

# NLP

(1) Basic Language Model

(2) Distributed " "

(3) Large language model (LLM) + Gen AI

2000-2014 → RNN/LSTM/GRU

2014 - Encoder-Decoder

Seq 2 Seq Learning with NN.

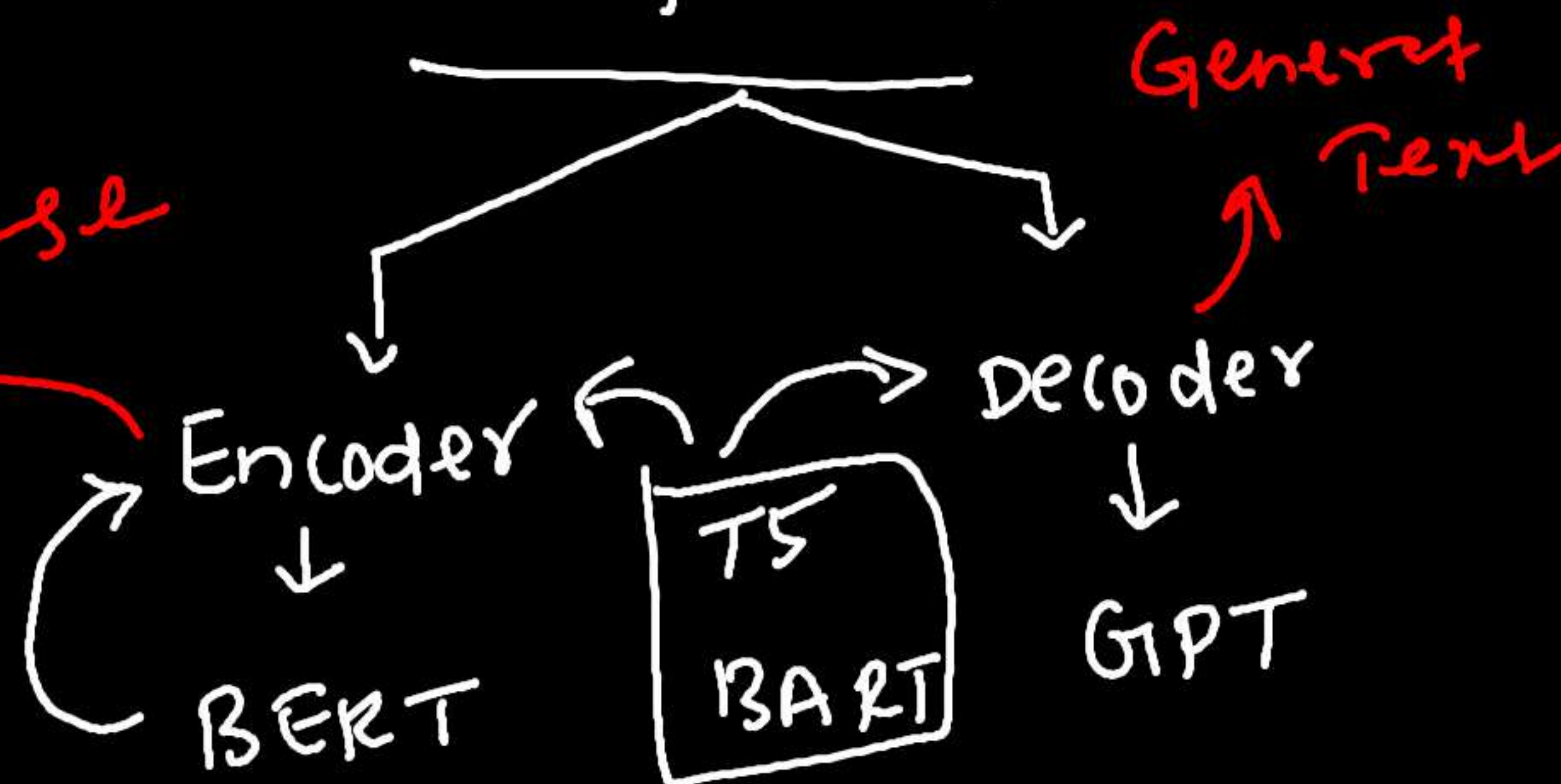
**LSTM**

Encoder-Decoder

Research Paper

## Transformers

*understand language*



## Transformers + BERT + GPT

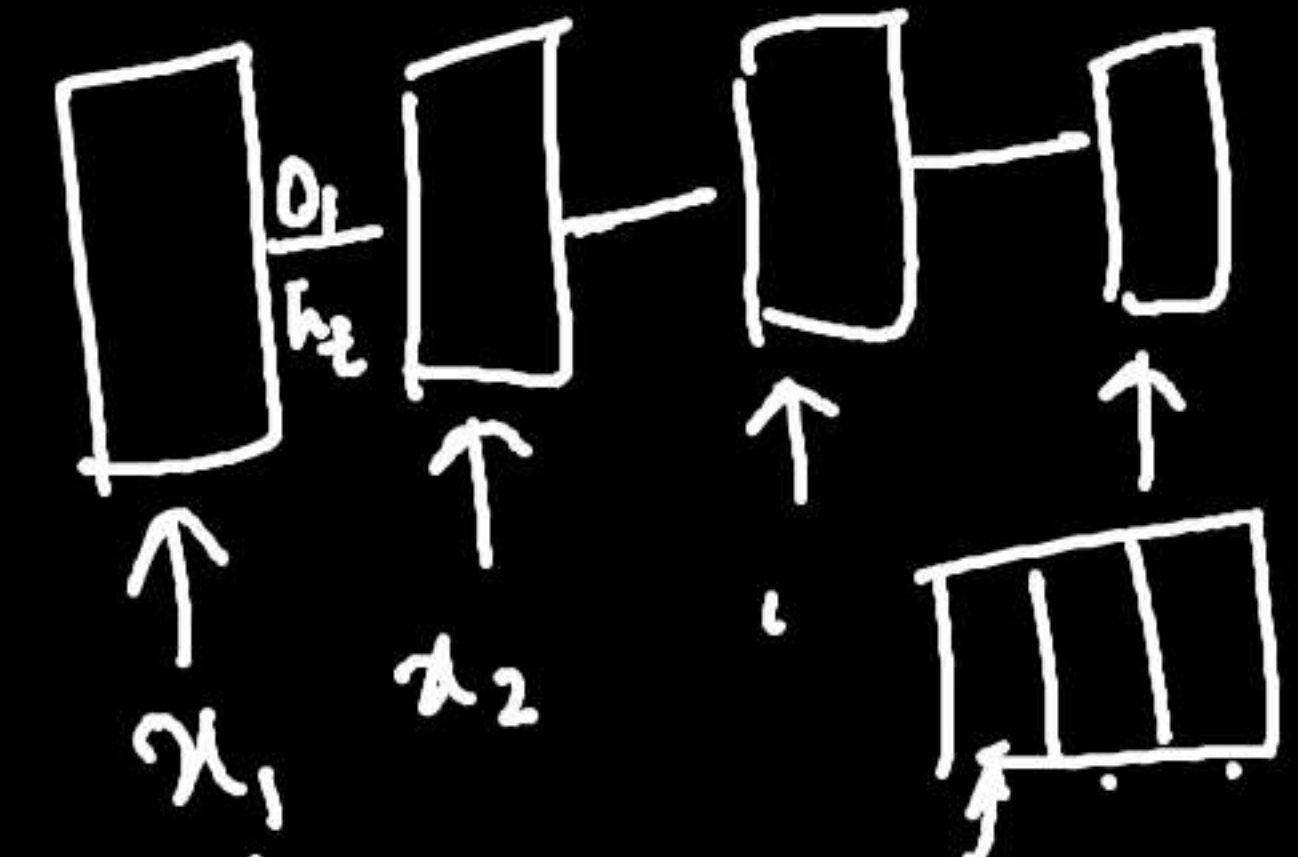
Theory + Case Study

Hugging Face - Arch  
 - models  
 - datasets  
 company + OpenSource communi..



NN

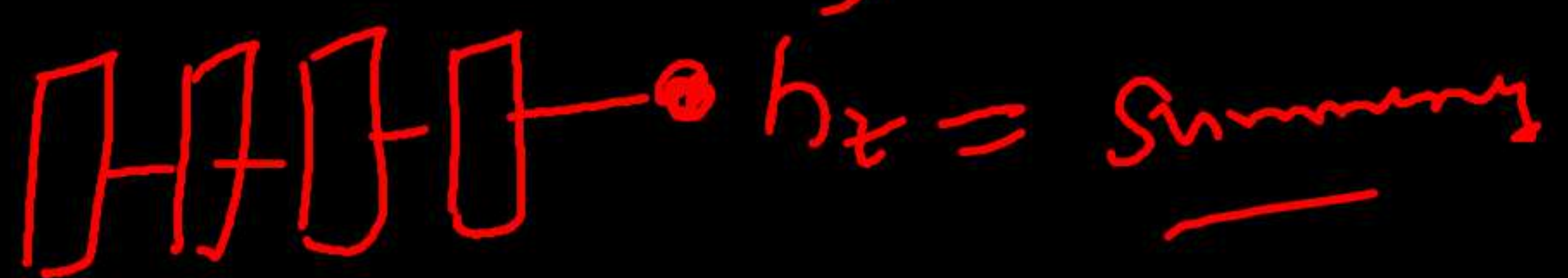
- ① ANN/MLP/DNN — Tabular data
- ② CNN / computer vision — Image
- ③ RNN/LSTM/GRU — Text (sequential data) → LSTM/GRU
- ④ Transformers — multimodal (1, 2, 3 + Reinforcement learning)



↳ seq 2 seq data

Self-attention model | parallel processing

↓  
very fast



vice to met you



2015 - Bahdanau - Neural Machine Translation

→ Attention model → very famous

Architecture

LSTM

sequence model

→ teacher

→ slow

→ huge data

→ Problem

Att -

30%

seq →

$$C_i = \sum \alpha_{ij} h_j$$

attention weights.

$1 \times a = 10$   
 $2 \times b = 10$   
 $100 \times c = 10 = (d)$   
 $500 \times d = 1, 2, 4$  → weight



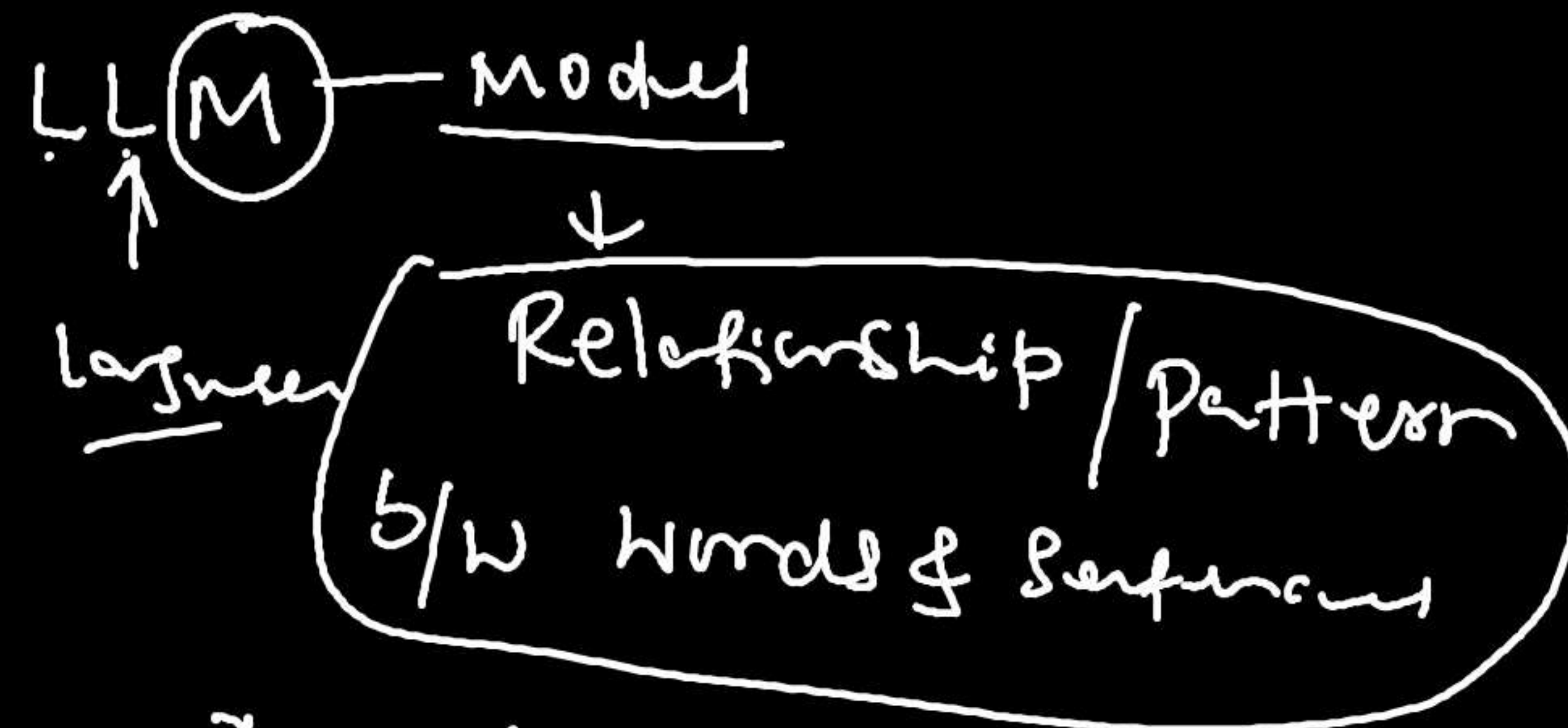
V.V. Significant

2017 - Transformers

- Attention Is All You Need

① What is model?

② What is language model?



x      input → output  
y

Language  
data — 

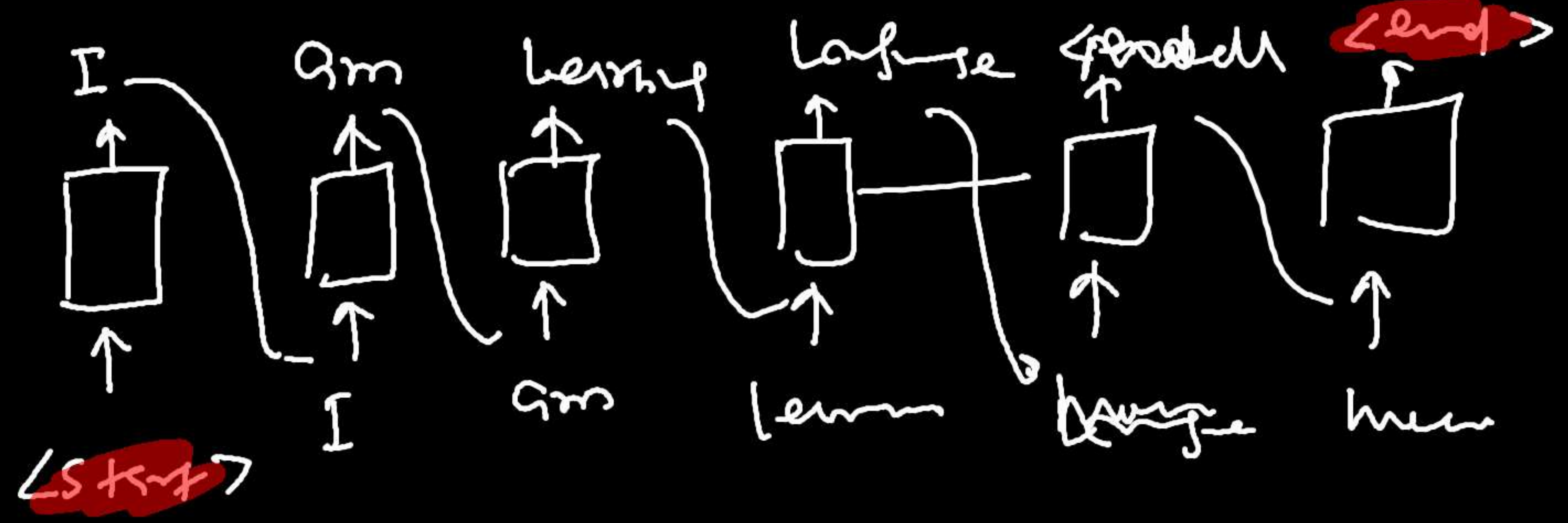


# ④ Language Modeling Techniques? uni-directional

① Autoregressive Task → A mask is applied to

✓ ② Auto Encoding task fill the seq  
→ Bi-directional

<S> I am learning Language modeling <S>



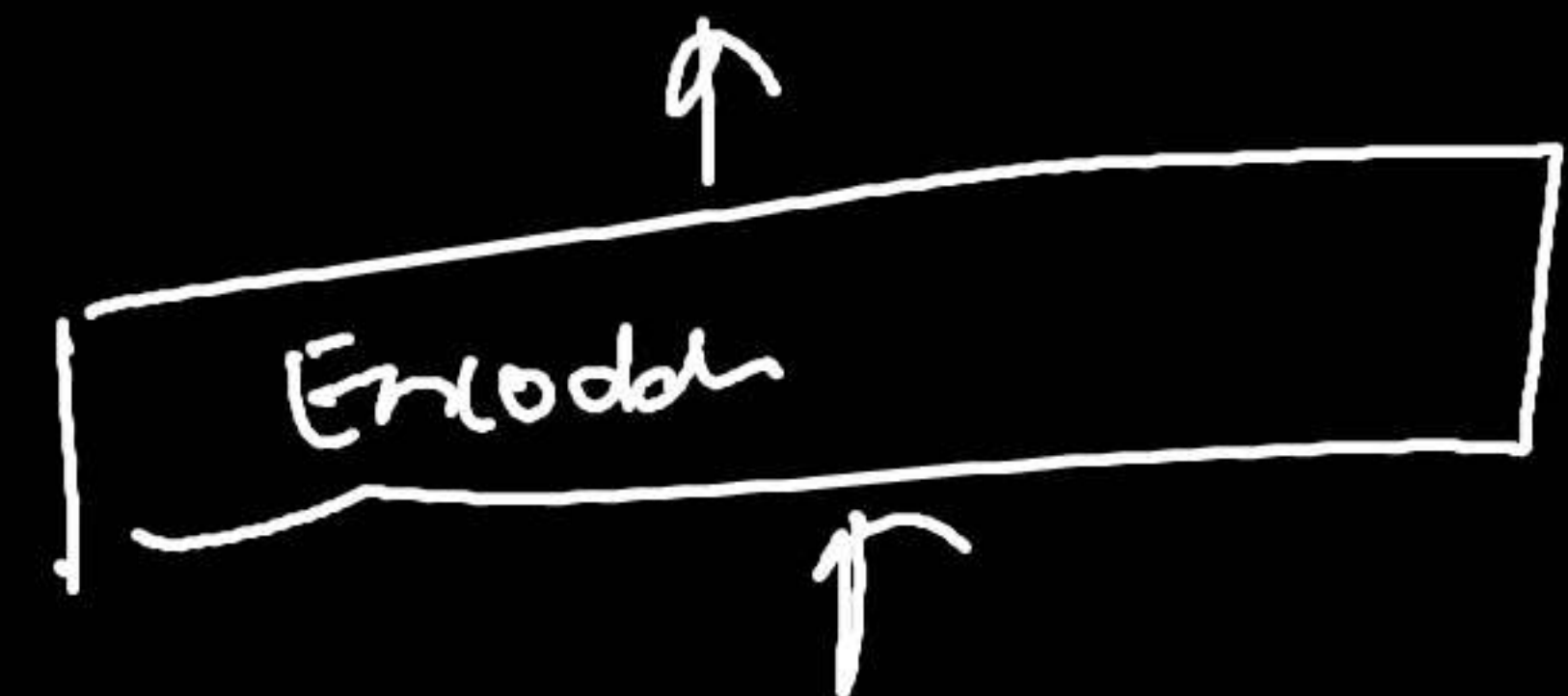


Auto Regressive  $\rightarrow$  Predict one token at a time

Auto Encoder  $\rightarrow$  Generate the entire target seq ~~token~~ together

Aim to  $\max^m$  the likelihood of the next token given the previous token

$\min^m$  Reconstruction error.





## Summary

2014 → Encoder-Decoder Architecture

2015 → Bahdanau Attention Mechanism

2017 :- Transformers Architecture  
↳ Attention is all you need  
↳ universal model

2018 → LLM → Fine-tuned model



Paper :- Attention is all you need

2017 → Proposed Transformer Architecture  
which used —

- Encoder-Decoder Architecture
- Self-attention Mechanism
- Positional Encoding
- Multihead self-attention model
- Cross Attention with Multihead.



$w_2 \sim 50 \text{ to } 300$

## Building Blocks

— please note, this is a Simplified

View of Transformer

- Encoder-Decoder
  - Attention mechanism
  - Parallel processing
- make huge improvement

is good at generating text

$N \times$

adds the notion of self-information

Positional Encoding

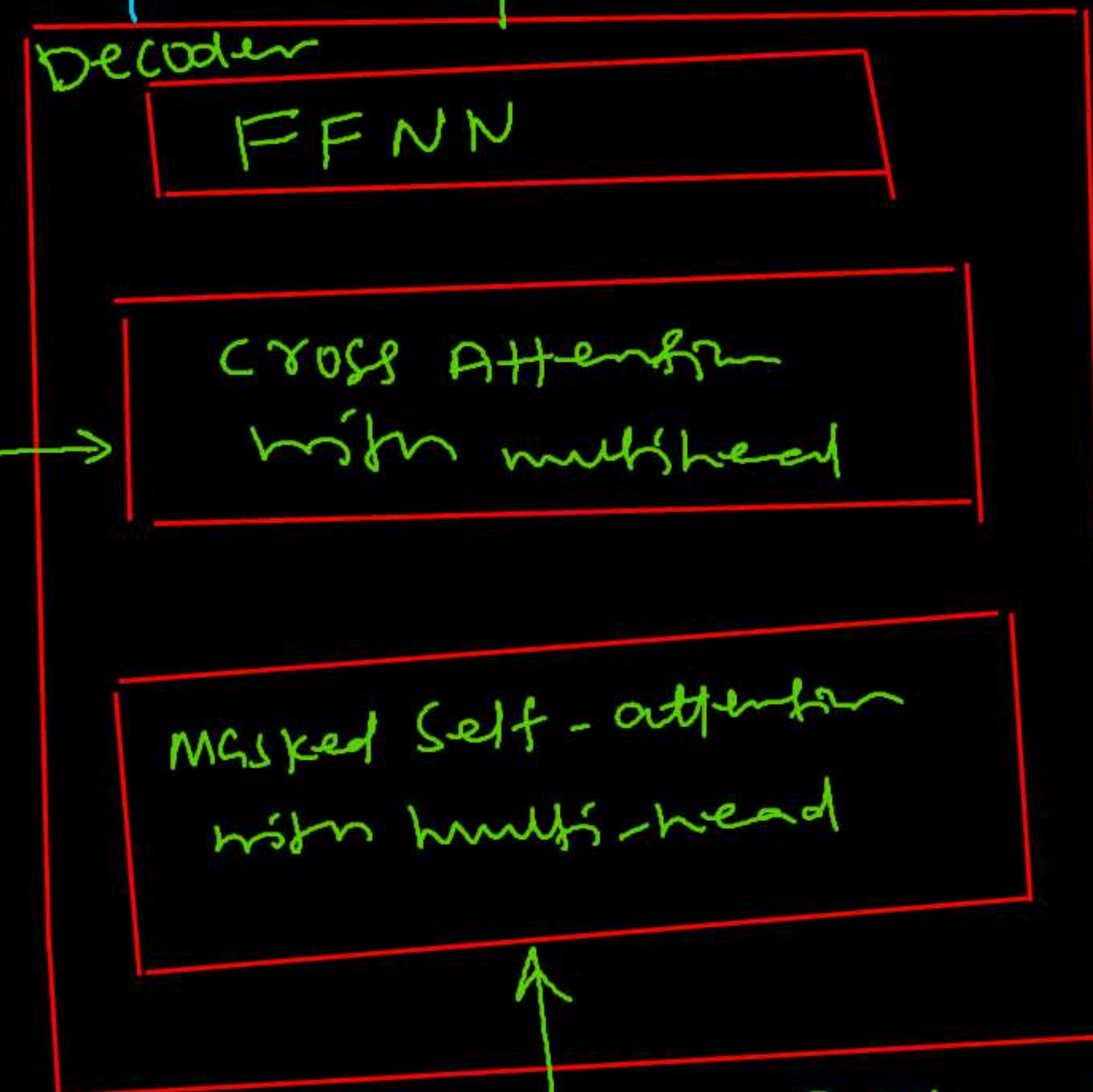
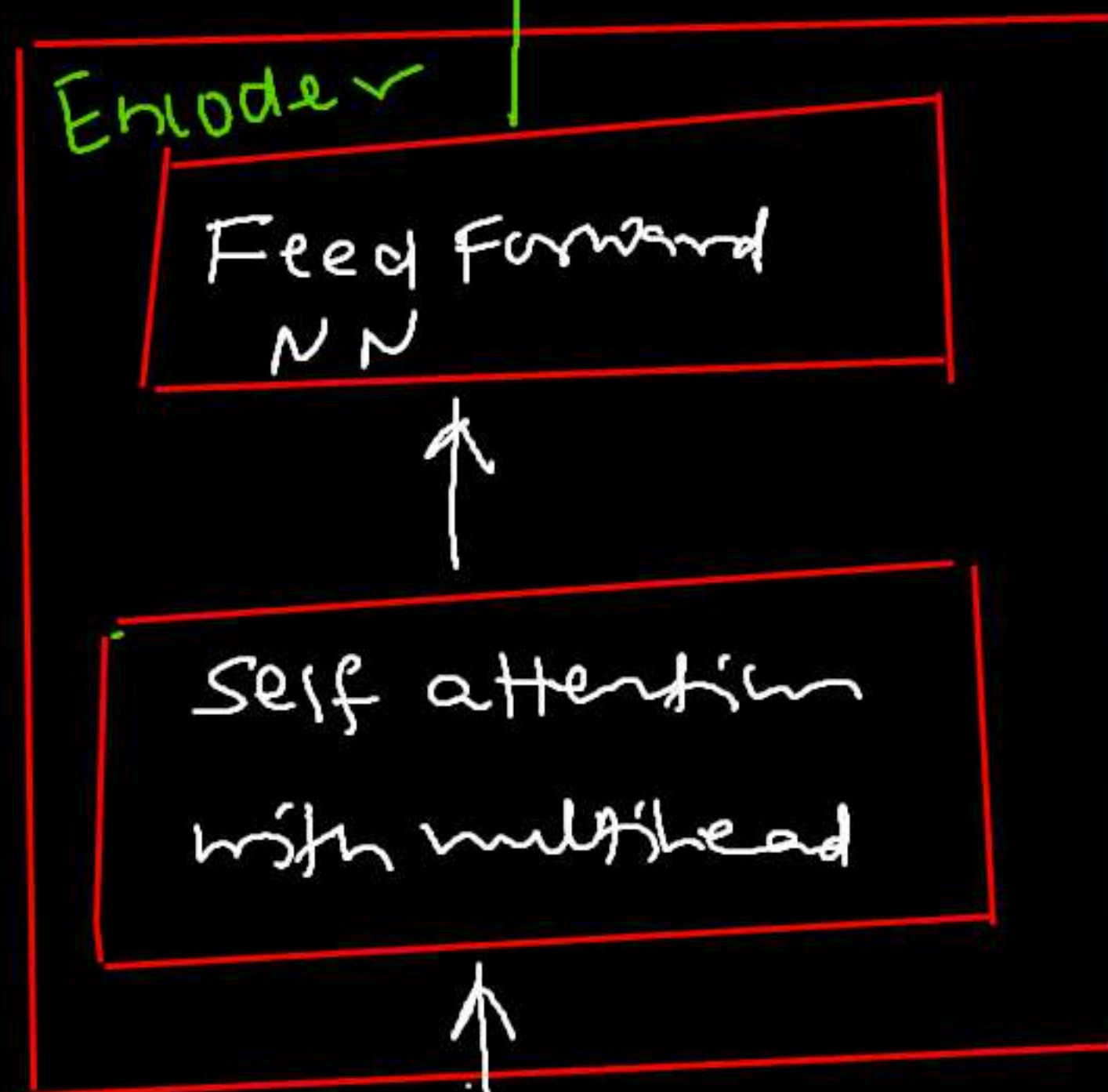
+

Input Encoding

Input

→ understand the long term

→ Encoder is good at understanding "text"



+

Input Encoding

Input

→ Positional Encoding

→ Hindi  
S  
↓  
G  
 $\times N = 6$

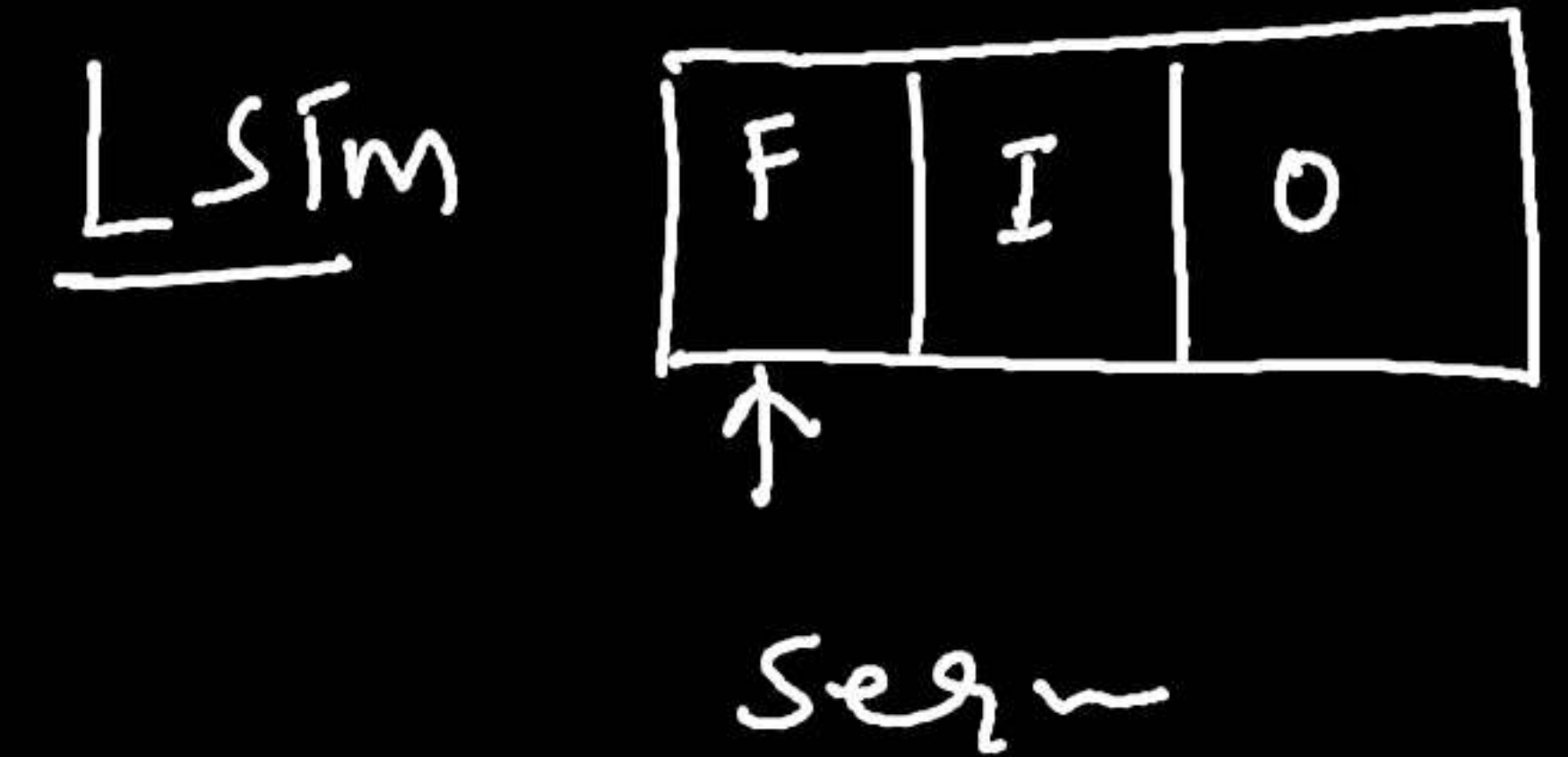
Output



## Self-attention mtr multthead

① Attention replaced need to recurrence & make the process of ~~embedding calculation parallel~~.

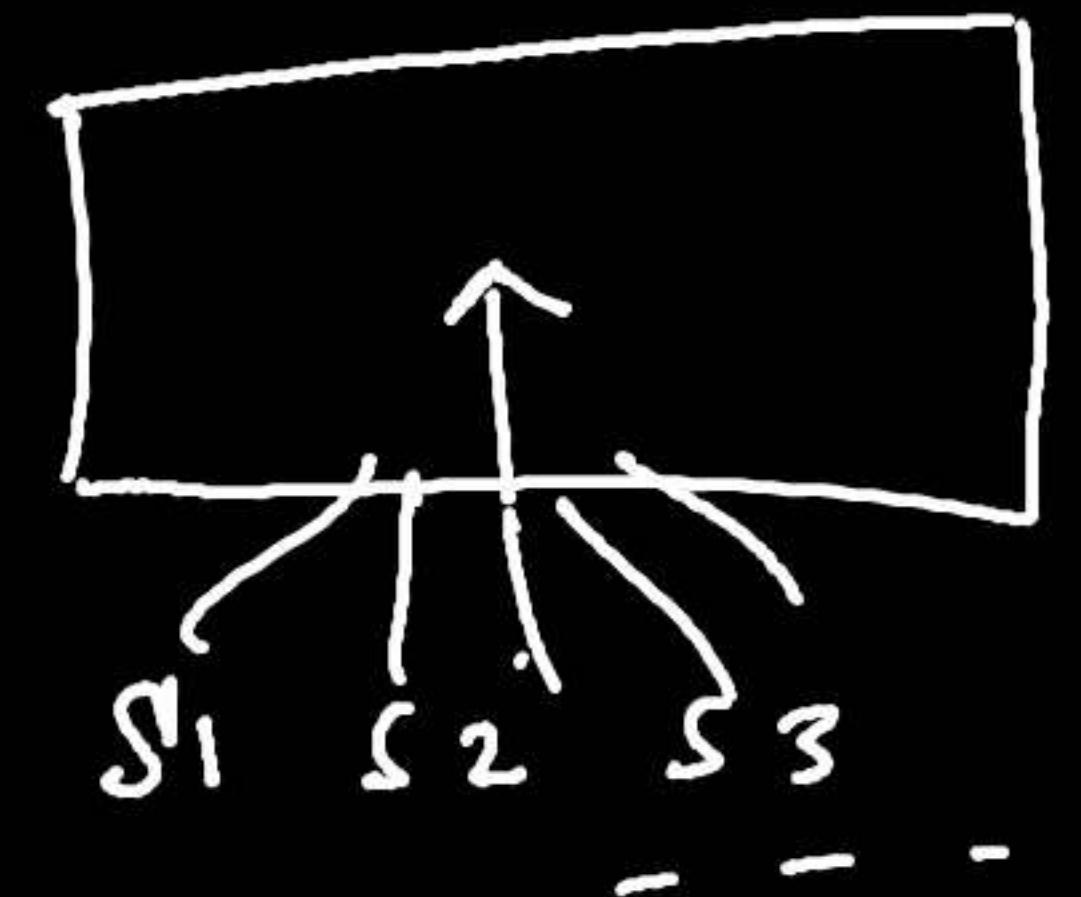
② Attention mechanism will find how each word related to all other words in sequence & ~~generate contextual embedding~~



\* How to find word similarity?

→ Attention will run dot product b/w word vectors determine the strongest relationship of a word with all other words.

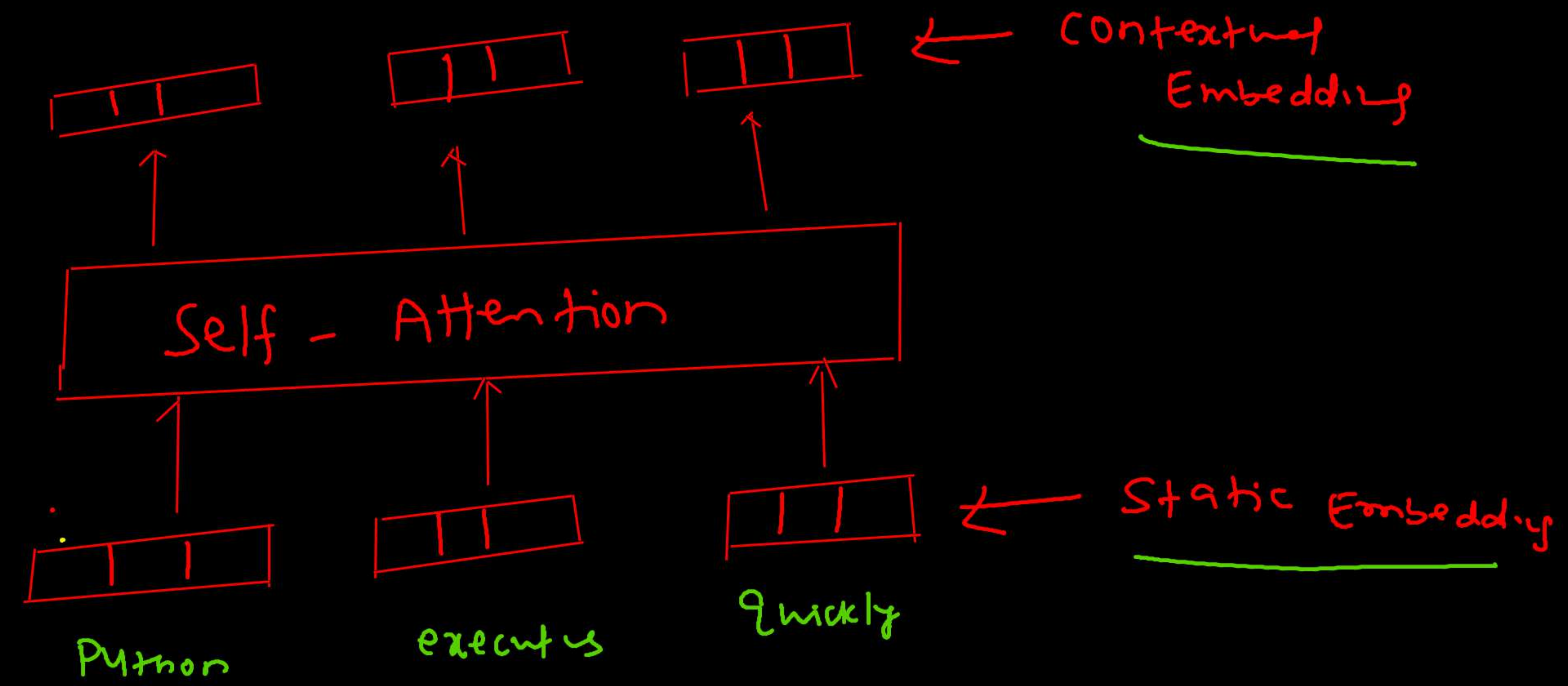
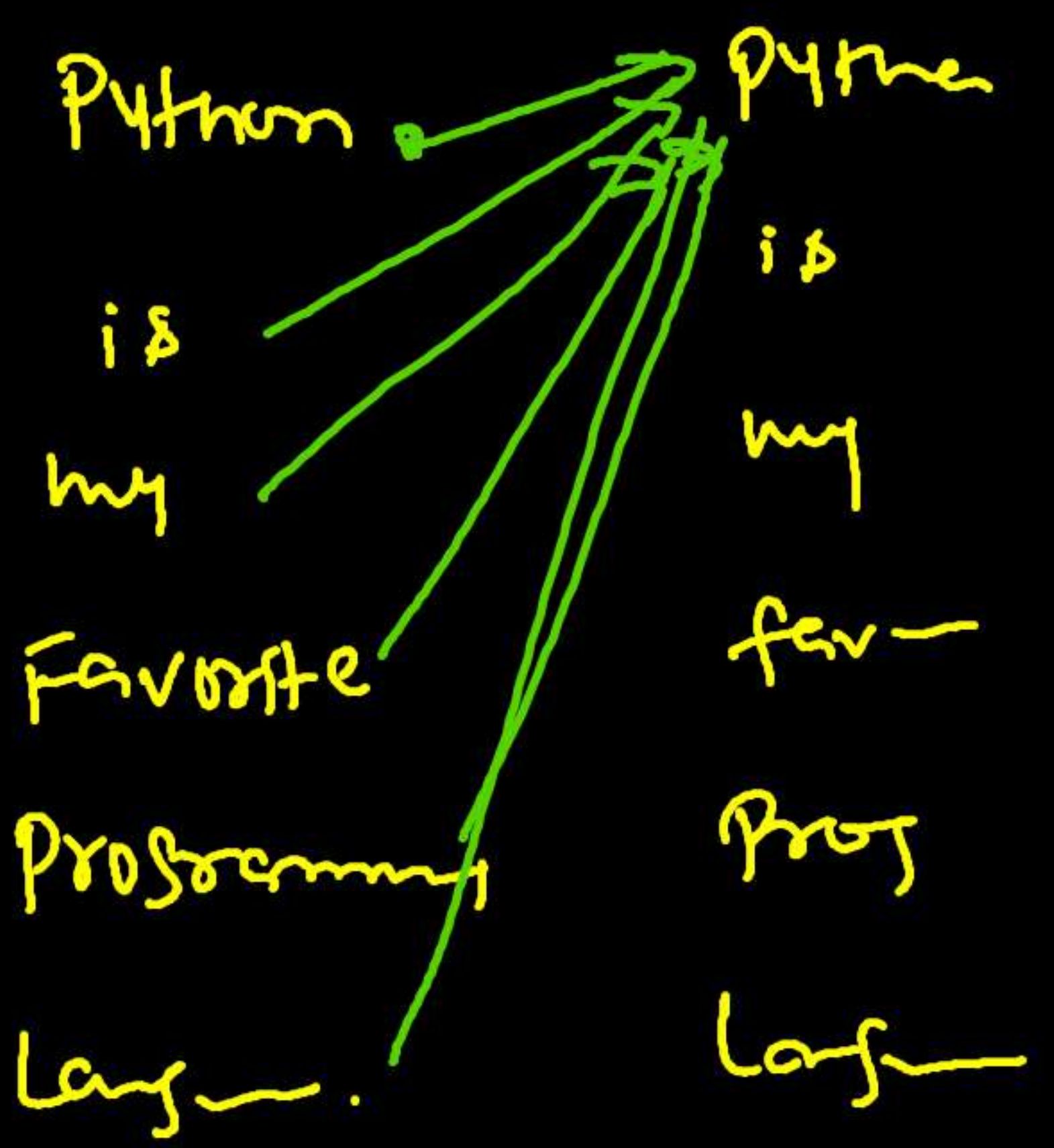
Self-Att





eg :- python is my favorite programming language.

# Self - Attention





dance = 5

W2V

$$\begin{aligned} \text{python} &= [0.1 \ 0.2 \ 0.9 \ 0.7 \ 0.6] = \underline{\underline{0.6}} \\ \text{executes} &= [0.2 \ 0.3 \ 0.4 \ 0.6 \ 0.3] = \underline{\underline{0.5}} \\ \text{quicks} &= [0.3 \ 0.1 \ 0.8 \ 0.7 \ 0.6] = \underline{\underline{0.3}} \end{aligned}$$

Avg

$e_{\text{python}}$        $e_{\text{executes}}$        $e_{\text{quicks}}$

$$\begin{aligned} \text{python} &= 0.6 * \text{python} + 0.5 * \text{executes} + 0.3 * \text{quicks} \\ \text{executes} &= 0.2 * \text{python} + 0.7 * \text{executes} + 0.1 * \text{quicks} \\ \text{quicks} &= 0.1 * \text{python} + 0.1 * \text{executes} + 0.8 * \text{quicks} \end{aligned}$$



$$z_{python} = \overset{e_{python} \times e_{python}}{\boxed{0.5}} \times e_{python} + \overset{e_{python} \times e_{exec}}{\boxed{0.3}} \times e_{exec} + \overset{e_{python} \times e_{quick}}{\boxed{0.2}} \times e_{quick}$$

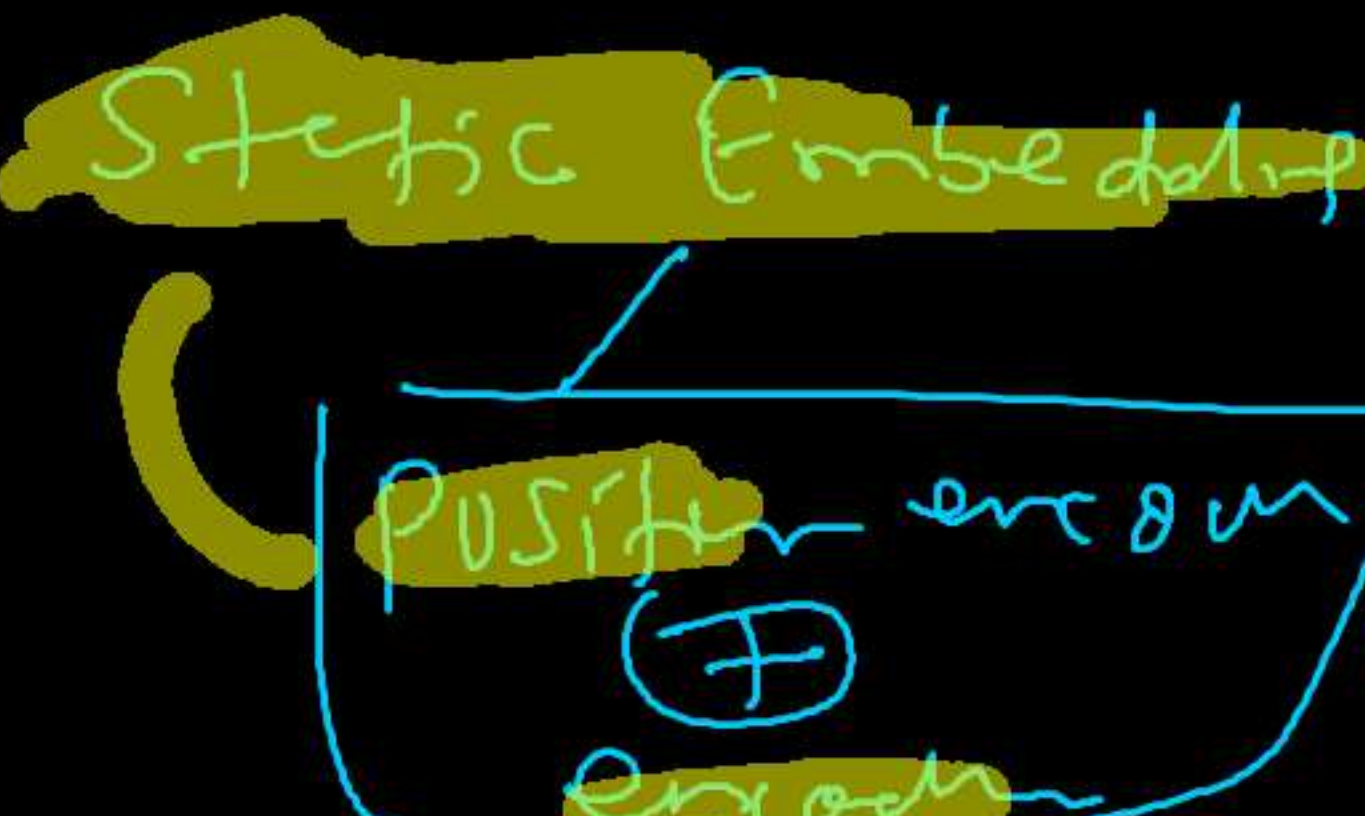
$$z_{exec} = \overset{e_{exec} \times e_{python}}{\boxed{0.2}} \times e_{python} + \overset{e_{exec} \times e_{exec}}{\boxed{0.7}} \times e_{exec} + \overset{e_{exec} \times e_{quick}}{\boxed{0.1}} \times e_{quick}$$

$$z_{quick} = \underset{e_{quick} \times e_{python}}{\boxed{0.1}} \times e_{python} + \underset{e_{quick} \times e_{python}}{\boxed{0.1}} \times e_{exec} + \underset{e_{quick} \times e_{quick}}{\boxed{0.8}} \times e_{quick}$$




 You are screen sharing
 

 Stop share





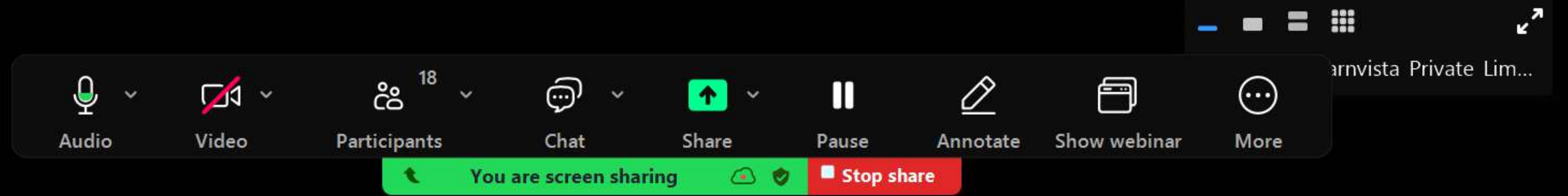
V.V.V. Forp  
Intern

Q(1) What is Q, K & V?

Q(2) - Why Q, K & V?

Q - Query - What information the model is trying to extract.

K - Key :- Content or features that the model can pay attention to

A screenshot of a Zoom meeting toolbar. It includes icons for Audio, Video, Participants (18), Chat, Share, Pause, Annotate, Show webinar, and More. Below these icons is a green bar that says "You are screen sharing" and a red button that says "Stop share".



→ token embedding directly without transform  
in  $Q, K$  &  $V$ , the model would be limited to  
just one fixed way of computing relationship  
between tokens.



① init the 3 weights:  $W_q$   $W_k$   $W_v$

② cal:  $Q = XW_q$   $K = XW_k$   $V = XW_v$

$$[Q]_{T \times d} \quad [K]_{T \times d} \quad [V]_{T \times d}$$

③ cal

$$\text{Soft} \left( \frac{Q * K^T}{\sqrt{d_k}} \right) V$$