



Purpose - The purpose of this memo is explain the intuition, nuts and bolts of the transformer architecture visually.

What is the problem it solves?

- We are given a sequence of inputs $x_1, x_2, x_3 \dots x_n$ and we want to predict a sequence of output $y_1, y_2 \dots y_n$
- This generalized abstraction can be used to solve many problems as illustrated below:
 - given a sentence in english, what would it mean in french
 - given a review of a product, what is the customer sentiment
 - given an email as input, what should be the response
 - given a sequence of user interactions with a product (TikTok videos), what are they likely to watch next.
 - given a sequence of words describing an image, what would the image look like.

This abstraction is not limited to the modality of text, and extends to image, audio as well thereby enabling us to solve problems in the field of perception as well.

Why do we need CONTEXT to solve these problems?

- Let's take an example.

"Elon Musk is a revolutionary leader. He makes
Gold Gets happen"

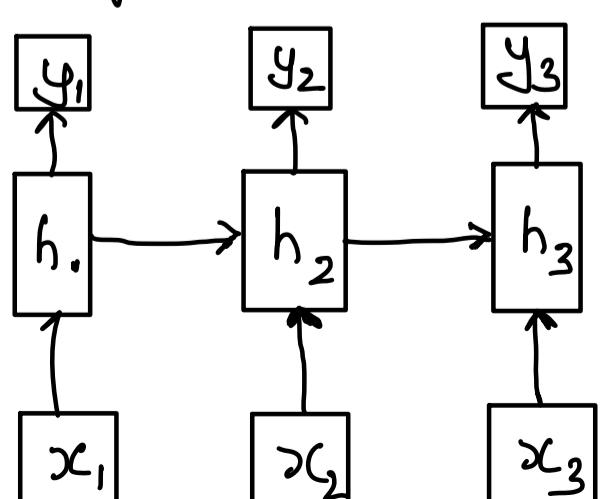
As seen in the example above, to understand the meaning of word "He" we need to retain **CONTEXT** of the last few words, and also pay **ATTENTION** to the words **ELON** and **LEADER**.

Our brain has memory and we don't understand how but it automatically figures out the right words to pay **ATTENTION** to find the right meaning.

How do we get computers to do so?

- for computers this can then be modelled as, $y = \text{some function } F(x_1, x_2, \dots, x_n)$.
i.e. to learn some function over a sequence of inputs x_1, \dots, x_n that can predict y using the ability to retain **CONTEXT** and **ATTENTION**.

What was the older way to solve this problem?
- have a neural network that can carry forward the past inputs x_1, \dots, x_i when it is working to process x_{i+1} as shown in the figure below.



INSIGHTS

- y_2 uses the influence of x_2 (directly) and x_1 (indirectly) and so on.
- H_1, \dots, H_3 are **VECTORS** that carry forward the inputs.
- H_i are also called as **HIDDEN/LATENT states** (for techies)

let's talk about VECTORS briefly

most quantities cannot be represented with a single number. Naive examples include lat, lng which is 2 dimensional or velocity which has speed and direction.

These representations are called as vectors.

As an example to find the distance between SFO & LA given their LAT, LNG vectors.

$$SFO = [\text{lat}, \text{lng}], \quad LA = [\text{lat}, \text{lng}].$$

we can vector difference between the two

$$\Delta \text{lat} = \text{Lat}_{SFO} - \text{Lat}_{LA}.$$

$$\Delta \text{lng} = \text{lng}_{SFO} - \text{lng}_{LA}.$$

and then use Δlat , Δlng in the Haversine formula.

In our example from the past "He depends on the word Elon" so it clearly has to be represented

using a vector.

ONE-HOT ENCODING

$[1, 0, 0, 0, \dots]$

$[0, 1, 0, 0, \dots]$

:

Elon
is
a
great
leader

He
makes
revolutionary
things

$[0, 0, 0, 0, \dots, 1]$

VECTORS WITH CONTEXT

$[3.0, 2.14, 1.789, \dots]$

$[1.34, 1.92, 2.13, \dots]$

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

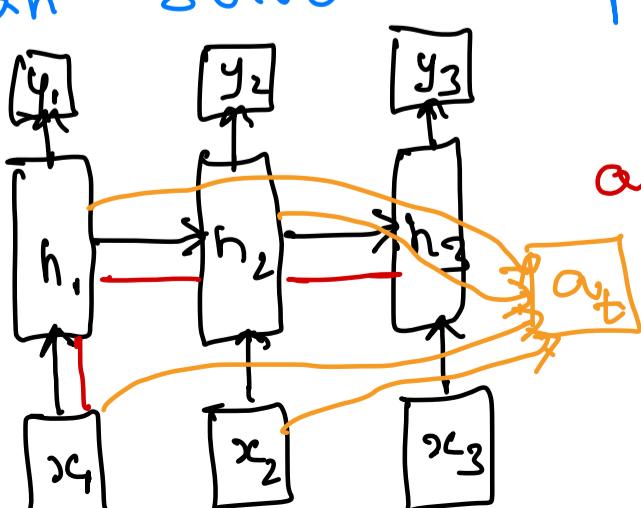
$[2.5], [8.35], \dots]$

- In 1-hot encoding, we represent every word in a unique position. It's a simple way but we can't capture any context.
- In the other vector representation, we capture CONTEXT and ATTENTION. This is more rich and useful.

- Now the Problem statement is how to learn these meaningful vector representation over a sequence of DATA in this case text.?

Things like why neural networks work or the elements that help train them are outside the scope of this memo for now. Lets assume we can train networks on input data, see the error and adjust the weights in the right direction to reduce them. These weights represent a high-dimensional non linear function that learns the pattern in the data.

if we had a network that retains memory and context and can attend to different past inputs then we can solve this problem.



- α_i influence is carried forward through a vector representation h_i .

- α_i - is attention. This nothing but weights to each term (past/future) that help us decide how much attention we should give.

This memory based network is called as RNN (Recurrent Neural Network), LSTM (long short term memory).

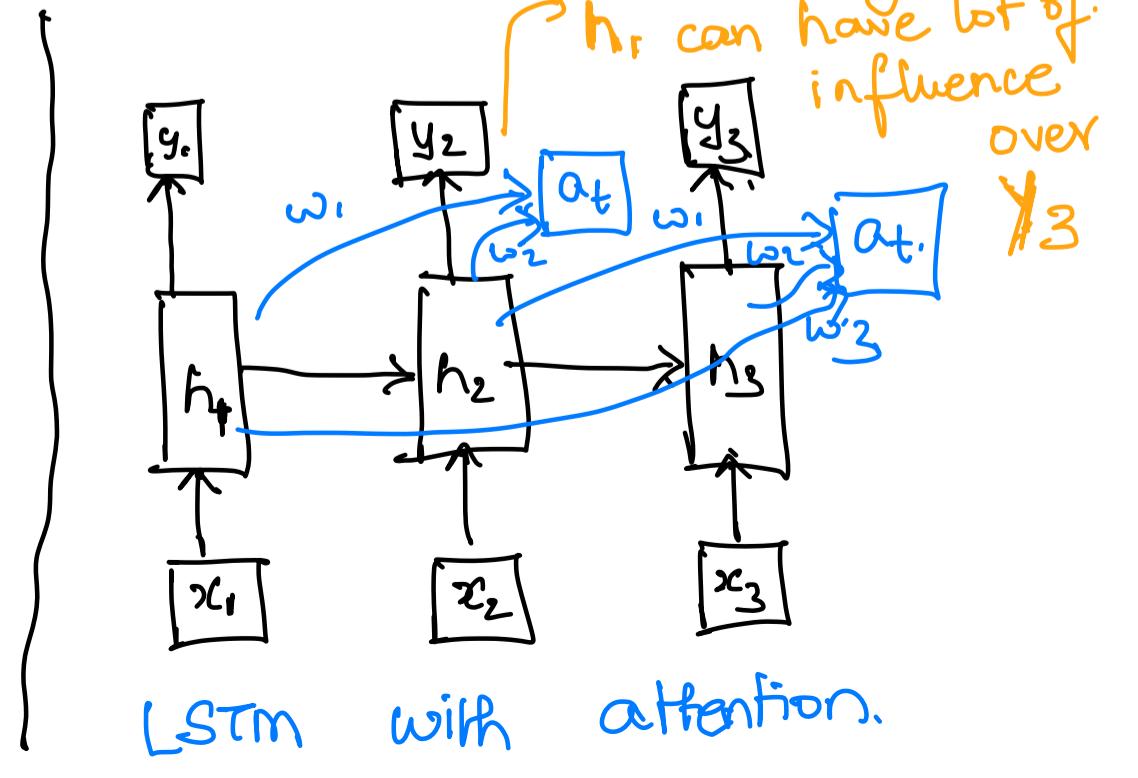
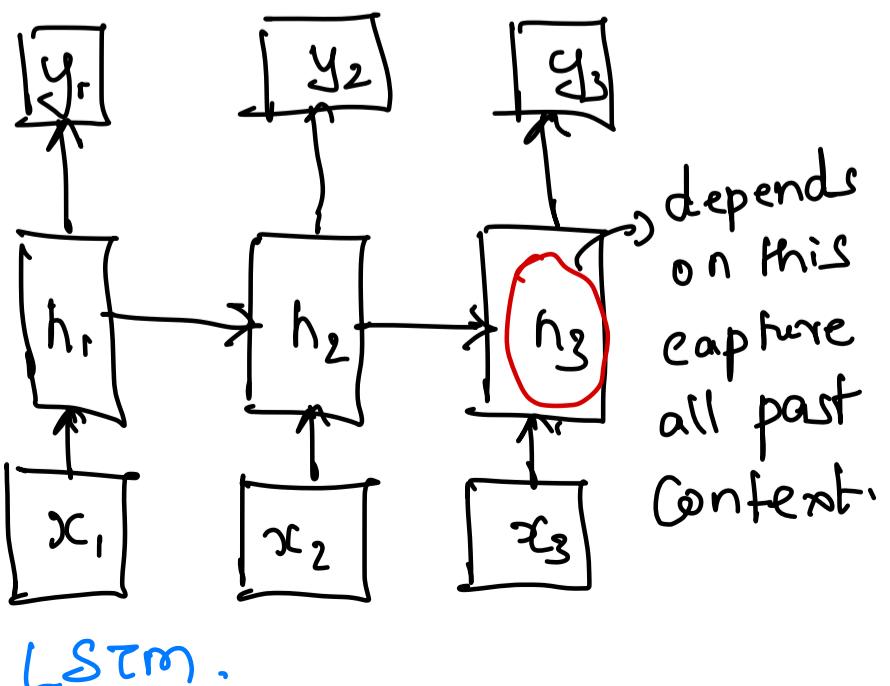
Downside of RNN|LSTM.

- cannot have large context window (gradient explosion while training - for techies).
- cannot be parallelised for training as you make the network deeper.

Key insight - The concept of attention was invented before transformers. It is a simple mechanism by which different states (hidden in LSTM) and inputs (transformed) are weighted and combined together. These weights are also learned as a part of the training process.

Difference Between LSTM and attention.

- Instead of depending on last hidden state to represent all inputs (x_i until now), if attention will consider all hidden states. This should lead to better context.



How do we define the weights in attention?

We learn them depending on the relative importance of the term as a part of the learning process.

We cannot go very far without being able to retain large enough CONTEXT. - Why?

Imagine giving the entire book of Shakespeare as input or a thriller novel where the context from the first few & final chapters is needed to figure out the climax. How do we get neural networks to have such large CONTEXT.

enter \downarrow TRANSFORMERS.

KEY INSIGHT - What if instead of paying attention to hidden states, we can directly pay attention to inputs (x_i) then we have solved the problem of trying to maintain a large enough context.

\downarrow
Let's see HOW.

Let's consider the sequence

$$\text{Elon} \rightarrow x_1 = [x_{11}, x_{12}, x_{13}, \dots, x_{1,512}]_{1 \times 512}$$

$$\text{is} \rightarrow x_2 = [x_{21}, x_{22}, x_{23}, \dots, x_{2,512}]_{1 \times 512}$$

$$\text{revolutionary} \rightarrow x_3 = [x_{31}, x_{32}, x_{33}, \dots, x_{3,512}]_{1 \times 512}$$

$$\text{leader} \rightarrow x_4 = [x_{41}, x_{42}, x_{43}, \dots, x_{4,512}]_{1 \times 512}$$

Why we have chosen 512 as size of vector?
- can it capture all content for x_4 . This is more than science

Each of the input vectors $x_1 \dots x_5$ can be represented as a matrix as shown below. The values of vectors are learned as a part of the learning process.

x_{Elon}	\rightarrow	$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1,512} \\ 2.21 & 1.65 & 2.413 & & 8.46 \end{bmatrix}$	$\xrightarrow{\text{Sample values.}}$
x_{is}	\rightarrow	$\begin{bmatrix} x_{21} & x_{22} & x_{23} & \dots & x_{2,512} \end{bmatrix}$	
x_{a}	\rightarrow	$\begin{bmatrix} x_{31} & x_{32} & x_{33} & \dots & x_{3,512} \end{bmatrix}$	
$x_{\text{revolutionary}}$	\rightarrow	$\begin{bmatrix} x_{41} & x_{42} & x_{43} & \dots & x_{4,512} \\ 3.33 & 9.62 & 9.87 & & 18.737 \end{bmatrix}$	
x_{leader}	\rightarrow	$\begin{bmatrix} x_{51} & x_{52} & x_{53} & \dots & x_{5,512} \\ 1.11 & 2.15 & 4.36 & & 5.99 \end{bmatrix}$	5×512

Input matrix X is also known as

$$X_{\text{Elon}} = [x_{11} \quad x_{12} \quad x_{13} \quad \dots \quad x_{1,512}]$$

and so on

Vectors are also called as **TENSORS**. and generally represented as **COLUMN MATRICES**.

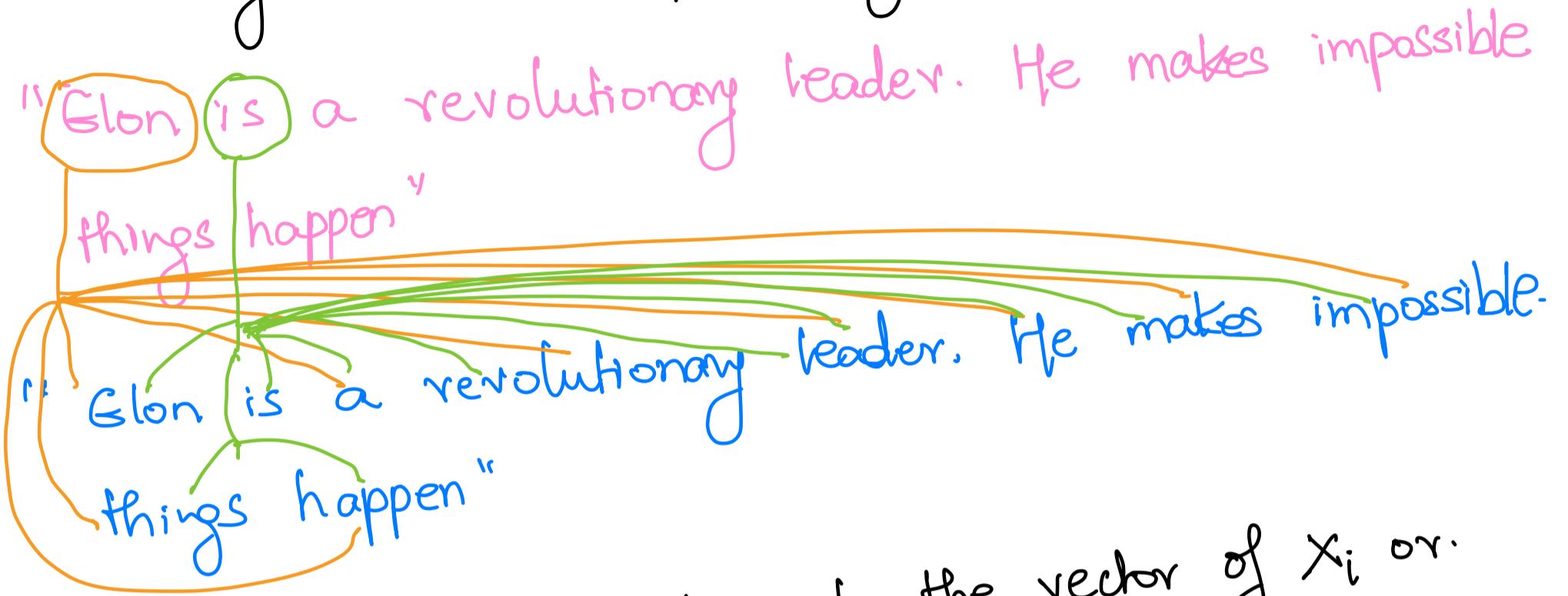
$$X_{\text{Elon}} = \begin{bmatrix} x_{11} \\ 2.21 \\ 1.65 \\ 2.413 \\ \vdots \\ 8.46 \end{bmatrix} \quad \Rightarrow \quad X = \begin{bmatrix} x_{1(\text{Elon})}^T \\ x_{2(\text{is})}^T \\ x_{3(a)}^T \\ x_{4(\text{revolutionary})}^T \\ x_{5(\text{leader})}^T \end{bmatrix}$$

Transpose in matrix represented by T , it flips columns & columns to rows. each column represents a word here.

note change in dimens ion.

512×5

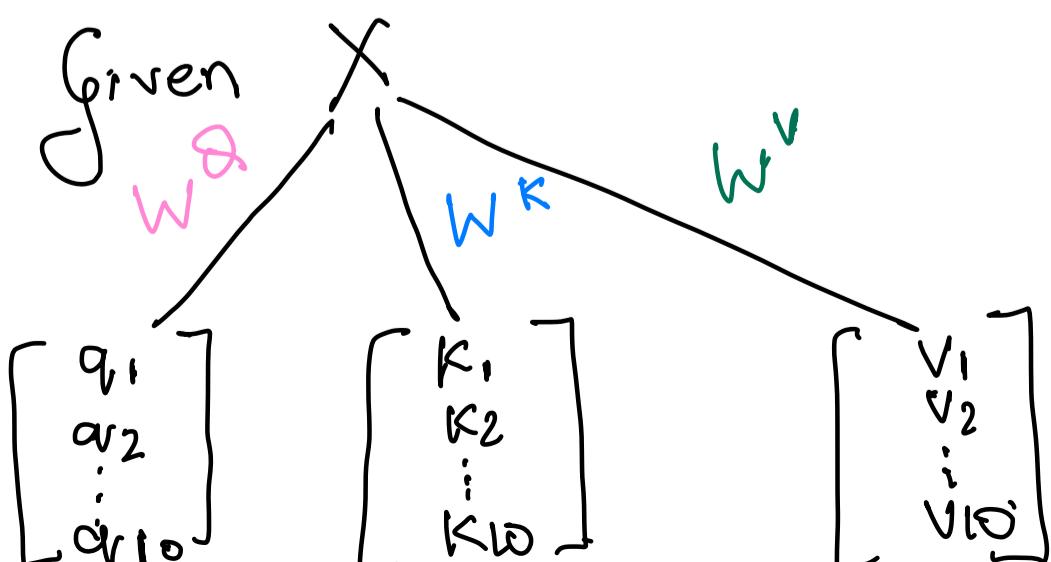
Given these input representations, we have to learn a way that how much **ATTENTION** we should give to every other word for a given word.



- Mathematically, we have to the vector of x_i or x_{Elon} and find how similar it is to all the other vectors of all words in x_i . This is visible in the pink and green arrows shown here.
- Instead of directly working on x_i , we will work on a non-linear representation of x_i . Let's define these as.

$$\text{Query} = Q \cdot (\text{Pink representation of } x \text{ multiplied by some weights})$$

$$\text{Key} = K \cdot (\text{Blue representation of } x \text{ multiplied by some weights})$$



- we know that the best way to learn weights is a fully connected neural network.

let's take x_i from X which the vector for word "Elon".
 we know x_i is 512 rows and 1 column vector = (512×1) .
 Let's say we want to compress the dimensionality of.
 the representations of $X(Q, K, V)$ to 64. I think the.
 intuition here is from auto-encoder (outside scope).
 where the network is forced to learn. Then.

$$x_i = (512 \times 1)$$

$$\text{and } q_{i1} = (64 \times 1)$$

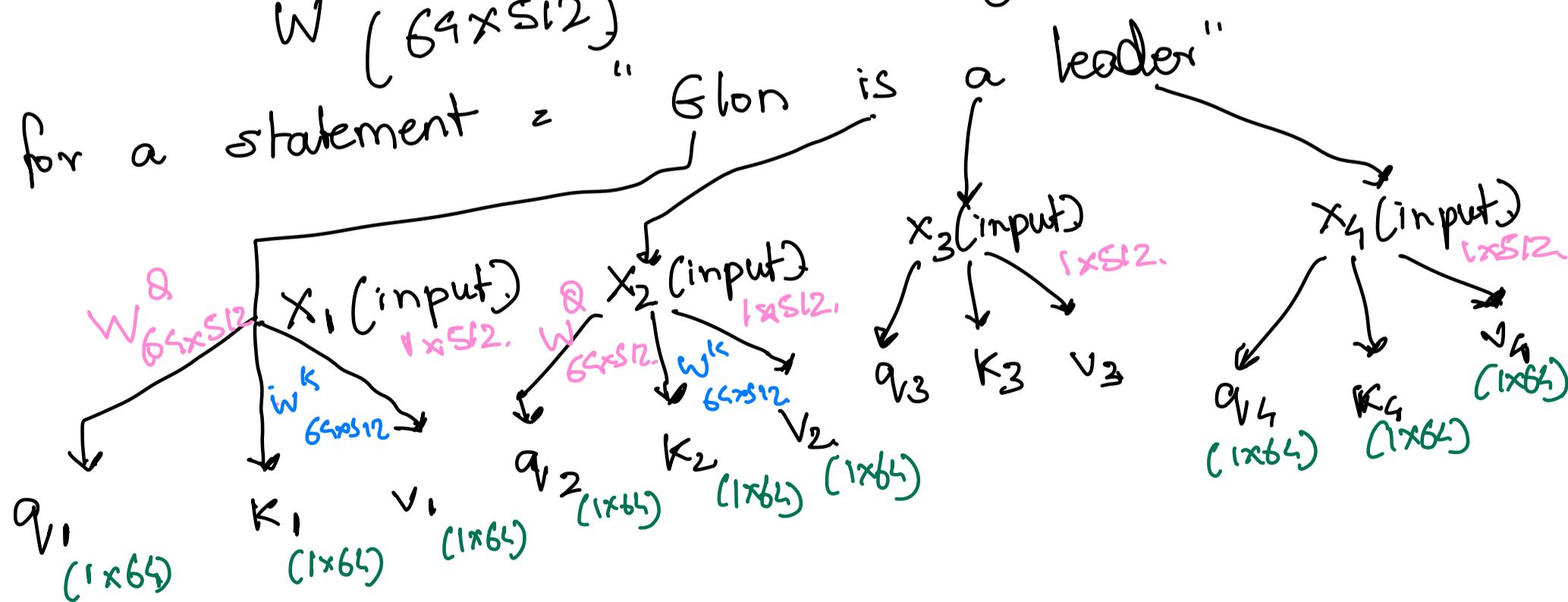
so the weight matrix of q_{i1} i.e. $W^Q = 64 \times 512$.

$$\text{Since } q_{i1}(1 \times 64) = (W^Q \cdot x_i)_{(64 \times 512)}^T$$

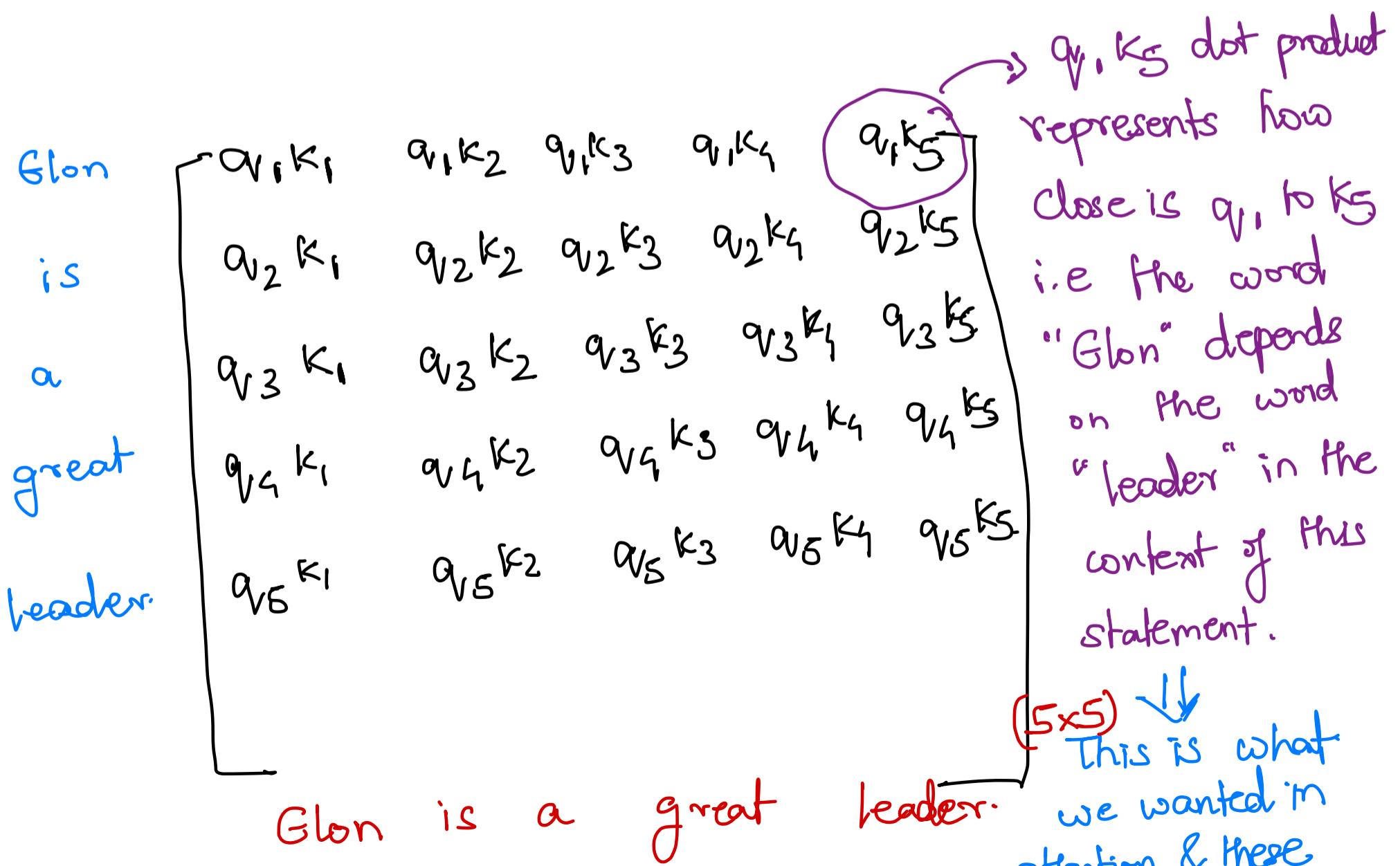
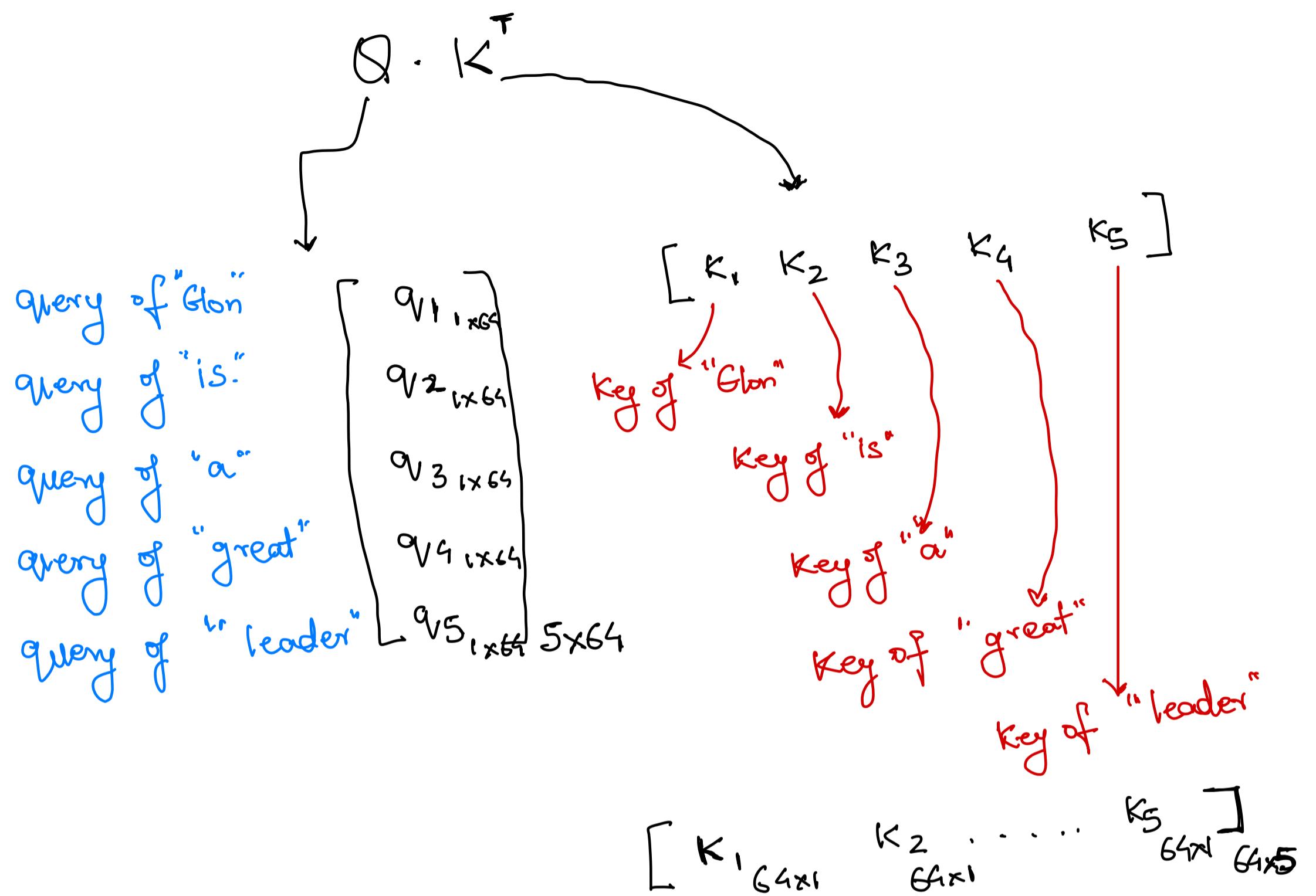
Similarly, we will learn a separate weight
 matrix to get key and value representations
 They will be defined as.

$$W^K \quad \text{and} \quad W^V \quad (64 \times 512)$$

for a statement = "Elon is a leader"



Now, let's calculate "SELF ATTENTION"



- Using this attention matrix, we can know how much attention we have to give to every other word for a given word.

The size of the attention matrix in this case is 5×5 . We took two matrices of size 5×64 and 64×5 to get this.

Now, why do we need to normalize the attention matrix.

- The dot product between K and Q vectors represent the similarity within them. As the dimensionality of these vectors (in this example 64) increases, the magnitude of the dot product increases simply because more terms are contributing to the sum. Without normalization, this increase in magnitude can lead to large values in attention scores. This would result in large variance in attention scores. This step we will
- Another reason is that after this step we will apply a function called as SOFTMAX to convert these attention scores into probabilities. That function is sensitive to large variations in inputs values, it can lead to gradients that are close to zero, and slow down the learning process due to vanishing gradient problem.

- The choice of $\sqrt{d_K}$ (key vector) as normalization factor comes from empirical observations. It was found that it prevents softmax function from entering regions of small gradients.

After normalization our attention matrix will look like.

$$\frac{1}{\sqrt{d_K}} \begin{bmatrix} q_{V1} K_1 & \dots & \dots & \dots & q_{V1} K_5 \\ \vdots & & & & \vdots \\ q_{V5} K_1 & \dots & \dots & \dots & q_{V5} K_5 \end{bmatrix}$$

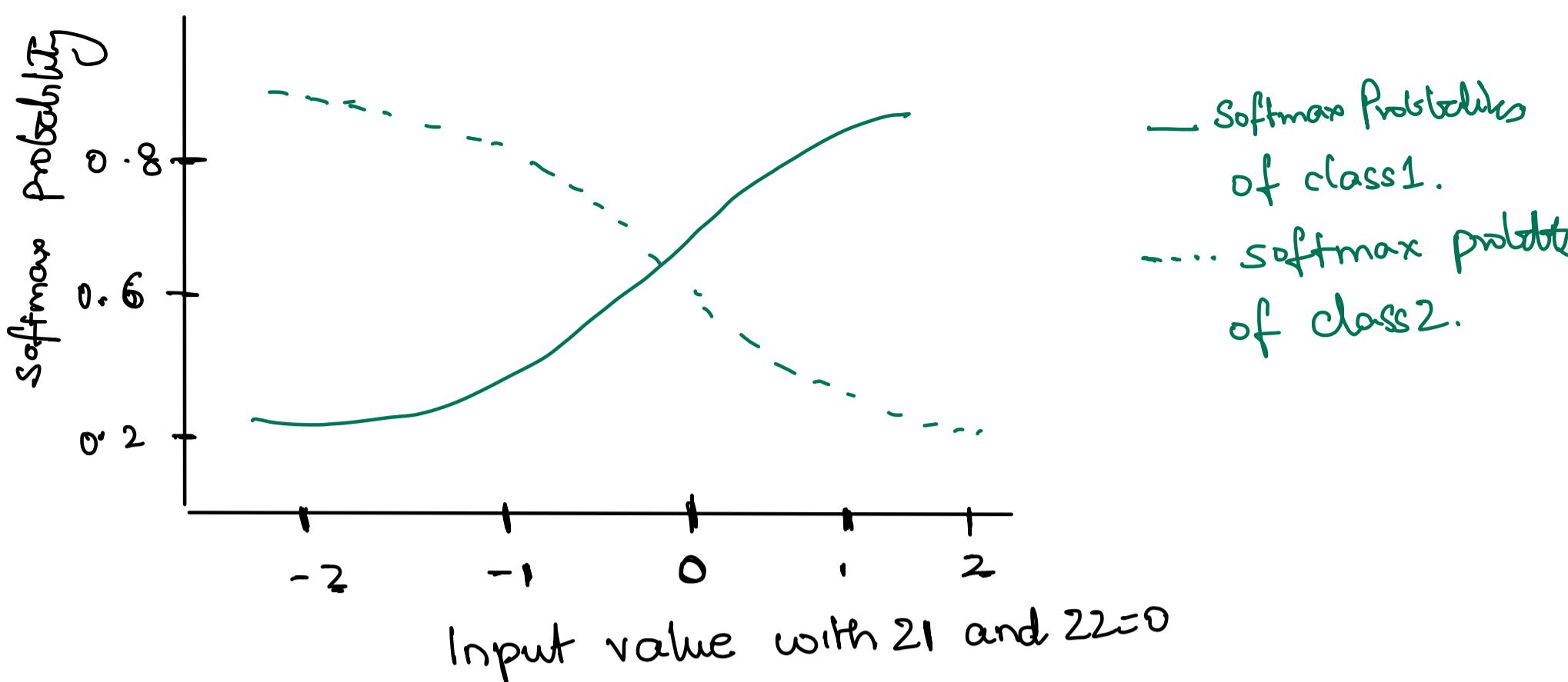
Remember this was $Q \cdot K^T$

Let's discuss how do we convert these attention scores into probabilities and why?

Why? - In the first row of this matrix, it represents how likely is the word "Glon" related to all other words in the given sentence. The other words are represented as columns. This is why we need a way to calculate probability distribution for every row. The other way to think about this is that we are comparing Query of one word with the key of all other words.

$$\text{SOFTMAX} \left(\frac{1}{\sqrt{d_K}} \begin{bmatrix} q_{V1} K_1 & \dots & \dots & \dots & q_{V1} K_5 \\ \vdots & & & & \vdots \\ q_{V5} K_1 & \dots & \dots & \dots & q_{V5} K_5 \end{bmatrix} \right)$$

These scores produce a probability distribution. softmax is applied in the horizontal direction.



In this example we vary the input value z_1 keeping z_2 fixed at zero. As you can see as the input value z_1 increases, the probability of class 1 increases and approaches 1, and at the same time probability of class 2 approaches ZERO.

Softmax takes model outputs and converts them into a probability distribution across multiple classes.

The curve shows that even small changes in inputs can lead to significant changes in probabilities especially when inputs are close to each other. This is a feature of softmax - HOW?

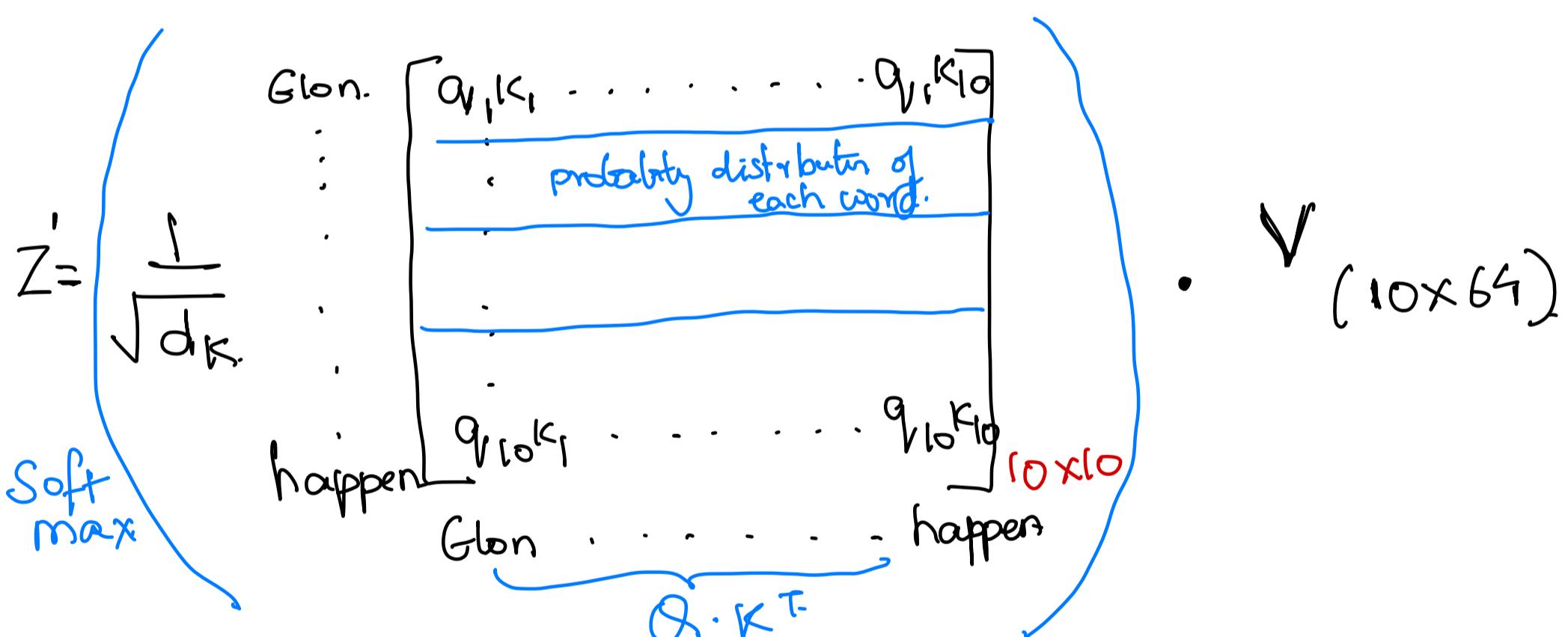
$$\text{softmax}(z) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

It applies exponential to each value of z_i (VECTOR). This ensures all values are positive and amplifies differences between input values (larger values become significantly larger than smaller values). Beneficial for highlighting most likely class in classifications.

To summarize, softmax does the following.

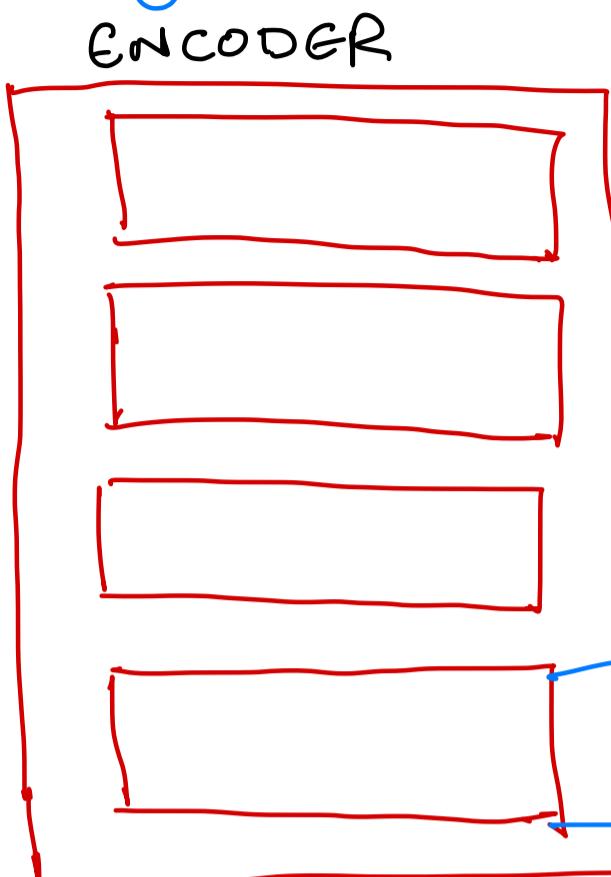
$$\text{out} = \begin{bmatrix} P(Y = \text{class1} | x) \\ P(Y = \text{class2} | x) \\ \vdots \\ P(Y = \text{classN} | x) \end{bmatrix}$$

Finally we take the probability distribution of attention scores and multiply it to \mathbf{V} (the third non-linear representation of x after \mathbf{Q} & \mathbf{K}) to apply attention to the given sentence.



$$Z = \begin{bmatrix} (z_1) \rightarrow 1 \times 64 \\ \vdots \\ \vdots \\ (z_{10}) \rightarrow 1 \times 64 \\ 10 \times 64 \end{bmatrix}$$

Putting it all together -



$$Z = \text{Softmax} \left(\frac{1}{\sqrt{dk}} (Q \cdot K^T) \cdot V \right)$$

x_1, x_2, x_n (inputs)

...

...

MHA

FFN

Feed forward network used to learn
 w^Q, w^K, w^V

multi head

attention (running in

z_1, \dots, z_n SELF ATTENTION

outputs.

N transformer blocks stacked
on top of each other

↳ What is the intuition behind

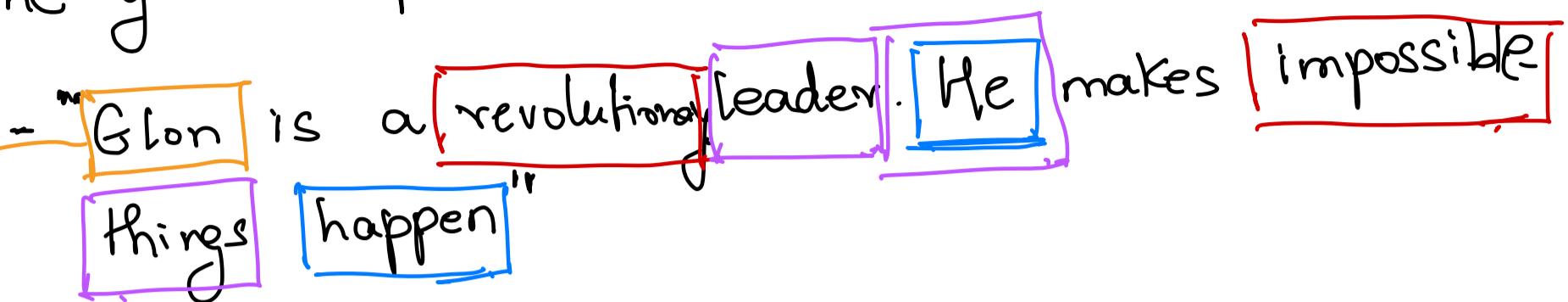
N blocks

- from my intuition of training, CNN's & Deep networks, as you stack up layers, the layers deeper in the network learn higher order abstractions. As an example, if the input is an image, the first layer learns where eyes, nose, ears are in a face. The further down layers in a network use these as features and start to learn how a face looks like (eyes are 6cm apart and so on). With language also, as you read a sentence then the entire paragraph, chapter, chapters, you learn higher order semantics about the language.

- This is visible in the magic of GPT - trained to learn next word prediction but learns all possible things on top

Why do we need MULTIGEAD ATTENTION?

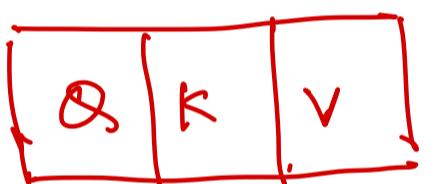
In the given input statement



MULTI-HEAD ATTENTION WORKS.

Head1. will have it's own Q, K, V matrices
→ for a given word/token, let's take "ELON"
and maybe paying attention to some words
in the statement. as shown in BLUG.
"He" and "Happens"

Head2



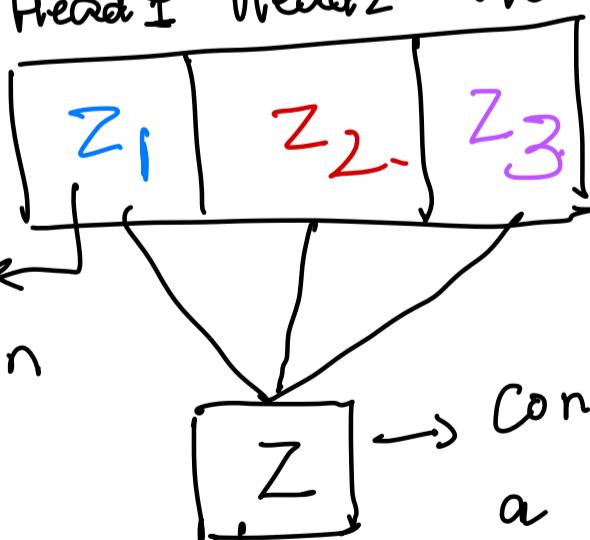
with it's Q, K, V would be paying attention
to "revolutionary" and "impossible"

Head3.



with it's Q, K, V would be paying attention
to "leader", "He" and "things".

all heads will generate their own output 'Z' matrix.



64 in dimension

64 in dimensions.

Concat all z_1, z_2, z_3 into Z using
a linear layer transformations

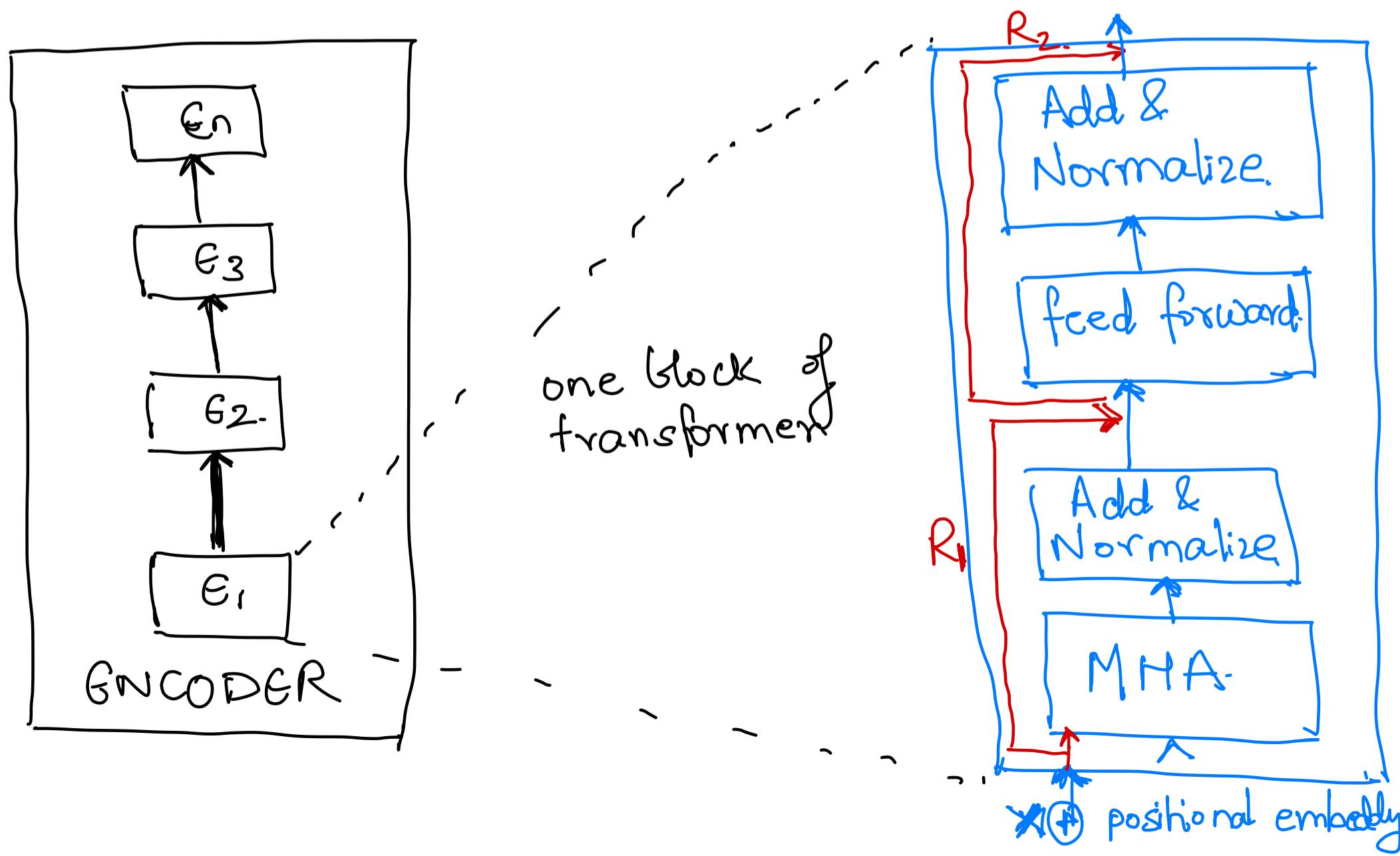
Key Insights

- 1.) Include various contexts in training.
- 2.) each head learns it's own context.
- 3.) Multi head leads to a better representation of context.

- Capture local &
global dependencies.

How are keeping the order of words/tokens in the SEQUENCE?

- In the current form of encoder if you give the statement "Glon is a revolutionary leader" in any order of words to the ENCODER, the output will be same.
We need some way to ensure that the knowledge that "Glon" is the first word is captured. This will retain the order of sequence and change if we pass a different permutation of input.
- if we change X (input vector) such that we add $X_i + P_i$ where P is a vector that captures the position of token X_i in the sequence then every different permutation of X will have a different vector. This is how we taken care of sequencing.



- R1 & R2 are residual connections.

The intuition of residual connections comes from RESNET but can be thought of as if deeper layers learn higher order abstractions of images and that then through residual layers we ensure that they still have access to the original input in the deeper layers.

- The way they are handled in code is by doing vector addition (element wise) for the vector z_i from layer j of encoder to x_i & then pass this as input to layer j of encoder.
- MHA contains self attention blocks (more than one to capture deeper context).

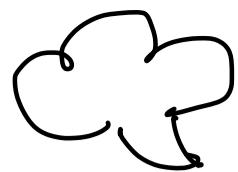
- Normalization used here is layer normalization. Instead of working on batches (like in NN) it works on normalizing every input sample w.r.t. its features. (Goto Chat GPT for a better explanation).

What can we build with this transformer
based ENCODER?

- Given text or images as inputs, it will use attention to learn vector representations of these inputs. Using deeper architecture, it will capture higher order abstractions in the inputs

Using these embeddings

- 1) Build embeddings on amazon review dataset and classify the sentiment of users as positive or negative.
 - 2) Build embeddings of keywords being used in search & products to recommend. Use Nearest neighbour search to find which keywords are more relevant to which product. Build a good enough search experience.
 - 3) Take images of content posted on social media and build embeddings and then classify content as acceptable or not.
- ⋮



How do we get to the magical state that it starts to generate content like GPT?

- How do we solve the following problems.

- a) While typing an email, complete the next few words.
- b) Build something like Chatgpt. - given a prompt of text, complete the next sentence with a response that is useful and meaningful to a fellow HUMAN.

Let's enter DECODER'S.

We need an architectural tweak so that the following things get achieved-

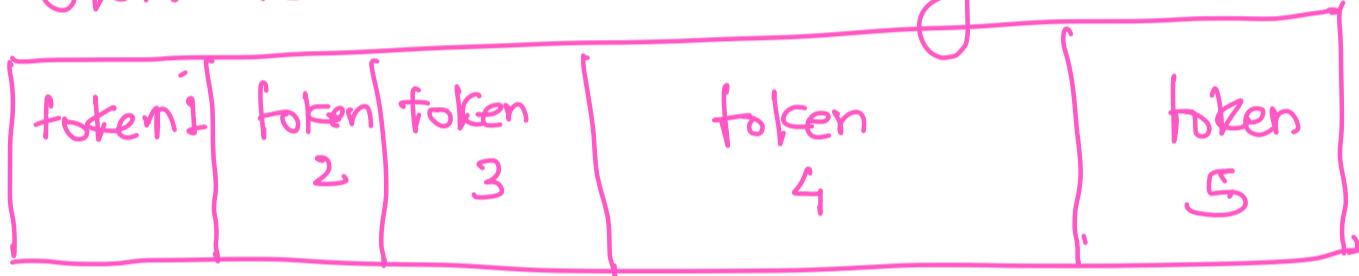
- 1) When we are calculating attention for input x_i , we do not look at any token from x_{i+1} and only look at tokens upto x_i .
- 2) as we generate an output, we feed it back as input to use it as an extended prompt to generate the probability distribution of next word.

Turns out to achieve this, we need just one tweak to our encoder architecture, and specifically the attention blocks.

for input Elon is a revolutionary leader.



Elon is a revolutionary leader.



when processing token 3 (α), it will only consider token 1/2/3 for attention and not token 4 & 5.

This allows the decoder to predict next words based on attention scores of previous tokens (1..3).

This is called as MASKED SELF ATTENTION.

we achieve this through the following simple way -
Put negative infinity in positions that don't need attention

Elon	$q_{1,1} K_1$	\dots	\dots	$q_{1,5} K_5$
is	:			
a	:			
revolutionary	:			
leader	$q_{1,1} K_1$	\dots	\dots	$q_{1,5} K_5$

Elon is a revolutionary leader.

$$- \begin{bmatrix} 0 & -\infty & -\infty & -\infty & -\infty \\ 0 & 0 & -\infty & -\infty & -\infty \\ 0 & 0 & 0 & -\infty & -\infty \\ 0 & 0 & 0 & 0 & -\infty \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Mask matrix.

ATTENTION MATRIX.

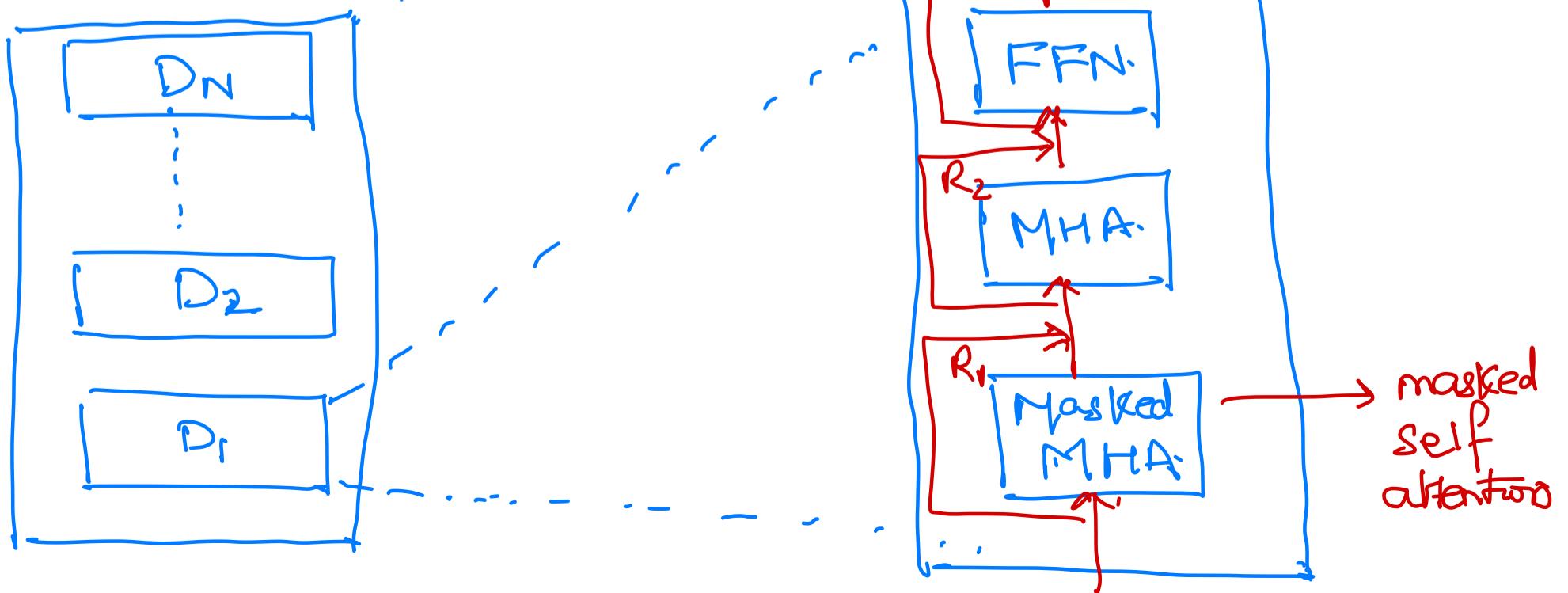
Softmax to this matrix.

\downarrow
Exponent of negative infinity is zero.

MSA
(masked self attention)

$$\begin{bmatrix} q_{v_1} k_1 & 0 & 0 & 0 & 0 \\ q_{v_2} k_1 & q_{v_2} k_2 & 0 & 0 & 0 \\ q_{v_3} k_1 & q_{v_3} k_2 & q_{v_3} k_3 & 0 & 0 \\ q_{v_4} k_1 & q_{v_4} k_2 & q_{v_4} k_3 & q_{v_4} k_4 & 0 \\ q_{v_5} k_1 & q_{v_5} k_2 & q_{v_5} k_3 & q_{v_5} k_4 & q_{v_5} k_5 \end{bmatrix}$$

This can produce the auto regressive behavior DGRs are famous for.



The magic of ChatGPT?

Turns out if keep training a bigger decoder to predict the next word then it starts to get better at all kinds of language tasks

Bigger \Rightarrow More tokens (dataset)

More blocks of transformers. (deeper)

Bigger context window to calculate.

ATTENTION. The size of the matrix.

\rightarrow This is the key difference between GPT-2/3/3.5/4 (also has images i.e multimodal).

Training for next word prediction on a bigger context window, dataset and model forces the model to learn higher order representations of data. and become good at all language tasks - sentiment, summarization, Q&A, paraphrase and much more.

\rightarrow They have also applied Reinforcement learning with human feedback on top of GPT answers to learn a policy that will pass answers that humans are likely to accept.

How does language translation work then?

This is achieved through an architectural tweak, wherein we put an encoder and decoder together.

English \rightarrow Hello how are you.

French \rightarrow Bonjour comment allez-vous

The English text is given to encoder to create the embeddings (latent representation). The decoder uses a mechanism called as **CROSS ATTENTION** to achieve this. It is depicted below visually.

