Name:- Mohammad Shaique Solanur ; 7062950

---

Q. 04-1

Lemma:  Let s$ be a string of length $|s\$| \geq n$
$$\sqrt{\frac{n}{2}}$$

and let T be its compressed suffix tree.
Then the following holds about T :-

(i) T has exactly n leaves.

(ii) T has atleast n-1 inner nodes other than the root.

(iii) T has atmost 2(n-1) edges.

Proof:-

Let  $\lambda = |$ No of leaves in T$|$

$I = |$ No. of inner nodes of T$|$

$E = |$ Edges of T$|$

W·K·T (we know that) : T is the suffix tree
of s$; There is exactly one leaf for
each of the n suffixes of s$.

$$\therefore \boxed{\lambda = n} \quad - ①$$

Now: Counting edges by incoming edges.

Every node except the root has exactly
one incoming edges. Thus

$$E = |\text{no. of Non-root nodes}| = (I + L) - 1$$

$$\left(\begin{array}{c} \text{For root} \\ \text{node} \end{array}\right)$$

$$\Rightarrow \boxed{E = I + L - 1} \quad - ②$$

**Now; Counting edges by outgoing edges.**

Here, without loss of generality; we can
assume that T is a compressed suffix tree;
i-e No nodes have exactly one child,
or Every inner node has atleast 2
children. $(\geq 2)$. And leaves have 0 children.

$$\therefore \boxed{E = \sum_{i \in T} (\text{No. of children } (i)) \geq \sum_{i \in T} 2 = 2I} - ③$$

Combining ② & ③ we get

$$I + L - 1 = E \geq 2I$$

$\Rightarrow \quad I + L - 1 \geq 2I$

$\Rightarrow \quad L - 1 \geq I$

substituting ① $(L = n)$ we get

$$\boxed{I \leq n - 1} - ④$$

There are atmost $(n-1)$ inner nodes in $T$.

Now, Substituting ④ & ① in ② we get

$$E = I + 2 - 1 \leq (n-1) + n-1) = 2n-2 = 2(n-1)$$

$$\therefore \boxed{E = 2(n-1)}$$

Hence proved.

---

## Q.4.2

Using the hint from the question;

We make relations b/w suffix trees and supermaximal repeats as follows.

1. Inner Nodes $\longleftrightarrow$ repeats.

    a) In suffix trees ; an inner node $v$ has atleast 2 children.

b) Each child is basically a distinct suffix that shares pam label of $v$.

c) Therefore pam label substring w occurs at least twice in the Text.

Every inner node Labels a repeat of the text.

2. Maximal repeats $\longleftrightarrow$ branching + left-maximality

a) Right maximality is taken care of off any innernode because we cannot extend substring w one more character to the right and still cover all of its occurrences with a single node.

(b) To capture left maximality:-
We need to look at the character to the immidiate left of substring w.
In suffix tree, we can do this by looking

at the edges that go to the leaves
at inner node $v$.

(c) If there are at least two distinct
left-contexts, we cannot extend substring
$w$ one more character to the left
and still cover all occurrences.

$v$ is a maximal-repeat node iff it has $\geq 2$ children
and $\geq 2$ distinct left contexts

## 3. Super-maximal repeats

(a) If we have any other maximal
repeat $w'$ and if $w \in w'$ then
$w$ is not supermaximal. Which means
that a super maximal repeat should
not appears inside any other maximal
repeats.

(b) From Tree's perspective it means the node of any other maximal repeat $w'$ is parent or parent-of-parent of the node of substring $w$.
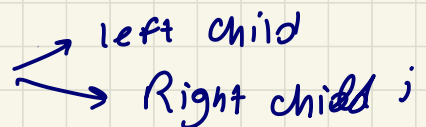
(c) So, for a supermaximal repeat substring $w$, the inner node $v$ should be top-most in the tree. i.e. There should not be any other node above $v$ in the hierarchy which is a node of another maximal substring.

(d) Vice-versa is that there must be no inner node; which is heirarchically lower than $v$; which is maximal.

$v$ is a supermaximal iff it has no descendant which is maximal.

Using Above conclusions, above, we will sketch out an algorithm.

### Step 1 : Preprocessing

(i) Create a Suffix Array of String s and store starting positions in array. [Linear Time]

(ii) Compute lcp array (Linear Time)

(iii) Build child table → left child
→ Right child ;

such that each interval in lcp can list its children.

(iv) Build left context Array :- LC array

for each suffix in Suffix Array, look one character to the left in the

original string. Correspond this Array
with the Suffix Array.

# Step 2
Check for Maximality.

Conditions:-

for any interval [i----j] in Suffix Array ;

check whether the LC table b/w [i---j]
have size ≥ 2. (i.e. those substrings
have atleast two differents characters
immidiately before it)

And.

The interval [i---j] in Suffix Array
should cover atleast two suffixes

i.e. $\boxed{j - i + 1 \geq 2}$

If both conditions hold, then we can say
that the interval [i---j] corresponds to maximal repeat.

## Step 3 Check for Super maximality

This can be done by post-order DFS

(i) **Process** on each **child.**

For each child interval given by [i---j],
we perform DFS. If any child DFS
discovers any maximal repeats in its
subtree; initiate a variable "max_present"
and make it "True"

$\Rightarrow$ max_present == True.

(ii) **Check** current interval

Define a function "max_status" which
if the current interval [i---j] is

maximal or not.

The conditions of this are

$$j - i + 1 < 2 \ \&\& \ len(LC[i:j+1]) \geq 2$$

As provided is <u>step 2</u>

If the condition is True, set "max_status"
== True.

If none of the descedants of the node
given by this interval $[i---j]$ has
"max_present" == False ; then this interval
must by super maximal repeat.

(iii) Go up the heirarchy:-

Return wheather "this subtree has any
maximal repeat at or below $[i---j]$".
That way the parent if one of its descendants

was already maximal or not.

## Step-4

Start this Algorithm at root with
interval [0...n-1]
When initial call finishes, we will have
exactly those repeats which are maximal
but not contained in any larger repeats.

## Conclusion:

1. We preprocess the string into
   Suffix Array, LCA array, Child Array
   and Left Context Array.
   All of these can be done in $O(n)$ time

2. We define a set of conditions to check whether an interval [i---j] is maximal repeat

3. We do one post-order DFS query, and when we hit a maximal repeat that no deeper node has already seen, we output it as supermaximal repeat.

— ✗ End of Algorithm ✗—