

**CSCI 5832**  
**Natural Language**  
**Processing**

---

Lecture 3  
Jim Martin

1/24/08 1

---

---

---

---

---

---

---

---

**Today 1/22**

---

- Regexs, FSAs and languages
  - Determinism and Non-Determinism
- Combining FSAs
- English Morphology

1/24/08 2

---

---

---

---

---

---

---

---

**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 what we'll be doing all semester.
- They also conveniently (?) correspond closely to what linguists say we need for morphology and parts of syntax.
  - *Coincidence?*

1/24/08 3

---

---

---

---

---

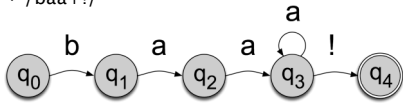
---

---

---

## FSA as Graphs

- Let's start with the sheep language from the text
  - /baa+!/



1/24/08

4

---

---

---

---

---

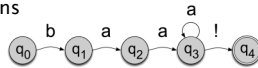
---

---

---

## Sheep FSA

- We can say the following things about this machine
  - It has 5 states
  - b, a, and ! are in its alphabet
  - q0 is the start state
  - q4 is an accept state
  - It has 5 transitions



1/24/08

5

---

---

---

---

---

---

---

---

## More Formally

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

1/24/08

6

---

---

---

---

---

---

---

---

## 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.

1/24/08

7

---

---

---

---

---

---

---

---

## 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

1/24/08

8

---

---

---

---

---

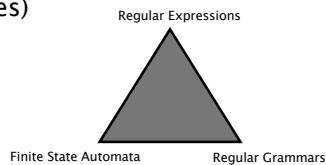
---

---

---

## Three Views

- Three equivalent formal ways to look at what we're up to (not including tables)



1/24/08

9

---

---

---

---

---

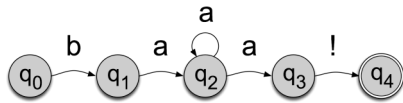
---

---

---

## But note

- There are other machines that correspond to this same language



- More on this one later

1/24/08

10

---

---

---

---

---

---

---

---

## About Alphabets

- Don't take that 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.

1/24/08

11

---

---

---

---

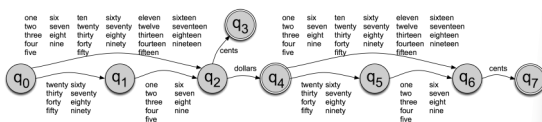
---

---

---

---

## Dollars and Cents



1/24/08

12

---

---

---

---

---

---

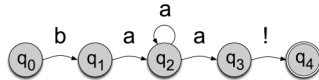
---

---

## QxΣ → Q

- The guts of FSAs can ultimately be represented as tables

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



1/24/08

13

---

---

---

---

---

---

---

---

---

---

## 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 defined by the machine
- Or... it's the process of determining if a regular expression matches a string
- Those all amount to the same thing in the end

1/24/08

14

---

---

---

---

---

---

---

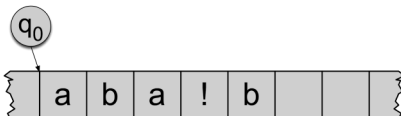
---

---

---

## Recognition

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



1/24/08

15

---

---

---

---

---

---

---

---

---

---

## 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.

1/24/08

16

---

---

---

---

---

---

---

---

## D-Recognize

```
function D-RECOGNIZE(tape, machine) returns accept or reject
  index ← Beginning of tape
  current-state ← 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
    elsif transition-table[current-state, tape[index]] is empty then
      return reject
    else
      current-state ← transition-table[current-state, tape[index]]
      index ← index + 1
  end
```

1/24/08

17

---

---

---

---

---

---

---

---

## Key Points

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

1/24/08

18

---

---

---

---

---

---

---

---

## 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 to an interpreter

1/24/08

19

---

---

---

---

---

---

---

---

## 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
- Its trivial because?

1/24/08

20

---

---

---

---

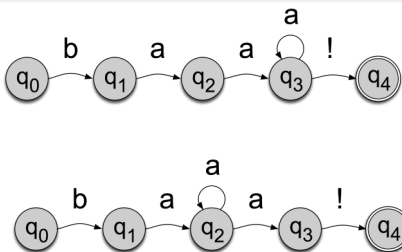
---

---

---

---

## Non-Determinism



1/24/08

21

---

---

---

---

---

---

---

---

## Non-Determinism

- Yet another technique
  - ♦ Epsilon transitions
  - ♦ Key point: these transitions do not examine or advance the tape during recognition

```

graph LR
    start(( )) --> q0((q0))
    q0 -- b --> q1((q1))
    q1 -- a --> q2((q2))
    q2 -- a --> q3((q3))
    q3 -- ! --> q4(((q4)))
    q3 -- ε --> q2
  
```

1/24/08 22

---

---

---

---

---

---

---

---

## 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 and can not accept

1/24/08 23

---

---

---

---

---

---

---

---

## ND Recognition

- Two basic approaches (used in all major implementations of Regular Expressions)
  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).

1/24/08 24

---

---

---

---

---

---

---

---



## Implementations

Program	(Original) Author	Version	Regex Engine
awk	Alex. Witteveen, Koenigsmann	generic	DFA
new_awk	Brian Kernighan	generic	DFA
GNU awk	Arnold Robbins	recent	Mostly DFA, some NFA
MS awk	Monette Kern Systems		POSIX NFA
asawk	Mike Brennan	all	POSIX NFA
agrep	Alfred Aho	generic	DFA
MS agrep	Monette Kern Systems		POSIX NFA
GNU Emacs	Richard Stallman	all	Trad. NFA (POSIX NFA available)
Expect	Don Libes	all	Traditional NFA
egrep	Dick Hight	generic	Traditional NFA
grep	Ken Thompson	generic	Traditional NFA
GNU grep	Mike Haertel	Version 2.0	Mostly DFA, but some NFA
GNU find	GNU		Traditional NFA
lex	Mike Lesk	generic	DFA
flex	Vern Paxson	all	DFA
flex	Monette Kern Systems		POSIX NFA
move	Eric Schlotzrood	generic	Traditional NFA
flex	Mark Nudelman		Variable (usually Trad. NFA)
Perl	Larry Wall	all	Traditional NFA
Python	Guido van Rossum	all	Traditional NFA
sed	Lee McMahon	generic	Traditional NFA
TED	John Ousterhout	all	Traditional NFA
tr	Bill Joy	generic	Traditional NFA

1/24/08

25

## 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.

1/24/08

26

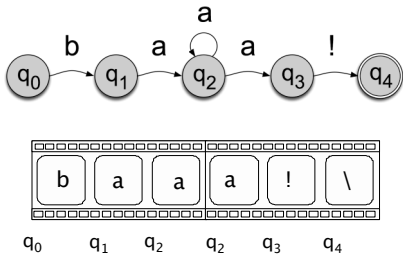
## Non-Deterministic Recognition

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

1/24/08

27

## Example



1/24/08

28

---

---

---

---

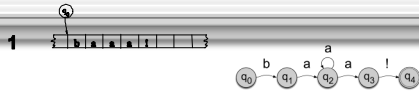
---

---

---

---

## Example



1/24/08

29

---

---

---

---

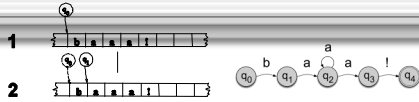
---

---

---

---

## Example



1/24/08

30

---

---

---

---

---

---

---

---

### Example

$q_0 \xrightarrow{b} q_1 \xrightarrow{a} q_2 \xrightarrow{a} q_2 \xrightarrow{!} q_3 \xrightarrow{\cdot} q_4$

1/24/08 31

---

---

---

---

---

---

---

---

### Example

1/24/08 32

---

---

---

---

---

---

---

---

### Example

1/24/08 33

---

---

---

---

---

---

---

---

### Example

1/24/08 34

---

---

---

---

---

---

---

---

### Example

1/24/08 35

---

---

---

---

---

---

---

---

### Example

1/24/08 36

---

---

---

---

---

---

---

---

## 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.

1/24/08

37

---

---

---

---

---

---

---

---

## ND-Recognize

```
function ND-RECOGNIZE(tape, machine) returns accept or reject
  agenda ← {(Initial state of machine, beginning of tape)}
  current-search-state ← NEXT(agenda)
  loop
    if ACCEPT-STATE?(current-search-state) returns true then
      return accept
    else
      agenda ← agenda ∪ GENERATE-NEW-STATES(current-search-state)
    if agenda is empty then
      return reject
    else
      current-search-state ← NEXT(agenda)
  end
```

1/24/08

38

---

---

---

---

---

---

---

---

## Infinite Search

- If you're not careful such searches can go into an infinite loop.
- How?

1/24/08

39

---

---

---

---

---

---

---

---

## Why Bother?

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

1/24/08

40

---

---

---

---

---

---

---

---

## Compositional Machines

- Formal languages are just sets of strings
- Therefore, we can talk about various set operations (intersection, union, concatenation)
- This turns out to be a useful exercise

1/24/08

41

---

---

---

---

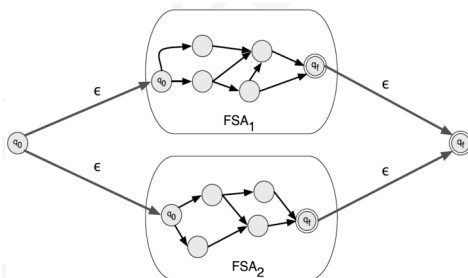
---

---

---

---

## Union



1/24/08

42

---

---

---

---

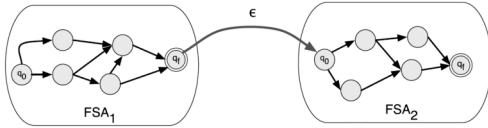
---

---

---

---

## Concatenation



1/24/08

43

---

---

---

---

---

---

---

---

## Negation

- Construct a machine M2 to accept all strings not accepted by machine M1 and reject all the strings accepted by M1
  - ♦ Invert all the accept and not accept states in M1
- Does that work for non-deterministic machines?

1/24/08

44

---

---

---

---

---

---

---

---

## Intersection

- Accept a string that is in both of two specified languages
- An indirect construction...
  - ♦  $A \cap B = \sim(\sim A \cup \sim B)$

1/24/08

45

---

---

---

---

---

---

---

---

## Motivation

- Consider the expression  
*Let's have a meeting on Thursday, Jan 26<sup>th</sup>*
  - Writing an FSA to recognize English date expressions is not terribly hard.
  - Except for the part about rejecting invalid dates.
  - Write two FSAs: one for the form of the dates, and one for the calendar arithmetic part
  - Intersect the two machines

1/24/08

46

---

---

---

---

---

---

---

---

## Next Time

- Finish Chapter 3

1/24/08

47

---

---

---

---

---

---

---

---