

mid sem.

NLP

### core challenges of NLP

Ambiguity  
(ek shabd do arth)

contextual understanding

Pragmatic  
(tone of sentence)

low resource language  
(lack of training data)

ethical concern  
(Bias in training data)

- Lexical  
(1 word 2 meaning)
- Syntactic  
(1 word multiple grammar)
- semantic  
(contextual meaning)

Multilingual  
(1 sentence  
2 language)  
e.g. Hindi-English

(challenges ka positive lih do).

### Future of NLP

zero shot

decision making

Multilanguages

Integration  
w/ other models

no dataset bias

### Language modeling (Backbone of NLP).

- It predicts next word kya h

?

#### Grammar

- It uses grammar rules of the language.

Statistical model.

- Based on probability.

#### Models

1) Unigram - each word independent P

2) Bigram - one prev. word P

1) Top down - start w/ sentence ka root word

2) Bottom up - small units

#### Parsing Technique

4) break into small unit

$$S \rightarrow NP VP$$

2) Bottom up - words ko jod ke  
sentence banata hai

$$NP VP \rightarrow S$$

3) CYK Algo

Dynamic prog. use karta hai  
for parsing CFG.

3) Trigram - 2 prev. word

P

4) Neural language - uses deep learning arch.  
eg - RNN  
limit - huge computational resources.

## Probabilistic approaches

1) MLE -

PC

### smoothing technique

1) Add one - Adds 1 to word count  
(Laplace)

2) Good-Turing - Adjust probability for unseen words.

3) Back off

4) Interpolation

## Evaluation metrics

### ① Perplexity (PPL)

- Measures how well a model predicts seq. of words

$$\text{Performance} \propto \frac{1}{\text{PPL}}$$

### ② Cross Entropy

- Measures avg bits to encode a message.

## Finite state Automata (FSA).

- Abstracted computational model. used to recognise patterns

State ( $Q$ ) - set of finite states

Alphabet ( $\Sigma$ ) - input alphabet

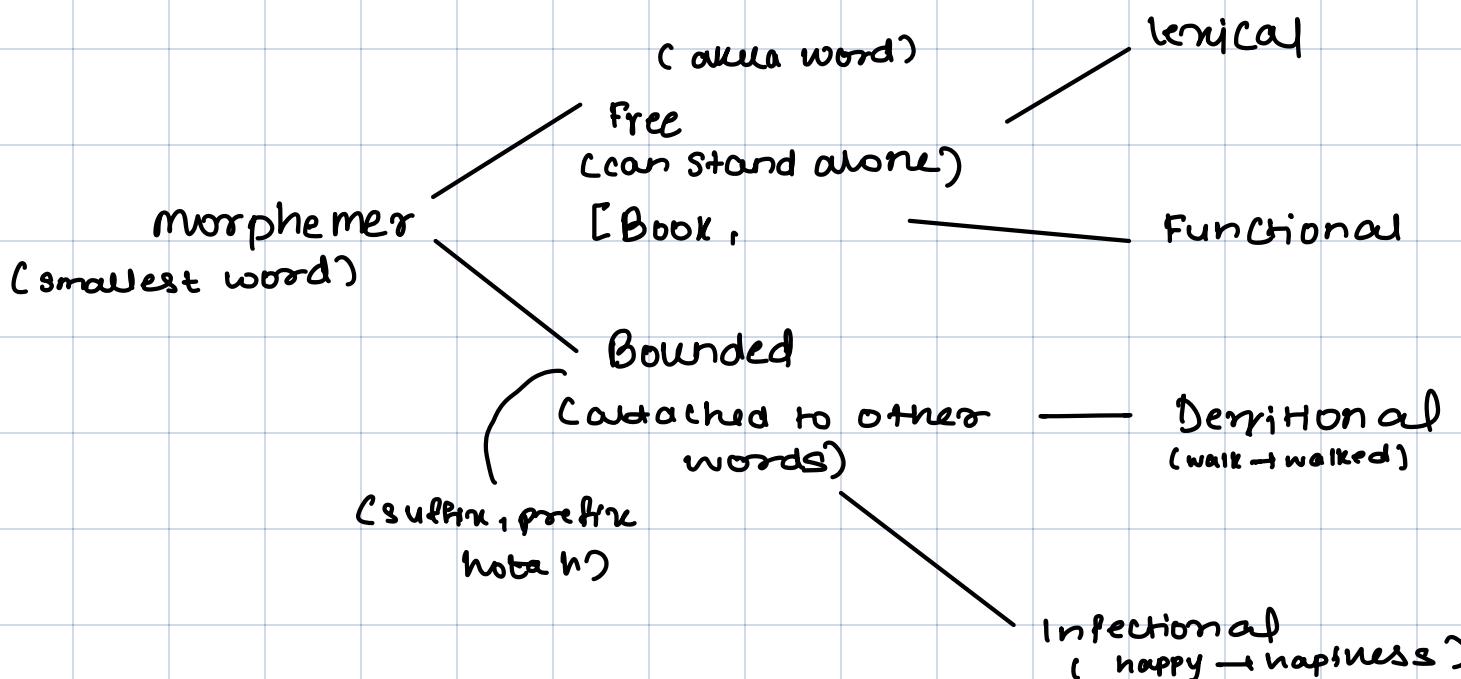
start state ( $Q_0$ ) - initial state

Accepting state ( $F$ ) - Accept the string (validate)

Transition function ( $\delta$ ) - Rule defining b/w transitions.

## English morphology

- focuses on structure & formation of words



morphological parsing - Analysing words into morphology

- Tokenisation - text → words
- morphological segmentation - words → morphemes
- morpheme Tagging - identifying morpheme & root
- Reconstructing meaning - combining meaning of morphemes.

## Tools

- Porter stemmer
- Spacy
- Lancaster stemmer

## challenges .

- Ambiguity
- complex
- irregular form

## Finite state Transducer

( FSA - accept/reject input/output  
FST - maps output to input )

### FST structure

- Q - start ka set
- $\Sigma_{in}$  - input alphabet
- $\Sigma_{out}$  - output alphabet
- $q_0$  - start state
- F - final state
- $\delta$  - Transition function

### Transition notation

$q_i$ : current state  
 $a$ : input symbol  
 $b$ : output symbol  
 $q_j$ : Next state

### Tokenisation

(Text to small units mein break karta h)

### Types

- word - Text to word [ "NLP is cool"  $\rightarrow$  "NLP" "is" "cool" ]
- sentence - text to sentence [ "Hi. How are you"  $\rightarrow$  "Hi , "How are you" ]
- subword - word to word [ "unhappiness"  $\rightarrow$  "un" "happiness" ]
- character - text  $\rightarrow$  individual char. [ "NLP"  $\rightarrow$  "N" "L" "P" ]
- special - special symbol handle karta hai [ "visit www... com"  $\Rightarrow$  "visit" "www" "com" ]

### Techniques

- Rule - Pre defined rules pe tokenize karta hai
- Statistical - Training data ke basis pe probability
- Neural - DL na use karta hai

## FSA

Accepts / rejects strings  
only yes/no

state + transition

spell checking . pattern  
matching

e.g. check if "cat" is valid  
- state

(current, input) → next-state

## FST.

converts i/p string to o/p  
produces new o/p.

state + transition + i/p, o/p pairs

morphological analysis (breaking word)  
text translation

e.g. change running → "run+ing"

(current.state, i/p:o/p) → next-state

## NLP Pipeline

Raw Text <sup>removal</sup>



Text cleaning → lower casing → Tokenisation → Removing stop words



Dependency NER → POS → Lemmatisation → Stemming  
Parsing → Tagging

## Vectorisation

[ just like Raw → Cleaning → Writing → Cooking → Tasty ]  
e.g. <sup>food</sup>

Regex → (like search pattern)

• Helps find, match, replace  
words / numbers / patterns

+ <sup>n</sup> - re.search(), re.findall(), re.replace()  
first match all matches replaces

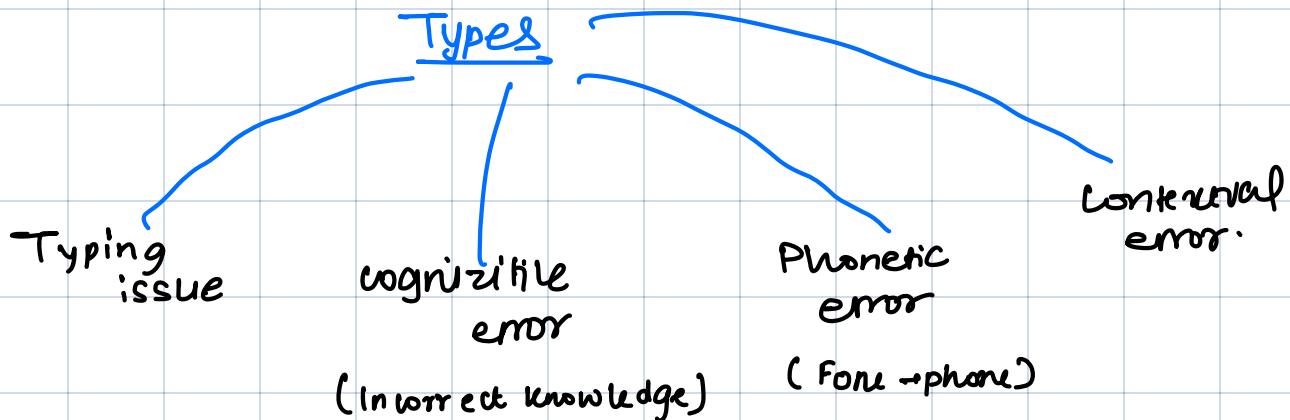
symbols - . ^ \$ \* + ? { } [ ]  
any char start of string end of string repeat or more repeat or more specific any char no-of-rep char.

Eg - import re  
match = re.search ('cat', 'I have a cat') <sup>y</sup> To find 1 word  
print (match.group())

matches = re.findall ('a', 'banana') <sup>y</sup> To find all occurrences

## Error detection.

Identifying wrong word & suggesting right word.



### Techniques for error detection

- Dictionary - pre defined dictionary (contextual error detect nhi kar payega)
- N-Gram - word sequence ka pattern check karta hai
- Edit distance - minimum possible changes ka count karta hai
- Neural - DL models such as transformers
- Probabilistic - Error ko ek noisy commun. channel ke through generate karta hai.

### Technique for error correction

- MED - ham changes karta hai
- contextual - Nearest words ki correction hisab se
- Noisy - Bayesian inference channel model ka use karte correction
- Neural spell - Transformer, correction - seq2seq jaise DL kein use karta hai

(Group of N-words)

## unit-2.

N-Gram (Help in predicting the next item)

unsmoothed

Smoothed.

(It can yield 0 probab.  
for unseen events)

(Improve by redistributing  
probability)

$P(\text{Phone} | \text{Please turn off your cell})$

uni -  $P(\text{Phone})$

Bigram:  $P(\text{Phone} | \text{cell})$  Tri ( $\text{Phone} | \text{your cell}$ )

How to handle unseen probability :-

Smoothing.

(Agar koi data training mein ho aur wo testing aa jaye  
toh wo unseen probab. nota hai)

D) Laplace (Add-on)

- Has word ki frequency mein +1 add kiya jata hai.

Good Turing - Adjust probability of unseen events.  
Kneser-Ney - Probability redistribute across n-gram levels.

Back off

- Agar Trigram na ho paye, toh bigram / unigram use

- Agar Bigram na ho toh unigram

(Basically back off krke ek level niche)

## Interpolation

- combine probability from tri, bi, unigram to get a better guess.

**POS tagging** (Assigns word classes to token for syntactic analysis).

↳ Helps us understand sentence structure, disambiguate word meaning.

- \* penn tree-bank common POS tagging technique.

Appm - machine translation.

Method of POS tagging .

① Rule Based - Pre-defined grammar rules

② stochastic - Probability-based models

③ Transformation - rules & stochastic ka combination

# POS category

- ① Noun , Verb , Adjective , Adverb  
proper common Main Auxi  
larry
- ↳ open  
(new word  
add no suffix  
nai)
- ② Preposition , Pronoun , Article , conjunction  
↳ closed  
(fixed words)

Rule based PBS tagging

Stage 1 :- uses dictionary

Stage 2 :- uses Large list of handcrafted rules.

## Tagging rule

Rule 1 - Agar previous word article ho toh saare verb tag nikalo

Rule 2 - Agar next verb tag ho saare verb nikalo.

## Hidden Markov model (HMM). (Generative model)

- statistical model used for sequences.

① Hidden state - POS Tags (hidden) observe layer se

② Observation state - (visible) hidden ko predict karne  
has

Parameter  
③ Transition - one state to another  
④ Emission - Hidden se observation ke chances  
⑤ Initial - starting state ki prob. distribution  
emission prob.

Transition prob.

$$P(\text{tag at } t \mid \text{tag at } t-1)$$

$$= P(t_i)$$

$$P(\text{word at } t \mid \text{tag at } t)$$

$$= P(w_i)$$

Goal - Find the most likely sequence of tags given the words

Maximum entropy markov model. Assigns probability without any bias & only on known data.  
 $y_{\text{model}} = P(\text{tag} \mid \text{context})$

\* No strong independence assumption like in HMM.

## Advantage

- Can use many overlapping features
- Better for complex lang. patterns

Eg -

If previous word 'the' next word likely noun.

## Appn

App

- 1) Text Analysis
- 2) Speech recog
- 3) Text to speech
- 4) Machine translation
- 5) DNA Analysis
- 6) Music

### 3 important problem in HMM -

10.

- decoding (used to find hidden state)
- ① Viterbi Algo - uses DP #  $O(N^2t)$  time complexity
- ② Initial probability assign
- ③ Transition & emission probab calculate
- ④ Best path find out Urkle output seq find out karo.
- ⑤ (Likelihood) ✓ (Forward Algo) /  $TN^2$   
 (used to est. probab. of observation sequence).  
 # calculate total probab. of different possible state sequence of HMM.

### ⑥ forward | Backward Algo (Learning).

# uses supervised L  
 ↳ (we use MLE)

steps

① Initial prob est.

② calculate transition / Emission probab.

③ Reestimate prob.

④ Repeat till convergence

# unsupervised huatok

it is true HMM case

↳ (we use forward | Backward algo).

## Unit-3

(CFU context free grammar)

- set of rules that help build sentences (like a recipe)

(replace homage )      S - sentence      VP - verb phrase

↖ Non terminals - NP - noun phrase

Actual words - cat etc

↙ Terminals

(final word) Production rules - Rules that show non terminals break

$S \rightarrow NP\ VP$  (noun followed by verb)

$NP \rightarrow Det\ N$  (noun phrase made of determiner & noun)

Treebank

syntactic - syntactic focus kinda has

semantic - Focus on word meaning & relationship.

- collection of sentences that've been parsed into tree structures

. Penn Treebank is a famous treebank . it contains thousands of sentences, each drawn as a tree with NP, VP etc.

## CNF

- rewriting CFU so that every production rule is simple

①  $A \xrightarrow{=} BC$  : A non terminal produces two non terminal

②  $A \rightarrow \alpha$  : A non terminal produces a terminal

\* useful because it makes rule binary & predictable

How to convert CFG to CNF

(  
eliminate  
empty space)

remove  
unit productions

RHS ka start sym  
nikalo  
null value nikalo  
a terminal nikalo  
Non T  $\rightarrow$  T

\* Binorize the  
rules

\* Replace mixed  
rules

Eg:-

CFG

$S \rightarrow \$ a A B$

Useless  
null  
empty  
Terminal  
long RHS

$R \rightarrow \alpha$   
 $X \rightarrow AB$

$S \rightarrow R X$   
CNF

(mix rule)

1)  $R \rightarrow a$      $S \rightarrow R A B$

(converting terminal to non T)

2)  $X \rightarrow AB$      $S \rightarrow R X$   
G CNF

(Binorizing)

CNF

(Kreibach normal form)

Every rule must start with terminal by non terminals

$A \rightarrow \alpha B$   
Terminal

why better? - Every sentence starts with a real word.

e.g

$A \rightarrow BC$

①  $A \rightarrow a D$   
Terminal      (BC) (or any other Terminal)

CYK Alg - Refer to ppt

## \* Parsing Techniques

### Top-down Parsing -

start from S to generate i/p string

root  $\rightarrow$  terminal

S  $\rightarrow$  NP VP

NP  $\rightarrow$  det N      VP  $\rightarrow$  NP

Bottom up parsing - start from the input string  
Build tree from word upward.

Det, N, V

Dynamic Parsing. After parsing result store it in mem.  
to avoid redundant computation.

Shallow Parsing. (chunking) - Build parse tree ke  
sikha extract meaningful info. (Fast)

Parsing

Lexical Parsing - Accurate (Because ye direct lexicons  
ko parse karta hai).

Collins Parser - Probability parser using **lexicalised**  
**CRF rule**. Enhance - Enhancing accuracy.

# CKY Parsing Algorithm

## Example: Parsing a Simple Sentence

We have the following context-free grammar (CFG) in Chomsky Normal Form (CNF):

$$\begin{aligned} S &\rightarrow NP VP \\ NP &\rightarrow Det N \\ VP &\rightarrow V NP \\ Det &\rightarrow "the" \\ N &\rightarrow "cat" \mid "dog" \\ V &\rightarrow "chased" \mid "saw" \end{aligned}$$

CKY ALGO

It is a bottom up algo  
Requires grammar in CNF.

## CKY Parsing Algorithm

### Input Sentence

Let's parse the sentence:  
"the cat chased the dog"

### Step 1: Construct the CYK Table

We create a triangular table for the sentence:

Word	the	cat	chased	the	dog
1	Det	N	V	Det	N
2	NP	-	VP	NP	-
3	-	S	-	S	-
4	-	-	S	-	-
5	-	-	-	S	-

## Probabilistic CKY Algorithm

Probabilistic CKY ALGO

### Step 1: CYK Table Initialization

Each cell  $(i, j)$  stores:

- Possible non-terminals
- The maximum probability of generating the substring

Word	John	chased	the	dog
1	NP (0.2)	V (0.6)	Det (0.6)	N (0.5)
2	-	VP ( $0.7 \times 0.5 = 0.35$ )	NP ( $0.8 \times 0.5 = 0.4$ )	-
3	-	-	S ( $1.0 \times 0.2 \times 0.35 = 0.07$ )	-
4	-	-	-	S ( $1.0 \times 0.2 \times 0.7 \times 0.4 = 0.056$ )

### Step 2: Parse Tree Extraction

We extract the best parse tree using backtracking.

$$S \rightarrow NP(John) \quad VP(V : chased, NP(Det : the, N : dog))$$

## CKY Algorithm

### Grammar in CNF:

- $S \rightarrow NP VP$
- $NP \rightarrow Det N$
- $VP \rightarrow V NP$
- $Det \rightarrow "the" \mid "a"$
- $N \rightarrow "cat" \mid "dog"$
- $V \rightarrow "sees"$

### Probabilistic Grammar (PCFG):

- $S \rightarrow NP VP (1.0)$
- $NP \rightarrow Det N (0.9) \mid NP \rightarrow N (0.1)$
- $VP \rightarrow V NP (0.85) \mid VP \rightarrow V (0.15)$
- $Det \rightarrow "the" (0.6) \mid Det \rightarrow "a" (0.4)$
- $N \rightarrow "cat" (0.5) \mid N \rightarrow "dog" (0.5)$
- $V \rightarrow "sees" (1.0)$

	The	Cat	Sees	a	Dog
1	Det(0.6) )	N(0.5)	V(1.0)	Det(0.4) )	N(0.5)
2	NP(0.9 *0.5)		V(1.0)	NP(0.9 *0.5)	
3	NP(0.9 *0.5)		VP(0.8 5 *0.5)		
4			S(1.0 *0.9*0. 5*0.85 *0.5)		
5					S

The cat sees a dog.

## Dependency Grammar

- word-to-word relation not phrase
- each word depends on another word

eg. "Economic news had little effect"

sub verb

? obj

had → verb → it is the root. All others depend on it.

## Ambiguity -

when one sentence has multiple parse trees

For - Relationship b/w two things

unit - 14

\* Shreyan likes coffee

Shreyan(x)  $\Rightarrow$  likes (x, coffee)  
relation

(All Girls) like coffee.

( $\forall x$ ) Girls(x)  $\Rightarrow$  likes (x, coffee)

## Unit - 4

logical components -

$\wedge$   
Symbol  
 $\vee$

AND  
OR  
b connecting

NOT

IMPLIES

IF AND ONLY IF

FOL - Nikhil(x)

$\Rightarrow$  Likes(x, ted)  
FOL about

LHS

RHS

$\forall x \text{ Boys}(x) \Rightarrow \text{Likes}(\underline{x}, \underline{\text{ted}})$   
 $\exists x \text{ Boy}_3(x)$

DOL

## First Order Logic (FOL)

### Components

1) constants

specific objects

(John, Appu)

2) Variables

represent arbitrary obj

(u, v, z)

3) Predicates

Relationships b/w them

Human(John)  $\rightarrow$  John is a human.

### 4) Functions

Function takes i/p & give o/p

Father(John) = David (David is

John's Father)

### 5) Quantifiers

1) Universal ( $\forall$ ) (For all)

2) Existential ( $\exists$ ) (For some)

$\wedge \vee \rightarrow \leftrightarrow$

some boys are intelligent

$\exists x : \text{boys}(x) \wedge \text{intelligent}(x)$

every gardener likes the sun

$(\forall x) \text{gardener}(x) \rightarrow \text{Likes}(x, \text{sun})$

{ some } object { and }

All obj

Implies for { objects }

You can fool all of the people some of the time

$(\exists t) (\forall x) \text{person}(x) \wedge \text{time}(t) \Rightarrow \text{can-be-fooled}(x, t)$

$\wedge x \quad \forall t$

## Requirement of Knowledge rep in NLP.

- expressiveness
    - ↳ complex ideas easily.
  - efficiency
    - ↳ reason logical deduction
  - clarity
    - ↳ clear
  - scalability
    - ↳ large data
  - consistency
    - ↳ no contradiction allowed
  - Modularity
    - ↳ knowledge can be updated easily.

FOL for semantic meaning

Natural lang are ambiguous ↗ multiple meanings

FOL gives a mathematical structure

↳ helps machine understand & reason better.

## Thematic agents

## Syntax - driven semantic analysis

- assigning meaning to sentence based on syntactic structure

"John saw many with book"  
we know who has the book

# Semantic Attachment

- Attach semantic meaning  
Buy → trigger financial decision

## semantic

- meaning of words
- deontic verbal underst.
- . Bank → financial inst.
- . word meaning
- used in machine trans.

## Pragmatic

- Based on context
- conversational understanding
- can you pass the salt?  
Urge, not a ques
- Used in chatbot

## Reasoning + Winograd

↳ process of  
drawing conclusion  
logically

↓  
need for machine  
to think like human

↓  
rule      } winograd  
cause  
effect

↳ solves ambiguity using real world knowledge  
not just patterns

✓ Test for common sense in AI.  
Harder than tuning.

## Description Logic (DL)

- used to combine aspects of both FOL & set theory.

## Key Features

### ① Concept (classes)

- collection of individual objects
- Relationship b/w ind.

### ② roles

- (Terminological Box)  
e.g. hasChild,  
worksAt

### ③ Axioms

- (Humans subclass of Animal)

Existential quantif. ( $\exists$ )

$\exists \text{ } R \cdot C$   $\rightsquigarrow$  class

$\exists$  haschild. human

universal quantif.

$\forall R \cdot C$

$\forall$  haschild. human

(FOL & Set theory combination)

DL

concept

(classes)

(

collection  
of objects

role  
(relationship)

properties

Individual  
Object)

1  
relationships

# Axioms,

T-box - Rules & definitions

A-Box - Assertions

## First order L

- General rep. of facts & rules
- Very expressive
- complex, hard to automate
- uses variables, quantifiers
- $\forall x (Dog(x) \rightarrow Animal(x))$
- Harder to decide
- General AI, expert sys

## Description L

- Special knowledge for concepts & relationship
- less expressive
- simpler, optimal.
- uses classes, roles, individuals
- $Dog \sqsubseteq Animal$
- easier, automated
- semantic web

## Unit-5

→ splitting text into meaningful sections

Discourse segmentation. (A grp of sentences that are logically connected)

- each sentence = a coherent unit

e.g. = news article → "intro, event details"

### Challenges

- Hidden Boundaries
- Shifting Topic
- Different writing style
- Ambiguous sentence
- Context Understanding.

### case study. (fake news detection)

Problem: Automatically detecting fake news.

- challenges faced:
- lack of progression
  - mix unrelated events
  - contains contradicted limited context

example: " vaccine hr side effects . Also a scientist in Norway discovered Alien DNA in vaccine".

→ not logically connected . It requires discourse analysis -

Coherence .

It is used to test the output quality.

- It should show relation b/w utterances
- Relationship b/w entity.

## Discourse Segmentation

① unsupervised

② supervised

### ① Unsupervised Discourse

- no levels
- finds topic change naturally
- uses lexical cohesion
  - related words
  - sticking together

### ② Supervised

- Labelled training data.
- It shows where exactly topic changes.
- use cue words
  - "however" "because"

## Text coherence

Result - "Rahul is late. He will be punished"

Explanation - "Rahul fought w/ friend. He was drunk"

Parallel - "He wants food. She want money"

Both want smth at same time

Rohan is from delhi. Ram is from delhi"

Both mean smth but no related and give clear explanation.

Occasion - "Rahul took money, he gave to rohan"  
S1 inferred from S0-

## Building Hierarchical Discourse Structure

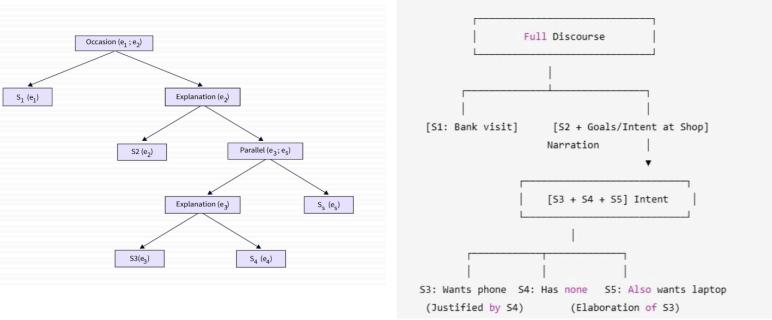


Let us now try to build a hierarchical discourse structure with the help of a group of statements. We generally create the hierarchical structure among the coherence relations to get the entire discourse in NLP.

Let us consider the following phrases and serially number them.

- S1: Rahul went to the bank to deposit money.
- S2: He then went to Ronan's shop.
- S3: He wanted a phone.
- S4: He did not have a phone.
- S5: He also wanted to buy a laptop from Rohan's shop.

Now the entire discourse can be represented using the below hierarchical discourse structure.



## Reference Resolution -

For example, look at the below sentences.

- Rahul went to the farm.
- He cooked food.
- His farm was very big.

Referent - whom we referring to

Antecedant - first time it being mentioned

No-ref - a word referring to same.

## Word Sense

→ words having multiple meanings)

e.g. - Bank → financial  
                        riverside.

why important?

→ correct understanding, improves translation, better search result, accurate sentiment analysis.

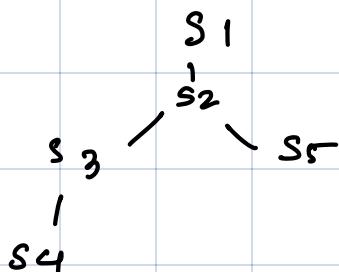
S<sub>1</sub> & S<sub>2</sub>: occasion

S<sub>2</sub> leads to  
↓

S<sub>3</sub> : explanation

S<sub>5</sub> : parallel goal

S<sub>4</sub> explanation S<sub>3</sub>.



(Big Act → detail → sub det → emp)

when we see "he", "it"  
what does it refer to?

He, His → Rahul

How to parse ?

- dictionary . supervised
- unsupervised . semi-supervised

# Tools

## • Porter Stemmer

- stemming algo

=

(chops off suffix)

"-ing", "-ly", "ed"

eg- playing → play  
faster → fast

Very fast but  
not accurate.

## Code

```
from nltk.stem import PorterStemmer
```

```
stemmer = PorterStemmer()
```

```
words = [" " ]
```

```
stems = [stemmer.stem(word) for  
word in words]
```

## • Lemmatizer

- PoS Tagging algo

- reducing word to its dictionary form.

eg- "am" "be" "is" → be  
"running" → run

More accurate but  
slower.

## Code -

```
from nltk.stem import WordNetLemmatizer
```

```
lemmatizer = WordNetLemmatizer()
```

```
words = [ " " ]
```

```
lemmas = [lemmatizer.lemmatize(word,  
pos='v') for word in words]
```

## Wordnet

- Large lexical database of eng
- Group words into synsets (set of synonyms)
- Provide relationships (synonym, antonym etc.).

## Brown corpus

- First text corpora in eng
- 1 million words categorised by topics
- Used for lang modelling
- PoS Tagging

## FrameNet

- Semantic Project
- How word evokes frames
- q. 'commerce-buy'  
includes buyer, seller, goods, price