



HAL
open science

Fixed-parameter Tractability of the Maximum Agreement Supertree Problem

Sylvain Guillemot, Vincent Berry

► **To cite this version:**

Sylvain Guillemot, Vincent Berry. Fixed-parameter Tractability of the Maximum Agreement Supertree Problem. RR-07005, 2007, pp.17. lirmm-00130406v1

HAL Id: lirmm-00130406

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00130406v1>

Submitted on 27 Feb 2007 (v1), last revised 27 Feb 2007 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Fixed-Parameter Tractability of the Maximum Agreement Supertree Problem^{*}

EXTENDED ABSTRACT

Sylvain Guillemot¹ and Vincent Berry¹

Equipe *Méthodes et Algorithmes pour la Bioinformatique*, LIRMM, C.N.R.S.-Université Montpellier II.
{sguillem,vberry}@lirmm.fr

Abstract. Given a ground set L of labels and a collection of trees whose leaves are bijectively labelled by some elements of L , the Maximum Agreement Supertree problem (SMAST) is the following: find a tree T on a largest label set $L' \subseteq L$ that homeomorphically contains every input tree restricted to L' . The problem finds applications in phylogenetics, databases and data mining. In this paper we focus on the parameterized complexity of this NP-hard problem. We consider different combinations of parameters for SMAST as well as particular cases, providing both FPT algorithms and intractability results.

1 Introduction

Motivation. Supertree construction consists in building trees on a large set of labels from smaller trees covering parts of the label set. This task finds application in bioinformatics where trees represent phylogenies, but also in other fields [1, 2]. In phylogenetics, the labels are bijectively associated with the leaves of the trees and represent current organisms, while internal nodes represent hypothetical ancestors. The topological information in the input trees consists in the groupings of labels induced by internal nodes, representing related sets of organisms such as species, orders, families, etc. The goal is to build a supertree complying as much as possible with the topological information of the source trees. The task is relatively easy when the input trees agree on the relative positions of the labels. In this case, it is possible to find in polynomial time a supertree that contains any input tree as an induced subtree, hence that incorporates all topological information provided by the data [1]. However, in practice several input trees usually disagree on the position of some leaves with respect to other leaves.

Related work. Some methods aim at producing supertrees incorporating as much input information as possible under the constraint that they do not contradict any input tree: they avoid disagreements between the input trees by collapsing some of their edges [3, 4] or by excluding some of their leaves, i.e. labels. The Maximum Agreement Supertree (SMAST) method [5–7] is appared to the latter kind. Given a collection \mathcal{T} of k

^{*} This paper was supported by the *Action incitative* BIOSTIC-LR.

trees of maximum degree d with labels taken in a ground set L of size n , an agreement supertree for \mathcal{T} is a tree T on a subset $L' \subseteq L$ such that each tree of \mathcal{T} restricted to L' is included in T . The SMAST problem consists in finding an agreement supertree containing the maximum number of labels from L .

This problem is NP-hard in general as it generalizes the MAST problem [8]. SMAST remains NP-hard when d is unrestricted for $k \geq 3$ input trees [6] and for trees of degree $d \geq 2$ when k is unrestricted [5]. Moreover, [6, 5] have also considered the complement problem, which is a minimization problem where the measure is the number p of labels missing in an agreement supertree. This problem can not be approximated in polynomial time within a constant factor, unless $P = NP$ [5]. The corresponding decision problem parameterized in p is W[2]-hard [5].

For the particular case of $d = 2$, [6] gave an $O(n^{k^2})$ algorithm for SMAST. For $k = 2$ both [6, 5] shown that SMAST can be solved in polynomial time, by reduction to MAST.

Our results. In this paper, we focus on the particular case where $d = 2$. Note that in phylogenetics, the input trees of SMAST will often be binary as a result of the optimization algorithms used to analyze raw molecular data. We improve on previous results in several ways.

First, we show that SMAST on k rooted binary trees on a label set of size n can be solved in $O((2k)^p kn^2)$. This algorithm is only exponential in p , that roughly represents the extent to which the input trees disagree. Thus, the algorithm will be reasonably fast when dealing with collections of trees obtained for genes displaying a low level of homoplasy. Then, we provide an $O((8n)^k)$ algorithm, independent of p , and significantly improving on the $O(n^{k^2})$ algorithm of [6]. This algorithm shows that SMAST is tractable for a small number of trees, extending in some sense the previously known results for $k = 2$ trees [6, 5, 9]. We also obtain some fixed-parameter intractability results for various combinations of parameters of SMAST.

We then consider SMAST on collections of rooted triples (binary trees on 3 leaves), focusing on the complexity of this variant parameterized in p . Since this problem is equivalent to SMAST in its general setting [9], it is W[2]-hard. However, we show here that an FPT algorithm can be achieved for *complete* collections of rooted triples, i.e., when there is at least one rooted triple for each set of 3 labels in L . This results from the fact that conflicts between the input trees can be circumvented to small sets of labels, leading to $O(4^p n^3)$ and $O(3 \cdot 12^p + n^4)$ algorithms.

2 Definitions

We consider rooted trees which are bijectively leaf-labelled. Let T be such a tree, we identify its leaf set with its label set, denoted by $L(T)$. The *size* of T is $|T| := |L(T)|$.

The node set of T is denoted by $N(T)$, and $r(T)$ stands for the root of T . We use a parenthesized notation for trees: if x is a label, then x denotes the tree whose root is a leaf labelled by x ; if T_1, \dots, T_k are trees, then (T_1, \dots, T_k) stands for the tree whose root is connected to the child subtrees T_1, \dots, T_k .

If x is a node of T , $T(x)$ stands for the subtree of T rooted at x , and $L(x)$ for the label set of this subtree. If x, y are two nodes of T , then $x <_T y$ means that x is a descendant of y in T . The upper bound of two nodes x, y of T w.r.t. $<_T$ is called the lowest common ancestor of x, y , and is denoted by $\text{lca}_T(x, y)$. If x, y are two nodes of T s.t. $x <_T y$, denote by $\text{child}_T(x, y)$ the child of y along the path joining y to x in T . If x is an internal node of T , the set of children of x in T is denoted by $\text{children}_T(x)$.

Given a tree T and a label set L , the *restriction* of T to L , denoted by $T|L$, is the tree homeomorphic to the smallest subtree of T connecting leaves of L . Let T, T' be two trees. We say that T *embeds* in T' , denoted by $T \leq T'$, iff $T = T'|L(T)$. We say that T *partially embeds* in T' , denoted by $T \bowtie T'$, iff $T|L(T') = T'|L(T)$. A *collection* is a family $\mathcal{T} = \{T_1, \dots, T_k\}$ of trees, the label set of the collection is $L(\mathcal{T}) = \cup_{i=1}^k L(T_i)$. Given a label set L , the *restriction* of \mathcal{T} to L is the collection $\mathcal{T}|L = \{T_1|L, \dots, T_k|L\}$. See Figure 1 for an example of a collection.

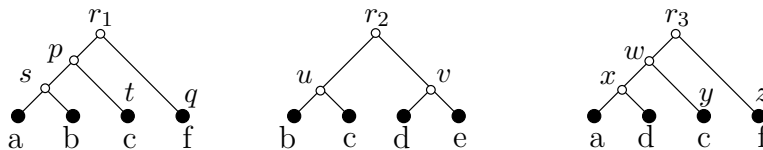


Fig. 1. A collection \mathcal{T} of 3 input trees.

An *agreement supertree* for \mathcal{T} is a tree S s.t. $L(S) \subseteq L(\mathcal{T})$ and for each $i \in [k]$, $S \bowtie T_i$. We say that S is a *total agreement supertree* for \mathcal{T} if additionally $L(S) = L(\mathcal{T})$. The collection \mathcal{T} is *compatible* iff there exists a total agreement supertree for \mathcal{T} . A *conflict* between \mathcal{T} is a set $C \subseteq L(\mathcal{T})$ s.t. $\mathcal{T}|C$ is incompatible. For instance, $T = (((a, b), c), (e, f))$ is an agreement supertree for the collection \mathcal{T} of Figure 1, and $C = \{a, b, c, d\}$ is a conflict between \mathcal{T} .

Given a collection \mathcal{T} , we define $\text{SMAST}(\mathcal{T})$ as the set of agreement supertrees for \mathcal{T} . The MAXIMUM AGREEMENT SUPERTREE problem (SMAST) asks: given a collection \mathcal{T} , find an agreement supertree for \mathcal{T} with the largest size. Equivalently, it amounts to seek a largest set $L \subseteq L(\mathcal{T})$ s.t. $\mathcal{T}|L$ is compatible. The size of such an optimal solution is denoted by $\#\text{SMAST}(\mathcal{T})$. We also denote by P-SMAST the parameterized version of SMAST, which asks: given a collection \mathcal{T} and a parameter p , can \mathcal{T} be made compatible by removing at most p labels?

3 Solving Smast on binary trees

Throughout this section, we consider a fixed collection $\mathcal{T} = \{T_1, \dots, T_k\}$ of binary trees, we let n denote the size of the label set and k the number of trees.

If T is a tree, we define $N^\perp(T) := N(T) \cup \{\perp\}$. We extend the notation $T(u)$ to $u \in N^\perp(T)$, s.t. if $u = \perp$ then $T(u)$ is the empty tree. We extend the relation \leq_T to $N^\perp(T)$ s.t. $\perp \leq_T x$ for each $x \in N^\perp(T)$.

A *position* in \mathcal{T} is a tuple $\pi = (u_1, \dots, u_k)$, where each $u_i \in N^\perp(T_i)$. For $i \in [k]$, the i th component of π is denoted $\pi[i]$. We set $I(\pi) = \{i \in [k] : \pi[i] \text{ is an internal node of } T_i\}$. We define the *initial position* $\pi_\top = (r_1, \dots, r_k)$, where each r_i is the root of T_i . We define the *final position* $\pi_\perp = (\perp, \dots, \perp)$. We let $\Pi(\mathcal{T})$ denote the set of positions in \mathcal{T} .

3.1 Solving Smast in $O((2k)^p \times kn^2)$ time

In this section, we describe an algorithm deciding the compatibility of a collection in $O(kn^2)$ time, and returning a conflict of size $\leq 2k$ in case of incompatibility. This yields an FPT algorithm for P-SMAST with $O((2k)^p \times kn^2)$ running time.

The well-known BUILD algorithm [1, 10] decides the compatibility of a collection but doesn't give a conflict in case of incompatibility. Like BUILD, the algorithm presented here builds the supertree using a recursive top-down approach. Each step constructs a graph where the connected components correspond to the subtrees hanging from the root of the supertree. However, we replace the graph used in BUILD with a graph that when connected yields a conflict of size $\leq 2k$, identified thanks to a spanning tree.

We begin with some additional definitions. A position is *reduced* iff each component is either \perp or an internal node; to any position π , we associate a reduced position $\pi \downarrow$ by replacing by \perp any component of π that is a leaf. In the following, we will assume that π is a reduced position in \mathcal{T} . We set $\mathcal{T}(\pi) := \{T_1(u_1), \dots, T_k(u_k)\}$.

We define the graph $G(\mathcal{T}, \pi)$ as follows: (i) its vertex set is $V = \cup_{i \in I(\pi)} \text{children}_{T_i}(\pi[i])$; (ii) two vertices $x, y \in V$ are adjacent iff $L(x) \cap L(y) \neq \emptyset$. In other terms, $G(\mathcal{T}, \pi)$ is the intersection graph of the set system $\{L(x) : x \in V\}$. See Figure 2 for an example of such graphs.

If $V' \subseteq V$, we define the position $\text{Succ}_{V'}(\pi)$ as the position π' s.t.

- if $\pi[i] = \perp$, then $\pi'[i] = \perp$;
- if $\pi[i]$ is an internal node of T_i , with children v_i, v'_i , then one of the following holds:
 - (i) $\pi'[i] = v_i$ if $v_i \in V', v'_i \notin V'$, (ii) $\pi'[i] = v'_i$ if $v_i \notin V', v'_i \in V'$, (iii) $\pi'[i] = u_i$ if $v_i \in V', v'_i \in V'$, (iv) $\pi'[i] = \perp$ if $v_i \notin V', v'_i \notin V'$.

We set $\text{succ}_{V'}(\pi) = \text{Succ}_{V'}(\pi) \downarrow$. Given $V' \subseteq V$, set $L(V') = \cup_{x \in V'} L(x)$. Set $L(\pi) = L(\mathcal{T}(\pi))$. We will repeatedly use the following simple observations:

Lemma 1. *Let $V' \subseteq V$. Then: $L(\text{Succ}_{V'}(\pi)) \subseteq L(V')$.*

Lemma 2. *Two sets $V_1, V_2 \subseteq V$ are connected in $G(\mathcal{T}, \pi)$ iff $L(V_1) \cap L(V_2) \neq \emptyset$.*

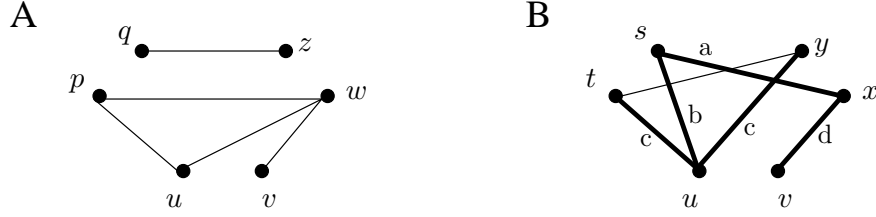


Fig. 2. A. The graph $G(\mathcal{T}, \pi_{\top})$ of the position $\pi_{\top} = (r_1, r_2, r_3)$ for the collection of trees of Figure 1. This graph is disconnected, the two connected components indicate the two successor positions $\pi_1 = (p, r_2, w)$ and $\pi_2 = (q, \perp, z)$ of π_{\top} . B. The graph $G(\mathcal{T}, \pi_1)$ is connected. Choosing a spanning tree (bold edges) of the graph and an arbitrary label shared by the two subtrees corresponding to the extremities of each edge of this tree identifies a conflict $C = \{a, b, c, d\}$.

Given a position π in \mathcal{T} , we say that π is *compatible* iff $\mathcal{T}(\pi)$ is compatible. Note that π_{\perp} is compatible. Observe also that: π is compatible iff $\pi \downarrow$ is compatible. Let us now consider a reduced position $\pi \neq \pi_{\perp}$. Then $|L(\pi)| \geq 2$. We have the following recursive characterization of compatibility for this case:

Lemma 3. *Suppose that $|L(\pi)| \geq 2$. The following are equivalent:*

- π is compatible;
- there exists a partition V_1, V_2 of V s.t. (i) V_1, V_2 are disconnected in $G(\mathcal{T}, \pi)$, (ii) $\text{succ}_{V_1}(\pi), \text{succ}_{V_2}(\pi)$ are compatible.

Moreover, if the graph $G(\mathcal{T}, \pi)$ turns out to be connected, a spanning tree of this graph yields a small conflict between \mathcal{T} :

Lemma 4. *Suppose that $G(\mathcal{T}, \pi)$ is connected, and let $T = (V, F)$ be a spanning tree of $G(\mathcal{T}, \pi)$. For each edge $e = \{u, v\} \in F$, choose $l_e \in L(u) \cap L(v)$. Then $C = \{l_e : e \in F\}$ is a conflict between \mathcal{T} .*

Lemmas 3 and 4 give rise to an algorithm for deciding the compatibility of a collection, and obtaining a conflict of small size in case of incompatibility:

Theorem 1. *There is an algorithm which, in $O(kn^2)$ time, decides if \mathcal{T} is compatible, or returns a conflict of size $\leq 2k$.*

Proof. We define a procedure $\text{ISCOMPATIBLE}(\pi)$ which takes as input a reduced position, decides if π is compatible, or returns a conflict of size $\leq 2k$ in case of incompatibility. The procedure is as follows: (i) answer ("yes") if $\pi = \pi_{\perp}$; (ii) if $\pi \neq \pi_{\perp}$, test whether $G(\mathcal{T}, \pi)$ is connected.

- If the graph is connected, then let $T = (V, F)$ be a spanning tree of $G(\mathcal{T}, \pi)$, choose $l_e \in L(u) \cap L(v)$ for each edge $e = \{u, v\} \in F$, construct $C = \{l_e : e \in F\}$, and return ("no", C).

- If the graph is not connected, then let V_1, V_2 be a partition of V in two disconnected sets, and construct the positions π_1, π_2 where $\pi_i = \text{succ}_{V_i}(\pi)$. Call $\text{ISCOMPATIBLE}(\pi_1)$, let R_1 be its result; if $R_1 = (\text{"yes"})$ then call $\text{ISCOMPATIBLE}(\pi_2)$ and return its result R_2 , else return R_1 .

To decide if \mathcal{T} is compatible, we simply call $\text{ISCOMPATIBLE}(\pi_{\top} \downarrow)$.

We now justify the correctness and the running time of the algorithm. The correctness of the procedure ISCOMPATIBLE follows from lemmas 3 and 4. For the running time, we rely on the fact that using appropriate data structures, we can ensure that a call to ISCOMPATIBLE takes $O(kn)$ time (see Appendix for details). By lemmas 1 and 2, the total number of calls to ISCOMPATIBLE is $O(n)$, therefore the total running time of the algorithm is $O(kn^2)$. \square

The algorithm of Theorem 1 yields a simple FPT algorithm for P-SMAST using the bounded search tree technique:

Theorem 2. *The P-SMAST problem can be solved in $O((2k)^p \times kn^2)$ time.*

3.2 Solving Smast in $O((8n)^k)$ time

In this section, we describe an algorithm which solves SMAST in $O((8n)^k)$ time. The algorithm uses dynamic programming, and is somewhat similar in spirit to the algorithm described in [8] for solving MAST on two trees.

We first give an alternative definition of the \bowtie relation in terms of *partial embeddings*. Let T, T' be two trees, say that a partial embedding of T into T' is a function $\phi : N(T) \rightarrow N(T') \cup \{\perp\}$ such that:

- for any x leaf of T , we have $\phi(x) = \perp$ if $x \notin L(T')$, or $\phi(x) = x$ otherwise,
- for any x internal node of T with children u_1, \dots, u_p , let $V = \{j : \phi(u_j) \neq \perp\}$, then (i) either $V = \emptyset$, and $\phi(x) = \perp$, (ii) either $V = \{i\}$ and $\phi(x) = \phi(u_i)$, (ii) or $|V| \geq 2$ and $\phi(x) >_T \phi(u_i)$ for each $i \in V$, and the nodes $\{\text{child}_T(\phi(u_i), \phi(x)) : i \in V\}$ are pairwise distinct.

Then: $T \bowtie T'$ iff there exists a partial embedding of T into T' (or equivalently, a partial embedding of T' into T).

Given a collection \mathcal{T} , the algorithm computes values $\#SMAST(\pi)$ for each position π . Let $SMAST(\pi)$ denote the set of trees T s.t. (i) T is an agreement supertree for \mathcal{T} , (ii) for each i , the partial embedding $\phi_i : T \rightarrow T_i$ is such that $\phi_i(r(T)) \leq_{T_i} \pi[i]$. We denote by $\#SMAST(\pi)$ the size of a largest tree of $SMAST(\pi)$.

We now define two values $\#SMAST_1(\pi)$ and $\#SMAST_2(\pi)$, from which $\#SMAST(\pi)$ is computed. We first define $\#SMAST_1(\pi)$. Say that a position π' is a *successor* of π

iff there exists $i \in [k]$ s.t. $\pi'[i]$ is a child of $\pi[i]$ and $\pi'[j] = \pi[j]$ for each $j \neq i$. Let $S(\pi)$ denote the set of successors of π . Then:

$$\#SMAST_1(\pi) = \max_{\pi' \in S(\pi)} \#SMAST(\pi'). \quad (1)$$

We now define $\#SMAST_2(\pi)$. Say that a pair of positions (π_1, π_2) is a *decomposition* of π iff (i) $\pi_1 \neq \pi, \pi_2 \neq \pi$, (ii) for each $i \in [k]$, the following holds:

- either $\pi[i] = \perp$, in which case $\pi_1[i] = \pi_2[i] = \perp$;
- either $\pi[i]$ is a leaf x , in which case we have $\{\pi_1[i], \pi_2[i]\} = \{\perp, x\}$;
- either $\pi[i]$ is an internal node u with two children v, v' , in which case we have either $\{\pi_1[i], \pi_2[i]\} = \{\perp, x\}$ or $\{\pi_1[i], \pi_2[i]\} = \{v, v'\}$.

Let $D(\pi)$ denote the set of decompositions of π . Then:

$$\#SMAST_2(\pi) = \max_{(\pi_1, \pi_2) \in D(\pi)} (\#SMAST(\pi_1) + \#SMAST(\pi_2)). \quad (2)$$

We define the relation $\leq_{\mathcal{T}}$ on $\Pi(\mathcal{T})$ by: $\pi \leq_{\mathcal{T}} \pi'$ iff for each $i \in [k]$, $\pi[i] \leq_{T_i} \pi'[i]$. We observe that:

Lemma 5. (i) If $\pi' \in S(\pi)$, then $\pi' <_{\mathcal{T}} \pi$, (ii) if $(\pi_1, \pi_2) \in D(\pi)$ then $\pi_i <_{\mathcal{T}} \pi$ for $i \in \{1, 2\}$.

Lemma 6. If $\pi' \leq_{\mathcal{T}} \pi$, then $SMAST(\pi') \subseteq SMAST(\pi)$.

We now give a recurrence relation for computing $\#SMAST(\pi)$. Say that a position π is *terminal* if for each $i \in [k]$, $\pi[i]$ is a leaf or \perp .

We first consider terminal positions. Given $x \in L(\mathcal{T})$, let $Ind(x) = \{i \in [k] : x \in L(T_i)\}$. Given a terminal position π , and given $x \in L(\pi)$, define $Ind(x, \pi) = \{i \in [k] : \pi[i] = x\}$; observe that $Ind(x, \pi) \subseteq Ind(x)$. Say that an element $x \in L(\pi)$ is *nice* iff $Ind(x, \pi) = Ind(x)$, and let $Nice(\pi)$ denote the set of nice elements of $L(\pi)$. Then:

Lemma 7. Suppose that π is terminal. Then: $\#SMAST(\pi) = |Nice(\pi)|$.

Lemma 8. Suppose that π is not terminal. Then:

$$\#SMAST(\pi) = \max(\#SMAST_1(\pi), \#SMAST_2(\pi)).$$

Proof. We first prove that $\#SMAST_1(\pi) \leq \#SMAST(\pi)$. Let $S \in SMAST(\pi')$ for some $\pi' \in S(\pi)$, s.t. $|S|$ is maximal. Since $\pi' <_{\mathcal{T}} \pi$ by Lemma 5, we have $S \in SMAST(\pi)$ by Lemma 6, and the result follows.

We now prove that $\#SMAST_2(\pi) \leq \#SMAST(\pi)$. Let $(\pi_1, \pi_2) \in D(\pi)$, and let S_1, S_2 s.t. $S_j \in SMAST(\pi_j)$, $|S_j|$ maximal. If one of the S_j 's is empty, say S_1 , then $\#SMAST(\pi_1) = 0$, and we obtain $\#SMAST_2(\pi) = |S_2| = \#SMAST(\pi_2) \leq \#SMAST(\pi)$

by lemmas 5 and 6. Suppose now that S_1, S_2 are not empty. For $j \in \{1, 2\}$, since $S_j \in \text{SMAS}T(\pi_j)$, there exists partial embeddings $\phi_{j,i} : S_j \rightarrow T_i$ s.t. $\phi_{j,i}(r(S_j)) \leq_{T_i} \pi_j[i]$ for each $i \in [k]$. Let $S = (S_1, S_2)$, we claim that $S \in \text{SMAS}T(\pi)$. Indeed, define $\phi_i : S \rightarrow T_i$ as follows. Set $\phi_i(x) = \phi_{j,i}(x)$ if x is a node of S_j , and $\phi_i(x) = \text{lca}_{T_i}(\phi_{1,i}(r(S_1)), \phi_{2,i}(r(S_2)))$ if x is the root of S . It is clear that: (i) $L(S_1) \cap L(S_2) = \emptyset$, hence S is well-defined, (ii) ϕ_i is a partial embedding of S into T_i , (iii) $\phi_i(r(S)) \leq_{T_i} \pi[i]$ (proof in Appendix). We conclude that $\#\text{SMAS}T_2(\pi) = |S_1| + |S_2| = |S| \leq \#\text{SMAS}T(\pi)$.

Finally, we show that $\#\text{SMAS}T(\pi) \leq \max(\#\text{SMAS}T_1(\pi), \#\text{SMAS}T_2(\pi))$. Let $S \in \text{SMAS}T(\pi)$ s.t. $|S|$ is maximal. Then there exists partial embeddings $\phi_i : S \rightarrow T_i$ s.t. $\phi_i(r(S)) \leq_{T_i} \pi[i]$ for each $i \in [k]$. Let $u_i = \phi_i(r(S))$ for each i . We consider two cases.

First case: there exists $i \in [k]$ s.t. $u_i <_{T_i} \pi[i]$. This case holds in particular if $|S| \leq 1$. Define π' from π by setting the i th component to $\text{child}_{T_i}(u_i, \pi[i])$, then $\pi' \in S(\pi)$. We verify that $S \in \text{SMAS}T(\pi')$: indeed, ϕ_i is a partial embedding of S into T_i s.t. $\phi_i(r(S)) \leq_{T_j} \pi'[j]$ for each j . We conclude that $|S| = \#\text{SMAS}T(\pi) \leq \#\text{SMAS}T(\pi') \leq \#\text{SMAS}T_1(\pi)$.

Second case: $u_i = \pi[i]$ for each $i \in [k]$. In this case, we have $|S| \geq 2$, hence $S = (S_1, S_2)$. Let u be the root of S , let v_i be the root of S_i in S , then $\pi = (\phi_1(u), \dots, \phi_k(u))$. For $j \in \{1, 2\}$, define π_j as follows: given $i \in [k]$, (i) if $\phi_i(v_j) = \phi_i(u)$, set $\pi_j[i] = \phi_i(u)$, (ii) if $\phi_i(v_j) = \perp$, set $\pi_j[i] = \perp$, (iii) if $\phi_i(v_j) <_{T_i} \phi_i(u)$, set $\pi_j[i] = \text{child}_{T_i}(\phi_i(v_j), \phi_i(u))$. It is clear that $(\pi_1, \pi_2) \in D(\pi)$ (proof in Appendix). We now show that $S_j \in \text{SMAS}T(\pi_j)$: indeed, ϕ_i is a partial embedding of S_j into T_i , and by definition of π_j we have $\phi_i(r(S_j)) \leq_{T_i} \pi_j[i]$ for each $i \in [k]$. We conclude that $|S| = \#\text{SMAS}T(\pi) = |S_1| + |S_2| \leq \#\text{SMAS}T(\pi_1) + \#\text{SMAS}T(\pi_2) \leq \#\text{SMAS}T_2(\pi)$. \square

Lemmas 7 and 8 yield an algorithm for computing $\#\text{SMAS}T(\mathcal{T})$:

Theorem 3. $\#\text{SMAS}T(\mathcal{T})$ can be computed in $O((8n)^k)$ time.

Proof. Using dynamic programming, the algorithm computes the values $\#\text{SMAS}T(\pi)$ for each position π , using the recurrence relations stated in lemmas 7 and 8. The correctness of the algorithm follows from the lemmas, and the termination of the algorithm is ensured by Lemma 5 and the fact that $<_{\mathcal{T}}$ is an order relation on $\Pi(\mathcal{T})$.

We now consider the space and time requirements for the algorithm. First observe that the number of positions π in \mathcal{T} is $\leq (2n)^k$: indeed, a component $\pi[i]$ has $\leq 2n$ possible values (one of the $\leq 2n - 1$ nodes of T_i , or the value \perp). It follows that the space complexity is $O((2n)^k)$. We claim that the time complexity is $O((8n)^k)$. Indeed, consider the time required to compute $\#\text{SMAS}T(\pi)$, assuming that the values $\#\text{SMAS}T(\pi')$ for $\pi' <_{\mathcal{T}} \pi$ are available. Testing if π is terminal requires $O(k)$ time. If π is terminal, computing $|\text{Nice}(\pi)|$ takes $O(k)$ time. If π is nonterminal, then we need to compute $\#\text{SMAS}T_1(\pi)$ and $\#\text{SMAS}T_2(\pi)$, which respectively require $O(k)$ and $O(4^k)$ time. Thus, $\#\text{SMAS}T(\pi)$ is computed in $O(4^k)$ time, hence the total running time of the algorithm is $O((8n)^k)$. \square

3.3 Hardness results

The parameterized complexity of the SMAST problem on binary trees is considered w.r.t. the following parameters: k denotes the number of input trees, l denotes an upper bound on the maximum size of the input trees, p (resp. q) denotes an upper (resp. lower) bound on the number of labels to remove (resp. conserve) in order to obtain compatibility of the collection. Our complexity results for several combinations of the parameters are summarized in Theorem 4:

Theorem 4. *We have the following hardness results for SMAST:*

Parameters	Complexity of SMAST
q	W[1]-complete (even for $l = 3$)
q, k	W[1]-complete (even for $l = 3$)
p	W[2]-hard (even for $l = 3$)
k, p	FPT by a $O((2k)^p \times kn^2)$ time algorithm
k	XNL-hard, solvable in $O((8n)^k)$ time

4 Solving Smast on complete collection of triples

A *rooted triple* (or *triple* for short) is a binary tree T s.t. $|L(T)| = 3$. A *collection of triples* is a collection $\mathcal{R} = \{t_1, \dots, t_k\}$ where each t_i is a triple. \mathcal{R} is *complete* iff each set of three labels in $L(\mathcal{R})$ is present in at least one t_i . To a binary tree T , we associate a complete collection of triples $rt(T)$ formed by the triples $t_i \leq T$. For a complete collection \mathcal{R} , we say that \mathcal{R} is *treelike* iff there exists a tree T s.t. $\mathcal{R} = rt(T)$; then we say that \mathcal{R} *displays* T .

Let P-SMAST-CR denote the restriction of P-SMAST to complete collections of triples. We can show that non-treelike collections have conflicts of size ≤ 4 , a result similar to that known on quartets [11]. This allows to solve P-SMAST-CR in $O(n^4 + 3.12^p)$ time by reduction to 4-HITTING SET [12]. and also in $O(4^p n^4)$ time by bounded search (see, e.g. [13]). In the following, we describe a faster algorithm with $O(4^p n^3)$ running time. We first present an algorithm to decide treelikeness in linear $O(n^3)$ time (Proposition 1 and Theorem 5).

Proposition 1. *There is an algorithm INSERT-LABEL-OR-FIND-CONFLICT(\mathcal{R}, X, x, T) which takes a complete collection of triples \mathcal{R} , a set $X \subseteq L(\mathcal{R})$, an element $x \in L(\mathcal{R}) \setminus X$ and a tree T s.t. $\mathcal{R}|_X$ displays T , and in $O(n^2)$ time decides if $\mathcal{R}' = \mathcal{R}|(X \cup \{x\})$ is treelike. Additionally, the algorithm returns the tree T' displayed by \mathcal{R}' in case of positive answer, or returns a conflict C between \mathcal{R}' with $|C| \leq 4$ in case of negative answer.*

Proof. In a first step, the algorithm checks whether \mathcal{R} contains two different triples on the same set of three labels x, ℓ, ℓ' . In such a case, they form a conflict of size 3 which is

then returned by the algorithm. Suppose now that no such conflict is found. Let u be an internal node of T , and let v, v' be the two children of u . An u -fork is a pair $\{l, l'\}$ where $l \in L(v), l' \in L(v')$. Each u -fork $\{l, l'\}$ will propose a status $s_{l,l'}$ for the positioning of x w.r.t. u in T , where $s_{l,l'}$ is computed from \mathcal{R} as follows: $s_{l,l'} = L$ if $lx|l' \in \mathcal{R}$, $s_{l,l'} = R$ if $l'x|l \in \mathcal{R}$, $s_{l,l'} = U$ if $ll'|x \in \mathcal{R}$ (L, R, U respectively stand for left, right and up).

The second step of the algorithm consists in successively considering each internal node u . For a given node u , it checks that the different u -forks propose the same status. This verification is performed as follows: (i) for each u -fork $\{l, l'\}$, determine the status $s_{l,l'}$; (ii) if there exists l, l'_1, l'_2 s.t. $s_{l,l'_1} \neq s_{l,l'_2}$, then $C = \{x, l, l'_1, l'_2\}$ is a conflict; (iii) if there exists l_1, l_2, l' s.t. $s_{l_1,l'} \neq s_{l_2,l'}$, then $C = \{x, l_1, l_2, l'\}$ is a conflict; (iv) else, all values $s_{l,l'}$ are identical, let s_u denote this value.

In a third step, the algorithm checks that the different statuses are compatible. They are compatible iff for each edge u, v of T with u above v , we have: (i) if v is the left child of u , then $s_u = R \Rightarrow s_v = U$, (ii) if v is the right child of u , then $s_u = L \Rightarrow s_v = U$, (iii) if v is a child of u , then $s_u = U \Rightarrow s_v = U$. If one pair of nodes u, v does not meet the above requirements, then by considering $\{l, l'\}$ v -fork and $\{l, l''\}$ u -fork, we obtain a conflict $C = \{x, l, l', l''\}$. Otherwise, consider the sets of nodes u s.t. $s_u \neq U$, they form a (possibly empty) path in T starting at the root and ending at a node v . Then $\mathcal{R}|(X \cup \{x\})$ is treelike, and displays the tree obtained from T by inserting x above v , which is returned by the algorithm.

We now justify the running time of the algorithm. The first step trivially takes $O(n^2)$ time. Consider the second step. Given a node u , let F_u be the set of u -forks, then an internal node u is processed in time $O(|F_u|)$. Therefore, the time required by the second step is $\sum_u O(|F_u|) = O(n^2)$. Now consider the third step. The algorithm checks that for each edge u, v of T , Conditions (i)-(ii)-(iii) hold: for a given edge, checking the conditions or finding a conflict is done in constant time, hence the time required by this step is $O(n)$. It follows that the total time required by the algorithm is $O(n^2)$. \square

Theorem 5. *There is an algorithm FIND-TREE-OR-CONFLICT(\mathcal{R}) which takes a complete collection of triples \mathcal{R} , and in $O(n^3)$ time decides if \mathcal{R} is treelike, returns a tree T displayed by \mathcal{R} in case of positive answer, or a conflict C between \mathcal{R} with $|C| \leq 4$ in case of negative answer.*

Proof. We use the procedure INSERT-LABEL-OR-FIND-CONFLICT to decide treelikeness as follows. We iteratively insert each label, starting from an empty tree, until: (i) either every label has been inserted, in which case the collection is treelike and the displayed tree is returned, (ii) or a conflict is found and returned. \square

Using bounded search, we obtain:

Theorem 6. *The P-SMAST-CR problem can be solved in $O(4^p n^3)$ time.*

References

1. Aho, A.V., Sagiv, Y., Szymanski, T.G., Ullman, J.D.: Inferring a tree from lowest common ancestors with an application to the optimization of relational expressions. *SIAM Journal on Computing* **10**(3) (1981) 405–421
2. Xia, Y., Yang, Y.: Mining closed and maximal frequent subtrees from databases of labeled rooted trees. *IEEE Transactions on Knowledge and Data Engineering* **17**(2) (2005) 190–202
3. Gordon, A.G.: Consensus supertrees: the synthesis of rooted trees containing overlapping sets of labelled leaves. *Journal of Classification* **3** (1986) 335–348
4. Ranwez, V., Berry, V., Criscuolo, A., Guillemot, S., Douzery, E.: Vote or veto: desirable properties for supertree methods. submitted to *syst. biol.*, LIRMM (2007)
5. Berry, V., Nicolas, F.: Maximum agreement and compatible supertrees. In Sahinalp, S.C., Muthukrishnan, S., Dogrusoz, U., eds.: *Proceedings of CPM*. Volume 3109 of LNCS. (2004) 205–219
6. Jansson, J., Ng, J.H.K., Sadakane, K., Sung, W.K.: Rooted maximum agreement supertrees. In: *Proceedings of the 6th Latin American Symposium on Theoretical Informatics (LATIN)*. (2004) (in press).
7. Kao, M.Y.: *Encyclopedia of algorithms*. <http://refworks.springer.com/algorithms/> (2007)
8. Steel, M., Warnow, T.: Kaikoura tree theorems: computing the maximum agreement subtree. *Information Processing Letters* **48**(2) (1993) 77–82
9. Berry, V., Nicolas, F.: Maximum agreement and compatible supertrees. *Journal of Discrete Algorithms* (in press) (2007)
10. Henzinger, M., King, V., Warnow, T.: Constructing a Tree from Homeomorphic Subtrees, with Applications to Computational Evolutionary Biology. *Algorithmica* **24**(1) (1999) 1–13
11. Bandelt, H., Dress, A.: Reconstructing the shape of a tree from observed dissimilarity data. *Advances in Applied Mathematics* **7** (1986) 309–343
12. Fernau, H.: *Parameterized algorithmics: A graph-theoretic approach*. Habilitationsschrift, Universität Tübingen, Germany (2005)
13. Gramm, J., Niedermeier, R.: A fixed-parameter algorithm for minimum quartet inconsistency. *Journal of Computer and System Sciences* **67**(4) (2003) 723–741
14. Guillemot, S., Berry, V.: Fixed-parameter tractability of the maximum agreement supertree problem. Technical report, LIRMM (2007)

5 Appendix

We give below proofs for some results stated in the paper. Proofs for other results can be found in [14]. We first establish a preliminary lemma.

Lemma 9. *Suppose that $\mathcal{T} = \{T_1, \dots, T_k\}$ is compatible. Then:*

1. *for each $L \subseteq L(\mathcal{T})$, $\mathcal{T}|L$ is compatible.*
2. *for each position π in \mathcal{T} , $\mathcal{T}(\pi)$ is compatible.*

Proof. Let S be a total agreement supertree for \mathcal{T} .

To prove Point 1, consider $L \subseteq L(\mathcal{T})$. Then for each $i \in [k]$ we have $T_i \leq S \Rightarrow T_i|L \leq S|L$. Hence $S|L$ is a total agreement supertree for $\mathcal{T}|L$.

To prove Point 2, consider a position π in \mathcal{T} . Remark that $\mathcal{T}|L(\pi)$ is compatible by Point 1. Let S be a total agreement supertree for $\mathcal{T}|L(\pi)$, we claim that S is a total agreement supertree for $\mathcal{T}(\pi)$. Indeed, we have $L(S) = L(\pi) = L(\mathcal{T}(\pi))$, and given $i \in [k]$ we have $T_i(u_i) \leq T_i|L(\pi) \leq S$. \square

Proof of Lemma 3 (\Rightarrow). Suppose that π is compatible. Let S be a total agreement supertree for $\mathcal{T}(\pi)$. Since $|L(\pi)| \geq 2$, we have $S = (S_1, S_2)$. Since S is a total agreement supertree for $\mathcal{T}(\pi)$, we have $T_i(\pi[i]) \leq S$ for each $i \in [k]$. Define a partition V_1, V_2 of V as follows. Let $u_i = \pi[i]$, and suppose that u_i is an internal node of T_i , with children v_i, v'_i . Then $T_i(u_i) = (T_i(v_i), T_i(v'_i))$. Together with $T_i(u_i) \leq S$, this yields $T_i(v_i) \leq S_1$ or $T_i(v_i) \leq S_2$: add v_i to V_1 in the first case, to V_2 in the second case. Proceed similarly for v'_i .

We first prove Point (i). To see that V_1, V_2 are disconnected in $G(\mathcal{T}, \pi)$, observe that $L(V_1) \subseteq L(S_1)$ and $L(V_2) \subseteq L(S_2)$. Indeed, if $x \in V_j$ with $x \in \{v_i, v'_i\}$ then $T_i(x) \leq S_j$, hence $L(x) \subseteq L(S_j)$. It follows that V_1, V_2 are disconnected by Lemma 2. We now prove Point (ii). Let $\pi_j = \text{succ}_{V_j}(\pi)$, then π_j is a position in $\mathcal{T}(\pi)$, and since $\mathcal{T}(\pi)$ is compatible by assumption, it follows that π_j is compatible by Point 2 of Lemma 9.

(\Leftarrow). Suppose that there exists a partition V_1, V_2 of V satisfying Points (i), (ii). Let $\pi_j = \text{succ}_{V_j}(\pi)$. Since $\text{succ}_{V_j}(\pi)$ is compatible, it follows that π_j is compatible. Hence, there exists a total agreement supertree S_j for $\mathcal{T}(\pi_j)$, which thus satisfies: $T_i(\pi_j[i]) \leq S_j$ for each i . By Lemma 1, we then have $L(S_j) = L(\pi_j) \subseteq L(V_j)$. Since V_1, V_2 are disconnected in $G(\mathcal{T}, \pi)$, it follows that $L(V_1) \cap L(V_2) = \emptyset$ by Lemma 2. Therefore we have $L(S_1) \cap L(S_2) = \emptyset$, and we can define the tree $S = (S_1, S_2)$. We show that S is a total agreement supertree for $\mathcal{T}(\pi)$: to this end, we need to show that $T_i(\pi[i]) \leq S$ for each $i \in [k]$.

Fix such an i , let $u_i = \pi[i]$. If $u_i = \perp$, then the relation holds obviously. Suppose now that u_i is an internal node of T_i , and let v_i, v'_i be its two children. We consider three cases. If $v_i, v'_i \in V_1$: then $\pi_1[i] = u_i$, therefore we have $T_i(u_i) \leq S_1$, and we conclude that $T_i(u_i) \leq S$. If $v_i, v'_i \in V_2$: then $\pi_2[i] = u_i$, therefore we have $T_i(u_i) \leq S_2$, and we conclude that $T_i(u_i) \leq S$. If $v_i \in V_1, v'_i \in V_2$: then $\pi_1[i] = v_i$, which implies that $T_i(v_i) \leq S_1$, and $\pi_2[i] = v'_i$, which implies that $T_i(v'_i) \leq S_2$. It is easy to see that $T_i(v_i) \leq S_1$ and $T_i(v'_i) \leq S_2$ imply that $T_i(u_i) = (T_i(v_i), T_i(v'_i)) \leq (S_1, S_2) = S$. \square

Proof of Lemma 4: We show that $\mathcal{T}' = \mathcal{T}|C$ is incompatible. For each $i \in I(\pi)$, let $u_i = \pi[i]$, and let v_i, v'_i be its two children in T_i . By definition of C , the sets $L(v_i) \cap C, L(v'_i) \cap C$ are not empty, hence to the nodes v_i, v'_i, u_i there corresponds nodes $\tilde{v}_i, \tilde{v}'_i, \tilde{u}_i$ in $T_i|C$. Define the position π' by setting $\pi'[i] = \perp$ if $i \notin I(\pi)$, $\pi'[i] = \tilde{u}_i$ if $i \in I(\pi)$. Consider the graph $G(\mathcal{T}', \pi')$, then by definition of C for each edge $\{x, y\}$ of T , the edge $\{\tilde{x}, \tilde{y}\}$ is present in $G(\mathcal{T}', \pi')$, therefore the tree T' formed of these edges is a spanning tree of $G(\mathcal{T}', \pi')$, hence the graph is connected. By Lemma 3, we conclude that π' is an incompatible position of \mathcal{T}' , therefore \mathcal{T}' is incompatible (by Point 2 of Lemma 9). \square

End of the proof of Theorem 1. We show that the procedure ISCOMPATIBLE can be implemented as a $O(kn)$ time algorithm. Obviously, (i) testing if $\pi = \pi_\perp$ is done in $O(k)$ time, (ii) given T spanning tree of $G(\mathcal{T}, \pi)$, constructing C is done in $|T| = O(k)$ time, provided we have stored a label l_e for each edge of T , (iii) given V_1, V_2 partition of V , constructing the positions π_1, π_2 is done in $O(k)$ time. We now justify that in $O(kn)$ time we can perform a connexity test on $G(\mathcal{T}, \pi)$.

The crucial point is that the algorithm tests the connexity of the graph, by working on the intersection model of $G := G(\mathcal{T}, \pi)$ provided by the sets $\{L(x) : x \in V\}$. In this way, we avoid constructing the adjacency matrix of G , which would require $O(k^2n)$ time. We thus need to describe a connexity test for a graph $G = (V, E)$ given by an intersection model $\{S_v : v \in V\}$, where the S_v are subsets of a base set S . We will justify that the algorithm has running time $O(kn)$, where $k = |V|$ and $n = |S|$.

The algorithm proceeds as follows. It performs a traversal of the graph, by starting at an arbitrary vertex $u \in V$, and maintains the following information during the traversal: (i) the set U of nodes already visited, (ii) a set F of edges forming a spanning tree of $G[U]$. At each step, the algorithm seeks a *transverse edge*, which is an edge $e = \{u, v\} \in E$ with $u \in U, v \in \bar{U}$. If such an edge is found, then v is added to U , and e is added to F . If no such edge exists, the algorithm stops, and the graph is connected iff $U = V$.

We show that using appropriate data structures, a step of the algorithm can be done in $O(n)$ time. For each $x \in S$, let $V_x = \{v \in V : x \in S_v\}$. We maintain for each $x \in S$, two lists representing the sets $U_x = V_x \cap U$ and $\bar{U}_x = V_x \cap \bar{U}$. Initializing these lists at the beginning of the algorithm is done in $O(kn)$ time. Moreover, at a given step of the algorithm: (i) we can find a transversal edge in $O(n)$ time, (ii) we can update the structures in $O(n)$ time. To justify Point (i), observe that finding a transversal edge amounts to find an element $x \in S$ s.t. U_x, \bar{U}_x are non empty; if such an x is found then by choosing $u \in U_x, v \in \bar{U}_x$ we obtain a transverse edge $\{u, v\}$; clearly, these operations can be performed in $O(n)$ time. To justify Point (ii), observe that when visiting a new vertex v , we need, for each $x \in S_v$, to add v to U_x and to remove v from \bar{U}_x , which can be performed in $O(n)$ time by using appropriate linkage. \square

End of the proof of Lemma 8.

(i) $L(S_1) \cap L(S_2) = \emptyset$: indeed, if there was $x \in L(S_1) \cap L(S_2)$ then we would have $x \in L(T_i)$ for some $i \in [k]$; then $x = \phi_{j,i}(x) \leq_{T_i} \pi_j[i]$. Since $x \in L(T_i)$, we must have $\pi_1[i], \pi_2[i] \neq \perp$; then they are equal to distinct children of $\pi[i]$, impossible.

(ii) ϕ_i is a partial embedding of S into T_i :

- if $x \in L(S)$, then $x \in L(S_j)$. We conclude using the fact that $\phi_{j,i}$ is a partial embedding and that $\phi_i(x) = \phi_{j,i}(x)$.
- if x is an internal node of S with children x', x'' , then:
 - if $x \in N(S_j)$, we conclude using the fact that $\phi_{j,i}$ is a partial embedding and that $\phi_i(x) = \phi_{j,i}(x)$.
 - if $x = r(S)$, with children $x' = r(S_1), x'' = r(S_2)$: then
 - * either $\phi_{1,i}(x') = \phi_{2,i}(x'') = \perp$, in which case $\phi_i(x) = \perp$;
 - * either $\phi_{1,i}(x') \neq \perp, \phi_{2,i}(x'') = \perp$, in which case $\phi_i(x) = \phi_{1,i}(x')$;
 - * either $\phi_{1,i}(x') = \perp, \phi_{2,i}(x'') \neq \perp$, in which case $\phi_i(x) = \phi_{2,i}(x'')$;
 - * either $\phi_{1,i}(x') \neq \perp, \phi_{2,i}(x'') \neq \perp$, in which case these are nodes y', y'' s.t. $y' \leq_{T_i} \pi_1[i], y'' \leq_{T_i} \pi_2[i]$. Since $(\pi_1, \pi_2) \in D(\pi)$, it follows that $\pi_1[i], \pi_2[i]$ are $\neq \perp$ and are distinct children of $\pi[i]$, hence $\phi_i(x) = \pi[i]$, which implies that $\phi_{1,i}(x') <_{T_i} \phi_i(x), \phi_{2,i}(x'') <_{T_i} \phi_i(x)$.

(iii) $\phi_i(r(S)) \leq_{T_i} \pi[i]$: follows from the definition of $\phi_i(r(S))$ and from the fact that $\phi_i(r(S_j)) \leq_{T_i} \pi_j[i]$.

We show that $(\pi_1, \pi_2) \in D(\pi)$. Indeed, (i) we have $\pi_j \neq \pi$ since if we had $\pi_1[i] = \pi[i]$ for each i , this would imply $\pi_2[i] = \perp$ for each i , but given $x \in L(S_2)$ there exists i s.t. $x \in L(T_i)$, impossible; (ii) fix $i \in [k]$:

- if $\pi[i] = \perp$, then $\phi_i(u) = \perp$, and we then have $\phi_i(v_1) = \phi_i(v_2) = \perp$ by definition of a partial embedding, hence $\pi_1[i] = \pi_2[i] = \perp$;
- if $\pi[i] \neq \perp$, then $\phi_i(u)$ is a node of T_i , and we have:
 - either $\phi_i(v_1), \phi_i(v_2) \neq \perp$, in which case the nodes $\text{child}_{T_i}(\phi_i(v_1), \phi_i(u)), \text{child}_{T_i}(\phi_i(v_2), \phi_i(u))$ are distinct, which implies that $\pi_1[i], \pi_2[i]$ are distinct children of $\pi[i]$;
 - or one of $\phi_i(v_1), \phi_i(v_2)$ is equal to \perp , in which case the other must be equal to $\phi_i(u)$, which implies that $\pi_1[i] = \pi[i], \pi_2[i] = \perp$ or the symmetric case.

□