

## Constructing Computer Virus Phylogenies

Leslie Ann Goldberg

*University of Warwick, Coventry CV4 7AL, United Kingdom\**

Paul W. Goldberg

*University of Warwick, Coventry CV4 7AL, United Kingdom<sup>†</sup>*

Cynthia A. Phillips

*Sandia National Labs, P.O. Box 5800, Albuquerque, New Mexico 87185<sup>‡</sup>*

and

Gregory B. Sorkin

*IBM T. J. Watson Research Center, P.O. Box 218, Yorktown Heights, New York 10598*

Received March 26, 1996; revised August 7, 1997

There has been much recent algorithmic work on the problem of reconstructing the evolutionary history of biological species. Computer virus specialists are interested in finding the evolutionary history of computer viruses—a virus is often written using code fragments from one or more other viruses, which are its immediate ancestors. A phylogeny for a collection of computer viruses is a directed acyclic graph whose nodes are the viruses and whose edges map ancestors to descendants and satisfy the property that each code fragment is “invented” only once. To provide a simple explanation for the data, we consider the problem of constructing such a phylogeny with a minimum number of edges. This optimization

\* Part of this work was performed at Sandia National Laboratories and was supported by the U.S. Department of Energy under contract DE-AC04-76AL85000. Part of this work was supported by the ESPRIT Basic Research Action Programme of the EC under contract 7141 (project ALCOM-IT).

<sup>†</sup> Part of this work was performed at Sandia National Laboratories and was supported by the U.S. Department of Energy under contract DE-AC04-76AL85000.

<sup>‡</sup> This work was performed under U.S. Department of Energy contract DE-AC04-76AL85000.

problem is NP-hard, and we present positive and negative results for associated approximation problems. When tree solutions exist, they can be constructed and randomly sampled in polynomial time. © 1998 Academic Press

## 1. INTRODUCTION

There are now several thousand different computer viruses in existence, with new ones being written at a rate of 3 to 4 per day. Most of these are based upon previous ones: someone copies and modifies a virus, or creates a new virus with subroutines borrowed from one or more **ancestors**.

For most purposes, a computer virus can be regarded as a fixed string of bytes, each byte consisting of eight bits. If one virus is based on another, long substrings of the ancestor, say 20 bytes or more, will appear in the descendant. Using probability models similar to those employed in speech recognition it is possible to estimate the probability that a given byte string occurs in several viruses *by chance* [15]; if the probability is low but the string does occur in several viruses then we conclude that it was written for one virus, and it was copied into the others.

We wish to infer an evolutionary or phylogenetic history for a set of computer viruses. As most phylogenetic literature to date has been based upon biological evolution, we adopt that terminology. Thus, the viruses in the input set  $\mathcal{S} = \{s_1, \dots, s_n\}$  are called **species**. The species are defined by a set of **binary characters**  $\mathcal{E} = \{c_1, \dots, c_k\}$ . A binary character is a function  $c: \mathcal{S} \rightarrow \{0, 1\}$ . (In general, the range of a character can be arbitrary, but the presence or absence of byte strings can be modeled with binary characters.) Each character  $c$  corresponds to a byte string, with  $c(s) = 1$  if the string occurs in species  $s$  and  $c(s) = 0$ , otherwise. If  $c(s) = 1$ , we say that species  $s$  *has* or *contains* character  $c$ . In analogy with terminology from the logic synthesis area of computer circuit design, we define the **on-set**  $S_c$  of a character  $c$  to be the set of all species on which its value is 1:  $S_c = \{s \in \mathcal{S} \mid c(s) = 1\}$ . A character  $c$  is **trivial** if  $|S_c| \leq 1$ . A trivial character can be ignored because it imposes no constraints on possible solutions.

We assume that the input species are all related: that the bipartite graph joining characters to species that have them is connected. Otherwise, the connected components can be considered independently.

We also assume that each code fragment is invented only once. For sufficiently long fragments this is justified by differences in programming style, the many possible orderings of unconstrained events, etc. We model the evolution of a set of viral species with a directed graph in which an edge  $s_i \rightarrow s_j$  indicates that species  $s_i$  is an ancestor of species  $s_j$  (i.e.,  $s_j$  inherited some character(s) from  $s_i$ ).

**DEFINITION 1.** A **phyloDAG** for input species  $\mathcal{S}$  and characters  $\mathcal{E}$  is a directed acyclic graph (DAG) with node set  $\mathcal{S}$ . For each character  $c \in \mathcal{E}$ , the subgraph induced by on-set  $S_c$  is connected, in the sense that from a single species  $a_c \in S_c$ , which we call the **archetype**, there is a directed path, within  $S_c$ , to every other  $s \in S_c$ .

The phyloDAG model allows the possibility that a species may be derived from several ancestors rather than from a single ancestor. We will explain the motivation behind this new degree of freedom right after some brief comments on the mathematics of the model.

A phyloDAG exists for any inputs  $(\mathcal{S}, \mathcal{E})$ : for any chronology ascribed to the species (i.e., any total ordering of the species set), the directed graph with edges from each species to all later species is a phyloDAG. However, every pair of species is related by an edge in this graph. Because most virus species presumably have few ancestors, we seek a **minimum phyloDAG**, one with a minimum number of directed edges.

We assume that the input is given in the following compact format: for each species  $s \in S$ , we are given a list of the characters  $c$  for which  $c(s) = 1$ .

**DEFINITION 2.** The **input length**  $l = l(\mathcal{S}, \mathcal{E}) = \sum_{c \in \mathcal{E}} |S_c|$ . The **size** is  $n = |\mathcal{S}|$ . The number of characters is  $k = |\mathcal{E}|$ .

Our approach to the evolution problem corresponds to a so-called *restricted* model of evolution: one in which we are not allowed to introduce hypothetical species outside of the input set. This model is well suited to computer viruses, where because of good worldwide communications, sharing of data between antivirus organizations, and the brief history involved, there are likely to be very few gaps in our viral database—a situation quite different from that in biology. Previous work on restricted models of evolution will be discussed in Section 1.4. For our model, if additional species could be introduced into a phyloDAG, there would always be a trivial sparse phyloDAG: a star graph with the center an added species  $s$  such that  $c(s) = 1$  for  $c \in \mathcal{E}$ .

### 1.1. Problem Motivation

Sorkin's study of computer virus evolution [18] motivated our study of the phyloDAG model. There are about 6000 computer virus species in existence, of which many are simple modifications of predecessors. The Jerusalem, Vienna, and Blackjack virus families, for instance, each contain from scores to hundreds of related species. The author of a computer virus can equally well incorporate computer code (instructions) from several existing viruses, which is how multiple ancestry arises. Experts disagree as to the frequency with which this occurs, and one of our eventual aims is to

resolve this issue. (Another form of multiple ancestry is well established, but not addressed here. It comes from virus “toolkits”: collections of mix-and-match software components from which viruses can be assembled.)

The evolutionary classification of computer viruses can be helpful in several ways. First, a taxonomy provides a natural organization for the sizeable libraries of computer viruses that antivirus organizations must maintain. Second, new viruses must be analyzed to tailor countermeasures, in a process that can be partly but not completely automated. If a new virus is related to one that was previously analyzed, the analysis may be simplified.

The most practical application of evolutionary information may be in increasing the efficiency of virus scanners. In a slightly simplified mathematical view, each of the 6000 computer virus species is represented as a byte string, typically 2000 bytes long. When antivirus programs “scan” for infected files (and antivirus programs do more than just this) they use a “signature” of about 20 bytes to stand in for each virus: the signature must always occur in the corresponding virus, and must never occur in legitimate computer code. If one signature can be used for several viruses, savings (in space more than time) can be achieved: the scanner requires only a minimum-sized set of signatures which together “cover” all the computer virus species.

In fact, the characters we will use to form a basis for computer virus phylogenies are such shared signatures. They are defined as, say, all strings of 20 bytes or more that occur in at least 2 viruses but in no legitimate programs. They can be found, using linear space and time, by straightforward application of suffix trees [7]. All viral and legitimate strings are concatenated together, separated by a special character, and a suffix tree is constructed. Its leaves represent all suffixes of the input string, and its internal nodes—viewed as paths from root partway to leaf—denote prefixes of suffixes, which is to say substrings of the input string. Depth-first search can be used to propagate, from leaves to root, the number of times each substring appears, or in fact the number of times it appears in viruses and (separately) in legitimate strings.

## 1.2. *Biological Application*

Beyond the computer virus realm for which it was conceived, the phyloDAG is also a plausible model for evolution of bacterial populations. Bacteria reproduces through simple cell division. A single cell divides into two daughter cells which each receive an exact copy of the parent cell’s genetic information (other than mutations that occur in transcription). However, there are at least three known methods whereby bacteria of different populations can exchange genetic information: transformation, transduction, and conjugation [13].

In *transformation*, a bacterium transports exogenous (outside the cell) DNA into the cell, where it can become incorporated into the bacterium's DNA. The exogenous DNA can come from another bacterium that has lysed (broken apart) and released its DNA into the medium. Only certain types of bacteria can do this and only under certain circumstances; some bacteria only bring in DNA that is quite similar to their own, while others will bring in any DNA, but will incorporate it only if it is suitably similar.

*Transduction* involves the transfer of genes from one bacterium to another via a bacteriophage (a virus that infects bacteria). Normally a virus infects a cell by binding to the cell and injecting its DNA. The virus then takes over the cell and forces it to make many more viruses. The infected cell then lyses (breaks apart), releasing the new virus particles. There are two mechanisms whereby viruses transmit genetic information. The first is generalized transduction: sometimes when the bacterial cell is producing new viruses, the viral package is filled with DNA from the host bacterium rather than the viral DNA. The process is random and so any piece of DNA can be packaged this way. When this "virus" is released, it can "infect" a cell by injecting its contents, but these contents are just bacterial DNA. This DNA will not kill the cell, and can become incorporated into the new host's DNA. The second mechanism is specialized transduction via lysogenic viruses. These viruses, upon infecting a bacterium, insert their DNA into the host DNA at a particular spot and coexist. When given the proper stimuli, the viral DNA is excised from the host DNA to carry out the normal infection cycle. Sometimes this excision is not done correctly, and pieces of the host DNA are excised as well. They are then packaged into the new viruses and transmitted to new hosts. Only genes near the attachment sites are transmitted this way, but the transmission is very efficient.

*Conjugation* involves the direct contact of two bacteria and the transmission of plasmids from one (donor) to the other (recipient). Plasmids are rings of DNA that are much smaller than the bacterial genome. They exist in the bacterial cell independently from the genome and are capable of replicating when the cell divides. Conjugative plasmids encode the proteins, etc., necessary for conjugation, thus engineering their own transmission. Conjugative plasmids can bring other genes with them into new cells, and can also allow the transmission of arbitrary plasmids. These plasmids can become incorporated into the cell DNA; for example, the genetic material of *E. Coli*'s F plasmid, which allows sexual conjugation, is incorporated into the host genome at a rate of  $10^{-5}$  per cell division. This is an important mechanism, because it is the primary way bacteria transfer drug resistance.

Because these mechanisms allow arbitrary exchange of genes from one population to another, bacterial evolution does not seem to follow the

“divergent evolution” implied by a tree: populations can evolve from multiple sources. Bacteria reproduce very rapidly and some region of their genome mutate frequently. Therefore, characters based on single-site mutations may not have a single archetype. However, for genes with sufficiently large mutation differences from any genes seen previously, it is reasonable to assume that as a rule there is unique evolution, and therefore a unique archetype.

### 1.3. Paper Organization and Results

We will show in Section 3.3 that the minimum phyloDAG problem is “hard”: in polynomial time, it cannot be solved exactly unless  $P = NP$ , nor can it be approximated to within better than a logarithmic factor unless  $NP \subseteq DTIME(n^{O(\log \log n)})$ . In fact, we know of no way to approximate minimum phyloDAG to within a logarithmic factor: Section 3.3 shows that various natural greedy strategies (including randomized ones) do not even approximate within a factor of  $cn$ .

Because of the difficulty of the phyloDAG problem, we consider two variants. In the first variant, we require that each species have just one ancestor, so that the phyloDAG is an arborescence (a tree with edges directed away from a root). If the arborescence’s vertices are labeled with the values of one character, the postulate that no character is “invented” twice corresponds to an assertion that there is at most one directed edge labeled  $0 \rightarrow 1$ . Thus the sequence of labels along any source-to-leaf path is described by the regular expression  $0^*1^*0^*$ , that is, zero or more 0s, followed by zero or more 1s, and finally zero or more 0s again. In Section 2 we define a **0–1–0 phylogeny** to be an arborescent phyloDAG’s underlying undirected tree. Species  $\mathcal{S}$  and characters  $\mathcal{C}$  may be consistent with zero, one, or multiple 0–1–0 phylogenies. We give two polynomial-time algorithms to randomly sample 0–1–0 phylogenies if any exist.

The first **atomic-set** algorithm (Section 2.1) computes a concise data structure that represents all 0–1–0 phylogenies for the input data and can be used to select a phylogeny uniformly at random in time  $O(nl)$ . When no solution exists the algorithm returns a **witness set**: a concise indication of *why* there can be no phylogenetic tree.

The second **minimum spanning tree** algorithm (Section 2.2) characterizes a 0–1–0 phylogeny of the input species set as a minimum spanning tree (MST) of a particular undirected edge-weighted graph. With it, 0–1–0 phylogenies can be constructed in deterministic time  $O(ln + n^2 \log n)$  or (with high probability) in randomized time  $O(ln)$ , and sampled uniformly at random in time  $O(ln + M(n))$ , where  $M(n)$  is the time needed to multiply two  $n \times n$  matrices. It does not produce a concise witness when there is no 0–1–0 phylogeny.

The second variant of phyloDAG is simply its undirected analogue. A **phylograph** for species  $\mathcal{S}$  and characters  $\mathcal{C}$  is an undirected graph with vertex set  $\mathcal{S}$ , with the property that the subgraph induced by the on-set of each character  $c \in \mathcal{C}$  is connected. The **Minimum Phylograph** problem is to find a phylograph with the minimum number of edges. Theorem 17 shows that it is hard to approximate Minimum Phylograph within a factor less than  $\frac{1}{4} \ln l$ , while Theorem 18 shows that approximating it within a factor of  $\ln l$  is easy.

The model of computation used in this paper is the uniform-cost random-access machine.

#### 1.4. Related Work

Previous work in phylogeny focused on constructing phylogenetic *trees*. However, the problem of modeling computer virus evolution is more suited to phylographs and phyloDAGs, in which undirected cycles may arise. As far as we know, ours is the first phylogenetic work that allows cycles.

There is substantial literature on character-based phylogenies where each subgraph induced by all species sharing a state for a character is required to be connected. This problem is called the **perfect phylogeny** problem, and is NP-complete for the “unrestricted” case (where putative species may be added) with general characters [3, 19]. For the unrestricted case with binary characters Gusfield gives an elegant  $O(nk)$  algorithm [12], and for the restricted case with general characters Goldberg *et al.* [11] give an algorithm analogous to the MST algorithm of Section 2.2. (To clarify the relationship between perfect phylogeny and our problem, note that the restricted perfect phylogeny problem with binary characters could be called the “0–1 phylogeny problem.”)

Our 0–1–0 phylogeny problem is similar to a restricted version of the **general character compatibility** problem of Benham *et al.* [2]. There a character  $c$  maps each species  $s$  to a *subset*  $c(s) \subseteq \{0, 1, 2\}$  rather than to a single value; the *leaves* of the tree are the species  $\mathcal{S}$ ; for each  $c$  and  $s$  a single value from  $c(s)$  is chosen as a label; and the goal is to find a rooted perfect phylogeny in which the sequence of labels along any root-to-leaf path is of the form  $0 \rightarrow 1 \rightarrow 2$ . The problem is NP-hard [2].

A preliminary version of this article appeared as [10].

## 2. COMPUTING A 0–1–0 PHYLOGENY

The case in which each species has only one ancestor is of special interest, and corresponds to cases in which the phyloDAG is an arborescence—a tree with all edges directed away from some root. There is a

straightforward  $n:1$  correspondence between arborescences and undirected trees: the undirected graph underlying an arborescence is a tree; and each of the  $n$  possible rootings of a tree is an arborescence.<sup>1</sup> Therefore we concentrate on undirected 0–1–0 phylogenies:

**DEFINITION 3.** An (undirected) **0–1–0 phylogeny**, or **phylogenetic tree**, is a tree  $T$  on species  $\mathcal{S}$  with characters  $\mathcal{E}$  such that each on-set  $S_c$  induces a subtree of  $T$ .

The requirement that each on-set induces a subtree corresponds to the requirement that each code fragment is invented only once. We allow the possibility that the code fragment may be dropped when one virus is used to create a new virus. Thus, we are interested in 0–1–0 phylogenies and not in perfect phylogenies.

If  $\mathbf{T}$  is a phyloDAG whose underlying graph is a tree  $T$ , then  $T$  is a 0–1–0 phylogeny as previously defined: as each on-set  $S_c$  was connected in  $\mathbf{T}$ , it is connected in  $T$ . Also, if  $T$  is a 0–1–0 phylogeny, any arborescence based on  $T$  is a phyloDAG: the archetype of any character  $c$  is the species in  $S_c$  closest to the root. In this section, we will show how to generate 0–1–0 phylogenies, and how to generate them uniformly at random. Given a uniformly random phylogenetic tree, choosing a root uniformly at random generates a uniformly random arborescent phyloDAG.

Because an arborescence can be rooted anywhere, a 0–1–0 phylogeny alone does not determine an evolutionary chronology, but it can be useful in combination with external information. For example, if the first species' identity is known, the rest of the evolutionary history follows.

### 2.1. The Atomic-Set Algorithm for Computing 0–1–0 Phylogenies

As described in the Introduction, our atomic-set algorithm produces a data structure, an AS-tree, which concisely represents all 0–1–0 phylogenies for species  $\mathcal{S}$  and characters  $\mathcal{E}$ , and can be used to generate an arbitrary solution or a solution chosen uniformly at random.

Generalizing the definition of the on-set of a character, define the on-set of a collection of characters to be the species having *all* those characters:  $S_C = \bigcap_{c \in C} S_c$ .

**DEFINITION 4.** Let  $\hat{C} \subseteq \mathcal{E}$  be a maximal (not necessarily *maximum*) set of characters for which  $|S_{\hat{C}}| \geq 2$ . Then  $A = S_{\hat{C}}$  is an **atomic set** with **defining characters**  $\hat{C}$ .

<sup>1</sup> There exist phyloDAGs whose underlying graphs are trees but which are not arborescences. An example, for species with characters  $(a)$ ,  $(ab)$ , and  $(b)$ , is  $(a) \rightarrow (ab) \leftarrow (b)$ . But because such phyloDAGs simply multiple ancestors for some species, they are not especially interesting.

LEMMA 5. For any atomic set  $A$  and character  $c$ , either  $S_c \supseteq A$  ( $c$  is a defining character), or  $|S_c \cap A| = 1$  ( $c$  is a **nondefining** character **owned** by the sole species  $s \in S_c \cap A$ ), or  $S_c \cap A = \emptyset$  ( $c$  is an **avoiding** character).

*Proof.* The only logical possibility missing is that  $|S_c \cap A| \geq 2$  but  $S_c \cap A \neq A$ , which would contradict the maximality of  $A$ 's set of defining characters. ■

An atomic set can be constructed in time  $O(kn)$ : start with  $\hat{C} = \emptyset$  (so  $S_{\hat{C}} = \mathcal{S}$ ), sweep through all characters  $c \in \mathcal{E}$  in turn, reject  $c$  if  $|S_{\hat{C}} \cap S_c| \leq 1$ , but otherwise add  $c$  to the defining set, so  $\hat{C} := \hat{C} \cup \{c\}$ . An  $O(l)$ -time implementation of this algorithm is described in the Appendix.

LEMMA 6. Suppose all species in  $\mathcal{S}$  are connected, i.e., the bipartite graph joining characters to species that have them is connected. Then if  $s_1, s_2 \in \mathcal{S}$  have no characters in common, no phylogeny contains the edge  $(s_1, s_2)$ .

*Proof.* Suppose a phylogenetic tree  $T$  contained  $(s_1, s_2)$ , and delete  $(s_1, s_2)$  to create a forest  $T'$ , consisting of two trees. For any character  $c$  and any  $s, s' \in S_c$ ,  $T$  has a path  $s, \dots, s'$  within  $S_c$ . The path does not include the edge  $(s_1, s_2)$ , because not both  $s_1$  and  $s_2$  can be in  $S_c$ , so  $T'$  contains the same path. Thus in  $T'$  there is a path from any species having character  $c$  to any other. Given the connectedness of the species-character graph, a series of such paths joins any species in  $\mathcal{S}$  to any other, contradicting the fact that  $T'$  is not a connected graph. ■

LEMMA 7. If  $A$  is an atomic set, then in any 0-1-0 phylogeny,  $A$ 's induced subgraph is a subtree.

*Proof.* In a 0-1-0 phylogenetic tree  $T$ , the on-set of any character  $c \in \mathcal{E}$  induces a connected subgraph, therefore a subtree.  $A$  is the intersection of the subtrees corresponding to  $A$ 's defining characters, and the intersection of subtrees is itself a subtree. ■

LEMMA 8. For any phylogeny  $T$  and atomic set  $A$ , if the subtree  $T_A$  is replaced by any other tree  $T'_A$  on the set  $A$ , the resultant overall tree  $T'$  is also a phylogeny.

*Proof.* For any character  $c$  and species  $s, s' \in S_c$ , consider the (unique) path  $s, \dots, s'$  in  $T$ . If  $S_c \cap A = \emptyset$ , the path never enters  $A$ , so it is unaffected (i.e., the identical path exists in  $T'$ ). If  $|S_c \cap A| = 1$ , the path touches at most one vertex in  $A$ , hence no edges within  $A$ , and is unaffected. Otherwise (by Lemma 5)  $S_c \supseteq A$ , and if the path through  $T$  included any subpaths through  $T_A$  (in fact there can be at most one), those sections could be replaced by subpaths through  $T'_A$  (and thus still within  $S_c$ ). So connectedness of all characters in  $T$  implies the same for  $T'$ , and  $T'$  is a phylogeny. ■

LEMMA 9. For any phylogeny  $T$  and atomic set  $A$ , if  $A$  is **collapsed**—replaced by a single species “ $a$ ” having all defining and nondefining characters of  $A$  (but not its avoiding characters), and the subtree  $T_A$  is contracted to the single species  $a$ , then the resultant overall tree  $T'$  is a phylogeny for  $\mathcal{S}' = (S \setminus A) \cup \{a\}$ .

*Proof.* Same as the previous one.

LEMMA 10. If  $(\mathcal{S}, \mathcal{E})$  has an atomic set  $A$ , with species  $s_1, s_2 \in A$  owning nondefining characters  $c_1, c_2$  respectively, and if  $S_{c_1} \cap S_{c_2} \neq \emptyset$ , then there is no 0–1–0 phylogeny for  $\mathcal{S}$ .

*Proof.* Suppose there is a phylogeny  $T$  for  $\mathcal{S}$ . Root  $T$  at any  $s_3 \in S_{c_1} \cap S_{c_2}$ , and let  $s_x$  be the lowest common ancestor of  $s_1$  and  $s_2$ . Then the path (all paths in a tree are unique) from  $s_1$  to  $s_2$  passes through  $s_x$ ; the path from  $s_3$  to  $s_1$  passes through  $s_x$  (because  $s_x$  is an ancestor of  $s_1$ ); and the path from  $s_3$  to  $s_2$  passes through  $s_x$  (because  $s_x$  is an ancestor of  $s_2$ ). By Lemma 7,  $A$  induces a subtree, so  $s_1, s_2 \in A$  implies that the  $s_1$ – $s_2$  path is contained in  $A$ , and in particular  $s_x \in A$ . Similarly,  $s_1, s_3 \in S_{c_1}$  implies  $s_x \in S_{c_1}$ , and  $s_2, s_3 \in S_{c_2}$  implies  $s_x \in S_{c_2}$ . Therefore  $s_x \in A \cap S_{c_1} \cap S_{c_2}$ . But  $c_1$  and  $c_2$  are nondefining characters with distinct owners, so  $A \cap S_{c_1} \cap S_{c_2} = \emptyset$ , a contradiction. ■

If the hypotheses of Lemma 10 are satisfied, we say that the atomic set  $A$ , characters  $c_1, c_2$ , and species  $s_1, s_2$  provide a **witness** attesting to the nonexistence of any 0–1–0 phylogenetic tree.

LEMMA 11. Let  $A$  be an atomic set, and suppose that no  $s_1, s_2, c_1, c_2$  satisfy the conditions of Lemma 10. As before, “collapse”  $A$  to the single species  $a$  having all defining and nondefining characters of  $A$ . If  $\mathcal{S}' = (\mathcal{S} \setminus A) \cup \{a\}$  has a phylogeny, so does  $\mathcal{S}$ .

*Proof.* Let  $T'$  be a phylogeny for  $\mathcal{S}'$ . Delete  $a$  and its incident edges, and replace them with the st  $A$  and any tree on  $A$ . Additionally, replace each edge  $(s, a)$  with a single edge as follows.

By Lemma 6,  $s$  and  $a$  must share some character(s), which (because  $a$  has them) must be defining or nondefining characters of  $A$ . If  $s$  and  $a$  share any *nondefining* characters, those characters must have a single owner  $s'$  (or else  $A$ , these characters, and their owners are a negative witness), in which case add the edge  $(s, s')$ . Otherwise,  $s$  and  $a$  only share defining characters of  $A$ , in which case add any edge  $(s, s')$  with  $s' \in A$ .

Replacement of each edge  $(s, a)$  with an edge  $(s, s')$ ,  $s' \in A$ , means that the tree components created by  $a$ 's deletion are all connected to the tree on  $A$ , creating a single tree  $T$ . Using arguments similar to those in Lemma 8, all characters induce connected components in  $T$  as they did in  $T'$ . ■

In fact, the constructive nature of the proof of Lemma 11 immediately suggests the **atomic-set algorithm**. Starting from  $\mathcal{S}_0 := \mathcal{S}$ , repeatedly, find an atomic set  $A_i$  and check for a witness as in Lemma 10. If one is found, terminate negatively. Otherwise, collapse  $A_i$  to a single new species  $a_i$ , and redefine the species set to be  $\mathcal{S}_i := (\mathcal{S}_{i-1} \setminus A_i) \cup \{a_i\}$ . Because each atomic set contains at least two species, this reduces the number of species, and needs to be performed at most  $n - 1$  times.

We construct the AS-tree during this contraction phase. The leaves of the AS-tree are the species in  $\mathcal{S}$ , and all elements of any set  $A_i$  have  $a_i$  as their parent. Equivalently, the final  $a_i$  is the root of the AS-tree, and each  $a_j$  has all species in  $A_j$  as children. This tree concisely represents all possible phylogenies.

Now, starting at the root of the AS-tree, we expand any node  $a_i$  whose parent is already expanded using the method suggested by the proof of Lemma 11: Replace  $a_i$  with  $A_i$  and form any tree  $T_i$  on  $A_i$ . For each old edge  $(s, a_i)$ , if  $s$  has a nondefining character  $c$  of  $A_i$ , add edge  $(s, \text{owner}_{A_i}(c))$ ; otherwise  $s$  must have only defining characters, in which case add any edge  $(s, s')$ ,  $s' \in A_i$ .

**THEOREM 12.** *The preceding algorithm produces a phylogeny for  $\mathcal{S}, \mathcal{E}$  if one exists, and otherwise produces a negative witness. The AS-tree that it constructs represents all possible 0–1–0 phylogenies. Thus, if the algorithm is implemented to choose trees  $T_i$  uniformly at random, and to choose  $s' \in A_i$  uniformly at random for defining-character edges  $(s, s')$ , then it produces a uniformly random undirected 0–1–0 phylogeny.*

*Proof.* The first assertion follows directly from the preceding sequence of lemmas. If we detect a negative witness, we correctly terminate negatively by Lemma 10 coupled with Lemma 9. Otherwise, by Lemmas 9 and 11, we can collapse the atomic set, solve the problem on the new set, and “expand” the collapsed set to a 0–1–0 phylogeny. Lemmas 8 and 9 and the proof of Lemma 11 show that all 0–1–0 phylogenies can be produced from the AS-tree obtained by the algorithm. the choices made in the expansion phase are independent and lead to different phylogenies. The uniform generation of phylogenies follows from this one-to-one correspondence between phylogenies, and choices in the algorithm. ■

Note that the order in which atomic sets are chosen by the algorithm affects the final AS-tree that is obtained, but that *any* AS-tree obtained by the algorithm is a concise representation of *all* possible 0–1–0 phylogenies. Furthermore, any AS tree obtained by the algorithm can be used to randomly sample the 0–1–0 phylogenies of the input.

The atomic-set algorithm produces an AS-tree in time  $O(nl)$ : in each of the  $O(n)$  collapsing iterations, we find an atomic set, check for a witness, and collapse the set, each such operation taking time  $O(l)$ . (See the Appendix.)

The expansion can be completed in time  $O(nl)$ . There are  $O(n)$  expansions. To expand node  $a_i$ , we can produce a random tree on the set  $A_i$  in time  $O(|A_i|)$ , because a labeled tree on  $r$  nodes can be selected uniformly at random in time  $O(r)$ . (See, for example, [16]). If we store pointers to owners of nondefining characters when constructing the AS-tree, we can connect this tree to its neighbors in time  $O(l)$ .

## 2.2. The Minimum Spanning Tree Algorithm

In this section we give a second algorithm for computing 0–1–0 phylogenies. It is very simple, and it is based on the observation that 0–1–0 phylogenies for species  $\mathcal{S}$  and characters  $\mathcal{E}$  correspond to minimum-weight spanning trees (MSTs) of a particular undirected edge-weighted graph  $G(\mathcal{S}, \mathcal{E})$ . (This observation was also used in [11] to obtain an algorithm finding restricted perfect phylogenies.)

The graph  $G(\mathcal{S}, \mathcal{E})$  is a complete graph on  $\mathcal{S}$ , with edge weights  $w(s_1, s_2) = k - |\{c \in \mathcal{E} \mid c(s_1) = c(s_2) = 1\}|$ . It can be constructed in  $O(ln)$  time.

**THEOREM 13.** *0–1–0 phylogenies for  $(\mathcal{S}, \mathcal{E})$  are spanning trees of  $G(\mathcal{S}, \mathcal{E})$  with weight  $nk - l$ . Furthermore,  $G(\mathcal{S}, \mathcal{E})$  has no spanning trees of smaller weight.*

*Proof.* Every spanning tree of  $G(S)$  has weight at least  $nk - l$ , because the contribution of each character  $c$  to the total weight is at least  $(n - 1) - (|S_c| - 1)$ . Spanning trees of  $G(S)$  with weight  $nk - l$  correspond to trees in which each on-set  $S_c$  is connected (see [11]).

Because of this correspondence, phylogenies can be constructed (or randomly sampled) by established algorithms for constructing (or randomly sampling) MSTs. Prim's algorithm [17, 9] constructs an MST of  $G$  in  $O(m \log m)$  time, where  $m$  is the number of edges in  $G$ , and  $m \binom{n}{2}$  for  $G = G(\mathcal{S}, \mathcal{E})$ . If a faster algorithm is required, Karger, Klein, and Tarjan's randomized algorithm constructs an MST, with high probability, in  $O(m)$  time [14]. (Their model of computation is a unit-cost random-access machine with the restriction that the only operation allowed on edge weights are binary comparisons. See also the other algorithms discussed in [14].)

Given an unweighted  $n$ -vertex graph, an algorithm of Colbourn, Myrvold, and Neufeld [5] selects a spanning tree uniformly at random in

$O(M(n))$  time.<sup>2</sup> (Here  $M(n) = O(n^{2.376})$  is the time needed to multiply two  $n \times n$  matrices [6].) Colbourn and Jerrum [4] note that the algorithm can be used to select an MST of a weighted graph  $G$  uniformly at random in  $O(M(n))$  time: construct a random spanning tree on each connected component of the subgraph of  $G$  induced by the edges of minimum weight, put the spanning trees' edges into the final solution, contract the spanning trees, and repeat.

Compared with the atomic-set algorithm, the MST approach has the advantage of using an unusually widely understood and simple paradigm, a benefit echoed in the availability and efficiency of computer programs. However, it does not supply a structural representation of all possible phylogenies, nor a concise witness when no phylogeny exists.

### 3. PHYLOGRAPHS AND PHYLODAGS

Having considered the problem of constructing phylogenetic trees, we now turn to phylogenies that are not trees. In particular, we consider the phylograph and phyloDAG problems that were defined in the Introduction. In Section 3.1 we prove that it is hard to approximate the optimal phylograph within better than a logarithmic factor, and in Section 3.2 that the natural greedy algorithm gives an approximation within such a factor. In Section 3.3 we show both that it is hard to approximate the optimal phyloDAG within better than a logarithmic factor, and that in this case the natural greedy algorithm can perform very badly, even on average.

#### 3.1. *Hardness of Approximation of Phylograph*

Hardness results for Minimum Phylograph follow from those known for **Minimum Set Cover** and problems equivalent to it in terms of approximation ratio, notably **Minimum Dominating Set**.

**DEFINITION 14.** The **neighborhood** of a vertex  $v$  of a graph  $G = (V, E)$  is the set  $N(v) = \{v\} \cup \{w : (v, w) \in E\}$ . A **dominating set** of  $G$  is a set of vertices  $D \subseteq V$  whose neighborhoods cover the graph:  $\bigcup_{d \in D} N(d) = V$ .

It is well known and it is easily proved that the natural greedy algorithm for Minimum Dominating Set (or the related problems) is a  $\ln n$  approximation algorithm: for a graph  $G = (V, E)$ , the dominating set produced by

<sup>2</sup> Another randomized algorithm, due to Wilson [20], as an expected running time equal to the *mean hitting time* of the graph; this is often smaller than  $M(n)$ , but can be larger.

the greedy algorithm is at most  $\ln|V|$  times larger than the minimum dominating set. In [8], Feige shows that this is a threshold:

**THEOREM 15 (Feige).** *Let  $c$  be a constant in the range  $0 < c < 1$ . Unless  $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$ , there is no polynomial-time algorithm that takes as input a graph  $G$  and outputs a dominating set  $D$  of  $G$  such that  $|D|$  is within a factor of  $c \ln|V|$  of the minimum possible value.*

Feige's is the latest in (and contains a good review of) a sequence of works on this problem. Another which is relevant here, because of its weaker hypothesis, is that of Bellare *et al.* [1]:

**THEOREM 16 (Bellare, Goldwasser, Lund, and Russell).** *Unless  $\text{P} = \text{NP}$ , there is no polynomial-time algorithm that approximates Minimum Dominating Set to within any constant factor.*

From these results we can show that Minimum Phylograph cannot be approximated to within any constant factor unless  $\text{P} = \text{NP}$ , and cannot be approximated to better than a logarithmic factor unless  $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$ .

**THEOREM 17.** *Unless  $\text{P} = \text{NP}$ , for any constant  $c > 0$ , there is no polynomial-time algorithm that takes as input species  $\mathcal{S}$  and characters  $\mathcal{E}$  and outputs a phylograph  $G = (\mathcal{S}, E)$  such that  $|E|$  is within a factor of  $c$  of the minimum possible value.*

*Similarly, for any  $c < 1/4$ , unless  $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$ , there is no polynomial-time algorithm approximating Minimum Phylograph within a factor  $c \ln l$  (where  $l$  is the input length defined earlier).*

*Proof.* We use an approximation-preserving reduction from Minimum Dominating Set to Minimum Phylograph. Given an input  $G = (V, E_G)$  with  $|V| = \nu$ , construct an instance  $P$  to Minimum Phylograph as follows: The species set is  $\mathcal{S} = V \cup X$  where  $X$  is a set of  $\nu^3$  "auxiliary vertices." For each pair of vertices  $\{v_1, v_2\} \in V^{(2)}$ , define a character with on-set  $\{v_1, v_2\}$ . Thus any phylograph for  $P$  contains each edge in the complete graph on  $V$ . In addition, for each pair of vertices  $(v, x) \in V \times X$  we define a character with on-set  $\{x\} \cup N(v)$ .

If  $P_0 = \{\mathcal{S}, E_0\}$  is an optimal phylograph for  $P$ , and  $D_0$  is a minimum dominating set for  $G$ , then  $|E_0| = \binom{\nu}{2} + |X||D_0|$ . To see this, observe that the complete graph on  $V$  added to  $X \times D_0$  is a phylograph for  $P$ , so  $|E_0| \leq \binom{\nu}{2} + |X||D_0|$ . On the other hand, every phylograph for  $P$  has at least  $\binom{\nu}{2}$  edges connecting species in  $V$  and has at least  $|D_0|$  edges adjacent to each  $x \in X$ .

Suppose we had an algorithm  $A$  that could produce a phylograph  $(\mathcal{S}, E_A)$  for  $P$  with  $|E_A| \leq c \ln(l)|E_0|$  edges. By the construction of  $P$ ,

some vertex  $x \in X$  is connected to a dominating set  $D$  for  $G$  with  $|D| \leq |E_A|/|X| \leq c \ln(l)|E_0|/|X|$  edges. Because  $|E_0| = \binom{\nu}{2} + |X||D_0|$ , we have

$$|D| \leq c \ln(l) \frac{\binom{\nu}{2} + |X||D_0|}{|X|}.$$

Thus (because  $|X| = \nu^3$  and  $|D_0| \geq 1$ ),  $|D| \leq c(1 + o(1))\ln(l)|D_0|$ . Now note that  $l = \nu(\nu - 1) + 2\nu|X| + 2|E_G||X| = O(\nu^5)$ . Thus,  $|D| \leq 5c(1 + o(1))\ln(\nu)|D_0|$ , which is contrary to Theorem 15 if  $c < \frac{1}{5}$ , unless  $\text{NP} \subseteq \text{DTIME}(n^{O(\log \log n)})$ . Using  $|X| = \nu^{2+\epsilon}$  instead of  $\nu^3$  gives the constant  $\frac{1}{4}$ .

Similarly, a constant-approximation algorithm is forbidden by Theorem 16, unless  $\text{P} = \text{NP}$ . ■

### 3.2. Greedy Algorithm for Phylograph

There is a natural greedy algorithm for the Minimum Phylograph problem. In a phylograph, every character's induced subgraph consists of a single connected component, so the greedy algorithm "grows" a solution by iteratively adding an edge that maximally reduces the number of connected components.

The same notation needed to define the algorithm more precisely can be used in the proof of its quality. Given species  $\mathcal{S}$  and characters  $\mathcal{C}$ , and a set of edges  $E \subseteq \mathcal{S}^{(2)}$  define the "cost" of  $E$  to be

$$f(E) = \sum_{c \in \mathcal{C}} \text{components}(S_c) - |\mathcal{C}|,$$

where  $\text{components}(S_c)$  denotes the number of connected components in the subgraph of  $(\mathcal{S}, E)$  induced by the on-set of  $c$ . Thus  $f(\emptyset) = \sum_{c \in \mathcal{C}} |S_c| - |\mathcal{C}| = l - |\mathcal{C}|$ , and if  $E$  is a phylograph,  $f(E) = \sum_{c \in \mathcal{C}} 1 - |\mathcal{C}| = 0$ .

For any edge set  $E$  and any edge  $e$ , let  $\Delta_E(e) = f(E) - f(E \cup \{e\})$  be the amount by which  $e$  decreases the cost  $f$ . The greedy algorithm begins with each species an isolated vertex, and iteratively adds the edge which maximally decreases the cost, until the cost is 0. In pseudocode:

```

Let  $i := 0$  and  $E_G(0) := \emptyset$ 
While  $f(E_G(i)) > 0$  do
begin
  Let  $i := i + 1$ 
  Let  $e$  be an edge maximizing  $\Delta_{E_G(i-1)}(e)$ 
  Let  $E_G(i) := E_G(i-1) \cup \{e\}$ 
end
Return the set  $E_G = E_G(i)$ 

```

**THEOREM 18.** *Suppose that for species  $S$  and characters  $C$ , of total input length  $l$ , the minimum phylograph  $\{e(1), \dots, e(r)\}$  has cardinality  $r$ . Then the greedy algorithm produces a phylograph  $E_G$  of size  $|E_G| \leq r \ln(l - |C|)$ .*

*Proof.* If we have any partial solution, adding in all  $r$  edges of a minimum phylograph will certainly yield a phylograph. Because  $r$  more edges are enough to complete the job, some edge (one of these, even) must take care of at least  $1/r$ th of the cost. If the initial cost was  $f(\emptyset)$ , and the greedy algorithm reduces it by  $1 - 1/r$  at each step, after  $r \ln f(\emptyset)$  steps the cost must be reduced below 1, and the algorithm must have terminated.<sup>3</sup>

More formally, for any edge set  $E(0)$  define a series of sets  $E(0) \subseteq \dots \subseteq E(r)$ , where  $E(i) = E(0) \cup \{e(1), \dots, e(i)\}$  and the edges  $e(i)$  are those of the minimum phylograph. Note that  $E(r)$  is a phylograph, because it contains the minimum phylograph. Because components (with respect to any character) only become more connected as  $i$  increases, for any  $e$ , if  $i \leq j$  then  $\Delta_{E(i)}(e) \geq \Delta_{E(j)}(e)$ . Thus for any starting set  $E_0$ ,

$$\begin{aligned} r \cdot \max_{e \in \mathcal{A}^{(2)}} \Delta_{E(0)}(e) &\geq \sum_{i=1}^r \Delta_{E(0)}(e(i)) \\ &\geq \sum_{i=1}^r \Delta_{E(i-1)}(e(i)) \\ &= \sum_{i=1}^r [f(E(i-1)) - f(E(i))] \\ &= f(E(0)) - f(E(r)) \\ &= f(E(0)). \end{aligned}$$

Comparing the first and last quantities, we conclude that there always exists an edge  $e$  for which  $\Delta_{E(0)}(e) \geq f(E_0)/r$ .

Therefore the greedy algorithm reduces the cost by a factor  $1 - 1/r$  at each step. Because the initial cost is  $l - |C|$ , the cost after  $r \ln(l - |C|)$  steps of the greedy algorithm is at most  $(1 - 1/r)^{r \ln(l - |C|)}(l - |C|) \leq 1$ . The greedy algorithm therefore terminates within  $r \ln(l - |C|)$  steps, producing a phylograph of the same size.

This complements the result of Theorem 17: Minimum Phylograph is apparently hard to approximate to better than a factor of  $\frac{1}{4} \ln l$ , but easy to approximate to a factor  $\ln(l - |C|) \leq \ln l$ . It would be of some interest to

<sup>3</sup> The same approach will not work for phyloDAG. Because directed cycles are forbidden, chosen edges constrain the addition of future ones, and even if there was a solution of size  $r$  initially, there may not be once some edges are chosen suboptimally.

derive better bounds on the constant  $c$ ,  $\frac{1}{4} < c \leq 1$ , for which  $(c \ln l)$ -approximability is possible.

### 3.3. PhyloDAGs

We begin by observing that a phyloDAG cannot always be obtained by directing the edges of a phylograph. Consider four species with  $s_1$  defined by characters  $(b, c, d)$ ,  $s_2$  by  $(a, c, d)$ ,  $s_3$  by  $(a, b, d)$ , and  $s_4$  by  $(a, b, c)$ . The cycle  $s_1, s_2, s_3, s_4, s_1$  is a 4-edge phylograph, but there is no way to direct the edges of the cycle to obtain a phyloDAG: any acyclic orientation will create two archetypes for some character's on-set.

We now prove the following theorem, which is analogous to Theorem 17.

**THEOREM 19.** *Unless  $P = NP$ , for any constant  $c > 0$ , there is no polynomial-time algorithm that takes as input species  $\mathcal{S}$  and characters  $\mathcal{E}$  and outputs a phyloDAG  $G = (\mathcal{S}, E)$  such that  $|E|$  is within a factor of  $c$  of the minimum possible value.*

*Similarly, for any  $c < \frac{1}{4}$ , unless  $NP \subseteq DTIME(n^{O(\log \log n)})$ , there is no polynomial-time algorithm approximating within a factor  $c \ln l$ .*

*Proof.* The proof uses the same reduction as the proof of Theorem 17. Let  $E_0$  be the edge set in an optimal phyloDAG for  $P$ . We must show that  $|E_0| = \binom{v}{2} + |X||D_0|$ . The direction that differs from the proof of Theorem 17 is showing that given a dominating set  $D_0$ , we can construct a phyloDAG of size  $\binom{v}{2} + |X||D_0|$ . To do so, first construct a phylograph (as in the proof of Theorem 17). Then direct edges having both end points in  $V$  according to a total order on the vertices in  $V$ , and direct all remaining edges from vertices in  $V$  toward vertices in  $X$ . The resulting digraph has no directed cycles and each character has a unique archetype. Therefore, it is a phyloDAG. The rest of the proof is identical to that of Theorem 17. ■

As already noted, the natural greedy algorithm does not work well for phyloDAGs: the phyloDAG problem seems to be more difficult because the prohibition of cycles means that it is possible for the greedy algorithm to add a “bad” edge which prevents other “good” edges from being added later. In the remainder of this section, we give an example of a species set for which various natural greedy approaches for constructing a phyloDAG lead to an  $\Omega(n)$  ratio between the size (number of edges) of the constructed phyloDAG and the size of the optimal phyloDAG. A randomized strategy has an  $\Omega(n)$  expected ratio and has a ratio of  $\Omega(n/\log n)$  with high probability.

We construct a species set as follows. There are  $n$  species  $s_1, \dots, s_n$ , and two distinguished species  $s'$  and  $s''$ . Now we add

- $2n$  characters shared by  $s'$  and  $s''$ ;
- $2$  characters shared by  $s'$  and  $s_i$ , for  $i = 1, \dots, n$
- $1$  character shared by  $s''$ ,  $s_i$ , and  $s_j$ , for  $1 \leq i, j \leq n$ ,  $i \neq j$ .

Duplicating characters forces the order in which a greedy algorithm connects species. We hide this duplication from an algorithm that checks for it by adding a set  $S_d$  of dummy species, where  $|S_d| = \lceil \log(4n) \rceil$ . There are  $2^{\lceil \log(4n) \rceil} \geq 4n$  distinct subsets of  $S_d$ . We add one such subset to each of the  $4n$  nonunique characters.<sup>4</sup> An optimal solution has  $O(n)$  edges, consisting of an edge from  $s'$  to  $s''$ , edges from  $s'$  and  $s''$  to each of the  $s_i$ , and edges from  $s'$  to each species in  $S_d$ .

A phyloDAG has exactly one archetype for each character. A greedy algorithm begins with each species an isolated node, thus an archetype for each character it contains. A natural edge to add in a greedy fashion is one that maximally reduces the number of archetypes (over all characters). Of course, we may not introduce directed cycles.

There may be times where we can choose the direction of the edge to be introduced (for example at the first iteration) and we show that the algorithm performs badly for any of the following strategies:

- The direction is chosen arbitrarily.
- The direction is chosen uniformly at random. (The expected performance of the algorithm is bad for this example, and the example can be modified so that the bad performance occurs with high probability.)
- The edge is directed out from the node with the larger number of characters. (This is a natural way of breaking ties, because we expect ancestral nodes to have many characters.)

A greedy algorithm starts by putting an edge between  $s'$  and  $s''$ , and an edge between  $s'$  (or possibly  $s''$ ) and each species in  $S_d$ . Then it adds edges between  $s'$  and the  $s_i$ . If direction are chosen arbitrarily we may assume that these edges are from  $s'$  to  $s''$ , and from each of the  $s_i$  to  $s'$ . Hence it is now impossible to add edges from  $s''$  to any of the  $s_i$ , because they would create directed cycles. This means that in order to prevent there being two archetypes for a character shared by  $s''$ ,  $s_i$ , and  $s_j$ , species  $s_i$  must be connected to  $s_j$  by an edge. This results in  $\binom{n}{2}$  edges.

Now consider the variant where the direction of an edge is chosen uniformly at random whenever it is equally good to direct it either way.

<sup>4</sup> An algorithm may also check for *domination*, where  $s_d$  contains a subset of the characters contained by  $s$ . We can remove the dominated species  $s_d$  from the instance and later direct an edge from  $s$  to  $s_d$  in the phylogeny for the reduced set. To avoid this situation here, we add a character  $\{s_d(i), s_d(i+1)\}$  for  $i = 1, \dots, [4n]$ , which chains the dummies together. This does not change the asymptotic size of the optimal solution.

With high probability (i.e., with complement probability that is exponentially small in  $n$ ), there will be at least  $n/4$  edges directed from the  $s_i$ s to  $s'$ . If the edge between  $s'$  and  $s''$  is directed the wrong way (i.e., from  $s'$  to  $s''$ ) then these  $s_i$  nodes will have to be connected in a clique, resulting in a quadratic number of edges. If we now consider a species set consisting of  $\alpha \log(n)$  copies of the species set as described (for a positive constant  $\alpha$ ), we see that the optimal solution has  $\Theta(n \log n)$  nodes and edges, and with probability at least  $1 - n^{-\alpha}$ , at least one of those copies will have the edge between  $s'$  and  $s''$  directed the wrong way, resulting in  $\Theta(n^2)$  edges.

If edges are directed away from nodes with higher numbers of characters, then the algorithm can be forced to take the “wrong” direction for the edges by adding dummy characters at the nodes from which we want the edges to be directed.

## ACKNOWLEDGMENTS

We are grateful to Phil MacKenzie, Tom Martin, Baruch Schieber, Madhu Sudan, Luca Trevisan, and David Wilson for useful discussions, and to Luca for directing us to [8]. We also thank the reviewers and referees for their helpful comments.

## 4. APPENDIX

### 4.1. An $O(l)$ -time Algorithm to Compute an Atomic Set

We presume that each species is described by a sorted list of the characters it contains. From this, construct a description of each character, as a sorted list of the species containing it. This can be done in linear time: loop through species  $i$ ; loop through characters  $j$  on  $i$ ; add species  $i$  to character  $j$ 's list.

Now the basic algorithm is:

Let  $A_0 := \mathcal{S}$  (0th atomic set contains all species)

Let  $D := \emptyset$  (set of defining characters is initially empty)

Loop through characters  $i$ , and consider the list  $S_i$  of species having character  $i$ :

- (1) size :=  $|A_{i-1} \cap S_i|$
- (2) If size  $< 2$  then  $A_i := A_{i-1}$
- (3) If size  $\geq 2$  then  $A_i := A_{i-1} \cap S_i$ , and  $D := D \cup \{i\}$
- (4) next  $i$

We now show how to compute the intersection size (step 1) and the intersection itself (step 3) in linear time. This implementation gives a

linear-time algorithm overall. Let all the set  $A$  be represented doubly, as both a sorted linked list, and as a binary array of length  $k$  (with 1s for species present in  $A$ , 0s for species absent from  $A$ ).

Computing the size of  $A \cap S_i$  can be done in time  $|S_i|$ : Over species  $s \in S_i$ , sum up the binary array elements  $A[s]$ . Thus all iterations of step 1, together, take time of order  $\sum_i |S_i|$ .

Computing  $A' := A \cap S_i$  can be done in time  $|A| + |S_i|$ : The ordered list for  $A'$  is constructed by stepping through the ordered lists for  $A$  and  $S_i$  in synchrony, advancing in the list with the smaller current value, and augmenting the list for  $A'$  when the lists for  $A$  and  $S_i$  have the same current value. The binary array for  $A'$  is formed by modifying that of  $A$ , which is no longer needed for any other purpose; the list for  $A$  is used to set all 1s in the array back to 0, and then the list for  $A'$  is used to set 1s appropriately.

Thus over values  $i$  where step 3 is executed, the total time consumed is of order

$$\begin{aligned} \sum_i |A_{i-1}| + \sum_i |S_i| &\leq |A_0| + \sum_i |A_i| + \sum_i |S_i| \\ &\leq |\mathcal{S}| + 2 \sum_i |S_i|, \end{aligned}$$

because the execution of step 3 implies that  $|A_i| \leq |S_i|$ . Thus the total time consumed by all steps of the algorithm is at most of order  $|\mathcal{S}| + \sum_i |S_i| = O(I)$ .

## REFERENCES

1. M. Bellare, S. Goldwasser, C. Lund, and A. Russell, Efficient probabilistically checkable proofs and applications to approximation, in "Proceedings of the Twentyfifth Annual ACM Symposium on the Theory of Computing, 1993," pp. 294–304.
2. C. Benham, S. Kannan, M. Paterson, and T. Warnow, Hen's teeth and whale's feet: Generalized characters and their compatibility, *J. Math. Biol.* **2**(4) (1995), 515–525.
3. H. Bodlaender, M. Fellows, and T. Warnow, Two strikes against perfect phylogeny, in "Proceedings of the Nineteenth International Colloquium on Automata, Languages, and Programming, 1992," Lecture Notes in Computer Science, Springer-Verlag, Berlin/New York, pp. 273–283.
4. C. Colbourn and M. Jerrum, 1995, personal communication.
5. C. Colbourn, W. Myrvold, and E. Neufeld, Two algorithms for unranking arborescences, *J. Algorithms*, to appear.
6. D. Coppersmith and S. Winograd, Matrix multiplication via arithmetic progressions, *J. Symbolic Comput.* **9** (1990), 251–280.
7. M. Corchemore and W. Rytter, "Text Algorithms," Oxford Univ. Press, London, 1994.
8. U. Feige, A threshold of  $\ln n$  for approximating set cover, in "Proceedings of the Twentyeighth Annual ACM Symposium on the Theory of Computing, 1996," pp. 286–293.

9. A. Gibbons, "Algorithmic Graph Theory," Cambridge Univ. Press, Cambridge, U.K., 1985.
10. L. Goldberg, P. Goldberg, C. Phillips, and G. Sorkin, Constructing computer virus phylogenies, in "Proceedings of Combinatorial Pattern matching, 1996," pp. 253–270.
11. L. Goldberg, P. Goldberg, C. Phillips, E. Sweedyk, and T. Warnow, Computing the phylogenetic number to find good evolutionary trees, *Discrete Appl. Math.* **71**, 1–3 (1996), 111–136.
12. D. Gusfield, Efficient algorithms for inferring evolutionary trees, *Networks* **21** (1991), 12–28.
13. W. Jiklik, H. Willett, D. Amos, and C. Wilfert, Eds. "Zinsser Microbiology," Appleton & Lange, East Norwalk, CT, 20th ed., 1992.
14. D. Karger, P. Klein, and R. Tarjan, A randomized linear-time algorithm to find minimum spanning trees, *J. Assoc. Comput. Mach.* **42**(2) (1995).
15. J. Kephart and W. Arnold, Automatic extraction of computer virus signatures, in "Proceedings of the Fourth Virus Bulletin International Conference, 1994" (R. Ford, Ed.), Virus Bulletin Ltd., pp. 179–194.
16. A. Nijenhuis and H. Wilf, "Combinatorial Algorithms for Computers and Calculators," Academic Press, 2nd ed., 1978.
17. R. Prim, Shortest connection networks and some generalizations, *Bell Syst. Tech. J.* **36** (1957), 1389–1401.
18. G. Sorkin, Grouping related computer viruses into families, in "Proceedings of the IBM Security ITS," October 1994.
19. M. Steel, The complexity of reconstructing trees from qualitative characters and subtrees, *J. Classification* **9** (1992), 91–116.
20. D. Wilson, Generating random spanning trees more quickly than the cover time, submitted for publication.