

STRING MATCHING ALGORITHMS

Tanjina Helaly

STRING MATCHING

- Also known as
 - Substring matching
 - Pattern matching
- Find pattern of length M in a text of length N .

Diagram illustrating the KMP algorithm. The **Text** string is "A A C A B A B A A B A B A A B". The **Pattern** string is "A B A A A B A A".



APPLICATIONS

- Particular patterns in DNA sequences.
- Text editors
- Search engines
 - web crawling: finding strings inside other strings
- Spam detection
 - Look for pattern e.g. profit & bankaccount, lose weight
- Screen scraping
- Plagiarism Detection



FORMAL DEFINITION

- We formalize the string-matching problem as follows.
 - We assume that
 - **Text** $T[1\dots n]$ \rightarrow an array of length n .
 - **Pattern** $P[1\dots m]$ \rightarrow an array of length $m \leq n$.
 - The elements of P and T are characters drawn from a finite alphabet Σ .
 - For example, we may have $\Sigma = \{0,1, \dots\}$ or $\Sigma = \{a, b, \dots\}$ or $\Sigma = \{\text{ASCII values}\}$
 - Given strings T (text) and P (pattern), the pattern matching problem consists of finding a substring of T equal to P .



WHAT IS SHIFT

- The pattern P *occurs with shift* s in text T
 - if $0 \leq s \leq n-m$ and $T[s+1 \dots s+m] = P[1 \dots m]$
 - (that is, if $T[s+j] = P[j]$, for $1 \leq j \leq m$).
- If P occurs with shift s in T , then we call s a *valid shift*;
- The *string-matching problem* is the problem of finding **all valid shifts** with which a given pattern P occurs in a given text T .

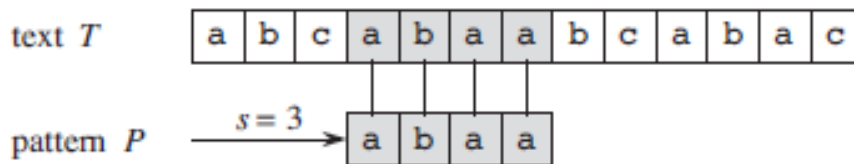


Figure 32.1 An example of the string-matching problem, where we want to find all occurrences of the pattern $P = \text{abaa}$ in the text $T = \text{abcabaabcabac}$. The pattern occurs only once in the text, at shift $s = 3$, which we call a valid shift. A vertical line connects each character of the pattern to its matching character in the text, and all matched characters are shaded.

PATTERN MATCHING ALGORITHMS

- Brute-force algorithm
- Rabin-Karp algorithm
- Boyer-Moore algorithm
- Knuth-Morris-Pratt algorithm

Algorithm	Preprocessing time	Matching time
Naive	0	$O((n - m + 1)m)$
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$
Finite automaton	$O(m \Sigma)$	$\Theta(n)$
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(n)$

Figure 32.2 The string-matching algorithms in this chapter and their preprocessing and matching times.



THE NAIVE STRING-MATCHING ALGORITHM

- Brute force
- finds all valid shifts using a loop that checks the condition $P[1\dots m] = T[s+1 \dots s+m]$ 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  $P[1\dots m] == T[s+1 \dots s+m]$   
5     print "Pattern occurs with shift"  $s$ 
```



THE NAIVE STRING-MATCHING ALGORITHM

- COMPLEXITY

- Complexity $\rightarrow O((n-m+1)m)$
 - 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.
 - Why need to compare all m even if there is a mismatch before m ?
 - Has some room to improve
 - Worst case. See the example below.
 - $T = \text{aaaaaaaa}$
 - $P = \text{aaaa}$



SCOPE OF IMPROVEMENT

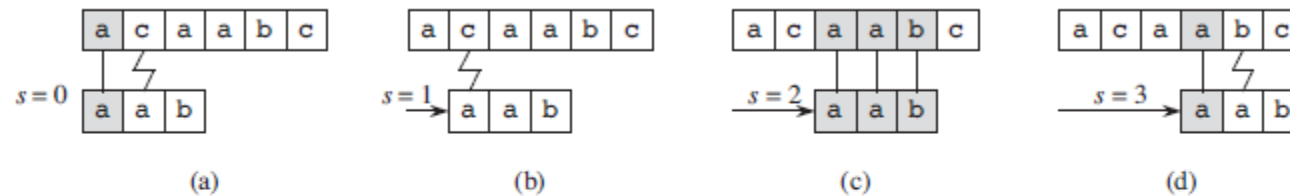


Figure 32.4 The operation of the naive string matcher for the pattern $P = \text{aab}$ and the text $T = \text{ac aabc}$. We can imagine the pattern P as a template that we slide next to the text. (a)–(d) The four successive alignments tried by the naive string matcher. In each part, vertical lines connect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



THE NAIVE STRING-MATCHING ALGORITHM - WEAKNESS

- 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 $T[4] = \text{b}$.



THE RABIN-KARP ALGORITHM



THE RABIN-KARP ALGORITHM

- Rabin and Karp proposed a string-matching algorithm that
 - performs well in practice and
 - 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 $\Theta((n-m+1)m)$.
 - however, its average-case running time is better.



THE RABIN-KARP ALGORITHM

- Assume we have to find for a match of the following patten in the text below.
 - Text – BALLTHEBALL length $\rightarrow N = 11$
 - Pattern – BALL length $\rightarrow M = 4$ ($M \leq N$)
- Assume $\text{hash}(\text{BALL}) = x$
- Now
 - Take a window of size M in the Text from beginning
 - Calculate the hash
 - If the hash of the window and the pattern is same.
 - Compare each character of the window and the pattern.
 - Slide the window 1 position right.
 - Repeat until the window contains $<M$ character.



THE RABIN-KARP ALGORITHM

- The solution of working modulo q is not perfect, however:

$hash(T_s) \equiv hash(P)$ does not imply $Ts = p$

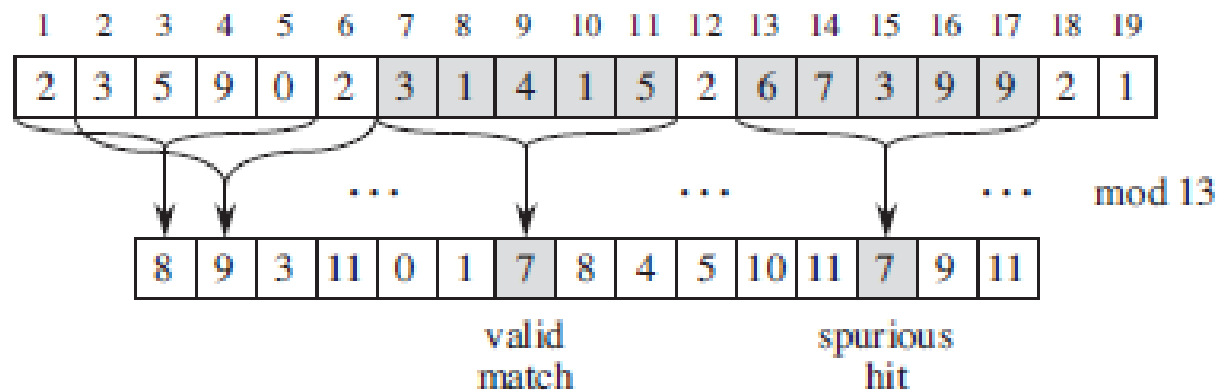
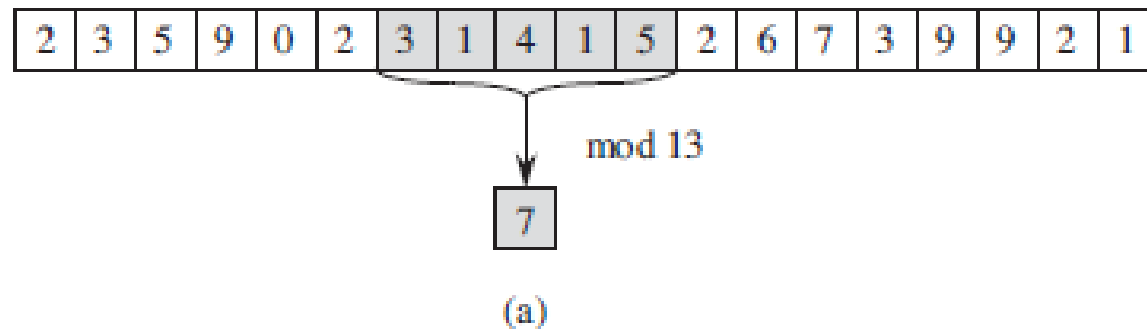
but

$hash(T_s) \not\equiv hash(P)$ imply $Ts \neq p$

- We can thus use the test $hash(T_s) \equiv hash(P)$ as a fast heuristic test to rule out invalid shifts s .
- Any shift s for which $hash(T_s) \equiv hash(P)$ must be tested further to see whether s is really valid or we just have a ***spurious hit***.



THE RABIN-KARP ALGORITHM



CAN WE IMPROVE?

- Use the hash value to current window calculate the hash value of next window.
 - Will reduce the time complexity of hash function from $\Theta(m)$ to $\Theta(1)$



PATTERN CALCULATION

- Pattern – “31415”
- Value = $3 \times 10^4 + 1 \times 10^3 + 4 \times 10^2 + 1 \times 10^1 + 5 \times 10^0$
- Or $((((3 \times 10 + 1) \times 10 + 4) \times 10 + 1) \times 10 + 5)$
- **Algorithm:** (for calculating the value of current window)
- $p=0$, $P="31415"$, $d=10$, $m=5$

for i from 1 to m
 $p = (d * p + P[i])$



HASH CALCULATION

- Some Mod properties
 - $(A + B) \bmod C = ((A \bmod C) + (B \bmod C)) \bmod C$
 - $(A \times B) \bmod C = ((A \bmod C) \times (B \bmod C)) \bmod C$
- So, Algorithm for **Mod(Hash)** calculation without any previous hash

for i from 1 to m
 $p = (d * p + P[i]) \bmod q$

- Complexity = $O(m)$



HASH CALCULATION – IMPROVED

2	3	5	9	0	2	3	1	4	1	5	2	6	7	3	9	9	2	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

- $T[s+1 .. s+m] = T_s = 31415$
- $T[s+2 .. s+m+1] = T_{s+1} = 10(31415 - 10000 \times 3) + 2 = 14152$
 $= d(T[s+1..s+m] - d^{m-1} \times T[s+1]) + T[s+m+1]$

Similar expression for hash calculation

- $t_{s+1} = (d(t_s - T[s+1]h) + T[s+m+1]) \bmod q$ [Complexity = O(1)]
where $h = d^{m-1} \bmod q$
 t_s = hash value of $T[s+1 .. s+m]$ (at shift s)
- So, we are using the hash value to current window (t_s)
calculate the hash value of next window (t_{s+1}).



IMPROVEMENT

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

○ Where

- t_s = hash value at shift s
- t_{s+1} = hash value at shift $s+1$ (next window)
- d = radix. (number of character in the alphabet system)
 - For numerical $d = 10$
 - For binary $d = 2$
 - For ASCII $d = 256$
- q = a prime number
- $h = d^{m-1} \bmod q$
- $T[s+1]$ = leading character of current window
- $T[s+m+1]$ = trailing character that will be brought in in the next window.



THE RABIN-KARP ALGORITHM

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]) mod 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 RABIN-KARP ALGORITHM-TIME COMPLEXITY

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

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



TIME COMPLEXITY

- Best Case (No matching) : $\Theta(n)$ [$m+n-m+1$]
- Worst Case : $\Theta(mn)$ [$m + (n-m+1)m$]
 - Text = “aaaaaaaaa”
 - Pattern = “aaaa”;
- Average Case: $\Theta(m+n)$ [realistic]



EXAMPLE

- $T = \text{BALLTHEBALL}$, $P = \text{BALL}$, $q = 29$ (prime number)
- Calculate hash, p
 - Repeat for each character and calculate $p = dp + P[i]$
 - $p = (dp + P[1]) \bmod q = 256 * 0 + 66 \bmod 29 = 8$
 - $p = (dp + P[2]) \bmod q = 256 * 8 + 65 \bmod 29 = 25$
 - $p = (dp + P[3]) \bmod q = 256 * 25 + 76 \bmod 29 = 9$
 - $p = (dp + P[4]) \bmod q = 256 * 9 + 76 \bmod 29 = 2$
- Calculate hash, t_0
 - Repeat for each character and calculate $t_0 = dt_0 + T[i]$
 - $t_0 = (dt_0 + T[1]) \bmod q = 256 * 0 + 66 \bmod 29 = 8$
 - $t_0 = (dt_0 + T[2]) \bmod q = 256 * 8 + 65 \bmod 29 = 25$
 - $t_0 = (dt_0 + T[3]) \bmod q = 256 * 25 + 76 \bmod 29 = 9$
 - $t_0 = (dt_0 + T[4]) \bmod q = 256 * 9 + 76 \bmod 29 = 2$



EXAMPLE

- $P = 2, t_0 = 2, m = 4, n = 11,$
- $q = 29, h = 256^3 \bmod 29 = 20$

Window (red text)	t_s	2 hashes equals	Text matches ?
B ALLTHEBALL	2	Yes	Yes
B A LLTHEBALL	$(256 * (2 - 66 * 20) + 84) \bmod 29 = 4$	No	NA
BALL T HEBALL	$(256 * (4 - 65 * 20) + 72) \bmod 29 = 27$	No	NA
BALLTHE B ALL	$(256 * (27 - 76 * 20) + 69) \bmod 29 = 23$	No	NA
BALLTHEB A LL	$(256 * (23 - 76 * 20) + 66) \bmod 29 = 11$	No	NA
BALLTHEB A LL	$(256 * (11 - 84 * 20) + 65) \bmod 29 = 0$	No	NA
BALLTHEB A LL	$(256 * (0 - 72 * 20) + 76) \bmod 29 = 26$	No	NA
BALLTHEB A LL	$(256 * (26 - 69 * 20) + 76) \bmod 29 = 2$	Yes	Yes



- Advantages:

- Can be extended to 2d pattern.
- Can be extended to find multiple patterns.

- Disadvantages:

- Arithmetic operation is slower than character comparison.

- Try at home:

- Extend Rabin-Karp to find any one of \mathbf{P} possible patterns in a text of size \mathbf{N} .



REFERENCES

- Chapter 32 (32.1, 32.2) - Cormen

