

Sunanda
Asst. Professor
CSE, DCE, BHU

UNIT- 01.

NLP approaches

(P) Top

1) Rationalist → person who based their approach on reason

2) Empiricist → experience derived from sensations

* Levels of processing

1) Lexical analysis → Analysis of words
Word-level process

2) Syntactic analysis → Sequence of words

3) Semantic analysis → Meaning of the language
Interpret the meaning of
larger units.

4) Discourse Analysis → deals with purposeful use of
sentences in situation

Challenges

* Not possible to write/mimic as done by humans.

* Multiple meaning associated with the same word.

Ex: bank, can, ball, still,

↳ Need to develop models/algo's to resolve them.

* Machine will not understand the feelings / expression of the stats.

* Taj → Indian / others

Ambiguity: beard of tea / hotel / monument

- Applicability:
 ① Machine translation (grammars of lang)
 ② Speech synthesis (with NL interfaces to DB) ③ Information retrieval
 ④ Question answering ⑤ Chatbot.
 ⑥ Information Extraction
 ⑦ Text summarization.

Language Modelling

entity
description of some complex entity
(or) process.

Language Model → description of language.

approaches in

Grammar-based LM Statistical language Modelling.

Grammar-based LM: Uses the grammar of a language to create its model.

↳ Consists of hand-coded rules to define the structure & ordering of various constituents appearing in unit.

Statistical LM: Ex: LFG (Left lexical function grammar) creates LM by training it from a corpus.

* In order to capture regularities of language, training corpus needs to be large.
Ex: CFG

Various grammar-based language Models:

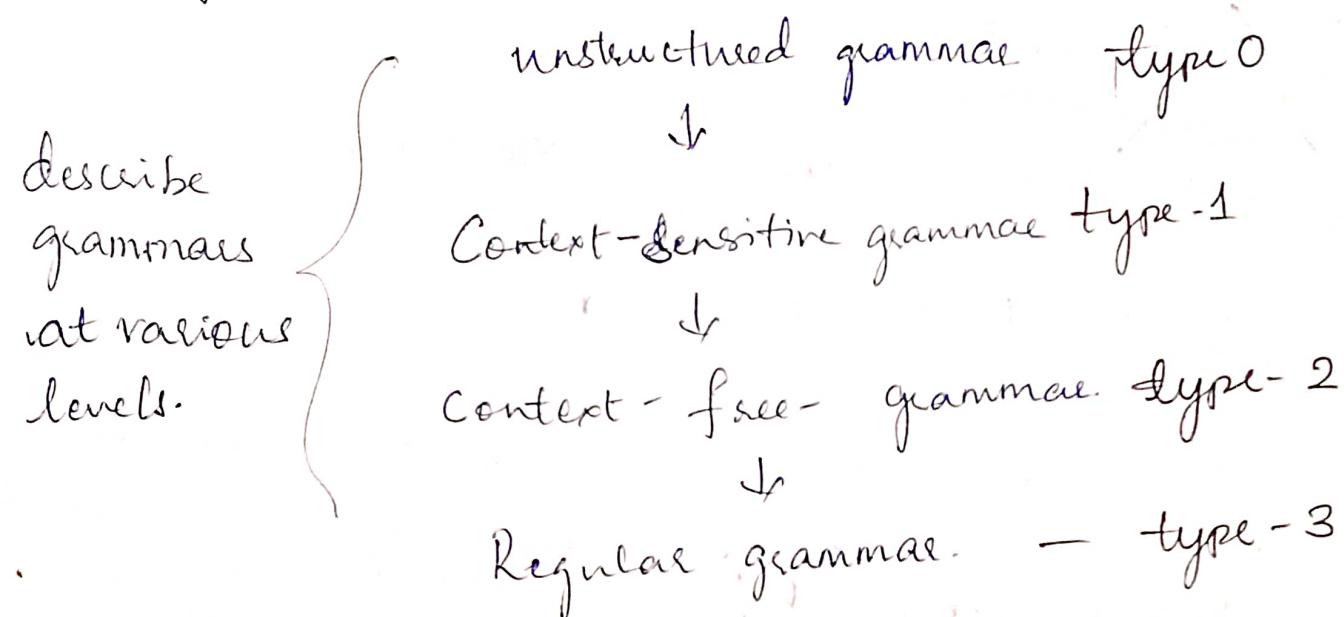
1. Generative grammars:

* generate sentences in a language by known collection of words & rules in that language.

* only those sentences that can be generated as per the rules in the language are grammatical.

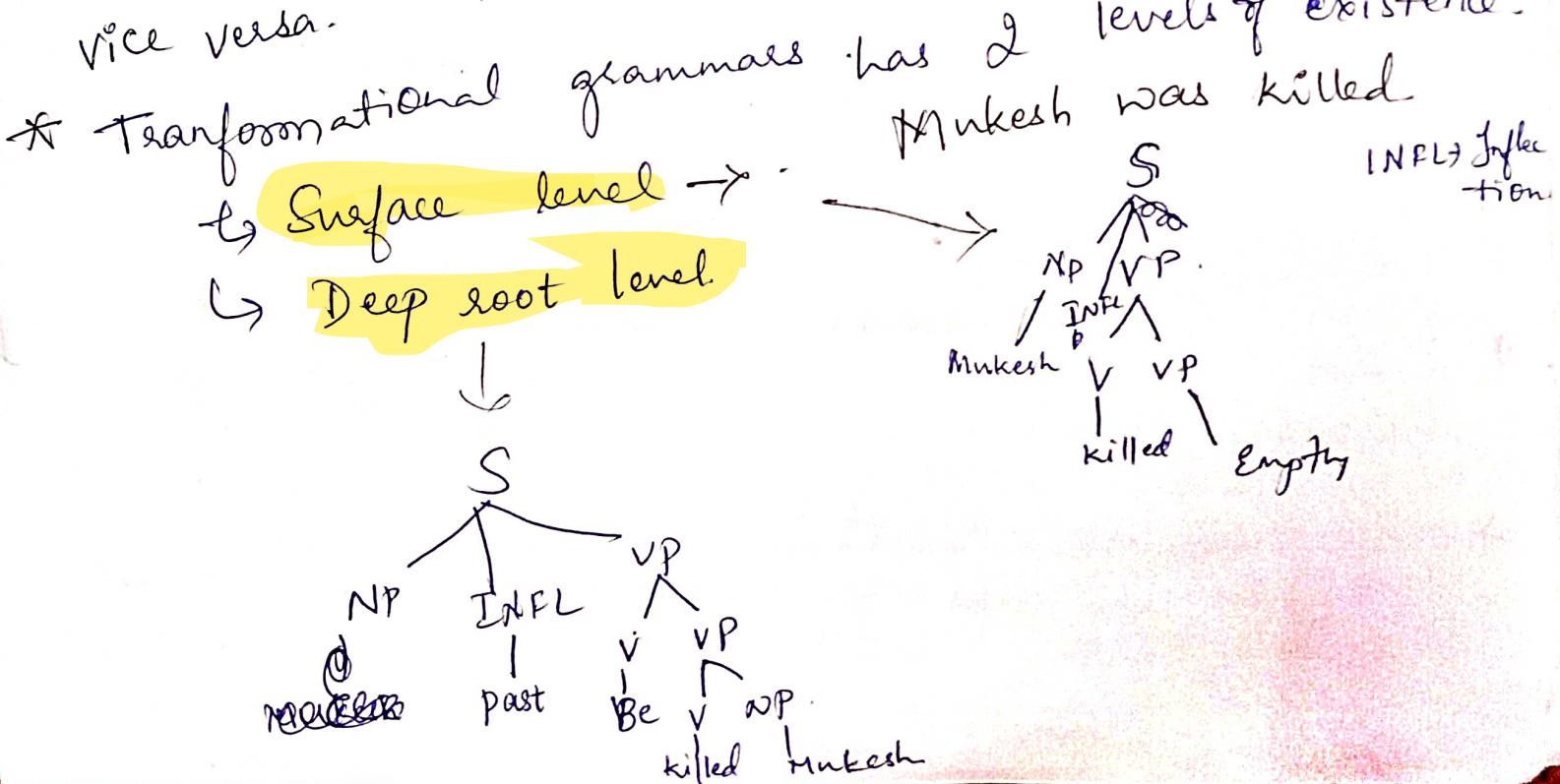
2) Hierarchical Grammars:

- * described classes of grammars in a hierarchical manner where top layer contained the grammars represented by its sub classes.



3) Government & Binding:

- * How well the sentences are formed by its meaning,
- * ~~Sentences~~ Sentences are given at the syntactical level and the transformation from meaning to syntax (or vice versa).



1) Components of GR

1) Lexical Functional grammar (LFG) Model.

- * LFG represents sentences at syntactic level.
 - ↳ Constituent structure (C-structure)
 - ↳ Functional structure (F-structure)
- * LFG is a formal system for grammatical representation in 1982.
- * Lexical Functional grammar.
 - Composed of 2 terms
 - the fact that the lexical rules can be formulated to help define the given structure of a sentence.
 - grammatical functions, such as Subject + object, (S-O) roles played by various arguments in a sentence.
- * C-structure → derived from usual phrase of sentence structure syntax.
- * F-structure → final product obtained from C-structure, by functional specifications are annotated.

Statistical Language Model

- * It is the probability distribution $P(s)$ over all possible word sequences.
- Ex: n-gram model.
- * It is used to estimate the probability (likelihood) of a sentence^(SL) Add-one smoothing

2b

n-gram Model:

- * To estimate the probability of a sentence, decompose sentence probability into a product of conditional probabilities using chain rule:

$$P(S) = P(w_1, w_2, w_3, \dots, w_n)$$

$$= P(w_1) P(w_2|w_1) P(w_3|w_1, w_2) P(w_4|w_1, w_2, w_3) \dots \dots (P(w_n|w_1, w_2, \dots, w_{n-1}))$$

$$= \boxed{\prod_{i=1}^n P(w_i | h_i)} \quad (h_i \xrightarrow{\text{on}} \begin{matrix} \text{history of word } w_i \\ w_1, w_2, \dots, w_{i-1} \end{matrix})$$

- * To get sentence probability, we need to calculate the word probability which is preceding it in a sequence.

- * The ~~time~~ history (h_i) to the previous one word only.

↳ bi-gram ($n=1$)

* previous 2 words → tri-gram ($n=2$)

$$P(S) \approx \prod_{i=1}^n P(w_i | w_{i-1}) \rightarrow \text{Bi-gram}$$

$$P(S) \approx \prod_{i=1}^n P(w_i | w_{i-2}, w_{i-1}) \rightarrow \text{Tri-gram}$$

- * Ex: $P(\text{east} | \text{The Arabian Knights are fairy tales of the})$

↳ $P(\text{east} | \text{the})$ — Bi-gram

↳ $P(\text{east} | \text{of the})$ — tri gram

- * pseudoword ($\langle s \rangle$) is introduced to mark the beginning of the sentence

→ Bi-gram $\langle s_1 \rangle$ · tri-gram → $\langle s_1 \rangle \& \langle s_2 \rangle$.

Word level Analysis

One way to analyse natural language text is by breaking it down into words, identifying morphological variants, detecting & correcting misspelled words & identify correct part-of-speech of word.

Regular Expressions

pattern-matching standard for string parsing & replacement.

- * They are useful tools for the design of language compiler & used for tokenization; describing lexicons, morphological analysis etc..
- * R.E are case sensitive. $[A-Z]$ / $[a-z]$ is different.
- * Single symbol expression $\rightarrow /a/$ [set contains 'a' match.]

RE	Match	Ex
$[abc]$	Match any of a, b & c	Refresher course will start from
$[A-Z]$	Any character b/w A & Z	the course end in Dec 20
$[^A-Z]$	Any char. other than Upper letter	TREC Conference
$[abc]$	any char. not a, b, c	• NLP is good.
$[a^]$	not a (on 'a')	↑ has 3 diff uses.

RE

Match

Example

* [Kleene *]

Zero or more occurrences of previous char

 $|b^*| \rightarrow b, bb, bbb$
 $|ab|^*| \rightarrow aa, ab, abab, ababab, ...$

+ [Kleene +]

One or more —||—

 $|a^+| \rightarrow aa, a, aaa, aaaa, ...$

?

Exactly one occurrence of previous char

/example?/ → example examples.

• (dot)
[Wildcardchar]

Any single char.

/•.al/ → cal, bat, mat
...

^ [Caret]

Match at the beginning of a line

/^The nature \\$/
↳ exactly for the line

\$

Match at the end of line

/\\$The nature --- - - -

Some special characters

RE	Expansion	Description	Example
\d \D	[0-9]	any digit	Party of 5
\. D	[^0-9]	any non-digit	Blue Moon
\w	[a-zA-Z0-9_]	any alphanumeric space	Daiyu
\W	[^\\w]	non-alphanumeric	!!!
\s	[\\t\\n\\f]	white space (space; tab)	
\S	[^\\s]	non-white space	in Concoed
\	[*\\+\\?\\\$\\]	escape special characters	what is this?

Finite State Automata

(A)

- * Machine has no. of states are finite & the alphabet of input symbol is 'finite' [limited/not infinite]
- * Machine moves automatically i.e., the change of state is governed by the input.
- * FSA recognize the regular languages represented by regular expressions.
- * Ex Sheep Talk : /baat!/
- * FSA is a 5-tuple

Q : set of states $\{q_0, q_1, q_2, \dots, q_n\}$

Σ : alphabet of symbols $\{a, b, !, .\}$ (set of input symbols)

q_0 : A start / initial state

F : set of final states in Q $\{q_f\}$

$\delta(q, i)$: Transition function mapping $Q \times \Sigma \rightarrow Q$

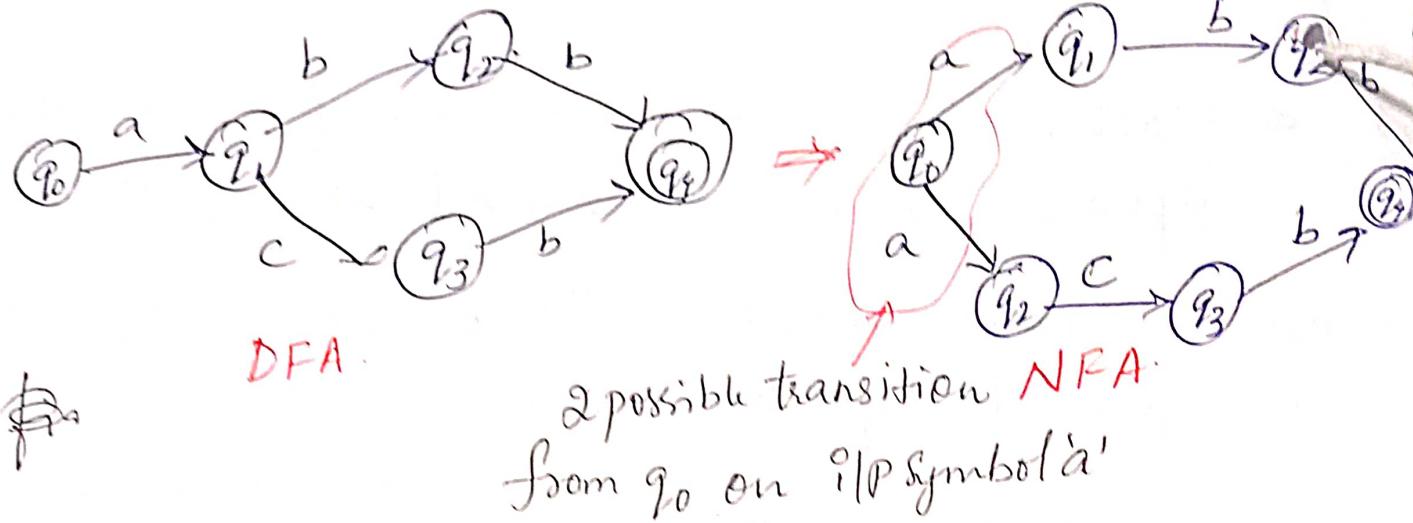
* FSA can be deterministic (or non-deterministic) (NFA)

DFA: Only one transition possible from a state on

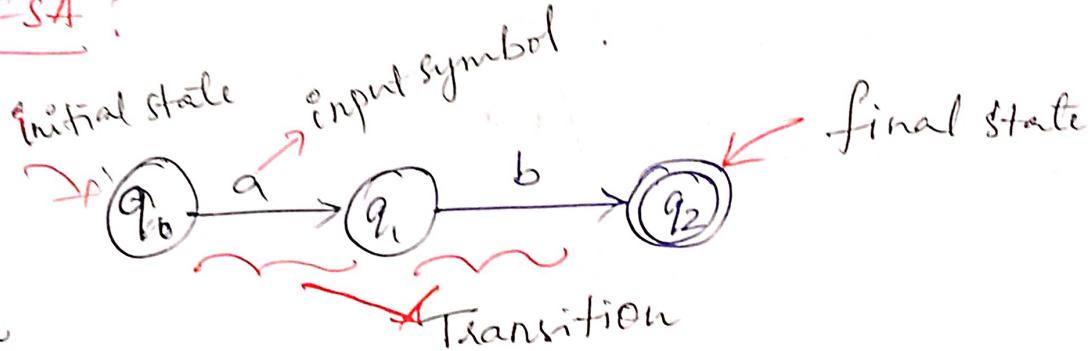
NFA: for each state there can be more than one symbol.

one transition on a given symbol, each leading to a different state.

Ex:



Ex for FSA:



$$4 \quad Q = \{q_0, q_1, q_2\}$$

$$\Sigma = \{a, b\}$$

q_0 = initial state

$$F = \{q_2\}$$

$$S(q_0, a) \Rightarrow q_1, \quad S(q_1, b) = q_2.$$

State-Transition table: list of transitions in table form

where \rightarrow row \rightarrow states

column \rightarrow input

Entries in table \rightarrow transition corresponds to input

$\rightarrow \emptyset \rightarrow$ Missing transition.

Ex:

State \downarrow	Inputs	
	a	b
State: q_0	q_1	\emptyset
q_1	\emptyset	q_2
final: q_2	\emptyset	\emptyset

Morphological parsing

(taking some input & producing some sort of structures for it) (5)

"Morphology" is the study of word structure & the formation of words from the smaller units (morphemes).

* Morphemes → Smallest meaning-bearing units in language.

Ex: eggs → has 2 morphemes:
egg s.

Morphological parser tells us that 'eggs' is the plural form of noun stem egg.

9 types classes of Morphemes:

↳ base word → Stem → Main morphem. (Central meaning) (Stems or lemmas)

↳ Affixes → Modify the stem meaning.

→ Prefix: before stem

→ Suffix: end of stem

→ Infix: inside a stem

→ Circumfix: Either end of the stem

3 ways of word formation:

V V
Run - Running

↳ Inflection: root word is combined to yield same class of stem

↳ Derivation: Combines a word stem with morpheme yield different class of word.

↳ Compounding:

Eg: Compounding noun from verb.

noun → Computation from verb → Computer

process of merging 2 or more words to form a new word.

Eg: Personal Computer, Desktop, overlook.

* Morphological parsing takes an input word in a text, as ~~of~~ it produces the parsed form consisting of 'canonical form' (or lemma) of the word & a 'set of tags'.

Ex: Part of speech, inflectional properties (gender, number, person, tense...).

Requirements for morphological parsers

- 1) Lexicon
 - 2) Morphotactics
 - 3) Orthographic rules

Lexicons lists stems + affixes together with basic information about them.

doing → $\underbrace{\text{do}}_{\text{stem}} + \underbrace{\text{V} + \text{PP}}_{\text{Morphological infn}} (\text{Present participle})$

books → book + N + PL (plural noun)

Morphotactics:

Ordering of Morphemes, (2)

describes the way morphemes are combined to form valid English words.
Ex: rest-ness-less ← invalid.

Orthographic rules

Graphographic rules: Spelling rules that specify the changes that occur when 2 morphemes are combined.

Ex: $y \rightarrow \text{ice}$, $\text{ice} \rightarrow y$, $E \rightarrow S$

- * two-level morphological model, ⁿ
- ↳ In this model, a word is represented as a correspondence between its lexical level form & its surface level form.
- * Surface level → represents actual spelling of the word.
- * lexical level → represents the concatenation of its constituent morphemes.
- * Morphological parsing is viewed as a mapping from the surface level into morpheme & feature sequences on the lexical level.
- * Ex: Surface form ~~is~~ represented for 'playing' is

P	L	a	y	i	n	g
---	---	---	---	---	---	---

 Surface level:

P	L	a	y	+V	PL
---	---	---	---	----	----

 Lexical level:

book	N	+V	PL
------	---	----	----

Similarly for "books" → Surface level form.
 "book+N+PL" → Lexical level form.
- * This Model is implemented with kind of FSA called Finite-State Transducer (FST).
- * Transducer maps set of symbols to another.
- * FST parses over the i/p string & through o/p strings in the form of symbol.

* FST is a 6-tuple $(\Sigma_1, \Sigma_2, Q, S, F)$.

Q → set of states

S → initial state

F → set of final states ($F \subseteq Q$)

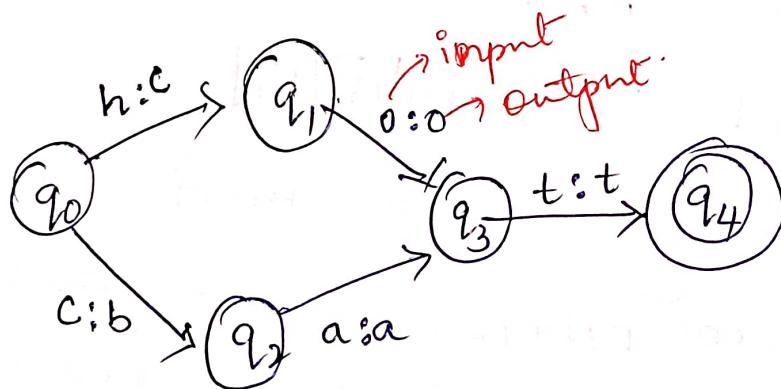
Σ_1 → input alphabet

Σ_2 → output alphabet

δ → Function mapping, $Q \times (\Sigma_1 \cup \{\epsilon\}) \times (\Sigma_2 \cup \{\epsilon\})$ to subset of power set of Q .

* FST transitions labelled with symbol from $\Sigma_1 \times \Sigma_2$ (input & output alphabets).

Ex: FST for that accepts 2 IP strings,
hot & cat maps them onto cat & bat respectively.

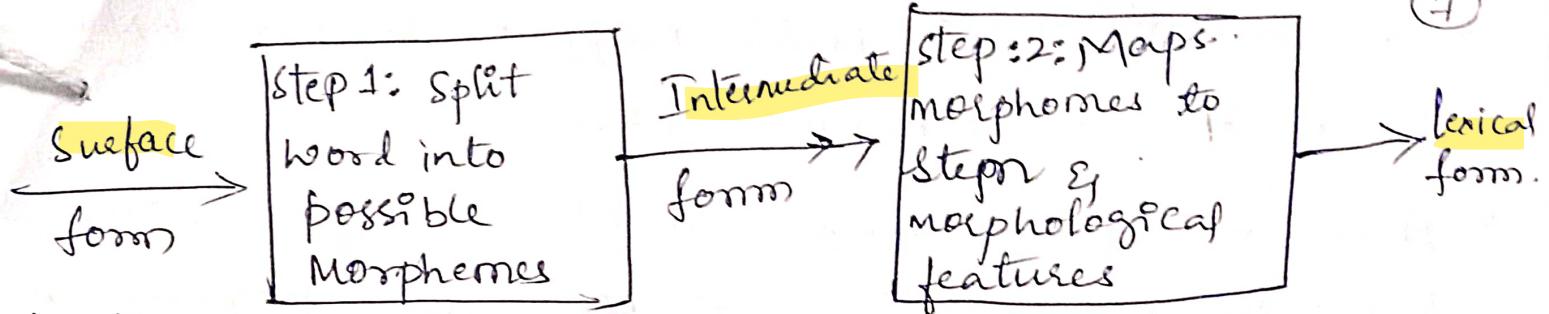


* 2-level morphology using FST.

Soft

(P.T.O)

(7)



birds
birds
boxes
goose
geese

birds
bird+s
box+s
goose
geese

bird+N+SG
bird+N+PL
box+N+PL
goose+N+SG
geese+N+PL

fig: Two-Step-Morphological parser

- * first splits word into its possible components.
- * Here we also consider spelling rules.

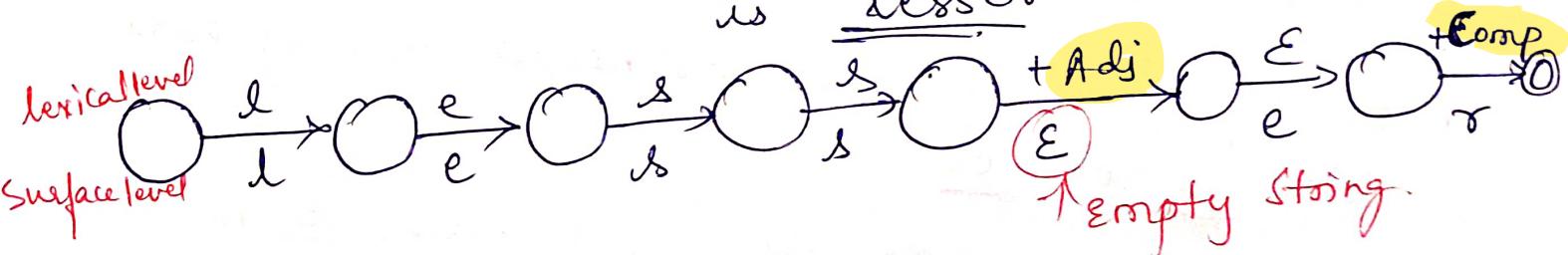
Ex ① birds ↗ bird + s

② boxes ↗ box + s
boxe + s

- * The o/p of the 1st step is concatenation of morphemes i.e stems + affixes.

There can be more than one representation for a word.

- * FST for adjective less → ~~less~~ Comparative form is "less"



- * In Second Step, birds mapped to bird+N+PL, if boxes to \Rightarrow box+N+PL.
- * Now we found 'boxe' is not a legal stem. This tells tells splitting boxes into box+es is incorrect & therefore discarded.
- * In other cases, spouses (or) parser
 ↓ ↓
 spouses (or) parser+s is correct.
- * Use Orthographic rules to handle spelling Variations
Ex: add e after -s, -z, -x, -ch, -sh before s.
 [dish \rightarrow dishes, box \rightarrow boxes].
- * Each of these steps can be implemented with the help of transducers.
- * Thus, we need to build 2 transducers,
 - \hookrightarrow 1. To maps the surface form to intermediate form
 - 2. To maps intermediate form to lexical form
- * FST ~~for~~ based Morphological parser for singular & plural noun in English:
- * plural form of regular nouns usually ends with 's' (or) 'es'
- * plural form of however, word ending 's' need not necessarily be the plural form of a word.
Ex: miss, ass \rightarrow singular form.

* deletion of 'e' ~~with~~ when introducing a morpheme boundary.

* deletion is usually required for words ending in ~~s~~, ~~es~~, ~~z~~ (e.g. suffixes & boxes). is shown below.

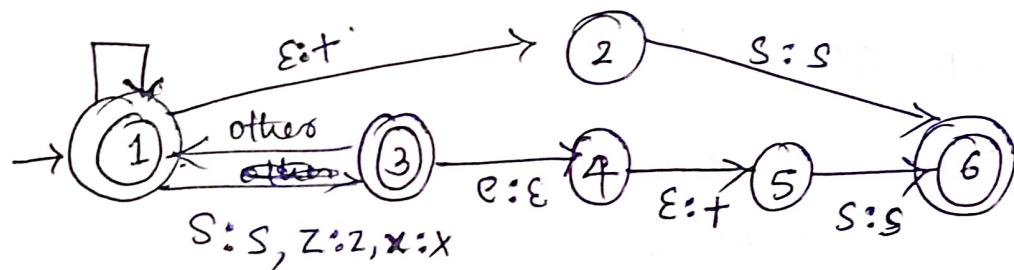
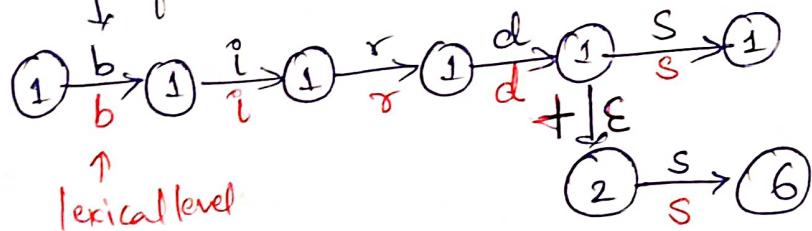


fig: Simplified FST, mapping English nouns to the intermediate form.

* Possible state sequences for the given surface form ~~birds or books~~ ~~use it's input is~~.

surface level.



OLP: birds

OLP: bird+s

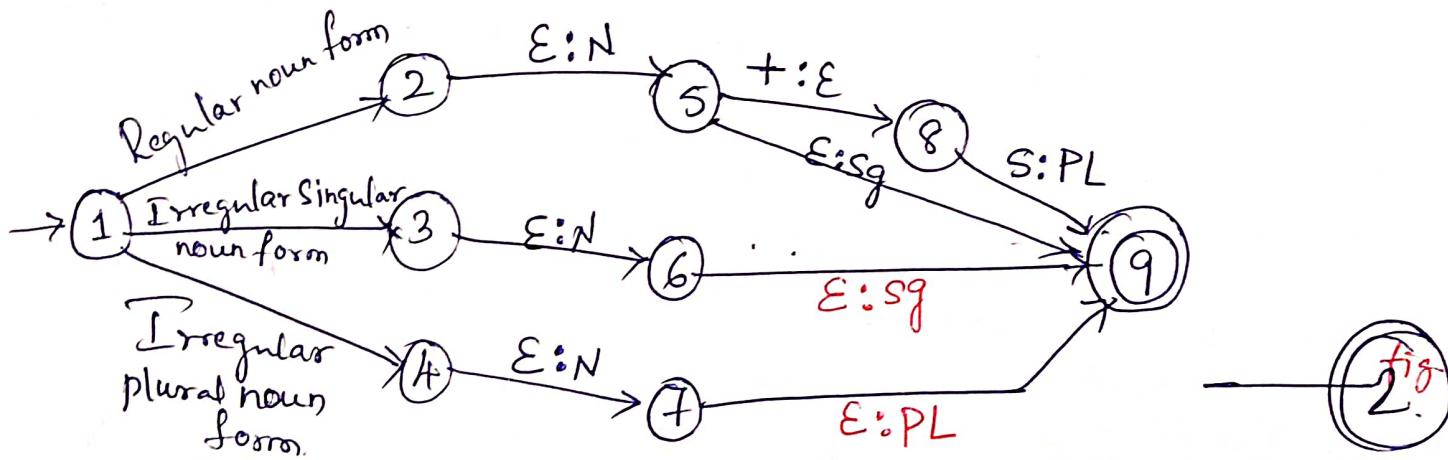


fig: Transducer for Step 2.

- * The next step is to develop a transducer that maps from intermediate level to lexical level. The IP to the transducer has one of the following forms
 - ↳ Regular noun stem (eg: bird, cat)
 - ↳ Regular noun stem + S (eg: bird + s)
 - ↳ Singular irregular noun stem (eg. goose)
 - ↳ Plural irregular noun stem (eg. geese)

* The general structure of this transducer shown in fig ②

- * Mapping of surface form geese (irregular noun) to its correct stem ·goose· as:
- | | | | | | |
|-----|-----|-----|-----|-----|--|
| g:g | e:o | e:o | s:s | e:e | $\begin{bmatrix} g & o & o & s & e \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ g & e & e & s & e \end{bmatrix}$ |
|-----|-----|-----|-----|-----|--|

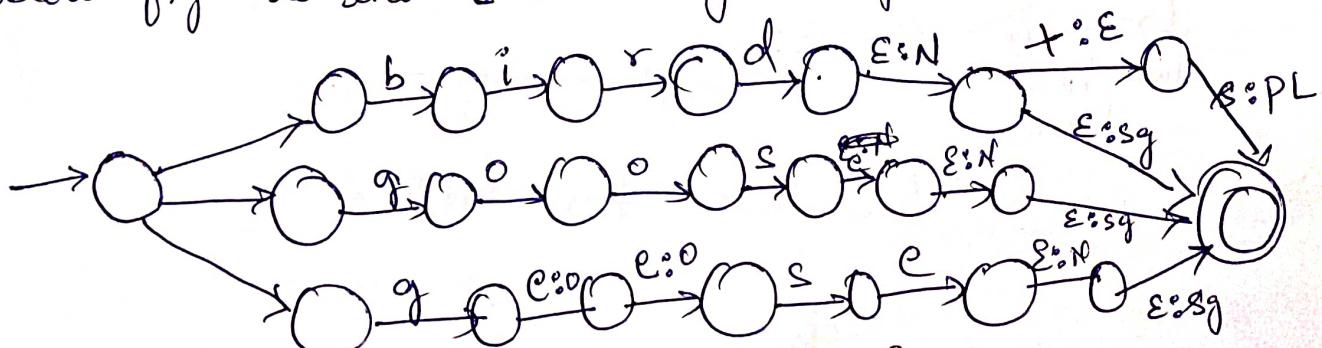
* bird → b:b i:i r:r d:d

- * Surface word "bird" maps to morphological ~~"bird+N+Pl"~~ "bird+N+Pl" as follows.

b:b i:i r:r d:d + ε:N + s:PL

[Each letter maps to itself but 'ε' maps to N + PL if 's' maps to PL]

- * Below fig is the resulting Composed transducer.



Spelling error detection & Correction (Refer PPT)