

en

بخش inputs



Inputs

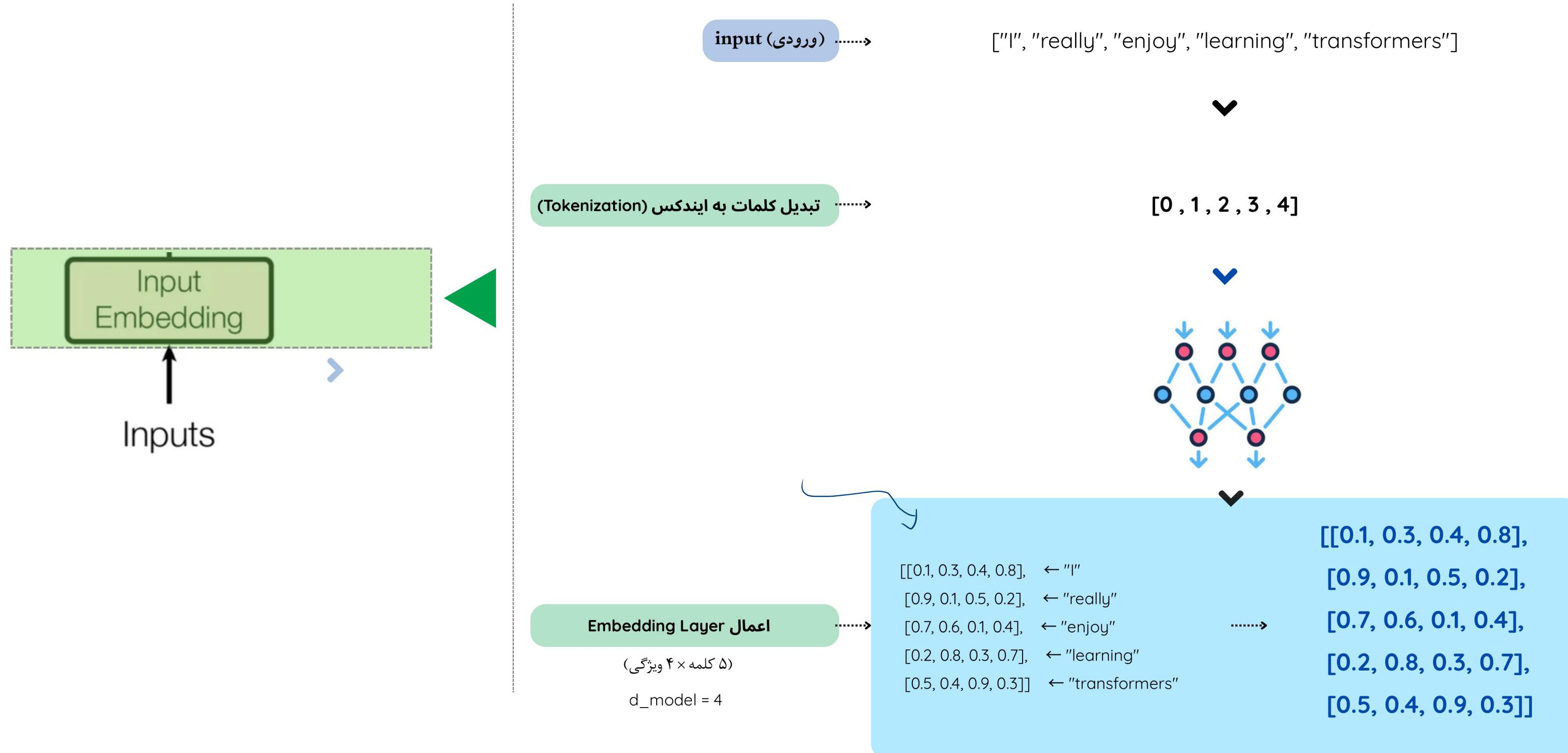


I really enjoy learning transformers



["I", "really", "enjoy", "learning", "transformers"]

بخش inputs



Word Embedding

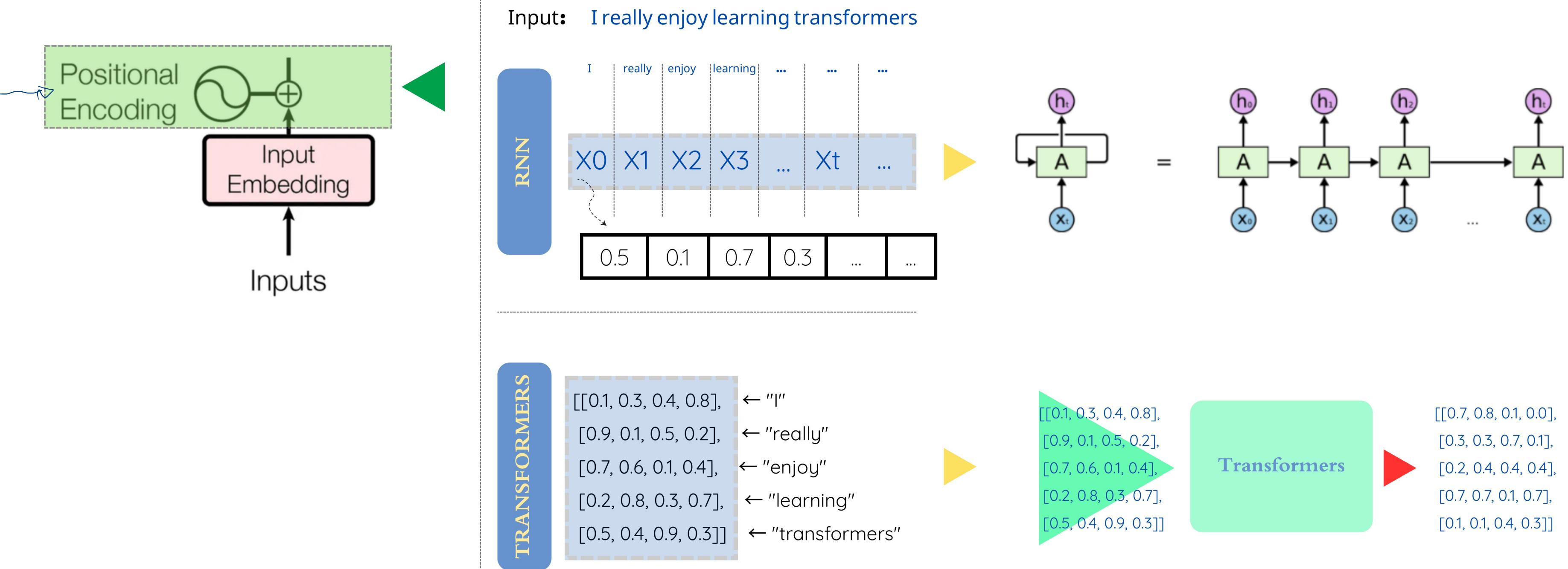
- One-Hot Encoding
- TF-IDF, Bag-of-Words
- Contextual Embedding
- Trainable Embedding Layer

چرا position encoding



برخلاف RNN‌ها که توکن‌ها را به ترتیب پردازش می‌کنند، Transformer کل جمله را یکجا (موازی) پردازش می‌کند.

به صورت ذاتی نمی‌داند ترتیب کلمات چیست!

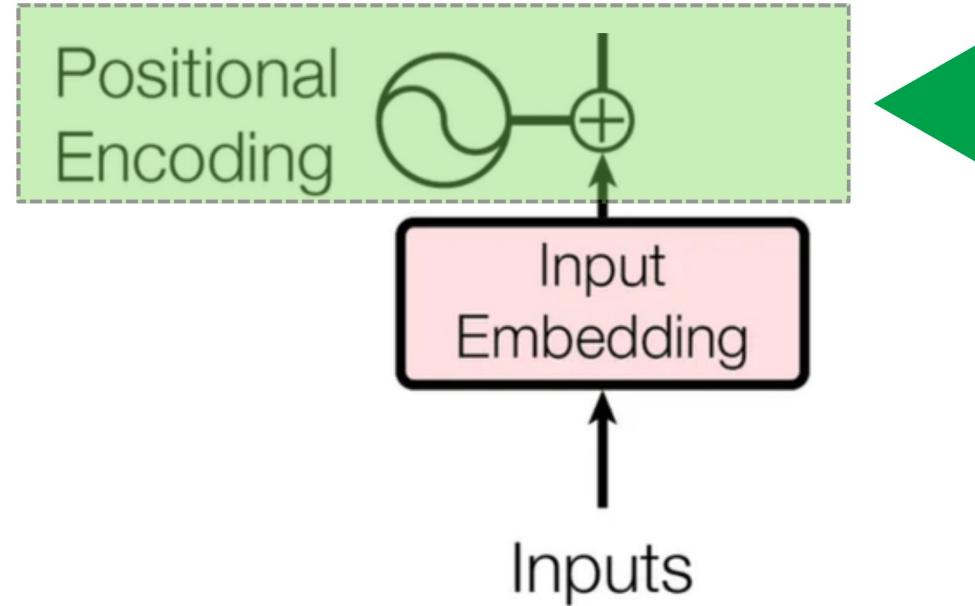


چگونه (sinusoidal positional encoding) ؟ positional encoding

برای اینکه مدل بتواند ترتیب توکن‌ها را تشخیص دهد، به هر توکن یک بردار موقعیتی (Positional Encoding) اضافه می‌شود.

دادن اطلاعات «موقعیت» به هر توکن به طوری که مدل بفهمد "I" اول جمله است و "آخر" آخر "transformers"

:Positional Encoding



For each position pos and dimension i :

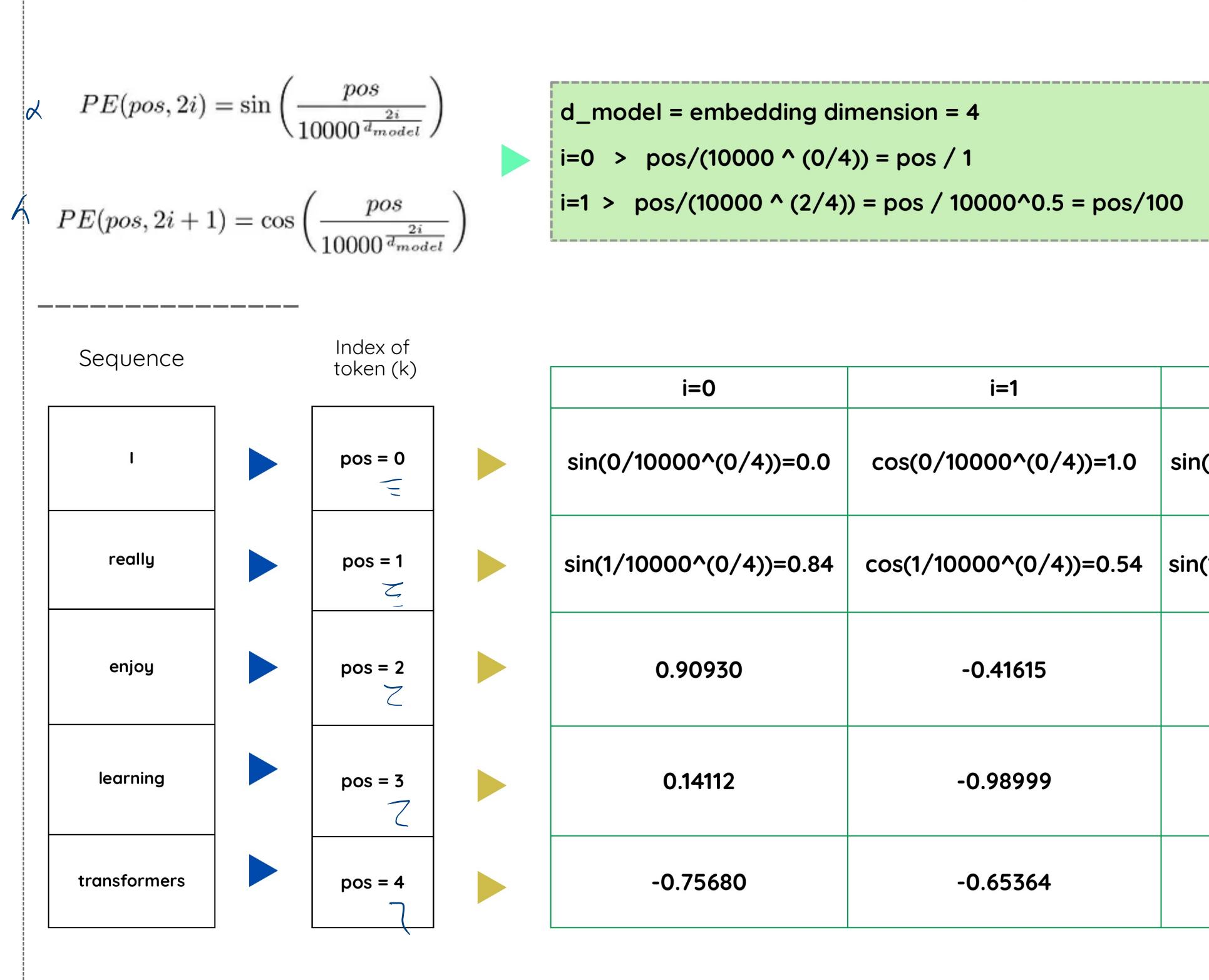
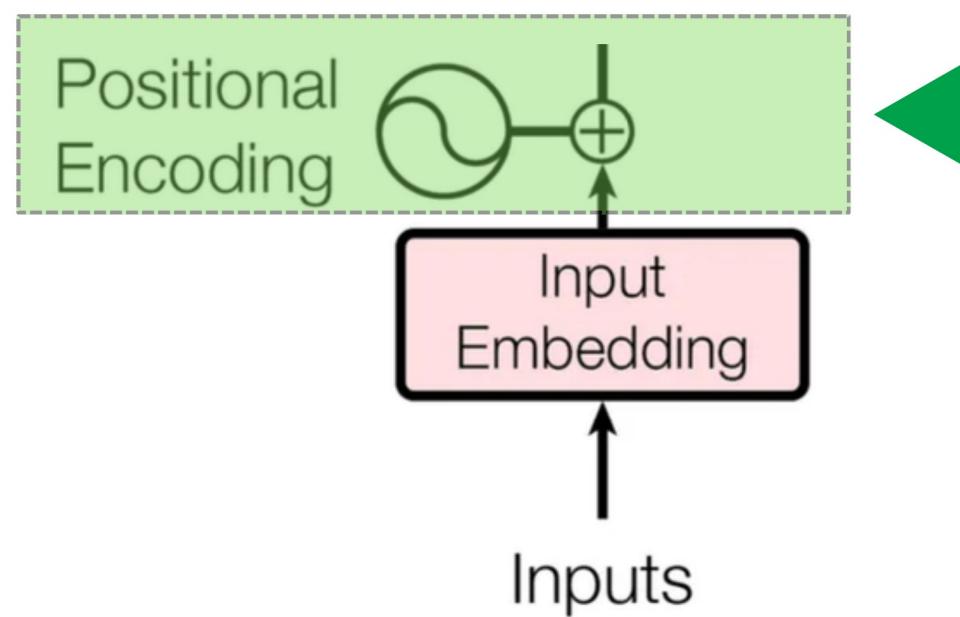
$$\text{PE}(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

$$\text{PE}(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

Here, d_{model} is the dimension of the model.

d_{model} = embedding dimension

چگونه position encoding

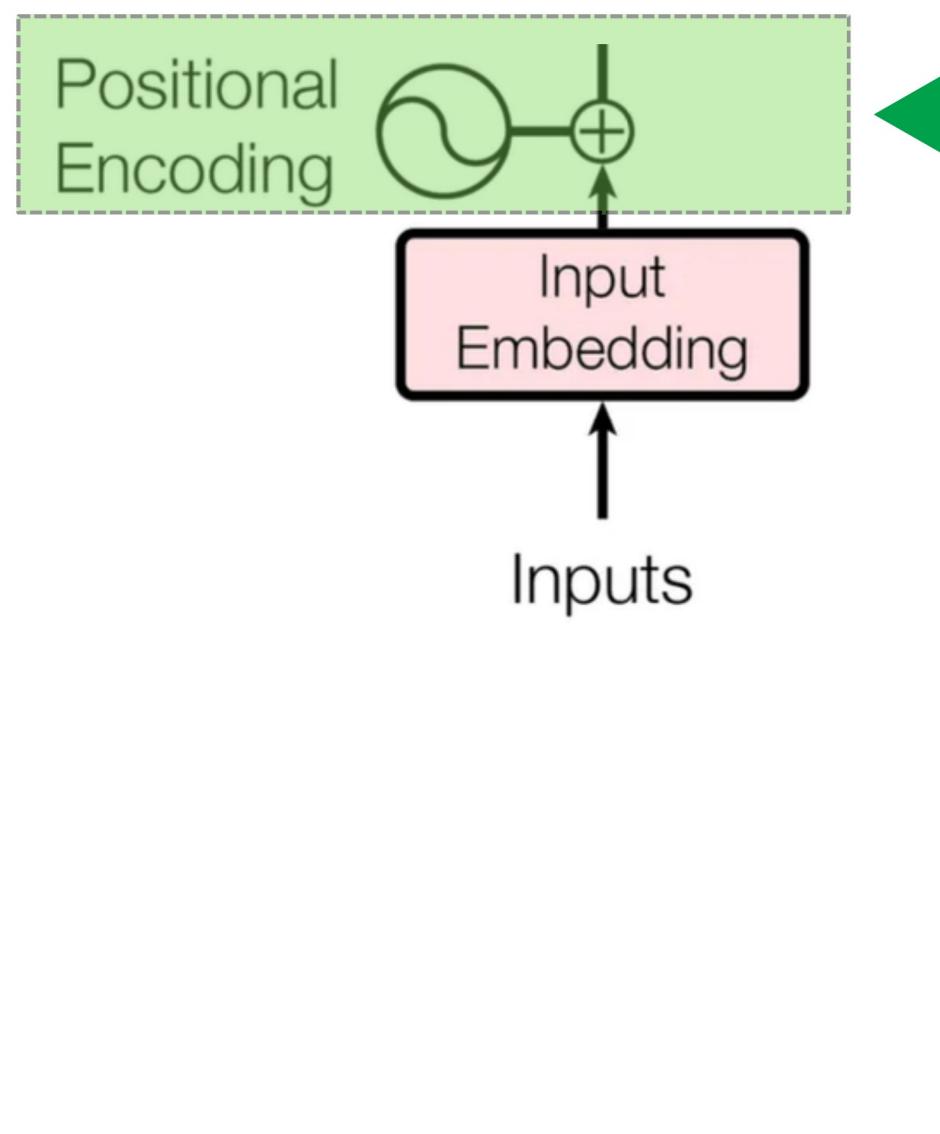


نحوه ترکیب اطلاعات ترتیبی (از مرحله ای Input Embedding) و بردار هر توکن (کلمه) (از مرحله ای Positional Encoding)



برای هر توکن (کلمه)، بردار position encoding با بردار embedding جمع می شود.

$$X = E + PE$$



I		0.1	0.3	0.4	0.8
really	0.9	0.1	0.5	0.2	
enjoy	0.7	0.6	0.1	0.4	
learning	0.2	0.8	0.3	0.7	
transformers	0.5	0.4	0.9	0.3	

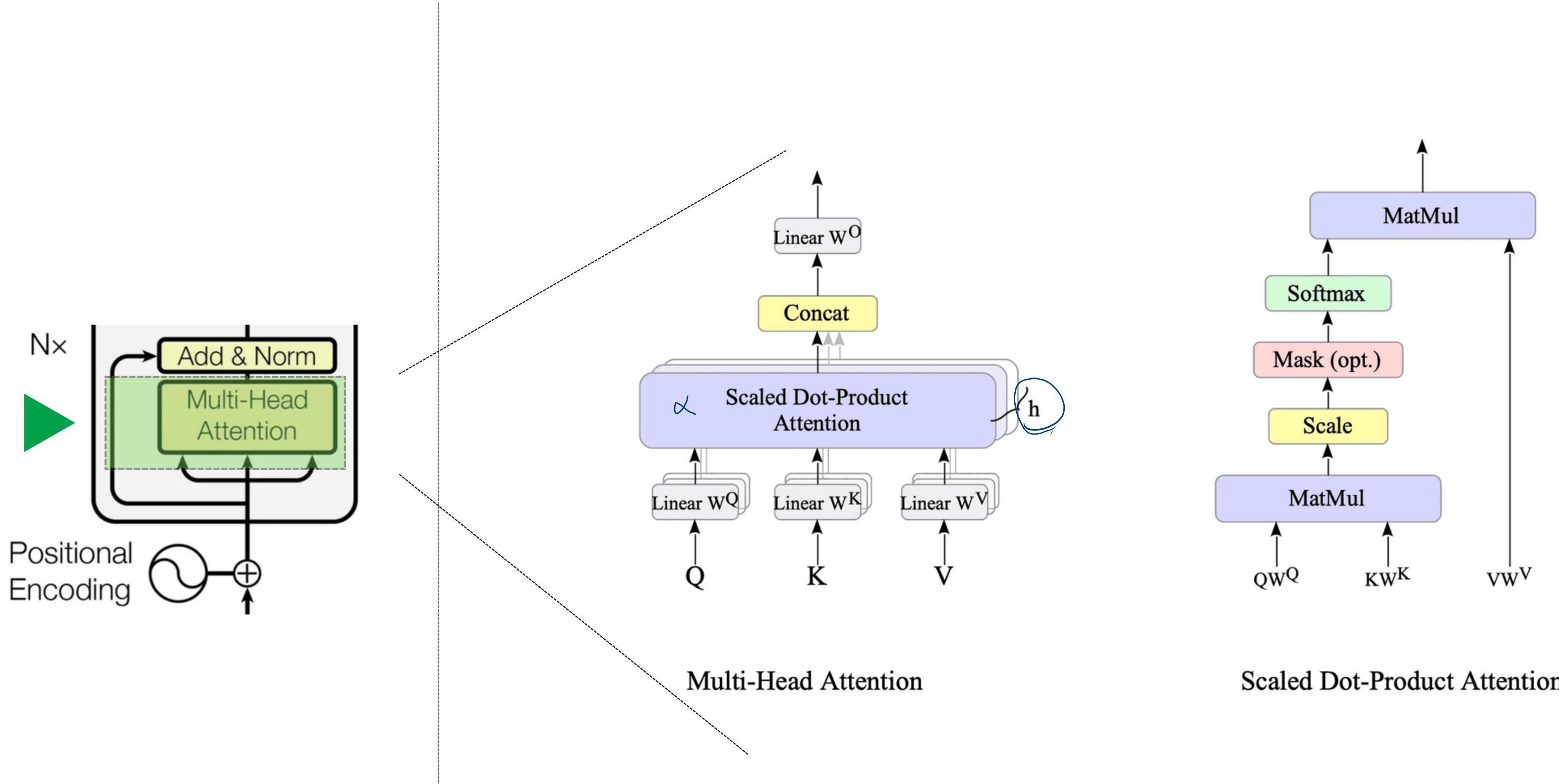
0.0 X	1.0	0.0	1.0
0.84	0.54	0.01	0.99
0	-0.41	0.02	0.99
0.14	-0.98	0.029	0.99
-0.75	-0.65	0.039	0.999

0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

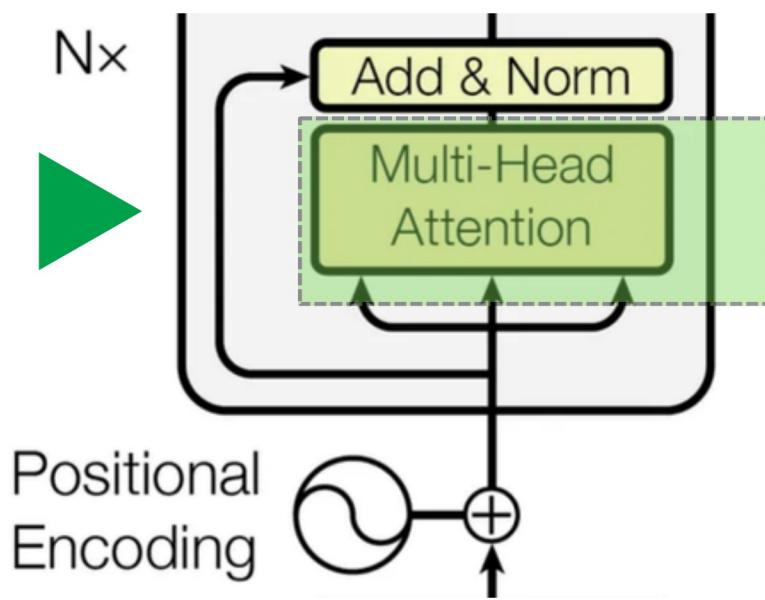
از این ماتریس (خروجی بخش (Positional Encoding مدل چگونه ترتیب کلمات را می فهمد؟

	X			
I	0.10	1.30	0.40	1.80
really	1.74	0.64	0.51	1.19
enjoy	1.60	0.18	0.12	1.39
learning	0.34	-0.18	0.32	1.69
transformers	-0.25	-0.25	0.93	1.29

نحوه ورود خروجی Multi Head Attention & Positional Encoding



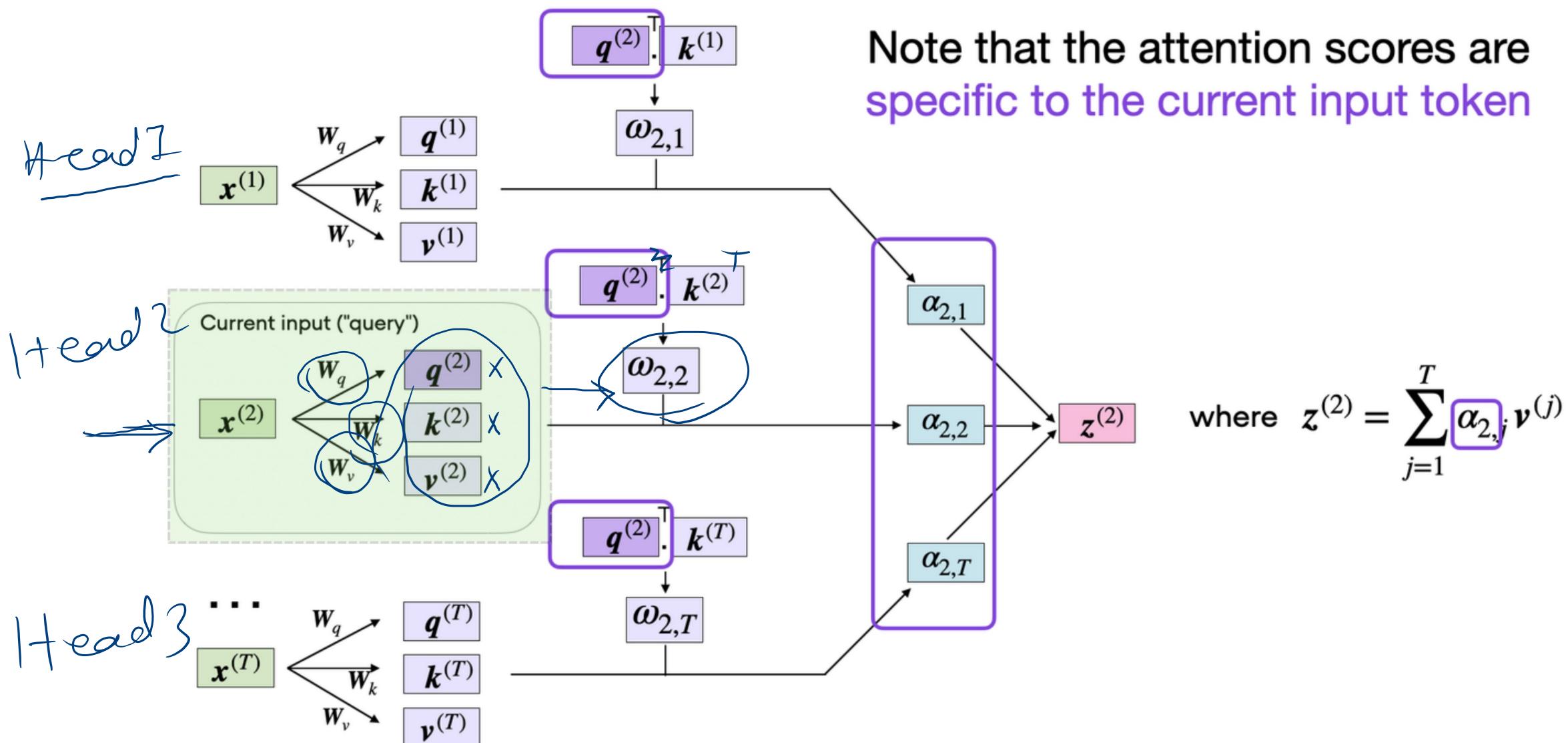
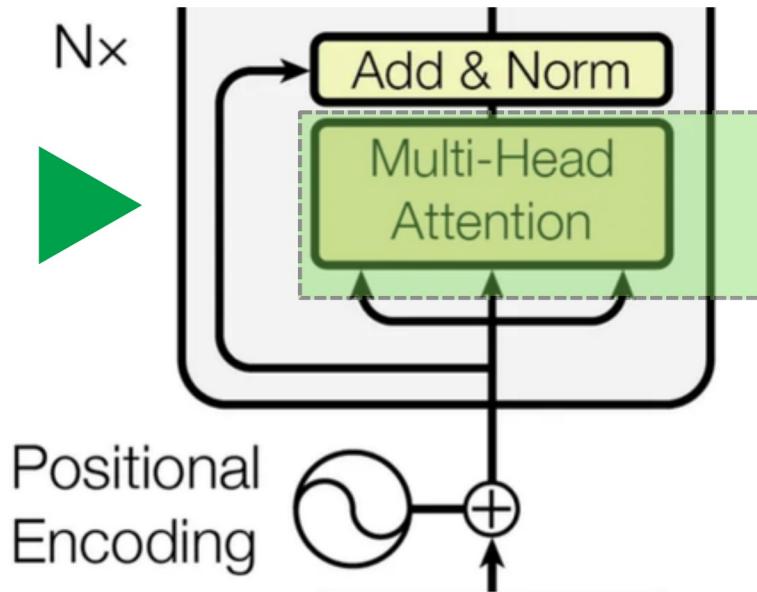
فرمول مکانیزم

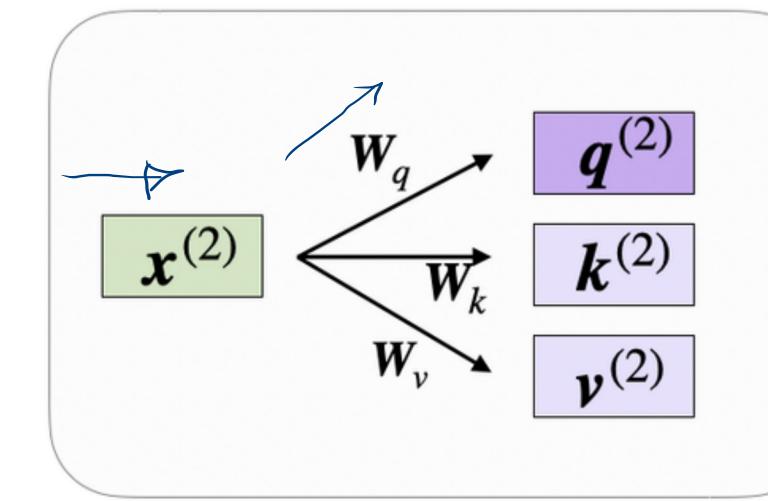
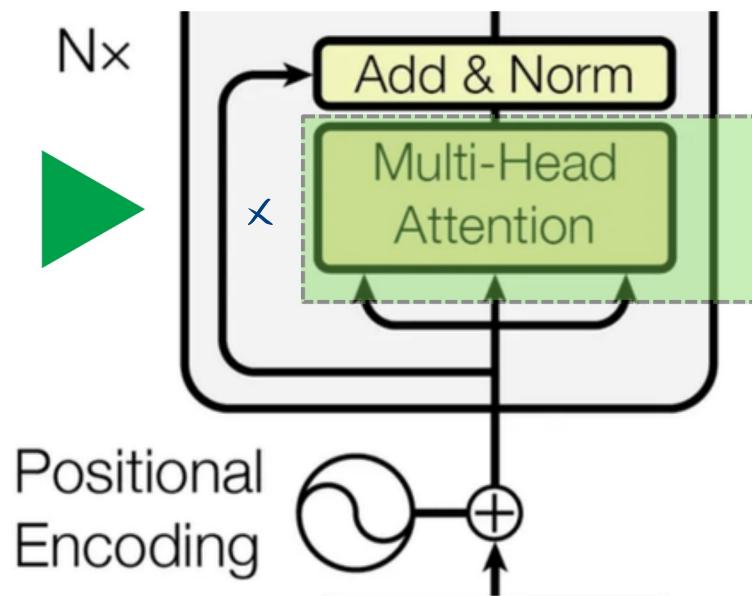


$$\text{Attention}(Q, K, V) = \underbrace{\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)}_{\alpha} V$$

where Q, K and V are Query, Key and Value vectors.

نمایش دیگری از معماری مکانیزم





X			
0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

X			
0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

X			
0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

محاسبه بردارهای Query, Key, Value مکانیزم Self Attention

$$\begin{matrix} X \\ * \\ W_q \\ = \\ Q \end{matrix}$$

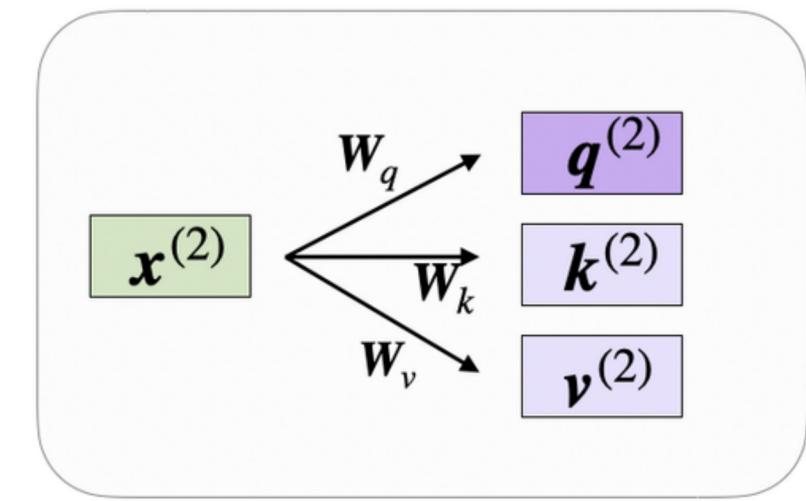
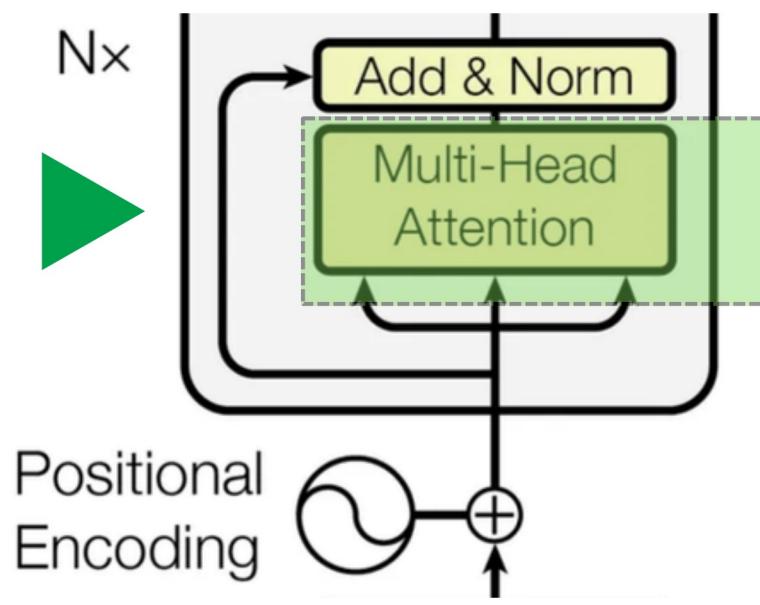
0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

$$\begin{matrix} X \\ * \\ W_k \\ = \\ K \end{matrix}$$

0.5	2.89	1.45	-0.44
3.96	1.08	1.31	-1.6
4.35	1.67	0.89	-1.0
4.35	1.82	0.22	0.43
4.39	0.43	-0.06	0.72

$$\begin{matrix} X \\ * \\ W_v \\ = \\ V \end{matrix}$$

0.17	0.01	-3.48	1.8
-1.9	-3.93	-3.23	2.01
-2.04	-3.75	-2.59	1.52
-1.75	-2.12	-2.74	2.08
-1.73	-1.7	-3.04	2.88



Q

0.10	1.30	0.40	1.80
1.74	0.64	0.51	1.19
1.60	0.18	0.12	1.39
0.34	-0.18	0.32	1.69
-0.25	-0.25	0.93	1.29

نمایش دیگری از معماری مکانیزم Self Attention



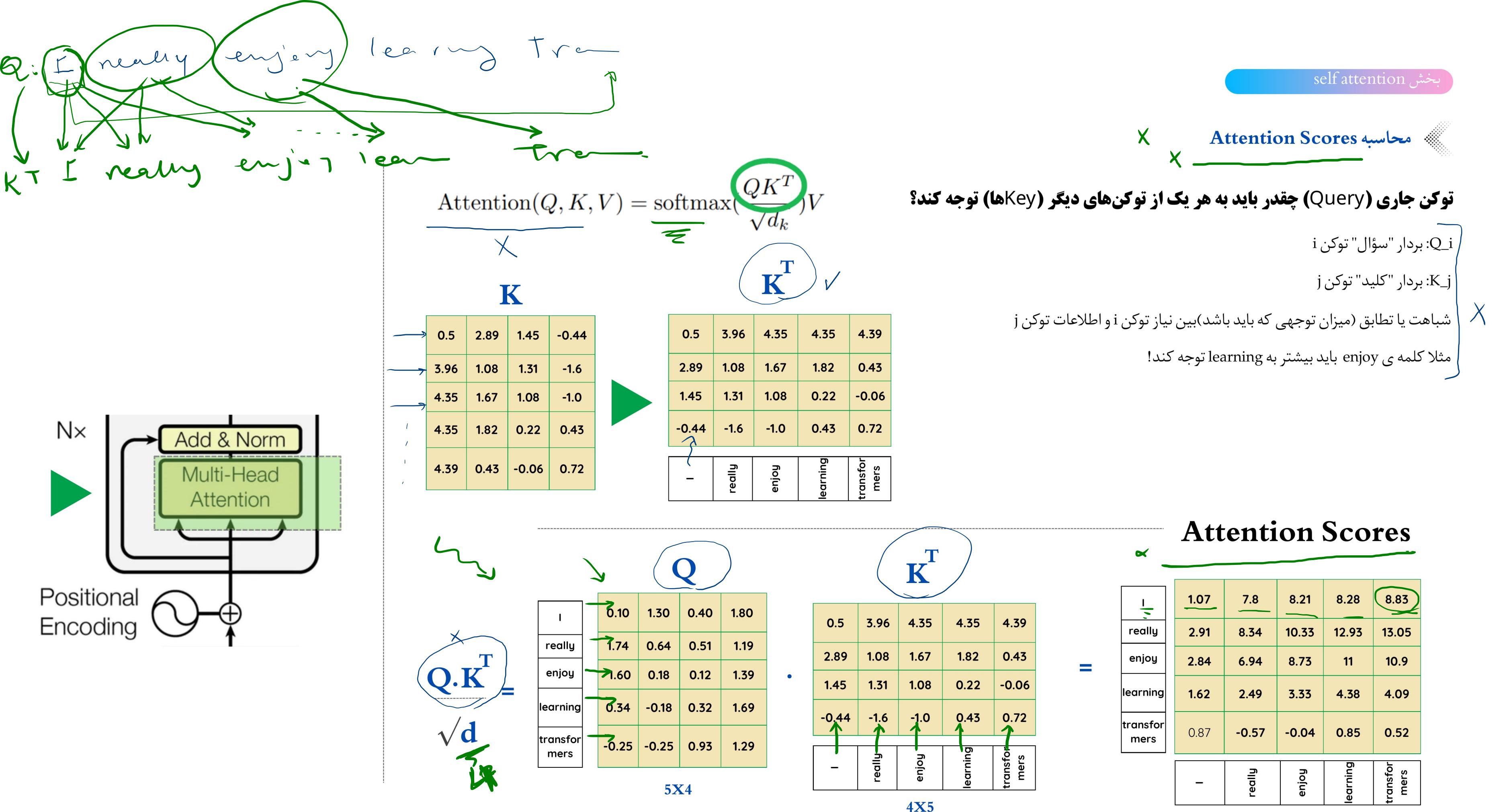
K

0.5	2.89	1.45	-0.44
3.96	1.08	1.31	-1.6
4.35	1.67	0.89	-1.0
4.35	1.82	0.22	0.43
4.39	0.43	-0.06	0.72

V

0.17	0.01	-3.48	1.8
-1.9	-3.93	-3.23	2.01
-2.04	-3.75	-2.59	1.52
-1.75	-2.12	-2.74	2.08
-1.73	-1.7	-3.04	2.88

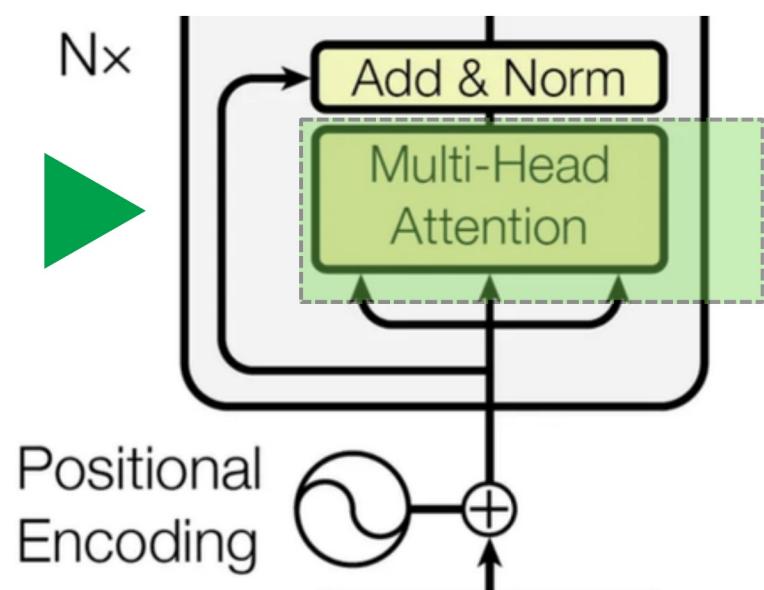
(5x4)



محاسبه Attention Scores

توکن جاری (Query) چقدر باید به هر یک از توکن‌های دیگر (Key‌ها) توجه کند؟

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



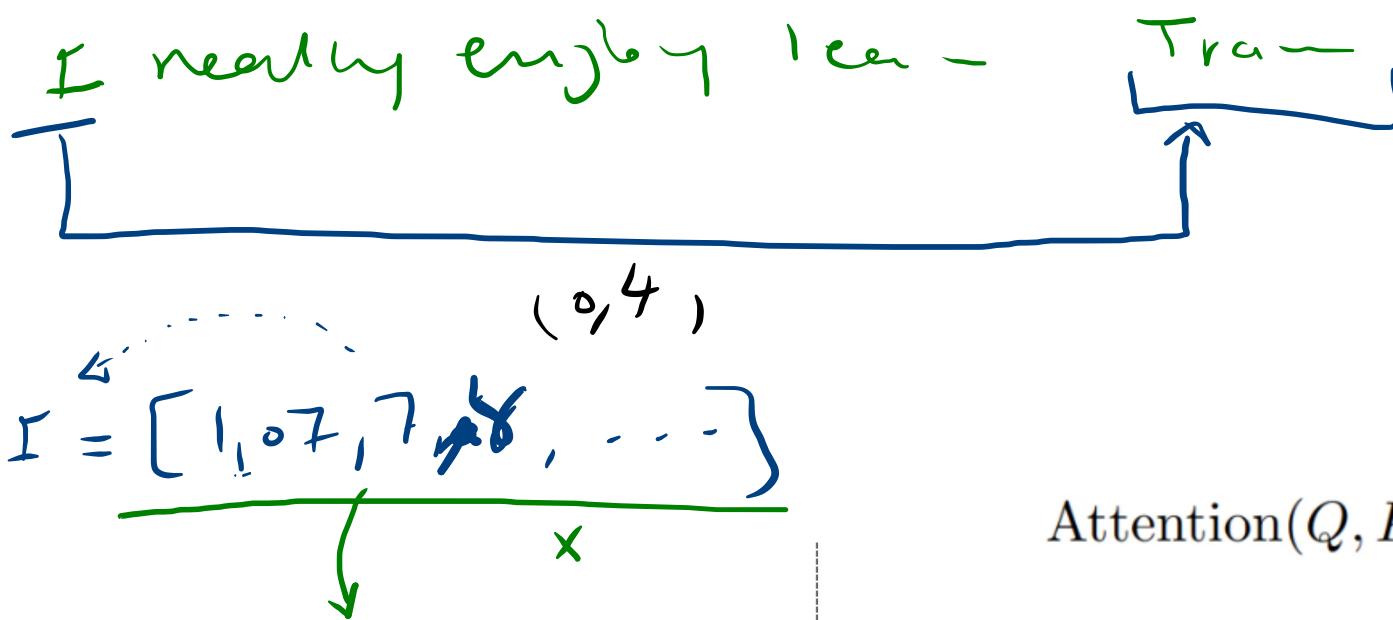
$$\text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d}}\right) = \text{softmax}(A)$$

1.07	7.8	8.21	8.28	8.83
2.91	8.34	10.33	12.93	13.05
2.84	6.94	8.73	11	10.9
1.62	2.49	3.33	4.38	4.09
0.87	-0.57	-0.04	0.85	0.52

Attention Weights

1	2	3	4	5
0.0002	0.1443	0.2185	0.234	0.403
0.0	0.0046	0.0335	0.4509	0.5111
0.0001	0.0085	0.0508	0.495	0.4456
0.0274	0.0655	0.1517	0.4332	0.3222
0.3	0.0716	0.1215	0.2945	0.2124

Context Query Position

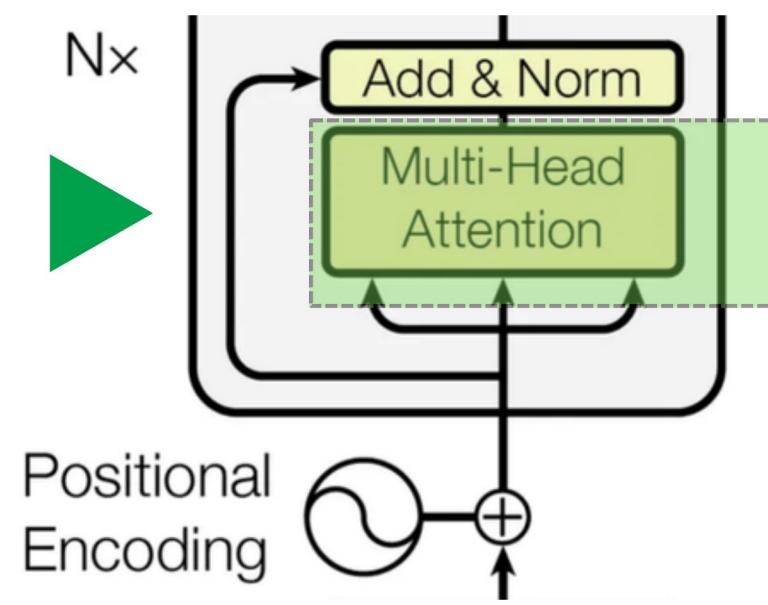


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

محاسبه Attention Output

نتیجه نهایی مکانیزم Self-Attention

خروجی Attention برای هر توکن، یک بردار جدید است که ترکیبی از اطلاعات بقیه توکن‌هاست، وزن دار شده با توجهی که آن توکن به دیگران داشته است.



Attention Weights

Σ

$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V =$

Σ	0.0002	0.1443	0.2185	0.234	0.403
Γ	0.0	0.0046	0.0335	0.4509	0.5111
Σ	0.0001	0.0085	0.0508	0.495	0.4456
Σ	0.0274	0.0655	0.1517	0.4332	0.3222
Σ	0.3	0.0716	0.1215	0.2945	0.2124

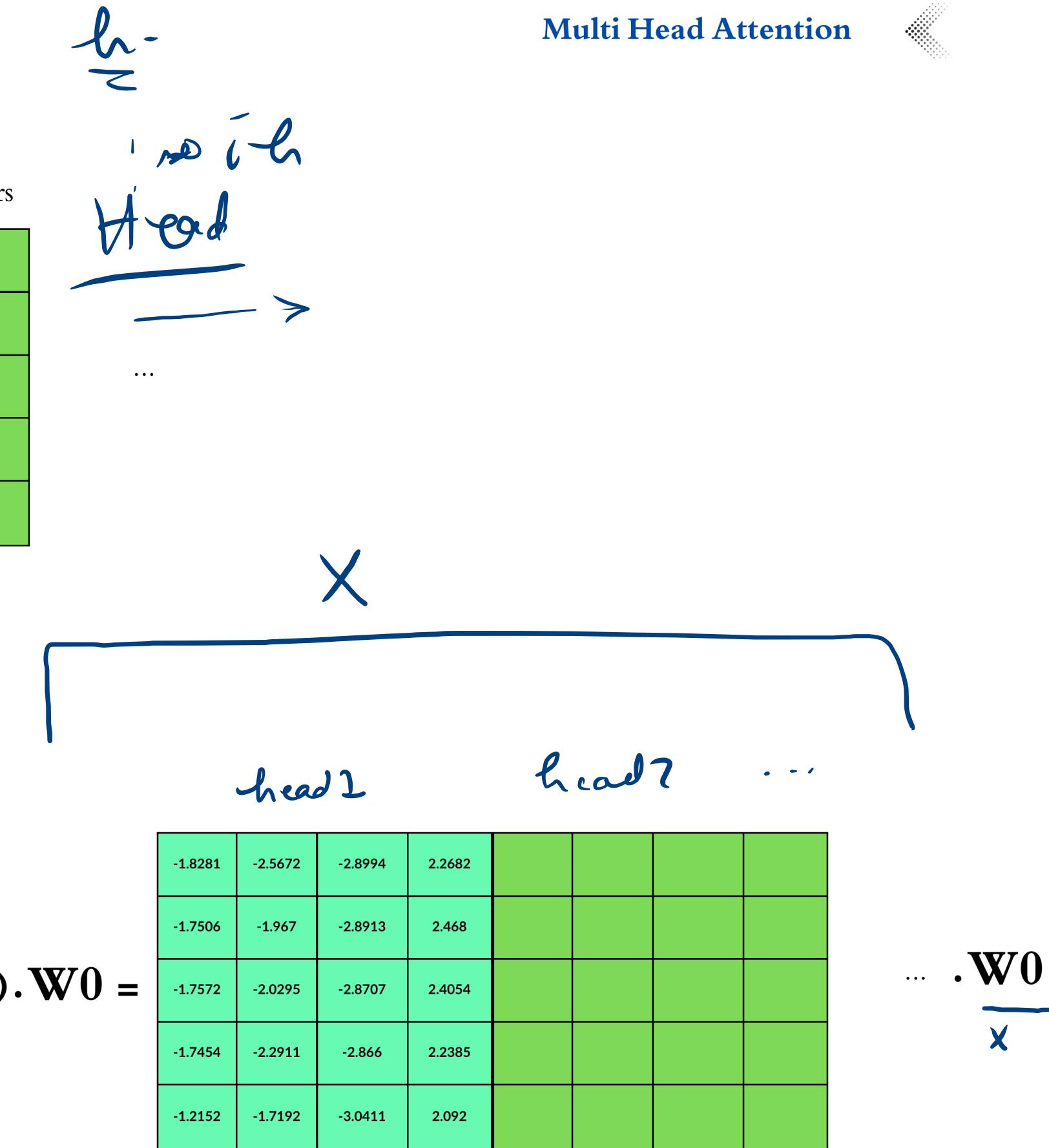
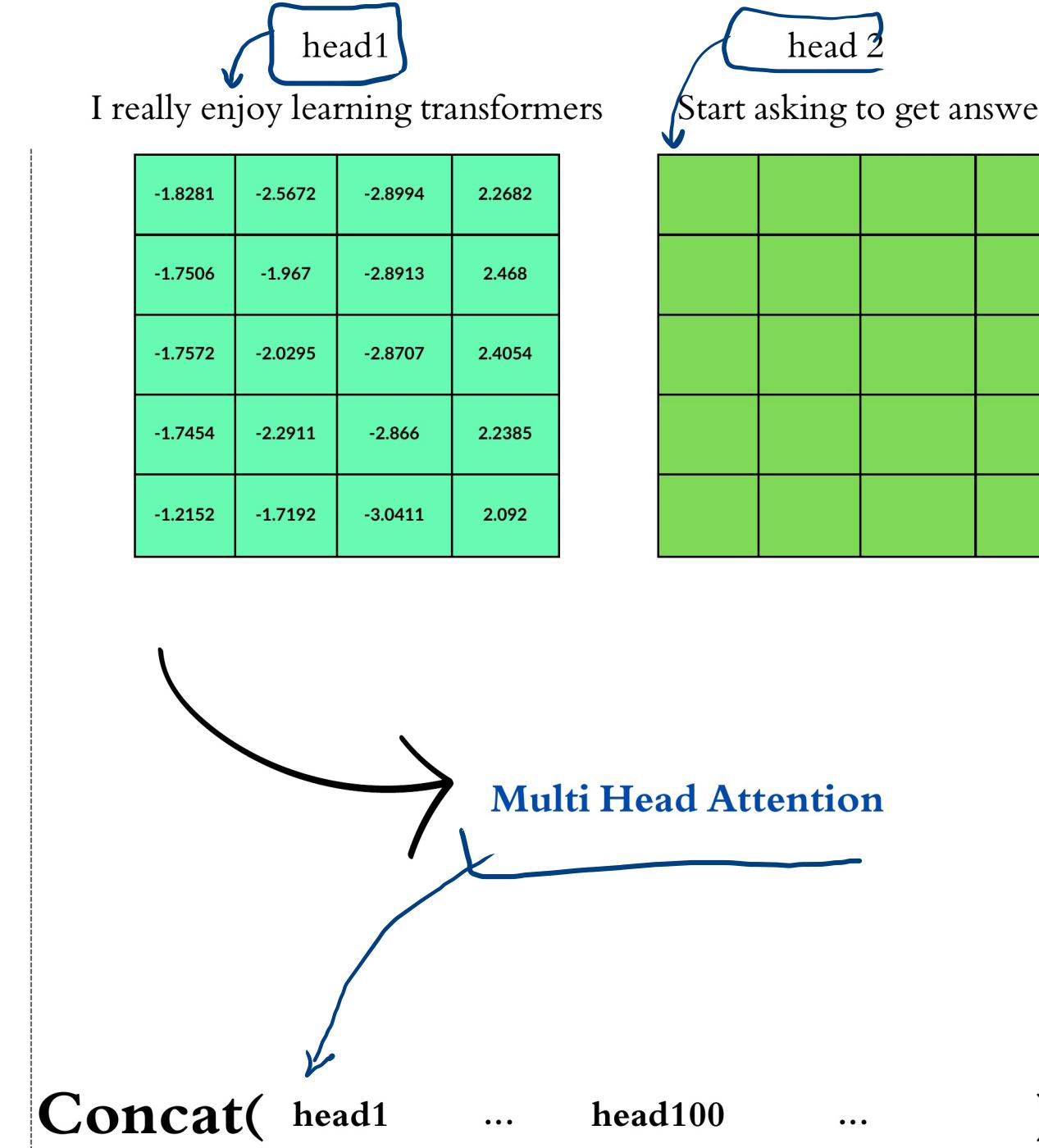
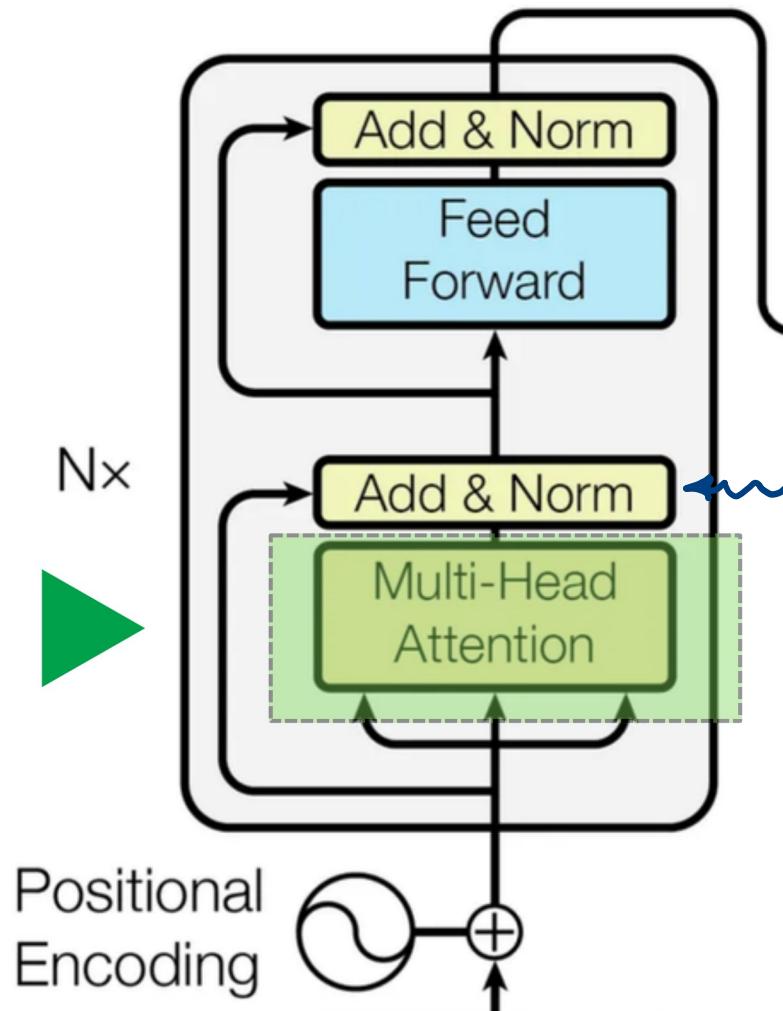
$\text{Train} \times V$

X	0.17	0.01	-3.48	1.8
X	-1.9	-3.93	-3.23	2.01
X	-2.04	-3.75	-2.59	1.52
X	-1.75	-2.12	-2.74	2.08
X	-1.73	-1.7	-3.04	2.88

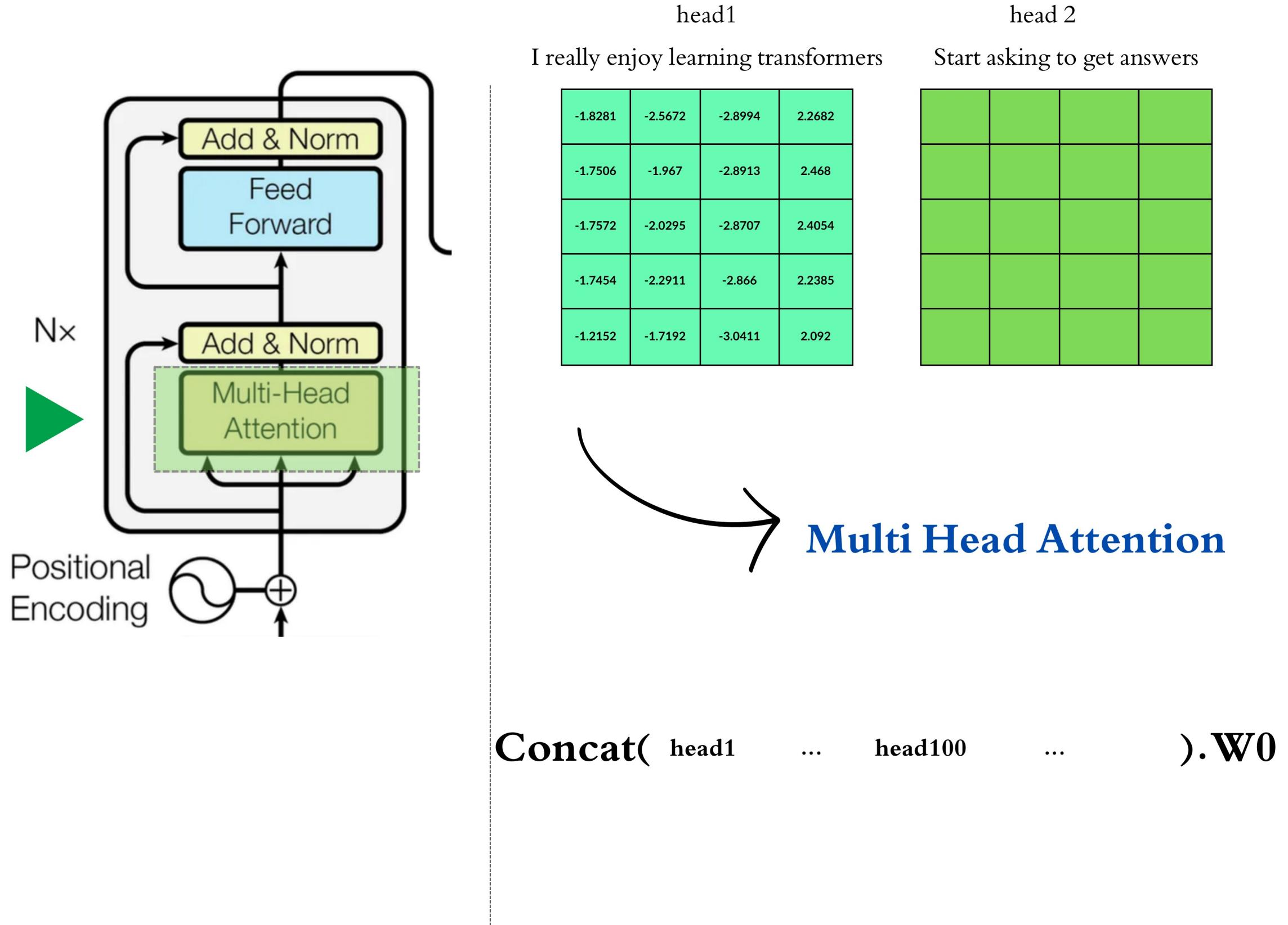
=

V	-1.8281	-2.5672	-2.8994	2.2682
V	-1.7506	-1.967	-2.8913	2.468
V	-1.7572	-2.0295	-2.8707	2.4054
V	-1.7454	-2.2911	-2.866	2.2385
V	-1.2152	-1.7192	-3.0411	2.092

attention output



Multi Head Attention



-1.8281	-2.5672	-2.8994	2.2682				
-1.7506	-1.967	-2.8913	2.468				
-1.7572	-2.0295	-2.8707	2.4054				
-1.7454	-2.2911	-2.866	2.2385				
-1.2152	-1.7192	-3.0411	2.092				

... .W0

