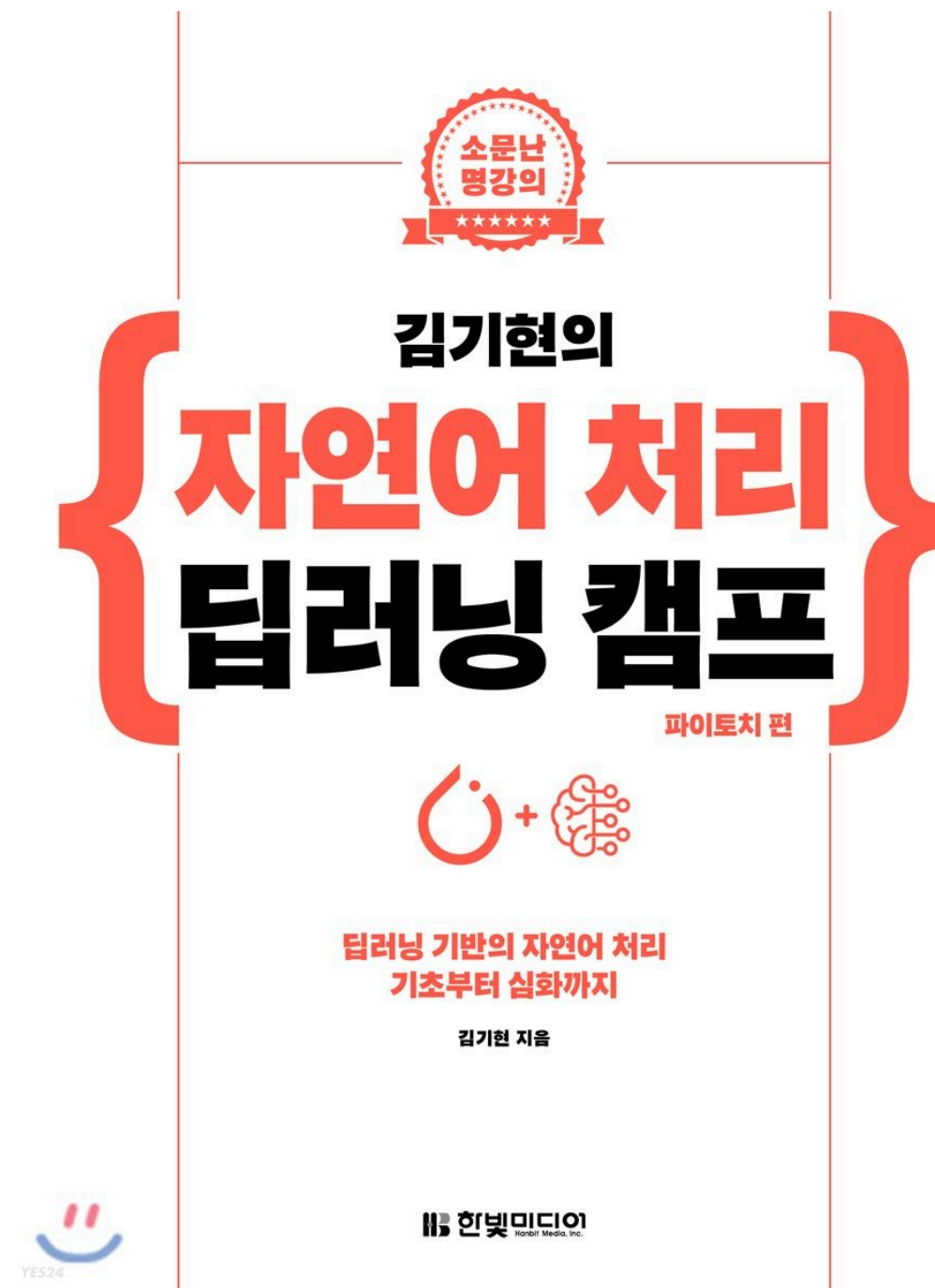


자연어처리 책 정리 발표 2주차



2022.03.25(금)

목차

1. n-gram
2. NNLM
3. Transformer

1. n-gram

1. n-gram

주어진 문장에 대해 어떻게
확률을 구할 수 있을까?

1. n-gram

어쩔, TV

1. n-gram

어쩔, TV

$P(\text{어쩔, TV})$

1. n-gram

어쩔, TV

$P(\text{어쩔, TV})$

$$P(\text{어쩔, TV}) = P(\text{어쩔})P(\text{TV}|\text{어쩔})$$

$$\text{because } P(\text{TV} | \text{어쩔}) = \frac{P(\text{어쩔, TV})}{P(\text{어쩔})}$$

1. n-gram

$$P(\text{어쩔}, TV, \dots, \text{쿠쿠루빙봉}) = P(\text{어쩔})P(TV | \text{어쩔})P(\text{어쩔} | \text{어쩔}, TV) \cdots P(\text{쿠쿠루빙봉} | \text{어쩔}, TV, \text{어쩔}, \text{냉장고})$$

1. n-gram

$$\underline{P(\text{어쩔}, TV, \text{어쩔}, \text{냉장고})} = \underline{P(\text{냉장고} | \text{어쩔}, TV, \text{어쩔})} P(\text{어쩔}, TV, \text{어쩔})$$

1. n-gram

$$\begin{aligned} P(\text{어쩔}, TV, \text{어쩔}, \text{냉장고}) &= P(\text{냉장고} \mid \text{어쩔}, TV, \text{어쩔}) \underline{P(\text{어쩔}, TV, \text{어쩔})} \\ &= P(\text{냉장고} \mid \text{어쩔}, TV, \text{어쩔}) \underline{P(\text{어쩔} \mid \text{어쩔}, TV)} P(\text{어쩔}, TV) \end{aligned}$$

1. n-gram

$$\begin{aligned} P(\text{어쩔}, TV, \text{어쩔}, \text{냉장고}) &= P(\text{냉장고} \mid \text{어쩔}, TV, \text{어쩔})P(\text{어쩔}, TV, \text{어쩔}) \\ &= P(\text{냉장고} \mid \text{어쩔}, TV, \text{어쩔})P(\text{어쩔} \mid \text{어쩔}, TV)\underline{P(\text{어쩔}, TV)} \\ &= P(\text{냉장고} \mid \text{어쩔}, TV, \text{어쩔})P(\text{어쩔} \mid \text{어쩔}, TV)\underline{P(TV \mid \text{어쩔})P(\text{어쩔})} \end{aligned}$$

1. n-gram

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{<i})$$

1. n-gram

$$\log P(w_1, w_2, \dots, w_n) = \sum_{i=1}^n \log P(w_i | w_{<i})$$

1. n-gram

$$P(\text{냉장고} \mid BOS, \text{어쩔}, TV, \text{어쩔}) \approx \frac{\text{Count}(BOS, \text{어쩔}, TV, \text{어쩔}, \text{냉장고})}{\text{Count}(BOS, \text{어쩔}, TV, \text{어쩔})}$$

1. n-gram

$$P(\text{냉장고} \mid BOS, \text{어쩔}, TV, \text{어쩔}) \approx \frac{\text{Count}(BOS, \text{어쩔}, TV, \text{어쩔}, \text{냉장고})}{\text{Count}(BOS, \text{어쩔}, TV, \text{어쩔})}$$

희소성 문제

1. n-gram

마르코프 가정 (Markov assumption)

1. n-gram

$$P(x_i | x_1, x_2, \dots, x_{i-1}) \approx P(x_i | x_{i-k}, \dots, x_{i-1})$$

1. n-gram

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{<i})$$

1. n-gram

$$\begin{aligned} P(w_1, w_2, \dots, w_n) &= \prod_{i=1}^n P(w_i | w_{<i}) \\ &\approx \prod_{i=1}^n P(x_i | x_{i-k}, \dots, x_{i-1}) \end{aligned}$$

1. n-gram

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{<i})$$

$$\approx \prod_{i=1}^n P(x_i | x_{i-k}, \dots, x_{i-1})$$

$$\log P(w_1, w_2, \dots, w_n) = \sum_{i=1}^n \log P(w_i | w_{i-k}, \dots, x_{i-1})$$

1. n-gram

k	n-gram	명칭
0	1-gram	uni-gram
1	2-gram	bi-gram
2	3-gram	tri-gram

1. n-gram

k	n-gram	명칭
0	1-gram	uni-gram
1	2-gram	bi-gram
2	3-gram	tri-gram

1. n-gram

$$P(x_i | x_{i-2}, x_{i-1}) = \frac{\textit{Count}(x_{i-2}, x_{i-1}, x_i)}{\textit{Count}(x_{i-2}, x_{i-1})}$$

1. n-gram

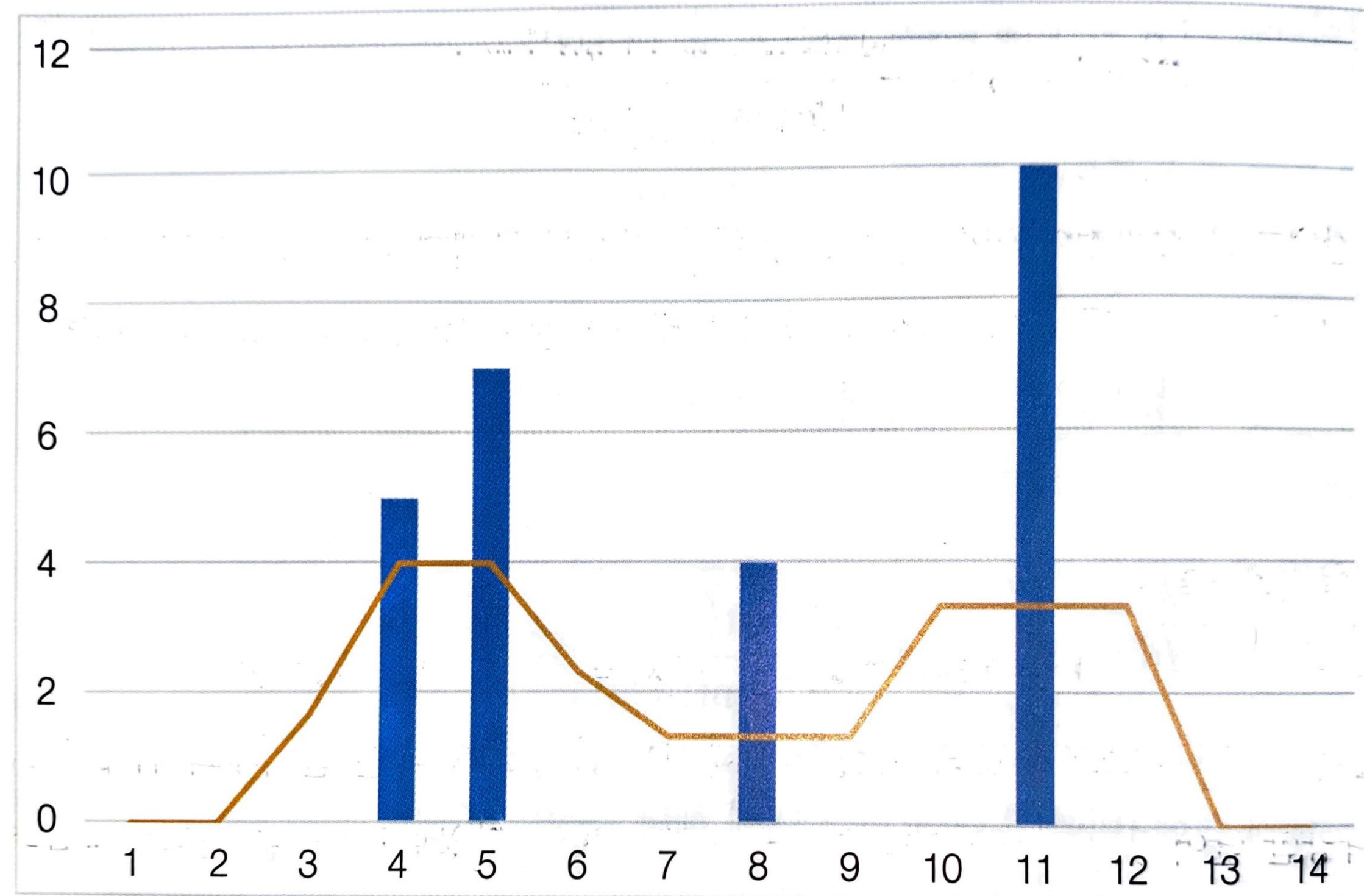
일반화

1. n-gram

출현 횟수를 단순히 확률값으로
추정할 경우 문제가 없을까?

1. n-gram

스무딩과 디스카운팅



▶ 스무딩을 통해 분포의 모양을 좀 더 평탄하게 만들 수 있습니다.

문장 출현 빈도.

1. n-gram

$$P(w_i | w_{<i}) \approx \frac{\textit{Count}(w_{<i}, w_i) + \underline{1}}{\textit{Count}(w_{<i}) + \underline{V}}$$

1. n-gram

$$\begin{aligned} P(w_i | w_{<i}) &\approx \frac{\text{Count}(w_{<i}, w_i) + k}{\text{Count}(w_{<i}) + kV} \\ &\approx \frac{\text{Count}(w_{<i}, w_i) + (m/V)}{\text{Count}(w_{<i}) + m} \end{aligned}$$

1. n-gram

$$P(w_i | w_{<i}) \approx \frac{\textit{Count}(w_{<i}, w_i) + m \underline{P(w_i)}}{\textit{Count}(w_{<i}) + m}$$

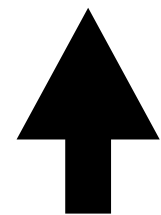
1. n-gram

Kneser-Ney(KN) 디스카운팅

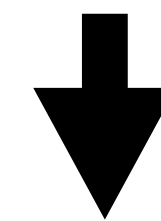
machine learning

deep learning

learning



laptop



1. n-gram

$$Score_{continuation}(w) \propto | \{ v : Count(v, w) > 0 \} |$$

w 와 함께 나타난 v 들의 집합의 크기

1. n-gram

$$Score_{continuation}(w) = \frac{|\{v : Count(v, w) > 0\}|}{\sum_{w' \in W} |\{v : Count(v, w') > 0\}|}$$

w 와 함께 나타난 v 들의 집합의 크기

전체 단어 집합에서 샘플링한 w' 이 v 와 함께 나타난 집합의 크기 합

1. n-gram

출현 빈도를 확률로 근사한 방법과 유사

다양한 단어 뒤에 나타나는 단어의 점수

$$P_{KN}(w_i | w_{i-1}) = \frac{\max(\text{Count}(w_{i-1}, w_i) - d, 0)}{\text{Count}(w_{i-1})} + \lambda(w_{i-1}) \times \text{Score}_{\text{continuation}}(w_i),$$

$$\text{where } \lambda(w_{i-1}) = \frac{d}{\sum_v \text{Count}(w_{i-1}, v)} \times |\{w : c(w_{i-1}, v) > 0\}|$$

1. n-gram

인터폴레이션(보간) interpolation

1. n-gram

$$\tilde{P}(w_n | w_{n-k}, \dots, w_{n-1}) = \underline{\lambda} P_1(w_n | w_{n-k}, \dots, w_{n-1}) + \underline{(1 - \lambda)} P_2(w_n | w_{n-k}, \dots, w_{n-1})$$

where $0 < \lambda < 1$

1. n-gram

일반 영역

- $P(\text{진정제} \mid \text{준비, 된}) = 0.00001$
- $P(\text{사나이} \mid \text{준비, 된}) = 0.01$

1. n-gram

일반 영역

- $P(\text{진정제} \mid \text{준비, 된}) = 0.000001$
- $P(\text{사나이} \mid \text{준비, 된}) = 0.01$

특화 영역

- $P(\text{진정제} \mid \text{준비, 된}) = 0.09$
- $P(\text{약} \mid \text{준비, 된}) = 0.04$

1. n-gram

일반 영역

- $P(\text{진정제} \mid \text{준비, 된}) = 0.000001$
- $P(\text{사나이} \mid \text{준비, 된}) = 0.01$

특화 영역

- $P(\text{진정제} \mid \text{준비, 된}) = 0.09$
- $P(\text{약} \mid \text{준비, 된}) = 0.04$

인터폴레이션 결과

$$- P(\text{진정제} \mid \text{준비, 된}) = \underbrace{0.5}_{\lambda} * 0.09 + \underbrace{(1 - 0.5)}_{1 - \lambda} * 0.000001 = 0.0450005$$

1. n-gram

백오프(back-off)

1. n-gram

$$\begin{aligned}\tilde{P}(w_n | w_{n-k}, \dots, w_{n-1}) = & \lambda_1 P(w_n | w_{n-k}, \dots, w_{n-1}) \\ & + \lambda_2 P_2(w_n | w_{n-k}, \dots, w_{n-1}) \\ & + \dots \\ & + \lambda_k P(w_n),\end{aligned}$$

$$\text{where } \sum_i \lambda_i = 1$$

2. NNLM

2. NNLM

고양이는 좋은 반려동물입니다.

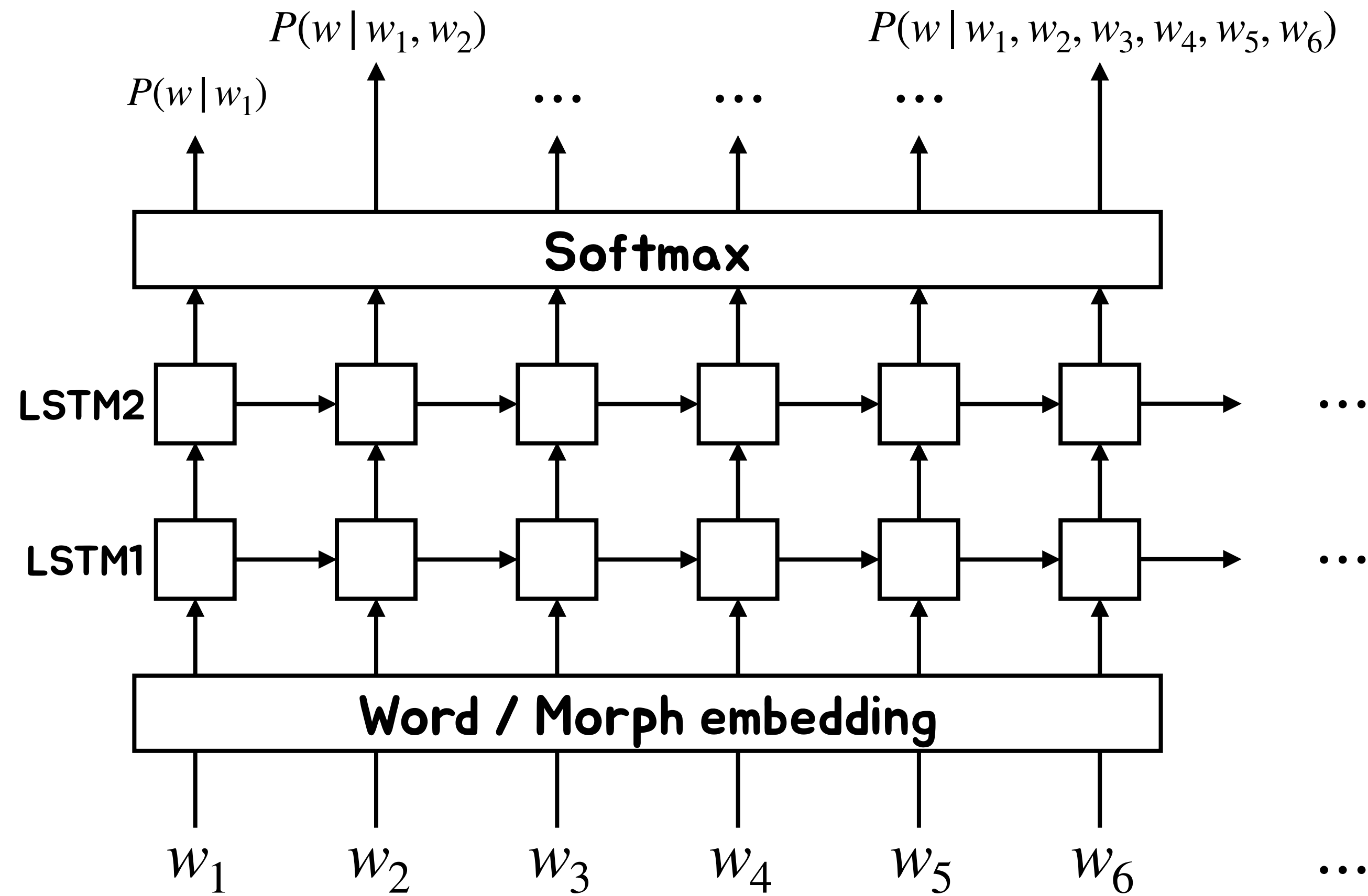
2. NNLM

고양이는 좋은 반려동물입니다.

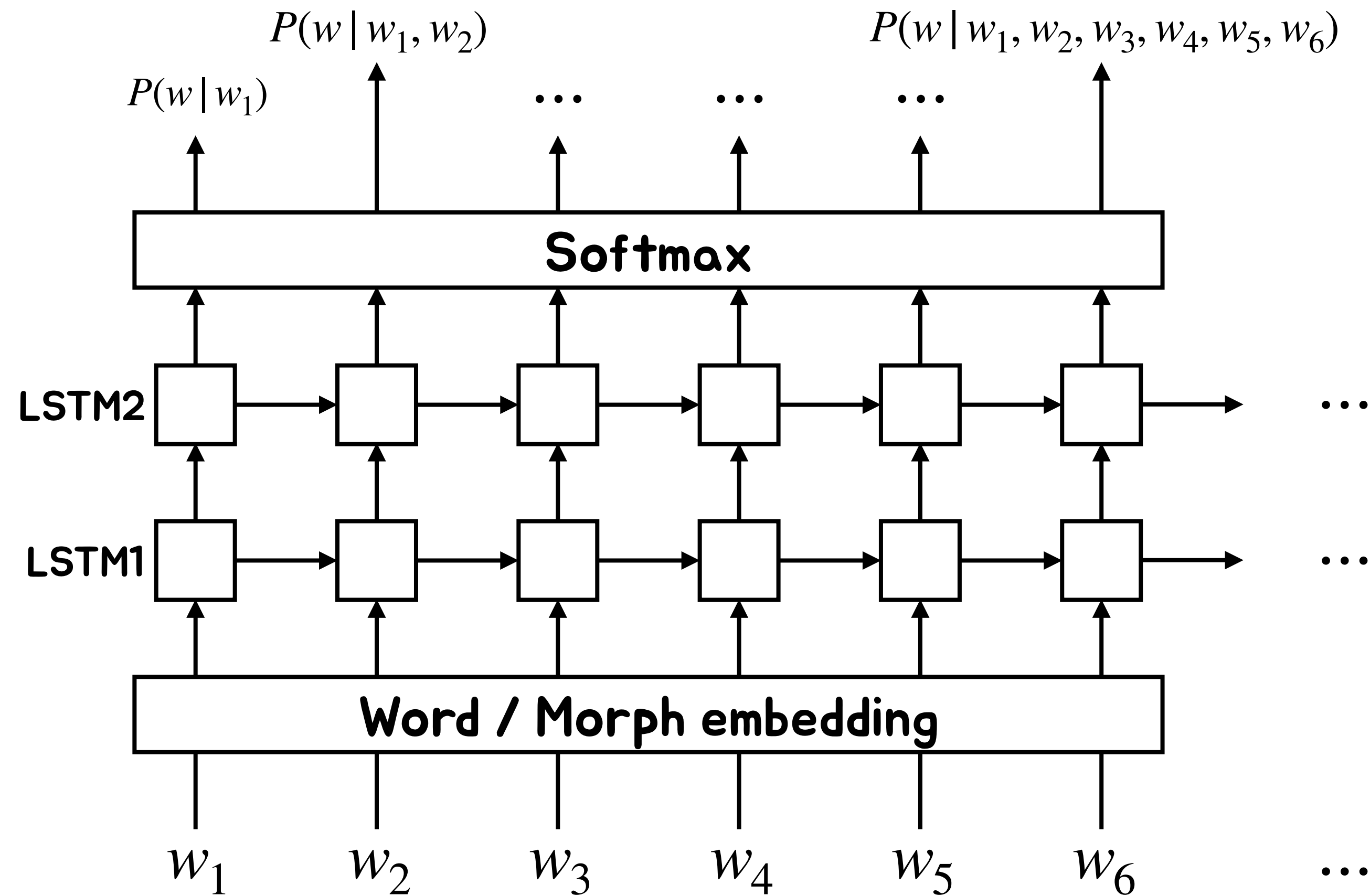
$P(\text{반려동물} \mid \text{강아지는, 좋은})$

$P(\text{반려동물} \mid \text{자동차는, 좋은})$

2. NNLM

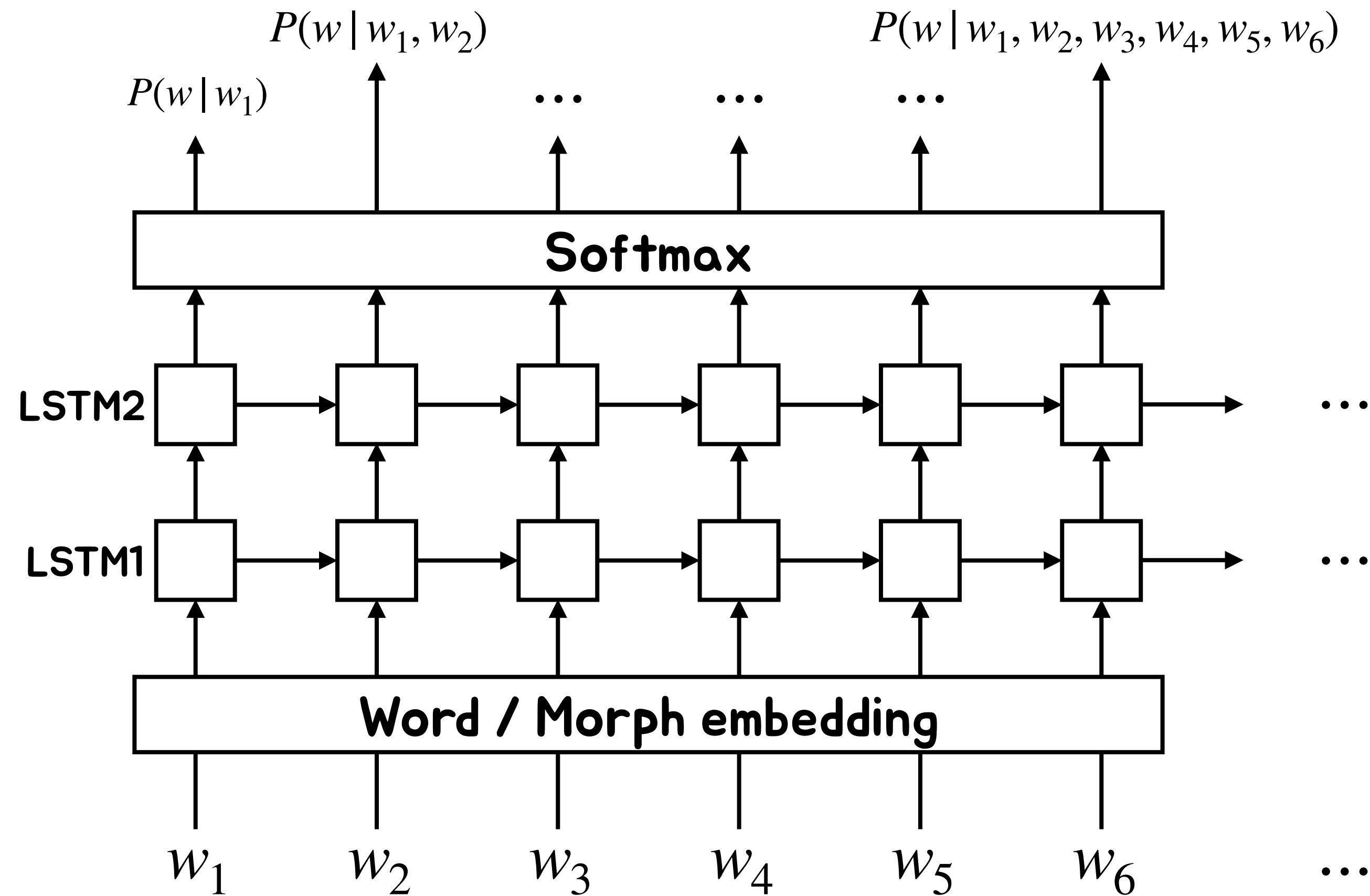


2. NNLM



$$P(w_1, w_2, \dots, w_k) = \prod_{i=1}^k P(w_i | w_{<i})$$

2. NNLM



$$\log P(w_1, w_2, \dots, w_k) = \sum_{i=1}^k \log P(w_i | w_{<i})$$

2. NNLM

$$x_{1:n} = \{x_0, x_1, \dots, x_n, x_{n+1}\}$$

where $x_0 = BOS$ and $x_{n+1} = EOS$

$$\hat{x}_{i+1} = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_i))))$$

$$\hat{x}_{1:n}[1:] = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_{1:n}[: - 1])))),$$

$$\text{linear}_{d_1 \rightarrow d_2}(x) = Wx + b \text{ where } W \in \mathbb{R}^{d_1 \times d_2} \text{ and } b \in \mathbb{R}^{d_2},$$

and hidden_size is dimension of hidden state and $|V|$ is size of vocabulary.

2. NNLM

$$x_{1:n} = \{x_0, x_1, \dots, x_n, x_{n+1}\}$$

where $x_0 = BOS$ and $x_{n+1} = EOS$

$$\hat{x}_{i+1} = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_i))))$$

$$\hat{x}_{1:n}[1:] = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_{1:n}[: - 1])))),$$

$$\text{linear}_{d_1 \rightarrow d_2}(x) = Wx + b \text{ where } W \in \mathbb{R}^{d_1 \times d_2} \text{ and } b \in \mathbb{R}^{d_2},$$

and hidden_size is dimension of hidden state and $|V|$ is size of vocabulary.

2. NNLM

$$x_{1:n} = \{x_0, x_1, \dots, x_n, x_{n+1}\}$$

where $x_0 = BOS$ and $x_{n+1} = EOS$

$$\hat{x}_{i+1} = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_i))))$$

$$\hat{x}_{1:n}[1:] = \text{softmax}(\text{linear}_{\text{hidden_size} \rightarrow |V|}(\text{RNN}(\text{emb}(x_{1:n}[: - 1])))),$$

$$\text{linear}_{d_1 \rightarrow d_2}(x) = Wx + b \text{ where } W \in \mathbb{R}^{d_1 \times d_2} \text{ and } b \in \mathbb{R}^{d_2},$$

and *hidden_size* is dimension of hidden state and $|V|$ is size of vocabulary.

2. NNLM

$$x_{1:n}[: - 1] = \{x_0, x_1, \dots, x_n\}$$

$$x_{emb} = emb(x_{1:n}[1 : - 1])$$

EOS 제외한 BOS + 입력문장(n)의 길이

where $|x_{1:n}[: - 1]| = (batch_size, \underline{n + 1})$

and $|x_{emb}| = (batch_size, n + 1, word_vec_dim)$

2. NNLM

$$h_{0:n} = RNN(x_{emb})$$

Where $|h_{0:n}| = (batch_size, n + 1, hidden_size)$

2. NNLM

$$\hat{x}_{1:n} = \text{softmax} \left(\text{linear}_{\text{hidden_size} \rightarrow |V|}(h_{0:n}) \right)$$

Where $|\hat{x}_{1:n}| = (\text{batch_size}, n + 1, |V|)$

and $x_{1:n}[1 :] = \{x_1, x_2, \dots, x_{n+1}\}$

2. NNLM

$$\mathcal{L}(\hat{x}_{1:n}, x_{1:n}[1 :]) = -\frac{1}{n} \sum_{i=1}^m \sum_{j=1}^{n+1} x_j^i \log \hat{x}_j^i$$

Where x_j^i is one – hot vector

2. NNLM

n-gram vs NNLM



2. NNLM

n-gram vs NNLM

 n-gram 성능  NNLM

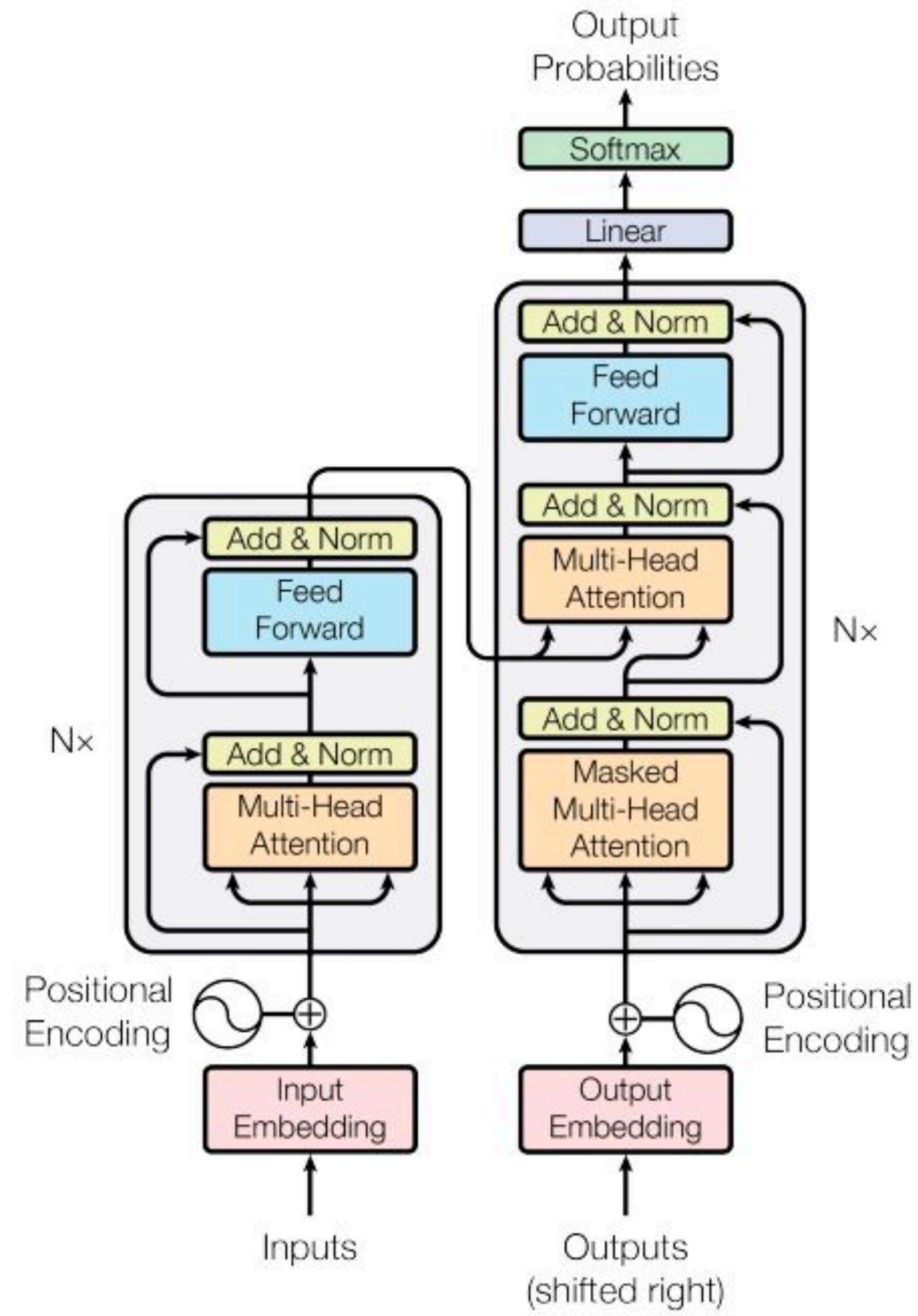
  n-gram 연산량  NNLM

3. Transformer

3. Transforemr

Attention is all you need

3. Transforemr



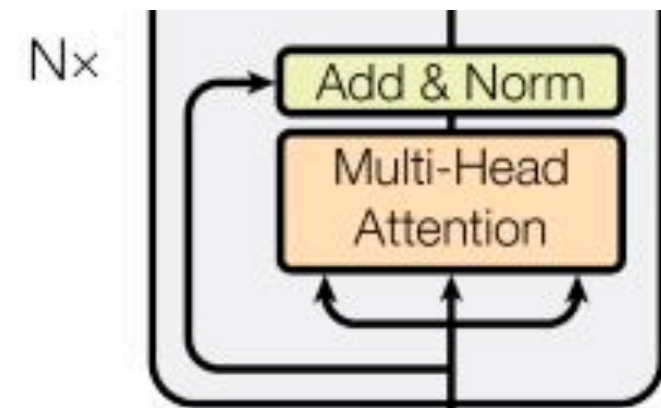
3. Transforemr

Sub Module

3. Transforemr

Sub Module

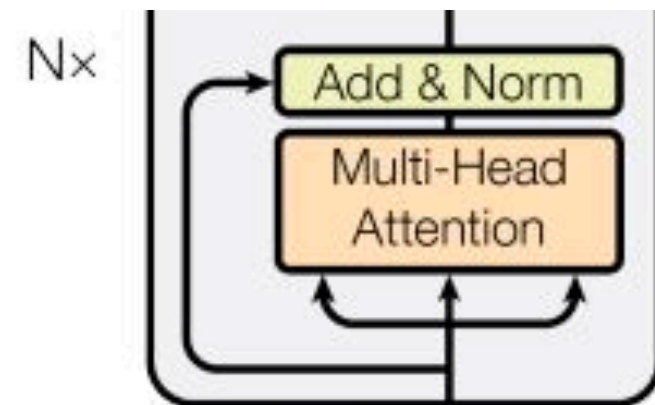
self-attention



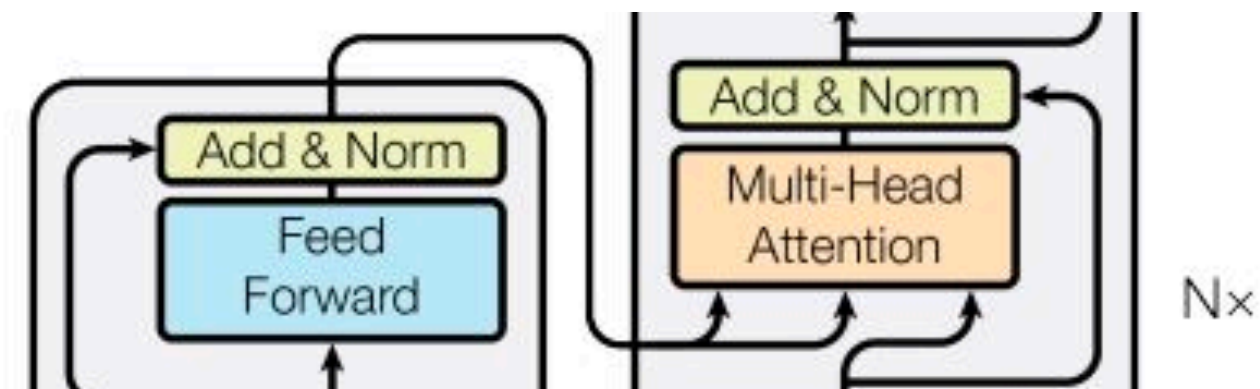
3. Transforemr

Sub Module

self-attention



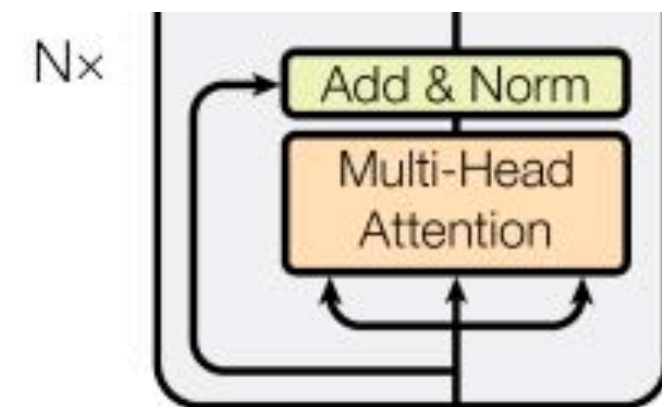
attention



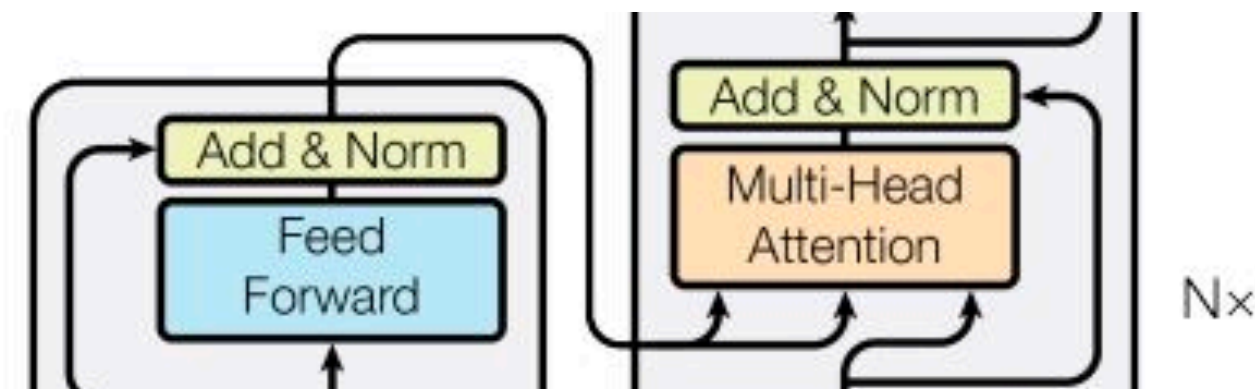
3. Transforemr

Sub Module

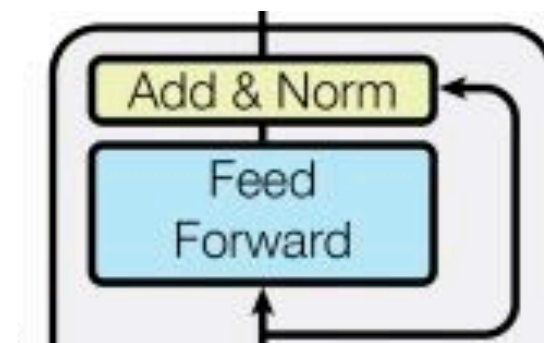
self-attention



attention

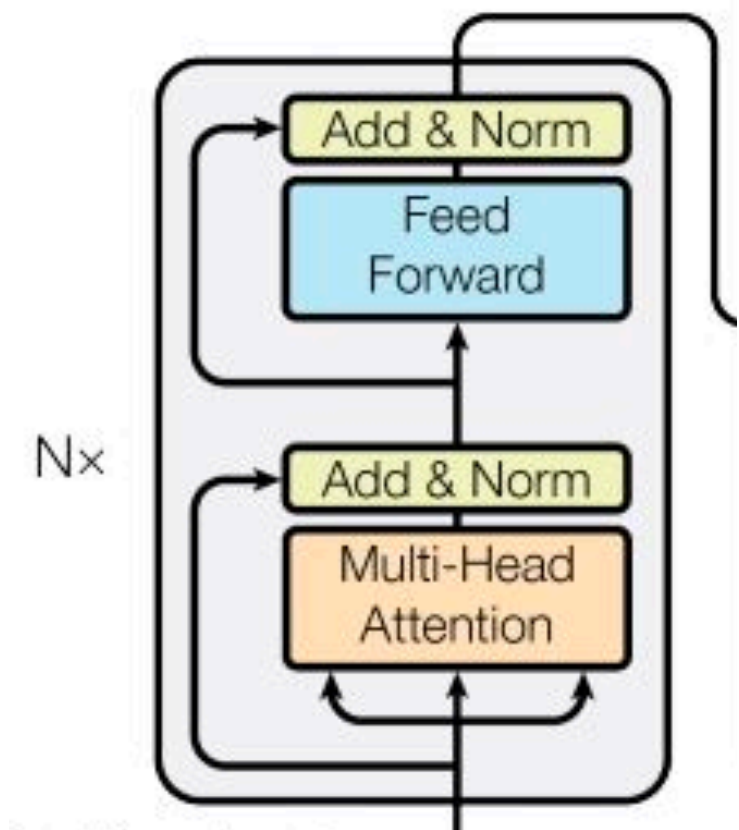


Feed Forward Layer



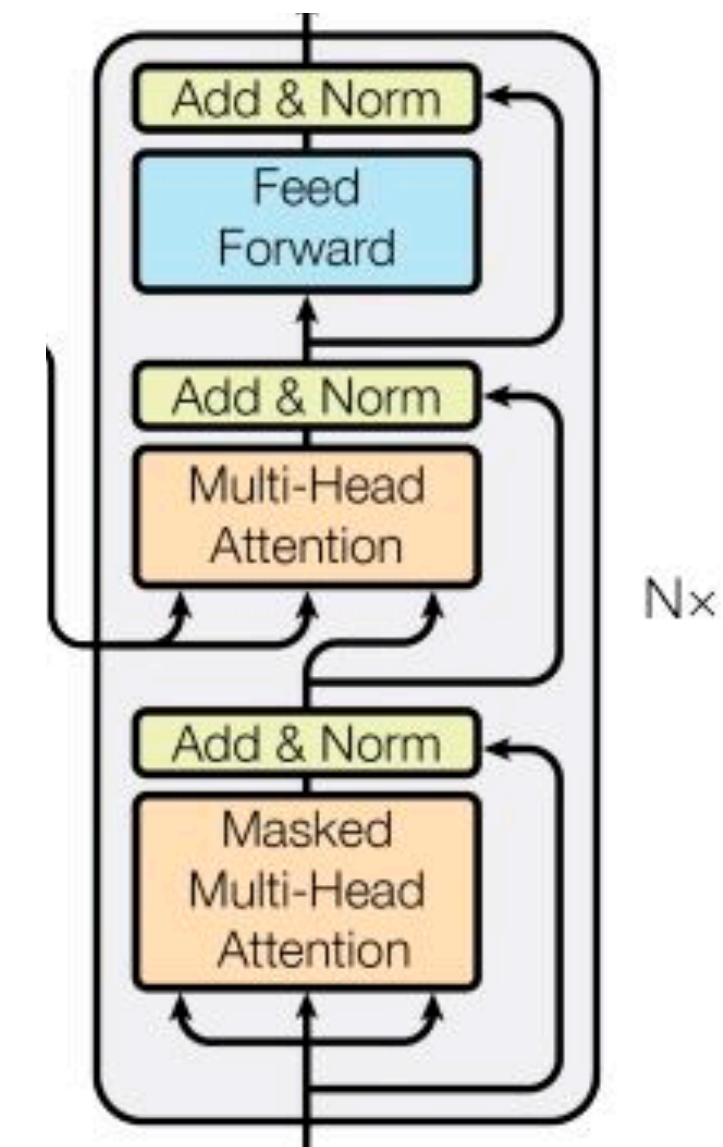
3. Transforemr

Encoder



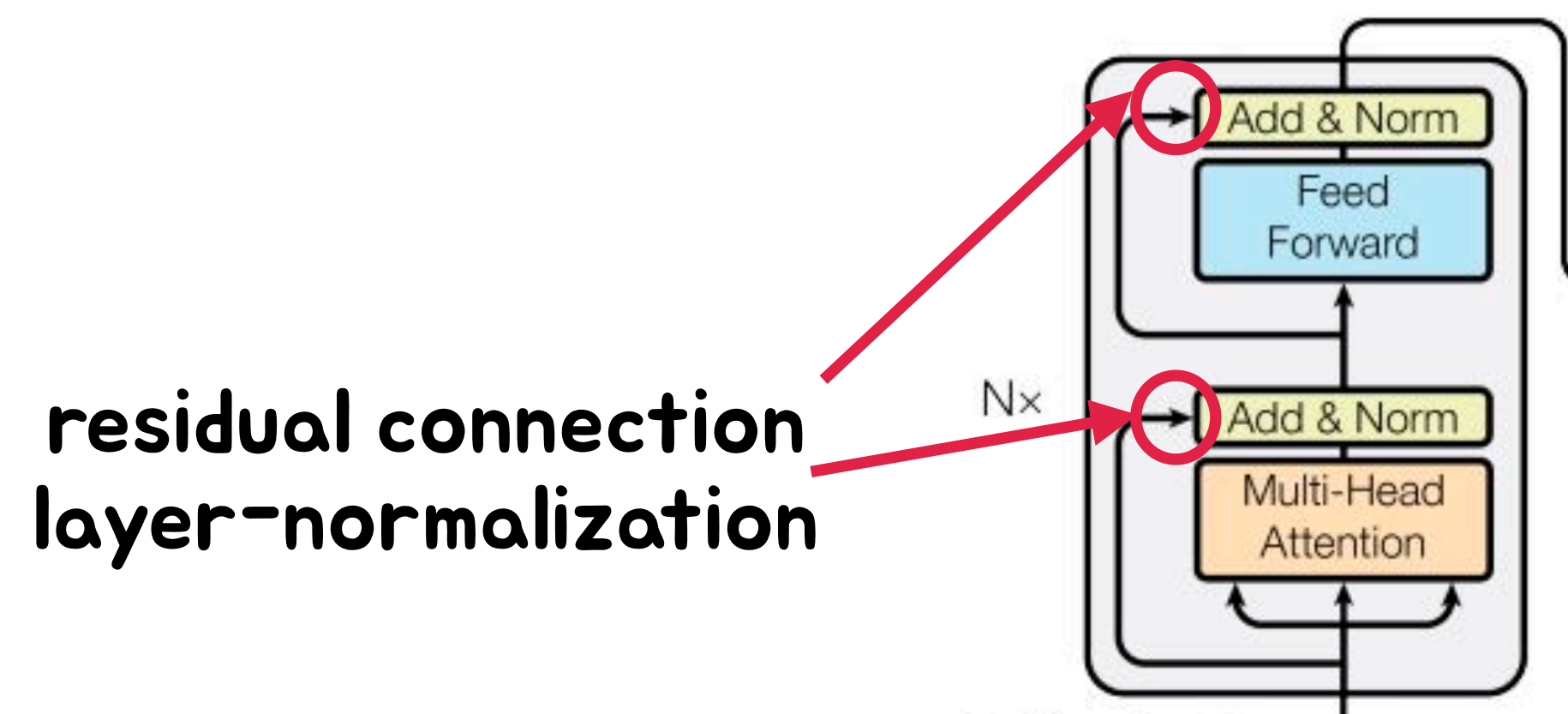
3. Transforemr

Decoder

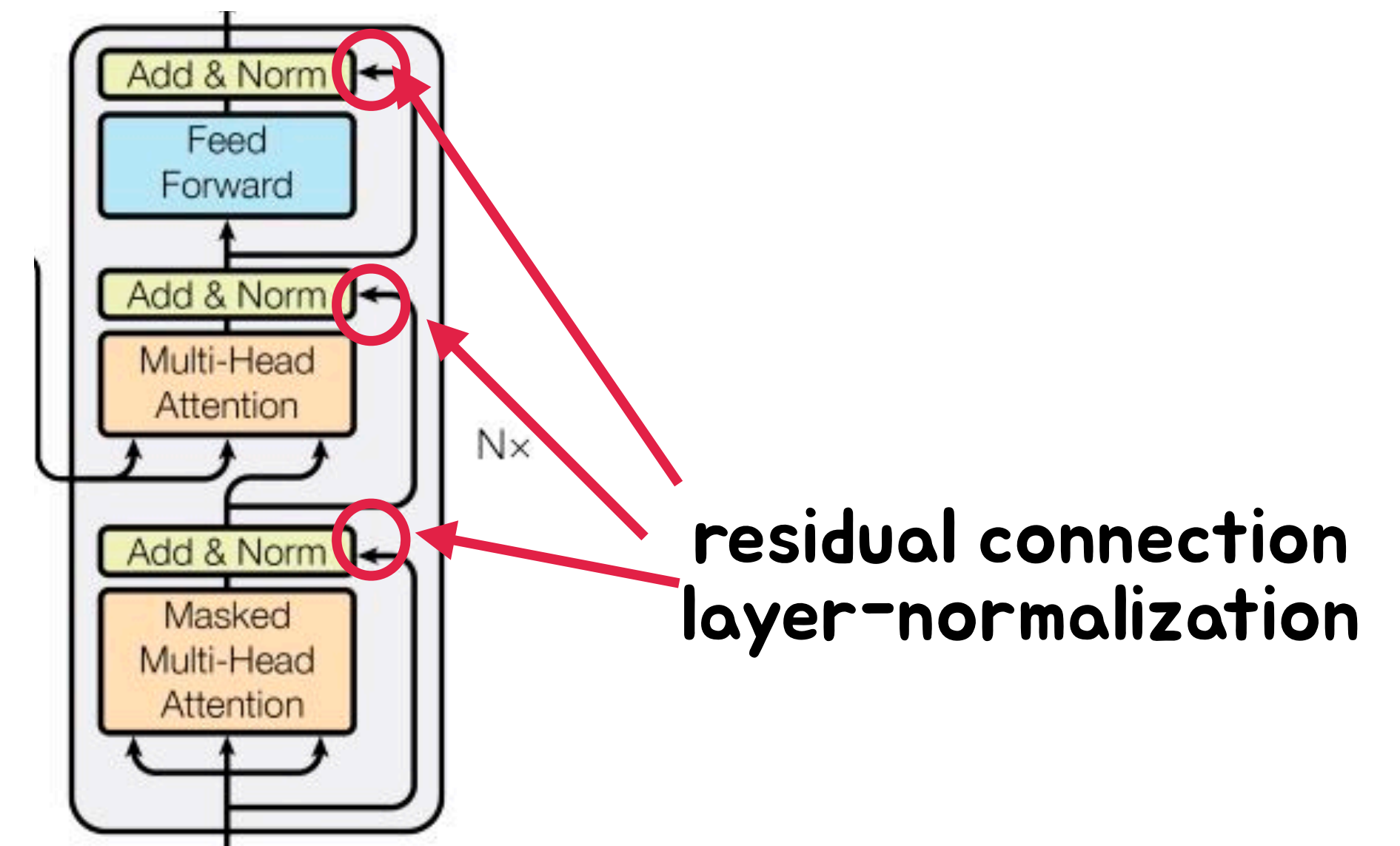


3. Transforemr

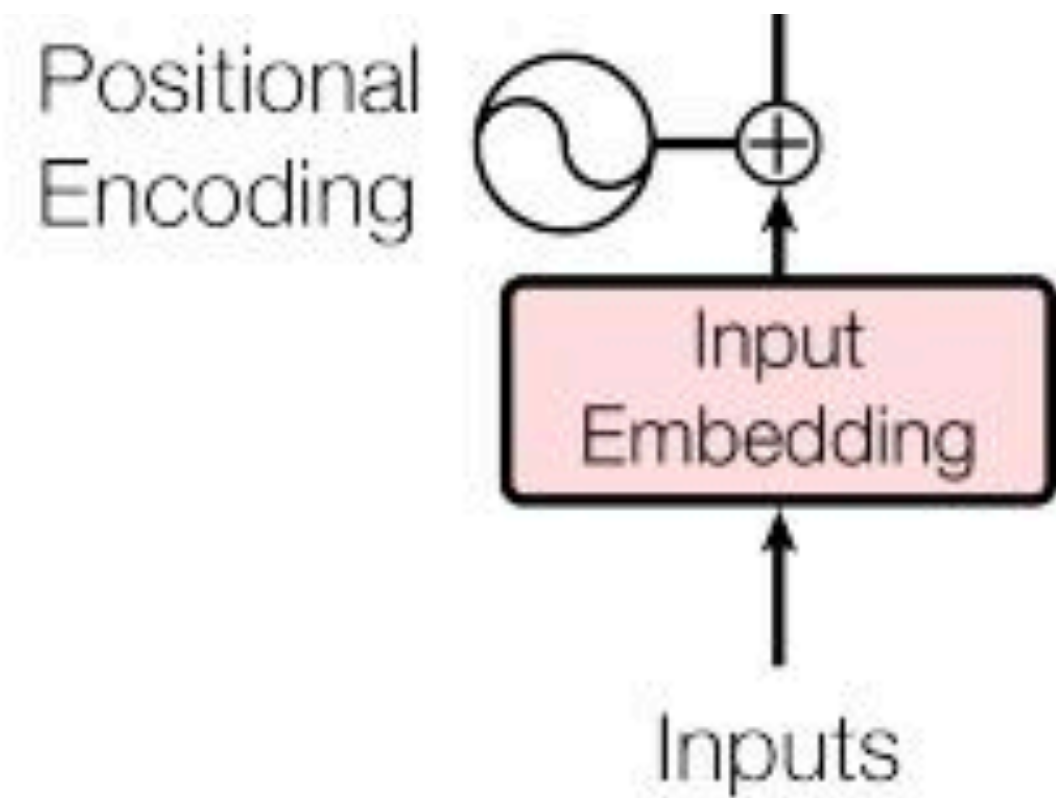
Encoder



Decoder



Positional Encoding



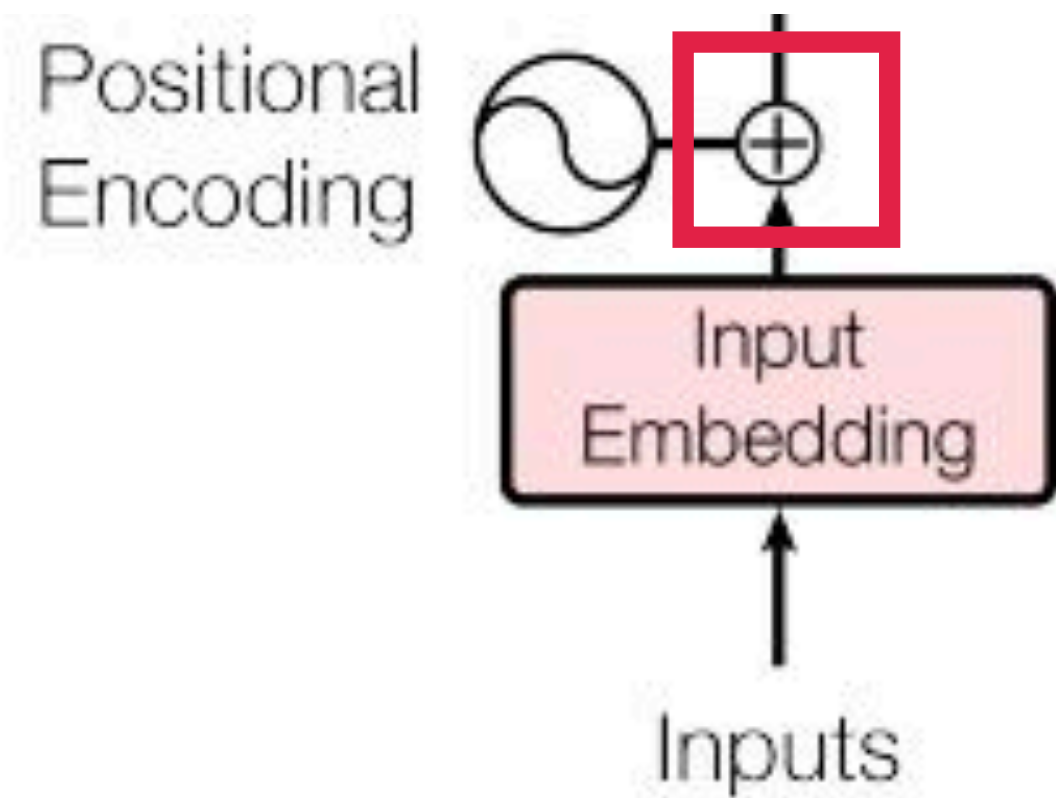
3. Transforemr

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

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

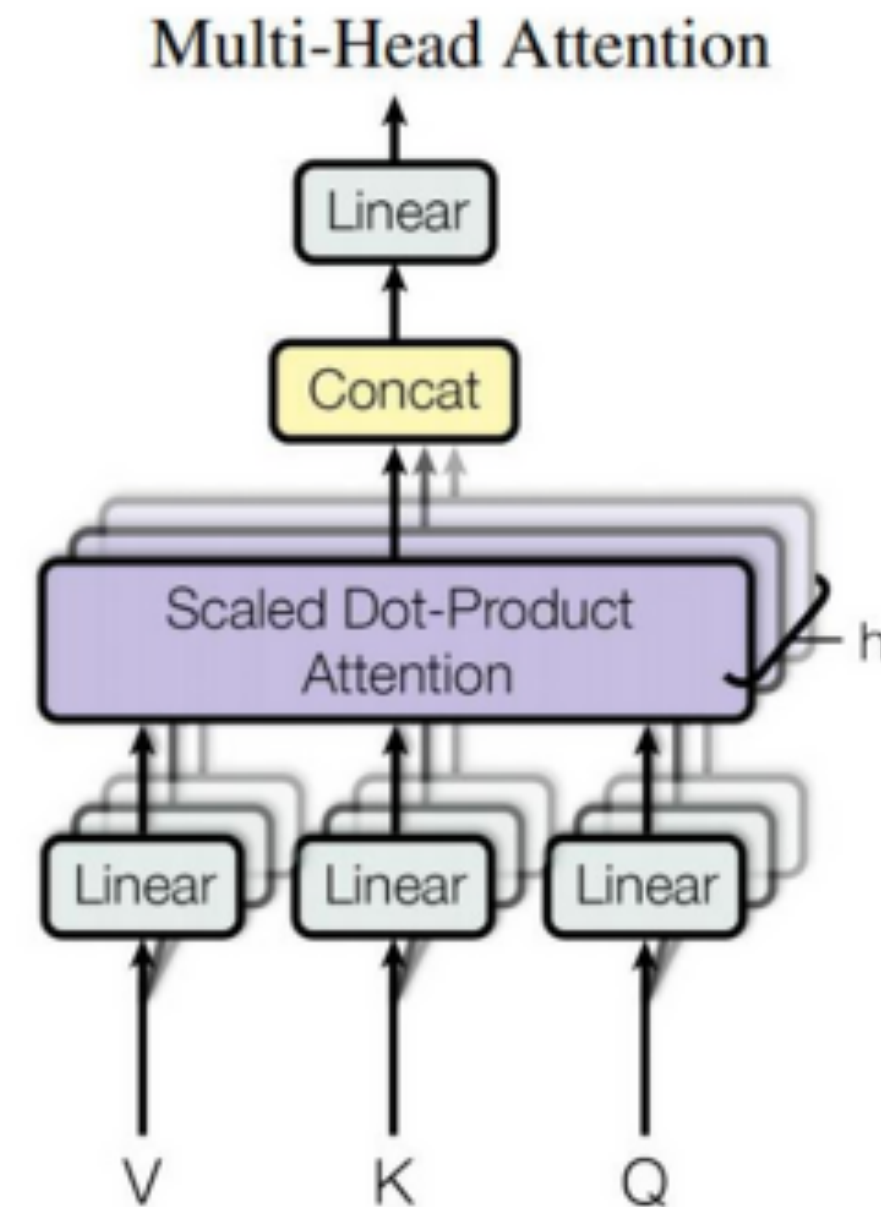


Positional Encoding



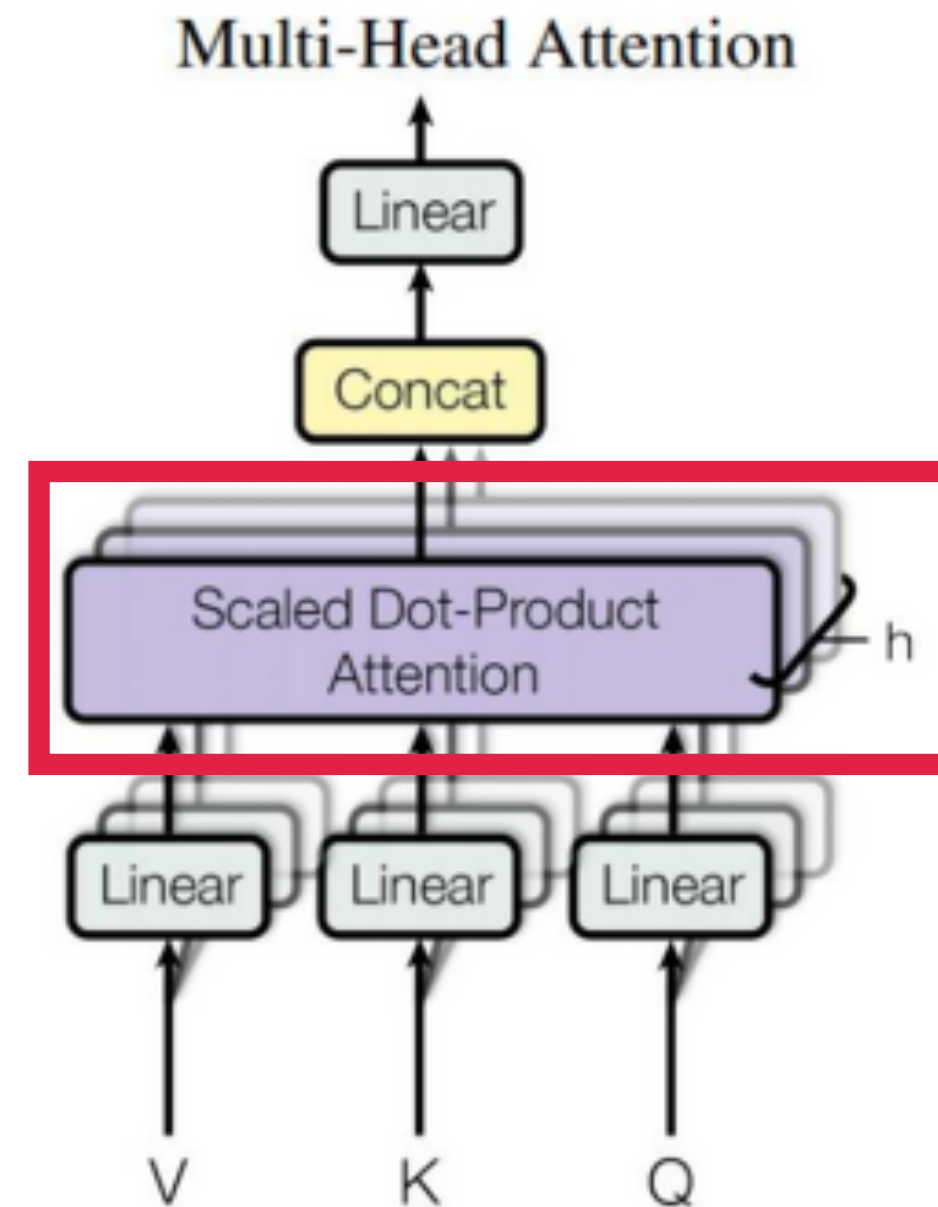
3. Transforemr

Multi-Head Attention



3. Transforemr

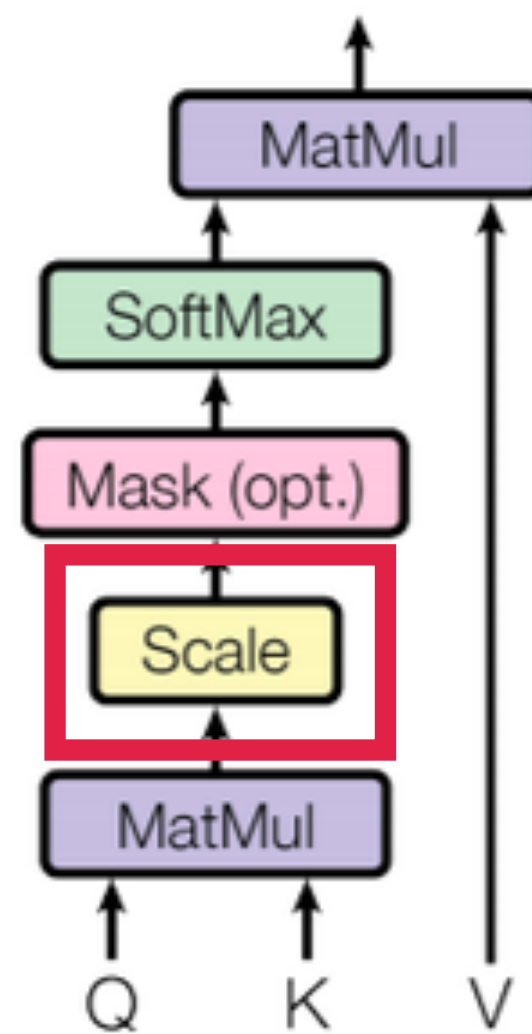
Multi-Head Attention



3. Transforemr

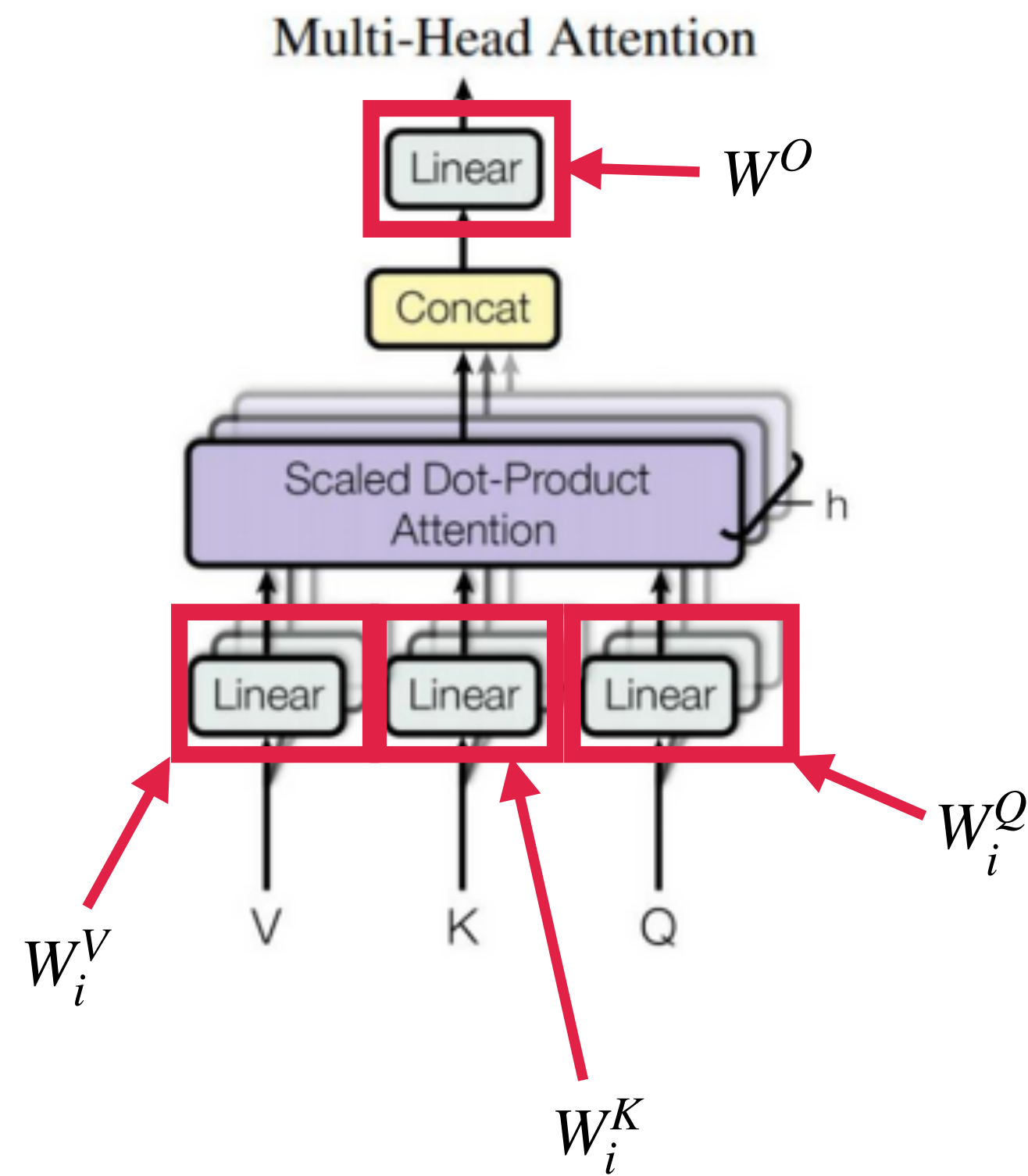
Scaled Dot-Product Attention

Scaled Dot-Product Attention



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

3. Transforemr



concat 을 의미

$$MultiHead(Q, K, V) = [head_1; head_2; \cdots; head_h]W^O$$

where $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$

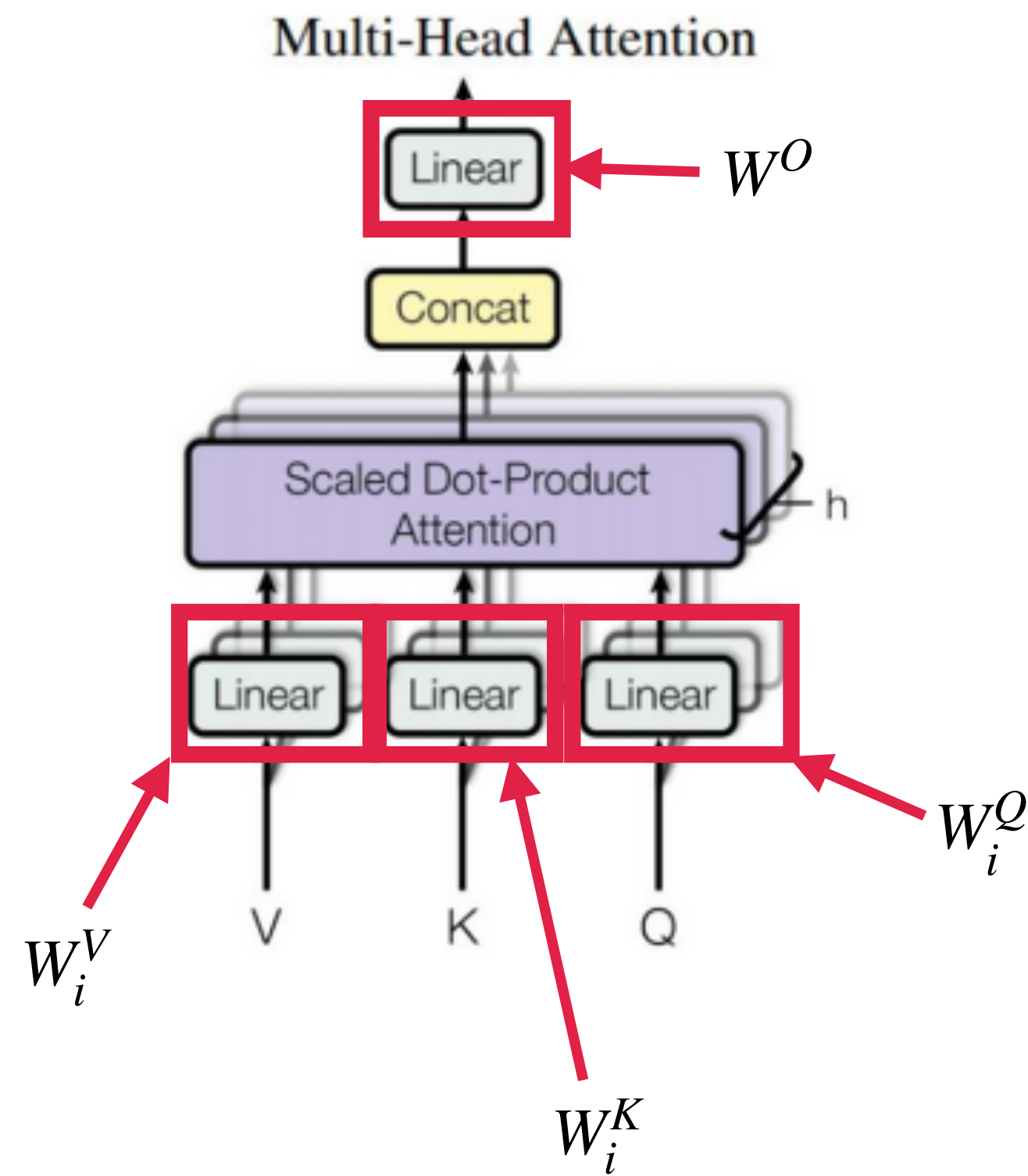
3. Transforemr

$$|Q| = (batch_size, \underline{m}, hidden_size)$$

$$|K| = |V| = (batch_size, \underline{n}, hidden_size)$$

where n is length of source sentence, and m is length of target sentence .

3. Transforemr



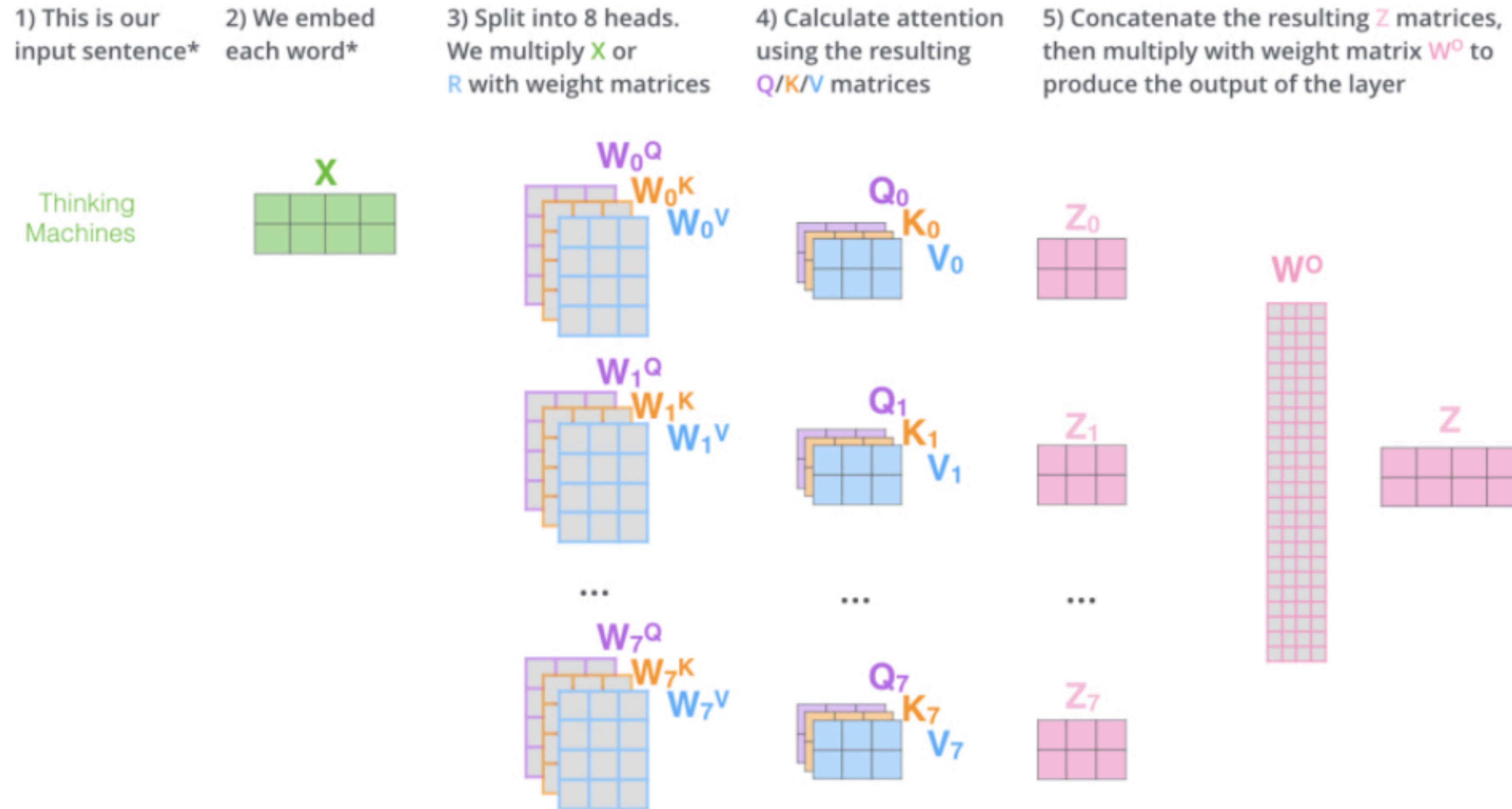
$$|W_i^Q| = |W_i^K| = |W_i^V| = (hidden_size, head_size)$$

$$|W^O| = (head_size \times h, hidden_size)$$

where $hidden_size = 512$, $h = 8$ and $hidden_size = head \times h$

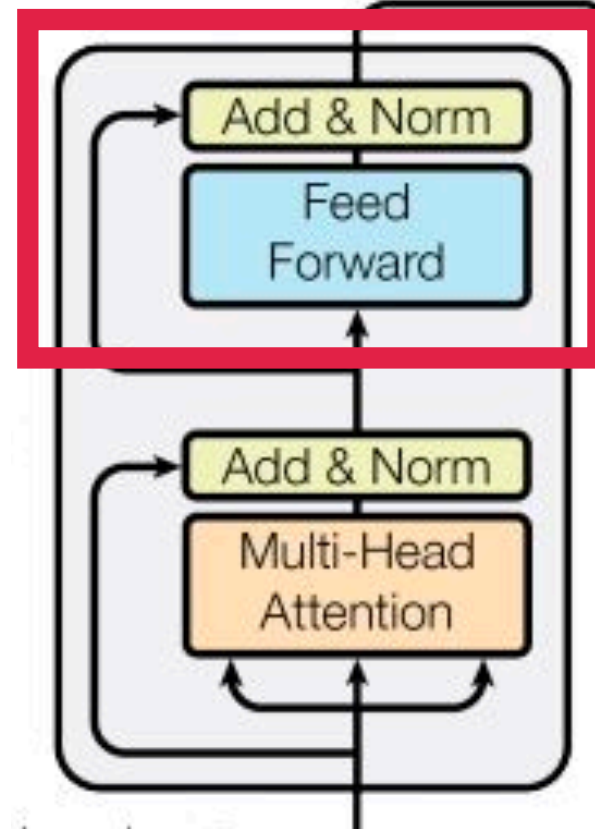
3. Transforemr

Multi-Head Attention



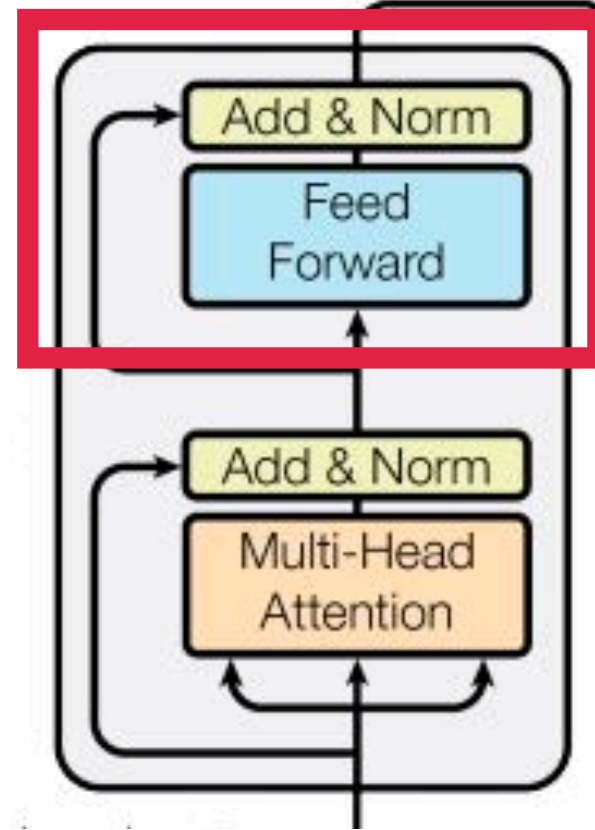
3. Transforemr

Feed Forward Layer



3. Transforemr

Feed Forward Layer



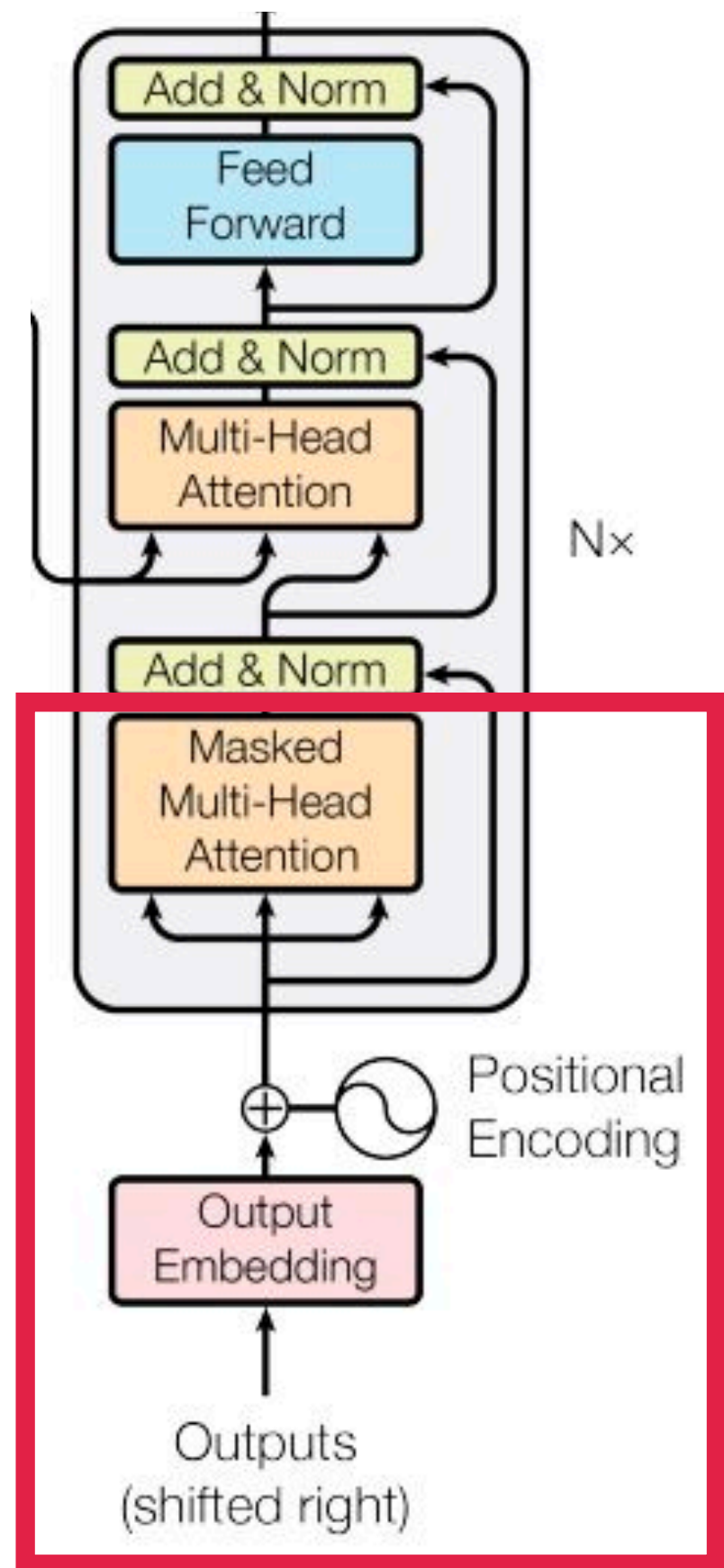
$$FFN(x) = ReLU(xW_1 + b_1)W_2 + b_2$$

where $|x| = (batch_size, n, hidden_size)$

and $w_1 \in \mathbb{R}^{hidden_size \times d_{ff}}$, $W_2 \in \mathbb{R}^{d_{ff} \times hidden_size}$ and $d_{ff} = 2048$

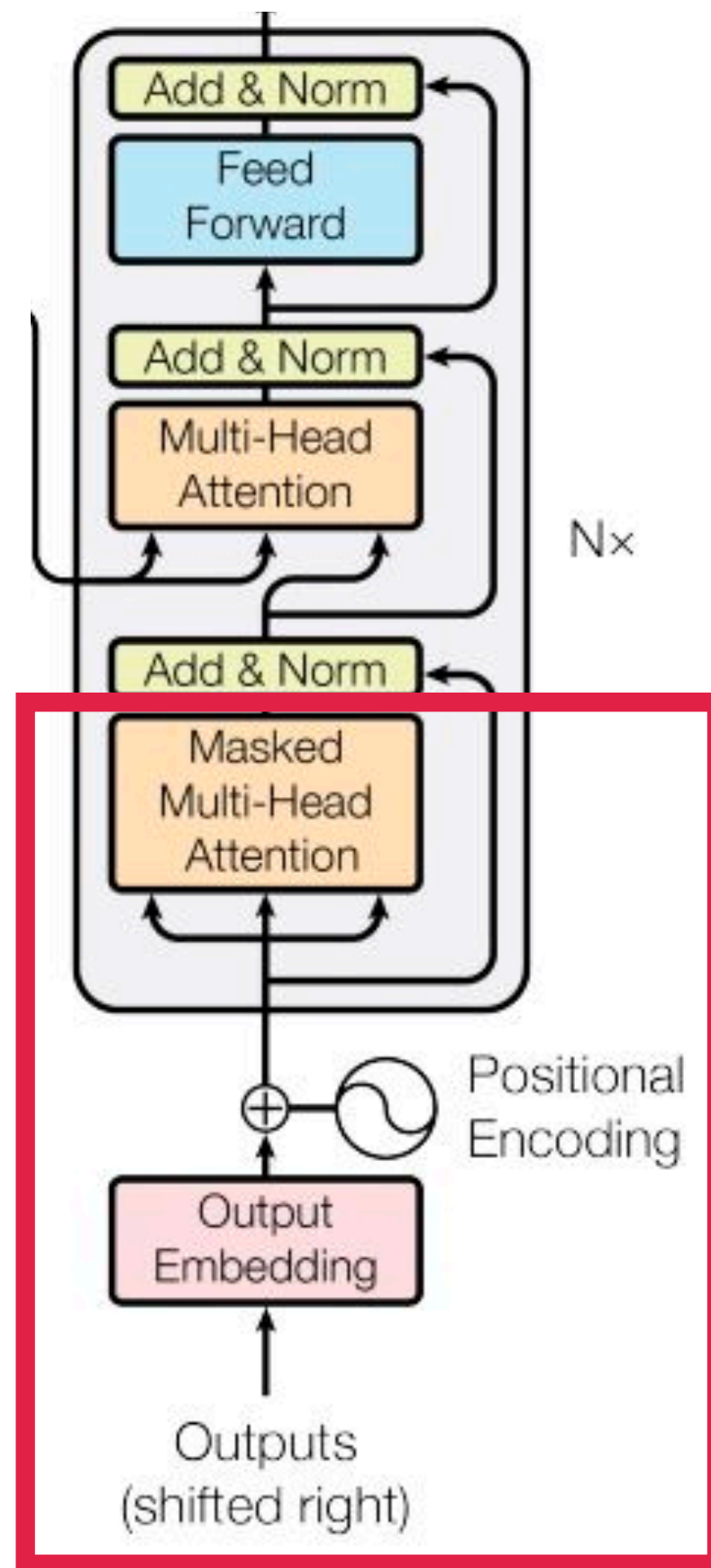
3. Transforemr

Masked Multi-Head Attention



3. Transforemr

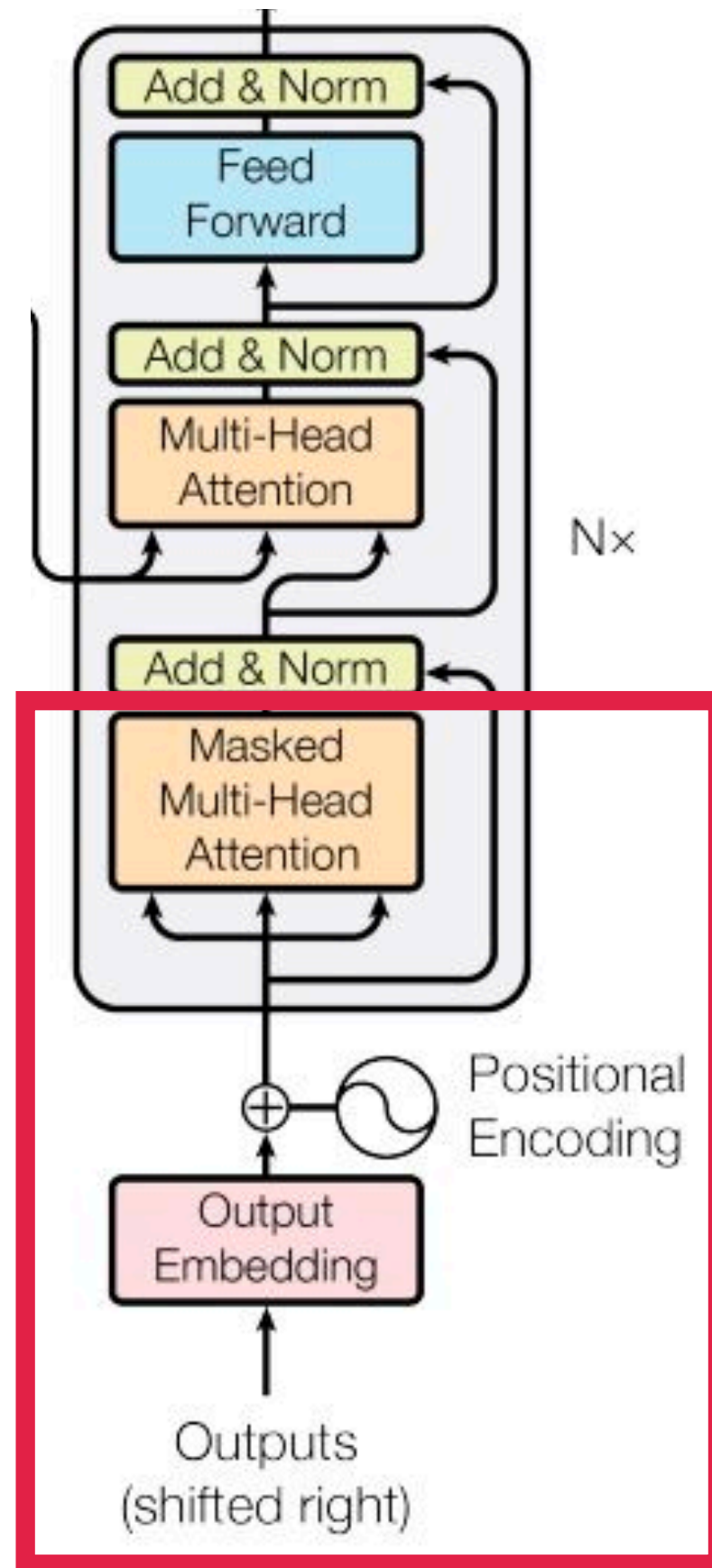
Masked Multi-Head Attention



문제 발생!!

3. Transforemr

Masked Multi-Head Attention



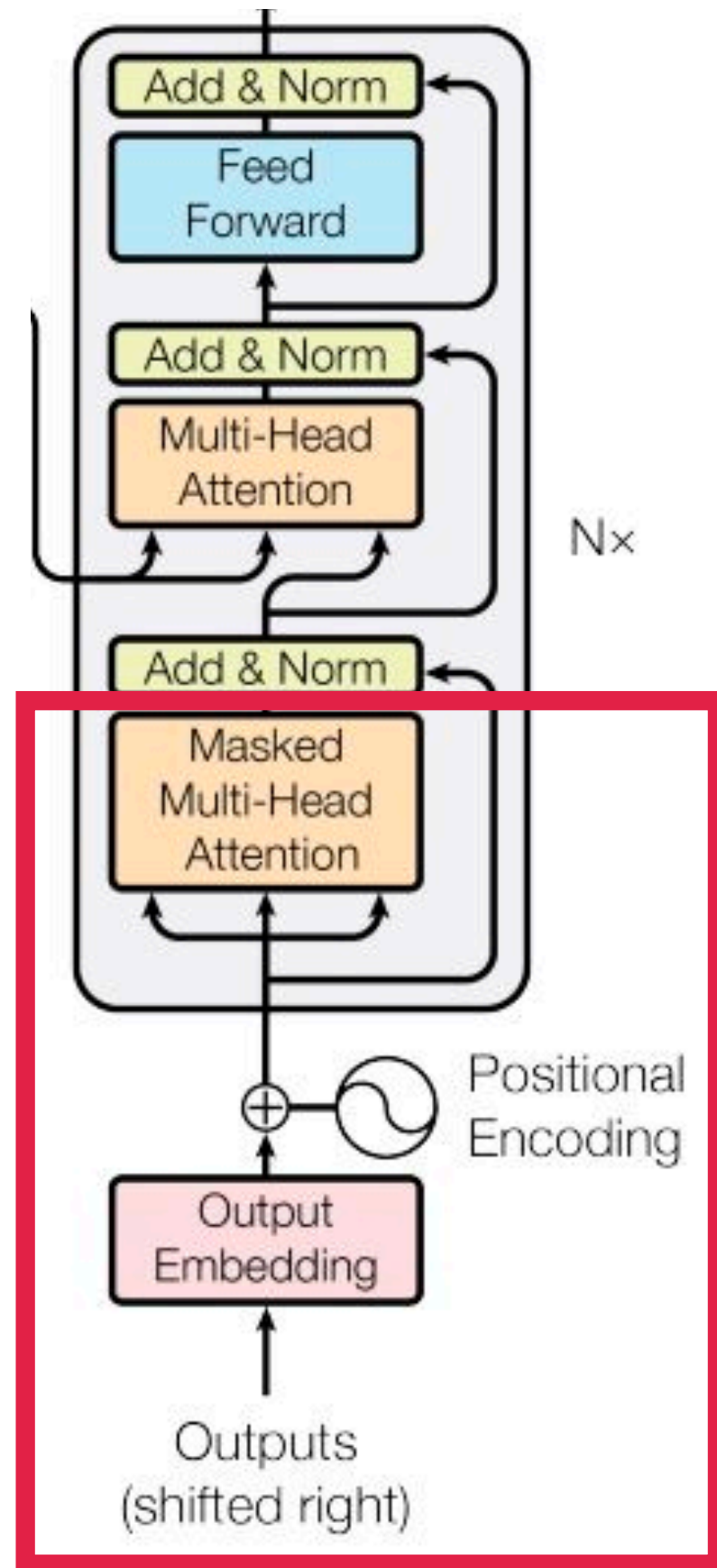
문제 발생!!

$$\begin{matrix} & Q \\ & \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} < sos > \\ je \\ suis \\ étudiant \end{matrix} & \begin{matrix} \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \end{matrix} \times K^T \begin{matrix} & \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \end{matrix} = \begin{matrix} & \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} & & & & \end{matrix} \end{matrix} \begin{matrix} < sos > & je & suis & étudiant \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{matrix}$$

Attention Score Matrix

3. Transforemr

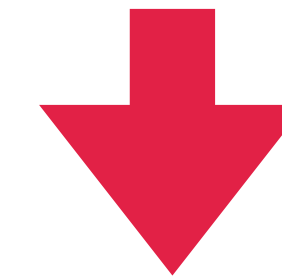
Masked Multi-Head Attention



문제 발생!!

$$\begin{matrix} & Q \\ & \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} < sos > \\ je \\ suis \\ \acute{e}tudiant \end{matrix} & \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \end{matrix} \times \begin{matrix} & K^T \\ & \begin{matrix} & & & & \end{matrix} \\ \begin{matrix} < sos > je suis \acute{e}tudiant \end{matrix} & \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \end{matrix} = \begin{matrix} & \begin{matrix} < sos > je suis \acute{e}tudiant \end{matrix} \\ \begin{matrix} < sos > \\ je \\ suis \\ \acute{e}tudiant \end{matrix} & \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \end{matrix}$$

Attention Score Matrix

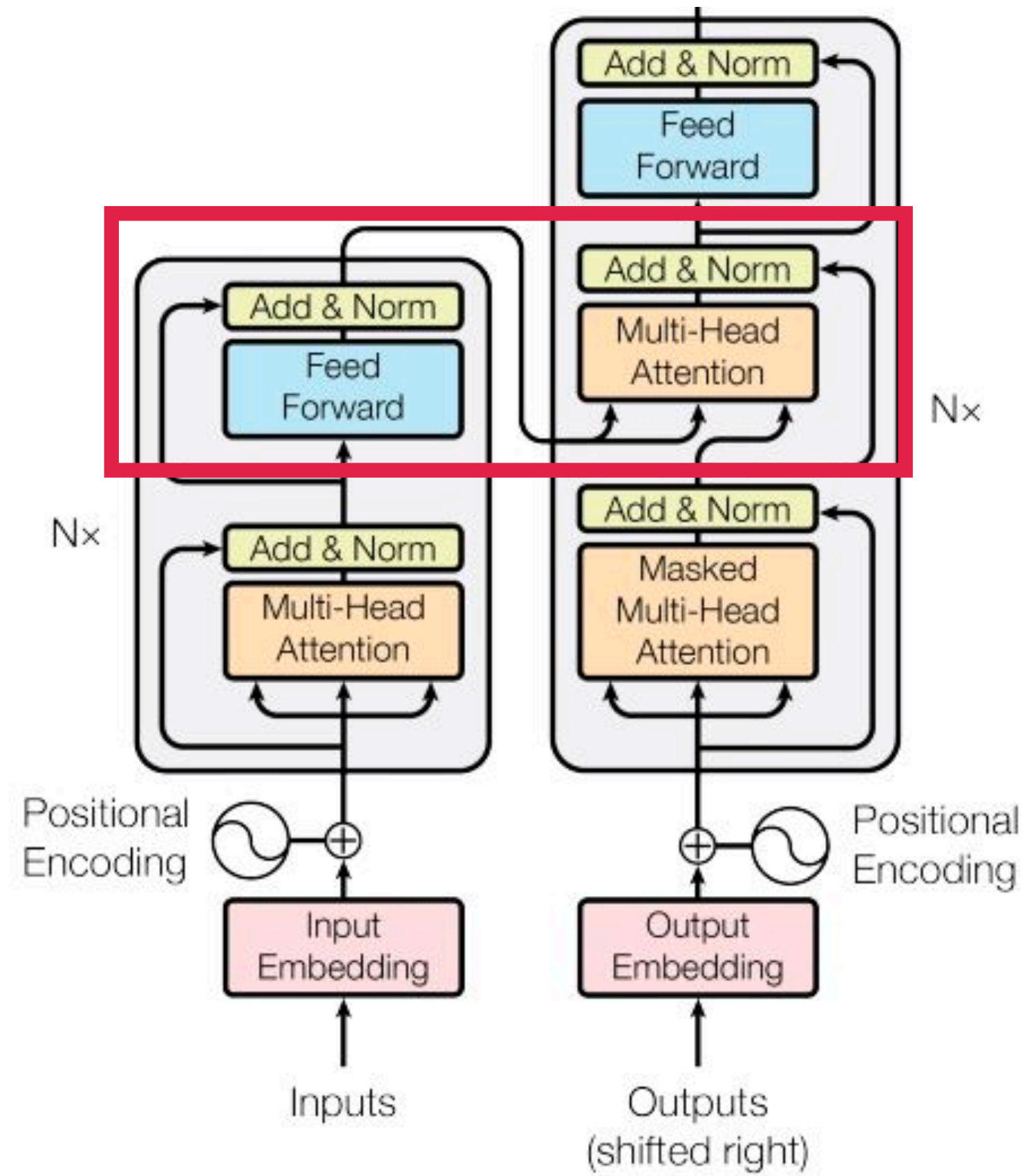


$$\begin{matrix} & \begin{matrix} < sos > je suis \acute{e}tudiant \end{matrix} \\ \begin{matrix} < sos > \\ je \\ suis \\ \acute{e}tudiant \end{matrix} & \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \end{matrix}$$

Attention Score Matrix

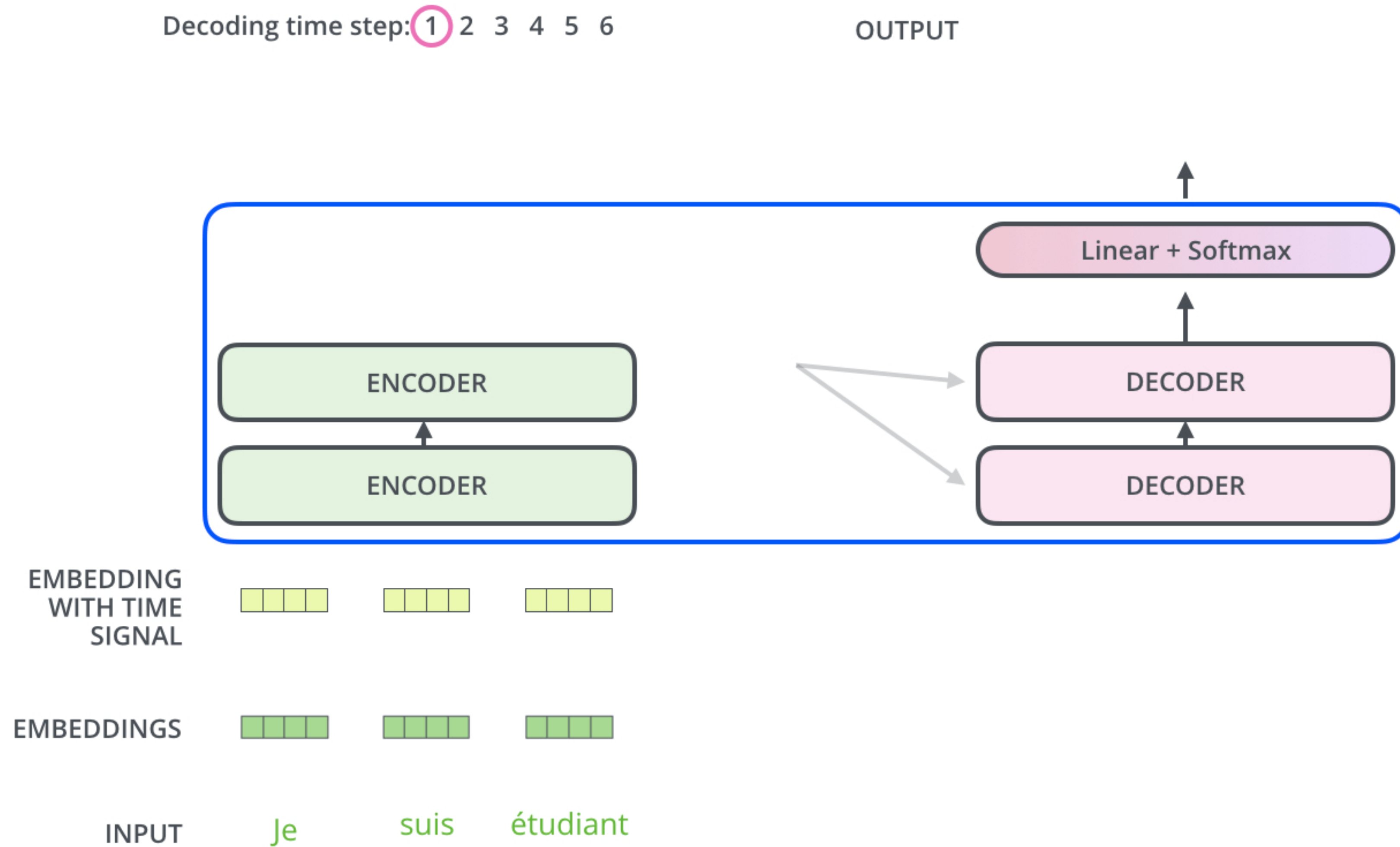
3. Transforemr

Encoder-Decoder Attention



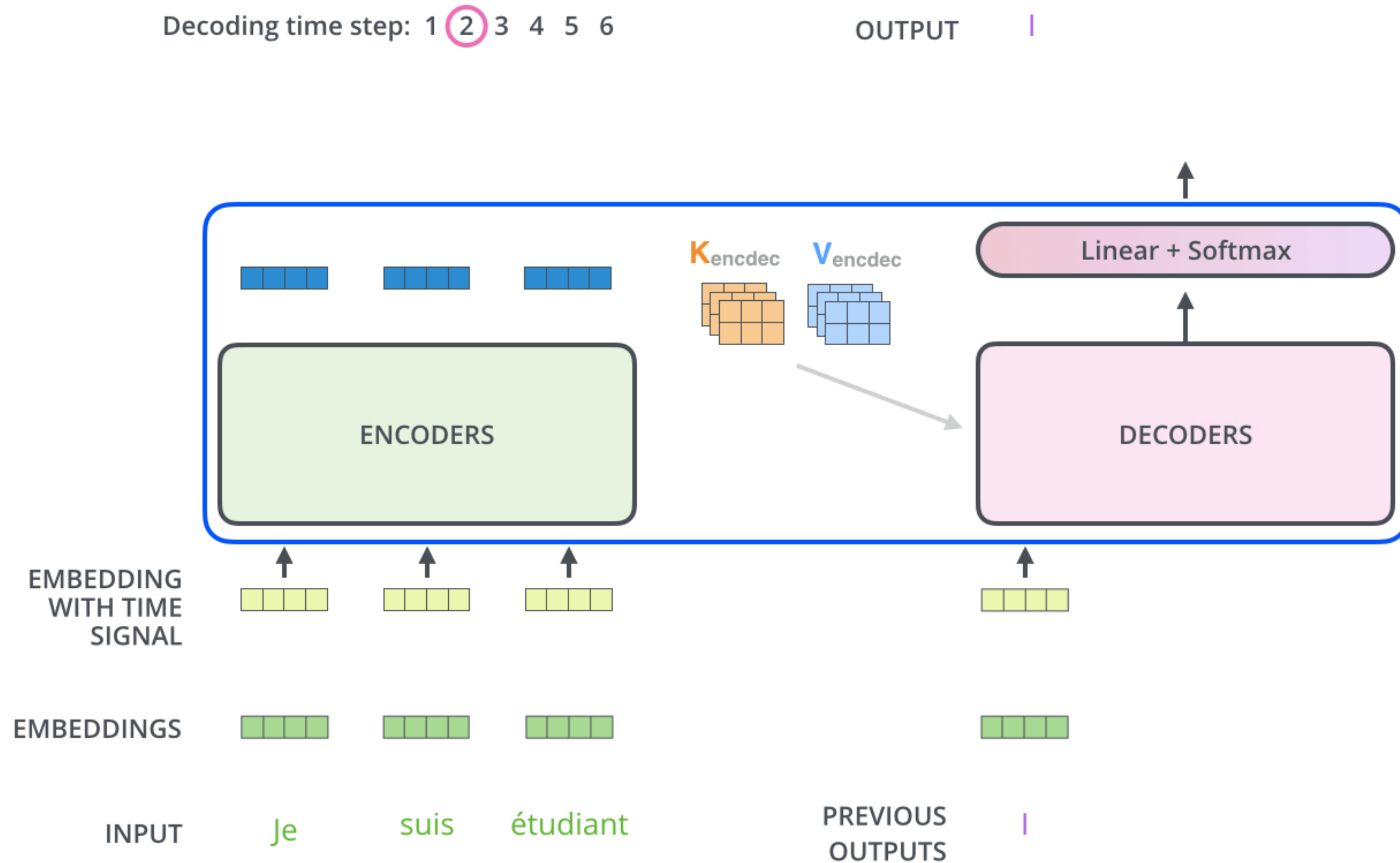
3. Transforemr

Encoder-Decoder Attention



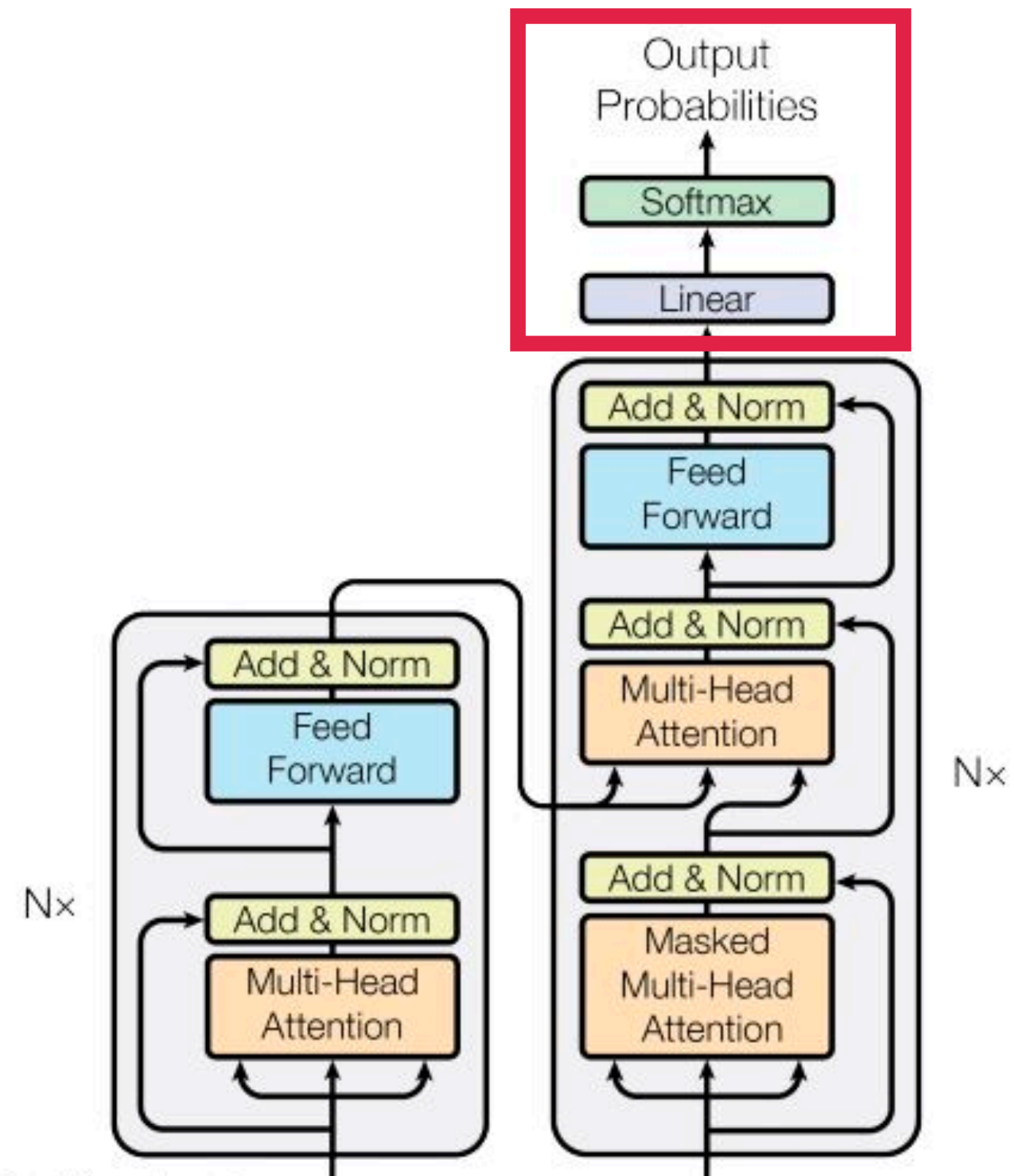
3. Transforemr

Encoder-Decoder Attention



3. Transforemr

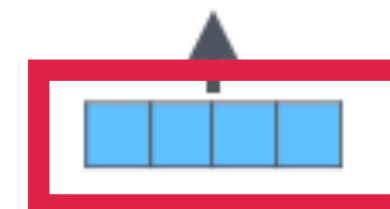
Linear Layer and Softmax Layer



3. Transforemr

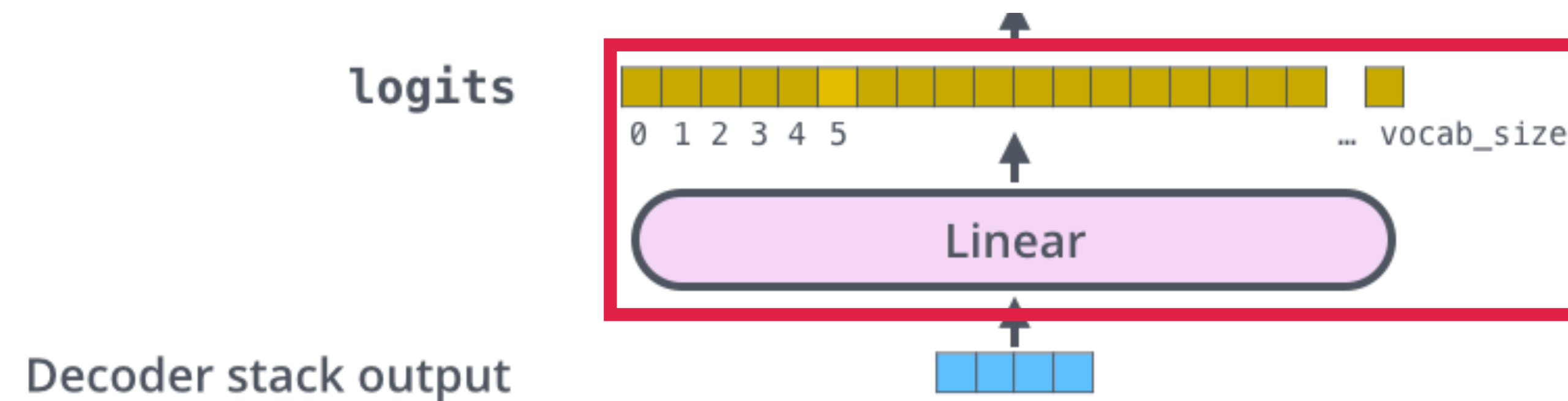
Linear Layer and Softmax Layer

Decoder stack output



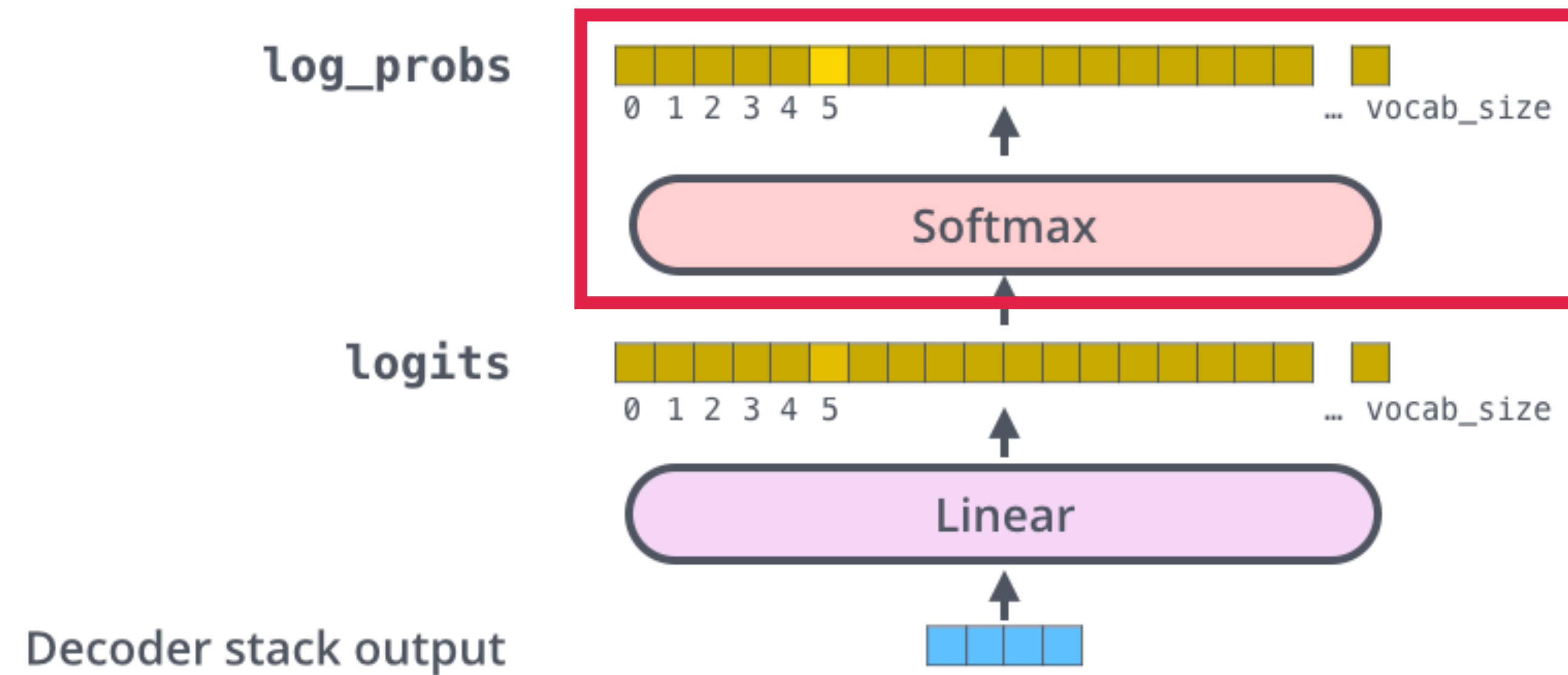
3. Transforemr

Linear Layer and Softmax Layer



3. Transforemr

Linear Layer and Softmax Layer

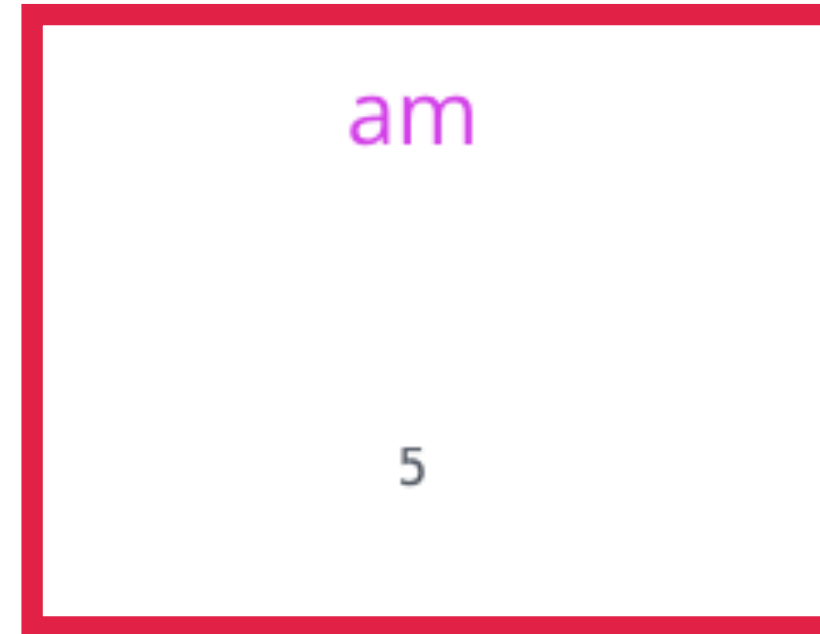


3. Transforemr

Linear Layer and Softmax Layer

Which word in our vocabulary
is associated with this index?

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



log_probs



Softmax

logits

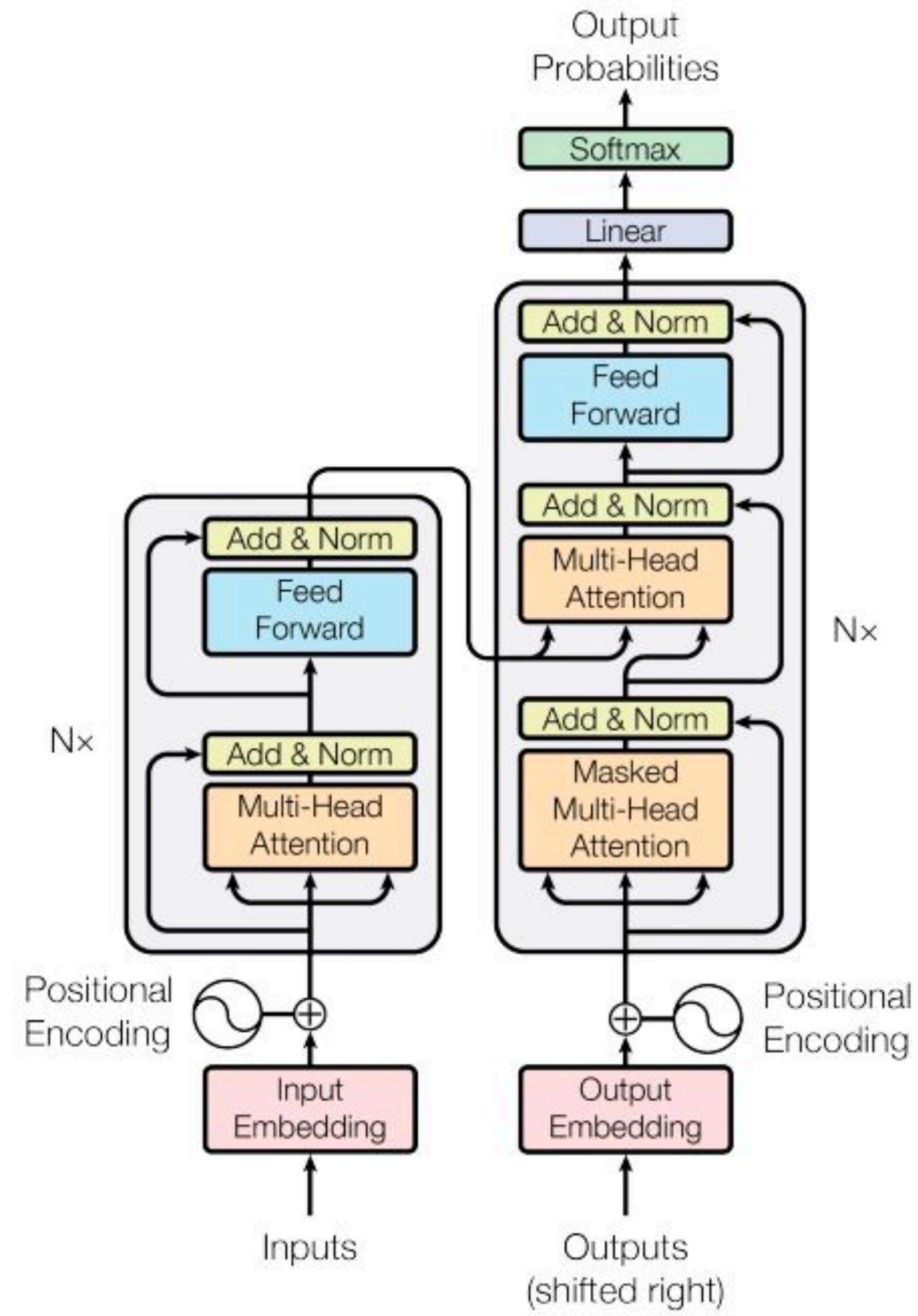


Linear

Decoder stack output



3. Transforemr



Q&A