

# Space Efficient Algorithms for Ordered Tree Comparison

Lusheng Wang · Kaizhong Zhang

Received: 27 January 2006 / Accepted: 27 June 2006 / Published online: 15 January 2008  
© Springer Science+Business Media, LLC 2008

**Abstract** In this paper we present techniques to significantly improve the space complexity of several ordered tree comparison algorithms without sacrificing the corresponding time complexity. We present new algorithms for computing the constrained ordered tree edit distance and the alignment of (ordered) trees. The techniques can also be applied to other related problems.

**Keywords** Space efficient algorithms · Constrained tree edit distance · Alignment of trees

## 1 Introduction

Ordered labeled trees are trees whose nodes are labeled and in which the left to right order among siblings is significant. Comparing such trees with distance and/or similarity measures has applications in several diverse areas such as computational molecular biology [5, 8, 14], computer vision, pattern recognition, programming compilation and natural language processing [2].

Algorithms have been developed for ordered labeled tree comparisons [5, 9, 10, 12, 15]. The degree-one edit distance was introduced by Selkow [9] in which insertions and deletions are restricted to the leaves of the trees. The degree-two edit distance was introduced by Zhang et al. [13, 16] in which insertions and deletions

---

L. Wang (✉)  
Dept. of Computer Science, City University of Hong Kong, 83 Tat Chee Avenue, Kowloon,  
Hong Kong  
e-mail: [lwang@cs.cityu.edu.hk](mailto:lwang@cs.cityu.edu.hk)

K. Zhang  
Dept. of Computer Science, University of Western Ontario, London, ON N6A 5B7, Canada  
e-mail: [kzhang@csd.uwo.ca](mailto:kzhang@csd.uwo.ca)

are restricted to tree nodes with zero or one child. Zhang [12] introduced the constrained edit distance which is a bit more general than degree-two edit distance. The (general) edit distance between ordered labeled trees was introduced by Tai [10]. His algorithm has been improved by several authors [3, 6, 15]. The alignment of trees was introduced by Jiang et al. in [5].

Given two ordered trees  $T_1$  and  $T_2$ , the degree-one edit distance, the degree-two edit distance, and the constrained edit distance can all be computed in  $(|T_1||T_2|)$  time and space, where  $|T|$  is the number of nodes in  $T$  [9, 12, 13]. Richter [7] presented an algorithm for the constrained edit distance with  $O(|T_1||T_2|\deg(T_1)\deg(T_2))$  time and  $O(|T_1|\deg(T_2)\deg(T_2))$  space, where  $\deg(T)$  denotes the depth of  $T$  and  $\deg(T)$  denotes the degree of  $T$ . For small degree and low depth trees, this is a space improvement. According to a recent survey on tree edit distance by Bille [1], the algorithms of Zhang [12] and Richter [7] are currently the best for the constrained tree edit distance. In this paper, we present an algorithm for the constrained tree edit distance with  $O(|T_1||T_2|)$  time and  $O(\log(|T_1||T_2|))$  space. The techniques used can also be used for the degree-one and the degree-two edit distance to achieve the same time and space complexities.

For alignment of trees, the algorithm in [5] runs in  $O(|T_1||T_2|(\deg(T_1) + \deg(T_2))^2)$  time and needs  $O(|T_1||T_2|(\deg(T_1) + \deg(T_2)))$  space. Recently, there is a strike to reduce the space [11]. The space required for the new algorithm is  $O(\log(|T_1||T_2|(\deg(T_1) + \deg(T_2))\deg(T_1)))$ . However, the running time is increased to  $O(|T_1|^2|T_2|(\deg(T_1) + \deg(T_2))^2)$ . In this paper, we proposed a new algorithm that runs in time  $O(|T_1||T_2|(\deg(T_1) + \deg(T_2))^2)$  and requires the same space as in [11]  $O(\log(|T_1||T_2|(\deg(T_1) + \deg(T_2))\deg(T_1)))$ .

## 2 Constrained Edit Distance between Ordered Trees

### 2.1 Notations

The nodes in an ordered tree of size  $n$  are numbered from 1 to  $n$  according to the postorder, where the left siblings are ordered before the right siblings. Given an ordered labeled tree  $T$ , the  $i$ th node of tree  $T$  is represented as  $t[i]$ , the subtree of  $T$  rooted at node  $t[i]$  is represented as  $T[i]$ , and the subforest obtained by deleting  $t[i]$  from  $T[i]$  is represented as  $F[i]$ . The parent of node  $t[i]$  in  $T$  is denoted as  $t[p(i)]$ . The number of nodes in a tree  $T$  is denoted as  $|T|$ .

Let  $\theta$  be the empty tree, and  $\lambda$  a space.  $\gamma(a, \lambda)$ ,  $\gamma(\lambda, a)$ , and  $\gamma(a, b)$  denote the cost of deleting  $a$ , inserting  $a$ , and substituting  $a$  with  $b$ , respectively.

Consider two trees  $T_1$  and  $T_2$  to compare. Suppose that the degrees of node  $t_1[i]$  and node  $t_2[j]$  are  $m_i$  and  $n_j$ , respectively. Denote the children of  $t_1[i]$  from left to right as  $t_1[i_1], \dots, t_1[i_{m_i}]$  and children of  $t_2[j]$  from left to right as  $t_2[j_1], \dots, t_2[j_{n_j}]$ .

The constrained edit distance metric is based on a constrained edit mapping allowed between two trees. We will first give the definition of constrained editing mapping and then use it to define the constrained editing distance metric.

Formally we define a triple  $(M, T_1, T_2)$  to be a constrained edit mapping from  $T_1$  to  $T_2$ , where  $M$  is any set of pairs of integers  $(i, j)$  satisfying:

- (1)  $1 \leq i \leq |T_1|, 1 \leq j \leq |T_2|$ ;
- (2) For any pair of  $(i_1, j_1)$  and  $(i_2, j_2)$  in  $M$ ,
  - (a)  $i_1 = i_2$  iff  $j_1 = j_2$  (one-to-one)
  - (b)  $t_1[i_1]$  is to the left of  $t_1[i_2]$  iff  $t_2[j_1]$  is to the left of  $t_2[j_2]$  (sibling order preserved);
  - (c)  $t_1[i_1]$  is an ancestor of  $t_1[i_2]$  iff  $t_2[j_1]$  is an ancestor of  $t_2[j_2]$  (ancestor order preserved);
- (3) For any triple  $(i_1, j_1), (i_2, j_2)$  and  $(i_3, j_3)$  in  $M$ , let  $\text{lca}()$  represent least common ancestor function,  $t_1[\text{lca}(i_1, i_2)]$  is a proper ancestor of  $t_1[i_3]$  iff  $t_2[\text{lca}(j_1, j_2)]$  is a proper ancestor of  $t_2[j_3]$ .

This definition adds constraint (3) to the definition of the edit mapping [15]. This is why we call it a constrained edit mapping. The intuitive idea behind this definition is that two separate subtrees of  $T_1$  should be mapped to two subtrees of  $T_2$  and vice versa.

We will use  $M$  instead of  $(M, T_1, T_2)$  if there is no confusion. Let  $M$  be a constrained editing mapping from  $T_1$  to  $T_2$ . We can define the cost of  $M$ :

$$\gamma(M) = \sum_{(i,j) \in M} \gamma(t_1[i], t_2[j]) + \sum_{i \notin M} \gamma(t_1[i], \lambda) + \sum_{j \notin M} \gamma(\lambda, t_2[j]).$$

Let  $F_1[p]$  and  $F_2[q]$  be forests of  $T_1$  and  $T_2$ . The cost of  $M$  with restrict to  $F_1[p]$  and  $F_2[q]$  is

$$\begin{aligned} \gamma(M, F_1[p], F_2[q]) = & \sum_{(i,j) \in M \& t_1[i] \in F_1[p] \& t_2[j] \in F_2[q]} \gamma(t_1[i], t_2[j]) \\ & + \sum_{i \notin M \& t_1[i] \in F_1[p]} \gamma(t_1[i], \lambda) + \sum_{j \notin M \& t_2[j] \in F_2[q]} \gamma(\lambda, t_2[j]). \end{aligned}$$

The constrained edit distance between  $T_1$  and  $T_2$  is defined as:

$$D(T_1, T_2) = \min_M \{\gamma(M) \mid M \text{ is a constrained edit mapping between } T_1 \text{ to } T_2\}.$$

Similarly, the constrained edit distance between  $F_1[p]$  and  $F_2[q]$  is defined as

$$D(F_1[p], F_2[q]) = \min_M \{\gamma(M, F_1[p], F_2[q]) \mid M \text{ is a constrained edit mapping between } T_1 \text{ to } T_2\}.$$

## 2.2 A Simple Algorithm

We now present an algorithm for computing the constrained edit distance between two ordered trees. This algorithm is given in [12] and is the basis of our new space efficient algorithm.

The algorithm in [12] for the constrained edit distance between two ordered trees is given in Fig. 1. While we do not show the details of the correctness of the algorithm,

Input:  $T_1$  and  $T_2$

Output:  $D(T_1[i], T_2[j])$ , where  $1 \leq i \leq |T_1|$  and  $1 \leq j \leq |T_2|$

$D(\theta, \theta) = 0$ ;

**for**  $i = 1$  **to**  $|T_1|$

$D(F_1[i], \theta) = \sum_{k=1}^{n_i} D(T_1[i_k], \theta)$

$D(T_1[i], \theta) = D(F_1[i], \theta) + \gamma(t_1[i], \lambda)$

**for**  $j = 1$  **to**  $|T_2|$

$D(\theta, F_2[j]) = \sum_{k=1}^{n_j} D(\theta, T_2[j_k])$

$D(\theta, T_2[j]) = D(\theta, F_2[j]) + \gamma(\lambda, t_2[j])$

**for**  $i = 1$  **to**  $|T_1|$

**for**  $j = 1$  **to**  $|T_2|$

Compute  $E(m_i, n_j)$

$$D(F_1[i], F_2[j]) = \min \begin{cases} D(\theta, F_2[j]) + \min_{1 \leq t \leq n_j} \{D(F_1[i], F_2[j_t]) - D(\theta, F_2[j_t])\} \dots (2.1) \\ D(F_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(F_1[i_s], F_2[j]) - D(F_1[i_s], \theta)\} \dots (2.2) \\ E(m_i, n_j) \dots (2.3) \end{cases}$$

$$D(T_1[i], T_2[j]) = \min \begin{cases} D(\theta, T_2[j]) + \min_{1 \leq t \leq n_j} \{D(T_1[i], T_2[j_t]) - D(\theta, T_2[j_t])\} \dots (3.1) \\ D(T_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(T_1[i_s], T_2[j]) - D(T_1[i_s], \theta)\} \dots (3.2) \\ D(F_1[i], F_2[j]) + \gamma(t_1[i], t_2[j]) \dots (3.3) \end{cases}$$

**Fig. 1** A simple algorithm, Algorithm 1

we now explain the ideas of the algorithm. Let  $t_1[i]$  and  $t_2[j]$  be two nodes.  $T_1[i]$  and  $T_2[j]$  are the two subtrees rooted at  $t_1[i]$  and  $t_2[j]$ , respectively. We use  $m_i$  and  $n_j$  to denote the numbers of children for  $t_1[i]$  and  $t_2[j]$ , respectively.

The computation of  $D(T_1[i], T_2[j])$  has several cases. In one case,  $T_1[i]$  is matched to a child subtree of  $T_2[j]$ . See (3.1) in Fig. 1. In a symmetric case,  $T_2[j]$  is matched to a child subtree of  $T_1[i]$ . See (3.2) in Fig. 1. We denote

$$G_i = D(T_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(T_1[i_s], T_2[j]) - D(T_1[i_s], \theta)\}$$

and will use it in the new algorithm later. In the last case,  $t_1[i]$  is matched with  $t_2[j]$  and therefore  $F_1[i]$  is matched with  $F_2[j]$  which means that we need to know  $D(F_1[i], F_2[j])$ . See (3.3) in Fig. 1.

Similarly, the computation of  $D(F_1[i], F_2[j])$  has several cases. In the first case,  $F_1[i]$  is matched to  $F_2[j_t]$ ,  $1 \leq t \leq n_j$ , a subforest of a child of  $F_2[j]$ . See (2.1) in Fig. 1. In the second case,  $F_2[j]$  is matched to  $F_1[i_s]$ ,  $1 \leq s \leq m_i$ , a subforest of a child of  $F_1[i]$ . See (2.2) in Fig. 1. We denote

$$F_i = D(F_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(F_1[i_s], F_2[j]) - D(F_1[i_s], \theta)\}.$$

In the third case, there is a matching between  $T_1[i_1], \dots, T_1[i_{m_i}]$  and  $T_2[j_1], \dots, T_2[j_{n_j}]$ . The optimal matching between these subtrees are computed using the sequence edit distance algorithm treating each subtree as an unit. For a given pair of nodes  $t_1[i]$  and  $t_2[j]$ ,  $E(m_i, n_j)$  is the cost for an optimal matching between  $T_1[i_1], \dots, T_1[i_{m_i}]$  and  $T_2[j_1], \dots, T_2[j_{n_j}]$ . The code for the computation of  $E(m_i, n_j)$  is shown in Fig. 2.

**Fig. 2** The code to compute  $E[m_i, n_j]$

```

 $c(0, 0) = 0$ 
for  $s = 1$  to  $m_i$ 
     $c(s, 0) = c(s - 1, 0) + D(T_1[i_s], \theta)$ 
for  $t = 1$  to  $n_j$ 
     $c(0, t) = c(0, t - 1) + D(\theta, T_2[j_t])$ 
for  $s = 1$  to  $m_i$ 
    for  $t = 1$  to  $n_j$ 
        
$$c(s, t) = \min \begin{cases} c(s, t - 1) + D(\theta, T_2[j_t]) \\ c(s - 1, t) + D(T_1[i_s], \theta) \\ c(s - 1, t - 1) + D(T_1[i_s], T_2[j_t]) \end{cases}$$

 $E[m_i, n_j] = c(m_i, n_j).$ 

```

The time and space complexity for computing  $E(m_i, n_j)$  is obviously  $O(m_i \times n_j)$ . The complexity of computing  $D(T_1[i], T_2[j])$  is bounded by  $O(m_i + n_j)$ . The complexity of computing  $D(F_1[i], F_2[j])$  is bounded by  $O(m_i + n_j)$ .

Hence for any pair  $i$  and  $j$ , the complexity of computing  $D(T_1[i], T_2[j])$  and  $D(F_1[i], F_2[j])$  is bounded by  $O(m_i \times n_j)$ . Therefore the complexity of the algorithm is

$$\sum_{i=1}^{|T_1|} \sum_{j=1}^{|T_2|} O(m_i \times n_j) = O\left(\sum_{i=1}^{|T_1|} m_i \times \sum_{j=1}^{|T_2|} n_j\right) = O(|T_1| \times |T_2|).$$

### 2.3 Reducing the Space Complexity of Algorithm 1

In practice, the space required for Algorithm 1 might be a bottleneck. In this section, we propose a method to modify Algorithm 1 so that the space required is reduced to  $O(\log(|T_1|) \cdot |T_2|)$ .

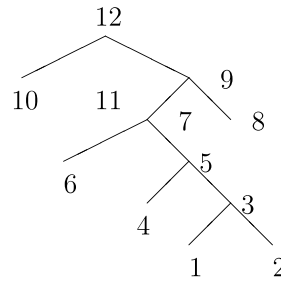
The basic idea is straightforward. In order to compute the constrained edit distance, some computed values are no longer useful after they were used. We simply release the space that are no longer useful. To achieve our goal, we use two tricks: (1) We use a new order of nodes for  $T_1$ ; and (2) We also modify the computation of node  $t_1[i]$  slightly. Let  $t_1[i_1], t_1[i_2], \dots, t_1[i_{m_i}]$  be the children of  $t_1[i]$ . When the computation for node  $t_1[i_1]$  is completed, we immediately start the computation of  $E[m_i, n_j]$ ,  $F_i$  and  $G_i$  without waiting for the computation of  $t_1[i_2]$ .

*Using a New Order of Nodes in  $T_1$*  Here we give a more restricted order for the nodes in  $T_1$ . Consider a node  $t_1[i]$  in  $T_1$ . Let  $t_1[i_1], t_1[i_2], \dots, t_1[i_{m_i}]$  be the children of  $t_1[i]$  in  $T_1$  and  $|T_1[i_k]|$  be the number of nodes in the subtree  $T_1[i_k]$ . The *computational order* of the children of  $t_1[i]$  is that the child with largest number of nodes, which we refer to as the *favorable* child, would be the first. For the rest of children of  $t_1[i]$ , the order is from left to right. This order is applied to all the nodes and if  $t_1[u]$  and  $t_1[v]$  are siblings in  $T_1$  and  $t_1[u]$  is ordered before  $t_1[v]$  then all nodes in  $T_1[u]$  are ordered before any node in  $T_1[v]$ . Figure 3 gives an example of the new order.

Let us consider a new algorithm, called Algorithm 1.1. Algorithm 1.1 is identical to Algorithm 1 except that the order for  $T_1$  is new. Figure 4 gives the sketch.

Considering the computation for nodes  $t_1[i]$  in  $T_1$  and  $t_2[j]$  in  $T_2$ . In order to compute  $E[m_i, n_j]$ , we have to keep the computed values for nodes  $t_1[i_1], t_1[i_2]$ ,

**Fig. 3** An example of the new order



Input:  $T_1$  and  $T_2$

Output:  $D(T_1[i], T_2[j])$ , where  $1 \leq i \leq |T_1|$  and  $1 \leq j \leq |T_2|$

**for**  $i = 1$  **to**  $|T_1|$  (using the new order)

**for**  $j = 1$  **to**  $|T_2|$

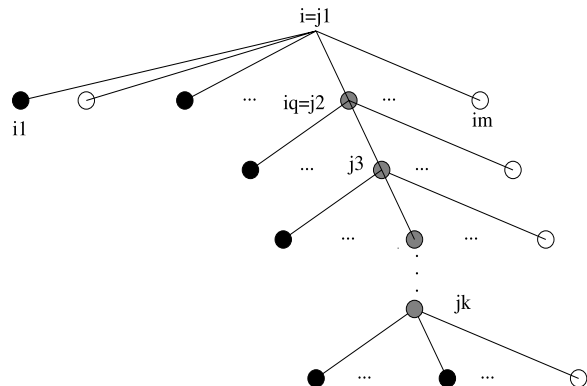
        Compute  $E(m_i, n_j)$

$$D(F_1[i], F_2[j]) = \min \begin{cases} D(\theta, F_2[j]) + \min_{1 \leq t \leq n_j} \{D(F_1[i], F_2[j_t]) - D(\theta, F_2[j_t])\} \dots (2.1) \\ D(F_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(F_1[i_s], F_2[j]) - D(F_1[i_s], \theta)\} \dots (2.2) \\ E(m_i, n_j) \dots (2.3) \end{cases}$$

$$D(T_1[i], T_2[j]) = \min \begin{cases} D(\theta, T_2[j]) + \min_{1 \leq t \leq n_j} \{D(T_1[i], T_2[j_t]) - D(\theta, T_2[j_t])\} \dots (3.1) \\ D(T_1[i], \theta) + \min_{1 \leq s \leq m_i} \{D(T_1[i_s], T_2[j]) - D(T_1[i_s], \theta)\} \dots (3.2) \\ D(F_1[i], F_2[j]) + \gamma(t_1[i], t_2[j]) \dots (3.3) \end{cases}$$

**Fig. 4** Algorithm 1.1

**Fig. 5** The nodes required to be stored in order to compute the values on node  $t_1[i]$  in  $T_1$



$\dots, t_1[i_{m_i}]$ . To get the values on  $t_1[i_k]$ , we have to know the values on the children of  $t_1[i_k]$ . In general, in the computation for node  $t_1[i]$ , we have to store all the dark nodes as shown in Fig. 5. Since we use the new order, we know that the height of the tree in Fig. 5 is at most  $O(\log |T_1|)$ . Thus, the space required is  $O((\log |T_1|) \deg(T_1) |T_2|)$ , where  $\deg(T_1)$  is the degree of  $T_1$ . To achieve a better space complexity, we use the second trick.

Input:  $T_1$  and  $T_2$

Output:  $D(T_1, T_2[j])$ , where  $1 \leq j \leq |T_2|$

Initialization of variables

**for**  $i = 1$  **to**  $|T_1|$  (the new order)

**for**  $j = 1$  **to**  $|T_2|$

$$D(F_1[i], F_2[j]) = \min \begin{cases} D(\theta, F_2[j]) + \min_{1 \leq t \leq n_j} \{D(F_1[i], F_2[j_t]) - D(\theta, F_2[j_t])\} \\ F_i \\ E(m_i, n_j) \end{cases}$$

$$D(T_1[i], T_2[j]) = \min \begin{cases} D(\theta, T_2[j]) + \min_{1 \leq t \leq n_j} \{D(T_1[i], T_2[j_t]) - D(\theta, T_2[j_t])\} \\ G_i \\ D(F_1[i], F_2[j]) + \gamma(t_1[i], t_2[j]) \end{cases}$$

Let  $t_1[p(i)]$  be the parent of  $t_1[i]$  and  $t_1[p_f]$  be the favorable child of  $t_1[p(i)]$

**Case 1:** if  $t_1[i]$  is the favorable child of  $t_1[p(i)]$  **then**

$$E_{p(i)}(0) = 0;$$

$$F_{p(i)} = D(F_1[p(i)], \theta) + D(F_1[i], F_2[j]) - D(F_1[i], \theta);$$

$$G_{p(i)} = D(T_1[p(i)], \theta) + D(T_1[i], T_2[j]) - D(T_1[i], \theta);$$

**for**  $t = 1$  **to**  $n_j$

$$E_{p(i)}(t) = E_{p(i)}(t-1) + D(\theta, T_2[j_t]);$$

**Case 2:** if  $t_1[i]$  is not the favorable child of  $t_1[p(i)]$  or  $t_1[i]$  is the leftmost child of  $t_1[p(i)]$  **then**

$$F_{p(i)} = \min\{F_p, D(F_1[p(i)], \theta) + D(F_1[i], F_2[j]) - D(F_1[i], \theta)\}$$

$$G_{p(i)} = \min\{G_p, D(T_1[p(i)], \theta) + D(T_1[i], T_2[j]) - D(T_1[i], \theta)\}$$

**for**  $t = 0$  **to**  $n_j$

$$E_{p(i)}^0(t) = E_p(t);$$

$$E_{p(i)}(0) = E_{p(i)}^0(0) + D(T_1[i], \theta);$$

**for**  $t = 1$  **to**  $n_j$

$$E_{p(i)}(t) = \min \begin{cases} E_{p(i)}(t-1) + D(\theta, T_2[j_t]) \\ E_{p(i)}^0(t) + D(T_1[i], \theta) \\ E_{p(i)}^0(t-1) + D(T_1[i], T_2[j_t]) \end{cases}$$

**Case 3:** if  $t_1[i]$  is the left sibling of the favorable child  $t_1[p_f]$  of  $t_1[p(i)]$  **then**

**for**  $t = 0$  **to**  $n_j$

$$E_{p(i)}^0(t) = E_{p(i)}(t);$$

$$E_{p(i)}(0) = E_{p(i)}^0(0) + D(T_1[p_f], \theta);$$

**for**  $t = 1$  **to**  $n_j$

$$E_{p(i)}(t) = \min \begin{cases} E_{p(i)}(t-1) + D(\theta, T_2[j_t]) \\ E_{p(i)}^0(t) + D(T_1[p_f], \theta) \\ E_{p(i)}^0(t-1) + D(T_1[p_f], T_2[j_t]) \end{cases}$$

**Fig. 6** A more space efficient algorithm, Algorithm 2

**Changing the Computation on Node  $t_1[i]$**  Now, we slightly change the computation on node  $t_1[i]$ . When the computation for one of the children of  $t_1[i]$ , say,  $t_1[i_1]$ , is completed, we immediately start the computation of  $E[m_i, n_j]$ ,  $F_i$  and  $G_i$  without waiting for the computation of  $t_1[i_2]$ . In this way, we can avoid storing the values on the dark nodes as in Fig. 5. Instead, we only have to keep the values on the grey nodes in the path from  $t_1[i]$  to the bottom of the tree. Thus, the space complexity is reduced to be  $O(\log |T_1| \times |T_2|)$ .

The modified algorithm, Algorithm 2, is shown in Fig. 6. In the algorithm, we will use this new order for  $T_1$  and the same post order as in Algorithm 1 for  $T_2$ . For each

node  $t_1[i]$ , we will first compute  $D(F_1[i], F_2[j])$  and  $D(T_1[i], T_2[j])$  for all nodes  $t_2[j]$  in  $T_2$  assuming  $E[m_i, n_j]$ ,  $F_i$ , and  $G_i$  have been computed.

*Computation of  $E[m_i, n_j]$ ,  $F_i$  and  $G_i$*  Suppose that  $t_1[p(i)]$  is the parent of  $t_1[i]$ . We can immediately use  $D(F_1[i], F_2[j])$  and  $D(T_1[i], T_2[j])$  for the computation of  $E[m_{p(i)}, n_j]$ ,  $F_{p(i)}$  and  $G_{p(i)}$ . There are several cases.

**Case 1:**  $t_1[i]$  is the favorable child of  $t_1[p(i)]$  and it is not the leftmost child of  $t_1[p(i)]$ . In this case, we can start computing  $F_{p(i)}$  and  $G_{p(i)}$  by setting

$$F_{p(i)} = D(F_1[p(i)], \theta) + D(F_1[i], F_2[j]) - D(F_1[i], \theta)$$

and

$$G_{p(i)} = D(T_1[p(i)], \theta) + D(T_1[i], T_2[j]) - D(T_1[i], \theta).$$

We also start the computation of  $E[m_{p(i)}, n_j]$ . The computation for  $E[m_i, n_j]$  must be from left to right (see Fig. 1) and thus we have to keep the values on node  $t_1[i]$  for later use since  $t_1[i]$  is not the leftmost child of  $t_1[p(i)]$ . Again, as shown in Fig. 1, computing  $E[m_i, n_j]$  needs the values for the  $m_i \times n_j$  cells  $c(s, t)$  for  $s = 1, 2, \dots, m_i$  and  $j = 1, 2, \dots, n_j$ . There are two indices  $s$  and  $t$ . Here in the new algorithm, we also need the values for the  $m_i \times n_j$  cells. However, we do not simultaneously store all the  $m_i \times n_j$  cells. To get the value of  $c(s, t)$ , we need the values of the three cells  $c(s-1, t-1)$ ,  $c(s-1, t)$  and  $c(s, t-1)$ . We use  $E_{p(i)}^0(t)$  to store  $c(s-1, t)$  and use  $E_{p(i)}(t)$  to store  $c(s, t)$ . In this way, we do not have to store all  $c(1, t)$ ,  $c(2, t)$ ,  $\dots$ ,  $c(m_i, t)$  simultaneously. Instead, we just need to keep  $E_{p(i)}^0(t)$  and  $E_{p(i)}(t)$ . Therefore at this stage, we set the initial values  $E_{p(i)}(0) = 0$  and  $E_{p(i)}(t) = E_{p(i)}(t-1) + D(\theta, T_2[j_t])$  for  $t = 1, 2, \dots, n_j$ .

**Case 2:**  $t_1[i]$  is not the favorable child or it is the favorable child and also the leftmost child. In this case, we can use it immediately for the computation of  $E(p(i), j)$ . We first pass the value  $E_{p(i)}(t)$  to  $E_{p(i)}^0(t)$  and then compute the new values for  $E_{p(i)}(t)$  as

$$E_p(0) = E_p^0(0) + D(T_1[i], \theta);$$

**for**  $t = 1$  **to**  $n_j$

$$E_{p(i)}(t) = \min \begin{cases} E_{p(i)}(t-1) + D(\theta, T_2[j_t]) \\ E_{p(i)}^0(t) + D(T_1[i], \theta) \\ E_{p(i)}^0(t-1) + D(T_1[i], T_2[j_t]) \end{cases}$$

Therefore for each node  $t_1[i]$  in  $T_1$  such that at least its favorable child's computation is completed, we at most will need to keep two sets of values: the values for its favorable child and the partial computation of  $E(m_i, n_j)$ . Notice that for any node  $t_1[i]$  in  $T_1$ , if the computation of its favorable child is not completed, we do not need to store any value and if the computation of all its children have been completed, then in the next step the computation for node  $t_1[i]$  is completed and we can release the space used.



**Fig. 7** Initialization of the variables

```

 $D(\theta, \theta) = 0;$ 
for  $i = 1$  to  $|T_1|$ 
     $D(F_1[i], \theta) = \sum_{k=1}^{n_i} D(T_1[i_k], \theta)$ 
     $D(T_1[i], \theta) = D(F_1[i], \theta) + \gamma(t_1[i], \lambda)$ 

for  $j = 1$  to  $|T_2|$ 
     $D(\theta, F_2[j]) = \sum_{k=1}^{n_j} D(\theta, T_2[j_k])$ 
     $D(\theta, T_2[j]) = D(\theta, F_2[j]) + \gamma(\lambda, t_2[j])$ 
    
```

**Case 3:**  $t_1[i]$  is the left sibling of the favorable child  $t_1[p_f]$  of  $t_1[p(i)]$ . The computation for node  $t_1[i]$  has been done in Case 2. Now, we do computation for node  $t_1[p_f]$  to calculate  $E(m_i, n_j)$ . See Fig. 6. By using such a new order for  $T_1$  and the above straightforward idea, we now have an algorithm with much less space. The initialization of variables is shown in Fig. 7. The modified algorithm is called Algorithm 2 and is given in Fig. 6.

The following lemma shows that the space complexity of Algorithm 2 is reduced to  $O(\log(|T_1|)|T_2|)$ .

**Lemma 1** *The time and space complexities of Algorithm 2 are  $O(|T_1||T_2|)$  and  $O(\log(|T_1|)|T_2|)$ .*

*Proof* From the algorithm, it is clear that for each  $i$  of  $T_1$ , the number of steps is bounded by  $O(\sum_{j=1}^{|T_2|} n_j)$ . Therefore the time complexity of Algorithm 2 is  $O(\sum_{i=1}^{|T_1|} \sum_{j=1}^{|T_2|} n_j) = O(|T_1||T_2|)$ .

For the space complexity, let us consider the status of a node in  $T_1$  with respect to the space requirement. We say a node  $t_1[i]$  in  $T_1$  is active if the computation of its favorable child is completed and the computation for node  $t_1[i]$  itself has not been completed. Therefore a node is inactive if the computation of its favorable child has not been completed or if the computation for node  $i$  is completed.

If a node is inactive because its computation is completed and this node is the favorable child of its parent, then we cannot immediately release the space used since these values will be used for the computation of its parent. However in this situation, we can contribute the space requirement of the favorable child to its parent, which is active, when counting the space requirement. Therefore, we can assume that if a node is inactive, then we can release the space used for that node. This means that we only have to check maximum space used by active nodes during the execution of the algorithm.

For any active node  $t_1[p]$ , assuming that its favorable child is  $t_1[p_f]$ , then we need  $O(|T_2|)$  to store  $D(F_1[p_f], F_2[j])$  and  $D(T_1[p_f], T_2[j])$  for  $1 \leq j \leq |T_2|$ , and  $O(\sum_{j=1}^{|T_2|} n_j) = O(|T_2|)$  to store  $E_{p_f}(t)$  where  $1 \leq t \leq n_j$  for all  $j$  in  $T_2$ . Therefore for each active node, the space required is  $O(|T_2|)$ . Because of the new order, if two nodes are active at the same time, then one has to be the ancestor of the other. Therefore all active nodes of  $T_1$  are on a path in  $T_1$ . Let  $t_1[s]$  and  $t_1[t]$  be two neighboring active nodes on this path and  $t_1[s]$  is an ancestor of  $t_1[t]$ , then  $|T_1[s]| \geq 2|T_1[t]|$  since the computation of  $t_1[s_f]$ , the favorable child of  $t_1[s]$  with maximum size, is already completed and therefore  $t_1[t]$  is not a descendant of  $t_1[s_f]$ .

This means that there are at most  $O(\log(|T_1|))$  active nodes. Therefore the space complexity is  $O(\log(|T_1|)|T_2|)$ .  $\square$

## 2.4 Finding the Optimal Constrained Edit Mapping between Two Trees

In previous section we show that  $D(T_1, T_2)$  can be computed in  $O(|T_1||T_2|)$  time and  $O(\log(|T_1|)|T_2|)$  space. However in real application the optimal mapping that achieves  $D(T_1, T_2)$  maybe required.

Finding the optimal mapping in  $O(\log(|T_1|)|T_2|)$  space is in fact a more difficult task. This is similar to the situation of computing edit distance between two sequences. Computing the edit distance in linear space is easy, but finding the optimal edit script in linear space is more involved. Hirschberg [4] presented a clever way to do this.

In this section we will present a method that can find the optimal mapping between two ordered trees in  $O(|T_1||T_2|)$  time and  $O(\log(|T_1|)|T_2|)$  space. When the input is two sequences (trees such that each node only has one child), then our method produces the optimal edit script for sequence edit distance in  $O(|T_1||T_2|)$  time and  $O(|T_2|)$  space.

The main idea is as the follows. Given  $T_1$  and  $T_2$ , there is a unique node, called key node,  $t_1[k]$  in  $T_1$  such that in  $O(|T_1||T_2|)$  time and  $O(\log(|T_1|)|T_2|)$  space we can determine, in the optimal mapping, what subtree and subforest in  $T_2$  that  $T_1[k]$  and  $F_1[k]$  would match to produce  $D(T_1, T_2)$ . With this information, we then decompose the optimal mapping into several components mapping such that for each component the subtree or subforest of  $T_1$  involved has a size less than or equal to half of  $|T_1|$ . This means that in the next step if we repeat the same process for each component then the total cost for all the components is  $O(0.5|T_1||T_2|)$  time using  $O(\log(|T_1|)|T_2|)$  space. Therefore, repeating this process at most  $\log(|T_1|)$  times, we can compute the optimal mapping in  $O(|T_1||T_2|)$  time and  $O(\log(|T_1|)|T_2|)$  space.

Given  $T_1$ , the unique key node  $t_1[k]$ , with children  $t_1[k_1], t_2[k_2], \dots, t_1[k_{m_k}]$ , is a node satisfying the following properties.

- $|T_1[k]| > 0.5|T_1|$ ,
- $|T_1[k_s]| \leq 0.5|T_1|$  for  $1 \leq s \leq m_k$ .

With a small modification of Algorithm 2, when computing  $D(T_1, T_2)$ , we can also compute two integers  $t_1$  and  $t_2$  for subtree  $T_1[k]$  and two integers  $f_1$  and  $f_2$  for subforest  $F_1[k]$ .

We now define the meaning of  $t_1$  and  $t_2$ . If in the optimal mapping  $T_1[k]$  is deleted, then  $t_1 = t_2 = 0$ . If in the optimal mapping subtree  $T_1[k]$  is matched with subtree  $T_2[l]$  such that in this mapping  $T_1[k_s]$  is matching with  $T_2[l]$  and the rest of  $T_1[k]$  is deleted, then  $t_1 = k_s$  and  $t_2 = l$ . If in the optimal mapping node  $t_1[k]$  is matched with node  $t_2[l]$ , then  $t_1 = k$  and  $t_2 = l$ . If in the optimal mapping  $t_1[k]$  is deleted and  $F_1[k]$  is matched with  $F_2[l]$ , then  $t_1 = k + 1$  and  $t_2 = l$ .

$f_1$  and  $f_2$  are defined similarly. Note that  $f_1$  and  $f_2$  are only meaningful when  $t_1 = k$  or  $t_1 = k + 1$ . Otherwise  $f_1 = f_2 = -1$ . If in the optimal mapping  $F_1[k]$  is deleted, then  $f_1 = f_2 = 0$ . If in the optimal mapping  $F_1[k]$  is matched with  $F_2[l]$  such that in this mapping  $F_1[k_s]$  is matching with  $F_2[l]$  and the rest of  $F_1[k]$  are

deleted, then  $f_1 = k_s$  and  $f_2 = l$ . If in the optimal mapping  $F_1[k]$  is matched with  $F_2[l]$  and  $D(F_1[k], F_2[l]) = E(m_k, n_l)$ , then  $f_1 = k$  and  $f_2 = l$ .

We now give the formulae for computing  $t_1$ ,  $t_2$ ,  $f_1$ , and  $f_2$ . Since those formulae are self explanatory, we omit their proofs. We use  $t\_T(i, j)$  to represent  $(t_1, t_2)$  in the optimal mapping of  $D(T_1[i], T_2[j])$ . We use  $f\_T(i, j)$  to represent  $(f_1, f_2)$  in the optimal mapping of  $D(T_1[i], T_2[j])$ . We use  $t\_F(i, j)$  to represent  $(t_1, t_2)$  in the optimal mapping of  $D(F_1[i], F_2[j])$ . We use  $f\_F(i, j)$  to represent  $(f_1, f_2)$  in the optimal mapping of  $D(F_1[i], F_2[j])$ .

**Lemma 2** *Let  $t_1[k]$  be the key node of  $T_1$ , then*

$$f\_F(k, j) = \begin{cases} (k, j) & \text{if } D(F_1[k], F_2[j]) = E(m_k, n_j), \\ (k_s, j) & \text{if } D(F_1[k], F_2[j]) = D(F_1[k], \theta) + D(F_1[k_s], F_2[j]) \\ & \quad - D(F_1[k_s], \theta), \\ f\_F(k, j_t) & \text{if } D(F_1[k], F_2[j]) = D(\theta, F_2[j]) + D(F_1[k], F_2[j_t]) \\ & \quad - D(\theta, F_2[j_t]), \end{cases}$$

$$t\_T(k, j) = \begin{cases} (k, j) & \text{if } D(T_1[k], T_2[j]) = D(F_1[k], F_2[j]) + \gamma(t_1[k], t_2[j]), \\ (k_s, j) & \text{if } D(T_1[k], T_2[j]) = D(T_1[k], \theta) + D(T_1[k_s], T_2[j]) \\ & \quad - D(T_1[k_s], \theta), \\ t\_T(k, j_t) & \text{if } D(T_1[k], T_2[j]) = D(\theta, T_2[j]) + D(T_1[k], T_2[j_t]) \\ & \quad - D(\theta, T_2[j_t]), \end{cases}$$

$$f\_T(k, j) = \begin{cases} f\_F(k, j) & \text{if } D(T_1[k], T_2[j]) = D(F_1[k], F_2[j]) + \gamma(t_1[k], t_2[j]), \\ (-1, -1) & \text{if } D(T_1[k], T_2[j]) = D(T_1[k], \theta) + D(T_1[k_s], T_2[j]) \\ & \quad - D(T_1[k_s], \theta), \\ f\_T(k, j_t) & \text{if } D(T_1[k], T_2[j]) = D(\theta, T_2[j]) + D(T_1[k], T_2[j_t]) \\ & \quad - D(\theta, T_2[j_t]). \end{cases}$$

**Lemma 3** *Let  $t_1[i]$  be the parent of  $t_1[k]$ . Let  $x$  be  $t$  or  $f$  in the following formula for  $x\_T(i, j)$ .*

$$t\_F(i, j) = \begin{cases} (0, 0) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[k] \text{ is deleted,} \\ t\_T(k, j_t) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[k] \text{ is matched} \\ & \quad \text{to } T_2[j_t], \\ (0, 0) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[i_s], F_2[j]) \\ & \quad - D(F_1[i_s], \theta), i_s \neq k, \\ (k+1, j) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[k], F_2[j]) \\ & \quad - D(F_1[k], \theta), \\ t\_F(i, j_t) & \text{if } D(F_1[i], F_2[j]) = D(\theta, F_2[j]) + D(F_1[i], F_2[j_t]) \\ & \quad - D(\theta, F_2[j_t]), \end{cases}$$

$$f\_F(i, j) = \begin{cases} (0, 0) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[k] \text{ is deleted,} \\ f\_T(k, j_t) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[k] \text{ is matched} \\ & \text{to } T_2[j_t], \\ (0, 0) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[i_s], F_2[j]) \\ & \quad - D(F_1[i_s], \theta), i_s \neq k, \\ f\_F(k, j) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[k], F_2[j]) \\ & \quad - D(F_1[k], \theta), \\ f\_F(i, j_t) & \text{if } D(F_1[i], F_2[j]) = D(\theta, F_2[j]) + D(F_1[i], F_2[j_t]) \\ & \quad - D(\theta, F_2[j_t]), \end{cases}$$

$$x\_T(i, j) = \begin{cases} x\_F(i, j) & \text{if } D(T_1[i], T_2[j]) = D(F_1[i], F_2[j]) + \gamma(t_1[i], t_2[j]), \\ (0, 0) & \text{if } D(T_1[i], T_2[j]) = D(T_1[i], \theta) + D(T_1[i_s], T_2[j]) \\ & \quad - D(T_1[i_s], \theta), i_s \neq k, \\ x\_T(k, j) & \text{if } D(T_1[i], T_2[j]) = D(T_1[i], \theta) + D(T_1[k], T_2[j]) \\ & \quad - D(T_1[k], \theta), \\ x\_T(i, j_t) & \text{if } D(T_1[i], T_2[j]) = D(\theta, T_2[j]) + D(T_1[i], T_2[j_t]) \\ & \quad - D(\theta, T_2[j_t]). \end{cases}$$

**Lemma 4** Let  $t_1[i]$  be a proper ancestor of  $t_1[p(k)]$  and  $t_1[i_k]$  be the child of  $t_1[i]$  which is an ancestor of  $t_1[k]$ . Let  $x$  be  $t$  or  $f$  in the following formulae for  $x\_F(i, j)$  and  $x\_T(i, j)$ .

$$x\_F(i, j) = \begin{cases} (0, 0) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[i_k] \text{ is deleted,} \\ x\_T(i_k, j_t) & \text{if } D(F_1[i], F_2[j]) = E(m_i, n_j) \text{ and } T_1[i_k] \text{ is matched} \\ & \text{to } T_2[j_t], \\ (0, 0) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[i_s], F_2[j]) \\ & \quad - D(F_1[i_s], \theta), s \neq k, \\ x\_F(i_k, j) & \text{if } D(F_1[i], F_2[j]) = D(F_1[i], \theta) + D(F_1[i_k], F_2[j]) \\ & \quad - D(F_1[i_k], \theta), \\ x\_F(i, j_t) & \text{if } D(F_1[i], F_2[j]) = D(\theta, F_2[j]) + D(F_1[i], F_2[j_t]) \\ & \quad - D(\theta, F_2[j_t]), \end{cases}$$

$$x\_T(i, j) = \begin{cases} x\_F(i, j) & \text{if } D(T_1[i], T_2[j]) = D(F_1[i], F_2[j]) + \gamma(t_1[i], t_2[j]), \\ (0, 0) & \text{if } D(T_1[i], T_2[j]) = D(T_1[i], \theta) + D(T_1[i_s], T_2[j]) \\ & \quad - D(T_1[i_s], \theta), s \neq k, \\ x\_T(i_k, j) & \text{if } D(T_1[i], T_2[j]) = D(T_1[i], \theta) + D(T_1[i_k], T_2[j]) \\ & \quad - D(T_1[i_k], \theta), \\ x\_T(i, j_t) & \text{if } D(T_1[i], T_2[j]) = D(\theta, T_2[j]) + D(T_1[i], T_2[j_t]) \\ & \quad - D(\theta, T_2[j_t]). \end{cases}$$

From these formulae, it is clear that we can compute  $t\_T(|T_1|, |T_2|)$  and  $f\_T(|T_1|, |T_2|)$  with a small modification of Algorithm 2.

We now consider how to use  $t\_T(|T_1|, |T_2|) = (t_1, t_2)$  and  $f\_T(|T_1|, |T_2|) = (f_1, f_2)$  to decompose the optimal mapping into smaller component mappings.

**Case 1:**  $t_1 = t_2 = 0$ . In this case, since  $T_1[k]$  is deleted, it is not in the optimal mapping for  $D(T_1, T_2)$ . Therefore the optimal mapping for  $D(T_1, T_2)$  would be the op-

timal mapping for  $D(T_1/T_1[k], T_2)$  where  $T_1/T_1[k]$  is  $T_1$  with  $T_1[k]$  deleted. Note that  $|T_1/T_1[k]| < 0.5|T_1|$ .

**Case 2:**  $t_1 = k_s$  and  $t_2 = l$ . In this case,  $T_1[k]$  is matched with  $T_2[l]$  and in this matching,  $T_1[k_s]$  is matched with  $T_2[l]$  and the rest of  $T_1[k]$  is deleted. Therefore the optimal mapping for  $D(T_1, T_2)$  has two component mappings. One component is the optimal mapping for  $D(T_1/F_1[k], T_2/F_2[l])$  with additional condition that  $t_1[k]$  has to match with  $t_2[l]$  though  $(k, l)$  is not in the optimal mapping of  $D(T_1, T_2)$  since  $t_1[k]$  is deleted. The other component is the optimal mapping of  $D(T_1[k_s], T_2[l])$ . Note that  $|T_1/F_1[k]| \leq 0.5|T_1|$  and  $|T_1[k_s]| \leq 0.5|T_1|$ .

Computing  $D(T_1/F_1[k], T_2/F_2[l])$  with condition that some of the leaves of the two trees are forced to match would not be more difficult than computing  $D(T_1/F_1[k], T_2/F_2[l])$  without any condition. In fact the condition will force the ancestors of these matched leaves of the two trees to match making the computation less expensive.

**Case 3:**  $t_1 = k$  and  $t_2 = l$ ,  $f_1 = f_2 = 0$ . In case 3, case 4 and case 5,  $t_1 = k$  and  $t_2 = l$ . This means that in the optimal mapping for  $D(T_1, T_2)$ ,  $t_1[k]$  is matched with  $t_2[l]$ . Therefore one component mapping (for case 3–5) is the optimal mapping for  $D(T_1/F_1[k], T_2/F_2[l])$  with additional condition that  $t_1[k]$  matched with  $t_2[l]$ . Note that  $|T_1/T_1[k]| < 0.5|T_1|$ .

In this case, since  $f_1 = f_2 = 0$  there is no additional component.

**Case 4:**  $t_1 = k$  and  $t_2 = l$ ,  $f_1 = k_s$  and  $f_2 = l$ . In this case, one component mapping is the same as that in case 3 and the other component mapping is the optimal mapping for  $D(F_1[k_s], F_2[l])$ . Note that  $|F_1[k_s]| \leq 0.5|T_1|$ .

**Case 5:**  $t_1 = k$  and  $t_2 = l$ ,  $f_1 = k$  and  $f_2 = l$ . In this case, one component mapping is the same as that in case 3 and the other component mapping is the optimal mapping for  $D(F_1[k], F_2[l])$ . Since in this case,  $D(F_1[k], F_2[l]) = E(m_k, n_l)$ , the optimal mapping for  $D(F_1[k], F_2[l])$  can be further decomposed into the optimal mappings of the matching subtree pairs of  $E(m_k, n_l)$ . Note that  $|T_1[k_s]| \leq 0.5|T_1|$  for  $1 \leq s \leq m_k$ .

**Case 6–8:**  $t_1 = k + 1$  and  $t_2 = l$ . This means that in the optimal mapping for  $D(T_1, T_2)$ ,  $t_1[k]$  is deleted and  $F_1[k]$  is matched with  $F_2[l]$ . Therefore one component mapping (for case 6–8) is the optimal mapping for  $D(T_1/F_1[k], T_2/F_2[l])$  with additional condition that  $t_2[l]$  matched with a proper ancestor of  $t_1[k]$ . Note that  $|T_1 - T_1[k]| < 0.5|T_1|$ .

There is no additional component for case 6. The other component for case 7 is exactly the same as that of case 4. The other components for case 8 are exactly the same as that of case 5.

In all the above cases, with  $t_1$ ,  $t_2$ ,  $f_1$ , and  $f_2$ , we can determine components needed except matching subtree pairs of  $E(m_k, n_l)$  for case 5 and case 8. Finding the matching subtree pairs of  $E(m_k, n_l)$  can be done in  $O(|F_1[k]| |F_2[l]|)$  time and  $O(\log(|F_1[k]|) |F_2[l]|)$  space. First we find  $T_1[k_s]$  such that  $\sum_{i=1}^{s-1} |T_1[k_i]| \leq 0.5|F_1[k]|$  and  $\sum_{i=1}^s |T_1[k_i]| > 0.5|F_1[k]|$ . We then determine  $T_2[l_t]$ , in  $O(|F_1[k]| |F_2[l]|)$  time and  $O(\log(|F_1[k]|) |F_2[l]|)$  space, such that either  $T_1[k_s]$  matched with  $T_2[l_t]$  in  $E(m_k, n_l)$  or  $T_1[k_s]$  is deleted,  $T_1[k_1], \dots, T_1[k_{s-1}]$  are matched with  $T_2[l_1], \dots, T_2[l_t]$ , and  $T_1[k_{s+1}], \dots, T_1[k_{m_k}]$  are matched with  $T_2[l_{t+1}], \dots, T_2[l_{n_l}]$  in  $E(m_k, n_l)$ . In both cases, since  $\sum_{i=1}^{s-1} |T_1[k_i]| \leq 0.5|F_1[k]|$  and

$\sum_{i=s+1}^{m_k} |T_1[k_i]| \leq 0.5|F_1[k]|$ , we can repeat and get all the subtree matching pairs for  $E(m_k, n_l)$  in  $O(|F_1[k]||F_2[l]|)$  time and  $O(\log(|F_1[k]|)|F_2[l]|)$  space.

Therefore using  $c|T_1||T_2|$  time for some constant  $c$  and  $O(\log(|T_1|)|T_2|)$  space we can decompose the optimal mapping for  $D(T_1, T_2)$  in to some components such that for each component the involved subtree or subforest of  $T_1$  has a size less or equal to  $0.5|T_1|$ .

In the next step we will use the same algorithm for all the components and the total time is bounded by  $c0.5|T_1||T_2|$ .

The above analysis means that for our algorithm the space complexity is bounded by  $O(\log(|T_1|)|T_2|)$  and the time complexity is bounded by  $c|T_1||T_2| + c\frac{1}{2}|T_1||T_2| + c\frac{1}{4}|T_1||T_2| + \dots \leq 2c|T_1||T_2| = O(|T_1||T_2|)$ .

**Theorem 1** *Given two ordered trees  $T_1$  and  $T_2$ , the constrained edit distance and the optimal constrained mapping between them can be computed in  $O(|T_1||T_2|)$  time and  $O(\log(|T_1|)|T_2|)$  space.*

### 3 Alignment of Trees

The alignment of tree is another measure for comparison of ordered trees [5]. The definition of *alignment of trees* is as follows: Inserting a node  $u$  into  $T$  means that for some node  $v$  (could be a leaf) in  $T$ , we make  $u$  the parent of the consecutive subsequence of the children of  $v$  (if any) and then  $v$  the parent of  $u$ . We also allow to directly add/insert a node as a child of a leaf in the tree. Given two trees  $T_1$  and  $T_2$ , an alignment of the two trees can be obtained by first inserting nodes labeled with spaces into  $T_1$  and  $T_2$  such that the two resulting tree  $T'_1$  and  $T'_2$  have the same structure, (i.e., they are identical if labels are ignored) and then overlaying  $T'_1$  and  $T'_2$ . A score is defined for each pair of labels. The value of an alignment is the total score of all the opposing labels in the alignment. The problem here is to find an alignment with the optimal (maximum or minimum) value.

**Theorem 2** *The algorithm in [11] for alignment of trees runs in  $O(|T_1||T_2|(\deg(T_1) + \deg(T_2))^2)$  time and requires  $O(\log(|T_1|)|T_2|(\deg(T_1) + \deg(T_2))\deg(T_1))$  space.*

The basic ideas of the algorithm are similar to that for constrained edit distance. We first use the method in [11] to compute the cost of an optimal alignment. The remaining task is to get the alignment.

Let  $q$  be the node in  $T_1$  such that  $|T_1[q]| \geq 0.5|T_1|$  and for any child  $q_i$  of node  $q$ ,  $|T_1[q_i]| < 0.5|T_1|$ .  $q$  is referred to as a *cutting point* in  $T_1$ . By definition, there always exists a cutting point for any  $T_1$ .

By the definition of the alignment, node  $q$  is either aligned with a node  $T_2[j]$  or aligned with an inserted node (labeled with empty). We can modify the algorithm for computing the cost of an optimal alignment so that when the cost of an optimal alignment is computed, we immediately know the configuration of node  $q$  in the alignment. Thus, we can decompose the alignment of the whole two tree  $T_1$  and  $T_2$  into two parts (1) the alignment of subtree  $T_1[q]$  with a subtree or subforest of  $T_2$ , and (2) the alignment of the remaining parts of  $T_1$  and  $T_2$ .

Repeating the process, we can get the alignment. Similar to Theorem 1, we can show that the time required is still  $O(|T_1||T_2|(\deg(T_1) + \deg(T_2))^2)$ . In fact, the running time is twice of that for computing the cost of an optimal alignment.

## 4 Conclusion

We have presented space efficient algorithms for the computation of constrained edit distance and tree alignment for ordered rooted trees. The techniques can also be applied to other tree comparison problems such as degree-one and degree-two edit distance between ordered trees.

**Acknowledgements** Kaizhong Zhang is partially supported by the Natural Sciences and Engineering Research Council of Canada under Grant No. OGP0046373. Lusheng Wang is fully supported by a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China [Project No. CityU 120905].

## References

1. Bille, P.: A survey on tree edit distance and related problems. *Theor. Comput. Sci.* **337**, 217–239 (2005)
2. Chodorow, M., Klavans, J.L.: Locating syntactic patterns in text corpora. Manuscript, Lexical systems, IBM Research, T.J. Watson Research Center, Yorktown Heights, New York (1990)
3. Dulucq, S., Touzet, H.: Decomposition algorithm for tree editing distance. *J. Discrete Algorithms* (2004)
4. Hirschberg, D.S.: A linear space algorithm for computing maximal common subsequences. *Commun. ACM* **18**, 341–343 (1975)
5. Jiang, T., Wang, L., Zhang, K.: Alignment of trees—an alternative to tree edit. *Theor. Comput. Sci.* **143**(1), 137–148 (1995)
6. Klein, P.: Computing the edit-distance between unrooted ordered trees. In: *Proceedings of 6th European Symposium on Algorithms*, pp. 91–102 (1998)
7. Richter, T.: A new measure of the distance between ordered trees and its applications. Technical Report 85166-cs, Department of Computer Science, University of Bonn (1997)
8. Shapiro, B., Zhang, K.: Comparing multiple RNA secondary structures using tree comparisons. *Comput. Appl. Biosci.* **6**(4), 309–318 (1990)
9. Selkow, S.M.: The tree-to-tree editing problem. *Inf. Process. Lett.* **6**, 184–186 (1977)
10. Tai, K.C.: The tree-to-tree correction problem. *J. ACM* **26**, 422–433 (1979)
11. Wang, L., Zhao, J.: Parametric alignment of ordered trees. *Bioinformatics* **19**, 2237–2245 (2003)
12. Zhang, K.: Algorithms for the constrained editing distance between ordered labeled trees and related problems. *Pattern Recognit.* **28**(3), 463–474 (1995)
13. Zhang, K.: Efficient parallel algorithms for tree editing problems. In: *Proceedings of the Seventh Symposium on Combinatorial Pattern Matching*, Laguna Beach, California, June 1996. *Lecture Notes in Computer Science*, vol. 1075, pp. 361–372. Springer, Berlin (1996)
14. Zhang, K.: Computing similarity between RNA secondary structures. In: *Proceedings of IEEE International Joint Symposia on Intelligence and Systems*, Rockville, Maryland, pp. 126–132 (May 1998)
15. Zhang, K., Shasha, D.: Simple fast algorithms for the editing distance between trees and related problems. *SIAM J. Comput.* **18**(6), 1245–1262 (1989)
16. Zhang, K., Wang, J.T.L., Shasha, D.: On the editing distance between undirected acyclic graphs. *Int. J. Found. Comput. Sci.* **7**(1), 43–57 (1996)