

Seq2Seq에서 이전 모든 data를 특정 크기의 dimension을 갖는 vector로 embedding이 거기 붙는다.

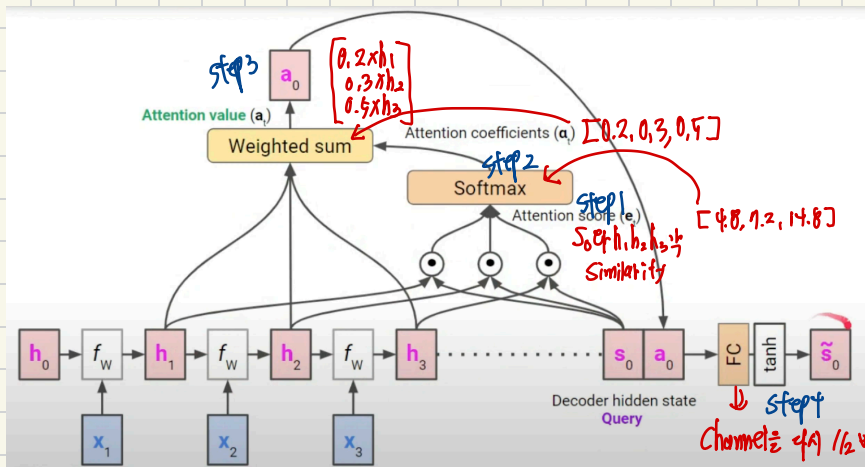
decoding이 encoder의 최종 vector만 사용하지 않고, encoder 중간 중간의 output을 참조하자.

< Attention Idea > 어디에 집중할지 결정한다.

$Q, K$ 를 사용하여 weight (가중치)를 생성하고 Value를 곱해주는 연산.

$Q$  : decoder의 time  $t$ 의 현재상태 (가중)  
 $K$  : encoder의 hidden state 들  
 $V$  : "

현재 state encoder에 집중?



위의 과정을 모든  $S_n$ 에 대해 진행한다.

$$h_1, h_2, \dots, h_n \in \mathbb{R}^h$$

$$s_t \in \mathbb{R}^h \quad \# \text{ 현재 state.}$$

$$e_t = [s_t^T h_1, s_t^T h_2, \dots, s_t^T h_T] \quad \# \text{ encoder-} \rightarrow \text{hidden state} \\ \text{들과 유사도 계산}$$

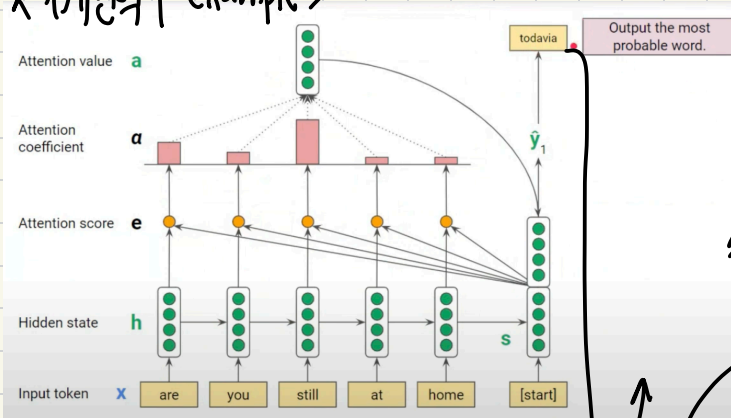
$$\alpha_t = \text{softmax}(e_t) \in \mathbb{R}^T \quad \# \text{ softmax}$$

$$a_t = \sum_{i=1}^T [\alpha_t]_i \cdot h_i \in \mathbb{R}^h \quad \# \text{ value를 곱한다}$$

$$[a_t, s_t] \in \mathbb{R}^{2h} \quad \# \text{ state에 append 한다.}$$

# Similarity는 cosine sim이 간단하면서도 성능이 좋아서 많이 사용.

<가계보예시 example>



생각할 때는 max length 설정.

while result = <eos> 인 경우

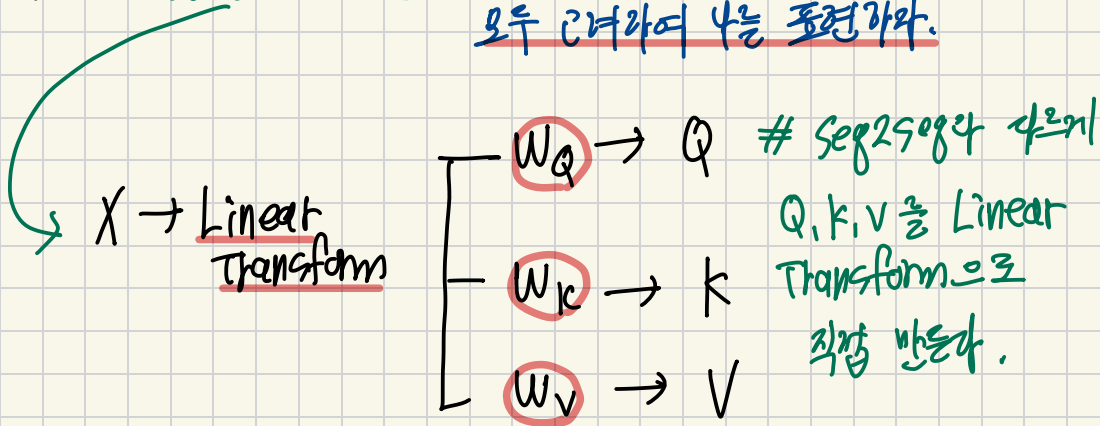
피드백 (teacher force)

# < Transformers >

NN방 CNN방 원칙 동일하다.  $\Rightarrow \hat{y} = W_2 \sigma(W_1 x)$  결국  $x$ 들의 weighted sum을 사용한 모델

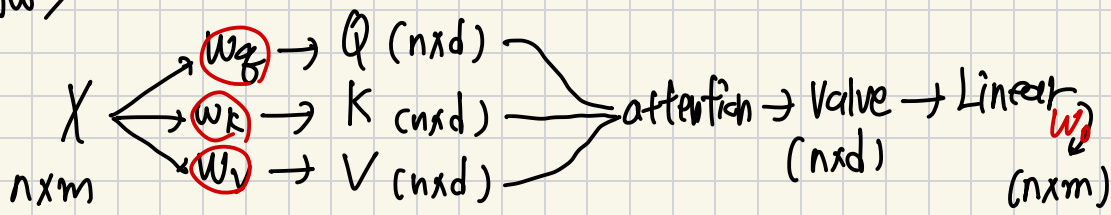
RNN  $\Rightarrow$  애드 결국  $\tanh(W_2 \cdot \tanh(W_1 \cdot x_i))$ 처럼  $x$ 들의 weighted sum.

$x$ 를 data를 두어 각  $x$ 를 서로 상호작용하여  $x$ 를 표현  
 함.  $\Rightarrow$  "Self attention":  $x$ 를 표현시 나를 알아 주변 data들도 모두 고려하여  $x$ 를 표현하라.

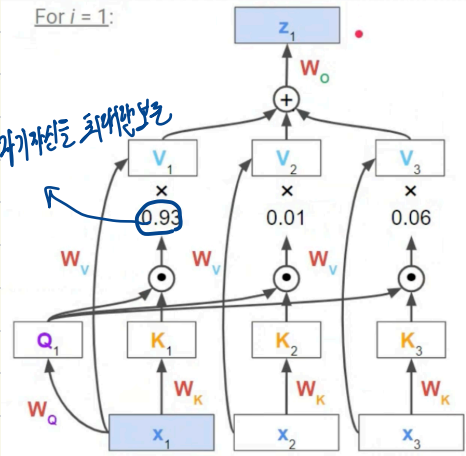


$\times W_Q, W_K, W_V$ 은 모든 input에 same 가다  
 learnable parameter라서 어떻게 Q, K, V 역할 잘 수행하도록 학습한다.

## < Flow >



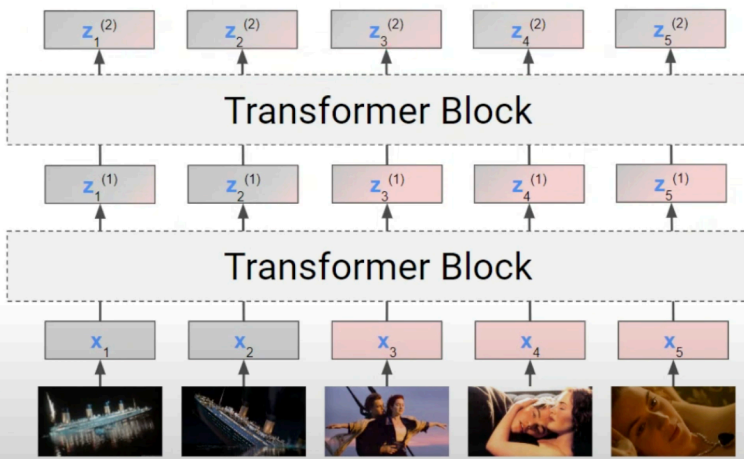
For  $i = 1$ :



↳ " $x_1$ 을  $x_2, x_3$ 과의 상호작용으로 새롭게 표현하여  $z_1$  이 된다"

$z_1$ 은  $x_1$ 의 상호작용된 새로운 자성이 된다.

즉, 모든  $x$ 들은 주변 data들과의 상호작용으로 재표현 (Transform) 한다. ★  
"contextualize"



즉, <sup>self</sup> Attention은  
자기자신을 문맥에 맞게  
재표현하는 방식이다.