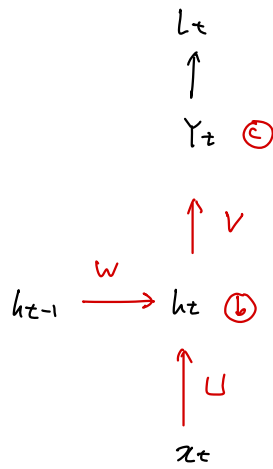


Vanilla RNN

$$\begin{cases} h_t = \sigma(W h_{t-1} + U x_t + b) \\ y_t = V h_t + c \end{cases}$$

$$\begin{cases} L_t = \frac{1}{2} (y_t - z_t)^2 \\ L = \sum L_t \end{cases}$$



$$\frac{\partial L}{\partial V} = \sum_t \underbrace{\left(\frac{\partial L_t}{\partial y_t} \right)}_{y_t - z_t} \cdot \underbrace{\left(\frac{\partial y_t}{\partial V} \right)}_{h_t}$$

$$\frac{\partial L}{\partial c} = \sum_t \underbrace{\left(\frac{\partial L_t}{\partial y_t} \right)}_{y_t - z_t} \cdot \underbrace{\left(\frac{\partial y_t}{\partial c} \right)}_1$$

Backprop Through Time

$$\frac{\partial L}{\partial w} = \left(\frac{\partial L_1}{\partial h_1} \cdot \frac{\partial h_1}{\partial w} \right) + \left(\frac{\partial L_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial w} \right) + \left(\frac{\partial L_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial w} \right)$$

$$\frac{\partial h_t}{\partial w} = \sum_{q=1}^L \left(\frac{\partial h_{t+1}}{\partial h_t} \dots \underbrace{\left(\frac{\partial h_t}{\partial h_{t-1}} \right)}_{(1-h_{t-1}^2) \cdot W} \cdot \underbrace{\left(\frac{\partial h_t}{\partial w} \right)}_{(1-h_{t-1}^2) h_{t-1}} \right)$$

Note

$$(\tanh(x))' = 1 - \tanh^2(x)$$

$$\left| \frac{\partial h_t}{\partial h_{t-1}} \right| > 1 \quad \text{then} \quad \text{clipping} : \text{threshold} \cdot \frac{g}{\|g\|}$$

$$\left| \frac{\partial h_t}{\partial h_{t-1}} \right| < 1 \quad \text{then} \quad \text{clipping} : 0 \rightarrow 0$$

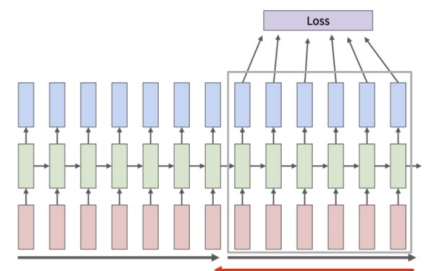
$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_{t+1}} \cdot \underbrace{\left(\frac{\partial h_{t+1}}{\partial h_t} \right)}_{\approx I \text{ then}}$$

grad can flow

LSTM
GRU

Note.

Truncated Backprop.



Vanilla RSTM

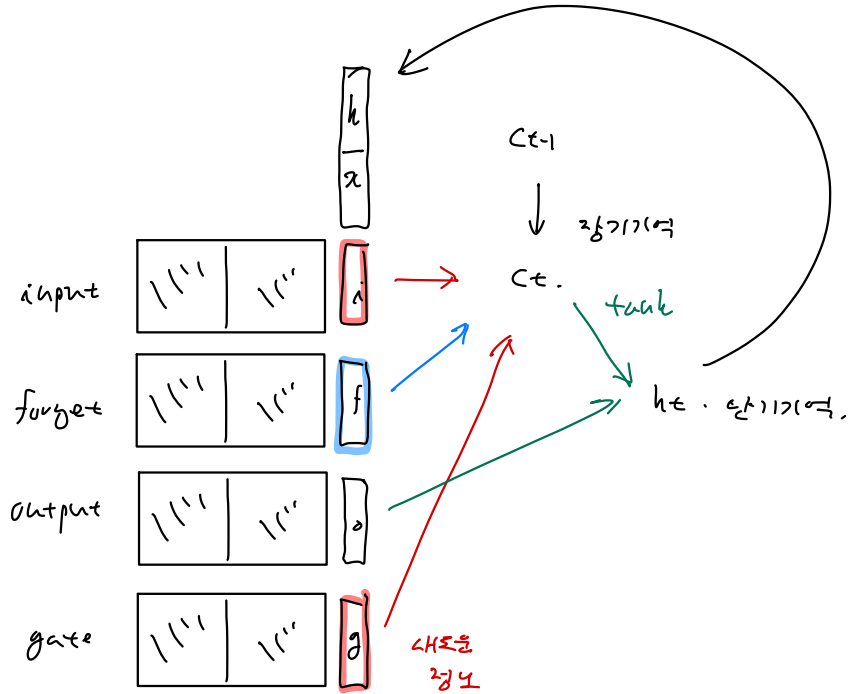
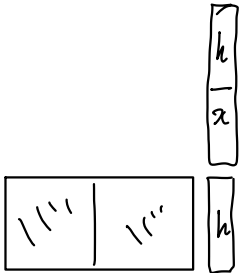
$$h_t = \tanh(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix})$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$



$$\frac{\partial L}{\partial c_t} = \frac{\partial L}{\partial c_{t+1}} \frac{\partial c_{t+1}}{\partial c_t}$$

$$c_{t+1} = f_{t+1} \odot c_t + i_{t+1} \odot g_{t+1}$$

$$\frac{\partial c_{t+1}}{\partial c_t} = f_{t+1} + (\dots)$$

$$\frac{\partial L}{\partial c_t} \approx \frac{\partial L}{\partial c_{t+1}} \frac{\partial c_{t+1}}{\partial c_t} \approx \frac{\partial L}{\partial c_{t+1}} f_{t+1}$$

uninterrupted
grad flow

$$h_t = \text{LSTM}(z_t, h_{t-1}, c_{t-1}; \theta)$$

Note. One-Hot Encoding

$$z_t = W_{emb} \hat{x}_t$$

one-hot

"a" [0 0 1]

"man" [0 1 0]

$$\hat{x}_{t+1} = W_d h_t + b_d$$

"다음 단어로 등장할 확률"

training

$$\sum_{t=1}^T \mathcal{L}(\hat{x}_t, x_t)$$

cross-entropy loss $-\sum p(i) \log f(i)$

inference

$$\left[\begin{array}{c} \text{Greedy} \\ \text{argmax } \mathcal{L}(x_t) \end{array} \right] \quad \text{or} \quad \left[\begin{array}{c} \text{stochastic} \\ x_t \sim \mathcal{L}(x_t) \end{array} \right]$$

... <EOS> \Rightarrow 생성 중단.

img to text

$$f_{\theta}(y | f_{\phi}(x))$$

CNN

RNN

"use feature vector as first word"

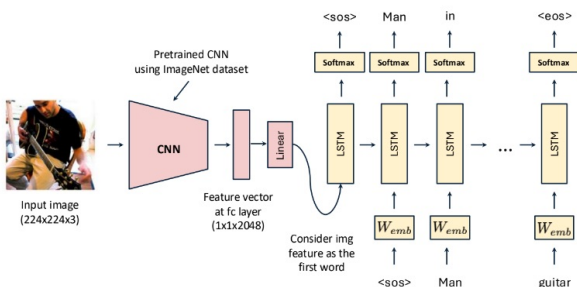
"이 정보를 읽어서 무슨 단어를?"

개념 ① 매 순간 단어 예측이나 할게 있네

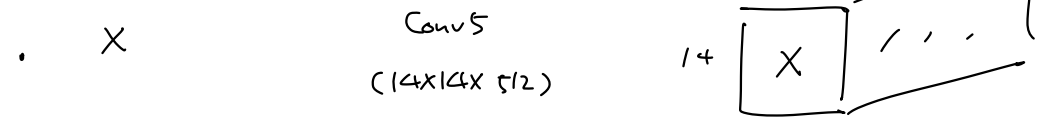
개념 ② Attention

: 단어를 생성할 때

관련된 이미지의 특징 벡터를 활용!



Spatial Attention



Score

$e^{(t)} = f_{att}(X, h_{t-1})$

각 구간은 "a bird"와 얼마나 관련 있는가?

$\alpha^{(t)}$

Normalize

$\alpha^{(t)} = \text{softmax}(e^{(t)})$

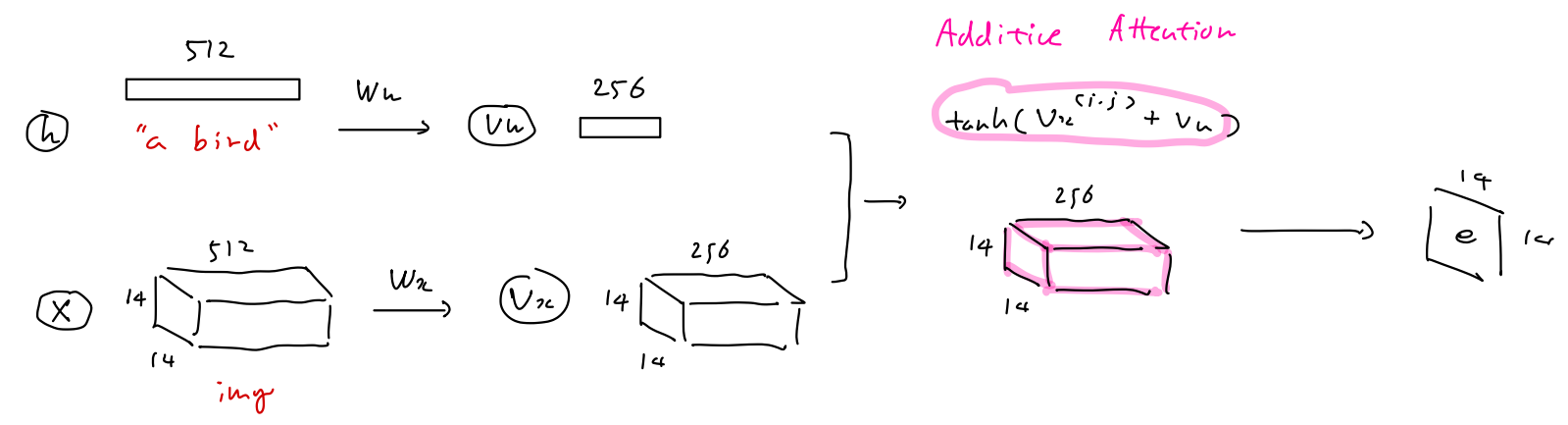
Context

$z^{(t)} = \sum_{i,j} \alpha^{(t)}_{i,j} x_{i,j}$

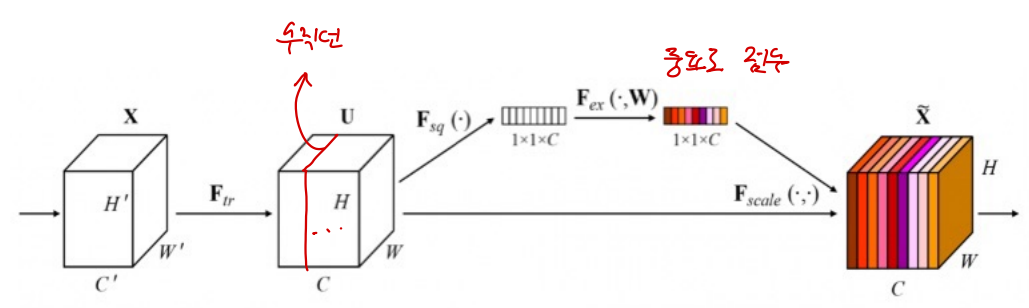
$z^{(t)}$

$h^{(t)} = \text{LSTM}(h_{t-1}, \text{Embed } y_t, z^{(t)})$

Simple Spatial Attention



Squeeze-and-Excitation Net : Channel Attention



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

$$\text{softmax} \frac{e^{z_j}}{\sum e^{z_j}}$$

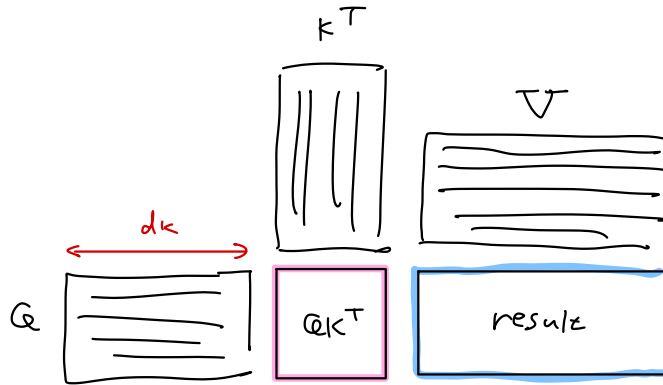
train



$$\text{Query} = XW_Q : n \times d_k$$

$$\text{Key} = XW_K : n \times d_k$$

$$\text{Value} = XW_V : n \times d_v$$



Self-Attention

vs

Cross-Attention

$$Q, K, V \Rightarrow X$$

$$Q \Rightarrow Y \text{ (decoder)}$$

$$K, V \Rightarrow X \text{ (encoder)}$$

가계년역

(역시)

이러지 않음

Q : 한영문

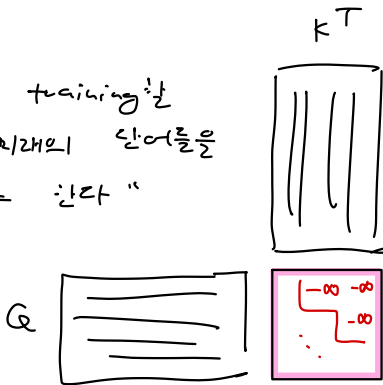
Q : h-1

K, V : 원문

K, V : 이터리 블록

Masked Self-Attention

"decoder를 training할 때에 는 미래의 단어들을 가려야 한다"



i^{th} row < j^{th} column
then $- \infty$

positional encoding

$$: x_j + p_j \text{ or } \text{concat}(x_j, p_j)$$

$$2d \begin{bmatrix} \sin(w_1 \times pos) \\ \cos(w_1 \times pos) \\ \vdots \\ \sin(w_d \times pos) \\ \cos(w_d \times pos) \end{bmatrix}$$

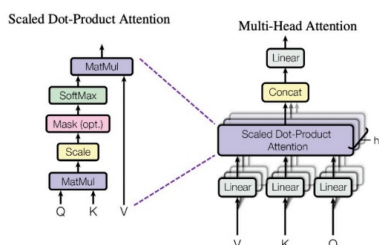
$$\text{where } w_k = \frac{1}{10000^{k/d}}$$

Multi-Head

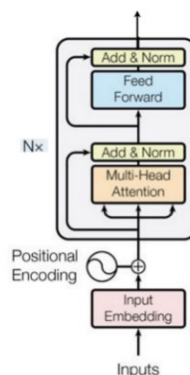
h개의 다중

Attention

W_Q, W_K, W_V 를 이용



Transformer Encoder



$$\rightarrow \text{LayerNorm}(h + \text{FF}(h))$$

$$\rightarrow \text{LayerNorm}(h + \text{Attn}(h, h, h))$$

$$\rightarrow \text{FF}(x_i)$$

"Attention is All You Need"

Attention & Transformer 은이 이이

새끼를 얻는 Attention

$$Q \begin{bmatrix} \text{어려운} \end{bmatrix} \quad K \begin{bmatrix} \text{어려운} & \text{어려운} & \text{어려운} \end{bmatrix} \quad V \begin{bmatrix} \text{어려운} & \dots & \text{어려운} \end{bmatrix}$$

★

수식어 얻는 Attention

$$Q \begin{bmatrix} \text{어려운} \end{bmatrix} \quad K \begin{bmatrix} \text{어려운} & \text{어려운} & \text{어려운} \end{bmatrix} \quad V \begin{bmatrix} \text{어려운} & \dots & \text{어려운} \end{bmatrix} \Rightarrow \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \\ \vdots \end{bmatrix}$$

★

다음 어휘를
이동시킬래.

Linearization softmax $\left(\frac{QK^T}{\sqrt{d_k}} \right)$ is too big!

$$\text{sim}(f, k) = \exp\left(\frac{f^T k}{\sqrt{D}}\right) \rightarrow \text{keu}(f, k) = \phi(f)^T \phi(k)$$

$$f_{\text{len}} \quad \tilde{V}_i = \frac{\sum_{j=1}^N \text{keu}(f_i, k_j) \cdot v_j}{\sum_{j=1}^N \text{keu}(f_i, k_j)} = \frac{\phi(f_i)^T \sum_{j=1}^N \phi(k_j) \cdot v_j}{\phi(f_i)^T \sum_{j=1}^N \phi(k_j)}$$

2차원 공간
 미리 계산 (v x d_k)
 2차원 공간
 미리 계산 (d_k)

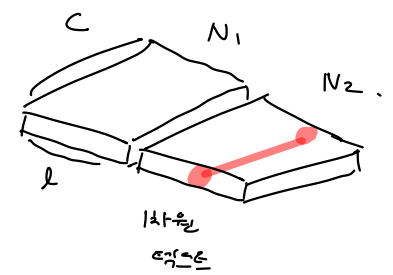
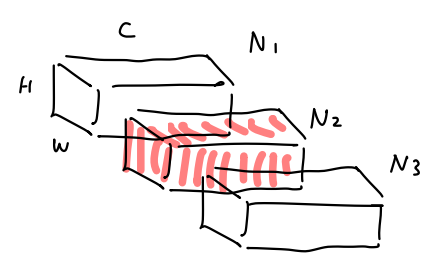
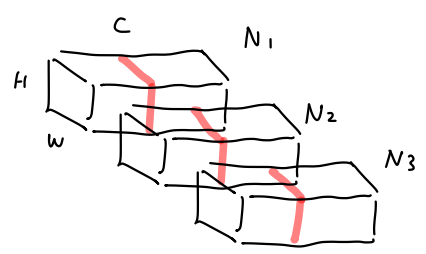
Normalization -

배치 차원 채널 높이/너비
 $\downarrow \quad \downarrow \quad \downarrow$
 $N \times C \times (H \times W)$
 $X \in \mathbb{R}$ or L

$$\text{LN}(x_i) = \beta \cdot \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \checkmark$$

값은 가변적일 수 있다!

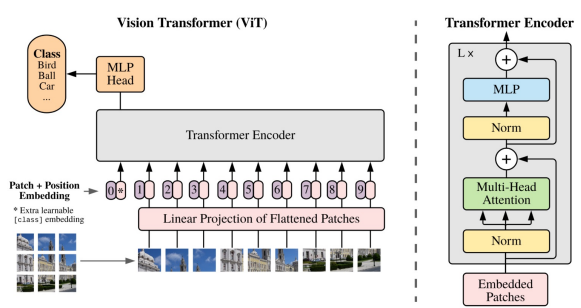
- ① Batch Norm (CNN)
- ② Layer Norm (CNN)
- ③ Layer Norm (Transformer)



$$\mu_c = \frac{1}{NHW} \sum_{n,h,w} x_{n,c,h,w}$$

$$\mu_n = \frac{1}{CHW} \sum_{c,h,w} x_{n,c,h,w}$$

$$\mu_{n,l} = \frac{1}{C} \sum_c x_{n,l,c}$$



이제 이를 N개로 잘라낸다. (patching)

$$Z_0 = \begin{bmatrix} \overbrace{x_{cls}}^{D_{cls}} ; \overbrace{x_p^1 E}^{D_{p^1}} ; \dots ; \overbrace{x_p^N E}^{D_{p^N}} \end{bmatrix} + E_{pos}$$

$$\begin{cases} Z'_l = \text{MSA}(\text{LN}(Z_{l-1})) + Z_{l-1} \\ Z_l = \text{MLP}(\text{LN}(Z'_l)) + Z'_l \end{cases}$$

각번 앞에 있는 D차원 정도인 가려냄.

GELU = $x \cdot \Phi(x)$
 가우시안 분포의 누적분포함수

CNN

: 작은 filter가

sliding window로

feature map을 만든다.

✓ 계층적 구조

: 엡지, 선 \rightarrow 눈, 코, 나뭇

✓ 지역적 처리

: 필터 크기(에) 의존

ViT

: 이미지를 patch로 나눌 수

$\begin{cases} \text{flattening} \\ \text{position encoding} \end{cases}$ 을 처리

transformer에 넣는다.

✓ 전역적 처리

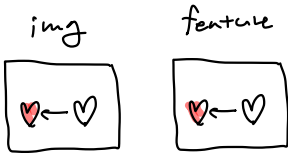
: 전체 간의 관계

관선에 관계

강한 사전가정

• locality

• Translation Equivariance



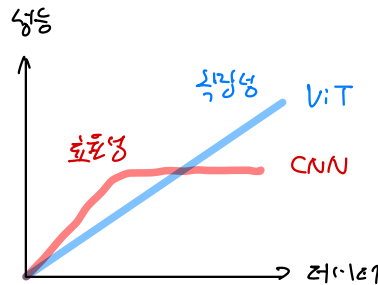
$$F(T(x)) = T(F(x))$$

약한 사전가정

• ~~locality~~ locality

• ~~Translation~~ Translation Equivariance

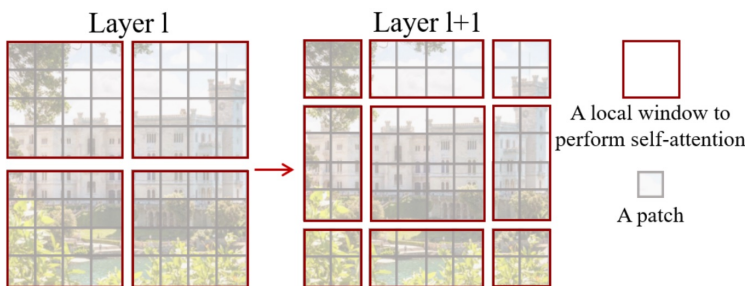
\Leftarrow position encoding



Note.

translation invariance

$$F(T(x)) = F(x)$$



① ViA vs W-MSA.

Window 기반의 Attention.

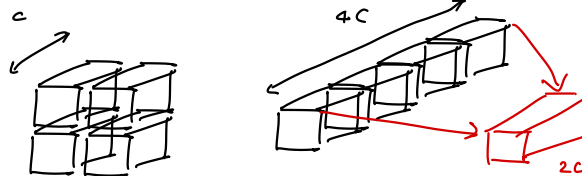
$$\begin{cases} \Omega(\text{MSA}) = 4hwc^2 + 2(hw)^2 \cdot c \\ \Omega(\text{W-MSA}) = \dots + 2M^2 hwc \end{cases}$$

② Shifted window

: 윈도우 크기 고정

$M/2$ 만큼 겹칠.

③ Patch Merging



like pooling

in CNN.