

Why Transformer?

1. All the tokens in parallel
2. Handling long Set of tokens was challenging with other architecture
3. A word can be represented in different context in different place

I went to bank to withdraw money
I went to river bank to get ~~some~~ fresh air

I have interest in biking
The interest in the loan is increasing

The embeddings may not be 100%
Correct representation

we update embedding based on
the input passed

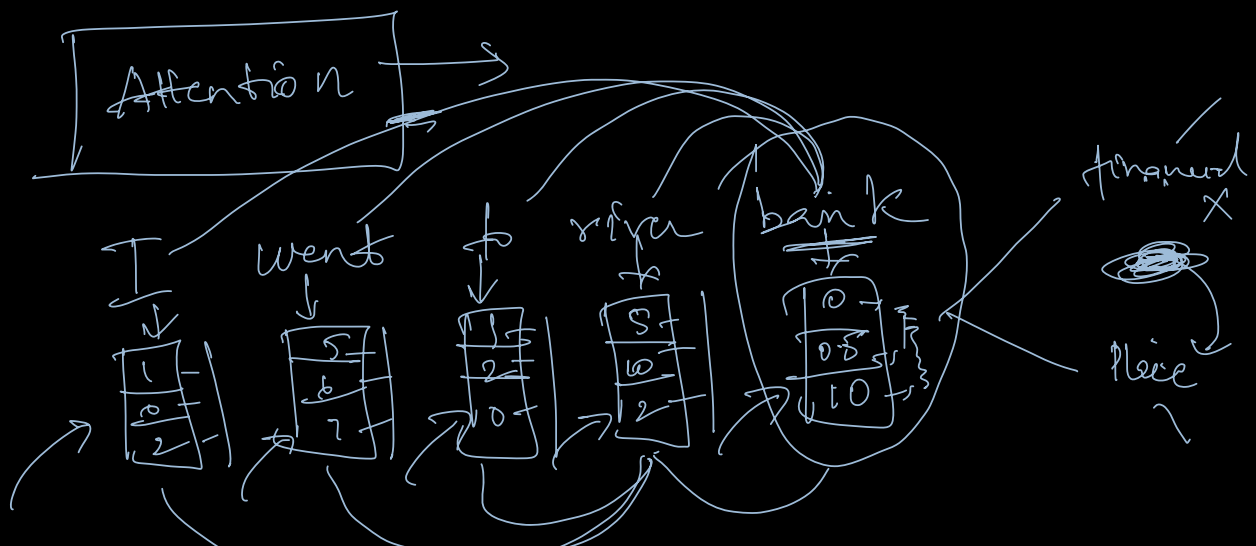
4.

Diffusion model ←

Image → Image
text → Image
Image → Text
Voice → Text
Text → Voice

2017 → Attention is all you need ←

Transformers



Updating the embeddings based on the Input passed by comparing the other tokens

How?

Similar to cosine similarity

Scaled dot product

$$X_1 = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}, X_2 = \begin{bmatrix} 0 \\ 1 \\ 5 \end{bmatrix}$$

$$X_1 \cdot X_2 = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 5 \end{bmatrix} = 0 + 0 + 10 = 10$$

$$\Rightarrow \frac{10}{\sqrt{3}}$$

Similarity score to tokens

$\text{len}(X_1) = 3$

	1	went	to	river	bank
1	1	0.5	0.6	0.1	2.0
went		1			
to			1		
river				1	
bank					1

$$I = 0.5 \text{ went} + 0.6 \text{ to} + 10 \text{ river} + 20 \text{ bank} = 1$$

Softmax \rightarrow

\downarrow

$= 0.01 \text{ went} + 0.01 \text{ to} + 0.28 \text{ river} + 0.7 \text{ bank}$

\rightarrow to represent how each word is influencing the other word

Update each word embedding based on the influencing factor

How?

Key Queries Values

Multi Head Attention

why?

Decision tree
Random Forest
Stacking

Linearly transform the embedding

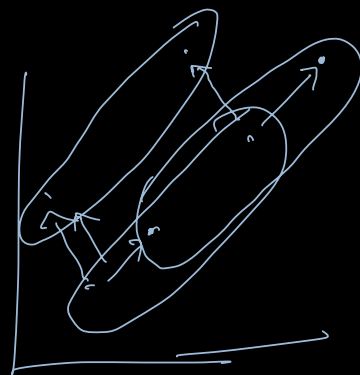
0
1
2

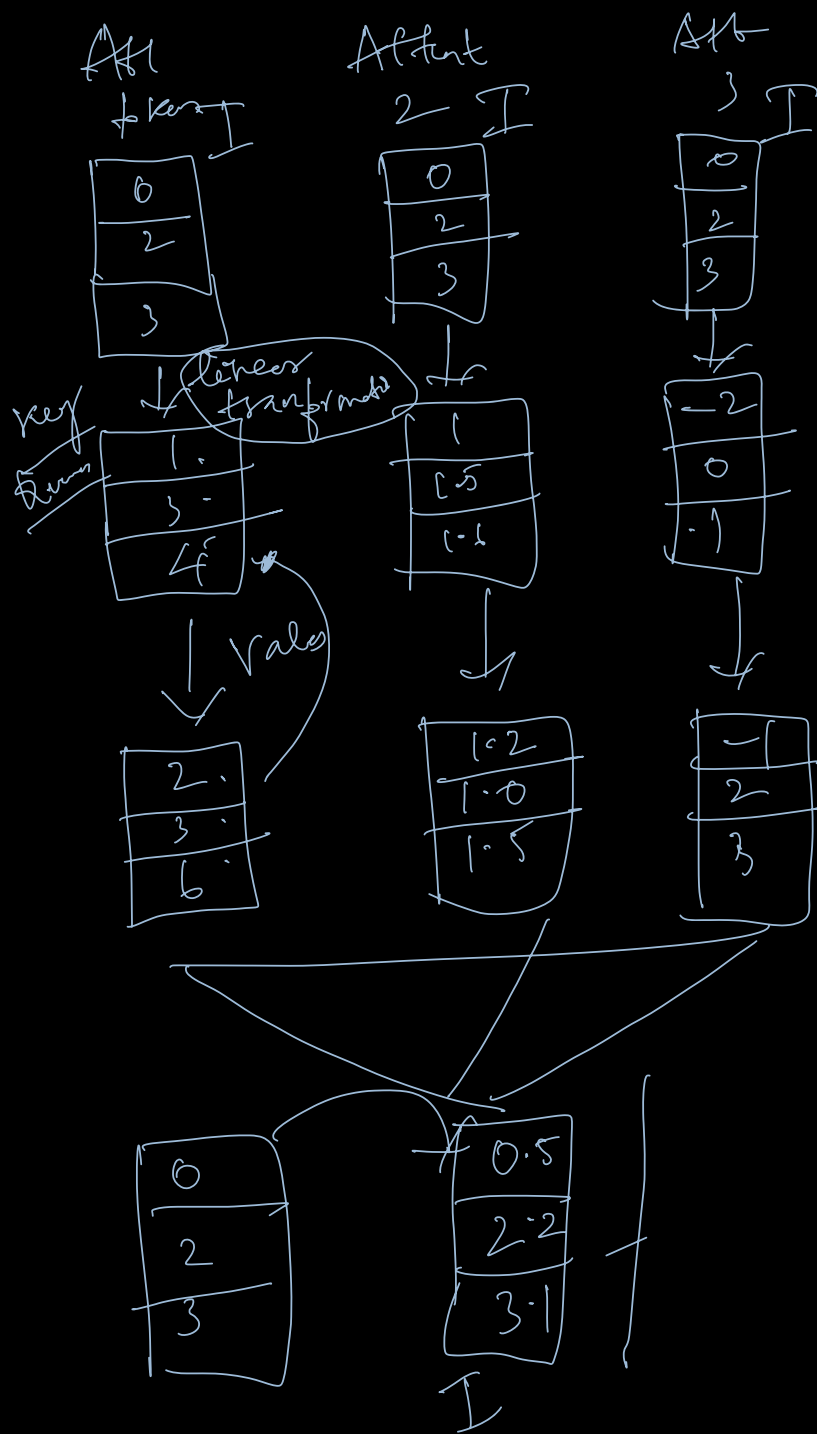
+1

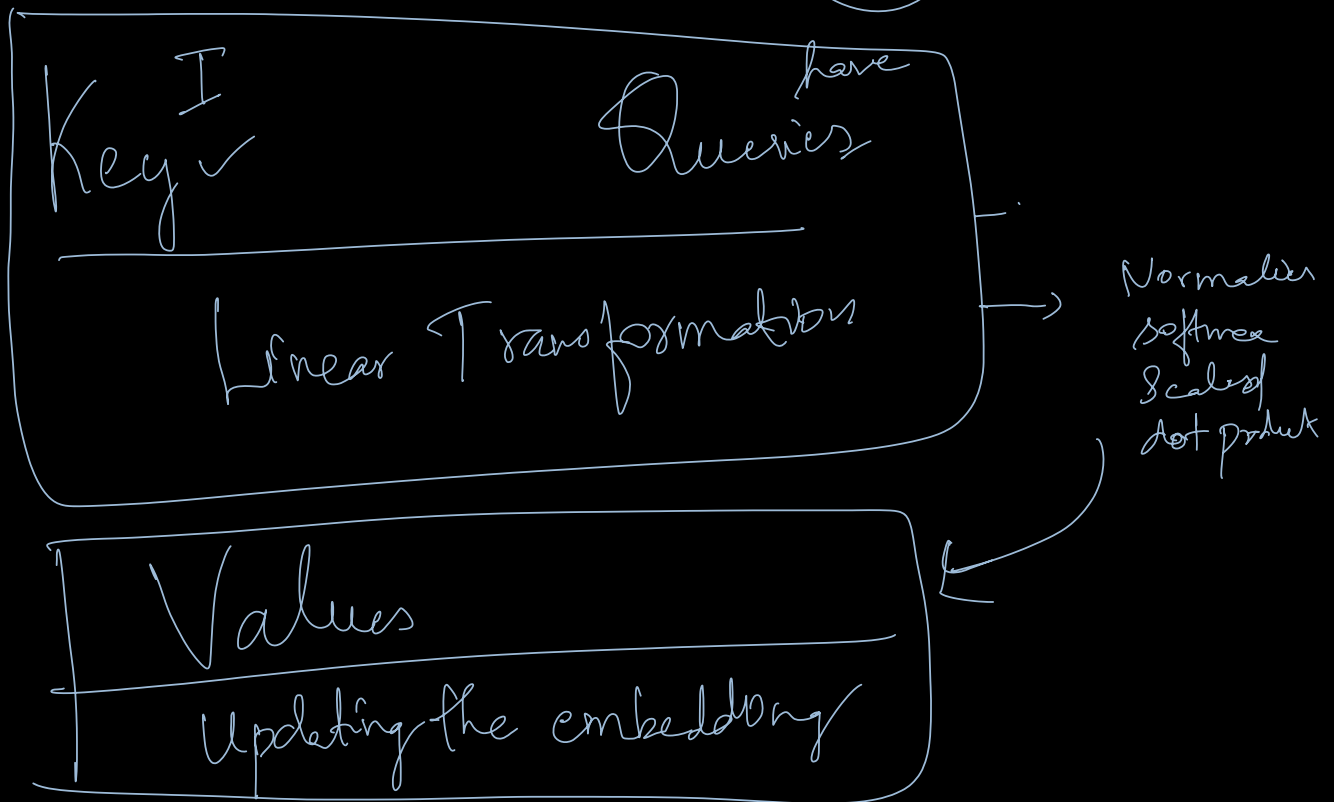
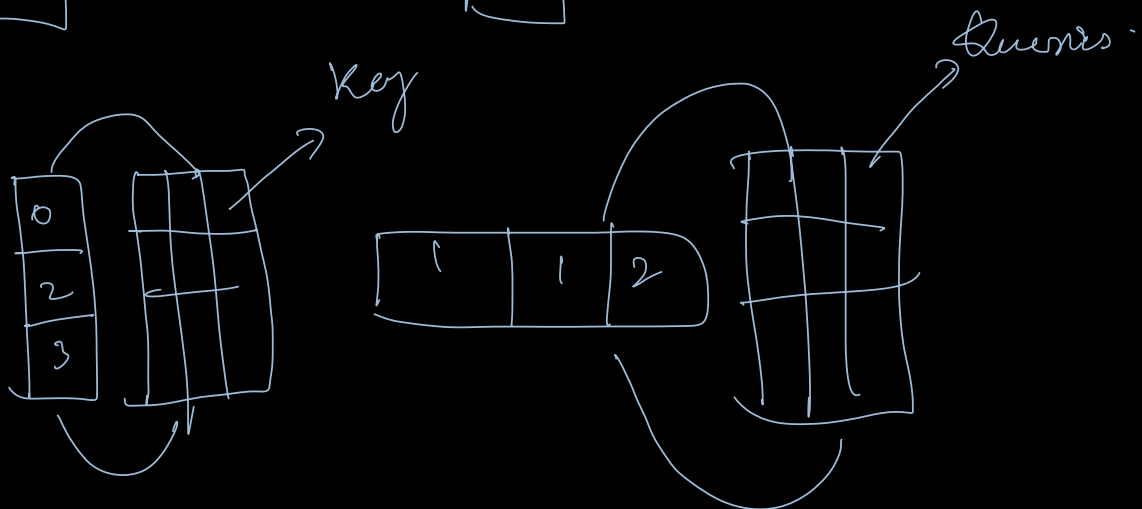
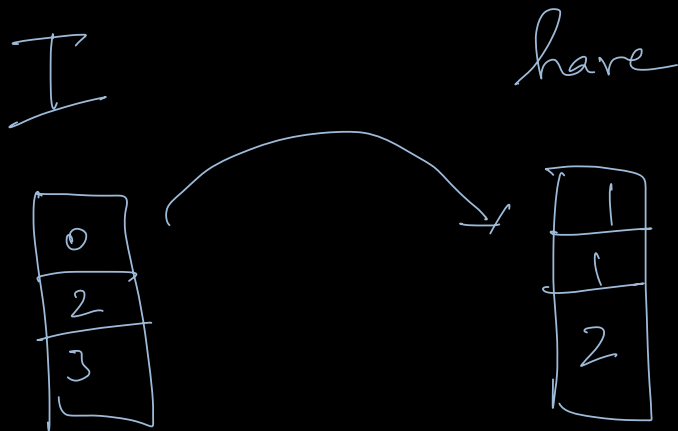
+2

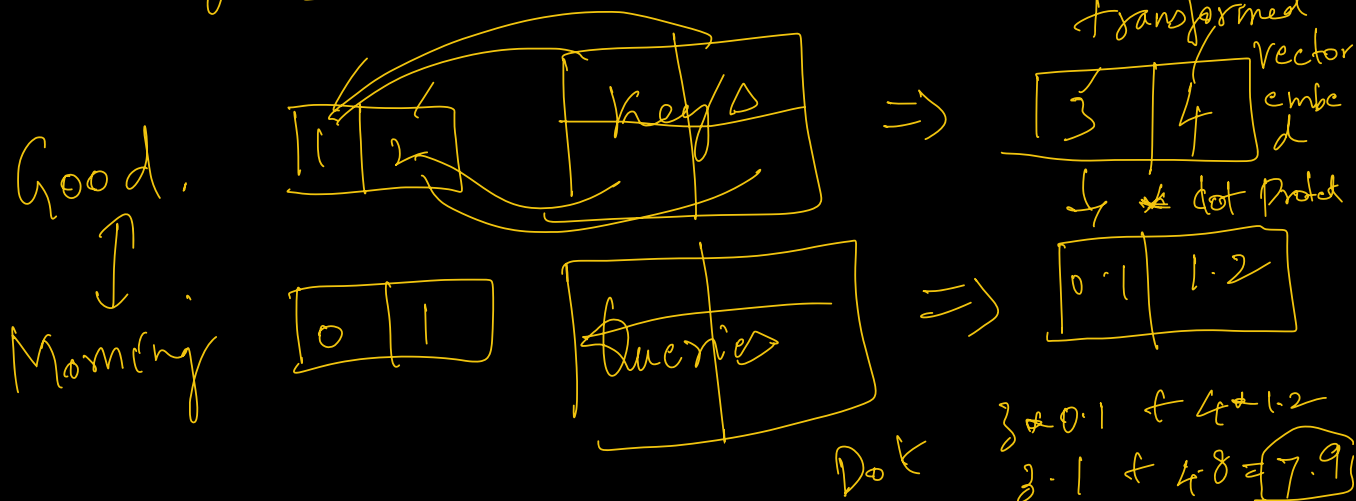
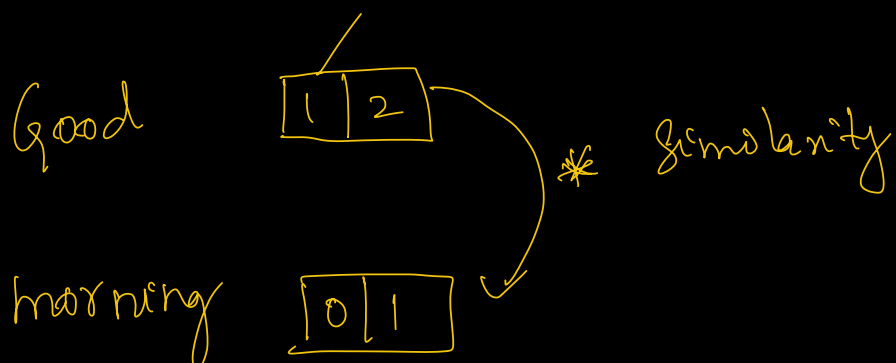
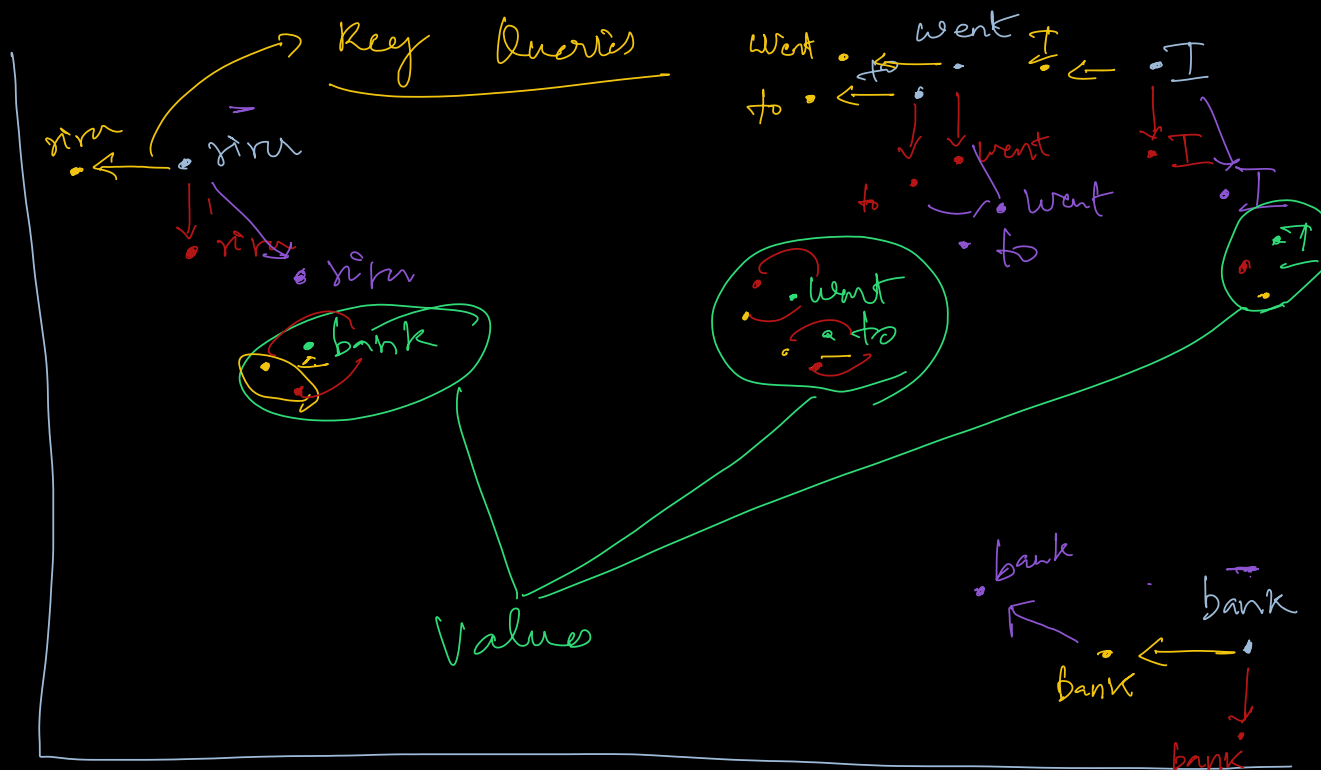
+3

1
2
3









2x2

	good	morning	
good	1	6.8	Scaled Dot product = $\frac{7.9}{\sqrt{2}} = 6.8$ ✓
morning	6.8	1	

Score

	good	morning
good	0.3	0.7
morning	0.7	0.3

of max.

Normalize the similarity with addition of vector becomes

①

Multiply the values to get new embedding

0.3 0.7
0.7 0.3

Values

1

Good Morning	<div>1 3</div>	<div>3 9</div>	am	Good	<div>1 3</div>	<div>15 45</div>
	Good	morning	.	l	an	Good
Good	1	2	.	3	4	5
morning						
l						
am						
Good						

Positional encoding

odd 1, 3, 5, 7, ... cos

even 2, 4, 6, 8, ... sin

odd $\rightarrow \frac{3-i}{2}$ and position
even $\rightarrow \frac{2-i}{2}$

Good

1
3