

#### Genome Rearrangements Finding preserved genes

• We have seen before how to rearrange a genome to obtain another one based on:

- Reversals

- Knowledge of preserved blocks (or genes)

 Now we are concerned in determining the preserved genes, or more generally, given two strings x and y, determine all their possible maximal matches



#### A simple way: k-mers

Another possibility is to find common *k*-mers. Here's one way:

Algorithm:

- Denote a k-mer of x by (w, 0, i) where  $w = x_{i}...x_{i+k-1}$
- Denote a k-mer of y by (z, 1, j) where  $z = y_j \dots y_{j+k-1}$
- Sort them lexicographically
- Deduce all k long matches between x and y



#### Disadvantages

• Worst case running time still O(mn)

- e.g. O(m) k-mers in x and O(n) k-mers in y are identical.

- Thought: aren't we supposed to find these anyway and, therefore, the O(mn) bound is not necessarily bad?
- No, they might be small insignificant matches, we are interested in maximal matches
- Making *k* larger reduces the running time because it results in a shortest list of matches, but we might miss significant matches

Saad Mneimneh



- We will use an efficient data structure called suffix tree that stores all suffixes of a string and supports fast lookup
- Definition: A suffix tree *T* of a string *s* of length *m* is a rooted tree such that:
  - It has exactly m leaves numbered 1 to m
  - Each internal node other than the root has at least two children
  - Each edge is labeled by a substring of s
  - No two edge labels out of a node start with the same character
  - For any leaf i the concatenation of the edge labels on the path from the root to i spells out the suffix  $x_{\!\!p}..x_{\!\!m}$

Saad Mr











A suffix tree satisfies the following:

- |E| = |V| − 1 (tree)
- Number of leaves = m + 1 (now |s| = m + 1)
- Since each internal node has at least two children, the number of edges |*E*| = *O*(# leaves) = *O*(*m*)
- Any sub-tree with k leaves satisfies |E| = O(k)



## Building a suffix tree

Here's a simple algorithm:

given  $x = x_1 \dots x_m$ 

insert a special character \$ at the end of s initialize the tree  ${\cal T}$  to one root

**for** j = 1 to m + 1

find the longest match of  $\mathbf{x}_j...\mathbf{x}_m$  in T starting from the root and following a unique path

split the edge where the match stops, add a new node w

add an edge  $(w_{i})$  (j is the new leaf) and label it with the remaining unmatched characters of  $x_{j}\ldots x_{m}$ 



### Analysis

- Running time: O(m<sup>2</sup>)
   Each suffix requires O(m) time to update the tree
- But, there exists an O(m) time suffix tree algorithm
- Space: O(m)
  - How? Each label has O(m) characters and we have |E| = O(m) labels!
  - Solution: do not explicitly store labels, but store the indices [*i,j*] of a label

Saad Mn

Saad Mr

#### Now what?

- How can we use the suffix tree data structure to identify all maximal matches between two strings *x* and *y*?
- Consider first the following problem: given a string *x*, determine all locations where another string *y* occurs.
- This can be solved efficiently as described next.

## string matching

Find all occurrences of y in x

Algorithm	
build a suffix tree <i>T</i> for <i>x</i>	O(n
Match the characters of y along the unique path in <i>T</i> until (case 1) either y is exhausted or (case 2) no more matches are possible	O(r
if (case 2) y does not occur in x	O(1
else the <i>k</i> leaves in the sub-tree below the point of the last match give the <i>k</i> location of <i>y</i> in <i>x</i> (traverse the tree in linear time)	0(
	Saad Mneimn

#### Correctness

- Why is the string matching algorithm correct?
- If *y* occurs in *x* at position *i*, then the *i*<sup>th</sup> suffix of *x* must start with *y*
- Therefore, leaf *i* must be reached by the path determined by *y*



Saad Mn

#### Finding maximal matches

 Given x and y, we would like to find all maximal matches between x and y

- $x_{j} \dots x_{i+1} = y_j \dots y_{j+1}$
- Cannot extend  $x_{i}...x_{i+1}$  and  $y_{j}...y_{j+1}$  and obtain a match
- We will find all matches starting at y<sub>j</sub> that cannot be extended to the right
  - Build a suffix tree for x (do this only once)
  - Find the path in *T* determined by the longest possible prefix of the suffix y<sub>j</sub>...y<sub>n</sub> (it could stop in the middle of an edge e in that case e is part of the path)
  - Let  $v_k, k = 1...p$  be an internal node on this path and  $T_k$  be the sub-tree rooted at  $v_k$  that excludes  $v_{k+1}$
  - Identify the leaves in each sub-tree



#### What about left?

- Given a leaf *i*, let *left*(*i*) be the character *x*<sub>*i*-1</sub>
- If  $left(i) \neq y_{j-1}$ , then *i* represents a maximal match
- Therefore, we obtain all maximal matches between *x* and *y* in *O*(*mn*) time by repeating the previous algorithm for every suffix of *y*





# Generalized suffix tree for a set of strings

- We can build a suffix tree for a set of strings s<sub>1</sub>, s<sub>2</sub>, ..., s<sub>n</sub>
  - Append a different end of string marker to each string in the set
  - concatenate all the strings together
  - build a suffix tree for the concatenated string
- The resulting suffix tree will have a leaf for each suffix of the concatenated string and is build in time proportional to the sum of all lengths
- The leaf numbers can be easily converted to two numbers, one identifying a string *s*<sub>i</sub> and the other a starting position in *s*<sub>i</sub>









#### Suffix tree for x and y

- Therefore, given two strings *x* and *y*, we can build a suffix tree for both in O(*m* + *n*) (i.e. linear) time.
- Each leaf in the tree represents
  - Either a suffix from x
  - Or a suffix from y
- Mark each internal node v with x (y) if there is a leaf in the sub-tree
  of v representing a suffix from x (y). This can be done in linear time
  by a bottom up traversal of the tree from leaves to the root.
- Note that if v is marked x (y), all ancestors of v are marked x (y).







#### Left diverse node

An internal node v is *left diverse* iff it has two children v<sub>1</sub> and v<sub>2</sub> with a leaf *i* for x in v<sub>1</sub>'s sub-tree and a leaf *j* for y in v<sub>2</sub>'s sub-tree, such that *left*(*j*)  $\neq$  *left*(*j*)

(assume  $x_0$  and  $y_0$  are different and distinct from any other character)

If  $u\alpha p$  is a substring of x and  $w\alpha q$  is a substring of y for  $u \neq w$  and  $p \neq q$ , then  $\alpha$  corresponds to a left diverse node v and vice-versa.

Proof: similar to previous proof

Call such an  $\boldsymbol{\alpha}$  a maximal common substring



#### Compact representation

- Therefore, we have only O(m + n) maximal common substrings for x and y (but each maximal common substring might appear in multiple locations)
- If we identify left diverse nodes in linear time, we need only O(m + n) time and space to come up with this compact representation of all maximal common substrings
- A maximal match can be represented as (p<sub>1</sub>, p<sub>2</sub>, l) where p<sub>1</sub> and p<sub>2</sub> are the positions of a maximal common substring of length l in x and y respectively
- We can obtain all maximal matches in O(m + n + k) where k is their number (we will not present the algorithm)



