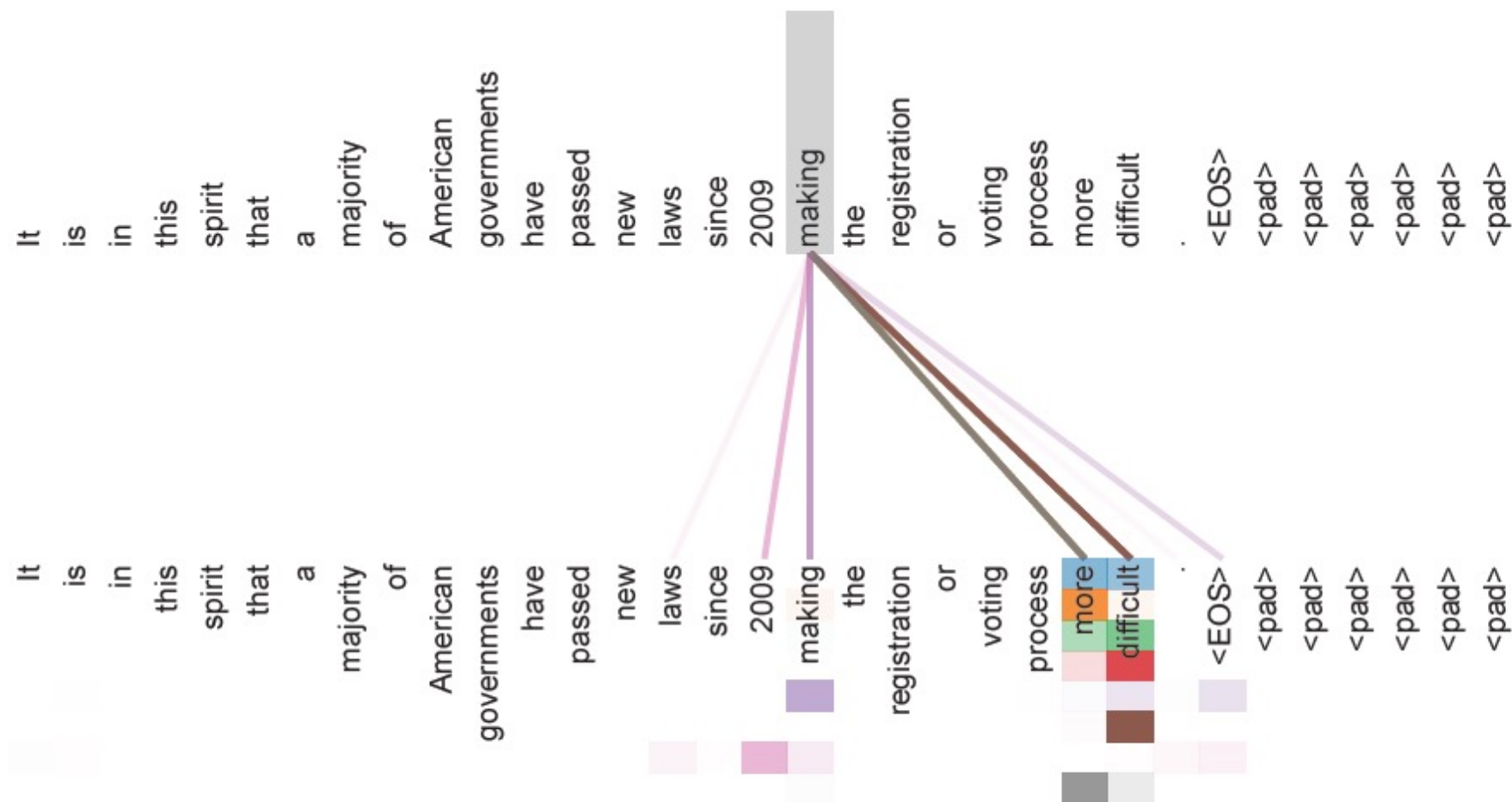


Figure 3: An example of the attention mechanism following long-distance dependencies in the encoder self-attention in layer 5 of 6. Many of the attention heads attend to a distant dependency of the verb 'making', completing the phrase 'making...more difficult'. Attentions here shown only for the word 'making'. Different colors represent different heads. Best viewed in color.

# Neural Machine Translation by Jointly Learning to Align and Translate(2014, Bahdanau et al.)

$$p(y_i | y_1, \dots, y_{i-1}, x) = g(y_{i-1}, s_i, c_i)$$

each conditional probability at decoder rnn cell

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

s\_i is an RNN hidden state for time i

$$(h_1, \dots, h_{T_x})$$

encoder maps the input sentence into a sequence of annotations  
context vector c\_i depends on this.

**★ each annotation  $h_i$  contains information about the whole input sequence with a strong focus on the parts surrounding the  $i$ -th word of the input sequence.**

# Neural Machine Translation by Jointly Learning to Align and Translate(2014, Bahdanau et al.)

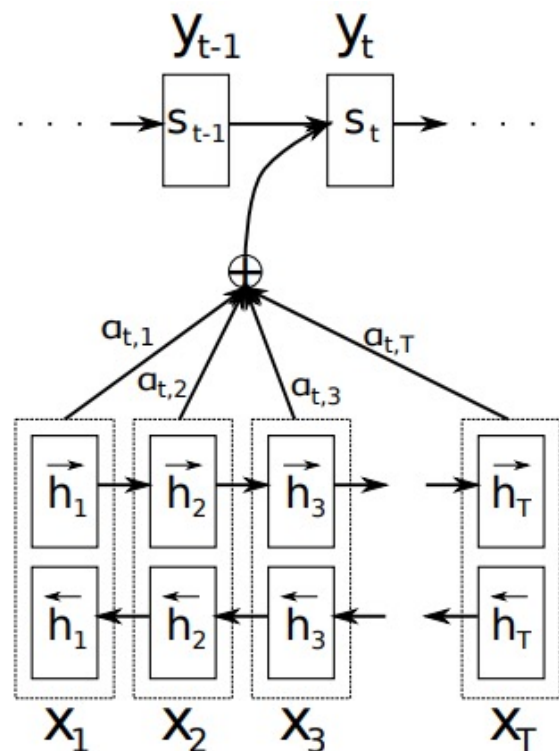


Figure 1: The graphical illustration of the proposed model trying to generate the  $t$ -th target word  $y_t$  given a source sentence  $(x_1, x_2, \dots, x_T)$ .

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

context vector  $c_i$  can be described as weighted sum of annotations  $h_i$ .

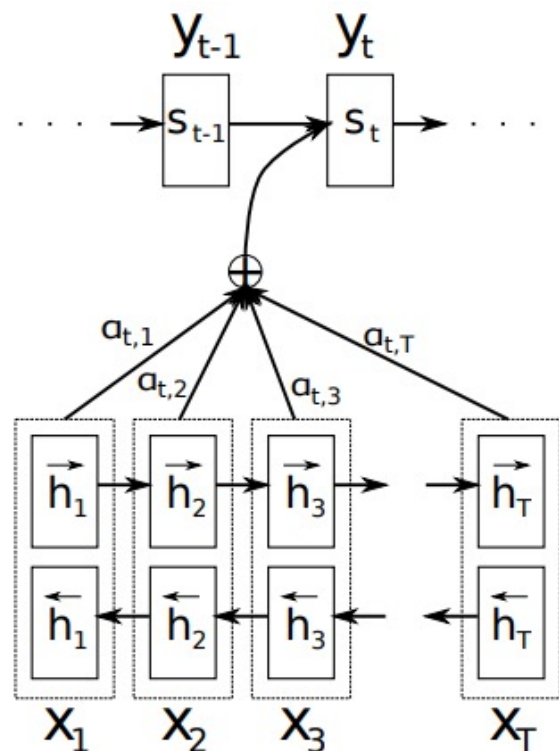
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$i$  = decoder  $i$ -th time step  
 $j$  = input sequences'  $j$ -th annotation

$$e_{ij} = a(s_{i-1}, h_j)$$

$a$  is an alignment model.  
 scores how well the inputs around position  $j$  and the output at position  $i$  match.

# Neural Machine Translation by Jointly Learning to Align and Translate(2014, Bahdanau et al.)



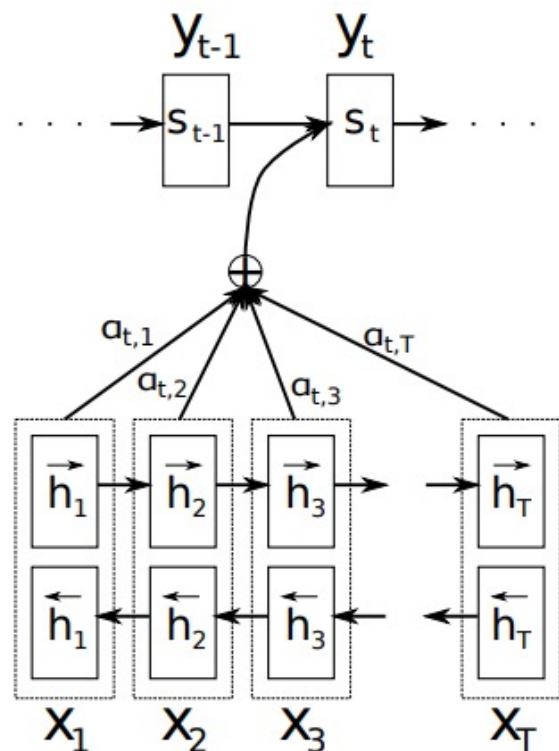
*alignment model  $a$*

$$e_{ij} = a(s_{i-1}, h_j)$$

***feed-forward neural network!***

Figure 1: The graphical illustration of the proposed model trying to generate the  $t$ -th target word  $y_t$  given a source sentence  $(x_1, x_2, \dots, x_T)$ .

# Neural Machine Translation by Jointly Learning to Align and Translate(2014, Bahdanau et al.)



*bi-directional RNN*

*→ not only the preceding words, but also the following words.*

$$(\vec{h}_1, \dots, \vec{h}_{T_x})$$

forward hidden state

$$(\overleftarrow{h}_1, \dots, \overleftarrow{h}_{T_x})$$

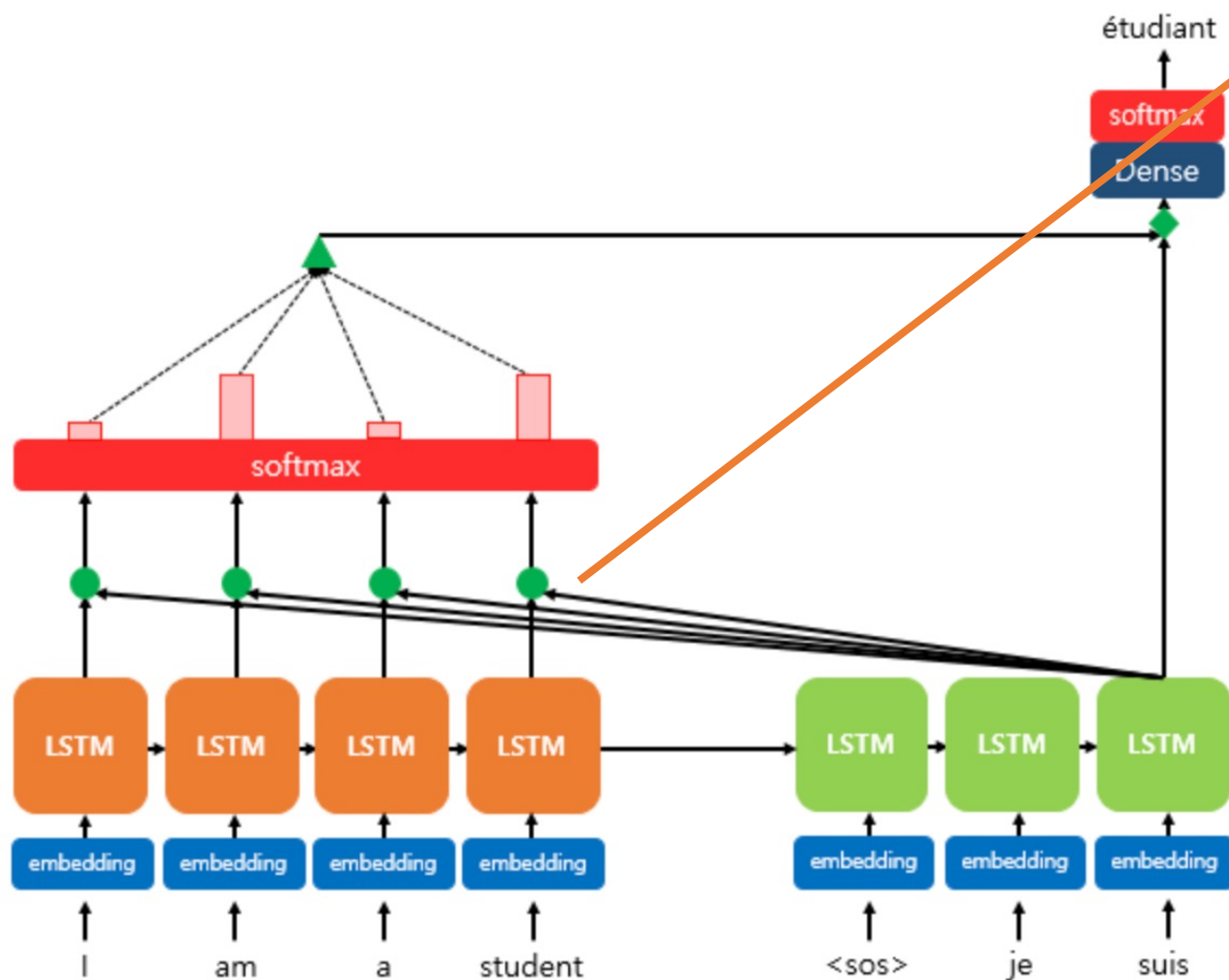
backward hidden state

$$h_j = \left[ \vec{h}_j^\top; \overleftarrow{h}_j^\top \right]^\top$$

concatenation

Figure 1: The graphical illustration of the proposed model trying to generate the  $t$ -th target word  $y_t$  given a source sentence  $(x_1, x_2, \dots, x_T)$ .

# Attention Mechanism



$$e_{ij} = a(s_{i-1}, h_j)$$

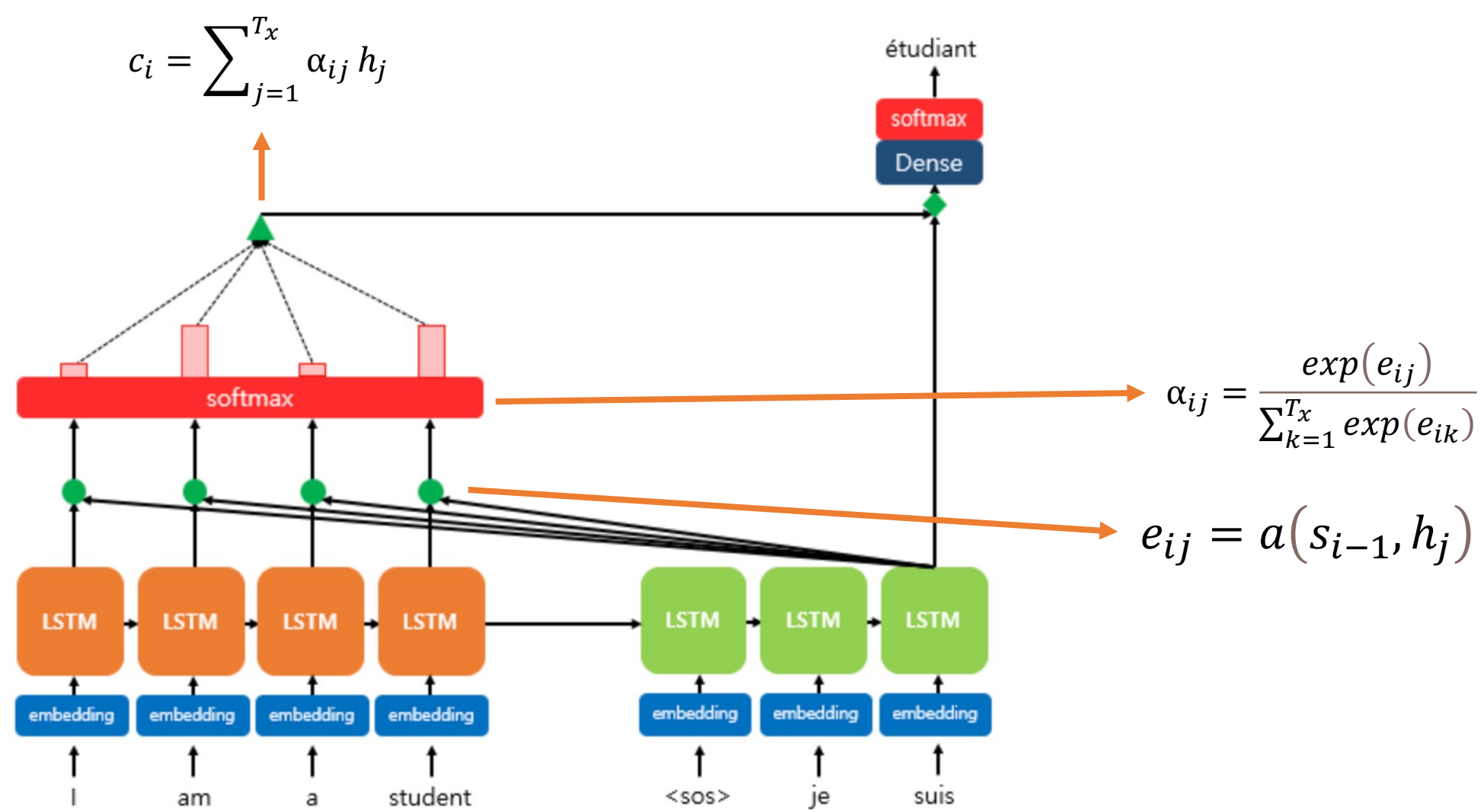
dot-product attention score

The diagram shows a horizontal green vector labeled  $s_t^T$  multiplied by a vertical orange vector labeled  $h_i$ . The result is a single value, representing the dot-product attention score.

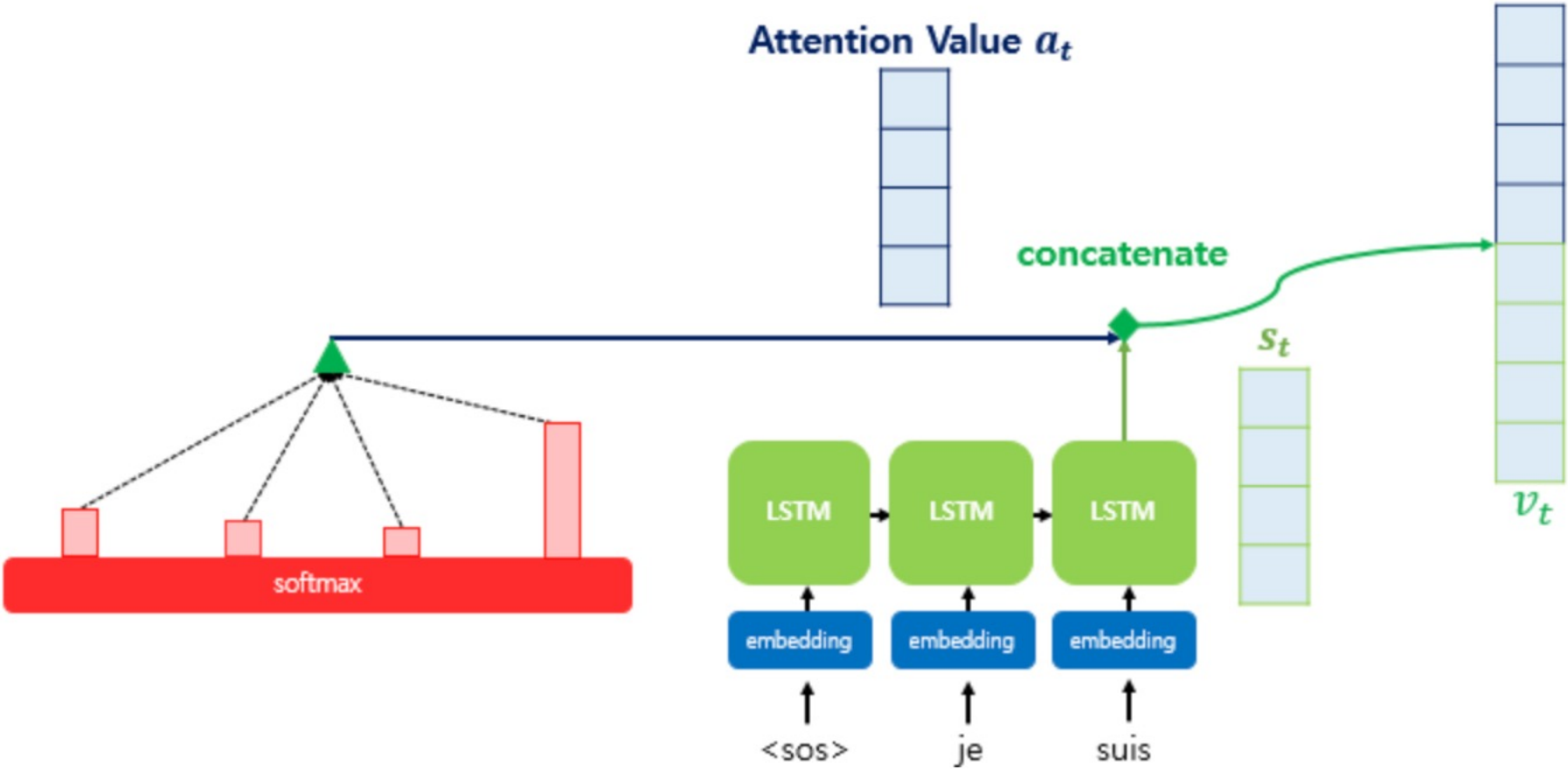
$$a(s_{i-1}, h_j) = s_{i-1}^T h_j$$



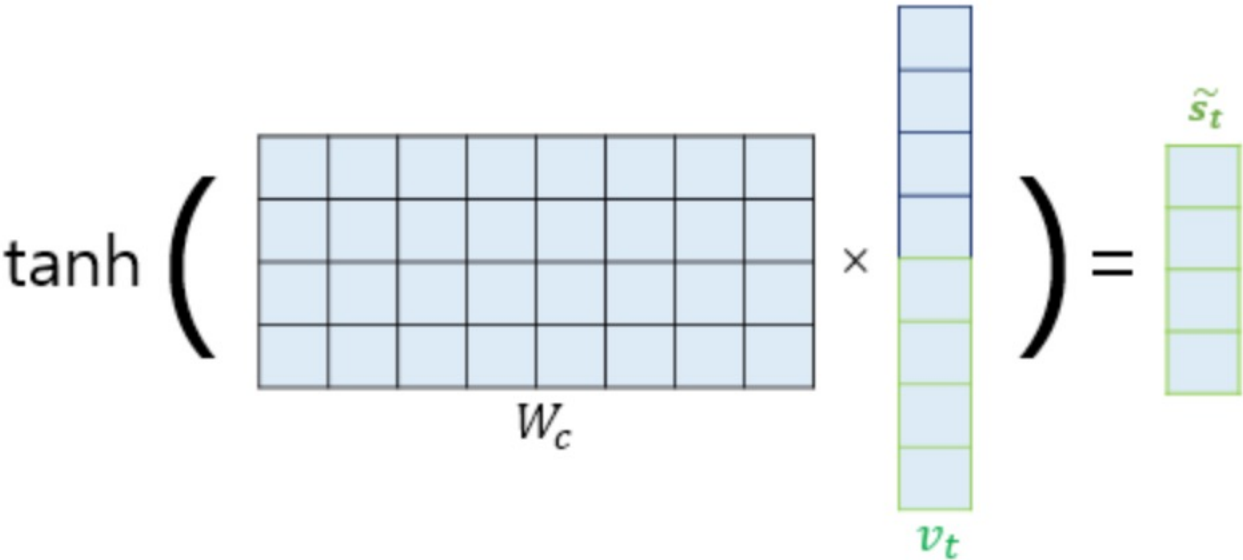
# Attention Mechanism



# Attention Mechanism



# Attention Mechanism



$$\tilde{s}_t = \tanh(W_c v_t + b_c)$$

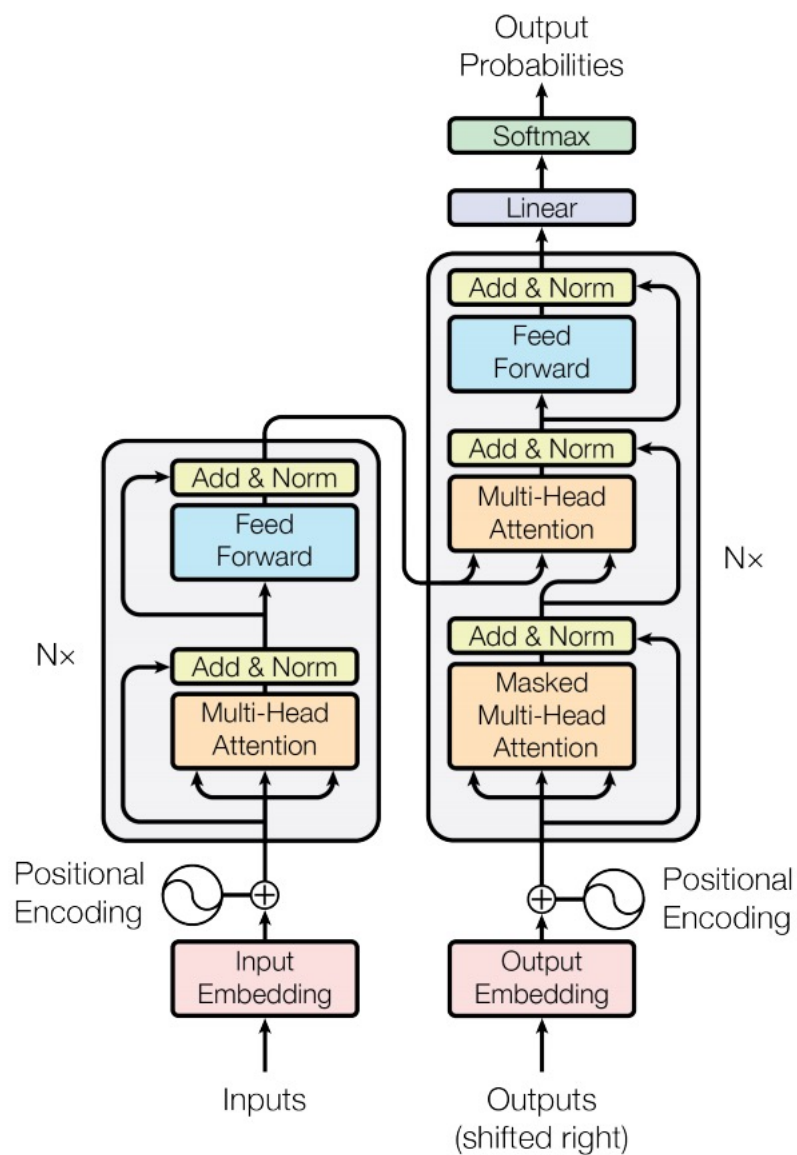


$$\hat{y}_t = \text{Softmax}(W_y \tilde{s}_t + b_y)$$

## Attention Score function

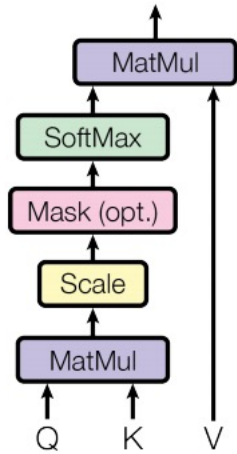
이름	스코어 함수	Defined by
<i>dot</i>	$score(s_t, h_i) = s_t^T h_i$	Luong et al. (2015)
<i>scaled dot</i>	$score(s_t, h_i) = \frac{s_t^T h_i}{\sqrt{n}}$	Vaswani et al. (2017)
<i>general</i>	$score(s_t, h_i) = s_t^T W_a h_i$ // 단, $W_a$ 는 학습 가능한 가중치 행렬	Luong et al. (2015)
<i>concat</i>	$score(s_t, h_i) = W_a^T \tanh(W_b[s_t; h_i])$ $score(s_t, h_i) = W_a^T \tanh(W_b s_t + W_c h_i)$	Bahdanau et al. (2015)
<i>location - base</i>	$\alpha_t = softmax(W_a s_t)$ // $\alpha_t$ 산출 시에 $s_t$ 만 사용하는 방법.	Luong et al. (2015)

# Attention Is All You Need (2017, Vaswani et al.)



# Attention Is All You Need (2017, Vaswani et al.)

Scaled Dot-Product Attention



Query: *decoder hidden state at time step t*

Keys: *(all) encoder cells' hidden state*

Values: *(all) encoder cells' hidden state*

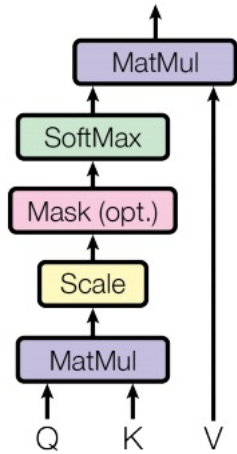
$$\text{score}(s_t, h_i) = \frac{s_t^T h_i}{\sqrt{n}}$$

$$\alpha_{ij} = \frac{\text{score}(s_t, h_i)}{\sum_{k=1}^{T_x} \text{score}(s_t, h_k)}$$

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

# Attention Is All You Need (2017, Vaswani et al.)

Scaled Dot-Product Attention



Query: *decoder hidden state at time step t*

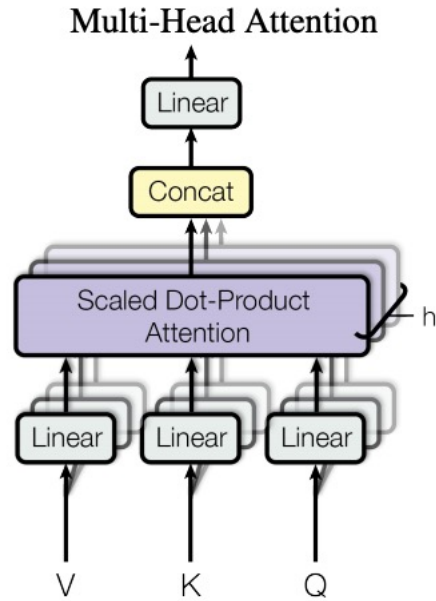
Keys: *(all) encoder cells' hidden state*

Values: *(all) encoder cells' hidden state*

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



## Attention Is All You Need (2017, Vaswani et al.)



$$\text{Multihead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

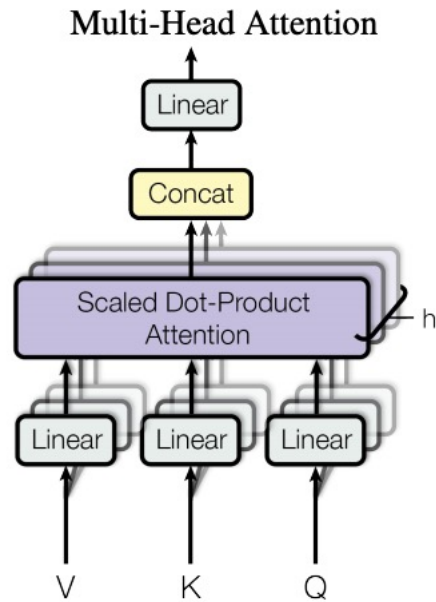
$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$$

$$W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$$

$$d_k = d_v = d_{\text{model}}/h = 64$$

## Attention Is All You Need (2017, Vaswani et al.)



$$QW_i^Q = [d_Q \times d_{model}] \times [d_{model} \times d_k] = [d_Q \times d_k]$$

$$KW_i^K = [d_K \times d_{model}] \times [d_{model} \times d_k] = [d_K \times d_k]$$

$$VW_i^V = [d_V \times d_{model}] \times [d_{model} \times d_v] = [d_V \times d_v]$$

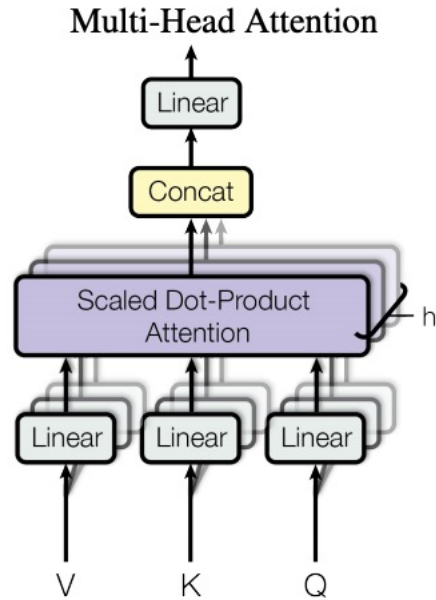


$$Attention(QW_i^Q, KW_i^K, VW_i^V) = [d_V \times d_v]$$



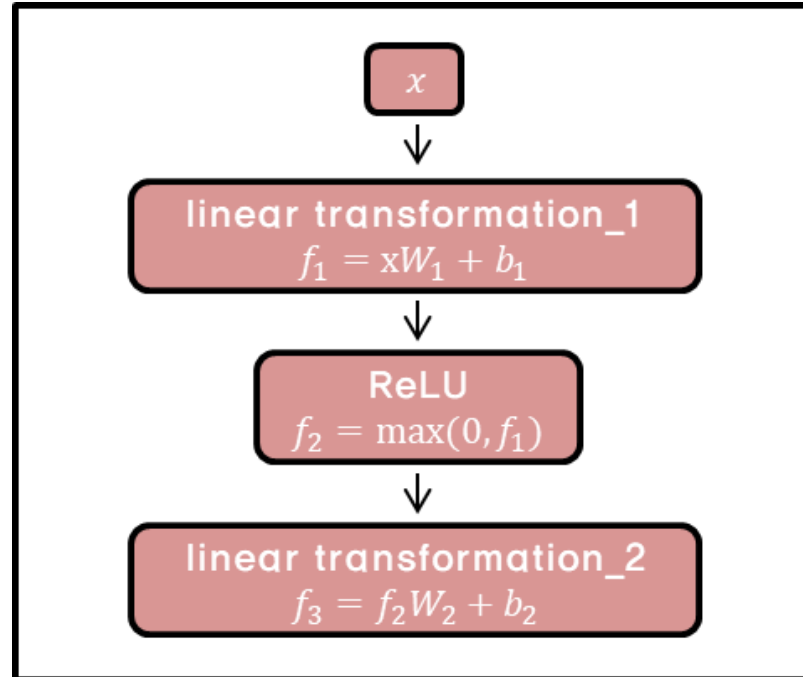
$$Concat(QW_i^Q, KW_i^K, VW_i^V)W^O = [d_V \times h d_v] \times [h d_v \times d_{model}] = [d_V \times d_{model}]$$

# Attention Is All You Need (2017, Vaswani et al.)



## Position-wise Feed-Forward Networks

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



## Attention Is All You Need (2017, Vaswani et al.)

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

# Attention Is All You Need (2017, Vaswani et al.)

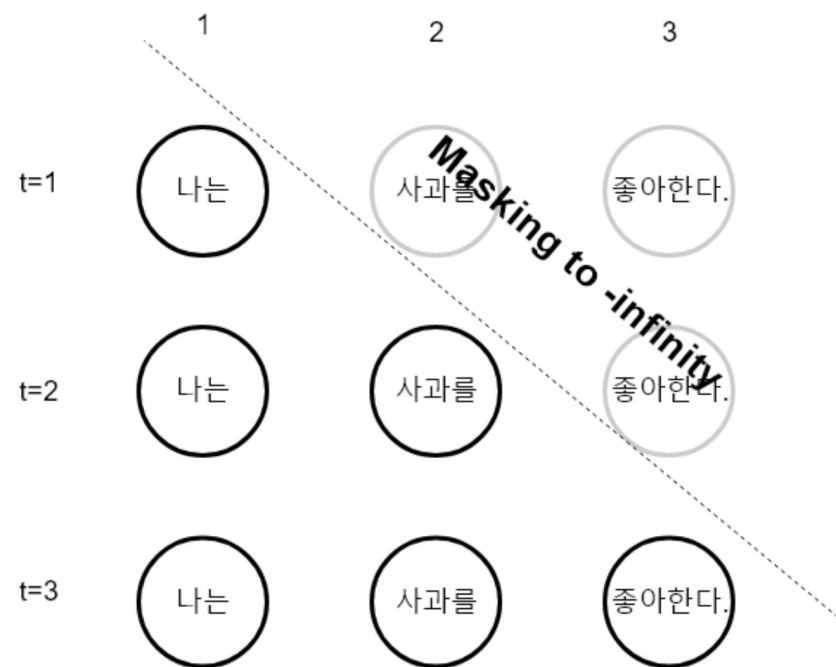
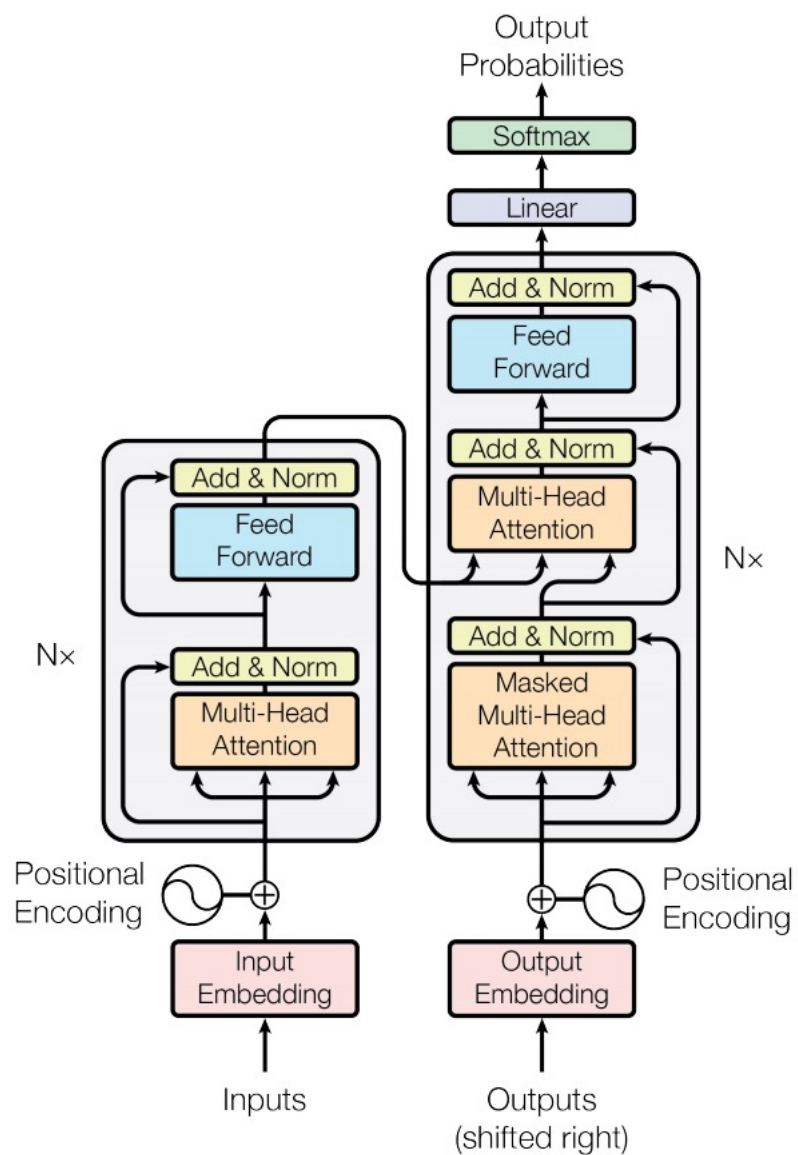
positional encoding

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

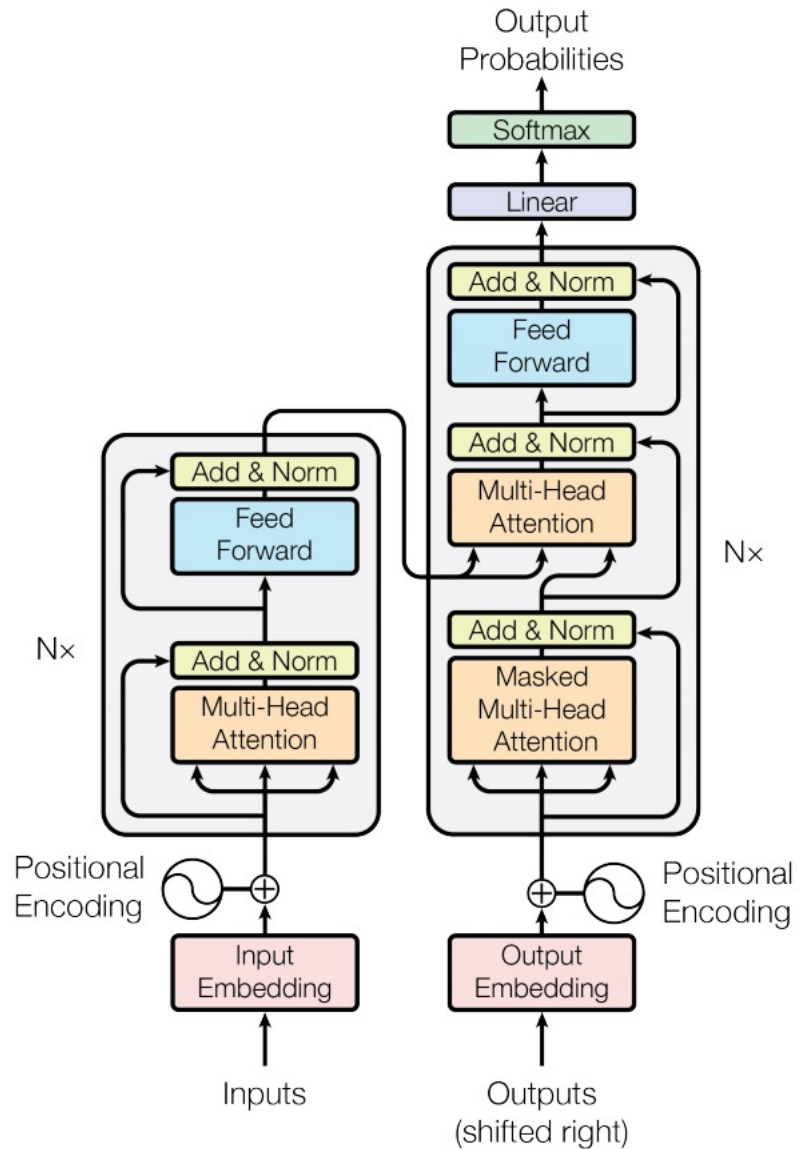
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- allow the model to easily learn to attend by relative positions.
- $PE_{pos+k}$  can be represented as a linear function of  $PE_{pos}$ .
- sinusoidal version allow the model to extrapolate to sequence lengths longer than the ones encountered when training.

# Attention Is All You Need (2017, Vaswani et al.)



# Attention Is All You Need (2017, Vaswani et al.)



- no more long time dependency
- can emphasize values which are “close” to query
- parallelization (multi-head attention)



# 과제

- 예시 코드에서 데이터, rnn cell, attention 바꿔 가며 학습해 보기 (nmae score 비교하기)
- 다른 데이터 써도 됨
- scaled dot attention 구현