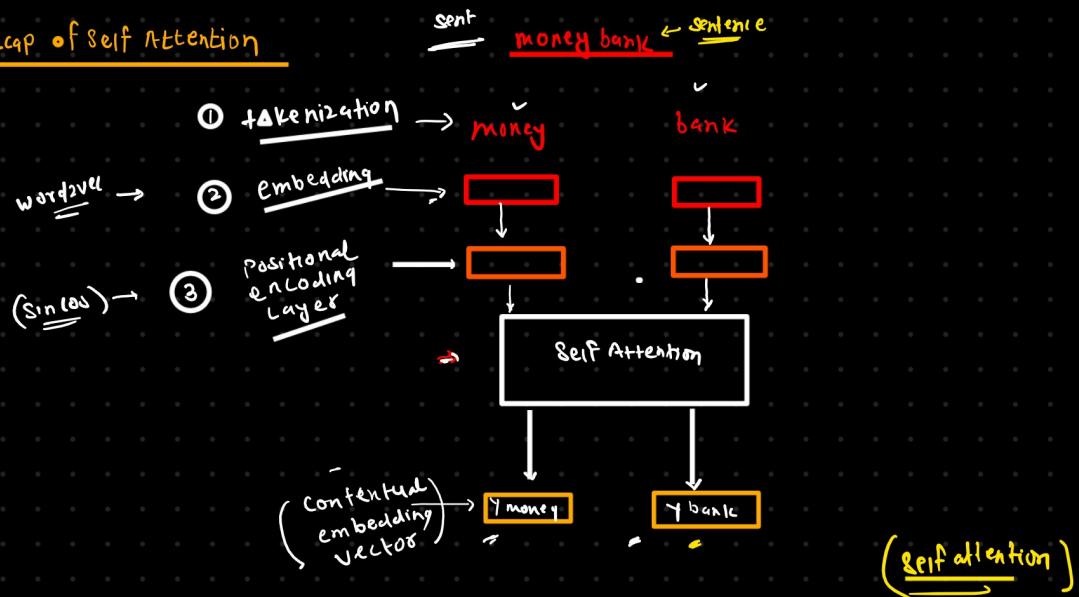
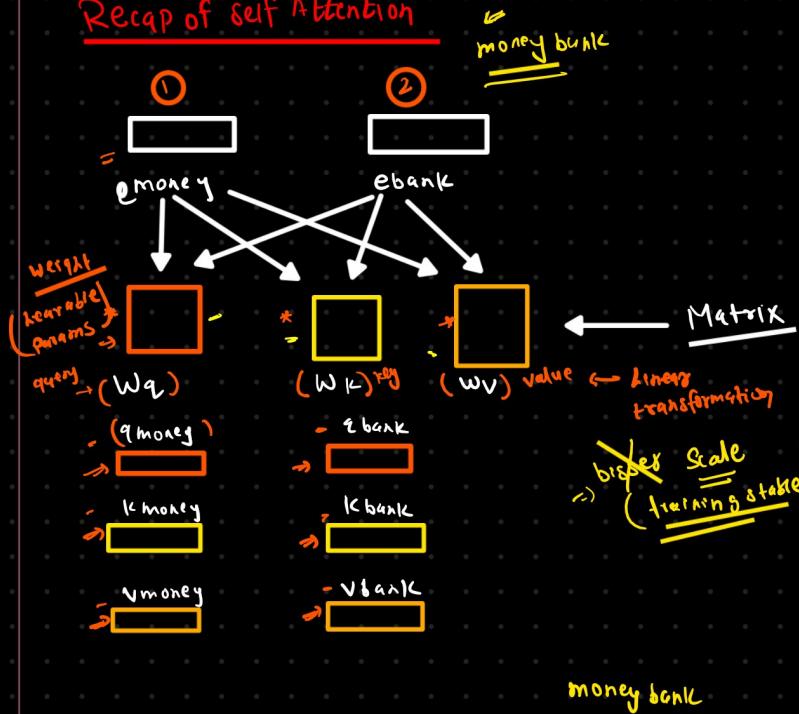


(Multhead, Cross & masked self Attention)

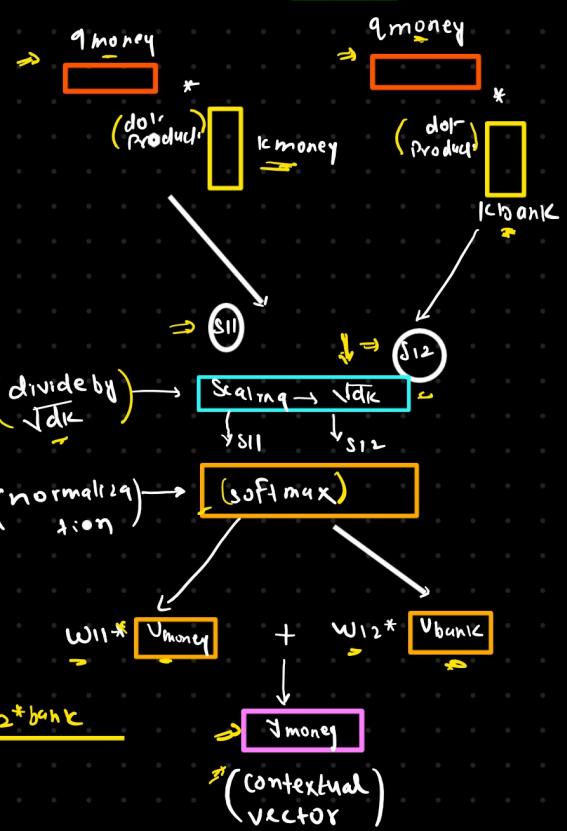
Recap of Self Attention



Recap of Self Attention



Let's look into money word



ambiguity, 1 sent \rightarrow interpretation \rightarrow multiple way

required

Now why we understand Multhead self Attention

- ① We use it for ambiguous sentence NLP
- ② self attention with single head only can capture one perspective of sentence
- ③ in NLP one sentence have multiple Perspective.

For ex

① (Summarization) \rightarrow sent \Rightarrow multiple \rightarrow (decode)

② (Few other sentence) \rightarrow The chicken is ready to eat \rightarrow man saw the astronomer with a telescope

In electrical engineering, a transformer is a passive component that transfers electrical energy from one electrical circuit to another circuit, or multiple circuits. A varying current in any coil of the transformer produces a varying magnetic flux in the transformer's core, which induces a varying electromotive force (EMF) across any other coils wound around the same core.

Electrical energy can be transferred between separate coils without a metallic (conductive) connection between the two circuits. Faraday's law of induction, discovered in 1831, describes the induced voltage effect in any coil due to

↓ transformer → multiple → (single attention)

(Summary 1) person 1

multihead

→ entire context of data X

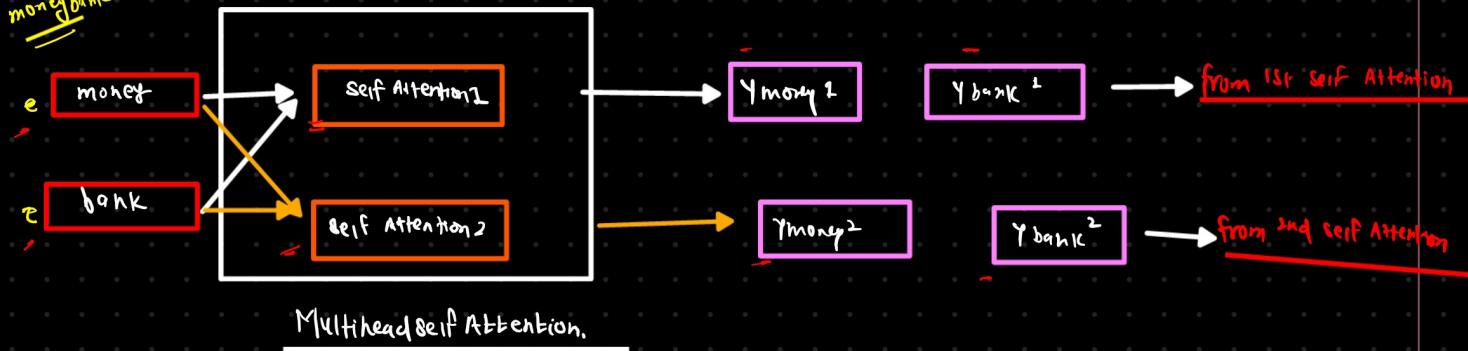
A transformer is a passive electrical component that transfers energy between circuits via electromagnetic induction. A varying current in one coil induces a magnetic flux in the core, generating an electromotive force (EMF) in other coils.

Faraday's law of induction explains this phenomenon, enabling energy transfer without direct conductive connection.

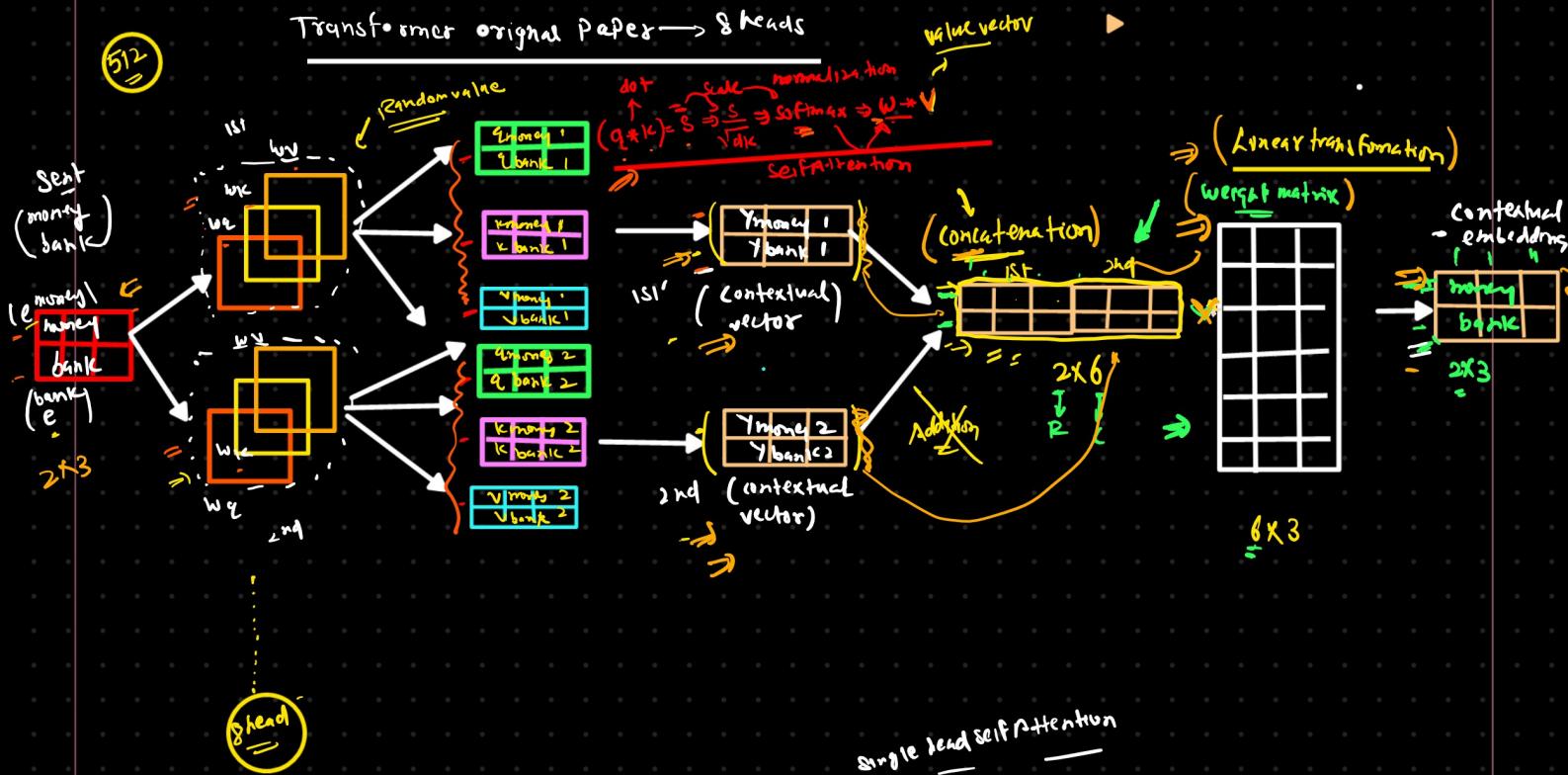
Summary 2 person 2

Transformers use electromagnetic induction to transfer electrical energy between circuits. A changing current in one coil creates a magnetic flux, inducing voltage in another coil. This process, based on Faraday's law, allows energy transfer without physical connections.

memory bank



Transformer original paper → 8 heads



8 head = single head self attention

① single head self attention

② multihead self attention

③ multiples of matrix

④ more than 1

⑤ heads ⇒ hyperparameter

one head → single set of K, Q, V matrix

Original research paper → 8 head ⇒

= $\begin{bmatrix} q \\ k \\ v \end{bmatrix}$ } single set one head

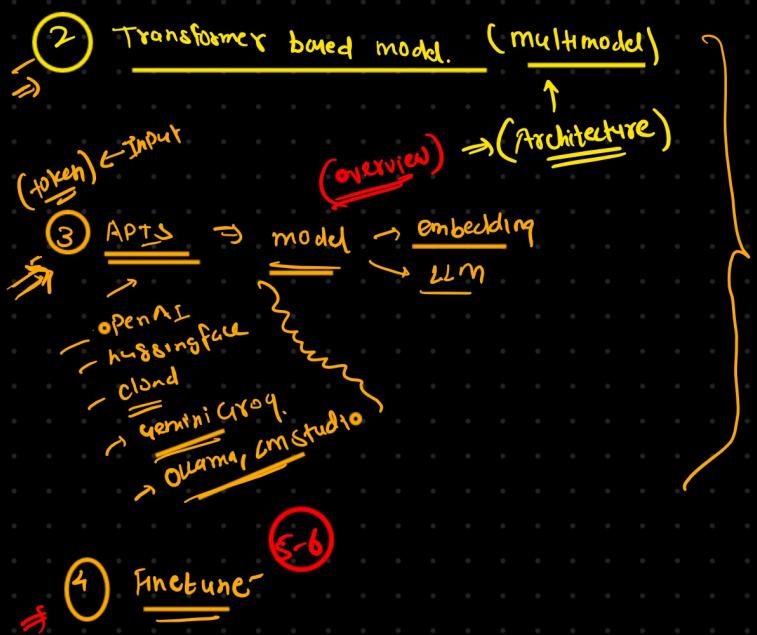
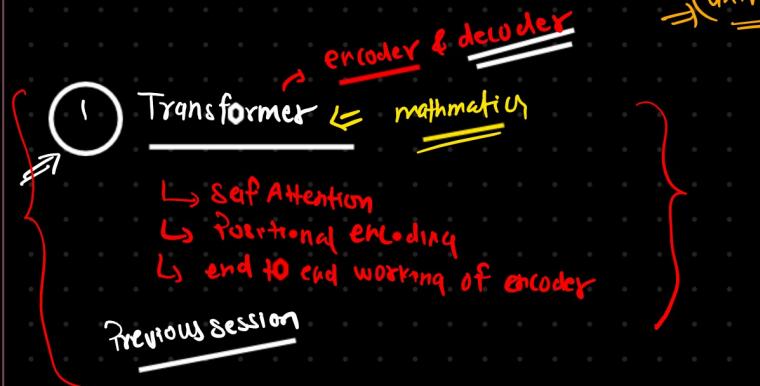
? multihead self attention

↓ multiples of matrix

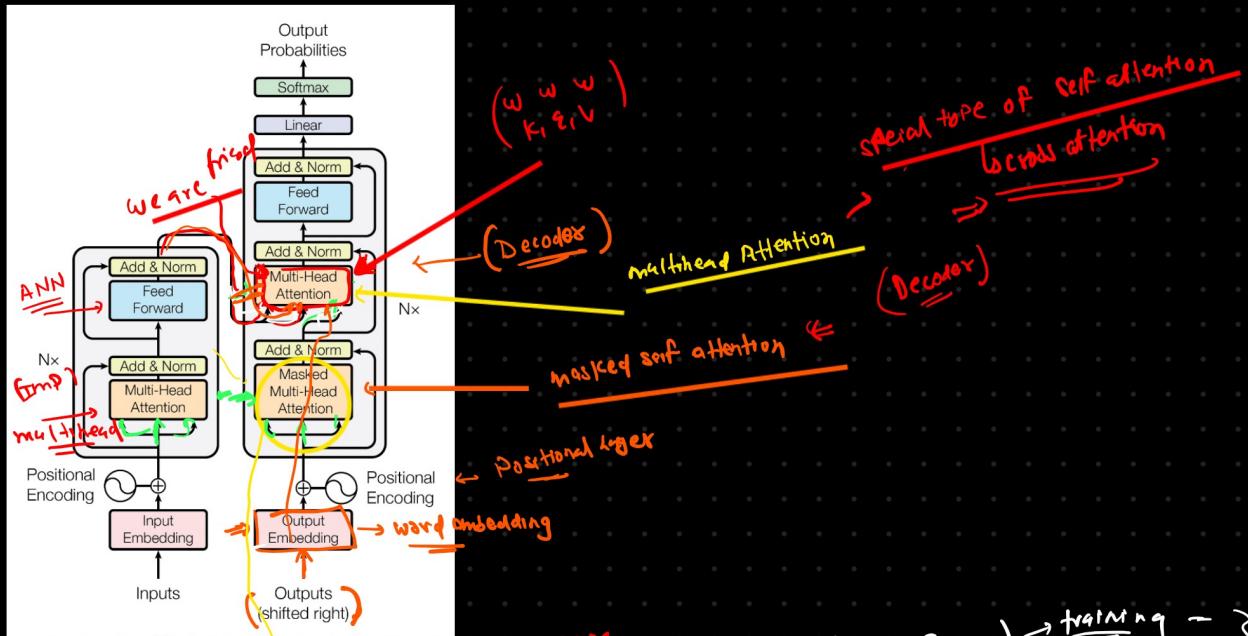
more than 1

→ heads ⇒ hyperparameter

Lets understand about the cross encoder



- 2 Cross Attention -
3 Masked attention -

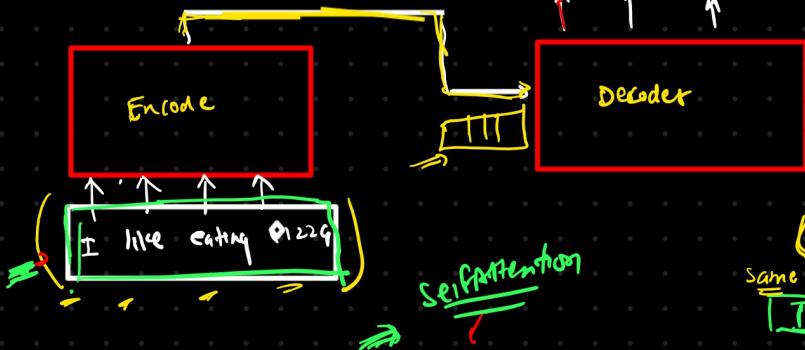


Transformers \Rightarrow translation \Rightarrow en. (english) \rightarrow hi. (hindi)

(I like eating pizza) \rightarrow प्रेत (प्रेत)

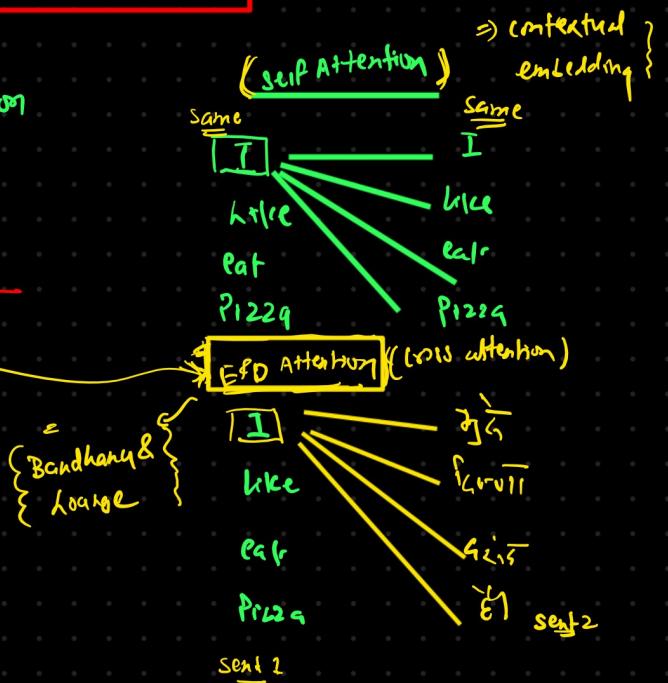
प्रेत \rightarrow next word (next word)

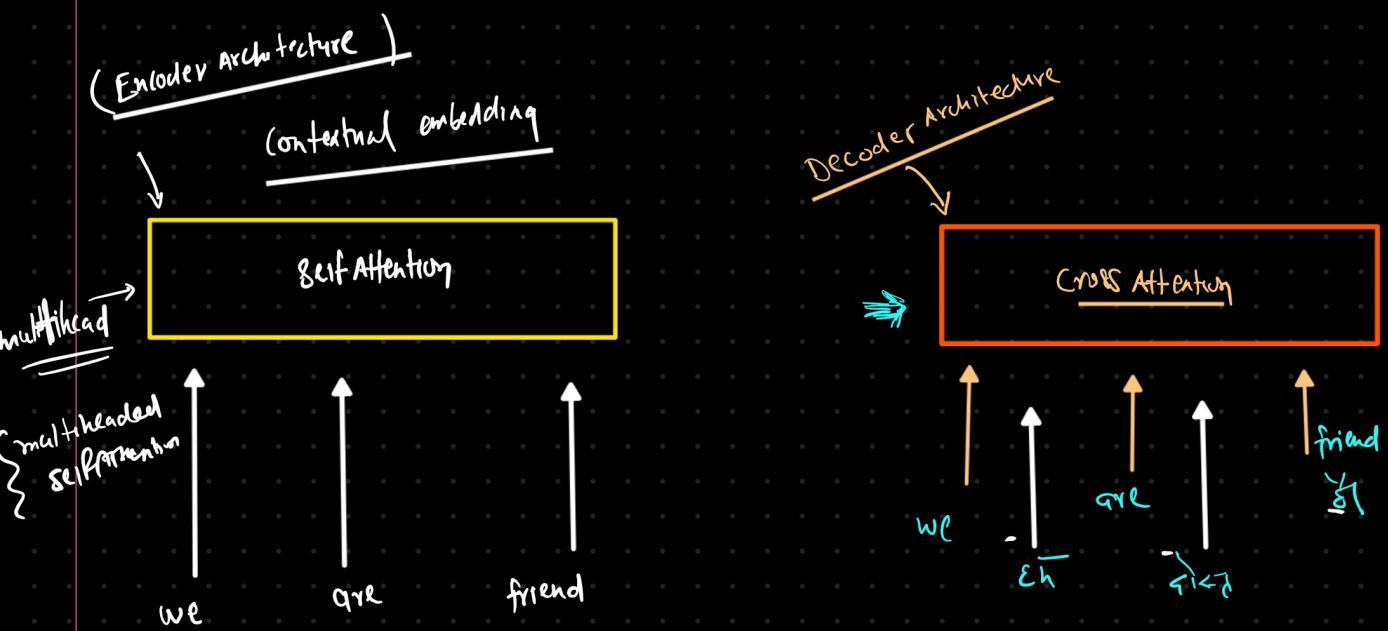
प्रेत \rightarrow बिना ?



- 1) what generated till now
2) input sentence

CROSS attention
transformer





SelfAttention

$$\begin{aligned}
 ce_we &= w_{11} * e_we + w_{12} * e_are + w_{13} * e_friend \\
 ce_are &= w_{21} * e_we + w_{22} * e_are + w_{23} * e_friend \\
 ce_friend &= w_{31} * e_we + w_{32} * e_are + w_{33} * e_friend
 \end{aligned}$$

selfattention

	we	are	friend
we	*	*	*
are	*	*	*
friend	*	*	*

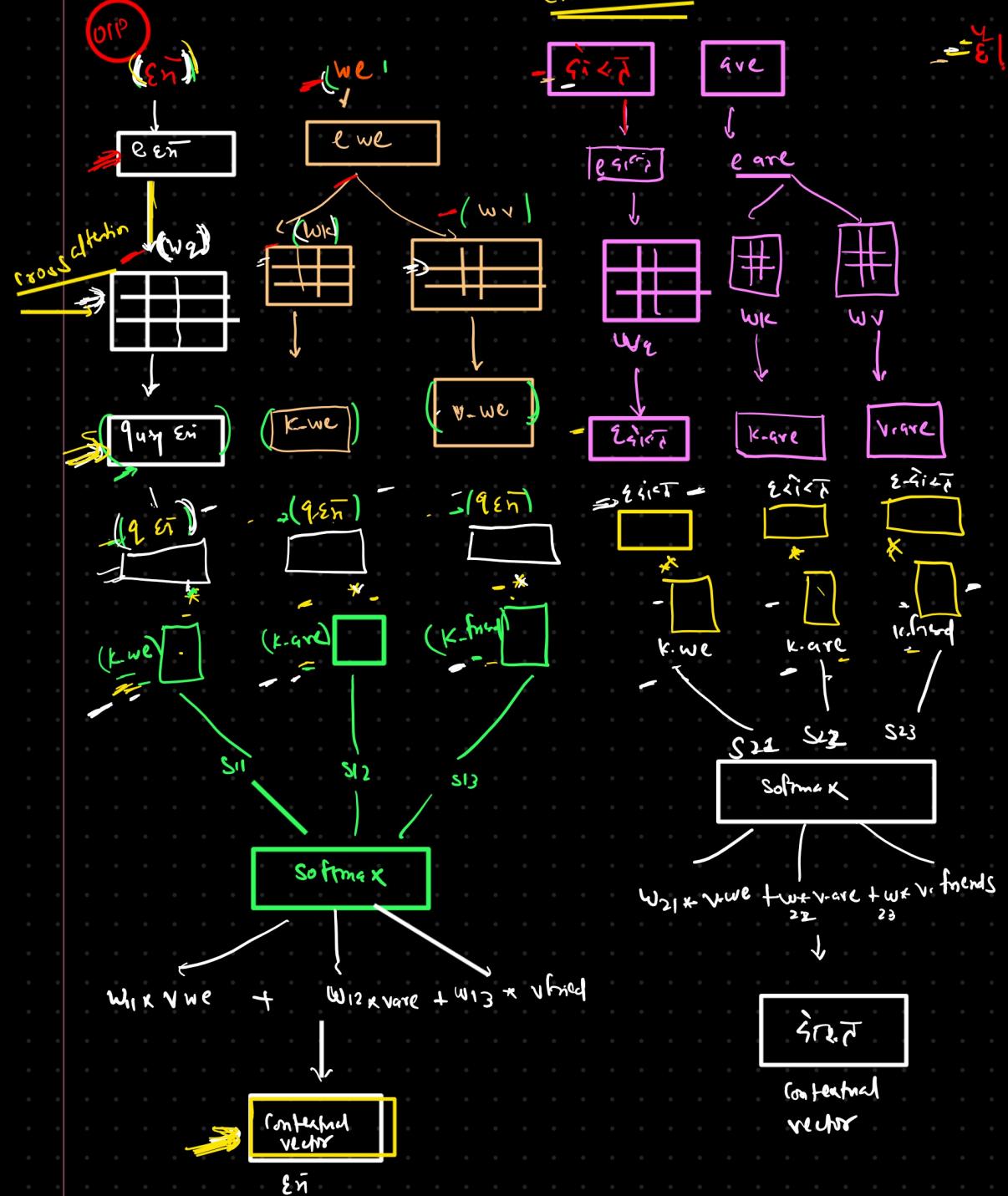
3 input \rightarrow query, key, value

$\Rightarrow Q, K, V \rightarrow$ matrix (weight matrix)

query \rightarrow output seq (e_n grid \vec{v}) -

key & value \rightarrow input seq (we are friends)

Cross attention

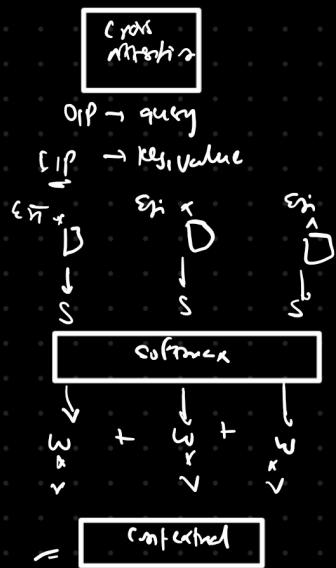
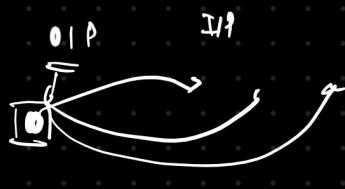
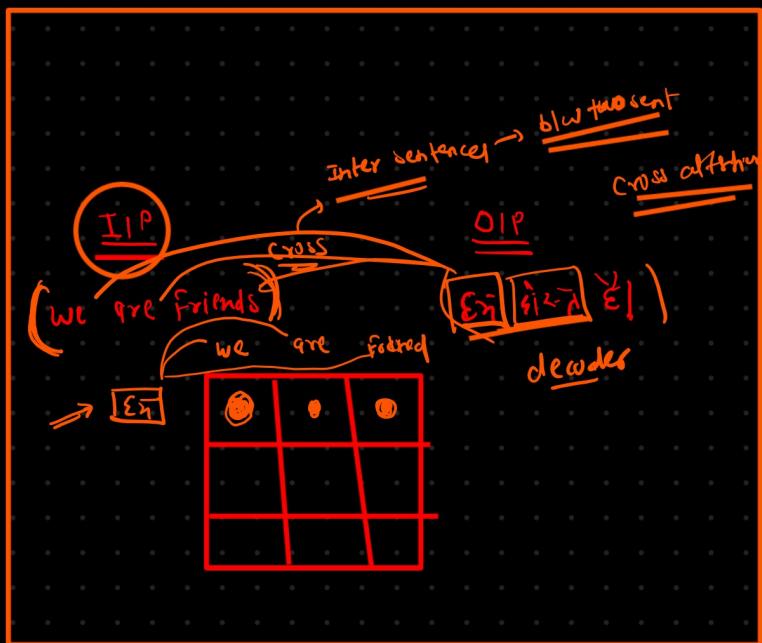
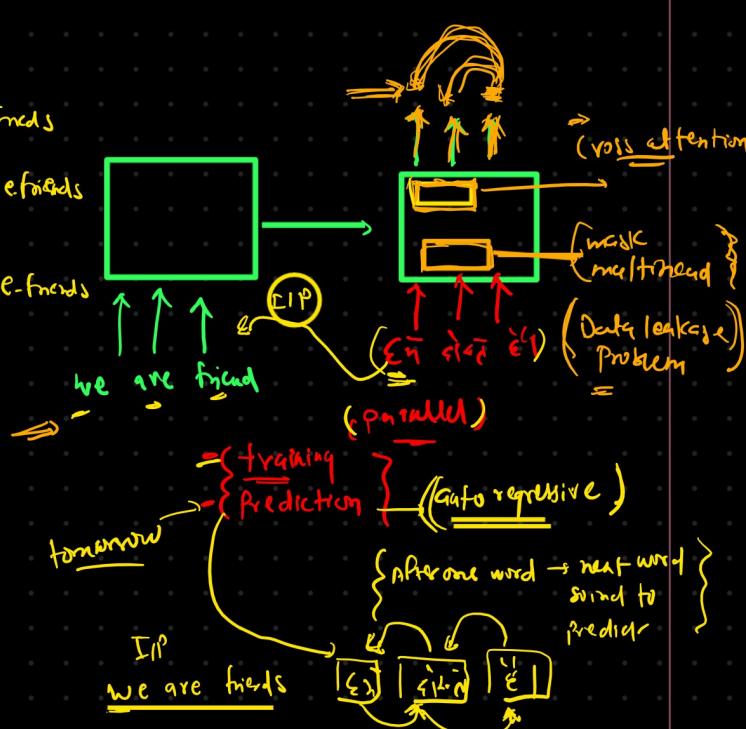


WA QIP friends

e_n	*	*	*
e_{in}	*	*	*
\vec{v}_1	*	*	*

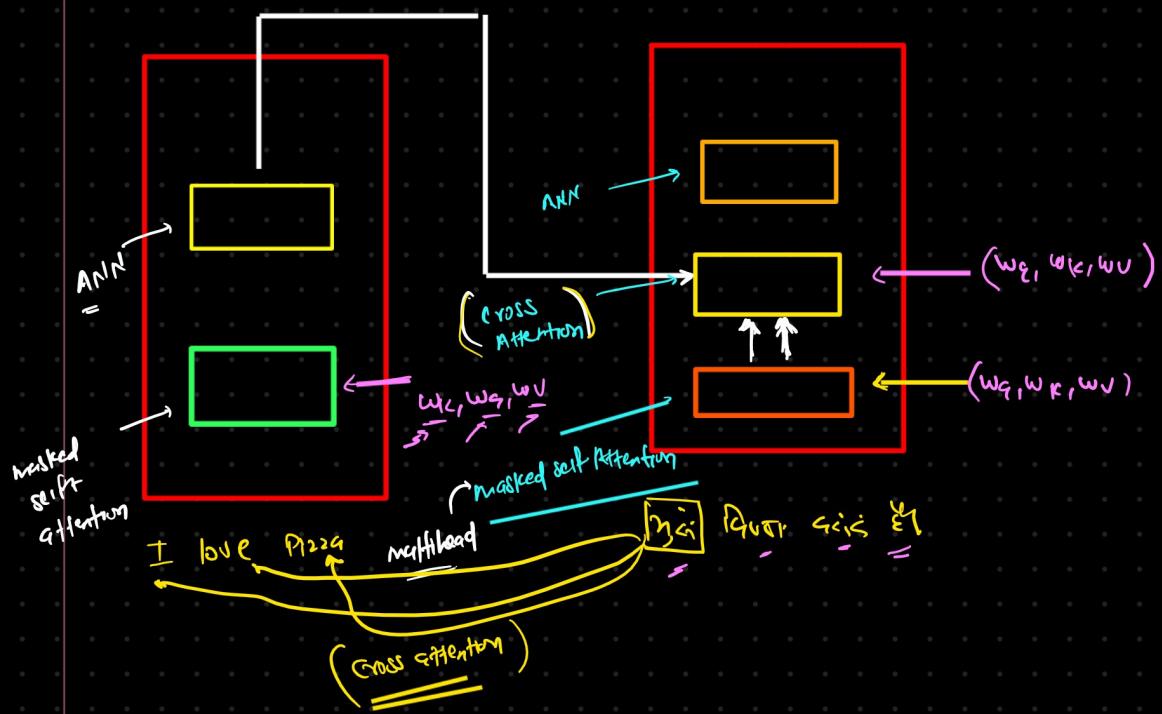
~~Cross~~ Complex

$$\begin{aligned} \underline{e_{e-\bar{x}}} &= w_{11} \cdot e_we + w_{12} \cdot e_are + w_{13} \cdot e_friends \\ \underline{e_{e-\bar{y} \bar{x}}} &= w_{21} \cdot e_we + w_{22} \cdot e_are + w_{23} \cdot e_friends \\ \underline{e_{e-\bar{z}}} &= w_{31} \cdot e_we + w_{32} \cdot e_are + w_{33} \cdot e_friends \end{aligned}$$

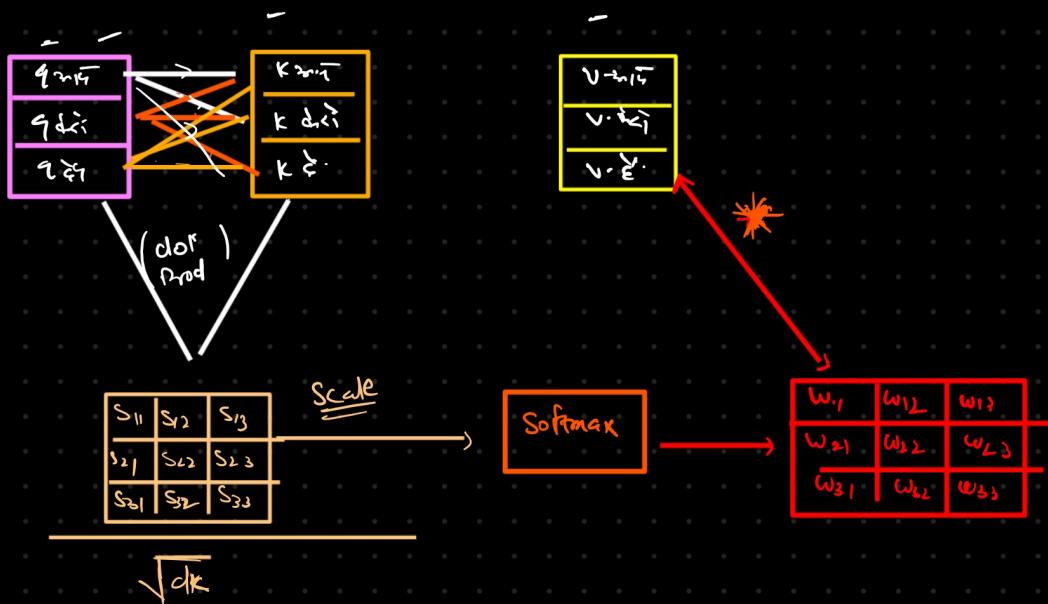
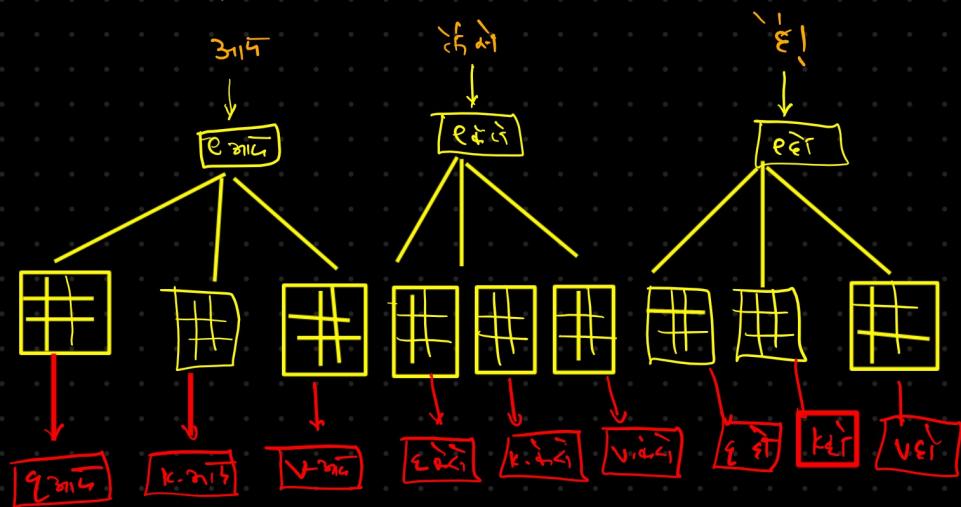


- 1 self attention -
- 2 multihed -
- 3 cross -
- 4 masked attention

	y_{11}	y_{12}	y_{13}	y_{14}
I	-	-	-	-
<u>Like</u>	-	-	-	-
<u>snow</u>	-	-	-	-



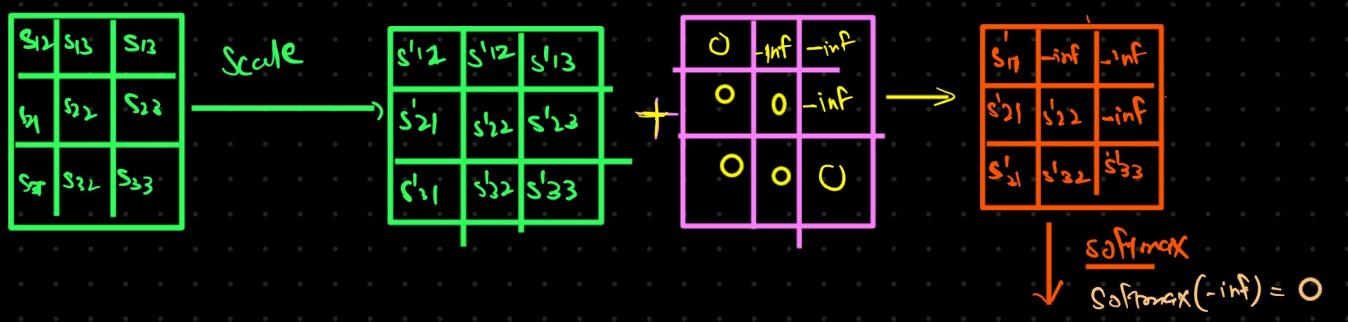
Now are you
bitting like me!



$$\begin{aligned}
 \text{[31] ce} &= w_{11} * v_{\text{31}} + w_{12} * v_{\text{32}} + w_{13} * v_{\text{33}} \\
 \text{[22] ce} &= w_{21} * v_{\text{21}} + w_{22} * v_{\text{22}} + w_{23} * v_{\text{23}} \\
 \text{[11] co} &= w_{31} * v_{\text{31}} + w_{32} * v_{\text{32}} + w_{33} * v_{\text{33}}
 \end{aligned}$$

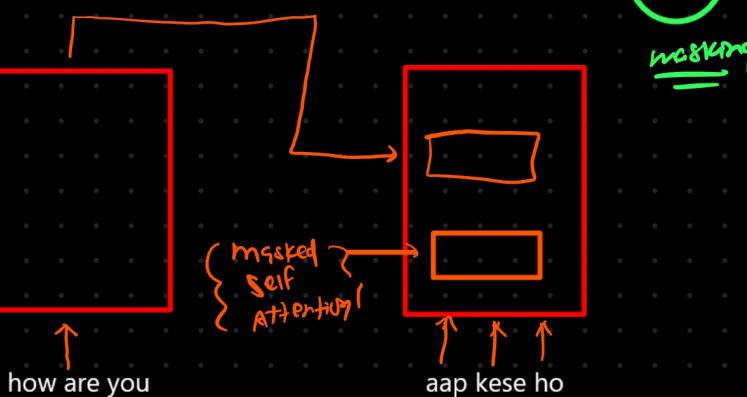
MASKed Self Attention

Data leakage \Rightarrow overfitting \Rightarrow (cheating)



mathematically

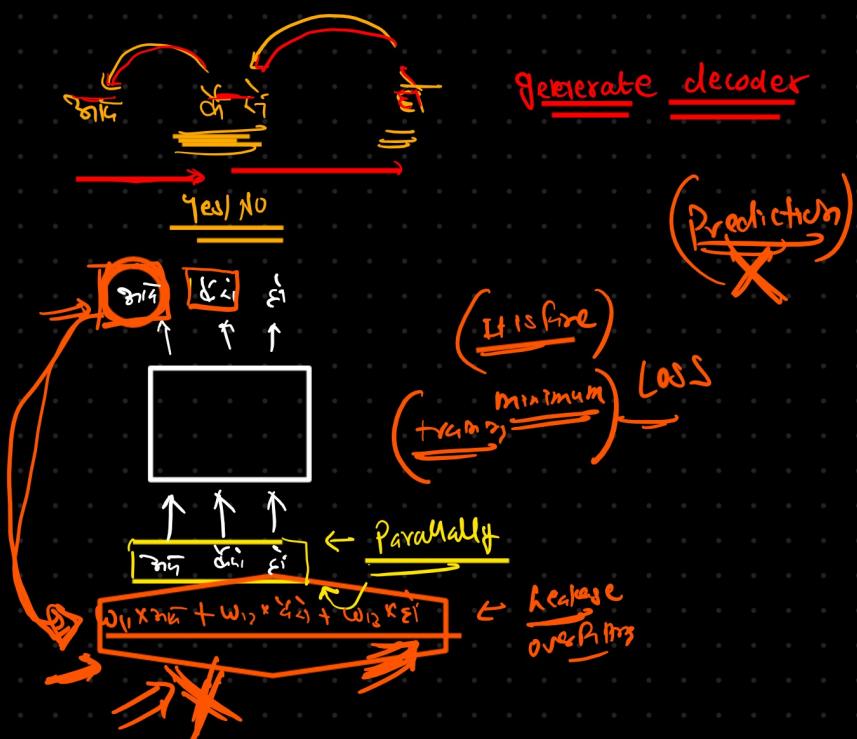
Translation Problem



$$\hat{y}_{ce} = w_{11} \cdot \hat{m}_e + w_{12} \cdot \hat{d}_e + w_{13} \cdot \hat{\epsilon}_e$$

$$\hat{y}_{ce} = w_{21} \cdot \hat{m}_e + w_{22} \cdot \hat{d}_e + w_{23} \cdot \hat{\epsilon}_e$$

$$\hat{y}_{ce} = w_{31} \cdot \hat{m}_e + w_{32} \cdot \hat{d}_e + w_{33} \cdot \hat{\epsilon}_e$$



Current token at the time of training

You **revealed** the future token

Current token \hat{m}_1

$$\Rightarrow (w_{11} \cdot \hat{m}_1 + w_{12} \cdot \hat{d}_1 + w_{13} \cdot \hat{\epsilon}_1)$$

{Training it is fine}
{but when prediction
it will fail}

\hat{m}_1, \hat{d}_1

o

$$w_{21} \cdot \hat{m}_1 + w_{22} \cdot \hat{d}_1 + w_{23} \cdot \hat{\epsilon}_1$$

$$(\hat{m}_1, \hat{d}_1, \hat{\epsilon}_1)$$

no

$$w_{31} \cdot \hat{m}_1 + w_{32} \cdot \hat{d}_1 + w_{33} \cdot \hat{\epsilon}_1$$

