

## Gen AI

→ generate new data based on training data.

Unstructured data → Understand pattern & learn. → generate content - Audio, Video, Image, text

Input → discriminative model → Cat / Dog.

Sample data → generative model → [ ] → new data.

- LLM → Large language model (foundational models) use deep learning algorithm to understand natural language.

→ text to text, text to image, text to video.  
Summarizer, translator, code generator.

## • pipeline -

- 1> Data acquisition
- 3> Feature engineering
- 5> Evaluation

- 2> Data preparation -
- 4> Modelling
- 6> Deployment
- 7> Monitoring.

## ① Data acquisition -

- Available data (csv, txt, pdf, xlsx, docs)
- Other data (Internet, DB, API, Scrapping)
- No data (create own data)  
+ LLM to generate data.

Note: less data → Data augmentation.

I am Anu my name is Anu

- + replace with synonyms.
- + Biagram flip (अक्षरों को बदलना)
- + Back translate
- + Add additional noise.

## ② Data Preprocessing -

① Cleanup : HTML removal, emoji, spelling correction

② Basic preprocessing

③ Advance preprocessing

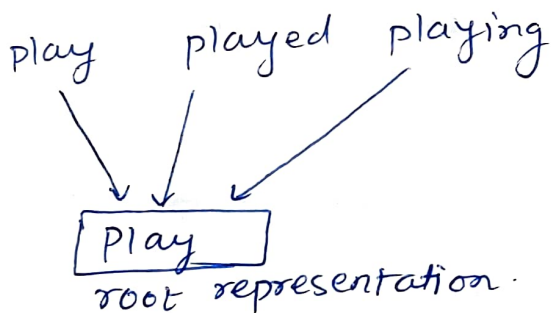
# Basic Preprocessing :

→ tokenization — { word, sentence

# Optional preprocessing :

→ Stop word removal → Lemmatization.  
→ stemming → punctuation removal.  
→ lower case → language detection

### Stemming



### Lemmatization

→ root representation readable.

# Advance preprocessing -

→ parts of speech tagging

→ Parsing (Parse tree)

→ Coreference resolution.

Anamika is a good girl, she is talented also.

## ③ Feature Engineering -

→ text Vectorization

- TFIDF
- Bag of word.
- word2vec
- one Hot
- Transformers model.

## ④ Modelling -

→ Choose models

open source model  
LLM

paid

## ⑤ Evaluation -

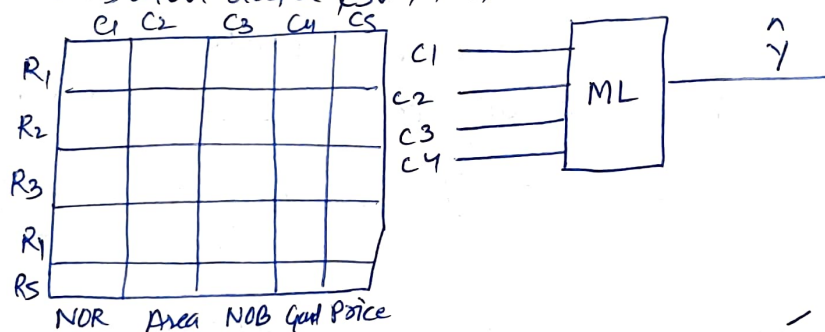
- ① Intrinsic  $\rightarrow$  metrics  $\rightarrow$  Gen AI eng.
- ② Extrinsic  $\rightarrow$  after doing deployment  
 $\rightarrow$  production

③ Deployment  
monitoring + Retraining

## ⑥ Data Representation -

- ① Feature extraction from text/Image -  
vector representation from text/Image  
 $\downarrow$   
Numbers.

#. ML: tabular data csv, xlsx



#. CV: Image, Videos. - Pixels. (0-255)  
white - Black.

# Audio: frequency.

#. text: unstructured data.  
 $\rightarrow$  dimensionality issue

① My name is Anne. X

② I am Anamika. X  
4 5 6

0  
0 0  
0 0 0  $\rightarrow \hat{y}$   
0 0

techniques:

- ① One - Hot encoding
- ② Bag of word (BOW)



# One-Hot Encoding

P <sub>1</sub>	I am Anamika.
P <sub>2</sub>	Anamika is studying.
P <sub>3</sub>	Anamika is learning AI.
P <sub>4</sub>	She is a smart girl.

(i) Corpus: Entire data.

$$P_1 + P_2 + P_3 + P_4.$$

→ find unique words.

$$n = 10 \text{ unique words.}$$

I, am, Anamika, is, Studying, learning, AI, She, a, Smart girl.

I	am	is	Anamika	Studying	learn	AI	She	a	Smart	girl.
1	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0	0	0

(word find kro 1 likho, aur 0 dgaao.)

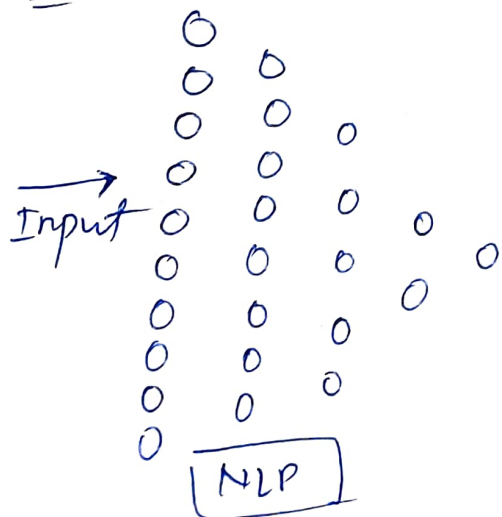
↓  
make a matrix.

P<sub>1</sub> = for each paragraph make 2-D matrices.

$\begin{bmatrix} [1 & 0 & 0 & 0 & 0 \dots] \\ [ - & - & - & - & - ] \end{bmatrix}$  > finding the unique words in the sentence.

P<sub>2</sub> = . Perform the same operation.

(ii) make a neural network of unique words (n)



## \* Drawbacks

→ Sparsity (lots of zeroes)

→ No fixed size.

→ out of vocabulary

→ not capturing semantic relationship.

## Bag of Words

i) Corpus - whole data (unique)

$n=6$

P <sub>1</sub>	I am Anu.
P <sub>2</sub>	I am a girl.
P <sub>3</sub>	I am a Student.

I | am | Anu | a | girl | Student

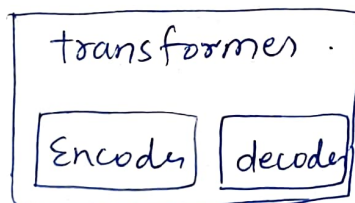
→ It will see how many times each word present in any sentence and make matrix.

	I	am	Anu	a	girl	Student
1	1	1	1	0	0	0
1	1	1	0	1	1	0
1	1	1	0	1	0	1

→ pandas  
→ sklearn  
|  
countvectorizer.

→ Semantic picture not capture.

## Transformer Tree



BERT  
XLNet

GPT

### Open Source

+ mistral  
+ Llama  
+ Gemini  
+ Falcon  
+ Claude  
+ StableLM

## Prompt Designing

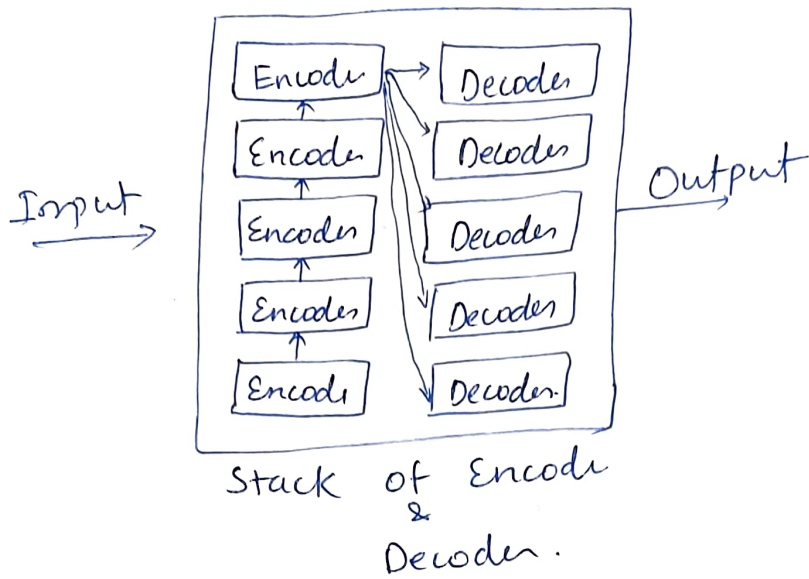
Zero shot Learning - using a single command.

Few shot Learning - command + Example.

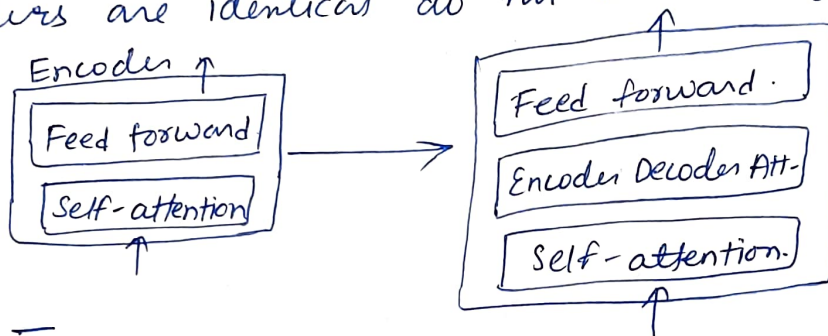
## Transformer Architecture

Input → [transformer] → Output.

Input → [Encoder] → [Decoder] → Output



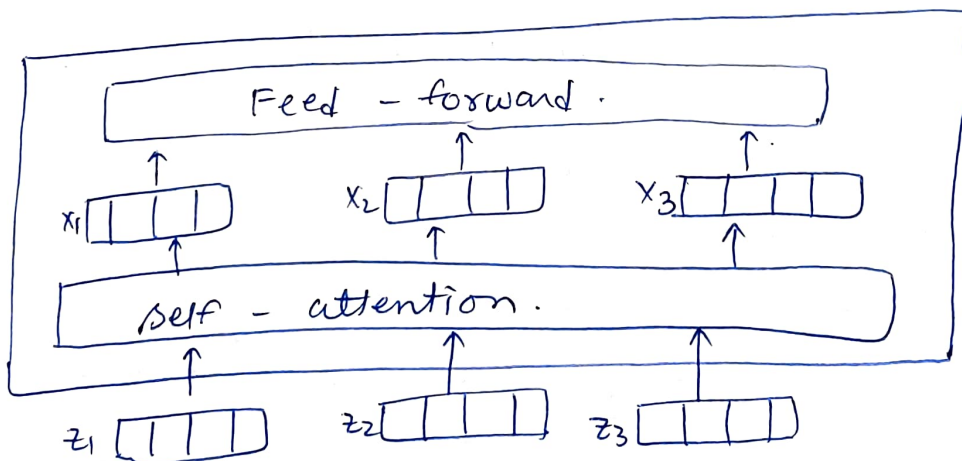
→ Encoders are identical do not share weights -



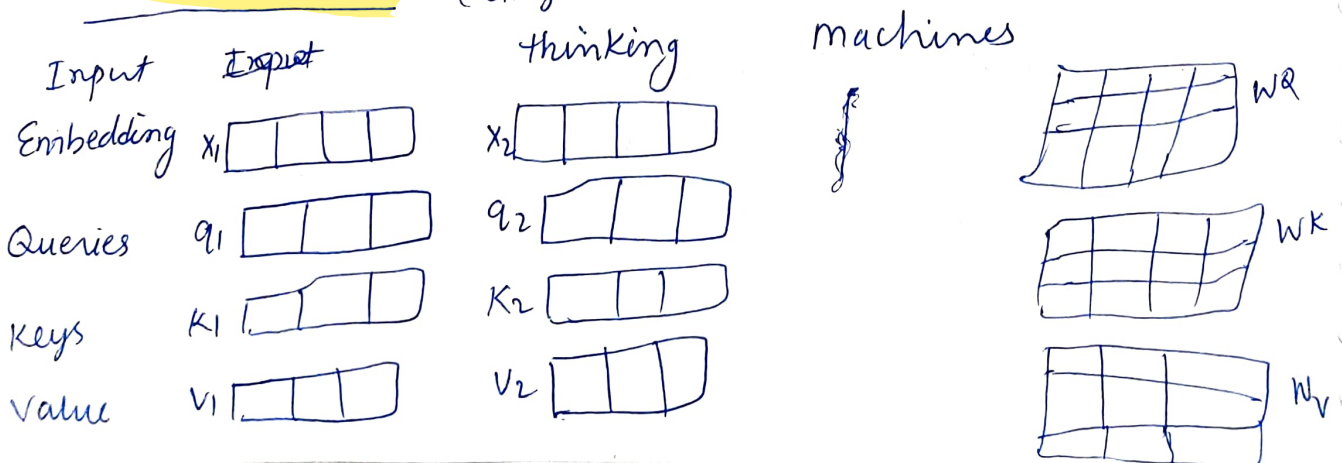
Working

I am Anamika.

Word 2 Vec.  $z_1$   $z_2$   $z_3$



• Self-attention - (Single head attention)





Score

$$q_1 \cdot k_1 =$$

$$q_1 \cdot k_2 =$$

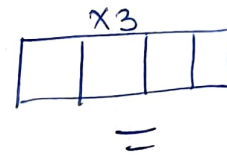
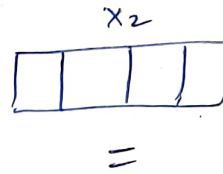
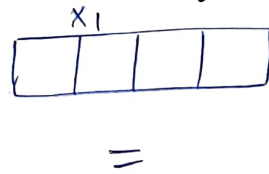
Softmax  $\rightarrow$  activation function

$$\text{Softmax} \left( \frac{\begin{matrix} Q \\ \begin{bmatrix} \square & \square & \square \end{bmatrix} \end{matrix} \times \begin{matrix} K^T \\ \begin{bmatrix} \square & \square \end{bmatrix} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} V \\ \begin{bmatrix} \square & \square \end{bmatrix} \end{matrix}$$

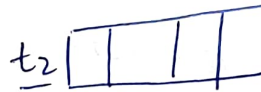
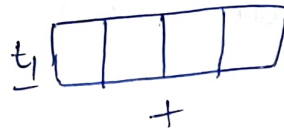
$$= \begin{bmatrix} \square & \square & \square \end{bmatrix} z$$

## #-Positional Encoding -

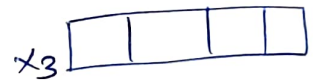
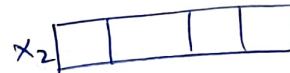
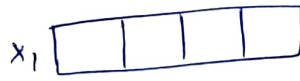
Embedding  
with time  
signal



Positional  
Encoding



Embedding

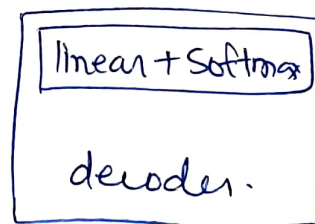
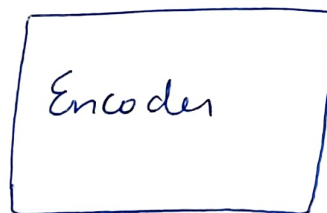


Input

I

am

Anamika.



$\rightarrow$  transformer.  
Architecture

Linguistics - study of language.

- Clear instruction
- Adopt a persona
- Specify the format
- Avoid leading the answer
- Limit the scope

### Vector DataBase

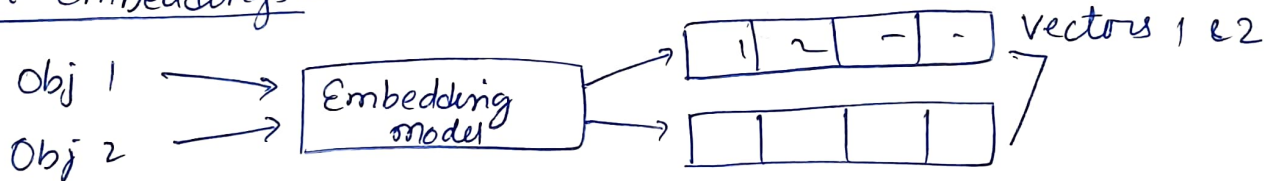
Unstructured data.

High Dimensional data. (pdf, Image, word)

↓  
High dimensional data vector 

↓  
Vector database

### • Vector embeddings -



Ex - Unstructure data as text -

Embedding model.

Feature	King	Queen	man	Women	monkey
gender	1	0	1	0	1
Wealth	1	1	0.5	0.3	0
Power	1	1	0.5	0.3	0
Weight	1	0.5	0.5	0.3	0.2

male  $\rightarrow$  1  
female  $\rightarrow$  0  
yes  $\rightarrow$  1  
NO  $\rightarrow$  0

King  $\Rightarrow [1, 1, 1, 1]$

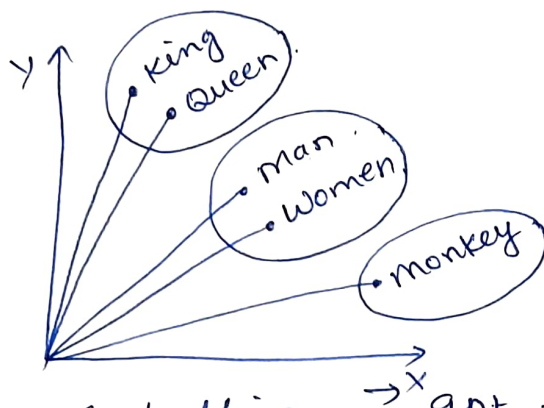
Queen  $\Rightarrow [0, 1, 1, 0.5]$

man  $\Rightarrow [1, 0.5, 0.5, 0.5]$

Women  $\Rightarrow [0, 0.3, 0.3, 0.3]$

monkey  $\Rightarrow [1, 0, 0, 0.2]$





Prince — King  
Queen.

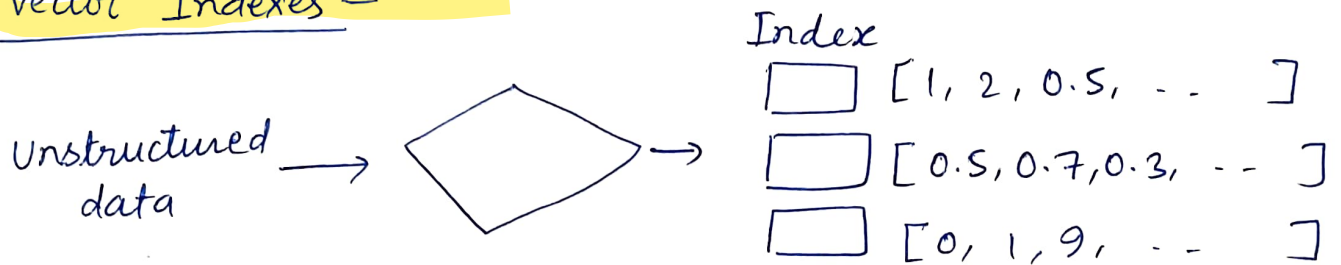
→ Understand the essence of other words.

- Open AI Embedding → gpt model
- Hugging face Embedding → open Source.
- llama 2 Embedding → facebook
- Google Palm embedding → Google

#### Vector DB

- Pinecone
- Neo4j
- Weaviate
- ChromaDB
- FAISS

#### • Vector Indexes —



- It will check the new data is falling under which cluster and then search only in those indexes.
- Long term memory for LLMs.
- Use for Semantic Search based on meaning of context.
- Similarity Search.
- Recommendation engine as well.