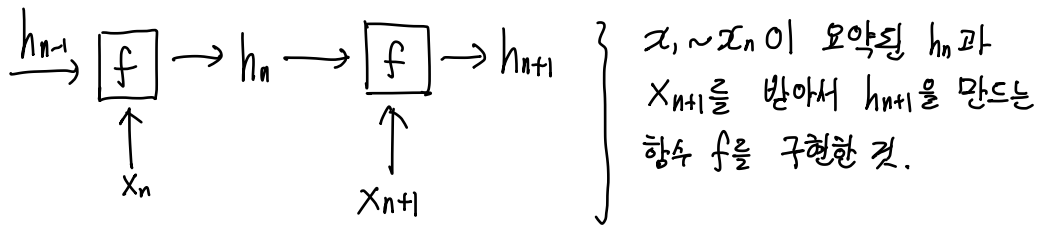


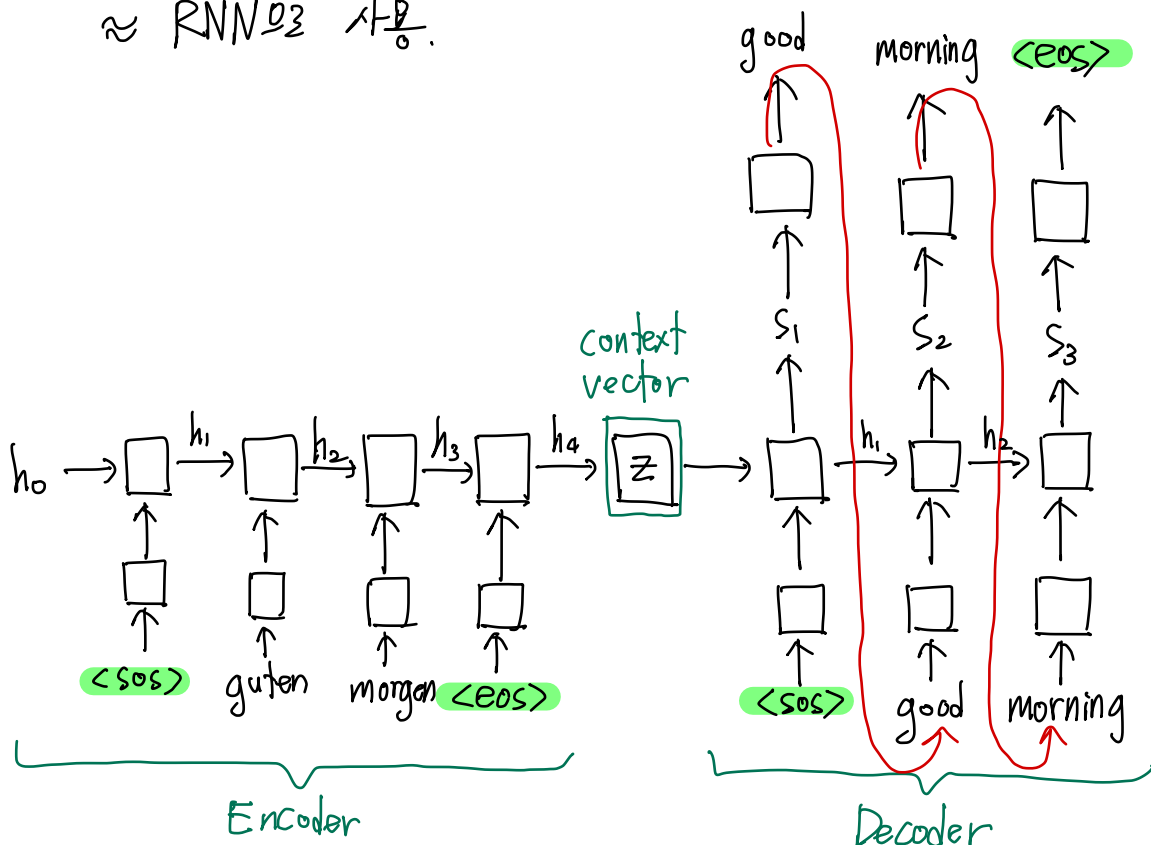
# RNN



$$f(x_n, h_{n-1})$$

## Seq2Seq

- 번역 machine.
- 입력/출력 문장의 단어 수 다름.
- RNN encoder, RNN decoder 사용.
- : Special token
- 단점
  - ① Hidden state가 멀리있는 내용을 기억해내기 못함. (Context vector)
  - ② 출력한 것을 다시 입력함.
- LSTM: 장기/단기 기억을 구분해서 관리하는 RNN.  
 $\approx$  RNN으로 사용.



## Embedding Layer

·  $\mathbb{R}^{\text{사전차원}} \rightarrow \mathbb{R}^{500}$

· 사전 차원으로 one-hot encoding 돼 있는 벡터를  $W$  곱해서 차원 축소.

ex)  $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{\text{사전}} \rightarrow W \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{500}$

Word2vec

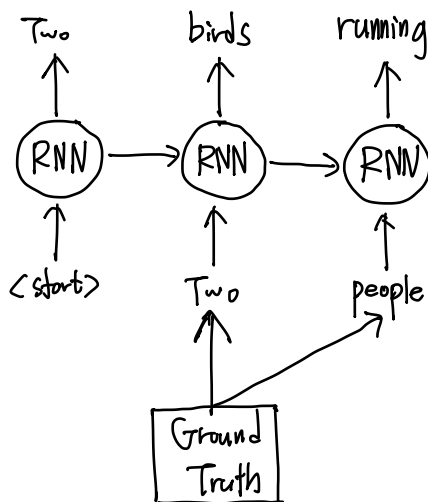
## Teacher Forcing

· Decoder 에서 사용하는 기법.

· 예측값 말고 정답 입력.

· 다음 단어 맞추기

$\Rightarrow$  Language Model



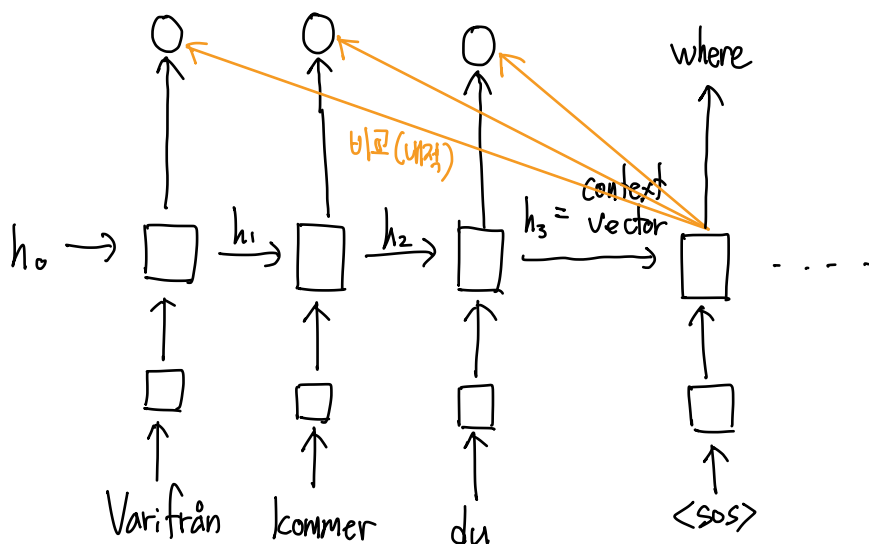
## Attention

· Decoder에 다음 입력되는건 context vector라는 그냥 1개의 벡터.

$\Rightarrow$  빈약해서 앞의 내용 잘 기억 못함.

· Decoder가 출력하기 전에 Encoder의 출력들인 hidden state들과 내적을 통해 비교

$\Rightarrow$  볼에너전에 상관관계를 비교해서 붙여넣다!



- $\{V_1, V_2, \dots, V_n\}$ 이  $\mathbb{R}^n$ 의 orthonormal basis라면  
 $\rightarrow$  수렴, 크기 1

$\rightarrow$  Query

$$X = C_1 V_1 + C_2 V_2 + \dots + C_n V_n$$

$$= (X \cdot V_1) V_1 + (X \cdot V_2) V_2 + \dots + (X \cdot V_n) V_n$$

$\downarrow$  key     $\downarrow$  value     $\downarrow$  key     $\downarrow$  value     $\downarrow$  key     $\downarrow$  value

변형 ①

$$\rightarrow ((W_Q X) \cdot (W_K V_1)) \underbrace{W_V V_1}_{\tilde{V}_1} + ((W_Q X) \cdot (W_K V_2)) \underbrace{W_V V_2}_{\tilde{V}_2} + \dots + ((W_Q X) \cdot (W_K V_n)) \underbrace{W_V V_n}_{\tilde{V}_n}$$

$\alpha_1$      $\alpha_2$      $\alpha_n$

$$\tilde{X} = \alpha_1 \tilde{V}_1 + \alpha_2 \tilde{V}_2 + \dots + \alpha_n \tilde{V}_n$$

변형 ②

$$\rightarrow \frac{e^{\alpha_1}}{\sum_{i=1}^n e^{\alpha_i}} \tilde{V}_1 + \frac{e^{\alpha_2}}{\sum_{i=1}^n e^{\alpha_i}} \tilde{V}_2 + \dots + \frac{e^{\alpha_n}}{\sum_{i=1}^n e^{\alpha_i}} \tilde{V}_n$$

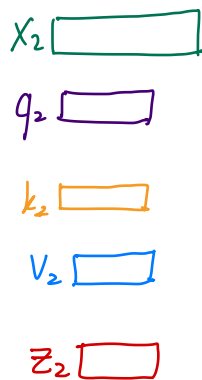
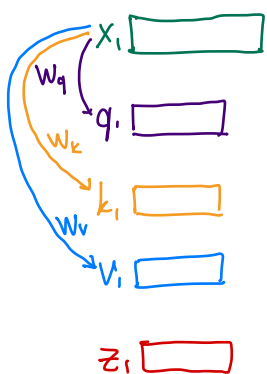
softmax ; skip-gram 처럼.

## Self Attention

- 위의 attention 식에서  $x$ 가  $V_1$  or  $V_2$  or  $\dots$  or  $V_n$ .
- $V_1, V_2, \dots, V_n$ 이  $\begin{cases} \text{Query} \\ \text{key} \\ \text{value} \end{cases}$  전부 다 됨.

Thinking

Machine



$$\Rightarrow z_1 = \text{softmax}(q_1 \cdot k_1) V_1 + \text{softmax}(q_1 \cdot k_2) V_2$$

(temperature 사용한 softmax)

## Multi-headed Attention

- $Q, K, V$  조합  $\rightarrow$   $n$  여러개 만들어서 다양한 관점 반영하는 것.
- Original transformer model은 head 8개.
- $W^O$ 은 각 head의 가중치 반영.
- parameter :  $\begin{cases} W^Q \\ W^K \\ W^V \\ W^O \end{cases}$

## Positional Encoding & Embedding

subword) word-piece, BPE

$$\begin{aligned} \text{unaffordable} &= \text{un} + \text{afford} + \text{able} \\ &= [19, 28, 42] \leftarrow \text{정수 index로 tokenizing.} \end{aligned}$$

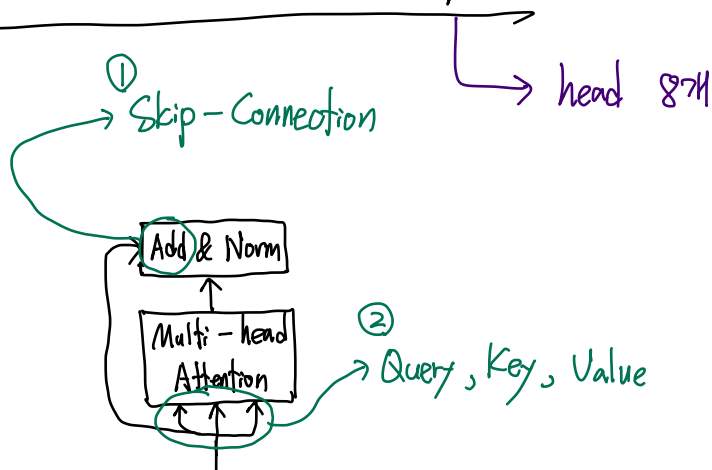
$\Rightarrow$  30000 차원으로 차원 축소 : pretrained

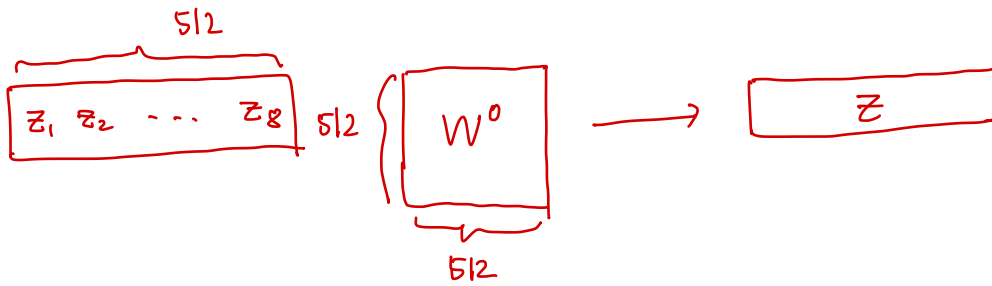
$$\begin{aligned} &\downarrow W (\text{Embedding Matrix}) : \text{parameter.} \\ W \begin{bmatrix} \end{bmatrix} &\in \mathbb{R}^{512} \text{ 차원 축소} \end{aligned}$$

Position Encoding)  $\sin / \cos$  함수 사용. (짝/홀)

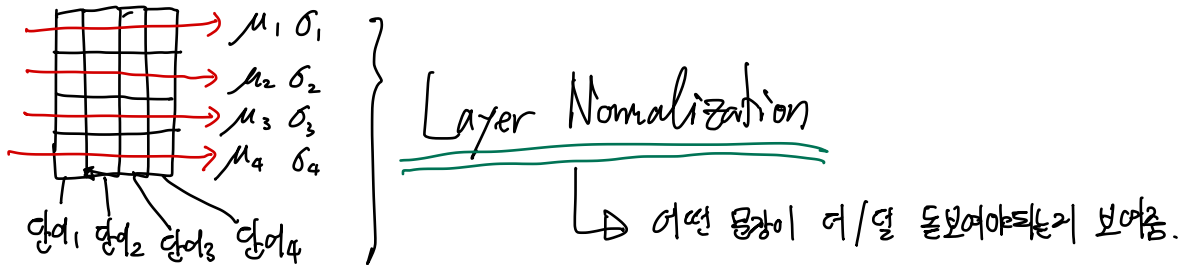
같은 위치  $\Rightarrow$  주기 같게

## Multi-head Self Attention layer



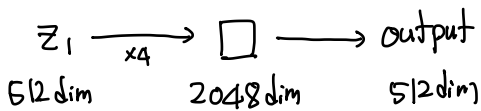


## [Residual Connection]



## Pointwise FFN

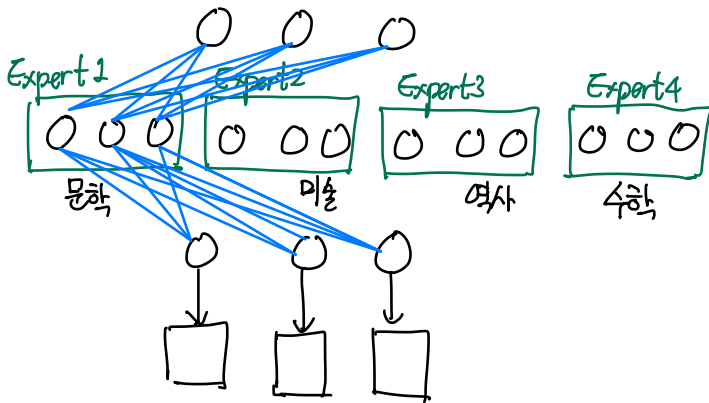
- $z_1, z_2$  마다 목적에 맞게 FFN에 넣고 돌림.



- 대부분의 parameter 잡아먹음.

## Mixture of Expert (MOE)

- 각각의 일에 맞게 Expert 중 선택함.



## Masked Self Attention

- Self-Attention (Encoder)에서는 전체 문장 알고있지만  
Decoder (번역)에서는 뒤에 어떤 단어 묻지 미리 알수 없다.

→ Masked Self-Attention.

자기자신까지 끝에서 계산한다.

→  $n+1$ 개 까지 알면 그걸 다 넣어서 반복.

\* Key, Value는 Query에 영향 준다.