

CSCI3220 2018-19 First Term Assignment 2

I declare that the assignment here submitted is original except for source material explicitly acknowledged, and that the same or closely related material has not been previously submitted for another course. I also acknowledge that I am aware of University policy and regulations on honesty in academic work, and of the disciplinary guidelines and procedures applicable to breaches of such policy and regulations, as contained in the following websites.

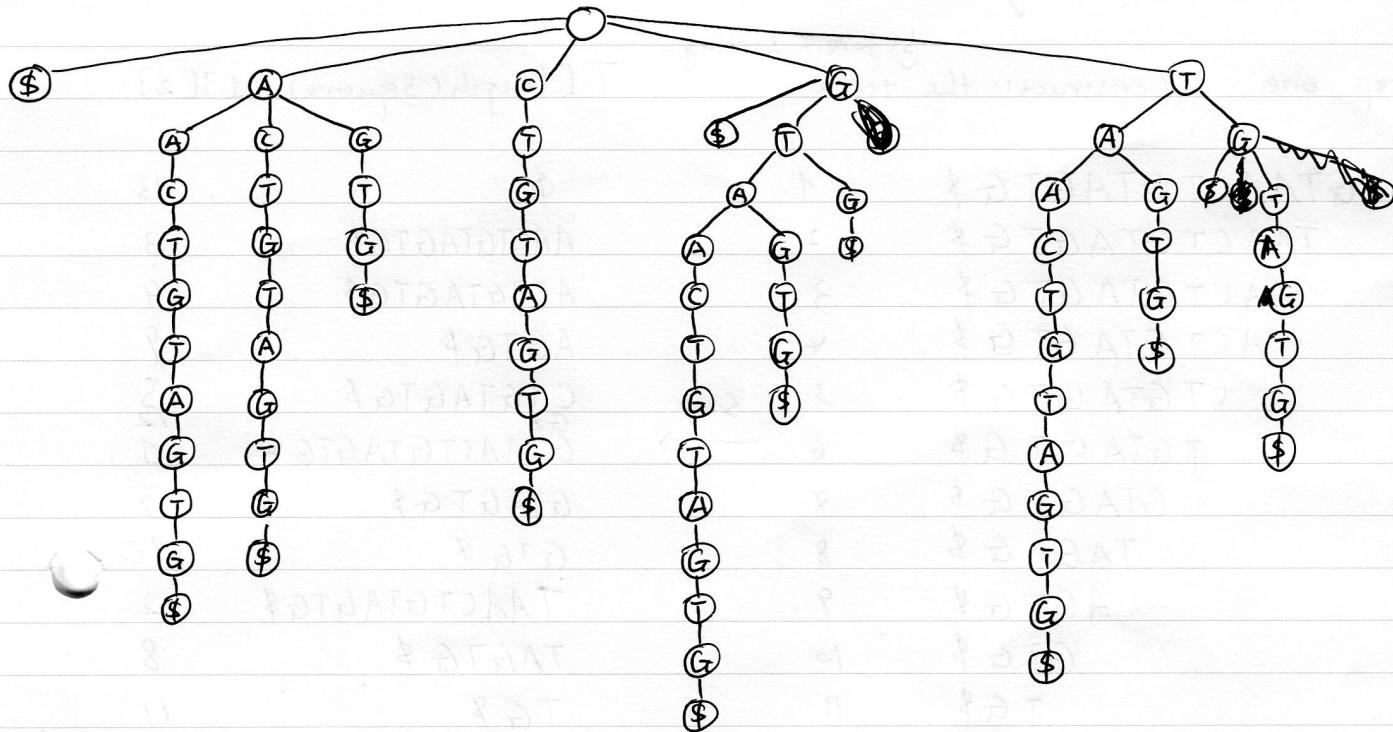
University Guideline on Academic Honesty:  
<http://www.cuhk.edu.hk/policy/academichonesty/>

Student Name: Zuowen Wang  
Student ID : 1155123906

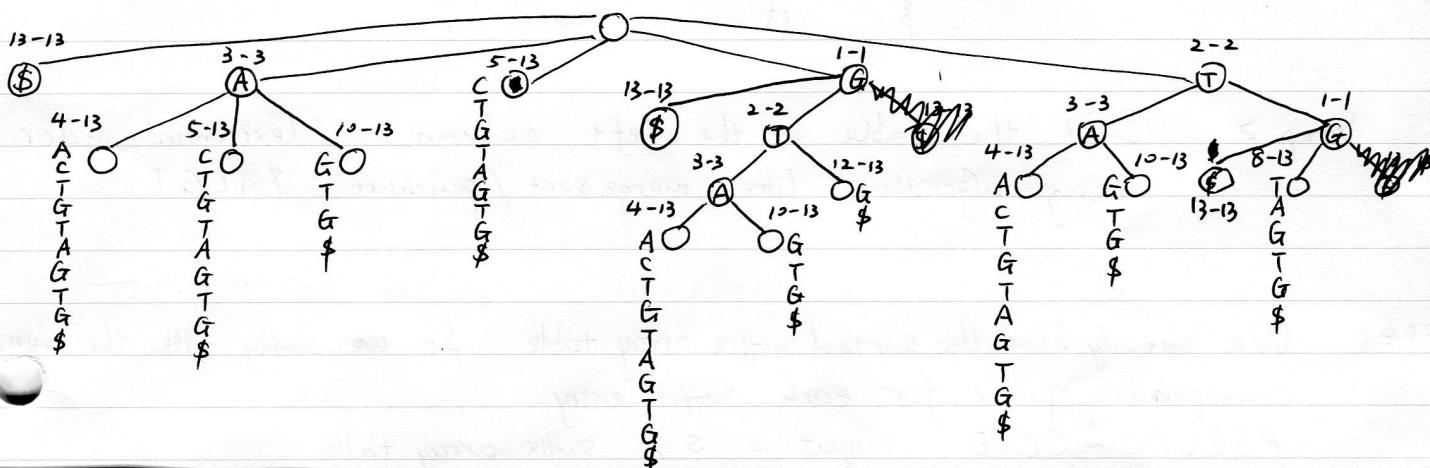
1. (a)

S = GTAACTGTAGTG \$

P1



(b)



(C) We traverse the suffix tree (compressed) with index using DFS. When we encounter a leaf node, we output (its lower index - depth of its parent) to a dynamic array (for convenience). Notice the root has depth 0.

pseudo code: Input: Tree G, Output: Suffix Array SA.

Main: call DFS(G, root, -1)

DFS:  $\text{DFS}(G, v, \text{depth})$  // depth is depth of  $v$ 's parent  
 $\text{depth}++;$   
 $\text{foreach child } u \text{ of } v \text{ in lexicographic order}$   
 $\text{if } u \text{ is a leaf.}$   
 $\text{SA.append}(u.\text{index}.low - \text{depth})$   
 $\text{else}$   
 $\text{DFS}(G, u, \text{depth})$

Using this algorithm we can get suffix array from suffix tree

SA: 13, 4, 4, 9, 5, 12, 1, 7, 10, 2, 8, 11, 6

(d) Input :  $s = \text{GTAACGTAGTG\$}$

Goal: Output : The Suffix Array

A

Step 1. Create an array of String. The length of array is  $s.length$ . And store all suffixes (from longest to the \\$) in this array.

Step 2. Sort this array A using lexicographic order \$ACGT

Better use in place merge sort since suffix array method was created for memory limitation. or heap sort.

Step 3 Output the corresponding location for each entry of A.

The location can be computed as  $(s.length - A[i].length + 1)$

A

$\#A'$

$$= 13 - 1 + 1 = 13$$

GTAACGTAGTG \$	\$	13
TAACGTAGTG \$	AACTGTAGTG \$	3
AACTGTAGTG \$	ACTGTAGTG \$	4
ACTGTAGTG \$	AGTG \$	9
CTGTAGTG \$	CTGTAGTG \$	Step 3 5
TGTAGTG \$	G \$	Step 3 Compute Location 12
GTAGTG \$	GTAACGTAGTG \$	1
TA GTG \$	GTAGTG \$	7
AGTG \$	GTG \$	10
GTG \$	TAACGTAGTG \$	2
TG \$	TAGTG \$	8
G \$	TG \$	11
\$	TGTAGTG \$	6

(e) Assume we have suffix array Sa. Now we want to get the corresponding BWT b of s.

Step 1. (actually there is only one step...)

We can construct b directly from sa and s using the following pseudo code: initially  $b = ""$  // empty string.

for i from 1 to  $s.length$

$b = b + \boxed{\text{_____}} s.charAt(\underbrace{sa[i] - 1})$

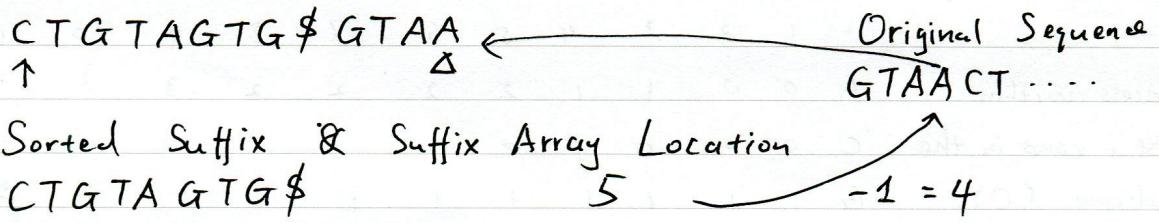
if this == 0, then return \$!

Output b.

Our b is : GTATAT\$TAGGGC

The reason why we can do that is because the last character in Sorted rotation should be right before the character at Location in the original Sequence . For instance

### Sorted Rotation



~~Name of Suffix~~

- (f) We use rotation matrix and sort the rows in lexicographic order . Then we output the sequence in the last column as b

Step 1. Construct the rotation matrix

Step 2 . Sort by row . in lexi- order

Step 3 . Output b .

$\checkmark \text{ GTAACTGTAGTG\$}$	$\$ \text{ GTAACTGTAGTG}$	$b = \begin{array}{ c c } \hline G & G \\ \hline \end{array}$
$\checkmark \text{ TAACTGTAGTG\$G}$	$\text{AACTGTAGTG\$GT}$	$\begin{array}{ c c } \hline T & \\ \hline \end{array}$
$\checkmark \text{ AACTGTAGTG\$GT}$	$\text{ACTGTAGTG\$GTA}$	$\begin{array}{ c c } \hline A & A \\ \hline \end{array}$
$\checkmark \text{ ACTGTAGTG\$GTA}$	$\text{AGTG\$GTAACTGT}$	$\begin{array}{ c c } \hline T & T \\ \hline \end{array}$
$\checkmark \text{ CTGTAGTG\$GTAA}$	$\text{CTGTAGTG\$GTAA}$	$\begin{array}{ c c } \hline A & A \\ \hline \end{array}$
$\checkmark \text{ TGTAGTG\$GTAA}$	$\text{TAACTGTAGTG\$G}$	$\begin{array}{ c c } \hline G & T \\ \hline \end{array}$
$\checkmark \text{ TAACTGTAGTG\$G}$	$\text{TAG} \dots \text{G}$	$\begin{array}{ c c } \hline G & \$ \\ \hline \end{array}$
$\checkmark \text{ GTAGTG\$GTAACT}$	$\text{TG\$} \dots \text{G}$	$\begin{array}{ c c } \hline G & T \\ \hline \end{array}$
$\checkmark \text{ TAGTG\$GTAACTG}$	$\text{TGT} \dots \text{C}$	$\begin{array}{ c c } \hline C & A \\ \hline \end{array}$
$\checkmark \text{ AGTG\$GTAACTGT}$	$\text{G\$} \dots \text{T}$	$\begin{array}{ c c } \hline T & G \\ \hline \end{array}$
$\checkmark \text{ GTG\$GTAACTGTGA}$	$\text{GTAA} \dots \text{\$}$	$\begin{array}{ c c } \hline \$ & G \\ \hline \end{array}$
$\checkmark \text{ TG\$GTAACTGTGA}$	$\text{GTAG} \dots \text{T}$	$\begin{array}{ c c } \hline T & G \\ \hline \end{array}$
$\checkmark \text{ G\$GTAACTGTAGT}$	$\text{GTG\$} \dots \text{A}$	$\begin{array}{ c c } \hline A & C \\ \hline \end{array}$
$\checkmark \text{ \$GTAACTGTAGTG}$		

(g) Last

1 2 3 4 5 6 7 8 9 10 11 12 13

(g) Last Column (b) G T A T A T \$ T A G G G C

(Using this we can construct First Column:)

△ First Column \$ A A A C G G G G T T T T

First Occurrence (F)

in First Column

A C G T

2 5 6 10

1 2 3 4 5 6 7 8 9 10 11 12 13

Occurrences within A 0 0 1 1 2 2 2 2 3 3 3 3 3

the first i rows in the C 0 0 0 0 0 0 0 0 0 0 0 0 1

last column (O) G 1 1 1 1 1 1 1 1 2 3 4 4 4

T 0 1 1 2 2 3 3 4 4 4 4 4 4 4

SA

13 3 4 9 5 12 1 7 10 2 8 11 6

△ Notice that we don't need to store First Column in order to get F since we can compute F simply based on b.

Now using b, F, O, we can construct the original sequence S and the suffix array SA. I only give the first several steps here since its also on the lecture notes and all the steps follow the same pattern.

① We start from the first character G in b. Since it must be right before \$ in S. We see O[1, G] = 1, the check F[G] + O[1, G] - 1 = 6 then jump to 6th character in b, b[6] = T, and then check O[6, T] = 3, F[T] + O[6, T] - 1 = 12, so we got next character is b[12] = G ...

We can represent these steps as follow:

(\*) \$ → G<sub>1</sub> → T<sub>3</sub> → G<sub>4</sub> → A<sub>3</sub> → T<sub>2</sub> → G<sub>3</sub> → T<sub>4</sub> → C<sub>1</sub> → A<sub>2</sub> → A<sub>1</sub> → T<sub>1</sub> → G<sub>2</sub>  
So the original sequence S is GTAACGTAGTGTG\$

② Now we can construct SA, which is very useful and will be stored on hard disk since it will be accessed not very frequent in text search. (only once)

If we construct it from original S it will cost O(n log n)

But if we hash the characters in String (\*) above, using its position as value, character as key. we can get SA in O(n)

(\*) \$ → G<sub>1</sub> → T<sub>3</sub> → G<sub>4</sub> → A<sub>3</sub> → T<sub>2</sub> → G<sub>3</sub> → T<sub>4</sub> → C<sub>1</sub> → A<sub>2</sub> → A<sub>1</sub> → T<sub>1</sub> → G<sub>2</sub>

hash into 13 12 11 10 5 8 7 6 5 4 3 2 1

Then we get the hash value of b (with occurrences in O)

b: G T A T A T \$ T A G G G C

O(b): G<sub>1</sub> T<sub>1</sub> A<sub>1</sub> T<sub>2</sub> A<sub>2</sub> T<sub>3</sub> \$ T<sub>4</sub> A<sub>3</sub> G<sub>2</sub> G<sub>3</sub> G<sub>4</sub> C<sub>1</sub>

SA 13 3 4 9 5 12 1 7 10 2 8 11 6

= hash.get(O(b)) + 1

if (14), change count to 1 be done with and (14) case.

Now we have constructed all tables we need for text search. Let's start the searching

To search GTA we first search the occurrence of A then T then G.

1.  $O[13, A] = 3$  occurrences of A in total
2. In F these A's appear on row  $F[A] = 2$  to row  $F[A] + 3 - 1 = 4$
3. on row  $2-1=1$ ,  $O[2-1, T] = 0$  T's have appeared in last column
4. on row 4,  $O[4, T] = 2$  T's have appeared in last column.
5. Therefore these rows cover the  $(O+1) = 1st$  to the  $2nd$  occurrences of ~~A~~ T
6. In the first column these T's appear on row  $F[T] + \frac{1}{2} - 1 = 10$  and  $F[T] + 2 - 1 = 11$

Then we check G. (actually GT)

7. on row  $10-1=9$ ,  $O[9, G] = 1$  G's has appeared in last col.
8. on row  $11-\cancel{10}=1$ ,  $O[11, G] = 3$  G's has appeared in last column.
9. therefore these rows cover the  $(1+1) = 2nd$  to the  $3rd$  occurrences of G.
10. In the first column these G's appear on row  $F[G] + 2 - 1 = 7$  and  $F[G] + 3 - 1 = 8$
11. These GTA's appear at position  $SA[7] = 1$  and  $SA[8] = 8$  of the input string s.

(h) Well apparently not equal in each of the  $n+1$  positions.  
since in BWT output b, \$ can never appear in the first place of b, since \$ is smallest in lexicographical order, so the first line of ~~the~~ sorted rotation matrix is always \$ -----