

**Exercise 2.3.5:** In the only-if portion of Theorem 2.12 we omitted the proof by induction on  $|w|$  that if  $\delta_D(q_0, w) = p$  then  $\delta_N(q_0, w) = \{p\}$ . Supply this proof.

**Exercise 2.3.6:** In the box on "Dead States and DFA's Missing Some Transitions," we claim that if  $N$  is an NFA that has at most one choice of state for any state and input symbol (i.e.,  $\delta(q, a)$  never has size greater than 1), then the DFA  $D$  constructed from  $N$  by the subset construction has exactly the states and transitions of  $N$  plus transitions to a new dead state whenever  $N$  is missing a transition for a given state and input symbol. Prove this contention.

**Exercise 2.3.7:** In Example 2.13 we claimed that the NFA  $N$  is in state  $q_i$ , for  $i = 1, 2, \dots, n$ , after reading input sequence  $w$  if and only if the  $i$ th symbol from the end of  $w$  is 1. Prove this claim.

## 2.4 An Application: Text Search

In this section, we shall see that the abstract study of the previous section, where we considered the "problem" of deciding whether a sequence of bits ends in 01, is actually an excellent model for several real problems that appear in applications such as Web search and extraction of information from text.

### 2.4.1 Finding Strings in Text

A common problem in the age of the Web and other on-line text repositories is the following. Given a set of words, find all documents that contain one (or all) of those words. A search engine is a popular example of this process. The search engine uses a particular technology, called *inverted indexes*, where for each word appearing on the Web (there are 100,000,000 different words), a list of all the places where that word occurs is stored. Machines with very large amounts of main memory keep the most common of these lists available, allowing many people to search for documents at once.

Inverted-index techniques do not make use of finite automata, but they also take very large amounts of time for crawlers to copy the Web and set up the indexes. There are a number of related applications that are unsuited for inverted indexes, but are good applications for automaton-based techniques. The characteristics that make an application suitable for searches that use automata are:

1. The repository on which the search is conducted is rapidly changing. For example:
  - (a) Every day, news analysts want to search the day's on-line news articles for relevant topics. For example, a financial analyst might search for certain stock ticker symbols or names of companies.

## 2.4 AN APPLICATION: TEXT SEARCH

(b) A "shopping robot" wants to search for the current prices charged for the items that its clients request. The robot will retrieve current catalog pages from the Web and then search those pages for words that suggest a price for a particular item.

2. The documents to be searched cannot be cataloged. For example, Amazon.com does not make it easy for crawlers to find all the pages for all the books that the company sells. Rather, these pages are generated "on the fly" in response to queries. However, we could send a query for books on a certain topic, say "finite automata," and then search the pages retrieved for certain words, e.g., "excellent" in a review portion.

### 2.4.2 Nondeterministic Finite Automata for Text Search

Suppose we are given a set of words, which we shall call the *keywords*, and we want to find occurrences of any of these words. In applications such as these, a useful way to proceed is to design a nondeterministic finite automaton, which signals, by entering an accepting state, that it has seen one of the keywords. The text of a document is fed, one character at a time to this NFA, which then recognizes occurrences of the keywords in this text. There is a simple form to an NFA that recognizes a set of keywords.

1. There is a start state with a transition to itself on every input symbol, e.g. every printable ASCII character if we are examining text. Intuitively, the start state represents a "guess" that we have not yet begun to see one of the keywords, even if we have seen some letters of one of these words.
2. For each keyword  $a_1 a_2 \dots a_k$ , there are  $k$  states, say  $q_1, q_2, \dots, q_k$ . There is a transition from the start state to  $q_1$  on symbol  $a_1$ , a transition from  $q_1$  to  $q_2$  on symbol  $a_2$ , and so on. The state  $q_k$  is an accepting state and indicates that the keyword  $a_1 a_2 \dots a_k$  has been found.

**Example 2.14:** Suppose we want to design an NFA to recognize occurrences of the words web and ebay. The transition diagram for the NFA designed using the rules above is in Fig. 2.16. State 1 is the start state, and we use  $\Sigma$  to stand for the set of all printable ASCII characters. States 2 through 4 have the job of recognizing web, while states 5 through 8 recognize ebay.  $\square$

Of course the NFA is not a program. We have two major choices for an implementation of this NFA.

1. Write a program that simulates this NFA by computing the set of states it is in after reading each input symbol. The simulation was suggested in Fig. 2.10.
2. Convert the NFA to an equivalent DFA using the subset construction. Then simulate the DFA directly.



$p_i$ 's on symbol  $x$ . On all symbols  $x$  such that there are no transitions out of any of the  $p_i$ 's on symbol  $x$ , let this DFA state have a transition on  $x$  to that state of the DFA consisting of  $q_0$  and all states that are reached from  $q_0$  in the NFA following an arc labeled  $x$ .

For instance, consider state 135 of Fig. 2.17. The NFA of Fig. 2.16 has transitions on symbol  $b$  from states 3 and 5 to states 4 and 6, respectively. Therefore, on symbol  $b$ , 135 goes to 146. On symbol  $e$ , there are no transitions of the NFA out of 3 or 5, but there is a transition from 1 to 5. Thus, in the DFA, 135 goes to 15 on input  $e$ . Similarly, on input  $w$ , 135 goes to 12.

On every other symbol  $x$ , there are no transitions out of 3 or 5, and state 1 goes only to itself. Thus, there are transitions from 135 to 1 on every symbol in  $\Sigma$  other than  $b$ ,  $e$ , and  $w$ . We use the notation  $\Sigma - b - e - w$  to represent this set, and use similar representations of other sets in which a few symbols are removed from  $\Sigma$ .  $\square$

## 2.4.4 Exercises for Section 2.4

Exercise 2.4.1: Design NFA's to recognize the following sets of strings.

- \* a)  $abc, abd, \text{ and } aacd$ . Assume the alphabet is  $\{a, b, c, d\}$ .
- b)  $0101, 101, \text{ and } 011$ .
- c)  $ab, bc, \text{ and } ca$ . Assume the alphabet is  $\{a, b, c\}$ .

Exercise 2.4.2: Convert each of your NFA's from Exercise 2.4.1 to DFA's.

## 2.5 Finite Automata With Epsilon-Transitions

We shall now introduce another extension of the finite automaton. The new "feature" is that we allow a transition on  $\epsilon$ , the empty string. In effect, an NFA is allowed to make a transition spontaneously, without receiving an input symbol. Like the nondeterminism added in Section 2.3, this new capability does not expand the class of languages that can be accepted by finite automata, but it does give us some added "programming convenience." We shall also see, when we take up regular expressions in Section 3.1, how NFA's with  $\epsilon$ -transitions, which we call  $\epsilon$ -NFA's, are closely related to regular expressions and useful in proving the equivalence between the classes of languages accepted by finite automata and by regular expressions.

### 2.5.1 Uses of $\epsilon$ -Transitions

We shall begin with an informal treatment of  $\epsilon$ -NFA's, using transition diagrams with  $\epsilon$  allowed as a label. In the examples to follow, think of the automaton as accepting those sequences of labels along paths from the start state to an accepting state. However, each  $\epsilon$  along a path is "invisible", i.e., it contributes nothing to the string along the path.

## 2.5. FINITE AUTOMATA WITH EPSILON-TRANSITIONS

Example 2.16: In Fig. 2.18 is an  $\epsilon$ -NFA that accepts decimal numbers consisting of:

1. An optional + or - sign.
2. A string of digits,
3. A decimal point, and
4. Another string of digits. Either this string of digits, or the string (2) can be empty, but at least one of the two strings of digits must be nonempty.

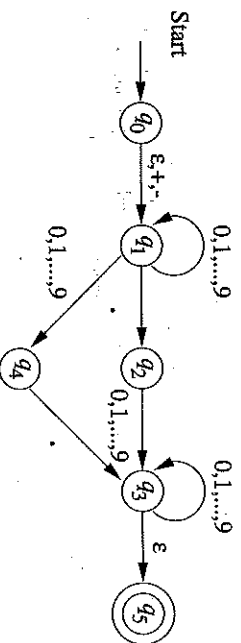


Figure 2.18: An  $\epsilon$ -NFA accepting decimal numbers

Of particular interest is the transition from  $q_0$  to  $q_1$  on any of  $\epsilon$ ,  $+$ , or  $-$ . Thus, state  $q_1$  represents the situation in which we have seen the sign if there is one, and perhaps some digits, but not the decimal point. State  $q_3$  represents the situation where we have just seen the decimal point, and may or may not have seen prior digits. In  $q_4$  we have definitely seen at least one digit, but not the decimal point. Thus, the interpretation of  $q_3$  is that we have seen a decimal point and at least one digit, either before or after the decimal point. We may stay in  $q_3$  reading whatever digits there are, and also have the option of "guessing" the string of digits is complete and going spontaneously to  $q_5$ , the accepting state.  $\square$

Example 2.17: The strategy we outlined in Example 2.14 for building an NFA that recognizes a set of keywords can be simplified further if we allow  $\epsilon$ -transitions. For instance, the NFA recognizing the keywords `web` and `ebay`, which we saw in Fig. 2.16, can also be implemented with  $\epsilon$ -transitions as in Fig. 2.19. In general, we construct a complete sequence of states for each keyword, as if it were the only word the automaton needed to recognize. Then, we add a new start state (state 9 in Fig. 2.19), with  $\epsilon$ -transitions to the start-states of the automata for each of the keywords.  $\square$