Optimal Completion and Comparison of Incomplete Phylogenetic Trees Under Robinson-Foulds Distance

4 Keegan Yao 🖂

5 Department of Computer Science and Engineering, University of Connecticut, Storrs, USA

6 Mukul S. Bansal 🖂

7 Department of Computer Science and Engineering, University of Connecticut, Storrs, USA

8 — Abstract -

The comparison of phylogenetic trees is a fundamental task in phylogenetics and evolutionary biology. In many cases, these comparisons involve trees inferred on the same set of leaves, and many distance 10 measures exist to facilitate such comparisons. However, several applications in phylogenetics require 11 the comparison of trees that have non-identical leaf sets. The traditional approach for handling 12 such comparisons is to first restrict the two trees being compared to just their common leaf set. An 13 alternative, conceptually superior approach that has shown promise is to first *complete* the trees by 14 adding missing leaves so that the completed trees have identical leaf sets. This alternative approach 15 requires the computation of optimal completions of the two trees that minimize the distance between 16 them. However, no polynomial-time algorithms currently exist for this optimal completion problem 17 under any standard phylogenetic distance measure. 18

In this work, we provide the first polynomial-time algorithms for the above problem under the widely used Robinson-Foulds (RF) distance measure. This hitherto unsolved problem is referred to as the RF(+) problem. We (i) show that a recently proposed linear-time algorithm for a restricted version of the RF(+) problem is a 2-approximation for the RF(+) problem, and (ii) provide an exact $O(nk^2)$ -time algorithm for the RF(+) problem, where n is the total number of distinct leaf labels in the two trees being compared and k, bounded above by n, depends on the topologies and leaf set overlap of the two trees. Our results hold for both rooted and unrooted binary trees.

We implemented our exact algorithm and applied it to several biological datasets. Our results show that completion-based RF distance can lead to very different inferences regarding phylogenetic similarity compared to traditional RF distance. An open-source implementation of our algorithms is freely available from https://compbio.engr.uconn.edu/software/RF_plus.

 $_{32}$ $\,$ Mathematics of computing \rightarrow Trees

Keywords and phrases Phylogenetic tree comparison, Robinson-Foulds Distance, Optimal tree
 completion, Algorithms, Dynamic programming

- 35 Digital Object Identifier
- ³⁶ Funding Keegan Yao: University of Connecticut Summer Undergraduate Research Fund
- 37 Mukul S. Bansal: US National Science Foundation award IIS 1553421

38 **1** Introduction

³⁹ *Phylogenetic trees*, or simply *phylogenies*, are leaf-labeled trees that depict the evolutionary

- $_{40}$ $\,$ relationships between different species, genes, or other biological entities such as cells in an
- $_{\rm 41}$ $\,$ organism or individuals from a population. In phylogenetic trees, leaf nodes represent extant
- ⁴² entities while internal nodes represent hypothetical ancestors. Many different methodologies,
- ⁴³ algorithms, and data types exist for estimating phylogenies, and there is often considerable
- 44 uncertainty and error in their inference, with different methods or data types suggesting



different evolutionary relationships between the same extant entities. Many distance (or 45 similarity) measures have therefore been developed for systematically comparing different phylogenetic trees, including the widely used Robinson-Foulds distance [29], triplet and 47 quartet distances [14, 17], nearest neighbor interchange (NNI) and subtree prune and regraft 48 (SPR) distances [31, 18, 34], maximum agreement subtrees [19, 2, 15], nodal distance [9], 49 geodesic distance [23] and others. However, these distance measures implicitly assume that 50 the two trees being compared have identical leaf sets, an assumption that is often violated 51 in practice. Indeed, several applications, such as supertree construction [24, 6, 10, 32, 1], 52 phylogenetic database search [28, 30, 11, 25], and clustering of phylogenetics [20, 35], require 53 the computation of distances between trees with partially overlapping leaf sets. 54

The traditional approach to comparing two trees with only partially overlapping leaf sets 55 is to first restrict (i.e., prune down) both trees to their shared leaf set. This restriction based 56 approach, though simple to conceptualize and compute, can result in the loss of valuable 57 topological information through scrapping of leaves that are not common to both trees. 58 An alternative approach to comparing trees with non-identical leaf sets is to *complete* or 59 fill in each of the input trees to the union of their leaf sets in a way which minimizes the 60 distance between them, and then compute their distance. This approach, though conceptually 61 more complex, successfully incorporates all topological information in both the trees being 62 compared. In addition to its more complete use of topological information, the completion 63 based approach also has the benefit of a larger range of attainable values due to comparisons 64 over larger extended trees rather than smaller induced trees. Despite these advantages, 65 no polynomial-time algorithms currently exist for completion based comparison under any 66 standard phylogenetic distance measure. In this work, we provide the first polynomial-time 67 algorithms for optimal completion and comparison of incomplete phylogenetic trees under 68 the widely used Robinson-Foulds (RF) distance measure. Following existing literature [4], 69 we refer to completion based RF distance as RF(+), the traditional restriction based RF 70 distance as RF(-), and the problem of computing the RF(+) distance between two trees as 71 the RF(+) problem. Figure 1 illustrates the difference between RF(-) and RF(+) distances. 72 Previous work. The idea of completion based Robinson-Foulds distance arose at least 73 decade ago when Cotton and Wilkinson introduced majority-rule supertrees [13] and \mathbf{a} 74 defined two variants, majority-rule(-) and majority-rule(+) supertrees, based on RF(-) and 75 RF(+), respectively. Completion based majority-rule(+) supertrees and some variants were 76 subsequently shown to have many desirable properties [16]. Later, Kupczok [22] characterized 77 the RF(+) distance for the restricted special case where the leaf set of one tree is a subset 78 of the leaf set of the other in terms of incompatible splits between the two trees. For this 79 restricted special case, referred to as the One Tree RF(+) (OT-RF(+)) problem [4], an 80 $O(n^2)$ -time algorithm was proposed by Christensen et. al. in 2017 [12], where n is the 81 total number of distinct leaf labels in the two trees being compared. More recently, Bansal 82 proposed an optimal O(n)-time algorithm for this OT-RF(+) problem [3, 4]. Bansal also 83 proposed a restricted formulation of the RF(+) problem, called the *Extraneous-Clade-Free* 84 RF(+) (EF-RF(+)) problem, which is based on computing optimal completions that avoid 85 the creation of any subtrees formed by joining together two subtrees unique to each one of the 86 two input trees. Essentially, the EF-RF(+) problem disallows certain types of completions; 87 specifically, it ignores how subtrees exclusive to one input tree impact the overall optimal 88 position where subtrees from the other input tree should be added. Bansal showed that the 89 EF-RF(+) problem can be solved in O(n) time [4]. These linear-time algorithms for the 90 OT-RF(+) and EF-RF(+) problems are applicable to both rooted and unrooted trees. 91

92 Our Contributions. In this work, we provide the first polynomial-time algorithms for



Figure 1 RF(-) and RF(+) distances. The figure shows a "base" tree *S* and two other trees *U* and *V*, with Le(U) = Le(V), being compared to *S*. S_*, U_* and V_* represent the trees *S*, *U* and *V*, respectively, when restricted to the common leaf set. U^* and V^* are the optimal RF(+) completions of *U* and *V* with respect to *S*. S_U^* and S_V^* are the optimal RF(+) completions of *S* with respect to *U* and *V*, respectively. Filled in nodes represent matched nodes (Definition 2.2). Here, $RF(S_*, U_*) = 2$ and $RF(S_*, V_*) = 4$ while $RF(S_U^*, U^*) = 8$ and $RF(S_V^*, V^*) = 4$. Thus, in this example, *U* is closer to *S* than *V* under RF(-) but *V* is closer to *S* than *U* under RF(+).

the RF(+) problem for both rooted and unrooted trees. Specifically, we make the fol-93 lowing contributions: First, we show that the EF-RF(+) distance between two trees is a 94 2-approximation for the RF(+) distance between those trees. Since the EF-RF(+) problem 95 can be solved in O(n) time, this yields a linear time 2-approximation algorithm for the RF(+)96 problem. Second, we provide an $O(nk^2)$ -time exact algorithm for the RF(+) problem, where 97 k, bounded above by n, is the number of maximal subtrees exclusive to one input tree. And 98 third, we perform an extensive experimental study which demonstrates that the use of RF(+)99 distance can lead to very different inferences regarding phylogenetic similarity compared 100 to RF(-) distance. We also find that, in practice, EF-RF(+) distances are often very close 101 to RF(+) distances, suggesting that the linear-time algorithm for computing EF-RF(+)102 distances could be an excellent heuristic for estimating RF(+) distances between large trees. 103 The rest of this manuscript is organized as follows: Preliminaries and problem definitions

The rest of this manuscript is organized as follows: Preliminaries and problem definitions appear in the next section. We describe the linear time 2-approximation algorithm in Section 3, and the exact algorithm in Section 4. Section 5 shows how our algorithms can be extended to unrooted trees, and Section 6 describes the results of our experimental study. ¹⁰⁸ Concluding remarks appear in Section 7. Proofs of all lemmas and theorems from Sections 3
 ¹⁰⁹ and 4 appear in the Appendix.

110 2 Definitions and Preliminaries

We follow basic definitions and problem formulations from [4]. All trees will be unordered. 111 Given a tree T, we denote its node set, edge set, and leaf set by V(T), E(T), and Le(T), 112 respectively. The set of all non-leaf (i.e., internal) nodes of T is denoted by I(T). If T is 113 rooted, the root node of T is denoted by rt(T), the parent of a node $v \in V(T)$ by $pa_T(v)$, 114 its set of children by $Ch_T(v)$, and the (maximal) subtree of T rooted at v by T(v). If two 115 nodes in T have the same parent, they are called *siblings* of each other. If $pa_T(v)$ has exactly 116 two children, then we will denote the sibling of v as $sib_T(v)$. The least common ancestor, 117 denoted $lca_T(L)$, of a set $L \subseteq Le(T)$ in T is defined to be the node $v \in V(T)$ such that 118 $L \subseteq Le(T(v))$ and $L \not\subseteq Le(T(u))$ for any child u of v. For convenience, given a collection of 119 vertices a_1, \ldots, a_m in T, we will define $lca_T(a_1, \ldots, a_m) = lca_T(Le(T(a_1)) \cup \cdots \cup Le(T(a_m)))$. 120 Given a rooted tree T and $a, b \in V(T)$, we say that $a \leq b$ if $a \in V(T(b))$, and a < b if 121 $a \in V(T(b))$ and $a \neq b$. A rooted tree is *binary* if all of its internal nodes have exactly 122 two children, while an unrooted tree is *binary* if all its nodes have degree either 1 or 3. 123 Throughout this work, the term *tree* refers to binary trees with uniquely labeled leaves. 124

Let T be a rooted or unrooted tree. Given a set $L \subseteq Le(T)$, let T_L be the minimal subtree of T with leaf set L. We define the *leaf induced subtree* T[L] of T on leaf set L to be the tree obtained from T_L by successively removing each non-root node of degree two and adjoining its two neighbors.

▶ Definition 2.1 (Completion of a tree). Given a tree T and a set L' such that $Le(T) \subseteq L'$, a completion of T on L' is a tree T' such that Le(T') = L' and T'[Le(T)] = T.

If T is a rooted tree, for each node $v \in V(T)$, the clade $C_T(v)$ is defined to be the set 131 of all leaf nodes in T(v); i.e. $C_T(v) = Le(T(v))$. We denote the set of all clades of a rooted 132 tree T by Clade(T). This concept can be extended to unrooted trees as follows. If T is an 133 unrooted tree, each edge $(u, v) \in E(T)$ defines a partition of the leaf set of T into two disjoint 134 subsets $Le(T_u)$ and $Le(T_v)$, where T_u is the subtree containing node u and T_v is the subtree 135 containing node v, obtained when edge (u, v) is removed from T. The partition induced by 136 any edge $(u, v) \in E(T)$ is called a *split* and is represented by the set $\{Le(T_u), Le(T_v)\}$. The 137 set of all splits in an unrooted tree T is denoted by Split(T). 138

▶ Definition 2.2 (Matched and mismatched nodes). Given rooted trees S and T, and a node $v \in V(S)$, we call v a matched node with respect to T if $C_S(v) \in Clade(T)$, and a mismatched node otherwise. Analogously, $C_S(v)$ is called a matched clade if $C_S(v) \in Clade(T)$, and a mismatched clade otherwise.

The symmetric difference of two sets A and B, denoted by $A\Delta B$, is the set $(A \setminus B) \cup (B \setminus A)$. We now define the Robinson-Foulds distance and the two problems that we solve in this paper.

¹⁴⁶ ► Definition 2.3 (Robinson-Foulds distance). The Robinson-Foulds (RF) distance, RF(S,T), ¹⁴⁷ between two trees S and T is defined to be $|Clade(S)\Delta Clade(T)|$ if S and T are rooted trees, ¹⁴⁸ and $|Split(S)\Delta Split(T)|$ if S and T are unrooted trees.

Problem 1 (Rooted RF(+) (R-RF(+))). Given two rooted binary trees S and T, compute a binary completion S^{*} of S on $Le(S) \cup Le(T)$ and a binary completion T^{*} of T on $Le(S) \cup Le(T)$ such that $RF(S^*, T^*)$ is minimized. ¹⁵² ► Problem 2 (Unrooted RF(+) (U-RF(+))). Given two unrooted binary trees S and T, ¹⁵³ compute a binary completion S^* of S on $Le(S) \cup Le(T)$ and a binary completion T^* of T on ¹⁵⁴ $Le(S) \cup Le(T)$ such that $RF(S^*, T^*)$ is minimized.

These problems can equivalently be viewed as maximizing the number of matched clades 155 or minimizing the number of mismatched clades between completions of the input trees. 156 Our algorithms for the problems above rely on first computing exact solutions for restricted 157 variants of those problems. These restricted variants of R-RF(+) and U-RF(+) were first 158 proposed and defined in [4] and are referred to as the Extraneous-Clade-Free R-RF(+) (EF-159 R-RF(+) and Extraneous-Split-Free U-RF(+) (EF-U-RF(+)) problems. These restricted 160 variants are based on computing optimal completions that do not contain any subtrees 161 formed by joining together two subtrees unique to each one of the two input trees. Next, 162 we first define extraneous clades and extraneous splits, and then state the EF-R-RF(+) and 163 EF-U-RF(+) problems. 164

▶ Definition 2.4 (Extraneous clade [4]). Suppose S and T are rooted trees. Given completions S' and T' of S and T, respectively, on $Le(S) \cup Le(T)$, we define a clade of S' or T' to be an extraneous clade if it contains leaves from both S and T but no leaves that are common to S and T.

An extraneous split is simply the analogous notion for unrooted trees and we refer the reader to [4] for a formal definition. The corresponding problem variants can now be defined as follows:

▶ Problem 3 (Extraneous-Clade-Free R-RF(+) (EF-R-RF(+)) [4]). Given two rooted trees S and T, compute a completion S' of S on $Le(S) \cup Le(T)$ and a completion T' of T on $Le(S) \cup Le(T)$ such that S' and T' do not contain any extraneous clades and RF(S',T') is minimized.

▶ Problem 4 (Extraneous-Split-Free U-RF(+) (EF-U-RF(+)) [4]). Given two unrooted trees S and T such that $|Le(S) \cap Le(T)| \ge 2$, compute a completion S' of S on $Le(S) \cup Le(T)$ and a completion T' of T on $Le(S) \cup Le(T)$ such that S' and T' do not contain any extraneous splits and RF(S', T') is minimized.

Figure 2 provides examples of completions with and without extraneous clades. Both the EF-R-RF(+) and EF-U-RF(+) problems can be solved optimally in linear time [4].

¹⁸² Note. In the remainder of this section, as well as in Sections 3 and 4 we focus on only the ¹⁸³ rooted version of RF(+), i.e., on the R-RF(+) problem, and implicitly assume that the two ¹⁸⁴ trees being compared, S and T, are rooted.

Node coloring scheme for rooted trees. For ease of presentation, we assign a color to some of the nodes of the two rooted input trees as follows. These node colorings can also be used to define red and green subtrees.

▶ Definition 2.5 (Red and Green Nodes). Let S and T be two arbitrary rooted trees. A node $v \in V(S)$ is called a red node (with respect to T) if $Le(S(v)) \subseteq Le(S) \setminus Le(T)$. Analogously, a node $v \in V(T)$ is called a green node (with respect to S) if $Le(T(v)) \subseteq Le(T) \setminus Le(S)$.

▶ Definition 2.6 (Red and Green Subtrees). A subtree S(u), where $u \in V(S)$, is called a red subtree of S if u is a red node. A subtree T(u), where $u \in V(T)$, is called a green subtree of T if u is a green node. A subtree S(u), where $u \in V(S)$, is called a maximal red subtree of S if S(u) is a red subtree and either u = rt(S) or $pa_S(u)$ is not red. A subtree T(u), where $u \in V(T)$, is called a maximal green subtree of T if T(u) is a green subtree and either



Figure 2 EF-RF(+) and RF(+) completions. S', T' are optimal EF-R-RF(+) completions (without extraneous clades) of S and T, respectively, and completions S^*, T^* are optimal RF(+) completions. Nodes labeled with downward and upward pointing triangles are red and green nodes, respectively, as defined in Definition 2.5. Filled in nodes correspond to matched clades.

¹⁹⁶ u = rt(T) or $pa_T(u)$ is not green. Note that all nodes in a red (green) subtree must be red ¹⁹⁷ (green).

¹⁹⁸ Under this node coloring, completing a tree S with respect to tree T entails adding all ¹⁹⁹ the green leaves of T into S and completing a tree T with respect to tree S entails adding, or ²⁰⁰ grafting, all the red leaves of S into T. Importantly, as we show later in Theorem 3.1, under ²⁰¹ R-RF(+) problem, there exist optimal completions of S and T in which all grafted subtrees ²⁰² are maximal red or green subtrees. In other words, to optimally complete S we must only ²⁰³ add the maximal green subtrees of T to S, and vice versa.

Notational conventions. S and T will denote the two given (input) trees to be completed/compared. Going forward, we will generally use S' and T' to represent completions (optimal or non-optimal) with *no* extraneous clades, and S^* and T^* to represent completions that *may* include extraneous clades.

3 EF-R-RF(+) is a 2-Approximation for R-RF(+)

Observe that any optimal pair of R-RF(+) completions can be modified into a pair of (not 209 necessarily optimal) EF-R-RF(+) completions by breaking apart any existing extraneous 210 clades and reinserting the red/green leaves in a manner that avoids forming extraneous 211 clades. In this section, we will show how to perform such a modification of optimal R-RF(+)212 completions so that the resulting increase in RF distance is appropriately bounded. This 213 will establish that EF-R-RF(+) distance is a 2-approximation for R-RF(+) distance and 214 will yield a linear-time 2-approximation algorithm for the R-RF(+) problem. We will first 215 establish the presence of *canonical* optimal R-RF(+) completions that satisfy some desirable 216 structural properties. 217

Notation and terminology. Given completions S^* and T^* of S and T, if there exists an extraneous clade $C_{T^*}(v)$ for some vertex $v \in T^*$, then we will call the subtree $T^*(v)$ an extraneous subtree. If the children s and t of v satisfy $C_{T^*}(s) \in Clade(S)$ and $C_{T^*}(t) \in$ Clade(T), then we will denote the extraneous subtree by $\{s,t\}$. To simplify notation, we will

²²¹ Clade(T), then we will denote the extraneous subtree by $\{s, t\}$. To simplify notation, we will ²²² write $pa_{T^*}\{s, t\}$ to express the parent $pa_{T^*}(lca_{T^*}(s, t))$ of the root node of the extraneous

write $pa_{T^*}\{s,t\}$ to express the parent $pa_{T^*}(lca_{T^*}(s,t))$ of the root node of the extraneous subtree $\{s,t\}$ in completion T^* . Likewise, we will write $sib_{T^*}\{s,t\}$ to express $sib_{T^*}(lca_{T^*}(s,t))$,

i.e., the sibling of the root node of extraneous subtree $\{s,t\}$ in T^* .

Next, we show that there always exists an optimal pair of R-RF(+) completions in which

all extraneous clades are of the form $\{s,t\}$, and any such extraneous clade appears in both

²²⁷ completions. We refer to such optimal R-RF completions S^* and T^* of S and T as canonical ²²⁸ optimal R-RF(+) completions.

Theorem 3.1. Let S and T be rooted binary trees. Then, there exist optimal completions S^* and T^* under the R-RF(+) problem with the following properties:

- 1. Every subtree inserted into S^* is a maximal green subtree of T, and every subtree inserted into T^* is a maximal red subtree of S,
- 233 2. Every extraneous subtree in S^* and T^* is of the form $\{s,t\}$, where s is the root of a 234 maximal red subtree in S and t is the root of a maximal green subtree in T,
- **3.** Every extraneous subtree $\{s,t\}$ which is a subtree of S^* is also a subtree of T^* and vice versa.

Decomposition of canonical optimal R-RF(+) completions. Given an extraneous subtree $\{s, t\}$ in canonical optimal R-RF(+) completions S^*, T^* of S and T, where $s \in V(S)$ and $t \in V(T)$, we define a *decomposition* of the extraneous subtree $\{s, t\}$ as a modification of the completions S^* and T^* , yielding new completions S' and T' with strictly fewer extraneous subtrees, as follows:

- 1. If either none or both of the nodes $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ are matches (in S^* and T^*), then the decomposition occurs as described below.
- In tree T^* , prune out the grafted subtree S(s) and regraft it at the parent edge of node $sib_{T^*}\{s,t\}$.
- ²⁴⁶ In tree S^* , prune out the grafted subtree T(t) and regraft it at the parent edge of ²⁴⁷ node $pa_{S^*}\{s,t\}$. If $pa_{S^*}\{s,t\} = rt(S^*)$, then create a new root node with children t²⁴⁸ and $pa_{S^*}\{s,t\}$.

249 2. Otherwise, if exactly one of the nodes $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ is a matched node (in S^* 250 and T^*), then the decomposition occurs as described below. Without loss of generality, 251 assume that $pa_{S^*}\{s,t\}$ is a match and $pa_{T^*}\{s,t\}$ a mismatch.

- ²⁵² In tree S^* , prune out the grafted subtree T(t) and regraft it at the parent edge of node ²⁵³ $sib_{S^*}\{s,t\}.$
- In tree T^* , prune out the grafted subtree S(s) and regraft it at the parent edge of that unique node $u \in V(T^*)$ for which $C_{T^*}(u) = C_{S^*}(pa_{S^*}\{s,t\})$. If $u = rt(S^*)$, then create a new root node with children s and $pa_{S^*}\{s,t\}$. Note that u must exist since $pa_{S^*}\{s,t\}$ is a matched node.

This decomposition is illustrated in Figure 3. The following lemma characterises how the RF distance between S^* and T^* is impacted as their extraneous subtrees are decomposed.

Lemma 3.2. Let S' and T' denote the trees obtained by decomposing extraneous subtree $\{s,t\}$ in completions S^{*} and T^{*}, respectively.

262 1. If $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ are both matched nodes then $RF(S',T') = RF(S^*,T^*)$.

263 2. If exactly one of $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ is a matched node then $RF(S',T') = RF(S^*,T^*)$.

3. If neither $pa_{S^*}\{s,t\}$ nor $pa_{T^*}\{s,t\}$ is a matched node then $RF(S',T') = RF(S^*,T^*) + 2$.

The 2-approximation now follows by appropriately bounding the number of extraneous subtrees $\{s, t\}$ that fall in category 3 of the above lemma.

▶ **Theorem 3.3.** Let S^* and T^* represent optimal completions of S and T, respectively, under the R-RF(+) problem. Let S' and T' represent optimal completions of S and T respectively under the EF-R-RF(+) problem. Then, $RF(S', T') \leq 2 \cdot RF(S^*, T^*)$.



Figure 3 Decomposition of extraneous clades. Shown here is a decomposition of completions S^* and T^* into completions S' and T'. Nodes labeled with downward and upward pointing triangles are red and green nodes, respectively. Extraneous subtree $\{b, g\}$ is of type 1 where both parents match, extraneous subtree $\{d, h\}$ is of type 1 where neither parent is a match, and extraneous subtree $\{e, i\}$ is of type 2. Matches between corresponding completions are denoted by filled in nodes.

²⁷⁰ 4 An Efficient Exact Algorithm for R-RF(+) Distance

As shown above, optimal EF-R-RF(+) completions 2-approximate RF(+) distance. We now show how to construct optimal R-RF(+) completions by modifying optimal EF-R-RF(+) completions.

Notation and terminology. We refer to EF-R-RF(+) completions resulting from the 274 Two Tree Completion Algorithm of [4] as canonical EF-R-RF(+) completions. This is due to 275 the way that maximal red and green subtrees are topologically well placed in such completions. 276 We will refer to the placement of a maximal colored subtree under the Two Tree Completion 277 Algorithm as a canonical EF-R-RF(+) position. The placement of each maximal red subtree 278 R of S, rooted at r, in canonical EF-R-RF(+) completion T' of T has the useful property 279 that all leaves $a \in Le(S) \cap Le(T)$ where $lca_S(a, r) = pa_S(r)$ also satisfy $lca_{T'}(a, r) = pa_{T'}(r)$, 280 and all leaves $b \in Le(S) \cap Le(T)$ where $lca_{T'}(b, r) > pa_{T'}(r)$ also satisfy $lca_S(b, r) > pa_S(r)$. 281 By Theorem 3.1, we know that there exists an optimal pair of R-RF(+) completions 282 where the only extraneous subtrees are of the form $\{s, t\}$. We will first show that a canonical 283 pair of R-RF(+) completions can be constructed by taking a canonical pair of EF-R-RF(+)284 completions and pairing up extraneous subtrees of the form $\{s, t\}$ in an optimal manner. We 285 will then design a recurrence relation which computes the best possible change to the RF 286 distance caused by pairing up extraneous subtrees of the form $\{s,t\}$, and show that this 287 change to the RF distance can be computed in near linear time depending on the leaf-set 288 overlap between the input trees. 289

Lemma 4.1. There exist canonical R-RF(+) completions S^* and T^* of rooted binary trees S and T such that every subtree grafted into S^* and T^* is either in an extraneous subtree or in its canonical EF-R-RF(+) position.

In the remainder of this section, let S', T' and S^*, T^* represent canonical EF-R-RF(+)



Figure 4 The tree T''. This figure shows the relationship between T', T'', and T^* . In this example, observe that there is exactly one extraneous subtree $\{s,t\}$ in the optimal completions S^* and T^* , and that $RF(S',T'') = RF(S^*,T^*) + 2$. Moreover, T'' in this example cannot be a completion of T since the green leaf i has been regrafted. But constructing T'' is simply an intermediary step for constructing completions S^* and T^* . Matches are denoted by filled in nodes.

and R-RF(+) completions of S and T, respectively. We will soon define the subproblems that are the basis of our dynamic programming algorithm. Before doing so, we motivate the dynamic programming recurrence relation with the following lemma, which describes a new useful tree T'' that is easier to construct from T' and preserves the important topological structure of T^* . Our dynamic programming algorithm actually constructs T'', and we can then easily use T'' to generate S^* and T^* .

▶ Lemma 4.2. Let T'' be the tree obtained by taking T^* and regrafting every extraneous subtree $\{s,t\}$ along the parent edge of $lca_{T^*}(lca_{T^*}(Le(sib_S(s))),t)$. Then $RF(S',T'') = RF(S^*,T^*) + 2m$, where m is the number of extraneous subtrees $\{s,t\}$ contained in T^* .

Note that T'' itself may not be a completion of T. In particular, in the construction of T'', pruning and regrafting the maximal green subtree T(t) is necessary if the extraneous subtree $\{s,t\}$ is formed and $lca_{T'}(s,t) \neq pa_{T'}(t)$. Moving any subtree of T in T' changes T'to no longer be a completion of T. Figure 4 shows a concrete example.

▶ **Definition 4.3.** Let the colors red and green be associated with the binary values 0 and 307 1, respectively. For $v \in V(T')$ and $c \in \{0,1\}$, let cMax(c,v) be the total number of maximal 308 subtrees of color c in T'(v). Moreover, let m be an integer such that $0 \le m \le cMax(c, v)$. We 309 define Cost(v, m, c) to be $\min_{\widehat{T}}(RF(S', \widehat{T}) - 2p - RF(S', T'))$, where \widehat{T} is obtained from T' by 310 regrafting maximal red and green subtrees in T'(v) under the constraint that each extraneous 311 subtree $\{s,t\}$ is grafted along the parent edge of $lca_{T'(v)}(s,t)$ and exactly m maximal c-colored 312 subtrees in T'(v) have been regrafted along the parent edge of v, excluding extraneous subtrees 313 (see Figure 5 for an example), and p denotes the number of extraneous subtrees of the form 314 $\{s,t\}$ in T. 315

In the trivial case when v is the root of a maximal c-colored subtree, we will say that it is possible to push one red subtree up to the parent edge of v or down from the parent edge of v.

Note that the *Cost*() subproblem builds the optimal RF(+) distance. However, the cost is defined based on Lemma 4.2 by constructing T'' and subtracting out the extraneous subtrees



Figure 5 Illustration of tree \hat{T} . The figure shows an example of what the tree \hat{T} might look like after computing Cost(u, 2, 1), where c and d have both been regrafted *iteratively* along the parent edge of u and not regrafted into an extraneous subtree. Note that the extraneous subtree $\{e, f\}$ has also been regrafted along the parent edge of u, though it does not contribute to the value of m = 2. In particular, $u = lca_{T'}(e, f)$, so the extraneous subtree $\{e, f\}$ will appear at the same position in \hat{T} and T''. Moreover, f is not included as one of the two maximal green subtrees grafted onto the parent edge of u since it is a part of an extraneous subtree. For each choice of vertex v, integer m and color c implying to the minimum Cost(rt(T'), 0, 0) value, the corresponding optimal \hat{T} provides the topological structure of T'' when restricted to the subtree rooted at v.

as they are produced. Moreover, we subtract the constant term RF(S', T') to express the cost as the *change* in RF distance.

We point out that the choice of \widehat{T} implying Cost(rt(T'), 0, 0) is exactly T'' by Lemmas 4.1 322 and 4.2. Furthermore, for any internal node v in T', and for the choice of m, c which imply 323 the optimal cost value of Cost(rt(T'), 0, 0) via the upcoming recurrence relation, the tree 324 $\widehat{T}(v)$ which admits Cost(v, m, c) is exactly equal to T''(v). In this sense, each \widehat{T} captures an 325 entire subtree of T''. Note that on a local scale, in any specific \overline{T} there may be a red or green 326 subtree regrafted outside of an extraneous subtree and outside of its canonical EF-R-RF(+)327 position. However, it can be concluded that either eventually these red and green subtrees 328 will be paired in extraneous subtrees for some later T, or the particular cost value does not 329 imply the optimal Cost(rt(T'), 0, 0). 330

The next lemma provides a recurrence relation that can compute each Cost(v, m, c)efficiently. In this recurrence relation, a subscript of L or R denotes the *left* or *right* child, respectively. For example, if a vertex v is an internal node in T then v_L is the left child of v, and if c is a color associated with vertex v then c_L is a color associated with vertex v_L . Note that the trees are unordered, so we use "left" and "right" here only to distinguish between the two children of an internal node.

▶ Lemma 4.4. Let $f(m_i, v_i, c_i)$ equal 2 when $m_i > 0$ and v_i is a match with color other than c_i , and 0 otherwise. Let $g_c(m_L, m_R, c_L, c_R)$ equal $2 \cdot \min\{m_L, m_R\}$ when $c_L \neq c_R$, and 0 when $c_L = c_R = c$. Then,

$$Cost(v, m, c) = \min_{m_L, m_R, c_L, c_R} \left\{ \begin{array}{c} Cost(v_L, m_L, c_L) + Cost(v_R, m_R, c_R) \\ + f(m_L, v_L, c_L) + f(m_R, v_R, c_R) - g_c(m_L, m_R, c_L, c_R) \end{array} \right\}$$

if v is an internal node of T', and Cost(v, m, c) = 0 if v is a leaf of T', where:

- 342 (a) $c, c_L, c_R \in \{0, 1\}$, and either $c_L \neq c_R$ or $c_L = c_R = c$,
- 343 **(b)** $0 \le m \le cMax(c, v),$
- (c) If $c_L \neq c_R$, then $m_i m_j = m$ for $i, j \in \{L, R\}, i \neq j$ satisfying $c_i = c_i$
- 345 (d) If $c_L = c_R = c$, then $m_L + m_R = m$

The functions f and g_c from Lemma 4.4 both track local changes in matched and 346 mismatched nodes. In particular, f tracks a local change between RF(S', T') and RF(S', T'')347 while g_c tracks a local change between RF(S', T'') and $RF(S^*, T^*)$. We now provide our 348 dynamic programming algorithm for computing the R-RF(+) distance between S and T. 349 **Algorithm** Compute-R-RF+(S,T)350 1: Compute the EF-R-RF(+) completions S' and T' of S and T. 351 2: for v in T' in postorder do 352 if v is a leaf then 3: 353 Set Cost(v, 0, 0) = Cost(v, 0, 1) = 0. 4: 35 if v is the root of a maximal red (0) or green (1) subtree then 5: 355 Set $Cost(v, 1, c_v) = 0$, where c_v is the color of v. 6: 356 7: else 357 8: for each color c and value $0 \le m \le cMax(c, v)$ do 358

9: Compute Cost(v, m, c) using the recurrence relation from Lemma 4.4

360 10: return RF(S', T') + Cost(rt(T'), 0, 0)

The algorithm above can be easily augmented to compute optimal completions by backtracking and determining the optimal values of m and c at each vertex of T' implying Cost(rt(T'), 0, 0). Using these optimal m and c values, we can determine when opposite colored subtrees converge and construct T''. From T'', we simply move each extraneous subtree $\{s, t\}$ into the canonical EF-R-RF(+) position for T(t) to build T^* and form the same extraneous subtrees in S' to build S^* .

Theorem 4.5. The RF(+) distance between two rooted binary trees S and T can be computed in $O(nk^2)$ time, where $n = |Le(S) \cup Le(T)|$ and k is the number of maximal red and green subtrees in S and T.

5 Extension to Unrooted Trees

³⁷¹ Our algorithm for the R-RF(+) problem can be easily adapted for the U-RF(+) problem. ³⁷² Specifically, the following algorithm computes the unrooted RF(+) distance between two ³⁷³ unrooted input trees S and T with at least one leaf in common.

- 374 Algorithm Compute-U-RF+(S, T)
- 1: Let l be any leaf from $Le(S) \cap Le(T)$. Produce two rooted trees \widehat{S} and \widehat{T} by rooting Sand T, respectively, on the edge which connects l to the rest of each tree.

2: Compute the RF(+) distance d between \widehat{S} and \widehat{T} using Algorithm Compute-R-RF+(S,T). 378 3: Return d

The correctness of this algorithm is easy to establish based on the well-understood association between rooted and unrooted RF distances [10, 4], and further technical details and proofs are therefore omitted. This yields the following two theorems.

Theorem 5.1. The U-RF(+) problem can be solved in $O(nk^2)$ time, where $n = |Le(S) \cup Le(T)|$ and k is the number of maximal red and green subtrees in the corresponding EF-U-RF(+) completion of S or T.

▶ Theorem 5.2. Let S^* and T^* represent optimal completions of unrooted trees S and T, respectively, under the U-RF(+) problem. Let S' and T' represent optimal completions of Sand T, respectively, under the EF-U-RF(+) problem. Then, $RF(S',T') \leq 2 \cdot RF(S^*,T^*)$.



Figure 6 Fraction of conflicting triples for different leaf-overlap ratios. The figure contains three plots, one for each dataset, which each show the fraction of triples of type-1, type-2, and type-3 for different ranges of leaf-overlap ratio, among all triples of trees within the same leaf-overlap ratio range in that dataset. The dotted line represents the total number of conflicting triples (i.e., all triples of type 1, 2 or 3). *x*-axis labels denote the center of each interval of size 0.1. Each leaf-overlap ratio range is a closed interval and *includes* the boundary, e.g., *x*-axis label 0.15 represents the range [0.1 - 0.2].

6 Experimental Evaluation

We implemented our exact algorithm and performed experiments to assess the impact of using RF(+) distance instead of RF(-) distance on inferences related to tree similarity. We also conducted experiments to assess how well the linear-time algorithm for computing EF-RF(+) distances approximates RF(+) distances in practice. All our experiments were performed using real biological phylogenetic tree datasets on marsupials [8] (158 trees), legumes [33] (22 trees), and placental mammals [7] (726 trees).

Experiment 1: Conflicts between RF(+) and RF(-). Given two trees S and T, let $RF^+(S,T)$ and $RF^-(S,T)$, respectively, denote the RF(+) and RF(-) distances between them. We used the above datasets to measure the number of times that for any "base" tree S, there is a tree T_1 which is closer to S than T_2 under one of RF(+) or RF(-) but not closer under the other distance measure. This motivates the following definitions to describe each possible case of a change in order.

⁴⁰¹ **Type-1 Triples:** Triple (S, T_1, T_2) is Type-1 if $RF^-(S, T_1) < RF^-(S, T_2)$ but $RF^+(S, T_1) >$ ⁴⁰² $RF^+(S, T_2)$, or $RF^-(S, T_2) < RF^-(S, T_1)$ but $RF^+(S, T_2) > RF^+(S, T_1)$. A Type-1 triple ⁴⁰³ indicates when the ordering of T_1 and T_2 by distance from S strictly changes as the distance ⁴⁰⁴ function changes between RF(-) and RF(+).

⁴⁰⁵ **Type-2 Triples:** Triple (S, T_1, T_2) is Type-2 if $RF^-(S, T_1) = RF^-(S, T_2)$ but $RF^+(S, T_1) \neq$ ⁴⁰⁶ $RF^+(S, T_2)$. A Type-2 triple indicates when T_1 and T_2 have equal distance to S under RF(-) ⁴⁰⁷ but not under RF(+).

⁴⁰⁸ **Type-3 Triples:** Triple (S, T_1, T_2) is Type-3 if $RF^-(S, T_1) \neq RF^-(S, T_2)$ but $RF^+(S, T_1) = RF^+(S, T_2)$. A Type-3 triple indicates when T_1 and T_2 have equal distance to S under RF(+)⁴¹⁰ but not under RF(-).

⁴¹¹ Observe that the magnitude of difference between RF(+) and RF(-) distances depends ⁴¹² on the level of overlap between the trees being compared. To account for this effect, we ⁴¹³ define the *leaf-overlap ratio* of a pair of trees (S, T) to be the following ratio: $|Le(S) \cap Le(T)|$

K. Yao and M. S. Bansal

divided by min{|Le(S)|, |Le(T)|}, and the leaf-overlap ratio of a triple of trees S, T_1 , and T_2 to be the minimum pairwise leaf-overlap ratio between (S, T_1) and (S, T_2) .

We performed this experiment for each subset of three trees from each dataset, and Figure 416 6 shows its results. As the figure shows, the proportion of conflicting triples (type-1, 2, or 3) 417 tends to increase as the triple leaf-overlap ratio increases. In particular, at least 10% of all 418 triples show a conflict (either of type-1, 2, or 3) when the leaf-overlap ratio is 0.7 or greater. 419 Even for leaf-overlap ratio as small as 0.4, we find that 5% of all triples show a conflict. 420 This demonstrates that RF(+) and RF(-) frequently differ starkly in their assessments of 421 relative similarities between trees. Observe that the results on the Legumes dataset are vastly 422 different from the results on the other two datasets. This is mainly because the Legumes 423 dataset consists of only 22 trees, which is significantly smaller than the 158 tree and 726 tree 424 datasets. For instance, the number of triples within each leaf overlap ratio range (interval 425 size 0.1) is between 8,214,518 and 50,815,687 for the placental mammals dataset, between 426 3,287 and 1,652,701 for the Marsupials dataset, but only 6, 16, 5, and 0, respectively, for the 427 Legumes dataset for leaf overlap ratio ranges [0.5 - 0.6], [0.6 - 0.7], [0.7 - 0.8], and [0.8 - 0.9].428



Figure 7 Difference between sets of closest trees under RF(+) and RF(-). Plots in the left column show the number of query trees where the set of closest trees with a minimum leaf-overlap ratio of 0.7 differ under RF(+) and RF(-) distances for each of the three biological data sets. Plots in the right column show the number of query trees where the set of closest 10% of trees with a minimum leaf-overlap ratio of 0.5 differ under RF(+) and RF(-) distances. Results are presented for varying levels of difference between the sets (labels on the *x*-axes). The sizes of the datasets, in order from top to bottom, are 158 trees, 22 trees and 726 trees. Each tree in each of these datasets was used as a query tree for this analysis.

Experiment 2: Impact on phylogenetic database search and clustering. Next, we assessed the potential impact of using RF(+) distance on applications related to phylogenetic database search and clustering. Specifically, we assessed how, for each "query" tree in each dataset, the sets of the "closest" trees to it differed under RF(+) and RF(-). Specifically, we measured how the sets of (i) the most similar trees and (ii) the most similar 10% of trees (i.e., top 10% closest matches) differed when using RF(+) and RF(-) distances. To avoid any ambiguity in defining these sets, we include all trees with equal distance, even if that results

in sets of different sizes under RF(+) and RF(-).

For our comparison of the most similar trees, we found that the sets of closest trees 437 under RF(+) and RF(-) all had a distance of 0 to the query tree and were identical, for 438 all choices of the query tree in all datasets. To perform a more meaningful comparison, we 439 therefore required a minimum leaf-overlap ratio of 0.7, i.e., only those trees with a minimum 440 leaf-overlap ratio of 0.7 with the query tree could be compared with the query tree. Likewise, 441 for our comparison of the most similar 10% of trees, we found that the sets of closest 10%442 of trees were generally identical under RF(+) and RF(-) if no minimum leaf-overlap ratio 443 was imposed. We therefore imposed a minimum leaf-overlap ratio of 0.5 for the analysis, 444 which was the smallest ratio for which a non-negligible fraction of query trees returned 445 differing sets under RF(+) and RF(-). Figure 7 shows the results of both these analyses. We 446 find that there are several query trees in each dataset for which there is a large difference 447 (normalised symmetric difference greater than, say, 0.4) between their sets of closest trees 448 under RF(+) and RF(-). For the sets of closest 10% of trees, we find that over 25% of trees 449 in the marsupials dataset, 18% of trees in the legumes dataset, and almost 15% of trees in 450 the placental mammals dataset return different sets of closest 10% of trees under RF(+)451 and RF(-) distances. These results demonstrate how using RF(+) distance can substantially 452 impact phylogenetic database search and phylogenetic tree clustering, especially when the 453 trees under consideration have a sufficiently large overlap in their leaf sets. 454

Experiment 3: Comparison of EF-RF(+) and RF(+). Finally, we used simulated 455 and real datasets to compare the distances inferred under EF-RF(+) and RF(+), and to 456 study the runtime and scalability of our implementation. For our analysis with simulated 457 data, we generated two datasets of random trees using the birth-death model implemented 458 in SaGePhy [21] (specific parameter values: height of tree = 1.0, birth rate = 5.0 and 459 death rate = 0.05). The first simulated dataset consisted of 100 randomly generated trees, 460 each with between 200 and 300 leaves. The second simulated dataset also consisted of 100 461 randomly generated trees, but each with between 900 and 1000 leaves. The average leaf-set 462 sizes for these two datasets were 244.95 and 941.14, respectively, and the average pairwise 463 leaf-overlap ratio for both datasets was approximately 0.5. For each pair of trees in each 464 dataset, we measured how close the EF-RF(+) distance is to the RF(+) distance for that 465 pair. Figure 8 plots the distribution of the ratio of RF(+) distance to EF-RF(+) distance for 466 the two datasets. As that figure shows, the ratio of RF(+) distance to EF-RF(+) distance is 467 approximately 0.92, on average, and roughly follows a Gaussian distribution. 468



Figure 8 Comparison of EF-RF(+) and RF(+) distances on simulated trees. The two plots show the distribution of the ratio of RF(+) distance to EF-RF(+) distance for the two simulated datasets consisting of randomly generated birth-death trees. Each dataset contains 100 trees and results are shown for all $\binom{100}{2}$ pairs of trees in each dataset.

For the three biological datasets, we found that the ratio of RF(+) distance to EF-RF(+)distance was equal to one for an overwhelmingly large proportion of pairs of trees within all

three datasets. Specifically, for the marsupials, legumes, and placental mammals datasets, 471 the average ratios of RF(+) distance to EF-RF(+) distance were 0.998, 0.993, and 0.995, 472 respectively. In fact, 99.07%, 93.81%, and 96.82% of the pairs in these datasets, respectively, 473 had identical EF-RF(+) and RF(+) distances. Even when the trees being compared were 474 restricted to have at least 0.4 leaf-overlap ratio, 95.97%, 78.79%, and 95.59% of the pairs in 475 marsupials, legumes, and placental mammals datasets, respectively, had identical EF-RF(+) 476 and RF(+) distances. This discrepancy between results for simulated data and real data is 477 not surprising since we expect any pair of randomly generated trees to have smaller maximal 478 red and green subtrees and greater RF(-) distance, presenting more opportunities to improve 479 the distance by creating extraneous clades. Together, these results on simulated and real 480 datasets show that EF-RF(+) distance, which is linear-time computable, is generally very 481 close to RF(+) distance in practice. 482

Runtime comparison. We also measured the runtimes of the two algorithms and found that, on average, computing EF-RF(+) distances took 0.06 seconds for the first simulated dataset and 0.25 seconds for the second simulated dataset. Corresponding average runtimes for computing RF(+) distances were 0.17 seconds and 1.04 seconds, respectively. All timed experiments were run on a single core of a 2.1 GHz Intel Xeon processor.

488 7 Conclusion

Completion based comparison of incomplete phylogenetic trees is an emerging, promising 489 approach for tree comparison. In this work, we developed the first polynomial-time exact 490 algorithm for the RF(+) problem. We also established a linear-time 2-approximation 491 algorithm for the problem. These algorithms allow for more principled comparison of 492 incomplete phylogenetic trees than was hitherto possible, and our experimental analysis 493 shows that RF(+) distance can lead to very different inferences regarding phylogenetic 494 similarity compared to traditional RF distance. Moreover, our results suggest that the linear-495 time 2-approximation algorithm for the RF(+) problem almost always computes optimal or 496 near-optimal RF(+) distances in practice. 497

In addition to their utility for improved tree comparison and clustering, our solutions for 498 the RF(+) problem also have implications for phylogenomics. Many modern phylogenomics 499 methods for reconstructing evolutionary histories and understanding genome-scale patterns 500 of evolution are designed to work with complete phylogenies from genomic loci across 501 the genomes of the considered species [5, 26, 27, 20, 12], and loci that yield incomplete 502 phylogenies are often discarded, resulting in only a fraction of the available genomic sequence 503 information being used for the phylogenomic analysis. Thus, problems related to optimal 504 completion of incomplete phylogenies (i.e., imputation of complete phylogenies) arise naturally 505 in phylogenomics. Our algorithms for the RF(+) problem may provide a principled way to 506 impute such complete phylogenies. 507

The current work is restricted to comparison of binary trees under the Robinson-Foulds 508 metric, and it can be extended in many useful ways. A possible next step could include 509 consideration of non-binary trees in computing distances between incomplete trees. Fu-510 ture work could also entail development of similar completion based methods under other 511 distance/similarity measures such as triplet/quartet distances [14, 17], nearest neighbor 512 interchange (NNI) and subtree prune and regraft (SPR) distances [31, 18, 34], and nodal 513 distance [9]. Furthermore, the idea of computing optimal completions could be extended 514 to multi-labeled trees, which arise frequently in practice due to evolutionary events such as 515 gene duplication. 516

 Wasiu A. Akanni, Mark Wilkinson, Christopher J. Creevey, Peter G. Foster, and Davide Pisa Implementing and testing bayesian and maximum-likelihood supertree methods in phylogene- ics. Royal Society Open Science, 2(8), 2015. URL: http://rsos.royalsocietypublishir org/content/2/8/140436, doi:10.1098/rsos.140436. Amihood Amir and Dmitry Keselman. Maximum agreement subtree in a set of evolutiona trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656–1669, 199 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209–226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robins son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	_
 Implementing and testing bayesian and maximum-likelihood supertree methods in phylogenetics. Royal Society Open Science, 2(8), 2015. URL: http://rsos.royalsocietypublishir org/content/2/8/140436, doi:10.1098/rsos.140436. Amihood Amir and Dmitry Keselman. Maximum agreement subtree in a set of evolutional trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656-1669, 199 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	ni.
 ics. Royal Society Open Science, 2(8), 2015. URL: http://rsos.royalsocietypublishir org/content/2/8/140436, doi:10.1098/rsos.140436. Amihood Amir and Dmitry Keselman. Maximum agreement subtree in a set of evolutiona trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656-1669, 199 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robins son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	et-
 ⁵²¹ org/content/2/8/140436, doi:10.1098/rsos.140436. ⁵²² 2 Amihood Amir and Dmitry Keselman. Maximum agreement subtree in a set of evolutiona trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656-1669, 199. doi:10.1137/S0097539794269461. ⁵²⁵ 3 Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. ⁵²⁹ 4 Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	ıg.
 Amihood Amir and Dmitry Keselman. Maximum agreement subtree in a set of evolutiona trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656-1669, 198 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	0
 trees: Metrics and efficient algorithms. SIAM Journal on Computing, 26(6):1656-1669, 199 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	rv
 doi:10.1137/S0097539794269461. Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In <i>Comparative Genomics - 16th International Conference</i> <i>RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings</i>, pag 209–226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robin son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020.)7.
 Mukul S. Bansal. Linear-time algorithms for some phylogenetic tree completion problem under robinson-foulds distance. In <i>Comparative Genomics - 16th International Conference</i> <i>RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings</i>, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robin son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	
 ⁵²⁵ Intuition Difference and the algorithms for some phylogenetic free completion problem under robinson-foulds distance. In Comparative Genomics - 16th International Conference RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209–226, 2018. doi:10.1007/978-3-030-00834-5_12. ⁵²⁹ 4 Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robinson-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	ns
 RECOMB-CG 2018, Magog-Orford, QC, Canada, October 9-12, 2018, Proceedings, pag 209-226, 2018. doi:10.1007/978-3-030-00834-5_12. Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robi son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020. 	.15 C.C.
 ¹¹ ¹¹ ¹¹ ¹² ¹² ¹² ¹² ¹²	es.
4 Mukul S. Bansal. Linear-time algorithms for phylogenetic tree completion under robi son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020.	
son-foulds distance. Algorithms for Molecular Biology, 15:6, 2020.	n-
5 Mukul S Bansal Guy Banay Timothy I Harlow I Peter Cogarten and Ron Sham	ir
Systematic inference of highways of horizontal gene transfer in prokarvotes <i>Bioinformati</i>	
⁵³³ 29(5):571–579, 2013.	,
6 Mukul S. Bansal, I. Cordon Burlaigh, Oliver Fulanetain, and David Fernández Baca. Robinse	m
foulds supertrees Algorithms for Molecular Biology 5(1):18 Feb 2010	11-
7 Bobin Bock Olaf Bininda Emonds Marcal Cardillo, Fu Guo Liu, and Andy Purvis	Δ
higher level MRP supertree of placental mammals <i>BMC Eval Biol</i> 6(1):03 2006 de	A i·
$10 \ 1186/1471-2148-6-93$	1.
 Marcel Cardille Olef P. D. Dininde Emends Elizabeth Boelies and Andr. Durris A specie 	20
⁵³⁹ o Marcel Cardino, Olar R. P. Binnida-Emolids, Enzabeth Boakes, and Andy Purvis. A specie	38-
⁵⁴⁰ level phylogenetic superfield of marsuphals. <i>Journal of Zoology</i> , 204.11–51, 2004.	
⁵⁴¹ 9 Gabriel Cardona, Merce Liabres, Francesc Rossello, and Gabriel Vallente. Nodal distances I	or T.
542 rooted phylogenetic trees. Journal of Mathematical Biology, 61(2):253–276, Aug 2010. UR	т:
⁵⁴⁴ 10 Ruchi Chaudhary, J Gordon Burleign, and David Fernandez-Baca. Fast local search F	or
⁵⁴⁵ unrooted robinson-rounds supertrees. <i>IEEE/ACM Transactions on Computational Biology an</i> <i>Displaymenting (TCPP)</i> , 0(4):1004, 1012, 2012	ia
 Diomjormanics (TCDD), 9(4):1004–1013, 2012. Diomonia Chen, L.Conder, Durchick, Michael C. Densel, and Densid Farmán des Dava. Dhalafin des 	
547 11 Dunong Chen, J Gordon Burleign, Mukul S Bansal, and David Fernandez-Baca. Phylonnae	er:
an intelligent search engine for phylogenetic tree databases. <i>DNC Evolutionary Diology</i> , 8(1):	<i>i</i> 0,
12 Carel Christener Evin I/ Meller Derniel Verlageneti and Tende Werner Ortinal Co	
550 12 Sarah Christensen, Erin K. Molloy, Pranjal Vachaspati, and Tandy Warnow. Optimal Col	n-
⁵⁵¹ pletion of incomplete Gene Trees in Polynomial Time Using OUTAL. In Russell Schwar	τz
(WA PL 2017) volume SS of Leibniz International Proceedings in Information (IIPIss) page	cs
27:1-27:14 Dagstuhl Cormany 2017 Schloss Dagstuhl-Loibniz Zontrum fuor Informatik	es
 12 James A. Cotton Mark Willingen and Mile Steel Majority rule supertroop. Customed 	
⁵⁵⁵ 15 James A. Cotton, Mark Winkinson, and Mike Steel. Majority-rule superfrees. Systemat	<i>ic</i>
$_{556}$ $_{556}$ $_{560}$ $_$	1:
14 Develop E Critchlery Dennig K Deerl Churlin Gion and Daniel Esith. The triples distant	
for rooted bifurenting phylogenetic trees. <i>Custometric Diclose</i> , 45(2),222–224, 1006 UD	се т.
⁵⁵⁹ for rooted billinearing phylogenetic trees. Systematic Biology, 45(5):525–554, 1990. UK	L:
15 Damian M de Vienne Tetione Gineud and Olivier G Martin A commune i de Catati	
topological similarity between troop <i>Bioinformatica</i> 22(22),2110, 2124, 2007 UDL (1)++	ng.
$\frac{1000}{1000}$ topological similarity between trees. <i>Diotiljorniutics</i> , 25(25).5119–5124, 2007. URL: fift	ŀ٠
16 Jianwarg Dang and David Farmander Dasa Dramatics of maiority rule current and current of the second seco	tic
564 IU Jiamong Dong and David remandez-Daca. Properties of majority-rule supertrees. Systemat Biology 58(3):360-367-2000 UDL: thttp://dx.doi.org/10.1002/graphic/cm-030_daid	<i>.ic</i>
1093/sysbio/syb032.	

K. Yao and M.S. Bansal

567	17	George F. Estabrook, F. R. McMorris, and Christopher A. Meacham. Comparison of undirected
568		phylogenetic trees based on subtrees of four evolutionary units. Systematic Zoology, 34(2):193–
569		200, 1985. URL: http://www.jstor.org/stable/2413326.
570	18	J. Felsenstein. Inferring Phylogenies. Sinauer Assoc., Sunderland, Mass, 2003.
571	19	C. R. Finden and A. D. Gordon. Obtaining common pruned trees. <i>Journal of Classification</i> .
572		2(1):255-276, Dec 1985. doi:10.1007/BF01908078.
573	20	Kevin Gori, Tomasz Suchan, Nadir Alvarez, Nick Goldman, and Christophe Dessimoz. Cluster-
574		ing genes of common evolutionary history. <i>Molecular Biology and Evolution</i> , 33(6):1590–1605.
575		2016. URL: http://dx.doi.org/10.1093/molbev/msw038, doi:10.1093/molbev/msw038.
576	21	Soumva Kundu and Mukul S Bansal. Sagephy: An improved phylogenetic simulation framework
577		for gene and subgene evolution. <i>Bioinformatics</i> , 35(18):3496–3498, 2019.
578	22	Anne Kupczok. Split-based computation of majority-rule supertrees. BMC Evolutionary
579		Biology, 11(1):205, Jul 2011. URL: https://doi.org/10.1186/1471-2148-11-205.
580	23	Anne Kupczok, Arndt Von Haeseler, and Steffen Klaere. An exact algorithm for the geodesic
581		distance between phylogenetic trees. Journal of Computational Biology, 15(6):577–591, 2008.
582	24	Harris T. Lin, J. Gordon Burleigh, and Oliver Eulenstein. Triplet supertree heuristics for the
583		tree of life. BMC Bioinformatics, 10(1):S8, Jan 2009. doi:10.1186/1471-2105-10-S1-S8.
584	25	Michelle M. McMahon, Akshay Deepak, David FernÄindez-Baca, Darren Boss, and Michael J.
585		Sanderson. Stbase: One million species trees for comparative biology. <i>PLOS ONE</i> , 10(2):1–17,
586		02 2015. doi:10.1371/journal.pone.0117987.
587	26	S. Mirarab, R. Reaz, Md. S. Bayzid, T. Zimmermann, M. S. Swenson, and T. Warnow.
588		ASTRAL: genome-scale coalescent-based species tree estimation. <i>Bioinformatics</i> , 30(17):i541–
589		i548, 2014. URL: http://dx.doi.org/10.1093/bioinformatics/btu462, doi:10.1093/
590		bioinformatics/btu462.
591	27	Siavash Mirarab, Md. Shamsuzzoha Bayzid, Bastien Boussau, and Tandy Warnow. Statistical
592		binning enables an accurate coalescent-based estimation of the avian tree. Science, 346(6215),
593		2014. URL: http://science.sciencemag.org/content/346/6215/1250463, doi:10.1126/
594		science.1250463.
595	28	William H Piel, MJ Donoghue, MJ Sanderson, and LUT Netherlands. Treebase: a database of
596		phylogenetic information. In Proceedings of the 2nd International Workshop of Species 2000,
597		2000.
598	29	D.F. Robinson and L.R. Foulds. Comparison of phylogenetic trees. Mathematical Bios-
599		ciences, 53(1):131 - 147, 1981. URL: http://www.sciencedirect.com/science/article/
600		pii/0025556481900432, doi:https://doi.org/10.1016/0025-5564(81)90043-2.
601	30	Jason TL Wang, Huiyuan Shan, Dennis Shasha, and William H Piel. Fast structural search in
602		phylogenetic databases. Evolutionary Bioinformatics, 2005(1):0–0, 2007.
603	31	M.S. Waterman and T.F. Smith. On the similarity of dendrograms. Journal of Theoretical
604		Biology, 73(4):789 - 800, 1978. URL: http://www.sciencedirect.com/science/article/
605		pii/0022519378901376, doi:https://doi.org/10.1016/0022-5193(78)90137-6.
606	32	Christopher Whidden, Norbert Zeh, and Robert G. Beiko. Supertrees based on the subtree
607		prune-and-regraft distance. Systematic Biology, 63(4):566-581, 2014. URL: +http://dx.doi.
608		org/10.1093/sysbio/syu023, doi:10.1093/sysbio/syu023.
609	33	M.F. Wojciechowski, M.J. Sanderson, K.P. Steele, and A. Liston. Molecular phylogeny of
610		the "Temperate Herbaceous Tribes" of Papilionoid legumes: a supertree approach. In P.S.
611		Herendeen and A. Bruneau, editors, Advances in Legume Systematics, volume 9, pages 277–298.
612		Royal Botanic Gardens, Kew, 2000.
613	34	Yufeng Wu. A practical method for exact computation of subtree prune and regraft distance.
614		<i>Bioinformatics</i> , 25(2):190-196, 2009. URL: +http://dx.doi.org/10.1093/bioinformatics/
615		btn606, doi:10.1093/bioinformatics/btn606.
616	35	Ruriko Yoshida, Kenji Fukumizu, and Chrysafis Vogiatzis. Multilocus phylogenetic ana-
617		lysis with gene tree clustering. Annals of Operations Research, Mar 2017. doi:10.1007/
618		s10479-017-2456-9.

619 Appendix

Proof of Theorem 3.1. Let S^* and T^* be arbitrarily chosen optimal completions of S and T under R-RF(+). We will modify S^* and T^* to be of the desired form. To do so, we first show that any maximal red subtree in S and any maximal green subtree of T can be made subtrees of S^* and T^* without increasing the RF distance between them (condition 1). Suppose there exist two maximal matched red subtrees R_1 and R_2 of S^* and T^* which neighbor each other in the original tree S. Let r_1 and r_2 be the roots of R_1 and R_2 . 1. Suppose both $C_{T^*}(pa_{T^*}(r_1)) \setminus C_{T^*}(r_1)$ and $C_{T^*}(pa_{T^*}(r_2)) \setminus C_{T^*}(r_2)$ contain some non-

green leaves. Observe that every matched clade in T^* containing $C_{T^*}(r_1) \cup C_{T^*}(r_2)$ must also contain $C_{T^*}(lca_{T^*}(r_1, r_2))$ because R_1 and R_2 neighbor each other in S by assumption. Therefore, we can regraft R_2 to neighbor R_1 in T^* without increasing the RF distance between S^* and T^* . Moreover, if there are any green subtrees inserted along the path from R_1 to R_2 in S^* , then they can be regrafted along the parent edge of $lca_{S^*}(r_1, r_2)$ without increasing the Robinson-Foulds distance.

2. Suppose, without loss of generality, that $C_{T^*}(pa_{T^*}(r_2)) \setminus C_{T^*}(r_2)$ contains only green 633 leaves. That is, suppose R_2 is contained in an extraneous subtree, whose root could 634 be a match without ancestoring R_1 . First, regraft R_2 in T^* to neighbor R_1 . Then, 635 regraft all green subtrees from the path in S^* connecting R_2 and R_1 to the parent edge 636 of $lca_{S^*}(r_1, r_2)$, preserving the topological structure of the green leaves. This does not 637 increase the RF distance between S^* and T^* . Notice that any originally matched clades 638 containing $Le(R_2)$ are mismatched. However, preserving the topological structure of the 639 green leaves from any matched clades containing $Le(R_2)$ also retains the same number of 640 matches except for one representing the smallest match containing R_2 . This is because 641 the only subtree removed (in both S^* and T^*) from these matched extraneous subtrees 642 is R_2 . Furthermore, the matched clade $Le(R_1) \cup Le(R_2)$ is formed in both S^* and T^* , 643 which counteracts this lost match. 644

If this is done iteratively for all such R_1 and R_2 , then we conclude that there exist optimal completions S^* and T^* where every maximal red subtree in S is also a subtree of S^* and T^* . The same argument applies for maximal green subtrees.

Now we will show that S^* and T^* can be modified to only contain extraneous subtrees of 648 the form $\{s, t\}$ without increasing the RF distance (condition 2). We will simultaneously 649 show that an extraneous subtree $\{s, t\}$ is a subtree of S^* if and only if it is a subtree of T^* 650 by construction (condition 3). Observe that if $Le(U) \cap Le(V) \cap Le(S) \neq \emptyset$ for two maximal 651 extraneous subtrees U and V of S^* and T^* respectively, then $Le(U) \cap Le(V) \cap Le(S) \subseteq Le(R)$ 652 for a single maximal red subtree R of S. Likewise if $Le(U) \cap Le(V) \cap Le(T) \neq \emptyset$, then 653 $Le(U) \cap Le(V) \cap Le(T) \subseteq Le(Y)$ for a single maximal green subtree Y of T. Therefore, every 654 maximal extraneous subtree in S^* or T^* satisfies one of the following two cases. 655

1. Without loss of generality, let U be a maximal extraneous subtree of S^* rooted at usuch that for *every* maximal extraneous subtree V of T^* , $Le(U) \cap Le(V) \cap Le(S) = \emptyset$ or $Le(U) \cap Le(V) \cap Le(T) = \emptyset$. Then, every *extraneous* clade contained in Le(U) must be a mismatch. Hence, every maximal green subtree of U can be regrafted along the parent edge of $pa_{S^*}(u)$ without increasing the Robinson-Foulds distance from T^* . This results in destroying *all* extraneous subtrees contained in U because $pa_{S^*}(u)$ is an ancestor of a maximal extraneous subtree and therefore possesses uncolored descendants.

⁶⁶³ 2. Let U and V be maximal extraneous subtree of S^* and T^* , rooted at u and v respectively, ⁶⁶⁴ satisfying $Le(U) \cap Le(V) \cap Le(S) \neq \emptyset$ and $Le(U) \cap Le(V) \cap Le(T) \neq \emptyset$. Then every ⁶⁶⁵ matched extraneous clade contained in Le(U) and Le(V) must contain elements of

 $Le(U) \cap Le(V) \cap Le(S)$ and $Le(U) \cap Le(V) \cap Le(T)$. Every maximal green subtree of U 666 with no leaves in $Le(U) \cap Le(V) \cap Le(T)$ can be regrafted along the parent edge of u 667 without increasing the RF distance. Likewise, every maximal red subtree of V with no 668 leaves in $Le(U) \cap Le(V) \cap Le(S)$ can be regrafted along the parent edge of v without 669 increasing the RF distance. Moreover, as described before, $Le(U) \cap Le(V) \cap Le(S) \subseteq Le(R)$ 670 and $Le(U) \cap Le(V) \cap Le(T) \subseteq Le(Y)$ for a single maximal red subtree R of S and a single 671 maximal green subtree Y of T. Hence, we are only left with the extraneous subtree 672 $\{rt_{S^*}(R), rt_{S^*}(Y)\}$ in S^* and $\{rt_{T^*}(R), rt_{T^*}(Y)\}$ in T^* . 673 Once every maximal extraneous subtree in S^* and T^* is handled according to the appropriate 674 case above, we are left with two optimal completions S^* and T^* of the desired form. 675 **Proof of Lemma 3.2.** Case 1: In this case, both $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ are matched 676 nodes. Here, we must have $Le(S^*(pa_{S^*}\{s,t\})) = Le(T^*(pa_{T^*}\{s,t\}))$. This holds because 677 $C_{S^*}(pa_{S^*}\{s,t\})$ and $C_{T^*}(pa_{T^*}\{s,t\})$ are both matches, and the smallest proper super-678 sets of $C_{T^*}(s) \cup C_{T^*}(t)$ in S^* and T^* respectively. By definition, the decomposition re-679 places the matched clades $C_{S^*}(s) \cup C_{S^*}(t)$ and $C_{T^*}(s) \cup C_{T^*}(t)$ with $C_{S^*}(pa_{S^*}\{s,t\}) \setminus$ 680 $C_{S^*}(t)$ and $C_{T^*}(pa_{T^*}\{s,t\}) \setminus C_{T^*}(t)$ in S^* and T^* , respectively. Since $Le(S^*(pa_{S^*}\{s,t\})) =$ 681 $Le(T^*(pa_{T^*}\{s,t\}))$, we conclude that $C_{S^*}(pa_{S^*}\{s,t\}) \setminus C_{S^*}(t)$ and $C_{T^*}(pa_{T^*}\{s,t\}) \setminus C_{T^*}(t)$ 682 are then matches in the resulting trees S' and T'. 683 Case 2: We now consider the case when exactly one of the nodes $pa_{S^*}\{s,t\}$ and $pa_{T^*}\{s,t\}$ is 684 a matched node. Without loss of generality, suppose $pa_{S^*}\{s,t\}$ is a match and $pa_{T^*}\{s,t\}$ 685 is not a match. For convenience, let x denote $pa_{S^*}\{s,t\}$, y denote $pa_{T^*}\{s,t\}$, and let u be 686 the element of $V(T^*)$ such that $C_{S^*}(x) = C_{T^*}(u)$. Then, observe that $C_{S^*}(x) \supset C_{T^*}(y)$, i.e., 687 y < u in T^{*}. Moreover, every node v along the path from y to u in T^{*} must be a mismatch 688 since $C_{T^*}(t) \subset C_{T^*}(v)$ and $C_{S^*}(t) \cap C_{S^*}(sib_{S^*}\{s,t\}) = \emptyset$ but $C_{T^*}(v) \cap C_{S^*}(sib_{S^*}\{s,t\}) \neq \emptyset$ 689 for arbitrary choice of v. Now, applying the decomposition of extraneous subtree $\{s,t\}$ to 690 S^* and T^* yields the modified trees S' and T'. Observe that this modification changes 691 exactly the $\{s, t\}$ clade, and all clades along the path from y to u in T^* . In S', the new clade 692 formed at the subtree rooted at $pa_{S'}(t)$ must be a matched node since $C_{S'}(pa_{S'}(t)) = C_{T'}(u)$. 693

- Moreover, in T', all clades $C_{T'}(v)$ along the path from y to u remain mismatches except for $C_{T'}(u)$ because it still holds that $C_{T'}(t) \subset C_{T'}(v)$ and $C_{S'}(t) \cap C_{S'}(sib_{S'}\{s,t\}) = \emptyset$ but $C_{T'}(v) \cap C_{S'}(sib_{S'}\{s,t\}) \neq \emptyset$ for arbitrary choice of v along the path. Thus, after the decomposition, the number of matched clades in S' (w.r.t. T') remains the same as the number of matched clades in S^* (w.r.t. T^*).
- *Case 3:* If neither $pa_{S^*}\{s,t\}$ nor $pa_{T^*}\{s,t\}$ is a matched node, then, following the same argument as in Case 1, S' will have one less matched node (w.r.t. T') than S^* (w.r.t. T^*). Namely, the clades $C_{S^*}(pa_{S^*}\{s,t\}) \setminus C_{S^*}(t)$ and $C_{T^*}(pa_{T^*}\{s,t\}) \setminus C_{T^*}(t)$ are mismatched clades in S' and T' respectively. Consequently, T' will have one less matched node as well. Thus, $RF(S', T') = RF(S^*, T^*) + 2$.

Proof of Theorem 3.3. Let $d = \frac{1}{2}RF(S^*, T^*)$ and let e be the number of extraneous clades in S^* . Then, we have that d is also the number of mismatches in S^* , or equivalently in T^* . Observe that at most d of the e extraneous clades have mismatched parent nodes in both trees. Thus, by Lemma 3.2, decomposing all e extraneous clades will *increase* the RF distance by at most $2d = RF(S^*, T^*)$. Therefore, the decomposed extraneous clade free completion will have an RF distance of at most $2 \cdot RF(S^*, T^*)$.

Proof of Lemma 4.1. Consider arbitrary canonical R-RF(+) completions S^* and T^* . We will show that any grafted subtree in S^* and T^* that is not in its canonical EF-R-RF(+) position or in an extraneous subtree can be regrafted into its canonical EF-R-RF(+) position without increasing the RF distance. Without loss of generality, suppose there exists a maximal red subtree R, with r denoting rt(R), in T^* such that R is neither in its canonical EF-R-RF(+) position nor in an extraneous subtree. Let u represent the canonical EF-R-RF(+) position of subtree R in completion T^* . Thus, $u \neq pa_{T^*}(r)$. Then, we have two possible cases: either $pa_{T^*}(r)$ is an ancestor of u or not $(pa_{T^*}(r) > u$ or $pa_{T^*}(r) \neq u)$.

1. Suppose $pa_{T^*}(r) > u$. We will prove that $pa_{T^*}(r)$ can be regrafted in position u without 718 increasing the RF distance. Since $pa_{T^*}(r) > u$, for any arbitrary node c on the path 719 from $pa_{T^*}(r)$ to u, there exists a subtree C of $T^*(c)$ rooted at one of the children of 720 c (the subtree not containing u) satisfying $pa_{T^*}(r) > c = lca_{T^*}(u, Le(C)) > u$ and 721 $pa_{S^*}(r) < lca_{S^*}(r, Le(C))$. Since $pa_{T^*}(r) > lca_{T^*}(u, Le(C)) > u$, we have that $pa_{T^*}(r) > lca_{T^*}(r) > lca_{T^*}(r)$ 722 $lca_{T^*}(Le(C), a) > a$ for all leaves $a \in Le(S) \cap Le(T)$ such that $a < pa_{S^*}(r)$. Since for 723 each such a, we have that $a < pa_{S^*}(r) < lca_{S^*}(a, Le(C))$ and $a < lca_{T^*}(a, Le(C)) = c < ca_{S^*}(a, Le(C))$ 724 $pa_{T^*}(r)$, every match containing Le(C) must also contain Le(R). In particular, c is not a 725 match. This is true for every node c along the path from $pa_{T^*}(r)$ to u. We can therefore 726 regraft R at position u without increasing the RF distance because every node along the 727 path from $pa_{T^*}(r)$ to u is already mismatched. 728

2. Now suppose $pa_{T^*}(r) \neq u$. We will prove that R can be regrafted along the parent 729 edge of $lca_{T^*}(pa_{T^*}(r), u)$ (equivalent position to u if u is an ancestor of $pa_{T^*}(r)$) without 730 increasing the RF distance. This will then reduce the case where $pa_{T^*}(r)$ is not an 731 ancestor of u to the previous case where $pa_{T^*}(r)$ is an ancestor of u. If $pa_{T^*}(r)$ is not an 732 ancestor of u, then there exist some $a_1, \ldots, a_k \in Le(S) \cap Le(T)$ such that $pa_{S^*}(r) > a_i$ 733 and $lca_{T^*}(pa_{T^*}(r), a_i) > pa_{T^*}(r)$ for all values of i. Therefore, $pa_{T^*}(r)$ is not a match, 734 as well as every node on the same path up to the node $lca_{T^*}(pa_{T^*}(r), a_1, \ldots, a_k)$ which 735 contains every a_i in its clade $C_{T^*}(pa_{T^*}(pa_{T^*}(r), a_1, \ldots, a_k))$. Then, regrafting R at the 736 parent edge of $lca_{T^*}(a_1,\ldots,a_k,pa_{T^*}(r)) = lca_{T^*}(pa_{T^*}(r),u)$ will not increase the RF 737 distance since there are no matches to become mismatched. 738 739

Proof of Lemma 4.2. For binary trees U and V, let \mathcal{M}_U^V denote the LCA map from U 740 to V. That is, on input $u \in V(U)$, $\mathcal{M}_U^V(u)$ returns $lca_V(C_U(u))$. We will show that 741 $RF(S',T'') - RF(S',T') = RF(S^*,T^*) - RF(S',T') + 2m$. Observe that the only changes 742 from S' and T' to S^*, T^* and T'' are the formations of the extraneous subtrees $\{s, t\}$. Then, 743 it suffices to confirm that for every extraneous subtree $\{s, t\}$, the number of mismatched 744 clades in $T''(pa_{T''}\{s,t\})$ equals the number of mismatched clades in $T^*(\mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\}))$ 745 plus the number of extraneous subtrees. For an arbitrary extraneous subtree $\{s,t\}$ in T^* , we 746 first count the mismatched clades in $T''(pa_{T''}\{s,t\})$. Then, we count the mismatched clades 747 in $T^*(\mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\}))$ and compare. 748

1. Suppose v lies along the path from $pa_{T''}\{s,t\}$ to the parent of the canonical EF-R-RF(+) 749 position for T(t) in T''. Moreover, suppose u lies along the path from $pa_{T''}\{s,t\}$ to the 750 parent of the canonical EF-R-RF(+) position for S(s) in T''. Then $C_{S'}(\mathcal{M}_{T''}^{S'}(v)) \supseteq$ 751 $C_{T''}(v) \cup C_{S'}(t)$ since v is an *ancestor* of the canonical EF-R-RF(+) position of T(t)752 in T'' and hence $\mathcal{M}_{T''}^{S'}(v)$ is an ancestor of the canonical EF-R-RF(+) position of T(t)753 in S'. Moreover, $C_{T''}(v) \cap C_{S'}(t) = \emptyset$ if $v \neq pa_{T''}\{s,t\}$ by construction of T''. Hence 754 if $v \neq pa_{T''}\{s,t\}$, then v is mismatched with respect to S'. Likewise, $C_{S'}(\mathcal{M}_{T''}^{S'}(u)) \supseteq$ 755 $C_{T''}(u) \cup C_{S'}(s)$ and $C_{T''}(u) \cap C_{S'}(s) = \emptyset$ if $u \neq pa_{T''}\{s,t\}$. Hence if $u \neq pa_{T''}\{s,t\}$, 756 then u is mismatched with respect to S'. Note that by construction, $C_{T''}(pa_{T''}\{s,t\}) =$ 757 $C_{T'}(lca_{T'}(s,t))$. Hence $pa_{T''}(s,t)$ is matched with respect to S' if and only if $lca_{T'}(s,t)$ 758 is, and every other node along either path is mismatched. 759

Note that the only remaining node impacted in the formation of $\{s, t\}$ is the root of the extraneous subtree in T''. This node must be mismatched with respect to S' since S' is an extraneous free completion.

2. Now suppose v lies along the path from $pa_{T^*}\{s,t\}$ (the canonical EF-R-RF(+) position for 763 T(t) in T^*) to $\mathcal{M}_{T^*}^{T^*}(pa_{S^*}\{s,t\})$ (the least common ancestor of the EF-R-RF(+) positions 764 in T^*). Moreover, suppose u lies along the path from $\mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\})$ to the parent of 765 the canonical EF-R-RF(+) position for S(s) in T^* . Observe that $\mathcal{M}_{T^*}^{S^*}(v)$ is an ancestor of 766 the extraneous subtree $\{s, t\}$ in S^* , and therefore $\mathcal{M}_{T^*}^{S^*}(v)$ is an ancestor of the canonical 767 EF-R-RF(+) position for S(s) in S^* . Then $C_{S^*}(\mathcal{M}_{T^*}^{S^*}(v)) \supseteq C_{T^*}(v) \cup C_{S^*}(sib_{S^*}\{s,t\}),$ 768 where $C_{T^*}(v) \cap C_{S^*}(sib_{S^*}\{s,t\}) = \emptyset$ if $v \neq \mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\})$. Additionally, notice 769 that $\mathcal{M}_{T^*}^{S^*}(u)$ is an ancestor of the canonical EF-R-RF(+) position for S(s) in S^* , and 770 therefore $\mathcal{M}_{T^*}^{T^*}(u)$ is an ancestor of the extraneous subtree $\{s,t\}$. Then $C_{S^*}(\mathcal{M}_{T^*}^{T^*}(u)) \supseteq$ 771 $C_{T^*}(u) \cup C_{S^*}(s)$, where $C_{T^*}(u) \cap C_{S^*}(s) = \emptyset$ if $u \neq \mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\})$. It follows that if 772 $u \neq \mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\})$, then u is a mismatched node. Likewise, if $v \neq \mathcal{M}_{S^*}^{T^*}(pa_{S^*}\{s,t\})$, 773 then v is a mismatched node. Furthermore, $C_{T^*}(\mathcal{M}_{T^*}^{T^*}(pa_{S^*}\{s,t\}))$ is a matched clade 774 with respect to S^* if and only if $C_{T'}(lca_{T'}(s,t))$ is a matched clade with respect to S'. 775 Note, again, that the only remaining node impacted in the formation of $\{s, t\}$ is the root 776 of the extraneous subtree $\{s, t\}$. Since S^* and T^* are canonical R-RF(+) completions, 777 we know that this node must be matched in S^* and T^* . 778 Now, observe that the union of paths connecting the canonical EF-R-RF(+) positions for 779

S(s) and T(t) to $pa_{T^*}\{s,t\}$ in T^* is the same size as the union of paths connecting the canonical EF-R-RF(+) positions for S(s) and T(t) to $pa_{T''}\{s,t\}$ in T''. Moreover, every node in each union of paths (except the common ancestor) is mismatched. Finally, the root of $\{s,t\}$ is mismatched in T'' but matched in T^* . Since the choice of $\{s,t\}$ was arbitrary, we conclude with $RF(S',T'') - RF(S',T') = RF(S^*,T^*) - RF(S',T') + 2m$, where m is the number of extraneous subtrees in T^* . Equivalently, $RF(S',T'') = RF(S^*,T^*) + 2m$.

Proof of Lemma 4.4. Let S, T be two input binary rooted trees, and let S', T' be their canonical EF-R-RF(+) completions. By the proof of Lemma 4.1, we observe two important points: First, it can only be beneficial to move a maximal red or green subtree if the maximal subtree is eventually paired in an extraneous subtree. And second, a maximal red or green subtree will increase the RF distance by a lower amount if it is paired in an extraneous subtree closer to the canonical EF-R-RF(+) position. The recurrence relation follows by induction.

Base Case: No extraneous clades can be formed at a leaf node and there are no matches
 to become mismatched. Hence, the cost at each leaf is indeed zero.

Inductive Step: Assume we have computed all $Cost(x, \cdot, \cdot)$ for all descendants x of an internal node v. Let $c \in \{0, 1\}$ and $0 \le m \le cMax(c, v)$ be arbitrarily given. We first show that twice the maximal number of new extraneous subtrees $\{s, t\}$ that can be formed at vgiven c_L, c_R, m_L and m_R is equal to $g_c(m_L, m_R, c_L, c_R)$. There are two cases to consider: 1. $c_L = c_R = c$ and 2. $c_L \ne c_R$ (at least one of c_L and c_R must equal c).

1. Suppose $c_L = c_R = c$ and let m_L, m_R be arbitrary nonnegative values such that $m_L + m_R =$ 800 m. Then by the first observation above, the condition $m_L + m_R = m$ is optimal to regraft 801 m subtrees of color $c_L = c_R = c$ along the parent edge of v. By the second observation 802 above, if there are any extraneous subtrees that can be paired at v then it is optimal to 803 do so at v. We cannot pair any maximal red and green subtrees at v because $c_L = c_R = c_r$ 804 which means that all m subtrees regrafted along the parent edge of v are the same color. 805 Hence, twice the number of new extraneous subtrees that can be formed at v is equal to 806 $g_c(m_L, m_R, c_L, c_R) = 0$ when $c_L = c_R = c$. 807

K. Yao and M. S. Bansal

2. Now suppose without loss of generality that $c_L \neq c_R$ and let m_L, m_R be arbitrary nonnegative values such that $|m_L - m_R| = m$. Then by the two observations above, the condition $|m_L - m_R| = m$ is optimal to regraft the $m_L + m_R$ subtrees on the parent edge of v. By the second observation above, if there are any extraneous subtrees that can be paired at v then it is optimal to do so at v. Note that since $c_L \neq c_R$, we can pair exactly $min\{m_L, m_R\}$ extraneous subtrees at v. Hence, twice the number of *new* extraneous subtrees that can be formed at v is equal to $g_c(m_L, m_R, c_L, c_R) = 2 \min\{m_L, m_R\}$.

We now show that, regardless of the choice of colors c_L and c_R , the new increase in RF 815 distance between S' and T' only by regrafting m_L and m_R subtrees from $T'(v_L)$ and $T'(v_R)$ at 816 the parent edge of v, respectively, is equal to $f(m_L, v_L, c_L) + f(m_R, v_R, c_R)$. Once a subtree is 817 regrafted at the parent edge of v_L , the only clade that can become mismatched by regrafting 818 the subtree on the parent edge of v is $C_{T'}(v_L)$. This clade only becomes mismatched if 819 it is a matched clade and it is not contained in a maximal c_L -colored subtree. Once the 820 clade is mismatched, regrafting all remaining m_L maximal subtrees on the parent edge of v821 cannot make v mismatched again. Therefore, the act of pruning and regrafting m_L maximal 822 c_L -colored subtrees from the parent edge of v_L to the parent edge of v increases the RF 823 distance between S' and T' by $f(m_L, v_L, c_L)$, one for each of S' and T' if a match becomes 824 mismatched. By symmetry, the new increase in RF distance between S' and T' from pruning 825 and regrafting m_R maximal c_R -colored subtrees from v_R to v is equal to $f(m_R, v_R, c_R)$. 826

We have determined that the maximal number of new extraneous subtrees which can be 827 formed is equal to $g_c(m_L, m_R, c_L, c_R)$, and the new increase in RF distance is $f(m_L, v_L, c_L) +$ 828 $f(m_R, v_R, c_R)$. Then the change in cost from v_L and v_R to v is equal to $f(m_L, v_L, c_L) + f(m_R, v_R, c_R)$ 829 $f(m_R, v_R, c_R) - q_c(m_L, m_R, c_L, c_R)$. Note if a maximal c_L -colored subtree of $T'(v_L)$ is 830 regrafted along the parent edge of v, it must first already be regrafted along parent edge 831 of v_L by construction. Then, the cost of regrafting m_L subtrees at the parent edge of v_L 832 must be $Cost(v_L, m_L, c_L)$. By symmetry, the right subtree adds a cost of $Cost(v_R, m_R, c_R)$. 833 Moreover, the cost values also subtract the number of extraneous subtrees formed in $T'(v_L)$ 834 and $T'(v_R)$. 835

Hence, the value of $RF(S', \hat{T}) - 2p - RF(S', T')$ given fixed c_L, c_R, m_L and m_R is $Cost(v_L, m_L, c_L) + Cost(v_R, m_R, c_R) + f(m_L, v_L, c_L) + f(m_R, v_R, c_R) - g_c(m_L, m_R, c_L, c_R)$. By definition, the cost Cost(v, m, c) is equal to the minimum over all methods of moving maximal colored subtrees in T'(v) while leaving m maximal c-colored subtrees regrafted along the parent edge of v and unpaired in an extraneous subtree. Then, taking the minimum over all possible c_L, c_R, m_L and m_R values provides the optimal cost value.

Proof of Theorem 4.5. We note that a pair of canonical extraneous free completions can be computed in O(n) time. To compute the optimal cost values at each vertex of an EF-R-RF(+) completion, Algorithm *Compute-R-RF+(S,T)* has a total of three nested for loops, over (1) the postorder traversal, (2) the values of c and m, and (3) the values of c_L, c_R, m_L and m_R when the recurrence relation is invoked. The total time complexity is then the product of the sizes of each nested loop. Note there are a constant number of colors.

⁸⁴⁸ **1.** The postorder traversal has O(n) nodes to parse.

2. Notice *m* must be bounded above by $\max\{cMax(0, v), cMax(1, v)\} \leq cMax(0, rt(T')) + cMax(1, rt(T')) = k$ for any vertex *v*. Hence, we have another multiplicative O(k) factor.

3. For each Cost(v, m, c) value, we observe that the number of possible values of m_L and m_R considered is again bounded above by k, adding another multiplicative O(k) factor. Thus, the total runtime to compute all cost values is $O(nk^2)$. Once all cost values are computed, the RF(+) distance can be computed in O(1) time.