

# Fast Algorithms for Computing Tree LCS

Shay Mozes<sup>1</sup> Dekel Tsur<sup>2</sup>  
Oren Weimann<sup>3</sup> Michal Ziv-Ukelson<sup>2</sup>

<sup>1</sup>Brown University

<sup>2</sup>Ben-Gurion University

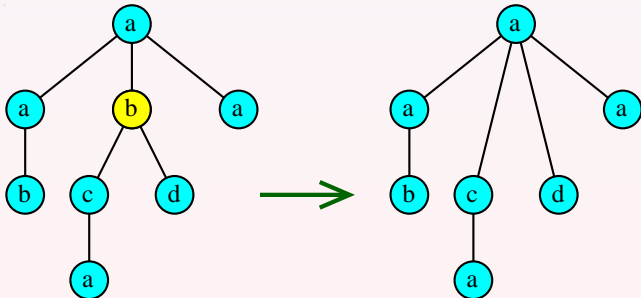
<sup>3</sup>Massachusetts Institute of Technology

# Definition of tree LCS

## Definition

The **LCS** of two trees is the size of largest forest that can be obtained from both trees by deleting nodes.

Deletion of a node  $v$  in a forest  $F$  means that the children of  $v$  become children of the parent of  $v$  (if it exists).

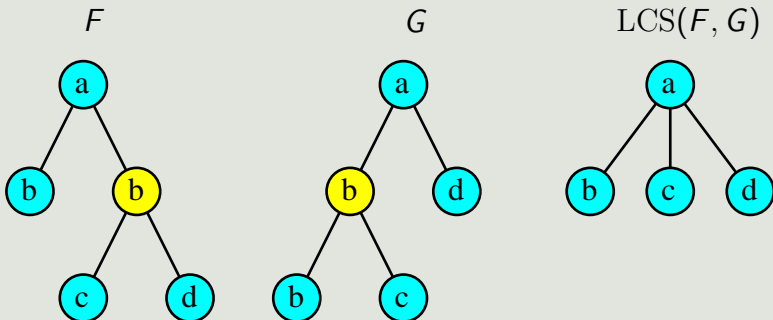


# Definition of tree LCS

## Definition

The **LCS** of two trees is the size of largest forest that can be obtained from both trees by deleting nodes.

## Example



# Known & new results

Let  $F$  and  $G$  be the input trees, with sizes  $n \geq m$ .

Known algorithms for tree edit distance/tree LCS:

- $O(nm \cdot \text{leaves}(F)^2 \cdot \text{leaves}(G)^2) = O(n^6)$  (Tai 79)
- $O(nm \cdot \text{cdepth}(F) \cdot \text{cdepth}(G)) = O(n^4)$  (Zhang & Shasha 89)
- $O(m^2 n \log n) = O(n^3 \log n)$  (Klein 98)
- $O(nm^2(1 + \log \frac{n}{m})) = O(n^3)$  (Demaine et al. 07)

# Known & new results

Let  $F$  and  $G$  be the input trees, with sizes  $n \geq m$ .

Known algorithms for tree edit distance/tree LCS:

- $O(nm \cdot \text{leaves}(F)^2 \cdot \text{leaves}(G)^2) = O(n^6)$  (Tai 79)
- $O(nm \cdot \text{cdepth}(F) \cdot \text{cdepth}(G)) = O(n^4)$  (Zhang & Shasha 89)
- $O(m^2 n \log n) = O(n^3 \log n)$  (Klein 98)
- $O(nm^2(1 + \log \frac{n}{m})) = O(n^3)$  (Demaine et al. 07)

New algorithms for tree LCS:

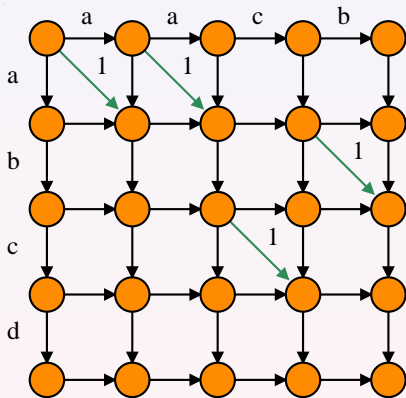
- $O(r \cdot \text{height}(F) \cdot \text{height}(G) \cdot \log \log m)$
- $O(Lr \log r \cdot \log \log m)$

$r$  = number of match pairs (pairs of vertices  $v \in F$ ,  $w \in G$  with equal labels)

$L = \text{LCS}(F, G)$



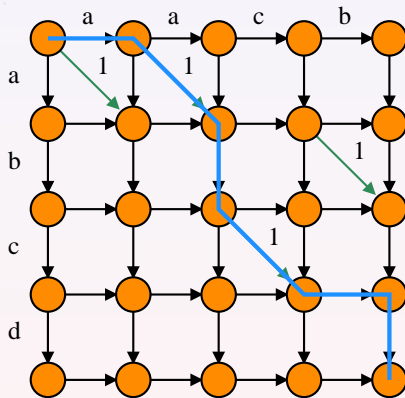
# Alignment graph for string LCS



The edges are

- Vertical and horizontal edges of weight 0 corresponding to deletion of character from  $S$  or  $T$ .
- Diagonal edges for every match pair of weight 1.

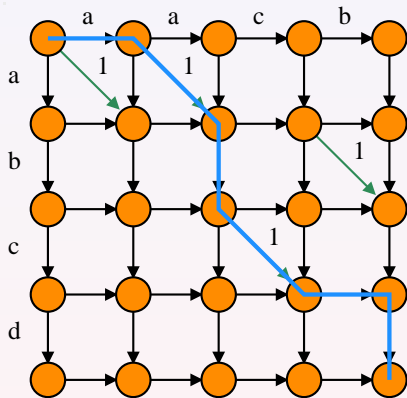
# Alignment graph for string LCS



- A path of weight  $k$  from  $(0, 0)$  to  $(n, m)$  corresponds to a common substring of size  $k$ .
- A common substring of size  $k$  corresponds to a path of weight  $k$  from  $(0, 0)$  to  $(n, m)$ .



# Alignment graph for string LCS

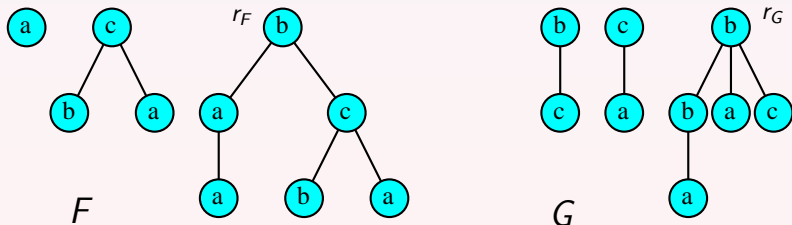


- The number of nonzero weight edges is  $r$ .
- The maximum weight path can be found in time  $O(r \log \log m)$  (Hunt-Szymanski 77).

# The algorithm of Zhang & Shasha

$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

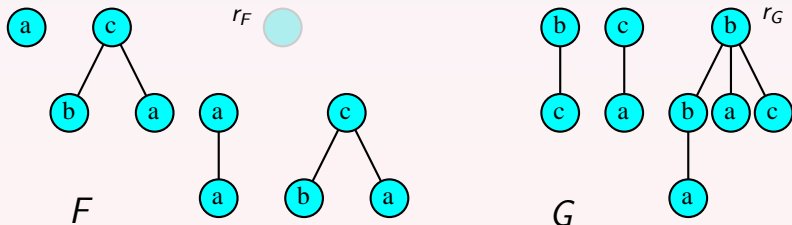
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Zhang & Shasha

$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

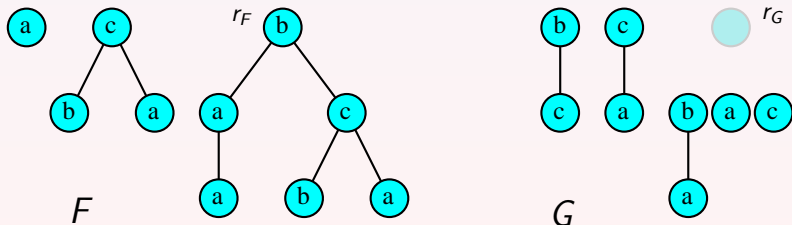
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Zhang & Shasha

$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

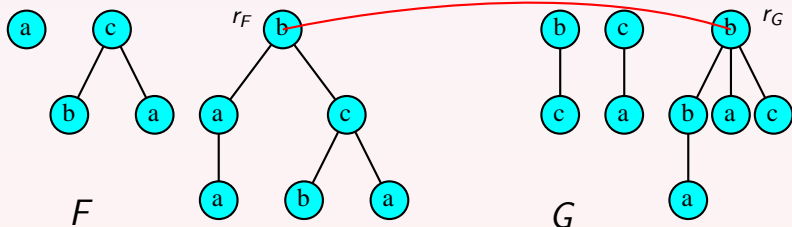
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Zhang & Shasha

$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

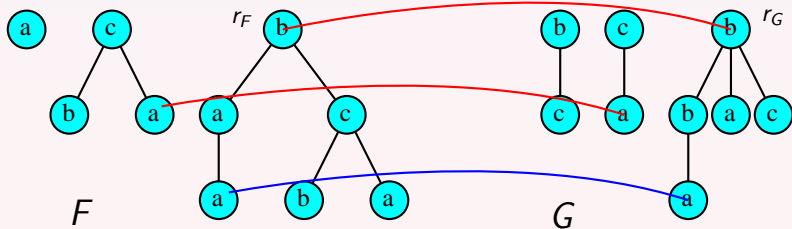
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Zhang & Shasha

$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

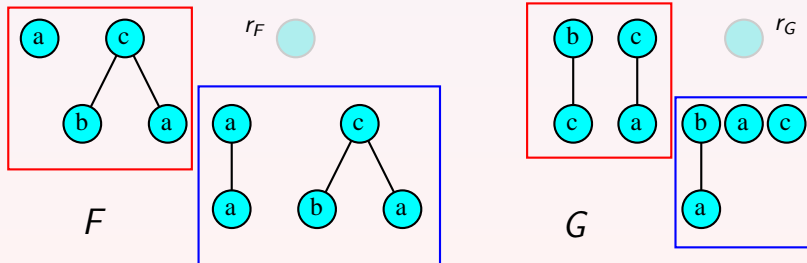
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Zhang & Shasha

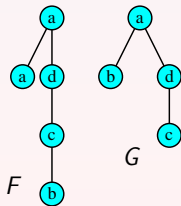
$R_F$  = rightmost tree in  $F$ .  $r_G$  = root of  $R_G$ .

$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# Relevant subproblems

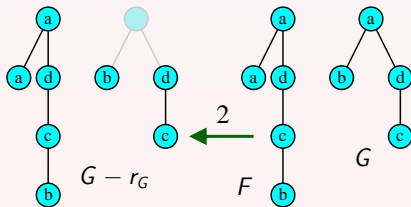
$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} 1 \\ 2 \end{array}$$





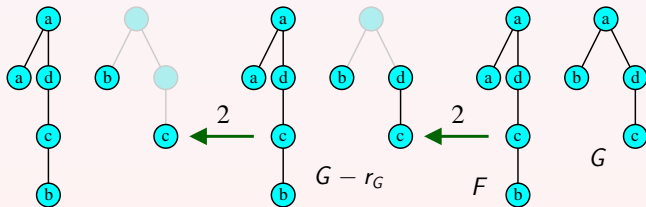
# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} 1 \\ 2 \end{array}$$



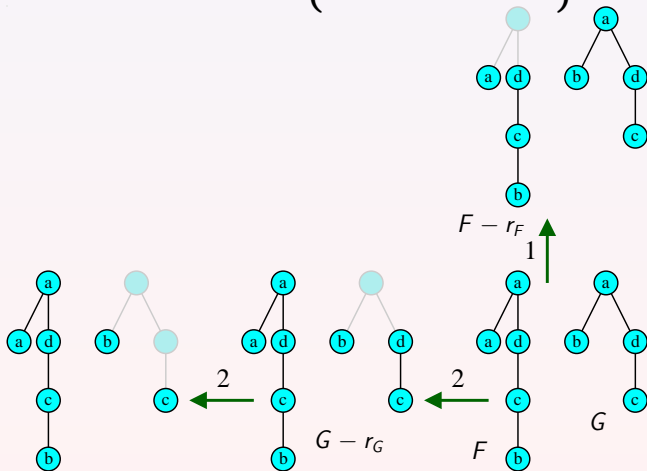
# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} 1 \\ 2 \end{array}$$



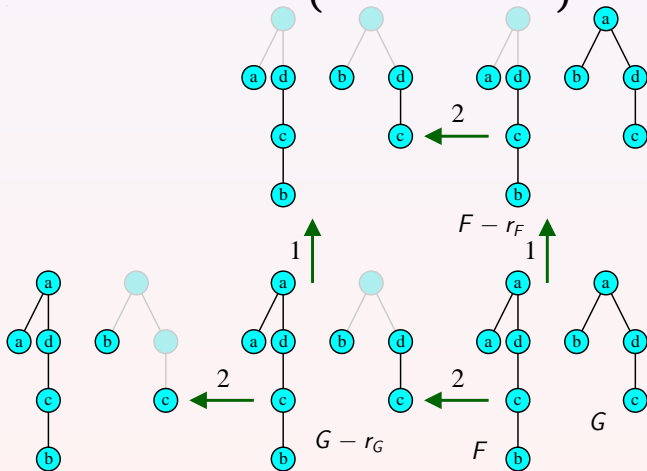
# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\}$$



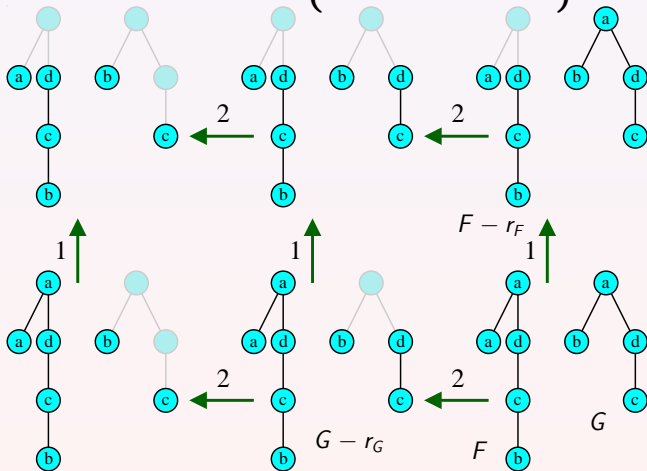
# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} 1 \\ 2 \end{array}$$



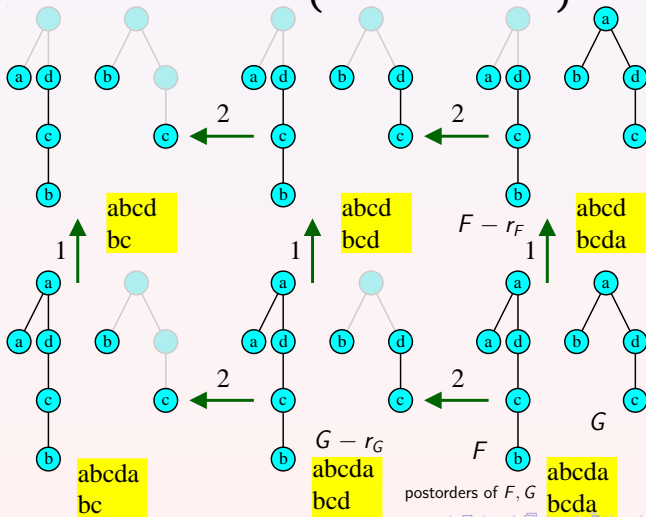
# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\}$$

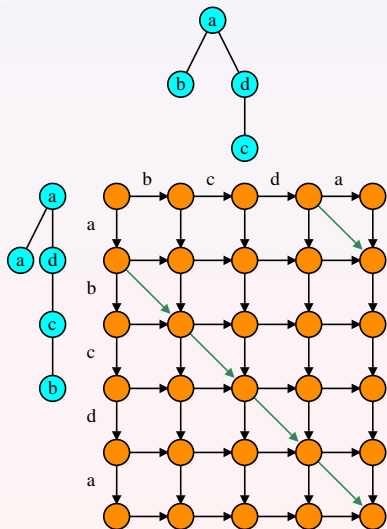


# Relevant subproblems

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} 1 \\ 2 \end{array}$$



# Alignment graph (Backofen et al., Touzet)



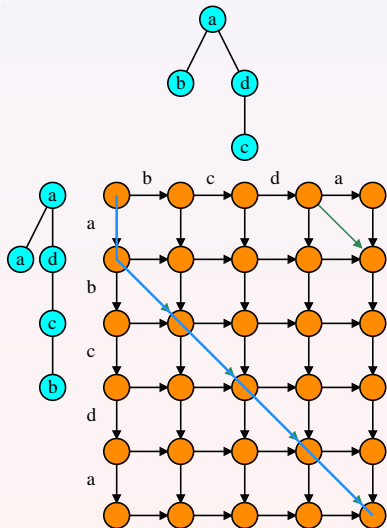
Consider the alignment graph for the two postorder strings of  $F$  and  $G$ .

A path of weight  $k$  **does not** correspond to a common subforest of size  $k$ .

# Alignment graph (Backofen et al., Touzet)

Consider the alignment graph for the two postorder strings of  $F$  and  $G$ .

A path of weight  $k$  **does not** correspond to a common subforest of size  $k$ .

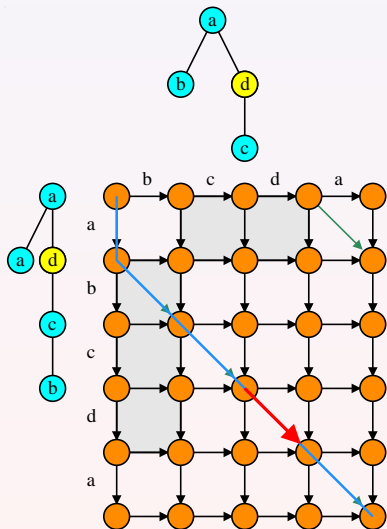




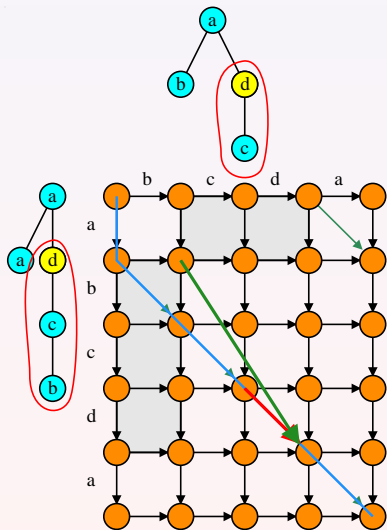
# Alignment graph (Backofen et al., Touzet)

Consider the alignment graph for the two postorder strings of  $F$  and  $G$ .

A path of weight  $k$  **does not** correspond to a common subforest of size  $k$ .



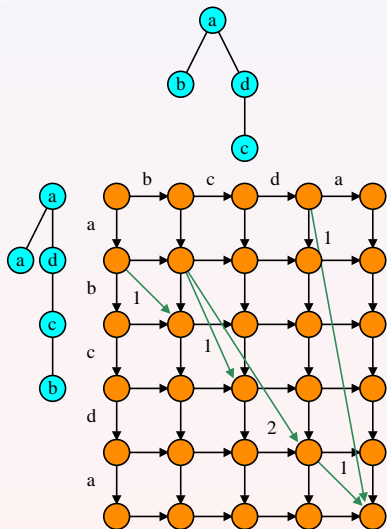
# Alignment graph (Backofen et al., Touzet)



Consider the alignment graph for the two postorder strings of  $F$  and  $G$ .

A path of weight  $k$  **does not** correspond to a common subforest of size  $k$ .

# Alignment graph (Backofen et al., Touzet)



For each match pair  $v, w$  (with  $v \in F, w \in G$ ) add an edge

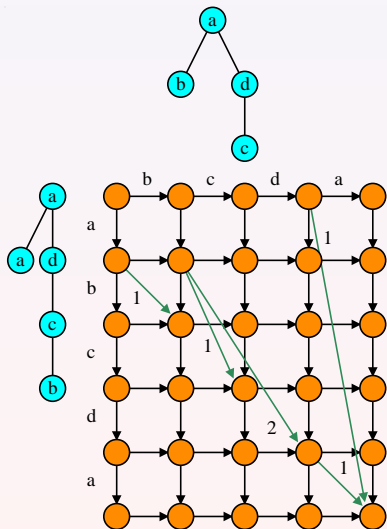
$$(i - |F_v|, j - |G_w|) \rightarrow (i, j)$$

of weight  $\text{LCS}(F_v, G_w)$ , where  $i$  and  $j$  are the ranks of  $v$  and  $w$  in the postorders of  $F$  and  $G$ .

For the match pair of the roots we just add an edge  $(n - 1, m - 1) \rightarrow (n, m)$  of weight 1.

The LCS values are computed recursively.

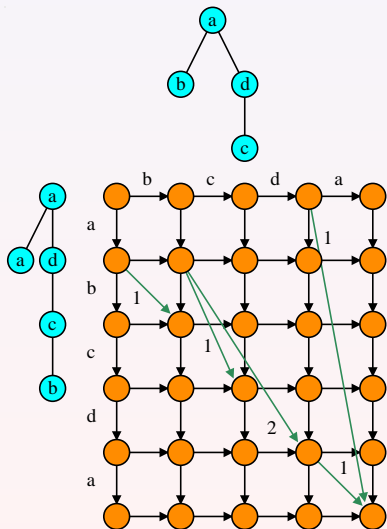
# Alignment graph (Backofen et al., Touzet)



Every path of weight  $k$  from  $(0, 0)$  to  $(n, m)$  corresponds to a common subforest of size  $k$ .

Every largest common subforest corresponds to a path of weight  $\text{LCS}(F, G)$ .

# Alignment graph (Backofen et al., Touzet)



There are exactly  $r$  edges with nonzero weight.

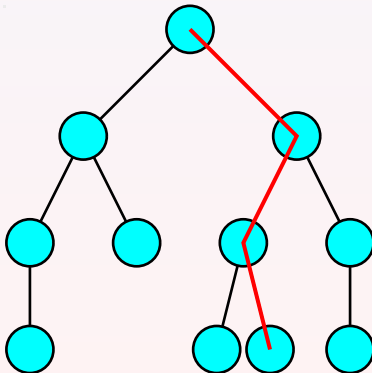
Computing the maximum weight path in the graph takes  $O(r \log \log m)$  time.

The time including the recursive calls is  $O(r \cdot \text{height}(F) \cdot \text{height}(G) \cdot \log \log m)$

# Heavy path

## Definition

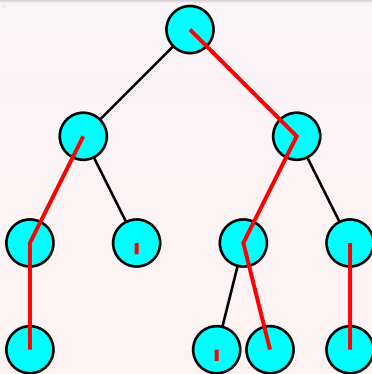
A **heavy path** of a tree is a path that starts at the root and goes from each vertex  $v$  to the child of  $v$  whose subtree contains more vertices.



# Heavy path decomposition

## Definition

A **heavy path decomposition** of a tree  $F$  is a collection of paths obtained by taking a heavy path of  $F$ , removing its vertices, and then building a heavy path decomposition for every remaining tree.



# The algorithm of Klein

$R_S$  = rightmost tree in  $S$ .  $r_S$  = root of  $R_S$ .

$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$



# The algorithm of Klein

$R_S$  = rightmost tree in  $S$ .  $r_S$  = root of  $R_S$ .

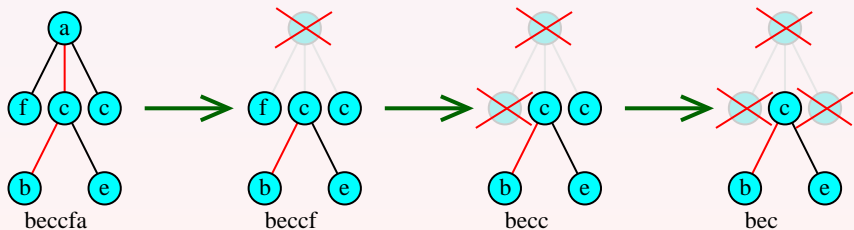
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G), \\ \text{LCS}(R_F - r_F, R_G - r_G) \\ \quad + \text{LCS}(F - R_F, G - R_G) + 1 \\ \quad \text{if } \text{label}(r_F) = \text{label}(r_G) \end{array} \right\}$$

$L_S$  = leftmost tree in  $S$ .  $l_S$  = root of  $L_S$ .

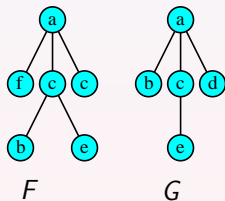
$$\text{LCS}(F, G) = \max \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G), \\ \text{LCS}(L_F - l_F, L_G - l_G) \\ \quad + \text{LCS}(F - L_F, G - L_G) + 1 \\ \quad \text{if } \text{label}(l_F) = \text{label}(l_G) \end{array} \right\}$$

# The algorithm of Klein

- Let  $F$  and  $G$  be the input trees, with  $|F| \geq |G|$ .
- Build a heavy path decomposition of  $F$ .
- Recursively compute  $LCS(F, G)$  using left & right rules.
- To compute  $LCS(F', G')$  for some  $F'$  and  $G'$ , let  $p$  be the highest path in decomposition that contain vertices of  $F'$ . Use the left rule if  $l_F$  is not on  $p$ . O/w use left rule.

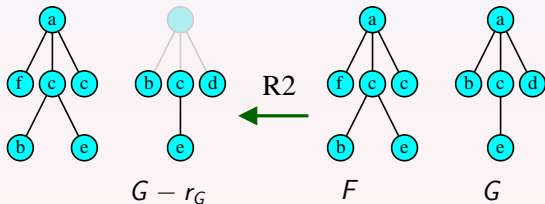


# Relevant subproblems



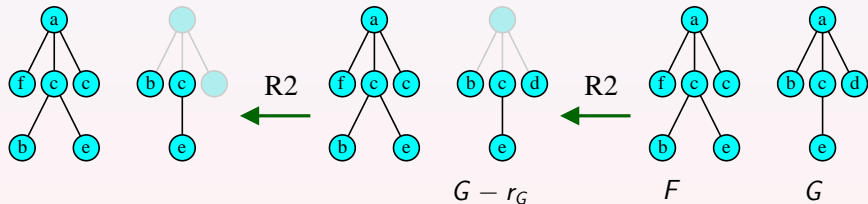
$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \end{array} \quad \text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \end{array}$$

# Relevant subproblems



$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \end{array} \quad \text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \end{array}$$

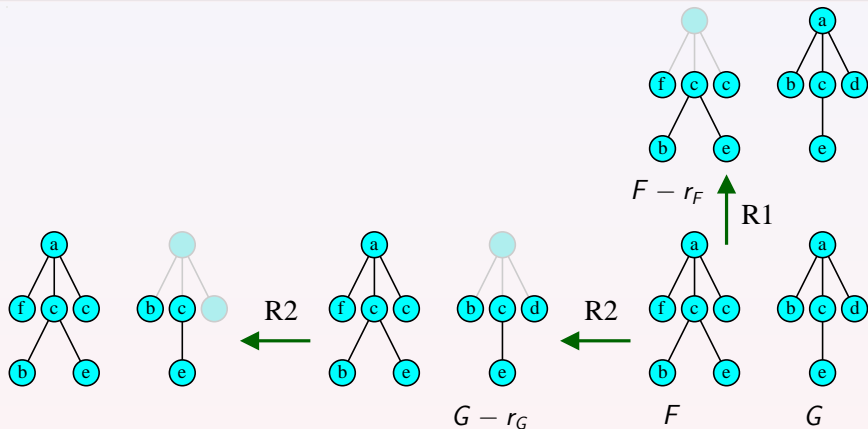
# Relevant subproblems



$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \\ \dots \end{array}$$

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \\ \dots \end{array}$$

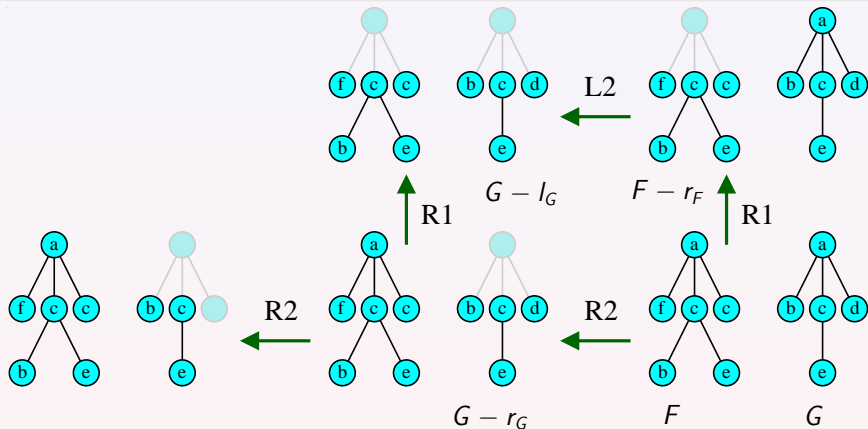
# Relevant subproblems



$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \\ \dots \end{array}$$

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \\ \dots \end{array}$$

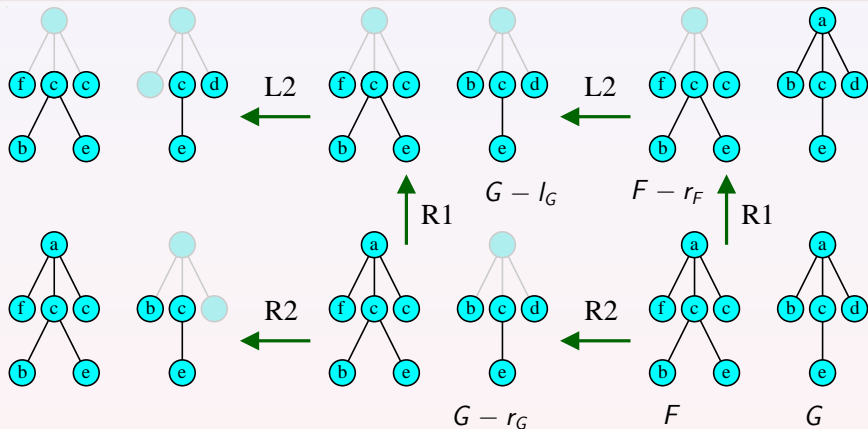
# Relevant subproblems



$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \\ \dots \end{array}$$

$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \\ \dots \end{array}$$

# Relevant subproblems

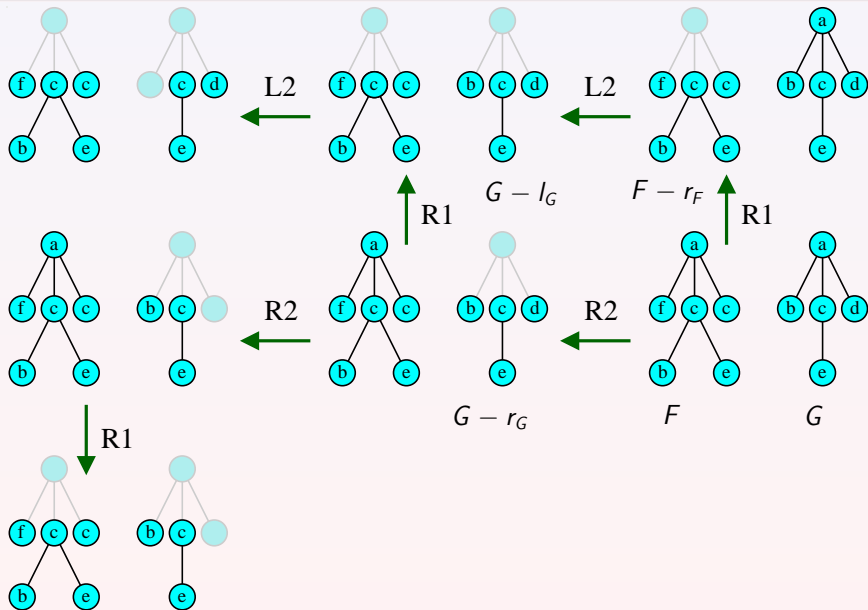


$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - l_F, G), \\ \text{LCS}(F, G - l_G) \\ \dots \end{array} \right\} \begin{array}{l} L1 \\ L2 \end{array}$$

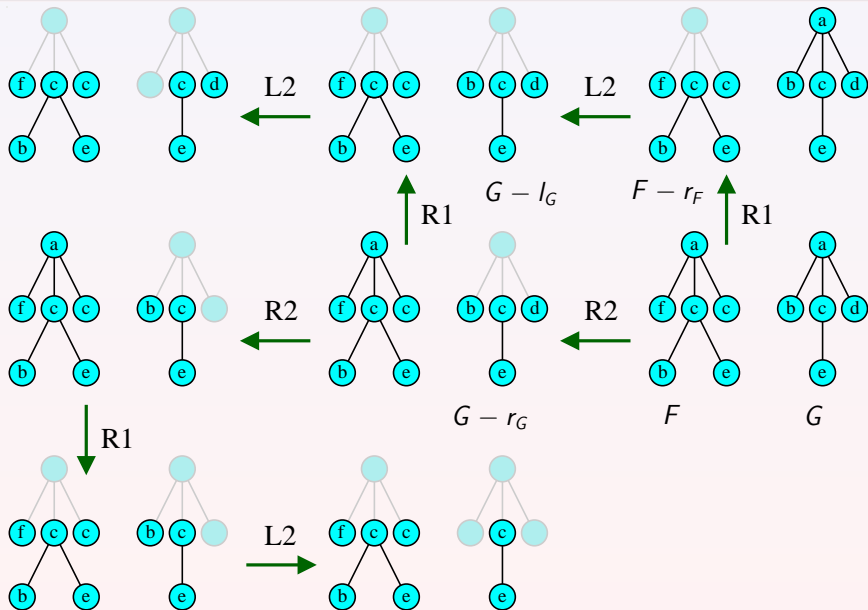
$$\text{LCS}(F, G) = \min \left\{ \begin{array}{l} \text{LCS}(F - r_F, G), \\ \text{LCS}(F, G - r_G) \\ \dots \end{array} \right\} \begin{array}{l} R1 \\ R2 \end{array}$$



# Relevant subproblems

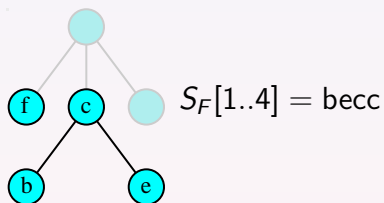
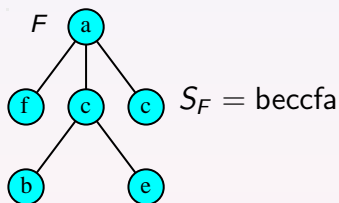


# Relevant subproblems

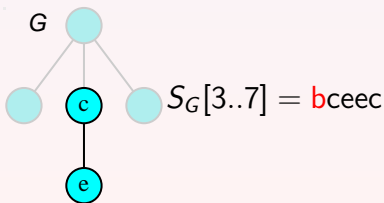
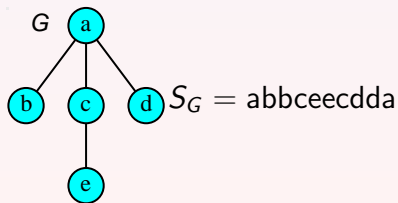


# Relevant subproblems

- $F$  represented by a string  $S_F$  according to deletion order.



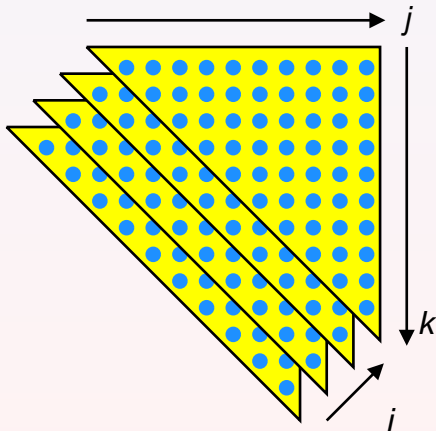
- $G$  represented by its Euler string  $S_G$ .



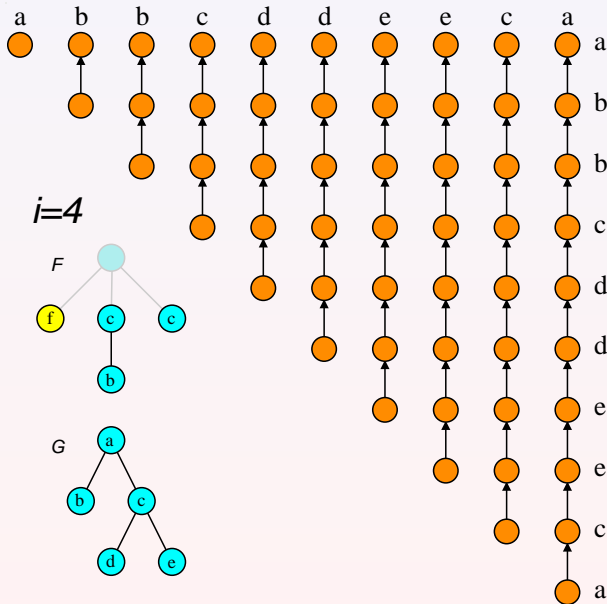
# Alignment graph

The alignment graph has vertex  $(i, j, k)$  for every  $0 \leq i \leq n$  and  $1 \leq j \leq k \leq 2m$ .

Vertex  $(i, j, k)$  corresponds to the pair  $S_F[1..i], S_G[j..k]$



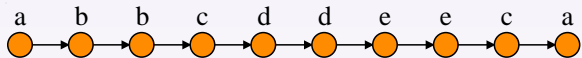
# Alignment graph



a Edges corresponding to deleting leftmost root of  $S_G[j..k]$ .

b Weights: 0

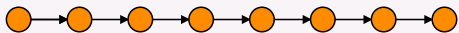
# Alignment graph



a Edges corresponding to deleting rightmost root of  $S_G[j..k]$ .



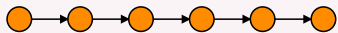
b



b Weights: 0



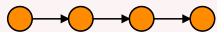
c



d



d



e



e

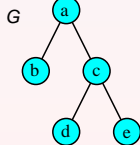
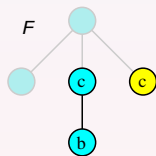


c

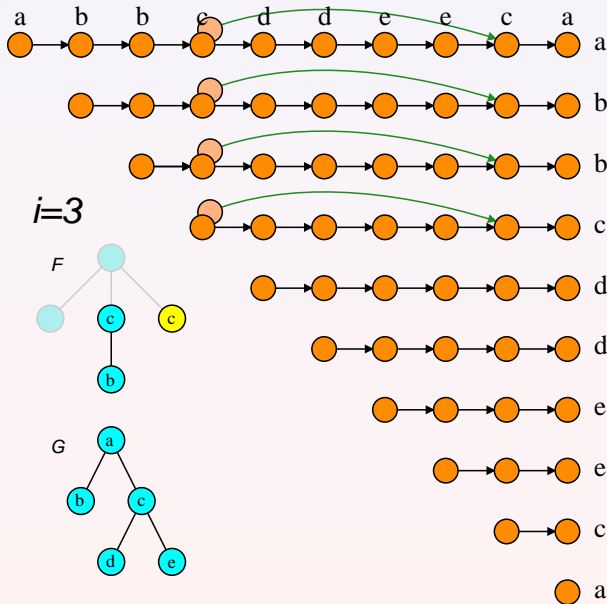


a

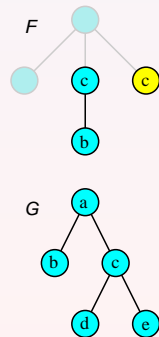
$i=3$



# Alignment graph



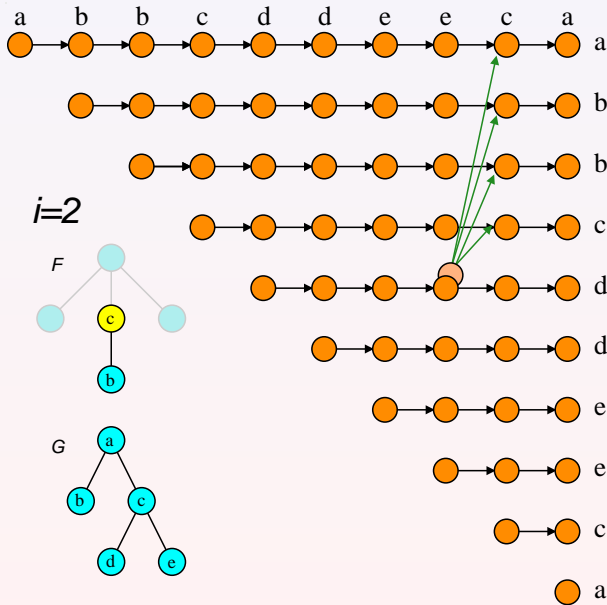
$i=3$



a Edges corresponding to matching rightmost roots of  $S_F[1..i]$  and  $S_G[j..k]$ .

Weights:  $\geq 1$  (LCS)

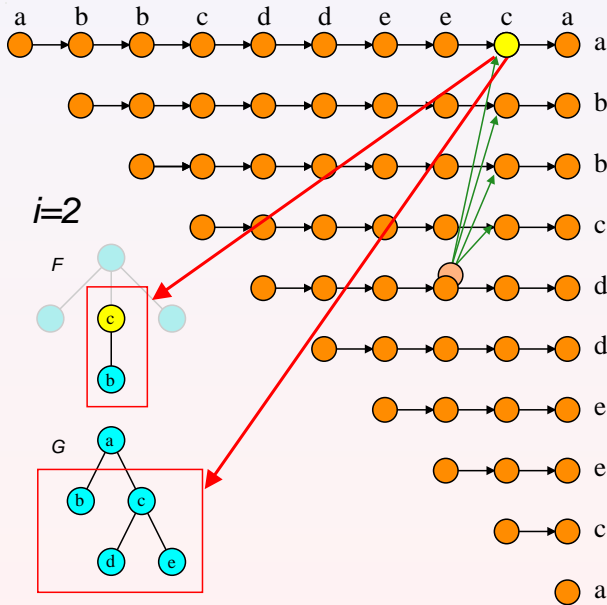
# Alignment graph



a If  $S_F[1..i]$  is a tree, the edges are different.  
Weights: 1



# Alignment graph



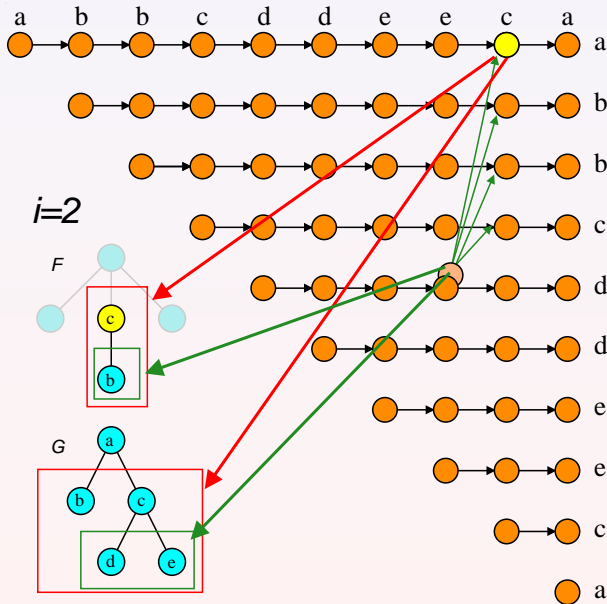
$a$  If  $S_F[1..i]$  is a tree, the edges are different.  
Weights: 1

$i=2$

$F$

$G$

# Alignment graph



a If  $S_F[1..i]$  is a tree, the edges are different.  
Weights: 1

# Summary

Algorithms for tree LCS:

- $O(r \cdot \text{height}(F) \cdot \text{height}(G) \cdot \log \log m)$
- $O(Lr \log r \cdot \log \log m)$

Algorithm for homeomorphic tree LCS

- $O(r \lg \lg m \cdot (\text{height}(F) + \text{height}(G)))$ .