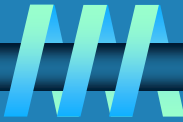


String Matching

Brute-Force String Matching



- pattern: a string of m characters to search for
- text: a (longer) string of n characters to search in
- problem: find a substring in the text that matches the pattern

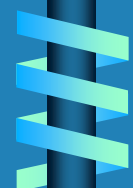
Brute-force algorithm

Step 1 Align pattern at beginning of text

Step 2 Moving from left to right, compare each character of pattern to the corresponding character in text until

- all characters are found to match (successful search); or
- a mismatch is detected

Step 3 While pattern is not found and the text is not yet exhausted, realign pattern one position to the right and repeat Step 2



Examples of Brute-Force String Matching

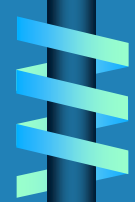


1. **Pattern:** 001011

Text: 10010101101001100101111010

2. **Pattern:** happy

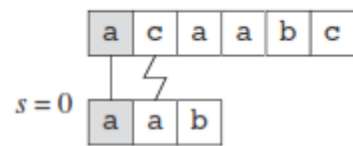
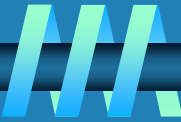
Text: It is never too late to have a happy
childhood.



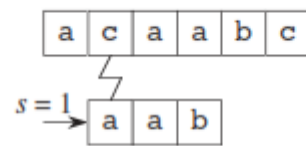
Pseudocode and Efficiency

ALGORITHM *BruteForceStringMatch*($T[0..n - 1]$, $P[0..m - 1]$)
//Implements brute-force string matching
//Input: An array $T[0..n - 1]$ of n characters representing a text and
// an array $P[0..m - 1]$ of m characters representing a pattern
//Output: The index of the first character in the text that starts a
// matching substring or -1 if the search is unsuccessful
for $i \leftarrow 0$ **to** $n - m$ **do**
 $j \leftarrow 0$
 while $j < m$ **and** $P[j] = T[i + j]$ **do**
 $j \leftarrow j + 1$
 if $j = m$ **return** i
return -1

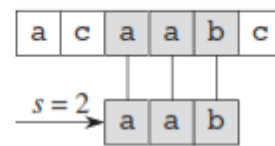
Efficiency:



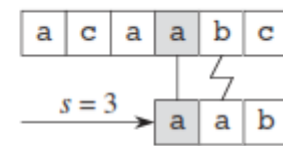
(a)



(b)



(c)



(d)

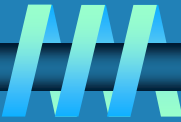


NONOTTONOTTONOTTONOTNONOT

NONOTTONOT

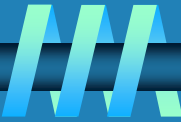
Figure 3.3 shows the first step of the brute-force algorithm. The pattern is aligned under the text, and the characters being compared are in bold. The alignment shows the pattern starting at index 0 of the text, with the first 'NON' matching, but the second 'NON' not matching because the text has 'OT' instead of 'OTTON'.

FIGURE 3.3 Example of brute-force string matching. (The pattern's characters that are compared with their text counterparts are in bold type.)



Shift amount	Text and pattern
0	GTAACAGTAAACG AAC
1	GTAACAGTAAACG AAC
2	GTAACAGTAAACG AAC
3	GTAACAGTAAACG AAC
4	GTAACAGTAAACG AAC
5	GTAACAGTAAACG AAC

Shift amount	Text and pattern
6	GTAACAGTAAACG AAC
7	GTAACAGTAAACG AAC
8	GTAACAGTAAACG AAC
9	GTAACAGTAAACG AAC
10	GTAACAGTAAACG AAC

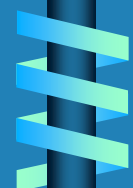


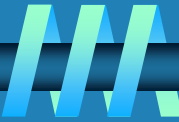
Determine the number of character comparisons made by the brute-force algorithm in searching for the pattern GANDHI in the text

THERE_IS_MORE_TO_LIFE_THAN_INCREASING_ITS_SPEED

Assume that the length of the text—it is 47 characters long—is known before the search starts.

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

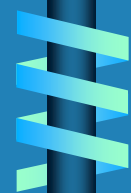




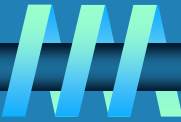
How many comparisons (both successful and unsuccessful) will be made by the brute-force algorithm in searching for each of the following patterns in the binary text of one thousand zeros?

- a.** 00001 **b.** 10000 **c.** 01010

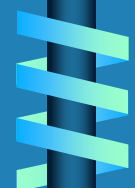
Give an example of a text of length n and a pattern of length m that constitutes a worst-case input for the brute-force string-matching algorithm. Exactly how many character comparisons will be made for such input?

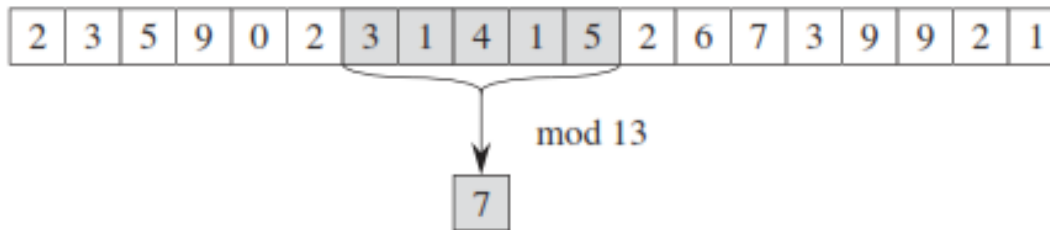
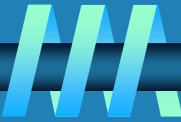


Rabin Karp string matching

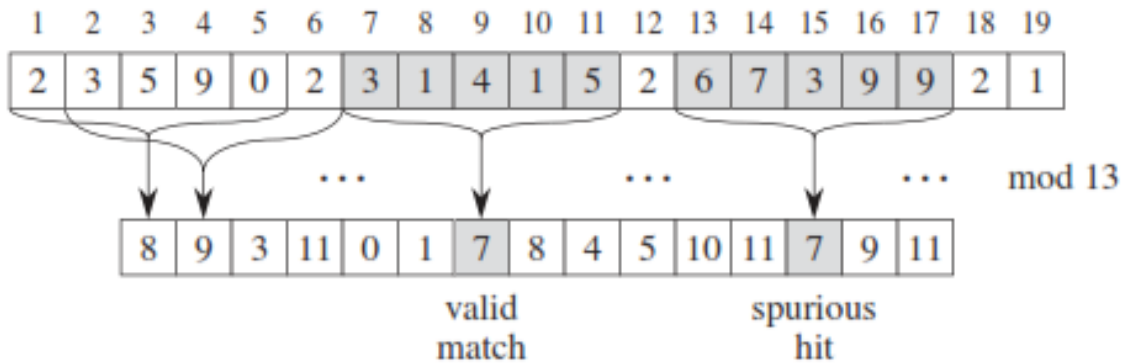


Given a pattern $P[1..m]$, let p denote its corresponding decimal value. In a similar manner, given a text $T[1..n]$, let t_s denote the decimal value of the length- m substring $T[s+1..s+m]$, for $s = 0, 1, \dots, n-m$. Certainly, $t_s = p$ if and only if $T[s+1..s+m] = P[1..m]$; thus, s is a valid 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,¹

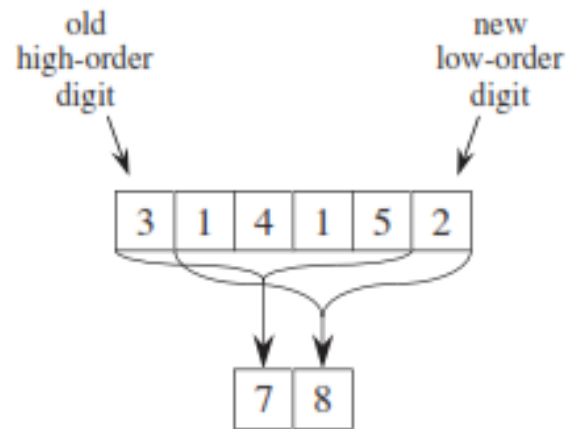
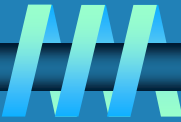




(a)



(b)



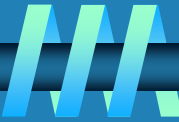
(c)

old high-order digit

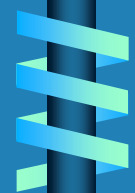
shift

new low-order digit

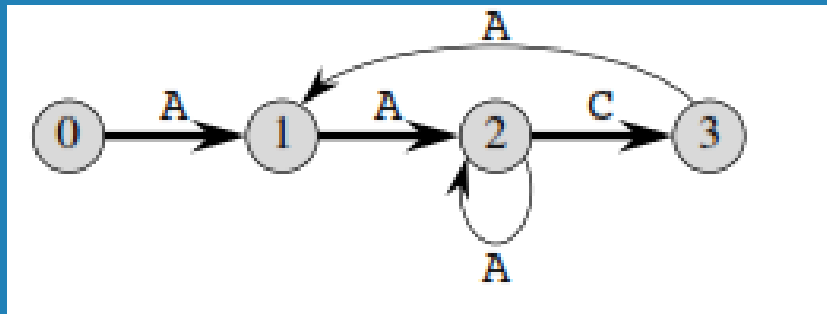
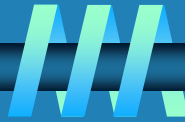
$$\begin{aligned} 14152 &\equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13} \\ &\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13} \\ &\equiv 8 \pmod{13} \end{aligned}$$



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

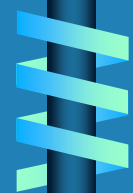


Finite Automata

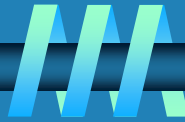


the pattern AAC does on the input text
GTAACAGTAAACG

state	0	0	0	1	2	3	1	0	0	1	2	2	3	0
character	G	T	A	A	C	A	G	T	A	A	A	C	G	



Knuth-Morris-Pratt algorithm



❧ Linear Time Algorithm

❧ txt = "AAAAABAAABA"

❧ pat = "AAAA"

❧ We compare first window of txt with pat

❧ txt = "AAAAABAAABA"

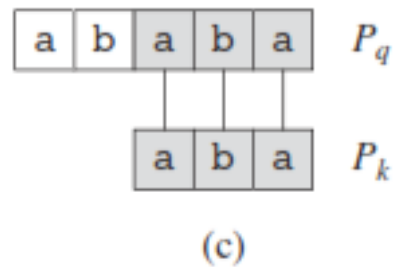
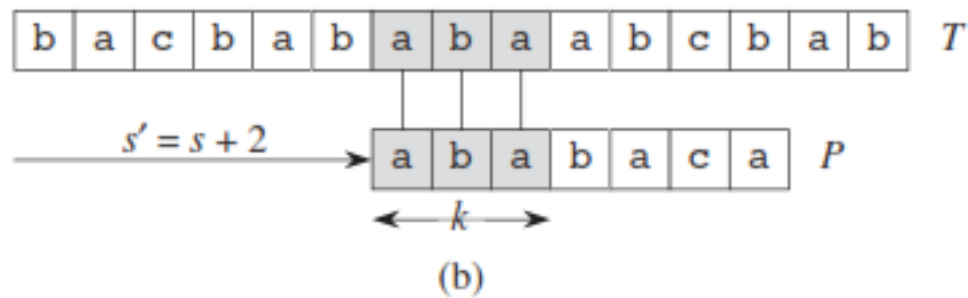
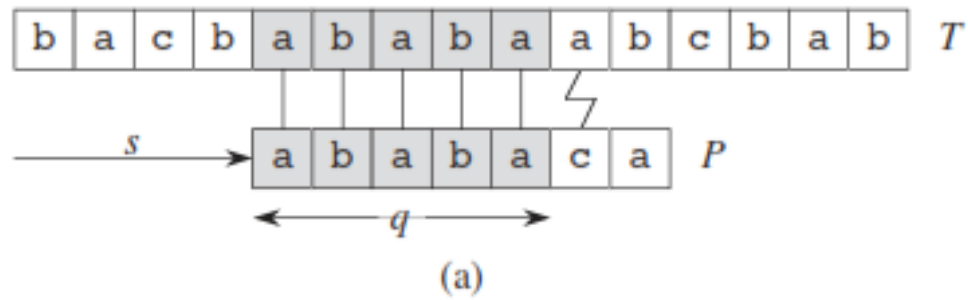
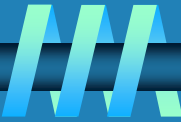
❧ pat = "AAAA" [Initial position]

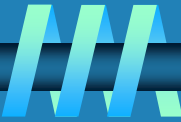
❧ We find a match. This is same as Naive String Matching.

❧ In the next step, we compare next window of txt with pat.

❧ txt = "AAAAABAAABA"

❧ pat = "AAAA" [Pattern shifted one position]





QUICKSORT(A, p, r)

```
1  if  $p < r$ 
2       $q = \text{PARTITION}(A, p, r)$ 
3      QUICKSORT( $A, p, q - 1$ )
4      QUICKSORT( $A, q + 1, r$ )
```

