

# Lec20-more-suffix-trees-arrays

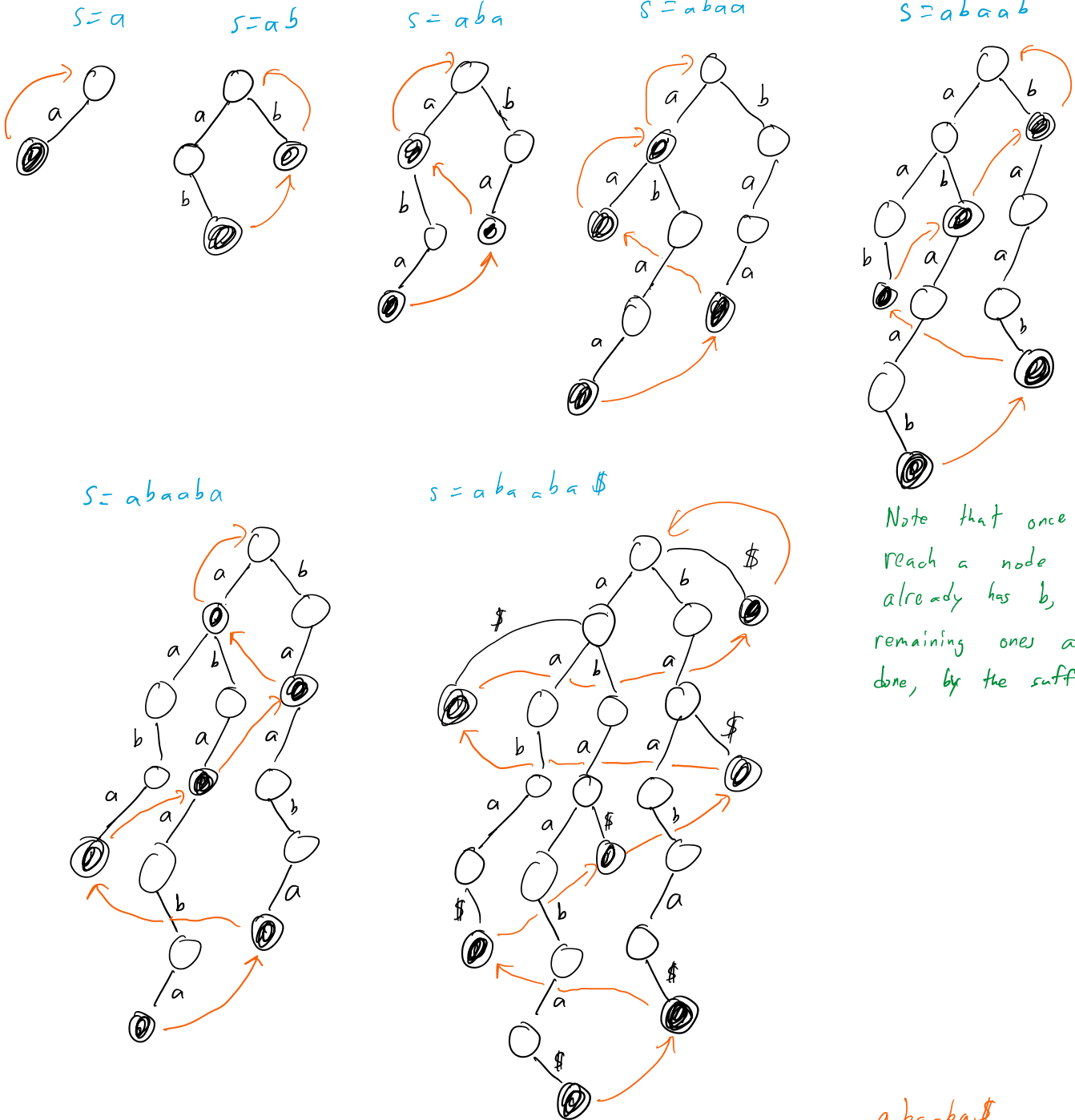
Sunday, October 20, 2024 10:14 PM

Last time we learned about suffix tries + trees.

Let's remind ourselves of suffix trie construction.

We build it up while walking the string from left to right.

$s = abaaba\$$



Note that once you reach a node that already has b, all remaining ones are already done, by the suffix link property

$abaaba\$$

Ex  $bbaa$

Longest Common Substring - L S

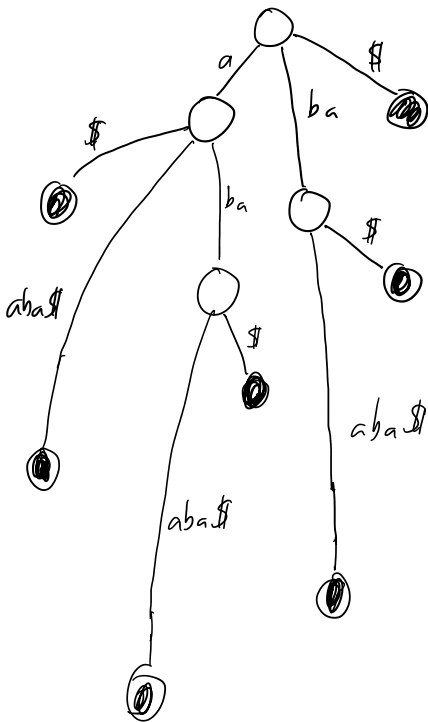
## Longest Common Substring $q \neq S$

Walk down  $q$  as far as possible

If dead end, save depth, follow suffix links, keep walking

When  $q$  exhausted, return longest substring found.

We can also compress it to a suffix tree.



Ex bbaa

Start at b

Dead end.

Start at  $\emptyset \rightarrow b \rightarrow a \rightarrow a$ .

Done.

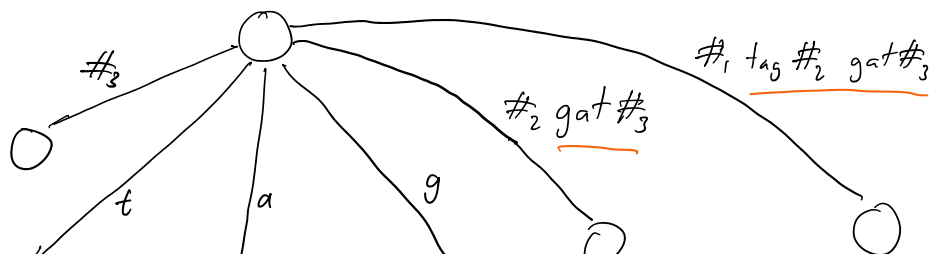
## Generalizing suffix trees to multiple strings (GST)

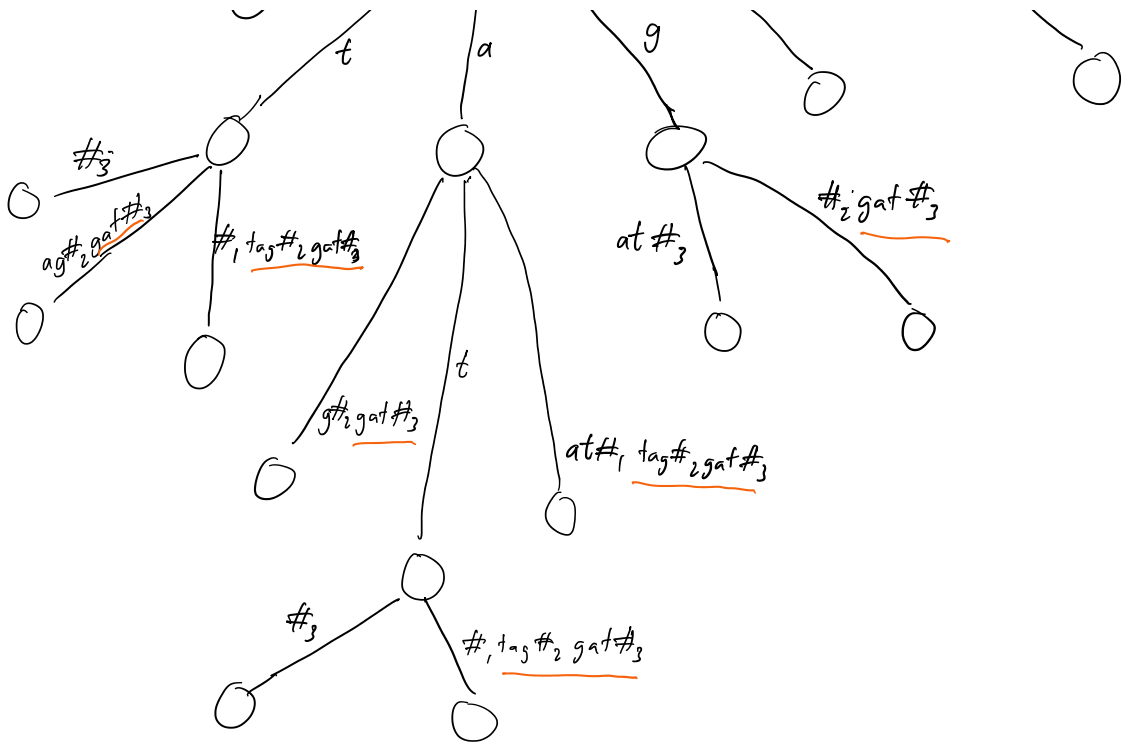
Goal: represent set of strings  $P = \{S_1, S_2, \dots, S_m\}$

Ex. aat, tag, gat

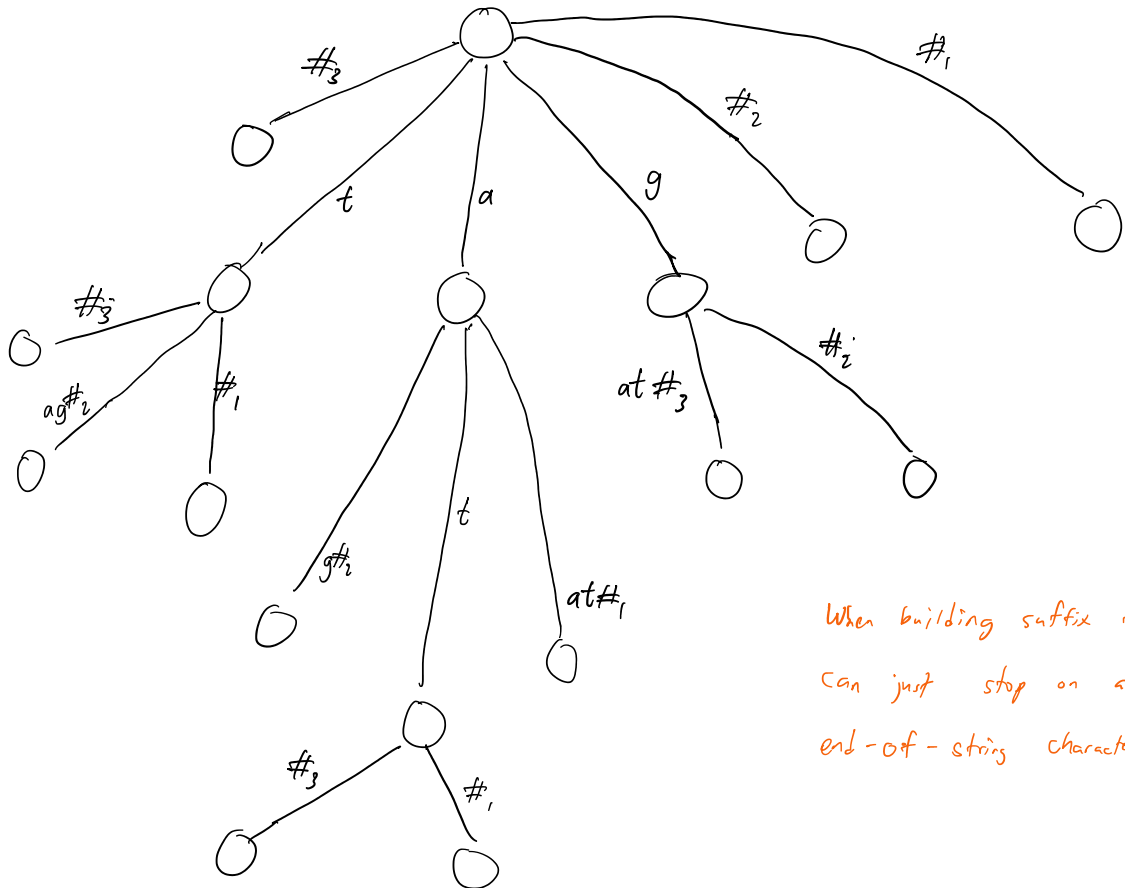
Sol. build suffix tree for aat#<sub>1</sub> tag#<sub>2</sub> gat#<sub>3</sub>

unique string terminators





Notice the #s are only on leaf nodes because they are unique  
 Remove any text after first # on leaf.



When building suffix tree,  
 can just stop on any  
 end-of-string character #.

## Applications:

Longest common substring of  $S, T$

Build generalized suffix tree for  $\{S, T\}$ .

Find deepest node that has descendants from both strings  
(containing both  $\#_1$  &  $\#_2$ )

Determine strings in database  $\{S_1, \dots, S_m\}$  that contain query  $q$ .

Build GST.

Follow path for  $q$ .

If we end at node  $u$ , traverse the tree below  $u$  & exit if  $\#_i$  is found.

Longest Common Extension at  $i, j$

Given strings  $S$  &  $T$ , want to find longest substring of  $S$  starting at  $i$  that matches a substring of  $T$  starting at  $j$ .

$S$  —————  
                   $LCE(i, j)$   
                  |  
                   $i$

$T$  —————  
                   $LCE(i, j)$   
                  |  
                   $j$

Ex. AAAATTTT CCCCGGGG  
                   $i=9$

TTTT CCCCGGGG AAAA  
                   $i=5$

Can solve by walking on both strings until mismatch, but we want to do this faster by preprocessing.

Sol. Build GST for  $S, T$ ,  $O(|S| + |T|)$  time.

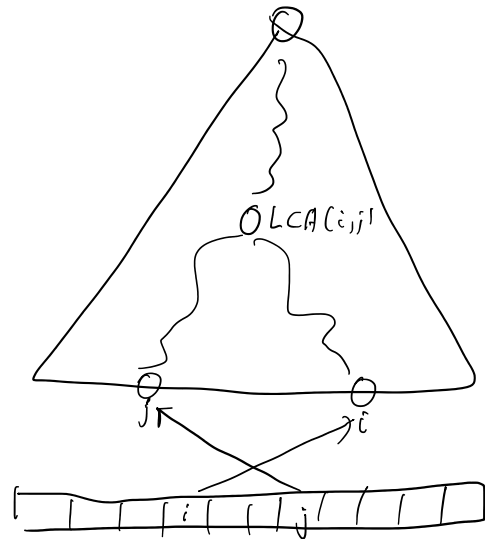
Preprocess tree so that lowest common ancestors can be found in constant time. } We are not covering LCA in 15351, but exists complicated data structure that takes  $O(|S| + |T|)$  time to build initially.  
In 15351/02613, we take for granted this LCA data structure.

Create array mapping suffix numbers to leaf nodes.  $O(|S| + |T|)$  time.

Given query  $(i, j)$ :

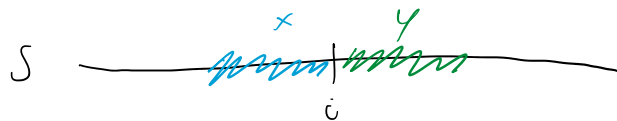
Find leaf nodes for  $i, j$ .  $O(1)$

Return string of LCA for  $i, j$ .  $O(1)$



### Using LCE to find palindromes

Maximal even palindrome at pos.  $i$ : the longest string to the left & right such that left half = reverse of right half.

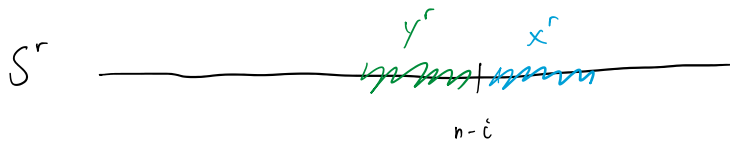


Ex. GATTAG

$y = \text{reverse}(x)$

Goal: find all maximal even palindromes.

Construct  $S^r = \text{reverse}(S)$ .  $O(|S|)$



Preprocess  $S$  &  $S^r$  so LCE queries take  $O(1)$  time.  $O(|S|)$

$LCE(i, n-i) = \text{length of longest palindrome centered at } i$

For pos  $i$ :

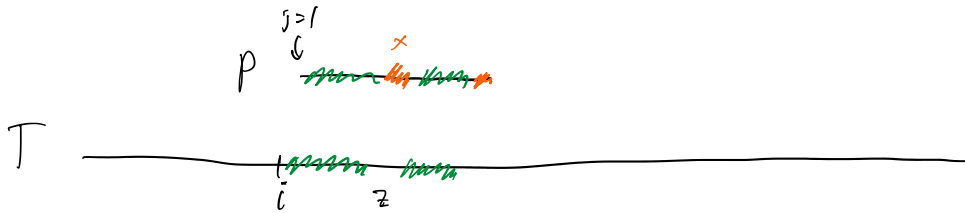
Compute  $LCE(i, n-i)$

$O(|S|)$   
 $O(1)$  }  $O(|S|)$

## k-mismatch using LCE

Given long reference string  $T$ , & a short string  $P$ , check if  $\exists$  a k-mismatch of  $P$  starting at pos  $i$  of  $T$ .

$\hookrightarrow$  k substitutions of characters in  $P$  to be exact match.



$j=1$  pos in  $P$   
 $c=0$  # mismatches so far

$O(k)$  { Repeat until  $c > k$ :  
 $j = j + \text{LCE}(i, j) + 1$   
 $i = i + \text{LCE}(i, j) + 1$   
If  $j \geq |P| + 1$ , return True. }  $O(1)$  to get  $\text{LCE}(i, j)$   
 $c = c + 1$  matched all of  $P$   
return False

Finding all k-mismatches takes  $O(k \cdot |T|)$  time

## Suffix arrays

Suffix trees are  $O(n)$  space, but the constant is large.

( $\approx 20$  bytes per char in good implementation)

10 Gb genome = 200 GB to store.

Suffix arrays are more efficient even than suffix trees & do most of the same ops.

Idea: Sort all suffixes of string lexicographically & keep the indices.

Idea: Sort all suffixes of string lexicographically & keep the indices.

$S = \text{cat}t\text{cat} \$$

- 1 cat $t$ cat $\$$
- 2 at $t$ cat $\$$
- 3 t $t$ cat $\$$
- 4 tcat $\$$
- 5 cat $\$$
- 6 at $\$$
- 7 t $\$$
- 8  $\$$

sorted  $\rightarrow$

- 8  $\$$
- 6 at $\$$
- 2 at $t$ cat $\$$
- 5 cat $\$$
- 1 cat $t$ cat $\$$
- 7 t $\$$
- 4 tcat $\$$
- 3 t $t$ cat $\$$

Search substring by binary search of suffixes.

Counting is easy because all suffixes starting with "at" are next to each other.

App: k-mer counting for all k-mers.

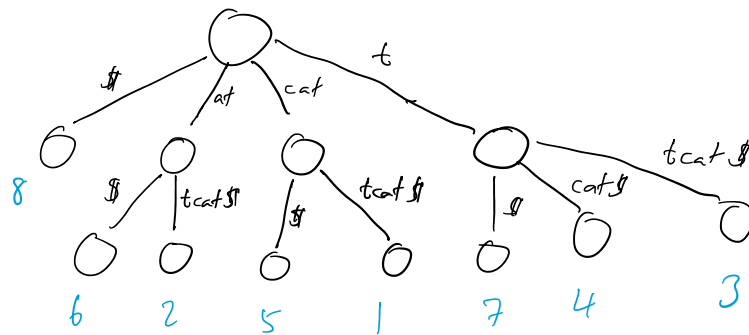
↳ only keep indices & original str  $S$ .

Construction:  $O(n^2 \log n)$  by sorting suffixes, where each comparison takes  $O(n)$  time.

But better  $O(n)$  algorithms, such as building suffix tree.

$\Sigma = \{ \$, a, c, t \}$

$S = \text{cat}t\text{cat} \$$   
1 2 3 4 5 6 7 8



If edges sorted lexicographically, leaf labels are exactly suffix array.