

## 15-853: Algorithms in the Real World

### Suffix Trees

15-853

Page 1

## Exact string searching

Given a text  $T$  of length  $n$  and pattern  $P$  of length  $m$   
"Quickly" find an occurrence (or all occurrences) of  $P$   
in  $T$

A Naïve solution:

Compare  $P$  with  $T[i..i+m]$  for all  $i$  —  $O(nm)$  time

How about  $O(n+m)$  time? (Knuth Morris Pratt)

How about  $O(n)$  preprocessing time and  
 $O(m)$  search time?

15-853

Page 2

Notation:

Capital letters for strings:  $A, B, S, T, \dots$

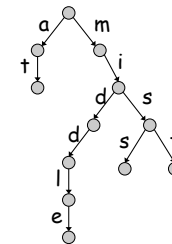
Lower-case letters for characters:  $a, b, c, \dots$

15-853

Page 3

## TRIEs

Dictionary = {at, middle, miss, mist}

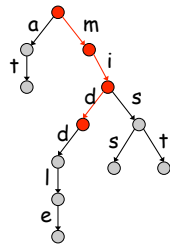


15-853

Page 4

## TRIEs (searching)

Consider searching a string of length  $m$  to see if it is a **prefix** of an element in the dictionary  
 Search time depends on implementation of nodes



e.g., Search("mid")

15-853

Page 5

## TRIEs (node implementation)

Consider searching a string of length  $m$  to see if it is a **prefix** of an element in the dictionary

Consider an  $n$ -node trie over alphabet  $\Sigma$ , with  $|\Sigma|=k$ .

Implementation choices:

- Array per node:  $O(nk)$  space,  $O(m)$  time search
- Tree per node:  $O(n)$  space,  $O(m \log k)$  time search
- Hash children:  $O(n)$  space,  $O(m)$  time

Don't need a separate table per node — can hash node pointer and child character

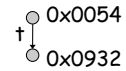


Table.Lookup((0x0054,t)) = 0x0932

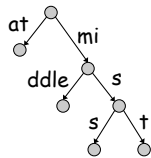
15-853

Page 6

## Compressed Tries

Also called PATRICIA tries/trees or radix trees.  
 All nodes with a single child are collapsed  
 Edges now represent substrings

Dictionary = {at, middle, miss, mist}



Takes less space in practice

15-853

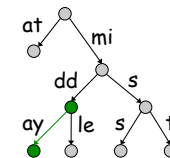
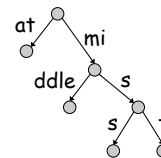
Page 7

## Insertion

Inserting string  $S$  into a compressed trie

- Find longest common prefix
- Split edge if necessary
- Add remaining suffix

Insert("midday")



Takes  $O(|S|)$  time

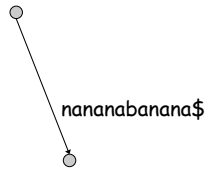
15-853

Page 8



## Suffix Tree Construction

S = nananabana\$  
↑  
i

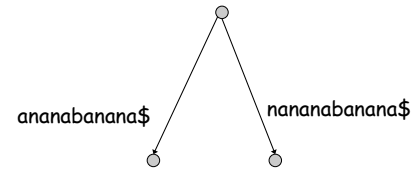


15-853

Page 13

## Suffix Tree Construction

S = nananabana\$  
↑  
i

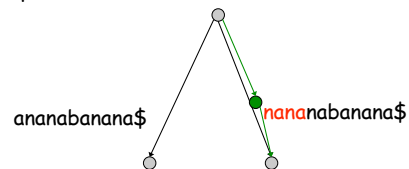


15-853

Page 14

## Suffix Tree Construction

S = nananabana\$  
↑  
i

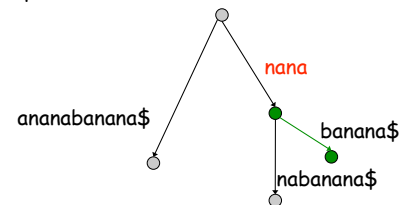


15-853

Page 15

## Suffix Tree Construction

S = nananabana\$  
↑  
i



inserting "nana\$" matches "nana": 4 comparisons  
inserting "ananabana\$" will match "ana": 3 comparisons  
then nabana\$ and abana\$...

15-853

Page 16

## First Improvement

Consider insert of "xAB" (nanabanana\$)

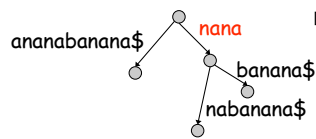
Suppose "xA" (nana) is a prefix in the tree

Then "A" (ana) is a prefix in the tree

⇒ Don't need to match A again — assume it's there.

S = nananabananana\$

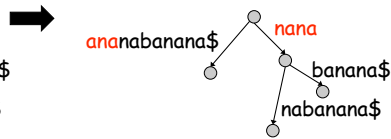
↑ ↑  
i j : end of match



15-853

S = nananabananana\$

↑ ↑  
i j

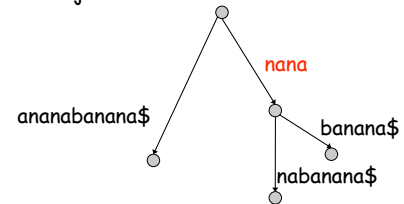


Page 17

## Suffix Tree Construction

S = nananabananana\$

↑ ↑  
i j



inserting "nabanana\$" matches "nana": 4 comparisons

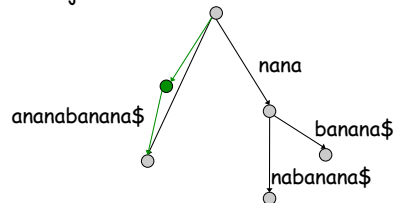
15-853

Page 18

## Suffix Tree Construction

S = nananabananana\$

↑ ↑  
i j



inserting "anabananana\$", "ana" already matched.  
just follow "a..." edge and split : 0 comparisons

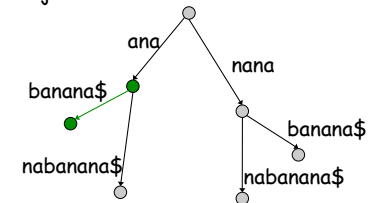
15-853

Page 19

## Suffix Tree Construction

S = nananabananana\$

↑ ↑  
i j



inserting "anabananana\$", "ana" already matched.  
just follow "a..." edge and split : 0 comparisons

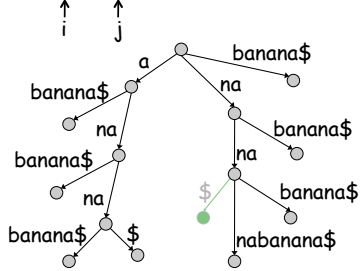
15-853

Page 20



## Suffix Tree Construction

S = nananab**anana**\$



Uh oh. Yes, "nana" is already there, but it takes extra work to follow intermediate nodes

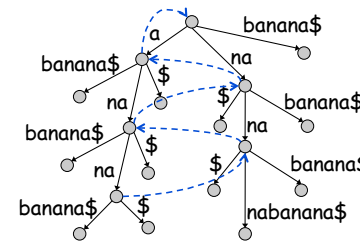
15-853

Page 25

## Second Improvement: Suffix Links

For every internal node "xA", keep a pointer to the node for "A"

S = nananabana\$



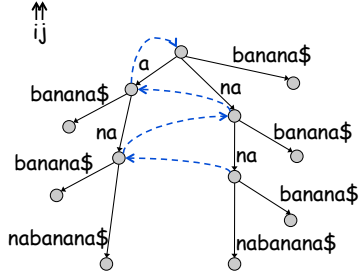
Getting from xA to A becomes cheaper — follow suffix links!

15-853

Page 26

## Suffix Tree Construction

S = nananabana\$

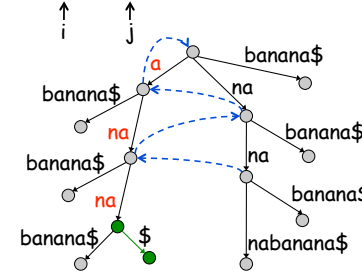


15-853

Page 27

## Suffix Tree Construction

S = nananab**anana**\$



15-853

Page 28







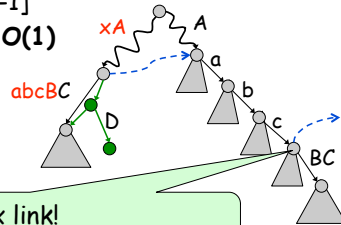
## Following Suffix Links

Suppose  $S[i..j-1]$  is a new node (i.e., resulting from an edge split)

1. Go to parent node  $S[i..k]$ , for some  $i \leq k < j-1$
2. Follow suffix link to  $S[i+1..k]$
3. Search down to  $S[i+1..j-1]$   
**not guaranteed to be  $O(1)$**

$S = ZxAabcBD$

↑  $i$                       ↑  $j$



Page 37

## Better Analysis

$S =$

current ↑                      ↑                      ↑  
suffix  $i$                       nearest  $k$                       matched  $j$   
   suffix link                      prefix

( $k$  not maintained by algorithm)

- $j$ : Each increment of  $j$  takes  $O(1)$  time
- $k$ : If searching from  $S[i+1..k]$  to  $S[i+1..j-1]$  follows  $r$  edges, then  $k$  increases by at least  $r$
- $i$ : Each iteration of  $i$  takes  $O(1)$  additional time
- $\Rightarrow$  total time is  $O(n)$  because  $i$ ,  $j$ , and  $k$  are each incremented at most  $n$  times.

15-853

Page 38

## Extending to multiple lists

Suppose we want to match a pattern with a dictionary of  $m$  texts with a total length of  $n$

Concatenate all texts (interspersed with special characters) and construct a common suffix tree

(Optional) truncate tree below special characters

Time:  $O(n+m) = O(n)$

or

Construct suffix tree on first text, then insert suffixes of second text and so on.

Time: also  $O(n)$

15-853

Page 39

## Longest Common Substring

Find the longest string that is a substring of both  $S_1$  and  $S_2$

Construct a common suffix tree for both texts

All leaves should be labeled with  $S_1$  and/or  $S_2$

Any node with a descendent labeled  $S_1$  and  $S_2$  is a common substring

The "deepest" such node is the longest common substring

Takes linear time with a tree traversal

15-853

Page 40

## Common substrings of m texts

Given m texts of total length n

Given some value  $1 \leq k \leq m$ , find the longest string that is a substring of at least k texts

Construct a common suffix tree, labeling each leaf with the text it comes from

For every internal node, find the number of distinctly labeled descendents :  $O(m)$  time per node.

Can report for all k with a single tree traversal  
time :  $O(mn)$  — not linear