

# 3 – Basic Text Processing



**CS 6320**

# Outline

---

- Regular Expressions
- Word Tokenization
- Word Morphology
- Sentence Segmentation

# Regular Expressions

---

- One of the unsung successes in standardization in computer science has been the regular expression (RE), *a language for specifying text search strings.*
  - This practical language is used in every computer language, word processor, and text processing tools like the Unix tools grep or Emacs.
- Formally, a regular expression is an algebraic notation for characterizing a set of strings.
- REs are particularly useful for searching in texts, when we have a pattern to search for and a corpus of texts to search through.
  - A regular expression search function will search through the corpus, returning all texts that match the pattern.
- The corpus can be a single document or a collection.
  - For example, the Unix command-line tool grep takes a regular expression and returns every line of the input document that matches the expression.

# Regular Expressions

---

- A formal language for specifying text strings
- How can we search for any of these?
  - woodchuck
  - woodchucks
  - Woodchuck
  - Woodchucks



# Regular Expressions: Disjunctions

---

- Letters inside square brackets []

Pattern	Matches
[wW] oodchuck	Woodchuck, woodchuck
[1234567890]	Any digit

- Ranges [A-Z]

Pattern	Matches	
[A-Z]	An upper case letter	Drenched Blossoms
[a-z]	A lower case letter	my beans were impatient
[0-9]	A single digit	Chapter <u>1</u> : Down the Rabbit Hole

# Regular Expressions: Negation in Disjunctions

---

- Negations [ ^Ss ]
  - Carat means negation only when first in []

Pattern	Matches	
[ ^A-Z ]	Not an upper case letter	O <u>y</u> fn pripetchik
[ ^Ss ]	Neither 'S' nor 's'	I have no exquisite reason"
[ ^e^ ]	Neither e nor ^	Look <u>h</u> ere
a^b	The pattern a carat b	Look up <u>a^b</u> now

# Regular Expressions: More Disjunctions

---

- Woodchucks is another name for groundhog!
- The pipe | for disjunction

Pattern	Matches
groundhog woodchuck	groundhog woodchuck
yours mine	yours mine
a b c	= [abc]
[gG] roundhog [Ww]oodchuck	Groundhog groundhog Woodchuck woodchuck



# More Examples

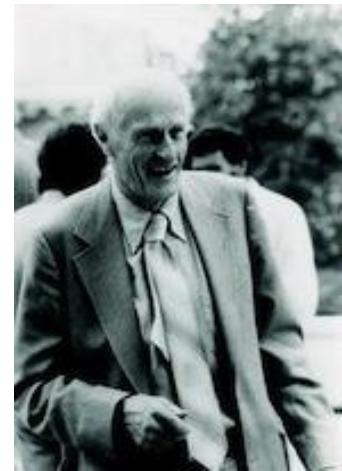
---

- Find all the instances of the word “the” in a text.
  - `/the/`  Misses capitalized examples
  - `/[tT]he/`  Incorrectly returns `other` or `Theology`
  - `/\b[tT]he\b/`

RE	Example Patterns Matched
<code>/woodchucks/</code>	“interesting links to <u>woodchucks</u> and lemurs”
<code>/a/</code>	“ <u>Mary</u> Ann stopped by <u>Mona</u> ’s”
<code>/Claire_says,/</code>	“ ‘Dagmar, my gift please,’ <u>Claire</u> says,”
<code>/DOROTHY/</code>	“SURRENDER <u>DOROTHY</u> ”
<code>/ ! /</code>	“You’ve left the burglar behind again!” said Nori

# Regular Expressions: ? \* +

Pattern	Matches	
colou?r	Optional previous char	<u>color</u> <u>colour</u>
oo*h!	0 or more of previous char	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
o+h!	1 or more of previous char	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
baa+		<u>baa</u> <u>baaa</u> <u>baaaa</u> <u>baaaaa</u>
beg.n		<u>begin</u> <u>begun</u> <u>begun</u> <u>beg3n</u>



Stephen C Kleene

## Operators

Kleene \*, Kleene +

# Anchors

---

- Kleene operators
  - Kleene \*: 0 or more occurrences of the previous char
    - `/[ab]*/` matches **ababab** or **bbbb**
  - Kleene +: 1 or more occurrences of the previous char
    - `/[0-9]+/` specifies a sequence of digits
- Special characters that anchor regular expressions to particular places in a string:
  - `^` specifies the **start** of a line
  - `$` matches the **end** of a line
  - `\b` matches a word boundary: `/\b35\b/` will match **There are 35 dresses in the room**, but it will not match **There are 356 destinations**.
  - `\B` matches a non-boundary

# What else???

- The use of caret ^ for **negation**.

RE	Match (single characters)	Example Patterns Matched
[ ^A-Z ]	not an upper case letter	“Oyfn pripetchik”
[ ^Ss ]	neither ‘S’ nor ‘s’	“I have no exquisite reason for’t”
[ ^\ . ]	not a period	“our resident Djinn”
[ e^ ]	either ‘e’ or ‘^’	“look up ^ now”
a^b	the pattern ‘a^b’	“look up a^ b now”

- The question mark ? Means **optionality** of the previous expression

RE	Match	Example Patterns Matched
woodchucks?	woodchuck or woodchucks	“woodchuck”
colou?r	color or colour	“colour”

- The use of the period . **To specify any character**

RE	Match	Example Patterns
/beg.n/	any character between <i>beg</i> and <i>n</i>	begin, beg’n, begun

# Advanced Operators-1

- Aliases for common set of characters

RE	Expansion	Match	Examples
\d	[ 0-9 ]	any digit	Party_of_5
\D	[ ^0-9 ]	any non-digit	Blue_moon
\w	[ a-zA-Z0-9_ ]	any alphanumeric/underscore	Daiyu
\W	[ ^\w ]	a non-alphanumeric	!!!
\s	[ \r\t\n\f ]	whitespace (space, tab)	
\S	[ ^\s ]	Non-whitespace	in_Concord

- Regular expression operators for counting

RE	Match
*	zero or more occurrences of the previous char or expression
+	one or more occurrences of the previous char or expression
?	exactly zero or one occurrence of the previous char or expression
{n}	$n$ occurrences of the previous char or expression
{n,m}	from $n$ to $m$ occurrences of the previous char or expression
{n,}	at least $n$ occurrences of the previous char or expression

# Advanced Operators-2

- Characters that need to be backslashed

RE	Match	Example Patterns Matched
\*	an asterisk “*”	“K\_A*P*L*A*N”
\.	a period “.”	“Dr. Livingston, I presume”
\?	a question mark	“Why don’t they come and lend a hand?”
\n	a newline	
\t	a tab	

- More examples: **ELIZA** simulating conversations:

User<sub>1</sub>: Men are all alike.

ELIZA<sub>1</sub>: IN WHAT WAY

User<sub>2</sub>: They’re always bugging us about something or other.

ELIZA<sub>2</sub>: CAN YOU THINK OF A SPECIFIC EXAMPLE

User<sub>3</sub>: Well, my boyfriend made me come here.

ELIZA<sub>3</sub>: YOUR BOYFRIEND MADE YOU COME HERE

User<sub>4</sub>: He says I’m depressed much of the time.

ELIZA<sub>4</sub>: I AM SORRY TO HEAR YOU ARE DEPRESSED

# How is Eliza working?

---

- **ELIZA** simulating conversations:

User<sub>1</sub>: Men are all alike.

ELIZA<sub>1</sub>: IN WHAT WAY

User<sub>2</sub>: They're always bugging us about something or other.

ELIZA<sub>2</sub>: CAN YOU THINK OF A SPECIFIC EXAMPLE

User<sub>3</sub>: Well, my boyfriend made me come here.

ELIZA<sub>3</sub>: YOUR BOYFRIEND MADE YOU COME HERE

User<sub>4</sub>: He says I'm depressed much of the time.

ELIZA<sub>4</sub>: I AM SORRY TO HEAR YOU ARE DEPRESSED

- It employs a cascade of regular expression substitutions

```
s/.* YOU ARE (depressed|sad) .*/I AM SORRY TO HEAR YOU ARE \1/  
s/.* YOU ARE (depressed|sad) .*/WHY DO YOU THINK YOU ARE \1/  
s/.* all .*/IN WHAT WAY/  
s/.* always .*/CAN YOU THINK OF A SPECIFIC EXAMPLE/
```

# Regular Expressions: Anchors ^ \$

---

- In a nutshell:

Pattern	Matches
[A-Z]	Palo Alto <u>P</u> alo Alto
^ [ ^A-Za-z ]	<u>1</u> "Hello" <u>H</u> ello
\ . \$	The end <u>.</u>
. \$	The end <u>?</u> The end <u>!</u>

\$ matches the end of a line

# Errors

---

- Two kinds of errors
  - Matching strings that we should not have matched  
(**there, then, other**)
    - False positives (Type I)
  - Not matching things that we should have matched  
(**The**)
    - False negatives (Type II)

# Errors cont.

---

- In NLP we are always dealing with these kinds of errors.
- Reducing the error rate for an application often involves two antagonistic efforts:
  - Increasing accuracy or precision (minimizing false positives)
  - Increasing coverage or recall (minimizing false negatives).

# Finite State Automata 1/4

- Regular Expressions can describe a FSA machine.
- FSAs are useful for NLP.
- FSA to recognize the “sheep language”  
/ baa+! /

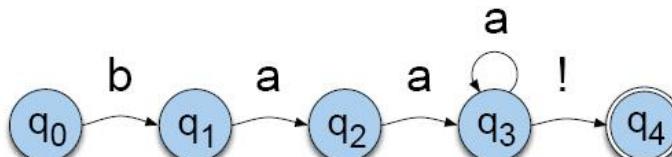


Figure 2.10 A finite state automaton for talking sheep.

- It has five states
- $q_0$  is the start state
- $q_4$  is the final state
- It has four transitions.

# Finite State Automata 2/4

- FSAs can be encoded as tables.

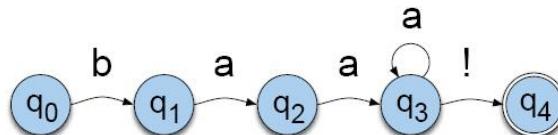


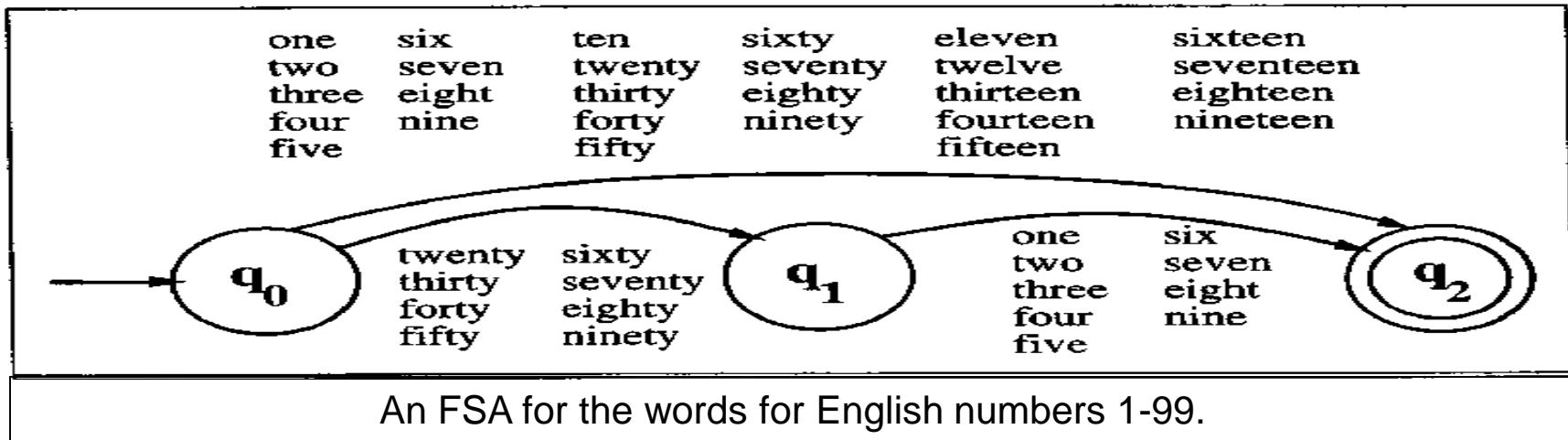
Figure 2.10 A finite-state automaton for talking sheep.

Input			
State	b	a	!
0	1	Ø	Ø
1	Ø	2	Ø
2	Ø	3	Ø
3	Ø	3	4
4:	Ø	Ø	Ø

state-transition table

# Finite State Automata 3/4

- An example: FSA to recognize amounts of money.
- Ten cents, three dollars, one dollar thirty five cents



# Finite State Automata 4/4

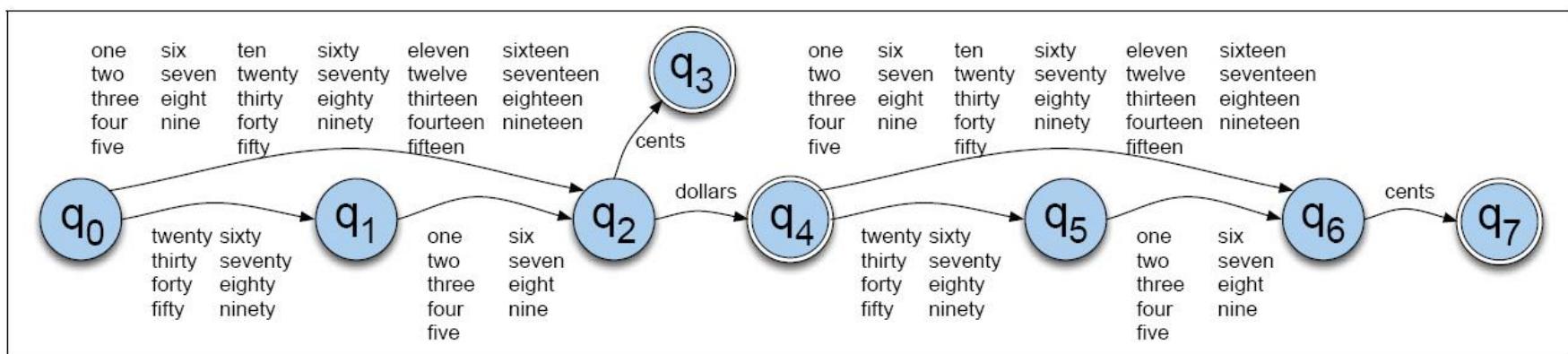


Figure 2.16 FSA for the simple dollars and cents

- FSAs can be formally specified as a 5 tuple:
  - The set of states:  $Q$
  - A finite alphabet:  $\Sigma$
  - A start state
  - A set of final states
  - A transition function that maps  $Q \times \Sigma$  to  $Q$ .

# More Formally again

---

- The definition of a finite automaton

$Q = q_0q_1q_2 \dots q_{N-1}$

a finite set of  $N$  **states**

$\Sigma$

a finite **input alphabet** of symbols

$q_0$

the **start state**

$F$

the set of **final states**,  $F \subseteq Q$

$\delta(q, i)$

the **transition function** or transition matrix between states. Given a state  $q \in Q$  and an input symbol  $i \in \Sigma$ ,  $\delta(q, i)$  returns a new state  $q' \in Q$ .  $\delta$  is thus a relation from  $Q \times \Sigma$  to  $Q$ ;

# Recognition

- Recognition is the process of determining whether or not a given input is accepted by a machine.
- In terms of REs, it is the process of determining whether or not a given input matches a particular RE.
- Recognition is viewed as processing an input written on a tape consisting of cells containing elements from the alphabet.

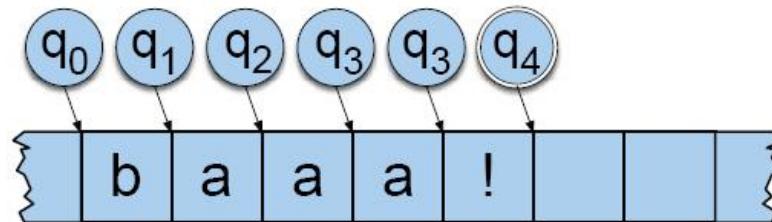


Figure 2.13 Tracing the execution of FSA #1 on some sheeptalk.

# Deterministic vs. Non-Deterministic FSAs

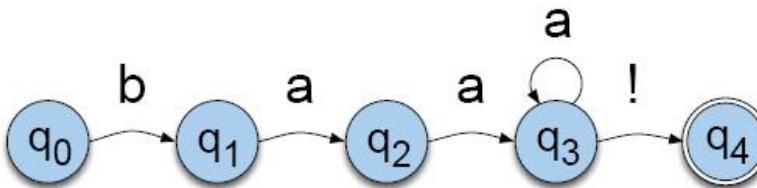


Figure 2.10 A finite-state automaton for talking sheep.

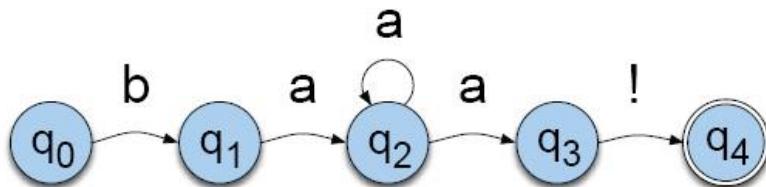


Figure 2.17 A non-deterministic finite-state automaton for talking sheep (NFSA #1).  
Compare with the deterministic automaton in Fig. 2.10

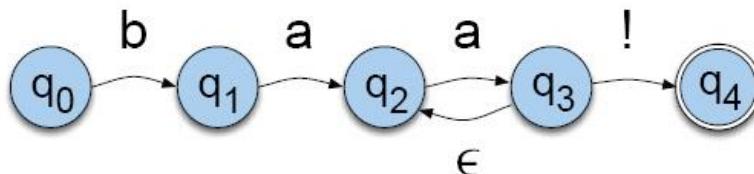


Figure 2.18 Another NFSA for the sheep language (NFSA #2). It differs from NFSA #1 in Fig. 2.17 in having an  $\epsilon$ -transition.

# ND Recognition as Search 1/3

---

- Idea:
  - A search state is a pairing of a single machine state with a position on the input tape.
  - By keeping track of not yet explored search states, a recognizer can systematically explore all possible paths through a machine given some input.

# ND Recognition as Search 2/3

- Depthfirst search or LIFO

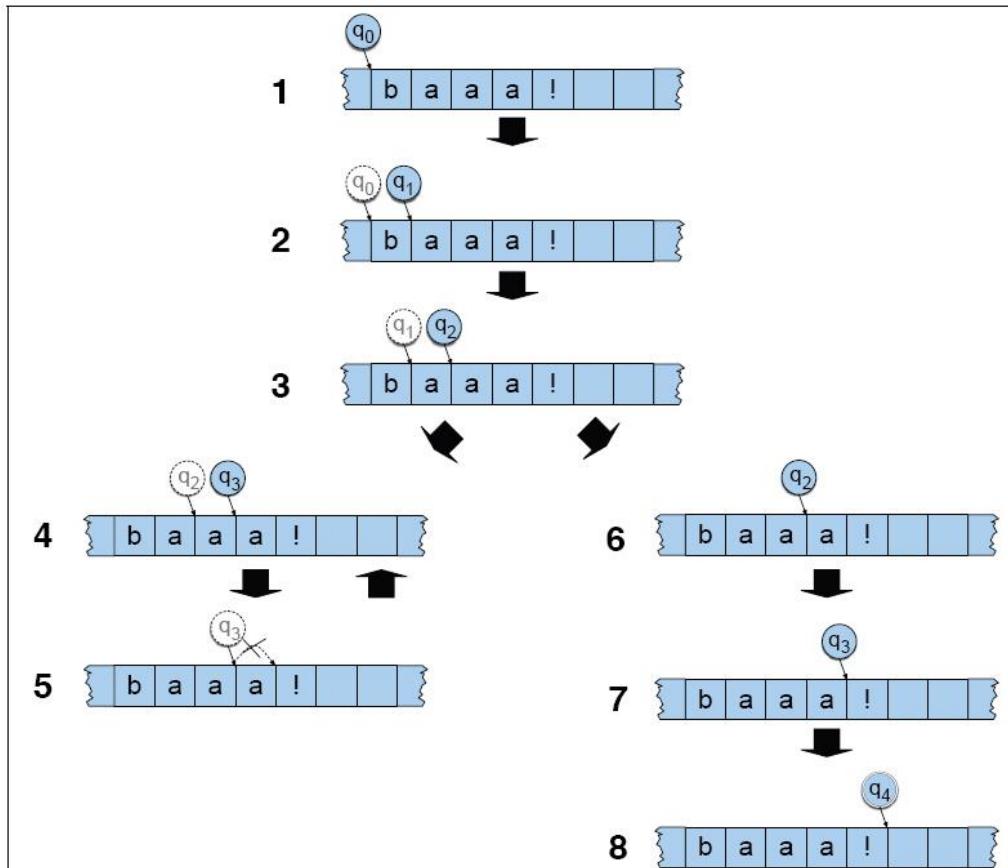


Figure 2.20 Tracing the execution of NFSA #1 (Fig. 2.17) on some sheeptalk.

# ND Recognition as Search 3/3

- Breadthfirst search or FIFO

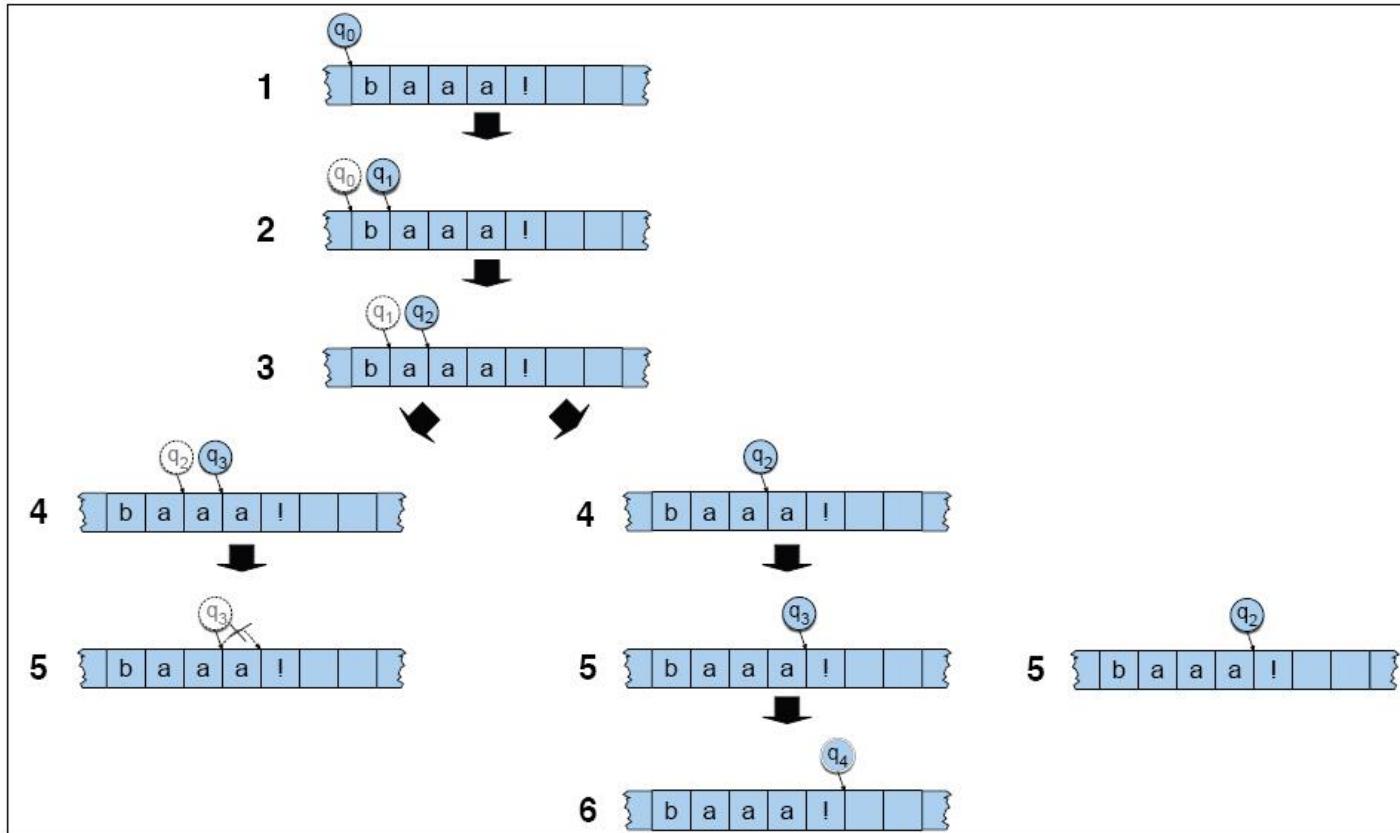


Figure 2.21 A breadth-first trace of FSA #1 on some sheeptalk

# Generative Grammars and Formal Languages

---

- A Formal Language is a set of strings composed of symbols from a finite set of symbols.
- FSAs (and REs) define formal languages without having to explicitly enumerate the set.
- The term generative refers to the idea that FSAs can be viewed as generators of formal languages as well as acceptors.
- To generate you traverse the machine and transition symbols on the tape rather than reading them.

# Summary

---

- Regular expressions play a surprisingly large role
  - Sophisticated sequences of regular expressions are often the first model for any text processing task
- For many hard tasks, we use machine learning classifiers
  - But regular expressions are used as features in the classifiers
  - Can be very useful in capturing generalizations

# Words

---

*they lay back on the San Francisco grass and looked at the stars and their*

- **Token:** an instance of that type in running text.
- How many?
  - 15 tokens (or 14)

# How many words?

---

**N** = number of tokens

**V** = vocabulary

Church and Gale (1990):  $|V| > O(N^{1/2})$

	<b>Tokens = N</b>	<b> V </b>
Switchboard phone conversations	2.4 million	20 thousand
Shakespeare	884,000	31 thousand
Google N-grams	1 trillion	13 million

# Issues in Tokenization

---

- Finland's capital → Finland Finlands Finland's ?
- what're, I'm, isn't → What are, I am, is not
- Hewlett-Packard → Hewlett Packard ?
- state-of-the-art → state of the art ?
- Lowercase → lowercase lowercase lower case ?
- San Francisco → one token or two?
- m.p.h., PhD. → ??

*In Natural Language Processing we care about punctuation – and do not discard it!!!!*

*Finland's capital* → *Finland ' s capital*

# Tokenization: language issues

---

- French
  - **L'ensemble** → one token or two?
    - **L** ? **L'** ? **Le** ?
    - Want **L'ensemble** to match with **un ensemble**
- German noun compounds are not segmented
  - **Lebensversicherungsgesellschaftsangestellter**
  - 'life insurance company employee'
  - German information retrieval needs **compound splitter**

**IMPORTANT: where can I find a good English Tokenizer???**

**Answer: Stanford NLP software:**

<https://nlp.stanford.edu/software/lex-parser.html>

# Word Morphology

---

- **Morphology** is the study of the way words are built from smaller units called morphemes
  - A morpheme is a minimal meaning-bearing unit in language.
  - Example: the word **DOG** has a single morpheme; the word **CATS** has two: (1) **CAT** and (2) **-S**
- There are two classes of morphemes: **stems** and **affixes**.
- Stems -- are the main morphemes of words.
- Affixes – add additional meaning to the stems to modify their meanings and grammatical functions

# Morphology and Finite-State Transducers

---

- There are four forms of affixes:
  1. **Prefixes** – precede the stem
  2. **Suffixes** – follow the stem
  3. **Infixes** – inserted inside the stem
  4. **Circumfixes** – both precede and follow the stem.
- Examples:
  - Prefixes - “un”, “a”
  - Suffixes - plurals, “ing”
  - Infixes - not common in English
  - Circumfixes : unbelievably  
un + believe + able + ly

# Kinds of Morphology

---

- The usage of prefixes and suffixes concatenated to the stem creates a **concatenative morphology**.
  - E.g. “door” and “doors”; “run” and “running”
- When morphemes are combined in more complex ways, we have a **non-concatenative morphology**.
  - E.g. “foot” and “feet”; “freeze” and “frozen”
- **Inflectional morphology** - is the combination of a word stem with a grammatical morpheme usually resulting in a word of the same class
  - Special case: **Templatic morphology**, also known as root-and-pattern morphology:
    - Used in Semitic languages, e.g. Arabic, Hebrew

# Example of Templatistic Morphology

---

- In Hebrew, a verb is constructed using two components:
  - A root (consisting usually of 3 consonants CCC) – carrying the main meaning
  - A template which gives the ordering of consonants and vowels and specifies more semantic information about the resulting verb, e.g., the voice (active, passive)
- Example: the tri-consonant root **lmd** (meaning: learn, study)
  - can be combined with a template **CaCaC** for active voice to produce the word **lamad** = “he studied”
  - can be combined with the template **CiCeC** for intensive to produce the word **limed** = “he taught”
  - can be combined with a template **CuCaC** for active voice to produce the word **lumad** = “he was taught”

# Producing words from morphemes

---

- By **inflection** – **Inflectional morphology** - is the combination of a word stem with a grammatical morpheme usually resulting in a word of the same class
  - E.g. “leaf” and “leaves”; “write” and “writes”
- By **derivation** – **Derivational Morphology** - is the combination of a word stem with a grammatical morpheme usually resulting in a word of a different class.
  - E.g. “leaf” and “leaflet”; “write” and “writer”
- By **compounding** – by combining multiple words together,
  - E.g. “doghouse”
- By **cliticization** – by combining a word stem with a clitic. A **clitic** is a morpheme that acts syntactically like a word but it is reduced in forms and it is attached to another word.
  - E.g. *I've*

# Clitization

---

- A **clitic** is a unit whose status lies between that of an affix and a word
  - The **phonological behavior** of clitics is like affixes: they tend to be short and unaccented
  - The **syntactic behavior** is more like words: they often act as pronouns, articles, conjunctions or verbs.

<b>Full Form</b>	<b>Critic</b>	<b>Full Form</b>	<b>Critic</b>
am	'm	have	've
are	're	has	's
is	's	had	'd
will	'll	would	'd

# Inflectional Morphology 1/2

- **Nouns:** have an affix for plural and an affix for possessive
  - Plural – the suffix **-s** and the alternative spelling **-es** for regular nouns

Singular	Plural
<i>dog</i>	<i>dogs</i>
<i>farm</i>	<i>farms</i>
<i>school</i>	<i>schools</i>
<i>car</i>	<i>cars</i>

Singular	Plural
<i>ibis</i>	<i>ibises</i>
<i>waltz</i>	<i>waltzes</i>
<i>box</i>	<i>boxes</i>
<i>butterfly</i>	<i>butterflies</i>

Irregular nouns:

Singular	Plural
<i>mouse</i>	<i>mice</i>
<i>ox</i>	<i>oxen</i>

Possessive - for words not ending in "s" the affix is **"'s"** (*children/ children's*) and for words ending in "s" the affix is **"'"** (*llama/llama's; llamas/llamas'*) (*Euripides/Euripides' comedies*)

# Inflectional Morphology 2/2

- Verbal inflection is more complex than nominal inflection
- There are three classes of verbs in English:
  - **Main verbs** (eat, sleep, walk)
  - **Modal verbs** (can, will, should)
  - **Primary verbs** (be, have, do)

Regular/ Irregular Verbs

Morphological Class	Regularly Inflected Verbs			
stem	walk	merge	try	map
-s form	walks	merges	tries	maps
-ing participle	walking	merging	trying	mapping
Past form or -ed participle	walked	merged	tried	mapped

Morphological Class	Irregularly Inflected Verbs		
stem	eat	catch	cut
-s form	eats	catches	cuts
-ing participle	eating	catching	cutting
preterite	ate	caught	cut
past participle	eaten	caught	cut

# Spanish Verb System

	<b>Present Indicative</b>	<b>Imperfect Indicative</b>	<b>Future</b>	<b>Preterite</b>	<b>Present Subjunctive</b>	<b>Conditional</b>	<b>Imperfect Subjunctive</b>	<b>Future Subjunctive</b>
1SG	amo	amaba	amaré	amé	ame	amaría	amara	amare
2SG	amas	amabas	amarás	amaste	ames	amarías	amaras	amares
3SG	ama	amaba	amará	amó	ame	amaría	amara	amáremo
1PL	amamos	amábamos	amaremos	amamos	amemos	amaríamos	amáramos	amáremos
2PL	amáis	amabais	amaréis	amasteis	améis	amaríais	amarais	amareis
3PL	aman	amaban	amarán	amaron	amen	amarían	amaran	amaren

“To love” in Spanish.

Some of the inflected forms of the verb “amar” in European Spanish  
50 distinct verb forms for each regular verb.

# Derivational Morphology

---

- Changes of word class
  - Nominalization: the formation of new nouns from verbs or adjectives.

Suffix	Base Verb/Adjective	Derived Noun
-ation	computerize (V)	computerization
-ee	appoint (V)	appointee
-er	kill (V)	killer
-ness	fuzzy (A)	fuzziness

- Adjectives can also be derived from verbs.

Suffix	Base Noun/Verb	Derived Adjective
-al	computation (N)	computational
-able	embrace (V)	embraceable
-less	clue (N)	clueless

# Morphological Parsing

- In general parsing means taking an input and producing a structure for it.
- Morphological parsing takes as an input words and produces a structure that reveals its morphological features.

	<b>English</b>		<b>Spanish</b>		
<b>Input</b>	<b>Morphological Parse</b>		<b>Input</b>	<b>Morphological Parse</b>	<b>Gloss</b>
cats	cat +N +PL		pavos	pavo +N +Masc +Pl	'ducks'
cat	cat +N +SG		pavo	pavo +N +Masc +Sg	'duck'
cities	city +N +Pl		bebo	beber +V +PInd +1P +Sg	'I drink'
geese	goose +N +Pl		canto	cantar +V +PInd +1P +Sg	'I sing'
goose	goose +N +Sg		canto	canto +N +Masc +Sg	'song'
goose	goose +V		puse	poner +V +Perf +1P +Sg	'I was able'
gooses	goose +V +3P +Sg		vino	venir +V +Perf +3P +Sg	'he/she came'
merging	merge +V +PresPart		vino	vino +N +Masc +Sg	'wine'
caught	catch +V +PastPart		lugar	lugar +N +Masc +Sg	'place'
caught	catch +V +Past				

# Morphology and FSAs

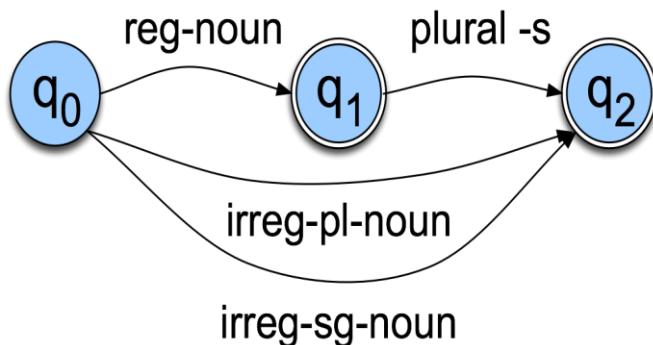
---

- We'd like to use the machinery provided by FSAs to capture these facts about morphology
  - **Accept strings** that are in the language
  - **Reject strings** that are not
  - And do so in a way that doesn't require us to in effect list all the words in the language

# Building a finite-state lexicon: simple rules

---

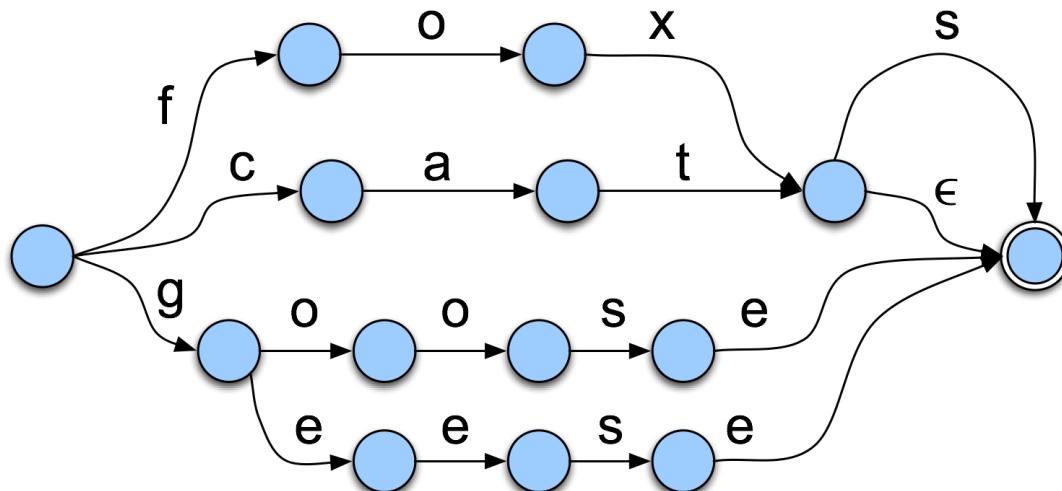
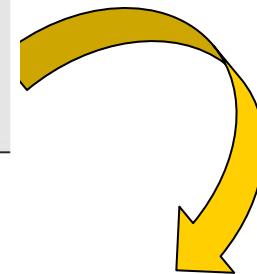
- A lexicon is a repository of words.
- Possibilities:
  - List all words in the language
  - Computational lexicons: list all stems and affixes of a language + representation of the morphotactics that tells us how they fit together.
    - A example of morphotactics: the finite-state automaton (FSA) for English nominal inflection



Nouns

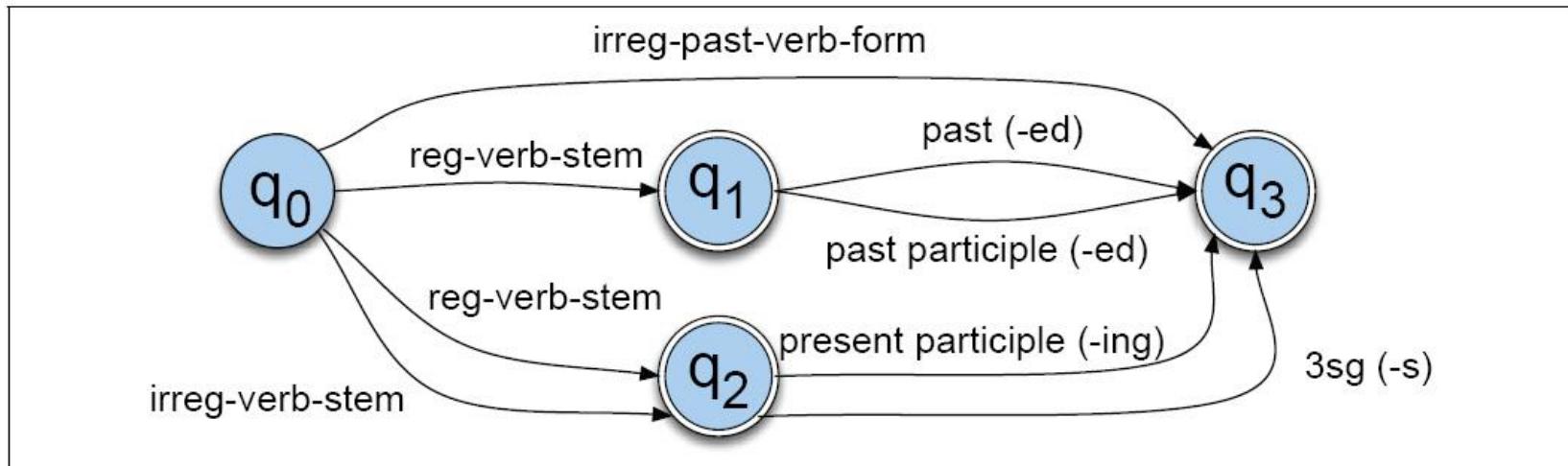
# Now Plug in the Words

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



# FSA for English Verb Inflection

reg-verb-stem	irreg-verb-stem	irreg-past-stem	past	past-part	pres-part	3sg
walk	cut	caught	-ed	-ed	-ing	-s
fry	speak	ate				
talk	sing	eaten				
impeach		sang				



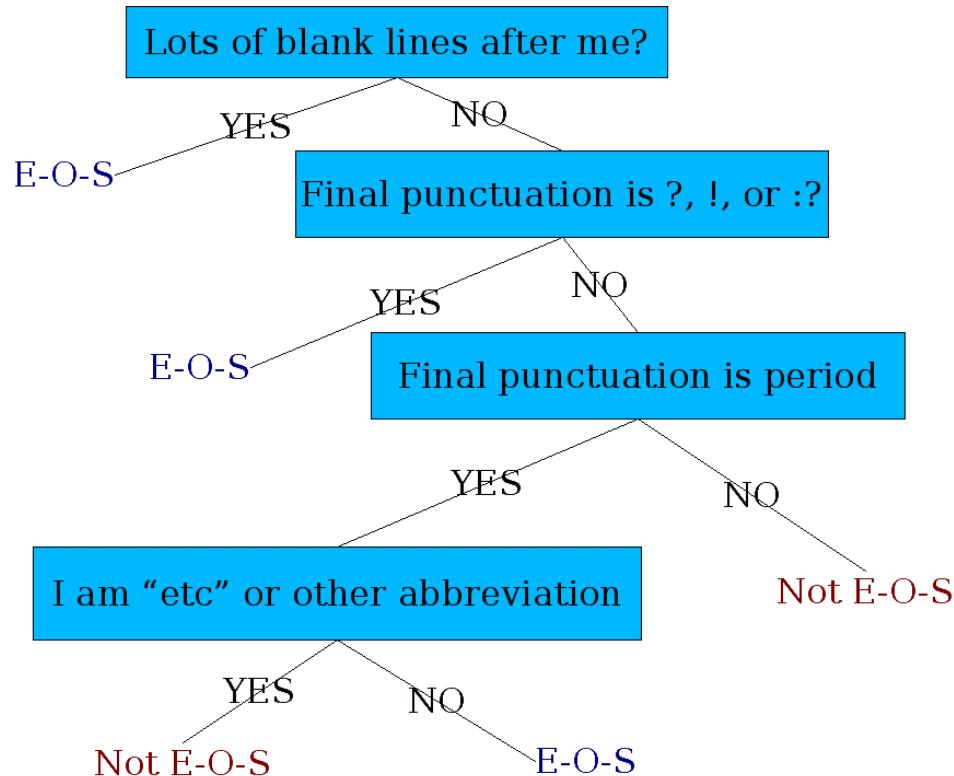
# Sentence Segmentation

---

- !, ? are relatively unambiguous
- Period “.” is quite ambiguous
  - Sentence boundary
  - Abbreviations like Inc. or Dr.
  - Numbers like .02% or 4.3
- Build a binary classifier
  - Looks at a “.”
  - Decides EndOfSentence/NotEndOfSentence
  - Classifiers: hand-written rules, regular expressions, or machine-learning

# Determining if a word is end-of-sentence: *a Decision Tree*

---



# More sophisticated decision tree features

---

- Case of word with ".": Upper, Lower, Cap, Number
- Case of word after ".": Upper, Lower, Cap, Number
- Numeric features
  - Length of word with “.”
  - Probability(word with “.” occurs at end-of-s)
  - Probability(word after “.” occurs at beginning-of-s)

# Implementing Decision Trees

---

- A decision tree is just an if-then-else statement
- The interesting research is choosing the features
- Setting up the structure is often too hard to do by hand
  - Hand-building only possible for very simple features, domains
    - For numeric features, it's too hard to pick each threshold
  - Instead, structure usually learned by machine learning from a training corpus

# Decision Trees and other classifiers

---

- We can think of the questions in a decision tree
- As features that could be exploited by any kind of classifier
  - Logistic regression
  - SVM
  - Neural Nets
  - etc.