

Module-2

Word Level Analysis

Prepared By

Prof. Suja Jayachandran

Reference :

1.Speech and Language Processing by Daniel Jurafsky,James and Martin

Word: Fundamental building block of any language . Also called orthographic tokens(ie separated by white space). Exception Chinese, and Japanese languages are not separated by white space.

Morphology:deals with syntax of complex words and parts of word called morphemes as well as what semantic property they convey.

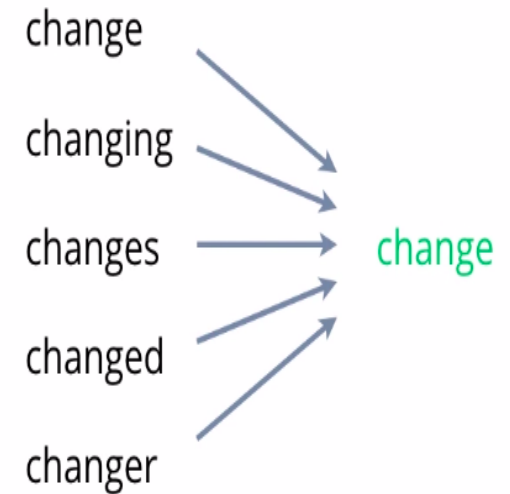
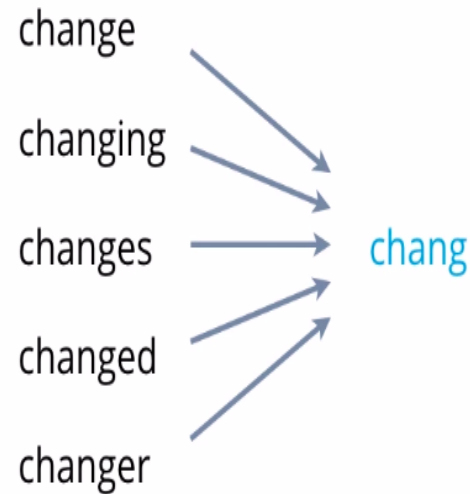
Inflectional vs Derivational Morphology		
More Information Online WWW.DIFFERENCEBETWEEN.COM		
	Inflectional Morphology	Derivational Morphology
DEFINITION	The study of the modification of words to fit into different grammatical contexts	The study of the formation of new words that differ either in syntactic category or in meaning from their base
MORPHEMES	Affixes that merely serve as grammatical markers and indicate some grammatical information about a word.	Affixes that are capable of either changing the meaning or the grammatical category of the word.
TYPE OF WORDS	Creates new forms of the same word	Creates new words

- **Tokenization:** Tokenization is the process of tokenizing or splitting a string, or text into a list of tokens. One can think of tokens as parts like a word is a token in a sentence, and a sentence is a token in a paragraph. It is important because the meaning of the text can be interpreted through analysis of the words present in the text.
 - Token: Non-unique words of a sentence
 - Type: Unique words of a sentence.
- **Stopword Removal:** Stop word removal is one of the most commonly used preprocessing steps across different NLP applications. The idea is simply to remove the words that occur commonly across all the documents in the corpus. Typically, articles and pronouns are generally classified as stop words. These words have no significance in some of the NLP tasks like information retrieval and classification.

- **Normalization** is the process of converting a token into its base form.
- It helps in :
 - Reducing the number of unique tokens present in the text.
 - Reducing variation of words in the text.
 - Removing redundant information.
- **Stemming**: It is an elementary rule-based process for removing inflectional forms from a token and output are the stem/root of the word.
- Eg: “Laughing”, “Laughed” , “Laughs” , “Laugh” will all become “laugh”
- It can produce words that are not part of the dictionary.
Eg: “His teams are not winning”.
After stemming: “Hi”, “team”, “are”, “not”, “winn”

- **Lemmatization**: It is a systematic step-by-step process for removing inflection forms of a word.
- It makes use of vocabulary, word structure, part of speech tags, and grammar relation.
- Output of lemmatization is the root word called a lemma.
- E.g. "Running", "Run", and "Ran" will be Run.

Stemming vs Lemmatization



Types of Stemming Algorithm

- **1. Porter Stemmer – PorterStemmer()**
 - Martin Porter invented the Porter Stemmer or Porter algorithm in 1980. Five steps of word reduction are used in the method, each with its own set of mapping rules. Porter Stemmer is the original stemmer and is renowned for its ease of use and rapidity. Frequently, the resultant stem is a shorter word with the same root meaning.
 - PorterStemmer() is a module in NLTK that implements the Porter Stemming technique.
- **2. Snowball Stemmer – SnowballStemmer()**
 - Martin Porter also created Snowball Stemmer. The method utilized in this instance is more precise and is referred to as “English Stemmer” or “Porter2 Stemmer.” It is somewhat faster and more logical than the original Porter Stemmer.
 - SnowballStemmer() is a module in NLTK that implements the Snowball stemming technique.
- **3. Lancaster Stemmer – LancasterStemmer()**
 - Lancaster Stemmer is straightforward, although it often produces results with excessive stemming. Over-stemming renders stems non-linguistic or meaningless.
 - LancasterStemmer() is a module in NLTK that implements the Lancaster stemming technique.
- **4. Regex Stemmer – RegexStemmer()**
 - Regex stemmer identifies morphological affixes using regular expressions. Substrings matching the regular expressions will be discarded.
 - RegexStemmer() is a module in NLTK that implements the Regex stemming technique.

- **Regular Expression**

- A regular expression (RE) is a language for specifying text search strings. RE helps us to match or find other strings or sets of strings, using a specialized syntax held in a pattern.
- Regular expression search requires a pattern that we want to search for, and a corpus of texts to search through.
- A regular expression search CORPUS function will search through the corpus returning all texts that contain the pattern
- Properties of Regular Expressions
 - American Mathematician Stephen Cole Kleene formalized the Regular Expression language.
 - RE is a formula in a special language, which can be used for specifying simple classes of strings, a sequence of symbols. In other words, we can say that RE is an algebraic notation for characterizing a set of strings.
 - Regular expression requires two things, one is the pattern that we wish to search and other is a corpus of text from which we need to search.

Searching for a particular word

Expression

share save flags

`/Mullen/g`

5 matches

Text




If you'd told Alex **Mullen** a few years ago that he was capable of memorising a whole pack of cards in 21.5 seconds, he would have said you were being ridiculous. His memory wasn't anything special – “below average” even. Fast-forward to today and **Mullen**, a medical student at the University of Mississippi, has just been crowned the World Memory Champion.

Mullen told me about a book he'd read called Moonwalking with Einstein. The book was written by Joshua Foer, a journalist who attended a US memory championship to write about what he thought would be “the Super Bowl of savants”. Instead, he found a group of people who had trained their memory using ancient techniques. Foer started practicing the techniques himself, and went on to win the competition the following year.

Mullen was spurred on to improve his own memory by Foer's story. “I definitely didn't have a great natural memory,” he said, “but in 2013 I started training using the techniques that Foer had talked about.” A year later, **Mullen** came second in the US memory championship. “It was really motivating, I kept practicing and eventually ended up at the 2015 world memory championship.”

- Searching a number

Expression

 share  save  flags

`/[0-9]+[.]?[0-9]/g`

3 matches

Text

If you'd told Alex Mullen a few years ago that he was capable of memorising a whole pack of cards in 21.5 seconds, he would have said you were being ridiculous. His memory wasn't anything special – “below average” even. Fast-forward to today and Mullen, a medical student at the University of Mississippi, has just been crowned the World Memory Champion.

Mullen told me about a book he'd read called Moonwalking with Einstein. The book was written by Joshua Foer, a journalist who attended a US memory championship to write about what he thought would be “the Super Bowl of savants”. Instead, he found a group of people who had trained their memory using ancient techniques. Foer started practicing the techniques himself, and went on to win the competition the following year.

Mullen was spurred on to improve his own memory by Foer's story. “I definitely didn't have a great natural memory,” he said, “but in 2013 I started training using the techniques that Foer had talked about.” A year later, Mullen came second in the US memory championship. “It was really motivating, I kept practicing and eventually ended up at the 2015 world memory championship.”

- Searching passages with number

Expression

[share](#) [save](#) [flags](#)

`/[Hh]e.*[0-9]+/g`

2 matches

Text

If you'd told Alex Mullen a few years ago that he was capable of memorising a whole pack of cards in 21.5 seconds, he would have said you were being ridiculous. His memory wasn't anything special – “below average” even. Fast-forward to today and Mullen, a medical student at the University of Mississippi, has just been crowned the World Memory Champion.

Mullen told me about a book he'd read called Moonwalking with Einstein. The book was written by Joshua Foer, a journalist who attended a US memory championship to write about what he thought would be “the Super Bowl of savants”. Instead, he found a group of people who had trained their memory using ancient techniques. Foer started practicing the techniques himself, and went on to win the competition the following year.

Mullen was spurred on to improve his own memory by Foer's story. “I definitely didn't have a great natural memory,” he said, “but in 2013 I started training using the techniques that Foer had talked about.” A year later, Mullen came second in the US memory championship. “It was really motivating, I kept practicing and eventually ended up at the 2015 world memory championship.”

- Named entity recognition

Expression

share

save

flags

/ [A-Za-z]+[0-9]+/g

28 matches

Text

The p53 gene family members p53, p73, and p63 display several isoforms derived from the presence of internal promoters and alternative splicing events. They are structural homologs but hold peculiar functional properties. p53, p73, and p63 are tumor suppressor genes that promote differentiation, senescence, and apoptosis. p53, unlike p73 and p63, is frequently mutated in cancer often displaying oncogenic "gain of function" activities correlated with the induction of proliferation, invasion, chemoresistance, and genomic instability in cancer cells. These oncogenic functions are promoted either by the aberrant transcriptional cooperation of mutant p53 (mutp53) with transcription cofactors (e.g., NF-Y, E2F1, Vitamin D Receptor, Ets-1, NF-kB and YAP) or by the interaction with the p53 family members, p73 and p63, determining their functional inactivation. The instauration of these aberrant transcriptional networks leads to increased cell growth, low activation of DNA damage response pathways (DNA damage response and DNA double-strand breaks response), enhanced invasion, and high chemoresistance to different conventional chemotherapeutic treatments. Several studies have clearly shown that different cancers harboring

RE: Build an FSA that modeled the subpart of English dealing with amounts of money. Such a formal language would model the subset of English consisting of phrases like ten cents, three dollars, one dollar thirty-five cents and so on

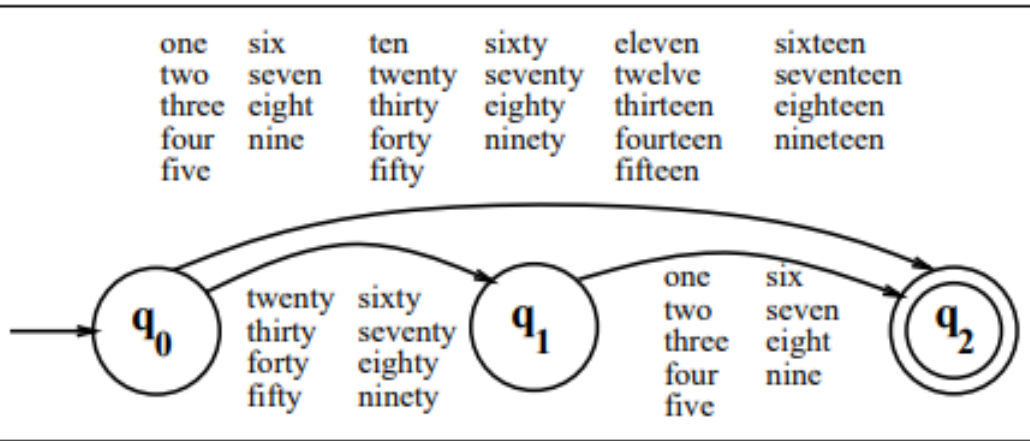


Figure 2.16 An FSA for the words for English numbers 1 – 99.

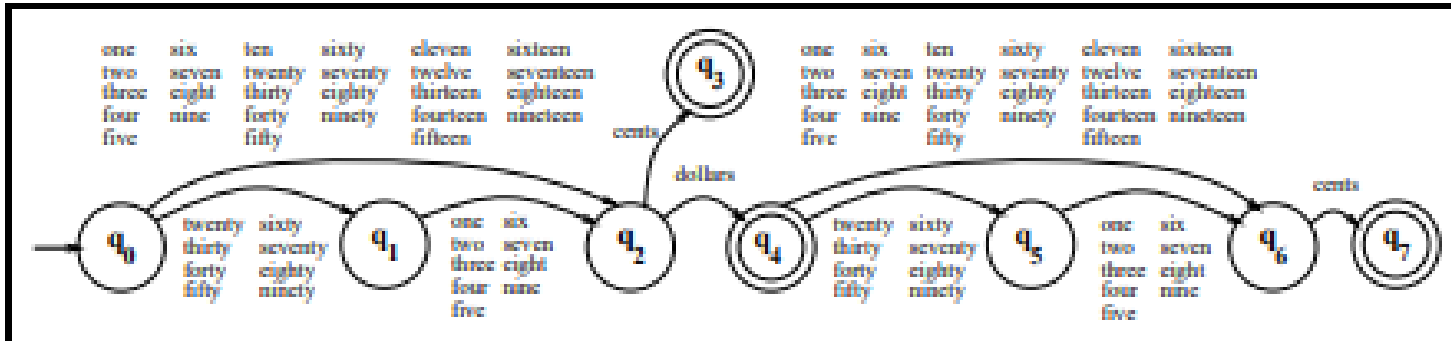


Figure 2.17 FSA for the simple dollars and cents.

Porter's Algorithm

- Porter (1980) algorithm we define a consonant as a letter other than A, E, I, O, and U, and other than Y preceded by a consonant. Any other letter is a vowel. (This is of course just an orthographic approximation.)
- Let c denote a consonant and v denote a vowel. C will stand for a string of one or more consonants, and V for a string of one or more vowels.
- Any written English word or word part can be represented by the following regular expression (where the parentheses $()$ are used to mark optional elements): $(C)(VC)^m(V)$
- For example the word troubles maps to the following sequence: troubles $CVCVC$ with no final V . We call the Kleene operator m the measure of any word or word part; the measure correlates very roughly with the number of syllables in the word or word part.

• Steps:

Step 1: Plural Nouns and Third Person Singular Verbs

The rules in this set do not have conditions:

SSES → SS	caresses → caress
IES → I	ponies → poni ties → ti
SS → SS	caress → caress
S → ε	cats → cat

Step 2a: Verbal Past Tense and Progressive Forms

(m > 1) EED → EE	feed → feed agreed → agree
(*v*) ED → ε	plastered → plaster bled → bled
(*v*) ING → ε	motoring → motor sing → sing

Step 2b: Cleanup

If the second or third of the rules in 2a is successful, we run the following rules (that remove double letters and put the E back on -ATE/-BLE)

AT → ATE	conflat(ed) → conflate
BL → BLE	troubl(ing) → trouble
IZ → IZE	siz(ed) → size
(*d & !(*L or *S or *Z)) → single letter	hopp(ing) → hop tann(ed) → tan fall(ing) → fall hiss(ing) → hiss fizz(ed) → fizz
(m=l & *o) → E	fail(ing) → fail fil(ing) → file

Step 3: Y → I

(*v*) Y → I	happy → happi sky → sky
-------------	----------------------------

Step 4: Derivational Morphology I: Multiple suffixes

(m > 0) ATIONAL → ATE	relational → relate
(m > 0) TIONAL → TION	conditional → condition
	rational → rational
(m > 0) ENCI → ENCE	valenci → valence
(m > 0) ANCI → ANCE	hesitanci → hesitance
(m > 0) IZER → IZE	digitizer → digitize
(m > 0) ABLI → ABLE	conformabli → conformable
(m > 0) ALLI → AL	radicalli → radical
(m > 0) ENTLI → ENT	differentli → different
(m > 0) ELI → E	vileli → vile
(m > 0) OUSLI → OUS	analogousli → analogous
(m > 0) IZATION → IZE	vietnamization → vietnamize
(m > 0) ATION → ATE	predication → predicate
(m > 0) ATOR → ATE	operator → operate
(m > 0) ALISM → AL	feudalism → feudal
(m > 0) IVENESS → IVE	decisiveness → decisive
(m > 0) FULNESS → FUL	hopefulness → hopeful
(m > 0) OUSNESS → OUS	callousness → callous
(m > 0) ALITI → AL	formaliti → formal
(m > 0) IVITI → IVE	sensitiviti → sensitive
(m > 0) BILITI → BLE	sensibiliti → sensible

Step 5: Derivational Morphology II: More multiple suffixes

(m > 0) ICATE → IC	triplicate → triplic
(m > 0) ATIVE → ε	formative → form
(m > 0) ALIZE → AL	formalize → formal
(m > 0) ICITI → IC	electricity → electric
(m > 0) FUL → ε	hopeful → hope
(m > 0) NESS → ε	goodness → good

Step 6: Derivational Morphology III: single suffixes

(m > 1) AL	→ ε	revival	→ reviv
(m > 1) ANCE	→ ε	allowance	→ allow
(m > 1) ENCE	→ ε	inference	→ infer
(m > 1) ER	→ ε	airliner	→ airlin
(m > 1) IC	→ ε	gyroscopic	→ gyroscop
(m > 1) ABLE	→ ε	defensible	→ defens
(m > 1) ANT	→ ε	irritant	→ irrit
(m > 1) EMENT	→ ε	replacement	→ replac
(m > 1) MENT	→ ε	adjustment	→ adjust
(m > 1) ENT	→ ε	dependent	→ depend
(m > 1) (*S or *T) & ION	→ ε	adoption	→ adopt
(m > 1) OU	→ ε	homologou	→ homolog
(m > 1) ISM	→ ε	communism	→ commun
(m > 1) ATE	→ ε	activate	→ activ
(m > 1) ITI	→ ε	angulariti	→ angular
(m > 1) OUS	→ ε	homologous	→ homolog
(m > 1) IVE	→ ε	effective	→ effect
(m > 1) IZE	→ ε	bowdlerize	→ bowdler

Step 7a: Cleanup

(m > 1) E → ε	probate → probat
	rate → rate
(m = 1 & ! *o) E → ε	cease → ceas

Step 7b: Cleanup

(m > 1 & *d *L) → [single letter]	controll → control
	roll → roll

Morphological Parsing

- **Morphology**-to understand how a word is formed
- **Morphological Parsing**-It is used to find morphemes form a word.
- Morphemes contain Stem(Root word) and Affix(Prefix(e.g. reform),infix(e.g. passerssby) and suffix(e.g. nationalist).
- Morphological Parser decides the order of words:
 - **Lexicons**-Stem,affix,part of speech(noun,adjective,verb)
 - **Morphotactics**:decides which morphemes should come based on rules.
 - **Orthographic rules**: lady+s=ladys(wrong)
lady+ies=ladies(true)

- **Types of Morphemes:**

1. **Free Morphemes:** Independent word having its own meaning(e.g. camera)

a) **Lexical Morphemes:**Adjective/Noun/Verb/Picture word(e.g. Yellow)

b) **Grammatical Morphemes:**Conjunction(e.g. and, or)

2. **Bound Morphemes:** No meaning of its own.(E.g. –ing (running))

a) **Inflectional Morphemes:**Words when combined with free morphemes it will not change the part of speech.IT will be always added as suffix.

Example: cat + s=cats

b) **Derivational Morphemes:**Words when combined with free morphemes it will change the part of speech.

Example: danger+ ous= dangerous

3. **Allomorphes:** antonym

Happy x unhappy

Rational x irrational

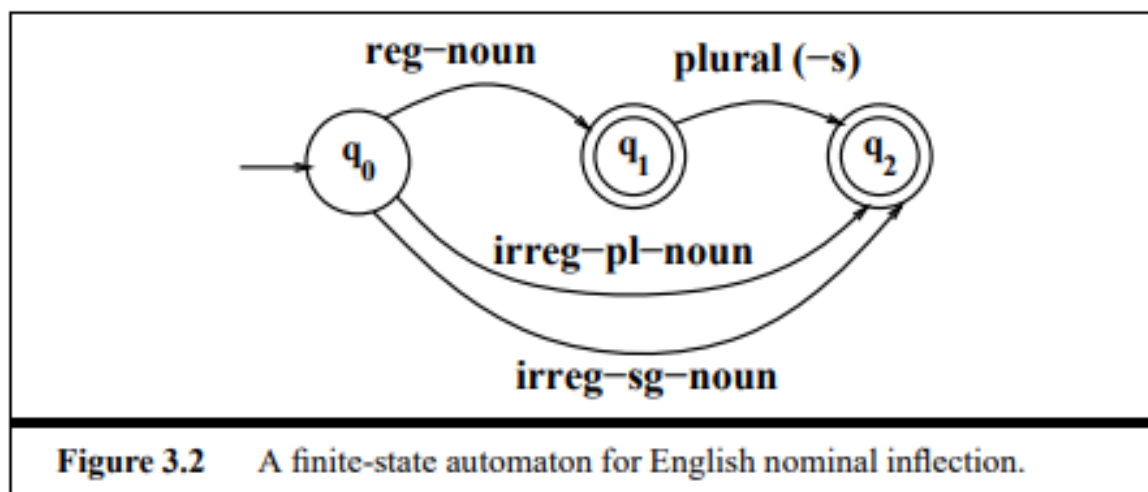
Possible x impossible

- Finite State Morphological Parsing

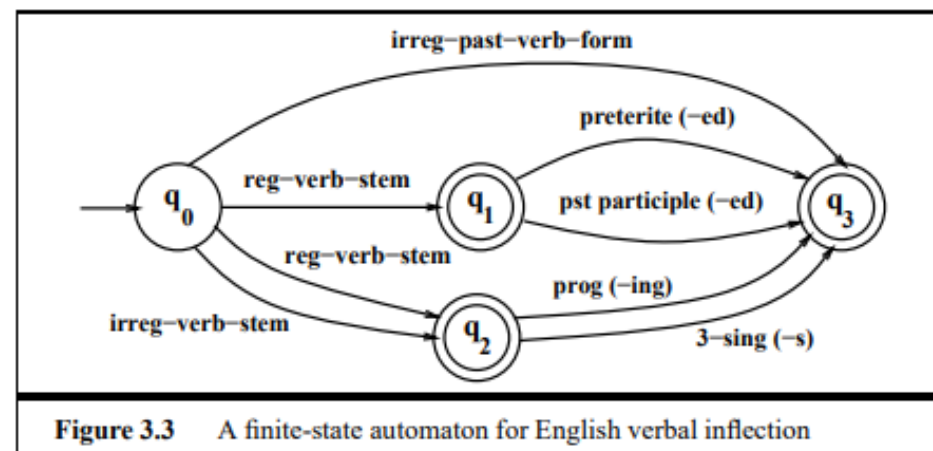
Let's now proceed to the problem of parsing English morphology. Consider a simple example: parsing just the productive nominal plural (-s) and the verbal progressive (-ing). Our goal will be to take input forms like those in the first column below and produce output forms like those in the second column.

Input	Morphological Parsed Output
cats	cat +N +PL
cat	cat +N +SG
cities	city +N +PL
geese	goose +N +PL
goose	(goose +N +SG) or (goose +V)
gooses	goose +V +3SG
merging	merge +V +PRES-PART
caught	(catch +V +PAST-PART) or (catch +V +PAST)

The second column contains the stem of each word as well as assorted morphological **features**. These features specify additional information about the stem. For example the feature +N means that the word is a noun; +SG means it is singular, +PL that it is plural. We will discuss features in Chapter 11; for now, consider +SG to be a primitive unit that means 'singular'. Note that some of the input forms (like *caught* or *goose*) will be ambiguous between different morphological parses.



reg-noun	irreg-pl-noun	irreg-sg-noun	plural
fox cat dog aardvark	geese sheep mice	goose sheep mouse	-s



Morphological Parsing with Finite-State Transducers(Pg 102 Jurafsky ebook)

- Two-level morphology, first proposed by Koskenniemi (1983).
- Two level morphology represents a word as a correspondence between a lexical level, which represents a simple concatenation of morphemes making up a word, and the surface level, which represents the actual spelling of the final word.
- Morphological parsing is implemented by building mapping rules that map letter sequences like cats on the surface level into morpheme and features sequences like cat +N +PL on the lexical level. Figure shows these two levels for the word cats. Note that the lexical level has the stem for a word, followed by the morphological information +N +PL which tells us that cats is a plural noun.

- The automaton that we use for performing the mapping between these two levels is the finite-state transducer or FST.
- A transducer maps between FST one set of symbols and another; a finite-state transducer does this via a finite automaton.
- Thus we usually visualize an FST as a two-tape automaton which recognizes or generates pairs of strings.
- The FST thus has a more general function than an FSA; where an FSA defines a formal language by defining a set of strings, an FST defines a relation between sets of strings. This relates to another view of an FST; as a machine that reads one string and generates another

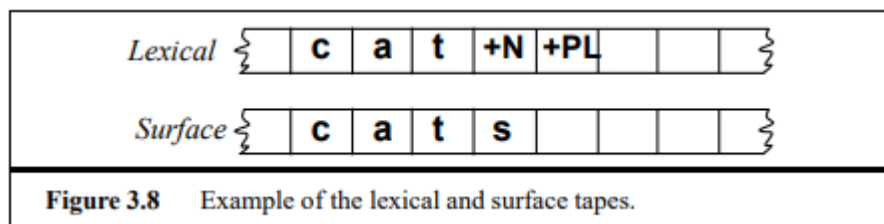
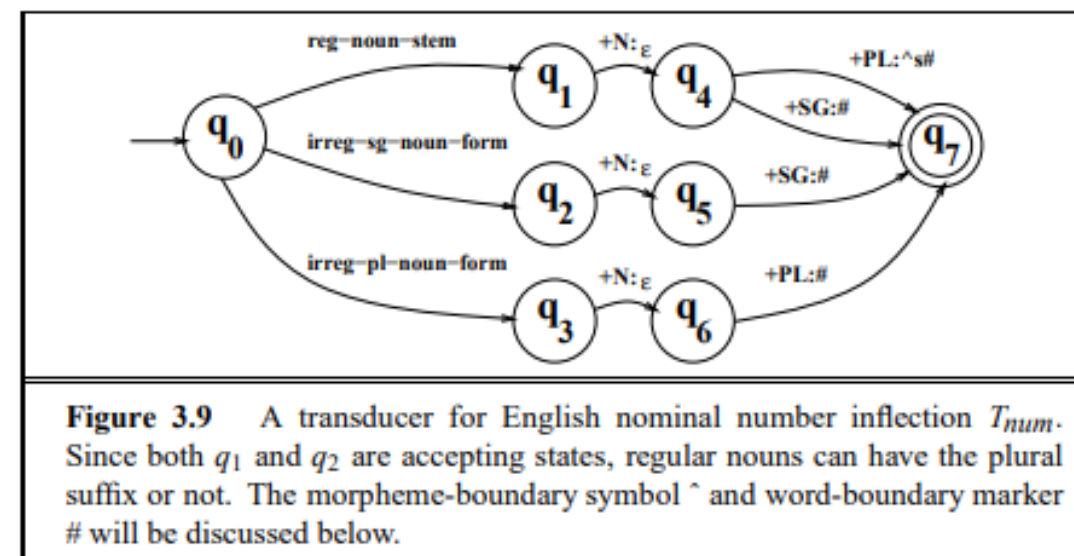


Figure 3.8 Example of the lexical and surface tapes.



That is, c maps to itself, as do a and t, while the morphological feature +N (recall that this means ‘noun’) maps to nothing (ϵ), and the feature +PL (meaning ‘plural’) maps to \hat{s} . The symbol $\hat{}$ indicates a morpheme boundary, while the symbol # indicates a word boundary.

This transducer will map plural nouns into the stem plus the morphological marker +PL, and singular nouns into the stem plus the morpheme +SG. Thus a surface *cats* will map to cat +N +PL as follows:

c:c a:a t:t +N: ϵ +PL: \hat{s} #

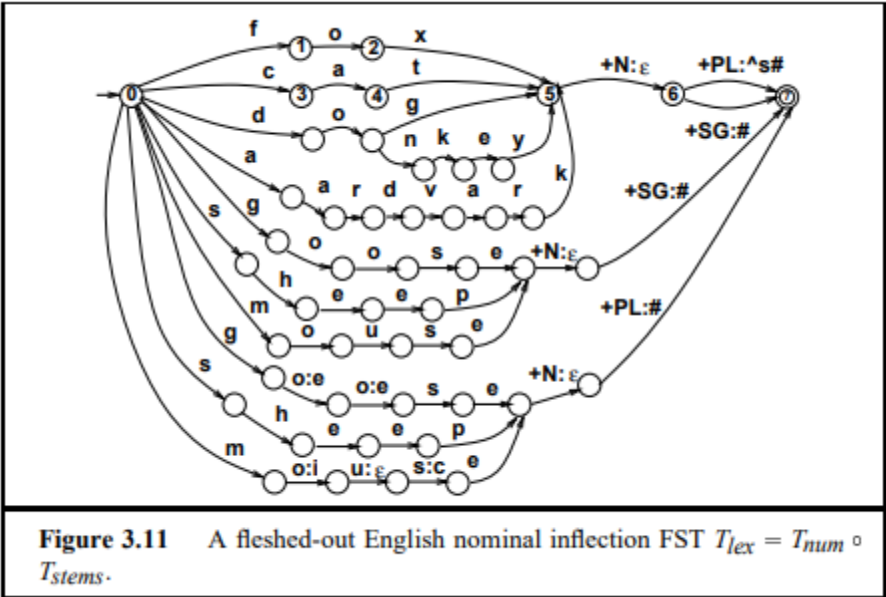


Figure 3.11 A fleshed-out English nominal inflection FST $T_{lex} = T_{num} \circ T_{stems}$.

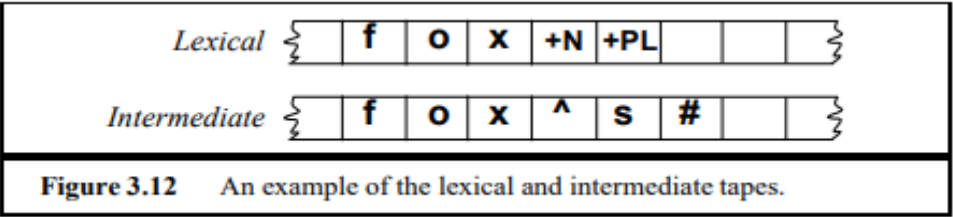


Figure 3.12 An example of the lexical and intermediate tapes.

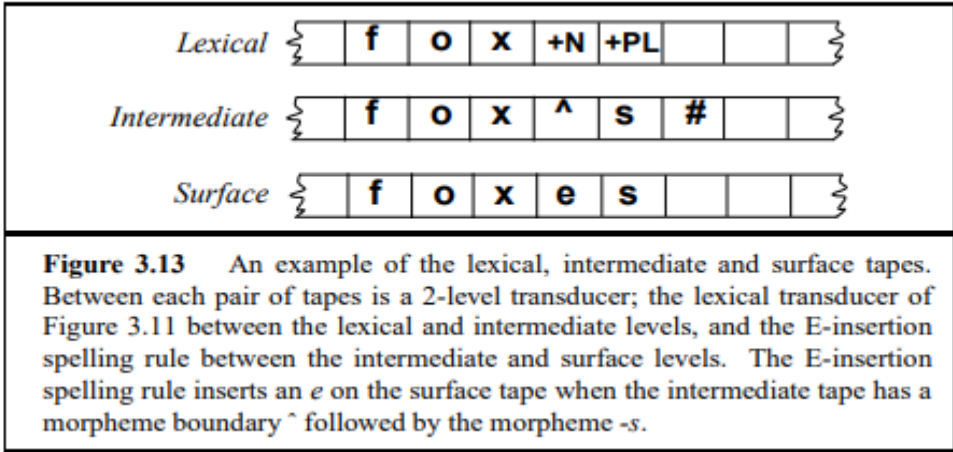


Figure 3.13 An example of the lexical, intermediate and surface tapes. Between each pair of tapes is a 2-level transducer; the lexical transducer of Figure 3.11 between the lexical and intermediate levels, and the E-insertion spelling rule between the intermediate and surface levels. The E-insertion spelling rule inserts an *e* on the surface tape when the intermediate tape has a morpheme boundary $\hat{}$ followed by the morpheme *-s*.

Orthographic Rules and Finite-State Transducers

Name	Description of Rule	Example
Consonant doubling	1-letter consonant doubled before <i>-ing/-ed</i>	beg/begging
E deletion	Silent e dropped before <i>-ing</i> and <i>-ed</i>	make/making
E insertion	e added after <i>-s,-z,-x,-ch, -sh</i> before <i>-s</i>	watch/watches
Y replacement	<i>-y</i> changes to <i>-ie</i> before <i>-s, -i</i> before <i>-ed</i>	try/tries
K insertion	verbs ending with <i>vowel + -c</i> add <i>-k</i>	panic/panicked

Language Model

- Language modeling is the way of determining the probability of any sequence of words. Language modeling is used in a wide variety of applications such as Speech Recognition, Spam filtering, etc.
- An N-gram language model predicts the probability of a given N-gram within any sequence of words in the language. A good N-gram model can predict the next word in the sentence i.e the value of $p(w|h)$.
- Two types of Language Modelings:
- **Statistical Language Modelings:** Statistical Language Modeling, or Language Modeling, is the development of probabilistic models that are able to predict the next word in the sequence given the words that precede. Examples such as N-gram language modeling.
- **Neural Language Modelings:** Neural network methods are achieving better results than classical methods both on standalone language models and when models are incorporated into larger models on challenging tasks like speech recognition and machine translation. A way of performing a neural language model is through word embeddings.

- **N-gram**
- N-gram can be defined as the contiguous sequence of n items from a given sample of text or speech. The items can be letters, words, or base pairs according to the application. The N-grams typically are collected from a text or speech corpus (A long text dataset).
- **N-gram Language Model:**
- An N-gram language model predicts the probability of a given N-gram within any sequence of words in the language. A good N-gram model can predict the next word in the sentence i.e the value of $p(w|h)$
- Example of N-gram such as unigram ("This", "article", "is", "on", "NLP") or bi-gram ('This article', 'article is', 'is on', 'on NLP')
- **Perplexity:** Perplexity is a measure of how good a probability distribution predicts a sample. It can be understood as a measure of uncertainty.
- $\text{Perplexity} = P(S)^{-1/n}$

- **Maximum Likelihood Estimate:** It is the method of estimating the parameter of an assumed probability distribution, given some observed data. It is the value that makes the observed data the “most probable”.

$$P(W_i | W_{i-1}) = \text{Count}(W_{i-1}, W_i) / \text{Count}(W_{i-1})$$

- **Laplace Smoothing:** Also called Add one smoothing. It is a smoothing technique that helps to tackle the problem of zero probability of a word in the text. When a particular bigram never occurred in our corpus data, then we get probability zero for that word. And when probability of any word is zero the overall effect is zero, which wastes the contribution of other words. To avoid this we can use Laplace smoothing.

- Practicenumeral based on N-gram model.

