

WHY NEIGHBOR-JOINING WORKS

RADU MIHAESCU, DAN LEVY, AND LIOR PACTER

ABSTRACT. We show that the neighbor-joining algorithm is a robust quartet method for constructing trees from distances. This leads to a new performance guarantee that contains Atteson’s optimal radius bound as a special case and explains many cases where neighbor-joining is successful even when Atteson’s criterion is not satisfied. We also provide a proof for Atteson’s conjecture on the optimal edge radius of the neighbor joining algorithm. The strong performance guarantees we provide also hold for the quadratic time fast neighbor-joining algorithm, thus providing a theoretical basis for inferring very large phylogenies with neighbor-joining.

1. INTRODUCTION

The widely used neighbor-joining algorithm [17] has been extensively analyzed and compared to other tree construction methods. Previous studies have mostly focused on empirical testing of neighbor-joining. Examples include the comparison of neighbor-joining with quartet [10] and maximum likelihood [9] methods, comprehensive comparisons of multiple programs [8, 11], and detailed testing of the limits of the neighbor-joining algorithm [12]. These studies have concluded that neighbor-joining is effective for many problems, and have recommended the algorithm. For example, in [10] it is remarked that “...quartet-based methods are much less accurate than the simple and efficient method of neighbor-joining”. In a recent study, Tamura et al. [20] conclude that there are “...bright prospects for the application of the NJ and related methods in inferring large phylogenies.”

Furthermore, new methods are now almost always compared with neighbor-joining to establish an improvement in performance [2, 6, 7, 14, 15, 16, 18]. In other words, neighbor-joining has become the standard by which new phylogenetic algorithms are compared, and continues to surface as an effective candidate method for constructing large phylogenies. This is remarkable, considering the simplicity of the neighbor-joining algorithm:

- (1) Given a dissimilarity map $\delta : X \times X \rightarrow \mathbb{R}$ (this is a map that satisfies $\delta(i, j) = \delta(j, i)$ and $\delta(i, i) = 0$), compute the *Q-criterion* for δ

$$Q_\delta(i, j) = (n - 2)\delta(i, j) - \sum_{k \neq i} \delta(i, k) - \sum_{k \neq j} \delta(j, k)$$

Then select a pair a, b that minimize Q_δ as motivated by the following theorem:

Theorem 1 (Saitou-Nei [17] and Studier-Keppler [19]). *Let δ_T be the tree metric corresponding to the tree T . The pair a, b that minimizes $Q_{\delta_T}(i, j)$ is a cherry in the tree.*

- (2) If there are more than three taxa, replace the putative cherry a and b with a leaf j_{ab} , and construct a new dissimilarity map where $\delta(i, j_{ab}) = \frac{1}{2}(\delta(i, a) + \delta(j, b))$. This is called the *reduction step*.
- (3) Repeat until there are three taxa.

Although the Q -criterion is easy to compute, the formula seems, at first glance, somewhat contrived and mysterious. However, the formulation of the Q -criterion is far from accidental and has many useful properties. For example, it is linear in the distances, is permutation equivariant (the input order of the taxa don't matter), and it is *consistent*, i.e., it correctly finds the tree corresponding to a tree metric. Bryant [3] has recently shown that the Q -criterion is in fact the unique selection criterion satisfying these properties. These results motivate the neighbor-joining algorithm, but they do not offer any insight into its performance with dissimilarity maps that are not tree metrics. More importantly, they do not address the central question of the behavior of neighbor-joining on dissimilarity maps that arise from maximum likelihood estimation of distances between sequences in multiple alignments. There is one result that addresses precisely these issues:

Theorem 2 (Atteson). *Neighbor-joining has l_∞ radius $\frac{1}{2}$.*

This means that if the distance estimates are at most half the minimal edge length of the tree away from their true value then neighbor-joining will reconstruct the correct tree. Atteson's theorem shows that neighbor-joining is *statistically consistent*. Informally, this means that neighbor-joining reconstructs the correct tree from dissimilarity maps estimated from sufficiently long multiple alignments. This has been a widely used justification for the observed success of neighbor-joining.

However, as noted in [13], Atteson's condition frequently fails to be satisfied even when neighbor-joining is successful. This is also remarked on in [4]: "In practice, most distances are far from being nearly additive. Thus, although, important, optimal reconstruction radius is not sufficient for an algorithm to be useful in practice." Our main result is an explanation of the observation that neighbor-joining is useful in practice, which we obtain by using a new consistency theorem (Theorem 17 in Section 4). Roughly speaking, our theorem states that neighbor-joining is successful when it works correctly for the quartets in the tree. Thus, Theorem 17 provides a crucial link between neighbor-joining and quartet methods, a connection that is first developed in Section 3. We also show that Atteson's theorem is a special case of our theorem.

In Section 5 we present a proof of Atteson's conjecture on the optimal edge radius of neighbor-joining. For a dissimilarity map δ whose l_∞ distance to a tree metric δ_T is less than $\frac{\epsilon}{4}$, we prove that the output T' of neighbor-joining applied to δ will contain all edges of T of length at least ϵ . We then say that neighbor-joining has l_∞ edge radius $\frac{1}{4}$. We say that T' contains the edge e if there exists an edge $e' \in T'$ such that the split of the taxa induced by removing e' from T' is the same as that induced by removing e

from T . This result is tight, as [1] provides a counterexample for edge radius larger than $\frac{1}{4}$. The methods we employ for this result are virtually the same as the ones used in the proof of our general consistency criterion. The key step is showing that an adaptation of the global criterion holds in the edge preservation scenario, suggesting that this may in fact be the "right way" of analyzing neighbor-joining.

2. THE SHIFTING LEMMA

In this section we provide a few preliminary observations which are crucial to the understanding of our results. We begin by observing that there is an alternative possible reduction step in the neighbor joining algorithm which is (2') If there are more than three taxa, replace the putative cherry a and b with a leaf j_{ab} , and construct a new dissimilarity map where $\delta(i, j_{ab}) = \frac{1}{2}(\delta(i, a) + \delta(j, b) - \delta(a, b))$.

Notice that the only difference between (2) and (2') consists of addition of the $-\delta(a, b)$ term. On a tree metric, the first version of the algorithm corresponds to replacing a cherry by a single node whose distance to the rest of the tree is the average of the distances of the two collapsed nodes. In the second version, the reduction step is equivalent with replacing the cherry by its "root", i.e., the point where the path between the cherry nodes connects to the rest of the tree. Although not immediately apparent, the two versions of the algorithm are equivalent, as we prove below.

Definition 3. We say the dissimilarity map δ' is a *shift* of δ if and only if there exists a fixed ϵ and a distinguished taxon a , such that $\delta'(a, x) = \delta(a, x) + \epsilon$ for all $x \neq a$ and $\delta'(x, y) = \delta(x, y)$ for all $x, y \neq a$.

Lemma 4 (The shifting lemma). *Shifting does not affect the outcome of neighbor-joining.*

Proof: This follows from the observation that shifting by ϵ around a changes the value of $Q(x, y)$ by exactly -2ϵ , for all pairs x, y . Moreover, the result of collapsing taxa x, y in the shifted dissimilarity map is the same as the result of collapsing the same taxa in the initial dissimilarity map, or an ϵ shift of it. By these two observations, at each step neighbor-joining will collapse the same pairs as before the shift. \square We note that the proof holds for both versions of the algorithm.

Corollary 5. *The two versions of neighbor-joining are equivalent: they always produce the same tree.*

Proof: Collapsing x, y by the second reduction method gives a dissimilarity map that is a $\frac{\delta(x, y)}{2}$ -shift of the one produced by collapsing by the first type of reduction step. \square

These results are well known. For example, they are implicit in the results of [3] on the uniqueness of the Q -criterion.

It is important to notice that the operation of shifting a real tree metric δ_T by ϵ around taxon a corresponds exactly to modifying the length of the leaf edge of T corresponding to a by exactly ϵ . So in effect we can allow negative leaf edges since shifting around all vertices by a large enough constant will make them positive, while the outcome of the algorithm is the same. Of course, the statement of our edge radius results does not make

sense in the case of negative edge lengths. However, the only negative edges are the leaf edges, and in that case the statements we make are vacuous. They hold trivially since neighbor-joining reconstructs the bipartition consisting of one taxon vs. the rest of the taxa set, correctly, regardless of the input.

In the remainder of the paper, by a tree metric we will mean a shift of a tree metric, i.e. a metric corresponding to a tree where leaf edges are allowed to be negative.

3. QUARTETS AND NEIGHBOR-JOINING

We now show that for four taxa, neighbor joining is equivalent to the *four point method* [5]. We will use the notation $(ij : kl)$ for a quartet where i, j and k, l are cherries (two leaves adjacent to the same vertex) in a tree.

Proposition 6. *Let $X = \{i, j, k, l\}$ and $\delta : X \times X \rightarrow \mathbb{R}$ be a dissimilarity map. The neighbor-joining algorithm will return a tree $(ij : kl)$ where $\delta(i, j) + \delta(k, l) \leq \min(\delta(i, k) + \delta(j, l), \delta(i, l) + \delta(j, k))$.*

This result can be easily derived using the Q -criterion, but we prefer to motivate it using an alternative formulation of the neighbor-joining criterion.

For a dissimilarity map δ , let

$$w_\delta(ij : kl) = \frac{1}{2} (\delta(i, k) + \delta(i, j) + \delta(j, k) + \delta(j, l)) - \delta(i, j) - \delta(k, l).$$

Note that for a quartet $(ij : kl)$ in a tree T with corresponding tree metric δ_T , $w_{\delta_T}(ij : kl)$ is double the length of the internal edge in the quartet.

Theorem 7. *If δ_T is the tree metric corresponding to a tree T and*

$$Z_{\delta_T}(i, j) = \sum_{k, l \in X \setminus \{i, j\}} w_{\delta_T}(ij : kl)$$

then the pair a, b that maximizes $Z_{\delta_T}(i, j)$ is a cherry in the tree.

Proof: Observe that

$$Z_{\delta_T}(i, j) = -T - \frac{(n-1)}{2} Q_{\delta_T}(i, j)$$

where $T = \sum_{x, y \in T} \delta_T(x, y)$ does not depend on i or j . The theorem now follows directly from Theorem 1. \square

Although the naive computation of the Z -criterion requires quadratic time, the formulation of the Q -criterion shows that each entry in the Z -matrix is just a sum of a linear number of distances. One may wonder why the Z -criterion is worth mentioning at all. We outline a number of reasons why the Z -criterion may be a more natural way to formulate the neighbor-joining selection criterion. For example, note that in the case of four taxa i, j, k, l , the Z -criterion is just $Z_\delta(i, j) = w_\delta(ij : kl)$ and Proposition 6 follows immediately from Theorem 7. The Z -criterion also highlights the fact that for a tree metric, the neighbor-joining selection criterion does not depend on the length of edges adjacent to leaves. This is remarked on in the proof of the consistency of neighbor-joining in [3]. Most importantly, the Z -criterion highlights the connection

between neighbor-joining and quartet methods. Recall that the naive quartet method consists of choosing a quartet for each four taxa using the four point method (Proposition 6), and then returning the tree consistent with all the quartets (if such a tree exists). This leads us to

Definition 8. A dissimilarity map δ is *quartet consistent* with a tree T if for every $(ij : kl) \in T$, $w_\delta(ij : kl) > \max(w_\delta(ik : jl), w_\delta(il : jk))$.

The naive quartet method will reconstruct a tree T from a dissimilarity map δ if δ is quartet consistent with T . In the next section we will prove the following extension of Proposition 6:

Theorem 9. If $4 \leq |X| \leq 7$ and $\delta : X \times X \rightarrow R$ is a dissimilarity map that is quartet consistent with a binary tree T then the neighbor-joining algorithm applied to δ will construct a tree with the same topology as T . Furthermore, if $5 \leq |X| \leq 7$ then there exists $\epsilon > 0$ such that if $\|\tilde{\delta} - \delta\|_\infty < \epsilon$, neighbor-joining applied to $\tilde{\delta}$ will reconstruct a tree with the same topology as T .

As we have pointed out, neighbor-joining is equivalent to the naive quartet method for $|X| = 4$. Theorem 9 states that neighbor-joining is at least as good as the naive quartet method for trees with at most 7 taxa, and is in fact *robust* to small changes in the metric. The following example illustrates this.

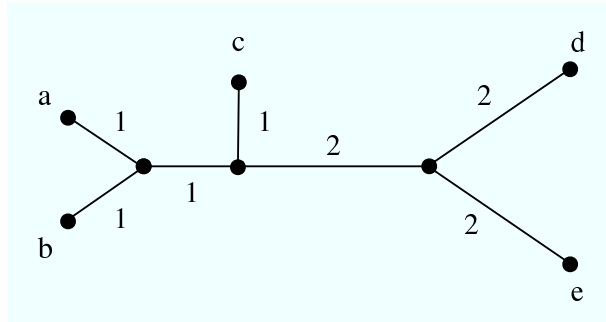


FIGURE 1. A five leaf tree.

Example 10. Let T be the 5 leaf tree shown in Figure 1 that corresponds to the tree metric δ_T , and consider the distorted dissimilarity map δ

$$\delta_T = \begin{matrix} & a & b & c & d & e \\ \begin{matrix} a \\ b \\ c \\ d \\ e \end{matrix} & \begin{pmatrix} 0 & 2 & 3 & 6 & 6 \\ 2 & 0 & 3 & 6 & 6 \\ 3 & 3 & 0 & 5 & 5 \\ 6 & 6 & 5 & 0 & 4 \\ 3 & 6 & 5 & 4 & 0 \end{pmatrix} \end{matrix}, \quad \delta = \begin{matrix} & a & b & c & d & e \\ \begin{matrix} a \\ b \\ c \\ d \\ e \end{matrix} & \begin{pmatrix} 0 & 2 & 3 & 6 & 3 \\ 2 & 0 & 3 & 6 & 6 \\ 3 & 3 & 0 & 5 & 5 \\ 6 & 6 & 5 & 0 & 4 \\ 3 & 6 & 5 & 4 & 0 \end{pmatrix} \end{matrix}.$$

Note that δ is not quartet consistent with T , because $\delta(a, e) + \delta(b, c) < \delta(a, b) + \delta(c, e)$. However it is easy to verify that neighbor-joining constructs a tree with the same topology as T . The example shows that neighbor-joining can construct the correct tree even when the naive quartet method fails.

The next example show that Theorem 9 fails for trees with more than 7 taxa.

Example 11. Let T be the 8 leaf tree shown below and let δ_T be its corresponding tree metric.

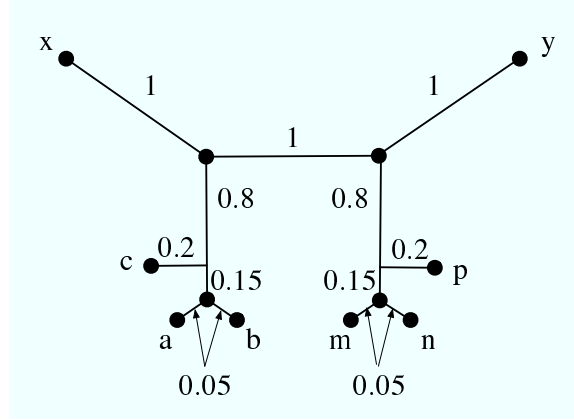


FIGURE 2. An eight leaf tree.

$$\delta_T = \begin{matrix} & \begin{matrix} x & y & a & b & c & m & n & p \end{matrix} \\ \begin{matrix} x \\ y \\ a \\ b \\ c \\ m \\ n \\ p \end{matrix} & \begin{pmatrix} 0 & 3 & 2 & 2 & 2 & 3 & 3 & 3 \\ 3 & 0 & 3 & 3 & 3 & 2 & 2 & 2 \\ 2 & 3 & 0 & 0.1 & 0.4 & 3 & 3 & 3 \\ 2 & 3 & 0.1 & 0 & 0.4 & 3 & 3 & 3 \\ 2 & 3 & 0.4 & 0.4 & 0 & 3 & 3 & 3 \\ 3 & 2 & 3 & 3 & 3 & 0 & 0.1 & 0.4 \\ 3 & 2 & 3 & 3 & 3 & 0.1 & 0 & 0.4 \\ 3 & 2 & 3 & 3 & 3 & 0.4 & 0.4 & 0 \end{pmatrix} \end{matrix},$$

$$\delta = \begin{matrix} & \begin{matrix} x & y & a & b & c & m & n & p \end{matrix} \\ \begin{matrix} x \\ y \\ a \\ b \\ c \\ m \\ n \\ p \end{matrix} & \begin{pmatrix} 0 & 2.7 & 2.6 & 2.6 & 2.6 & 4.4 & 4.4 & 4.4 \\ 2.7 & 0 & 4.4 & 4.4 & 4.4 & 2.6 & 2.6 & 2.6 \\ 2.6 & 4.4 & 0 & 0.1 & 0.4 & 2.7 & 2.7 & 2.7 \\ 2.6 & 4.4 & 0.1 & 0 & 0.4 & 2.7 & 2.7 & 2.7 \\ 2.6 & 4.4 & 0.4 & 0.4 & 0 & 2.7 & 2.7 & 2.7 \\ 4.4 & 2.6 & 2.7 & 2.7 & 2.7 & 0 & 0.1 & 0.4 \\ 4.4 & 2.6 & 2.7 & 2.7 & 2.7 & 0.1 & 0 & 0.4 \\ 4.4 & 2.6 & 2.7 & 2.7 & 2.7 & 0.4 & 0.4 & 0 \end{pmatrix} \end{matrix}.$$

Consider the distorted dissimilarity map δ which is quartet consistent with T . It is easy to see that $Q_\delta(x, y) = -31.2$, while $Q_\delta(a, b) = Q_\delta(m, n) = -30.2$. Therefore the function

Q_δ is not minimized at one of the cherries of T , and the neighbor-joining algorithm applied to δ outputs a tree different from T , in which x, y form a cherry.

Fortunately, there is a single extra condition which ensures that neighbor-joining correctly reconstructs a tree. In what follows we say that a leaf x is *interior to a quartet* $(ij : kl)$ in a tree T if neither $(ik : xl)$ nor $(ik : xj)$ are quartets in T .

Definition 12. A dissimilarity map $\delta : X \times X \rightarrow \mathbb{R}$ is *quartet additive* with a tree T if for every $(ij : kl) \in T$ with x interior to $(ij : kl)$, and y not interior to $(ij : kl)$ such that $(ij : xy)$ is not a quartet, we have $w(kl : xy) > w(ij : xy)$.

We conclude this section by introducing the abbreviation $w_\delta(ij : k)$ for $w_\delta(ij : kl)$ where $k = l$. Note that

$$w_\delta(ij : k) = w(ij : kk) = \delta(i, k) + \delta(j, k) - \delta(i, j).$$

In other words, for a tree metric δ_T , $w_{\delta_T}(ij : k)$ is the length of the path between k and the path between i and j . We can reformulate the neighbor-joining criterion in terms of such path estimates.

Theorem 13. Let δ_T be the tree metric corresponding to a tree T . The pair a, b that maximizes

$$S_{\delta_T}(i, j) = 2\delta_T(i, j) + \sum_{k \neq i, j} w_{\delta_T}(ij : k)$$

is a cherry in T .

The terms $w(ij : k)$, together with the *S-criterion* and the two Lemmas below simplify the proof of our main theorem in the next section.

Lemma 14. If δ is quartet consistent with a tree T , and x is a leaf not interior to the quartet $(ij : kl)$, but $(ij : kx)$ and $(ik : lx)$ are quartets, then $w(ij : x) > w(kl : x)$.

Proof: Note that

$$\delta(i, x) + \delta(j, x) + \delta(k, l) > \delta(i, x) + \delta(j, k) + \delta(x, l) > \delta(i, j) + \delta(x, k) + \delta(x, l)$$

which is equivalent to

$$w(ij : x) > w(kl : x).$$

□

Lemma 15. If δ is quartet additive with a tree T , and $(ij : kl)$ is a quartet in T with x interior to $(ij : kl)$, and y not interior to $(ij : kl)$ such that $(ij : xy)$ is not a quartet then

$$w(kl : x) + w(kl : y) > w(ij : x) + w(kl : y).$$

Proof: Note that

$$\begin{aligned} w(kl : xy) &= \frac{1}{2}(w(kl : x) + w(kl : y)) - \delta(x, y) \\ w(ij : xy) &= \frac{1}{2}(w(ij : x) + w(ij : y)) - \delta(x, y) \end{aligned}$$

so $w(kl : xy) > w(ij : xy)$ is equivalent to $w(kl : x) + w(kl : y) > w(ij : x) + w(ij : y)$. \square

Remark 16. Quartet consistency and additivity are invariant with respect to the shifting operation. We leave this as an exercise to the reader.

4. A CONSISTENCY THEOREM FOR NEIGHBOR-JOINING

Theorem 17. *If $\delta : X \times X \rightarrow \mathbb{R}$ is quartet consistent and quartet additive with a tree T , then neighbor-joining applied to δ will construct a tree with the same topology as T .*

Lemma 18. *The reduction step maintains quartet consistency and additivity with T .*

Proof: Without loss of generality, we may assume that we are collapsing taxa x, y into taxon z by using the first type of reduction step. We need to prove that the new dissimilarity map is quartet consistent and additive with the new tree T' . In other words, for every set of taxa $X = \{z, x_1, \dots, x_k\}$, where $k = 3$ or $k = 5$, we must have $L(z, x_1, \dots, x_k) > 0$, where L is a linear expression depending on the topology of the induced tree on X . But since x and y formed a cherry, the topology of the tree induced on X by T' is the same as that induced by T on $X - \{z\} \cup \{x\}$ and on $X - \{z\} \cup \{y\}$. But by induction, $L(X - \{z\} \cup \{x\}) > 0$ and $L(X - \{z\} \cup \{y\}) > 0$. Since L is linear and we are using the first reduction step, $L(X) = \frac{1}{2}(L(X - \{z\} \cup \{x\}) + L(X - \{z\} \cup \{y\})) > 0$. \square

Proof of Theorem 17: By the above lemma, it suffices to prove that at any step, the pair of taxa which maximize the S -criterion form a cherry. We argue by contradiction. Let us consider a pair δ, T that are a minimum size counterexample to the above statement. Let i, j be the pair of taxa which maximize S_δ , but that do not form a cherry. There are three steps in the proof. First we show that neither i nor j form part of a cherry, then we show that they are separated by exactly three edges, and finally we reach a contradiction by finding a cherry k, l for which $S_\delta(k, l) > S_\delta(i, j)$.

Suppose, without loss of generality, that leaf i forms a cherry with leaf $k \neq j$. Then

$$S_\delta(i, j) = 2\delta(i, j) + w_\delta(ij : k) + \sum_{l \neq i, j, k} w_\delta(ij : l) = \delta(i, j) + \delta(i, k) + \delta(j, k) + \sum_{l \neq i, j, k} w_\delta(ij : l).$$

Similarly,

$$S_\delta(i, k) = 2\delta(i, k) + w_\delta(ik : j) + \sum_{l \neq i, j, k} w_\delta(ik : l) = \delta(i, j) + \delta(i, k) + \delta(j, k) + \sum_{l \neq i, j, k} w_\delta(ik : l).$$

But since i and k form a cherry and δ is quartet consistent, for any leaf l we have that $\delta(l, k) + \delta(i, j) > \delta(i, k) + \delta(j, l)$ so $\delta(l, i) + \delta(l, k) - \delta(i, k) > \delta(l, j) + \delta(l, i) - \delta(i, j) \Rightarrow w_\delta(ik : l) > w_\delta(ij : k)$. Therefore, $S_\delta(i, k) > S_\delta(i, j)$ and i, j cannot be part of a cherry. It follows that the path joining i with j in T has at least 3 edges.

Now suppose that i and j are separated by four edges or more. Let $T_1, \dots, T_k, k \geq 3$ be the subtrees adjacent to vertices on the path between i and j , with T_1 closest to i and T_k closest to j . Note that by the first part of the proof, T_1 and T_k have at least two leaves each. Let T' be the tree T with T_2 removed and δ' be the dissimilarity map δ restricted to the leaves in T' . Since T is a minimum counterexample, $S_{\delta'}$ must be maximized at a

cherry in T' . Since neither i nor j are part of a cherry in T' , the maximizing cherry a, b has to be in one of the subtrees T_1, T_3, \dots, T_k . Note that $S_{\delta'}(a, b) > S_{\delta'}(i, j)$ but

$$\begin{aligned} S_{\delta'}(a, b) &= 2\delta(a, b) + w_{\delta}(ab : i) + w_{\delta}(ab : j) + \sum_{l \neq a, b, i, j; l \notin T_2} w_{\delta}(ab : l) \\ &= w_{\delta}(ab : ij) + 2\delta(a, b) + 2\delta(i, j) + \sum_{l \neq a, b, i, j; l \notin T_2} w_{\delta}(ab : l). \end{aligned}$$

Similarly,

$$\begin{aligned} S_{\delta'}(i, j) &= 2\delta(i, j) + w_{\delta}(ij : a) + w_{\delta}(ij : b) + \sum_{l \neq a, b, i, j; l \notin T_2} w_{\delta}(ij : l) \\ &= w_{\delta}(ab : ij) + 2\delta(a, b) + 2\delta(i, j) + \sum_{l \neq a, b, i, j; l \notin T_2} w_{\delta}(ij : l). \end{aligned}$$

Therefore

$$S_{\delta'}(a, b) - S_{\delta'}(i, j) = \sum_{l \neq a, b, i, j; l \notin T_2} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) > 0.$$

But

$$S_{\delta}(a, b) - S_{\delta}(i, j) = \sum_{l \neq a, b, i, j} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) < 0,$$

and thus it must be that

$$\sum_{l \in T_2} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) < 0.$$

By Lemma 14, $w_{\delta}(ab : l) > w_{\delta}(ij : l)$ for any $l \in T_2$ and we have a contradiction. It follows that $k = 2$ and the path between i and j has exactly three edges.

Suppose, without loss of generality, that T_2 has at least as many leaves as T_1 . Let a, b be a cherry in T_1 . Set $s = |T_1| - 2$, and group the leaves of T_2 into two sets: S_1 of cardinality s and S_2 of cardinality $|T_2| - s$. There exists a bijection f between the leaves of $T_1 - \{a, b\}$ and S_1 . As before, we have

$$\begin{aligned} S_{\delta}(a, b) - S_{\delta}(i, j) &= \sum_{l \neq a, b, i, j} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) \\ &= \sum_{l \in T_1; l \neq a, b} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) + \sum_{l \in S_1} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) \\ &\quad + \sum_{l \in S_2} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) \\ &= \sum_{l \in T_1; l \neq a, b} (w_{\delta}(ab : l) - w_{\delta}(ij : l)) + (w_{\delta}(ab : f(l)) - w_{\delta}(ij : f(l))) \\ &\quad + \sum_{l \in S_2} (w_{\delta}(ab : l) - w_{\delta}(ij : l)). \end{aligned}$$

The first sum is positive by Lemma 15, and the second is positive by Lemma 14. Therefore $S_\delta(a, b) > S_\delta(i, j)$ and we have a contradiction. \square

Remark 19. Theorem 9 follows from the observation that quartet consistency suffices in the proof of Theorem 17 when $4 \leq |X| \leq 7$.

Corollary 20. *Neighbor-joining has l_∞ radius $\frac{1}{2}$.*

Proof: It suffices to show that if δ_T is a tree metric and δ is a metric with $\max_{i,j} |\delta_T(i, j) - \delta(i, j)| < \frac{1}{2} \min_{e \in E(T)} l(e)$ where $l(e)$ is the length of edge e in T , then δ is quartet consistent and quartet additive with T . \square

The next corollary extends the visibility lemma from [4]. A cherry a, b is *visible* from a with respect to a dissimilarity map δ if

$$b = \operatorname{argmax}_{x \neq a} Z_\delta(a, x).$$

Corollary 21. *If δ is quartet consistent with a tree T and a is a leaf with a, b a cherry, then (a, b) is visible from a with respect to δ .*

Proof: This follows directly from the first step of the proof of Theorem 17. \square

The visibility lemma is the key to developing a fast neighbor joining algorithm (FNJ) that has optimal run time complexity $O(n^2)$ [4]. In fact, we can conclude that

Corollary 22. *If δ is quartet consistent and quartet additive with respect to a tree T then FNJ will reconstruct T from δ with optimal run time complexity $O(n^2)$.*

5. THE EDGE RADIUS OF THE NEIGHBOR-JOINING ALGORITHM

In this section we prove a strengthening of a conjecture proposed by Atteson in [1] about the edge radius of the neighbor-joining algorithm. The methods used are very similar to the ones used in the previous section. For ease of exposition, in what follows we will drop the requirement that the input trees are binary. We also continue to allow negative leaf edges. Note that an internal node of degree at least 3 corresponds to an internal edge of length 0 in a binary tree. We first state our theorem in the form of Atteson's conjecture.

Theorem 23. *Let T be a tree with associated tree metric δ_T , and let e be an edge of T of length $l(e)$. If δ is a tree metric whose l_∞ distance to δ_T is less than $\frac{l(e)}{4}$, then neighbor-joining applied to δ will reconstruct the edge e correctly, i.e., the tree T' output by neighbor-joining will contain an edge e' which induces the same split in the tree T' as e induces in T .*

Since the necessary l_∞ error bound needed for the neighbor-joining algorithm to reconstruct an edge correctly is $\frac{1}{4}$ of the length of the edge, we say that the edge radius of neighbor-joining is $\frac{1}{4}$. This result is optimal. In [1], Atteson presents an example where the theorem fails for l_∞ error larger than $\frac{l(e)}{4}$.

Definition 24. Let T be a tree and e a non-leaf edge of length l in T , corresponding to the split $P|Q$ of X . We say that a dissimilarity map δ is $P|Q$ -additive with respect to the tree metric δ_T if the following conditions hold

- $\delta(x, y) - \delta_T(x, y) < \frac{l}{4}$ for all pairs $x, y \in P$ and all pairs $x, y \in Q$,
- $|\delta(x, y) - \delta_T(x, y)| < \frac{l}{4}$ for all pairs $x \in P$ and $y \in Q$.

We are ready to state the main theorem of the section

Theorem 25. *If δ is a $P|Q$ -additive dissimilarity map with respect to the tree T and the tree metric δ_T , then the neighbor-joining algorithm applied to δ will output a tree T' which contains $P|Q$ among its edge-induced splits.*

This is a strengthening of Atteson's initial conjecture since we are not imposing a lower bound on the estimated distances between taxa situated on the same side of the split.

Proof of Theorem 25: The proof proceeds similarly to the proof in of Theorem 17. We define $\tilde{\delta} = \delta - \delta_T$, in other words $\tilde{\delta}(x, y)$ is the error in the $\delta(x, y)$ estimate of $\delta_T(x, y)$.

Lemma 26. *Given $P|Q$ -additive dissimilarity map δ , with respect to a tree T , collapsing a pair of taxa $x, y \in P$ will result in a metric δ' which is $P'|Q$ -additive with respect to a tree T' . Here P' is the set of taxa obtained by replacing x, y by the collapsed node z in P .*

Proof: As before, we assume without loss of generality that we are using the first variant of the reduction step. Now shift the new dissimilarity map δ' by $\frac{\delta_T(x, y)}{2}$ around the new taxon z . Again, this can be done without affecting the outcome of the algorithm. In effect, this is equivalent to defining the distances with respect to z in the following manner:

$$\delta'(z, a) = \frac{1}{2}(\delta(x, a) + \delta(y, a) - \delta_T(x, y)).$$

Now let e be the edge in T corresponding to the $P|Q$ split. Let l be its length. Let p be the path in T that joins x and y . Then $e \notin p$ and let o be the internal point of T where a path from e reaches p . Consider the new tree T' where the taxa x and y are removed and the new taxon z is placed exactly at the internal point o . $\delta_{T'}$ and P' are defined in the obvious way and δ' is defined by collapsing x, y in δ according to the reduction described above.

Let $T_0 \dots T_k$ be the subtrees of T hanging off the path p and let T_0 be the one that contains e . We now observe that $\delta'(a, b) = \delta(a, b)$ and $\delta_{T'}(a, b) = \delta_T(a, b)$ for $a, b \neq z$. In this case the errors remain the same. For $a \in Q$, and therefore $a \in T_0$, we observe that

$$\delta_{T'}(z, a) = \delta_T(o, a) = \frac{1}{2}(\delta_T(x, a) + \delta_T(y, a) - \delta_T(x, y)).$$

Therefore

$$|\delta'(z, a) - \delta_{T'}(a, z)| = |(\delta(x, a) - \delta_T(x, a)) + (\delta(y, a) - \delta_T(y, a))| \frac{1}{2} < \frac{l}{4}.$$

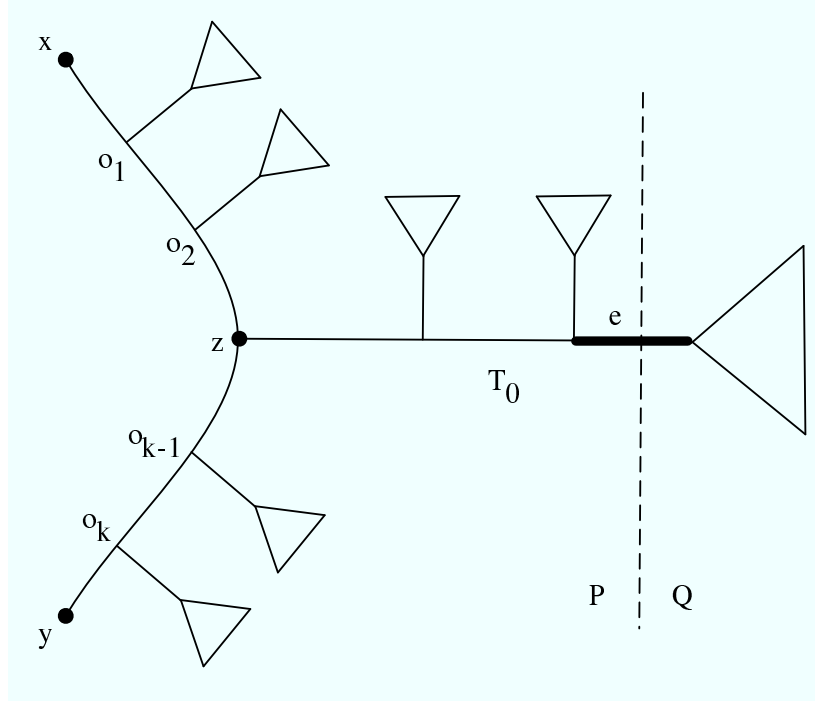


FIGURE 3. The collapsing lemma.

Now for all i let o_i be the point of the path p that is the root of the subtree T_i . Then for any $a \in P$, so $a \in T_i$ for $i \neq 0$ or $a \in T_0 \cap P$,

$$\begin{aligned} \delta'(z, a) &= \delta_{T'}(a, o_i) + [(\delta(x, a) - \delta_T(x, a)) + (\delta(y, a) - \delta_T(y, a))]/2 = \\ \delta_{T'}(a, z) &+ [(\delta(x, a) - \delta_T(x, a)) + (\delta(y, a) - \delta_T(y, a))]/2 - \delta_T(o, o_i) < \\ &\delta_{T'}(a, z) + \frac{l}{4} - \delta_T(o, o_i). \end{aligned}$$

Therefore δ' is $P|Q$ -additive with respect to T' and T' "contains" the edge e . \square

By the previous lemma, it is enough to prove that the first step in the neighbor-joining algorithm, when applied to a $P|Q$ -additive dissimilarity map δ , is to collapse two taxa on the same side of the split $P|Q$. In other words, the criterion $S(\cdot, \cdot)$ cannot be maximized at some pair $x \in P$ and $y \in Q$.

Suppose that we have an instance of a dissimilarity map δ which is $P|Q$ additive with respect to T and δ_T . As before, we let e be the edge inducing $P|Q$ in T and we let l be its length. Suppose that the criterion $S(\cdot, \cdot)$ is maximized at the pair x, y with $x \in P$ and $y \in Q$. For ease of exposition, we define the dissimilarity map $\tilde{\delta} = \delta - \delta_T$. In other words $\tilde{\delta}(a, b)$ is the error in the estimate $\delta(a, b)$ of the true distance $\delta_T(a, b)$.

Let T_1, \dots, T_k be the subtrees of T that hang off of the path p , enumerated in order from the one closest to x to the one closest to y . Let T_j and T_{j+1} be the two trees that anchor at the endpoints of e . Let us call these endpoints u and v . Now consider the operation

of removing the subtrees T_i , for $i \neq j, j + 1$ from their "roots" and attaching them at the nodes u for $i < j$ and v for $i > j + 1$. We maintain the same edge lengths otherwise and we call the new tree T' . We define the dissimilarity map $\delta'(a, b) = \delta_{T'}(a, b) + \tilde{\delta}(a, b)$ for all a, b . In other words we maintain the same errors and change the underlying tree by re-grafting the subtrees hanging off the path p to the endpoints of e . It is trivial to see that δ' is still $P|Q$ additive with respect to e and T' . Moreover, $S_{\delta'}(x, y) = S_{\delta}(x, y)$ and $S_{\delta'}(a, b) \leq S_{\delta}(a, b)$ for $a, b \in P$ or $a, b \in Q$. This is true because the distance between taxa in different T_i 's can only be reduced and moreover $w(a, b : c)$ can also only be reduced for $a, b \in P$ or $a, b \in Q$. Trivially then, proving that there is a pair $a, b \in P$ or $a, b \in Q$ such that $S_{\delta'}(a, b) \geq S_{\delta'}(x, y)$ implies that for the same pair a, b , $S_{\delta}(a, b) \geq S_{\delta}(x, y)$. We can therefore replace T and δ with T' and δ' .

We can now assume that p consists of only three edges, namely e and the two leaf edges corresponding to x and y . We let T_x and T_y be the subtrees of T hanging off the path p at the endpoints of e on x 's and y 's side respectively.

We now proceed in the same manner as in the previous section.

Lemma 27. (*Modified four point condition*) *Suppose we have $a, i \in P$ and $b, j \in Q$. Then $w_{\delta}(a, i : b) > w_{\delta}(i, j : b)$.*

Proof: This is a trivial verification based on the fact that the middle edge of the quartet $(a, i : b, j)$ is at least as long as e . \square

Corollary 28. *Neither of the vertices x and y can be part of a cherry.*

Proof: Suppose x is in a cherry with vertex a . Then $P = \{a, x\}$. The statement follows from applying the previous lemma to the quartet $(a, x : y, j)$ for all $j \in Q - \{y\}$ and summing. \square

Lemma 29. *Let $a, b \in T_x$ and $c \in T_y$. Then $w_{\delta}(a, b : c) > w_{\delta}(x, y : c)$.*

Proof:

$$\begin{aligned} w_{\delta}(a, b : c) - w_{\delta}(x, y : c) &= \delta(a, c) + \delta(b, c) - \delta(a, b) - \delta(x, c) - \delta(y, c) + \delta(x, y) = \\ &= w_{\delta_T}(a, b : c) - w_{\delta_T}(x, y : c) + \tilde{\delta}(a, c) + \tilde{\delta}(b, c) - \tilde{\delta}(a, b) - \tilde{\delta}(x, c) - \tilde{\delta}(y, c) + \tilde{\delta}(x, y) > \\ &> 2l - 6\left(\frac{l}{4}\right) = \frac{l}{2} > 0 \end{aligned}$$

\square

Corollary 30. *T_x and T_y have at least three leaves each.*

Proof: Suppose T_x has less than three leaves. It cannot have only one by the previous corollary. Then T_x has exactly two leaves, a, b . Applying the previous lemma to a, b and all c 's in T_y and summing up gives $S_{\delta}(a, b) \geq S_{\delta}(x, y)$, a contradiction. \square

Now let us suppose that T_x has fewer leaves than T_y . Let s be the number of leaves in T_x . As before, we split the edges of T_y into a set Y' of $s - 2$ leaves and $Y'' = l(T_y) - Y'$.

Let $a, b, c \in T_x$ and $d \in Y'$. Then

$$(1) \quad \begin{aligned} & 2w_\delta(a, b : c, d) - 2w_\delta(x, y : c, d) = \\ & 2w_{\delta_T}(a, b : c, d) - 2w_{\delta_T}(x, y : c, d) + 2w_{\tilde{\delta}}(a, b : c, d) - 2w_{\tilde{\delta}}(x, y : c, d) = \\ & 2w_{\delta_T}(a, b : c, d) + 2l + 2w_{\tilde{\delta}}(a, b : c, d) - 2w_{\tilde{\delta}}(x, y : c, d) \end{aligned}$$

But

$$\begin{aligned} 2w_{\tilde{\delta}}(a, b : c, d) &= \tilde{\delta}(a, c) + \tilde{\delta}(b, c) + \tilde{\delta}(a, d) + \tilde{\delta}(b, d) - 2\tilde{\delta}(a, b) - 2\tilde{\delta}(c, d) > \\ & \tilde{\delta}(a, c) + \tilde{\delta}(b, c) - 2\tilde{\delta}(a, b) - 2\tilde{\delta}(c, d) - \frac{l}{2} \end{aligned}$$

and

$$\begin{aligned} 2w_{\tilde{\delta}}(x, y : c, d) &= \tilde{\delta}(x, c) + \tilde{\delta}(y, c) + \tilde{\delta}(x, d) + \tilde{\delta}(y, d) - 2\tilde{\delta}(x, y) - 2\tilde{\delta}(c, d) < \\ & \frac{3}{2}l - 2\tilde{\delta}(c, d). \end{aligned}$$

Subtracting the two inequalities above we obtain

$$(2) \quad 2w_\delta(a, b : c, d) - 2w_\delta(x, y : c, d) > 2w_{\delta_T}(a, b : c, d) + \tilde{\delta}(a, c) + \tilde{\delta}(b, c) - 2\tilde{\delta}(a, b)$$

But notice that

$$w_{\delta_T}(a, b : c, d) + w_{\delta_T}(a, c : b, d) + w_{\delta_T}(b, c : a, d) = 0,$$

therefore by summing 2 for permutations of $\{a, b, c\}$, we obtain

$$w_\delta(a, b : c, d) + w_\delta(a, c : b, d) + w_\delta(b, c : a, d) > w_\delta(x, y : c, d) + w_\delta(x, y : b, d) + w_\delta(x, y : a, d).$$

Summing over all possibilities of $a, b, c \in T_x$ and $d \in Y'$ yields

$$(3) \quad \sum_{a, b, c \in T_x, d \in Y'} w_\delta(a, b : c, d) - w_\delta(x, y : c, d) > 0.$$

But remember that

$$2w_\delta(a, b : c, d) = w_\delta(a, b : c) + w_\delta(a, b : d) - 2\delta(c, d),$$

so

$$w_\delta(a, b : c, d) - w_\delta(x, y : c, d) = \frac{1}{2}(w_\delta(a, b : c) + w_\delta(a, b : d) - w_\delta(x, y : c) - w_\delta(x, y : d)).$$

By rearranging in the sum above, the following holds

$$\begin{aligned} & \sum_{a, b, c \in T_x, d \in Y'} w_\delta(a, b : c, d) - w_\delta(x, y : c, d) = \\ & \sum_{a, b \in T_x} \sum_{c \in T_x - \{a, b\}} \sum_{d \in Y'} [w_\delta(a, b : c) - w_\delta(x, y : c) + w_\delta(a, b : d) - w_\delta(x, y : d)] = \\ & = (s-2) \sum_{a, b \in T_x} \sum_{c \in T_x \cup Y' - \{a, b\}} [w_\delta(a, b : c) - w_\delta(x, y : c)] > 0. \end{aligned}$$

Moreover, by Lemma 29, for all pairs $a, b \in T_x$ and $c \in Y''$,

$$\sum_{a,b \in T_x} \sum_{c \in Y''} [w_\delta(a, b : c) - w_\delta(x, y : c)] > 0.$$

Summing the above two relations over all pairs a, b gives

$$\begin{aligned} \sum_{a,b \in T_x} \sum_{c \in T_x \cup T_y - \{a,b\}} [w_\delta(a, b : c) - w_\delta(x, y : c)] = \\ \sum_{a,b \in T_x} [S_\delta(a, b) - S_\delta(x, y)] > 0. \end{aligned}$$

Thus the average of the S_δ criterion for all pairs in T_x is larger than $S_\delta(x, y)$, so there is at least one pair $a, b \in T_x$ such that $S_\delta(a, b) > S_\delta(x, y)$. Therefore the neighbor-joining algorithm will collapse a, b before x, y . We conclude that for a $P|Q$ -additive δ , the neighbor-joining algorithm will always collapse a pair of taxa on the same side of $P|Q$. By Lemma 26, this is enough to prove that the split $P|Q$ is preserved in the output of neighbor-joining, which concludes the proof of our theorem. \square

We conclude by stating that our analysis holds trivially for the fast neighbor-joining algorithm of [4]. This follows from the observation that for a $P|Q$ -additive dissimilarity map δ , no pair x, y with $x \in P$ and $y \in Q$ can maximize the $S(\cdot, \cdot)$ criterion, and therefore the maximizing pair, which has to be visible from both of its members, has both taxa on the same side of the partition $P|Q$.

Corollary 31. *If δ is a $P|Q$ -additive dissimilarity map with respect to a tree T , then FNJ applied to δ will reconstruct a tree T' which contains $P|Q$ among its set of edge-induced splits.*

6. ACKNOWLEDGMENTS

Radu Mihaescu was supported by the Fannie and John Hertz Foundation. Lior Pachter was partially supported by NIH grant R01HG2362 and NSF grant CCF0347992. Dan Levy was supported by NIH grant GM68423.

REFERENCES

1. K Atteson, *The performance of neighbor-joining methods of phylogenetic reconstruction*, *Algorithmica* **25** (1999), 251–278.
2. WJ Bruno, ND Socci, and AL Halpern, *Weighted neighbor-joining: a likelihood-based approach to distance-based phylogeny reconstruction*, *Molecular Biology and Evolution* **17** (2000), no. 1, 189–197.
3. D Bryant, *On the uniqueness of the selection criterion in neighbor-joining*, *Journal of Classification* **22** (2005), no. 1, 3–15.
4. I Elias and J Lagergren, *Fast neighbor joining*, *Proceedings of the International Colloquium on Automata, Languages and Programming (ICALP '05)*, 2005.
5. PL Erdős, MA Steel, LA Székely, and TJ Warnow, *A few logs suffice to build (almost) all trees. I*, *Random Structures and Algorithms* **14** (1999), no. 2, 153–184.
6. JS Farris, VA Albert, M Källersjö, D Lipscomb, and AG Kluge, *Parsimony jackknifing outperforms neighbor-joining*, *Cladistics* **12** (1996), 99–124.
7. O Gascuel, *BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data*, *Molecular Biology and Evolution* **14** (1997), no. 7, 685–695.
8. BG Hall, *Comparison of the accuracies of several phylogenetic methods using protein and DNA sequences*, *Molecular Biology and Evolution* **22** (2005), no. 3, 792–802.
9. J Huelsenbeck and D Hillis, *Success of phylogenetic methods in the four-taxon case*, *Systematic Biology* **42** (1993), no. 3, 247–264.
10. K St. John, T Warnow, B Moret, and L Vawter, *Performance study of phylogenetic methods: (unweighted) quartet methods and neighbor joining*, *Journal of Algorithms* **48** (2003), 174–193.
11. MK Kuhner and J Felsenstein, *A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates*, *Molecular Biology and Evolution* **11** (1994), no. 459–468.
12. S Kumar and SR Gadagker, *Efficiency of the neighbor-joining method in reconstructing evolutionary relationships in large phylogenies*, *Journal of Molecular Evolution* **51** (2000), 544–553.
13. D Levy, R Yoshida, and L Pachter, *Beyond pairwise distances: neighbor joining with phylogenetic diversity estimates*, *Molecular Biology and Evolution*, in press (2006).
14. GJ Olsen, H Matsuda, R Hagstrom, and R Overbeek, *fastDNAm1: A tool for construction of phylogenetic trees of DNA sequences using maximum likelihood*, *Computational Applied Biosciences* **10** (1994), 41–48.
15. S Ota and WH Li, *NJML: A hybrid algorithm for the neighbor-joining and maximum likelihood methods*, *Molecular Biology and Evolution* **17** (2000), no. 9, 1401–1409.
16. V Ranwez and O Gascuel, *Improvement of distance-based phylogenetic methods by a local maximum likelihood approach using triplets*, *Molecular Biology and Evolution* **19** (2002), no. 11, 1952–1963.
17. N Saitou and M Nei, *The neighbor joining method: a new method for reconstructing phylogenetic trees*, *Molecular Biology and Evolution* **4** (1987), no. 4, 406–425.
18. K Strimmer and A von Haeseler, *Quartet puzzling: A quartet maximum likelihood method for reconstructing tree topologies*, *Molecular Biology and Evolution* **13** (1996), 964–969.
19. JA Studier and KJ Keppler, *A note on the neighbor-joining method of saitou and nei*, *Molecular Biology and Evolution* **5** (1988), 729–731.
20. K Tamura, M Nei, and S Kumar, *Prospects for inferring very large phylogenies by using the neighbor-joining method*, *Proceedings of the National Academy of Sciences* **101** (2004), 11030–11035.

DEPARTMENT OF MATHEMATICS, UC BERKELEY

E-mail address: {mihaescu,levyd,lpachter}@math.berkeley.edu