# Transformers

ANN → Tabular Data

RNN → Sequence data
(Text)

CNN → Image Data

Transformer → Sequence
to Sequence
task

Impact of Transformers

```
Revolution
in
NLP

Demoraltisiy
AI
  ↓    ↓
Bert, GPT

Multi model
capability

Accerlation
of
GenAI

Unification
of
deep learning
```

# Self Attention

## The what

NLP → Words to number [vectorization]
                    Very important

## Methods.

1. OHE

mat   cat   mat

cat   nat   nat

$[1\,0\,0]$   $[0\,1\,0]$   $[1\,0\,0]$

|     | mat | cat | nat |
| --- | --- | --- | --- |
| mat | 1 | 0 | 0 |
| cat | 0 | 1 | 0 |
| nat | 0 | 0 | 1 |

2) Bag of words (BOW)

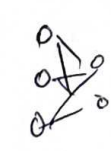|       | mat | rat | cat |
| ----- | --- | --- | --- |
| $S_1$ | 2 | 0 | 1 | → Cat 1 time in first sentence

↳ mat 2 times in first sentence

| $S_2$ | 0 | 2 | 1 |

↳ 2nd sentence

3) Word embedding → Semantic Meaning

Training Data send (Very Large) → ↘ → $\begin{bmatrix} \circ \\ \circ \\ \circ \\ \circ \end{bmatrix}$    Each word converted into n-dim Vector

↓
NN

↳ n-dim Vector

let say, we have 5-dim vector

king → [0.6   0.2   1.0   0.9]

queen → [0.3   0.2   0.4   1.0].

If both words are similar or same than Vector of the both words will be similar because of we are finding Semantic.

The problem of " Average Meaning "

1> An apple a day keeps the doctor away
2> Apple is healthy.
3> Apple is better than orange.
4> Apple makes great phones

:
:

→ [ ] 〰 → [ : ]

                                                                2 dim

[x  y]
taste ↙  ↘ technology

1st line → Apple Khane/taste → [x    y]
         ki baat ho rhi        [0.6   0]

2nd line → Mone sure ]
          Apple khane/taste  [0.7   0]
          ki baat ho rhi

3rd line → Apple khane/taste    [0.8   0]
          ki baat ho rhi

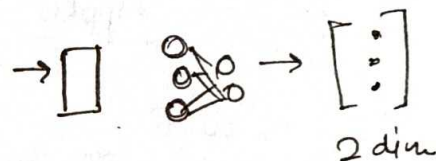4th line → Technology ki    [0.8   0.2]
          bat ho rhi

Total sentence → 10000

    9000 sentence        1000 sentence ]     → overall vector
         ↓                    ↓                  [0.9   0.3]
    Fruits/taste          Tech             move tilted to taste
                                           than tech

Tech



taste

Data → tilted toward taste than Tech

Data → tilted toward tech than taste

## Problem

Word embedding create one time and use many time.

└→ Static

└→ ek baar embedding ban gya to har bar wahi use karn erai.

eg:

└→ Eng to hindi Translation

Apple launched a new phone while I was eating an orange.

\* But our data is tilted toward taste so this Apple treated as a fruit not technology. [0.9    0.3]

But we want to change value based on content.
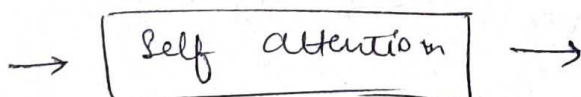
So, this problem solve self Attention.

e apple                       Y apple

e launch   →   [ Self attention ]   →   Y launch

e phone                       Y phon

e orange                       Y orange

└→ embeddings               └→ New embedding

(Smart conceptual embedding)

Humans                    Love                    Smartphones
  ↓                        ↓                          ↓
→ [ | | ] $e_1$      → [ | | ] $e_2$        → [ | | ] $e_3$

  ↓                        ↓                          ↓

→ [ Calculation →      Self    Attention                        ]

  ↓                        ↓                          ↓

→ [ | | ]            → [ | | ]              → [ | | ]
   $Y_1$                  $Y_2$                    $Y_3$

## First principle Approach

| money   bank   grows |                    | rivers   bank   flows |

                ↓                                          ↓

bank → 0.3 money + 0.7 bank +              bank → 0.5 river + 0.4 bank
          0.1 grows                                  + 0.1 flows

So, bank word not only
made with bank but also
made with other words too

money = 0.7 money + 0.2 bank + 0.1 grows

bank = 0.25 money + 0.7 bank + 0.05 grow

grows = 0.1 money + 0.2 bank + 0.7 grows

$$\text{river} = 0.8 \text{ river} + 0.15 \text{ bank} + 0.05 \text{ flows}$$

$$\text{bank} = 0.2 \text{ river} + 0.78 \text{ bank} + 0.02 \text{ flows}$$

$$\text{flows} = 0.4 \text{ river} + 0.01 \text{ bank} + 0.59 \text{ flows}$$

$n$ dim vector [ ][ ][ ]

$$e_{money}^{(new)} = \boxed{0.7} \; e_{money} + \boxed{0.2} e_{bank} + \boxed{0.1} e_{grows}$$

$$e_{bank}^{(new)} = 0.25 \; e_{money} + 0.7 \; e_{bank} + 0.05 \; e_{grows}$$

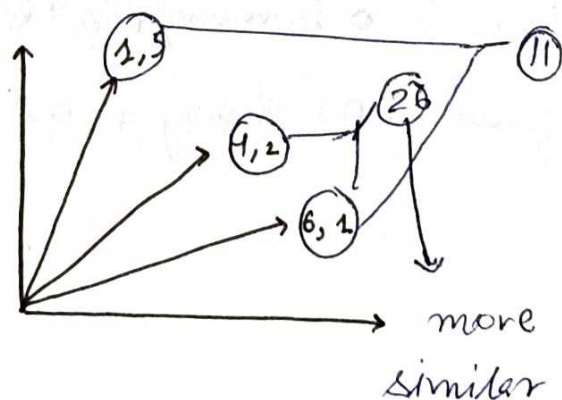$$e_{grows}^{(new)} = 0.1 \; e_{money} + 0.2 \; e_{bank} + 0.7 \; e_{grows}$$

$$
\begin{bmatrix}
0.7 \text{ times old embedding money se bna hai} \\
0.2 \text{ times old embedding bank se bna hai} \\
0.1 \text{ time old embedding grows se bna hai}
\end{bmatrix}
$$

a we can say

$$
\begin{bmatrix}
0.7 \rightarrow \text{similarity bet}^n \text{ money embeddig and money embeddy} \\
0.2 \rightarrow \text{similarity bet}^n \text{ money embeddig and bank embdig} \\
0.1 \rightarrow \text{similarity bet}^n \text{ money embeddig and grows embeddig}
\end{bmatrix}
$$

And All three embeddig are vector. And Dot product bet$^n$ vectors is known as similarity

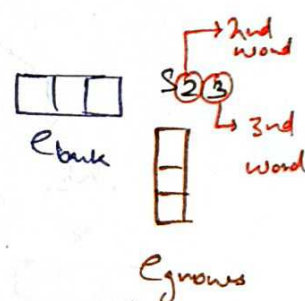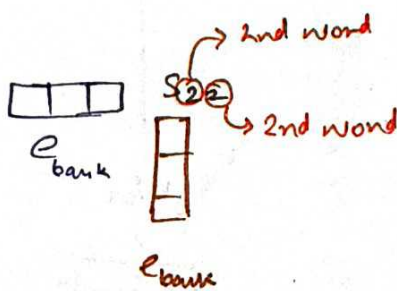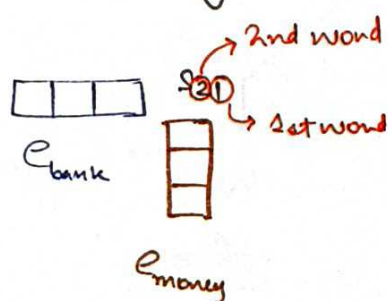$$e_{bank}^{(new)} = \boxed{0.25}\ e_{money} + \boxed{0.7}\ e_{bank} + \boxed{0.05}\ e_{grows}$$

$$e_{bank}^{(new)} = \left[ e_{bank} \cdot e_{money}^T \right] e_{money} + \left[ e_{bank} \cdot e_{bank}^T \right] e_{bank} +$$

$$\left[ e_{bank} \cdot e_{grows}^T \right] e_{grows}$$

Similarity

$$0.25 + 0.7 + 0.05 = 1$$
↳ Normalized



Similarity diagrams:

→ 2nd word
$S_{21}$
↳ 1st word
$e_{bank}$ ... $e_{money}$

→ 2nd word
$S_{22}$
↳ 2nd word
$e_{bank}$ ... $e_{bank}$

→ 2nd word
$S_{23}$
↳ 3rd word
$e_{bank}$ ... $e_{grows}$

Soft Max

$W_{21}$
$$W_{21} = \frac{e^{S_{21}}}{e^{S_{21}} + e^{S_{22}} + e^{S_{23}}}$$
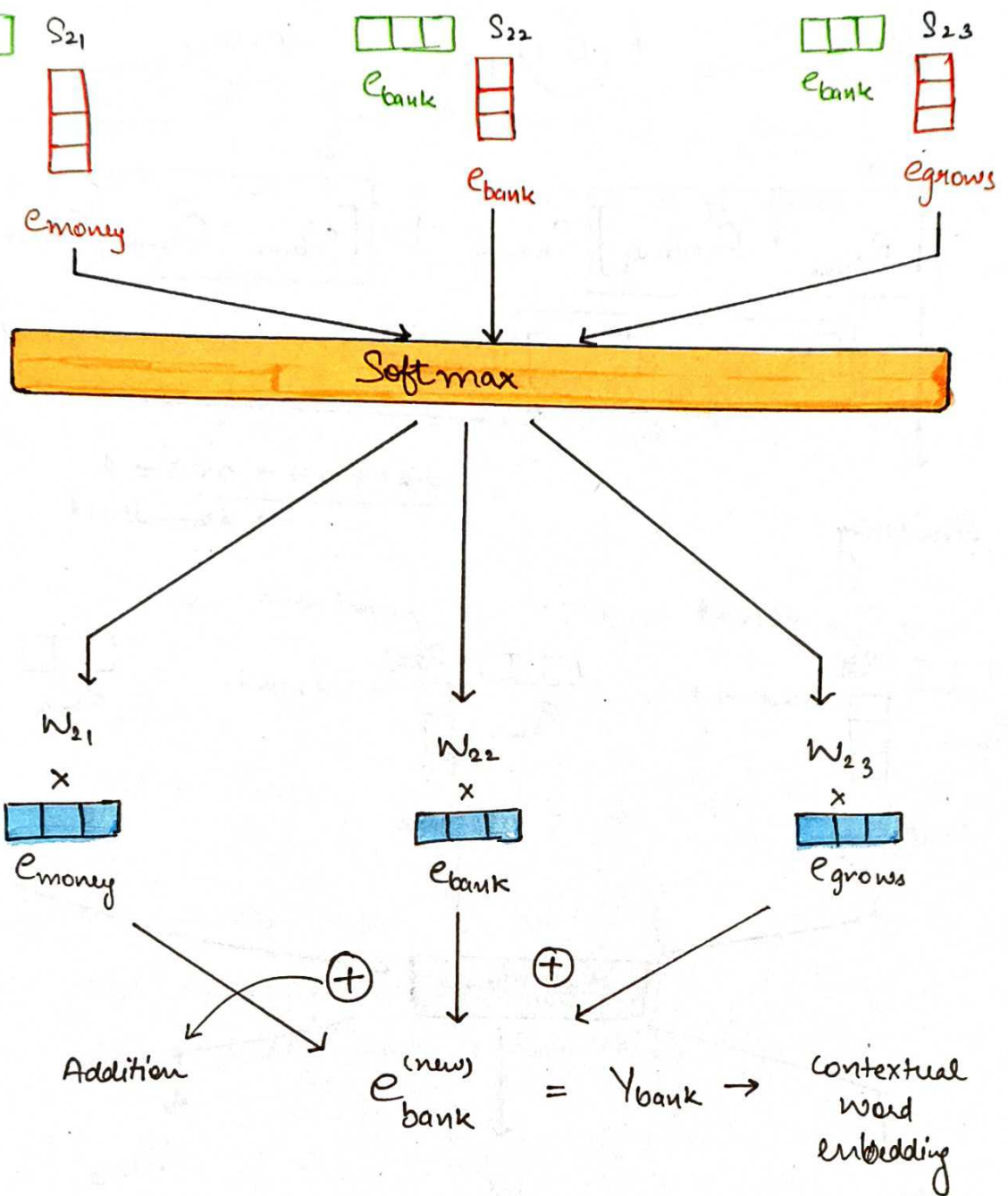
$W_{22}$
$$W_{22} = \frac{e^{S_{22}}}{e^{S_{21}} + e^{S_{22}} + e^{S_{23}}}$$

$W_{23}$
$$W_{23} = \frac{e^{S_{23}}}{e^{S_{21}} + e^{S_{22}} + e^{S_{23}}}$$
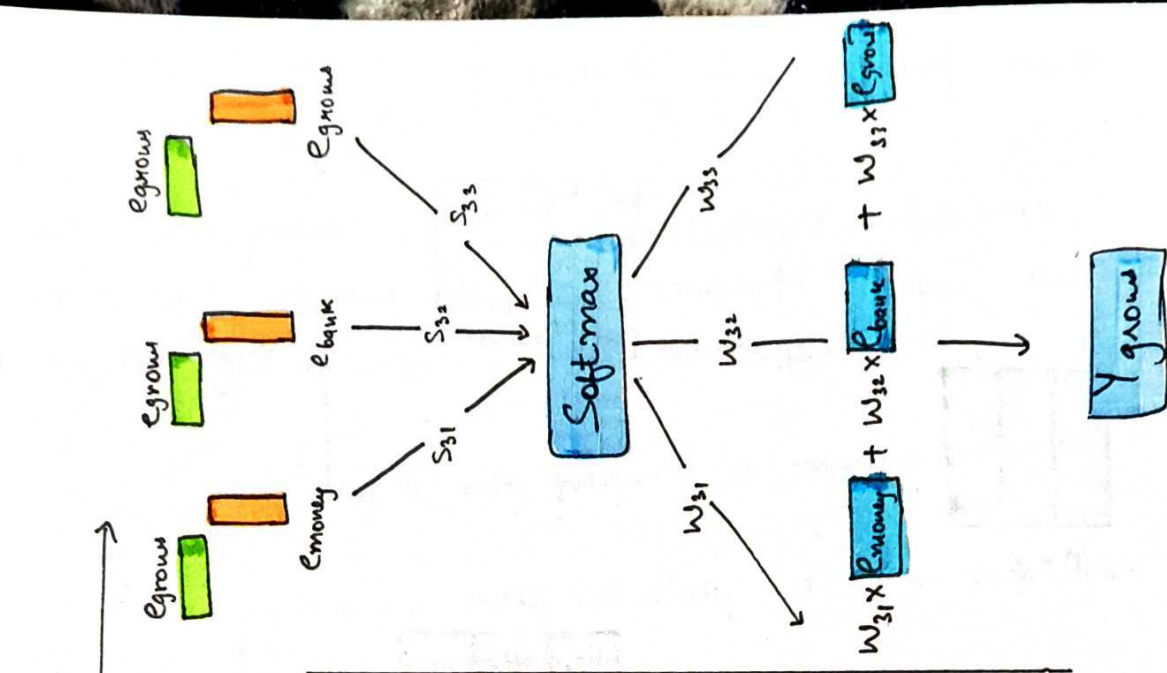
$W_{21}$, $W_{23}$ and $W_{22}$ is Normalized because sum of $W_{21} + W_{23} + W_{22}$ is 1.

$S_{21}$    $e_{bank}$    $e_{money}$

$S_{22}$    $e_{bank}$    $e_{bank}$

$S_{23}$    $e_{bank}$    $e_{grows}$

Softmax

$W_{21}$
$\times$
$e_{money}$

$W_{22}$
$\times$
$e_{bank}$

$W_{23}$
$\times$
$e_{grows}$

Addition

$e_{bank}^{(new)}$ $=$ $Y_{bank}$ $\rightarrow$ Contextual word embedding

This diagram is only for $Y_{bank}$
and $Y_{bank}$ is Contextual word embedding

**Top section (grows):**

$e_{grows}$, $e_{bank}$, $e_{money}$

$S_{33}$, $S_{32}$, $S_{31}$ → Softmax

$W_{33}$, $W_{32}$, $W_{31}$

$$W_{31} \times e_{money} + W_{32} \times e_{bank} + W_{33} \times e_{grows}$$

$Y_{grows}$

**Middle section (bank):**

grows, bank, money

$e_{grows}$, $e_{bank}$, $e_{money}$

$S_{23}$, $S_{22}$, $S_{21}$ → Softmax

$W_{23}$, $W_{22}$, $W_{21}$

$$W_{21} \times e_{money} + W_{22} \times e_{bank} + W_{23} \times e_{grows}$$

$Y_{bank}$

**Bottom section (money):**

$e_{money}$, $e_{bank}$, $e_{grows}$

$S_{23}$, $S_{12}$, $S_{11}$ → Softmax

$W_{13}$, $W_{12}$, $W_{11}$

$$W_{11} \times e_{money} + W_{12} \times e_{bank} + W_{13} + e_{grows}$$

$Y_{money}$

# How to ttrue process in Parallel?

| e money |
|---------|
| e bank |
| e grows |

$3 * n$

embedding vector size

No. of words

$n * 3$

| $S_{11}$ | $S_{12}$ | $S_{13}$ |
|----------|----------|----------|
| $S_{21}$ | $S_{22}$ | $S_{23}$ |
| $S_{31}$ | $S_{32}$ | $S_{33}$ |

Softmax

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
|----------|----------|----------|
| $W_{21}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}$ |

$3 * 3$

| Y money |
|---------|
| Y bank |
| Y grows |

$3 * n$

| e money |
|---------|
| e bank |
| e grows |

$3 * n$

→ There is not Learning Parameter.

## Problem

Our Approach is general contexual embedding Not task specific contextual embedding.

eg:-    piece   of   cake → केक का टुकड़ा

          ↓       ↓    ↓

          $e_1$   $e_2$  $e_3$
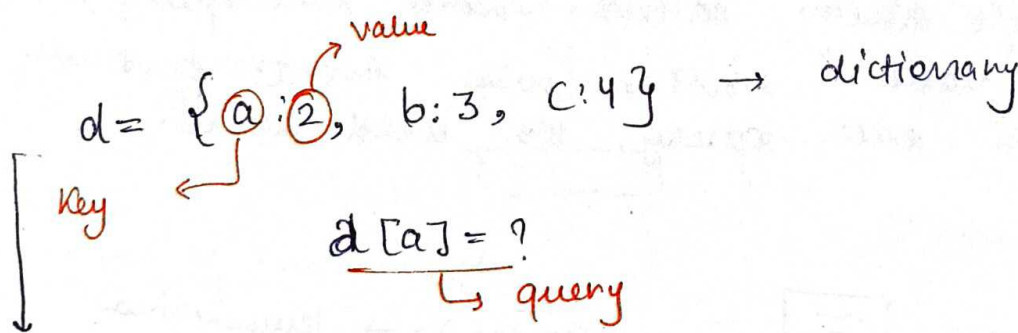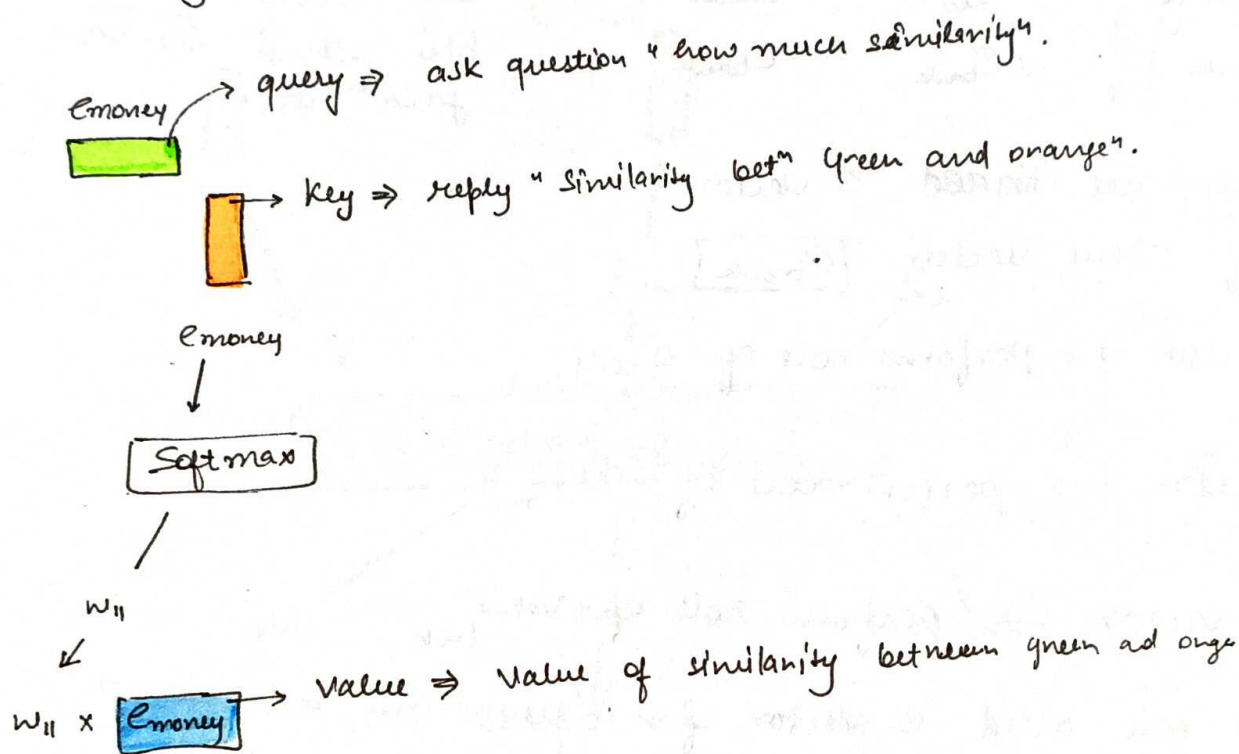
          ↓       ↓    ↓

          $Y_1$   $Y_2$  $Y_3$

* Because we are using general contextual embedding our output is बहुत /आसान काम-
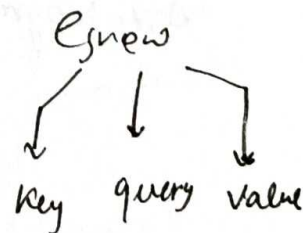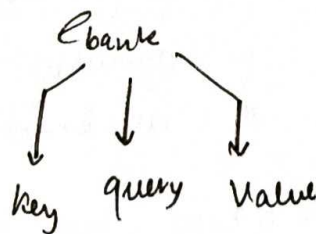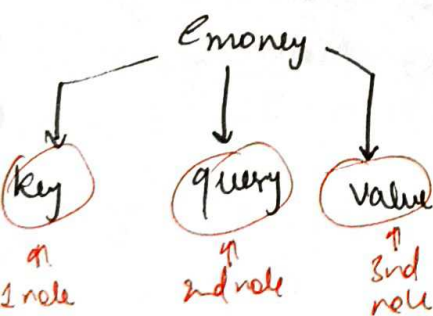~~काम्म~~ "केक का टुकड़ा". kabhi bhi "बहुत आसान काम" output nhi kar sakte.

but in my data

piece of cake → बहुत आसान काम

If Task specific contextual embedding use then might be
output is "बड़े आसमान काम".

At some point general contextual embedding will fail.
If I am doing sentiment analysis then embedding →
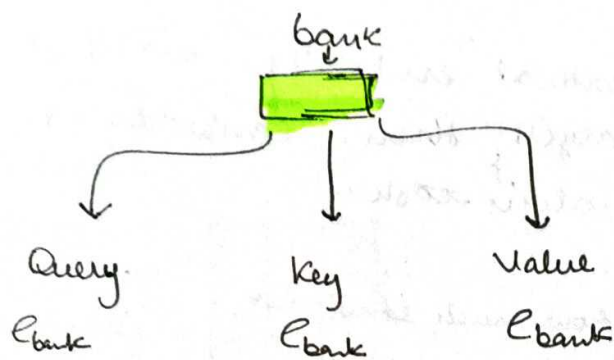accordingly to sentiment analysis task.



query ⇒ ask question " how much similarity".

$E_{money}$

key ⇒ reply " Similarity bet$^n$ green and orange".

$e_{money}$
↓
Softmax
/
$W_{11}$
↙
$W_{11} \times E_{money}$ → value ⇒ value of similarity between green ad orge

$$d = \{ @:2, \ b:3, \ c:4 \} \to dictionary$$

value

key

$$d[a] = ?$$
↳ query

→ doing same thing on diagram.
→ every embedding play 3 roles.



$E_{money}$

Key          query          value
1st role     2nd role       3rd role

$E_{bank}$

key    query    value

$E_{new}$

key   query   value

# Query , Key & Value Vectors



**Problem**

Ye query bhi khud hi ban jata hai. Key bhi and value bhi khud hi ban jata hai.

* ideally we need 3 vector of this vector. [Cbank]
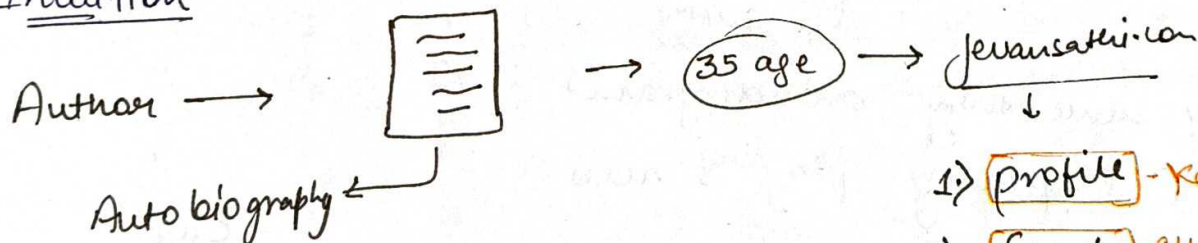
1 vector → perform role of query bank

2 vector → perform role of key bank

3 vector → perform role of Value bank

Why we need 3 vector for query, key and value?

because single vector cannot work like query, key and value. That is why create particular vector for each query, key and value.
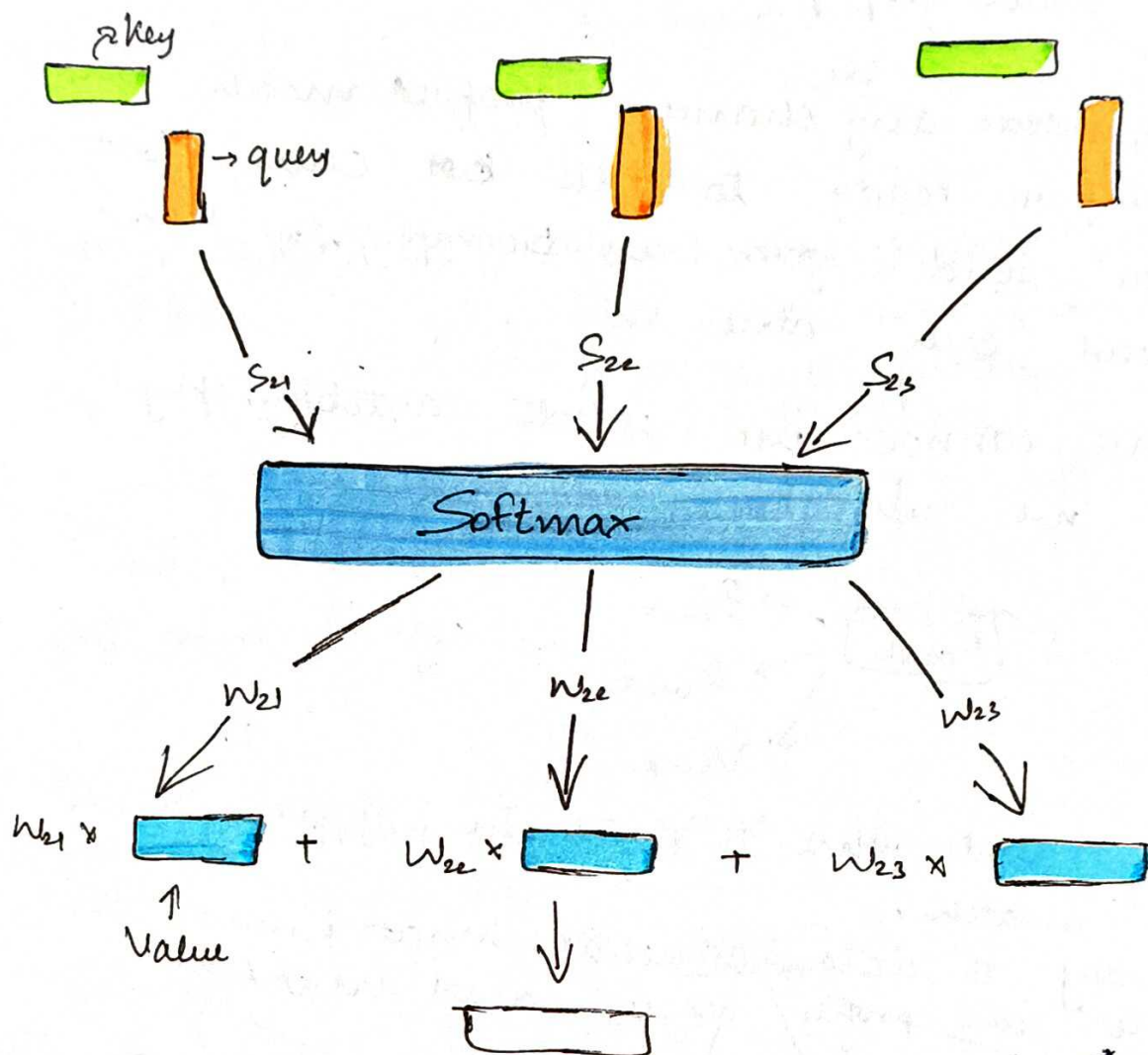
## Intuition

Author →  → (35 age) → jevansathi.com

Auto biography

1) [profile] - Key

2) [Search] query →

3) [Match] → value

query → qualification of girl, state of living, hobbies

key → My profile ⇒ girl know about me.

value → After match ⇒ Start conversation.



Conclusion of this example:

Embedding → key
→ query
→ value

1. If I want to share my profile to girl then I'll share personalities, good things, hobbies. If I'll store my Autobiography instead of Simple profile because in autobiography already written so, it will create awkwardness to other girl.

2. If I want search girl and I write all autobiography in search bar because In autobiography also wrote

about which type of girls I like to marry. Wrote whole biography in search bar and may be I give wrong girls option because of long autobiography.

suggest

3. By chance any chance profiles match and Girl is ready to talk but chat but you send your autobiography eto know about each other.

So, we cannot use whole autobiography. That's we use key, query, value.

$$E_{bank} \rightarrow Q_{bank}$$
$$\searrow K_{bank}$$
$$\searrow V_{bank}$$

How I decide what to write in Profile, Search and match.

→ Accordiy to data I decide what i want to write in profile search and match.

How?

→ First I write → I wrote political books Suggestion Those In profile → girls who are interested in politics

According to data, I understood that I have to change profile from Political book writer to Authors / writer

but I don't like girls who are interested in politics

→ **First 9 write in Search** → Working proffesion / Non-Working ⟶

Decision ← After conuersation with of both type of girls,
|
Data

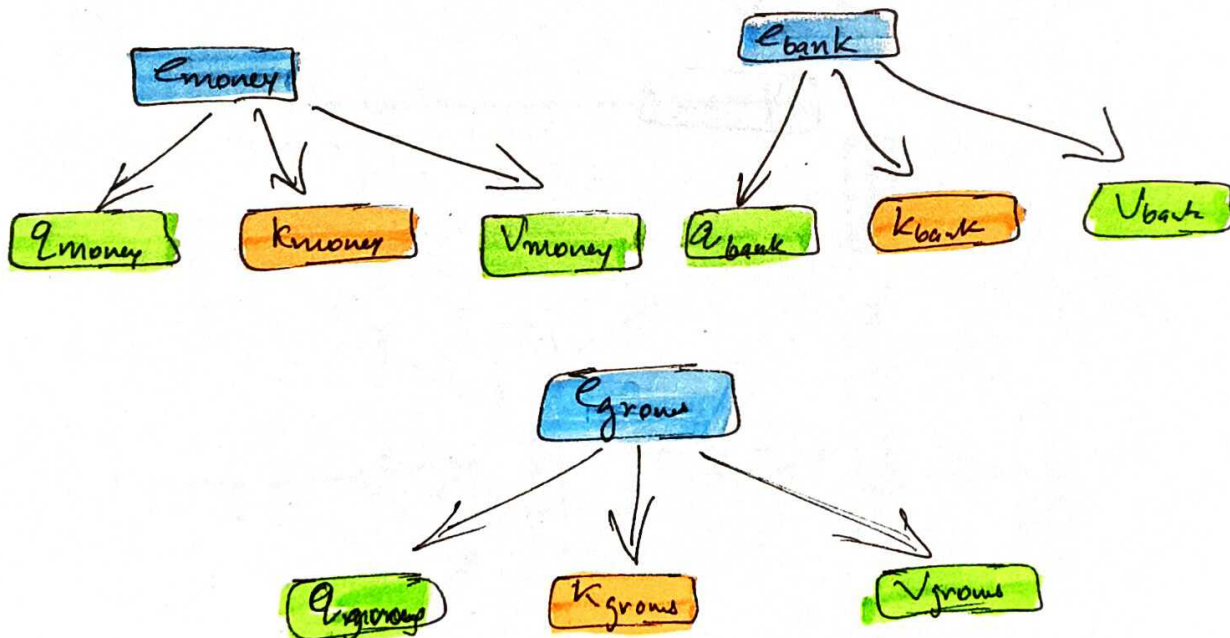9 like only working girls

→ **First 9 write in Match** → After Match 9 talk with girls with more enthusiast And most of the girl not like it. ⟩—Data
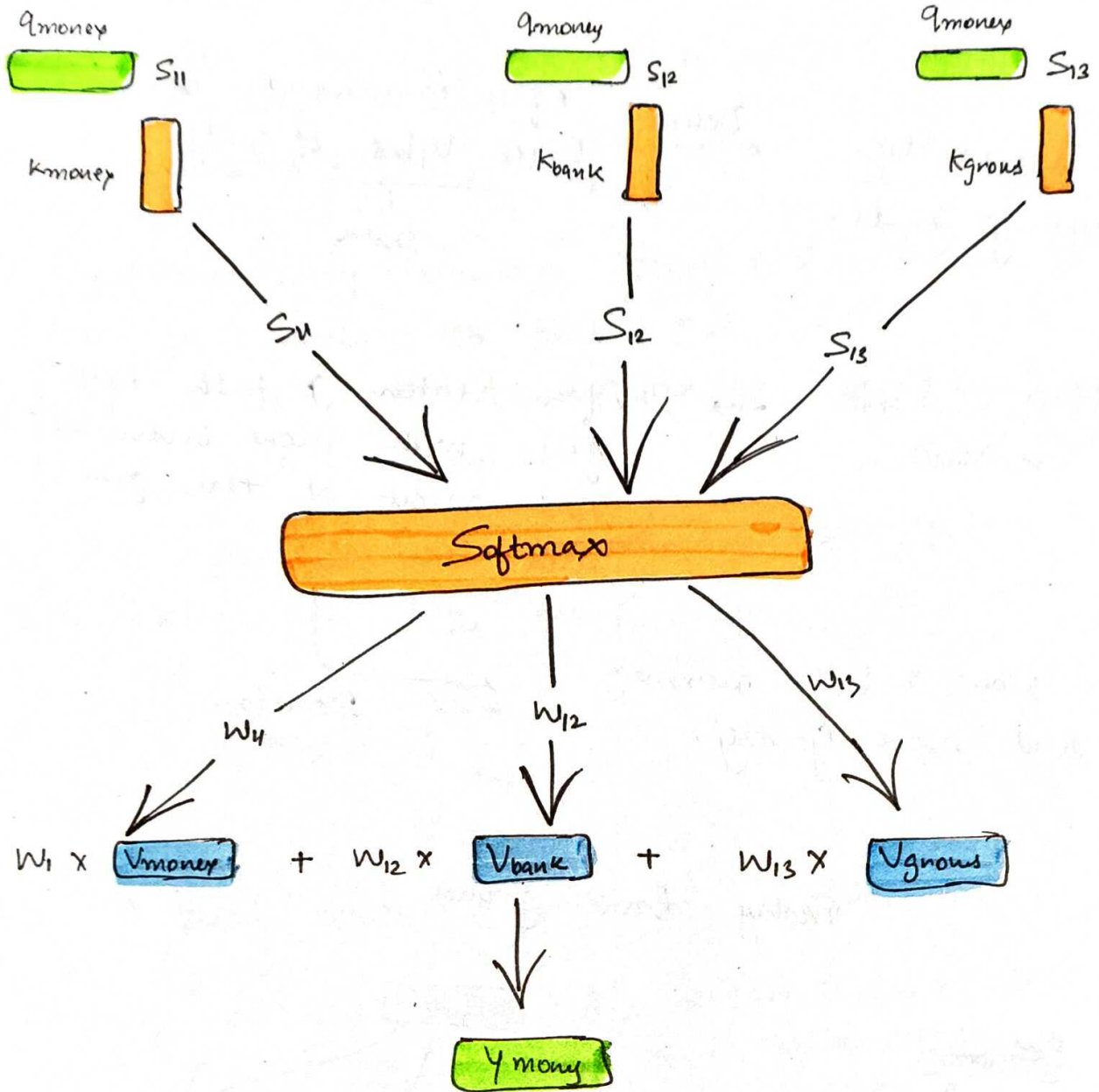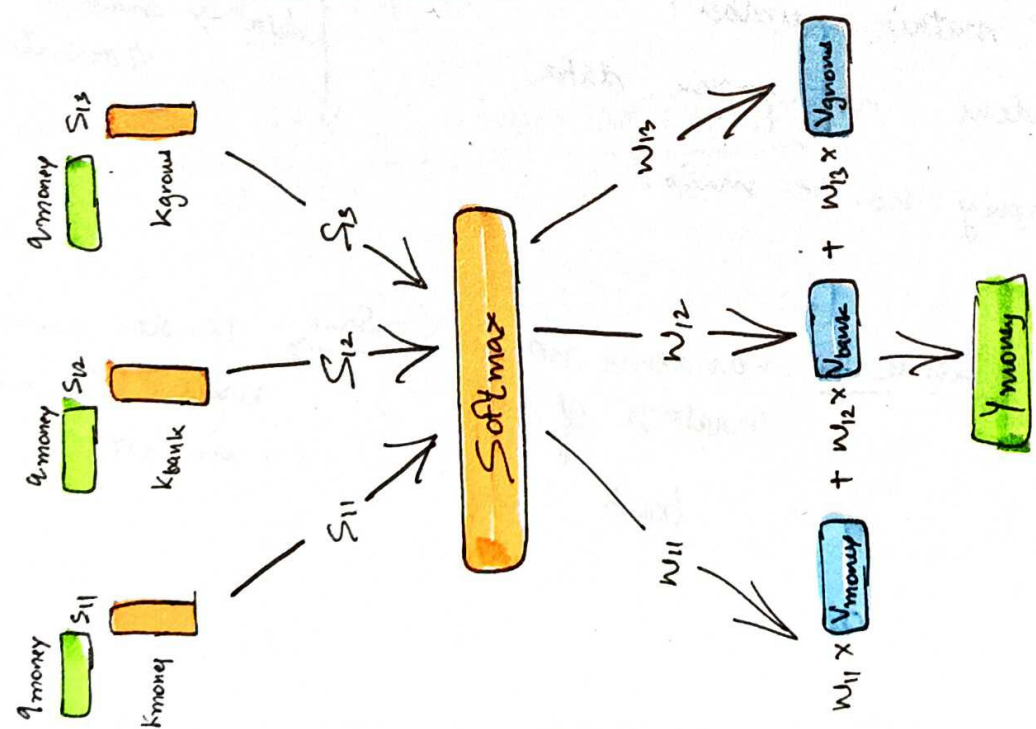
Now, 9 talk normaly and more Gently. ← Decision

Money bank grows

Cmoney
→ 9money  Kmoney  Vmoney

Cbank
→ Abank  Kbank  Vbank

Cgrows
→ 9grows  Kgrows  Vgrows

# Money

$q_{money}$    $S_{11}$

$k_{money}$

$q_{money}$    $S_{12}$

$k_{bank}$

$q_{money}$    $S_{13}$

$k_{grous}$

$S_{11}$     $S_{12}$     $S_{13}$

## Softmax

$W_{11}$    $W_{12}$    $W_{13}$

$W_1 \times$ $V_{money}$ $+$ $W_{12} \times$ $V_{bank}$ $+$ $W_{13} \times$ $V_{grous}$

$Y_{mony}$

**Block 1**

q_money    q_money    q_money

S_11    S_12    S_13

k_money    k_bank    k_grow

S_11    S_12    S_3  → Softmax

W_11    W_12    W_13

V_money    V_bank    V_ground

$$W_{11} \times V_{money} + W_{12} \times V_{bank} + W_{13} \times V_{ground}$$

Y_money

---

**Block 2**

q_book    q_book    q_book

S_11    S_12    S_13

k_mony    k_bank    k_grow

S_11    S_12    S_13  → Softmax

W_11    W_12    W_13

V_mony    V_bank    V_grow

$$W_{11} \times V_{money} + W_{12} \times V_{bank} + W_{13} \times V_{grow}$$

Y_bank

---

**Block 3**

q_grow    q_grow    q_grow

S_11    S_12    S_13

k_mony    k_bank    k_grow

S_11    S_12    S_13  → Softmax

W_11    W_12    W_13

V_money    V_bank    V_...

$$W_{11} \times V_{money} + W_{12} \times V_{bank} + W_{13} \times V_{...}$$

Y_ground

How to make Key vector, query vector, value vector from embedding vector.

query   Key   $e_{bank}$
value

(i) magnitude (Scaling)

(ii) Linear transform

$$\begin{bmatrix} & \\ & \end{bmatrix} \rightarrow$$ Multiply with Matrix

$e_{bank}$

matrix

$W_q$

$q_{bank}$

$W_k$

$K_{bank}$

$W_v$

$V_{bank}$

eg:
data → machine translation

How to find matrix number?
→ first random no. from data
→ Then improvey no. in matrix.

random no. in   Same   random no. in   Same   random no. in
matrix of $W_q$  →     matrix of $W_q$   →    matrix of $W_v$
   Money                   bank                  in green

random no. of matrix $\xrightarrow{\text{Same}}$ random no. of matrix in $\xrightarrow{\text{Same}}$ random no. of matrix in Wk (money) $\qquad$ $W_k$ (bank) $\qquad$ $W_k$ (grows)

random no. of matrix in $W_v$ $\xrightarrow{\text{Same}}$ random no. of matrix in $W_v$ bank $\xrightarrow{\text{Same}}$ random no. of matrix in $W_v$ grows

money
↓
emoney

bank
↓
ebank

grows
↓
egrows

Qmoney  Kmoney  Vmoney

Qbank   Kbank   Vbank

Qgrows  Kgrows  Vgrows

All this process are parallely.

Mathematical formula

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