



Transformers

Prof. Khaled Mostafa El-Sayed

khaledms@fci-cu.edu.eg

Faculty of Computers and Artificial Intelligence

Cairo University

+

Aug 2020

Transformers

Processing Sequential Data

Transformer Building Blocks

Self Attention Mechanism

Multi-Head Attention

Parallelization Implementation Issues

Positional Information

Building Transformer Block

Transformer Decoder

+

Transformer Encoder as a Classifier

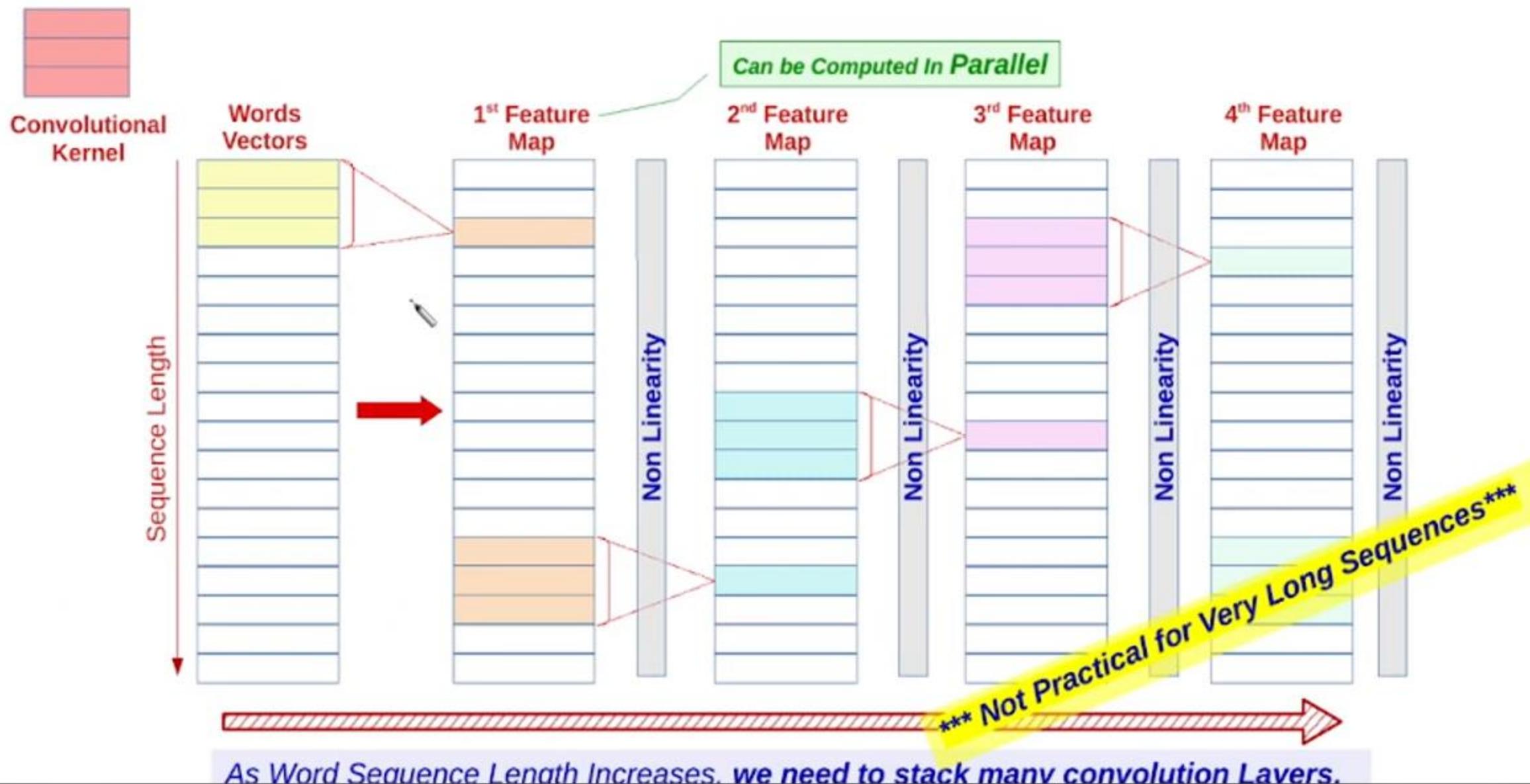
What is Next ?

Processing Sequential Data

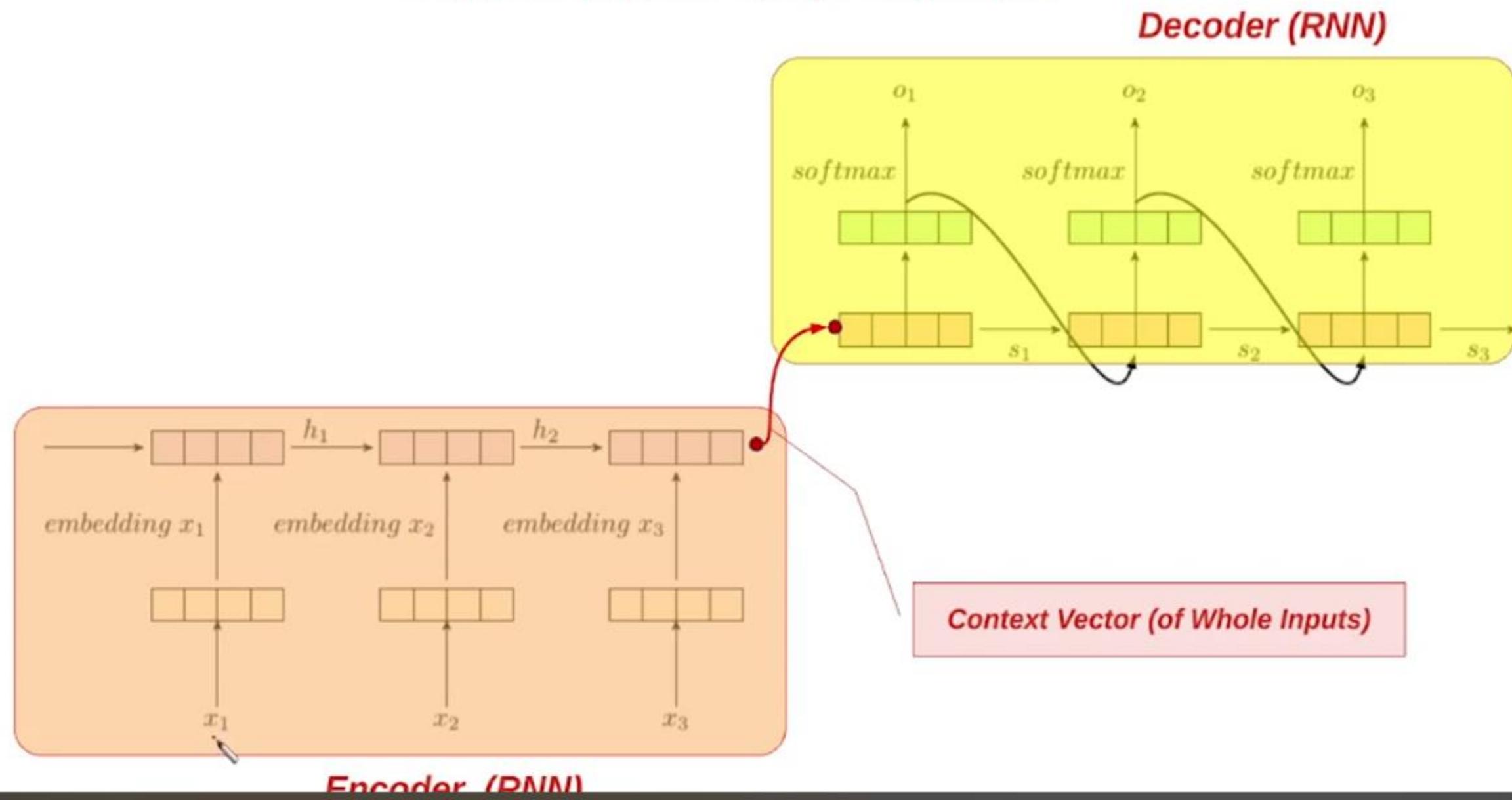
1-D Convolutional Neural Network (1-D CNN)

Recurrent Neural Network (RNN)

1-D Convolutional Neural Network



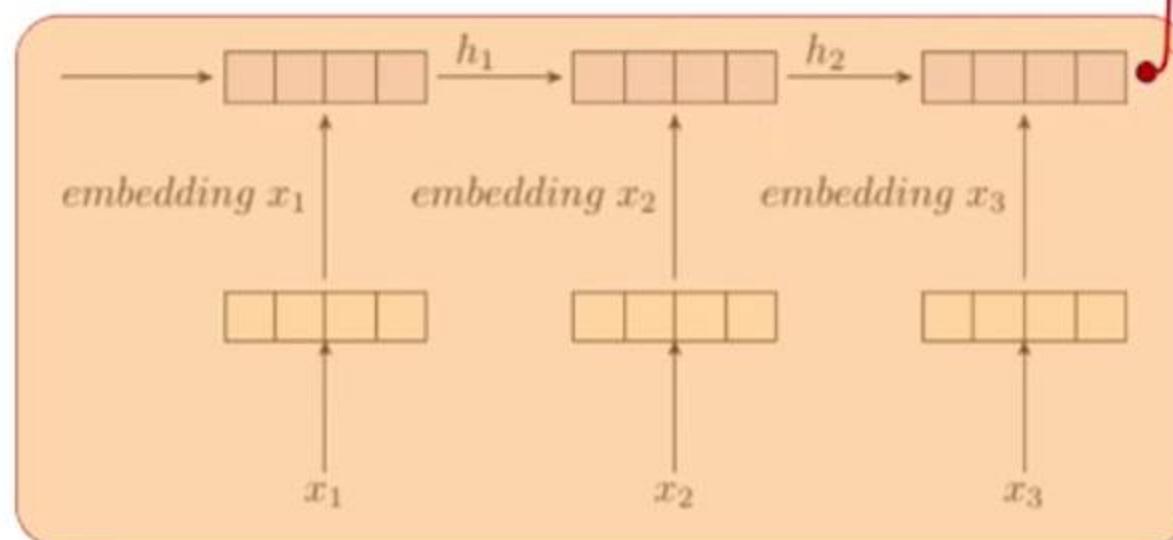
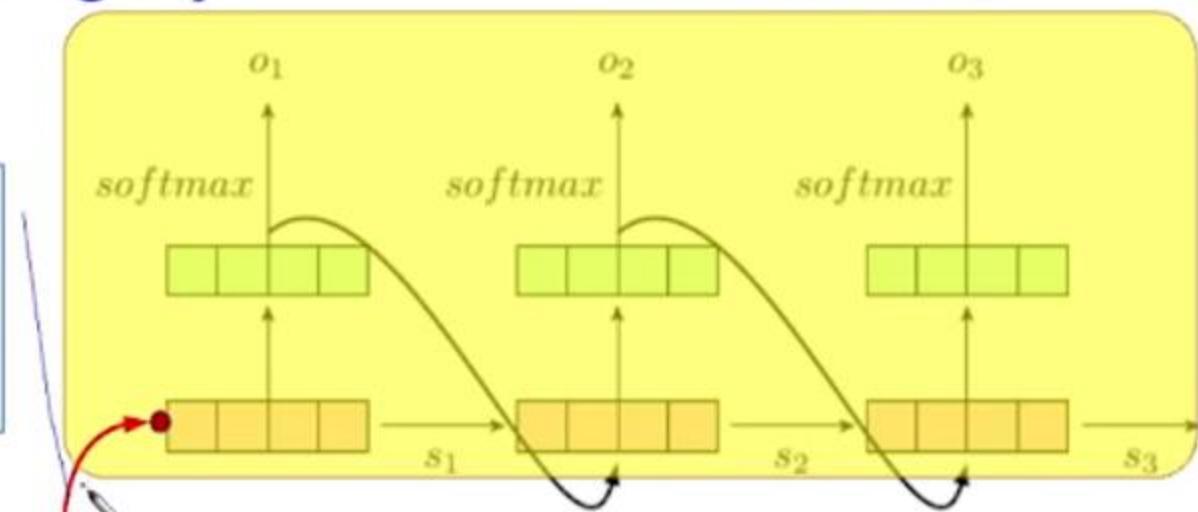
Traditional RNN Seq2Seq model



Traditional RNN Seq2Seq model (Challenges)

Decoder (RNN)

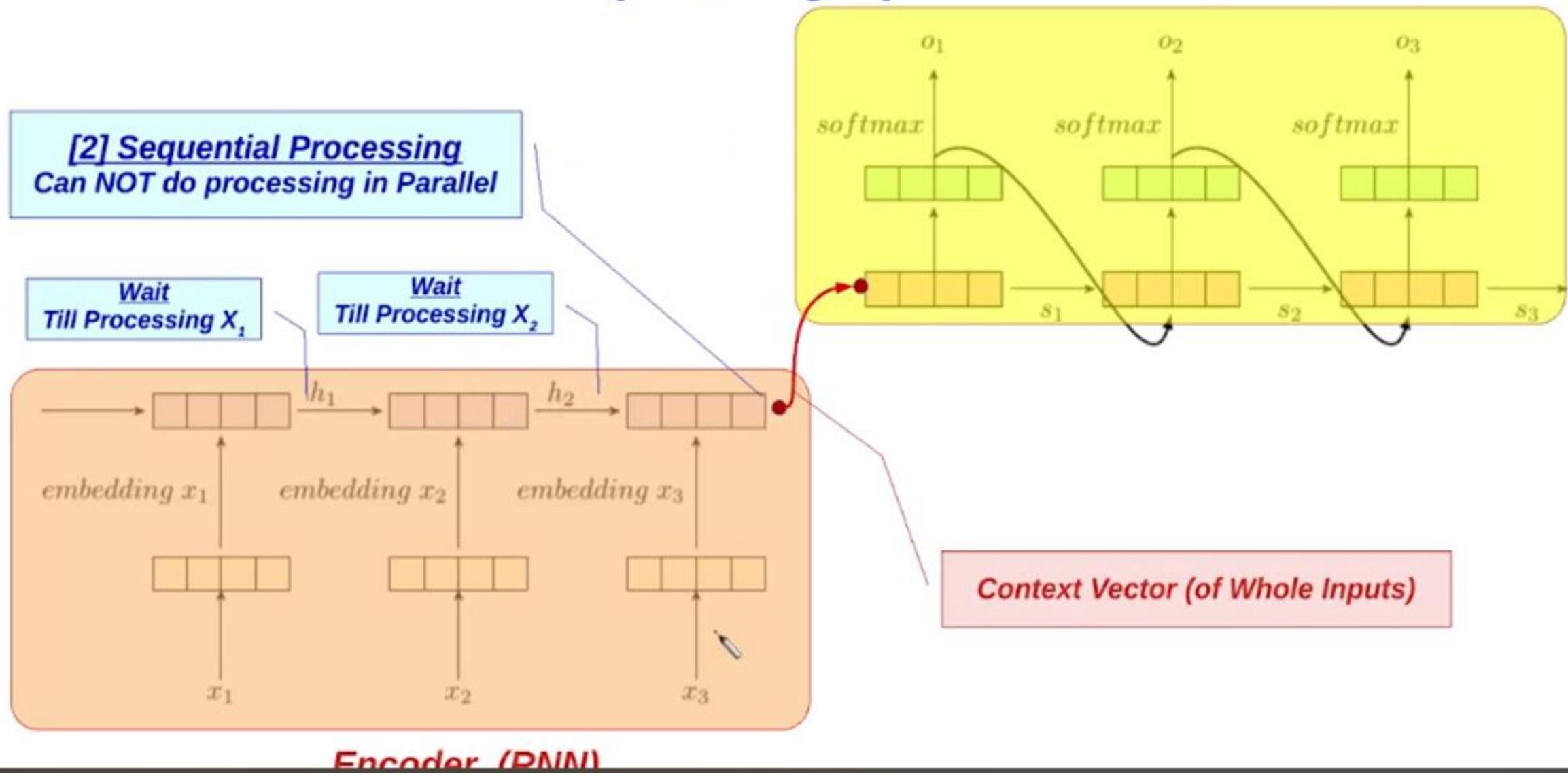
[1] Bottleneck
*The meaning of the entire input sequence is
 Expected to be captured by a single context vector
 with fixed dimensionality*



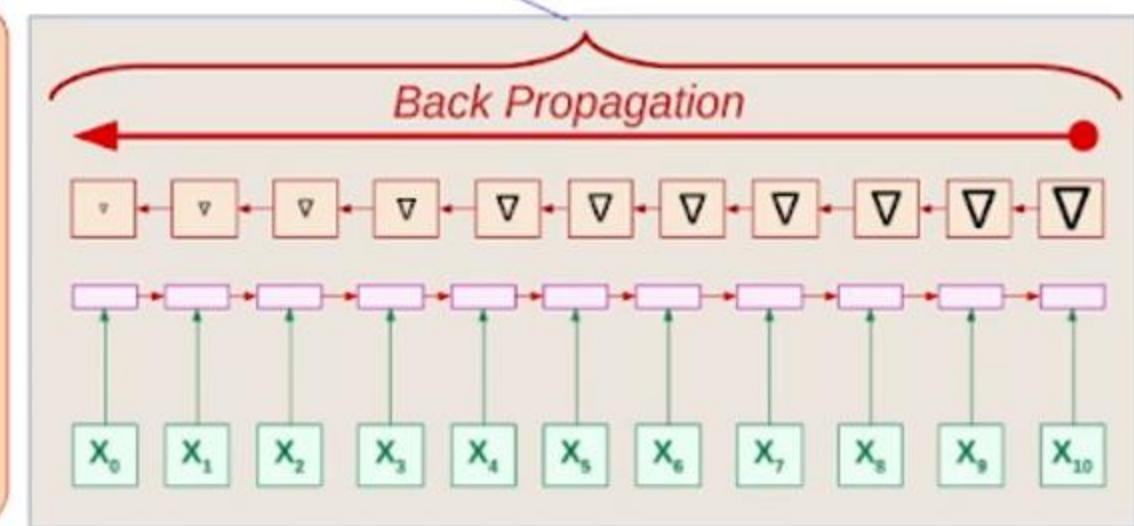
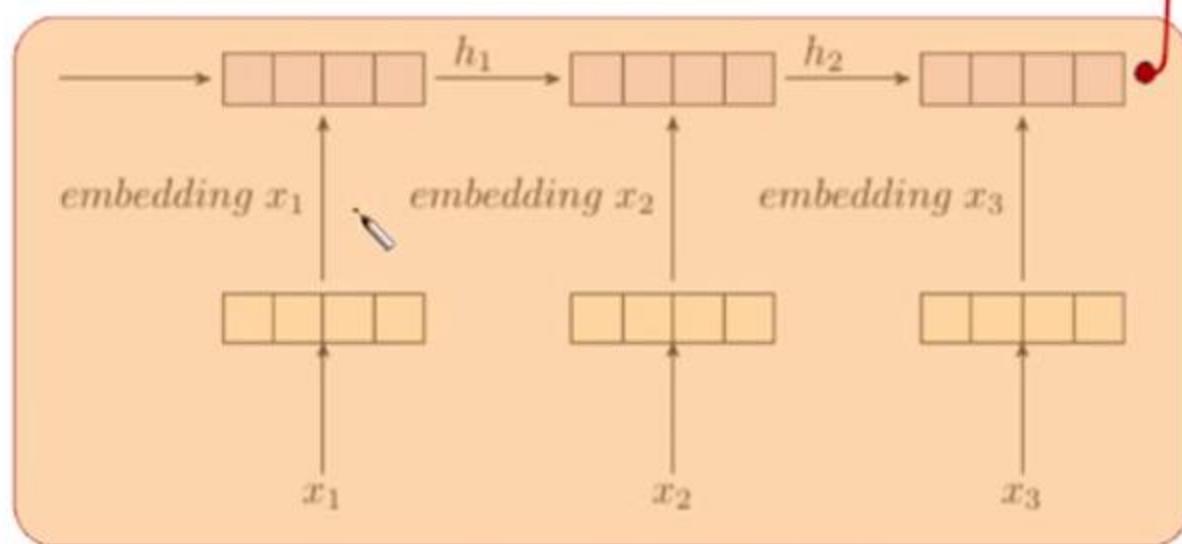
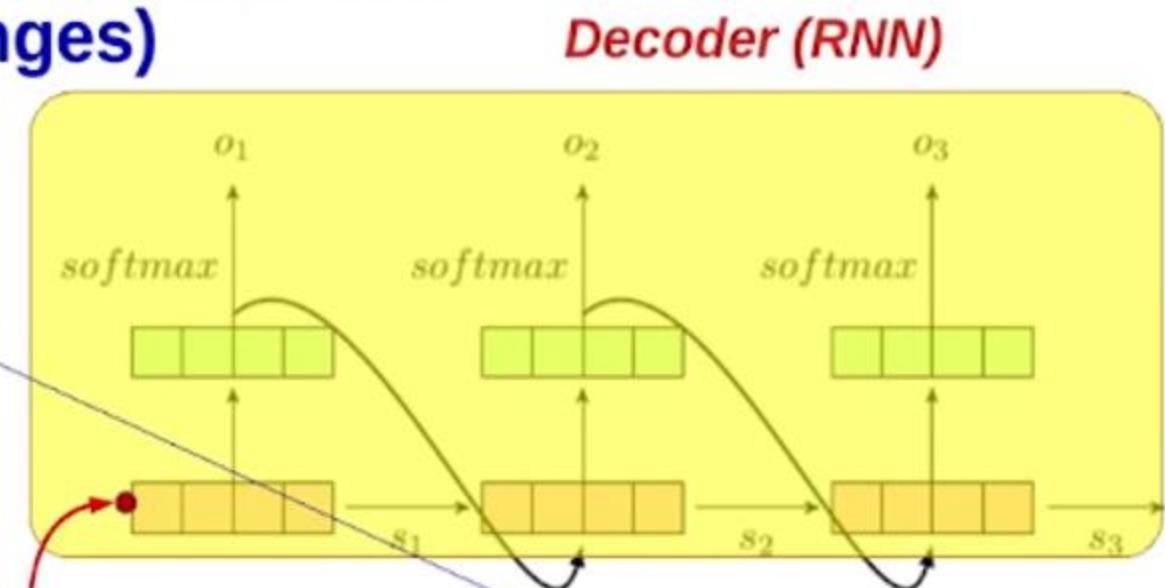
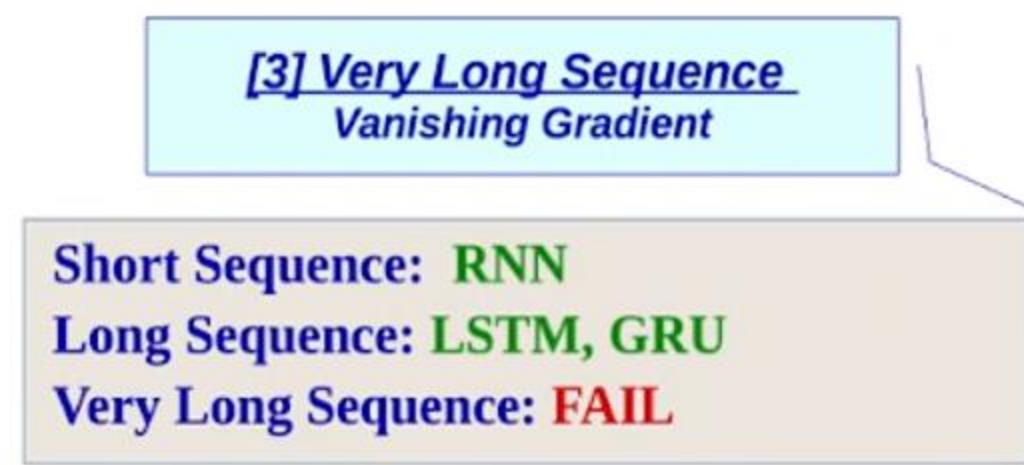
Encoder (RNN)

Context Vector (of Whole Inputs)

Traditional RNN Seq2Seq model (Challenges)



Traditional RNN Seq2Seq model (Challenges)



RNN Seq2Seq model [with Attention]

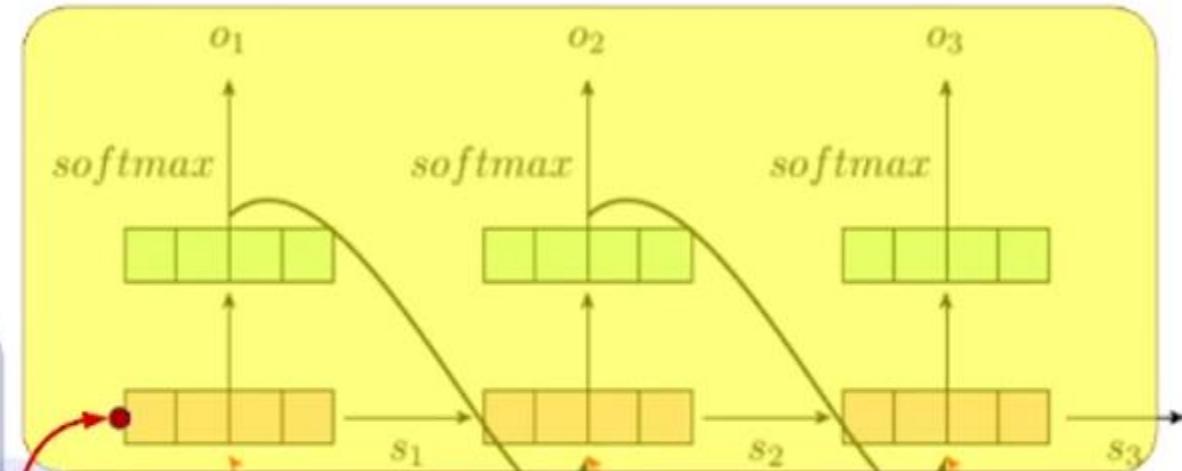
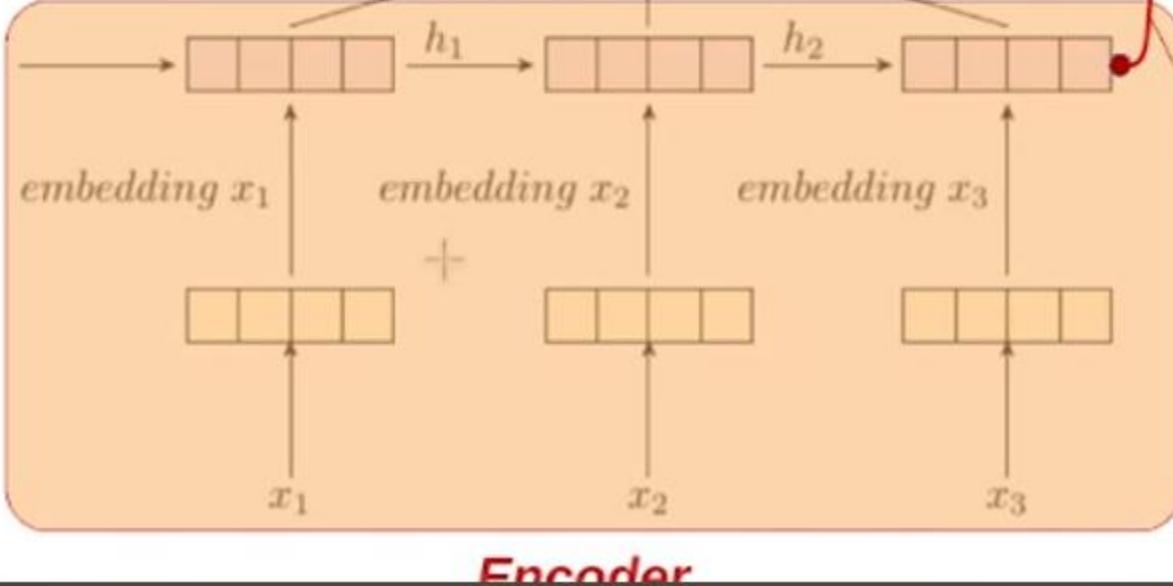
[To overcome Bottleneck challenge]

Decoder

weighted sum of the hidden states

Attention

attention



Attention

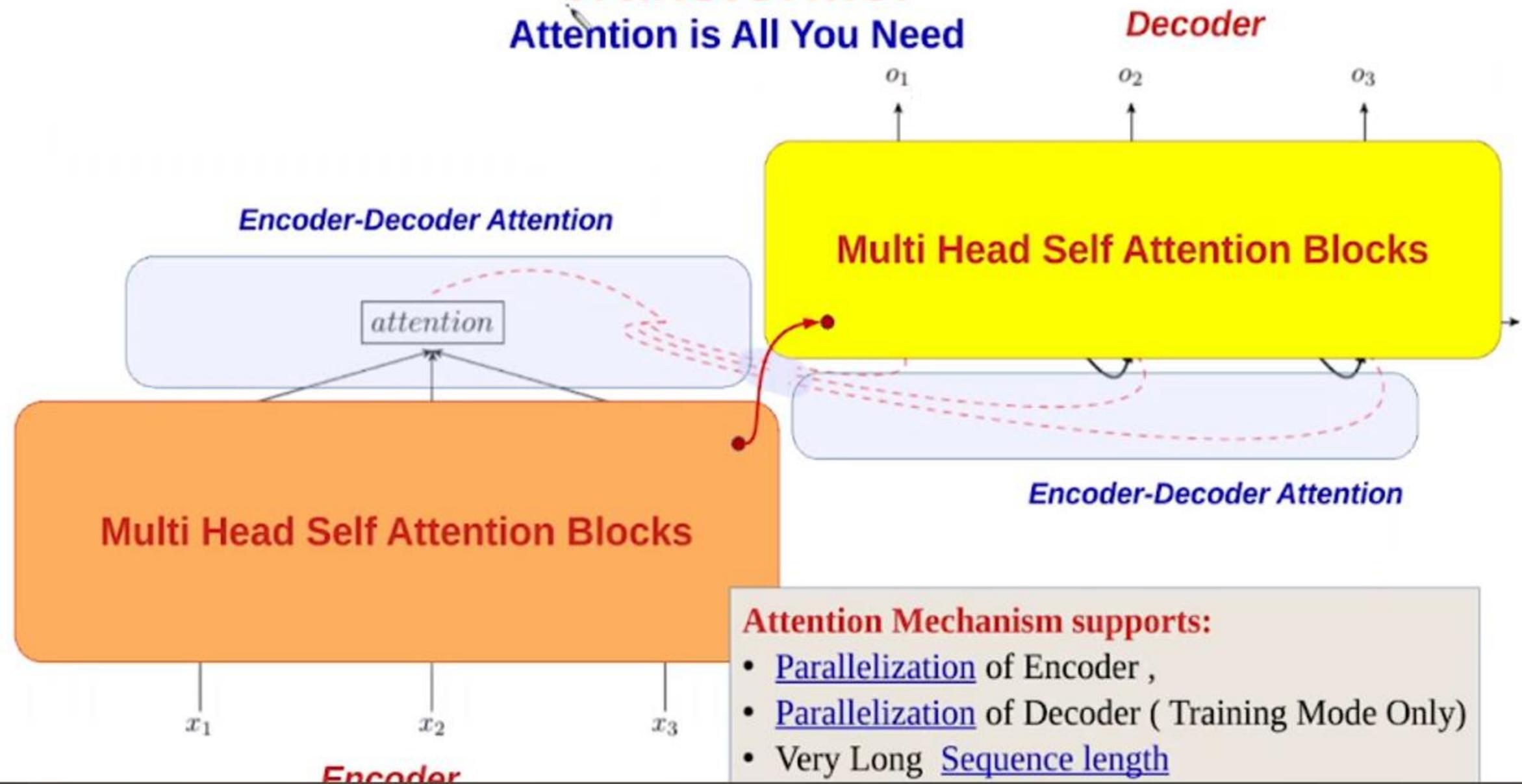
Context Vector (of Whole Inputs)

Decoder Utilizes:

- Context Vector
- Weighted sum of hidden states

Transformer

Attention is All You Need



Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez*[†]

University of Toronto

aidan@cs.toronto.edu

Lukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin*[‡]

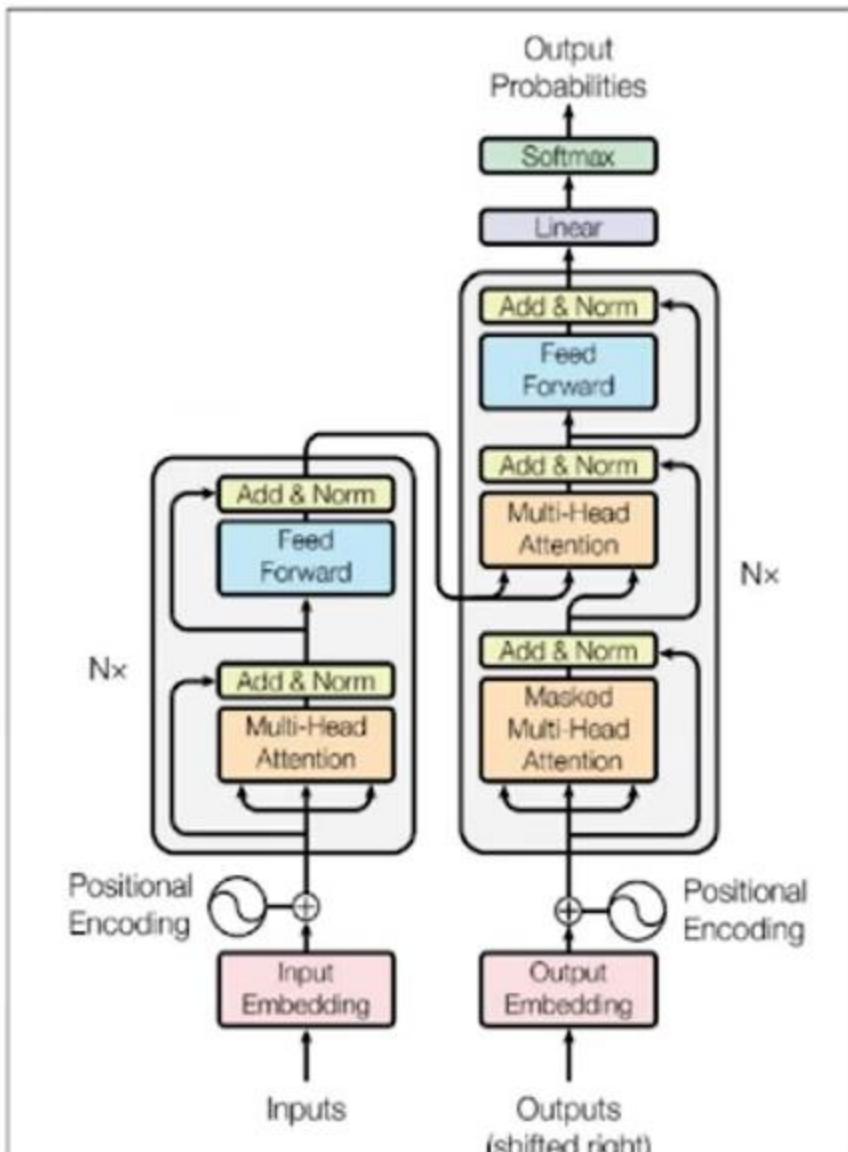
illia.polosukhin@gmail.com

Abstract

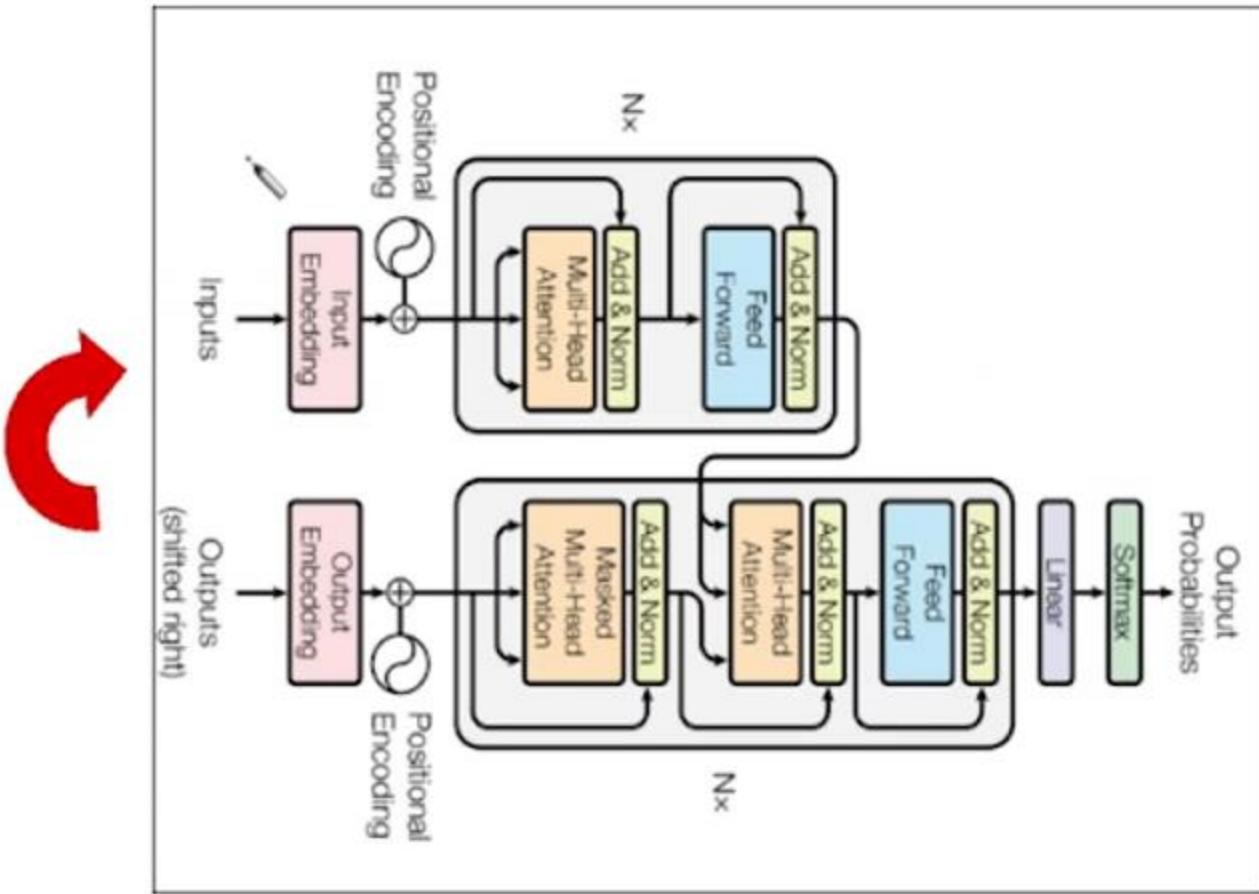
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the

Transformer Building Blocks

Encoder – Decoder Transformer

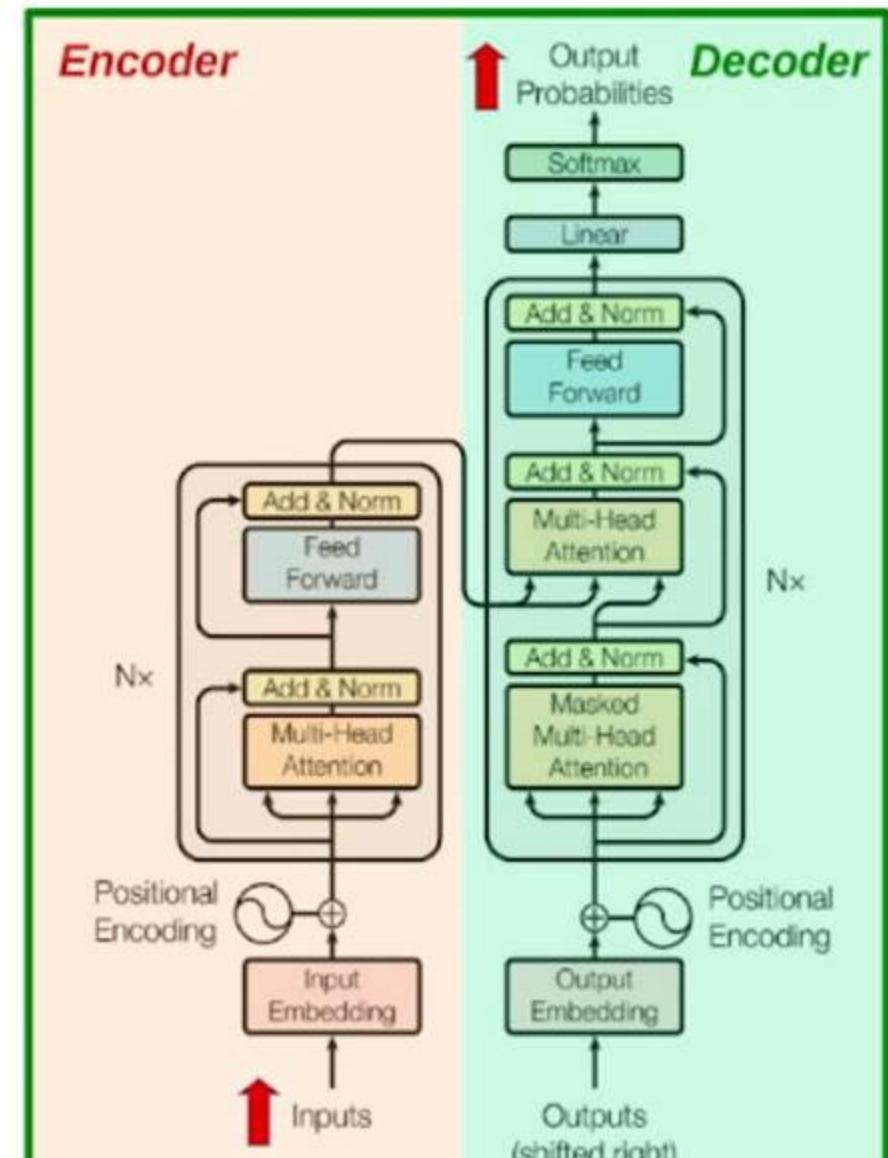
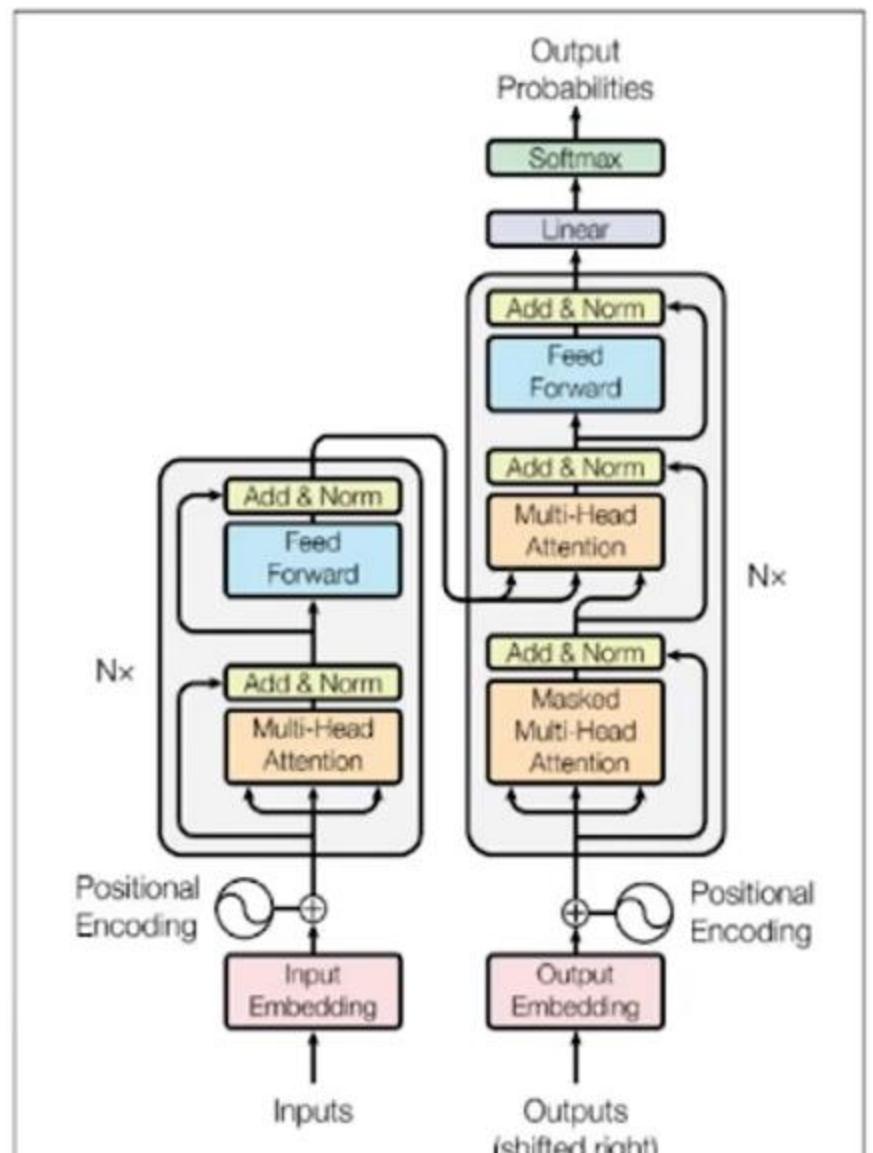


Orientation is not the Issue

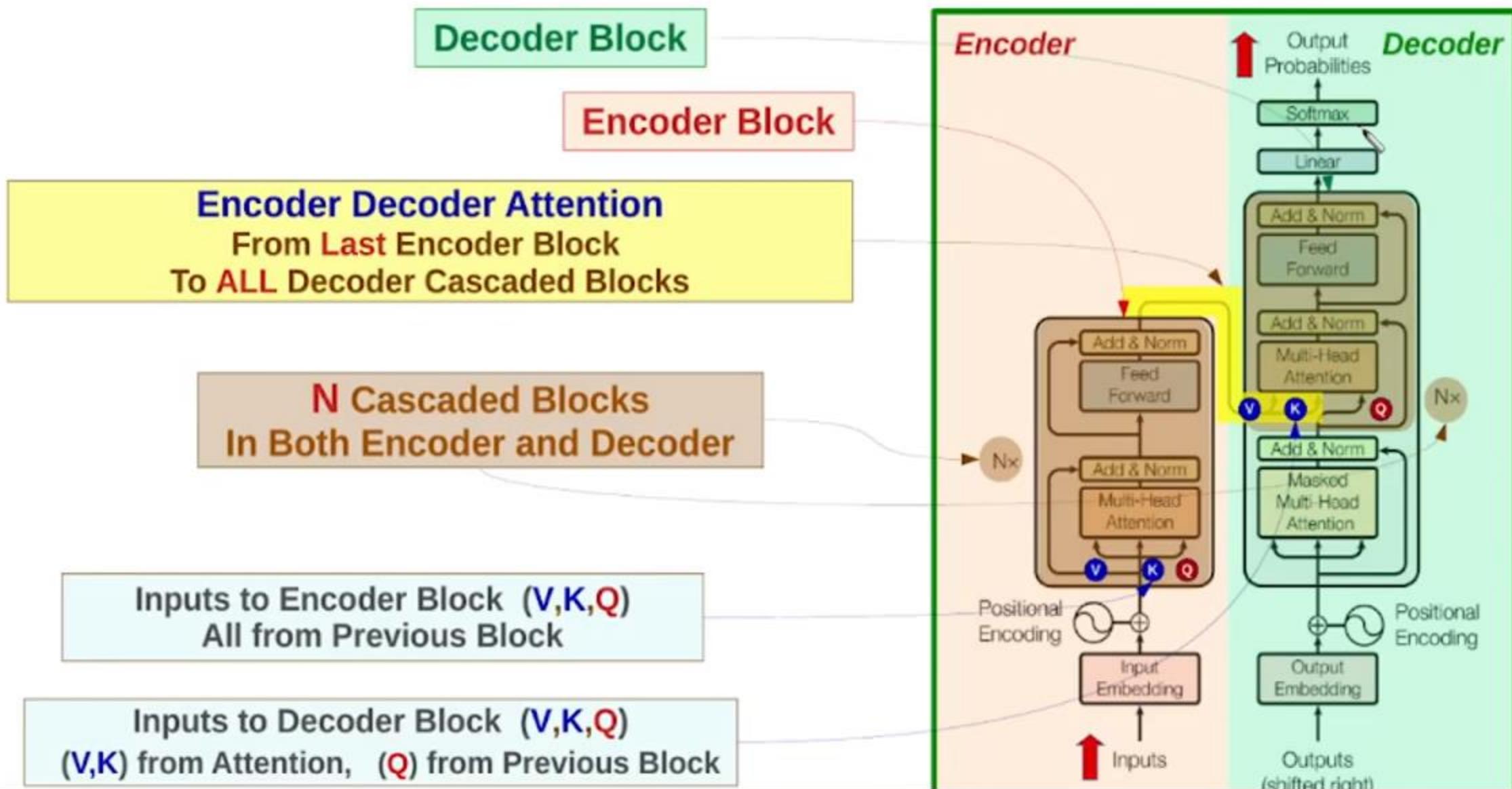


Let us focus on the Building blocks

Encoder – Decoder Transformer



Encoder – Decoder Transformer



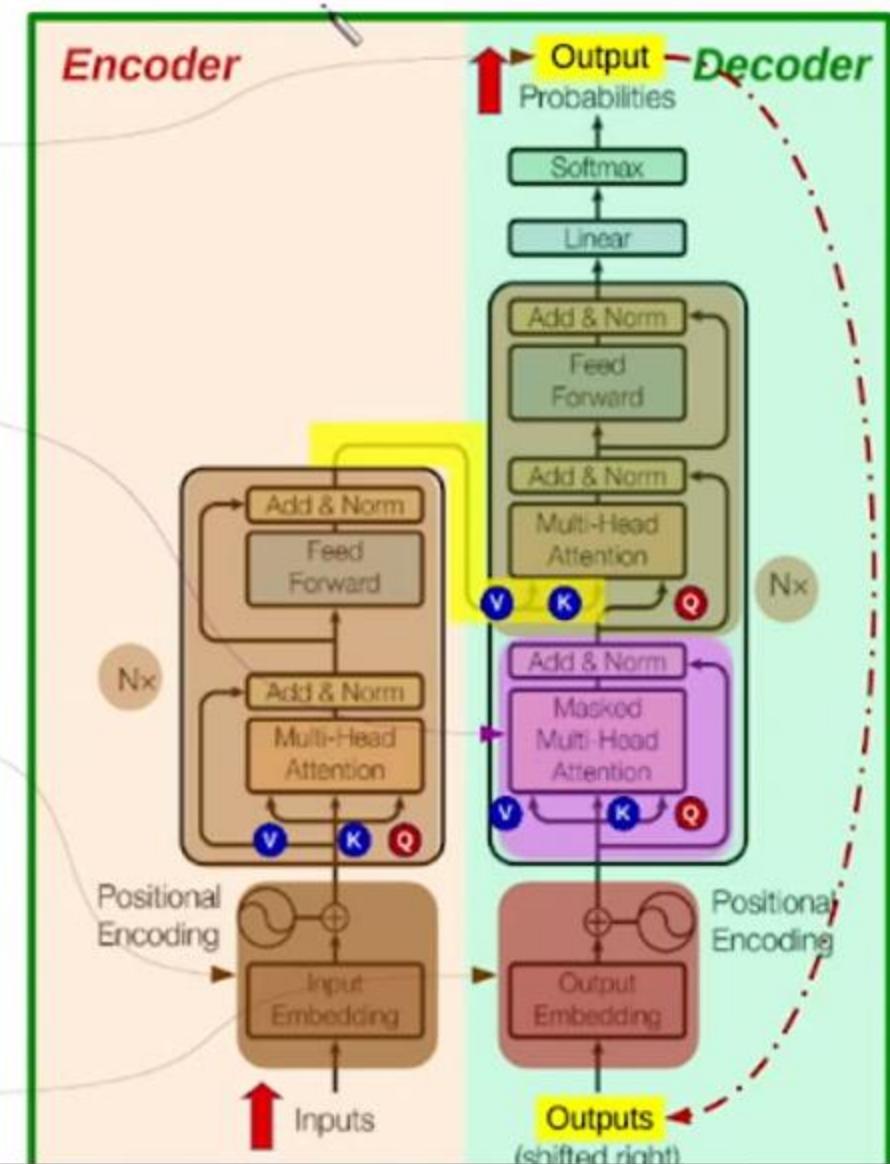
Encoder – Decoder Transformer

Decoder Output (at step “ $t-1$ ”)
Is supplied as Decoder Input (at step “ t ”)

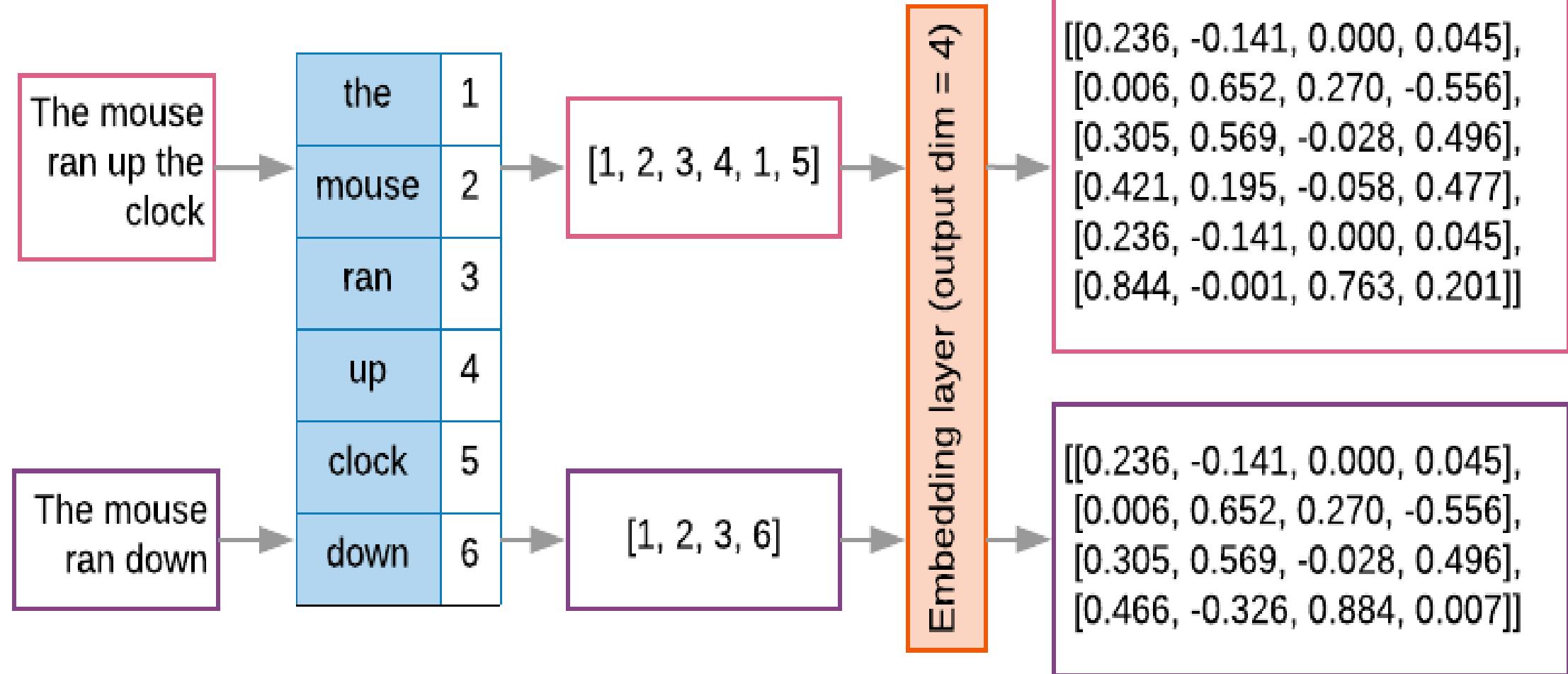
First Decoder Block
No Input from Encoder Attention

Input Positional Encoder
(Position of Input Word within the Input Sequence)

Output Positional Encoder
(Position of Output Word within the Output Sequence)



Token and embedding



$$d=4$$

good

$$\begin{array}{c} \boxed{} \\ \boxed{} \end{array} + \boxed{} = \boxed{}$$

$[d_{11}]$ $[d_{11}]$ $[d_{11}]$

$$\begin{array}{cc} good & boy \\ p=0 & p=1 \end{array}$$

$$pe(p, 2i) = \sin\left(\frac{p}{10000} 2i\pi\right)$$

$$pe(p, 2i+1) = \cos\left(\frac{p}{10000} 2i\pi\right)$$

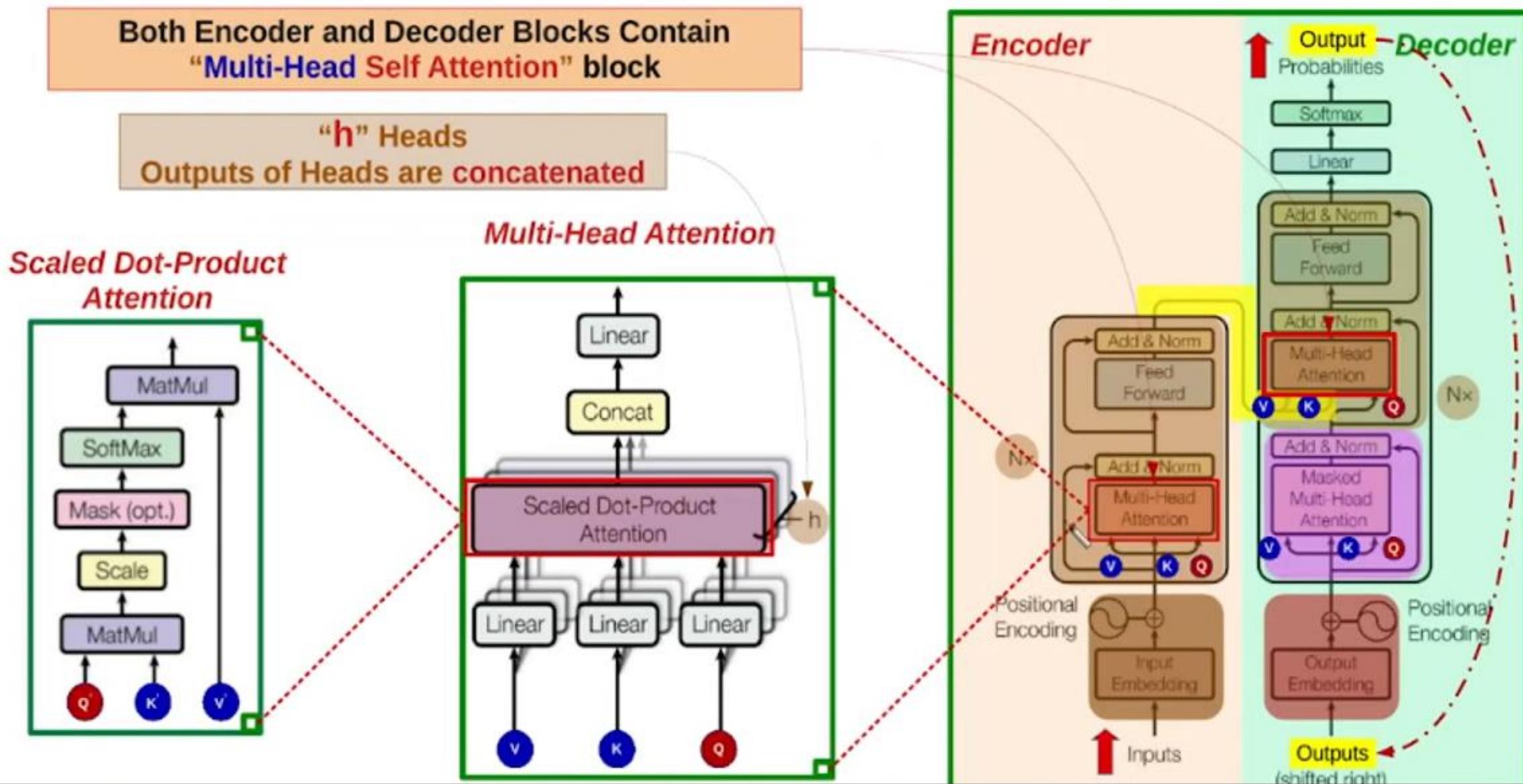
$$i = \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \end{array} \quad /2 = \begin{array}{c} b \\ 0 \\ 1 \\ 1 \end{array} = \begin{array}{c} \sin\left(\frac{0}{10000} 2 \cdot 0 \pi\right) \\ \cos\left(\frac{0}{10000} 2 \cdot 0 \pi\right) \\ \sin\left(\frac{0}{10000} 2 \cdot 1 \pi\right) \end{array}$$

$$i = \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \end{array} \quad /2 = \begin{array}{c} 0 \\ 0 \\ 1 \\ 1 \end{array} = \begin{array}{c} boy(p=1) \\ \sin\left(\frac{1}{10000} 2 \cdot 1 \pi\right) \\ \cos\left(\frac{1}{10000} 2 \cdot 1 \pi\right) \\ \sin\left(\frac{1}{10000} 2 \cdot 1 \pi\right) \end{array}$$

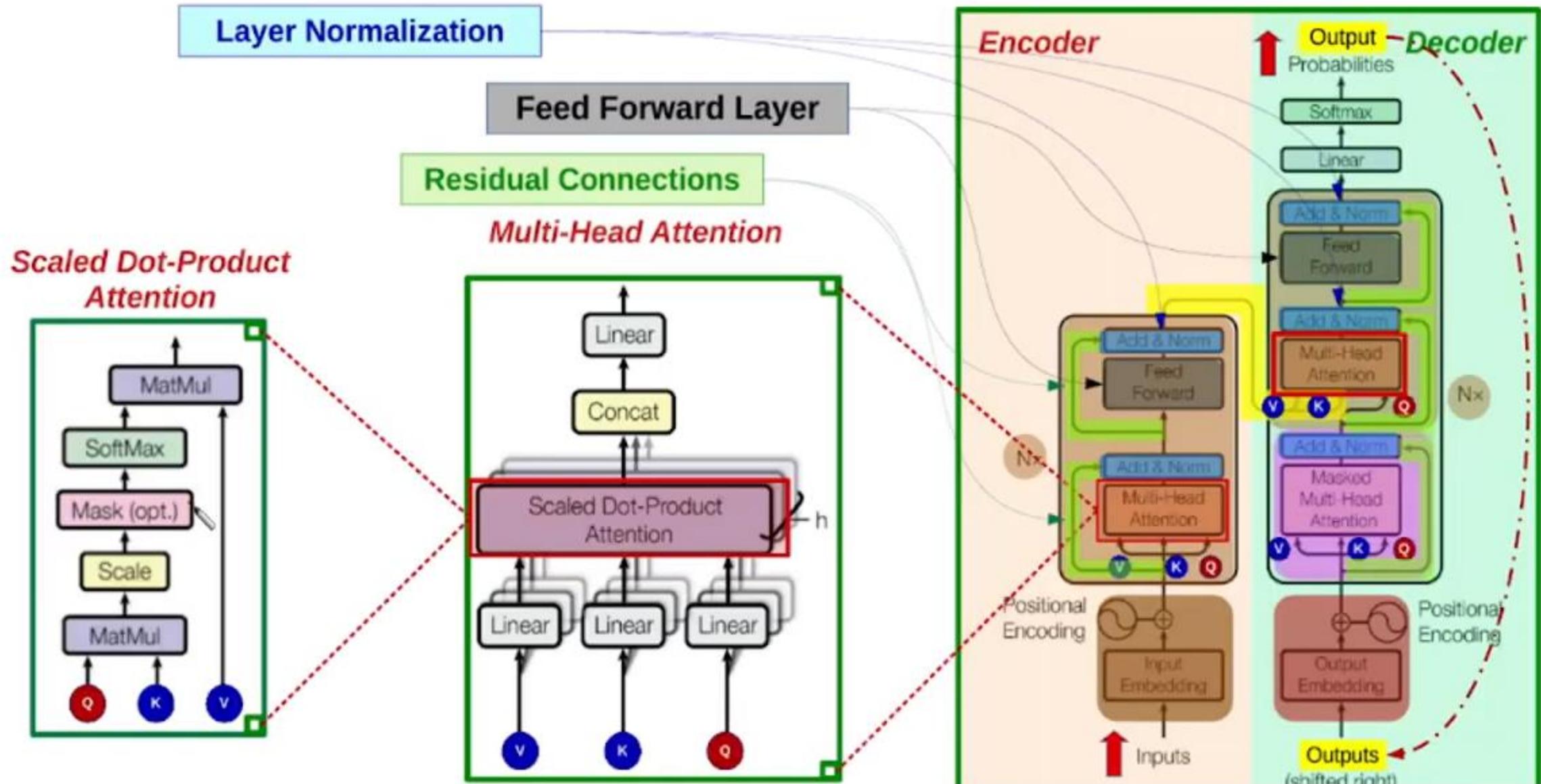
$$\begin{array}{c} \boxed{} \\ \boxed{} \end{array} + \boxed{} = \boxed{}$$

$[d_{11}]$ $[d_{11}]$ $[d_{11}]$

Encoder – Decoder Transformer

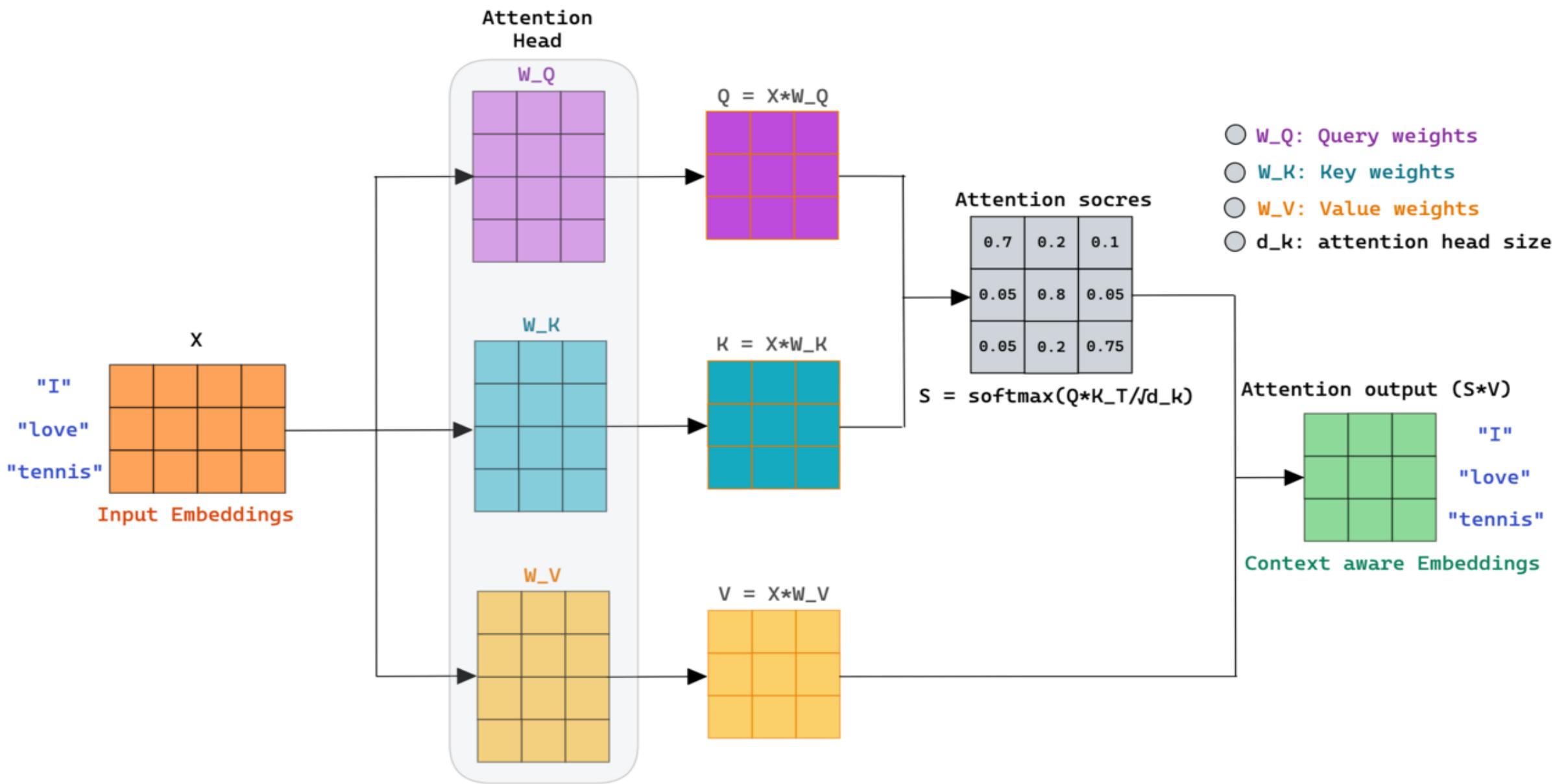


Encoder – Decoder Transformer

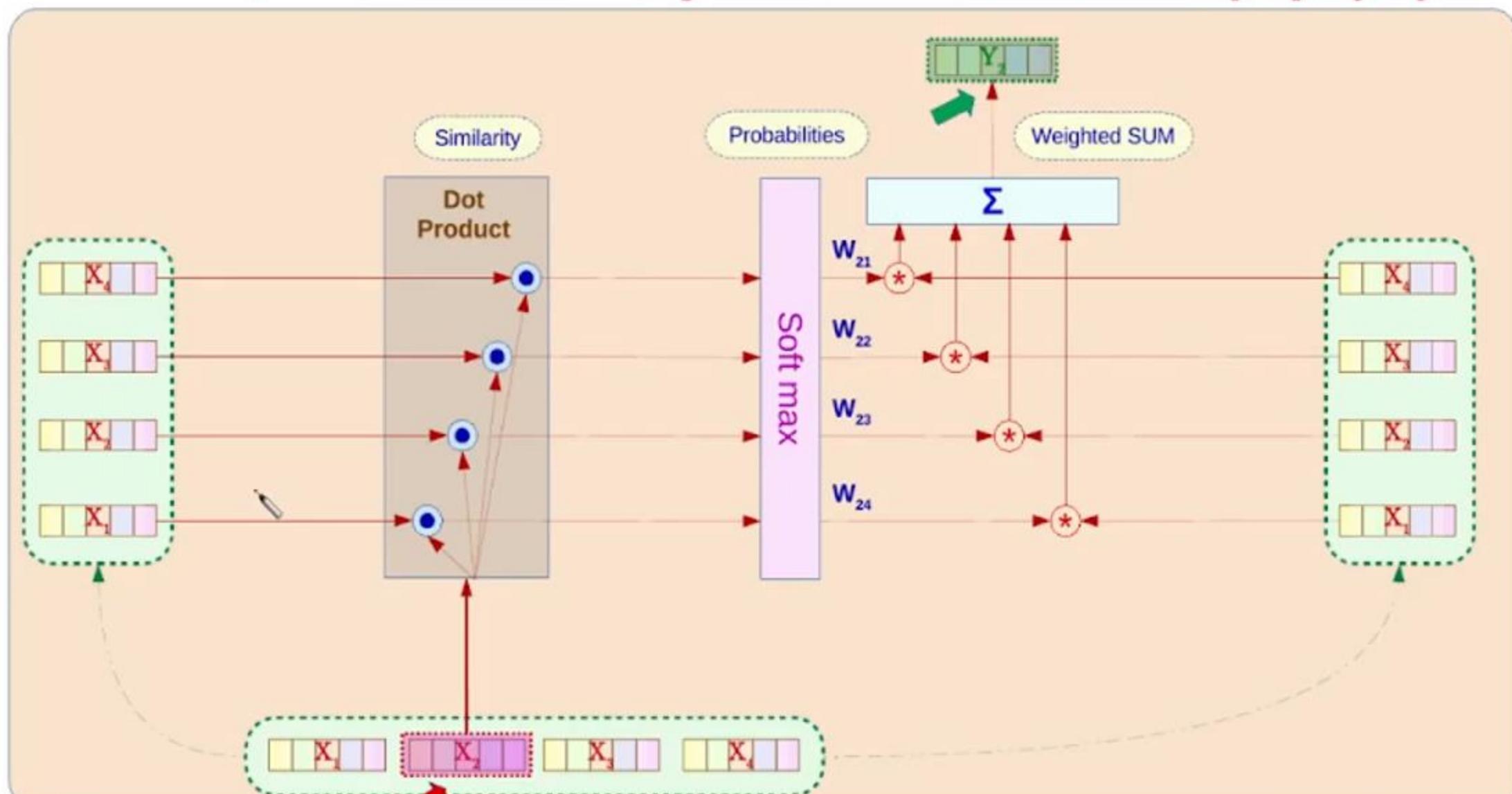


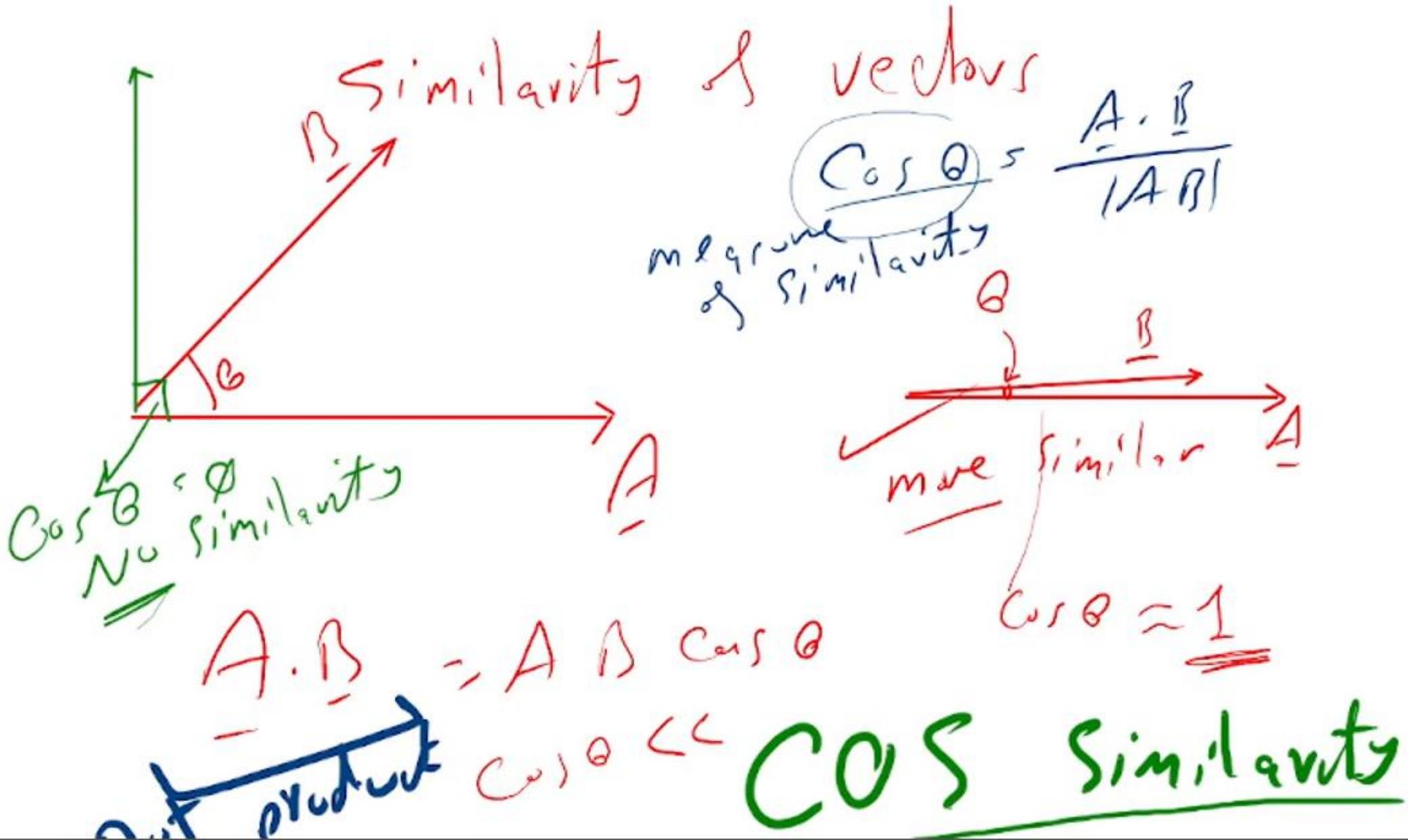
Self Attention Mechanism

Self-attention



Self Attention Mechanism (Fundamental Operation)
(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)





$X_2 \rightarrow Y_2 \Rightarrow$ related to context words
(around X_2)

$$Y_2 = \underbrace{w_1 X_1 + w_2 X_2 + w_3 X_3 + w_4 X_4}_{\text{sum}} \text{ weighted}$$

similarity of $X_1, X_2 \rightsquigarrow w_1$

" $X_3, X_2 \rightsquigarrow w_3$

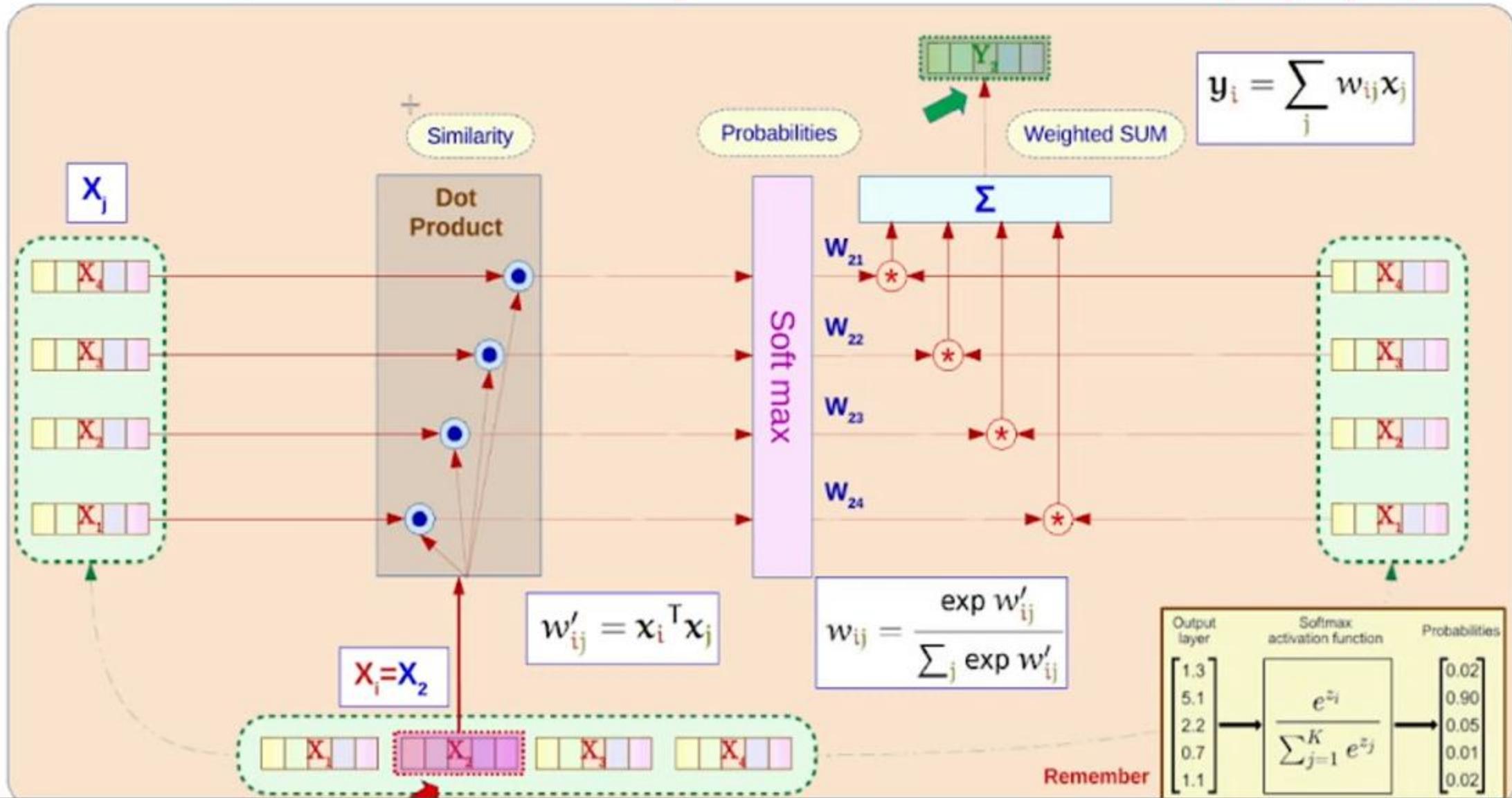
" $X_4, X_2 \rightsquigarrow w_4$

" $X_2, X_2 \rightsquigarrow \underline{\underline{w_2}}$

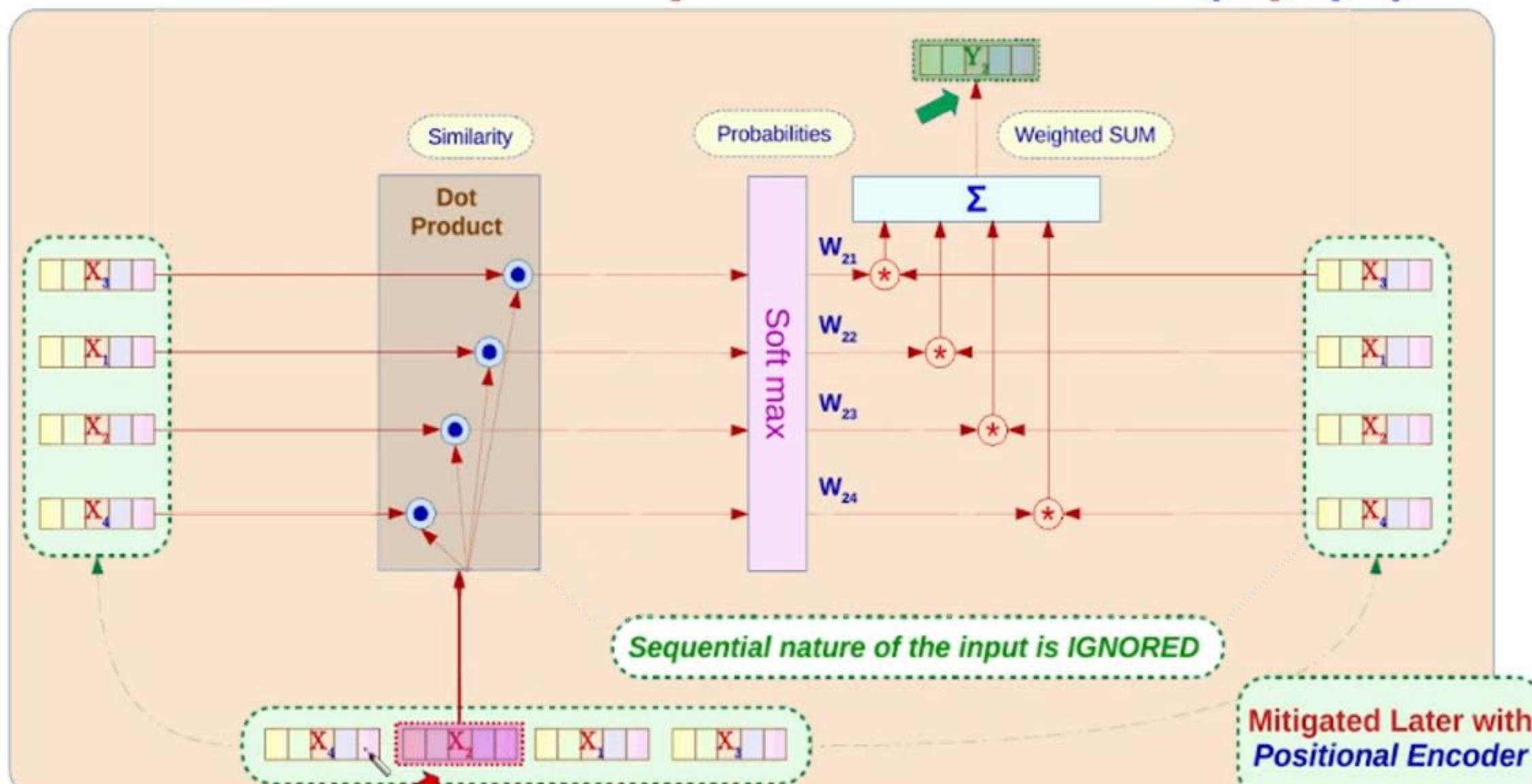
sum = $\underline{\underline{1}}$

Self Attention Mechanism (Fundamental Operation)

(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)

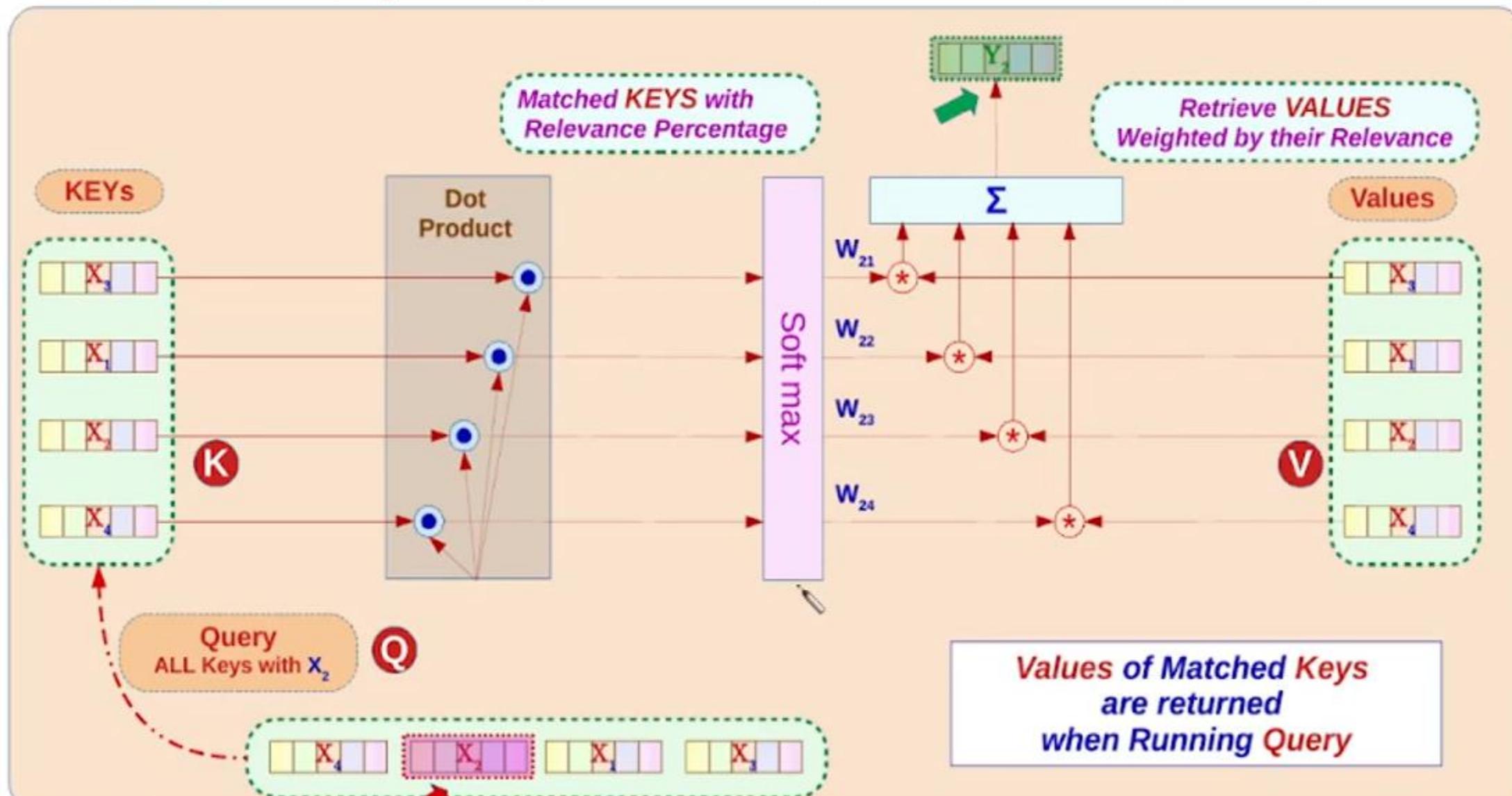


Self Attention Mechanism (Not Sensitive to Words Order)
(Self Attention of word X_2 w.r.t All words in Sentence (X_4, X_2, X_1, X_3)



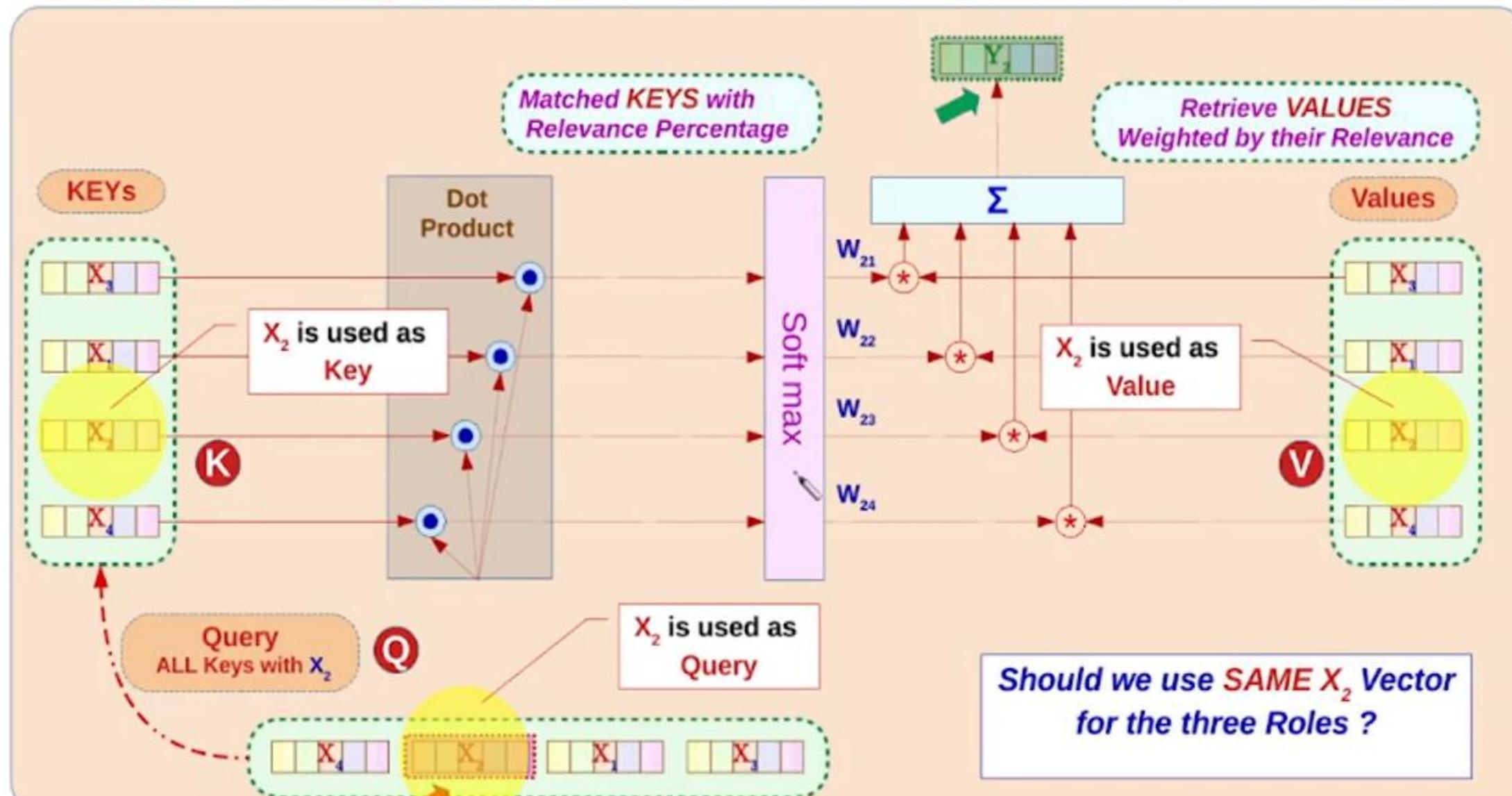
Self Attention Mechanism (as Query and {Key-Value} Dictionary)

(Query using X_2 All Keys in Dictionary and retrieve Corresponding Values



Self Attention Mechanism (Input Word plays Different Roles)

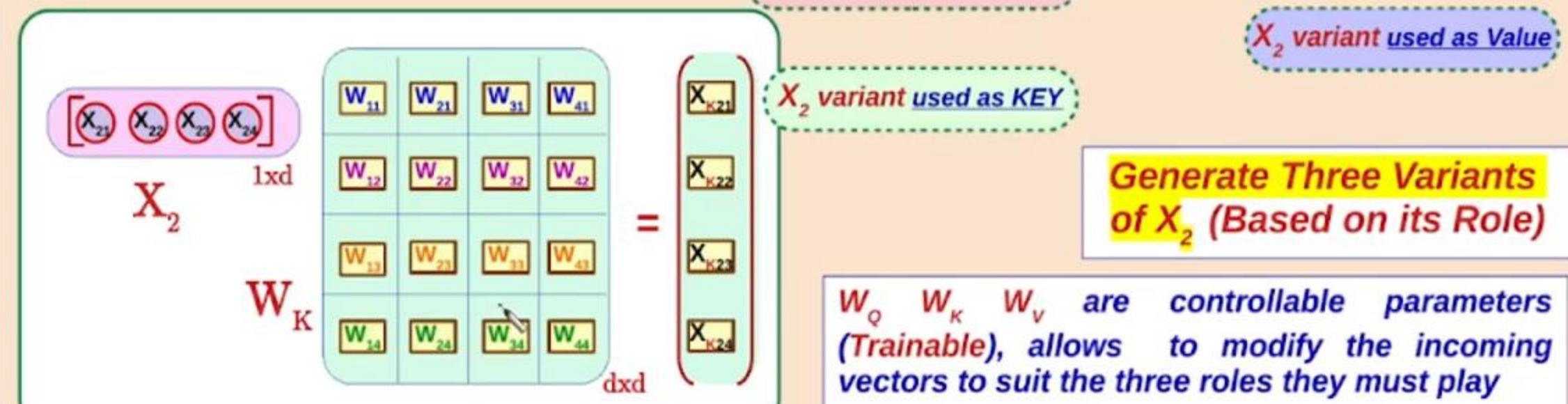
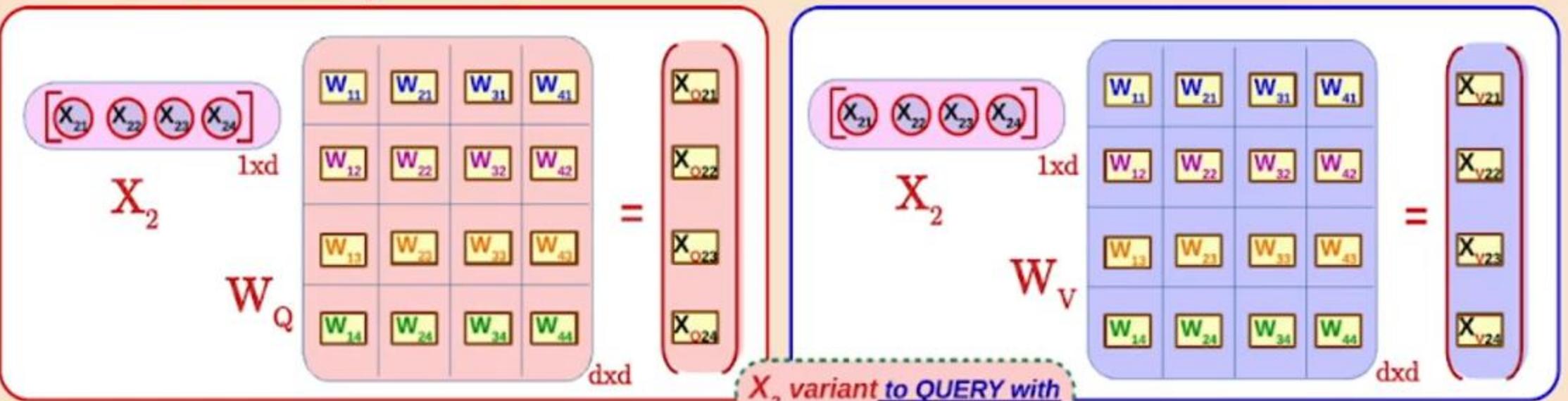
(Query using X_2 , All Keys in Dictionary and retrieve Corresponding Values



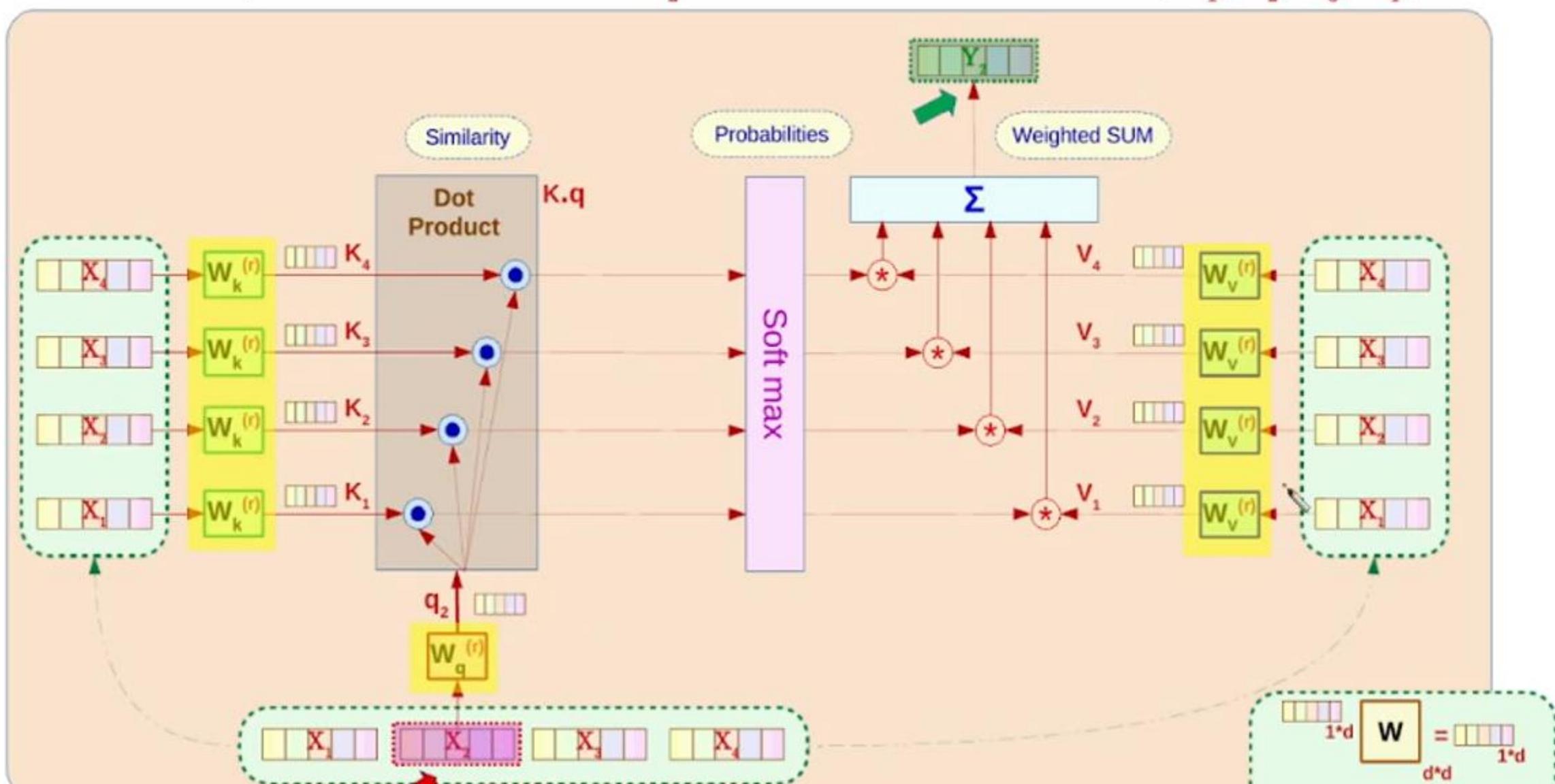
Self Attention Mechanism (Linear Transformation of Inputs)

[26]

(Query using X_2 All Keys in Dictionary and retrieve Corresponding Values



Self Attention Mechanism (Linear Transformation of Inputs)
(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)



Self Attention Mechanism

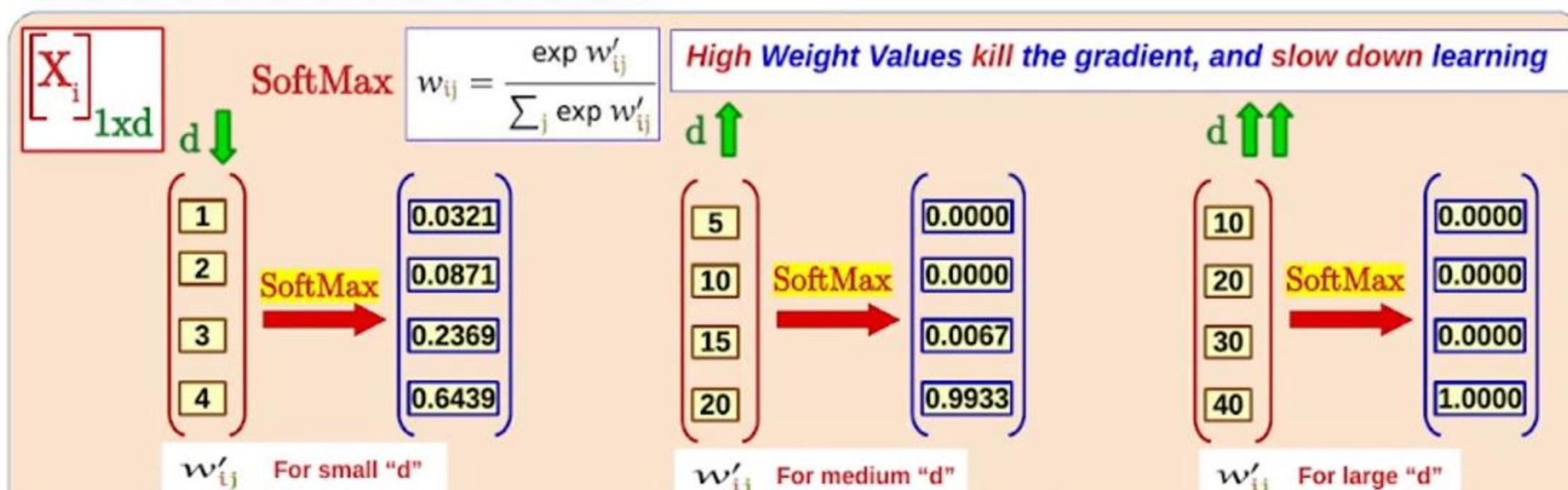
(Embedding Vector Dimension and Softmax Sensitivity)

$$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

$\left[\mathbf{X}_i\right]_{1 \times d}$

$d \uparrow$ $\uparrow w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$

*As Dimension of Embedding "d" Increases
the Calculated weight value "w" increases*



Self Attention Mechanism

(Embedding Vector Dimension and Softmax Sensitivity)

$$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$$

$$\left[\mathbf{X}_i \right]_{1 \times d}$$

$d \uparrow$

$w'_{ij} = \mathbf{x}_i^T \mathbf{x}_j$

As Dimension of Embedding " d " Increases
the Calculated weight value " w' " increases

Scale down "Weights" Values
before applying SoftMax

Scaling Factor should be related to
Embedding Dimension " d "

$$\left[\mathbf{X}_i \right]_{1 \times d}$$

SoftMax

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

1
2
3
4

SoftMax

0.0321
0.0871
0.2369
0.6439

w'_{ij} For small " d "

High Weight Values kill the gradient, and slow down learning

$d \uparrow$

5
10
15
20

SoftMax

0.0000
0.0000
0.0067
0.9933

w'_{ij} For medium " d "

$d \uparrow\uparrow$

10
20
30
40

SoftMax

0.0000
0.0000
0.0000
1.0000

w'_{ij} For large " d "

Self Attention Mechanism

(Scaling based on Square Root of Embedding Vector Dimension)

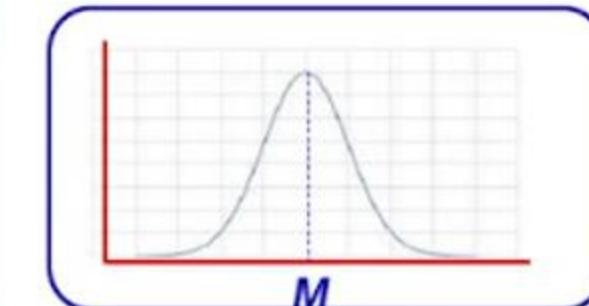
$$\begin{array}{c} V_1 \quad V_2 \\ \left[\begin{array}{c} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \\ X_7 \\ X_8 \\ X_9 \\ X_{10} \\ X_{11} \\ X_{12} \end{array} \right] \quad \left[\begin{array}{c} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \end{array} \right] \\ 1 \times d_1 \qquad \qquad \qquad 1 \times d_2 \end{array}$$

Average Values of elements of $V_1 = M$

Average Values of elements of $V_2 = M$

Vector Length of V_1 (all Ms) = $V_1^T V_1 = M\sqrt{d_1}$

Vector Length of V_2 (all Ms) = $V_2^T V_2 = M\sqrt{d_2}$



\sqrt{d} is a good choice for scaling

$4\sqrt{d}$ Some times is used

$$\text{SoftMax } \left[\begin{array}{c} w_{ij} \end{array} \right]$$

$$\text{SoftMax } \left[\frac{w_{ij}}{\sqrt{d}} \right]$$

$$d=100, \sqrt{d}=10$$

Remember:
 "d" is
Dimension of Embedding
 Not
 length of input Sequence

$$\begin{pmatrix} 10 \\ 20 \\ 30 \\ 40 \end{pmatrix} \xrightarrow{\text{SoftMax}} \begin{pmatrix} 0.0000 \\ 0.0000 \\ 0.0000 \\ 1.0000 \end{pmatrix}$$

$$\begin{pmatrix} 10 \\ 20 \\ 30 \\ 40 \end{pmatrix} / \sqrt{d} \xrightarrow{\text{Scaled}} \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} \xrightarrow{\text{SoftMax}} \begin{pmatrix} 0.0321 \\ 0.0871 \\ 0.2369 \\ 0.6439 \end{pmatrix}$$

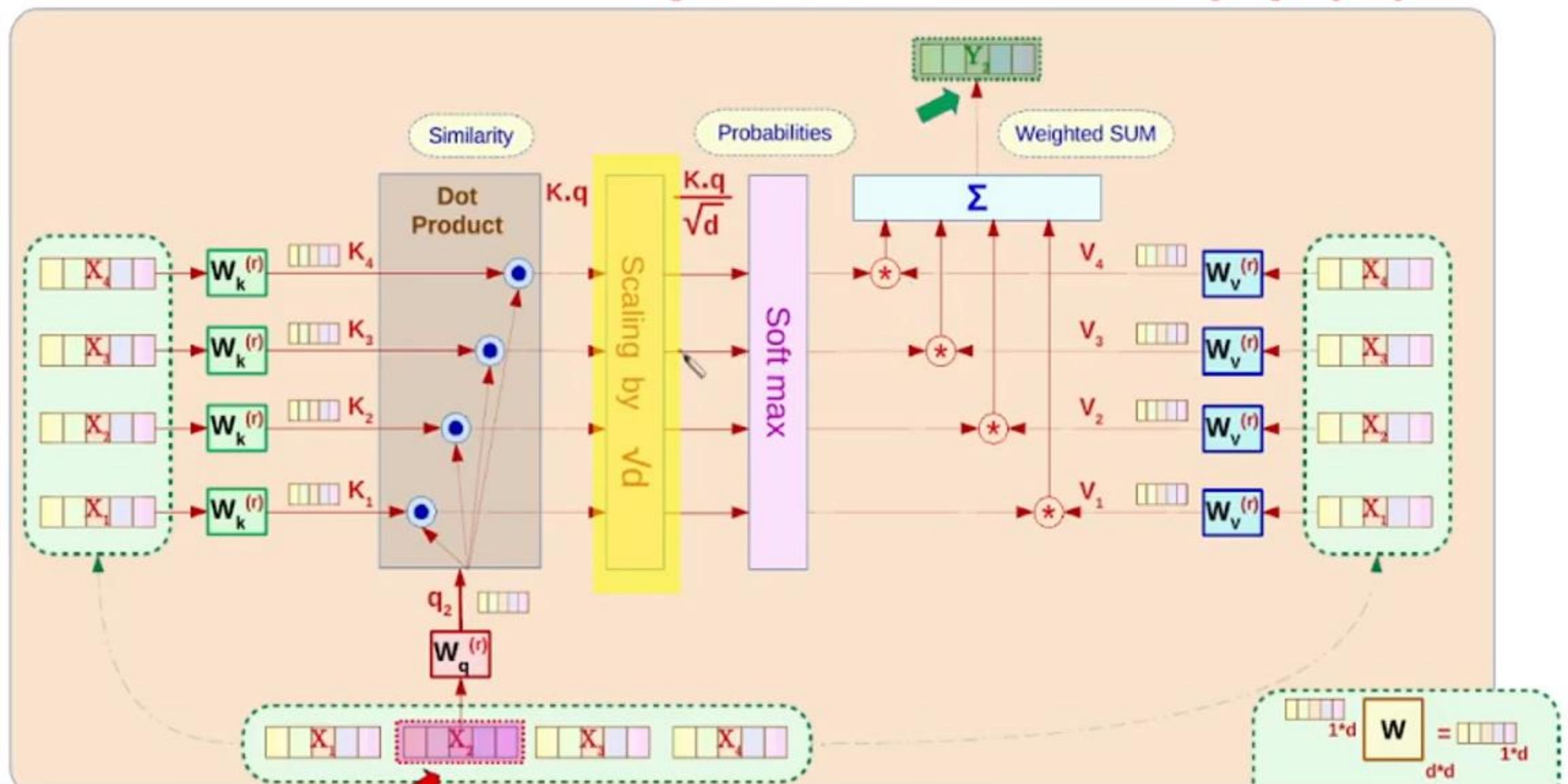
$$\begin{bmatrix} s & s & s & s \end{bmatrix} \begin{bmatrix} s \\ s \\ s \\ s \end{bmatrix} = \boxed{s^2 + s^2 + s^2 + s^2} = \boxed{4s^2}$$



$$\begin{bmatrix} s & s & s & s & s \end{bmatrix} \begin{bmatrix} s \\ s \\ s \\ s \\ s \end{bmatrix}_{100} = \boxed{s^2 + s^2 + \dots + s^2} = \boxed{100 \times s^2} = \boxed{s \sqrt{100}} \rightarrow d$$

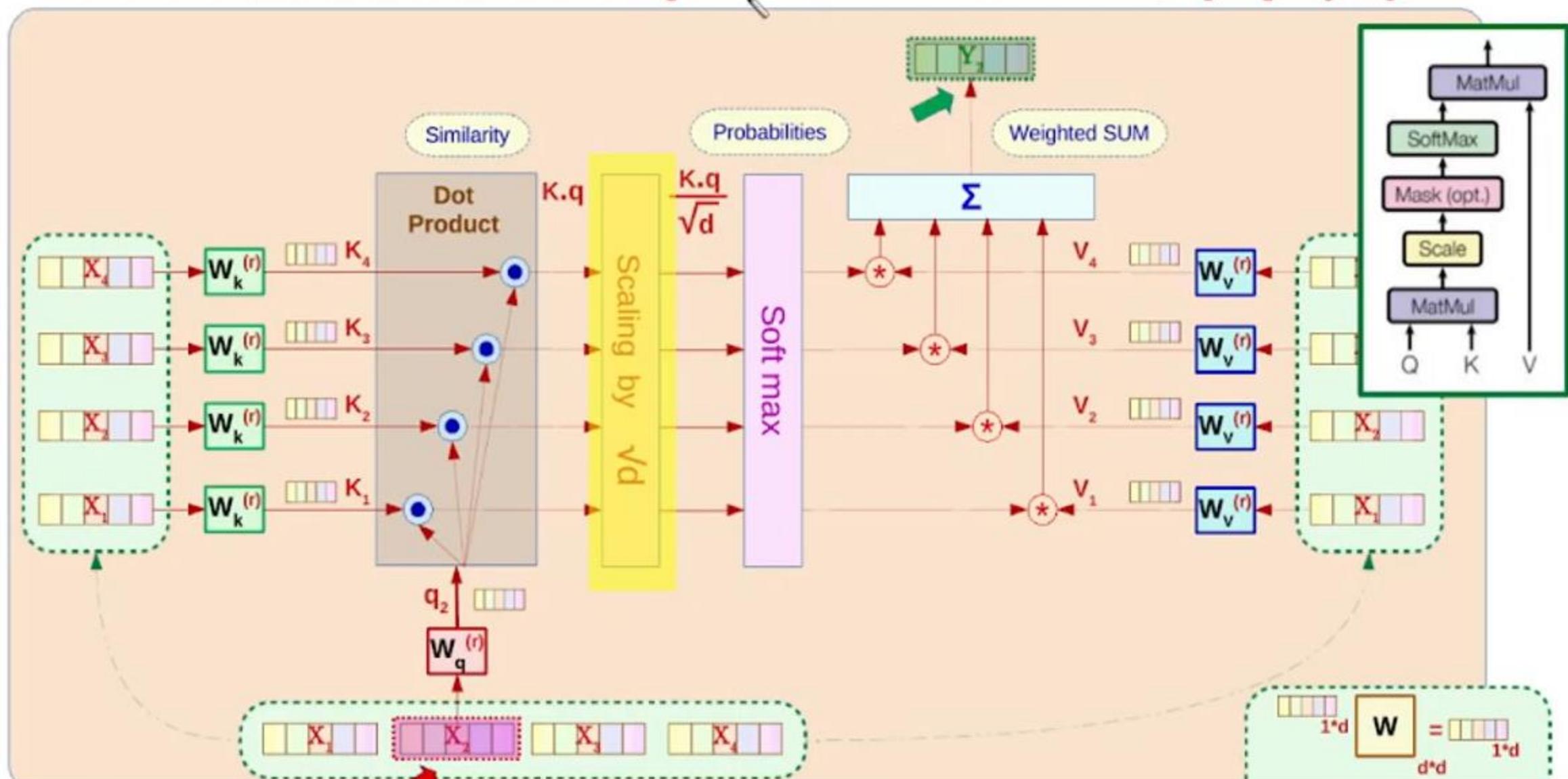
Self Attention Mechanism (Adding Scaling)

(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)



Self Attention Mechanism (Adding Scaling)

(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)



$$\begin{bmatrix} \{x_1\} \\ \{x_2\} \\ \{x_3\} \\ \{x_4\} \end{bmatrix} * \begin{bmatrix} (0) & (x_1) & (x_2) & (x_3) \end{bmatrix} = \begin{bmatrix} (0) & (0) & (0) \\ (0) & (0) & (0) \\ (0) & (0) & (0) \\ (0) & (0) & (0) \end{bmatrix}$$

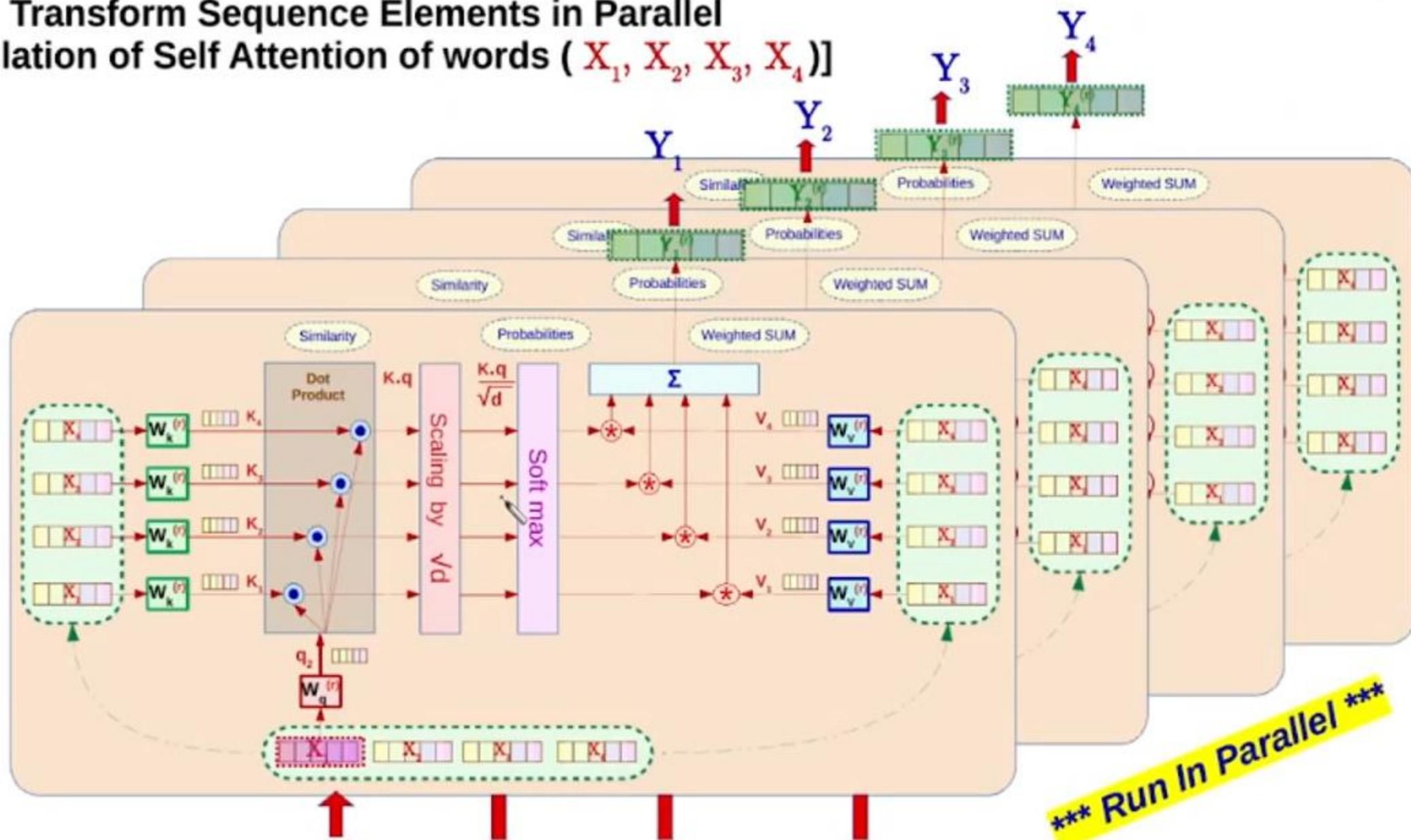
Inst matrix keys
 Mat Mul

To get similarities with other word
 of All words At one time

Parallel Processing

Transform Sequence Elements in Parallel

[Calculation of Self Attention of words (X_1, X_2, X_3, X_4)]



Multi-Head Attention

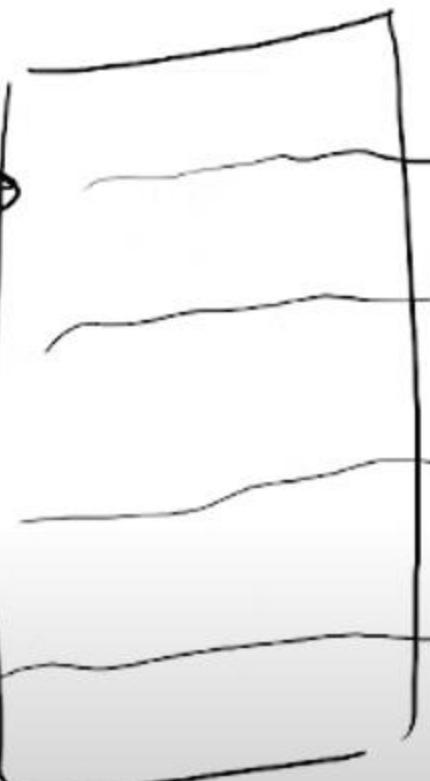


Wide Architecture

Narrow Architecture

Individual word vectors

$x_1 \rightarrow$
 $x_2 \rightarrow$
 $x_3 \rightarrow$
 $x_4 \rightarrow$



word vectors that embed the context words

$$y_i, \{x_1, x_2, x_3, x_4\}$$

y_1

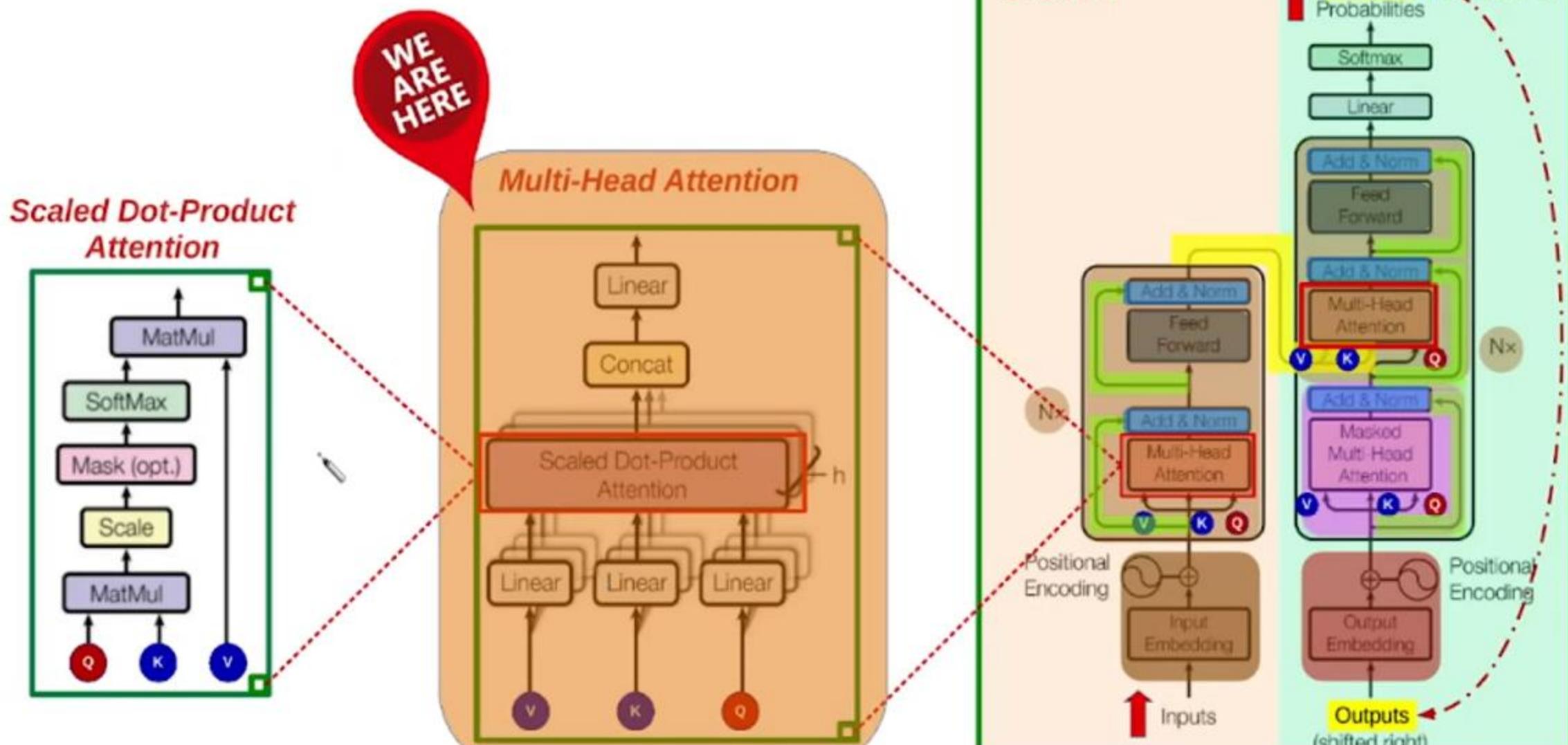
y_2

y_3

y_4

all Done in parallel

Encoder – Decoder Transformer



Multi-Head Self Attention

Combining several self attention mechanisms give the self attention greater power of discrimination

Each Head is indexed with “r”

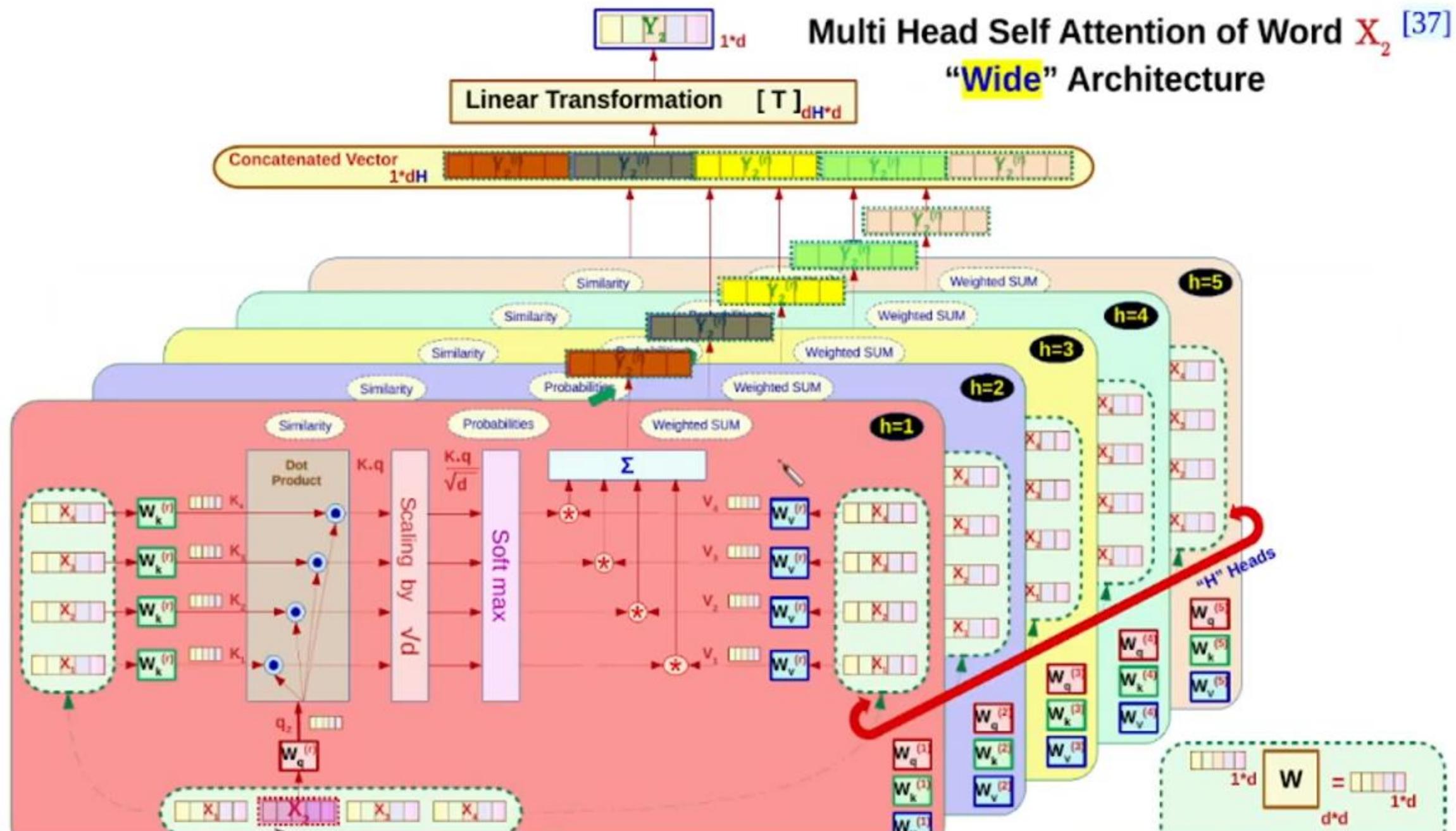
Each head has its own matrices \mathbf{W}_q^r , \mathbf{W}_k^r , \mathbf{W}_v^r .

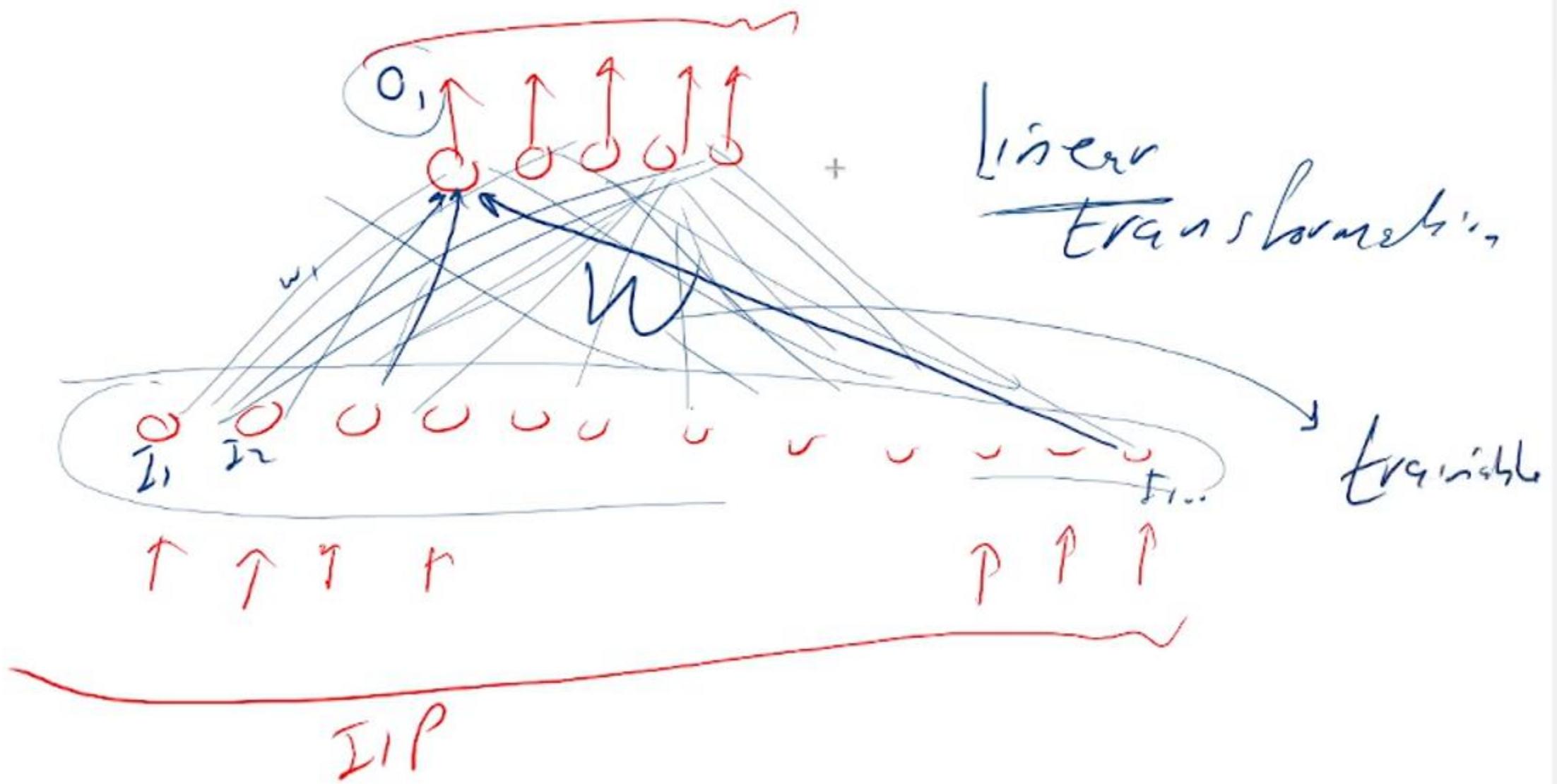
For input \mathbf{X} each attention head produces a different output vector \mathbf{Y}^r .

All Output vectors are concatenated , and pass through a linear transformation to reduce the dimension back to original \mathbf{X} dimension.

Multi Head Self Attention of Word X_2 [37]

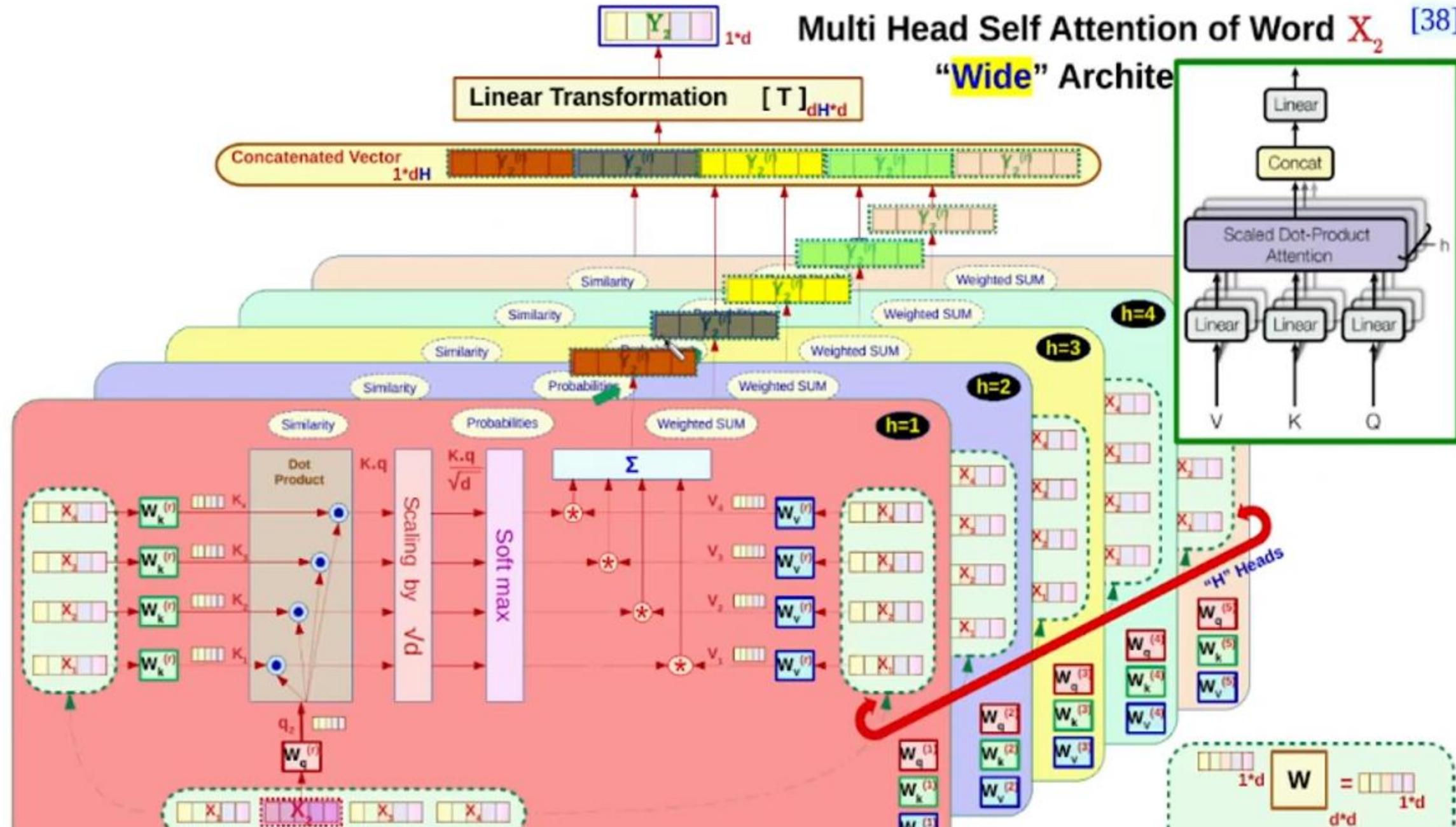
“Wide” Architecture





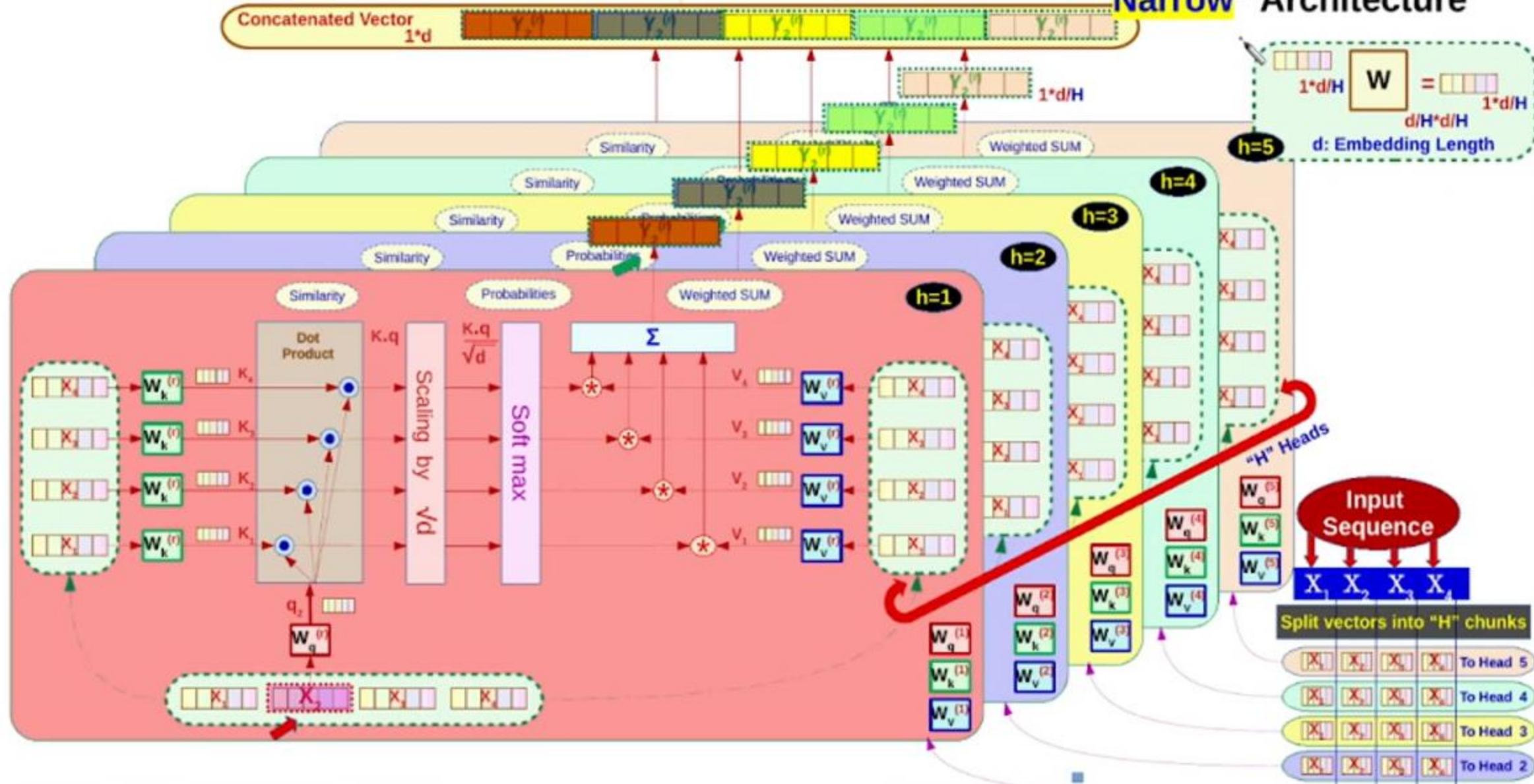
Multi Head Self Attention of Word X_2 [38]

"Wide" Archite



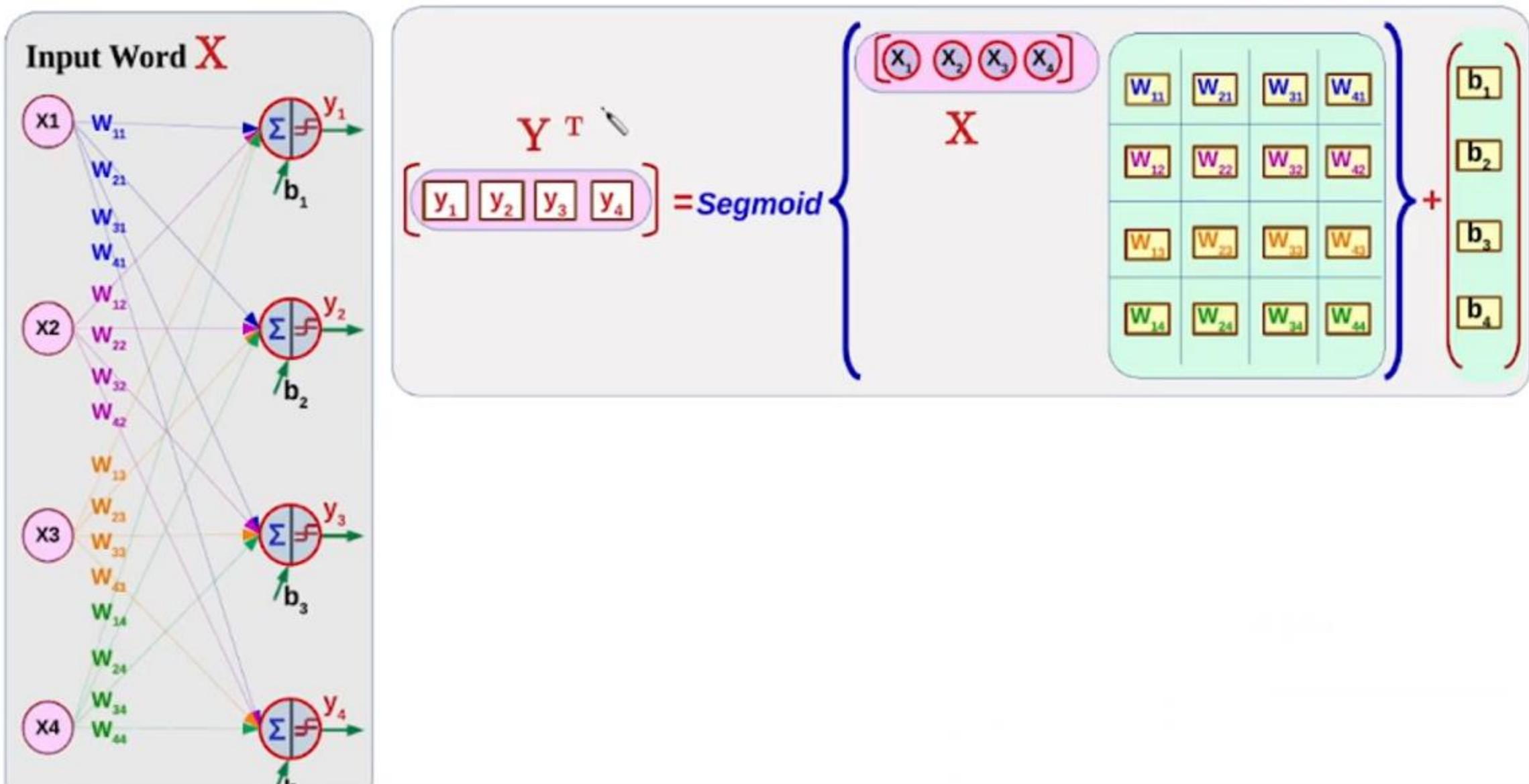
Multi Head Self Attention of Word X_2 [39]

“Narrow” Architecture



Parallelization Implementation Issues

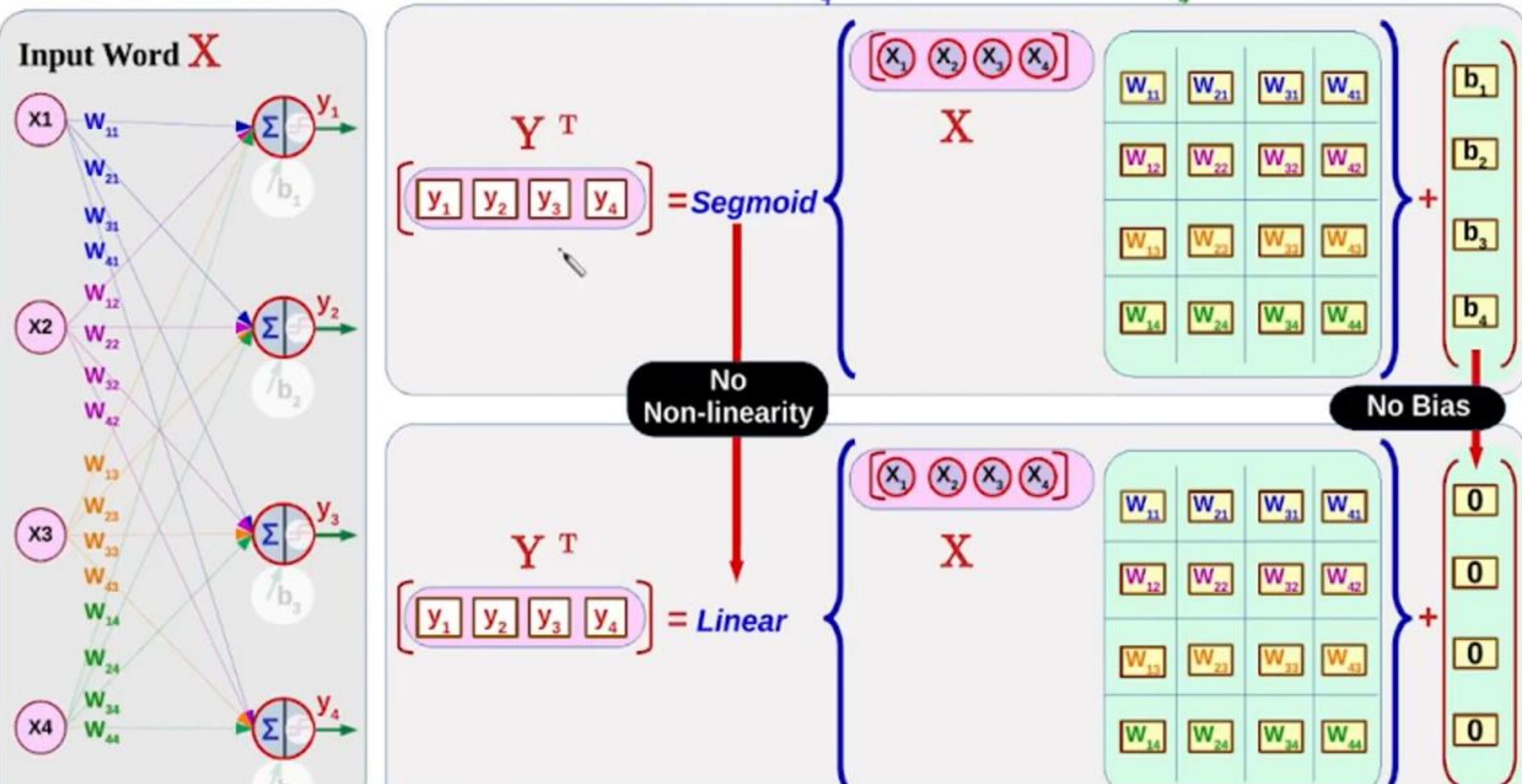
Using FF NN for Multiplication of a Vector by a Matrix



Using FF NN for Multiplication of ONE Vector by a Matrix

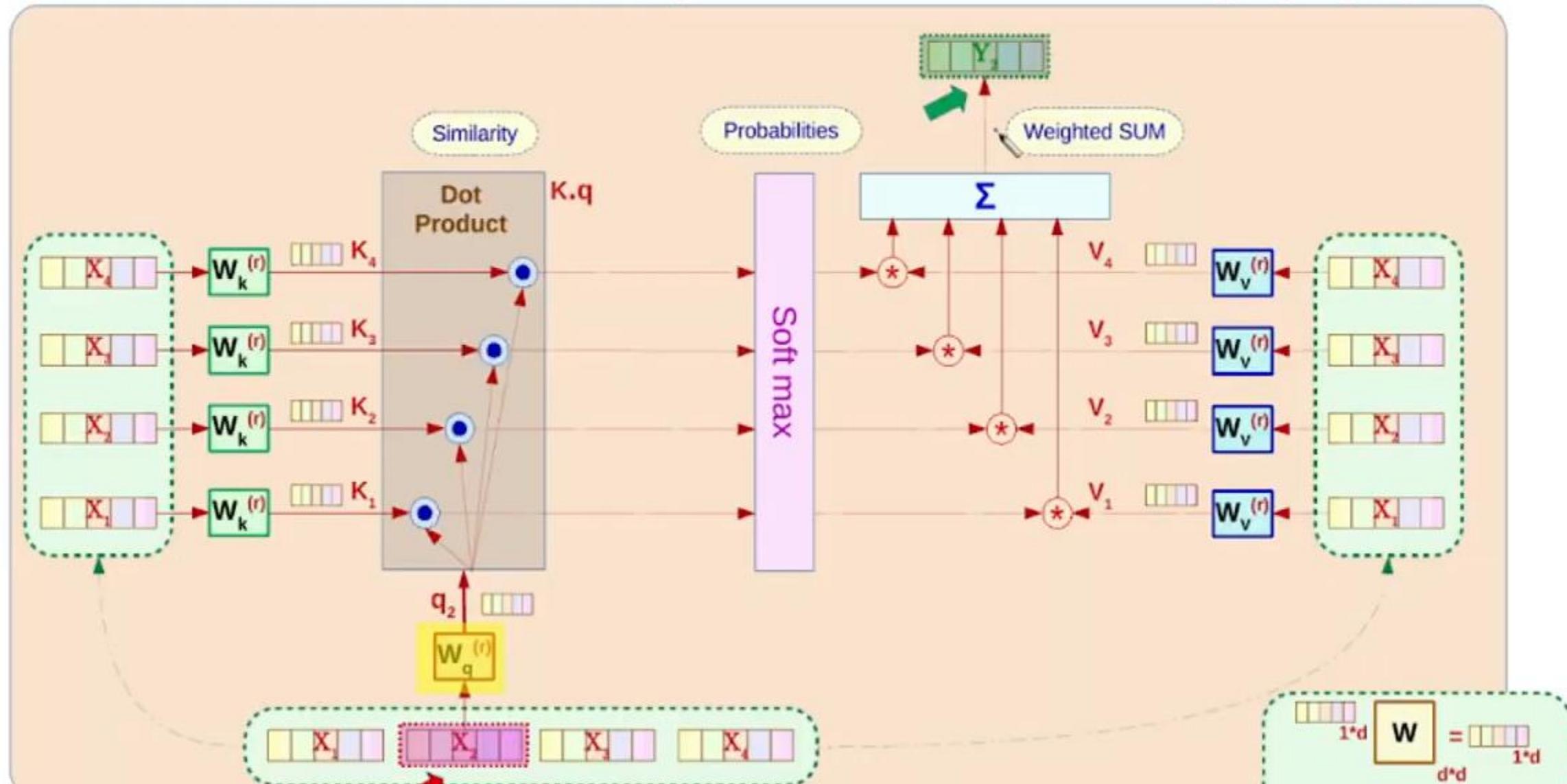
[42]

Example: Multiply Vector of Current Word X by W_q (Query) {Remember: W_q is Squared Matrix}



Self Attention Mechanism (Linear Transformation of Inputs)

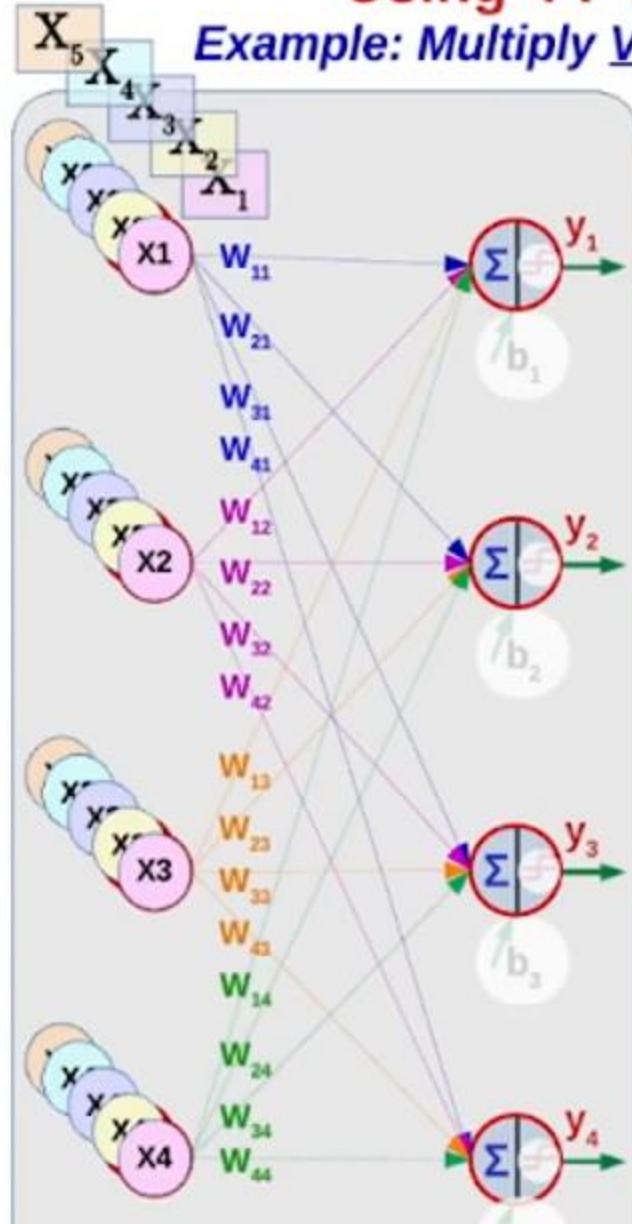
(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)



Using FF NN for Multiplication of Multiple Vectors by a Matrix

[44]

Example: Multiply Vectors of ALL Words by W_k (Key) or W_v (Value) { W_k, W_v is Squared Matrix }

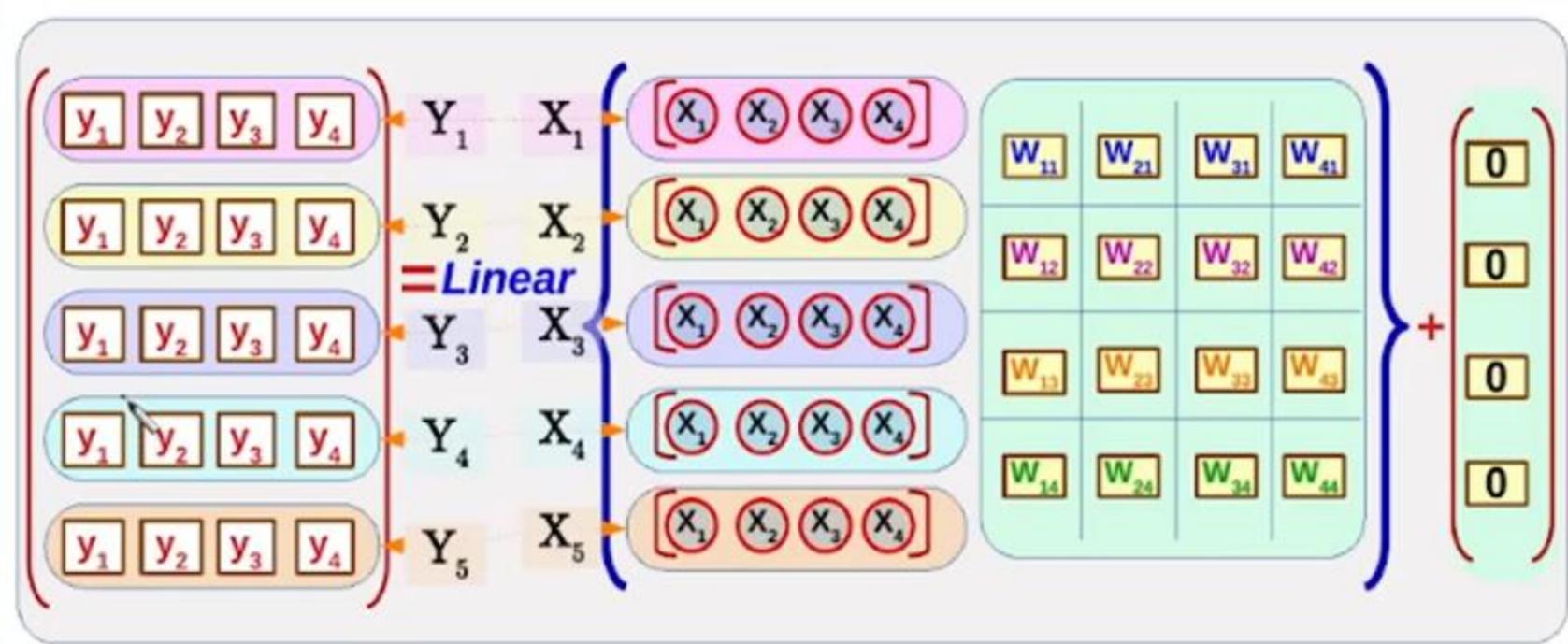


Design FF NN with the following Considerations:

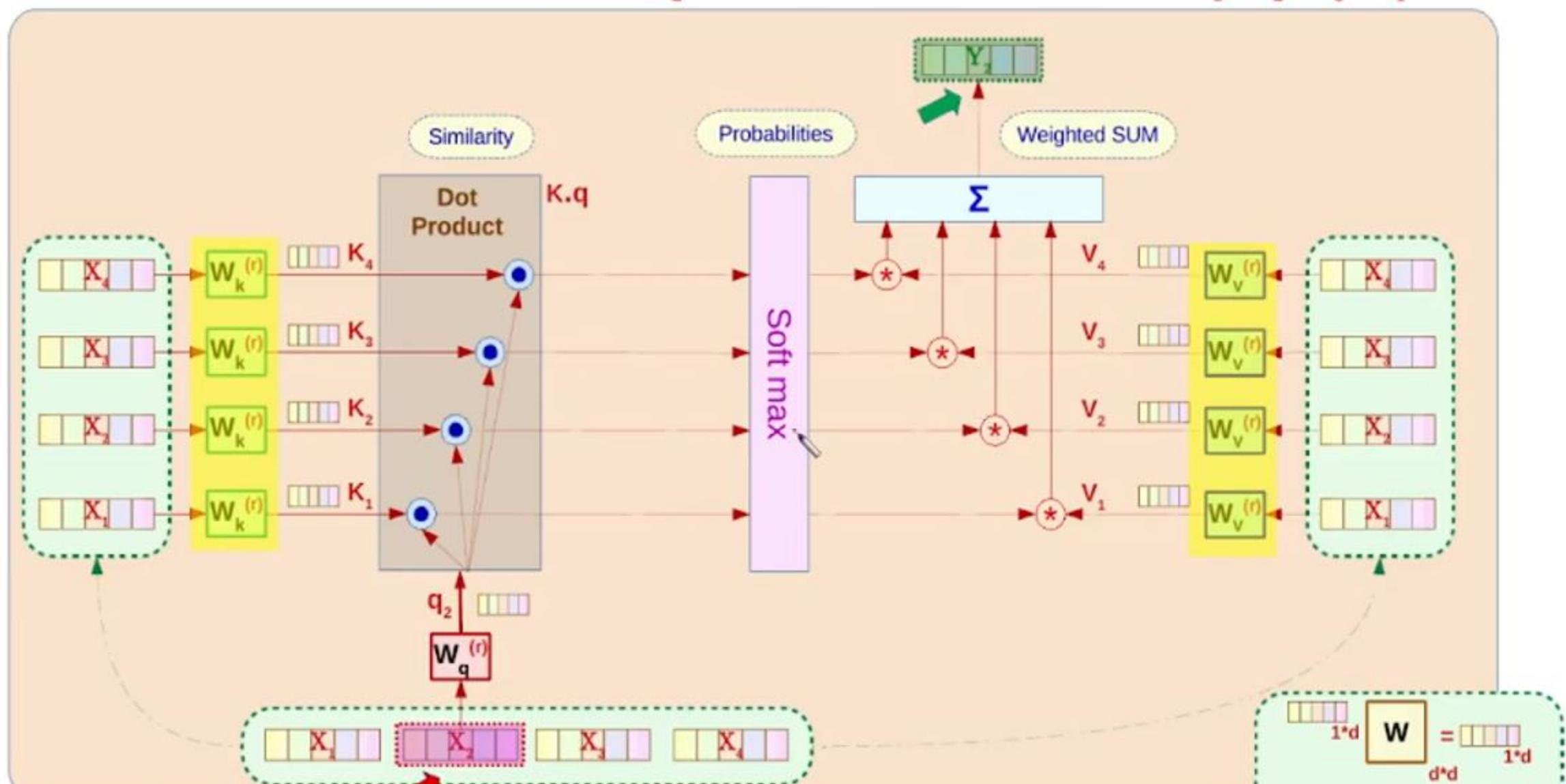
Activation: Linear

Bias: Zeros

Inputs: All Words Within Same Paragraph (X_1, X_2, X_3, X_4, X_5)

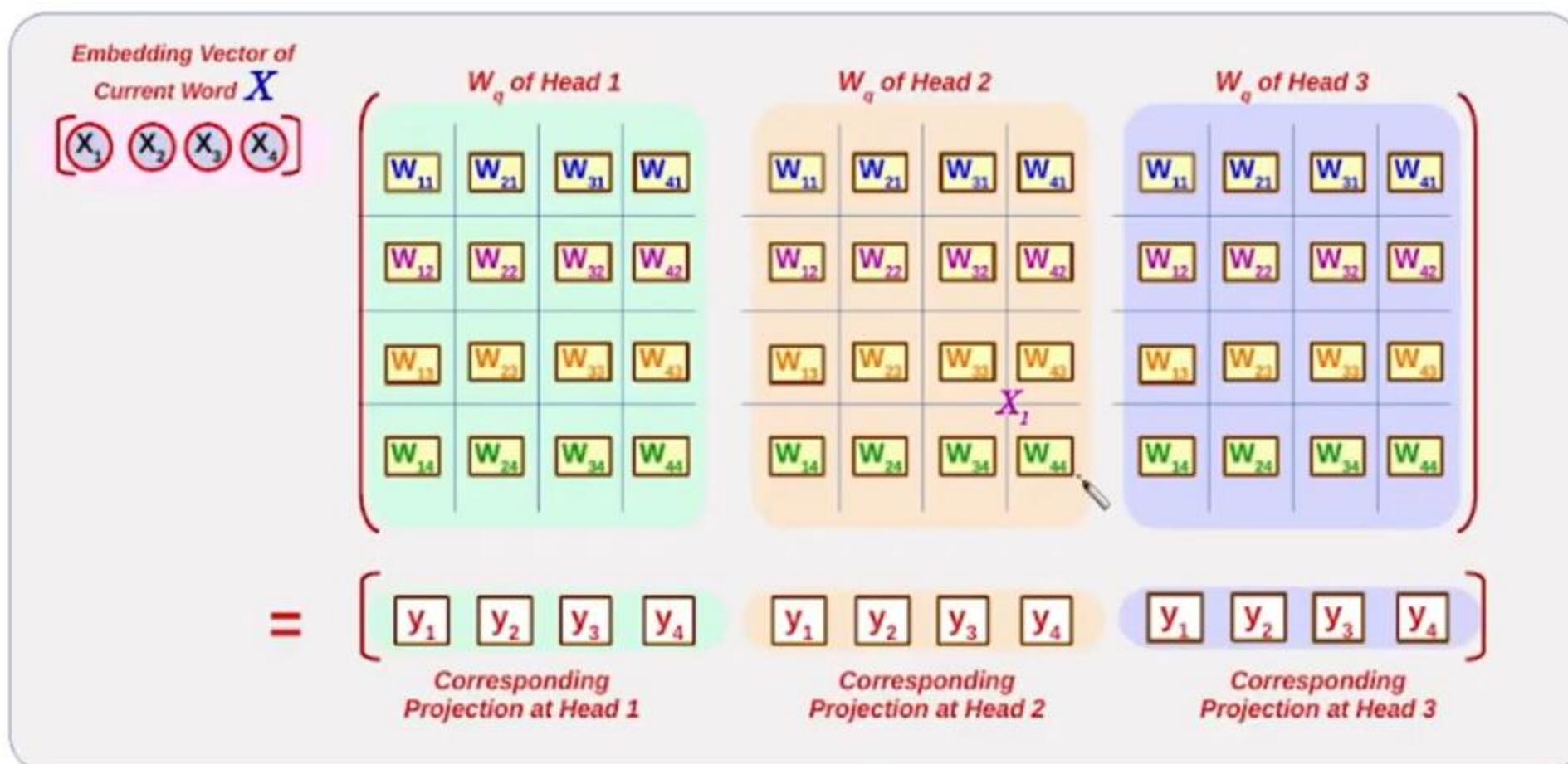


Self Attention Mechanism (Linear Transformation of Inputs)
(Self Attention of word X_2 w.r.t All words in Sentence (X_1, X_2, X_3, X_4)



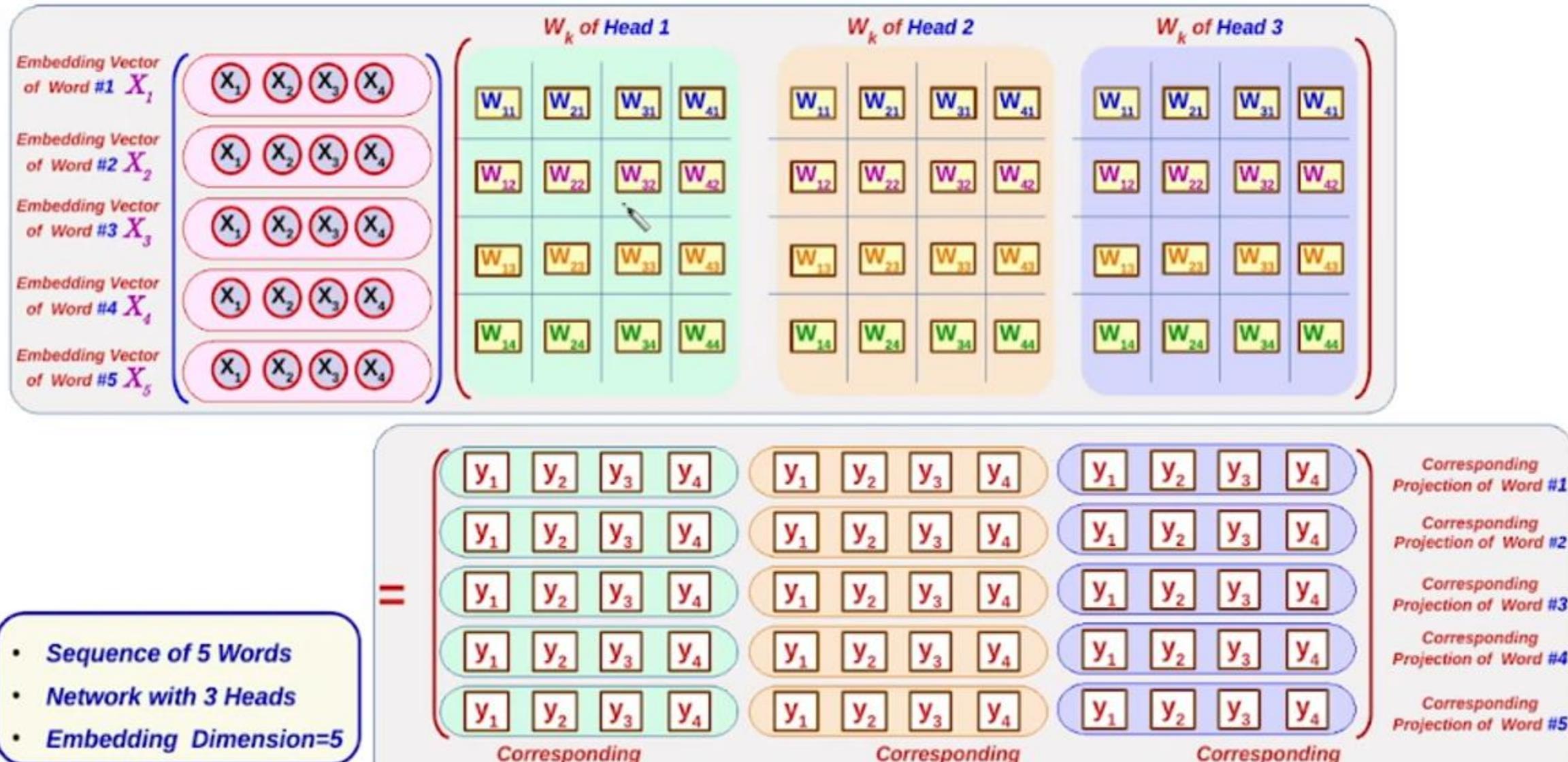
Using FF NN for Multiplication of Vector by a Multiple Matrices

Example: Multiply Vector of Current Word by Multiple (different) W_q (Queries of Multiple Heads)

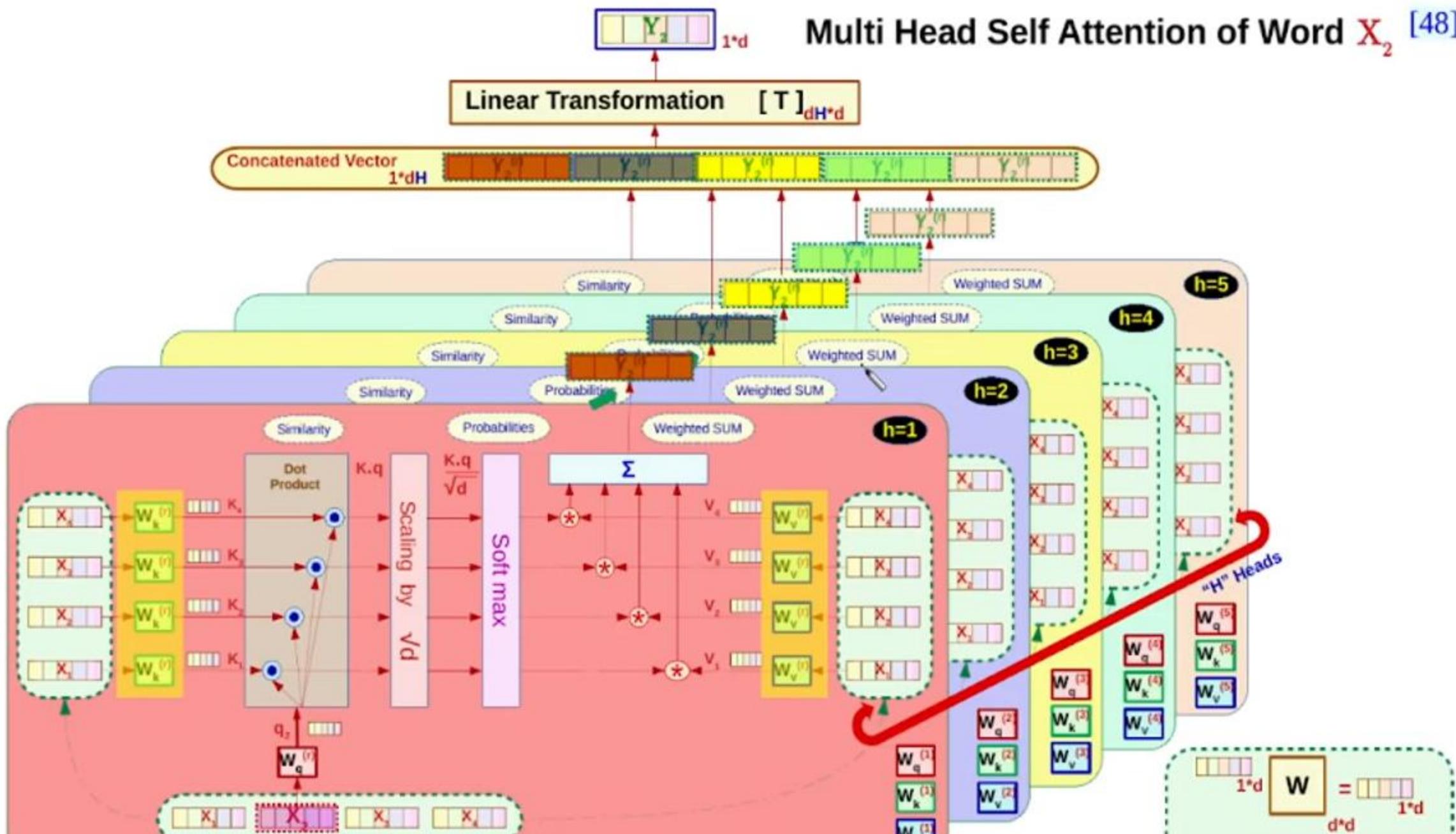


Using FF NN for Multiplication of Multiple Vectors by Multiple Matrices

Usage: Multiply Vectors of All Words by Multiple (different) W_k (Keys of Multiple Heads)



Multi Head Self Attention of Word X_2 [48]



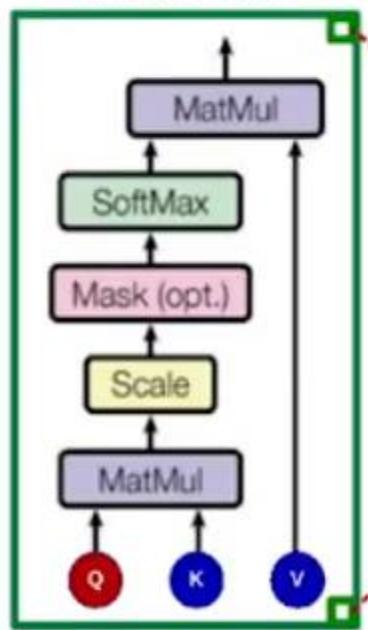
Adding Position Information

[1] Position Encoding

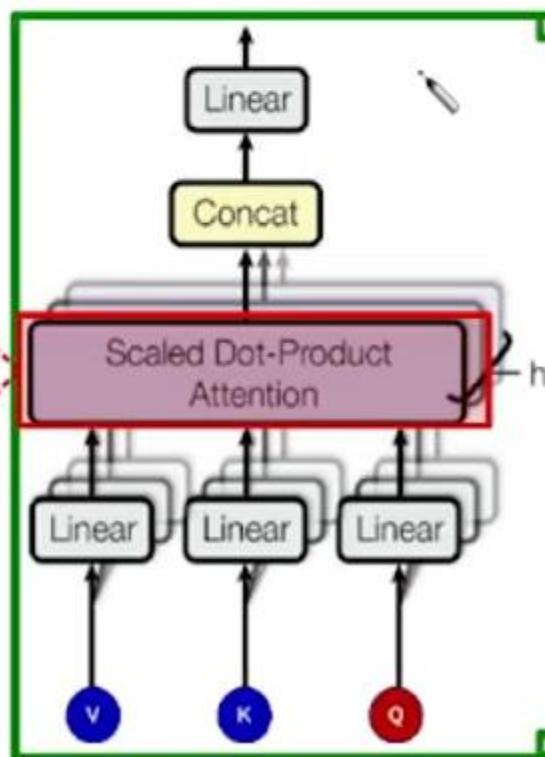
[2] Position Embedding

Encoder – Decoder Transformer

Scaled Dot-Product Attention



Multi-Head Attention



Encoder

