

Computing the Quartet Distance Between Evolutionary Trees in Time $O(n \log^2 n)$

Gerth Stølting Brodal^{1,*}, Rolf Fagerberg^{1,*}, and Christian N. S. Pedersen^{1,*}

BRICS[†], Department of Computer Science, University of Aarhus, Ny Munkegade, DK-8000 Århus C, Denmark. E-mail: {gerth,rolf,cstorm}@brics.dk

Abstract Evolutionary trees describing the relationship for a set of species are central in evolutionary biology, and quantifying differences between evolutionary trees is an important task. One previously proposed measure for this is the quartet distance. The quartet distance between two unrooted evolutionary trees is the number of quartet topology differences between the two trees, where a quartet topology is the topological subtree induced by four species. In this paper, we present an algorithm for computing the quartet distance between two unrooted evolutionary trees of n species in time $O(n \log^2 n)$. The previous best algorithm runs in time $O(n^2)$.

1 Introduction

The evolutionary relationship for a set of species is commonly described by an evolutionary tree. This is a rooted tree where the leaves correspond to the species, and the internal nodes correspond to speciation events, i.e. the points in time where the evolution has diverged in different directions. The direction of the evolution is described by the location of the root, which corresponds to the most recent common ancestor for all the species, and the rate of evolution is described by assigning lengths to the edges. The true evolutionary tree for a set of species is rarely known, hence estimating it from obtainable information about the species, e.g. genomic data, is of great interest. The problem of computationally estimating aspects of the true evolutionary tree requires a model describing how to use the available information about the species in question. Given a model, the problem of estimating certain aspects of the true evolutionary tree is often referred to as constructing the evolutionary tree in that model. Many models and methods for constructing evolutionary trees have been presented, see [10, Chap. 17] for an overview.

An important aspect of the true evolutionary tree is the undirected tree topology induced by ignoring the location of root and the length of the edges. Many models and methods are concerned with estimating this tree topology,

* Partially supported by the IST Programme of the EU under contract number IST-1999-14186 (ALCOM-FT).

† Basic Research in Computer Science, www.brics.dk, funded by the Danish National Research Foundation.

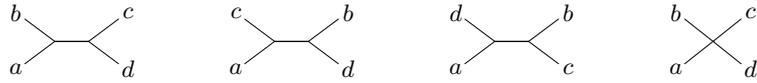


Figure 1. The four possible quartet topologies of species $a, b, c,$ and d

usually under the further assumption that all internal nodes have degree three. We say that such models and methods are concerned with constructing the unrooted evolutionary tree of degree three for a set of species. For the remainder of this paper an evolutionary tree denotes an unrooted evolutionary tree of degree three.

Different models and methods often yield different estimates of the evolutionary tree for the same set of species. The same model and method can also give rise to different evolutionary trees for the same set of species when applied to different information about the species, e.g. different genes. To study such differences in a systematic manner, one must be able to quantify differences between evolutionary trees using well-defined and efficient methods.

One approach for comparing two evolutionary trees is to determine a consensus tree (or forest) that reflects common traits of the two trees, e.g. the maximum agreement subtree. Much work has been concerned with developing efficient methods for computing the maximum agreement subtree of two or more evolutionary trees, see e.g. [2]. Another approach for comparing two evolutionary trees is to define a distance measure between two trees and compare the two trees by computing the distance. Several distance measures have been proposed, e.g. the symmetric difference metric [12], the nearest-neighbor interchange metric [16], the subtree transfer distance [1], the Robinson and Foulds metric [13], and the quartet metric [8]. Each distance measure has different properties and reflects different aspects of biology, e.g. the subtree transfer distance is related to the number of recombinations between the two sets of species. The quartet metric has several attractive properties. Bryant *et al.* in [5] discuss the properties of the quartet metric and conclude that it does not suffer from drawbacks of the other distance measures. For example, measures based on transformation operations, e.g. the subtree transfer distance, do not distinguish between transformations that affect a large number of leaves and transformations that affect a small number of leaves.

In this paper, we study the quartet metric. For an evolutionary tree, the *quartet topology* of four species is the topological subtree induced by these species. In general, the possible quartet topologies of four species are the four shown in Fig. 1. Of these, the right-most cannot occur if we assume that all internal nodes have degree three. It is well-known that the complete set of quartet topologies is unique for a given tree and that the tree can be uniquely recovered from its set of quartet topologies in polynomial time [6]. If the tree has degree three, then, as observed in [11], it can be recovered from its set of quartet topologies in time $O(n \log n)$ using methods [4,9,11] for constructing an evolutionary tree in the experiment model in time $O(n \log n)$.

Given two evolutionary trees on the same set of n species, the *quartet distance* between them is the number of sets of four species for which the quartet topologies differ in the two trees. Since there are $\binom{n}{4}$ sets of four species, the quartet distance can be calculated in time $O(n^4)$ by examining the sets one by one. Steel and Penny in [14] present an algorithm for computing the quartet distance in time $O(n^3)$. Bryant *et al.* in [5] present an algorithm that computes the quartet distance in time $O(n^2)$. In this paper, we present an algorithm that computes the quartet distance in time $O(n \log^2 n)$, making it possible to compare much larger evolutionary trees. Our solution is based on two techniques: the smaller-half trick, also used by methods for finding tandem repeats in strings, see e.g. [15], and a data structure related to the data structure for dynamic expression trees [7].

The rest of the paper is organized as follows. In Sect. 2, we introduce quartets and present our algorithm for computing the quartet distance between two unrooted evolutionary trees. In Sect. 3, we describe a hierarchical decomposition of unrooted trees which is an essential part of the data structure used by our algorithm. In Sect. 4, we present the details of our data structure.

2 The Algorithm

As mentioned, we in this paper by an *evolutionary tree* mean an unrooted tree where all nodes are either leaves (i.e. have degree one) or have degree three, and where the leaves are uniquely labeled by the elements of a set S of species. Let n denote the size of S .

For an evolutionary tree T , the *quartet topology* of four species a, b, c , and d is the topological subtree of T induced by these species. In general, the possible quartet topologies for species a, b, c, d are the four shown in Fig. 1. Of these, the right-most does not occur in our setting, due to the assumption about all internal nodes having degree three. Hence, the quartet topology is a pairing of the four species into two pairs, defined by letting a and b be a pair if among the three paths in T from a to respectively b, c , and d , the path to b is the first to separate from the others.

Given two evolutionary trees T_1 and T_2 on the same set S of species, the *quartet distance* between the two trees is the number of four-sets $\{a, b, c, d\} \subseteq S$, for which the quartet topologies in T_1 and T_2 differ. As there are $\binom{n}{4}$ different four-sets in S , the quartet distance can also be calculated as $\binom{n}{4}$ minus the number of four-sets for which the quartet topologies in T_1 and T_2 are identical. In this paper, we give an algorithm for finding this number in time $O(n \log^2 n)$.

To facilitate the counting of identical quartet topologies in the two trees, we view the quartet topology of a four-set $\{a, b, c, d\}$ as two *oriented* quartet topologies given by the two possible orientations of the “middle edge” of the topology. Figure 2 shows the two oriented quartet topologies arising from one unoriented quartet topology.

Clearly, the number of identical oriented quartet topologies between the trees T_1 and T_2 is twice the number of identical unoriented quartet topologies.

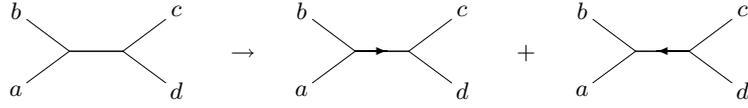


Figure 2. The two orientations of a quartet topology

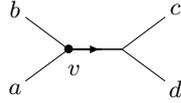


Figure 3. A generic quartet

The goal of our algorithm is to count identical oriented quartet topologies. For brevity, we in the rest of this paper let the word *quartet* denote an oriented quartet topology of a four-set.

We associate quartets to internal nodes in T_1 as follows: Consider the generic quartet in Fig. 3, where the orientation is from the pair $\{a, b\}$ to the pair $\{c, d\}$. There is a unique node v in T_1 where the paths from a and b to c (and d) meet. We associate the quartet of Fig. 3 with the node v . This partitions the $2^{\binom{n}{4}}$ quartets into $n - 2$ disjoint sets, as there are $n - 2$ internal nodes in a tree of n leaves, when all internal nodes have degree three.

For an internal node v in T_1 , we by the subtrees *incident* to v mean the three subtrees which arise if v and its three incident edges are removed from T_1 . These are shown in Fig 4, denoted by A , B , and C . The number of quartets associated with v is given by the expression

$$\binom{|A|}{2} \cdot |B| \cdot |C| + \binom{|B|}{2} \cdot |C| \cdot |A| + \binom{|C|}{2} \cdot |A| \cdot |B|,$$

where $|T|$ denotes the number of leaves in subtree T .

The strategy of the algorithm is for each internal node v in T_1 to count how many of the quartets associated with v which also are quartets of T_2 . The sum over all nodes in T_1 of these counts then gives the required number of identical quartets in T_1 and T_2 .

To do this, the algorithm colors the elements of S using the three colors \mathcal{A} , \mathcal{B} , and \mathcal{C} . The coloring is maintained via the data structure described in Sect. 4. When v is an internal node in T_1 , we say that the elements of S are colored *according* to v if the labels of the leaves of one of the three subtrees incident to v all have color \mathcal{A} , the labels of the leaves of another of the subtrees all have color \mathcal{B} , and the labels of the leaves of the remaining subtree all have color \mathcal{C} .

The central feature of the data structure is that if the elements of S are colored according to a node v in T_1 , then it can return in constant time the number of quartets associated with v which also are quartets in T_2 . The data

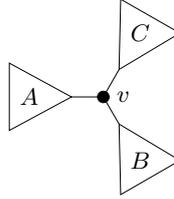


Figure 4. Subtrees incident to an internal node v

structure also allows the color of an element to be changed in time $O(\log n)$, given a pointer to the element.

The algorithm starts by rooting T_1 at an arbitrary leaf. It then calculates the size $|v|$ of each node v in T_1 during a postorder traversal starting at the root, where $|v|$ denotes the number of leaves below v , and stores this information in the nodes. It also colors all elements of S by the color \mathcal{C} .

The algorithm then calculates the desired sum of the counts for all internal nodes of T_1 in a recursive fashion, starting at the single child of the root of T_1 . To achieve the claimed complexity, the algorithm at a node v will recurse first on its smaller child, then on its larger child, and finally add the count for v to the sum calculated so far.

In Fig. 5, the algorithm is described in pseudo-code as a recursive procedure $\text{Count}(v)$. A call to $\text{Count}(v)$ returns the sum of the counts for the internal nodes of T_1 which are below v . Initially, it is called with v set to the single child of the root of T_1 . The two routines $\text{Small}(v)$ and $\text{Large}(v)$ return the child of v having smallest, respectively largest, size. The routine $\text{NodeCount}(v)$ is a call to the data structure of Sect. 4, returning the count for the node v . The routine $\text{ColorLeaves}(v, \mathcal{X})$ colors by the color \mathcal{X} all elements in the data structure which are labels of leaves below v in T_1 . This is done by a traversal of the subtree in T_1 rooted at v . By maintaining bi-directional pointers between elements of S in the data structure and the leaves in T_1 and T_2 which they label, this takes time $O(|v| \cdot \log n)$.

Theorem 1. *Let T_1 and T_2 be two unrooted evolutionary trees on the same set S of species, and let all internal nodes in the trees have degree three. Then the quartet distance between T_1 and T_2 can be found in time $O(n \log^2 n)$.*

Proof. We here assume the existence of the data structure discussed above. This existence is proven in Sect. 4. By induction on the number of calls to $\text{Count}(v)$, it follows that the algorithm above maintains the invariants:

1. At the *beginning* of the execution of an instance of $\text{Count}(v)$, all elements in S are colored by the color \mathcal{C} .
2. At the *end* of the execution of an instance of $\text{Count}(v)$, all elements in S which are labels of leaves *below* v in T_1 are colored by the color \mathcal{A} , and all other elements in S are colored by the color \mathcal{C} .

```

Procedure Count( $v$ )
  if  $v$  is a leaf then
    color  $v$  by the color  $\mathcal{A}$ 
    return 0
  else
     $x = \text{Count}(\text{Small}(v))$ 
    ColorLeaves( $\text{Small}(v)$ ,  $\mathcal{C}$ )
     $y = \text{Count}(\text{Large}(v))$ 
    ColorLeaves( $\text{Small}(v)$ ,  $\mathcal{B}$ )
     $z = \text{NodeCount}(v)$ 
    ColorLeaves( $\text{Small}(v)$ ,  $\mathcal{A}$ )
    return  $x + y + z$ 

```

Figure 5. The algorithm

The invariants imply that when a call to $\text{NodeCount}(v)$ takes place, labels of leaves in the subtree of $\text{Small}(v)$ are labeled by the color \mathcal{B} , labels of leaves in the subtree of $\text{Large}(v)$ are labeled by the color \mathcal{A} , and the remaining elements are labeled by the color \mathcal{C} . In other words, the elements of S are colored according to v . Correctness of the algorithm follows.

For complexity, note that the work incurred by an instance of $\text{Count}(v)$, not counting recursive calls made during this instance, is $O(|\text{Small}(v)| \cdot \log n)$. Let this work be accounted for by charging each leaf below $\text{Small}(v)$ in T_1 (or v itself, if it is a leaf) an amount $O(\log n)$ of work. For a given leaf, this charging can only happen at nodes v on the path from the leaf to the root where the path goes from $\text{Small}(v)$ to v . As the size of v is at least twice as large as the size of $\text{Small}(v)$, this can only happen $\log n$ times. Hence, each leaf is at most charged $O(\log^2 n)$ work in total, and the result follows. \square

3 Hierarchical Decomposition

An essential part of the data structure in Sect. 4 is a *hierarchical decomposition* of the evolutionary tree T_2 . Given an unrooted tree T where all nodes have degree at most three, we in the following describe how to obtain a hierarchical decomposition of T with logarithmic height. Our decomposition is very similar to the decompositions used for solving the parallel and dynamic expression tree evaluation problems [3,7], but in our setting the underlying tree is considered to be unrooted.

We base our hierarchical decomposition on the notion of *components*. We define a component C in T to be one of the following:

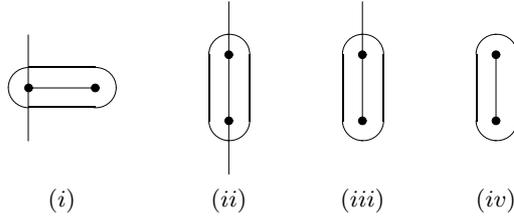


Figure 6. The four possible types of compositions of components

1. A set consisting of a single node of T .
2. A connected subset of the nodes of T , such that at most two nodes in C are connected by an edge to nodes in $T \setminus C$.

In other words, a component is either a set consisting of a single node, or a connected subset of nodes such that the cut defined by the subset is of size at most two. The *external edges* of a component C of T are the edges in T connecting nodes in C and $T \setminus C$. The *degree* of a component is the number of external edges of the component. By the second condition above, a component with two or more nodes can have degree at most two.

Each node of T (including leaves) constitutes a component of type 1. Components of type 2 are formed as the union of two adjacent components C' and C'' , where C' and C'' are said to be adjacent if there exist an edge (u, v) in T such that $u \in C'$ and $v \in C''$. We call such a union a *composition*. We only allow the four compositions depicted in Fig. 6. Nodes represent contracted components and ovals represent component compositions. Types (i), (iii), and (iv) are the cases where a component with degree one is composed with a component of degree three, two, and one respectively. Type (ii) is the case where two components with degree two are composed into a new component with degree two. Note that each composition of two components corresponds to a unique edge in the tree T , namely the edge connecting the two components.

A hierarchical decomposition of an unrooted tree T is a rooted binary tree, in the following denoted $H(T)$. Each node of $H(T)$ represents a component in T . Leaves of $H(T)$ represent components of type 1, and there is a one-to-one mapping between these components and the leaves of $H(T)$. An internal node v of $H(T)$ represent a component of type 2 formed by the composition of the two components represented by the children of v .

Lemma 1. *For every unrooted tree with n nodes and all nodes having degree at most three, there exists a hierarchical decomposition tree with height $O(\log n)$. The decomposition can be computed in time $O(n)$.*

Proof. Given a tree with n nodes, we construct a hierarchical decomposition bottom-up in $O(\log n)$ steps. Initially we start with each node being a component

by itself. In each step we greedily select an arbitrary maximal set of independent compositions, using time linear in the number of remaining components.

Let n denote the number of components at the beginning of a step. A composition of type (iv) will occur if and only if $n = 2$. If $n \geq 3$, let n_1 , n_2 , and n_3 denote the number of components of degree one, two and three respectively. We have $n = n_1 + n_2 + n_3$ and $n_3 = n_1 - 2$. Since $n \geq 3$, there are n_1 possible compositions of types (i) and (iii). We observe that the only edges not corresponding to legal compositions are edges connecting a component of degree three with a component of degree two or three. Since there are at most $3n_3$ such edges, the number of possible compositions is at least $n - 1 - 3n_3 = n - 3n_1 + 5$. If $n_1 < n/4$, then this bound is at least $n/4$. It follows that there are always at least $n/4$ possible compositions. Since each possible composition can conflict with at most two other compositions, any maximal set of non-conflicting compositions has size at least $n/12$.

After k steps, at most $n(11/12)^k$ components will remain. In particular, at most one component will remain after at most $\lceil \log_{12/11} n \rceil$ steps, so the height of the hierarchical decomposition tree is bounded by $\lceil \log_{12/11} n \rceil$. Since the number of components decreases geometrically for each step, the total time becomes $O(n)$. \square

4 Counting Quartets in Components

Given a coloring of the elements in S with the colors \mathcal{A} , \mathcal{B} , and \mathcal{C} , and given a quartet oriented as in Fig. 3 from the pair $\{a, b\}$ to the pair $\{c, d\}$, we say that the quartet is *compatible* with the coloring if a and b have different colors, and c and d both have the remaining color. Let T be an evolutionary tree for S , and let $H(T)$ be the hierarchical decomposition tree for T , as defined in Sect. 3.

Lemma 2. *When S is colored according to a choice of v in T , then the set of quartets compatible with the coloring is exactly the quartets associated with v .*

Proof. Follows from the definitions of quartets being compatible with a coloring and quartets being associated with a node. \square

We now describe how to decorate the nodes of $H(T)$ with information such that the number of quartets of T which are compatible with a given coloring of S can be returned in constant time. Furthermore, for a given coloring, the information can be generated in $O(n)$ time, and if one element of S changes color, the information can be updated in time $O(\log n)$.

For each node of $H(T)$, we store a tuple (a, b, c) of integers and a function F . Recall that a node in $H(T)$ represents a component in T . The integers a , b , and c of the tuple are the number of elements at the leaves contained in this component which are colored \mathcal{A} , \mathcal{B} , and \mathcal{C} , respectively. A component has k external edges for k between zero and three (the case of zero external edges occurs only at the root of $H(T)$). The function F has three variables for each of the external edges of the component. For a component with at least one external edge, we number

these edges arbitrarily from 1 to k and denote the three variables corresponding to edge i by \mathbf{a}_i , \mathbf{b}_i , and \mathbf{c}_i . If an external edge were removed from T , two subtrees of T would arise, of which one does not contain the component in question. We call this subtree the subtree *induced* by the external edge. The variables \mathbf{a}_i , \mathbf{b}_i , and \mathbf{c}_i denote the number of elements in leaves from the subtree induced by edge i which are colored \mathcal{A} , \mathcal{B} , and \mathcal{C} , respectively. Finally, F states, as a function of the variables \mathbf{a}_i , \mathbf{b}_i , and \mathbf{c}_i for $1 \leq i \leq k$, the number of the quartets which are both associated (in the sense defined in Sect. 2) with nodes in the component *and* are compatible with the given coloring. It will turn out that F is actually a polynomial of total degree at most four.

The root of $H(T)$ represents a component which comprises the entire tree T , i.e. the component has no external nodes, so the function F stored there is actually a constant. Hence, the number of quartets of T which are compatible with a given coloring of S is part of the information stored at the root.

Lemma 3. *The tree $H(T)$ can be decorated with the information described above in time $O(n)$.*

Proof. The information is computed in a bottom up fashion during a traversal of $H(T)$. We first describe how the information for leaves in $H(T)$ is generated, i.e. for nodes representing single node components. Recall that a node in T is either a leaf and has degree one, or is an internal node and has degree three.

For a component consisting of a single leaf with an element colored \mathcal{A} , \mathcal{B} , or \mathcal{C} , the tuple is $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$, respectively. The function F is identically zero, as quartets are only associated with internal nodes of T , not with leaves of T .

For a component consisting of a single degree three node u , the tuple is $(0, 0, 0)$, as no leaves of T are contained in the component. The function F should count the number of quartets which are both compatible with the coloring and associated with u in T . A quartet oriented from the pair $\{a, b\}$ to the pair $\{c, d\}$ fulfills this requirement precisely when c and d are contained in one of the three subtrees induced by the external edges of the component, and they have the same color, and a and b each are in one of the remaining two induced subtrees and each have one of the remaining two colors. For the case that c and d are in the subtree induced by edge number one and have color \mathcal{A} , the number of quartets fulfilling this is

$$\binom{\mathbf{a}_1}{2} \cdot (\mathbf{b}_2 \mathbf{c}_3 + \mathbf{b}_3 \mathbf{c}_2).$$

Summing over all $3 \cdot 3 = 9$ choices of the induced subtree and color for c and d , we get:

$$\begin{aligned}
F(\mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1, \mathbf{a}_2, \mathbf{b}_2, \mathbf{c}_2, \mathbf{a}_3, \mathbf{b}_3, \mathbf{c}_3) \\
&= \binom{\mathbf{a}_1}{2} \cdot (\mathbf{b}_2 \mathbf{c}_3 + \mathbf{b}_3 \mathbf{c}_2) + \binom{\mathbf{a}_2}{2} \cdot (\mathbf{b}_1 \mathbf{c}_3 + \mathbf{b}_3 \mathbf{c}_1) + \binom{\mathbf{a}_3}{2} \cdot (\mathbf{b}_2 \mathbf{c}_1 + \mathbf{b}_1 \mathbf{c}_2) \\
&+ \binom{\mathbf{b}_1}{2} \cdot (\mathbf{a}_2 \mathbf{c}_3 + \mathbf{a}_3 \mathbf{c}_2) + \binom{\mathbf{b}_2}{2} \cdot (\mathbf{a}_1 \mathbf{c}_3 + \mathbf{a}_3 \mathbf{c}_1) + \binom{\mathbf{b}_3}{2} \cdot (\mathbf{a}_2 \mathbf{c}_1 + \mathbf{a}_1 \mathbf{c}_2) \\
&+ \binom{\mathbf{c}_1}{2} \cdot (\mathbf{b}_2 \mathbf{a}_3 + \mathbf{b}_3 \mathbf{a}_2) + \binom{\mathbf{c}_2}{2} \cdot (\mathbf{b}_1 \mathbf{a}_3 + \mathbf{b}_3 \mathbf{a}_1) + \binom{\mathbf{c}_3}{2} \cdot (\mathbf{b}_2 \mathbf{a}_1 + \mathbf{b}_1 \mathbf{a}_2)
\end{aligned}$$

We now turn to the generation of the information stored in the internal nodes of $H(T)$. Consider the composition of two components C' and C'' . Let (a', b', c') and F' , and (a'', b'', c'') and F'' be the information stored at the nodes representing the components C' and C'' . The information stored at the node representing the composition C of C' and C'' is $(a' + a'', b' + b'', c' + c'')$ and F , where F depends on the type of composition. If the component composition is of type *(ii)*, we consider the case where the numbering of external edges of components is such that the first external edge of C' and C'' is the edge connecting C' and C'' , and the second external edge of C' is the first external edge of C , and the second external edge of C'' is the second external edge of C . The remaining cases of numbering of external edges are obtained by appropriate changes of the arguments to F' and F'' .

$$\begin{aligned}
F(\mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1, \mathbf{a}_2, \mathbf{b}_2, \mathbf{c}_2) \\
&= F'(\mathbf{a}_2 + \mathbf{a}'', \mathbf{b}_2 + \mathbf{b}'', \mathbf{c}_2 + \mathbf{c}'', \mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1) \\
&+ F''(\mathbf{a}_1 + \mathbf{a}', \mathbf{b}_1 + \mathbf{b}', \mathbf{c}_1 + \mathbf{c}', \mathbf{a}_2, \mathbf{b}_2, \mathbf{c}_2)
\end{aligned}$$

Component compositions of type *(iii)* and *(iv)* are identical to type *(ii)*, except that the definition of F is simpler. For type *(iii)* we have (assuming that C'' is the component of degree one)

$$F(\mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1) = F'(\mathbf{a}'', \mathbf{b}'', \mathbf{c}'', \mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1) + F''(\mathbf{a}_1 + \mathbf{a}', \mathbf{b}_1 + \mathbf{b}', \mathbf{c}_1 + \mathbf{c}') ,$$

and for type *(iv)* we have

$$F = F'(\mathbf{a}'', \mathbf{b}'', \mathbf{c}'') + F''(\mathbf{a}', \mathbf{b}', \mathbf{c}') .$$

Note that for type *(iv)* compositions, F is a constant.

Finally, we for type *(i)* compositions get the following expression for F , assuming C' has degree one and the first and second external edges of C are the second and third external edges of C'' , respectively.

$$\begin{aligned}
F(\mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1, \mathbf{a}_2, \mathbf{b}_2, \mathbf{c}_2) \\
&= F'(\mathbf{a}_1 + \mathbf{a}_2 + \mathbf{a}'', \mathbf{b}_1 + \mathbf{b}_2 + \mathbf{b}'', \mathbf{c}_1 + \mathbf{c}_2 + \mathbf{c}'') \\
&+ F''(\mathbf{a}', \mathbf{b}', \mathbf{c}', \mathbf{a}_1, \mathbf{b}_1, \mathbf{c}_1, \mathbf{a}_2, \mathbf{b}_2, \mathbf{c}_2)
\end{aligned}$$

By structural induction on the definition of the F functions stored at components, it follows that F is always a polynomial of total degree at most four.

Polynomials with total degree at most four and at most nine variables can be stored in constant space by storing the coefficients of the polynomials, and they can be manipulated in constant time, e.g. when adding or composing two polynomials. We conclude that for a component C which is the composition of two components C' and C'' , the information to be stored at C can be computed in constant time, provided that the information stored at C' and C'' is known. It follows that $H(T)$ can be decorated in time $O(n)$. \square

Lemma 4. *The decoration of $H(T)$ can be updated in $O(\log n)$ time when the color of an element in S changes.*

Proof. From the proof of Lemma 3 we know that the decoration of a node in $H(T)$ only depends on the decoration of the children of the node in $H(T)$, i.e. the only decorations that need to be updated in $H(T)$ while changing the color of an element in S are the ancestors of the leaf in $H(T)$ corresponding to the element. Since $H(T)$ has height $O(\log n)$ and the decoration of a node takes constant time to compute knowing the decoration of the children, it follows that the decoration of $H(T)$ can be updated in time $O(\log n)$. \square

Lemma 5. *When S is colored according to a choice of v in T_1 , then the set of quartets compatible with the coloring is exactly the quartets associated with v .*

Proof. Follows from the definitions of the colors and compatible quartets. \square

Corollary 1. *If the above construction is done with T_2 for T , and the coloring of S is according to a choice of v in T_1 , then the quartets in T_2 compatible with the coloring are exactly the quartets which are in both T_1 and T_2 . Furthermore, the number of such quartets is exactly the value of the constant function F stored at the root of $H(T_2)$.*

References

1. B. L. Allen and M. Steel. Subtree transfer operations and their induced metrics on evolutionary trees. *Annals of Combinatorics*, 5:1–13, 2001.
2. 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, 1997.
3. R. P. Brent. The parallel evaluation of general arithmetic expressions. *Journal of the ACM*, 21(2):201–206, Apr. 1974.
4. G. S. Brodal, R. Fagerberg, C. N. S. Pedersen, and A. Östlin. The complexity of constructing evolutionary trees using experiments. In *Proc. 28th International Colloquium on Automata, Languages, and Programming*, volume 2076 of *Lecture Notes in Computer Science*, pages 140–151. Springer-Verlag, 2001.
5. D. Bryant, J. Tsang, P. E. Kearney, and M. Li. Computing the quartet distance between evolutionary trees. In *Proceedings of the Eleventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 285–286, N.Y., Jan. 9–11 2000. ACM Press.

6. P. Buneman. The recovery of trees from measures of dissimilarity. *Mathematics in Archeological and Historical Sciences*, pages 387–395, 1971.
7. R. F. Cohen and R. Tamassia. Dynamic expression trees. *Algorithmica*, 13(3):245–265, 1995.
8. G. Estabrook, F. McMorris, and C. Meacham. Comparison of undirected phylogenetic trees based on subtrees of four evolutionary units. *Syst. Zool.*, 34(2):193–200, 1985.
9. M. Farach, S. Kannan, and T. J. Warnow. A robust model for finding optimal evolutionary trees. *Algorithmica*, 13(1/2):155–179, 1995.
10. D. Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge University Press, 1997.
11. A. Lingas, H. Olsson, and A. Östlin. Efficient merging, construction, and maintenance of evolutionary trees. In *Proc. 26th Int. Colloquium on Automata, Languages and Programming*, volume 1644 of *Lecture Notes in Computer Science*, pages 544–553. Springer-Verlag, 1999.
12. D. F. Robinson and L. R. Foulds. Comparison of weighted labelled trees. In *Combinatorial mathematics, VI (Proc. Sixth Austral. Conf., Univ. New England, Armidale, 1978)*, Lecture Notes in Mathematics, pages 119–126. Springer, Berlin, 1979.
13. D. F. Robinson and L. R. Foulds. Comparison of phylogenetic trees. *Math. Biosci.*, 53(1-2):131–147, 1981.
14. M. Steel and D. Penny. Distribution of tree comparison metrics—some new results syst. *Syst. Biol.*, 42(2):126–141, 1993.
15. J. Stoye and D. Gusfield. Simple and flexible detection of contiguous repeats using a suffix tree. In *Proc. 9th Annual Symposium on Combinatorial Pattern Matching*, volume 1448 of *Lecture Notes in Computer Science*, pages 140–152. Springer-Verlag, 1998.
16. M. S. Waterman and T. F. Smith. On the similarity of dendrograms. *Journal of Theoretical Biology*, 73:789–800, 1978.