## Suffix Tree and Suffix Array

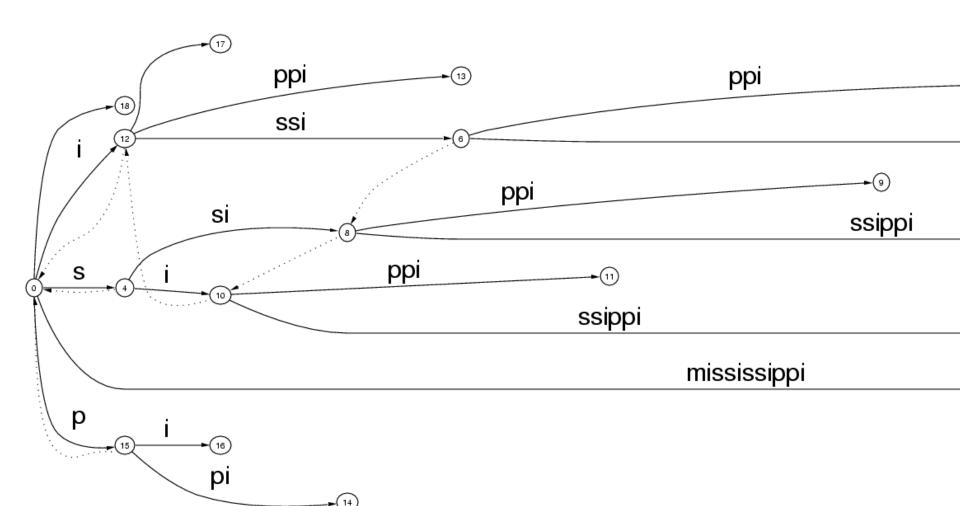
R92922025	Brain Chen
R92548028	Pluto Chang

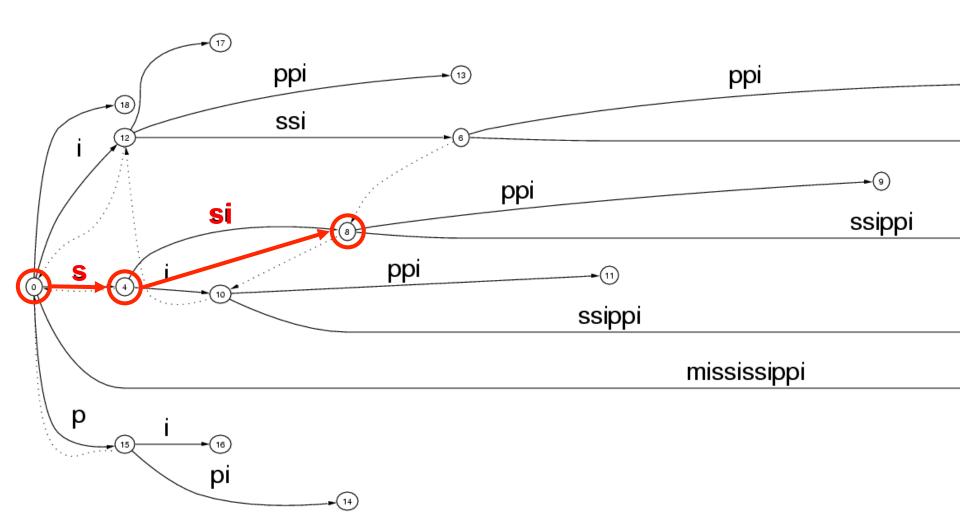
# Outline

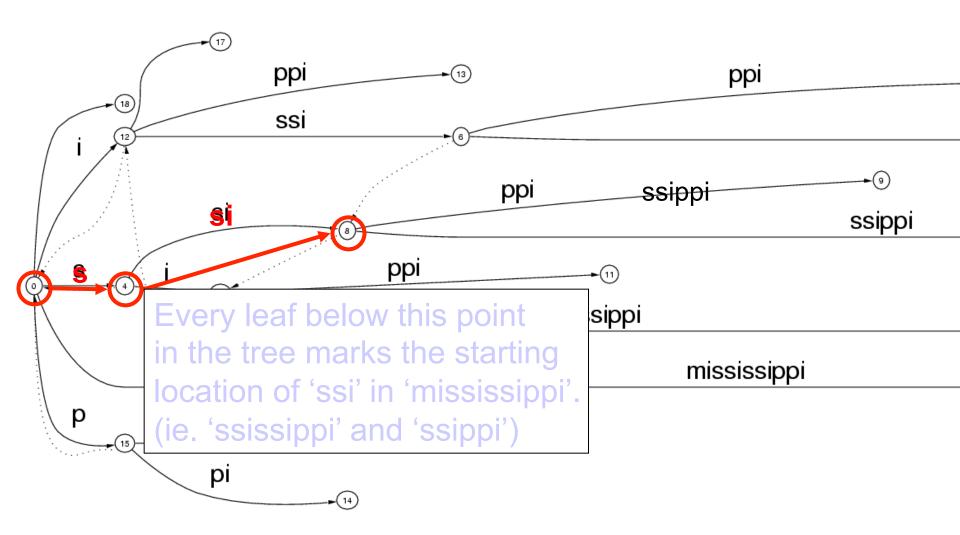
- Motivation
- Exact Matching Problem
- Suffix Tree
  - Building issues
- Suffix Array
  - Build
  - Search
  - Longest common prefixes
- Extra topics discussion
- Suffix Tree VS. Suffix Array

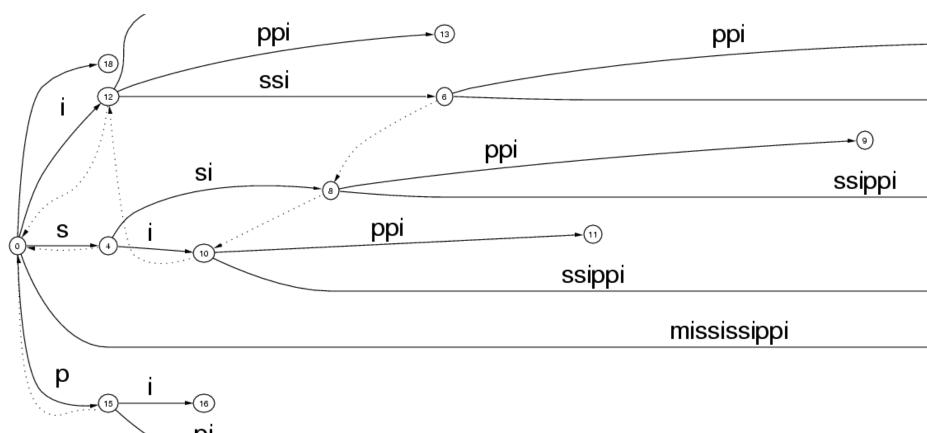
#### Motivation

- Text search
  - Need fast searching algorithm(with low space cost)
- •DNA sequences and protein sequences are too large to search by traditional algorithms
- Some improved algorithms perform efficiently
  - •KMP, BM algorithms for string matching
  - Suffix Tree with linear construction and searching time
  - Suffix Array with Suffix Tree based construction

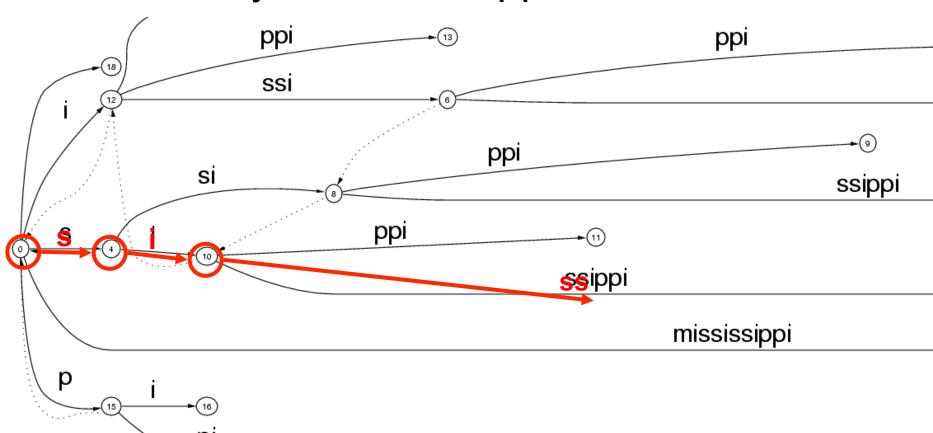




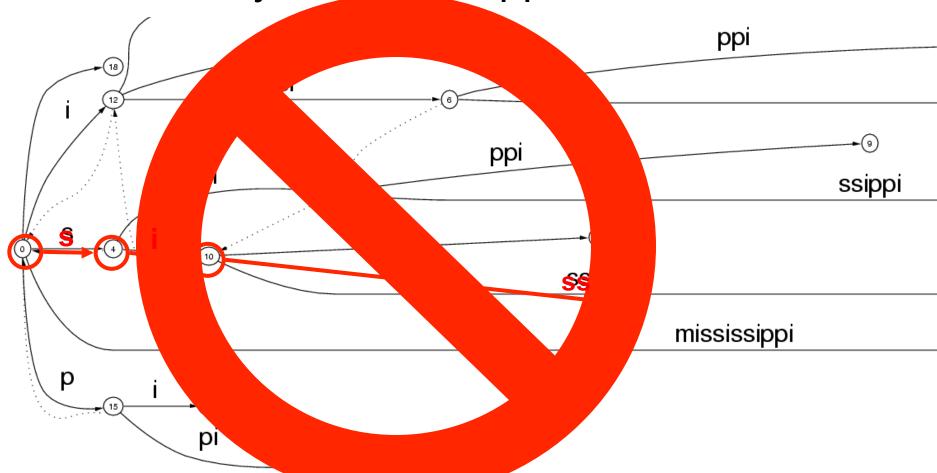












- So what? Knuth-Morris-Pratt and Boyer-Moore both achieve this worst case bound.
  - *O*(*m*+*n*) when the text and pattern are presented together.
- Suffix trees are <u>much</u> faster when the text is fixed and known first while the patterns vary.
  - $\bigcirc$  O(m) for single time processing the text, then only O(n) for each new pattern.
- Aho-Corasick is faster for searching a number of patterns at one time against a single text.

#### Boyer-Moore Algorithm

- For string matching(exact matching problem)
- Time complexity O(m+n) for worst case and O(n/m) for absense
- Method: backward matching with 2 jumping arrays(bad character table and good suffix table)

#### What are suffix arrays and trees?

- Text indexing data structures
- not word based
- allow search for patterns or
- computation of statistics

#### **Important Properties**

- Size
- Speed of exact matching
- Space required for construction
- Time required for construction



# Suffix Tree

#### Properties of a Suffix Tree

- Each tree edge is labeled by a substring of S.
- Each internal node has at least 2 children.
- Each  $S_{(i)}$  has its corresponding labeled path from root to a leaf, for  $1 \le i \le n$ .
- There are n leaves.
- No edges branching out from the same internal node can start with the same character.

How do we build a suffix tree?

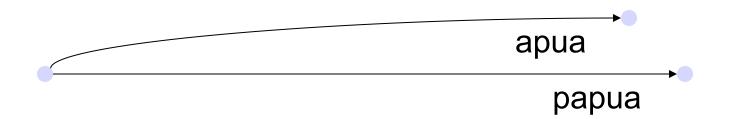
while suffixes remain:

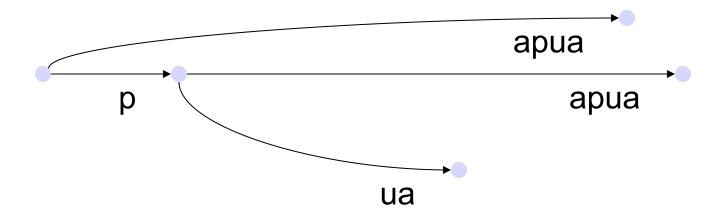
add next shortest suffix to the tree

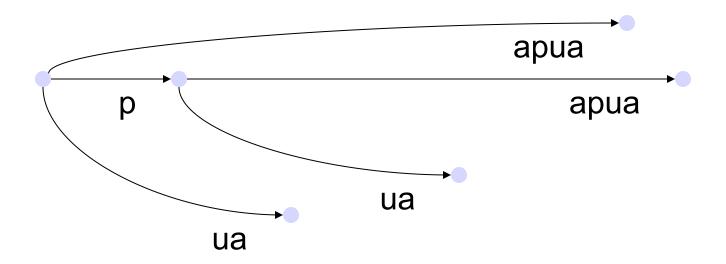


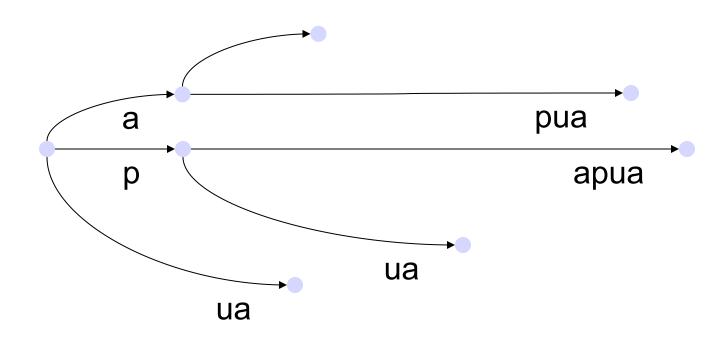
papua

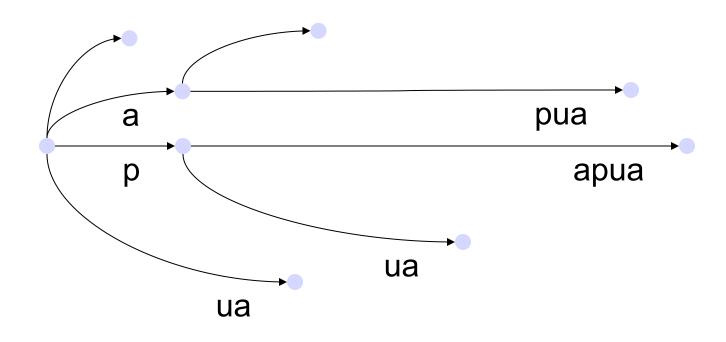












How do we build a suffix tree?

while suffixes remain:

add next shortest suffix to the tree

Naïve method -  $O(m^2)$  (m = text size)

#### Building the Suffix Tree in O(m) Time

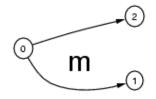
- In the previous example, we assumed that the tree can be built in O(m) time.
- Weiner showed original O(m) algorithm (Knuth is claimed to have called it "the algorithm of 1973")
- More space efficient algorithm by McCreight in 1976
- Simpler 'on-line' algorithm by Ukkonen in 1995

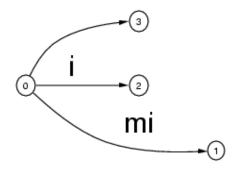
#### Ukkonen's Algorithm

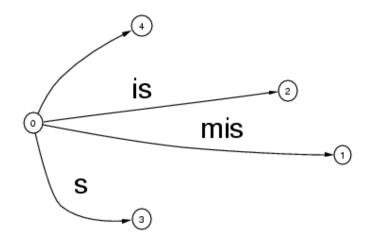
- Build suffix tree T for string S[1..m]
  - OBuild the tree in m phases, one for each character. At the end of phase i, we will have tree  $T_i$ , which is the tree representing the prefix S[1..i].
    - In each phase i, we have i extensions, one for each character in the current prefix. At the end of extension j, we will have ensured that S[j..i] is in the tree T<sub>i</sub>.

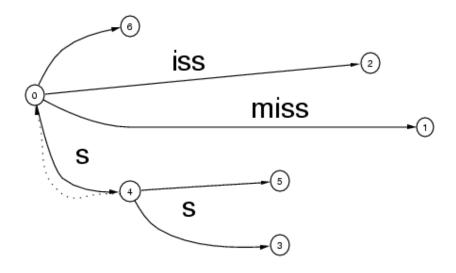
#### Ukkonen's Algorithm

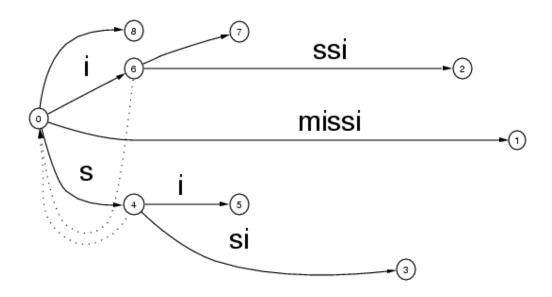
- 3 possible ways to extend S[j..i] with character i +1.
- 1. S[j..i] ends at a leaf. Add the character i+1 to the end of the leaf edge.
- 2. There is a path through *S[j..i]*, but no match for the *i*+1 character. Split the edge and create a new node if necessary, then add a new leaf with character *i*+1.
- 3. There is already a path through *S[j..i+1]*. Do nothing.

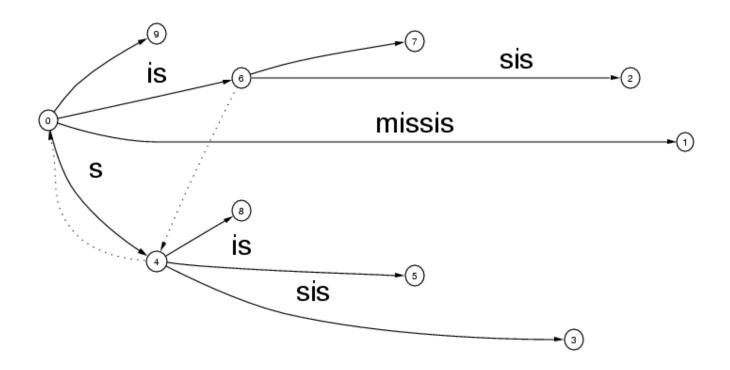


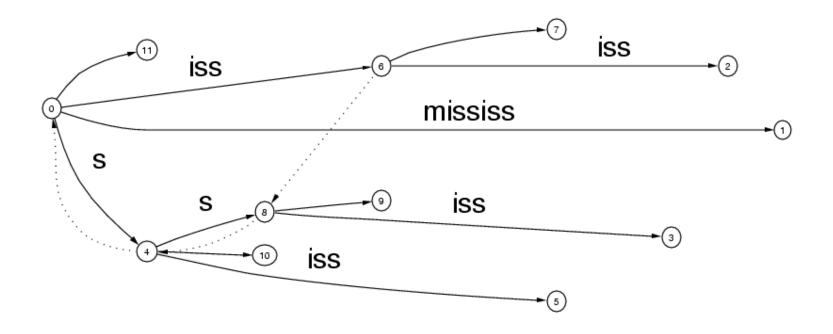


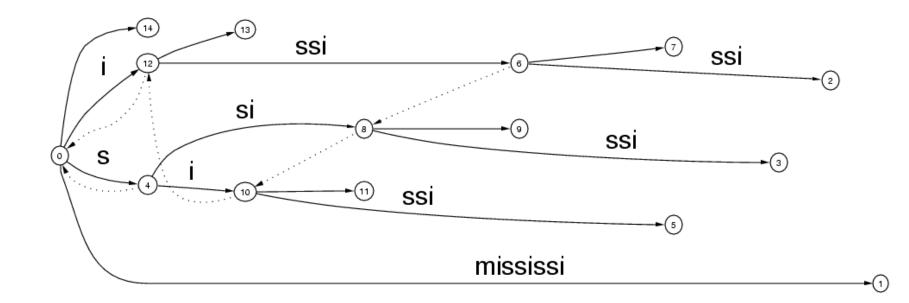


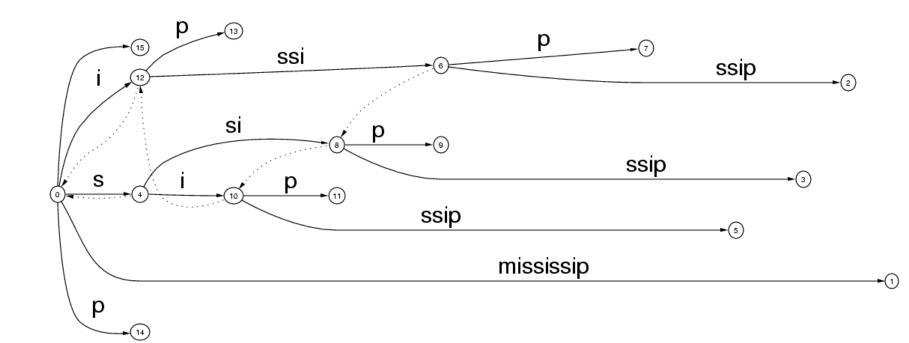


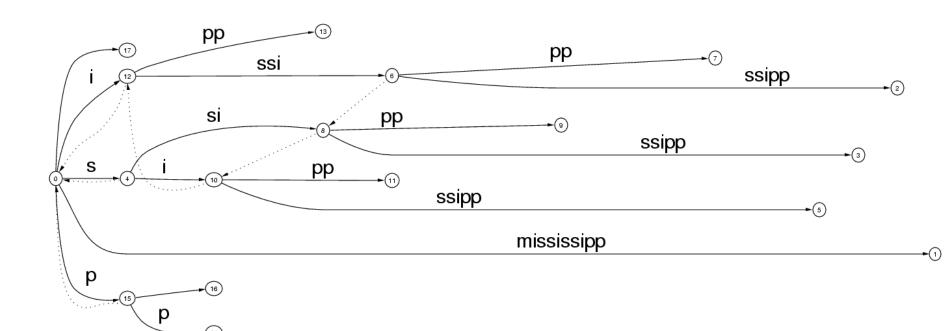


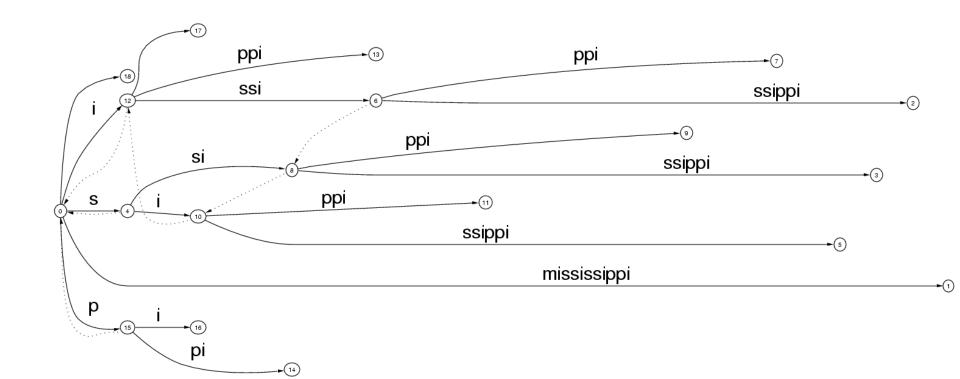












# Ukkonen's Algorithm

- In the form just presented, this is an  $O(m^3)$  time,  $O(m^2)$  space algorithm.
- We need a few implementation speed-ups to achieve the O(m) time and O(m) space bounds.



# Suffix Array

# The Suffix Array

**Definition:** Given a string **D** the <u>suffix</u> array **SA** for this string is the sorted list of pointers to all suffixes of **D**.

(Manber, Myers 1990)

# The Suffix Array

http://par.cse.nsysu.edu.tw/~cbyang/

- In a suffix array, all suffixes of S are in the non-decreasing lexical order.
- For example, S="ATCACATCATCA"

i	0	1	2	3	4	5	6	7	8	9	10	11
A	11	3	8	0	5	10	2	7	4	9	1	6

3	ATCACATCATCA	S <sub>(0)</sub>
10	TCACATCATCA	S <sub>(1)</sub>
6	CACATCATCA	S <sub>(2)</sub>
1	ACATCATCA	S <sub>(3)</sub>
8	CATCATCA	S <sub>(4)</sub>
4	ATCATCA	S <sub>(5)</sub>
11	TCATCA	S <sub>(6)</sub>
7	CATCA	S <sub>(7)</sub>
2	ATCA	S <sub>(8)</sub>
9	TCA	S <sub>(9)</sub>
5	CA	S <sub>(10)</sub>
0	А	S <sub>(11)</sub>

0	А	S <sub>(11)</sub>
1	ACATCATCA	S <sub>(3)</sub>
2	ATCA	S <sub>(8)</sub>
3	ATCACATCATCA	S <sub>(0)</sub>
4	ATCATCA	S <sub>(5)</sub>
5	CA	S <sub>(10)</sub>
6	CACATCATCA	S <sub>(2)</sub>
7	CATCA	S <sub>(7)</sub>
8	CATCATCA	S <sub>(4)</sub>
9	TCA	S <sub>(9)</sub>
10	TCACATCATCA	$S_{(1)}$
11	TCATCA	S <sub>(6)</sub>



fin

#### How do we build it?

- Build a suffix tree
- Traverse the tree in DFS, lexicographically picking edges outgoing from each node and fill the suffix array.
- O(n) time
- Suffix tree construction loses some of the advantage that the suffix array has over the suffix tree

#### Direct suffix array construction algorithm

Unfortunately, it is difficult to solve this problem with the suffix array Pos alone because Pos has lost the information on tree topology. In direct algorithm, the array Height (saving lcp information) has the information on the tree topology which is lost in the suffix array P

"Linear-Time Longest-Common-Prefix Computation in Suffix Arrays and Its Applications"

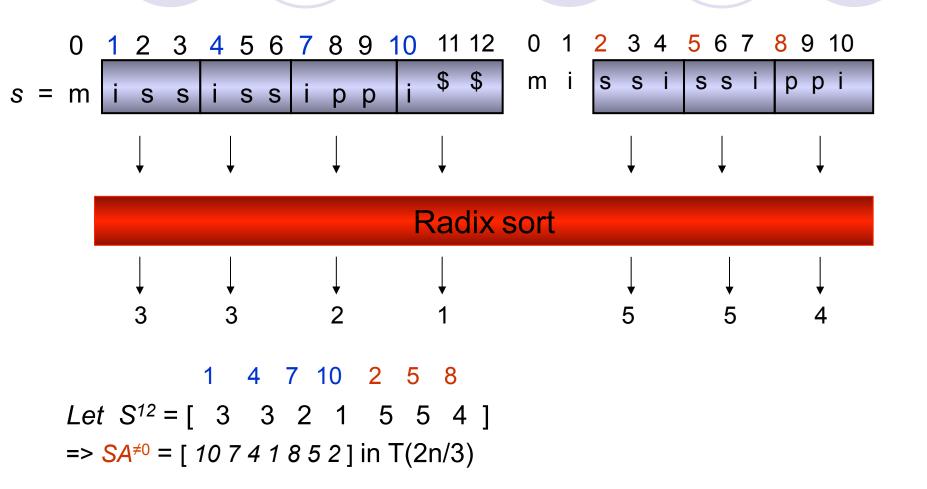
# Skew-algorithm

- Step 1:
  - $SA^{\neq 0}$  = sort the suffixes starting at position  $i \neq 0$  mod 3.
- Step 2:
  - $SA^{=0}$  = sort the suffixes starting at position i = 0 mod 3.
- Step 3:

 $SA = \text{merge } SA^{=0} \text{ and } SA^{\neq 0}.$ 

```
0 1 2 3 4 5 6 7 8 9 10
s = m i s s i s s i p p i
```

Step 1:  $SA^{\neq 0}$  = sort the suffixes starting at position  $i \neq 0 \mod 3$ .



#### 1 4 7 10 2 5 8

$$s^{12} = [ 3 \ 3 \ 2 \ 1 \ 5 \ 5 \ 4 ]$$
 $s^{12}_{1} = 3 \ 3 \ 2 \ 1 \ 5 \ 5 \ 4$ 
 $s^{12}_{4} = 3 \ 2 \ 1 \ 5 \ 5 \ 4$ 
 $s^{12}_{7} = 2 \ 1 \ 5 \ 5 \ 4$ 
 $s^{12}_{10} = 3 \ 5 \ 5 \ 4$ 
 $s^{12}_{10} = 5 \ 5 \ 4$ 
 $s^{12}_{2} = 5 \ 5 \ 4$ 
 $s^{12}_{5} = 5 \ 4$ 
 $s^{12}_{8} = 5 \ 4$ 

$$s = m i s s i s s i p p i$$
 $s_1 = i s s i s s i p p i$ 
 $s_4 = i s s i p p i$ 
 $s_7 = i p p i$ 
 $s_{10} = i$ 
 $s_2 = i s s i s s i p p i$ 
 $s_5 = i s s i p p i$ 
 $s_8 = i p p i$ 

$$SA^{\neq 0} = [10741852],$$

It suffices to show that  $S^{12}_{i} < S^{12}_{j} <=> s_{i} < s_{j}$ .

### Compare $S_i$ and $S_j$ where i = 0, $j \neq 0$ mod 3:

- case 1:  $j = 1 \mod 3$ 
  - :  $i + 1 = 1 \mod 3$ ,  $j+1 = 2 \mod 3$
  - ... compare (s[i],  $S_{i+1}$ ) with (s[j],  $S_{j+1}$ ) s[i] = the character at i in s, s, in constant time.
- case 2:  $j = 2 \mod 3$ 
  - :  $i + 2 = 2 \mod 3$ ,  $j+2 = 1 \mod 3$
  - ... compare (s[i], s[i+1],  $S_{i+2}$ ) with (s[j], s[j+1],  $S_{i+2}$ ) in constant time

$$S^{12}_{i} < S^{12}_{j} <=> S_{i} < S_{j}$$

Case 1:  $i = j \mod 3$ 

1 4 7 10 2 5 8 
$$s^{12} = [3 \ 3 \ 2 \ 1 \ 5 \ 5 \ 4]$$

$$s^{12}_{4} = \begin{bmatrix} 4 & 7 & 10 & 2 & 5 & 8 \\ 3 & 2 & 1 & 5 & 5 & 4 \end{bmatrix}$$

$$s^{12}_{4} < s^{12}_{1}$$

$$s_4 = \begin{bmatrix} 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ i & s & s & i & p & p & i & $ & $ & $ \end{bmatrix}$$

$$s_1 = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ i & s & s & i & s & s & i & p & p & i & $ & $ & $ \end{bmatrix}$$

$$s_4 < s_1$$

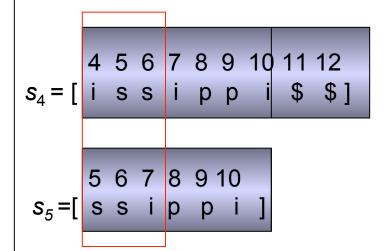
$$S^{12}_{i} < S^{12}_{j} <=> s_{i} < s_{j}$$

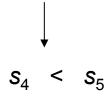
Case 2:  $i \neq j \mod 3$ 

Ex:  

$$s^{12}_{4} = \begin{bmatrix} 4 & 7 & 10 & 2 & 5 & 8 \\ 3 & 2 & 1 & 5 & 5 & 4 \end{bmatrix}$$
  
 $s^{12}_{5} = \begin{bmatrix} 5 & 8 \\ 4 \end{bmatrix}$ 

$$s^{12}_{4} < s^{12}_{5}$$





Step 2:  $SA^{=0}$  = sort the suffixes starting at position i = 0 mod 3.

- The rank of s<sub>j</sub> among {s<sub>k</sub> | k ≠ 0 mod 3 } was determined in Step1 for all j ≠ 0 mod 3.
- $SA^{=0} = \text{radix sort } \{ (s[i], S_{i+1}) \mid i = 0 \mod 3 \}.$

```
0 1 2 3 4 5 6 7 8 9 10
                 s = m i s s i s s i p p i
    (s[i], S_{i+1})
0: (m, ississippi)
                                 9: (p, i)
                                                                  0: (m, ississippi)
                                                   Radix sort
                     Step 1
3: (s, issippi)
                                 6: (s, ippi)
                                                                  9: (p, i)
6: (s, ippi)
                                 3: (s, issippi)
                                                                  6: (s, ippi)
9: (p, i)
                                 0: (m, ississippi)
                                                                  3: (s, issippi)
```

Step 3:  $SA = \text{merge } SA^{=0} \text{ and } SA^{\neq 0}$ .

- $SA^{=0} = [s_0 \ s_9 \ s_6 \ s_3]$
- $SA^{\neq 0} = [s_{10} \ s_7 \ s_4 \ s_1 \ s_8 \ s_5 \ s_2]$
- $SA = \text{merge } SA^{=0} \text{ and } SA^{\neq 0}$ =  $[s_{10} s_7 s_4 s_1 s_0 s_9 s_8 s_6 s_3 s_5 s_2]$ 
  - = [107410986352]

It is in time O(n) if we can determine the relative order of  $S_i \in SA^{=0}$  and  $S_j \in SA^{\neq 0}$  in constant time.

#### Time complexity analysis

- Step1: O(n) + T(2n/3)
- Step2: O(n)
- Step3: O(n)
- T(n) = O(n) + T(2n/3) = O(n)

# Exact matching using a Suffix Array ABAABBAABBAC

#### **SUFFIX ARRAY SA:**

SA = 203691584710

Basic Idea: 2 binary searches in SA Search for leftmost position Search for rightmost position

#### ABAABBABBAC

2 0 3 6 9 1 5 8 4 7 10 0 1 2 3 4 5 6 7 8 9 10

#### ABAABBABBAC

Continue binary search in the right (larger) half of SA

#### ABAABBAAC

More occurences of BB left of this one possible!

#### ABAABBABBAC

leftmost position of BB is pointed to by SA[8]

#### ABAABBABBAC

More occurences of BB right of this one possible!

#### ABAABBABBAC

rightmost position of BB is pointed to by SA[9]

#### Results of search for: BB

#### ABAABBAABC

2 0 3 6 9 1 5 8 4 7 10 0 1 2 3 4 5 6 7 8 9 10

leftmost position of BB is pointed to by SA[8] rightmost position of BB is pointed to by SA[9]

=>All occurences of the pattern BB are pointed to by SA[8..9]

# **Important Properties**

for |SA| = n and m = length of pattern:

- Size : 1 Pointer per Letter (4 Byte if n < 4Gb)
- Speed of exact matching:
  - O(log n) binary search steps
  - # of compared chars is O(mlogn)
     can be reduced to O(m + log n)

### Longest common prefixes

- Definition: Icp(i,j) is the length of the longest common prefix of the suffixes beginning at SA[i] and SA[j].
- Mississippi Example
  - $\bigcirc$ SA[2] = 4 (issippi)
  - $\bigcirc$ SA[3] = 1 (ississippi)
  - Olcp(2, 3) = 4

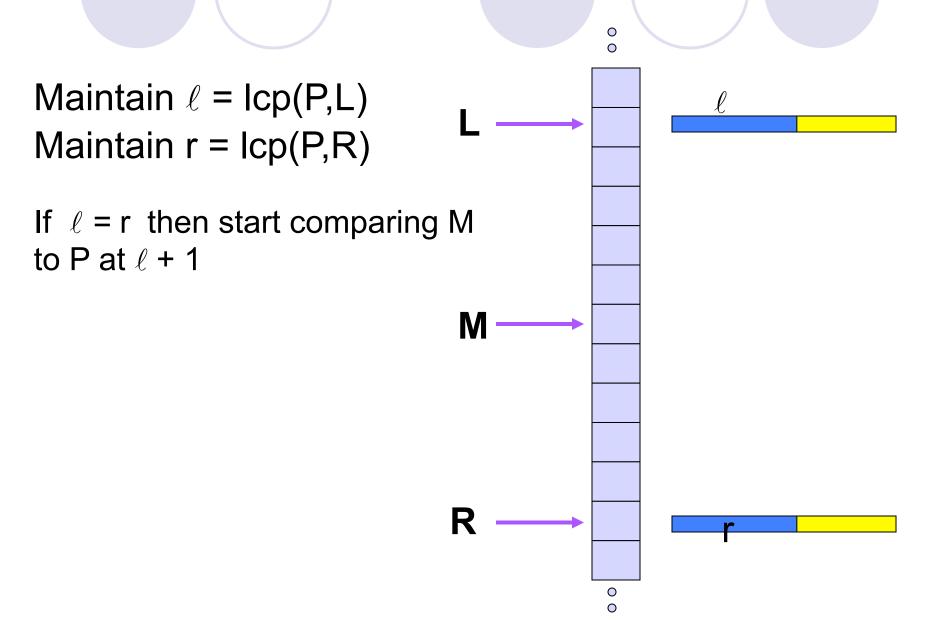
```
s = m i s s i s s i p p i
SA = [10 7 4 1 0 9 8 6 3 5 2]
```

# Example

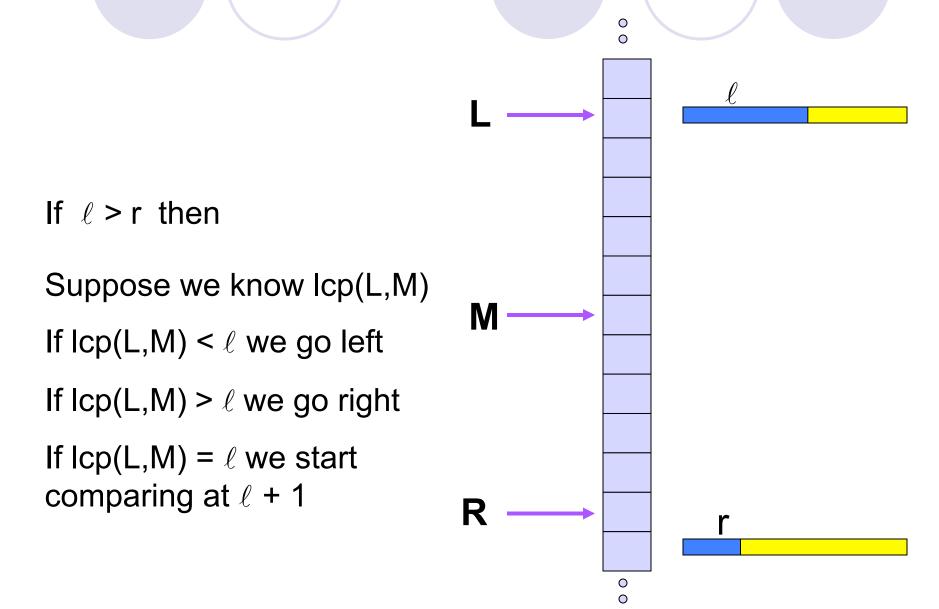
Let S = mississippi

10 ippi issippi Let P = issaississippi mississippi p M ppi 8 sippi 6 sisippi 3 ssippi ssissippi

#### How do we accelerate the search?



#### How do we accelerate the search?



# Analysis of the acceleration

If we do more than a single comparison in an iteration then  $\max(\ell, r)$  grows by 1 for each comparison  $\rightarrow$  O(logn + m) time

# Complicated Sorting Algorithm

- Using radix sort for each characters, totally O(N²)
- Using radix sort for each H characters, and for 2H, 4H, 8H etc. →O(NlogN)

#### Precomputed LCP Array Construction

- Compute Icps between suffixes that are consecutive in the sorted Pos array:
- Range Minimum Query Theorem:
- Given H-bucket lcps, compute 2H-bucket lcps
- still require too much time

### Precomputed LCP Array Construction

- Using height(i) = Icp(A<sub>Pos[i-1]</sub>, A<sub>Pos[i]</sub>)
- Using Hgt[i] to record height(i) when it is correct
- For b-th iteration
  - Oif height(i) ≤ (b-1)H and height(i) < bH, then
    Hgt[i] = height(i)</pre>
  - Otherwise, Hgt[i] = N+1 (undefined)

### Precomputed LCP Array Construction

- Constructing interval tree
  - O(N)-space height balanced tree structure that records the minimum pairwise lcp over a collection of intervals of the suffix array
- Compute min(  $Hgt[k] : k \leftarrow [i, j]$  )
- Takes O(log N) time
- overall O(NlogN) time

### Linear Time Expected-case Variations

- Require additional O(N) structure
- Longest Repeated Substring
  - $\bigcirc$ 2log<sub>| $\Sigma$ |</sub>N+O(1)
  - Sorting algorithm => O(N log log N)
- Linear Time Algorithm
  - Perform RadixSort on T-symbols of each suffix
  - Improve both sorting algorithm and Icp computation

### Constant Time Icp Construction

- LCP[i] = lcp(SA[i], SA[i+1])
- $Lcp(i, j) = min_{i <= k < j} LCP[k]$
- Case 1:
  - $\bigcirc$ j mod 3 = 1, k mod 3 = 2 => adjacent
  - $\bigcirc$ j' = (j-1)/3, k' = (n+k-2)/3 => adjacent
  - $OI = Icp^{12}(j', k') = LCP^{12}[SA^{12}[j']-1]$
  - $\bigcirc$ LCP[i] = lcp(j, k) = 3l + lcp(j+3l, k+3l) <= 2
  - Constant time

# Constant Time Icp Construction

#### Case 2:

- $\bigcirc$  J mod 3 = 0, k mod 3 = 1 (or k mod 3 = 2)
- $\bigcirc$  If s[j]  $\neq$ s[k], LCP[i] = 0
- Otherwise, LCP[i] = 1 + <u>lcp(j+1, k+1)</u> ← Case 1
- Olcp(j+1, k+1) = 3I + lcp(j+1+3I, k+1+3I), if SA[j+1], SA[k+1] are adjacent
- If not adjacent, perform range minimum query
- No suffix is involved in more that two lcp queries at the top level of the extended skew algorithm
- Constant time

- LCP[i] = lcp(SA[i], SA[i+1])
- j = SA[i], k = SA[i+1]
- Case 1:
  - $\bigcirc$ j mod 3 = 1, k mod 3 = 2
  - $\bigcirc$ j' = (j-1)/3, k' = (n+k-2)/3 => adjacent in SA<sup>12</sup>
  - $\mathcal{L} = \text{Icp}^{12}(j', k') = \text{LCP}^{12}[SA^{12}[j']]$
  - $\bigcirc$ LCP[i] = lcp(j, k) = 3 $\ell$  + lcp(j+3 $\ell$ , k+3 $\ell$ ) <= 2
  - Constant time

```
0 1 2 3 4 5 6 7 8 9 0
m i s s i s s i p p i
s^{12} = [ 3 3 2 1 5 5 4 ]
SA^{12} = [ 3 2 1 0 6 5 4 ]
LCP^{12} = [ 0 0 1 0 0 1
0 ]
```

 LCP<sup>12</sup> is used to decide triple-lcps (groups of lcps of 3 characters)

- To answer range minimum queries on LCP<sup>12</sup> needs O(n) time
- Lemma: No suffix is involved in more than two lcp queries at the top level of the extended skew algorithm
  - A suffix can be involved in lcp queries only with its two lexicographically nearest neighbors that have the same preceding character

- LCP<sup>12</sup> construction algorithm
  - LCP<sup>12</sup> array is divided into blocks of size log(n)
  - For each block [a, b], precompute and store the following data:
    - For all i ← [a, b], Qi identifies all j ← [a, i] such that LCP<sup>12</sup>[j] < min<sub>k ←[i+1, i]</sub> LCP<sup>12</sup>[k]
    - For all i ← [a, b], the minimum values over the ranges [a, i] and [i, b]
    - The minimum for all ranges that end just before or begin just after
       [a, b] and contain exactly a power of two full blocks
  - [i, j] is completely inside a block
    - Its minimum can be found with the help of Qj in constant time
  - [i, j] is covered with some ranges whose minimun is stored
    - Its minimum is the smallest of those minima

- LCP[i] =  $lcp(j, k) = 3\ell + lcp(j+3\ell, k+3\ell) <= 2$
- l represents the number of triple-lcps
- 3l represents the number of characters of lcp triples
- The rest is non-triple lcps, which have length at most 2
- Applying character comparison, they can be done in constant time (at most 2 comparisons)
- Computing LCP[i] is O(1) for case 1

- Case 2:
  - $\bigcirc$  J mod 3 = 0, k mod 3 = 1
  - $\bigcirc$  If s[j]  $\neq$ s[k], LCP[i] = 0
  - Otherwise, LCP[i] = 1 + lcp(j+1, k+1) ← Case 1
  - Olcp(j+1, k+1) = 3I + lcp(j+1+3I, k+1+3I), if SA[j+1], SA[k+1] are adjacent
  - If not adjacent, perform range minimum query
  - No suffix is involved in more that two lcp queries at the top level of the extended skew algorithm
  - Constant time

#### **Applications of Suffix Trees and Suffix Arrays**

- Exact String Match
- The Exact Set Matching Problem
  - The problem of finding all occurrences from a set of strings P in a text T, where the set is input all at once.
- The Substring Problem for a Database of Patterns
  - A set of strings, or a database, is first known and fixed. Later sequence of strings will be presented and for each presented string S, the algorithm must find all the strings in the database containing S as a substring.

### **Applications of Suffix Trees and Suffix Arrays**

- Longest Common Substring of Two Strings
- Recognizing DNA Contamination
- Common Substrings of More Than Two Strings
- Building a Smaller Directed Graph for Exact Matching
  - how to compress a suffix tree into a directed acyclic graph(DAG) that can be used to solve the exact matching problem (and others) in linear time but that uses <u>less space</u> than the tree.

#### **Applications of Suffix Trees and Suffix Arrays**

- A Reverse Role for Suffix Trees, and Major Space Reduction
  - Operine *ms(i)* to be the <u>length</u> of the longest substring of T starting at position i that matches a substring somewhere (but we don't know where) in P. These values are called the *matching statistics*.
- Space-Efficient Longest Common Substring Algorithm
- All-Pairs Suffix-Prefix Matching
  - O Given two string Si and Sj, and suffix of Si that matches a prefix of Sj is called a *suffix-prefix match* of Si,Sj.

# Suffix Trees and Suffix Arrays

#### Suffix

- Each position in the text is considered as a text suffix.
  - A string that does from that text position to the end to the text

### Advantage

• They answer efficiently more complex queries.

#### Drawback

- Costly construction process
- The text must be readily available at query time
- The results are not delivered in text position order.

NLP Laboratory of Hanshin University http://infocom.chonan.ac.kr/~limhs/

### Compression

- Suffix trees can be compressed almost to size of suffix arrays
- Suffix arrays can't be compressed (almost random), but can be constructed over compressed text
  - instead of Huffman, use a code that respects alphabetic order
  - almost the same compression
- Signature files are sparse, so can be compressed
  - ratios up to 70%

# Compression

- Suffix trees and suffix arrays
  - Suffix arrays are very hard to compress further.
    - Because they represent an almost perfectly random permutation of the pointers to the text.
  - Suffix arrays on compressed text
    - The main advantage is that both index construction and querying almost double their performance.
      - Construction is faster because more compressed text fits in the same memory space and therefore fewer text blocks are needed.
      - Searching is faster because a large part of the search time is spent in disk seek operations over the text area to compare suffixes.

### Where have suffix trees been used?

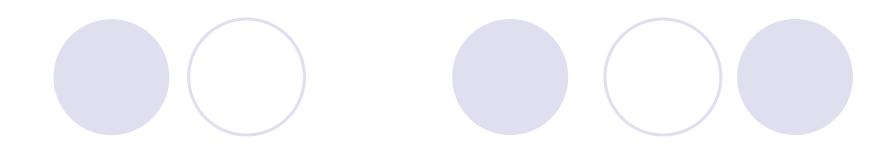
- Problems
  - Olinear-time longest common substring
  - Oconstant-time least common ancestor
  - omaximally repetitive structures
  - all-pairs suffix-prefix matching
  - compression
  - Oinexact matching
  - conversion to suffix arrays

### Where have suffix trees / arrays been used?

- Applications
  - The Human Genome Project (see Skiena)
  - Omotif discovery (see Arabidopsis genome project)
  - OPST probabilistic suffix trees
  - SVM string kernels
  - chromosome-level similarities and rearrangements

### When have suffix trees / arrays been used?

- When they solve your problem.
- When you need results fast!
- When you have memory to spare.
- ...more caveats.



fin

