

- ① Create a Token Object. with loc no where it originated.
  - ② Internally each token is represented by a single character for the comparison.
- 

- ① Comparison of token strings
- ① Algorithm - Greedy String Tilting.
- ② Aim: Find set of substrings that are the same and satisfy following rules:-
  - ① Match only 1 token of A string with 1 token of B string.
  - ② Substrings to be found independent of their position in their string.
  - ③ Long substring matches are preferred over short ones.

# Robin Karp Algorithm for String Matching

Text  
Pattern

① Preprocessing text pattern    ② Matching

Let us assume  $\Sigma = \{0, 1, 2, 3, 4, \dots, 9\}$

→ each character is a decimal digit.

We can view a string of  $k$  digits as length  $k$  decimal number

string =  $\frac{31456}{\text{string}}$  = decimal number.  $\frac{31456}{\text{number}}$

1 2 3 4 5 7 2  
substring + Number.

a - - - z

A - - - z

Radix  $d = 52$  number.

Radix  $d$  no system.

[7898].

(2)

2246      7898

Pattern Matching

$p$  = decimal value.

Horner's Rule.



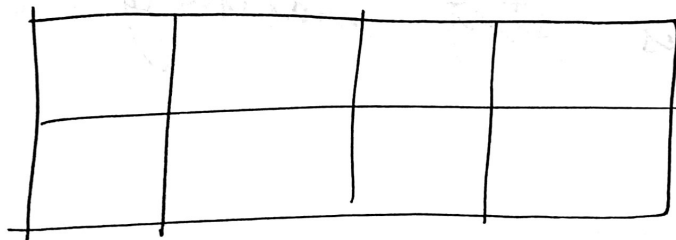
$t_s =$

$$3 \times 10^4 + 2 \times 10^3 + 7 \times 10^2 + 3 \times 10^1 + 2 \times 10^0$$

$$= [3^5 \times 10^4 + 2 \times 10^3 + 7 \times 10^2 + 3 \times 10^1 + 1 \times 10^0]$$

$$= 1 \cdot 0$$

Tiling Problem:



$$\text{count}(n) = \text{count}(n-1) + \text{count}(n-2)$$

$O(\log n)$  - Fibonacci sequence.

Text: "GEEKS FOR GEEKS"

length =  $N$

Pattern: "GEEK"

length =  $M$

index = 0, 8

Slide pattern - match hash value of pattern with hash value of current substring of text.

If hash value matches start matching individual characters

hashing converts text to numeric value.

Algorithm uses the fact that if two strings are equal - their hash values are also equal.

Assume length of text  $>$  length of pattern  $M$ .

hash("GEEK") =  $X$ .

window size =  $M$ .

Hash Function:

Requirement: Hash at the next shift must be efficiently computable ( $O(1)$ ) from current hash value and next character in text.

Rehashing:  $\text{hash}(\text{txt}[s+1 \dots s+m]) = d(\text{hash}(\text{txt}[s \dots s+m-1]) - \text{txt}[s] * h) + (\text{txt}[s+m]) \bmod q$

hash next = hash current \*  $d$  - trailing current char + leading current char

$d$  = No of characters in alphabet = 256

$$h = d^{(m-1)}$$

$q$  = prime number

#define d 256

// pat → pattern

// txt → text

// q → prime no.

void search (char pat[], char txt[], int q)

{  
int M = strlen (pat);

int N = strlen (txt);

int i, j;

int p = 0; // hash value for pattern

int t = 0; // hash value for text.

int h = 1;

for (i = 0; i < M - 1; i++)

h = (h \* d) % q;

for (i = 0; i < M; i++)

h = (d \* h + pat[i]) % q;

t = (d \* t + txt[i]) % q; }

text = "ABCD B"

pat = "DC"

$$q = 11$$

$$M = 2 \quad N = 5$$

$$h = (1 * 256) \% 11 = 3$$

Iter 1  $p = (256 * 0 + 68) \% 11 = 2$

$$t = (256 * 0 + 65) \% 11 = 10$$

Iter 2  $p = (256 * 2 + 67) \% 11 = 7$

$$t = \dots$$

Slide over-patterns on text.

Repeat iterations with +ve.