

UNIT-6

string matching

Naive string matching

- * string matching means finding all occurrences of a pattern in a given text [either number series or Alphabets]
- * Naive pattern searching is the simplest method among other pattern searching algorithm. It checks for all character of the main string to the pattern. If finding one of all exact occurrences of a pattern in a text.
- * This algorithm is helpful to smaller texts. We can find substring by checking once for the string.
- * The naive approach tests all the possible placement of pattern $p[1 \dots m]$ relative to Text $T[1 \dots n]$.

We try to shift ($s = 0, 1, 2, \dots, n-m$) successively and for each

shift s .

Compare $T[s+1, \dots, s+m]$ to $p[1 \dots m]$.

1. $n \leftarrow \text{length}[T]$

2. $m \leftarrow \text{length}[p]$

3. for $s \leftarrow 0$ to $n-m$

4. do If $p[1 \dots m] = T[s+1, \dots, s+m]$

5. then point "pattern occurs with shift"

Analysis:- for loop executes $n-m+1$ times; Time complexity = $O(n-m)$

Q1 Text $T = a c a a b c$

pattern $P = a a b$; the pattern is found at which index using naive string matching algorithm.

SQ

Text $T = \boxed{a} \boxed{c} \boxed{a} \boxed{a} \boxed{b} \boxed{c}$

pattern $P = \boxed{a} \boxed{a} \boxed{b}$

Step1 :- matching starting from 0 index, initially shift $s=0$

$\boxed{a} \boxed{c} \boxed{a} \boxed{a} \boxed{b} \boxed{c}$

$s=0$

$\boxed{a} \boxed{a} \boxed{b}$

string not matching, then shifting to next index.

Step2 :- pattern is shifting to next index. $s=1$

$\boxed{a} \boxed{c} \boxed{a} \boxed{a} \boxed{b} \boxed{c}$

$s=1$

$\boxed{a} \boxed{a} \boxed{b}$

not matching, then pattern is shifting to next index.

Step3 :- pattern is shifting to next index $s=2$

$\boxed{a} \boxed{c} \boxed{a} \boxed{a} \boxed{b} \boxed{c}$

$s=2$

$\boxed{a} \boxed{a} \boxed{b}$

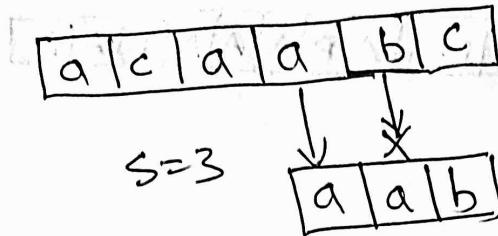
$T[2] \rightarrow a$

in the given pattern

1, 2, 3 character are matching with the given text

∴ pattern is found at index 2

Step 4:- Again pattern is shifting to next index

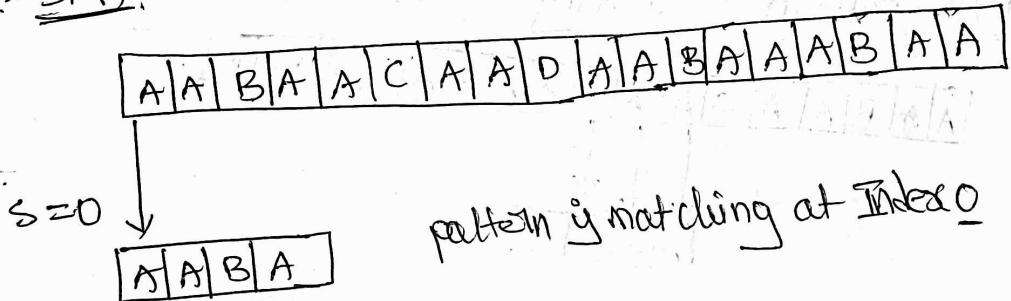


In a given pattern 1st character is matching, 2nd character is not matching with given Text

\therefore pattern aab is found at Index 2

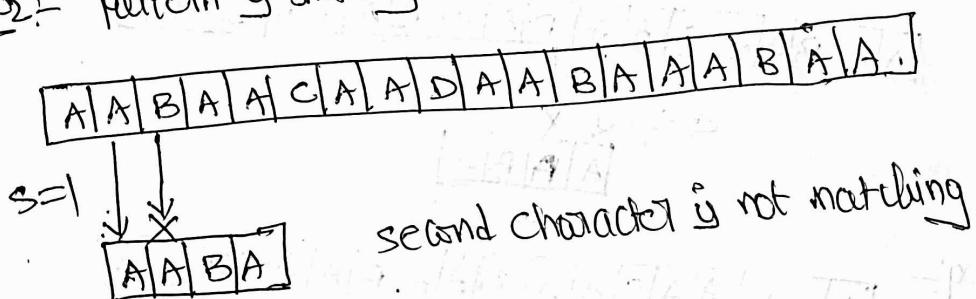
Sol 2:- Text $T = A A B A A C A A D A A B A A A B A A$
pattern $P = A A B A$

Sol :- Step 1:-



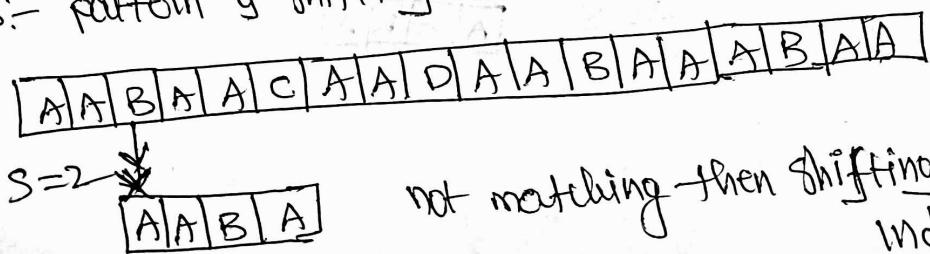
pattern is matching at Index 0

Step 2:- pattern is shifting to next index



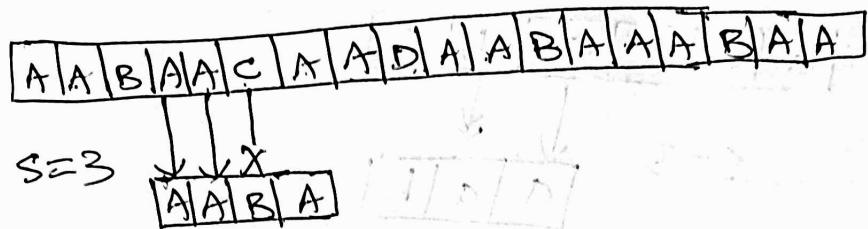
second character is not matching

Step 3:- pattern is shifting to next index

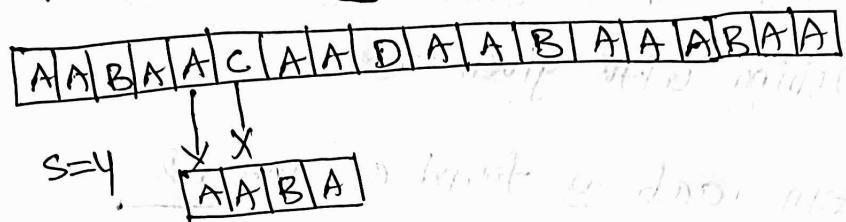


not matching then shifting to next index.

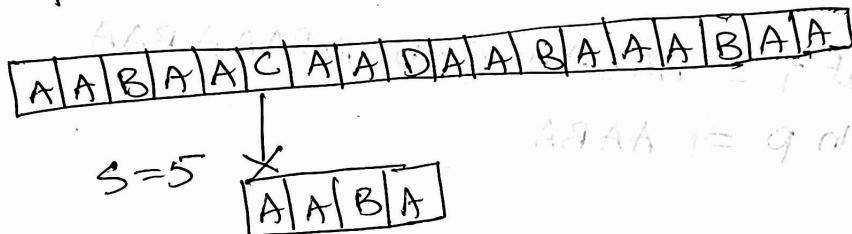
Step4:- pattern is shifting to next index.



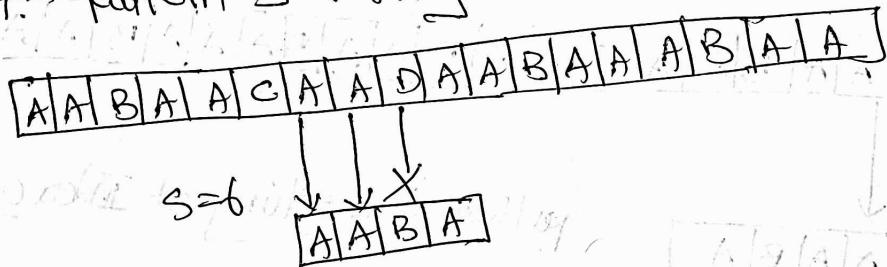
Step5:- pattern is shifting to next index



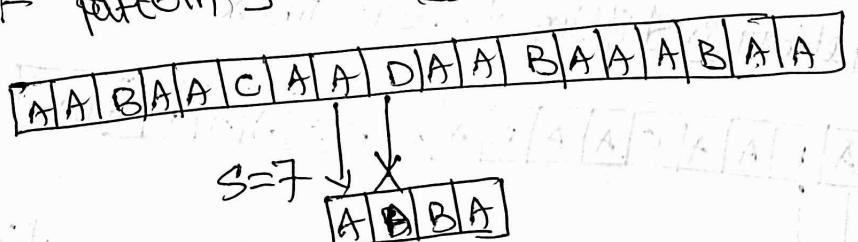
Step6:- pattern is shifting to next index



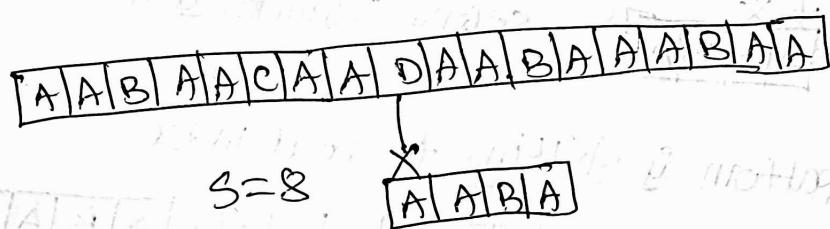
Step7:- pattern is shifting to next index



Step8:- pattern is shifting to next index



Step9:-



Step 0:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=9$

$\boxed{\text{A A B A}}$

pattern matching at
Index - 9

Step 1:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=10$

$\boxed{\text{A A B A}}$

pattern not matching

Step 2:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=11$

$\boxed{\text{A A B A}}$

pattern not matching

Step 3:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=12$

$\boxed{\text{A A B A}}$

pattern not matching

Step 4:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=13$

$\boxed{\text{A A B A}}$

pattern matching
at index 13

Step 5:- $\boxed{\text{A A B A A C A A D A A A B A A A B A A}}$

$S=14$

$\boxed{\text{A A B A}}$

pattern not
matching

∴ pattern matching at index 0

" "

" "

Index 9

" "

" "

Index 13

∴ Pattern matching is successful at index 0 & 13

Algorithm

```

void search (char *pat, char *txt)
{
    int M = strlen (pat);
    int N = strlen (txt);
    for (int i=0; i<=N-M; i++) // A loop to slide pattern one by one
    {
        int j;
        for (j=0; j<M; j++) // current index i, check for pattern
        match.
        if (txt[i+j] != pat[j])
            break;
        if (j == M) // If pat[0...M-1] = txt[i, i+1 ... i+M-1]
            printf ("pattern found at index %d", i);
    }
}

```

Best Case:- $\text{text} = \text{"AABCCAAADDEE"}$
 $\text{pattern} = \text{"FAA"}$

the first character of the pattern is not present in text at all.
 the pattern is just move the window, the time complexity is the

: Time Complexity = $O(n)$. text window size.

Worst Case:- All characters of the text and pattern are same

$\text{text} = \text{"AAAAA A A A A A A"}$

$\text{pattern} = \text{"AAAA"}$

In each window we need to compare all the characters until the last window.

: Time Complexity = $O(m \times (n - m + 1))$

(8)

If the last character is different even though we need to check all characters in the window.

$\text{Text} = \text{"A A A A A A A A A A B"}$
 $\text{pattern} = \text{"A A A B"}$

* Note:- pattern searching is an important problem in computer science. When we do search in a string in notepad/word file browser (Q) database (Q) google search, pattern searching algorithms are used to show the search results

Tries string matching

- * Trie is a tree that stores a set of strings. How we are going to store a string in a tree. A Trie is a tree-like information retrieval data structure whose nodes store the letters of an alphabet.
- * Every node except root node will store a letter in the alphabet. It is known as Radix Tree or Digital Tree or prefix tree

$$\Leftrightarrow S = \{ \text{bear}, \text{bell}, \text{bed}, \text{bul}, \text{stop}, \text{stock} \}$$

○ root

- * Trie is an efficient information retrieval data structure. It is also called Digital Tree and sometimes called Radix Tree or prefix tree.

- * Trie is pronounced as "try" it comes from the word "retrieval".

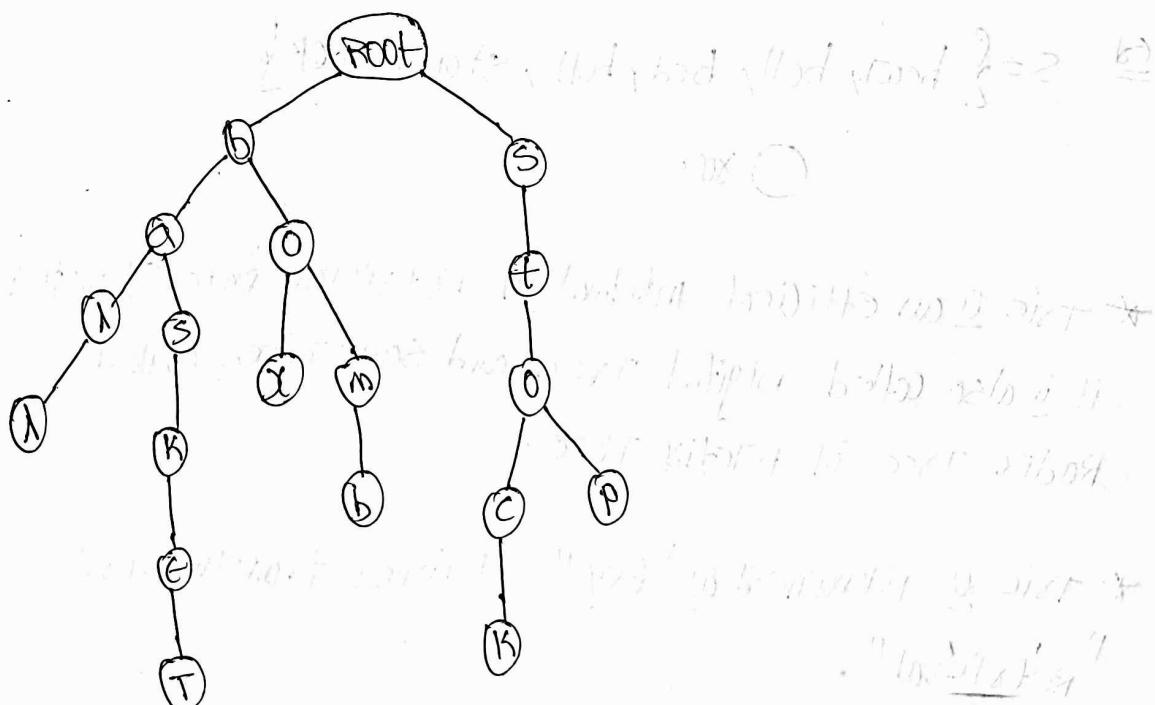
- * There are 3 types of tries

1. Standard Trie
2. Compressed Trie
3. Suffix Trie

standard Trie:-

- * it is an ordered Tree like Data Structure.
- * each node [except the Root node] in a standard trie is labeled with a character.
- * the children of a node are in Alphabetical Order.
- * Each node of branch represents a possible character or keys of words.
- * Each node of branch may have multiple character branches.

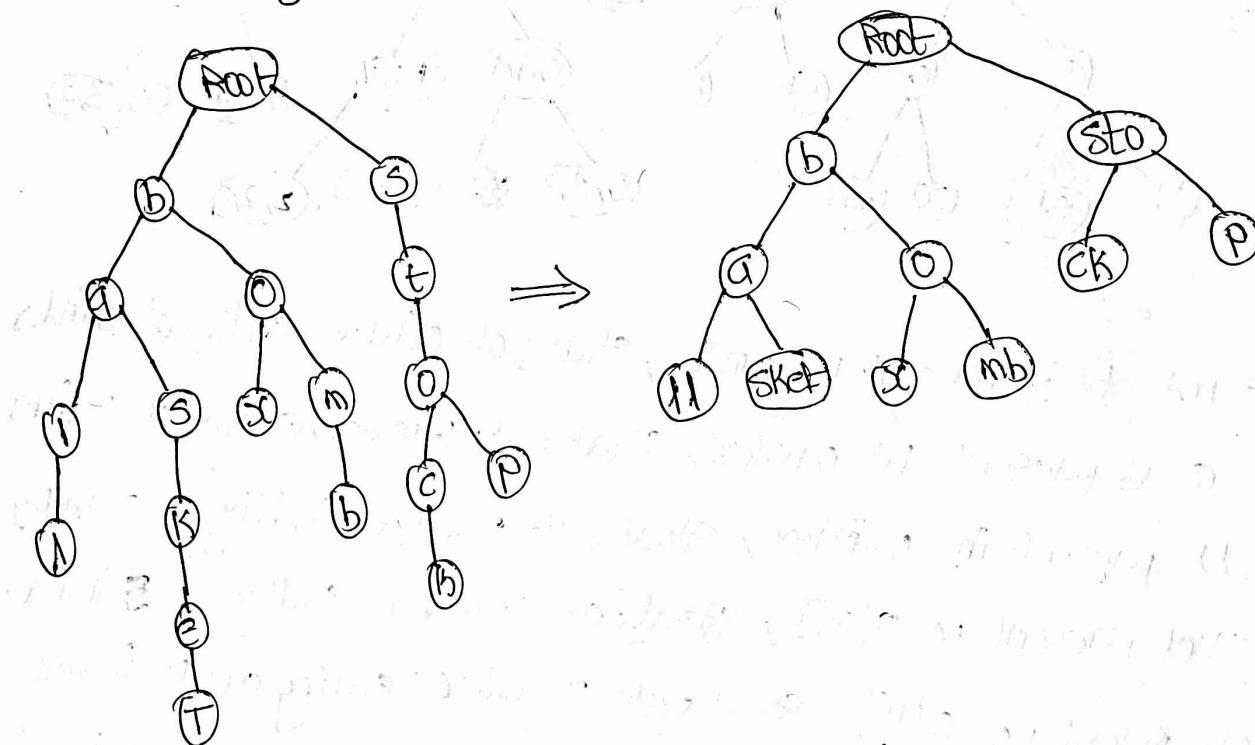
Ex string S = {ball, box, bomb, basket, stock, stop}



Compressed Trie

- * it's a advanced version of standard tries. If we have a node with one child, we try to merge them to get the easy form of the Trie.
- * each node [except the leaf nodes] have at least 2 children.
- * it is used to achieve space optimization.

Given string $S = \{ball, box, bomb, basket, stock, stop\}$



* Representation of a Compressed Trie

Index	0	1	2	3	4	5
-------	---	---	---	---	---	---

$S[0] = ball$

$S[1] = box$

$S[2] = bomb$

$S[3] = basket$

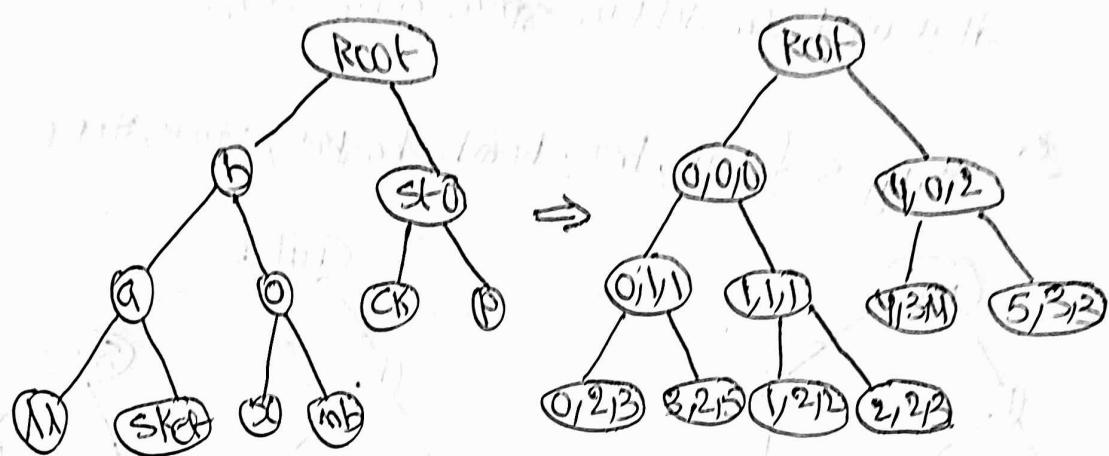
$S[4] = stock$

$S[5] = stop$

* we represent a compressed tree in three tuple relation

role(i, j, k)

- i represents index at which the word is present
- j represents the index at which the prefix starts
- k represents the index at which the prefix ends



- node b is present in 0 -index, starts at 0 -index, ending at 0 -index
- a is present in 0 -index, starts at 1 -index, ending at 1 -index
- 11 is present in 0 -index, starts at 2 -index, ending at 3 -index
- $sket$ present in $s[3]$, starts at 2 -index, ending at 5 -index
- ck parent in $s[4]$, starts at 3 -index, ending at 4 -index
- p present in $s[5]$, starts at 3 -index, ending at 3 -index.
- $s[0]$ present in $s[4]$, ^{prefix} starts at 0 -index; ending at 2 -index
- o present in $s[1]$, prefix starts at 1 -index, ending at 1 -index
- x present in $s[1]$, prefix starts at 2 -index, ending at 2 -index
- mb present in $s[2]$, prefix starts at 2 -index, ending at 3 -index

Suffix Trie

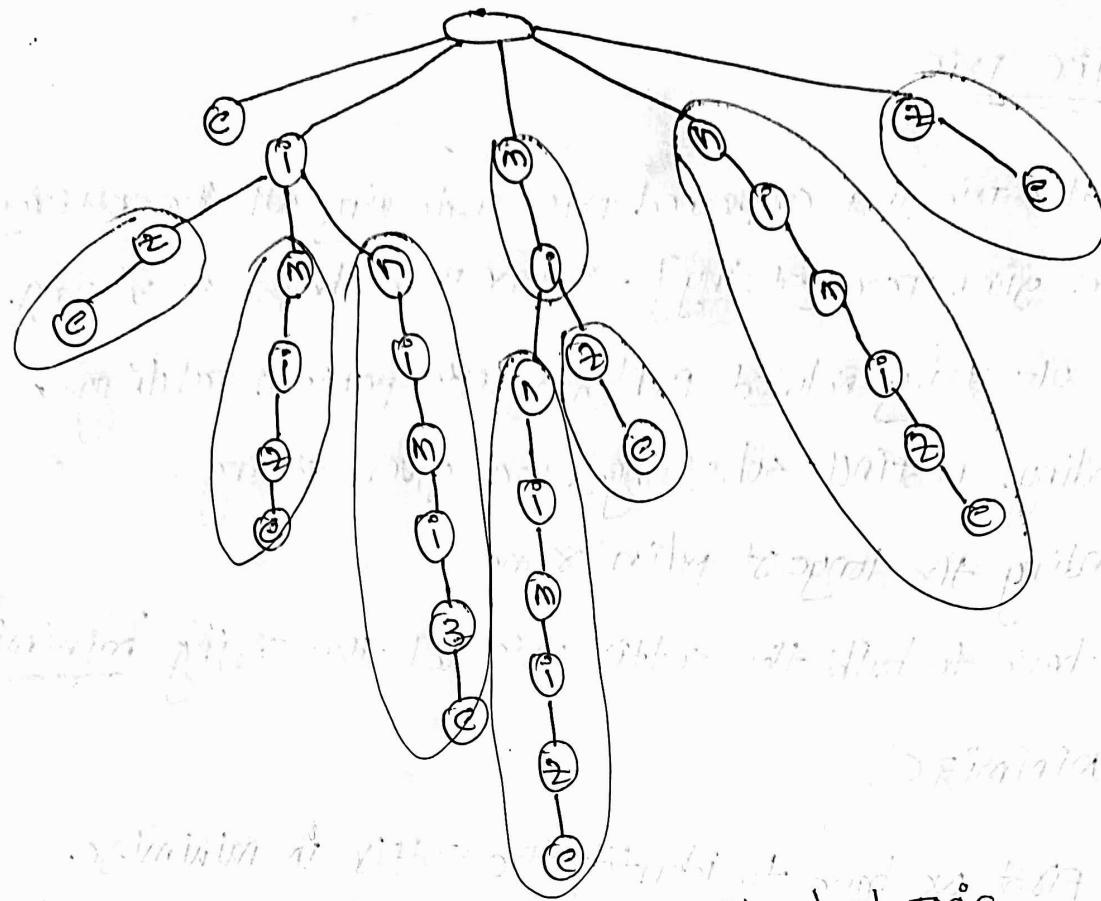
- * Suffix Trie is a compressed Trie containing all the suffixes of the given Text [string]. Suffix Trie helps in solving a lot of string related problems like pattern matching, finding distinct substrings in a given string, finding the longest palindrome.
- * We have to build the suffix Trie for the string "minimize".

Q minimize

A: - First we have to identify the suffix in minimize.
suffix is "e", then include "z" becomes "ze", then include "i" becomes "ize".
then include "m" becomes "mize"
then include "i" becomes "imize"
then include "n" becomes "nimize"
then include "i" becomes "inimize"
then include "m" becomes "minimize"

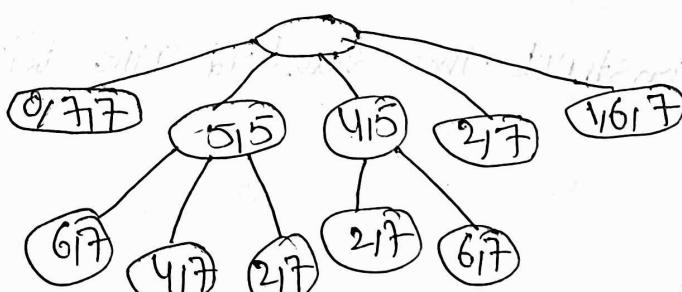
e
ze
ize
mize
inimize
minimize

Then we can construct the standard Trie for this string



Compressed Trie by the above standard Trie

the compressed Trie can be represented
in the tuple representation



~~$s[0] = e$
 $s[1] = z$
 $s[2] = i$
 $s[3] = m$
 $s[4] = i$
 $s[5] = n$
 $s[6] = m$
 $s[7] = i$~~

Rabin-Karp Algorithm [Robin Karp in 1987]

- * Robin-Karp Algorithm is used for searching (or) matching patterns in the Text [string] using a hash function.
- * the hash function is a tool to map a larger input value to a smaller output value. This output is called the hash value.
- * A sequence of characters is taken and checked for the possibility of the presence of the required string.
- * If the possibility is found then character matching is performed.
- * pattern hash value and Text window hash value both are matching then character matching is performed. otherwise no character matching is performed.

Ex. Text = A B C C D D A E F G and
pattern = C C D

so! — n is the length of the pattern and

m is the length of the Text

n=3 and m=10

and is the number of characters in the input set Text

d = 10.

Text window

A	B	C
---	---	---

 C D D A E F G
pattern C C D

* calculate the hash value of the pattern "CCD"

$$\text{hash value for pattern } (P) = \sum_{j=0}^{m-1} (v * d^{m-j}) \bmod 13$$

A - 1
B - 2
C - 3
D - 4
E - 5
F - 6
G - 7
H - 8
I - 9

: v is the ASCII value of the character and
choose any prime number [here we use 13]
in such a way that we can perform all the
calculations

$$\text{pattern } (P) = ((3 \times 10^3 + 4 \times 10^2 + 4 \times 10^1) \bmod 13)$$

$$= (3 \times 10^2 + 4 \times 10^1 + 4 \times 10^0) \bmod 13$$

$$\text{pattern } (P) = (300 + 40 + 4) \bmod 13$$

$$\text{pattern } (P) = 344 \bmod 13 = \underline{\underline{6}}$$

* calculate the hash value of the text window of size m.

for the first text window ABC

$$\text{hash value for Text } (T) = (1 \times 10^3 + 2 \times 10^2 + 3 \times 10^1) \bmod 13$$

$$= (100 + 20 + 3) \bmod 13$$

$$= 123 \bmod 13 = \underline{\underline{6}}$$

* compare the hash value of the pattern with the hash value of the text. If both are matching then only we can check the character-matching. otherwise no need to check character-matching.

* In the example, the hash value of the first window and hash value of the pattern both are matching then we can check the character-matching

Text FIRST window =

A	B	C
---	---	---

 D E F G F G
pattern =

C	D	D
---	---	---

not matching then go to the next Text window.

Step2: If not matching then Text index is incremented (i++)
pattern index remains same

Text second window = A

B	C	D
---	---	---

 E A E F G
pattern =

C	D	D
---	---	---

* hash value of the next window TEXT is

$$\begin{aligned} \text{pattern}(T) &= (2 \times 10^3 + 3 \times 10^2 + 4 \times 10^1) \bmod 13 \\ &= (200 + 30 + 4) \bmod 13 \\ &= 234 \bmod 13 = \underline{12} \end{aligned}$$

$$\begin{aligned} \text{pattern}(T) &= (234) \bmod 13 = \underline{12} \\ &= 233 \bmod 13 = \underline{12} \end{aligned}$$

* hash value of the pattern and the hash value of the Text second window [6 ≠ 12] both are not matching.

Then Text widow is changing. widow is moving to next window.

Step 3 :- Text window = A B C ~~C~~ D DA E F G
 pattern = $\boxed{C \mid D \mid P}$

* hash value of the Text window $y =$

$$= (3 \times 10^3 + 9 \times 10^{3-2} + 4 \times 10^{3-3}) \bmod 13$$

$$= (3 \times 10^2 + 9 \times 10^1 + 4 \times 10^0) \bmod 13$$

$$= (300 + 90 + 4) \bmod 13$$

$$= (394 - 2) \bmod 13$$

$$= 332 \bmod 13 = \underline{\underline{7}}$$

hash

* hash value of the Text window and the value of the pattern both are not matching (6 ≠ 7)

∴ then window is moving one step [i value is incremented]

Step 4 :- Text next window = A B C $\boxed{C \mid D \mid D}$ A E F G.

pattern =

$\boxed{C \mid D \mid D}$

* hash value of the next Text window y

$$\text{hash value of } \del{\text{pattern}}{P} = (3 \times 10^3 + 4 \times 10^{3-2} + 4 \times 10^{3-3}) \bmod 13$$

$$= (3 \times 10^2 + 4 \times 10^1 + 4 \times 10^0) \bmod 13$$

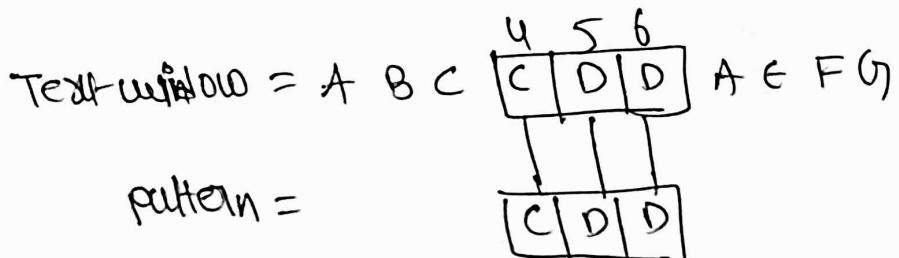
$$= (300 + 40 + 4) \bmod 13$$

$$= (344 - 3) \bmod 13$$

$$\text{hash}(P) = 331 \bmod 13$$

$$= \underline{\underline{6}}$$

\therefore hash value of the next Text window and the hash value of the pattern ($6=6$) both are same. Then we can check the character matching.



\therefore matching found at 4th index.

Continue the procedure for the next characters in the text.

Step 5:-

Knuth morris pratt [Kmp] Algorithm

- * KMP Algorithm is used to find a "pattern" in a "text". This algorithm compares character by character from left to right.
- * But whenever a mismatch occurs, it uses a preprocessed table called "prefix Table" to skip characters comparison while matching.
- * The basic pattern matching Algorithm doesn't study the pattern & what is there in the pattern. w/ it is blindly it is checking every alphabet; just shifting this to the right in case of mismatch.
- * The basic idea of KMP Algorithm is to, just study the pattern and reduce some time, we can see the idea of KMP algorithm.

Ex pattern: a b c d a b c
 1 2 3 4 5 6 7

prefix: a, ab, abc, abcd ..

suffix: c, bc, abc, dabc ..

The idea of this KMP is to ^{inside the string} is there prefix same as suffix.

pattern: a b c d a b c
 | |
 same

The idea is that the beginning of the pattern is appearing again anywhere else in the pattern or not. That is part of IPS or prefix Table. Now see how to prepare Prefix Table.

Eg: P1: a b c d a b e a b f
 0 0 0 0 1 2 0 1 2 0

at 1st index, a is appearing in two times

at 2nd index, b is appearing in two times.

G1 a b c d e a b f a b c
 0 0 0 0 0 1 2 0 1 2 3

↓

→ appearing with 3rd index character

↓

→ appearing with 2nd index character

↓

→ appearing with 1st index character.

G2 a a b c a d a a b e

0 1 0 0 1 0 1 1 2 0

↓

→ appearing with 2nd index character

↓

→ appearing with 1st index character

XG a a a a b a a c d

1 2 3 4 1 2 3 4

appearing with 1st index

↓

→ appearing with 3rd index.

2nd index

Ex:- Text: ab a b a b d
 pattern: ab a b d

Sol:- First prepare the prefix Table (Π) $\pi(i)$ Table for the given pattern.

j	0	1	2	3	4	5
	a	b	a	b	d	
	0	0	1	2	0	

i	0	1	2	3	4	5	6	7
	a	b	a	b	a	b	d	

Now compare $j+1 = i$ then j is incremented and i is also incremented.

Step:- $j+1 = i \Rightarrow 0+1 = 1 \Rightarrow a = a$ then j, i values increment

j	0	1	2	3	4	5	6	7
	a	b	a	b	a	b	d	

i	0	1	2	3	4	5
	a	b	a	b	d	

$b = b$ matching

then i, j values are incremented

Step2:-

1	2	3	4	5	6	7
a	b	a	b	a	b	d

$$j+1 = i$$

$$2+1 = 3$$

$$a = a$$

matching then i, j values are incremented.

0	1	2	3	4	5
a	b	a	b	d	
0	0	1	2	3	

Step3

1	2	3	4	5	6	7
a	b	a	b	a	b	d

0	1	2	3	4	5
a	b	b	b	d	
0	0	1	2	0	

$$\therefore j+1 = i$$

$$3+1 = 4$$

$$4 = 4$$

$$b = b$$

matching then i, j values are incremented.

Step4

1	2	3	4	5	6	7
a	b	a	b	a	b	d

0	1	2	3	4	5
a	b	a	b	d	
0	0	1	2	0	

$$\therefore j+1 = i$$

$$4+1 = 5$$

$$5 = 5$$

$d \neq a$, not matching
then present y is at 4th index.
that corresponding index is 2.

then j is moving to corresponding index: 2.

Step5

1	2	3	4	5	6	7
a	b	a	b	a	b	d

0	1	2	3	4	5
a	b	a	b	d	
0	0	1	2	0	

$$\therefore j+1 = i$$

$$2+1 = 3$$

$$3 = 3$$

$$a = a$$
 matching

then i, j values are incremented

Step 6:-

	1	2	3	4	5	6	7
a	b	a	b	a	b	d	

0	1	2	3	4	5
a	b	a	b	d	
0	0	1	2	3	

$$\begin{aligned}\therefore j+1 &= i \\ 3+1 &= 6 \\ 4 &= 6 \\ b &= b\end{aligned}$$

matching then i, j values are incremented.

Step 7:-

	1	2	3	4	5	6	7
a	b	a	b	a	b	d	

0	1	2	3	4	5
a	b	a	b	d	
0	0	1	2	3	

$$\begin{aligned}\therefore j+1 &= i \\ 4+1 &= 7 \\ 5 &= 7 \\ d &= \underline{d}\end{aligned}$$

matching then i, j values are incremented.

Q1 Text ~~pattern~~: a b a b c a b c a b a b a b d
pattern : ababd.

Sol) First prepare the prefix Table.

0	1	2	3	4	5
	a	b	a	b	d
	0	0	1	2	0

we can add prefix 0th index for performing operation

Boyer moore Algorithm

* Boyer moore Algorithm is the fastest way to find the pattern (P) matching the pattern in a given string (T) Text.

* this algorithm is finding the pattern in a string.

it comparing each character of pattern to find a word of the same character in to the string.

* when character don't match, the search jumps to the next matching position in the pattern by the value indicated in the Bad match Table.

* the Bad match Table indicates how many jumps should it move from current position to the next.

* how to create the Bad match Table, first step is calculate the value of each letter of the pattern to create the Bad match Table using the formula

$$\text{value} = [\text{length of pattern} - \text{index of each letter in the pattern} - 1]$$

(1) $\begin{matrix} 0 & 1 & 2 & 3 \\ a & b & c & d \end{matrix}$

$$\text{value}(a) = [4 - 0 - 1] = 3$$

$$\text{value}(b) = [4 - 1 - 1] = 2$$

$$\text{value}(c) = [4 - 2 - 1] = 1$$
 value of lost letter and other letters

$$\text{value}(d) = [4 - 3 - 1] = 0$$
 that not in substring will be the length of pattern [4]

Letter	a	b	c	d	*
value	3	2	1	0	4

* After that, you can compare the substring [pattern] and the string [Text].

* you start from the index of the end letter in the pattern. In this case the letter is "d".

* If the letter matches, then compare with the proceeding letter "c" in this example.

* if doesn't match, check its value in the Bad Match Table.

* then skip the number of positions that the table value indicates.

table value indicates.

Ex: Text T = This is a Test

Pattern P = Test

Sol: First construct Bad Match

Table Using the pattern:

* the pattern containing Repeated character 'T'.

* Bad match Table doesn't contain Repeated character

* values can be calculated by using the equation

Letter	T	E	S	*
Value	3	2	1	4

value = [length of pattern] - index of each letter in the pattern - 1

$$\text{value}(T) = [4 - 0 - 1] = 3$$

$$\text{value}(E) = [4 - 1 - 1] = 2$$

$$\text{value}(S) = [4 - 2 - 1] = 1$$

Step 2:

Text T =

T	H	I	S	I	S	A	T	E	S	T
---	---	---	---	---	---	---	---	---	---	---

Pattern P =

T	E	S	T
---	---	---	---

 not matching

But 'S' is in Bad Match Table with the value '1'.

means pattern is moving one position to the Right

Step 3:

Text =

T	H	I	S	I	S	A	T	E	S	T
---	---	---	---	---	---	---	---	---	---	---

T	e	S	T
---	---	---	---

 not matching.

* If it's empty character, then consider lost character '*' in a Bad match Table. Empty character is not there in a Bad match Table, then consider lost character '*'. That indicates value 4.

* then pattern is moving 4 positions to the Right.

Step4: Text:

T	H	I	S	I	S	A	T	G	S	T
---	---	---	---	---	---	---	---	---	---	---

T	E	S	T
---	---	---	---

not matching.

and character A is also not there in Bad match Table. Then consider last character 'T' in a Bad match Table, the value is 4. Then pattern is moving 4 positions to the right.

Step5:

Text:

T	H	I	S	I	S	A	T	G	S	T
---	---	---	---	---	---	---	---	---	---	---

T	E	S	T
---	---	---	---

X

not matching.

But is in Bad match Table with the value '1'.

means pattern moving one step to Right.

Step6:

Text:

T	H	I	S	I	S	A	T	G	S	T
---	---	---	---	---	---	---	---	---	---	---

T	E	S	T
---	---	---	---

matching

then compare the prefixes one by one.

∴ pattern is matching at 10th position in a given Text.

