

Journal of Computational Biology

Journal of Computational Biology: <http://mc.manuscriptcentral.com/liebert/jcb>

Quantifying Hybridization in Realistic Time

Journal:	<i>Journal of Computational Biology</i>
Manuscript ID:	JCB-2009-0166
Manuscript Type:	Original Paper
Date Submitted by the Author:	17-Jul-2009
Complete List of Authors:	Collins, Joshua; University of Canterbury, Mathematics and Statistics Linz, Simone; University of California, Davis, Computer Science Semple, Charles; University of Canterbury, Mathematics and Statistics
Keyword:	combinatorics, evolution, NP-completeness, PHYLOGENETIC TREES, algorithms
Abstract:	<p>Recently, numerous practical and theoretical studies in evolutionary biology aim at calculating the extent to which reticulation---for example horizontal gene transfer, hybridization, or recombination---has influenced the evolution for a set of present-day species. It has been shown that inferring the minimum number of hybridization events that is needed to simultaneously explain the evolutionary history for a set of trees is an NP-hard and also fixed-parameter tractable problem. In this paper, we give a new fixed-parameter algorithm for computing the minimum number of hybridization events for when two rooted binary phylogenetic trees are given. This newly developed algorithm is based on interleaving---a technique using repeated kernelization steps that are applied throughout the exhaustive search part of a fixed-parameter algorithm. To show that our algorithm runs efficiently to be applicable to a wide range of practical problem instances, we apply it to a grass data set and highlight the significant improvements in terms of running times in comparison to an algorithm which has previously been implemented.</p>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



For Peer Review

Quantifying Hybridization in Realistic Time

Joshua Collins¹, Simone Linz^{2,*}, and Charles Semple³

¹Biomathematics Research Centre

Department of Mathematics and Statistics

University of Canterbury

Private Bag 4800

Christchurch 8140, New Zealand

Email: j.collins@math.canterbury.ac.nz

Telephone: +64 3 3642987 ext 3858

Fax: +64 3 3642587

²Department of Computer Science

University of California, Davis

1 Shields Ave.

Davis 95616, USA

Email: linzs@cs.ucdavis.edu

Telephone: +1 530 7523785

Fax: +1 530 7524767

³Biomathematics Research Centre

Department of Mathematics and Statistics

University of Canterbury

Private Bag 4800

Christchurch 8140, New Zealand

Email: c.semple@math.canterbury.ac.nz

Telephone: +64 3 364 2987 ext 8349

Fax: +64 3 3642587

*Corresponding author.

Running head. Quantifying Hybridization in Realistic Time

Key words: reticulate evolution, hybridization, agreement forests, interleaving, fixed-parameter tractability.

1 Abstract

1
2
3
4
5
6
7
8 Recently, numerous practical and theoretical studies in evolutionary biology aim at calcu-
9
10 lating the extent to which reticulation—for example horizontal gene transfer, hybridiza-
11
12 tion, or recombination—has influenced the evolution for a set of present-day species. It
13
14 has been shown that inferring the minimum number of hybridization events that is needed
15
16 to simultaneously explain the evolutionary history for a set of trees is an NP-hard and
17
18 also fixed-parameter tractable problem. In this paper, we give a new fixed-parameter
19
20 algorithm for computing the minimum number of hybridization events for when two
21
22 rooted binary phylogenetic trees are given. This newly developed algorithm is based on
23
24 interleaving—a technique using repeated kernelization steps that are applied throughout
25
26 the exhaustive search part of a fixed-parameter algorithm. To show that our algorithm
27
28 runs efficiently to be applicable to a wide range of practical problem instances, we apply
29
30 it to a grass data set and highlight the significant improvements in terms of running
31
32 times in comparison to an algorithm which has previously been implemented.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

2 Introduction

Molecular evolution (phylogenetics) is a lively field of research that is affected by a variety of scientific disciplines. Viewing it from the perspective of computer science, the NP-hardness of many fundamental problems in phylogenetics makes it a challenging subject to study. The theoretically well-analyzed and widely-applied tree reconstruction methods of maximum parsimony and maximum likelihood are prominent examples of such NP-hard problems (Foulds and Graham, 1982; Chor and Tuller, 2005; Roch, 2006). In this paper, we consider a particular NP-hard problem that is fundamental in the study of reticulate evolution.

Evolutionary biologists often observe inconsistencies amongst phylogenetic trees that represent the evolution of different parts of present-day species genomes. Such inconsistencies can essentially be caused either by reticulation events like horizontal gene transfer, hybridization, and recombination, or by non-biological processes like sequencing errors or signals in the data that may yield trees whose branching patterns do not always represent the correct evolutionary history. Here, we assume that hybridization (as a representative of reticulation) has led to the observed inconsistencies. In this case, it may be more appropriate to represent the evolutionary history of a set of species by a phylogenetic network rather than a phylogenetic tree since the parents of a hybrid taxa belong to two different species. It is consequently desirable to calculate a hybridization network that simultaneously explains the evolutionary histories for a given set of trees and minimizes the number of hybridization events. The reason for the latter is that it quantifies the significance of hybridization for the evolution of the species under consideration. However, computing this minimum number is NP-hard even for two trees (Bordewich and Semple, 2007a). Known as MINIMUM HYBRIDIZATION, it is this two-tree problem that is the focus of this paper.

To overcome the computational burden of NP-hard problems, one frequently resorts to approximation algorithms, heuristics, or solving polynomial-time restrictions of the

1
2
3
4 problem. However, the solutions obtained from these approaches are not always accept-
5
6 able; for example, this may be due to complex and expensive processes that were needed
7
8 to generate certain data sets. In such cases, fixed-parameter algorithms have proven to
9
10 be a valuable tool to calculate the exact solution of a computationally-hard problem.
11
12 Mathematically speaking, a problem of size n , parameterized by k , is *fixed-parameter*
13
14 *tractable* if it can be computed in $O(f(k)p(n))$, where f is a computable function and p
15
16 is a fixed polynomial. The importance of this running time is in the separation of the
17
18 variables k and n . Thus, if k is small, the problem may be tractable in reasonable time
19
20 even if n is large. For a more detailed description of fixed-parameter tractability (FPT),
21
22 we refer the interested reader to Downey and Fellows (1998) and Flum and Grohe (2006).
23
24

25
26 MINIMUM HYBRIDIZATION and various other problems in computational biology are
27
28 known to be fixed-parameter tractable (for example, Ávila et al., 2006; Bordewich and
29
30 Semple, 2007b; Gramm et al., 2008). However, practical fixed-parameter algorithms
31
32 that have been applied to biological data sets rarely exist. Recently, Bordewich et al.
33
34 (2007) implemented a fixed-parameter algorithm for MINIMUM HYBRIDIZATION. By
35
36 applying this algorithm to a grass data set, the authors subsequently showed that many
37
38 problem instances were computable within a couple of minutes. However, there were
39
40 several instances to which the algorithm did not return the exact answer in reasonable
41
42 time. In particular, for three tree pairs, the running time to calculate the exact solution
43
44 was at least two days. Other studies in computational biology that have introduced
45
46 fixed-parameter algorithms and applied them to biological or synthetic data sets are for
47
48 example described in Dehne et al. (2006), and Gramm and Niedermeier (2002, 2003).
49
50

51
52 To keep up with the constant progress in molecular biology, which primarily originates
53
54 from the development of efficient DNA sequencing technologies, it is of importance to
55
56 develop new and to further improve existing fixed-parameter algorithms such that they
57
58 can cope with an increasing data set size. A common technique for obtaining a fixed-
59
60 parameter algorithm for a problem is *kernelization*. The aim of this technique is to shrink
the size of the initial problem instance to its difficult core by quickly resolving those

1
2
3 parts of the problem that are easily dealt with. Kernelization is the technique used by
4 Bordewich et al. (2007) in their implementation for solving MINIMUM HYBRIDIZATION.
5 To complement kernelization, *interleaving* has been introduced as a new method in the
6 design of fixed-parameter algorithms (Niedermeier and Rossmanith, 2000). Interleaving
7 refers to repeated kernelization steps while one systematically processes a bounded search
8 tree. Apart from Abu-Khzam et al. (2006) and Dehne et al. (2006), where the authors
9 showed that interleaving has a positive impact on the overall running time of a fixed-
10 parameter algorithm, this technique has so far attracted more attention in theoretical
11 analyses concerned with FPT than in practical studies.
12
13
14
15
16
17
18
19
20
21
22

23 Making use of interleaving, we present a fixed-parameter algorithm for solving MIN-
24 IMUM HYBRIDIZATION that performs significantly better than that given by Bordewich
25 et al. (2007). This improvement is highlighted by the fact that all instances of the grass
26 data set described above can be solved in less than eleven minutes. As an example of
27 the new algorithm's performance, an instance for which the previously implemented al-
28 gorithm did not return the solution within two days, can now be calculated in less than
29 a minute.
30
31
32
33
34
35
36
37

38 The new algorithm—called HYBRIDINTERLEAVE—has been implemented in Java and
39 is available for application at
40
41
42

43 <http://www.math.canterbury.ac.nz/~c.semler/software.shtml>
44
45
46

47 or
48
49

50 <http://wwwcsif.cs.ucdavis.edu/~linzs/>.
51
52
53

54 To start a calculation with HYBRIDINTERLEAVE, the algorithm requires the two input
55 trees to be given in Newick format. As output, the program provides the user with the
56 minimum number of hybridization events that explain the two input trees.
57
58
59
60

This paper is organized as follows. The next section contains some preliminaries and

1
2
3
4 formally states the decision problem MINIMUM HYBRIDIZATION for which a previously
5 established fixed-parameter algorithm is summarized in Section 4. The new algorithm
6 HYBRIDINTERLEAVE and its proof of correctness are given in Section 5, while Section 6
7 analyzes the running times of HYBRIDINTERLEAVE when applied to a grass data set
8 and compares it with the running times of the recently implemented algorithm HYBRID-
9 NUMBER (Bordewich et al., 2007). We end the paper with some concluding remarks in
10 Section 7. Unless otherwise stated, the notation and terminology follows Semple and
11 Steel (2003).
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

3 Preliminaries

This section provides preliminary definitions which are used throughout the rest of this paper and formally states the decision problem MINIMUM HYBRIDIZATION.

3.1 Phylogenetic Trees

A *rooted binary phylogenetic X -tree* \mathcal{T} is a rooted tree whose root has degree two while all other interior vertices have degree three. The leaf set X is the label set of \mathcal{T} and frequently denoted by $\mathcal{L}(\mathcal{T})$. Furthermore, a subset A of X is a *cluster* of \mathcal{T} if there is a vertex v whose set of descendants is precisely A . We view v as an ancestor and descendant of itself.

In the course of this paper, two types of subtrees play an important role. Let X' be a subset of X , and let \mathcal{T} be a rooted phylogenetic X -tree. The minimal rooted subtree of \mathcal{T} that connects all leaves in X' is denoted by $\mathcal{T}(X')$. Furthermore, the subtree obtained from $\mathcal{T}(X')$ by contracting all non-root degree-2 vertices is the *restriction of \mathcal{T} to X'* and is denoted by $\mathcal{T}|X'$. Lastly, a subtree is *pendant* if it can be detached from \mathcal{T} by deleting a single edge.

3.2 Hybridization Networks

A *hybridization network* \mathcal{H} on a set X is a rooted acyclic digraph with root ρ such that the following properties are satisfied:

- (i) X is the set of vertices of out-degree 0,
- (ii) the out-degree of ρ is at least 2, and
- (iii) for all vertices v with $d^+(v) = 1$, we have $d^-(v) = 2$,

where $d^+(v)$ and $d^-(v)$ denote the out-degree and in-degree of v , respectively. The set X represents a collection of present-day taxa, and vertices of in-degree two represent an exchange of genetic information between their parents. Generically, we call these vertices *hybridization vertices*.

Figure 1

To quantify the number of hybridization events, the *hybridization number* of \mathcal{H} is the total number of hybridization vertices. Observe that if \mathcal{H} is a rooted binary phylogenetic tree, then $h(\mathcal{H}) = 0$. Ignoring the thickness of the arcs, the left-hand-side of Figure 1 shows a hybridization network \mathcal{H} whose hybridization number is three.

Now let \mathcal{H} be a hybridization network on X , and let \mathcal{T} be a rooted binary phylogenetic X' -tree with $X' \subseteq X$. We say that \mathcal{T} is *displayed* by \mathcal{H} if \mathcal{T} can be obtained from \mathcal{H} by deleting a subset of its edges and any resulting degree-0 vertices, and then contracting edges. For example, the rooted binary phylogenetic tree \mathcal{T} shown in Figure 1 is displayed by the hybridization network \mathcal{H} of the same figure. Intuitively, if \mathcal{H} displays \mathcal{T} , then all of the ancestral relationships visualized by \mathcal{T} are visualized by \mathcal{H} . Extending the definition of the hybridization number to two rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} , we set

$$h(\mathcal{S}, \mathcal{T}) = \min\{h(\mathcal{H}) : \mathcal{H} \text{ is a hybridization network that displays } \mathcal{S} \text{ and } \mathcal{T}\}.$$

With the above definition, we now formally state MINIMUM HYBRIDIZATION.

Problem: MINIMUM HYBRIDIZATION($\mathcal{S}, \mathcal{T}, k$)

Instance: Two rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} , and an integer k .

Question: Is $h(\mathcal{S}, \mathcal{T}) < k$?

3.3 Agreement Forests

Figure

2

Originating from an idea in Hein et al. (1996), Bordewich and Semple (2007a) showed that MINIMUM HYBRIDIZATION is NP-complete by using a characterization of the problem in terms of agreement forests. Such forests play a fundamental role in this paper. For the purpose of the upcoming definitions, we regard the root of a rooted binary phylogenetic X -tree \mathcal{T} as a vertex ρ at the end of a pendant edge adjoined to the original root. For an example of two such trees, see Figure 2. Furthermore, we view ρ as an element of the label set of \mathcal{T} ; thus $\mathcal{L}(\mathcal{T}) = X \cup \{\rho\}$. Now, let \mathcal{S} and \mathcal{T} be two rooted binary phylogenetic X -trees. An *agreement forest* $\mathcal{F} = \{\mathcal{L}_\rho, \mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_k\}$ for \mathcal{S} and \mathcal{T} is a partition of $\mathcal{L}(\mathcal{S})$ such that $\rho \in \mathcal{L}_\rho$ and the following conditions are fulfilled:

- (1) For all $i \in \{\rho, 1, 2, \dots, k\}$, we have $\mathcal{S}|_{\mathcal{L}_i} \cong \mathcal{T}|_{\mathcal{L}_i}$.
- (2) The trees in $\{\mathcal{S}(\mathcal{L}_i) : i \in \{\rho, 1, 2, \dots, k\}\}$ and $\{\mathcal{T}(\mathcal{L}_i) : i \in \{\rho, 1, 2, \dots, k\}\}$ are vertex-disjoint subtrees of \mathcal{S} and \mathcal{T} , respectively.

As an example, two agreement forests for the two rooted binary phylogenetic trees \mathcal{S} and \mathcal{T} shown in Figure 2 are $\mathcal{F} = \{\{\rho, 7\}, \{1, 2, 3\}, \{4, 5, 6\}\}$ and $\mathcal{F}' = \{\{\rho, 1, 2, 3, 7\}, \{4\}, \{5\}, \{6\}\}$.

A characterization of the minimum number $h(\mathcal{S}, \mathcal{T})$ of hybridization events in terms of agreement forests requires an additional condition. Without going into details, this condition avoids the possibility of species inheriting genetic material from their own descendants. Let $\mathcal{F} = \{\mathcal{L}_\rho, \mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_k\}$ be an agreement forest for \mathcal{S} and \mathcal{T} . Let $G_{\mathcal{F}}$ be the directed graph that has vertex set \mathcal{F} and an arc from \mathcal{L}_i to \mathcal{L}_j precisely if $i \neq j$, and either

- (1) the root of $\mathcal{S}(\mathcal{L}_i)$ is an ancestor of the root of $\mathcal{S}(\mathcal{L}_j)$ in \mathcal{S} or
- (2) the root of $\mathcal{T}(\mathcal{L}_i)$ is an ancestor of the root of $\mathcal{T}(\mathcal{L}_j)$ in \mathcal{T} .

Figure

3

1
2
3
4 We call \mathcal{F} an *acyclic-agreement forest* for \mathcal{S} and \mathcal{T} if $G_{\mathcal{F}}$ has no directed cycle. Moreover,
5
6 if \mathcal{F} contains the smallest number of parts over all acyclic-agreement forests for \mathcal{S} and
7
8 \mathcal{T} , we say that \mathcal{F} is a *maximum-acyclic-agreement forest* for \mathcal{S} and \mathcal{T} , in which case, we
9
10 denote this number minus one by $m_a(\mathcal{S}, \mathcal{T})$. Figure 3 shows the two digraphs $G_{\mathcal{F}}$ and
11
12 $G_{\mathcal{F}'}$ that are associated with the agreement forests \mathcal{F} and \mathcal{F}' , respectively, for the two
13
14 rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} depicted in Figure 2. Note that, as $G_{\mathcal{F}'}$ is
15
16 acyclic, \mathcal{F}' is an acyclic-agreement forest for \mathcal{S} and \mathcal{T} while \mathcal{F} is no such forest. Indeed,
17
18 \mathcal{F}' is a maximum-acyclic-agreement forest for \mathcal{S} and \mathcal{T} . Baroni et al. (2005) established
19
20 the following characterization.
21

22
23 **Theorem 1.** *Let \mathcal{S} and \mathcal{T} be two rooted binary phylogenetic X -trees. Then*
24

$$25 \quad h(\mathcal{S}, \mathcal{T}) = m_a(\mathcal{S}, \mathcal{T}).$$

26
27
28
29

30 It is this characterization that was used to show that MINIMUM HYBRIDIZATION is
31
32 NP-complete.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

4 Overview of a Known Fixed-Parameter Algorithm for MINIMUM HYBRIDIZATION

In this section, we summarize the basic ideas of the first fixed-parameter algorithm for MINIMUM HYBRIDIZATION. This algorithm is based on an earlier result that showed, for a pair of rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} , MINIMUM HYBRIDIZATION is fixed-parameter tractable with $h(\mathcal{S}, \mathcal{T})$ being the parameter (Bordewich and Semple, 2007b). In establishing this result, the authors used two reductions—called the subtree and chain reduction—that kernelize \mathcal{S} and \mathcal{T} to two smaller trees whose number of leaves is linear in $h(\mathcal{S}, \mathcal{T})$.

Before detailing these reductions, we need some additional definitions. Let \mathcal{T} be a rooted binary phylogenetic X -tree. An n -chain of \mathcal{T} is an ordered tuple (a_1, a_2, \dots, a_n) of elements in X such that the parent of a_1 is either the same as the parent of a_2 or a child of the parent of a_2 and, for all $i \in \{2, 3, \dots, n-1\}$, the parent of a_i is a child of the parent of a_{i+1} . For example, referring to Figure 2, $(1, 2, 3, 4, 5, 6)$ and $(4, 5, 6)$ are chains of \mathcal{S} . Furthermore, let \mathcal{S} and \mathcal{T} be two rooted binary phylogenetic X -trees, and let P be a disjoint collection of 2-element subsets of X such that each pair $\{a, b\} \in P$ is a common 2-chain of \mathcal{S} and \mathcal{T} . Let the weight function $w : P \rightarrow \mathbb{Z}^+$ assign each element of P a positive integer weight. We refer to \mathcal{S} and \mathcal{T} with an associated weight function w as a *pair of weighted phylogenetic trees on X* . The set P is referred to as the set of *weighted 2-chains* of \mathcal{S} and \mathcal{T} . Unless otherwise stated, we will use w to denote the weighting of the 2-chains in P .

We now describe the reductions. Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees with weighted 2-chain set P .

Subtree reduction. Replace a maximal pendant subtree that is common to \mathcal{S} and \mathcal{T} by a single leaf with a new label. Furthermore, delete all members in P whose elements label leaves of the pendant subtree under consideration.

Chain reduction. Replace a maximal n -chain (a_1, a_2, \dots, a_n) with $n > 2$ that occurs identically in \mathcal{S} and \mathcal{T} by a 2-chain with new labels a and b . Furthermore, add the 2-element set $\{a, b\}$ to P with an associated weight of

$$w(\{a, b\}) = n - 2 + \sum_{\substack{\{a_i, a_j\} \in P \\ \text{and } a_i, a_j \in \\ \{a_1, \dots, a_n\}}} w(\{a_i, a_j\}),$$

and delete all members in P whose elements are in $\{a_1, a_2, \dots, a_n\}$. An explicit example of the chain reduction is shown in Figure 4, where the two rooted binary phylogenetic trees \mathcal{S}' and \mathcal{T}' have been obtained from \mathcal{S} and \mathcal{T} , which are shown in Figure 2, by replacing the 3-chain $(1, 2, 3)$ with the 2-chain (a, b) , and similarly, the 3-chain $(4, 5, 6)$ with (c, d) . Note that $w(\{a, b\}) = w(\{c, d\}) = 1$.

The correspondence between the trees resulting from repeated applications of the subtree and chain reductions, and the initial two trees is given in the next lemma. This correspondence is done via a notion of agreement forests that extends acyclic-agreement forests. An agreement forest \mathcal{F} for two rooted weighted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} with weighted 2-chain set P is called *legitimate* if it is acyclic and the following property holds:

- (P) For each $\{a, b\} \in P$, either $\{a\}$ and $\{b\}$ are elements in \mathcal{F} , or there exists an element of \mathcal{F} , say \mathcal{L} , such that $\{a, b\} \subseteq \mathcal{L}$.

Let \mathcal{F} be a legitimate-agreement forest for \mathcal{S} and \mathcal{T} . Set

$$w_c(\mathcal{F}, P) = \sum_{\substack{\{a, b\} \in P; \\ \{a\}, \{b\} \in \mathcal{F}}} w(\{a, b\}),$$

we define the *weight* of \mathcal{F} as

$$w(\mathcal{F}) = |\mathcal{F}| - 1 + w_c(\mathcal{F}, P)$$

and set $f(\mathcal{S}, \mathcal{T})$ to be the minimum weight of a legitimate-agreement forest for \mathcal{S} and \mathcal{T} . Note that we always have $f(\mathcal{S}, \mathcal{T}) \geq h(\mathcal{S}, \mathcal{T})$, and $f(\mathcal{S}, \mathcal{T}) = h(\mathcal{S}, \mathcal{T})$ whenever P is empty. The next two lemmas are central to showing that MINIMUM HYBRIDIZATION is fixed-parameter tractable (Bordewich and Semple, 2007b).

Lemma 2. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees, and let \mathcal{S}' and \mathcal{T}' be two weighted rooted binary phylogenetic X' -trees that have been obtained from \mathcal{S} and \mathcal{T} , respectively, by applying the subtree or chain reduction. Then $f(\mathcal{S}, \mathcal{T}) = f(\mathcal{S}', \mathcal{T}')$.*

Lemma 3. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees whose weighted 2-chain set P is empty. Furthermore, let \mathcal{S}' and \mathcal{T}' be two weighted rooted binary phylogenetic X' -trees that have been obtained from \mathcal{S} and \mathcal{T} , respectively, by repeatedly applying the subtree and chain reduction until no further reduction is possible. Then $|X'| \leq 14h(\mathcal{S}, \mathcal{T})$.*

Cluster reduction. Besides repeatedly applying the subtree and chain reductions to kernelize a problem instance of MINIMUM HYBRIDIZATION before exhaustively calculating a legitimate-agreement forest of minimum weight, we can use a third reduction that breaks a problem instance of MINIMUM HYBRIDIZATION into two smaller subproblems. This reduction is depicted in Figure 5 and can be repeatedly intertwined with the other two reductions before the inevitable exhaustive search part of the algorithm. How the two smaller problem instances relate to the original instance is described in the next corollary. Due to Linz (2008), this corollary generalizes the unweighted version given by Baroni et al. (2006).

Corollary 4. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees with weighted 2-chain set P , and let A be a common minimal cluster of both \mathcal{S} and \mathcal{T} with $|A| \geq 2$. Then,*

$$f(\mathcal{S}, \mathcal{T}) = f(\mathcal{S}|A, \mathcal{T}|A) + f(\mathcal{S}_a, \mathcal{T}_a),$$

where \mathcal{S}_a and \mathcal{T}_a are the trees obtained from \mathcal{S} and \mathcal{T} , respectively, by replacing the pendant subtree whose label set is precisely A with a new leaf labeled a .

In the last corollary, the weighted 2-chain sets, P_A and P_a say, of $\mathcal{S}|A$ and $\mathcal{T}|A$, and \mathcal{S}_a and \mathcal{T}_a , respectively, are

$$P_A = \{\{\ell, \ell'\} : \{\ell, \ell'\} \subseteq P \text{ and } \ell, \ell' \in A\}$$

and

$$P_a = \{\{\ell, \ell'\} : \{\ell, \ell'\} \subseteq P \text{ and } \ell, \ell' \notin A\},$$

where, for both sets, the weight of each element $\{\ell, \ell'\}$ is equal to the weighting of $\{\ell, \ell'\}$ in the weight function associated with P .

Remarks.

- (i) Note that the cluster reduction can repeatedly be applied to break \mathcal{S} and \mathcal{T} into as many smaller tree pairs as possible by setting A to be a minimal common cluster of \mathcal{S} and \mathcal{T} with $|A| \geq 2$, and resetting \mathcal{S} and \mathcal{T} to be \mathcal{S}_a and \mathcal{T}_a , respectively, before applying this reduction again until $\mathcal{S} \cong \mathcal{T}$.
- (ii) We impose *maximality* on a common pendant subtree and a common n -chain and *minimality* on a common cluster to guarantee that the corresponding label set of any such common feature intersects each member of P in either both elements or neither.

5 A New Algorithm for MINIMUM HYBRIDIZATION

In this section, we present our fixed-parameter algorithm HYBRIDINTERLEAVE for MINIMUM HYBRIDIZATION. It makes use of the subtree, chain, and cluster reductions, but importantly, in terms of obtaining significantly decreased running times (see Section 6), it additionally uses interleaving.

Before outlining HYBRIDINTERLEAVE and giving its pseudocode, we state two lemmas that are central to its correctness and description. Let \mathcal{T} be a rooted binary phylogenetic X -tree, and let ℓ and ℓ' be elements of X . To ease reading in this section, we use $\mathcal{T}[-\ell]$ to denote $\mathcal{T}|(\mathcal{L}(\mathcal{T}) - \{\ell\})$ and $\mathcal{T}[-\ell, \ell']$ to denote $\mathcal{T}|(\mathcal{L}(\mathcal{T}) - \{\ell, \ell'\})$. Furthermore, let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees with weighted 2-chain set P . If ℓ is contained in a member of P , we say that ℓ *crosses* P .

Lemma 5. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees with weighted 2-chain set P . Suppose that \mathcal{S} and \mathcal{T} have no common pendant subtree whose leaf set size is at least 2. Then, for each $\ell \in X$, we have*

$$f(\mathcal{S}, \mathcal{T}) \leq f(\mathcal{S}[-\ell, \ell'], \mathcal{T}[-\ell, \ell']) + 2 + w(\{\ell, \ell'\})$$

if ℓ crosses P with $\{\ell, \ell'\} \in P$, and

$$f(\mathcal{S}, \mathcal{T}) \leq f(\mathcal{S}[-\ell], \mathcal{T}[-\ell]) + 1$$

otherwise.

Proof. First assume that ℓ crosses P in an element $\{\ell, \ell'\}$. Let \mathcal{F}_ℓ be a legitimate-agreement forest for $\mathcal{S}[-\ell, \ell']$ and $\mathcal{T}[-\ell, \ell']$ of minimum weight. Then it is easily checked that

$$\mathcal{F} = \mathcal{F}_\ell \cup \{\{\ell\}, \{\ell'\}\}$$

is a legitimate-agreement forest for \mathcal{S} and \mathcal{T} . Moreover, we have $|\mathcal{F}| = |\mathcal{F}_\ell| + 2$ and

$w_c(\mathcal{F}_\ell, P - \{\ell, \ell'\}) = w_c(\mathcal{F}, P) - w(\{\ell, \ell'\})$. Hence,

$$\begin{aligned}
 f(\mathcal{S}[-\ell, \ell'], \mathcal{T}[-\ell, \ell']) + 2 + w(\{\ell, \ell'\}) & \quad (1) \\
 &= |\mathcal{F}_\ell| - 1 + w_c(\mathcal{F}_\ell, P - \{\ell, \ell'\}) + 2 + w(\{\ell, \ell'\}) \\
 &= |\mathcal{F}| - 1 + w_c(\mathcal{F}, P) \\
 &\geq f(\mathcal{S}, \mathcal{T}).
 \end{aligned}$$

Now assume that ℓ does not cross P . Let \mathcal{F}_ℓ be a legitimate-agreement forest for $\mathcal{S}[-\ell]$ and $\mathcal{T}[-\ell]$ of minimum weight. Again, it is clear that

$$\mathcal{F} = \mathcal{F}_\ell \cup \{\ell\}$$

is a legitimate-agreement forest for \mathcal{S} and \mathcal{T} . Moreover, we have $|\mathcal{F}| = |\mathcal{F}_\ell| + 1$ and $w_c(\mathcal{F}_\ell, P) = w_c(\mathcal{F}, P)$. Thus

$$\begin{aligned}
 f(\mathcal{S}[-\ell], \mathcal{T}[-\ell]) + 1 &= |\mathcal{F}_\ell| - 1 + w_c(\mathcal{F}_\ell, P) + 1 & (2) \\
 &= |\mathcal{F}| - 1 + w_c(\mathcal{F}, P) \\
 &\geq f(\mathcal{S}, \mathcal{T}).
 \end{aligned}$$

Inequalities (1) and (2) establish the lemma. \square

Lemma 6. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees with weighted 2-chain set P . Suppose that \mathcal{S} and \mathcal{T} have no common pendant subtree whose leaf set size is at least 2. Then there exists an element $\ell \in X$ such that either ℓ crosses P with $\{\ell, \ell'\} \in P$ and*

$$f(\mathcal{S}, \mathcal{T}) = f(\mathcal{S}[-\ell, \ell'], \mathcal{T}[-\ell, \ell']) + 2 + w(\{\ell, \ell'\}),$$

or ℓ does not cross P and

$$f(\mathcal{S}, \mathcal{T}) = f(\mathcal{S}[-\ell], \mathcal{T}[-\ell]) + 1.$$

Proof. Let $\mathcal{F} = \{\mathcal{L}_\rho, \mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_k\}$ be a legitimate-agreement forest for \mathcal{S} and \mathcal{T} of minimum weight. First observe that, as $G_{\mathcal{F}}$ is acyclic, it has a vertex \mathcal{L}_i with $i \in$

1
2
3
4 $\{\rho, 1, 2, \dots, k\}$ whose out-degree is zero. Furthermore, since \mathcal{S} and \mathcal{T} have no common
5
6 pendant subtree whose leaf set size is at least 2, \mathcal{L}_i is a singleton in \mathcal{F} . Since ρ is never
7
8 a singleton in \mathcal{F} by Lemma 1 of Baroni et al. (2005), we may assume that $\mathcal{L}_i = \{\ell\}$,
9
10 where $\ell \in X$.

11
12 First assume that ℓ crosses P in an element $\{\ell, \ell'\}$. Since \mathcal{F} is legitimate, $\{\ell'\} \in \mathcal{F}$
13
14 as $\{\ell\} \in \mathcal{F}$, and so

$$15 \quad \mathcal{F}_\ell = \mathcal{F} - \{\{\ell\}, \{\ell'\}\}$$

16
17 is a legitimate-agreement forest for $\mathcal{S}[-\ell, \ell']$ and $\mathcal{T}[-\ell, \ell']$. Furthermore, we have $|\mathcal{F}| =$
18
19 $|\mathcal{F}_\ell| + 2$ and $w_c(\mathcal{F}, P) = w_c(\mathcal{F}_\ell, P - \{\ell, \ell'\}) + w(\{\ell, \ell'\})$. It now follows that

$$20 \quad \begin{aligned} 21 \quad f(\mathcal{S}, \mathcal{T}) &= |\mathcal{F}| - 1 + w_c(\mathcal{F}, P) & (3) \\ 22 \quad &= |\mathcal{F}_\ell| + 2 - 1 + w_c(\mathcal{F}_\ell, P - \{\ell, \ell'\}) + w(\{\ell, \ell'\}) \\ 23 \quad &\geq f(\mathcal{S}[-\ell, \ell'], \mathcal{T}[-\ell, \ell']) + 2 + w(\{\ell, \ell'\}). \end{aligned}$$

24
25
26
27
28
29
30
31
32
33 Second assume that ℓ does not cross P . Since \mathcal{F} is a legitimate-agreement forest for
34
35 \mathcal{S} and \mathcal{T} ,

$$36 \quad \mathcal{F}_\ell = \mathcal{F} - \{\ell\}$$

37
38 is such a forest for $\mathcal{S}[-\ell]$ and $\mathcal{T}[-\ell]$. Furthermore, we have $|\mathcal{F}| = |\mathcal{F}_\ell| + 1$ and $w_c(\mathcal{F}, P) =$
39
40 $w_c(\mathcal{F}_\ell, P)$. It now follows that

$$41 \quad \begin{aligned} 42 \quad f(\mathcal{S}, \mathcal{T}) &= |\mathcal{F}| - 1 + w_c(\mathcal{F}, P) & (4) \\ 43 \quad &= |\mathcal{F}_\ell| + 1 - 1 + w_c(\mathcal{F}_\ell, P) \\ 44 \quad &\geq f(\mathcal{S}[-\ell], \mathcal{T}[-\ell]) + 1. \end{aligned}$$

45
46
47
48
49
50
51
52
53
54 Combining (3) and (4) with Lemma 5 gives the lemma. \square

55
56
57
58 We next give a brief outline of the algorithm HYBRIDINTERLEAVE before detailing
59
60 its pseudocode. The algorithm takes as input two rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} , and an integer k , and outputs $h(\mathcal{S}, \mathcal{T})$ precisely if $h(\mathcal{S}, \mathcal{T}) < k$. It starts

with initializing the variable P , the set of weighted 2-chains, that has previously been obtained by applying a chain reduction to \mathcal{S} and \mathcal{T} . Recall that w is the weight function associated with P . HYBRIDINTERLEAVE then directly calls the subroutine INTERLEAVE which contains the key features of this algorithm.

If $k > 0$, INTERLEAVE initially finds all maximal pendant subtrees that are common to \mathcal{S} and \mathcal{T} . If the resulting two trees have a label set size of at most 3, then, as $\rho \in \mathcal{L}(\mathcal{S})$, they are identical. Consequently, INTERLEAVE directly returns 0 as the minimum weight for a legitimate-agreement forest of \mathcal{S} and \mathcal{T} . Otherwise, the algorithm proceeds with replacing each maximal common n -chain, where $n \geq 3$, with a 2-chain. Resetting \mathcal{S} and \mathcal{T} to be the reduced weighted trees, they always have a cluster A with $2 \leq |A| < |\mathcal{L}(\mathcal{S})|$ in common which allows for an application of the cluster reduction. This reduction returns two new tree pairs. The second pair \mathcal{S}'' and \mathcal{T}'' has been obtained from \mathcal{S} and \mathcal{T} by replacing $\mathcal{S}(A)$ and $\mathcal{T}(A)$, respectively, by a new leaf while the first pair is $\mathcal{S}' = \mathcal{S}|A$ and $\mathcal{T}' = \mathcal{T}|A$ (viewing the root of \mathcal{S}' and \mathcal{T}' , respectively, as a vertex ρ' adjoined to the original root by a pendant edge). With $P' = \{\{\ell, \ell'\} \in P : \{\ell, \ell'\} \subseteq A\}$ whose associated weight function is w' , the algorithm next checks whether there exists a legitimate-agreement forest for \mathcal{S}' and \mathcal{T}' with $f(\mathcal{S}', \mathcal{T}') < k$, where w' is obtained from w by restricting its domain to members that are subsets of A . To this end, the subroutine branches into $|A|$ computational paths, where each path corresponds to an element of A and a call to INTERLEAVE. This guarantees that an element is found for which Lemma 6 holds. Furthermore, for each $\ell \in A$, the algorithm successively resets the variable h , which was originally initialized with k , to the minimum of the current value of h and the return value of the associated recursive call to INTERLEAVE increased by $2 + w'(\{\ell, \ell'\})$ if ℓ crosses P' with $\{\ell, \ell'\}$, or increased by 1 otherwise. Thus, at each step, h equals k or it contains the minimum weight over all legitimate-agreement forests for \mathcal{S}' and \mathcal{T}' that have previously been considered. After at most k iterations, INTERLEAVE($\mathcal{S}, \mathcal{T}, k$) declares $h + \text{INTERLEAVE}(\mathcal{S}'', \mathcal{T}'', w'', k - h)$, where w'' is obtained from w by restricting its domain to members that are not subsets of A . Eventually, HYBRIDINTERLEAVE

1
2
3
4 either returns $h(\mathcal{S}, \mathcal{T})$ if $h(\mathcal{S}', \mathcal{T}') = h < k$ and $h(\mathcal{S}'', \mathcal{T}'') < k - h$, or it returns k .
5
6

7 The pseudocode for HYBRIDINTERLEAVE is given below. The pseudocodes for the
8 subtree, chain, and cluster reductions are given in Bordewich et al. (2007). Because of
9 this and the description given earlier, we have omitted their respective pseudocodes.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Algorithm: HYBRIDINTERLEAVE($\mathcal{S}, \mathcal{T}, k$)

procedure INTERLEAVE($\mathcal{S}, \mathcal{T}, w, k$)

if $k \leq 0$

then return (k)

$(\mathcal{S}, \mathcal{T}, w) \leftarrow$ SUBTREEREDUCTION($\mathcal{S}, \mathcal{T}, w$)

if $|\mathcal{L}(\mathcal{S})| \leq 3$

then return (0)

$(\mathcal{S}, \mathcal{T}, w) \leftarrow$ CHAINREDUCTION($\mathcal{S}, \mathcal{T}, w$)

$(\mathcal{S}', \mathcal{T}', w', \mathcal{S}'', \mathcal{T}'', w'') \leftarrow$ CLUSTERREDUCTION($\mathcal{S}, \mathcal{T}, w$)

$h \leftarrow k$

for each $\ell \in \mathcal{L}(\mathcal{S}') - \{\rho'\}$

do $\left\{ \begin{array}{l} \text{if } \exists \ell' \in \mathcal{L}(\mathcal{S}') - \{\rho', \ell\} \text{ such that } \{\ell, \ell'\} \in \text{domain } w' \\ \text{then } \left\{ \begin{array}{l} \mathcal{S}' \leftarrow \mathcal{S}'[-\ell, \ell'] \\ \mathcal{T}' \leftarrow \mathcal{T}'[-\ell, \ell'] \\ h \leftarrow \min\{h, \text{INTERLEAVE}(\mathcal{S}', \mathcal{T}', w', h - \\ w'(\{\ell, \ell'\}) - 2) + 2 + w'(\{\ell, \ell'\})\} \end{array} \right. \\ \text{else } \left\{ \begin{array}{l} \mathcal{S}' \leftarrow \mathcal{S}'[-\ell] \\ \mathcal{T}' \leftarrow \mathcal{T}'[-\ell] \\ h \leftarrow \min\{h, \text{INTERLEAVE}(\mathcal{S}', \mathcal{T}', w', h - 1) + 1\} \end{array} \right. \end{array} \right.$

return ($h + \text{INTERLEAVE}(\mathcal{S}'', \mathcal{T}'', w'', k - h)$)

main

$P \leftarrow \emptyset$

$w : P \rightarrow \mathbb{Z}^+$

$k \leftarrow \text{INTERLEAVE}(\mathcal{S}, \mathcal{T}, w, k)$

return (k)

Remark. Let X' be the label set of the two rooted binary phylogenetic trees that result from repeated applications of the subtree and chain reduction in the *first* call to INTERLEAVE. By Lemma 3, we freely assume for the rest of the paper that the algorithm directly returns k if $|X'| > 14k$.

Figure

6

To illustrate the reductions that are performed throughout the algorithm HYBRID-INTERLEAVE, consider Figure 6 and the call HYBRIDINTERLEAVE($\mathcal{S}_1, \mathcal{T}_1, k$) with $k > 2$. The algorithm first replaces the common pendant subtree $\mathcal{S}_1|_{\{1,2\}}$ of \mathcal{S}_1 and \mathcal{T}_1 with a single leaf labeled a . Since the resulting two trees \mathcal{S}_2 and \mathcal{T}_2 cannot be reduced further under any of the three reductions, INTERLEAVE is recursively called, say for $\mathcal{S}_2[-7]$ and $\mathcal{T}_2[-7]$. Next interleaving comes into play as \mathcal{S}_3 and \mathcal{T}_3 have a common 3-chain $(4, 5, 6)$ that is replaced with the 2-chain (b, c) . Again, \mathcal{S}_4 and \mathcal{T}_4 are fully reduced and INTERLEAVE is called, say for $\mathcal{S}_4[-3]$ and $\mathcal{T}_4[-3]$. The two obtained trees \mathcal{S}_5 and \mathcal{T}_5 are identical and the algorithm ultimately returns $h = 2$ for the described computational path of the search tree.

We next establish the correctness of HYBRIDINTERLEAVE.

Theorem 7. *Let \mathcal{S} and \mathcal{T} be a pair of rooted binary phylogenetic X -trees. Then the output of HYBRIDINTERLEAVE($\mathcal{S}, \mathcal{T}, k$) is $h(\mathcal{S}, \mathcal{T})$ if and only if $h(\mathcal{S}, \mathcal{T}) < k$; otherwise it is k .*

Proof. The proof is by induction on k . If $k = 0$, then INTERLEAVE immediately returns 0, and so the theorem holds. Now suppose that $k \geq 1$ and that the theorem holds whenever the input parameter is at most $k - 1$. Because of the structure of HYBRIDINTERLEAVE and Corollary 4, to establish this part of the induction, it suffices to show that the first call to the **for each** loop correctly returns $h + \text{INTERLEAVE}(\mathcal{S}'', \mathcal{T}'', w'', k - h)$ with $h = f(\mathcal{S}', \mathcal{T}')$ if and only if $f(\mathcal{S}', \mathcal{T}') < k$, otherwise with $h = k$.

By Lemma 6, there is an $\ell \in \mathcal{L}(\mathcal{S}') - \{\rho'\}$ such that one of the following holds:

(a) If ℓ does not cross P , then

$$f(\mathcal{S}'[-\ell], \mathcal{T}'[-\ell]) = f(\mathcal{S}', \mathcal{T}') - 1.$$

(b) If ℓ crosses P with $\{\ell, \ell'\} \in P$, then

$$f(\mathcal{S}'[-\ell, \ell'], \mathcal{T}'[-\ell, \ell']) = f(\mathcal{S}', \mathcal{T}') - 2 - w(\{\ell, \ell'\}).$$

Moreover, by Lemma 5, for all $\ell \in \mathcal{L}(\mathcal{S}') - \{\rho'\}$, we have

$$f(\mathcal{S}'[-\ell], \mathcal{T}'[-\ell]) \geq f(\mathcal{S}', \mathcal{T}') - 1$$

if ℓ does not cross P and

$$f(\mathcal{S}'[-\ell, \ell'], \mathcal{T}'[-\ell, \ell']) \geq f(\mathcal{S}', \mathcal{T}') - 2 - w(\{\ell, \ell'\})$$

if ℓ crosses P with $\{\ell, \ell'\} \in P$. It now follows by the induction assumption and Lemma 5 that if $f(\mathcal{S}', \mathcal{T}') \geq k$, then the first call to the **for each** loop correctly returns $k + \text{INTERLEAVE}(\mathcal{S}'', \mathcal{T}'', w'', 0)$. Furthermore, by the induction assumption and Lemma 6, if $f(\mathcal{S}', \mathcal{T}') < k$, then the first call to the **for each** loop correctly returns $h + \text{INTERLEAVE}(\mathcal{S}'', \mathcal{T}'', w'', k - h)$, where $h = f(\mathcal{S}', \mathcal{T}')$. This completes the proof of the theorem. \square

We end this section by analyzing the running time of `HYBRIDINTERLEAVE` and comparing it with the time complexity of a previous implemented algorithm to solve `MINIMUM HYBRIDIZATION`.

Proposition 8. *Let \mathcal{S} and \mathcal{T} be two weighted rooted binary phylogenetic X -trees whose weighted 2-chain set P is empty. Furthermore, let k be an integer. Then the running time of `HYBRIDINTERLEAVE`($\mathcal{S}, \mathcal{T}, k$) is $O((14k)^k n^3)$, where $n = |X|$.*

Proof. By repeatedly applying the subtree and chain reductions to \mathcal{S} and \mathcal{T} until no further reduction is possible, it follows from Lemma 3 that the leaf set size of the obtained weighted rooted binary phylogenetic X' -trees \mathcal{S}' and \mathcal{T}' is at most $14h(\mathcal{S}, \mathcal{T})$.

1
2
3
4 Furthermore, while the subtree and chain reduction can be computed in $O(n^3)$ (Bor-
5 dewich and Semple, 2007b), a single application of the cluster reduction results in an
6 $O(2n)$ algorithm. Thus, calling all three reductions takes time $O(n^3)$.
7
8

9
10 Since HYBRIDINTERLEAVE directly returns k if $|X'| > 14k$, we may assume that
11 $|X'| \leq 14k$. The remaining part of this proof is by induction on k . If $k = 0$, then the
12 algorithm returns k in constant time. Now suppose that the running time of HYBRID-
13 INTERLEAVE is $O((14k')^{k'}n^3)$ for all $0 \leq k' < k$. Let A be a minimal common cluster
14 of \mathcal{S}' and \mathcal{T}' . As $14k \geq |A|$, the algorithm makes at most $14k$ calls to INTERLEAVE for
15 the tree pair $\mathcal{S}'|A$ and $\mathcal{T}'|A$ with parameter of at most $k - 1$. Thus the running time is
16 $O(n^3 + 14k(14(k - 1))^{k-1}n^3)$ which is $O((14k)^kn^3)$ as claimed. \square
17
18
19
20
21
22
23
24
25
26

27 **Remark.** The running time of HYBRIDINTERLEAVE given in Proposition 8 is an im-
28 provable upper bound since it only considers the kernelization of the two inputted phy-
29 logenetic trees and not applications of the subtree and chain reductions to previously
30 reduced trees for which INTERLEAVE is recursively called from within the **for each** loop.
31 As a comparison, $O((2 \cdot 14k)^k + n^3)$ is the theoretical worst-case running time of HY-
32 BRIDNUMBER (Bordewich and Semple, 2007b). Thinking of the exhaustive search part
33 of both algorithms HYBRIDINTERLEAVE and HYBRIDNUMBER as successively deleting
34 a set of edges to calculate the hybridization number, the difference in the two running
35 times is because the bounded search tree in HYBRIDINTERLEAVE is only based on the
36 deletion of pendant edges of the trees under consideration and not all of the edges as
37 in HYBRIDNUMBER. The fact that one needs only consider pendant edges follows from
38 Lemmas 5 and 6. An indication of how much better the theoretical worst-case running
39 time of HYBRIDINTERLEAVE could possibly be is highlighted in the next section when
40 we compare the running times of both algorithms on a biological data set.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

6 Experimental Results

To evaluate the performance of HYBRIDINTERLEAVE, we applied it to a grass (*Poaceae*) data set that has previously been used for running-time analyses in the context of calculating the hybridization number (Bordewich et al., 2007) and the rooted subtree prune and regraft distance (for details, see Section 7) which is frequently used to calculate the dissimilarities between two phylogenies for when reticulation is not assumed to be its major cause. The *Poaceae* data set was originally provided by the Grass Phylogeny Working Group (2001) and contains DNA sequences for six genetic loci, each with up to 65 taxa. Details about this data set and how a gene tree was reconstructed for each locus can be found in Bordewich et al. (2007) (and references therein). Species of the *Poaceae* family are subject to numerous natural hybridization events (Ellstrand et al., 1996). Therefore, the conflicting signals in this data set are more likely to be due to hybridization than to other processes causing inconsistencies.

For each of the 15 tree pairs, we restricted the two associated phylogenies to taxa that are common to both (second column of Table 1) and calculated the hybridization number of the resulting trees. The results are summarized in Table 1, where—beside the hybridization numbers—the running times for HYBRIDNUMBER and HYBRIDINTERLEAVE are compared for each tree pair. A detailed description of the former algorithm, is given by Bordewich et al. (2007). Note that we reran HYBRIDNUMBER to guarantee consistency among the obtained running times for both algorithms. While HYBRIDNUMBER computes the hybridization number for eight tree pairs within a couple of minutes, HYBRIDINTERLEAVE does so for all instances of the *Poaceae* data set and performs significantly faster. The latter algorithm successfully completes each program run in less than 8 minutes and calculates hybridization numbers as high as 19 for gene tree pairs with up to 46 taxa. This seems remarkably quick since HYBRIDNUMBER cannot calculate the exact solution for three tree pairs (*ndhF* and *ITS*, *rbcL* and *ITS*, and *rpoC2* and *ITS*) within 48 hours. The running time of HYBRIDINTERLEAVE mostly depends

1
2
3
4 on the exhaustive search part of this algorithm since the reductions can be computed
5
6 in polynomial-time. Clearly, the running time primarily decreases with an increase in
7
8 the number of taxa that can be reduced by any of the three reductions. On the other
9
10 hand, if the reductions have little effect because the trees only share a limited amount
11
12 of common features such as subtrees, chains, or clusters, then the running time greatly
13
14 increases with the hybridization number.
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

7 Concluding Remarks

In this paper, we presented the new algorithm HYBRIDINTERLEAVE that exactly calculates the hybridization number for two rooted binary phylogenetic trees. The algorithm can be applied to answer questions that consider the extent to which hybridization has influenced evolution and, therefore, shaped the current diversity of species. However, from a biological point of view, the results should carefully be interpreted since the algorithm is based on the assumption that hybridization is the only cause of gene tree inconsistencies. Moreover, it is possible that the real number of hybridization events for two trees is underestimated because HYBRIDINTERLEAVE minimizes this number and the true biological scenario might be less parsimonious. Of course, calculating a hybridization network that realizes this number is a desirable extension of our work. Indeed, given a maximum-acyclic-agreement forest, there is straightforward algorithm for doing this (Semple, 2007). However, such a network is typically not unique and it is part of ongoing research to implement an algorithm that calculates all possible hybridization networks that display two rooted binary phylogenetic trees and minimizes the hybridization number. Nevertheless, we believe that HYBRIDINTERLEAVE provides an important first step towards analyzing the occurrence of hybridization within a data set and, additionally, is remarkably quick.

We have shown that interleaving is an advantageous technique to speed-up the previously implemented fixed-parameter algorithm HYBRIDNUMBER. Referring back to the running time results summarized in Table 1, it is likely that HYBRIDINTERLEAVE can also compute the exact hybridization number in a reasonable short amount of time for problem instances that either contain bigger trees or have a greater hybridization number than those of the *Poaceae* data set. In conclusion, interleaving has proven to be most effective for our purpose of providing an exact algorithm to compute the hybridization number for two phylogenies of biologically relevant size, and we look forward to seeing whether interleaving has the same positive impact when applied to other fixed-parameter

tractable problems.

We end this paper with a remark on how interleaving can also be applied to calculate the rooted subtree prune and regraft (rSPR) distance. Loosely speaking, the graph-theoretic operation of rSPR cuts (prunes) a subtree and reattaches (regrafts) it to another part of the tree. The *rSPR distance* between two arbitrary rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} is the smallest number of rSPR operations that transforms \mathcal{S} into \mathcal{T} . We denote this distance by $d_{\text{rSPR}}(\mathcal{S}, \mathcal{T})$ and note that it is well-defined since one can always transform \mathcal{S} into \mathcal{T} via a sequence of single rSPR operations. Like MINIMUM HYBRIDIZATION, calculating $d_{\text{rSPR}}(\mathcal{S}, \mathcal{T})$ is NP-hard and fixed-parameter tractable (Bordewich and Semple, 2004). Furthermore, the following theorem was central to obtaining these results.

Theorem 9. *Let \mathcal{S} and \mathcal{T} be two rooted binary phylogenetic X -trees, and let $m(\mathcal{S}, \mathcal{T})$ denote the smallest number of elements among all agreement forests for \mathcal{S} and \mathcal{T} minus one. Then*

$$d_{\text{rSPR}}(\mathcal{S}, \mathcal{T}) = m(\mathcal{S}, \mathcal{T}).$$

Given the strong similarities between the characterizations of $h(\mathcal{S}, \mathcal{T})$ and $d_{\text{rSPR}}(\mathcal{S}, \mathcal{T})$ (see Theorems 1 and 9), it is not surprising that interleaving can also be applied to calculate the latter distance. However, while it is sufficient to exclusively consider 1-element subsets in the **for each** loop of HYBRIDINTERLEAVE, for calculating the rSPR distance, we need to iterate through all proper subsets of the label set under consideration, and thus, subsequently apply analogous subtree, chain, and cluster reductions to possibly more than one tree pair. This is due to the missing acyclic property in the context of calculating $d_{\text{rSPR}}(\mathcal{T}, \mathcal{T}')$. A detailed description of how interleaving can be applied in order to speed-up the computation of $d_{\text{rSPR}}(\mathcal{T}, \mathcal{T}')$ is given by Collins (2009).

8 Acknowledgement

We thank Michael Langston for introducing and highlighting to us the use of interleaving in fixed-parameter algorithms, and Daniel Huson for constructive comments on an earlier version of this paper. We thank the New Zealand Marsden Fund and the Department of Mathematics and Statistics at the University of Canterbury for their financial support. This research was partially supported by NSF grants SEI-BIO 0513910 and IIS-0803564.

9 Disclosure Statement

No competing financial interests exist.

REFERENCES

- 1
2
3
4
5
6
7 Abu-Khzam, F.N., Langston, M.A, Shanbhag, P., and Symons, C.T. 2006. Scalable
8 parallel algorithms for FPT problems. *Algorithmica* 45, 269–284.
- 9
10
11 Ávila, L.F., García, G., Serna, M., and Thilikos, D.M. 2006. A list of parameterized
12 problems in bioinformatics, *Technical report LSI-06-24-R*, Technical University of Cat-
13 alonia.
- 14
15
16
17
18
19 Baroni, M., Grünewald, S., Moulton, V., and Semple, C. 2005. Bounding the number of
20 hybridisation events for a consistent evolutionary history. *J. Math. Biol.* 51, 171–182.
- 21
22
23 Baroni, M., Semple, C., and Steel, M. 2006. Hybrids in real time. *Syst. Biol.* 55, 46–56.
- 24
25
26
27 Bordewich, M. and Semple, C. 2004. On the computational complexity of the rooted
28 subtree prune and regraft distance. *Ann. Combin.* 8, 409–423.
- 29
30
31 Bordewich, M. and Semple, C. 2007a. Computing the minimum number of hybridization
32 events for a consistent evolutionary history. *Discrete Appl. Math.* 155, 914–928.
- 33
34
35
36 Bordewich, M. and Semple, C. 2007b. Computing the hybridization number of two phy-
37 logenetic trees is fixed-parameter tractable. *IEEE Trans. Comput. Biol Bioinf.* 4, 458–
38 466.
- 39
40
41
42
43 Bordewich, M., Linz, S., St. John, K., and Semple, C. 2007. A reduction algorithm for
44 computing the hybridization number of two trees. *Evol. Bioinform.* 3, 86–98.
- 45
46
47
48
49 Chor, B. and Tuller, T. 2005. Maximum likelihood of evolutionary trees is hard, 296–310.
50 *In* Miyano, S., Mesirov, J., Kasif, S., Istrail, S., Pevzner, P., and Waterman, M., eds.,
51 *9th Annual International Conference, RECOMB 2005, Lecture Notes in Computer*
52 *Science*, Springer, Berlin.
- 53
54
55
56
57 Collins, J. 2009. *Rekernelisation algorithms in hybrid phylogenies*. MSc Thesis, University
58 of Canterbury, Christchurch, New Zealand.
- 59
60

- 1
2
3
4 Dehne, F., Langston, M.A., Luo, X., Pitre, S., Shaw, P., and Zhang, Y. 2006. The
5 cluster editing problem: Implementations and experiments, 13–24. *In* Bodlaender, H.
6 L. and Langston, M. A., eds., *Parameterized and Exact Computation, Lecture Notes*
7 *in Computer Science*, Springer, Berlin.
- 8
9
10
11
12 Downey, R. and Fellows, M. 1998. *Parameterized Complexity (Monographs in Computer*
13 *Science)*. Springer, Berlin.
- 14
15
16
17
18 Ellstrand, N.C., Whitkus, R., and Rieseberg, L.H. 1996. Distribution of spontaneous
19 plant hybrids. *P. Natl. Acad. Sci. USA* 93, 5090-5093.
- 20
21
22
23 Flum, J. and Grohe, G. 2006. *Parameterized Complexity Theory*. Springer, Berlin.
- 24
25
26
27 Foulds, L.R. and Graham, R.L. 1982. The Steiner problem in phylogeny is NP-complete.
28 *Adv. Appl. Math.* 3, 43–49.
- 29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
- Gramm, J. and Niedermeier, R. 2002. Breakpoint medians and breakpoint phylogenies:
A fixed-parameter approach. *Bioinformatics* 18, S128–S139.
- Gramm, J. and Niedermeier, R. 2003. A fixed-parameter algorithm for minimum quartet
inconsistency. *J. Comput. Syst. Sci.* 67, 723–741.
- Gramm, J., Nickelsen, A., and Tantau, T. 2008. Fixed-parameter algorithms in phylo-
genetics, *The Computer Journal* 51, 79–101.
- Grass Phylogeny Working Group. 2001. Phylogeny and subfamilial classification of the
grasses (*Poaceae*). *Ann. Mo. Bot. Gard.* 88, 373–457.
- Hein, J., Jing, T., Wang, L., and Zhang, K. 1996. On the complexity of comparing
evolutionary trees. *Discrete Appl. Math.* 71, 153–169.
- Linz, S. 2008. *Reticulation in evolution*. PhD Thesis, Heinrich-Heine-University,
Düsseldorf, Germany.
- Niedermeier, R. and Rossmanith, P. 2000. A general method to speed up fixed-parameter-
tractable algorithms. *Inform. Process. Lett.* 73, 125–129.

1
2
3
4 Roch, S. 2006. A short proof that phylogenetic tree reconstruction by maximum likeli-
5 hood is hard. *IEEE Trans. Comput. Biol Bioinf.* 3, 1545–5963.
6
7

8
9 Semple, C. 2007. Hybridization networks, 277–314. In Gascuel, O. and Steel, M., eds.,
10 *Reconstructing Evolution: New Mathematical and Computational Advances*, Oxford
11 University Press.
12
13
14

15
16 Semple, C. and M. Steel. 2003. *Phylogenetics*. Oxford University Press.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Table 1: Running time comparison of HYBRIDINTERLEAVE with HYBRIDNUMBER (Bordewich et al., 2007) for the *Poaceae* data set (Grass Phylogeny Working Group, 2001).

Pairwise combination	#Taxa	Hybridization number	RT ^a of HYBRIDNUMBER	RT ^a of HYBRIDINTERLEAVE
<i>ndhF</i> <i>phyB</i>	40	14	5.9 h	23 s
<i>ndhF</i> <i>rbcL</i>	36	13	5.3 h	3 s
<i>ndhF</i> <i>rpoC2</i>	34	12	13 h	6 s
<i>ndhF</i> <i>waxy</i>	19	9	150 s	< 1 s
<i>ndhF</i> <i>ITS</i>	46	19	> 48 h	258 s
<i>phyB</i> <i>rbcL</i>	21	4	< 1 s	< 1 s
<i>phyB</i> <i>rpoC2</i>	21	7	90 s	< 1 s
<i>phyB</i> <i>waxy</i>	14	3	< 1 s	< 1 s
<i>phyB</i> <i>ITS</i>	30	8	10 s	< 1 s
<i>rbcL</i> <i>rpoC2</i>	26	13	15.2 h	8 s
<i>rbcL</i> <i>waxy</i>	12	7	132 s	< 1 s
<i>rbcL</i> <i>ITS</i>	29	14	> 48 h	612 s
<i>rpoC2</i> <i>waxy</i>	10	1	< 1 s	< 1 s
<i>rpoC2</i> <i>ITS</i>	31	15	> 48 h	57 s
<i>waxy</i> <i>ITS</i>	15	8	330 s	< 1 s

^aRunning time (RT) on a 2.66 GHz CPU, 2 GB RAM machine measured in seconds (s) and hours (h), respectively.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 1: A hybridization network \mathcal{H} with $h(\mathcal{H}) = 3$, and a rooted binary phylogenetic tree \mathcal{T} that is displayed by \mathcal{H} (indicated by the thicker arcs in \mathcal{H}).

For Peer Review

Figure 2: Two rooted binary phylogenetic X -trees \mathcal{S} and \mathcal{T} with their roots labeled ρ .

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 3: Two agreement forests \mathcal{F} and \mathcal{F}' and their associated digraphs $G_{\mathcal{F}}$ and $G_{\mathcal{F}'}$, respectively, for \mathcal{S} and \mathcal{T} shown in Figure 2.

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 4: Two rooted binary phylogenetic X -trees \mathcal{S}' and \mathcal{T}' that have been obtained from \mathcal{S} and \mathcal{T} depicted in Figure 2 by repeated applications of the chain reduction.

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 5: A cluster reduction applied to the two rooted binary phylogenetic trees \mathcal{S} and \mathcal{T} , where \mathcal{S}_a and \mathcal{T}_a have been obtained from \mathcal{S} and \mathcal{T} , respectively, by replacing the pendant subtree whose label set is A with a new leaf labeled a .

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12 Figure 6: The intermediate tree pairs of a computational path in the search tree generated
13 by a call to $\text{HYBRIDINTERLEAVE}(\mathcal{S}_1, \mathcal{T}_1, k)$ with $k > 2$. Note that the cluster reduction
14 is omitted here since the second tree pair which is returned from this reduction always
15 consists of two identical trees. For details, see text.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

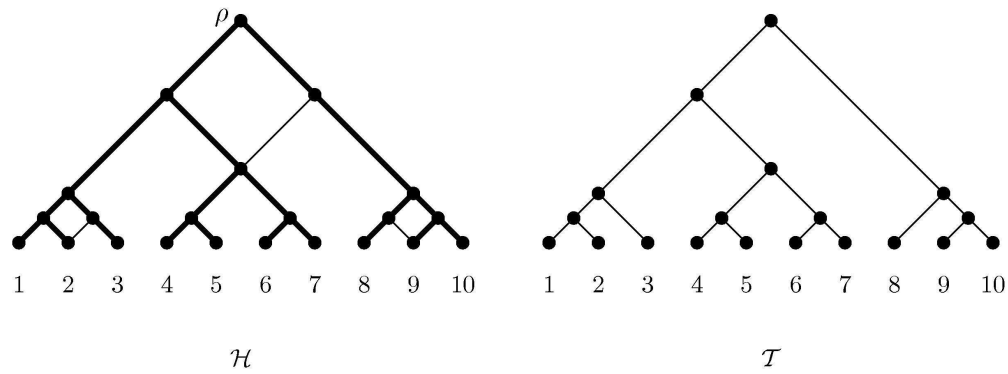


Figure 1
102x37mm (600 x 600 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

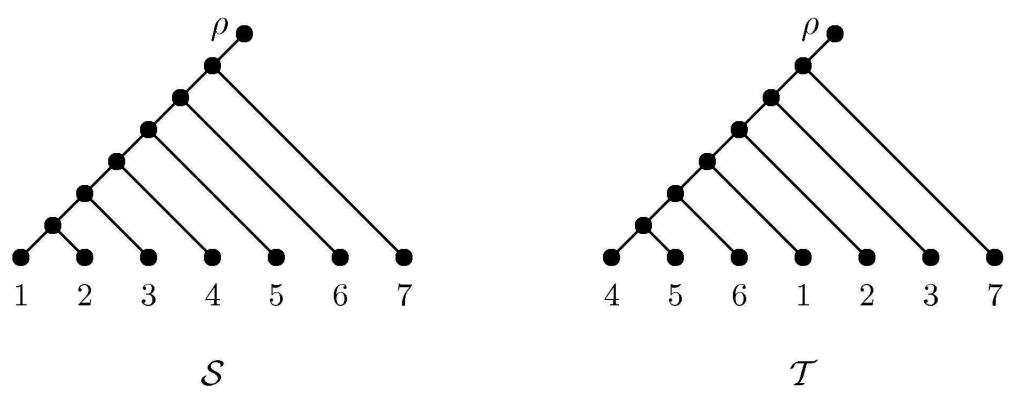


Figure 2
79x29mm (600 x 600 DPI)

Peer Review

$$\mathcal{F} = \{\{\rho, 7\}, \{1, 2, 3\}, \{4, 5, 6\}\} \quad \mathcal{F}' = \{\{\rho, 1, 2, 3, 7\}, \{4\}, \{5\}, \{6\}\}$$

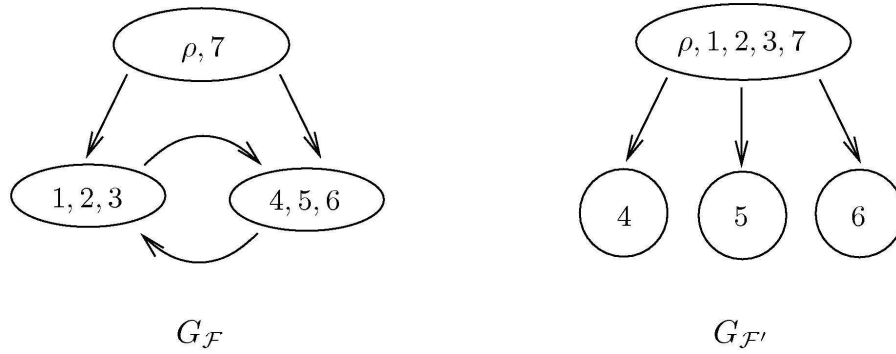


Figure 3
86x40mm (600 x 600 DPI)

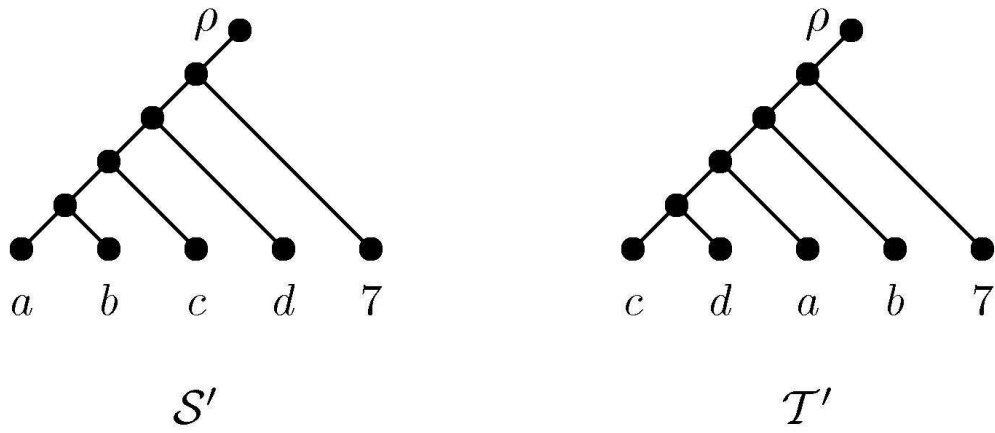


Figure 4
58x24mm (600 x 600 DPI)

Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

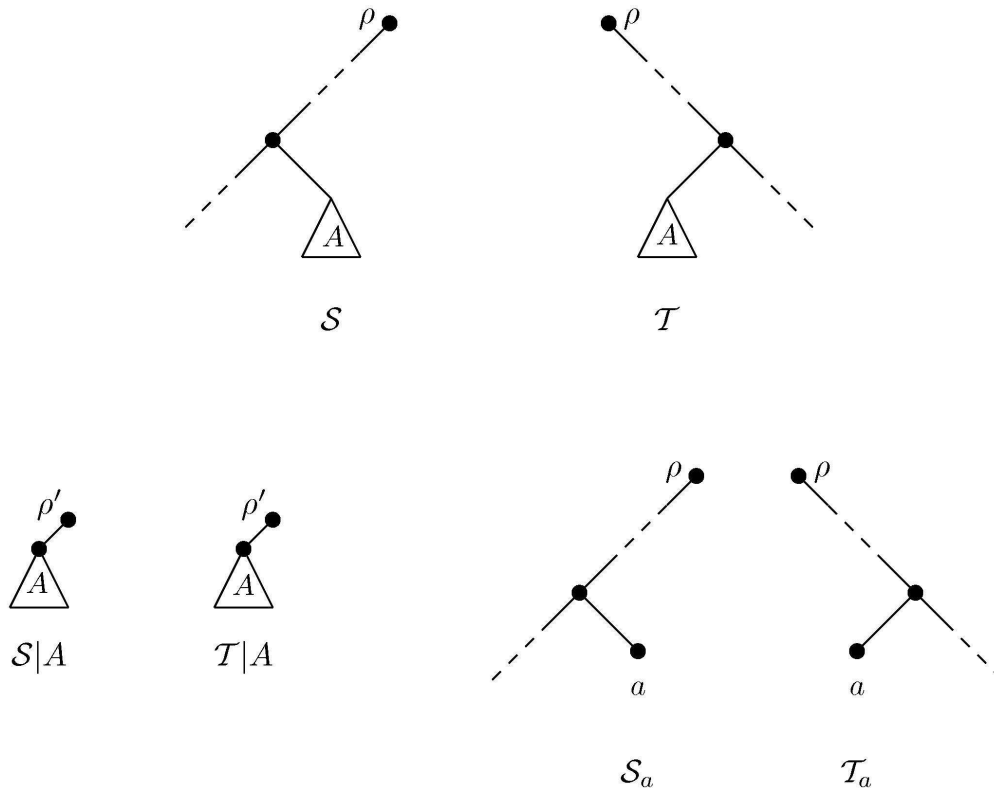


Figure 5
87x68mm (600 x 600 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

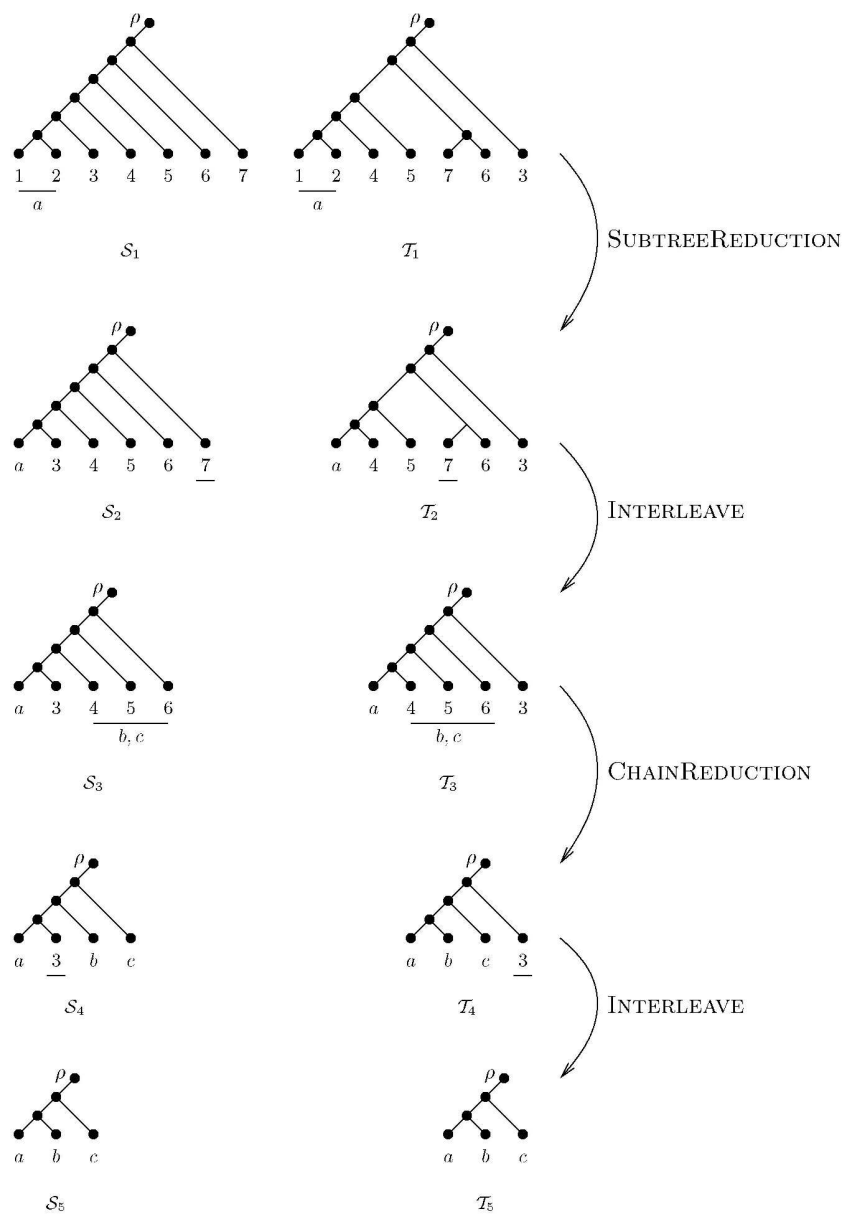


Figure 6
113x163mm (600 x 600 DPI)