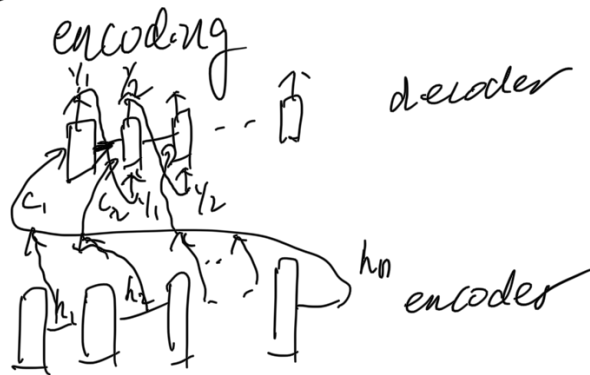
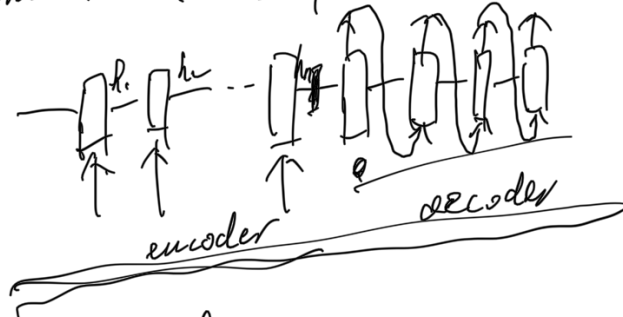


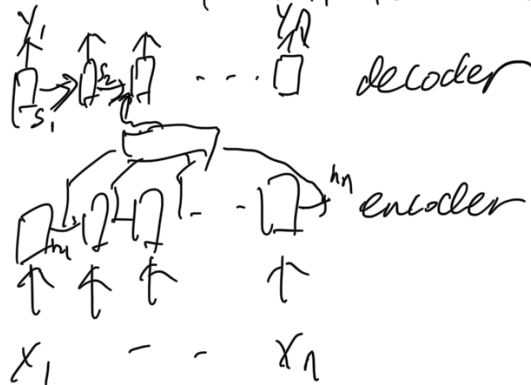
Lecture 15

Attention and Transformer



Input sequence $X = [x_1, x_2, \dots, x_n]$, length = n

output sequence $Y = [y_1, y_2, \dots, y_m]$, length = m



$$y_i = f(s_i), \quad s_i = g(s_{i-1}, c_i, y_{i-1})$$

\uparrow \leftarrow
 $x_1 \quad \dots \quad x_n$

$$\rightarrow c_i = \sum_{k=1}^n a_{i,k} h_k$$

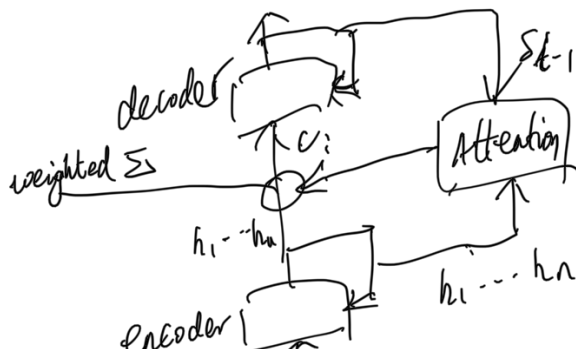
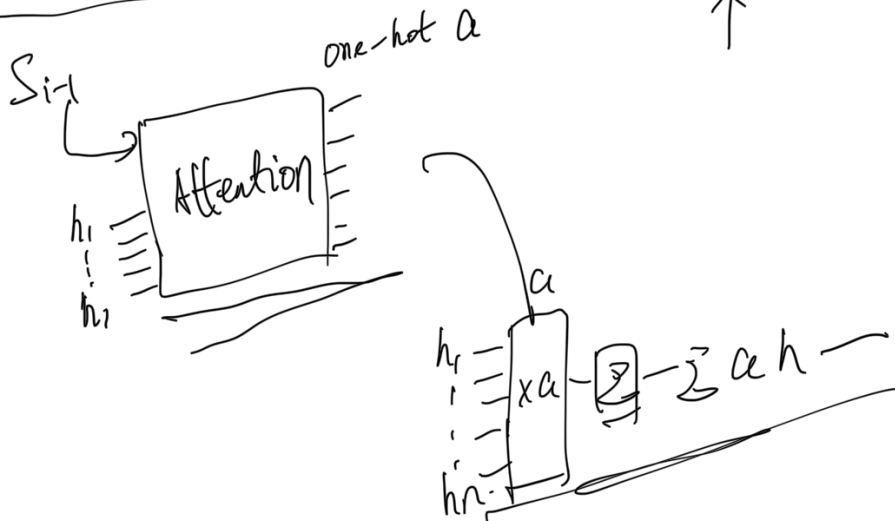
$k=1$ $\frac{1}{n}$
weights

O1 - learn $a_{i,k}$ directly from data, then A is fixed. think A is a matrix of $a_{i,k}$.

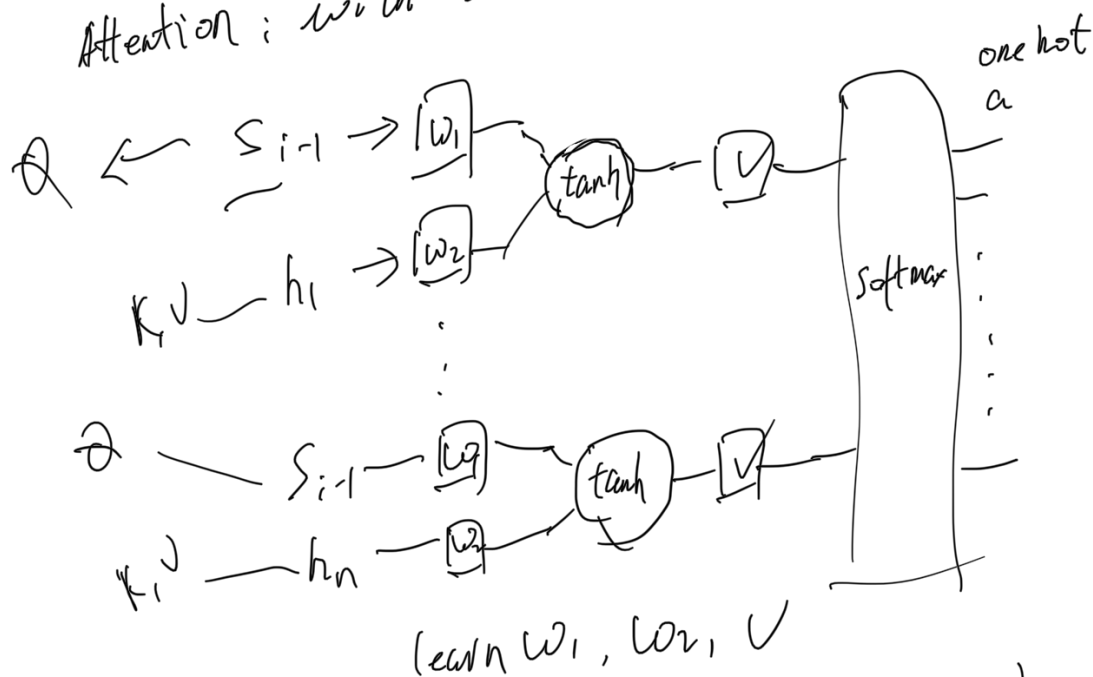
✓ O2 - relate $a_{i,k}$ to some variables, A is a function

$$a_{i,k} = \text{align}(y_i, x_k) \\ = \frac{\exp(\text{score}(s_{i-1}, h_k))}{\sum \dots} \quad \checkmark \text{softmax}$$

$$\text{score}(s_{i-1}, h_k) = v \tanh(w_1 s_{i-1} + w_2 h_k) \leftarrow$$



Attention: with additive



Several attention mechanisms, (score functions)

Content-based attention $\text{score}(S_{i-1}, h_k) = \text{cwise}(S_{i-1}, h_k)$

General $\text{score}(S_{i-1}, h_k) = S_{i-1}^T W h_k$
 \uparrow to be learned

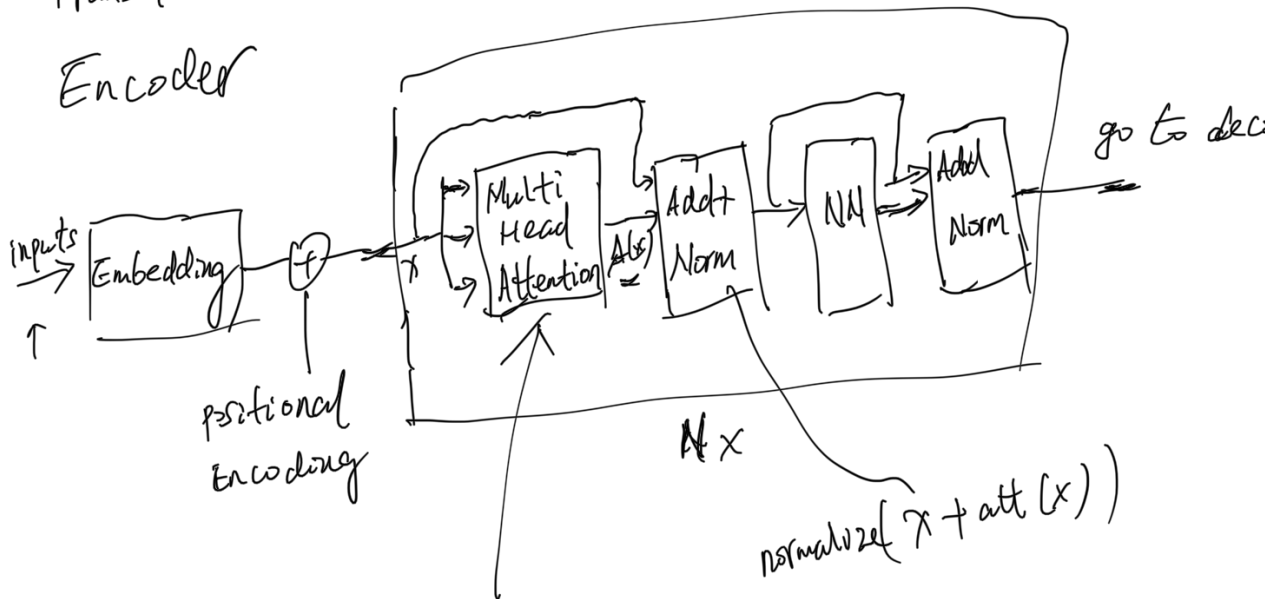
Dot-Product $\text{score}(S_{i-1}, h_k) = S_{i-1}^T h_i$

Scaled Dot product: $\text{score}(S_{i-1}, h_k) = \frac{S_{i-1}^T h_i}{\sqrt{n}}$

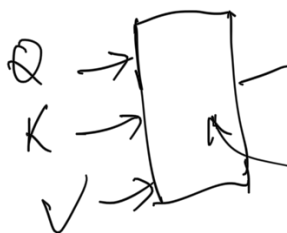
Global/soft $\rightarrow h_1 \dots h_n$

Local/Hard $\rightarrow h_i \dots h_j$ subset of $1, \dots, n$

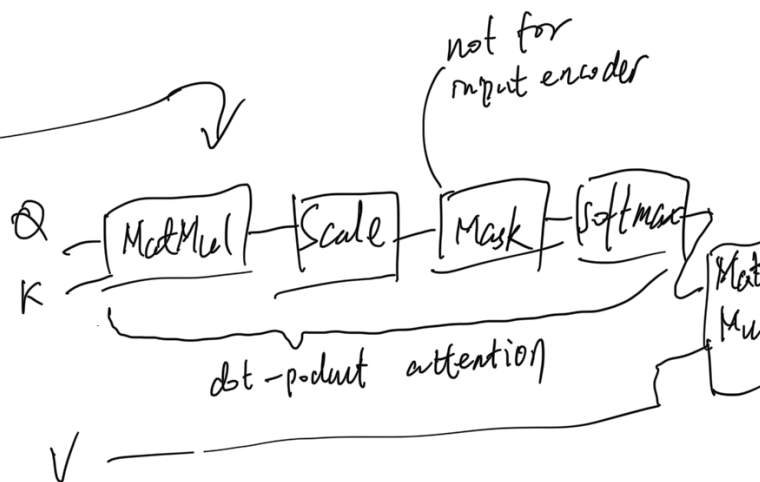
Transformer Encoder



self-Attention,

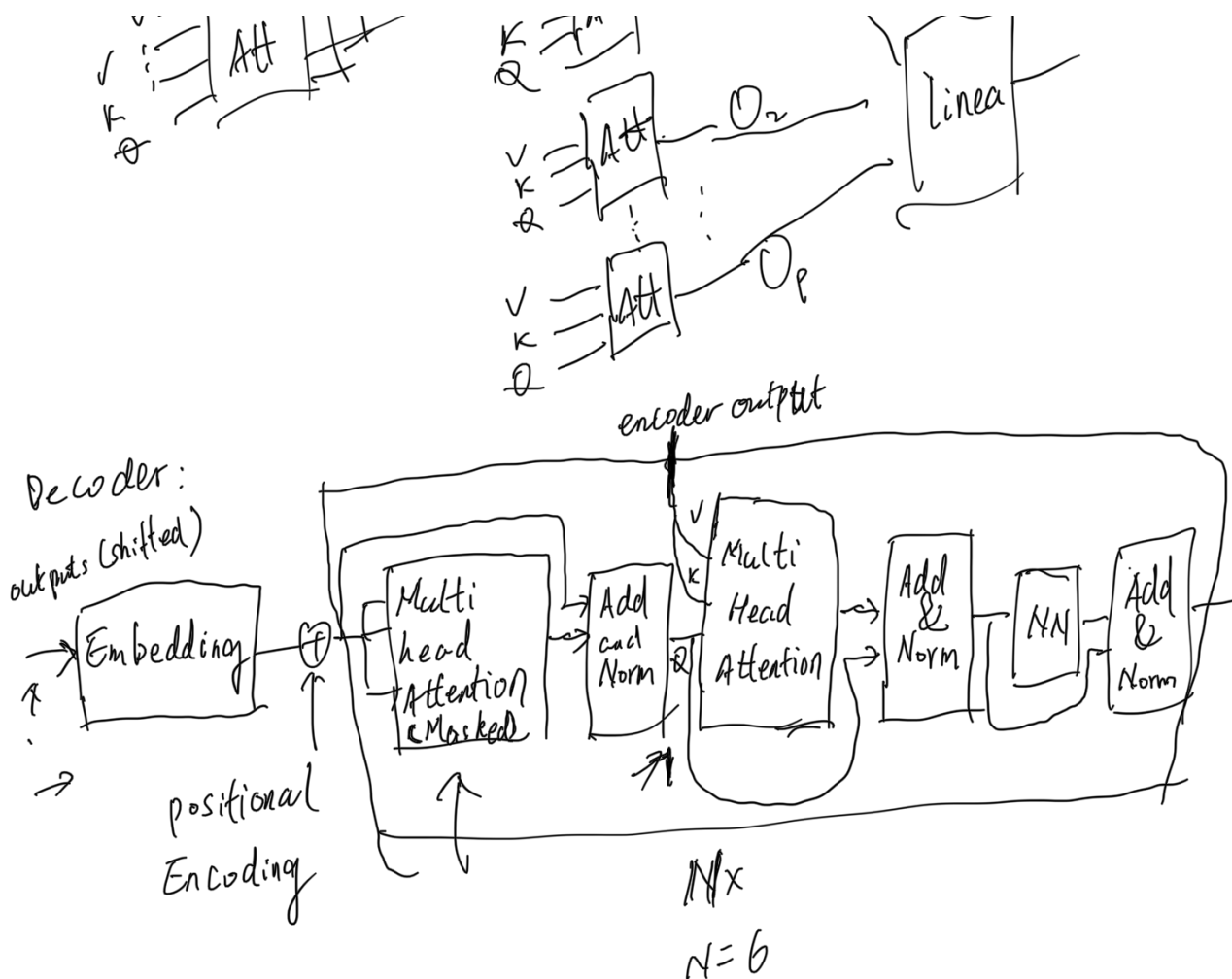


Query, Key, Value.



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





Suppose the output is "I live in Tampa"

↑

I live in Tampa
1 2 3 4 5

Mask:

For 1st input

0 0 0 0

2nd

1 0 0 0

3rd

1 1 0 0

positional encoding

$$f(t) = \begin{cases} \sin(\omega_k t) & \text{if } i = 2k \\ \cos(\omega_k t) & \text{if } i = 2k+1 \end{cases}$$

$$\begin{bmatrix} \sin(\omega_1 t) \\ \cos(\omega_1 t) \\ \sin(\omega_2 t) \\ \cos(\omega_2 t) \\ \vdots \\ \sin(\omega_{d/2} t) \\ \cos(\omega_{d/2} t) \end{bmatrix}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$