

SAN097-0907C,
CONF-970683--1

Approximation Algorithms for the Fixed-Topology Phylogenetic Number Problem

Mary Cryan* Leslie Ann Goldberg[†] Cynthia A. Phillips[‡]

Abstract. In the ℓ -phylogeny problem, one wishes to construct an evolutionary tree for a set of species represented by characters, in which each state of each character induces no more than ℓ connected components. We consider the fixed-topology version of this problem for fixed-topologies of arbitrary degree. This version of the problem is known to be \mathcal{NP} -complete for $\ell \geq 3$ even for degree-3 trees in which no state labels more than $\ell + 1$ leaves (and therefore there is a trivial $\ell + 1$ phylogeny). We give a 2-approximation algorithm for all $\ell \geq 3$ for arbitrary input topologies and we give an optimal approximation algorithm that constructs a 4-phylogeny when a 3-phylogeny exists. Dynamic programming techniques, which are typically used in fixed-topology problems, cannot be applied to ℓ -phylogeny problems. Our 2-approximation algorithm is the first application of linear programming to approximation algorithms for phylogeny problems. We extend our results to a related problem in which characters are polymorphic.

MASTER

* maryc@dcs.warwick.ac.uk. Department of Computer Science, University of Warwick, Coventry CV4 7AL, United Kingdom. This work was partly supported by ESPRIT LTR Project no. 20244 — ALCOM-IT.

[†] leslie@dcs.warwick.ac.uk. Department of Computer Science, University of Warwick, Coventry CV4 7AL, United Kingdom. Part of this work took place during a visit to Sandia National Laboratories which was supported by University of Warwick Research and Teaching Innovations Grant 0951CSA and by the U.S. Department of Energy under contract DE-AC04-76AL85000. Part of this work was supported by ESPRIT LTR Project no. 20244 — ALCOM-IT.

[‡] caphill@cs.sandia.gov. Sandia National Laboratories, Albuquerque, NM. This work was performed under U.S. Department of Energy contract number DE-AC04-76AL85000.

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED

un

1. Introduction. The evolutionary biologist collects information on extant species (and fossil evidence) and attempts to infer the evolutionary history of a set of species. Most mathematical models of this process assume divergent evolution, meaning that once two species diverge, they never share genetic material again. Therefore, evolution is modelled as a tree (*phylogeny*), with extant species as leaves and (extant, extinct, or hypothesized) ancestors as internal nodes. Species have been modelled in several ways, depending upon the nature of available information and the mechanism for gathering that information. Based upon these representations, differing measures of evolutionary distance and objective function are used to evaluate the goodness of a proposed evolutionary tree.

In this paper we assume that input data is character-based. Let S be an input set of n species. A *character* c is a function from the species set S to a set R_c of *states*. If we are given a set of characters c_1, \dots, c_k for S , each species is a vector from $R_{c_1} \times \dots \times R_{c_k}$, and any such vector can represent a hypothesized ancestor. Characters can be used to model biomolecular data, such as a column in a multiple sequence alignment, but in this paper, we think of characters as morphological properties such as coloration or the ability to fly.

Character-based phylogenies are typically evaluated by some *parsimony-like* measure, meaning that the total evolutionary change is somehow minimized. In this paper, we consider the ℓ -*phylogeny* metric introduced in [8]. Given a phylogenetic tree, a character c_i and a state $j \in R_{c_i}$, let l_{ij} be the number of connected components in $c_i^{-1}(j)$ (the subtree induced by the species with state j in character i). A phylogeny is an ℓ -phylogeny if $\max_{c_i, j \in R_{c_i}} l_{ij} \leq \ell$. The ℓ -*phylogeny problem* is to determine if an input consisting of a species set S and a set of characters c_1, \dots, c_k has an ℓ -phylogeny. The *phylogenetic number problem* is to determine the minimum ℓ such that it has an ℓ -phylogeny. The classic parsimony problem is to find a tree that minimizes $\sum_{c_i, j \in R_{c_i}} l_{ij}$. The compatibility problem is to maximize $|\{c_i : l_{ij} = 1 \text{ for all } j \in R_{c_i}\}|$. A 1-phylogeny is called a *perfect phylogeny*. All three problems (ℓ -phylogeny for $\ell \geq 1$, parsimony, compatibility) are \mathcal{NP} -complete [1, 8, 4, 6, 13]. Parsimony, ℓ -phylogeny, and compatibility all allow states of a character to evolve multiple times. However, both parsimony and compatibility allow some characters to evolve many times. The ℓ -phylogeny metric requires *balanced* evolution, in that no one character can pay for most of the evolutionary changes. Thus, ℓ -phylogeny is a better measure than parsimony or compatibility in biological situations in which all characters are believed to evolve slowly.

In this paper we consider the *fixed-topology* variant of the ℓ -phylogeny problem, where in addition to the species set and characters, we are also given a tree T in which internal nodes are unlabelled, each leaf is labelled with a species $s \in S$ and each species $s \in S$ is the label of exactly one leaf of T . The *fixed-topology ℓ -phylogeny problem* is the problem of determining labels for the internal nodes so that the resulting phylogeny is an ℓ -phylogeny, or determining that such a labelling does not exist.

Fixed-topology algorithms can be used as *filters*. Current phylogeny-producing software can generate thousands of trees which are (approximately) equally good under some metric such

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, make any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.



DISCLAIMER

Portions of this document may be illegible in electronic image products. Images are produced from the best available original document.



as maximum likelihood or parsimony. We can think of these outputs as proposed topologies. One way to differentiate these hypotheses is to see which topologies also have low phylogenetic number. For example, the original trees can be generated by biomolecular sequence data, and they can then be filtered using morphological data with slowly-evolving traits.

It will be convenient to allow a node to remain unlabelled in one or more characters in a fixed topology. In this case, the node disagrees with all of its neighbors on all unlabelled characters. We can easily extend such a labelling to one in which every node is labelled without increasing ℓ_{ij} for any i or j : for any character j , for each connected component of nodes which are not labelled, choose any neighbouring node v which is labelled i_v and label the entire component with the state i_v . This does not introduce any extra component for i_v , nor does it break components of any other state that weren't already broken.

In the fixed-topology setting, optimal trees for the parsimony and compatibility metrics can be found in polynomial time [7]. The fixed-topology ℓ -phylogeny problem can be solved in polynomial time for $\ell \leq 2$, but is \mathcal{NP} -complete for $\ell \geq 3$ even for degree-3 trees in which no state labels more than $\ell + 1$ leaves (and therefore there is a trivial $\ell + 1$ phylogeny) [8]. Fitch's algorithm for parsimony uses dynamic programming. Dynamic programming also gives good algorithms in some cases for finding phylogenies when characters are polymorphic [2].

Jiang, Lawler, and Wang [10] consider the fixed-topology tree-alignment problem, where species are represented as biomolecular sequences, the cost of an edge in the tree is the edit distance between the labels at its endpoints, and the goal is to minimize the sum of the costs over all edges. They give a 2-approximation for bounded-degree input topologies and extend this to obtain a polynomial-time approximation scheme (PTAS). In Lemma 3 of [10], they prove that the best lifted tree (in which the label of each internal node is equal to the label of one of its children) is within a factor of 2 of the best tree with arbitrary labels. The proof only uses the triangle inequality (it does not use any other facts about the cost measure). Therefore, the result holds for several other cost measures, including ℓ -phylogeny, parsimony, and the minimum-load cost measure for phylogenies with polymorphic characters which was introduced in [2]. It also holds for the variant of ℓ -phylogeny in which ℓ_i is specified for each character c_i . This variant was introduced in [8]. We refer to it as the *generalized ℓ -phylogeny problem*. In fact, Lemma 3 of [10] holds for the fixed-topology problem with *arbitrary* input topologies, though the authors do not state this fact since they do not use it. Despite the applicability of Lemma 3, the algorithmic method of Jiang, Lawler and Wang does not seem to be useful in developing approximation algorithms for the fixed-topology ℓ -phylogeny problem (or for related problems). Jiang et al. use dynamic programming to find the minimum-cost lifted tree. Dynamic programming is not efficient for the more global metric of ℓ -phylogeny. The dynamic programming proceeds by computing an optimal labelling for a subtree for each possible labelling of the root of the subtree. For metrics where cost is summed over edges (such as parsimony or tree alignment), one only needs to find the lowest-cost labelling for a given root label. For the ℓ -phylogeny problem, the cost of a tree depends upon how many times each state is broken for a given character. One

cannot tell *a priori* which state will be the limiting one. Therefore, instead of maintaining a single optimal tree for each root label, we must maintain all trees whose cost (represented as a vector of components for each state) is undominated. This number can be exponential in r , the number of states, even for bounded-degree input trees. This is a common theme in combinatorial optimization: the more global nature of minimax makes it harder to compute than summation objectives, but also more useful.

Gusfield and Wang [14] take the approach of [10] a step further by proving that the best uniform lifted tree (ULT) is within a factor of 2 of the best arbitrarily-labelled tree. In a uniform lifted tree on each level, all internal nodes are labeled by the same child (e.g. all nodes at level one take the label of their leftmost child). This proof also extends to the ℓ -phylogeny metric. If the input tree is a complete binary tree, then there are only n ULTs, and exhaustive search is efficient, giving an algorithm which is faster than ours and has an equivalent performance bound. However, when the input tree isn't complete (even if it is binary), Gusfield and Wang use dynamic programming to find the minimum-cost ULT, and so their method fails when it is applied to the ℓ -phylogeny problem. Wang, Jiang, and Gusfield recently improved the efficiency of their PTAS for tree alignment [15], but still use dynamic programming.

We give a simple 2-approximation for the fixed-topology ℓ -phylogeny problem that works for arbitrary input topologies. It is based on rounding the linear-programming relaxation of an integer programming formulation for fixed-topology ℓ -phylogeny. To our knowledge, this is the first application of linear-programming technology to phylogeny problems.

As we described earlier, ℓ -phylogeny is most appropriate for slowly-evolving characters. It is most restrictive (and hence most different from parsimony) when ℓ is small. Therefore, we look more closely at the first *NP*-hard case: $\ell = 3$. For this case, we give an algorithm based upon the structure of a 3-phylogeny that will construct a 4-phylogeny if the input instance has a 3-phylogeny.

The remainder of our paper is organized as follows: in section 2, we give the 2-approximation algorithm for the ℓ -phylogeny problem. In section 3 we give the optimal approximation algorithm for inputs with 3-phylogenies. Section 4 discusses extensions of our results to a problem with polymorphic characters.

2. A 2-approximation algorithm for the fixed-topology phylogenetic number problem. The interaction between characters in phylogeny problems affects the choice of the topology, but it does not affect the labelling of the internal nodes once the topology is chosen. Thus, for this problem, we can consider each character separately.

Let $c : S \rightarrow \{1, \dots, r\}$ be a character and let T be a tree with root q and leaves labelled by character states $1, \dots, r$. For each state i , let T_i be the subtree of T consisting of all the leaves labelled i and the minimum set of edges connecting these leaves. Let $L(T_i)$ be the set of leaves of T_i , and let rt_i , the root of T_i , be the node of T_i closest to the root of T . The *important nodes* of T_i are the leaf nodes and the nodes of degree greater than 2. An *i -path* p of T_i is a sequence of edges of T_i that connects two important nodes of T_i , but does not pass through any other

important nodes. The two important nodes are referred to as the *endpoints* of p , and the other nodes along the i -path are said to be *on* p (an i -path need not have any nodes on it). Although the edges of the tree T are undirected, we will sometimes use the notation $(v \rightarrow w)$ for an edge or i -path with endpoints v and w , to indicate that v is nearer to the root of T than w (v is the *higher* endpoint and w is the *lower* endpoint); otherwise we will write edges and i -paths as (v, w) . If the label of the upper endpoint v or some node on the i -path $p = (v \rightarrow w)$ differs from the label for the lower endpoint w , then we say that p *breaks* state i .

Given a tree T with each node labeled from the set $\{1, \dots, r\}$, we need a way to count the number of components induced by the nodes labeled i . Since the tree is rooted, we can assign each connected component a root, namely the node closest to the root of T . We then count the number of roots for components labelled i . A node is the root of its component if its label differs from that of its parent. The root q , which has no parent, is also the root of its component. Therefore we have the following:

OBSERVATION 2.1. *Let T be a tree with its leaves and internal nodes labelled by elements of $\{1, \dots, r\}$. For each i , let T_i be defined as above, and let q be the root of tree T . Then the number of connected components induced by the nodes labelled i is $|\{e = (v \rightarrow w) : c(v) \neq i, c(w) = i\}| + Y_i$, where $Y_i = 1$ if q is labelled i and 0 otherwise.*

We now define an integer linear program (ILP) which solves the fixed-topology ℓ -phylogeny problem. The linear-programming relaxation of this ILP is the key to our 2-approximation algorithm. The integer linear program \mathcal{I} uses the variables $X_{v,i}$, for each state $i \in \{1, \dots, r\}$, and each node v in the tree T , and the variables $X_{p,i}$ and $cost_{p,i}$, for each state i , and each i -path p of T_i . These variables have the following interpretation:

$$\begin{aligned} X_{v,i} &= \begin{cases} 1 & \text{if node } v \text{ is labelled } i \\ 0 & \text{otherwise} \end{cases} \\ X_{p,i} &= \begin{cases} 1 & \text{if all nodes on } p \text{ are labelled } i \\ 0 & \text{otherwise} \end{cases} \\ cost_{p,i} &= \begin{cases} 1 & \text{if the lower endpoint of } p \text{ is the root of a component of state } i \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

ILP \mathcal{I} is defined as follows:

minimize ℓ

subject to

- (1) $X_{v,i} = 1$ for each leaf $v \in T_i$, $i = 1, \dots, r$
- (2) $X_{v,i} = 0$ if $v \notin T_i$, $i = 1, \dots, r$
- (3) $\sum_{i=1}^r X_{v,i} \leq 1$ $\forall v \in T$
- (4) $X_{p,i} = X_{v,i}$ $i = 1, \dots, r, \forall p \in T_i, \forall v \in p$

- (5) $X_{p,i} \leq X_{v,i} \quad i = 1, \dots, r, \forall p \in T_i, \text{ each endpoint } v \in p$
(6) $cost_{p,i} \geq X_{v,i} - X_{p,i} \quad i = 1, \dots, r, \forall p \in T_i \text{ with lower endpoint } v$
(7) $\sum_p cost_{p,i} + X_{rt_i,i} \leq \ell \quad i = 1, \dots, r$
(8) $X_{v,i}, X_{p,i}, cost_{p,i} \in \{0, 1\}$

Constraint (8) assures that the cost ($cost_{v,i}$), i -path ($X_{p,i}$), and vertex ($X_{v,i}$) variables serve as indicator variables in accordance with their interpretation. Constraint (1) labels the leaves in accordance with the input. Constraint (2) prohibits labelling a node v with a state i when v is not in T_i (the number of components labelled i could not possibly be reduced by this labelling). Constraint (3) ensures that each internal node will have no more than one label. Constraints (4) and (5) ensure that for each tree T_i , nodes on paths are taken all-or-none; if any node on an i -path p (including endpoints) is lost to a state i , then it does no good to have any of the other nodes on the path (though it may be beneficial to maintain one or both endpoints). Constraint (6) computes the path costs, and constraint (7) ensures that each state has no more than ℓ connected components. This is an implementation of Observation 2.1. Since there is no i -path in T_i with rt_i as its lower endpoint, we must explicitly check the root of each tree T_i , just as we checked the global root in Observation 2.1.

Integer program \mathcal{I} solves the fixed-topology ℓ -phylogeny problem. We only require, however, that the optimal value of ℓ given by \mathcal{I} is a lower bound on the phylogenetic number of tree T with the given leaf labelling (proof omitted):

PROPOSITION 2.2. *If there exists an ℓ -phylogeny for tree T with a given leaf labelling, then there is a feasible solution for the integer linear program for this value of ℓ .*

Integer linear programming is \mathcal{NP} -hard in general [3, 9, 11], so we cannot solve it directly in polynomial time. (In fact, doing so would solve the fixed-topology ℓ -phylogeny problem, which we know to be \mathcal{NP} -hard for $\ell \geq 3$ from [8].) However, we can solve the linear-programming relaxation \mathcal{L} of \mathcal{I} , which consists of all the constraints of \mathcal{I} except that Constraint (8) is replaced by the constraint $0 \leq X_{v,i}, X_{p,i}, cost_{p,i} \leq 1$ (8').

THEOREM 2.3. *If there is a solution for the linear program \mathcal{L} for a fixed topology T with leaves labelled with states from $\{1, \dots, r\}$, then we can assign states to the internal nodes of T such that no state $i \in \{1, \dots, r\}$ has more than 2ℓ components.*

Proof. The 2ℓ phylogeny for the character $c : S \rightarrow \{1, \dots, r\}$ on T is constructed by assigning states to the nodes of each tree T_i based on the $X_{v,i}$ values. For each state $i \in \{1, \dots, r\}$, consider each internal node v of T_i . A node v is labelled i if and only if $X_{v,i} > 1/2$, and there is path $v, w_1, w_2, \dots, w_k, v^*$ through tree T_i to a leaf v^* of T_i where $X_{w_j,i} > 1/2$ for all $j = 1, \dots, k$. If $X_{v,i} > 1/2$, but there is no such path, then node v is *isolated*, and by our procedure remains unlabelled. A node v also remains unlabelled if $X_{v,i} \leq 1/2$ for all states i .

To show that the labelling is a 2ℓ -phylogeny, we show that each component of state i adds at least $1/2$ to the sum $(\sum_p cost_{p,i}) + X_{rt_i,i}$. From Observation 2.1, the number of connected components for the state i is $|\{e = (v \rightarrow w) : c(v) \neq i, c(w) = i\}| + Y_i$, where Y_i is 1 if q has

state i (and therefore $q = rt_i$) and 0 otherwise. Constraints (5) and (4) ensure that if the edge $e = (v \rightarrow w)$ has $c(v) \neq i$ and $c(w) = i$ then either w is the root of T_i , or w must be an endpoint node with $X_{w,i} > 1/2$, and that either $X_{v,i} \leq 1/2$ or v is isolated. However, since w is labelled i , w must not be isolated, and therefore v would not be isolated if $X_{v,i}$ was greater than $1/2$. So $X_{v,i} \leq 1/2$, and $X_{p,i} \leq 1/2$ for the i -path p with lower endpoint w . Therefore we need only calculate the number of i -paths p with a lower endpoint w such that $X_{p,i} \leq 1/2$, $X_{w,i} > 1/2$, and w is not isolated.

Suppose $p = (v \rightarrow w)$ is such an i -path. Since w is not isolated and the node above w is not labelled i , there is a sequence $p_1 = (w \rightarrow v_1)$, $p_2 = (v_1 \rightarrow v_2), \dots, p_j = (v_{j-1} \rightarrow v_j)$ of i -paths of T_i such that $X_{p,i} > 1/2$ for every $p \in \{p_1, \dots, p_j\}$ and $X_{v,i} > 1/2$ for every $v \in \{v_1, \dots, v_j\}$, and v_j is a leaf of T_i . Calculating $cost_{p,i} + cost_{p_1,i} + \dots + cost_{p_j,i} = (X_{w,i} - X_{p,i}) + (X_{v_1,i} - X_{p_1,i}) + \dots + (X_{v_j,i} - X_{p_j,i}) = -X_{p,i} + (X_{w,i} - X_{p_1,i}) + (X_{v_1,i} - X_{p_2,i}) + \dots + (X_{v_{j-1},i} - X_{p_j,i}) + X_{v_j,i}$, we know by constraints (5) that $X_{w,i} - X_{p_1,i} \geq 0$, $X_{v_1,i} - X_{p_2,i} \geq 0$, \dots , $X_{v_{j-1},i} - X_{p_j,i} \geq 0$. So $cost_{p,i} + cost_{p_1,i} + \dots + cost_{p_j,i} \geq X_{v_j,i} - X_{p,i} = 1 - X_{p,i} \geq 1/2$.

Note also that for any two i -paths $p = (v \rightarrow w)$ and $p' = (v' \rightarrow w')$ which break i , the i -labelled paths to leaves are disjoint (because they are in separate components of i). Therefore each i -path p which breaks i in our construction contributes at least $1/2$ to the sum $(\sum_p cost_{p,i})$. If rt_i is labelled i (and hence the root of a component of i), then $X_{rt_i,i} > 1/2$ (corresponding to an edge $(v \rightarrow rt_i)$ in T or to the case $Y_i = 1$). So $2 \times ((\sum_p cost_{p,i}) + X_{rt_i,i}) \geq |\{e = (v \rightarrow w) : c(v) \neq i, c(w) = i\}| + Y_i$, and therefore $2\ell \geq |\{e = (v \rightarrow w) : c(v) \neq i, c(w) = i\}| + Y_i$.

□

Recall that we have considered each character separately in our 2-approximation algorithm. Thus, our work applies to the generalized ℓ -phylogeny problem (and not just to the ordinary ℓ -phylogeny problem). In particular, we have the following theorem.

THEOREM 2.4. *There is a 2-approximation algorithm for the generalized ℓ -phylogeny problem.*

3. 4-phylogeny algorithm. In this section we give an algorithm which takes a fixed-topology phylogeny instance with arbitrary topology and, as long as it has a 3-phylogeny, finds a 4-phylogeny for the instance.

We use the following definitions, in addition to those that we used for the 2-approximation. We will maintain a forest F_i for every state i . A **branch point** of F_i is a node in F_i with degree 3. We say that a node $v \in F_i$ is *claimed* by state i if it is not in F_j for any $j \neq i$.

The algorithm generalizes the fixed-topology 2-phylogeny algorithm of [8]. It consists of a *forced phase* and then an *approximation phase*. The forced phase produces a partial labelling which can still be extended to a 3-phylogeny; it makes no labelling decisions that are not forced if one is to have a 3-phylogeny. The approximation phase removes all remaining contention for labels, but it can break some states into four pieces. Because finding a fixed-topology 3-phylogeny is \mathcal{NP} -complete [8], this is an optimal approximation algorithm for phylogeny instances with 3-phylogenies.

3.1. The Forced Phase of the Algorithm. Initially, for every state i we will have $F_i = T_i$. During the forced phase of the algorithm, nodes will be removed from the forests F_i . The invariant during the forced phase of the algorithm is that there is a 3-phylogeny in which every node v is assigned a state j such that $v \in F_j$. The forced phase applies the following rules in any order until none can be applied. If any forest F_i is broken into more than three components by the application of these rules, then the instance has no 3-phylogeny and the algorithm terminates.

1. For any i -path (v, w) , let S be the set containing v and w and the nodes on the i -path. If S contains two or more branch points of F_j then every node on the i -path is removed from F_i . (Note that in the updated copy of F_i (after the rule is applied), v and w will have lower degree than in the original F_i . Furthermore, if v has degree 2 in the updated F_i then the i -path containing it will consist of nodes from two different i -paths in the original F_i . Similarly, i -paths can be merged as a result of the following rules.)
2. If F_i has $C(F_i)$ connected components and F_i contains a node v of degree at least $5 - C(F_i)$ then in every forest F_j with $j \neq i$, v and all nodes on j -paths adjacent to v are removed from F_j (i.e., i claims node v).
3. If v is a branch point of F_i and two i -paths (v, w_1) and (v, w_2) adjacent to v contain branch points of F_j (for some $j \neq i$), then in every forest F_k with $k \neq i$, every branch point $w \notin \{v, w_1, w_2\}$ of F_i and every node on every k -path adjacent to w is removed from F_k (i.e., F_i claims all branchpoints except v, w_1 , and w_2).

Rule 1 is justified by observing that in any 3-phylogeny, each forest F_i gives up at most two disjoint i -paths, or a single branchpoint with the i -paths adjacent to it. In the setting in which rule 1 is applied, if F_i were to claim the path in question, then F_j would lose two branchpoints and necessarily be in at least four components. Therefore, in any 3-phylogeny for the input, F_i cannot have that i -path. Note that once any node on an i -path is lost to F_i , then F_i has no reason to claim any other nodes on the i -path.

Rule 2 is justified by the following observations. If there is a node of degree at least 4 in tree T_i , then it must be labelled i in any 3-phylogeny (losing it will break state i into at least 4 pieces). Once F_i has been forced to give up an i -path, it cannot give up another branchpoint. Finally, once F_i has been forced into three pieces, then it must claim all remaining nodes in F_i .

Rule 3 is applied when we isolate a region where a break in F_i must occur, but do not yet know exactly where the break will occur. If two paths adjacent to a branchpoint of F_i contain branchpoints of F_j , then by the previous argument for rule 1, F_i cannot keep both of those paths. Therefore, outside of the affected region (those two i -paths), F_i can act as though the forest has been cut into at least two pieces, and can claim all branchpoints.

3.2. The Approximation Phase of the Algorithm. In the following, *releasing* a degree-2 node $v \in F_i$ removes all nodes on its i -path from F_i . Releasing a higher-degree node $v \in F_i$ removes v and all nodes on i -paths adjacent to v from F_i . The approximation phase consists of the following steps.

1. For each connected component C of F_i , if the root of C is unclaimed then F_i releases the root of C . Also, if this root has degree 2, F_i releases any unclaimed branch points at the ends of i -paths adjacent to this root.
2. If, after the forced phase, F_i is in a single component with exactly one unclaimed branch point, w , then it releases w .
3. If, after the forced phase, F_i is in a single component with exactly two unclaimed branch points, w_1 and w_2 which are the two endpoints of an i -path, and the path from the root to w_2 passes through w_1 , then F_i releases w_2 .
4. If, after the forced phase, F_i is in a single component with exactly three unclaimed branch points, w_1 , w_2 and w_3 where there is an i -path from w_1 to w_2 and an i -path from w_2 to w_3 , then F_i releases w_2 .

3.3. The Proof of Correctness. The proof requires the following observation (proof omitted), and follows from Lemma 3.2 and Lemma 3.7.

OBSERVATION 3.1. *If F_i is in one component and it releases two branchpoints w_1 and w_2 which share an i -path, then the resulting forest F_i has at most 4 components.*

LEMMA 3.2. *At the end of the approximation phase, every forest F_i has at most 4 connected components.*

Proof. The forest F_i can be in at most three components at the end of the forced phase. If F_i is in three components at the end of the forced phase, then, by Rule 2 of the forced phase, every remaining node in F_i is claimed during the forced phase, so nothing is removed from F_i during the approximation phase. If F_i is in two components after the forced phase, then, again by Rule 2, all branch points of F_i are claimed during the forced phase, so no branch points are removed from F_i during the approximation phase. Step 1 of the approximation phase, therefore, will remove at most one path from each component (when the root has degree 2, since degree-3 roots are claimed) and therefore breaks F_i into at most four components. In this case Steps 2–4 of the approximation phase do not apply, and at most one of Steps 2–4 can apply to each of the remaining cases.

If F_i is in one component after the forced phase, and Steps 2–4 do not apply then we have two cases. If the root is degree three, Step 1 results in at most 3 components. If the root is degree two, then F_i could release the two branchpoints on either end of this i -path, resulting in at most four components by Observation 3.1.

Suppose F_i is in one component with exactly one unclaimed branchpoint after the forced phase. If that branchpoint is released by Step 1, then F_i is in at most 3 components after that step (only that branchpoint and its adjacent i -paths are removed from F_i), and Step 2 is redundant. Otherwise, Step 1 releases only the i -path through the root, since its endpoints are claimed, resulting in two components, and the subsequent application of Step 2 adds at most two more for a total of four.

Suppose Step 3 can be applied to F_i . If w_1 is not an endpoint of the i -path through the root of F_i , then Step 1 results in two components (both endpoints are claimed since neither is

w_1 or w_2). Subsequently removing w_2 by Step 3 results in at most two more components for a total of four. If w_1 is an endpoint of the i -path containing the root of F_i , then both w_1 and w_2 are released (and nothing more). Since they share an i -path, this results in at most four components by Observation 3.1.

Finally, suppose Step 4 can be applied to F_i . If none of w_1, w_2 or w_3 is the root or is an endpoint of the i -path adjacent to the root, then Step 1 will result in additional components only if the root has degree two. Since both of the endpoints of this i -path are claimed in this case, removing this i -path and w_2 (by Step 4) results in at most four components. Otherwise, the combined application of Steps 1 and 4 requires the release of w_2 , and possibly one of w_1 and w_3 as well (but not both, since if w_2 is the root, neither of the other branchpoints will be released). By Observation 3.1 this will result in at most four components. \square

The following lemmas use this fact:

FACT 3.3. ([8]) *The intersection of two subtrees of a tree is connected and contains the root of at least one of the subtrees.*

LEMMA 3.4. *If F_i and F_j are each in two components after the forced phase, then after the approximation phase, there is no node that is in F_i and in F_j .*

Proof. This proof is similar to the correctness proof of the fixed-topology 2-phylogeny algorithm in [8]. Let C_i be a component of F_i after the forced phase, and let C_j be a component of F_j after the forced phase. Since F_i and F_j are both split in two components during the forced phase, all branch points in C_i and C_j are claimed during the forced phase (by Rule 2), and their intersection is a path in the fixed topology (i.e., all nodes are degree 2). Furthermore, the root of C_i or C_j is in the intersection. Therefore, the contention is cleared in Step 1 of the approximation phase of the algorithm. \square

LEMMA 3.5. *If F_i is in one component after the forced phase, and F_j is in two components after the forced phase, then there is no node that is in F_i and in F_j after the approximation phase.*

Proof. First note that no branch point of T_j is part of T_i . (Since F_j is in only 2 pieces, it gave up only degree-2 nodes in the forced phase, and subsequently claimed all branch points. None of these are in F_i , since F_i was unbroken in the forced phase). Thus, the intersection of T_i and T_j is a path in T_j . We conclude that the intersection of F_i and F_j is a path in F_j and contains at most one branch point of F_i (otherwise, the path would be removed from F_j during the forced phase by Rule 1). If the intersection contains the root of F_j , then the contention will be removed during Step 1 of the approximation phase. Otherwise, the intersection contains the root of F_i . Thus, the single branch point of F_i that is contained in the intersection of F_i and F_j is either the root of F_i or it is an endpoint of the i -path containing the root of F_i . In either case, the contention will be removed in Step 1 of the approximation phase. \square

LEMMA 3.6. *If F_i and F_j are each in one component after the forced phase, then there is no node that is in F_i and in F_j after the approximation phase.*

Proof. (sketch) The proof is by case analysis. Case (α, β, γ) represents the situation in

which the intersection of F_i and F_j after the forced phase contains α branch points of F_i and β branch points of F_j , γ of which are shared.

Cases $(0, 0, 0)$, $(1, 0, 0)$, and $(1, 1, \gamma)$ use arguments similar to those in Lemmas 3.4 and 3.5. Cases $(2, 0, 0)$, $(2, 1, 1)$, $(3, 3, 0)$, $(\alpha > 1, \beta > 1, \gamma > 0)$, and $(3, \beta \leq 1, \gamma)$ are all forbidden by rule 1 of the forced phase. We now give a proof of case $(2, 1, 0)$. The remaining cases $(2, 2, 0)$ and $(3, 2, 0)$ use similar arguments.

Case $(2, 1, 0)$: By Rule 1, after the forced phase, the intersection of F_i and F_j has the branch-point of F_j (w_2), between the two branchpoints for F_i (w_1 and w_3). Let w_4 be the other endpoint of the j -path adjacent to w_2 that contains w_1 , and let w_5 be the other endpoint of the j -path adjacent to w_2 that contains w_3 . By Rule 3 of the forced phase, all branch points of F_j except w_2 , w_4 and w_5 are claimed. Suppose that w_2 is not released by F_j during the approximation phase. Then by Rules 2 and 4 of the approximation phase, exactly one of $\{w_4, w_5\}$ is unclaimed after the forced phase. Suppose without loss of generality this is w_4 . Because w_2 is not released by F_j in the approximation phase, the root of F_j is not w_2 or on any j -path adjacent to it. Therefore the root of F_i is in the intersection. If the root of F_i is on the i -path between w_1 and w_3 then by Step 1, F_i will release both w_1 and w_3 , and the contention will be cleared. If the root of F_i was on one of the other i -paths adjacent to w_1 , then F_j would have released w_2 by step 3. Finally, if the root is on an i -path adjacent to w_3 (but not w_1), F_i will release w_3 by Step 1, clearing the i -path from w_3 to w_1 (not including w_1). By Step 3 F_j will release w_4 , clearing the j -path from w_2 to w_4 (not including w_2). Therefore the contention is removed. \square

LEMMA 3.7. *After the approximation phase, every node is in at most one forest F_i .*

Proof. The lemma follows from Lemmas 3.4, 3.5 and 3.6. \square

4. Approximating Polymorphism. A **polymorphic character** (see [12]) allows more than one state per character per species. This type of character has strong application in linguistics [2, 16]. If there are r states, a polymorphic character is a function $c : S \rightarrow (2^{\{1, \dots, r\}} - \emptyset)$. For a given set of species, the *load* is the maximum number of states for any character for any species.

We consider a simple model where the *loss* of a state from parent to child costs nothing, and *mutation* of a state (changing from parent to child), costs a fixed amount c_m . Given a tree labelled uniquely by species with polymorphic characters, we wish to determine labels for the internal nodes which minimize the total mutation in the tree (sum over all edges of the number of mutations across the edge), and achieves a given load bound ℓ . Bonet et. al. [2], considered more complicated models motivated by [12], but they proved that even this simple case is \mathcal{NP} -complete even if the input tree is binary. We can extend the results of section 2 to obtain, for any $\alpha > 1$, an $(\alpha, \frac{\alpha}{\alpha-1})$ -approximation algorithm for this problem with arbitrary input topology, meaning that if there is a load- ℓ phylogeny of cost c , we can find a load- $\alpha\ell$ phylogeny with cost at most $(\alpha/(\alpha-1))c$. (Note that taking $\alpha = 2$ gives a $(2, 2)$ -approximation algorithm.)

REFERENCES

- [1] H. Bodlaender, M. Fellows, T. Warnow, "Two Strikes Against Perfect Phylogeny", *procs. of the 19th International Congress on Automata, Languages and Programming (ICALP)*, pp. 273-287, Springer-Verlag Lecture Notes in Computer Science, 1992.
- [2] M. Bonet, C. Phillips, T.J. Warnow and S. Yooseph, Constructing Evolutionary Trees in the Presence of Polymorphic Characters, *Proceedings of the 28th Annual ACM Symposium on the Theory of Computing* (1996).
- [3] I. Borosh and L.B. Treybig, Bounds on positive integral solutions of linear Diophantine equations, *Proceedings of the American Mathematical Society*, Vol 55 (1976).
- [4] W.H.E. Day, Computationally difficult parsimony problems in phylogenetic systematics, *Journal of Theoretical Biology*, Vol 103 (1983).
- [5] W.H.E. Day, D.S. Johnson and D. Sankoff, The computational complexity of inferring phylogenies by parsimony, *Mathematical biosciences*, Vol 81 (1986).
- [6] W.H.E. Day and D. Sankoff, "Computational complexity of inferring phylogenies by compatibility", *Systematic Zoology*, 35(2): 224-229, 1986.
- [7] W. Fitch, Towards defining the course of evolution: minimum change for a specified tree topology, *Systematic Zoology*, Vol 20 (1971).
- [8] L.A. Goldberg, P.W. Goldberg, C.A. Phillips, E. Sweedyk and T. Warnow, Minimizing phylogenetic number to find good evolutionary trees, *Discrete Applied Mathematics*, to appear.
- [9] M.R. Garey and D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-completeness*, W.H. Freeman and Company (1979).
- [10] T. Jiang, E.L. Lawler and L. Wang, Aligning Sequences via an Evolutionary Tree: Complexity and Approximation, *Proceedings of the 26th Annual ACM Symposium on the Theory of Computing* (1994).
- [11] R.M. Karp, Reducibility among combinatorial problems, *Complexity of Computer Computations*, eds. R.E. Miller and J.W. Thatcher, Plenum Press (1972).
- [12] M. Nei, *Molecular Evolutionary genetics*, Columbia University Press, New York (1987).
- [13] M.A. Steel, "The complexity of reconstructing trees from qualitative characters and subtrees", *Journal of Classification*, 9 91-116, 1992.
- [14] L. Wang and D. Gusfield, Improved Approximation Algorithms for Tree Alignment, *Proceedings of CPM 1996*, 220-233.
- [15] L. Wang, T. Jiang, and D. Gusfield, "A more efficient approximation scheme for tree alignment", To appear in *Proceedings of the of First Annual International Conference on Computational Molecular Biology*, Jan. 1997.
- [16] T. WARNOW, D. RINGE AND A. TAYLOR, *A character based method for reconstructing evolutionary history for natural languages*, Tech Report, Institute for Research in Cognitive Science, 1995, and *Proceedings 1996 ACM/SIAM Symposium on Discrete Algorithms*.