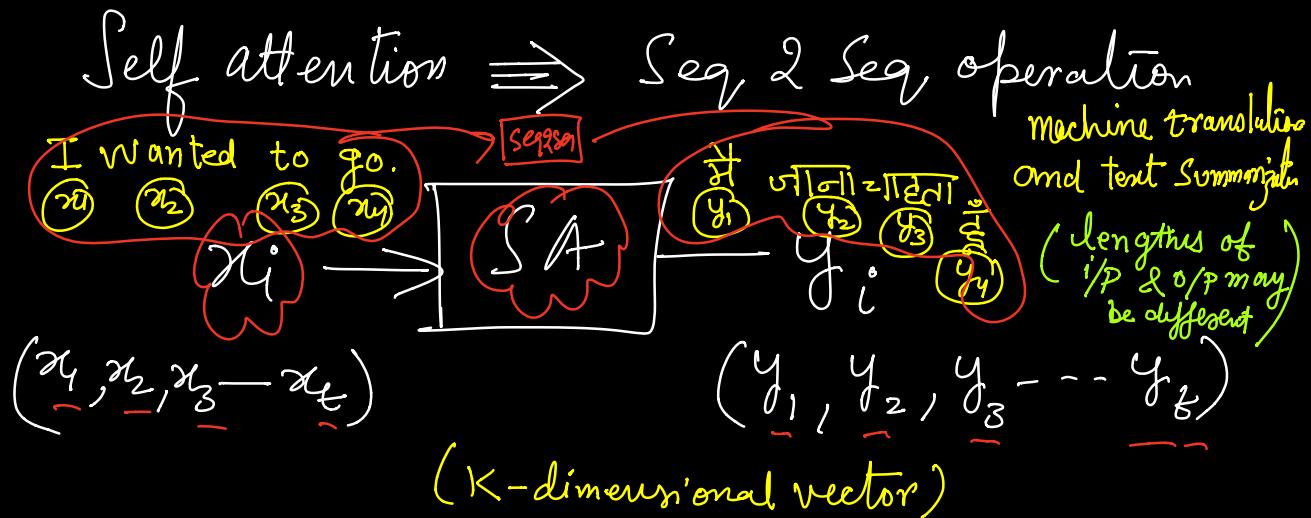


Transformer N/L

• www.peterbloem.nl/blog/transfomers



$$x_i \in K\text{-dim vec} \rightarrow \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} & \dots & x_{1K} \\ x_{21} & x_{22} & x_{23} & x_{24} & \dots & x_{2K} \\ x_{31} & x_{32} & x_{33} & x_{34} & \dots & x_{3K} \\ \vdots & \vdots & \vdots & \vdots & & \vdots \\ x_{t1} & x_{t2} & x_{t3} & x_{t4} & \dots & x_{tK} \end{pmatrix}$$

\times (t \times K)
rows \rightarrow columns

Any $y_i = \sum_{j=1}^t w_{ij} x_j$ (Similarly we have Y)

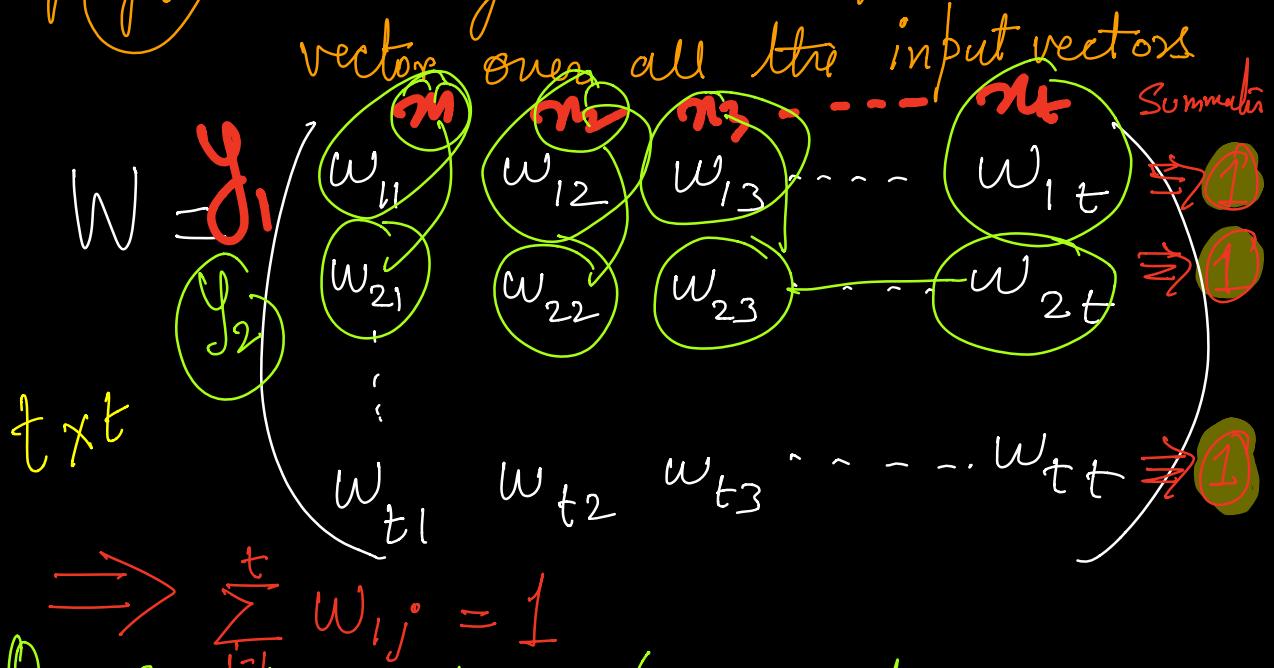
Matrix weighted average over all input vector

$$\therefore y_i = w_{i1} x_1 + w_{i2} x_2 + w_{i3} x_3 + \dots + w_{it} x_t$$

$\forall i = 1 \dots t$ Scales vector Vector addition

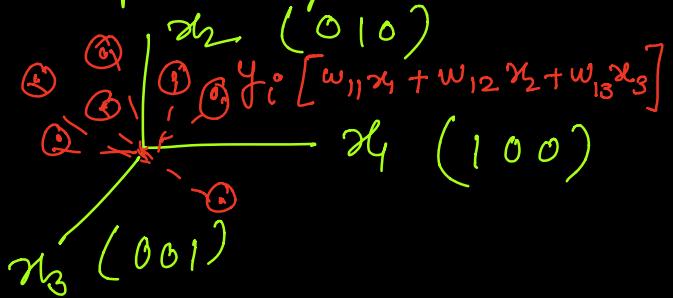
• Just like X & Y matrix
 \therefore there is a W matrix.

Any y_i is nothing but the weighted averaged



$$\Rightarrow \sum_{j=1}^t w_{ij} = 1$$

One can say that x_i 's can be seen as the basis and the space spanned by them will hold the y_i 's



This weight matrix $W = (w_{ij})$ is $t \times t$

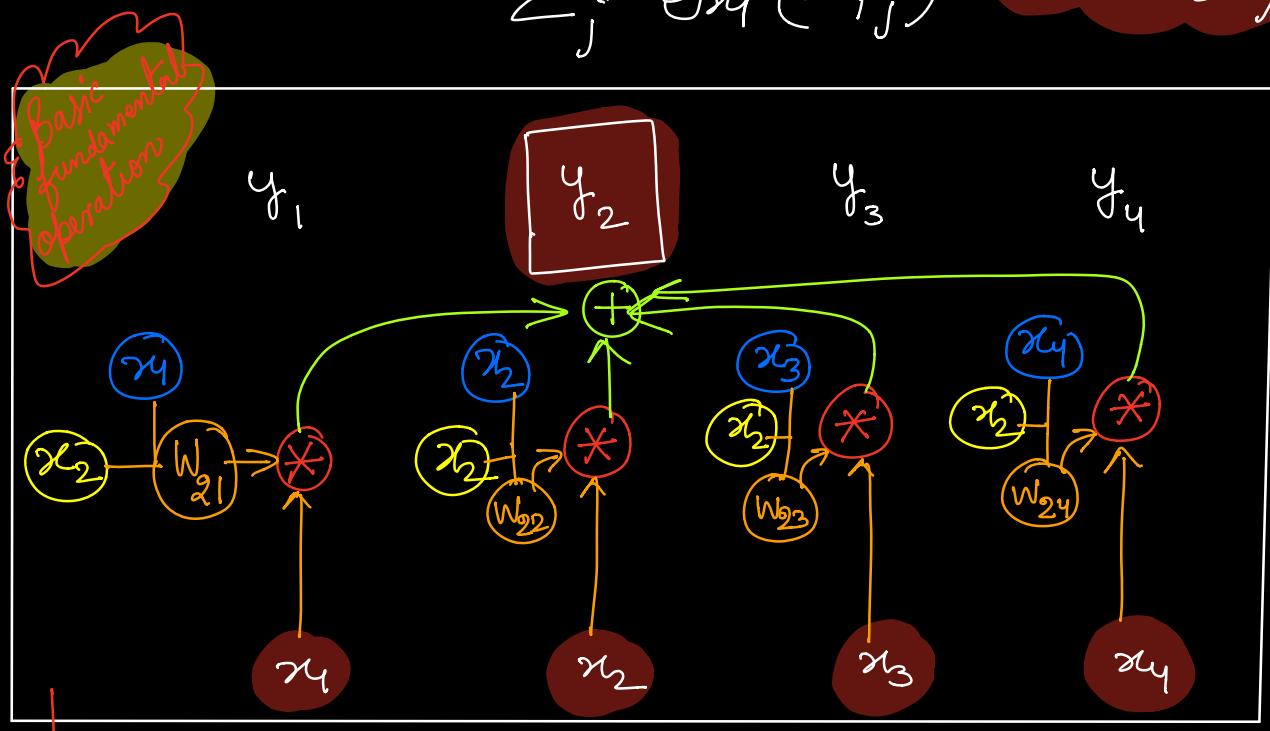
$$w_{ij} = f(x_i^{\circ}, x_j^{\circ}) = \frac{x_i^{\circ T} x_j^{\circ}}{\|x_i^{\circ}\| \|x_j^{\circ}\|}$$

Note: How related 2 vectors are. This is cosine similarity

Its value can range from $(-\infty, \infty)$
 hence in order to normalize it apply Softmax.

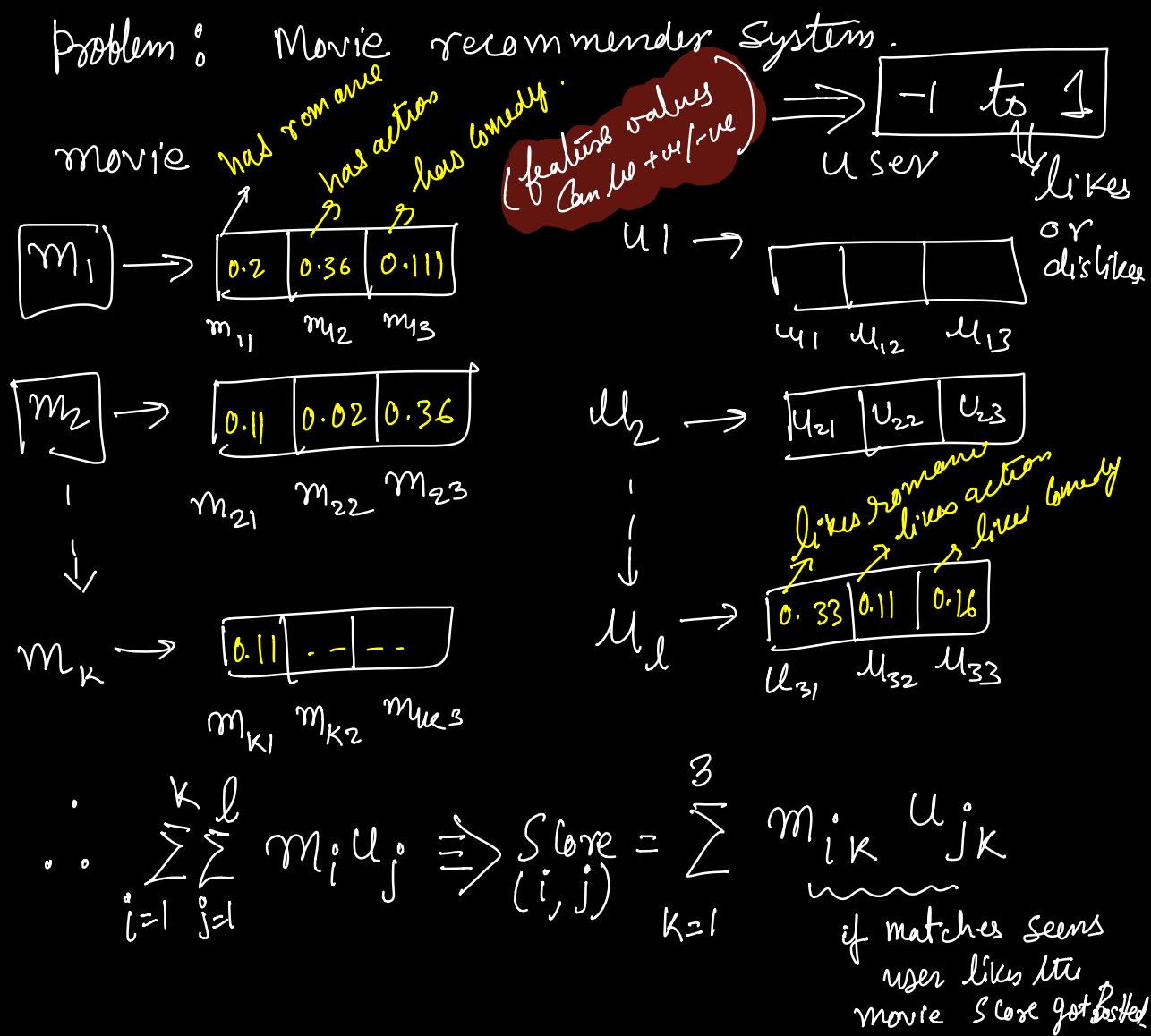
$$w_{ij}^* = \frac{\exp(w_{ij}^*)}{\sum_j \exp(w_{ij}^*)}$$

Row-wise normalization



Only operation in transformers that propagates information between the vectors. Other than self attention all operations are applied on individual inputs.

Why do self attention works ?



Basically each movie can be seen as a feature vector as well as each user " " " " " of same dim. dot product \Rightarrow cosine similarity is being computed.

\Rightarrow But such feature annotations need to be done manually as for million of movies and atleast extraction of K features \Rightarrow [Impractical]

Instead of manual annotations we will let it to be learned automatically. fix up the hyperparameters and learn these parameters.

\hookrightarrow Movie & user \Rightarrow feature need to be learned

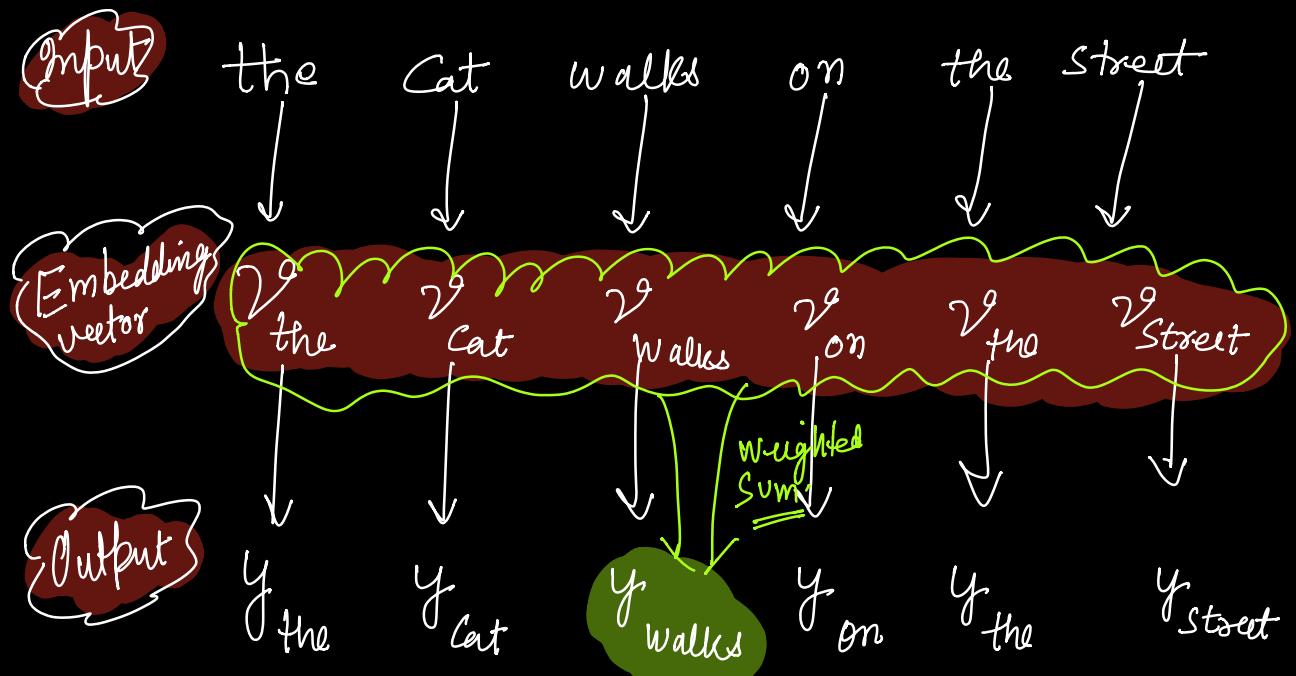
(*) taking it towards Supervised learning

(*) Get a small dataset with the class/annotation and learn the feature that can be seen as the best representation.

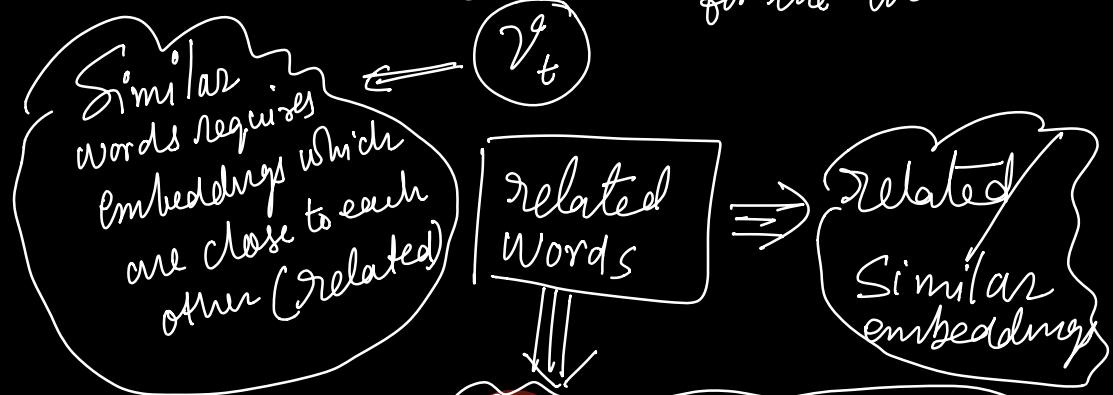
Learning the embedding

Can also be done using basic matrix factorization

Converting into a supervised optimization problem
Self attention Basic



\Rightarrow We need to learn $[v_{word}]$ embedding vector for the words



The problem in hand

But relatedness completely depends upon the optimization taste

Relatedness captured using dot product

O/p is weighted sums over the whole input sequence with weights determined by these

Dot product \Rightarrow How much these two vectors (x_1 & x_2)
 $(x_1 \cdot x_2)$ are related

Currently we are trying to do [seq 2 seq] operations

\hookrightarrow No parameters learned for Self attention

Only word encoding are determined

till now and are used to get the
"relatedness" \Rightarrow using dot product.

\hookrightarrow Self attention sees its input as a set
"NOT" \Rightarrow Sequence. (even though it is
sequence)

$[x_1 \ x_2 \ x_3 \ x_4 \ x_5]$ \rightarrow all belongs to
the input set (I)

This means that order
does not matter and

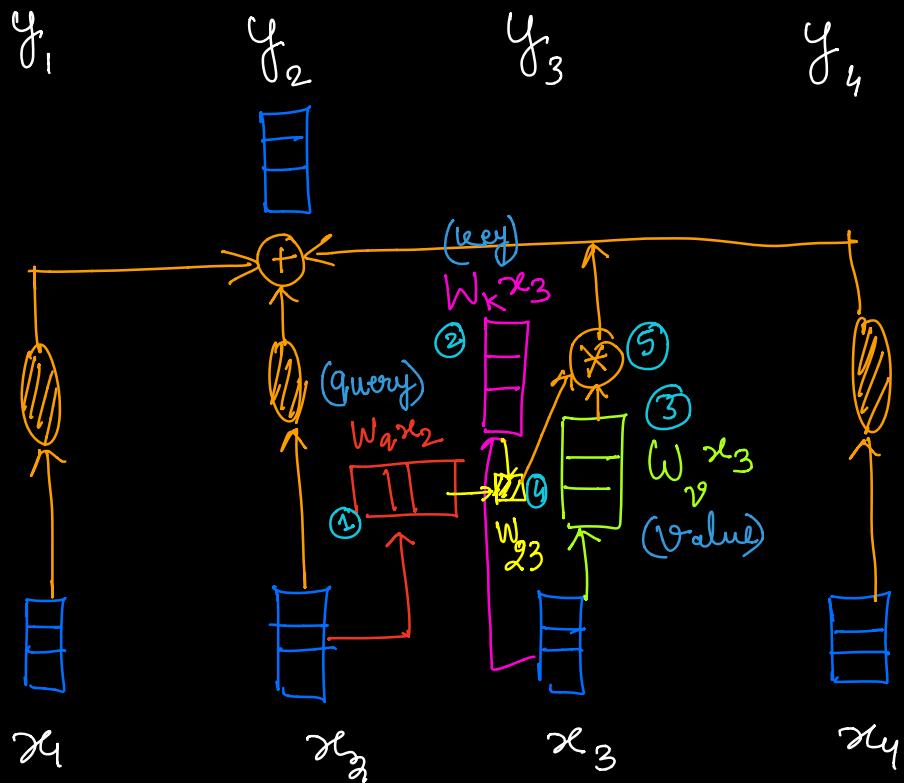
any permuted input sequence
will also generate the same
O/P.

Permutation
Equivariant

This view is addressed
in full transfer
Now introducing additional weights
to be learnt for performance (modern transformers) and
tricks

A) Query, Keys and values (Trick-01)

One can see that in the above transformer network any input vector (x_i) is used in ③ ways.



\Rightarrow Weighted sum of $[x_1, x_2, x_3, x_4]$
 but weights (w_i) for each (x_i) depends upon
 how much it is related (dot product) to x_2 .

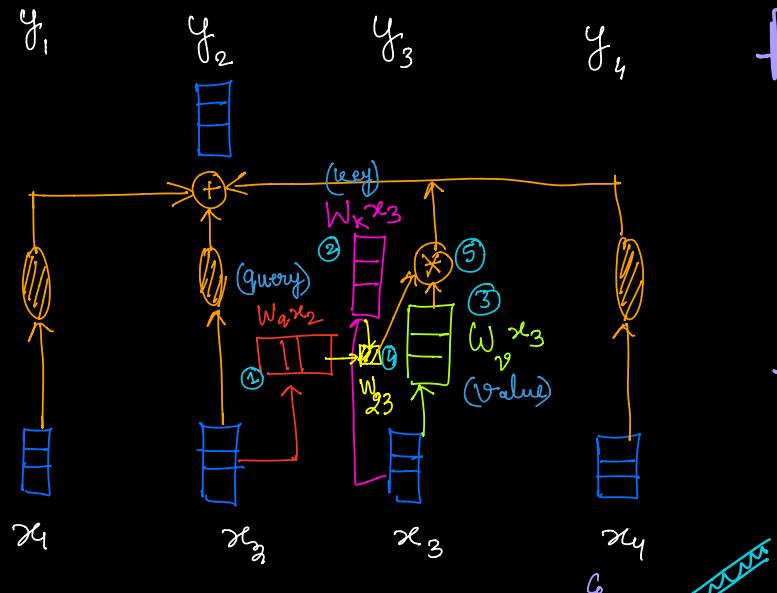
One can see that each x_i (Say x_3) is used

in ③ different ways in self attention -



* Key \div In order to participate in the generation of all $[y_{j \neq i}]$, again $[x_i]$ is used as a key and compared with all $[x_{j \neq i}]$ to get its weight contribution.

* Value \div Finally $[x_i]$ is also converted into a value vector that can be multiplied with the weight to compute each o/p weight.



For $(\mathcal{Y},)$ Computation

a) x_4 query

for relatedness
using $[w_q x_i]$

With the key of
all the inputs
i.e $\begin{bmatrix} W_k \\ x_i \end{bmatrix}$ it is

b) Now once all the related-ness of (x_i) with all other $(x_i's)$ got computed in
 $[W_{1,j}] \forall j=1-4$

[α_j] used
Converted into key and
multiplied to

$$\textcircled{O} [W_{i,j}] \quad \forall j=1 \text{ to } 4$$

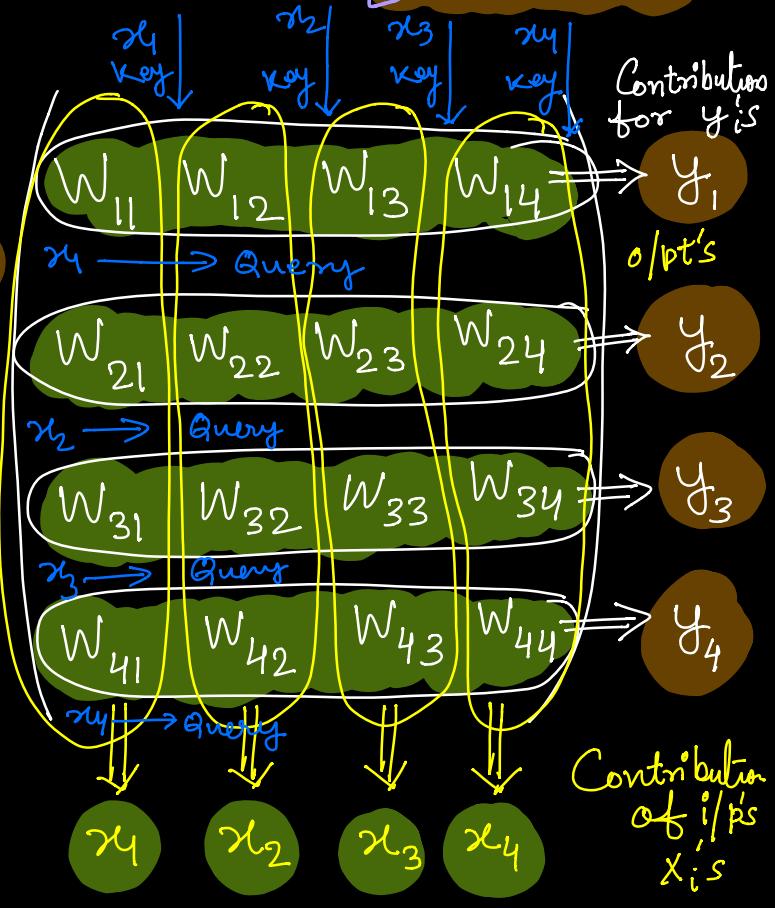
where

$$W_{ii} \leftarrow \text{Contributions of } x_i \text{ in } y_i \rightarrow y_1$$

$[W_1]$ of x_1 into y_2

$$[W_3] \xrightarrow{\text{distrinv}} y_3$$

[W₄₁] of 24 in \bar{g}_4



\therefore Horizontally $\Rightarrow \text{Row}(i) \Rightarrow \text{Query of } [x_i]$

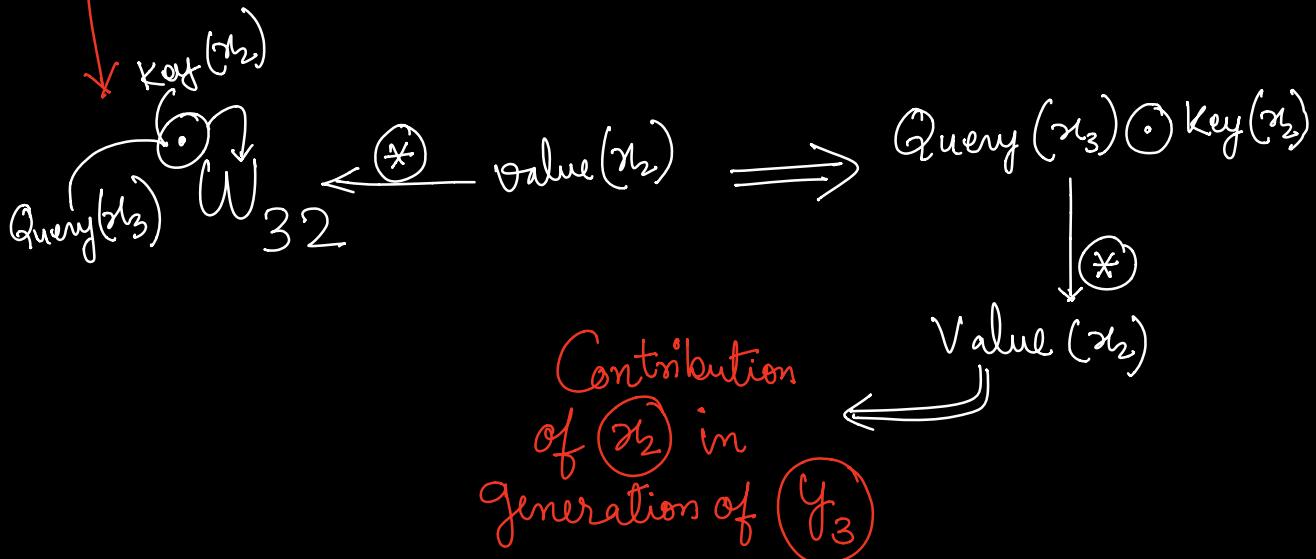
Vertically $\Rightarrow \text{Col}(j) \Rightarrow \text{Key of } [x_j]$

$$y_1 = W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + W_{14}x_4$$

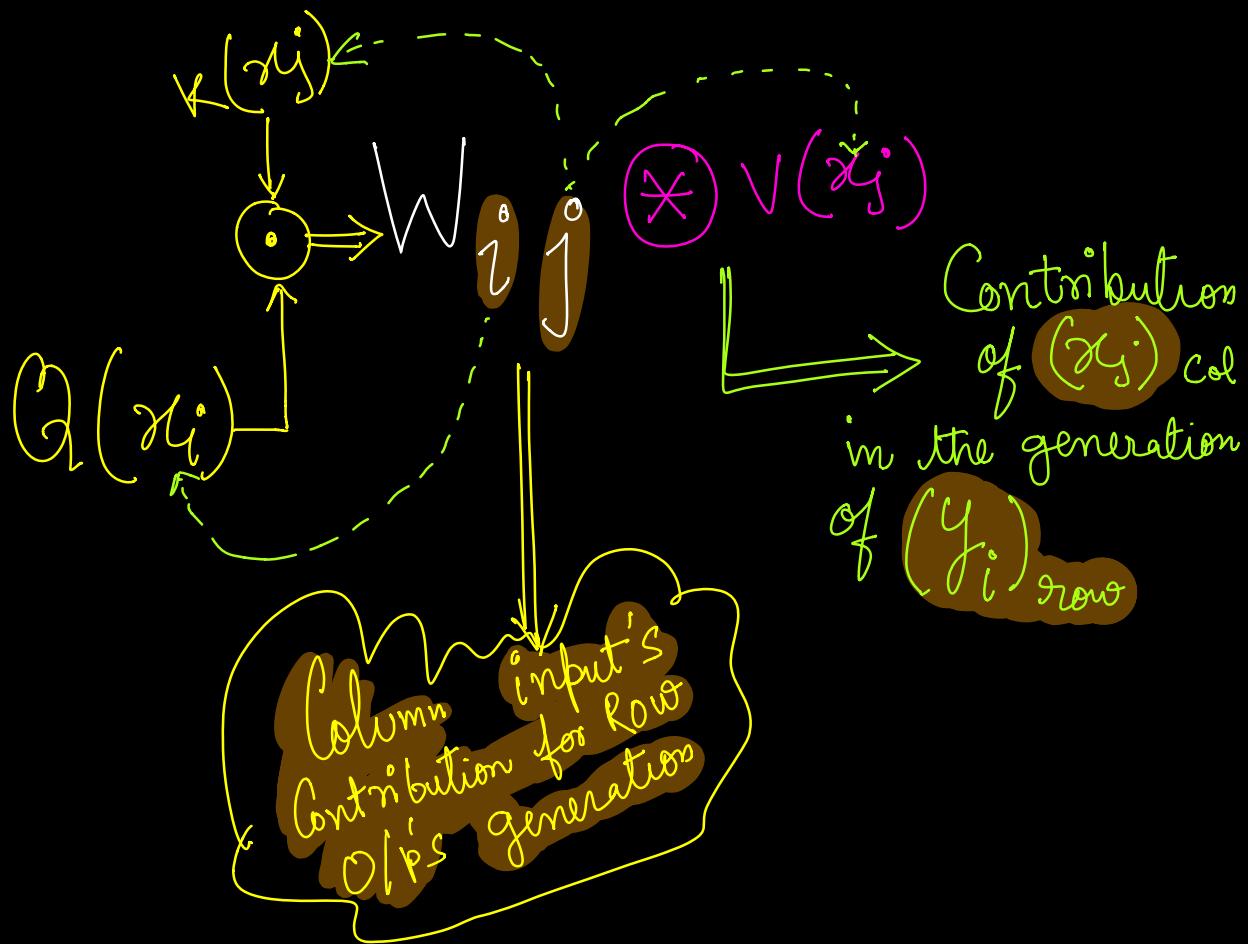
$$y_2 = W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + W_{24}x_4$$

$$y_3 = W_{31}x_1 + W_{32}x_2 + W_{33}x_3 + W_{34}x_4$$

$$y_4 = W_{41}x_1 + W_{42}x_2 + W_{43}x_3 + W_{44}x_4$$



In general :



These W_q, W_k, W_v are all $K \times K$ matrix

and can be seen as a linear transformation of
and K -dim (x_i) input vector from \mathbb{R}^K to some
other K -dim transformed vector. These 3 transformation
need to be learned in order to transform any input
 $[x_i]$, so that it can adapt to play its ③ different
roles appropriately for SA modules.

These $[W_q, W_k, W_v]$ can be used in many different ways to realize the attention. One such way is

$$\textcircled{1} \quad \text{Query}(x_i) = W_q x_i \quad \left\| \begin{array}{l} y_i = F(x_1, x_2, \dots, x_n) \\ \text{all input} \end{array} \right.$$

$$\textcircled{2} \quad \text{Key}(x_i) = W_k x_i$$

$$\textcircled{3} \quad \text{Value}(x_i) = W_v x_i$$

$$\textcircled{4} \quad \sum_{j=1}^n w'_{ij} = \text{Query}(x_i)^T * \text{Key}(x_j)$$

$$\textcircled{5} \quad \sum_{j=1}^n w_{ij} = \text{Softmax}(w'_{ij})$$

$$\textcircled{6} \quad y_i = \sum_{j=1}^n w_{ij} * \text{Value}(x_j)$$

The self attention layer got parameterized using (W_q, W_k, W_v)

(B) Scaling the Dot product (Trick-02)

$$\forall_{i,j} w'_{ij} = \frac{\text{Query}(x_i)^T * \text{Key}(x_j)}{\sqrt{K}}$$

Average value of the dot product keeps on growing with the embeddings dimensions $[K]$. As finally

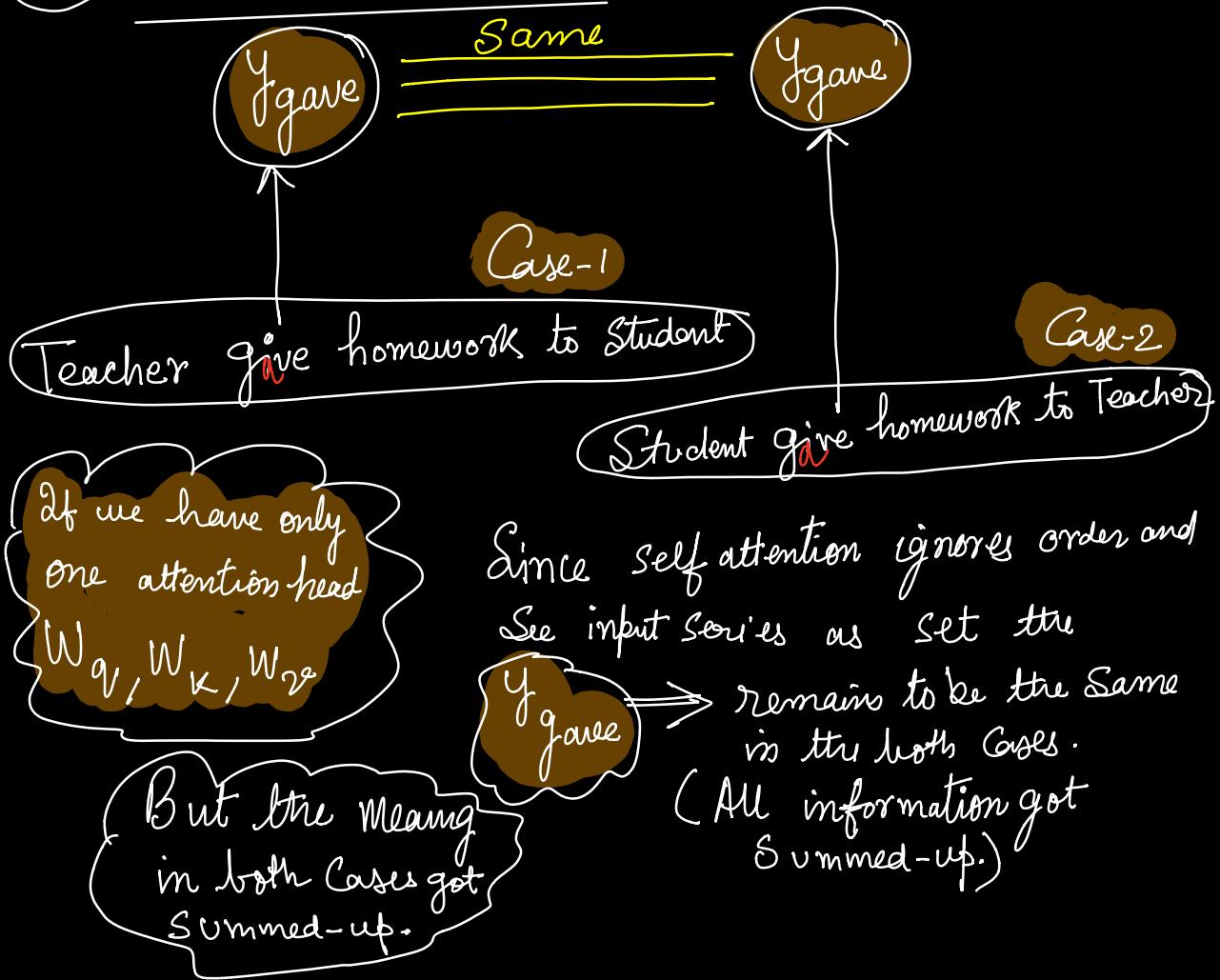
$[w'_{ij}] \Rightarrow$ got applied by Softmax this growth got exponentiated and may be effected adversely.

Hence we need to normalize it. Since our input vector $\in \mathbb{R}^K$ (Say all are $[C] \rightarrow \begin{bmatrix} c \\ c \\ \vdots \\ c \end{bmatrix}$)

$$\therefore C \xrightarrow[\text{(Growth is } \sqrt{K})]{\text{K-dim}} \sqrt{K} * C \quad \text{Norm} \Rightarrow [\sqrt{K} * c]$$

Hence normalizing it by (\sqrt{K}) .

C) Multi-Head attention



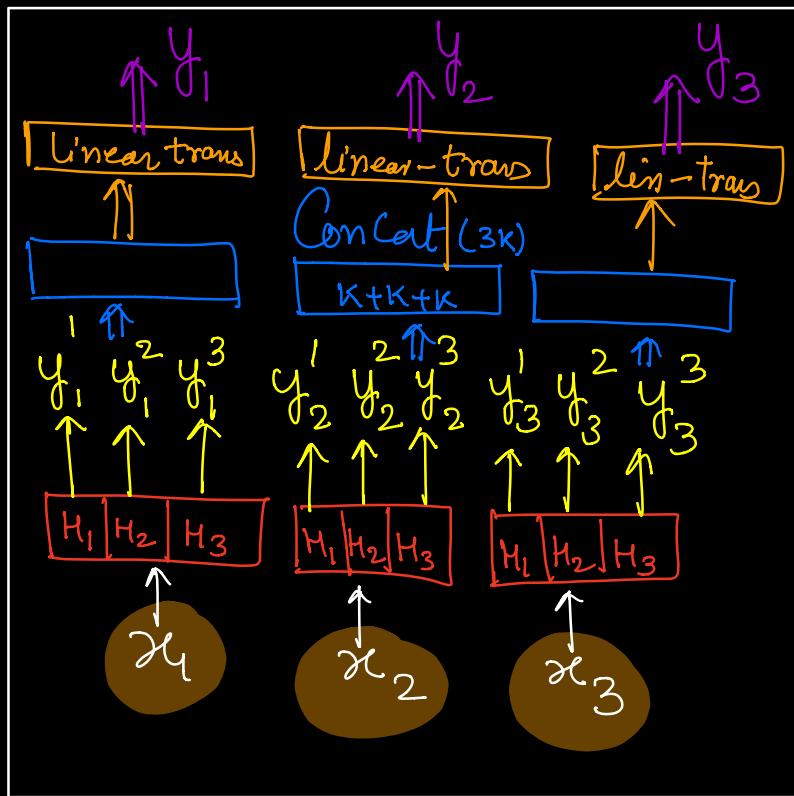
Same word mean differently to different neighbourhood/Context.

One Self-attention head \rightarrow information got Summed up (Averaged)

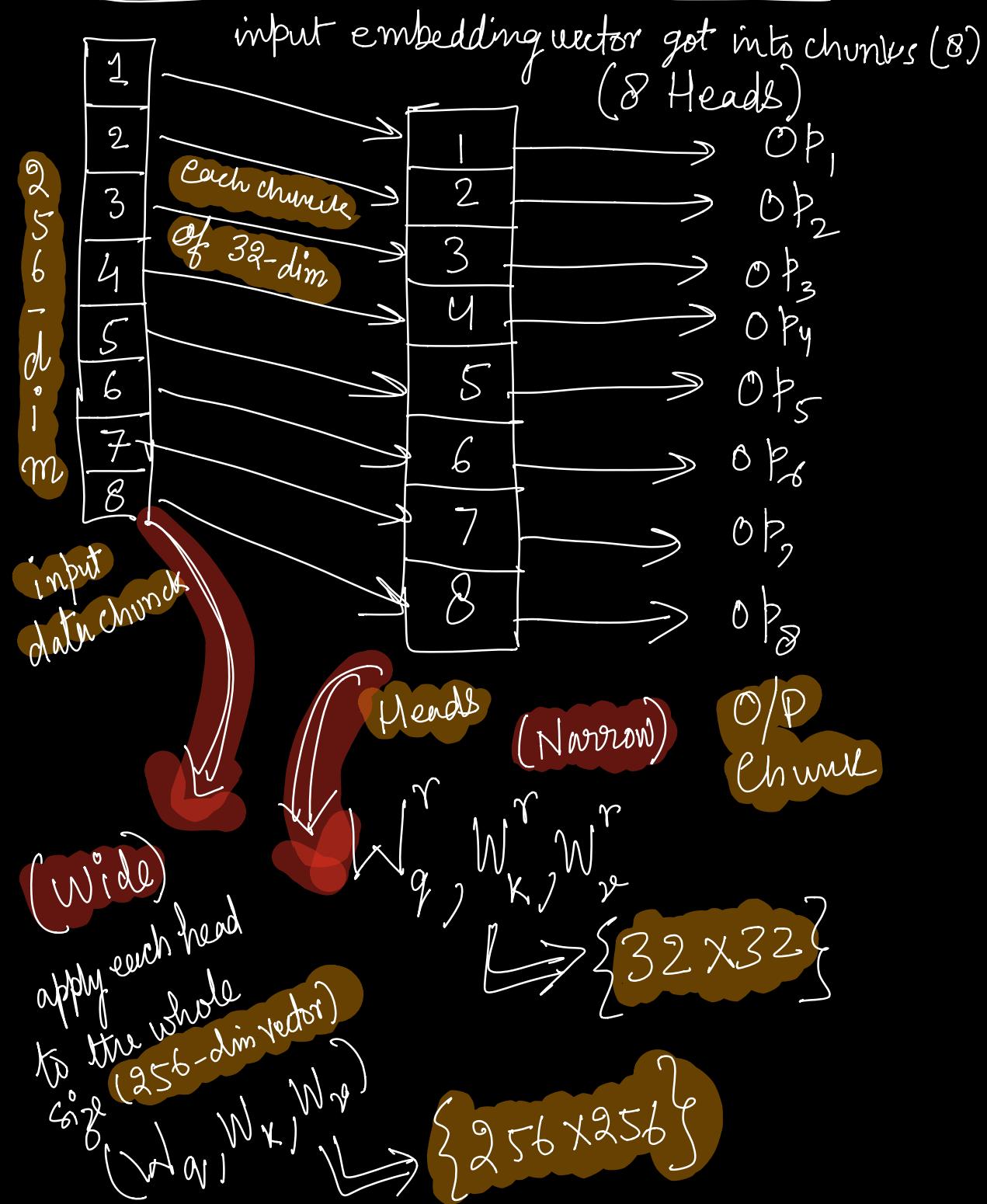
But adding more and more heads will provide more power to SA to address Contextual meaning rather than fixed

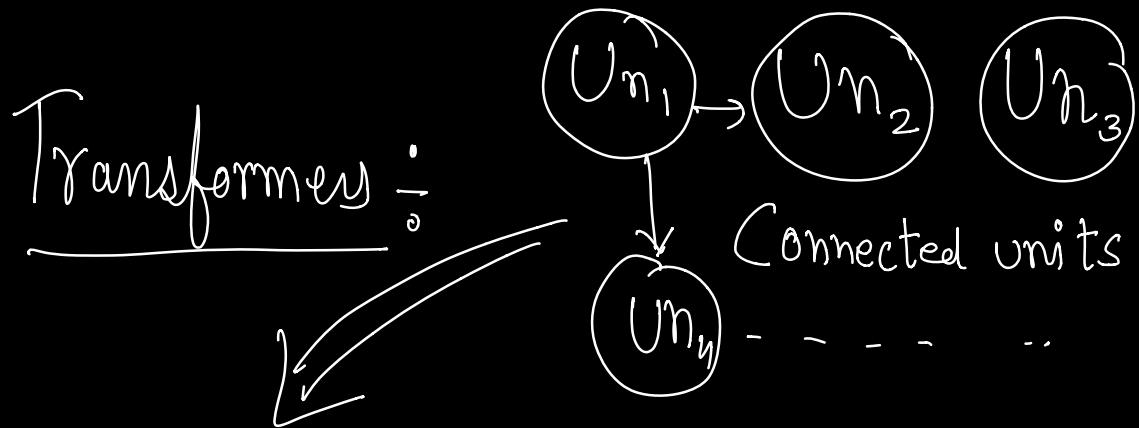
$H_1 \rightarrow W_q^1, W_k^1, W_v^1$
 $H_2 \rightarrow W_q^2, W_k^2, W_v^2$
 \vdots
 $H_r \rightarrow W_q^r, W_k^r, W_v^r$

Multiple attention heads.

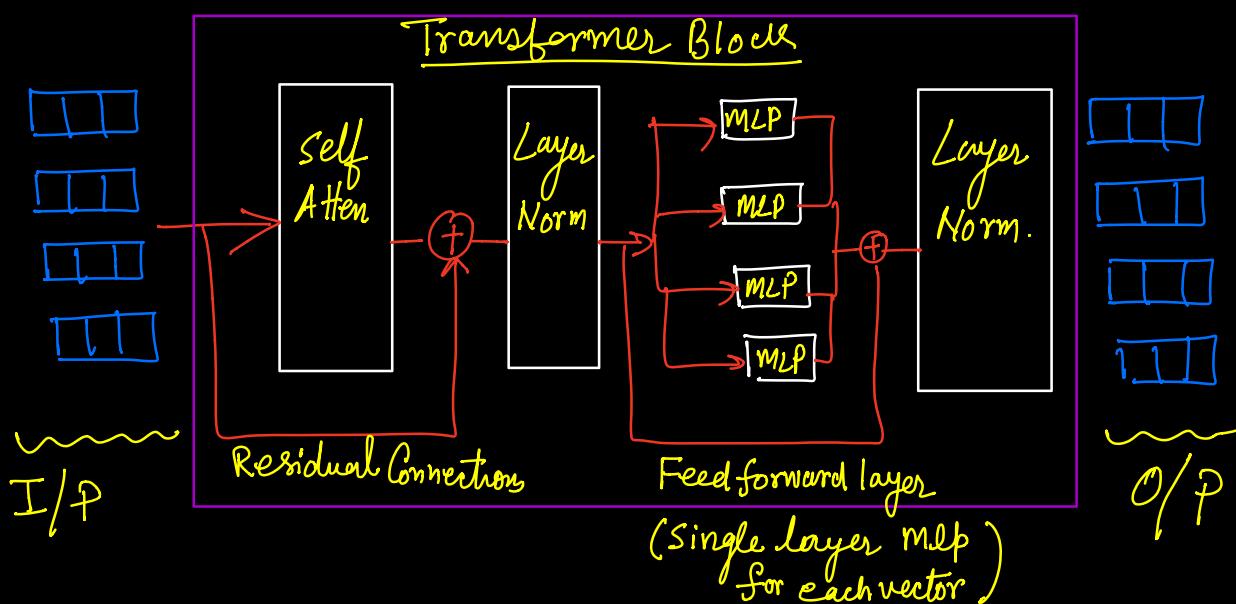


D) Narrow and Wide Self-Attention





An architecture that can process a set of units (Sequence of tokens/pixels) where units can only interact via self attention.
 {Repeated self-attention}

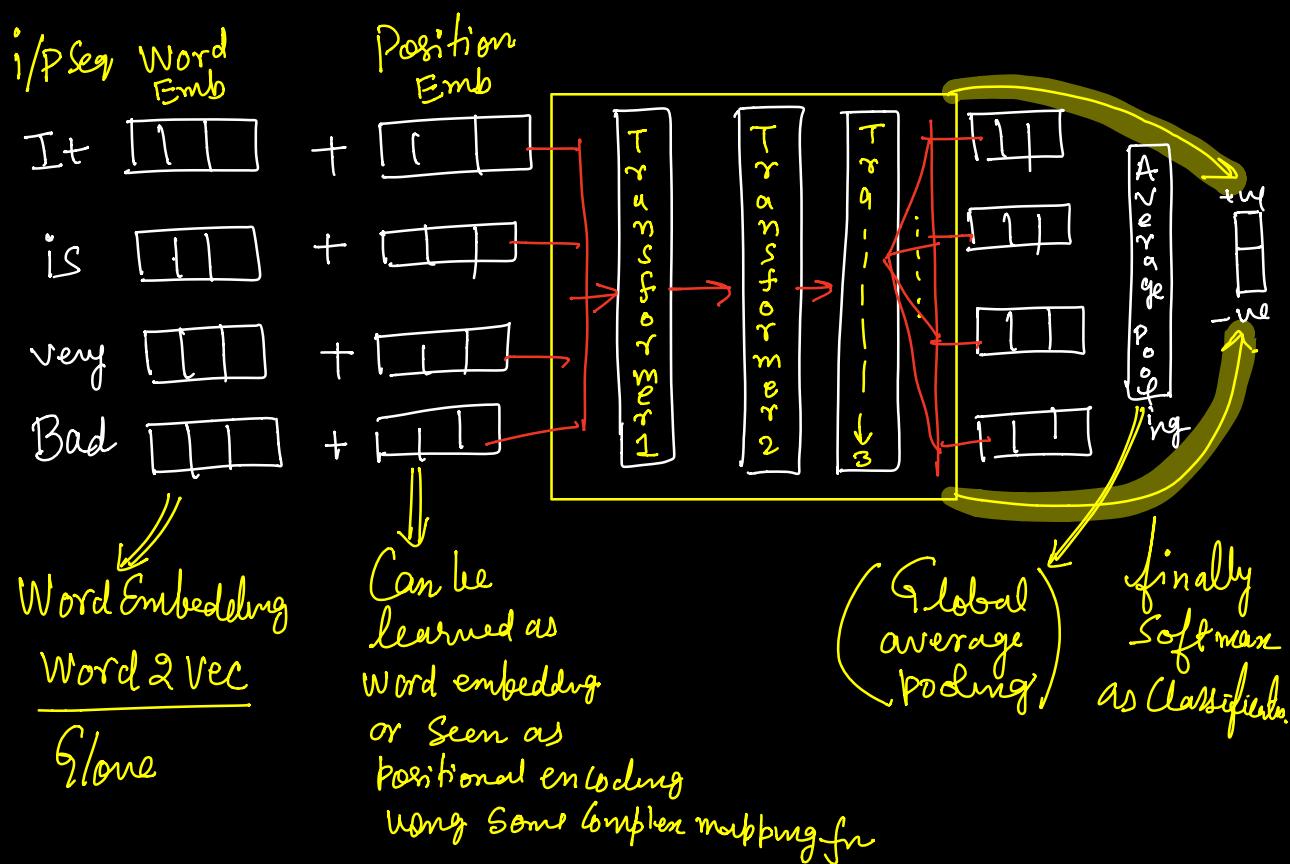


⊗ Normalization and Residual learning are Standard tricks to train deep networks.

I Designing Classification Transformer

IMDb Sentiment classification dataset \Rightarrow tokenized in words
 (make a sequence) for movie reviews
 Classifier $\downarrow \downarrow$
 Binary Classification
 $\{+ve/-ve\}$

Transformer based classification



Requirement of positional embeddings :

The position of any word (w_i) in a sentence (S) plays a very crucial role

Any sentence can be seen as a complex time series with a well defined semantic associated to it

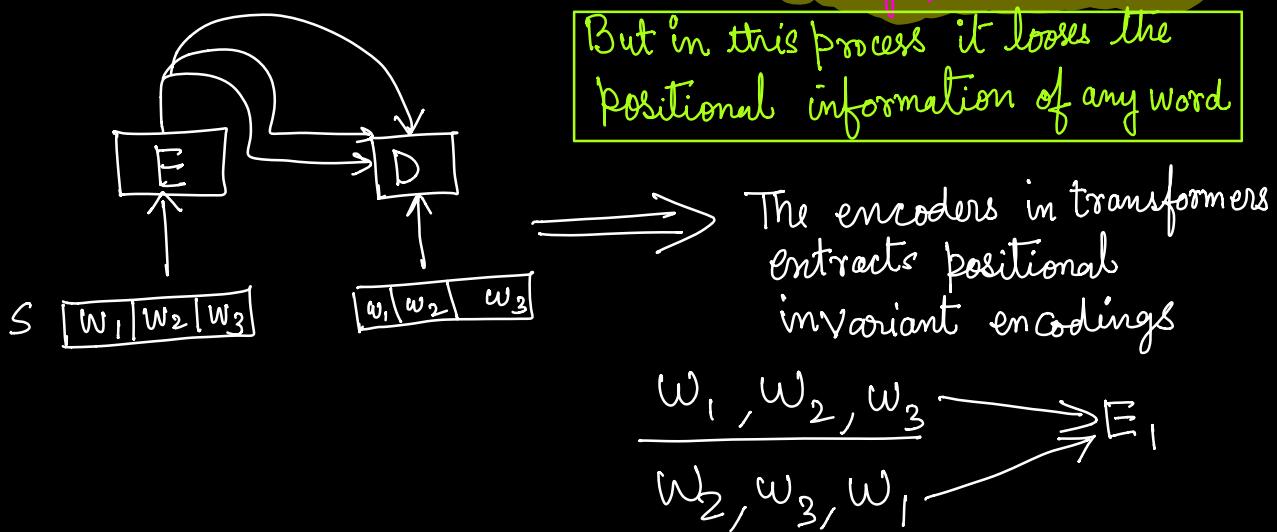
In RNN, the input got fed into the RNN Unit (Sequentially). Hence the ordering information implicitly got embedded

Sequential Computation (slow)

In Transformers multi-head self attention Completely ignores the word ordering.
(Hence can capture longer dependency & fast)

Massively parallel (fast)

But in this process it loses the positional information of any word



Some basic options are :

- ① Decimal encoding for each (t).

$$t = 0, 1, 2, 3, \dots \dots \dots 10,000$$

(Some big No)

↓
Unbounded
Do not generalize
for longer length
Sequences
testing sentences
may be longer
 Δt is
consistent

(Normalize it b/w 0-1)
Bounded

But each i/p is of different size
therefore $\Delta t \Rightarrow$ Gap between

② Timestep is not consistent
across all sentences

Basic Criteria's for encoding

- ① Position of a word in any sentence must have to be unique. (unique time step encodings)
- ② Distance b/w any two timesteps (Δt) must have to be consistent across all length sequence
- ③ Encoding need to be bounded and generalize well for longer sequence
- ④ It has to be deterministic and fixed for any given i/p.
(Should not be a random fn).

Hence, for each timestep (t), in any i/p sentence (S) we require a function $f: \mathbb{N} \rightarrow \mathbb{R}^d$

$$t \in \mathbb{N}$$

d -dimensional
positional encoding

$$\vec{p}_t^{(i)} = f(t)^{(i)}$$

i^{th} element of
the positional
embedding for
time (t).

function f
takes t as an
i/p and gives
a d -dim p.e

$$\vec{P}_t \in \mathbb{R}^d$$

$$\sin(\omega_t k)$$

when $i = 2k$
i.e $i = \text{even}$

[i] ranges from : $1, 2, 3, 4, \dots, d$

[k] ranges : $k=1, k=2, k=3, \dots, k=d/2$

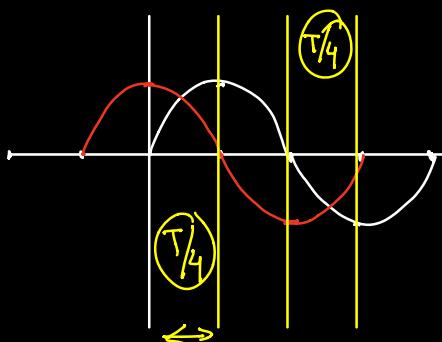
$$\cos(\omega_t k)$$

when $i = 2k+1$
i.e $i = \text{odd}$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 t) \\ \cos(\omega_1 t) \\ \sin(\omega_2 t) \\ \cos(\omega_2 t) \\ \sin(\omega_3 t) \\ \vdots \\ \vdots \\ \sin(\omega_{d/2} t) \\ \cos(\omega_{d/2} t) \end{bmatrix}$$

- * Basically utilizing sinusoidal function with a varying (ω)
- * One (ω) encodes 2 dimensions using (\sin) & (\cos) each.
- * Therefore a total of ($d/2$) frequencies are used to encode (d) dimensions.

Q1: Why do you require more than one sinusoid function.



only $(T/4) + (T/4) \Rightarrow$ any one sinusoid can give us unique values.
 $(T/2)$
 after that values keep on repeating.

\Rightarrow Instead of a single sinusoid if one can use 2 sinusoids may be at same freq then there combination may lead to produce unique encoding for each timestep
 But \Rightarrow it will be (2-D) encoding.

$$\omega = 2\pi f$$

$$f = \frac{1}{T} \quad \& \quad \omega = \frac{2\pi}{T} \quad \text{also} \quad f = \frac{T}{2\pi}$$

Hence $\sin(\omega_1 t)$ & $\cos(\omega_1 t) \Rightarrow$ cover upto $\left(\frac{2\pi}{\omega_1}\right) = T_1$
Time period w/o repeating:

Adding more sinusoids will

increase the time period w/o repeating the sequence.

As suggested in paper

$$\omega_k = \frac{1}{(10,000)^{2k/d}}$$

$$K=0, \quad \omega_0=1, \quad f_0 = \frac{1}{2\pi}, \quad T_0 = 2\pi \quad \left(\text{one } \sin \text{ & one } \cos \text{ for each } K \right)$$

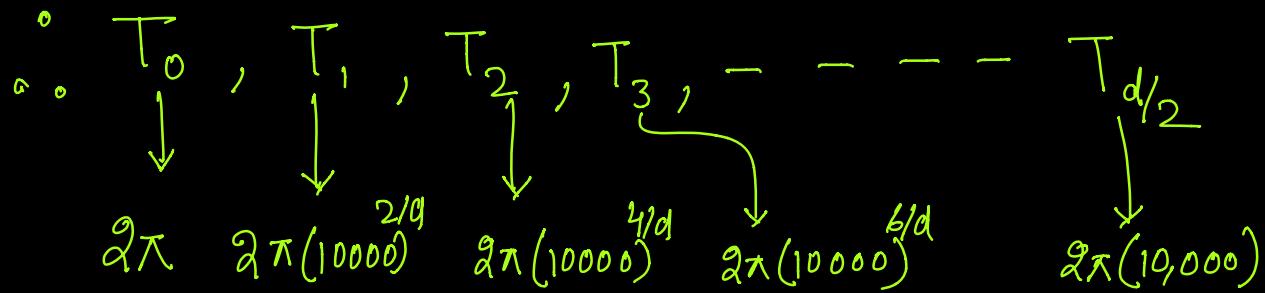
$$K=1, \quad \omega_1 = \frac{1}{(10,000)^{2/d}}, \quad f_1 = \frac{1}{2\pi(10,000)^{2/d}}, \quad T_1 = 2\pi(10,000)^{2/d}$$

$$K=2, \quad \omega_2 = \frac{1}{(10,000)^4/d}, \quad f_2 = \frac{1}{2\pi(10,000)^4/d}, \quad T_2 = 2\pi(10,000)^4/d$$

$$K=3, \quad \omega_3 = \frac{1}{(10,000)^6/d}, \quad f_3 = \frac{1}{2\pi(10,000)^6/d}, \quad T_3 = 2\pi(10,000)^6/d$$

finally $K = d/2$

$$K = d/2, \omega_{d/2} = \frac{1}{(10,000)}, f_{d/2} = \frac{1}{2\pi(10,000)}, T_{d/2} = 2\pi(10,000)$$



- ✳ Systematically increasing the time period of the sinusoid from $2\pi \longrightarrow 2\pi(10,000)$ [Equi-distant]
- ✳ It is a GP with $\tau = (10,000)^{2/d}$

✳ The d -dim encoding is going to be repeated after

$$\text{lcm}(T_0, T_1, T_2, \dots, T_{d/2}) = T_{d/2} = 2\pi(10,000)$$

lcm of a GP is its final term.

$$\begin{array}{ccccccccc} t=3 & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 \\ t=2 & 1 & 2 & 1 & 2 & 1 & 2 & 1 & 2 & 1 \end{array}$$

after (6)

$$\begin{array}{ccccccccc} t=3 & 1 & 2 & 3 & 1 & 2 & 3 & 1 & 2 & 3 & 1 \\ t=4 & 1 & 2 & 3 & 4 & 1 & 2 & 3 & 4 & 1 & 2 & 3 & 4 \end{array}$$

after (12)

Q2: Why dividing $2\pi \longrightarrow 2\pi(10,000)$ into $(d/2)$
 when even lcm of $(2\pi, 2\pi(10,000)) = (2\pi)(10,000)$.
 Instead of taking a few why to take $(d/2)$ frequencies.

Before that let us see how to use these positional encodings.

for any word say (w_t) in any sentence $S = \{w_1, \dots, w_n\}$

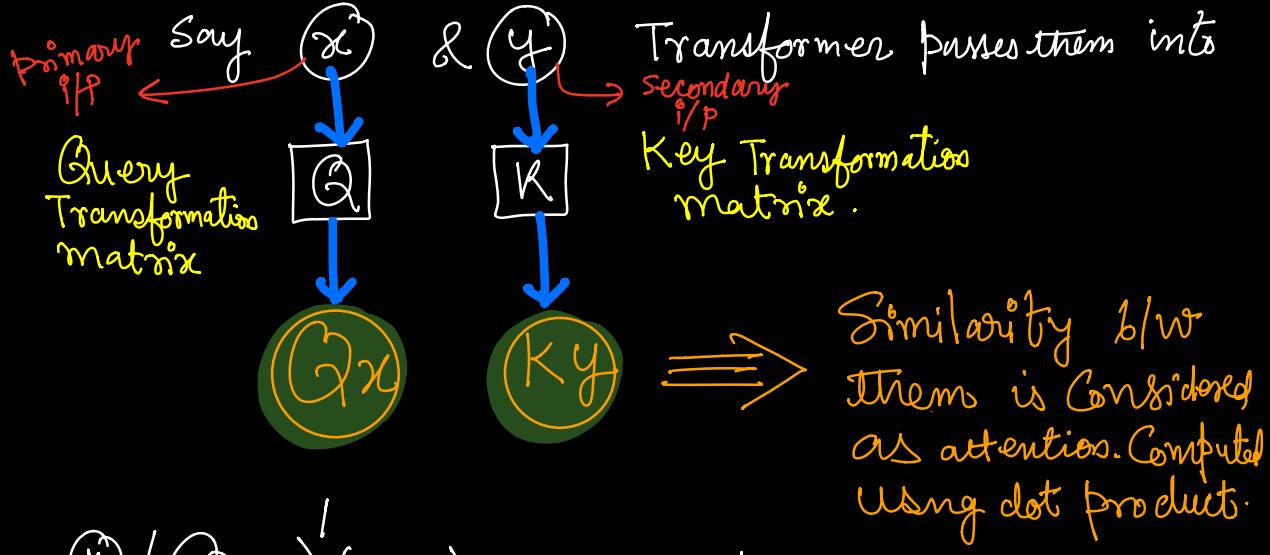
$$\Psi'(w_t) = \underbrace{\Psi(w_t)}_{\text{Word embedding}} + \underbrace{\vec{P}_t}_{\text{position embedding}}$$

final embedding fed to the N/W

They can also be concatenated instead of addition (that is more intuitive) \iff Few reasons supporting additions :

- ① It requires less memory.
- ② With max-timescale = 10,000, (10,001) Word has the same encoding as first word.
- ③ Assuming any seq of length 20 (which is common) with Word embedding of 512 dims, Most of the position embedding dimensions will be used/dominated by word embedding (WE).
- ④ $(WE + PE)$ got trained from scratch in transformer.
↳ This enable not to encode important information in lower dimensions of WE.
↳ for smaller sequence this is almost like concatenation happening.

⑤ In order to learn the attention b/w 2 given embeddings

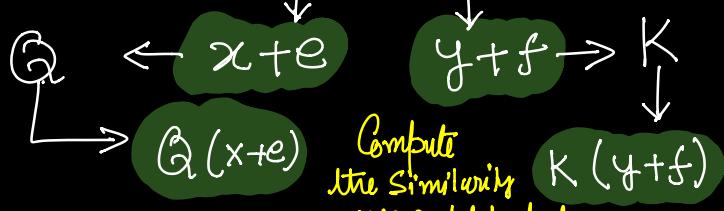


$$*(Qx)'(Ky) = x'(Q'K)y$$

How much attention
Should we pay to Word (x)
Given word (y).

Learning a single transformation matrix ($Q'K$)
It transforms secondary i/p (y) into a space where it can be directly compared with (x).

\Rightarrow Now adding pe's to x & y



$$\begin{aligned} (Q(x+e))' (K(y+f)) &= (Qx + Qe)' (Ky + Kf) \\ &= (Qx)' Ky + (Qx)' Kf + (Qe)' Ky + (Qe)' Kf \end{aligned}$$

$$\Rightarrow x' (Q'K)y + x' (Q'K)f + e' (Q'K)y + e' (Q'K)f$$

Assuming $(Q'K)$ as a single transform $\Rightarrow \mathbb{T}$

$$(Q(x+e)) \cdot (K(y+f)) = x' T y + x' T f + e' T y + e' T f$$

\downarrow \downarrow \downarrow \downarrow
 $AT(x||y)$ $AT(x||f)$ $AT(y||e)$ $AT(e||f)$

Attention to word x given word y

Attention to word x given the position f of word y

Attention to word y given the position e of word x

Attention to the position e of word y given the position f of word y .

Hence learning \mathbb{T} with $(WE + PE)$ need to look for these 4 tasks simultaneously.

to accomplish

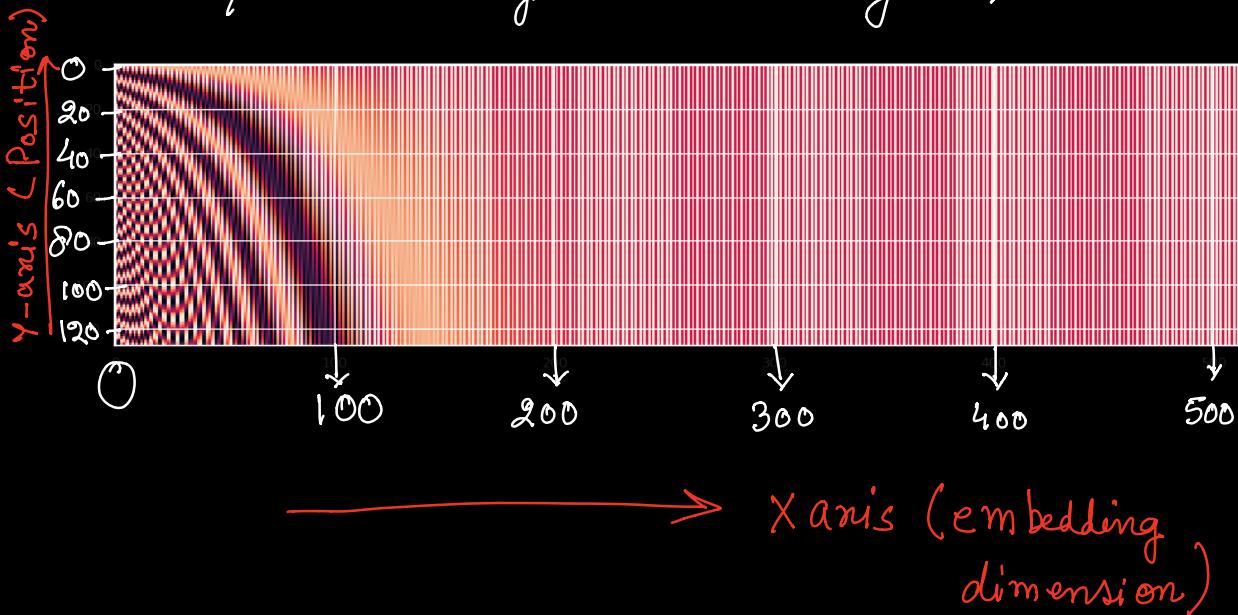
Subspaces
in that high dim.
Space \Rightarrow can be
approximately orthogonal



* Concatenation may give orthogonality/separation between (WE) & (PE) but as both are high dimensional approximate orthogonality is naturally being achieved. (w/o concatenation Cost of more parameters)

Q2 (Answer): Since we are going to add both WE & PE we need to make their dimensions "Same".

Assuming any seq of length Q which is common with word embedding of 512 dim, Most of the position embedding dimensions will be used/dominated by word embedding (WE).



Q3: How does it address relative positions.

(i.e) PE for t & $t+\Delta t$ need to be consistent.

\therefore for some fixed (K) (PE_{pos}) & (PE_{pos+K}) need to be linearly related.

As we all know that for each frequency (ω_k) we have $\sin(\omega_k t)$ and $\cos(\omega_k t)$ components.

Let us assume a linear transformation $(M) \in \mathbb{R}^{2 \times 2}$ independent of (t) as follows:

$$M = f(\Delta t)$$

$$M * \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k(t + \Delta t)) \\ \cos(\omega_k(t + \Delta t)) \end{bmatrix}$$

This ensures that the PE at (t) for (ω_k) frequency is linearly related to that of $(t + \Delta t)$ and is independent of (t) (only depends on (Δt))

Let us compute such an (M) .

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} * \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k(t + \phi)) \\ \cos(\omega_k(t + \phi)) \end{bmatrix}$$

$$\begin{bmatrix} u_1 \sin(\omega_k t) + v_1 \cos(\omega_k t) \\ u_2 \sin(\omega_k t) + v_2 \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \cos(\omega_k \Delta t) \sin(\omega_k t) + \sin(\omega_k \Delta t) \cos(\omega_k t) \\ -\sin(\omega_k \Delta t) \sin(\omega_k t) + \cos(\omega_k \Delta t) \cos(\omega_k t) \end{bmatrix}$$

$$u_1 = \cos(\omega_k \Delta t) \quad v_1 = \sin(\omega_k \Delta t)$$

$$u_2 = -\sin(\omega_k \Delta t) \quad v_2 = \cos(\omega_k \Delta t)$$

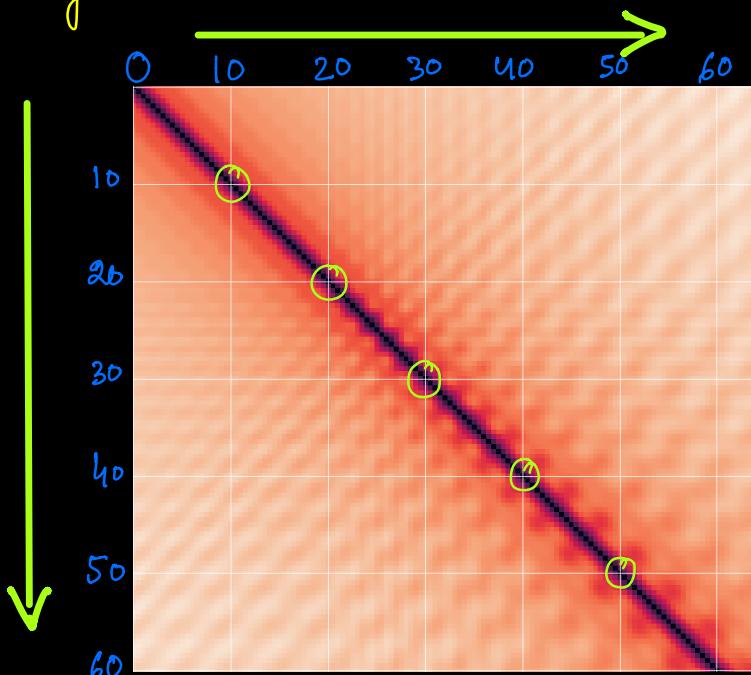
$$\therefore M_{\Delta t, k} = \begin{bmatrix} \cos(\omega_k \Delta t) & \sin(\omega_k \Delta t) \\ -\sin(\omega_k \Delta t) & \cos(\omega_k \Delta t) \end{bmatrix}$$

*Same as
Rotation
matrix*

\Rightarrow Similarly we can have $M_{\Delta t, k}$ (for all k 's)

Hence such a transformation can eventually relate (\vec{P}_t) & $(\vec{P}_{t+\phi})$ linearly

ensuring model learns to attend relative positions.



Showing the
dot product of
PE's for all
time steps.

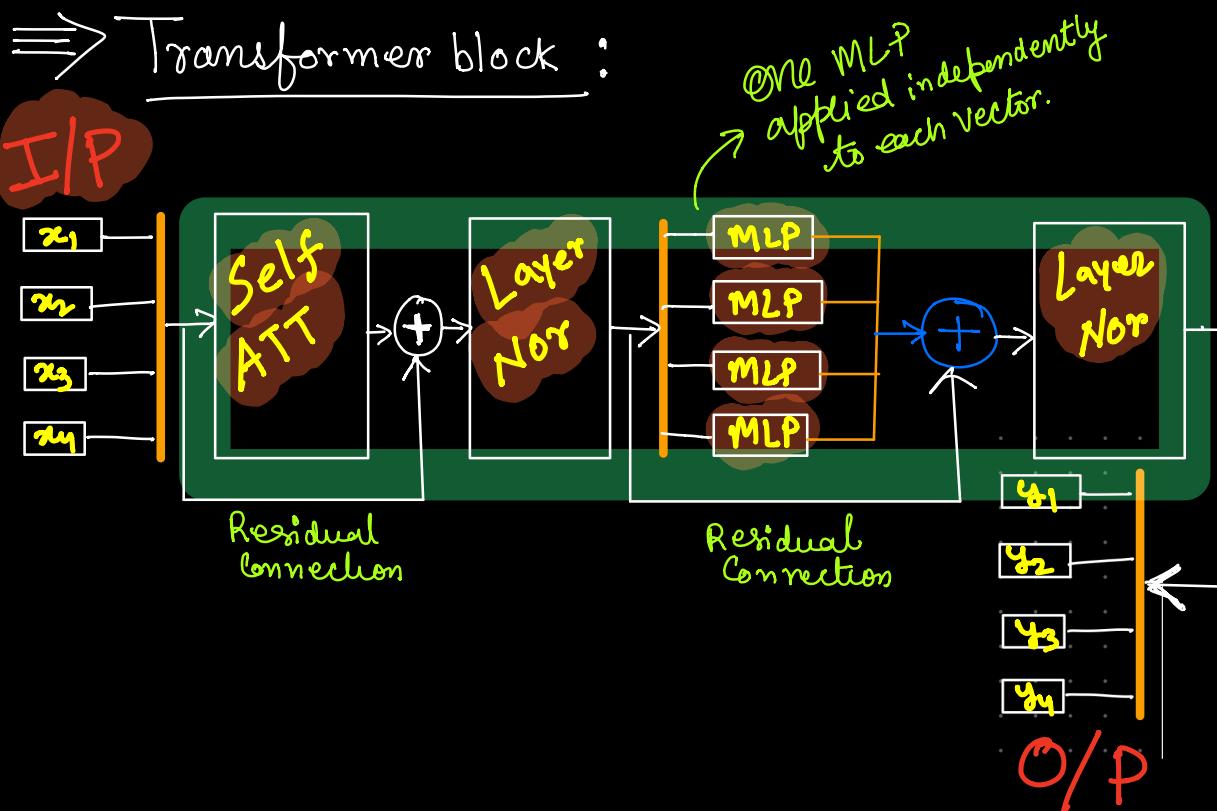
⊗ Distance b/w
neighbouring time
steps are symmetric

⊗ Decay linearly
with time.

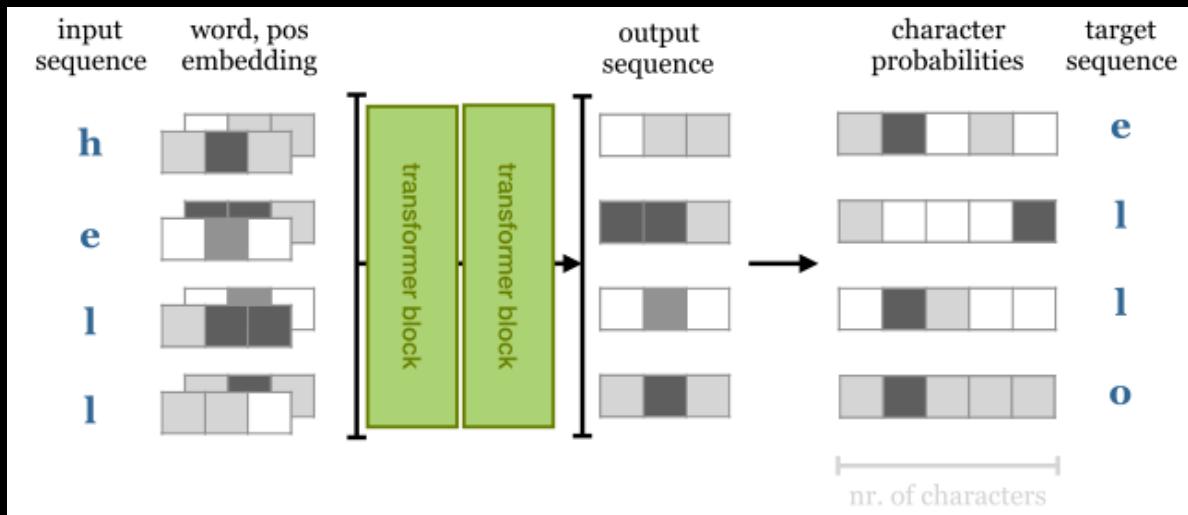
Diagonals are zero and off diagonals are symmetric.

Transformer Architecture uses Self-Attention.

↓
It has units that can process tokens, sequences, set of pixels but units can interact via Self-attention.

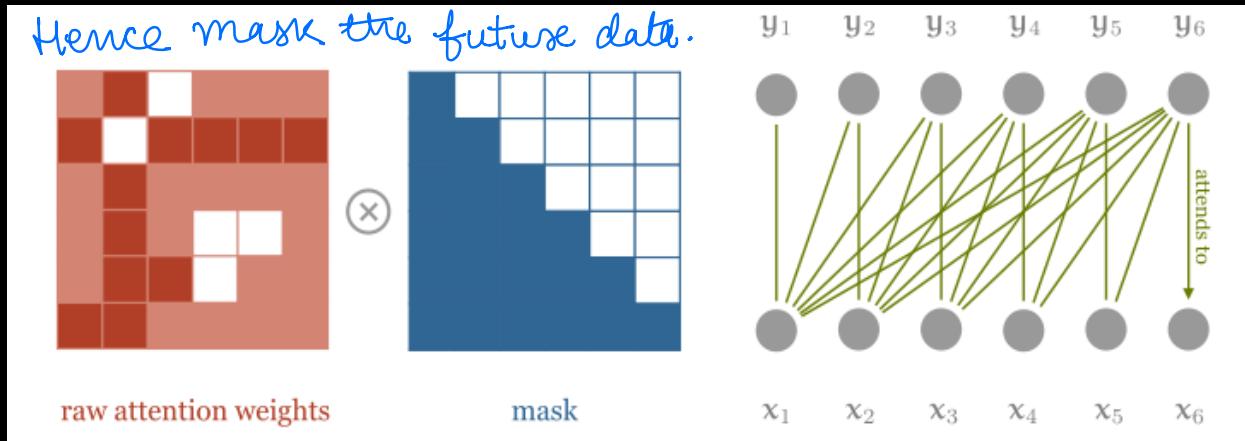


How to generate Text using Transformers



O/P depends upon entire input
so next character prediction is just
(Trivial Step) choosing.

Hence mask the future data.



Only allow to attend the previous data.

