

# A Reduction Algorithm for Copying and Hybridization Networks

Magnus Bordewich<sup>1</sup>, Simone Linz<sup>2,3</sup>, Katherine St. John<sup>4</sup>,  
and Charles Semple<sup>2</sup>

<sup>1</sup> Department of Computer Science, Durham University, Durham, DH 1 1TA, United Kingdom

<sup>2</sup> Bioinformatics Research Centre, Department of Mathematics and Statistics, University of Canterbury, Christchurch, New Zealand

<sup>3</sup> Department of Bioinformatics, Heinrich Heine University, Düsseldorf, Germany

<sup>4</sup> Department of Mathematics and Computer Science, Lehman College, City University of New York, USA

**Correspondence:** Simone Linz, Bioinformatics Research Centre, Department of Mathematics and Statistics, University of Canterbury, Private Bag 4800, Christchurch, New Zealand. Tel: +64 3 364 2639, Email: linz@cs.uni-duesseldorf.de

**Running head:** A Reduction Algorithm for Hybridization

**Key words:** hybridization networks, reticulate evolution, green forest

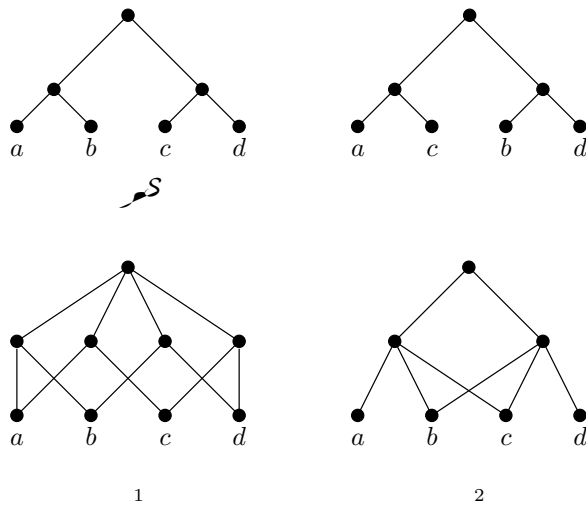
## 1 Abstract

Hybridization is an important evolutionary process for many groups of species. Thus conflicting signals in a dataset may not be the result of sequencing or coding errors, but due to the fact that hybridization has played a significant role in the evolution of the species.



When the initial collection consists of rooted binary phylogenetic trees (Bordewich and Sempeck, 2005). Consequently, as a result of this computation, currently most research considers the tree problem. There are no known algorithms for approximating this latter problem. However, some of these algorithms are either algorithms solving a restricted version of the problem (e.g. Heled and Lerner, 2005; Huson et al., 2005; Nishihara et al., 2005) or polynomial time heuristics with no guarantee of the correctness of their solution (e.g. Nishihara et al., 2005).

In this paper, we describe new and recently implemented algorithm for solving the tree problem with no restrictions based on three reductions that preserve the amount of hybridization. All of these reductions use the use of similarities between the trees. It has recently been shown that two of the reductions are enough to guarantee that the algorithm is a parameter tractable where the parameter is the smallest number of hybridizations to appear in the initial trees (Bordewich and Sempeck, 2005). This means that the algorithm runs efficiently when this smallest number is bounded. The remaining reduction is



**Figure 1.** Two rooted binary phylogenetic trees and two hybridization networks  $\mathcal{S}_1$  and  $\mathcal{S}_2$  which appear in both trees.

combination of the reduced parameter result described in Bordewich and Semple, whose proof of correctness is given by Proposition 4 of that paper and the cluster reduction described in Bordewich et al., whose proof of correctness is given by Theorem 1 in that paper. For simplicity in this paper we only describe the algorithms. For the reader interested in the finer details we refer the reader to the original papers.

### 3 Reduction Algorithm for Hybridization

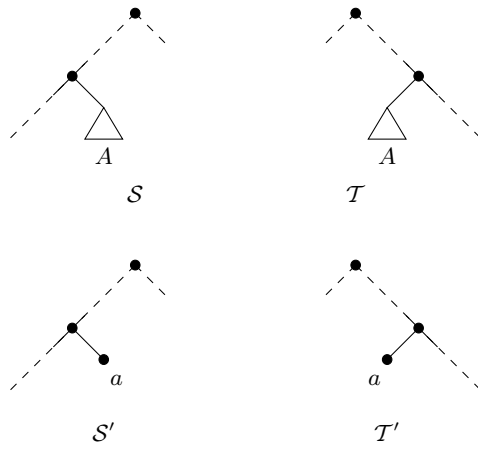
begin with for description of the tree problem. A **rooted binary phylogenetic  $X$ -tree** is a rooted tree that has leaf set  $X$  and whose root has degree two and all other interior vertices have degree three. A **cluster** of  $X$  is a subset of  $X$  that contains precisely the elements that are descendants of some vertex of  $T$ .

A **rooted acyclic digraph** is a digraph with no directed cycles. Each such digraph has a distinguished vertex  $\rho$  whose in-degree is zero and has the property that there is a directed path from  $\rho$  to every other vertex. For a vertex  $v$  in a digraph  $D$  denote the **in-degree** of  $v$  the number of edges directed into  $v$  by  $d^-(v)$  and the **out-degree** of  $v$  the number of edges directed out of  $v$  by  $d^+(v)$ . A **hybridization network** on  $X$  is a rooted acyclic digraph with root  $\rho$  in which

- i)  $X$  is the set of vertices of out-degree zero
- ii)  $d^+(\rho) = 2$  and



such network, the size of the resulting green forest for  $S$  and here the size of forest is the number of trees in the forest. On the other hand, if we are given a green forest for  $S$  and then one can reverse this process to construct hybridization network that explains  $S$  and



**Figure 2.** Two rooted phylogenetic trees  $S$  and  $T$  reduced under the subtree reduction rule. The triangle





Furthermore, the correctness of the cluster reduction rule follows from Proposition 4 of Bordewich and Sempe, 2004.

- ii) Bordewich and Sempe, 2004, showed that the subtree reduction and cluster reductions, by themselves, are enough to characterize the problem and give a polynomial algorithm for **Hybridization Number**. The cluster reduction provides a new, very useful tool for reducing the problem into a number of smaller problems that is required is that the subtrees have identical leaf sets, the topologies of the two subtrees can be completely different.
- iii) Without going into details, the cluster reduction has a similar flavor to the Decomposition Theorem in Huson **et al.**, 2005. This theorem describes a one-to-one correspondence between the overlapping cycles of an unrooted network  $\mathcal{N}$ , the connected components of the incompatibility graph of the splits generated by  $\mathcal{N}$ , and the netted components of the splits graph of the splits generated by  $\mathcal{N}$ . However, since this theorem yields an algorithm for minimizing the number of hybridization vertices amongst restricted classes of networks, it is important to note that it does not give a general strategy for minimizing this number amongst hybridization networks. In fact, there is no guarantee that such a reduction leads to an optimal solution. In contrast, Bonini **et al.**, 2007, showed that such a strategy, in particular, the cluster reduction, works for rooted trees. It is an interesting open

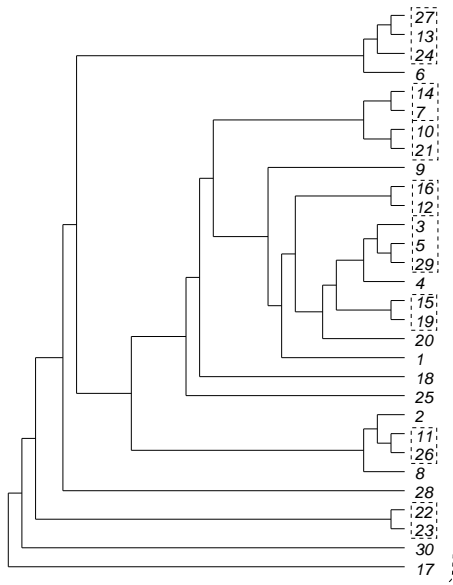


**Table 2.** Results for the **Poaceae** dataset.

pairwise combination	# taxa	hybridization number	run time <sup>a</sup>
<i>nd</i> - <i>yB</i>	40	14	11 h
<i>nd</i> - <i>cL</i>	36	13	11.8 h
<i>nd</i> - <i>oC</i>	34	12	26.3 h
<i>nd</i> - <i>xy</i>	19	9	320 s
<i>nd</i> - <i>⚡</i>	46	at least 15	2 d
<i>yB</i> - <i>cL</i>	21	4	1 s
<i>yB</i> - <i>oC</i>	21	7	180 s
<i>yB</i> - <i>xy</i>	14	3	1 s
<i>yB</i> - <i>⚡</i>	30	8	19 s
<i>cL</i> - <i>oC</i>	26	13	29.5 h
<i>cL</i> - <i>xy</i>	12	7	230 s
<i>cL</i> - <i>⚡</i>	29	at least 9	2 d
<i>oC</i> - <i>xy</i>	10	1	1 s
<i>oC</i> - <i>⚡</i>	31	at least 10	2 d
<i>xy</i> - <i>⚡</i>	15	8	620 s

<sup>a</sup>run time on a 2000 MHz CPU, 2 GB RAM machine measured in seconds (s), hours (h), and days (d), respectively

post sequence phytochrome B **phyB** and the nuclear sequence of the inter-nuclear transcribed spacer of ribosomal DNA **ITS** which have been overlapping this set of present day species see the round indicated by the grey circle



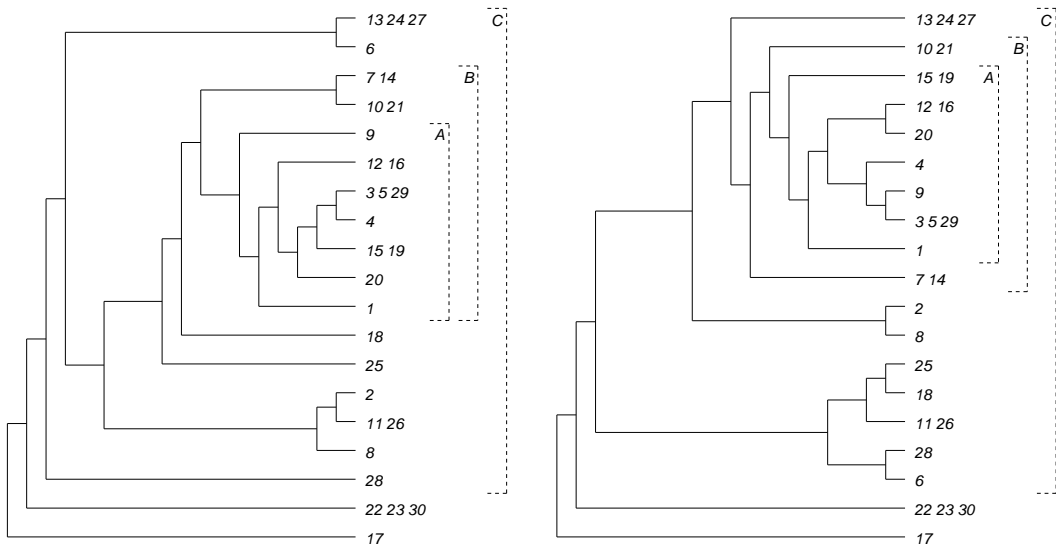
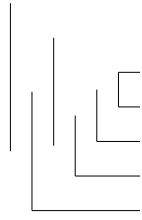


Figure 6.



subtrees, clades, or clusters which is likely for any biologically plausible tree. The algorithm performs recombination and the hybridization number are calculated as found in Robinson et al.

Note that **HybridNumber** calculates a lower bound for the number of hybridization events to explain the differences between two phylogenetic gene trees assuming that hybridization is the only cause of incongruence between the two trees. It is possible that the real number of hybridization events that happened during the evolution of the collection of present day species under consideration is underestimated. Indeed, it is possible that some hybridization events are never recognized. Nevertheless, the algorithm provides the important results.





Nehel Ruths and Ling LS. RIATA HGT: a fast and accurate heuristic for reconstructing horizontal gene transfer. In **Proceedings of the Eleventh International Computing and Combinatorics Conference (COCOON 05)** Lecture Notes in Computer Science vol. 3196 Springer, 2004.

Nehel, Linder CR et al. Reconstructing reticulate evolution in species theory and practice. **Journal of Computational Biology** 12

Olsen G, Madsen H, Hestrov B et al. A fast DNA algorithm for construction of phylogenetic trees of DNA sequences using parallel processing. **Comput Appl Biosci** 10

Rieseberg LH, Rymond Q, Rosenthal DM et al. Major ecological transitions in island sunflowers facilitated by hybridization. **Science** 301

Schmidt HA. Phylogenetic trees from ordered sets. PhD thesis Heinrich Heine University, Düsseldorf.

Sepe C and Steel M. **Phylogenetics**. Oxford University Press.

## Appendix

### A Pseudocode

Here we present the pseudocode of **HybridNumber**. For a rooted binary phylogenetic  $X$  tree and subset  $A$  of  $X$  we denote the initial tree of connecting the elements in  $A$  by  $A$ . Further we denote the tree formed by replacing cluster  $A$  with the new set  $C$  by  $A \rightarrow C$ . If  $B$  is a subset of  $X$  we use  $B$  to denote the phylogenetic tree obtained from  $X$  by deleting each of the elements in  $B$  and suppressing any resulting degree two vertices. Finally,  $E$  denotes the forest obtained from the tree  $X$  by deleting the edges in the set  $E$ . Because of the chain reduction rule the input to **HybridNumber** includes a weight function  $w$  on pairs of taxa. This can be taken to be zero for pairs in the initial input.

```
Algorithm A.1: HybridNumber  $\mathcal{S}, w$   
 $\mathcal{S}, w \leftarrow \text{SubtreeReduction } \mathcal{S}, w$   
 $\mathcal{S}, w \leftarrow \text{ChainReduction } \mathcal{S}, w$   
if  $\text{initial condition on cluster } C \text{ of } \mathcal{S}_{\text{nd}}$  and  
   $C < \text{number of t of } \mathcal{S}$   
do  $\left\{ \begin{array}{l} \mathcal{S}_{1, 1}, w_1, \mathcal{S}_{2, 2}, w_2 \leftarrow \text{ClusterReduction } \mathcal{S}, w \\ h_1 \leftarrow \text{ExhaustiveSearch } \mathcal{S}_{1, 1}, w_1, \mathcal{S}_{2, 2}, w_2 \end{array} \right.$ 
```

```

Algorithm A.4: ClusterReduction  $\mathcal{S}, w$ 
 $C$  ← initialization of cluster of  $\mathcal{S}$ 
 $\mathcal{S}_1, \mathcal{S}_2 \leftarrow \mathcal{S}$ 
 $i \leftarrow 1$ 
 $C \leftarrow \mathcal{S}_i$ 
 $w_1 \leftarrow w$  restricted to pairs of trees in  $C$ 
 $w_2 \leftarrow w$  restricted to pairs of trees not in  $C$ 
return  $\mathcal{S}_1, w_1, \mathcal{S}_2, w_2$ 

```

```

Algorithm A.5: ExhaustiveSearch  $\mathcal{S}, w$ 
if  $\mathcal{S}$  return
 $h \leftarrow$  number of leaves of  $\mathcal{S}$ 
 $i \leftarrow$ 
repeat
  for each  $E$  subset of the edges of  $\mathcal{S}$  such that  $|E| = i$ 
    do
       $\mathcal{S}' \leftarrow \mathcal{S} \setminus E$ 
      if  $\mathcal{S}'$  is a non-cyclic green forest of  $\mathcal{S}$ 
        do
           $P \leftarrow$  pairs  $(a, b)$  of isolated trees in  $\mathcal{S}'$ 
           $h' \leftarrow i - \sum_{(a,b) \in P} w(a, b)$ 
          if  $h' < h$ 
            do  $h \leftarrow h'$ 
       $i \leftarrow i - 1$ 
until  $i = h$ 
return  $h$ 

```

**Remarks**

- The actual implemented algorithms contain various simplifications compared to the pseudocode in order to improve running time. In particular, these changes do not affect the theoretical worst case running time. In practice they are efficient. An example is that no green forest has isolated internal vertices, hence in the exhaustive search we do not need to consider subsets of edges of size  $i$  to delete from  $\mathcal{S}$  which contain the three edges incident with particular vertices.
- In **HybridNumber** for applying cluster reduction, the cluster reduction cannot be reduced any further using the reductions in which case we immediately call **ExhaustiveSearch**. However, it may not be possible to further reduce the remainder of the trees and so we call **HybridNumber**.

