



TUTORIAL 10: TRANSFORMER

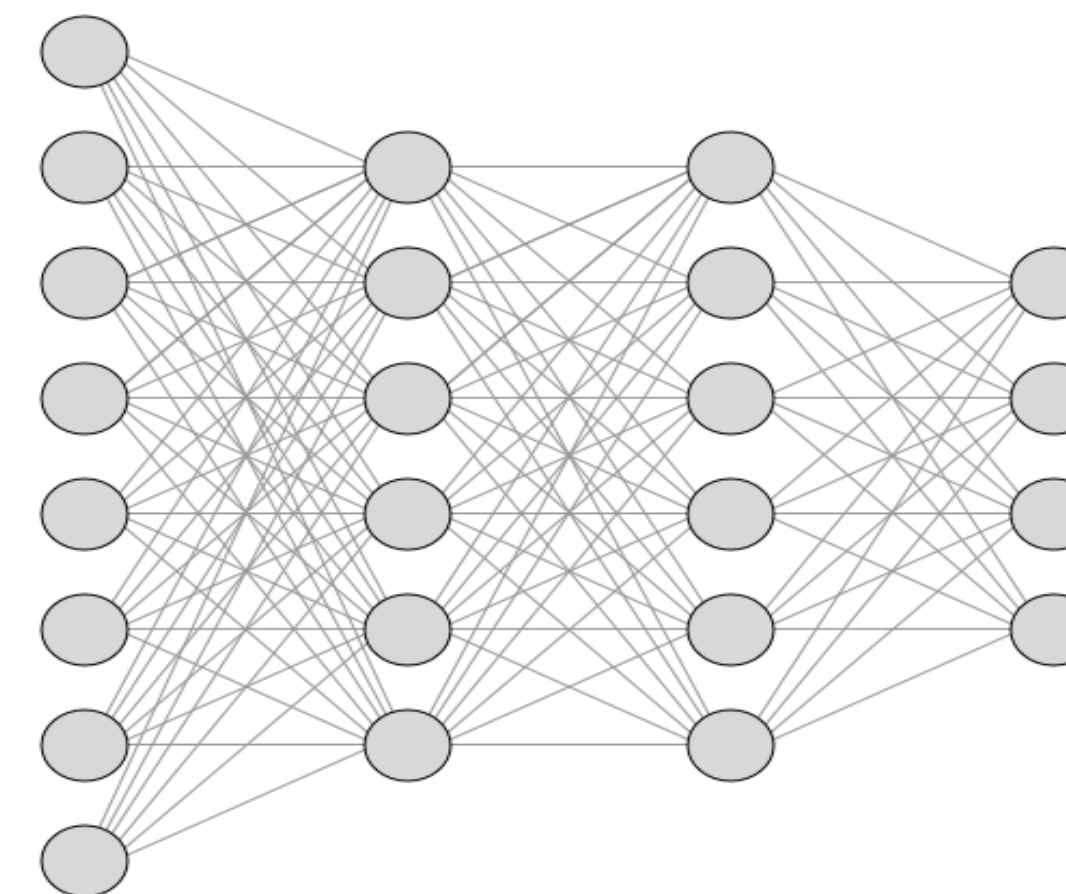
PART 1: ENCODER PART, ATTENTION

Instructor: Prof. Hsing-Kuo Pao
TA: Zolnamar Dorjsembe (Zola)

Date: April 30, 2024

Contents

- Recap: LSTM
- Attention
- Transformer's encoder part in detail





LET'S RECAP: RNN & LSTM

Join at www.kahoot.it
or with the Kahoot! app

Game PIN:

526 8431



Kahoot!

Waiting for players...



Remember the last week's activity



Solving Transformer by Hand: A Step-by-Step Math Example

Source: <https://levelup.gitconnected.com/understanding-transformers-from-start-to-end-a-step-by-step-math-example-16d4e64e6eb1>



Step 1. Defining Our Dataset

Dataset (corpus)

I drink and I know things.

When you play the game of thrones, you win or you die.

The true enemy won't wait out the storm, He brings the storm.

Our entire dataset containing only three sentences

Step 2. Finding Vocab Size

$$\text{vocab size} = \text{count}(\text{set}(N))$$

vocab_size formula where N is total number of words



Dataset (Corpus)

I drink and I know things.

When you play the game of thrones, you win or you die.

The true enemy won't wait out the storm, He brings the storm.

$$N = \begin{bmatrix} \text{I, drink, and, I, Know, things,} \\ \text{When, you, play, the, game, of, thrones, you, win, or, you, die,} \\ \text{The, true, enemy, won't, wait, out, the, storm, He, brings, the, storm} \end{bmatrix}$$

calculating variable N



$$\text{vocab size} = \text{count}(\text{set}(N))$$

$$\begin{aligned} &\rightarrow \text{set} \left(\begin{array}{l} \text{I, drink, and, I, Know, things,} \\ \text{When, you, play, the, game, of, thrones, you, win, or, you, die,} \\ \text{The, true, enemy, won't, wait, out, the, storm, He, brings, the, storm} \end{array} \right) \\ &\rightarrow \text{count} \left(\begin{array}{l} \text{I, drink, and, Know, things, When, you, play, the, game, of,} \\ \text{thrones, win, or, die, true, enemy, won't, wait, out, storm, He,} \\ \text{brings} \end{array} \right) \\ &\rightarrow = 23 \end{aligned}$$

finding vocab size

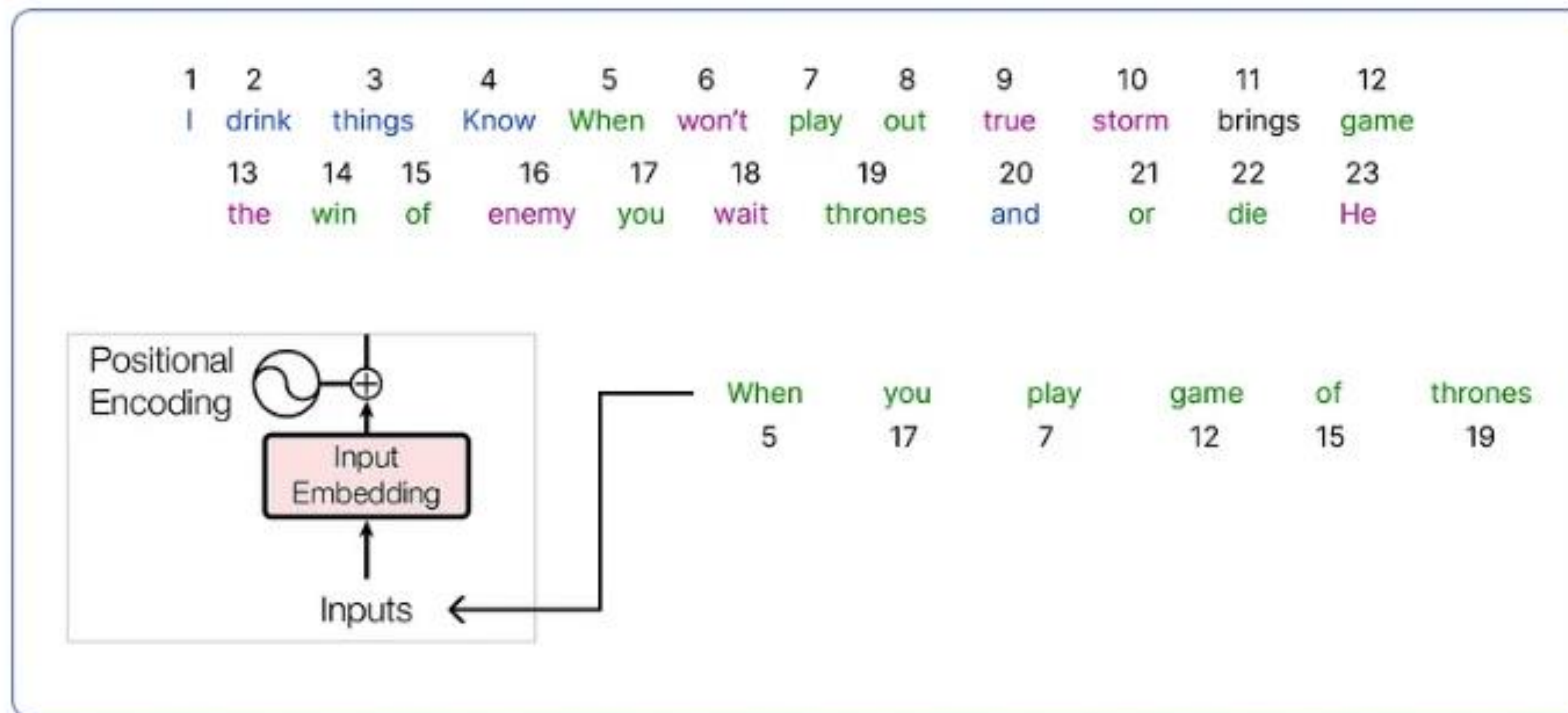


Step 3. Encoding

1	2	3	4	5	6	7	8	9	10	11	12
I	drink	things	Know	When	won't	play	out	true	storm	brings	game
13	14	15	16	17	18	19	20	21	22	23	
the	win	of	enemy	you	wait	thrones	and	or	die	He	

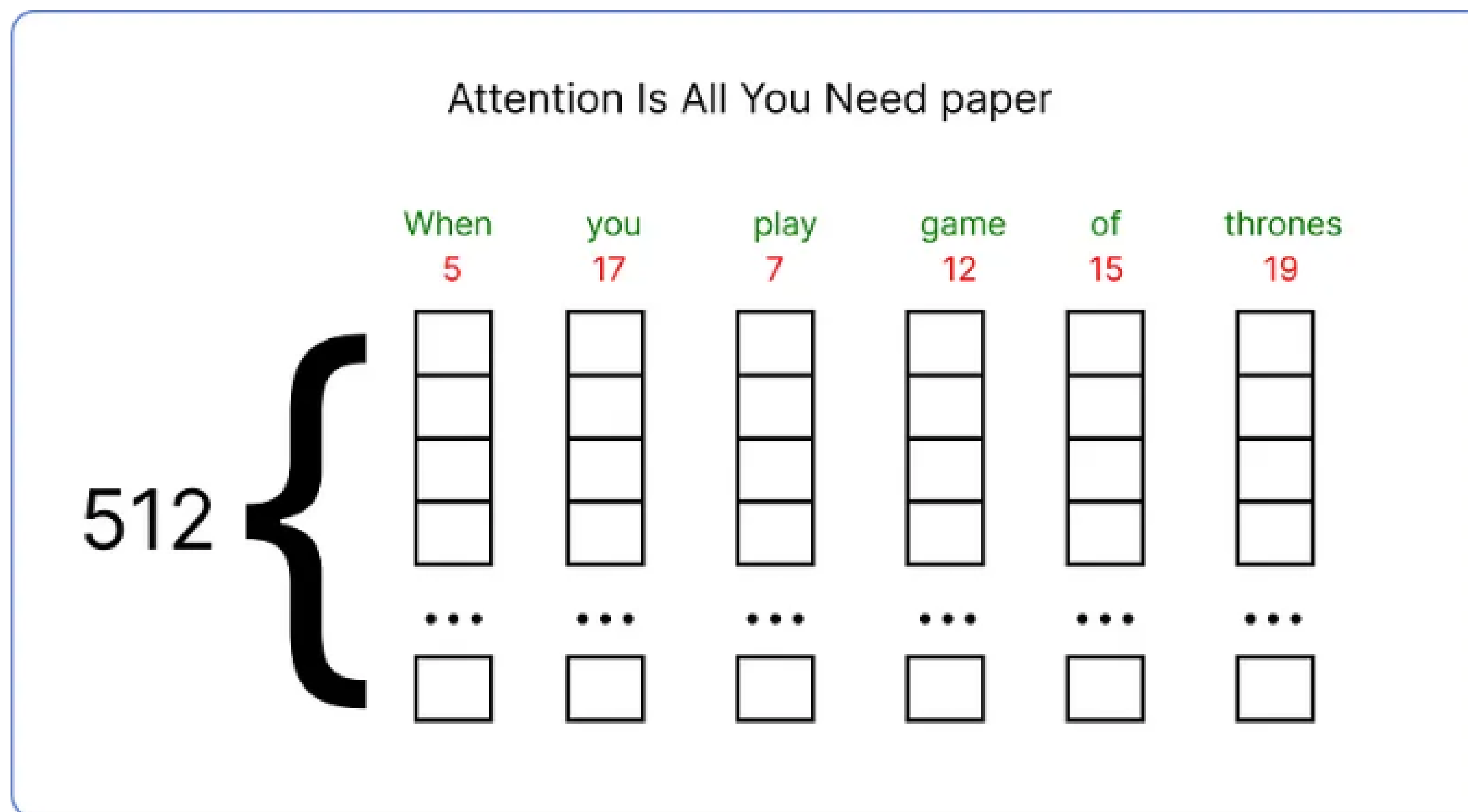
encoding our unique words

Step 4. Calculating Embedding



Input sentence for transformer

Step 4. Calculating Embedding



Original Paper uses 512 dimension vector

Step 4. Calculating Embedding

(For demonstration, use an embedding vector with a dimension of 6)



Embedding vectors of our input

Step 5. Calculating Positional Embedding

Embedding vector for any word

even position
odd position
even position
odd position
even position

...

For even position

$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i / d_{\text{model}}})$$

For odd position

$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i / d_{\text{model}}})$$

Positional Embedding formula

Step 5. Calculating Positional Embedding

When
5

i	e1	Position	Formula	p1
0	0.79	Even	$\sin(0/10000^{(2*0/6)})$	0
1	0.6	Odd	$\cos(0/10000^{(2*1/6)})$	1
2	0.96	Even	$\sin(0/10000^{(2*2/6)})$	0
3	0.64	Odd	$\cos(0/10000^{(2*3/6)})$	1
4	0.97	Even	$\sin(0/10000^{(2*4/6)})$	0
5	0.2	Odd	$\cos(0/10000^{(2*5/6)})$	1

d (dim) 6

POS 0

Positional Embedding for word: **When**

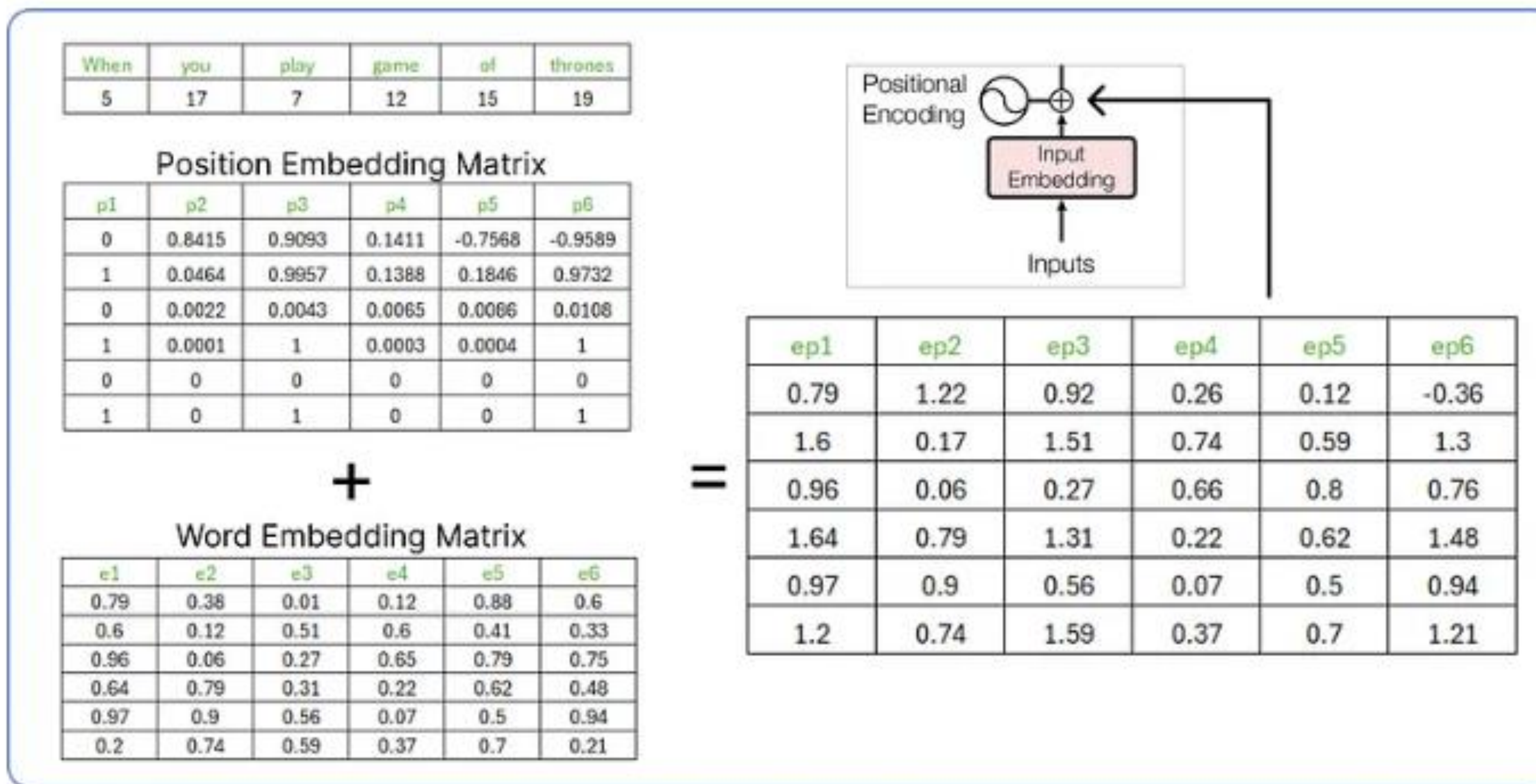
Step 5. Calculating Positional Embedding

	When	you	play	game	of	thrones
	5	17	7	12	15	19

i	p1	p2	p3	p4	p5	p6
0	0	0.8415	0.9093	0.1411	-0.7568	-0.9589
1	1	0.0464	0.9957	0.1388	0.1846	0.9732
2	0	0.0022	0.0043	0.0065	0.0086	0.0108
3	1	0.0001	1	0.0003	0.0004	1
4	0	0	0	0	0	0
5	1	0	1	0	0	1
d (dim)	6	6	6	6	6	6
POS	0	1	2	3	4	5

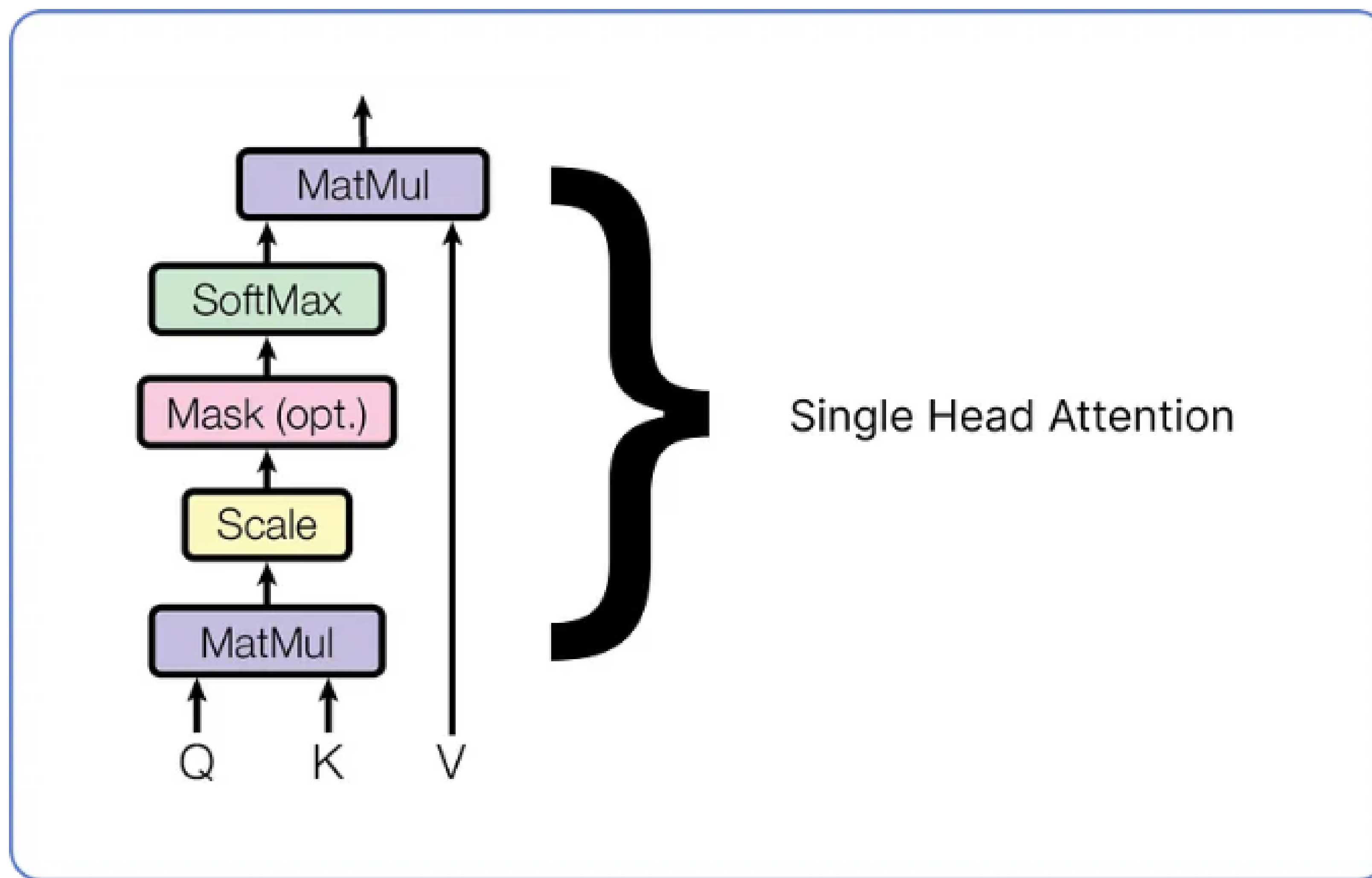
Calculating Positional Embeddings of our input **(The calculated values are rounded)**

Step 6. Concatenating Positional and Word Embeddings



concatenation step

Step 7. Multi Head Attention



Single Head attention in Transformer

Step 7. Multi Head Attention: Query matrix

Word Embedding + Positional Embedding

When	0.79	1.6	0.96	1.64	0.97	1.2
you	1.22	0.17	0.06	0.79	0.9	0.74
play	0.92	1.51	0.27	1.31	0.56	1.59
game	0.26	0.74	0.66	0.22	0.07	0.37
of	0.12	0.59	0.8	0.62	0.5	0.7
thrones	-0.36	1.3	0.76	1.48	0.94	1.21

6 x 6

Linear weights for query

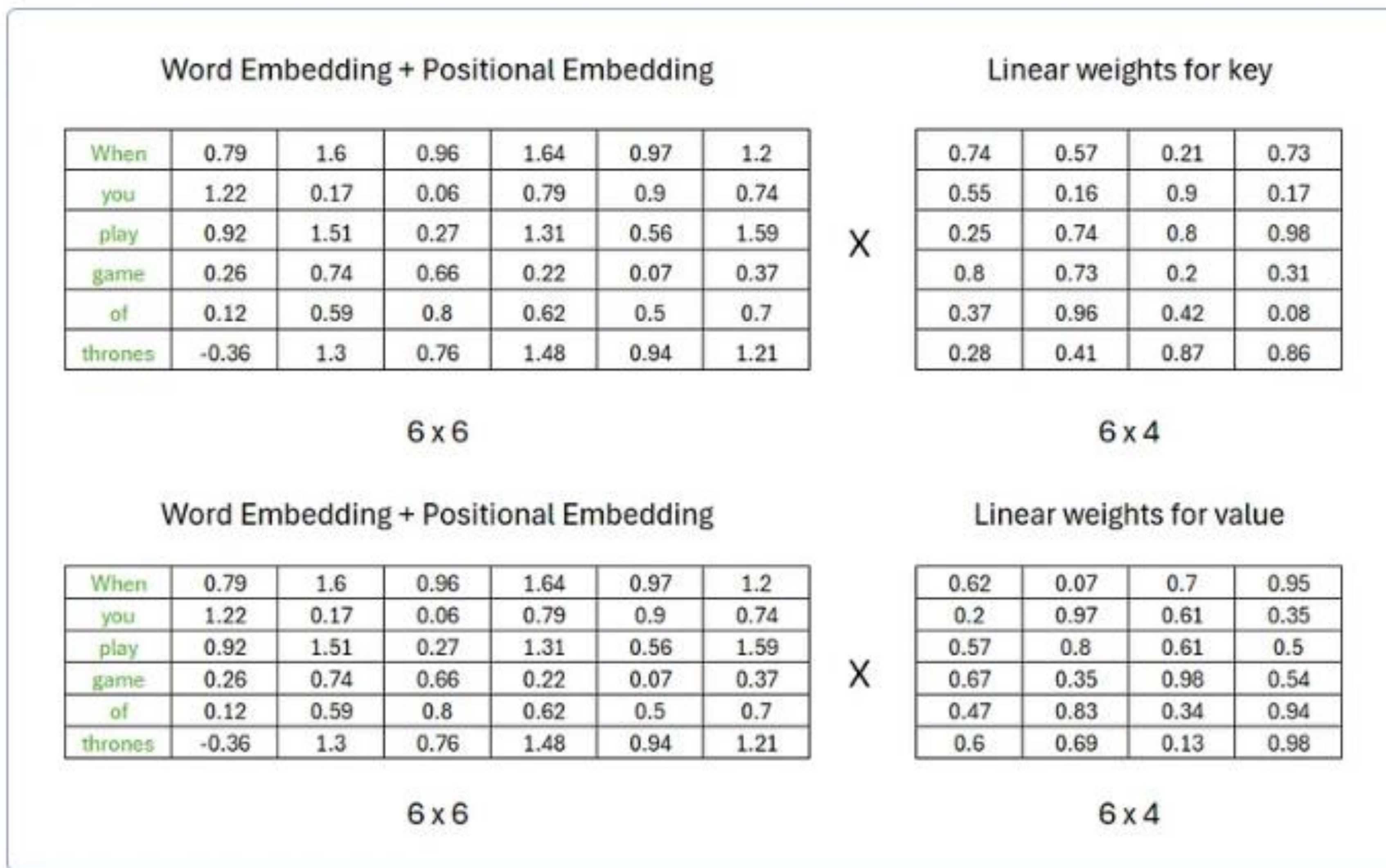
0.52	0.45	0.91	0.69
0.05	0.85	0.37	0.83
0.49	0.1	0.56	0.61
0.71	0.64	0.4	0.14
0.76	0.27	0.92	0.67
0.85	0.56	0.57	0.07

6 x 4

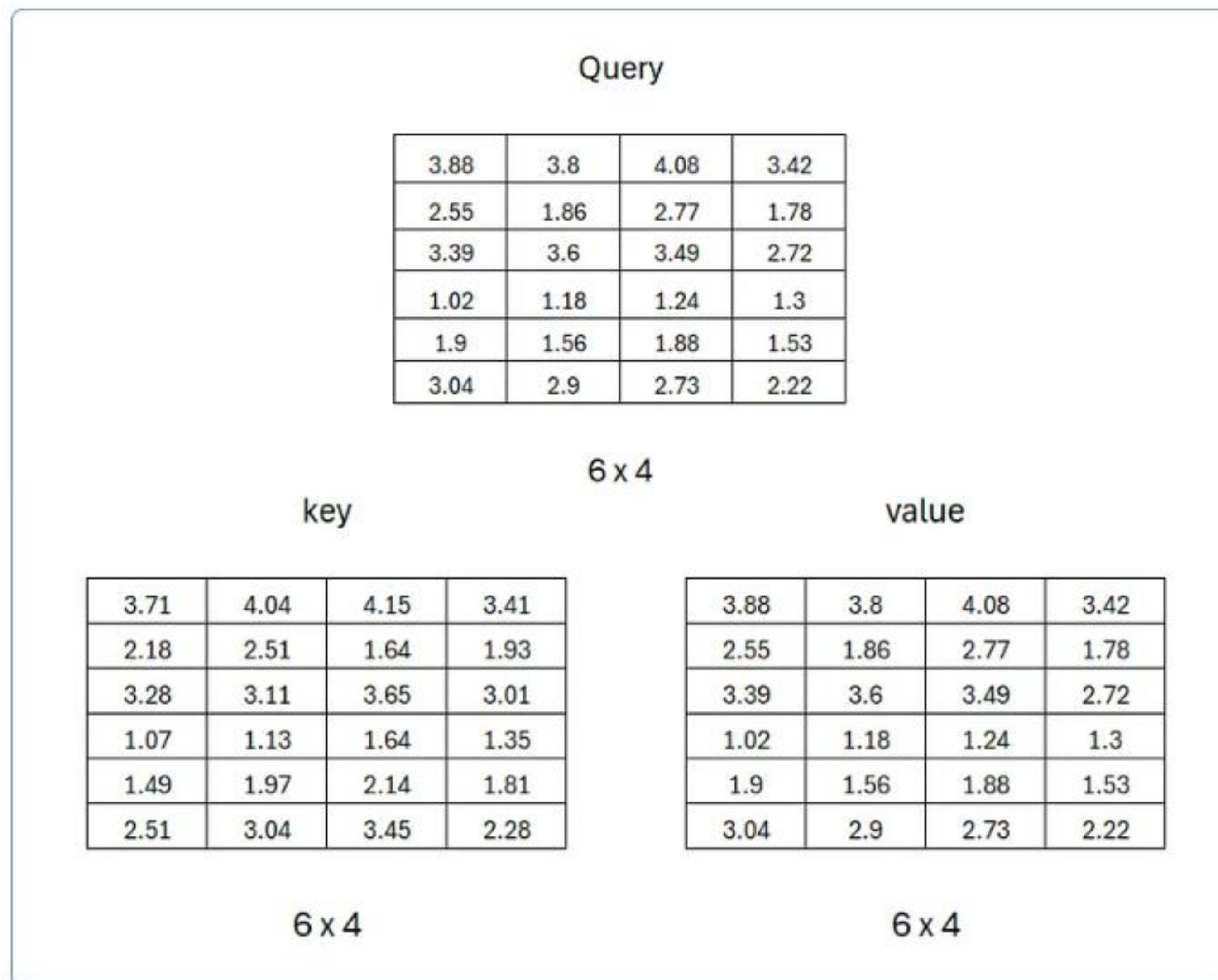
X

calculating Query matrix

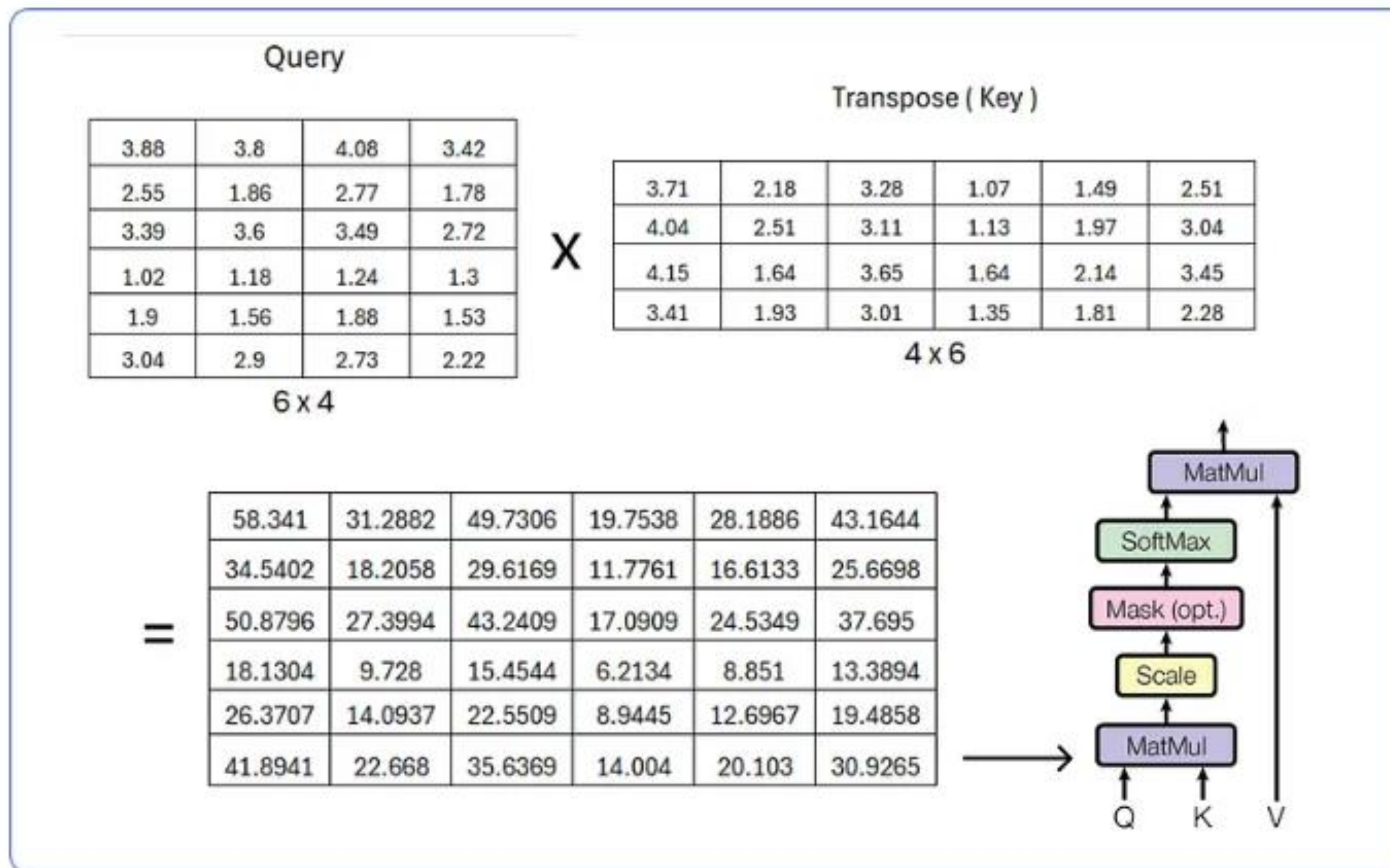
Step 7. Multi Head Attention: Key and Value Matrices



Step 7. Multi Head Attention: Query, Key and Values



Step 7. Multi Head Attention: Calculating single-head attention



matrix multiplication between Query and Key

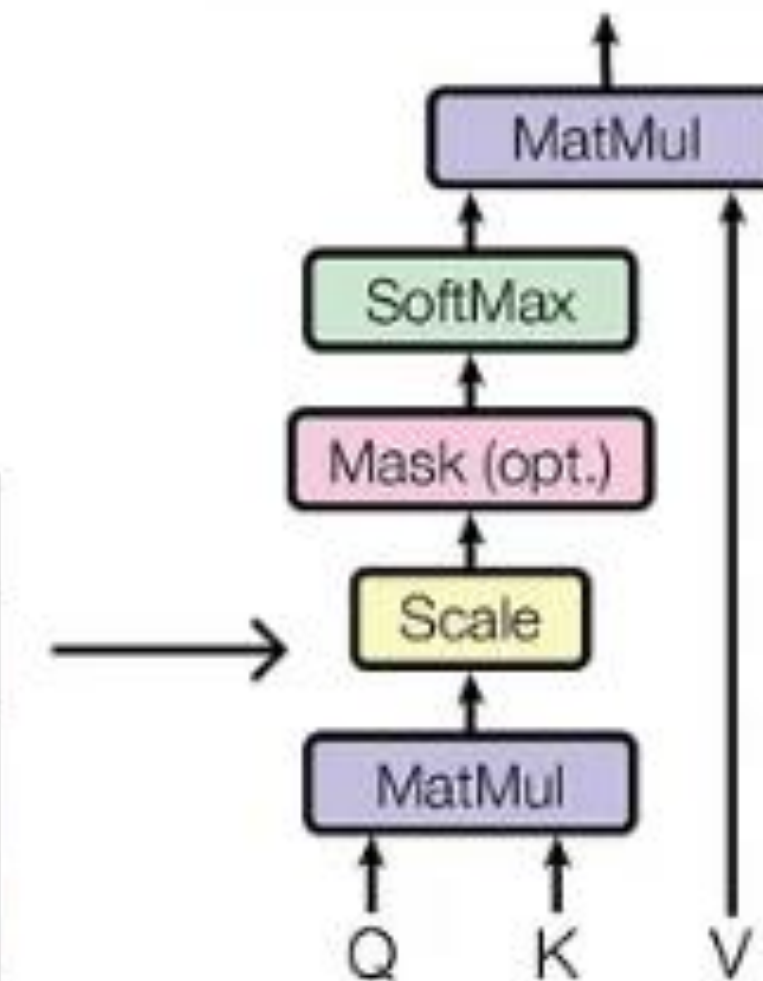
Step 7. Multi Head Attention: Calculating single-head attention

58.341	31.2882	49.7306	19.7538	28.1886	43.1644
34.5402	18.2058	29.6169	11.7761	16.6133	25.6698
50.8796	27.3994	43.2409	17.0909	24.5349	37.695
18.1304	9.728	15.4544	6.2134	8.851	13.3894
26.3707	14.0937	22.5509	8.9445	12.6967	19.4858
41.8941	22.668	35.6369	14.004	20.103	30.9265

$\sqrt{d_k}$ where d (dimension) is 6

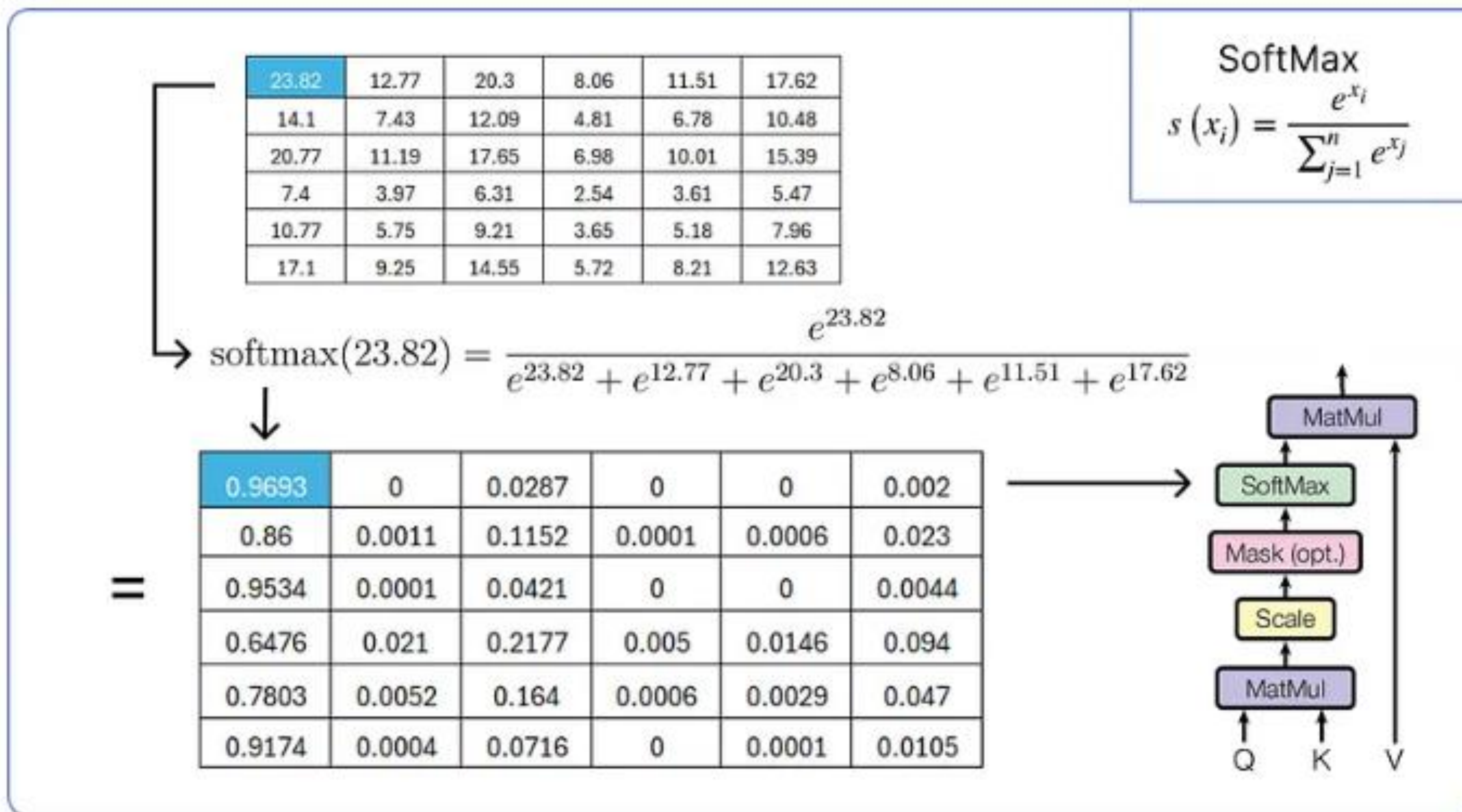
=

23.81721	12.77409	20.30219	8.062904	11.50852	17.62
14.1009	7.434201	12.09231	4.809165	6.781004	10.47973
20.77167	11.186	17.65266	6.976963	10.01433	15.39096
7.401542	3.972256	6.307436	2.535222	3.612997	5.466445
10.76551	5.752218	9.205999	3.64974	5.184753	7.956759
17.10152	9.254989	14.54997	5.715476	8.205791	12.62712

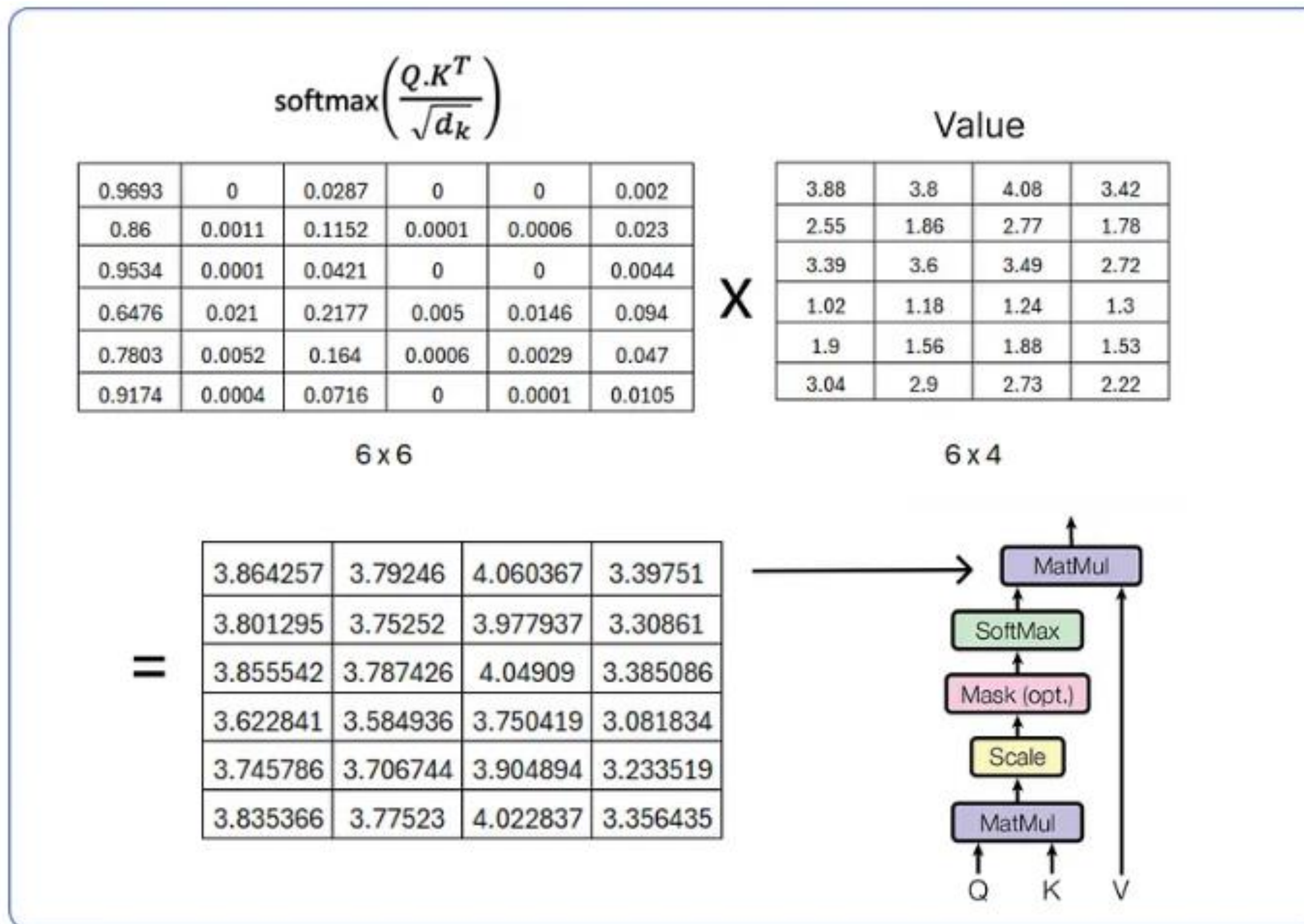


scaling the resultant matrix with dimension 5

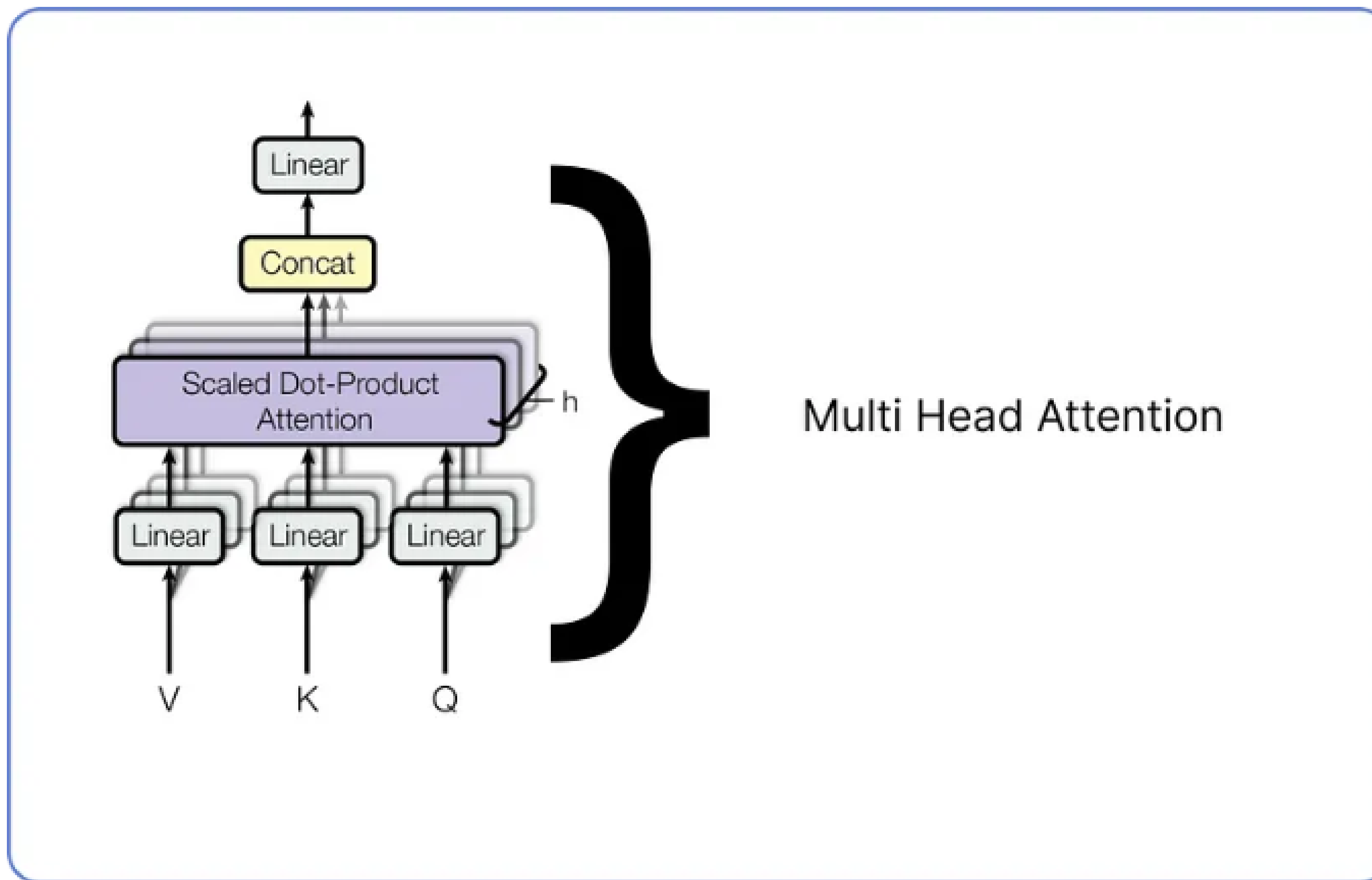
Step 7. Multi Head Attention: Calculating single-head attention



Step 7. Multi Head Attention: Calculating single-head attention



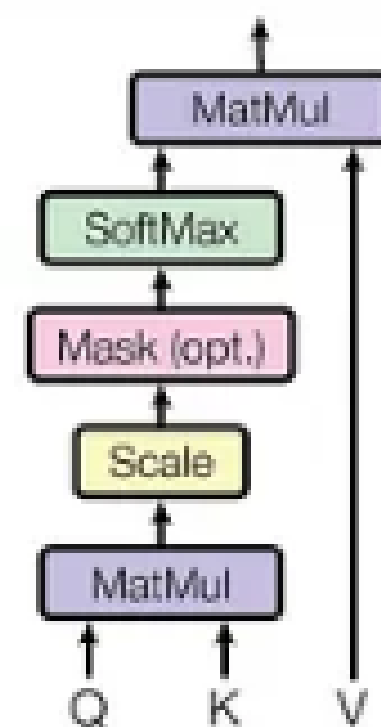
Step 7. Multi Head Attention:



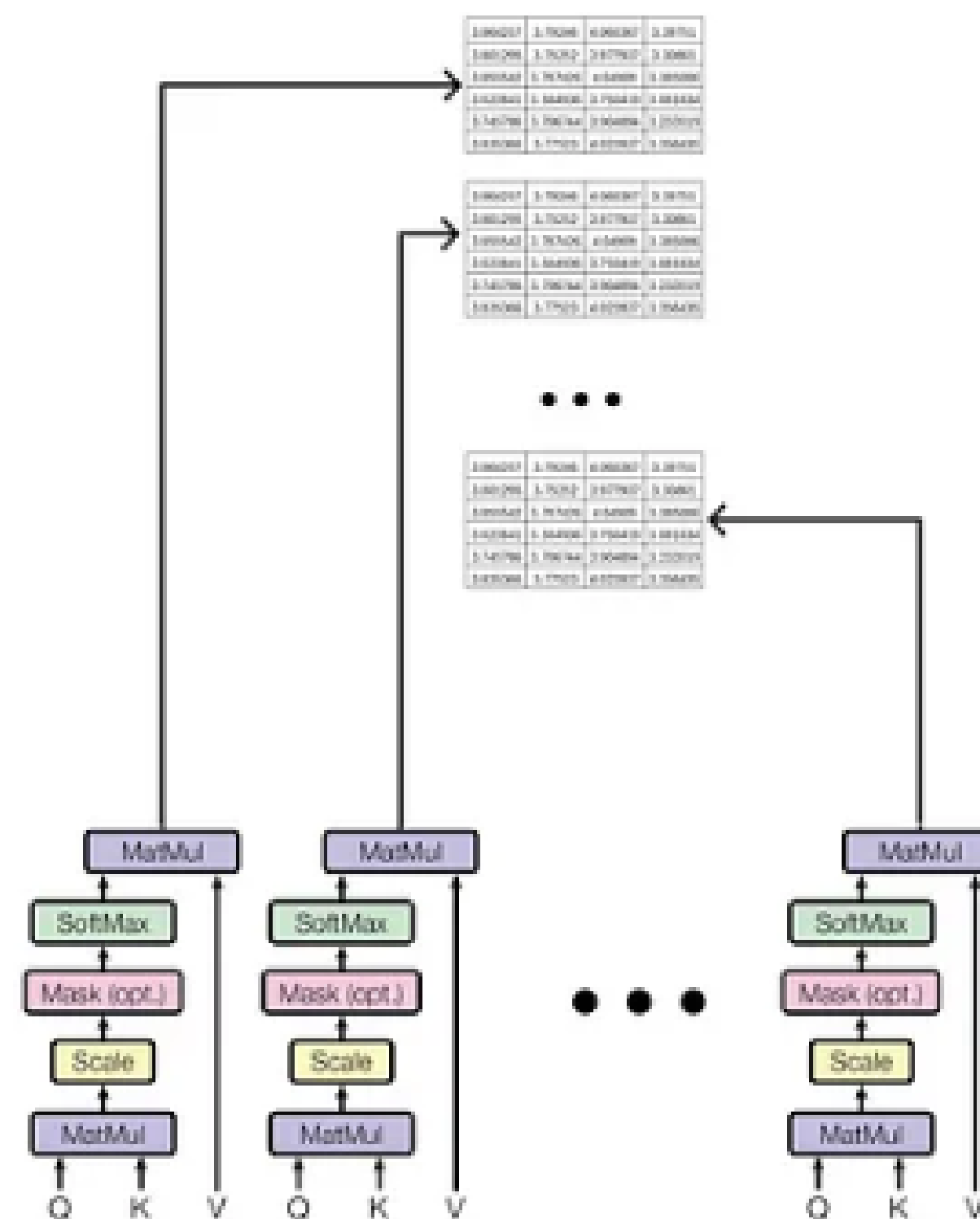
Step 7. Multi Head Attention: Single-head vs Multi-head attention

Single Head Attention
Our Case

3.864257	3.79246	4.060367	3.39751
3.801295	3.75252	3.977937	3.30861
3.855542	3.787426	4.04909	3.385086
3.622841	3.584936	3.750419	3.081834
3.745786	3.706744	3.904894	3.233519
3.835366	3.77523	4.022837	3.356435



Multi Head Attention (N Heads)
Real world Case Concatenation



Step 7. Multi Head Attention: Normalizing single-head attention

$$\begin{array}{c}
 \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) \quad \times \quad \text{Value} \\
 \begin{array}{|c|c|c|c|}
 \hline
 3.86 & 3.79 & 4.06 & 3.4 \\
 \hline
 3.8 & 3.75 & 3.98 & 3.31 \\
 \hline
 3.86 & 3.79 & 4.05 & 3.39 \\
 \hline
 3.62 & 3.58 & 3.75 & 3.08 \\
 \hline
 3.75 & 3.71 & 3.9 & 3.23 \\
 \hline
 3.84 & 3.78 & 4.02 & 3.36 \\
 \hline
 \end{array} \\
 6 \times 4
 \end{array}
 \times
 \begin{array}{c}
 \text{Linear weights} \\
 \text{columns length must be} \\
 \text{(embedding+positional) matrix columns length} \\
 \begin{array}{|c|c|c|c|c|c|}
 \hline
 0.8 & 0.34 & 0.45 & 0.54 & 0.07 & 0.53 \\
 \hline
 0.85 & 0.74 & 0.78 & 0.5 & 0.75 & 0.55 \\
 \hline
 0.53 & 0.81 & 0.55 & 0.59 & 0.49 & 0.14 \\
 \hline
 0.7 & 0.6 & 0.12 & 0.42 & 0.29 & 0.87 \\
 \hline
 \end{array} \\
 4 \times 6
 \end{array}
 =
 \begin{array}{|c|c|c|c|c|c|}
 \hline
 10.84 & 9.45 & 7.33 & 7.8 & 6.09 & 7.66 \\
 \hline
 10.65 & 9.28 & 7.22 & 7.67 & 5.99 & 7.51 \\
 \hline
 10.83 & 9.43 & 7.33 & 7.79 & 6.08 & 7.65 \\
 \hline
 10.08 & 8.77 & 6.85 & 7.25 & 5.67 & 7.09 \\
 \hline
 10.48 & 9.12 & 7.11 & 7.54 & 5.89 & 7.38 \\
 \hline
 10.77 & 9.38 & 7.29 & 7.75 & 6.05 & 7.6 \\
 \hline
 \end{array}$$

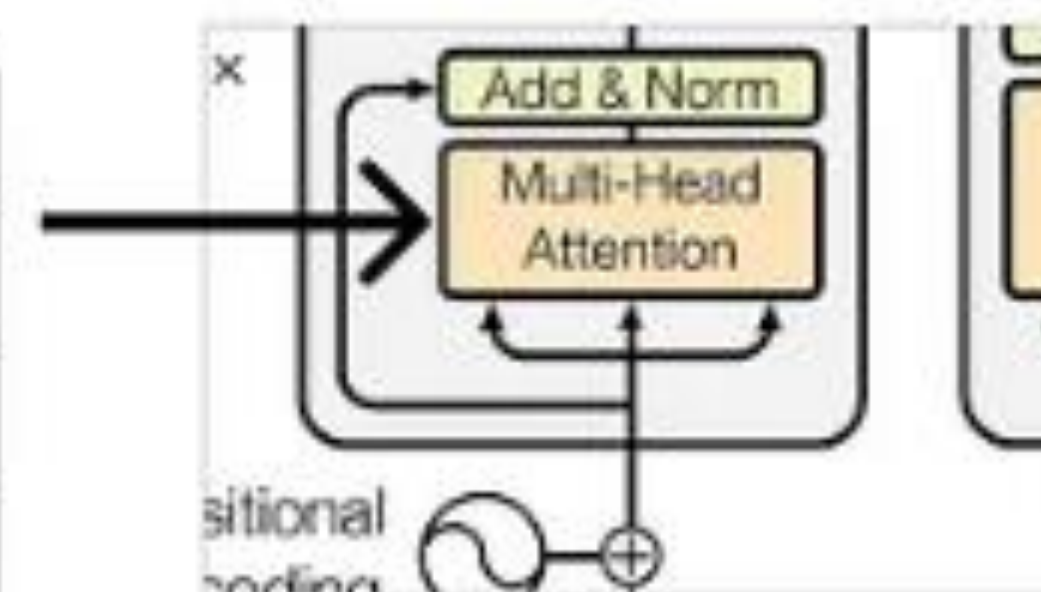
normalizing single head attention matrix

Step 7. Multi Head Attention:

Output of Multi Head attention

10.84	9.45	7.33	7.8	6.09	7.66
10.65	9.28	7.22	7.67	5.99	7.51
10.83	9.43	7.33	7.79	6.08	7.65
10.08	8.77	6.85	7.25	5.67	7.09
10.48	9.12	7.11	7.54	5.89	7.38
10.77	9.38	7.29	7.75	6.05	7.6

6 x 6



Output matrix of multi head attention

Step 8. Adding and Normalizing:

Word Embedding + Positional Embedding

When	0.79	1.6	0.96	1.64	0.97	1.2
you	1.22	0.17	0.06	0.79	0.9	0.74
play	0.92	1.51	0.27	1.31	0.56	1.59
game	0.26	0.74	0.66	0.22	0.07	0.37
of	0.12	0.59	0.8	0.62	0.5	0.7
thrones	-0.36	1.3	0.76	1.48	0.94	1.21

6 x 6

+

Output of Multi Head attention

10.84	9.45	7.33	7.8	6.09	7.66
10.65	9.28	7.22	7.67	5.99	7.51
10.83	9.43	7.33	7.79	6.08	7.65
10.08	8.77	6.85	7.25	5.67	7.09
10.48	9.12	7.11	7.54	5.89	7.38
10.77	9.38	7.29	7.75	6.05	7.6

6 x 6

=

11.63	11.05	8.29	9.44	7.06	8.86
11.87	9.45	7.28	8.46	6.89	8.25
11.75	10.94	7.6	9.1	6.64	9.24
10.34	9.51	7.51	7.47	5.74	7.46
10.6	9.71	7.91	8.16	6.39	8.08
10.41	10.68	8.05	9.23	6.99	8.81

Adding matrices to perform add and norm step

Step 8. Adding and Normalizing:

$$mean = \frac{\sum_{i=1}^N X_i}{N}$$

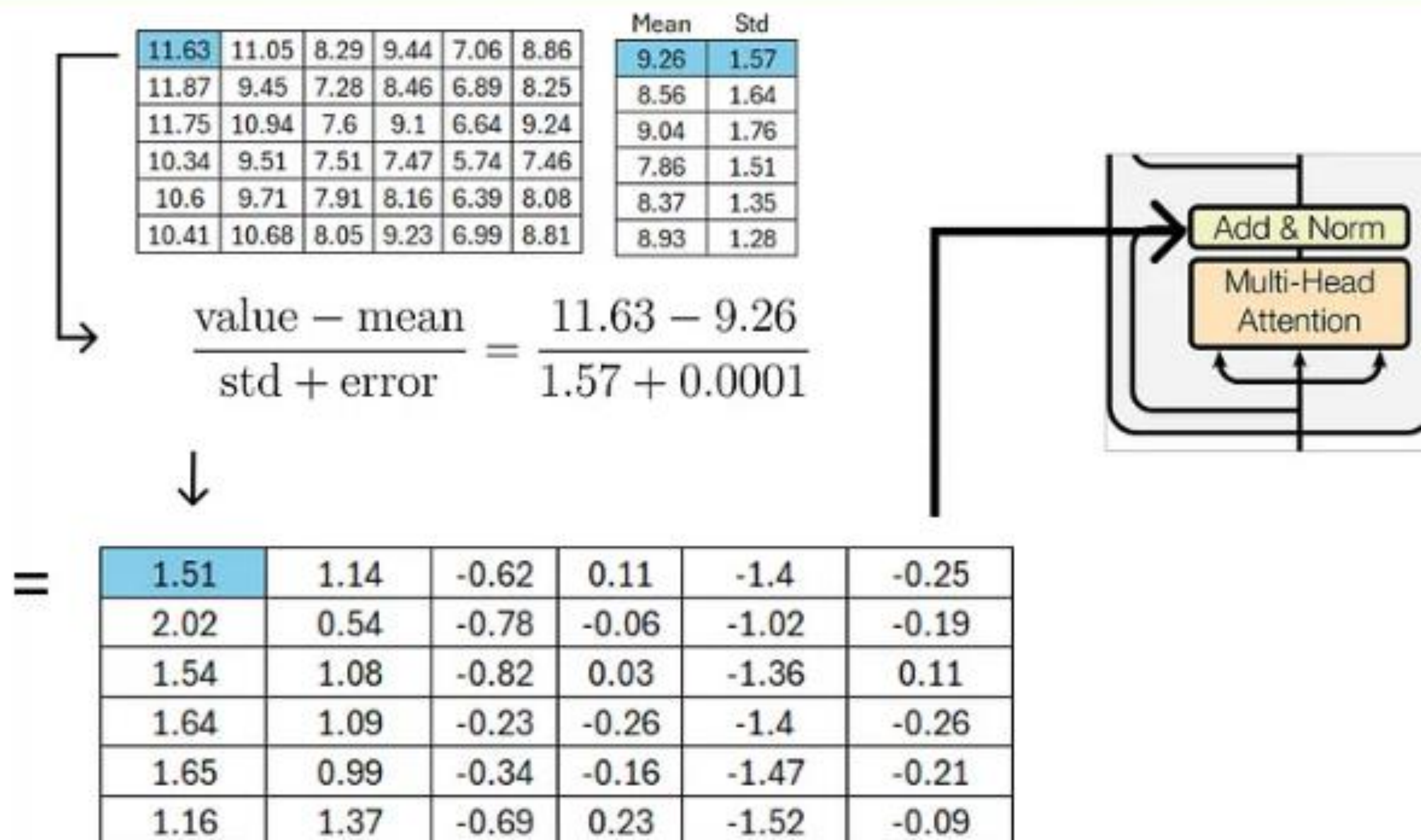
$$standard\ dev. = \sqrt{\frac{\sum_{i=1}^N (X_i - \mu)^2}{N}}$$

Row Wise Implementation

						Mean	Standard Deviation
11.63	11.05	8.29	9.44	7.06	8.86	9.26	1.57
11.87	9.45	7.28	8.46	6.89	8.25	8.56	1.64
11.75	10.94	7.6	9.1	6.64	9.24	9.04	1.76
10.34	9.51	7.51	7.47	5.74	7.46	7.86	1.51
10.6	9.71	7.91	8.16	6.39	8.08	8.37	1.35
10.41	10.68	8.05	9.23	6.99	8.81	8.93	1.28

calculating meand and std.

Step 8. Adding and Normalizing:

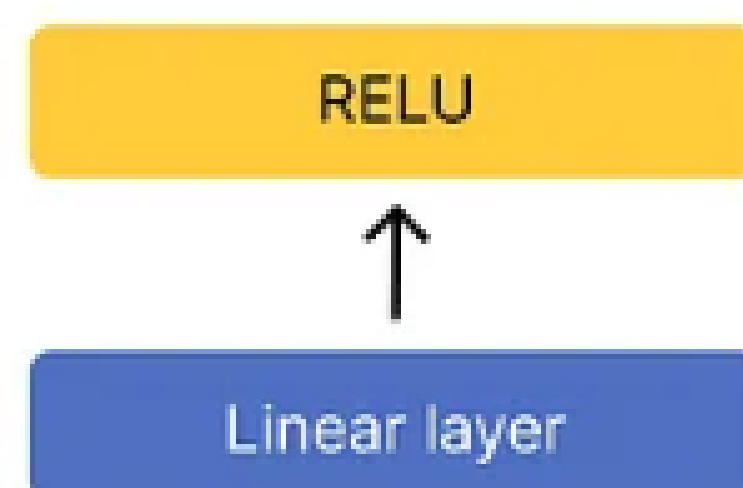


Step 9. Feed Forward Network:

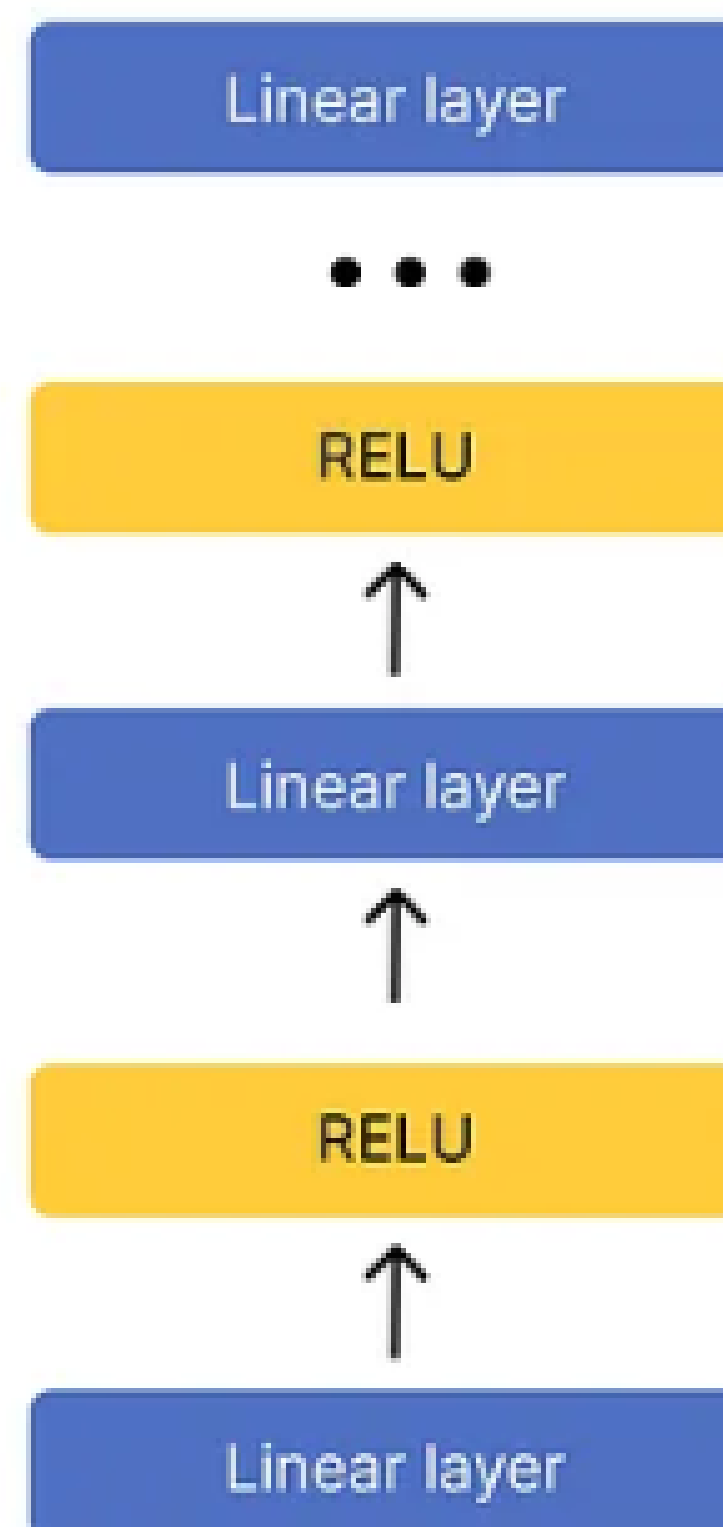
$$\text{ReLU}(x) = \max(0, x)$$

$$\text{Linear Layer} = X \cdot W + b$$

our case (one linear layer)



Real world case
(multiple layers)



Step 9. Feed Forward Network:

$$\begin{array}{c}
 \begin{array}{c} \text{Matrix after add and norm step} \\ \begin{array}{|c|c|c|c|c|c|} \hline 1.51 & 1.14 & -0.62 & 0.11 & -1.4 & -0.25 \\ \hline 2.02 & 0.54 & -0.78 & -0.06 & -1.02 & -0.19 \\ \hline 1.54 & 1.08 & -0.82 & 0.03 & -1.36 & 0.11 \\ \hline 1.64 & 1.09 & -0.23 & -0.26 & -1.4 & -0.26 \\ \hline 1.65 & 0.99 & -0.34 & -0.16 & -1.47 & -0.21 \\ \hline 1.16 & 1.37 & -0.69 & 0.23 & -1.52 & -0.09 \\ \hline \end{array} \\ 6 \times 6 \end{array} \\
 \times \\
 \begin{array}{c} W \\ \begin{array}{|c|c|c|c|c|c|} \hline 0.5 & 0.05 & 0.97 & 0.22 & 0.56 & 0.02 \\ \hline 0.17 & 0.52 & 0.63 & 0.48 & 0.06 & 0.6 \\ \hline 0.53 & 0.87 & 0.47 & 0.1 & 0.31 & 0.79 \\ \hline 0.83 & 0.58 & 0.38 & 0.09 & 0.64 & 0.25 \\ \hline 0.81 & 0.85 & 0.74 & 0.35 & 0.31 & 0.53 \\ \hline 0.25 & 0.31 & 0.22 & 0.77 & 0.57 & 0.85 \\ \hline \end{array} \\ 6 \times 6 \end{array} \\
 \\
 \begin{array}{c} X \cdot W \\ \begin{array}{|c|c|c|c|c|c|} \hline 0.49 & 1.07 & 0.84 & 0.14 & 0.22 & 0.7 \\ \hline 0.24 & 1.26 & 1.11 & 0.12 & 0.46 & 0.97 \\ \hline 0.53 & 1.18 & -0.82 & 0.39 & 0.33 & 0.59 \\ \hline 0.53 & 0.97 & 0.98 & 0.15 & 0.16 & 0.52 \\ \hline 0.56 & 1.11 & -0.87 & 0.11 & 0.2 & 0.64 \\ \hline 0.62 & 1.02 & 0.61 & 0.26 & 0.14 & 0.52 \\ \hline \end{array} \\ 6 \times 6 \end{array} \\
 + \\
 \begin{array}{c} \text{Bias} \\ \begin{array}{|c|c|c|c|c|c|} \hline b1 & b2 & b3 & b4 & b5 & b6 \\ \hline 0.42 & 0.18 & 0.25 & 0.42 & 0.35 & 0.45 \\ \hline \end{array} \end{array} \\
 \\
 = \\
 \begin{array}{c} \begin{array}{|c|c|c|c|c|c|} \hline 0.91 & 1.25 & 1.09 & 0.56 & 0.57 & 1.15 \\ \hline 0.66 & 1.44 & 1.36 & 0.54 & 0.81 & 1.42 \\ \hline 0.95 & 1.36 & -0.57 & 0.81 & 0.68 & 1.04 \\ \hline 0.95 & 1.15 & 1.23 & 0.57 & 0.51 & 0.97 \\ \hline 0.98 & 1.29 & -0.62 & 0.53 & 0.55 & 1.09 \\ \hline 1.04 & 1.2 & 0.86 & 0.68 & 0.49 & 0.97 \\ \hline \end{array} \\ 6 \times 6 \end{array}
 \end{array}$$

Step 9. Feed Forward Network:

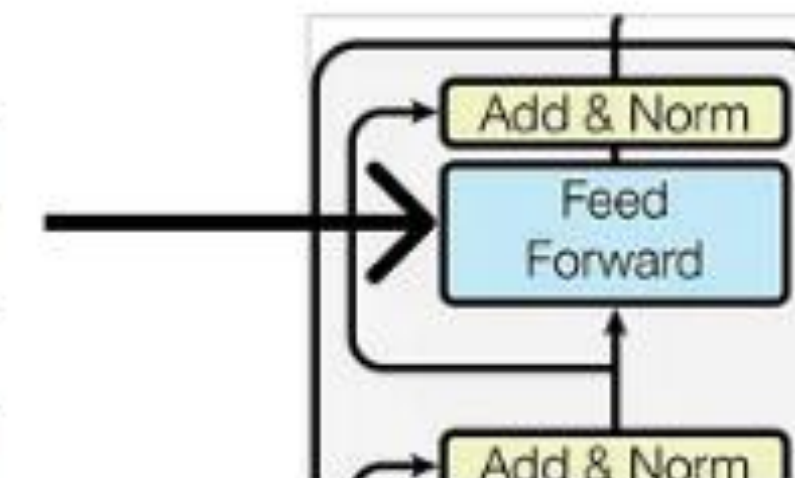
$$\text{ReLU}(x) = \max(0, x)$$

0.91	1.25	1.09	0.56	0.57	1.15
0.66	1.44	1.36	0.54	0.81	1.42
0.95	1.36	-0.57	0.81	0.68	1.04
0.95	1.15	1.23	0.57	0.51	0.97
0.98	1.29	-0.62	0.53	0.55	1.09
1.04	1.2	0.86	0.68	0.49	0.97

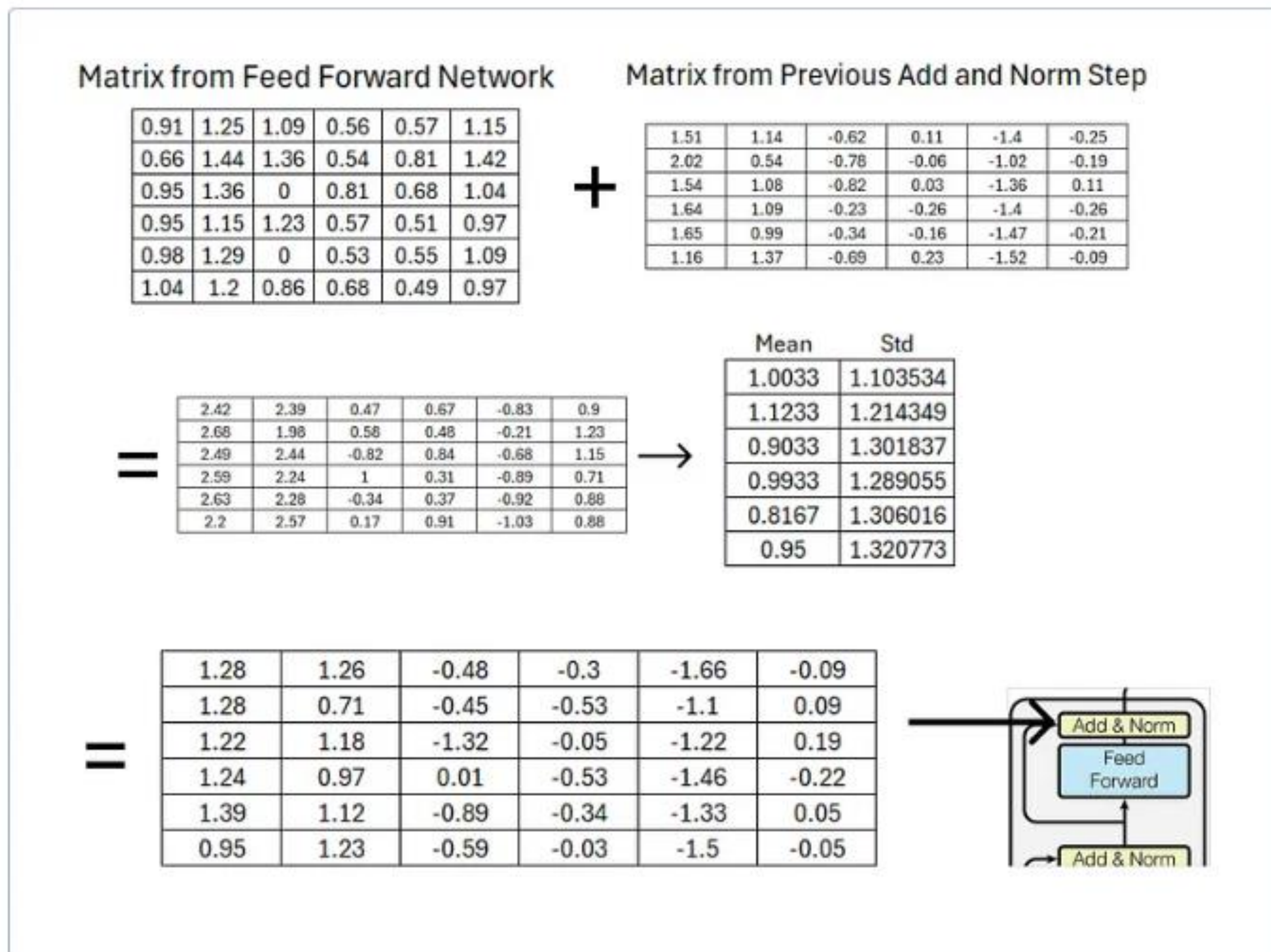
6 x 6

$$\rightarrow \max(0, 0.91)$$

0.91	1.25	1.09	0.56	0.57	1.15
0.66	1.44	1.36	0.54	0.81	1.42
0.95	1.36	0	0.81	0.68	1.04
0.95	1.15	1.23	0.57	0.51	0.97
0.98	1.29	0	0.53	0.55	1.09
1.04	1.2	0.86	0.68	0.49	0.97

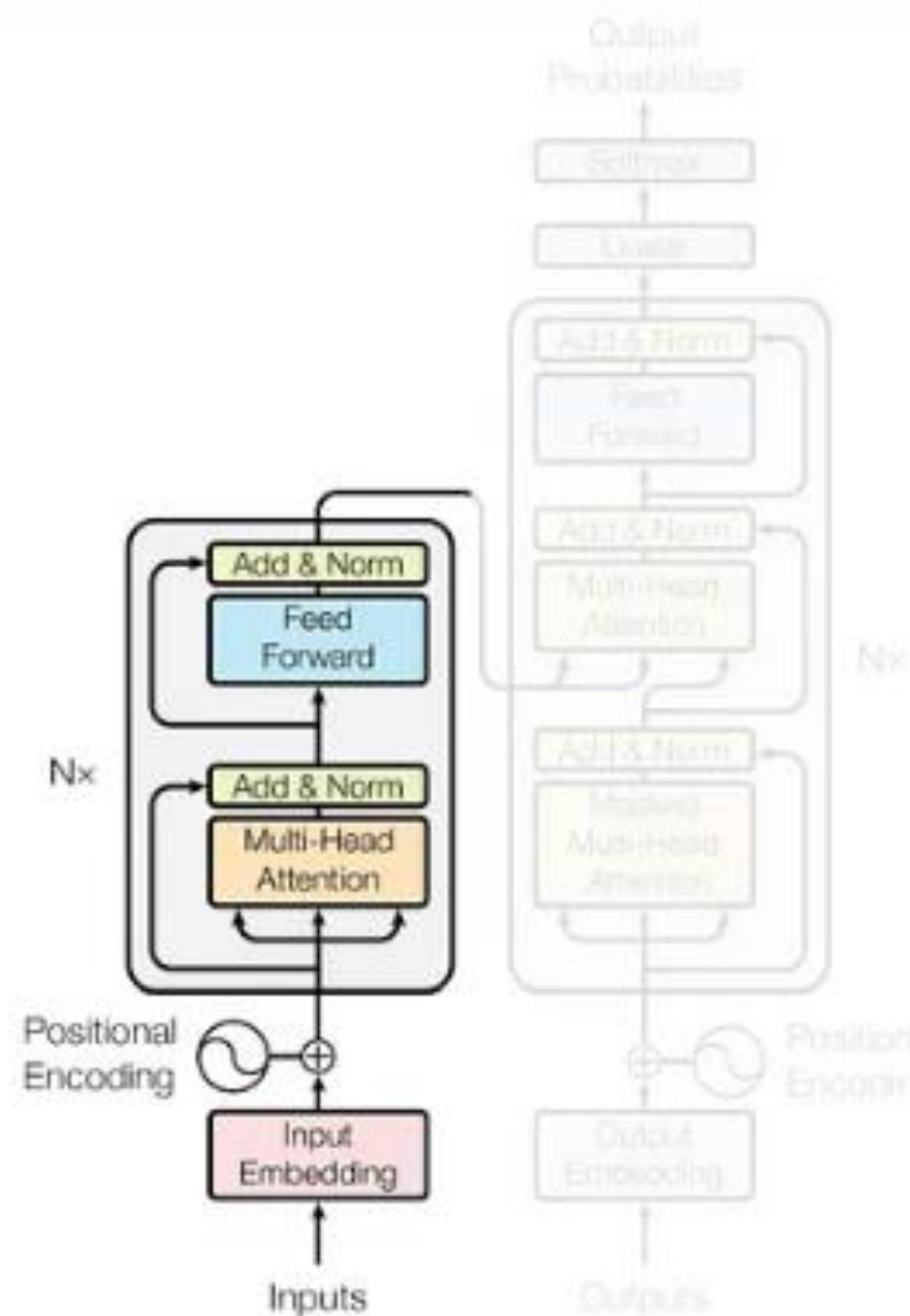


Step 10. Adding and Normalizing Again:



Step 11. Decoder Part:

What we have covered so far



What we have to cover

