

**32.2-4**

Alice has a copy of a long  $n$ -bit file  $A = \langle a_{n-1}, a_{n-2}, \dots, a_0 \rangle$ , and Bob similarly has an  $n$ -bit file  $B = \langle b_{n-1}, b_{n-2}, \dots, b_0 \rangle$ . Alice and Bob wish to know if their files are identical. To avoid transmitting all of  $A$  or  $B$ , they use the following fast probabilistic check. Together, they select a prime  $q > 1000n$  and randomly select an integer  $x$  from  $\{0, 1, \dots, q-1\}$ . Then, Alice evaluates

$$A(x) = \left( \sum_{i=0}^{n-1} a_i x^i \right) \bmod q$$

and Bob similarly evaluates  $B(x)$ . Prove that if  $A \neq B$ , there is at most one chance in 1000 that  $A(x) = B(x)$ , whereas if the two files are the same,  $A(x)$  is necessarily the same as  $B(x)$ . (*Hint*: See Exercise 31.4-4.)

---

**32.3 String matching with finite automata**

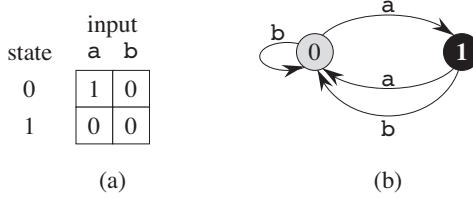
Many string-matching algorithms build a finite automaton—a simple machine for processing information—that scans the text string  $T$  for all occurrences of the pattern  $P$ . This section presents a method for building such an automaton. These string-matching automata are very efficient: they examine each text character *exactly once*, taking constant time per text character. The matching time used—after preprocessing the pattern to build the automaton—is therefore  $\Theta(n)$ . The time to build the automaton, however, can be large if  $\Sigma$  is large. Section 32.4 describes a clever way around this problem.

We begin this section with the definition of a finite automaton. We then examine a special string-matching automaton and show how to use it to find occurrences of a pattern in a text. Finally, we shall show how to construct the string-matching automaton for a given input pattern.

**Finite automata**

A *finite automaton*  $M$ , illustrated in Figure 32.6, is a 5-tuple  $(Q, q_0, A, \Sigma, \delta)$ , where

- $Q$  is a finite set of *states*,
- $q_0 \in Q$  is the *start state*,
- $A \subseteq Q$  is a distinguished set of *accepting states*,
- $\Sigma$  is a finite *input alphabet*,
- $\delta$  is a function from  $Q \times \Sigma$  into  $Q$ , called the *transition function* of  $M$ .



**Figure 32.6** A simple two-state finite automaton with state set  $Q = \{0, 1\}$ , start state  $q_0 = 0$ , and input alphabet  $\Sigma = \{a, b\}$ . **(a)** A tabular representation of the transition function  $\delta$ . **(b)** An equivalent state-transition diagram. State 1, shown blackend, is the only accepting state. Directed edges represent transitions. For example, the edge from state 1 to state 0 labeled **b** indicates that  $\delta(1, b) = 0$ . This automaton accepts those strings that end in an odd number of **a**'s. More precisely, it accepts a string  $x$  if and only if  $x = yz$ , where  $y = \varepsilon$  or  $y$  ends with a **b**, and  $z = a^k$ , where  $k$  is odd. For example, on input **abaaaa**, including the start state, this automaton enters the sequence of states  $\{0, 1, 0, 1, 0, 1\}$ , and so it accepts this input. For input **abbbaa**, it enters the sequence of states  $\{0, 1, 0, 0, 0, 1, 0\}$ , and so it rejects this input.

The finite automaton begins in state  $q_0$  and reads the characters of its input string one at a time. If the automaton is in state  $q$  and reads input character  $a$ , it moves (“makes a transition”) from state  $q$  to state  $\delta(q, a)$ . Whenever its current state  $q$  is a member of  $A$ , the machine  $M$  has **accepted** the string read so far. An input that is not accepted is **rejected**.

A finite automaton  $M$  induces a function  $\phi$ , called the **final-state function**, from  $\Sigma^*$  to  $Q$  such that  $\phi(w)$  is the state  $M$  ends up in after scanning the string  $w$ . Thus,  $M$  accepts a string  $w$  if and only if  $\phi(w) \in A$ . We define the function  $\phi$  recursively, using the transition function:

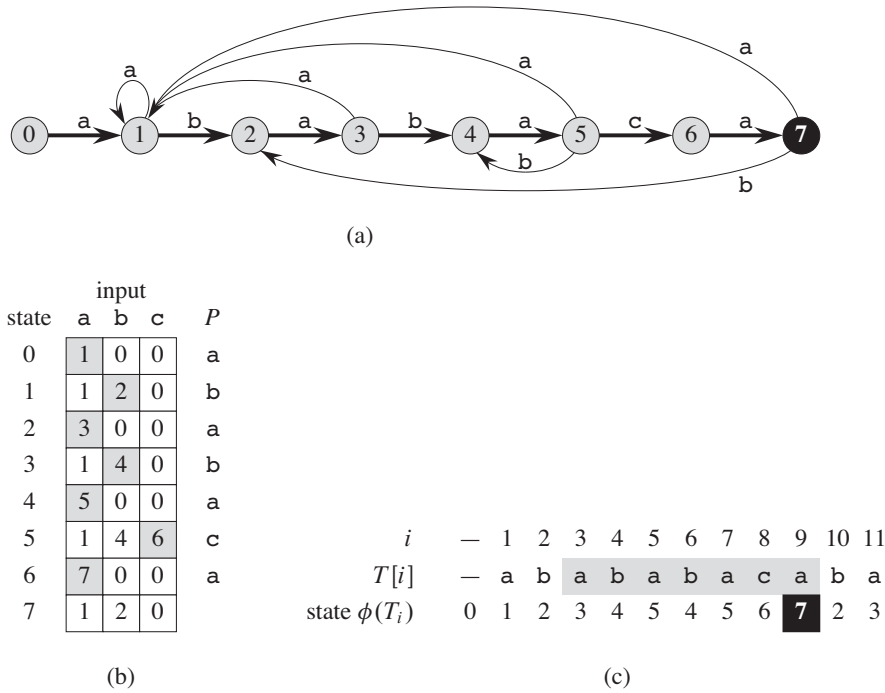
$$\begin{aligned}\phi(\varepsilon) &= q_0, \\ \phi(wa) &= \delta(\phi(w), a) \quad \text{for } w \in \Sigma^*, a \in \Sigma.\end{aligned}$$

### String-matching automata

For a given pattern  $P$ , we construct a string-matching automaton in a preprocessing step before using it to search the text string. Figure 32.7 illustrates how we construct the automaton for the pattern  $P = ababaca$ . From now on, we shall assume that  $P$  is a given fixed pattern string; for brevity, we shall not indicate the dependence upon  $P$  in our notation.

In order to specify the string-matching automaton corresponding to a given pattern  $P[1..m]$ , we first define an auxiliary function  $\sigma$ , called the **suffix function** corresponding to  $P$ . The function  $\sigma$  maps  $\Sigma^*$  to  $\{0, 1, \dots, m\}$  such that  $\sigma(x)$  is the length of the longest prefix of  $P$  that is also a suffix of  $x$ :

$$\sigma(x) = \max \{k : P_k \sqsubset x\} . \quad (32.3)$$



**Figure 32.7** (a) A state-transition diagram for the string-matching automaton that accepts all strings ending in the string **ababaca**. State 0 is the start state, and state 7 (shown blackened) is the only accepting state. A directed edge from state  $i$  to state  $j$  labeled  $a$  represents  $\delta(i, a) = j$ . The right-going edges forming the “spine” of the automaton, shown heavy in the figure, correspond to successful matches between pattern and input characters. The left-going edges correspond to failing matches. Some edges corresponding to failing matches are omitted; by convention, if a state  $i$  has no outgoing edge labeled  $a$  for some  $a \in \Sigma$ , then  $\delta(i, a) = 0$ . (b) The corresponding transition function  $\delta$ , and the pattern string  $P = \mathbf{ababaca}$ . The entries corresponding to successful matches between pattern and input characters are shown shaded. (c) The operation of the automaton on the text  $T = \mathbf{abababacaba}$ . Under each text character  $T[i]$  appears the state  $\phi(T_i)$  that the automaton is in after processing the prefix  $T_i$ . The automaton finds one occurrence of the pattern, ending in position 9.

The suffix function  $\sigma$  is well defined since the empty string  $P_0 = \varepsilon$  is a suffix of every string. As examples, for the pattern  $P = \mathbf{ab}$ , we have  $\sigma(\varepsilon) = 0$ ,  $\sigma(\mathbf{ccaca}) = 1$ , and  $\sigma(\mathbf{ccab}) = 2$ . For a pattern  $P$  of length  $m$ , we have  $\sigma(x) = m$  if and only if  $P \sqsubset x$ . From the definition of the suffix function,  $x \sqsubset y$  implies  $\sigma(x) \leq \sigma(y)$ .

We define the string-matching automaton that corresponds to a given pattern  $P[1..m]$  as follows:

- The state set  $Q$  is  $\{0, 1, \dots, m\}$ . The start state  $q_0$  is state 0, and state  $m$  is the only accepting state.
- The transition function  $\delta$  is defined by the following equation, for any state  $q$  and character  $a$ :

$$\delta(q, a) = \sigma(P_q a) . \quad (32.4)$$

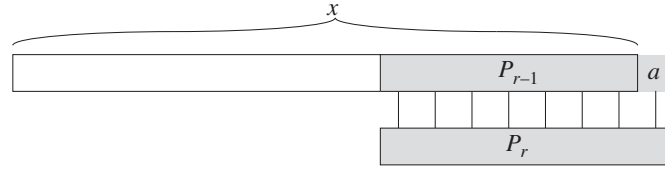
We define  $\delta(q, a) = \sigma(P_q a)$  because we want to keep track of the longest prefix of the pattern  $P$  that has matched the text string  $T$  so far. We consider the most recently read characters of  $T$ . In order for a substring of  $T$ —let’s say the substring ending at  $T[i]$ —to match some prefix  $P_j$  of  $P$ , this prefix  $P_j$  must be a suffix of  $T_i$ . Suppose that  $q = \phi(T_i)$ , so that after reading  $T_i$ , the automaton is in state  $q$ . We design the transition function  $\delta$  so that this state number,  $q$ , tells us the length of the longest prefix of  $P$  that matches a suffix of  $T_i$ . That is, in state  $q$ ,  $P_q \sqsubset T_i$  and  $q = \sigma(T_i)$ . (Whenever  $q = m$ , all  $m$  characters of  $P$  match a suffix of  $T_i$ , and so we have found a match.) Thus, since  $\phi(T_i)$  and  $\sigma(T_i)$  both equal  $q$ , we shall see (in Theorem 32.4, below) that the automaton maintains the following invariant:

$$\phi(T_i) = \sigma(T_i) . \quad (32.5)$$

If the automaton is in state  $q$  and reads the next character  $T[i + 1] = a$ , then we want the transition to lead to the state corresponding to the longest prefix of  $P$  that is a suffix of  $T_i a$ , and that state is  $\sigma(T_i a)$ . Because  $P_q$  is the longest prefix of  $P$  that is a suffix of  $T_i$ , the longest prefix of  $P$  that is a suffix of  $T_i a$  is not only  $\sigma(T_i a)$ , but also  $\sigma(P_q a)$ . (Lemma 32.3, on page 1000, proves that  $\sigma(T_i a) = \sigma(P_q a)$ .) Thus, when the automaton is in state  $q$ , we want the transition function on character  $a$  to take the automaton to state  $\sigma(P_q a)$ .

There are two cases to consider. In the first case,  $a = P[q + 1]$ , so that the character  $a$  continues to match the pattern; in this case, because  $\delta(q, a) = q + 1$ , the transition continues to go along the “spine” of the automaton (the heavy edges in Figure 32.7). In the second case,  $a \neq P[q + 1]$ , so that  $a$  does not continue to match the pattern. Here, we must find a smaller prefix of  $P$  that is also a suffix of  $T_i$ . Because the preprocessing step matches the pattern against itself when creating the string-matching automaton, the transition function quickly identifies the longest such smaller prefix of  $P$ .

Let’s look at an example. The string-matching automaton of Figure 32.7 has  $\delta(5, c) = 6$ , illustrating the first case, in which the match continues. To illustrate the second case, observe that the automaton of Figure 32.7 has  $\delta(5, b) = 4$ . We make this transition because if the automaton reads a **b** in state  $q = 5$ , then  $P_q b = ababab$ , and the longest prefix of  $P$  that is also a suffix of **ababab** is  $P_4 = abab$ .



**Figure 32.8** An illustration for the proof of Lemma 32.2. The figure shows that  $r \leq \sigma(x) + 1$ , where  $r = \sigma(xa)$ .

To clarify the operation of a string-matching automaton, we now give a simple, efficient program for simulating the behavior of such an automaton (represented by its transition function  $\delta$ ) in finding occurrences of a pattern  $P$  of length  $m$  in an input text  $T[1..n]$ . As for any string-matching automaton for a pattern of length  $m$ , the state set  $Q$  is  $\{0, 1, \dots, m\}$ , the start state is 0, and the only accepting state is state  $m$ .

FINITE-AUTOMATON-MATCHER( $T, \delta, m$ )

```

1   $n = T.length$ 
2   $q = 0$ 
3  for  $i = 1$  to  $n$ 
4       $q = \delta(q, T[i])$ 
5      if  $q == m$ 
6          print "Pattern occurs with shift"  $i - m$ 
```

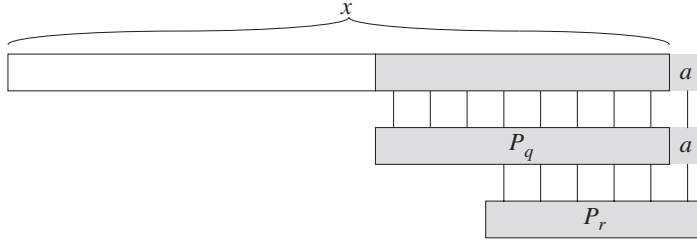
From the simple loop structure of FINITE-AUTOMATON-MATCHER, we can easily see that its matching time on a text string of length  $n$  is  $\Theta(n)$ . This matching time, however, does not include the preprocessing time required to compute the transition function  $\delta$ . We address this problem later, after first proving that the procedure FINITE-AUTOMATON-MATCHER operates correctly.

Consider how the automaton operates on an input text  $T[1..n]$ . We shall prove that the automaton is in state  $\sigma(T_i)$  after scanning character  $T[i]$ . Since  $\sigma(T_i) = m$  if and only if  $P \sqsubset T_i$ , the machine is in the accepting state  $m$  if and only if it has just scanned the pattern  $P$ . To prove this result, we make use of the following two lemmas about the suffix function  $\sigma$ .

**Lemma 32.2 (Suffix-function inequality)**

For any string  $x$  and character  $a$ , we have  $\sigma(xa) \leq \sigma(x) + 1$ .

**Proof** Referring to Figure 32.8, let  $r = \sigma(xa)$ . If  $r = 0$ , then the conclusion  $\sigma(xa) = r \leq \sigma(x) + 1$  is trivially satisfied, by the nonnegativity of  $\sigma(x)$ . Now assume that  $r > 0$ . Then,  $P_r \sqsubset xa$ , by the definition of  $\sigma$ . Thus,  $P_{r-1} \sqsubset x$ , by



**Figure 32.9** An illustration for the proof of Lemma 32.3. The figure shows that  $r = \sigma(P_q a)$ , where  $q = \sigma(x)$  and  $r = \sigma(xa)$ .

dropping the  $a$  from the end of  $P_r$  and from the end of  $xa$ . Therefore,  $r-1 \leq \sigma(x)$ , since  $\sigma(x)$  is the largest  $k$  such that  $P_k \sqsubset x$ , and thus  $\sigma(xa) = r \leq \sigma(x) + 1$ . ■

**Lemma 32.3 (Suffix-function recursion lemma)**

For any string  $x$  and character  $a$ , if  $q = \sigma(x)$ , then  $\sigma(xa) = \sigma(P_q a)$ .

**Proof** From the definition of  $\sigma$ , we have  $P_q \sqsubset x$ . As Figure 32.9 shows, we also have  $P_q a \sqsubset xa$ . If we let  $r = \sigma(xa)$ , then  $P_r \sqsubset xa$  and, by Lemma 32.2,  $r \leq q + 1$ . Thus, we have  $|P_r| = r \leq q + 1 = |P_q a|$ . Since  $P_q a \sqsubset xa$ ,  $P_r \sqsubset xa$ , and  $|P_r| \leq |P_q a|$ , Lemma 32.1 implies that  $P_r \sqsubset P_q a$ . Therefore,  $r \leq \sigma(P_q a)$ , that is,  $\sigma(xa) \leq \sigma(P_q a)$ . But we also have  $\sigma(P_q a) \leq \sigma(xa)$ , since  $P_q a \sqsubset xa$ . Thus,  $\sigma(xa) = \sigma(P_q a)$ . ■

We are now ready to prove our main theorem characterizing the behavior of a string-matching automaton on a given input text. As noted above, this theorem shows that the automaton is merely keeping track, at each step, of the longest prefix of the pattern that is a suffix of what has been read so far. In other words, the automaton maintains the invariant (32.5).

**Theorem 32.4**

If  $\phi$  is the final-state function of a string-matching automaton for a given pattern  $P$  and  $T[1..n]$  is an input text for the automaton, then

$$\phi(T_i) = \sigma(T_i)$$

for  $i = 0, 1, \dots, n$ .

**Proof** The proof is by induction on  $i$ . For  $i = 0$ , the theorem is trivially true, since  $T_0 = \varepsilon$ . Thus,  $\phi(T_0) = 0 = \sigma(T_0)$ .

Now, we assume that  $\phi(T_i) = \sigma(T_i)$  and prove that  $\phi(T_{i+1}) = \sigma(T_{i+1})$ . Let  $q$  denote  $\phi(T_i)$ , and let  $a$  denote  $T[i + 1]$ . Then,

$$\begin{aligned}
 \phi(T_{i+1}) &= \phi(T_i a) && \text{(by the definitions of } T_{i+1} \text{ and } a) \\
 &= \delta(\phi(T_i), a) && \text{(by the definition of } \phi) \\
 &= \delta(q, a) && \text{(by the definition of } q) \\
 &= \sigma(P_q a) && \text{(by the definition (32.4) of } \delta) \\
 &= \sigma(T_i a) && \text{(by Lemma 32.3 and induction)} \\
 &= \sigma(T_{i+1}) && \text{(by the definition of } T_{i+1}) \quad \blacksquare
 \end{aligned}$$

By Theorem 32.4, if the machine enters state  $q$  on line 4, then  $q$  is the largest value such that  $P_q \sqsupseteq T_i$ . Thus, we have  $q = m$  on line 5 if and only if the machine has just scanned an occurrence of the pattern  $P$ . We conclude that FINITE-AUTOMATON-MATCHER operates correctly.

### Computing the transition function

The following procedure computes the transition function  $\delta$  from a given pattern  $P[1..m]$ .

COMPUTE-TRANSITION-FUNCTION( $P, \Sigma$ )

```

1   $m = P.length$ 
2  for  $q = 0$  to  $m$ 
3      for each character  $a \in \Sigma$ 
4           $k = \min(m + 1, q + 2)$ 
5          repeat
6               $k = k - 1$ 
7          until  $P_k \sqsupseteq P_q a$ 
8           $\delta(q, a) = k$ 
9  return  $\delta$ 

```

This procedure computes  $\delta(q, a)$  in a straightforward manner according to its definition in equation (32.4). The nested loops beginning on lines 2 and 3 consider all states  $q$  and all characters  $a$ , and lines 4–8 set  $\delta(q, a)$  to be the largest  $k$  such that  $P_k \sqsupseteq P_q a$ . The code starts with the largest conceivable value of  $k$ , which is  $\min(m, q + 1)$ . It then decreases  $k$  until  $P_k \sqsupseteq P_q a$ , which must eventually occur, since  $P_0 = \varepsilon$  is a suffix of every string.

The running time of COMPUTE-TRANSITION-FUNCTION is  $O(m^3 |\Sigma|)$ , because the outer loops contribute a factor of  $m |\Sigma|$ , the inner **repeat** loop can run at most  $m + 1$  times, and the test  $P_k \sqsupseteq P_q a$  on line 7 can require comparing up

to  $m$  characters. Much faster procedures exist; by utilizing some cleverly computed information about the pattern  $P$  (see Exercise 32.4-8), we can improve the time required to compute  $\delta$  from  $P$  to  $O(m |\Sigma|)$ . With this improved procedure for computing  $\delta$ , we can find all occurrences of a length- $m$  pattern in a length- $n$  text over an alphabet  $\Sigma$  with  $O(m |\Sigma|)$  preprocessing time and  $\Theta(n)$  matching time.

### Exercises

#### 32.3-1

Construct the string-matching automaton for the pattern  $P = \text{aabab}$  and illustrate its operation on the text string  $T = \text{aaababaabaababaab}$ .

#### 32.3-2

Draw a state-transition diagram for a string-matching automaton for the pattern  $\text{ababbabbababbababbabb}$  over the alphabet  $\Sigma = \{a, b\}$ .

#### 32.3-3

We call a pattern  $P$  *nonoverlappable* if  $P_k \sqsubset P_q$  implies  $k = 0$  or  $k = q$ . Describe the state-transition diagram of the string-matching automaton for a nonoverlappable pattern.

#### 32.3-4 ★

Given two patterns  $P$  and  $P'$ , describe how to construct a finite automaton that determines all occurrences of *either* pattern. Try to minimize the number of states in your automaton.

#### 32.3-5

Given a pattern  $P$  containing gap characters (see Exercise 32.1-4), show how to build a finite automaton that can find an occurrence of  $P$  in a text  $T$  in  $O(n)$  matching time, where  $n = |T|$ .

---

## ★ 32.4 The Knuth-Morris-Pratt algorithm

We now present a linear-time string-matching algorithm due to Knuth, Morris, and Pratt. This algorithm avoids computing the transition function  $\delta$  altogether, and its matching time is  $\Theta(n)$  using just an auxiliary function  $\pi$ , which we precompute from the pattern in time  $\Theta(m)$  and store in an array  $\pi[1..m]$ . The array  $\pi$  allows us to compute the transition function  $\delta$  efficiently (in an amortized sense) “on the fly” as needed. Loosely speaking, for any state  $q = 0, 1, \dots, m$  and any character