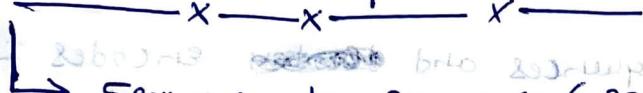


5 Sequence To Sequence Models:-



→ Sequence-to-Sequence (Seq2Seq) Model are a class of Models designed to handle tasks where Input and Output are sequences of varying lengths.

→ They are widely used in Natural Language Processing (NLP) tasks such as :- { Applications }

Following are the applications :-

1. Machine translation :- Convert Text from one language to another.

2. Text summarization

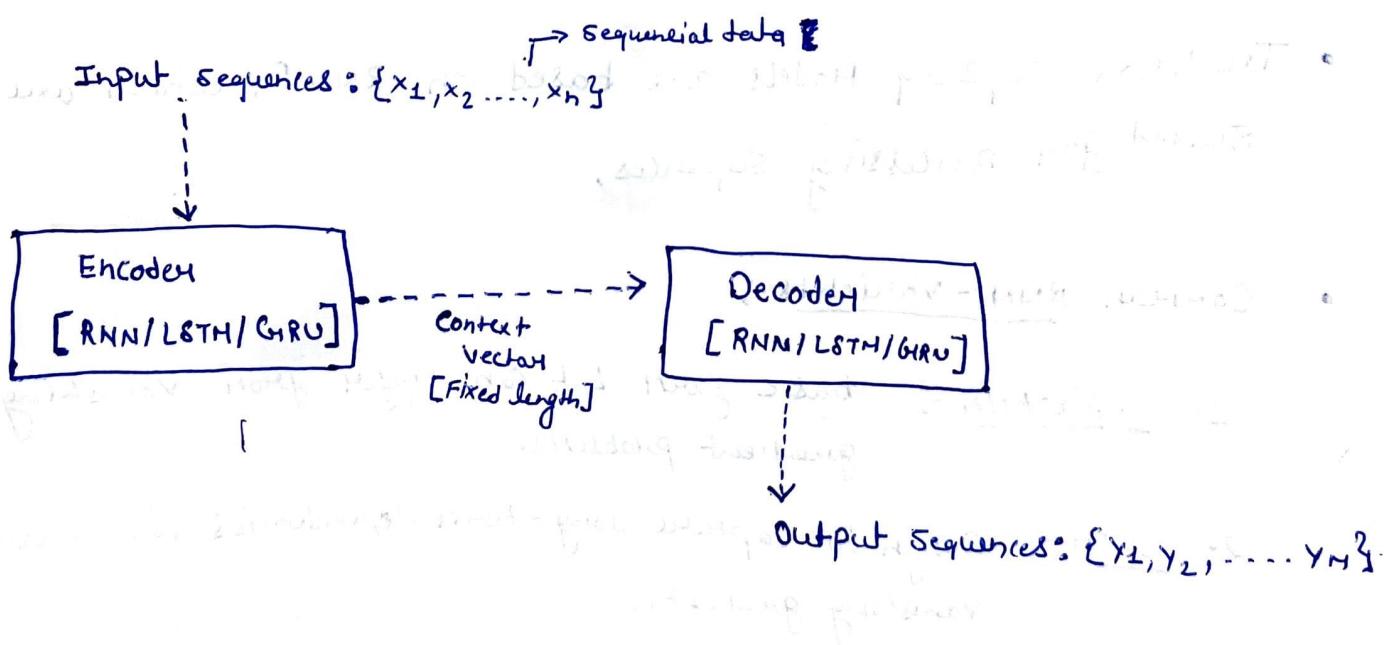
3. Dialogue systems :- Building conversational agents.

4. Text-generation :- Generating human-like text based on a given prompt. OA chatbot.

5. speech-to-text.

Basic Structure of Seq2Seq

→ Encoder-Decoder Architecture, the first Seq2Seq Model.



* NOTE :- It is Not necessarily that the length of the input-output seq is same.

Encoder - Decoder Architecture.

Encoder:

↳ Processes the input sequences and ~~encodes~~ encodes it into a fixed-length context vector (thought vector). This context vector captures the essence of the input sequences.

Decoder:

↳ Takes the context vector from the encoder and generates the output sequences.

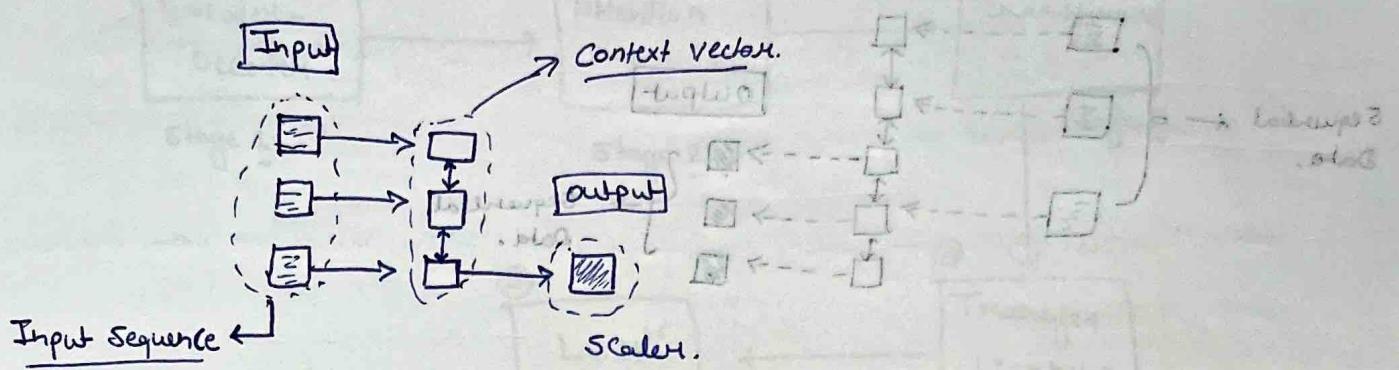
The output is produced token by token, usually using the previously generated token as input for the next step.

RNN-Based Seq2Seq Models:

- Traditional seq2seq Models are based on RNN's, which are well-suited for processing sequences.
- Common RNN Variants:
 - 1. Simple-RNN: basic form but can suffer from vanishing gradient problems.
 - 2. LSTM: Design to capture long-term dependencies and avoid vanishing gradients.
 - 3. GRU: A simpler alternative to LSTM with similar performance.

Types of RNN :-

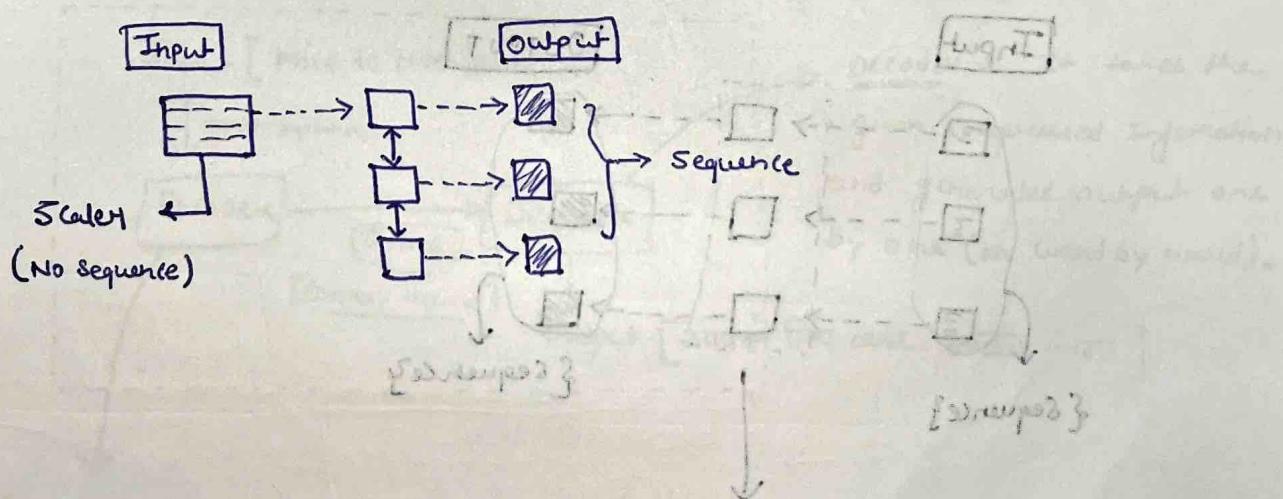
①. Many to ONE :-



Ex - Sentiment Analysis \Rightarrow Input \rightarrow Movie Review.

sent per two two steps \rightarrow output sent per two two steps \rightarrow predict (positive, Negative, Neutral).

②. ONE TO MANY :-

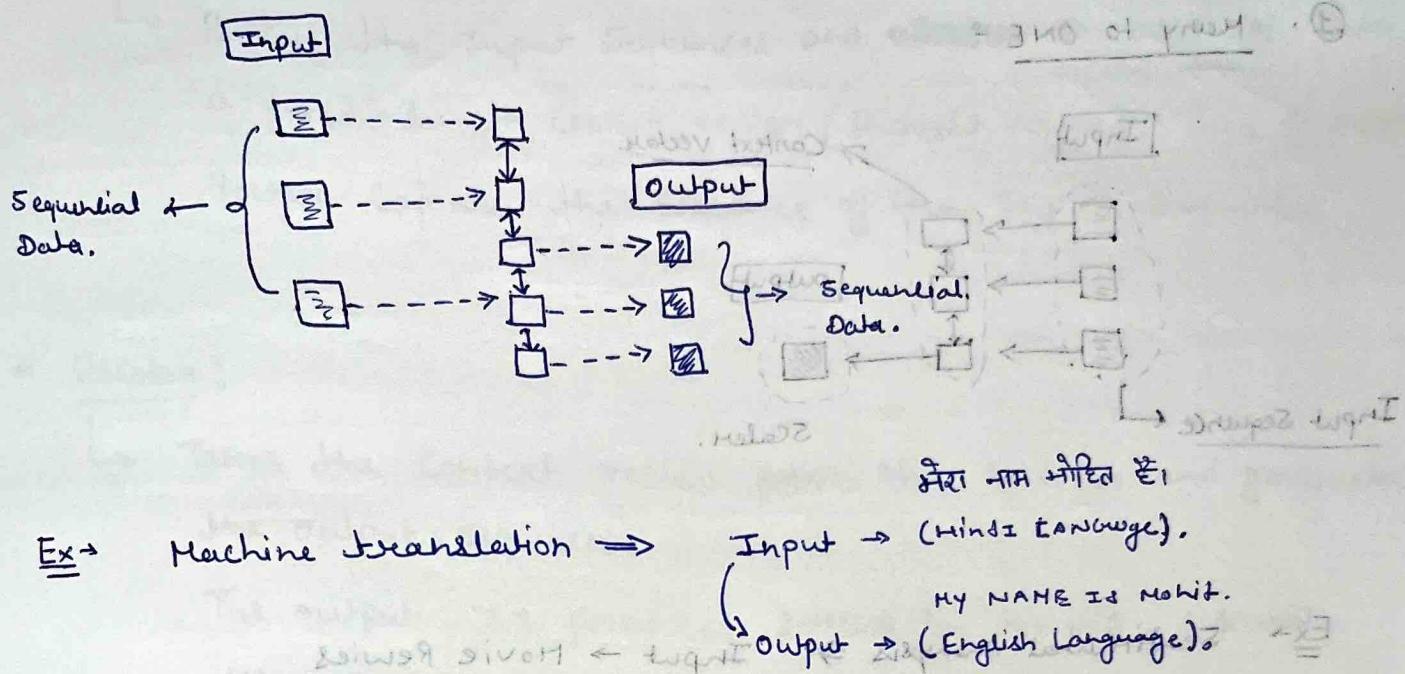


Ex - Image Captioning \Rightarrow Input \rightarrow Image (Image of something)

nothing we will show

Output \rightarrow Description \Rightarrow "Man IS riding Horse!"

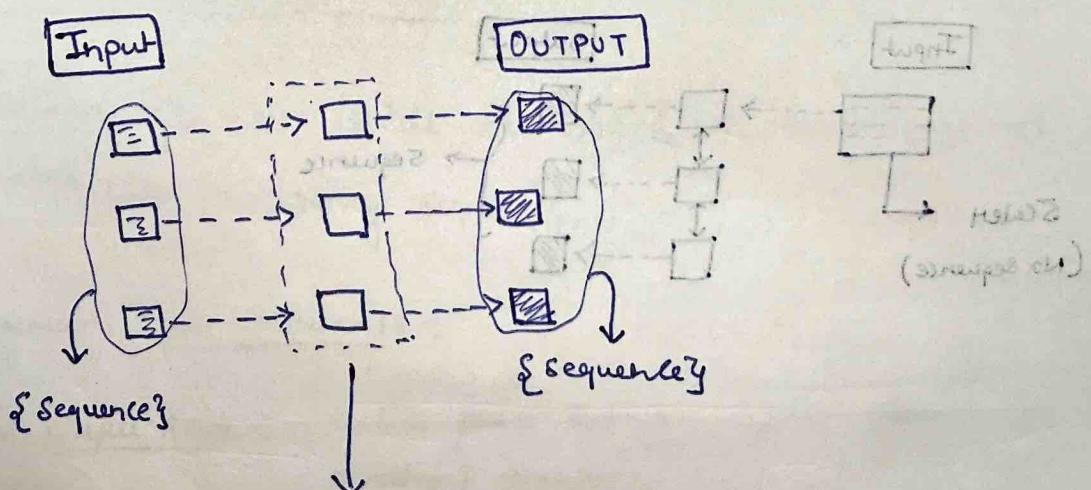
3. MANY to MANY \therefore { ASYNCHRONOUS }
 (having input & output of size, output & output)



NOTE \rightarrow Input may have different sequence and output may have different sequence.

\rightarrow Sequence to sequence Model designed to solve this many to many problem.

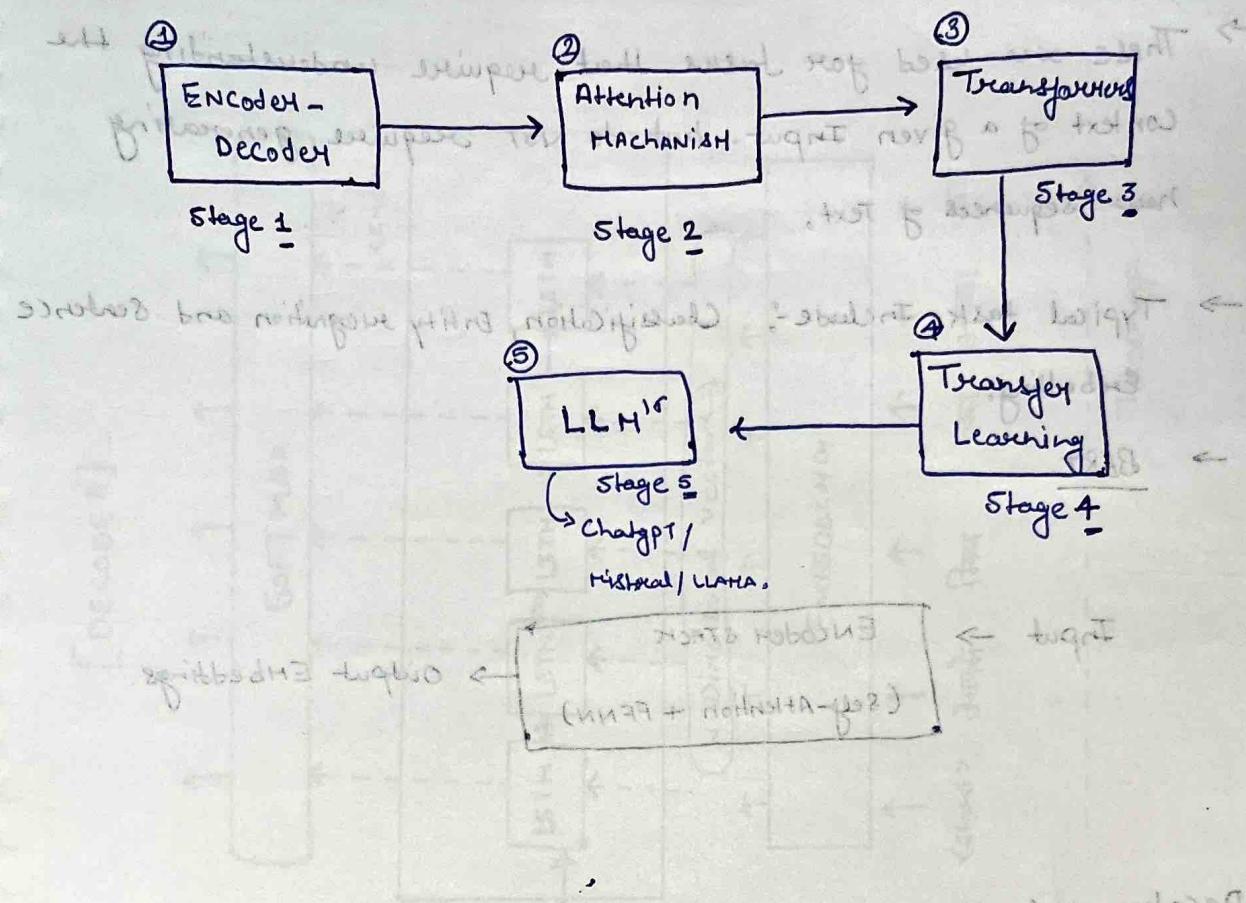
4. MANY to MANY \therefore { SYNCHRONOUS }



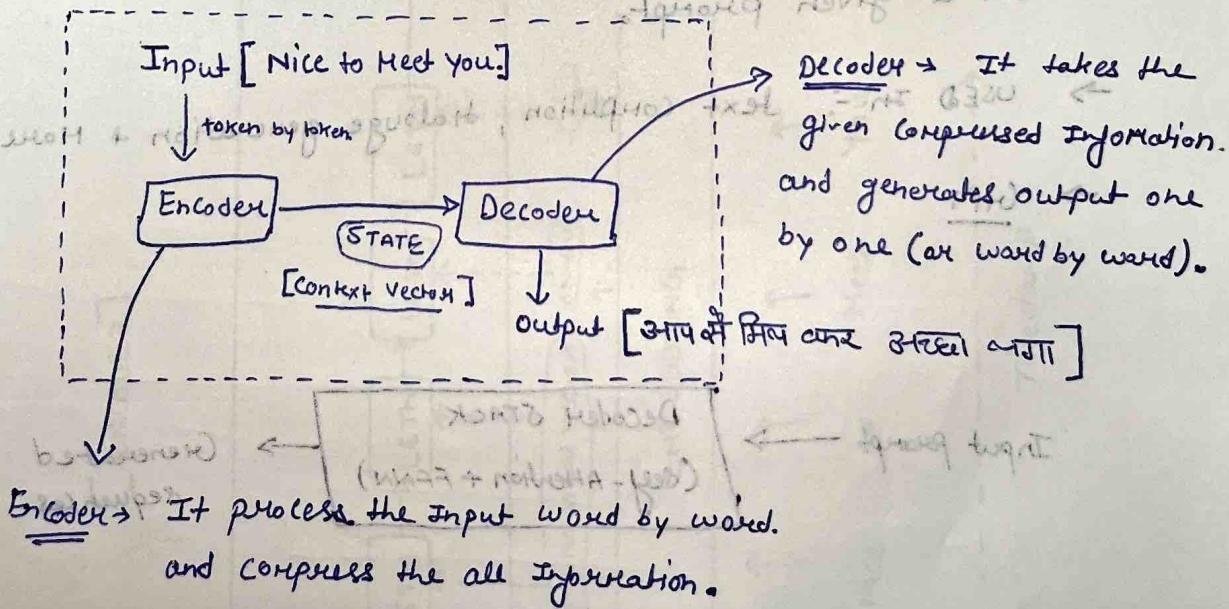
USING :- ①. Parts of speech tagging.

②. Name entity recognition.

History of Seq to Seq Models :-



STAGE-1 : Encoder - Decoder Architecture



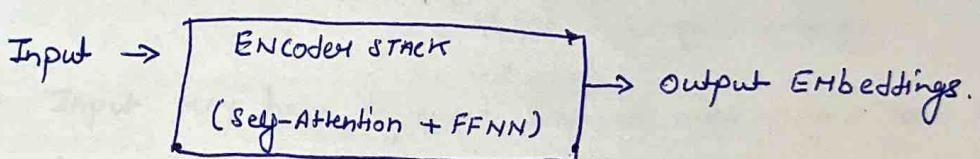
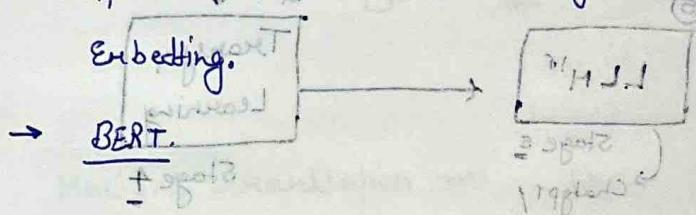
NOTE →

↳ ~~Encoder~~ → Extract features of input to predict

①. Encoder-only Models :-

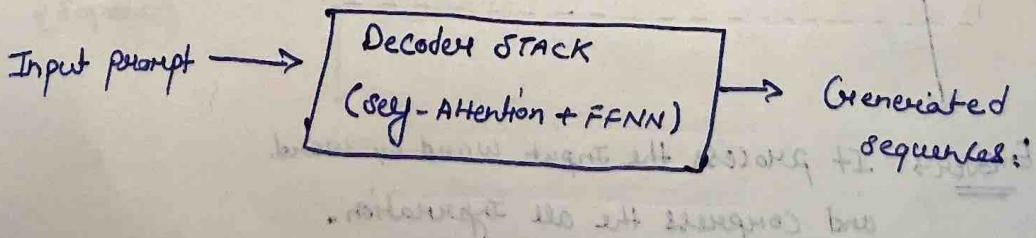
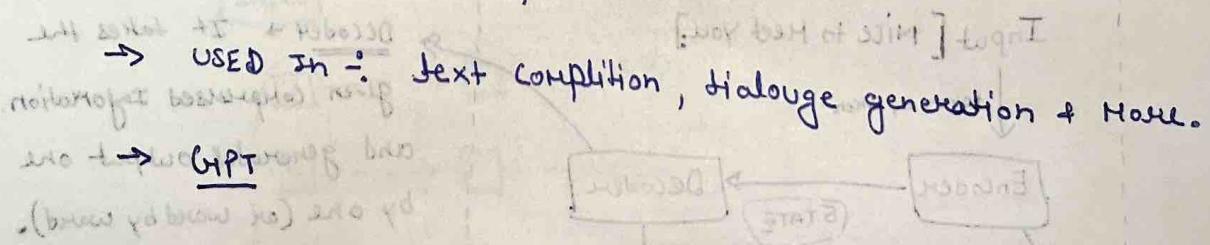
→ These are used for tasks that require understanding the context of a given Input, but do NOT require generating new sequences of Text.

→ Typical tasks include :- Classification, Entity recognition and Sentence



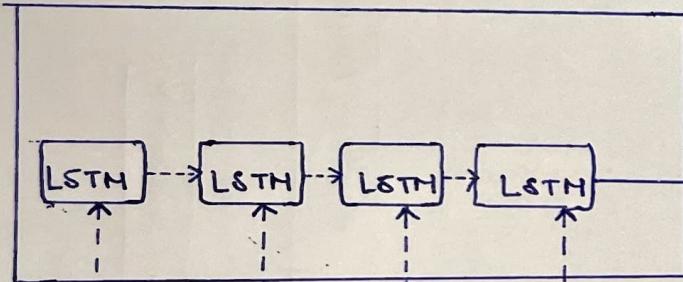
②. Decoder-only Models :-

→ They are used for generating Contextually text based output on a given prompt.



Encoder - Decoder Architecture :-

[ENCODER]

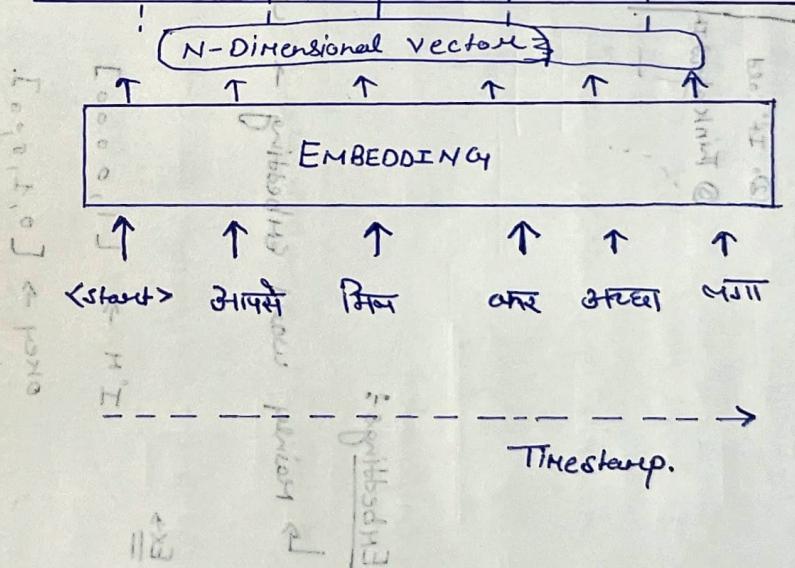
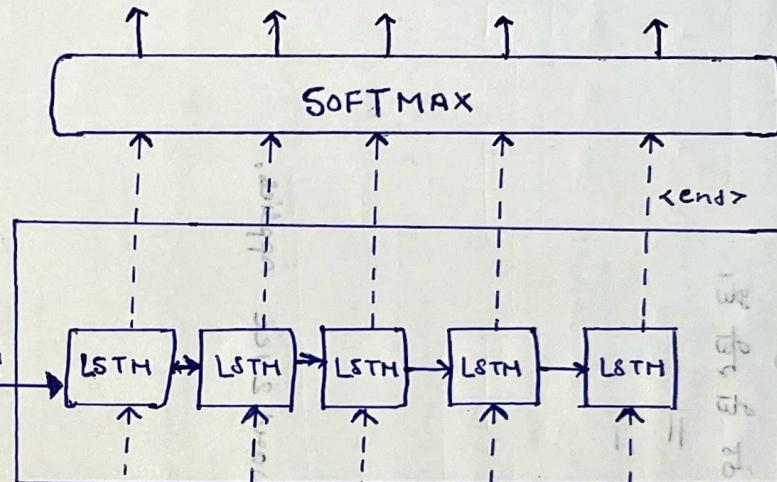


You can see it. In this case it is 3-D vector.

Nice to Meet You

----->
Timestamp

[DECODER]



STAGE-2 :- ATTENTION MECHANISM.

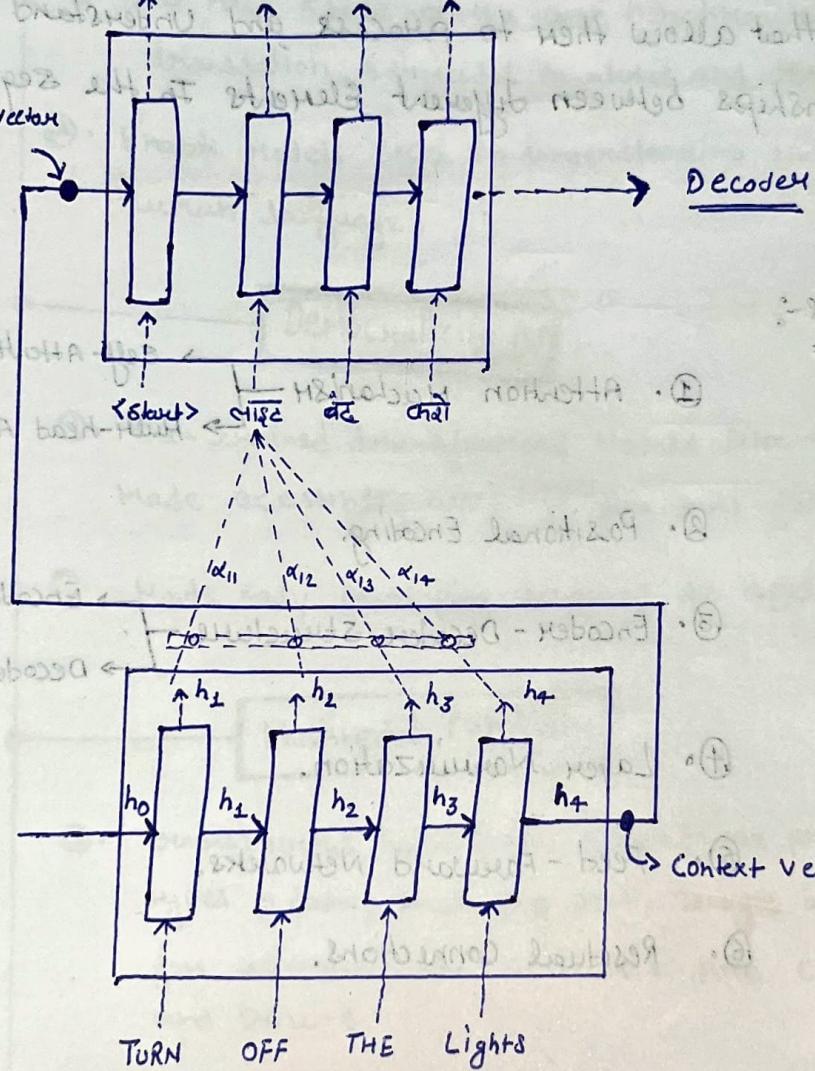
- Attention Mechanism is a technique in ML + NLP that allows models to focus on different parts of the input sequence when producing an output.
- It helps model weigh the importance of different input elements dynamically, rather than treating all input elements equally.

Key Components :-

- . query :- Represents the item for which attention is being computed.
- . key :- Represents the items in the input sequences to which the model is attending.
- . value :- Represents the actual information or content of the input sequences.

Formulas :-

- . Attention score calculation \Rightarrow $\boxed{\text{Score}(q, k) = q \cdot k^T}$
- . Attention Weights \Rightarrow $\boxed{\text{Weights} = \text{softmax}(\text{Score}(q, k))}$
- . Context Vector \Rightarrow $\boxed{\text{Context} = \text{Weights} \cdot V}$



* STAGE-3 :- Transformer

↳ Transformers are a type of neural network architecture designed to handle sequential data, such as text, by using mechanisms that allow them to process and understand the relationships between different elements in the sequence.



Key Concepts:-

- ①. Attention Mechanism
 - Self-Attention.
 - Multi-head Attention.

- ②. Positional Encoding.

- ③. Encoder - Decoder Structure
 - Encoder.
 - Decoder.

- ④. Layer Normalization.

- ⑤. Feed - Forward Networks.

- ⑥. Residual Connections.

Applications:- NLP, Vision, speech.

- # ①. Advantages ⇒ Parallelization + Scalability.

- ②. Disadvantages ⇒ Computationally Expensive + Data Hungry.

Impact of Transformers :-

INSIGHT

Revolution In NLP

- ① Set new state-of-the-art benchmarks in tasks like:-
Translation, Sentiment Analysis and Text-generation.
- ② Enable Models Excel in understanding and generating
human language.

Democratizing AI

- ① Pre-trained Transformers Models like → GPT and BERT are made accessible, allowing fine-tune for specific tasks.
- ② Made easy developing advanced AI applications.

Multimodal Capability

- ③ Transformers have been adapted for processing multiple types of data, including text, images and audio, leading to advancements in models like CLIP (for text-to-image) and DALL-E.

Acceleration of GenAI

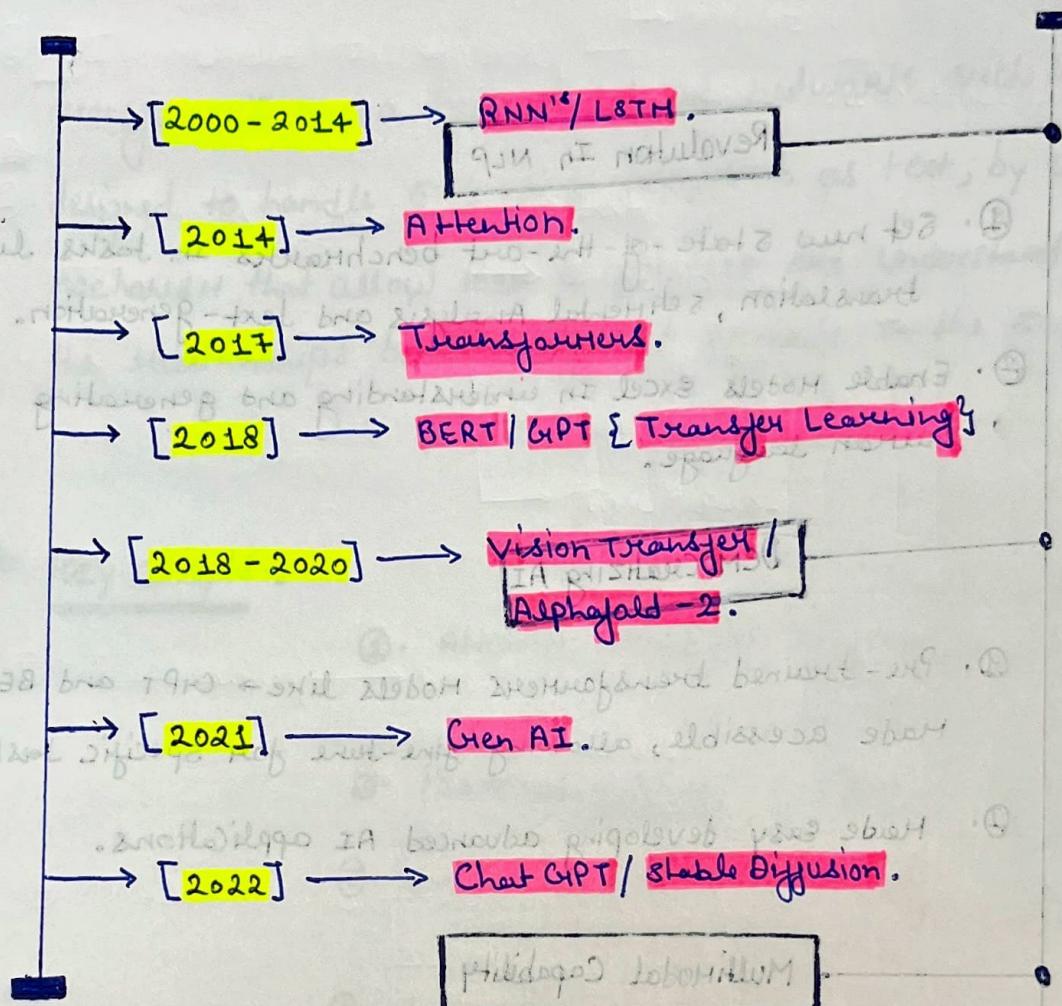
- ④ Transformers have driven rapid progress in GenAI, enabling the creation of high-quality text, images, music and more.
- ⑤ Models like GPT-3 and its successors show potential use of LLM's.

Unification of deep learning

- ⑥ Transformer Architecture can be applied across various domains, from NLP to vision, and even to scientific data, pushing towards a more unified deep learning framework.

The Timeline:

from now to tonight



AlphaFold

AlphaFold

AlphaFold

Advantages :-

- ①. Scalability.
- ②. Transfer Learning.
- ③. Multimodal.
- ④. flexible Architecture.
 - Encoder only (BERT).
 - Decoder only (GPT).
- ⑤. Ecosystem.
 - huggingface.
 - video / blogs.
- ⑥. Integrated AI Media.

Disadvantages :-

- ①. High Computation (GPU).
- ②. Large Data { text → unsupervised labeling }
- ③. overfitting.
- ④. Energy Consumption.
- ⑤. Interpreter.
- ⑥. Bias → Ethical Concerns.

Future :-

①. Improvements In efficiency. → Position,
→ Quantization.
→ Knowledge Distribution.

②. Multimodal Capabilities → text, general report
→ Image / speech, laboratory
→ Sensory data.
→ Bio-Heuristic feedback.

③. Responsibility → Bias.
→ Ethical Concerns. IA botany IIT

④. Domain specific → Doctor GPT.
→ Teacher GPT.
→ Legal GPT.

⑤. Multi-lingual → English, Hindi, French, etc.

⑥. Interpretability → Critical domains like finance, medical etc.

Till now Transformers are Black-box Model.

①

Attention Mechanism

A) Self-Attention :-

↳ Self-Attention is a mechanism used in neural networks,

Especially in models like transformers to help the network focus on different parts of the input text more effectively.

Simple Explanation :- Imagine you're reading a sentence, and you want to understand the meaning of a word by looking at the surrounding words.

Self-Attention allows the model to look at each word. Ex - "Cat", the model might pay attention to "sat" and "on" because they are more relevant to understanding the context around "cat".

Steps In Self-Attention :-

①. Input Words :- Let's say we have a sentence with words as input.

②. Create Attention Scores :- For each word, the model calculates a score (or weight) that shows how much attention to give to every other word.

③. Weighted sum :- The model then multiplies each word's embedding by its corresponding attention score and sums them up to get a new representation of the word.

Workflow :-

Hierarchical architecture

Step 1 :- Input Words \Rightarrow A sentence like "The Cat Sat" Is Input the Model.

Step 2 :- Attention Scores \Rightarrow A Table Shows the attention scores for each pair of words.

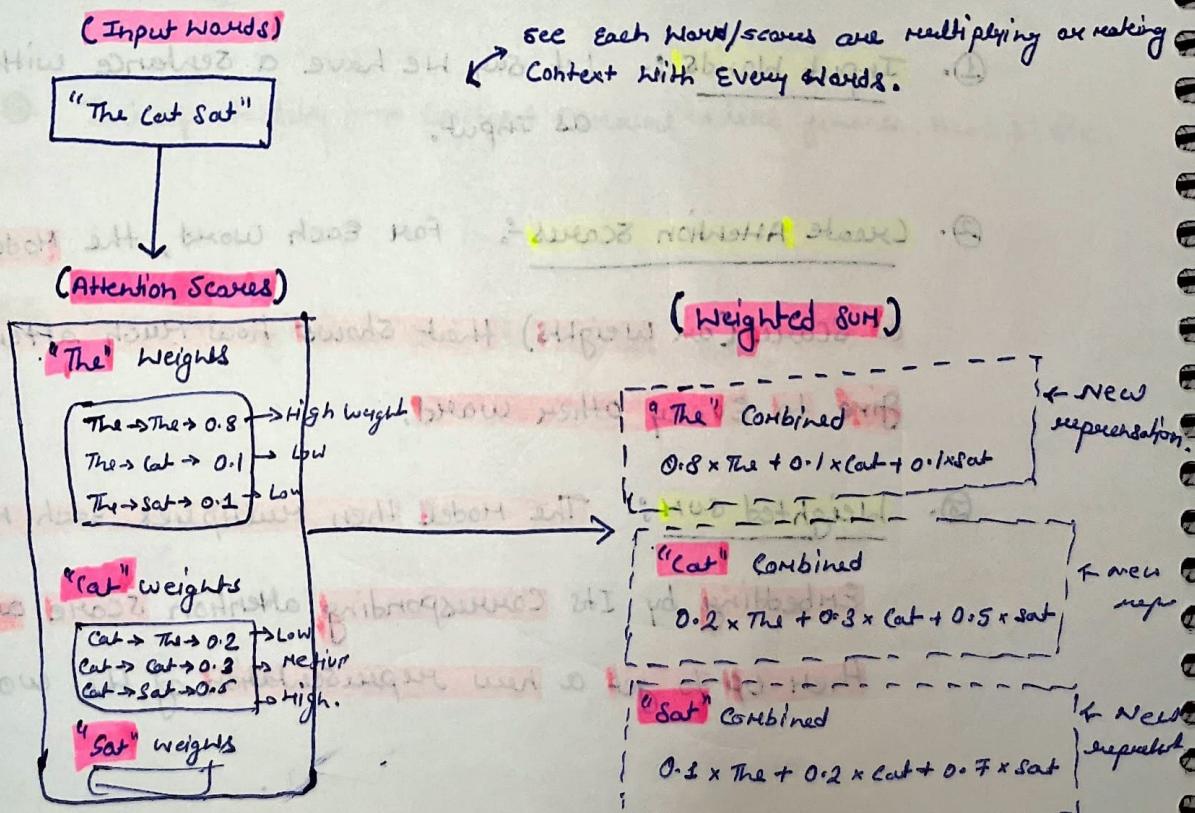
Ex -

①. "The" pays most attention to "The" and some to "Cat" and "Sat".

②. "Cat" pays more attention to "Sat" than "The".

Step 3 :- Weighted sum \Rightarrow These scores are used to create a new representation of each word, which takes into account the context provided by other words.

Diagram :-



First principle Approach :-

MONEY BANK grows

river bank flows

- * Notice both sentence have "bank" word with different meaning.

$$\text{bank} \rightarrow 0.3 \text{ Money} + 0.7 \text{ bank} + 0.1 \text{ grows}$$



$$\text{bank} \rightarrow 0.5 \text{ river} + 0.4 \text{ bank} + 0.1 \text{ flows}$$



- * See both embedding of bank now different because of Contextual Embeddings.

Ex :-

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

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

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



Word Embedding (S_1)

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

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

$$\text{flows} = 0.4 \text{ river} + 0.03 \text{ bank} + 0.59 \text{ flows}$$



Word Embedding (S_2)

n -dimension (here 3 words so $(3 \times 3) \rightarrow 3$ -dimension)

(Chew)

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



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

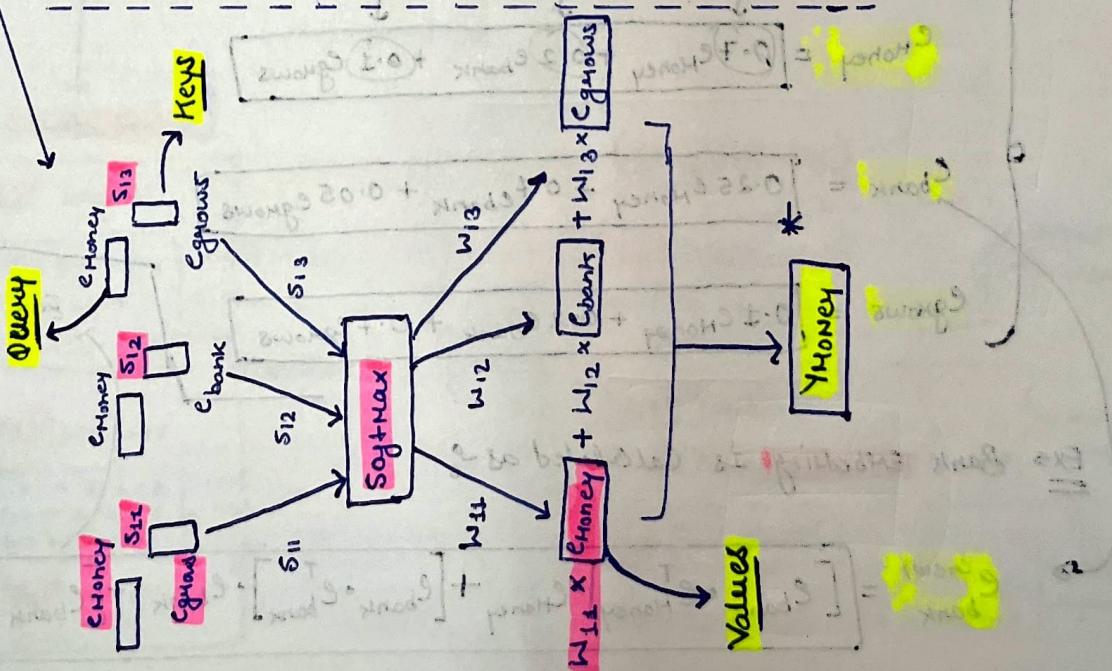
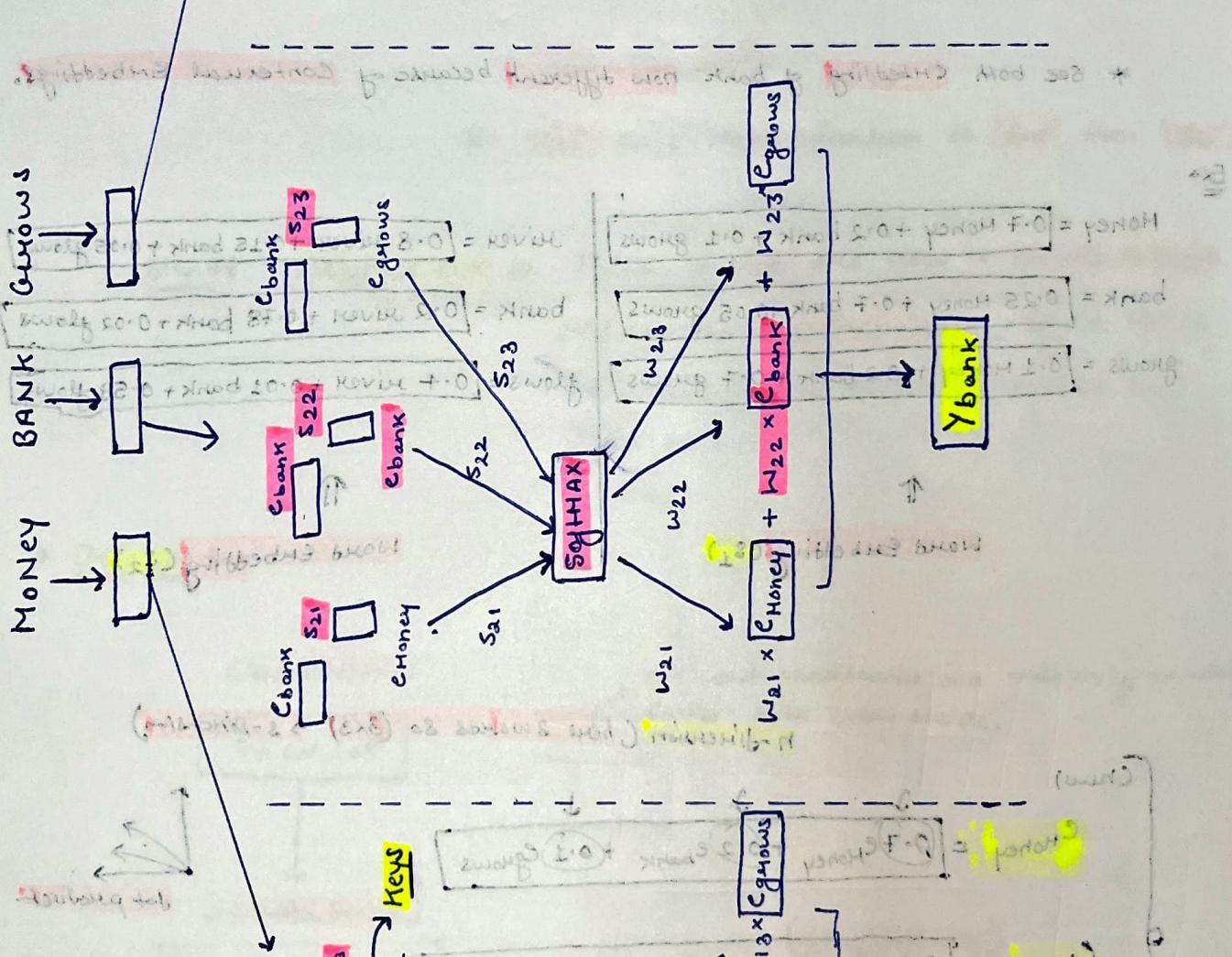
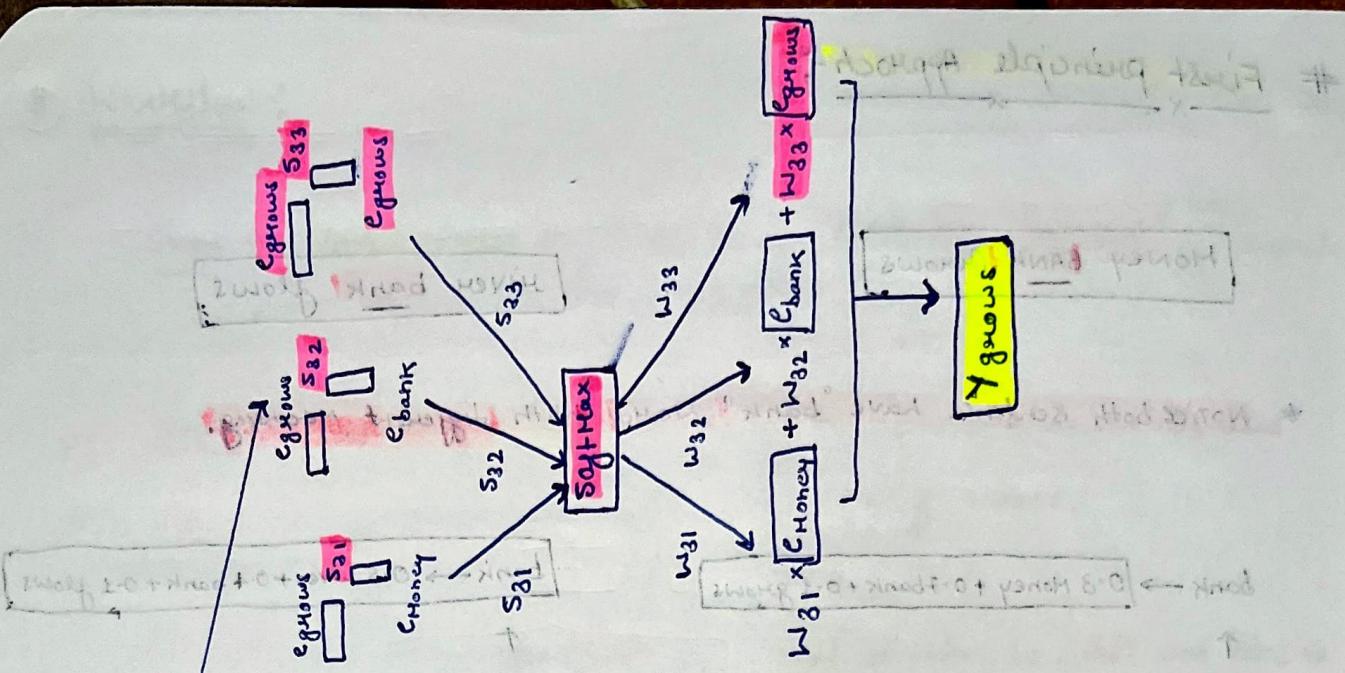
similarity

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



Ex :- Bank Embedding is calculated as :-

$$e_{\text{bank}}^{\text{new}} = [e_{\text{bank}} \cdot e_{\text{Money}}^T] \cdot e_{\text{Money}} + [e_{\text{bank}} \cdot e_{\text{bank}}^T] \cdot e_{\text{bank}} + [e_{\text{bank}} \cdot e_{\text{grows}}^T] \cdot e_{\text{grows}}$$



Key operations :-

V \rightarrow Input \rightarrow Output

→ The main steps in the self-Attention Mechanism are :-
1. Multi-head Self-Attention
2. Layer Norm
3. Residual Connection + Layer Norm

①. "Query", "Key" & "Value" Matrices :- Each token's embedding is linearly transformed into three vectors.

(A). Query vector (Q) :-

→ Represents the current token being processed.

(B). Key vector (K) :-

→ Represents other tokens to compare with.

(C). Value vector (V) :-

→ Represents the actual token embeddings used in the weighted sum.

Weighted SUM.

②. Dot product b/w Q + K :-

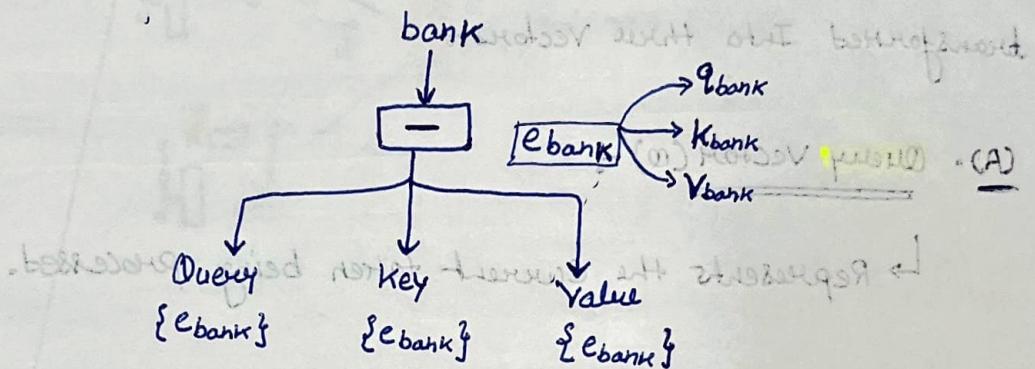
→ The similarity between the query and key vectors is computed using a dot product.

③. SoftMax :- The resulting dot products are normalized using a softmax function to compute attention weights.

→ This helps determine how much attention should be paid to each key.

④ • Weighted sum of V :-

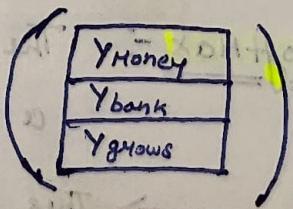
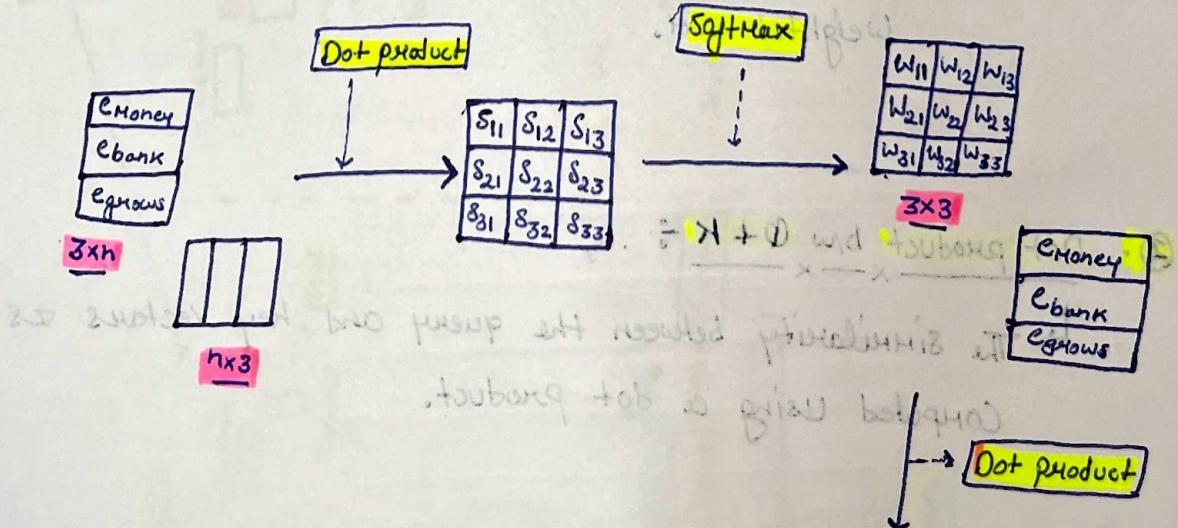
↳ The attention weights are applied to the value vectors to get the final representation of each token.



→ (A) Query → (B)

→ (C) Value

Visualization :-



3xh

Parallelism :-

↳ Self-Attention have ability to process operations In parallel, making transformers well-suited for large-scale training on GPU's or TPU.

①. Matrix Multiplications :-

↳ Instead of processing tokens one by one, The self-attention mechanism can operate on the entire sequences of tokens In parallel.

→ This is possible because Inputs (queries, keys and values) are Matrices and operations are dot-products.

②. Attention score Calculation:-

↳ Dot product are multiplying by Query Matrix (Q) with the transpose of the Key Matrix (K^T).

$$\text{Attention Score} = Q \cdot K^T$$

③. Sigmoid + Scaling :-

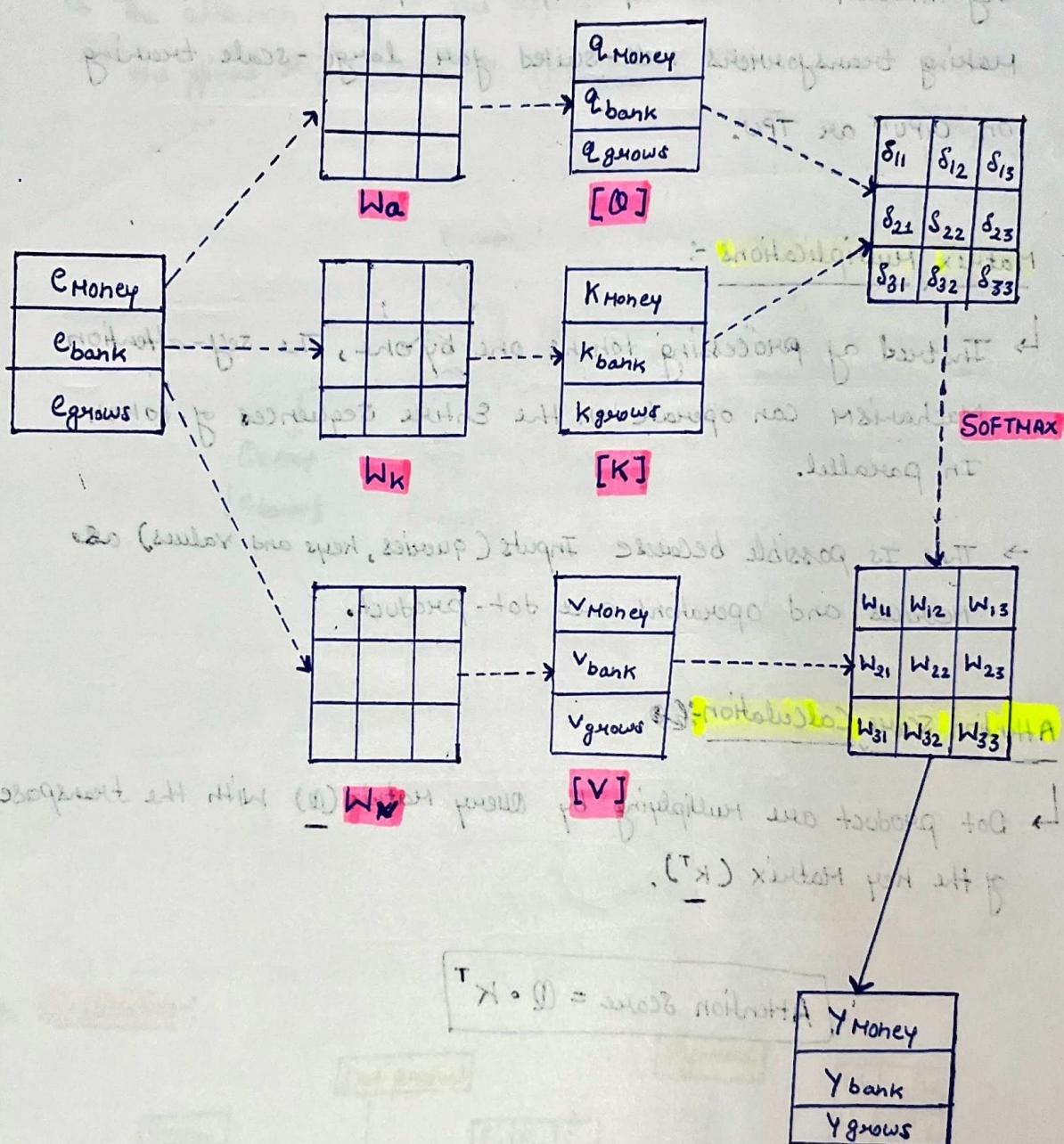
↳ The attention scores for all tokens can also be scaled and normalized In parallel.

$$\text{Normalized Weights} = \text{sigmoid}[Q \cdot K^T]$$

④. Weighted sum of values :- The weighted sum is calculated for all tokens In one go by multiplying the sigmoid output matrix with the Value matrix (V).

$$\text{Context vector} = [(Q \cdot K^T) \cdot V]$$

Complete Seq-Attention Structure :-



$$[Q \cdot K^T] \cdot V = \text{Context Vector}$$

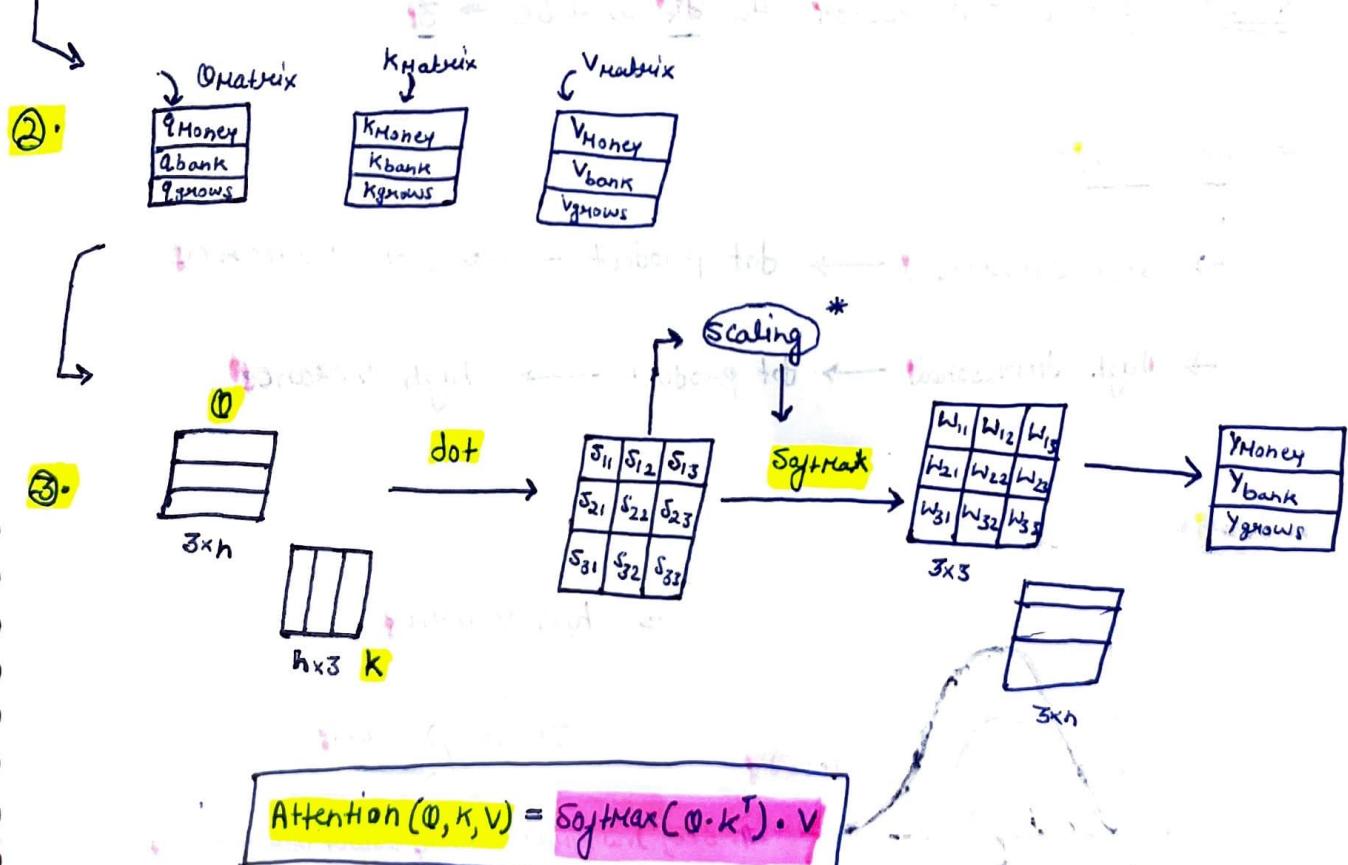
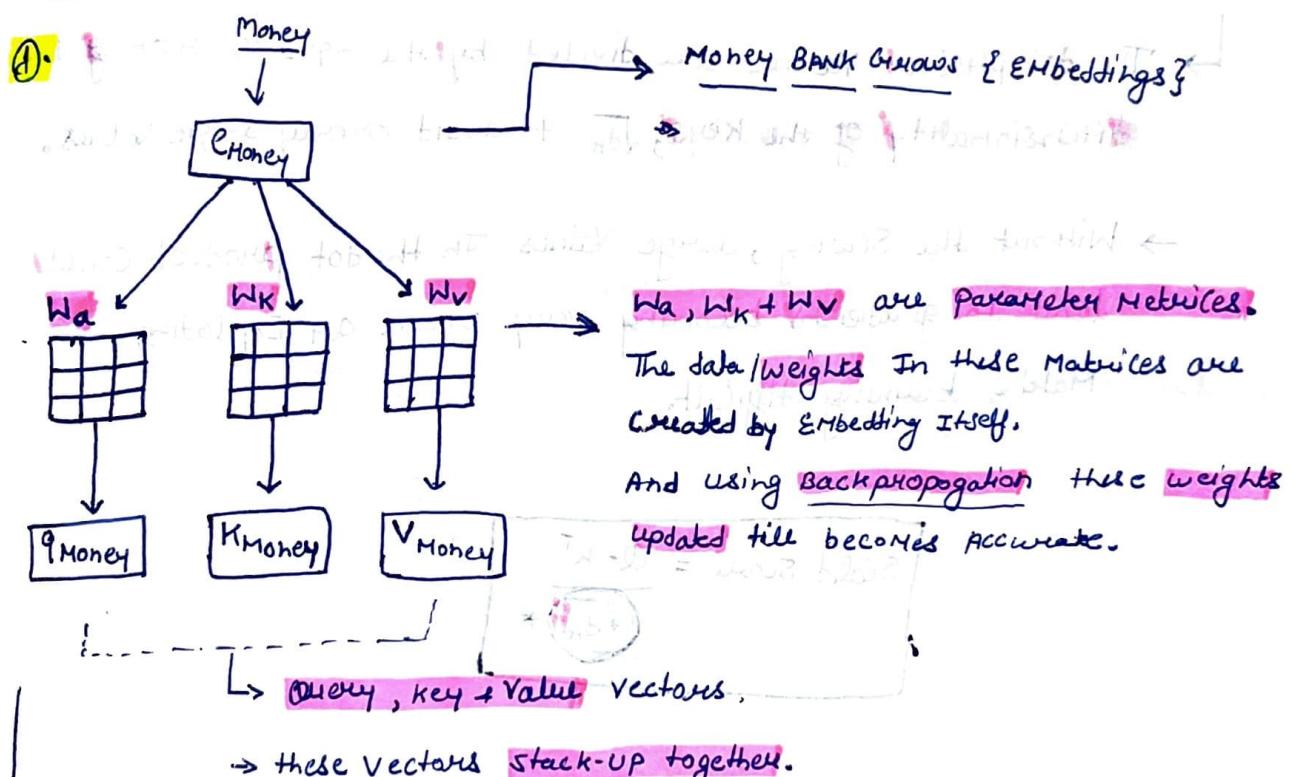
Context also helps to calculate the weighted sum of words in context with attention weights.

$$[V \cdot (Q \cdot K^T)] = \text{Context Vector}$$

Scaled Dot Product - Attention :

→ Scaled Dot-product Attention is a key mechanism used in transformer Models, especially in NLP tasks like, machine translation, text generation and summarization.

Flow :



Scaled DOT-Product :- Formula.

$$\text{Attention}(Q, K, V) = \text{SoftMax}\left(\frac{Q \cdot K^T}{\sqrt{dk}}\right) \cdot V$$

Attention score is the dot product of query and key, scaled by \sqrt{dk} .

* Scaling (\sqrt{dk}) :-

→ The dot product values are divided by the square root of the dimensionality of the keys, \sqrt{dk} to avoid overly large values.

→ Without the scaling, large values in the dot product could lead to gradients becoming very small or exploding, making training difficult.

$$\text{Scaled Score} = \frac{Q \cdot K^T}{\sqrt{dk}} *$$

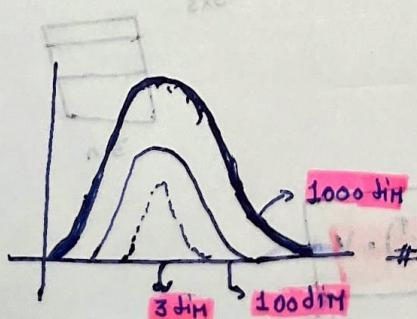
NOTE → for a 3-D vector, the dk will be $\rightarrow 3$

Key points :-

→ low dimensional → dot product → low variance

→ high dimensional → dot product → high variance.

Graph →



It's a problem.

See, low dimension have lower variance and vice-versa.

Multi-Head Attention :-

→ Multi-head attention is a key mechanism in the Transformer architecture that allows the model to capture different types of relationships between tokens in a sequence.

→ Instead of using a single attention function, multi-head attention applies multiple attention mechanisms (or heads) in parallel, allowing the model to focus on various parts of the sequences simultaneously.

* How It Works:-

①. Each attention head computes self-attention independently using different learned linear projections of the original queries, keys and values.

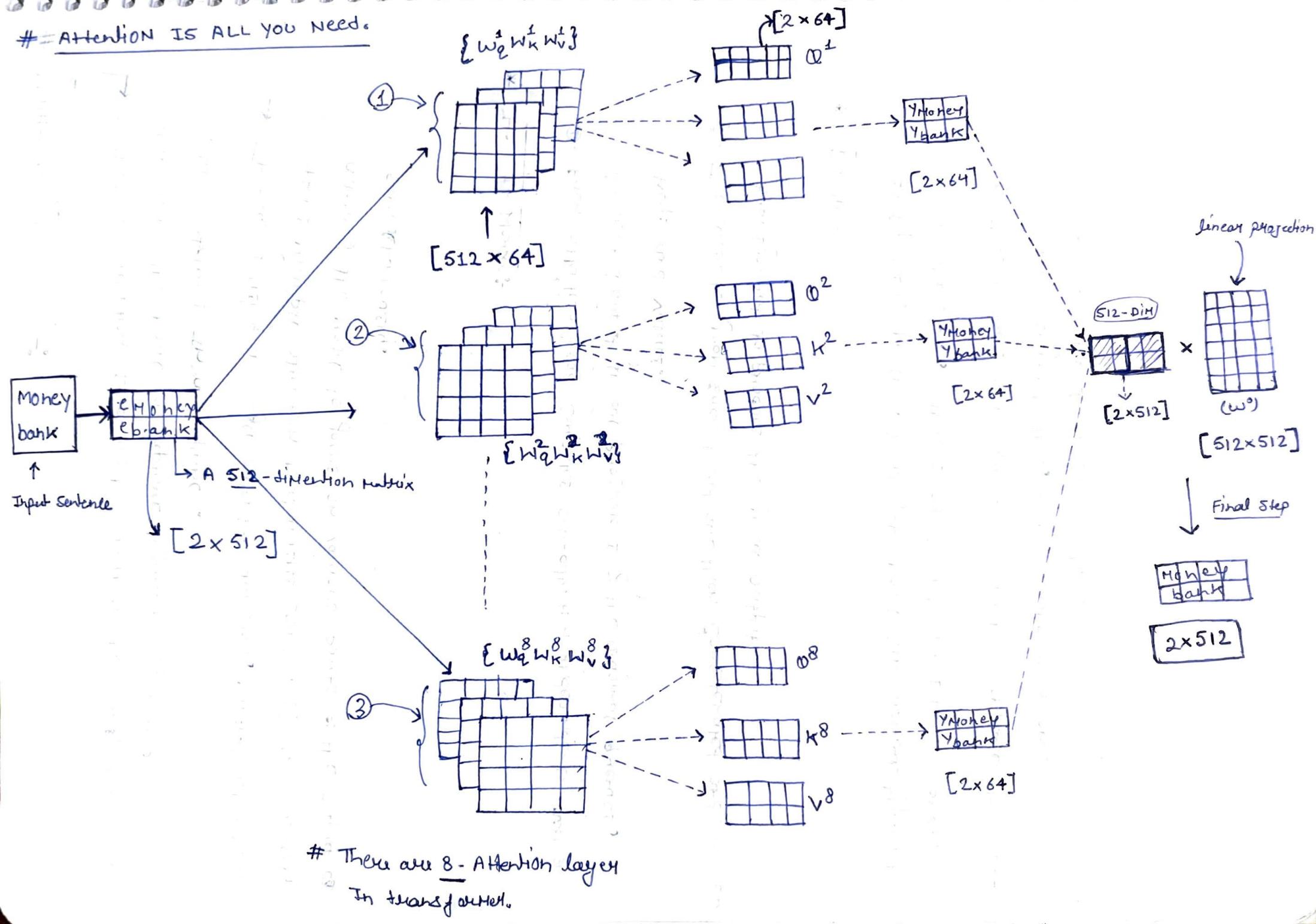
②. The results from each head are then concatenated and linearly transformed to produce the final output.

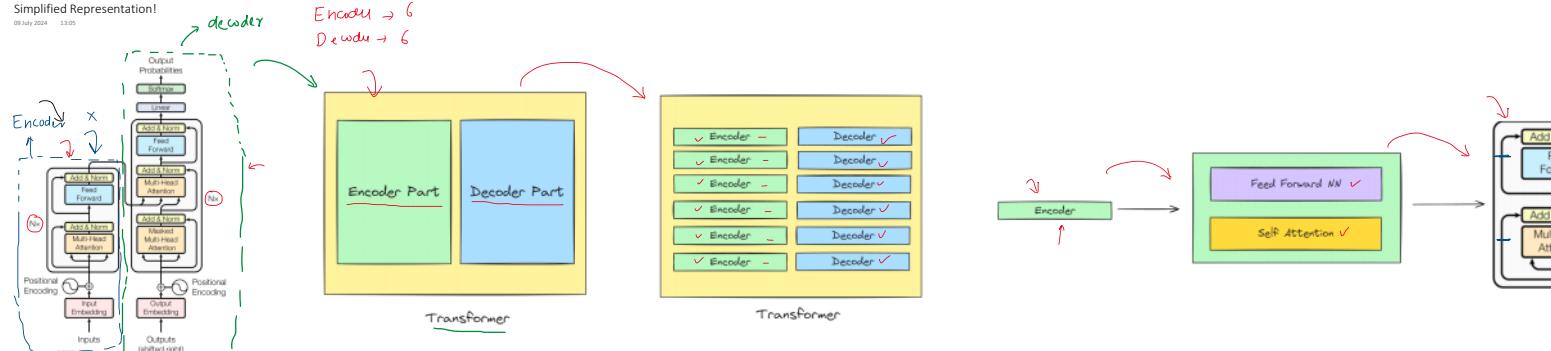
* Benefits:-

①. By using multiple heads, the model can capture different relationship patterns or dependencies between words.

②. This provides a richer representation compared to single-head self-attention, which might focus on only one aspect of the token relationships.

ATTENTION IS ALL YOU NEED





Encoder Architecture

09 July 2024 16:31

