

Regular Expressions and Finite State Automata

Mausam

(Based on slides by Jurafsky & Martin,
Julia Hirschberg)

Regular Expressions and Text Searching

- Everybody does it
 - ◆ Emacs, vi, perl, grep, etc..
- Regular expressions are a compact textual representation of a set of strings representing a language.

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

Figure 2.1

RE	Match	Example Patterns
<code>/[wW]oodchuck/</code>	Woodchuck or woodchuck	“ <u>W</u> oodchuck”
<code>/[abc]/</code>	‘a’, ‘b’, <i>or</i> ‘c’	“In uomini, in soldat <u>i</u> ”
<code>/[1234567890]/</code>	any digit	“plenty of <u>7</u> to 5”

Figure 2.2

RE	Match	Example Patterns Matched
/ [A-Z] /	an upper case letter	“we should call it ‘ <u>D</u> renched Blossoms’ ”
/ [a-z] /	a lower case letter	“ <u>m</u> y beans were impatient to be hoed!”
/ [0-9] /	a single digit	“Chapter <u>1</u> : Down the Rabbit Hole”

Figure 2.3

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

Figure 2.4

RE	Match	Example Patterns Matched
woodchucks?	woodchuck or woodchucks	“ <u>woodchuck</u> ”
colou?r	color or colour	“ <u>colour</u> ”

Figure 2.5

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

Parenthesis	()
Counters	* + ? { }
Sequences and anchors	the ^my end\$
Disjunction	

Figure 2.6

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

Figure 2.8

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	

Example

- Find all the instances of the word “the” in a text.
 - ♦ `/the/`
 - ♦ `/[tT]he/`
 - ♦ `/\b[tT]he\b/`

Errors

- The process we just went through was based on **two fixing 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

- We'll be telling the same story for many tasks, all semester. 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).

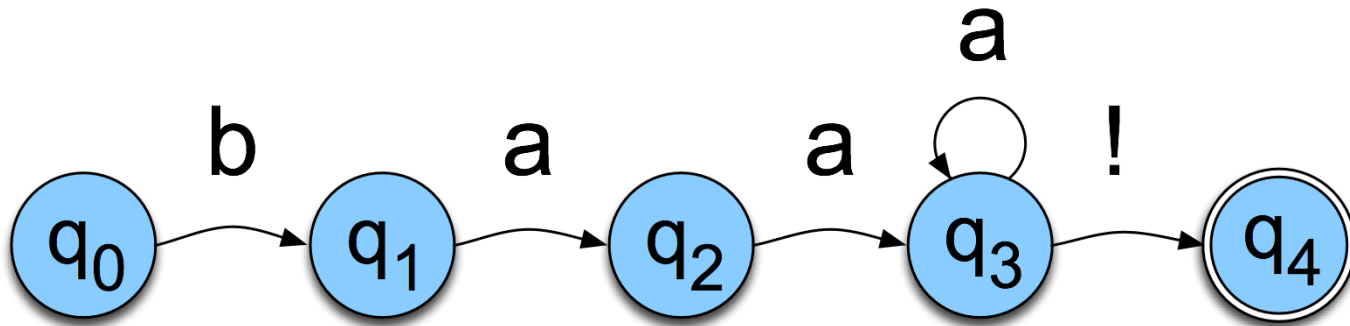
Finite State Automata

- Regular expressions can be viewed as a textual way of specifying the structure of finite-state automata.
- FSAs and their probabilistic relatives are at the core of much of what we'll be doing all semester.
- They also capture significant aspects of what linguists say we need for **morphology** and parts of **syntax**.

FSAs as Graphs

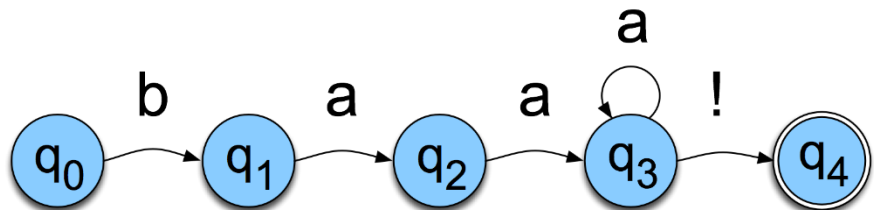
- Let's start with the sheep language from Chapter 2

♦ /baa+!/



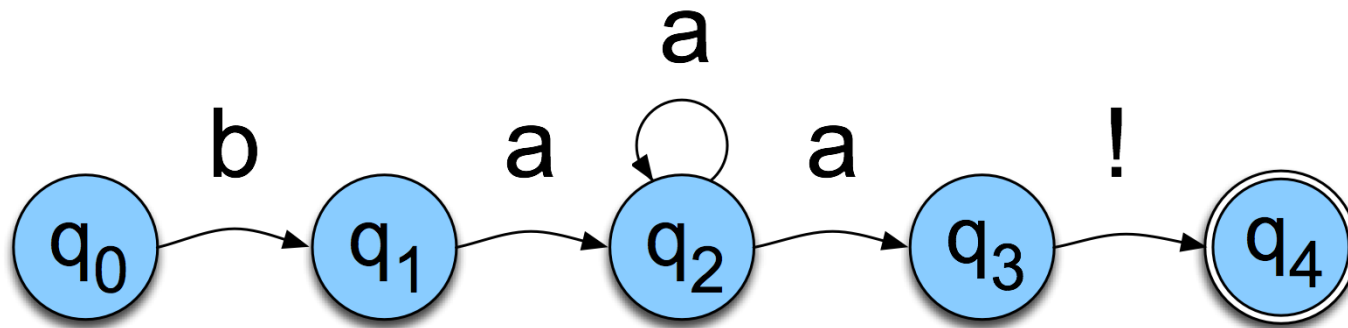
Sheep FSA

- We can say the following things about this machine
 - ♦ It has 5 states
 - ♦ **b**, **a**, and **!** are in its alphabet
 - ♦ q_0 is the start state
 - ♦ q_4 is an accept state
 - ♦ It has 5 transitions



But Note

- There are other machines that correspond to this same language



- More on this one later

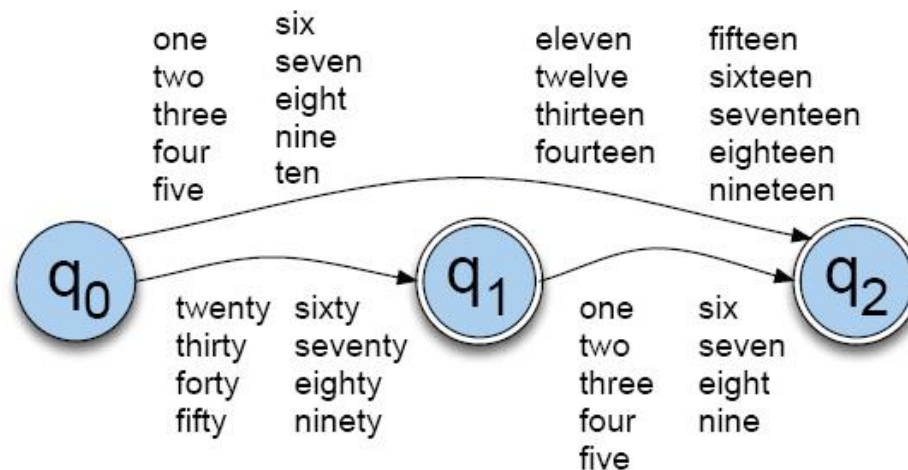
More Formally

- You can specify an FSA by enumerating the following things.
 - ♦ The set of states: Q
 - ♦ A finite alphabet: Σ
 - ♦ A start state
 - ♦ A set of accept/final states
 - ♦ A transition function that maps $Q \times \Sigma$ to Q

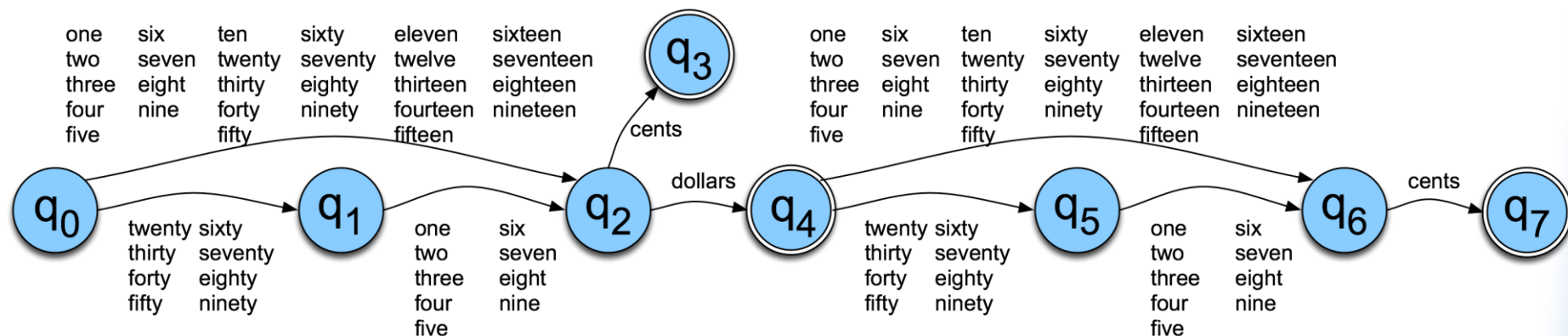
About Alphabets

- Don't take term *alphabet* word too narrowly; it just means we need a finite set of symbols in the input.
- These symbols can and will stand for bigger objects that can have internal structure.

Dollars and Cents



Dollars and Cents

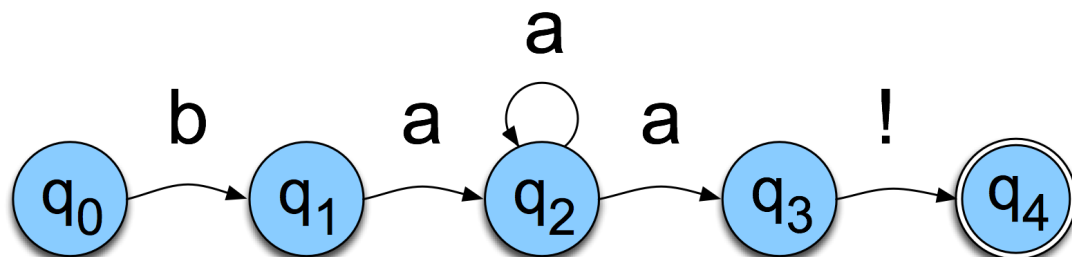


Yet Another View

- The guts of FSAs can ultimately be represented as tables

If you're in state 1 and you're looking at an a, go to state 2

	b	a	!	e
0	1			
1		2		
2		2,3		
3			4	
4				

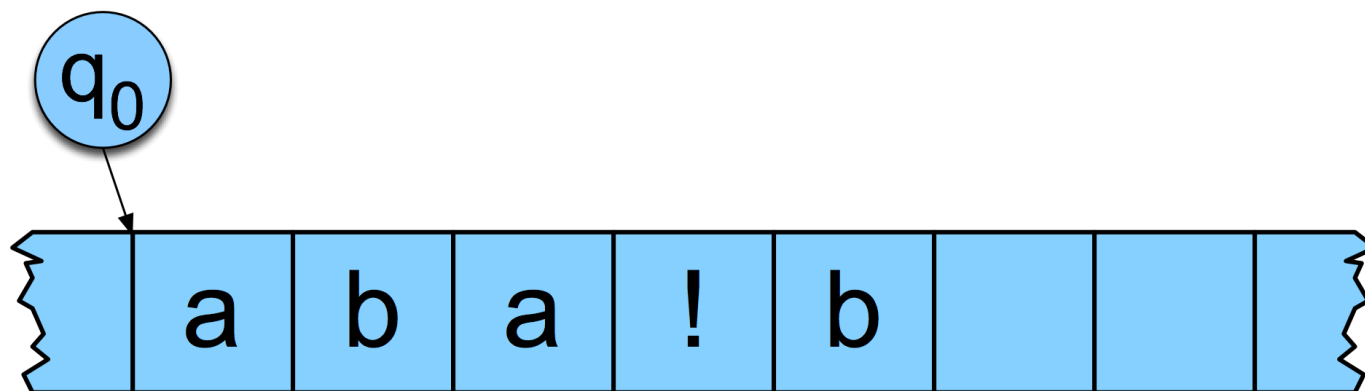


Recognition

- Recognition is the process of determining if a string should be accepted by a machine
- Or... it's the process of determining if a string is in the language we're defining with the machine
- Or... it's the process of determining if a regular expression matches a string
- Those all amount the same thing in the end

Recognition

- Traditionally, (Turing's notion) this process is depicted with a tape.



Recognition

- Simply a process of starting in the start state
- Examining the current input
- Consulting the table
- Going to a new state and updating the tape pointer.
- Until you run out of tape.

D-Recognize

function D-RECOGNIZE(*tape, machine*) **returns** accept or reject

index \leftarrow Beginning of tape

current-state \leftarrow Initial state of machine

loop

if End of input has been reached **then**

if *current-state* is an accept state **then**

return accept

else

return reject

elseif *transition-table*[*current-state*,*tape*[*index*]] is empty **then**

return reject

else

current-state \leftarrow *transition-table*[*current-state*,*tape*[*index*]]

index \leftarrow *index* + 1

end

Key Points

- Deterministic means that at each point in processing there is always one unique thing to do (no choices).
- D-recognize is a simple table-driven interpreter
- The algorithm is universal for all unambiguous regular languages.
 - ♦ To change the machine, you simply change the table.

Key Points

- Crudely therefore... matching strings with regular expressions (ala Perl, grep, etc.) is a matter of
 - ♦ translating the regular expression into a machine (a table) and
 - ♦ passing the table and the string to an interpreter

Recognition as Search

- You can view this algorithm as a trivial kind of *state-space search*.
- States are pairings of tape positions and state numbers.
- Operators are compiled into the table
- Goal state is a pairing with the end of tape position and a final accept state
- It is trivial because?

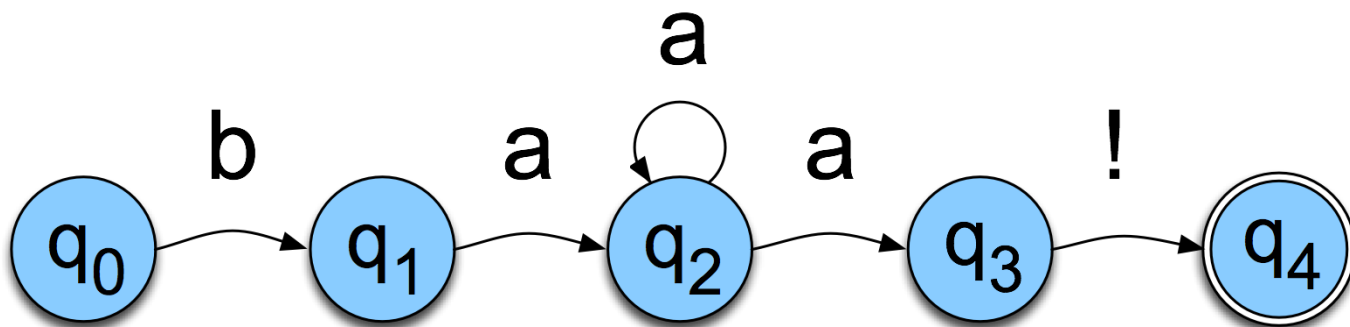
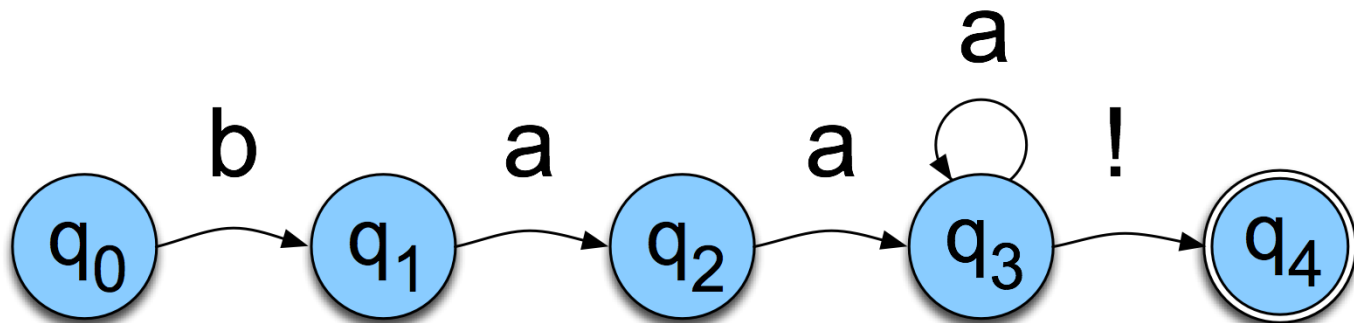
Generative Formalisms

- *Formal Languages* are sets of strings composed of symbols from a finite set of symbols.
- Finite-state automata define formal languages (without having to enumerate all the strings in the language)
- The term *Generative* is based on the view that you can run the machine as a generator to get strings from the language.

Generative Formalisms

- FSAs can be viewed from two perspectives:
 - ◆ Acceptors that can tell you if a string is in the language
 - ◆ Generators to produce *all and only* the strings in the language

Non-Determinism



Equivalence

- Non-deterministic machines can be converted to deterministic ones with a fairly simple construction
- That means that they have the same power; non-deterministic machines are not more powerful than deterministic ones in terms of the languages they can accept

ND Recognition

- Two basic approaches (used in all major implementations of regular expressions, see Friedl 2006)
 1. Either take a ND machine and convert it to a D machine and then do recognition with that.
 2. Or explicitly manage the process of recognition as a state-space search (leaving the machine as is).

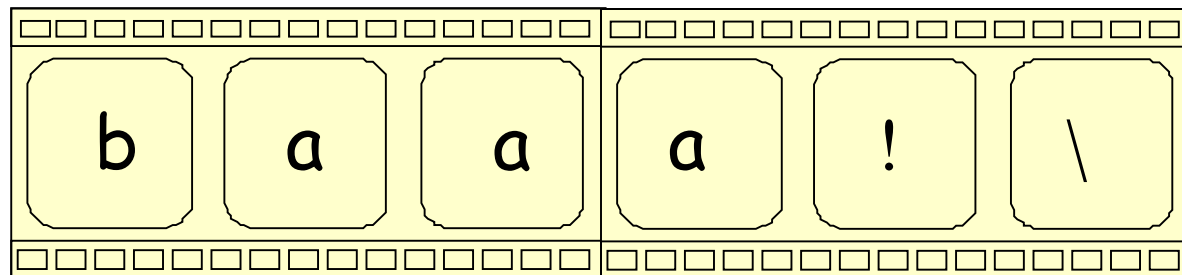
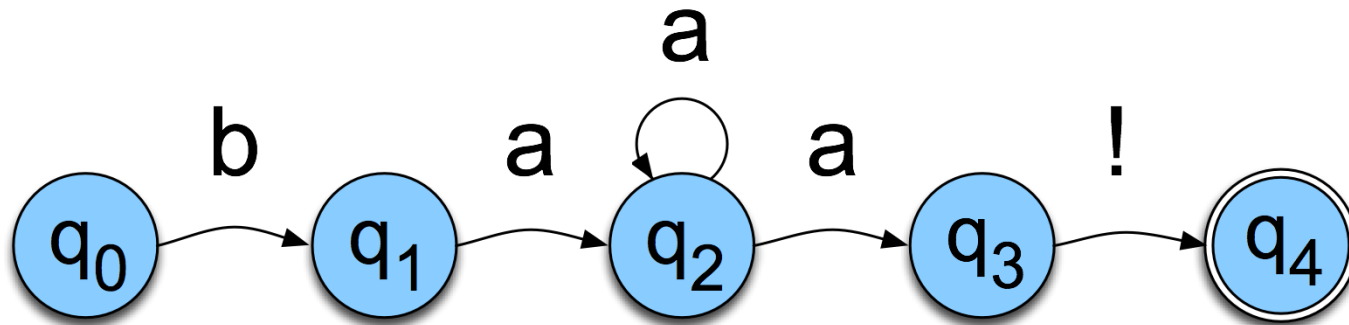
Non-Deterministic Recognition: Search

- In a ND FSA **there exists at least one path** through the machine for a string that is in the language defined by the machine.
- **But not all paths** directed through the machine for an accept string lead to an accept state.
- **No paths** through the machine lead to an accept state for a string not in the language.

Non-Deterministic Recognition

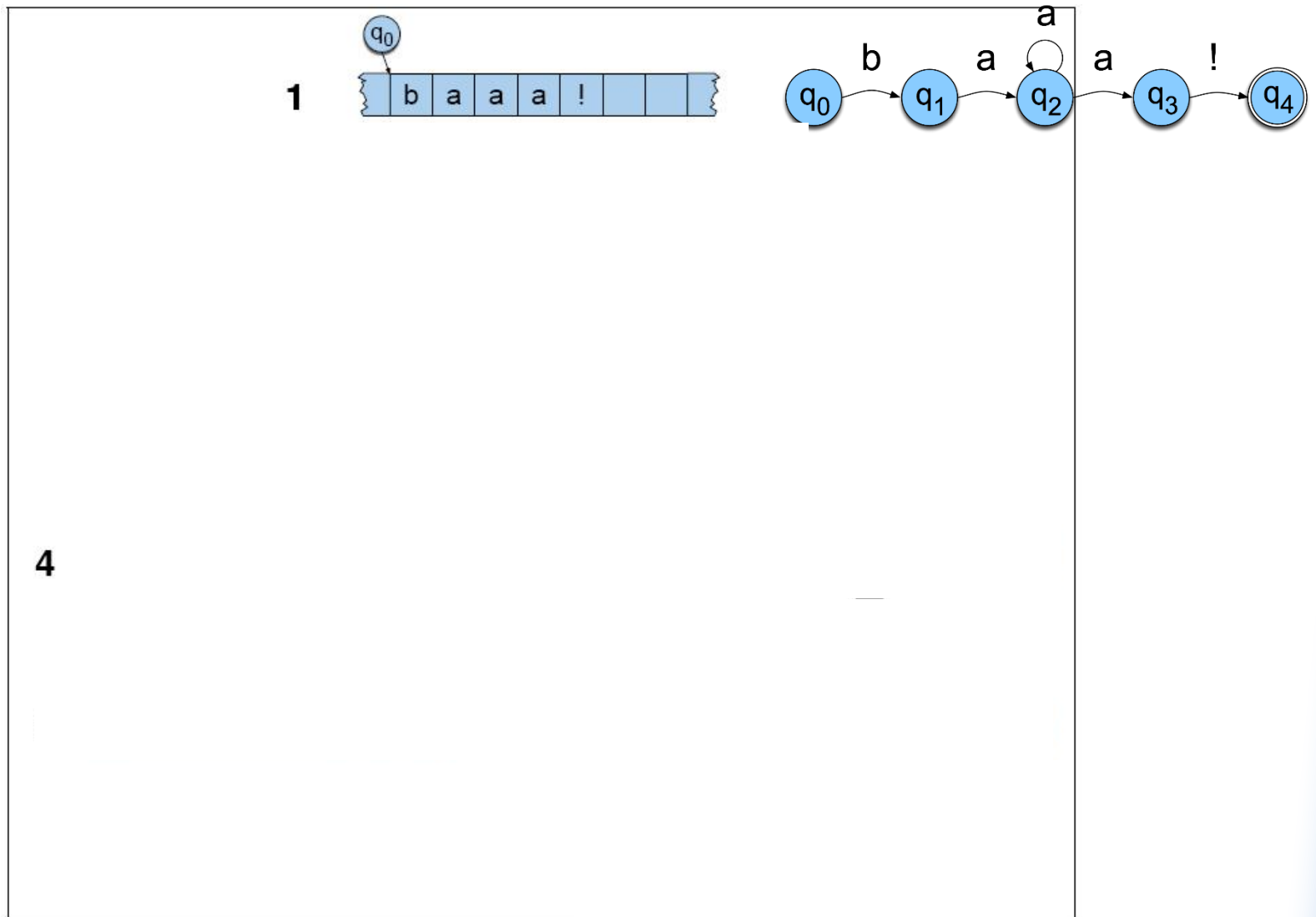
- So **success** in non-deterministic recognition occurs when a path is found through the machine that ends in an accept.
- **Failure** occurs when **all** of the possible paths for a given string lead to failure.

Example



q_0 q_1 q_2 q_2 q_3 q_4

Example



Key Points

- States in the search space are **pairings of tape positions and states** in the machine.
- By keeping track of **as yet unexplored states**, a recognizer can systematically explore all the paths through the machine given an input.

Why Bother?

- Non-determinism doesn't get us more formal power and it causes headaches so why bother?
 - ♦ More natural (understandable) solutions

FSTs (Contd)

FST-based Tokenization

```
#!/usr/bin/perl

$letternumber = "[A-Za-z0-9]";
$notletter = "[^A-Za-z0-9]";
$alwayssep = "[\\?!()\\\";\\/\\'|]";
$clitic = "('|:|-|'S|'D|'M|'LL|'RE|'VE|N'T|'s|'d|'m|'ll|'re|'ve|n't)";

$abbr{"Co."} = 1; $abbr{"Dr."} = 1; $abbr{"Jan."} = 1; $abbr{"Feb."} = 1;

while ($line = <>){ # read the next line from standard input

    # put whitespace around unambiguous separators
    $line =~ s/$alwayssep/ $& /g;

    # put whitespace around commas that aren't inside numbers
    $line =~ s/([0-9]),/$1 , /g;
    $line =~ s/,([0-9])/ , $1/g;

    # distinguish singlequotes from apostrophes by
    # segmenting off single quotes not preceded by letter
    $line =~ s/'/'$& /g;
    $line =~ s/($notletter)'/ $1 '/g;

    # segment off unambiguous word-final clitics and punctuation
    $line =~ s/$clitic$/ $&/g;
    $line =~ s/$clitic($notletter)/ $1 $2/g;

    # now deal with periods. For each possible word
    @possiblewords=split(/\s+/, $line);
    foreach $word (@possiblewords) {
        # if it ends in a period,
        if (($word =~ /$letternumber\./))
            && !($abbr{$word}) # and isn't on the abbreviation list
            # and isn't a sequence of letters and periods (U.S.)
            # and doesn't resemble an abbreviation (no vowels: Inc.)
            && !($word =~
                /\^[A-Za-z]\.([A-Za-z]\.)+|[A-Z][bcdafhj-nptvxz]+\.\.)/) {
            # then segment off the period
            $word =~ s/\.$/ \./;
        }
        # expand clitics
        $word =~ s/'ve/have/;
        $word =~ s/'m/am/;
        print $word, " ";
    }
    print "\n";
}
```

Porter Stemmer (1980)

- Common algorithm for stemming English
- Conventions + 5 phases of reductions
 - ◆ phases applied sequentially
 - ◆ each phase consists of a set of commands
 - ◆ sample convention: *Of the rules in a compound command, select the one that applies to the longest suffix.*

Porter Stemmer (1980)

- Standard, very popular and usable **stemmer** (IR, IE) – identify a word's **stem**
- Sequence of cascaded rewrite rules, e.g.
 - ♦ $IZE \rightarrow \varepsilon$ (e.g. unionize \rightarrow union)
 - ♦ $CY \rightarrow T$ (e.g. frequency \rightarrow frequent)
 - ♦ $ING \rightarrow \varepsilon$, if stem contains vowel (motoring \rightarrow motor)
- Can be implemented as a lexicon-free FST (many implementations available on the web)
- <http://text-processing.com/demo/stem/>

Summing Up

- Regular expressions and FSAs can represent subsets of natural language as well as regular languages
 - ♦ Both representations may be difficult for humans to use for any real subset of a language
 - ♦ Can be hard to scale up: e.g., when many choices at any point
 - ♦ But quick, powerful and easy to use for small problems
 - ♦ AT&T Finite State Toolkit **does** scale
- Finite state transducers and rules are common ways to incorporate linguistic ideas in NLP for small applications