

# NLP+ Deep Learning(Generative AI)

# Natural Language Processing – Deep Learning Modules



Word Embedding

LSTM-Long Short -Term Memory

Sequence 2 Sequence

Transformer

# Abdul Kalam is a Great man

Abdul  
Kalam  
is  
a  
Great  
man

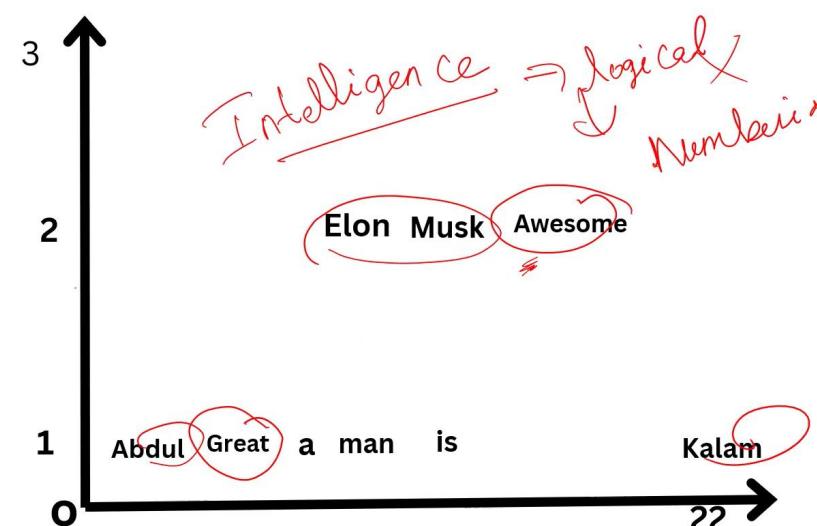
- 1
- 22
- 13
- 4
- 5
- 6

# Elon Musk is Awesome

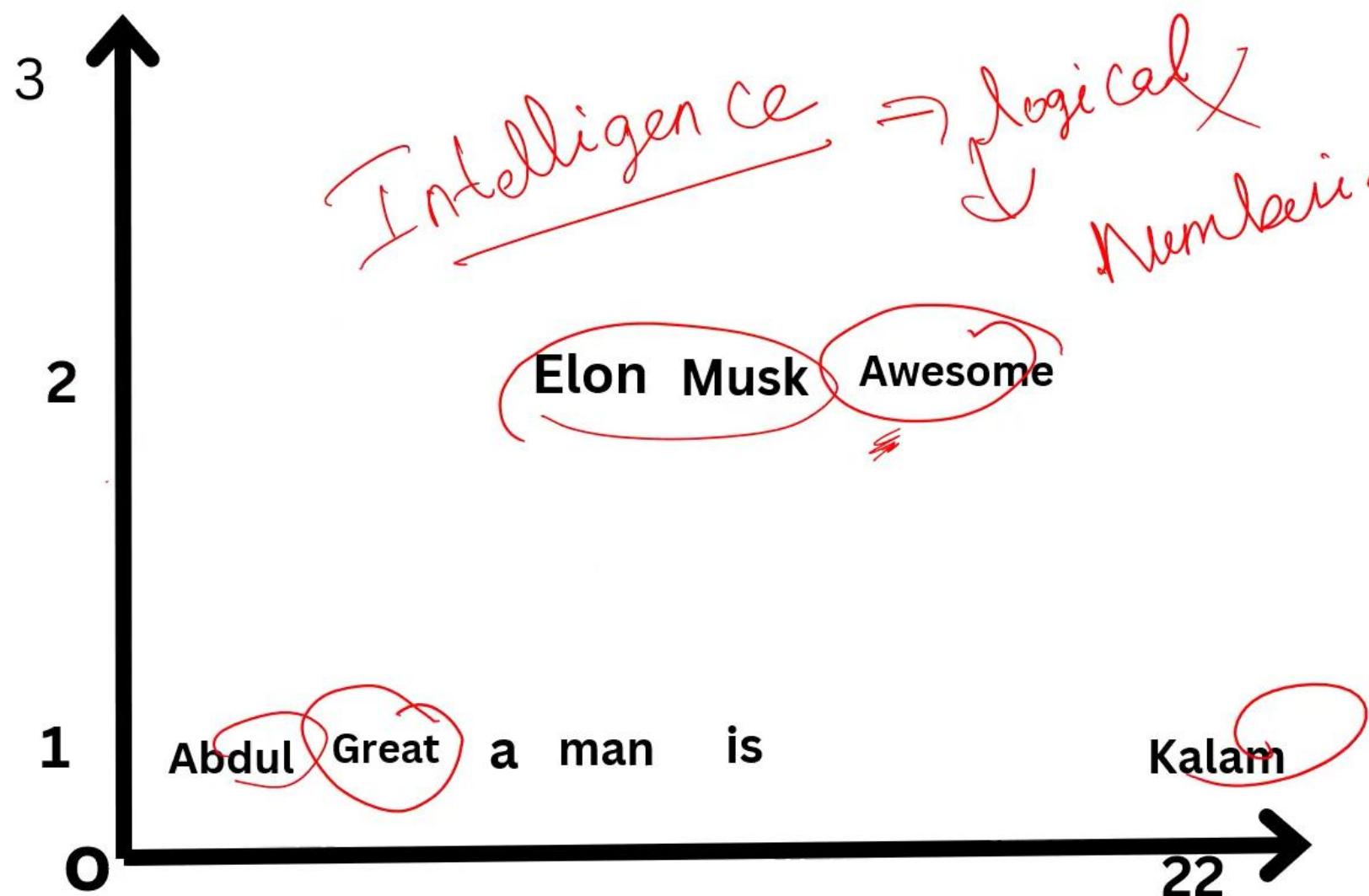
Elon  
Musk  
is  
Awesome

- 11
- 12
- 13
- 19

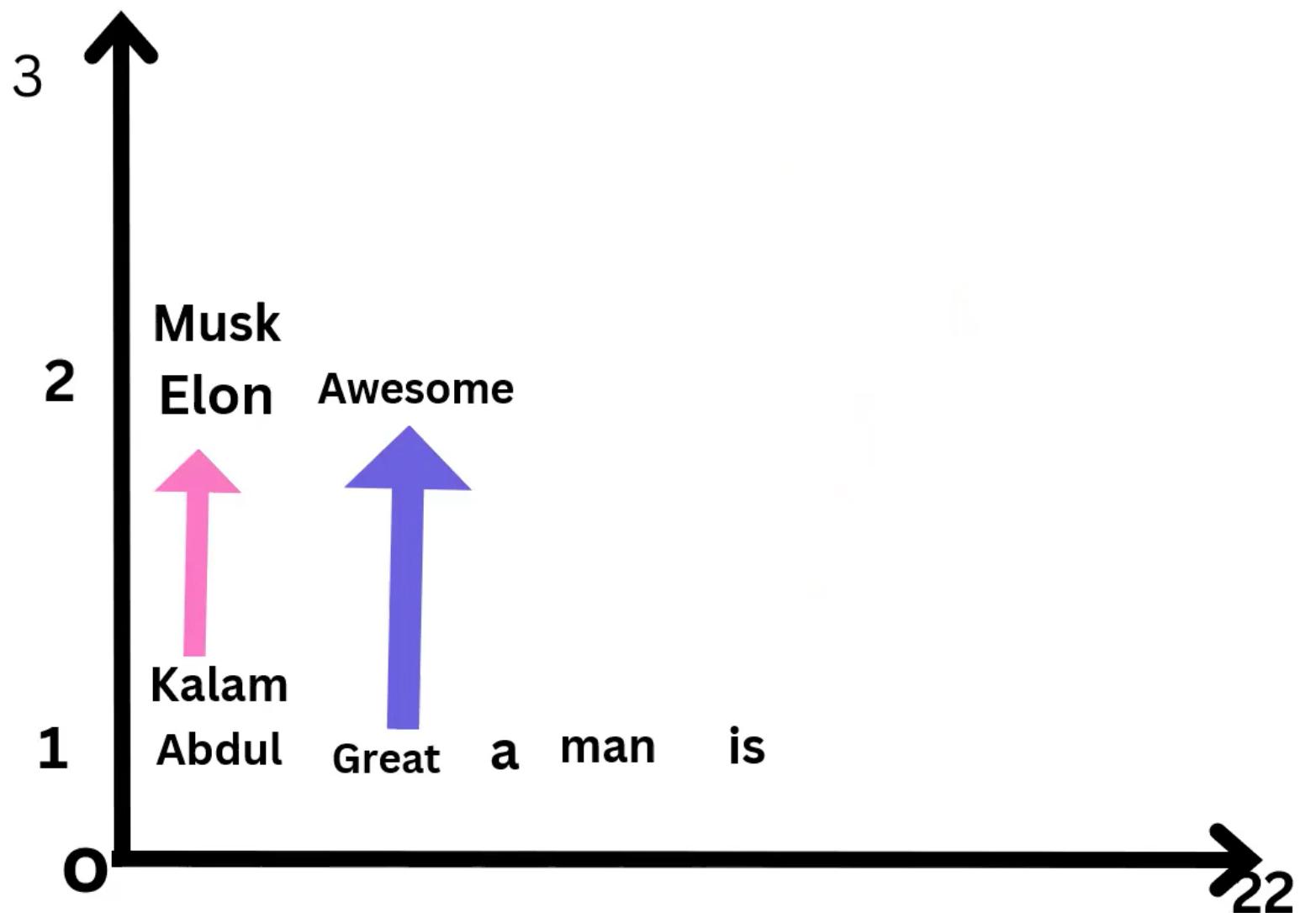
Why Word Embedding is needed?



## Why Word Embedding is needed?

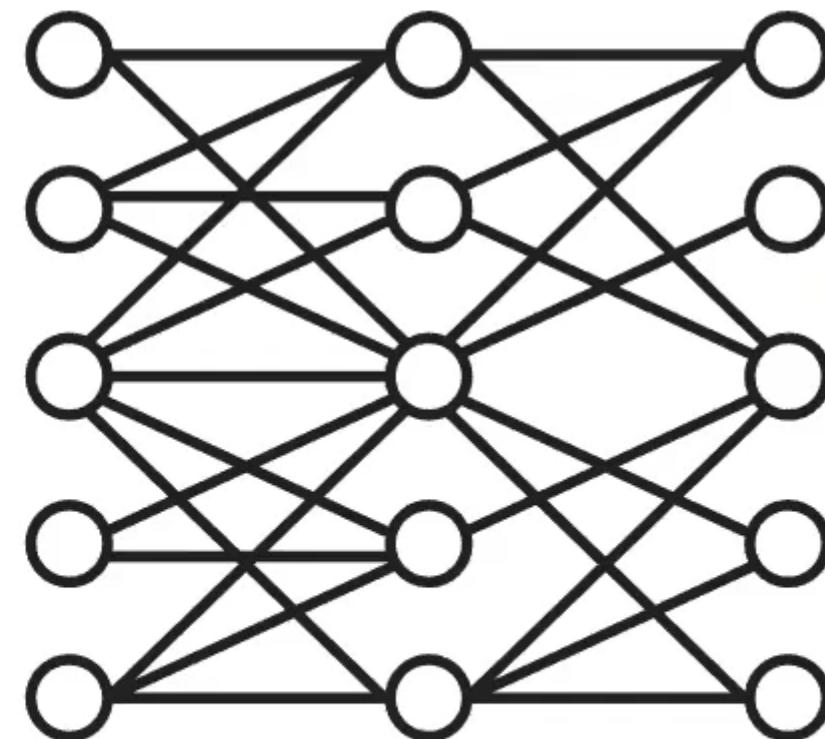


## Why Word Embedding is needed?



# Word Embeddings using Neural Network

Using Neural  
Network



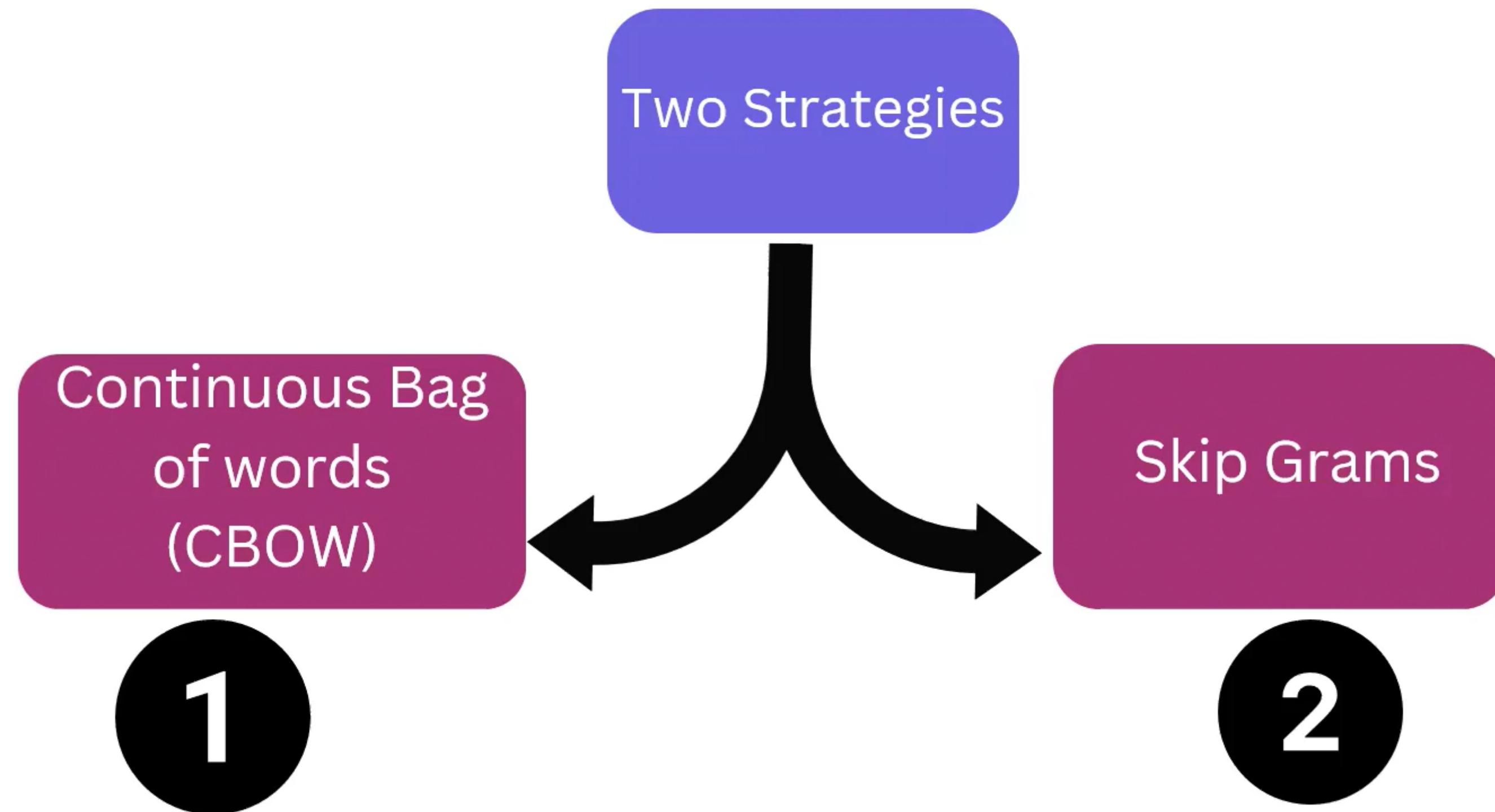
We can create  
meaningful  
numbers

word  $\rightarrow$  Vec

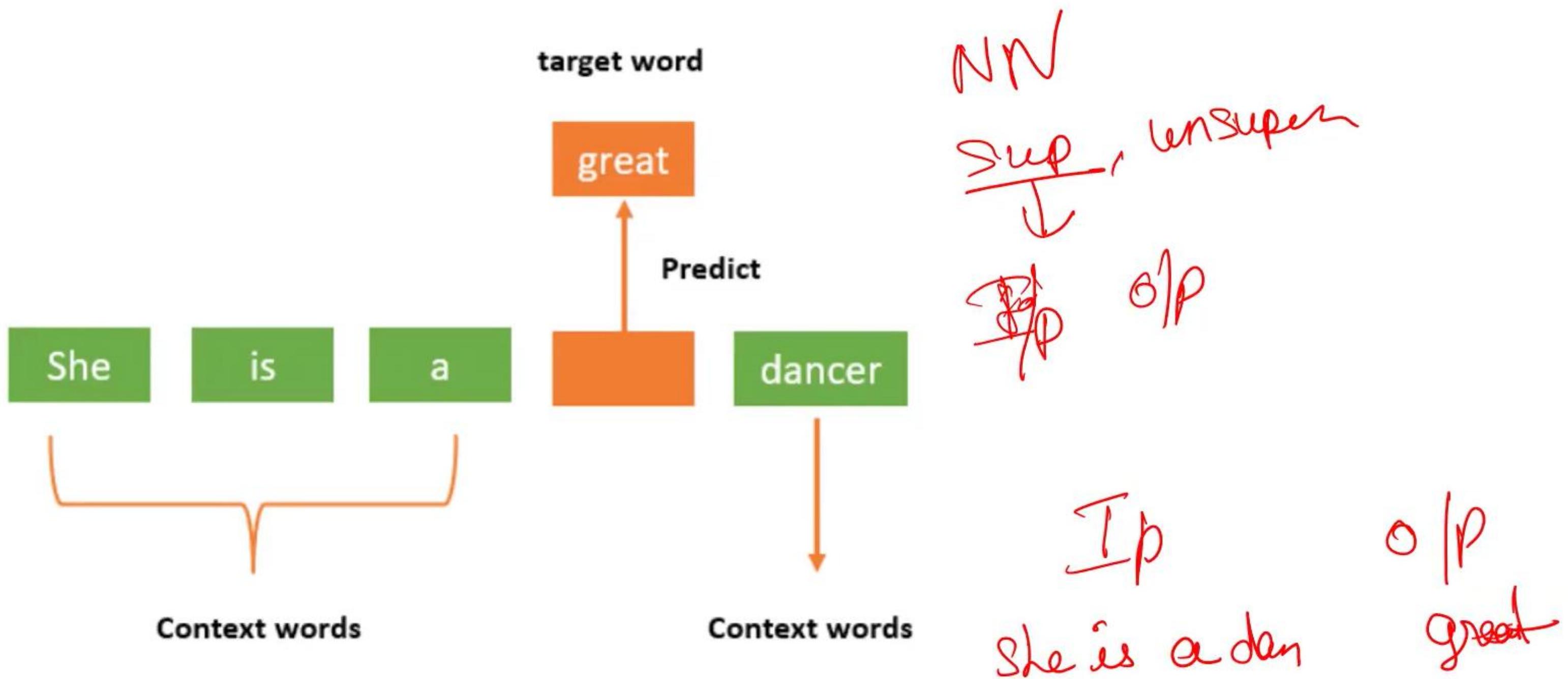
fastest      by  $\rightarrow$  pain  
                pre

words  $\rightarrow$  meanings

# Word Embedding-Word2Vec



# Word Embedding-Word2Vec-CBOW



## Continuous Bag of words (CBOW)

~~The quick brown fox jumps over the lazy dog~~

Context Word, Target Word, Window Size

Context Word: (-2, Target Word, +2)

The quick brown fox jumps over the lazy dog

(-2) (-1) (T) (1) (2)

Window 3:

Context words ( 'The', 'quick', 'fox', 'jumps' ) -  
Target word ( 'brown' )

The quick brown fox jumps over the lazy dog  
(-2)(-1) (T) (1) (2)

Window 1:

Context words ('N/A', 'N/A', quick, brown) -  
Target word ('The') ~~top~~

Window 2

Context words (NA, the, brown, fox)  
Target word ("quick")

I/P

O/P

The quick brown fox jumps over the lazy dog  
(-2) (-1) (T) (1) (2)

N/A,N/A, quick, brown  
N/A, The, brown, fox  
The, quick, fox, jumps

The  
quick  
brown

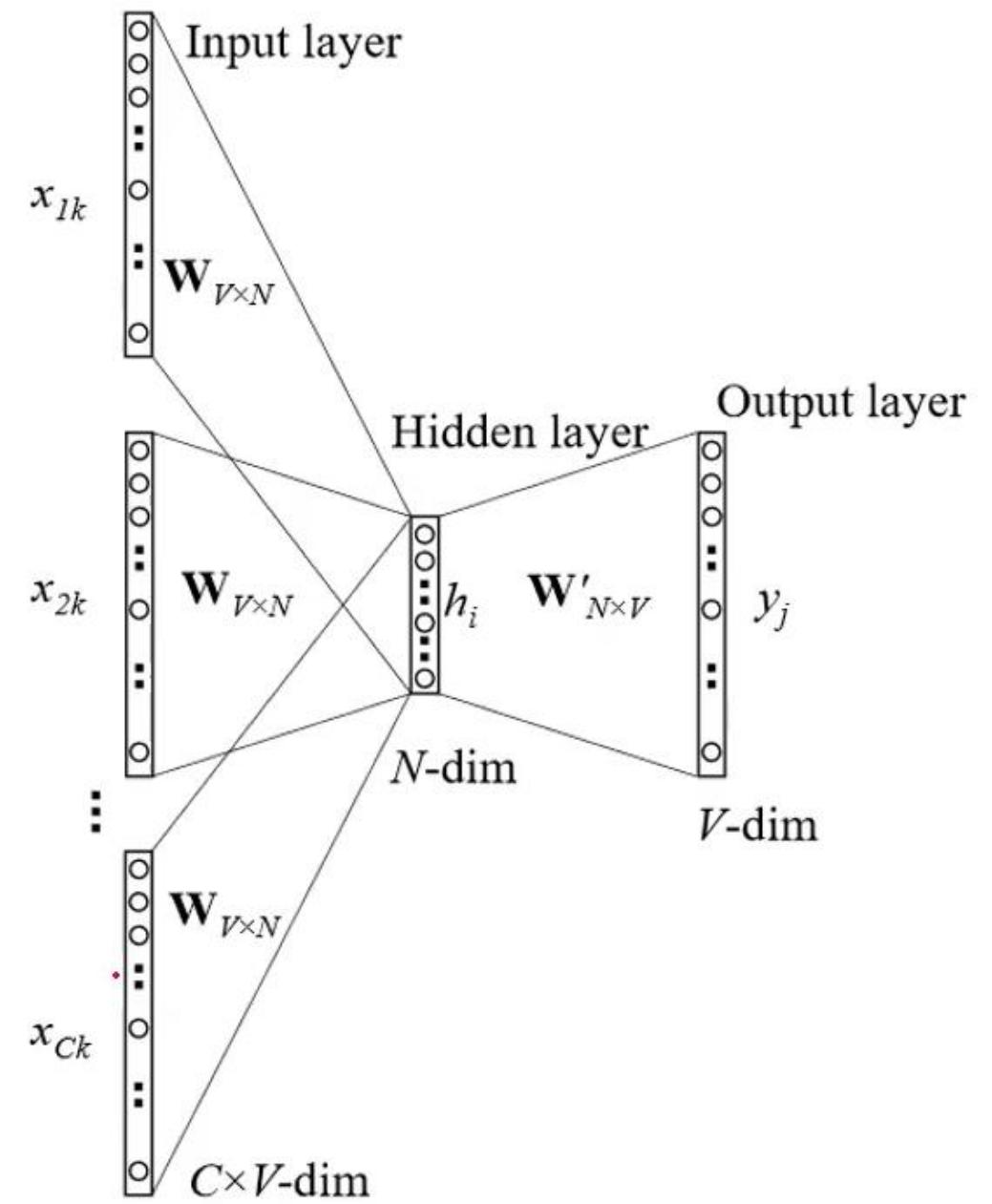
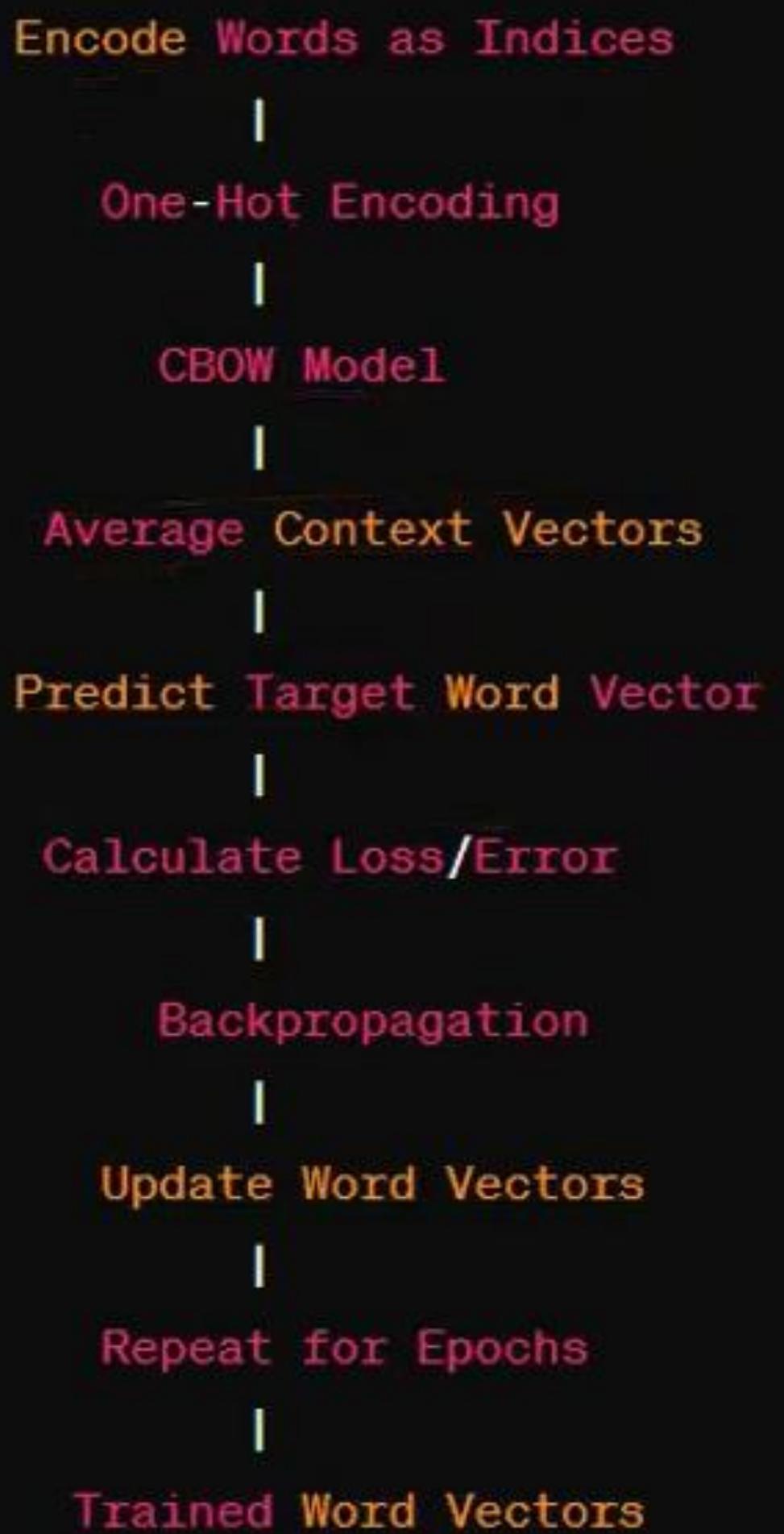
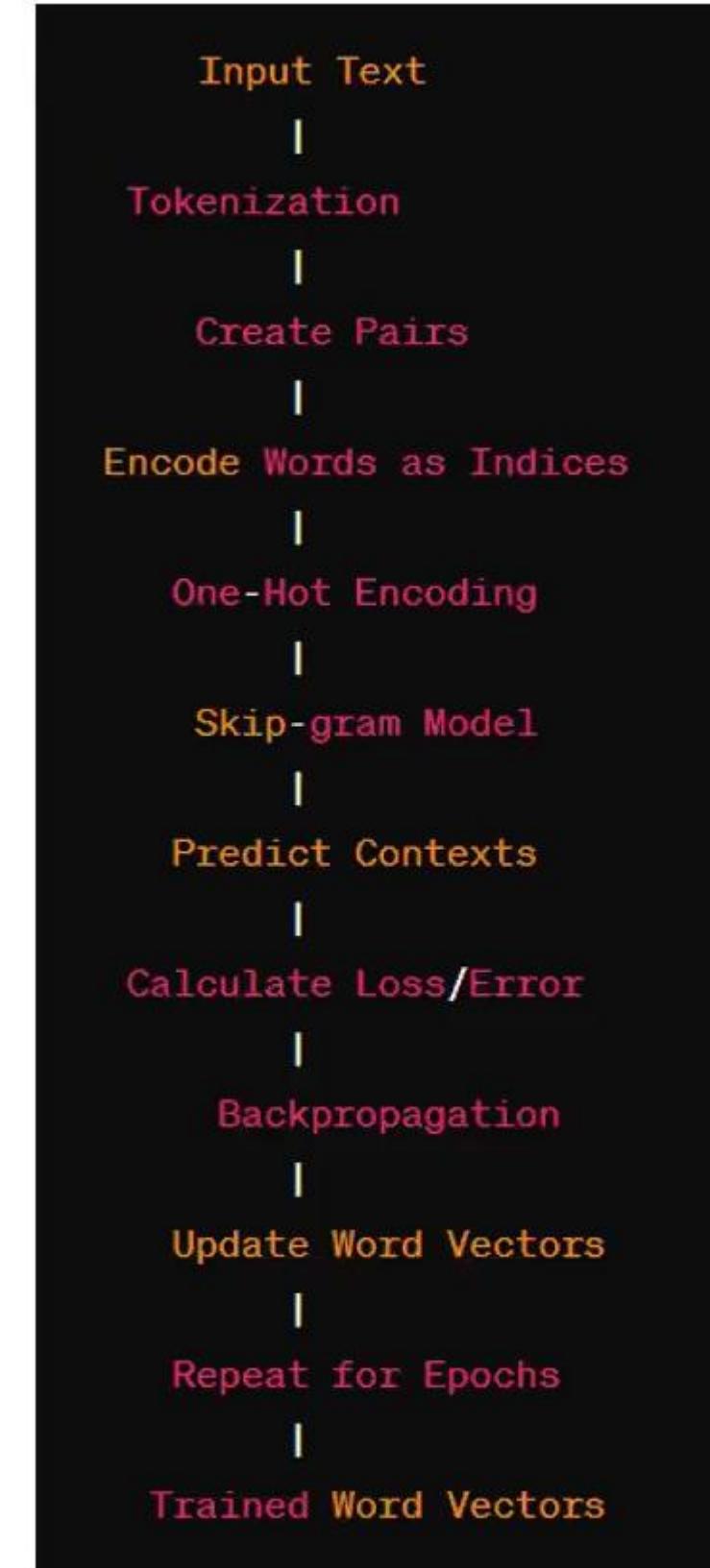
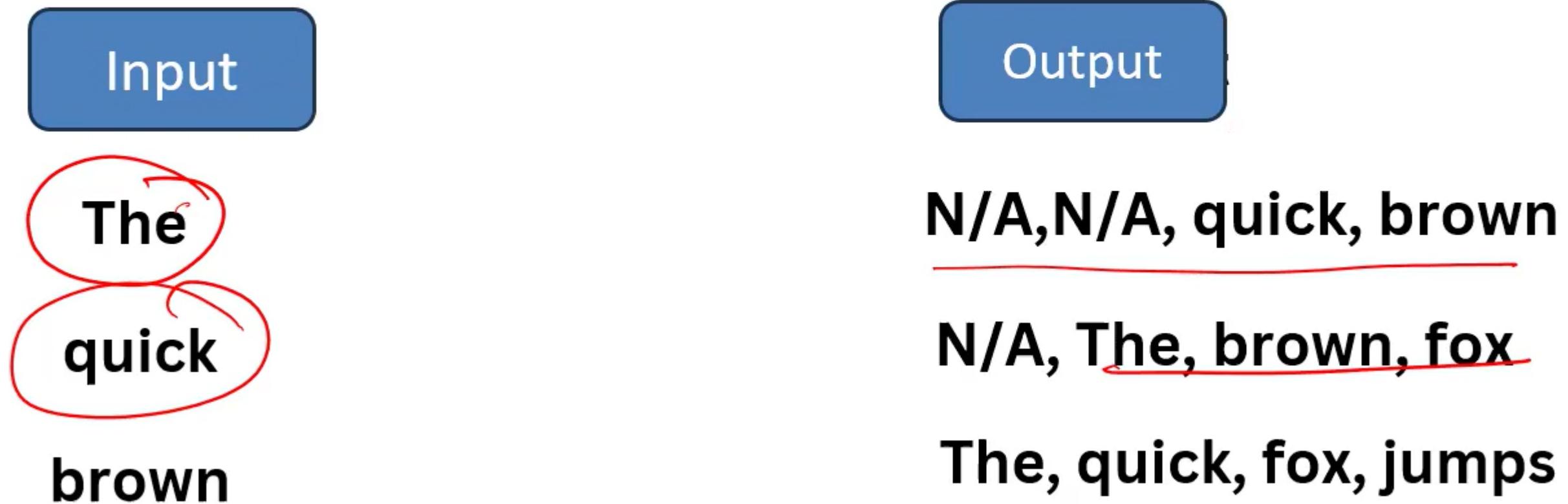


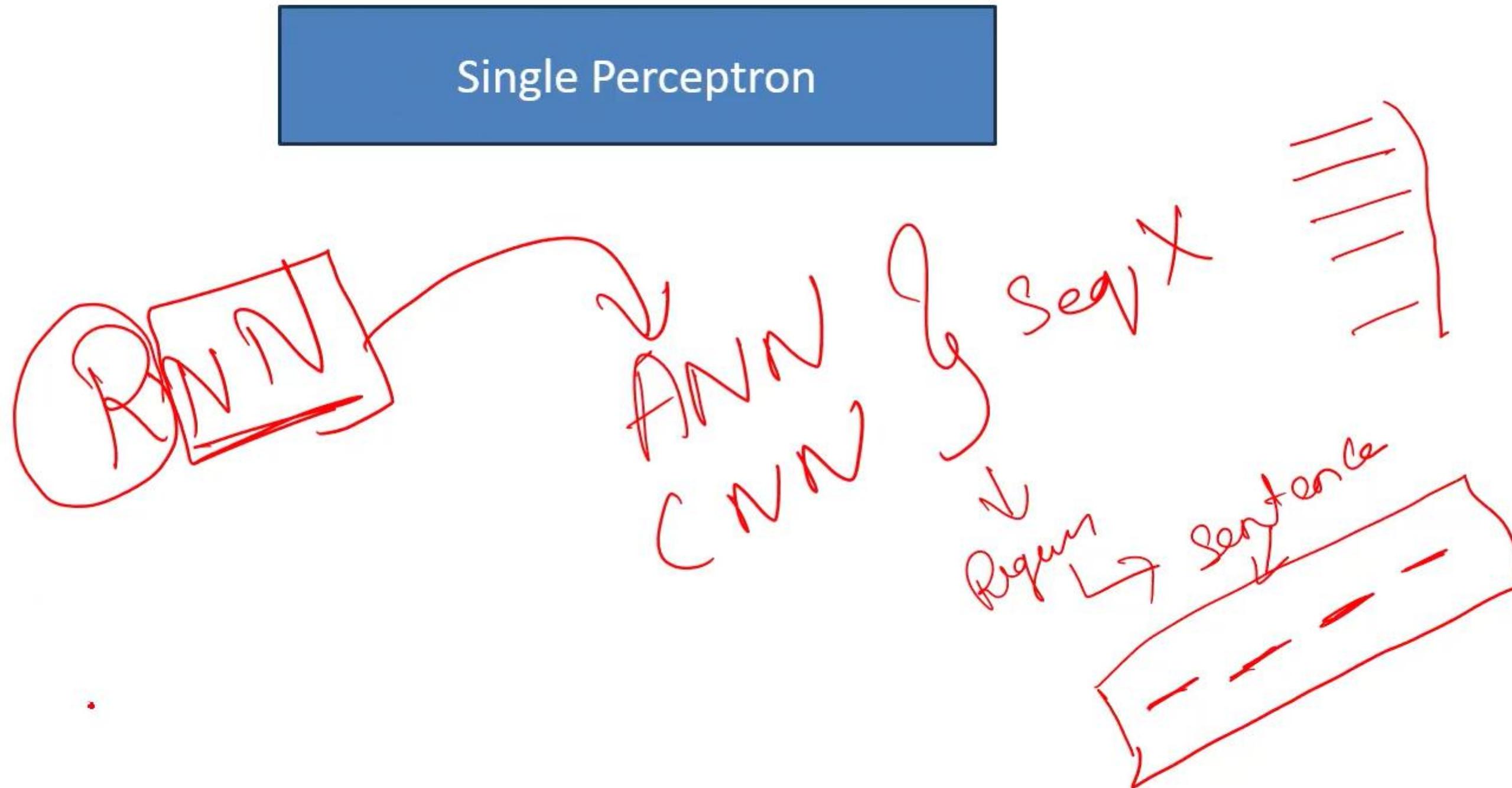
Figure 2: Continuous bag-of-word model



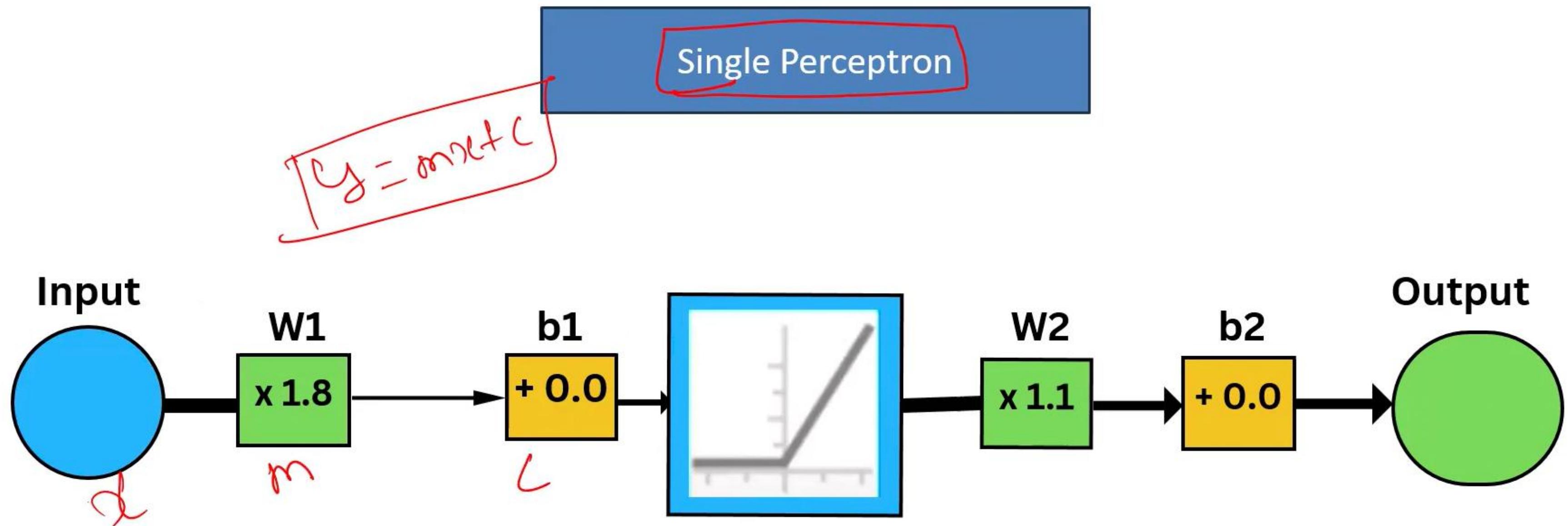
# Word Embedding-Word2Vec-Skip Gram

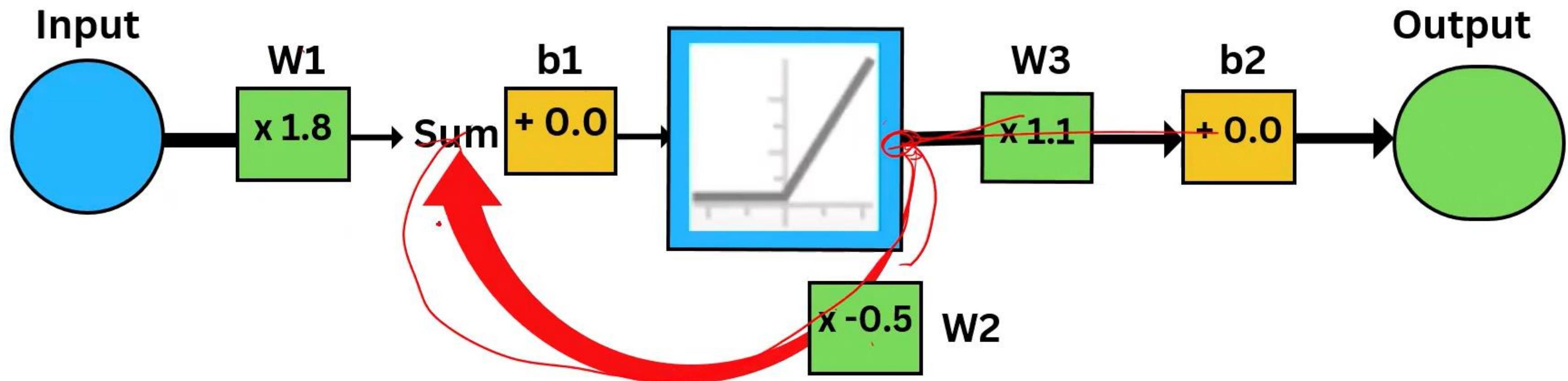


# Long Short-Term Memory- (LSTM)

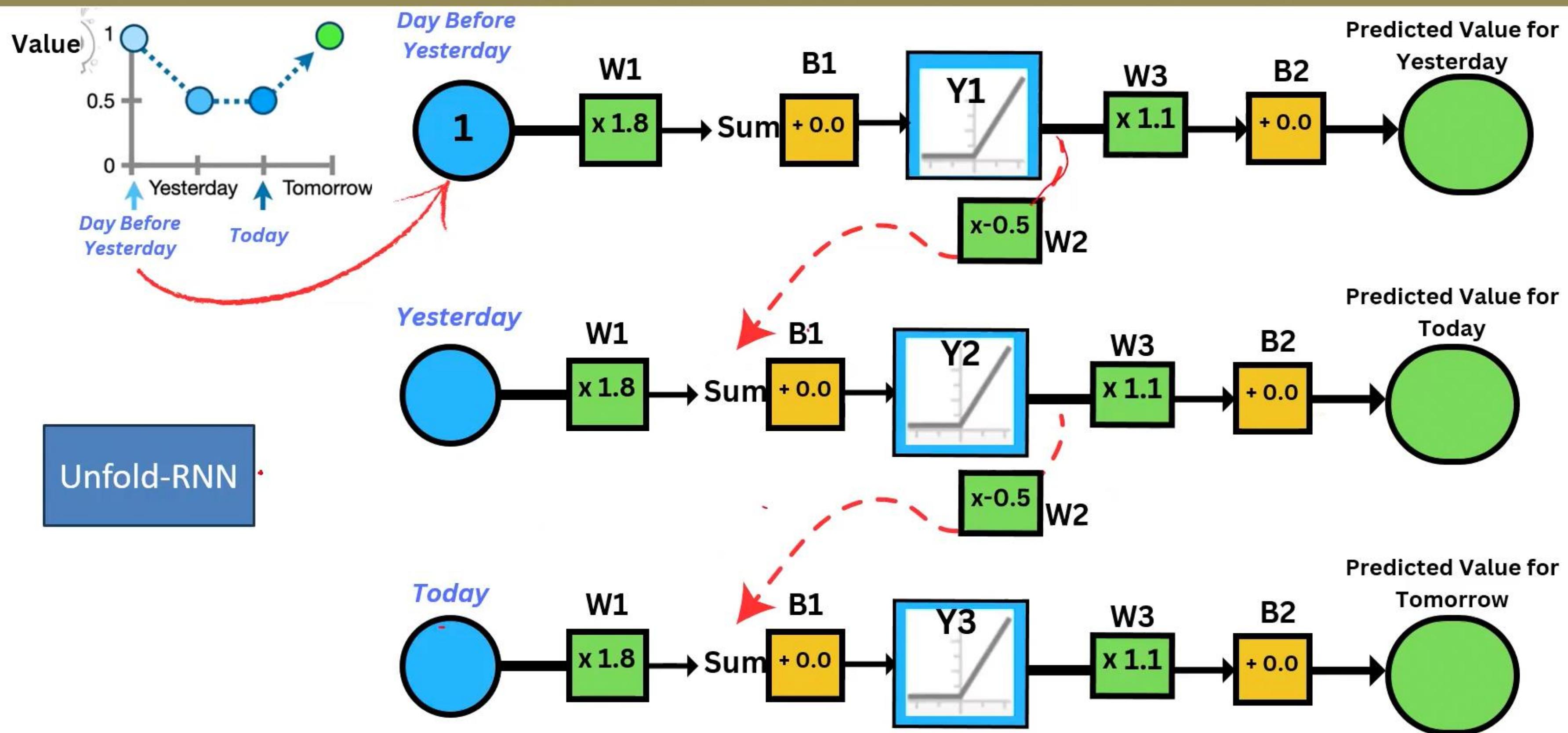


# Long Short-Term Memory- (LSTM)

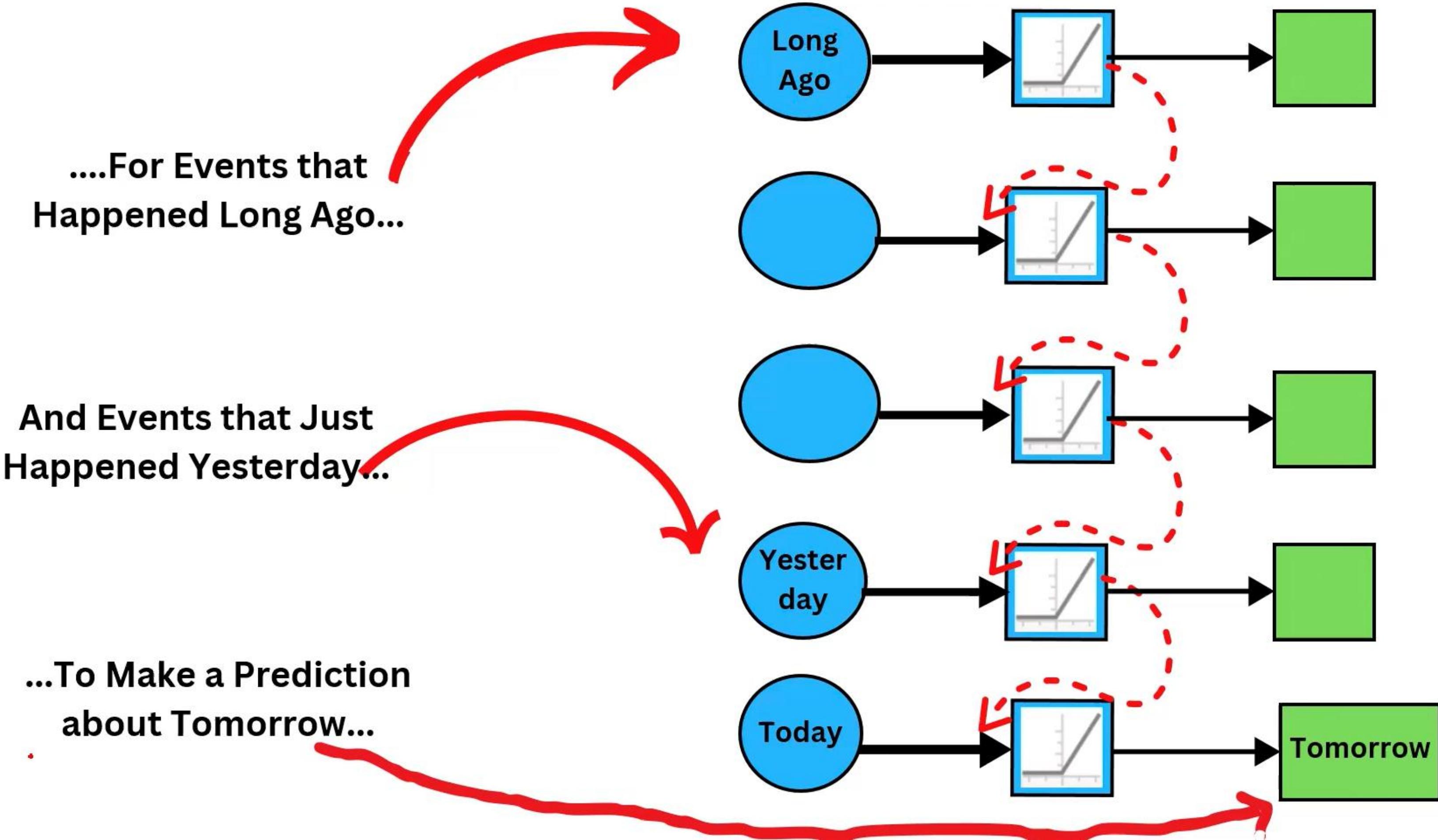




# Long Short-Term Memory- (LSTM)



# Long Short-Term Memory- (LSTM)

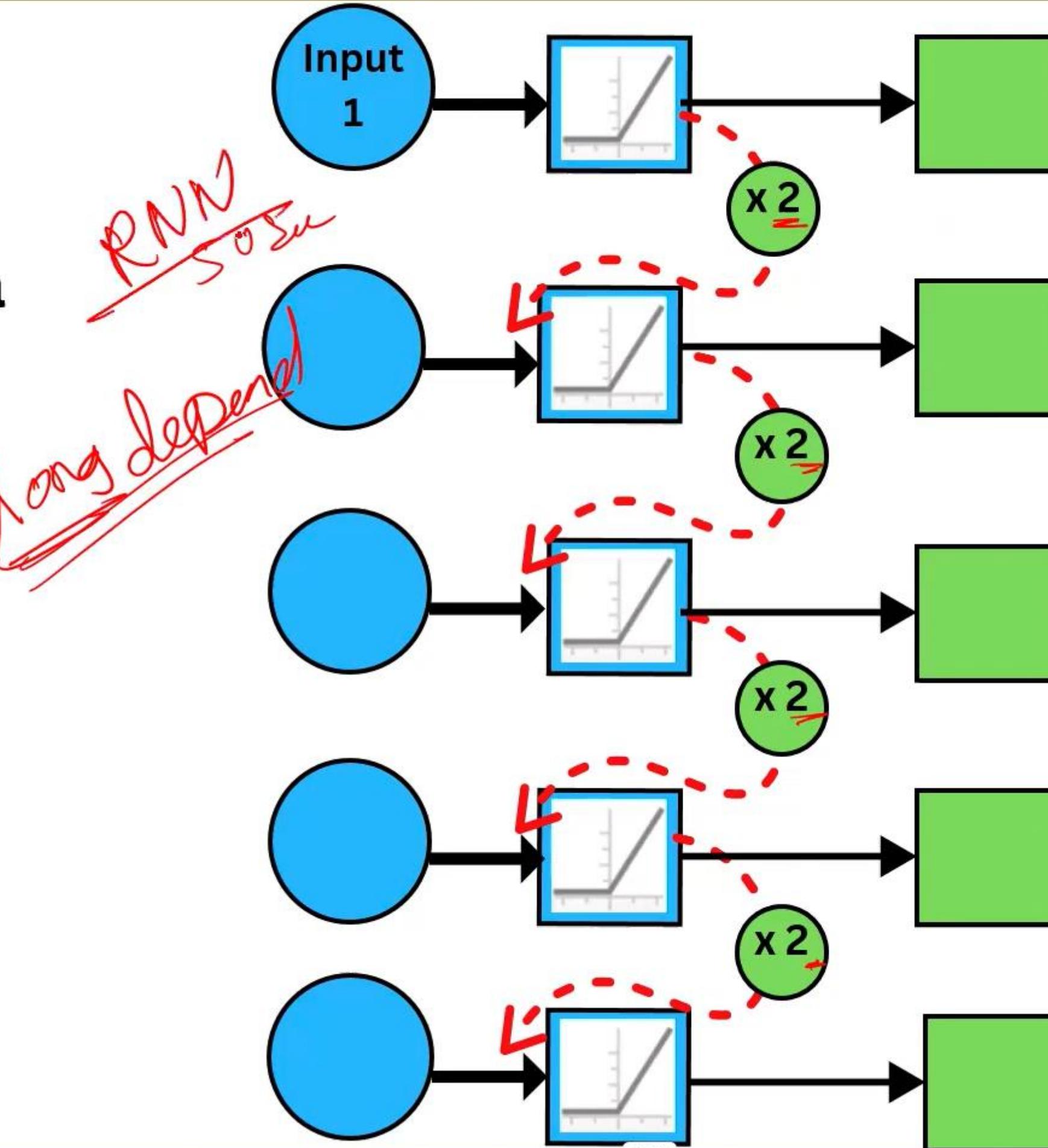


Exploding

2 to the 50<sup>th</sup> Power is a Huge Number

$$= \text{Input}_1 \times 2^{50}$$

$$= \text{Input}_1 \times \text{A Huge Number}$$



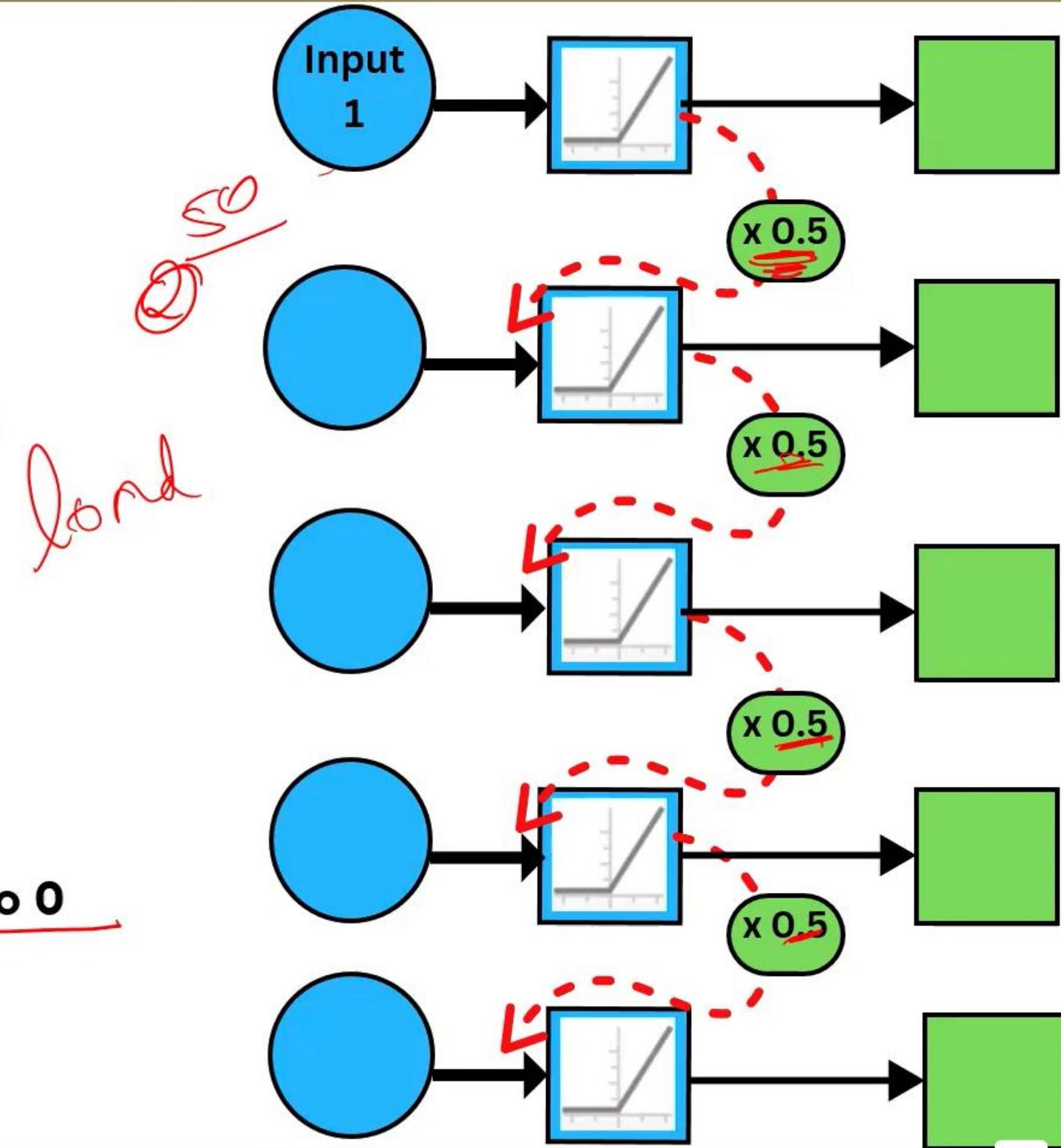
# Long Short-Term Memory- (LSTM)

Vanishing

This Number Super Close to 0  
would Cause the Gradient,  
Which We need for Gradient  
Descent

$$= \text{Input 1} \times 0.5$$

$$= \text{Input 1} \times \text{A Number Super Close to 0}$$



## ◆ 1. ReLU (Rectified Linear Unit) — *Most Popular*

$$f(x) = \max(0, x)$$

- Used in CNN, Deep Neural Networks
- Fast & simple
- Problem: “Dying ReLU” for negative values
- Use case: Hidden layers in almost all modern models

## ◆ 2. Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

- Output between 0 and 1
- Good for binary classification
- Vanishing gradient problem
- Use case: Output layer for binary classification

## ◆ 3. Tanh (Hyperbolic Tangent)

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- Output between -1 and +1
- Zero-centered
- Still suffers from vanishing gradients
- Use case: RNN, LSTM, NLP models

## tanh – Positive Value

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(2) = \frac{e^2 - e^{-2}}{e^2 + e^{-2}}$$

**=0.96**

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

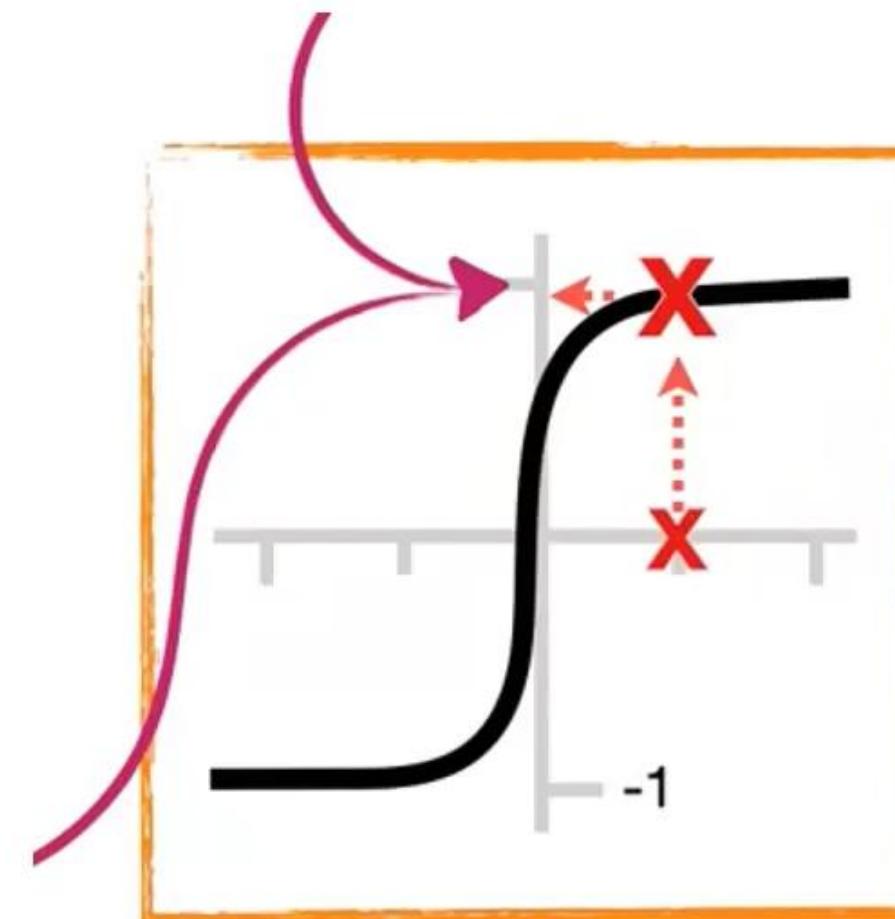
$$f(-5) = \frac{e^{-5} - e^5}{e^{-5} + e^5}$$

**=-1**

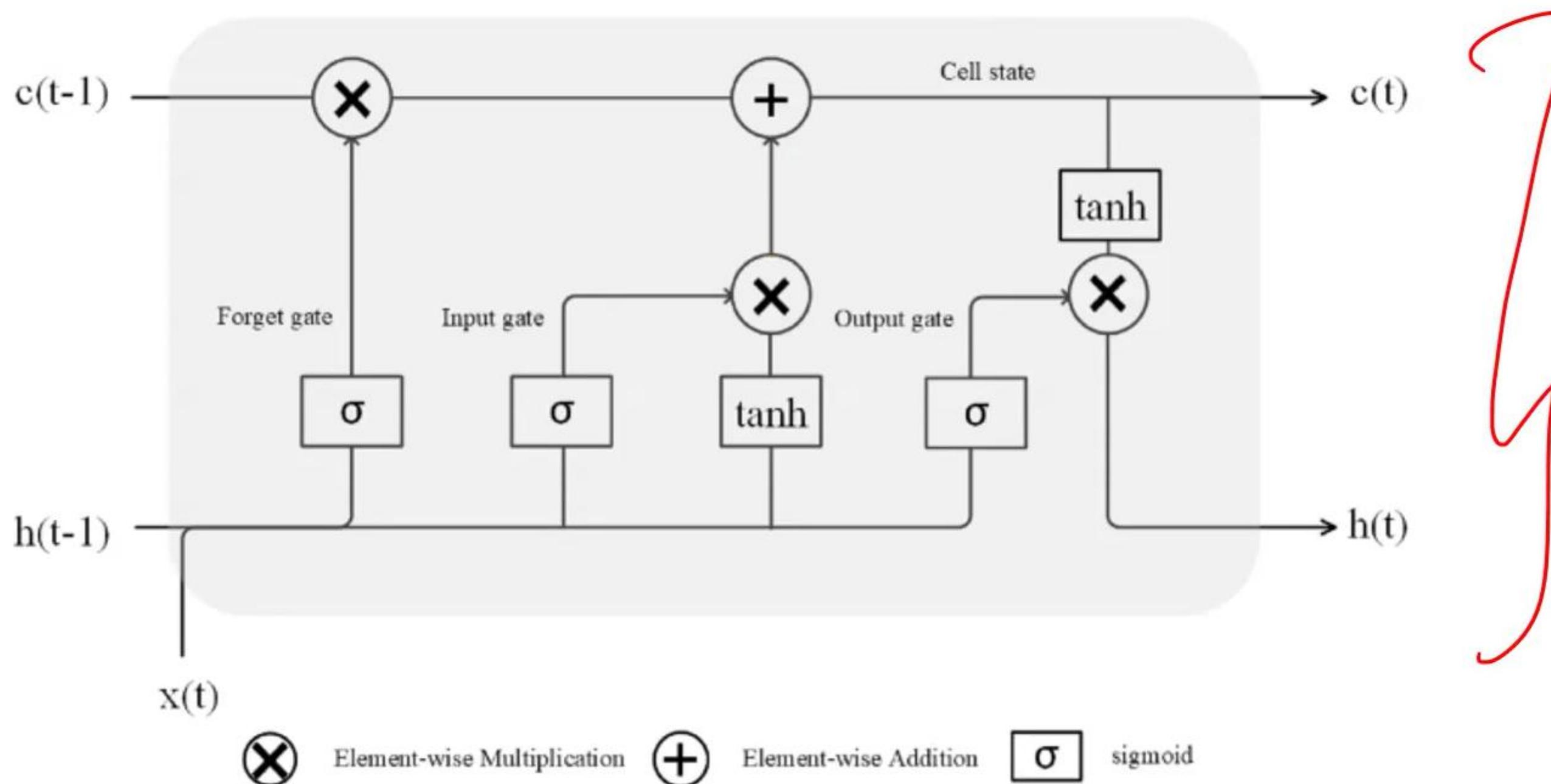
- 3. Tanh (Hyperbolic Tangent)

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

- ✓ Output between -1 and +1
- ✓ Zero-centered
- ✗ Still suffers from vanishing gradients
- ✳ Use case: RNN, LSTM, NLP models



## Long Short-Term Memory- (LSTM)-Case -1-Positive Value

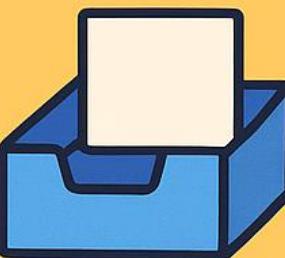


# LSTM MEMORY BOX



## 1. Forget Old Stuff

“Throw away things I don’t need anymore!”



## 2. Add New Stuff

“Learn something new!”



## 3. Mix Old + New

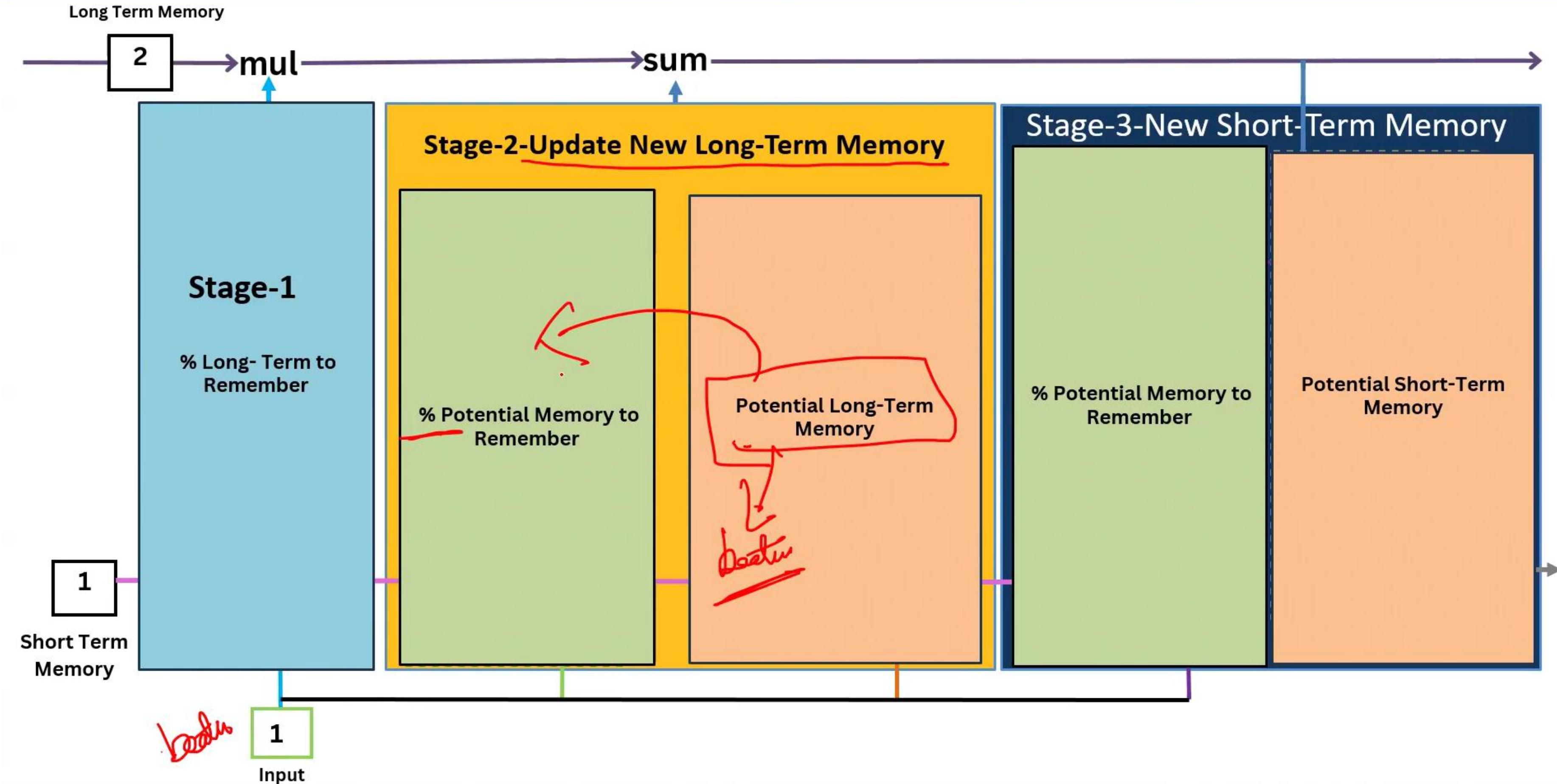
“Update my brain!”



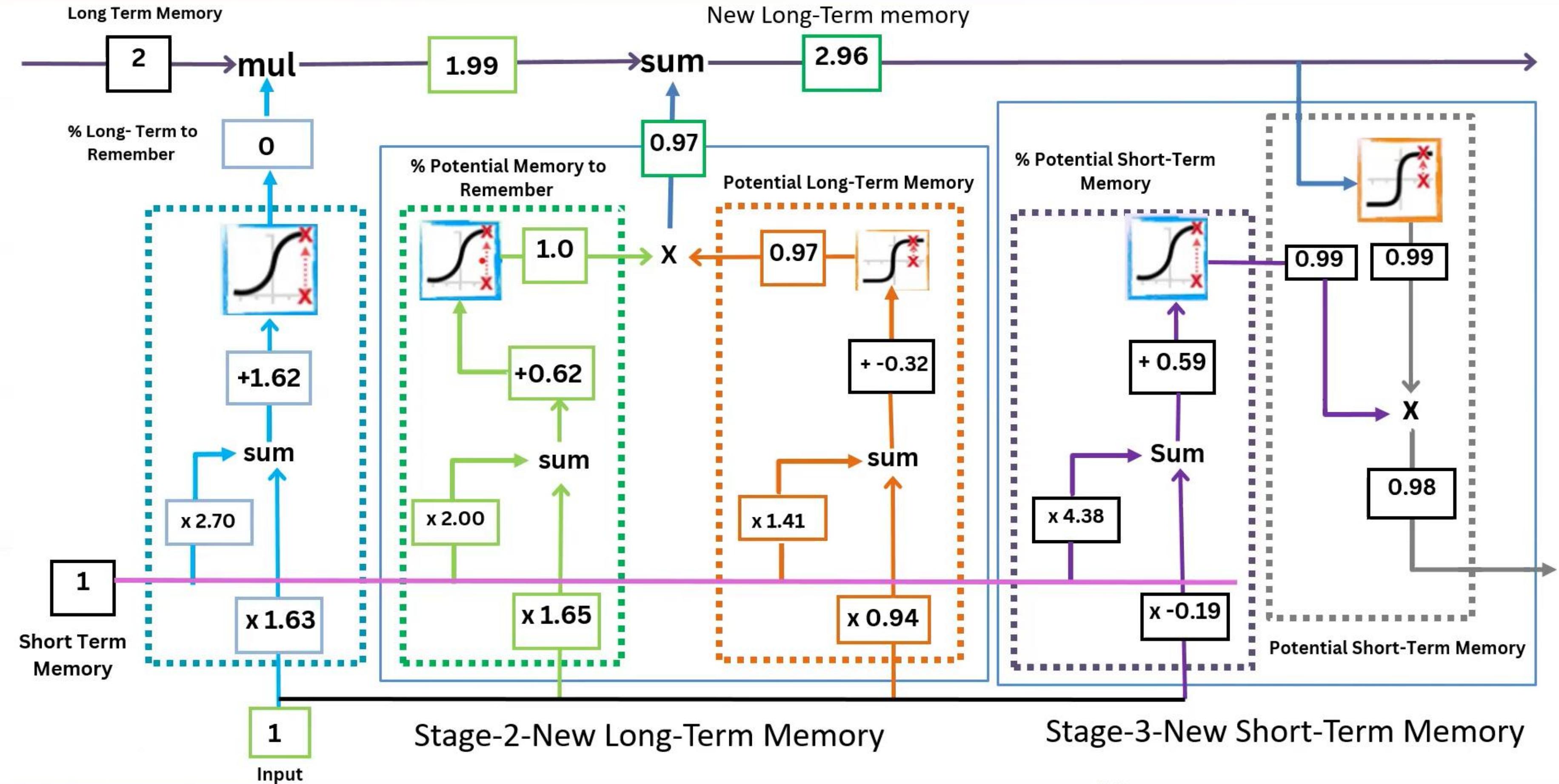
## 4. Output Answer

“Say something useful now!”

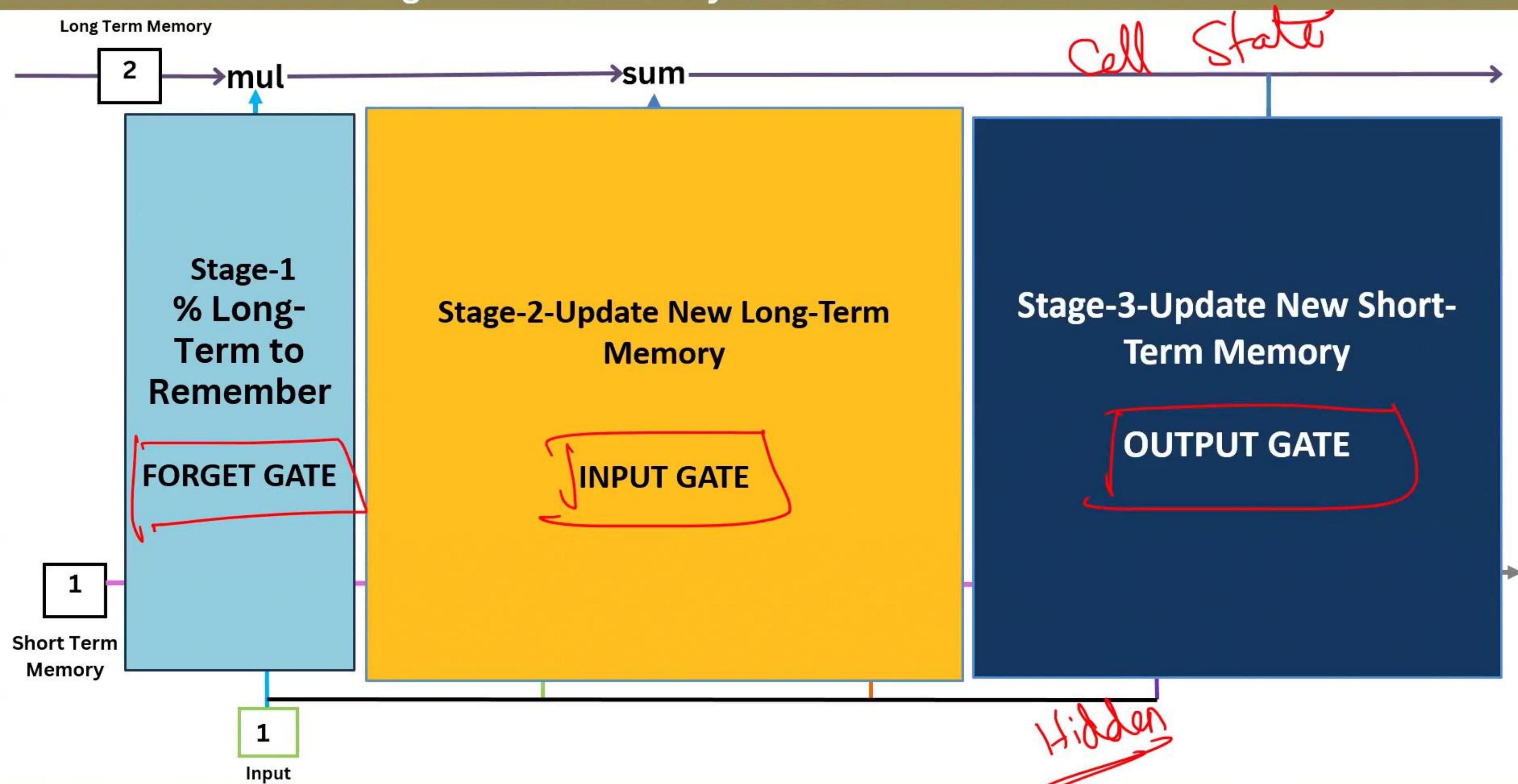
# Long Short-Term Memory- (LSTM)-Case -1-Positive Value



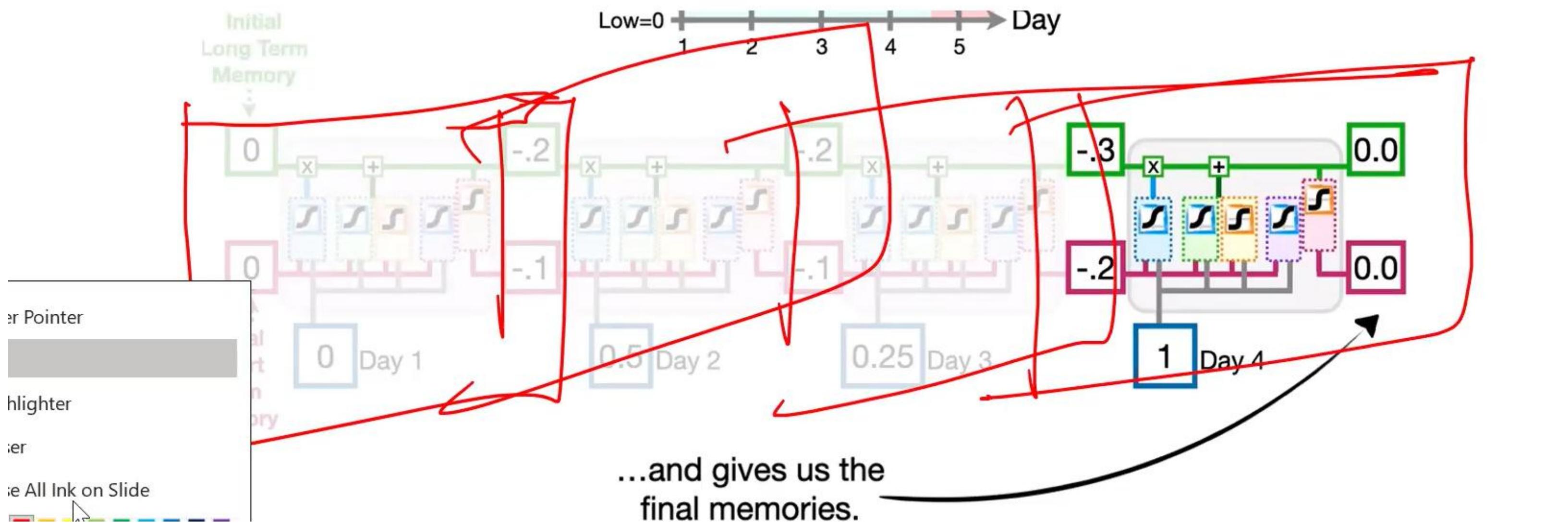
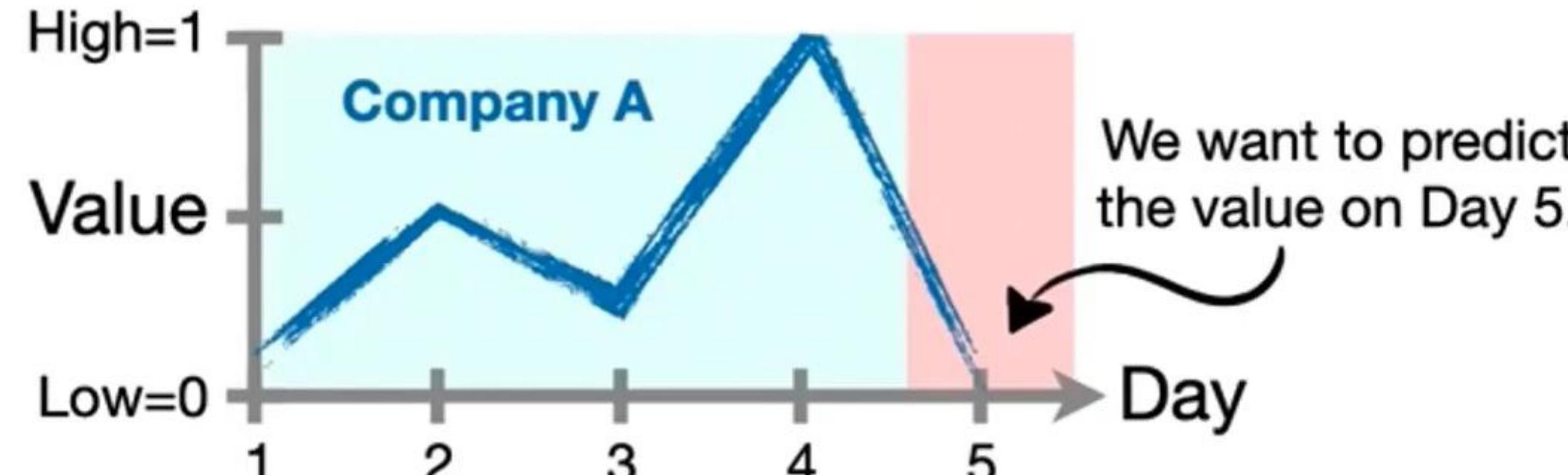
# Long Short-Term Memory- (LSTM)-Case -1-Positive Value



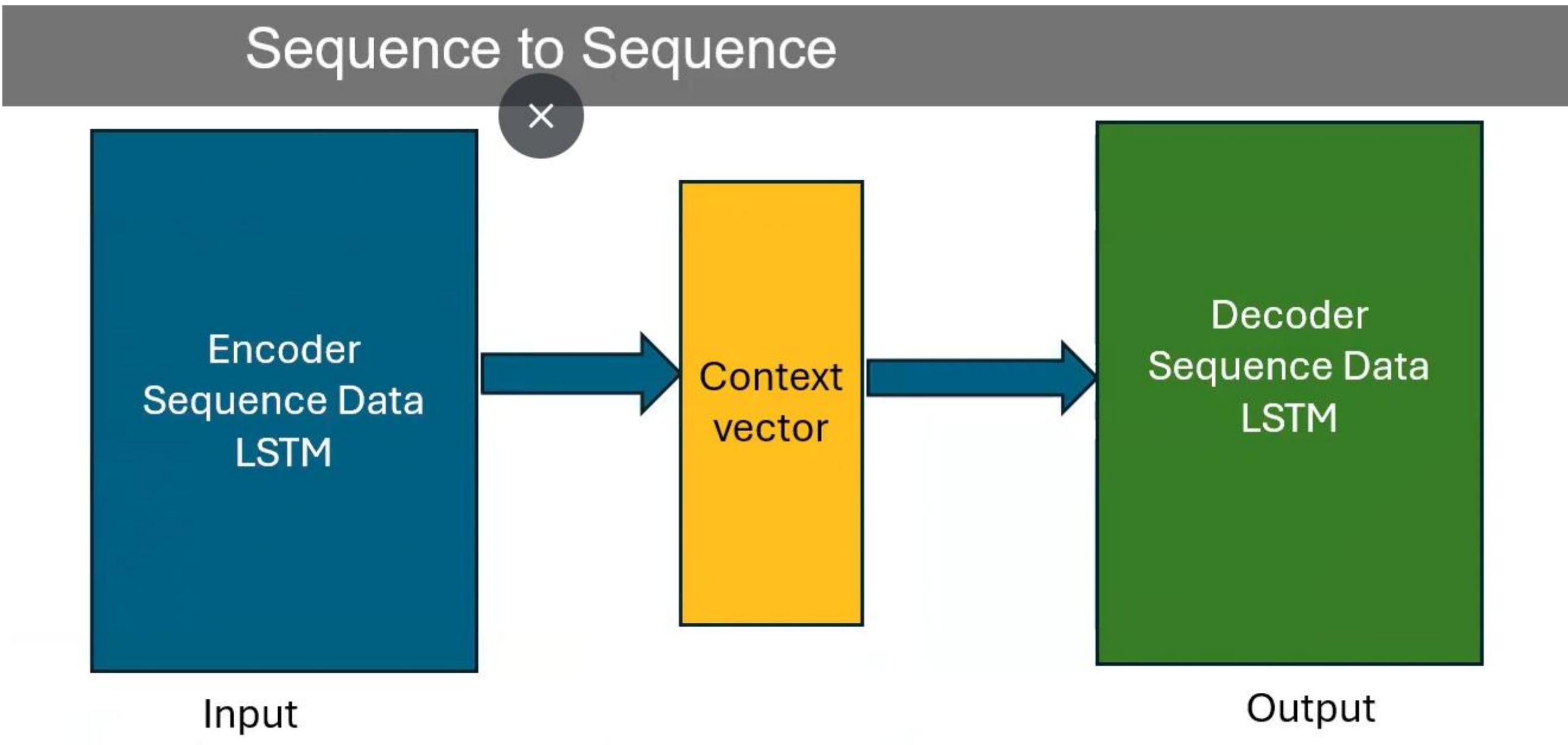
## Long Short-Term Memory- (LSTM)-Case -1-Positive Value



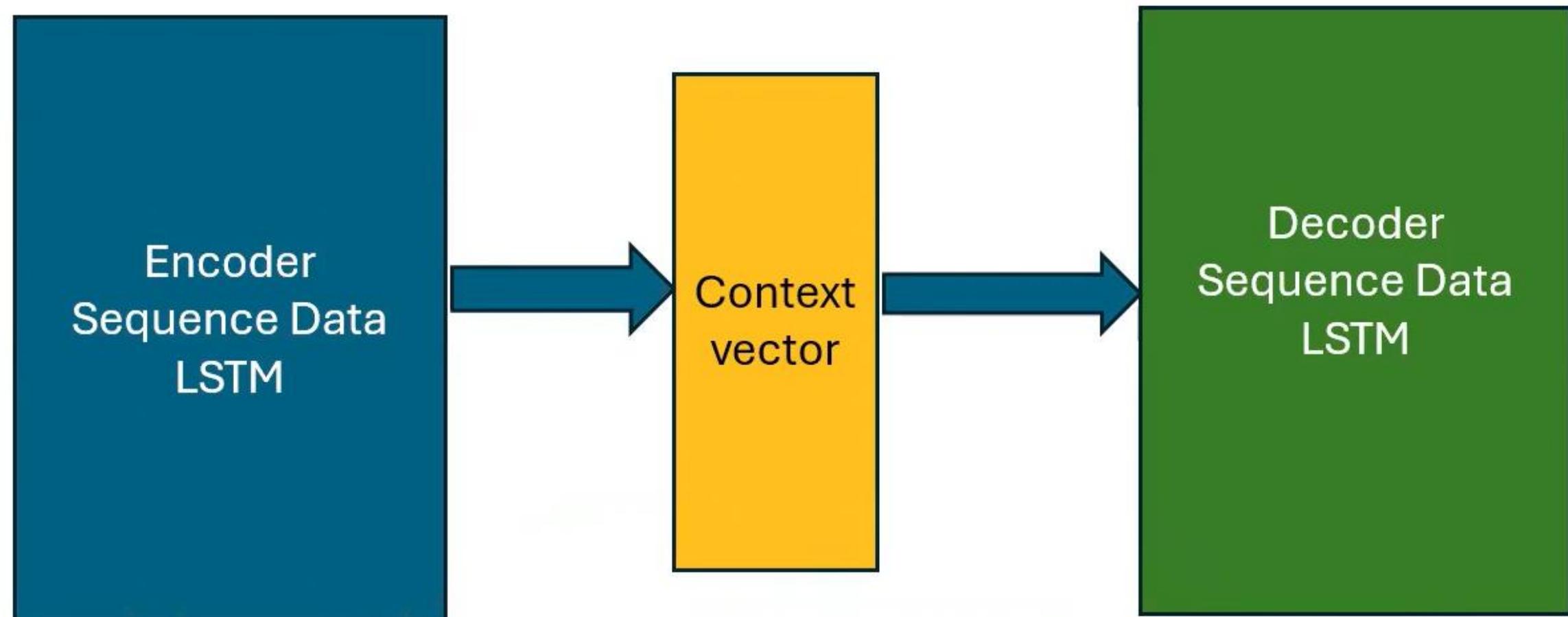
## Long Short-Term Memory- (LSTM)



## Sequence to Sequence



# Sequence to Sequence



Input

Output

**Machine Translation**

English

Tamil

**Text Summarization**

Paragraphs

Summarization

**Question Answering**

Questions

Answer

# Tamil Cinema Famous Dialog Completion

input

நா ஒரு தடவை சொன்னா

ஒரு வாட்டி முடிவு பண்ணிட்டா

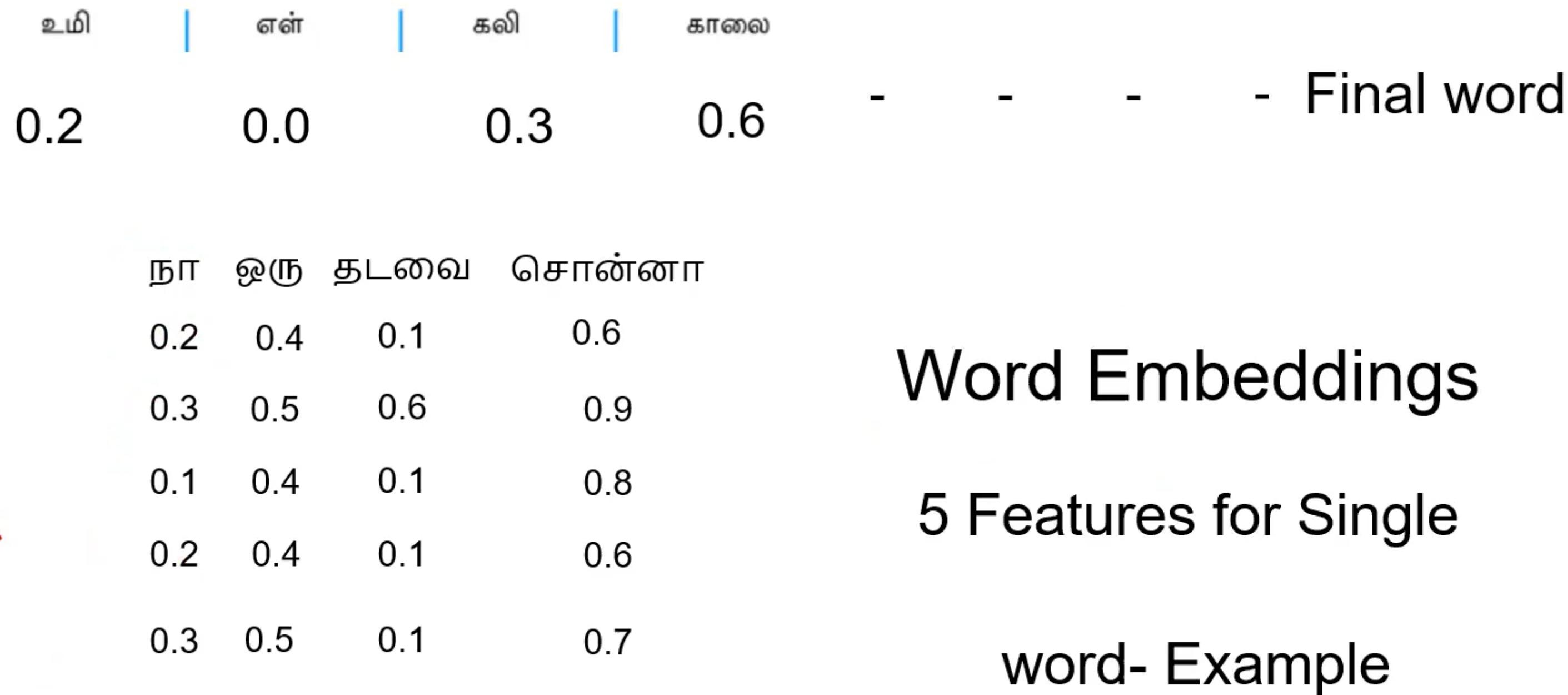
Output

நாறு தடவா சொன்ன மாறி

என் பேச்சை நானே

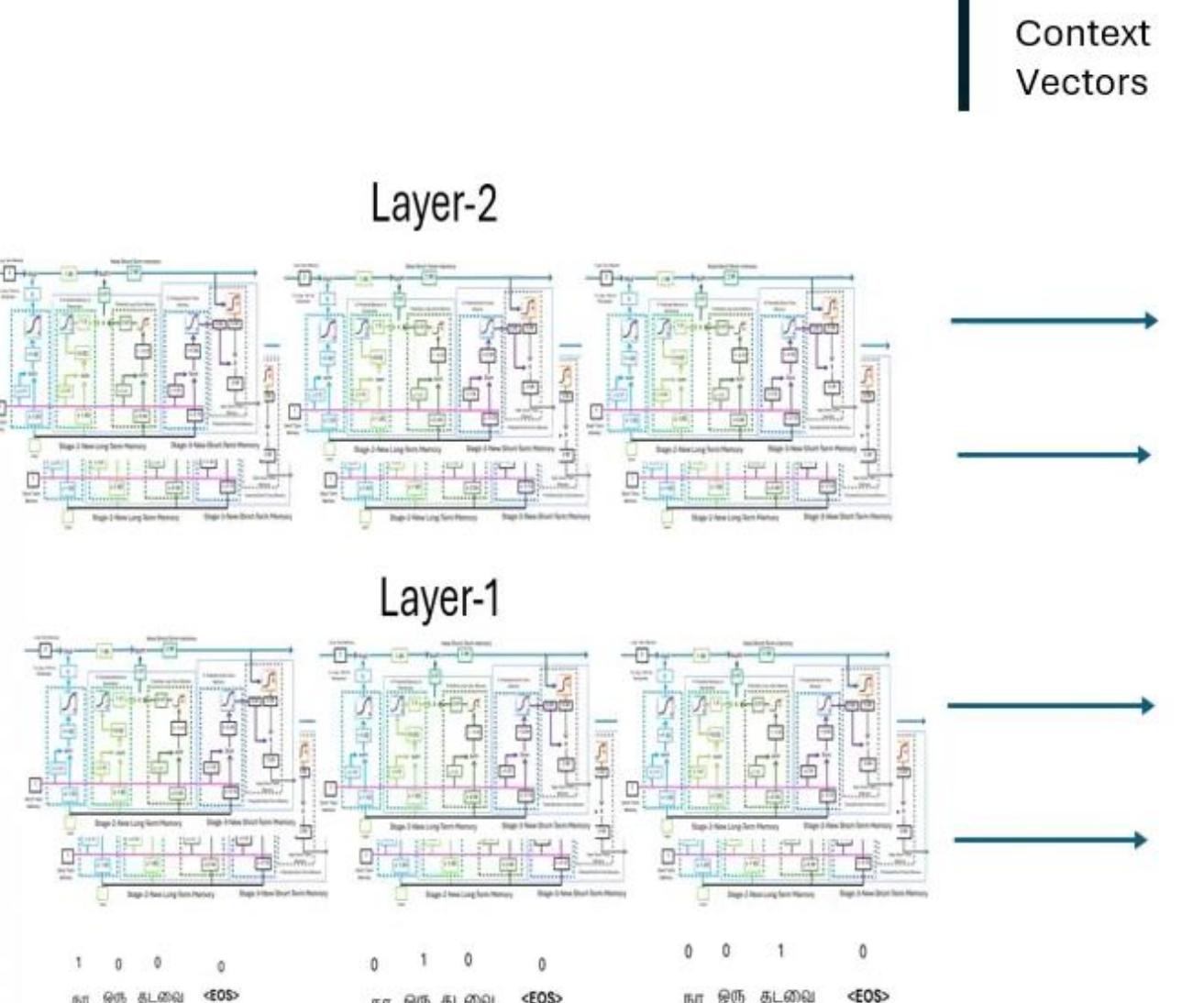
கேட்கமாட்டேன்

# Sequence to Sequence

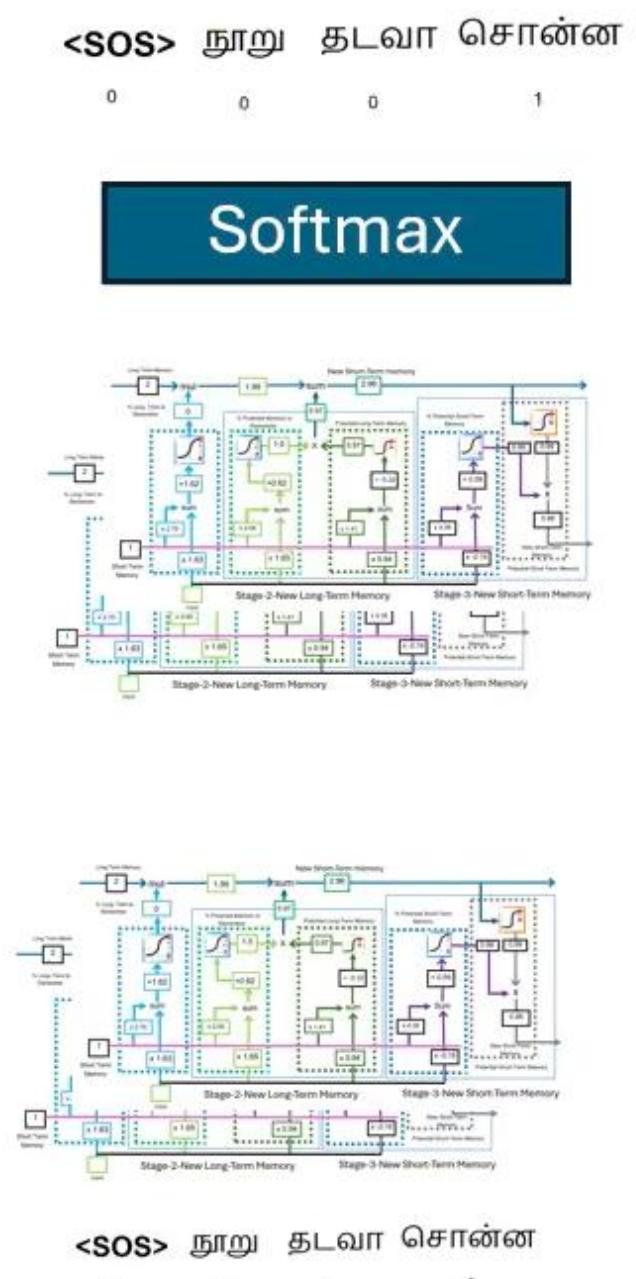
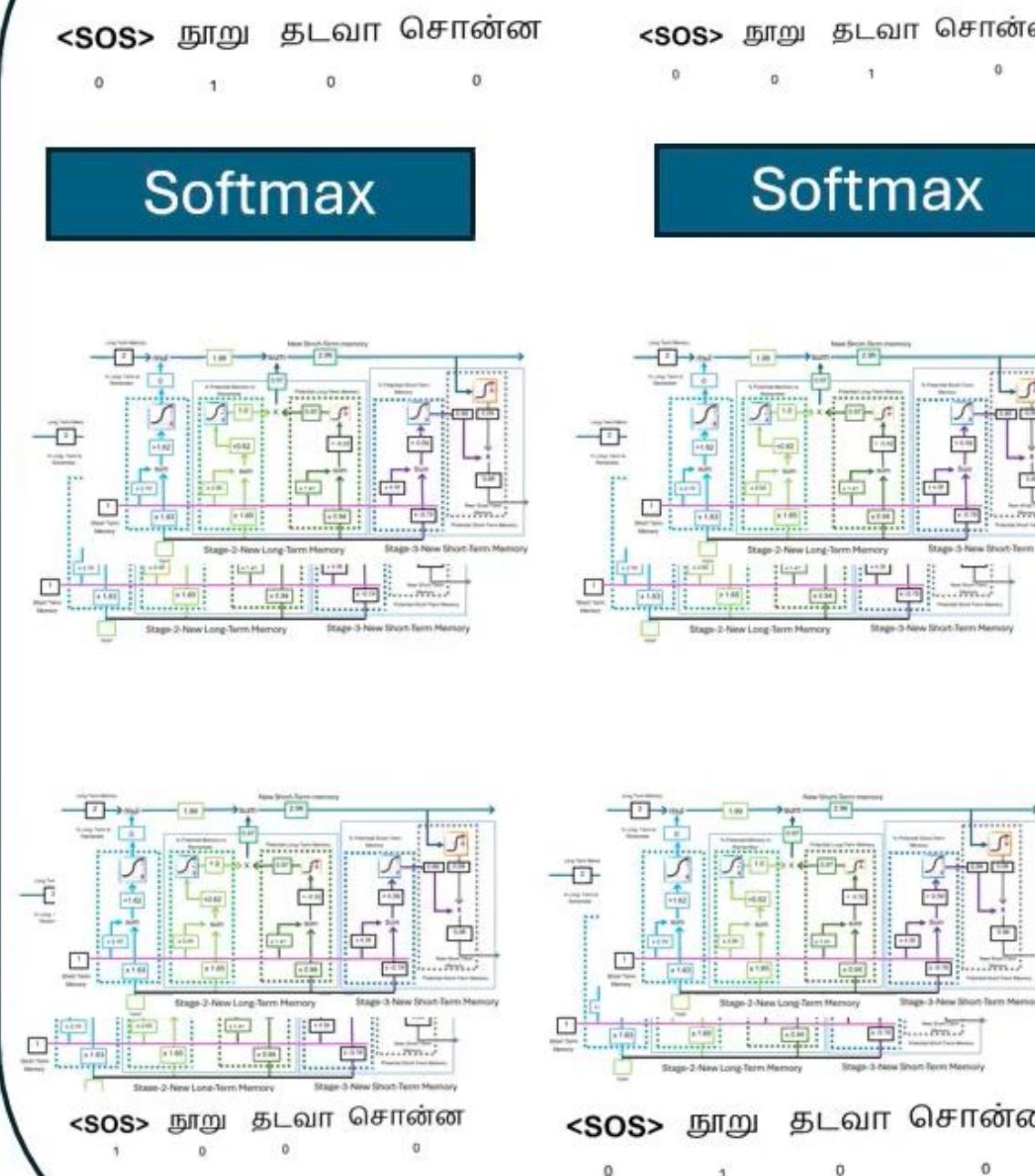


# Sequence to Sequence

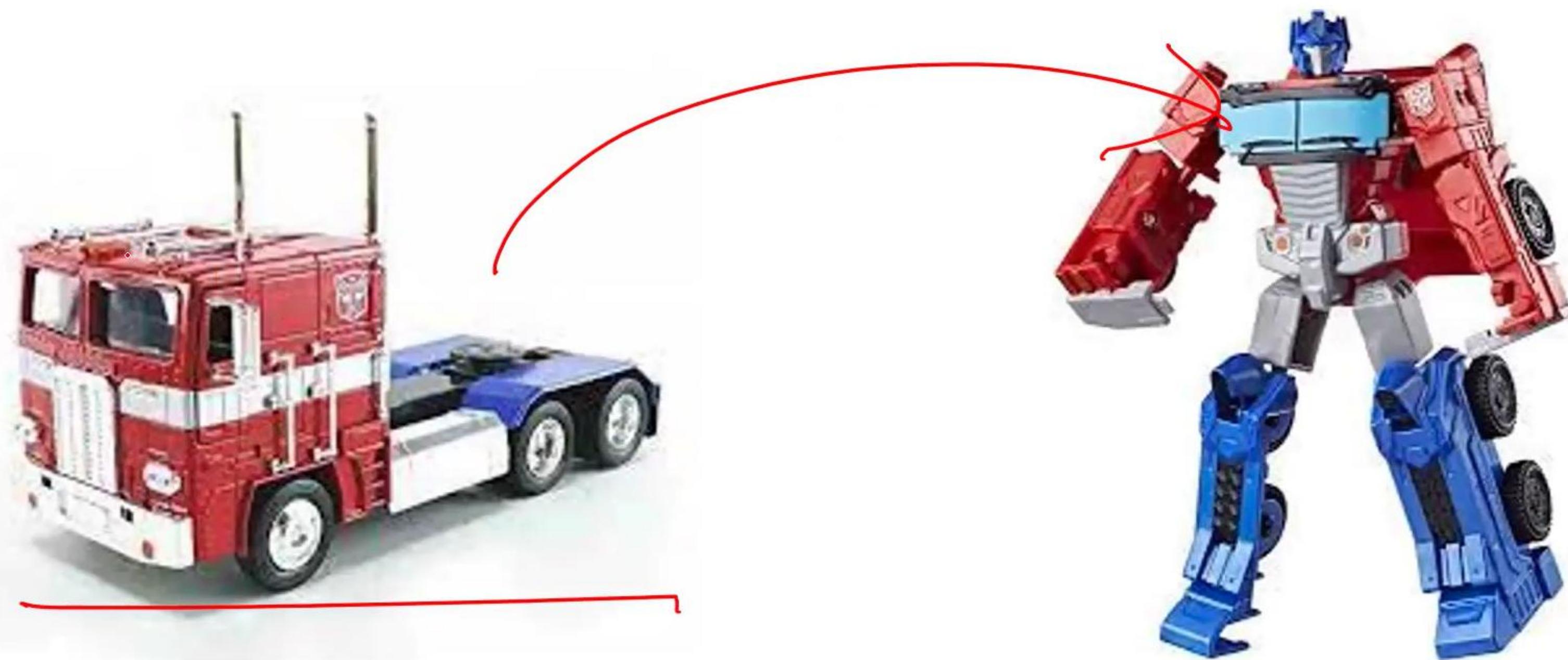
## Encoder



## Decoder



# Why Transformer?

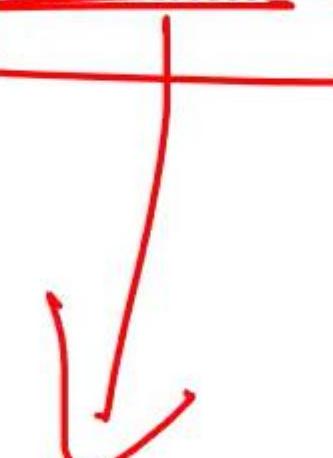




You

Complete the sentence " You are beautiful...."

IP



C.

→ (6)



ChatGPT

"You are beautiful just the way you are, inside and out."





You

write python to handle large number rows



ChatGPT

Handling large numbers of rows in Python often involves using efficient libraries like pandas or Dask for data manipulation and processing. Here's a basic example using pandas:

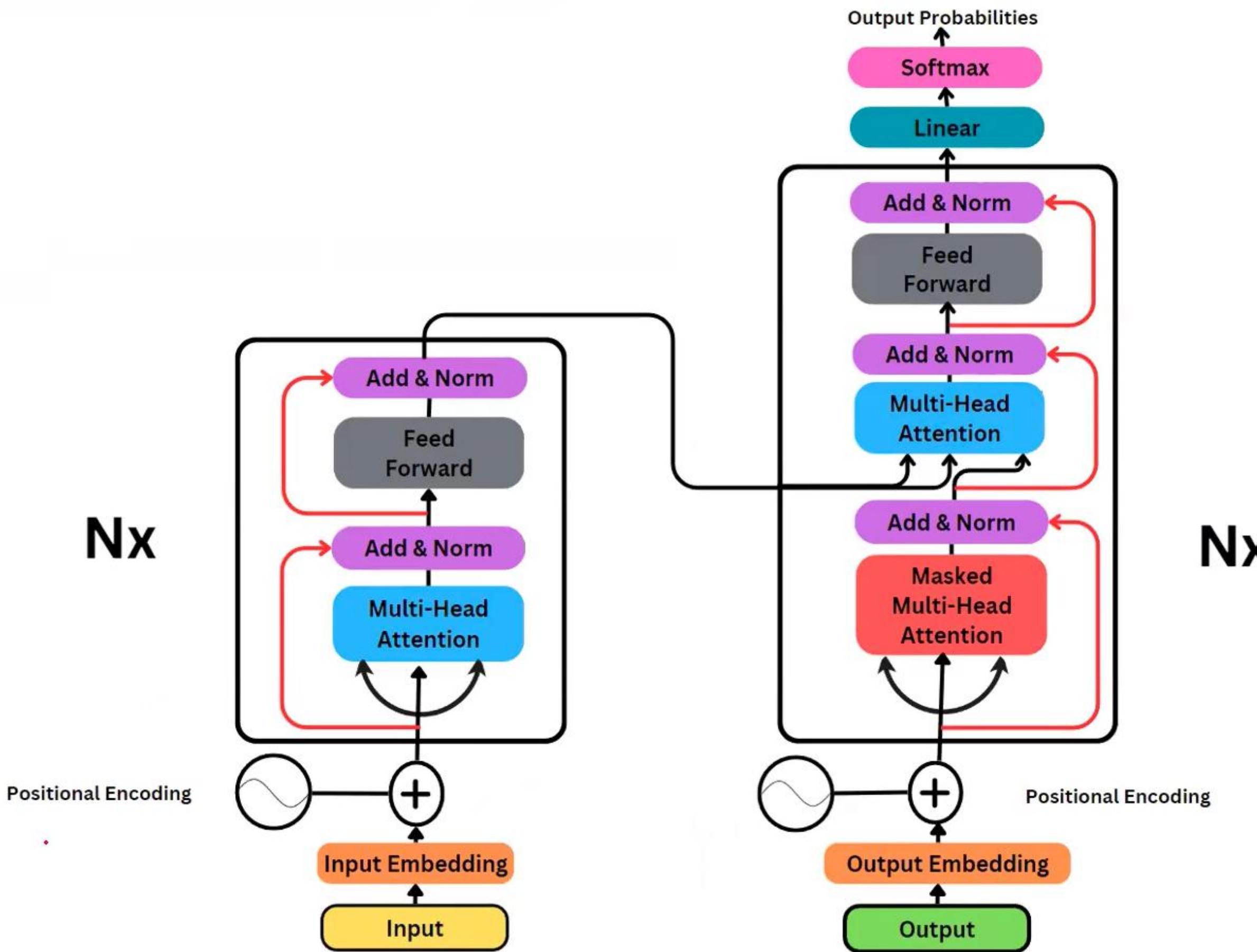
python

Copy code

```
import pandas as pd

# Assuming you have a large CSV file named 'data.csv'
# Load the data into a pandas DataFrame in chunks
chunk_size = 100000 # Adjust the chunk size as needed
data_chunks = pd.read_csv('data.csv', chunksize=chunk_size)
```

# Why Transformer?



# Why Transformer?

Supervised  
Learning

Unsupervised  
Learning

## Why Transformer?

# Tamil Cinema Famous Dialog Completion

**input**

நா ஒரு தடவை சொன்னா

ஒரு வாட்டி முடிவு பண்ணிட்டா

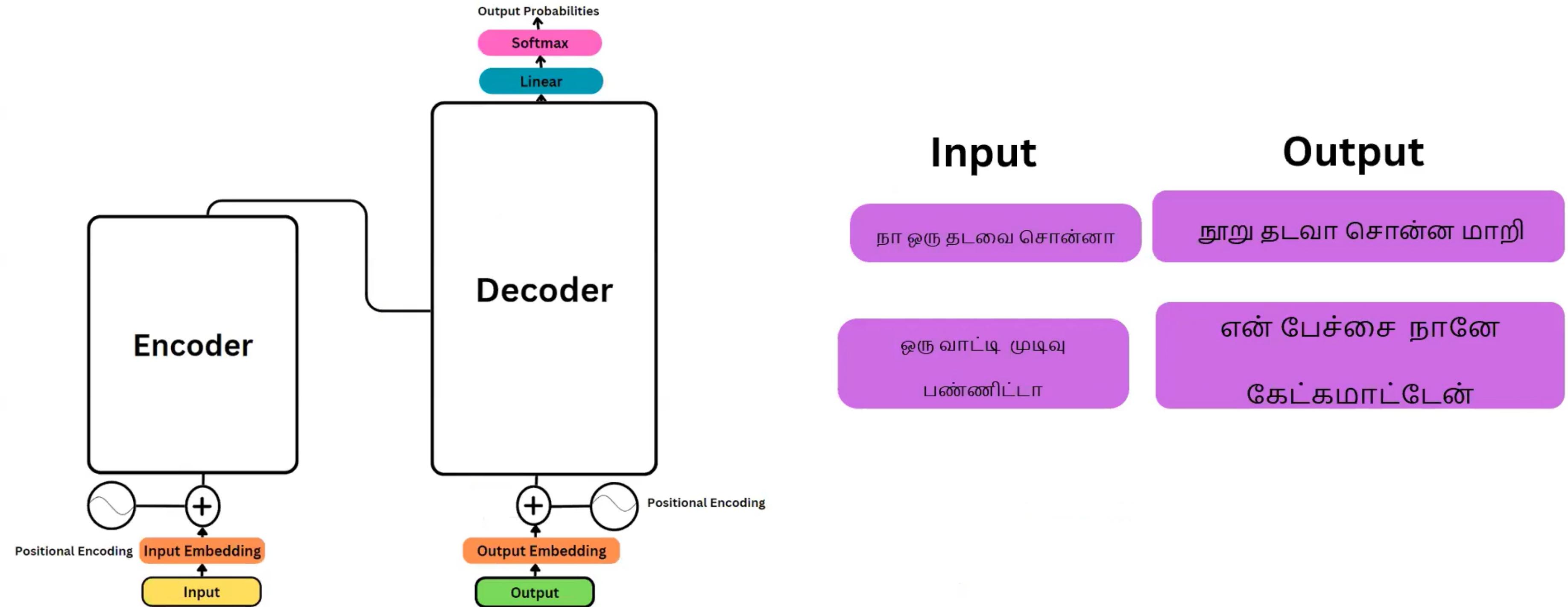
**Output**

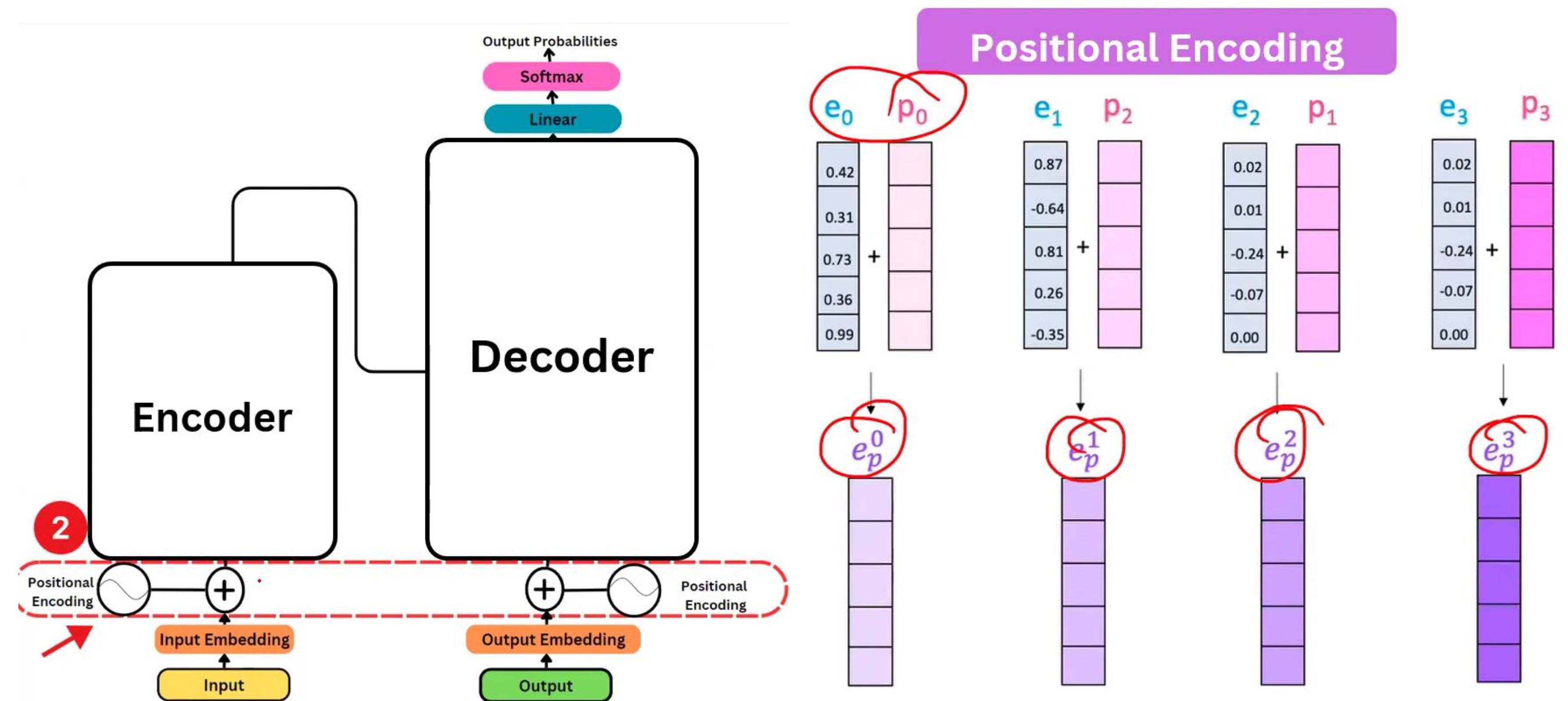
நூறு தடவா சொன்ன மாறி

என் பேச்சை நானே

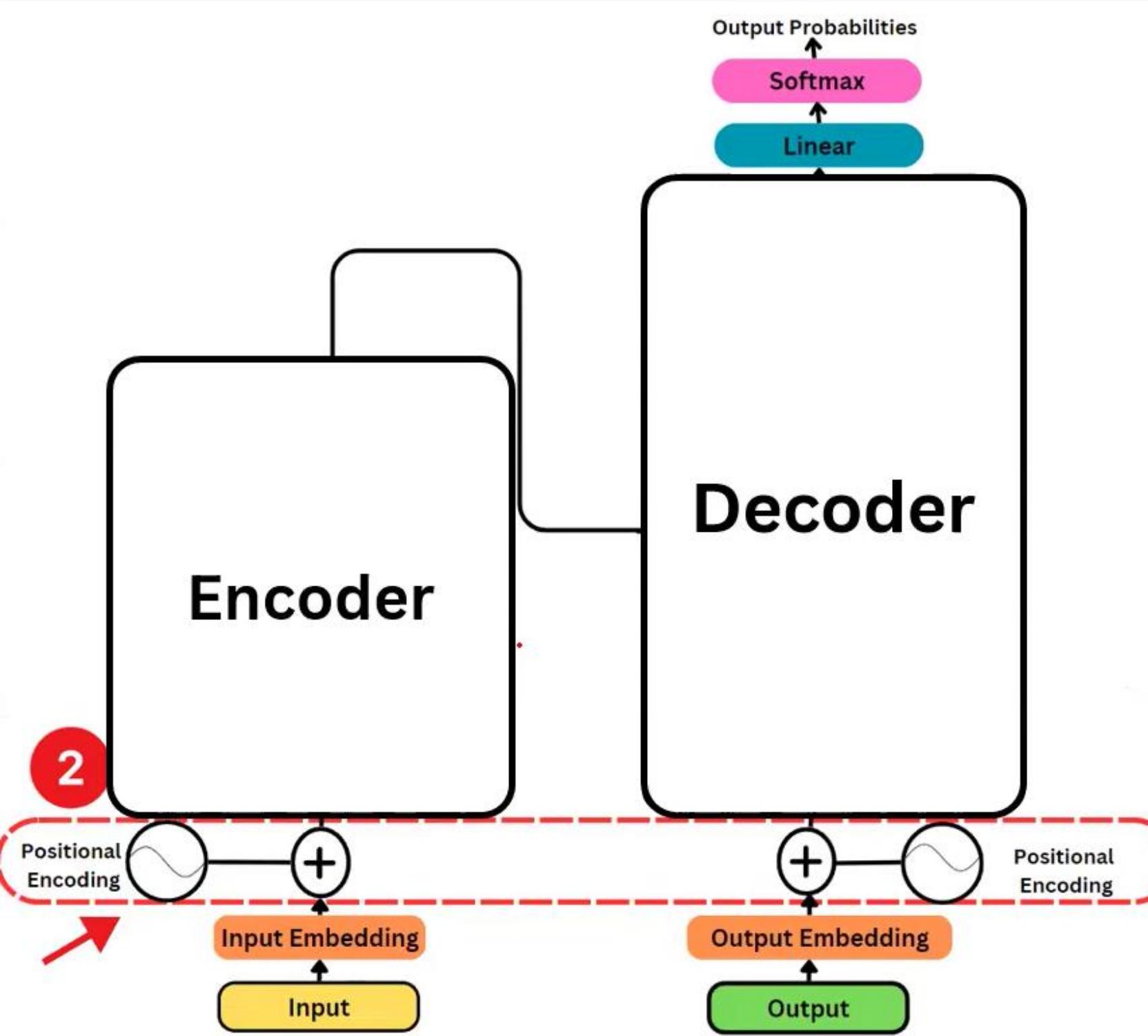
கேட்கமாட்டேன்

# Why Transformer?





## Transformer-Positional Encoding



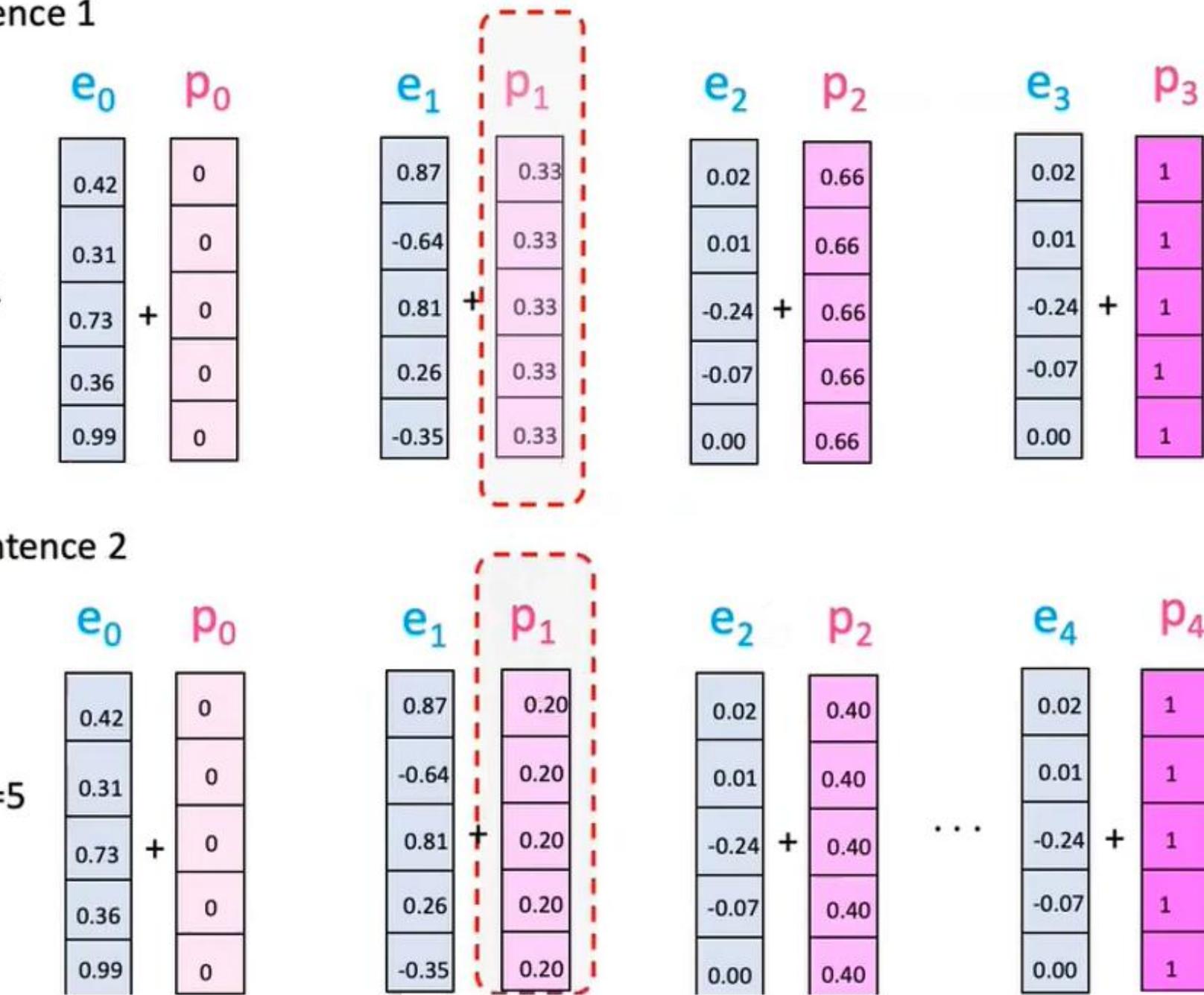
## Positional Encoding -Option 2

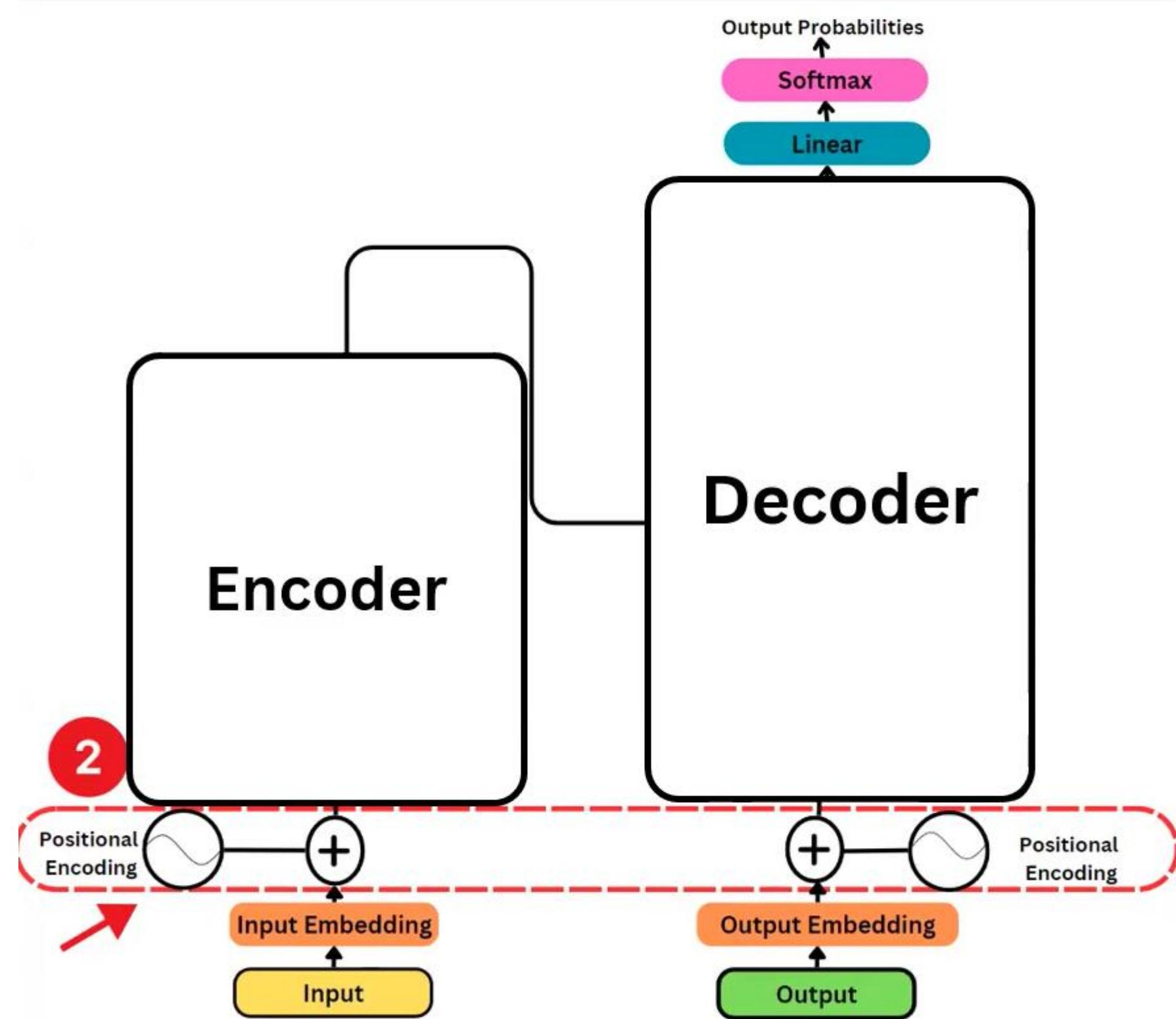
## Sentence:

N=4

### Sentence 2

N





### Positional Encoding -Option-3

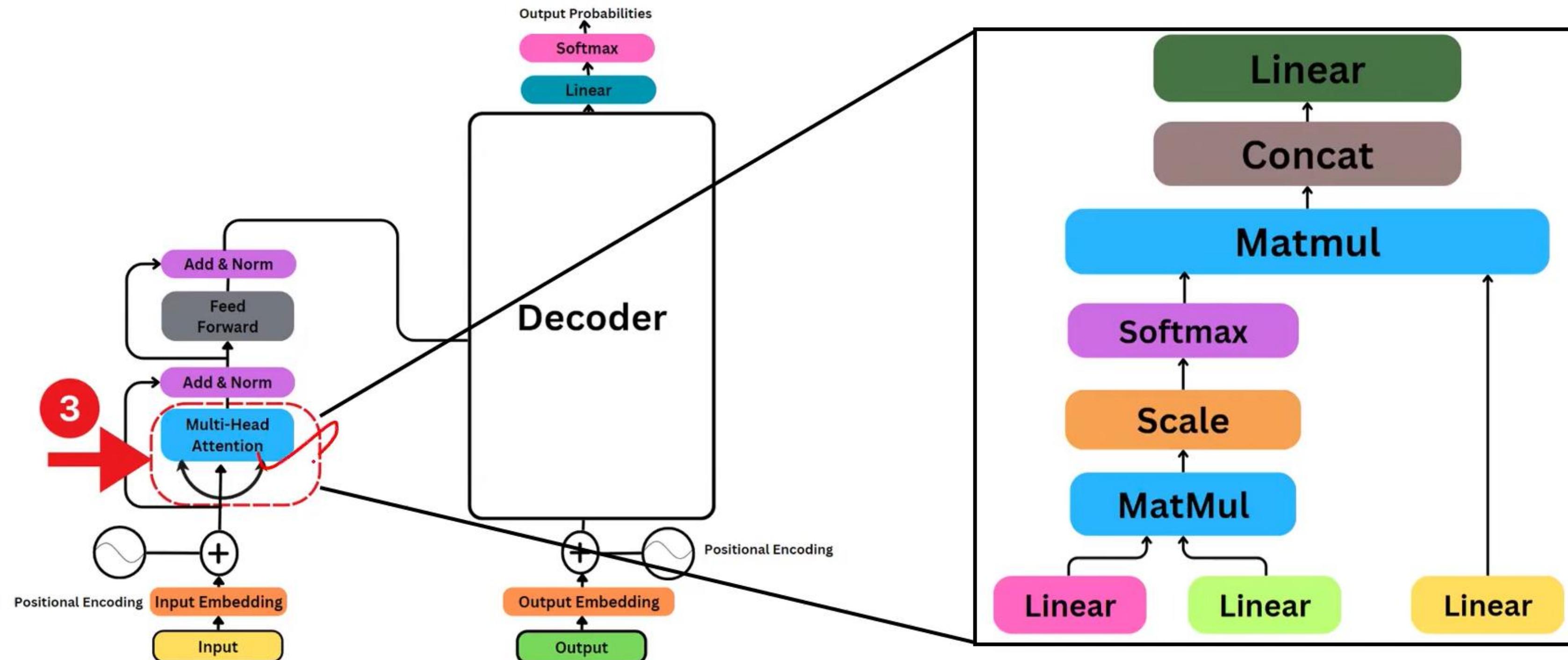
$$PE_{(pos,2i)} = \frac{\sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)}{?}$$

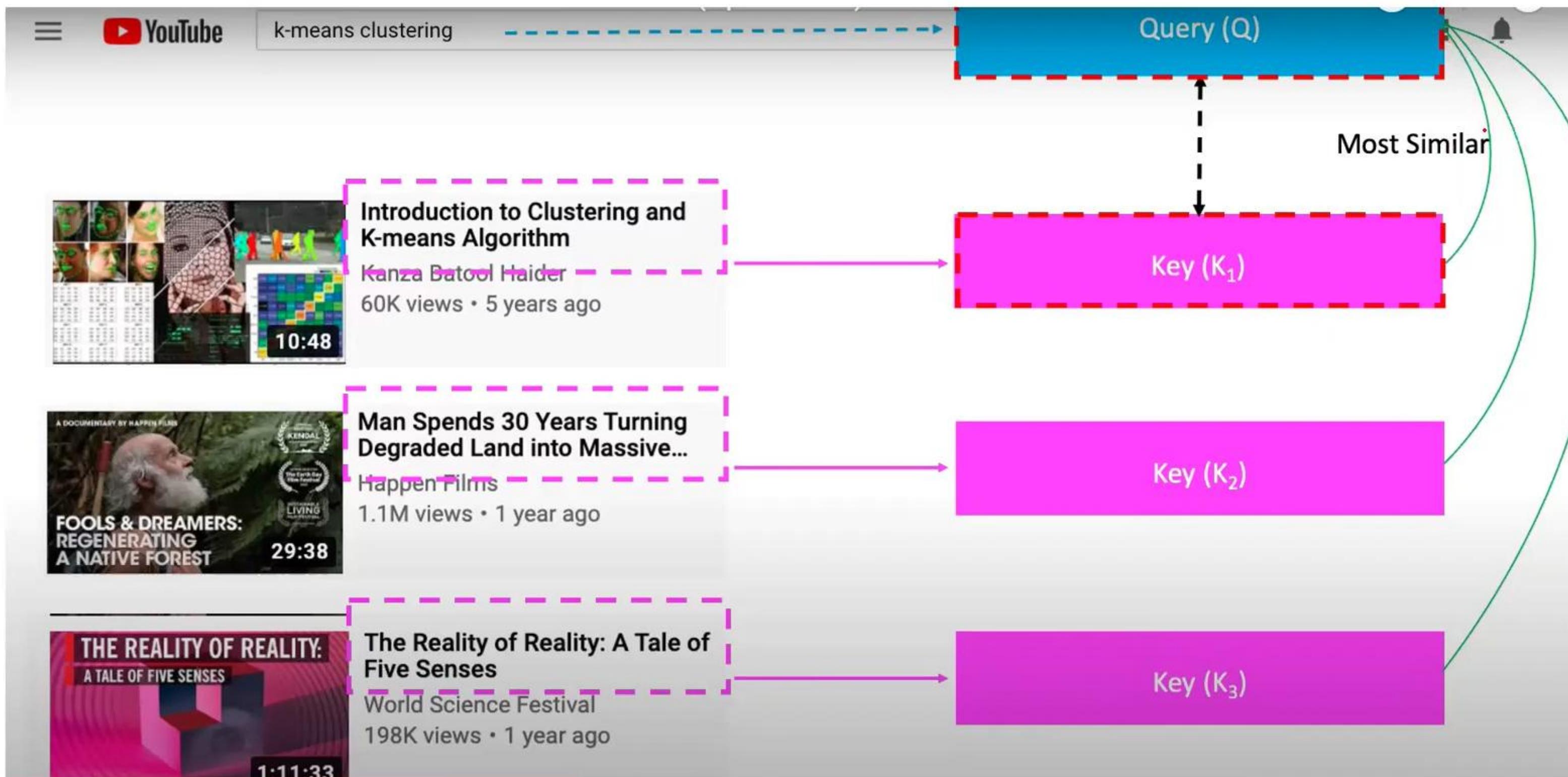
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right) ?$$

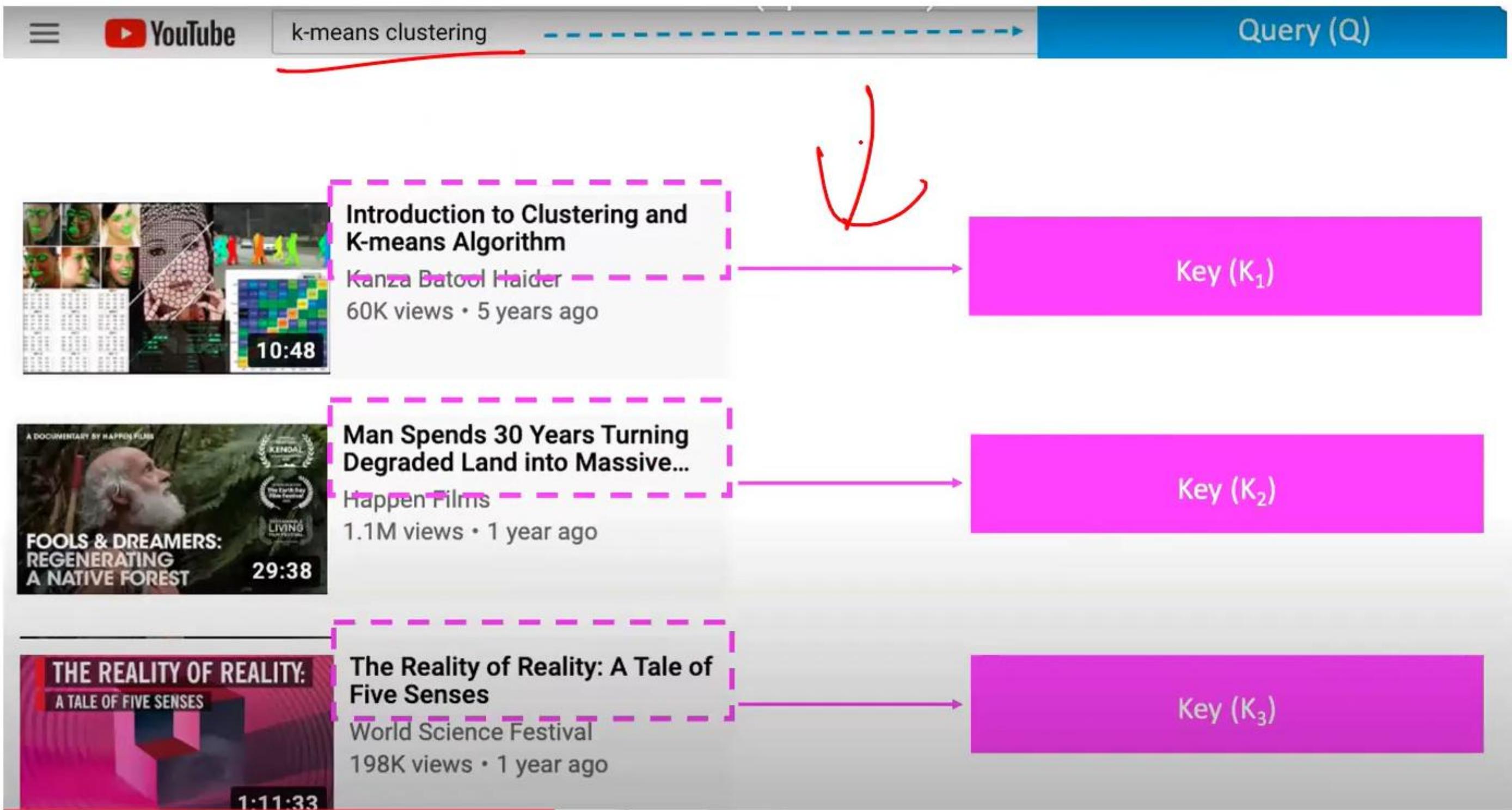
odd

?

even

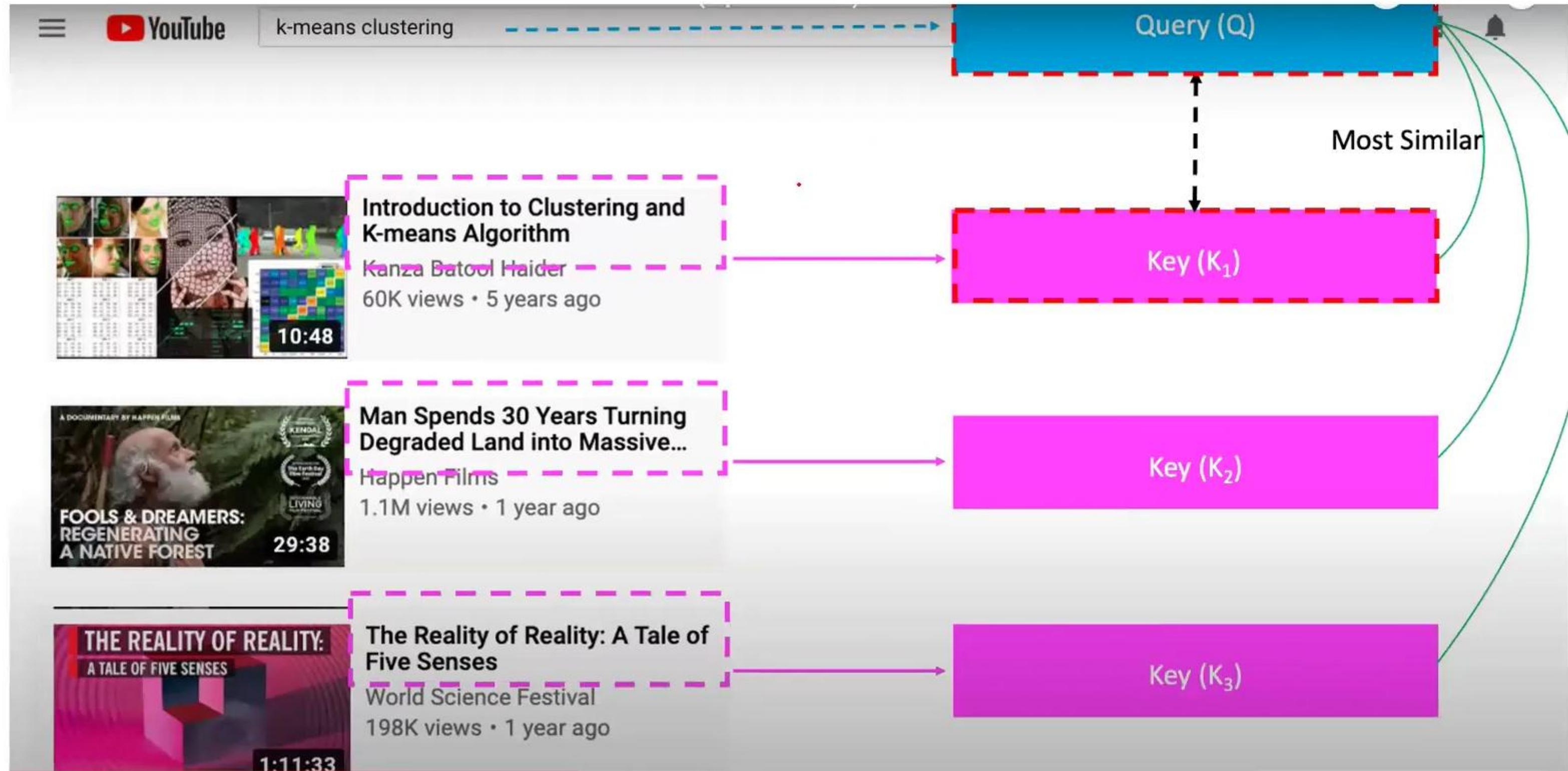




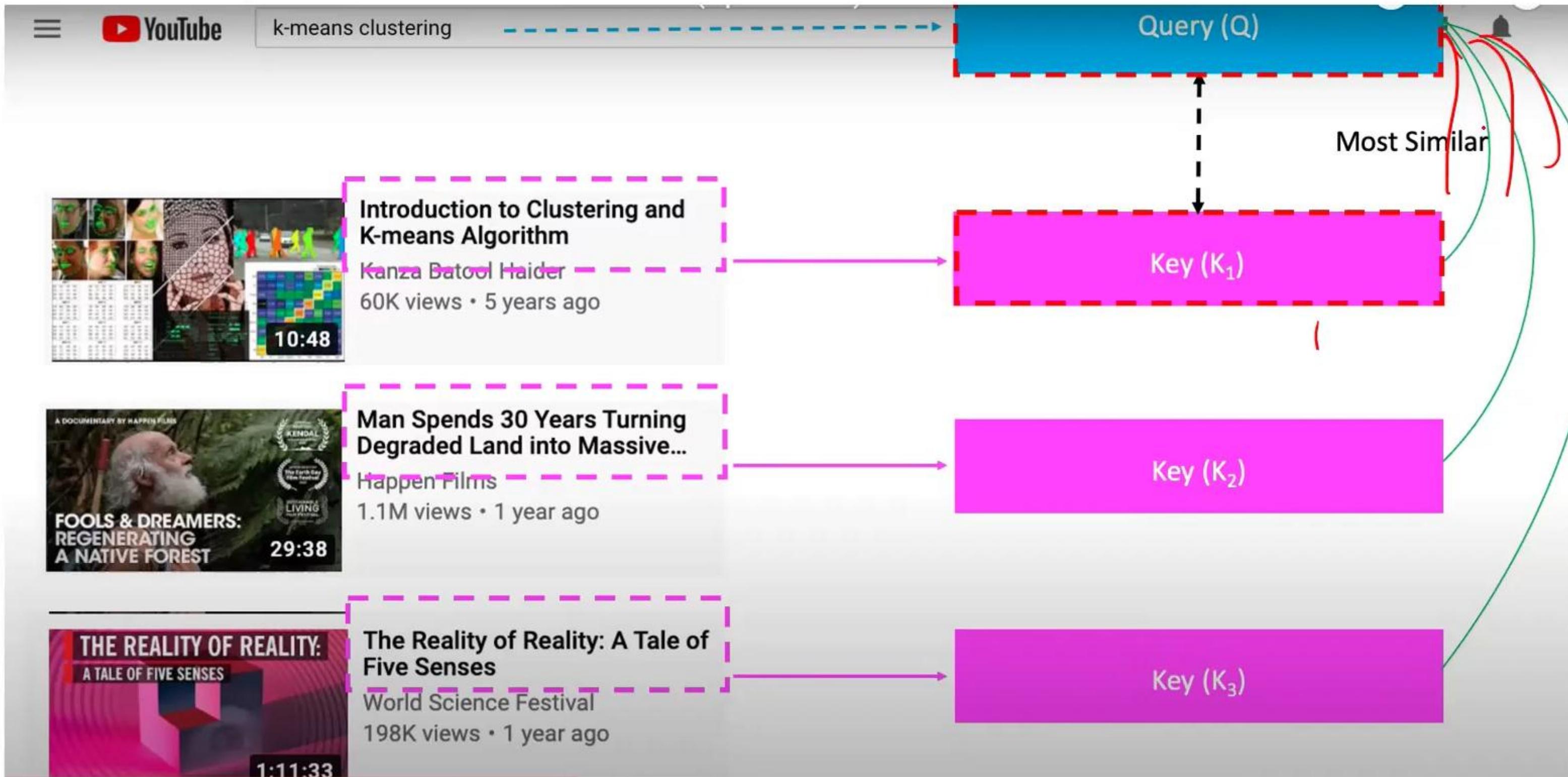


# Multi Head Attention

3



# Multi Head Attention



# Multi Head Attention

X

YouTube k-means clustering

Query (Q)

Key ( $K_1$ )

Value ( $V_1$ )

Introduction to Clustering and K-means Algorithm  
Kanza Batool Haider  
60K views • 5 years ago

10:48

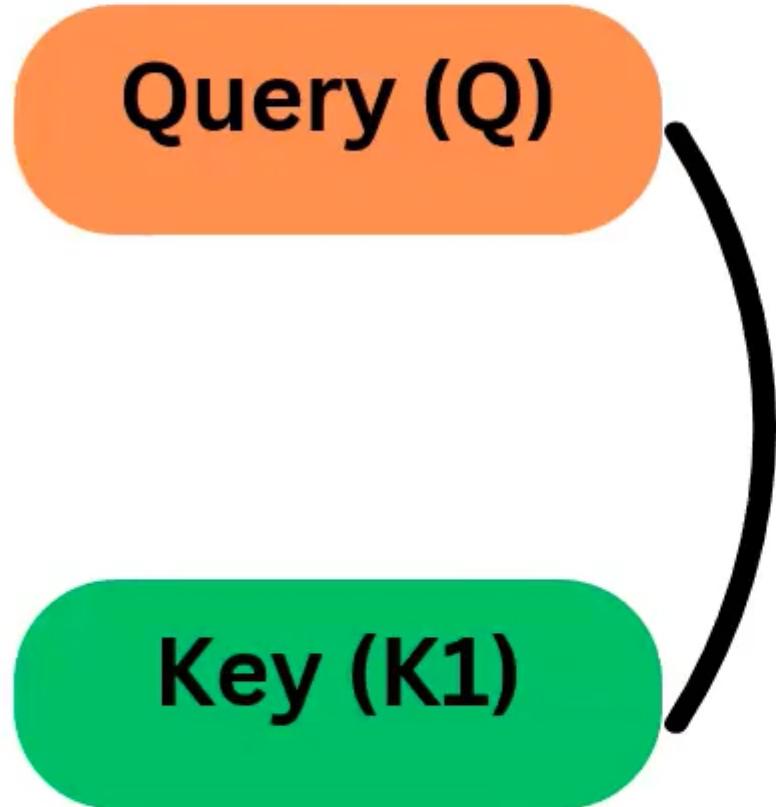
A diagram illustrating the Multi-Head Attention mechanism. On the left, a YouTube search result for "k-means clustering" is shown. A dashed orange box encloses a video thumbnail titled "Introduction to Clustering and K-means Algorithm" by Kanza Batool Haider, which has 60K views and is 10:48 long. An orange arrow points from this video thumbnail to a pink box labeled "Key ( $K_1$ )". Another orange arrow points from the same video thumbnail to an orange box labeled "Value ( $V_1$ )". The top right of the slide shows a blue bar labeled "Query (Q)".

# Why Transformer?

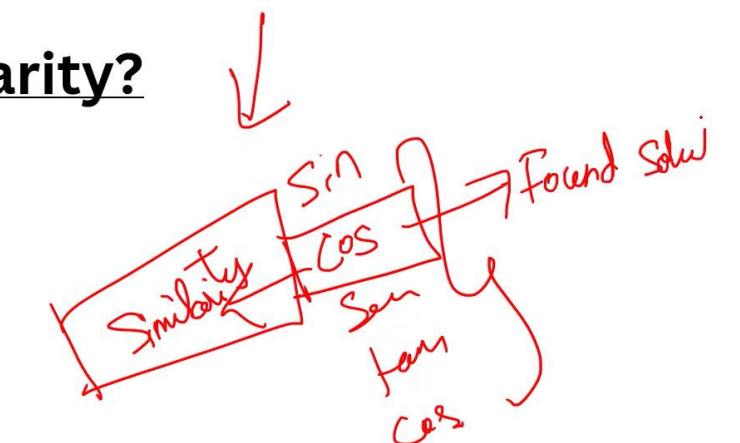
Ramisha codes Python Programming

# Similarity?

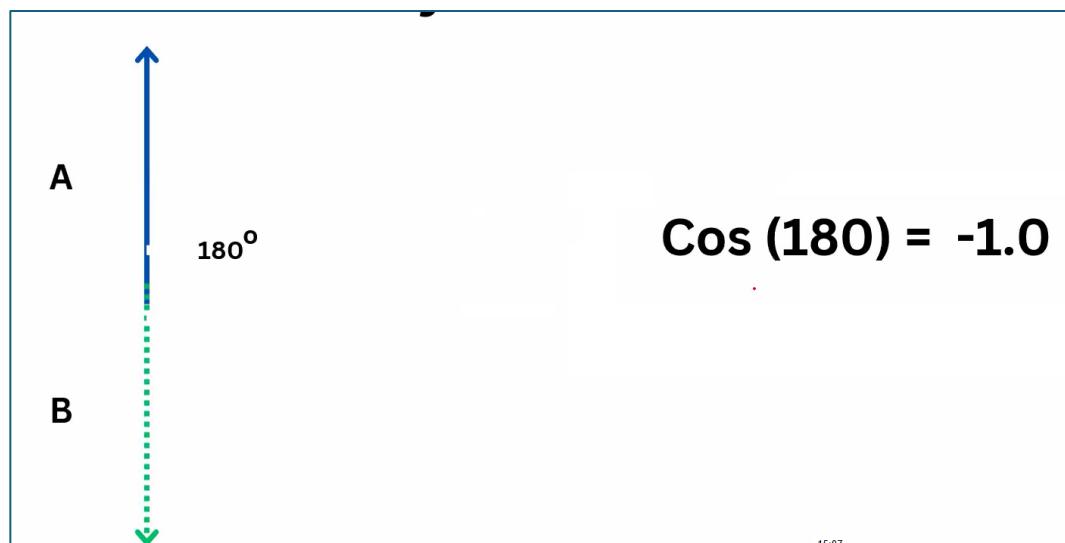
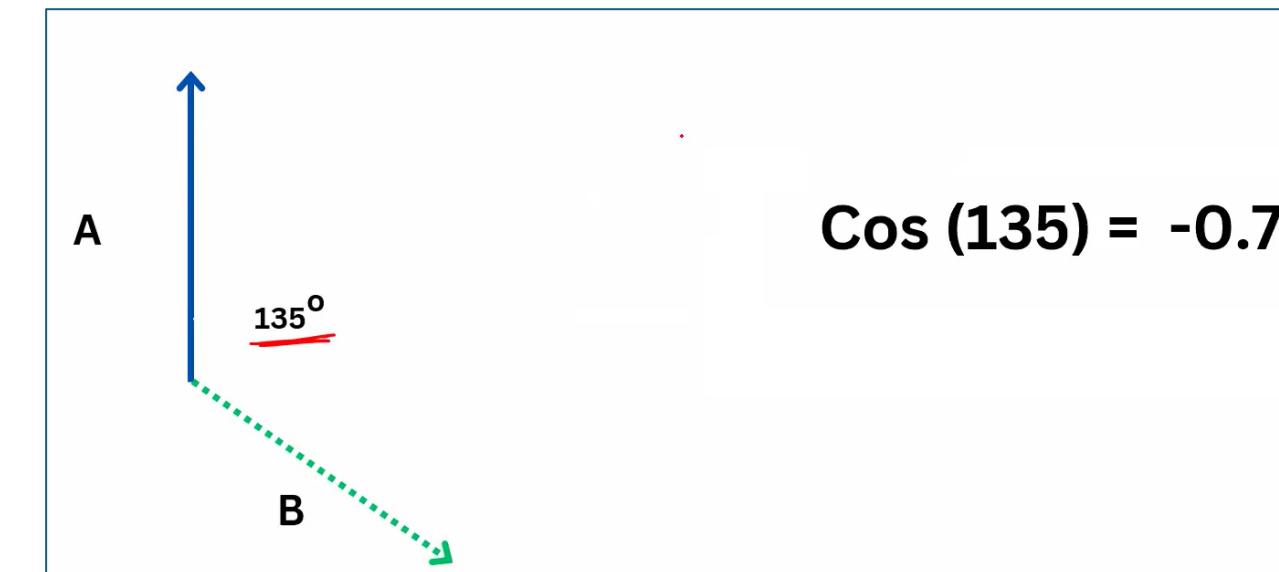
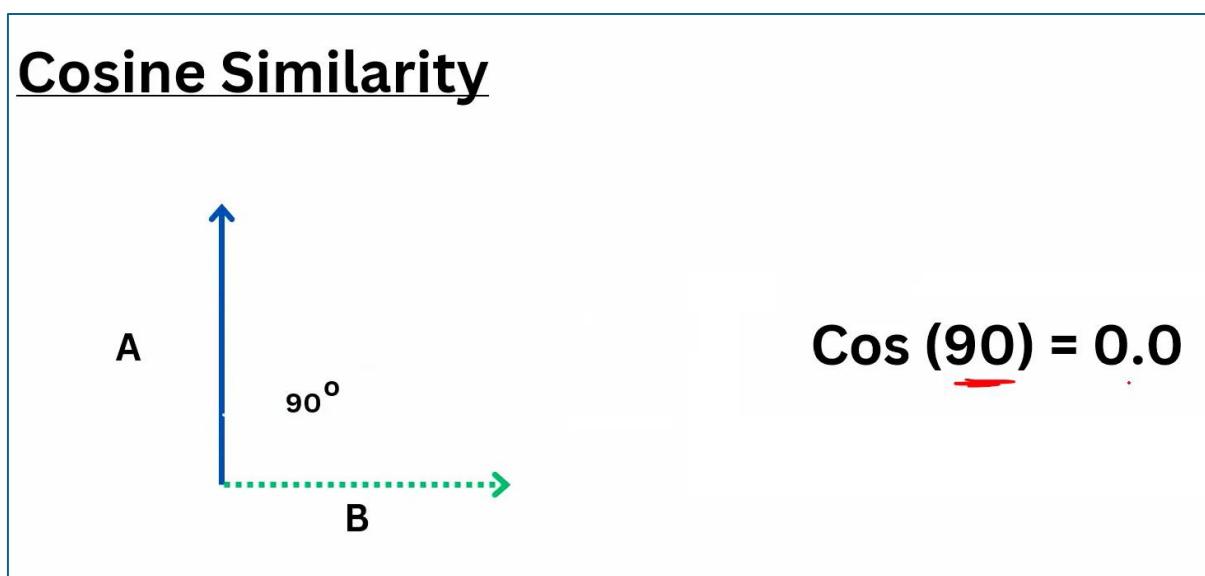
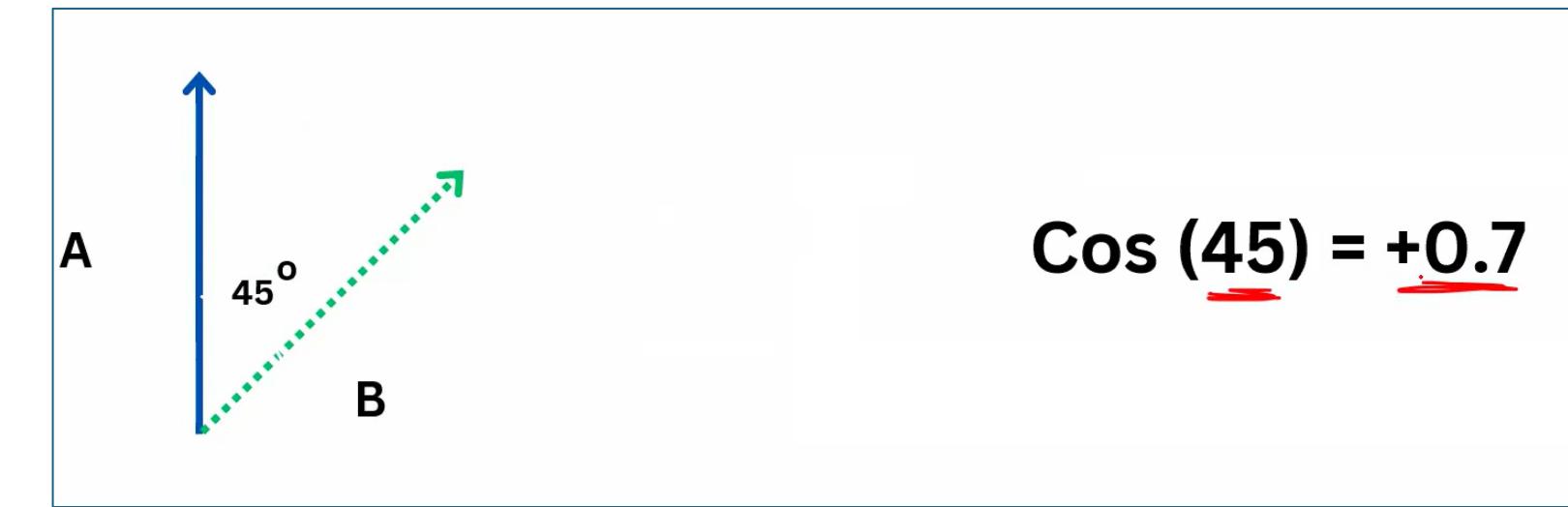
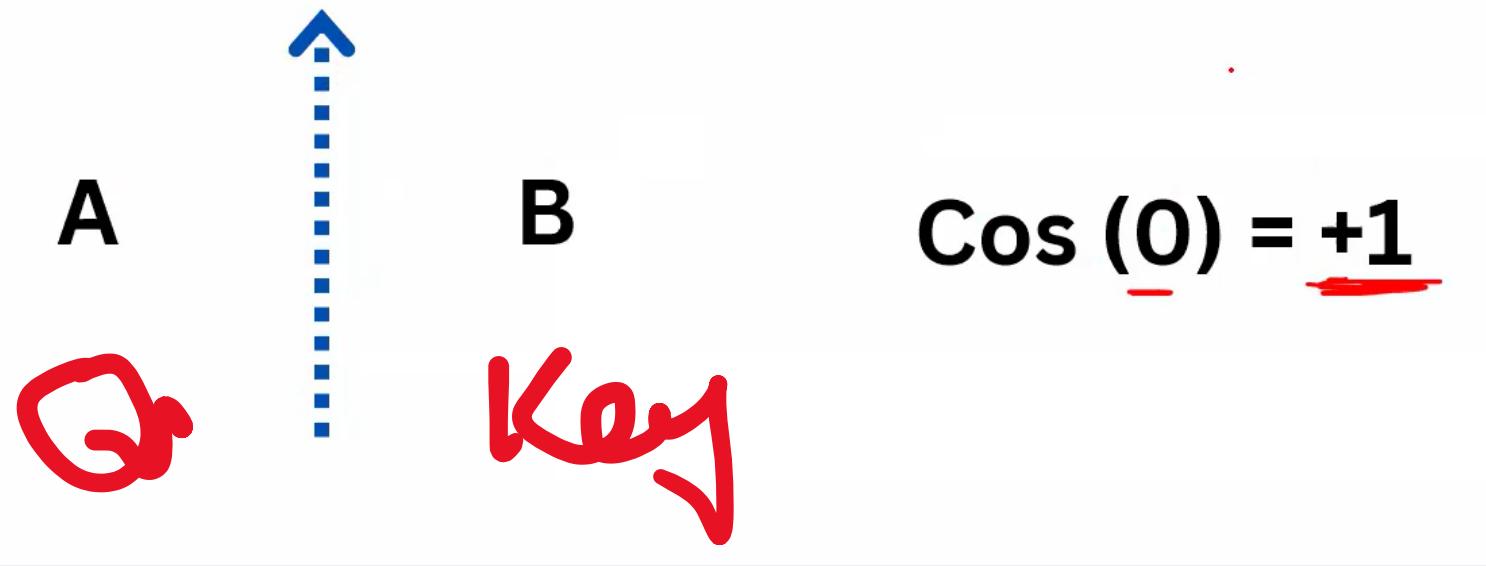
-1      to      +1  
Most Dissimilar      Most similar



Similarity?



*Similarity?*



-1 to +1  
Most Dissimilar Most similar

$$\sin(0) = 0 \times 1$$

$$\sin(90) = 1$$

$$\sin(180) = 0 \times (-1)$$

$$\cos(0) = 1$$

$$\cos(90) = 0$$

$$\cos(180) = -1$$

$$\cos(A, B) = \frac{A \cdot B}{\|A\| \|B\|}$$

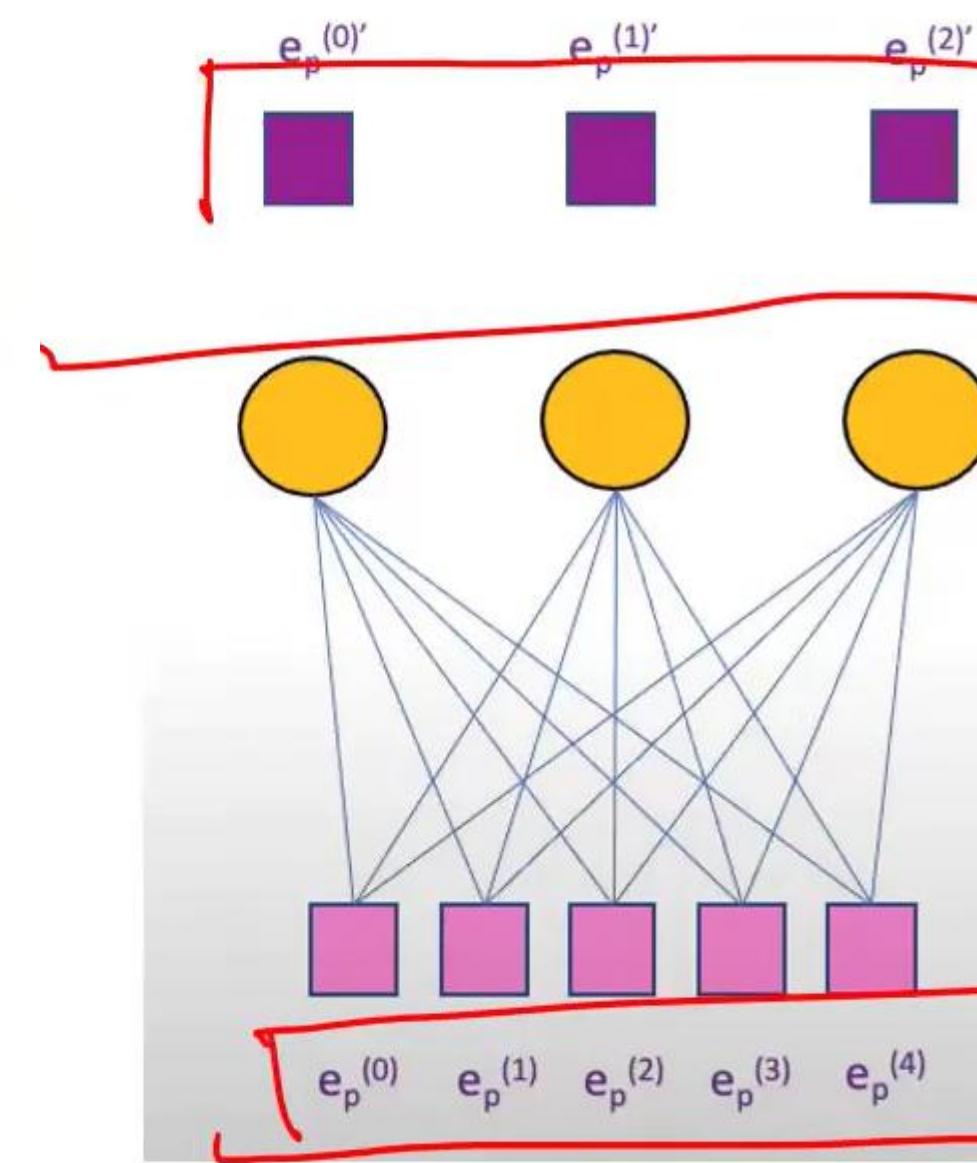
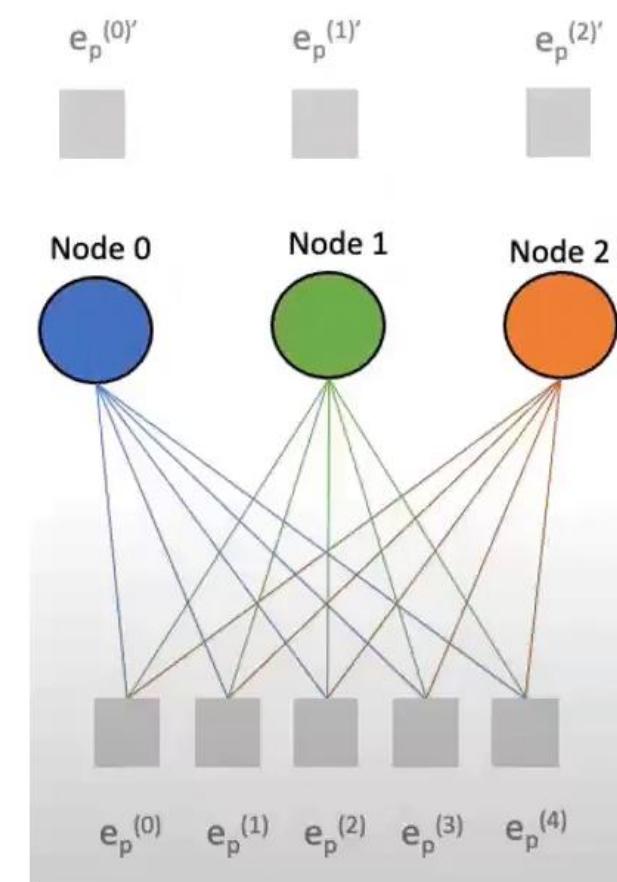
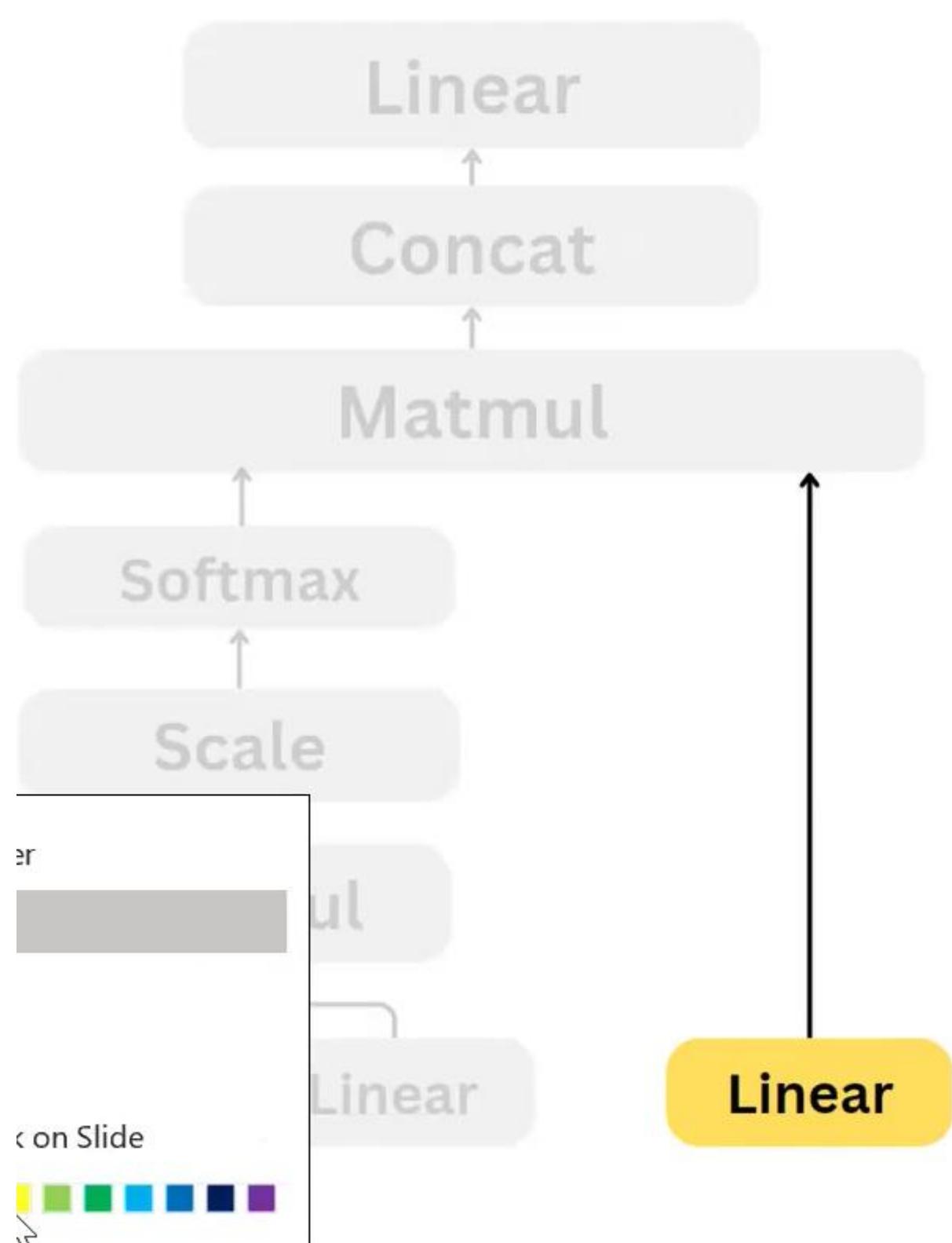
Similarity b/w Vectors



$$\text{Similarity}(A, B) = \frac{A \cdot B}{\text{Scaling}}$$

Similarity b/w Query and Key

$$\text{Similarity}(Q, K) = \frac{Q \cdot K^T}{\text{Scaling}}$$

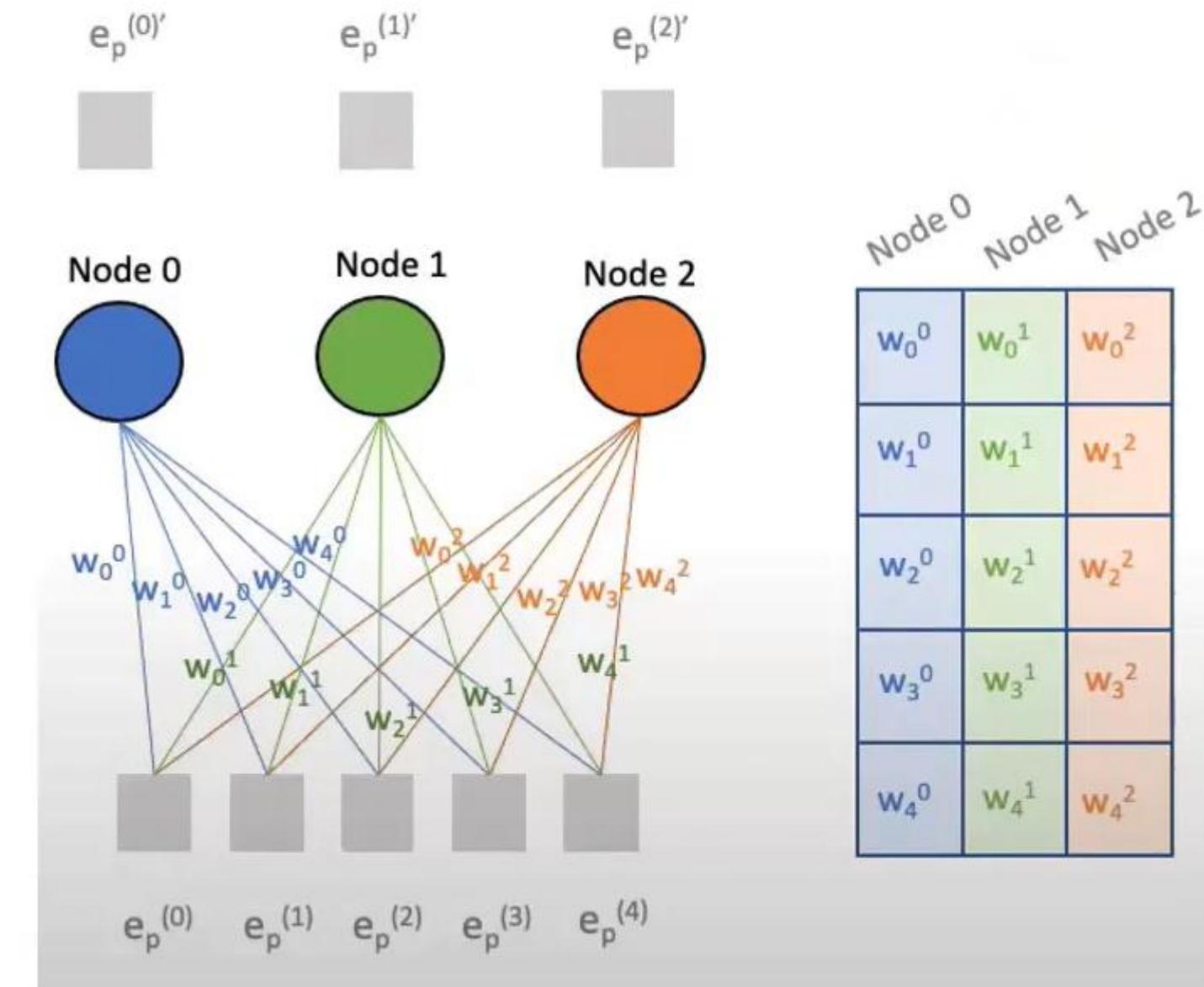


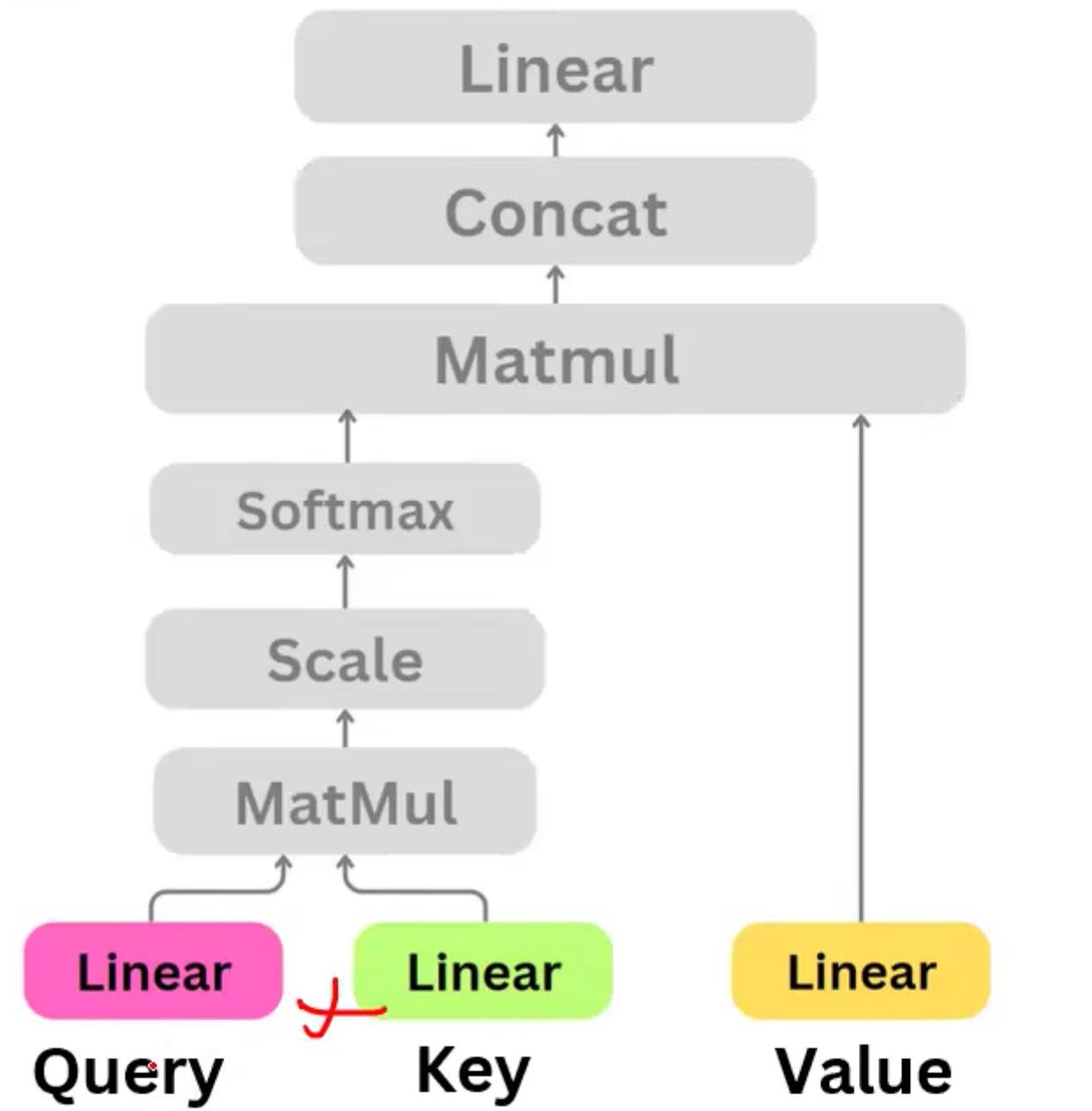
r  
t  
il



*Acti.*

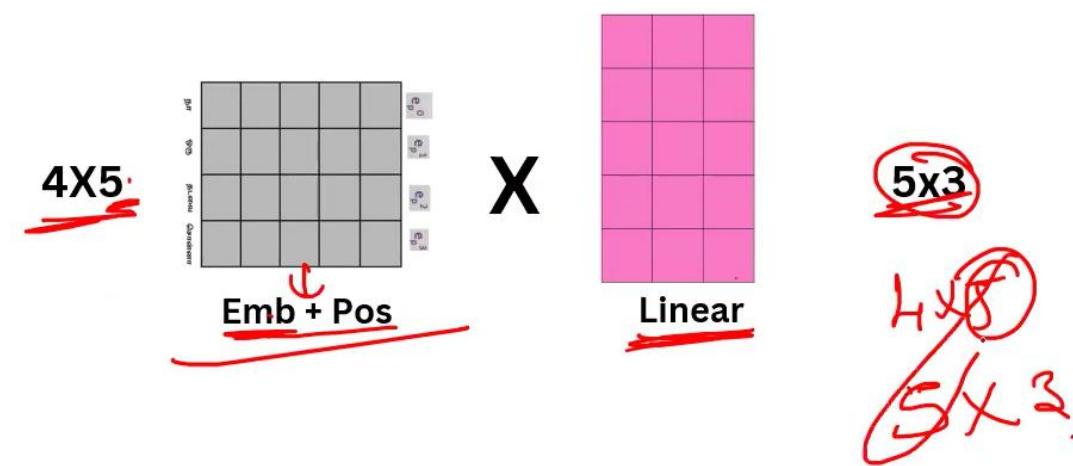
**Linear**

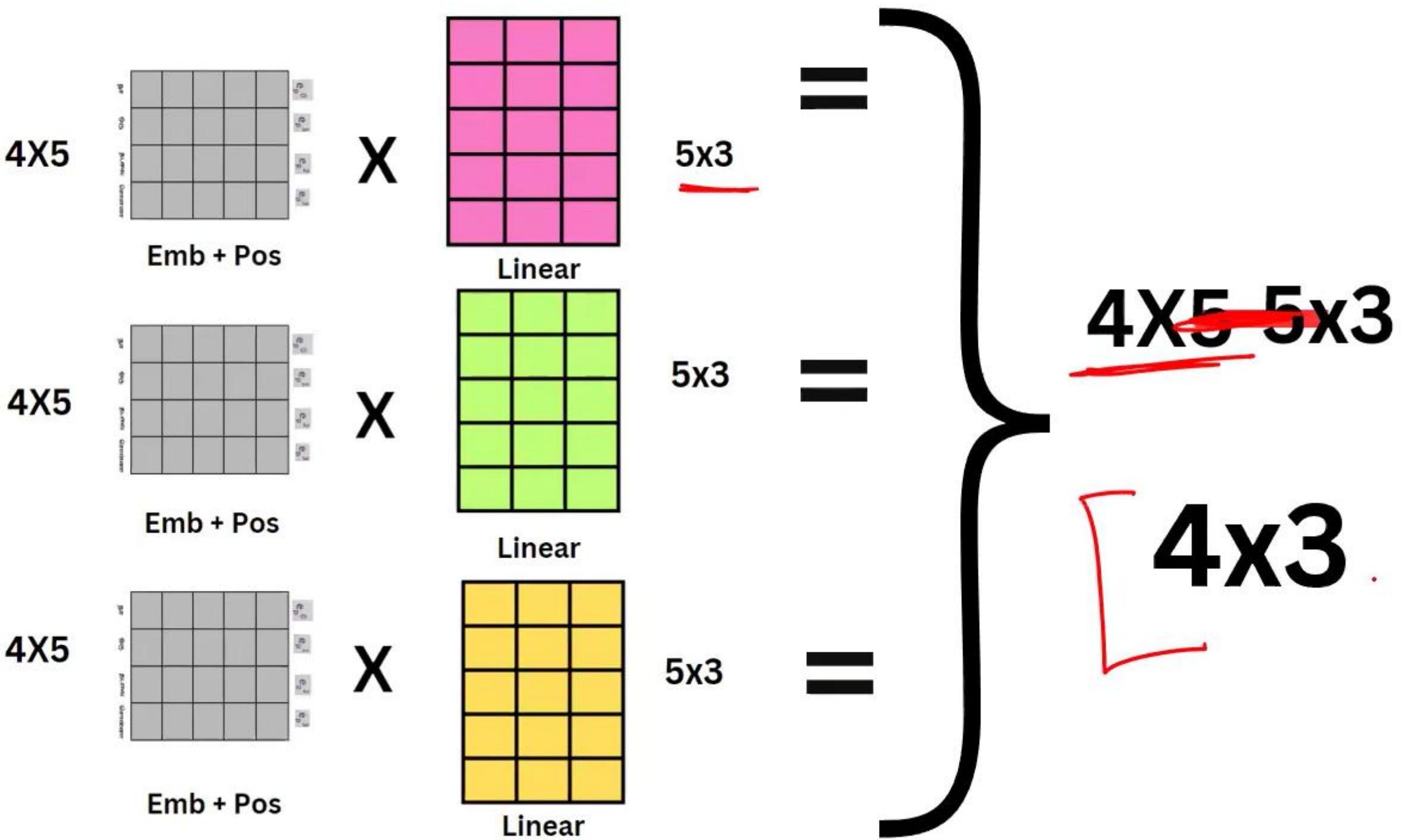
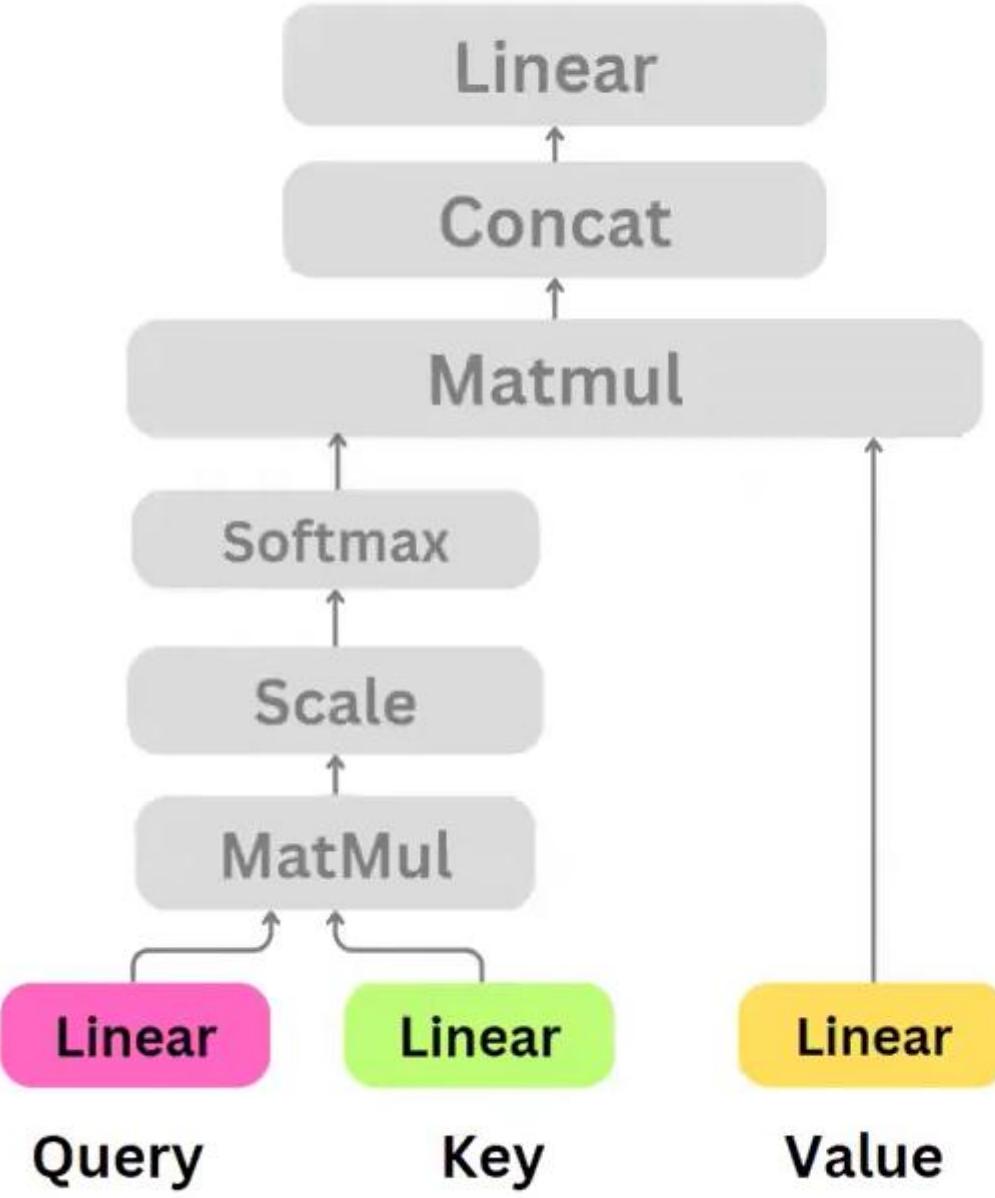


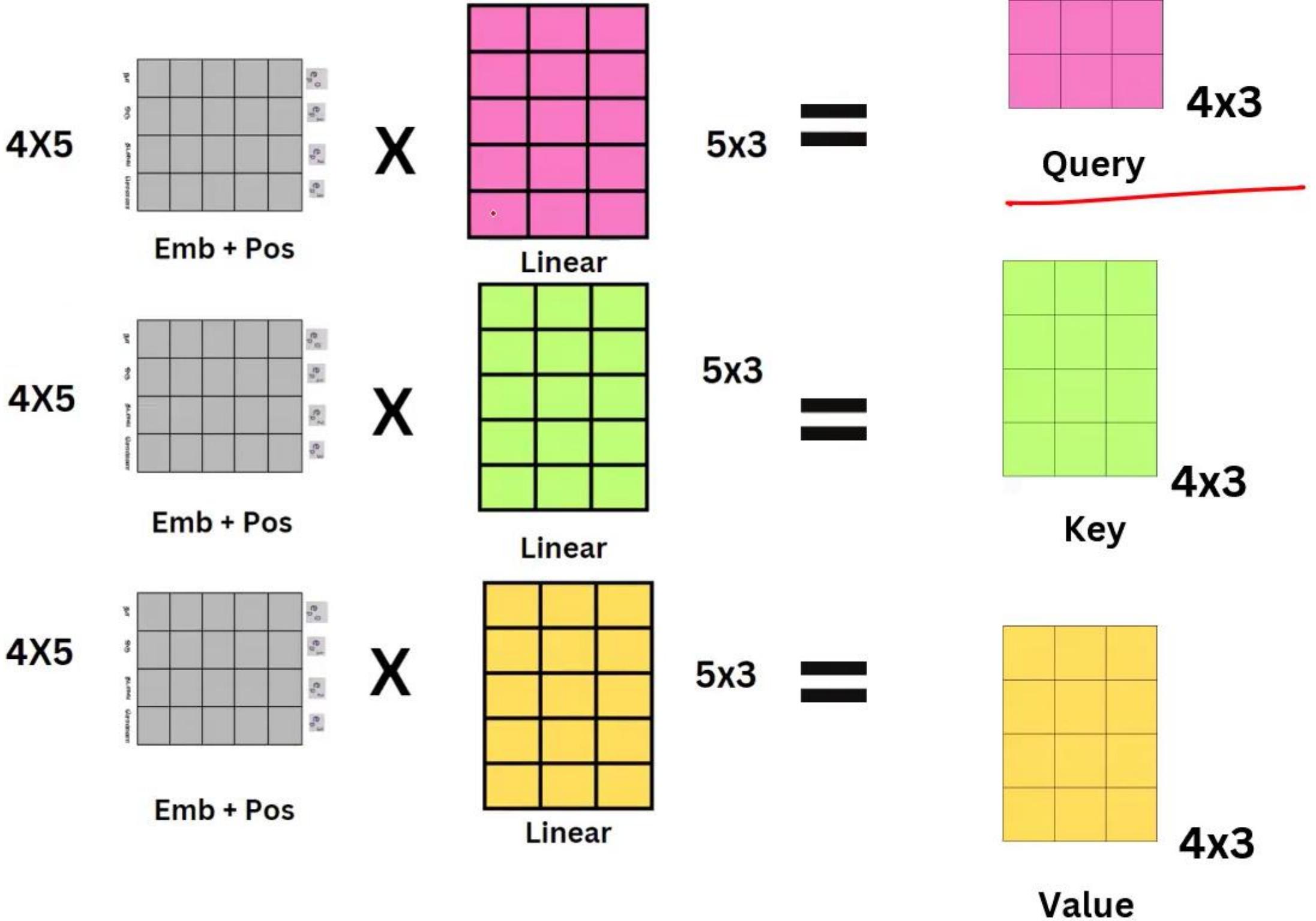
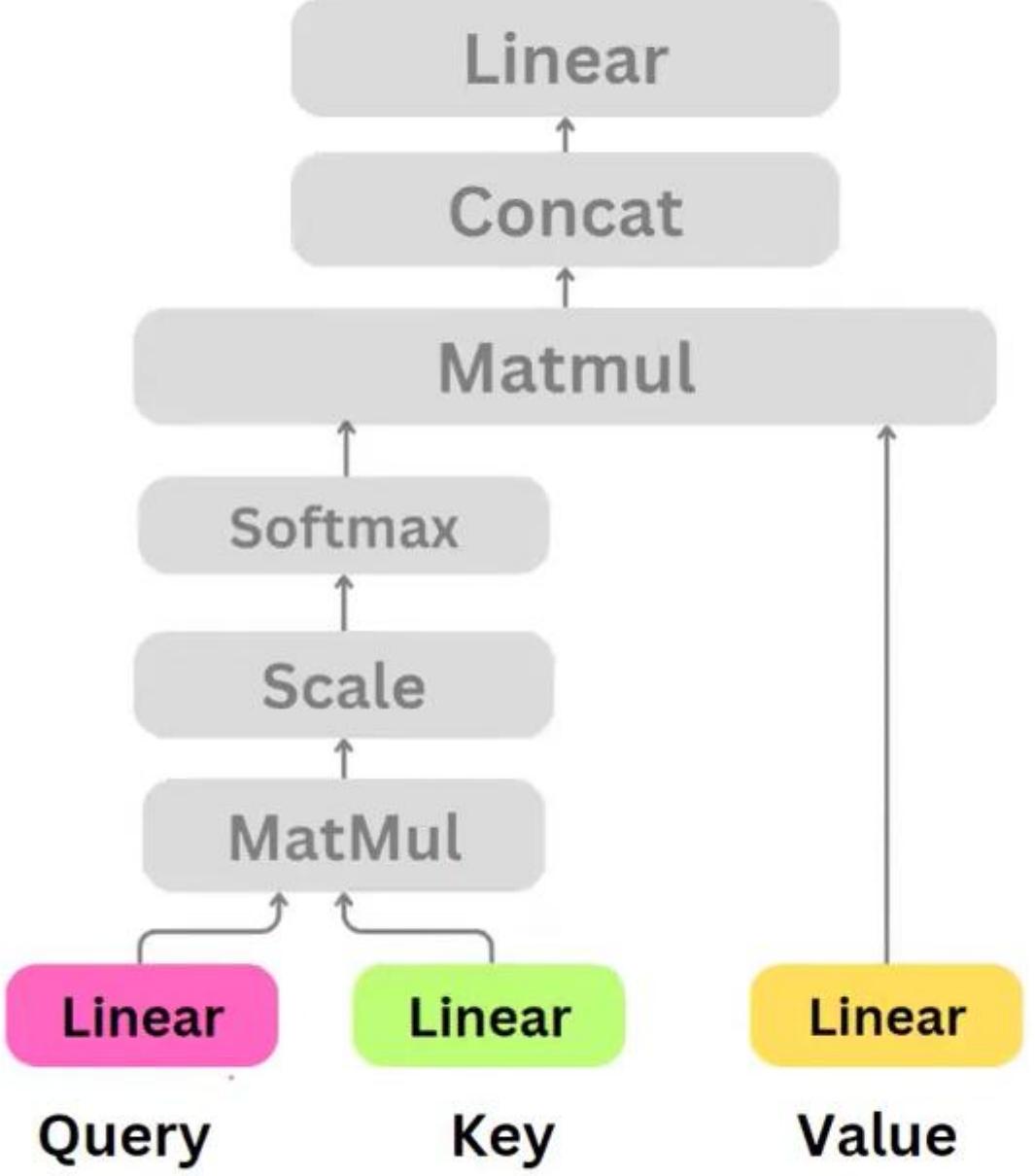


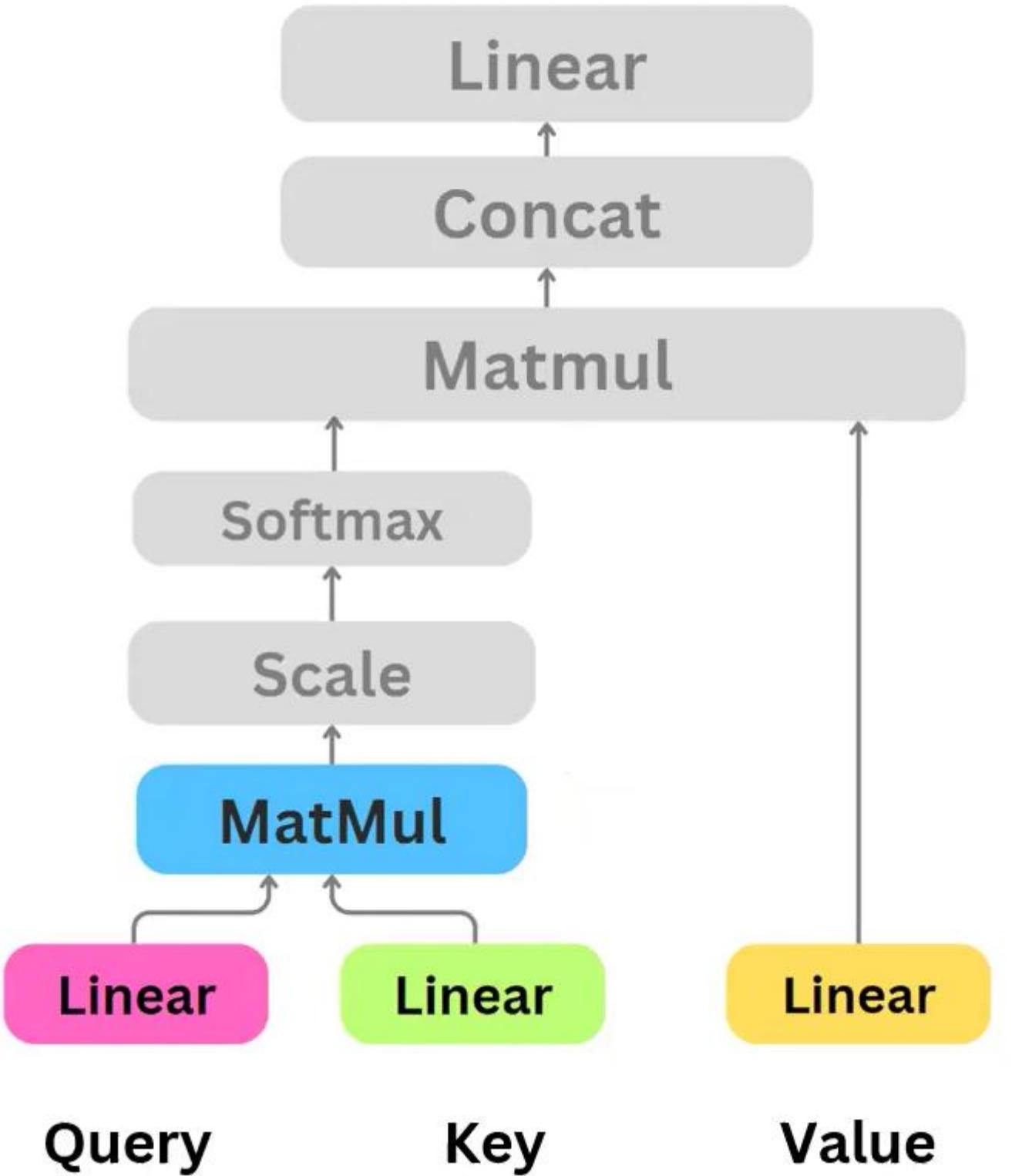
$e_p^0$	$e_p^1$	$e_p^2$	$e_p^3$
பொ	புரு	நடவெ	சொல்லை
நா	புரு	நடவெ	சொல்லை
பொ	புரு	நடவெ	சொல்லை
நா	புரு	நடவெ	சொல்லை

$e_p^0$	$e_p^1$	$e_p^2$	$e_p^3$
பொ	புரு	நடவெ	சொல்லை
நா	புரு	நடவெ	சொல்லை
பொ	புரு	நடவெ	சொல்லை
நா	புரு	நடவெ	சொல்லை









$$\text{Query} \quad 4 \times 3$$

\*

$$\text{Key} \quad 3 \times 4$$

=

$$4 \times 4$$

$$\mathbf{Q} \cdot \mathbf{K}^T$$

**Self Attention**

நா	ஓரு	தட்டை	சொன்னா	
நா	0.3	0.4	0.9	0.1
ஓரு	0.2	0.6	0.8	0.4
தட்டை	0.4	0.3	0.5	0.3
சொன்னா	0.1	0.2	0.3	0.2

4x4

4x4



Attention Filter

\*



Value/Original

=



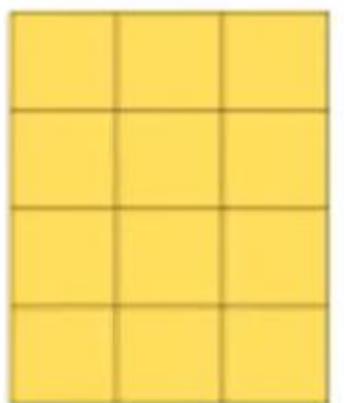
Filtered Important Features

	நா	இரு	தடவை	சொன்னா
நா	0.3	0.4	0.9	0.1
இரு	0.2	0.6	0.8	0.4
தடவை	0.4	0.3	0.5	0.3
சொன்னா	0.1	0.2	0.3	0.2

**4x4**

Attention Filter

\*



Value

**4x3**

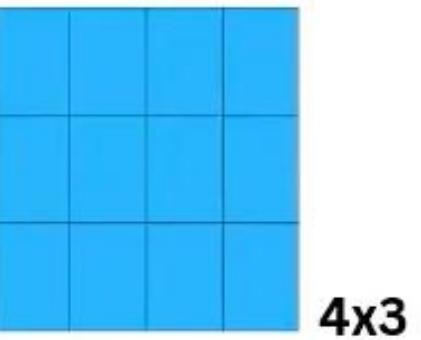
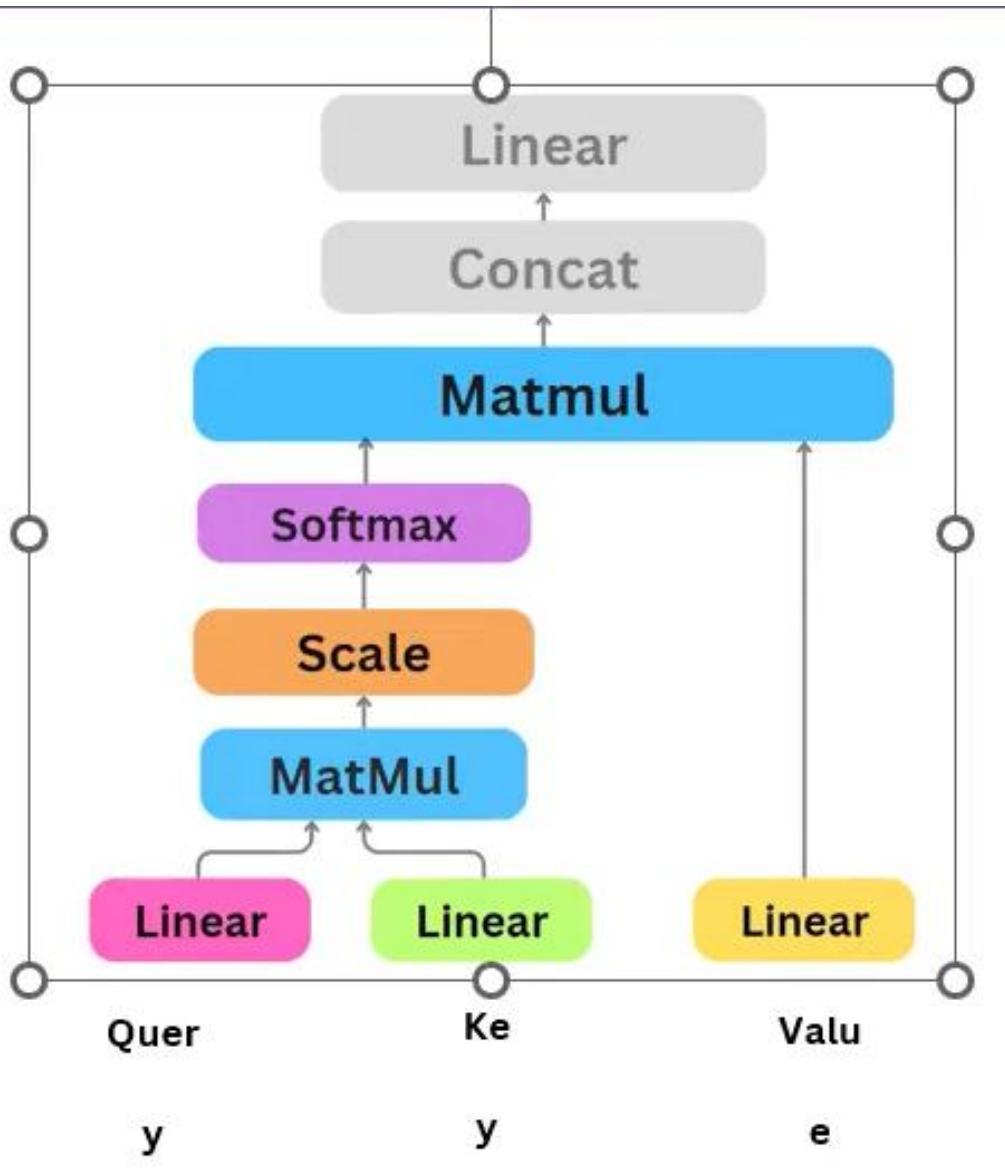
=



Filtered Important

7

Features



Filtered Important  
Features

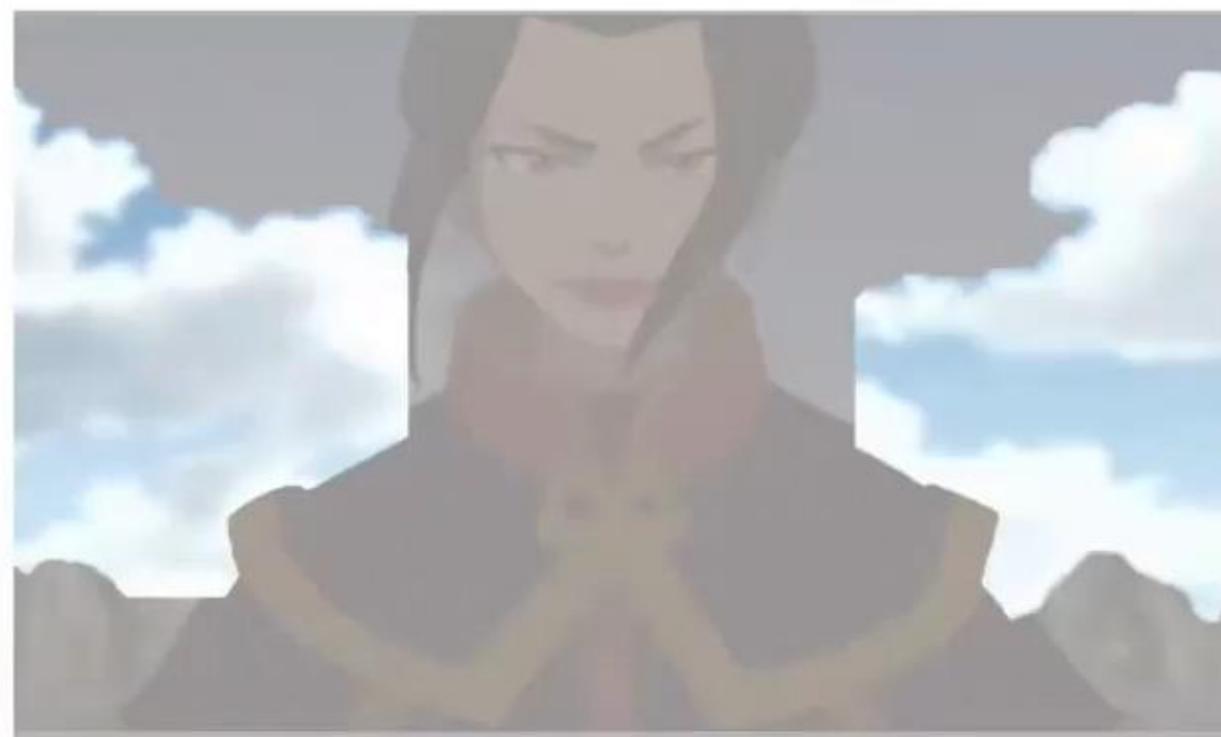
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right)V$$

## Multi Head Attention-Similarity

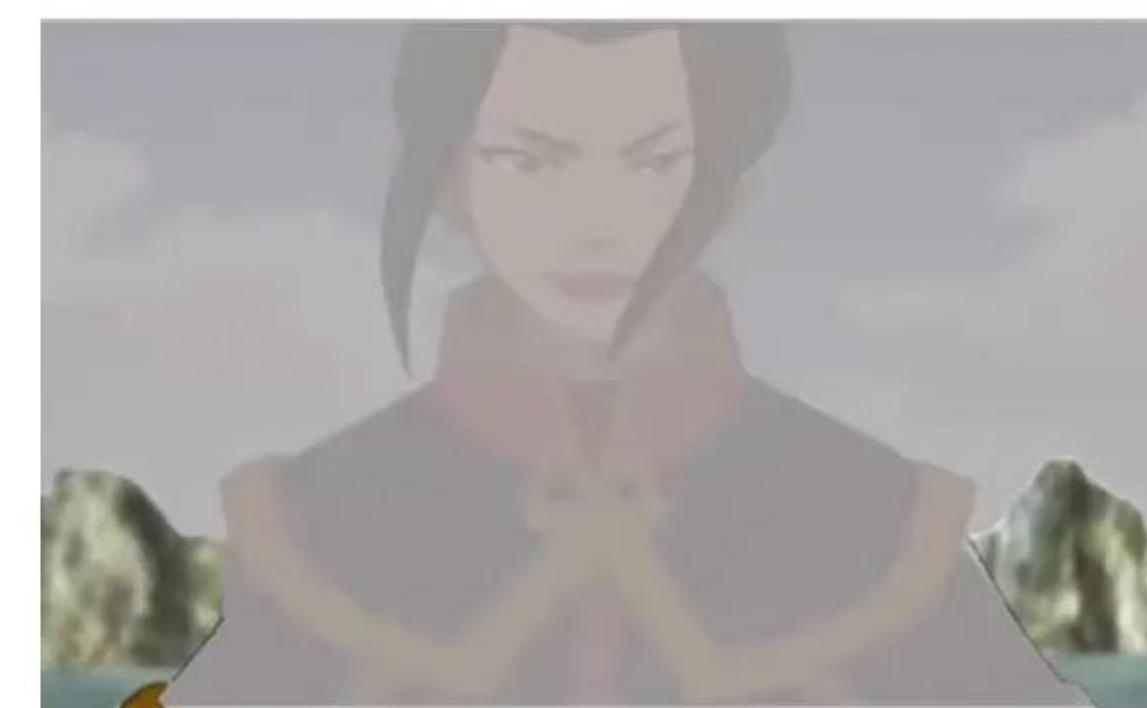
Attention Filter 1



Attention Filter 2

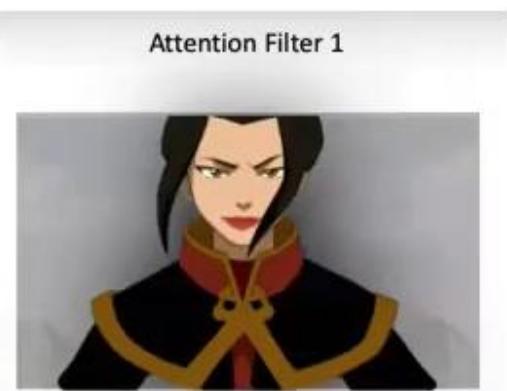
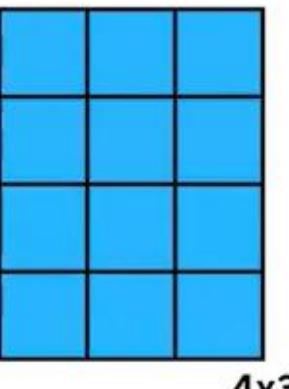


Attention Filter 3



# Multi Head Attention-Similarity

**Head- 1**

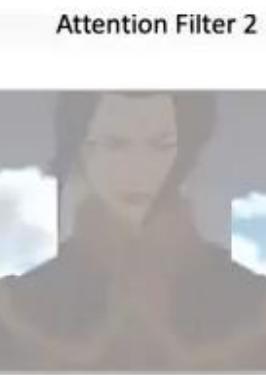
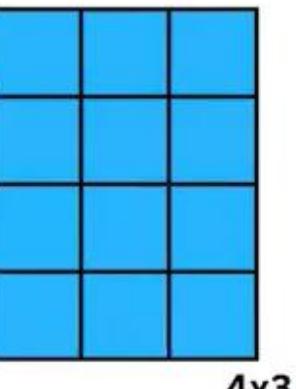


Filtered Important  
Features

$$\begin{matrix} \text{BF} & \text{GB} & \text{Blue} & \text{Green} \\ \text{BF} & 0.3 & 0.4 & 0.9 & 0.1 \\ \text{GB} & 0.2 & 0.6 & 0.8 & 0.4 \\ \text{Blue} & 0.4 & 0.3 & 0.5 & 0.3 \\ \text{Green} & 0.1 & 0.2 & 0.3 & 0.2 \end{matrix} \quad * \quad \begin{matrix} \text{Value} \\ 4 \times 3 \end{matrix}$$

Attention Filter 4x4

**Head- 2**

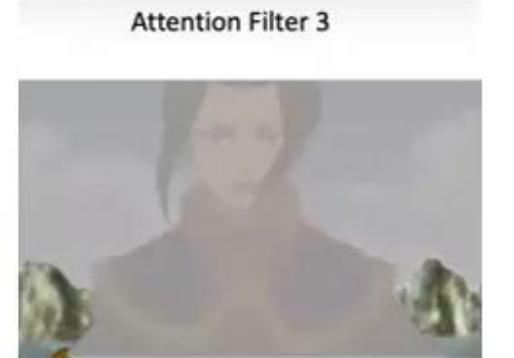


Filtered Important  
Features

$$\begin{matrix} \text{BF} & \text{GB} & \text{Blue} & \text{Green} \\ \text{BF} & 0.3 & 0.4 & 0.9 & 0.1 \\ \text{GB} & 0.2 & 0.6 & 0.8 & 0.4 \\ \text{Blue} & 0.4 & 0.3 & 0.5 & 0.3 \\ \text{Green} & 0.1 & 0.2 & 0.3 & 0.2 \end{matrix} \quad * \quad \begin{matrix} \text{Value} \\ 4 \times 3 \end{matrix}$$

Attention Filter 4x4

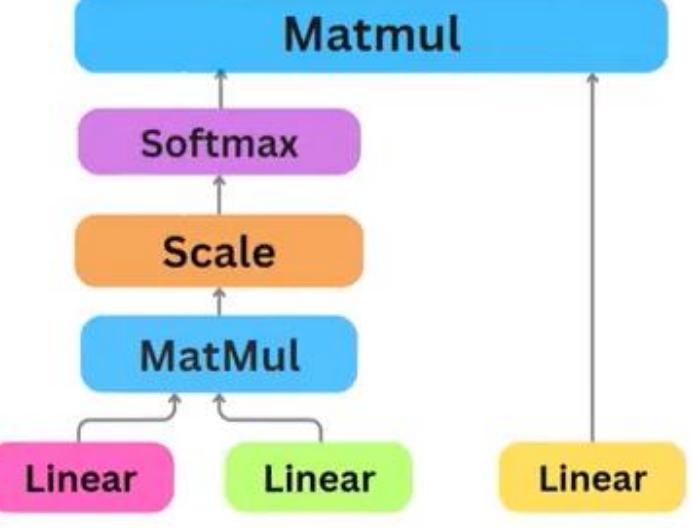
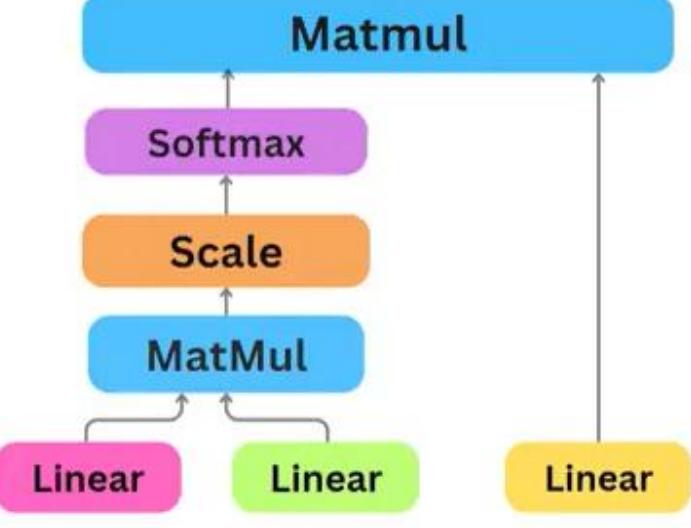
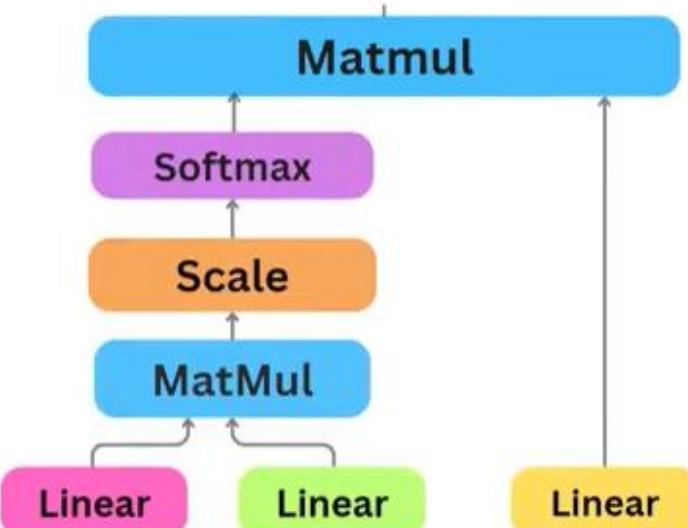
**Head-3**

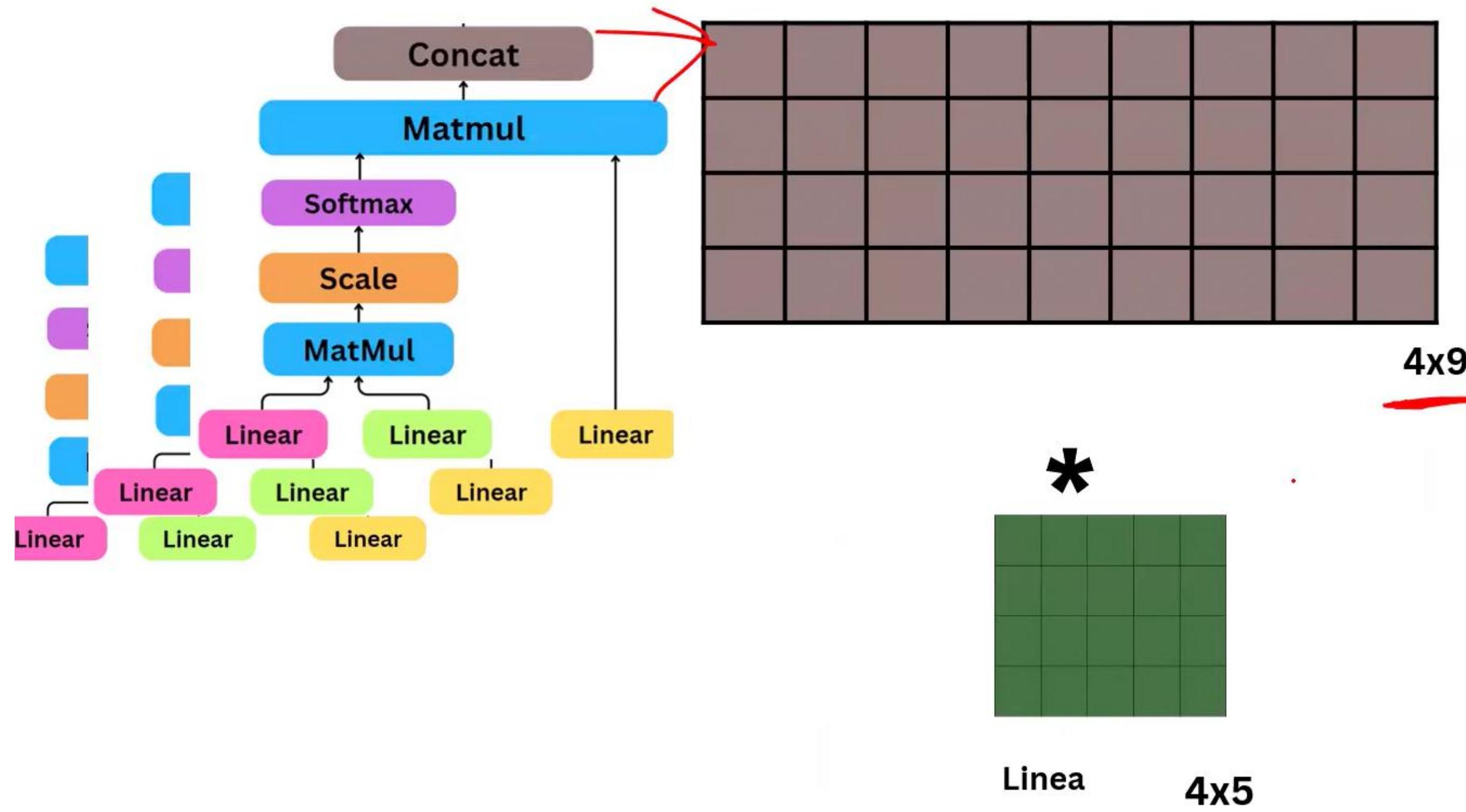


Filtered Important  
Features

$$\begin{matrix} \text{BF} & \text{GB} & \text{Blue} & \text{Green} \\ \text{BF} & 0.3 & 0.4 & 0.9 & 0.1 \\ \text{GB} & 0.2 & 0.6 & 0.8 & 0.4 \\ \text{Blue} & 0.4 & 0.3 & 0.5 & 0.3 \\ \text{Green} & 0.1 & 0.2 & 0.3 & 0.2 \end{matrix} \quad * \quad \begin{matrix} \text{Value} \\ 4 \times 3 \end{matrix}$$

Attention Filter 4x4





## Multi Head Attention-Similarity

