

09 Mar 2025

# SELF ATTENTION

- The self attention mechanism enables each token in a sequence to attend to every other token, dynamically adjusting its representation based on its relationship ~~with~~ to other tokens.
- This flexibility allows the transformers model to handle long-range dependencies and parallelize computations efficiently.

lets take a two sentence

(i) money bank grows

(ii) river bank flaws.

for 1st sentence

$$\text{money} = 0.7 \text{money} + 0.2 \text{bank} + 0.1 \text{grows}$$

$$\text{bank} = 0.25 \text{money} + 0.7 \text{bank} + 0.05 \text{grows}$$

$$\text{grows} = 0.1 \text{money} + 0.2 \text{bank} + 0.7 \text{grows}$$

or  
in the term of embedding

for 2nd sentence

$$\text{river} = 0.8 \text{river} + 0.15 \text{bank} + 0.05 \text{flaws}$$

$$\text{bank} = 0.2 \text{river} + 0.78 \text{bank} + 0.02 \text{flaws}$$

$$\text{flaws} = 0.4 \text{river} + 0.01 \text{bank} + 0.59 \text{flaws}$$

or  
in the term of embedding

$$e_m(\text{new}) = 0.7 e_m + 0.2 e_b + 0.1 e_g \quad \text{--- (a)}$$

$$e_b(\text{new}) = 0.25 e_m + 0.7 e_b + 0.05 e_g \quad \text{--- (b)}$$

$$e_g(\text{new}) = 0.1 e_m + 0.2 e_b + 0.7 e_g \quad \text{--- (c)}$$

→ it means.

$e_m(\text{new})$  is makes with 0.7 time  $e_m$ , 0.2 time  $e_b$  and 0.1 time  $e_g$

or we can say that coefficients (e.g. 0.7) is similarity between the  $e_m(\text{new})$  and  $e_m$



Here,

$e_m$  - embedding of money

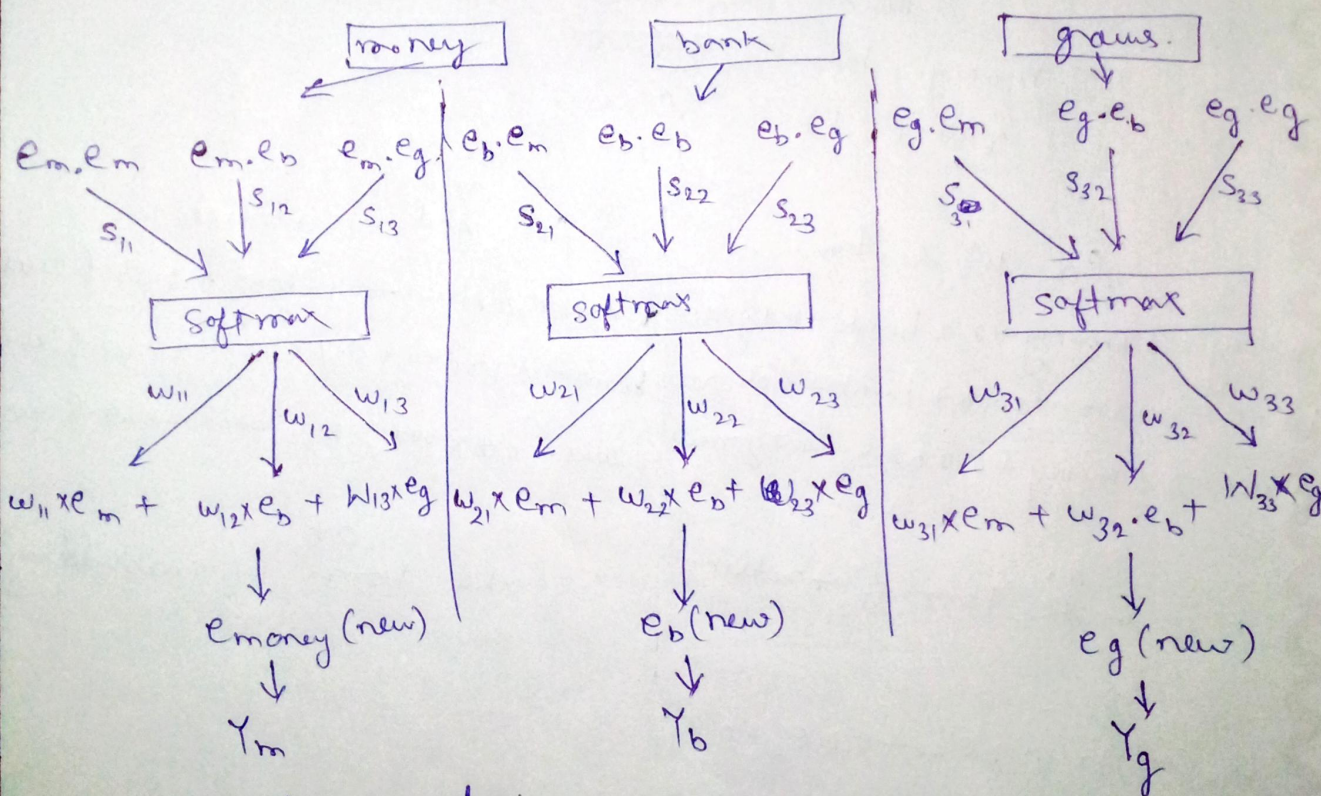
$e_b$  - embedding of bank

$e_g$  - embedding of grams

Here these embeddings are vectors in high dimension  
therefore find the similarity between two vectors  
we use the dot product

$$e_b(\text{new}) = [e_b \cdot e_m^T] e_m + [e_b \cdot e_b^T] e_b + [e_b \cdot e_g^T] e_g$$

from equation (b)



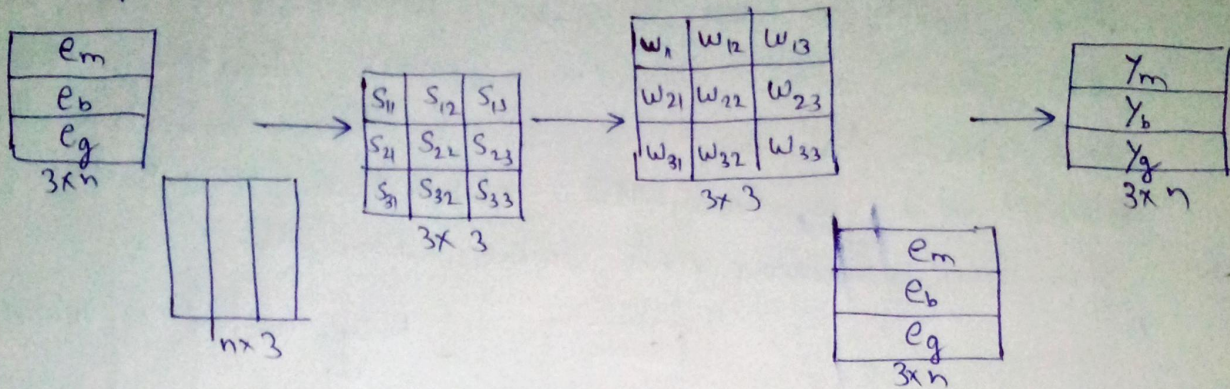
point to consider:

- This operation is a parallel operation
- There are no parameters involved
- This is generate the general contextual embedding not a task specific contextual embedding



(i) architecture of general contextual embedding of each word

OR



Here 3 number of word in sentence  
n is size of embedding vector

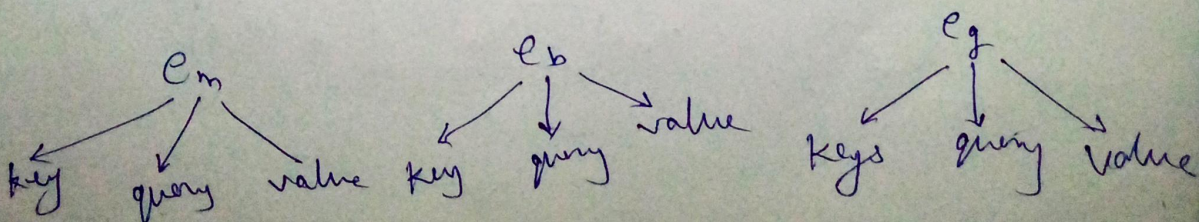
example:

in general contextual embedding  
piece of case - केक का टुकड़ा  
break a leg - टंग तोड़ दो

in task specific contextual embedding  
piece of case - बहुत आसान काम  
break a leg - शुभकामनाएँ

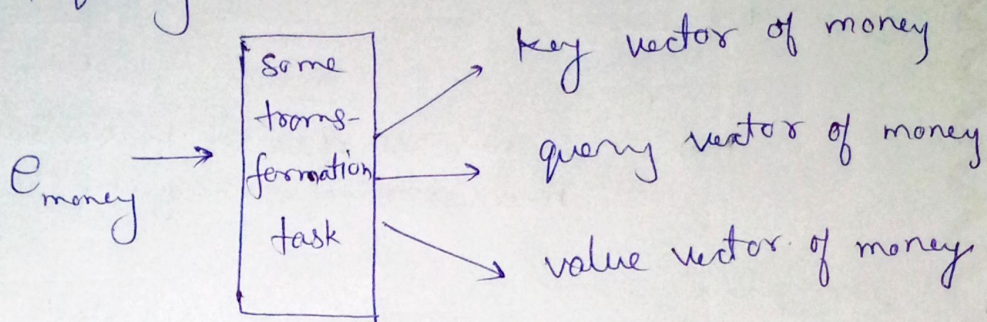
Now, due to drawback of general contextual embedding we have to improve the architecture of ~~the~~ the finding the contextual embedding with the help of some learnable parameters.

in the above word embedding of any word is use a three ways.

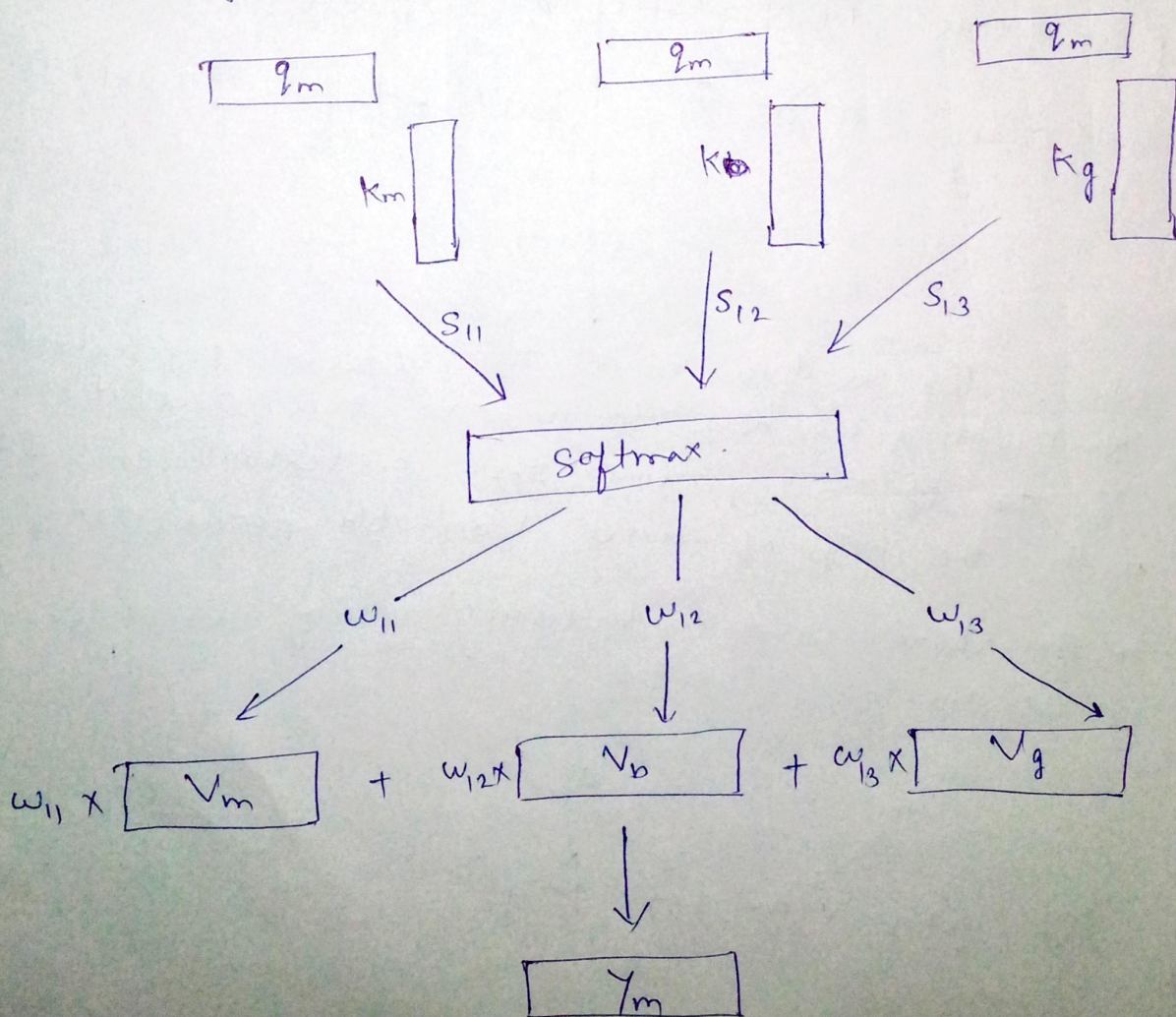




- In given sentence the embedding of given word have contain lot of given information to use the given word for use a specific function then we have to perform some transformation on word and produce the subvector. (three) example for key, query and value. example



- now find the contextual embedding of "money"





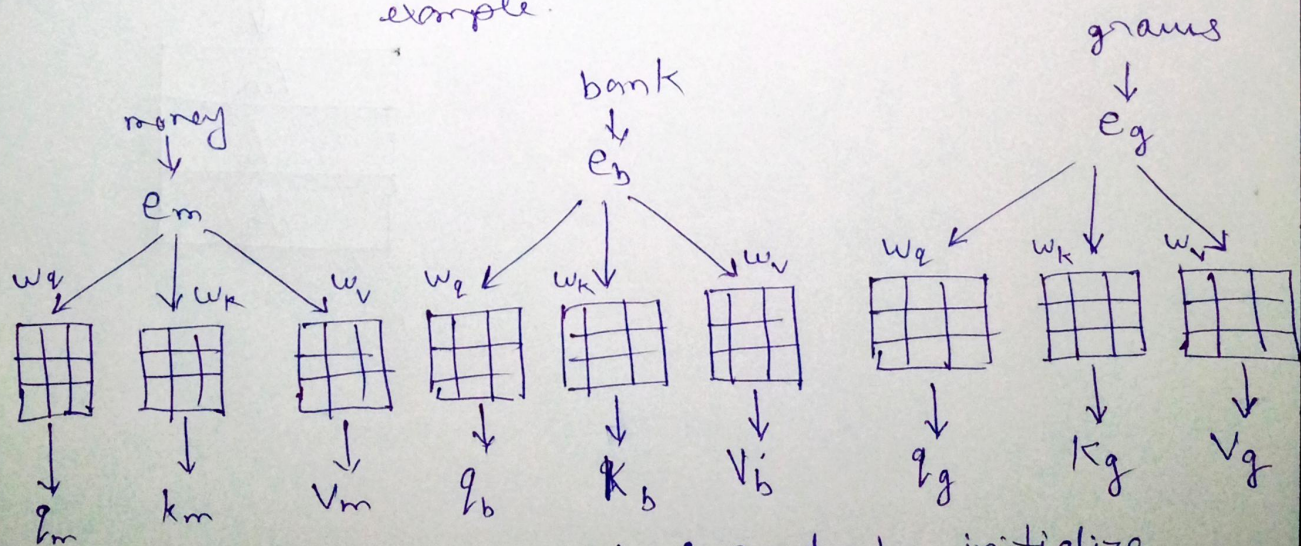
- similar way of structure of other other words also ("bank", "grow").

Here,

$q_m$  - query's vector of money word  
 $k_m$  - key's vector of money word.  
 $k_b$  - key's vector of bank word.  
 $k_g$  - key's vector of grow word.  
 $v_m$  - value's vector of money word.  
 $v_b$  - value's vector of ~~bank~~ bank word.  
 $v_g$  - value's vector of grow word.

- Now let's find what type of transformation/operation is perform to find the these vectors of each word (key, query, value).

- magnitude (scaling)
- linear transformation ✓  
(multiply vector ( $e_m, e_b, e_g$ ) with matrix)  
example.



inside the matrix is randomly initialize.  
value they will update the value for further iterations



(ii) Architecture of task specific contextual embedding of each word

