

Suffix Tree Construction

Friday, March 5, 2021 10:11 AM

Goal: Given a string $S \in \{1..n\}$ over alphabet Σ ,
Construct the suffix tree of S (STCS).

Naive Algorithm:

Build ST Naive ($S \in \{1..n\}$) {

Init: $T_0 \leftarrow \text{root}$

$O(n)$ for ($i = 1$ to n) do {

$O(n) \rightarrow T_i \leftarrow \text{Insert } \text{suff}_i \text{ into } T_{i-1}$

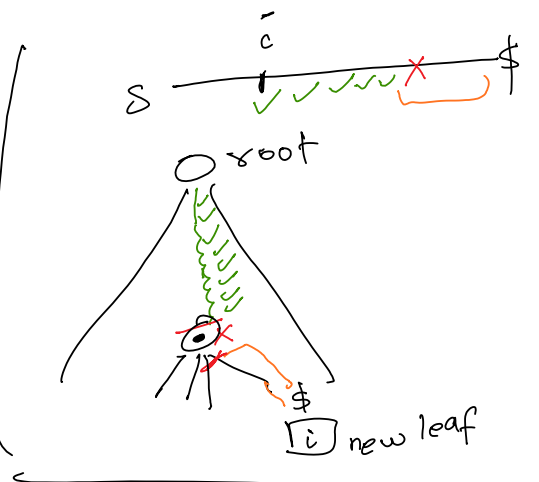
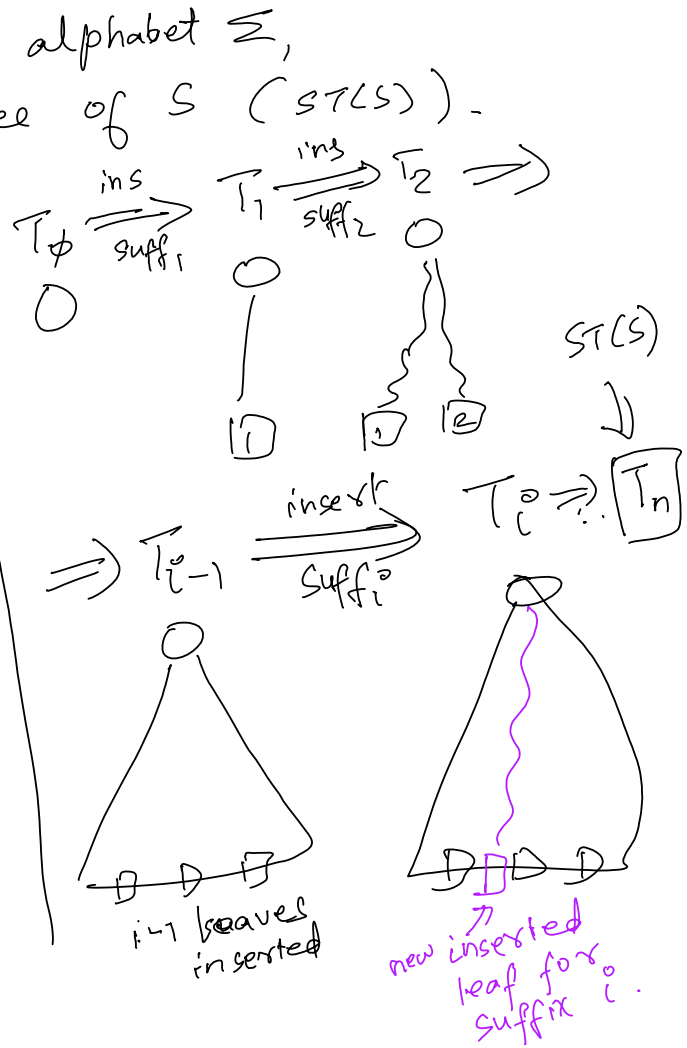
}

STCS $\leftarrow T_n$ // return

}

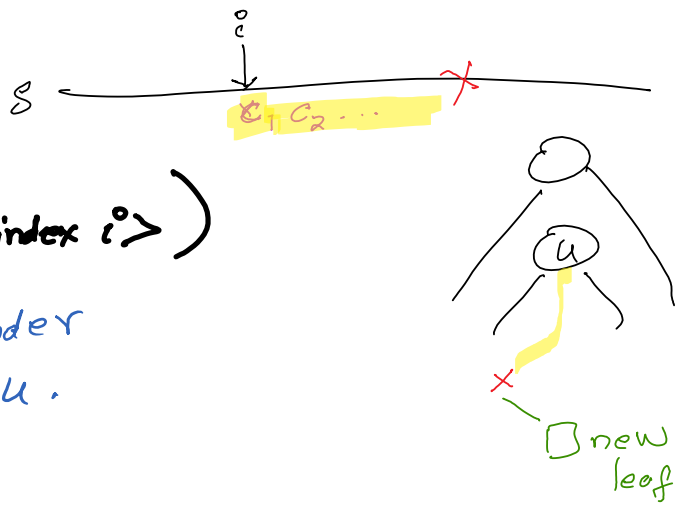
Find Path ($\text{root}, S[i..n]$)

Complexity: worst case = $O(n^2)$ time



FindPath function

Monday, March 1, 2021 10:50 AM



Find Path (Tree node u , \langle string s , index i \rangle)

//Goal: to insert suffix $s[i..]$ under node u .

{
 0) Initialize: $v \leftarrow u$, $x \leftarrow s[i..end]$

1) REPEAT {

a) Find the child branch that starts with $x[i]$

b) if no such branch exists {
 - insert new leaf for s under u

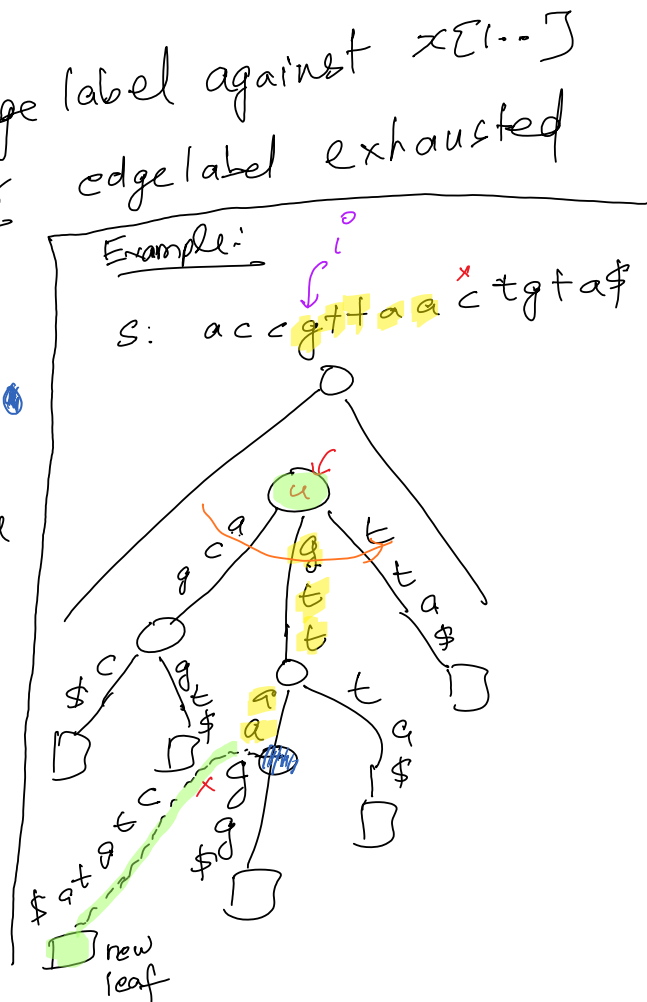
} - return.

c) Compare characters of the edge label against $x[i..]$ until first mismatch or edge label exhausted

d) if mismatch {
 - break edge
 - create a new internal node
 - create new leaf for s under that node
 - return

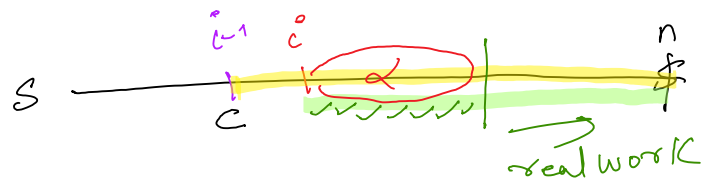
} else {
 $v \leftarrow$ next internal node
 $x \leftarrow x[i+1..]$

} UNTIL (leaf inserted)

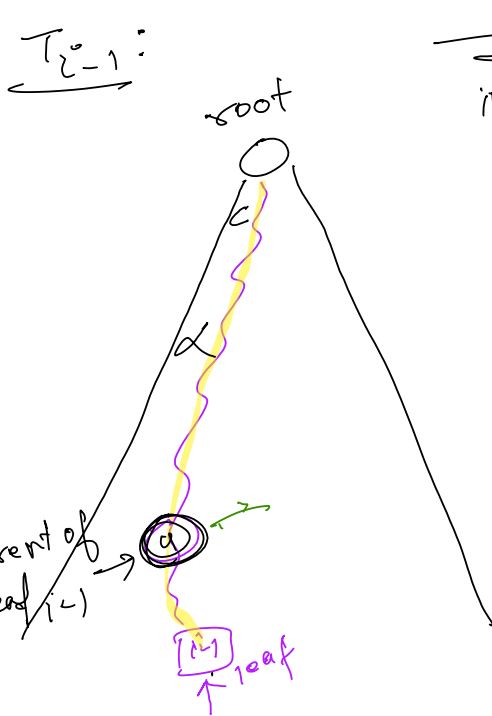


Suffix Links

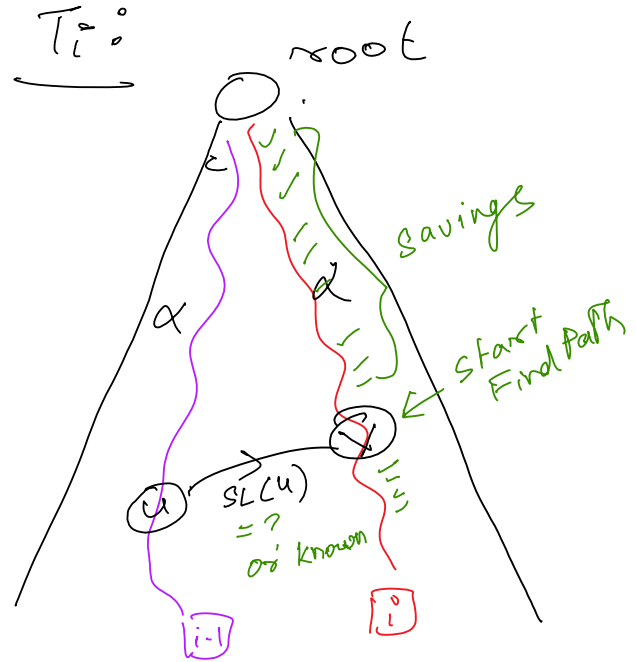
Friday, March 5, 2021 10:21 AM



Q) How to improve Insertion of suff_i into T_{i-1} ?



insert suff_i

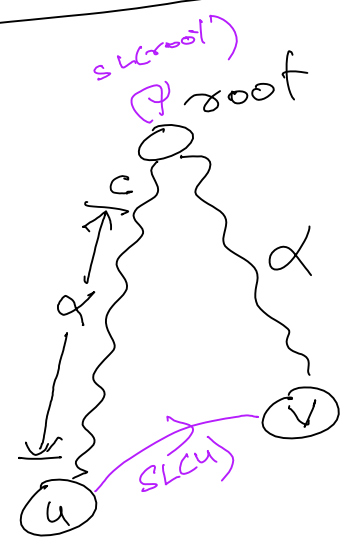


Suffix Link: Let u be an internal node with the path label $c\alpha$.

Then, Suffix link of u ($SL(u)$) is the internal node v with the path label α .

i.e., $SL(u) = v$, for some v pathlabel(v) = α

Special case: $SL(\text{root}) = \text{root}$



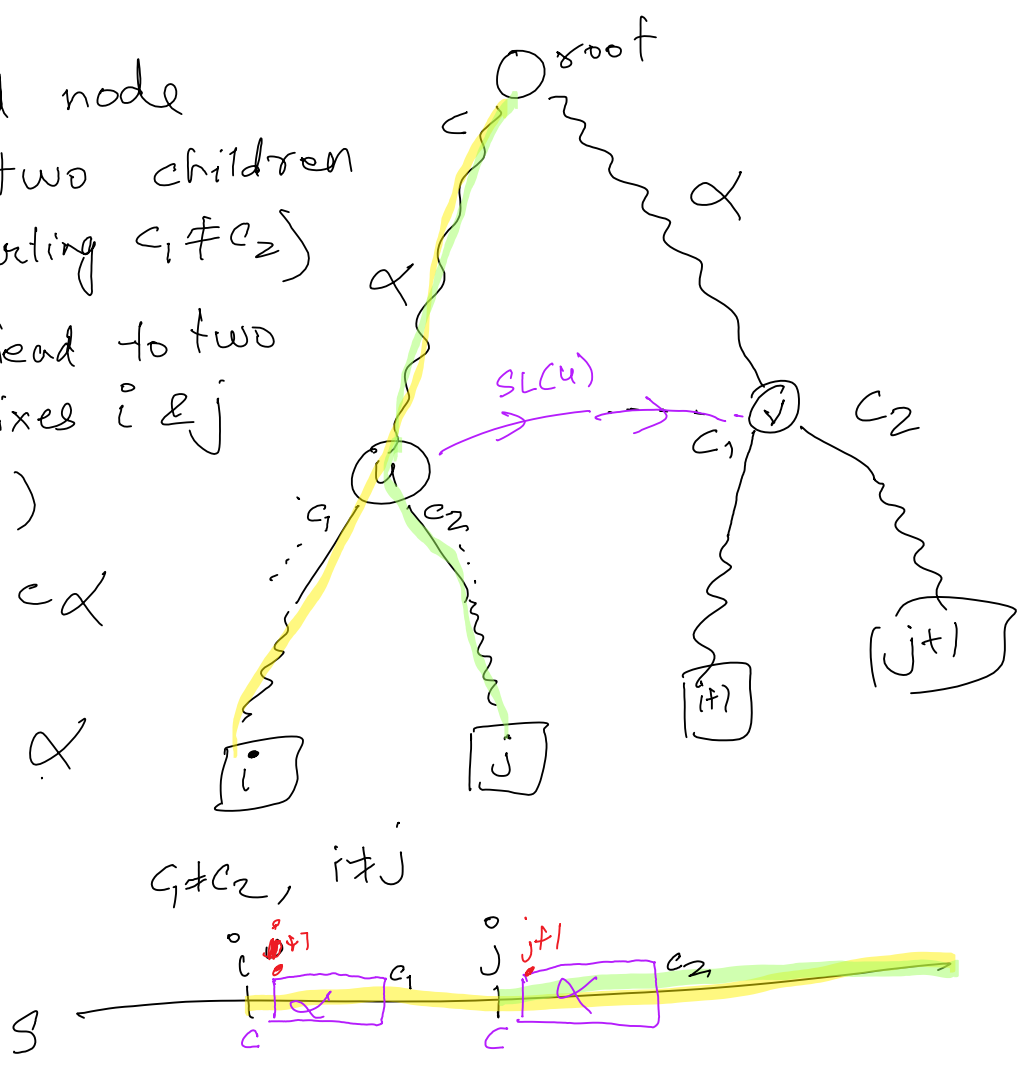
Suffix Links are defined on a suffix tree for all nodes

Friday, March 5, 2021 10:22 AM

Lemma: The suffix link of any internal node u is always defined.
 i.e., $SL(u) = v$, where v will always be a valid internal node in STCS.

Proof:

Since u is an internal node
 \Rightarrow it has at least two children
 (with edge labels starting $c_1 \neq c_2$)
 \Rightarrow those two children lead to two different leaves, suffixes i & j
 ($i \neq j$)
 $\Rightarrow lcp(\text{suffix}_i, \text{suffix}_j) = c_1$
 $\Rightarrow lcp(\text{suffix}_{i+1}, \text{suffix}_{j+1}) = \alpha$
 \Rightarrow there has to exist an internal node with path label α . \square



(on the timing of the Suffix Link of a node will become available for use)

Suffix Link lemma (during construction)

Monday, March 8, 2021 10:16 AM

- Suffix links provide a shortcut, to insert suff_i into T_{i-1}

Q) However as the tree is being constructed (incrementally), is there a guarantee (of some kind) about the timing by which suffix links will be established for newly created internal nodes? \Rightarrow Answer: Yes (see below lemma)

Lemma: At the end of iteration $i-1$,

let u be the parent of the leaf corresponding suff_{i-1} .

If u was newly created during iteration $\#i-1$, then $SL(u)$ will be established latest by the end of iteration $\#i$.

Proof: (let path label $(u) = c\alpha$)

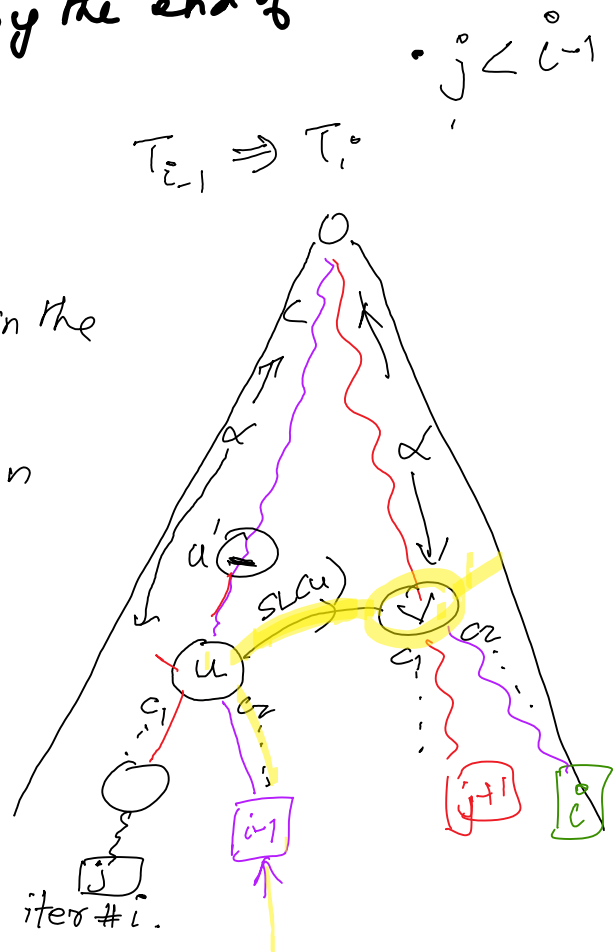
u was created while inserting suff_{i-1}

$\Rightarrow \exists \text{suff}_j$, where $j < i-1$, already in the tree and under u

$\Rightarrow \text{suff}_{j+1}$ should have also already been inserted in the tree

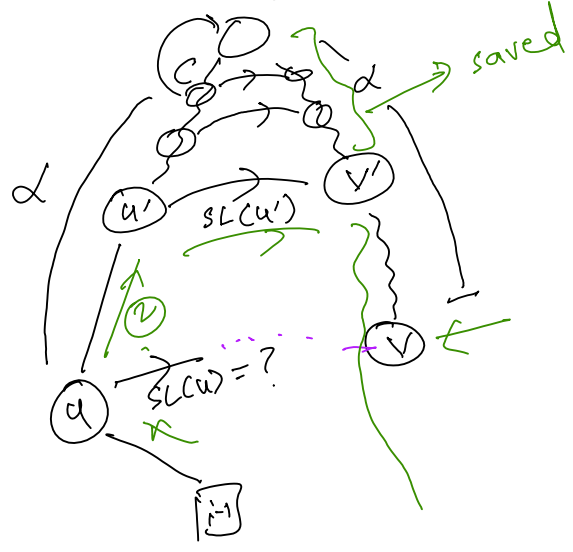
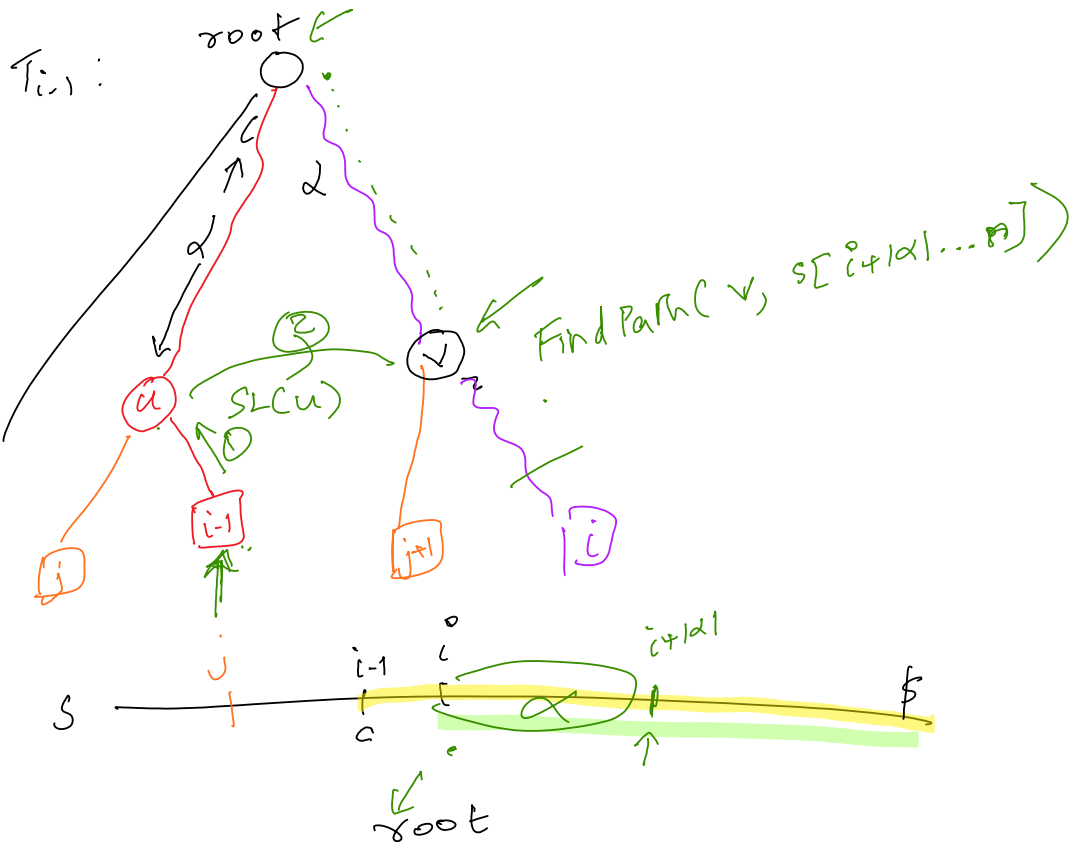
\Rightarrow while inserting the next suffix, suff_i into T_{i-1} , internal node v with path label α must be created (if it didn't already exist).

$\Rightarrow SL(u) = v$ can be established during iter $\#i$.



Recap

Wednesday, March 10, 2021 10:57 AM



McCreight's algorithm for linear time suffix tree construction (1976)

Monday, March 8, 2021 10:22 AM

Algorithmic Pseudocode:

Build ST($S[1..n]$, ϵ)

$T_0 \leftarrow \text{root}$

for ($i = 1$ to n) { // for each suffix

Let $u \leftarrow$ parent of leaf.

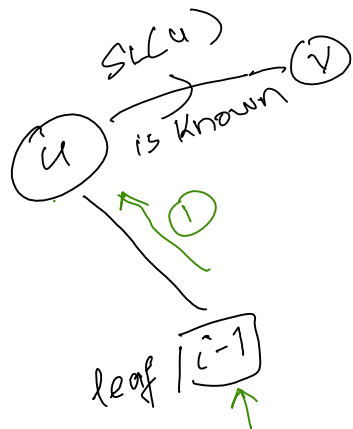
Cases

$\rightarrow u$ already exists ($SL(u) = \text{known}$)

$\rightarrow u$ is new ($\therefore SL(u) = \text{unknown}$)

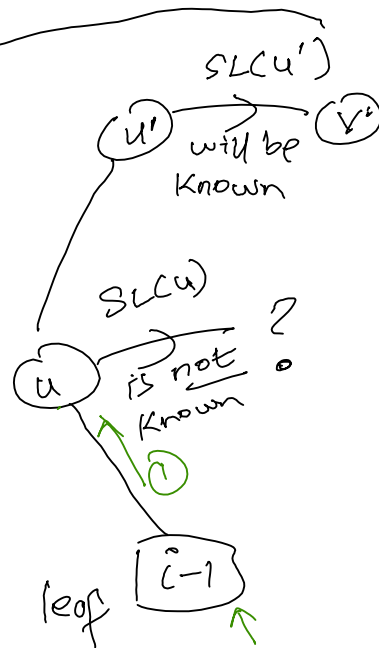
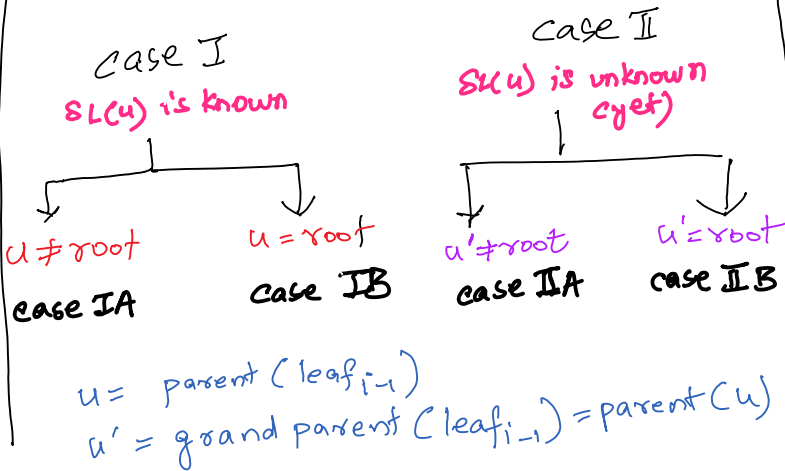
insert suffix into $T_{i-1} \Rightarrow T_i$

} Output $T_n \Rightarrow \text{ST}(S)$



Case I

Case Taxonomy



Case II

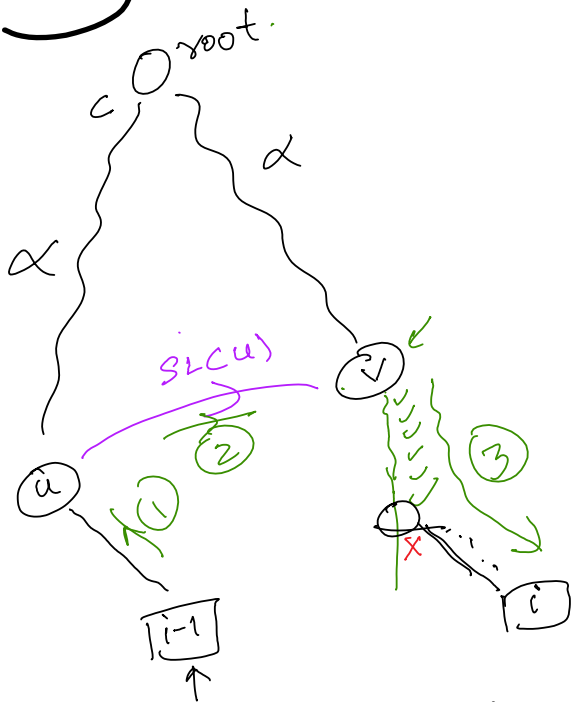
McCreight's algorithm: Cases

Monday, March 8, 2021 10:23 AM

$$|\alpha| = v.\text{string depth.}$$

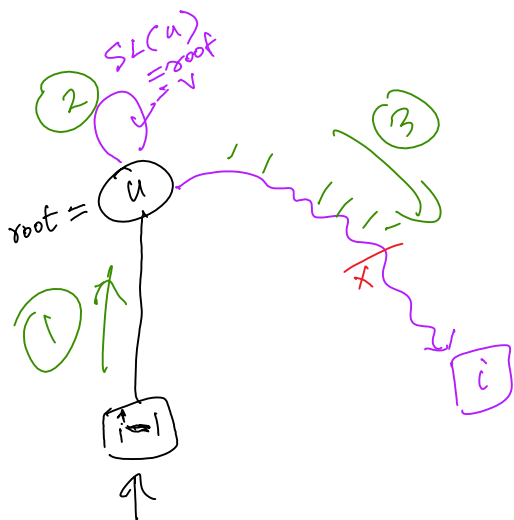
Let $u \leftarrow$ parent of leaf i_{i-1} .

Case IA) $SL(u)$ is known & u is not the root



- ① go to parent u
- ② Take $SL(u)$ to v
- ③ FindPath(v , $S[i+|\alpha| \dots n]$)

Case IB) $SL(u)$ is known & u is the root



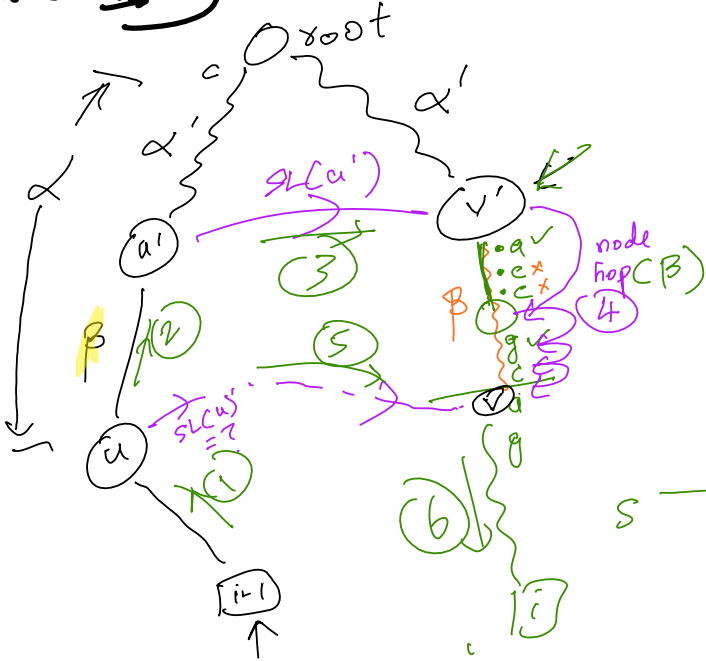
- ① Go to parent u ($= \text{root}$)
- ② Take $SL(\text{root}) \Rightarrow (v = \text{root})$
- ③ FindPath($v = \text{root}$, $S[i \dots n]$)

$$c \alpha = c \alpha' \beta$$

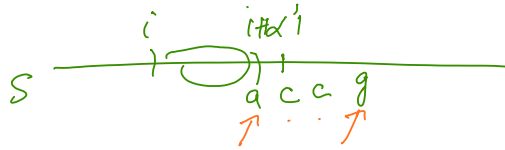
McCreight's algorithm: cases

Monday, March 8, 2021 10:26 AM

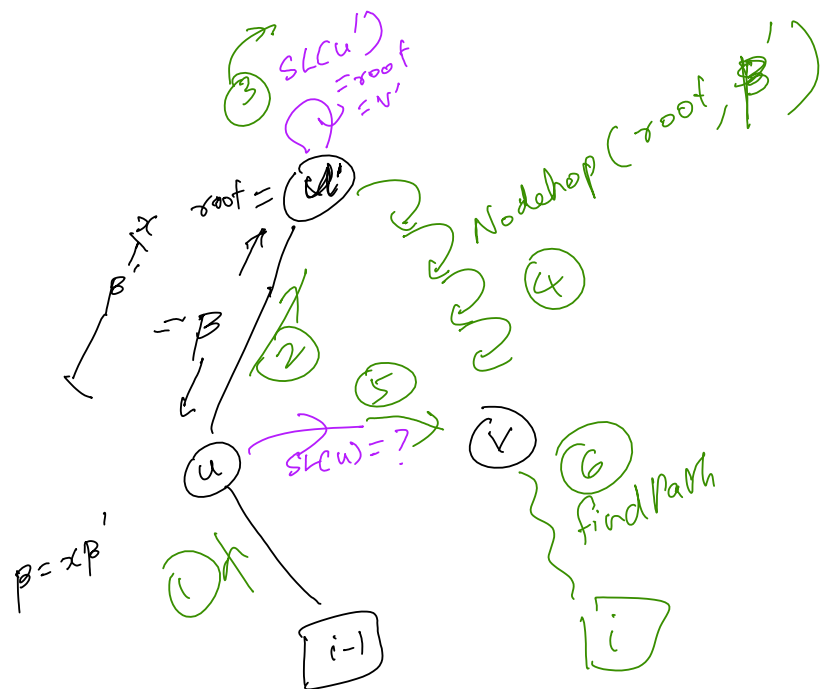
Case II A) $SL(u)$ is unknown and u' is not the root.



- ① Go to parent u
- ② Go to grand parent u'
 $\beta \leftarrow \text{edge label } (u' \text{ to } u)$
- ③ Take $SL(u')$ to v'
- ④ Nodehops ($v', SL[i+|\alpha'|] \dots$) to v
create v if necessary - arry
- ⑤ Establish $SL(u) = v$
- ⑥ FindPath($v, \beta[i+|\alpha'| \dots n]$)



Case II B) $SL(u)$ is unknown and u' is the root.



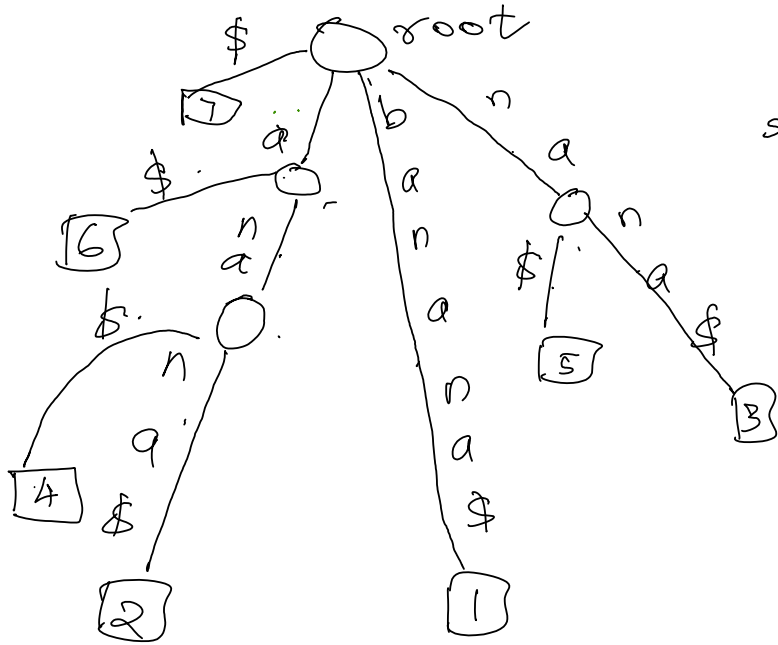
- ① go to parent u
- ② go to grand parent u'
- let edge label (u' to u) be $\beta = x \beta'$
- ③ Take $SL(u')$ to root $= v'$
- ④ Node Hops(root, β') to v
- create v if necessary
- ⑤ Establish $SL(u) = v$
- ⑥ Find Path($v, \beta[i+|\beta'| \dots n]$)

Project 2: Testing your ST construction code

Monday, March 8, 2021 2:18 PM

Use small examples:

$s = \text{banana}\$$



$\$ < a < c < g < t$

1 2 3 4 5 6 7
s: b a n a n a \$

① Basic Statistics

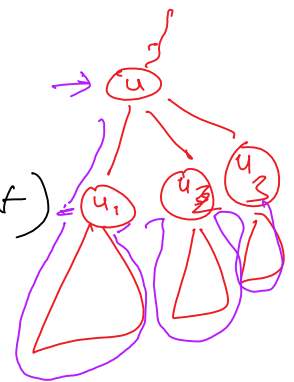
#nodes
 { # internal nodes (4) # leaves (7) }
 #edges = 10

some other examples to try:

mississippi\$
 applemapleapply\$
 anantharamankalyanaraman\$

② DFS disp (node *T) Display function

```
{
  // print contents of curr. node (T)
  For every child of curr node (from left to right)
    DFS disp (child)
}
```



Project 2: Testing your ST construction code

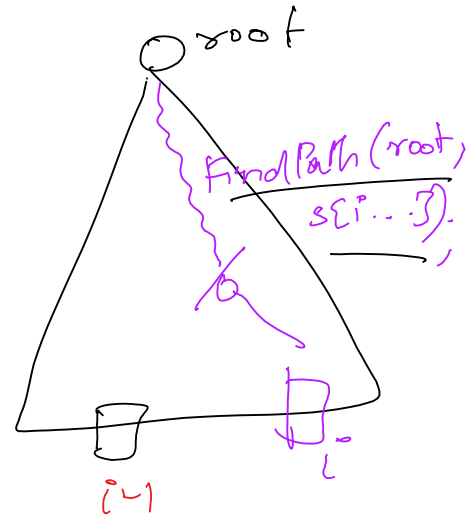
Monday, March 8, 2021 2:18 PM

Build ST Naive ($s[1..n]$) // i.e., implement a no-suffix link version for testing

$T_0 \leftarrow \text{root}$
 For ($i = 1$ to n)

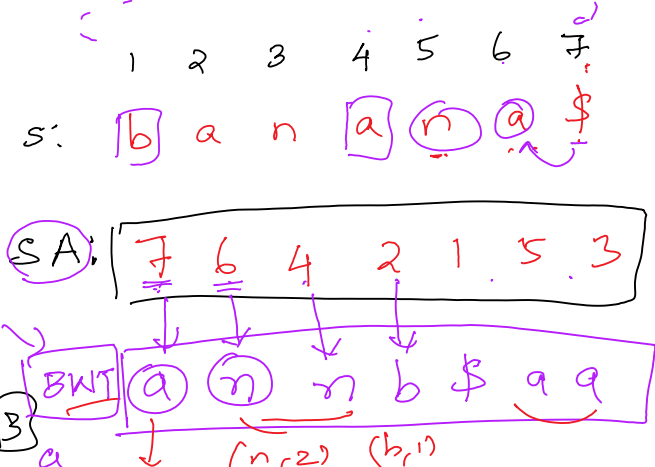
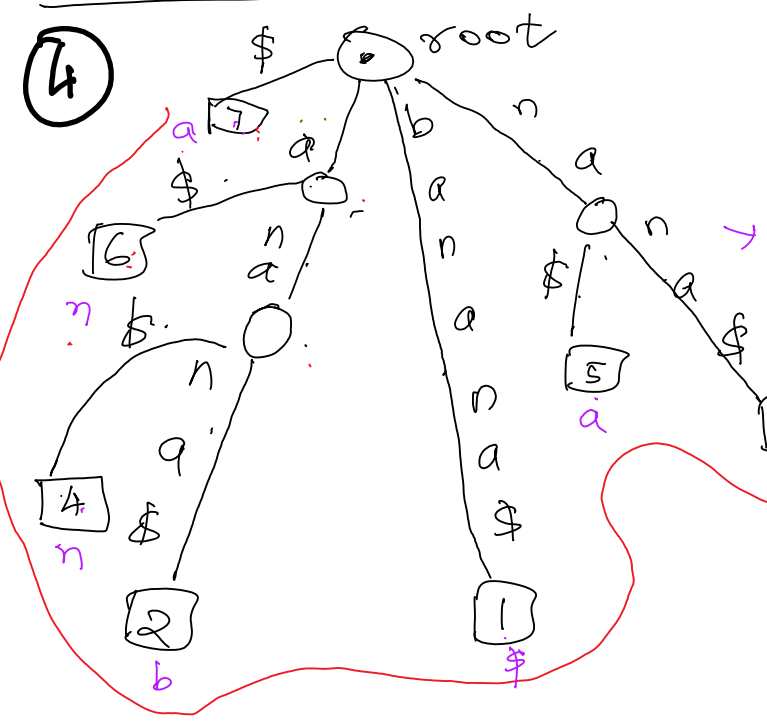
$T_i \leftarrow$ "Insert" suff_i into T_{i-1}

FindPath ($\text{root}, s[1..i]$)



3
 Naive

4



BWT ← Burrows Wheeler Transform

DFS-BWT

a
 n
 n
 b
 \$
 a
 a

If BWT matches, then code most likely right.