

Finding a Maximum Likelihood Tree Is Hard

BENNY CHOR AND TAMIR TULLER

Tel-Aviv University, Tel-Aviv, Israel

Abstract. Maximum likelihood (ML) is an increasingly popular optimality criterion for selecting evolutionary trees [Felsenstein 1981]. Finding optimal ML trees appears to be a very hard computational task, but for tractable cases, ML is the method of choice. In particular, algorithms and heuristics for ML take longer to run than algorithms and heuristics for the second major character based criterion, maximum parsimony (MP). However, while MP has been known to be NP-complete for over 20 years [Foulds and Graham, 1982; Day et al. 1986], such a hardness result for ML has so far eluded researchers in the field.

An important work by Tuffley and Steel [1997] proves quantitative relations between the parsimony values of given sequences and the corresponding log likelihood values. However, a direct application of their work would only give an *exponential time* reduction from MP to ML. Another step in this direction has recently been made by Addario-Berry et al. [2004], who proved that *ancestral maximum likelihood* (AML) is NP-complete. AML “lies in between” the two problems, having some properties of MP and some properties of ML. Still, the AML proof is not directly applicable to the ML problem.

We resolve the question, showing that “regular” ML on phylogenetic trees is indeed intractable. Our reduction follows the vertex cover reductions for MP [Day et al. 1986] and AML [Addario-Berry et al. 2004], but its starting point is an approximation version of vertex cover, known as GAP VC. The crux of our work is not the reduction, but its correctness proof. The proof goes through a series of tree modifications, while controlling the likelihood losses at each step, using the bounds of Tuffley and Steel [1997]. The proof can be viewed as correlating the value of any ML solution to an arbitrarily close approximation to vertex cover.

Categories and Subject Descriptors: F.2.m [Analysis of Algorithms and Problem Complexity]: Miscellaneous

General Terms: Theory, Algorithms

Additional Key Words and Phrases: Maximum likelihood, tree reconstruction, maximum parsimony, intractability, approximate vertex cover

This research was supported by ISF grant 418/00.

An extended abstract of this work was submitted to RECOMB2005 on October 29, 2004, and published in the conference proceedings on May 14, 2005.

Authors’ address: School of Computer Science, Tel-Aviv University, Ramat-Aviv 69978, Israel, e-mail: {bchor,tamirtul}@post.tau.ac.il.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701, USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2006 ACM 0004-5411/06/0900-0722 \$5.00

1. Background

Molecular data, and even complete genomes, are being sequenced at an increasing pace. This newly accumulated information should make it possible to resolve long standing questions in evolution, such as reconstructing the phylogenetic tree of placental mammals and estimating the times of species divergence. The analysis of this data flood requires sophisticated mathematical tools and algorithmic techniques. Most tree reconstruction methods are either distance-based or character-based. Of the later, two are widely used in practice: MP (*maximum parsimony* [Fitch 1971]) and ML (*maximum likelihood* [Felsenstein 1981]). It is known that ML is *consistent*, namely with high probability, for long enough input sequences, the correct tree is the tree maximizing the likelihood [Felsenstein 2004, ch. 16]. Consistency does not hold for MP, and in fact for certain families of trees (the so-called *Felsenstein zone* [Felsenstein 1996]). MP will reconstruct the *wrong* trees, even for arbitrarily long input sequences. The two methods are known to be computationally intensive, and exact algorithms are limited to just about 20 sequences or less. This forces practitioners to resort to heuristics. For both exact algorithms and heuristics, ML seems a *harder* problem than MP.

In the absence of concrete lower-bound techniques, the major tool for demonstrating computational intractability remains NP hardness proofs. Both MP and ML have well-defined objective functions, and the related decision problems (or at least discretized versions of them) are in the complexity class NP. It has been known for over 20 years that MP is NP-complete [Foulds and Graham 1982; Day et al. 1986; Day 1987; Day and Sankoff 1986; Steel 1992], (see also Wareham [1993] and references). The proof of Day et al. [1986] employs an elegant reduction from vertex cover (VC). However, no such result has been found for ML to date. This is particularly frustrating in light of the intuition among practitioners that ML is harder than MP.

Tuffley and Steel [1997] have investigated the quantitative relations between MP and ML. In particular, they showed that if the m sequences are padded with sufficiently many zeroes, the ML and MP trees coincide. Since parsimony is invariant under padding by zeroes, this approach could in principle lead to a reduction from MP to ML. Unfortunately, the upper bound provided in Tuffley and Steel [1997] on the padding length is *exponential* in m . A step in a different direction was taken by Addario-Berry et al. [2004]. They studied the complexity of AML (ANCESTRAL MAXIMUM LIKELIHOOD) [Koshi and Goldstein 1996; Yang et al. 1995]. This variant of ML is “between” MP and ML in that it is a likelihood method (like ML) but it reconstructs sequences for internal vertices (like MP). They showed that AML is NP-complete, using a reduction from (exact) VERTEX COVER.

Our NP hardness proof of ML uses ingredients from both Tuffley and Steel [1997] and Addario-Berry et al. [2004], as well as new insights on the behavior of the likelihood function on trees. The reduction itself is essentially identical to that given for MP by Day et al. [1986], and also used in the AML paper Addario-Berry et al. [2004]. However, our starting point is not *exact* VC but the *gap* version of it [Berman and Karpinski 1999; Karpinski 2001]. The proof of correctness for this reduction relative to ML is different, and substantially more involved. We define a family of *canonical trees*. Every such tree is associated with a unique cover in the original graph. We show that if L is the likelihood of the canonical tree, n is the number of vertices in the original graph, m is the number of edges in the original

graph, and c is the size of the associated cover, then as $n \rightarrow \infty$,

$$\frac{-\log(L)}{(m+c)\log(n)} \rightarrow 1.$$

In particular, this gives an inverse relation between likelihood and cover size: Larger L implies smaller c , and vice-versa.

When proving the correctness of the reduction, we want to establish two directions: (\Rightarrow) If the original graph has a small cover, then there is a tree with high likelihood, and (\Leftarrow) that the existence of a tree with high likelihood implies the existence of a small cover. The first direction is easy, using the canonical tree related to the small vertex cover. It is the other direction that is hard, because there is no obvious relation between the log likelihood of a *noncanonical* tree and the size of any cover. What we do, starting from any ML tree, is to apply a sequence of modifications that leads it to a *canonical tree*. The whole series of modifications may actually *decrease* the likelihood of the resulting, canonical tree vs. the original, ML one. We use the techniques of Tuffley and Steel [1997] to infer likelihood properties from parsimony ones. In particular, we combine Tuffley and Steel [1997] and the degree bound of the original graph to show that in every step, the log likelihood decreases by at most $O(\log n)$ bits. Finally, we show that the total number of modifications is not too large—at most $n/\log \log n$. This allows us to show that the overall loss in log likelihood is at most $O(n \log n / \log \log n)$. We also show that the log likelihood of the final, canonical tree is $\theta(n \log n)$. This implies the ratio of the log likelihood of the last, canonical tree, and the log likelihood of the ML tree, approaches 1 as $n \rightarrow \infty$. This proves that log ML is tightly related to an approximate vertex cover, establishing the NP hardness of ML.

2. Overview of Proof

In this section we give a high level description of the hardness proof. The reduction is from the GAP VERTEX COVER problem on graphs whose degree is at most three, a problem proved NP-hard in Berman and Karpinski [1999] and Karpinski [2001].

Given an undirected graph $G = (V, E)$ of max degree 3 with $n = |V|$ nodes and $m = |E| \leq 1.5n$ edges, we construct an ML instance, consisting of $m + 1$ binary strings of length n . The ML problem is to find a tree with the $m + 1$ sequences at its leaves, and an assignment of substitution probabilities to the edges of that tree (edge lengths), such that the likelihood of generating the given sequences is maximized. The proof relates the approximate max log likelihood value to the size of a vertex cover in G . This approximation is tight enough to enable the solution of the original gap problem. Our reduction follows the one for maximum parsimony given by Day et al. [1986] and for ancestral ML, given by Addario-Berry et al. [2004]. Both reductions were from the (exact) VERTEX COVER problem. In this reduction, we generate one string, denote 0^n consisting entirely of zeros, and m strings, called *edge strings* that contain exactly two ones, and naturally encodes an edge.

Consider all unrooted weighted trees with $m + 1$ leaves that have the given sequences at their leaves. We say that such a tree is in *canonical form* if the following properties hold (see Figure 1):

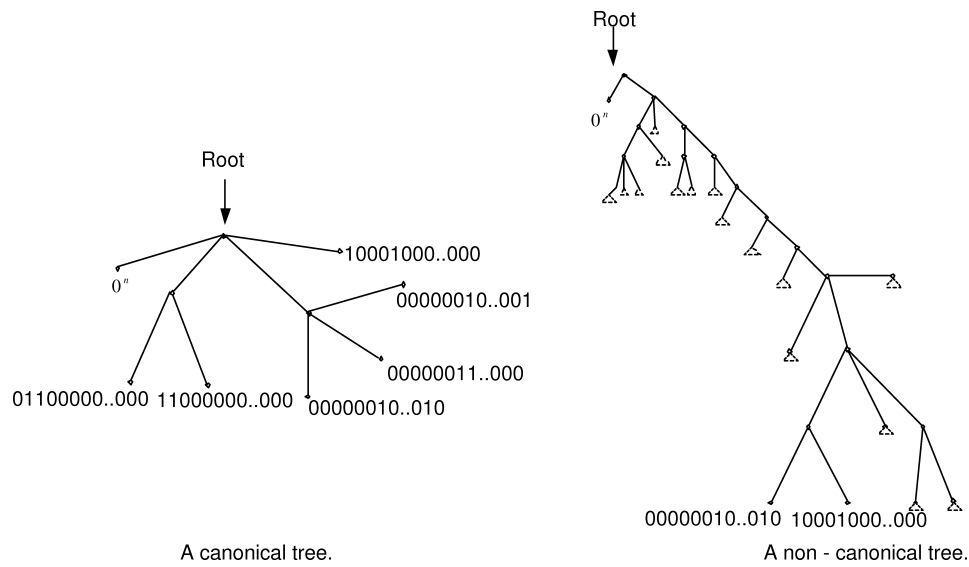


FIG. 1. Canonical (left) and noncanonical (right) trees.

Definition 2.1 (Canonical Form).

- (1) The tree has an internal node, the root, such that the leaf with label 0^n is a child of the root, and the edge connecting that leaf to the root has length 0.
- (2) The tree has height at most two (namely all leaves are one or two tree edges away from the root).
- (3) All internal nodes except the root have degree two or three.
- (4) For every internal vertex v , except the root, there is a position such that the labels for all the children of v have a 1 in that position.

Canonical trees uniquely define a vertex cover, where the leaves of each subtree corresponds to one, two, or three original edges that are covered by one node. (For the subtrees with one leaf, the covering vertex can correspond to either end point, while for size two and three subtrees, the covering vertex is uniquely defined.) Consequently, given a tree in canonical form, we can quantify the size of the corresponding vertex cover of the original graph. The reason we force the root to be connected to the all zero leaf with an edge of weight 0 is that this way the root itself is “forced” to have label 0^n in all trees with nonzero likelihood. This enables us to express the likelihood of a canonical form tree as a product of the likelihoods of its subtrees. In particular, there is no influence, or dependency, between different subtrees.

The major part of the proof is showing that given any ML tree, T_{ML} , with the given $m + 1$ “reduction sequences” at its leaves, there is a series of local modifications on trees with the given sequences at their leaves, such that in each modification the log likelihood of the resulting tree is decreased by at most $O(\log n)$ per step, and the final tree, T_{CA} , is in canonical form. The number of modifications is $o(n)$, which is small enough to establish a tight ratio $1 - o(1)$ between the max log likelihood and the log likelihood of the final, canonical tree. In each step, we transform one tree to another. We identify a small subforest, containing between $\log \log n$ and $2 \log \log n$ leaves.

Such a subforest is a union of disjoint subtrees that hang off a common internal node (not the “root”). Using the bound on the degree of the original graph, we show that the parsimony score of this subforest when its root is labeled by the all zero string can be worse by at most a constant $B < 8^3$ than the score with any other root labeling. Using the results of Tuffley and Steel [1997], and the small size of the subtree, it is possible to unroot this subforest, rearrange it, and connect it directly to the root in a “canonical way,” such that the overall log likelihood of the whole tree decreases by at most $B \log n + o(\log n)$. Over the series of $m/\log \log n$ modifications, the overall decrease is at most $Bn \log n/\log \log n + o(n \log n/\log \log n) = O(n \log n/\log \log n)$. We show that the log likelihood of the final canonical tree, T_{CA} , is $\theta(n \log n)$. This is sufficiently large to show that despite such decrease,

$$\frac{\log Pr(S|T_{CA})}{\log Pr(S|T_{ML})} = 1 - o(1).$$

Every tree in canonical form naturally corresponds to a vertex cover in the original graph. The tight relation between $\log Pr(S|T_{ML})$ and $\log Pr(S|T_{CA})$ implies a tight relationship between the size of an approximate vertex cover in the original graph and the maximum likelihood tree on the given sequences, and establishes the NP hardness of maximum likelihood on phylogenetic trees.

3. Model, Definitions and Notations

In this section, we describe the model and basic definitions that we will use later. These definitions include phylogenetic trees and characters, the parsimony score, Neyman’s two state model, and the likelihood function. In most of this paper, we assume that characters are in one of two states, 0 or 1. Let $S = [s(1), s(2), s(3), \dots, s(m)] \in \{0, 1\}^{m \times n}$ be the observed sequences of length n over m taxa (m leaves). Given such sequences, both the maximum parsimony and the maximum likelihood criteria aim at finding the tree (or trees) that “best explain” this data. Each uses a different objective function. In this section, both are defined and explained.

Definition 3.1 (Subtrees and Subforests). Let $T = (V(T), E(T))$ be a rooted tree. We assume that each nonroot internal vertex of T has a degree at least three. The root has degree 2 or higher. For a nonroot internal node v , the node on the unique path from v to the root that is adjacent to v is called v ’s father. The other (two or more) nodes adjacent to v are called v ’s children. Let u_1, \dots, u_k be these children, and let T_1, \dots, T_k be the subtrees rooted at them, respectively. The union of subset of T_1, \dots, T_k together with v is called the subtree rooted at v . If $\{i_1, \dots, i_a\}$ and $\{j_1, \dots, j_b\}$ are two disjoint subsets of $\{1, \dots, k\}$, we say that the subforest $\{T_{i_1}, \dots, T_{i_a}\}$ and $\{T_{j_1}, \dots, T_{j_b}\}$ rooted at v are *disjoint*.

Definition 3.2 (Phylogenetic Trees, Characters, Labellings Tuffley and Steel 1997). A phylogenetic tree with m leaves is a tree $T = (V(T), E(T))$ such that each leaf (degree one vertex) is given a unique label from $[m] = \{1, \dots, m\}$. A function $\lambda : [m] \rightarrow \{0, 1\}$ is called a state function for T . A function $\hat{\lambda} : V(T) \rightarrow \{0, 1\}$ is called an extension of λ on T if it coincides with λ on the leaves of T . In a similar way, we define a function $\lambda^n : [m] \mapsto \{0, 1\}^n$ and an extension

$\hat{\lambda}^n : V(T) \mapsto \{0, 1\}^n$. This later function is called a labelling of T . If $\hat{\lambda}^n(v) = s$ we say that the string s is the labelling of the vertex v . Given a labelling $\hat{\lambda}^n$, and an edge $e = (u, v)$ let $d_e(\hat{\lambda}^n)$ denote the number of positions at which $\hat{\lambda}^n(u)$ and $\hat{\lambda}^n(v)$ differ.

Definition 3.3 (Maximum Parsimony Score). Let T be a phylogenetic tree with m leaves, and S be a set of m binary strings, all of length n . Let $\lambda_{pars}^n : [m] \rightarrow \{0, 1\}^n$ be the function taking i to the i th string S_i in S . Let $\hat{\lambda}_{pars}^n : V(T) \mapsto \{0, 1\}^n$ be an extension of λ_{pars}^n that minimizes the expression $\sum_{e \in E(T)} d_e(\hat{\lambda}^n)$. We define $pars(S, T, \lambda^n)$, the parsimony score for S, T, λ^n , as the value of this sum. A maximum parsimony tree (or trees) for the set of binary strings, S , is a tree (or trees) with leaves $1, 2, \dots, m$ that minimizes $pars(S, T, \lambda^n)$. The value of the sum on this tree is called the parsimony score for the set of strings S .

When the labeling $\hat{\lambda}^n$ is clear, we simply use d_e instead of $d_e(\hat{\lambda}^n)$. In the likelihood setting, we endow edges with “mutation probabilities.” For a rooted tree T , let $\mathbf{p} = [p_e]_{e \in E(T)}$ be the edge probabilities. We use the Neyman two states model [Neyman 1971]. Given labels of length n , each position $j \in \{1, \dots, n\}$ is called a *site*. According to this model:

- Leaf labels are strings from $\{0, 1\}^n$.
- There is uniform distribution of states at the root, namely the probability of each string from $\{0, 1\}^n$ at the root is $(1/2)^n$.
- The “edge probability”, p_e , satisfies $0 \leq p_e \leq \frac{1}{2}$.
- The probability of a net change of state (from ‘1’ to ‘0’ or vice versa) occurring across an edge e (a “mutation event”) is given by p_e . This probability is also called the “length”, or “weight”, of edge e .
- Mutation (change) events on different edges are independent.
- Different sites mutate independently.

Following Edwards [1972], $Pr(T, \mathbf{p}|S)$, the likelihood of a tree T with edge length \mathbf{p} and distribution of strings at the root, given the data S , is defined as the conditional probability that the model T, \mathbf{p} produces the observed data S , $L(T, \mathbf{p}|S) \triangleq Pr(S|T, \mathbf{p})$.

Definition 3.4 (Maximum Likelihood Score). Let T be a rooted phylogenetic tree with edge length \mathbf{p} , and with m leaves. Let S be a set of m binary strings, all of length n . The likelihood of a tree T with edge length \mathbf{p} such that each leaf in the tree get a unique labeling from S is $Pr(S|T, \mathbf{p})$. A maximum likelihood tree (or trees) for the set of binary strings, S , is a tree (or trees) with leaves $1, 2, \dots, m$ that, together with the optimal edge length \mathbf{p} maximize $Pr(S|T, \mathbf{p})$.

Let $S \in \{0, 1\}^{m \times n}$ denote a set of m sequences of length n . Let \mathbf{a} ranges over all combinations of assigning labels (length n 0 or 1 strings) to the r internal nodes of T with $r \leq m - 2$ internal nodes, and edge probabilities \mathbf{p} . Let the term $M(p_e, S_i, a_i)$ denote either p_e or $(1 - p_e)$, depending on whether in the i th site of S and \mathbf{a} , the two endpoints of e are assigned different characters states (and then $M(p_e, S_i, a_i) = p_e$) or the same characters states (and then $M(p_e, S_i, a_i) = 1 - p_e$), and assume uniform distribution at the root, the likelihood of T, \mathbf{p} given the data,

S is:

$$L(T, \mathbf{p}|S) \equiv Pr(S|T, \mathbf{p}) = \frac{1}{2^n} \prod_{i=1}^n \sum_{\mathbf{a} \in \{0,1\}^r} \prod_{e \in E(T)} M(p_e, S_i, a_i), \quad (1)$$

The $(1/2)^n$ is a constant factor, common to all trees and sequences. We drop it in the sequel. This notion of ML is termed maximum *average* likelihood in Steel and Penny [2000]. The ML solution (or solutions) for a specific tree T is the point (or points) in the edge space $\mathbf{p} = [p_e]_{e \in E(T)}$ that maximizes the expression $L(T, \mathbf{p}|S)$. The global ML solution is the pair (or pairs) (T, \mathbf{p}) , maximizing the likelihood over all trees T of M leaves, labeled by S , and all edge probabilities \mathbf{p} (for more details, see Felsenstein [1981], Steel [1994], and Tuffley and Steel [1997]). By the independence of sites, an equivalent way to define the likelihood of observing S in the tree T is:

$$L(T, \mathbf{p}|S) \equiv Pr(S|T, \mathbf{p}) = \frac{1}{2^n} \sum_{\lambda^n \in \{0,1\}^{n \times r}} \prod_{e \in E(T)} p_e^{d_e(\lambda^n)} \cdot (1 - p_e(\lambda^n))^{n - d_e(\lambda^n)} \quad (2)$$

In the rest of the article, we use this definition for likelihood. Our model of substitution is reversible. This implies that the likelihood (and our results) is invariant under the choice of the root, and hence is the same for the unrooted tree as well. This is Felsenstein's pulley principle [Felsenstein 1981].

4. Properties of Maximum Likelihood Trees

In this section, we prove some useful properties of ML trees. We start with properties of general trees and continue with canonical ones.

4.1. GENERAL PROPERTIES OF ML TREES. In our NP-hardness proof, we want to show that the ML tree for a set of reduction strings, defined in the next section, have log likelihood arbitrarily close to the log likelihood of some canonical tree. We achieve this by a sequence of pruning subforests that satisfy certain conditions, and rearranging them in a canonical way around the "root". We will bound the decrease in log likelihood resulting from such arrangements. The following lemma is used several times in the rest of this article.

LEMMA 4.1. *Let T be a phylogenetic tree with edge probabilities \mathbf{p} , let S (set of binary string of length n) denote the labelling for the leaves of the tree. Suppose F_1 and F_2 are two disjoint forests that partition T , and have the node x as their common root. Let S_1 and S_2 be the leaf labelings of F_1 and F_2 , respectively, and let $\mathbf{p}_1, \mathbf{p}_2$ be the induced, corresponding edge probabilities. Then, the likelihood of observing S given T and \mathbf{p} equals*

$$Pr(S|T, \mathbf{p}) = \sum_{s \in \{0,1\}^n} Pr(S_1, \lambda^n(x) = s|F_1, \mathbf{p}_1) \cdot Pr(S_2, \lambda^n(x) = s|F_2, \mathbf{p}_2).$$

PROOF. Follows directly from Eq. (2). \square

For "standard" phylogenetic trees, the internal nodes do not have any specified labeling, while leaves are labelled by sequences of length n . In the course of our modifications, we could also have a leaf with no labeling (see Figure 2). The natural

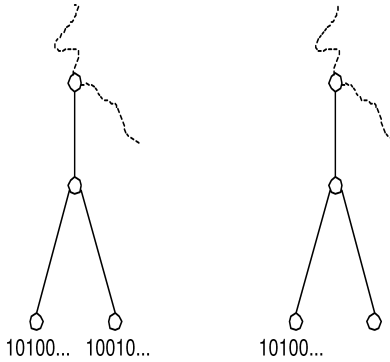


FIG. 2. A standard tree (labeled leaves) and a non-standard tree (some leaves unlabeled).

way to define the likelihood of a tree with such leaves is to treat them as internal nodes, namely summing over all their possible labellings. The next lemma states that such “unlabeled” leaves can be pruned without effecting the likelihood.

LEMMA 4.2. *Let T be a phylogenetic tree with an unlabelled leaf. By pruning this leaf (and the edge connecting to it), we get a tree, T' , with equal likelihood.*

PROOF. Let $S = \{S_i\}$ be the set of leaf labels (binary strings of length n). Let h be the unlabeled leaf, and let h' be its neighbor in T . According to the definition:

$$\begin{aligned} Pr(S | T, \mathbf{p}) &= \sum_{s \in \{0,1\}^n} Pr(S, \lambda^n(h) = s | T, \mathbf{p}) \\ &= \sum_{s \in \{0,1\}^n} \sum_{r \in \{0,1\}^n} Pr(S, \lambda^n(h') = r | T', \mathbf{p}) \cdot Pr(\lambda^n(h) = s, \lambda^n(h') = r | \mathbf{p}) \\ &= \sum_{r \in \{0,1\}^n} Pr(S, \lambda^n(h') = r | T', \mathbf{p}) \\ &= Pr(S | T', \mathbf{p}) \end{aligned}$$

since $\sum_{s \in \{0,1\}^n} Pr(\lambda^n(h) = s, \lambda^n(h') = r | \mathbf{p}) = 1$. \square

LEMMA 4.3. *Let T be a phylogenetic tree with an internal node, h , of degree two, and let g_1, g_2 be its neighbors. Then, h can be eliminated to create an (g_1, g_2) edge without changing the likelihood of T .*

PROOF. Let p_{h,g_1} and p_{h,g_2} be the mutation probabilities of the edges (h, g_1) and (h, g_2) , respectively. Set $p_{g_1,g_2} = p_{h,g_1}(1 - p_{h,g_2}) + p_{h,g_2}(1 - p_{h,g_1})$. It is easy to see that for $0 \leq p_{h,g_1}, p_{h,g_2} \leq 1$, we get $0 \leq p_{g_1,g_2} \leq 1$, and that the mutation probability across the path from g_1 to g_2 does not change. \square

Our NP completeness proof heavily uses a sequence of tree modifications. In each modification, we pruned a subforest, rearrange it, and graft it on the root of the tree. The following theorem establishes a connection between the likelihood of the original and the rearranged subforest, and the change in the total likelihood of the tree.

Definition 4.4 (*Uprooting, Regrafting and Rearranging a Subforest*). Let T be a phylogenetic tree, rooted at the internal node termed *root*. Let h be an internal

node of T . Suppose T_1, \dots, T_j make up a subforest, rooted at h . Denote this subforest by F . Suppose v_1, \dots, v_j are the roots of T_1, \dots, T_j , respectively. Let e_1, \dots, e_j denote the edges connecting h to these roots v_1, \dots, v_j , respectively. The operation of uprooting F and regrafting it on the root consists of deleting the j edges e_1, \dots, e_j , and adding j new edges, $(root, v_1), \dots, (root, v_j)$. Pictorially, this creates a new tree with the j subtrees T_1, \dots, T_j hung off the root. Finally, we take all the leaves' labels of this subforest and rearrange them in a new subforest F_{new} . This rearrangement involve a new edge structure and edge probabilities. This last step is termed rearrangement (see Figure 3). We denote the tree resulting from these operations by $T_{arranged}$.

THEOREM 4.5. *Let T be a phylogenetic tree, with edge probabilities \mathbf{p} , a set of m labels S (each of length n) in its leaves, such that one of its leaves is labeled by the all zero sequence. Let root denote an internal node on T that is at distance 0 from this leaf. Suppose T_1, \dots, T_j is a subforest of T , rooted at h . Denote by F^- the original subforest, by \mathbf{p}^- be its edge probabilities, and by S^- the labels of its leaves. Suppose we uproot, regraft, and rearrange F^- (Definition 4.4). Let \mathbf{p}_{new} denote the edge probabilities of the rearranged forest F_{new} . Let $T_{arranged}$ denote the resulting tree, where the edge probabilities in the ‘‘old’’ part are as in \mathbf{p} . For every labeling s of the node h , the probability $Pr(S^-, \lambda^n(h) = s | F^-, \mathbf{p}^-)$ is well defined. Suppose there is a $W > 0$ such that for every labeling s of h ,*

$$Pr(S^-, \lambda^n(h) = 0 | F_{new}, \mathbf{p}_{new}) \geq W \cdot Pr(S^-, \lambda^n(h) = s | F^-, \mathbf{p}^-).$$

Then

$$Pr(S | T_{arranged}, \mathbf{p}_{arranged}) \geq W \cdot Pr(S | T, \mathbf{p}).$$

PROOF. The probability of S , given the initial tree, T , and \mathbf{p} , is $Pr(S | T, \mathbf{p})$. By Lemma 4.1,

$$\begin{aligned} Pr(S | T, \mathbf{p}) &= \sum_{s \in \{0,1\}^n} Pr(S^-, \lambda^n(h)) \\ &= s | F^-, \mathbf{p}^- \cdot Pr(S \setminus S^-, \lambda^n(h) = s | T \setminus F^-, \mathbf{p} \setminus \mathbf{p}^-). \end{aligned}$$

The probability of S , given $T_{arranged}$ and $\mathbf{p}_{arranged}$ equals (By Lemma 4.1 again),

$$\begin{aligned} &Pr(S | T_{arranged}, \mathbf{p}_{arranged}) \\ &= \sum_{s \in \{0,1\}^n} Pr(S^-, \lambda^n(h) = 0^n | F_{new}, \mathbf{p}_{new}) \\ &\quad \cdot Pr(S \setminus S^-, \lambda^n(h) = s | T_{arranged} \setminus F^-, \mathbf{p}_{arranged} \setminus \mathbf{p}_{new}) \\ &= \sum_{s \in \{0,1\}^n} Pr(S^-, \lambda^n(h) = 0^n | F_{new}, \mathbf{p}_{new}) \\ &\quad \cdot Pr(S \setminus S^-, \lambda^n(h) = s | T \setminus F^-, \mathbf{p} \setminus \mathbf{p}^-). \end{aligned}$$

By our assumption, for each $s \in \{0, 1\}^n$,

$$Pr(S^-, \lambda^n(h) = 0^n | F_{new}, \mathbf{p}_{new}) \geq W \cdot Pr(S^-, \lambda^n(h) = s | F^-, \mathbf{p}^-),$$

and thus $Pr(S | T_{arranged}, \mathbf{p}_{arranged}) \geq W \cdot Pr(S | T, \mathbf{p})$, as desired. \square

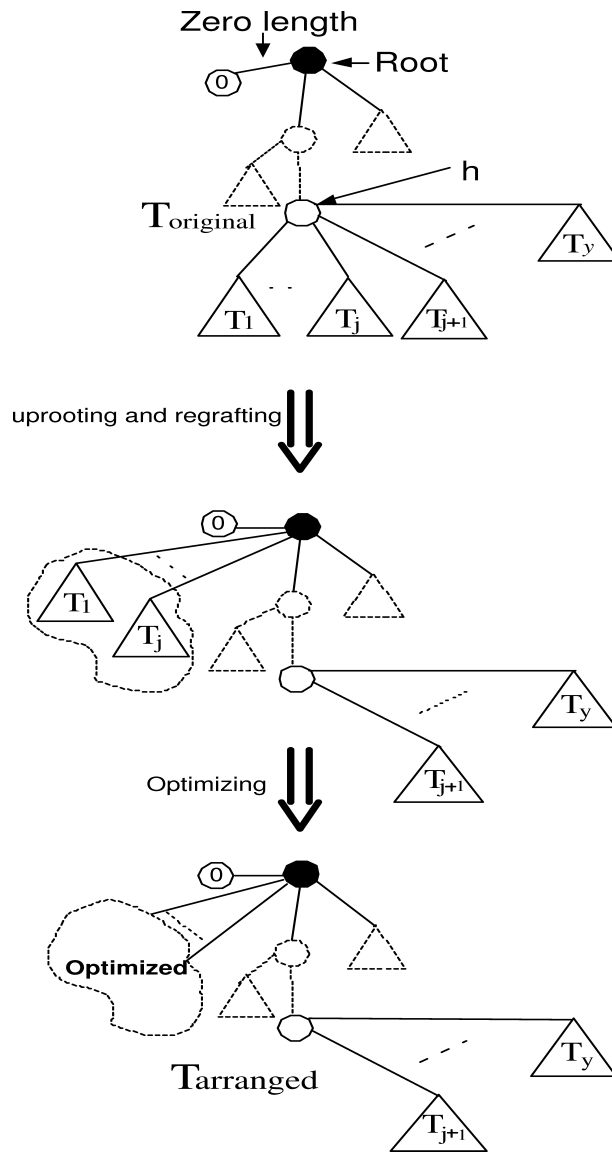


FIG. 3. Uprooting, regrafting and optimizing the subforest T_1, \dots, T_j .

The following corollary is a direct result of Theorem 4.5, and uses the same notation.

COROLLARY 4.6. *Let S^- denote the strings at the leaves of F^- , and let s denote a labelling of h . Let $T^*(s), \mathbf{p}^*(s)$, respectively, denote the structure and edge lengths of a tree that maximize the likelihood of the strings $S^- \cup \{s\}$. Let s^* be a string that maximizes this likelihood, namely for all $s \in \{0, 1\}^n$, $Pr(S^- \cup \{s\} | T^*(s), \mathbf{p}^*(s)) \leq Pr(S^- \cup \{s^*\} | T^*(s^*), \mathbf{p}^*(s^*))$. If $Pr(S^-, \lambda^n(h) = 0^n | F_{\text{new}}, \mathbf{p}_{\text{new}}) \geq W \cdot Pr(S^- \cup \{s^*\} | T^*(s^*), \mathbf{p}^*(s^*))$, then $Pr(S | T_{\text{arranged}}, \mathbf{p}_{\text{arranged}}) \geq W \cdot Pr(S | T, \mathbf{p})$.*

The proof of the following lemma is immediate.

LEMMA 4.7. Let T_{ML}^S, \mathbf{p} denote the structure and edges lengths of an ML tree for the set of strings S . For any $s \in \{0, 1\}^n$, $Pr(S \cup \{s\} | T_{ML}^{S \cup \{s\}}, \mathbf{p}') \leq Pr(S | T_{ML}^S, \mathbf{p})$, where $T_{ML}^{S \cup \{s\}}$ and \mathbf{p}' are optimal structure and edge lengths for $S \cup \{s\}$, respectively.

Definition 4.8 (Rearrangement and Root Labelling, RRL). Let s^* denote the labeling at the root, h , of a subforest F^- , with a set of leaves S^- , that maximized the probability of $S^- \cup \{s^*\}$ given a structure for F^- and a labeling for h , when the maximization is over all the labeling of h , and over all structures for the subforests with the set of leaves S^- . Namely, we are allowed to optimize the root labeling and the structure of the subforest, but are interested only in the resulting root labeling.

In the following lemma, we focus on the subforest F^- , and ignore the rest of the tree. We are allowed to change the structure of F^- , and want to find an optimal labeling s^* of its root. We denote such a labeling s^* of the root, h , an optimal root labeling. By adding one edges of length zero, and one new node, a subforest with a label in its root can easily be translated to a subtree without effecting its same conditional probability. Thus, we use here subtrees and not subforests.

LEMMA 4.9. Suppose there are sites where the value of all the strings in S^- is 0. Then, there is an optimal labeling of a root of F^- , s^* , whose value in each of these sites is also 0.

PROOF. Let us disconnect F^- from the rest of the tree, this makes the subforest a tree, T^{S^-} , with h being one of its leaves. Let s_1 , be one of the sequences in S^- . By the pulley principle, we can root the tree at new node with zero distance to the leaf of s_1 , without changing the conditional probability $Pr(S^- | T^{S^-})$. Removing the leaf h will not lower this probability, by Lemma 4.7. The result is a tree with root label s_1 , which clearly satisfies the required condition. \square

For any tree, T , on m leaves and any observed sequences S , we denote by \mathbf{p}^* the edge probabilities that maximize $Pr(S | T, \mathbf{p})$. The following theorem is a restatement of Theorem 7 from Tuffley and Steel [1997].

THEOREM 4.10. Let S be a set of m binary strings of length $n = n_c + n_{nc}$, where n_c is the number of constant characters in S (i.e. positions that have the same value for all the strings). Let T be a tree on m leaves. Let $pars(S, T)$ denote the parsimony score of S on the tree T . Then

$$2^{-\log(n_c) \cdot pars(S, T) - C_{T, pars(S, T)}^d} \leq Pr(S | T) \stackrel{\Delta}{=} Pr(S | T, \mathbf{p}^*) \leq 2^{-\log(n_c) \cdot pars(S, T) - C_{T, pars(S, T)}^u}$$

and

$$\lim_{n_c \rightarrow \infty} \frac{-\log(Pr(S | T, \mathbf{p}^*))}{\log(n_c)} = pars(S, T),$$

where

$$C_{T, pars(S, T)}^u, C_{T, pars(S, T)}^d = O(n_{nc} m + pars(S, T) \log pars(S, T)).$$

If we hold m fixed and pad the strings in S , then n_c increases, but $pars(S, T)$ remains invariant. The first terms in the exponents of both the upper and lower

bounds become dominant. This establishes the limit, and furthermore the fact that the ML tree “converges” to an MP tree.

COROLLARY 4.11. *Let S contain m binary sequences of length n . Let T_a and T_b be two trees with the strings of S in their leaves. Let \mathbf{p}_a^* and \mathbf{p}_b^* denote the optimal edge lengths for S on these two trees, respectively. Suppose that the strings in S have n_c constant sites. Then, there is n_c large enough (it should be at least $2^{C_{T_i, pars(S, T)}^u} = o(n_c)$) such that $pars(S, T_a) < pars(S, T_b)$ implies $Pr(S|\mathbf{p}_a^*, T_a) > Pr(S|\mathbf{p}_b^*, T_b)$.*

We remark that in general, equality in the parsimony score does not imply equality in the likelihood. The next corollary generalizes the previous one to general trees with one internal node that is labeled. (The definition of likelihood for such trees is easily generalized, and is omitted.) Suppose S is a set containing length k strings, which share n_c constant positions. The likelihood, $Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*)$, of a subforest F with r subtrees T_1, \dots, T_r , sets of labelings of the subtrees’ leaves S_1, \dots, S_r , optimal edge lengths \mathbf{p}^* , and with a label $\lambda^n(h) = s$ at the root of the subforest (see Figure 3) is

$$Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*) = \prod_{i=1}^r Pr(S_i, \lambda^n(h) = s|\mathbf{p}^*, T^i).$$

Therefore

$$Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*) \geq \prod_{i=1}^r 2^{-\log(n_c) \cdot pars(S_i \cup \{s\}, T_i) - C_{T_i, pars(S_i \cup \{s\}, T_i)}^d}, \text{ and}$$

$$Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*) \leq \prod_{i=1}^r 2^{-\log(n_c) \cdot pars(S_i \cup \{s\}, T_i) - C_{T_i, pars(S_i \cup \{s\}, T_i)}^u},$$

where $C_{T_i, pars(S_i \cup \{s\}, T_i)}^u$ and $C_{T_i, pars(S_i \cup \{s\}, T_i)}^d$ are the functions defined in Theorem 4.10. Let $pars(S \cup \{s\}, F) = \sum_i pars(S_i \cup \{s\}, T_i)$, and $C_{S \cup \{s\}, F}^u = \sum_i C_{T_i, pars(S_i \cup \{s\}, T_i)}^u$ and let $C_{S \cup \{s\}, F}^d = \sum_i C_{T_i, pars(S_i \cup \{s\}, T_i)}^d$. Summing up the exponents, we get

COROLLARY 4.1.2

$$\begin{aligned} 2^{-\log(n_c) \cdot pars(S \cup \{s\}, F) - C_{S \cup \{s\}, F}^d} &\leq Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*) \\ &\leq 2^{-\log(n_c) \cdot pars(S \cup \{s\}, F) - C_{S \cup \{s\}, F}^u}. \end{aligned}$$

PROOF. The proof follows directly from Theorem 4.10 and the properties of our model. \square

Let S denote the set of labeling of the leaves of a subforest $F = (V, E)$. Let n_{nc} denote the number of nonconstant sites in $S \cup \{s\}$, where $\lambda^n(h) = s$ is the labeling at the F ’s root, h . The “reduction strings” we will deal with have the property that the number of nonconstant sites of the labelings of the subforest F is small, namely $n_{nc} \leq 2|F|$.

In this case, by Theorem 4.10 and Corollary 4.1.2, we get the following relationship between the log likelihood and the parsimony:

$$\begin{aligned} \log Pr(S, \lambda^n(h) = s|F, \mathbf{p}^*) &= O(pars(S, F) \cdot \log(n_c)) \\ &\quad + O(pars(S, F) \cdot \log(pars(S, F))) + O(|F|^2). \end{aligned}$$

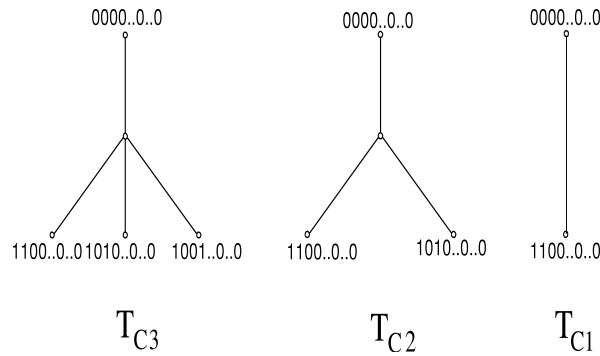


FIG. 4. Building blocks of the maximum likelihood tree for our reduction.

4.1.1. *Properties of Canonical ML Trees.* In this section, we study properties related to canonical ML trees (Definition 2.1), properties that play an important role in our reduction. Throughout this section, the strings we deal with are binary “reduction strings” of length n , originating from a graph of n nodes and m edges.

Definition 4.12. Let T_{C_i} ($i = 1, 2$, or 3) be a tree with $i + 1$ leaves and one internal node (i.e. T_{C_i} has the star topology), such that one of the strings in the leaves is the all zero string (of length n). The other i strings are all of weight 2 (two 1s), and for $i > 1$ they all share one “1” position (see Figure 4). Let $ML_i(n)$ be the log ML score of T_{C_i} for the optimal edge lengths of the tree. Let S_{C_i} denote the strings in the leaves of tree T_{C_i} .

It is easy to see that $ML_i(n)$ does not depend on the specific choice of strings in T_{C_i} .

LEMMA 4.13. *There are constants $C^d \leq C^u$ such that for large enough n ,*

$$-(i + 1) \cdot \log(n) + C^d \leq ML_i(n) \leq -(i + 1) \cdot \log(n) + C^u .$$

PROOF. The proof follows from Theorem 4.10 and direct calculations. \square

THEOREM 4.14. *Let S be a set of “reduction strings”, corresponding to a graph with n nodes and m edges. Let T be a canonical trees with $m + 1$ leaves, labelled by S . Let d denote the degree of its root. Let \mathbf{p}^* be optimal edges weights for S with respect to this tree. Then, for the constants $C^d \leq C^u$ of Lemma 4.13, and for n large enough,*

$$-(d + m) \cdot \log n + dC^d \leq \log Pr(S|T, \mathbf{p}^*) \leq -(d + m) \cdot \log n + dC^u .$$

PROOF. In accordance with Lemma 4.13, $ML_i(n)$, the log likelihood of a subtree T_{C_i} with i nonzero leaves and n long labelings, that is hung off the root of the canonical tree, satisfies

$$-(i + 1) \log(n) + C^d \leq ML_i(n) \leq -(i + 1) \log(n) + C^u .$$

Since T ’s root is effectively labeled by the all zero vector, the log likelihood of S given T is obtained by summing all log likelihoods of its d subtrees. All subtrees together have m leaves, other than the all zero leaf. So the log likelihood of S given T satisfies the inequalities given in the statement of the theorem. \square

Let S be a set of “reduction strings”, corresponding to a graph with n nodes and m edges. Let T_a and T_b be two canonical trees with $m + 1$ leaves, labeled by S . Let d_a and d_b denote the degrees of the roots of these trees, correspondingly. Let p_a^* and p_b^* be optimal edges weights for S with respect to these trees. Then, for the constants $C^d \leq C^u$ of Lemma 4.13, the log likelihood ratio of these trees satisfies

$$\frac{-(d_a + m) \cdot \log n + d_a \cdot C^d}{-(d_b + m) \cdot \log n + d_b \cdot C^u} \leq \frac{\log \Pr(S|T_a, \mathbf{p}_a^*)}{\log \Pr(S|T_b, \mathbf{p}_b^*)} \leq \frac{-(d_a + m) \cdot \log n + d_a \cdot C^u}{-(d_b + m) \cdot \log n + d_b \cdot C^d}.$$

Since $d_a, d_b \leq n$, both the left-hand side and the right-hand side converge to $(d_a + m)/(d_b + m)$ as n grows. This implies that for large enough n ,

COROLLARY 4.15. *For any arbitrarily small ε , there is n_0 large enough such that for all $n \geq n_0$*

$$\frac{d_a + m}{d_b + m} \cdot (1 - \varepsilon) \leq \frac{\log \Pr(S|p_a^*, T_a)}{\log \Pr(S|p_b^*, T_b)} \leq \frac{d_a + m}{d_b + m} \cdot (1 + \varepsilon).$$

And in particular if $d_a = d_b$, then

$$\lim_{n \rightarrow \infty} \frac{\log \Pr(S|T_a, p_a^*)}{\log \Pr(S|T_b, p_b^*)} = 1.$$

5. NP-Hardness of Maximum Likelihood

Building upon the ML machinery developed so far, we now turn to the proof that ML reconstruction on trees is NP hard. We start by formally defining the decision version of maximum likelihood, and then of the gap version of vertex cover we use.

PROBLEM 5.1 (MAXIMUM LIKELIHOOD (ML))

Input. S , A set of binary strings, all of the same length, and a negative number L .

Question. Is there a tree, T , such that $\log \Pr(S|T, \mathbf{p}^*(S, T)) > L$?

A gap vertex cover problem is the following:

PROBLEM 5.2 (GAP PROBLEM FOR VERTEX COVER, $gap - VC[c_1, c_2]$)

Input. A graph, $G = (V, E)$, two positive numbers, c_1 and c_2 .

Task. Does G have a vertex cover smaller than c_1 , or is the size of each vertex cover larger than c_2 ? (If the minimum vertex cover is in the intermediate range, there is no requirement.)

Our proof implies a reduction from the gap version of vertex cover, restricted to degree 3 graphs (undirected graph of degree at most 3 in each node). We rely on the following hardness result of Karpinski and Berman [1999].

THEOREM 5.3 ([BERMAN AND KARPINSKI 1999]). *The following problem,¹ $gap - VC_3[\frac{144}{284} \cdot n, \frac{145}{284} \cdot n]$, is NP-hard: Given a degree 3 graph, G on n nodes, is the minimum VC of G smaller than $\frac{144}{284} \cdot n$, or is it larger than $\frac{145}{284} \cdot n$?*

We reduce that specific version of $gap - VC_3$ above to ML.

¹ We could also use the deep gap VC results of Hästad [2001] and Dinur and Safra [2005]. However, their graphs are of bounded degree greater than 3 and it seems that the modification to bounded degree 3 graphs would yield smaller gaps (not effecting the hardness of ML, though).

5.1. REDUCTION AND PROOF OUTLINE. Given an instance $G = (V, E)$ of $gap-VC_3$, denote $|V| = n$, $|E| = m$, $c_1 = \frac{144}{284} \cdot n$ and $c_2 = \frac{145}{284} \cdot n$. We construct an instance $\langle S, L \rangle$ of ML such that S is a set of $m + 1$ strings, each string of length n , and $L = -(m + \frac{c_1+c_2}{2}) \cdot \log n$.

The first string in S consists of all zeros (the all zeros string), $\underbrace{00\dots 0}_{n}$, and for every edge $e = (i, j) \in E$ there is a string, $S(e) = \underbrace{00\dots 0}_{i-1} \underbrace{100\dots 0}_{j-i-1} \underbrace{100\dots 0}_{n-j} \underbrace{00\dots 0}_n$ where the i th and the j th positions are set to 1, and all the rest are set to 0. These m strings are called “edge strings”. From now on, the trees we refer to have leaves whose labels are generated by this construction.

We use asymptotic properties of likelihood of trees, so most claims will hold when the input graph is large enough (i.e., $n = |V|$ is large enough). In our proof, we deal with small size subtrees or subforests, containing at most $2 \cdot \log \log n$ leaves.

We will need the following relation for the expressions in the likelihood of the subforests to hold (see Corollary 4.1.2):

$$\frac{C_{S \cup \{s\}, F}^d}{\log(n_c) \cdot \text{pars}(S \cup \{s\}, F)}, \frac{C_{S \cup \{s\}, F}^u}{\log(n_c) \cdot \text{pars}(S \cup \{s\}, F)} \xrightarrow{n \rightarrow \infty} 0.$$

By Lemma 4.9, we can assume s have 0 in positions were all the subforest’s strings have 0. The parsimony score (and n_{nc} , the number of nonconstant sites) of such a subforest is no more than $4 \cdot \log \log n$, thus $C_{S \cup \{s\}, F}^d, C_{S \cup \{s\}, F}^u = O((\log \log n)^2)$. Since $n_c = O(n)$ we can use the quantitative relations between parsimony and likelihood as proved in Corollary 4.1.2, our proof is strongly relies on these relations.

5.2. FROM ML TO CANONICAL TREES. In this section, we show that for every $\varepsilon > 0$, there is an $n_0 > 0$ such that for $n > n_0$, the ratio between the log likelihood and the maximum log likelihood of some canonical tree is upper bounded by $(1 + \varepsilon)$.

Given an ML tree, T , if it is in canonical form, we are done. Otherwise, we locate subtrees of T , T_1, T_2, \dots, T_ℓ with a common root, such that the number of leaves in $\bigcup_{i=1}^\ell T_i$ is in the interval $[\log \log n, 2 \cdot \log \log n]$. Notice that this is a subforest as there may be other subtrees rooted at the same node. It is easy to show that such a subforest always exists (Lemma 5.4).

LEMMA 5.4. *Suppose T is a rooted tree and v is an internal node such that the number of leaves below v is at least q . Then, v has a descendent, u , such that u is the root of a subforest consisting of ℓ subtrees T_1, T_2, \dots, T_ℓ ($\ell \geq 1$), and the number of leaves in the subforest $\bigcup_{i=1}^\ell T_i$ is in the range $[q, 2 \cdot q]$.*

The next lemmata, we show that the ratio of the log-likelihood of such a subforest when the all zero labelling is placed in its root, and the log-likelihood of the same subforest with the best labelling in its root, is close to 1.

LEMMA 5.5. *Let u be an internal node or the root h in F , whose degree is $r \geq 9$, and let $s \in \{0, 1\}^n$. Consider an assignment of labels to internal nodes of F , where h is assigned s . Among such assignments, those that optimize the parsimony score label u with 0^n .*

PROOF. We assume here $u \neq h$, the case were $u = h$ can be proved in a similar way. It suffices to prove the claim for every position separately. The internal node u have $r - 1$ subtrees below it, and one edge “above” it, leading to h . Out of these

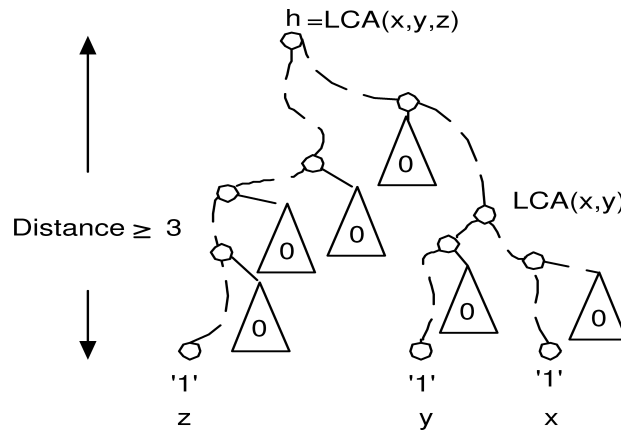


FIG. 5. Lemma 5.6, case (2).

subtrees, at most three have “1” in the position of interest (since our graphs are of degree 3). For the other $r - 4 > 4$ subtrees, since their leaves have 0 in the position, the most parsimonious assignment will label all their nodes with 0, as can be seen by running Fitch algorithm [Fitch 1971]. Therefore, u has at least five neighbor nodes with 0 in this position, and at most four with 1. Any parsimonious assignment will thus label u with 0. \square

LEMMA 5.6. *Let h be the root of a subforest F (h has at least two children in F) whose leaves are labelled by a subset of the reduction strings, S' . Suppose that in each position, the leaves labeled with “1” are at distance at least 3 from h . Then, the maximum parsimony score for S' on F is achievable with the all zero labeling in h .*

PROOF. Consider an arbitrary position. Since the reduction strings emanate from a degree 3 graph, there are at most three leaves x, y, z with “1” in this position. Let $LCA(x, y), LCA(x, y, z)$ denote the least common ancestors of x, y and x, y, z , respectively (see Figure 6). Then either $LCA(x, y)$ is equal to $LCA(x, y, z)$ or is below it in F . For any node j in F , we denote by $pa(j)$ the parent of j . There are three cases:

- (1) $h = LCA(x, y, z)$ and $LCA(x, y) = LCA(x, y, z)$: There are at least two intermediate nodes on each of the paths from h to x, y , and z , respectively. The leaves below each of these intermediate nodes, other than x, y, z , all contains “0” in this position. Thus, by Fitch algorithm [Fitch 1971], the best assignment to the nodes in the paths from z, y , and x to h is “0”. Therefore, if we assign “1” to h we lose 1 on each edge leading to h , a total loss of 3. If we assign “0” to h , we lose nothing on the node immediately below to h .
- (2) $h = LCA(x, y, z)$ and $LCA(x, y) \neq LCA(x, y, z)$ (see Figure 5): The path from $LCA(x, y)$ to h joins the path from the path from z to h at h and not below h , for otherwise we’d have $h \neq LCA(x, y, z)$. Now there are at least 2 internal nodes on the path from z to h , none of which is a root to a subtree with a “1” in this position. Therefore, by Fitch algorithm [Fitch 1971], the best assignment to the internal nodes on the path from z to h is “0”. Thus, if we assign “0” to h we may cost 1 in the score due to the node just below h in the path between

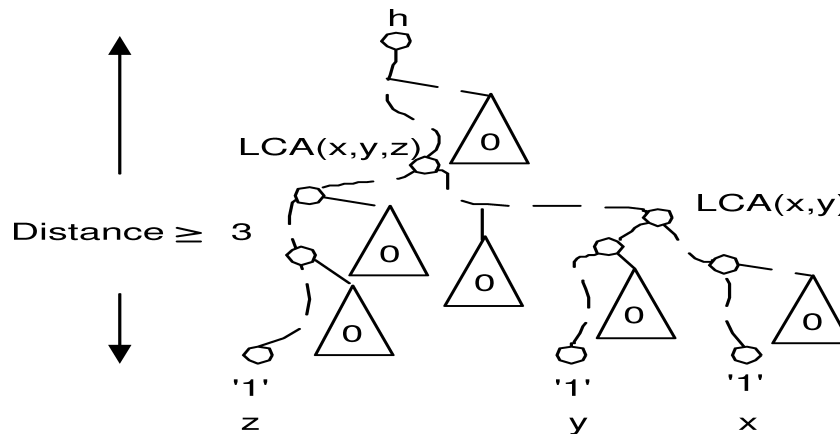


FIG. 6. Lemma 5.6, case (3).

$LCA(x, y)$ and h , but cost nothing in the score due to the edges to the other children of h . On the other hand, if we assign “1” to h we certainly pay 1 in the score, due to the node before last in the path from z to h , and may pay 1 in the score due to the node just below h in the path between $LCA(x, y)$ and h . Thus, the “0” labeling to h is not worse than the “1” labeling.

- (3) $h \neq LCA(x, y, z)$ (see Figure 6): Since h has at least two children, all the leaves under one of these children are “0”, so the algorithm of Fitch assigns “0” to this node. Since $LCA(x, y, z)$ is below h , Fitch’s algorithm will assign “1” to at most one of h children. Thus, the “0” labelling to h is not worse than the “1” labeling. \square

COROLLARY 5.7. *Let h be a root of a subforest F whose leaves are labeled by a set of reduction strings, S . Consider a specific position, and suppose all leaves having “1” in the position are either at distance ≥ 3 from the root, or have an internal node of degree ≥ 9 on the path to the root. Then, the parsimony score of F when labeling the root h with 0 at this position is at least as good as when labelling the root with 1.*

PROOF. In accordance with Lemma 5.5, if there is an internal node of degree ≥ 9 , it is better to assign 0 to this node. This enables us to disregard the “1” leaves of distance smaller than 3 from h with a high degree node on their path to h . The other leaves with “1” are at distance at least 3 from h , and we can now apply Lemma 5.6. \square

THEOREM 5.8. *Let T be a tree whose leaves are labeled by a subset of the reduction strings that are nonzero (i.e., the labeling of the leaves do not include the all zero string). Let F be a subforest of T , rooted at h , and let $s \in \{0, 1\}^n$ be a label of h . The parsimony score of F with the 0^n labeling at the root can be worse by at most $8^3 - 1$ than the parsimony score of F with label s at its root.*

PROOF. If the degree of h is larger than 8, then by Lemma 5.5 the best assignment to h is 0^n and we are done. Otherwise, we say that a leaf in F is *dangerous* if its distance from h less than 3, and all its ancestors on the path to h have degree ≤ 8 . Simple counting shows that the number of dangerous leaves in F is smaller than

$8 + 8^2$. Every dangerous leaf has 2 positions where it is labeled “1”. Each such position can be “1” in at most 3 leaves because the graph G is of degree 3. Therefore for any of these positions, changing the label at h from 1 to 0 will worsen the parsimony score by at most 3. There are at most $2 \cdot (8 + 8^2)$ such positions. So changing to 0 at h will cause at most $2 \cdot 3 \cdot (8 + 8^2) < 8^3$ parsimony degradations. In accordance of Lemma 5.7, in all other positions, the “0” label at h is optimal. \square

The following lemma was proved by Day et al. [1986] and Addario-Barry et al. [2004].

LEMMA 5.9. *Let $S' \subseteq S$ be a subset of the reduction strings, which contains the all zero string. The structure of the best parsimony tree for S' is canonical.*

THEOREM 5.10. *For every $\varepsilon > 0$, there is an n' such that for all $n \geq n'$, if $S \subseteq \{0, 1\}^n$ is a set of $m + 1$ reduction strings corresponding to a degree 3 graph with n nodes and m edges, then the following holds: Let T_{ML} denote an ML tree for S , and let \mathbf{p}_{ML}^* and be an optimal edges length for this ML tree. Then, there is a canonical tree for S , T_{CA} , with optimal edges length \mathbf{p}_{CA}^* , such that*

$$\frac{\log Pr(S|T_{ML}, \mathbf{p}_{ML}^*)}{\log Pr(S|T_{CA}, \mathbf{p}_{CA}^*)} < 1 + \varepsilon .$$

PROOF. We start from any ML tree, T_{ML} , and show how to transform it to a canonical tree, T_{CA} , with “close enough” log likelihood, in a sequence of up to $n / \log \log(n)$ steps. Each step involves a small, local change to the current tree. We identify a subforest (disjoint subtrees with a common root), whose number of leaves is in the interval $[\log \log(n), 2 \log \log(n)]$. By Lemma 5.4, if the root of the whole tree has a subtree with more than $\log \log(n)$ leaves, we can find such a subforest. In such a case, we first prune this subforest, and then regraft it under the tree’s root. Let S_F be the restriction of the set S to the labeling of F ’s leaves. By Theorem 5.8, the parsimony score of S_F on F , with the 0^n label at the root, is larger by at most $B \triangleq 8^3$ than the parsimony score of S_F on F with any $s \in \{0, 1\}^n$ label at its root. We will show how this implies that for every s , the likelihood $Pr(S_F, \lambda^n(h) = s|F, \mathbf{p}_s^*)$ is not much larger than $Pr(S_F, \lambda^n(h) = 0^n|F, \mathbf{p}_s^*)$. To prove this, we use the results of Tuffley and Steel [1997] (Corollary 4.1.2), the properties of our strings, and the small size of the subforest.

Let s^*, F^* be a string and a subforest structure, respectively, for which $Pr(S_F, \lambda^n(h) = s|F, \mathbf{p}_s^*)$ is maximized. By Lemma 4.9, such an s^* that has 1s only in positions where some $s \in S_F$ has a 1 exists.

Clearly, instead of bounding the likelihood difference for every $s \in \{0, 1\}^n$, it suffices to bound $Pr(S_F, \lambda^n(h) = s^*|F^*, \mathbf{p}_s^*) - Pr(S_F, \lambda^n(h) = 0^n|F, \mathbf{p}_s^*)$.

The parsimony score of $S_F \cup \{0^n\}$ on F , $pars(S_F \cup \{0^n\}, F)$, is no larger than the number of “1” entries in $S_F \cup \{0^n\}$. There are at most two “1”s per string, and the number of nonzero strings is at most $2 \log \log n$. Therefore, $pars(S_F \cup \{0^n\}, F) \leq 4 \log \log n$. Since s^* can have a “1” only in positions where some $s \in S_F$ has a “1”, the total number of “1”s in $S_F \cup \{s^*\}$ is at most twice the number in $S_F \cup \{0^n\}$. Therefore, $pars(S_F \cup \{s^*\}, F^*) \leq 8 \log \log n$.

Let us denote by n_{s^*} the number of constant sites with s^* at the root of F^* , and by k_0 the number of constant sites with 0^n at the root of F , by par_{s^*} the parsimony score of F^* with s^* at the root, and by par_0 the parsimony score of F with 0^n at

the root. Then, we show that $par_0 - par_{s^*} \leq B$, and know that n_{s^*}, n_0 are both $n - O(\log \log n)$, implying $\log(n_0) - \log(n_{s^*}) = \theta(1)$. Finally, by Theorem 4.10, the order of magnitude of both $C_{S_F \cup \{0^n\}, F}^u$ and $C_{S_F \cup \{s^*\}, F^*}^d$ is $O(n_{nc}m + \text{pars} \cdot \log \text{pars})$, where $n_{nc} = O(\log \log n)$ is the number of nonconstant sites in the set of strings, $m \leq 2 \log \log n$ is the number of strings, and $\text{pars} = O(\log \log n)$ is the parsimony value. All by all, in our case $C_{S_F \cup \{0^n\}, F}^u - C_{S_F \cup \{s^*\}, F^*}^d$ are both $O((\log \log n)^2)$, and so is their difference. These inequalities imply

$$\begin{aligned}
& \log(\Pr(S_F, \lambda^n(h) = s^* | F^*, \mathbf{p}_{s^*}^*)) - \log(\Pr(S_F, \lambda^n(h) = 0^n | F, \mathbf{p}_0^*)) \\
& \leq -\log(n_{s^*}) \cdot par_{s^*} - C_{S_F \cup \{s^*\}, F^*}^u - (-\log(n_0) \cdot par_0 - C_{S_F \cup \{0^n\}, F}^d) \\
& = \log(n_0) \cdot (par_0 - par_{s^*}) + (\log(n_0) - \log(n_{s^*})) \cdot par_{s^*} \\
& \quad + (C_{S_F \cup \{0^n\}, F}^u - C_{S_F \cup \{s^*\}, F^*}^d) \\
& \leq B \log(n_0) + \theta(par_{s^*}) + (C_{S_F \cup \{0^n\}, F}^u - C_{S_F \cup \{s^*\}, F^*}^d) \\
& \leq B \log(n_0) + O(\log \log(n)) + O((\log \log n)^2) \\
& \leq B \log(n) + o(\log n)
\end{aligned}$$

To summarize, for each $s \in \{0, 1\}^n$, we get

$$\begin{aligned}
& \log(\Pr(S_F, \lambda^n(h) = s | F, \mathbf{p}_s^*)) - \log(\Pr(S_F, \lambda^n(h) = 0^n | F, \mathbf{p}_0^*)) \\
& \leq B \log(n) + o(\log n).
\end{aligned}$$

Let $T_{\text{arranged}}, \mathbf{p}_{\text{arranged}}^*$ denote the tree resulting from uprooting and regrafting the subforest F , and its optimal edge lengths. The conditions of Theorem 4.5 apply, so we conclude

$$\log \Pr(S | T, \mathbf{p}^*) - \log \Pr(S | T_{\text{arranged}}, \mathbf{p}_{\text{arranged}}^*) \leq B \log(n) + o(\log n),$$

Namely each single uprooting decreases the overall log likelihood of S by no more than $B \log(n) + o(\log n)$. All uprootings therefore decrease the log likelihood by at most $Bn \log(n) / \log n + o(n \log n / \log \log n)$.

After a sequence of up to $n / \log \log n$ such uprootings, we get a tree having no subtrees with $\log \log n$ or more leaves. To get the desired canonical tree, we *separately* “canonize” each small subtree, namely rearrange it in an optimal canonical form. In accordance with Lemma 5.9, we can rearrange such a subforest in a canonical form, with the all zero root, such that its parsimony score does not deteriorate. Let F_c denote such canonical rearrangement, and let \mathbf{p}_c^* denote the optimal edge lengths for the rearrangement, and n_c the number of constant sites in the strings set. By Corollary 4.1.2,

$$\begin{aligned}
& \log \Pr(S, \lambda^n(h) = 0^n | F, \mathbf{p}_0^*) - \log \Pr(S, \lambda^n(h) = 0^n | F_c, \mathbf{p}_c^*) \\
& \leq -\log(n_c) \cdot \text{pars}(S \cup 0^n, F) - C_{S \cup \{0^n\}, F}^u \\
& \quad - (-\log(n_c) \cdot \text{pars}(S \cup 0^n, F_c) - C_{S \cup \{0^n\}, F_c}^d) \\
& = \log(n_c)(\text{pars}(S \cup 0^n, F_c) - \text{pars}(S \cup 0^n, F)) + (C_{S \cup \{0^n\}, F_c}^d - C_{S \cup \{0^n\}, F}^u) \\
& \leq C_{S \cup \{0^n\}, F_c}^d - C_{S \cup \{0^n\}, F}^u \\
& = O((\log \log n)^2)
\end{aligned}$$

Therefore, each such rearrangement can decrease the log likelihood of S given T by at most $O((\log \log n)^2)$. The minimal size of a subforest that needs rearrangement is 2, so here are no more than $n/2$ subforests to be rearranged. Overall, the decrease in log likelihood due to the rearrangements is $O(n(\log \log n)^2)$. Taking into both uprootings and rearrangements, the total log likelihood loss of the process is $Bn \log(n)/\log \log n + o(n \log n/\log \log n) + O(n(\log \log n)^2) = o(n \log n)$.

In accordance with Theorem 4.14, the log-likelihood of all canonical trees is larger than $-n \log(n)$. We just showed the existence of a canonical tree whose log likelihood differs from the log likelihood of any ML tree by less than $Bn \log(n)/\log \log(n) + o(n \log(n)/\log \log(n))$. Thus, there must be a constant $K > 0$ such that the log-likelihood of any ML tree is at most $-K \cdot n \log(n)$, and consequently there is a canonical tree such that the ratio between the log likelihood of the ML tree and this tree is

$$\frac{-K \cdot n \log(n)}{-K \cdot n \log(n) - O(n \cdot \log n/\log \log(n))} = 1 + O\left(\frac{1}{\log \log n}\right)$$

implying that for every $\varepsilon > 0$ there is an n' such that for all $n > n'$

$$Pr(S|T_{ML})/Pr(S|T_{CA}) < 1 + \varepsilon. \quad \square$$

We remark that the size of the subforests could be chosen to be different than $\theta(\log \log n)$ and still get similar result, provided they are neither too small nor too large.

5.3. VALIDITY OF THE REDUCTION. In this section, we complete the proof, by showing that indeed we have a reduction from $GAP - VC_3$ to ML. We show that if G has a small enough cover, then the likelihood of the corresponding canonical tree is high (this is the easy direction), and if the likelihood is high, then there is a small cover (this is the harder direction). The translation of sizes, from covers to log likelihood, and vice versa, is not sharp, but introduces some slack. This is why a gap version of vertex cover, instead of exact vertex cover, is required as our starting point.

The next lemma establishes a connection between MP and VC , and was used in the NP-hardness proof of MP .

LEMMA 5.11 ([DAY ET AL. 1986; ADDARIO-BARRY ET AL. 2004]). $G = (V, E)$ has a vertex cover of size c if and only if there is a canonical tree with parsimony score $c + m$, where c is the degree of the root.

THEOREM 5.12. For every $0 < \varepsilon$, there is an n' such that for every $n \geq n'$, if G is a degree 3 graph on n nodes and m edges, with a cover of size at most c , then there is a tree T such that the log-likelihood of S satisfies

$$\log(Pr(S|T, \mathbf{p}^*(S, T))) > -(1 + \varepsilon)(m + c) \log n.$$

On the other hand, if the size of every cover is $\geq c$, then the log likelihood of S given T satisfies

$$\log(Pr(S|T, \mathbf{p}^*(S, T))) < -(1 - \varepsilon)(m + c) \log n.$$

PROOF. Suppose G has a vertex cover of size $\leq c$. Since G 's is of bounded degree 3, c satisfies $m/3 \leq c \leq m$, and $n \leq m \leq 1.5n$. In accordance with Lemma 5.11, there is a canonical tree, T_{CA} , with parsimony score $c + m$, such

that the degree of its root is c . In accordance with Theorem 4.14, the log likelihood of S , given this tree is, $-(c + m) \log(n) + \theta(n)$. Since $m, c = \theta(n)$, $\log(\Pr(S|T_{CA}, \mathbf{p}^*(S, T_{CA}))) = -(c + m) \log(n) + \theta(n)$ implies that for every $\varepsilon > 0$ and large enough n ,

$$\log(\Pr(S|T_{CA}, \mathbf{p}^*(S, T_{CA}))) > -(m + c) \log(n)(1 + \varepsilon).$$

For the other direction, suppose the size of every cover of G is greater or equal to c . In accordance with Lemma 5.11, the parsimony score of each canonical tree is at least $m + c$. Thus, by Theorem 4.14, the likelihood of S , given any tree, is at most $-(m + c) \log(n) + cC^u$ (where C^u is the constant from the theorem). Since $m, c = \theta(n)$ we get that for every $\varepsilon_1 > 0$ and large enough n ,

$$-(m + c) \log(n) + cC^u < -(m + c) \log(n)(1 - \varepsilon_1).$$

In accordance with Theorem 5.10, this implies that the likelihood of S with respect to any ML tree satisfies

$$\log(\Pr(S|T_{ML}, \mathbf{p}^*(S, T_{ML}))) < -(m + c) \log(n)(1 - \varepsilon_1)(1 - \varepsilon_2),$$

where $\varepsilon_1, \varepsilon_2$ are arbitrarily small, and n is large enough. Thus, for every ε there is n' such that for $n > n'$ the likelihood of the best trees satisfies

$$\log(\Pr(S|T_{ML}, \mathbf{p}^*(S, T_{ML}))) < -(m + c) \log(n)(1 - \varepsilon). \quad \square$$

THEOREM 5.13. *ML reconstruction on trees is NP-hard.*

PROOF. Let $G = (V, E)$ be an instance of *gap-VC₃*. Denote $|V| = n, |E| = m$, $c_1 = \frac{144}{284} \cdot n$ and $c_2 = \frac{145}{284} \cdot n$. Recall that in the reduction, we construct an instance $\langle S, L \rangle$ of *ML* such that S is a set of $m + 1$ strings, each string is of length n , and the threshold is $L = -(m + \frac{c_1 + c_2}{2}) \cdot \log n$.

Suppose $G \in \text{gap-VC}_3$. Then G has a cover of size $\leq c_1$. According to Theorem 5.12, for every $\varepsilon > 0$ and large enough n , $\log(\Pr(S|T_{ML}, \mathbf{p}^*(S, T_{ML}))) > -(m + c_1) \log(n)(1 + \varepsilon)$. Thus, in order to show that $\langle S, L \rangle \in \text{ML}$, it suffices to show the existence of $\varepsilon > 0$ so that $(m + c_1)(1 + \varepsilon) < m + (c_1 + c_2)/2$. Since $c_1 < n$ and $m \leq 1.5n$, and $(c_2 - c_1)/2 = n/568$, simple arithmetic shows that by taking $\varepsilon = 1/1420$, the inequality is satisfied.

Suppose $G \notin \text{gap-VC}_3$. Then every cover of G is of size $\geq c_2$. In accordance with Theorem 5.12, for every $\varepsilon > 0$ and large enough n , $\log(\Pr(S|T_{ML}, \mathbf{p}^*(S, T_{ML}))) < -(m + c_2) \log(n)(1 - \varepsilon)$. In order to show that $\langle S, L \rangle \notin \text{ML}$, it suffices to show the existence of $\varepsilon > 0$ so that $(m + c_2)(1 - \varepsilon) > m + (c_1 + c_2)/2$. Since $c_2 < n$ and $m \leq 1.5n$, and $(c_2 - c_1)/2 = n/568$, simple arithmetic shows that by taking $\varepsilon = 1/1420$ again, the inequality is satisfied. \square

5.4. OTHER SUBSTITUTION MODELS. Our NP hardness result was stated in Neyman's two states model of substitution. What about 4 states DNA, or proteins? It turns out that proving that such an extension of the main theorem also holds for these situations is not difficult. In this section, we prove NP-hardness of maximum likelihood reconstruction under the Jukes-Cantor model [Jukes and Cantor 1969]. This model is a special case of Kimura 2 parameter and 3 parameter models, and of more elaborate models of DNA substitution. The same holds for protein sequences as well.

Suppose we have a c state alphabet, Σ (for DNA sequences, $c = 4$). Let α_e denote a substitution parameter associated with the edge e . In the JC model, there

is a certain probability $1 - p_e$ that a character does not change across the edge e . If it does change, the probabilities of changing to any one of the other $c - 1$ characters are equal, $p_e/(c - 1)$. The likelihood of S given a tree under this model is defined in a way similar to Eq. (2):

$$Pr(S|T, \mathbf{p}) = \sum_{\lambda^n \in \Sigma^{n \times r}} \prod_{e \in E(T)} \frac{p_e}{c - 1}^{d_e(\lambda^n)} (1 - p_e(\lambda^n))^{n - d_e(\lambda^n)}. \quad (3)$$

According to Tuffley and Steel 1997], we get for this model relations that are similar to the theorems, lemmata and corollaries in Section 4 (with $C_{T, pars(S, T)}^u$ and $C_{T, pars(S, T)}^d$ that are different but have the same order of magnitude). Thus, our reduction holds for the JC model, and consequently for all models extending the JC model.

6. Concluding Remarks and Further Research

In this work, we proved that ML reconstruction of phylogenetic trees is computationally intractable. We used the simplest model of substitution—the Neyman two states model [Neyman 1971]. This NP-hardness proof generalizes to the Jukes–Cantor model [Jukes and Cantor 1969], and then to the Kimura and other models of DNA and protein substitution.

Since the extended abstract of this work was submitted and published in RECOMB05, the results in this article have been extended in three directions: We have shown that ML remains hard even under the assumption of molecular clock. We proved an initial $1 + \varepsilon$ hardness result of approximation for log likelihood (for a rather small ε). We developed an approximation algorithm for log likelihood for special, biologically interesting, sets of inputs. These results were presented in ISMB 2005. Starting the reduction from vertex cover, and going through canonical trees, were crucial in our hardness proof of ML under molecular clock [Chor and Tuller 2005]. After the conference version of this paper was submitted (and widely distributed), Roch [2006] published a different proof for hardness of ML. His proof reduce directly from the hardness of MP. Consequently, it is shorter, but seems not to yield the hardness under molecular clock, mentioned above.

Vertex cover, which is the starting point for our reduction, has a simple 2-approximation algorithm. Maximum parsimony has 2-approximation (and better) algorithms. What about *any constant* approximation algorithms for log likelihood? So far, no constant factor approximations are known. It will be interesting to find a b approximation of log likelihood for some constant $b > 1 + \varepsilon$ (for the above ε), or to prove that no such efficient algorithm exists (unless $P = NP$). Finally, it would be nice to identify regions of inputs where ML is *tractable*. In this context, we note that it is not even known what is the complexity of *small ML*, where the sequences and the unweighted tree are given, and the goal is to find optimal edge lengths. In practice, local search techniques such as EM or hill climbing seem to perform well, but no proof of performance is known, and multiple maxima [Steel 1994; Chor et al. 2000] shed doubts even on the (worst-case) correctness of this approach.

ACKNOWLEDGMENTS. We wish to thank Isaac Elias for helpful discussions, Sagi Snir for reading early drafts of the manuscript, and the anonymous referees for their thorough job.

REFERENCES

- ADDARIO-BERRY, L., CHOR, B., HALLETT, M., LAGERGREN, J., PANCONESI, A., AND WAREHAM, T. 2004. Ancestral maximum likelihood of evolutionary trees is hard. *J. Bioinf. Comput. Biol.* 2, 2, 257–271.
- BERMAN P., AND KARPINSKI, M. 1999. On some tighter inapproximability results. In *Proceedings of the 25th International Colloquium on Automatic Languages, and Programing*. Lecture Notes in Computer Science. Springer-Verlag, New York.
- CHOR, B., HENDY, M. D., HOLLAND, B. R., AND PENNY, D. 2000. Multiple maxima of likelihood in phylogenetic trees: An analytic approach. *Mol. Biol. Evol.* 17, 10, 1529–1541.
- CHOR B. AND TULLER, T. 2005. Maximum likelihood of evolutionary trees: Hardness and approximation. *Bioinformatics* 21, 1, i97–i106.
- DAY, W. 1987. The computational complexity of inferring phylogenies from dissimilarity matrix. *Bull. Math. Biol.* 49, 4, 461–467.
- DAY, W., JOHNSON, D., AND SANKOFF, D. 1986. The computational complexity of inferring rooted phylogenies by parsimony. *Mathematical Biosciences* 81, 33–42.
- DAY W., AND SANKOFF, D. 1986. The computational complexity of inferring phylogenies by compatibility. *Syst. Zool.* 35, 2, 224–229.
- DINUR I., AND SAFRA, S. 2005. On the importance of being biased (1.36 hardness of approximating vertex-cover). *Ann. Math.* 162, 439–485.
- EDWARDS, A. W. F. 1972. *Likelihood*. Cambridge University Press, Cambridge, UK.
- FELSENSTEIN, J. 1981. Evolutionary trees from DNA sequences: A maximum likelihood approach. *J. Mol. Evol.* 17, 368–376.
- FELSENSTEIN, J. 1996. Inferring phylogenies from protein sequences by parsimony, distance, and likelihood methods. *Meth. Enzym.* 266, 419–427.
- FELSENSTEIN, J. 2004. *Inferring Phylogenies*. Sinauer Associates, Sunderland, MA.
- FITCH, W. M. 1971. Toward defining the course of evolution: Minimum change for specified tree topology. *Syst. Zool.* 20, 406–416.
- FOULDS, L., AND GRAHAM, R. 1982. The Steiner problem in phylogeny is NP-complete. *Adv. Appl. Math.* 3, 43–49.
- HÅSTAD, J. 2001. Some optimal inapproximability results. *J. ACM* 48, 798–859.
- JUKES, T. H. AND CANTOR, C. R. 1969. Evolution of protein molecules. In *Mammalian Protein Metabolism*. H. N. Munro, Ed. Academic Press, New York, pp. 21–132.
- KARPINSKI, M. 2001. Approximating bounded degree instances of np-hard problems. In *Proceedings of the 42nd IEEE Symposium on Foundations of Computer Science*. IEEE Computer Society Press, Los Alamitos, CA.
- KOSHI, M., AND GOLDSTEIN, R. 1996. Probabilistic reconstruction of ancestral nucleotide and amino acid sequences. *J. Molec. Evol.* 42, 313–320.
- NEYMAN, J. 1971. Molecular studies of evolution: A source of novel statistical problems. In *Statistical Decision Theory and Related Topics*, S. Gupta and Y. Jackel, Eds. Academic Press, New York, pp. 1–27.
- ROCH, S. 2006. A short proof that phylogenetic tree reconstruction by maximum likelihood is hard. *IEEE/ACM Trans. Comput. Biol. Bioinform.* 3, 1, 92–94. accepted.
- STEEL, M. 1992. The complexity of reconstructing trees from qualitative characters and subtrees. *J. Classif.* 9, 71–90.
- STEEL, M. 1994. The maximum likelihood point for a phlogenetic tree is not unique. *Syst. Biol.*, 43, 560–564.
- STEEL, M., AND PENNY, D. 2000. Parsimony, likelihood and the role of models in molecular phylogenetics. *Mol. Biol. Evol.* 17, 839–850.
- TUFFLEY, C., AND STEEL, M. 1997. Link between maximum likelihood and maximum parsimony under a simple model of site substitution. *Bull. Math. Biol.* 59, 3, 581–607.
- WAREHAM, T. 1993. On the computational complexity of inferring evolutionary trees. Tech. Rep. 93–01. Dep. Computer Science, Memorial University of Newfoundland.
- YANG, Z., KUMAR, S., AND NEI, M. 1995. A new method of inferring of ancestral nucleotide and amino acid sequences. *Genetics* 141, 1641–1650.

RECEIVED JULY 2005; REVISED FEBRUARY 2006; ACCEPTED AUGUST 2006