

# Statistically Consistent Rooting of Species Trees Under the Multispecies Coalescent Model

Yasamin Tabatabaee<sup>1</sup>, Sébastien Roch<sup>2</sup>, and Tandy Warnow<sup>1(⊠)</sup>

 $^{1}\,$  Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, USA

warnow@illinois.edu

 $^{2}\,$  Department of Mathematics, University of Wisconsin-Madison, Madison, USA

**Abstract.** Rooted species trees are used in several downstream applications of phylogenetics. Most species tree estimation methods produce unrooted trees and additional methods are then used to root these unrooted trees. Recently, Quintet Rooting (QR) (Tabatabaee et al., ISMB and Bioinformatics 2022), a polynomial-time method for rooting an unrooted species tree given unrooted gene trees under the multispecies coalescent, was introduced. QR, which is based on a proof of identifiability of rooted 5-taxon trees in the presence of incomplete lineage sorting, was shown to have good accuracy, improving over other methods for rooting species trees when incomplete lineage sorting was the only cause of gene tree discordance, except when gene tree estimation error was very high. However, the statistical consistency of QR was left as an open question. Here, we present QR-STAR, a polynomial-time variant of QR that has an additional step for determining the rooted shape of each quintet tree. We prove that QR-STAR is statistically consistent under the multispecies coalescent model, and our simulation study shows that QR-STAR matches or improves on the accuracy of QR. QR-STAR is available in open source form at https://github.com/ytabatabaee/Quintet-Rooting.

**Keywords:** Species Tree Estimation  $\cdot$  Rooting  $\cdot$  Statistical Consistency  $\cdot$  Multispecies Coalescent

#### 1 Introduction

Inferring rooted species trees is important for many downstream applications of phylogenetics, such as comparative genomics [7,11] and dating [25]. These estimations use different loci from across the genomes of the selected species, and so are referred to as multi-locus analyses. If rooted gene trees can be accurately inferred, then the rooted species tree can be estimated from them [14]; however, this is not a reliable assumption [29]. Hence, the standard approach is to first

The full paper is available at https://doi.org/10.1101/2022.10.26.513897 and includes the Supplementary Materials.

<sup>©</sup> The Author(s) 2023

H. Tang (Ed.): RECOMB 2023, LNBI 13976, pp. 41–57, 2023.

estimate an unrooted species tree using multi-locus datasets, and then root that estimated tree.

The estimation of the unrooted species tree is challenged by biological processes, such as incomplete lineage sorting (ILS) or gene duplication and loss (GDL), that can result in different parts of the genome having different evolutionary trees. When ILS or GDL occur, statistically consistent estimation of the unrooted species tree requires techniques that take the source of heterogeneity into consideration [15,22]. The case of ILS, as modeled by the multispecies coalescent (MSC) model [10], is the most well-studied, and there are several methods for estimating unrooted species trees that have been proven statistically consistent under the MSC (see [22] for a survey of such methods).

The general problem of rooting a species tree (or indeed even a gene tree) is of independent interest, but presents many challenges. A common approach is the use of an outgroup taxon (i.e., the inclusion of a species that is outside the smallest clade containing the remaining species), so that the resultant tree is rooted on the edge leading to the outgroup [16]. However, outgroup selection has its own difficulties: if the outgroup is too distant, then it may be attached fairly randomly to the tree containing the remaining species, and if it is too close, it may even be an ingroup taxon [5,6,9,13]. Other approaches use branch lengths estimated on the tree to find the root based on specific optimization criteria; however, these approaches tend to degrade in accuracy unless the strict molecular clock holds (which assumes that all sites along the genome evolve under a constant rate) [8,18,31].

Quintet Rooting (QR) [30] is a recently introduced method that is designed to root a given species tree using the unrooted gene tree topologies, under the assumption that the gene trees can differ from the species tree due to ILS. QR is based on mathematical theory established by Allman, Degnan, and Rhodes [2], which showed that the rooted species tree topology is identifiable from the unrooted gene tree topologies whenever the number of species is at least five. In [30], QR was shown to provide good accuracy for rooting both estimated and true species trees in the presence of ILS, compared to alternative methods.

However, QR was not proven to be statistically consistent for locating the root. Thus, we do not have a proof that the root location selected by QR, when given the true species tree topology, will converge to the correct location as the number of gene trees in the input increases. Although much attention has been paid to establishing statistical consistency for unrooted species tree estimation methods and many methods, such as ASTRAL [19], SVDQuartets [32] and BUCKy [12], have been proven to be statistically consistent estimators of the unrooted species tree topology under the MSC, to the best of our knowledge, no prior study has addressed the statistical consistency properties of methods for rooting species trees.

In this paper, we argue that QR is not guaranteed to be statistically consistent under the MSC, but we also present a modification to QR, which we call QR-STAR, that we prove statistically consistent. Moreover, QR-STAR, like QR, runs in polynomial time. We also provide results of a simulation study comparing QR to QR-STAR. Due to space limitations, most of the proofs and results from the simulation study are presented in the full version of the paper available at https://doi.org/10.1101/2022.10.26.513897.

## 2 Background

We present the theory from [2] first, which establishes identifiability of the rooted species tree from unrooted quintet trees, and then we describe Quintet Rooting (QR), our earlier method for rooting species trees. Together these form the basis for deriving our new method, QR-STAR, which we present in the next section.

## 2.1 Allman, Degnan, and Rhodes (ADR) Theory

Allman, Degnan, and Rhodes (ADR) [2] established that the unrooted topology of the species tree is identifiable from four-leaf unrooted gene trees under the MSC, a result that is well known and used in several "quartet-based" methods for estimating species trees under the MSC [12,17,19,24]. ADR also proved that the rooted species tree topology is identifiable from unrooted five-leaf gene tree topologies; this result is much less well known, but was recently used in the development of QR for rooting species trees.

ADR have described the probability distribution of unrooted gene tree topologies under each 5-taxon MSC model species tree. On a given set of five taxa, there exist 105 different rooted binary trees, labeled with  $R_1, \ldots, R_{105}^{-1}$ , that can be categorized into three groups based on their (unlabeled) rooted shapes: caterpillar, balanced and pseudo-caterpillar [27]. An example of a tree from each category is shown in Fig. 1. Each 5-taxon model species tree defines a specific probability distribution over the 15 different unrooted gene tree topologies on the same leafset, shown with  $T_1, \ldots, T_{15}$ . Theorem 9 in [2] states that this distribution uniquely determines the rooted tree topology and its internal branch lengths for trees with at least five taxa.

To prove this identifiability result, the ADR theory specifies a set of linear invariants (i.e., equalities) and inequalities that must hold between the probabilities of unrooted 5-taxon gene trees, for any choice of the parameters of the model species tree. These linear invariants and inequalities define a partial order on the probabilities of the topologies of the different 5-taxon unrooted gene trees. In other words, two gene tree probabilities  $u_i = \mathbb{P}(T_i)$  and  $u_j = \mathbb{P}(T_j)$  can have one of four possible relationships:  $u_i > u_j$ ,  $u_j > u_i$ ,  $u_i = u_j$ , or  $u_i$  and  $u_j$  are not comparable.

<sup>&</sup>lt;sup>1</sup> The labeling of rooted and unrooted trees in this paper is consistent with the notations and leaf-labeling used in Tables 4–5 in [2] as well as in [30].

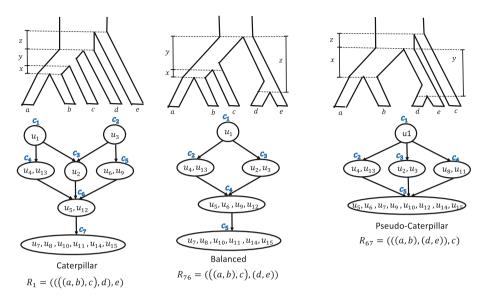


Fig. 1. ADR invariants and inequalities for different rooted topological shapes. The invariants and inequalities found by ADR define, for each rooted 5-taxon model tree topology, a partial order on the probabilities of the 15 unrooted 5-leaf gene trees; importantly, the partial order depends only on the "rooted shape" of the 5-taxon model species trees (i.e., caterpillar, balanced and pseudo-caterpillar). Thus, the topology of any 5-leaf rooted binary species tree is uniquely determined by the partial order, and so can be determined from the true distribution on unrooted 5-leaf gene trees (i.e., it is identifiable, as established by ADR).

Figure 1 shows examples of these partial orders with a Hasse diagram for a particular leaf labeling of trees from each rooted shape. Note that some probabilities are members of the same set (e.g., for  $R_1$ , set  $c_4$  contains both  $u_4$  and  $u_{13}$ , indicating that  $u_4 = u_{13}$ ), and so we refer to the sets  $c_i$  as equivalence classes on these probabilities. Furthermore, we will denote the set of equivalence classes associated with a 5-taxon rooted tree R with  $C_R$ . As can be seen in Fig. 1, the number of equivalence classes depends on the shape of the rooted species tree, with caterpillar, balanced and pseudo-caterpillar trees having 7, 5 and 5 equivalence classes, respectively.

Each directed edge between two equivalence classes in these Hasse diagrams defines an inequality, so that all gene tree probabilities in class  $c_a$  at the source of an edge are greater than all gene tree probabilities in class  $c_b$  at the target, and we denote this by  $c_a > c_b$ . The exact values of the unrooted gene tree probabilities depend on the internal branch lengths of the model tree, and ADR provide a set of formulas that relate the model tree parameters to the probability distribution of the unrooted gene trees in Appendix B of [2], which will be used in our proofs.

## 2.2 Quintet Rooting

The input to QR is an unrooted species tree T with n leaves and a set  $\mathcal{G}$  of k single-copy unrooted gene trees where the gene trees draw their leaves from the leafset of T, denoted by  $\mathcal{L}(T)$ . Given this input, QR searches over all possible rootings of T and returns a tree most consistent with the distribution of quintets (i.e., 5-taxon trees) in the input gene trees.

QR approaches this problem by selecting a set Q of quintets of taxa from  $\mathcal{L}(T)$  (called the "quintet sampling" step; refer to Supplementary Materials Sec. A for details), and scoring all rooted versions of T based on their induced trees on these quintets. The subtree  $T_{|q}$ , T restricted to taxa in quintet set q, can be rooted on any of its seven edges. In a preprocessing step, QR computes a score for each of these seven different rootings for all trees induced on the quintets in set Q, based on a cost function. This results in  $7 \times |Q|$  computations, and therefore the preprocessing step takes O(k(|Q|+n)). Next, for every rooted version of T, QR sums up the costs of all its induced rooted trees on quintets in Q using the scores computed in the preprocessing step, and returns the rooting with the minimum overall cost. Since T can be rooted on any of its 2n-3 edges, the scoring step takes O(n+|Q|) time. Therefore, the overall runtime of QR when using an O(n) sampling of quintets is O(nk). Figure 2 shows the pipeline of QR and its individual steps.

Thus, QR provides an exact solution to the optimization problem with the following input and output:

- **Input:** An unrooted tree topology T, a set of k unrooted gene tree topologies  $\mathcal{G} = \{g_1, g_2, \ldots, g_k\}$ , a set Q containing quintets of taxa from leafset  $\mathcal{L}(T)$  and a cost function  $\text{Cost}(r, \vec{u})$ .
- **Output:** Rooted tree R with topology T such that  $\sum_{q \in Q} \operatorname{Cost}(R|_q, \vec{u}_q)$  is minimized, where  $\vec{u}_q$  is the distribution of unrooted gene tree quintets in  $\mathcal{G}|_q = \{g_1|_q, g_2|_q, \ldots, g_k|_q\}.$

Cost Function. The cost function  $\operatorname{Cost}(R|_q,\vec{u}_q)$  measures the fitness of the rooted quintet tree  $R|_q$  with the distribution of the unrooted gene trees restricted to q (i.e.,  $\vec{u}_q$ ), according to the linear invariants and inequalities derived from the ADR theory. In particular, this cost function is designed to penalize a rooted tree  $R|_q$  if the estimated quintet distribution  $\vec{u}_q$  violates some of the inequalities or invariants in its partial order. To this end, a penalty term was considered for each invariant and inequality in the partial order of a 5-taxon rooted tree that is violated in a quintet distribution. The cost function was defined based on a linear combination of these penalty terms, and had the following form, where r is a 5-taxon rooted tree and  $\vec{u}$  is an estimated quintet distribution:

$$\operatorname{Cost}(r, \vec{\hat{u}}) = \underbrace{\sum_{c \in C_r} \frac{1}{|c|} \sum_{u_a, u_b \in c} |\hat{u}_a - \hat{u}_b|}_{\text{Invariants Penalty}} + \underbrace{\sum_{c > c' \in C_r} \frac{1}{|c'|} \sum_{u_a \in c, u_b \in c'} \max(0, \hat{u}_b - \hat{u}_a)}_{\text{Inequalities Penalty}}.$$

(1)

The normalization factors  $\frac{1}{|c|}$  and  $\frac{1}{|c'|}$  were used to reduce a topological bias that arose from differences in the sizes of the equivalence classes for each tree shape.

## 3 QR-STAR

QR-STAR is an extension to QR that has an additional step for determining the rooted shape (i.e., the rooted topology without the leaf labels) of a quintet tree, as well as an associated penalty term in its cost function. This penalty term compares the rooted shape of the 5-taxon tree, denoted by S(r), with the rooted shape inferred by QR-STAR from the given quintet distribution, denoted by  $\hat{S}(\hat{u})$ . The motivation for this additional preprocessing step is that, as we argue in Supplementary Materials Sec. C, the cost function of QR does not guarantee statistical consistency. The cost function of QR-STAR takes the following general form

$$\operatorname{Cost}^*(r, \vec{\hat{u}}) = \underbrace{\sum_{c \in C_r} \sum_{u_a, u_b \in c} \alpha_{a,b} |\hat{u}_a - \hat{u}_b|}_{\text{Invariants Penalty}} + \underbrace{\sum_{c > c' \in C_r} \sum_{u_a \in c, u_b \in c'} \beta_{a,b} \max(0, \hat{u}_b - \hat{u}_a)}_{\text{Inequalities Penalty}} + \underbrace{C1|S(r) \neq \hat{S}(\hat{u})|}_{\text{Shape Penalty}}$$
(2)

where  $\alpha_{a,b} \geq 0$  and  $\beta_{a,b}, C > 0$  are constant real numbers for all  $a, b^2$ . Let  $\alpha_{\max} = \max_{a,b}(\alpha_{a,b})$  and  $\beta_{\min} = \min_{a,b}(\beta_{a,b})$  where a, b ranges over all pairs of indices a, b used in the penalty terms in Eq. 2.

Each of the 105 rooted binary trees on a given set of 5 leaves have a unique set of inequalities and invariants that can be derived from the ADR theory. The cost function in Eq. 2 considers a penalty term for these inequalities and invariants as well as the shape of the tree, so that  $\operatorname{Cost}^*(r, \vec{u})$  is minimized for a rooted 5-taxon tree r that best describes the given estimated quintet distribution.

#### 3.1 Determining the Rooted Shape

Model 5-taxon species trees with different rooted shapes (i.e., caterpillar, balanced, pseudo-caterpillar) define equivalence classes with different class sizes on the unrooted gene tree probability distribution. These class sizes can be used to determine the unlabeled shape of a rooted tree, when given the *true* gene tree probability distribution. For example, the size of the equivalence class with the smallest gene tree probabilities is 8 for the pseudo-caterpillar trees and 6 for balanced or caterpillar trees. Therefore, the size of the equivalence class corresponding to the minimal element in the partial order can differentiate a pseudo-caterpillar tree from other tree shapes. Moreover, both caterpillar and

<sup>&</sup>lt;sup>2</sup> Refer to Remark 1 for why  $\alpha_{a,b}$  does not need to be strictly positive.

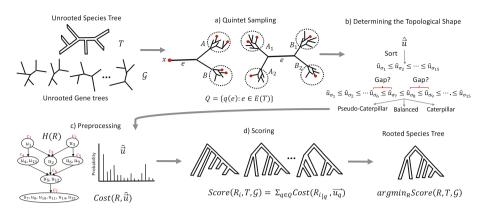


Fig. 2. Pipeline of QR and QR-STAR. The input is an unrooted species tree T and set of unrooted gene trees  $\mathcal G$  on the same leafset. a) The sampling step selects a set Q of quintets from the leafset of T (shown is the linear encoding sampling). b) An additional step in QR-STAR that determines the rooted shape for each selected quintet. c) The preprocessing step computes a cost for each of the seven possible rootings of each selected quintet. d) The scoring step computes a score for each rooted tree in the search space based on the costs computed in the preprocessing step, and returns a rooting of T with minimum score.

balanced trees have a unique class with the second smallest probability, which is of size 2 for caterpillar trees and 4 for balanced trees and this can be used to differentiate a caterpillar tree from a balanced tree. This approach is used in Theorem 9 in [2] for establishing the identifiability of rooted 5-taxon trees from unrooted gene trees.

However, given an *estimated* gene tree distribution, it is likely that none of the invariants derived from the ADR theory exactly hold, and so the class sizes cannot be directly determined and the approach above cannot be used as is to infer the shape of a rooted quintet. Here we propose a simple modification for determining the rooted shape of a tree from the estimated distribution of unrooted gene trees, by looking for significant gaps between quintet gene tree probabilities.

Let T be the unrooted species tree with  $n \geq 5$  leaves given to QR-STAR and q be a quintet of taxa from  $\mathcal{L}(T)$ . Let  $\vec{u}$  be the quintet distribution estimated from input gene trees induced on taxa in set q. QR-STAR first sorts  $\vec{u}$  in ascending order to get  $\hat{u}_{\sigma_1} \leq \hat{u}_{\sigma_2} \leq \cdots \leq \hat{u}_{\sigma_{15}}$ . Let  $A(k) = \sqrt{\frac{2}{k} \ln(30|Q|k)}$  (refer to Lemma 4 for the derivation of A(k)), where k is the number of input gene trees and |Q| is the size of the set of sampled quintets in QR-STAR (this depends on the number n of taxa and is assumed fixed), and note that  $\lim_{k\to\infty} A(k) = 0$ . The first step of QR-STAR computes an estimate of the rooted shape of a quintet q, denoted by  $\hat{S}(\hat{u})$  in Eq. 2, as follows:

- estimate the rooted shape  $\hat{S}(\hat{u})$  as pseudo-caterpillar if  $\hat{u}_{\sigma_7} \hat{u}_{\sigma_6} < A(k)$ ;
- estimate the rooted shape  $\hat{S}(\hat{u})$  as balanced if  $\hat{u}_{\sigma_7} \hat{u}_{\sigma_6} \geq A(k)$  and  $\hat{u}_{\sigma_9} \hat{u}_{\sigma_8} < A(k)$ ;
- estimate the rooted shape  $\hat{S}(\hat{u})$  as caterpillar if  $\hat{u}_{\sigma_7} \hat{u}_{\sigma_6} \geq A(k)$  and  $\hat{u}_{\sigma_9} \hat{u}_{\sigma_8} \geq A(k)$ .

The runtime of QR-STAR is the same as QR, as determining the topological shape for each quintet is done in constant time, so that the overall runtime remains O(nk) when a linear sampling of quintets is used.

## 4 Theoretical Results

In this section, we provide the main theoretical results, starting with a series of lemmas and theorems that will be used in the proof of statistical consistency of QR-STAR in Theorem 2. Throughout this paper, we assume that discordance between species trees and gene trees is solely due to ILS. In establishing statistical consistency, we assume that input gene trees are true gene trees and, thus, have no gene tree estimation error. If not otherwise specified, all trees are assumed to be fully resolved (i.e., binary). Due to space constraints, the proofs are provided in Supplementary Materials Section B. We begin with some definitions and key observations.

**Definition 1 (Path length parameter).** Let R be an MSC model species tree. Let f(R) be the length of the shortest internal branch of R and g(R) be the length of the longest internal path (i.e., a path formed from only the internal branches) of R. We define the path length parameter of R as

$$h(R) = \frac{1}{18}e^{-3g(R)}(1 - e^{-f(R)})^2$$
(3)

Note that  $h(R) \in (0, \frac{1}{18})$  since  $\exp(-x) \in (0, 1)$  for all x > 0 and the branch lengths have positive values. The formula for Eq. 3 is derived from the proof of Lemma 2 in Supplementary Materials Sec. B.

**Lemma 1.** Let R be an MSC model species tree with  $n \geq 5$  leaves and q be an arbitrary set of 5 leaves from  $\mathcal{L}(R)$ . Then  $h(R|_q) \geq h(R)$  where  $R|_q$  is the rooted tree R restricted to taxa in set q.

**Lemma 2.** Let R be an MSC model species tree with 5 leaves and internal branch lengths x, y, and z. Let  $\vec{u}$  be the probability distribution that R defines on the unrooted 5-taxon gene tree topologies. If  $\vec{u}$  is an estimate of  $\vec{u}$  such that given  $\epsilon > 0$ , we have  $|\hat{u}_i - u_i| < \epsilon$  for all  $1 \le i \le 15$ , then the following inequality holds:

$$\forall_{c>c'\in C_R} \forall_{u_a\in c, u_b\in c'} : \hat{u}_a - \hat{u}_b > h(R) - 2\epsilon. \tag{4}$$

**Definition 2.** For a 5-taxon rooted tree R, we define  $I_R$  as the set of ordered pairs (i,j),  $1 \le i \ne j \le 15$ , corresponding to inequalities in the form  $u_i > u_j$  defined according to the partial order of R. The inequalities that are a result of transitivity (i.e.  $u_i > u_j$  and  $u_j > u_k$  implies  $u_i > u_k$ ) are not included in  $I_R$ .

**Definition 3.** Let V(R, R') be the set of violated inequalities of two rooted 5-taxon trees R and R', i.e., all pairs  $\{i, j\}$  such that  $(i, j) \in I_R$  and  $(j, i) \in I_{R'}$ .

Figure 3a shows an example of V(R,R') computed for caterpillar trees and Fig. 3b is a heatmap showing the function |V(R,R')| computed for the seven possible rootings of an unrooted quintet tree. The set V(R,R') can be easily computed from  $I_R$  and  $I_{R'}$  for all pairs of rooted 5-taxon trees, and  $I_R$  is derived from the ADR theory for all 105 5-taxon rooted trees in the Supplementary Materials, Sec. S2 in [30].

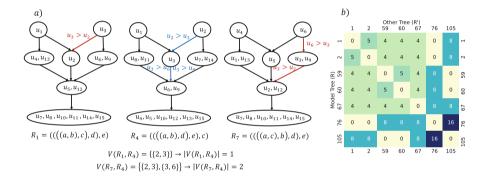


Fig. 3. Conflicting inequality penalty terms between rooted 5-taxon species trees. a) Set of violated inequality penalty terms in the partial orders of  $R_1$  and  $R_7$  with respect to  $R_4$ , which are all caterpillar trees. The red edges show violations of inequalities in tree  $R_4$ , highlighted in blue. b) Heatmap showing the number of pairwise violated penalty terms (function |V(R, R')|) of seven possible rooted trees having unrooted topology with bipartitions ab|cde and abc|de. The dark colors indicate more violations, and the lightest color corresponds to no violations (|V(R, R')| = 0). (Color figure online)

**Lemma 3.** (a) For 5-taxon binary rooted trees R and R' with the same rooted shape, the set V(R, R') is always non-empty. (b) For each balanced tree B, there exist two caterpillar trees  $C_1$  and  $C_2$  such that  $V(B, C_i) = \emptyset$ .

#### 4.1 Statistical Consistency

In this section, we establish statistical consistency for QR-STAR under the MSC and provide the sufficient condition for a set of sampled quintets to lead to consistency. That is, we prove that as the number of input true gene trees increases,

the probability that QR-STAR and its variants correctly root the given unrooted species tree converges to 1. We first prove statistical consistency for QR-STAR when the model tree has only five taxa in Theorem 1 and then extend the proofs to trees with arbitrary numbers of taxa in Theorem 2. The main idea of the proof of consistency for 5-taxon trees is that we show as the number of input gene trees increases, the cost of the true rooted tree becomes arbitrarily close to zero, but the cost of any other rooted tree is bounded away from zero, where the bound depends on the path length parameter of the model tree, h(R).

**Lemma 4.** Let R be an MSC model species tree with  $n \ge 5$  leaves and Q be a set of quintets of taxa from  $\mathcal{L}(R)$ . Given  $\delta > 0$  and k > 0 unrooted gene tree topologies, the following inequality holds, where  $A_{\delta}(k) = \sqrt{\frac{2}{k} \ln(\frac{30|Q|}{\delta})}$ 

$$\mathbb{P}\left(\forall_{q\in Q}\forall_{1\leq i\leq 15}|(\hat{u}_q)_i - (u_q)_i| < \frac{A_\delta(k)}{2}\right) \ge 1 - \delta. \tag{5}$$

Setting  $\delta = \frac{1}{k}$  in Eq. 5, we get  $A(k) = \sqrt{\frac{2}{k} \ln(30|Q|k)}$ , which is the bound that is used for determining the rooted shape of each quintet in the first step of QR-STAR as well as the proofs of statistical consistency. When R has only five taxa, A(k) becomes  $\sqrt{\frac{2}{k} \ln(30k)}$ , as Q can only contain one quintet.

**Lemma 5 (Correct determination of rooted shape).** Let R be a 5-taxon model species tree and  $\vec{u}$  be the probability distribution that it defines on the unrooted 5-taxon gene tree topologies. There is an integer k>0 such that if we are given at least k unrooted gene trees drawn i.i.d. from the distribution  $\vec{u}$ , the first step of QR-STAR will correctly determine the rooted shape of R with probability at least  $1-\frac{1}{k}$ .

Lemma 6 (Upper bound on the cost of the model tree). Let R be a 5-taxon model species tree and  $\vec{u}$  be the probability distribution that it defines on the unrooted 5-taxon gene tree topologies. There is an integer k > 0 such that if we are given at least k unrooted gene trees drawn i.i.d. from distribution  $\vec{u}$ , then  $\text{Cost}^*(R, \overline{\hat{u}})$  is less than  $31\alpha_{\max}A(k)$  with probability at least  $1 - \frac{1}{k}$ .

Theorem 1 (Statistical Consistency of QR-STAR for 5-taxon trees). Let R be a 5-taxon model species tree and  $\vec{u}$  be the distribution that it defines on the unrooted 5-taxon gene tree topologies. Given a set  $\mathcal G$  of unrooted true quintet gene trees drawn i.i.d. from  $\vec{u}$ , QR-STAR is a statistically consistent estimator of R under the MSC.

Remark 1. Note that when  $\alpha_{\max}=0$ , meaning that the invariant penalty terms are removed from the cost function, the cost of the true tree would become exactly zero according to the proof of Lemma 6, and the cost of any other tree would be positive when k is large enough so that the conditions of Theorem 1 hold. Hence in this case, the condition in Eq. 7 (see full version of the paper) will reduce to  $A(k) < \frac{1}{2}h(R)$ .

Remark 2. Note that Lemma 3(a) holds for all pairs of 5-taxon rooted trees with the same rooted shape and with different permutations of the leaf-labeling, regardless of whether they have the same unrooted topology or not. Due to this property, it is possible to differentiate all pairs of 5-taxon rooted trees in a statistically consistent manner with the cost function of QR-STAR, without prior knowledge about the unrooted tree topology, and hence Theorem 1 does not assume that the unrooted topology is given as input.

The next lemma and theorem extend the proof of statistical consistency to trees with n > 5 taxa. The linear encoding of a tree T by quintets is defined in Supplementary Materials Section A.

Lemma 7 (Identifiability of the root from the linear encoding). Let R and R' be rooted trees with unrooted topology T and distinct roots. Let  $Q_{LE}(T)$  be the set of quintets of leaves in a linear encoding of T. There is at least one quintet of taxa  $q \in Q_{LE}(T)$  so that  $R_{|q}$  and  $R'_{|q}$  have different rooted topologies.

Lemma 7 states that no two distinct rooted trees with topology T induce the same set of rooted quintet trees on quintets of taxa in set  $Q_{LE}(T)$ . Clearly, the same is true for any superset Q such that  $Q_{LE}(T) \subseteq Q$ , including the set  $Q_5$  of all quintets of taxa on the leafset of T. There might also be other quintet sets that are not a superset of  $Q_{LE}(T)$ , but have the property that no two rooted versions of T define the same set of rooted quintets on their elements. We generalize the proof of consistency to all set of sampled quintets with this property.

**Definition 4.** Let T be an unrooted tree and Q be a set of quintets of taxa from  $\mathcal{L}(T)$ . We say Q is "root-identifying" if every rooted tree R with topology T is identifiable from T and the set of rooted quintet trees in  $\{R|_q: q \in Q\}$ , i.e., no two rooted trees with topology T induce the same set of rooted quintet trees on Q.

**Theorem 2 (Statistical Consistency of QR-STAR).** Let R be an MSC model species tree with  $n \geq 5$  leaves and let T denote its unrooted topology. Given T and a set G of unrooted true gene trees on the leafset L(T), QR-STAR is a statistically consistent estimator of the rooted version of T under the MSC, if the set of sampled quintets Q is root-identifying.

## 5 Experimental Study

We performed an experimental study on simulated datasets to explore the parameter space of QR-STAR on a training dataset, and then compared its accuracy to QR on a test dataset. We used the 101-taxon simulated datasets from [34] as our training data, which had model conditions characterized by four levels of gene tree estimation error (GTEE) ranging from 0.23 to 0.55 (measured in terms of normalized Robinson-Foulds (RF) [26] distance between true and estimated gene trees) for 1000 genes. The normalized RF distance between the model species tree and

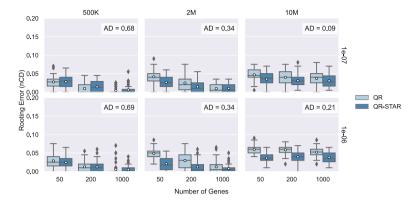


Fig. 4. Rooting the model species tree with estimated gene trees on S200 datasets. Comparison between QR and QR-STAR in terms of rooting error (nCD) for rooting the true unrooted species tree topology using estimated gene trees (GTEE levels vary from 0.22 (for low ILS) to 0.49 (for high ILS)) on the 201-taxon datasets from [20] with 50 replicates in each model condition. The columns show tree height (500K for high ILS, 2M for moderate ILS, and 10M for low ILS), and the rows show speciation rate (1e-06 or 1e-07).

true gene trees (denoted average distance, or AD) in this dataset was 0.46, which indicates moderate ILS. For the test data (see Table D1 in the Supplementary Materials for empirical statistics), we used a set of 201-taxon simulated datasets from [20]; these are characterized by two different speciation rates and three tree heights (500K for high ILS, 2M for moderate ILS, and 10M for low ILS) and three numbers of genes for each combination of speciation rate and tree height. GTEE levels on the test data varied from 0.22 (for low ILS) to 0.49 (for high ILS). The AD levels ranged from 0.09 (for the 10M, 1e–07 condition) to 0.69 (for the 500K, 1e–06 condition). The number of replicates for each model condition for both the training and test datasets was 50.

We measured the error in the rooted species tree in terms of average normalized clade distance (nCD) [30], which is an extension of RF error for rooted trees. For our training experiment, we only rooted the true species tree topology to directly observe the rooting error. In our test experiments, we rooted both the model species tree and estimated species tree, as produced by ASTRAL, using both true and estimated gene trees (which were estimated using FastTree [23]). Additional information about the simulation study, datasets, and software commands are provided in Supplementary Materials Section D.

In our training experiments, we explored the impact of the shape coefficient C and the ratio  $\frac{\alpha_{max}}{\beta_{min}}$  (that describes the relative impact of invariants and inequalities) on the accuracy of QR-STAR. Results for the training experiments (provided in Supplementary Materials Sec. E1) show that there are wide ranges of settings for the algorithmic parameters that provide the best accuracy. We used these training results and theoretical considerations related to sample

complexity of QR-STAR to set the algorithmic parameters to C = 1e-02 and  $\frac{\alpha_{max}}{\beta} = 0$ .

 $\frac{\overline{\beta_{min}}}{\beta_{min}} = 0.$ Figure 4 shows a comparison between QR and QR-STAR in terms of rooting error for rooting the model species tree topology using estimated gene trees on the test datasets. Increasing the ILS level (by reducing tree height) decreases the rooting error, and increasing the number of genes also generally reduces rooting error (although much less under the lowest ILS level where tree height is 10M). To understand the impact of ILS in Fig. 4, note that the true species tree is being rooted and so ILS will not impact species tree estimation accuracy. However, the level of ILS impacts information about rooting location, which comes from the distribution of gene tree topologies. Thus, with lower ILS, it is likely that many gene trees that have low probability of appearing will not appear in the input. In this case, some estimates of quintet probabilities would become zero, and it may not be possible to differentiate some of the rooted quintets using the inequality and invariants derived from the ADR theory. In the extreme case, when there is no discordance, there will be only one quintet gene tree with non-zero probability, and the identifiability theorem in [2] would not hold and it becomes impossible to find the root. This trend can be compared to the impact of ILS level on the problem of estimating the unrooted topology of the species tree, where increases in ILS generally lead to increases in error [19–21].

A comparison between QR and QR-STAR shows that QR-STAR generally matched or improved on QR; the only exception was for the high ILS conditions, where the two methods were very close but with perhaps a small advantage to QR. On these high ILS conditions, however, GTEE is also large, and QR-STAR is more accurate than QR when used with true gene trees, even under high ILS (Supplementary Materials Sec. E). Hence, the issue is likely to be high GTEE rather than high ILS, suggesting that QR-STAR is slightly more affected by GTEE compared to QR.

#### 6 Conclusion

In this work we presented QR-STAR, a polynomial time statistically consistent method for rooting species trees under the multispecies coalescent model. QR-STAR is an extension to QR, a method for rooting species trees introduced in [30]. QR-STAR differs from QR in that it has an additional step for determining the topological shape of each unrooted quintet selected in the QR algorithm, and incorporates the knowledge of this shape in its cost function, alongside the invariants and inequalities previously used in QR. We also showed that the statistical consistency for QR-STAR holds for a larger family of optimization problems based on cost functions and sampling methods.

To the best of our knowledge, this is the first work that established the statistical consistency of any method for rooting species trees under a model that incorporates gene tree heterogeneity. It remains to be investigated whether other rooting methods can also be proven statistically consistent under models of gene evolution inside species trees, such as the MSC or models of GDL. For example,

STRIDE [4] and DISCO+QR [33] are methods that have been developed for rooting species trees from gene family trees, where genes evolve under gene duplication and loss (GDL); however, it is not known whether these methods are statistically consistent under any GDL model.

Our simulation study showed as well that QR-STAR generally improved on QR in a wide range of model conditions. Given that QR itself improved on other methods for rooting species trees (as shown in [30]), this experimental study suggests that QR-STAR may be a useful tool for rooting species trees when gene tree discordance due to ILS is present.

This study suggests several directions for future research. For example, we proved statistical consistency for one class of cost functions, which was a linear combination of the invariant, inequality and shape penalty terms; however, cost functions in other forms could also be explored and proven statistically consistent. The proof of Theorem 1 suggests that the sample complexity of QR-STAR depends on the function h(R), which is based on both the length of the shortest branch and the longest path in the model tree. This suggests that having very short or very long branches can both confound rooting under ILS, which is also suggested in previous studies [1,2]. This is unlike what is known for species tree estimation methods such as ASTRAL, where the sample complexity is only affected by the shortest branch of the model tree [3,28], and trees with long branches are easier to estimate.

Another theoretical direction is the construction of the rooted species tree directly from the unrooted gene trees. As explained in Remark 2, the proof of consistency of QR-STAR for 5-taxon trees does not depend upon the knowledge of the unrooted tree topology; this suggests that it is possible to estimate the rooted topology of the species tree in a statistically consistency manner directly from unrooted gene tree topologies. Future work could focus on developing statistically consistent methods for this problem, which is significantly harder than the problem of rooting a given tree.

There are also directions for improving empirical results. An important consideration in designing a good cost function is its empirical performance, as many cost functions can lead to statistical consistency but may not provide accurate estimations of the rooted tree in practice (see Figures E1 and E2 in the Supplementary Materials). One potential direction is to incorporate estimated branch lengths, whether of the gene trees or the unrooted species tree, into the rooting procedure.

**Acknowledgments.** SR was supported by NSF grants DMS-1902892, DMS1916378 and DMS-2023239 (TRIPODS Phase II), as well as a Vilas Associates Award. TW was supported by the Grainger Foundation. SR thanks Cécile Ané and her group for helpful discussions. YT thanks Mohammed El-Kebir for helpful suggestions on an earlier version of this work. The authors thank the reviewers for their feedback.

## References

- Alanzi, A.R., Degnan, J.H.: Inferring rooted species trees from unrooted gene trees using approximate Bayesian computation. Mol. Phylogenet. Evol. 116, 13–24 (2017). https://doi.org/10.1016/j.ympev.2017.07.017
- Allman, E.S., Degnan, J.H., Rhodes, J.A.: Identifying the rooted species tree from the distribution of unrooted gene trees under the coalescent. J. Math. Biol. 62(6), 833–862 (2011). https://doi.org/10.1007/s00285-010-0355-7
- 3. Chan, Y., Li, Q., Scornavacca, C.: The large-sample asymptotic behaviour of quartet-based summary methods for species tree inference. J. Math. Biol. **85**(3), 1–22 (2022). https://doi.org/10.1007/s00285-022-01786-4
- 4. Emms, D.M., Kelly, S.: STRIDE: species tree root inference from gene duplication events. Mol. Biol. Evol. **34**(12), 3267–3278 (2017)
- Felsenstein, J.: Cases in which parsimony or compatibility methods will be positively misleading. Syst. Zool. 27(4), 401–410 (1978)
- Graham, S.W., Olmstead, R.G., Barrett, S.C.: Rooting phylogenetic trees with distant outgroups: a case study from the commelinoid monocots. Mol. Biol. Evol. 19(10), 1769–1781 (2002)
- Skarp-de Haan, C.: Comparative genomics of unintrogressed Campylobacter coli clades 2 and 3. BMC Genomics 15(1), 1–14 (2014). https://doi.org/10.1186/1471-2164-15-129
- 8. Hess, P.N., De Moraes Russo, C.A.: An empirical test of the midpoint rooting method. Biol. J. Linn. Soc. **92**(4), 669–674 (2007)
- Holland, B., Penny, D., Hendy, M.: Outgroup misplacement and phylogenetic inaccuracy under a molecular clock-a simulation study. Syst. Biol. 52(2), 229–238 (2003). https://doi.org/10.1080/10635150390192771
- Hudson, R.R.: Testing the constant-rate neutral allele model with protein sequence data. Evolution 203–217 (1983). https://doi.org/10.2307/2408186
- Jun, S.R., et al.: Ebolavirus comparative genomics. FEMS Microbiol. Rev. 39(5), 764–778 (2015)
- Larget, B.R., Kotha, S.K., Dewey, C.N., Ané, C.: BUCKy: gene tree/species tree reconciliation with Bayesian concordance analysis. Bioinformatics 26(22), 2910–2911 (2010). https://doi.org/10.1093/bioinformatics/btq539
- 13. Li, C., Matthes-Rosana, K.A., Garcia, M., Naylor, G.J.: Phylogenetics of Chondrichthyes and the problem of rooting phylogenies with distant outgroups. Mol. Phylogenet. Evol. **63**(2), 365–373 (2012)
- 14. Liu, L., Yu, L., Edwards, S.V.: A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. BMC Evol. Biol. **10**(1), 1–18 (2010). https://doi.org/10.1186/1471-2148-10-302
- Maddison, W.P.: Gene trees in species trees. Syst. Biol. 46(3), 523–536 (1997). https://doi.org/10.1093/sysbio/46.3.523
- Maddison, W.P., Donoghue, M.J., Maddison, D.R.: Outgroup analysis and parsimony. Syst. Biol. 33(1), 83–103 (1984)
- 17. Mahbub, M., Wahab, Z., Reaz, R., Rahman, M.S., Bayzid, M.S.: wQFM: highly accurate genome-scale species tree estimation from weighted quartets. Bioinformatics 37(21), 3734–3743 (2021). https://doi.org/10.1093/bioinformatics/btab428
- 18. Mai, U., Sayyari, E., Mirarab, S.: Minimum variance rooting of phylogenetic trees and implications for species tree reconstruction. PLoS ONE 12(8), e0182238 (2017)
- 19. Mirarab, S., Reaz, R., Bayzid, M.S., Zimmermann, T., Swenson, M.S., Warnow, T.: ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics **30**(17), i541–i548 (2014). https://doi.org/10.1093/bioinformatics/btu462

- 20. Mirarab, S., Warnow, T.: ASTRAL-II: coalescent-based species tree estimation with many hundreds of taxa and thousands of genes. Bioinformatics **31**(12), i44–i52 (2015). https://doi.org/10.1093/bioinformatics/btv234
- Molloy, E.K., Warnow, T.: To include or not to include: the impact of gene filtering on species tree estimation methods. Syst. Biol. 67(2), 285–303 (2018)
- 22. Posada, D.: Phylogenomics for systematic biology. Syst. Biol. **65**(3), 353–356 (2016). https://doi.org/10.1093/sysbio/syw027
- 23. Price, M.N., Dehal, P.S., Arkin, A.P.: FastTree 2-approximately maximum-likelihood trees for large alignments. PLoS ONE 5(3), e9490 (2010)
- 24. Rabiee, M., Mirarab, S.: QuCo: quartet-based co-estimation of species trees and gene trees. Bioinformatics **38**(Supplement\_1), i413–i421 (2022)
- Renner, S.S., Grimm, G.W., Schneeweiss, G.M., Stuessy, T.F., Ricklefs, R.E.: Rooting and dating maples (Acer) with an uncorrelated-rates molecular clock: implications for North American/Asian disjunctions. Syst. Biol. 57(5), 795–808 (2008). https://doi.org/10.1080/10635150802422282
- 26. Robinson, D.F., Foulds, L.R.: Comparison of phylogenetic trees. Math. Biosci. **53**(1–2), 131–147 (1981). https://doi.org/10.1016/0025-5564(81)90043-2
- 27. Rosenberg, N.A.: Counting coalescent histories. J. Comput. Biol. **14**(3), 360–377 (2007). https://doi.org/10.1089/cmb.2006.0109
- Shekhar, S., Roch, S., Mirarab, S.: Species tree estimation using ASTRAL: how many genes are enough? IEEE/ACM Trans. Comput. Biol. Bioinf. 15(5), 1738– 1747 (2017). https://doi.org/10.1109/TCBB.2017.2757930
- Simmons, M.P., Springer, M.S., Gatesy, J.: Gene-tree misrooting drives conflicts in phylogenomic coalescent analyses of palaeognath birds. Mol. Phylogenet. Evol. 167, 107344 (2022). https://doi.org/10.1016/j.ympev.2021.107344
- Tabatabaee, Y., Sarker, K., Warnow, T.: Quintet Rooting: rooting species trees under the multi-species coalescent model. Bioinformatics 38(Supplement\_1), i109– i117 (2022). https://doi.org/10.1093/bioinformatics/btac224
- Tria, F.D.K., Landan, G., Dagan, T.: Phylogenetic rooting using minimal ancestor deviation. Nat. Ecol. Evol. 1(1), 1–7 (2017). https://doi.org/10.1038/s41559-017-0193
- 32. Wascher, M., Kubatko, L.: Consistency of SVDQuartets and maximum likelihood for coalescent-based species tree estimation. Syst. Biol. **70**(1), 33–48 (2021). https://doi.org/10.1093/sysbio/syaa039
- 33. Willson, J., Tabatabaee, Y., Liu, B., Warnow, T.: DISCO+QR: rooting species trees in the presence of GDL and ILS. Bioinform. Adv. **3**(1), vbad015 (2023). https://doi.org/10.1093/bioadv/vbad015
- 34. Zhang, C., Rabiee, M., Sayyari, E., Mirarab, S.: ASTRAL-III: polynomial time species tree reconstruction from partially resolved gene trees. BMC Bioinform. 19(6), 15–30 (2018). https://doi.org/10.1186/s12859-018-2129-y

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (http://creativecommons.org/licenses/by/4.0/), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

