

Seq2Seq Data

Camp 65X

1) Tabular Data

CGPA, IQ, placement

train it on ANN → DNN
add more layers
Artificial Deep
importance

2) Image Data

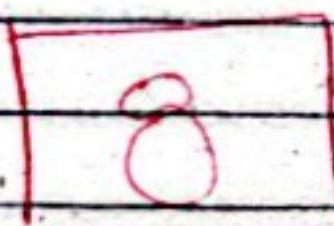


image classification

↓
is made of pixel

2D structure

has some meaning
that can't understand
ANN

then ↓

CNN

3) Sequential Data

→ textual data

→ time series based data

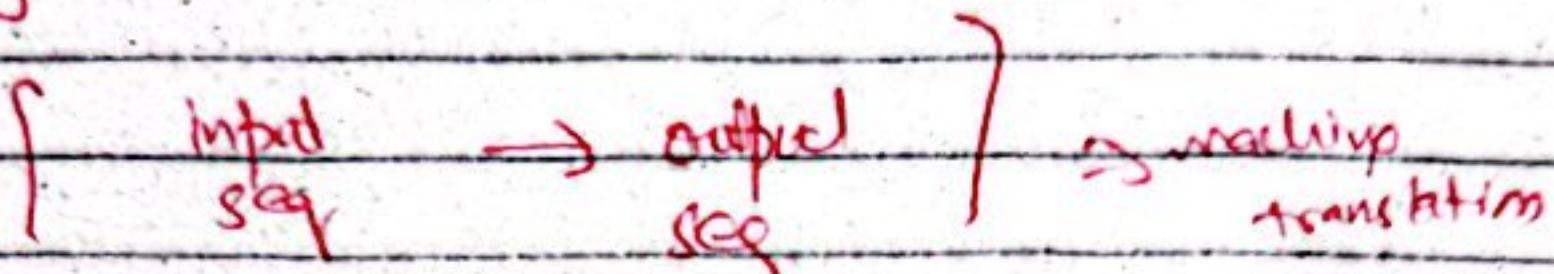
sequence (order) has some important
but unfortunately ANN and CNN
can't understand

→ RNN

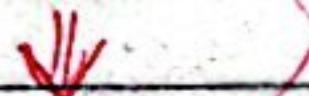
LSTM GRUs

4) Seq2Seq Data:-

Such data where input is also sequence and output is also sequence
e.g machine translation



nice to meet you (4 word)



(b) los / je ~~wie~~ i (16 words)

We made such neural n/w architecture that handle Seq2Seq data.

Q Why it is difficult to solve Seq2Seq Data?

1) input → sentence in some language
↳ variable length.

2) output → sentence in some language
↳ variable length

3) not guaranteed input & output Seq2Seq data at same length

~~with dealing with~~
it mean ^{with dealing with} variable length is difficult

input \rightarrow output
6 word 600 words

variable length

↳ handle we learn in
LSTM & GRU

but only for input

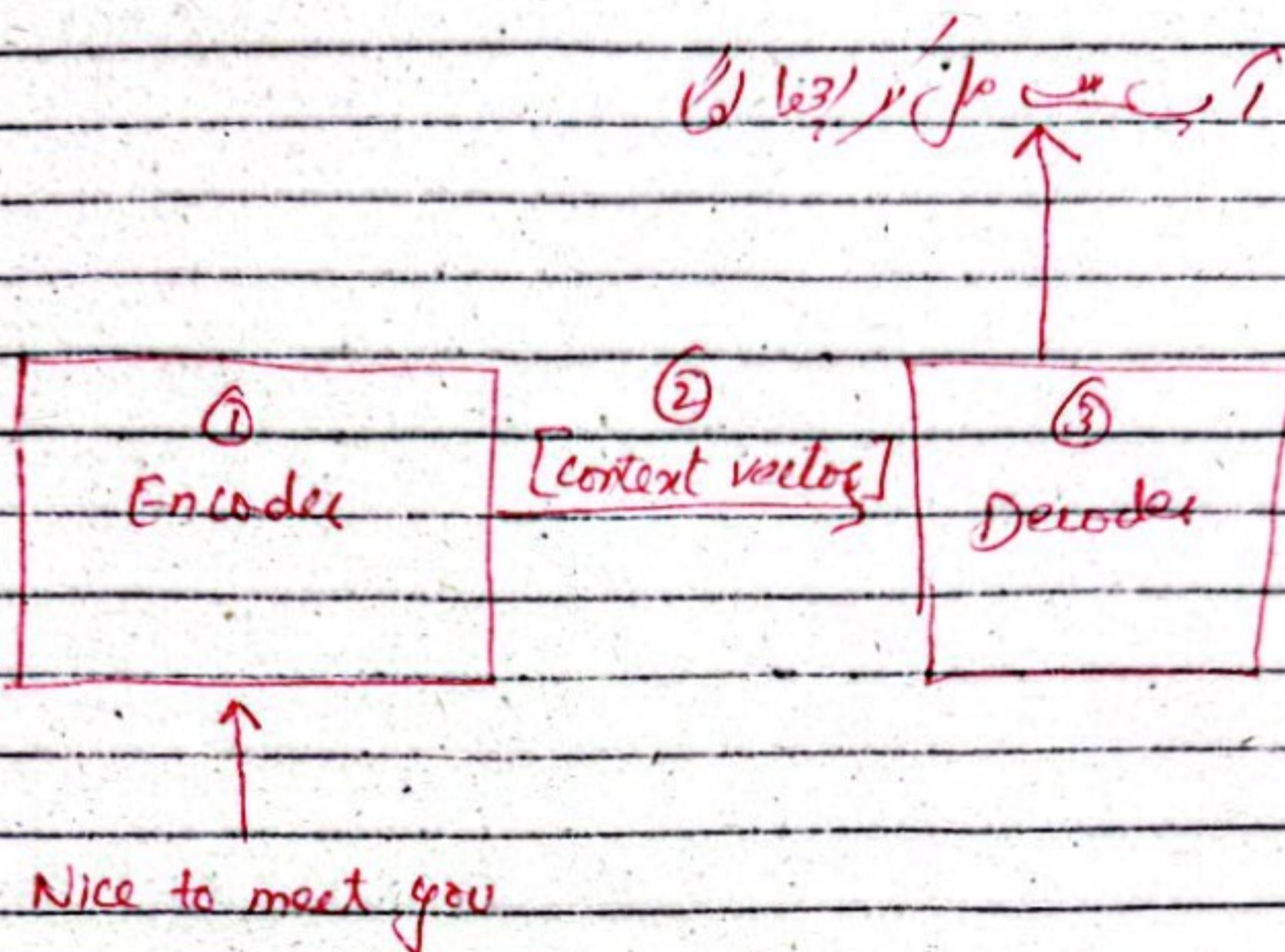
but not for output and what we
learn to solve seq to seq problem.

Encoder Decoder Architecture

Pre-requisite

\rightarrow RNN / LSTM

\rightarrow Sequence to sequence with neural n/w
research paper



Encoder :-

is basically that section of architecture where input sequence give.

↳ to make sure input given token by token (word by word)

Encoder → to understand complete sentence
 ↳ to capture its essence
 ↳ and try to summarize

↳ output
 vector (set of numbers)
 that called (context vector:)

Content vector is basically summary of this sentence.

that made from encoder intelligent

that content vector give to ~~decoder~~

to understand content vector (summary)

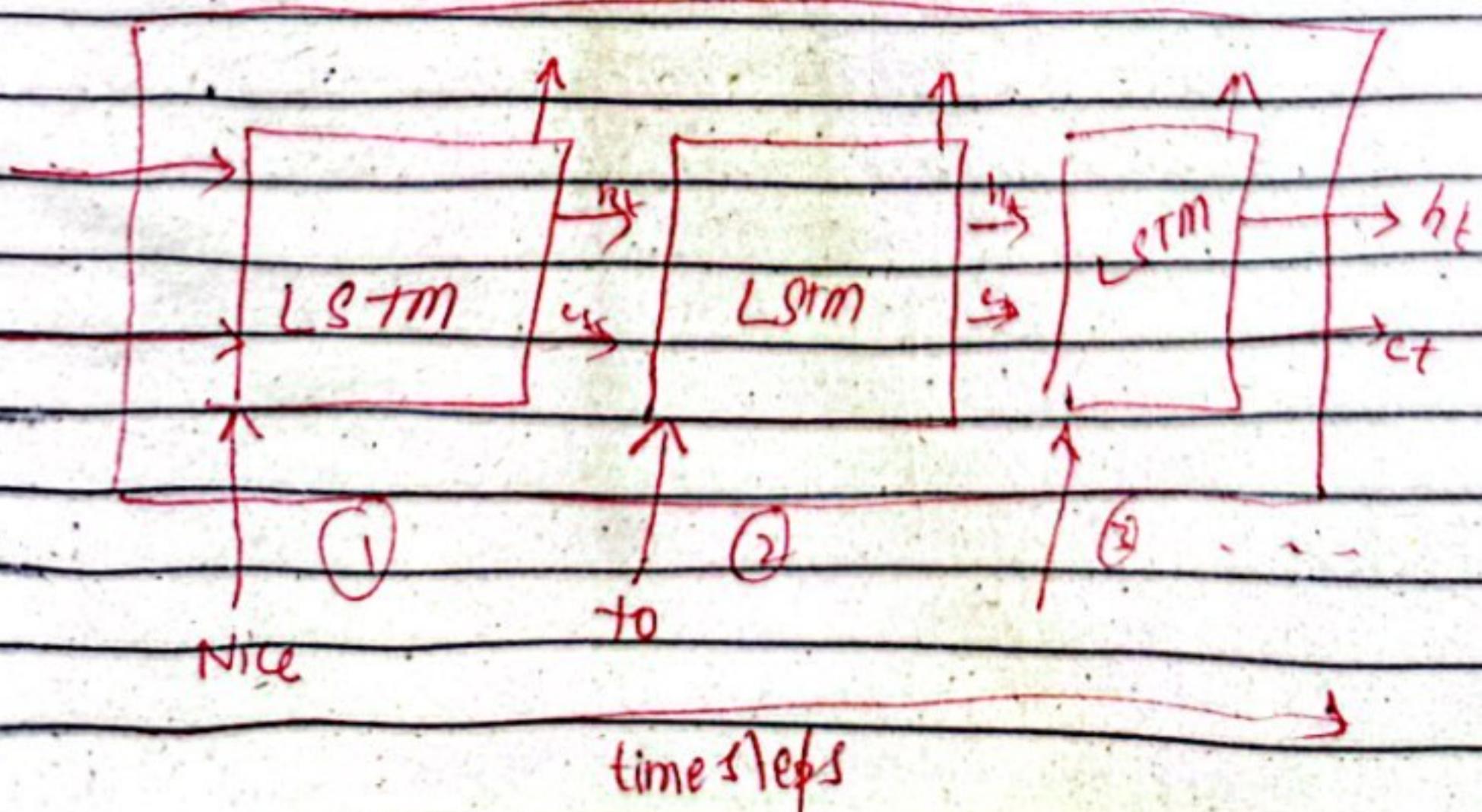


and then taken by layers (word to word) output.

What's under the encoder and decoder

↳ we want such thing - that process sequence in both encoder & decoder.

Encoder \rightarrow  Encoder is one LSTM cell, we unfold over time



In whole process LSTM understand sentence and send its summarization through hidden layers.

After final step, we get h_t and c_t that is our final representation of sentence.

That become the context vector.

That h_t and c_t pass to the decoder.

In summary Encoder is made of one LSTM.

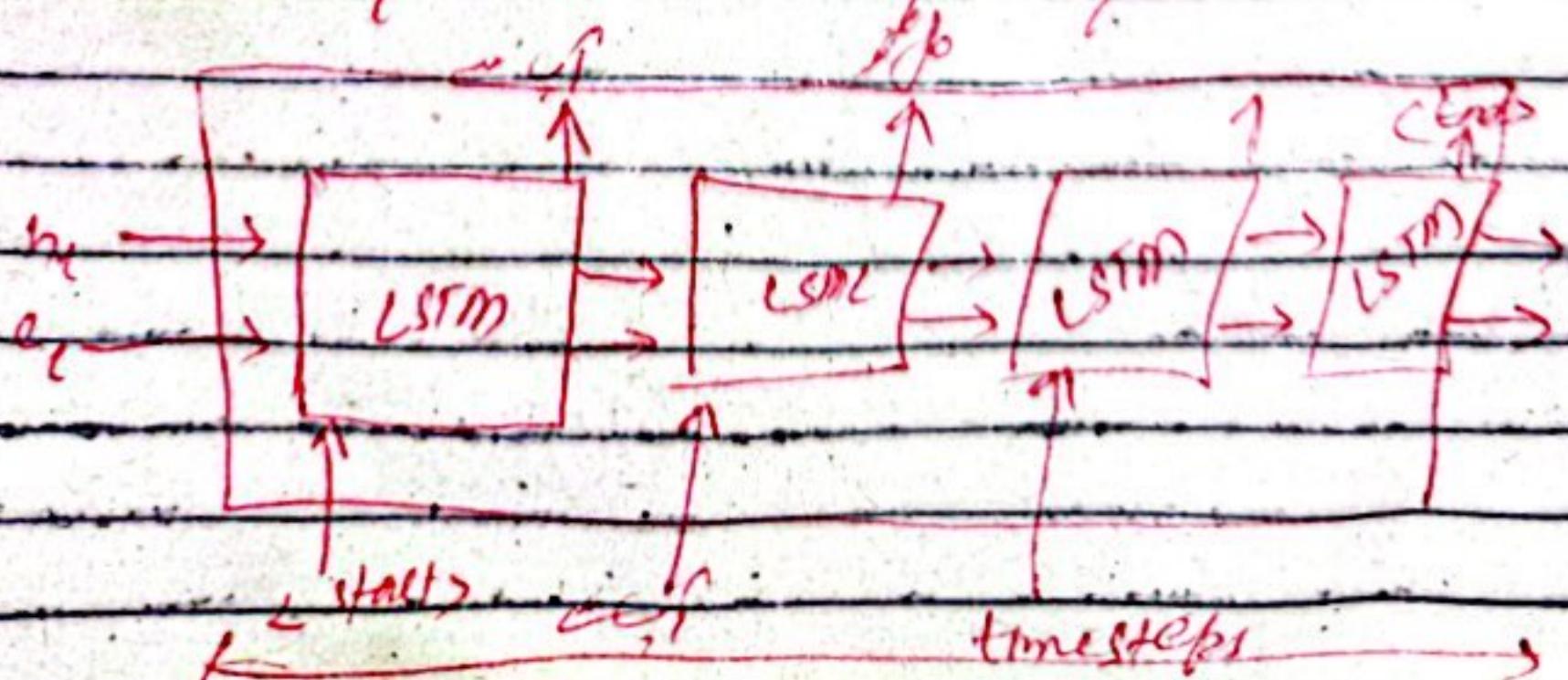
It is not necessary to used LSTM you also used GRU.

RNN generally not used because it has vanishing problem.

Decoder :-

also has one LSTM

→ here LSTM works it in every time step some output produced.



h_t and c_t that we got from encoder
that pass it to decoder through
initial state.

and in that way decoder get
content of input sequence.

at timestep 1 you not only pass content
vector but also pass special symbol
in input which we call <start>.

By seeing <start> decoder understand
now it ^{start} produce output.

and it produce any one word output
at timestep 1 and that output
we pass at input at timestep 2
along with internal state and so on
until <END> token ~~produced~~ produced.
And it stop output produced.

Training

Dataset

Machine translation

English-Hindi-Dataset (Kaggle)

15

45

Encoder

English

- 1) think about it
- 2) come in

WAV

1 2 3
4 5 6

first convert into number.

→ first convert into ASCII

[think, about, it] → [72, 109, 105]
[come, in] → [105, 115, 110]

Now convert into number

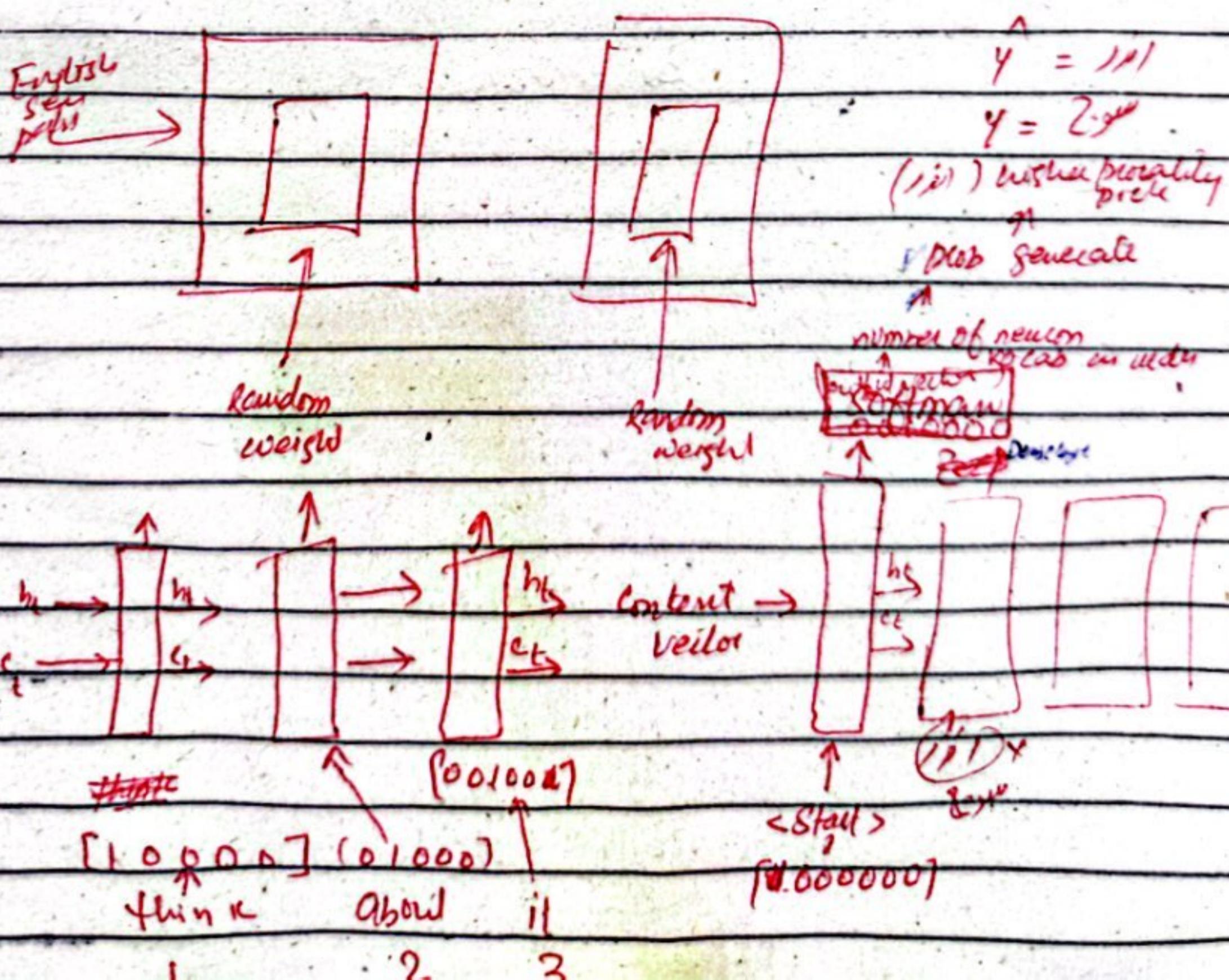
→ one hot encoding
eng [one hot encoding]

	think	about	it	come	in	wav	15	45
1	1 0 0 0 0	0 1 0 0 0	0 0 1 0 0	0 0 0 1 0	0 0 0 0 1	<starts>	1 0 0 0 0 0	0 1 0 0 0 0
2	0 0 0 0 1	1 0 0 0 0	0 1 0 0 0	0 0 1 0 0	0 0 0 1 0	<ends>	0 1 0 0 0 0	0 0 1 0 0 0
3	0 0 0 0 0	0 0 1 0 0	1 0 0 0 0	0 1 0 0 0	0 0 0 1 0	2-5"	0 0 1 0 0 0	0 0 0 1 0 0
4	0 0 0 0 0	0 0 0 1 0	0 0 1 0 0	0 0 0 1 0	1	d	0 0 0 1 0 0	0 0 0 0 1 0

Decoder.

h_t	0 0 0 0 1 0 0
c_t	0 0 0 0 0 1 0
h_{t-1}	0 0 0 0 0 0 1

Row 1 pass [think about it $\rightarrow 1 \text{ or } 2$]



in decoder.

→ pass content vector and start symbol as input.

→ it generate some output, softmax layer. (three prob here)

→ suppose it generate $|j_1\rangle$

→ But we pass that output as input
but it is mixed.

so researcher think if we pass encoded
as input (Training become worst)

so it decide

whether it got encoded output.

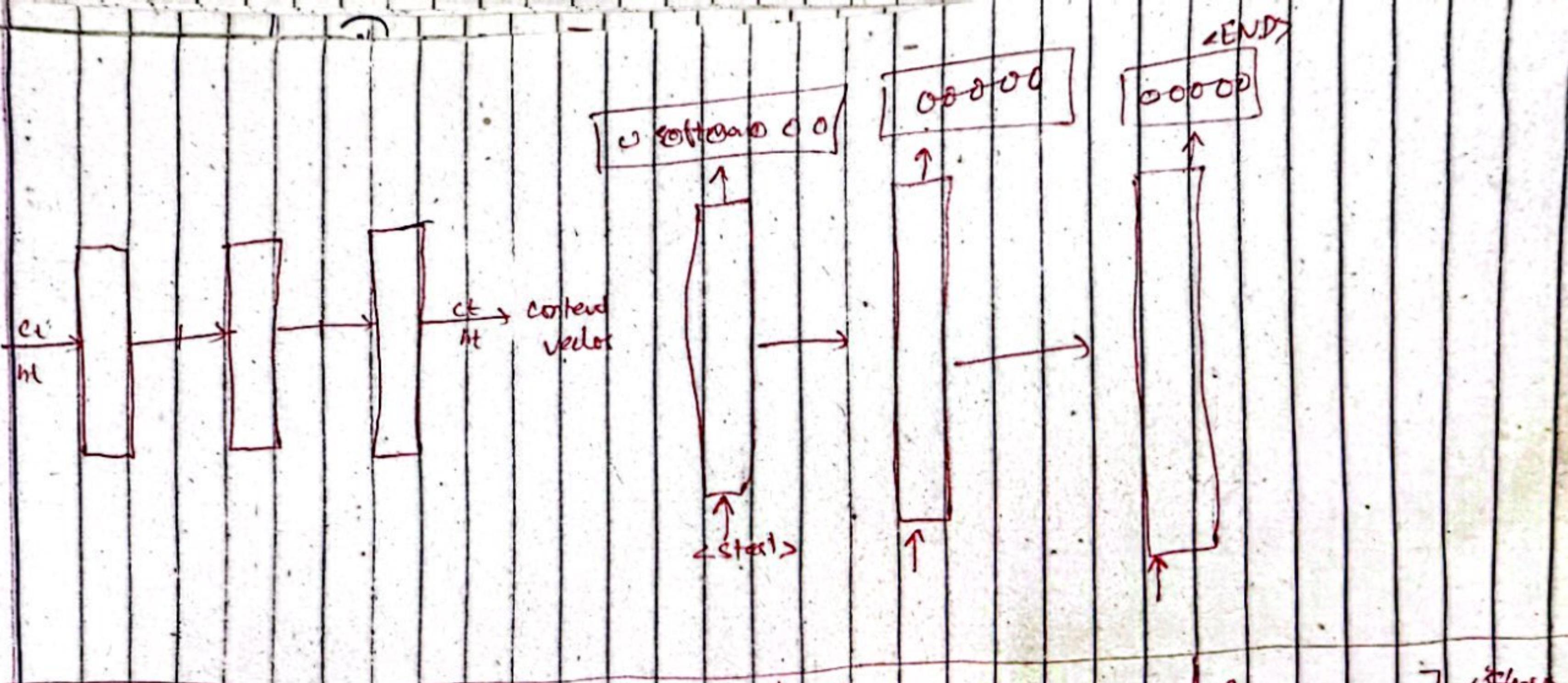
it give correct word as input

That's is called Teacher Forcing

its benefit is training fast
convergence

→ total loss

→ average loss



151

\hat{y} -predicted

193

Category loss euphony

$$L = - \sum_{i=1}^7 y_i \log(\hat{y}_i^{\text{pred}})$$

~~L_1~~ $L_1 = -1 \log(0.01)$

$$L_2 = 1$$

$$C_3 = -1 \cos(0.04) \\ = 0.99$$

real loss = 2.39

avg loss = 0.7

initial weight \rightarrow forward pass

output
compare
with
correct output

loss

Gradient

↓
weight

Weight update

using some
optimizer
ADAM

new
weight

↑

Learning
rate

batch
size
(speed)

1) forward

2) loss

3) Backpropagation

4) gradient

5) update

parameters

measure any one
particular parameter
how it contribute
in loss function
and that parameter
adjust it which
direction where
loss less.

→ Stochastic GD,
ADAM

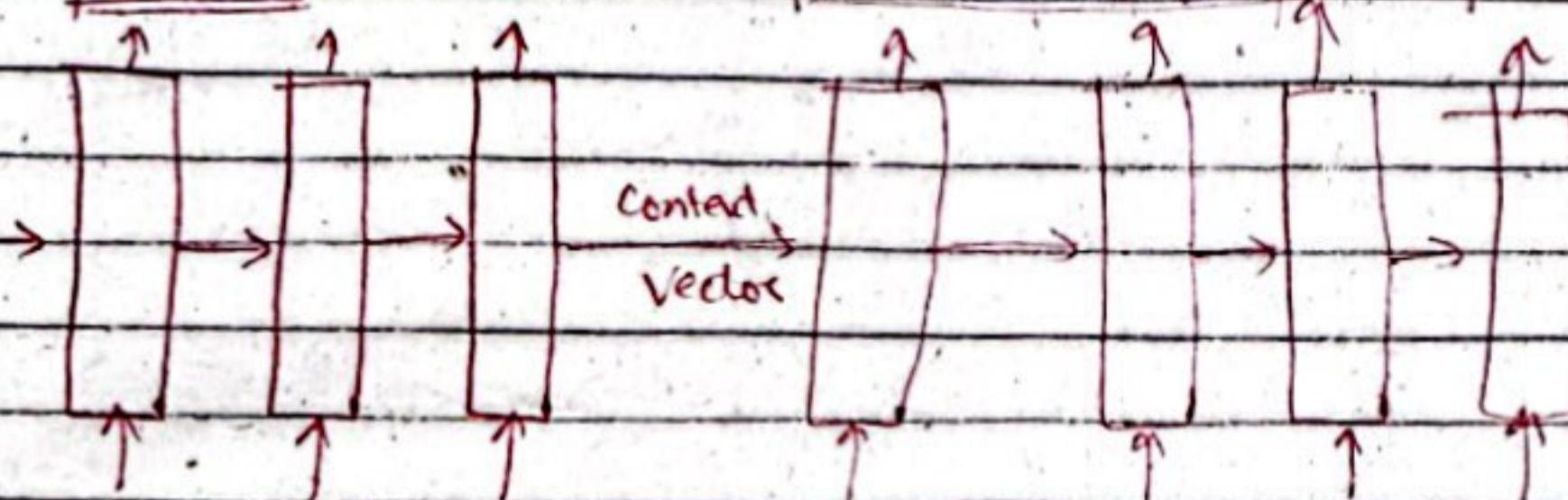
Optimizer adjust
weight in that direction
where loss less
and that adjustment
we based on
learning rate
hyperparameter
batch size

Prediction

① ② ③ ④ ⑤ ⑥ ⑦
 <start> 2.90 1.11 7.30 <END>

2.90 3.0 2 <END>

Prediction



think about it

<start> 2.90 3.0 2

[1000000] [0100000] [0000001] [0000100]

input \rightarrow [think about it]

2.90 3.0 2

the prediction

Improvement 1 [Embedding]

one-hot encoding



If 1 loc it means 1-loc data
it does not need



word embedding used

low dimension of embedding

dense representation

size: 30

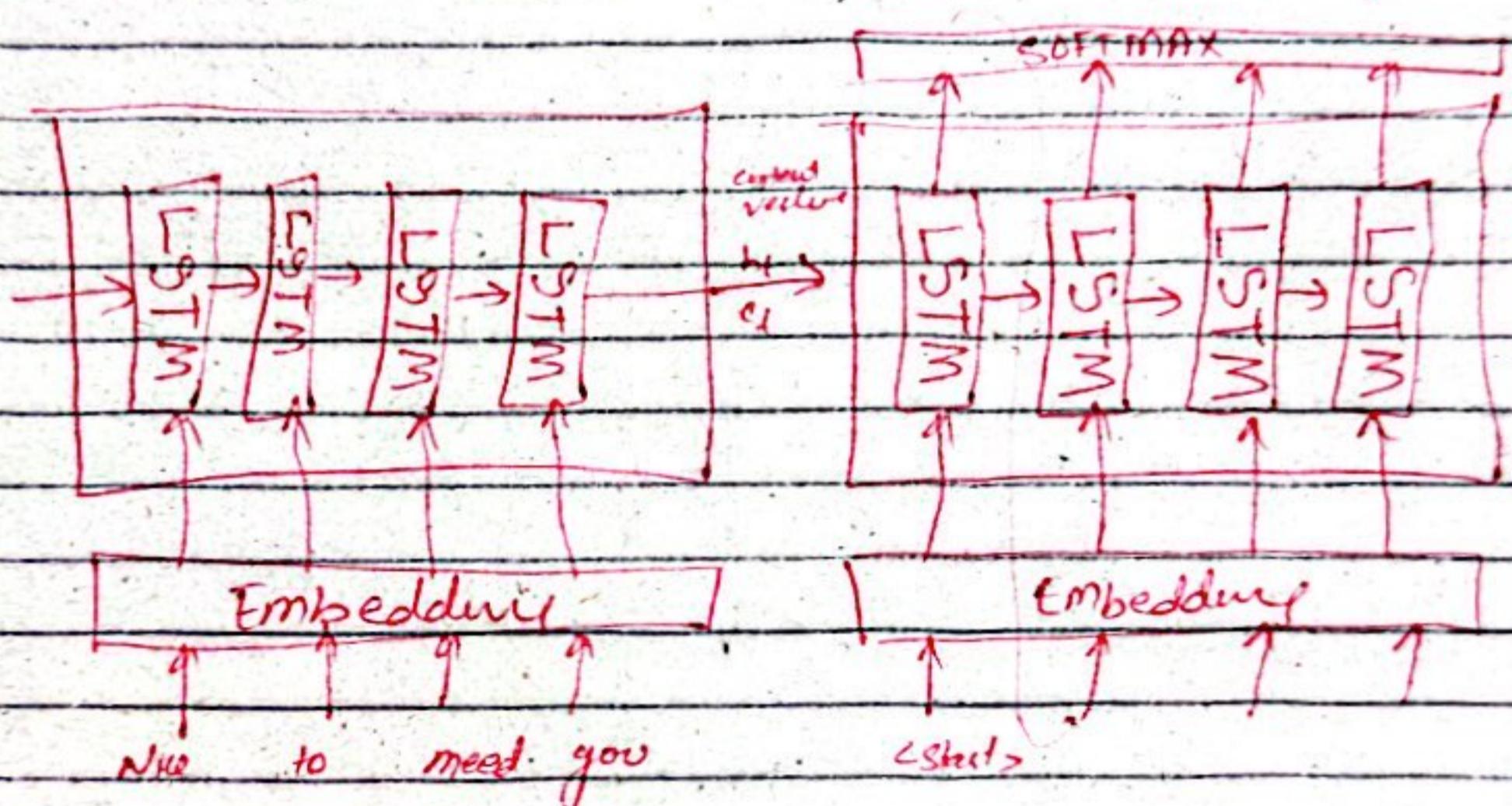
summary of word

that can't fit in
one-hot vector

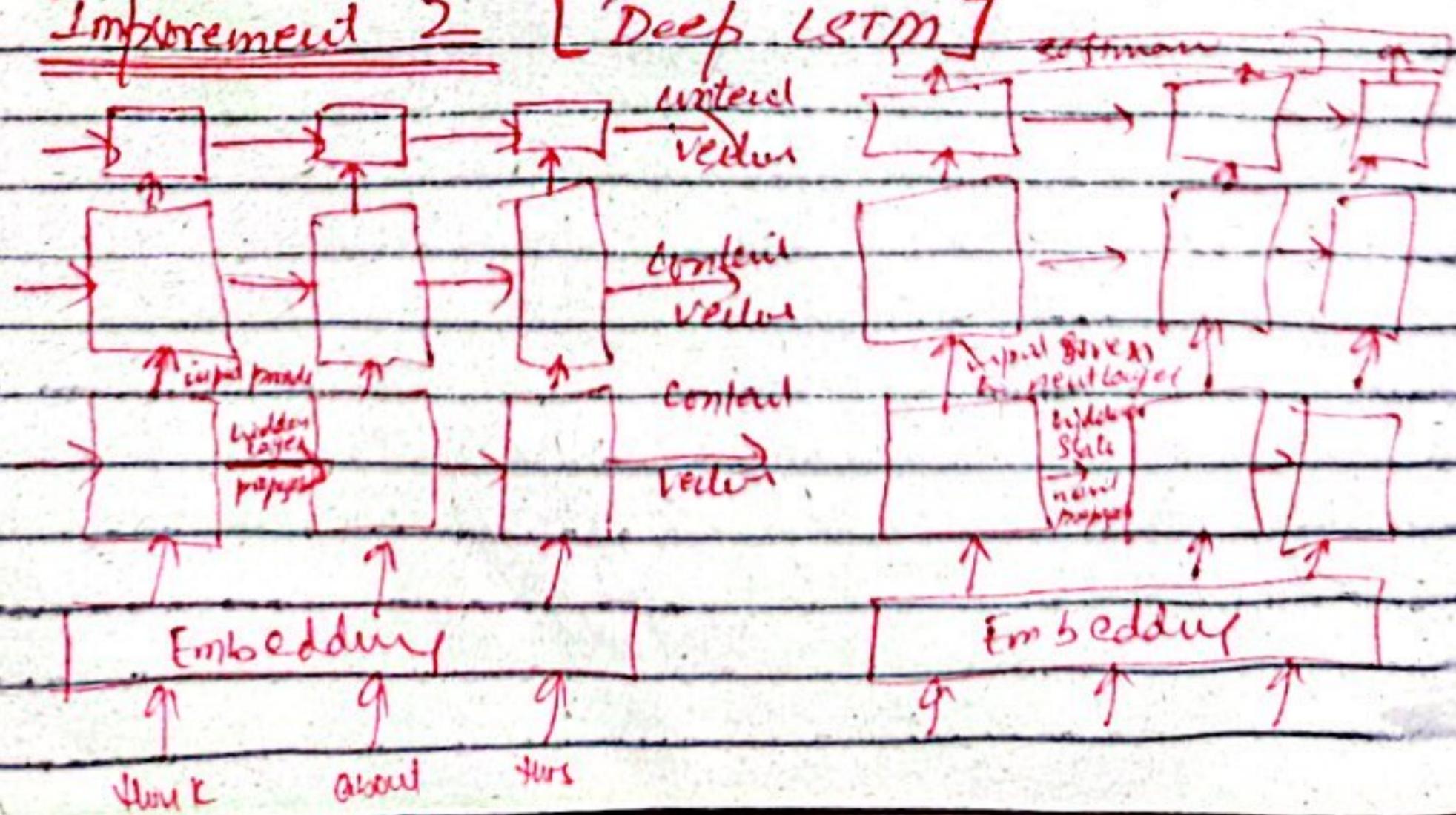
Embedding

pretrained
word2vec/Glove train
(own dataset)
captions

Encoder & Decoder both used Embedding



Improvement 2 [Deep LSTM]



3 Reason why you do this

1) long term dependencies
more efficiently handle
(long paragraph)

Single layer performance in long paragraph
is not good.

In deep Encoding and decoding

1) context vector get from multiple layers
we have more capability to store/
represent that summary

if one context vector is not efficiently
give summary of single sentence

2) layered representation :-

Hierarchical data more understand

→ lower level ^{LSTM} understand

(word level) →
→ singular
→ plural
→ which
→ pronoun

→ middle level LSTM

(sentence level) understand
↳ context understanding

→ top level LSTM

paragraph level understand
↳ full paragraph content

as The phone battery is bad But the overall
phone is good

③ Deep learning NN

Whatever you increase DNN parameters,
its learning capability increase
mean it capture minute minute
until not overfit

capacity increase \rightarrow small thing more
capture

\rightarrow better generalization

if sufficient data provide not
overfit and new data it better
give results.

not

\rightarrow input variable also better capture

Improvement 3

Reversing the effect

benefit

still words are more closer
and it pick good content
gradient backpropagate less effect.

Drawbacks

last words more far

It work on such language

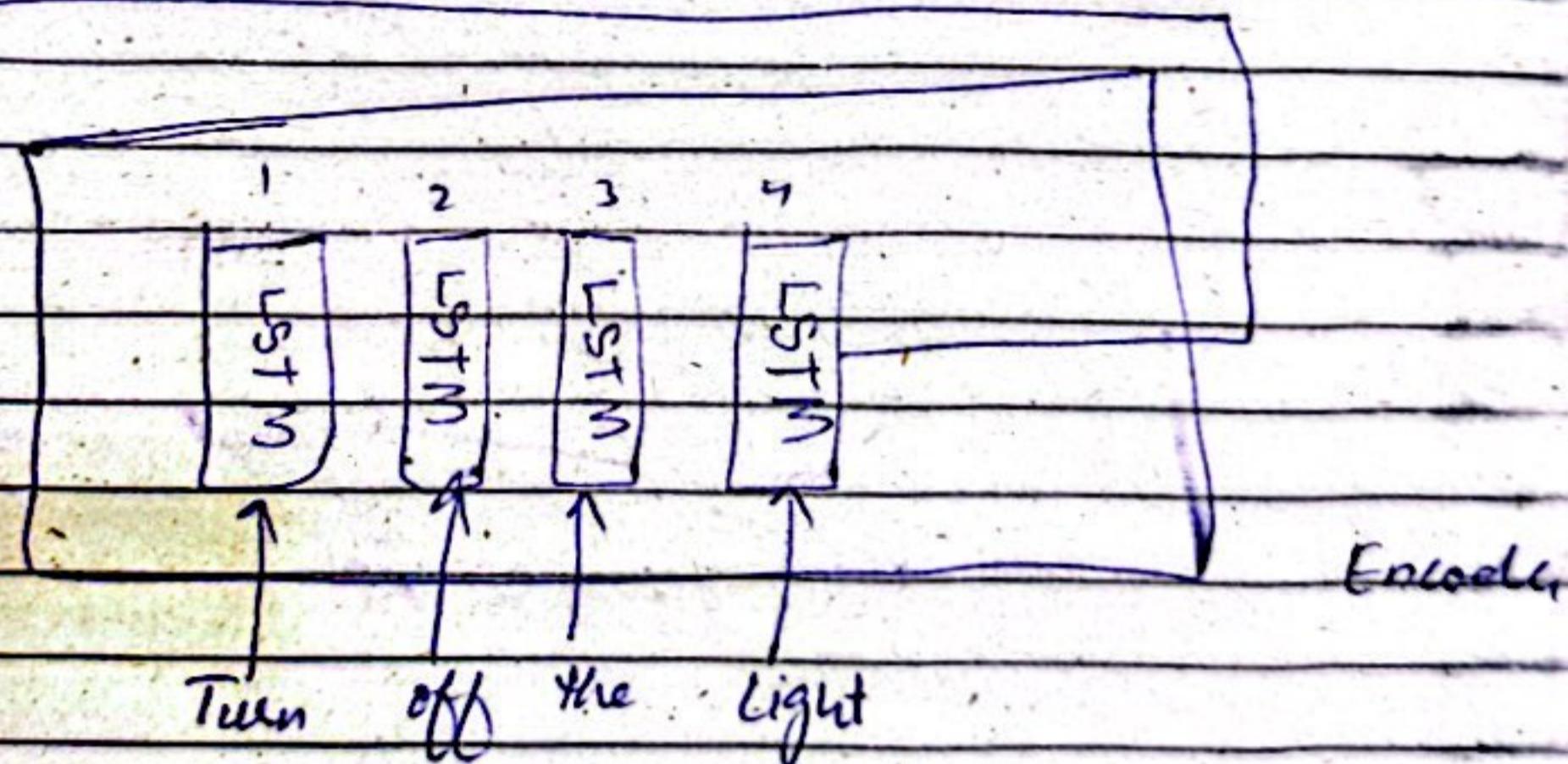
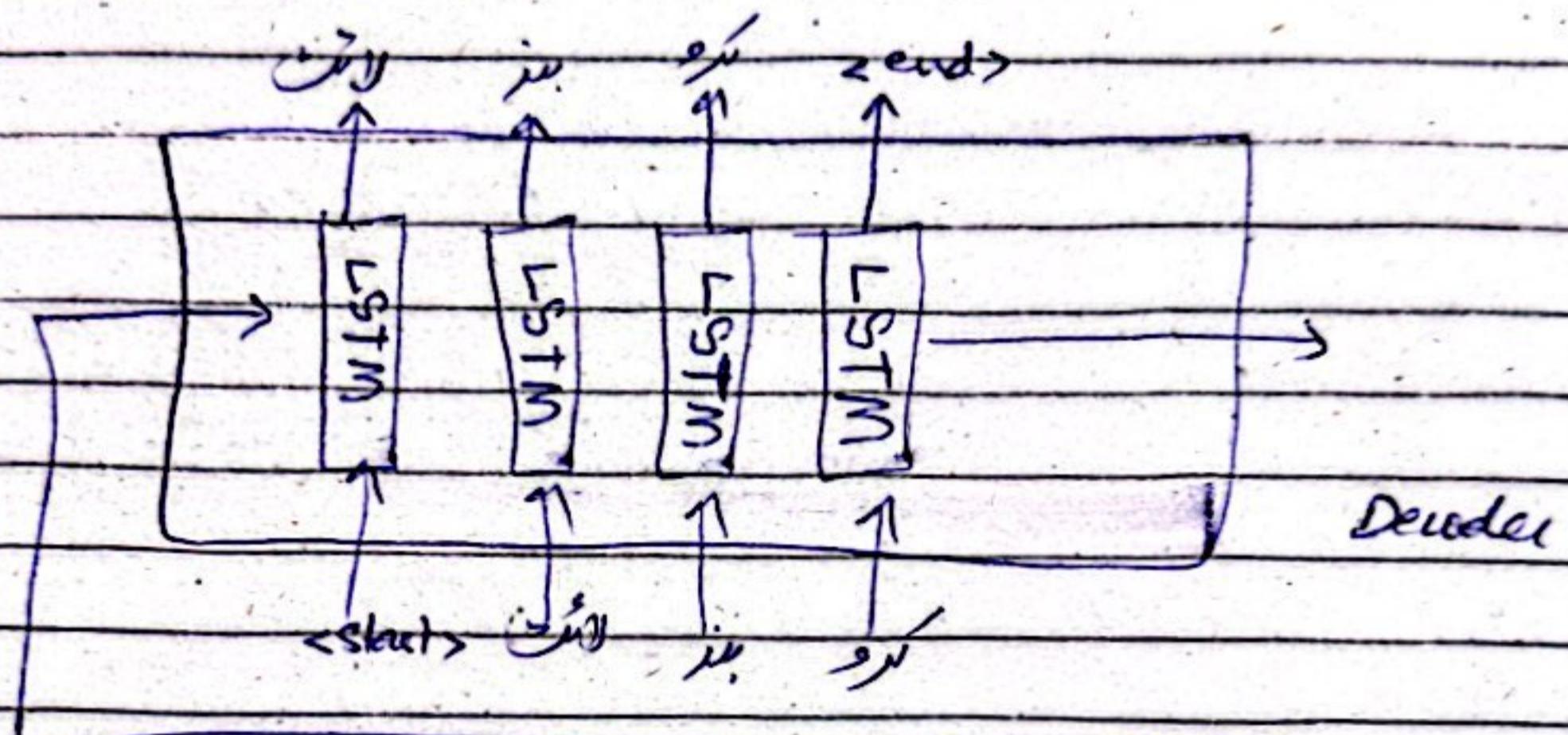
↳ initial content
more information (English → French)

but not for many languages.

Attention

Why we need attention?

Encoder Decoder Architecture



You provide one sentence to encoder. Encoder step by step process.

one time it process whole then it generate summary & context vector & pass that vector as input in decoder. decoder step by step output sentence.

problem of Encoder & Decoder :-

so word sentence, and one time summarize and then translate is very difficult.

problem is word > 25 context vector load is too very huge.

very small set of numbers you give very huge responsibility.

(\hookrightarrow) unnecessary pressure create on Encoder
→ at a time whole sentence load and learn and tell it decoder

(\hookrightarrow) at any time step at decoder, we don't need input sentence to translate we need some specific word / set of words

but problem in this architecture. whole input sentence give at very time stamp go and figure out.

which word you need or your decision

→ because of static representation, decoder face problem to decode.

best if dynamic if it is, any particular output point only specific part attention

Solution

human being if translate long sentence it only focus particular part of sentence and translate it. other part is blur.

This concept we introduce in our architecture.

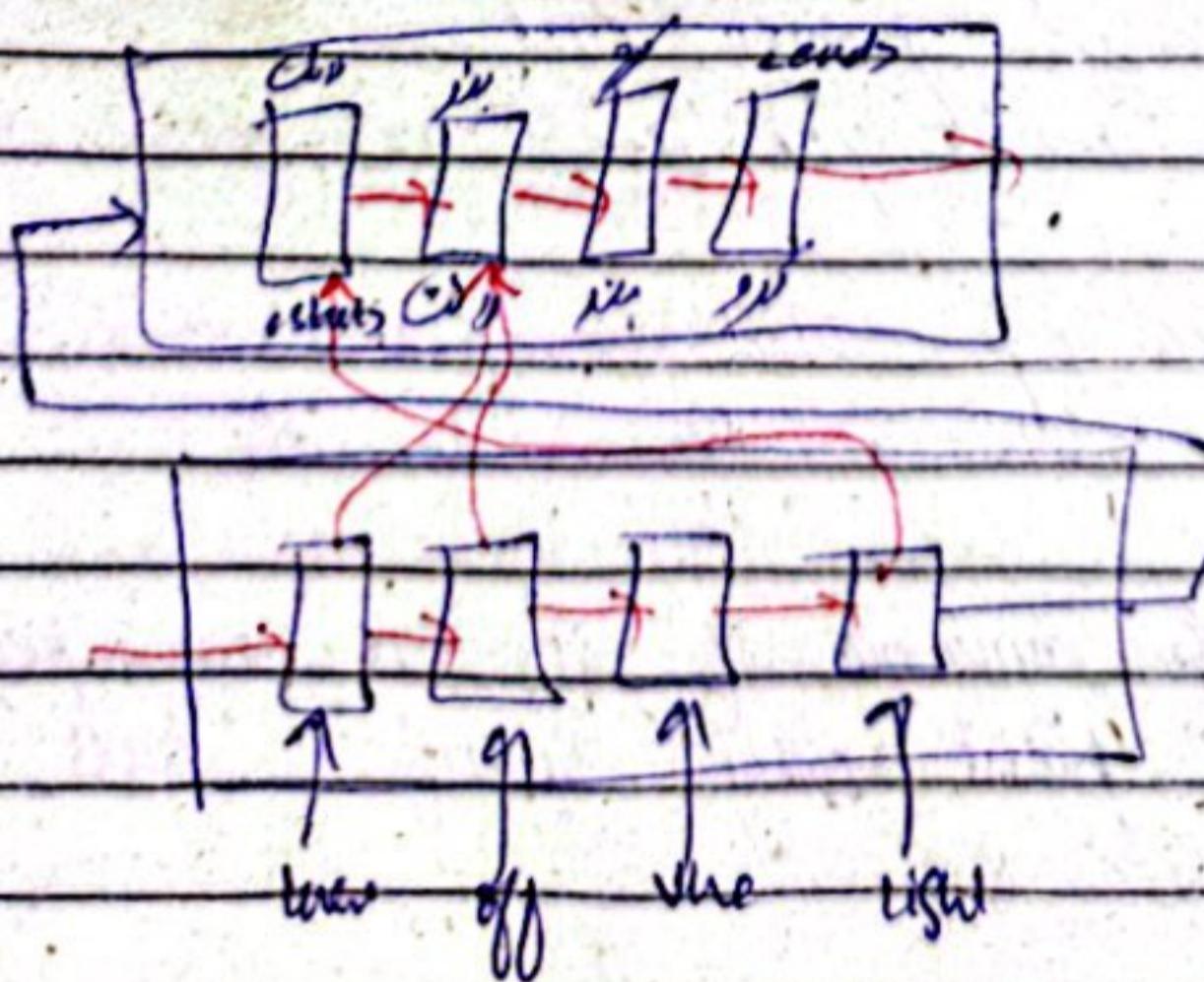
In decode

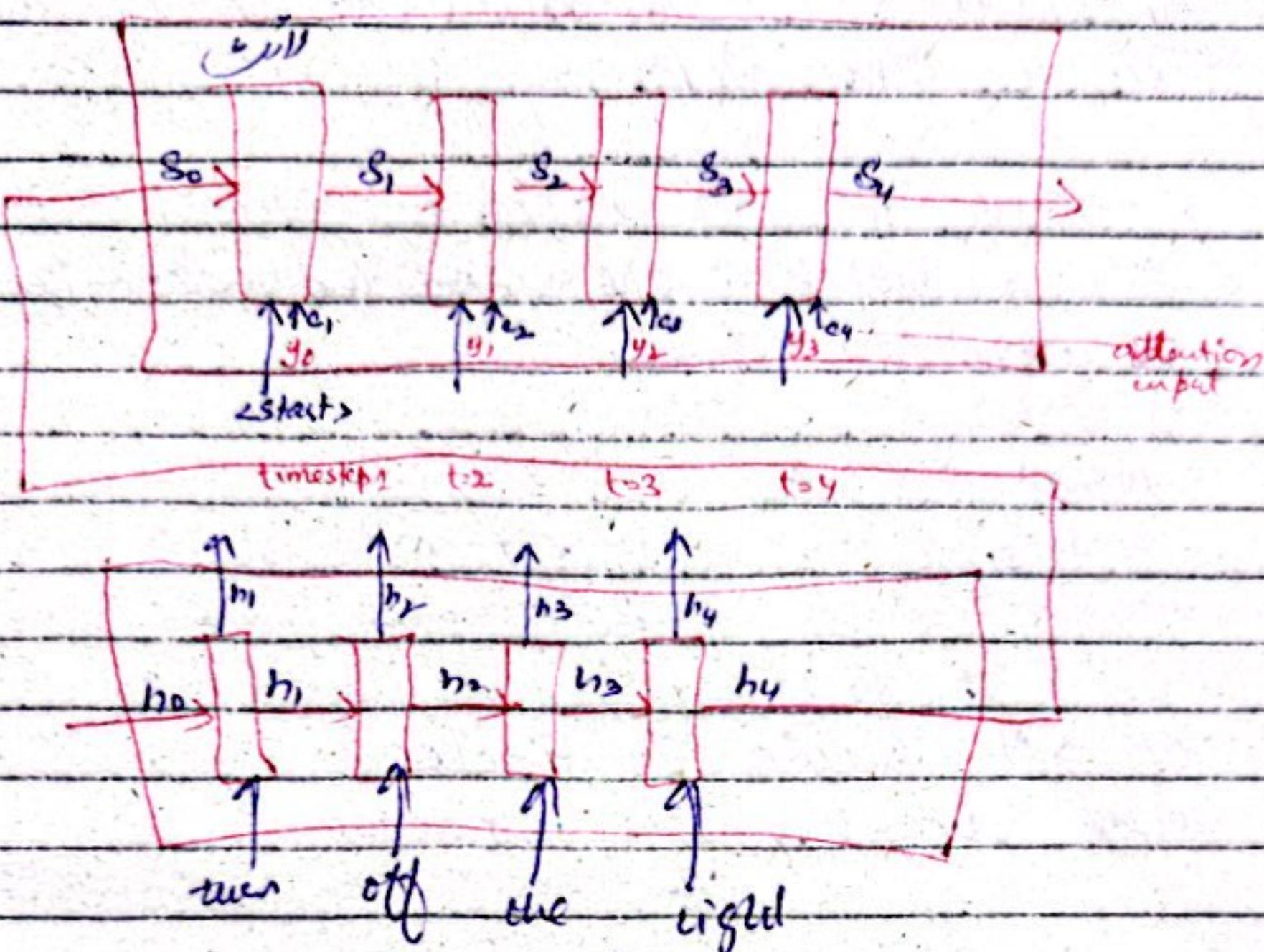
at timestamp t to print decoder. we tell information in any way of decoder that encoder timestep is very important to print. $t-1$ to print in encoder timestep 1 and 2 are important and so on.

We dynamically generate information at every time stamp of every decoder, which encode timestamp or set of time steps important currently.

and that mechanism is called

Attention mechanism





we know that, in attention mechanism
we pass out information which encodes time
step is important for decode each
time step. teacher forcing
output hidden

Decoder $t = 2$ (y_t, s_{t+1}) provide input (Vanilla)

In attention you also pass

$t = 2$

y_{t-1}, s_{t-1}, c_t

its purpose to
tell decoder
at current timestep
which encoded
hidden state
is useful
attention
input

what c_t is look like?

- 1 - vector } to answer what question
- 2 - scalar } first code what is
- 3 - matrix } the purpose of c_t ?

~~encoder~~ ^{encoder} first time step which hidden state of h_1 , hidden state is useful to point c_i at decoder first step.

If someone tell me megirally which hidden state is important for decoder each time step c_i work is only h_1, c_i

Decoder

$$t=2$$

$$h_1 \text{ and } h_2$$

$$t=1$$

$$h_1$$

c_i^* \Rightarrow it should be vector

which dimension of c_i .

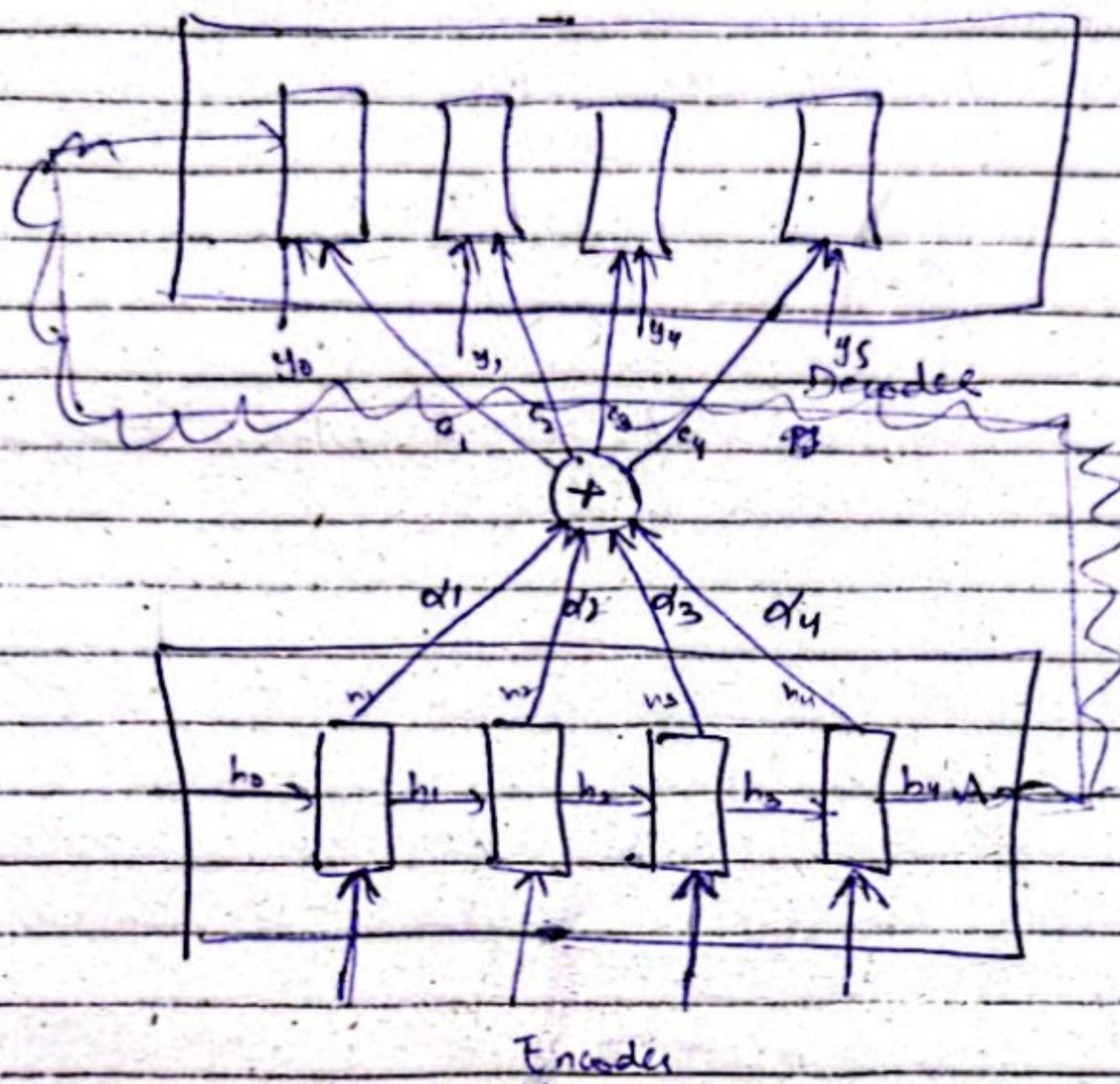
$$\begin{bmatrix} 1 & 0.5 & 0.6 & 0.3 \end{bmatrix} \Rightarrow h_1$$

$$\begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \Rightarrow h_2$$

$$\begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \Rightarrow h_3$$

c_i exactly same as h_i .

but what if, I will send two hidden as we weighted sum. and send it.



Decoder time step "i"
 Encoder time step "j"

In attention mechanism:-

It assign weight of all hidden state if
 said

$$c_{i,j} = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3 + \alpha_4 h_4$$

vector ^{hidden}

c_i = 4 dimension

If we find c_i , we find weight for each
 hidden state, which α_i is dominate its
 impact is highest.

for each time step ... for decoder
 α_i is different for each time step

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

$$\begin{aligned} C_1 &= \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4 \\ C_2 &= \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4 + \alpha_{25} h_5 \end{aligned}$$

$$\text{total } \alpha = i + j = 4 \times 4 = 16$$

How "α" are created? how it calculates
because it has information which hidden state is important.

$\alpha_{21} \rightarrow$ alignment / similarity score

which tell decoder 2nd time

decoder 2nd time output which output is present, it
decoder 1st output encoder first time step
how much hole.

α_{21} which quantities it depend?

$h_1 s_1$ ↓
previous hidden
state of decoder
how?

given translation based
second step output
how much hole?

$\alpha_{21} \rightarrow f(h_1, s)$

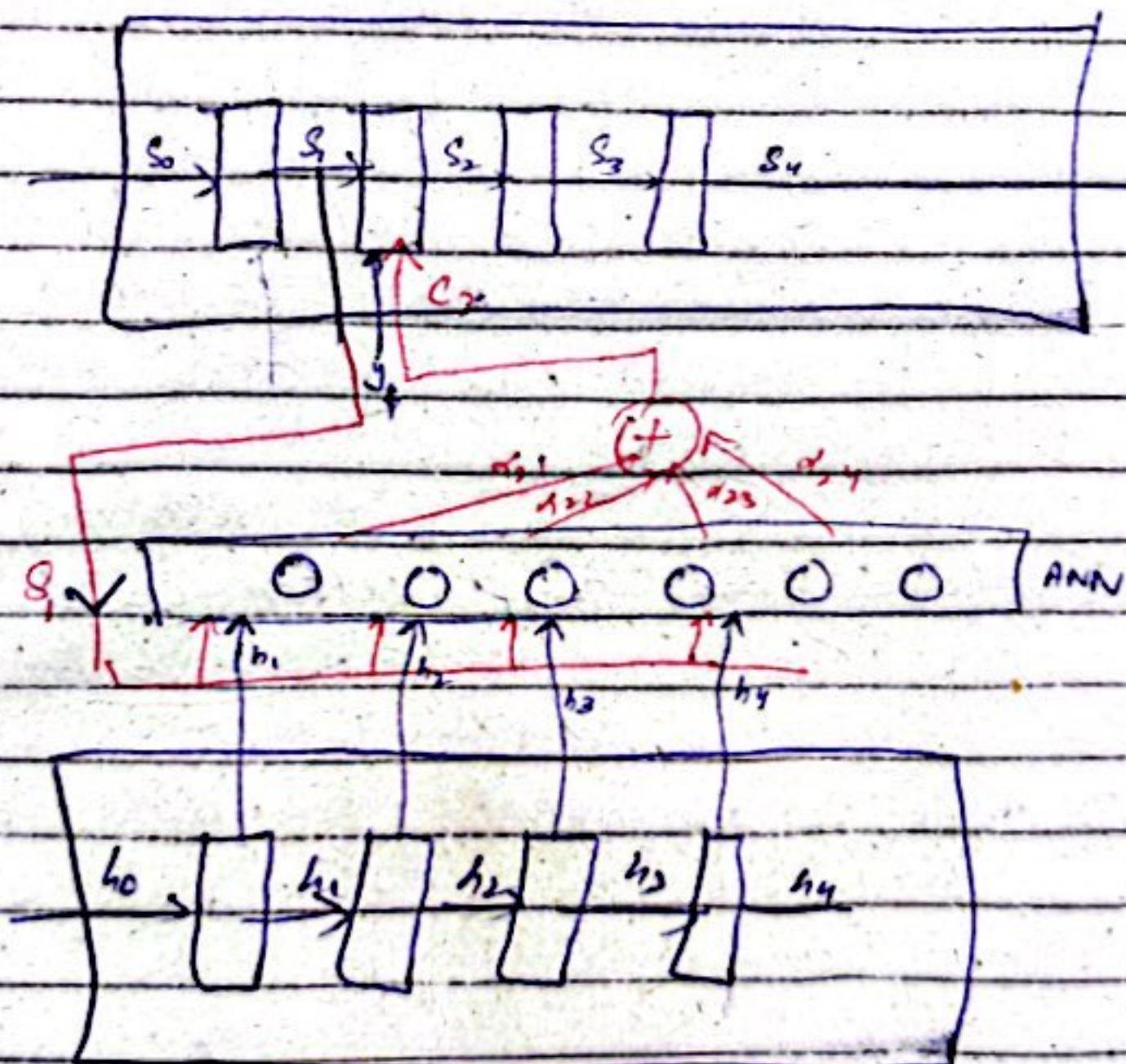
$$d_{ij} = f(h_j, s_{i-1})$$

$$d_{23} = f(h_3, s_1)$$

what function is used?

researches - think why it have functions.

ANN \rightarrow universal function Approximator.
we can directly used neural n/w.



$$c_2 = d_{21} h_1 + d_{22} h_2 + d_{23} h_3 + d_{24} h_4$$

BLEU Score \Rightarrow translate Quality

Sentence length > 30 BLEU score

maintain/stable in attention mechanism

and other BLEU score > 30 decrease

attention lead = $\alpha \Rightarrow$ plot.

man off the light

in	d_{11}	d_{12}	d_{13}	d_{14}
is	d_{21}	d_{22}	d_{23}	d_{24}
it's	d_{31}	d_{32}	d_{33}	d_{34}
end	d_{41}	d_{42}	d_{43}	d_{44}

english to french translation. (paper)
its tell benefit of attacking attention.

} encoder \rightarrow bi-direction LSTM

birectional and forward
content

decoder \rightarrow LSTM

Bahdanau Attention vs Luong Attention

Attention

Machine translation :- NMT (Neural Machine Translation)
turn off the light \rightarrow 灯关了

Encoder & Decoder

→ LSTM, bi-LSTM, stacked LSTM used as $\text{hidden} \rightarrow \text{out}$

Step by step each word sentence, encoder
goal it to generate summary context
vector) and that context vector passed
to decoder.

Decoder work is step by step
try to print output.

Problem :-

For > 30 paragraph, document
encoder face difficulty to generate
content vector (set of numbers how
much data you create store. Content
vector can't create efficiently and
that's why decoder can't generate
good quality translation.

To eliminate that problem new mechanism
came what is **Attention**.

* particular

Decide any timestep, you print a word,
you don't need full sentence, you
need ~~some~~ words or subset of words.

light
turn off.

c_i

i

Attention mechanism said, don't need
of send whole content vector.

we dynamically select h_i.

at each time step of decoder, content
vector send. c₁, c₂, c₃, c₄

i) point {h₁} want
ii) point {h₁, h₂}

* if we want any particular c_i = $\sum \alpha_j h_j$

$$c_1 = \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4$$

$$c_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$$

Total alphas } input sentence word x
alignment store } output sentence word

how we find α (word to word similarity values score)

- Q -
- 1 - Bahadonov additive attention
 - 2 - Luong attention

Bahadonov attention (Additive Attention)

α = (alignment score) how much importance of decoder time step encoder which time step how much important.

$$\alpha_{ij} = f(h_j, \underbrace{s_{i-1}}_{\substack{\text{previous decoder} \\ \text{time step}}})$$

α tell given translation current decoder output how much it depend on encoder which time step.

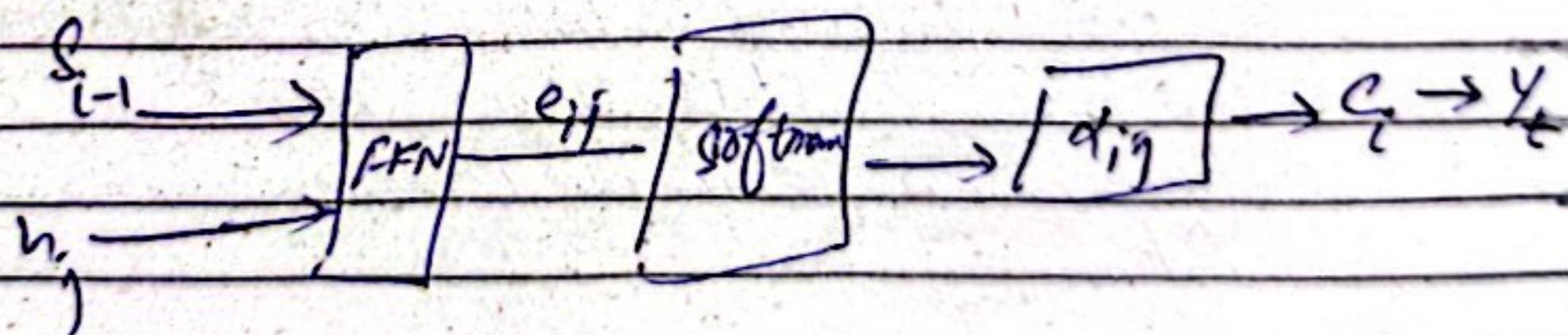
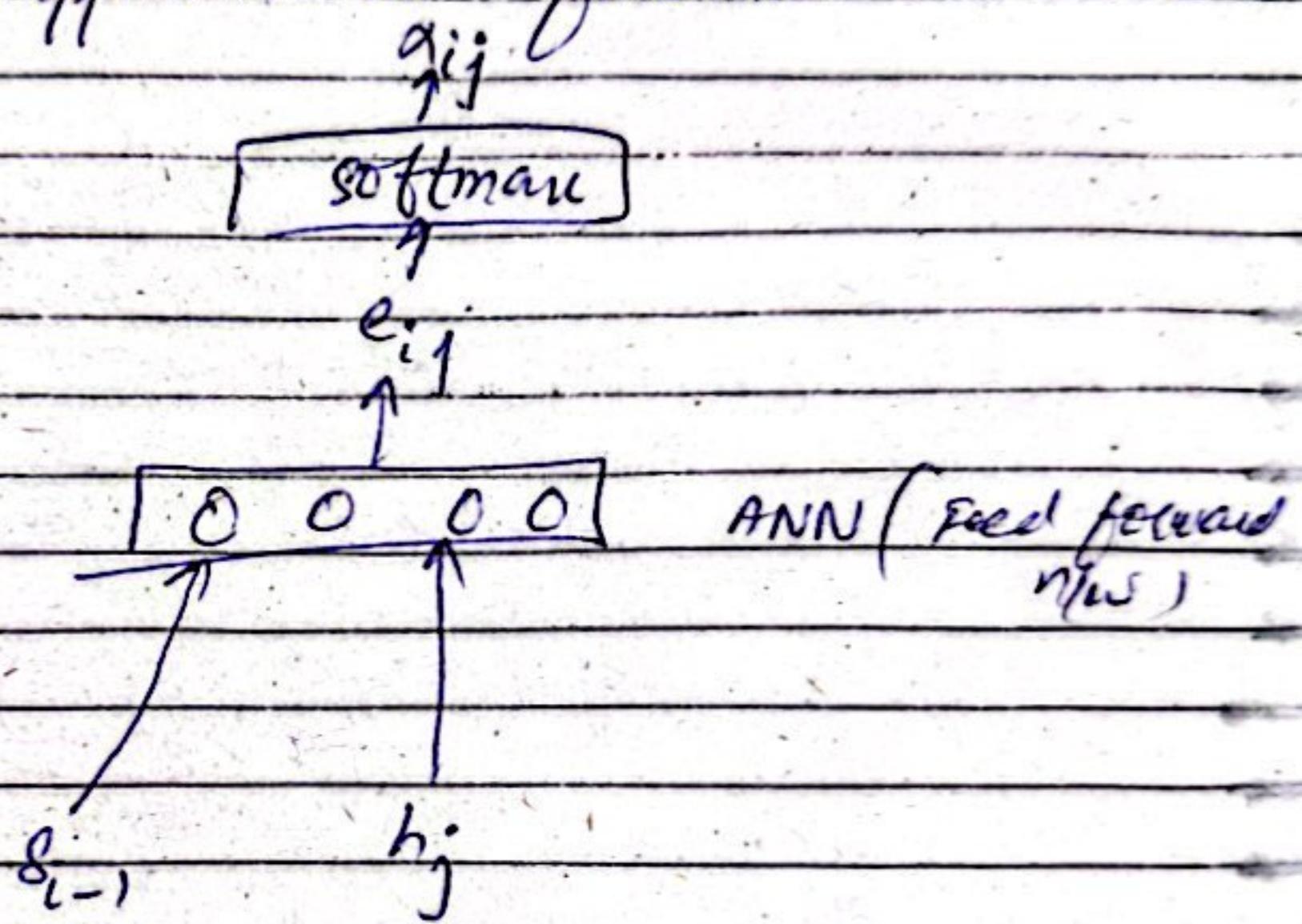
Alignment score will be mathematic function of two things

- 1- encoder hidden state
- 2- decoder previous ^{sup}hidden state

Q which mathematic we used
A Bahadonov Attention said, why we word hard to find that function.

instead we approximate that function
with the help of ANN.

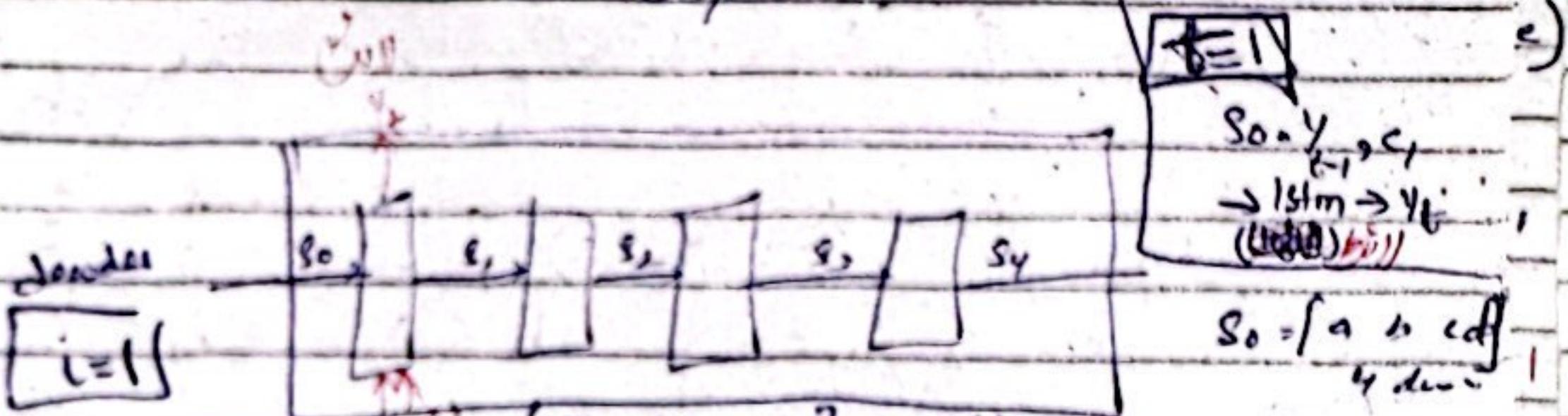
ANN \Rightarrow approximation function.



Architecture

Encoder is same, but in decoder
you add FFNN (Feed forward neural net)

hidden layer $\{0, 0, 0\}$ three neurons
+ output



$$f(1 \times 4) \quad [0, 0, 0, 0] \rightarrow \text{softmax}$$

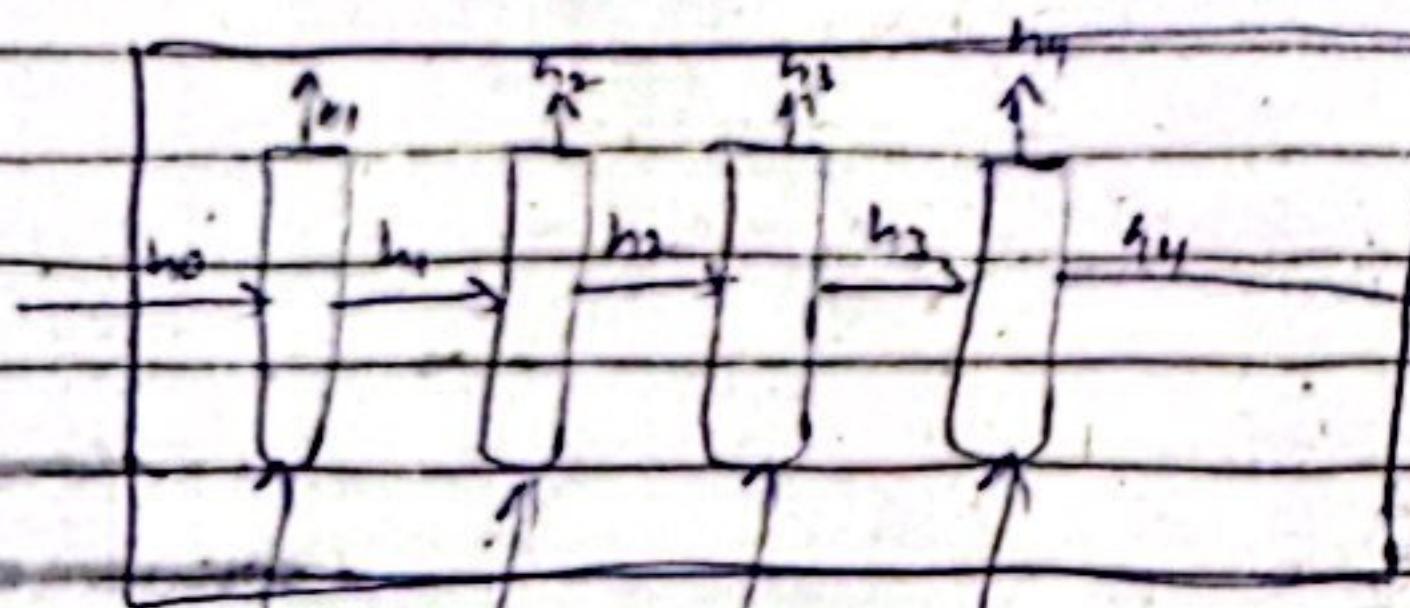
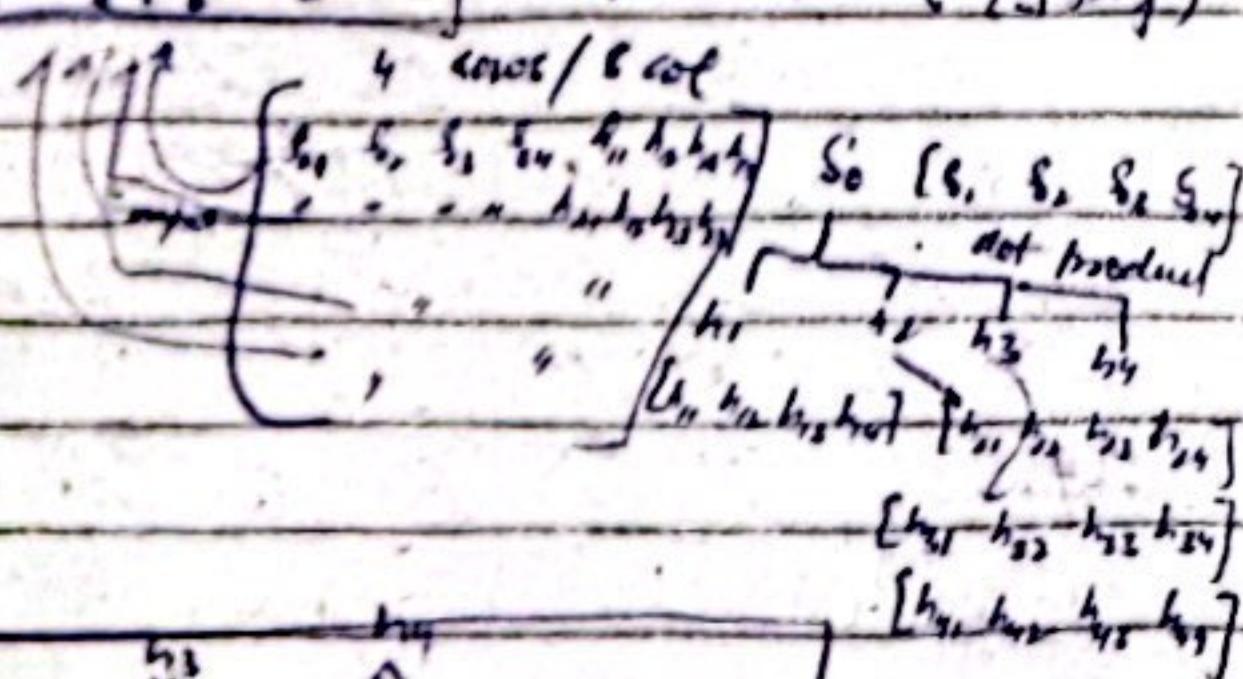
$$(3 \times 1) \quad (3 \times 1) \quad \xrightarrow{\text{weight}} [0] \rightarrow \text{output}$$

$$4 \times 3 \quad 3 \times 3$$

$$(2 \times 3) \quad (3 \times 2) \quad \xrightarrow{\text{softmax}} [0, 0, 0] \rightarrow \text{output}$$

4x8 RBGs $[0, 0, 0, 0, 0, 0, 0, 0]$ input $a(h_1, h_2)$

\rightarrow you send one row
or batch of 4



turn off the light

$$h_t = f(a + c_0)$$

dim

↳ encoder Sentence send

~~the , h₁, h₂, h₃, h₄~~ came

↳ Decoder t=1

$$S_0 [s_{01} s_{02} s_{03} s_{04}]$$

$$\begin{matrix} h_1 & h_2 & h_3 & h_4 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ [h_{11} h_{12} h_{13} h_{14}] & [h_{21} h_{22} h_{23} h_{24}] \end{matrix}$$

$$\begin{bmatrix} s_{01} s_{02} s_{03} s_{04} h_{11} h_{12} h_{13} h_{14} \\ s_{01} s_{02} s_{03} s_{04} h_{21} h_{22} h_{23} h_{24} \\ s_{01} s_{02} s_{03} s_{04} h_{31} h_{32} h_{33} h_{34} \\ s_{01} s_{02} s_{03} s_{04} h_{41} h_{42} h_{43} h_{44} \end{bmatrix}$$

$$[h_{31} h_{32} h_{33} h_{34}]$$

$$[h_{41} h_{42} h_{43} h_{44}]$$

batch
of 4

$$\begin{array}{cccc} a_{11} & a_{12} & a_{13} & a_{14} \\ \hline \text{softmax} & & & \end{array}$$

$$[e_{11} e_{12} e_{13} e_{14}]$$

4x1

$$\boxed{0}$$

4x1

N

[3x17]

[3x1]

$$\boxed{0 \ 0 \ 0 \ 1}$$

2x3

4x3

$$\boxed{1 \ 1 \ 1}$$

[2x3]

4x8

$$\boxed{0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ -1}$$

1x8

feature
of one

feature
of one

$$C_1 = a_{11} h_1 + a_{12} h_2 + a_{13} h_3 + a_{14} h_4$$

Decoder $t=2$

$$S_t = [s_{11} \ s_{12} \ s_{13} \ s_{14}]$$

$$s_2 = g_1 h_1 + g_2 h_2 + g_3 h_3 + g_4 h_4$$

e)

$$\begin{bmatrix} s_{11} & s_{12} & s_{13} & s_{14} & h_{11} & h_{12} & h_{13} & h_{14} \\ s_{11} & s_{12} & s_{13} & s_{14} & h_{21} & h_{22} & h_{23} & h_{24} \\ s_{11} & s_{12} & s_{13} & s_{14} & h_{31} & h_{32} & h_{33} & h_{34} \\ s_{11} & s_{12} & s_{13} & s_{14} & h_{41} & h_{42} & h_{43} & h_{44} \end{bmatrix}$$

$$s_1 \ s_2 \ y_i \rightarrow \text{LSTM} \rightarrow y_i \rightarrow p/s_i$$

weight same ~~for~~ until all sentence
print

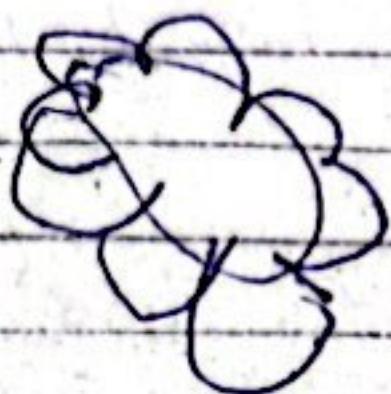
then ~~new~~ back propagate
then ~~sentence~~ send.

⇒ decoder hidden state
change only
~~because that's why~~
I change each
time step.)

time distributed ~~feed~~ feed connected
neural N/W.

parameter update
when ~~back~~
propagate

A B C D E F G H I J K L O P Q



$$c_i = \sum_j a_{ij} h_j$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^n \exp(e_{ik})}$$

Softmax

$$e_{ij} = V \cdot \tanh(W [s_{i-1} \cdot h_j] + b)$$

Bahdanau Attention

also called additive attention

Neural N/W \Rightarrow alignment node.

Luong Attention

(Multiplicative Attention)

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

Luong Attention

$$\alpha_{ij} = f(s_i, h_j)$$

more
dynamically
output
adjust

current
decoder
hidden
state

not used Complex NLP

Simple $s_i \cdot h_j$ (dot product)
take.

Bahdanau Attention

$$\alpha_{ij} = f(s_{i-1}, h_j)$$

Previous
decoder
hidden
state

FFNN

unnecessary

parameter add/
calculate

$$s_i = [a \ b \ c \ d]$$

$$h_j = [c \ d \ e \ f]$$

two vector
similar
dot product

two vector
dissimilar
dot product

$$s_i = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

$$\cdot [c \ d \ e \ f] = [ae \ bd \ ce \ df]$$

softmax

$$\alpha_{ij} = e^{e_{ij}}$$

softmax

α_{ij}

less no
of param

these basic funda is not

how good you find approximation

function, you have main

funda is encoder which

hidden state is useful

that also do dot product

→ dot product fast

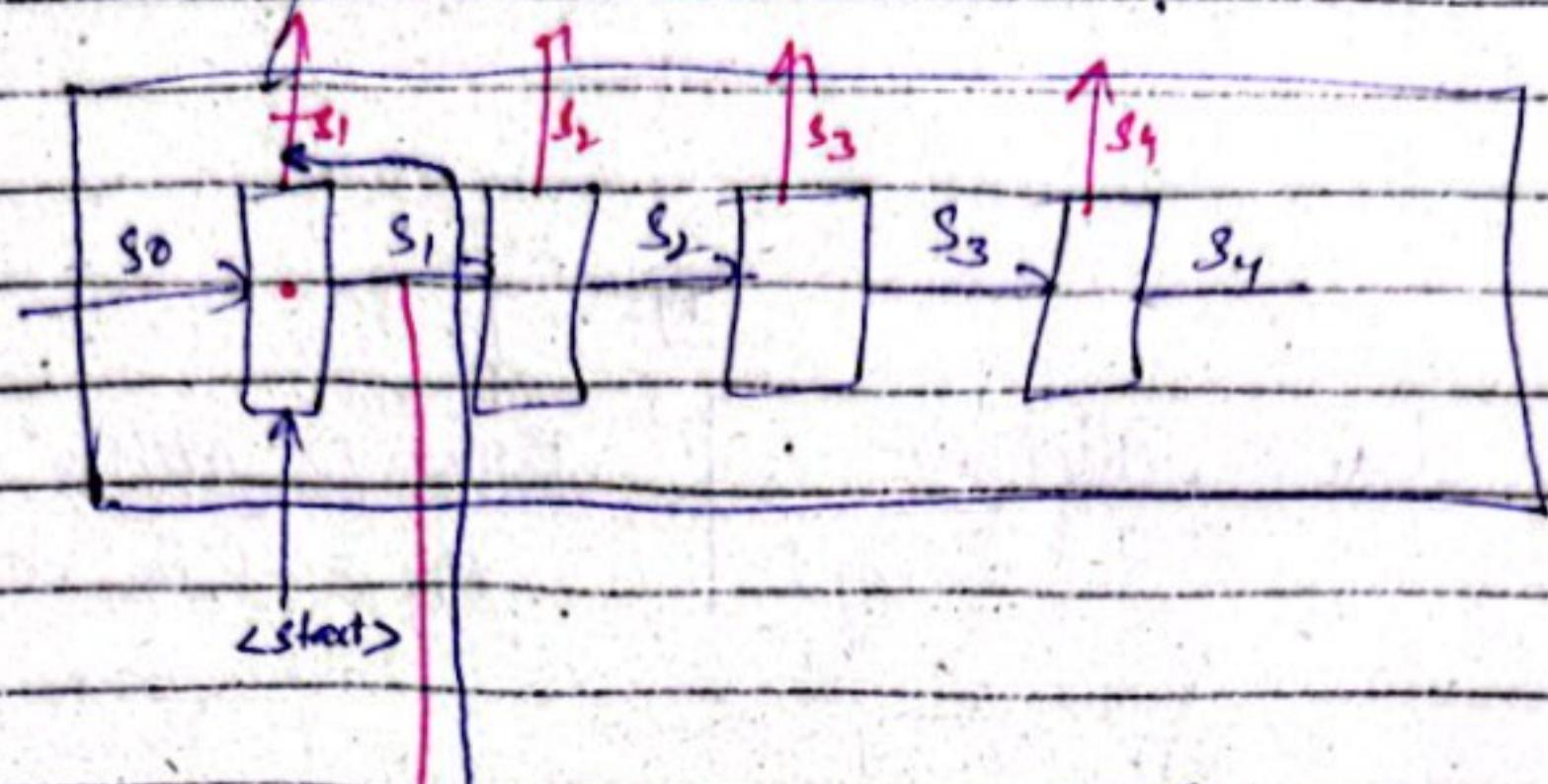
→ recent decoder hidden state performance

updated information

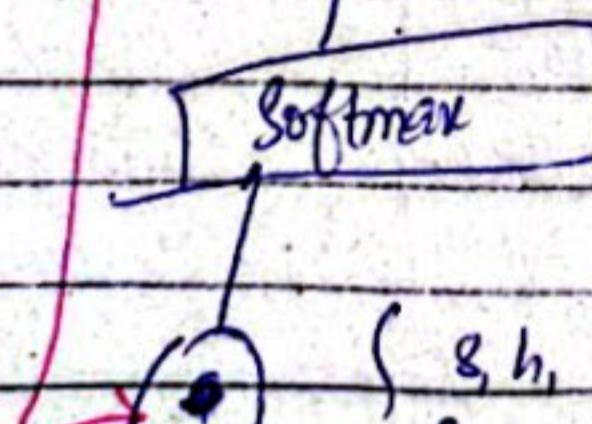
→ recent decoded
→ dot product

$$S_i \cdot C_i = \tilde{S}_i \rightarrow \text{softmax} \rightarrow \text{softmax}$$

numerical n/w



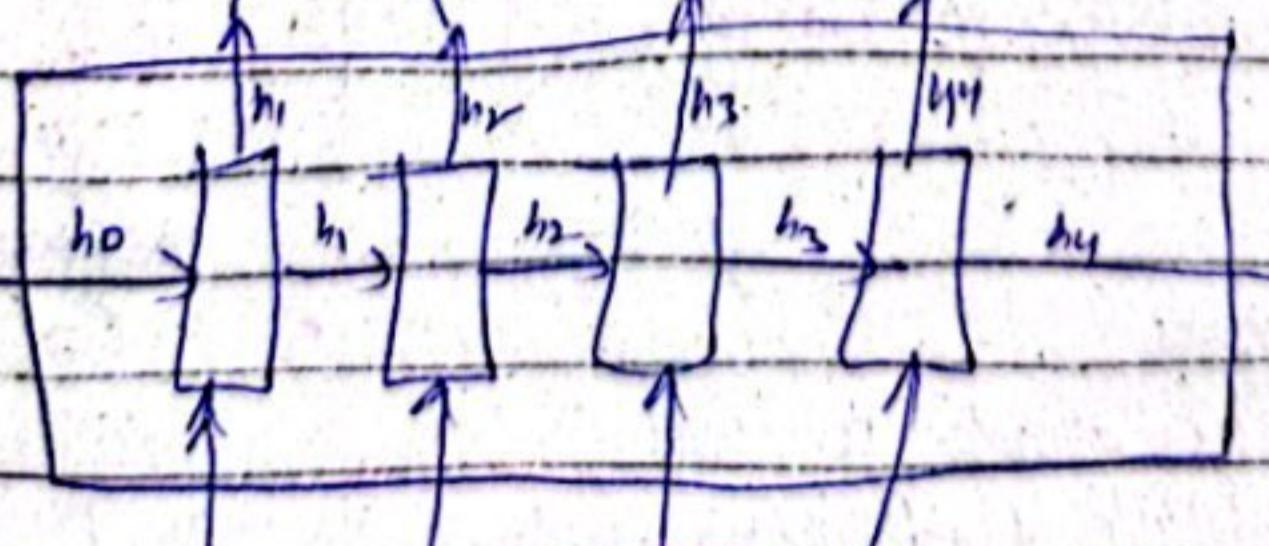
$$\alpha_{ij}, \alpha_{i1}, \alpha_{i2}, \alpha_{i3}, \alpha_{i4} \Rightarrow c_i = \sum \alpha_{ij} h_j$$



$$\{g, h_1, e, h_2, g, h_3, g, h_4\}$$

$$\{e_1, e_2, e_3, e_4\}$$

$$\{\alpha_{i1}, \alpha_{i2}, \alpha_{i3}, \alpha_{i4}\} \rightarrow c_i$$



turn off the light

→ performance better

→ fastly

Self Attention

NLP → Application

→ sentiment

→ machine translation

→ Named Entity Recognition

important requirement

how words converted into number (vectorize)
efficiently

word → number (vectorization)

1- One hot Encoding

→ mat cat mat → [1 0 0] [0 1 0] [1 0 0]

→ cat rat eat → [0 1 0] [0 0 1] [0 0 1]

↳ unique word find

mat

cat.

rat

	mat	cat	eat
mat	1	0	0
cat	0	1	0
rat	0	0	1

what way is efficient

2) Bag Of words

→ how times a word occurs

→ word occur or not

)

unique word

	mat	sat	cat
Sent 1	1	2	0
Sent 2	0	1	7

3) TF-IDF

4) Word Embedding

↳ word better way to represent numbers.

↳ because it has capability of semantic meaning capture.

↳ take large learning data

↳ neural net understand and which word is used in which context and every word represented into n-dm vector

doggy → [0.6 0.2 0.1 0.7]

king → [0.3 0.2 0.4 0.07]

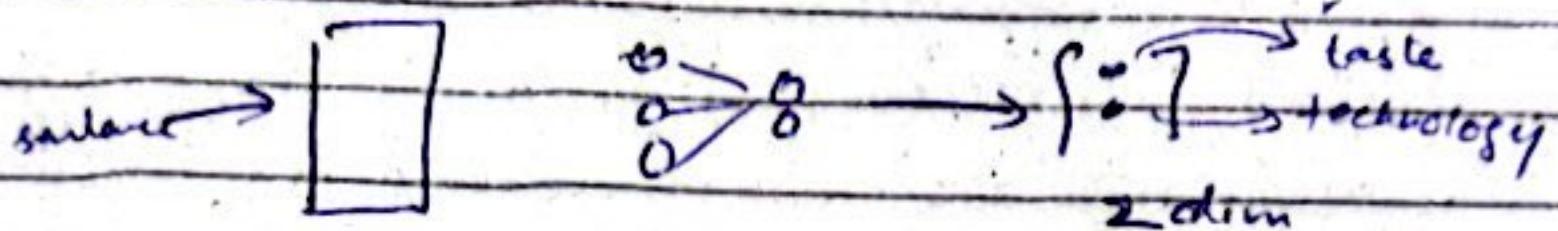
two word similar → vector similar
each dim. represent some aspect (feature)
about data we can't see

each dim. ^{some} meaning and each dim. similarity see

word meaning
somehow 'changes' into vector

The problem in embedding "average embedding"

- An apple a day keeps the doctor away
- Apple is healthy
- Apple is better than orange
- Apple makes great phone



our word embedding does not capture meaning, it captures average meaning

Apple	Taste	Technology	on an average that particular word in your data has particular meaning it used
1)	[0.6 0]		
2)	[0.7 0]		
3)	[0.8 0]		
4)	[0.8 0.2]		
Total	10,000 soul		

- Given sentence apple as a fruit came
- 1000 " " " " a company

taste Technology
[0.9 0.3]

tech → taste is more tilted toward taste
in comparison to technology
taste

word embedding how looks, it depend on data.

word embedding basically any particular word "average meaning" represent in vector.

word embedding will one time and used so many times.
it is static.

apple [0.9 0.3]

used so many
times

It is prediction

e.g application

translation

Eng → Ukr

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

first you convert each word its respective word embedding

[0.9 0.3] [] [] []

↳ but its accept where apple used

as a fruit

but in this other sentence, it

is used as a company

ideally instead of static embedding
used Contextual embedding

that value dynamically
change based on content.

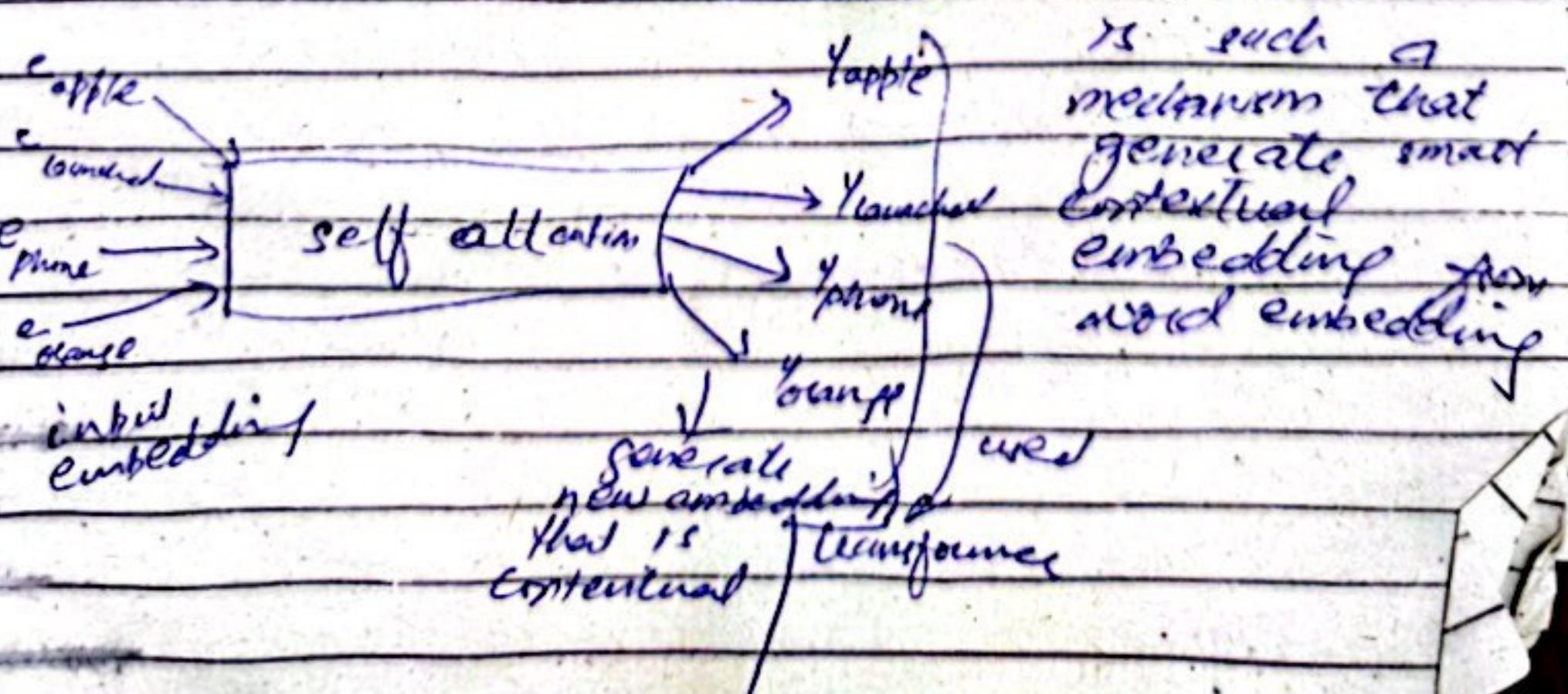
means If see in sentence word
"launched" or "plane", it automatically
increase their embedding.

$$\begin{bmatrix} 0.9 & 0.3 \\ 0.3 & 0.8 \end{bmatrix}$$

ideally it can't confuse on "orange" word.

but our word embedding can't do
because it is static

So solve this problem Self attention
came.



SELF - Attention

Any NLP application you make

↳ How represent word into vector

Hello → [0.3 0.9 1.0]

One-hot encoding
Bow
Embedding
word semantic meaning change into words

but problem in word embedding is,
it is static. Once it create, you can't.

↪ financial institution

Money Bank	Rice Bank
[0.6 0.2 0.01 0.7]	[0.6 0.2 0.01 0.7]

place ↑

But bank embedding is exactly same. It is problem.

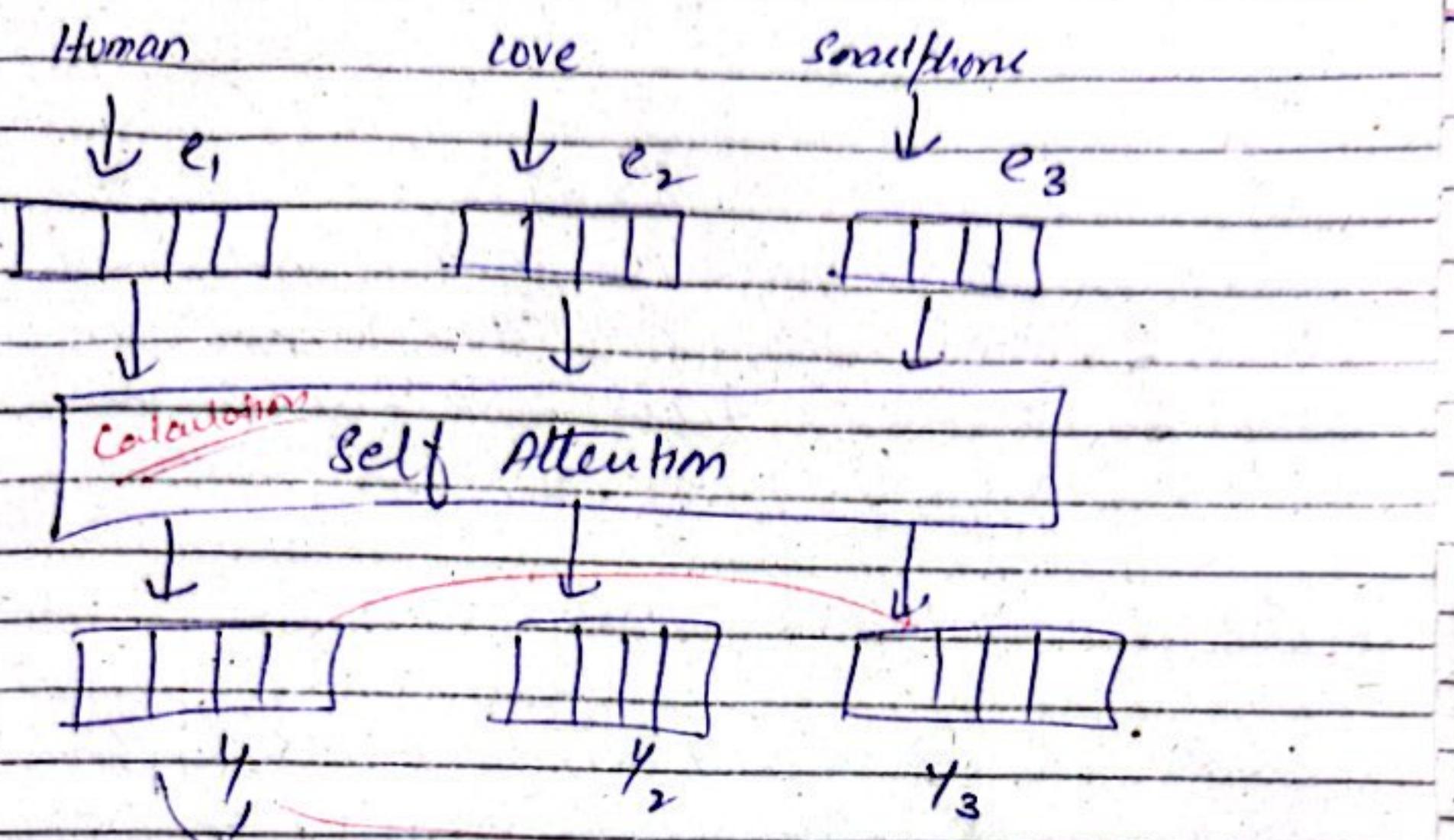
So solve this problem, we need contextual embedding.

- ↳ dynamically change
- ↳ context-aware so content it change.

↳ self attention solve this problem

causal

static embedding → dynamic
Contextual
embedding



It understand current word context
Content it need.

First principle Approach

Sen 1 - money bank grows
Sen 2 - river bank flows

} bank is same word but used in different context
it's meaning is not same but in embedding it is same

Sen 1
bank ≠ bank

$$\text{bank} = 0.3 \text{ money} + 0.7 \text{ bank} + 0.1 \text{ grow}$$

Sen 2

$$\text{bank} = 0.5 \text{ river} + 0.4 \text{ bank} + 0.1 \text{ flows}$$

Now bank meaning in different context
is different.

Sent 1

$$\text{money} = 0.7 \text{ money} + 0.2 \text{ bank} + 0.1 \text{ grows}$$

$$\boxed{\text{bank} = 0.25 \text{ money} + 0.7 \text{ bank} + 0.05 \text{ grows}}$$

$$\text{grows} = 0.01 \text{ money} + 0.02 \text{ bank} + 0.07 \text{ grows.}$$

Sent 2

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

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

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

It's benefit is, if you focus on bank
on both sentence.

L.H.S is same but R.H.S is different
in both sentence.

⇒ Now do little bit change
words are change into embedding.

$$e_{\text{money}}^{\text{new}} = 0.7 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grows}}$$

$$e_{\text{bank}}^{\text{new}} = 0.25 e_{\text{money}} + 0.7 e_{\text{bank}} + 0.05 e_{\text{grows}}$$

vector →
n-dim

vector → n-dims

Now what represent number 0.7, 0.2, 0.1

$$e_{\text{new}} = 0.7 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grow}}$$

new embedding is made up of

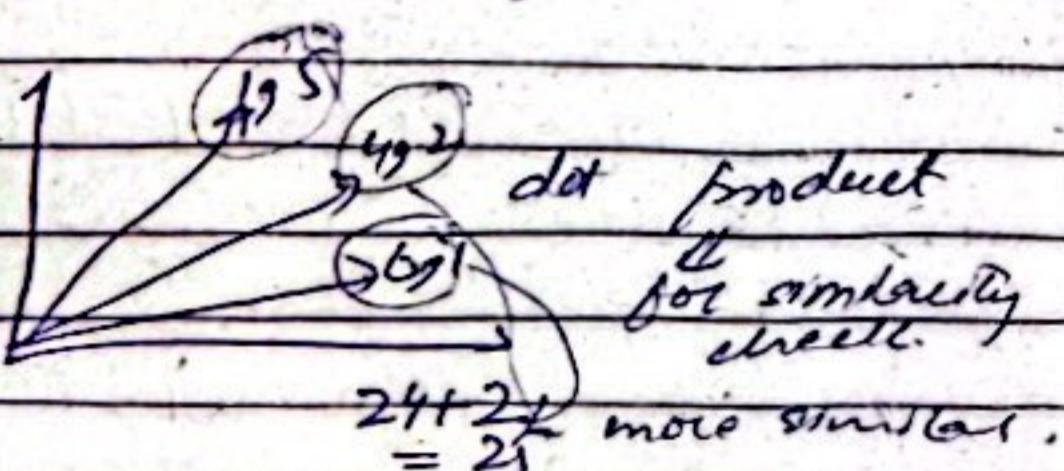
- 0.7 embedding of old money
- 0.2 " " " bank
- 0.1 " " " grow

actually it represent
similarity b/w words embeddings

$$\begin{aligned} 0.7 &\rightarrow e_{\text{money}} & e_{\text{money}} \\ 0.2 &\rightarrow e_{\text{money}} & e_{\text{bank}} \\ 0.1 &\rightarrow e_{\text{money}} & e_{\text{grow}} \end{aligned}$$

embedding is vector in high dimension
space

dot product
b/w word
embedding



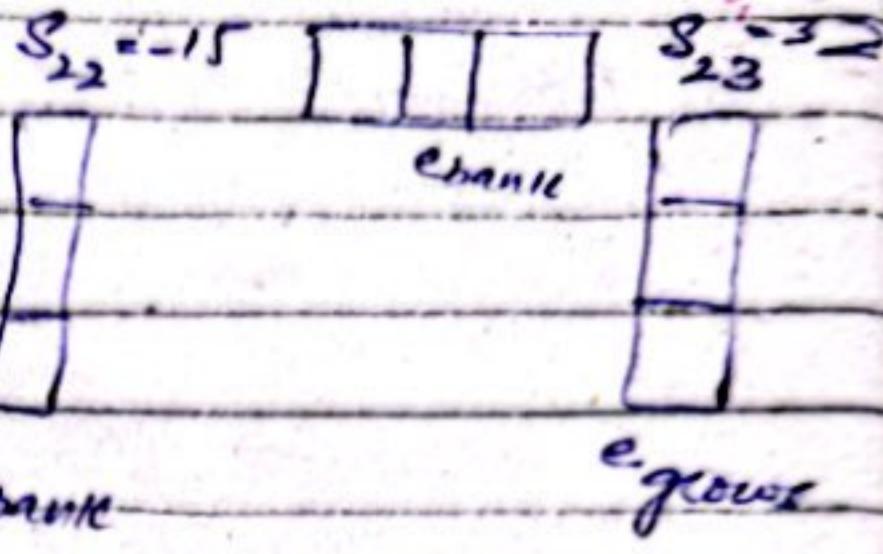
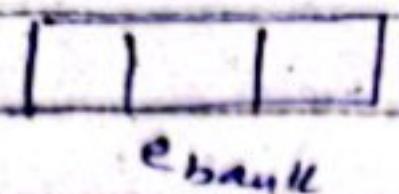
$$0.7 + 0.2 = 0.9$$

$$\begin{aligned} e_{\text{bank}}^{\text{new}} &= [e_{\text{bank}} \cdot e_{\text{money}}] \cdot e_{\text{money}} + \\ &\quad [e_{\text{bank}} \cdot e_{\text{bank}}] \cdot e_{\text{bank}} + \\ &\quad [e_{\text{bank}} \cdot e_{\text{grow}}] \cdot e_{\text{grow}}. \end{aligned}$$

after dot product and -

bank \rightarrow 2nd word
money \rightarrow 2nd word

sum = 1



money

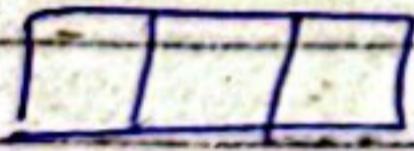
bank

grows

after softmax

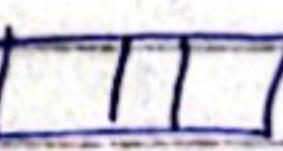
$w_{2,1}$

x

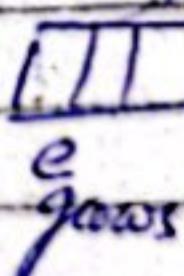


e_{bank}

+



e_{bank}



e_{grows}

$$\Rightarrow w_{2,1} + w_{2,2} + w_{2,3} = 1$$

$$w_{2,1} = \frac{e^{s_{2,1}}}{e^{s_{2,1}} + e^{s_{2,2}} + e^{s_{2,3}}}$$

$e_{\text{bank}}^{\text{new}}$

v_{bank}

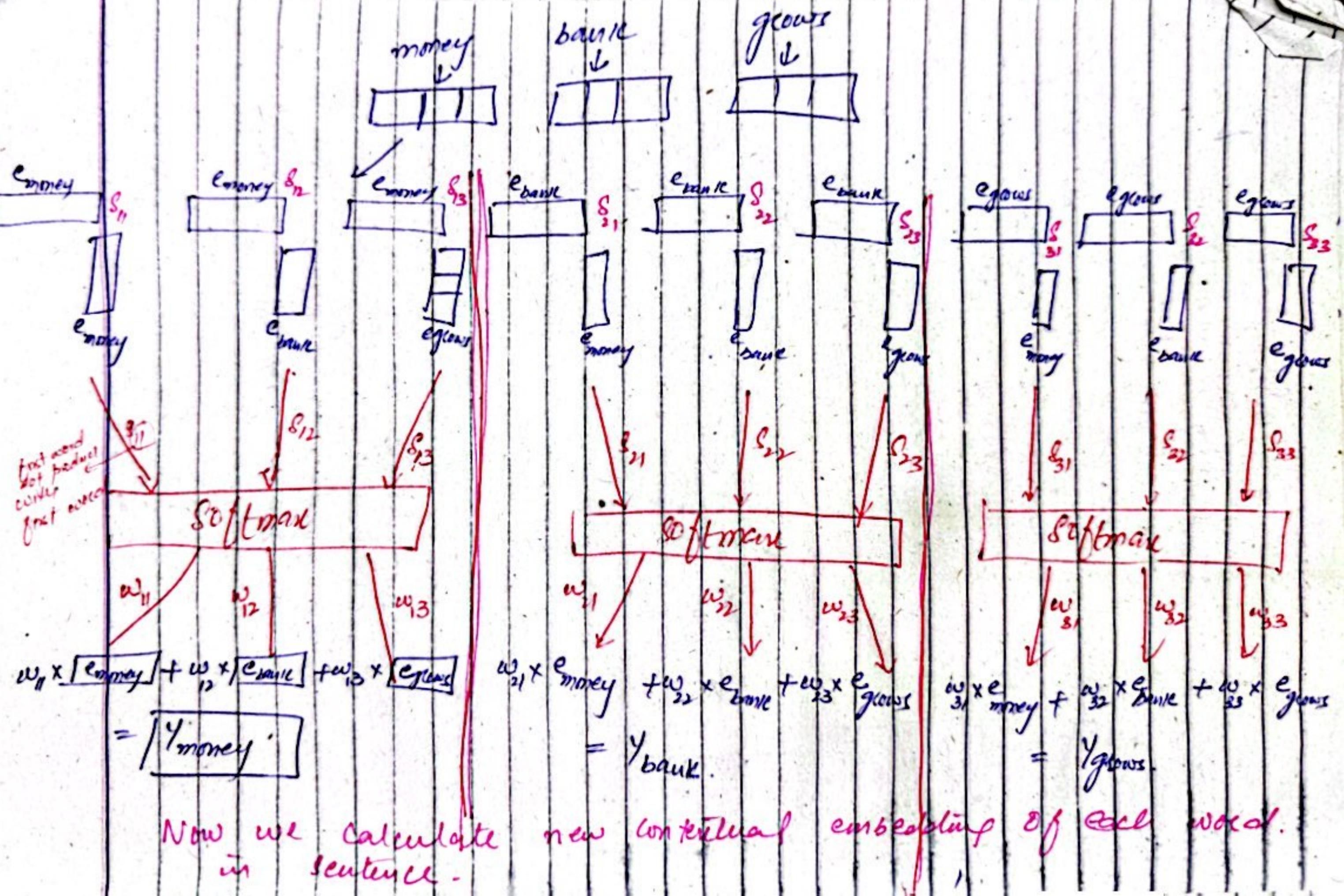
$$w_{2,2} = \frac{e^{s_{2,2}}}{e^{s_{2,1}} + e^{s_{2,2}} + e^{s_{2,3}}}$$

✓

contentual

word embedding

it is weighted sum of these embeddings



Point to consider :-

→ This is parallel operation.

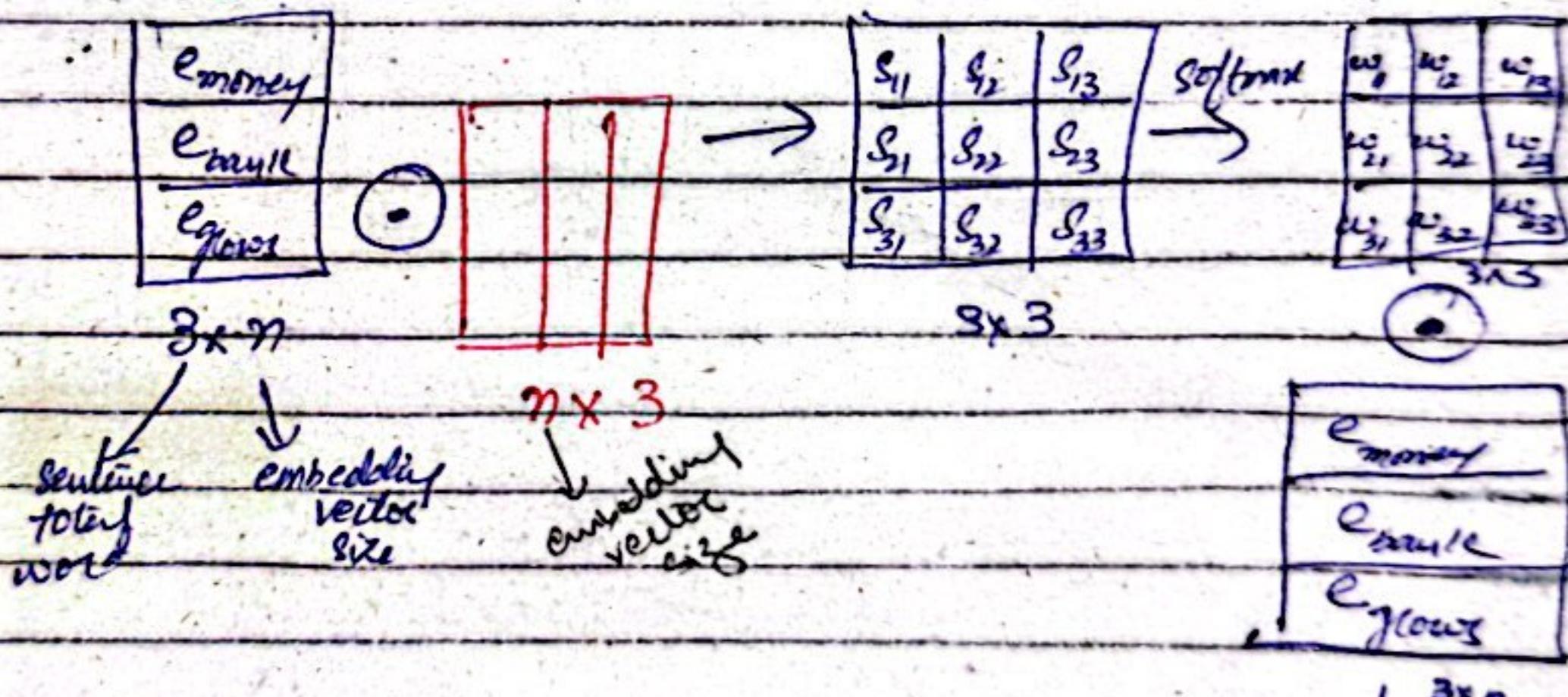
↳ no one depend on each other.

3rd not depend on 2 or 1 word.

→ There are no ^{learning} parameters involved.

↳ no learning based on task specific.

↳ it is general



Your sentence have 30 or 300

words, you at the same time

Contextual embedding generate
as many as you want.

entire thing is parallel

spread width

$3 \times n$

Ymoney

Ymusic

Ygames

$3 \times n$

Disadvantage

→ sequence ^{info} loss
first bank come or money.

hello → e, → y,
^{contentual}

↳ it can't learn anything
from data. ~~ans | under~~
~~how are you | we are~~
it learn from ^{current} sentence. ~~I am good | you are also~~

it is general contentual embedding
not task contentual embedding)

piece of cake $\hookrightarrow 1 \text{st}^1 2^2$

$\downarrow \quad \downarrow \quad \downarrow$ ~~but in our~~
 $e_1 \quad e_2 \quad e_3$ state

$\downarrow \quad \downarrow \quad \downarrow$
 $y_1 \quad y_2 \quad y_3$

~~16(16)~~ \hookrightarrow
it can't generate
task contentual
embedding.

it is general, it work on many
cases. but it can't work where task
contentual.

break a leg \rightarrow \rightarrow \rightarrow
 \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow

only contextual embedding generation is sufficient.

We want task contextual embedding.

Sentiment problem \Rightarrow based on sentiment problem if generate contextual embedding

translation problem \Rightarrow based on translation problem, if generate contextual embedding

we add some weight & bias, self attention problem learn, in this data that words how used.

What learning based on data not only one sentence.

Problems (Revise)

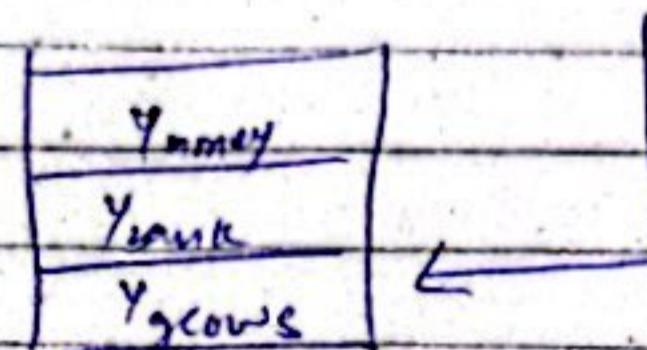
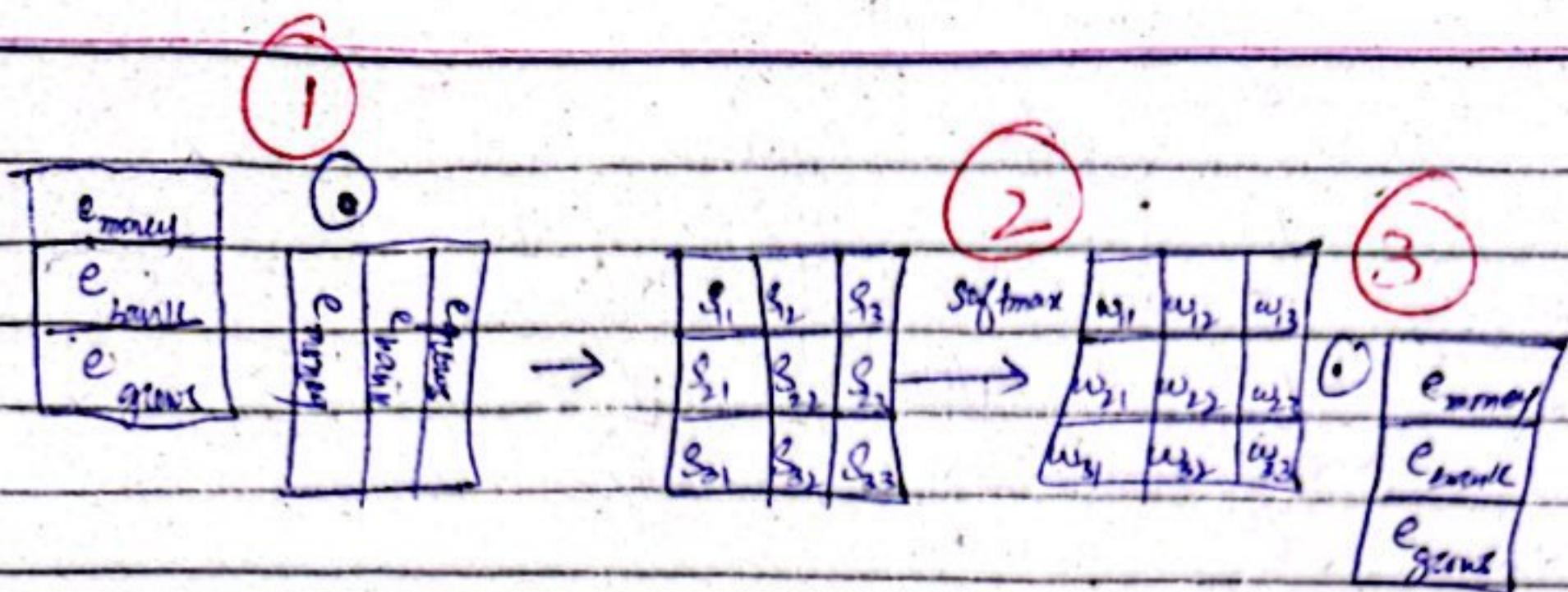
word \rightarrow number \rightarrow embedding (semantic)

| not capture
contentual
e.g. river bank
money bank

↓
Contentual embedding

| \rightarrow self attention

| general - task
no parameter parameter



where learnable parameters add?

- 1) dot product ?
- 2) softmax \times no weight introduce
- 3) dot product ?

only 1 and 3 where weight are introduce

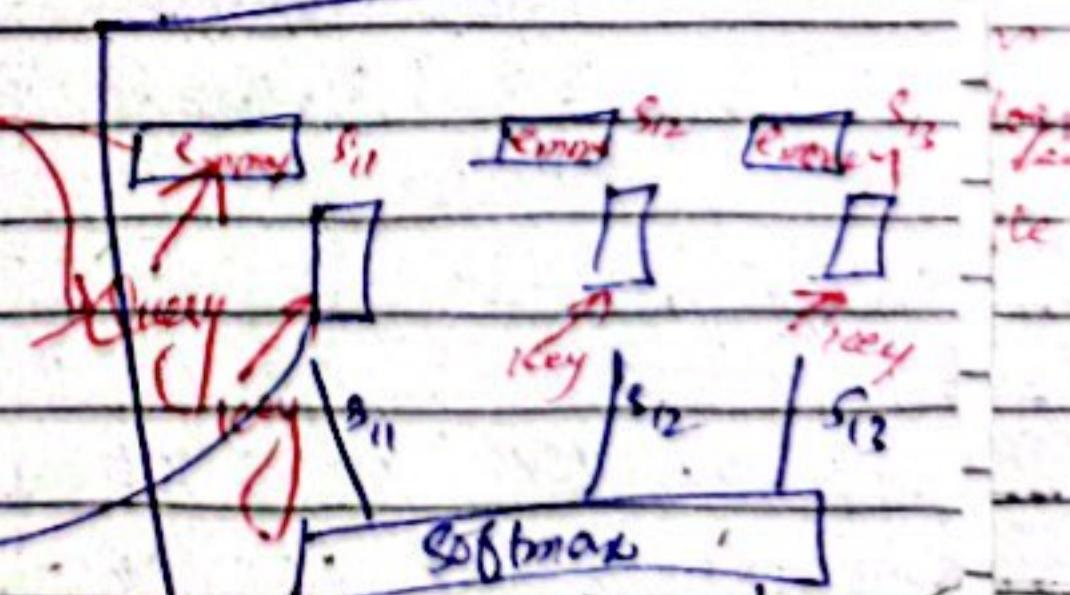
when we see "emmy" only when we "emmy" only

$$\boxed{e_{\text{emmy}}} e_{\text{emmy}} + e_{\text{emmy}} \downarrow \quad \boxed{e_{\text{bank}}} e_{\text{bank}} + e_{\text{bank}}$$

it ^{are} remaining all
embedding

how much similarity
b/w you and me.

Query



when query ask question, reply

answer, that is "key"

that. it answer its "values"

$$\begin{aligned} \text{weighted sum} &= w_{11} \cdot e_{\text{emmy}} + w_{12} \cdot e_{\text{bank}} + w_{13} \cdot e_{\text{goats}} \\ &= \boxed{w_{11} \cdot e_{\text{emmy}}} + \boxed{w_{12} \cdot e_{\text{bank}}} + \boxed{w_{13} \cdot e_{\text{goats}}} \end{aligned}$$

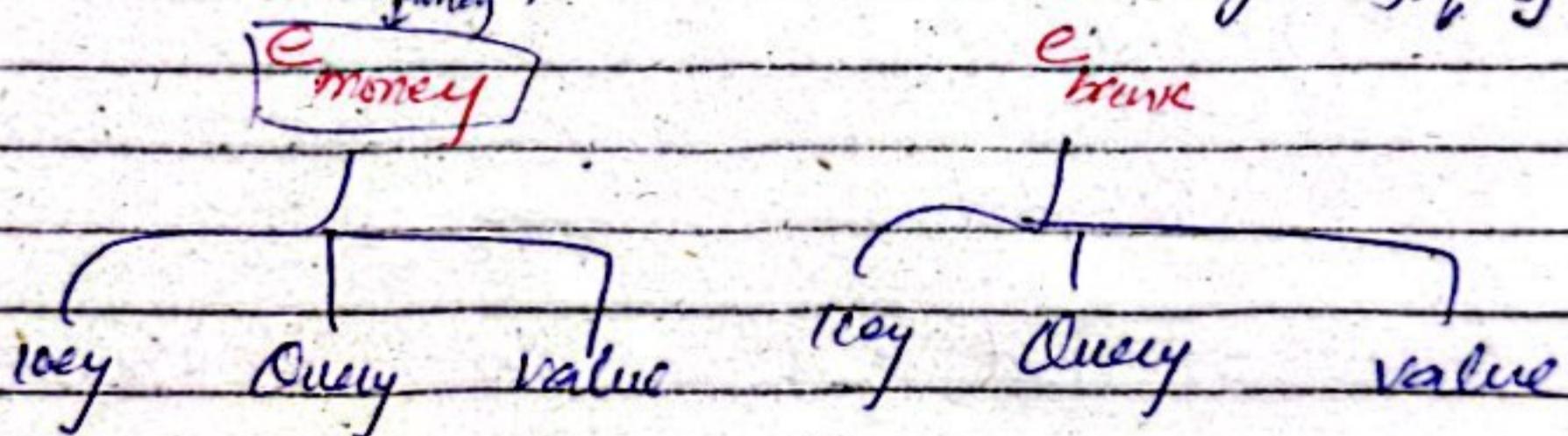
values.

Dict value value value
 $d = \{ a: 2, b: 3, c: 4 \}$
 ↓ ↓ ↓
 key key key

$d[a] = ?$

Query

so embedding we used in three ways key, query, value



During

Contentual embedding calculation process, every embedding vector three roles play

1- Key

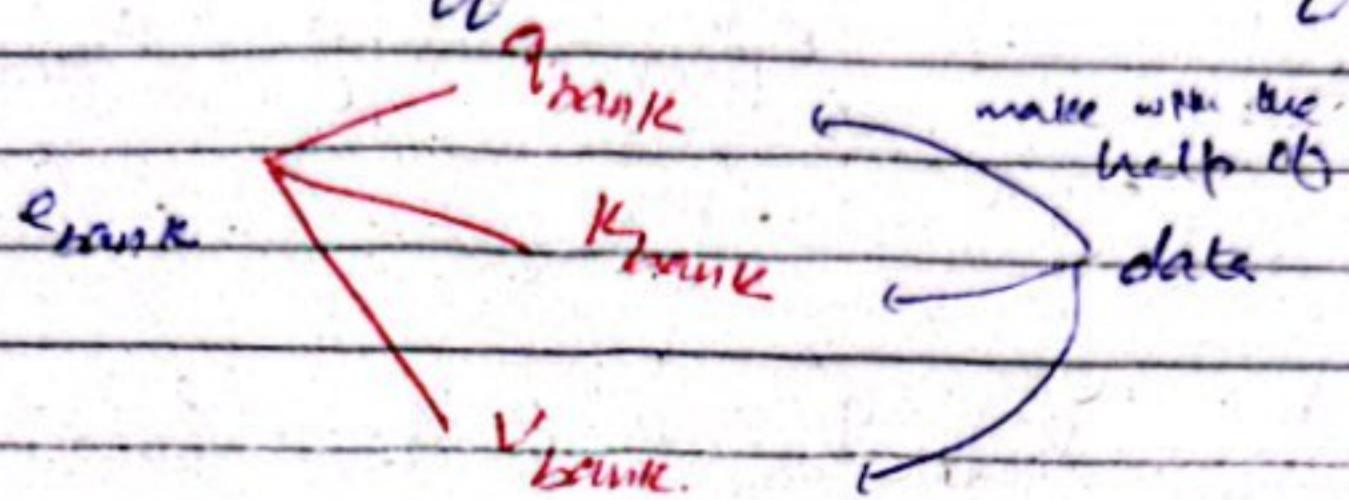
2- Query

3- Value

problems

Key :- that vector or not make sense
 value is a " " " "
 query is " " " " "

ideally three different vector made from it



Q what benefit from it?

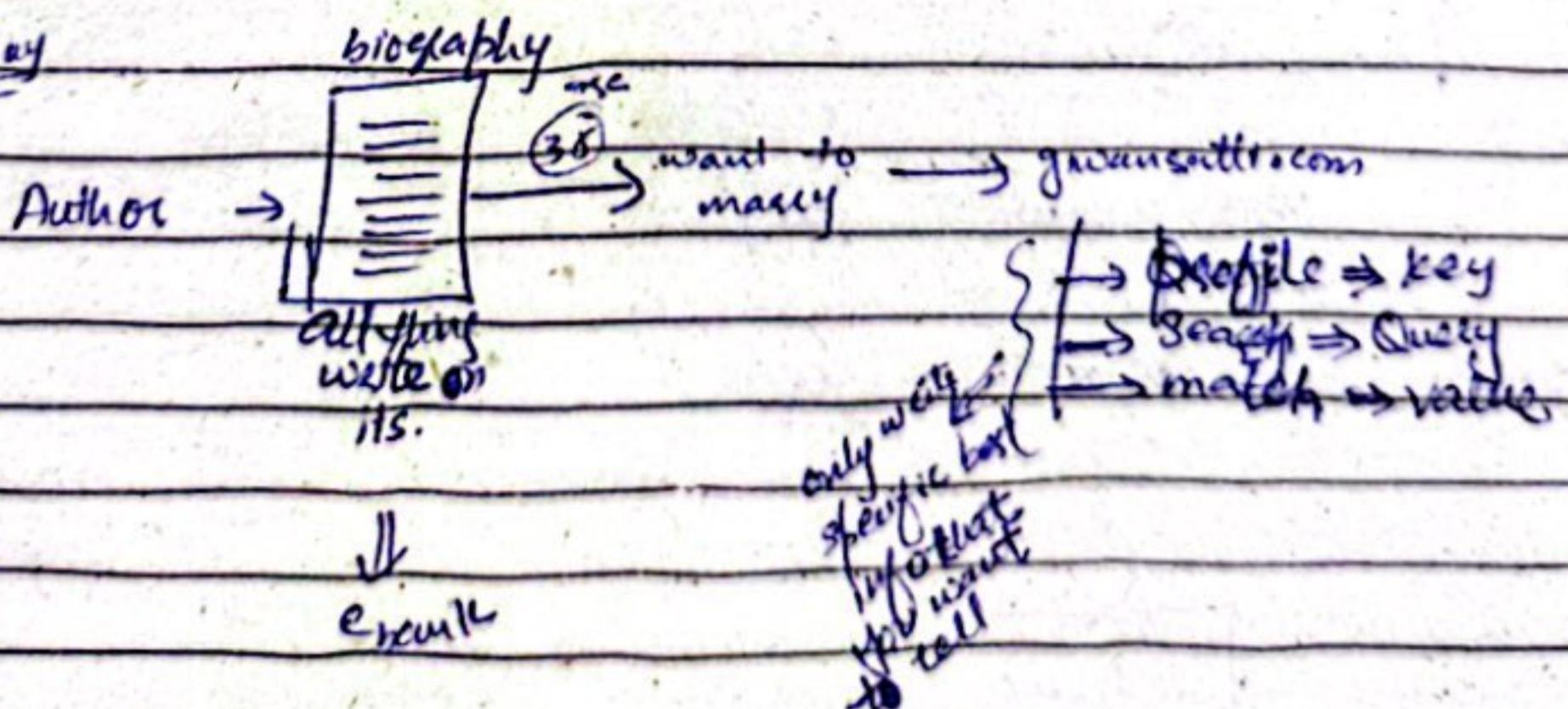
A: Separation of concerns.

if I put all work on one vector, it can't perform perfectly.

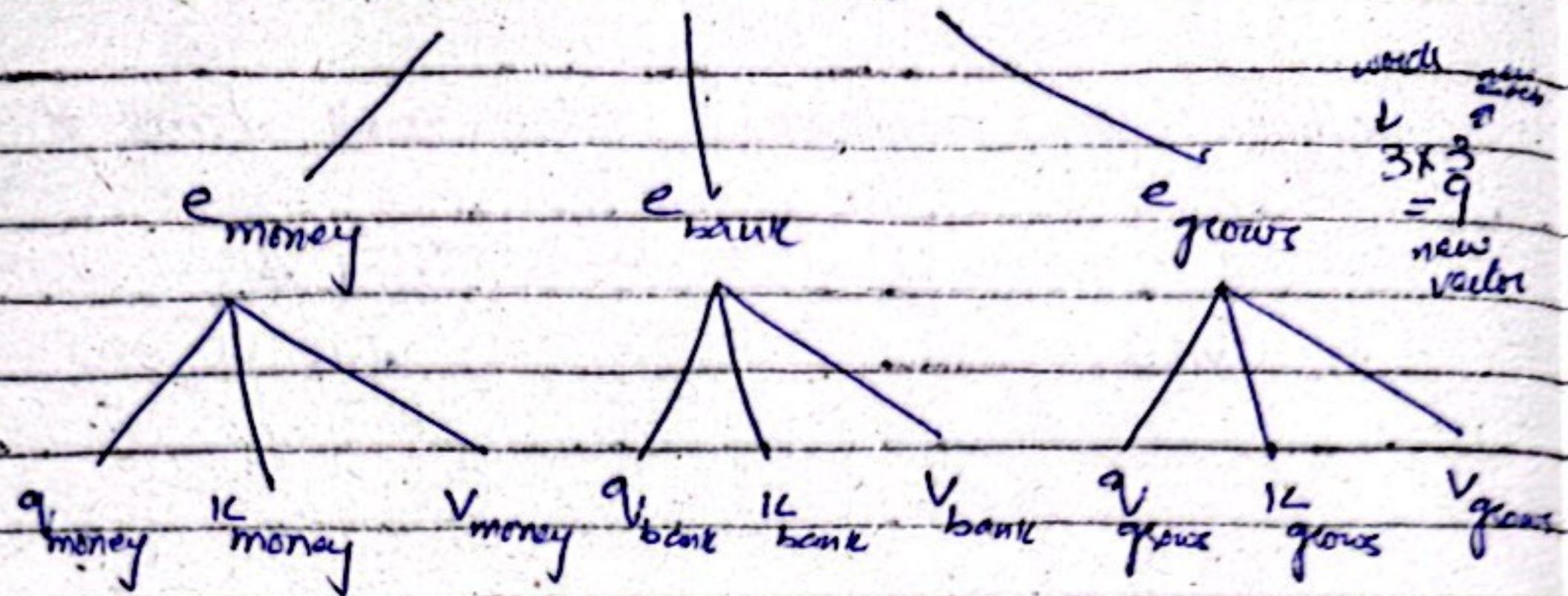
but if I transform it three different vector, it will do its own work perfectly.

example

say



money bank grows



money

continual Embedding

$$w_{\text{money}} \quad s_1 \\ \boxed{v_{\text{money}}} \quad \boxed{s_2} \\ q_{\text{money}} \quad s_3$$

$$k_{\text{money}} \quad k_{\text{bank}} \quad k_{\text{grows}}$$

$$s_{11}$$

$$s_{12}$$

$$s_{13}$$

softmax

$$w_{11}$$

$$w_{12}$$

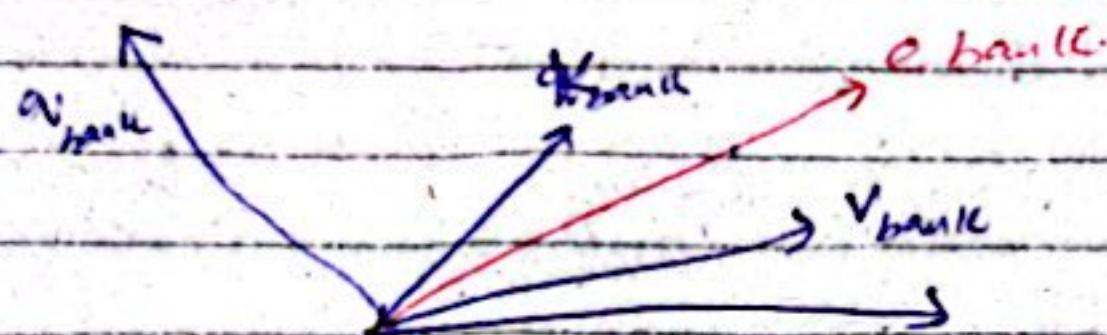
$$w_{13}$$

$$w_{11} \times v_{\text{money}} + w_{12} \times v_{\text{bank}} + w_{13} \times v_{\text{grows}}$$

y_{money}

~~How~~ given embedding, how query, key and value vector make.

→ we know we take help of dot product.



→ To make new vectors from one vector

① magnitude (increase or decrease vector)

② linear transformation.

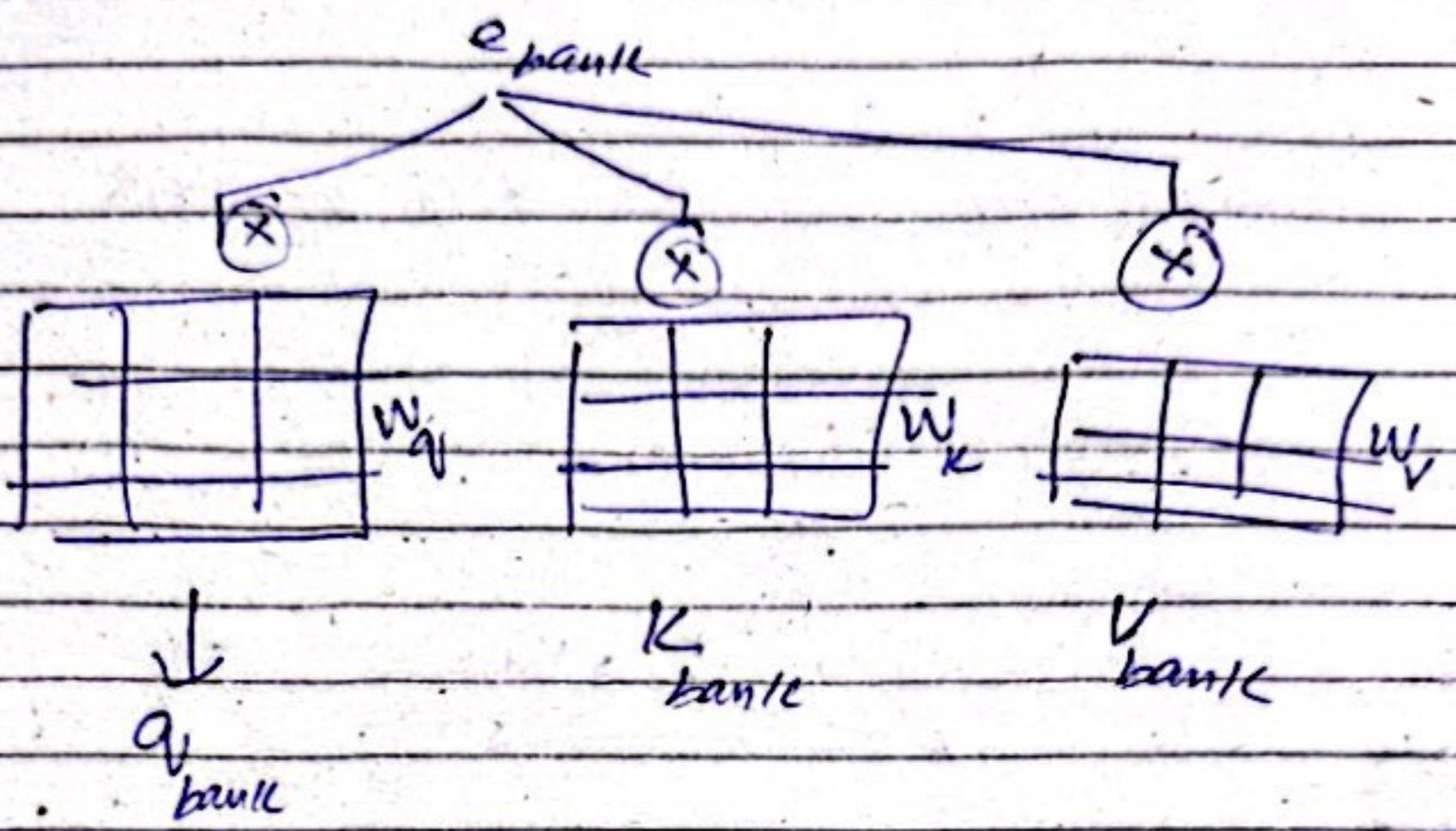
vector $x \begin{bmatrix} & \\ & \end{bmatrix}$

matrix

↳ magnitude and direction
both change

Scaling is not best, so linear transformation is best in our case.

~~Embedding for TPR~~



What values are in matrix?

w_q, w_k and w_v

\Rightarrow we decide new vector of mode from original embedding by the help of data.

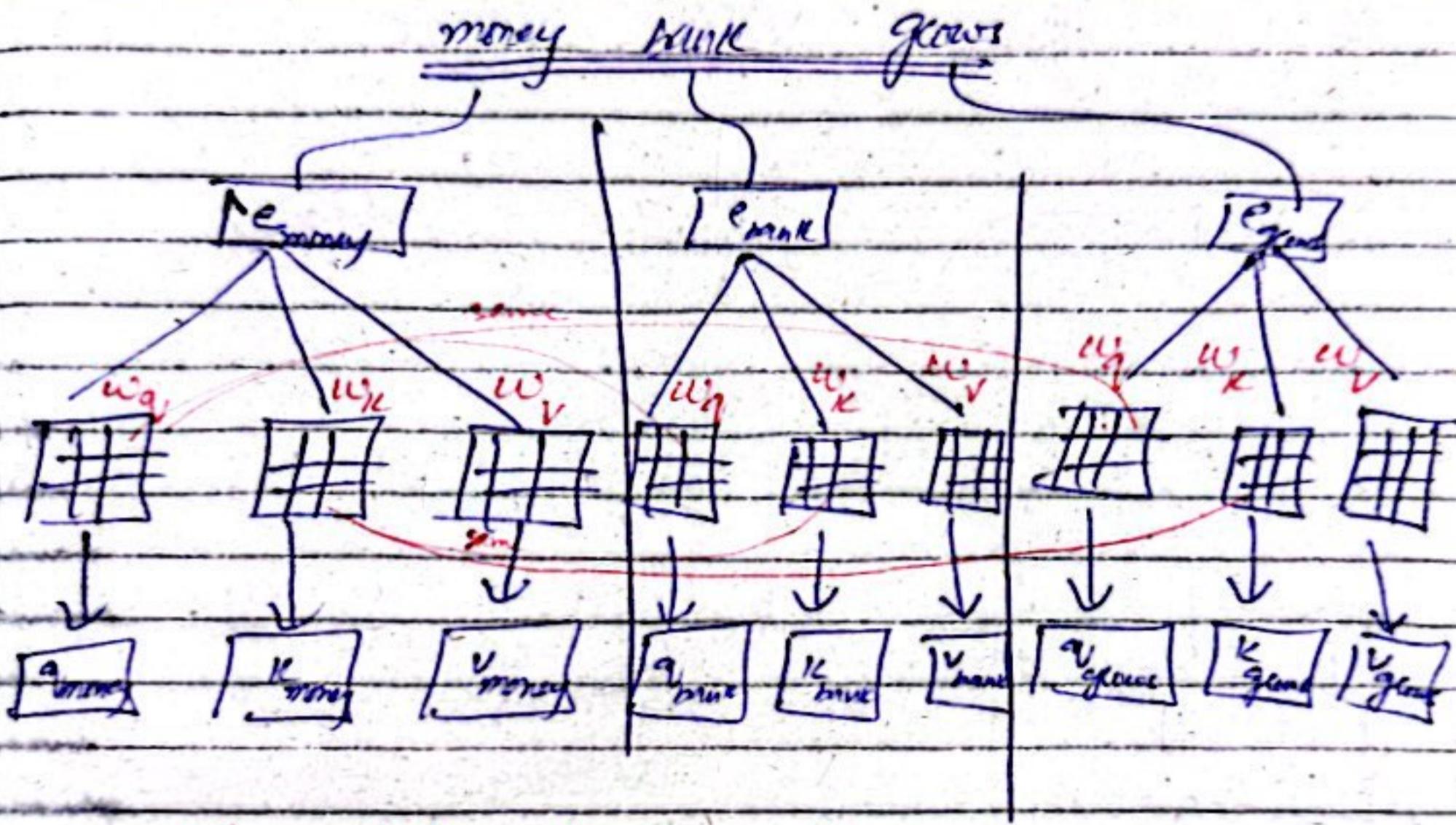
\Rightarrow we not decide what number is matrix, training process decide

what we will do?

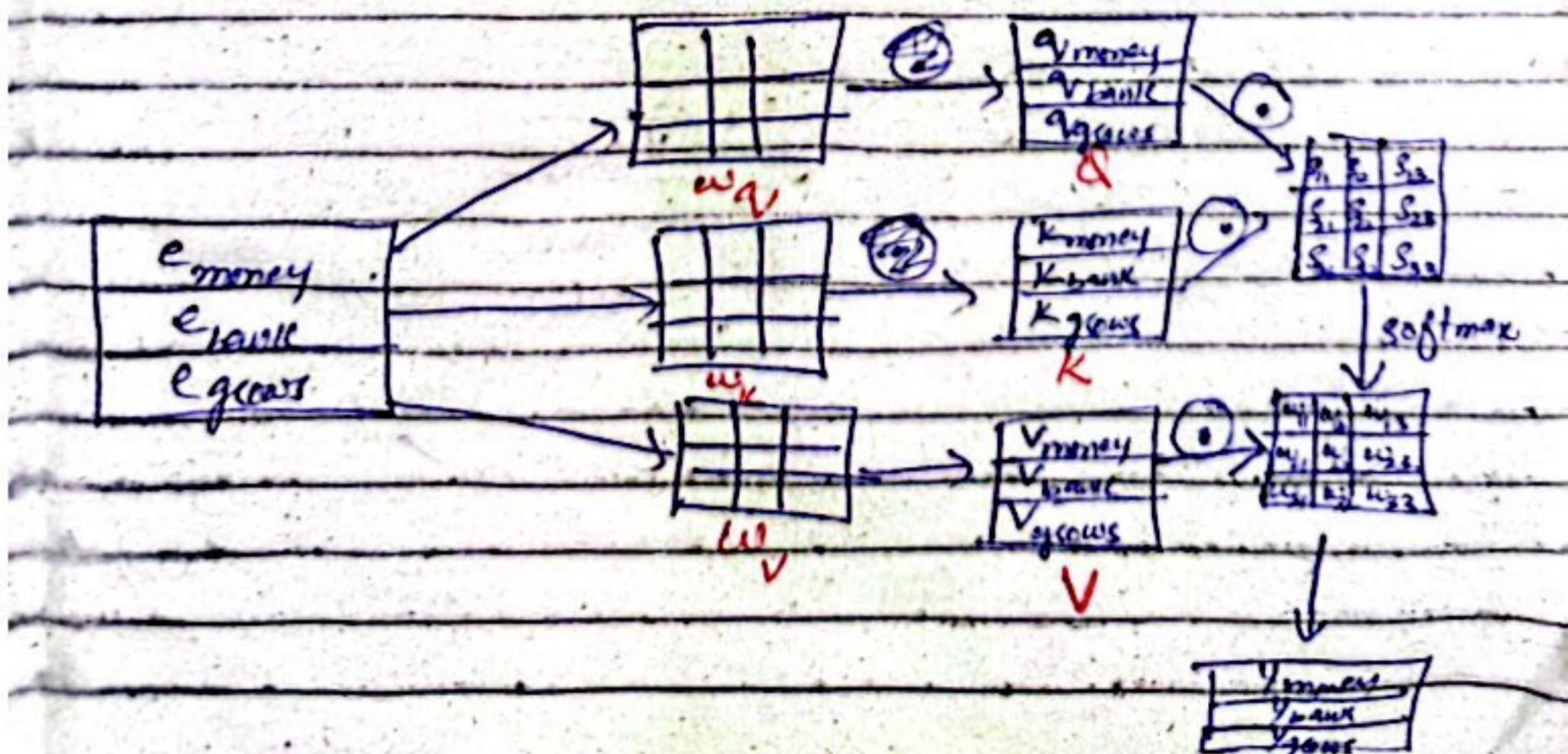
We will start with random value (weights)

and it generate query, key, and value that this 'translation' occurs some mistake occur, back propagate weight update and next sentence goes and so on.

and after learning it gain coded value.

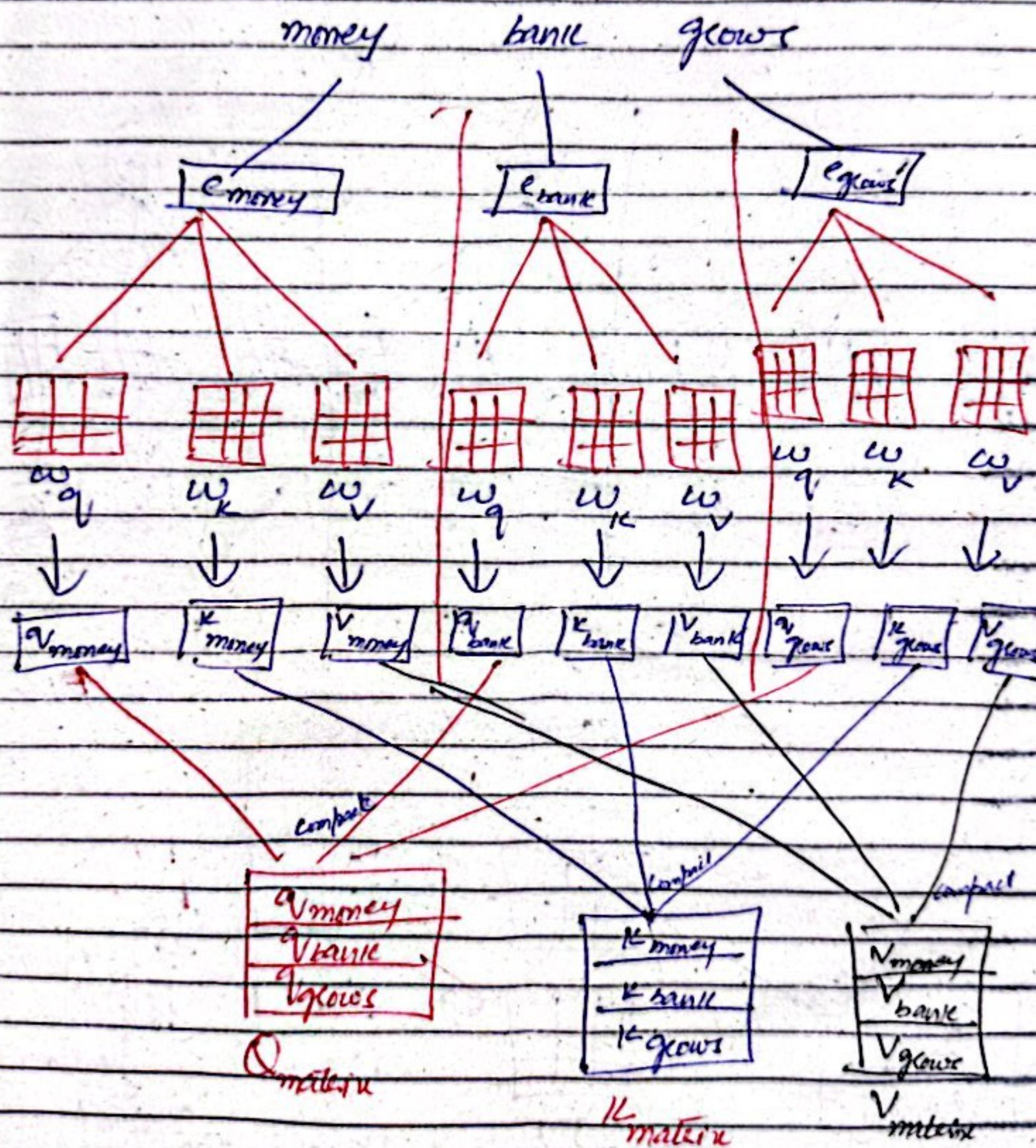


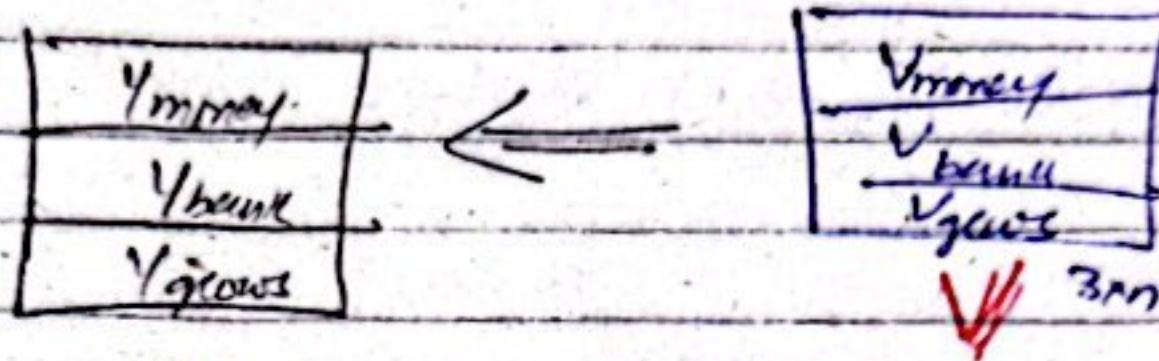
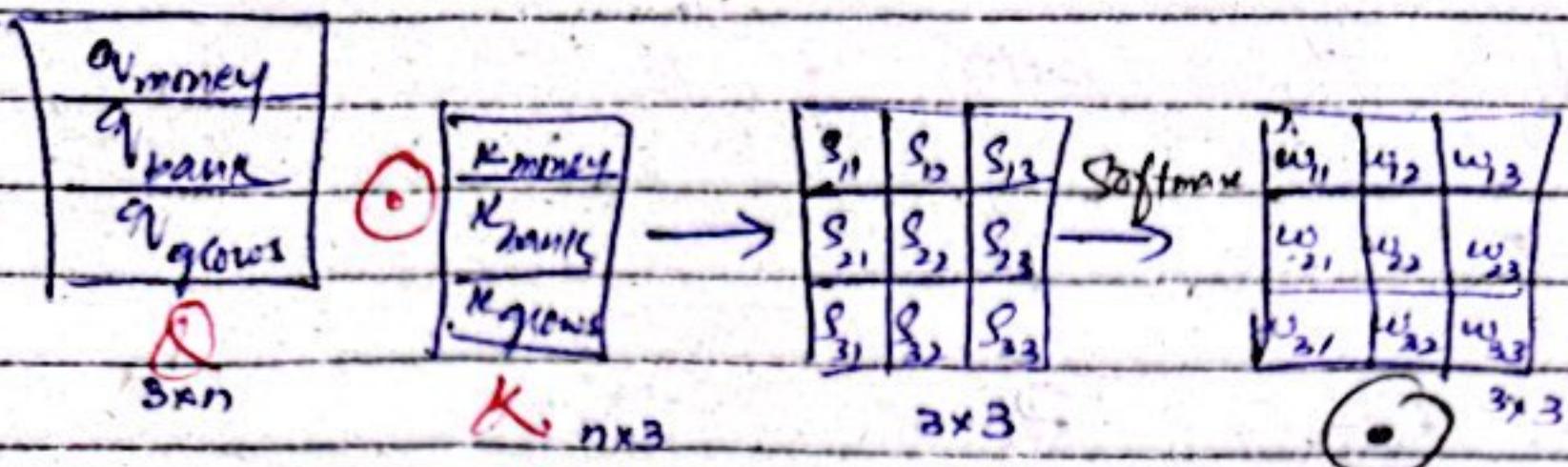
all words have same
 w_q, w_k, w_v . These are parameter learnable
 during training
it also called ceaser



Scaled Dot product Attention

Recap





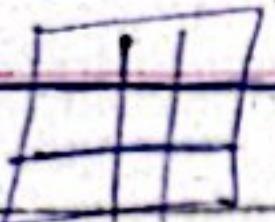
$$\text{attention}_{(Q,K,V)} = \left[\text{softmax}((Q \cdot K^T)) \right] \cdot V$$

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

Scaled dot product production

$$\text{softmax} \left(\frac{Q \cdot K'}{\sqrt{d_K}} \right) \cdot V$$

Scale \downarrow Unstable gradient creation / Numerical issues
 Why scale and why for d_K ?



$$K \quad d_K = 3$$

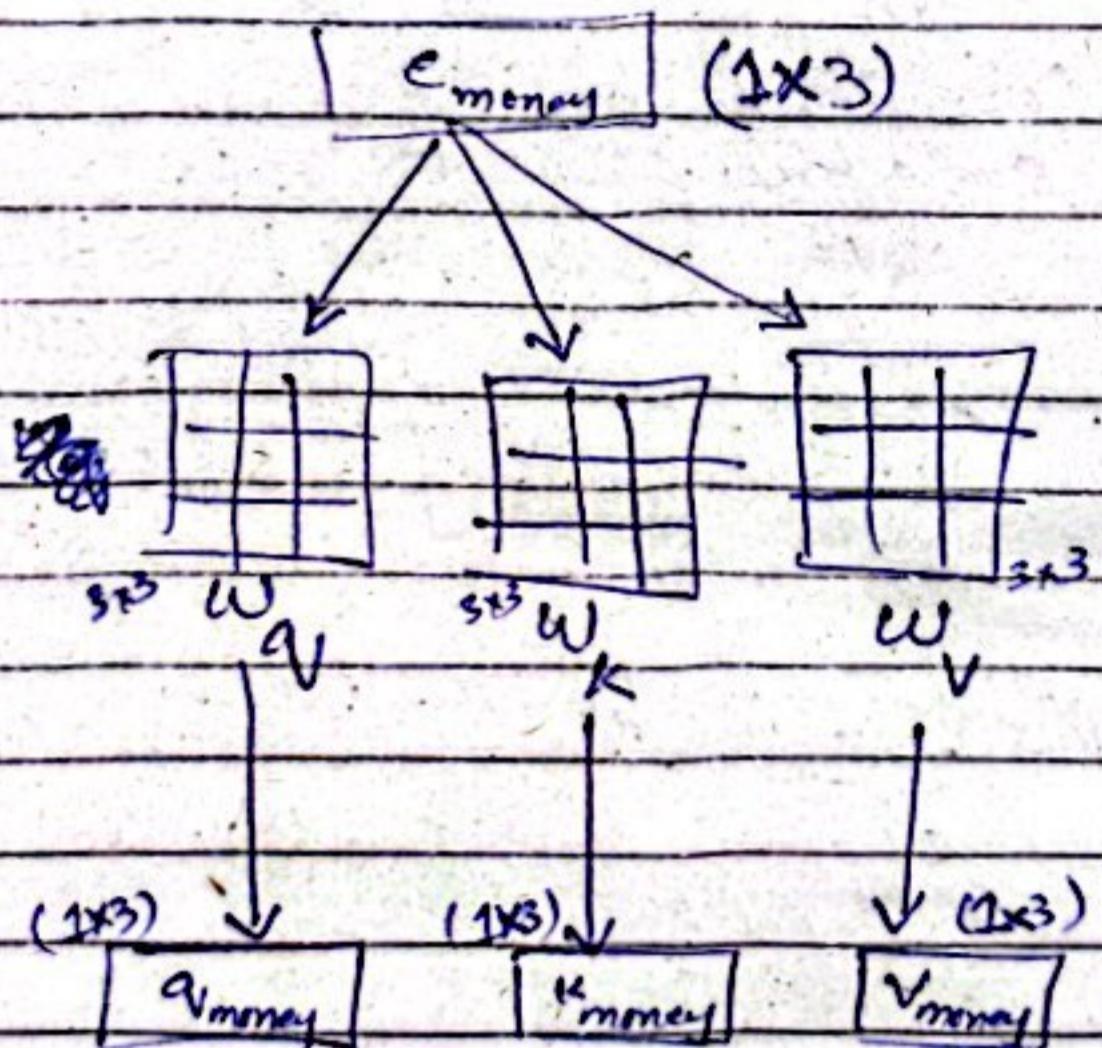
d_K = dimensionality of K -vector

→ optional
→ encoding
→ transform
→ activation

e_{money} 3 dim (1×3)

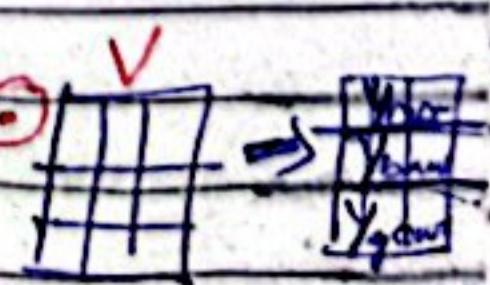
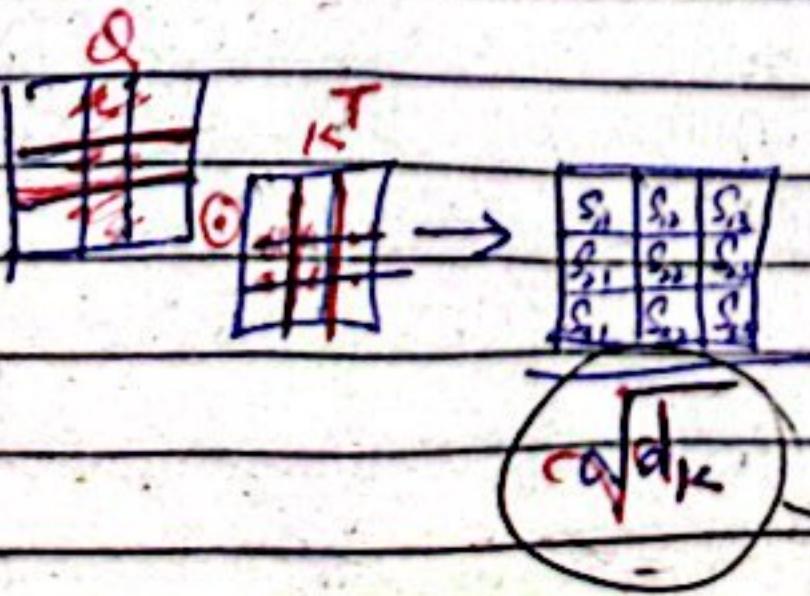
$$d(w_q) = d(w_{i_0}) = d(w)$$

it also different
 $d(w_q) \neq d(w_x) \neq d(w_v)$



$$d_q = d_K = d_v$$

in case
out it may
be different



$\alpha \cdot \frac{1}{d_K}$

why scale?

Nature of Dot-product

mode of

$Q \rightarrow 3$ vectors

$K \rightarrow 3$ "

$3 \times 3 \Rightarrow 9$ vector
dot product
9 numbers

↳
mean (μ)
variance (σ^2)

Nature of dot product

low dim \rightarrow dot product. \rightarrow low variance

high " " " " \rightarrow high "

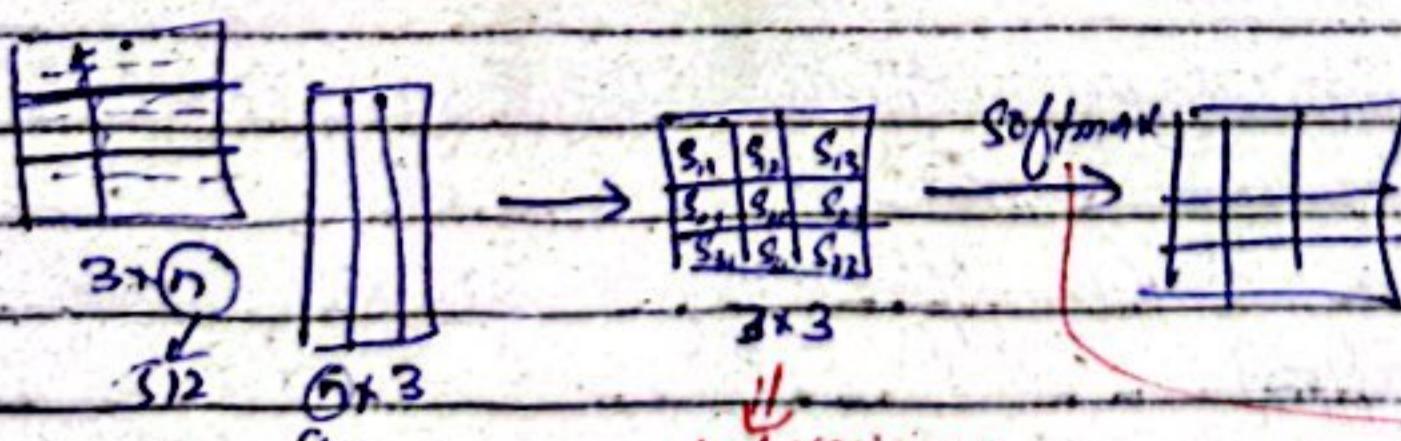
2-dim

$$\begin{array}{l} 1) [1 \ 2] \odot [3, 2] \\ 2) \quad \quad \quad = b \\ 3) \quad \quad \quad = c \\ 4) \quad \quad \quad = d \\ 5) \quad \quad \quad = e \end{array} \quad \left. \begin{array}{l} -9 \\ = 5 \\ = c^2 \\ = d^2 \\ = e^2 \end{array} \right\}$$

3-dim

$$\begin{array}{l} 1) [- - -] \odot [- - -] = f \\ 2) \quad \quad \quad = g \\ 3) \quad \quad \quad = h \\ 4) \quad \quad \quad = i \\ 5) \quad \quad \quad = j \end{array} \quad \left. \begin{array}{l} = f^2 \\ = g^2 \\ = h^2 \\ = i^2 \\ = j^2 \end{array} \right\}$$

$$g_2 > 6$$



↳
but variance
is very
high

some numbers too large
" " " small

→ convert numbers
into probability

↳
if you give
too large, it
prob is high
if you give
too small number,
it prob is low.

$$\text{softmax}(4, 5) \quad \text{softmax}(1, 10)$$

$$\approx (26.89, 73.11)$$

$$\approx (0.01, 99.99)$$

If some number is too large, softmax give that number high prob and some number give too low prob

gab too many exist, problem come, we train whole process, it mean back propagate, gradient calculate and that basis weight update.

Whatever it see large number, it whole training try to correct that number and smaller number ignore, vanishing gradient problem occur, parameters never update, training process 'deadly'.

To reduce variance, automatically when apply softmax, such prob assign that is approximately equal and training correctly occur.

To solve that problem?

↳ you can't take high dimension, take low dimension embedding but in case you can't access more useful information.

↳ high dimension \rightarrow high variance

↓ problem
low variance

$$x = 10, 20, 30, 40, 50, 60, 70$$

$$\text{val}(A) = 400$$

$$x = \frac{1}{10}, \frac{2}{10}, \frac{3}{10}, \frac{4}{10}, \frac{5}{10}, \frac{6}{10}, \frac{7}{10}$$

$$x = 1, 2, 3, 4, 5, 6, 7$$

$$\text{val}(A) = 7$$

why used d_n ?

- 1) The magnitude of dot matrix product grows with d_n .
- 2) Large dot products can cause (softmax activation) to push values into regions of small gradients, making optimization harder.



$$v_1 \begin{pmatrix} v_4 \\ v_5 \\ v_6 \end{pmatrix}$$

The softmax function has smaller changes (or gradients) when the inputs are large.

\Rightarrow one value becomes much larger than the other.

Small gradients make it harder for the model to adjust and improve during training.

* Softmax will treat it most important and ignore smaller values.

* Gradient becomes very small.

Variance shift $\rightarrow \text{var}(X) \neq \text{var}(Y)$

1-dim

Random variable

$$(a) \rightarrow (b) \rightarrow a \cdot b$$

expected value

$$(c) \rightarrow (d) \rightarrow a \cdot c$$

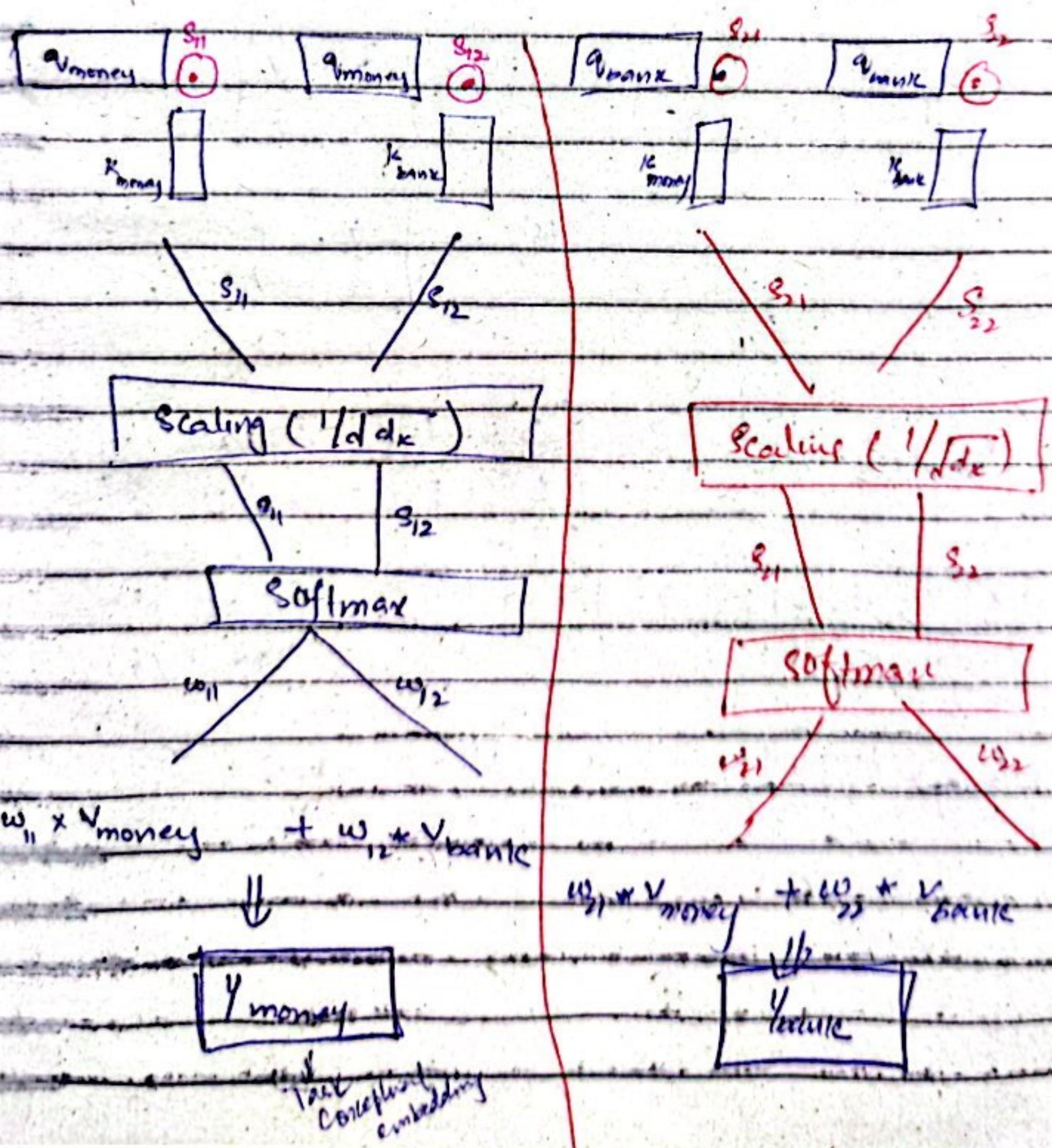
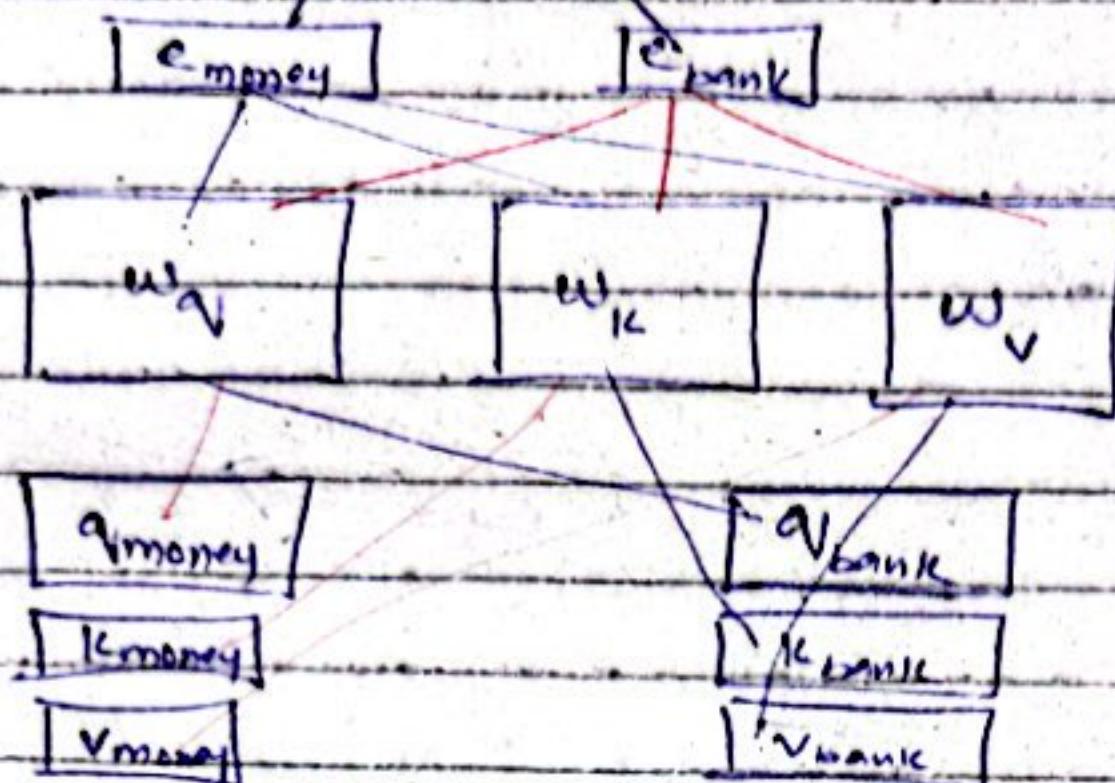
2-dim

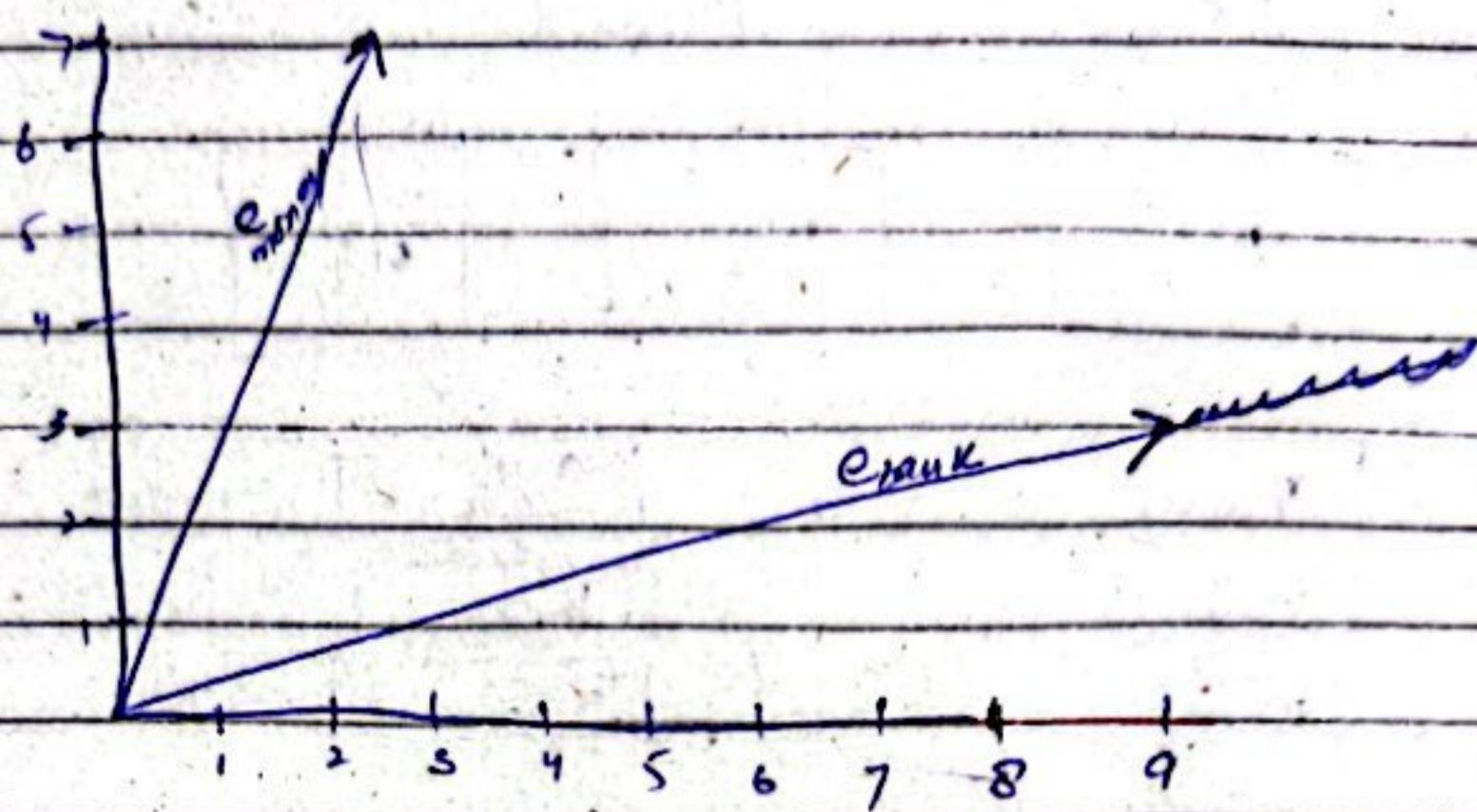
$$(a, b) \rightarrow a \cdot c + b \cdot d \quad y \rightarrow \text{var}_y(Y)$$

$$(e, f) \rightarrow a \cdot e + b \cdot f$$

$$(g, h) \rightarrow a \cdot g + b \cdot h$$

money bank



$e_{\text{money}} [2, 4)$ $e_{\text{carrie}} [9, 3)$ 

$$w_9 = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad w_{10} = \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix}.$$

$$w_v = \begin{bmatrix} 9 & 1 \\ 2 & 1 \end{bmatrix}$$

$$e_{\text{money}} \cdot w_9 = \{2 \ 7\} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \{4+7 \ 2+4\} = \{11 \ 8\}$$

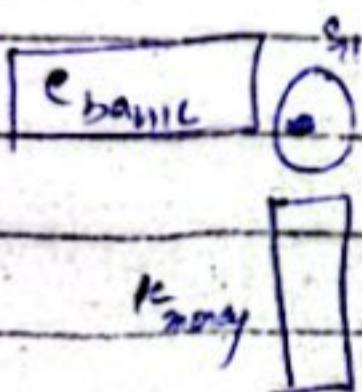
$$e_{\text{money}} \cdot w_{10} = \{2 \ 7\} \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} = \{6+35 \ 8+7\} = \{41 \ 15\}$$

$$e_{\text{money}} \cdot w_v = \{2 \ 7\} \begin{bmatrix} 9 & 1 \\ 2 & 1 \end{bmatrix} = \{8+14 \ 2+7\} = \{22 \ 9\}.$$

$$e_{\text{basic}} \cdot w_q = [9, 3] \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 18+3 & 3+6 \\ 2+6 & 9 \end{bmatrix} = \begin{bmatrix} 21 & 9 \\ 8 & 9 \end{bmatrix}$$

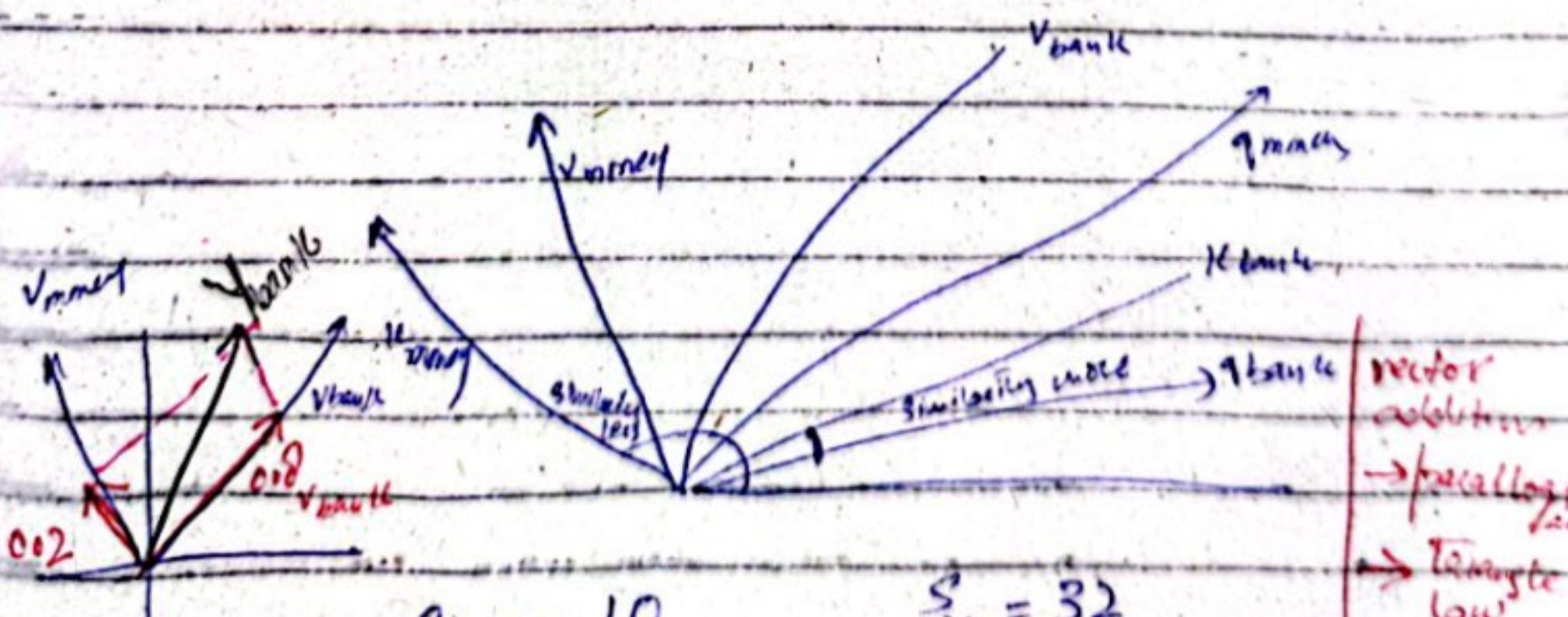
$$e_{\text{basic}} \cdot w_R = [9, 3] \begin{bmatrix} 3 & 4 \\ 5 & 1 \end{bmatrix} = \begin{bmatrix} 27+15 & 36+3 \\ 42+5 & 9 \end{bmatrix} = \begin{bmatrix} 42 & 39 \\ 47 & 9 \end{bmatrix}$$

$$e_{\text{basic}} \cdot w_V = [9, 3] \begin{bmatrix} 4 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 36+6 & 9+2 \\ 42+2 & 9 \end{bmatrix} = \begin{bmatrix} 42 & 11 \\ 44 & 9 \end{bmatrix}$$



$$e_{\text{basic}} \cdot k_{\text{mach}}^T = \begin{bmatrix} 2 & 3 \end{bmatrix} \begin{bmatrix} 3 & 1 \\ 1 & 5 \end{bmatrix}$$

$$e_{\text{basic}} = \begin{bmatrix} 82+95 \\ - \end{bmatrix} = \begin{bmatrix} - \end{bmatrix}$$



$$g_{21} = 10 \quad g_{12} = 32$$

Scaling

$$g'_{21} = \frac{10}{2} = 7.09$$

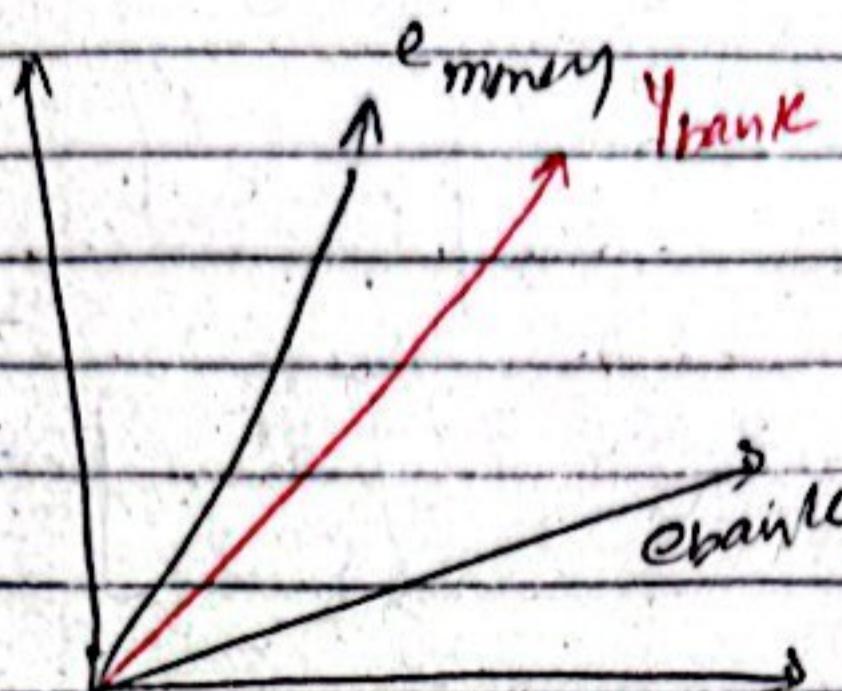
$$g'_{12} = 22.69$$

softmax?

$$w_{21} = 0.2$$

$$w_{12} = 0.8$$

money bank



bank embedding
is too far
then money
embedding
When we
apply self
attention then

Self attention on the basis
of content embedding
generate

Y bank continual
embedding if it
too closest to
money embedding

self attention acts like
gravity

~~self attention~~
~~self is content~~
~~usage~~

money word pull more toward
itself.

Why self attention called "self"?

why. self attention called "attention"?

Recap of Attention :-

language translation task solve using neural net.

→ that type of task (sequence to sequence)
use Encoder Decoder Architecture.

→ Encoder - feed sentence word by word, and
Encoder ^{sentence} process and it maintains hidden
state and the end hidden state
are (content vector) and that
content vector send to decoder.

→ content vector contain summary of
sentence.

Decoder pick content vector and steps
by step pass on output part.

This approach is good. but one big
problem in it.

you whole full burden on content
vector. (content vector is nothing but
just a set of numbers)

but in some way your sentence size is
greater than 30. (in that case, you
can't effectively summarize your sentence
in content vector).

then researcher notice sentence > 30,
translation quality bad. Then

attention mechanism came.

Attention :-

Instead after seeing whole sentence, its summary only put in one vector and send off it of decoder if you.

If decoder each timestep, If you send different different context vector and that content vector has info in which input words is useful to print that word.

With the help of hidden state we make that content vector

$$\alpha_i = \sum_{j=1}^n \alpha_{ij} h_j$$

to tell ~~at~~ this current ^{word} ~~current~~ final which input hidden state is useful

How much - total alpha's
= total input word x total output words

How α are calculate?

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

Living attention :-

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

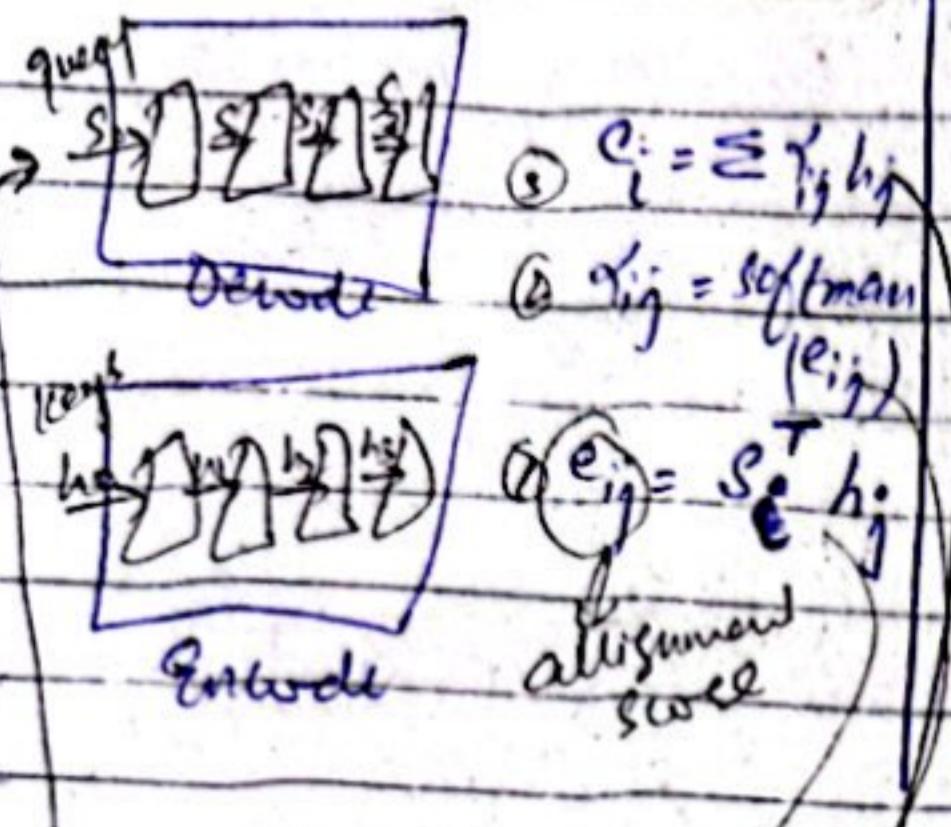
alignment score

$$c_i = \sum_{j=1}^n \alpha_{ij}$$

$$\alpha_{ij} = \frac{e_{ij}}{\sum_j e_{ij}}$$

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

Self attention



Turn off the light

e_1, e_2, e_3, e_4

y_1, y_2, y_3, y_4

contextual embedding

turn off the light

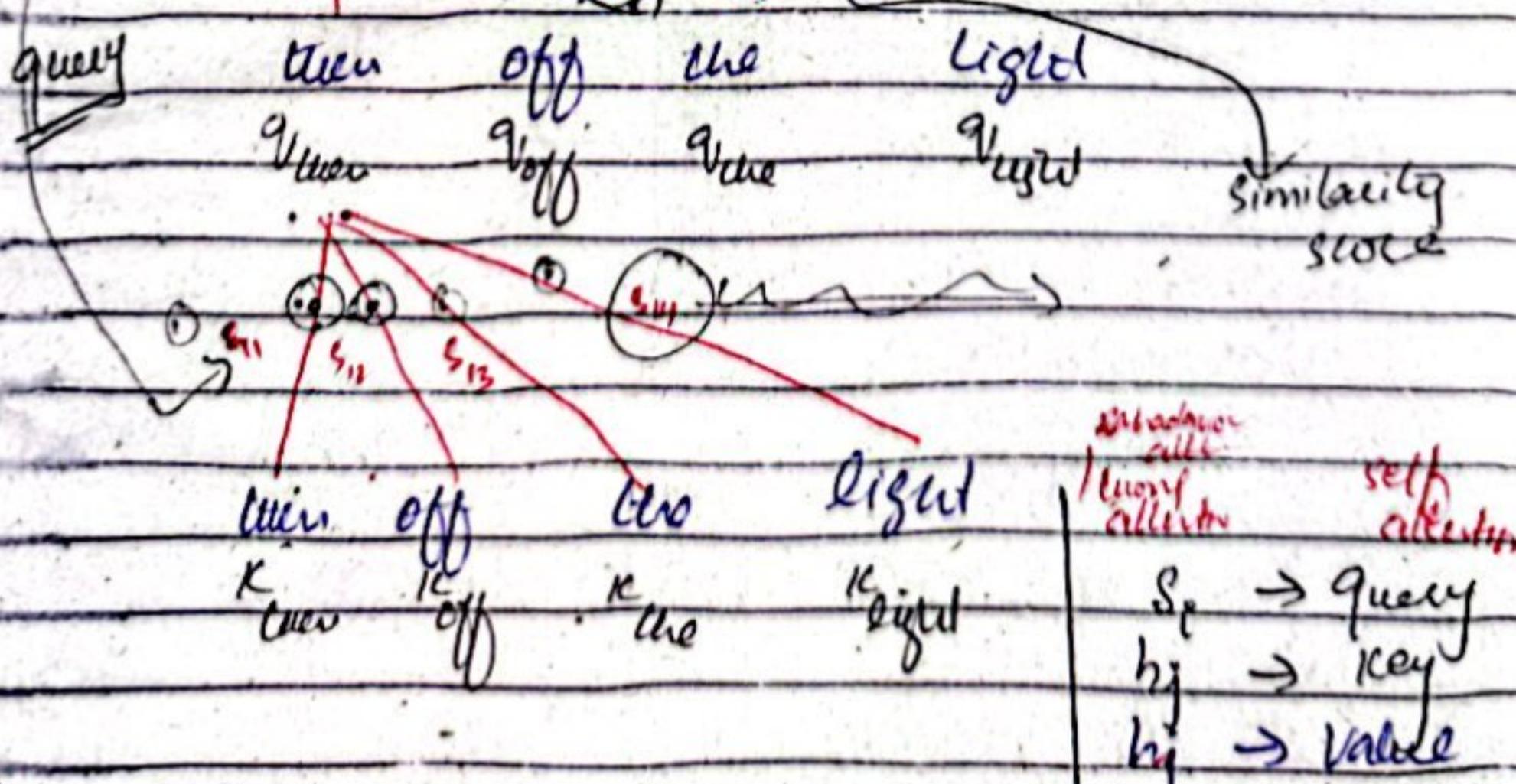
e_1, e_2, e_3, e_4

turn view off off off the the the

Similar

$$y_{\text{view}} = w_{11} \cdot v_{\text{view}} + w_{12} \cdot v_{\text{off}} + w_{13} \cdot v_{\text{the}} + w_{14} \cdot v_{\text{light}}$$

$$\Rightarrow \text{softmax}(\dots) \rightarrow w_{11}, w_{12}, w_{13}, w_{14}$$



$s_0, s_1, s_2, s_3 \Rightarrow \text{query}$

That's self attention called attention.

long attention find two different sequence alignment score deduce

Self attention (same sequence)

Intrasequence

enc

enc

similarity score

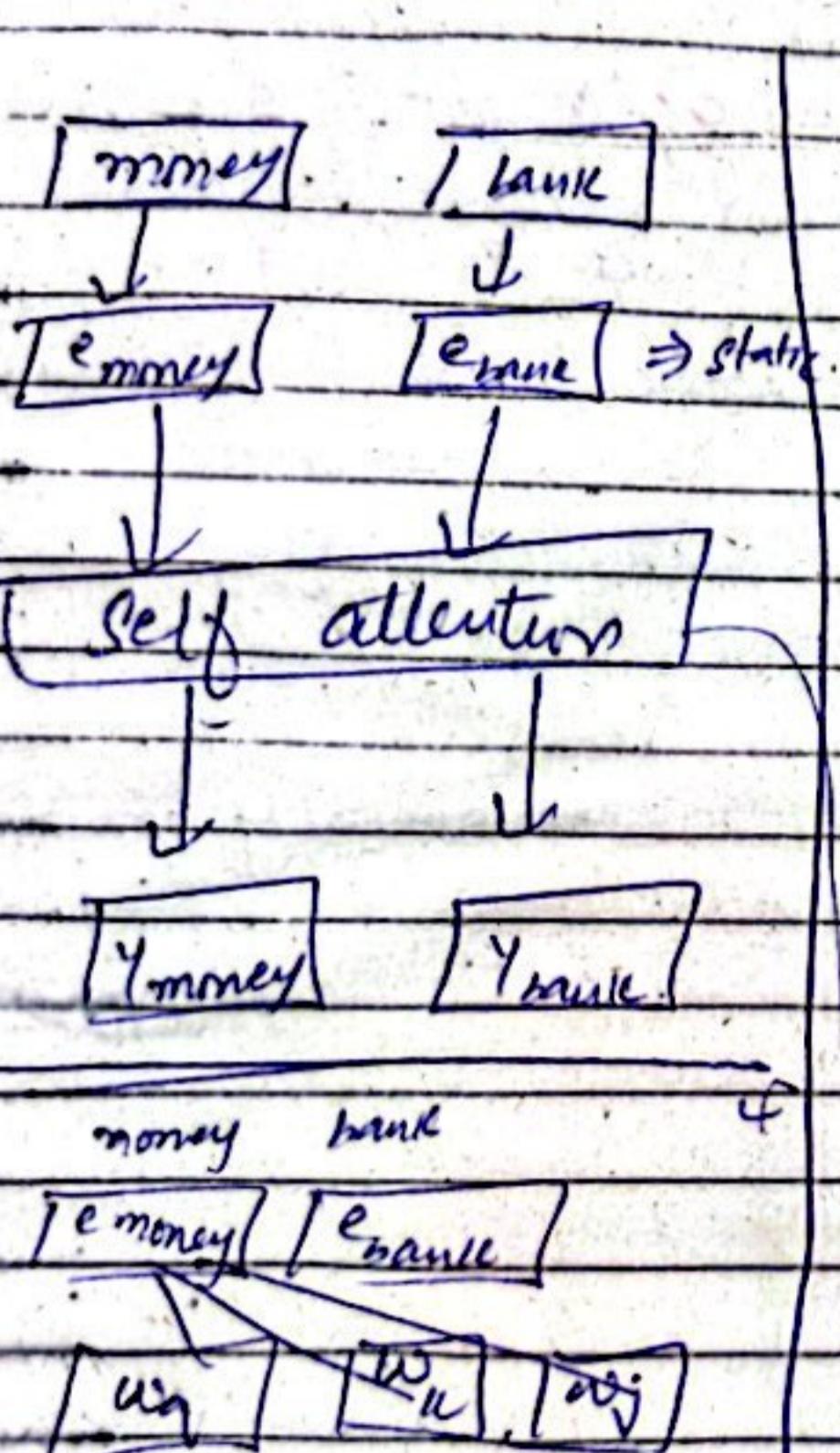
- multihead
- position encoding
- transformer architecture

Multi-head Attention :-

Recap of self-attention :-

Self-attention is a technique which you used to generate contextual embedding.

money bank $\xrightarrow{\text{represent}}$ number



$$\begin{cases} q_{\text{money}} & q_{\text{bank}} \\ k_{\text{money}} & k_{\text{bank}} \\ v_{\text{money}} & v_{\text{bank}} \end{cases}$$

- 1- one hot encoding
- 2- BOW
- 3- TF-IDF
- 4- Embedding
(semantic meaning capture)

money bank
 \rightarrow problem
 \Rightarrow static
 \hookrightarrow money bank
 \hookrightarrow give bank

both sentence bank
 embedding is same but
 its meaning is
 different in
 two sentence.

To solve that
 problem.

- 5- Self embedding
- \hookrightarrow generate
Contextual
 embedding
 (which context it came,
 it change)

Problem in self attention

The man saw the astronomer with a telescope

⇒ two perspective of sentence

↳ Tele^{scope} of astronomer is seen by man

↳ Telescope (seen by man) of astronomer is

This sentence is ambiguous.

but unfortunately self attention only

Capture one meaning/perspective

Self attention can't capture multiple perspective.

Working of self attention

↳ embedding of each word find

↳ query, key, value find of each word.

↳ then similarly check below query,
key, value vector

1) ~~possible meaning~~

With similarity score

man , telescope ,

—

caus , "

—

2) ~~same meaning~~

astronomer , telescope

✓

saw , astronomer

✓

doc summarization

→ If that paragraph depends on self
attention. it generate one
perspective summary.

We want multiple perspective summary.

similarity score

The man saw the officer with a telescope.

The
man
saw
the
officer
with
a
telescope

Self attention only make one table.
multiple meaning/perspective can't capture.

To solve that problem multi-head attention came.

Multi-head attention

→ multiple self-attention used are Multi-head attention.

money bank

↓ ↓

e_{money} e_{bank}

↓ ↓

self attention

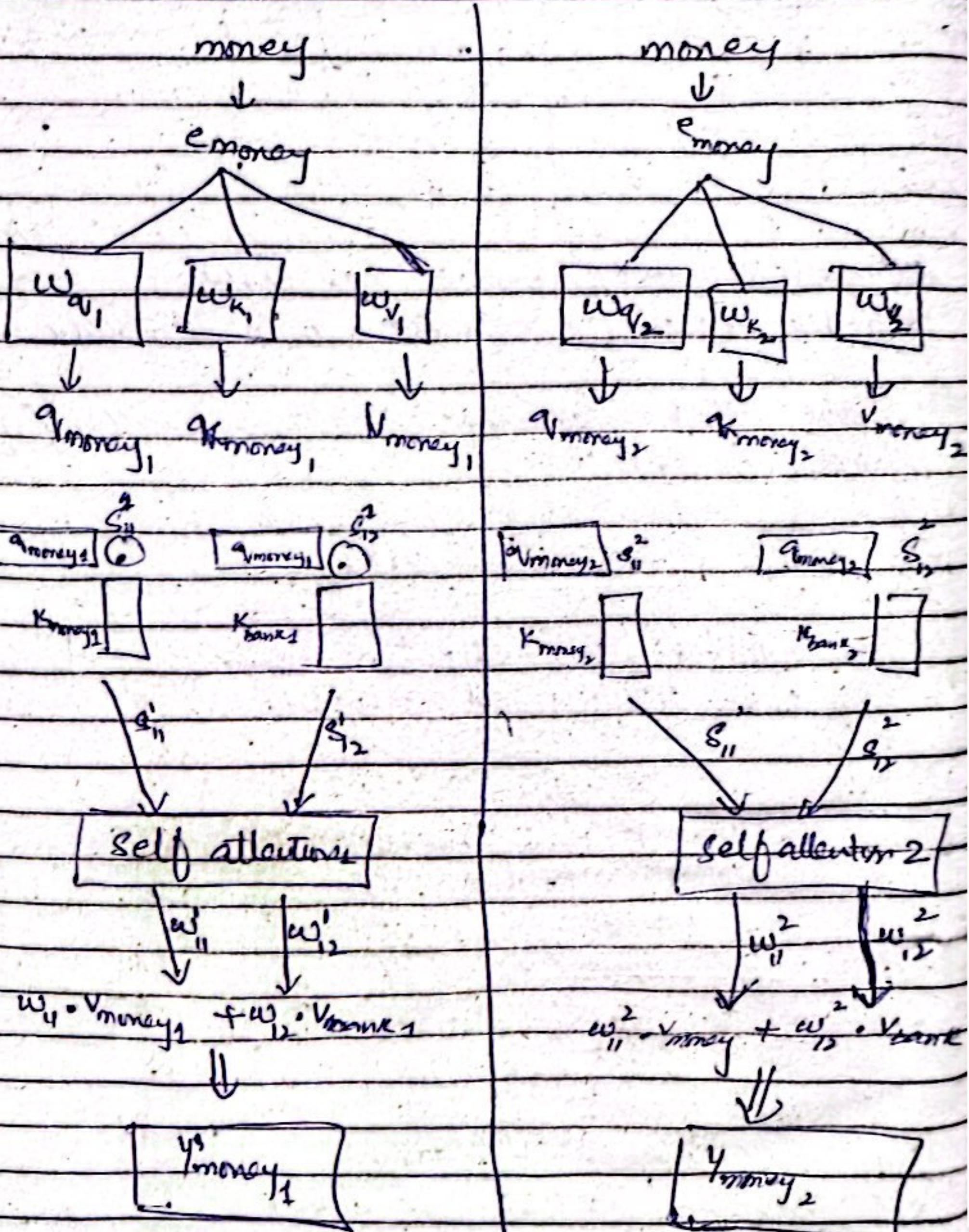
self attention 2

$y_{\text{money}3}$
 $y_{\text{bank}3}$

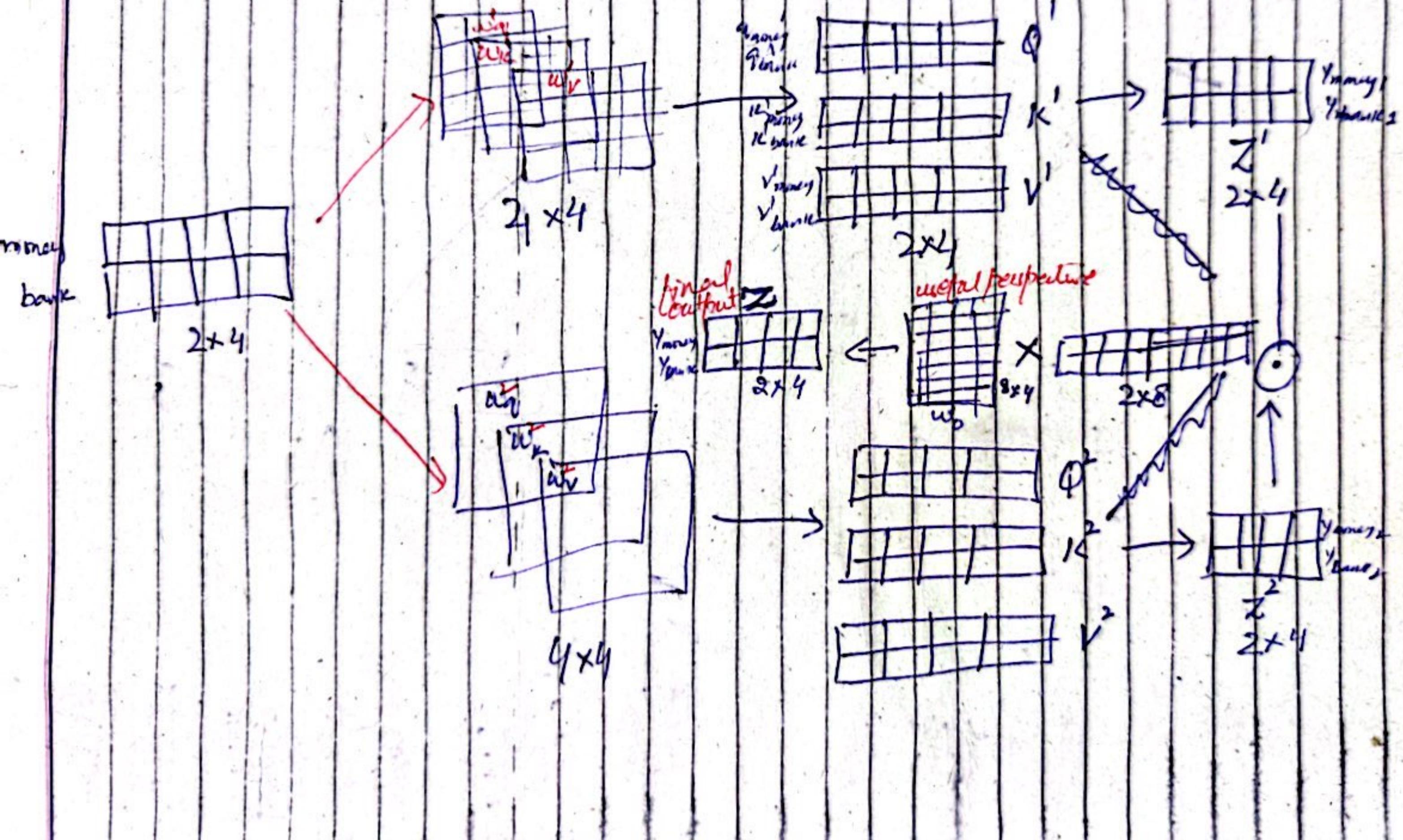
$y_{\text{money}2}$
 $y_{\text{bank}2}$

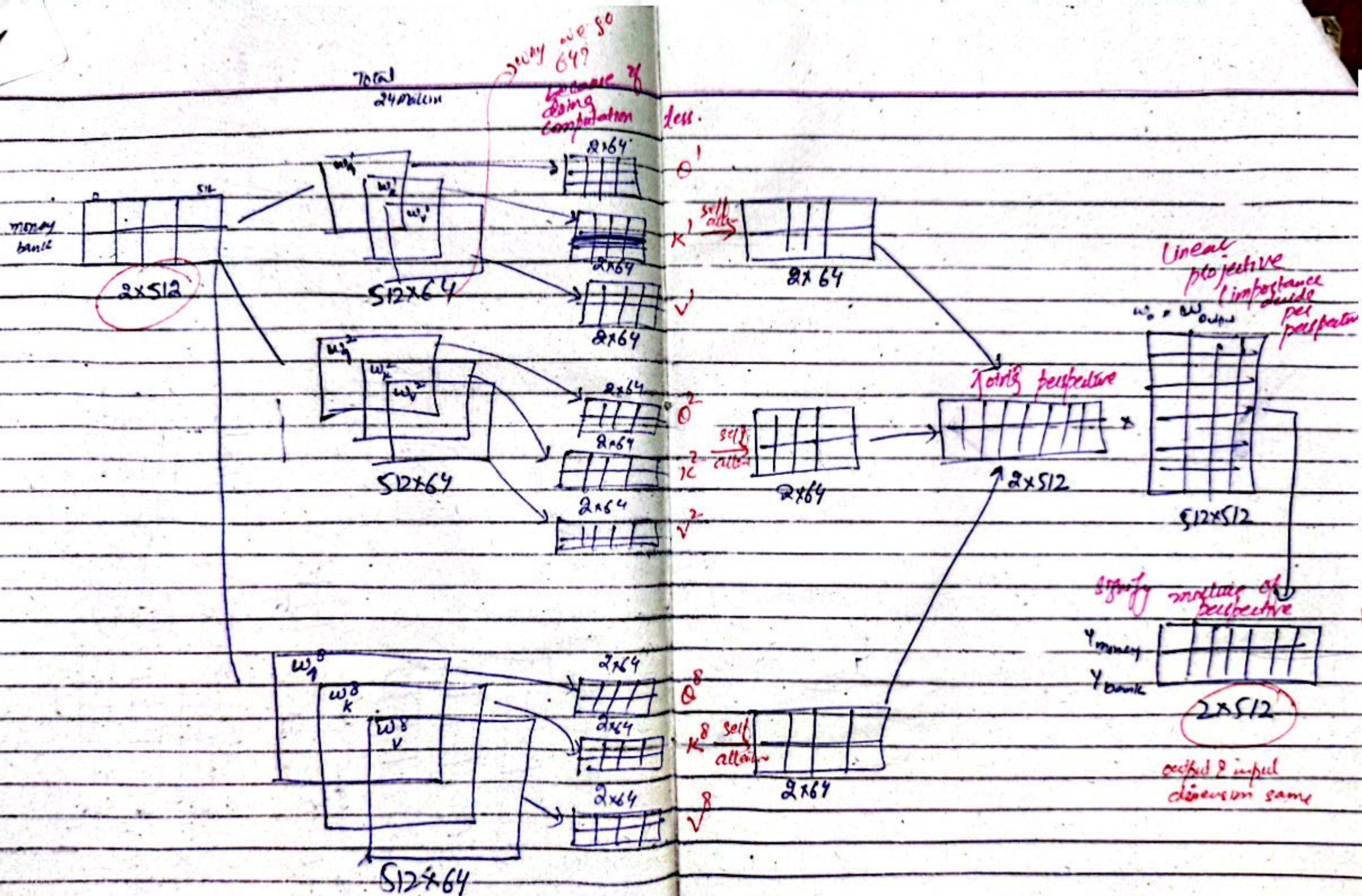
It parallelly do, capture different perspective.

sec. - money bank



How multihead apply same time





overall configuration is same as
→ single one attention used with 52
dimension.

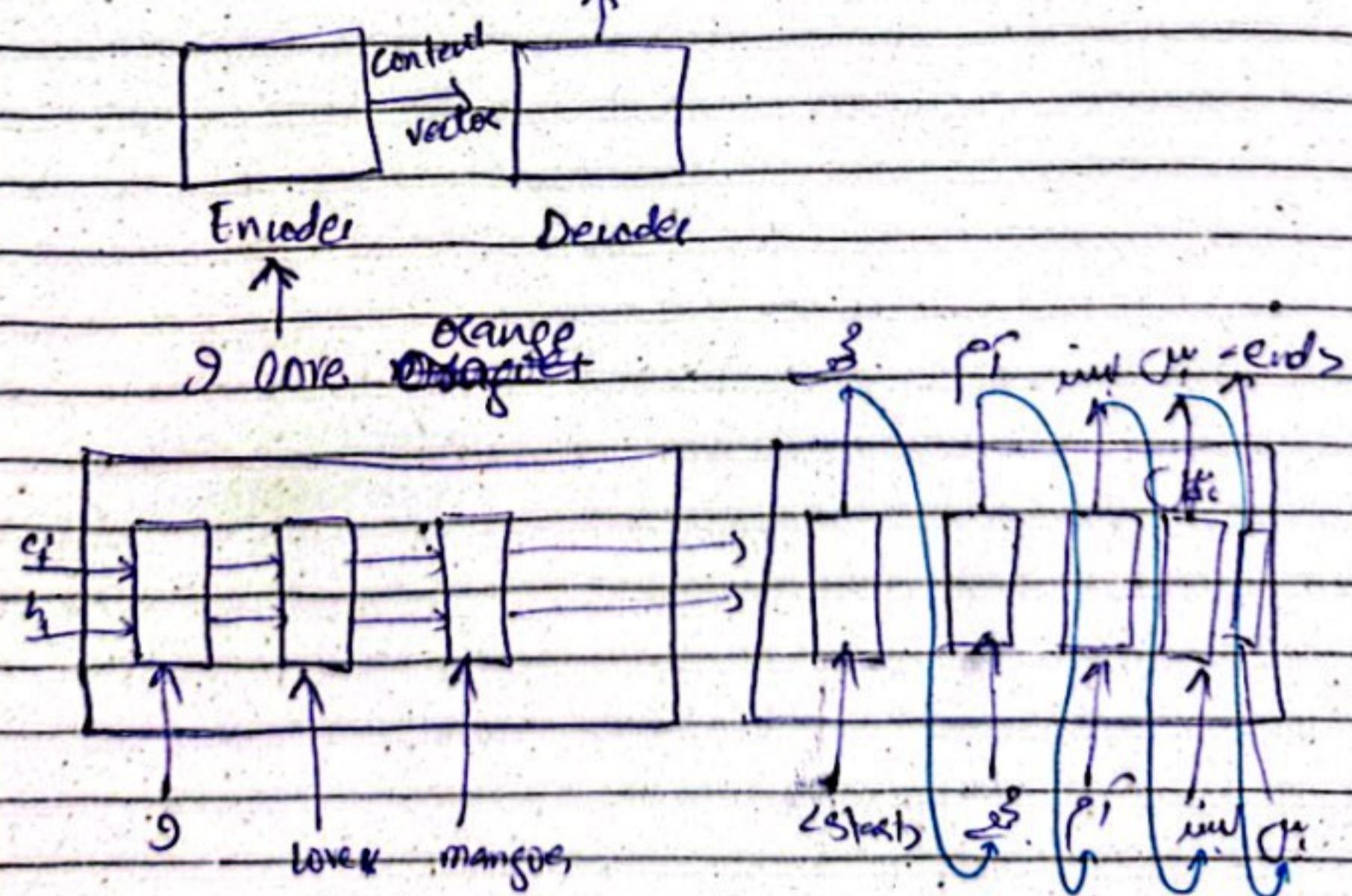
↳ multiple perspective come

Masked Multi-head Attention

The transformer decoder is autoregressive at inference time (prediction / testing) and non-autoregressive at training time.

Autoregressive \Rightarrow

next prediction based on previous sequence
 $i, i+1, \dots, i+k$.



time step what output.

at any time-step, if you predict output you need previous

You called that model called autoregressive.

Why seq to seq are autoregressive.

When we predict next word you need to know what is predict before because next word based on previous words.

That's why we can't do that thing in parallel.

→ Transformer in training time not autoregressive.

If in training we used teacher forcing, we don't give wrong word in training.

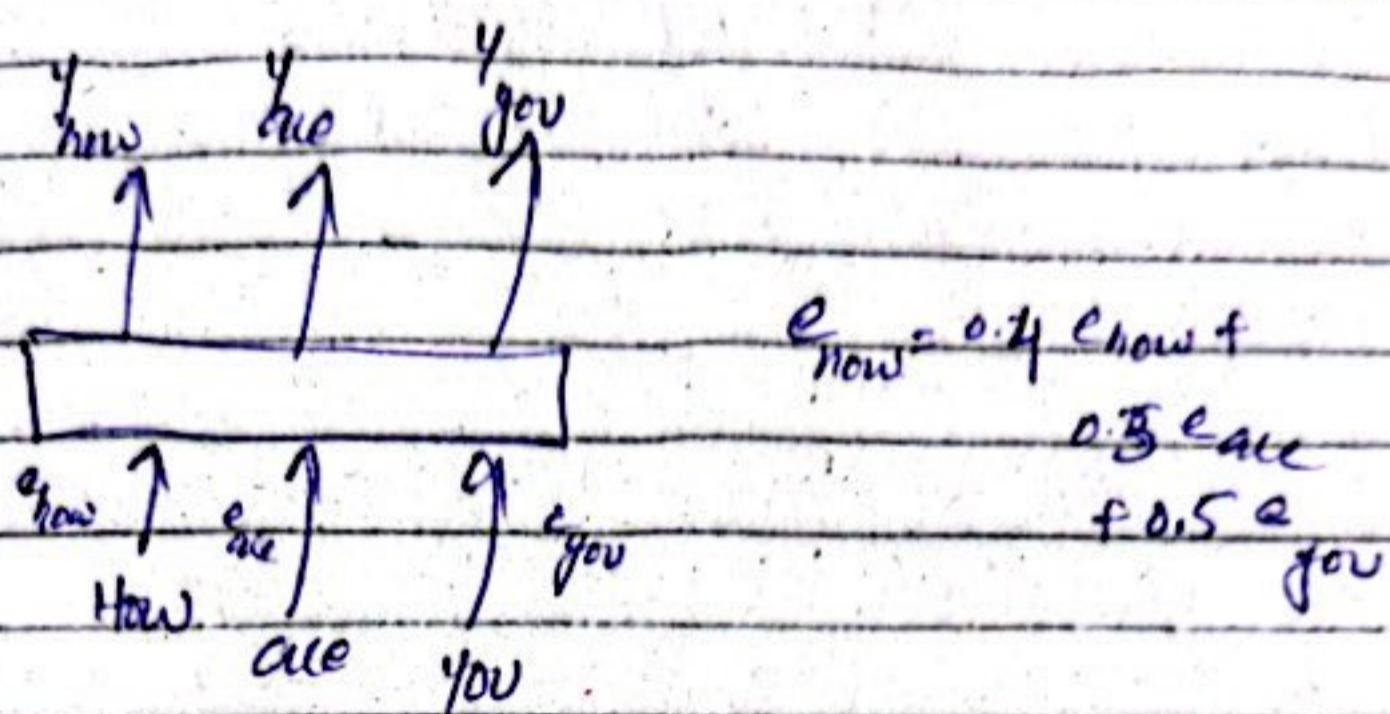
In training, we know the output of sentence. Why we^{not} only give the output ~~time~~ of new word from dataset without seeing what is produced.

→ Is there any drawback of making autoregressive? (training)

If 300 words, it means 301 times separate run decoder.

→ You do that thing in parallel way also during training because you know the next word from dataset.

Soft attention \rightarrow next word contextual embeddings
you know the next word.



in telling e will make from future prediction word value and itself also.

\Rightarrow you don't know the next word how to solve that eg.

$$e_{now} = 0.2 e_{now} + 0.1 e_{are} + 0.7 e_{you}$$

$$e_{are} = 0.3 e_{now} + 0.2 e_{are} + 0.5 e_{you}$$

$$e_{you} = 0.7 e_{now} + 0.3 e_{are} + 0.1 e_{you}$$

Autoregressive

- \rightarrow no date leakage
- \rightarrow slow

Non-Autoregressive

- \rightarrow fast
- \rightarrow date leakage

Knows
Grace
you

Knows
Grace
you

V_{now}
V_{ace}
V_{you}

O

K

V

dot product

mapped

$$\begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix}$$

Scale

$$\sqrt{d_k}$$

$$\begin{bmatrix} S'_{11} & S'_{12} & S'_{13} \\ S'_{21} & S'_{22} & S'_{23} \\ S'_{31} & S'_{32} & S'_{33} \end{bmatrix}$$

$$U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

softmax

softmax

Knows
Grace
you

V_{now}
V_{ace}
V_{you}

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$y_{\text{now}} = w_{11} V_{\text{now}} + [w_{12} V_{\text{ace}} + w_{13} V_{\text{you}}]$$

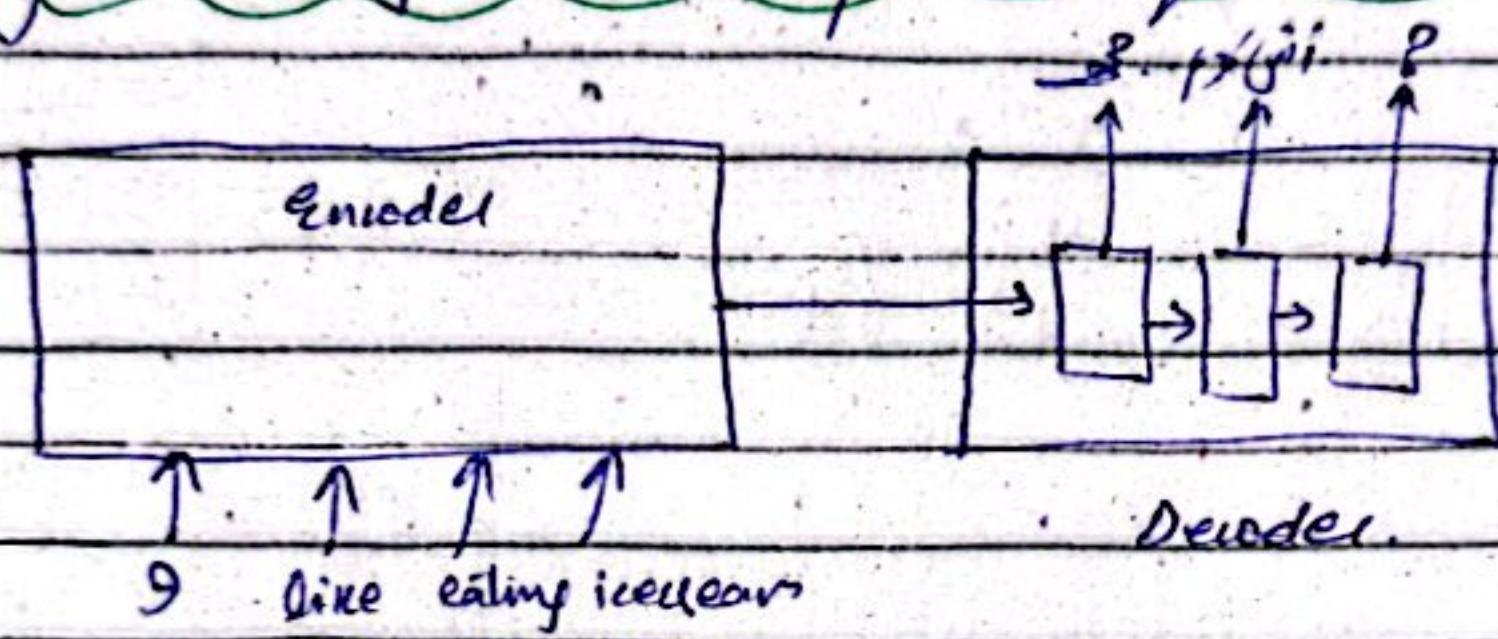
$$y_{\text{ace}} = w_{21} V_{\text{now}} + w_{22} V_{\text{ace}} + w_{23} V_{\text{you}}$$

$$y_{\text{you}} = w_{31} V_{\text{now}} + w_{32} V_{\text{ace}} + w_{33} V_{\text{you}}$$

Cross Attention

Cross Attention is a mechanism used in transformer architectures, particularly in tasks involving sequence-to-sequence able like translation or summarization.

It allows a model to focus on different parts of an input sequence when generating an output sequence.



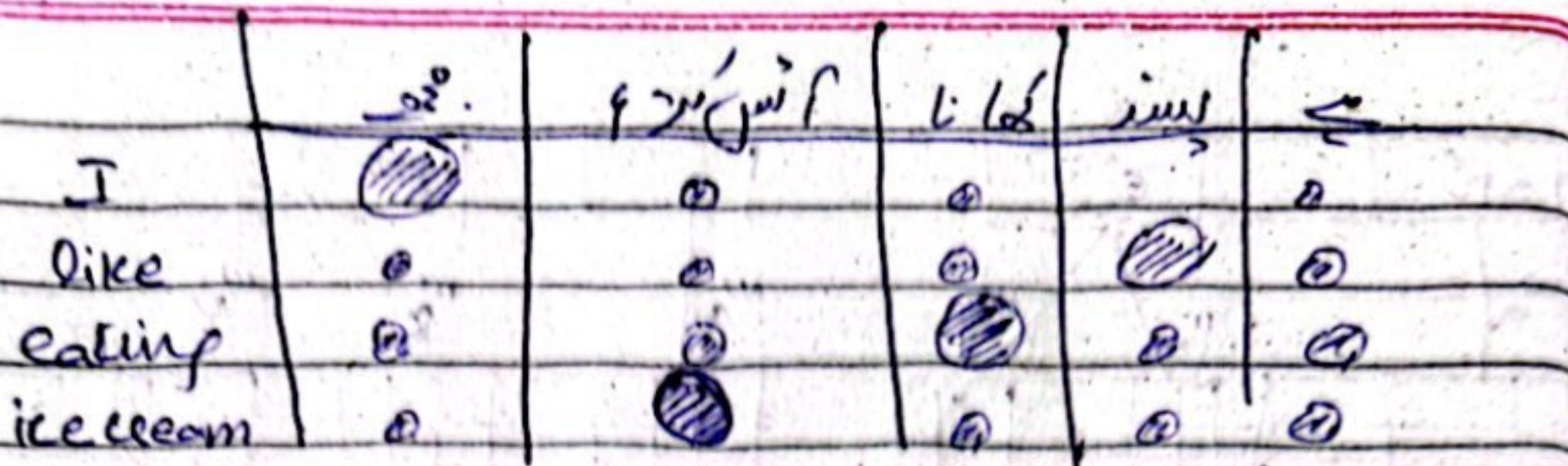
Encoder \rightarrow To understand sentence in very good way
that summary pass to decoder.

Decoder \rightarrow decoder get context vector and start predicting/translation
 \rightarrow seq to seq it predict/constraint

Decoder time 3 \Rightarrow what next word

Self attention (same sentence context generate) \rightarrow decoder what generate until till now
what is relationship b/w input and output sentence

Encoder information



different
two sequence what relationship

↙
cross self attention

same sequence what relationship

↙
self attention

Cross attention is conceptually similar
to self attention

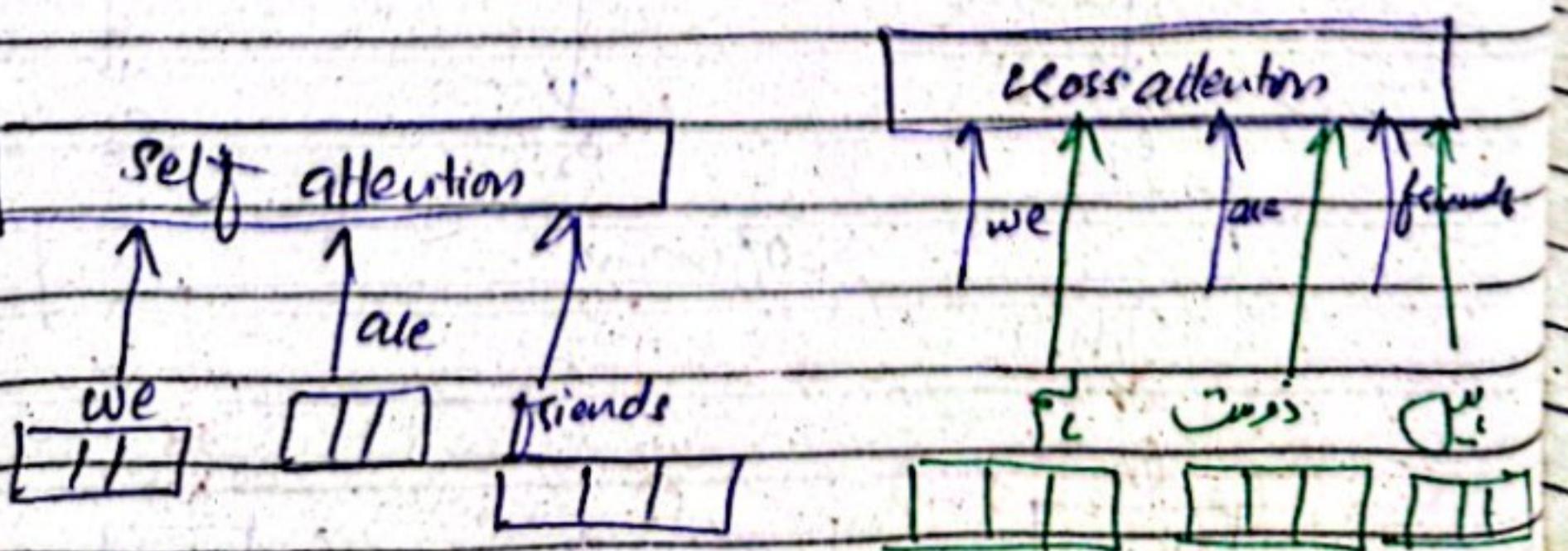
- input
- processing
- output

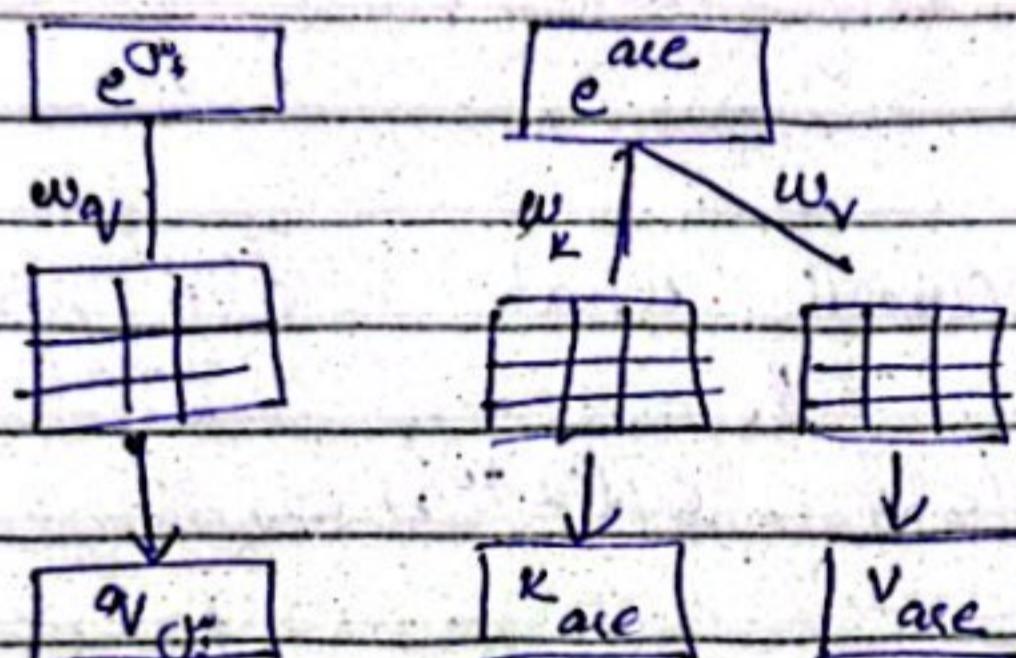
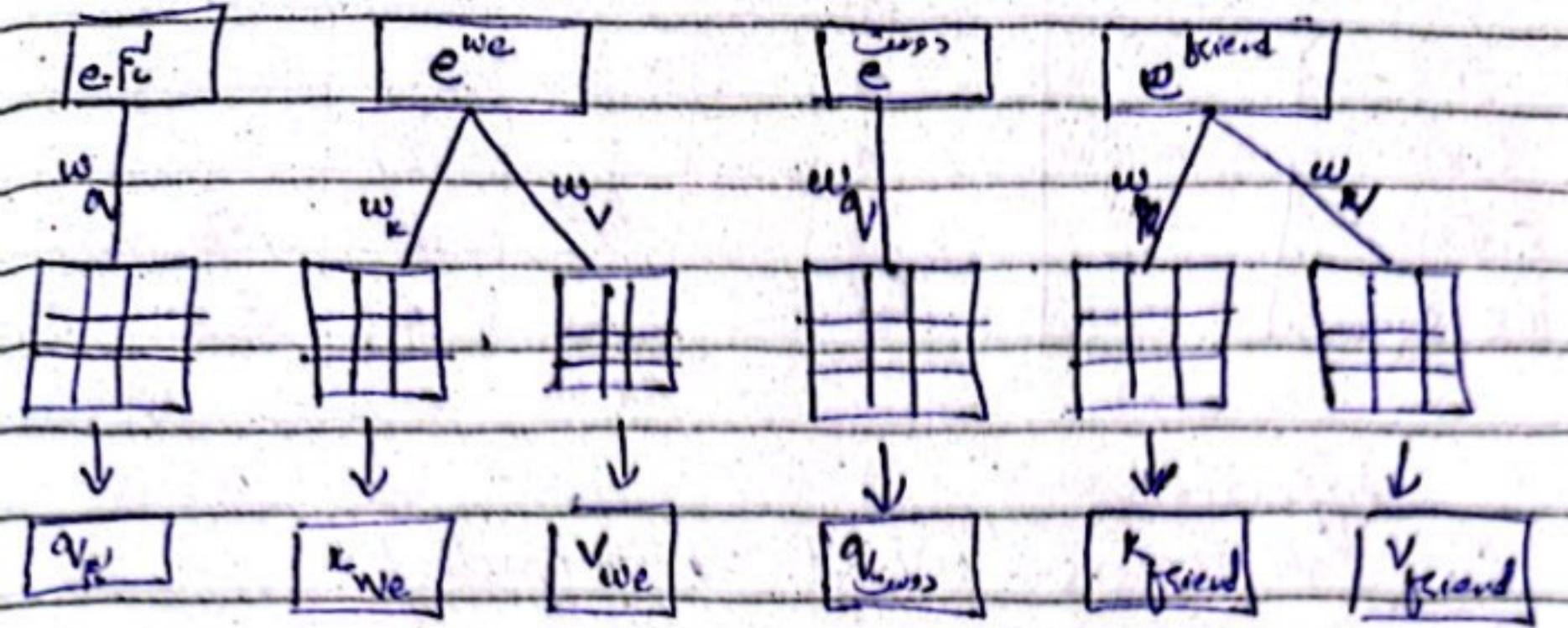
self

cross

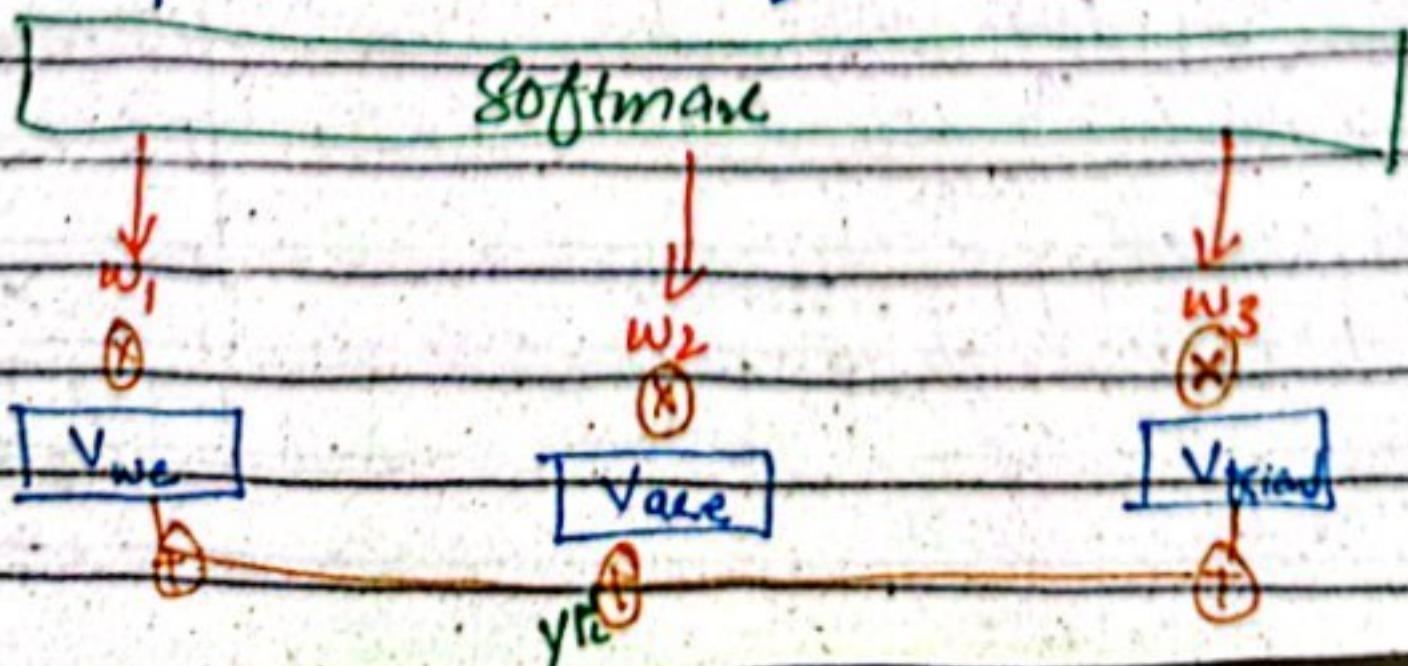
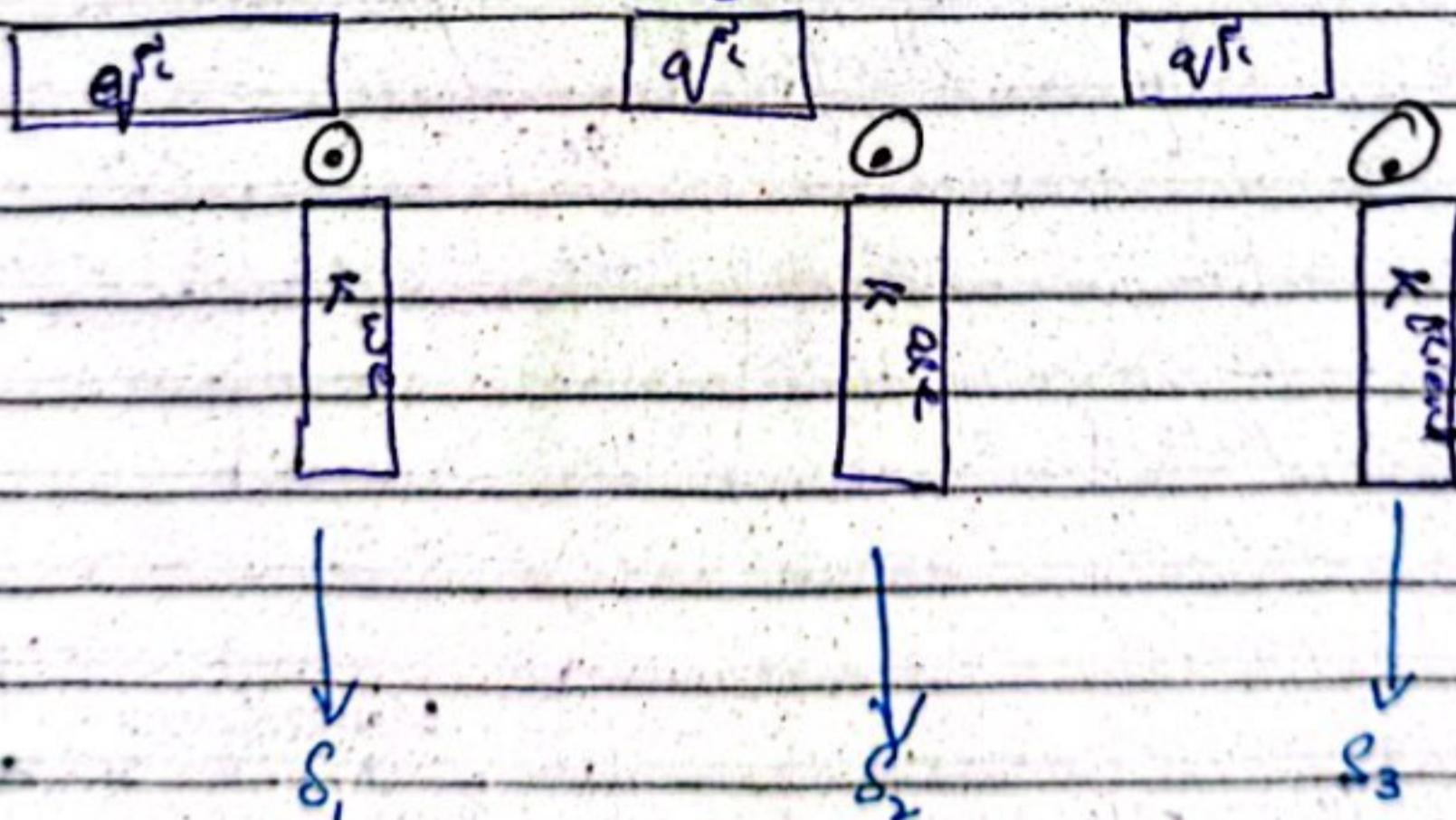
single sequence embedding

two different sequence embeddings



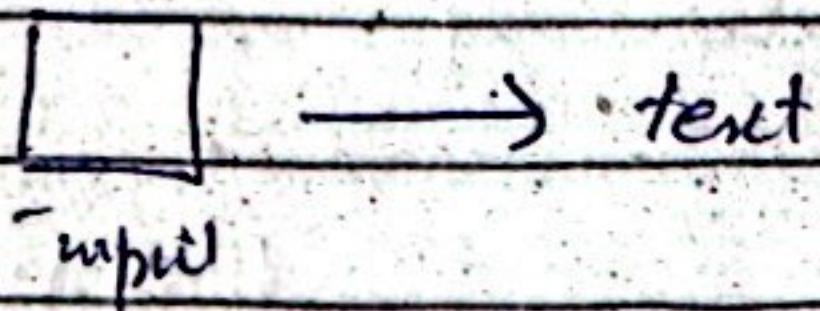


e^{f_i}



Use cases

image captioning



text → image

text → speech