

3-4 geek symbols.

A
C
G
T

gene is a part of DNA which comes out as RNA - converted to protein

3 of the characters make one protein molecule (amino acid molecule).

$\overline{4} \overline{4} \overline{4}$

64 possible

20 in nature (redundancy) Many to one mapping

cell's behaviour - mediated by protein

↓
set of amino acid molecules.

DNA is a string
RNA is also a 1D string of characters
Seq is protein. (3D structure, van Der Waal forces)
It gets folded.
(lowest energy state).

Its shape determines which reaction it will catalyze.

1000-5000 characters - 1 gene.

30,000 genes in a human.

10% of genome is useful, some genes don't come out as RNA.

some change in nucleotides.
dominant gene, recessive gene.

same molecule can be generated by different genes.
64 \rightarrow 20.

sequence one full DNA string - 2001.
sequence sequence eg: HIV virus
 \downarrow
changes itself every 3 days.

Biological Warfare
Gene editing

Given a big DNA - finding a gene in it.
Bigger string if chota string chudna.

KMP algorithm starts

13

```

for (i=0 ; i<n ; i++)
{
    for (j=0 ; j<m ; j++)
    {
        if (a[i+j] != b[j])
            break ;
    }
    if (j==m)
        printf ("found string at %d", i) ;
}

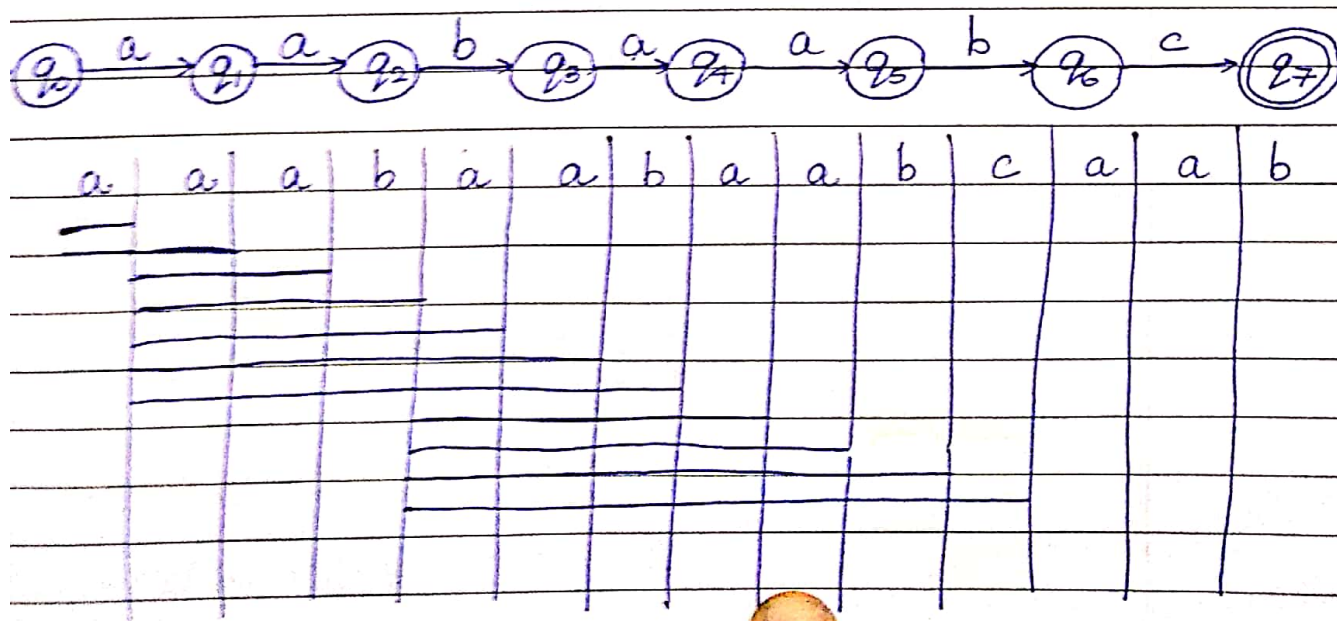
```

$O(nm)$

Construct a finite automata for it.

str to search = "aabaabc"

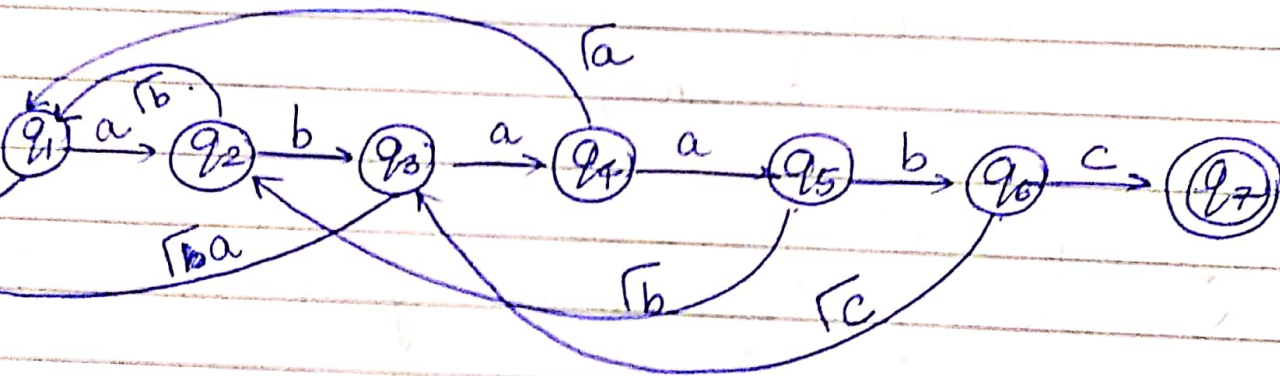
str given = "aaa baabaa bca aab"



If we fail at some stage, which state should we jump to?

what is max^m initial part of string possible at the point of failure.

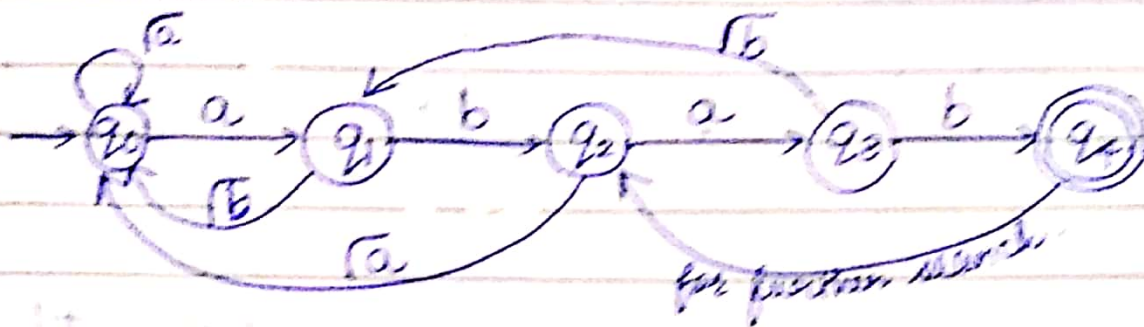
maintain failure table - failure lead to which state?
maintain index to check last matched.



Building of failure table
KMP matching

Pattern = "abob"

Text = "abobooooob"



1	2	3	4
a	b	a	b

1	2	3	4
0	0	1	2

failure states

For this algo, array starts from 1 and string from 1 to m.

Failure Table Algorithm
O/P = Array π

Compute - failure - table (string p)

$m = \text{length}(p)$

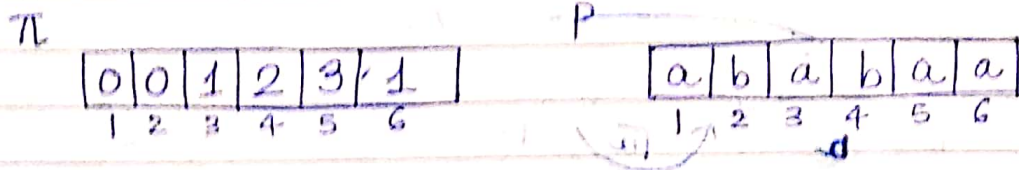
$\pi[1] = 0$

$k = 0$

for $q = 2$ to m

while ($k > 0$ && $p[k+1] \neq p[q]$)

$O(m)$ - Amortized Analysis



Text: "a b b a b a e"
"a b a c a b a b a a"

$\pi[j]$ - failure at j , go to $\pi[j]$

$\pi[j]$ - failure at $j+1$, start comparison of text string at next position with pattern at $\pi[j]$. ?

KMP_matcher (string T , string P)

{

$m = \text{length}(T)$

$n = \text{length}(P)$

$\pi = \text{compute_failure_table}(P)$

$q = 0$

for ($i = 1$ to n)

{

while ($q > 0$ & $P[q+1] \neq T[i]$)

{

$q = \pi[q]$

}

return ($\pi[i]$)

Labels - suffix array

Page No.	
Date	

12	11	2	5	2	1	10	3	7	4	6	3
i	i	i	i	m	p	p	s	s	s	s	(alphabetical length increasing
\$	p	s	s	i	i	p	i	s	i	s	
	p	s	s	s	\$	i	p	i	s	i	
	i	i	i	s		\$	p	p	s	s	everything
	\$	p	s	i		p	i	p	i	s	is
		p	s	s		\$	p	i	p	i	lexicogra- phical
		i	i	s				\$	p	p	
		\$	p	i					i	\$	
			p	p							
			p	p							
			\$	\$							

suffix array is labels on top
now write common elements between two
adjacent strings

LCP array → common prefix
↓
longest common prefix

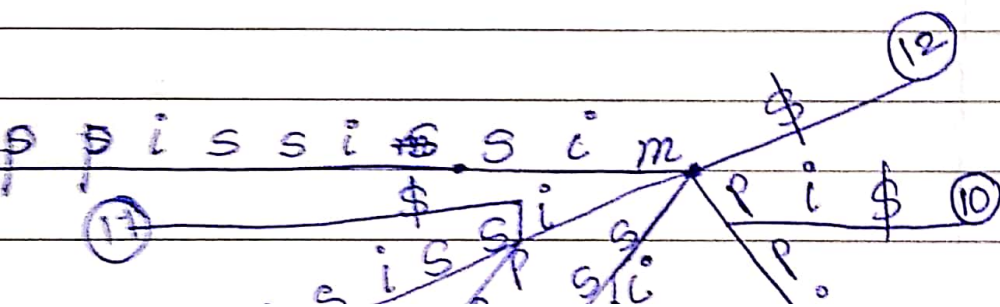
Spina

m i s s i s s i p p i s

suffixes

diffuse tone

Suffix Trees



Boerwies - Wheeler Transformation

Every substring of S is a prefix of some suffix of S .

check whether q is a substring of T .

- follow path of q starting from root.
- If q is exhausted, then q in T .

check whether q is suffix of T .

- follow path of q starting from root.
- if we end at leaf node, then q is a suffix of T .

or $q(\text{string})$

check # of occurrences of a character in T .

- follow the path of q starting from root
- number of leaves under the node reached after exhausting q is # of occurrences

longest repeat in T

- deepest node that has at least 2 leaf nodes

lexicographically first suffix.

→

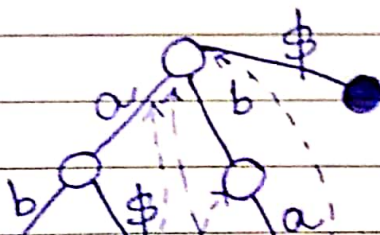
suffix links

- reach a node
- remove first letter on this path
- link to that suffix's node.

Eg) a a b linked to a b
 b a a linked to a a.

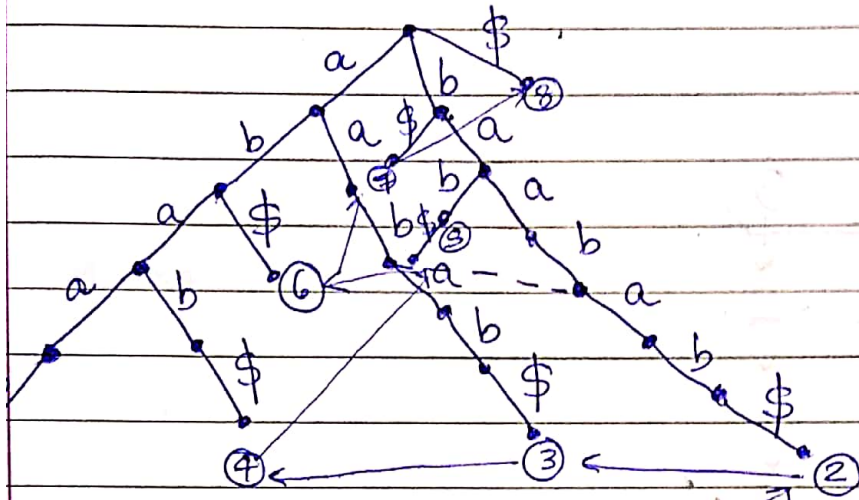
Walk down the tree following q .
If you hit a dead end, save current depth and follow suffix link from the current node.
When you exhaust q , return the longest substring found.

a b a \$.



Ukkonen's Algorithm
↓
constructing suffix tree

Suffix tree for $abaabab\$ - T$



creates suffix
links for every
possible suffix

$ababab\$$
 $baabab\$$
 $oabab\$$
 $abab\$$
 $bab\$$
 $ab\$$
 $b\$$
 $\$$

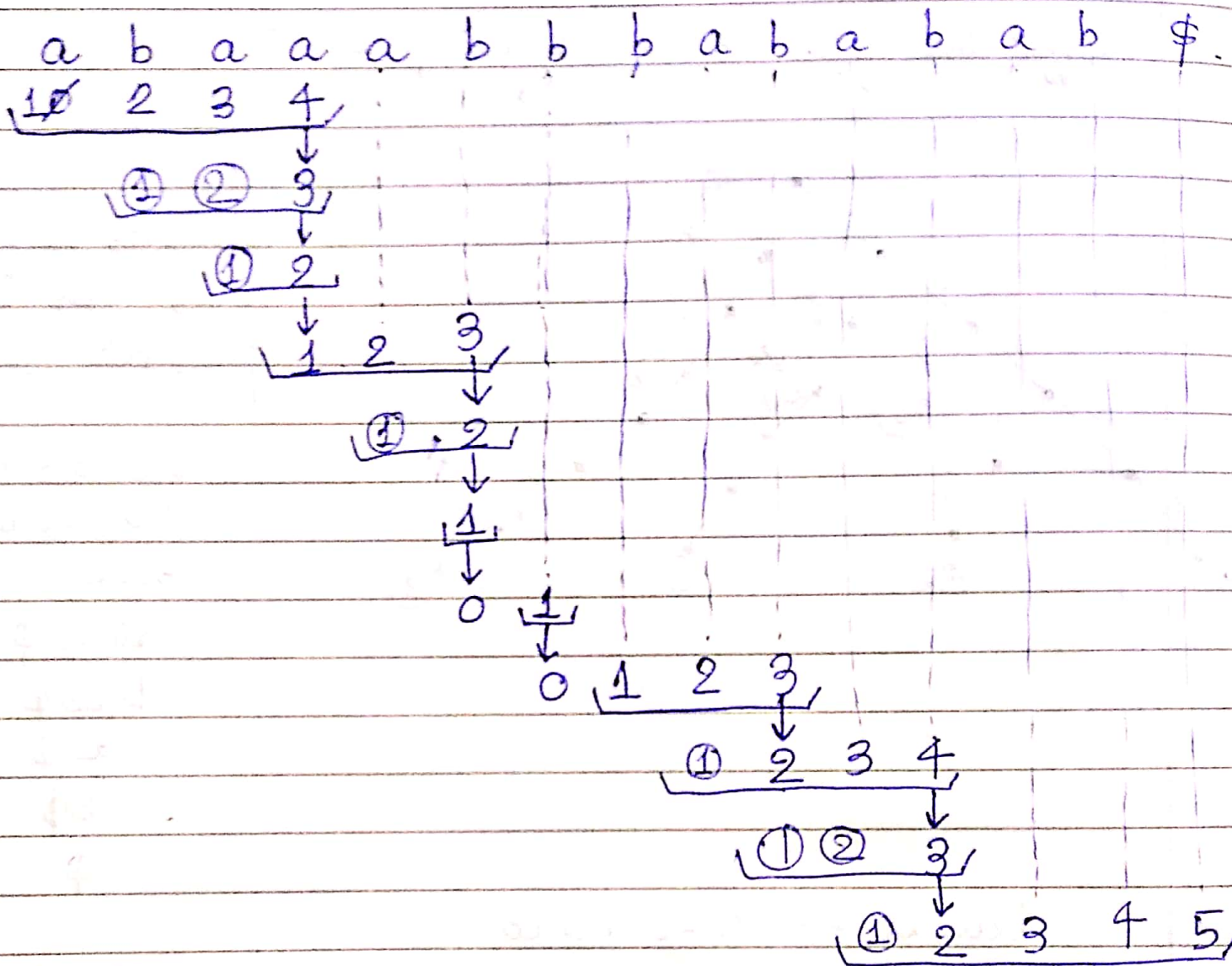
baab. — take baab
remove first character
if we have remaining string in path,
we connect baab to aab.

aab . ab ko aab .
ab b ko ab .

Connect every node to corresponding suffix link node

$abaaa bbbababab\$ - q$

Find longest common substring of q and T .



We have

a b a a

aab

b

bab.

abab

abab\$

baa

aa

ab

How many links?

n characters in string

$$n + (n-1) + (n-2) + \dots = O(n^2)$$

Building the tree - $O(n^2)$

Tree with links - $O(n^2)$

Just the tree - $O(n)$

If string 1 $\rightarrow n$ (build suffix tree)

string 2 $\rightarrow m$

\Rightarrow max^m depth we can go is m in the tree and come back

$$\therefore O(2m)$$

$$\Rightarrow O(m)$$

Uses

How closely two gene strings are related.

$O(m)$ - Amortized Analysis

 π

0	0	1	2	3	1
1	2	3	4	5	6

-P

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

Text: "a b a b a b a a"
 "a b a c a b a b a a"

$\pi[j]$ - failure at j , go to $\pi[j]$.

$\pi[j]$ - failure at $j+1$, start comparison of text string at next position with pattern at $\pi[j]$. ?

KMP matcher (string T , string P)

{

$n = \text{length}(T)$

$m = \text{length}(P)$

$\pi = \text{compute_failure_table}(P)$

$q = 0$

for ($i = 1$ to n)

{

while ($q > 0$ && $P[q+1] \neq T[i]$)

{

$q = \pi[q]$

{

if ($P[q+1] == T[i]$)

$q++$

if ($q == m$)

{

printf ("found substring")

$q = \pi[q]$

{

{

}