

# Computing Matching Statistics on Repetitive Texts

---

Younan Gao

Faculty of Computer Science, Dalhousie University, Canada

# Matching Statistics

The matching statistics  $MS$  of a pattern  $P[1..m]$  with respect to a text  $T[1..n]$  is an array of integers  $MS[1..m]$  such that the  $i$ -th entry  $MS[i]$  stores the length of the longest prefix of  $P[i..m]$  that occurs in  $T$ .

For example, given that  $T[1..8] = \text{"aaabbbcc"}$  and  $P[1..5] = \text{"ccabb"}$ , the matching statistics  $MS[1..5] = \{2, 1, 3, 2, 1\}$ .

## Related Work

Space	Time	Reference
$O(n)$	$O(m \lg \sigma)$	Textbook
$(n \lg \sigma + o(n \lg \sigma))$ bits	$O(m \lg \sigma)$	Enno et al.
$O(r + S(n))$	$O(m \cdot f(n))$	Bannai et al.
$O(z + S(n))$	$O(m^2 \lg^\epsilon z + m \cdot f(n))$	New
$O(z \lg z + S(n))$	$O(m^2 + m \lg z \lg \lg z + m \cdot f(n))$	New
$O(z \lg z + \frac{z}{\log_\sigma n} \lg^{2\epsilon+1} z + S(n))$	$O(m^2 \lg \lg \sigma + m \cdot f(n))$	New
$O(z \lg z + S(n))$	$O(m^2 + m \cdot f(n))$	$\sigma$ is constant

- $z$  is the num of phrases in the Lempel-Ziv Parsing, while  $r$  is the num of runs in BWT.
- Assume that there is a data structure of  $S(n)$  words of space to support retrieving any substring  $T[i..i + \ell]$  in  $O(f(n) + \ell)$  time.
- $r = O(z \lg^2 n)$ .

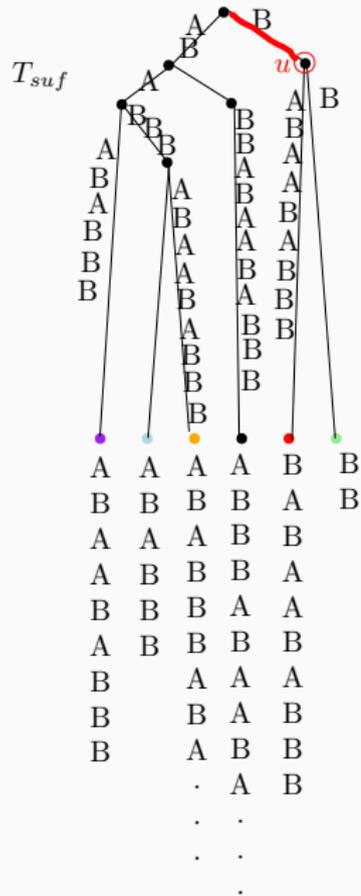
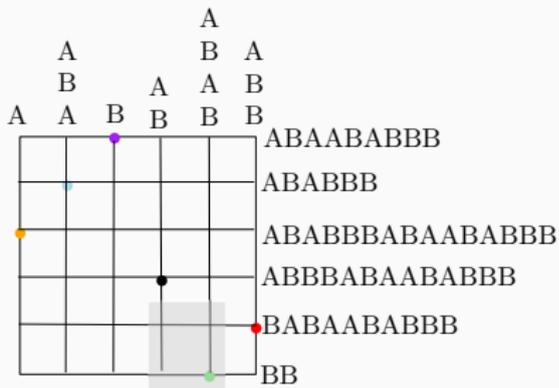
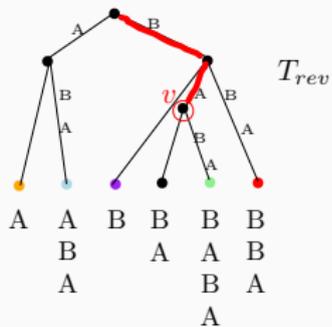
# Preliminaries

$Text[1..16] = A|AB|ABB|B|ABA|ABAB|BB$

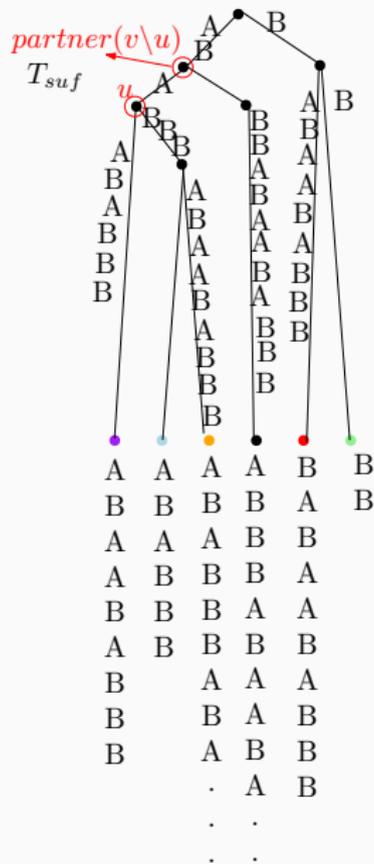
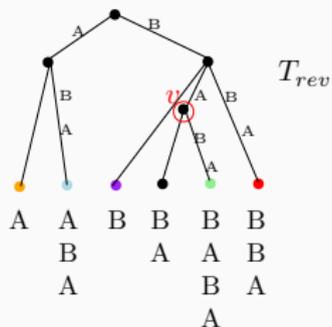
Phrases	Phrases <sub>rev</sub>	Suffixes
<i>A</i>	<i>A</i>	<i> AB ABB B ABA ABAB BB</i>
<i>AB</i>	<i>BA</i>	<i> ABB B ABA ABAB BB</i>
<i>ABB</i>	<i>BBA</i>	<i> B ABA ABAB BB</i>
<i>B</i>	<i>B</i>	<i> ABA ABAB BB</i>
<i>ABA</i>	<i>ABA</i>	<i> ABAB BB</i>
<i>ABAB</i>	<i>BABA</i>	<i> BB</i>



# Preliminaries: Induced Relationship

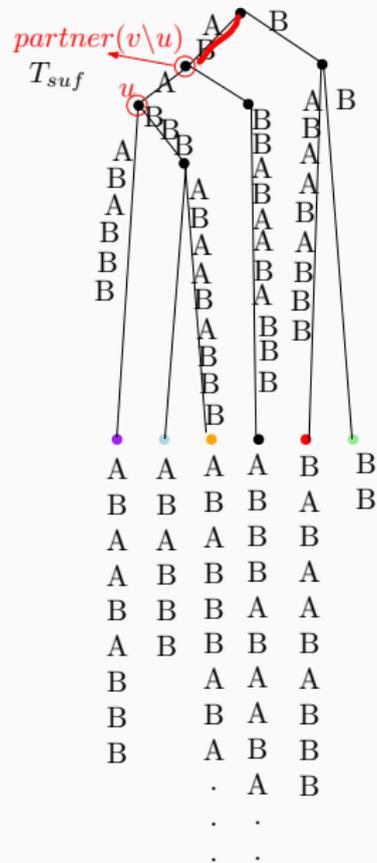
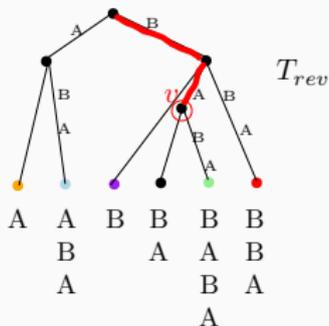


# Preliminaries: Partner Finding



- Operation  $partner(v \setminus u)$  can be implemented by 2D orthogonal range succ/prec queries.
- String ABAB appears in the text, AABABBBABAABABBB, but string ABABA does not;

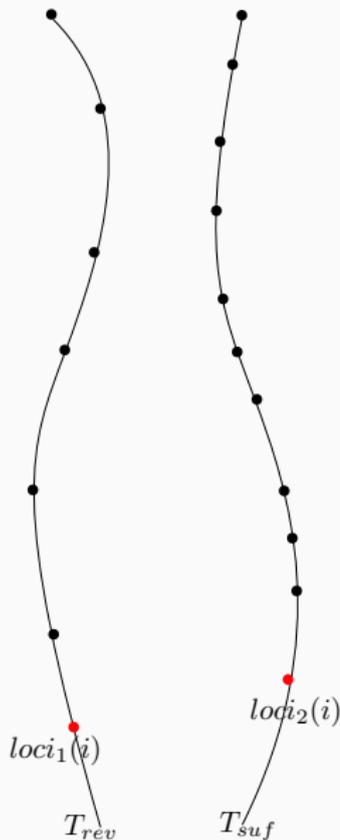
# Preliminaries: Partner Finding



- Operation  $partner(v \setminus u)$  can be implemented by 2D orthogonal range succ/prec queries.
- String ABAB appears in the text, AABABBBABAABABBB, but string ABABA does not;

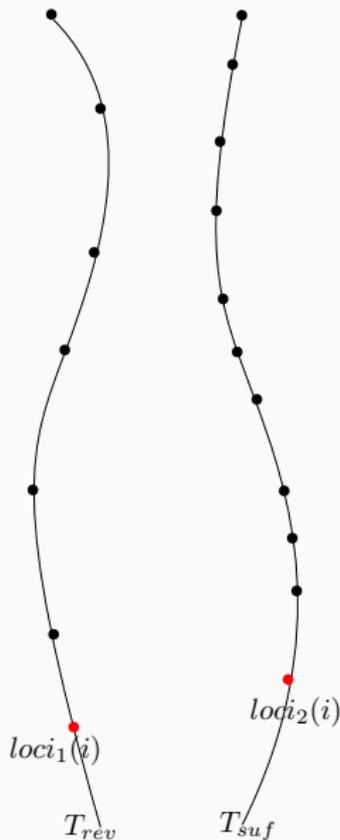
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



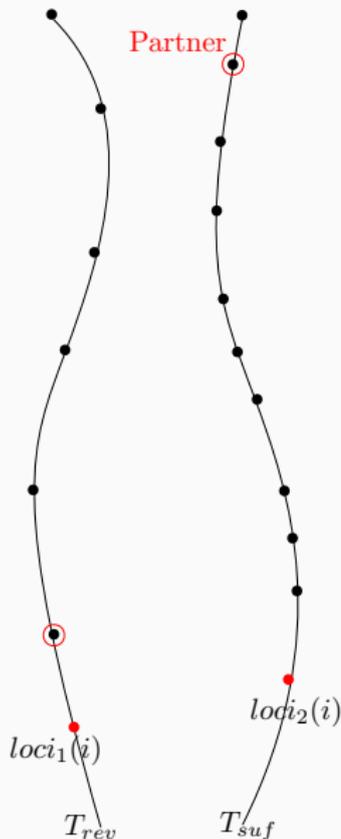
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



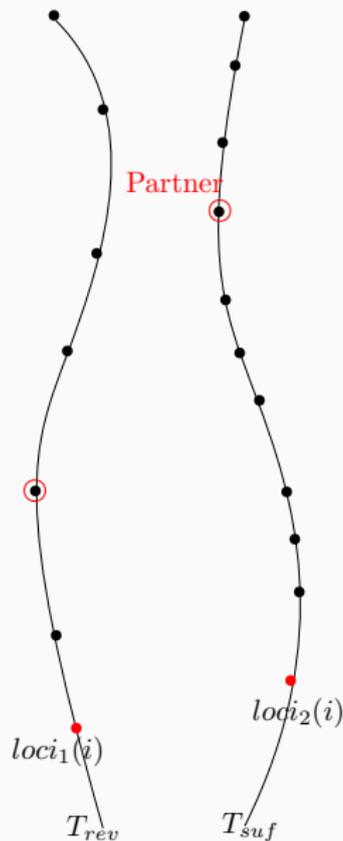
## Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



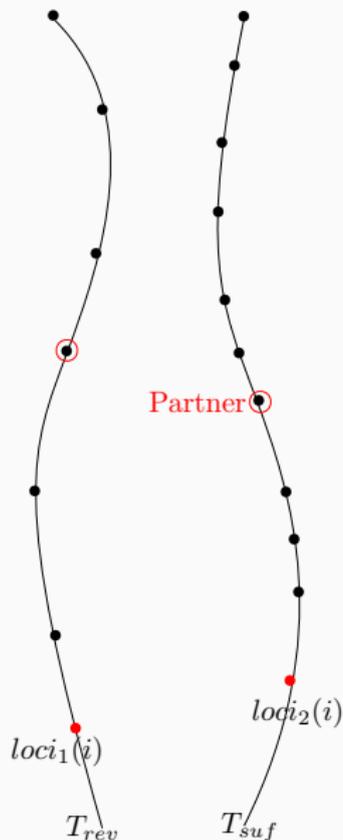
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



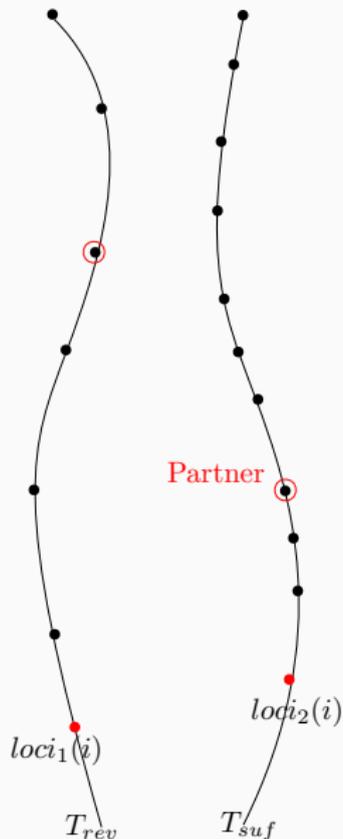
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



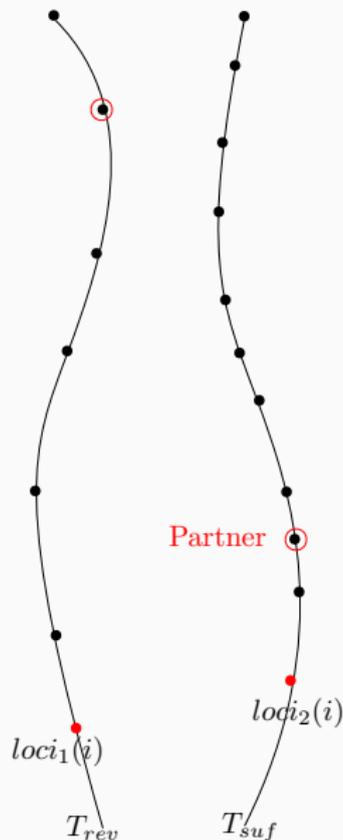
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



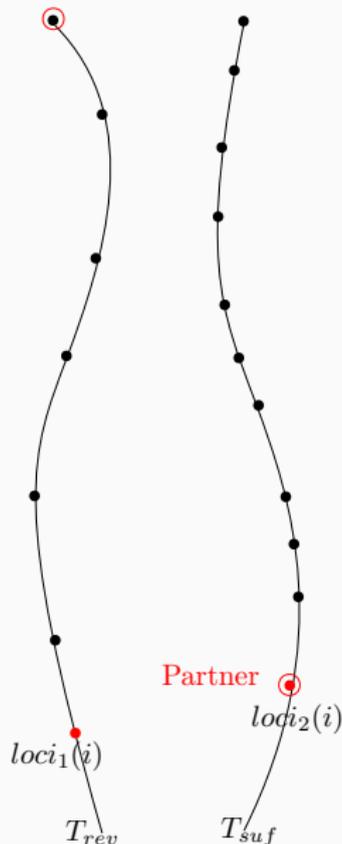
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



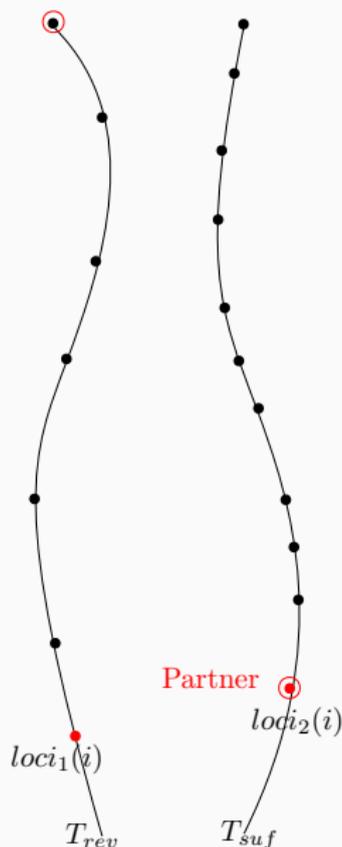
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



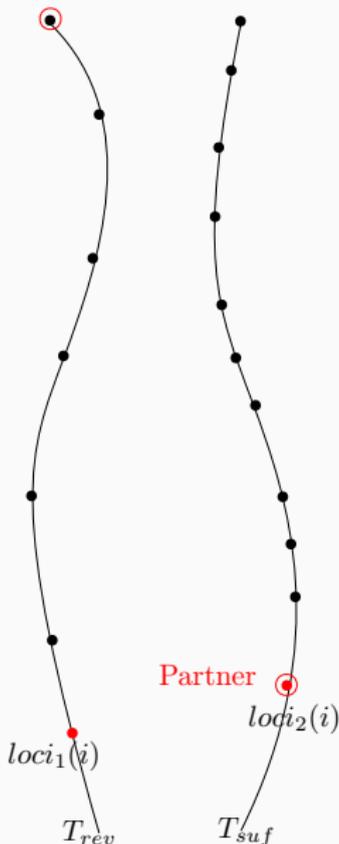
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



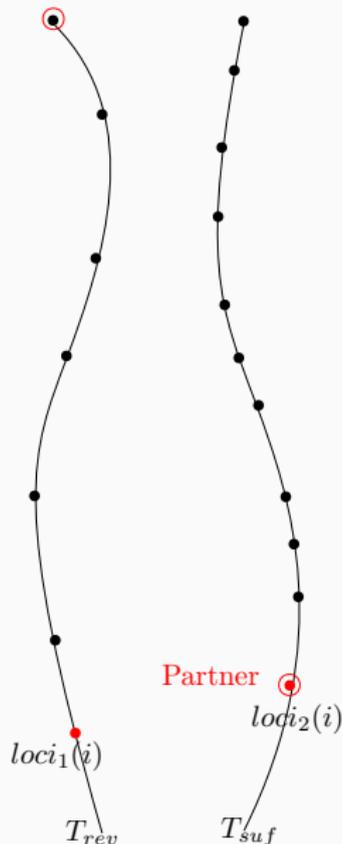
# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



# Naive Method for Computing $MS[1..m]$

- Query pattern  $P[1..m]$  is split into  $m - 1$  pairs of prefix and suffix pairs.
- Consider the  $i$ -th pair,  $P[1..i]$  and  $P[i + 1, m]$ ;
- For each ancestor,  $u$ , of  $loc_1(i)$ , find  $partner(u \setminus loc_2(i))$
- Overall, we have  $\sum_{i=1}^{m-1} i = O(m^2)$  partner (or LCP) queries.
- Query Time:  $O(m^2 + mf(n) + m^2 \lg^\epsilon z)$ .
- Space Cost:  $O(z + S(n))$  words of space.



## Second Method: LPMEM

$P[1..i] = abwwz$   
↓  
 $loci_1(i)$

$P[i+1..m] = xyycd\dots$   
↓  
 $loci_2(i)$

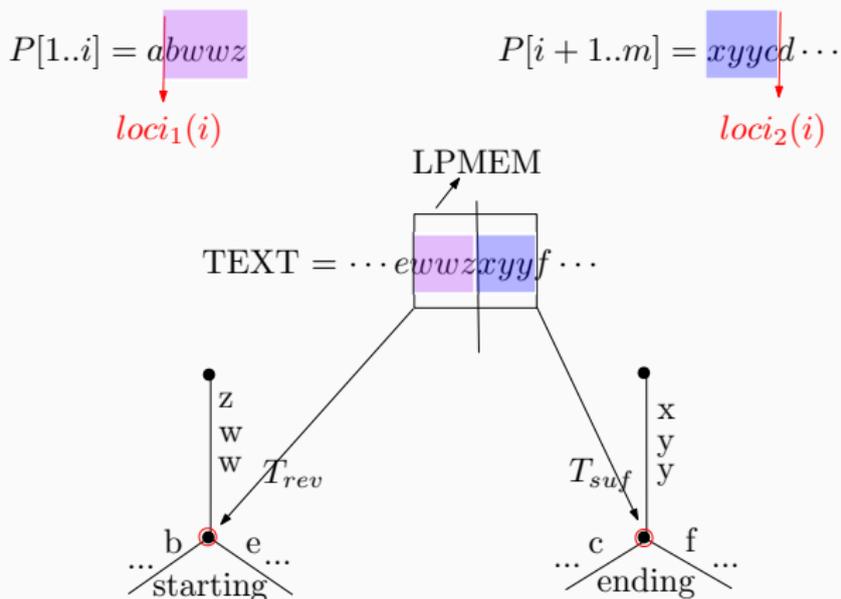
## Second Method: LPMEM

$P[1..i] = abwwz$   
↓  
 $loc_1(i)$

$P[i+1..m] = xyycd\dots$   
↓  
 $loc_2(i)$

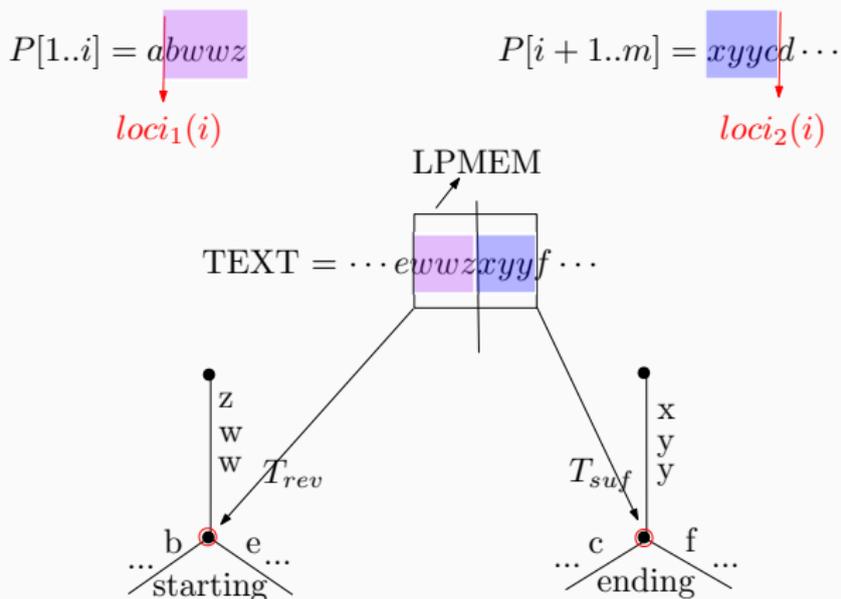
LPMEM  
↑  
TEXT =  $\dots ewwzxyyf \dots$

## Second Method: LPMEM



- The number,  $occ$ , of LPMEMs for  $P[1..m]$  is at most  $\binom{m}{2}$ .

## Second Method: LPMEM

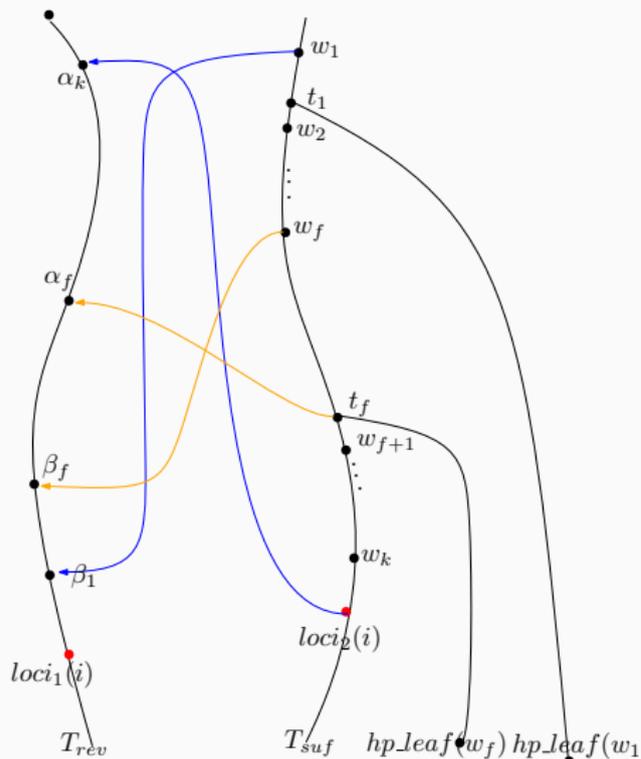


- The number,  $occ$ , of LPMEMs for  $P[1..m]$  is at most  $\binom{m}{2}$ .
- MS can be computed in  $O(m + occ)$  time —Algorithm 2.



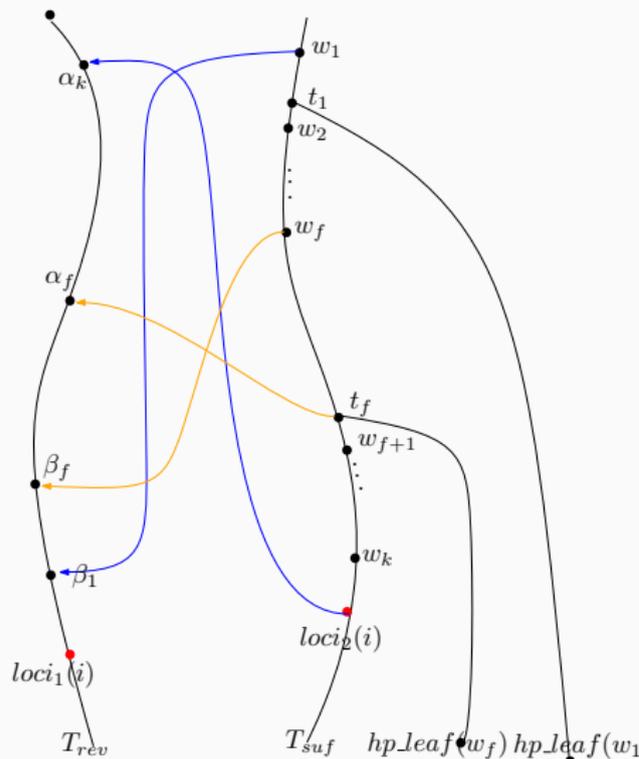
## Second Method: Query Algorithms

- Apply heavy path decomposition on  $T_{suf}$ .
- For each  $1 \leq f \leq k = O(\lg z)$ , set  $\alpha_f = \text{partner}(t_f \setminus \text{loci}_1(i))$  and  $\beta_f = \text{partner}(w_f \setminus \text{loci}_1(i))$ ;
- $\alpha_f$  and  $\text{partner}(\alpha_f \setminus \text{loci}_2(i))$  induce a LPMEM; so do  $\beta_f$  and  $\text{partner}(\beta_f \setminus \text{loci}_2(i))$ : **Lemma 6**;
- If  $u$  and  $v$  are induced together, and if  $v$  stays between  $w_f$  and  $t_f$ , then  $u$  stays between  $\alpha_f$  and  $\beta_f$ : **Lemma 7**.



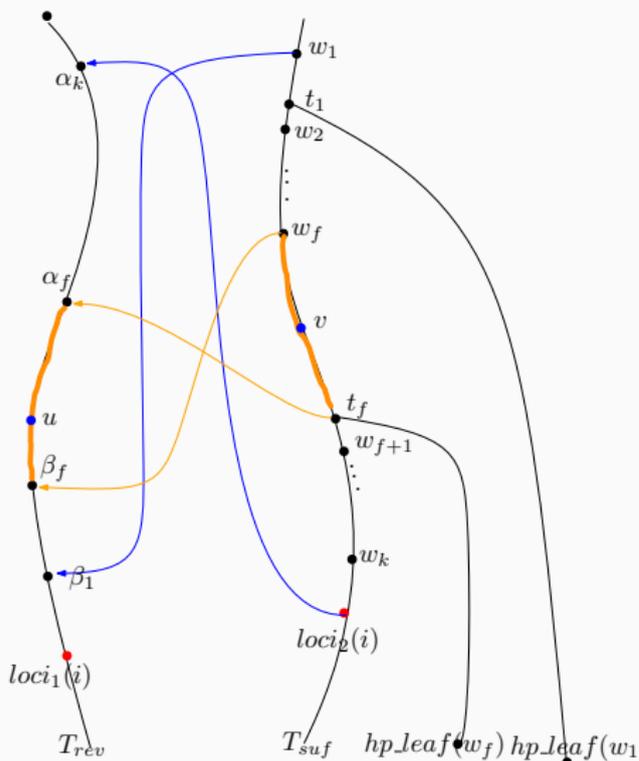
## Second Method: Query Algorithms

- Apply heavy path decomposition on  $T_{suf}$ .
- For each  $1 \leq f \leq k = O(\lg z)$ , set  $\alpha_f = \text{partner}(t_f \setminus \text{loci}_1(i))$  and  $\beta_f = \text{partner}(w_f \setminus \text{loci}_1(i))$ ;
- $\alpha_f$  and  $\text{partner}(\alpha_f \setminus \text{loci}_2(i))$  induce a LPMEM; so do  $\beta_f$  and  $\text{partner}(\beta_f \setminus \text{loci}_2(i))$ : **Lemma 6**;
- If  $u$  and  $v$  are induced together, and if  $v$  stays between  $w_f$  and  $t_f$ , then  $u$  stays between  $\alpha_f$  and  $\beta_f$ : **Lemma 7**.



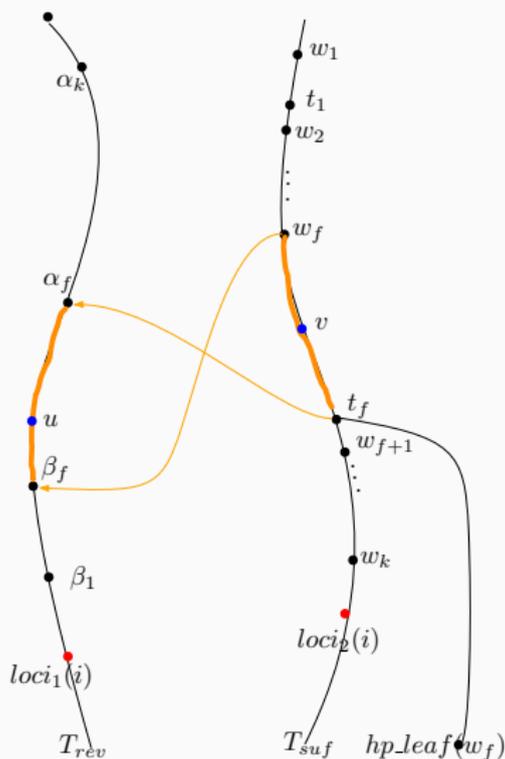
## Second Method: Query Algorithms

- Apply heavy path decomposition on  $T_{suf}$ .
- For each  $1 \leq f \leq k = O(\lg z)$ , set  $\alpha_f = \text{partner}(t_f \setminus \text{loc}_1(i))$  and  $\beta_f = \text{partner}(w_f \setminus \text{loc}_1(i))$ ;
- $\alpha_f$  and  $\text{partner}(\alpha_f \setminus \text{loc}_2(i))$  induce a LPMEM; so do  $\beta_f$  and  $\text{partner}(\beta_f \setminus \text{loc}_2(i))$ : **Lemma 6**;
- If  $u$  and  $v$  are induced together, and if  $v$  stays between  $w_f$  and  $t_f$ , then  $u$  stays between  $\alpha_f$  and  $\beta_f$ : **Lemma 7**.



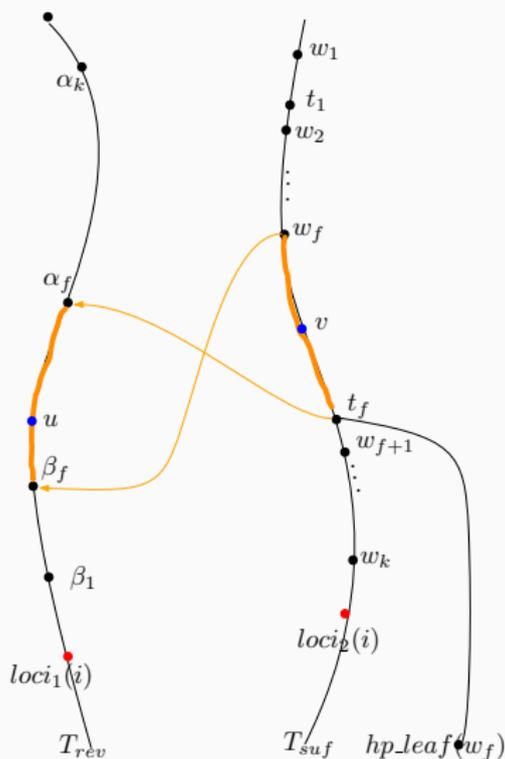
## Second Method: Preprocessing

- $w_f$  and  $hp\_leaf(w_f)$  are known during the preprocessing stage;
- $t_f$  and  $loci_2(i)$  are unknown, but  $partner(u \setminus loci_2(i)) = partner(u \setminus hp\_leaf(w_f)) = partner(u \setminus t_f)$ : **Lemma 8**;
- $\alpha_f$  and  $\beta_f$  are unknown, but we can use the induced subtree  $T_{rev}(w_f)$ .



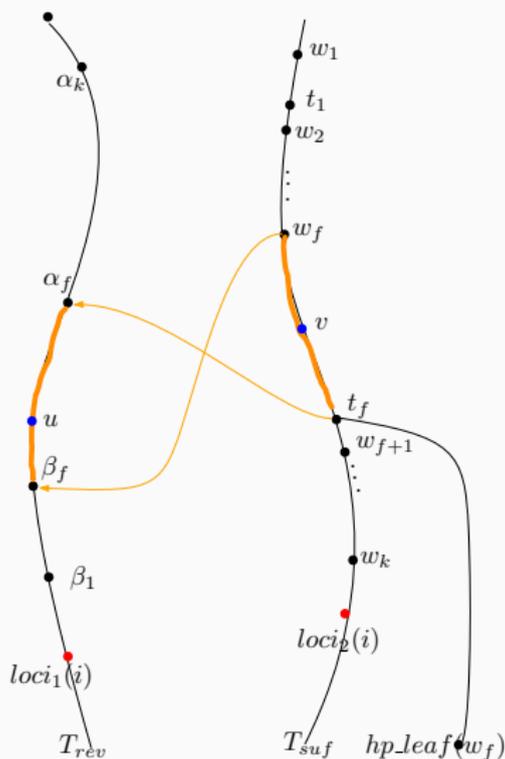
## Second Method: Preprocessing

- $w_f$  and  $hp\_leaf(w_f)$  are known during the preprocessing stage;
- $t_f$  and  $loci_2(i)$  are unknown, but  $partner(u \setminus loci_2(i)) = partner(u \setminus hp\_leaf(w_f)) = partner(u \setminus t_f)$ : **Lemma 8**;
- $\alpha_f$  and  $\beta_f$  are unknown, but we can use the induced subtree  $T_{rev}(w_f)$ .

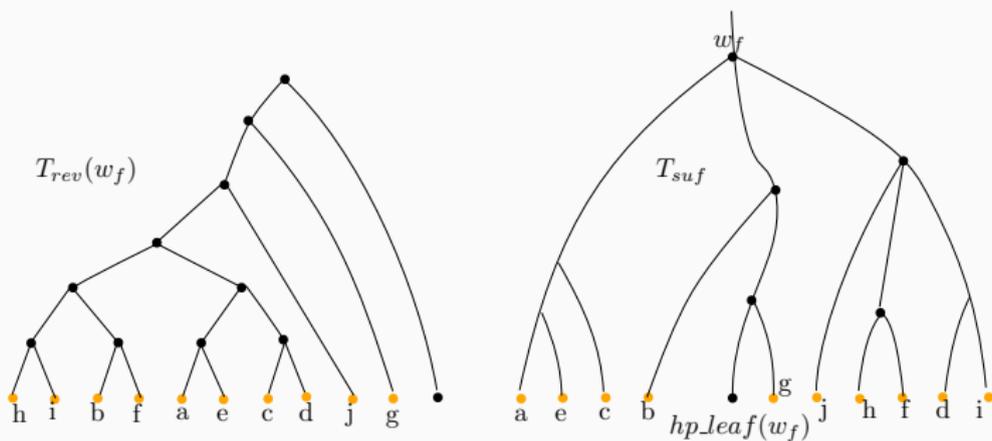


## Second Method: Preprocessing

- $w_f$  and  $hp\_leaf(w_f)$  are known during the preprocessing stage;
- $t_f$  and  $loci_2(i)$  are unknown, but  $partner(u \setminus loci_2(i)) = partner(u \setminus hp\_leaf(w_f)) = partner(u \setminus t_f)$ : **Lemma 8**;
- $\alpha_f$  and  $\beta_f$  are unknown, but we can use the induced subtree  $T_{rev}(w_f)$ .

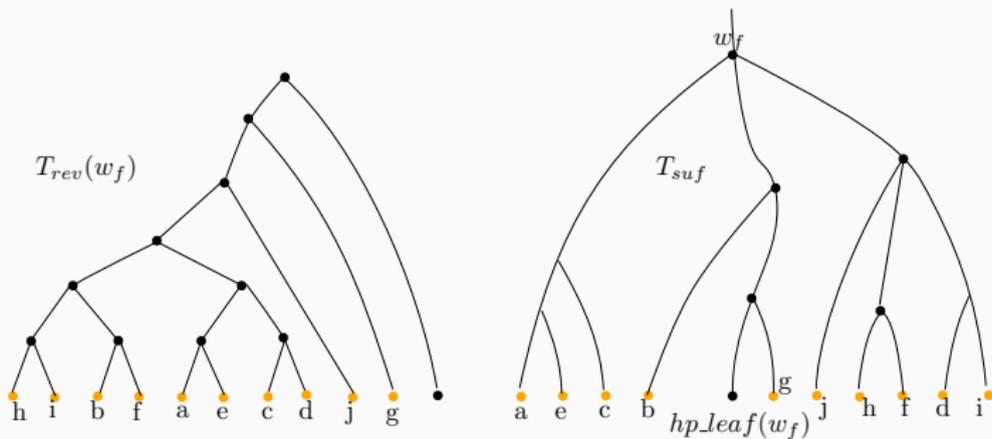


## Second Method: Induced Subtree $T_{ref}(w_f)$



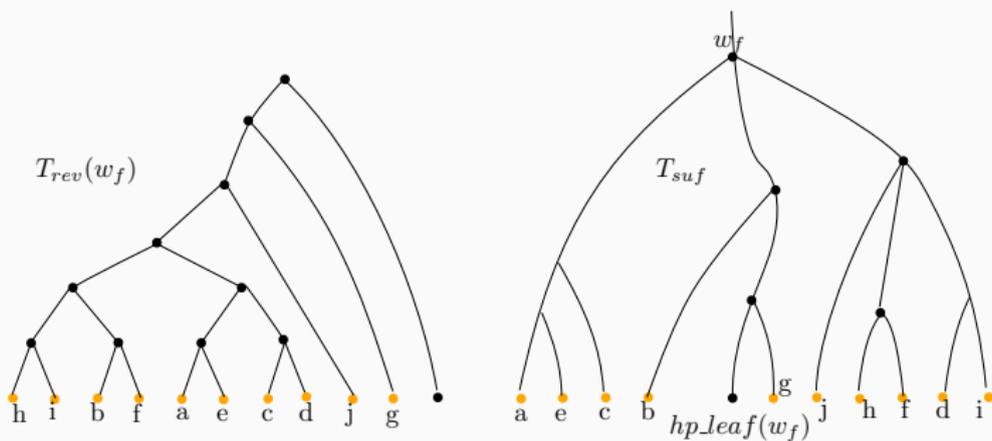
- Leaves in  $T_{rev}$  that are induced with  $w_f$  are called special leaves.
- $T_{rev}(w_f)$  contains the special leaves and their LCAs (special nodes).
- $\sum_{w_f} |special(w_f)| = O(z \lg z)$ .
- The left endpoint of a LPMEM always stays at a special internal node: **Lemma 9**.

## Second Method: Induced Subtree $T_{ref}(w_f)$



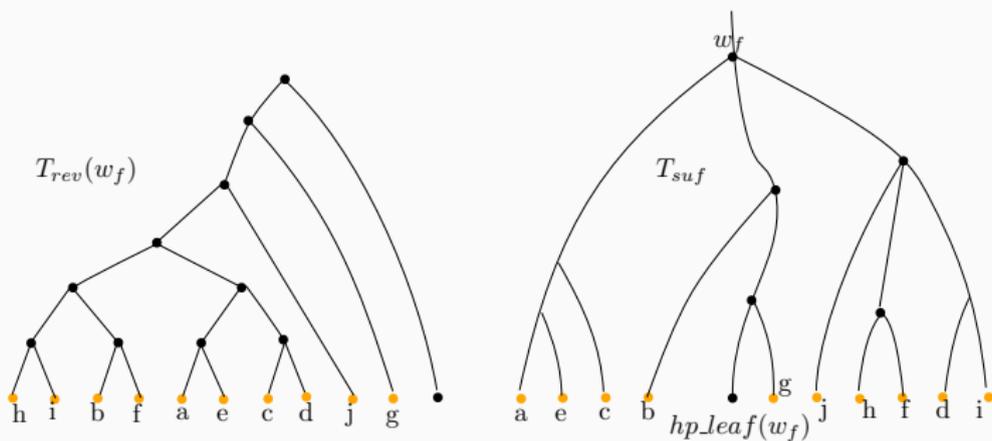
- Leaves in  $T_{rev}$  that are induced with  $w_f$  are called special leaves.
- $T_{rev}(w_f)$  contains the special leaves and their LCAs (special nodes).
- $\sum_{w_f} |special(w_f)| = O(z \lg z)$ .
- The left endpoint of a LPMEM always stays at a special internal node: **Lemma 9**.

## Second Method: Induced Subtree $T_{ref}(w_f)$



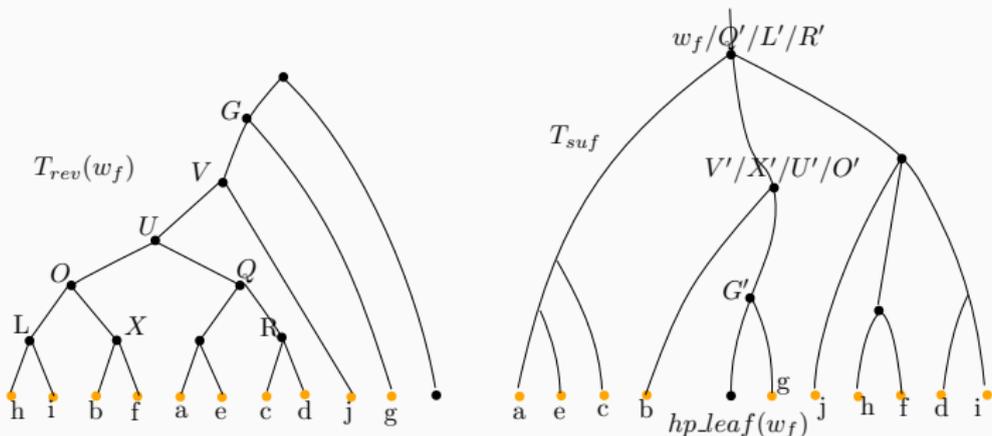
- Leaves in  $T_{rev}$  that are induced with  $w_f$  are called special leaves.
- $T_{rev}(w_f)$  contains the special leaves and their LCAs (special nodes).
- $\sum_{w_f} |special(w_f)| = O(z \lg z)$ .
- The left endpoint of a LPMEM always stays at a special internal node: **Lemma 9**.

## Second Method: Induced Subtree $T_{ref}(w_f)$



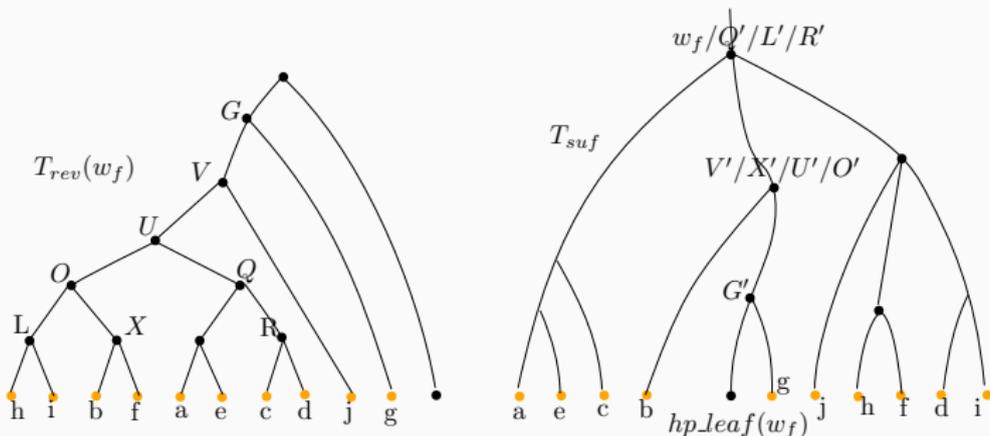
- Leaves in  $T_{rev}$  that are induced with  $w_f$  are called special leaves.
- $T_{rev}(w_f)$  contains the special leaves and their LCAs (special nodes).
- $\sum_{w_f} |special(w_f)| = O(z \lg z)$ .
- The left endpoint of a LPMEM always stays at a special internal node: **Lemma 9**.

## Second Method: Induced Subtree $T_{ref}(w_f)$



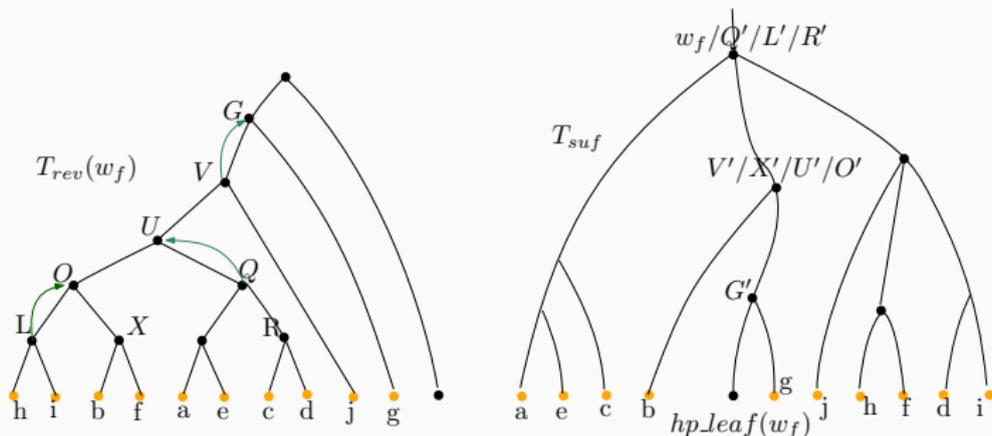
- Each special internal node  $U$  in  $T_{rev}(w_f)$  stores a pointer,  $e_2$ , pointing to  $U' = partner(U/hp\_leaf(w_f))$  in  $T_{suf}$ .
- Each special internal node  $U$  stores another pointer,  $e_0$ :
  - Given a pair of parent and child nodes,  $O$  and  $L$ , if  $L$  does not induce with  $O'$ , then pointer  $e_0$  of  $L$  points to  $O$ .
  - Otherwise, its  $e_0$  points to the same node as of its lowest ancestor that has a valid  $e_0$ .

## Second Method: Induced Subtree $T_{ref}(w_f)$



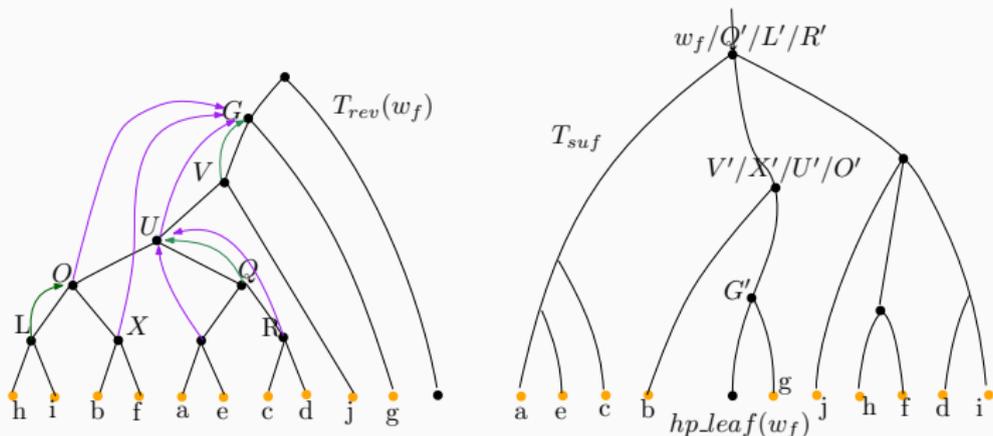
- Each special internal node  $U$  in  $T_{rev}(w_f)$  stores a pointer,  $e_2$ , pointing to  $U' = partner(U/hp\_leaf(w_f))$  in  $T_{suf}$ .
- Each special internal node  $U$  stores another pointer,  $e_0$ :
  - Given a pair of parent and child nodes,  $O$  and  $L$ , if  $L$  does not induce with  $O'$ , then pointer  $e_0$  of  $L$  points to  $O$ .
  - Otherwise, its  $e_0$  points to the same node as of its lowest ancestor that has a valid  $e_0$ .

## Second Method: Induced Subtree $T_{ref}(w_f)$



- Each special internal node  $U$  in  $T_{rev}(w_f)$  stores a pointer,  $e_2$ , pointing to  $U' = partner(U/hp\_leaf(w_f))$  in  $T_{suf}$ .
- Each special internal node  $U$  stores another pointer,  $e_0$ :
  - Given a pair of parent and child nodes,  $O$  and  $L$ , if  $L$  does not induce with  $O'$ , then pointer  $e_0$  of  $L$  points to  $O$ .
  - Otherwise, its  $e_0$  points to the same node as of its lowest ancestor that has a valid  $e_0$ .

## Second Method: Induced Subtree $T_{ref}(w_f)$



- Each special internal node  $U$  in  $T_{rev}(w_f)$  stores a pointer,  $e_2$ , pointing to  $U' = partner(U/hp\_leaf(w_f))$  in  $T_{suf}$ .
- Each special internal node  $U$  stores another pointer,  $e_0$ :
  - Given a pair of parent and child nodes,  $O$  and  $L$ , if  $L$  does not induce with  $O'$ , then pointer  $e_0$  of  $L$  points to  $O$ .
  - Otherwise, its  $e_0$  points to the same node as of its lowest ancestor that has a valid  $e_0$ .

## Second Method & Open Problems

- The second method:
  - Query Time:  $O(m^2 + mf(n) + m \lg z \lg \lg z)$ .
  - Space Cost:  $O(z \lg z + S(n))$  words of space.
- Open Problems:
  - Query Time:  $O(m \cdot f(n))$
  - Space:  $O(z \lg z + S(n))$  or even  $O(z + S(n))$ ?