

### 32.1 The naive string-matching algorithm

The naive algorithm finds all valid shifts using a loop that checks the condition  $P[1..m] = T[s+1..s+m]$  for each of the  $n - m + 1$  possible values of  $s$ .

NAIVE-STRING-MATCHER( $T, P$ )

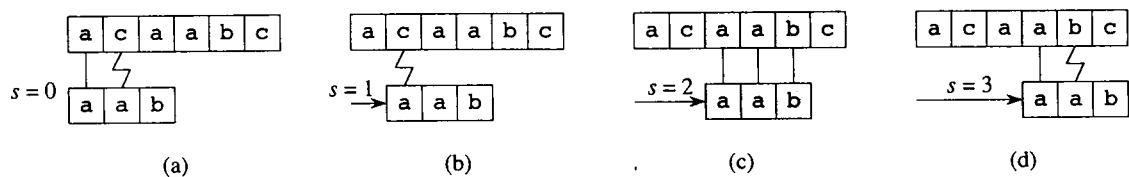
```

1   $n = T.length$ 
2   $m = P.length$ 
3  for  $s = 0$  to  $n - m$ 
4      if  $P[1..m] == T[s+1..s+m]$ 
5          print "Pattern occurs with shift"  $s$ 

```

Figure 32.4 portrays the naive string-matching procedure as sliding a “template” containing the pattern over the text, noting for which shifts all of the characters on the template equal the corresponding characters in the text. The **for** loop of lines 3–5 considers each possible shift explicitly. The test in line 4 determines whether the current shift is valid; this test implicitly loops to check corresponding character positions until all positions match successfully or a mismatch is found. Line 5 prints out each valid shift  $s$ .

Procedure NAIVE-STRING-MATCHER takes time  $O((n - m + 1)m)$ , and this bound is tight in the worst case. For example, consider the text string  $a^n$  (a string of  $n$  a’s) and the pattern  $a^m$ . For each of the  $n - m + 1$  possible values of the shift  $s$ , the implicit loop on line 4 to compare corresponding characters must execute  $m$  times to validate the shift. The worst-case running time is thus  $\Theta((n - m + 1)m)$ , which is  $\Theta(n^2)$  if  $m = \lfloor n/2 \rfloor$ . Because it requires no preprocessing, NAIVE-STRING-MATCHER’s running time equals its matching time.



**Figure 32.4** The operation of the naive string matcher for the pattern  $P = aab$  and the text  $T = acaabc$ . We can imagine the pattern  $P$  as a template that we slide next to the text. (a)–(d) The four successive alignments tried by the naive string matcher. In each part, vertical lines connect corresponding regions found to match (shown shaded), and a jagged line connects the first mismatched character found, if any. The algorithm finds one occurrence of the pattern, at shift  $s = 2$ , shown in part (c).

As we shall see, NAIVE-STRING-MATCHER is not an optimal procedure for this problem. Indeed, in this chapter we shall see that the Knuth-Morris-Pratt algorithm is much better in the worst case. The naive string-matcher is inefficient because it entirely ignores information gained about the text for one value of  $s$  when it considers other values of  $s$ . Such information can be quite valuable, however. For example, if  $P = \text{aaab}$  and we find that  $s = 0$  is valid, then none of the shifts 1, 2, or 3 are valid, since  $T[4] = \text{b}$ . In the following sections, we examine several ways to make effective use of this sort of information.

### Exercises

#### 32.1-1

Show the comparisons the naive string matcher makes for the pattern  $P = 0001$  in the text  $T = 000010001010001$ .

#### 32.1-2

Suppose that all characters in the pattern  $P$  are different. Show how to accelerate NAIVE-STRING-MATCHER to run in time  $O(n)$  on an  $n$ -character text  $T$ .

#### 32.1-3

Suppose that pattern  $P$  and text  $T$  are *randomly* chosen strings of length  $m$  and  $n$ , respectively, from the  $d$ -ary alphabet  $\Sigma_d = \{0, 1, \dots, d-1\}$ , where  $d \geq 2$ . Show that the *expected* number of character-to-character comparisons made by the implicit loop in line 4 of the naive algorithm is

$$(n - m + 1) \frac{1 - d^{-m}}{1 - d^{-1}} \leq 2(n - m + 1)$$

over all executions of this loop. (Assume that the naive algorithm stops comparing characters for a given shift once it finds a mismatch or matches the entire pattern.) Thus, for randomly chosen strings, the naive algorithm is quite efficient.

#### 32.1-4

Suppose we allow the pattern  $P$  to contain occurrences of a *gap character*  $\diamond$  that can match an *arbitrary* string of characters (even one of zero length). For example, the pattern  $\text{ab}\diamond\text{ba}\diamond\text{c}$  occurs in the text  $\text{cabccbacbacab}$  as

```
c ab cc ba cba c ab
  ab   $\diamond$  ba   $\diamond$  c
```

and as

```
c ab ccbac ba c ab.
  ab   $\diamond$  ba   $\diamond$  c
```

Note that the gap character may occur an arbitrary number of times in the pattern but not at all in the text. Give a polynomial-time algorithm to determine whether such a pattern  $P$  occurs in a given text  $T$ , and analyze the running time of your algorithm.

## 32.2 The Rabin-Karp algorithm

Rabin and Karp proposed a string-matching algorithm that performs well in practice and that also generalizes to other algorithms for related problems, such as two-dimensional pattern matching. The Rabin-Karp algorithm uses  $\Theta(m)$  preprocessing time, and its worst-case running time is  $\Theta((n - m + 1)m)$ . Based on certain assumptions, however, its average-case running time is better.

This algorithm makes use of elementary number-theoretic notions such as the equivalence of two numbers modulo a third number. You might want to refer to Section 31.1 for the relevant definitions.

For expository purposes, let us assume that  $\Sigma = \{0, 1, 2, \dots, 9\}$ , so that each character is a decimal digit. (In the general case, we can assume that each character is a digit in radix- $d$  notation, where  $d = |\Sigma|$ .) We can then view a string of  $k$  consecutive characters as representing a length- $k$  decimal number. The character string 31415 thus corresponds to the decimal number 31,415. Because we interpret the input characters as both graphical symbols and digits, we find it convenient in this section to denote them as we would digits, in our standard text font.

Given a pattern  $P[1..m]$ , let  $p$  denote its corresponding decimal value. In a similar manner, given a text  $T[1..n]$ , let  $t_s$  denote the decimal value of the length- $m$  substring  $T[s + 1..s + m]$ , for  $s = 0, 1, \dots, n - m$ . Certainly,  $t_s = p$  if and only if  $T[s + 1..s + m] = P[1..m]$ ; thus,  $s$  is a valid shift if and only if  $t_s = p$ . If we could compute  $p$  in time  $\Theta(m)$  and all the  $t_s$  values in a total of  $\Theta(n - m + 1)$  time,<sup>1</sup> then we could determine all valid shifts  $s$  in time  $\Theta(m) + \Theta(n - m + 1) = \Theta(n)$  by comparing  $p$  with each of the  $t_s$  values. (For the moment, let's not worry about the possibility that  $p$  and the  $t_s$  values might be very large numbers.)

We can compute  $p$  in time  $\Theta(m)$  using Horner's rule (see Section 30.1):

$$p = P[m] + 10(P[m - 1] + 10(P[m - 2] + \dots + 10(P[2] + 10P[1]) \dots)) .$$

Similarly, we can compute  $t_0$  from  $T[1..m]$  in time  $\Theta(m)$ .

<sup>1</sup>We write  $\Theta(n - m + 1)$  instead of  $\Theta(n - m)$  because  $s$  takes on  $n - m + 1$  different values. The "+1" is significant in an asymptotic sense because when  $m = n$ , computing the lone  $t_s$  value takes  $\Theta(1)$  time, not  $\Theta(0)$  time.

To compute the remaining values  $t_1, t_2, \dots, t_{n-m}$  in time  $\Theta(n - m)$ , we observe that we can compute  $t_{s+1}$  from  $t_s$  in constant time, since

$$t_{s+1} = 10(t_s - 10^{m-1}T[s + 1]) + T[s + m + 1]. \quad (32.1)$$

Subtracting  $10^{m-1}T[s + 1]$  removes the high-order digit from  $t_s$ , multiplying the result by 10 shifts the number left by one digit position, and adding  $T[s + m + 1]$  brings in the appropriate low-order digit. For example, if  $m = 5$  and  $t_s = 31415$ , then we wish to remove the high-order digit  $T[s + 1] = 3$  and bring in the new low-order digit (suppose it is  $T[s + 5 + 1] = 2$ ) to obtain

$$\begin{aligned} t_{s+1} &= 10(31415 - 10000 \cdot 3) + 2 \\ &= 14152. \end{aligned}$$

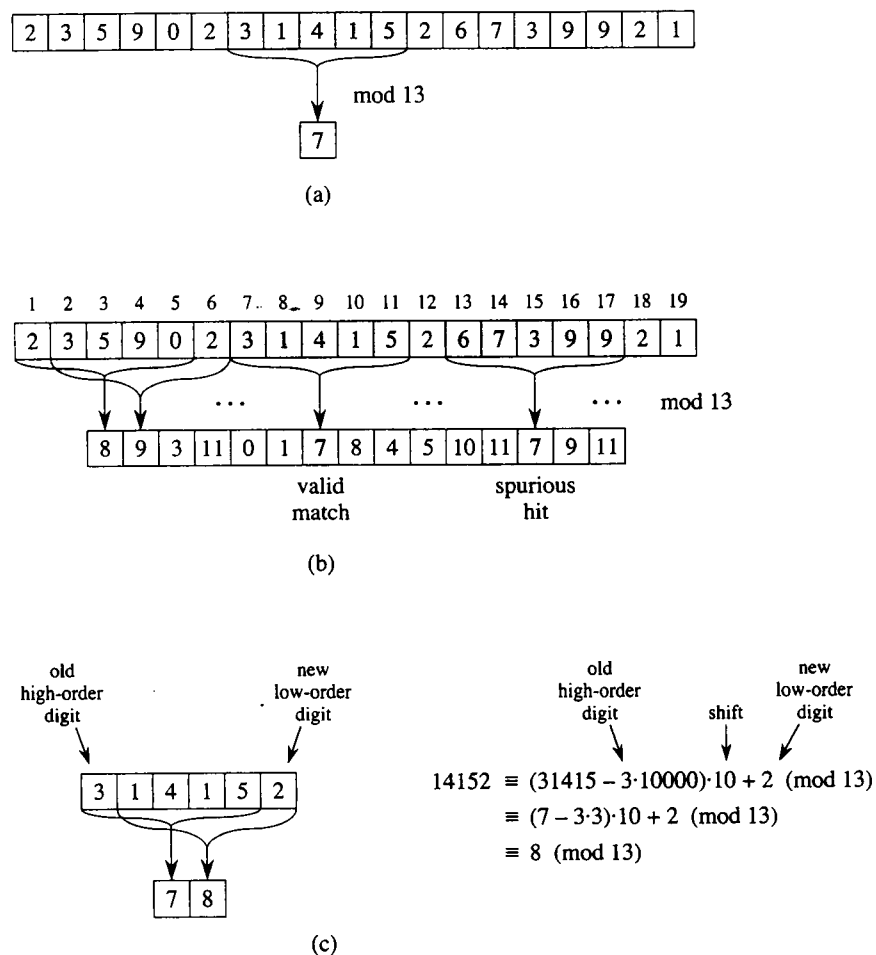
If we precompute the constant  $10^{m-1}$  (which we can do in time  $O(\lg m)$  using the techniques of Section 31.6, although for this application a straightforward  $O(m)$ -time method suffices), then each execution of equation (32.1) takes a constant number of arithmetic operations. Thus, we can compute  $p$  in time  $\Theta(m)$ , and we can compute all of  $t_0, t_1, \dots, t_{n-m}$  in time  $\Theta(n - m + 1)$ . Therefore, we can find all occurrences of the pattern  $P[1..m]$  in the text  $T[1..n]$  with  $\Theta(m)$  preprocessing time and  $\Theta(n - m + 1)$  matching time.

Until now, we have intentionally overlooked one problem:  $p$  and  $t_s$  may be too large to work with conveniently. If  $P$  contains  $m$  characters, then we cannot reasonably assume that each arithmetic operation on  $p$  (which is  $m$  digits long) takes “constant time.” Fortunately, we can solve this problem easily, as Figure 32.5 shows: compute  $p$  and the  $t_s$  values modulo a suitable modulus  $q$ . We can compute  $p$  modulo  $q$  in  $\Theta(m)$  time and all the  $t_s$  values modulo  $q$  in  $\Theta(n - m + 1)$  time. If we choose the modulus  $q$  as a prime such that  $10q$  just fits within one computer word, then we can perform all the necessary computations with single-precision arithmetic. In general, with a  $d$ -ary alphabet  $\{0, 1, \dots, d - 1\}$ , we choose  $q$  so that  $dq$  fits within a computer word and adjust the recurrence equation (32.1) to work modulo  $q$ , so that it becomes

$$t_{s+1} = (d(t_s - T[s + 1]h) + T[s + m + 1]) \bmod q, \quad (32.2)$$

where  $h \equiv d^{m-1} \pmod{q}$  is the value of the digit “1” in the high-order position of an  $m$ -digit text window.

The solution of working modulo  $q$  is not perfect, however:  $t_s \equiv p \pmod{q}$  does not imply that  $t_s = p$ . On the other hand, if  $t_s \not\equiv p \pmod{q}$ , then we definitely have that  $t_s \neq p$ , so that shift  $s$  is invalid. We can thus use the test  $t_s \equiv p \pmod{q}$  as a fast heuristic test to rule out invalid shifts  $s$ . Any shift  $s$  for which  $t_s \equiv p \pmod{q}$  must be tested further to see whether  $s$  is really valid or we just have a *spurious hit*. This additional test explicitly checks the condition



**Figure 32.5** The Rabin-Karp algorithm. Each character is a decimal digit, and we compute values modulo 13. (a) A text string. A window of length 5 is shown shaded. The numerical value of the shaded number, computed modulo 13, yields the value 7. (b) The same text string with values computed modulo 13 for each possible position of a length-5 window. Assuming the pattern  $P = 31415$ , we look for windows whose value modulo 13 is 7, since  $31415 \equiv 7 \pmod{13}$ . The algorithm finds two such windows, shown shaded in the figure. The first, beginning at text position 7, is indeed an occurrence of the pattern, while the second, beginning at text position 13, is a spurious hit. (c) How to compute the value for a window in constant time, given the value for the previous window. The first window has value 31415. Dropping the high-order digit 3, shifting left (multiplying by 10), and then adding in the low-order digit 2 gives us the new value 14152. Because all computations are performed modulo 13, the value for the first window is 7, and the value for the new window is 8.

$P[1..m] = T[s+1..s+m]$ . If  $q$  is large enough, then we hope that spurious hits occur infrequently enough that the cost of the extra checking is low.

The following procedure makes these ideas precise. The inputs to the procedure are the text  $T$ , the pattern  $P$ , the radix  $d$  to use (which is typically taken to be  $|\Sigma|$ ), and the prime  $q$  to use.

RABIN-KARP-MATCHER( $T, P, d, q$ )

```

1   $n = T.length$ 
2   $m = P.length$ 
3   $h = d^{m-1} \bmod q$ 
4   $p = 0$ 
5   $t_0 = 0$ 
6  for  $i = 1$  to  $m$  // preprocessing
7       $p = (dp + P[i]) \bmod q$ 
8       $t_0 = (dt_0 + T[i]) \bmod q$ 
9  for  $s = 0$  to  $n - m$  // matching
10     if  $p == t_s$ 
11         if  $P[1..m] == T[s+1..s+m]$ 
12             print "Pattern occurs with shift"  $s$ 
13     if  $s < n - m$ 
14          $t_{s+1} = (d(t_s - T[s+1])h) + T[s+m+1] \bmod q$ 
```

The procedure RABIN-KARP-MATCHER works as follows. All characters are interpreted as radix- $d$  digits. The subscripts on  $t$  are provided only for clarity; the program works correctly if all the subscripts are dropped. Line 3 initializes  $h$  to the value of the high-order digit position of an  $m$ -digit window. Lines 4–8 compute  $p$  as the value of  $P[1..m] \bmod q$  and  $t_0$  as the value of  $T[1..m] \bmod q$ . The for loop of lines 9–14 iterates through all possible shifts  $s$ , maintaining the following invariant:

Whenever line 10 is executed,  $t_s = T[s+1..s+m] \bmod q$ .

If  $p = t_s$  in line 10 (a “hit”), then line 11 checks to see whether  $P[1..m] = T[s+1..s+m]$  in order to rule out the possibility of a spurious hit. Line 12 prints out any valid shifts that are found. If  $s < n - m$  (checked in line 13), then the for loop will execute at least one more time, and so line 14 first executes to ensure that the loop invariant holds when we get back to line 10. Line 14 computes the value of  $t_{s+1} \bmod q$  from the value of  $t_s \bmod q$  in constant time using equation (32.2) directly.

RABIN-KARP-MATCHER takes  $\Theta(m)$  preprocessing time, and its matching time is  $\Theta((n - m + 1)m)$  in the worst case, since (like the naive string-matching algorithm) the Rabin-Karp algorithm explicitly verifies every valid shift. If  $P = a^m$

and  $T = a^n$ , then verifying takes time  $\Theta((n-m+1)m)$ , since each of the  $n-m+1$  possible shifts is valid.

In many applications, we expect few valid shifts—perhaps some constant  $c$  of them. In such applications, the expected matching time of the algorithm is only  $O((n-m+1) + cm) = O(n+m)$ , plus the time required to process spurious hits. We can base a heuristic analysis on the assumption that reducing values modulo  $q$  acts like a random mapping from  $\Sigma^*$  to  $\mathbb{Z}_q$ . (See the discussion on the use of division for hashing in Section 11.3.1. It is difficult to formalize and prove such an assumption, although one viable approach is to assume that  $q$  is chosen randomly from integers of the appropriate size. We shall not pursue this formalization here.) We can then expect that the number of spurious hits is  $O(n/q)$ , since we can estimate the chance that an arbitrary  $t_s$  will be equivalent to  $p$ , modulo  $q$ , as  $1/q$ . Since there are  $O(n)$  positions at which the test of line 10 fails and we spend  $O(m)$  time for each hit, the expected matching time taken by the Rabin-Karp algorithm is

$$O(n) + O(m(v + n/q)),$$

where  $v$  is the number of valid shifts. This running time is  $O(n)$  if  $v = O(1)$  and we choose  $q \geq m$ . That is, if the expected number of valid shifts is small ( $O(1)$ ) and we choose the prime  $q$  to be larger than the length of the pattern, then we can expect the Rabin-Karp procedure to use only  $O(n+m)$  matching time. Since  $m \leq n$ , this expected matching time is  $O(n)$ .

### Exercises

#### 32.2-1

Working modulo  $q = 11$ , how many spurious hits does the Rabin-Karp matcher encounter in the text  $T = 3141592653589793$  when looking for the pattern  $P = 26$ ?

#### 32.2-2

How would you extend the Rabin-Karp method to the problem of searching a text string for an occurrence of any one of a given set of  $k$  patterns? Start by assuming that all  $k$  patterns have the same length. Then generalize your solution to allow the patterns to have different lengths.

#### 32.2-3

Show how to extend the Rabin-Karp method to handle the problem of looking for a given  $m \times m$  pattern in an  $n \times n$  array of characters. (The pattern may be shifted vertically and horizontally, but it may not be rotated.)