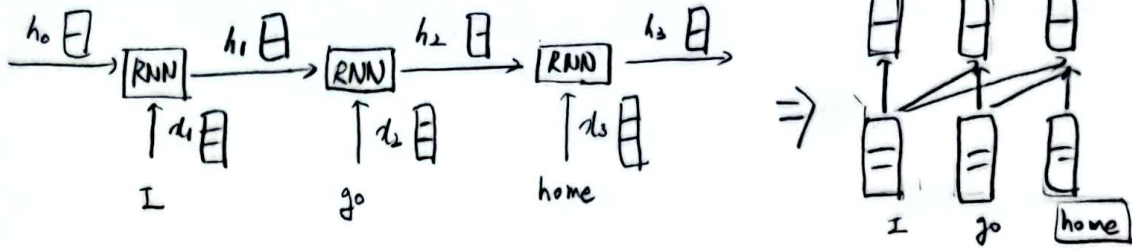
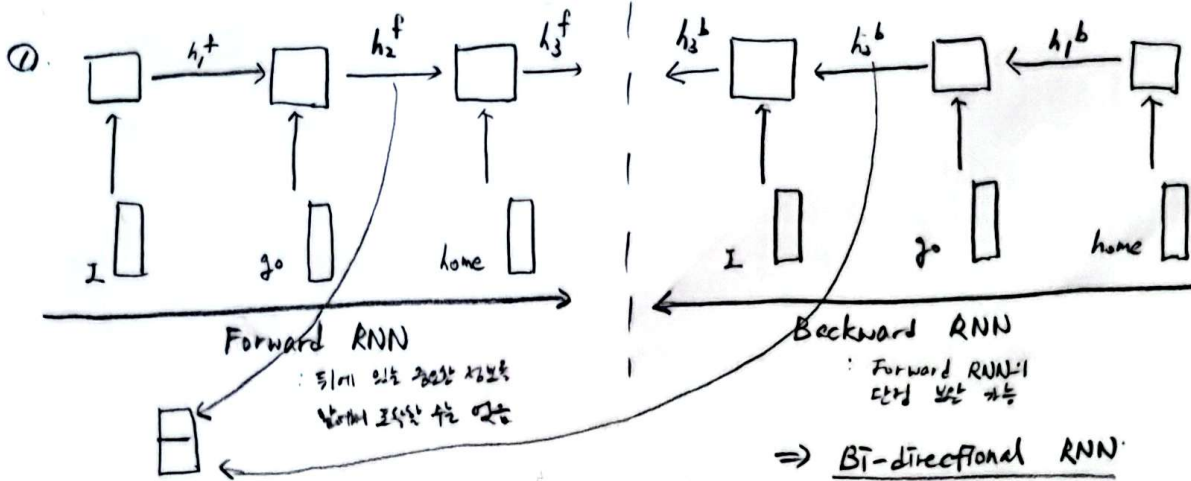


# RNN: Long-term Dependency.



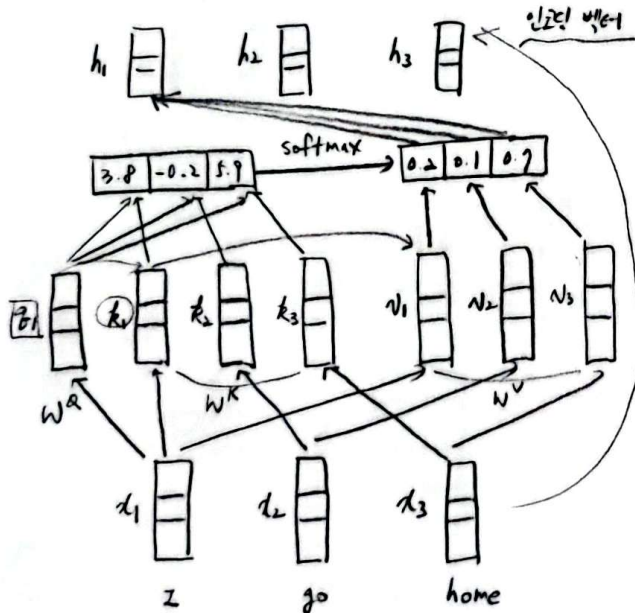
문장이 길어지면 멀리 있는 step까지 전달되는 과정에서 정보 유실 발생 가능성 있음

해결책)



hidden state를 RNN마다 각각 가져다와 concat을 해서 인코딩 벡터를 만든다.

## ② Transformer → sequence 길이나 길이에 time step이 지니는 정보 손실이 없게 가질 수 있는



Input	단어 1	단어 2
Embedding	$x_1$	$x_2$
Queries	$q_1$	$q_2$
Keys	$k_1$	$k_2$
Values	$v_1$	$v_2$

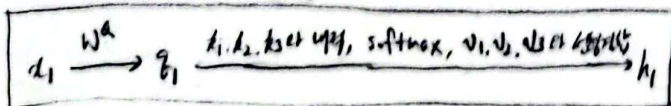
↑  
각각의 query에 softmax해 나온 가중치를 곱해서  
서열 정보에 사용 되는 정보 벡터.  
(Key 벡터와 유사한 이점)

Decoder hidden state vector 얻는  
주요한 vector 세로의 유사도를 구하는 방법 사용.

기준 벡터 → Query

(Self-attention으로 Encoder hidden state vector로  
유사도 구함)

ex)



$$\begin{matrix} x_1 \\ x_2 \end{matrix} \begin{matrix} \times \\ \times \end{matrix} \begin{matrix} W^Q \\ W^K \end{matrix} = \begin{matrix} q_1 \\ q_2 \end{matrix}$$

$$\begin{matrix} W^V \end{matrix} \begin{matrix} v_1 \\ v_2 \end{matrix} = \begin{matrix} v_1 \\ v_2 \end{matrix}$$

# Scaled Dot-Product Attention

Inputs:  $q$ : query,  $(k, v)$ : set of key-value pairs ( $q, k, v$ 는 모두 vector)

Output은 Values의 가중평균

각 Value의 가중치는 query와 각 value에 대응하는 key의 내적으로 계산

Query와 Key는 같은 차원  $d_k$ . Value는  $d_v$  (이때  $d_k$ 과  $d_v$ 가 같아도 되는 것 아님)

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot k_i)}{\sum_j \exp(q \cdot k_j)} v_i$$

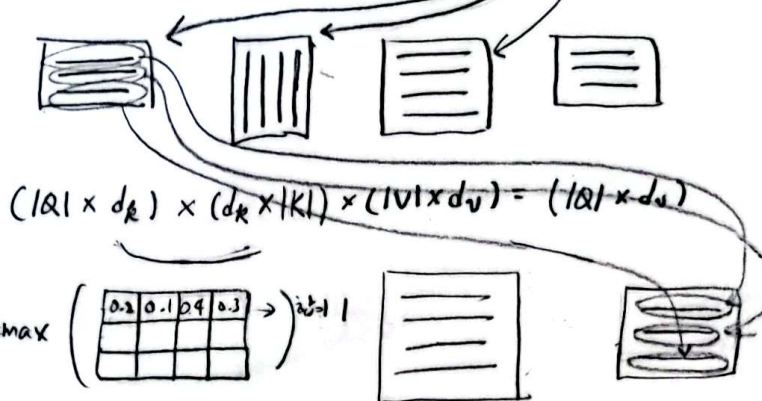
attention module  
각 key와 query vector의 내적 계산  
softmax

$k_1 \rightarrow 0.2$   
 $k_2 \rightarrow 0.8$

만약 query가 여러 개이면,

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

$$K = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$



실제 Transformer 구현 상으로는 동일한 shape로 mapping된  $Q, K, V$ 가 사용되어 각 matrix의 shape는 모두 동일

ex)  $X \times W^Q = Q$

$X \times W^K = K$

$X \times W^V = V$

$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

$$= Z$$

문제점)

$d_k$ 가 커지면  $QK^T$ 의 분산이 증가  
일부 값이 커질 수 있으며, softmax 값이  
한쪽으로 쏠릴 수 있고, 학습 시에 gradient  
값이 작아져 가짜 소문 가능성이 커진다

$$\begin{bmatrix} a & b \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = a + 2b$$

$q$   $k$

$a, b$ 의 값을 0, 분산: 1인 확률분포를 할 때  
(4점 득점)

$a + 2b$ 의 분산: 0,  $a$ 의 분산: 1,  $b$ 의 분산: 1  
 $\rightarrow a + 2b$ 의 분산: 2

내적 값

$$\begin{bmatrix} 1 & 1 & 0.8 & 0.7 \end{bmatrix}$$

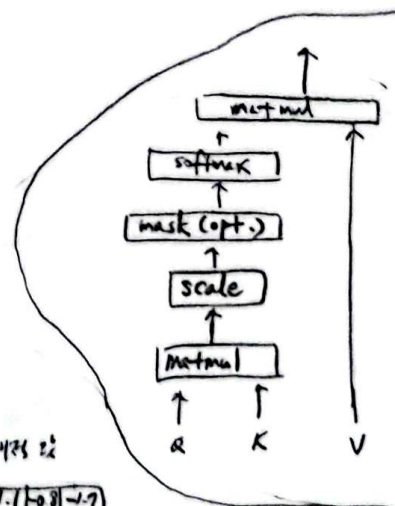
해결책)

만약 각각 10차원이라면 분산: 0, 분산: 100

$\rightarrow$  표준편차: 10  $\rightarrow \begin{bmatrix} 0 & -11 & 7 \end{bmatrix}$

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

분산을 일정하게 유지





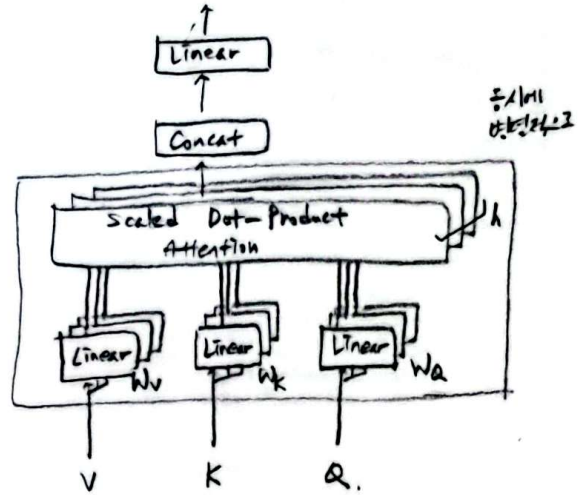
# Multi-Head Attention

동일한 sequence가 주어져도 특정한 query word에 대해 서로 다른 기준으로 여러 측면에서의 정보를 추출해야 하는 장점이 있다.

ex) 여러 문장으로 이루어져 있지만, 특정한 query word의 sequence를 볼 수 있다.

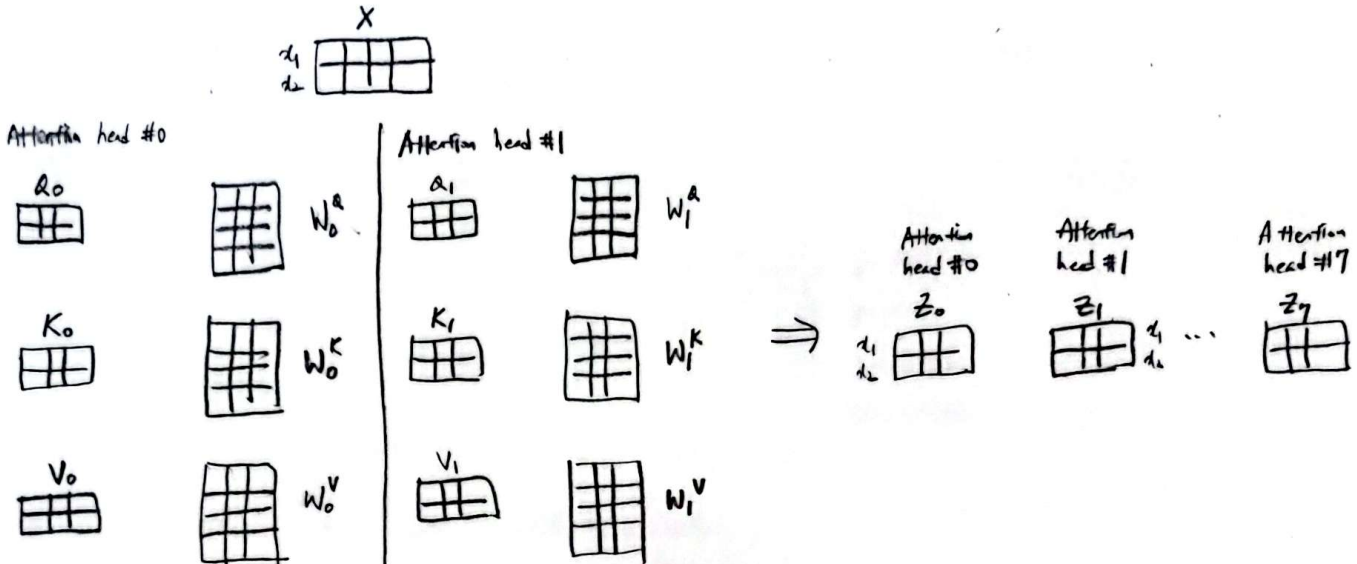
- I went to school
- I studied hard.
- I came back home.
- I took the rest.

여기서 'I' query word에 대한 정보는 어떻게 가져오는가?  
 ① 'I'가 한 문장을 기준으로 ② 'I'가 있는 단어의 번호를 뽑을 수 있다.

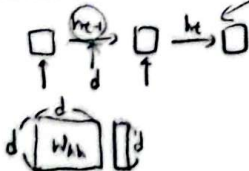


$$\text{Multihead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

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



- $n$ : sequence 길이
- $d$ : vector의 차원 (query, key, val dim)
- $k$ : convolutional kernel size
- $r$ : restricted self-attention의 neighborhood의 크기



$$\text{Self-attention} = \frac{(QK^T)}{\sqrt{d_k}} V$$

	Complexity per layer	Sequential operations	Maximum path length
Self-attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$

long-term dependencies를 갖는 경우  
 long-term dependencies 때문.

$n$ : 계산할 수 없는 메모리  
 $d$ : 정보의 양

Self-attention은 Recurrent에 비해 계산량이 많고 메모리가 많이 필요하다.  
 Self-attention은 GPU에서 실행되는 반면 Recurrent는 CPU에서 실행된다.  
 Recurrent는 두 가지 방향으로 진행되어야 한다 (백워드와 포워드)

# Block-Based Model

## Residual connection

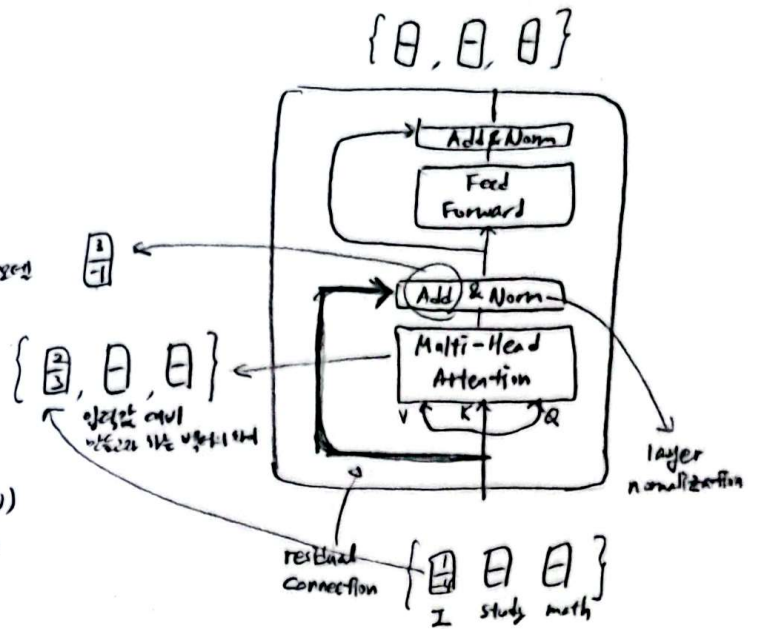
- CV 분야에 있어 깊은 layer의 NN(Neural Net)을 만들 때, Gradient Vanishing 문제를 해결하기 위한 방법 중 하나로 layer를 쌓아올릴 때 이 층은 선형층 보다는 비선형층이 더 중요

입력 벡터와 Embedding의 output에 선형층은 출력 벡터의 차이를 동일하게 만들 수 있음

각 Block은 2개의 sub-layer (Multi-head attention, 2-layer feed-forward NN (with ReLU))

각각은 추가로 2개의 step (Residual connection, Layer normalization, Layer Norm (1 + sublayer(x)))

을 거친다.



## Layer Normalization

→ 각 layer에서 input의 평균: 0, 분산: 1이 되도록 한다

$$\mu^L = \frac{1}{H} \sum_{i=1}^H x_i^L, \quad \sigma^L = \sqrt{\frac{1}{H} \sum_{i=1}^H (x_i^L - \mu^L)^2}, \quad h_i = f\left(\frac{g_i}{\sigma_i}(x_i - \mu_i) + b_i\right)$$

$$\mu = 2.5, \quad \sigma = 1.11, \quad \mu = 0.5, \quad \sigma = 2.06$$

4	2
2	1
3	-3
5	2

thinking machines



$$\mu = 0, \quad \sigma = 1, \quad \mu = 0, \quad \sigma = 1$$

0.65	0.7
-0.85	-1.5
-1.35	-0.3
1.25	1.1

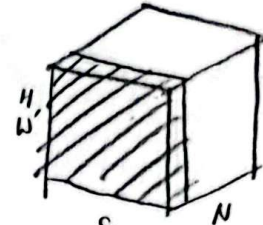
thinking machines



2.95	3.1
0.55	-0.5
1.35	0.3
-0.5	-1.2

thinking machines

$$\begin{aligned} y &= 3x + 1 \\ y &= x + 1 \\ y &= -1 + 0 \\ y &= -2x + 2 \end{aligned}$$



① 각 word vector를 평균 0, 분산 1이 되도록 normalization

② 각 sequence vector에 대해 가능한 parameter에 affine transformation 적용 (평균, 분산을 주입해 줌)

$$y = ax + b$$

## Positional Encoding

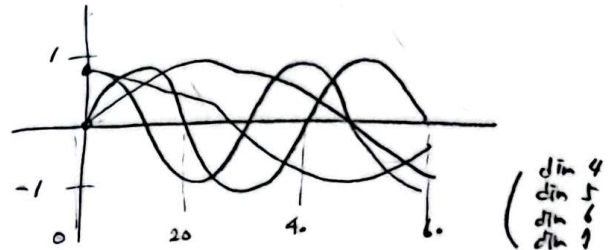
RNN과 달리 self-attention은 교환법칙이 성립하므로 순서를 무관하다  
즉, sequence는 2차원에서 인코딩해야 하는 값이 정해진 순서로 인코딩하게 된다  
이를 보완하기 위해 positional encoding이 나왔다

→ 서로 다른 순서간에 구별이 되는 특정한 sin/cos 값을 임의 벡터에 더한다  
sin, cos 등의 주기함수를 사용함, 주기를 다르게 해서 여러 값을 만든다

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

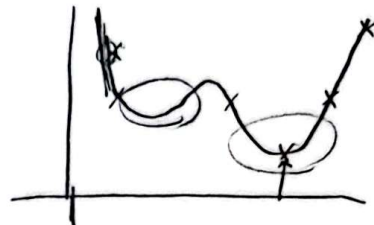
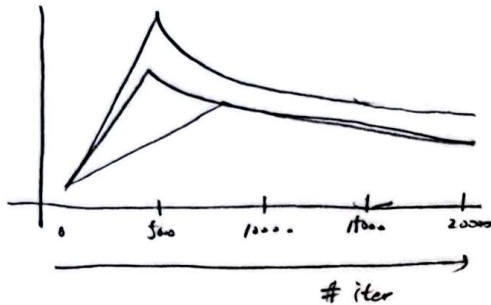
$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i/d_{model}})$$

→  $d_{model}$ 에 따라 주기가 달라진다



## Learning Rate Scheduler

learning rate를 학습 중에 적절히 변경할 수 있다



$$\text{learning rate} = d_{model}^{-0.5} \cdot \min(\#step^{-0.5}, \#step \cdot \text{warmup-steps}^{-1.5})$$

초반에는 너무 많은 보폭이 반영되어 강도를  $\text{lr}$ 을 작게 한다

어느 정도 step이 진행되어 안정해지면 실제 도달해야 하는 과정에서 아직 멀 수 있어서

동향을 주기 위해  $\text{lr}$ 을 iteration 수에 비례해서 증가시켜 준다

그래서 학습률 근처에 도달하면 (보이 크면 수렴해 있고 멈출 수 있으므로)  $\text{lr}$ 을 점차 줄여준다

