

$$\$ < a < \dots < z$$

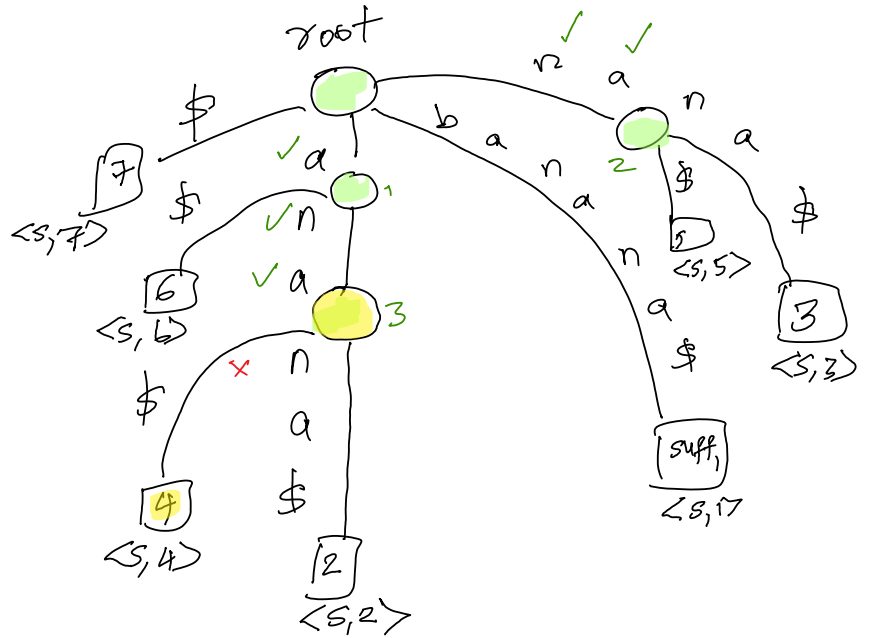
Suffix Trees

Wednesday, March 3, 2021 10:55 AM

Definition: Suffix tree:

Given a string S of length n , over alphabet Σ , the suffix tree of S is the PATRICIA tree for all suffixes of S .

	1	2	3	4	5	6	7
S :	b	a	n	a	n	a	\$
<u>suffix#</u>							
1:	b	a	n	a	n	a	\$
2:	a	n	a	n	a	\$	
3:	n	a	n	a	\$		
→ 4:	a	n	a	\$			
→ 5:	n	a	\$				
6:	a	\$					
7:	\$						



ST: Suffix Tree:

ST properties:

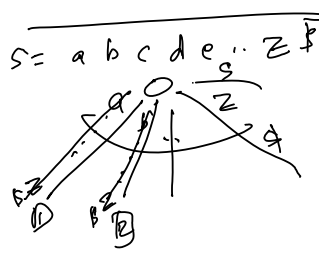
$$\#edges = \#nodes - 1$$

- No. leaves = n
- No. internal nodes $\leq n - 1$

- The pathlabel of leaf i = suffix i of S

Claim: The path-label of every internal node is a repeating substring. (proof in Lemma 1)

Space complexity = $O(n) + O(n^2)$
 ↑ for leaves & internal nodes ↑ for storing all edge labels.
 ← can be reduced to $O(n)$ if we store each edge label as an integer pair.



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

(Think of a worstcase example for this)

Suffix Tree Implementation

Thursday, March 4, 2021 2:40 PM

Q) How to implement the tree data structure (in a space-efficient way)?

Node structure:

```

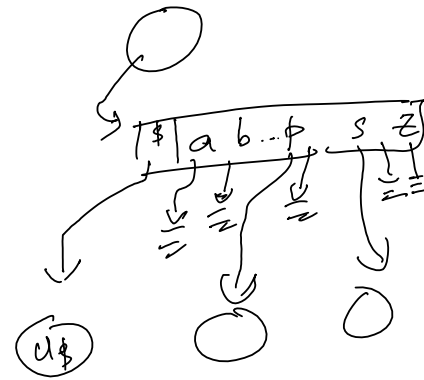
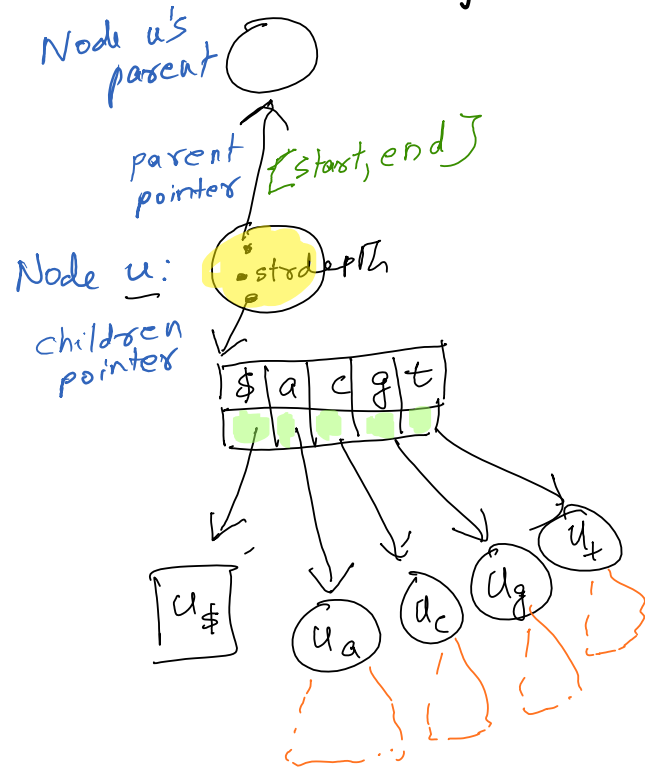
struct/class Node {
4 bytes ← int id;
4 ← int string-depth;
8 ← struct Node *parent;
8 ← pair<int, int> <start, end>;
40 ← struct Node *child [1 ≤ i ≤ 1];
    
```

64 bytes per node
 $\#nodes < 2n$

$n = 10^6$
 $2n \times 64 = 128n$
 space constant

Input string size

$n = 10^6 \Rightarrow 128MB$ tree space
 $10^9 \Rightarrow 128GB$ tree space

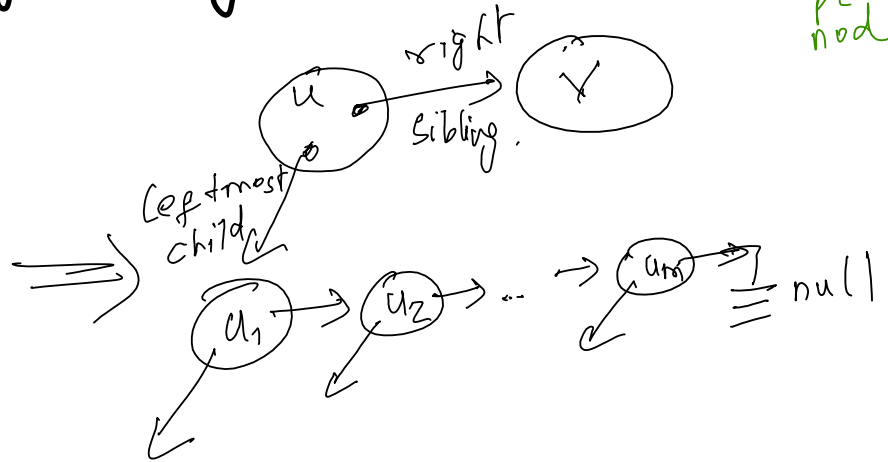
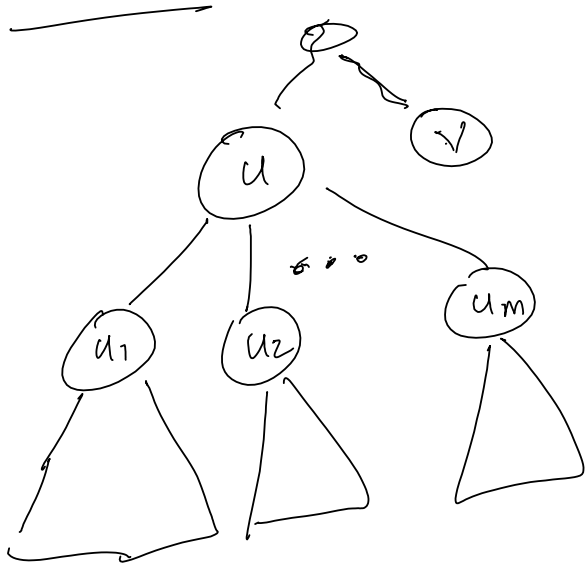


Space efficient Tree Implementation Tricks

Thursday, March 4, 2021 10:45 AM

Given an arbitrary m -way tree T , there are a couple of ways to reduce the space consumption. (memory footprint).

Trick #1 left-child, right-sibling approach



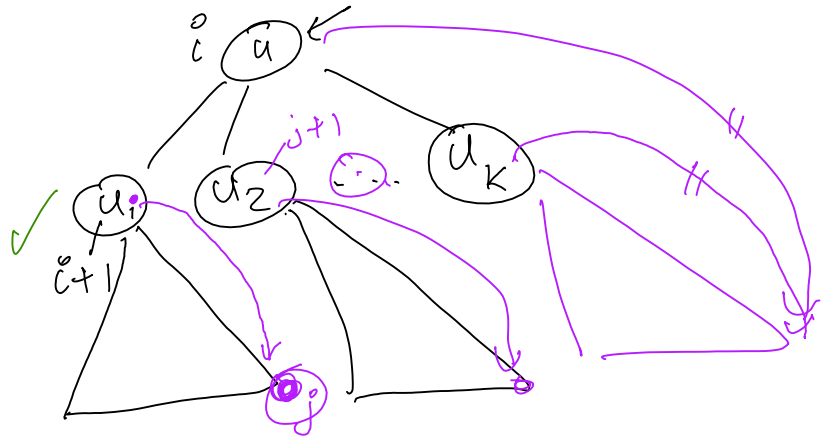
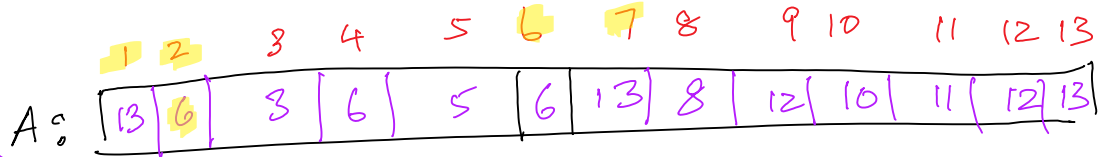
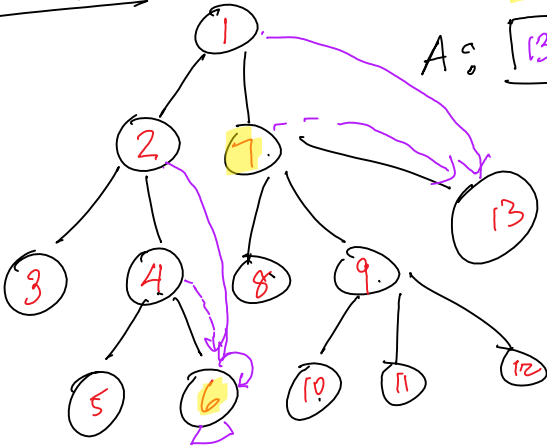
2 pointers per node

1 pointer per node

Trick #2 The Right most Leaf Approach:

- store nodes in DFS order
- Each node stores only one pointer (the pointer to the rightmost leaf under its subtree).

Example:



Definition: A "repeat" of a string s is a substring that occurs at least twice in s .

The Longest Repeat Problem

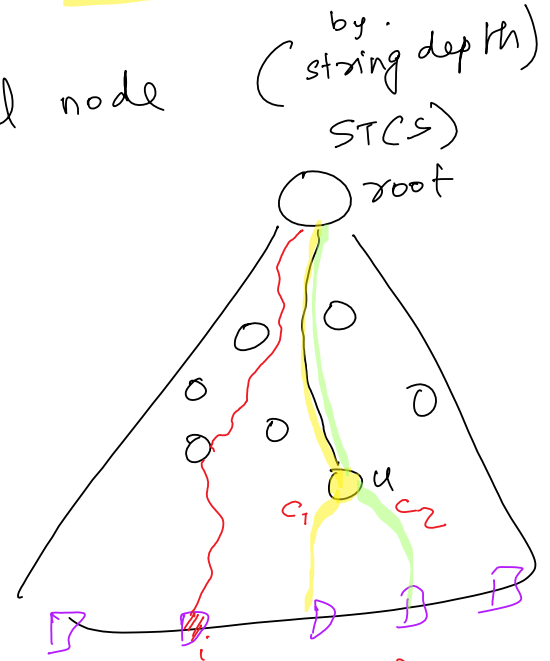
Wednesday, March 3, 2021 11:05 AM

Q1) Given $s [1..n]$, find a longest repeating substring ("repeat")

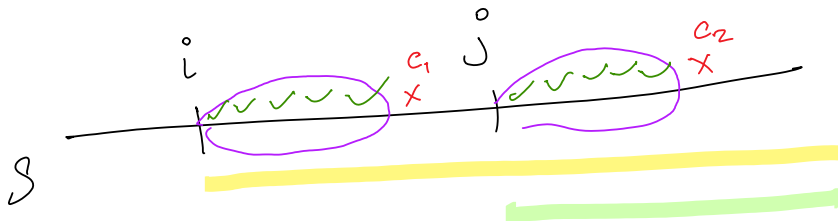
$b \text{ a n } \underline{\text{a n a}} \Rightarrow \text{"ana"}$

Algo:

- 1) Build STCS
- 2) Look for any "deepest" internal node u
→ call this node u
- 3) Output pathlabel (u)



pathlabel (leaf l)
= $s_{\text{uff}} l$



Claim: The algorithm above will output a longest repeat of s always

To prove this claim, we will prove

(i.e., will not miss it) it in two parts

(see next page)

Internal Nodes and Repeating Substrings

Thursday, March 4, 2021 9:51 AM

Lemma 1. Let u be any internal node in $STCS$.
Then the path-label (u) is a repeat in S .

Proof:

Since u is an internal node
 $\Rightarrow u$ has at least two children
 \Rightarrow those two children lead to at least two leaves, say i & j .

\Rightarrow Since $\text{pathlabel}(\text{leaf}_i) = \text{suff}_i$
 and $\text{pathlabel}(\text{leaf}_j) = \text{suff}_j$:

\Rightarrow suff_i and suff_j share
 $(\alpha =) \text{pathlabel}(u)$ as a common prefix

$\Rightarrow \alpha$ is the $\text{lcp}(\text{suff}_i, \text{suff}_j)$ because $c_1 \neq c_2$
 (right maximality)

\Rightarrow And, since $i \neq j$,
 α has to be a (right-maximal) repeat in S . \square

Corollary: If u is a deepest internal node in $STCS$,
 it implies that $\alpha (= \text{pathlabel}(u))$
 is also a longest repeat in S . \square

