$$\text{Var}(z) > \text{Var}(y) > \text{Var}(x)$$

$$\text{Var}(z) \simeq 3 \text{Var}(x)$$

If $d$ dim then $d_m \to d \text{Var}(x)$

* Somehow
$$\begin{bmatrix} ac + bd \\ ae + bf \\ \vdots \end{bmatrix}$$
aise number se
← ise divide karna hai
ki every time $\text{Var}(x)$ hi aaye

$$\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} c & d \end{bmatrix}$$
$$\frac{}{\sqrt{2}}$$
$$\begin{bmatrix} e & f \end{bmatrix}$$
$$\frac{}{\sqrt{2}}$$
dim
$$\begin{bmatrix} g & h \end{bmatrix}$$
$$\frac{}{\sqrt{2}} \to \sqrt{d_k}$$

$\left. \phantom{xxx} \right\} 2 \text{Var}(x)$
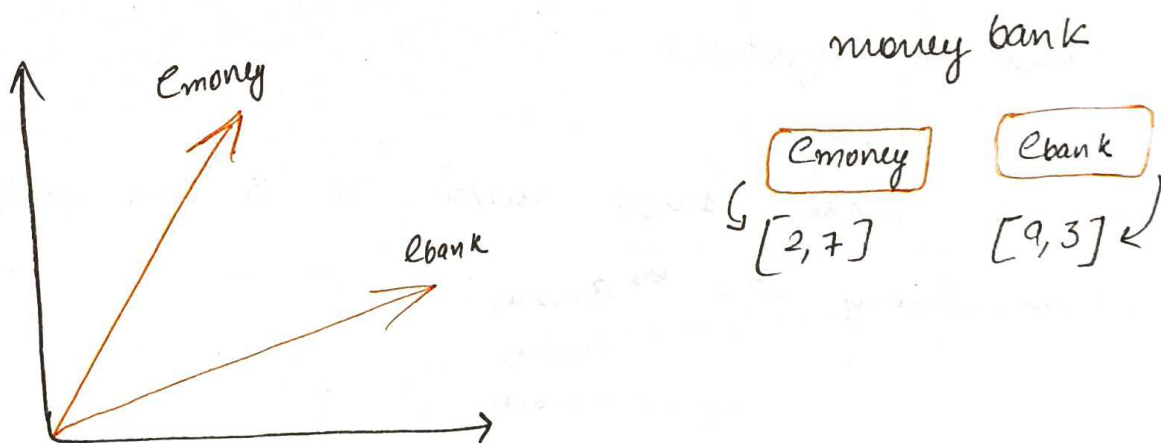
$\frac{1}{2} \text{Var}(y) \to \frac{1}{2} 2 \text{Var}(x) = \text{Var}(x)$

$d_k \to$ dimension

<u>**Self Attention Geometric Intuition**</u>
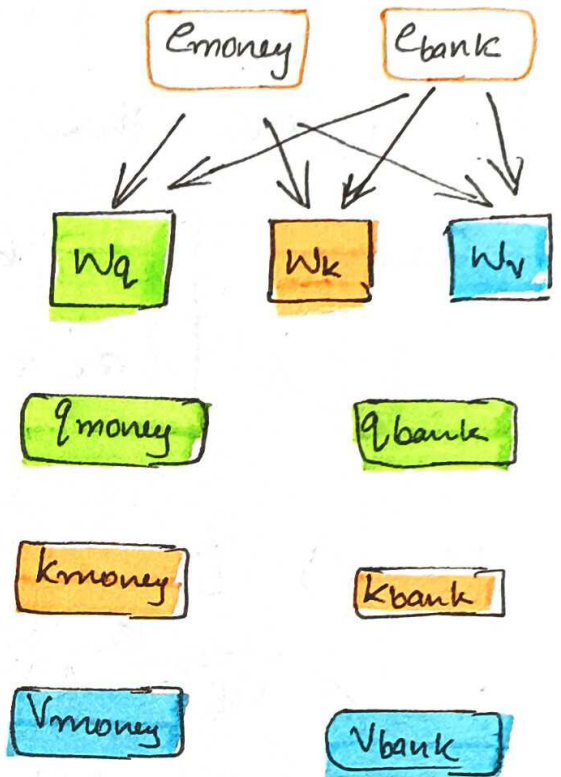


money bank

$e_{money}$  $e_{bank}$

$\hookrightarrow [2, 7]$  $[9, 3] \hookleftarrow$

money bank

dot product $e_{money}$

$W_q$       $W_k$       $W_v$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} \quad \begin{bmatrix} 4 & 1 \\ 2 & 1 \end{bmatrix}$$

2×2      2×2      2×2

$V_{money}$

$K_{money}$

$q_{money}$

$e_{money}$     $e_{bank}$

$W_q$       $W_k$       $W_v$

$q_{money}$         $q_{bank}$

$K_{money}$        $K_{bank}$

$V_{money}$        $V_{bank}$

$V_{money}$     $V_{bank}$

$q_{money}$

$K_{bank}$

$K_{money}$

$q_{bank}$

dot product $e_{bank}$

$W_q$       $W_k$       $W_v$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} \quad \begin{bmatrix} 4 & 1 \\ 2 & 1 \end{bmatrix}$$
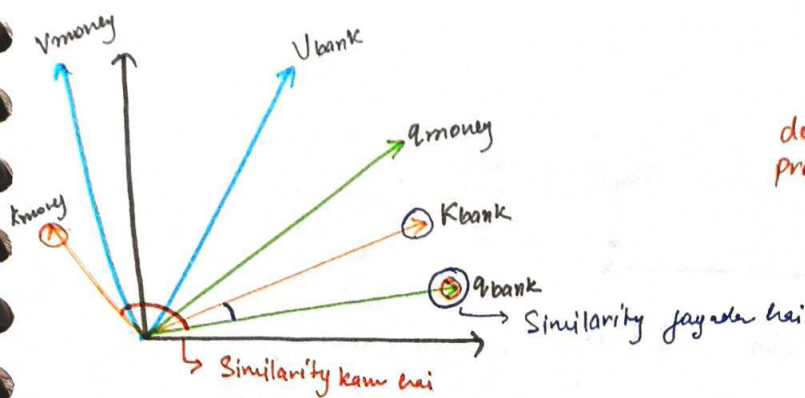
\* All values are hypothetial

Here we make single vector to 3 new vector

from $e_{mony}$ $\xrightarrow{to}$ $q_{money}$

                 $K_{mony}$

                 $V_{mony}$

Similarity jayada hai
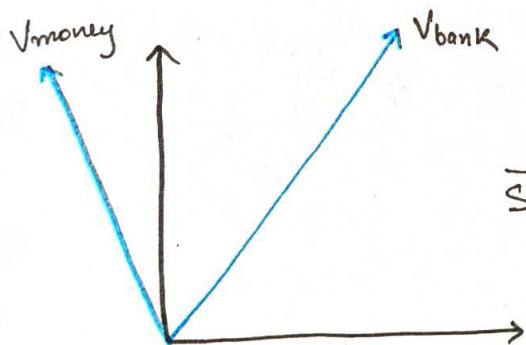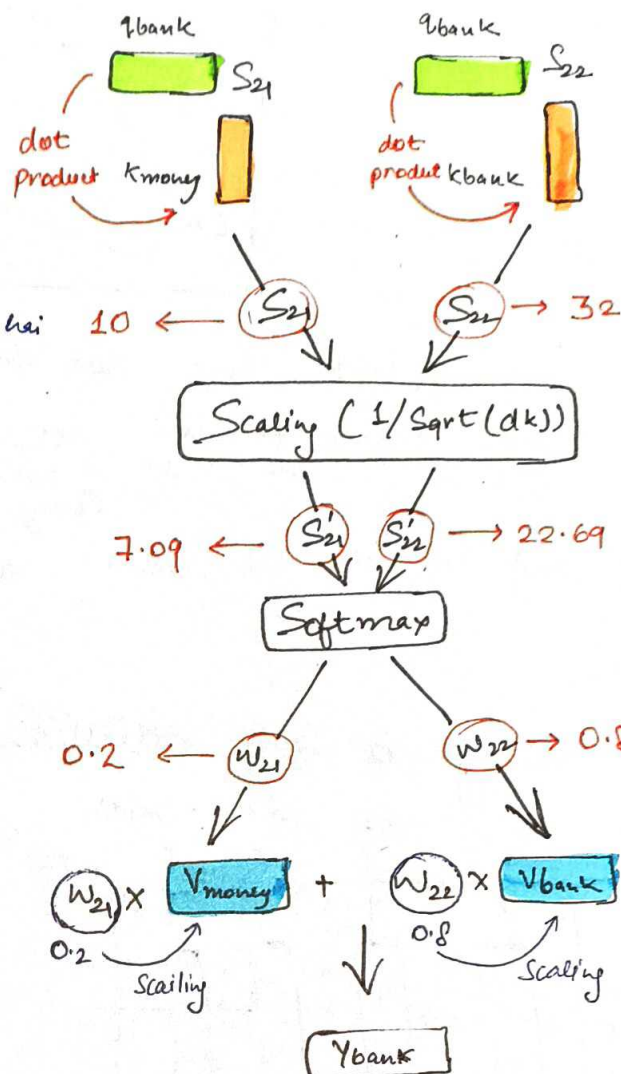
Similarity kam hai

## Dot Product

$$S_{21} = 10 \qquad S_{22} = 32$$
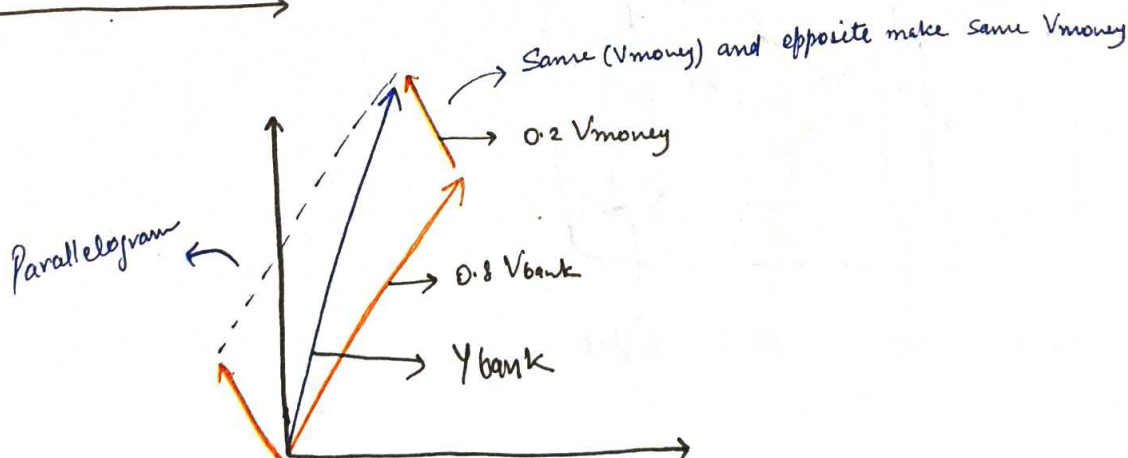
## Scaling

$$S'_{21} = \frac{10}{\sqrt{2}} = 7.09 \qquad S'_{22} = \frac{32}{\sqrt{2}} = 22.69$$

## Softmax

$$W_{21} = 0.2 \qquad W_{22} = 0.8$$



Scaling



10 ← $S_{21}$

$S_{22}$ → 32

Scaling ( 1/Sqrt (dk))

7.09 ← $S'_{21}$ $S'_{22}$ → 22.69

Softmax

0.2 ← $W_{21}$

$W_{22}$ → 0.8

$W_{21}$ × Vmoney + $W_{22}$ × Vbank

0.2 Scaling 0.8 Scaling

Ybank



0.8 Vbank

0.2 Vmoney

Same (Vmoney) and opposite make same Vmoney



Same (Vmoney) and opposite make same Vmoney

0.2 Vmoney

0.8 Vbank

Ybank

Parallelogram

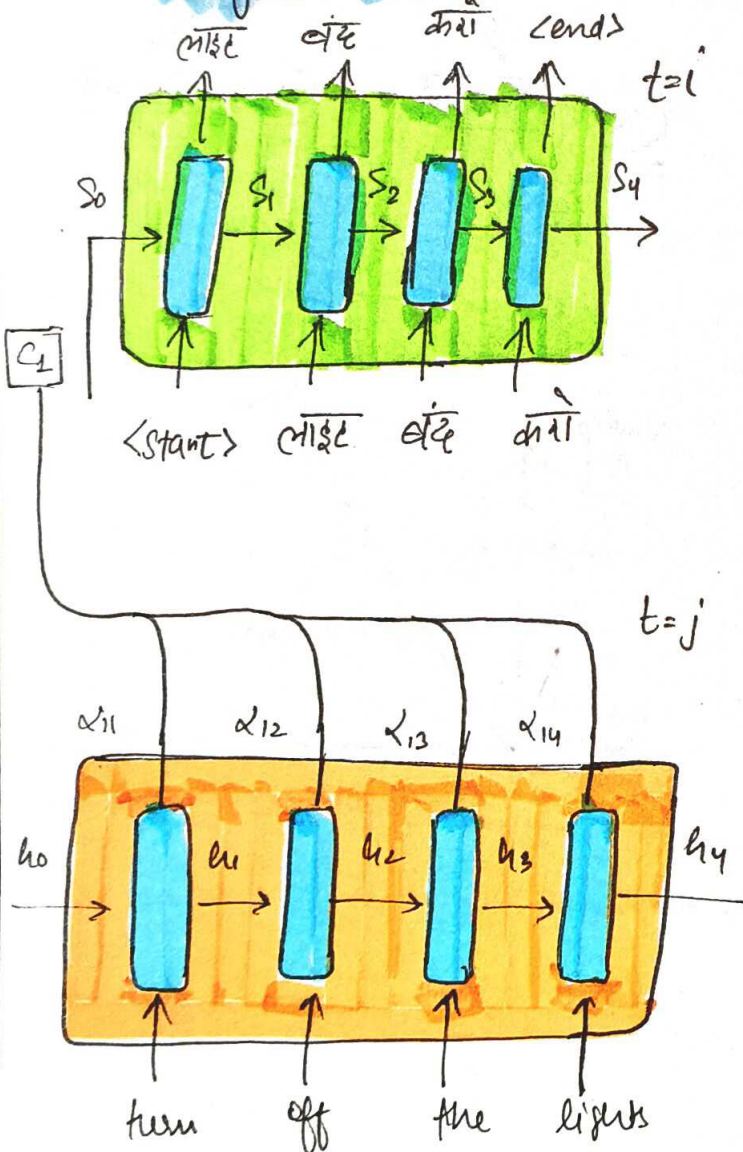* Embedding bank bohot dur tha Embedding money starting mai.
* Now, distance between <u>money</u> and <u>bank</u> is decreased.
  
  $e_{money}$       $y_{bank}$ (new)
* Self Attention is content aware.

## Why is self Attention called "self"?



लाइट   बंद   करो   <end>

$t=i$

$<start>$   लाइट   बंद   करो

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

$$\alpha_{ij} = softmax(e_{ij})$$

$$e_{ij} = S_i^T h_j$$

$t=j$

$\alpha_{11}$   $\alpha_{12}$   $\alpha_{13}$   $\alpha_{14}$

turn   off   the   lights

# Self Attention

[Turn] off the lights

$\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$

$e_1$   $e_2$   $e_3$   $e_4$

$q_{turn}$   $k_{turn}$   $V_{turn}$

$Y_{turn}$

[Turn] off the lights

$q_{turn}$   $q_{off}$   $q_{the}$   $q_{lights}$

$Y_{turn} = W_{11} V_{turn} + W_{12} V_{off} + W_{13} V_{the} + W_{14} V_{lights}$

$(S_{11})$   $(S_{12})$   $(S_{13})$   $(S_{14})$

$\rightarrow$ softmax $(S_{ij}) \Rightarrow W_{11} = $ Softmax $(S_{11})$

$W_{12} = $ Softmax $(S_{12})$

Turn off the lights

$k_{turn}$   $k_{off}$   $k_{the}$   $k_{lights}$

* $Y_{turn}$ similar to $C_i$

* $\alpha_{ij} = $ similar to $W_{11} = $ Softmax $(S_{11})$

  Softmax $(e_{ij})$

* $\begin{bmatrix} S_i \rightarrow \text{query} \\ h_j \rightarrow \text{key} \\ h_j \rightarrow \text{value} \end{bmatrix}$

* Because of these three similarity we called Attention (Luong Attention similar with self Attention)
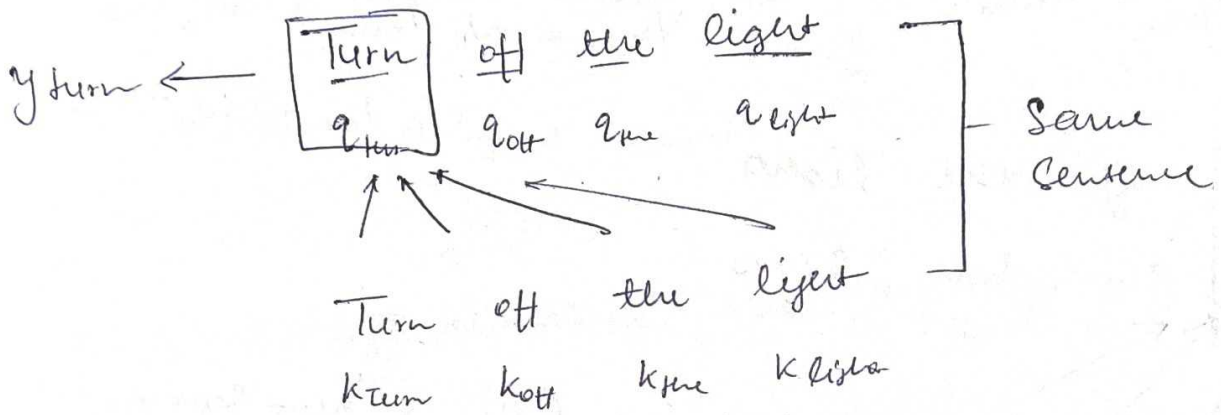
Why "Self" called in Self Attention?

In Luoy Attention calculate between two different sequences.

In Self Attention calculate in a single sequence like calculate similarity between same sentence

Y turn ← $\boxed{\text{Turn}}$ off the light

$q_{turn}$   $q_{off}$   $q_{the}$   $q_{light}$   — Same Sentence

Turn off the light

$K_{turn}$   $K_{off}$   $K_{the}$   $K_{light}$

# Problem with Self Attention

[The man saw the astronomer with a telescope]

Meaning:
1. बंदे ने दूरबीन पर एक एस्ट्रोनॉमर देखा |
2. बंदे ने एस्ट्रोनॉमर को देखा उसके पास दूरबीन था

So, this Sentence is Ambiguity. 2 meaning of single sentence.

But Self Attention capture only one meaning from both meaning (1 and 2).

* How Self Attention works?

The  man,  saw  the  astronomers  with  a  telescope

```
The
 |
Find embeding
 |
generate
 / | \
q  k  v   ← vectors
```

Calculate
Similarity Scores

* Same process for all words.

There is a chance to capture first meaning.

Apply self Attention ⤵

The man saw the astronomer with a telescope

Similarity score
Saw and telescope
high

Similarity score bet$^n$ man and
telescope is very high.

So, Meaning is : Man telescope लगा के देख रहा है

Apply self Attention ⤵

The man saw the astronomer with a telescope

Similarity score
bet$^n$ saw and astronomer
is very

Similarity
Score bet$^n$ astronomer
and telescope is very
high

So Meaning: Man ने Astronomers को देखा And uske हाथ
में telescope tha.

Problem : Self Attention figure out only

Single perspective . If Multiple perspective are
present then self Attention cannot capture.

In NLP, there are multiple scenario to capture multiple perspective.

example → Document Summorization tool.

# Multi-head Attention

The man saw the astronomer with a telescope

↳ 2 meaning

So, we use two self Attention.

money bank

↓      ↓

Emoney      Ebank

Matrics ⟹    $W_q$    $W_k$    $W_v$   →

$W_q^1 \quad W_k^1 \quad W_v^1$

$W_q^2 \quad W_k^2 \quad W_v^2$

if self Attention is three

$W_q^3 \quad W_k^3 \quad W_v^3$

$q_{money}$      $q_{bank}$

$k_{money}$      $k_{bank}$

$V_{money}$      $V_{bank}$

↳ ② set

② set

because two self attention