

# Fast recovery of evolutionary trees with thousands of nodes

Miklós Csűrös  
Department of Computer Science  
Yale University  
New Haven, CT 06520  
miklos.csuros@yale.edu

## ABSTRACT

We present a novel distance-based algorithm for evolutionary tree reconstruction. Our algorithm reconstructs the topology of a tree with  $n$  leaves in  $\mathcal{O}(n^2)$  time using  $\mathcal{O}(n)$  working space. In the general Markov model of evolution, the algorithm recovers the topology successfully with  $(1 - o(1))$  probability from sequences with polynomial length in  $n$ . Moreover, for almost all trees, our algorithm achieves the same success probability on polylogarithmic sample sizes. The theoretical results are supported by simulation experiments involving trees with 500, 1895, and 3135 leaves. The topologies of the trees are recovered with high success from 2000 bp DNA sequences.

## 1. INTRODUCTION

What is the largest evolutionary tree we can derive today? The limits of large-scale phylogeny reconstruction are determined by the availability of useful molecular sequences, and by the availability of useful reconstruction methods. With current advances in bioinformatics, DNA sequencing is now both fast and reliable enough that efficiency is becoming a major concern for large-scale problems in phylogeny reconstruction. Ambitious projects such as the Green Plant Phylogeny (GPP) project and the Ribosomal Database Project (RDP) [22] involve phylogenies with hundreds and thousands of homologous DNA sequences. When reconstructing a large tree, primary considerations for efficiency are computational speed and statistical accuracy. For instance, algorithms with exponential running time in the tree size cannot be used with trees that have more than a few tens of leaves. In fact, even algorithms that build trees with  $n$  leaves in  $\mathcal{O}(n^4)$  time may be too slow if  $n$  is in the order of thousands. On the other hand, algorithms that fail to extract topology information efficiently enough from the input sequences may require inordinately large amounts of data, preventing successful reconstruction of large trees. Recent theoretical results on the statistical efficiency of distance-

based algorithms [11, 9, 17] make them ideal candidates for large-scale phylogeny reconstruction. Nonetheless, simulation studies corroborating the theoretical predictions for trees of sizes comparable with those in GPP and RDP are still needed.

This paper has two goals. First, it presents a novel distance-based algorithm with provably high statistical and computational efficiency. Secondly, it reports the results of experiments conducted with large, biologically-motivated model trees with various ranges of mutation probabilities. In the experiments, we simulated DNA sequence evolution in the Jukes-Cantor [18] model on trees with 500, 1895, and 3135 leaves. These trees are unusually large for simulation studies. To our knowledge, the largest trees reconstructed from simulated data have 256 leaves [19]. The methods we compare include Neighbor-Joining [26], BioNJ [15], Weighbor [3], our Harmonic Greedy Triplets algorithm, and parsimony. Our results establish that even such large trees can be successfully recovered from DNA sequences with 2000 nucleotides. The experimental results support our theoretical results on the efficiency of our algorithm, called Harmonic Greedy Triplets with the Four-Point Condition (HGT/FP).

### 1.1 The general Markov model of sequence evolution

Mathematical models of sequence evolution play a fundamental role in providing a framework for developing evolutionary tree reconstruction algorithms, and for analyzing the algorithms' computational and statistical characteristics. A widely studied model is the general Markov model [27], in which sequence characters evolve independently. The model is formulated as follows. Let  $\mathcal{A} = \{1, 2, \dots, r\}$  be a finite alphabet of size  $r \geq 2$ , and for every  $\ell > 0$ , let  $\mathcal{A}^\ell$  denote the set of sequences over  $\mathcal{A}$  with length  $\ell$ . An evolutionary tree  $T$  is defined by an underlying tree and a mutation model. The underlying tree is a rooted binary tree representing the evolutionary ancestor-descendant relationships between its nodes. For every length  $\ell > 0$ , the mutation model randomly associates sequences of  $\mathcal{A}^\ell$  with the nodes. The vector formed by the characters in the same position of the sequences is called a *site*. In the general Markov model, the sites are independent and identically distributed.

Let  $E$  be the set of tree edges, let  $V = \langle u_1, \dots, u_k \rangle$  be the ordered set of tree nodes, and let  $\xi = \langle \xi^{(u_1)}, \dots, \xi^{(u_k)} \rangle$  be a site. For each node  $u \in V$ , the random variable  $\xi^{(u)}$  is called

the label of  $u$ . The distribution of  $\xi$  is determined by a *root label distribution*  $\pi = \langle \pi_1, \dots, \pi_r \rangle$  and by a set of *mutation matrices*  $\{\mathbf{M}_e: e \in E\}$  assigned to the tree edges. For each edge  $e$ ,  $\mathbf{M}_e$  is an  $r \times r$  stochastic matrix. The root label is distributed according to  $\pi$ . For all nodes  $u, v \in V$ , if  $v$  is the child of  $u$  on edge  $e$ , then for all characters  $i, j \in V$ ,

$$\mathbb{P}\{\xi^{(v)} = j \mid \xi^{(u)} = i\} = \mathbf{M}_e[i, j].$$

In other words, labels evolve along the tree from the root towards the leaves. For simplicity's sake, we assume that for every node  $u$  and character  $i \in \mathcal{A}$ ,  $\mathbb{P}\{\xi^{(u)} = i\} \neq 0$ , and that for every edge  $e$ ,  $0 < |\det \mathbf{M}_e| < 1$ .

## 1.2 Efficient evolutionary tree reconstruction

For an evolutionary tree  $T$ , the topology  $\Psi(T)$  is the unrooted binary tree, which is obtained from the underlying tree by removing the direction of the edges, and by replacing the root and its two incident edges with one single edge connecting the root's children. The problem of evolutionary tree reconstruction is that of finding  $\Psi(T)$  from sequences associated with the leaf set  $L$ , called a *sample*. An evolutionary tree reconstruction algorithm outputs an unrooted tree  $\Psi^*$  with the same leaf set  $L$  for an input sample. The algorithm succeeds if  $\Psi^* = \Psi(T)$ , i.e., if the topology is recovered. The *success rate* of an algorithm on sample sequences of length  $\ell$  is the probability that  $T$  generates a sample for which the algorithm succeeds. The minimum sample length required to achieve a given success rate  $(1 - \delta)$ , where  $0 < \delta < 1$  is the error probability sought, defines the *statistical efficiency* [19] of the algorithm. Let  $n$  denote the number of leaves in  $\Psi(T)$ . An algorithm is statistically efficient if the minimum sample length is polynomial in  $n$  and  $(1/\delta)$ . An algorithm is *computationally efficient* if its running time is polynomial in  $n$  and  $\ell$ .

Most popular evolutionary tree reconstruction algorithms today fall short of achieving computational or statistical efficiency. The HGT/FP algorithm, however, is both computationally and statistically efficient. In fact, it is the fastest statistically efficient algorithm to date, running in  $\mathcal{O}(n^2)$  time. Previous theoretical results on efficiency include those of Erdős *et al.* [12, 11], who started the study of statistical efficiency in the context of topology recovery and who also devised the first algorithms with provable computational and statistical efficiency. Their algorithms run in  $\mathcal{O}(n^5 \log n)$  and  $\mathcal{O}(n^4 \log n)$  time. Farach and Kannan [13] introduced the study of sample sizes required by evolutionary tree reconstruction algorithms in probabilistic models of sequence evolution, but their problem is slightly different from ours since their primary focus was to estimate the distribution of leaf labels based on a sample. Cryan *et al.* [7] gave a polynomial-time solution to the problem and proved that their algorithm is statistically efficient. The recently developed Disc Covering Method of Huson *et al.* [17] is statistically efficient but needs to solve an NP-hard problem at its core, and its heuristic implementation runs in  $\mathcal{O}(n^4)$  time. The statistically efficient Fast Harmonic Greedy Triplets algorithm [9] also runs in  $\mathcal{O}(n^2)$  time but its efficiency is contingent upon knowing of the highest mutation rate on a tree edge. A major advantage of HGT/FP is that it does not rely on any such input parameter. Warnow *et al.* [30] describe an algorithm that turns a statistically efficient algorithm need-

ing such a parameter into an algorithm that is statistically efficient without it. The transformation, however, increases the running time of the original algorithm by  $\mathcal{O}(n^4)$ .

## 1.3 Distance-based algorithms

A number of evolutionary tree reconstruction algorithms calculate an  $n \times n$  matrix in a preprocessing step from the input sequences, and build the tree based on the matrix. The input matrix estimates a matrix of *evolutionary distances* between leaves, defined as follows. Let  $T$  be an evolutionary tree with  $n$  leaves. Evolutionary distances between the nodes of  $\Psi(T)$  arise by equipping the edges of  $\Psi(T)$  with positive weights. The edge weights are also called *edge lengths*. The distance between two nodes is the sum of the edge lengths on the path between them. A *tree metric*  $\mathbf{D}$  over  $\Psi(T)$  is the  $n \times n$  matrix of pairwise distances between leaves. The tree metric  $\mathbf{D}$  is a functional of the site distribution and uniquely determines the topology  $\Psi(T)$ .

Common tree metrics in the general Markov model include paralinear distance [20] and LogDet distance [27, 21]. Define the  $r \times r$  matrix  $\mathbf{M}_{uv}$  for all nodes  $u, v$  of  $\Psi(T)$  by its entries as

$$\mathbf{M}_{uv} = \left[ \mathbb{P}\{\xi^{(v)} = j \mid \xi^{(u)} = i\} : i, j \in \mathcal{A} \right].$$

The paralinear distance is defined as

$$\mathbf{D}[u, v] = -\ln \sqrt{(\det \mathbf{M}_{uv})(\det \mathbf{M}_{vu})}, \quad (1)$$

for all leaves  $u, v$ . Since the expression on the right-hand side is additive along any path in  $\Psi(T)$ , the paralinear distance is a tree metric. Specifically, if the labels form a time-reversible Markov chain along any path in  $\Psi(T)$ , which is a frequent assumption in molecular evolutionary studies [29], then the paralinear distance is realized by setting the length of each edge  $e$  to  $(-\ln |\det \mathbf{M}_e|)$ . For all leaves  $u, v$ , let

$$\mathbf{J}_{uv} = \left[ \mathbb{P}\{\xi^{(u)} = i, \xi^{(v)} = j\} : i, j \in \mathcal{A} \right]$$

be the joint probability matrix of the leaf labels. The LogDet distance between leaves  $u$  and  $v$  is defined by  $\mathbf{D}[u, v] = -\ln |\det \mathbf{J}_{uv}|$ .

There are many other tree metrics within various subclasses of the general Markov model restricting the set of mutation matrices. For example, the Neyman-model [23] imposes that for every mutation matrix  $\mathbf{M}_e$  there exists a mutation probability  $0 < p_e < (1 - 1/r)$ , such that

$$\mathbf{M}_e[i, j] = \begin{cases} 1 - p_e & \text{if } i = j; \\ p_e/(r - 1) & \text{if } i \neq j. \end{cases}$$

Subsequently,  $\mathbf{D}[u, v] = -\ln(1 - \frac{r}{r-1} \mathbb{P}\{\xi^{(u)} \neq \xi^{(v)}\})$  is a tree metric. This tree metric is known as the Jukes-Cantor distance [18].

Distance-based algorithms thus have to estimate the tree metrics in a preprocessing step, which typically entails substituting probabilities in the tree metric's definition with relative frequencies calculated from the sample. We call such estimators *empirical tree metrics* and discuss them further in §2.3. An important feature of empirical tree metrics is

that estimation error increases with evolutionary distance between the leaves in question. In order to achieve statistical efficiency, a topology reconstruction algorithm has to strive to use leaves that are close to each other. The HGT/FP algorithm is designed with that goal in mind. The techniques we use to achieve that goal are based on analyzing the convergence rate of empirical tree metrics.

## 1.4 Recovering the topology from a tree metric

If an exact tree metric is known, then the problem of reconstructing the topology can be reduced to the problem of obtaining an unrooted tree with positive edge weights from distances between its leaves. A basic technique for that purpose uses *triplets*. A triplet  $uvw$  comprises three leaves  $u$ ,  $v$ , and  $w$  of  $\Psi(T)$ . Every triplet defines an internal node at which the three pairwise paths between the leaves intersect, with the four nodes forming a star. This internal node is the *center* of the triplet. Using the tree metric's definition, the distance between the center  $o$  and a leaf  $u$  in the triplet  $uvw$  can be calculated as

$$D_{uo} = \Delta(u, uvw) = \frac{\mathbf{D}[u, v] + \mathbf{D}[u, w] - \mathbf{D}[v, w]}{2},$$

where  $\Delta(u, uvw)$  denotes the triangle-star transformation formula on the right-hand side. This formula can be used repeatedly to reconstruct the topology with the edge lengths by adding one leaf and one internal node at a time [31]. The main idea of such a reconstruction is fairly simple. Let  $\Psi^*$  be a subtree of  $\Psi(T)$  spanned by a subset of the leaves, and let each edge of  $\Psi^*$  be weighted by the distance between its endpoints in  $\Psi(T)$ . If  $u$  and  $v$  are two leaves in  $\Psi^*$  and  $w$  is a leaf of  $\Psi(T)$  missing from  $\Psi^*$ , then the center  $o$  of  $uvw$  in  $\Psi^*$  is on the path  $P$  between  $u$  and  $v$ , and its exact location can be found by comparing  $\Delta(u, uvw)$  to the distances between  $u$  and the nodes on  $P$ . If that location falls properly on an edge  $e$  in  $\Psi^*$ , then  $o$  can be added on  $e$ , and  $w$  can be connected to it with an edge of length  $\Delta(w, uvw)$ . The edge lengths for the newly created edges between  $o$  and the endpoints of  $e$  can be calculated as the difference between  $\Delta(u, uvw)$  and the distances from  $u$  to the endpoints. This approach is complicated only by the fact that the center of  $uvw$  may be a node that is already in  $\Psi^*$ . In other words, there may be a node  $z$  on the path  $P$  between  $u$  and  $v$  that is at distance  $\Delta(u, uvw)$  from  $u$ . In that case  $w$  should be connected to  $\Psi^*$  through an internal node in the subtree rooted at  $w$  that contains neither  $u$  nor  $v$ , by using a different triplet. The reconstruction starts by selecting an arbitrary triplet  $uvw$  and initializing  $\Psi^*$  as the star formed by  $uvw$  and its center.

## 2. THE HGT/FP ALGORITHM

Using the algorithm outlined in §1.4 with an estimated tree metric  $\hat{\mathbf{D}}$  almost certainly leads to failure. The main reason is that  $\hat{\mathbf{D}}$  is usually not a tree metric, due to random estimation errors. As a consequence,  $\Psi(T)$  is not determined by  $\hat{\mathbf{D}}$ . We describe two specific measures to deal with the fact that  $\hat{\mathbf{D}}$  may not be a tree metric. These general measures are helpful for any algorithm following the outline of §1.4 and do not make assumptions about the exact way  $\hat{\mathbf{D}}$  is calculated. We address the problem of estimating edge lengths in §2.1. In §2.2 we address the problem of determining whether a triplet defines a new internal node in  $\Psi^*$ . These measures do not ensure statistical efficiency

on their own, and we analyze their error in §2.3 in conjunction with empirical tree metrics. The results of the analysis suggest a greedy selection of triplets, which is employed in HGT/FP. It is this greedy selection that leads not only to statistical efficiency but also to the  $\mathcal{O}(n^2)$  running time. The techniques for achieving the fast running time are described in §2.4.

## 2.1 Estimating edge lengths

The HGT/FP algorithm follows the general outline of the algorithm in §1.4 with specific techniques for dealing with estimated tree metrics. Let  $\hat{\mathbf{D}}$  be the estimated tree metric and let  $\hat{\Delta}(u, uvw) = (\hat{\mathbf{D}}[u, v] + \hat{\mathbf{D}}[u, w] - \hat{\mathbf{D}}[v, w])/2$  denote the corresponding triangle-star transformation formula for every triplet  $uvw$ . In order to prevent the accumulation of error in edge length estimates, the HGT/FP algorithm stores a triplet  $\text{def}(z)$  for each internal node  $z$  in  $\Psi^*$ , which is the triplet used for adding  $z$  to  $\Psi^*$ . For notational uniformity, let  $\text{def}(z) = \{z\}$  if  $z$  is a leaf. In order to add a new internal node  $o$  on an edge  $z_1z_2$  in  $\Psi^*$ ,  $o$  must be the center of a triplet  $u_1v_2w$  for which the following conditions hold:  $u_1 \in \text{def}(z_1)$ ,  $v_2 \in \text{def}(z_2)$ ,  $w$  is not in  $\Psi^*$ , and the edge  $z_1z_2$  is on the path between  $u_1$  and  $v_2$  in  $\Psi^*$ . Such an edge-triplet pair  $\langle z_1z_2, u_1v_2w \rangle$  is called *relevant*. Assume that  $z_1$  is an internal node, and  $\text{def}(z_1) = u_1v_1w_1$ . The value  $d_1 = |\hat{\Delta}(u_1, u_1v_1w_1) - \hat{\Delta}(u_1, u_1v_2w)|$  is an estimate of the distance between the centers of  $u_1v_1w_1$  and  $u_1v_2w$  in  $\Psi(T)$ . Similarly,  $d_2 = |\hat{\Delta}(v_2, u_2v_2w_2) - \hat{\Delta}(v_2, u_1v_2w)|$  estimates the distance between the centers of  $u_2v_2w_2$  and  $u_1v_2w$ . Let  $D_{z_1z_2}^*$  be the length of the edge  $z_1z_2$  in  $\Psi^*$ . The edge lengths for inserting  $o$  on  $z_1z_2$  are calculated by

$$\begin{aligned} D_{oz_1}^* &= (d_1 + D_{z_1z_2}^* - d_2)/2; \\ D_{oz_2}^* &= (d_2 + D_{z_1z_2}^* - d_1)/2. \end{aligned}$$

If  $z_i$  is a leaf for  $i = 1$  or for  $i = 2$ , then  $d_i = 0$  but otherwise the calculations are the same.

The theoretical importance of this procedure is that it results in edge length estimation errors that depend only on the error in estimating the center of individual triplets. Assume that  $\Psi^*$  is correct (i.e., it is topologically equivalent to the subtree of  $\Psi(T)$  spanned by the leaves of  $\Psi^*$ ) and the center of  $u_1v_1w$  falls onto the path between  $z_1$  and  $z_2$  in  $\Psi(T)$ . Assume further that the triplet centers are estimated within  $\epsilon$  error, i.e., that for every triplet  $x_1x_2x_3 \in \{u_1v_1w_1, u_2v_2w_2, u_1v_2w\}$  and leaf  $x_i$ ,

$$\left| \hat{\Delta}(x_i, x_1x_2x_3) - \Delta(x_i, x_1x_2x_3) \right| < \epsilon.$$

If  $|D_{z_1z_2}^* - D_{z_1z_2}| = 4\epsilon'$ , where  $D_{z_1z_2}$  is the distance between  $z_1$  and  $z_2$  in  $\Psi(T)$ , then  $|D_{oz_1}^* - D_{oz_1}| < 2\epsilon' + 2\epsilon$  and  $|D_{oz_2}^* - D_{oz_2}| < 2\epsilon' + 2\epsilon$ . Similar bounds hold if  $z_1$  or  $z_2$  is a leaf. If  $\epsilon' \leq \epsilon$ , then the error of the newly created edge lengths is bounded by  $4\epsilon$ . Consequently, if the maximum error in estimating triplet centers used by the algorithm is bounded by  $\epsilon$ , then all edge lengths are estimated within  $4\epsilon$  error, given that the topology is recovered correctly.

## 2.2 Finding triplet centers

While in the case of tree metrics, we can always tell whether a triplet defines a new internal node in the partially built topology  $\Psi^*$ , this is not so in the case of estimated tree metrics, where triplet centers may appear to define a new internal node due to estimation error. For example, even if the triplets  $uvw$  and  $uw'w'$  have the same center in  $\Psi(T)$ , it is possible that  $\hat{\Delta}(u, uvw) \neq \hat{\Delta}(u, uw'w')$  and thus the techniques in §1.4 for choosing triplets are likely to be inadequate.

A safeguarding measure in HGT/FP for dealing with an estimated tree metric is based on the *four-point condition* [5], which is defined as follows. An evolutionary tree with four leaves  $\{u, v, w, z\}$  has three possible topologies denoted by  $uv|wz$ ,  $uw|vz$  and  $uz|vw$ , depending on which leaf pairs are separated by the internal edge in the topology. The four-point condition states that by distance additivity, the topology is  $uv|wz$  if and only if

$$\mathbf{D}[u, v] + \mathbf{D}[w, z] < \mathbf{D}[u, w] + \mathbf{D}[v, z] = \mathbf{D}[u, z] + \mathbf{D}[v, w].$$

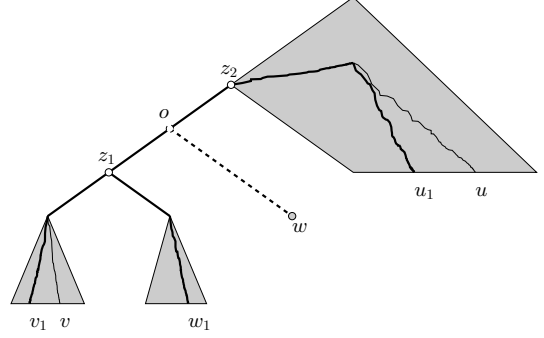
Since the equality of the two larger sums is unlikely when using an estimated tree metric, the HGT/FP algorithm employs the *relaxed four-point condition* [2], which for  $uv|wz$  is defined as

$$\begin{aligned} \hat{\mathbf{D}}[u, v] + \hat{\mathbf{D}}[w, z] &< \hat{\mathbf{D}}[u, w] + \hat{\mathbf{D}}[v, z]; \\ \hat{\mathbf{D}}[u, v] + \hat{\mathbf{D}}[w, z] &< \hat{\mathbf{D}}[u, z] + \hat{\mathbf{D}}[v, w]. \end{aligned} \quad (2)$$

Let  $\langle z_1 z_2, uvw \rangle$  be a relevant pair. The relaxed four-point condition is used to determine whether the center of  $uvw$  falls onto  $z_1 z_2$  in the following manner. Let  $z_1$  be an internal node in  $\Psi^*$ , let  $\text{def}(z_1) = u_1 v_1 w_1$ , and assume that  $z_2$  lies on the path between  $u_1$  and  $z_1$  in  $\Psi^*$  without loss of generality (see Figure 1). Recall that  $w$  is not a leaf in  $\Psi^*$ . HGT/FP tests whether the relaxed four-point condition holds for  $u_1 w | v_1 w_1$ . If so, then for the center  $o$  of  $uvw$ , the paths from  $z_1$  to  $o$  and to  $z_2$  overlap. The condition is used similarly with  $z_2$  if it is an internal node, in order to decide if the paths from  $z_2$  to  $o$  and to  $z_1$  overlap. If  $z_i$  is a leaf, then the condition for  $z_i$  is not tested. If the tested conditions hold for the pair  $\langle z_1 z_2, uvw \rangle$ , then it is called a *good relevant pair*. If  $\langle z_1 z_2, uvw \rangle$  is a good relevant pair, then HGT/FP concludes that the center of  $uvw$  can be inserted on the edge  $z_1 z_2$ . HGT/FP uses only good relevant pairs for adding new nodes. This way it tolerates some error in the estimated tree metrics, since Equation (2) may hold for the correct topology even if the distances between the leaves are estimated within a small error.

## 2.3 The Harmonic Greedy Triplets principle

The *Harmonic Greedy Triplets* (HGT) principle provides a guideline for the triplet selection mechanism when empirical tree metrics are used. The empirical tree metrics for the discussed distances are calculated as follows. The empirical Jukes-Cantor distance is computed by  $\hat{\mathbf{D}}[u, v] = -\ln(1 - \frac{r}{r-1} \hat{p}_{uv})$  where  $\hat{p}_{uv}$  is the relative frequency of the event  $\{\xi^{(u)} \neq \xi^{(v)}\}$  observed in the sample. (If the relative frequency is larger than  $(1 - 1/r)$ , then the distance is set to  $\infty$  or a large positive constant.) The empirical LogDet distance is calculated by computing the matrices  $\hat{\mathbf{J}}_{uv} = [\hat{p}_{uv,ij}]$  where  $\hat{p}_{uv,ij}$  is the relative frequency of the event  $\{\xi^{(u)} = i, \xi^{(v)} = j\}$  in the sample, and by setting



**Figure 1: Using the four-point condition with relevant triplets.**

$\hat{\mathbf{D}}[uv] = -\ln |\det \hat{\mathbf{J}}_{uv}|$ . Finally, in the case of empirical paralinear distance, we calculate  $\hat{\mathbf{M}}_{uv} = [\hat{p}_{uv,ij}/\hat{p}_{u,i}]$  where  $\hat{p}_{u,i}$  is the relative frequency of the event  $\{\xi^{(u)} = i\}$ , with the convention that if  $\hat{p}_{u,i} = 0$ , then  $\hat{\mathbf{M}}_{uv}[i, j] = 1$  for  $i = j$  and  $\hat{\mathbf{M}}_{uv}[i, j] = 0$  for  $i \neq j$ . The matrices  $\hat{\mathbf{M}}_{uv}$  and  $\hat{\mathbf{M}}_{vu}$  are used in place of  $\mathbf{M}_{uv}$  and  $\mathbf{M}_{vu}$  in Equation (1) to compute the empirical paralinear distance.

**LEMMA 1.** *Let  $\mathbf{D}$  be one of the tree metrics over  $\Psi(T)$  discussed, i.e., let  $\mathbf{D}$  be the paralinear, the LogDet, or the Jukes-Cantor distance. Let  $\hat{\mathbf{D}}$  be the corresponding empirical tree metric. There exist constants  $a, b > 0$  such that for all leaves  $u, v$  and  $0 < \epsilon < 1$ ,*

$$\begin{aligned} \mathbb{P}\left\{\mathbf{D}[u, v] - \hat{\mathbf{D}}[u, v] \geq -\ln(1 - \epsilon)\right\} &\leq ae^{-b\epsilon^2 S_{uv}^2}; \\ \mathbb{P}\left\{\mathbf{D}[u, v] - \hat{\mathbf{D}}[u, v] \leq -\ln(1 + \epsilon)\right\} &\leq ae^{-b\epsilon^2 S_{uv}^2}, \end{aligned} \quad (3)$$

where  $S_{uv} = e^{-\mathbf{D}[u, v]}$ .

**PROOF.** The lemma is proven for the Jukes-Cantor distance in [13, 9], and for the LogDet distance in [11]. For the paralinear distance, and for exact values of  $a$  and  $b$  in all cases, see [8].  $\square$

**DEFINITION 1.** *An estimated tree metric  $\hat{\mathbf{D}}$  for which Equation (3) holds is called an  $(a, b)$ -regular estimator for  $\mathbf{D}$ .*

The value  $S_{uv}$  in Lemma 1 is called the *similarity* between  $u$  and  $v$ . The HGT principle originates from Theorem 2 below, which relates the error in triplet center estimation with regular estimators to a harmonic average of similarities. For every triplet  $uvw$ , define the *average similarity*

$$S_{uvw} = \frac{3}{S_{uv}^{-1} + S_{uv}^{-1} + S_{uv}^{-1}} = \frac{3}{e^{\mathbf{D}[u, v]} + e^{\mathbf{D}[u, w]} + e^{\mathbf{D}[v, w]}}.$$

**THEOREM 2.** Let  $\hat{\mathbf{D}}$  be an  $(a, b)$ -regular estimator for the tree metric  $\mathbf{D}$ . For every triplet  $uvw$ , and  $0 < \epsilon < 1$ ,

$$\mathbb{P}\left\{\hat{\Delta}(u, uvw) - \Delta(u, uvw) \geq \frac{-\ln(1 - \epsilon)}{2}\right\} \leq 3a \exp\left(-\frac{b}{9}\ell\epsilon^2 S_{uvw}^2\right).$$

**PROOF.** We reported the theorem for the Jukes-Cantor distance in [10, 9]. For our detailed proof, see [8].  $\square$

The novel principle of HGT is that the selection of triplets in an algorithm following the outline of §1.4 with regular estimators should be a greedy selection of the triplet  $uvw$  with the largest *average estimated similarity* defined by

$$\hat{S}_{uvw} = \frac{3}{e^{\hat{\mathbf{D}}[u,v]} + e^{\hat{\mathbf{D}}[u,w]} + e^{\hat{\mathbf{D}}[v,w]}}.$$

## 2.4 Fast topology reconstruction

The HGT/FP algorithm uses good relevant pairs to add new nodes to  $\Psi^*$ . For every edge  $e \in \Psi^*$  and leaf  $w \notin \Psi^*$ , there are  $\mathcal{O}(1)$  relevant pairs of the form  $(e, uvw)$ . Maintaining the set of all  $\mathcal{O}(n^2)$  good relevant pairs while constructing  $\Psi^*$  would be possible: whenever a new leaf is added,  $\mathcal{O}(n)$  relevant pairs are eliminated,  $\mathcal{O}(n)$  new relevant pairs are created, and the new relevant pairs can be tested as described in §2.2. By the HGT principle, however, it is enough to consider one good relevant pair for every leaf  $w \notin \Psi^*$ , namely, the one in which the triplet has the largest average similarity. Denote the set of those pairs by  $\mathcal{R}$ . The HGT/FP algorithm maintains  $\mathcal{R}$  by updating it every time new nodes are added. The set  $\mathcal{R}$  is implemented as a vector of size  $n$  indexed by the leaves. Each entry of  $\mathcal{R}$  contains either null or a good relevant pair. In order to add a new internal node and a new leaf, the HGT/FP algorithm uses the relevant pair from  $\mathcal{R}$  in which the triplet has the largest average empirical similarity. The use of the HGT principle and relevant triplets results in the following theorem.

**THEOREM 3.** *The running time of the HGT/FP algorithm on a tree with  $n$  leaves is  $\mathcal{O}(n^2)$ . The algorithm uses  $\mathcal{O}(n)$  work space.*

**PROOF.** (*Sketch.*) The algorithm stores the tree  $\Psi^*$ , the vector  $\mathcal{R}$ , and  $\mathcal{O}(1)$  local variables, resulting in the  $\mathcal{O}(n)$  space requirement. Line F1 runs in  $\mathcal{O}(n^2)$  time. Since there are  $\mathcal{O}(1)$  relevant pairs for every edge  $e \in \Psi^*$  and leaf  $w \notin \Psi^*$ , Lines F3 and F8 take  $\mathcal{O}(n)$  time. Since there are  $\mathcal{O}(n)$  entries in  $\mathcal{R}$ , Lines F5 and F7 take  $\mathcal{O}(n)$  time also. Lines F2 and F6 update  $\Psi^*$  in  $\mathcal{O}(1)$  time. Thus, initialization in Lines F1–F3 takes  $\mathcal{O}(n^2)$  time, and the repeat loop of Line F4 is executed  $(n - 3)$  times, taking  $\mathcal{O}(n)$  time in each step, which results in the  $\mathcal{O}(n^2)$  total running time.  $\square$

The statistical efficiency of the HGT/FP algorithm is stated by the following theorem, for which the proof is sketched in the appendix.

**THEOREM 4.** Let  $T$  be an arbitrary evolutionary tree that has  $n$  leaves. Let  $\mathbf{D}$  be a tree metric over the topology  $\Psi(T)$ , and  $\hat{\mathbf{D}}$  be an  $(a, b)$ -regular estimator. Let  $D_{\min}$  be the minimum, and  $D_{\max}$  be the maximum distance between endpoints of edges in  $\Psi(T)$ , and define  $S_0 = e^{-D_{\max}}$ ,  $S_1 = 1 - e^{-D_{\min}}$ . For every error probability  $0 < \delta < 1$ , there exists

$$\ell = \mathcal{O}\left(\frac{\log \frac{a}{\delta} + \log n}{bS_0^{\mathcal{O}(\varrho)}S_1^2}\right), \quad (4)$$

with  $\varrho \leq 1 + \log_2(n - 1)$ , such that the success rate of HGT/FP is at least  $(1 - \delta)$  on samples of length  $\ell$ . Moreover, for almost every tree topology under the uniform or Yule-Harding distributions,  $\varrho = \mathcal{O}(\log \log n)$ .

*Remark.* The value  $\varrho$  in the theorem is the tree depth first studied in the context of evolutionary tree reconstruction by Erdős *et al.* [12] under different topology distributions.

## 3. SIMULATION EXPERIMENTS

We simulated DNA sequence evolution with 2000 nucleotides along three large trees in the Jukes-Cantor model. The 500-leaf tree has the topology of a seed plant phylogeny from Chase *et al.* [6]. The 1895-leaf tree is derived from the evolutionary tree of Eukaryotes in RDP [22]. The 3135-leaf tree is based on the subtree of Proteobacteria within the phylogeny of Prokaryotes in RDP. We scaled the edge lengths of the original trees using a linear transformation. We evaluated the accuracy of the reconstruction by the Robinson-Foulds error [25] RF%, which measures the percentage of misplaced internal edges in the tree.

The distance-based algorithms we used included BioNJ [15], Weighted Neighbor-Joining (Weighbor) [3], and Neighbor-Joining (NJ) [26]. The two former algorithms were recently developed, and are related to NJ. The NJ algorithm is arguably the most popular distance-based algorithm to date. All three algorithms run in  $\mathcal{O}(n^3)$  time [28], and their statistical efficiency is not proven. Atteson [1] derives sample length bounds for Neighbor-Joining and BioNJ similar to those of Equation (4), but the bounds use the diameter of the topology instead of  $\varrho$ . For Weighbor and BioNJ, we used the implementations provided by their authors; for NJ, we used its implementation in qclust [4]. We also included a heuristic parsimony method, called DNAPARS [14]. Parsimony algorithms aim at deriving a topology that gives rise to sample sequences with a minimal number of character changes along the edges, which is an NP-hard optimization problem [16]. It is known that the exact optimization is not statistically efficient for certain trees. In simulation experiments, however, they often perform very well [24].

Figure 3 shows the results of the simulations. In our experiments parsimony performs the best on the 500-leaf tree. Unfortunately, its running time increases rapidly with the tree size and the mutation probabilities, so that in some cases it takes several hours on a desktop computer<sup>1</sup> to recover the topology of the 500-leaf tree. In contrast, HGT/FP takes less than two minutes to reconstruct the topology of the

<sup>1</sup>We used a PC with Pentium III 500 MHz CPU and 256M memory, running Windows NT 4.0.

**Algorithm** Harmonic Greedy Triplets with Four Point Condition

**Input:** An  $n \times n$  estimated tree metric.

**Output:**  $\Psi^*$ .

F1 Select an arbitrary leaf  $u$  and find a triplet  $uvw$  with the maximum  $\hat{S}_{uvw}$ .

F2 Let  $\Psi^*$  be the star with three edges formed by  $uvw$  and its center  $o$ .

F3 Initialize  $\mathcal{R}$  using the good relevant pairs for edges  $uo$ ,  $vo$ ,  $wo$ .

F4 **repeat**

F5 Find  $\langle z_1z_2, uvw \rangle \in \mathcal{R}$  with the maximum  $\hat{S}_{uvw}$ .

F6 Add a new internal node  $z$  on  $z_1z_2$  and connect  $w$  to it.

F7 Delete the pairs from  $\mathcal{R}$  that contain the edge  $z_1z_2$ .

F8 Update  $\mathcal{R}$  using the good relevant pairs for edges  $z_1z$ ,  $z_2z$ ,  $wz$ .

F9 **until** all leaves are inserted to  $\Psi^*$ ; i.e., this loop has iterated  $(n - 3)$  times.

F10 Output  $\Psi^*$ .

**Figure 2: The HGT/FP algorithm.** Calculations pertaining to edge length estimation are sketched in §2.1. Good relevant pairs are discussed in §2.2. The set  $\mathcal{R}$  of good relevant pairs is discussed in §2.4.

3135-leaf tree. Weighbor proves to be even slower than parsimony despite its good asymptotic running time. It is also the least successful in recovering the topology among all algorithms considered. We omitted Weighbor and DNAPARS from the experiments with the larger trees because their running time increased to more than a day, and also omitted BioNJ because its performance is very similar to that of NJ. In the case of high mutation probabilities, HGT/FP performs better than NJ and BioNJ, while the neighbor-joining methods are better for low mutation probabilities, even though they do not recover the topology completely either.

Figure 4 shows the results of a different set of experiments comparing the effect of sample length on the recovery for HGT/FP and NJ. We simulated sequence evolution along the 1895-leaf tree with sample lengths ranging from 200 to 10000, for two edge scalings. In the case of high mutation probabilities, HGT/FP recovers the tree completely from 5000 bp sequences, while NJ misses more than 200 edges even for 10000 bp sequences. In the case of low mutation probabilities, NJ performs better than HGT/FP, but the difference is not so striking between the two algorithms as in the case of high mutation probabilities. In particular, the convergence rate of HGT/FP seems to be close to that of NJ.

## 4. CONCLUDING REMARKS

When working with trees with over one thousand leaves, the algorithms' running time becomes crucial. Existing  $\mathcal{O}(n^4)$ -time evolutionary tree building algorithms may take days to finish on today's desktop computers and slower algorithms are virtually unusable without having considerable insight into biological features of the data set at hand.

In addition to the computational issues, statistical characteristics of algorithms also become more stressed as one builds larger trees. Neighbor-Joining and most other algorithms have not been proven to require asymptotically polynomial sample sizes to correctly recover the topology, while HGT/FP is provably statistically efficient. Neighbor-Joining calculates edge lengths from an average that involves

distances between arbitrarily remote nodes in  $T$ , which may cause the estimation error to be very large. When the mutation probabilities are small, the averaging approach may be justified, as shown in the experiments with low mutation probabilities. On the other hand, the error committed while calculating the average is governed by the error in estimating the largest distance in the expression, which may be significant when mutation probabilities are large. In this case a greedy algorithm such as HGT/FP is more successful, as shown by our experimental results.

Many possible applications of evolutionary tree building algorithms may need to build large trees. Examples include large projects in evolutionary biology such as the ones cited and problems in molecular epidemiology. It will be of practical importance to determine which of the existing algorithms are the most suitable for the ranges of mutation probabilities and tree topologies defined by the application at hand.

## Acknowledgments

I am very thankful to Ming-Yang Kao and Dana Angluin for discussions leading to many of the results. I also thank the anonymous referees for helpful comments.

## 5. REFERENCES

- [1] K. Atteson. The performance of neighbor-joining methods of phylogeny reconstruction. *Algorithmica*, 25:251–278, 1999.
- [2] H.-J. Bandelt and A. Dress. Reconstructing the shape of a tree from observed dissimilarity data. *Advances in Applied Mathematics*, 7:309–343, 1986.
- [3] W. J. Bruno, N. D. Succi, and A. L. Halpern. Weighted Neighbor-Joining: a likelihood-based approach to distance-based phylogeny reconstruction. *Molecular Biology and Evolution*, 17:189–197, 2000.
- [4] J. Brzustowski. *qclust V0.2*, 1998. (<http://www.biology.ualberta.ca/jbrzusto/>).
- [5] P. Buneman. The recovery of trees from dissimilarity matrices. In F. R. Hodson, D. G. Kendall, and

- P. Tautu, editors, *Mathematics in the Archaeological and Historical Sciences*, pages 387–395. Edinburgh University Press, 1971.
- [6] M. W. Chase et al. Phylogenetics of seed plants: An analysis of nucleotide sequences from the plastid gene *rbcl*. *Annals of the Missouri Botanical Garden*, 80:528–580, 1993.
- [7] M. Cryan, L. A. Goldberg, and P. W. Goldberg. Evolutionary trees can be learned in polynomial time in the two-state general Markov model. In *39th Annual Symposium on Foundations of Computer Science*, pages 436–445. IEEE, 1998.
- [8] M. Csűrös. *Reconstructing Phylogenies in Markov Models of Sequence Evolution*. PhD thesis, Yale University, 2000. (<http://www.cs.yale.edu/~csuros-miklos/papers/>).
- [9] M. Csűrös and M.-Y. Kao. Provably fast and accurate recovery of evolutionary trees through Harmonic Greedy Triplets. *SIAM Journal on Computing*, 2001. To appear.
- [10] M. Csűrös and M.-Y. Kao. Recovering evolutionary trees through Harmonic Greedy Triplets. In *Proceedings of the Tenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 261–270, Baltimore, MD, 1999.
- [11] P. L. Erdős, M. A. Steel, L. A. Székely, and T. J. Warnow. A few logs suffice to build (almost) all trees (II). *Theoretical Computer Science*, 221:77–118, 1999.
- [12] P. L. Erdős, M. A. Steel, L. A. Székely, and T. J. Warnow. A few logs suffice to build (almost) all trees (I). *Random Structures and Algorithms*, 14:153–184, 1999.
- [13] M. Farach and S. Kannan. Efficient algorithms for inverting evolution. *J. ACM*, 46:437–449, 1999.
- [14] J. Felsenstein. *PHYLIP (Phylogeny Inference Package) version 3.5c*. Distributed by the author. University of Washington Department of Genetics, Seattle, Wash., 1993.
- [15] O. Gascuel. BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Molecular Biology and Evolution*, 14(7):685–695, 1997.
- [16] R. L. Graham and L. R. Foulds. Unlikelihood that minimal phylogenies for a realistic biological study can be constructed in reasonable computational time. *Mathematical Biosciences*, 60:133–142, 1982.
- [17] D. H. Huson, S. M. Nettles, and T. J. Warnow. Disk-Covering, a fast-converging method for phylogenetic tree reconstruction. *Journal of Computational Biology*, 6:369–386, 1999.
- [18] T. H. Jukes and C. R. Cantor. Evolution of protein molecules. In H. N. Munro, editor, *Mammalian Protein Metabolism*, volume III, chapter 24, pages 21–132. Academic Press, New York, 1969.
- [19] J. Kim. Large-scale phylogenies and measuring the performance of phylogenetic estimators. *Systematic Biology*, 47:43–60, 1998.
- [20] J. A. Lake. Reconstructing evolutionary trees from DNA and protein sequences: paralineal distances. *Proceedings of the National Academy of Sciences of the USA*, 91:1455–1459, 1994.
- [21] P. J. Lockhart, M. A. Steel, M. D. Hendy, and D. Penny. Recovering evolutionary trees under a more realistic model of sequence evolution. *Molecular Biology and Evolution*, 11:605–612, 1994.
- [22] B. L. Maidak et al. The RDP (Ribosomal Database Project) continues. *Nucleic Acids Research*, 28:173–174, 2000.
- [23] J. Neyman. Molecular studies of evolution: a source of novel statistical problems. In S. S. Gupta and J. Yackel, editors, *Statistical Decision Theory and Related Topics*, pages 1–27. Academic Press, New York, 1971.
- [24] K. Rice and T. Warnow. Parsimony is hard to beat! In *Computing and Combinatorics, Third Annual International Conference*, volume 1276 of *Lecture Notes in Computer Science*, pages 124–133. Springer-Verlag, 1997.
- [25] D. F. Robinson and L. R. Foulds. Comparison of phylogenetic trees. *Mathematical Biosciences*, 53:131–147, 1981.
- [26] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4(4):406–425, 1987.
- [27] M. A. Steel. Recovering a tree from the leaf colourations it generates under a Markov model. *Applied Mathematics Letters*, 7:19–24, 1994.
- [28] J. A. Studier and K. J. Keppler. A note on the neighbor-joining method of Saitou and Nei. *Molecular Biology and Evolution*, 5:729–731, 1988.
- [29] S. Tavaré. Some probabilistic and statistical problems in the analysis of DNA sequences. In *Lectures on mathematics in the life sciences*, volume 17, pages 57–86, Providence, RI, 1986. AMS.
- [30] T. Warnow, B. M. Moret, and K. S. John. Absolute convergence: true trees from short sequences. In *Twelfth Annual ACM-SIAM Symposium on Discrete Algorithms*, 2001. To appear.
- [31] M. S. Waterman, T. F. Smith, M. Singh, and W. A. Beyer. Additive evolutionary trees. *Journal of Theoretical Biology*, 64:199–213, 1977.

## APPENDIX

### A. PROOF OF THEOREM 4

This appendix outlines auxiliary lemmas leading to the proof of Theorem 4. In order to obtain the sample length bounds of the theorem, we bound the algorithm’s success probability. Let  $0 < S_{sm} < S_{lg} < 1$  be two threshold values on similarities with  $S_{sm} = S_{lg}/\sqrt{2}$  (we specify  $S_{lg}$  later). The

thresholds define three sets of triplets: every triplet  $uvw$  is either a *large* triplet, if  $S_{uvw} \geq S_{\text{lg}}$ , or a *medium* triplet if  $S_{\text{sm}} < S_{uvw} < S_{\text{lg}}$ , or a *small* triplet if  $S_{uvw} \leq S_{\text{sm}}$ . We show that with high probability, the HGT/FP algorithm recovers the tree correctly using only large and medium triplets. Let  $\Psi_k^*$  be the version of  $\Psi^*$  with  $k$  leaves at the beginning of the repeat loop in line F4 for  $k = 3, \dots, n-1$ , and let  $\Psi_n^* = \Psi^*$ , the algorithm's output. We prove for all  $k$  by induction that with high probability,  $\Psi_k^*$  is built correctly by using only large and medium triplets. Establishing the base case and the induction step relies on the following three arguments.

1. With high probability, the greedy selection favors large triplets over small triplets.
2. With high probability, HGT/FP correctly determines whether any relevant pair  $\langle z_1 z_2, uvw \rangle$ , for which  $uvw$  is not small is a good relevant pair.
3. There is always a relevant pair  $\langle z_1 z_2, uvw \rangle$  for which the triplet  $uvw$  is a large triplet and its center falls onto the edge  $z_1 z_2$ .

The first argument follows from Lemma 5, which states that if a regular estimator is used, then the probability of a small triplet appearing better than a large triplet is exponentially small in the sample sequence lengths.

LEMMA 5. *Assume that  $\hat{\mathbf{D}}$  is an  $(a, b)$ -regular estimator for  $\mathbf{D}$  calculated from sample sequences of length  $\ell$ . Let  $\mathcal{E}_3$  denote the random event that  $\hat{S}_{uvw} < \hat{S}_{u'v'w'}$  for every small triplet  $uvw$  and every large triplet  $u'v'w'$ .*

$$\mathbb{P}\{\mathcal{E}_3\} \geq 1 - \frac{a}{6} n^3 \exp\left(-b \frac{(\sqrt{2}-1)^2}{72} \ell S_{\text{lg}}^2\right).$$

PROOF. (*Sketch.*) Define the midpoint  $S_{\text{md}} = (S_{\text{sm}} + S_{\text{lg}})/2$ . For every small triplet  $uvw$ , the bound

$$\mathbb{P}\{\hat{S}_{uvw} \geq S_{\text{md}}\} \leq a \exp\left(-b \frac{(\sqrt{2}-1)^2}{72} \ell S_{\text{lg}}^2\right)$$

holds since  $\hat{\mathbf{D}}$  is  $(a, b)$ -regular. Similarly, for every large triplet  $u'v'w'$ ,

$$\mathbb{P}\{\hat{S}_{u'v'w'} \leq S_{\text{md}}\} \leq a \exp\left(-b \frac{(\sqrt{2}-1)^2}{72} \ell S_{\text{lg}}^2\right).$$

The two bounds imply the lemma, since there are  $\binom{n}{3} < n^3/6$  triplets.  $\square$

For the second argument, assume that  $\Psi_k^*$  has the correct topology, and that  $\text{def}(z)$  is a large or a medium triplet for every internal node  $z$ . Assume furthermore that  $\langle z_1 z_2, uvw \rangle$  is a relevant pair in  $\Psi_k^*$ . We take advantage of the fact that if  $uvw$  is not small, then the leaves for which the relaxed four-point condition is tested cannot be arbitrarily far from each other. Specifically, the distance between two leaves within the quartets is bounded from above by  $-2 \ln\left(\frac{2}{3} S_{\text{sm}}\right)$ , based on the fact that for every triplet  $uvw$  with center  $o$ ,  $S_{uo} \geq \frac{2}{3} S_{uvw}$ .

LEMMA 6. *A quartet is a short quartet if each pairwise distance between its leaves is less than  $-2 \ln\left(\frac{2}{3} S_{\text{sm}}\right)$ . Let  $\mathcal{E}_4$  denote the random event that for every short quartet  $uv|wz$  the relaxed four-point condition of Equation (2) holds.*

$$\mathbb{P}\{\mathcal{E}_4\} \geq 1 - an^2 \exp\left(-\frac{b}{81} \ell S_{\text{lg}}^4 S_1^2\right).$$

PROOF. Using the technique of [12, 11], the probability of the complementary event  $\bar{\mathcal{E}}_4$  is bounded by the probability that there is a leaf pair  $(u, v)$  for which  $|\mathbf{D}[u, v] - \hat{\mathbf{D}}[u, v]| \geq (-\ln(1 - S_1))/2$ . The lemma follows from the fact that  $\hat{\mathbf{D}}$  is  $(a, b)$ -regular, and that there are  $\binom{n}{2} < \frac{n^2}{2}$  leaf pairs.  $\square$

The third argument is based on Lemma 7. Lemma 7 depends on how large the defining triplets are for the internal nodes of  $\Psi(T)$ , and determines the value of  $S_{\text{lg}}$ .

DEFINITION 2. *Define the tree depth  $\rho$  as the smallest number such that for every edge  $e \in \Psi(T)$ , there is a path from each endpoint to a leaf with at most  $\rho$  edges, which does not go through  $e$ .*

LEMMA 7. *Assume that  $\Psi_k^*$  has the correct topology, and that for every internal node  $z$ ,  $\text{def}(z)$  is not small. If*

$$S_{\text{lg}} \leq \frac{3\sqrt{2}}{2} \left(\frac{\sqrt{2}-1}{\sqrt{2}+1}\right)^2 S_0^{2e+4} \quad \left(\approx \frac{S_0^{2e+4}}{16}\right),$$

*then the following statement holds. For every edge  $z_1 z_2 \in \Psi_k^*$ , if  $z_1$  and  $z_2$  are not connected by one edge in  $\Psi$ , then there exists a relevant pair  $\langle z_1 z_2, uvw \rangle$  such that  $uvw$  is large, and the center of  $uvw$  falls onto the path between  $z_1$  and  $z_2$  in  $\Psi(T)$ .*

PROOF. The full proof of this lemma can be found in [8]; a similar claim is proven in [9] for the Jukes-Cantor model.  $\square$

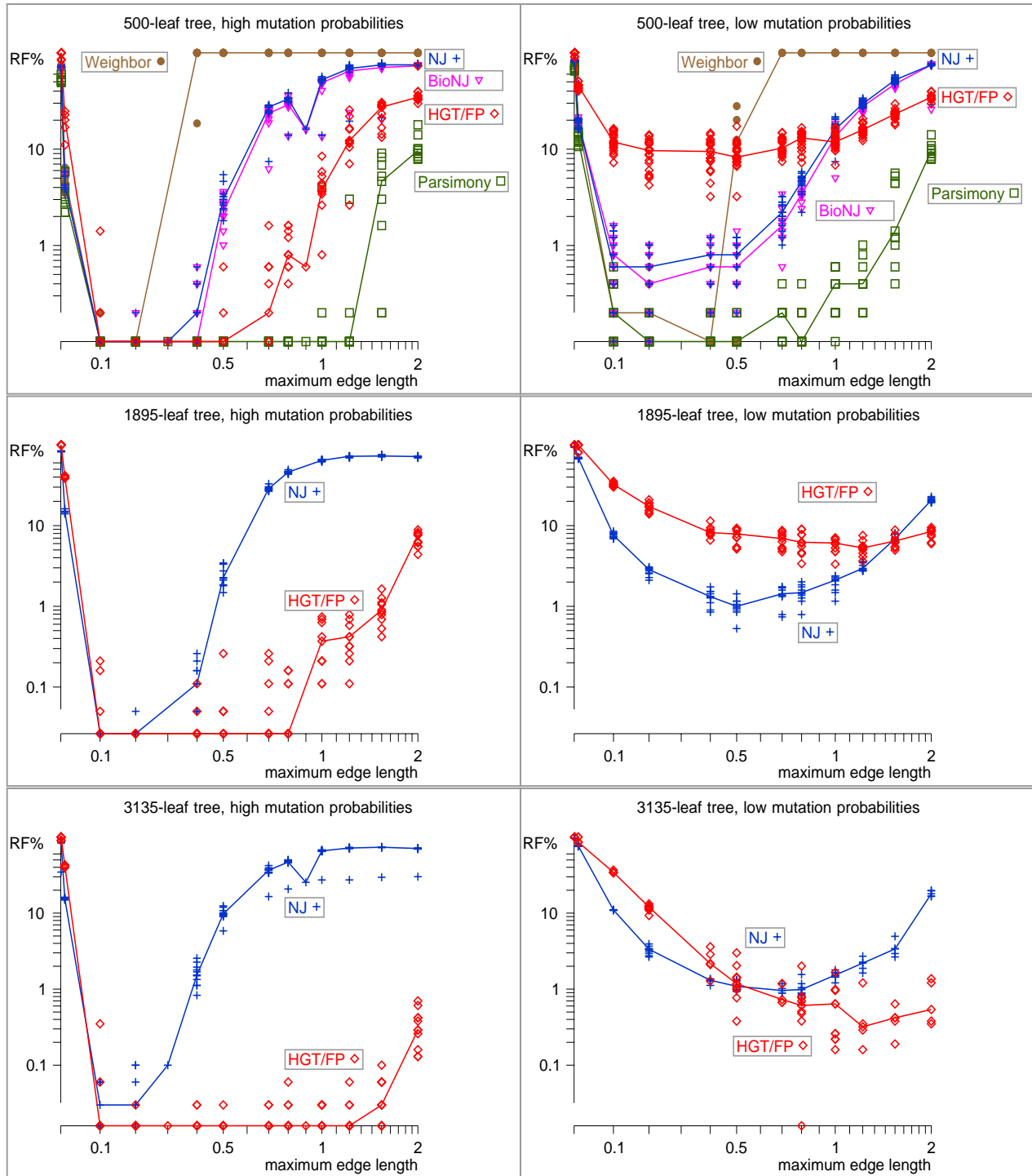
PROOF OF THEOREM 4. (*Sketch.*) We prove by induction that  $\mathcal{E}_3$  and  $\mathcal{E}_4$  imply that for all  $k$ ,  $\Psi_k^*$  has the correct topology and is built using only large and medium triplets. By  $\mathcal{E}_3$ , a medium or large triplet is selected to initialize  $\Psi^*$ , and thus the claim holds for  $k = 3$ . Assume that it holds for  $3 < k < n$ . By Lemma 7, there is a large triplet  $uvw$  in a relevant pair  $\langle z_1 z_2, uvw \rangle$  for which the center of  $uvw$  falls onto  $z_1 z_2$ . By  $\mathcal{E}_4$ ,  $\langle z_1 z_2, uvw \rangle$  is a good relevant pair. By  $\mathcal{E}_3$ , HGT/FP selects a good relevant pair with a medium or large triplet in Line F5, and by  $\mathcal{E}_4$ , the new nodes are added to  $\Psi_k^*$  correctly.

By Lemmas 5 and 6, if

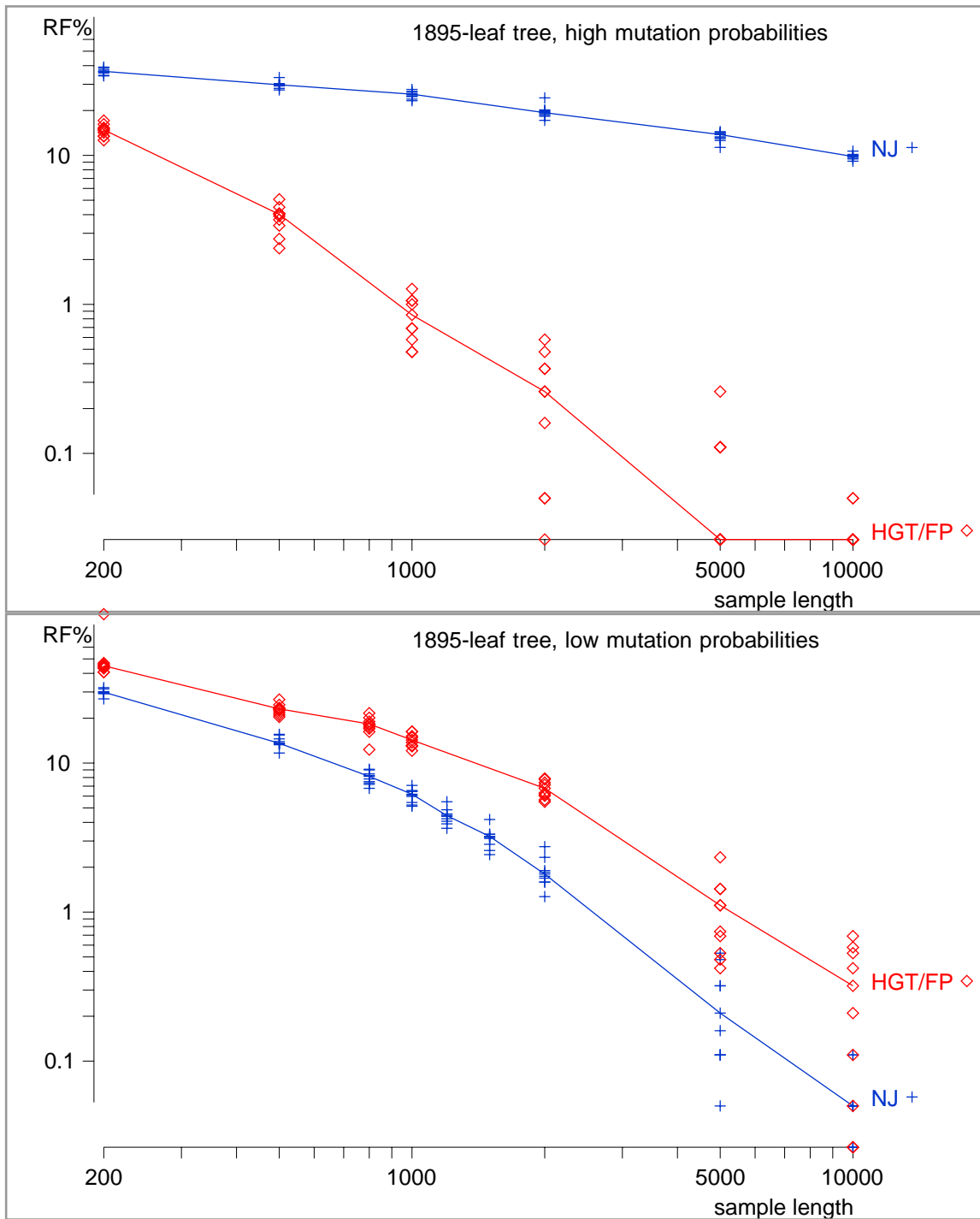
$$\ell \geq \max\left\{420 \frac{3 \ln n + \frac{a}{3\delta}}{b S_{\text{lg}}^2}, 81 \frac{2 \ln n + \frac{2a}{\delta}}{b S_1^2 S_{\text{lg}}^4}\right\},$$

then both  $\mathcal{E}_3$  and  $\mathcal{E}_4$  hold with probability at least  $(1 - \delta)$ . The theorem follows from setting  $S_{\text{lg}} = S_0^{2e+4}/17$  as suggested by Lemma 7.  $\square$





**Figure 3: Simulation of DNA sequence evolution in the Jukes-Cantor model along large trees with different mutation probabilities.** The plots show the percentage of misplaced internal edges (Robinson-Foulds error) as a function of the largest edge length  $D_{\max}$  in the tree after linear scaling. The graphs are calculated from generating ten set of samples with 2000 bp long sequences. The graphs go through the median values. On the left-hand side, the minimum edge length equals  $D_{\max}/10$ . On the right-hand side, it is set to  $D_{\max}/100$ . For reference,  $D_{\max} = 0.5$  corresponds to maximum mutation probability 0.30; on the left-hand side the minimum mutation probability equals 0.037 and on the right-hand side it equals 0.0037 for that scaling. Most edge lengths in the trees are very close to the minimum edge length.



**Figure 4: Simulation of DNA sequence evolution in the Jukes-Cantor model along the 1895-leaf tree with different sample lengths.** The plots show the percentage of misplaced internal edges (Robinson-Foulds error) as a function of the sample lengths. The graphs are calculated from generating ten sets of samples for each sequence length. The graphs go through the median values. The edge lengths on the top are linearly scaled to fall into the interval  $[0.1, 1]$ . The edge lengths on the bottom fall into the interval  $[0.01, 1]$ .