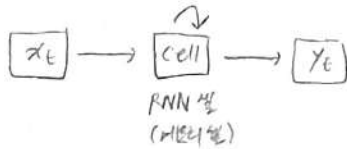


* tokenization : 주어진 코퍼스 (corpus) 에서 토큰 (token) 이라 불리는 단위로 나누는 작업

* 순환 신경망 (Recurrent Neural Network, RNN)



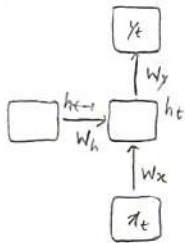
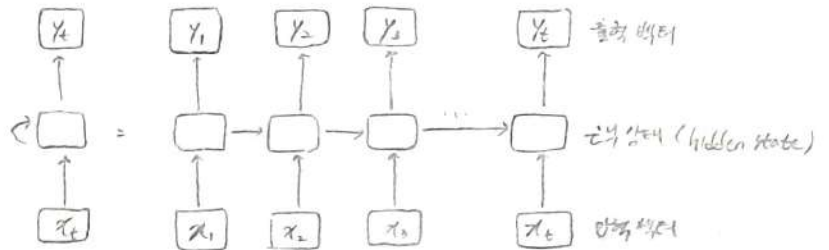
: 은닉층이 재귀형태

→ 입력과 출력의 길이 다르게 설계 가능

one-to-many ex. Image Captioning

many-to-one ex. sentiment classification, spam detection

many-to-many ex. 번역기, 기계적 영상 인식, 풀사 태깅



$$\text{은닉층} : h_t = \tanh(W_x x_t + W_h h_{t-1} + b)$$

$$\text{출력층} : y_t = f(W_y h_t + b) \quad (f: \text{비선형 활성화 함수, 시그모이드, 소프트맥스 등})$$

은닉 벡터의 차원이 d , 은닉 상태 크기 D_h 일때

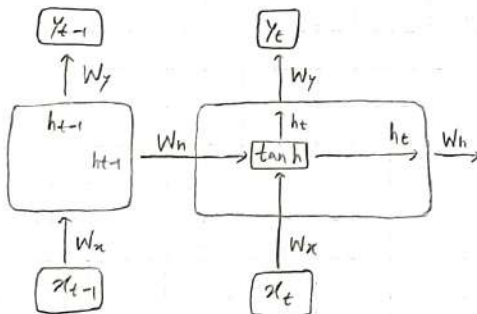
$$x_t: (d \times 1), W_x: (d \times D_h), W_h: (D_h \times D_h), h_{t-1}: (D_h \times 1), b: (D_h \times 1)$$

+ Deep Recurrent Neural Network, Bidirectional Recurrent Neural Network

* 장단기 메모리 (Long Short-Term Memory, LSTM)

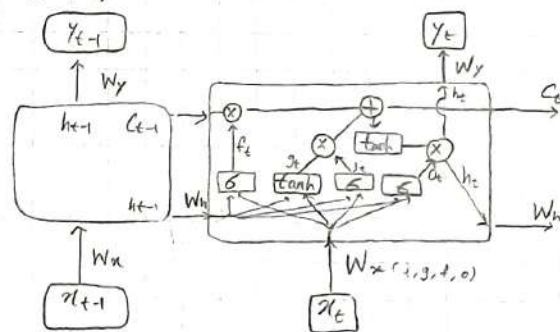
비선형 RNN의 한계 - 장기 의존성 문제 (the problem of Long-Term Dependencies) : 시간 (time step) 이 지날수록 x_t 정보량 손실

<RNN>



$$h_t = \tanh(W_x x_t + W_h h_{t-1} + b)$$

<LSTM>



C_t : 셀 상태 (cell state)

입력게이트 : 현재 정보를 저장하기 위한 게이트 $f_t = \sigma(W_{f1} x_t + W_{f2} h_{t-1} + b_f)$ / $i_t = \tanh(W_{i1} x_t + W_{i2} h_{t-1} + b_i)$

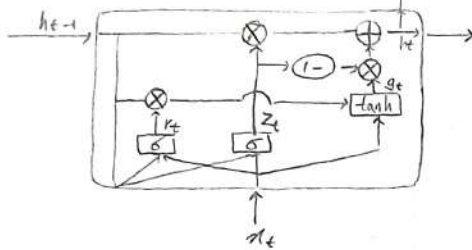
삭제게이트 : 지울 삭제하기 위한 게이트 $f_t = \sigma(W_{xf} x_t + W_{hf} h_{t-1} + b_f)$ 0과 1 사이 값. 0에 가까울수록 많이 삭제된 것

셀 상태 : 입력 게이트에서 선택된 정보를 삭제 게이트의 결과에 더한 것. $C_t = f_t \odot C_{t-1} + i_t \odot g_t$

삭제 게이트는 이전 시점의 입력을 얼마나 반영할 지 의미하고, 입력 게이트는 현재 시점의 입력을 얼마나 반영할 지 결정한다

출력게이트 : 현재 시점 t 의 은닉 상태를 출력하는 레이어. $o_t = \sigma(W_{xo} x_t + W_{ho} h_{t-1} + b_o)$ / $h_t = o_t \odot \tanh(C_t)$

* 게이트 순환 유닛 (Gated Recurrent Unit, GRU)



LSTM 구조 단원화 (매개변수 ↓)

업데이트 게이트 / 리셋게이트

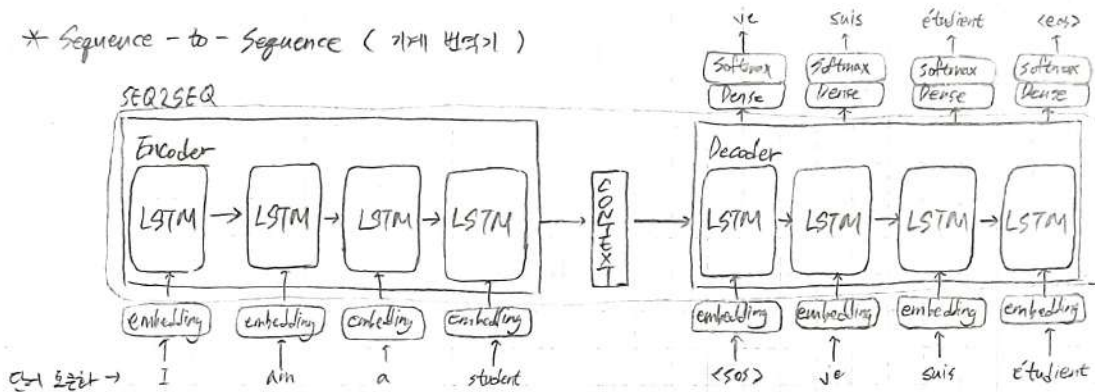
$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$g_t = \tanh(W_{hy}(r_t \odot h_{t-1}) + W_{xy}x_t + b_y)$$

$$h_t = (1 - z_t) \odot g_t + z_t \odot h_{t-1}$$

* Sequence-to-Sequence (기계 번역기)



인코더: 입력 문장의 정보가 하나의 컨텍스트 벡터로 압축, 디코더에 전송 / 디코더: 컨텍스트 벡터를 받아 번역된 단어를 한 개씩 순차적으로 출력

훈련 과정: Teacher forcing (교사 강요) - RNN의 모든 시점에서 이전 시점의 예측값 대신 실제값을 입력으로 주는 방법

테스트 과정: 인코더 RNN 실행 마지막 시점의 은닉 상태 (컨텍스트 벡터)를 디코더 RNN으로 넘긴다

→ 디코더는 문장의 시작을 나타내는 토큰 '<eos>'가 입력되면 다음에 등장하는 단어를 예측함, '<eos>'를 예측할 때까지 반복한다

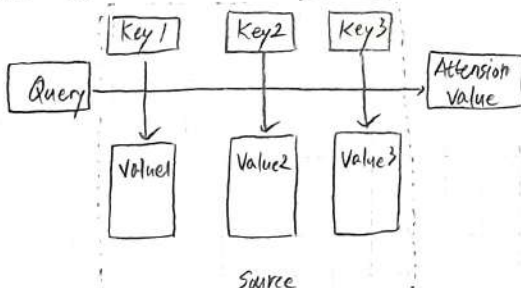
→ 디코더에서 각 시점의 RNN 셀에서 출력 벡터가 나오면 소프트맥스 함수를 통해 출력 시점의 각 단어별 확률값을 반환하고, 디코더는 출력 단어를 결정한다.

* Attention Mechanism

seq2seq 모델의 한계 - 1) 컨텍스트 벡터 정보 손실 2) 가변기 조절 문제

→ attention: 디코더에서 출력 단어를 예측하는 매 시점에서, 전체 입력 문장을 참고하여, 연관 부분을 더 집중해서 본다.

- Attention 함수



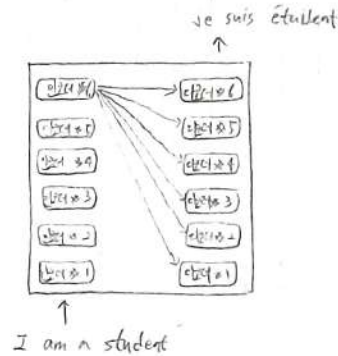
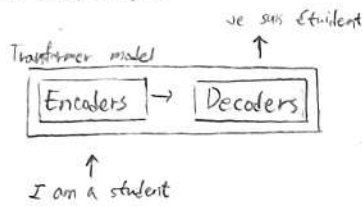
$$\text{Attention}(Q, k, V) = \text{Attention value}$$

Q = Query : + 시점의 디코더 셀에서 문장 상태

K = Keys : 모든 시점의 인코더 셀의 은닉 상태들

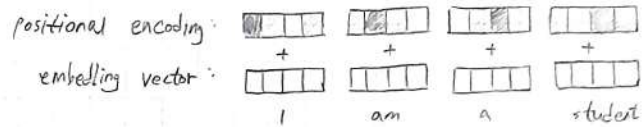
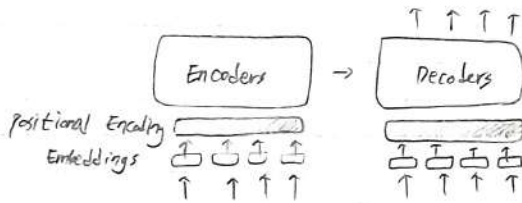
V = Value : 모든 시점의 인코더 셀의 은닉 상태들

* 트랜스포머 (Transformer)

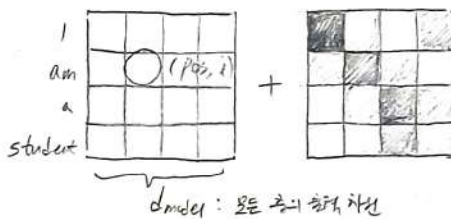


- 포지셔널 인코딩 (Positional Encoding)

포지셔널 인코딩: 트랜스포머의 입력에서, 단어의 위치 정보를 주기 위해 각 단어의 임베딩 벡터에 위치 정보를 더하여 모델의 입력으로 사용.



임베딩 벡터가 모여 만들어진 문장 행렬과 포지셔널 인코딩 행렬의 덧셈 연산을 통해 이루어진다



pos: 입력 문장에서 임베딩 벡터의 위치 / 1. 임베딩 벡터 내의 차원 인덱스

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

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

→ 주기적인 값을 생성하면서도 서로 다른 위치에 대해 고유한 값을 생성

pos에 대한 스케일링 역할

층수/축수 차이를 완화하여 차원 간 다양성을 확보해, 모델이 단어의 상대적 위치를 효과적으로 학습

- Attention

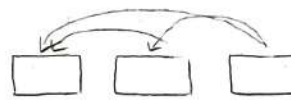
Encoder — Encoder self-Attention



(비트워의 출처)

$$Q = K = V$$

Decoder — Masked Decoder self-Attention



$$Q = K = V$$

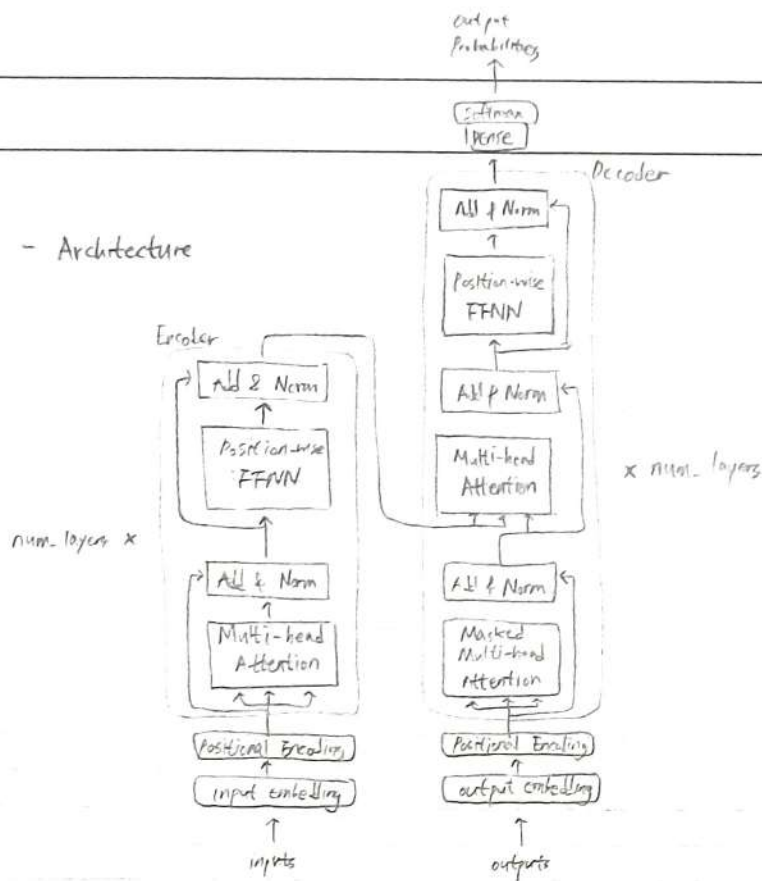
Encoder - Decoder Attention



$$Q = \text{decoder 벡터}$$

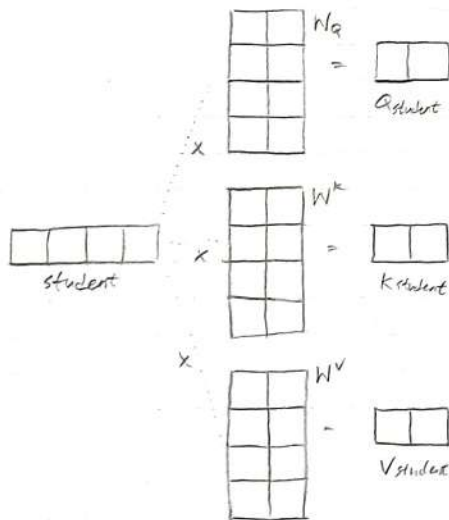
$$K = V = \text{인코더 벡터}$$

- Architecture



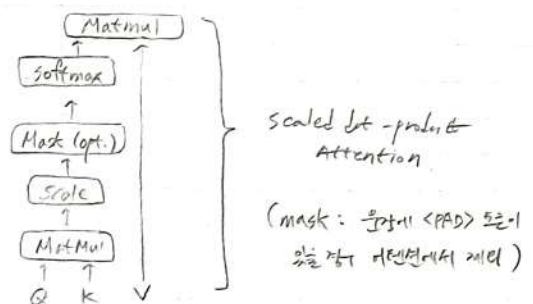
- Self Attention

$Q=K=V$: 입력 문장의 모든 단어 벡터를



Q, K, V 벡터들은 d_{model} 의 차원을 가지는 단어 벡터들보다 더 작은 차원을 가진다. 이때 파라미터 $num-head$ 로 결정된다

가중치 행렬은 $d_{model} \times (d_{model} / num-head)$ 의 크기를 가진다 모든 단어 벡터 각각에 Q, K, V 벡터를 곱한다



Scaled dot-product Attention

(mask: 문장에 <PAD> 토큰이 있을 경우 디펜던스에서 제외)

- Scaled dot-product Attention : dot-product 계산에서 스케일링 라인을 추가한 것

$$\text{Score}(q, k) = \frac{q \cdot k}{\sqrt{d_k}}, \quad (\text{논문에서 } n = d_{model} / num-head = d_k = d_v)$$

- 행렬 곱셈으로 원본 처리 - 문장 길이에 가중치 행렬을 곱하고, 각 단어에 적용된 값은 행렬 곱셈을 구한다

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

문장 길이는 $(seq-len, d_{model})$ 이고 d_k 는 Q 와 K 벡터의 차원, d_v 는 V 벡터의 차원이며 Q 와 K 행렬 크기는 $(seq-len, d_k)$, V 행렬의 크기는 $(seq-len, d_v)$

→ 가중치 행렬 W_Q, W_K 의 크기는 (d_{model}, d_k) , W_V 의 크기는 (d_{model}, d_v)

어떤 행렬 곱셈 d_v 의 크기는 $(seq-len, d_v)$

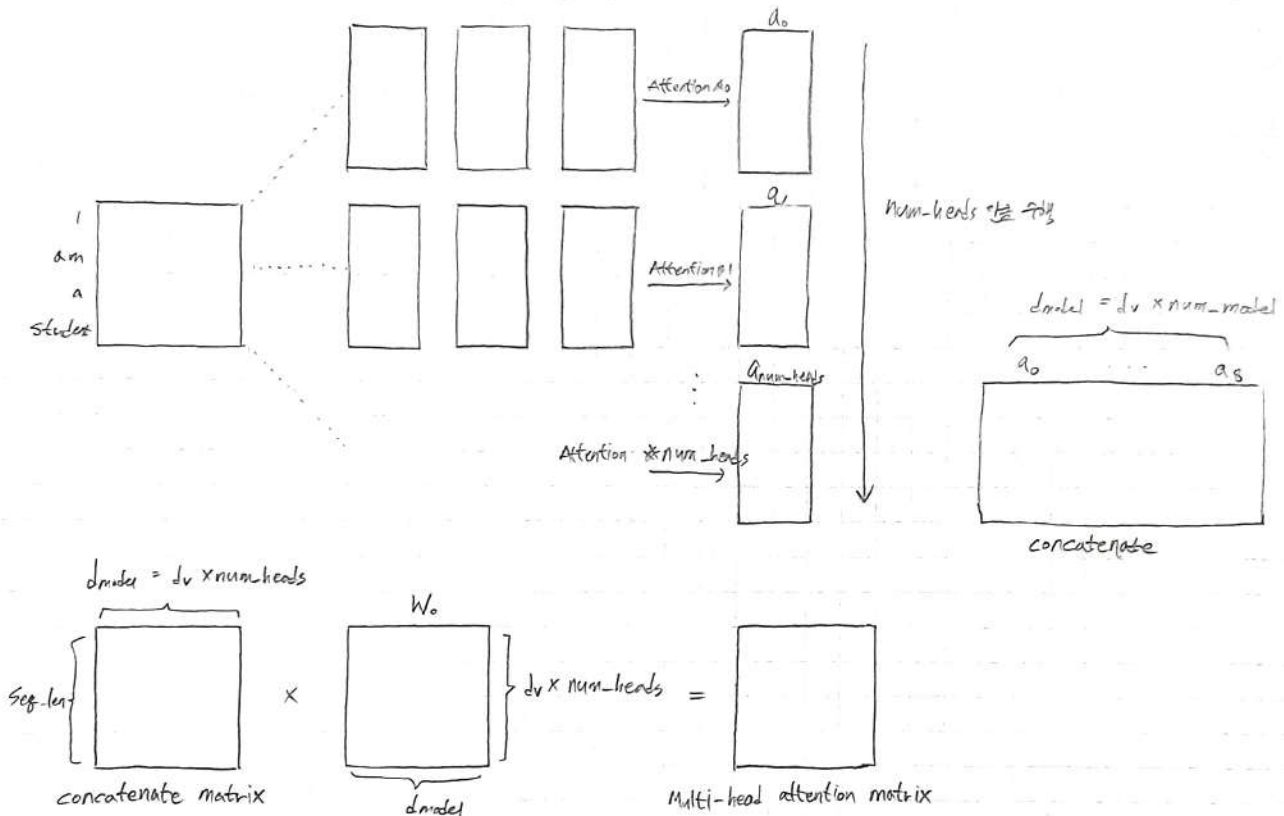
- Encoder

num_layers (하이퍼 파라미터) 개수의 인코딩 층을 쌓는다.

인코딩 층은 2개의 서브층(sub-layer) 멀티 헤드 어텐션 (self attention 병렬 사용)과 피드 포워드 신경망을 병렬로 사용

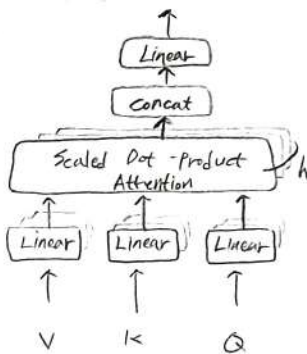
- Multi-Head Attention

병렬 어텐션 : 다른 서브층으로 정보를 추가하겠다는 것



병렬 어텐션 수행 → 모든 어텐션 헤드 연결 (concatenate) : (seq-len, d_{model})

→ 연결한 행렬을 가중치 W_0 와 곱함 → 멀티-헤드 어텐션 행렬 : (seq-len, d_{model})

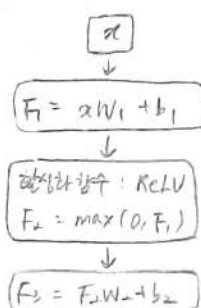
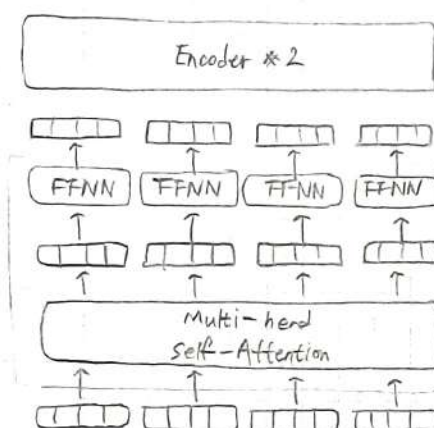


- 1) W_Q, W_K, W_V 에 해당하는 d_{model} 크기의 밀집층 (Dense layer)를 지내게 한다
- 2) 지정된 헤드 수 (num_heads) 만큼 나누자 (split)
- 3) Scaled Dot-Product Attention.
- 4) 나누어진 헤드들을 연결 (concatenate) 한다
- 5) W_0 에 해당하는 밀집층을 지내게 한다.

Position-wise FFNN (Fully-connected FFNN)

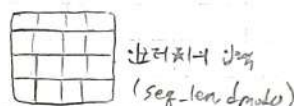
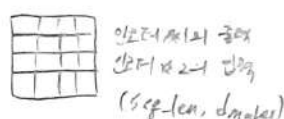
$$FFNN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

매개변수 W_1, b_1, W_2, b_2 는
하나의 인코더 층 내에서는 다른 문장, 다른 단어에
동일하게 적용, 인코더 층마다 다른 값



일터 헤드 어텐션의 결과로 나온
($seq-len, d_{model}$)의 크기를 차이는 해질
가능한 행렬 W_1 의 크기는 (d_{model}, d_{ff})

가능한 행렬 W_2 의 크기는 (d_{ff}, d_{model})



잔차 연결 (Residual Connection) 과 층 정규화 (Layer Normalization)

Add & Norm 기법

1) 잔차 연결

서브층의 입력과 출력을 더한다. $x + \text{Sublayer}(x)$

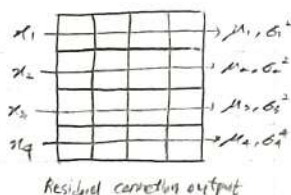
$$\text{ex. } H(x) = x + \text{Multi-head attention}(x)$$

$$\text{Multi-head attention input} + \text{Multi-head attention output} = \text{Residual Connection output}$$

2) 층 정규화

$$LN = \text{LayerNorm}(x + \text{Sublayer}(x))$$

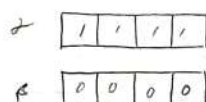
레이어의 마지막 처리에 대해서 표준화 부분을 추가, 이를 가지고 어떤 수치를 통해 값을 정규화하여 학습을 돕는다



1) 표준화 부분용에 정규화

$$\text{벡터 } x_i \text{의 각 } k \text{번째 값에 정규화} \quad \hat{x}_{i,k} = \frac{x_{i,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

2) 간헐적 베타 도입



$$\Rightarrow h_{i,j} = \gamma \hat{x}_{i,j} + \beta = \text{LayerNorm}(x_{i,j})$$

x_i 는 $h_{i,j}$ 값을 베타로 정규화된다.

- Decoder

↳ Masked Multi-head Self-Attention : 첫 번째 서브층

look-ahead mask : 트랜스포머는 문장 토큰으로 입력을 한 번에 받기에
현재 시점의 예측에서 현재 시점보다 미래에 있는 토큰들을 참조하지 않도록 하는 것.

이런 식으로 스코어 행렬에서 관심을 작성한다

↳ Encoder - Decoder Attention (Multi-head Attention) : 두 번째 서브층

Query : 디코더 토큰 / Key, Value : 인코더 토큰 (Self Attention x)

- 학습률

학습률 스케줄러 (Learning rate Scheduler) : 미리 학습 율을 정해두고 그 율에 따라 학습률이 조정되는 방법

step-num (단계) : 옵티마이저가 매개 변수를 업데이트 하는 한 번의 간격 값

step-num이 warmup-steps 보다 작을 경우 학습률을 선형적으로 증가시키다가, 소달하면 step-num이 커지면 점차적으로 감소시킨다.

$$lr_{rate} = lr_{init} \times \min(step_num^{-0.5}, step_num \times warmup_steps^{-1.5})$$