

# AN $O(N \log N)$ ALGORITHM FOR THE MAXIMUM AGREEMENT SUBTREE PROBLEM FOR BINARY TREES\*

RICHARD COLE<sup>†</sup>, MARTIN FARACH-COLTON<sup>‡</sup>, RAMESH HARIHARAN<sup>§</sup>, TERESA  
PRZYTYCKA<sup>¶</sup>, AND MIKKEL THORUP <sup>||</sup>

**Abstract.** The Maximum Agreement Subtree problem is the following: given two rooted trees whose leaves are drawn from the same set of items (e.g., species), find the largest subset of these items so that the portions of the two trees restricted to these items are isomorphic. We consider the case which occurs frequently in practice, i.e., the case when the trees are binary, and give an  $O(n \log n)$  time algorithm for this problem.

**1. Introduction.** Suppose we are given two rooted trees  $T_1$  and  $T_2$  with  $n$  leaves each. The internal nodes of each tree have at least two children each. The leaves in each tree are labeled with the same set of labels and further, no label occurs more than once in a particular tree. An *agreement subtree* of  $T_1$  and  $T_2$  is defined as follows. Let  $L_1$  be a subset of the leaves of  $T_1$  and let  $L_2$  be the subset of those leaves of  $T_2$  which have the same labels as leaves in  $L_1$ . The subtree of  $T_1$  induced by  $L_1$  is an agreement subtree of  $T_1$  and  $T_2$  if and only if it is *isomorphic* to the subtree of  $T_2$  induced by  $L_2$ . The Maximum Agreement Subtree problem (henceforth called *MAST*) asks for the largest agreement subtree of  $T_1$  and  $T_2$ .

We need to define the terms *induced subtree* and *isomorphism* used above. Intuitively, the subtree of  $T$  induced by a subset  $L$  of the leaves of  $T$  is the topological subtree of  $T$  restricted to the leaves in  $L$ , with branching information relevant to  $L$  preserved. More formally, for any two leaves  $a, b$  of a tree  $T$ , let  $\text{lca}_T(a, b)$  denote their lowest common ancestor in  $T$ . If  $a = b$ ,  $\text{lca}_T(a, b) = a$ . The *subtree*  $U$  of  $T$  induced by a subset  $L$  of the leaves is the tree with leaf set  $L$  and interior node set  $\{\text{lca}_T(a, b) \mid a, b \in L\}$  inheriting the ancestor relation from  $T$ , that is, for all  $a, b \in L$ ,  $\text{lca}_U(a, b) = \text{lca}_T(a, b)$ .

Intuitively, two trees are isomorphic if the children of each node in one of the trees can be reordered so that the leaf labels in each tree occur in the same order and the shapes of the two trees become identical. Formally, we say two trees  $U_1$  and  $U_2$  with the same leaf labels are *isomorphic* if there is a 1-1 mapping  $\mu$  between their nodes, mapping leaves to leaves with the same labels and such that for any two different leaves  $a, b$  of  $U_1$ ,  $\mu(\text{lca}_{U_1}(a, b)) = \text{lca}_{U_2}(\mu(a), \mu(b))$ .

**Motivation.** The *MAST* problem arises naturally in biology and linguistics as a measure of consistency between two evolutionary trees over species and languages, respectively. An evolutionary tree for a set of *taxa*, either species or languages, is a rooted tree whose leaves represent the taxa and whose internal nodes represent

---

\*Work supported by NSF grants CCR-9202900 and CCR-9503309, by Sloan and Department of Energy Postdoctoral Fellowship for Computational Biology, and by grant GM29458.

<sup>†</sup>Courant Institute of Mathematical Sciences, NYU.

<sup>‡</sup>DIMACS.

<sup>§</sup>Indian Institute of Science, Bangalore. Work done partly when the author was at Max Planck Institut für Informatik, Germany, and when visiting NYU.

<sup>¶</sup>Johns Hopkins School of Medicine.

<sup>||</sup>AT&T Labs-Research. Some of this work was done while at University of Copenhagen, Denmark, and while visiting MIT.

ancestor information. It is often difficult to determine the true phylogeny for a set of taxa, and one way to gain confidence in a particular tree is to have different lines of evidence supporting that tree. In the biological taxa case, one may construct trees from different parts of the DNA of the species. These are known as *gene trees*. For many reasons, these trees need not entirely agree, and so one is left with the task of finding a consensus of the various gene trees. The maximum agreement subtree is one method of arriving at such a consensus. Notice that a gene is usually a binary tree, since DNA replicates by a binary branching process. Therefore, the case of binary trees is of great interest.

Another application arises in automated translation between two languages [GY95]. The two trees are the parse trees for the same meaning sentences in the two languages. A complication that arises in this application (due in part to imperfect dictionaries) is that words need not be uniquely matched, i.e., a word at the leaf of one tree could match a number (usually small) of words at the leaves of the other tree. The aim is to find a maximum agreement subtree; this is done with the goal of improving context-using dictionaries for automated translation. So long as each word in one tree has only a constant number of matches in the other tree (possibly with differing weights), the algorithm given here can be used and its performance remains  $O(n \log n)$ . More generally, if there are  $m$  word matches in all, the performance becomes  $O((m + n) \log n)$ . Note however, that if there are two collections of equal meaning words in the two trees of sizes  $k_1$  and  $k_2$  respectively, they induce  $k_1 k_2$  matches.

**Previous Work.** Finden and Gordon [FG85] gave a heuristic algorithm for the *MAST* problem on binary trees which had an  $O(n^5)$  running time and did not guarantee an optimal solution. Kubicka, Kubicki and McMorris [KKM95] gave an  $O(n^{(.5+\epsilon) \log n})$  algorithm for the same problem. The first polynomial time algorithm for this problem was given by Steel and Warnow [SW93]; it had a running time of  $O(n^2)$ . Steel and Warnow also considered the case of non-binary and unrooted trees. Their algorithm takes  $O(n^2)$  time for fixed degree rooted and unrooted trees and  $O(n^{4.5} \log n)$  for arbitrary degree rooted and unrooted trees. They also give a linear reduction from the rooted to the unrooted case. Farach and Thorup gave an  $O(nc\sqrt{\log n})$  time algorithm for the *MAST* problem on binary trees; here  $c$  is a constant greater than 1. For arbitrary degree trees, their algorithm takes  $O(n^2 c\sqrt{\log n})$  time for the unrooted case [FT95] and  $O(n^{1.5} \log n)$  time for the rooted case [FT97]. Farach, Przytycka, and Thorup [FPT95a] obtained an  $O(n \log^3 n)$  algorithm for the *MAST* problem on binary trees. Kao [Ka95] obtained an algorithm for the same problem which takes  $O(n \log^2 n)$  time. This algorithm takes  $O(\min\{nd^2 \log d \log^2 n, nd^{\frac{3}{2}} \log^3 n\})$  for degree  $d$  trees. Finally, Cole and Hariharan [CR96] improved the algorithm from [FPT95a] to an  $O(n \log n)$  algorithm.

The *MAST* problem for more than two trees has also been studied. Amir and Keselman [AK97] showed that the problem is *NP*-hard for even 3 unbounded degree trees. However, polynomial bounds are known [AK97, FPT95b] for three or more bounded degree trees.

**Our Contribution.** This paper is the combined journal version of [FPT95a] and [CR96] and presents an  $O(n \log n)$  algorithm for the *MAST* problem for two binary trees.

The  $O(n \log^3 n)$  algorithm of [FPT95a] can be viewed as taking the following approach (although the authors do not describe it this way). It identifies two special

cases and then solves the general case by interpolating between these cases.

**Special Case 1:** The internal nodes in both trees form a path. The *MAST* problem reduces to essentially a size  $n$  Longest Increasing Subsequence Problem in this case. As is well known, this can be solved in  $O(n \log n)$  time.

**Special Case 2:** Both trees  $T_1$  and  $T_2$  are complete binary trees. For each node  $v$  in  $T_2$ , only certain nodes  $u$  in  $T_1$  can be usefully mapped to  $v$ , in the sense that the subtree of  $T_1$  rooted at  $u$  and the subtree of  $T_2$  rooted at  $v$  have a non-empty Agreement Subtree. There are  $O(n \log^2 n)$  such pairs  $(u, v)$ . This can be seen as follows. Note that for  $(u, v)$  to be such a pair, the subtree of  $T_1$  rooted at  $u$  and the subtree of  $T_2$  rooted at  $v$  must have a leaf-label in common. For each label, there are only  $O(\log^2 n)$  such pairs obtained by pairing each ancestor of the leaf with this label in  $T_1$  with each ancestor of the leaf with this label in  $T_2$ . The total number of interesting pairs is thus  $O(n \log^2 n)$ .

For each pair, computing the *MAST* takes  $O(1)$  time, as it is simply a question of deciding the best way of pairing their children.

The interpolation process takes a centroid decomposition of the two trees and compares pairs of centroid paths, rather than individual nodes as in the complete tree case. The comparison of a pair of centroid paths requires finding matchings with special properties in appropriately defined bipartite graphs, a non-trivial generalization of the Longest Increasing Subsequence problem. This process creates  $O(n \log^2 n)$  interesting  $(u, v)$  pairs, each of which takes  $O(\log n)$  time to process.

In [CR96] two improvements are given, each of which gains a  $\log n$  factor.

**Improvement 1.** The complete tree special case is improved to  $O(n \log n)$  time as follows. A pair of nodes  $(u, v)$ ,  $u \in T_1$ ,  $v \in T_2$ , is said to be *interesting* if there is an agreement subtree mapping  $u$  to  $v$ . As is shown below, for complete trees, the total number of interesting pairs  $(u, v)$  is just  $O(n \log n)$ . Consider a node  $v$  in  $T_2$ . Let  $L_2$  be the set of leaves which are descendants of  $v$ . Let  $L_1$  be the set of leaves in  $T_1$  which have the same labels as the leaves in  $L_2$ . The only nodes that may be mapped to  $v$  are the nodes  $u$  in the subtree of  $T_1$  induced by  $L_1$ . The number of such nodes  $u$  is  $O(\text{size of the subtree of } T_2 \text{ rooted at } v)$ . The total number of interesting pairs is thus the sum of the sizes of all subtrees of  $T_2$ , which is  $O(n \log n)$ .

This reduces the number of interesting pairs  $(u, v)$  to  $O(n \log n)$ . Again, processing a pair takes  $O(1)$  time (this is less obvious, for identifying the descendants of  $u$  which root the subtrees with which the two subtrees of  $v$  can be matched is non-trivial). Constructing the above induced subtree itself can be done in  $O(|L_1|)$  time, as will be detailed later. The basic tool here is to preprocess trees  $T_1$  and  $T_2$  in  $O(n)$  time so that least common ancestor queries can be answered in  $O(1)$  time.

**Improvement 2:** As in [FPT95a], when the trees are not complete binary trees, we take centroid paths and match pairs of centroid paths. The  $O(\log n)$  cost that the algorithm in [FPT95a] incurs in processing each such interesting pair of paths arises when there are large (polynomial in  $n$  size) instances of the generalized Longest Increasing Subsequence Problem. At first sight, it is not clear that large instances of these problems can be created for sufficiently many of the interesting pairs; unfortunately, this turns out to be the case. However, these problem instances still have some useful structure. By using (static) weighted trees we process pairs of interesting vertices in  $O(1)$  time per pair, on the average, as is shown by an appropriately parameterized analysis.

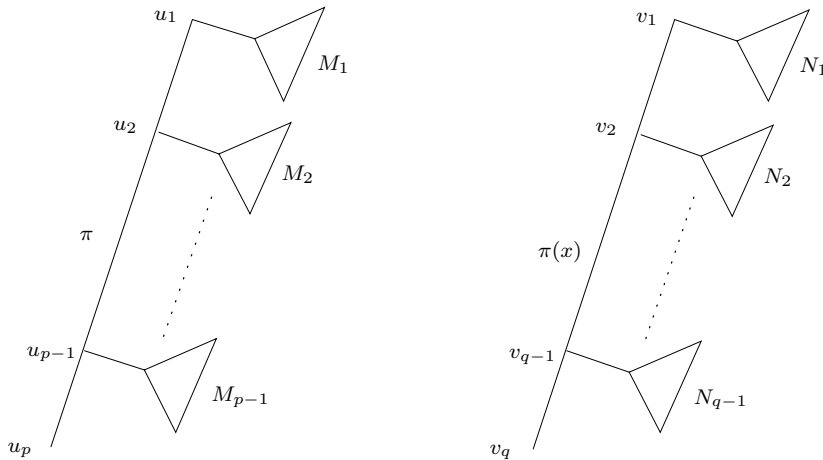


FIG. 3.1. Centroid paths in  $T_1$  and  $T_2$  with  $\pi$  starting in root of  $T_1$  and  $\pi(x)$  starting in  $x = v_1$  in  $T_2$ .

The paper is organized as follows. Section 2 gives some basic definitions and primitives. Section 3 outlines the algorithm. Section 4 provides further details of the algorithm and Section 5 gives the analysis. The remaining sections deal with problems raised in Sections 3-5.

**2. Definitions and Preliminaries.** All trees henceforth refer to binary trees whose internal nodes have exactly two children.

The tree  $T(x)$  denotes the subtree of  $T$  rooted at vertex  $x$ . The *size* of a tree  $T$ , denoted by  $|T|$ , is the number of leaves in it. In our problem,  $|T_1| = |T_2| = n$ .

Given a binary tree  $T$ , its *centroid decomposition* is a partitioning of its vertices into disjoint paths obtained as follows. First, for each internal node  $x$  in  $T$ , the edge to the child with the maximal number of leaves below it is called a *centroid edge*. Here ties are broken arbitrarily. Now, the centroid edges form a collection of disjoint paths, called *centroid paths*. The *beginning* of such a centroid path  $P$  is defined to be the vertex  $x$  closest to the root of  $T$ . Then, if we remove  $P$  from  $T(x)$ , we get a forest of trees called the *side trees* of  $T(x)$ , and then each side tree is of size at most  $|T(x)|/2$ . Note that the full centroid decomposition could also be found recursively by first finding the centroid path from the root, and then recursing on each side tree. Also note that the centroid decomposition of  $T$  can easily be found in  $O(n)$  time.

A tree  $T$  can be preprocessed in  $O(|T|)$  time so that given any subset  $L$  of its leaves in left to right order, the subtree induced by  $L$  can be computed in  $O(|L|)$  time. The details of this procedure are described in Section 8.

The set of labels at the leaves of  $T_1$  is identical to that at the leaves of  $T_2$ . For a leaf  $l$  in one of these trees, the leaf with the same label in the other tree is called its *twin*. Two subtrees, one from each tree, are said to *intersect* if and only if some leaf in one subtree has a twin in the other. The subtree of  $T_2$  induced by some subset of the leaves of  $T_1$  is the subtree of  $T_2$  induced by the twins of these leaves of  $T_1$ .

**3. Algorithm Outline.** We need some definitions to outline the algorithm.

**Definitions.** Let  $\pi$  be the centroid path containing the root of  $T_1$ . Let  $p = |\pi|$  and  $u_1, u_2, \dots, u_{p-1}, u_p$  be the vertices on this path in order from the root. Let  $M_1, M_2, \dots, M_{p-1}$  comprise the forest of side trees created by the removal of  $\pi$  from  $T_1$ . Let  $m_i = |M_i|$  be the number of leaves in  $M_i$  (see Fig.1). Recall that  $m_i \leq n/2$  as  $\pi$  is a centroid path. For technical reasons, we define  $M_p$  to be the tree consisting of the vertex  $u_p$  and set  $m_p = 1$ . Then,  $\sum_{i=1}^p m_i = n$ .

Given trees  $T_1, T_2$ , our aim is to determine the maximum agreement subtree of  $T_1$  and  $T_2$  efficiently. To do so, we will need to compute not just this agreement subtree but the maximum agreement subtree of  $T_1$  and  $T_2(w)$ , for each  $w$  in  $T_2$ . All these subtrees will be computed implicitly by a procedure  $Agree(T_1, T_2)$ .

$Agree(T_1, T_2)$  proceeds broadly as follows. See Fig.1. First, it recursively computes the maximum agreement subtree of  $M_i, T_2(w)$ , for each side tree  $M_i$  of  $T_1$  and each vertex  $w$  in  $T_2$ . This is not done explicitly though, as we will see shortly. Second, it uses the information gathered in this process to compute the maximum agreement subtree of  $T_1, T_2(w)$ , for each vertex  $w$  in  $T_2$ . Both steps together take  $O(n \log n)$  time.

We will describe  $Agree(T_1, T_2)$  in detail shortly. Two observations will be used in this description.

First, as in the discussion of interesting pairs in Improvement 1 from the introduction, note that it is not necessary to compute the maximum agreement subtrees of  $M_i, T_2(w)$  for all  $w \in T_2$ . The first reason is that the maximum agreement subtree of  $M_i, T_2(w)$  is empty if  $M_i$  does not intersect with  $T_2(w)$ . The second reason is that the maximum agreement subtrees of  $M_i, T_2(w)$  and  $M_i, T_2(\text{parent}(w))$  are identical if  $M_i$  does not intersect with the subtree of  $T_2$  rooted at  $w$ 's sibling. For these two reasons, it suffices to compute  $Agree(M_i, S_i)$  where  $S_i$  is the subtree of  $T_2$  induced by the leaves of  $M_i$ . This implicitly computes the maximum agreement subtrees of  $M_i, T_2(w)$  for all  $w \in T_2$ .

Second, note that if the maximum agreement subtree of  $T_1, T_2(w)$  does not have any vertex on the centroid path  $\pi$  of  $T_1$ , then all vertices in this agreement subtree belong to a single side tree  $M_i$  of  $T_1$ . For, if these vertices were distributed over two or more side trees of  $T_1$ , then these vertices would have a common ancestor on  $\pi$  which would be in the maximum agreement subtree as well. Thus, when the maximum agreement subtree of  $T_1, T_2(w)$  does not have any vertex on  $\pi$ , it can be determined using some  $Agree(M_i, S_i)$  from the previous paragraph,

**Algorithm Outline.**  $Agree(T_1, T_2)$  has three steps.

**Step 1.** The centroid decompositions of  $T_1$  and  $T_2$  are computed. This takes  $O(n)$  time.

**Step 2.** For each  $i$ ,  $1 \leq i \leq p-1$ ,  $Agree(M_i, S_i)$  is computed recursively, where  $S_i$  is the subtree of  $T_2$  induced by the leaves of  $M_i$ . We will inductively assume that this takes  $O(m_i \log m_i)$  for each  $i$ . Recall that if the maximum agreement subtree of  $T_1$  and  $T_2$  contains no vertex from  $\pi$  then it will be found in Step 2. Step 3 handles the other case.

**Step 3: Matching  $\pi$ .** For each  $w \in T_2$ , the largest agreement subtree for the trees  $T_1$  and  $T_2(w)$  is found using information from Step 2. Informally, we call this the process of ‘‘matching’’  $\pi$  at each of the vertices  $w$  of  $T_2$ . We will show how this is done in  $O(\sum_{i=1}^p m_i \log \frac{n}{m_i})$  time.

Clearly, the total time over all three steps is  $O(n \log n)$ . The following sections show how Step 3 is performed in  $O(\sum_{i=1}^p m_i \log \frac{n}{m_i})$  time. Starting in the next sec-

tion, we will define some bipartite graphs. Each graph will correspond to a particular centroid path in the centroid decomposition of  $T_2$  and will be used to match  $\pi$  as in Step 3 at all the vertices in this centroid path. These graphs will have the property that a particular kind of matching, which can be computed efficiently, will correspond to the relevant agreement subtrees.

**4. The Matching Graphs  $G(x)$  and the  $\pi$  Matching Algorithm. Definitions.** Recall that the centroid decomposition of  $T_2$  partitions its vertices into disjoint paths and that the beginning of such a path is defined to be the vertex closest to the root of  $T_2$  in that path. Let  $X$  denote the set of vertices in  $T_2$  at which paths in the above decomposition begin.

We define a number of bipartite graphs, one for each  $x \in X$ . The graph  $G(x)$  corresponding to vertex  $x$  is defined as follows.

**Vertices of  $G(x)$ .** The left vertex set  $L(x)$  of  $G(x)$  is a subset of  $\{u_1, \dots, u_p\}$ . Vertex  $u_i$ ,  $1 \leq i \leq p-1$ , is in the set if and only if  $M_i$  and  $T_2(x)$  intersect. Vertex  $u_p$  is in the set if and only if its twin is in  $T_2(x)$ .

The right vertex set  $R(x)$  of  $G(x)$  is exactly the set of vertices in the centroid path beginning at vertex  $x$ .

Since both sets of vertices are drawn from centroid paths, we order the vertices on each side in the order they occur on their respective centroid paths. The *topmost* vertex is the closest to the root and the *bottommost* is the farthest. Further, two edges  $(a, b)$  and  $(a', b')$  in  $G(x)$  are said to *cross* if  $a$  is above  $a'$  and  $b$  is below  $b'$ , or vice versa. In addition, edge  $(a, b)$  is said to *dominate*  $(a', b')$  in  $G(x)$  if  $a$  is above  $a'$  and  $b$  is above  $b'$ . The *topmost* edge in a set of edges, if any, is the edge which dominates all other edges in that set.

Before defining the edges of  $G(x)$ , we need the following definitions.

**Definitions.** Let  $\pi(x)$  be the centroid path containing  $x$ . Let  $q$  be the length of this path. Let  $v_1, v_2, \dots, v_q$  be the vertices on this path in order from the root. Let  $N_1, N_2, \dots, N_{q-1}$  comprise the forest of side trees of  $T_2(x)$  created by the removal of  $v_1, \dots, v_q$  from  $T_2(x)$ . Let  $n_i = |N_i|$ , for  $i = 1 \dots q-1$  (see Fig.1). For technical reasons, we define  $N_q$  to be the tree consisting of the vertex  $v_q$  and set  $n_q = 1$ . Then  $\sum_{i=1}^q n_i = |T_2(x)|$ .

**4.1. Motivation for Defining Edges of  $G(x)$ .** We first motivate the definition of edges of the graph  $G(x)$ . The purpose of  $G(x)$  is to determine maximum agreement subtrees of  $T_1$  and  $T_2(v_k)$ , for each  $k$ ,  $1 \leq k \leq q$ . The edges of  $G(x)$  will be defined so that maximum weight matchings of a certain kind (called *Agreement Matchings*) in  $G(x)$  will correspond to maximum agreement subtrees of  $T_1$  and  $T_2(v_k)$ ,  $1 \leq k \leq q$ . Clearly, the edges of  $G(x)$  and the agreement matchings must capture the structural properties of these maximum agreement subtrees. We outline these structural properties next and show how they lead to the edge definitions.

Consider the maximum agreement subtree  $\mathcal{A}$  of  $T_1, T_2(v_k)$ . Note that  $\mathcal{A}$  has the following properties.

**Case 1:** If  $\mathcal{A}$  has no vertices in  $\pi(x)$  then it must be the maximum agreement subtree of  $(T_1, N_j)$ , for some  $j$ ,  $k \leq j \leq q-1$ .

**Case 2:** Similarly, if  $\mathcal{A}$  has no vertices in  $\pi$  then it must be the maximum agreement subtree of  $(M_i, T_2(v_k))$ , for some  $i$ ,  $1 \leq i \leq p-1$ .

**Case 3:** Next, suppose  $\mathcal{A}$  has at least one vertex both from  $\pi$  and from  $\pi(x)$ . In other words, there is at least one vertex in  $\mathcal{A}$  which is in  $\pi$  and which maps to a vertex in  $\pi(x)$ .

Suppose vertex  $u_i \in \pi$  is one such vertex and it maps to vertex  $v_j \in \pi(x)$ . If  $u_i$  is not the bottommost such vertex in  $\pi$  and  $z$  is the unique child of  $u_i$  in  $\mathcal{A}$  which is not in  $\pi$ , then  $\mathcal{A}(z)$  must be the maximum agreement subtree of  $(M_i, N_j)$ .

Now if  $u_i$  is indeed the bottommost such vertex, we divide into subcases:

**Case 3.0:** If  $u_i$  is the leaf  $u_p$ , it must map to the other leaf  $v_q$  since leaves can only map to leaves. Then the subtree of  $\mathcal{A}$  rooted at  $u_i$  is just the vertex  $u_i$ , which is also the maximum agreement subtree of  $(T_1(u_i), N_q)$ .

For the remaining subcases we assume that  $u_i$  is not the leaf  $u_p$ .

**Case 3.1:** The subtrees of  $\mathcal{A}$  rooted at the two children of  $u_i$  in  $\mathcal{A}$  are the maximum agreement subtrees of the pairs  $(T_1(u_{i+1}), N_j)$  and  $(M_i, T_2(v_{j+1}))$ . The following facts will be of use in this case. If  $M_i$  and  $N_{j+1}$  do not intersect, then the maximum agreement subtree of  $(M_i, T_2(v_{j+1}))$  is identical to that of  $(M_i, T_2(v_{j'}))$ , where  $j' > j$  is the topmost vertex below  $v_j$  in  $\pi(x)$  such that  $M_i$  and  $N_{j'}$  intersect. And if  $M_{i+1}$  and  $N_j$  do not intersect, then the maximum agreement subtree of  $(T_1(u_{i+1}), N_j)$  is identical to that of  $(T_1(u_{i'}), N_j)$ , where  $i' > i$  is the topmost vertex below  $u_i$  in  $\pi$  such that  $M_{i'}$  and  $N_j$  intersect.

Note that if this subcase does not hold, then the subtree of  $\mathcal{A}$  rooted at one of the two children of  $u_i$  is just the maximum agreement subtree of  $(M_i, N_j)$ . In addition, all descendants of the other child of  $u_i$  in  $\mathcal{A}$  must come from either a single side tree below  $v_j$  in  $T_2$  or a single side tree below  $u_i$  in  $T_1$ . These situations are handled by the second and third cases, respectively.

**Case 3.2:** The subtrees of  $\mathcal{A}$  rooted at the two children of  $u_i$  in  $\mathcal{A}$  are the maximum agreement subtrees of  $(M_i, N_j)$  and  $(T_1(u_{i+1}), N_{j'})$ , for some  $j', j < j' \leq q$ .

**Case 3.3:** The subtrees of  $\mathcal{A}$  rooted at the two children of  $u_i$  in  $\mathcal{A}$  are the maximum agreement subtrees of  $(M_i, N_j)$  and  $(M_{i'}, T_2(v_{j+1}))$ , respectively, for some  $i', i < i' \leq p$ .

The following properties of  $\mathcal{A}$  can be inferred from the above case analysis.

**Property 1:**  $\mathcal{A}$  has a path (which is possibly empty) comprising vertices in  $\pi$  which map to vertices in  $\pi(x)$ , with the property that off-path subtrees are maximum agreement subtrees of the following three kinds: maximum agreement subtrees of  $(M_i, N_j)$  for some  $i, j$ , maximum agreement subtrees of  $(T_1(u_i), N_j)$  for some  $i, j$ , and maximum agreement subtrees of  $(M_i, T_2(v_j))$  for some  $i, j$ .

**Property 2:**  $\mathcal{A}$  contains at most one maximum agreement subtree of the second kind and at most one of the third kind. Also  $\mathcal{A}$  contains at least one maximum agreement subtree of either the second or the third kind. More specifically, in Case 3.1, there is one of each kind, in Cases 1 and 3.2, there is one of the second kind but none of the third kind, and in Cases 2 and 3.3, there is one of the third kind but none of the second kind. Finally, Case 3.0 can be viewed as either kind since  $M_p = T_1(u_p)$  and  $N_q = T_1(v_q)$ .

**Property 3:** In Case 3, there are zero or more maximum agreement subtrees of the first kind, all of which occur above subtrees of the second and third kinds. Here, by

above we refer to the relative positions of the nearest ancestors  $u_i, v_j$  on the centroid paths. In Cases 1 and 2, there are no subtrees of the first kind.

**Property 4:** If two maximum agreement subtrees of the first kind occur in  $\mathcal{A}$ , for instance, the maximum agreement subtrees of  $(M_i, N_j)$  and  $(M_{i'}, N_{j'})$ , then  $i < i'$  implies  $j < j'$ .

**Property 5:** If subtrees of both the second kind and the third kind exist in  $\mathcal{A}$ , for instance, the maximum agreement subtrees of  $(T_1(u_{i'}), N_j)$  and  $(M_i, T_2(v_{j'}))$ , respectively, then  $i < i'$  and  $j < j'$ .

To model these three different kinds of maximum agreement subtrees, we need three different kinds of edges in  $G(x)$ , namely white edges, red edges and green edges, respectively. The details of these edges are described next, followed by the definition of agreement matching which captures the above structural properties.

**4.2. Edges of  $G(x)$ .**  $G(x)$  is actually a multigraph, where each multiedge consists of three edges, a *white* edge, a *red* edge and a *green* edge, each of which has a distinct weight associated with it. A multiedge between  $u_i \in L(x)$ ,  $1 \leq i \leq p-1$  and  $v_j \in R(x)$ ,  $1 \leq j \leq q-1$ , exists if and only if  $M_i$  and  $N_j$  intersect. The white edge in this multiedge has weight equal to the size of the maximum agreement subtree of  $M_i$  and  $N_j$ . The red edge in this multiedge has weight equal to the size of the maximum agreement subtree of  $T_1(u_i)$  and  $N_j$ . The green edge in this multiedge has weight equal to the size of the maximum agreement subtree of  $M_i$  and  $T_2(v_j)$ . If  $u_p \in L(x)$  then there is a multiedge between  $u_p$  and  $v_j$  such that either  $j \neq q$  and  $u_p$ 's twin is in  $N_j$  or  $j = q$  and  $u_p$ 's twin is  $v_q$ ; all three edges in this multiedge have weight 1. In addition, there is a multiedge between  $u_i$  and  $v_q$  such that either  $i \neq p$  and  $v_q$ 's twin is in  $M_i$  or  $i = p$  and  $v_q$ 's twin is  $u_p$ ; all three edges in this multiedge have weight 1.

**4.3. Agreement Matchings in  $G(x)$ . Definitions.** We define a *proper crossing* in  $G(x)$  to be either a single red edge, a single green edge, or a red-green edge pair such that the two edges cross and further, the endpoint of the green edge in  $L(x)$  is above that of the red edge.

A matching in  $G(x)$  is an *agreement matching* if:

1. It has zero or more white edges and one proper crossing.
2. No white edge crosses any other edge; further, all white edges dominate the edges in the proper crossing.

See Fig.2. The weight of such a matching is just the sum of the weights of its edges. The following property of agreement matchings in  $G$  is crucial.

**The Key Property.** Each maximum weight agreement matching corresponds a maximum agreement subtree, and vice versa, as is made precise below.

LEMMA 4.1. *A maximum weight agreement matching  $\mathcal{M}$  containing only edges incident upon or below vertex  $w$  in  $R(x)$  corresponds to an agreement subtree  $\mathcal{A}$  of  $T_1, T_2(w)$ ,  $w \in \pi(x)$ , having the same weight.*

*Proof.* We associate with each white edge  $(u_i, v_j)$  the maximum agreement subtree of  $(M_i, N_j)$ . Similarly, we associate with each red edge  $(u_i, v_j)$  the maximum agreement subtree of  $(T_1(u_i), M_j)$  and with each green edge  $(u_i, v_j)$  the maximum agreement subtree of  $(N_i, T_2(v_j))$ . The tree  $\mathcal{A}$  is defined as follows.

$\mathcal{A}$  will have a path containing one vertex for each white edge (see Fig.2). These vertices occur in the same sequence from top to bottom as their corresponding edges.



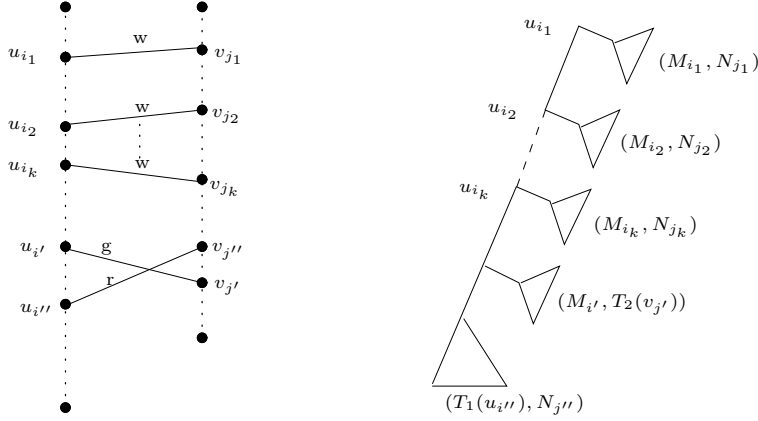


FIG. 4.1. An Agreement Matching with the Associated Agreement Tree.

The off-path children of these vertices will be the roots of the associated maximum agreement subtrees defined above. It remains to define the remaining child of the bottommost vertex on this path. This child will depend upon the nature of the proper crossing. We will define a tree  $\mathcal{A}'$  for the proper crossing as follows. This child will be the root of  $\mathcal{A}'$ .

If the proper crossing has a green edge and a red edge, then  $\mathcal{A}'$  is a tree whose left and right subtrees are the maximum agreement subtrees associated with the two edges. On the other hand, if the proper crossing has exactly one edge (red or green), then  $\mathcal{A}'$  is just the maximum agreement subtree associated with that edge.

And finally, if there are no white edges at all, then  $\mathcal{A} = \mathcal{A}'$ . The equivalence of weights is easy to check in all cases.  $\square$

LEMMA 4.2. A maximum agreement subtree  $\mathcal{A}$  of  $T_1$  and  $T_2(w)$ ,  $w \in \pi(x)$ , has a unique corresponding agreement matching which has the same weight and contains only edges incident upon or below  $w$  in  $R(x)$ .

*Proof.* We sketch how the agreement matching corresponding to  $\mathcal{A}$  is constructed. Recall Cases 1–3 and Properties 1–5 in Section 4.1.

If Case 1 occurs, then  $\mathcal{A}$  is of the maximum agreement subtree of  $(T_1, N_j)$ , for some  $v_j$  coinciding with or below  $w$  in  $\pi(x)$ . Let  $u_i$  be the topmost vertex in  $\pi$  such that  $M_i$  intersects with  $N_j$ . Then the maximum agreement subtree of  $(T_1(u_i), N_j)$  has the same weight as that of  $(T, N_j)$ . Further, there is a red edge between  $u_i$  and  $v_j$  in  $G(x)$ ; this red edge constitutes the agreement matching. Clearly, the weight of this matching is identical to that of  $\mathcal{A}$ .

If Case 2 occurs, a similar argument shows that the matching comprises a solitary green edge with the same weight as  $\mathcal{A}$ . In Case 3, a similar argument can be used to show that the matching contains a sequence of zero or more non-crossing white edges, lying above a proper crossing.  $\square$

Thus, in order to determine the maximum agreement subtree of  $T_1$  and  $T_2(w)$ ,  $w \in \pi(x)$ , it suffices to determine the maximum weight agreement matching in  $G(x)$  containing only edges incident upon or below vertex  $w$  in  $R(x)$ .

**4.4. Finding Maximum Weight Agreement Matching in  $G(x)$ .** To describe the algorithm for matching  $\pi$ , we need the following definitions followed by an

important theorem. The theorem itself will be proved in Section 7.

**Definitions.** The *degree* of a vertex in  $G(x)$  is defined as the number of white edges incident on it. Let  $d_x(u_i)$  denote the degree of  $u_i$ . A vertex in  $L(x)$  is called a *singleton* vertex if it has degree one; the white edge incident on it is called a *singleton* edge. Let  $nswe(x)$  denote the number of non-singleton white edges in  $G(x)$ . Let  $nsav(x)$  denote the number of vertices in  $R(x)$  which have at least one incident non-singleton white edge. Let  $SV(x)$  denote the set of singleton vertices in  $L(x)$ .

**THEOREM 4.3.** *Consider a particular  $x \in X$ . For each  $u_i \in L(x)$ , the largest weight agreement matching in  $G(x)$  containing only edges incident on or below  $u_i$  in  $L(x)$  can be found in time*

$$O\left(\sum_{i|d_x(u_i)>1} d_x(u_i) \log \frac{nsav(x)}{d_x(u_i)} + \sum_{(u_i, v_j) \in G(x) | d_x(u_i)=1} \log \frac{|T(x)|}{n_j}\right).$$

*Further, for each  $v_j \in R(x)$ , the largest weight agreement matching in  $G(x)$  containing only edges incident on or below  $v_j$  in  $R(x)$  can also be found in the same time.*

Theorem 4.3 is achieved by storing the vertices of  $R(x)$  in an appropriately weighted search tree. The construction is described in Section 7.

**Algorithm Outline for Matching  $\pi$ .** The matching graphs  $G(x)$  for all  $x \in X$  will be constructed in time proportional to the sum of the sizes of these graphs. This construction is described in Section 6. Then each matching graph  $G(x)$  is processed as follows (see Theorem 4.3). For each  $u_i \in L(x)$ , the largest weight agreement matching in  $G(x)$  containing only edges incident on or below  $u_i$  in  $L(x)$  is found. Further, for each  $v_j \in R(x)$ , the largest weight agreement matching in  $G(x)$  containing only edges incident on or below  $v_j$  in  $R(x)$  is also computed. This computation of agreement matchings is described in Section 7. For each  $w \in T_2$ , the largest agreement subtree of  $T_1$  and  $T_2(w)$  can be determined easily from the above information as it is given by the largest weight agreement matching in  $G(x)$  comprising only edges incident upon or below vertex  $w$  in  $R(x)$ . Section 5 shows that the total time taken above is  $O(\sum_{i=1}^p m_i \log \frac{n}{m_i})$ , as required.

**Inferring Maximum Agreement Subtrees.** Consider a vertex  $w \in T_2$ ; let  $x$  be the beginning of the centroid path in  $T_2$  containing  $w$ . Then  $w \in R(x)$ . The maximum agreement subtree of  $T_1$  and  $T_2(w)$  is given by the largest weight agreement matching in  $G(x)$  comprising only edges incident upon or below vertex  $w$  in  $R(x)$ .

**5. The Analysis.** We need the following preliminary lemmas before beginning the analysis.

**LEMMA 5.1.** *Consider graph  $G(x)$ . Then*

$$\sum_{i|d_x(u_i)>1} d_x(u_i) \log \frac{nsav(x)}{d_x(u_i)} \leq \sum_{i|d_x(u_i)>1} d_x(u_i) \log \frac{n}{m_i}.$$

*Proof.* Multiplying each side by  $\ln 2$ , we get the following equivalent inequality:

$$A = \sum_{i|d_x(u_i)>1} d_x(u_i) \ln \frac{nsav(x)}{d_x(u_i)} \leq \sum_{i|d_x(u_i)>1} d_x(u_i) \ln \frac{n}{m_i} = B.$$

Note that  $\sum_{i|d_x(u_i)>1} d_x(u_i) = nswe(x)$ . Let  $\alpha_i(x) > 0$  be such that  $d_x(u_i) = \alpha_i(x) \frac{m_i}{n} nswe(x)$ . Then  $\sum_{i|d_x(u_i)>1} \alpha_i(x) m_i = n$ . Also note that  $nsav(x) \leq nswe(x)$ .

Therefore,

$$A \leq B - \sum_{i|d_x(u_i)>1} \alpha_i(x) \frac{m_i}{n} \text{nswe}(x) \ln \alpha_i(x).$$

It suffices to show that

$$C = \sum_{i|d_x(u_i)>1} \alpha_i(x) m_i \ln \alpha_i(x) \geq 0.$$

We split  $C$  into two terms,

$$C_1 = \sum_{i|d_x(u_i)>1, \alpha_i(x) \geq 1} \alpha_i(x) m_i \ln \alpha_i(x)$$

and

$$C_2 = \sum_{i|d_x(u_i)>1, 0 < \alpha_i(x) < 1} \alpha_i(x) m_i \ln \alpha_i(x)$$

$$C_1 \geq \sum_{i|d_x(u_i)>1, \alpha_i(x) \geq 1} (\alpha_i(x) - 1) m_i.$$

Further,

$$\begin{aligned} & \sum_{i|d_x(u_i)>1, \alpha_i(x) \geq 1} (\alpha_i(x) - 1) m_i - \sum_{i|d_x(u_i)>1, \alpha_i(x) < 1} (1 - \alpha_i(x)) m_i \\ &= \sum_{i|d_x(u_i)>1} (\alpha_i(x) - 1) m_i = \sum_{i|d_x(u_i)>1} \alpha_i(x) m_i - \sum_{i|d_x(u_i)>1} m_i \geq n - n \geq 0. \end{aligned}$$

Therefore  $C \geq \sum_{i|d_x(u_i)>1, 0 < \alpha_i(x) < 1} (1 - \alpha_i(x) + \alpha_i(x) \ln \alpha_i(x)) m_i \geq 0$ .  $\square$

Recall that  $S_i$  denotes the subtree of  $T_2$  induced by the leaves of  $M_i$ .

LEMMA 5.2.  $\sum_{x \in X | d_x(u_i) > 1} d_x(u_i) = O(m_i)$ .

*Proof.* Consider  $G(x)$  such that  $d_x(u_i) > 1$ . Then all but one of the vertices of  $R(x)$  adjacent to  $u_i$  are also in  $S_i$ ; this is because  $M_i$  intersects both the right and the left subtrees of all but the bottommost of the vertices adjacent to  $u_i$  in  $R(x)$ . Since each vertex in  $S_i$  is in at most one matching graph  $G(x)$  and since  $|S_i| = m_i$ , the lemma follows.  $\square$

Consider a vertex  $u_i \in \pi$ . From Theorem 4.3 and Lemma 5.1, the following work is assigned to  $u_i$  when considering matching graph  $G(x)$ ,  $x \in X$ .

1. If  $M_i$  and  $T_2(x)$  do not intersect then no work is assigned to  $u_i$  as  $u_i$  is not in  $G(x)$ .
2. If  $d_x(u_i) = 1$  then the work assigned to  $u_i$  is  $O(\log \frac{|T_2(x)|}{n_j})$ , where  $v_j$  is the vertex in  $\pi(x)$  adjacent to  $u_i$ .
3. If  $d_x(u_i) > 1$  then the work assigned to  $u_i$  is  $O(d_x(u_i) \log \frac{n}{m_i})$ .

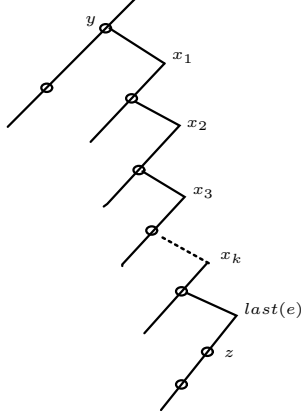


FIG. 5.1. Portion of  $T_2$  showing vertices in  $H(e)$  and  $last(e)$ .

The following is a corollary of Lemma 5.2 and the above bounds.

**COROLLARY 5.3.** *The work assigned to vertex  $u_i$  over all matching graphs  $G(x)$  with  $d_x(u_i) > 1$  is  $O(m_i \log \frac{n}{m_i})$ .*

It now suffices to account for the work assigned to vertex  $u_i$  over all matching graphs  $G(x)$  with  $d_x(u_i) = 1$ . Next, we show that this work is also  $O(m_i \log \frac{n}{m_i})$ . We use the tree  $S_i$  for this analysis.

**Analyzing over  $S_i$ .** Note that for each  $x \in X$  such that  $d_x(u_i) = 1$ ,  $x$  is not in  $S_i$ , i.e., either it lies on the path in  $T_2$  between the endpoints of some edge  $e$  in  $S_i$  or it lies on the path in  $T_2$  between the root of  $S_i$  and the root of  $T_2$ . In the former case,  $x$  is said to *lie* on edge  $e$  of  $S_i$ . To handle the latter case, we add a dummy edge to  $S_i$  connecting its root to a new node, representing the root of  $T_2$ . This new node now becomes the root of  $S_i$ . This addition of the dummy edge is only for the next few paragraphs, until the work done on that edge is accounted. Subsequent references to  $S_i$  will not have the dummy edge.

Consider the maximal subset  $H(e)$  of vertices  $x$  in  $X$  which lie on edge  $e$  of  $S_i$  and for which  $d_x(u_i) = 1$ . Let  $H(e) = \{x_1, x_2, \dots, x_k\}$ ; here the vertices appear in increasing order of distance from the root of  $T_2$ . Let  $e = (y, z)$ , with  $y$  the parent of  $z$  in  $S_i$ . Let  $first(e) = x_1$ . Let  $last(e)$  be the first vertex  $x$  in  $X$  such that  $d_x(u_i) > 1$  and  $x$  is on the path from  $x_k$  to  $z$  in  $T_2$ , if any; otherwise let  $last(e)$  be  $z$  (and then  $z$  is a leaf). See Fig.3.

**LEMMA 5.4.** *The sum of  $|T_2(first(e))|$  over any subset of edges of  $S_i$ , no two of which lie on the same root-to-leaf path in  $S_i$ , is  $O(n)$ .*

*Proof.* The above subtrees of  $T_2$  are all disjoint.  $\square$

**LEMMA 5.5.** *The work assigned to  $u_i$  on edge  $e$ , i.e., in processing graphs  $G(x)$ ,  $x \in H(e)$ , is  $O(\log \frac{|T_2(first(e))|}{|T_2(last(e))|})$ .*

*Proof.* The work assigned to  $u_i$  in processing  $G(x_j)$  is  $O(\log \frac{|T_2(x_j)|}{|T_2(x_{j+1})|})$  for  $1 \leq j < k$  and  $O(\log \frac{|T_2(x_k)|}{|T_2(last(e))|})$  for  $j = k$ . Thus the sum of the work assigned to  $u_i$  at the graphs  $G(x_j)$  is  $O(\log \frac{|T_2(first(e))|}{|T_2(last(e))|})$ .  $\square$

**COROLLARY 5.6.** *The work assigned to  $u_i$  on the dummy edge  $e$  is  $O(\log \frac{n}{m_i})$*

*Proof.*  $|T_2(last(e))| \geq m_i$ .  $\square$

It remains to analyze the work assigned to  $u_i$  over the non-dummy edges in  $S_i$ . From now onwards, we can ignore the dummy edge in  $S_i$ .

LEMMA 5.7. *Consider edges  $e, e' \in S_i$  such that  $e$  is on the path from  $e'$  to the root of  $S_i$ . If  $H(e), H(e')$  are non-empty, then  $|T_2(\text{last}(e))| \geq |T_2(\text{first}(e'))|$ .*

*Proof.*  $\text{first}(e')$  is a descendant of  $\text{last}(e)$  in  $T_2$ .  $\square$

We claim that sum of the work assigned to  $u_i$  over all the edges of  $S_i$  is  $O(m_i \log \frac{n}{m_i})$ . We show this next by applying tree contraction on  $S_i$ .

**Removal step.** First, remove all edges  $e$  in  $S_i$  incident upon leaves in  $S_i$ . The work done on these edges is bounded by the sum over all such edges  $e$  of  $O(\log |T_2(\text{first}(e))|)$ ; further, the number of such edges is at most  $m_i$ . By Lemma 5.4, this sum is at most  $O(m_i \log \frac{n}{m_i})$ .

**Contract step.** Next, contract all paths consisting only of degree two vertices in  $S_i$  into a single edge. The  $H$  set for such an edge  $e$  is defined to be the union of the  $H$  sets for the edges comprising the path which was contracted to give  $e$ . The work done on  $e$  is also defined to be the sum of the work done on the relevant edges.  $\text{first}(e)$  and  $\text{last}(e)$  are again defined as before. As is easily seen, Lemma 5.4, Lemma 5.5, and Lemma 5.7 hold for the new contracted tree as well. Further, this new tree has at most  $m_i/2$  leaves.

**Wrapping Up.**  $O(\log m_i)$  phases of the removal and contract steps are performed. In the  $j$ th phase, the work done on the edges removed is  $O(\frac{m_i}{2^{j-1}} \log \frac{n2^{j-1}}{m_i})$ . Summing up over all phases, we get

LEMMA 5.8. *The work assigned to vertex  $u_i$  over all matching graphs  $G(x)$  with  $d_x(u_i) = 1$  is  $O(m_i \log \frac{n}{m_i})$ .*

The following lemma is needed in the next section.

LEMMA 5.9. *The total number of edges incident on  $u_i$  over all matching graphs is  $O(m_i \log \frac{n}{m_i})$ .*

*Proof.* Recall that the work attributed to a singleton edge is of the form  $\log \frac{|T_2(x)|}{n_j} \geq 1$ . Thus, Lemma 5.8 implies that there are  $O(m_i \log \frac{n}{m_i})$  singleton edges incident to  $u_i$ , and by Lemma 5.2, there are  $O(m_i)$  non-singleton edges incident to  $u_i$ .  $\square$

THEOREM 5.10. *There is an algorithm for the Maximum Agreement Subtree Problem for two binary trees with an  $O(n \log n)$  running time.*

*Proof.* Corollary 5.3 and Lemma 5.8 imply that the total work assigned to  $u_i$  is  $O(m_i \log \frac{n}{m_i})$ . Hence the total matching cost from Step 3 in Section 3 is  $O(\sum m_i \log \frac{n}{m_i})$ . Further Step 1 takes linear time. Also, in Section 8, it will be shown that we can construct the recursive subproblems in Step 2 in linear time, so if we exclude the recursive calls, our cost is bounded by  $c \sum m_i \log \frac{n}{m_i}$  for some sufficiently large  $c$ . Inductively, we assume that each recursive call takes at most  $cm_i \log m_i$  time, and then the total time is at most

$$c \sum_i \left( m_i \log \frac{n}{m_i} + m_i \log m_i \right) = cn \log n,$$

as desired.  $\square$  The above analysis assumes (a) that we can construct the recursive subproblems in linear time, which will be done in Section 8, (b) that we can construct the matching graphs in time proportional to their sizes, which will be done in Sections 6 and 8, and (c) that Theorem 4.3 holds true, which will be proved in Sections 7 and 9.

**6. Constructing the Matching Graphs.** We show how all the matching graphs can be set up in time proportional to the sum of their sizes, which by Lemma 5.9 is  $O(\sum_{i=1}^{p-1} m_i \log \frac{n}{m_i})$ . First, we show how to set up the vertices and edges in each graph. Then we show how the weights on the edges are computed.

**Preprocessing.**  $T_2$  is preprocessed in linear time to compute a pointer from each vertex to the beginning of the centroid path containing it. It is also preprocessed to enable induced subtree computations in the same time bounds.

**6.1. Setting up Vertices and Edges..** The matching graphs in which vertex  $u_i$  appears along with the multiedges incident upon it in these graphs are determined in time proportional to the sum of the number of such multiedges over all such graphs as follows.

**Processing  $u_p$ .** First, consider the leaf  $u_p$  of  $T_1$ . The only matching graphs containing  $u_p$  are those which correspond to centroid paths beginning at vertices  $x$  of  $T_2$  such that  $x$  is an ancestor of the twin of  $u_p$  in  $T_2$ . Further, if  $u_p \in L(x)$  then there is a multiedge between  $u_p$  and vertex  $y \in R(x)$  if and only if  $y$  is the nearest ancestor of  $u_p$ 's twin in the centroid path beginning at  $x$ . Thus the matching graphs to which  $u_p$  belongs and the multiedges incident on  $u_p$  in these graphs can be determined in time proportional to the number of such graphs, given pointers from each vertex in  $T_2$  to the beginning of the centroid path containing it.

Lemmas 6.1 and 6.2 are needed for the next step.

**LEMMA 6.1.** *If vertex  $v_j$  in the centroid path beginning at vertex  $x$  of  $T_2$  is in  $S_i$  then  $u_i$  is adjacent to  $v_j$  in  $G(x)$ .*

*Proof.* Clearly, if  $j \neq q$  then  $M_i$  and  $N_j$  intersect and if  $j = q$  then  $v_j$ 's twin is in  $M_i$ .  $\square$

**LEMMA 6.2.** *If vertex  $v_j$  in the centroid path beginning at vertex  $x$  of  $T_2$  is not in  $S_i$  then  $u_i$  is adjacent to  $v_j$  in  $G(x)$  if and only if  $v_j \neq v_q$  and there exists some vertex  $y \in S_i$  which is in  $N_j$ .*

*Proof.* We assume that  $j \neq q$ . For if  $j = q$  then  $v_j = v_q$  is a leaf of  $T_2$  and since it does not appear in  $S_i$ , its twin is not in  $M_i$ , and therefore, there is no edge between  $u_i$  and  $v_j$ .

First, suppose some vertex  $y \in S_i$  is in  $N_j$ . Then, clearly,  $N_j$  intersects  $M_i$ . Therefore, there must be an edge between  $u_i$  and  $v_j$  in  $G(x)$ . Next, suppose that there is such an edge. Then  $N_j$  intersects  $M_i$ . Therefore, there exists some  $y \in S_i$  which is in  $N_j$ .  $\square$

**Processing  $u_i, 1 \leq i \leq p-1$ .** The subtrees of  $T_2$  induced by leaves of each  $M_i$  are computed in  $O(\sum_{i=1}^{p-1} m_i)$  time as described in Section 8. Let  $S_i$  denote this induced subtree. For each vertex  $z$  in  $S_i$ , perform the following in  $T_2$  until a vertex in the centroid path containing the parent of  $z$  in  $S_i$  is reached: repeatedly jump to the parent of the beginning of the centroid path in  $T_2$  containing the current vertex. By Lemmas 6.1 and 6.2, there is a multiedge from  $u_i$  to each vertex  $y$  of  $T_2$  (in the corresponding matching graph containing  $y$ ) encountered in this following procedure. Thus this procedure takes time proportional to the sum of the number of multiedges incident on  $u_i$  over all matching graphs it lies in, given pointers from each vertex in  $T_2$  to the beginning of the centroid path containing it.

**Remark.** For an edge between  $u_i$  and  $v_j, i \neq p, j \neq q$ , define  $map(i, j)$  as follows. If  $v_j \in S_i$  then  $map(i, j) = v_j$ . Otherwise, if  $v_j \notin S_i$ , then  $map(i, j)$  is that vertex in  $S_i$  which is closest to the root of  $S_i$  and a descendant of  $v_j$  in  $T_2$ . Note that  $map(i, j)$  can be easily computed in the course of the above procedure.

**6.2. Determining Edge Weights in  $G(x)$ .** Recall that for a multiedge between  $u_i$  and  $v_j$  in  $G(x)$ , we need to determine the sizes of the maximum agreement subtrees of the following pairs of trees.

1.  $M_i, N_j$ : white edge weight.
2.  $T_1(u_i), N_j$ : red edge weight.
3.  $M_i, T_2(v_j)$ : green edge weight.

Also recall that the multiedge itself indicates that  $M_i$  and  $N_j$  intersect.

Assume that the agreement matchings in graphs  $G(x')$  have already been determined, where  $x'$  is a descendant of  $x$  in  $T_2$ . Using this information and the information computed in Step 2, we show how the above required information can be computed for multiedge  $(u_i, v_j)$  in graph  $G(x)$  in constant time. Recall that in Step 2, the maximum agreement subtrees of  $M_i$  and the subtrees rooted at each vertex  $w$  of  $S_i$  were determined.

**White Edge Weight.** Let  $y = \text{map}(i, j)$ . If  $y \neq v_j$  then the maximum agreement subtree of  $M_i$  and the subtree of  $S_i$  rooted at  $y$  gives the desired information. Suppose  $y = v_j$ , i.e.,  $y \in S_i$ . Let  $z$  be the child of  $y \in S_i$  such that  $z \in N_j$ . The maximum agreement subtree of  $M_i$  and the subtree of  $S_i$  rooted at  $z$  gives the desired information in this case. This takes constant time.

**Green Edge Weight.** The maximum agreement subtree of  $M_i$  and the subtree of  $S_i$  rooted at  $y = \text{map}(i, j)$  gives the desired information in constant time.

**Red Edge Weight.** Let  $y$  be the root of  $N_j$ . Recall that the agreement matchings in graphs  $G(x')$  have already been determined, where  $x' \in X$  is a descendant of  $x$  in  $T_2$ . Since  $y \in X$ , agreement matchings in graph  $G(y)$  would already have been computed. Recall from Theorem 4.3 that for each vertex in  $L(y)$ , the maximum weight agreement matching containing only edges incident on or below that vertex in  $L(y)$  has been computed.

Note that since  $M_i$  intersects with  $T_2(y)$  (since a multiedge exists between  $u_i$  and  $v_j$ )  $u_i \in L(y)$ . The largest weight agreement matching in  $G(y)$  containing only edges incident on or below vertex  $u_i$  in  $L(y)$  gives the desired information. This information is computed as graph  $G(y)$  was processed, so it can be accessed in constant time now.

**7. Computing Agreement Matchings.** Consider graph  $G(x)$ . Recall that for each vertex in  $L(x)$ , we need to compute the largest weight agreement matching containing only edges incident on or below it in  $L(x)$ , and likewise for each vertex in  $R(x)$ . We outline the algorithm before giving details. The algorithm is similar to that in [FPT95a], but the data structure we use and the associated operations are different.

**Algorithm Outline.** First, a weight balanced binary search tree  $\mathcal{T}$  whose leaves are the vertices in  $R(x)$  is set up; here, the vertices in  $R(x)$  are given appropriate weights yet to be described. Next, the vertices in  $L(x)$  are considered in turn in bottom-to-top order. For each vertex  $u_i \in L(x)$ , the vertices adjacent to it in  $R(x)$  are searched for in  $\mathcal{T}$ ; the largest weight agreement matching with each white edge incident on  $u_i$  as topmost edge is found in the course of this search, as is the largest weight proper crossing for each green edge incident on  $u_i$ . From the above information, the largest weight agreement matching containing only edges incident on or below  $u_i$  in  $L(x)$  is easily found. Following the above search, the information stored in  $\mathcal{T}$  is updated. The time taken for processing  $u_i$  will be  $O(d_x(u_i) \log \frac{nsav(x)}{d_x(u_i)})$  if  $d_x(u_i) > 1$  and  $O(\log \frac{|T(x)|}{n_j})$  if  $d_x(u_i) = 1$  and  $u_i$  is adjacent to  $v_j$  in  $G(x)$ . After all vertices in

$L(x)$  have been processed, the vertices in  $R(x)$  are processed. For all such vertices  $v_j$ , the largest weight agreement matching containing only edges incident on or below  $v_j$  in  $R(x)$  are found in  $O(|R(x)|)$  time by a single scan of  $\mathcal{T}$ . The bounds in Theorem 4.3 follow.

**Weighted Search Tree  $\mathcal{T}$ .** Vertex  $v_j \in R(x)$  is given weight  $n_j + \frac{|T(x)|}{nsav(x)}$  if some non-singleton edge in  $G(x)$  is incident upon it, and weight  $n_j$ , otherwise. The sum of the weights of vertices in  $R(x)$  is at most  $2|T(x)|$ . The construction of  $\mathcal{T}$  using these weights is dealt with in Section 9.

Tree  $\mathcal{T}$  has the following three crucial characteristics.

1.  $\mathcal{T}$  can be constructed in  $O(|R(x)|)$  time.
2. Searching for  $v_j$  in  $\mathcal{T}$  takes  $O(\log \frac{|T(x)|}{n_j})$  time.
3. Searching for an ordered subset  $\{v_{j_1}, \dots, v_{j_k}\}$  of  $R(x)$ , each vertex in which has an incident non-singleton edge, takes  $O(k \log \frac{nsav(x)}{k})$  time. The procedure used here is to first search for  $v_{j_1}$  starting at the root, then search for  $v_{j_2}$  starting at  $v_{j_1}$  in the obvious way, and so on.

**Auxiliary Information in  $\mathcal{T}$ .** We maintain the following auxiliary information at each internal vertex in  $\mathcal{T}$ . Recall that we process the vertices of  $L(x)$  in order; an edge of  $G(x)$  is said to be *in*  $\mathcal{T}$  if its endpoint in  $L(x)$  has already been processed. Further, we say that an edge of  $G(x)$  is *in*  $\mathcal{T}(z)$  if it is in  $\mathcal{T}$  and its endpoint in  $R(x)$  is located in  $\mathcal{T}(z)$ . Let  $anc(z)$  denote the set of ancestors of  $z$  in  $\mathcal{T}$ ,  $z$  inclusive. For a leaf  $v_j \in \mathcal{T}$ ,  $lfringe(v_j)$  is the set of vertices in  $\mathcal{T}$  which are left children of vertices in  $anc(v_j)$  but not themselves in  $anc(v_j)$ .  $rfringe(v_j)$  is defined analogously.

The following information is maintained at each vertex  $z$  of  $\mathcal{T}$ .

1.  $g(z)$ : For each  $z$ ,  $\max_{z' \in anc(z)} g(z')$  will be the heaviest green edge in  $\mathcal{T}$  which forms a proper crossing with each red edge in  $\mathcal{T}(z)$ .
2.  $x(z)$ : This is largest weight proper crossing among the edges in  $\mathcal{T}(z)$ .
3.  $m(z)$ : This is largest weight agreement matching containing a white edge such that the topmost white edge is in  $\mathcal{T}(z)$ .
4.  $y(z)$ : This the largest weight proper crossing such that the green edge in this crossing is in  $\mathcal{T}$  but not in  $\mathcal{T}(z)$ , the red edge in this crossing is in  $\mathcal{T}(z)$ , and the green edge does not form a proper crossing with all the red edges in  $\mathcal{T}(z)$ .
5.  $r(z)$ : This is simply the heaviest red edge in  $\mathcal{T}(z)$ .

Next, we show how vertex  $u_i$  is processed, given that vertices below it in  $L(x)$  have been processed. For the moment assume that  $d_x(u_i) = 1$ . The case when  $d_x(u_i) > 1$  will be addressed later.

**Case 1.**  $d_x(u_i) = 1$ . Let  $v_j$  be the only vertex to which  $u_i$  is adjacent. First,  $v_j$  is found in  $\mathcal{T}$ ; this takes  $O(\log \frac{|T(x)|}{n_j})$  time. Next, the white, red, and green edges incident on  $u_i$  are processed as described below in the same time bound. An important fact to note is that in each case, the information in  $\mathcal{T}$  will be read and updated only at vertices in the set  $anc(v_j)$  and vertices which are children of vertices in this set; the time taken in this process will be proportional to the depth of  $v_j$ , i.e.,  $O(\log \frac{|T(x)|}{n_j})$ .

**Processing White Edge  $e = (u_i, v_j)$ .** First, the largest weight agreement matching with  $e$  as the topmost edge is determined. Then the  $m()$  values at vertices in  $anc(v_j)$  are updated according to the weight of this matching. All other information remains unchanged.

The above desired matching is computed as follows. There are two cases. Either this matching contains another white edge. The largest weight matching among all



such matchings is given by  $1 + \max_{z \in lfringe(v_j)} m(z)$ . The other case occurs when this matching contains only edge  $e$  plus a proper crossing. Thus, it suffices to compute the largest proper crossing containing edges dominated by  $e$ . This is given by  $\max\{\max_{z \in lfringe(v_j)} x(z), \max_{z \in lfringe(v_j)} (\max_{z' \in anc(z)} g(z')) + r(z), \max_{z \in lfringe(v_j)} y(z)\}$ . The first term here is the largest weight proper crossing in which both edges are in  $\mathcal{T}(z)$  for some  $z \in lfringe(v_j)$ . The second term is the largest weight proper crossing in which the red edge is in  $\mathcal{T}(z)$ , for some  $z \in lfringe(v_j)$ , the green edge is not in this subtree but it forms a proper crossing with each red edge in this subtree. The third term is the largest weight proper crossing in which the red edge is in  $\mathcal{T}(z)$ , for some  $z \in lfringe(v_j)$ , the green edge is not in this subtree and it does not form a proper crossing with some red edge in this subtree.

**Processing Red Edge**  $e = (u_i, v_j)$ . The  $m()$  and  $x()$  values remain unchanged in  $\mathcal{T}$ . Next, note that no green edge already in  $\mathcal{T}$  can form a proper crossing with  $e$ . This implies that the  $y()$  and  $g()$  values for  $z \in anc(v_j)$  need to be modified.

Consider  $y(z)$  first,  $z \in anc(v_j)$ . A green edge in  $\mathcal{T}$  which formed a proper crossing with all red edges in  $\mathcal{T}(z)$  does not do so any more. So  $y(z)$  is set to  $\max\{y(z), (\max_{z' \in anc(z)} g(z')) + r(z)\}$ .

Consider  $g(z)$  next,  $z \in anc(v_j)$ .  $g(z)$  is set to  $\phi$ . Before this is done,  $g(y)$  is updated to  $\max_{y' \in anc(y)} g(y')$  for each  $y \in lfringe(v_j)$  and  $y \in rfringe(v_j)$ . The invariant on  $g()$  is easily seen to be maintained.

Finally  $r(z)$  is set to  $\max\{r(z), wt(e)\}$ , for each  $z, z \in anc(v_j)$ .

**Processing Green Edge**  $e = (u_i, v_j)$ . Note that  $e$  can form a proper crossing with only those red edges in  $\mathcal{T}$  which are in  $\mathcal{T}(z)$ ,  $z \in rfringe(v_j)$ ; further,  $e$  forms a proper crossing with each such red edge. Therefore,  $g(z)$  is set to  $\max\{\max_{z' \in anc(z)} g(z'), wt(e)\}$  for each  $z \in rfringe(v_j)$ .

For each  $z \in anc(v_j)$ ,  $x(z)$  is then set to the larger of the current value and  $\max(wt(e) + r(z'))$ , the maximum being taken over all vertices  $z' \in rfringe(v_j)$  which are descendants of  $z$ . Also note that  $\max_{z \in rfringe(v_j)} (wt(e) + r(z))$  gives the largest weight proper crossing containing  $e$ .

**Case 2:**  $d_x(u_i) = k > 1$ . Suppose  $u_i$  is adjacent to  $v_{j_1}, v_{j_2}, \dots, v_{j_k}$ , in bottom to top order. Then these vertices are searched for in sequence in  $\mathcal{T}$ . This takes  $O(k \log \frac{nsav(x)}{k})$  time by the procedure mentioned earlier, i.e., first search for  $v_{j_1}$  starting at the root, then search for  $v_{j_2}$  starting at  $v_{j_1}$  in the obvious way, and so on. In the above process, all vertices in the set  $\{z | z \in (anc(v_{j_1}) \cup anc(v_{j_2}) \cup \dots \cup anc(v_{j_k}))\}$  are traversed. Again, as in Case 1, only information at vertices in the above set and at children of vertices in the above set needs to be read and updated. This takes time proportional to the size of the above set, which, in turn, is  $O(k \log \frac{nsav(x)}{k})$ .

**Processing  $R(x)$ .** It remains to show how, for each  $v_j \in R(x)$ , the largest weight agreement matching containing only edges incident on or below  $v_j$  in  $R(x)$ , is computed.

For each  $v_j \in R(x)$ , we find the largest weight agreement matching with some white edge incident upon  $v_j$  as the dominant edge and the largest weight proper crossing containing some red edge incident on  $v_j$ . This information clearly suffices. The first of the above two is given simply by  $m(v_j)$ . The second is given by  $\max\{y(z), \max_{z' \in anc(v_j)} g(z') + r(v_j)\}$ . Over all  $v_j \in R(x)$ , the computation of the above two values can be accomplished in a single pass of  $\mathcal{T}$  in  $O(|R(x)|)$  time.

**8. Computing Induced Subtrees.** We show how to preprocess a tree in  $O(|T|)$  time so that given any subset  $L$  of its leaves in order, the subtree induced by  $L$  can be computed in  $O(|L|)$  time. The construction is a generalization of the proof of Lemma 5.2 in [FT95].

$T$  is preprocessed for *Least Common Ancestor* (LCA) queries in  $O(|T|)$  time. This enables the computation of the LCA of any two leaves of  $T$  in constant time [HT84]. The distance of each vertex from the root of  $T$  is also computed; call this quantity the *depth* of a vertex.

Given the ordered set of leaves  $L = l_1, l_2, \dots, l_{|L|}$ , the following steps are executed. First, the LCA  $l'_i$  of each pair of consecutive leaves  $l_i, l_{i+1}$ ,  $1 \leq i \leq |L| - 1$ , is found; the  $l'_i$ 's will be the internal vertices in the subtree induced by  $L$ . Next, the edges between vertices are set up as follows.

For each vertex  $v$  in the sequence  $l_1, l'_1, l_2, l'_2, \dots, l_{|L|-1}, l'_{|L|-1}, l_{|L|}$ , two vertices  $v_{left}$  and  $v_{right}$  are computed.  $v_{left}$  is the nearest vertex to the left of  $v$ , if any, which has depth strictly less than  $v$ .  $v_{right}$  is defined analogously. This computation is easily accomplished in  $O(|L|)$  time. Finally, edges are put between  $v$  and one of  $v_{left}, v_{right}$ , whichever has greater depth. If exactly one of  $v_{left}, v_{right}$  is defined (this will happen only for vertices on the paths from the root to the leftmost and rightmost leaves in the induced subtree) then an edge is put between  $v$  and the vertex which is defined. The root of the induced tree will be the unique vertex for which both  $v_{left}$  and  $v_{right}$  are undefined; no edges need be put in this case.

**Step 2 of the Main Algorithm.** Step 2 (see Section 3) requires finding the induced trees  $S_i$  for each  $M_i$ ,  $1 \leq i \leq p - 1$ , in  $O(\sum_i^{p-1} m_i)$  time. This is done in two steps (essentially this procedure is described in [FT95]). First, the leaves of each  $M_i$  are sorted by the order in which their twins occur in  $T_2$ . This is done by bucket sorting all the leaves of  $T_1$  by the order in which their twins occur in  $T_2$ , and then bucket sorting them (in a stable way) by the order in which the trees  $M_i$  to which they belong occur in  $T_1$ . This takes  $O(\sum_i^{p-1} m_i)$  time.

Next, for each  $M_i$ ,  $1 \leq i \leq p - 1$ , the subtree of  $T_2$  induced by the leaves of  $M_i$  is found using the algorithm described above in  $O(m_i)$  time. The total time taken is  $O(\sum_i^{p-1} m_i)$ .

**9. The Weighted Search Tree.** We will now complete our solution to the maximum agreement subtree problem, by describing the weighted search trees from Section 7. Recall that we are given vertices  $v_1, v_2, \dots, v_{|R(x)|}$ , such that vertex  $v_j$  has weight  $w(v_j) = n_j + \frac{|T(x)|}{nsav(x)}$  if it has an incident non-singleton edge, and weight  $w(v_j) = n_j$ , otherwise. The sum of the vertex weights is bounded by  $2|T(x)|$ .

Theorem 9.1 from [Fre75, Meh77] shows that the weight balanced tree  $\mathcal{T}$  can be constructed in  $O(|R(x)|)$  time.

**THEOREM 9.1.** *Given weights  $w_1, \dots, w_n$  with sum  $W$ , a binary tree such that the depth of the  $i$ th leaf is  $O(1 + \log(W/w_i))$  can be constructed in  $O(n)$  time. In this tree, the total weight of all leaves in the subtree rooted at any node  $z$  will be at most half of the corresponding weight for the subtree rooted at the grandparent of  $z$ .*

It follows from Theorem 9.1 that the time to search for  $v_j$  in  $\mathcal{T}$  is  $O(1 + \log \frac{|2T(x)|}{w(v_j)}) = O(1 + \log \frac{|T(x)|}{n_j})$ .

Next consider the case when an ordered subset  $\{v_{j_1}, \dots, v_{j_k}\}$  of vertices is given, each having an incident non-singleton edge. The algorithm to search for these vertices is to first start from the root and search for  $v_{j_1}$ , then start from  $v_{j_1}$  and search for

$v_{j_2}$ , then start from  $v_{j_2}$  and search for  $v_{j_3}$ , and so on. Each search is performed in the obvious way. For, technical reasons, we return to the root at the end.

Each edge in  $\mathcal{T}$  is traversed at most twice during the above search, once in each direction.

Consider the topological subtree formed by the traversed edges. It has  $k$  leaves and  $k - 1$  internal nodes with two children. These nodes form a tree, with each edge in the tree corresponding to a path in the search tree. We associate with a node in the topological subtree the path in the search tree corresponding to the edge from the node to its parent in the topological subtree. The associated path for the root node of the topological subtree is the path from this node to the root of the search tree. We will give an upper bound on the total number of vertices on these paths excluding their endpoints. To do this, we give a lower bound on the total weight of the “off-path” subtrees for each path. (An “off-path” subtree for a node is simply the subtree which does not contain the continuation of the path.)

Let  $l$  be the number of internal vertices on one such path associated with node  $v$ . By Theorem 9.1, the sum of the weights of every second root of the off-path subtrees is at least  $(2w(v) - w(v)) + (4w(v) - 2w(v)) + \dots + (2^{\lfloor l/2 \rfloor} w(v) - 2^{\lfloor l/2 \rfloor - 1} w(v)) = (2^{\lfloor l/2 \rfloor} - 1)w(v)$ . But  $w(v) \geq \frac{|T(x)|}{nsav(x)}$  and the total weight is at most  $2T(x)$ . Simple calculus shows that the sum of these lower bounds on the path lengths is maximized when the terms  $w(v)$  are all at their minimum value and the path weights are all equal at  $2T(x)/(2k - 1)$ . This gives path lengths of  $O(\log \frac{nsav(x)}{2k-1})$  and hence a total path length of  $O(k \log \frac{nsav(x)}{k})$ .

**10. Concluding remarks.** We can generalize our technique to higher degree bounds  $d > 2$ , by combining it with techniques from [FT95, Section 2] for unbounded degrees. This appears to yield an algorithm with running time  $O(\min\{n\sqrt{d}\log^2 n, nd\log n \log d\})$ . We conjecture, however, that there is an algorithm with running time  $O(n\sqrt{d}\log n)$ .

**Acknowledgment.** We would like to thank the referees from SICOMP for some very thorough comments.

#### REFERENCES

- [AK97] A. AMIR, D. KESELMAN. *Maximum agreement subtree in a set of evolutionary trees*. SIAM Journal on Computing, 26(6), pp. 1656-1669, 1997.
- [CR96] R.. COLE, R. HARIHARAN. *An  $O(n \log n)$  algorithm for the maximum agreement subtree problem for binary trees*. Proc. of the 7th ACM-SIAM SODA, pp. 323-332, 1996.
- [FPT95a] M. FARACH, T. PRZYTYCKA, M. THORUP. *The maximum agreement subtree problem for binary trees*. Proc. of 2nd ESA, 1995.
- [FPT95b] M. FARACH, T. PRZYTYCKA, M. THORUP. *Agreement of many bounded degree evolutionary trees*. Information Processing Letters, 55(6), pp. 297-301, 1995.
- [FT95] M. FARACH, M. THORUP. *Fast comparison of evolutionary trees*. Information and Computation, 123(1), pp. 29-37, 1995.
- [FT97] M. FARACH, M. THORUP. *Sparse dynamic programming for evolutionary-tree comparison*. SIAM Journal on Computing, 26(1), pp. 210-230, 1997.
- [FG85] C. R. FINDEN, A. D. GORDON. *Obtaining common pruned trees*. Journal of Classification, 2, pp. 255-276, 1985.
- [Fre75] M. L. FREDMAN. *Two applications of a probabilistic search technique: sorting  $X + Y$  and building balanced search trees*. Proc. of the 7th ACM STOC, pp. 240-244, 1975.
- [GY95] R. GRISHMAN, R. YANGARBER. Private Communication, NYU, 1995.

- [HT84] D. HAREL AND R.E. TARJAN. *Fast algorithms for finding nearest common ancestor*. SIAM Journal on Computing, 13(2), pp. 338–355, 1984.
- [Ka95] M-Y. KAO. *Tree contractions and evolutionary trees*. SIAM Journal on Computing, 27(6), pp. 1592–1616, 1998.
- [KKM95] E. KUBICKA, G. KUBICKI, F. R. MCMORRIS. *An algorithm to find agreement subtrees*. Journal of Classification, 12, pp. 91–100, 1995.
- [Meh77] K. MEHLHORN. *A best possible bound for the weighted path length of binary search trees*. SIAM Journal on Computing, 6(2), pp. 235–239, 1977.
- [SW93] M. STEEL, T. WARNOW. *Kaikoura tree theorems: computing the maximum agreement subtree*. Information Processing Letters, 48, pp. 77–82, 1993.