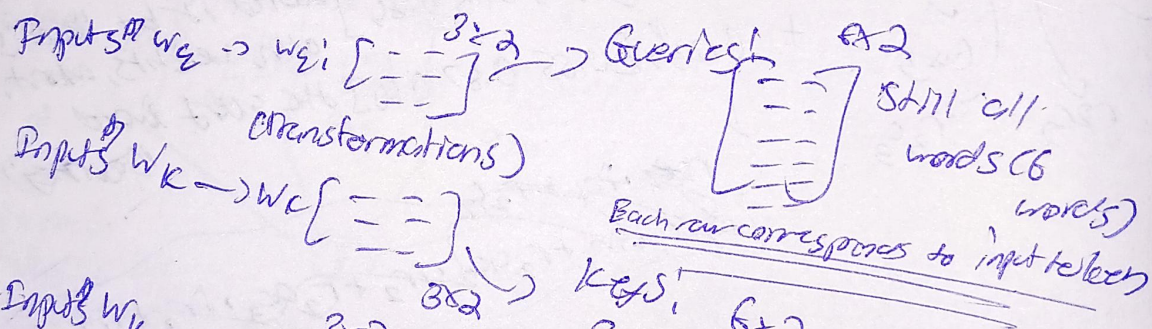


Self-Attention w/ learnable weights  
 Also, scaled dot-product attention

- ① Compute context vectors as weighted sums over input vectors specific to a certain input element
- ② Introduce weight matrices that are updated during model training
- ③ Crucial that model can learn to produce "good" context vectors
- ④  $W_q, W_k$  and  $W_v$  <sup>2, 1, 1</sup> query, key, value

These 3 matrices are used to project embedded input tokens  $x^{(i)}$  into query, key and value vectors

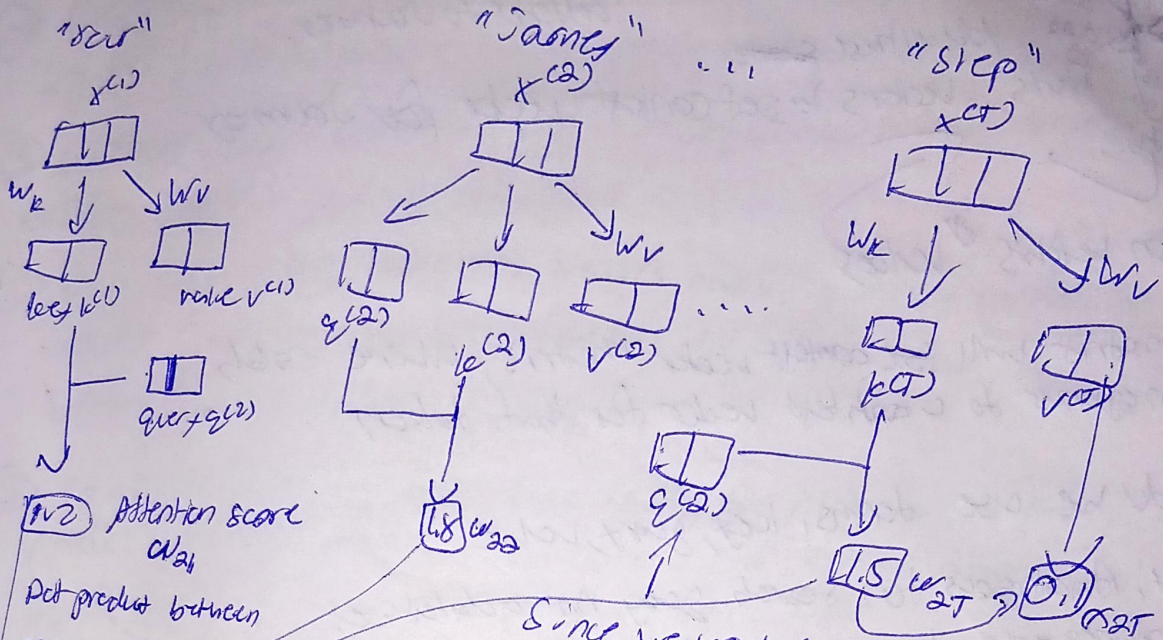


It means it undergoes optimization

After we get  $W_q, W_k, W_v$ , we do not need to look at input matrix again; has been transformed into 3 matrices.



⑤ Compute only one context vector  $z^{(2)}$  for instruction purposes  
 ↳ computing attention scores:



② Attention score  $a_{21}$

Dot product between  
query and key  
vectors

Queries ~~key~~  
keys

Since we want to compute the context vector for second input token,  $q^{(2)}$  is derived from that 2nd input token

↳ Then we compute Attention weights

③ Attention weights: (Scale by  $\sqrt{d_{keys}}$ ) (2) in this case

scale attention scores by  $\sqrt{d_{keys}}$  before softmax

Attn score  $\rightarrow$  scale by  $\sqrt{d_{keys}}$   $\rightarrow$  softmax  $\rightarrow$  Attention weights

④ Context vector

Attention weights  
↓  
context vector

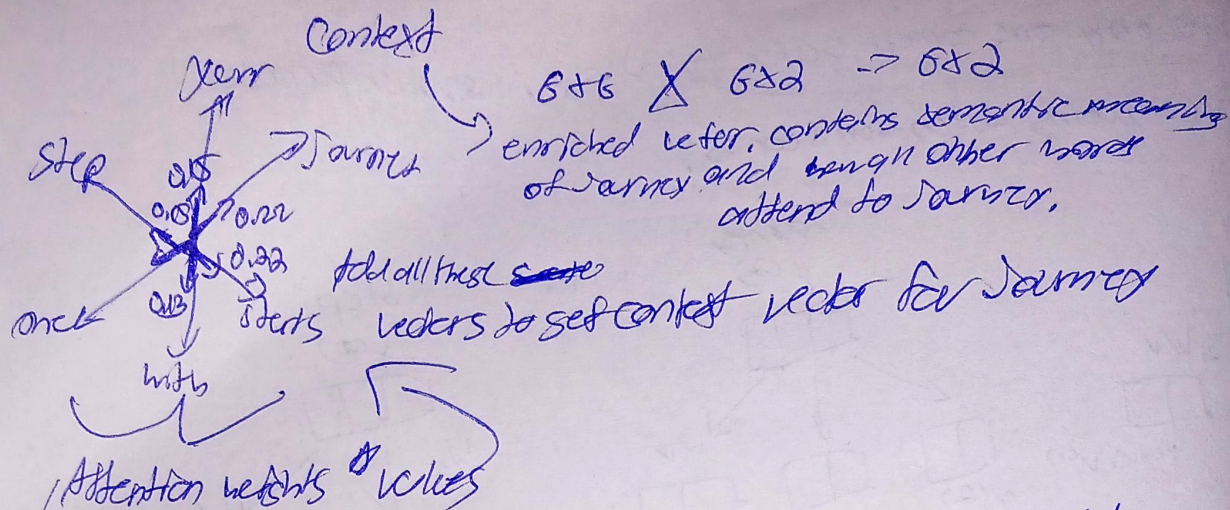
Next step is multiplying each value vector w/ its respective attention weight and then summing them to obtain the context vector

$[0.3 | 0.8]$

context vector  $z^{(2)}$

Multiply value vector w/ attention weight, which is then sent to scale context vector





\* Why do we use terms: key, query, value

Query: Analogous to search query in a database.

It represents the current token the model focuses on

key: In attention mechanism, each item in input sequence has a key. Keys are used to match w/ the query

value: It represents the output context or representation of the input items. Once model determines which keys (which parts of the input) are most relevant to the query (current focus item), it retrieves the corresponding values.