

# Finding Smallest Supertrees Under Minor Containment <sup>\*</sup>

Naomi Nishimura <sup>†</sup>   Prabhakar Ragde <sup>‡</sup>   Dimitrios M. Thilikos <sup>§</sup>

## Abstract

The diversity of application areas relying on tree-structured data results in a wide interest in algorithms which determine differences or similarities among trees. One way of measuring the similarity between trees is to find the smallest common superstructure or supertree, where common elements are typically defined in terms of a mapping or embedding. In the simplest case, a supertree will contain exact copies of each input tree, so that for each input tree, each vertex of a tree can be mapped to a vertex in the supertree such that each edge maps to the corresponding edge. More general mappings allow for the extraction of more subtle common elements captured by looser definitions of similarity.

We consider supertrees under the general mapping of minor containment. Minor containment generalizes both subgraph isomorphism and topological embedding; as a consequence of this generality, however, it is NP-complete to determine whether or not  $G$  is a minor of  $H$ , even for general trees. By focusing on trees of bounded degree, we obtain an  $O(n^3)$  algorithm which determines the smallest tree  $T$  such that both of the input trees are minors of  $T$ , even when the trees are assumed to be unrooted and unordered.

---

<sup>\*</sup>Research supported by the Natural Sciences and Engineering Research Council of Canada and Communications and Information Technology Ontario.

<sup>†</sup>Department of Computer Science, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1. Email: nishi@uwaterloo.ca

<sup>‡</sup>Department of Computer Science, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1. Email: pragde@uwaterloo.ca

<sup>§</sup>Departament de Llenguatges i Sistemes Informàtics, Universitat Politècnica de Catalunya, Campus Nord – Mòdul C5 – Desp. 211b, c/Jordi Girona Salgado, 1-3. E-08034, Barcelona, Spain. Email: sedthilk@lsi.upc.es Research supported by the Ministry of Education and Culture of Spain, Grant number MEC-DGES SB98 0K148809.

# 1 Introduction

The breadth of algorithmic research on trees stems from both the simplicity of the structure and the variety of application domains. When information about a data set can be derived from its tree structure, comparisons among two or more data sets can entail determining similarities among two or more trees. Algorithms of this type have been developed in areas such as compiler design, structured text databases, theory of natural languages, computer vision [18], and computational biology (the reader is directed to a previous paper on trees [10] for further references).

Comparisons of trees range from the classical tree pattern matching problem (finding an exact copy of one tree in another) to numerous variants, including problems on multiple trees and inexact matches. Each problem can be viewed as finding a way to relate trees by mappings, where trees are related if it is possible to map vertices to sets of vertices and edges to sets of edges subject to certain constraints. Researchers have considered different types of trees (ordered, unordered, labeled, unlabeled) and different mappings between pairs of trees (exact matching, approximate matching, subgraph isomorphism, topological embedding, minor containment) [3, 5, 9, 13, 14]. In addition, researchers have measured the similarity between trees by finding the largest common subtree or smallest common supertree under various constraints [1, 4, 7, 8, 10, 12, 19].

In this paper we consider the problem of finding the smallest common supertree under minor containment. Concisely, a graph  $G$  is a minor of a graph  $H$  if it is possible to map all the vertices in  $G$  to mutually disjoint connected subgraphs in  $H$  and there exists a bijection, from the edges of  $G$  to the edges of  $H$  that are not in any of these subgraphs, such that the images of the endpoints of any edge  $e$  in  $G$  contain the endpoints of the image of  $e$  through this bijection; equivalently we can view the mapping as taking edges to paths. Minor containment is of interest due to its generality; it encompasses both subgraph isomorphism and topological embedding and is fundamental in the work of Robertson and Seymour on graph minors [17]. However, due in large part to the generality, many problems which are tractable under subgraph isomorphism and topological embedding are NP-complete for minor containment. In particular, it is NP-complete to determine whether or not one tree is a minor of another [6], but this can be determined in polynomial time when there is a degree bound of  $O(\log n / \log \log n)$  [9]. We thus restrict our attention to trees of bounded degree, noting that the resultant supertree will also be of bounded degree (in contrast, a common subtree of two bounded degree trees may not have bounded degree).

Interest in supertrees under minor containment arises from their applications to editing, image clustering, genetics, chemical structure analysis, and evolution [12, 19]. Previous algo-

rithms to find supertrees have been limited to special cases: in ordered minor containment, there is an order imposed on the children of each node in each input tree, and this order must be preserved by the mapping [12]; for evolutionary trees, the leaves have distinct labels and are constrained to map to other leaves [19].

## 2 Preliminaries

Each input to our algorithm is a bounded-degree tree (an undirected graph with no cycles).  $V(T)$  denotes the vertices of  $T$  and  $E(T)$  the edges of  $T$ . A tree  $T$  may be rooted at a distinguished vertex  $r$ ; in this case we can view the rooted tree as a directed graph, with children and parents defined as in standard graph-theoretic references [2]. When processing rooted trees we will consider a *subtree*  $T_v$  of  $T$ , defined to be the subgraph of  $T$  induced by  $v$  and all its descendants. More generally, for  $A$  a subset of the children of some node  $v$ , we define  $T_A$  to be the subgraph induced by  $v$ , the vertices in  $A$ , and all descendants of nodes in  $A$ . For  $A$  an arbitrary subset of vertices,  $T[A]$  is defined to be the subgraph of  $T$  induced by  $A$ .

Given input trees  $Q$  and  $R$ , we wish to find a tree  $T$  such that both  $Q$  and  $R$  are minors of  $T$  and  $T$  is as small as possible. There are several equivalent definitions of minors; the most relevant one for our purposes is given below. Intuitively, a graph  $G$  is a minor of a graph  $H$  (or  $H$  is a major of  $G$ ) if  $H$  can be obtained from  $G$  by a series of vertex and edge deletions and edge contractions, where a contraction of an edge  $(u, v)$  in  $G$  is the operation that replaces  $u$  and  $v$  by a new vertex whose neighbors are the vertices that were adjacent to  $u$  or  $v$ . It is not difficult to see that, for trees, the following definition is equivalent:

**Definition:** A tree  $Q$  is a *minor* of a tree  $T$  if and only if there exists a surjection  $f : V(T) \rightarrow V(Q)$  such that

1. for each  $a \in V(Q)$ ,  $T[f^{-1}(a)]$  is connected;
2. for each pair  $a, b \in V(Q)$ ,  $f^{-1}(a) \cap f^{-1}(b) = \emptyset$ ; and
3. for  $S = \{(u, v) \in E(T) \mid f(u) \neq f(v)\}$ , there exists a bijection  $\xi : S \rightarrow E(Q)$  such that for each  $e = (s, t) \in S$ ,  $\xi(e) = (f(s), f(t))$ .

We call  $f$  a *minor embedding* of  $T$  into  $Q$ . Intuitively,  $f^{-1}(a)$  is the set of vertices of  $T$  contracted into  $a$ ; (2) captures the notion that each vertex of  $T$  corresponds to exactly one vertex of  $Q$ ; and (3) captures the notion that uncontracted edges of  $T$  are preserved in  $Q$ .

The problem we wish to solve is that of determining the smallest common acyclic major of  $Q$  and  $R$ , henceforth called the *smallest common tree major*. For  $\text{sctmj}(Q, R)$  the minimum number of vertices in a common tree major of  $Q$  and  $R$ , it is not difficult to see that  $\max\{|V(Q)|, |V(R)|\} \leq \text{sctmj}(Q, R) \leq |V(Q)| + |V(R)|$ . We observe that  $\text{sctmj}(Q, R) = |Q|$  if and only if  $R$  is a minor of  $Q$ . Duchet [6] proved that it is NP-complete to determine whether one tree is a minor of another. It is now easy to prove that deciding whether  $\text{sctmj}(Q, R) \leq k$  for two general trees  $Q, R$  is NP-complete. In view of this, we will restrict our attention to the case where the input graphs are both trees with maximum degree bounded by a fixed constant.

In the remainder of the paper we will make use of the following notational conventions. Since we will be finding a graph  $T$  such that  $Q$  and  $R$  are both minors of  $T$ , we will use  $f$  to denote the minor embedding of  $T$  into  $Q$  and  $g$  to denote the minor embedding of  $T$  into  $R$ . We will use letters near the beginning of the alphabet for vertices of  $Q$  and letters near the end of the alphabet for vertices of  $R$ .

### 3 Expansions

To facilitate understanding of the algorithm, it is beneficial to consider the mappings between  $Q, R$ , and a common tree major  $T$ . The edges of  $T$  correspond to edges in the input trees  $Q$  and  $R$ ; we distinguish between *strong* edges, which correspond to edges in both  $Q$  and  $R$ , and *weak* edges, each of which corresponds to an edge in only one of  $Q$  and  $R$ . For  $f$  and  $g$  the minor embeddings of  $T$  into  $Q$  and  $R$ , respectively,  $f^{-1}(a)$  and  $g^{-1}(u)$  describe connected subgraphs of  $T$ . Since for  $a \in V(G)$  each vertex in  $f^{-1}(a)$  is in  $g^{-1}(u)$  for some  $u \in V(R)$ , we can associate  $a$  with a set of vertices in  $V(R)$  with overlapping preimages. This notion is formalized in a graph called an *expansion* of  $Q$  and  $R$  consisting of edges between associated vertices. More formally:

**Definition:** For  $Q$  and  $R$  trees on disjoint sets of vertices, an *expansion* of  $Q$  and  $R$  is a bipartite graph  $\mathcal{E} = (V(\mathcal{E}), E(\mathcal{E}))$  with bipartition  $(V(Q), V(R))$  such that

1. the neighborhood in  $\mathcal{E}$  of any vertex of  $V(R)$  (respectively,  $V(Q)$ ) induces a connected subgraph of  $Q$  (respectively,  $R$ );
2.  $\mathcal{E}$  has no isolated vertices;
3. the neighborhoods in  $\mathcal{E}$  of two vertices in  $V(Q)$  (respectively,  $V(R)$ ) intersect in at most one vertex; and

4. for every edge  $(a, b)$  in  $E(Q)$ , either there are edges  $(a, u)$  and  $(b, u)$  in  $\mathcal{E}$  for some  $u \in V(R)$ , or there are edges  $(a, u)$  and  $(b, v)$  in  $\mathcal{E}$  for some edge  $(u, v) \in E(R)$  (and symmetrically for edges in  $R$ ).

Given an expansion  $\mathcal{E}$  of  $Q$  and  $R$ , we define  $T_{\mathcal{E}}$  to be a graph whose vertices are edges in  $\mathcal{E}$  and whose edges are formed by condition 4 in the definition above. For an edge  $(a, b) \in E(Q)$ , if there are edges  $(a, u)$  and  $(b, u)$  in  $\mathcal{E}$ , then  $\{(a, u), (b, u)\}$  is an edge in  $T_{\mathcal{E}}$ , and if there are edges  $(a, u)$  and  $(b, v)$  in  $\mathcal{E}$  for some  $(u, v) \in E(R)$ , and neither  $(a, v)$  nor  $(b, u)$  is in  $\mathcal{E}$ , then  $\{(a, u), (b, v)\}$  is in  $T_{\mathcal{E}}$ . Edges  $(u, v) \in E(R)$  define edges in  $T_{\mathcal{E}}$  in a similar fashion. In the former case we call the edge  $(a, b)$  *weak*; in the latter case,  $(a, b)$  and  $(u, v)$  are *strong*.

We will denote the weak (strong) edges of  $Q$  as  $\text{weak}(Q)$  ( $\text{strong}(Q)$ ) and we will use the analogous notation for  $R$  as well. Note that there is a natural bijection  $f_{\mathcal{E}}$  between strong edges in  $E(Q)$  and strong edges in  $E(R)$ . We call the edges of  $T_{\mathcal{E}}$  that are defined on the basis of weak edges of  $Q$  ( $R$ ),  $Q$ -*weak* ( $R$ -*weak*). If an edge of  $T_{\mathcal{E}}$  is not  $Q$  or  $R$ -weak, then we call it *strong*. There exist a natural bijection between the weak edges of  $Q$  ( $R$ ) and the  $Q$ -weak ( $R$ -weak) edges of  $T_{\mathcal{E}}$  and a natural bijection between the strong edges of  $Q$  ( $R$ ) and the strong edges of  $T_{\mathcal{E}}$ . Finally, as direct consequences of the definition of weak and strong edges,  $|\text{strong}(Q)| = |\text{strong}(R)|$  and  $|E(T_{\mathcal{E}})| = |\text{weak}(Q)| + |\text{weak}(R)| + |\text{strong}(Q)|$ . We define  $|E(T_{\mathcal{E}})|$  to be the *size of the expansion*  $\mathcal{E}$ .

For convenience, if  $\mathcal{E}$  is an expansion of two trees  $Q$  and  $R$ ,  $(a, b)$  is a strong edge of  $Q$ , and  $(u, v) = f_{\mathcal{E}}((a, b))$ , we will say that  $(a, b)$  and  $(u, v)$  are  $\mathcal{E}$ -*counterparts* of each other and conclude that  $(a, u), (b, v) \in \mathcal{E}$ . Finally, given a vertex  $t$  in  $T_{\mathcal{E}}$  which corresponds to an edge  $(a, u)$  of  $\mathcal{E}$  where  $a \in V(Q)$  and  $u \in V(R)$ ,  $a$  is the  $Q$ -*side* of  $t$  and  $u$  is the  $R$ -*side* of  $t$ .

The proof of the following lemma is a direct consequence of the definition of an expansion and is omitted.

**Lemma 3.1.** *Let  $\mathcal{E}_i$  be a minimum size expansion of two trees  $Q_i$  and  $R_i$  for  $i = 1, 2$ ,  $V(Q_1) \cap V(Q_2) = \{a\}$ ,  $V(R_1) \cap V(R_2) = \{u\}$ , and  $(a, u) \in \mathcal{E}_1 \cap \mathcal{E}_2$ . Then  $\mathcal{E}_1 \cup \mathcal{E}_2$  is an expansion of minimum size (among those containing  $(a, u)$ ) of  $Q_1 \cup Q_2$  and  $R_1 \cup R_2$ .*

The following two lemmas are direct applications of Lemma 3.1.

**Lemma 3.2.** *Let  $\mathcal{E}_i$  be a minimum size expansion of two trees  $Q_i$  and  $R_i$ ,  $a_i \in Q_i$  for  $i = 1, 2$ ,  $V(Q_1) \cap V(Q_2) = \emptyset$ ,  $V(R_1) \cap V(R_2) = \{u\}$ ,  $(a_1, u) \in \mathcal{E}_1$ , and  $(a_2, u) \in \mathcal{E}_2$ . Then  $\mathcal{E}_1 \cup \mathcal{E}_2$  is an expansion of minimum size (among those containing at least one of  $(a_1, u)$  and  $(a_2, u)$ ) of the graph with vertex set  $V(Q_1) \cup V(Q_2)$  and edge set  $E(Q_1) \cup E(Q_2) \cup \{(a_1, a_2)\}$  and  $R_1 \cup R_2$ .*

**Lemma 3.3.** *Let  $\mathcal{E}_i$  be a minimum size expansion of two disjoint trees  $Q_i$  and  $R_i$ ,  $a_i \in V(Q_i)$ ,  $u_i \in V(R_i)$ , and  $(a_i, u_i) \in \mathcal{E}_i$  for  $i = 1, 2$ . Then  $\mathcal{E}_1 \cup \mathcal{E}_2$  is an expansion of minimum size (among those containing at least one of  $(a_1, u_1)$  and  $(a_2, u_2)$ ) of the graph with vertex set  $V(Q_1) \cup V(Q_2)$  and edge set  $E(Q_1) \cup E(Q_2) \cup \{(a_1, a_2)\}$  and the graph with vertex set  $V(R_1) \cup V(R_2)$  and edge set  $E(R_1) \cup E(R_2) \cup \{(u_1, u_2)\}$ .*

Our algorithm will rely on relationships between neighborhoods of sets. We use  $N_G(v)$  to denote the neighborhood of the vertex  $v$  in the graph  $G$ . We say that two subsets  $S_1, S_2$  of the vertex set of a graph  $G$  are *touching* if either  $S_1 \cap S_2 \neq \emptyset$  or there exists an edge  $(v_1, v_2) \in E(G)$  for  $v_i \in S_i, i = 1, 2$ .

**Lemma 3.4.** *For any expansion  $\mathcal{E}$  of  $Q$  and  $R$  and any edge  $e = (a_1, a_2) \in E(Q)$  ( $e \in E(R)$ ),  $N_{\mathcal{E}}(a_1)$  and  $N_{\mathcal{E}}(a_2)$  are touching.*

**Proof.** By condition 4 of the definition of  $\mathcal{E}$ , for any edge  $(a_1, a_2) \in E(Q)$  there will exist either a vertex  $u$  in  $R$  where  $(a_1, u), (a_2, u) \in E(\mathcal{E})$  or there will exist an edge  $(u_1, u_2) \in E(R)$  such that  $(a_1, u_1), (a_2, u_2) \in E(\mathcal{E})$ . In the first case the connected graphs  $R[N_{\mathcal{E}}(a_1)]$  and  $R[N_{\mathcal{E}}(a_2)]$  have a common point  $u$  and in the second that they contain  $u_1$  and  $u_2$  respectively and  $(u_1, u_2) \in E(R)$ . Therefore, in both cases, their vertex sets are touching. ■

The lemma below is a useful tool in proving properties of expansions; it shows that if two pairs of nodes are related by an expansion, the paths joining the nodes are also related. In the remainder of the paper we use  $P_G(p_1, p_2)$  to denote the set of nodes in the (unique) path between vertices  $p_1$  and  $p_2$  in the tree  $G$ .

**Lemma 3.5.** *For any expansion  $\mathcal{E}$  of  $Q$  and  $R$ , if  $(a_i, u_i) \in \mathcal{E}, i = 1, 2$ , then any vertex in  $P_Q(a_1, a_2)$  has a neighbor in  $\mathcal{E}$  in  $P_R(u_1, u_2)$ .*

**Proof.** We will prove the lemma by contradiction, using induction on  $j$ , the size of  $P_Q(a_1, a_2)$ . Since the lemma holds trivially for  $j = 2$ , it suffices to show that the lemma holds for  $j = k$  assuming that it holds for all values  $j < k$ .

Suppose that there exist vertices in  $P_Q(a_1, a_2)$  whose sets of neighbors in  $\mathcal{E}$  do not intersect  $P_R(u_1, u_2)$ . We will call such vertices *bad* vertices and all other vertices in  $P_Q(a_1, a_2)$  *good*.

We first observe that if any interior vertex in  $P_Q(a_1, a_2)$  is a good vertex, then we can show that every vertex on the path has a neighbor in  $P_R(u_1, u_2)$ . That is, if  $b$  is a good vertex with neighbor  $v$  in  $P_R(u_1, u_2)$ , then we can apply the induction hypothesis on the smaller problem  $P_Q(a_1, b)$  and  $P_R(u_1, v)$  and also the smaller problem  $P_Q(b, a_2)$  and  $P_R(v, u_2)$  to reach our conclusion. We can now assume that every interior vertex in  $P_Q(a_1, a_2)$  is bad.

Furthermore, we can assume that there is no node in  $P_R(u_1, u_2)$  which is a neighbor of both  $a_1$  and  $a_2$ , since if there were such a node  $v$ , then by property 1 in the definition of  $\mathcal{E}$ , every node in  $P_Q(a_1, a_2)$  would also be in the neighborhood of  $v$ . Thus,  $N_{\mathcal{E}}(a_1) \cap N_{\mathcal{E}}(a_2) \cap P_R(u_1, u_2)$  is empty.

For each bad node  $a$ , we can define a vertex  $v(a)$  in  $P_R(u_1, u_2)$  which is the vertex in  $P_R(u_1, u_2)$  closest to  $N_{\mathcal{E}}(a)$  in  $R$ ; this vertex is unique due to property 1 in the definition of  $\mathcal{E}$ . We let  $b_i$  be the neighbor of  $a_i$  in  $P_Q(a_1, a_2)$  and show that  $v(b_i) \in N_{\mathcal{E}}(a_i) \cap P_R(u_1, u_2)$ . Suppose instead  $v(b_i) \notin N_{\mathcal{E}}(a_i) \cap P_R(u_1, u_2)$ . As  $R$  is a tree, we can partition the vertices of  $R \setminus P_R(u_1, u_2)$  into connected subgraphs on the basis of the closest vertex in  $P_R(u_1, u_2)$ . Since  $N_{\mathcal{E}}(b_i) \cap P_R(u_1, u_2) = \emptyset$ ,  $N_{\mathcal{E}}(b_i)$  must be contained entirely in one partition, namely that associated with  $v(b_i)$ . We observe that  $v(b_i)$  is a cutset separating  $N_{\mathcal{E}}(a_i)$  and  $N_{\mathcal{E}}(b_i)$  and not contained in either set. This contradicts Lemma 3.4, which states that since  $(a_i, b_i) \in E(Q)$ ,  $N_{\mathcal{E}}(a_i)$  and  $N_{\mathcal{E}}(b_i)$  are touching.

By a similar argument we can show that if  $a$  and  $b$  are bad neighbors in  $Q$ , then  $v(a) = v(b)$ . Since there is a path from  $b_1$  to  $b_2$ ,  $v(b_1) = v(b_2)$ . Since for  $i = 1, 2$ ,  $v(b_i) \in N_{\mathcal{E}}(a_i) \cap P_R(u_1, u_2)$ , then  $v(b_1) \in N_{\mathcal{E}}(a_1) \cap N_{\mathcal{E}}(a_2) \cap P_R(u_1, u_2)$ , which we proved to be empty. ■

**Lemma 3.6.** *If  $\mathcal{E}$  is an expansion of two trees  $Q$  and  $R$ , then  $T_{\mathcal{E}}$  is a common tree major of  $Q$  and  $R$ .*

**Proof.** We will prove first that  $T_{\mathcal{E}}$  is a tree. By property 1 of the definition of  $\mathcal{E}$ , for any vertex  $a$  in  $Q$ ,  $N_{\mathcal{E}}(a)$  induces in  $R$  a tree  $T^a$ , and hence the number of edges of  $\mathcal{E}$  with  $a$  as endpoint is equal to  $|E(T^a)| + 1$ . Moreover, all the edges in  $T^a$  are weak edges of  $R$  and any weak edge  $e$  of  $R$  is in some tree  $T^b$  where  $b$  is the vertex of  $Q$  adjacent to both endpoints of  $e$ . In addition, any edge of  $R$  belongs to only one tree  $T^a$  induced by the neighborhood, in  $\mathcal{E}$ , of some vertex  $a$  of  $Q$ . As a consequence of the above observations,

$$\begin{aligned} |V(T_{\mathcal{E}})| &= |E(\mathcal{E})| = \sum_{a \in V(Q)} (|\text{weak edges in } R[N_{\mathcal{E}}(a)]| + 1) \\ &= |V(Q)| + |\text{weak}(R)| = 1 + |E(Q)| + |\text{weak}(R)| \\ &= 1 + |\text{weak}(Q)| + |\text{strong}(Q)| + |\text{weak}(R)| \\ &= 1 + |E(T_{\mathcal{E}})| \end{aligned}$$

To show that  $T$  is a tree, it remains to show that  $T_{\mathcal{E}}$  is connected. Let  $t_1, t_2$  be two vertices in  $T_{\mathcal{E}}$  and let  $a_1$  and  $a_2$  be their  $Q$ -sides. We will use induction on  $j = |P_Q(a_1, a_2)|$ . Suppose first that  $j = 2$  and let  $(a_1, u_1)$  and  $(a_2, u_2)$  be the edges of  $\mathcal{E}$  corresponding to  $t_1$  and  $t_2$  respectively. By Lemma 3.4,  $N_{\mathcal{E}}(a_1)$  and  $N_{\mathcal{E}}(a_2)$  are touching. Therefore, there will

be in  $P_R(u_1, u_2)$  either a vertex  $u \in N_{\mathcal{E}}(a_1) \cap N_{\mathcal{E}}(a_2)$  or an edge  $(u, u')$  where  $u \in N_{\mathcal{E}}(a_1)$  and  $u' \in N_{\mathcal{E}}(a_2)$ . In the first case  $(a_1, u)$  and  $(a_2, u)$  and in the second  $(a_1, u)$  and  $(a_2, u')$  define two adjacent vertices  $t$  and  $t'$  of  $T_{\mathcal{E}}$ . In order to show that there exists a path in  $T_{\mathcal{E}}$  connecting  $t_1$  and  $t_2$ , we will prove that there exist two paths in  $T_{\mathcal{E}}$ , one connecting  $t_1$  with  $t$  and the other connecting  $t_2$  with  $t'$ . For any pair of edges  $(v_1, v_2), (v_2, v_3)$  of  $P_R(u_1, u)$ , there is a pair of edges  $(r_1, r_2)$  and  $(r_2, r_3)$  in  $E(T_{\mathcal{E}})$  where  $r_1, r_2, r_3$  correspond to  $(a_1, v_1), (a_1, v_2)$ , and  $(a_1, v_3)$  respectively. Using this observation, it is easy to see that  $t_1$  and  $t$  are connected in  $T_{\mathcal{E}}$ . The proof of the existence of a path connecting  $t$  and  $t_2$  in  $T_{\mathcal{E}}$  is similar and the base case of the induction holds.

Suppose now that the claim holds for  $j < k, k \geq 3$  and let  $t_1$  and  $t_2$  be two vertices in  $T_{\mathcal{E}}$  whose  $Q$ -sides are  $a_1$  and  $a_2$  and  $|P_Q(a_1, a_2)| = k$ . Let  $a'$  be the vertex in  $P_Q(a_1, a_2)$  that is adjacent to  $a_1$ . According to the definitions, there are two cases: (1)  $(a_1, a')$  is a strong edge with  $\mathcal{E}$ -counterpart  $(u^*, u')$  and thus there exist in  $T_{\mathcal{E}}$  two adjacent vertices  $r, t'$  corresponding to the edges  $(a_1, u^*)$  and  $(a', u')$  respectively, or (2)  $(a_1, a')$  is a weak edge whose endpoints are both connected to some vertex  $u^*$  in  $R$ , and there exist in  $T_{\mathcal{E}}$  two adjacent vertices  $r, t'$  corresponding to the edges  $(a_1, u^*)$  and  $(a', u^*)$ , respectively.

In either case, we can apply the induction hypothesis for  $t'$  and  $t_2$  since  $|P_Q(a', a_2)| < k$ . Therefore, there exists a path in  $T_{\mathcal{E}}$  connecting  $t'$  and  $t_2$  to which we can add edge  $(t', r)$  to form a path from  $r$  to  $t_2$  as well. It now remains to prove that there exists a path in  $T_{\mathcal{E}}$  connecting  $r$  and  $t_1$  in the case where  $r$  is different from  $t_1$ . The crucial property of  $r$  and  $t_1$  is that the edges of  $\mathcal{E}$  corresponding to them,  $(a_1, u_1)$  and  $(a_1, u^*)$ , both contain  $a_1$  as the  $Q$ -side. Since the neighborhood of  $a_1$  induces a tree  $R$ , there exists a path in this tree that connects  $u_1$  and  $u^*$ . Using the same arguments on  $t_1$  and  $r$  as we did for  $t_1$  and  $t$  in the base case, we can prove that there is a path in  $T_{\mathcal{E}}$  connecting  $t_1$  and  $r$  and therefore a path connecting  $t_1$  and  $t_2$ . Thus  $T_{\mathcal{E}}$  is connected and is a tree.

In order to prove that  $T_{\mathcal{E}}$  is a common major of  $Q$  and  $R$  we have to provide functions  $f$  and  $\xi$  as in definition 1. We define  $f : V(T_{\mathcal{E}}) \rightarrow V(Q)$ , such that  $f$  maps every vertex of  $T_{\mathcal{E}}$  to its  $Q$ -side and any edge in  $T_{\mathcal{E}}$  whose endpoints have different  $Q$ -sides to the edge of  $Q$  that connects them. The fact that condition 1 holds follows easily from the fact, observed above, that the vertices in  $T_{\mathcal{E}}$  with the same  $Q$ -side induce a connected subgraph of  $T_{\mathcal{E}}$ . Conditions 2 and 3 are direct consequences of the way  $T_{\mathcal{E}}$  is defined. The intuition behind the above definition of  $f$  is that a graph isomorphic to  $Q$  can be obtained from  $T_{\mathcal{E}}$  if we contract all the  $R$ -weak edges of  $T_{\mathcal{E}}$ . This proves that  $Q$  is a minor of  $T_{\mathcal{E}}$ . The proof that  $R$  is a minor of  $T_{\mathcal{E}}$  is symmetric. ■

**Lemma 3.7.** *For  $T$  a smallest common tree major of  $Q$  and  $R$ , there exists an expansion*



$\mathcal{E}$  such that  $T_{\mathcal{E}}$  is isomorphic to  $T$ .

**Proof:** Given minor embeddings  $f$  and  $g$  of  $T$  into  $Q$  and  $R$ , for each  $a \in V(Q)$  and each  $u \in V(R)$ ,  $|f^{-1}(a) \cap g^{-1}(u)| \leq 1$ , since otherwise the minor of  $T$  obtained after contracting the edges in the graph induced by  $\{f^{-1}(a) \cap g^{-1}(u)\}$  would be a smaller common tree major of  $Q$  and  $R$ . We define the expansion  $\mathcal{E}$  to be the set  $\{(a, u) : |f^{-1}(a) \cap g^{-1}(u)| = 1\}$ . It is straightforward to verify the claim that  $\mathcal{E}$  is an expansion of  $Q$  and  $R$ . ■

As a corollary of Lemmas 3.6 and 3.7, we can conclude that  $\text{sctmj}(Q, R)$  is the number of edges in the minimum expansion of  $Q$  and  $R$ . The following straightforward lemma reduces the problem to the computation of the rooted version of expansions.

**Lemma 3.8.** *For trees  $Q$  and  $R$  and for any  $a \in V(Q)$ ,  $\text{sctmj}(Q, R)$  is the minimum over all  $u \in V(R)$  of the number of edges in the smallest expansion  $\mathcal{E}$  of  $Q$  and  $R$  such that  $(a, u)$  is an edge in  $\mathcal{E}$ .*

We finish this section with the following useful observation:

**Lemma 3.9.** *For any trees  $Q$  and  $R$  where  $|E(Q)|, |E(R)| \geq 1$  and for any  $a \in V(Q)$  and  $u \in V(R)$ , the smallest expansion of  $Q$  and  $R$  that contains  $(a, u)$  as an edge has size smaller than  $|E(Q)| + |E(R)|$ .*

**Proof.** As  $|E(Q)|, |E(R)| \geq 1$ , there exist edges  $(a, b)$  and  $(u, v)$  with  $a$  and  $u$  as endpoints. Let  $Q_1$  and  $Q_2$  ( $R_1$  and  $R_2$ ) be the connected components of the graph formed by removing the edge  $(a, b)$  from  $Q$  (the graph formed by removing the edge  $(u, v)$  from  $R$ ) that contain  $a$  and  $b$  ( $u$  and  $v$ ) respectively. It is easy to verify that  $\mathcal{E} = (V(Q) \cup V(R), E)$  where

$$E = \{(c, u) \mid c \in V(Q_1)\} \cup \{(c, v) \mid c \in V(Q_2)\} \cup \\ \{(a, w) \mid w \in V(R_1)\} \cup \{(b, w) \mid w \in V(R_2)\}$$

is an expansion of  $Q$  and  $R$  containing  $(a, u)$ . Since  $|E| = |E(Q)| + |E(R)|$ ,  $|E(T_{\mathcal{E}})| = |E(Q)| + |E(R)| - 1$ . ■

## 4 Smallest common tree major algorithm

### 4.1 Algorithm overview

For algorithmic convenience, we construct a rooted tree major, where any node of either input tree could be associated with the root. We fix a root for one tree and then try all

possible rootings of the other tree; the following description concerns one possible choice of a root.

Our algorithm proceeds by dynamic programming, at each stage building tree majors of various subtrees of our inputs. After topologically sorting each tree with respect to the chosen root, we process each vertex  $a$  in  $V(Q)$  in order from leaves to root, pairing  $a$  with each  $u$  in  $V(R)$  in order from leaves to root.

For a given pair  $(a, u)$  we wish to determine the size of the largest common tree major  $T$  such that  $Q_a$  is a minor of  $T$  and  $R_u$  is a minor of  $T$  where for  $r$  the root of  $T$ ,  $f(r) = a$  and  $g(r) = u$ . We solve this problem using subproblems involving children of  $a$  and  $u$ , where in each subproblem we specify not only the roots of the subtrees of  $Q$  and  $R$ , but also the subsets of the children included thus far in the mapping.

Expansions, as defined in the previous section, give a convenient framework for expressing the progress of the algorithm, where expansions involving subgraphs of  $Q$  and  $R$  are augmented to form expansions of larger subgraphs of  $Q$  and  $R$ . The dynamic programming formulation of the problem relies on a set of subproblems at  $a \in V(Q)$  and  $u \in V(R)$ , where each subproblem corresponds to one choice of how the children of  $a$  and the children of  $u$  are related, assuming that  $(a, u)$  is to be an edge in the expansion and that all subproblems rooted at children have already been solved.

## 4.2 Technical lemmas

When processing  $(a, u)$ , we are assuming that  $(a, u) \in E(\mathcal{E})$  and attempting to see where subsets of the children of  $a$  and  $u$  can map. To build our intuition, we consider the process from the point of view of  $Q$  (viewing from  $R$  is symmetric and hence the reasoning identical). Each child  $b$  of  $a$  must eventually be involved in  $\mathcal{E}$ . There are four different cases for a child  $b$  of  $a$ , reflecting four different possible smaller expansions involving subtrees rooted at the children of  $a$  and  $u$  (for an illustration of the case analysis that follows, see Figure 1).

1. (epsilon child) The subtree rooted at  $b$  is not involved in any previous expansion. It will be included by creating an edge in  $\mathcal{E}$  from each vertex in the subtree to  $u$ .
2. (terminal child) The subtree rooted at  $b$  has been mapped to a subtree rooted at a child  $v$  of  $u$ , where  $(a, v)$  is not an edge in any previous expansion. In this case the edges  $(a, b)$  and  $(u, v)$  will be strong edges that are  $\mathcal{E}$ -counterparts.
3. (one-many child) The subtree rooted at  $b$  is mapped to subtrees rooted at a set of children of  $u$ , where  $(b, u)$  is an edge in a previous expansion. In this case  $(a, b)$  is a weak edge.

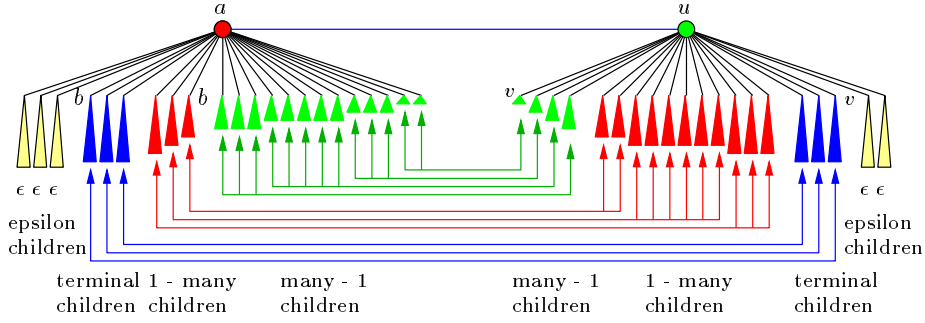


Figure 1: The different ways possible smaller expansions involving subtrees rooted at the children of  $a$  and  $u$  can be combined in a general expansion.

4. (many-one child) A set of subtrees rooted at children of  $a$  is mapped to a subtree rooted at a child  $v$  of  $u$ , where  $(a, v)$  is an edge in a previous expansion. In this case  $(u, v)$  is a weak edge.

We formalize the possible associations of children by a tuple for each possible pair of subsets  $A$  of children of  $a$  and  $X$  of children of  $u$  and each possible mapping among vertices.

**Definition:** Given two sets  $A, X$  we define  $\Pi(A, X)$  as the set containing all tuples

$$(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$$

that satisfy the following properties:

1.  $\{A^e, A^t, A^o, A^m\}$  is a partition of  $A$ ;
2.  $\{X^e, X^t, X^o, X^m\}$  is a partition of  $X$ ;
3.  $\tau : A^t \rightarrow X^t$  is a bijection;
4.  $\alpha : X^m \rightarrow A^o$  is a surjection; and
5.  $\chi : A^m \rightarrow X^o$  is a surjection.

It is not difficult to show that in the trees  $Q$  rooted at  $a$  and  $R$  rooted at  $u$ ,  $f_{\mathcal{E}}$  preserves the parent-child orientation of the strong edges. Suppose instead that  $(a_1, a_2)$  and  $(u_1, u_2)$  are  $\mathcal{E}$ -counterparts where  $a_1$  is in the path between  $a$  and  $a_2$  in  $Q$  but, in  $R$ ,  $u_2$  is in the path connecting  $u$  and  $u_1$ . Applying Lemma 3.5 for paths  $P_Q(a, a_2)$  and  $P_R(u, u_2)$ , we conclude  $a_1 \in P_Q(a, a_2)$  is adjacent, in  $\mathcal{E}$ , to a vertex in  $P_R(u, u_2)$ . But  $a_1$  is also adjacent to  $u_1$  (in

$\mathcal{E}$ ), and hence by property 1 of the definition of  $\mathcal{E}$   $a_1$  must be adjacent to  $u_2$ , a contradiction as  $(a_1, a_2)$  and  $(u_1, u_2)$  are strong edges. Using this observation we will always assume from now on that if  $(a_1, a_2)$  is strong and  $(u_1, u_2)$  is its  $\mathcal{E}$ -counterpart, then  $a_1$  ( $u_1$ ) is the endpoint closer to  $a$  ( $u$ ) in  $Q$  ( $R$ ). In general, whenever we mention an edge, the first endpoint of the pair will be the one that it is closer to the root of the tree to which it belongs.

We call two edges of a rooted tree *comparable* if one of them is in the path connecting the other with the root. If we have three mutually incomparable edges such that exactly two of them have a vertex different from the root as a common predecessor, we call the two edges the *close pair* of the triple.

**Lemma 4.1.** *For any expansion  $\mathcal{E}$  of two trees  $Q = Q_a$  and  $R = R_u$  such that  $(a, u) \in \mathcal{E}$ , if  $e_1, e_2$  are strong edges of  $Q$  and  $e'_1$  and  $e'_2$  are their  $\mathcal{E}$ -counterparts in  $R$ , then  $e_1$  is comparable with  $e_2$  if and only if  $e'_1$  is comparable with  $e'_2$ .*

**Proof.** We prove the lemma by contradiction. Without loss of generality,  $e_1 = (a_1, a_2)$  is in the path connecting  $e_2 = (a_3, a_4)$  and  $a$ , and  $e'_1 = (u_1, u_2)$  and  $e'_2 = (u_3, u_4)$  are not comparable. The incomparability of  $e'_1$  and  $e'_2$  means that  $u_2$  is not in  $P_R(u_1, u_4)$ . By applying Lemma 3.5 for paths  $P_Q(a_1, a_4)$  and  $P_R(u_1, u_4)$ ,  $a_2 \in P_Q(a_1, a_4)$  will be adjacent, in  $\mathcal{E}$ , to some vertex in  $P_R(u_1, u_4)$ . As  $a_2$  is also adjacent to  $u_2$  in  $\mathcal{E}$ , property 1 of the definition of  $\mathcal{E}$  requires that  $a_2$  be adjacent to  $u_1$ . Since  $e_1$  and  $e'_1$  are strong edges we have obtained a contradiction. The proof of the other direction is symmetric. ■

**Lemma 4.2.** *For any expansion  $\mathcal{E}$  of two trees  $Q = Q_a$  and  $R = R_u$  such that  $(a, u) \in \mathcal{E}$ , if  $e_1, e_2$  and  $e_3$  are strong mutually incomparable edges of  $Q$  and  $e'_1, e'_2$  and  $e'_3$  are their  $\mathcal{E}$ -counterparts in  $R$ , then  $e_1, e_2$  is the close pair of  $e_1, e_2, e_3$  if and only if  $e'_1, e'_2$  is the close pair of  $e'_1, e'_2, e'_3$ .*

**Proof.** In a proof by contradiction, we let  $e_i = (a_i, b_i), i = 1, 2, 3$ ,  $e'_i = (u_i, v_i), i = 1, 2, 3$ , and suppose that  $e_1$  and  $e_2$  form a close pair and  $e'_2$  and  $e'_3$  form a close pair. Let  $b$  be the common predecessor of  $e_1$  and  $e_2$  and  $v$  be the common predecessor of  $e'_2$  and  $e'_3$ . As a consequence of Lemma 3.5 on  $P_Q(b_1, b_2)$  and  $P_R(v_1, v_2)$ ,  $b$  must be adjacent in  $\mathcal{E}$  to some vertex in  $P_R(v_1, v_2)$ . Similarly, we can prove that  $b$  must be adjacent in  $\mathcal{E}$  to some vertex in  $P_R(v_1, v_3)$  and to some vertex in  $P_R(v_2, v_3)$ . It is not hard to see that as a consequence of these three facts and property 1 of the definition of  $\mathcal{E}$ ,  $b$  and  $v$  must be adjacent.

Using the same technique, by applying Lemma 3.5 to  $P_Q(b_2, b_3)$  and  $P_R(v_2, v_3)$ , we conclude that  $a \in P_Q(b_2, b_3)$  will be adjacent, in  $\mathcal{E}$ , to some vertex in  $P_R(v_2, v_3)$ . As  $a$  is also adjacent to  $u$ , by property 1 of the definition of  $\mathcal{E}$ ,  $a$  must be adjacent to  $v$ . By Lemma 3.5

for  $P_Q(b_1, b_2)$  and  $P_R(v_1, v_2)$ , by symmetry we can show that  $b$  is connected to  $u$  in  $\mathcal{E}$ . We have shown that  $(a, u), (a, v), (b, u), (b, v) \in E(\mathcal{E})$  which violates property 3 of the definition of  $\mathcal{E}$ . The proof of the other direction is symmetric. ■

Given a child  $b$  of  $a$  in  $Q_a$  we denote as  $\tilde{Q}_b$  the graph  $Q_b$  augmented with the edge  $(a, b)$ , and given a child  $v$  of  $u$  in  $R_u$  we denote as  $\tilde{R}_v$  the graph  $R_v$  augmented with the edge  $(u, v)$ .

**Lemma 4.3.** *For any expansion  $\mathcal{E}$  of two trees  $Q = Q_a$  and  $R = R_u$  such that  $(a, u) \in \mathcal{E}$ , there exists a tuple  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$  in  $\Pi(\text{children}(a), \text{children}(u))$ , such that the following hold:*

1. *there are no strong edges in  $Q_{A^e}$  or  $R_{X^e}$ ;*
2. *all edges from  $a$  to vertices in  $A^t$  and from  $u$  to vertices in  $X^t$  are strong;*
3. *all edges from  $a$  to vertices in  $A^o$  and from  $u$  to vertices in  $X^o$  are weak;*
4. *for all  $b \in A^t, v \in X^m$ , and  $c \in A^m$ ,  $f_{\mathcal{E}}$  maps  $(a, b)$  to  $(u, \tau(b))$ , the strong edges in  $Q_b$  to the strong edges in  $R_{\tau(b)}$ , the strong edges in  $R_v$  to the strong edges in  $Q_{\alpha(v)}$ , and the strong edges in  $Q_c$  to the strong edges in  $R_{\chi(c)}$ .*

**Proof.** We let  $E_a$  be the set of edges induced in  $Q$  by  $a$  and its children and let  $E_u$  be the set of edges induced in  $R$  by  $u$  and its children. To construct the desired partition, we first define sets  $A^e, A^t, X^e$ , and  $X^t$  as follows:  $A^e$  is the maximum subset of the children of  $a$  in  $Q$  with the property that for each  $b$  in  $A^e$ ,  $Q_b$  contains no strong edges;  $A^t$  consists of the children  $b$  of  $a$  such that  $(a, b)$  is the  $\mathcal{E}$ -counterpart of an edge  $(u, v)$  in  $E_u$ ;  $X^e$  and  $X^t$  are defined analogously. We form the bijection  $\tau$  by setting  $\tau(b) = v$  for  $(a, b)$  and  $(u, v)$   $\mathcal{E}$ -counterparts, for  $b \in A^t$ . We have now satisfied conditions 1 and 2, and it is straightforward to see that, for all  $b \in A^t$ ,  $f_{\mathcal{E}}$  maps  $(a, b)$  to  $(u, \tau(b))$ .

We now claim that for any  $b \in A^t$ , the strong edges of  $Q_b$  are mapped by  $f_{\mathcal{E}}$  to the strong edges of  $R_{\tau(b)}$ . Suppose instead that edge  $(a, b) \in E_a$ ,  $(c, d) \in E(Q_b)$ ,  $(u, v) = f_{\mathcal{E}}((a, b))$ , and  $(w, x) = f_{\mathcal{E}}((c, d)) \notin E(R_{\tau(b)})$  were a counterexample. Then, since  $(a, u), (d, x) \in E(\mathcal{E})$ , we can apply Lemma 3.5 to  $P_Q(a, d)$  and  $P_R(u, x)$  in order to conclude that the neighborhood of  $b$  in  $\mathcal{E}$  contains a vertex in  $P_R(u, x)$ . As  $(w, x)$  is not an edge of  $R_{\tau(b)}$ ,  $v = \tau(b)$  is not a vertex of this path. Since  $v$  is adjacent to  $b$  in  $\mathcal{E}$ , by property 1 of the definition of  $\mathcal{E}$ ,  $u$  must be a neighbor of  $b$  in  $\mathcal{E}$ . This results in a contradiction, as  $(a, u)$  and  $(b, v)$  are strong edges of  $Q$  and  $R$  respectively.

We now define  $A^o$  to include any child  $b$  of  $a$  for which  $\tilde{Q}_b$  contains strong edges whose  $\mathcal{E}$ -counterparts are in more than one of the trees  $\tilde{R}_w$  for children  $w$  of  $u$ . Notice that  $A^o$  and

$A^t$  are disjoint, as for any  $b \in A^t$ , the counterparts of the strong edges of  $\tilde{Q}_b$  are all in one tree  $\tilde{R}_w$ , namely  $\tilde{R}_{\tau(b)}$ .

We claim that for any  $b \in A^o$  the edge  $(a, b)$  is weak. Suppose instead that  $(a, b)$  were strong; since  $b \notin A^t$ , its  $\mathcal{E}$ -counterpart  $(x, w)$  must be in  $\tilde{Q}_v$  for some child  $v$  of  $u$ . Since  $b \in A^o$ , the  $\mathcal{E}$ -counterparts of the strong edges in  $\tilde{Q}_b$  are in more than one tree in  $\mathcal{R}_u$ , and thus there exists a tree  $\tilde{R}_{v'}$  different from  $\tilde{R}_v$  which contains at least one  $\mathcal{E}$ -counterpart  $(y, z)$  of a strong edge  $(c, d)$  in  $\tilde{Q}_b$ . Clearly  $(a, b)$  and  $(c, d)$  are comparable, contradicting Lemma 4.1 as  $(x, w)$  and  $(y, z)$  are incomparable. Therefore, all the edges connecting  $a$  with vertices in  $A^o$  are weak.

For  $b \in A^o$ , we let  $\tilde{R}_{v_1}, \dots, \tilde{R}_{v_r}$  be the trees that contain  $\mathcal{E}$ -counterparts of strong edges in  $\tilde{Q}_b$ .

We claim that the  $\mathcal{E}$ -counterparts of the strong edges in the  $\tilde{R}_{v_i}$ 's are all in  $\tilde{Q}_b$ . Suppose to the contrary that there exists a tree  $\tilde{R}_{v_i}$ , say  $\tilde{R}_{v_1}$ , containing a strong edge  $e'_1$  with its  $\mathcal{E}$ -counterpart  $e_1$  in  $\tilde{Q}_{b'}$  for some child  $b' \neq b$  of  $a$ . By definition,  $\tilde{R}_{v_1}$  contains a strong edge  $e'_2$  different from  $e'_1$  that is the counterpart of a strong edge in  $\tilde{Q}_b$ . In addition, also by definition,  $\tilde{Q}_b$  contains at least one strong edge  $e_3$  different from  $e_2$  whose  $\mathcal{E}$ -counterpart  $e'_3$  is in a tree  $\tilde{R}_{v_i}$  different from  $\tilde{R}_{v_1}$ . By Lemma 4.1,  $e_2$  and  $e_3$  ( $e'_1$  and  $e'_2$ ) are incomparable as their  $\mathcal{E}$ -counterparts  $e'_2$  and  $e'_3$  ( $e_1$  and  $e_2$ ) are incomparable. Moreover, the close pair of the first triple is  $e_2$  and  $e_3$  and the close pair of the second triple is  $e'_2$  and  $e'_1$ , violating Lemma 4.2. We can conclude that  $\mathcal{E}$ -counterparts of the strong edges in the  $\tilde{R}_{v_i}$ 's are all in  $\tilde{Q}_b$ .

We can now define  $X^m$  so that, for any  $b \in A^o$ ,  $X^m$  contains the children  $v_1, \dots, v_r$  of  $u$  such that  $\tilde{R}_{v_1}, \dots, \tilde{R}_{v_r}$  are the trees that contain  $\mathcal{E}$ -counterparts of strong edges in  $\tilde{Q}_b$ . The surjection  $\alpha$  maps any vertex  $v_i$  in  $X^m$  to the corresponding vertex  $b$  of  $A^o$ . Clearly,  $X^m$  and  $X^t$  are disjoint as for any  $v_i \in X^m$  the  $\mathcal{E}$ -counterparts of the strong edges of  $\tilde{R}_{v_i}$  belong to trees  $\tilde{Q}_b$  for children  $b$  of  $a$  such that  $(a, b)$  is weak. This completes the proofs of conditions 3 and 4 as far as sets  $X^m$  and  $A^o$  are concerned.

Working symmetrically, we can include in  $X^o$  all the children  $v$  of  $u$  such that the strong edges of  $\tilde{R}_v$  have  $\mathcal{E}$ -counterparts in more than one tree in  $\tilde{Q}_b$  for children  $b$  of  $a$ . As before,  $X^o$  and  $X^t$  are disjoint. Moreover,  $X^o$  and  $X^m$  are also disjoint as, according to the discussion above, for any  $v_i \in X_m$  the strong edges of  $\tilde{Q}_{v_i}$  are all in a single  $\tilde{Q}_b$ . Applying the same arguments as before, we can define the set  $A^m$  and surjection  $\chi : A^m \rightarrow X^o$  and verify that conditions 3 and 4 are satisfied for  $X^o$ ,  $A^m$ , and  $\chi$ .

The construction of the desired tuple is not yet complete. If  $b$  is a child of  $a$  that has not yet been classified as a member of  $A^e$ ,  $A^t$ ,  $A^o$ , or  $A^m$ , then the  $\mathcal{E}$ -counterparts of the strong edges of  $\tilde{Q}_b$  are *all* in exactly one tree  $\tilde{R}_v$  but  $(a, b)$  and  $(u, v)$  are not both strong edges.

We can make a similar claim for unclassified children  $v$  of  $u$ . Therefore, there is a bijection  $\sigma$  between the unclassified children of  $a$  and the unclassified children of  $u$  that allows us to classify each one of them arbitrarily in  $A^o$  and  $X^m$  respectively or in  $A^m$  and  $X^o$  respectively. For each such arbitrary choice  $\alpha$  or  $\chi$  is augmented by  $\sigma$  on the new pair of elements. By repeating the same arguments one can prove that, after this enhancement, the sets defined still satisfy properties 3 and 4 while  $\alpha$  and  $\chi$  remain surjections. In conclusion, the tuple  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$  satisfies properties 1-5 and the lemma holds. ■

In order to define the recurrence for our dynamic programming algorithm, we need to be able to decompose a minimum size expansion of two trees into minimum size expansions of pairs of subtrees. The following two lemmas do this for the subtrees needed when considering decompositions induced by removal of a strong edge or a weak edge, respectively.

**Lemma 4.4.** *If  $\mathcal{E}$  is a minimum size expansion of  $Q = Q_a$  and  $R = R_u$ ,  $(a, b)$  is a strong edge of  $Q$ , and  $(u, v) = f_\epsilon((a, b))$  is its  $\mathcal{E}$ -counterpart, then  $\mathcal{E}$  is the union of a minimum size expansion of  $Q \setminus Q_b$  and  $R \setminus R_v$  containing  $(a, u)$  and a minimum size expansion of  $Q_b$  and  $R_v$  containing  $(b, v)$ .*

**Proof.** We first claim that  $\mathcal{E}_1 = \mathcal{E}[V(Q_b) \cup V(R_v)]$  is an expansion of  $Q_b$  and  $R_v$  containing  $(b, v)$ . In order to prove this, it is enough to show that all the neighbors in  $\mathcal{E}$  of all the vertices in  $V(Q_b)$  ( $V(R_v)$ ) are in  $V(R_v)$  ( $V(Q_b)$ ). Suppose to the contrary that there exists an edge  $(c, w)$  in  $\mathcal{E}$  where  $c \in V(Q_b)$  and  $w \notin V(R_v)$ . If we now apply Lemma 3.5 for  $P_Q(a, c)$  and  $P_R(u, w)$ ,  $b \in P_Q(a, c)$  must be adjacent to some vertex not in  $R_v$ . As  $b$  is adjacent to  $v$  in  $\mathcal{E}$ , by property 1 of the definition of  $\mathcal{E}$ ,  $b$  must be adjacent to  $u$ , a contradiction as  $(a, b)$  and  $(u, v)$  are strong edges. By symmetry we can prove that  $\mathcal{E}_2 = \mathcal{E}[V(Q \setminus Q_b) \cup V(R \setminus R_v)]$  is an expansion of  $Q \setminus Q_b$  and  $R \setminus R_v$  containing  $(a, u)$ .

It now remains to prove that  $\mathcal{E}_1$  and  $\mathcal{E}_2$  are both minimum size expansions. Suppose instead that there is an expansion  $\mathcal{E}'$  of one of the pairs  $Q_b, R_v$  and  $Q \setminus Q_b, R \setminus R_v$ ,  $Q_b$  and  $R_v$ , that has size smaller than the one of  $\mathcal{E}_1$ . Then, by Lemma 3.3,  $\mathcal{E}' \cup \mathcal{E}_2$  is an expansion of  $Q$  and  $R$  with size smaller than  $\mathcal{E}$ , contradicting the minimality of  $\mathcal{E}$ . ■

For notational convenience, we will use short forms for various subgraphs of  $Q_a$  and  $R_u$ . For  $b$  a child of  $a$  in  $Q_a$ , we define  $Q_{-b}$  to be  $Q_a \setminus Q_b$ . For any subset  $X$  of children of  $u$  in  $R_u$ , we define  $R_{-X}$  to be  $R[(V(R) \setminus V(R_X)) \cup \{u\}]$ .

**Lemma 4.5.** *If  $\mathcal{E}$  is a minimum size expansion of  $Q = Q_a$  and  $R = R_u$ ,  $(a, b)$  is a weak edge of  $Q$ , and  $(a, u)$  and  $(b, u)$  are edges of  $\mathcal{E}$ , then*

1. *there exists a subset  $X$  of the children of  $u$  such that  $\mathcal{E}$  is the union of a minimum size expansion of  $Q_b$  and  $R_X$  containing  $(b, u)$  and a minimum size expansion of  $Q_{-b}$  and  $R_{-X}$  containing  $(a, u)$ , and*
2. *if  $Q_b$  contains only weak edges, then for any vertex  $c \in V(Q_b)$ ,  $N_{\mathcal{E}}(c) = \{u\}$ .*

**Proof.** We let  $X$  be the set containing any child  $v$  of  $u$  for which  $R_v$  contains a neighbor, in  $\mathcal{E}$ , of a vertex in  $Q_b$ . We then let  $\mathcal{E}_1 = \mathcal{E}[V(Q_b) \cup V(R_X)]$  and  $\mathcal{E}_2 = \mathcal{E}[V(Q_{-b}) \cup V(R_{-X})]$ . To show that  $\mathcal{E}_1$  is an expansion of  $Q_b$  and  $R_X$  and that  $\mathcal{E}_2$  is an expansion of  $Q_{-b}$  and  $R_{-X}$  by inheriting the properties of an expansion from  $\mathcal{E}$ , it will suffice to show that in  $\mathcal{E}$  all neighbors of vertices in  $V(Q_b)$  ( $V(R_X)$ ,  $V(Q_{-b})$ , and  $V(R_{-X})$ , respectively) are in  $V(R_X)$  ( $V(Q_b)$ ,  $V(R_{-X})$ , and  $V(Q_{-b})$ , respectively). The first of the four statements follows from the definition of  $X$ .

To prove the third claim by contradiction, suppose instead that a vertex  $c \in Q_{-b}$  is adjacent in  $\mathcal{E}$  to a vertex  $w$  outside of  $R_{-X}$ . Clearly,  $X \neq \emptyset$  and  $w$  is in one of the connected components of  $R_X - \{u\}$ . Let  $v$  be the vertex of  $X$  such that  $w \in R_v$ . In addition, by the definition of  $X$ ,  $Q_b$  contains at least one vertex  $d$  adjacent, in  $\mathcal{E}$ , to a vertex  $x$  in  $R_v$ . We now apply Lemma 3.5 to paths  $P_Q(d, c)$  and  $P_R(w, x)$  to conclude that  $a \in P_Q(d, c)$  is adjacent, in  $\mathcal{E}$ , to some vertex in  $R_v$ . Since  $a$  is also adjacent to  $u$  in  $\mathcal{E}$ , by property 1 of the definition of  $\mathcal{E}$ ,  $a$  must be adjacent to  $v$  in  $\mathcal{E}$ . Similarly, we can show that  $b$  is adjacent to  $v$ . Therefore, the neighborhoods of  $a$  and  $b$  have two vertices, i.e.  $u$  and  $v$ , in common. This contradicts property 3 of the definition of  $\mathcal{E}$  and hence the claim holds. The remaining claims can be proved in a similar manner.

To prove the first statement in the lemma, it now remains to show that  $\mathcal{E}_1$  and  $\mathcal{E}_2$  are both minimum size expansions. Suppose to the contrary that there is an expansion  $\mathcal{E}'$  of one of the pairs  $Q_b$  and  $R_X$  or  $Q_{-b}$  and  $R_{-X}$  that has size smaller than the one established above, say  $Q_b$  and  $R_X$  have an expansion  $\mathcal{E}'$  smaller than  $\mathcal{E}_1$ . Then, by Lemma 3.2,  $\mathcal{E}' \cup \mathcal{E}_2$  is an expansion of  $Q$  and  $R$  with size smaller than  $\mathcal{E}$ , contradicting the minimality of  $\mathcal{E}$ .

To prove the second statement in the lemma, it suffices to show that if  $Q_b$  contains only weak edges, then  $X = \emptyset$ . Suppose instead that  $|X| \geq 1$ . Then  $|E(R_X)| \geq 1$ . Since  $\mathcal{E}_2$  is a minimum size expansion of  $Q_b$  and  $R_X$ , and  $Q_b$  contains only weak edges,  $R_X$  contains only weak edges. But then  $|E(T_{\mathcal{E}_2})| = |E(Q_b)| + |E(R_X)|$ , contradicting Lemma 3.9. ■

The following lemma uses the structural information of Lemma 4.3, followed by repeated applications of Lemmas 4.4 and 4.5.

**Lemma 4.6.** *For any minimum size expansion  $\mathcal{E}$  of two trees  $Q = Q_a$  and  $R = R_u$  such that  $(a, u) \in \mathcal{E}$ , there exists a tuple  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$  in  $\Pi(\text{children}(a), \text{children}(u))$  such that  $(\mathcal{E}_e \cup \mathcal{E}_t \cup \mathcal{E}_a \cup \mathcal{E}_u)$  is a partition of  $\mathcal{E}$  where*



1.  $\mathcal{E}_e$  relates epsilon children;  $\mathcal{E}_e = \{(a, z) \mid z \in V(Q_{A^e})\} \cup \{(c, u) \mid c \in V(R_{X^e})\}$ .
2.  $\mathcal{E}_t$  relates terminal children;  $\mathcal{E}_t = \bigcup_{b \in A^t} \mathcal{E}_{t,b}$  where, for any vertex  $b \in A^t$ ,  $\mathcal{E}_{t,b}$  is a minimum expansion of  $Q_b$  and  $R_{\tau(b)}$  that contains  $(b, \tau(b))$ .
3.  $\mathcal{E}_a$  relates one-many children;  $\mathcal{E}_a = \bigcup_{b \in A^o} \mathcal{E}_{a,b}$  where, for any vertex  $b \in A^o$ ,  $\mathcal{E}_{a,b}$  is a minimum expansion of  $Q_b$  and  $R_{\alpha^{-1}(b)}$  that contains  $(b, u)$ .
4.  $\mathcal{E}_u$  relates many-one children;  $\mathcal{E}_u = \bigcup_{v \in X^o} \mathcal{E}_{u,v}$  where, for any vertex  $v \in X^o$ ,  $\mathcal{E}_{u,v}$  is a minimum expansion of  $Q_{\chi^{-1}(v)}$  and  $R_v$  that contains  $(a, v)$ .

**Proof.** We will prove the lemma by decomposing  $\mathcal{E}$  in groups of subexpansions of the four types described in Lemma 4.3. This decomposition will proceed step by step by applying inductively, Lemmas 4.4 and 4.5 as appropriate, depending on the type of subexpansion it is possible to extract.

Let  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi) \in \Pi(\text{children}(a), \text{children}(u))$  be as determined by Lemma 4.3. We will extract the decomposition of  $\mathcal{E}$  using induction on  $j = |A^e| + |X^e| + |A^t| + |A^o| + |X^o|$ . If  $j = 0$ , the result is trivial. We assume that it holds if  $j < k$  and we will prove that it also holds when  $j = k \geq 1$ . Let  $b \in A^e \cup X^e \cup A^t \cup A^o \cup X^o$ . We may assume that  $b$  is a vertex in  $A^e \cup A^t \cup A^o$ , as the case where  $b$  is a vertex in  $X^e \cup X^o$  is symmetric. We set  $\mathcal{E}_{\neg b} = \mathcal{E}[(V(Q \setminus Q_b) \cup V(R \setminus R_{\sigma(b)}))]$  where, if  $b \in A^e$ , (resp.  $b \in A^t$ ,  $b \in A^o$ ), then  $\sigma(b) = \emptyset$ , (resp.  $\sigma(b) = \tau(b)$ ,  $\sigma(b) = \alpha^{-1}(b)$ ).

We now claim that  $\mathcal{E}_{\neg b}$  is a minimum size expansion of  $Q \setminus Q_b$  and  $R \setminus R_{\sigma(b)}$  and that  $\mathcal{E}_b = \mathcal{E}[V(Q_b) \cup V(R_{\sigma(b)})]$  is a minimum size expansion of  $Q_b$  and  $R_{\sigma(b)}$ . When  $b \in A^t$ , the claim is a consequence of Lemma 4.4 and when  $b \in A^e \cup A^o$ , the claim is a consequence of Lemma 4.5.

We can now apply the induction hypothesis on  $\mathcal{E}_{\neg b}$  and derive the tuple  $(\{A_{\neg b}^e, A_{\neg b}^t, A_{\neg b}^o, A_{\neg b}^m\}, \{X_{\neg b}^e, X_{\neg b}^t, X_{\neg b}^o, X_{\neg b}^m\}, \tau_{\neg b}, \alpha_{\neg b}, \chi_{\neg b}) \in P(\text{children}(a) - \{b\}, \text{children}(u) - \sigma(b))$  and the corresponding partition  $(\mathcal{E}_{e,\neg b}, \mathcal{E}_{t,\neg b}, \mathcal{E}_{a,\neg b}, \mathcal{E}_{u,\neg b})$  of  $\mathcal{E}_{\neg b}$  satisfying conditions 1–4. If  $b \in A^t$ , and  $v$  is as defined in Lemma 4.4, for each member  $m$  of the tuple,  $m = m_{\neg b}$  with the following exceptions:  $A^t = A_{\neg b}^t \cup \{b\}$ ,  $X^t = X_{\neg b}^t \cup \{v\}$ , and  $\tau = \tau_{\neg b} \cup \{(b, v)\}$ . Suppose now that  $b \in A^e \cup A^o$  and  $X$  is as defined in Lemma 4.5. In this case, it is easy to see that if  $X = \emptyset$ , then  $A^e = A_{\neg b}^e \cup \{b\}$ , and for each other member  $m$  of the tuple  $m = m_{\neg b}$ . Finally, if  $X \neq \emptyset$ ,  $A^o = A_{\neg b}^o \cup \{b\}$ ,  $X^m = X_{\neg b}^m \cup X$ , and  $\alpha = \alpha_{\neg b} \cup \{(w, b) \mid w \in X\}$ , with all other members of the tuple unchanged. We construct the partition  $(\mathcal{E}_e \cup \mathcal{E}_t \cup \mathcal{E}_a \cup \mathcal{E}_u)$  of  $\mathcal{E}$  as follows.

If  $b \in A^e$ , then, by Lemma 4.5,  $\mathcal{E}_b = \{(c, u) \mid c \in V(Q_b)\}$ . We set  $\mathcal{E}_e = \mathcal{E}_b \cup \mathcal{E}_{e,\neg b}$ ,  $\mathcal{E}_t = \mathcal{E}_{t,\neg b}$ ,  $\mathcal{E}_a = \mathcal{E}_{a,\neg b}$ , and  $\mathcal{E}_u = \mathcal{E}_{u,\neg b}$ .

If  $b \in A^o$ , then, by Lemma 4.5,  $\mathcal{E}_b$  is a minimum expansion of  $Q_b$  and  $R_{\alpha^{-1}(b)} = R_{\sigma(b)}$ . We set  $\mathcal{E}_e = \mathcal{E}_{e,-b}$ ,  $\mathcal{E}_t = \mathcal{E}_{t,-b}$ ,  $\mathcal{E}_a = \mathcal{E}_b \cup \mathcal{E}_{a,-b}$ , and  $\mathcal{E}_u = \mathcal{E}_{u,-b}$ .

If  $b \in A^t$ , then, by Lemma 4.4,  $\mathcal{E}_b$  is a minimum expansion of  $Q_b$  and  $R_{\tau(b)} = R_{\sigma(b)}$ . We set  $\mathcal{E}_e = \mathcal{E}_{e,-b}$ ,  $\mathcal{E}_t = \mathcal{E}_b \cup \mathcal{E}_{t,-b}$ ,  $\mathcal{E}_a = \mathcal{E}_{a,-b}$ , and  $\mathcal{E}_u = \mathcal{E}_{u,-b}$ .

It now remains to verify that, in any case, the tuple

$$(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$$

along with the partition  $(\mathcal{E}_e \cup \mathcal{E}_t \cup \mathcal{E}_a \cup \mathcal{E}_u)$  satisfy conditions 1–4. This check is straightforward except for conditions 2 and 3 where we have to prove that the new expansions  $\mathcal{E}_t$  and  $\mathcal{E}_a$  are minimum. This follows from Lemmas 3.3 and 3.2 as, by their construction, they are the union of minimum expansions. ■

### 4.3 Algorithm details

**Procedure** Expansion( $Q, R, a, u$ )

*Input:* Two trees  $Q, R$  and two vertices  $a \in V(Q), u \in V(R)$ .

*Output:*  $\min\{|\mathcal{E}| : \mathcal{E} \text{ is an expansion of } Q \text{ and } R \text{ and } (a, u) \in \mathcal{E}\}$ .

```

1: Root  $Q$  and  $R$  at  $a$  and  $u$  respectively.
2: Topologically sort  $V(Q)$ , giving  $L_Q := \{a_1, \dots, a_{|V(Q)|}\}$  where  $a = a_{|V(Q)|}$ .
3: Topologically sort  $V(R)$ , giving  $L_R := \{u_1, \dots, u_{|V(R)|}\}$  where  $u = u_{|V(R)|}$ .
4: for  $i := 1 \dots |V(Q)|$  do
5:   for  $j := 1 \dots |V(R)|$  do
6:     if  $a_i$  and  $u_j$  are leaves then  $I(a_i, u_j, \emptyset, \emptyset) := 1$ 
7:     else
8:       for all  $X \subseteq \text{children}(u_j)$  and  $A \subseteq \text{children}(a_i)$  do
9:          $x := |V(Q)| + |V(R)|$ 
10:        for all  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi) \in \Pi(A, X)$  do
11:           $x := \min\{ x, |V(Q_{A^e})| + |V(R_{X^e})| - 1 +$  (i)
               $\sum_{b \in A^t} I(b, f_t(b), \text{children}(b), \text{children}(\tau(b))) +$  (ii)
               $\sum_{b \in A^o} I(b, u_i, \text{children}(b), \alpha^{-1}(b)) +$  (iii)
               $\sum_{v \in X^o} I(a_i, v, \chi^{-1}(v), \text{children}(v)) \}$  (iv)
12:           $I(a_i, u_j, A, X) := x$ 
13: return  $I(a, u, \text{children}(a), \text{children}(u))$ 

```

**Theorem 4.7.** *For any trees  $Q$  and  $R$  rooted at  $a$  and  $u$  respectively, Expansion  $(Q, R, a, u)$  returns the minimum number of edges in any expansion  $\mathcal{E}$  containing  $(a, u)$ .*

**Proof:** We prove that for  $Q$  and  $R$  rooted at  $a$  and  $u$  respectively, for any  $c \in V(Q)$ ,  $z \in V(R)$ , and any  $A \subseteq \text{children}(c)$  and  $X \subseteq \text{children}(z)$ , the quantity  $I(c, z, A, X)$  computed by the algorithm is the minimum number of edges over all expansions  $\mathcal{E}$  of  $Q_A$  and  $R_X$ , where  $(c, z) \in \mathcal{E}$ . The proof is by induction on the order of computation.

Consider the computation of  $I(c, z, A, X)$ . As  $L_Q$  and  $L_R$  are topological sorts of  $V(Q)$  and  $V(R)$  respectively, we can conclude that  $I(d, y, A_d, X_y)$  has already been computed in the following three cases, which cover the expressions on the right-hand side of step 11.

1.  $d \in \text{children}(c)$ ,  $y \in \text{children}(z)$ ,  $A_d \subseteq \text{children}(d)$ , and  $X_y \subseteq \text{children}(y)$ .
2.  $d \in \text{children}(c)$ ,  $y = z$ ,  $A_d \subseteq \text{children}(d)$ , and  $X_y \subseteq \text{children}(z)$ .
3.  $d = c$ ,  $y \in \text{children}(z)$ ,  $A_d \subseteq \text{children}(c)$ , and  $X_y \subseteq \text{children}(y)$ .

If we assume by the inductive hypothesis that the values  $I(d, y, A_d, X_y)$  are correct, then by Lemma 4.6 there is a choice of  $(\{A^e, A^t, A^o, A^m\}, \{X^e, X^t, X^o, X^m\}, \tau, \alpha, \chi)$  that results, at step 11, in  $x$  taking on the minimum number of edges in an expansion  $\mathcal{E}$  of  $Q_A$  and  $R_X$  containing  $(c, z)$ , as required. ■

**Theorem 4.8.** *For any pair of trees  $Q$  and  $R$  of bounded degree,  $\text{sctmj}(Q, R)$  can be computed in  $O(n^3)$  time where  $n = \max\{|V(Q)|, |V(R)|\}$ .*

**Proof:** The if-statement at step 6 is invoked  $O(n^2)$  times, and because the maximum degrees of  $Q$  and  $R$  are bounded by a constant, the loops at steps 8 and 10 result in a constant number of iterations of step 11. This quantity is multiplied by  $O(n)$ , the number of rootings to check. ■

## 5 Extensions of the algorithm

In this section we describe how our algorithm can be generalized to the problem of determining the edit distance (under certain conditions) of a pair of a edge-labeled, unrooted, unordered trees. The set of operations used is edge contraction, edge relabeling, and edge insertion. In this last operation, a vertex  $v$  is chosen, replaced by a pair of vertices  $v_1$  and  $v_2$  such that  $N(v_1)$  and  $N(v_2)$  partition  $N(v)$ , and a labeled edge  $(v_1, v_2)$  is inserted. The final

condition we impose on the edit sequence is that all insertions must be completed before any other operations.

Let  $Q$  and  $R$  be edge-labeled trees, i.e. for some given alphabet  $\Sigma$ , there exist two functions  $q : E(Q) \rightarrow \Sigma$  and  $r : E(R) \rightarrow \Sigma$ . We denote as  $g, h$  the cost functions where for any  $\sigma \in \Sigma$ ,  $g(\sigma)$  and  $h(\sigma)$  represent the cost of the contraction and the insertion respectively of an edge labeled with  $\sigma$ . Finally, for  $\sigma, \rho \in \Sigma$ , we denote as  $l(\sigma, \rho)$  the cost of changing the label of an edge from  $\sigma$  to  $\rho$ . We can now define as  $\text{dist}(Q, R)$  the smallest possible total cost of a sequence of operations which transforms  $Q$  to  $R$ , subject to the constraint that all insertions occur first.

Given such a sequence, we can reorder it (without altering its cost) so that the relabelings precede the contractions and follow the expansions. Let  $T_1$  be the tree after all expansions, and  $T_2$  the tree after all relabelings. Clearly, if labels are removed,  $T_1$  is isomorphic to  $T_2$ , and both are majors of both  $Q$  and  $R$ . Thus, for every edit sequence, there is a natural common supertree.

Conversely, let  $T$  be a common major of  $Q$  and  $R$  corresponding to some extension  $\mathcal{E}$  of  $Q$  and  $R$ . It is easy to see that  $Q$  can be transformed to  $R$  after the following sequence of operations: first insert in  $Q$  all the edges in  $E(T) - E(Q)$ , then relabel all the strong edges of  $T$  to the labelings they should have in  $R$ , and, finally, contract all the edges in  $E(T) - E(R)$ . Notice that, if  $S(T)$  contains the strong edges of  $T$ , the total cost of this sequence of operations is

$$\sum_{e \in E(T) - E(Q)} h(q(e)) + \sum_{e \in E(T) - E(R)} g(r(e)) + \sum_{e \in S(T)} l(q(e), r(e))$$

which, in turn, is equal to

$$\begin{aligned} & \sum_{e \in E(T)} h(q(e)) + \sum_{e \in E(T)} g(r(e)) + \sum_{e \in S(T)} l(q(e), r(e)) - C(R, Q) = \\ & \sum_{e \in E(T)} (h(q(e)) + g(r(e))) + \sum_{e \in S(T)} l(q(e), r(e)) - C(R, Q) \end{aligned}$$

where  $C(R, Q) = \sum_{e \in E(Q)} h(q(e)) + \sum_{e \in E(R)} g(r(e))$ . Therefore, in order to compute  $\text{dist}(Q, R)$  we have to find an expansion  $\mathcal{E}$  with major  $T$  where the quantity

$$Q(T) = \sum_{e \in E(T)} (h(q(e)) + g(r(e))) + \sum_{e \in S(T)} l(q(e), r(e))$$

is minimized. Following the methodology of the previous sections we set up a general version of  $I(c, z, A, X)$ , representing the minimum value of  $Q(T)$  over the  $T$ 's corresponding to all expansions  $\mathcal{E}$  of  $Q_A$  and  $Q_X$  where  $(c, z) \in \mathcal{E}$ . The only modification required for Procedure

Expansion( $Q, R, a, u$ ) concerns the way  $x$  is computed in line 11, which should change to the following:

$$\begin{aligned}
11: \quad x := \min\{ x, & \sum_{e \in E(Q_{A^e})} h(q(e)) + \sum_{e \in E(R_{X^e})} g(r(e)) + & \text{(i)} \\
& \sum_{b \in A^t} ( I(b, f_t(b), \text{children}(b), \text{children}(\tau(b))) + & \\
& \quad l(q(\{a_i, b\}), r(\{u_i, f_t(b)\})) ) + & \text{(ii)} \\
& \sum_{b \in A^o} ( I(b, u_i, \text{children}(b), \alpha^{-1}(b)) + h(q(\{a_i, b\})) ) + & \text{(iii)} \\
& \sum_{v \in X^o} ( I(a_i, v, \chi^{-1}(v), \text{children}(v)) + g(r(\{u_j, v\})) ) \} & \text{(iv)}
\end{aligned}$$

For completeness, in line 9,  $x$  should now be initialized as  $C(R, Q)$ .

Clearly, the above modifications do not require more time asymptotically, and we have the following:

**Theorem 5.1.** *The edit distance (under operations edge contraction, edge relabeling, and edge insertion, where all insertions come first) of any pair of edge-labeled trees  $Q$  and  $R$  of bounded degree, can be computed in  $O(n^3)$  time where  $n = \max\{|V(Q)|, |V(R)|\}$ .*

## 6 Conclusions and further work

We have shown an  $O(n^3)$  algorithm for finding the smallest common tree major of two trees  $Q$  and  $R$ , where both  $Q$  and  $R$  are unrooted and undirected, and have degree bounded by a fixed constant. The degree restriction can be relaxed to maximum degree  $O(\log n / \log \log n)$  while keeping the running time of the algorithm polynomial, since the multiplicative factor is  $d^{O(d)}$  for trees of maximum degree  $d$  (this factor arises from the number of tuples examined at line 10 of the algorithm). Our algorithm can be generalized to the problem of determining the edit distance (under the operations of edge contraction, edge relabeling, and edge insertion, where all insertions come first) of a pair of a edge-labeled, unrooted, unordered trees, by incorporating labels into the definition of the expansion. All of our algorithms can be implemented in NC using the technique of Brent restructuring to parallelize dynamic programming on trees [9]. Our work is also related to work on intertwines [15]: the value  $\text{sctmj}(Q, R)$  is the minimum size of an acyclic intertwine of  $Q$  and  $R$ .

Although the NP-completeness of minor containment for general trees suggests the intractability of finding the largest common subgraph under minors, there is hope for solving other related problems. The problem of determining whether or not  $G$  is a minor of  $H$  is solvable in polynomial time for  $G$  and  $H$  both bounded-degree partial  $k$ -trees [16] or for  $G$  and  $H$  both  $k$ -connected  $k$ -paths [11]; solving the largest common supergraph problem for

each of these graph classes would be an obvious extension to our work. Another obvious extension would be to solve the largest common tree major problem for three or more input trees.

## References

- [1] A. Amir and D. Keselman. Maximum agreement subtree in a set of evolutionary trees: metrics and efficient algorithms. *SIAM Journal on Computing*, 26(6):1656–1669, December 1997.
- [2] J. A. Bondy and U.S.R. Murty. *Graph Theory with Applications*. North-Holland, 1976.
- [3] M. J. Chung.  $O(n^{2.5})$  time algorithms for the subgraph homeomorphism problem on trees. *Journal of Algorithms*, 8:106–112, 1987.
- [4] Richard Cole and Ramesh Hariharan. An  $O(n \log n)$  algorithm for the maximum agreement subtree problem for binary trees. In *Proceedings of the Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 323–332, 1996.
- [5] M. Dubiner, Z. Galil, and E. Magen. Faster tree pattern matching. In *Proceedings of the 31st Annual Symposium on Foundations of Computer Science*, pages 145–150, 1990.
- [6] P. Duchet. Tree minors. Presentation at *AMS-IMS-SIAM Joint Summer Research Conference on Graph Minors*, 1991 (personal communication, A. Gupta).
- [7] M. Farach, T. Przytycka, and M. Thorup. On the agreement of many trees. *Information Processing Letters*, 55(6):297–301, 1995.
- [8] M. Farach and M. Thorup. Fast comparison of evolutionary trees. In *Proceedings of the Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 481–488, 1994.
- [9] A. Gupta and N. Nishimura. The parallel complexity of tree embedding problems. *Journal of Algorithms*, 18(1):176–200, 1995.
- [10] A. Gupta and N. Nishimura. Finding largest subtrees and smallest supertrees. *Algorithmica*, 21:183–210, 1998.
- [11] A. Gupta, N. Nishimura, A. Proskurowski, and P. Ragde. Embeddings of  $k$ -connected graphs of pathwidth  $k$ . Manuscript.

- [12] T. Jiang, L. Wang, and K. Zhang. Alignment of trees – an alternative to tree edit. In *Combinatorial Pattern Matching*, pages 75–86, 1994.
- [13] P. Kilpeläinen and H. Mannila. Ordered and unordered tree inclusion. *SIAM Journal on Computing*, 24(2):340–356, 1995.
- [14] S. R. Kosaraju. Efficient tree pattern matching. In *Proceedings of the 30th Annual Symposium on Foundations of Computer Science*, pages 178–183, 1989.
- [15] J. Lagergren. The size of an intertwiner. In *Proceedings of the 23rd International Colloquium on Automata, Languages, and Programming*, volume 820 of *Lecture Notes in Computer Science*, pages 520–531, 1994.
- [16] J. Matoušek and R. Thomas. On the complexity of finding iso- and other morphisms for partial  $k$ -trees. *Discrete Mathematics*, 108:343–364, 1992.
- [17] N. Robertson and P. Seymour. Graph minors II. Algorithmic aspects of tree-width. *Journal of Algorithms*, 7:309–322, 1986.
- [18] K. Siddiqi, A. Shokoufandeh, S. Dickinson, and S. Zucker. Shock graphs and shape matching. *International Journal of Computer Vision*, to appear.
- [19] T. Warnow. Tree compatibility and inferring evolutionary history. In *Proceedings of the Fourth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 382–391, 1993.