

# Lecture 9 - Trays, ~~Mismatches~~ and Wild cards, level ancestor

Last Time: Text  $T$  of length  $n$ , Pattern  $P$  of length  $m$ , Both over alphabet  $\Sigma$

## Suffix Trees (ST)

$O(n|\Sigma|)$  space or  $O(n)$  space  
 $O(m)$  query or  $O(m \log |\Sigma|)$  query

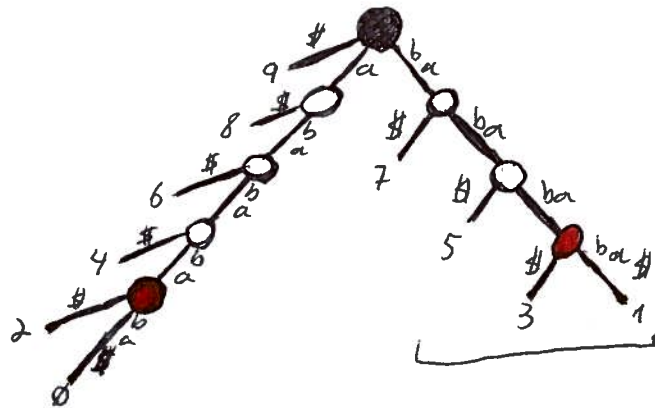
## Suffix Arrays (SA)

$O(n)$  space  
 $O(m + \log n)$  query

$\Rightarrow$  Suffix Tray [Cole, Kopolowitz, Lewenstein] IICALP 2006  
 $O(n)$  space  
 $O(m + \log |\Sigma|)$  query

$T = ababababa \#$   
 0 1 2 3 4 5 6 7 8 9

Idea: tree decomposition; we already saw heavy-path, preferred path decompositions. we will also see separator decomposition (L18) and ladder decomposition today.



SA: 9 8 6 4 2 0 7 5 3 1 = leaves of ST

- Internal nodes of ST correspond to an interval of SA
- Time to search on interval  $I$  is  $O(m + \log |I|)$ .

Definition:  $\Sigma$ -node is: (1) has at least  $|\Sigma|$  leaves in its subtree. (2) all its children do not.

# of leaves in subtree of a  $\Sigma$ -node  $\leq |\Sigma|^2$  so time to search in SA interval of a  $\Sigma$ -node is  $O(m + \log |\Sigma|)$

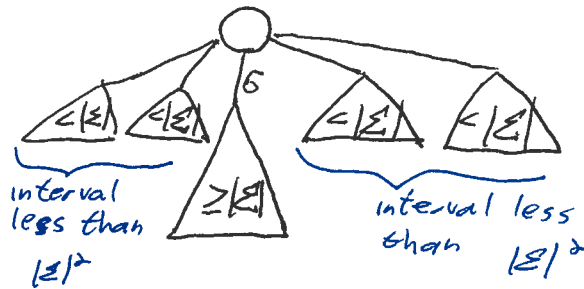
- This takes care of  $\Sigma$ -nodes

Remaining internal nodes partitioned into:

- (1) Branching- $\Sigma$ -node: has at least two children with subtree  $\geq |\Sigma|$
- (2) Others

# Branching- $\Sigma$ -nodes is  $O(n/|\Sigma|)$  so we can store an array of size  $|\Sigma|$  in everyone.

The remaining nodes:



- Navigation:
1.  $\Sigma$ -node: jump to SA.
  2. Branching- $\Sigma$ -node: look at  $\Sigma$ -array.
  3. Others: look at single character, walk to this child or jump to SA.

Search time:  $O(n) + O(m + \log|\Sigma|) = O(m + \log|\Sigma|)$

↑ search SA interval

to walk Branching- $\Sigma$ -nodes and single characters

Space:  $O(n)$

# Approximate String Matching

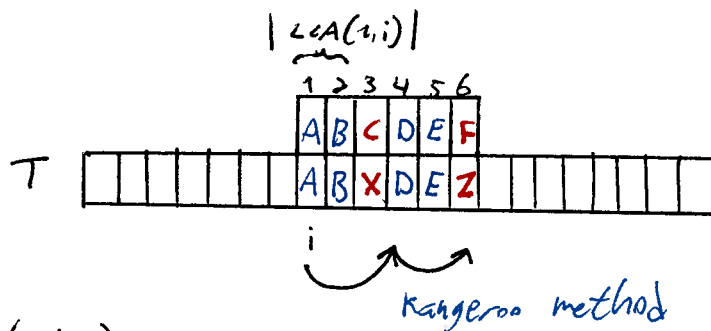
Given a text  $T$ , find occurrences of a query pattern  $p$  in  $T$  with error  $k$

- Hamming Distance: # of character mismatches - (Today)
- Edit distance: # of edits (insertions, deletions, mismatches) needed for an exact match - (Pset)

## "Online"

$P$  &  $T$  are given together.

- construct suffix tree + LCA on  $P\#T \Rightarrow \text{space } O(|T|)$
- For every location in  $T$ 
  - use  $k$  LCA queries to check if  $p$  occurs in this location.



Time:  $O(|T| \cdot k)$

## "Offline"

$T$  is given in advance,  $p$  is given online.

- Best known Bounds: [Cole, Gottlieb, Lewenstein] STOC 2004

Space and preprocessing:  $O(|T| (\lg |T|)^k / k!)$

query time:  $O(|p| + (c \lg |T|)^k / k! \lg \lg |T| + \underbrace{3^k}_{\text{only for edit distance}} \cdot \# \text{ occurrences})$

We focus on subproblem of searching with wildcards

- P contains at most k "don't care" characters (? wildcards)
- Find "exact matches" (? matches anything)

Solution:  $O(|T| \lg^k |T|)$  space and preprocessing

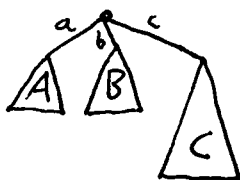
$O(2^k \lg \lg |T| + |P| \cdot \# \text{ occurrences})$  for query

Simple solution  $|Σ|^k \cdot |P|$  query

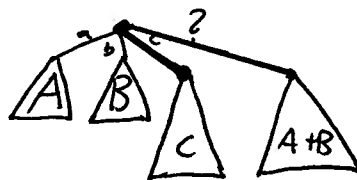
- walk down suffix tree of T as usual
- upon ? in P branch  $|Σ|$  ways

Better solution  $2^k |P|$  query

- Heavy-Light decomposition on suffix tree of T  
primary tree



recursively  
=>



Search branches 2 ways upon ?  
in P:  
1) heavy edge  
2) ? edge

- Each node in primary tree stores secondary tree on the union of light subtrees hanging off node, excluding first character

- recurs k times  $\rightarrow$  k+1 levels of trees
- light depth  $\leq \lg |T|$

$\Rightarrow$  each leaf appears in  $\leq \lg^k |T|$  trees

$\Rightarrow O(|T| \lg^k |T|)$  space & preprocessing

$2^k \lg \lg |T| + |P|$  solution uses "level ancestor" data structure.

as well as  
suffix tree + LCA

Determine in  $\lg \lg |T|$  time  
on which of the two edges we branch

## Level Ancestor Problem:

- Preprocess static tree
- level ancestor query: find  $l^{\text{th}}$  ancestor of node  $v$ .



[Berkman & Vishkin - JCSS 94, Dietz - WADS 97, Alstrup & Holm - ICALP 00  
Bender & Farach-Colton - TCS 04] ↖ dynamic tree

### Solutions:

① lookup table: on all possible queries

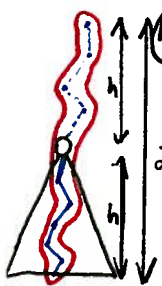
- $O(n^2)$  space
- $O(1)$  query

② Jump pointers: (skiplist) each node stores  $1^{\text{st}}, 2^{\text{nd}}, 4^{\text{th}}, 8^{\text{th}}, \dots$  ancestor

- $O(n \lg n)$  space
- $O(\lg n)$  query

③ long path decomposition: decompose tree by longest paths. store in the same as heavy-path decomp. each path as an array + parent pointer

- $O(n)$  space
- query might still visit  $\Theta(\sqrt{n})$  paths...
- Notice: node of height  $h$  is on path of length  $\geq h$ . So:



④ ladder decomposition: extend path of length  $l$  up by  $l$  levels.

- still  $O(n)$  space

$2h$  - height( $v$ ) =  $h \Rightarrow$  height (top node in  $v$ 's ladder)  $\geq 2h$

- query: as above, but using ladders

- each step at least doubles height  $\Rightarrow O(\lg n)$  query

⑤ combine jump pointers & ladder decomp.

- query:  $O(1)$ . 1 jump pointer  $\Rightarrow$  height( $v$ )  $> l/2$

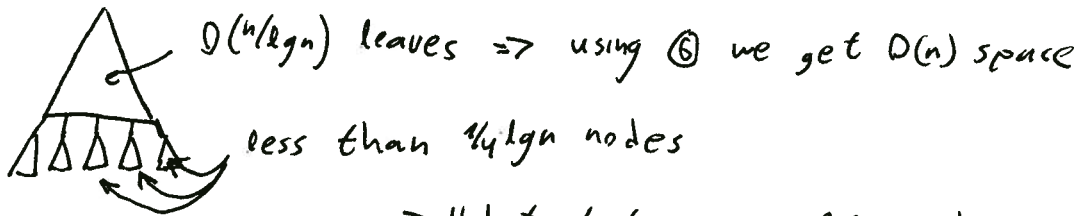
1 ladder step  $\Rightarrow$  ladder( $v$ ) extends  $> l/2$  above

$\Rightarrow$  contains  $l^{\text{th}}$  ancestor

- space -  $O(n \lg n)$

- ⑥ jump-pointer tuning: jump-pointers only at leaves
- start queries at leaves
  - depth- $d$  ancestor of  $v$  = depth- $d$  ancestor of leaf  $(v)$
  - $O(n + L \lg n)$  space, where  $L = \# \text{leaves}$
- $\uparrow$                        $\uparrow$   
 ladder                  jump  
 decomp                pointers

⑦ ART decomposition:



$$\begin{aligned} \Rightarrow \# \text{ distinct trees on } \frac{1}{4} \lg n \text{ nodes} \\ = \underbrace{C_{\frac{1}{4} \lg n}}_{\text{Catalan\#}} \leq 4^{\frac{1}{4} \lg n} = \sqrt{n} \end{aligned}$$

- lookup table for every possible ~~bottom~~ tree  $\Rightarrow O(\sqrt{n} \cdot \lg^2 n)$  space
- query: look in bottom tree, else get parent(root). query top  $\Rightarrow O(1)$  query

- \* In our "don't cares" data structure we need ancestors in a compressed trie  $\Rightarrow$  weights on the edges
  - total weight along any path  $\leq n$