

# Space-Efficient Construction Algorithm for Circular Suffix Tree

Wing-Kai Hon, Tsung-Han Ku, Rahul Shah and Sharma Thankachan

---

CPM2013

# Outline

---

- \* Preliminaries and Motivation
- \* Circular Suffix Tree
- \* Our Indexes and the Construction Algorithms

# Preliminaries

---

- \* Let  $P$  be a pattern of length  $|P|$ .
- \*  $P^\infty$  denote the string of infinite number of times of  $P$
- \*  $Q = P[i..|P|]P[1..i-1]$  be the  $i$ -th cyclic shift of  $P$ .

Ex:

Let  $P = abcd \Rightarrow P^\infty = abcdabcdabcdabcd\dots\dots$

Cyclic shifts of  $P$ :

$abcd, bcda, cdab, dabc$

# Circular Dictionary Matching

---

Given a set  $D = \{P_1, \dots, P_d\}$  of  $d$  patterns, the **circular dictionary matching** problem is to index  $D$  such that for any online query text  $T$ , we can quickly locate the occurrences of any pattern of  $D$  and all of its shifts within  $T$ .

In 2009, Hon et al. proposed a solution for the circular dictionary matching problem by using succinct circular suffix tree.

# Ordering

Let  $P$  and  $Q$  be two patterns. We say  $P$  is circularly larger than  $Q$ , or simply  $P$  larger than  $Q$ , if and only if  $P^\infty$  is lexicographically larger than  $Q^\infty$ . The notions of “smaller than”, or “equal to”, are defined analogously.

Example:

$P = abaab, Q = ab \rightarrow P < Q$

# Preliminaries

---

Lemma 1:

Let  $P$  and  $Q$  be two strings, and let  $m$  denote the maximum of  $|P|$  and  $|Q|$ . Then,  $P$  is circular smaller than  $Q$  if and only if  $P^\infty[1..2m]$  is lexicographically smaller than  $Q^\infty[1..2m]$ . In other words, to compare the “circular” lexicographical order of  $P$  and  $Q$ , we only need to compare  $P^\infty$  and  $Q^\infty$  directly up to length  $2m$ .

# Circular Suffix Tree

- \* The circular suffix tree  $ST_{circ}$  for  $D$  is a compact trie storing all circular suffixes of all  $P^\infty_j$ . In addition, the children of a node are arranged according to the lexicographical order of the labels in the incident edges.
- \* Consequently, the  $i$ -th leftmost leaf in the circular suffix tree will correspond to the  $i$ -th lexicographically smallest circular suffix, which is indexed by  $SA_{circ}[i]$ .

# Our Index

---

## Succinct Circular Suffix Tree:

- (1) The circular suffix array: for representing the leaf labels and the edge labels.
- (2) An *Hgt* array: for representing the length of the edge labels
- (3) A string of balanced parentheses: for encoding the tree structure.



# Circular Suffix Array

The circular suffix array  $SA_{circ}[1..n]$  for  $D$  is an arrangement of all circular suffixes of all  $P_j^\infty$  according to the lexicographical order.

Lemma 2: (Hon et al. 2012)

The circular suffix array can be stored succinctly in  $O(n \log \sigma)$  bits, and can be constructed in  $O(n \log n)$  time using  $O(n \log \sigma)$  bits working space. The functions  $SA_{circ}$  and  $SA_{circ}^{-1}$  can both be evaluated in  $O(\log n)$  time.

# *Hgt* Array

In  $Hgt[1..|P|]$  array,  $Hgt[i]$  stores the length of the longest common prefix between suffix  $S_i$  and its predecessor suffix  $S_{i-1}$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	-1
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	0
4	<i>bcdabc</i>	2	2
5	<i>c</i>	7	0
6	<i>cdabc</i>	3	1
7	<i>dabc</i>	4	0

An example of *Hgt* array for  $P = abcdabc$

# Kasai et al.'s algorithm

---

- \* Given the suffix array and inverse suffix array of a string  $S$ .
- \* Kasai et al.'s algorithm can compute the  $Hgt$  array of  $S$  in linear time.

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = \text{abcdabc}$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	
4	<i>bcdabc</i>	2	
5	<i>c</i>	7	
6	<i>cdabc</i>	3	
7	<i>dabc</i>	4	

$$h = 3$$

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = a**bc**dabc$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	
2	<i>abcdabc</i>	1	3
<b>3</b>	<b><i>bc</i></b>	6	
<b>4</b>	<b><i>bcdabc</i></b>	2	2
5	<i>c</i>	7	
6	<i>cdabc</i>	3	
7	<i>dabc</i>	4	

$$h = 2$$

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = abc**dabc**$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	
4	<i>bcdabc</i>	2	2
<b>5</b>	<b><i>c</i></b>	7	
<b>6</b>	<b><i>cdabc</i></b>	3	1
7	<i>dabc</i>	4	

$$h = 1$$

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = abc\dabc$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	
4	<i>bcdabc</i>	2	2
5	<i>c</i>	7	
6	<i>cdabc</i>	3	1
7	<i>dabc</i>	4	0

$$h = 0$$

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = abcdabc$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	-1
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	
4	<i>bcdabc</i>	2	2
5	<i>c</i>	7	
6	<i>cdabc</i>	3	1
7	<i>dabc</i>	4	0

$$h = 0$$



# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = abcda**bc**$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	-1
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	0
4	<i>bcdabc</i>	2	2
5	<i>c</i>	7	
6	<i>cdabc</i>	3	1
7	<i>dabc</i>	4	0

$$h = 0$$

# Kasai et al.'s algorithm

Compute the *Hgt* array of  $P = abcda**bc**$

Index	Suffix	SA	<i>Hgt</i>
1	<i>abc</i>	5	-1
2	<i>abcdabc</i>	1	3
3	<i>bc</i>	6	0
4	<i>bcdabc</i>	2	2
5	<i>c</i>	7	0
6	<i>cdabc</i>	3	1
7	<i>dabc</i>	4	0

$$h = 0$$

# *Hgt* Array of Circular Suffixes

- \* The *Hgt* array can naturally be defined for the circular suffix tree for  $D$ .
- \* We can also extend Kasai et al.'s algorithm to construct the *Hgt* array.
- \* Nevertheless, when we process the circular suffixes of  $P_j^\infty$  in round  $j$ , the total time required will become  $O(|P_j| + L_j)$ , where  $L_j$  denotes the maximum length of LCP between  $P_j$  and its predecessor circular suffix.
- \* In worst case,  $L_j = O(n)$ , the total running time will become  $O(dn)$

# *Hgt* Array of Circular Suffixes

Index	Suffix	$SA_{circ}$	<i>Hgt</i>
1	$(ab)^\infty$	(2, 2)	-1
2	$(abbb)^\infty$	(1, 4)	2
3	$(ba)^\infty$	(2, 1)	0
4	$(babb)^\infty$	(1, 3)	3
5	$(bbab)^\infty$	(1, 2)	1
6	$(bbba)^\infty$	(1, 1)	2
7	$(b)^\infty$	(3, 1)	3

An example of *Hgt* array for  $D = \{bbba, ba, b\}$

# *Hgt* Array of Circular Suffixes

---

- \* *Hgt'* stores the LCP between  $S_i$  and its successor  $S_{i+1}$ .
- $Hgt[i] = Hgt'[i-1]$
  
- \* In our algorithm, in each round, we compute  $Hgt[i]$  and  $Hgt'[i]$  such that when we compute the LCP between  $S_j$  and its predecessor, if  $Hgt'[j-1]$  is computed, we can set  $Hgt[j] = Hgt'[j-1]$  immediately.

# *Hgt* Array of Circular Suffixes

---

- \* Up to now, we are actually assuming the  $SA[i]$  and  $SA^{-1}[i]$  can be computed in constant time.
- \* Therefore, our algorithm for computing *Hgt* array is linear.
- \* Since we adopt Hon et al.'s succinct circular suffix array, the time complexity of our algorithm is  $O(n \log n)$ .

# Hgt Array of Circular Suffixes

- \* Define  $H[j, k] = \text{Hgt}[i]$ , if the  $k$ -th circular suffix of  $P_j^\infty$  is  $S_i$ .
- (1)  $\text{Hgt}[i] = H[j, k] = H[\text{SA}[i]]$
- (2)  $H[j, 1] \leq H[j, 2]+1 \leq H[j, 3]+2 \leq \dots$
- \* Since  $P_j^\infty$  is circular, its  $(|P_j|+1)$ -th circular suffix is the same to itself.
- $H[j, 1] \leq H[j, 2]+1 \leq \dots \leq H[j, |P_j|]+|P_j|-1 \leq H[j, |P_j|+1]+|P_j|$   
 $= H[j, 1]+|P_j|$

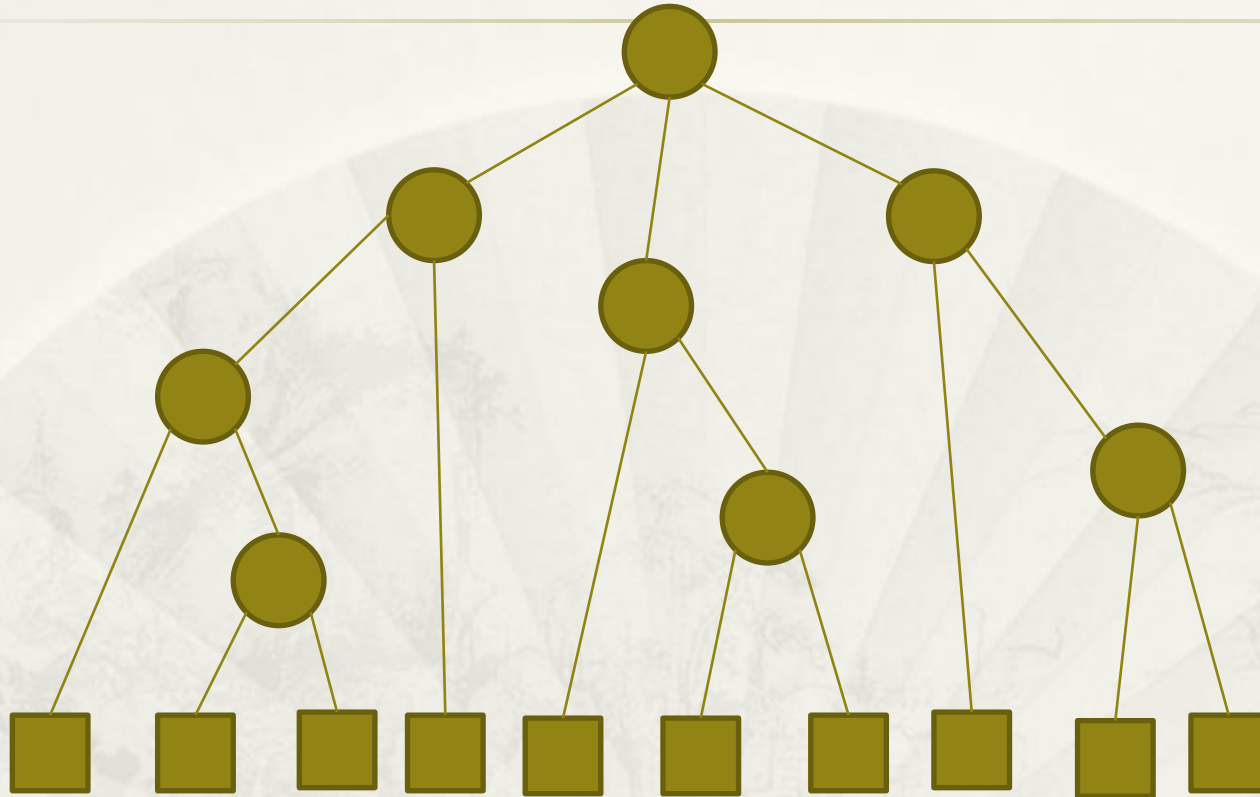
# *Hgt* Array of Circular Suffixes

---

- \* In stead of storing *Hgt* array, we store *H* array succinctly (i.e. store the differences).
- \* Thus, the space complexity of the *H* array is  $O(d \log n + n)$  bits.
- \* Therefore, the *Hgt* array can be constructed in  $O(n \log n)$  time using  $O(n \log \sigma + d \log n)$  bits of working space.

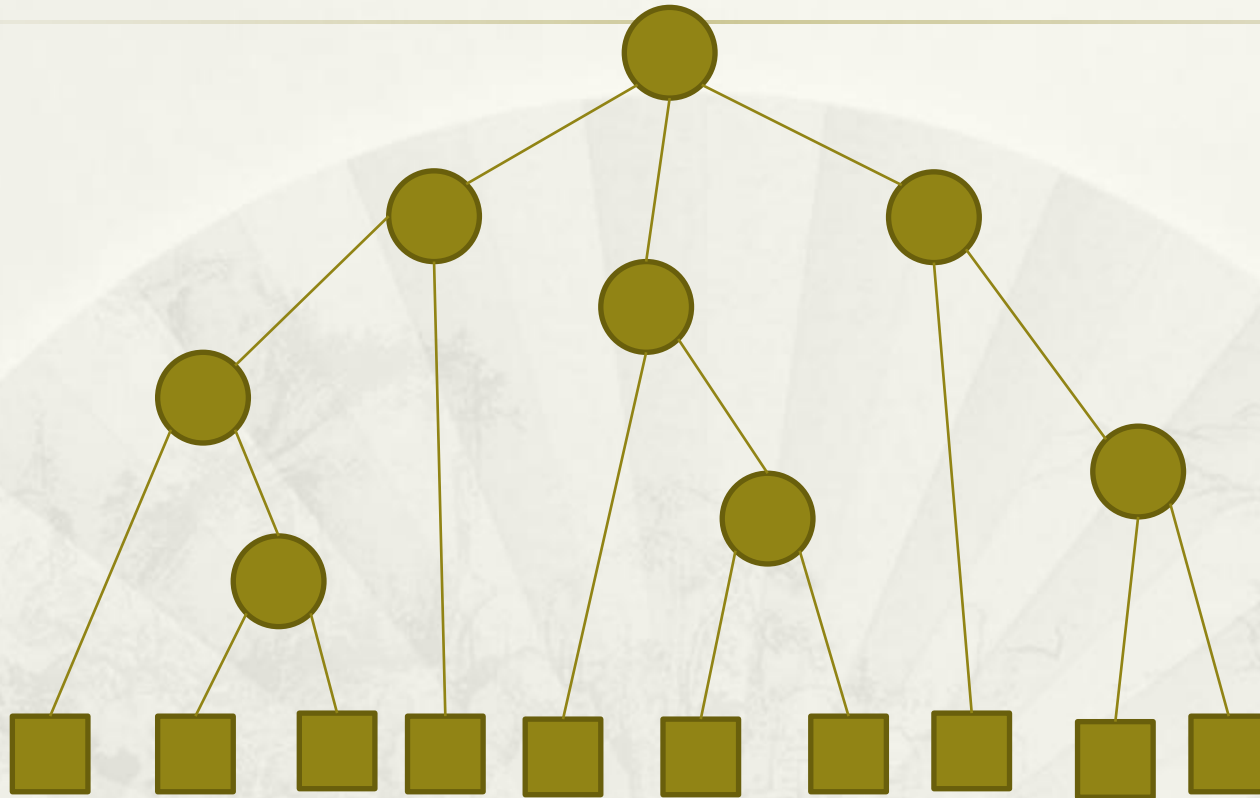


# Balanced parentheses representation



The output of the left to right bottom up traversal would be  
( ) ( ) ( ) ) ) ( ) ) ( ) ( ) ( ) ) ) ( ) ( ) ( ) ) ) )

# Balanced parentheses representation



The output of the right to left bottom up traversal would be  
 ( ( ( ( ) ( ( ) ( ) ( ) ( ( ) ( ( ) ( ) ( ( ) ( ( ) ( )

# Balanced parentheses representation

The output of the **left to right** bottom up traversal would be

( ) ( ) ( ) ) ) ( ) ) ( ) ( ) ( ) ) ) ( ) ( ) ( ) ) ) )

The output of the **right to left** bottom up traversal would be

( ( ( ( ) ( ( ) ( ) ( ) ( ( ) ( ( ) ( ) ( ( ) ( ( ) ( )

By using a linear time scanning algorithm, we can compute the balanced parentheses representation of the circular suffix tree.

(( ( ( ( ) ( ( ) ( ) ) ) ( ) ) ( ( ) ( ( ) ( ) ) ) ( ( ) ( ( ) ( ) ) ) )

Therefore, we can compute the balanced parentheses representation in  $O(n \log n)$  time by using  $O(n)$  bits working space.

---



Thank you