

Vector Embeddings

and

RAG

≈ 9:30 PM

- ① Intro to ML
- ② Languages and ML
- ③ How to create vector embeddings

↳ Word2Vec

↳ BERT

(Transformer Architecture)

- ④ How LMs work

- ⑤ RAG

- ⑥ Types of RAG

→ Text

- ⑦ Evaluation and adv topics

→ Embedding

- ⑧ Libraries

- ⑨ Retrieval Algos.

Intro to ML

Traditional Programming

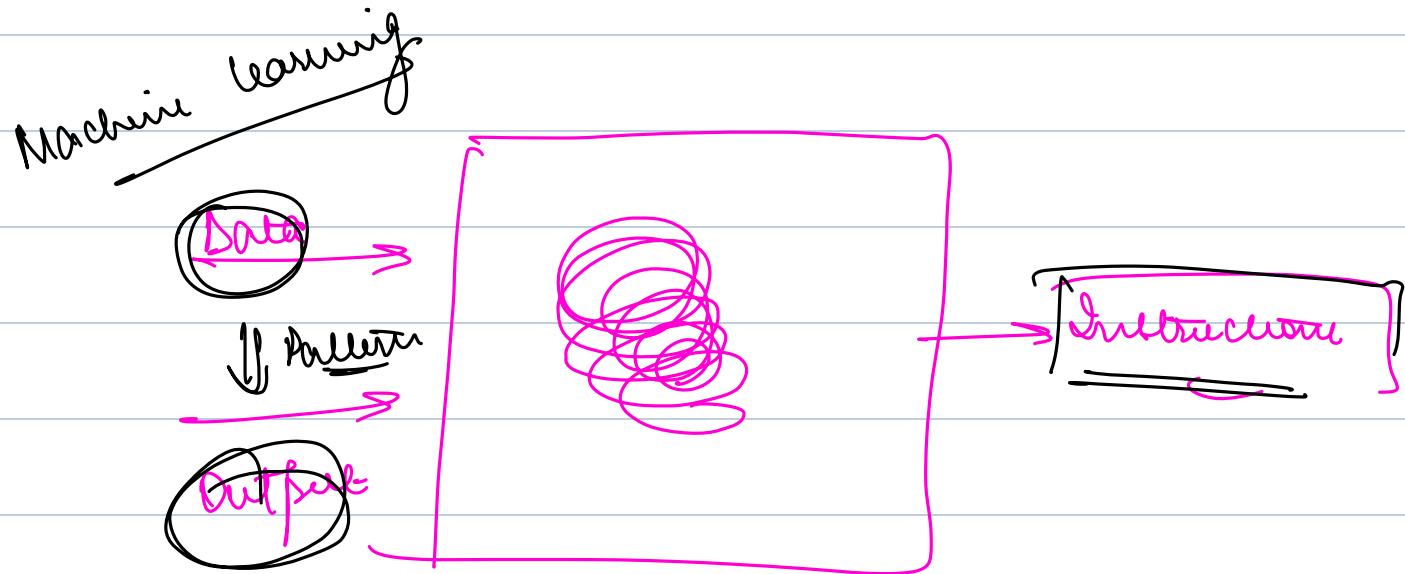
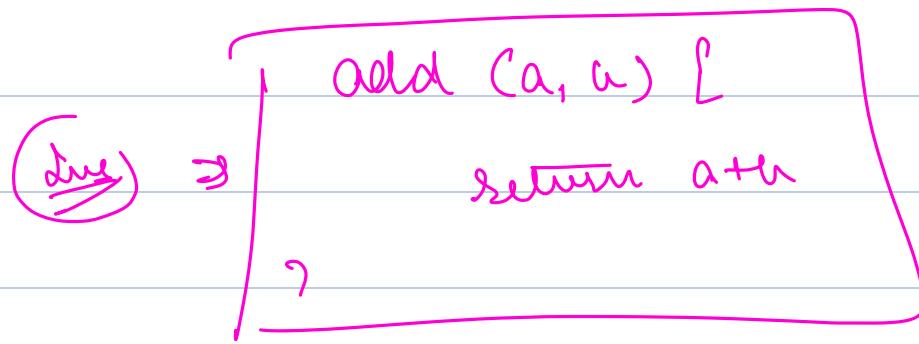
Instructions

+ Data

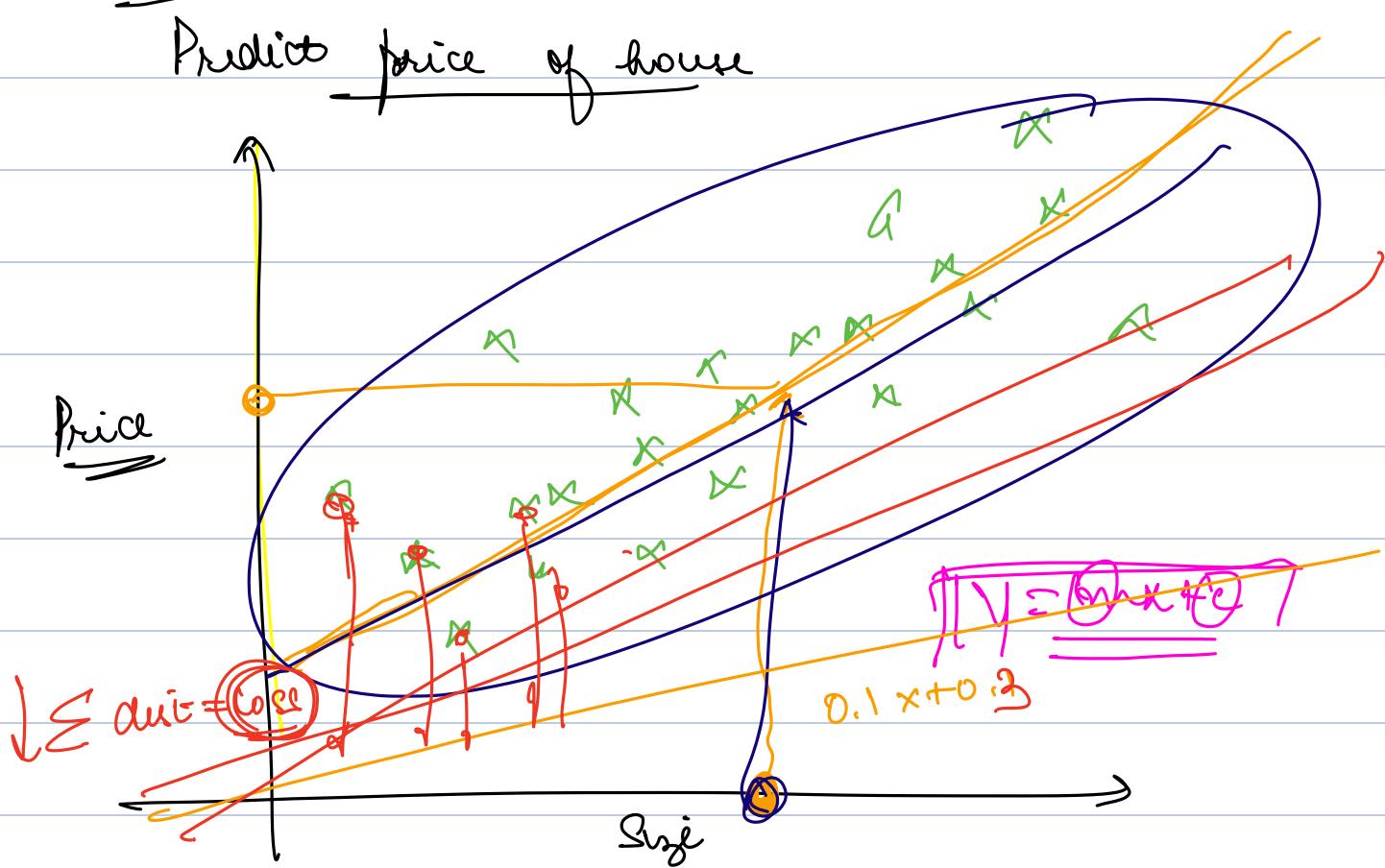
Comp

Output
(q)

2, 7



But How?



Assumption behind

ML

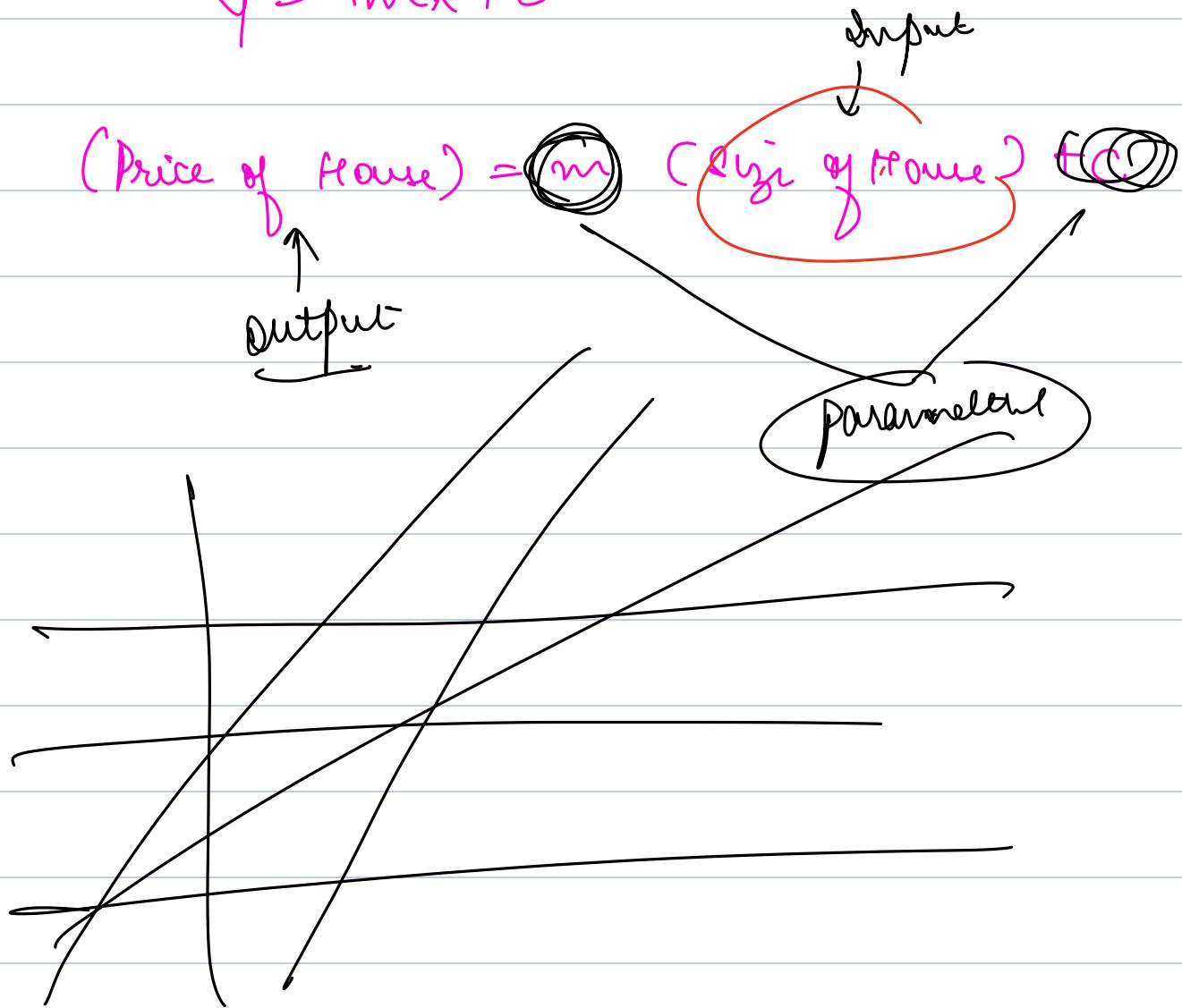
: the world works on

'defined pattern'

we just have to figure the pattern

Assume the given one can be rep
by a linear eqⁿ in x frame

$$y = mx + c$$



$$y = \alpha x_1 + \beta x_2 + c x_3 + \dots$$

assume an eqⁿ

$$y = \textcircled{m} x + \textcircled{c}$$

start by letting random value of
these param.

$$\Rightarrow \boxed{y = 0.1 x + 0.1}$$

- Iterate : \rightarrow till loss is minimized.
 \rightarrow calculate how far is your current
pattern from real data
 \rightarrow based on that your change variable

Gradient

Descent

$$\text{price} = \alpha_1 \cdot \text{Size} + \alpha_2 \cdot \# \text{bedrooms}$$

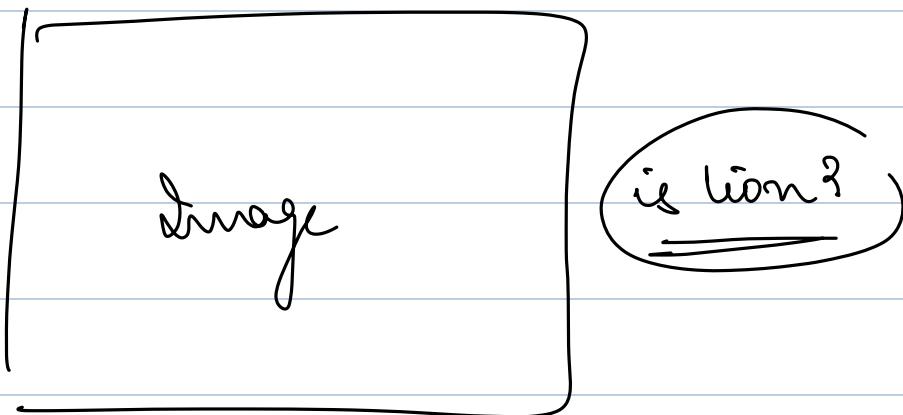
~~$\alpha_3 \cdot \# \text{rooms}$~~

Inputs \rightarrow Numbers

But! !

Every kind of input for a machine prediction task may not be numeric.

Eg: ~~Image~~ = Image

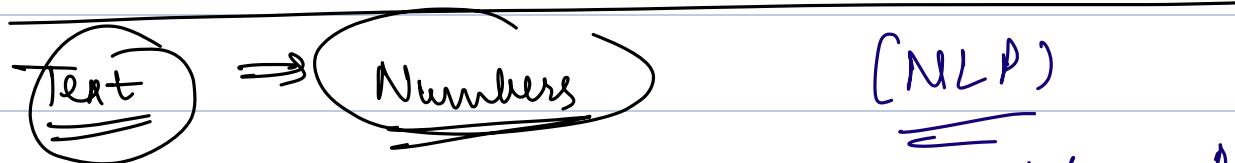
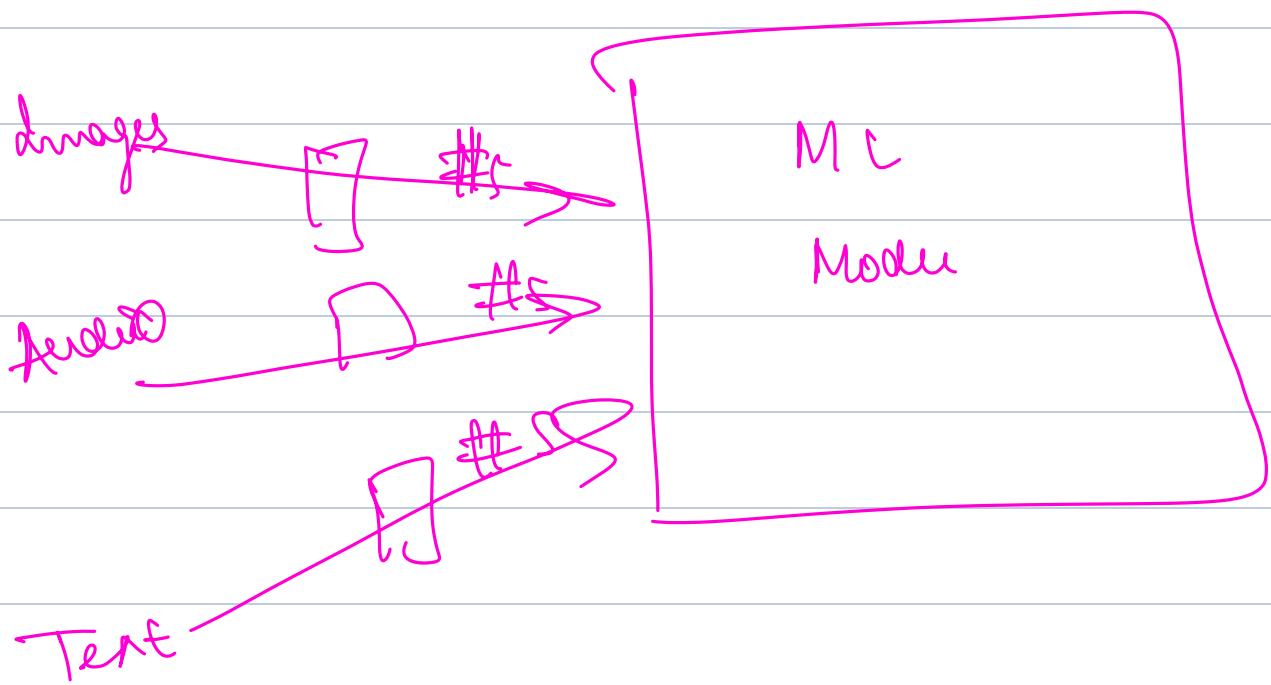


Input:
Hi How

\Rightarrow i | are | where

ML \Rightarrow





- ① Give every word a pos

Vocab \Rightarrow Set of all words that I know about

(dictionary)

\Rightarrow give rank to every word.

GPT3 \Rightarrow (52227 words)

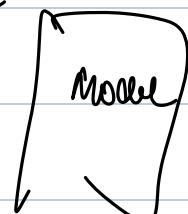
1 \rightarrow a

2 \rightarrow another

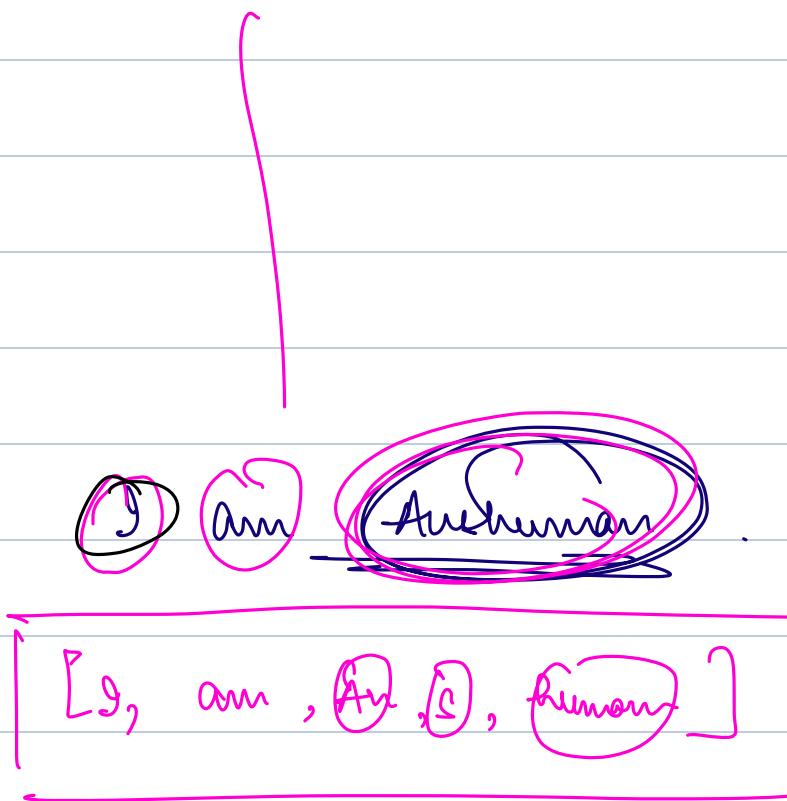
3 \rightarrow an

4 \rightarrow apple

5 have a apple
= [37, 29, 1, 4]



5 → ~~arrt~~



Token by Token

(1), 2 (2)

Input
↓
Tokens
↓
numeric rep for each

~~(200)~~
215

But!!

There is no logical reasoning behind the #s I have kept each word as.

(2).

(3).

a. very never learn at all.

→ Our rep of words should have some meaning associated to them

(eg)

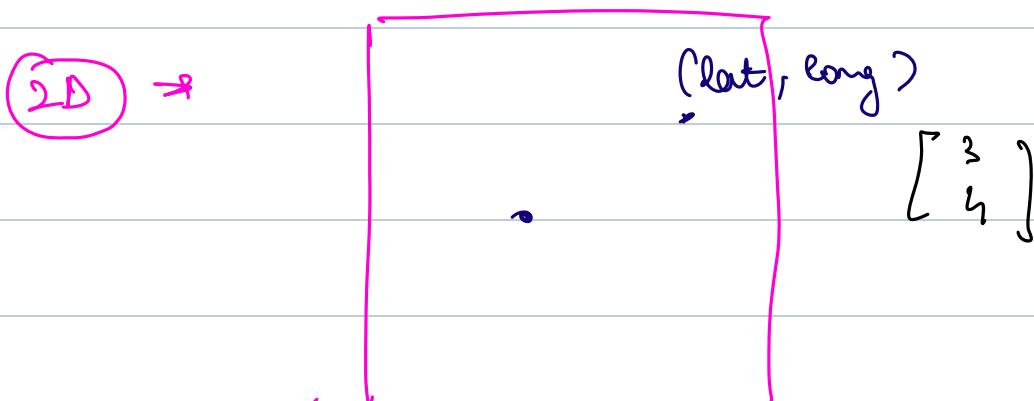
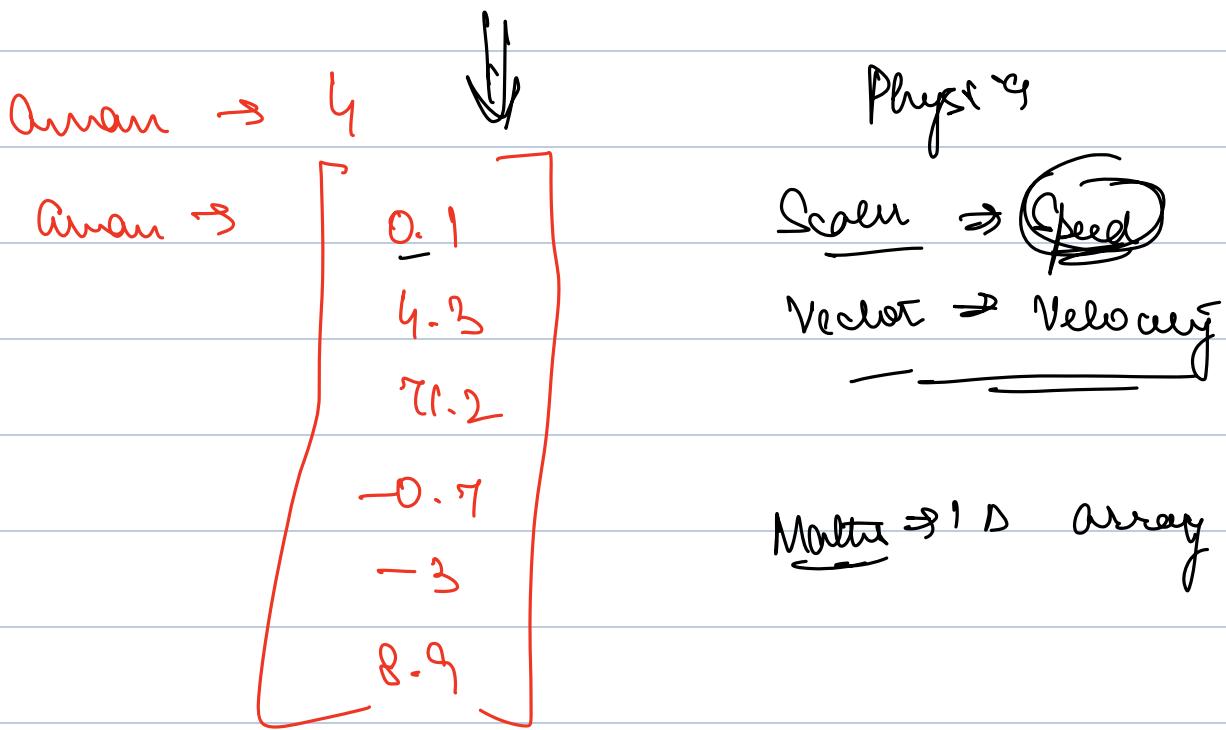
$$\text{Rep modi} - \text{Rep nida} = \text{Rep}_{\text{Meloni}} - \text{Rep}_{\text{Italy}}$$

$$\text{Rep King} - \text{Rep Queen} = \text{Rep Men} - \text{Rep women}$$

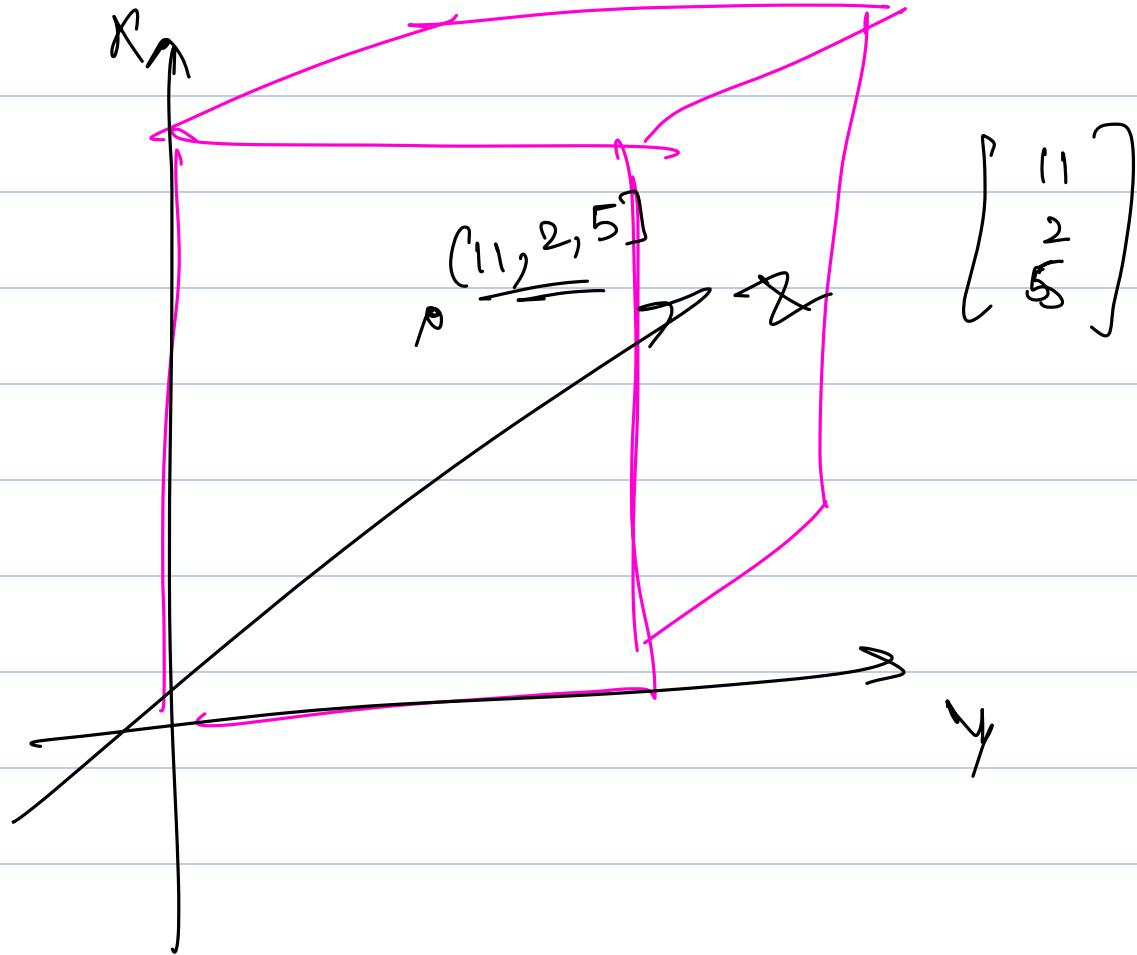
Am I do this?
Yes

Vector Embeddings

In vector embeddings, for each word in my dictionary, rather than figuring a numeric value, I try to figure a vector rep



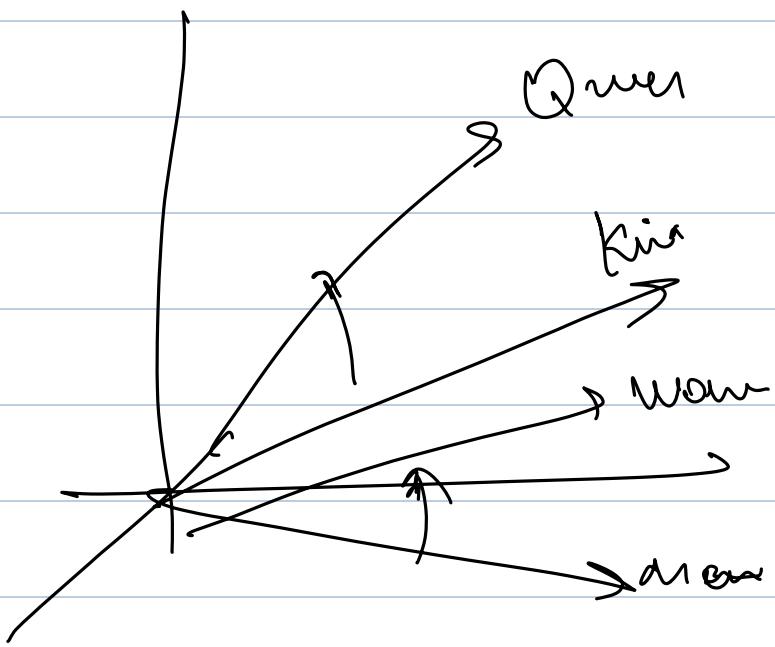
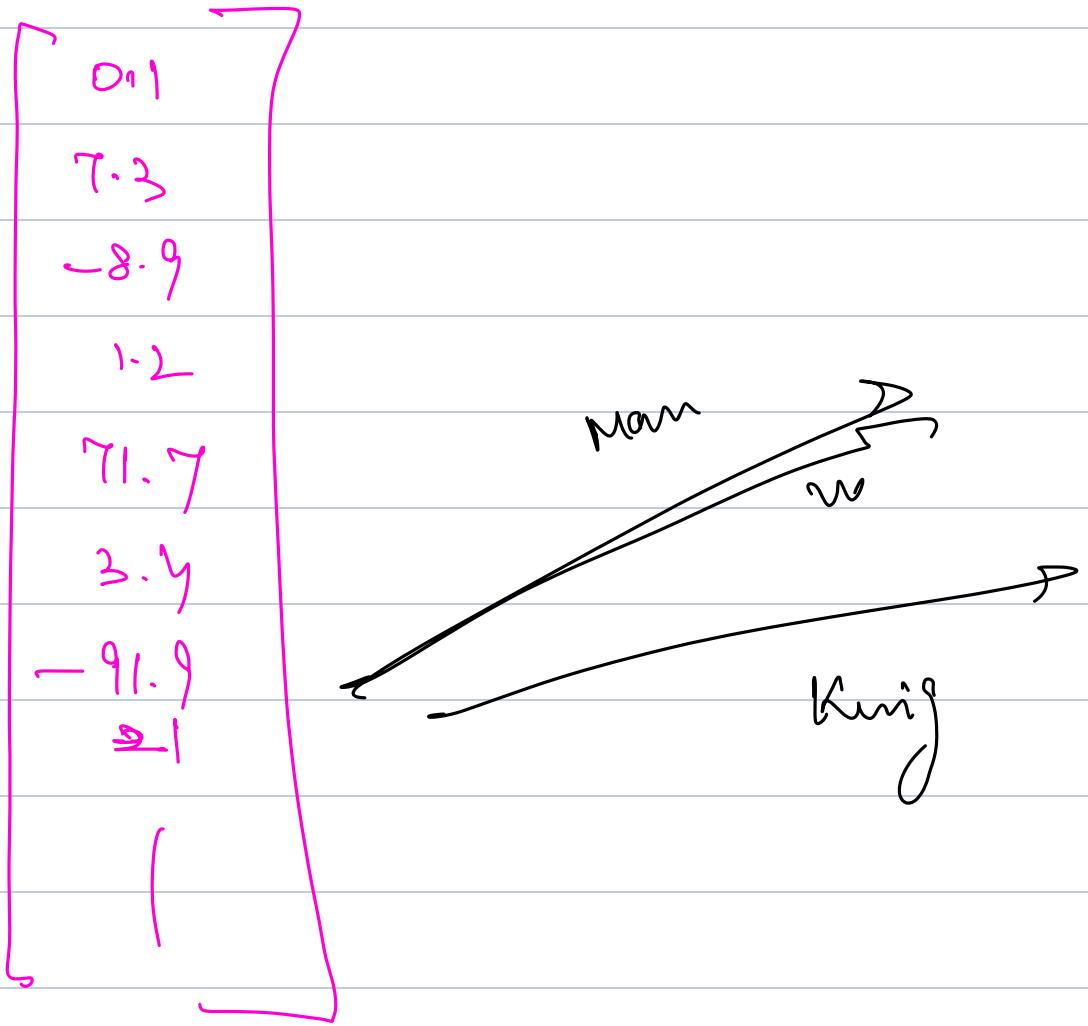
(3D)



$$\begin{bmatrix} 0.1 \\ 7 \\ -3 \\ 2 \\ 9 \end{bmatrix}$$

GPT 3

12096 dimension



Vector Embedding: Numeric Rep of words
in my dictionary

How? Cosine Similarity

$$= \frac{\text{distance } \vec{v}_1 - \vec{v}_2}{\|\vec{v}_1\| \|\vec{v}_2\|}$$

men

boy

S1

= There was a bear flying around

-

S2

+ liger was flying near my

{ if 2 words have similar words around them, probably they are similar.

① Word2Vec

↳ CBOW (Continuous Bag of Words)

↳ Skip Grams

CBOW ~~of~~ faster

whole wikipedia

۷

There was a bird flying near me

The was flying a plane yesterday

CBOW (

There was a plane

Predict ?

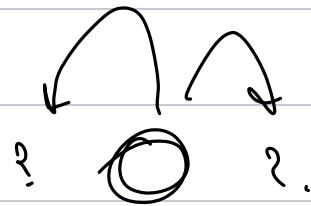
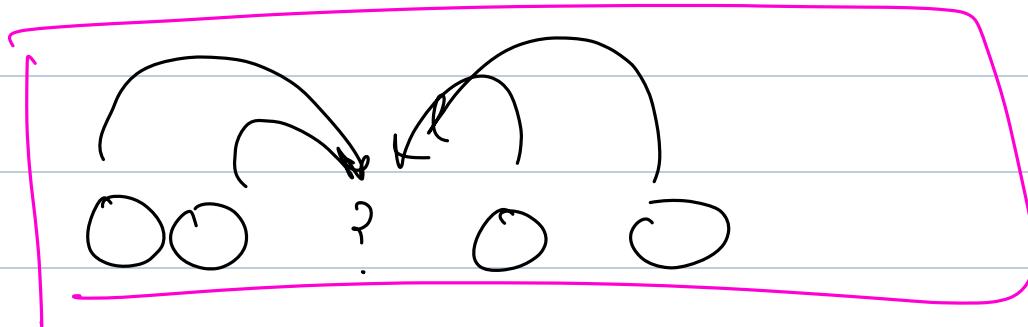
~~flynn yesterday Pigeon~~

Input
[the, who, - , -]

op

Skip Gram

predict the words around n



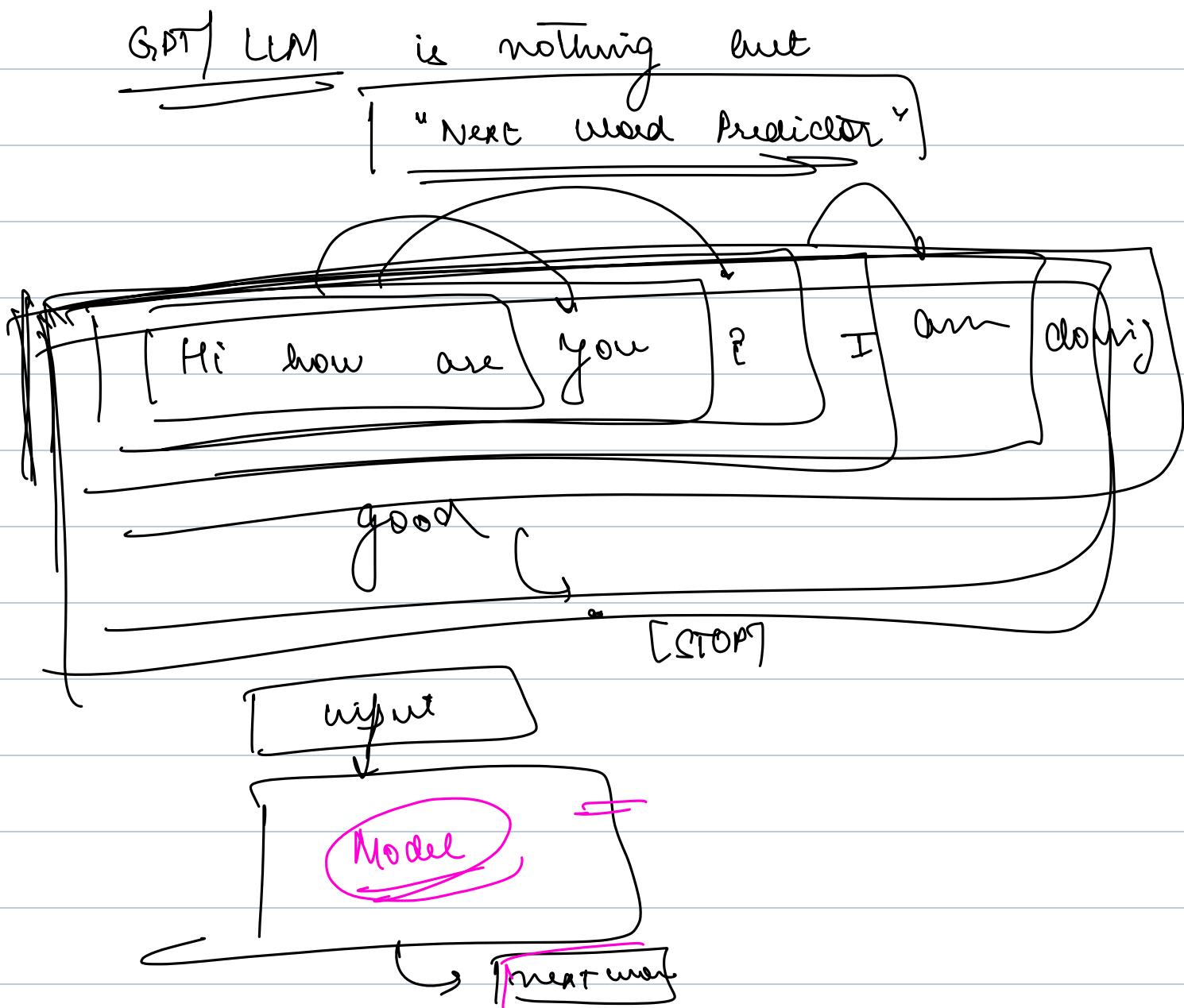
BERT

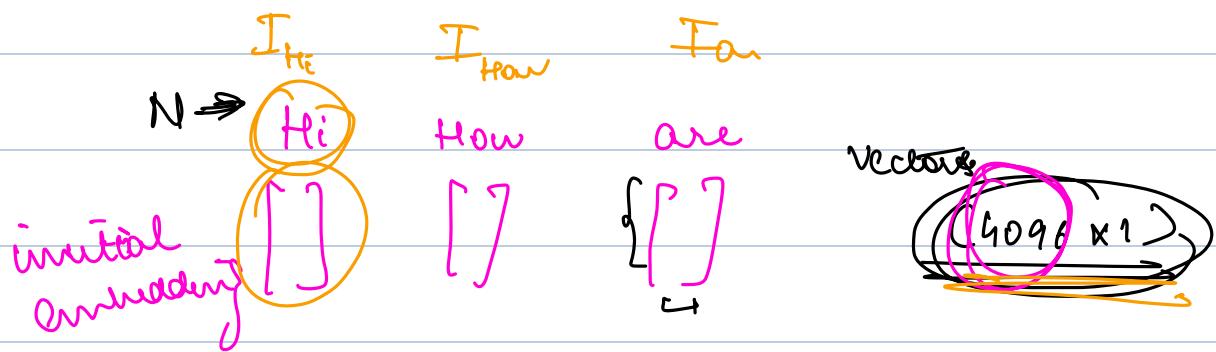
-

Intro to Transformer Architecture

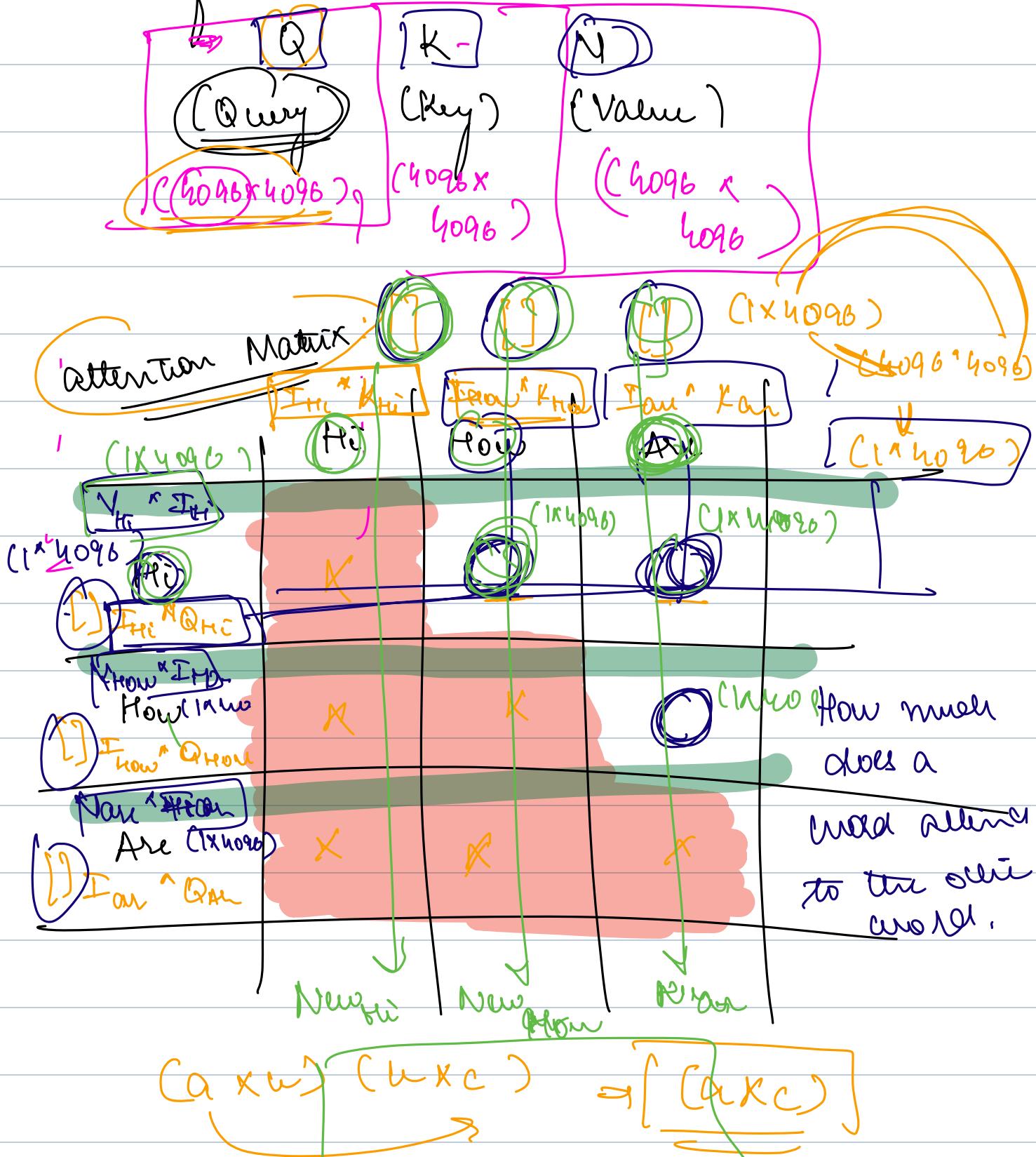
Assume an initial embedding for each token

(maybe start by embedding generated by word2vec)





In Transformer we learn 3 matrices:



There was a bee 



⇒ each preceding word may
impact the meaning of future
word

The was eating Apple

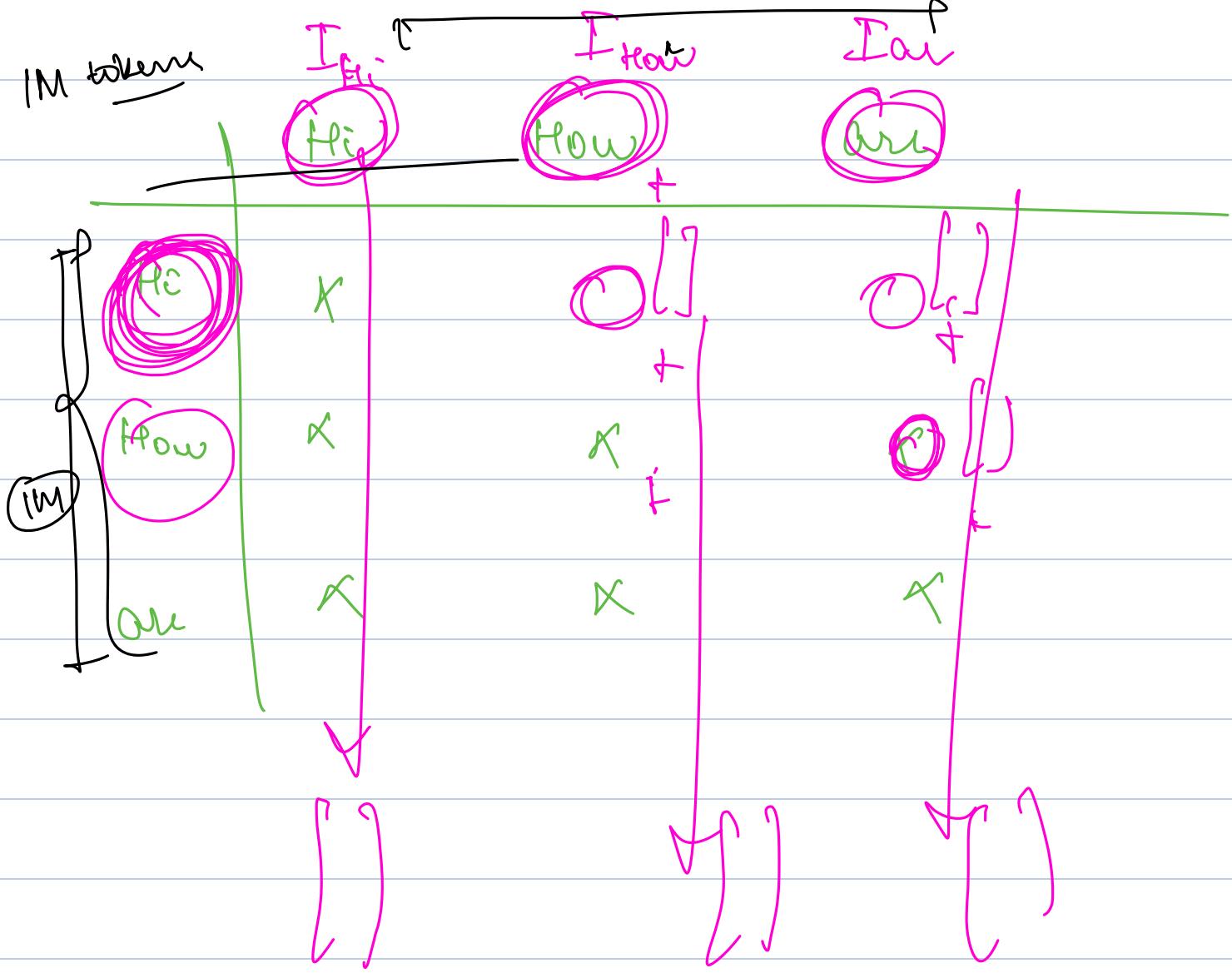
He was writing an Apple iphon

$$(1 \times 4096) \quad (1 \times 4096) =$$

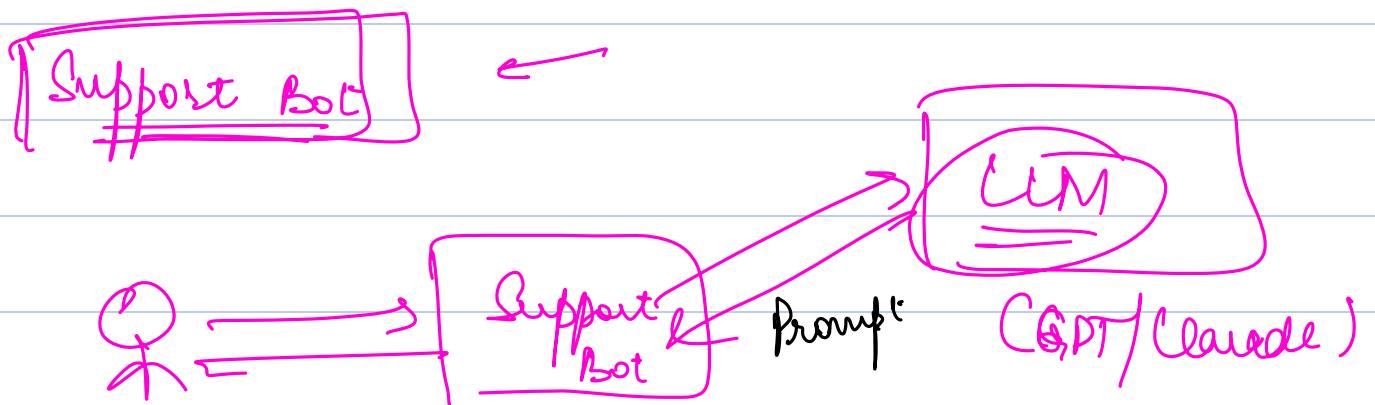
↓

$$\underline{(4096 \times 1)}$$

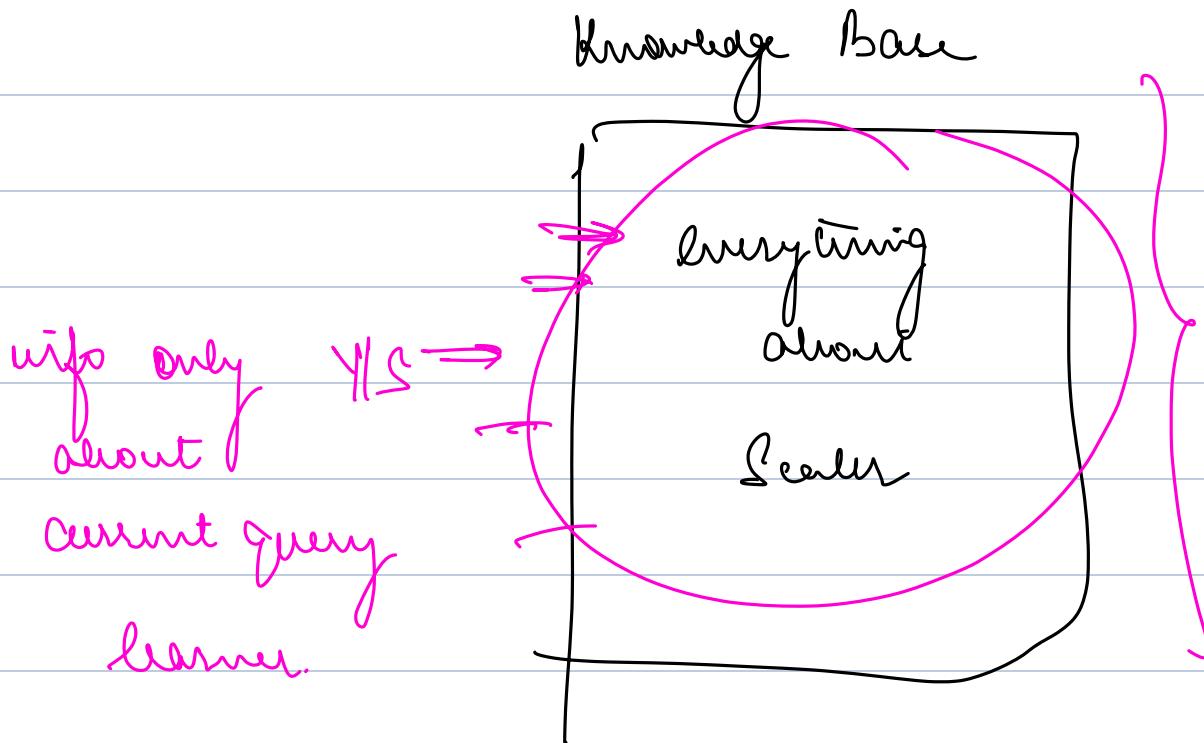
↓



RAG (Retrieval Augmented Generation)



You are — —



every LM \Rightarrow context window

\Rightarrow How much of input can it handle?

$$\underline{10^6 \times 10^6}$$

$$\rightarrow 10^{12}$$

$$4 \times 10^{12} \text{ By } \rightarrow 4 \text{ TB}$$

~~problem!~~

longer the prompt:

- ① More Costs
- ② Quality will also reduce.

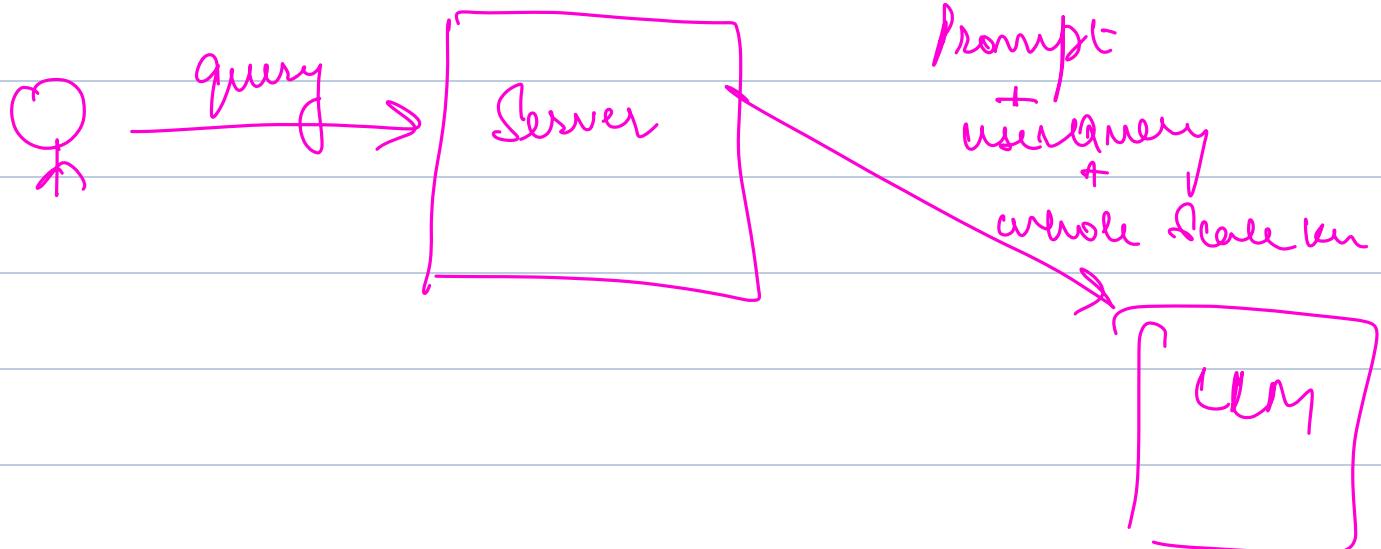
Use Case 1

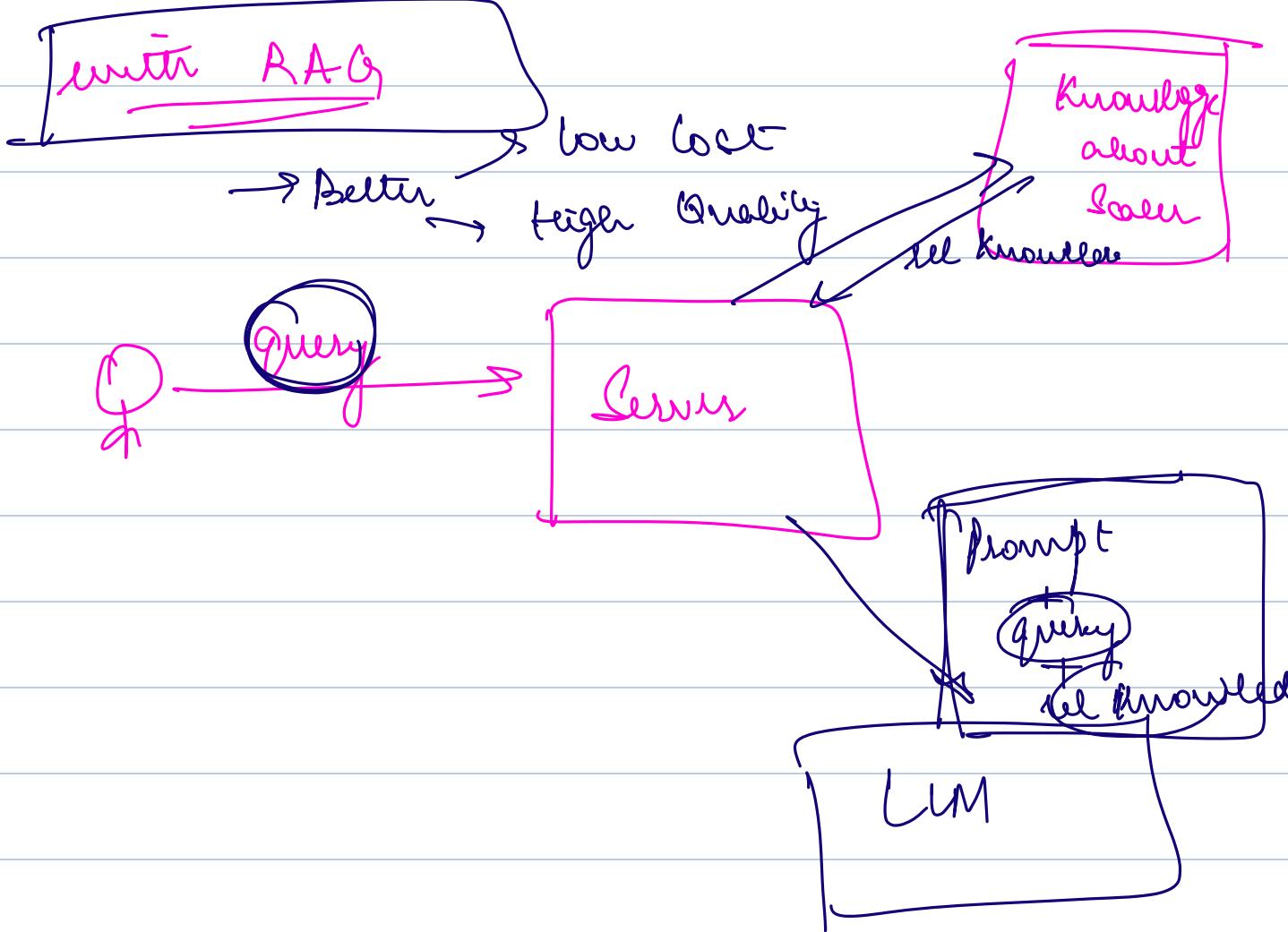
→ give only relevant info
to your model

↳ save cost

↳ better quality.

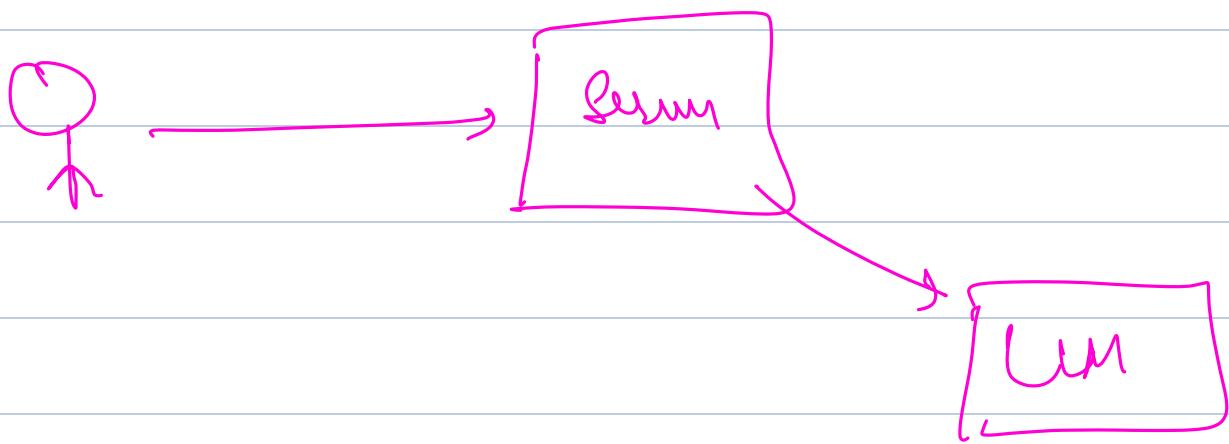
without PAG

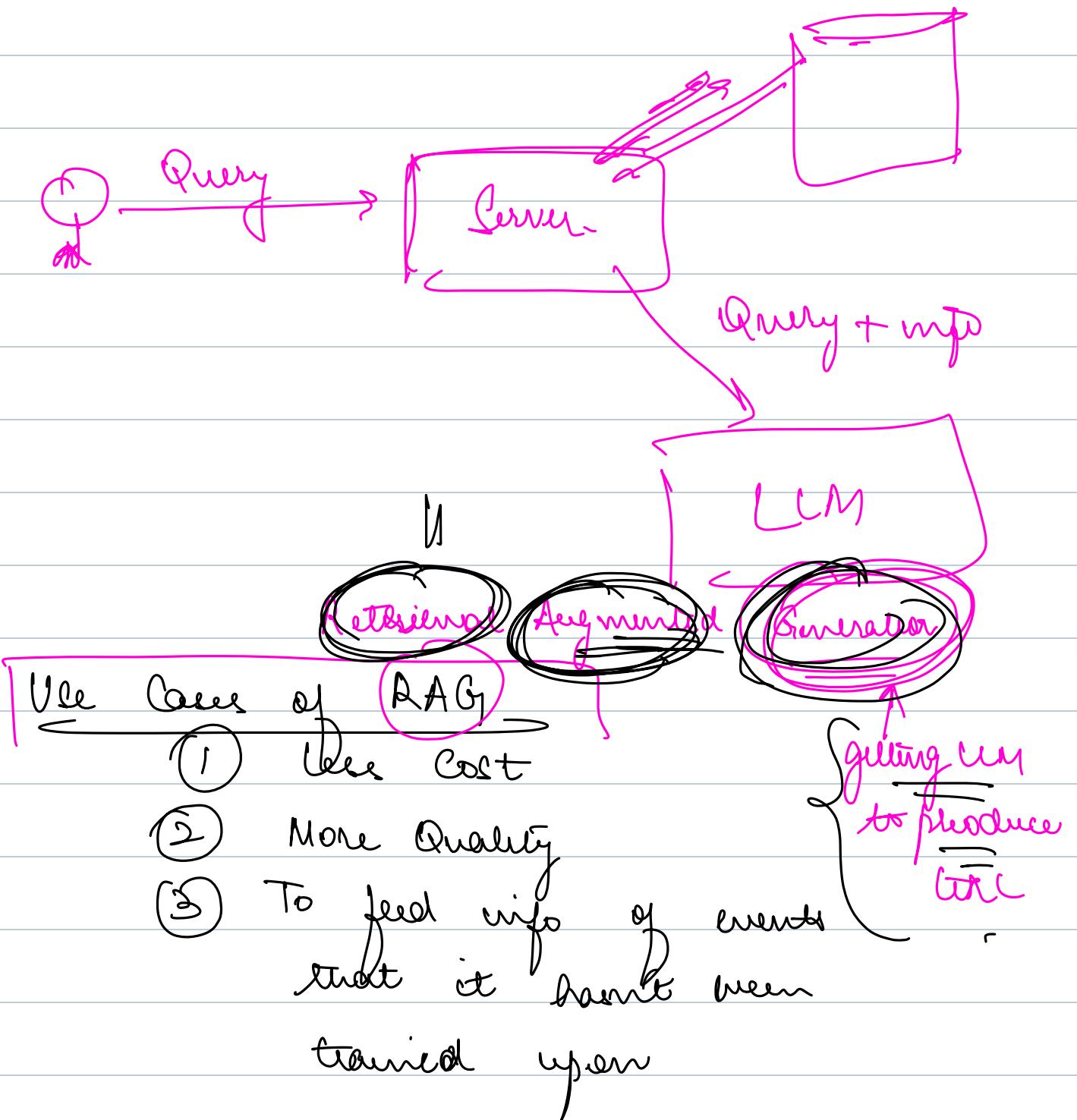




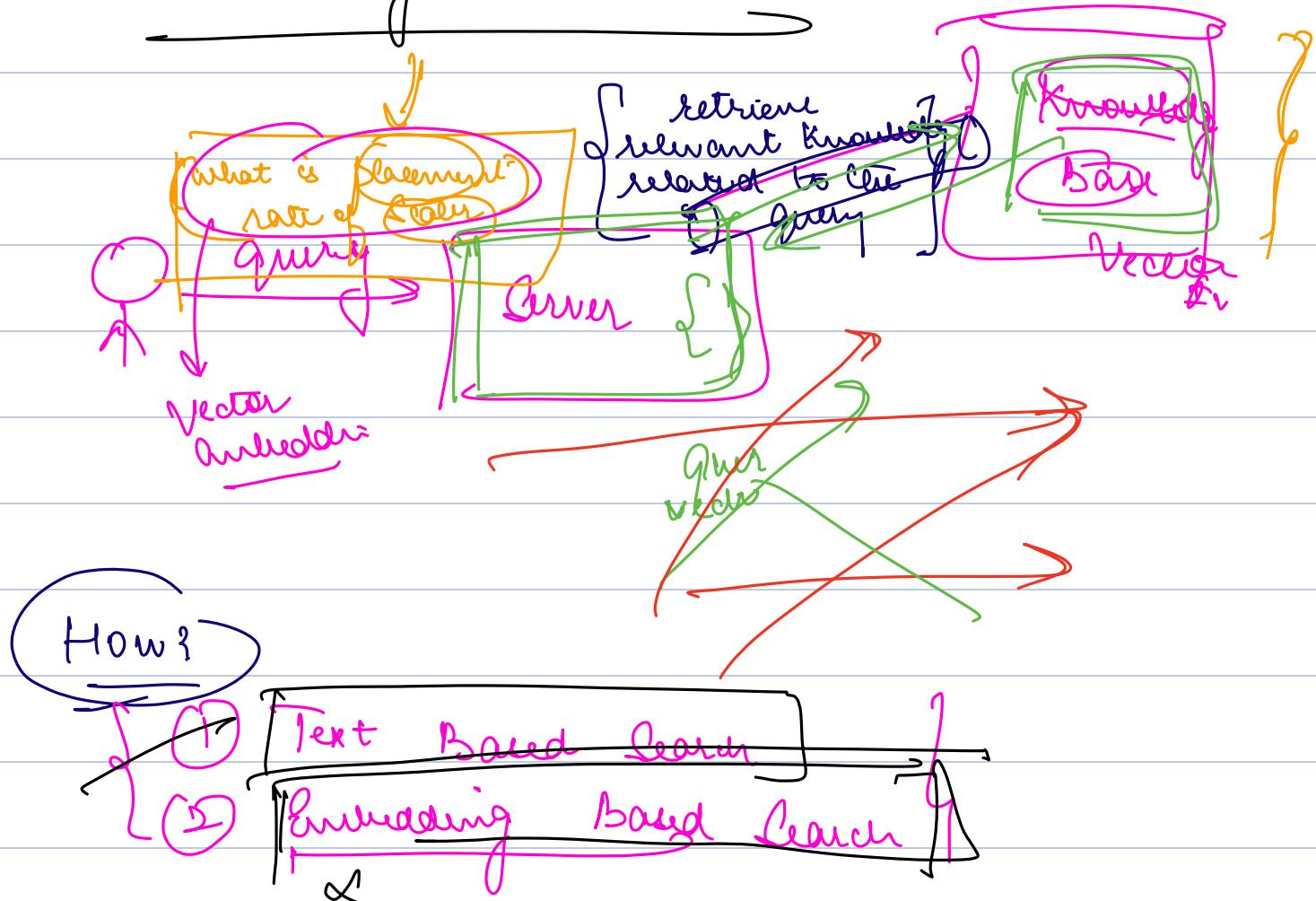
Issue 2

LLMs only know of info till the day were trained.
 ⇒ I might want to query about recent events



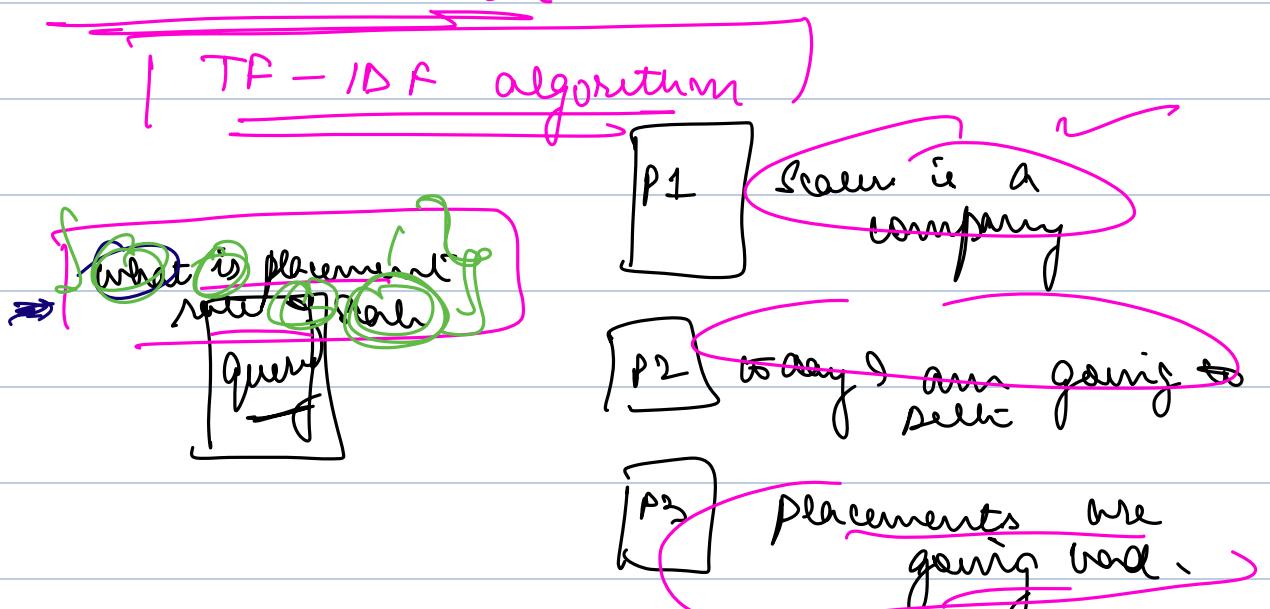


How do you do RAG



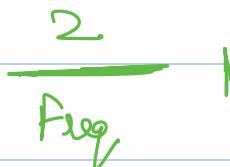
QOT.

Text Based Search



Find which webpage is most relevant to my query.

- ① presence of words in my query in the page.



A hand-drawn diagram showing two green ovals connected by a line. The left oval contains the text "Term Frequency". The right oval contains the text "Inverse Document Frequency". A green arrow points from "Term Frequency" to "Inverse Document Frequency".

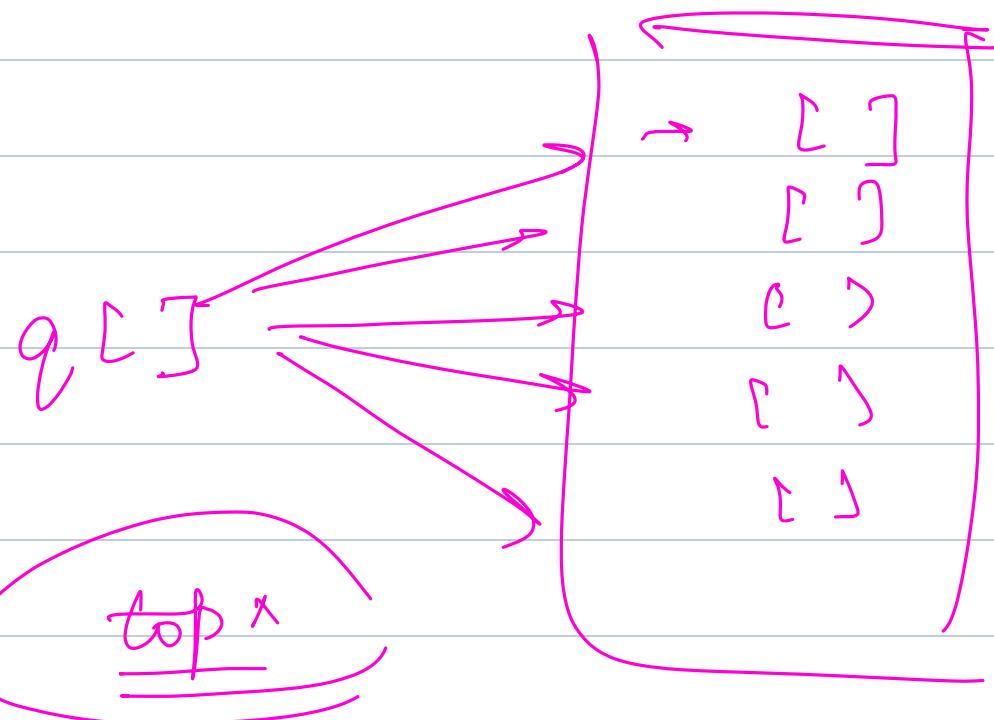
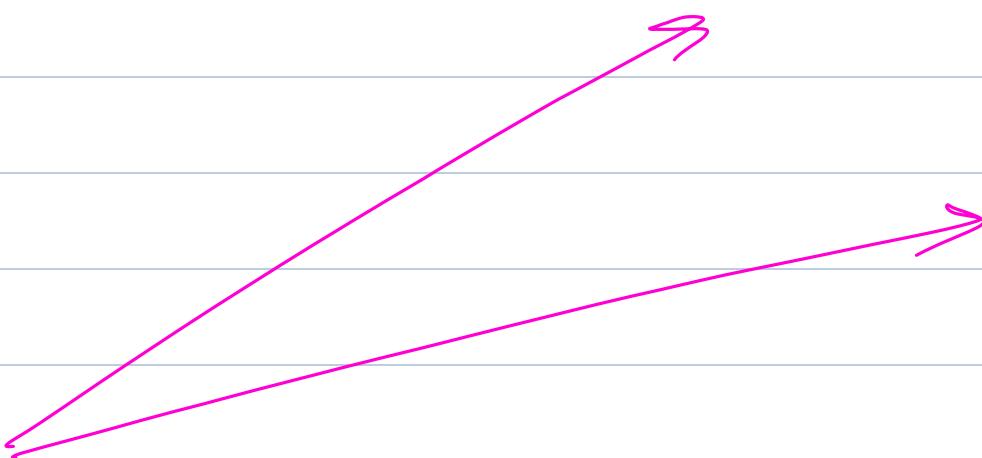
Embeddings Based Search

- ① In DB I store BB embedding for every QnA pair.

- ② when a user sends a query I find closest vector embedding in DB

③ I give top 10 of those to LM.

Cosine Similarity



Special Types of DBS

→ **vector DB**
[**pinecone, PQL vector module**]

Optimizing alarm

How to eval if your RAG is good

①

Precision

②

Recall

% of doc that are returned that match my query

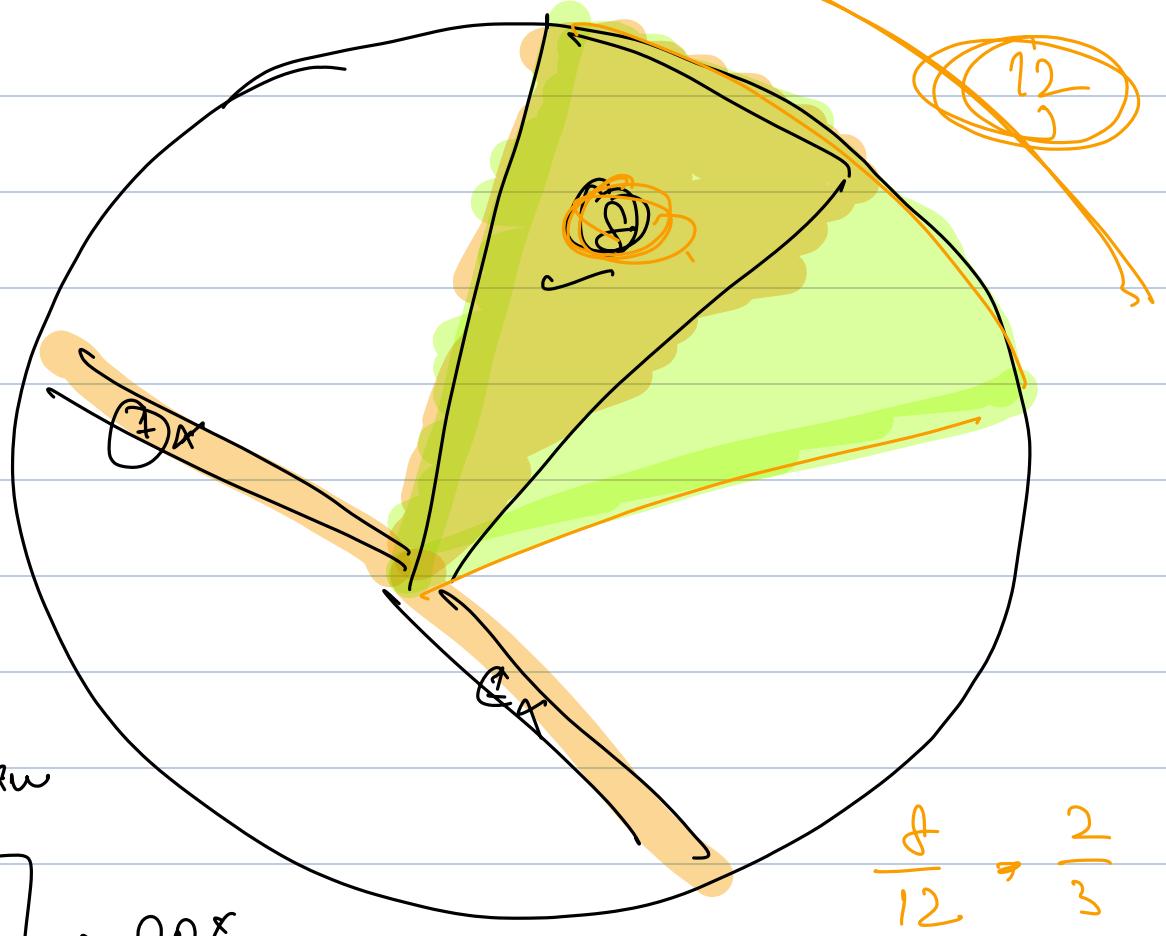
% of relevant docs that were retu

$$\approx 100\% \quad | \quad \boxed{10 \text{ docs}} \rightarrow \boxed{9} \quad \Rightarrow \quad \frac{\text{Precision}}{\rightarrow 90\%}$$

15 docs were relevant across dataset

I got

$$\frac{9}{15}$$



$$\boxed{\frac{8}{10}} \rightarrow 80\%$$

$$\frac{8}{12} \rightarrow \frac{2}{3}$$

66.
%

Lang Chanh

Lava Soden G