

Transformers → chapter 2

ANN → Tabular data

CNN → Image data

RNN → Sequential text

Neural
networks
Architecture

Transformer → Sequence to sequence task

ex) Translation

→ Question - Ans

→ Text summarisation

→ parallel processing

→ can train over a big data

→ Use self-attention

→ Induction in → Attention all you need

→ In 2017

→ Impact

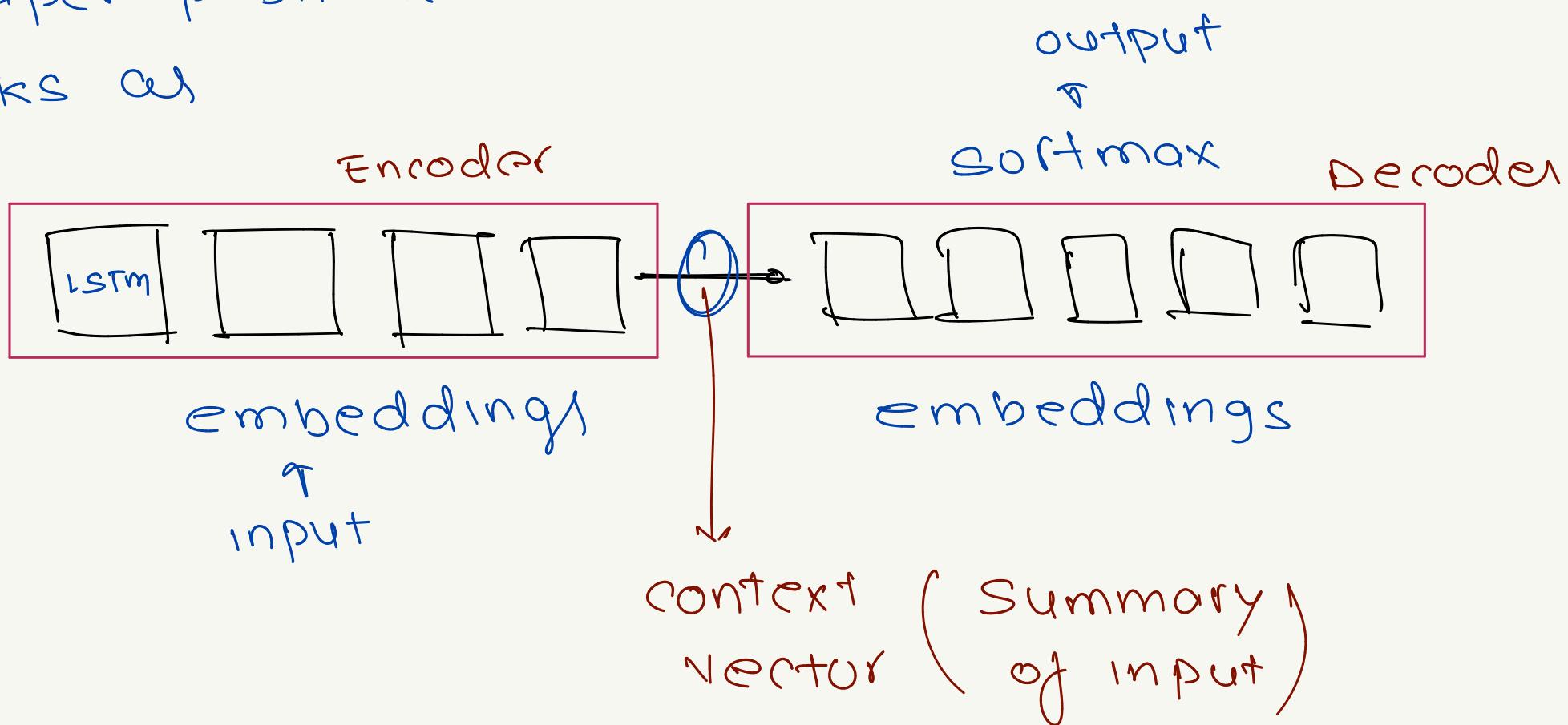
- Revolutionalized NLP (natural language processing)

- Democratizing AI
 - ↳ BERT
 - ↳ GPT
- Multimodel compatibility
 - ↳ can take text, image etc as input.
- Acceleration of Gen AI
 - generated text, image etc
- Unifying of deep learning.
 - ↳ Transformer can be used in place of NLP, Gen AI, Reinforcement.

Origin of Transformer → Papers

① Sequence to sequence Learning with Neural network

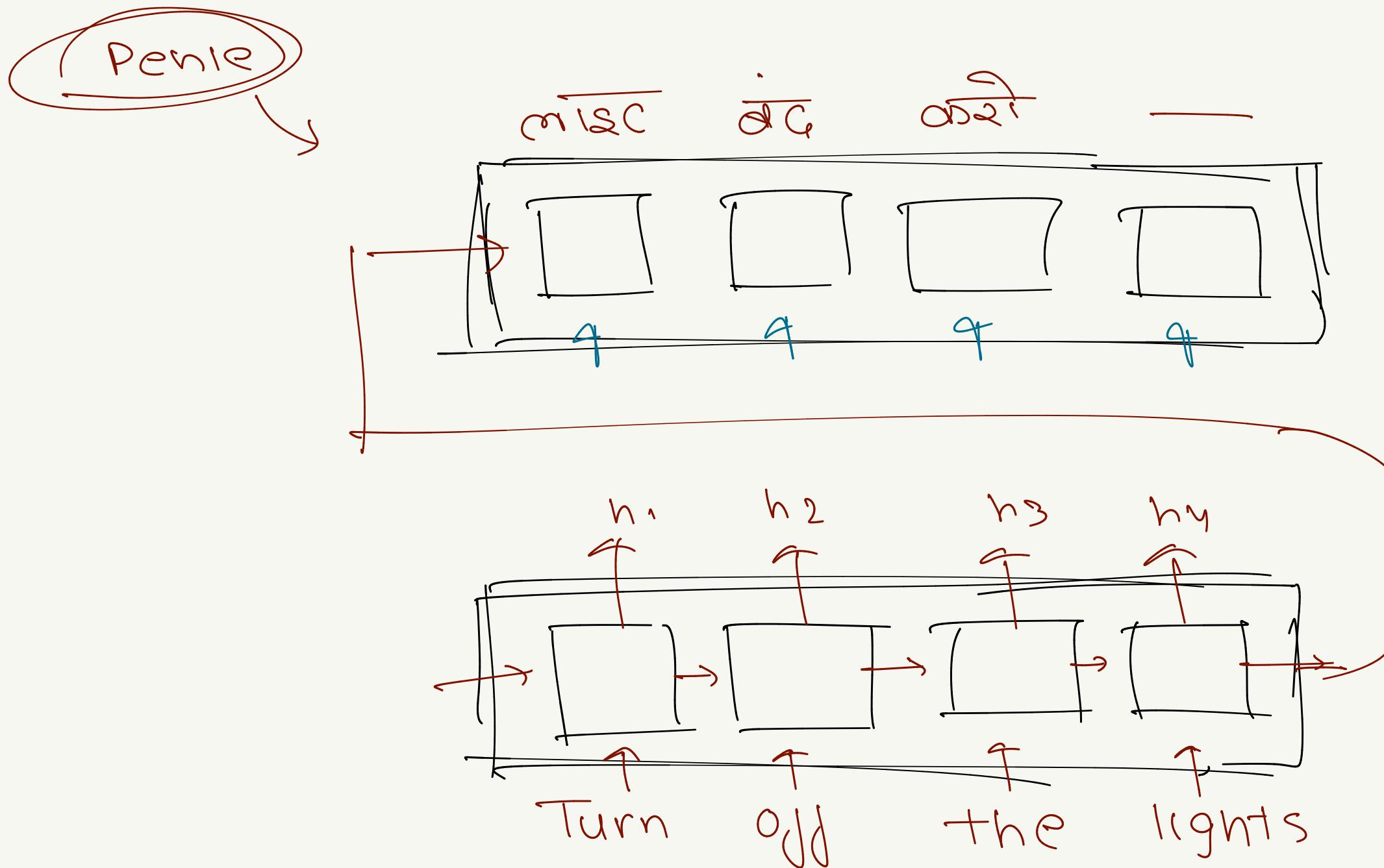
- A paper published in 2015 works as



- This architecture works perfectly upto 30 words but after that it fails to do properly.
- Flow → hum input ko context vector mein summezze karna chante hain par agar input bht bda ho gya to context vector sab acche se store nhi kar paayega.

② Neural machine translation by jointly learning to Align & Translate -

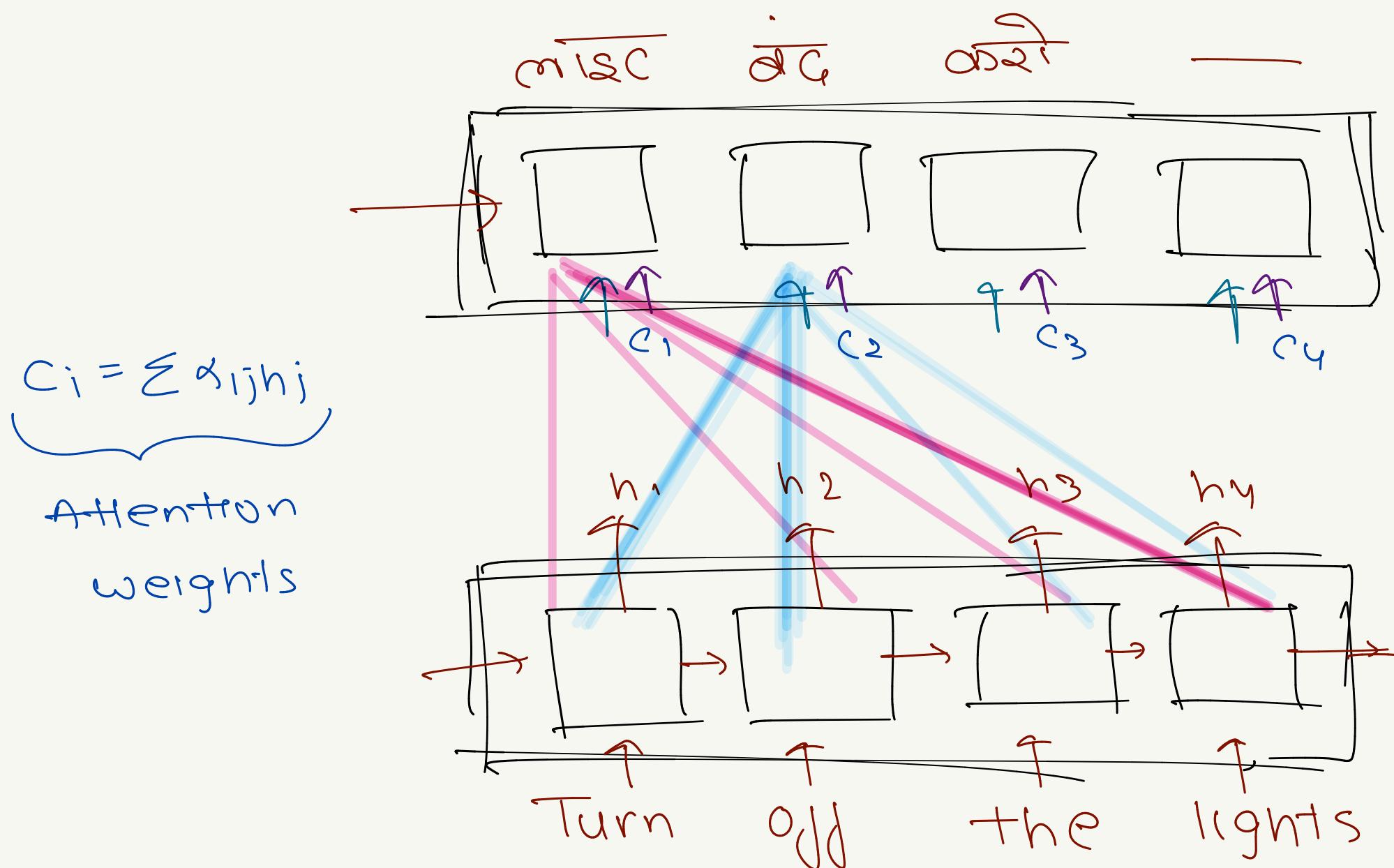
- first time introduced Attention



Penie hum saara ka saara line ke context vector bhejte the decoder mein

Abb

Hume saare input ek soath bhejne ki jarurat nahi hai.



$$c_i = \sum \alpha_{ij} h_j$$

Attention weights

Overall context vector is weightage sum of encoder ka hidden state ka.

Attention pe batata hai kaunse se input ka weightage jyada hogya kiss word ko print karne ke liye.

→ This architecture is perfectly working on even more than 30 words.

Flaw → sequential training
so very slow
can't train big data set. }
not parallel training
→ Bcz of ISPM

③

Attention all you need →

- In 2017
- No RNN or LSTM
- works on self-attention
- can be trained parallelly
- stable hyperparameters

2000 - 2014 → RNN / LSTM

2014 → Attention

2017 → Transformer

2018 → BERT / ViPT

2018 - 2020 → Vision Transfer / AlphaFold-2

2021 → Gen AI

2022 - 2024 → Chat ViPT / Stable diffusion

Generally research papers are incremental. A small change in previous research paper.

But

This paper was totally different from all older research papers.

Advantages of Transformer

- 1) highly scalable → parallel training
- 2) Transfer learning → *
- 3) multimodal Input/Output → image \leftrightarrow text
audio \leftrightarrow text
etc
image
- 4) Flexible architecture
 - ↳ can change small things
 - things can create our own architecture
 - BERT → encoder only
- 5) Ecosystem → many library (flair, hugging face)
 - & tools
 - & blogs & videos
- 6) Easily can be integrated in other AI.

Applications

① Chat GPT

② DALL-E2

③ AlphaGo/AlphaGo 2

④ openAI codeX

⑤ GitHub copilot

etc

Disadvantages

- ① High computational tools
- ② Data (requires lot of data)
- ③ Overfitting (maybe)
- ④ Energy consumption
- ⑤ Interpretability (its a blackbox model)
(not secure)
- ⑥ Bias → Ethical concerns

Self Attention

Q) What is self attention?

for NLP \Rightarrow words \rightarrow numbers (vectorization)

computer only understands
numbers

- OHE \Rightarrow One hot encoding

\rightarrow understands unique words

- Example \rightarrow
- 1) cat, Rat, cat
 - 2) mat, mat, Rat

	mat	rat	cat	
mat	1	0	0	$\rightarrow 100$
Rat	0	1	0	$\rightarrow 010$
cat	0	0	1	$\rightarrow 001$

So numerical representation of

$$(\text{cat } \text{Rat } \text{mat}) \rightarrow [100][010][001]$$

• BOW (Bag of words)

Example

1) cat mat mat

2) cat rat rat

→ Also counts number of occurrence.

BOW of 1) → [cat Rat mat]
 ↓ ↓
 [1 0 2]

BOW of 2) → [cat rat mat]
 ↓ ↓ ↓
 [1 2 0]

⇒ word embeddings

↳ have power to capture semantic meanings

what is the context & how
generally this word is used

→ How to create →

- collection of large amount of training data
- Sent to neural network
- It understands meaning & context
- & represent each word into n-dimensional vector. (n can be 64, 256, 512 etc)
- After this we can represent words into set of 5 numbers.

Example

$$\text{King} = [0.6 \quad 0.5 \quad 0.2 \quad 0 \quad 0.9] \quad \text{Similar}$$
$$\text{Queen} = [0.7 \quad 0.6 \quad 0.9 \quad 0.01 \quad 0.5]$$

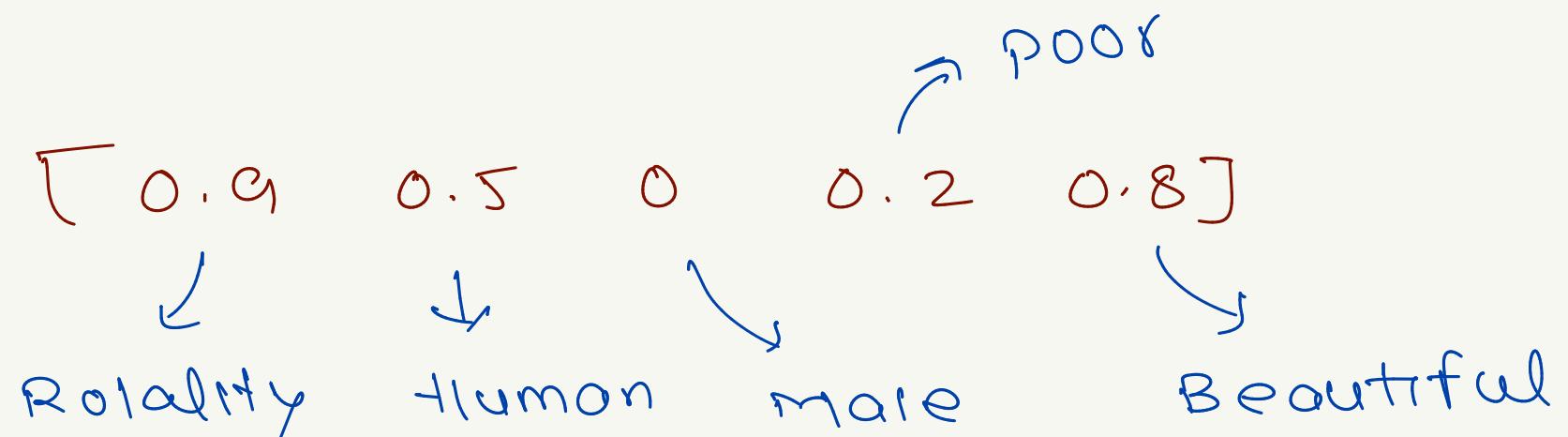
$$\text{Mouse} = [0.5 \quad 0.9 \quad 0.6 \quad 0.08 \quad 0.9] \quad \text{Different}$$

So, if two similar words have same vectors
then similar hence.

- vector ke har dimension ek parameter ko show karte hain.

Example

$$\text{Queen} = [0.9 \ 0.5 \ 0 \ 0.2 \ 0.8]$$



Just every dimension has a meaning.

- Embedding → word meaning → converts into vector

Problem

Example →

- 1) An apple a day keeps doc away.
- 2) Apple is healthy.
- 3) Apple is better than orange.
- 4) Apple makes good phone.

we will be
representing
all the words
into 2D $[x, y]$
 $[x, y] \rightarrow [Taste, Tech]$

word embedding captured → Average meaning

for sentence

- (1) $[0.6 \ 0]$
- (2) $[0.7 \ 0]$
- (3) $[0.8 \ 0]$
- (4) $[0.8 \ 0.2]$

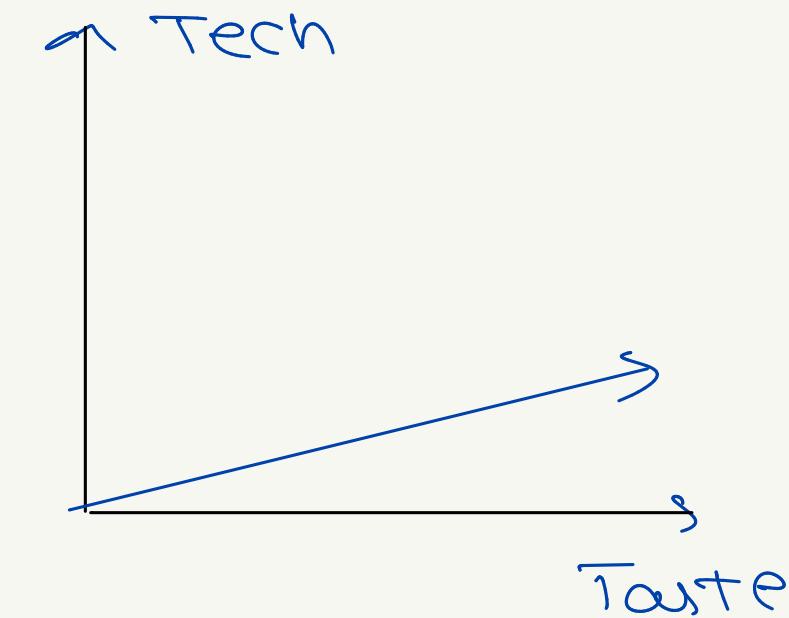
Let's assume out of 1000
900 \rightarrow Taste
100 \rightarrow Tech
overall vec $\rightarrow [0.9 \ 0.3]$

overall word embedding depends
on data set.

\rightarrow word embedding starts from book

bante hain phir static ho

goate hain for every further we.



Example

Apple launched a new phone while I was eating.

Convert this into Hindi.

$\rightarrow [0.9 \ 0.3]$ It will take a good
while it should be used as phone

So, instead of static embedding, it
should be \rightarrow contextual embedding.

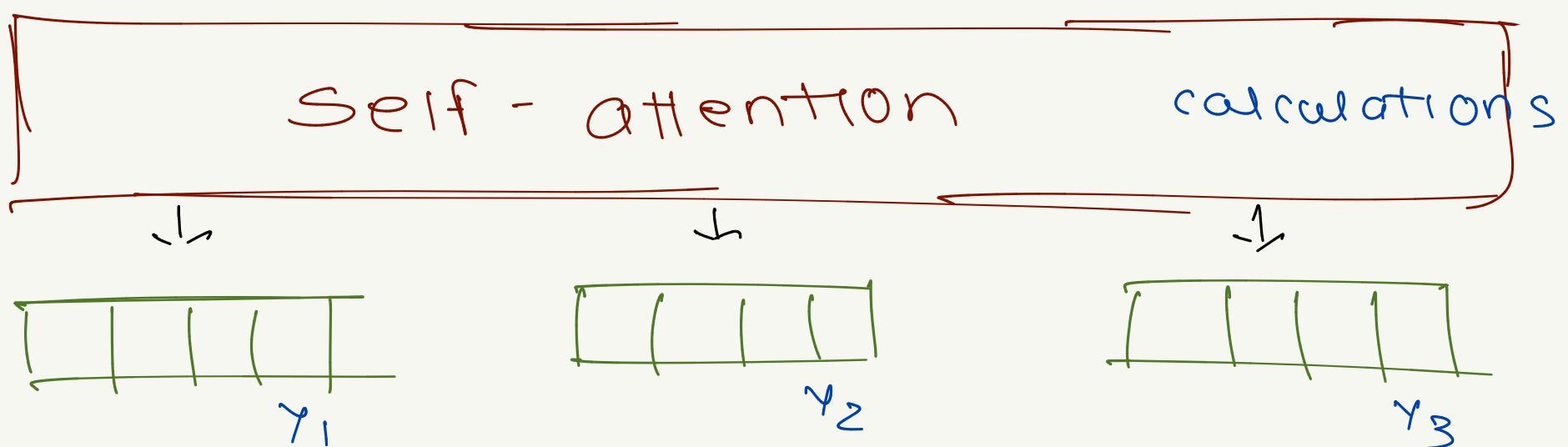
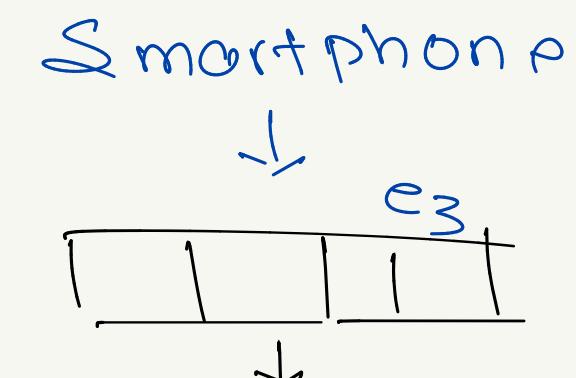
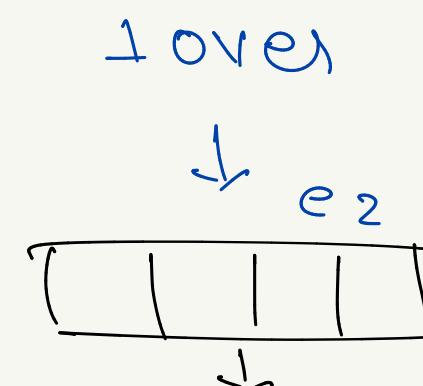
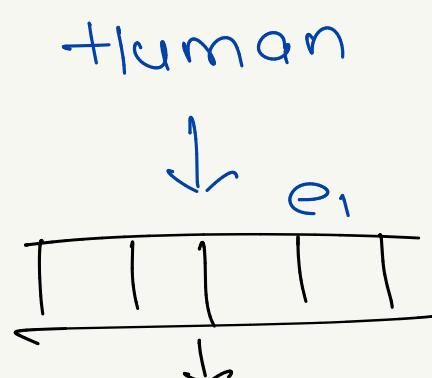
This problem is solved by Self-attention

It creates new embeddings based on context.
does not use static embeddings.

\rightarrow This is used in transformer.

How do self-attention works -

Example



money bank grows

bank \rightarrow 0.3 money

+ 0.7 bank

+ 0.1 grows

River bank flows

bank \rightarrow 0.5 river

+ 0.4 bank

+ 0.1 flows

money \rightarrow 0.7 money

+ 0.2 bank

+ 0.1 grows

river \rightarrow 0.8 river

+ 0.15 bank

+ 0.005 flows

grows \rightarrow 0.1 money

+ 0.2 bank

+ 0.7 grows

flows \rightarrow 0.4 river

+ 0.01 bank

+ 0.59 flows

we need to change money, bank, grow into embeddings too.

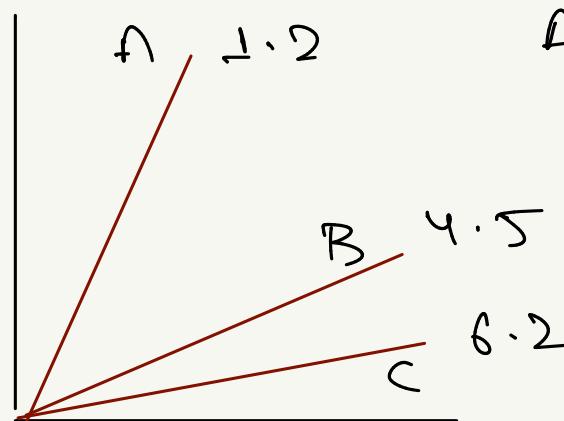
new embedding = 0.7(emoney) + 0.2(ebank) + 0.1(egrows)

of money (nD) (nD) (nD)

iska matlab ye hai ki new embedding of money
 0.7 time (money ke poorane embedding)
 + 0.2 time (bank " " "
 + 0.1 time (grow " " "
 se bna hai. (Representing similarities)

- Dot products shows similarity

for example



$$A \cdot C = 6 + 4 = 10$$

$$\begin{aligned} B \cdot C &= 24 + 10 \\ &= 34 \end{aligned}$$

As B & C are

more similar to

each other than A & C.

So the no.s in
 front of embeddings
 are actually

[DOT PRODUCTS]

So, we can write this eqn →

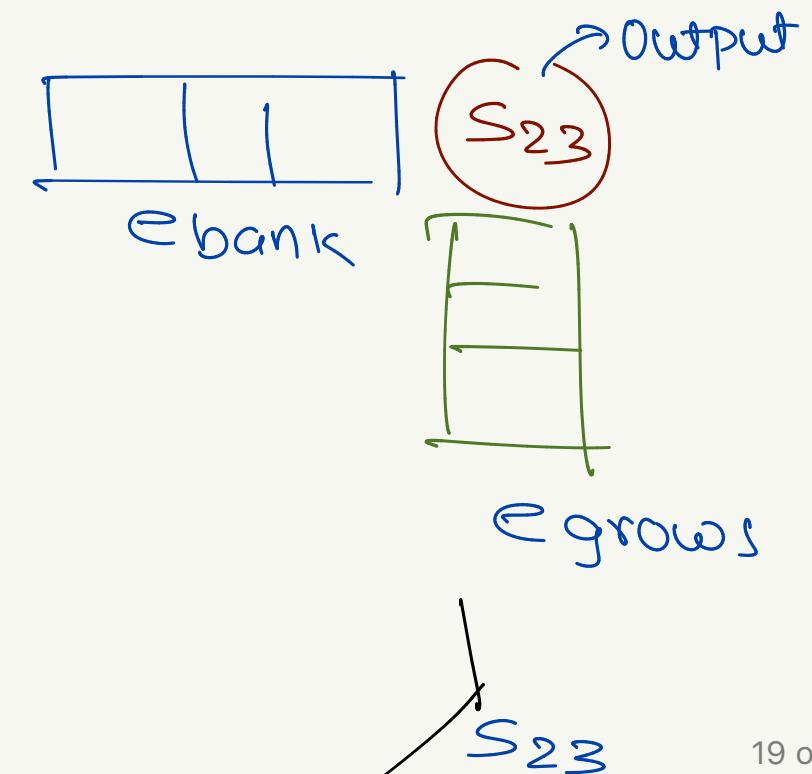
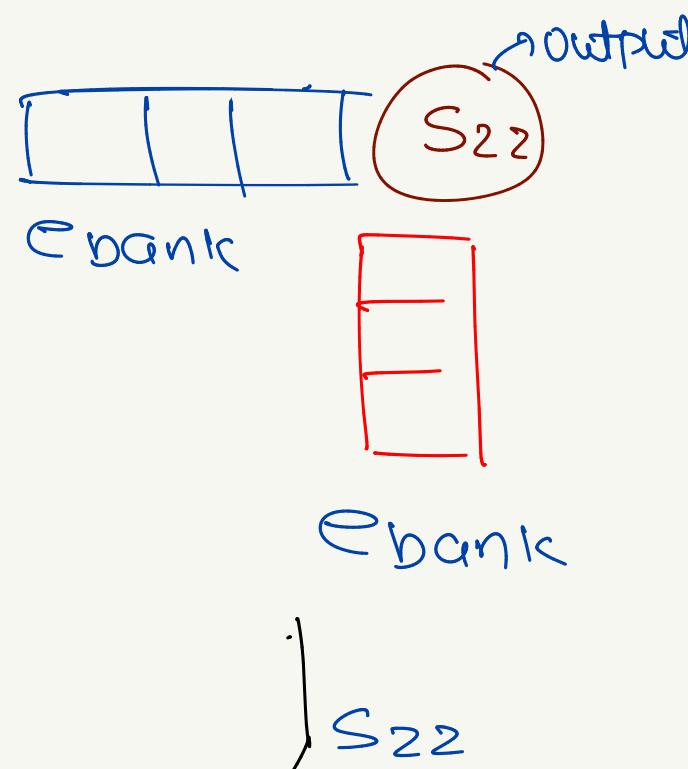
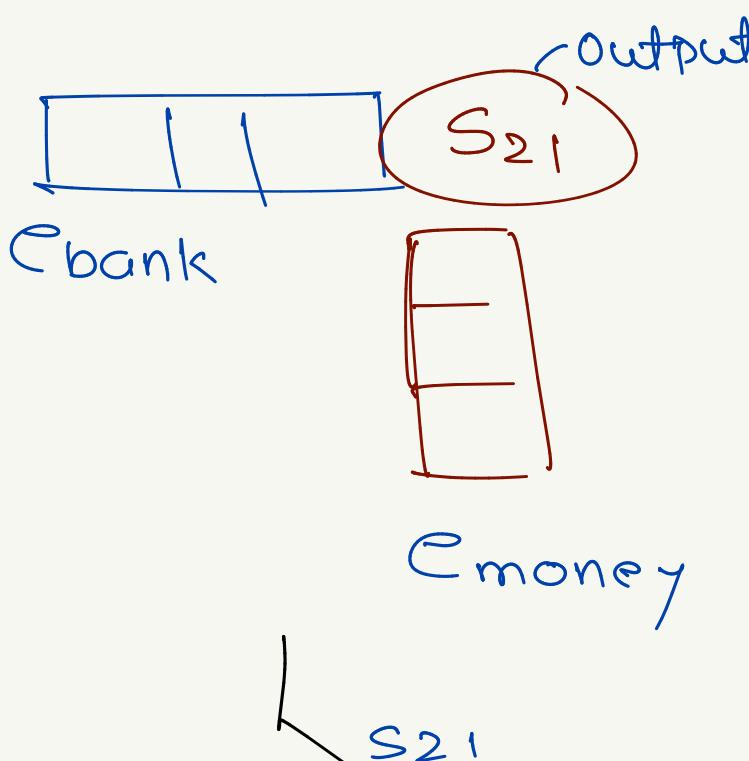
$$\begin{aligned}\underline{\text{bank}} &\rightarrow 0.3 \text{ money} \\ &+ 0.7 \text{ bank} \\ &+ 0.1 \text{ grows}\end{aligned}$$

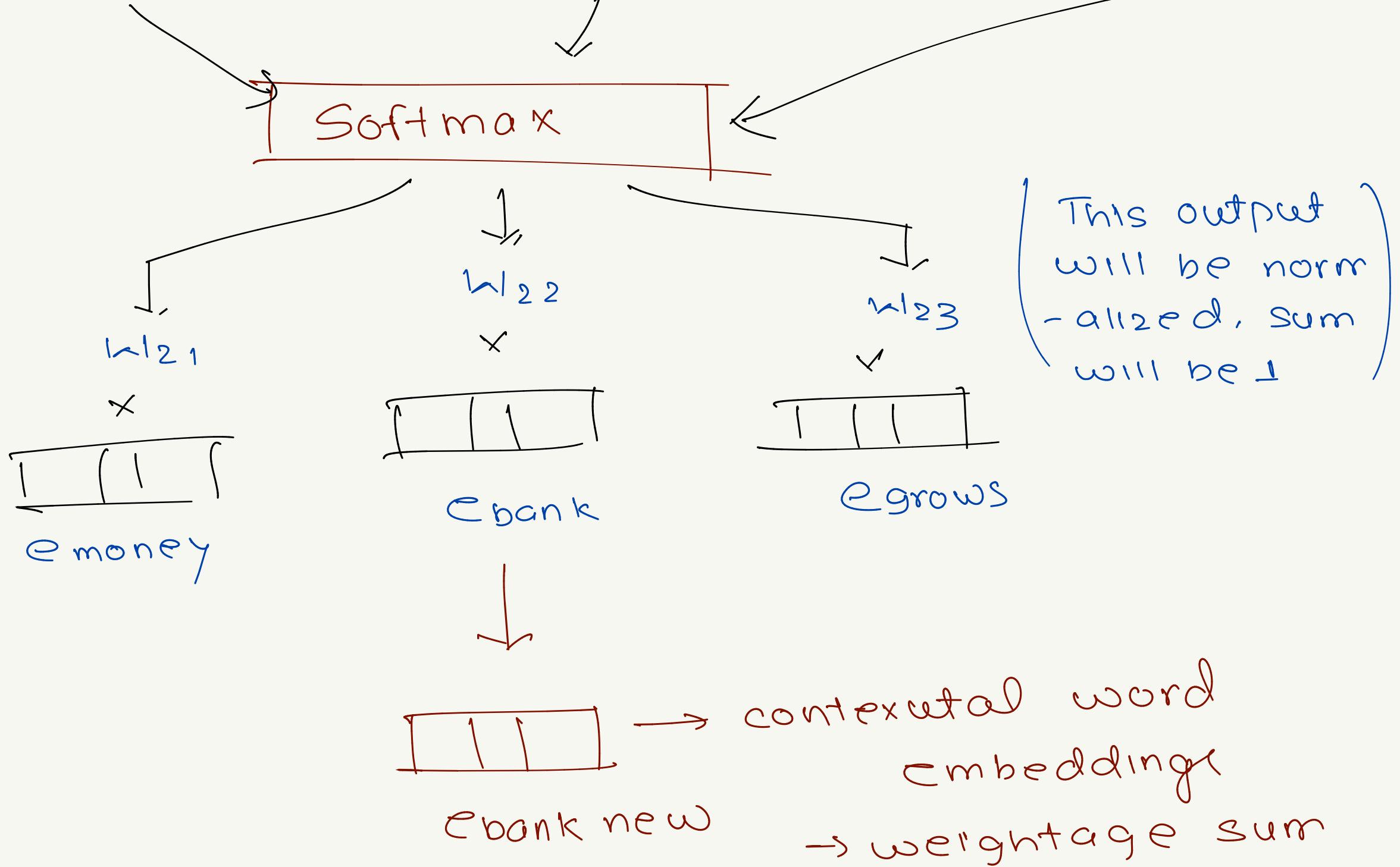
into → New embedding of $= T e_{\text{bank}} \circ e_{\text{money}}] e_{\text{money}}$

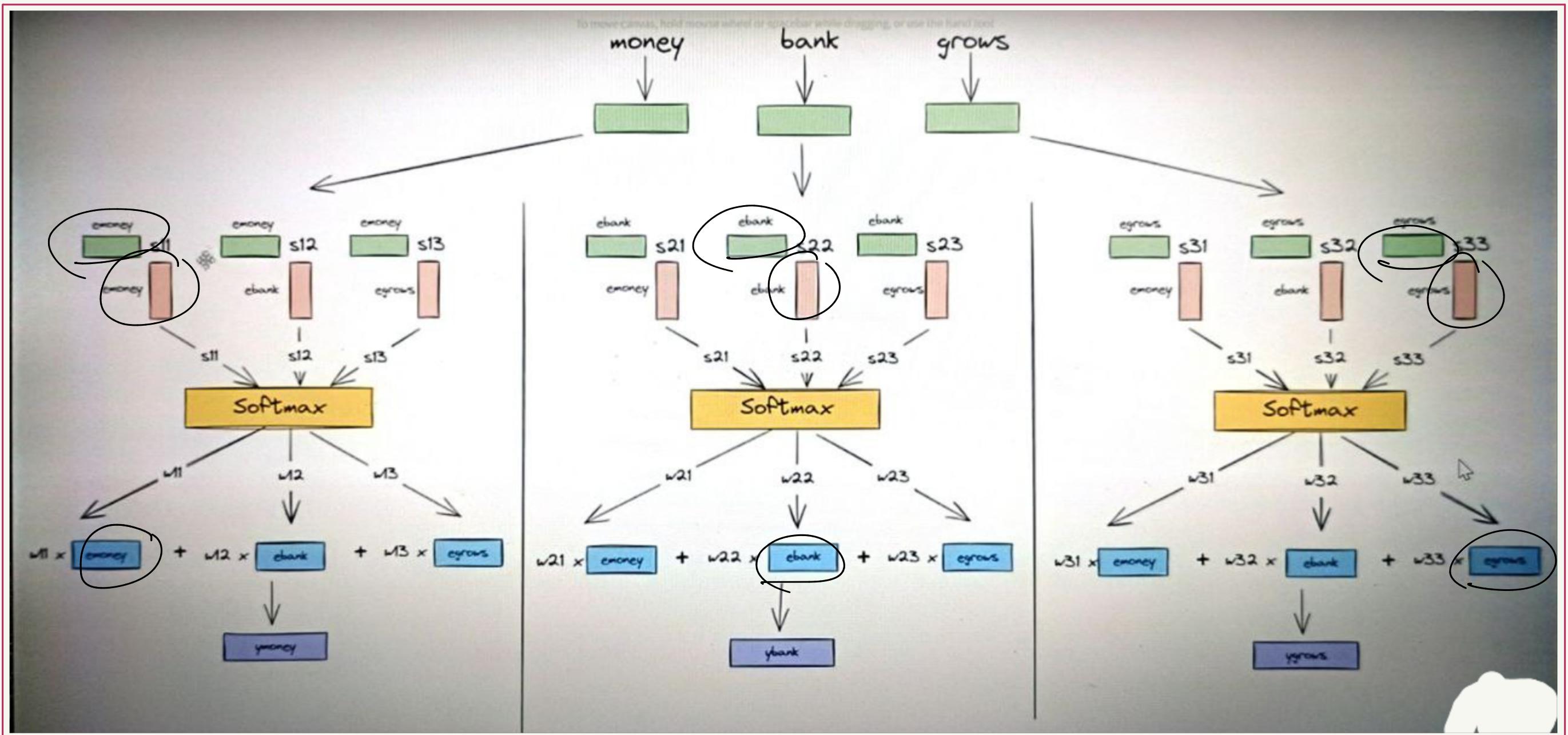
$$\begin{aligned}&\text{bank} \\ &+ [e_{\text{bank}} \circ e_{\text{bank}}] e_{\text{bank}} \\ &+ [e_{\text{bank}} \circ e_{\text{grows}}] e_{\text{grows}}\end{aligned}$$

visual of dot product →

Dot product







- This can be done parallelly.
(even there are 1000 words)
- we can generate contextual vector.

Disadvantage → loses sequential order.

→ There are no learning parameters.

means its not task specific, it just develops general contextual embedding.

Example

Piece of cake → but we want it as
केक और गोसी ↓
अदौत आसान (in data)
(generalized)
→ failed

Break a leg → रुक्खोंभिना
टिक टिक घेर ↓

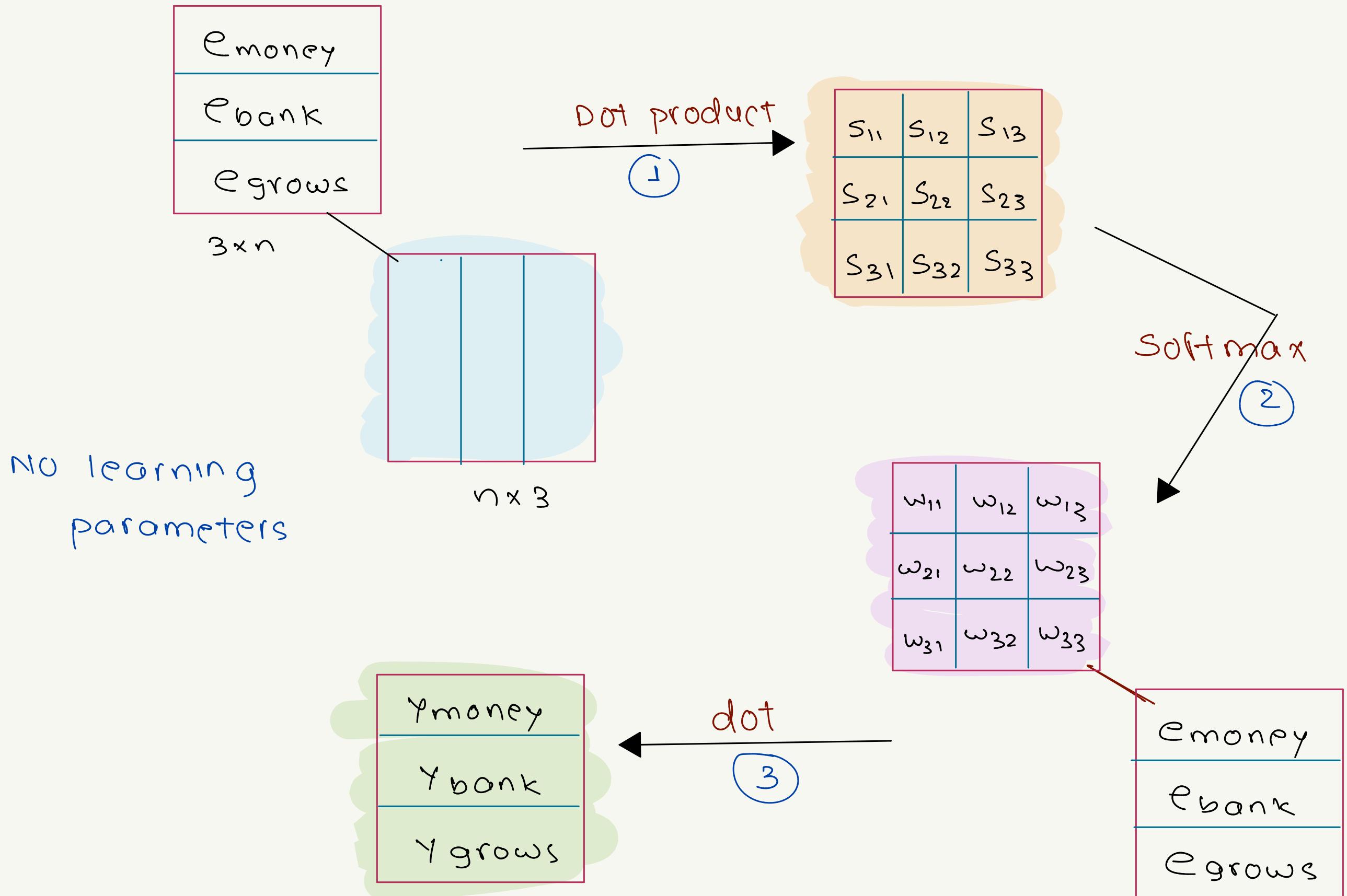
Overall → only contextual embedding is not perfect
we need task specific contextual embedding

General contextual
embedding

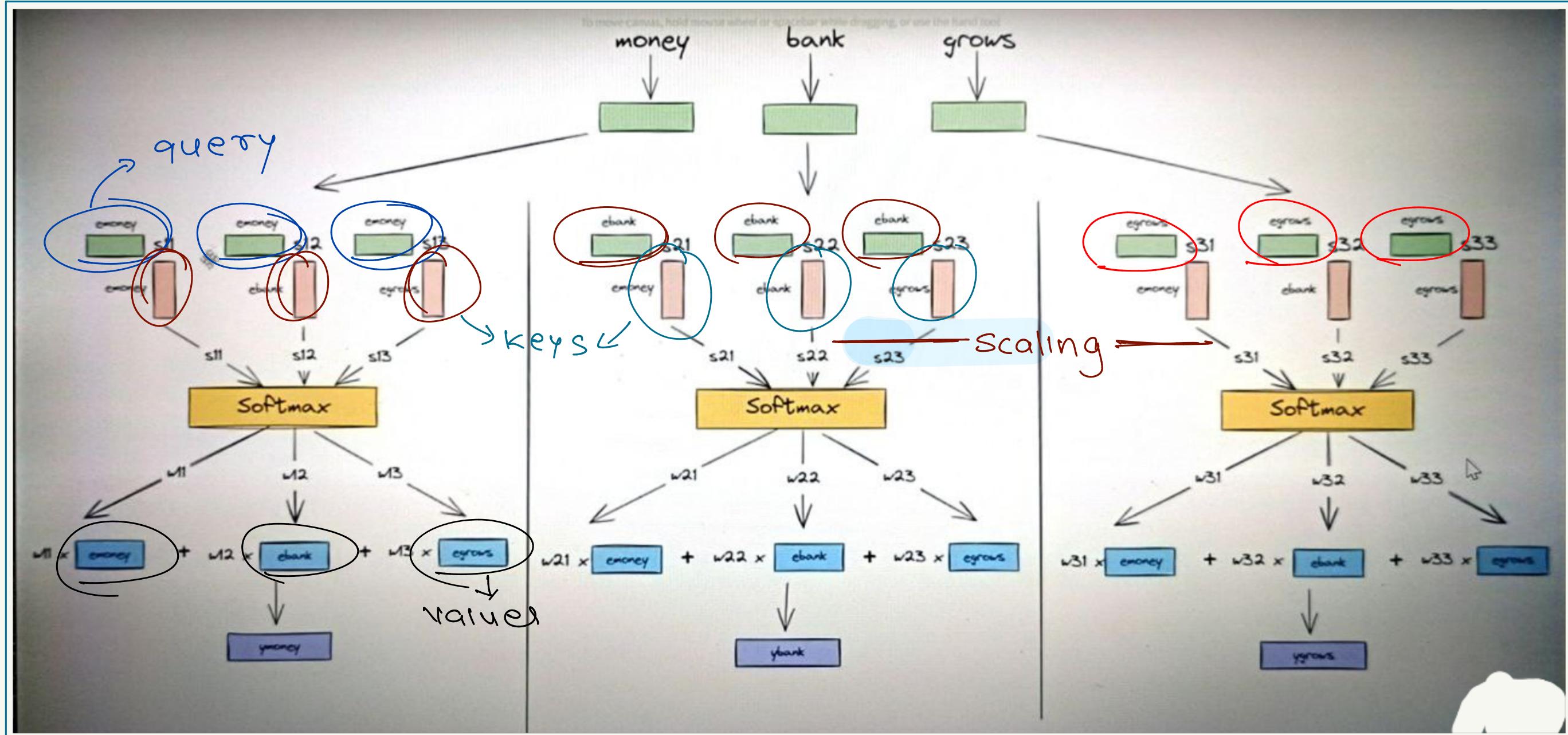
Add weights
& Bias

Task specific
contextual
embedding.

parallel embedding -



we can add weights in step① & ③ only.



Let's understand with the help of dictionary-

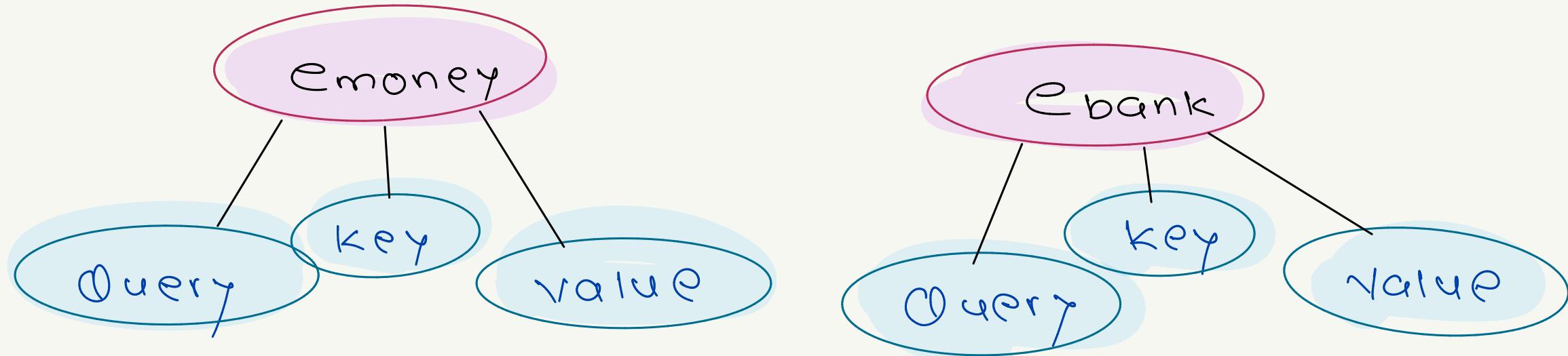
$$d = \{ \underset{\text{query}}{a} : 2, \underset{\text{key}}{b} : 3, c : 4 \}$$

↓ ↓
value

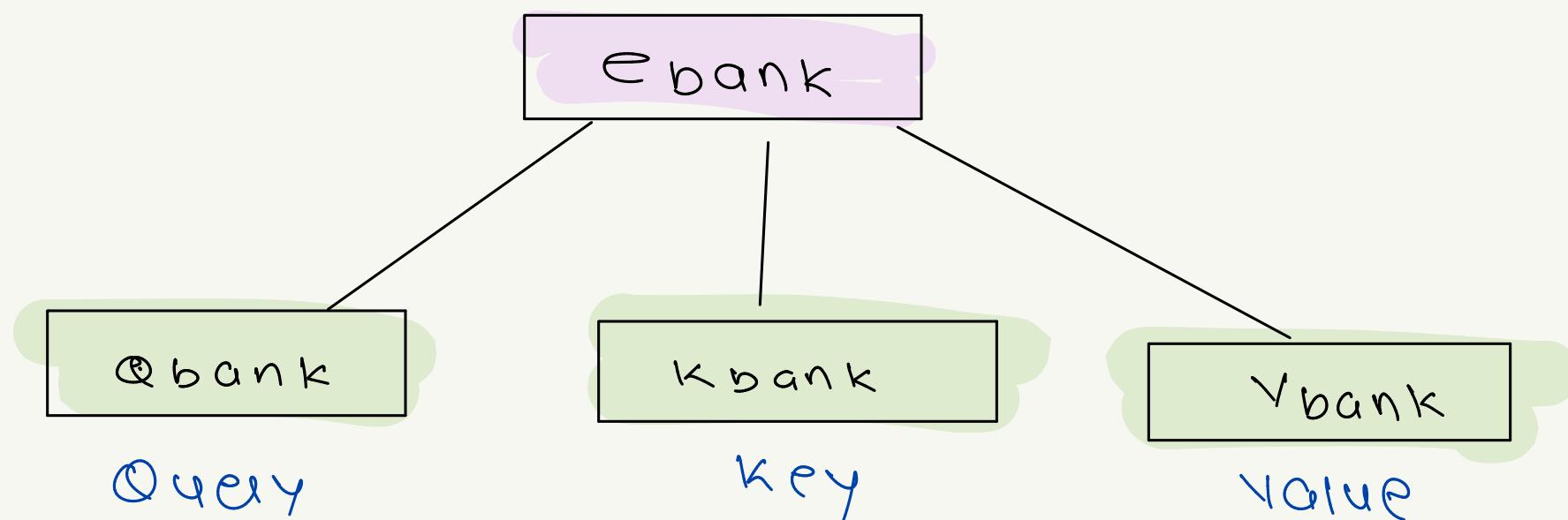
query key

If there is key named a, it will return value 2

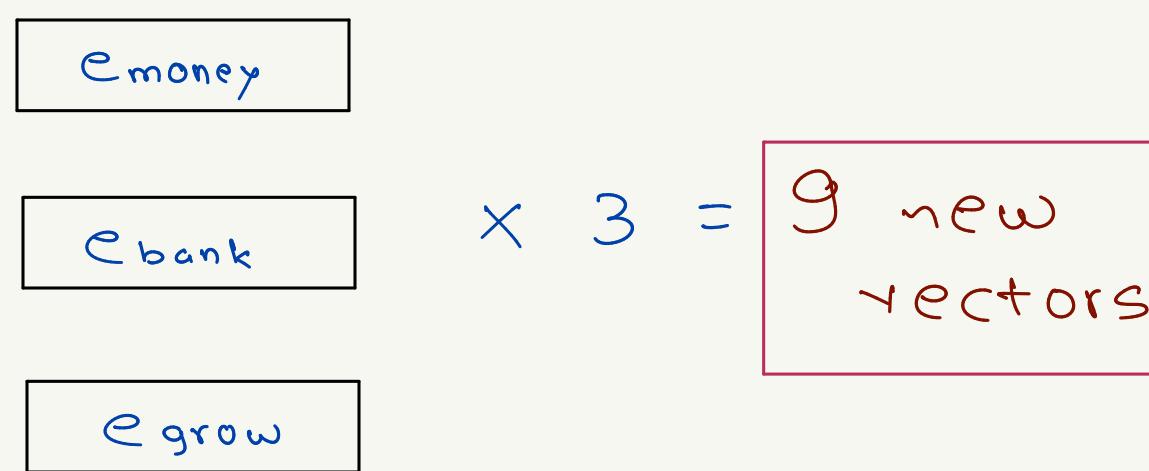
- Each of our word embedding is used in 3 ways



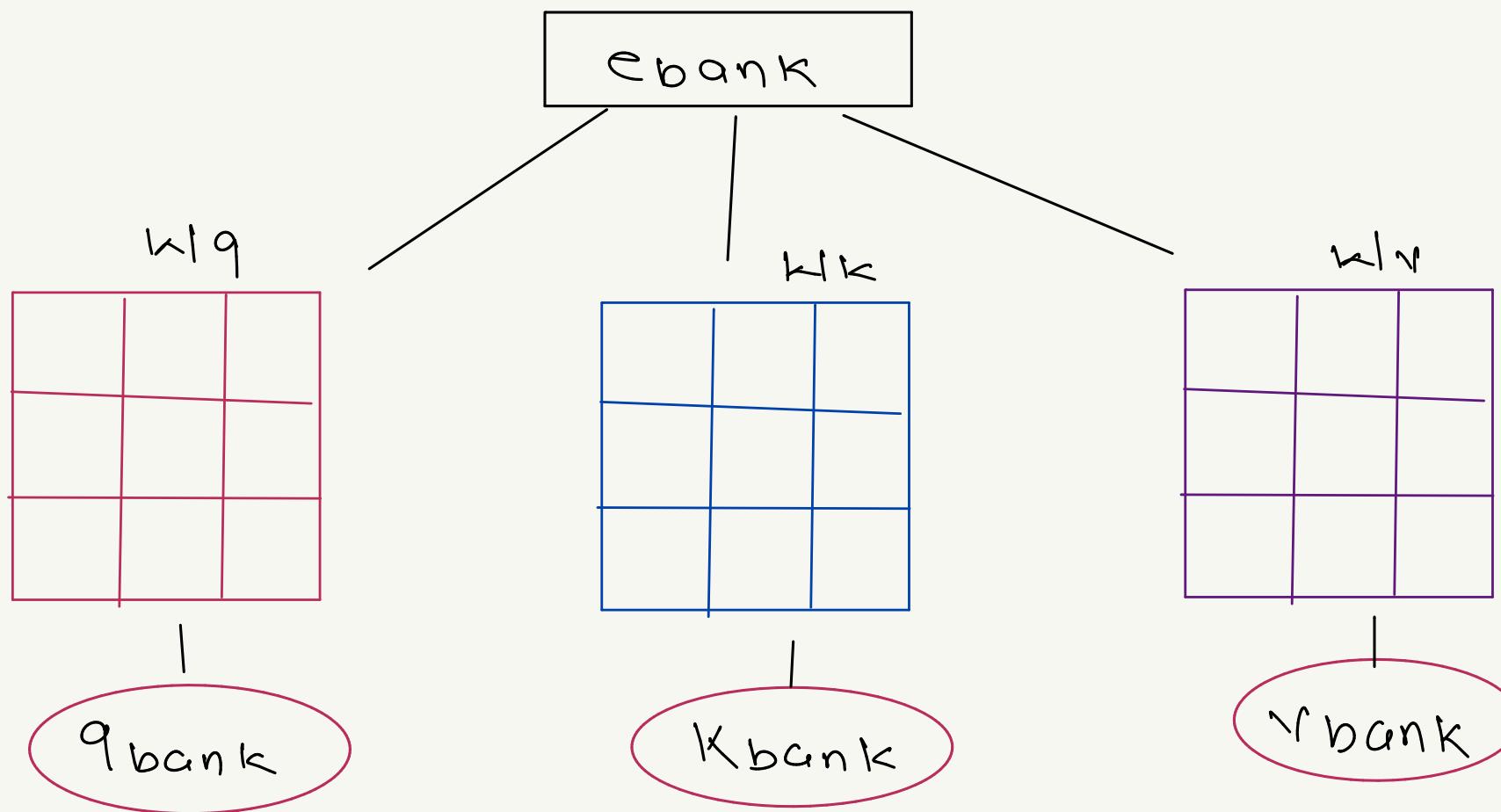
- So, we need to specify thing, we can use just word embedding as query, key & value.
- This Query, key and value vector will be formed with the help of data & we.



So, we have 3 words



By linear transformation
→ By multiplying original vector with diff matrix to get 3 new vectors.

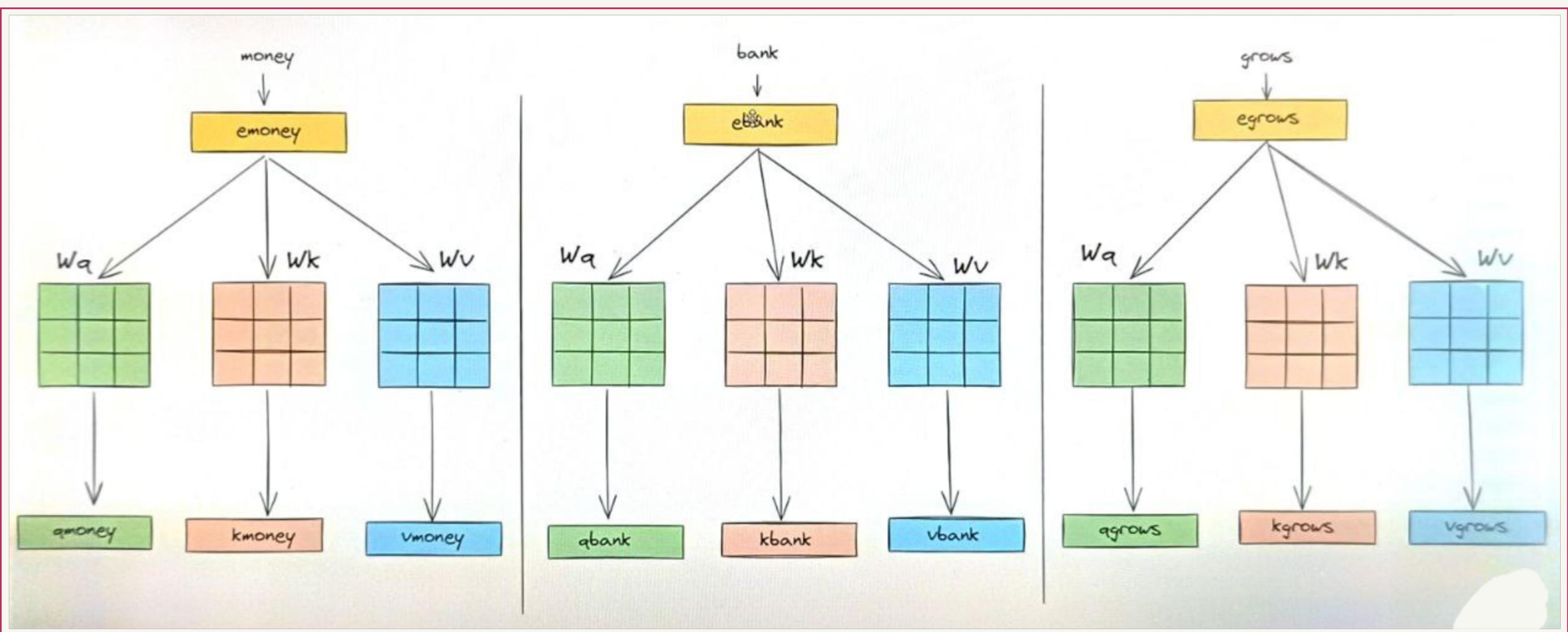


Ye matrix baar baar training
karke shi value gain kar lenge,

Start with random weights

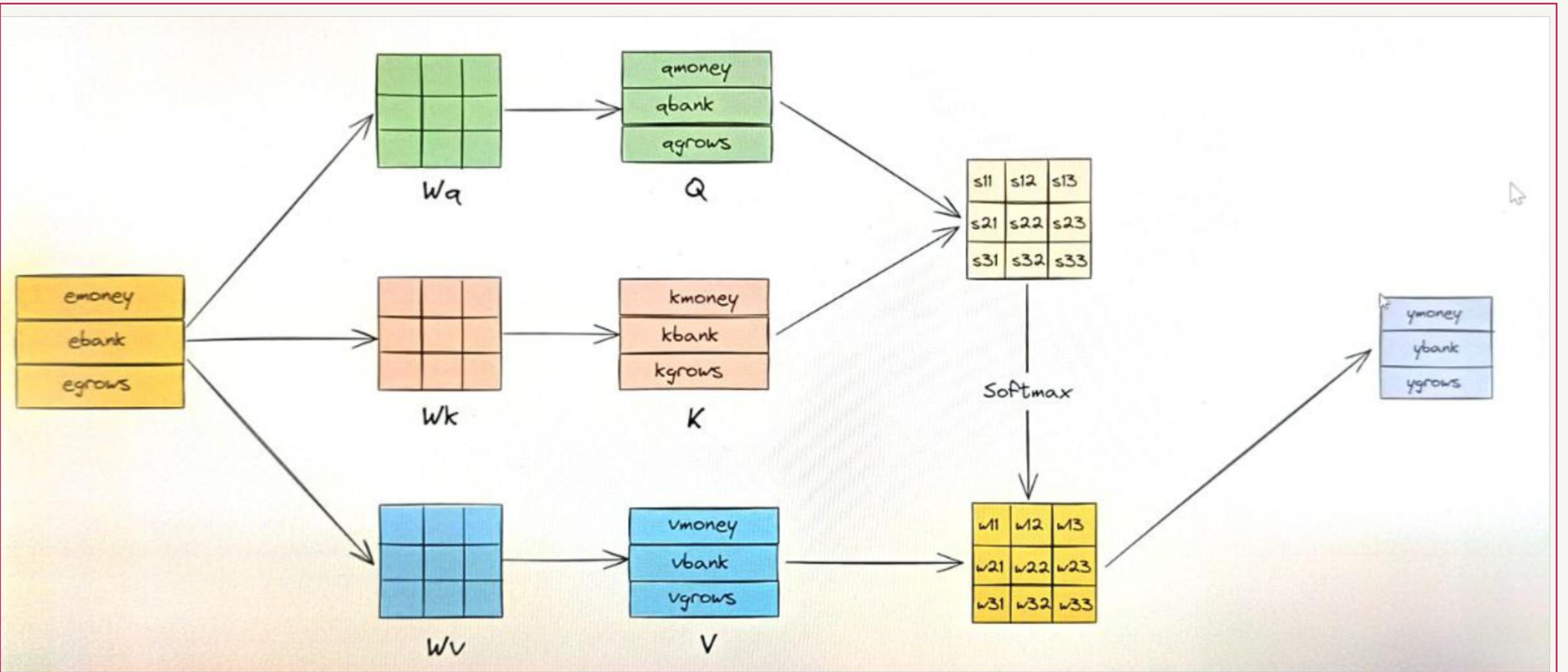
matrix will be decided by training process & data.

↓
we will



$$w_{q1} = w_{q2} = w_{q3} \quad \& \quad w_{k1} = w_{k2} = w_{k3} \quad \& \quad w_{v1} = w_{v2} = w_{v3}$$

& we finally got task contextual embedding



we can perform all the action parallelly with the help of powerful GPUs.

$$\text{Attention} = \text{Softmax}(Q \cdot K^T) \times V$$



$$\text{Softmax} \frac{(Q \cdot K^T) V}{\sqrt{dk}}$$

Scaled dot product
Attention

$d_k \rightarrow$ dimensionality of k vector.

to reduce unstable
gradiing.

\rightarrow dimension of key vector.

So, softmax lagore se pehle hum square value K^0
 $\sqrt{d_k}$ se divide karna padega.

why? \rightarrow Dot-product ka nature

$\underbrace{Q \cdot K^T}_{\text{Dot product of } Q \text{ matrix}}$

Actually we
are doing dot products
of multiple vectors.

when we have low dimensional dot product

\hookrightarrow low variance

when high dimensional dot product

\hookrightarrow high variance

problem

Softmax → converts given data into probability & normalize it.

→ depends on $e^n \rightarrow$ so bde no ko very high probability dega
aur chhote wale ko bht low.

so, overall Scaling se chhote aur bde no dono proper training hogi, kyoki scaling se variance kam noga.

Example

$$1) ([10, 20, 30, 40, 50, 60, 70])$$

$$\text{var} = 400$$

$$2) \left(\frac{10}{10}, \frac{20}{10}, \frac{30}{10}, \frac{40}{10}, \frac{50}{10}, \frac{60}{10}, \frac{70}{10} \right)$$

$$\rightarrow ([1, 2, 3, 4, 5, 6, 7])$$

$$\text{variance} \rightarrow 4.0$$

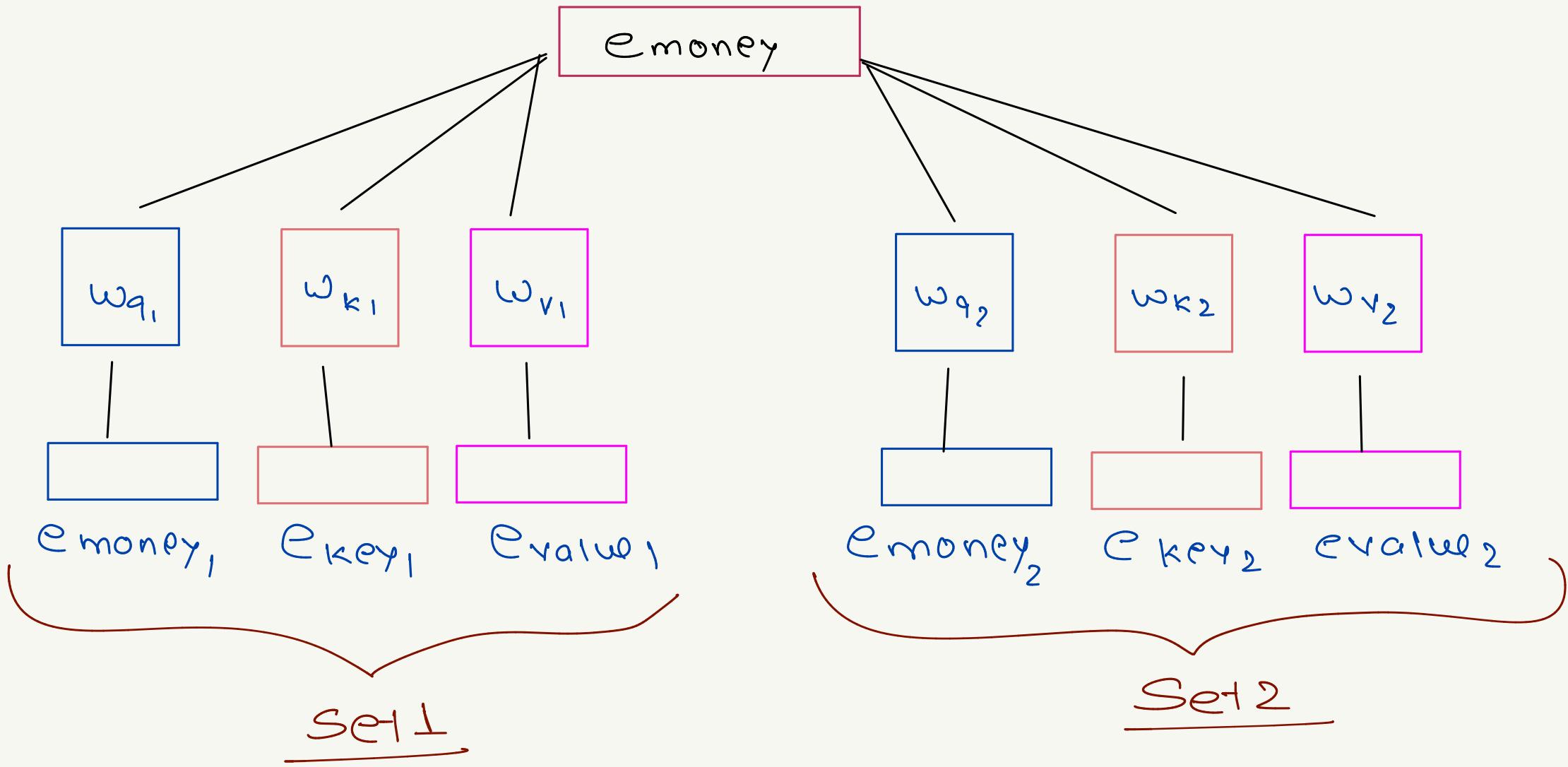
Spread
bw
no. in
dataset

Multi-head Attention

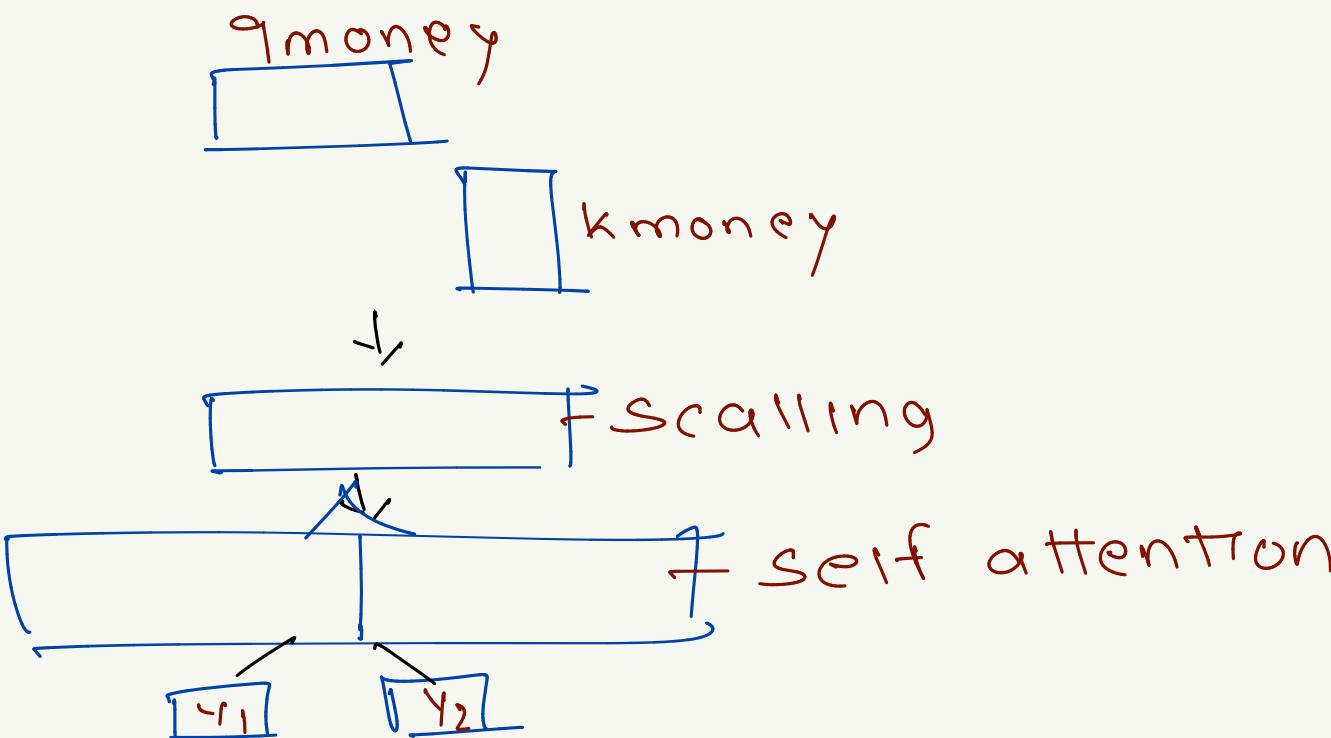
Example

- The man saw the astronomer with a telescope.
- have two meaning, but self-attention can only generate one meaning
 - Self attention can only consider single perception.
 - ⇒ Multihead attention solved this problem.
 - ↳ ek se jyada self attention mil hoga.

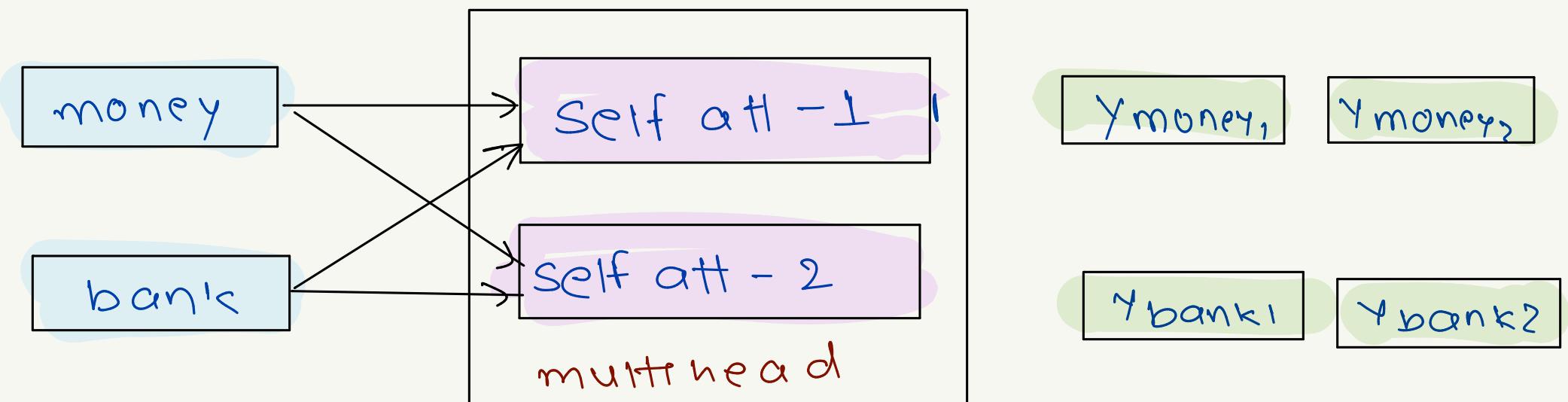
Example



- Aise hi bank & grows ke liye bhi 2-2 set aayenge



2 contextual text
vector aayega,



Positional Encoding

- we can generate contextual embeddings parallelly by self-attention.
- order of words lost ↓
for self attention
 - Harshit killed lion.
 - 2) Lion killed Harshit.

Both are same for self-attention.

Sol ① If we add one more dimension that will carry positional no. of word.

Prblm → If we want to process a book then last no will be a very big no. 10000 or more

1 Neural network hates large numbers
it prefers no bw $[-1 \text{ to } 1]$

↳ no. 1, 2, 3... are discrete nos while
neural network prefers smooth transitions
continuous.

we also have sol. for this.

we can divide all the positional no. with the
last no, so the no will be within -1 to 1.

Pblm

we can't normalize bcz neural network will
conclude → whether the second position is 2 or
0.2 or 0.00002. Position no. will not be consistent.

Pblm 3 → we can't capture relative positioning

Ex

horsht, is a good boy,
1 2 3 4 5 → absolute
positioning
Relative positioning) distance bw them.

Overall

- ① Unbounded
- ② Discrete
- ③ No relative position

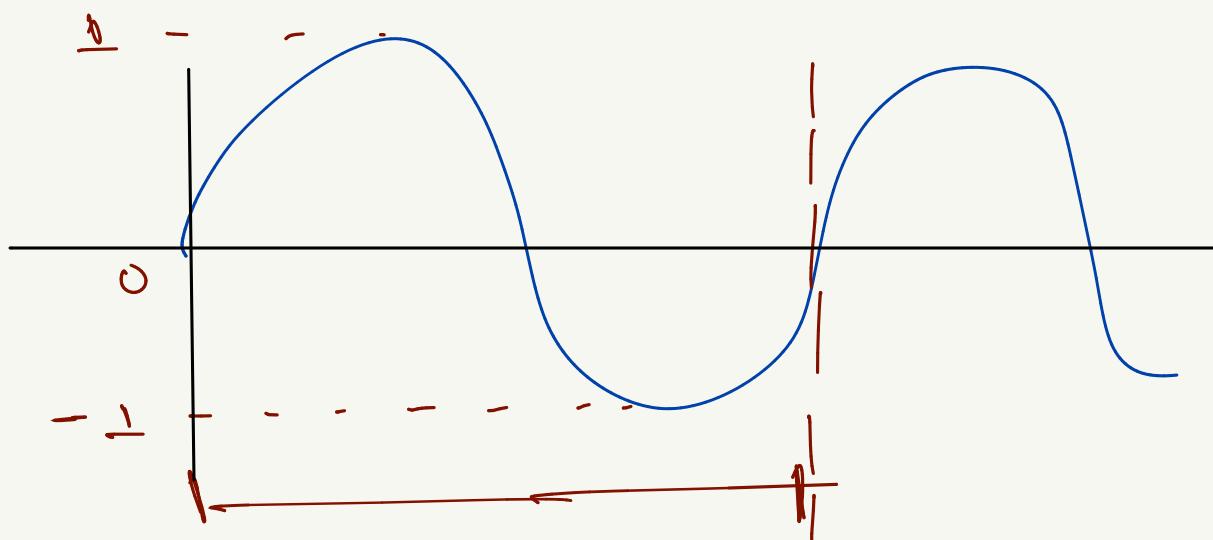
Bounded

continuous

periodic fun

All these criteria
is fulfilled by trigonometric fun.

- Sin fun
- continuous
- periodic
- bounded



Example

- 1 + harshit
- 2 killed
- 3 a
- 4 lion



Sci
f
Atte
ntr
m

$$y = \sin(\text{position})$$

Summary:

- $\sin(1) = 0.841$
- $\sin(2) = 0.909$
- $\sin(3) = 0.141$
- $\sin(4) = -0.757$

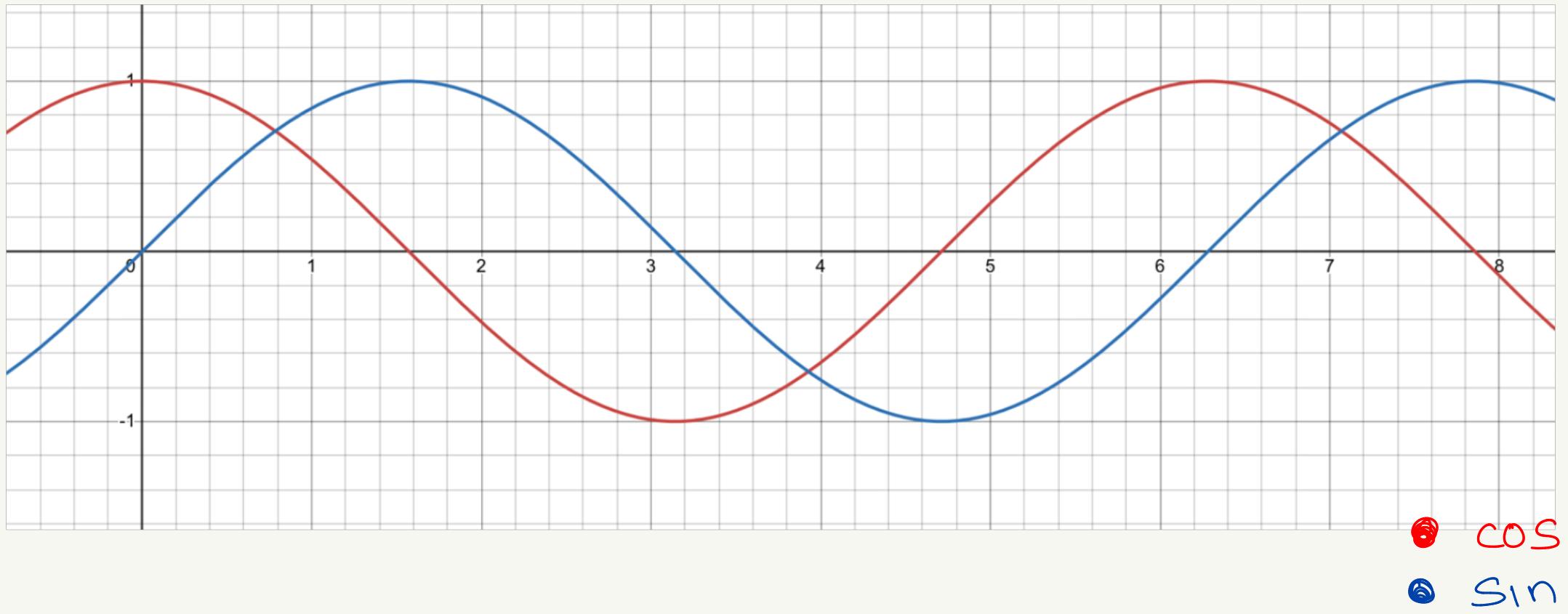
1	→ tarsnit	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>0.841</td></tr></table>	1	1	1	0.841
1	1	1	0.841			
2	killed	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>0.90</td></tr></table>	1	1	1	0.90
1	1	1	0.90			
3	a	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>0.14</td></tr></table>	1	1	1	0.14
1	1	1	0.14			
4	Lion	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>-0.75</td></tr></table>	1	1	1	-0.75
1	1	1	-0.75			

positional
encoding

Pblm

2 word ke positional encoding
same nhi ho sakta , par sin curve is periodical
so ye repeat karega & 2 words ke same encoding
dega.

Sol hum ek jn we korne ke badle 2 trigono
metric jn we kar sakte hai both sine & cos.
1 position encoding → vector



So, now we have 2 values for each no.

Sine Values

- $\sin(1) = 0.841$

- $\sin(2) = 0.909$

- $\sin(3) = 0.141$

- $\sin(4) = 0.656$

Cosine Values

- $\cos(1) = 0.5403$

- $\cos(2) = -0.4161$

- $\cos(3) = -0.9899$

- $\cos(4) = -0.6536$

Harshit

10.8	0.5
------	-----

Killed

	0.9 -0.4
--	------------

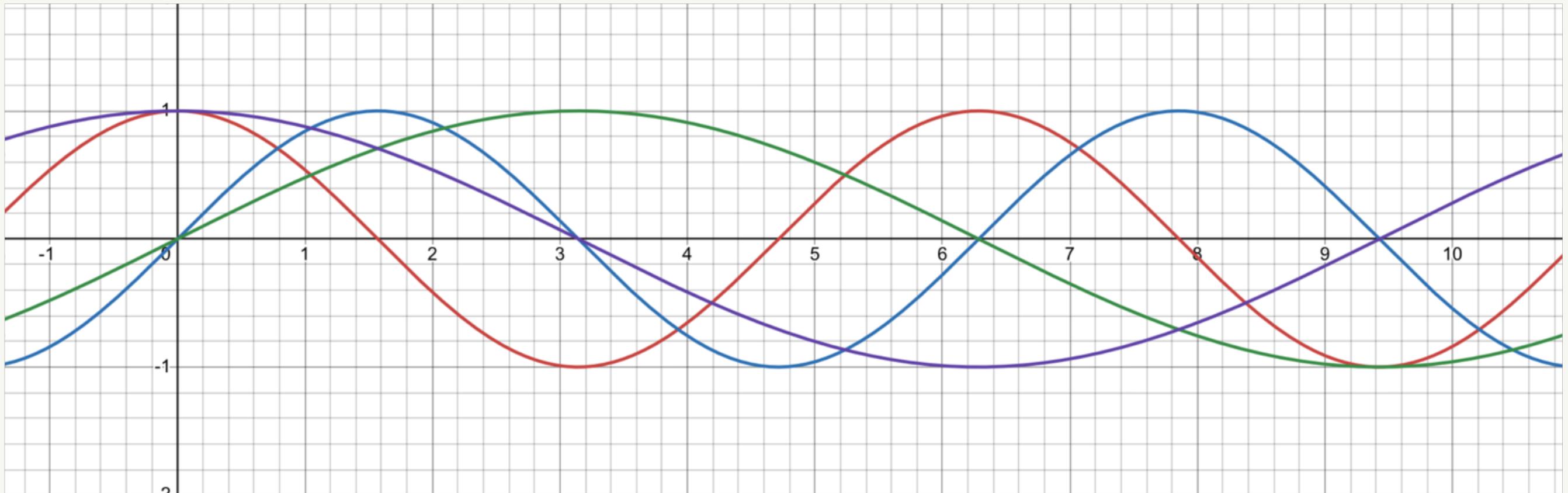
a

	0.1 0.9
--	-----------

10n

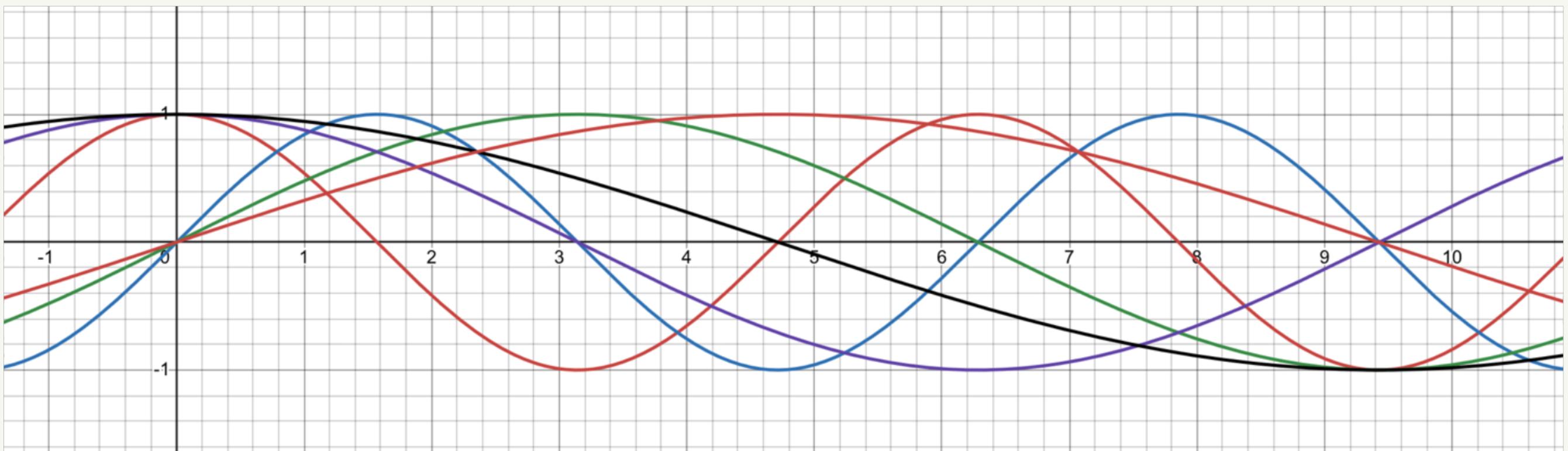
	0.6 0.6
--	-----------

There are still possibilities of matching, so we can add more functions $\cos(\text{position}/2)$ & $\sin(\text{position}/2)$

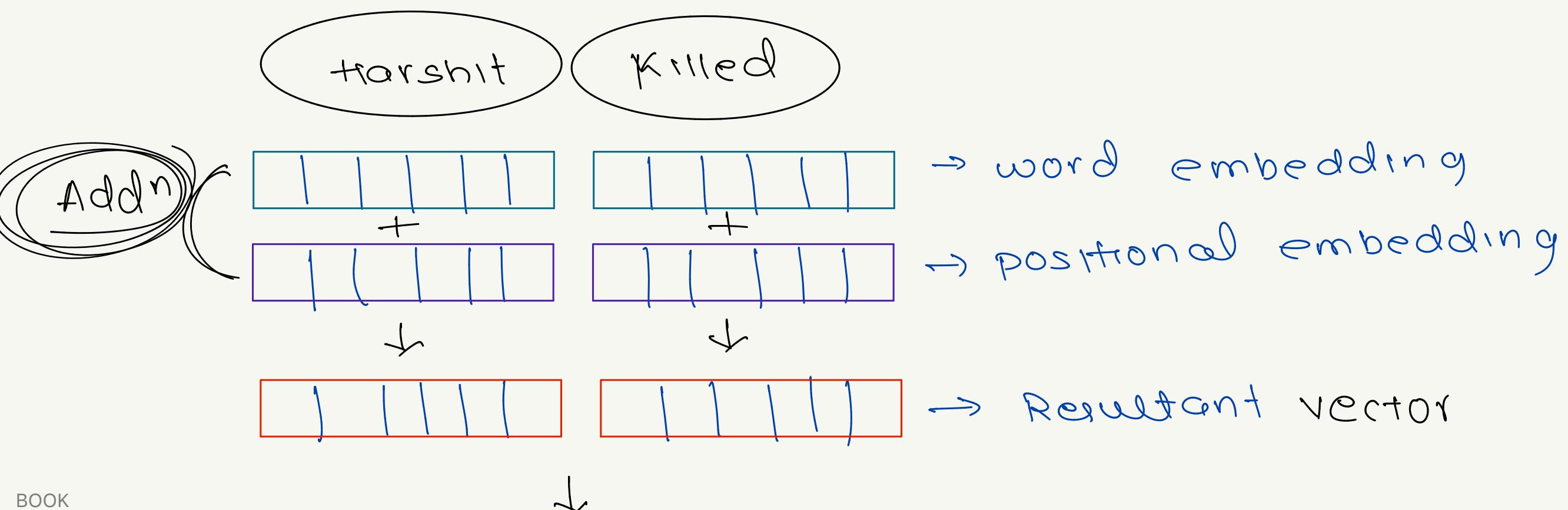


Now probability of same positional encoding has decreased more.

So, even now we can add more functions $\cos(\pi/4)$ & $\sin(\pi/4)$, $\sin(\pi/3)$ $\cos(\pi/3)$ etc



Har word ke positional embedding vector ka dimension exactly same hogा as dimension of embedding vector.



[Self Attention]

⇒ what's inside positional embedding.

Let's consider if positional vector
is of 6 dimension

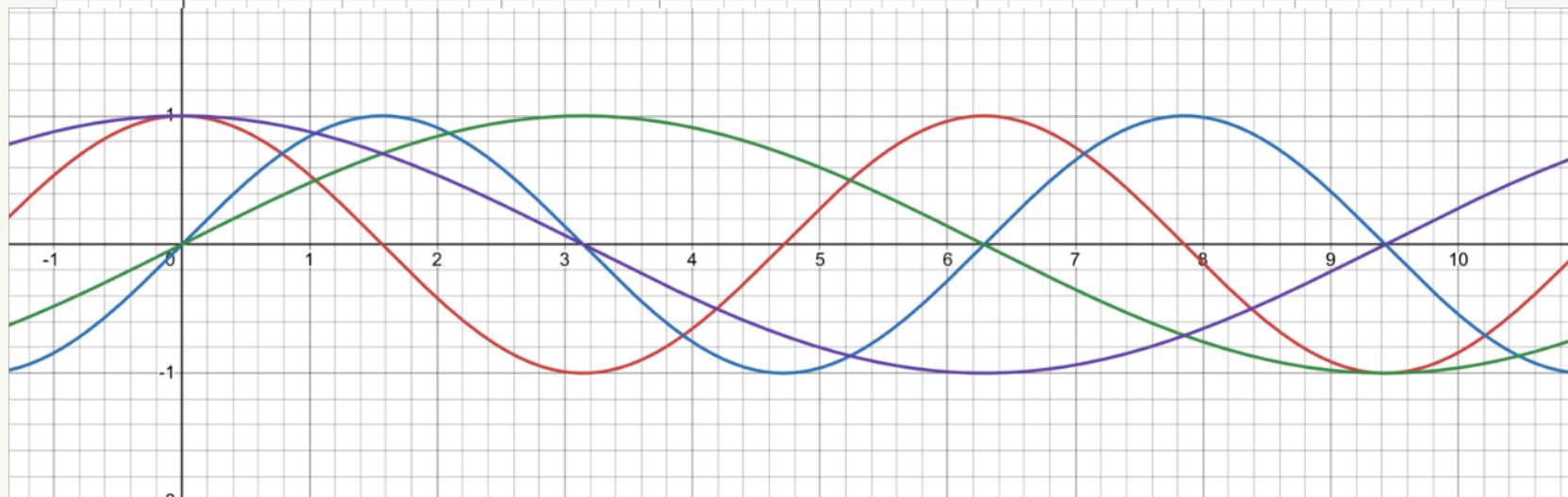
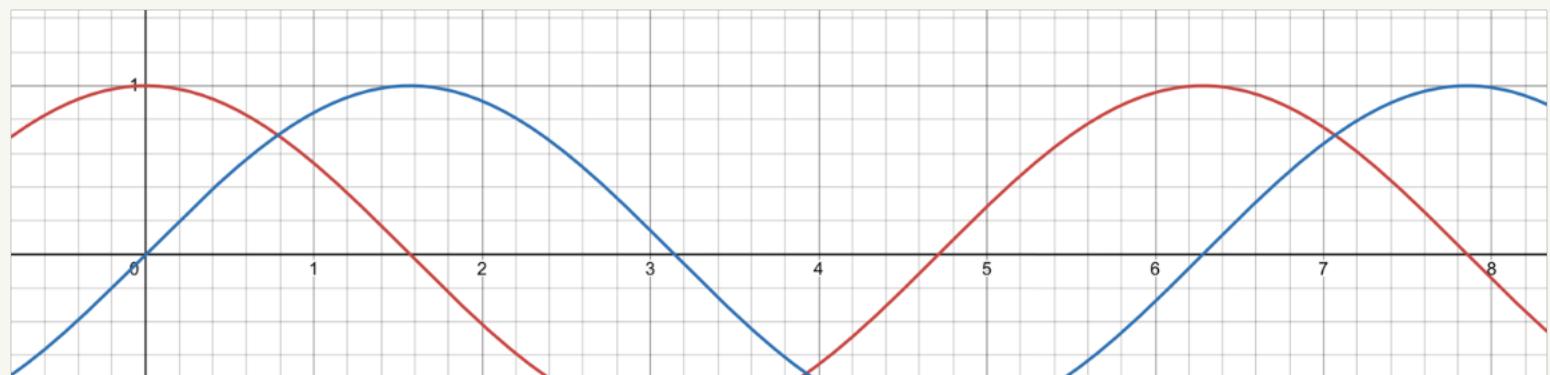
River

$\sin(1)$	$\cos(1)$	$\sin(1/2)$	$\cos(1/2)$	$\sin(1/3)$	$\cos(1/3)$
-----------	-----------	-------------	-------------	-------------	-------------

Bank

$\sin(2)$	$\cos(2)$	$\sin(2/2)$	$\cos(2/2)$	$\sin(2/3)$	$\cos(2/3)$
-----------	-----------	-------------	-------------	-------------	-------------

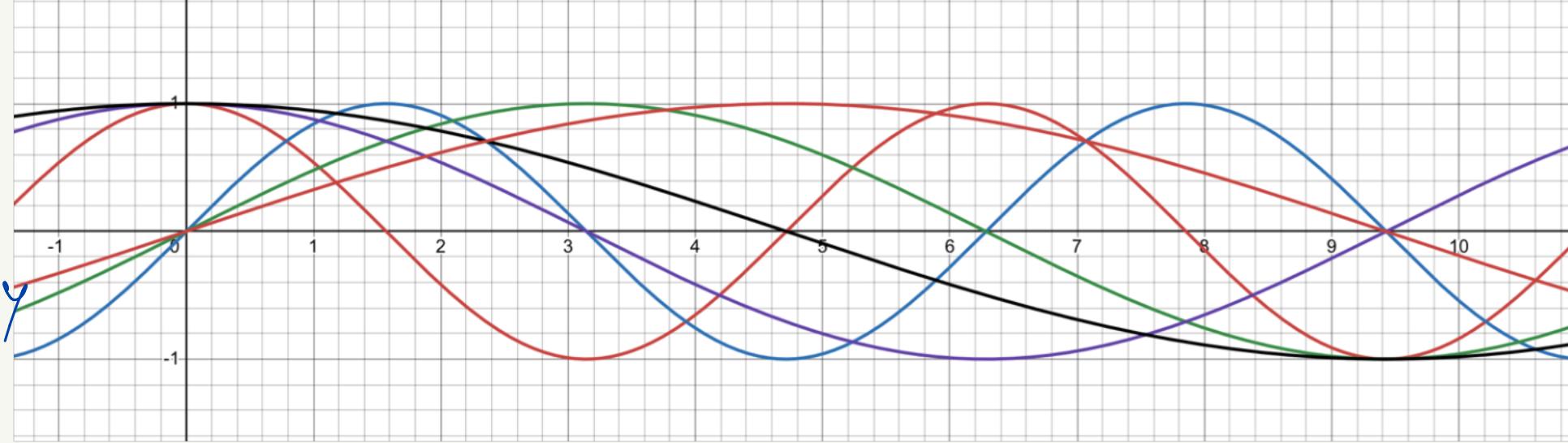
Jaise jaise aur
sin, cos curve
aayenge unki fre
quency kam hoti
jaoayegi.



Exact formula

pos → position
of word

dmodel → dimensionality
of word



$$PE(pos, 2i) = \sin(pos / 10000^{2i/dmodel})$$

$$PE(pos, 2i+1) = \cos(pos / 10000^{2i/dmodel})$$

Example

River

1	2	3	4	5	6
---	---	---	---	---	---

Bank

1	2	3	4	5	6
---	---	---	---	---	---

word embedding → 6D

positional embedding → 6D

pos → 0

dmodel → 6

pos → 1

dmodel → 6

$i \rightarrow 0 \text{ to } \frac{d_{\text{model}}}{2} - 1$

i.e. 0 - 2

$i = 0, 1, 2$

$i \rightarrow 0 \text{ to } \frac{d_{\text{model}}}{2} - 1$

i.e. 0 - 2

$i = 0, 1, 2$

For River

$$1D \rightarrow PE(0, 0) = \sin(0 / 10000^{\circ}) = 0$$

$$2D \rightarrow PE(0, 1) = \cos(0 / 10000^{\circ}) = 1$$

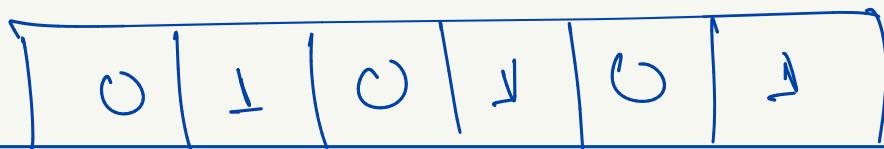
$$3D \rightarrow PE(0, 2) = \sin(0 / 10000^{\circ}/3) = 0$$

$$4D \rightarrow PE(0, 3) = \cos(0 / 10000^{\circ}/3) = 1$$

& so on

$$5D \rightarrow 0$$

$$6D \rightarrow 1$$



For Bank

using same formula -

$$1D \rightarrow 0.84$$

$$2D \rightarrow 0.54$$

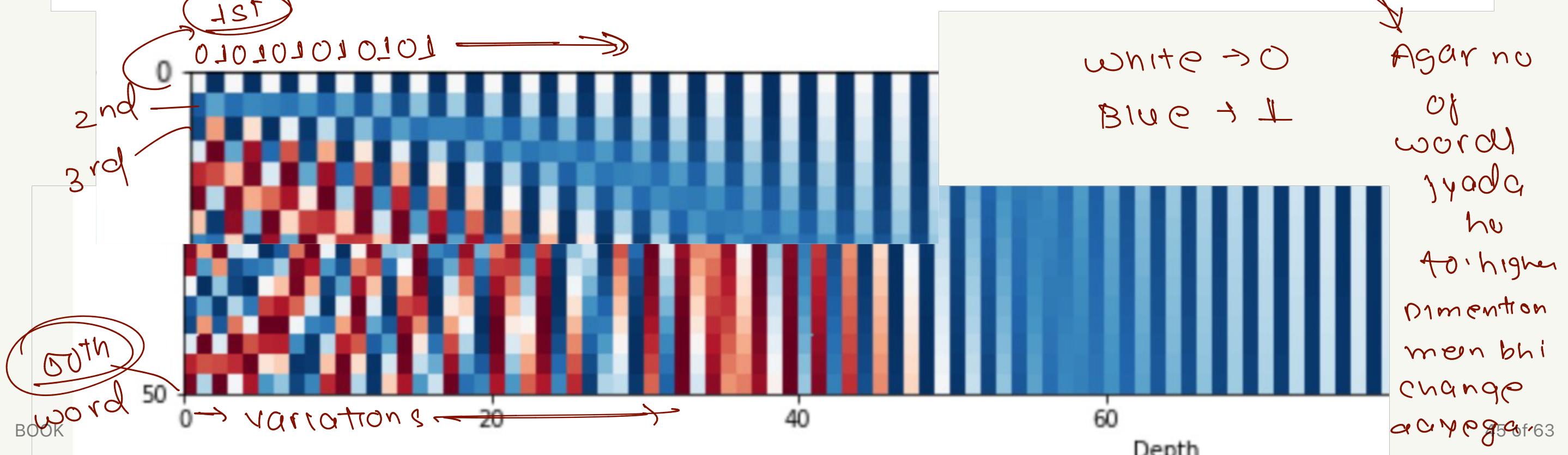
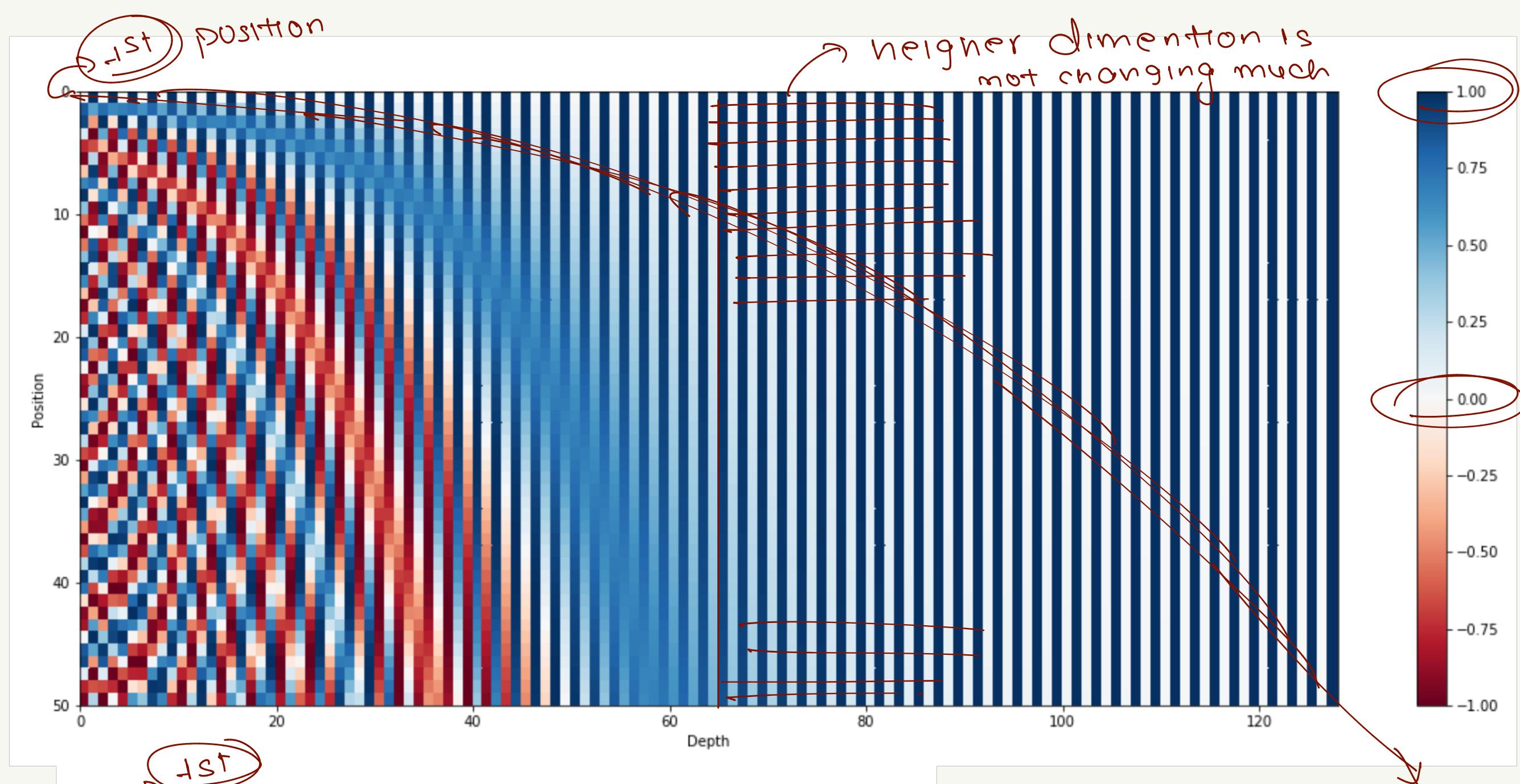
$$3D \rightarrow 0.04$$

$$4D \rightarrow 0.99$$

$$5D \rightarrow 0.00$$

$$6D \rightarrow 0.99$$

This graph shows positional encoding of 50 words each of 128 D.





ChatGPT

Masked Language Modeling (MLM) is a type of pre-training task used in training transformer-based language models, most notably BERT (Bidirectional Encoder Representations from Transformers). MLM is crucial for enabling these models to understand and generate human language in a context-aware manner. Here's a detailed look at MLM:

Key Concepts

1. Masked Tokens:

- During training, a percentage of the tokens in the input sequence are randomly replaced with a special [MASK] token.
- Typically, 15% of the tokens are masked in the training sequence.

2. Contextual Learning:

- The model is trained to predict the original tokens that have been replaced by the [MASK] tokens.
- This requires the model to use the context provided by the surrounding (unmasked) tokens to infer the masked token.