

String Matching Algorithms

Naive String-Matching 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
- takes time $O((n-m+1)*m)$, which is $\theta(n^2)$ if $m = \lfloor n/2 \rfloor$

Naive String-Matching Algorithm

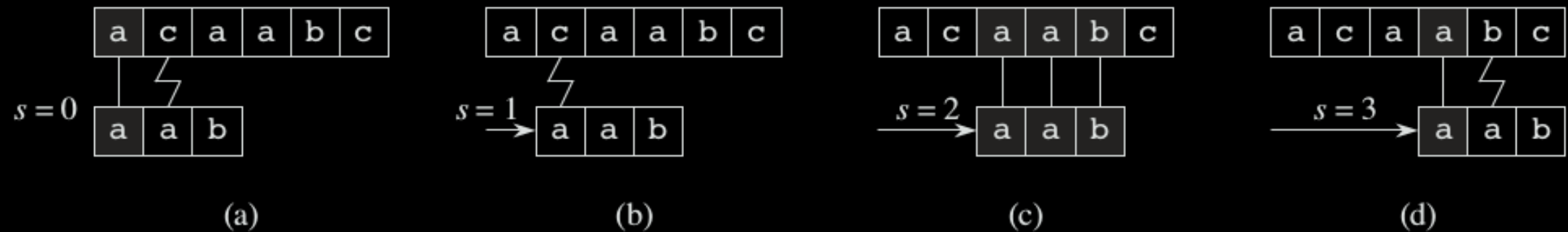


Figure 32.4 The operation of the naive string matcher for the pattern $P = \text{aab}$ and the text $T = \text{acaabc}$. We can imagine the pattern P as a template that we slide next to the text. (a)–(d) The four successive alignments tried by the naive string matcher. In each part, vertical lines connect 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).

Naive String-Matching Algorithm

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”  $s$ 
```

Naive String-Matching Algorithm

1

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

2

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

3

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

4

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

5

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

6

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

7

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

8

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

9

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

10

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1



Naive String-Matching Algorithm

11

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1


1	2	3	4
0	0	0	1

Naive String-Matching Algorithm

12

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

13

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

14

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

15

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

16

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1

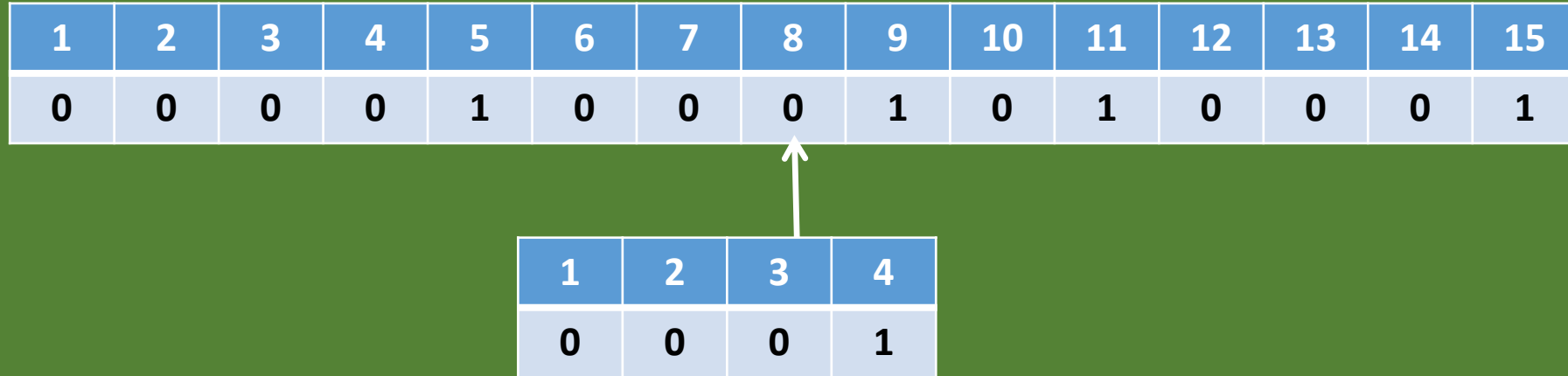


Naive String-Matching Algorithm

17

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

18

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

19

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

20

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

21

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

22

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

23

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

24

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

25

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

26

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

27

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

28

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

29

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

30

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1




Naive String-Matching Algorithm

31

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	1	0	0	0	1	0	1	0	0	0	1

1	2	3	4
0	0	0	1



Naive String-Matching Algorithm

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”  $s$ 
```

Rabin-Karp algorithm

- Performs well in practice and that also generalizes to other algorithms for related problems, such as two-dimensional pattern matching
- uses $\theta(m)$ preprocessing time, and its worst-case running time is $\theta((n-m+1)*m)$

Rabin-Karp algorithm

- Makes use of elementary number-theoretic notions such as the equivalence of two numbers modulo a third number
- Assume that $\Sigma = \{0, 1, 2, \dots, 9\}$, so that each character is a decimal digit
- Can then view a string of k consecutive characters as

Rabin-Karp algorithm

- Character string 31415 thus corresponds to the decimal number 31,415
- 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 decimal value of the length- m substring $T [s+1..s + m]$

Rabin-Karp algorithm

- Can compute p in time $\theta(m)$ using Horner's rule:

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

- Similarly compute t_0 from $T[1..m]$ in time $\theta(m)$
- t_1 is computed from $[2..m+1]$ and so on

Rabin-Karp algorithm

- Observe that we can compute t_{s+1} from t_s in constant

time

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

- 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
- $t_{s+1} = 10(31415 - 10000 * 3) + 2 = 14152$

Rabin-Karp algorithm

- 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
- 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.”

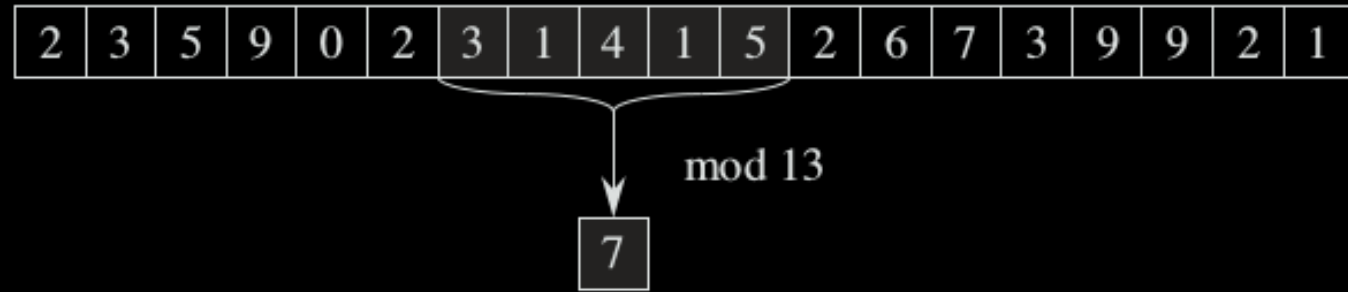
Rabin-Karp algorithm

- Fortunately, we can solve this problem easily, as Figure 32.5

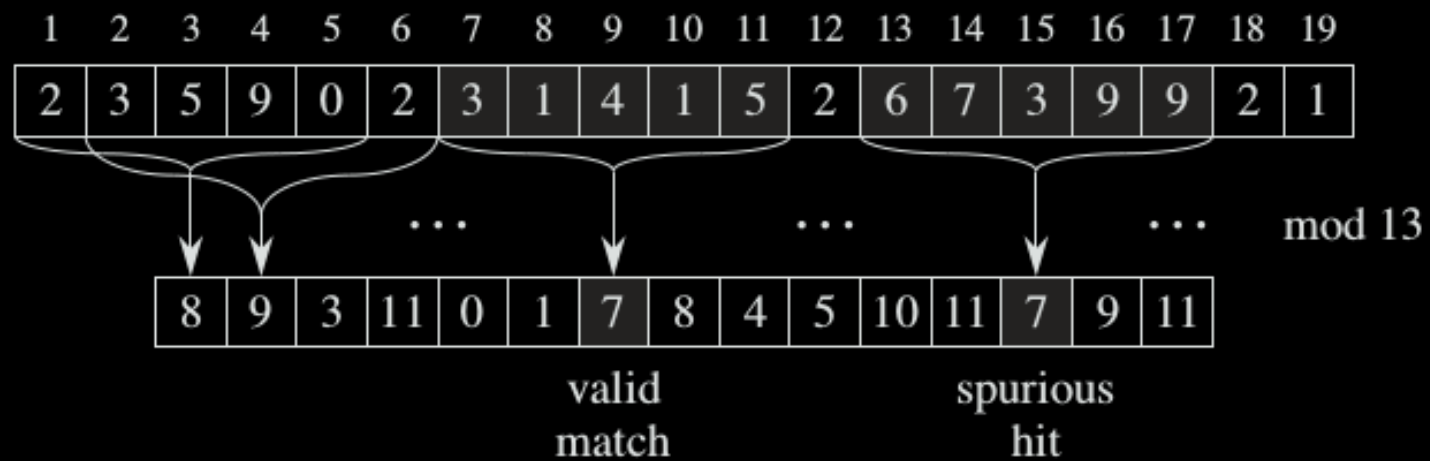
shows: compute p and the t_s values modulo a suitable modulus q

- Choose 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

Rabin-Karp algorithm



(a)



(b)

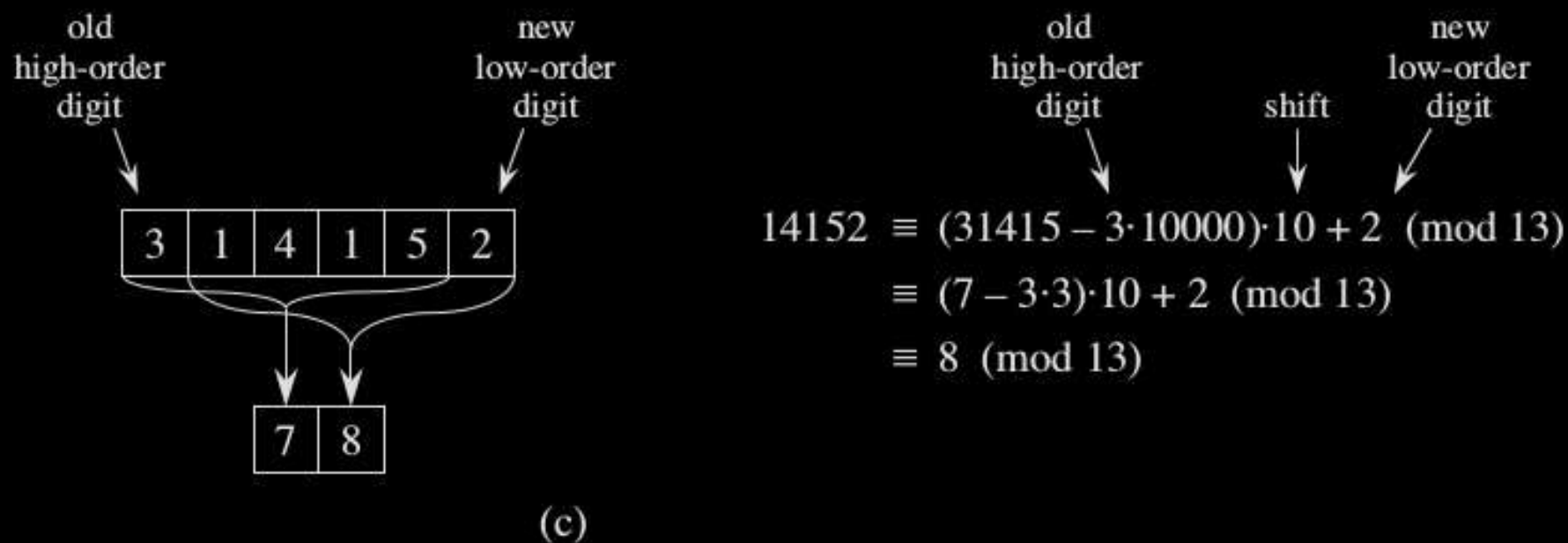


Figure 32.5 The **Rabin-Karp** algorithm. Each character is a decimal digit, and we compute values modulo 13. (a) A text string. A window of length 5 is shown shaded. The numerical value of the shaded number, computed modulo 13, yields the value 7. (b) The same text string with values computed modulo 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.

Rabin-Karp algorithm

- If q is large enough, then we hope that spurious hits occur infrequently enough that the cost of the extra checking is low.
- 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          $t_{s+1} = (d(t_s - T[s + 1]h) + T[s + m + 1]) \bmod q$ 
```