

Some Reviews :

How do we deal with time dependent data ?

① time as input = ??

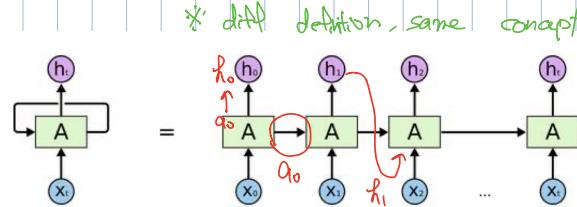
(Time 2 Vec ?, one-hot encoding ?)

② measure dependency between data

\Rightarrow to have memory property.

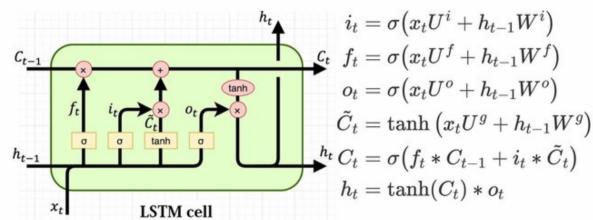
Recurrent Network:

Vanilla RNN :

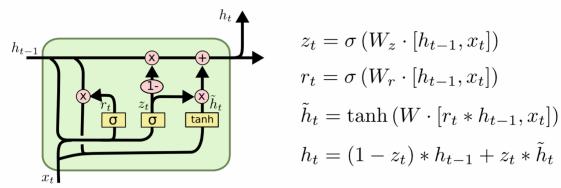


A: is the same network

LSTM :



GRU :



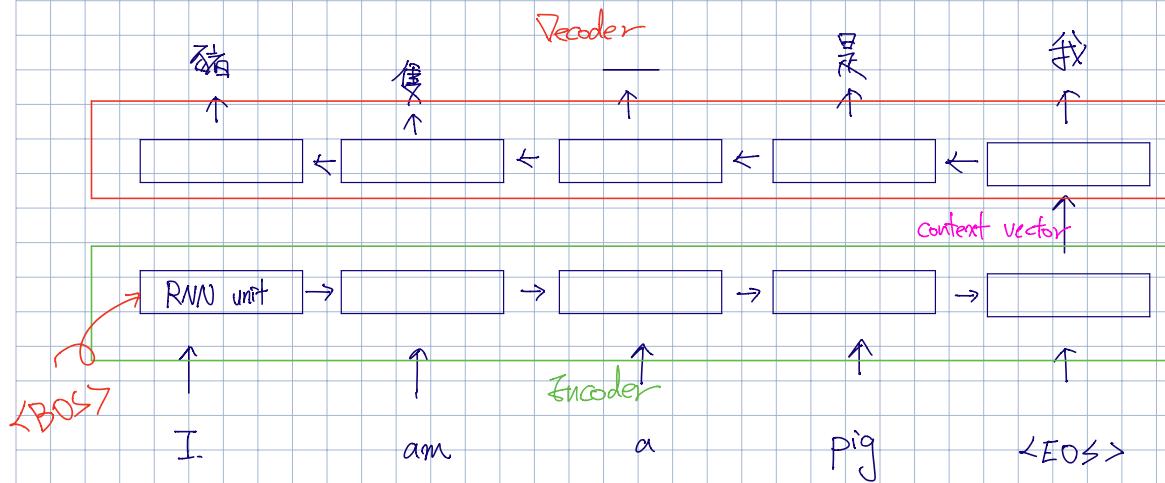
5 mins

How about translation tasks ?

Input : seq , Output : seq .

(length of seq cannot be pre-determined)

⇒ Encoder - Decoder Model.



* Encoder actions are trivial, but

Decoder: In training stage :

$h_{t-1} \neq h_t \Rightarrow$ Using truth answer as input of h_t

y_{t-1}

Bi directional LSTM instead of just RNN units :

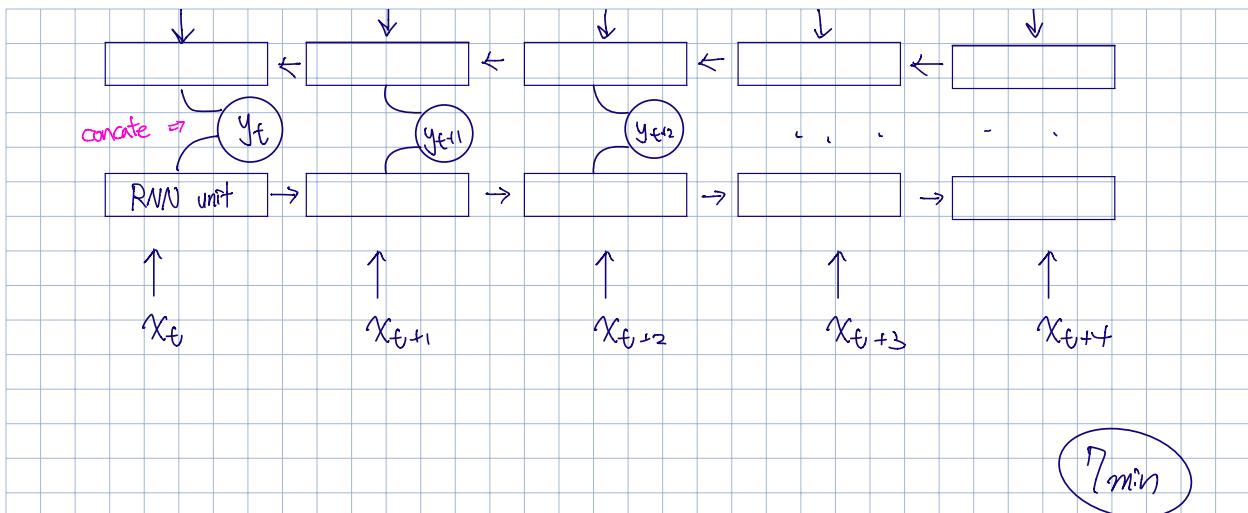
x_t

x_{t+1}

x_{t+2}

x_{t+3}

x_{t+4}



7 min

But still, we will have some problem.

① context vector \Rightarrow a fixed dim vector.

\Rightarrow no matter how long \Rightarrow the sequence

while use ReLU.

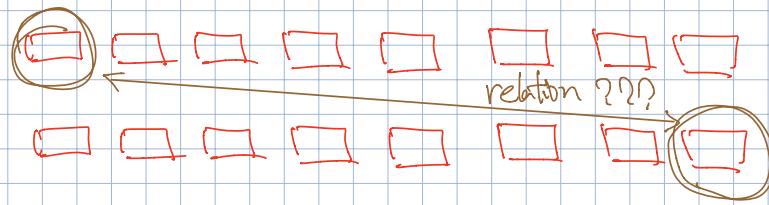
② RNN has the gradient vanishing, exploding problem.



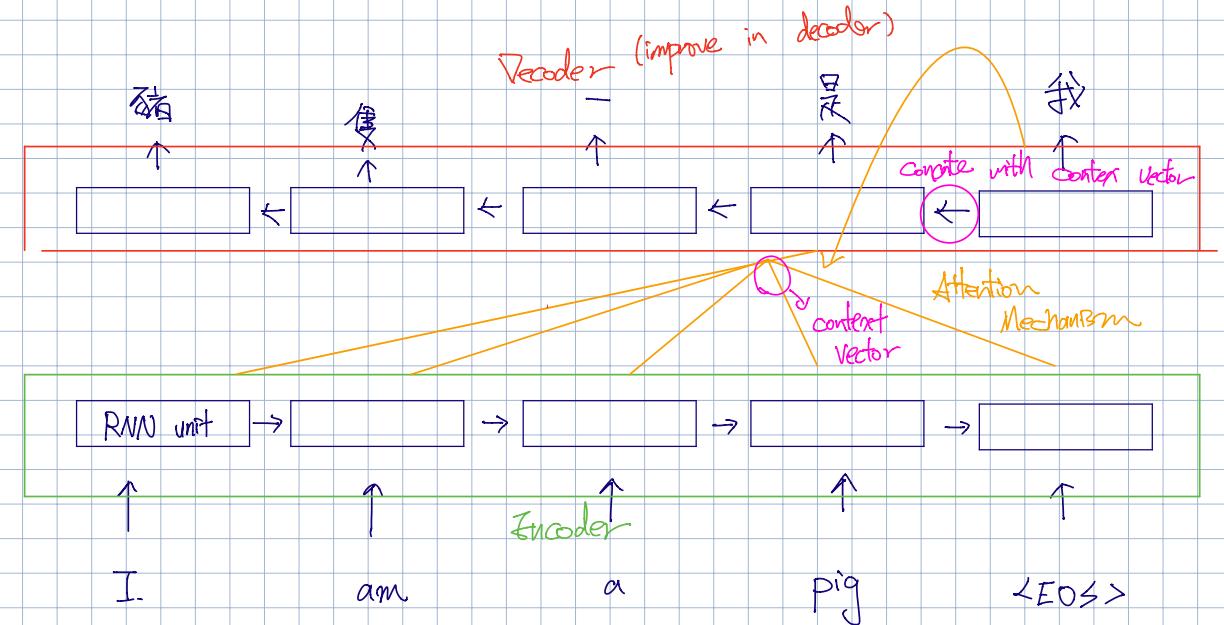
CNN \Rightarrow receptive field problem.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

③ As the seq being longer, the previous token has no effectness. Same as in Decoder.



Attention !



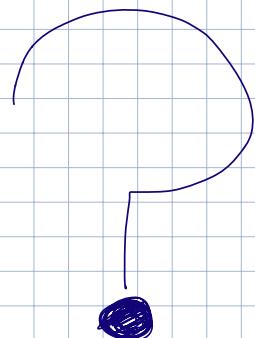
~~Calculate~~ Calculate each units' their own context vector

There are lots of ways to do "attention":

Ex: dot product attention, multiplicative attention, additive attention

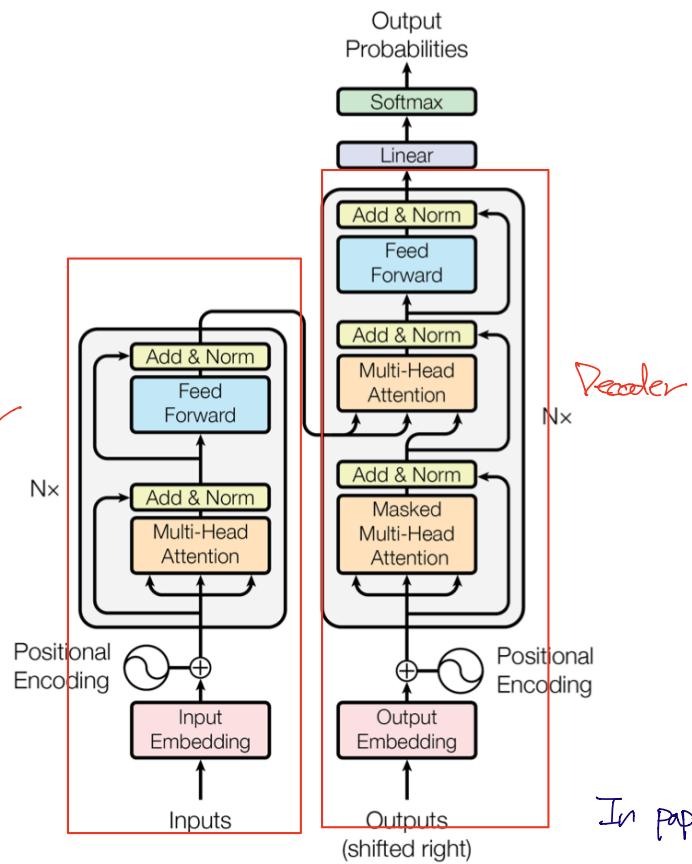
We now solved problem D. ②.

=> Is there anything else we can improve?



7 min

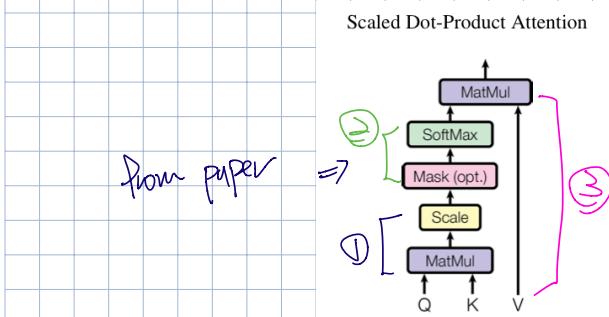
Transformer: self-attention
(Bye² RNN, CNN)



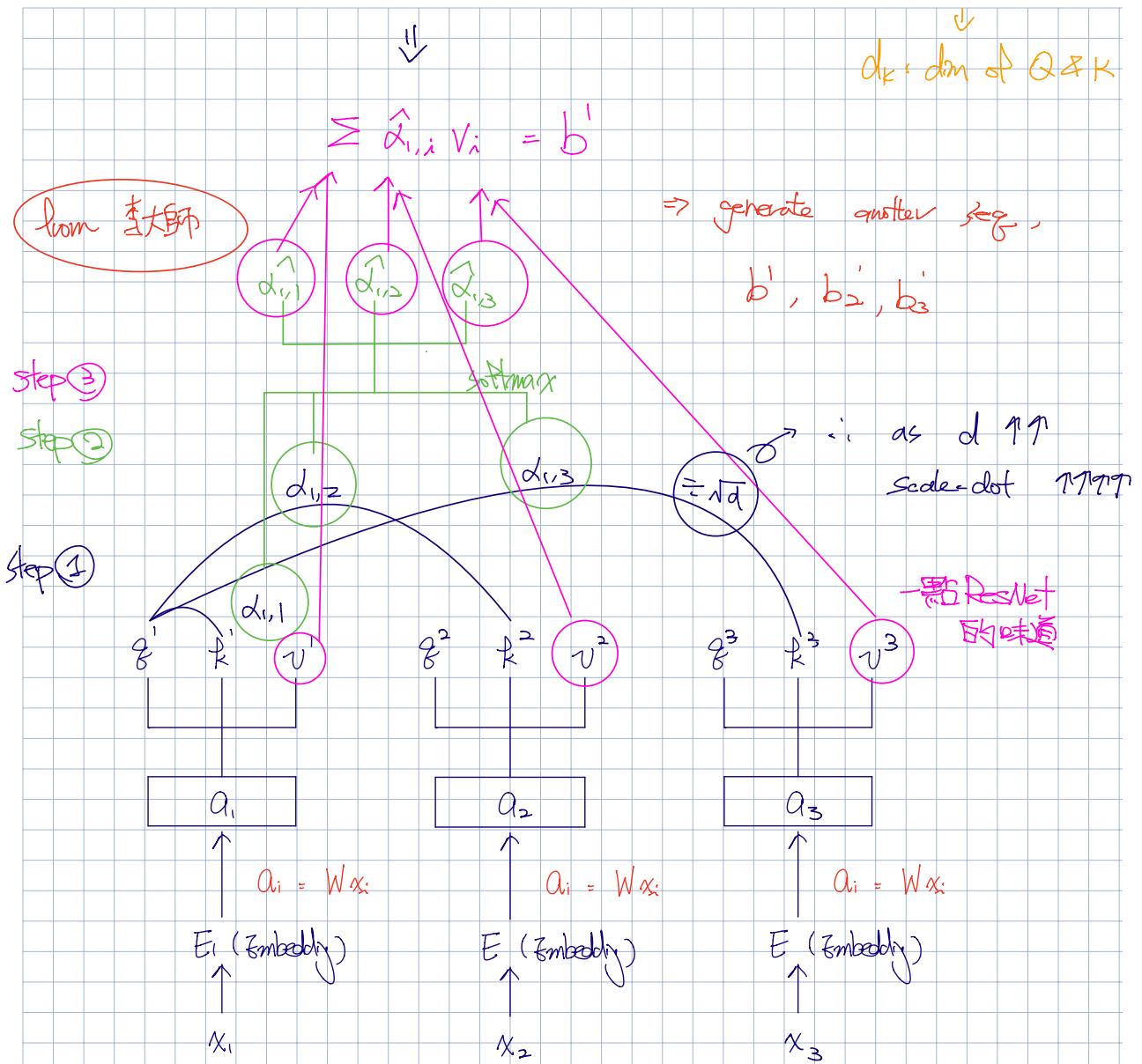
Don't worry! we will go through step by step.

1. Attention

(i) Scaled dot-product attention



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



Generalization form:

$$\begin{aligned} ① \quad W^Q a_i &= q^i \\ W^K a_i &= k^i \\ W^V a_i &= v^i \end{aligned} \Rightarrow$$

$$\begin{aligned} [W^Q] [a^1 a^2 a^3] &= [q^1 q^2 q^3] \\ [W^K] [a] &= [k] \\ [W^V] [a] &= [v] \end{aligned}$$

Similarly ,

②

$$\begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \end{bmatrix} \begin{bmatrix} g^1 & g^2 & g^3 \end{bmatrix}$$

$$A \downarrow \text{softmax}$$
$$\hat{A}$$

Similarly

③

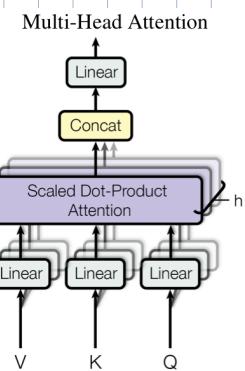
$$\begin{bmatrix} b_1 & b_2 & b_3 \end{bmatrix} = \begin{bmatrix} V_1 & V_2 & V_3 \end{bmatrix} \hat{A}$$

~~Star~~

Eventually , we are just calculating matrix multiplication !!

↳ speed up , parallel computation in GPU.

(ii) Multi-head Attention



MultiHead (Q, K, V)

= Concat (head₁, head₂, ..., head_n) W⁰

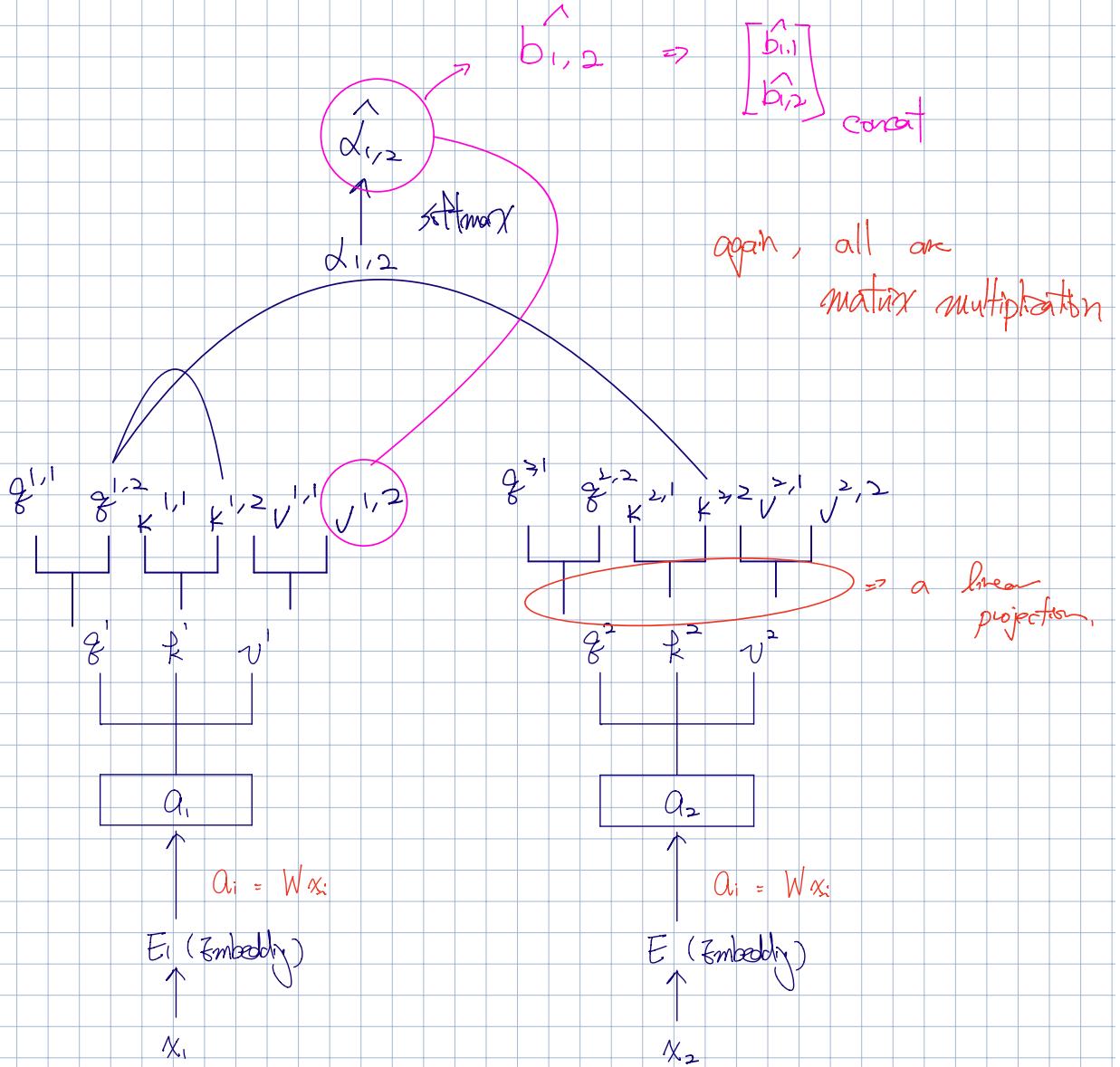
\Downarrow

$f \downarrow \quad 1 \downarrow \quad 1 \downarrow \quad \dots \downarrow \quad Q \quad \dots \quad K \quad \downarrow \quad V$

$$\text{head}_i = \text{Attention}(QW_i, KW_i, VW_i)$$

Similar as above, but heads ?

\Downarrow
GR
 $d_{\text{model}} \times d_k$



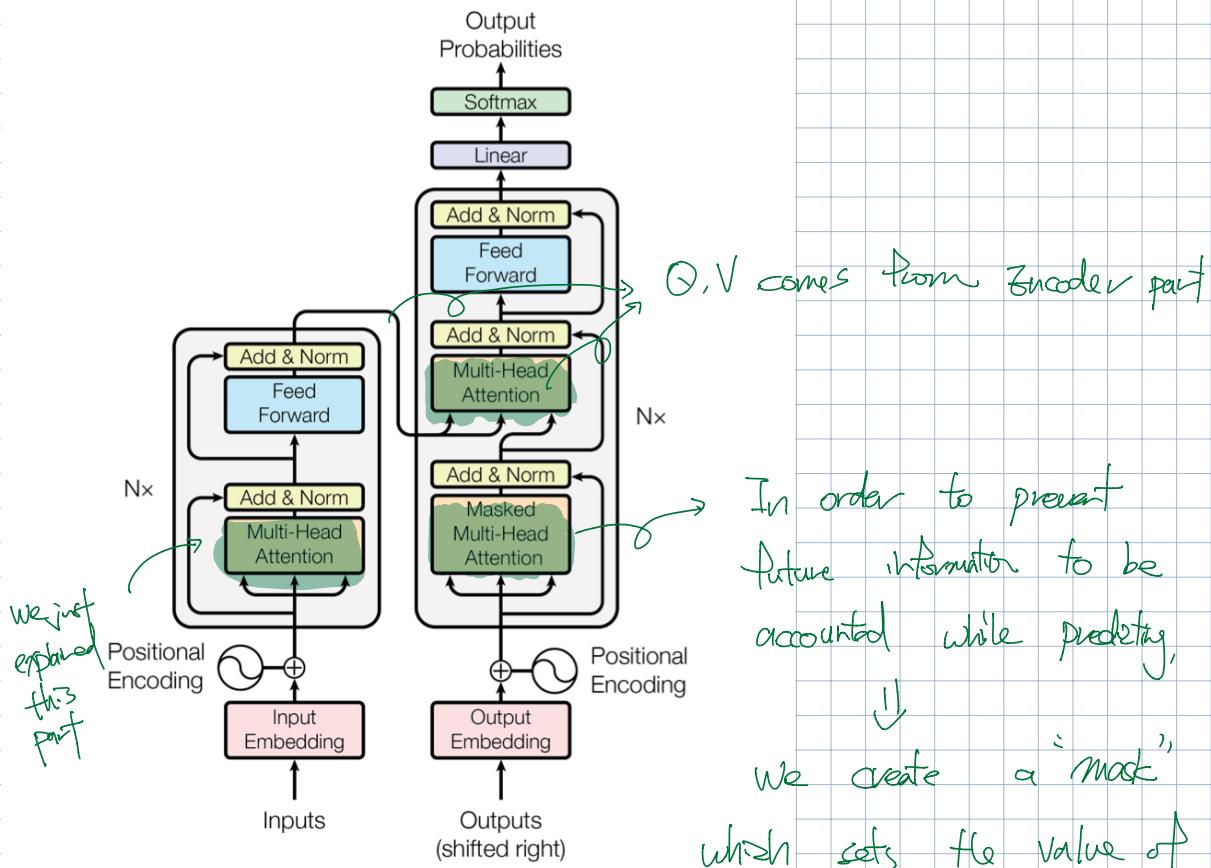
In paper, # heads = 8, $d_k, q, v = 64$, $d_{\text{model}} = 64 \times 8 = 512$

\hookrightarrow we can regard heads as different

Character of words, ex: grammar, syntax, ...

Back to Graphi

20min

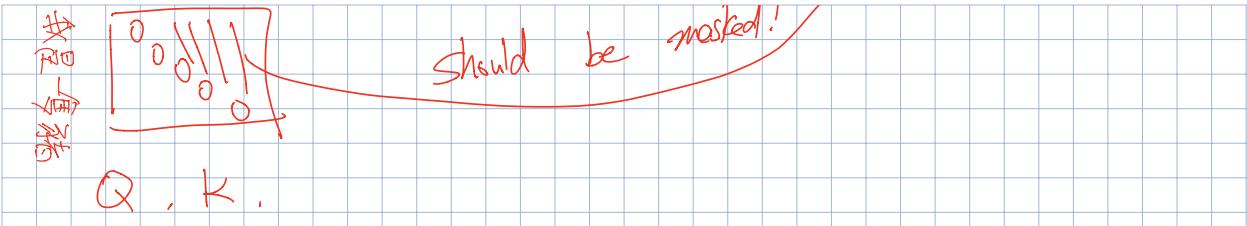


In order to prevent future information to be accounted while predicting,

We create a "mask" which sets the value of weight $\rightarrow -\infty$
(actually, -100000)

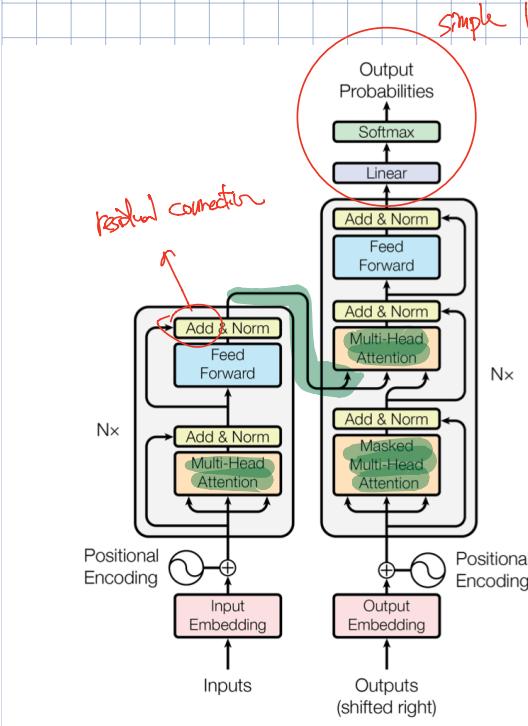
Ex: In Decoder -

我 (是) 一隻 猪. (Set to $-\infty$)
this word should not have attention with
(while doing this, it's not hard to imagine
we should set the QK^T 's $W.$ to -100000)
... \nearrow (mask upper triangle)
我 (是) 一隻 猪 !!



Note: there is another mask, we will talk later.

FFN & Norm:



Normalize before $\rightarrow \text{softmax}$
 $(\text{weights} \times \sqrt{d_{\text{model}}})$

simple linear transformation
 and softmax like output \Rightarrow another MN

$$\text{FFN}(x) = \text{Max}(0, xW_1 + b)W_2 + b_2$$

ReLU activation

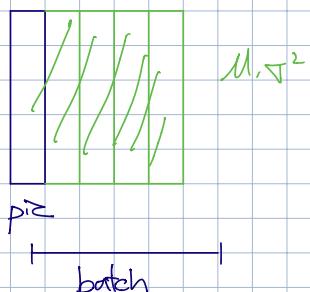
actually, it is just $2 \times \text{NN}$.

* in paper, d output = 512

$$d_{\text{inner}} = 2048$$

d ff.

Norm:



Layer Normalization:

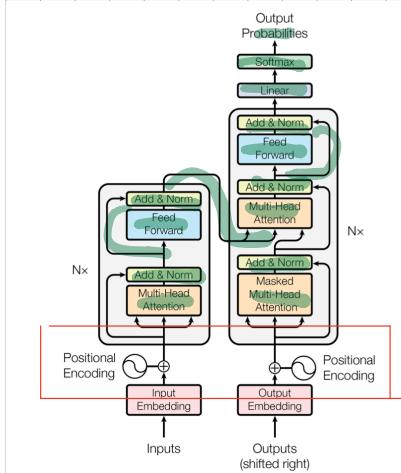
Usually use in RNN

The Last and most important idea !!

(In my view)

Positional Encoding:

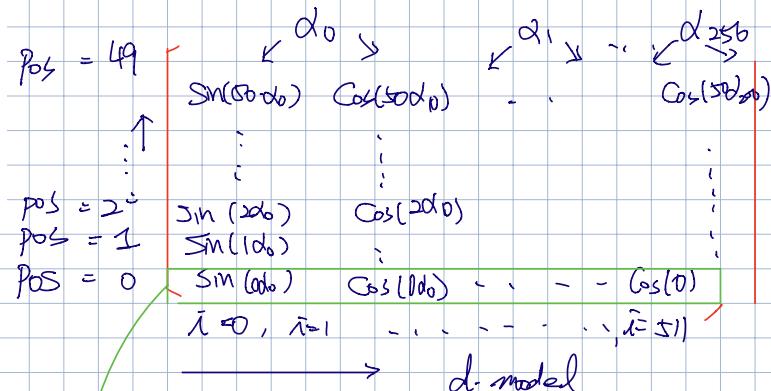
our architecture cannot tell the model
the order about the seq !!!



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

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

A matrix, let $\alpha_i = \frac{1}{10000^{2i/d_{\text{model}}}}$



→ and VR, ~~that~~ $PE(pos+k)$ can be rewritten as the linear combination of $PE(pos+k)$.

* In practice, we just create sin, cos np and concat

Why added instead of concate it? (讚嘆李大師)

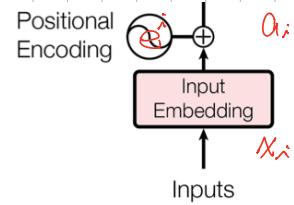
→ Suppose p_i represent one-hot position encoding.

x_i is the input

$$\begin{matrix} W^I & W^P \end{matrix} \quad \begin{matrix} x^i \\ p^i \end{matrix}$$

$$= \begin{matrix} W^I & x^i \end{matrix} \quad a^i$$

$$+ \begin{matrix} W^P & p^i \end{matrix} \quad e^i$$



↳ and the residual connection prevents this information loss !!!



Optimizer & Dropout:

We used the Adam optimizer [20] with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$. We varied the learning rate over the course of training, according to the formula:

$$lrate = d_{model}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$$

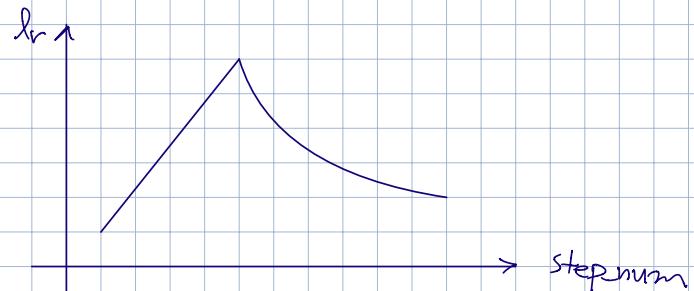
(3) num_steps

↳ linear ↑ while

24,000

start decay ↘
while
num_steps > 4000

In paper, Step_num = 4000



Dropout rate = 0.1 in every output of sublayer.

