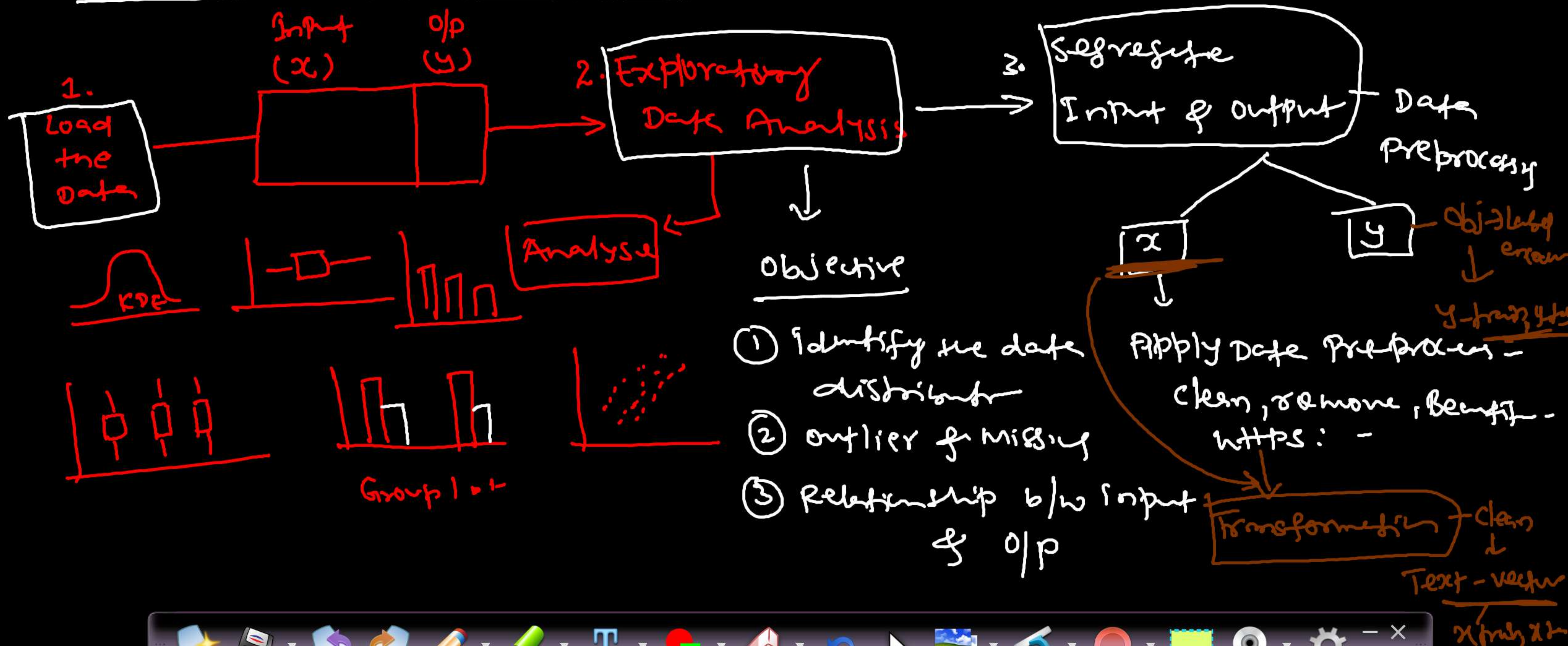


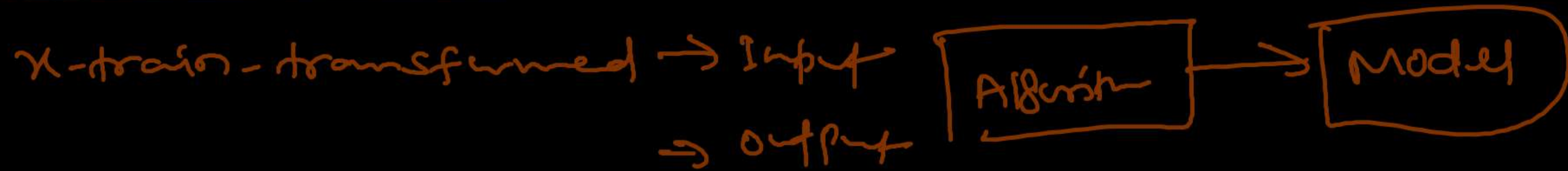
Introduction To NLP

Model Building Pipeline



* Learning Phase

Build the model



Testing Phase

(a) Predict on the test data



(b) Evaluation

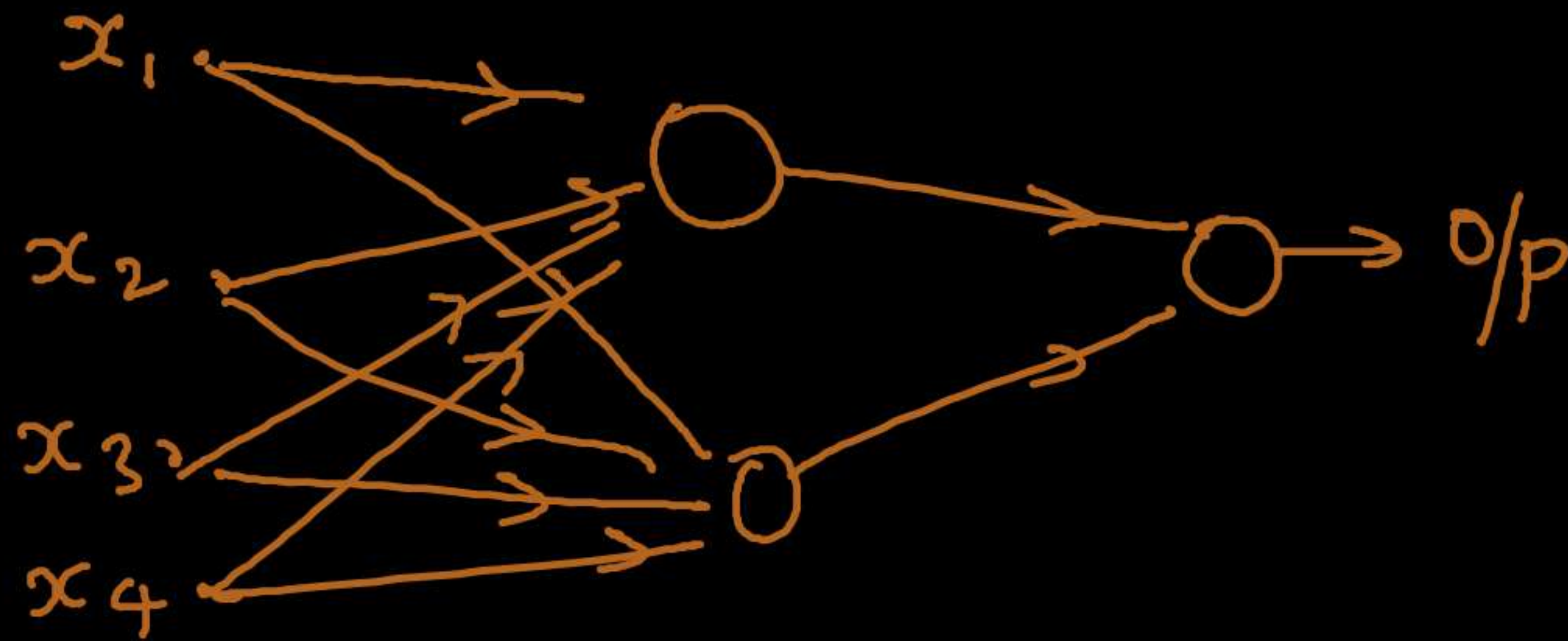
$y_{\text{-train}} \text{ vs } y_{\text{-train-pred}}$
 $y_{\text{-test}} \text{ vs } y_{\text{-test-pred}}$

} Score $\begin{cases} \text{high bsc} \\ \text{high variance} \end{cases}$

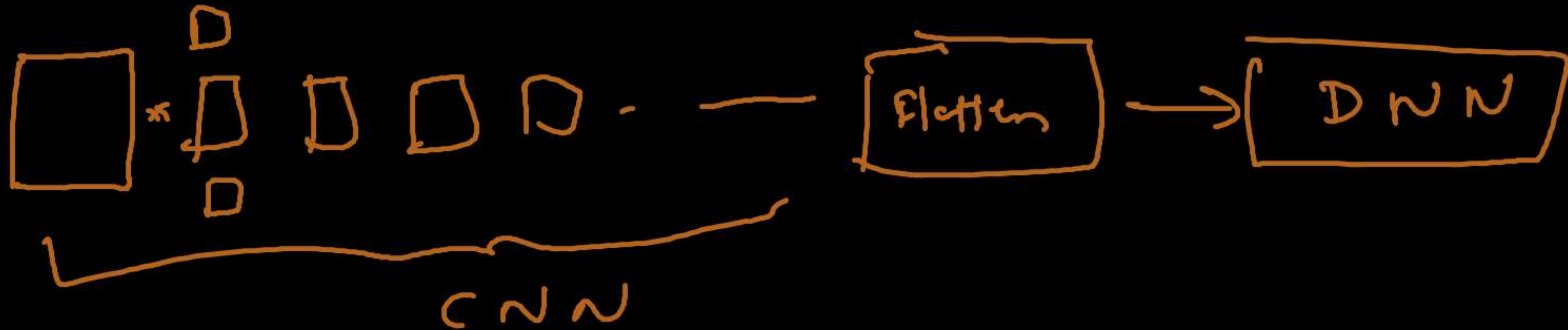
NLP in Deep Learning Perspective

① ANN → MLP/DNN → Tabular Data - Table $\begin{matrix} \text{Rows} \\ \text{columns} \end{matrix}$ - Label Data

Home size Age Nb. of rooms Location Price ?
 x_1



② CNN & Computer vision → Image classification, Object Detection



③ Sequential Data ^{Text Data}

Time series forecasting

NLP — Text Analytics / Text Generation

- (i) RNN
- (ii) LSTM RNN
- (iii) GRU RNN
- (iv) Encoder-Decoder
- (v) Seq-2-Seq

(vi) Transformer (Attention is all you need)

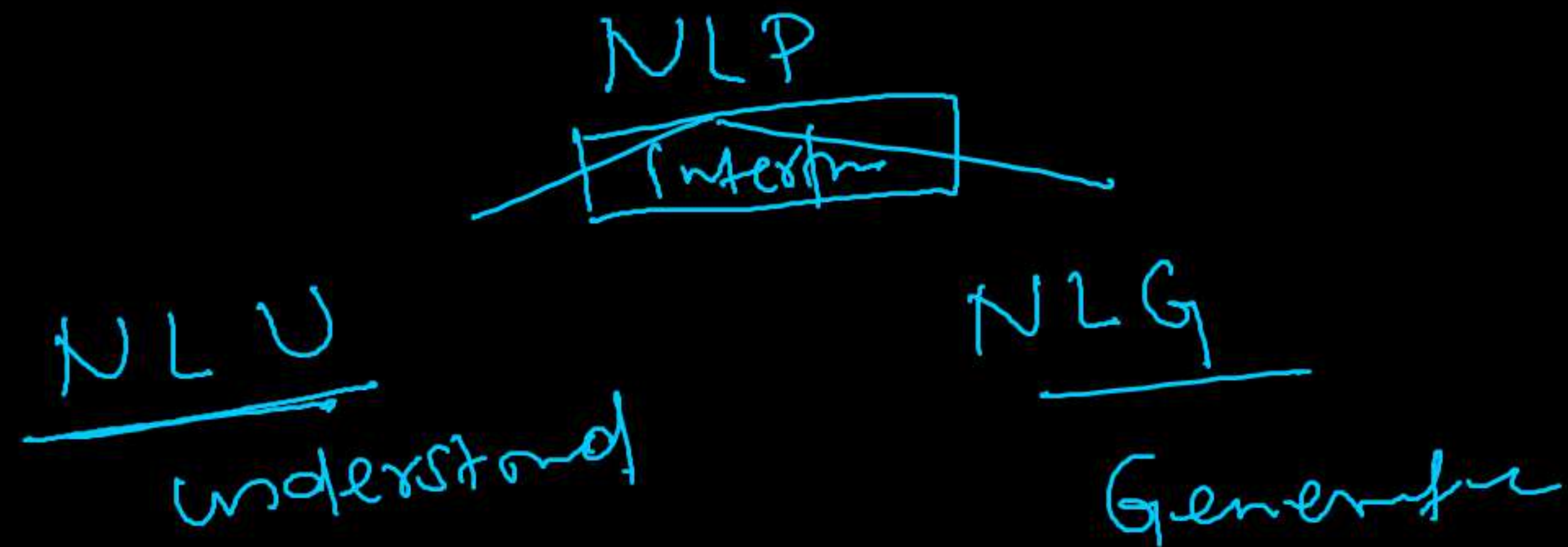
Adv :- Gen AI

- BERT
- GPT-3
- LLaMA 3.1
- LLM models

eg :-
 Chat Bot App → Q & A
 Language Translation
 Text correction / Generation

Interview Purpose

* What is NLP? → NLP is the branch of AI that involves computers to ~~understand~~, ~~interpret~~ & ~~generate~~ human language.



→ There are wide range of task in NLP :—

① Text Classification :- Spam detection, Sentiment Analysis.
Adult content filtering etc.

② Information Extraction → Named entity Recognition, (NER)
Parts of speech tagging

NER eg :- Apple is looking to launch my macbook for \$1500.
 ORG Product Price

eg :- She runs fast.
 Pronoun very Adverb
 POS → Action Describe
 SWS

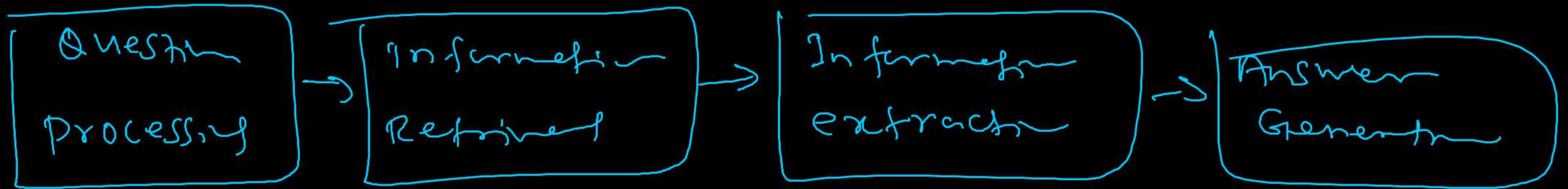
③ Information Retrieval :- Search & Retrieval

→ Documents, website, https, web search

④ Text Summarization → News feed creation, case report summarization in a big law firm

⑤ Machine Translation :- Translating one language to another

⑥ Question & Answering :-



DNN & CNN

Q:- Why is NLP Hard?

1. Ambiguity : - ie. uncertainty of meaning

eg. - The chicken is ready to eat.

eg. - The car hit the pole while it was moving.

2. Complexity of representation

like → poems, sarcasm, phrases etc.

eg. - You have a football game tomorrow, break a leg!

eg. → Yeah, right, because they worked so well last time.

Also called as

Text Representation (AKA vectorization or feature extraction)

1. Bag of words - Bow → capture word freq/count
2. Term Frequency - Inverse Document Frequency (TF-IDF) → capture the word importance
- Phase I → Basic Language representation
- Non-numeric repres. → to → Numeric representation
- Text to vect
Char to num
- ML
Eng

3. Word2vec (word to vector) → 2013 (Google)
4. Glove (Global vector) → 2015 (Stanford)
5. Fast Text → 2016 (Facebook)
- capture the word semantic meaning
- During DL era
- 2017
- Phase II → Distributed Language representation - 2009-2017

Phase 3

6. ELMo → 2017 (University of Washington)

7. Transformer (Attention all you need)

8. BERT → 2018 - It's just face

9. GPT-40 → 2024 (Open AI)

10. LLAMA-3 8.1 (July) - (Meta)

11. Gemini-1.5 - 2024 (Google)

12. Phi-1.5 - micro (2024)

During Gen AI

Exo

+ Chat Bot

+ RAG

+ Fine-tuning

+ Prompt engineering

+ Style defense

→ Contextual Language Representation (LLM) / Gen AI

Terminology

- (a) Document :- It is a single piece of text. It can be a sentence, paragraph, email, ~~message~~, article review etc - -
- (b) corpus :- collection of documents
- (c) Vocabulary :- set of all unique words that appear in a corpus.
- (d) Vectorization → It's a technique to convert text data into numerical vectors
- (e) Documents vectors → Numerical represent of documents

(f)

Document Term matrix :-

Rows

	Able	ability	thorough
1	0.1	0.1	0.5
0.1	0.2	0.3	0.8
0.2	0.3	1	1.3
1.2	2		1.9
2.3	3		2.7

* Text Cl

Able able

ability ability

Dog dog

Cat cat

- ① Removing
- ② converting
- ③ Removing
- ④ converting to root form

Webinar Chat

Vivek Shinde to Hosts and panelists

VS yes sir

chetan anand to Hosts and panelists

CA yes

Rohan to Hosts and panelists

R Yes

PRAMOD K. to Hosts and panelists

PK yes

SAHAS V SWAMY to Hosts and panelists

SV good to go sir!

Sneha to Hosts and panelists

S Yes Sir

Santoshkumar Pandit to Hosts and panelists

SP yes sir

going good

Who can see your messages? Recording on

To: Hosts and panelists

Type message here...