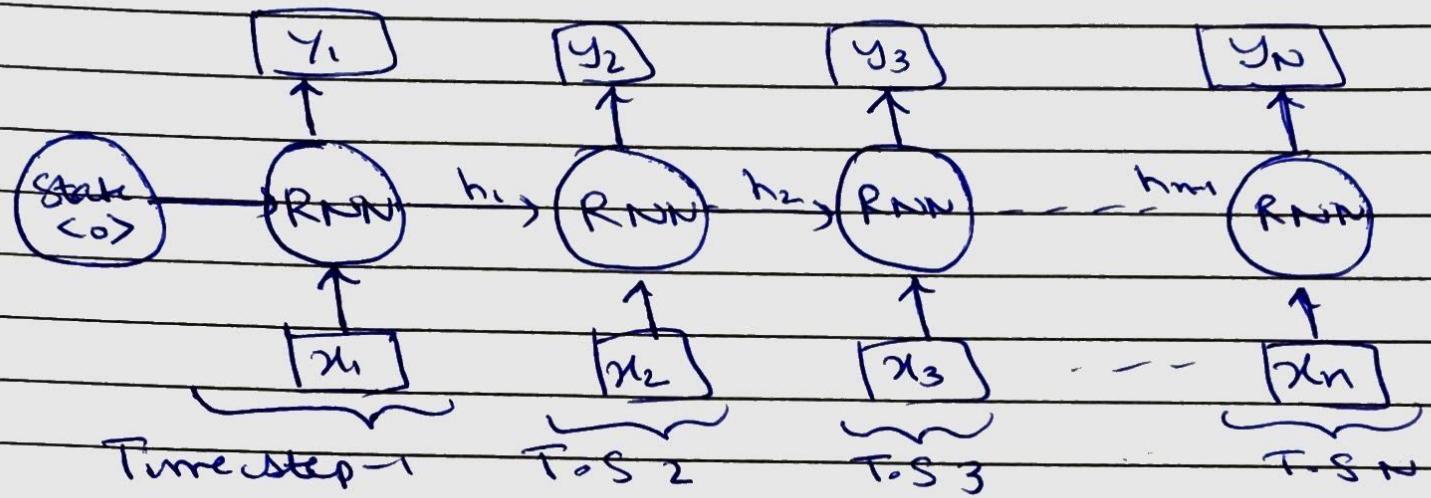# ATTENTION IS ALL YOU NEED

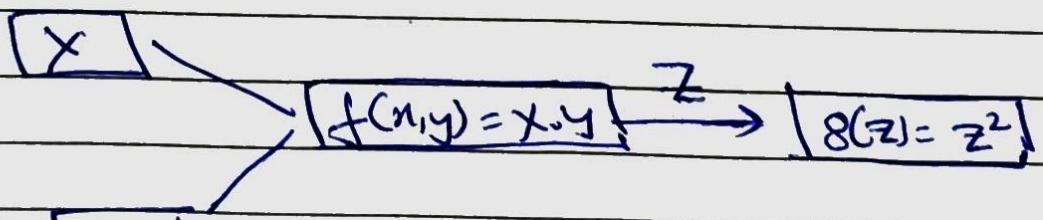## #RNNs and their problems:→



→ ~~The~~ Seq2 Seq model (map x to Y)
→ For n token, there will be n time stamps
→ Starts with 0, and each RNN has two inputs
   $x_i$ and $h_{i-1}$ (hidden state from previous
   computation)

## # Problems:→

(1) Slow computation for long sequences
    (since it is sequential)

(2) Vanishing and exploding gradients:
    eg. assume a simple graph



$$f(x,y) = x \cdot y \xrightarrow{z} g(z) = z^2$$

Here, $\dfrac{dy}{dx} = \dfrac{dy}{df} \cdot \dfrac{df}{dx}$

→ Thus, the longer the chain the more chances of vanishing / exploding gradients which is not desirable.!!

(3) **Difficulty in accessing information from long ago**
→ First token will have very less or negligible impact on the last token ( loss of context)
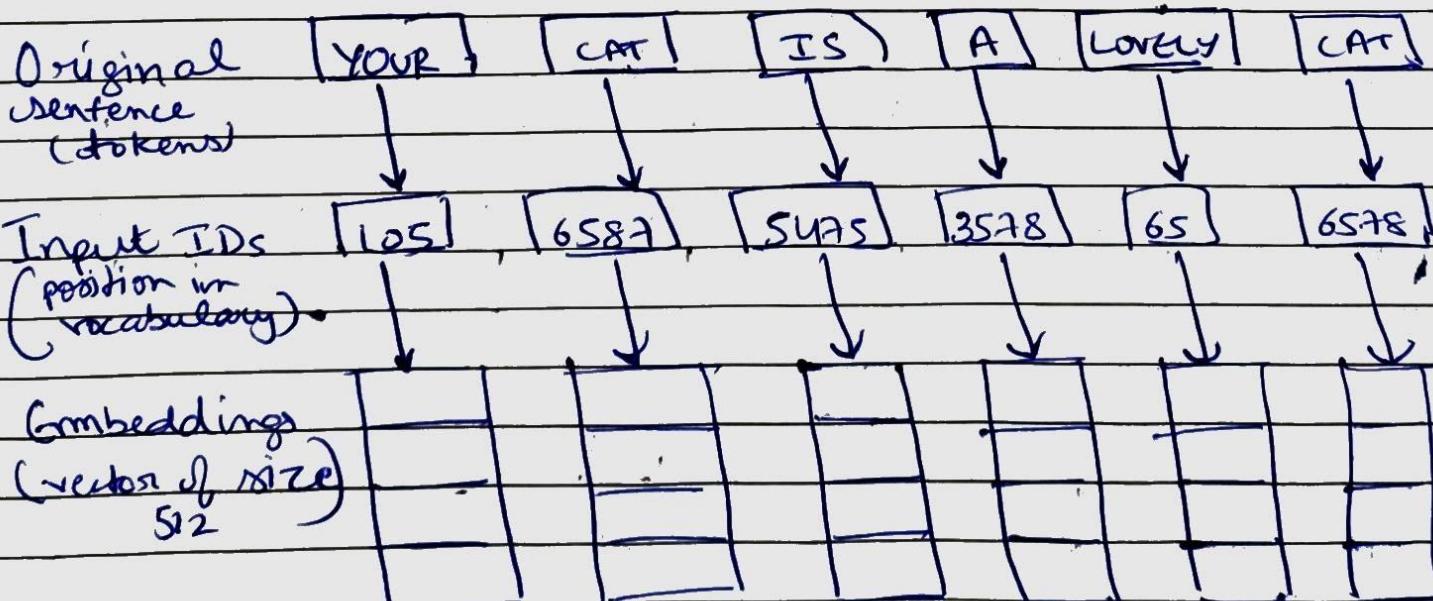→ This is undesirable!!

# ENCODER: →

**STEP 1 : Input embeddings : →**
eg. Your cat is a lovely cat

→ Tokenize it (i.e split into individual words)

→ Map these tokens into numbers which represent their position in our vocabulary

| Original sentence (tokens) | YOUR | CAT | IS | A | LOVELY | CAT |
|---|---|---|---|---|---|---|
| Input IDs (position in vocabulary) | 105 | 6587 | 5475 | 3578 | 65 | 6578 |
| Embeddings (vector of size) 512 | | | | | | |

→ However, note that these embeddings are learnable parameters and can change with time during training in order to satisfy the loss function!!

→ Also, we define $\boxed{d_{model} = 512}$ which represents the size of embedding vector of each word.

## STEP 2 : POSITIONAL ENCODINGS :→

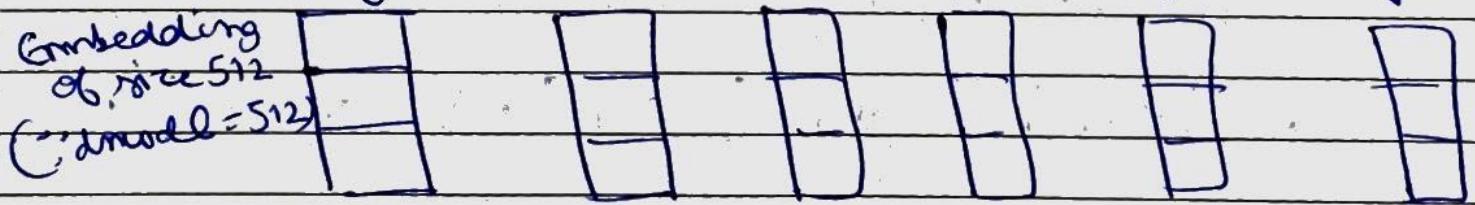→ We want each word to carry some information about it's spatial encoding in the sentence.
→ We want model to treat words that appear 'close' to each other as 'close' and those that are distant as 'distant'.
→ We want the PE to represent a pattern that can be learned by the model!! (HOW??)

eg.
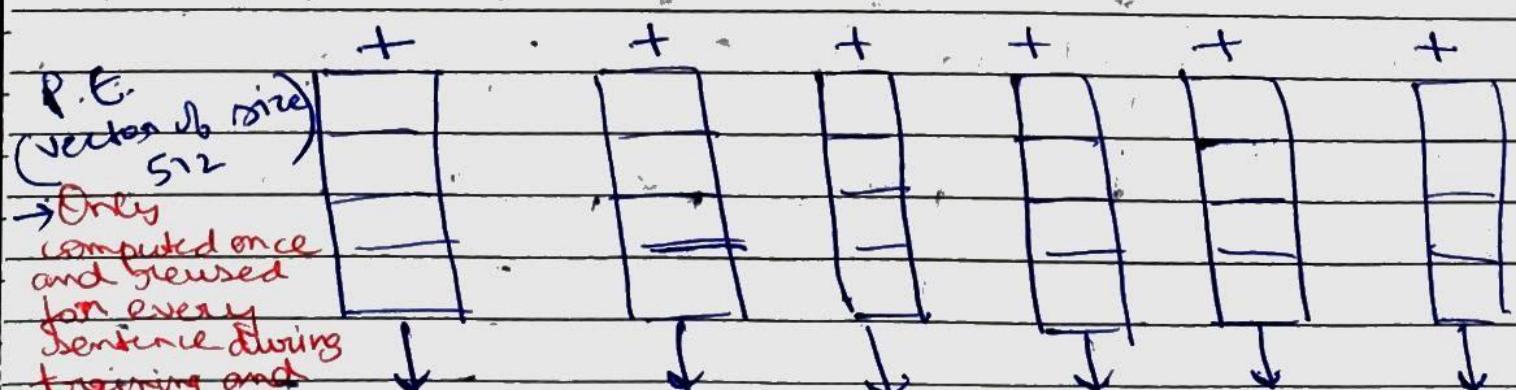original sentence | YOUR | | CAT | IS | A | LOVELY | | CAT |

Embedding of size 512 ($\therefore d_{model} = 512$)
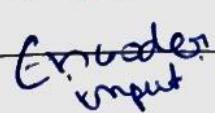
+     +     +     +     +     +

P.E. (vector of size 512)

→ Only computed once and reused for every sentence during training and inference

Encoder input

# How is PE calculated?

eg.

$$PE(pos, 2i) = \frac{\sin(pos)}{10000^{\frac{2i}{dmodel}}}$$

$$PE(pos, 2i+1) = \frac{\cos(pos)}{10000^{\frac{2i}{dmodel}}}$$

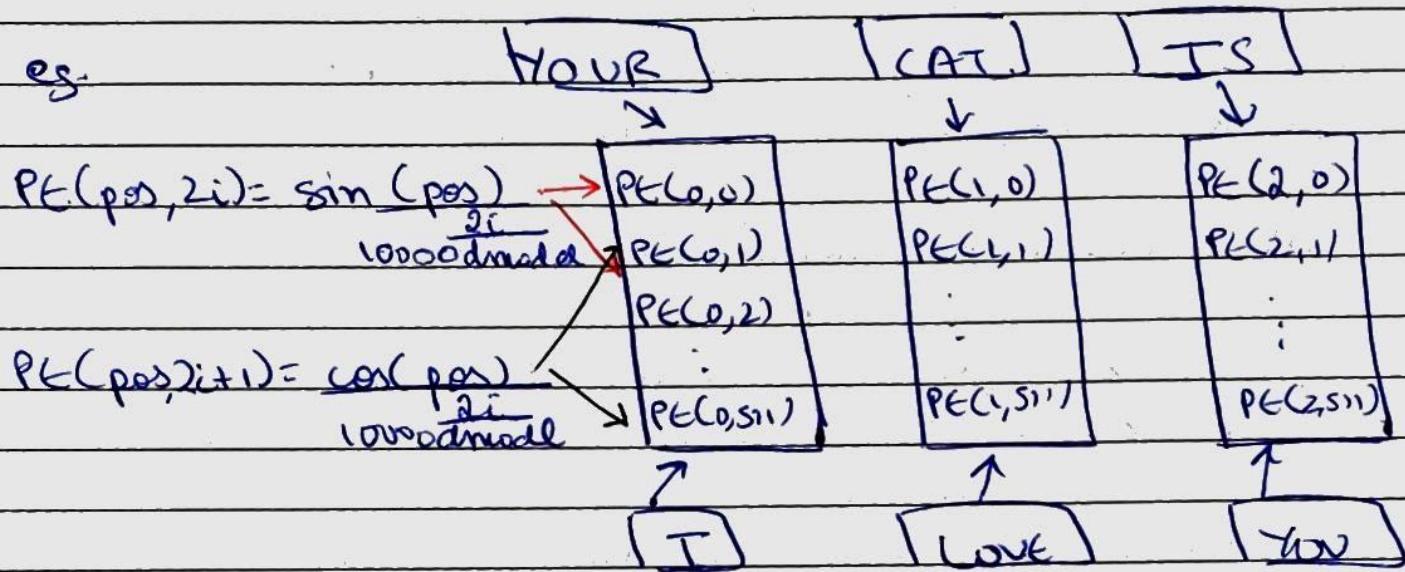| YOUR | CAT | IS |
|------|-----|-----|
| PE(0,0) | PE(1,0) | PE(2,0) |
| PE(0,1) | PE(1,1) | PE(2,1) |
| PE(0,2) | . | . |
| . | . | . |
| PE(0,511) | PE(1,511) | PE(2,511) |
| I | LOVE | YOU |

→ We can reuse same PE for diff sentences (as it is just a way of storing position of words) and is thus not learned or unique to diff. sentences.

→ Thus, they are just computed once and can be used again and again!!

Q =) Why trigonometric functions??

A =>                     like sin and cos naturally represent a pattern that the model can recognize as <u>continuous</u>, so relative positions are easier to see for the model. By watching the plot of these functions, we can also see a regular pattern, so we can hypothesize that the model will see it too.

# STEP-3     SELF-ATTENTION
## (with single head)

→ It existed before Transformers. The authors just changed it to Multi-Head!!
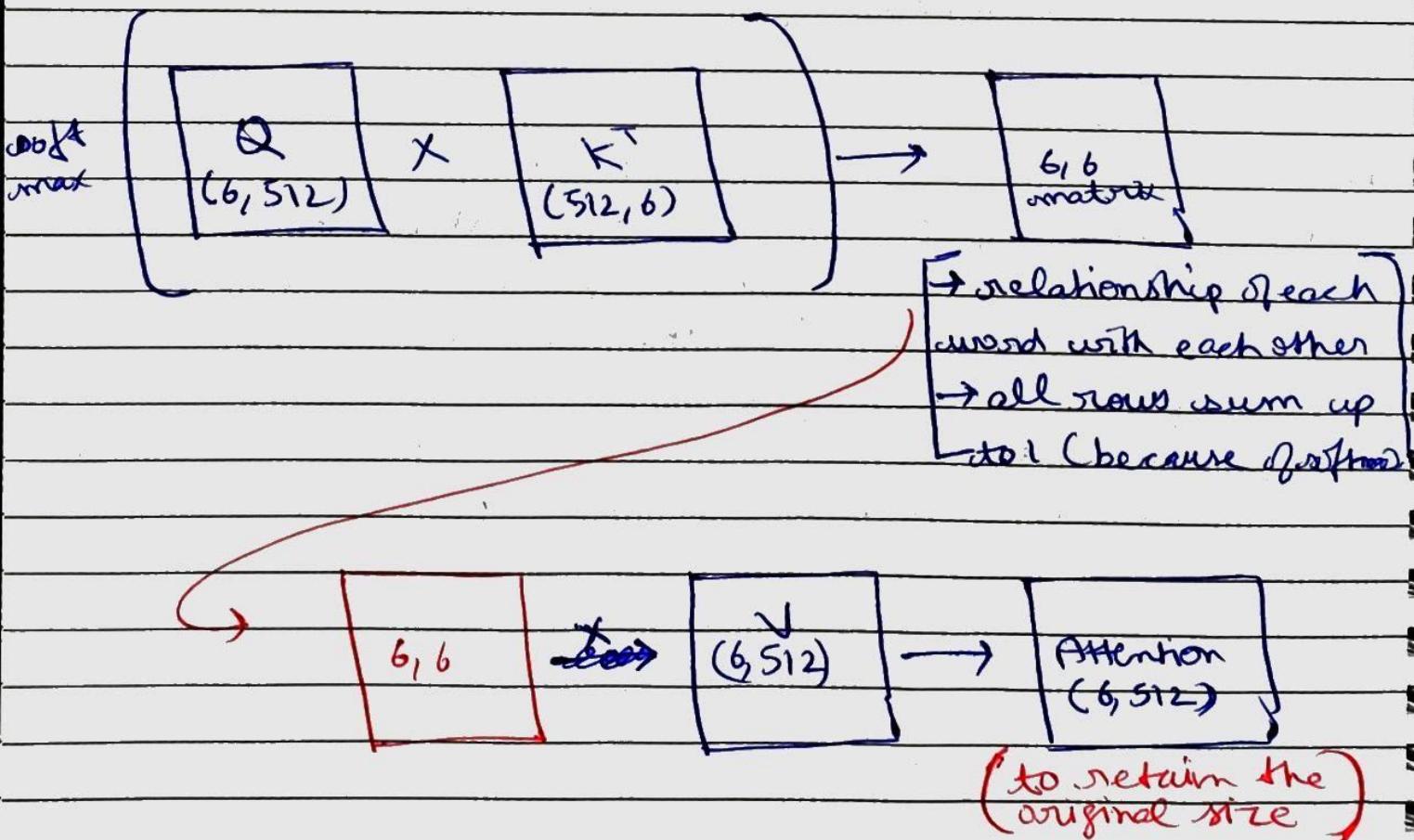
→ SA allows the model to relate words to each other

eg. IE → captures meaning of words

PE → give them info of position of words inside sentence

SA → relates words to each other.

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

Here, $Q, K, V$ are just the input sentence embeddings (after PE of course) of size $(6, 512)$
∵ our sequence is of length $= 6$, and $d_{model} = 512$



soft↑
max

| Q (6, 512) | X | $K^T$ (512, 6) | → | 6, 6 matrix |

→ relationship of each word with each other
→ all rows sum up to 1 (because of softmax)

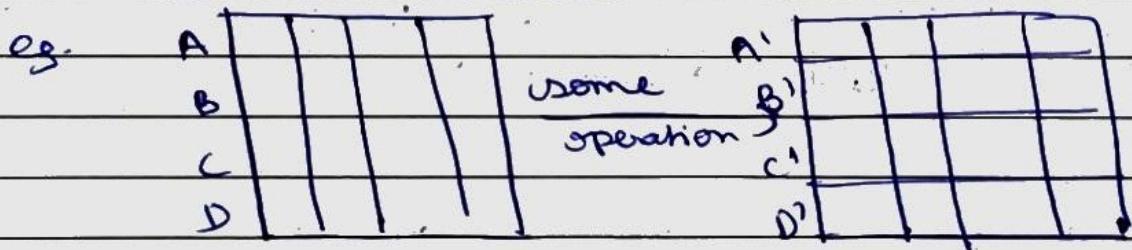| 6, 6 | ~~X~~ | V (6, 512) | → | Attention (6, 512) |

(to retain the original size)

→ Each row in this new attention matrix captures not only the meaning (given by the embeddings) or their position in sentence (given by PE) but now also represents each words interactions / relationship with other words!!

# Properties of self-attention:→

(1) Self attention is permutation-invariant (i.e swapping orders does not change the computed value)

→ This is desirable

eg.



if we swap B & C, it won't change the calculated value of B' and C'. It will just change their position in the matrix.



(2) It requires no parameters as of now (will change later when we discuss multi head attention)

(3) Values along diagonals are expected to be highest
  (i.e a word will have strongest relationship with itself)

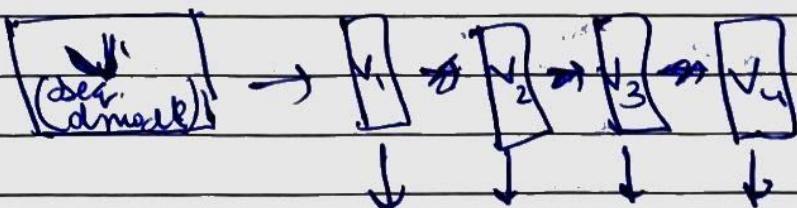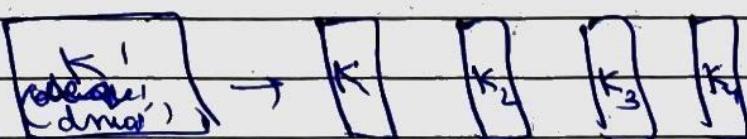(iv) If we don't want some positions to interact, we can always set their ~~position~~ value to $-\infty$ before applying softmax operation and the model will not learn those interactions? ($\because e^{-\infty}=0$). This is desirable when we study decoders!!

## STEP-4: MULTI-HEAD ATTENTION:

$$\underset{(seq, dmo)}{Q} \times \underset{(dmodel, dm)}{W^Q} = Q'$$

$$\underset{(seq, dmo)}{K} \times \underset{(dmo, dmod)}{W^K} = K'$$

Input (seq, dmodel)

$$\underset{(seq, dmod)}{V} \times \underset{(dmodel, dmod)}{V^K} = V'$$

$Q'$ (seq, dmo) $\rightarrow$ $Q_1$ $Q_2$ $Q_3$ $Q_4$ } seq

dmodel

$d_K$

$K'$ (seq, dmo) $\rightarrow$ $K_1$ $K_2$ $K_3$ $K_4$

$V'$ (seq, dmodel) $\rightarrow$ $V_1$ $V_2$ $V_3$ $V_4$

add

MHA

PE embed

I/P

$$\text{seq} \left\{ \begin{array}{|c|} \hline H\\E\\A\\D\\1 \hline \end{array} \right. \quad \begin{array}{|c|} \hline H\\E\\A\\D2 \hline \end{array} \quad \begin{array}{|c|} \hline H\\E\\A\\D3 \hline \end{array} \quad \begin{array}{|c|} \hline H\\E\\A\\D4 \hline \end{array} \xrightarrow{\text{concat}} \boxed{H} \times \boxed{W^0}$$

$\underbrace{\phantom{xxxxxxx}}_{dv}$

$(\text{seq}, \cancel{\phantom{x}} h \times d_v)$   $(h \times d_v, d_{model})$

$\underbrace{\phantom{xxxxxxxxxxxxxx}}_{d_{model}}$

$$= \boxed{\begin{array}{c} MHA \\ (\text{seq}, d_{model}) \end{array}}$$

→ Here, initially 3 copies of the input embeddings are made as Q, K, V (as previous)

→ These Q, K, V are multiplied by parameterized $W^Q$, $W^K$, $W^V$ to form $Q'$, $K'$ and $V'$

→ These $Q'$, $K'$ and $V'$ are split into multiple heads and then attention is performed on each head as per usual.

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

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

→ These heads are then concatenated to form H matrix and then multiplied by $W^0$ to give our final output.

$$\text{Multihead} (Q, K, V) = \text{concat}(\text{head}_1, \text{head}_2, \text{hed} \dots h)W^0$$

**Q4) But why split ???**

→ Here, each head$_i$ is watching the full sentence but a different aspect of ~~each word~~ the embedding of each word.

→ This is because, we want each head to watch diff. aspect of the same word.

eg. In chinese, a word can be a noun, verb, adverb, etc depending on the context.

→ So, what we want is that maybe one head relates it as a noun, another as verb and so on.. i-e, we want to look at all diff possibilities.

→ This attention can also be visualized !"

**Q⁵) Why Q, K, V ?**

⇒ Simple, it comes from python dictionary !"

**STEP-5 : ADD & NORM. ? →**

→ But, first we need to learn about layer normalisation.
→ Consider, three items in a batch.



ITEM 1          2          3

$\mu_1$         $\mu_2$      $\mu_3$
$\sigma_1^2$    $\sigma_2^2$  $\sigma_3^2$

→ We calculate mean & variance for each item separately and replace each value in item by

$$x_i = \frac{x - \mu_i}{\sqrt{\sigma_i^2 + \varepsilon}}$$

→ We also introduce two parameters, gamma and beta, that introduce some fluctuations in the data, because maybe having all values between 0 and 1 may be too restrictive for model. The network will learn to tune these two parameters to introduce fluctuations when necessary.

# DECODER :→
→ We have output embeddings (similar to IE in Encoder) and PE (similar to Encoder)
→ Next, we have, masked multihead attention
→ The K and V come from encoder wheras the query comes from decoder. Thus, it is no longer a self-attention but a cross-attention!!

STEP-1 · MASKED MULTIHEAD ATTENTION :→

→ Our goal is to make the model causal; ie; it means the output at a certain position can only depend on the words on the previous positions.
→ The model must not be able to see future words. (How??)
→ We do this by simply replacing the future words with $-\infty$.

| | YOUR | CAT | IS | A | LOVELY | CAT |
|---|---|---|---|---|---|---|
| YOUR | | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ |
| CAT | | | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ |
| IS | | | | $-\infty$ | $-\infty$ | $-\infty$ |
| A | | | | | $-\infty$ | $-\infty$ |
| LOVELY | | | | | | $-\infty$ |
| CAT | | | | | | |

# #TRAINING:→

→ We will consider a translation task.

    eg. I love you very much   (ENG)

$$\downarrow$$

      Ti Amo Malto        (ITALIAN)

→ We add special tokens to our sentence.



Label: Ti amo motto &lt;EOS&gt;

→ cross entropy loss

Softmax
(seq, vocab, size)

Linear
(seq, dmod) → (seq, vcab)

Decoder output
(seq, dmodel)

Encoder       K    V   → Decoder

                  ↑ Q

Encoder input
(seq, dmodel)

Decoder input
(seq, dmodel)

↑ IE + PE                ↑ OE + PE

&lt;SOS&gt; I love you very much &lt;EOS&gt;    &lt;SOS&gt; Ti amo molto
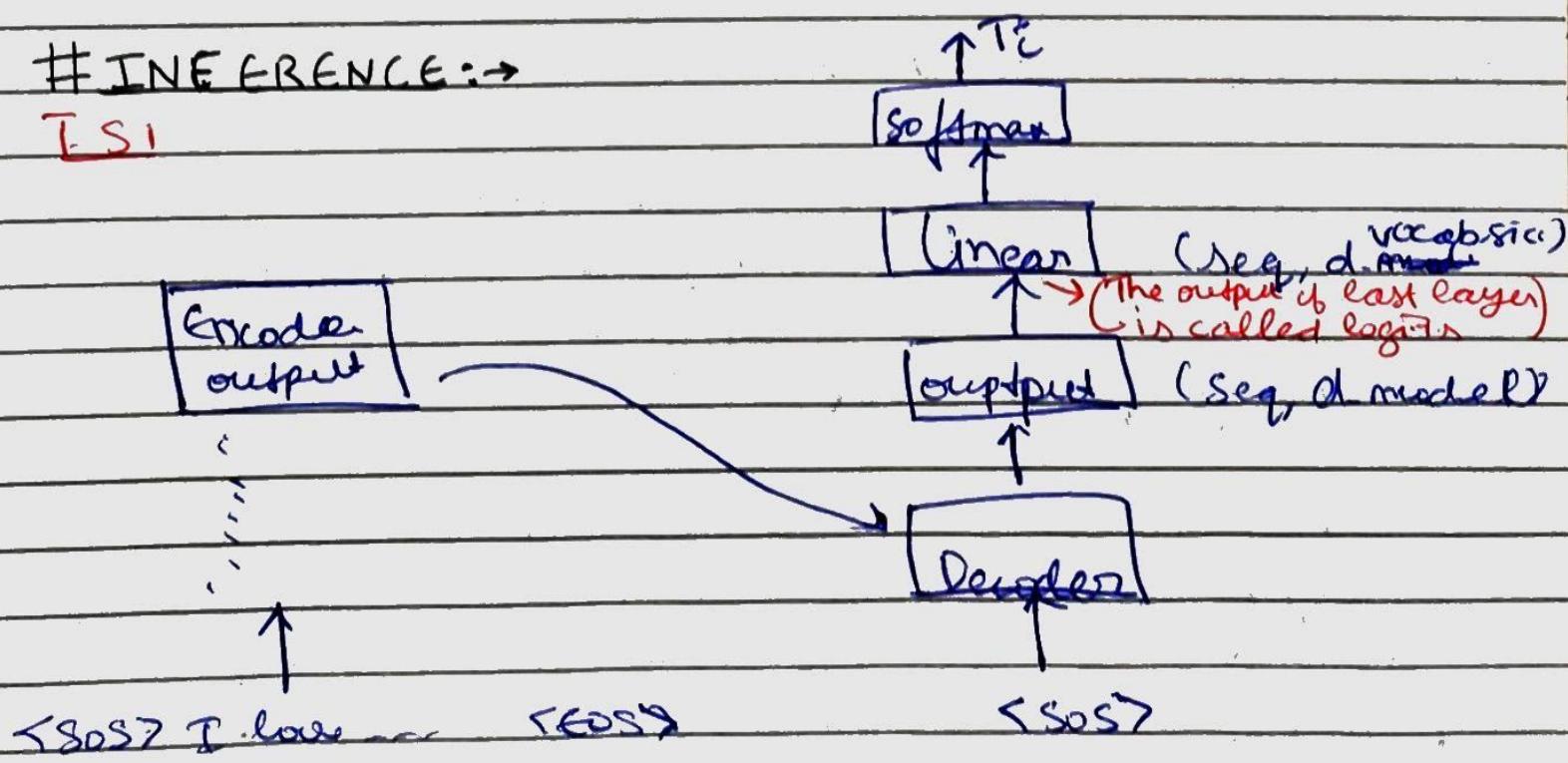
                                      (shifted right) by SOS

→ Next, we do JE + PE to create input embeddings for our encoder. IIly for decoder.

→ We get K and V from encoder and use Q from decoder (as it is masked) to get our decoder output

→ This output needs to be mapped to our vocabulary with a linear layer

→ Finally, we take the softmax to get our word predictions.

→ Once, we have our predicted label and the true label, we use cross entropy loss to update the weights.

→ This all happens in 1 time step (Time step = 1)
i.e all in one pass.

# INFERENCE :→
T.S1



↑ $T_i^c$
softmax
Linear    (seq, d.model → vocab_size)
          → (The output of last layer is called logits)
Encoder output    output    (seq, d.model)
Decoder

⟨SOS⟩ I love ...    ⟨EOS⟩          ⟨SOS⟩

→ We start with <SOS> and compute the logits as per decoder model.
→ We then perform softmax on these logits to get the next word which is Ti

Time step 2 :→
    <SOS> Ti
→ We don't compute encoder part again as our English sentence hasn't changed.
→ However, for decoder we go from
    <SOS> Ti    ⟶   Amo

TS 3 :→
    <SOS> Ti  ⟶  Molto

TS 4 :→
    <SOS> Ti Amo Molto ⟶ <EOS>

Thus; it is token by token until we generate <EOS>!!

# #INFERENCE STRATEGY :→

→ We selected, at every step, the word with the maximum probability / softmax value. This strategy is called greedy and usually does not perform well.

→ A better strategy is to use top B words and evaluate all possible next words for each of them and at each step, keeping the top B most probable sequences. This is Beam Search and generally performs better!!