

**Understanding Transformers:
A step-by-step hand simulation
to translate a sentence
from English to Hindi**

- How information changes through each layer of a basic transformer.
- Assumptions: weights are already optimized, so no backpropagation required

let's	1	0	0	0
to	0	1	0	0
go	0	0	1	0
<EOS>	0	0	0	1

English vocabulary = I_4 =

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

Word embedding:

$$I_4 \cdot \begin{bmatrix} 1.87 & 0.09 \\ -1.45 & 1.50 \\ -0.78 & 0.27 \\ 2.21 & -0.64 \end{bmatrix} = \begin{bmatrix} 1.87 & 0.09 \\ -1.45 & 1.50 \\ -0.78 & 0.27 \\ 2.21 & -0.64 \end{bmatrix} \begin{matrix} \text{let's} \\ \text{to} \\ \text{go} \\ \text{<EOS>} \end{matrix}$$

Sentence = [let's, go] = $\begin{bmatrix} 1.87 & 0.09 \\ -0.78 & 0.27 \end{bmatrix} \begin{matrix} \text{let's} \\ \text{go} \end{matrix}$

- Weights are optimised using backpropagation
- The process of optimising the weights is called training

But position matters!

[do , I , like , this]
[I , do , like , this]



different meanings, so just embedding words as vectors won't work

we need to somehow embed positions

<do, 0>; <I, 1>; <like, 2>; <this, 3>

<I, 0>; <do, 1>

so that the same words could be embedded differently if their position changes their meaning

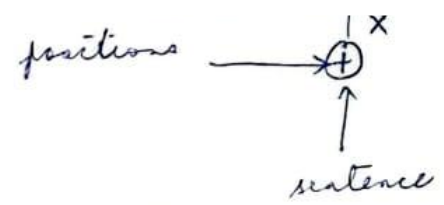
Positional embedding:

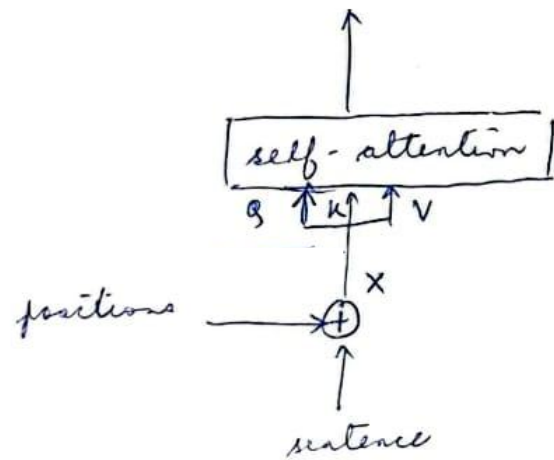
$$0 = \langle 0, 1 \rangle$$

$$1 = \langle -0.9, 0.4 \rangle$$

$$\begin{array}{l} X \\ = \\ = \end{array} \begin{array}{l} \text{sentence} \\ \begin{array}{|c|c|} \hline 1.87 & 0.09 \\ \hline -0.78 & 0.27 \\ \hline \end{array} \end{array} \begin{array}{l} \\ \text{let's} \\ \text{go} \end{array} \begin{array}{l} + \\ + \end{array} \begin{array}{l} \text{embedded positions} \\ \begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline -0.9 & 0.4 & 1 \\ \hline \end{array} \end{array}$$

$$\begin{array}{l} X \\ = \end{array} \begin{array}{|c|c|} \hline 1.87 & 1.09 \\ \hline -1.68 & 0.67 \\ \hline \end{array} \begin{array}{l} \langle \text{let's}, 0 \rangle \\ \langle \text{go}, 1 \rangle \end{array}$$





Self-attention keeps track of the relationships among words

The **stew** was cooked on the **stove**, and **it** tasted good

1. Encoder

Self-attention:

$X \cdot W_q = Q$: query : another way to represent X : what everyone is looking for

$X \cdot W_k = K$: key : yet another way to represent X : what everyone can offer

X		W_q		Q												
<table border="1"><tr><td>1.87</td><td>1.09</td></tr><tr><td>-1.68</td><td>0.67</td></tr></table>	1.87	1.09	-1.68	0.67	.	<table border="1"><tr><td>1.1</td><td>0.6</td></tr><tr><td>-2.8</td><td>2.4</td></tr></table>	1.1	0.6	-2.8	2.4	=	<table border="1"><tr><td>-0.995</td><td>3.74</td></tr><tr><td>-3.724</td><td>0.6</td></tr></table>	-0.995	3.74	-3.724	0.6
1.87	1.09															
-1.68	0.67															
1.1	0.6															
-2.8	2.4															
-0.995	3.74															
-3.724	0.6															
X		W_k		K												
<table border="1"><tr><td>1.87</td><td>1.09</td></tr><tr><td>-1.68</td><td>0.67</td></tr></table>	1.87	1.09	-1.68	0.67	.	<table border="1"><tr><td>-1.7</td><td>0.5</td></tr><tr><td>-1.4</td><td>0.9</td></tr></table>	-1.7	0.5	-1.4	0.9	=	<table border="1"><tr><td>-4.71</td><td>1.92</td></tr><tr><td>1.92</td><td>-0.24</td></tr></table>	-4.71	1.92	1.92	-0.24
1.87	1.09															
-1.68	0.67															
-1.7	0.5															
-1.4	0.9															
-4.71	1.92															
1.92	-0.24															

Self-attention:

Compare Q and K : $Q \cdot K^T$

why K^T ?: if $Q(4 \times 2)$ and $K(4 \times 2)$ then Q and K should be multipliable

$$Q \cdot K^T = \begin{bmatrix} 11.84 & -2.79 \\ 18.67 & -7.28 \end{bmatrix}$$

$$\text{softmax}(Q \cdot K^T) = \begin{bmatrix} 9.9e-1 & 1.05e-7 \\ 9.9e-1 & 2.01e-9 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} \text{ (approx.)}$$

Self-attention:

$X \cdot W_v = V$: value : yet another way to represent X : **what everyone is worth**

$$\begin{array}{c} X \\ \boxed{\begin{array}{cc} 1.87 & 1.09 \\ -1.68 & 0.67 \end{array}} \end{array} \cdot \begin{array}{c} W_v \\ \boxed{\begin{array}{cc} 1.5 & -1.0 \\ -0.3 & -0.2 \end{array}} \end{array} = \begin{array}{c} V \\ \boxed{\begin{array}{cc} 2.478 & -2.088 \\ -2.631 & 1.606 \end{array}} \end{array}$$

$\text{self_attention} = \text{softmax}(Q \cdot K^T) \cdot V$

$$= \begin{array}{c} \boxed{\begin{array}{cc} 1 & 0 \\ 1 & 0 \end{array}} \end{array} \cdot \begin{array}{c} \boxed{\begin{array}{cc} 2.478 & -2.088 \\ -2.631 & 1.606 \end{array}} \end{array} = \begin{array}{c} \boxed{\begin{array}{cc} 2.48 & -2.08 \\ 2.48 & -2.08 \end{array}} \end{array}$$

This could possibly suggest a one-word-translation from English to Hindi

The weights used to calculate self attention are the same for “lets” and “go”.

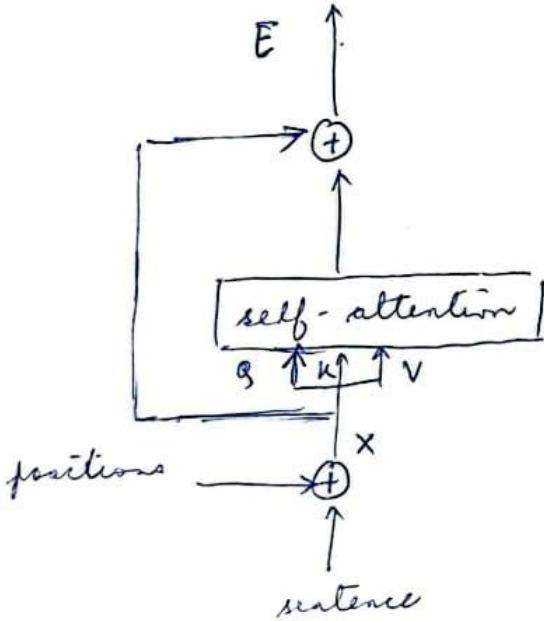
No matter how many words are input to the transformer, we reuse the same sets of weights for each word.

We can hence calculate all Q , K and V for all words at the same time. Was not possible with RNNs.

Residual connection:

$$E = \begin{bmatrix} 2.48 & -2.08 \\ 2.48 & -2.08 \end{bmatrix} + \begin{bmatrix} 1.87 & 1.09 \\ -1.68 & 0.67 \end{bmatrix} = \begin{bmatrix} 4.35 & -0.998 \\ 0.798 & -1.418 \end{bmatrix}$$

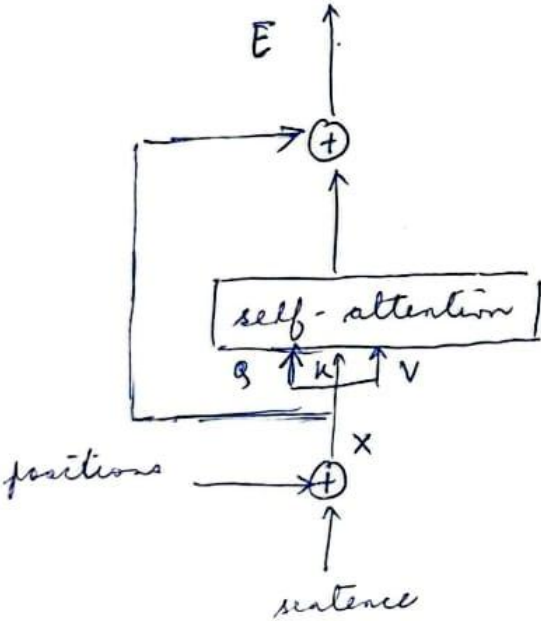
encoder output



Residual connection:

$$\begin{aligned}
 E &= \text{self_attention} + X \\
 &= \begin{bmatrix} 2.48 & -2.08 \\ 2.48 & -2.08 \end{bmatrix} + \begin{bmatrix} 1.87 & 1.09 \\ -1.68 & 0.67 \end{bmatrix} = \begin{bmatrix} 4.35 & -0.998 \\ 0.798 & -1.418 \end{bmatrix}
 \end{aligned}$$

encoder output



Residual connections help in training: SA layer can establish a relationship with input words without having to preserve all the functions that led to it.

2. Decoder

इधर (<i>idhar</i>)	1	0	0	0
चलो (<i>chalo</i>)	0	1	0	0
हम (<i>hum</i>)	0	0	1	0
<EOS>	0	0	0	1

Hindi vocabulary = I_4 =

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

Word embedding:

$I_4 \cdot$

-2.27	2.54
0.04	1.97
-0.77	-0.75
2.70	-1.34

 =

-2.27	2.54
0.04	1.97
-0.77	-0.75
2.70	-1.34

 इधर (*idhar*)
चलो (*chalo*)
हम (*hum*)
<EOS>

The decoder generates text iteratively, i.e., it predicts the next word based on the previous word(s). So unlike the encoder, where the entire sentence was processed together ([let's, go]), here we will process one word at a time.

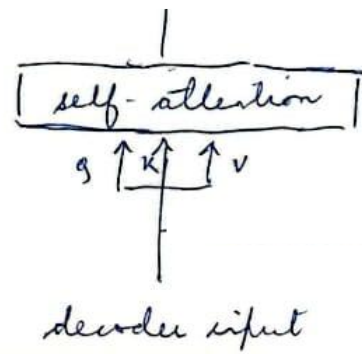
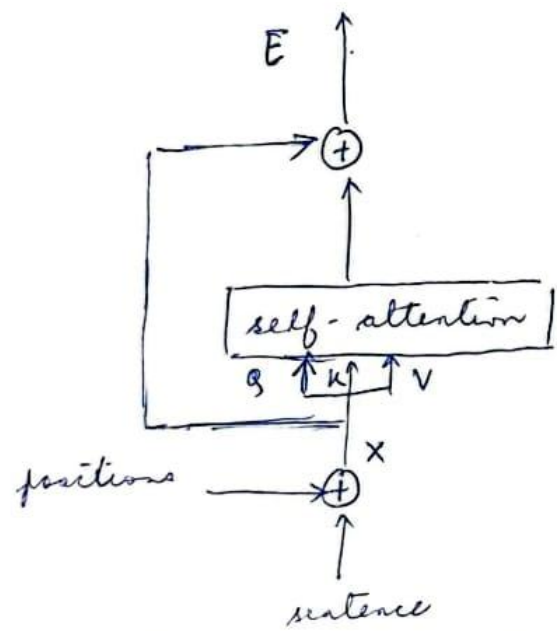
We assume the first word to be <EOS> in Hindi (or any target language), since that would be the last word of the (imaginary) previous sentence, preceding the first word of the current sentence.

-2.27	2.54	इधर (<i>idhar</i>)
0.04	1.97	चलो (<i>chalo</i>)
-0.77	-0.75	हम (<i>hum</i>)
2.70	-1.34	<EOS>

$$\begin{aligned}
 \langle \text{EOS} \rangle &= \langle 2.70, -1.34 \rangle \\
 \langle \text{EOS}, 0 \rangle &= \langle \text{EOS} \rangle + \langle 0 \rangle = \langle 2.70, -1.34 \rangle + \underbrace{\langle 0, 1 \rangle} \\
 &= \langle 2.70, -0.34 \rangle
 \end{aligned}$$

recall positional embedding? we use the same embeddings here.

0	1	0
-0.9	0.4	1



Decoder self-attention:

$$W_{q,sa} = \begin{bmatrix} 0.4 & 0.4 \\ -0.3 & 0.1 \end{bmatrix}$$

$$W_{k,sa} = \begin{bmatrix} 0.4 & -0.7 \\ -0.4 & -0.3 \end{bmatrix}$$

$$W_{v,sa} = \begin{bmatrix} -1.1 & -0.7 \\ -0.4 & 1.3 \end{bmatrix}$$

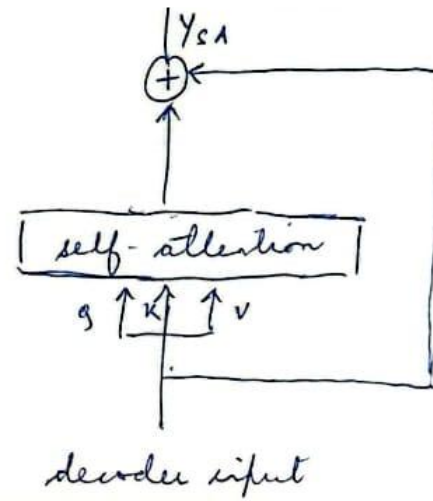
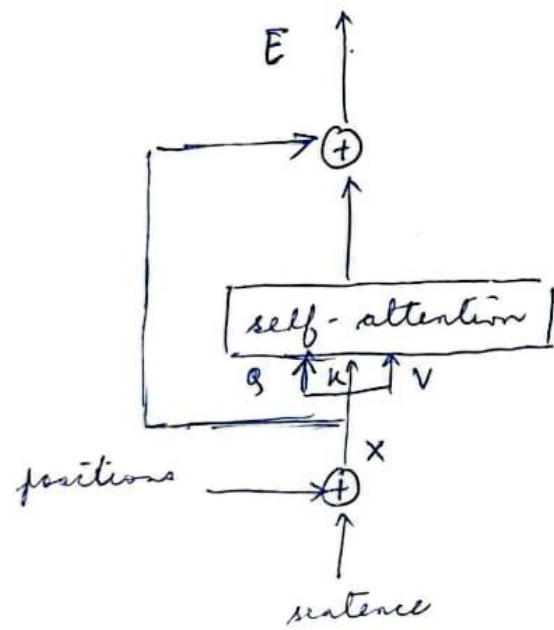
$$\langle \text{EOS}, 0 \rangle \cdot W_{q,sa} = Q_{sa}$$

$$\langle \text{EOS}, 0 \rangle \cdot W_{k,sa} = K_{sa}$$

$$\langle \text{EOS}, 0 \rangle \cdot W_{v,sa} = V_{sa}$$

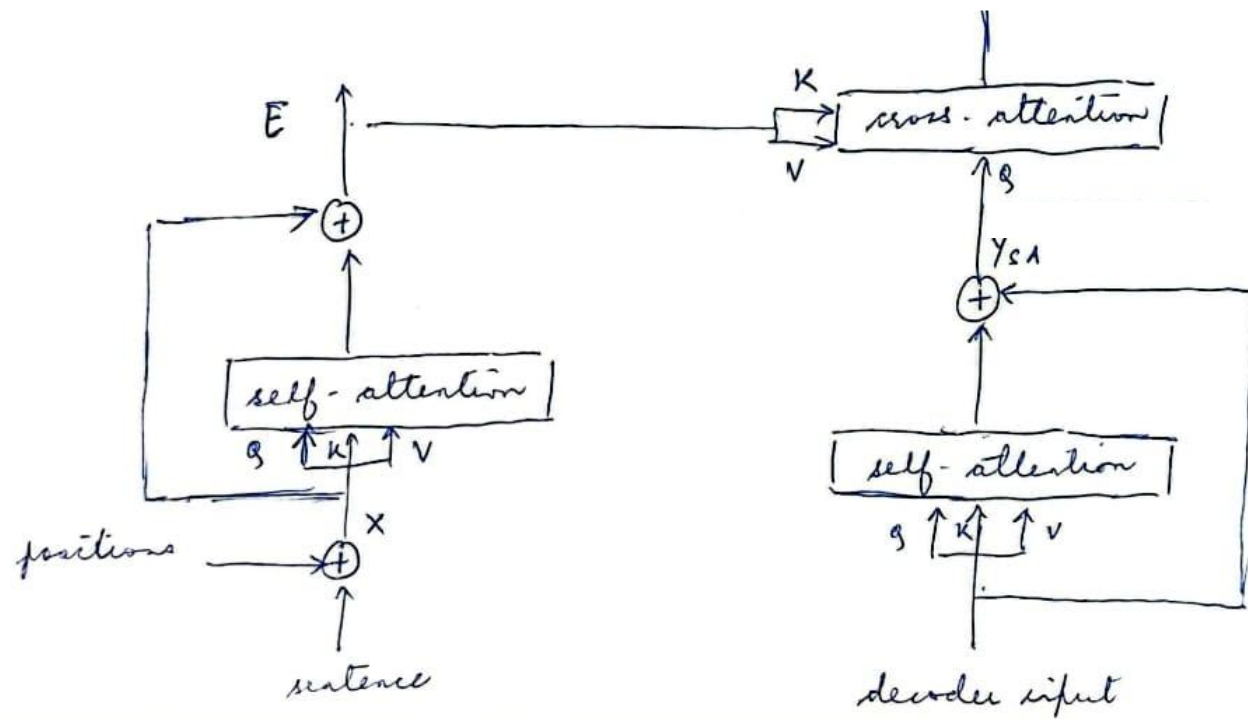
$$\text{decoder_self_attention} = \text{softmax}(Q_{sa} \cdot K_{sa}^T) \cdot V_{sa}$$

$$= \langle -2.834, -2.332 \rangle$$



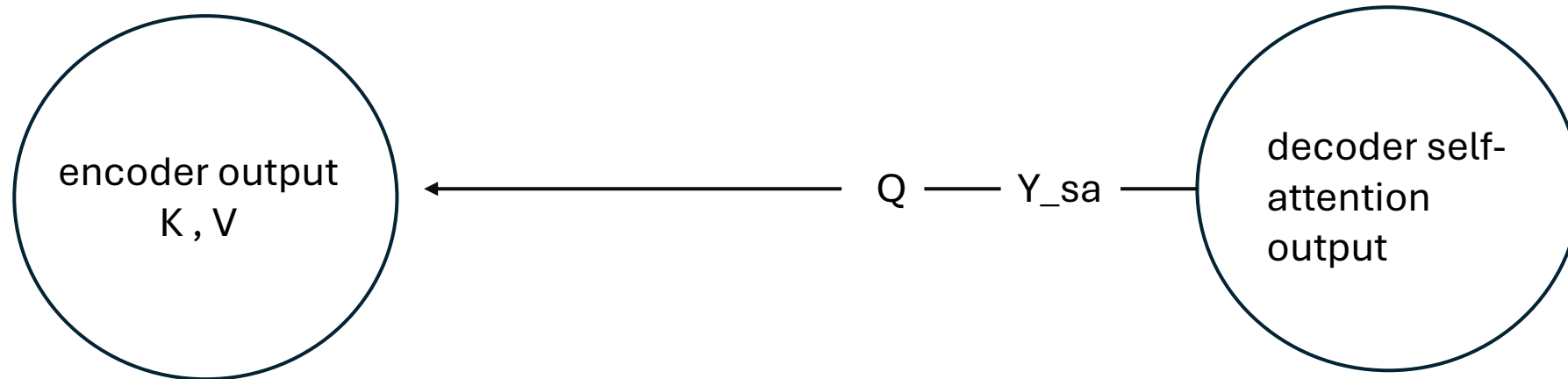
Residual network:

$$\begin{aligned} Y_{sa} &= \text{decoder_self_attention} + \langle \text{EOS}, 0 \rangle \\ &= \langle -2.834, -2.332 \rangle + \langle 2.70, -0.34 \rangle \\ &= \langle -0.134, -2.672 \rangle \end{aligned}$$



Decoder cross-attention:

We have now learnt the meaning of the current word w.r.t the other words in Hindi (through self attention)
But what does the current word mean w.r.t the English that was learnt by the encoder?



Decoder cross-attention:

$$W_{q,ca} = \begin{bmatrix} 1.5 & -0.3 \\ 0.3 & -1.0 \end{bmatrix}$$

$$W_{k,ca} = \begin{bmatrix} -1.1 & 0.3 \\ -0.3 & -0.8 \end{bmatrix}$$

$$W_{v,ca} = \begin{bmatrix} 1.1 & 0.6 \\ -1.2 & -0.5 \end{bmatrix}$$

Query from decoder self attention output:

$$Y_{sa} \cdot W_{q,ca} = Q_{ca} = \langle -1.0026, 2.7122 \rangle$$

Recall E, the encoder output?

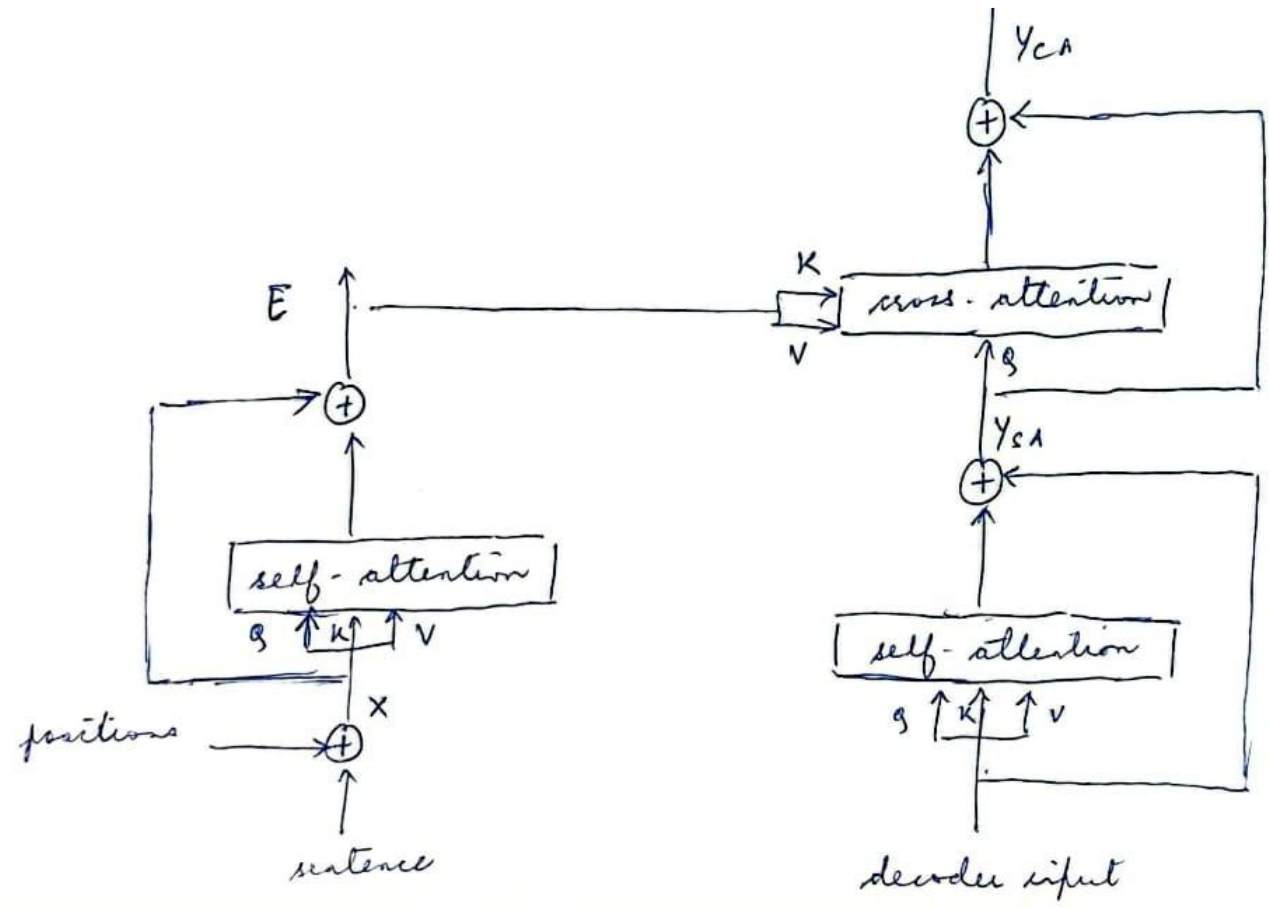
$$\begin{bmatrix} 4.35 & -0.998 \\ 0.798 & -1.418 \end{bmatrix}$$

Key and Value from encoder output:

$$E \cdot W_{k,ca} = K_{ca} = \begin{bmatrix} -4.48 & 2.10 \\ -0.45 & 1.37 \end{bmatrix}$$

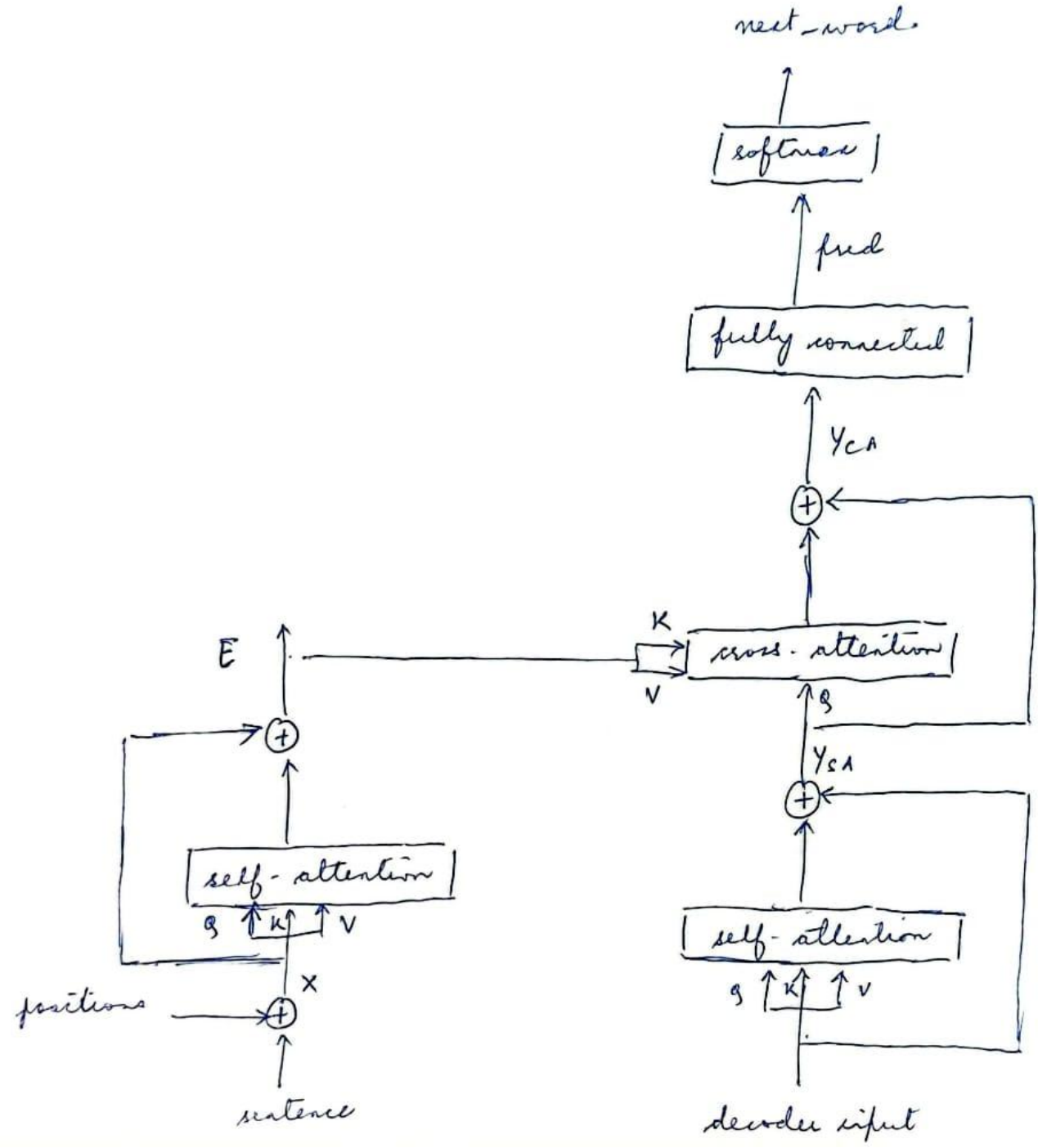
$$E \cdot W_{v,ca} = V_{ca} = \begin{bmatrix} 5.98 & 3.11 \\ 2.58 & 1.19 \end{bmatrix}$$

$$\text{decoder_cross_attention} = \text{softmax}(Q_{ca} \cdot K_{ca}^T) \cdot V_{ca} = \langle 5.97, 3.10 \rangle$$



Residual network:

$$\begin{aligned} Y_{ca} &= \text{decoder_cross_attention} + Y_{sa} \\ &= \langle 5.97, 3.10 \rangle + \langle -0.134, -2.672 \rangle \\ &= \langle 5.84, 0.43 \rangle \end{aligned}$$



Fully-connected layer

$$W = \begin{bmatrix} -0.6 & 0.8 & -0.1 & -1.0 \\ -2.0 & -0.9 & -1.1 & 1.6 \end{bmatrix}$$

$$\text{bias} = \langle -0.6, 1.4, -2.5, 0.5 \rangle$$

$$\begin{aligned} \text{pred} &= Y_{ca} \cdot W + \text{bias} \\ &= \langle -4.97, 5.68, -3.56, -4.65 \rangle \end{aligned}$$

$$\begin{aligned} \text{next_word} &= \text{softmax}(\text{pred}) \\ &= \langle 2.4e-5, 9.99e-1, 9.7e-5, 3.3e-5 \rangle = \langle 0, 1, 0, 0 \rangle \text{ (approx.)} \end{aligned}$$

चलो (*chalo*)

Now, चलो (*chalo*) enters the decoder, and the same decoder operations are repeated.

```
[46] chalo = hindi_words_embedded[1]
      chalo
```

```
↳ array([0.04, 1.97])
```

```
[47] chalo = chalo + W_pos_embedding[1] # 1st index
      chalo = chalo.reshape(1,2)
      chalo
```

```
↳ array([[ -0.86,  2.37]])
```

```
[48] Q_chalo = chalo @ W_q_dsa
      K_chalo = chalo @ W_k_dsa
      V_chalo = chalo @ W_v_dsa
```

```
[49] e = Q_chalo @ K_chalo.T
      softmax_output = np.exp(e) / np.sum(np.exp(e))
      softmax_output = np.array(softmax_output)
      self_attention_chalo = softmax_output @ V_chalo
      self_attention_chalo
```

```
↳ array([[ -2.000e-03,  3.683e+00]])
```

```
[50] chalo = self_attention_chalo + chalo
      chalo
```

```
↳ array([[ -0.862,  6.053]])
```

```
[51] Q_chalo = chalo @ W_q_dca
      dot_product_ca = Q_chalo @ K_encoder_output.T
      similarity = softmax(dot_product_ca)
```

```
[52] cross_attention = similarity @ V_encoder_output
```

```
[53] chalo = cross_attention + chalo
```

```
[54] pred = chalo @ W_fcl + bias
```

```
[55] np.argmax(softmax(pred))
```

```
↳ 3
```

3 = EOS

Next word predicted is <EOS>
Translation finished.