

Natural Language Processing

Basic Terminologies of NLP:-

① Corpus: A corpus is a large and structured set of texts used in NLP. It serves as a database for linguistic analysis and model training.

② Document:

Applications of NLP:-

① Sentiment Analysis:-

Analyze text data from social media, reviews and feedback to determine sentiment expressed (Positive, Negative, Neutral) towards the products, services or topics. This is widely used in market research to gauge public opinion.

② Chatbots and Virtual Assistants:-

Powering conversational agents that can understand and respond to human queries in natural language. Chatbots are used in customer service, while virtual assistants like Siri, Alexa, and Google Assistant help users perform tasks through voice commands.

③ Machine Translation:-

Automatically translating text from one language to another. Applications like Google Translate enable users to understand content in foreign language.

④ Named Entity Recognition: Identifying entities into pre-defined categories like people, organizations, etc. This is a type of extraction.

⑤ Speech Recognition

Converting speech into text.

⑥ Text Categorization

Categorizing text into predefined categories.

LEVELS OF NLP

① Morphology

② Syntax

③ Semantics

④ Discourse

⑤ Pragmatics

① Morphological

Morphology is the study of words. At the smallest level, a word is composed of morphemes, which are meaningful units of language.

④ Named Entity Recognition (NER):

Identifying and classifying key elements into text into predefined categories such as names of people, organisation, locations, dates and more.

This is useful in content classification, information extraction etc..

⑤ Speech Recognition

Converting spoken language into text.

⑥ Text Classification :-

Categorizing text into predefined groups. This is used in email filtering.

LEVELS OF NLP (or) STATES :-

① Morphological Analysis

② Syntactic Analysis

③ Semantic Analysis

④ Discourse Analysis

⑤ Pragmatic Analysis

① Morphological Analysis :-

Morphology is study of structure and formation of words. At this level, NLP systems analyze the smallest units of meaning within words, known as "morphemes". This includes underlying root words, prefixes, suffixes.

② Syntactic Analysis - (Parsing sentence)

Syntax is concerned with the arrangement of words in sentences and ~~grammatical~~ grammatical relationships between them. At syntactic level, analyses sentence structure to understand how words combine to form a sentence. This involves parsing the sentence to identify subjects, predators, objects.

③ Semantic Analysis,

Semantic Analysis deals with meaning of words, phrases, sentences and texts. At semantic level, NLP aims of understanding the intended meaning conveyed by a combination of words.

④ Discourse Analysis,

This analysis goes beyond the analysis of individual sentences. It involves analyzing the coherence and structure of longer texts.

⑤ Pragmatic Level,

Pragmatic focuses on use of language in context and how context influences the interpretation of meaning. It considers factors like the relationship b/w speaker and the listener, speaker's authority, speaker's intent, etc.



ited

STEPS IN NLP:-

- ① Tokenization
 - ② Lemmatization
 - ③ Stop Words
 - ④ Stemming
 - ⑤ Bag of Words
 - ⑥ TF-IDF [Term Frequency - Inverse Document Frequency]
 - ⑦ N-Gram
 - ⑧ Word2Vec → mapping semantically similar words to proximate point in vector space
 - ⑨ Avg Word2Vec.
- Vocabulary Create
Vectorization

I have been to
the college.

arrangement of words
metrical relationship

level, analysis

how words

involves persons,
objects, products,
events.

meaning of words,

Semantic level,

indicated means
words.

analysis of
analysis of the
texts

calculate the repetition of word

TF = $\frac{\text{No. of times term } t \text{ appears in document}}{\text{Total no. of terms in document}}$

calculate the importance of term

IDF = $\log \left(\frac{\text{Total no. of documents } D}{\text{No. of documents containing term } t} \right)$

Sent 1 → good boy
Sent 2 → good girl
Sent 3 → boy girl good

total score
no. of word repeat
no. of unique words
total words

	T.F sent 1	sent 2	sent 3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

	IDF
good	$\log(\frac{3}{3})$
boy	$\log(\frac{3}{2})$
girl	$\log(\frac{3}{2})$

T-F - T.D.F

	pool	long	mid
short	0	?	?
short	0	?	?
short	0	?	?

multiple
dimensions
of space

dimensions

Stream processing

Extreme processing

Volume

on the "vol-

UNIT-IV
Semantic & Discourse Analysis

Word Embedding :-

① Word2Vec :-

- * Word2Vec creates vectors of the words that are "distributed numerical representations of word features" — these words features could comprise of words that represent the context of individual words present in our vocabulary.
- * Word embeddings eventually help in establishing the association of word with another similar meaning word through the created vectors.
- * It is a two layer neural network to generate word embeddings given a text corpus.
- * Word Embeddings — Mapping of words in Vector Space

$\begin{matrix} 0.52 \\ 0.76 \\ 1.21 \\ 0.22 \\ -1.36 \\ 0.49 \\ -3.69 \\ -0.07 \end{matrix}$	$\xrightarrow{\text{Man}}$	$\begin{matrix} 0.73 \\ 0.59 \\ -1.67 \\ 1.32 \\ 0.36 \\ -1.49 \\ 2.71 \\ 0.05 \end{matrix}$
	Women \rightarrow	

$$\text{King} - \text{Man} + \text{Women} = \text{Queen}$$

Working

* The u
occu
tsp

Extr

The

Types

Two d

(a)

(b)

(a) CB

H

$V_{5 \times 1}, 0$

$V_{5 \times 1}, 0$
Vector
"Set"

Analysis

are "distributed" these words

represent
ent in our

the
meaning

word

vector

Working of Word2Vec :-

- * The word2Vec objective function causes the words that occur in similar contexts to have ~~similar~~ similar embeddings.

Ex:- The kid Said he would grow up to be superman
 The child Said she would grow up to be superman.

The words kid and child will have similar word vectors due to similar context.

Types of Word2Vec :-

Two different model architectures that can be used by Word2Vec to create the word embeddings are:

(a) Continuous Bag of Words (CBOW) Model.

(b) Skip Gram Model

\hookrightarrow predict context word from target.

\hookrightarrow predict target word from context

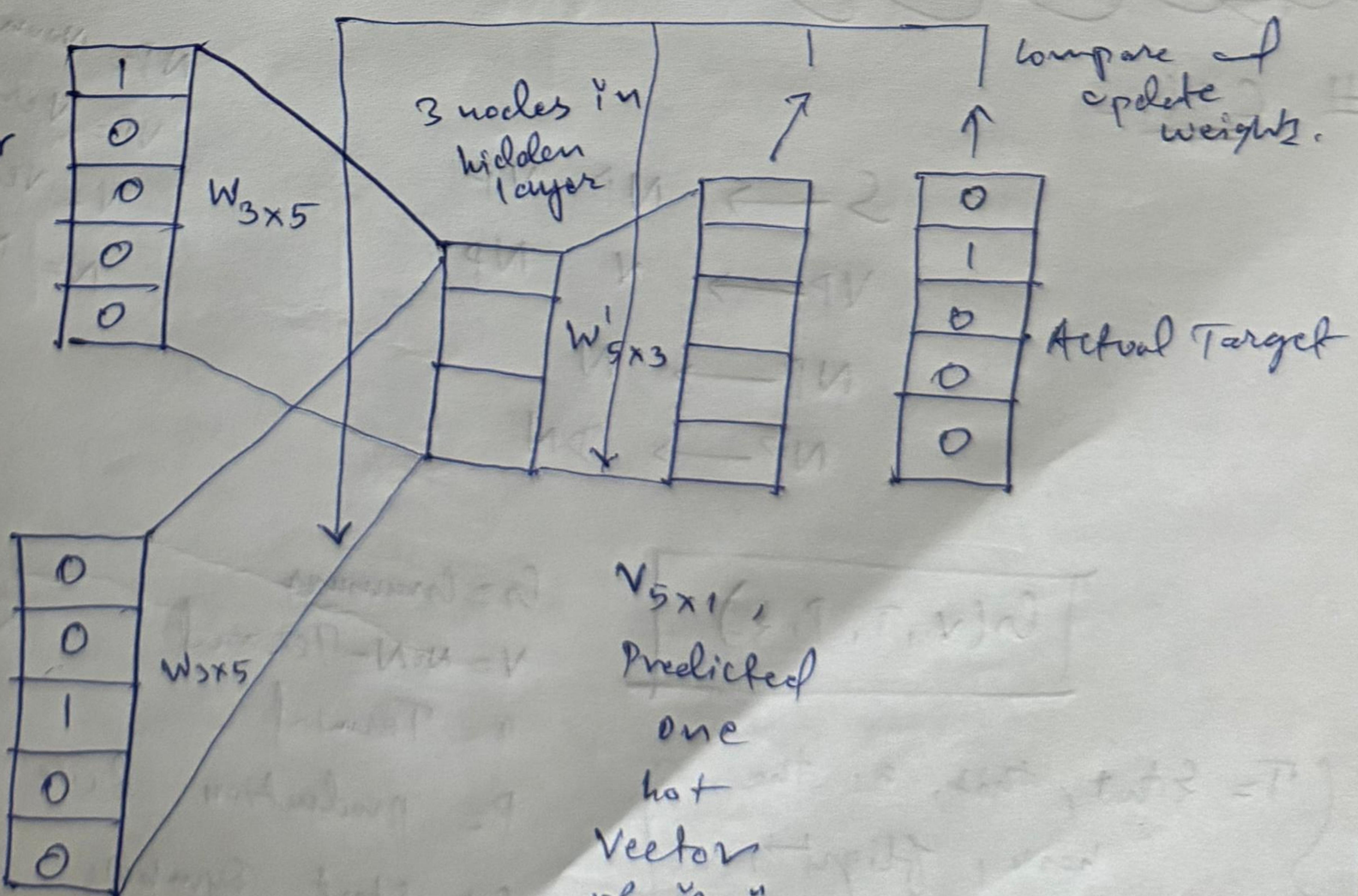
(a) CBOW :-

[Hope van Pet] You free

window size = 3

$V_{5 \times 1}$, one hot vector
of "Hope"

$V_{5 \times 1}$, one hot
vector of
"set"

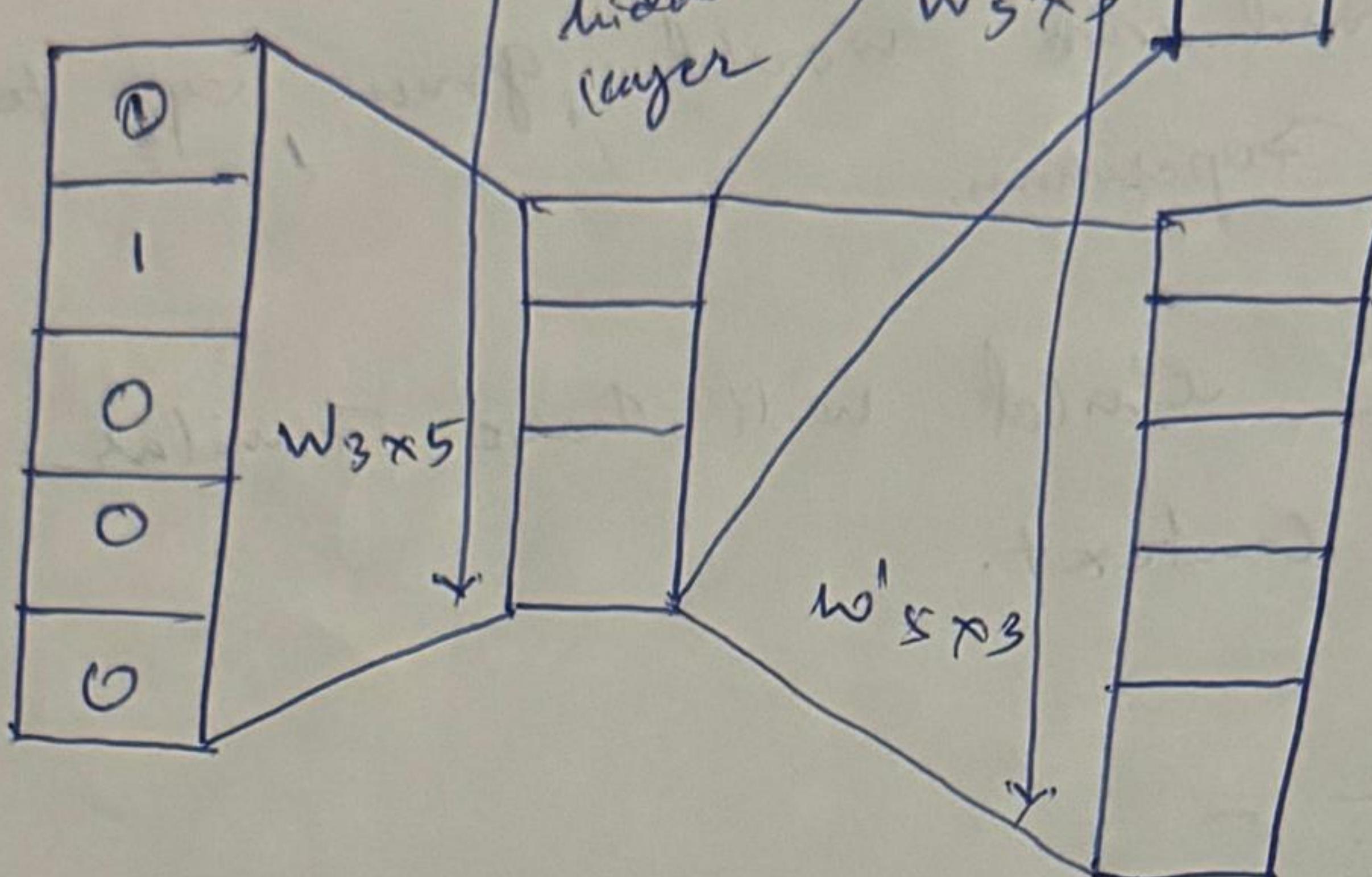


$V_{5 \times 1}$,
Predicted
one
hot
Vector
of "can"

(b) Skip Gram :-

Hope can set you free

$V_{5 \times 1}$, one hot
vector of "can"



Compare of
predicted weight
with weight

Actual Target

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

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

$V_{3 \times 1}$, predicted
vector of "set"

Skip Gram is used to predict the context word for a given target word. It's reverse of CBow algorithm. Here, target word is input while context words are output.

~~ADDED DOUBLE AND CONFIDENCE~~

CFG:-

$$\begin{array}{l} S \rightarrow NP \ NP \\ VP \rightarrow V \ NP \\ NP \rightarrow N \\ NP \rightarrow DN \end{array}$$

$NP = \text{Noun phrase}$
 $VP = \text{Verb phrase}$
 $N = \text{Verb}$
 $DN = \text{Noun}$

$$G(V, T, P, S)$$

$T = \{\text{that, this, a, the, book, flight, meat, does, Dead, Town, hall}\}$

$G = \text{Grammar}$
 $V = \text{Non-Terminal}$
 $T = \text{Terminal}$
 $P = \text{production Rule}$
 $S = \text{Start Symbol}$

$V = \{S, NP, Noun, VP, Det, Aux Verb\}$

$P = \{S \rightarrow NP \ NP, NP \rightarrow VP \ NP, VP \rightarrow V \ NP, NP \rightarrow N, NP \rightarrow DN\}$

Det →
Noun →
Verb →
Aux →

Rule :-

α
 β
 γ
 δ

Q.1 The ma

of
relative weight
Actual Page

$$P = \left\{ \begin{array}{l} S \rightarrow NP VP \\ S \rightarrow Aux NP VP \\ S \rightarrow VP \\ NP \rightarrow Det NOUN \\ VP \rightarrow Verb \\ VP \rightarrow Verb NP \end{array} \right.$$

Def → this / that / a / the
 Noun → book / flight / John / ball / meal
 Verb → book / include / Read.
 Aux → does / is

Rule 1 —

$$\alpha \rightarrow \beta$$

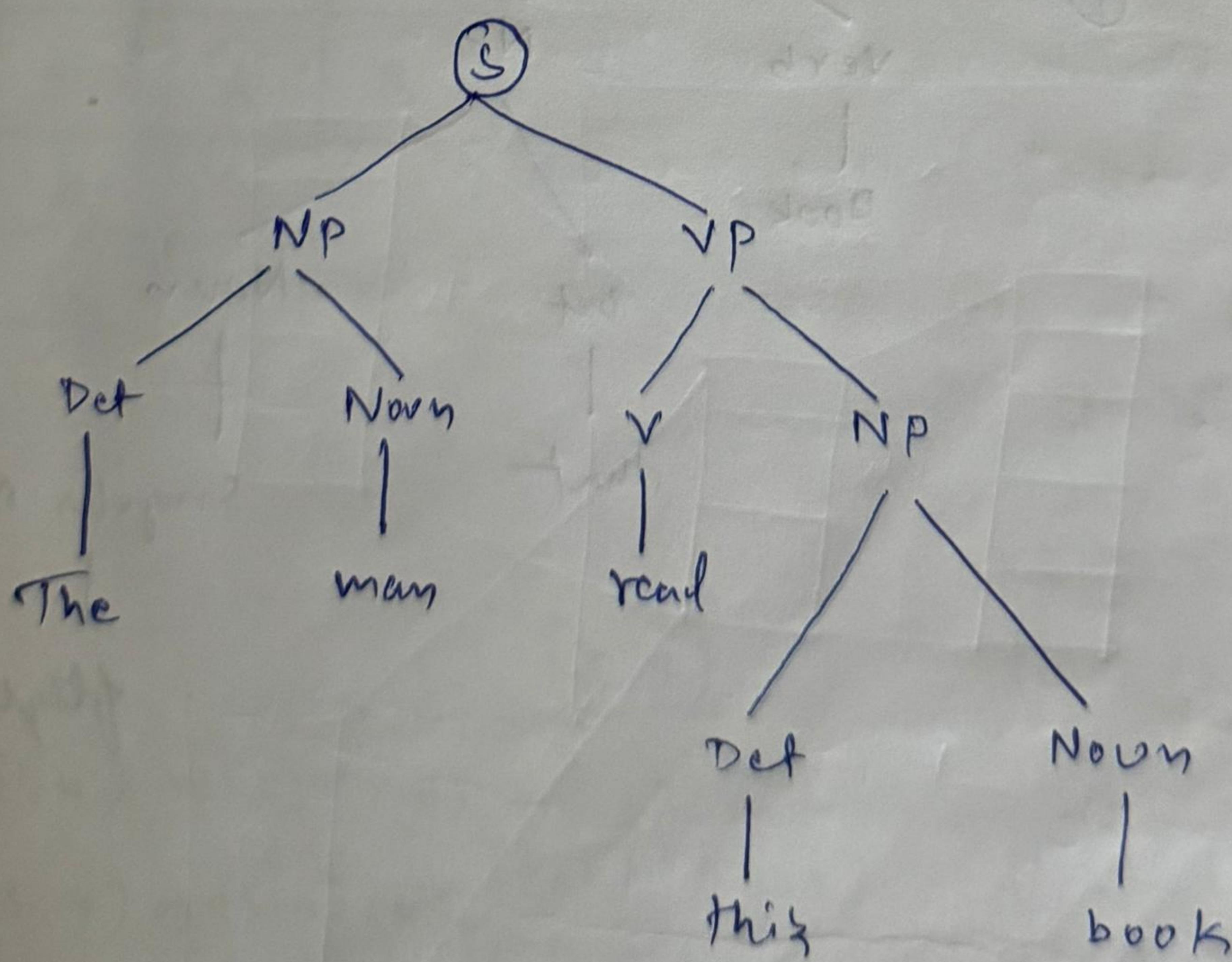
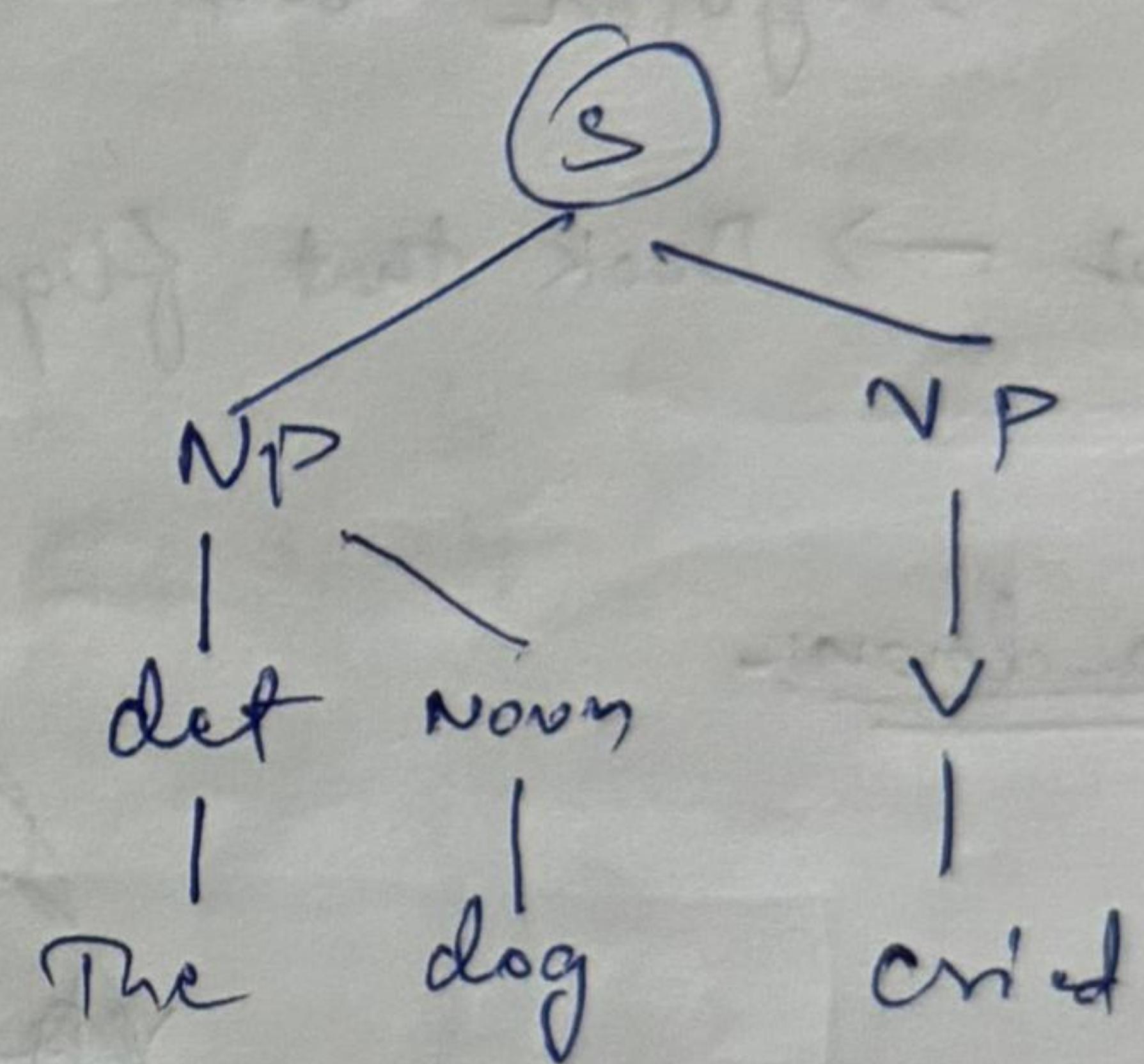
$\alpha \rightarrow$ Single variable

~~β ∈ (V+T)*~~

$$\beta \in (V+T)^*$$

Q.1 The man read this book

The dog cried
 det Noun Verb



Nouns,
Det,
Verb }

Parsing

PARSING

Top-down

Bottom up

②

Bottom up

Q. GRAMMER:

$S \rightarrow VP$

$VP \rightarrow \text{Verb } NP$

$NP \rightarrow \text{Def } Noun$

$ND \rightarrow \text{Def } Noun$

$\text{Def} \rightarrow \text{that}$

$\text{Noun} \rightarrow \text{Singular Noun}$

$\text{Verb} \rightarrow \text{Book}$

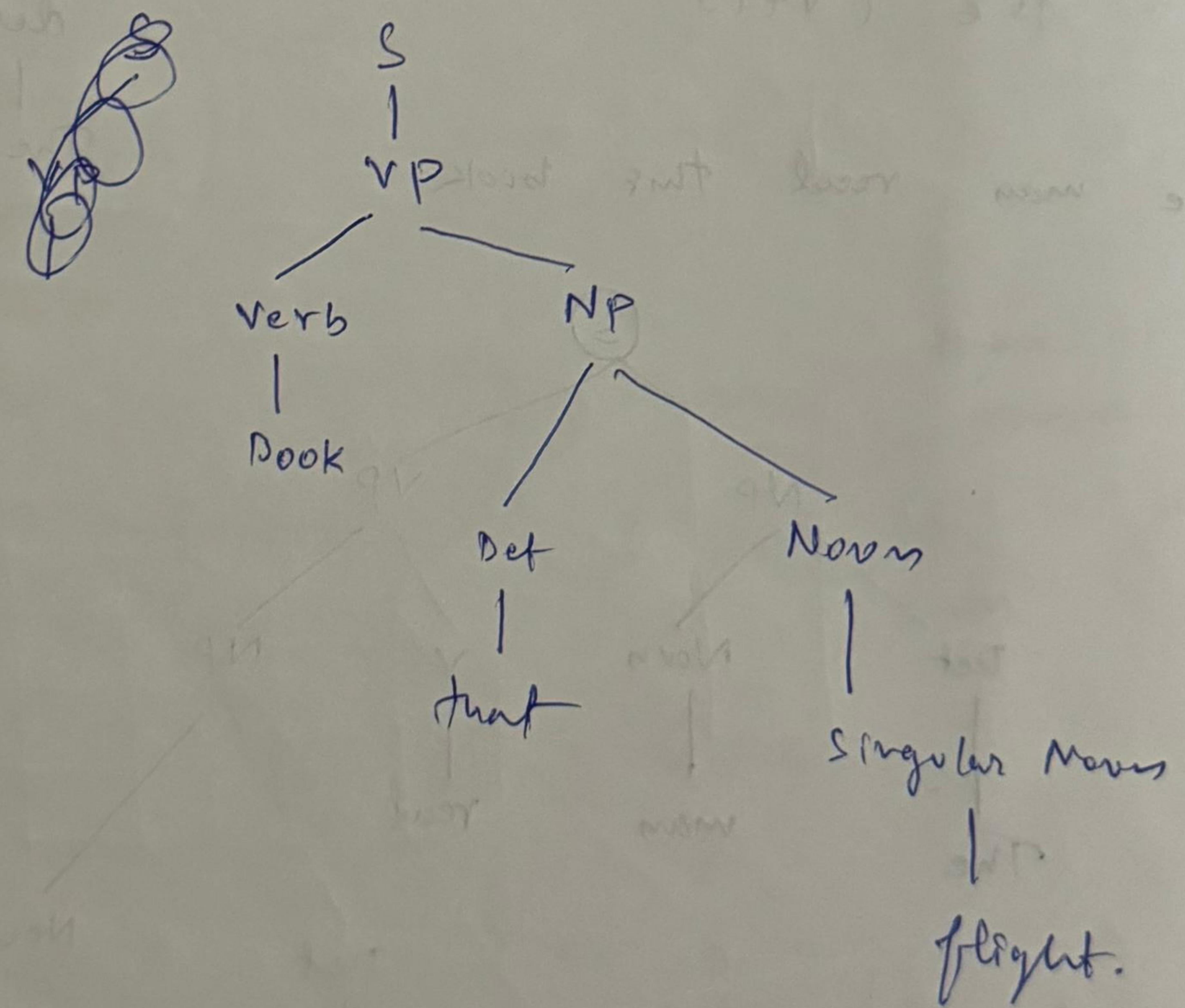
$\text{Singular noun} \rightarrow \text{flight.}$

Input \rightarrow Book that flight.

CKY Parse

Q. Check whether the following

① Top-down:



Done

1	1
2	5
3	1
4	1

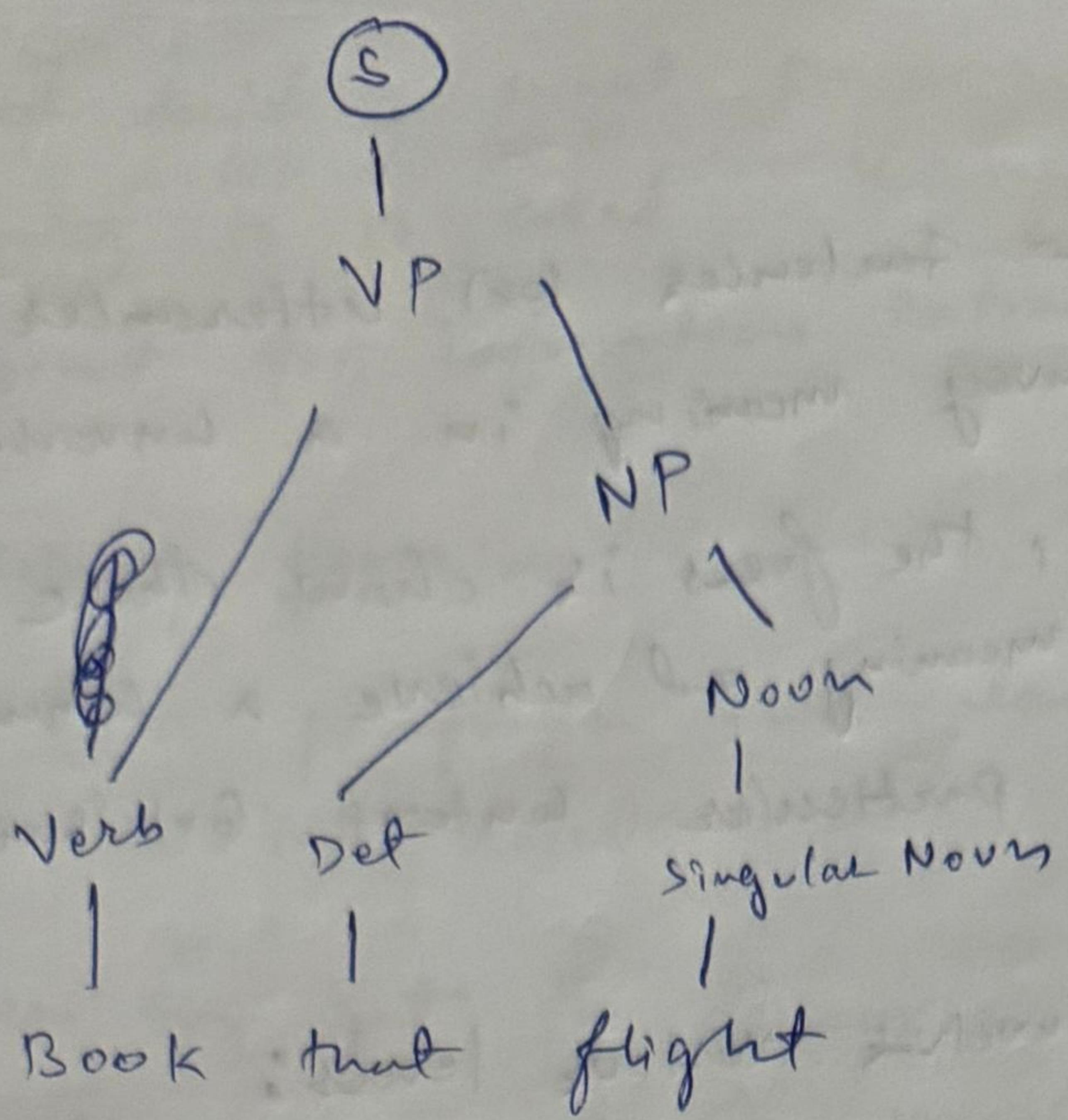
14

1 2 3
(1,1) (2,

(1,2) 13

(1,3) 14

②

Bottom UP# CKY Parsing :-

Q. Check whether a string "abbb" is a valid member of the following CFG

$$S \rightarrow AB$$

$$A \rightarrow BB/a$$

$$B \rightarrow AB/b$$

a b b b
1 2 3 4

Ans

	4	3	2	1
1	S, B	A	S, B	A
2	S, B	A	B	
3	A	B		
4	B			

~~(1,1)~~ 1 2
(1,1) (2,2)

A, B

(A, B)

S, B

2 3
(2,2) (3,3)

B, B

(B, B)

1 3
1 2 3
~~(1,1) (2,2) (3,3)~~

(1,1) (2,3)
A A = Ø

3 4

(1,2) (3,3) (3,3) (4,4)

(S, B) (B) B B

(S, B) (B, B) (B, B)

A

Ø A

Ø Ø Ø U A = A

A B AD

⇒ S, B

1 4
1 2 3 4 (Ø) (S, B) (2, 4)
(1, 1) (2, 4) ⇒ (A) (S, B) 2 3 4
(1, 2) (1 3, 4) ⇒ (S, B) (A) (2, 2) (3, 4)
(1, 3) (4, 4) ⇒ (A) (B) B A ATB
 (S, B) (2, 3) (4, 4)
 A B AD

DISCOURSE AND COHERENCE

(a) Discourse:-

Discourse refers to how sentences (or) utterances are connected and organized to convey meaning in a conversation (or) text.

- * In discourse analysis, the focus is how these sequences of language convey meaning and achieve a communicative purpose within a particular context (or) social setting.

Structure of discourse involves various levels:

(i) Cohesion:-

The linguistic devices used to link sentences and clauses within a text, such as pronouns, conjunctions, and lexical ties.

(ii) Coherence:-

The logical and semantic relationship that allows the text to be understood as a whole

(iii) Discourse markers,

words (or) phrases like "however", "furthermore" (or) "for example", which guide the reader (or) listener through the discourse, indicating relations b/w segments of discourse

(b) Coherence)

- * Coherence is about underlying connections that make a text (or) speech understandable and meaningful

- * It is the way components of the discourse fit together logically and semantically.

coherence

The goal
phrases

* It is imp

coherent

* If we re
paragraph

* Hence, we
if we
coherently
just a

entity-base

~~Relat~~

Discourse Seg

When we de

discourse /

* Segmentation

- ① info
- ② Te
- ③ D

Two types:-

① Unsuperv

② Supervise

* We ~~use~~

* We mat
mat

D

coreference resolution:

- * The goal of deciding what pronouns and other noun phrases refer to is called "coreference resolution".
- * It is important for information extraction, summarization.

coherent discourse texts:

- * If we read a paragraph, we can see that entire paragraph is interrelated.

- * Hence, we can say that the discourse is coherent, but if we only combine the newspaper headlines consecutively, then it is not discourse, it is just a group of sentences that is non-coherent.

entity-based coherence:

~~Relationship~~ Relationship b/w entities, is also coherent

Discourse Segmentation:-

When we determine the types of structures for a large discourse, we term its segmentation.

- * Segmentation is difficult, but necessary as:-

- ① info Retrieval
- ② Text summarization
- ③ Information extraction.

Two types:

① Unsupervised Discourse Segmentation

② Supervised Discourse Segmentation

- * We ~~not~~ deal with labeled boundary
- * We make use of discourse markers to do the segmentation

→ linear segmentation

→ classifications of
openly of
similar texts
with help of coherent
discourse

Discourse Markers :-

① Transition Words and Phrases :-

words like "however", "despite" and "as" indicate shifts in the narrative

② Punctuation :-

Periods(.) are often used to denote end of sentence and hence end of Segment. Not ~~all~~
all periods denote end of Segment.

③ Temporal and Spatial indicators :-

Phrases like "this time", "soon enough" and "as we" provide temporal (or) spatial context, which can signal the beginning (or) continuation of new segment.

④ Contextual Cues :-

The change of topic from one subject to another signifies significant cue for Segment boundary

Text Coherence :-

* To achieve coherent discourse, we must focus on coherence relations in specific

* As we know that coherence relation defines the possible connection b/w utterances in a discourse.

① Reunit

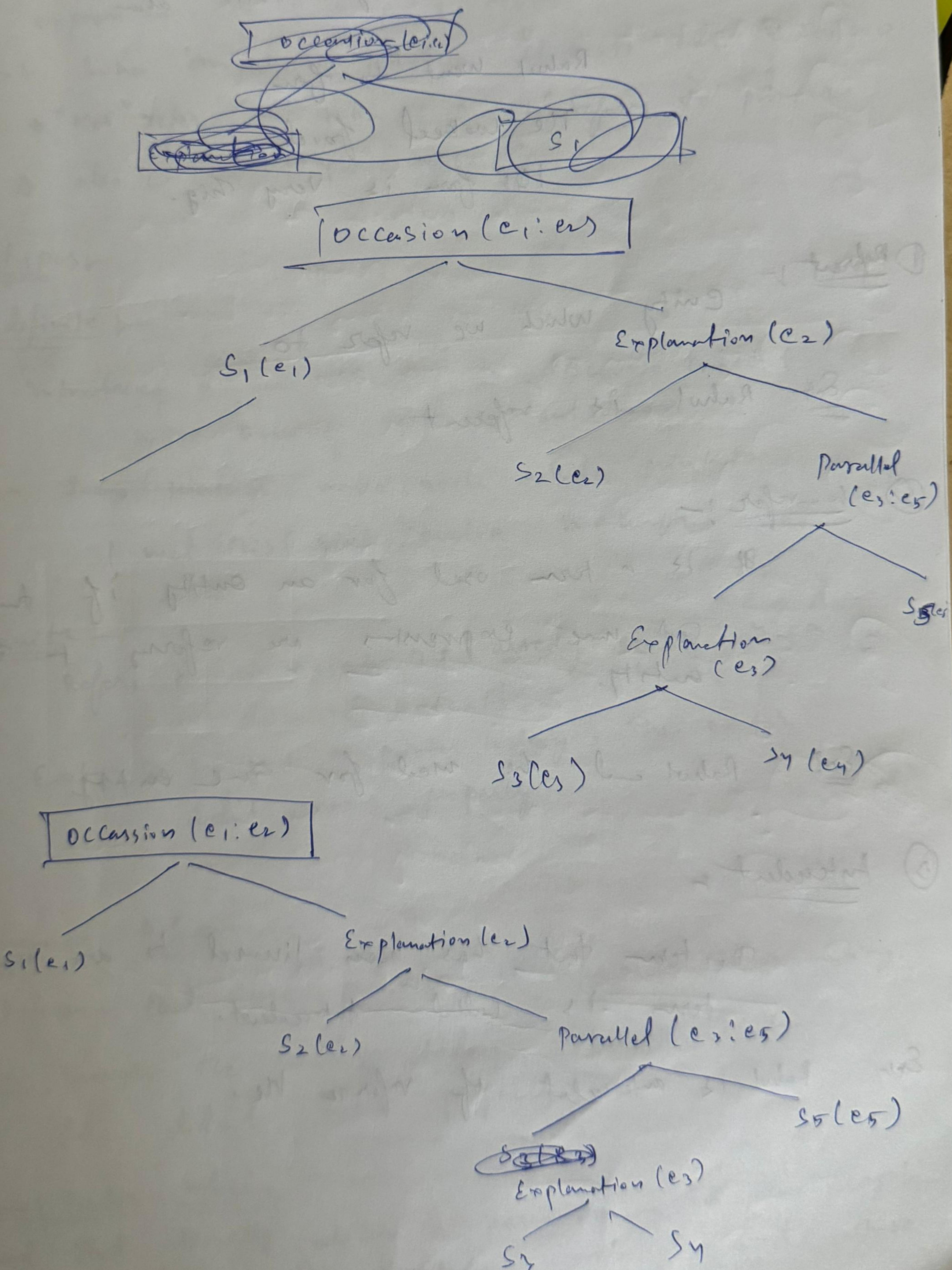
② Explainer

③ Parallel

④ Elevation

⑤ occasion

- s_1 = John went to bank to deposit his check
 s_2 = He then took a train to Bill's car dealership
 s_3 = He needed to buy a car
 s_4 = The company he works for is not near
 an public transportation
 s_5 = He also wanted to talk with Bill about
 soft ball team



Reference Resolution

- * The extraction of meaning (or) interpretation of the sentences of discourse is one of the most important tasks in natural language processing, to do so, we need to "what (or) who the entity" is.
- * ~~Def~~ Reference resolution means understanding the type of entity that is being talk about

Rahul went to farm
He cooked food
His farm is very big.

① Referent

Entity which we refer to

Def Rahul is referent

② Non-referent

If it is a term used for an entity if two or more expressions are referring to same entity.

Ex: Rahul and he used for the entity

③ Antecedent

The term that has been linked to the another term is called Antecedent.

Ex:- Rahul is antecedent of refernece He.

④ Answer

Reference

Five

④ Final

Ex:-

② Definition

⑤ Proper noun

Ref

Ex:-

⑥ Demonstrative A word

① Anaphora and Anaphoric:-

can be defined as term or reference used for an entity that has previously been introduced in the same sentence.

Reference :-

Linking a referring expression to another referring expression in surrounding text.

Ex:- Suhu brought a printer. It costs her ₹ 20,000
→ "Her" refers to Suhu and "it" refers to printer.
Also called "anaphoric reference".

Five types

① Indefinite :-

introduces a new object to discourse context.
most commonly with a/an some.

Ex:- Some printers make noise while printing
I met this girl earlier in course

② Definite :-

Refers to an object that already exists in a discourse context.

Ex:- I brought a printer today. The printer
didn't work properly

③ Pronominal Reference :-

Refers that use pronouns to refer to some entity.

Ex:- I brought the printer today. On installation,
it didn't work properly

④ Demonstrative Reference :-

A word that directly indicates a person / thing or few people
or few things. The demonstrative words are that, those

Ever think it's my bottle, live we that live
water bottle.

(5)

Ordinal Reference:

Use an ordinal like first, one etc..

Earlier I visited a computer store to buy a printer.

I have seen many and now I have to select
one.