

# THE DISTRIBUTIONS UNDER TWO SPECIES-TREE MODELS OF THE NUMBER OF ROOT ANCESTRAL CONFIGURATIONS FOR MATCHING GENE TREES AND SPECIES TREES

BY FILIPPO DISANTO<sup>1,a</sup>, MICHAEL FUCHS<sup>2,b</sup>, ARIEL R. PANINGBATAN<sup>3,c</sup> AND NOAH A. ROSENBERG<sup>4,d</sup>

<sup>1</sup>Department of Mathematics, University of Pisa, <sup>a</sup>[filippo.disanto@unipi.it](mailto:filippo.disanto@unipi.it)

<sup>2</sup>Department of Mathematical Sciences, National Chengchi University, <sup>b</sup>[mfuchs@nctu.edu.tw](mailto:mfuchs@nctu.edu.tw)

<sup>3</sup>Institute of Mathematics, University of the Philippines Diliman, <sup>c</sup>[arpaningbatan@math.upd.edu.ph](mailto:arpaningbatan@math.upd.edu.ph)

<sup>4</sup>Department of Biology, Stanford University, <sup>d</sup>[noahr@stanford.edu](mailto:noahr@stanford.edu)

For a pair consisting of a gene tree and a species tree, the *ancestral configurations* at a species-tree internal node are the distinct sets of gene lineages that can be present at that node. The enumeration of *root* ancestral configurations—ancestral configurations at the species-tree root—assists in describing the complexity of gene-tree probability calculations in evolutionary biology. Assuming that the gene tree and species tree match in topology, we study the distribution of the number of root ancestral configurations of a random labeled tree topology under the *uniform* and *Yule–Harding* models. We employ analytic combinatorics, considering ancestral configurations in the context of additive tree parameters and using singularity analysis to evaluate asymptotic growth of the coefficients of generating functions. For both models, we obtain asymptotic lognormal distributions for the number of root ancestral configurations. For Yule–Harding random trees, we also obtain the asymptotic mean ( $\sim 1.425^n$ ) and variance ( $\sim 2.045^n$ ) of the number of root ancestral configurations, paralleling previous results for the uniform model (mean  $(4/3)^n$ , variance  $\sim 1.822^n$ ). A methodological innovation is that to obtain the Yule–Harding asymptotic variance, singularity analysis is conducted from the Riccati differential equation that the generating function satisfies—without possessing the generating function itself.

**1. Introduction.** In the study of combinatorial properties of species trees (trees that describe evolutionary relationships among species) and gene trees (trees that describe evolutionary relationships among gene lineages for members of the species), one useful concept is that of an ancestral configuration. Given a gene tree, a species tree and a node of the species tree, an ancestral configuration is a list of the gene lineages that are present at the node of the species tree (Figure 1). Looking backward in time, or from the leaves of trees to the root, the fact that gene lineages only find their common ancestors once their associated species have found common ancestors produces conditions describing which ancestral configurations are present at a species tree node. These conditions enable the enumeration of the configurations. Ancestral configurations appear in recursive evaluations of the probabilities of gene tree topologies conditional on species tree topologies [46], so that enumerations of ancestral configurations assist in assessing the complexity of the computation.

When the node at which an ancestral configuration is considered is the root node of the species tree, ancestral configurations are termed *root ancestral configurations*, or *root configurations* for short. For matching gene trees and species trees—that is, if the species tree and gene tree have the same labeled topology—the number of root configurations is greater

---

Received December 2020; revised September 2021.

*MSC2020 subject classifications.* Primary 60C05, 92D15; secondary 05A15, 05A16, 05C05, 92B10.

*Key words and phrases.* Analytic combinatorics, gene trees, lognormal distribution, phylogenetics, Riccati equation, species trees.

than or equal to the number of ancestral configurations for any other species tree node. This property can be used to show that as the number of leaves increases, the total number of ancestral configurations for the gene tree and species tree—the sum of the number of ancestral configurations across all species tree nodes—has the same exponential growth as the number of root configurations ([14], Section 2.3.2). Hence, it suffices for investigations of the exponential growth of the total number of ancestral configurations for matching gene trees and species trees to focus on root configurations.

Disanto and Rosenberg [14] studied the number of root configurations for matching gene trees and species trees, considering the number of root configurations of families of increasingly large trees. They characterized the labeled tree topologies with the largest number of root configurations among trees with  $n$  leaves, showing that this number of root configurations lies between  $k_0^{n-1/4} - 1$  and  $k_0^n - 1$ , where  $k_0$  is a constant approximately equal to 1.5028 ([14], Proposition 4). They then studied the number of root configurations in trees selected uniformly at random from the set of labeled topologies with  $n$  leaves. Using techniques of analytic combinatorics, they showed that the mean number of root configurations grows with  $(4/3)^n$ , and the variance with  $\sim 1.8215^n$  ([14], Propositions 5 and 6).

Here, we extend these results on the distribution of the number of root configurations under a model imposing a uniform distribution on the set of labeled topologies. We review background results in Section 2. In Section 3, we describe correspondences between classes of trees, which we use in Section 4 to obtain an asymptotic normal distribution for the logarithm of the number of root configurations under the uniform model—and find that its mean, approximately  $0.272n$ , generates exponential growth  $e^{0.272n} \approx 1.313^n$ . In Section 5, we obtain similar results under the Yule–Harding model, including the asymptotic mean and variance of the number of root configurations and the asymptotic distribution of its logarithm. This set of computations also makes use of a correspondence between tree classes. The calculation of the asymptotic variance additionally employs a novel approach, in which asymptotic growth of the coefficients of a generating function that solves a Riccati equation is obtained without having the exact form of the generating function itself. We discuss the results in Section 6.

**2. Preliminaries.** We study ancestral configurations for rooted binary leaf-labeled trees. In Section 2.1, we introduce results on various classes of trees. In Section 2.2, we discuss the Yule–Harding distribution on labeled topologies. In Section 2.3, we recall properties of generating functions and analytic combinatorics. Following Wu [46], in Section 2.4 we define ancestral configurations, and we review enumerative results from Disanto and Rosenberg [14]. In Section 2.5, we relate ancestral configurations to additive tree parameters, which have been widely studied in the literature [27, 45].

2.1. *Classes of trees.* We will need to consider many classes of trees: labeled topologies, unlabeled topologies, ordered unlabeled topologies, labeled histories, unlabeled histories and ordered unlabeled histories. Many terms in the setting of evolutionary trees can be connected to concepts from settings that do not have a biological context [1, 4, 7]; our terminology generally follows that typical of mathematical studies of evolutionary trees [39].

2.1.1. *Labeled topologies.* We refer to a bifurcating rooted tree  $t$  with  $|t| = n$  labeled leaves as a *labeled topology* of size  $|t| = n$ , or a “tree” for short (Figure 1A); these trees are sometimes called phylogenetic trees or Schröder trees. They are *unordered* or *nonplane* in the sense that if left–right positions of two child nodes are exchanged in a labeled topology, then the same labeled topology is obtained. For the set  $\{a, b, c, \dots\}$  of possible labels for the leaves of a tree, we impose an alphabetical linear order  $a < b < c < \dots$ . The leaf labels of a tree of size  $n$  are the first  $n$  labels in the order  $<$ .

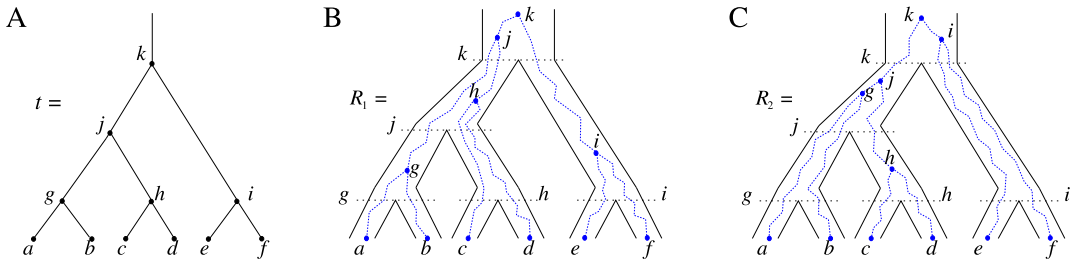


FIG. 1. A gene tree and species tree with matching labeled topology  $t$ . (A) A tree  $t$  of size 6, characterized by its shape and leaf labels. For convenience, we label the internal nodes of  $t$ , by  $g, h, i, j, k$  in this case, identifying each lineage (edge) by its immediate descendant node. For example, lineage  $h$  results from coalescence of lineages  $c$  and  $d$ . (B) A possible realization  $R_1$  of the gene tree in (A) (dotted lines) in the matching species tree (solid lines). The ancestral configurations at species tree nodes  $j$  and  $k$  are  $\{g, c, d\}$  and  $\{g, h, i\}$ , respectively. (C) A different realization  $R_2$  of the gene tree in (A) in the species tree. At species tree nodes  $j$  and  $k$ , the configurations are  $\{a, b, h\}$  and  $\{j, e, f\}$ , respectively. The figure is modified from Figure 1 of Disanto and Rosenberg [14] and Figure 1 of Disanto and Rosenberg [15].

We denote by  $T_n$  the set of trees of size  $n$ , with  $T = \bigcup_{n=1}^{\infty} T_n$  denoting the set of all trees. The number of trees of size  $n \geq 2$  is  $|T_n| = (2n - 3)!! = 1 \times 3 \times 5 \times \dots \times (2n - 3)$  [19], for  $n \geq 1$ ,

$$(1) \quad |T_n| = \frac{(2n - 2)!}{2^{n-1}(n - 1)!} = \frac{(2n)!}{2^n(2n - 1)n!}.$$

The exponential generating function for  $|T_n|$  is

$$T(z) = \sum_{t \in T} \frac{z^{|t|}}{|t|!} = \sum_{n=1}^{\infty} \frac{|T_n|z^n}{n!} = z + \frac{z^2}{2} + \frac{3z^3}{6} + \frac{15z^4}{24} + \dots,$$

given by Flajolet and Sedgewick ([23], Example II.19),

$$(2) \quad T(z) = 1 - \sqrt{1 - 2z}.$$

2.1.2. *Ordered unlabeled topologies.* An *orientation* of an unlabeled topology  $t$  is a plane embedding of  $t$  in which subtrees descending from the internal nodes of  $t$  are considered with a left–right orientation. For instance, the unlabeled topology underlying the labeled topology depicted in Figure 1A has exactly two different orientations, which are depicted in Figure 2A. An orientation of an unlabeled topology is called an *ordered unlabeled topology*, or a *plane unlabeled topology*. The set of all possible ordered unlabeled topologies of size  $n$  is enumerated by the Catalan number  $C_{n-1}$  ([38], Exercise 6.19d), where

$$(3) \quad C_n = \frac{1}{n + 1} \binom{2n}{n}.$$

The ordinary generating function is

$$C(z) = \sum_{n=0}^{\infty} C_n z^n = \frac{1 - \sqrt{1 - 4z}}{2z}.$$

With the leaves and associated incident edges stripped away so that only the tree connecting the internal nodes remains, an ordered unlabeled topology is also called a *Catalan tree* or *pruned binary tree*, for example, by Wagner [45] (see also Flajolet and Sedgewick [23], Example I.13).

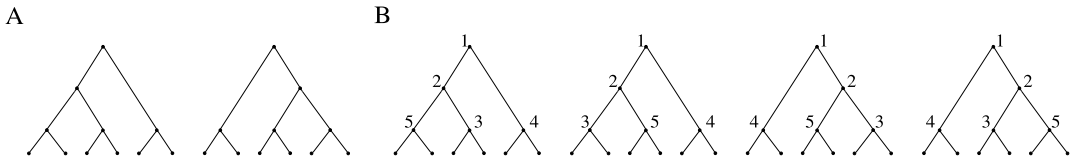


FIG. 2. *Ordered unlabeled topologies and histories. (A) The two orientations of the unlabeled topology that underlies the labeled topology of Figure 1A. (B) The four orientations of the unlabeled history underlying the labeled history in Figure 3A.*

2.1.3. *Labeled histories.* A labeled history is a labeled topology together with a temporal (linear) ordering of its internal nodes (Figure 3). Like a labeled topology, a labeled history is left–right unordered, or nonplane: if the left–right positions of two child nodes are interchanged in a labeled history, then the same labeled history is obtained. If  $t$  is a labeled history of size  $n$ , then we represent the time ordering of its  $n - 1$  bifurcations by bijectively associating each internal node of  $t$  with an integer label in the interval  $[1, n - 1]$ . The labeling is increasing in the sense that each internal node other than the root has a larger label than its parent node.

For a given label set of size  $n$ , the set of labeled histories is denoted  $H_n$ . Its cardinality is ([39], page 46)

$$(4) \quad |H_n| = \frac{n!(n - 1)!}{2^{n-1}}.$$

2.1.4. *Ordered unlabeled histories.* By removing leaf labels of a labeled history  $t$ , we obtain the unlabeled history underlying  $t$ . As we did for unlabeled topologies, we define an orientation of an unlabeled history  $t$  as a plane embedding of  $t$  in which child nodes are considered with a left–right orientation. Figure 2B shows the orientations of the unlabeled history underlying the labeled history of Figure 3A. We call each object so oriented an *ordered unlabeled history*, or a *plane unlabeled history*. The ordered unlabeled histories of size  $n$  are enumerated by  $F_{n-1}$  ([39], page 47),

$$(5) \quad F_n = n!.$$

Ordered unlabeled histories are also called *binary increasing trees* [3, 45] or *ranked oriented trees* [39].

2.2. *The Yule–Harding distribution.* Different labeled histories can share the same underlying labeled topology. For example, the labeled histories of Figure 3 have the underlying

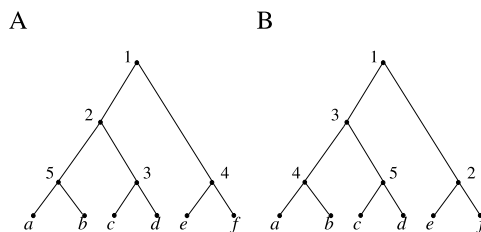


FIG. 3. *Labeled histories. (A) The labeled history of the labeled gene tree topology depicted in Figure 1B. The temporal ordering of the coalescence events in the gene tree is determined by the integer labeling of the internal nodes of the associated labeled topology. (B) The labeled history of the labeled gene tree topology depicted in Figure 1C.*

labeled topology depicted in Figure 1A. The number of labeled histories of size  $n$  with the same labeled topology  $t$  is

$$(6) \quad \frac{(n - 1)!}{\prod_{r=3}^n (r - 1)^{d_r(t)}}$$

where  $d_r(t)$  is the number of internal nodes of  $t$  from which exactly  $r$  leaves descend ([39], page 46). Equation (6) also appears as the so-called “shape functional” of binary search trees [20].

By summing the probability  $1/|H_n|$  of each uniformly distributed labeled history of size  $n$  with a given underlying labeled topology, the uniform distribution over the set  $H_n$  induces the Yule–Harding (or Yule) distribution over the set  $T_n$  of labeled topologies [6, 7, 17, 18, 25, 31, 32, 40, 48]. The probability of a labeled topology  $t$  is

$$(7) \quad P_{\text{YH}}(t) = \frac{2^{n-1}}{n! \prod_{r=3}^n (r - 1)^{d_r(t)}}.$$

Under this distribution, among all labeled topologies with size  $n$ , those with the largest number of labeled histories have the highest probability. For balanced labeled topologies, the product in the denominator of equation (7) tends to be smaller than for unbalanced topologies, resulting in a greater probability.

*2.3. Asymptotic growth and analytic combinatorics.* Our study concerns the growth of increasing sequences. A sequence of nonnegative numbers  $a_n$  is said to have exponential growth  $k^n$  or, equivalently, to be of exponential order  $k$ , if  $a_n = k^n s(n)$ , where  $s$  is subexponential, that is,  $\limsup_{n \rightarrow \infty} [s(n)^{1/n}] = 1$ . Sequence  $a_n$  grows exponentially in  $n$  if its exponential order exceeds 1.

If  $(a_n)$  has exponential order  $k_a$  and  $(b_n)$  has exponential order  $k_b < k_a$ , then the sequence of ratios  $b_n/a_n$  converges to 0 exponentially fast as  $(k_b/k_a)^n$ . If sequences  $a_n$  and  $b_n$  have the same exponential order, then we write  $a_n \asymp b_n$ . If in addition the ratio  $b_n/a_n$  converges to 1, then we write  $a_n \sim b_n$  and say that  $(a_n)$  and  $(b_n)$  have the same asymptotic growth.

Some results make use of techniques of analytic combinatorics (Flajolet and Sedgewick [23], Sections IV and VI). In particular, the entries of a sequence of integers  $(a_n)_{n \geq 0}$  can be interpreted as coefficients of the power series expansion  $A(z) = \sum_{n=0}^{\infty} a_n z^n$  at  $z = 0$  of a function  $A(z)$ , the generating function of the sequence. Considering  $z$  as a complex variable, the behavior of  $A(z)$  near its singularities—the points in the complex plane where  $A(z)$  is not analytic—can provide information on the growth of its coefficients. Under suitable conditions, a correspondence exists between the expansions  $A_\alpha(z)$ ,  $\alpha \in S$ , of the generating function  $A(z)$  near singularities in its set  $S$  of dominant singularities—that is, its singularities of smallest modulus—and the asymptotic growth of the coefficients  $a_n$ . In the simplest case, if  $\alpha$  is the only dominant singularity of  $A(z)$ , then the  $n$ th coefficient  $a_n$  of  $A(z)$  has asymptotic growth  $[z^n]A_\alpha(z)$ , that is, the  $n$ th coefficient of  $A_\alpha(z)$  (Theorem VI.4 of Flajolet and Sedgewick [23]). In symbols,

$$a_n \sim [z^n]A_\alpha(z).$$

The exponential order of sequence  $(a_n)$  is the inverse of the modulus of the dominant singularity  $\alpha$  of  $A(z)$  (Theorem IV.7 of Flajolet and Sedgewick [23]). That is,

$$a_n \asymp \alpha^{-n}.$$

As an example, sequence  $|T_n|/n!$ , with  $|T_n|$  as in equation (1), has exponential order 2 because  $\alpha = \frac{1}{2}$  is the only dominant singularity of the associated generating function in equation (2). Thus, as  $n \rightarrow \infty$ ,  $|T_n|/n!$  increases with a subexponential multiple of  $2^n$ .

2.4. *Ancestral configurations for matching gene trees and species trees.* In this section, following Disanto and Rosenberg [14], we review features of the objects on which our study focuses: the ancestral configurations of a gene tree  $G$  in a species tree  $S$ .

2.4.1. *Gene trees and species trees.* A *species tree* is a tree of evolutionary relationships among a set of species. A *gene tree* is a tree of evolutionary relationships among individual genetic lines of descent, or lineages, at a specific genomic site. Gene trees and species trees are typically viewed as objects evolving forward in time, from the root to the leaves, or backward in time, from the leaves to the root. They consist of both a labeled topology and a set of edge lengths, positive values that describe the lengths of time separating pairs of nodes.

In studies of gene trees and species trees, the leaf label set of a gene tree  $G$  is often taken to be a subset of the leaf label set of a species tree  $S$ , so that a gene tree evolves conditionally on the species tree. Here, because we consider only the combinatorial structure of gene trees and species trees, we are not concerned with numerical values of edge lengths. Hence, it is convenient to identify a gene tree or a species tree with its associated labeled topology; for ease of understanding, however, it is still said that a gene tree or species tree “has” a labeled topology rather than that it “is” a labeled topology. Because we are concerned with ancestor–descendant relationships, it is also convenient to retain a perspective that gene trees and species trees unfold over time.

We here examine the case that the leaf label sets of  $G$  and  $S$  are bijectively associated. In other words, a single genetic lineage is sampled from each species corresponding to a leaf of the species tree. We further restrict attention to the case in which  $G$  and  $S$  have the same labeled topology, so that the gene tree and species tree are said to be *matching*. With the perspective that a gene tree unfolds over time conditionally on a species tree, an instance of the evolutionary process that produces gene tree  $G$  on species tree  $S$  is a *realization* of  $G$  on  $S$ .

Looking backward in time, the lineages of  $G$  are traced back past nodes of  $S$  until the root of  $G$  is reached; at a given point in time, a lineage of  $G$  is associated with a label that contains information about which leaves descend from it. For convenience, a node of a gene tree or species tree is associated with its immediate ancestral edge, so that a node and its immediate ancestral edge are assigned the same label.

2.4.2. *Ancestral configurations.* An ancestral configuration can be viewed as a certain function of a realization of  $G$  on  $S$ , with  $G$  and  $S$  representing a gene tree and a species tree, respectively, and of a node of  $S$ . Suppose  $R$  is a realization of a gene tree  $G$  on a species tree  $S$ , where  $G = S = t$  (Figure 1). Looking backward in time, for node  $\eta$  of  $S$ , consider the set  $C(\eta, R)$  of genetic lineages—edges of  $G$ —that are present in  $S$  at the point in time just before node  $\eta$  is reached.

The set  $C(\eta, R)$  is the *ancestral configuration* of  $G$  at node  $\eta$  of  $S$ . For example, for tree  $t$  in Figure 1A, with the realization  $R_1$  of gene tree  $G = t$  in the species tree  $S = t$  in Figure 1B, just before the root node  $k$ , the gene lineages present in the species tree are lineages  $g, h$  and  $i$ . Hence, at species tree node  $k$ , the ancestral configuration is the set of gene lineages  $C(k, R_1) = \{g, h, i\}$ . Similarly, the ancestral configuration of the gene tree at species tree node  $j$  is  $C(j, R_1) = \{g, c, d\}$ . In Figure 1C, with a different realization  $R_2$  of the same gene tree, the ancestral configuration at the species tree root  $k$  is  $C(k, R_2) = \{j, e, f\}$ . The ancestral configuration at node  $j$  is  $C(j, R_2) = \{a, b, h\}$ .

Let  $\mathfrak{R}(G, S)$  be the set of realizations of gene tree  $G = t$  in species tree  $S = t$ . For a given node  $\eta$  of  $t$ , considering all possible elements  $R \in \mathfrak{R}(G, S)$ , the set of ancestral configurations is

$$(8) \quad C(\eta) = \{C(\eta, R) : R \in \mathfrak{R}(G, S)\}.$$



The associated number of ancestral configurations is

$$(9) \quad c_\eta = |C(\eta)|.$$

The quantity  $c_\eta$  counts the ways the lineages of  $G$  can reach the timepoint right before node  $\eta$  in  $S$ , considering all possible realizations of gene tree  $G$  in species tree  $S$ . Choosing  $t$  as in Figure 1A, we have  $C(g) = \{\{a, b\}\}$ ,  $C(h) = \{\{c, d\}\}$ ,  $C(i) = \{\{e, f\}\}$ ,  $C(j) = \{\{a, b, c, d\}, \{g, c, d\}, \{a, b, h\}, \{g, h\}\}$ , and

$$(10) \quad C(k) = \{\{j, i\}, \{j, e, f\}, \{g, h, i\}, \{g, h, e, f\}, \{a, b, h, i\}, \{a, b, h, e, f\}, \\ \{g, c, d, i\}, \{g, c, d, e, f\}, \{a, b, c, d, i\}, \{a, b, c, d, e, f\}\}.$$

For different realizations  $R_1, R_2 \in \mathfrak{R}(G, S)$  and an internal node  $\eta$ , it need not be true that  $C(\eta, R_1) \neq C(\eta, R_2)$ .

We say that a leaf or a 1-leaf tree has no ancestral configurations. The definition of an ancestral configuration at node  $\eta$ , by considering the timepoint right before node  $\eta$  in the species tree, excludes the case in which all gene tree lineages descended from gene tree node  $\eta$  have coalesced at species tree node  $\eta$ . Thus,  $\{\eta\} \notin C(\eta)$ .

Because we consider the case of  $G = S = t$ , the set  $C(\eta)$  and the quantity  $c_\eta$  in equations (8) and (9) depend only on node  $\eta$  and tree  $t$ . We use the term *configurations at node  $\eta$  of  $t$*  to denote elements of  $C(\eta)$ .

**2.4.3. Root and total configurations.** Our focus is on configurations at the root of  $t$ . Let  $N(t)$  be the set of nodes of a tree  $t$ , including both leaf nodes and internal nodes. With  $|t|$  leaf nodes and  $|t| - 1$  internal nodes in  $t$ ,  $|N(t)| = 2|t| - 1$ . Define the *total* number of configurations in  $t$  by

$$c = \sum_{\eta \in N(t)} c_\eta.$$

Let  $c_r$  be the number of configurations at the root  $r$  of  $t$ , or *root configurations* for short. Because  $c_r \geq c_\eta$  for each node  $\eta$  of  $t$ , we have

$$(11) \quad c_r \leq c \leq (2|t| - 1)c_r.$$

Quantities  $c$  and  $c_r$  are equal up to a factor that is at most polynomial in  $|t|$ , and they have the same exponential order when measured across families of trees of increasing size.

Selecting a tree of size  $n$  at random from the set of labeled topologies, inequality (11) gives  $\mathbb{E}_n[c_r] \leq \mathbb{E}_n[c] \leq 2n\mathbb{E}_n[c_r]$  and  $\mathbb{E}_n[c_r^2] \leq \mathbb{E}_n[c^2] \leq 4n^2\mathbb{E}_n[c_r^2]$ . In expectation  $\mathbb{E}$  and variance  $\mathbb{V}$ , exponential growth for total configurations follows that for root configurations:

$$(12) \quad \mathbb{E}_n[c] \asymp \mathbb{E}_n[c_r],$$

$$(13) \quad \mathbb{E}_n[c^2] \asymp \mathbb{E}_n[c_r^2],$$

$$(14) \quad \mathbb{V}_n[c] = \mathbb{E}_n[c^2] - \mathbb{E}_n[c]^2 \asymp \mathbb{E}_n[c_r^2] - \mathbb{E}_n[c_r]^2 = \mathbb{V}_n[c_r].$$

Equation (14) follows from the fact that the exponential growth of  $\mathbb{E}_n[c^2]$  is faster than that of  $\mathbb{E}_n[c]^2$ , as can be demonstrated from results in the next section (equations (17) and (19)), and the exponential growth of  $\mathbb{E}_n[c_r^2]$  is faster than that of  $\mathbb{E}_n[c_r]^2$  (equations (16) and (18)); we then have  $\mathbb{V}_n[c] \sim \mathbb{E}_n[c^2]$  and  $\mathbb{V}_n[c_r] \sim \mathbb{E}_n[c_r^2]$ , and equation (14) follows from equation (13).

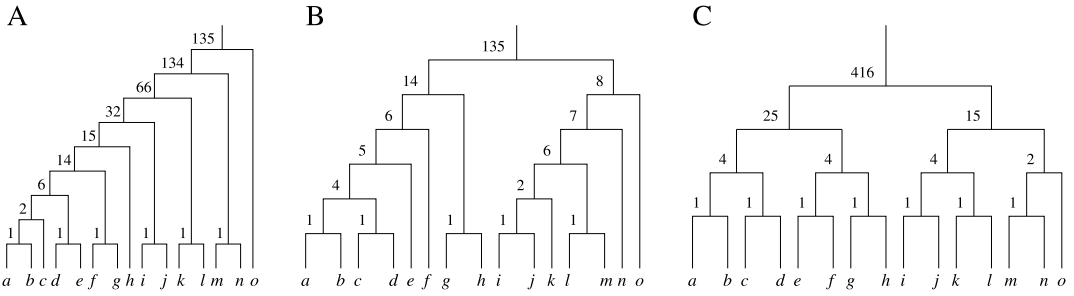


FIG. 4. The number of ancestral configurations at the internal nodes of three labeled topologies of size  $n = 15$ . (A, B) Two labeled topologies in which the number of root configurations is the mean number  $c_r = 135$  of root configurations calculated across the set of representative labelings of the unlabeled topologies of size 15. In this set, the labeled topologies in (A) and (B) have respectively the largest number 61776 and smallest number 14400 of labeled histories. (C) The labeled topology with 15 leaves that has the most root configurations (416) and the most labeled histories (2745600).

2.4.4. *Known results.* We recall some results of Disanto and Rosenberg [14] on the number of configurations possessed by a tree.

(i) For a given tree  $t$  with  $|t| > 1$ , let  $r$  denote the root node of  $t$ , with  $r_L$  and  $r_R$  being the two child nodes of  $r$ . The number  $c_r$  of possible configurations at  $r$  can be recursively computed as

$$(15) \quad c_r = (c_{r_L} + 1)(c_{r_R} + 1),$$

where we set  $c_r = 0$  if  $|t| = 1$ . Figure 4 illustrates the application of equation (15) successively from the leaves to the root of each of three labeled topologies of size  $n = 15$ .

(ii) Consider a representative labeling of each unlabeled topology of size  $n$ . Among these trees, the largest number of root configurations and the largest total number of configurations have exponential order  $k_0$ , where  $k_0 \approx 1.5028$ . The smallest number of root configurations and the smallest total number of configurations have polynomial growth with the tree size  $n$ . Furthermore, consider the balanced family of unlabeled topologies defined recursively by  $|t_1| = 1$  and  $t_n = (t_d, t_{n-d})$ , where  $d$  denotes the power of 2 nearest to  $\frac{n}{2}$ . Among the unlabeled topologies with  $n$  leaves,  $t_n$  has the largest number of root configurations. The maximally asymmetric caterpillar unlabeled topology has the smallest number of root configurations.

(iii) For a labeled topology of given size  $n$  selected uniformly at random, the mean number of root configurations  $c_r$  and the mean total number of configurations  $c$  grow asymptotically like

$$(16) \quad \mathbb{E}_n[c_r] \sim \sqrt{\frac{3}{2}} \left(\frac{4}{3}\right)^n,$$

$$(17) \quad \mathbb{E}_n[c] \asymp \left(\frac{4}{3}\right)^n.$$

The variances of  $c_r$  and  $c$  satisfy the asymptotic relations

$$(18) \quad \mathbb{V}_n[c_r] \sim \sqrt{\frac{7(11 - \sqrt{2})}{34}} \left[ \frac{4}{7(8\sqrt{2} - 11)} \right]^n,$$

$$(19) \quad \mathbb{V}_n[c] \asymp \left[ \frac{4}{7(8\sqrt{2} - 11)} \right]^n.$$



2.5. *Additive tree parameters and root configurations.* A quantity  $F(t)$  that can be computed for trees  $t$  and whose value can be calculated as

$$F(t) = F(t_L) + F(t_R) + f(t),$$

where  $t_L$  and  $t_R$  are the two root subtrees of  $t$ , is called an *additive tree parameter* with *toll function*  $f(t)$  [21, 27, 45]. Additive tree parameters and toll functions have been widely investigated ([27], Remark 1.16). We make use of results from Wagner [45]. For various tree families, Wagner [45] showed that an additive tree parameter  $F(t)$  is asymptotically normally distributed if the toll function  $f(t)$  is bounded and the mean of  $|f(t)|$ , considered over uniformly distributed trees of fixed size, goes to 0 exponentially fast as the tree size increases.

For a tree  $t$ , consider the quantity  $\log(c_r + 1)$ , that is, the natural logarithm of one more than the number of root configurations of  $t$ . From equation (15), a simple calculation yields for  $|t| \geq 2$ ,

$$(20) \quad \log(c_r + 1) = \log(c_{r_L} + 1) + \log(c_{r_R} + 1) + \log\left(1 + \frac{1}{c_r}\right).$$

In equation (20), if we set

$$F(t) = \log[c_r(t) + 1],$$

then the associated toll function is given for  $|t| \geq 2$  by

$$f(t) = \log\left[1 + \frac{1}{c_r(t)}\right].$$

We set  $f(t) = F(t) = \log(1) = 0$  if  $|t| = 1$ . We can therefore consider root configurations in the context of additive tree parameters.

**3. Equivalences for the distribution of the number of root configurations.** We prove a series of equivalences needed for analyzing distributional properties of the number of root configurations. In Section 3.1, we show that the distribution of the number of root configurations over uniformly distributed labeled topologies or labeled histories can be analyzed by considering equivalently the distribution of the number of root configurations over uniformly distributed ordered unlabeled topologies or ordered unlabeled histories, respectively. In Section 3.2, we obtain a correspondence between antichains of pruned binary trees and root configurations of ordered unlabeled topologies.

3.1. *Equivalences with ordered unlabeled topologies and histories.* Distributional properties of a tree parameter defined over the set of labeled topologies can in some cases be investigated by studying the same parameter over a different tree family. In particular, if the tree parameter under consideration depends only on tree topology, then its distribution can be equivalently analyzed over a different tree set taken under a probability model that induces or is induced by the probability model assumed for labeled topologies. In this direction, Blum et al. [4] derived a general framework for analyzing tree parameters of labeled topologies under a variety of probabilistic models defined over binary search trees.

In this section, we obtain results analogous to those of Blum et al. [4]. We show that the number of root configurations—or any other tree parameter that depends only on the branching structure of the tree—has the same distribution when considered over uniformly distributed labeled topologies or over uniformly distributed ordered unlabeled topologies of the same size (Lemma 3.1). Similarly, the number of root configurations has the same distribution over uniformly distributed labeled histories of size  $n$  as for uniformly distributed ordered unlabeled histories of size  $n$  (Lemma 3.3).

Moreover, because the uniform distribution over the set of labeled histories of size  $n$  induces the Yule–Harding distribution over the set of labeled topologies of size  $n$  (Section 2.2), as a direct consequence of Lemma 3.3 we have that the number of root configurations has the same distribution when considered over Yule–Harding-distributed labeled topologies or over uniformly distributed ordered unlabeled histories (Lemma 3.4). By using these facts, Propositions 3.2 and 3.5 give recursive formulas for the probabilities under the uniform and Yule–Harding probability models, respectively, that a random labeled topology of size  $n$  has  $c_r = \rho$  root configurations.

LEMMA 3.1. *The distribution of the number of root configurations over labeled topologies of size  $n$  selected uniformly at random matches the distribution of the number of root configurations over ordered unlabeled topologies of size  $n$  selected uniformly at random.*

PROOF. First, we note that the number of root configurations of a labeled topology or ordered unlabeled topology depends only on the underlying unlabeled topology. Thus, to prove the claim, it suffices to show that for each unlabeled topology  $t$  of size  $n$ , we have

$$(21) \quad \frac{\text{or}(t)}{C_{n-1}} = \frac{\text{lab}(t)}{|T_n|},$$

where  $\text{or}(t)$  and  $\text{lab}(t)$  are the number of orientations of  $t$  and the number of leaf labelings of  $t$ , respectively. Note from equations (3) and (1) that  $\text{or}(t)/C_{n-1}$  and  $\text{lab}(t)/|T_n|$  give the probability of the unlabeled topology  $t$  induced by the uniform distribution over the set of ordered unlabeled topologies and labeled topologies of  $n$  leaves, respectively.

By using  $C_{n-1} = \binom{2n-2}{n-1}/n$  and  $|T_n| = (2n-2)!/[2^{n-1}(n-1)!]$  from equations (3) and (1), equation (21) can be rewritten

$$\text{lab}(t) = \text{or}(t) \frac{n!}{2^{n-1}},$$

which we demonstrate by induction on the size of  $t$ . Let  $t_L$  and  $t_R$  be the two root subtrees of  $t$ , with sizes  $|t_L| = L$  and  $|t_R| = R$ . Thus, for  $n \geq 2$ ,

$$(22) \quad \text{lab}(t) = \text{lab}(t_L) \text{lab}(t_R) \binom{n}{L} \frac{1}{1 + \delta_{t_L=t_R}},$$

$$(23) \quad \text{or}(t) = \text{or}(t_L) \text{or}(t_R) \frac{2}{1 + \delta_{t_L=t_R}},$$

where  $\delta_{t_L=t_R} = 1$  if  $t_L = t_R$ , and  $\delta_{t_L=t_R} = 0$  otherwise. If we insert  $\text{lab}(t_L) = \text{or}(t_L)L!/2^{L-1}$  and  $\text{lab}(t_R) = \text{or}(t_R)R!/2^{R-1}$  into equation (22), then we find

$$\begin{aligned} \text{lab}(t) &= \text{or}(t_L) \text{or}(t_R) \frac{L!R!}{2^{n-2}} \binom{n}{L} \frac{1}{1 + \delta_{t_L=t_R}} \\ &= \text{or}(t_L) \text{or}(t_R) \frac{n!}{2^{n-1}} \frac{2}{1 + \delta_{t_L=t_R}} = \text{or}(t) \frac{n!}{2^{n-1}}, \end{aligned}$$

as desired.  $\square$

The proof shows that the ratio of orderings to labelings for an unlabeled topology is independent of the unlabeled topology. Hence, because the number of root configurations of a labeled topology or ordered unlabeled topology depends only on the underlying unlabeled topology, the probability that a labeled topology chosen uniformly at random has  $\rho$  root configurations equals the probability that an ordered unlabeled topology chosen uniformly at

random has  $\rho$  root configurations. We use Lemma 3.1 to calculate the probability that a labeled topology of size  $n$  selected under the uniform distribution has  $\rho$  root configurations as the probability that an ordered unlabeled topology of size  $n$  selected under the uniform distribution has  $\rho$  root configurations.

PROPOSITION 3.2. *Let  $R_n$  be the random variable that represents the number of root configurations in an ordered unlabeled topology of size  $n$  selected uniformly at random. (i) We have  $R_1 = 0$ , and for  $n \geq 2$ ,*

$$(24) \quad R_n \stackrel{d}{=} (R_{I_n} + 1)(R_{n-I_n}^* + 1),$$

where  $I_n$  is distributed over the interval  $[1, n - 1]$  with Catalan probability  $\mathbb{P}[I_n = j] = C_{j-1}C_{n-j-1}/C_{n-1}$ ,  $R_j^*$  is an independent copy of  $R_j$  for each  $j \in [1, n - 1]$ , and both  $R_j$  and  $R_j^*$  are independent of  $I_j$  for  $j \in [1, n - 1]$ . Furthermore, (ii) the probability that a random labeled topology of size  $n$  selected under the uniform distribution has  $c_r = \rho$  root configurations can be calculated as  $\mathbb{P}[c_r = \rho] = \mathbb{P}[R_n = \rho]$ , where  $\mathbb{P}[R_n = \rho]$  has recursive formula

$$(25) \quad \mathbb{P}[R_n = \rho] = \sum_{d \in \text{Div}(\rho)} \sum_{j=1}^{n-1} \mathbb{P}[I_n = j] \mathbb{P}[R_j = d - 1] \mathbb{P}\left[R_{n-j} = \frac{\rho}{d} - 1\right],$$

where  $\text{Div}(\rho)$  denotes the set of positive integers that divide  $\rho$ ,  $\mathbb{P}[I_n = j] = C_{j-1}C_{n-j-1}/C_{n-1}$  and  $\mathbb{P}[R_n = 0] = \delta_{n,1}$ .

PROOF. The recurrence in equation (24) follows from equation (15). Observe that for a random uniform ordered unlabeled topology  $t$  of  $n$  leaves, the probability that the left (or right) root subtree of  $t$  has size  $I_n = j$  is given by  $\mathbb{P}[I_n = j] = C_{j-1}C_{n-j-1}/C_{n-1}$ , where  $C_{j-1}$ ,  $C_{n-j-1}$  and  $C_{n-1}$  give the numbers of ordered unlabeled topologies of size  $j$ ,  $n - j$  and  $n$ , respectively (Section 2.1.2). This establishes (i).

For (ii), equation (25) is a direct consequence of Lemma 3.1 and equation (24).  $\square$

We now consider the equivalence between uniformly distributed labeled histories and uniformly distributed ordered unlabeled histories.

LEMMA 3.3. *The distribution of the number of root configurations over labeled histories of size  $n$  selected uniformly at random matches the distribution of the number of root configurations over ordered unlabeled histories of size  $n$  selected uniformly at random.*

PROOF. The proof is similar to that of Lemma 3.1: we show that for each unlabeled history  $t$  of size  $n$ , we have

$$(26) \quad \frac{\text{or}(t)}{F_{n-1}} = \frac{\text{lab}(t)}{|H_n|},$$

where  $\text{or}(t)$  and  $\text{lab}(t)$  are the number of orientations of  $t$  and the number of leaf labelings of  $t$ , respectively. In other words, we prove that the uniform distribution over the set of ordered unlabeled histories of size  $n$  and the uniform distribution over the set of labeled histories of size  $n$  both induce the same probability distribution over the set of unlabeled histories of  $n$  leaves. The same property has already been shown by Lambert and Stadler ([28], page 116) following a slightly different approach.

Using  $F_{n-1} = (n - 1)!$  and  $|H_n| = n!(n - 1)!/2^{n-1}$  from equations (5) and (4), equation (26) can be rewritten

$$\text{lab}(t) = \text{or}(t) \frac{n!}{2^{n-1}},$$

which we verify by induction on  $|t|$ . Let  $t_L$  and  $t_R$  denote the two root subtrees of  $t$ , with sizes  $|t_L| = L$  and  $|t_R| = R$ . Hence, for  $n \geq 2$  we have

$$(27) \quad \text{lab}(t) = \text{lab}(t_L) \text{lab}(t_R) \binom{n}{L},$$

$$(28) \quad \text{or}(t) = 2 \text{or}(t_L) \text{or}(t_R).$$

By setting  $\text{lab}(t_L) = \text{or}(t_L)L!/2^{L-1}$  and  $\text{lab}(t_R) = \text{or}(t_R)R!/2^{R-1}$  in equation (27), we find

$$\begin{aligned} \text{lab}(t) &= \text{or}(t_L) \text{or}(t_R) \frac{L!R!}{2^{n-2}} \binom{n}{L} \\ &= \text{or}(t_L) \text{or}(t_R) \frac{2n!}{2^{n-1}} = \text{or}(t) \frac{n!}{2^{n-1}}, \end{aligned}$$

as desired.  $\square$

Next, we describe implications of Lemma 3.3 for Yule–Harding-distributed labeled topologies.

LEMMA 3.4. *The distribution of the number of root configurations over labeled topologies of size  $n$  selected according to the Yule–Harding distribution matches the distribution of the number of root configurations over ordered unlabeled histories of size  $n$  selected uniformly at random.*

PROOF. The equivalence follows from Lemma 3.3 and the fact that the uniform distribution over labeled histories of size  $n$  induces the Yule–Harding distribution on the set of labeled topologies of size  $n$  (Section 2.2).  $\square$

By Lemma 3.4, we can calculate the probability that a labeled topology of size  $n$  selected under the Yule–Harding distribution has  $\rho$  root configurations as the probability that a random uniform ordered unlabeled history of size  $n$  has  $\rho$  root configurations. In particular, we have the following proposition.

PROPOSITION 3.5. *Let  $R_n$  be the random variable that represents the number of root configurations in an ordered unlabeled history of size  $n$  selected uniformly at random. (i) We have  $R_1 = 0$ , and for  $n \geq 2$ ,*

$$(29) \quad R_n \stackrel{d}{=} (R_{I_n} + 1)(R_{n-I_n}^* + 1),$$

where  $I_n$  is uniformly distributed over the interval  $[1, n - 1]$ ,  $R_j^*$  is an independent copy of  $R_j$  for each  $j \in [1, n - 1]$ , and both  $R_j$  and  $R_j^*$  are independent of  $I_j$  for  $j \in [1, n - 1]$ . Furthermore, (ii) the probability that a random labeled topology of size  $n$  selected under the Yule–Harding distribution has  $c_r = \rho$  root configurations can be calculated as  $\mathbb{P}[c_r = \rho] = \mathbb{P}[R_n = \rho]$ , where  $\mathbb{P}[R_n = \rho]$  has recursive formula

$$(30) \quad \mathbb{P}[R_n = \rho] = \sum_{d \in \text{Div}(\rho)} \sum_{j=1}^{n-1} \mathbb{P}[I_n = j] \mathbb{P}[R_j = d - 1] \mathbb{P}\left[R_{n-j} = \frac{\rho}{d} - 1\right],$$

where  $\text{Div}(\rho)$  denotes the set of positive integers that divide  $\rho$ ,  $\mathbb{P}[I_n = j] = \frac{1}{n-1}$  and  $\mathbb{P}[R_n = 0] = \delta_{n,1}$ .

PROOF. The formula in equation (29) follows directly from equation (15) when we observe that, for a random uniform ordered unlabeled history  $t$  of  $n$  leaves, the probability that the left (or right) root subtree of  $t$  has size  $I_n = j$  is

$$\mathbb{P}[I_n = j] = \frac{F_{j-1} F_{n-j-1} \binom{n-2}{j-1}}{F_{n-1}} = \frac{1}{n-1}.$$

Equation (30) is a direct consequence of Lemma 3.4 and equation (29).  $\square$

3.2. *Equivalences with antichains of pruned binary trees.* To use results of Wagner [45] to obtain probability distributions for root configurations, we must translate between root configurations for labeled topologies and nonempty antichains for pruned binary trees.

A pruned binary tree is an ordered unlabeled topology in which the external branches—those terminating in a leaf—have been removed. If a node of the initial ordered unlabeled topology has one incident external branch, then pruning renders the node of the pruned binary tree with only one immediate descendant; a node with two incident external branches is pruned to possess no immediate descendants. To illustrate the pruning operation, consider the ordered unlabeled topology depicted on the left of Figure 2A and assign arbitrary labels to all its nodes, as in Figure 1A. The leaf labels of the pruned binary tree resulting from this process can be described by the Newick format  $((g, h), i)$ . Note that pruned binary trees have their left–right orientation induced by the overlying ordered unlabeled topology.

If  $t$  is an ordered unlabeled topology of size  $n$  and  $\tilde{t}$  is its associated pruned binary tree of  $n - 1$  nodes, then we can consider  $\tilde{t}$  as the Hasse diagram of a partially ordered set with ground set given by the nodes of  $\tilde{t}$ —the *internal* nodes of  $t$ —and order relation determined by the descendant–ancestor relationship in  $\tilde{t}$ . An antichain of  $\tilde{t}$  is a subset of its nodes such that no two elements in the subset are comparable by the order relation. For instance, the two-element antichains of pruned binary tree  $((g, h), i)$  in Figure 1A are  $\{g, h\}$ ,  $\{g, i\}$ ,  $\{h, i\}$  and  $\{j, i\}$ .

The nonempty antichains of the pruned binary tree  $\tilde{t}$  bijectively correspond to the root configurations of the overlying ordered unlabeled topology  $t$ : omitting leaves from a root configuration of  $t$  yields an antichain of  $\tilde{t}$ , and adding leaves to an antichain of  $\tilde{t}$  so that each leaf of  $t$  is either represented or has one of its ancestral nodes represented yields a root configuration of  $t$ .

For instance, consider the set in equation (10) of the root configurations of the ordered unlabeled topology in Figure 1A. By omitting leaves from each configuration, we obtain the antichains of  $\tilde{t}$ :

$$\{\{j, i\}, \{j\}, \{g, h, i\}, \{g, h\}, \{h, i\}, \{h\}, \{g, i\}, \{g\}, \{i\}, \emptyset\}.$$

We make a substitution of the empty antichain  $\emptyset$  that emerges from the root configuration consisting of all the leaves by the antichain  $\{k\}$  consisting only of the root of  $\tilde{t}$ ; we have then bijectively paired all root configurations of  $t$  and all nonempty antichains of  $\tilde{t}$ . Using this correspondence, we have the next result.

LEMMA 3.6. *The distribution of the number of root configurations over labeled topologies of size  $n$  selected uniformly at random matches the distribution of the number of nonempty antichains over the set of  $(n - 1)$ -node pruned binary trees selected uniformly at random.*

PROOF. By Lemma 3.1, the number of root configurations has the same distribution when considered over uniformly distributed labeled topologies of size  $n$  or over uniformly distributed ordered unlabeled topologies of size  $n$ . By the correspondence between antichains

of pruned binary trees with  $n - 1$  nodes and root configurations of associated ordered unlabeled topologies of size  $n$ , the distribution of the number of root configurations over uniformly distributed ordered unlabeled topologies of size  $n$  matches the distribution of the number of nonempty antichains over uniformly distributed pruned binary trees with  $n - 1$  nodes.  $\square$

**4. Root configurations under the uniform distribution on labeled topologies.** Disanto and Rosenberg [14] determined the mean and variance of the number of root configurations for uniformly distributed labeled topologies of size  $n$  (Section 2.4.4). In this section, we use the correspondence with antichains given in Section 3.2 to show that the logarithm of the number of root configurations for uniformly distributed labeled topologies of size  $n$ , suitably rescaled, converges to a normal distribution.

Wagner ([45], Section 2.3.2) studied the number  $a(t)$  of nonempty antichains of a randomly selected pruned binary tree  $t$  of given size. For a pruned binary tree of  $n$  nodes selected uniformly at random, he considered  $\log a(t)$ , showing that  $(\log a - \mathbb{E}_n[\log a]) / \sqrt{\mathbb{V}_n[\log a]}$  converges to a standard normal distribution as  $n \rightarrow \infty$ , where  $\mathbb{E}_n[\log a] \sim \mu n$  and  $\mathbb{V}_n[\log a] \sim \sigma^2 n$ , with constants  $(\mu, \sigma^2) \approx (0.272, 0.034)$ .

By Lemma 3.6, Wagner’s variable  $\log a$  asymptotically has the same distribution as the variable  $\log c_r$  considered over uniformly distributed labeled topologies of size  $n + 1$ . We thus have the following result.

**PROPOSITION 4.1.** *The logarithm of the number of root configurations in a labeled topology of size  $n$  selected uniformly at random, rescaled as  $(\log c_r - \mathbb{E}_n[\log c_r]) / \sqrt{\mathbb{V}_n[\log c_r]}$ , converges to a standard normal distribution, where  $\mathbb{E}_n[\log c_r] \sim \mu n$  and  $\mathbb{V}_n[\log c_r] \sim \sigma^2 n$ ,  $(\mu, \sigma^2) \approx (0.272, 0.034)$ .*

The result gives an asymptotic lognormal distribution for the number of root configurations of a labeled topology of size  $n$  selected uniformly at random. Although we do not expect  $e^{\mathbb{E}_n[\log c_r]}$  and  $e^{\sigma_n[\log c_r]}$  to agree with  $\mathbb{E}_n[c_r]$  and  $\sigma_n[c_r]$ , for the mean we see that in the  $n \rightarrow \infty$  limit,  $e^{\mathbb{E}_n[\log c_r]} \approx e^{0.272n} \approx 1.313^n$ , numerically close to the exponential growth of  $\mathbb{E}_n[c_r]$ , or  $(4/3)^n$  (equation (16)). For the standard deviation,  $e^{\sigma_n[\log c_r]} \approx e^{\sqrt{0.034n}} \approx 1.202^n$  is not as close to the exponential growth of  $\sigma_n[c_r]$  from equation (18), which gives  $[2/\sqrt{7(8\sqrt{2} - 11)}]^n \approx 1.350^n$ .

For fixed  $n$ , we can compute the exact distribution of  $c_r$  and  $\log c_r$  under a uniform distribution across labeled topologies of size  $n$ , as described in Proposition 3.2(ii). Figure 5 shows the cumulative distribution  $\mathbb{P}[\log c_r \leq \mathbb{E}[\log c_r] + y\sigma[\log c_r]]$  as a function of  $y$ , when labeled topologies are selected uniformly at random among the  $2.13 \times 10^{14}$  labeled topologies with 15 leaves. To obtain the distribution, we can count root configurations for arbitrary labelings of each of the 4850 unlabeled topologies with 15 leaves, and then count labelings for each unlabeled topology ([39], page 47). Already for small tree size, the figure shows that the exact cumulative distribution is close to the cumulative distribution of a Gaussian random variable with mean 0 and variance 1.

**5. Root configurations under the Yule–Harding distribution on labeled topologies.** We next study distributional properties of the number of root configurations for labeled topologies selected under the Yule–Harding probability model. Section 2.2 noted that this model assigns higher probability to trees with a high degree of balance compared to that assigned by the uniform model; Section 2.4.4 noted that balanced trees have high numbers of root configurations relative to unbalanced trees. We therefore find that the mean number of root configurations for labeled topologies of size  $n$  grows exponentially faster under

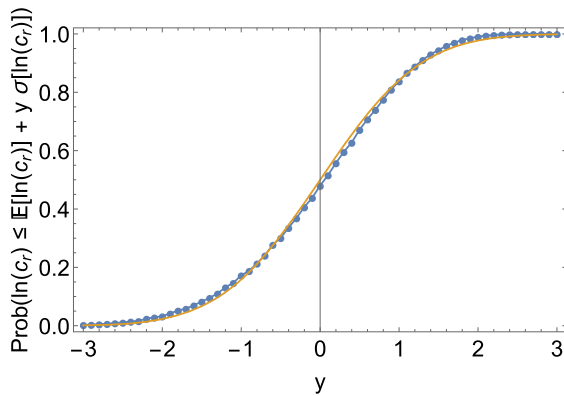


FIG. 5. Cumulative distribution of the natural logarithm of the number of root configurations for uniformly distributed labeled topologies of size  $n = 15$  (dotted line). Each dot has its abscissa determined by a value of  $y$  ranging in the interval  $y \in [-3, 3]$  in steps of 0.1. Given  $y$ , the quantity plotted is the probability that a labeled topology with  $n = 15$  chosen uniformly at random has a number of root configurations less than or equal to  $\exp(\mathbb{E}[\log c_r] + y\sigma[\log c_r])$ , where  $\mathbb{E}[\log c_r]$  and  $\sigma[\log c_r]$  are respectively the mean and standard deviation of the logarithm of the number of root configurations for uniformly distributed labeled topologies with  $n = 15$  leaves (Proposition 4.1). The solid line is the cumulative distribution of a Gaussian random variable with mean 0 and variance 1.

the Yule–Harding model than under the uniform model. The variance of the number of root configurations also has faster growth.

Note that in the main results of the section—Propositions 5.2, 5.4 and 5.5—expectations  $\mathbb{E}_n$  and variances  $\mathbb{V}_n$  are taken with respect to the Yule–Harding distribution.

5.1. *Lognormal distribution of the number of root configurations.* We begin the analysis of the number of root configurations under the Yule–Harding distribution by showing that the logarithm of the number of root configurations of a Yule–Harding random labeled topology of size  $n$ , when suitably rescaled, converges to a standard normal distribution.

The results in this section are obtained by considering root configurations over ordered unlabeled histories of given size selected under the uniform distribution. Owing to Lemma 3.4, we can demonstrate that the number of root configurations in a Yule–Harding random labeled topology of size  $n$  asymptotically follows a lognormal distribution by showing that the number of root configurations is asymptotically lognormally distributed when considered over the set of uniformly distributed ordered unlabeled histories of  $n$  leaves. We use a result of Wagner [45] for additive tree parameters of ordered unlabeled histories. We first must verify a technical condition for the mean of the random variable  $\log(1 + 1/c_r)$ , considered over uniformly distributed ordered unlabeled histories. This verification proceeds by considering cherry nodes [31], internal nodes whose two immediate descendant nodes are leaves.

LEMMA 5.1. *For uniformly distributed ordered unlabeled histories of size  $n$ , the mean value  $\mathbb{E}_n[\log(1 + 1/c_r)]$  of the random variable  $\log(1 + 1/c_r)$  converges to 0 exponentially fast as  $n$  increases. In particular,*

$$(31) \quad \mathbb{E}_n \left[ \log \left( 1 + \frac{1}{c_r} \right) \right] = \mathcal{O}(0.9^n).$$

PROOF. To show that  $\mathbb{E}_n[\log(1 + 1/c_r)]$  has exponential growth  $\mathcal{O}(0.9^n)$  for an ordered unlabeled history  $t$  of size  $n$  selected uniformly at random, we consider the mean value  $\mathbb{E}_n[2^{-\text{ch}}]$  of the random variable  $2^{-\text{ch}}$ —where  $\text{ch}$  is the number of cherries in  $t$ . We claim



that

$$(32) \quad \mathbb{E}_n[2^{-\text{ch}}] = \mathcal{O}(0.9^n).$$

For a tree  $t$  with  $|t| \geq 3$ ,  $c_r(t) \geq 2^{\text{ch}(t)}$ , as each cherry node generates a pair of ancestral configurations: the configuration corresponding to the node, and the configuration corresponding to its pair of leaves. At the root node, a root configuration can be obtained by choosing ancestral configurations at each of the cherry nodes and augmenting the configuration with leaves that do not descend from cherry nodes.

Noting  $\log(1 + x) \leq x$  for  $x > 0$ , for each ordered unlabeled history  $t$  with size  $|t| \geq 3$ , we have

$$\log\left[1 + \frac{1}{c_r(t)}\right] \leq \frac{1}{c_r(t)} \leq 2^{-\text{ch}(t)}.$$

By taking expectations, we see that equation (32) implies equation (31):

$$\mathbb{E}_n\left[\log\left(1 + \frac{1}{c_r}\right)\right] \leq \mathbb{E}_n[2^{-\text{ch}}].$$

It remains to verify equation (32). In their Theorem 2, Disanto and Wiehe [18] studied the generating function  $F(x, z)$  counting the number of unlabeled histories  $t$  of size  $n$  with a given number of cherries, where each unlabeled history  $t$  is weighted by its probability  $2^{n-1-\text{ch}(t)}/(n-1)!$  under the Yule–Harding distribution:

$$F(x, z) = \sum_t \frac{2^{n-1-\text{ch}(t)}}{(n-1)!} x^{\text{ch}(t)} z^n.$$

The sum proceeds over unlabeled histories (“ranked trees” in Disanto and Wiehe [18]). The coefficient of  $x^h z^n$  in  $F(x, z)$  gives the probability of  $h$  cherries in unlabeled histories of size  $n$  under the Yule–Harding distribution, or equivalently, the probability of  $h$  cherries in ordered unlabeled histories of size  $n$  selected uniformly at random. Hence, the expectation  $\mathbb{E}_n[2^{-\text{ch}}]$  is obtained from the coefficient of  $z^n$  in  $F(\frac{1}{2}, z)$ . From Disanto and Wiehe [18],

$$F\left(\frac{1}{2}, z\right) = f(z) = \frac{ze^{z\sqrt{2}} - z}{(\sqrt{2} - 2)e^{z\sqrt{2}} + 2 + \sqrt{2}}.$$

By Theorem IV.7 of Flajolet and Sedgewick [23] (see also Section 2.3),  $\mathbb{E}_n[2^{-\text{ch}}]$  grows exponentially like  $[z^n]f(z) \asymp \alpha^{-n}$ , where  $\alpha$  is the dominant singularity of  $f(z)$ . The value of  $\alpha$  is the solution of smallest modulus of the equation  $(\sqrt{2} - 2)e^{z\sqrt{2}} + 2 + \sqrt{2} = 0$ , whose left-hand side is the denominator of  $f(z)$ . Because

$$\alpha = \frac{1}{\sqrt{2}} \log\left(\frac{2 + \sqrt{2}}{2 - \sqrt{2}}\right) = \frac{\sqrt{2} \log(3 + 2\sqrt{2})}{2} \approx 1.246,$$

$\alpha^{-1} \approx 0.802$ , and thus, conservatively,  $\mathbb{E}_n[2^{-\text{ch}}] = \mathcal{O}(0.9^n)$ . Hence,  $\mathbb{E}_n[\log(1 + \frac{1}{c_r})]$  also decays to 0 as  $\mathcal{O}(0.9^n)$ .  $\square$

Considering as in Section 2.5 the additive tree parameter  $F(t) = \log[c_r(t) + 1]$ , by Lemma 5.1 we have demonstrated that the associated toll function  $f(t) = \log[1 + 1/c_r(t)]$  satisfies

$$(33) \quad \frac{\sum_t f(t)}{F_{n-1}} = \mathbb{E}_n\left[\log\left(1 + \frac{1}{c_r}\right)\right] = \mathcal{O}(0.9^n),$$

where the sum proceeds over all  $(n-1)!$  ordered unlabeled histories  $t$  of size  $n$  (equation (5)). Equation (33), together with the fact that  $f(t)$  is bounded because  $c_r(t) \geq 1$  for  $|t| \geq 2$ ,

show that the hypotheses of Theorem 4.2 of Wagner [45] are satisfied. By applying the theorem, we can conclude that for an ordered unlabeled history  $t$  of size  $n$  selected uniformly at random, the standardized version of the random variable  $F(t) = \log[c_r(t) + 1]$  converges asymptotically to a normal distribution with mean 0 and variance 1. By the same theorem, the mean and variance of  $F(t) = \log[c_r(t) + 1]$  grow respectively like  $\mu n$  and  $\sigma^2 n$ , for two constants

$$(34) \quad \mu = \sum_t \frac{2f(t)}{(|t| + 1)!} \approx 0.351,$$

$$(35) \quad \sigma^2 = \sum_t \frac{2f(t)[2F(t) - f(t)]}{(|t| + 1)!} - \mu^2 + \sum_{t_1} \sum_{t_2} \frac{4f(t_1)f(t_2)}{(|t_1| + 1)!(|t_2| + 1)!} \\ \times \left[ \frac{(|t_1| - 1)(|t_2| - 1)}{|t_1| + |t_2| - 1} - |t_1| - |t_2| + 2 + \frac{(|t_1| - 1)(|t_2| - 1)}{(|t_1| + |t_2|)(|t_1| + |t_2| + 1)} \right. \\ \left. + \frac{(|t_1| - 1)^2(|t_2| - 1)^2}{(|t_1| + |t_2| - 1)(|t_1| + |t_2|)(|t_1| + |t_2| + 1)} \right] \\ \approx 0.008.$$

Note that the sums in equations (34) and (35) are defined over all ordered unlabeled histories, but that the approximations have been calculated by disregarding histories of size strictly larger than 15 and 12 in the sums for  $\mu$  and  $\sigma^2$ , respectively. The equivalence of Lemma 3.4 between the distribution of the number of root configurations over uniformly distributed ordered unlabeled histories and the distribution of the number of root configurations over Yule–Harding distributed labeled topologies, coupled with the fact that the difference  $\log(c_r + 1) - \log c_r = \log(1 + 1/c_r)$  is small, finally yields the following proposition.

**PROPOSITION 5.2.** *The logarithm of the number of root configurations in a labeled topology of size  $n$  selected under the Yule–Harding distribution, rescaled as  $(\log c_r - \mathbb{E}_n[\log c_r]) / \sqrt{\mathbb{V}_n[\log c_r]}$ , converges to a standard normal distribution, where  $\mathbb{E}_n[\log c_r] \sim \mu n$  and  $\mathbb{V}_n[\log c_r] \sim \sigma^2 n$  for  $(\mu, \sigma^2) \approx (0.351, 0.008)$ .*

For fixed  $n$ , we can compute the exact distribution of  $c_r$  (and  $\log c_r$ ) under the Yule–Harding distribution across all labeled topologies of size  $n$  as in Proposition 3.5(ii). Similar to the computations in Figure 5, we can weight the counts of root configurations for unlabeled topologies by their Yule–Harding probabilities ([39], page 47). Figure 6 shows the cumulative distribution  $\mathbb{P}[\log c_r \leq \mathbb{E}[\log c_r] + y\sigma[\log c_r]]$  plotted as a function of  $y$ , when labeled topologies of size  $n = 15$  are selected under the Yule–Harding distribution. The distribution is close to the cumulative distribution of a Gaussian random variable with mean 0 and variance 1.

**5.2. Mean number of root configurations.** In Section 5.1, we have analyzed distributional properties of the logarithm of the number of root configurations considered over labeled topologies of given size selected under the Yule–Harding distribution. In this section, we study the mean number of root configurations under the Yule–Harding distribution.

From Lemma 3.4, the mean number of root configurations in a random labeled topology of size  $n$  selected under the Yule–Harding distribution is also the mean number of root configurations in a uniform random ordered unlabeled history of  $n$  leaves. To calculate this mean, we use the distributional recurrence in Proposition 3.5 for the variable  $R_n$ , and by applying generating functions and singularity analysis, we obtain the following result.

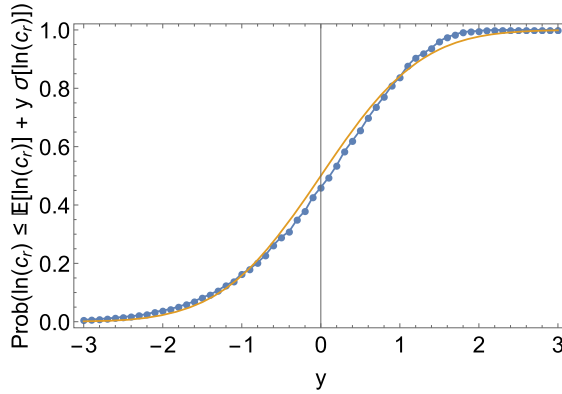


FIG. 6. Cumulative distribution of the natural logarithm of the number of root configurations for labeled topologies of size  $n = 15$  considered under the Yule–Harding distribution (dotted line). Each dot has its abscissa determined by a value of  $y$  ranging in the interval  $y \in [-3, 3]$  in steps of 0.1. Given  $y$ , the quantity plotted is the probability that a labeled topology with  $n = 15$  chosen at random under the Yule–Harding distribution has a number of root configurations less than or equal to  $\exp(\mathbb{E}[\log c_r] + y\sigma[\log c_r])$ , where  $\mathbb{E}[\log c_r]$  and  $\sigma[\log c_r]$  are respectively the mean and the standard deviation of the logarithm of the number of root configurations for Yule–Harding distributed labeled topologies of  $n = 15$  leaves (Proposition 5.2). The solid line is the cumulative distribution of a Gaussian random variable with mean 0 and variance 1.

PROPOSITION 5.3. The mean number of root configurations in an ordered unlabeled history of size  $n$  selected uniformly at random satisfies the asymptotic relation  $\mathbb{E}[R_n] \sim k_e^n$ , where  $k_e = 1/(1 - e^{-2\pi\sqrt{3}/9})$ .

PROOF. Set  $e_n \equiv \mathbb{E}[R_n]$ . Then  $\mathbb{E}[R_{I_n} R_{n-I_n}^*] = \sum_{j=1}^{n-1} \mathbb{P}[I_n = j] \mathbb{E}[R_j R_{n-j}^*] = \frac{1}{n-1} \times \sum_{j=1}^{n-1} \mathbb{E}[R_j] \mathbb{E}[R_{n-j}^*]$ . Proposition 3.5 yields for  $n \geq 2$  the recurrence

$$(36) \quad e_n = 1 + \frac{1}{n-1} \sum_{j=1}^{n-1} e_j e_{n-j} + \frac{2}{n-1} \sum_{j=1}^{n-1} e_j,$$

with initial condition  $e_1 = 0$ .

Defining the generating function

$$(37) \quad E(z) \equiv \sum_{n=1}^{\infty} e_n z^n = z^2 + 2z^3 + \frac{10}{3}z^4 + \frac{31}{6}z^5 + \dots,$$

the recurrence in equation (36) translates into the Riccati differential equation

$$(38) \quad zE'(z) = E(z)^2 + \frac{1+z}{1-z}E(z) + \frac{z^2}{(1-z)^2},$$

with initial condition  $E(0) = 0$ . To obtain the differential equation, we have multiplied both sides of equation (36) by  $(n-1)z^n$ , summed for  $n \geq 1$ , and then used the facts that  $\sum_{n=1}^{\infty} (n-1)e_n z^n = zE'(z) - E(z)$ ,  $\sum_{n=1}^{\infty} (n-1)z^n = z^2[1/(1-z)]' = z^2/(1-z)^2$ ,  $\sum_{n=1}^{\infty} (\sum_{j=1}^{n-1} e_j e_{n-j})z^n = E(z)^2$  and  $\sum_{n=1}^{\infty} (\sum_{j=1}^{n-1} e_j)z^n = E(z)[1/(1-z) - 1]$ .

Solving the differential equation yields

$$(39) \quad E(z) = \frac{2z \sin(\frac{\sqrt{3}}{2} \log(1-z))}{(z-1)[\sqrt{3} \cos(\frac{\sqrt{3}}{2} \log(1-z)) + \sin(\frac{\sqrt{3}}{2} \log(1-z))]}.$$

$E(z)$  has infinitely many singularities. The singularity of  $E(z)$  with smallest modulus occurs at  $z = \alpha \equiv 1 - e^{-2\pi\sqrt{3}/9} \approx 0.702$ . The singularity of smallest modulus is obtained by setting to 0 the factor

$$(40) \quad \sqrt{3} \cos\left[\frac{\sqrt{3}}{2} \log(1 - z)\right] + \sin\left[\frac{\sqrt{3}}{2} \log(1 - z)\right]$$

appearing in the denominator of equation (39). The expansion of  $E(z)$  at its dominant singularity  $z = \alpha$  looks like

$$E(z) \stackrel{z \rightarrow \alpha}{\sim} \frac{1}{1 - \frac{z}{\alpha}},$$

which can be obtained by plugging the Taylor expansion  $-\sqrt{3}e^{+2\pi\sqrt{3}/9}(z - \alpha)$  of the factor (40) in the denominator of equation (39). By Theorem VI.4 of Flajolet and Sedgewick [23] (see also Section 2.3), we finally obtain

$$[z^n]E(z) \sim [z^n]\left(\frac{1}{1 - \frac{z}{\alpha}}\right) = \alpha^{-n},$$

as  $n \rightarrow \infty$ .  $\square$

The next proposition follows immediately from Proposition 5.3 and Lemma 3.4.

**PROPOSITION 5.4.** *The mean number of root configurations in a labeled topology of size  $n$  selected at random under the Yule–Harding distribution has asymptotic growth  $\mathbb{E}_n[c_r] \sim k_e^n$ , where  $k_e = 1/(1 - e^{-2\pi\sqrt{3}/9}) \approx 1.42538682$ . Furthermore, the mean total number of configurations has asymptotic growth  $\mathbb{E}_n[c] \asymp \mathbb{E}_n[c_r]$ .*

For small tree size ( $n \leq 20$ ), Figure 7 plots the mean number of root configurations for a random tree of size  $n$  selected under the Yule–Harding distribution as a function of the corresponding mean under the uniform distribution. The plot provides a numerical visualization of the similar behavior of the numbers of root configurations under the Yule–Harding and uniform distributions. The mean is greater for the Yule–Harding distribution, but the two quantities are highly correlated, with Pearson’s correlation coefficient approximately 0.995.

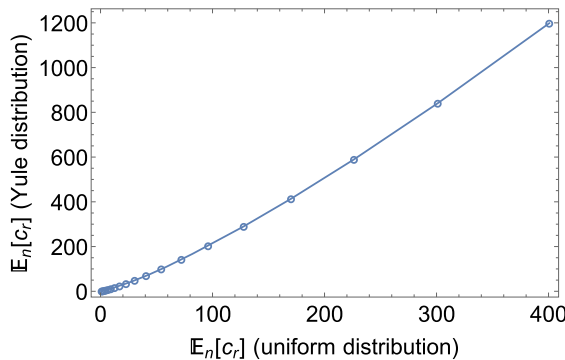


FIG. 7. Mean number of root configurations of labeled topologies of size  $n$  under the Yule–Harding and uniform distributions, for  $2 \leq n \leq 20$ . Values for the uniform distribution are computed from the power series expansion of equation (33) of Disanto and Rosenberg [14]; values for Yule–Harding are computed from the power series expansion of equation (39).

5.3. *Variance of the number of root configurations.* In this section, we analyze the asymptotic growth of the variance of the number of root configurations under the Yule–Harding distribution. In particular, by using Lemma 3.4, we study the variance of the number of root configurations in a uniform random ordered unlabeled history of size  $n$ .

Following Section 5.2 and squaring equation (29), we obtain a recurrence for  $s_n \equiv \mathbb{E}[R_n^2]$ . For  $n \geq 2$ ,

$$(41) \quad \begin{aligned} s_n = 1 + \frac{1}{n-1} \sum_{j=1}^{n-1} s_j s_{n-j} + \frac{2}{n-1} \sum_{j=1}^{n-1} s_j + \frac{4}{n-1} \sum_{j=1}^{n-1} s_j e_{n-j} \\ + \frac{4}{n-1} \sum_{j=1}^{n-1} e_j e_{n-j} + \frac{4}{n-1} \sum_{j=1}^{n-1} e_j, \end{aligned}$$

with initial condition  $s_1 = 0$ .

Starting from this recurrence, a symbolic calculation similar to that used to derive equation (38) shows that the generating function  $S(z) \equiv \sum_{n=1}^{\infty} s_n z^n = z^2 + 4z^3 + \frac{34}{3}z^4 + \frac{55}{2}z^5 \dots$  satisfies the Riccati differential equation

$$(42) \quad zS'(z) = S(z)^2 - S(z) \left[ \frac{1+z}{z-1} - 4E(z) \right] + \frac{[z - 2(z-1)E(z)]^2}{(z-1)^2}.$$

This equation can be written

$$(43) \quad S'(z) = g_2(z)S(z)^2 + g_1(z)S(z) + g_0(z)$$

by setting

$$(g_2(z), g_1(z), g_0(z)) \equiv \left( \frac{1}{z}, \left( 4E(z) - \frac{1+z}{z-1} \right) \frac{1}{z}, \frac{[z - 2(z-1)E(z)]^2}{z(z-1)^2} \right).$$

By substituting  $U(z) \equiv \exp[\int_0^z S(x)/(-x) dx]$ , we obtain  $S(z) = -zU'(z)/U(z)$ , and equation (43) can be rewritten as a second-order linear differential equation

$$(44) \quad U''(z) - \left( g_1(z) + \frac{g_2'(z)}{g_2(z)} \right) U'(z) + g_2(z)g_0(z)U(z) = 0.$$

The coefficients of equation (44) are analytic functions for  $|z| < 0.702$ , with a removable singularity at  $z = 0$  as the expansion (37) of  $E(z)$  starts with a quadratic nonzero term. Using existence results for the solutions of second-order ordinary differential equations,  $U(z)$  must be analytic for  $|z| < 0.702$ , the constant being the radius of convergence of  $E(z)$  as determined in the proof of Proposition 5.3. Therefore, also  $U'(z)$  is analytic for  $|z| < 0.702$ , and thus  $S(z)$  is a meromorphic function on this domain, being a quotient of two analytic functions. To analyze the singularities of a meromorphic function, one must locate the possible roots of its denominator function. In our case, the set of singularities of  $S(z)$  consists of the roots of  $U(z)$ . In particular, by studying in the Appendix the function  $U(z)$  in  $\mathcal{B} \equiv \{z \in \mathbb{C} : |z| \leq \frac{1}{2}\}$ , we find that  $S(z)$  has a unique dominant singularity  $\alpha \approx 0.4889986317$ , the unique and simple root of  $U(z)$  within  $\mathcal{B}$  (Proposition A.6).

As a consequence, we can write  $U(z) = (z - \alpha)\tilde{U}(z)$ , with  $\tilde{U}(\alpha) \neq 0$  and  $U'(\alpha) = (-\alpha)\tilde{U}'(\alpha) \neq 0$ . Therefore, for  $z \rightarrow \alpha$  the generating function  $S(z)$  admits the expansion

$$\begin{aligned} S(z) &= \frac{-zU'(z)}{U(z)} \underset{z \rightarrow \alpha}{\sim} \frac{(-\alpha)[U'(\alpha) + U''(\alpha)(z - \alpha) + \dots]}{U(\alpha) + U'(\alpha)(z - \alpha) + \dots} \underset{z \rightarrow \alpha}{\sim} \frac{(-\alpha)U'(\alpha)}{U'(\alpha)(z - \alpha)} \\ &= \frac{-\alpha}{z - \alpha} = \frac{1}{1 - \frac{z}{\alpha}}. \end{aligned}$$

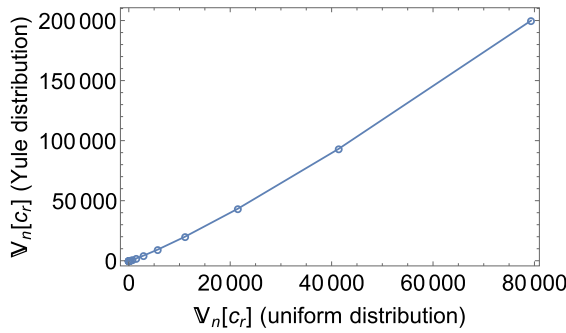


FIG. 8. Variance of the number of root configurations of labeled topologies of size  $n$  under the Yule–Harding and uniform distributions, for  $2 \leq n \leq 20$ . Values for the uniform distribution are computed from the power series expansion of equation (39) of Disanto and Rosenberg [14]; values for Yule–Harding are computed from equations (41) and (36).

From Theorem VI.4 of Flajolet and Sedgewick [23] (see also Section 2.3), we can thus recover the asymptotic growth of the associated coefficients

$$(45) \quad \mathbb{E}[R_n^2] = [z^n]S(z) \sim [z^n] \left( \frac{1}{1 - \frac{z}{\alpha}} \right) = \alpha^{-n},$$

and hence derive the asymptotic growth of the variance  $\mathbb{V}[R_n]$ . In particular, we have the following result.

**PROPOSITION 5.5.** *The variance of the number of root configurations in a labeled topology of size  $n$  selected at random under the Yule–Harding distribution has asymptotic growth  $\mathbb{V}_n[c_r] \sim k_v^n$ , where  $k_v \approx 2.0449954971$ . Furthermore, the variance of the total number of configurations has asymptotic growth  $\mathbb{V}_n[c] \asymp \mathbb{V}_n[c_r]$ .*

**PROOF.** For uniformly distributed ordered unlabeled histories of size  $n$ , equation (45) yields  $\mathbb{E}[R_n^2] \sim k_v^n$ ,  $k_v \equiv 1/\alpha \approx 2.0449954971$ . From Proposition 5.3,  $\mathbb{E}[R_n]^2 \sim (k_e^2)^n$ , with  $k_e^2 \approx 2.03$ . Because  $k_v > k_e^2$ , as  $n \rightarrow \infty$  we obtain

$$\mathbb{V}[R_n] = \mathbb{E}[R_n^2] - \mathbb{E}[R_n]^2 \sim k_v^n.$$

By Lemma 3.4, the variance of the variable  $R_n$  is the variance of the number of root configurations considered over labeled topologies of  $n$  leaves selected under the Yule–Harding distribution.  $\square$

As we did for the mean, we numerically visualize the similarity in variance of the number of root configurations for trees of size  $n$  selected at random under the Yule–Harding and uniform distributions. For small tree size ( $n \leq 20$ ), we plot in Figure 8 the variance of the number of root configurations for a random tree of size  $n$  selected under the Yule–Harding distribution as a function of the variance of the number of root configurations for a random uniform tree of the same size. As was true of the mean, the Yule–Harding and uniform distributions on labeled topologies give correlated variances (correlation coefficient 0.997).

**6. Discussion.** Considering gene trees and species trees with a matching labeled topology  $G = S = t$ , we have studied distributional properties of the number  $c_r$  of root ancestral configurations for labeled topologies  $t$  of fixed size under two probability models, the uniform model and the Yule–Harding model (Table 1). We have made use of techniques of analytic combinatorics, relying on equivalences across tree types (Section 3) and making particular

TABLE 1  
*Distributional properties of the number of root and total configurations*

Results		Uniform model		Yule–Harding model	
Root configurations	Mean	$\mathbb{E}_n[c_r] \sim 1.225 \cdot 1.333^n$	Equation (16)	$\mathbb{E}_n[c_r] \sim 1.425^n$	Proposition 5.4
	Variance	$\mathbb{V}_n[c_r] \sim 1.405 \cdot 1.822^n$	Equation (18)	$\mathbb{V}_n[c_r] \sim 2.045^n$	Proposition 5.5
	Lognormal distribution	$\mathbb{E}_n[\log c_r] \sim 0.272 \cdot n$ $\mathbb{V}_n[\log c_r] \sim 0.034 \cdot n$	Proposition 4.1 Proposition 4.1	$\mathbb{E}_n[\log c_r] \sim 0.351 \cdot n$ $\mathbb{V}_n[\log c_r] \sim 0.008 \cdot n$	Proposition 5.2 Proposition 5.2
Total configurations	Mean	$\mathbb{E}_n[c] \asymp 1.333^n$	Equation (17)	$\mathbb{E}_n[c] \asymp 1.425^n$	Proposition 5.4
	Variance	$\mathbb{V}_n[c] \asymp 1.822^n$	Equation (19)	$\mathbb{V}_n[c] \asymp 2.045^n$	Proposition 5.5

use of results of Wagner [45] on distributional properties of additive tree parameters for several families of trees.

Extending results of Disanto and Rosenberg [14], for the uniform model we have shown that the logarithm of the number of root configurations, when standardized, converges asymptotically to a standard normal distribution (Proposition 4.1). Under the Yule–Harding distribution, as is the case for uniformly distributed labeled topologies, the logarithm of the number of root configurations, when standardized, converges to a standard normal distribution (Proposition 5.2). The study produces the first results on asymptotic distributions under the uniform or Yule–Harding models for ancestral configurations, and further, for any of the recently studied combinatorial quantities that require consideration of both gene trees and species trees—ancestral configurations [14, 46], coalescent histories [2, 8, 11–13, 26, 33–35, 42], compact coalescent histories [16, 47], deep coalescence costs [29, 30, 41, 43, 44], history classes [36], nonequivalent ancestral configurations [15, 46] and ranked histories [9, 10, 37].

We have also determined the asymptotic growth of the mean and the variance of the number of root configurations, finding that under the Yule–Harding model,  $\mathbb{E}_n[c_r] \sim 1.425^n$  (Proposition 5.4) and  $\mathbb{V}_n[c_r] \sim 2.045^n$  (Proposition 5.5). As  $\mathbb{E}_n[c] \asymp \mathbb{E}_n[c_r]$  and  $\mathbb{V}_n[c] \asymp \mathbb{V}_n[c_r]$ , we also recover the exponential growth rate of the mean and the variance of the total number of configurations under the Yule–Harding model. These results were obtained by use of recursions to obtain Riccati differential equations for generating functions (equations (38) and (42)). For the case of the mean, the Riccati equation was solvable (equation (39)); for the variance, although the equation was not solvable, the asymptotic growth was nevertheless possible to obtain. Our method introduced for this case has potential for broader application, as many problems involving various types of trees and other combinatorial structures can lead to related Riccati equations [5, 22, 24].

Both the mean and the variance across labeled topologies of the number of ancestral configurations are empirically highly correlated between the uniform and Yule–Harding models (Figures 7 and 8). Alongside the results of Disanto and Rosenberg [14] for the uniform case, the larger values for Yule–Harding (Table 1) suggest a role for tree balance in predicting the number of root configurations. By considering a representative labeling for each unlabeled topology of size  $n = 15$ , in Figure 9 we plot on a logarithmic scale the number of root configurations as a function of the number of labeled histories, the latter calculated with equation (6). The numerical illustration in the figure shows that empirically, the two quantities are correlated: highly balanced labeled topologies—which tend to have larger numbers of labeled histories (Section 2.2)—in general have larger numbers of root configurations.

In particular, the largest number of root configurations is possessed by the balanced labeled topology depicted in Figure 4C, which also has the largest number of labeled histories, 2,745,600. The trend in this example is confirmed by our asymptotic results. Under the Yule–Harding probability model, which gives more weight to balanced labeled topologies



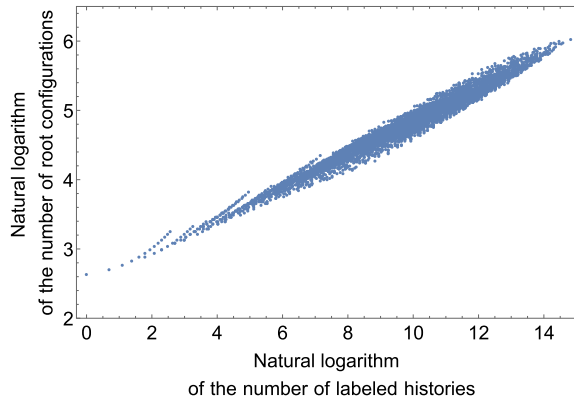


FIG. 9. Natural logarithm of the number of root configurations and natural logarithm of the number of labeled histories for a representative labeling of each unlabeled topology of size  $n = 15$ . The number of points plotted is 4850, the number of unlabeled topologies with  $n = 15$  leaves. The Pearson correlation is approximately 0.987 (0.784 without log scaling).

than does the uniform model, the mean number of root configurations and the mean total number of configurations grow exponentially faster than under the uniform distribution (Table 1). This differing behavior also accords with the proof of Disanto and Rosenberg [14] that balanced and caterpillar trees respectively possess the largest and smallest numbers of root configurations for fixed tree size (Section 2.4.3).

Several directions and extensions naturally arise from our work. First, we focused on root rather than total configurations; although some results for total configurations follow quickly (Table 1), we did not consider total configurations in detail. Second, we assumed that the gene tree and species tree had the same labeled topology, and we did not study nonmatching gene trees and species trees. The nonmatching case merits further analysis, as a nonmatching gene tree labeled topology can have more root and total configurations than the topology that matches the species tree [14]. Third, ancestral configurations can be considered up to an equivalence relationship that accounts for symmetries in gene trees [46]. The resulting equivalence classes—the nonequivalent ancestral configurations—are used for calculating probabilities of gene trees in STELLS [46], with computational complexity that depends on the number of these classes. Some investigation of this number has been carried out by Disanto and Rosenberg [15] for uniformly distributed matching gene trees and species trees. It would be of interest to see whether the techniques we have used could derive distributional properties of the number of nonequivalent ancestral configurations under the uniform and Yule–Harding probability models.

#### APPENDIX. THE FUNCTION $U(z)$ HAS A UNIQUE AND SIMPLE ROOT OF SMALLEST MODULUS

In this Appendix, we prove that the function  $U(z) \equiv \sum_{n=0}^{\infty} u_n z^n$ , which is analytic in the region  $|z| < 0.702$  and there satisfies the differential equation (44), has a unique and simple root  $\alpha$  of smallest modulus. We also calculate the first ten digits of  $\alpha \approx 0.4889986317$ . The calculation is performed without first solving the differential equation to obtain the function  $U(z)$ .

We start in Lemma A.1 by providing a recurrence for  $u_n$ , which is then used to find an upper bound of  $|u_n|$  in Lemma A.3. Next, we consider the set  $\mathcal{B} \equiv \{z \in \mathbb{C} : |z| \leq \frac{1}{2}\}$  in the complex plane and decompose  $U(z)$  into a sum  $U(z) = U_1(z) + U_2(z)$ , where  $U_1(z) = \sum_{n=0}^{100} u_n z^n$  is a polynomial and  $U_2(z) = \sum_{n=101}^{\infty} u_n z^n$ . The bound for  $|u_n|$  in Lemma A.3 yields a bound

for  $|U_1(z)|$  (Lemma A.4), which in turn implies that  $|U_1(z)| > |U_2(z)|$  if  $z \in \partial\mathcal{B}$ . Hence, by Rouché’s theorem we have that inside  $\mathcal{B}$ , the function  $U(z)$  has the same number of roots—considered with their multiplicity—as the polynomial  $U_1(z)$ . Lemma A.5 shows that  $U_1(z)$  has a unique and simple root inside  $\mathcal{B}$ , and in Proposition A.6 we conclude the proof of our claim by finding an approximation of the unique and simple root  $\alpha$  of  $U(z)$  inside  $\mathcal{B}$ —which turns out to be very close to the root of  $U_1(z)$  inside  $\mathcal{B}$ .

In  $U(z) = \sum_{n=0}^{\infty} u_n z^n$ , we have  $u_n \equiv [z^n]U(z)$ . From equation (44), we derive a recurrence for  $u_n$ . Recall that  $e_n$  gives the mean number of root configurations in an ordered unlabeled history of size  $n \geq 1$ .

LEMMA A.1. *For  $n \geq 2$ , we have*

$$(46) \quad \begin{aligned} u_n &= \frac{1}{n(n-1)} \sum_{k=0}^{n-1} (3n-k-3)u_k - \frac{4}{n(n-1)} \sum_{k=0}^{n-1} (n-2k-1)e_{n-k}u_k \\ &+ \frac{4}{n(n-1)} \sum_{k=0}^{n-1} \left( \sum_{j=0}^{n-k-1} e_j \right) u_k, \end{aligned}$$

with  $u_0 = 1$  and  $u_1 = 0$ .

PROOF. First, notice that for  $n \geq 0$ , the coefficient of  $z^n$  in each term of equation (44) can be written as

$$\begin{aligned} [z^n]U''(z) &= (n+2)(n+1)u_{n+2}, \\ -[z^n] \left( g_1 + \frac{g'_2}{g_2} \right) U'(z) &= - \sum_{k=0}^n (n-k+1)(4e_{k+1} + 2)u_{n-k+1}, \\ [z^n]g_2g_0U(z) &= \sum_{k=0}^n \left[ (k+1) + 4 \sum_{j=0}^k e_{j+1} + 4 \sum_{j=0}^{k+2} e_j e_{k-j+2} \right] u_{n-k}, \end{aligned}$$

where for convenience we set  $e_0 = 0$ .

Making a substitution to the index of summation, we have

$$-4 \sum_{k=0}^n (n-k+1)e_{k+1}u_{n-k+1} = -4 \sum_{k=0}^{n+1} ke_{n-k+2}u_k.$$

Hence, the sum for  $-[z^n](g_1 + g'_2/g_2)U'(z)$  can be simplified as

$$-[z^n] \left( g_1 + \frac{g'_2}{g_2} \right) U'(z) = -4 \sum_{k=0}^{n+1} ke_{n-k+2}u_k - 2 \sum_{k=0}^n (n-k+1)u_{n-k+1}.$$

The second sum in this equation together with the first sum  $\sum_{k=0}^n (k+1)u_{n-k}$  of  $[z^n]g_2g_0U(z)$  give

$$-2 \sum_{k=0}^n (n-k+1)u_{n-k+1} + \sum_{k=0}^n (k+1)u_{n-k} = \sum_{k=0}^{n+1} (n-3k+1)u_k.$$

Furthermore, by setting  $n = k + 2$  in equation (36), the inner sums of  $[z^n]g_2g_0U(z)$  can be rewritten as

$$4 \sum_{j=0}^k e_{j+1} + 4 \sum_{j=0}^{k+1} e_j e_{k-j+2} = 4(k+1)e_{k+2} - 4(k+1) - 4 \sum_{j=1}^{k+1} e_j.$$

Hence, the coefficient of  $z^n$  in equation (44) becomes

$$(n + 2)(n + 1)u_{n+2} - 4 \sum_{k=0}^{n+1} ke_{n-k+2}u_k + \sum_{k=0}^{n+1} (n - 3k + 1)u_k + \sum_{k=0}^n \left[ 4(k + 1)e_{k+2} - 4(k + 1) - 4 \sum_{j=1}^{k+1} e_j \right] u_{n-k}.$$

In this expression, we make two substitutions:

$$\begin{aligned} \sum_{k=0}^n 4(k + 1)e_{k+2}u_{n-k} &= \sum_{k=0}^{n+1} 4(n - k + 1)e_{n-k+2}u_k, \\ \sum_{k=0}^{n+1} (n - 3k + 1)u_k - 4 \sum_{k=0}^n (k + 1)u_{n-k} &= \sum_{k=0}^{n+1} (n - 3k + 1)u_k - 4 \sum_{k=0}^n (n - k + 1)u_k \\ &= \sum_{k=0}^{n+1} (-3n + k - 3)u_k, \end{aligned}$$

obtaining

$$(n + 2)(n + 1)u_{n+2} - 4 \sum_{k=0}^{n+1} ke_{n-k+2}u_k + \sum_{k=0}^{n+1} 4(n - k + 1)e_{n-k+2}u_k + \sum_{k=0}^{n+1} (-3n + k - 3)u_k + \sum_{k=0}^n \left( -4 \sum_{j=1}^{k+1} e_j \right) u_{n-k},$$

and thus

$$(n + 2)(n + 1)u_{n+2} + \sum_{k=0}^{n+1} 4(n - 2k + 1)e_{n-k+2}u_k + \sum_{k=0}^{n+1} (-3n + k - 3)u_k + \sum_{k=0}^n \left( -4 \sum_{j=1}^{k+1} e_j \right) u_{n-k}.$$

Finally, because  $e_0 = 0$ , in this expression we can substitute

$$\begin{aligned} \sum_{k=0}^n \left( -4 \sum_{j=1}^{k+1} e_j \right) u_{n-k} &= \sum_{k=0}^n \left( -4 \sum_{j=0}^{k+1} e_j \right) u_{n-k} = \sum_{k=0}^n \left( -4 \sum_{j=0}^{n-k+1} e_j \right) u_k \\ &= \sum_{k=0}^{n+1} \left( -4 \sum_{j=0}^{n-k+1} e_j \right) u_k, \end{aligned}$$

obtaining for  $n \geq 0$ ,

$$(n + 2)(n + 1)u_{n+2} + \sum_{k=0}^{n+1} 4(n - 2k + 1)e_{n-k+2}u_k - \sum_{k=0}^{n+1} (3n - k + 3)u_k - 4 \sum_{k=0}^{n+1} \left( \sum_{j=0}^{n-k+1} e_j \right) u_k = 0,$$

which rescaled is recurrence (46). The initial conditions  $u_0 = 1$  and  $u_1 = 0$ , follow from the fact that  $U(0) = 1$  and  $U'(0) = 0$  as  $U(z) = \exp[\int_0^z S(x)/(-x) dx]$ .  $\square$

In Lemma A.3, we use the recurrence to find an upper bound for  $|u_n|$ . First, we need an upper bound for  $e_n$ .

LEMMA A.2. For  $n \geq 0$ , we have  $e_n \leq (\frac{9}{10})(\frac{3}{2})^n$ .

PROOF. Using the recurrence (36), with the help of computing software we have shown that the inequality holds for  $0 \leq n \leq 41$ . We proceed by induction. Suppose the inequality holds for all  $k < n$  with  $n > 41$ . By equation (36),

$$\begin{aligned} e_n &\leq 1 + \frac{81}{100(n-1)} \sum_{j=1}^{n-1} \left(\frac{3}{2}\right)^j + \frac{9}{5(n-1)} \sum_{j=1}^{n-1} \left(\frac{3}{2}\right)^j \\ &= 1 + \frac{81}{100} \left(\frac{3}{2}\right)^n + \frac{18}{5(n-1)} \left(\frac{3}{2}\right)^n - \frac{27}{5(n-1)} \\ &= \frac{9}{10} \left(\frac{3}{2}\right)^n - \frac{9}{10} \left(\frac{1}{10} - \frac{4}{n-1}\right) \left(\frac{3}{2}\right)^n - \frac{27}{5(n-1)} + 1. \end{aligned}$$

In the last step, we can see that a positive number is subtracted from  $\frac{9}{10}(\frac{3}{2})^n$  for  $n > 41$ , as

$$\frac{9}{10} \left(\frac{1}{10} - \frac{4}{n-1}\right) \left(\frac{3}{2}\right)^n + \frac{27}{5(n-1)} - 1 > \frac{9}{10} \frac{1}{400} \left(\frac{3}{2}\right)^{42} - 1 > 0.$$

Thus, the claim is proved.  $\square$

LEMMA A.3. For  $n \geq 0$ , we have  $|u_n| \leq (\frac{9}{5})^n$ .

PROOF. Using recurrence (46), computing software verifies the inequality for  $0 \leq n \leq 25$ . We proceed by induction. Suppose that the inequality holds for all  $k < n$  with  $n > 25$ . For simplicity of computation, instead of the bound in Lemma A.2, we use the more conservative  $(\frac{3}{2})^n$  as a bound for  $e_n$ . With equation (46), we get

$$\begin{aligned} |u_n| &\leq \frac{3}{n} \sum_{k=0}^{n-1} \left(\frac{9}{5}\right)^k + \frac{4}{n} \sum_{k=0}^{n-1} \left(\frac{3}{2}\right)^{n-k} \left(\frac{9}{5}\right)^k + \frac{4}{n(n-1)} \sum_{k=0}^{n-1} \left(\sum_{j=0}^{n-k-1} \left(\frac{3}{2}\right)^j\right) \left(\frac{9}{5}\right)^k \\ &= \frac{15}{4n} \left(\frac{9}{5}\right)^n - \frac{15}{4n} + \frac{20}{n} \left(\frac{9}{5}\right)^n - \frac{20}{n} \left(\frac{3}{2}\right)^n + \frac{30}{n(n-1)} \left(\frac{9}{5}\right)^n - \frac{40}{n(n-1)} \left(\frac{3}{2}\right)^n \\ &\quad + \frac{10}{n(n-1)} \\ &= \frac{5(19n+5)}{4n(n-1)} \left(\frac{9}{5}\right)^n - \frac{20(n+1)}{n(n-1)} \left(\frac{3}{2}\right)^n - \frac{5(3n-11)}{4n(n-1)}. \end{aligned}$$

In the last step, we have  $|u_n| \leq (\frac{9}{5})^n$ , as for  $n > 25$ , the following two inequalities hold:

$$\begin{aligned} \frac{5(19n+5)}{4n(n-1)} &\leq 1, \\ -\frac{20(n+1)}{n(n-1)} \left(\frac{3}{2}\right)^n - \frac{5(3n-11)}{4n(n-1)} &\leq 0. \end{aligned}$$

Thus, the claim is proved.  $\square$

We now consider the set  $\mathcal{B} \equiv \{z \in \mathbb{C} : |z| \leq \frac{1}{2}\}$ , and the partition  $U(z) = \sum_{k=0}^{\infty} u_k z^k = U_1(z) + U_2(z)$ ,  $U_1(z) \equiv \sum_{k=0}^{100} u_k z^k$  and  $U_2(z) \equiv \sum_{k=101}^{\infty} u_k z^k$ . Using the bound for  $|u_n|$  from Lemma A.3, for each  $z \in \mathcal{B}$  we have

$$(47) \quad |U_2(z)| \leq \sum_{k=101}^{\infty} |u_k| |z|^k \leq \sum_{k=101}^{\infty} \left(\frac{9}{5}\right)^k \left(\frac{1}{2}\right)^k = 10 \left(\frac{9}{10}\right)^{101} \approx 0.0002390525900.$$

Next, we need a lower bound for  $|U_1(z)|$ .

LEMMA A.4. *We have  $\min_{z \in \partial \mathcal{B}} |U_1(z)| \geq \frac{3}{1000}$ .*

PROOF. We obtain the result by considering a function

$$G(t) \equiv \left[ \sum_{k=0}^{100} u_k \cos(kt) \left(\frac{1}{2}\right)^k \right]^2 + \left[ \sum_{k=0}^{100} u_k \sin(kt) \left(\frac{1}{2}\right)^k \right]^2.$$

$G(t)$  has period  $2\pi$ , with  $G(\pi - t) = G(\pi + t)$ , if  $t \in [0, \pi]$ . For  $|z| \in \partial \mathcal{B}$ , we can write  $z = \frac{1}{2}[\cos t + i \sin t]$  for  $t \in [0, 2\pi)$ , and thus

$$\begin{aligned} |U_1(z)| &= \left| \sum_{k=0}^{100} u_k \left[ \left(\frac{1}{2}\right) [\cos t + i \sin t] \right]^k \right| \\ &= \left| \sum_{k=0}^{100} u_k \cos(kt) \left(\frac{1}{2}\right)^k + i \sum_{k=0}^{100} u_k \sin(kt) \left(\frac{1}{2}\right)^k \right| = \sqrt{G(t)}. \end{aligned}$$

By using the bound in Lemma A.3, we have the following inequality:

$$\begin{aligned} (48) \quad |G'(t)| &= \left| 2 \left[ \sum_{k=0}^{100} u_k \cos(kt) \left(\frac{1}{2}\right)^k \right] \left[ - \sum_{k=0}^{100} k u_k \sin(kt) \left(\frac{1}{2}\right)^k \right] \right. \\ &\quad \left. + 2 \left[ \sum_{k=0}^{100} u_k \sin(kt) \left(\frac{1}{2}\right)^k \right] \left[ \sum_{k=0}^{100} k u_k \cos(kt) \left(\frac{1}{2}\right)^k \right] \right| \\ &\leq 2 \left| \sum_{k=0}^{100} u_k \cos(kt) \left(\frac{1}{2}\right)^k \right| \left| \sum_{k=0}^{100} k u_k \sin(kt) \left(\frac{1}{2}\right)^k \right| \\ &\quad + 2 \left| \sum_{k=0}^{100} u_k \sin(kt) \left(\frac{1}{2}\right)^k \right| \left| \sum_{k=0}^{100} k u_k \cos(kt) \left(\frac{1}{2}\right)^k \right| \\ &\leq 2 \left[ \sum_{k=0}^{100} |u_k| |\cos(kt)| \left(\frac{1}{2}\right)^k \right] \left[ \sum_{k=0}^{100} k |u_k| |\sin(kt)| \left(\frac{1}{2}\right)^k \right] \\ &\quad + 2 \left[ \sum_{k=0}^{100} |u_k| |\sin(kt)| \left(\frac{1}{2}\right)^k \right] \left[ \sum_{k=0}^{100} k |u_k| |\cos(kt)| \left(\frac{1}{2}\right)^k \right] \\ &\leq 4 \left[ \sum_{k=0}^{100} \left(\frac{9}{10}\right)^k \right] \left[ \sum_{k=0}^{100} k \left(\frac{9}{10}\right)^k \right] \approx 3598.862135. \end{aligned}$$

We set  $\mathcal{I} = \left\{ \frac{k\pi}{1,000,000} : k \in \mathbb{Z}, 0 \leq k \leq 1,000,000 \right\}$ . A numerical calculation shows that

$$(49) \quad \min_{t \in \mathcal{I}} G(t) = G(0) \approx 0.01949528529.$$

With these preparations complete, we prove our claim by showing that

$$(50) \quad \min_{t \in [0, \pi]} G(t) \geq \frac{9}{1,000,000}.$$

We prove equation (50) by contradiction. Suppose there exists  $t_0 \in [0, \pi]$  such that  $G(t_0) < \frac{9}{1,000,000}$ . Then we can find  $t_1 \in \mathcal{I}$  such that

$$(51) \quad |t_1 - t_0| \leq \frac{\pi}{2,000,000}.$$

By the mean value theorem, we can find  $c \in (t_0, t_1)$  such that  $G(t_1) - G(t_0) = G'(c)(t_1 - t_0)$ . From equations (48) and (51),

$$(52) \quad \frac{1800\pi}{1,000,000} \geq |G'(c)(t_1 - t_0)| = |G(t_1) - G(t_0)| \geq G(t_1) - G(t_0).$$

However, because  $t_1 \in \mathcal{I}$ , by equation (49), we have

$$G(t_1) - G(t_0) \geq G(0) - G(t_0) \geq \frac{1}{100} - \frac{9}{1,000,000} = \frac{9991}{1,000,000}.$$

This result contradicts the upper bound in equation (52). Thus, equation (50) holds and the claim has been proven.  $\square$

Next, we study the root of  $U_1(z)$  inside  $\mathcal{B}$ .

LEMMA A.5. *The polynomial  $U_1(z)$  has a unique (simple) root  $\beta$  inside  $\mathcal{B}$ , with  $\beta \approx 0.4889986317$ .*

PROOF. First, by the intermediate value theorem, there exists a real root  $\beta$  with  $0 < \beta < \frac{1}{2}$ , as we can numerically compute  $U_1(0)U_1(\frac{1}{2}) < 0$  for the polynomial  $U_1(z)$ . Thus, we must prove

$$\begin{aligned} \frac{U_1(z)}{z - \beta} &= \frac{U_1(z) - U_1(\beta)}{z - \beta} = \sum_{k=0}^{100} u_k \frac{z^k - \beta^k}{z - \beta} \\ &= \sum_{k=0}^{100} u_k \sum_{\ell=0}^{k-1} \beta^{k-1-\ell} z^\ell = \sum_{\ell=0}^{99} \left( \sum_{k=\ell+1}^{100} u_k \beta^{k-1-\ell} \right) z^\ell \end{aligned}$$

satisfies  $|U_1(z)/(z - \beta)| > 0$  in  $\mathcal{B}$ .

To do so, we first use the bisection method for root-finding to numerically approximate  $\beta$  by

$$\tilde{\beta} = \frac{1,101,127,027,820,569}{2,251,799,813,685,248} \approx 0.4889986317,$$

with the approximation error

$$(53) \quad |\beta - \tilde{\beta}| \leq \frac{1}{2^{50}}.$$

Then we define the polynomial

$$Q(z) \equiv \sum_{\ell=0}^{99} a_\ell z^\ell \quad \text{with } a_\ell \equiv \sum_{k=\ell+1}^{100} u_k \tilde{\beta}^{k-1-\ell},$$

through which we can write

$$\frac{U_1(z)}{z - \beta} = Q(z) + (\beta - \tilde{\beta})R(z),$$

$$R(z) \equiv \sum_{\ell=0}^{99} \left( \sum_{k=\ell+1}^{100} u_k \frac{\beta^{k-1-\ell} - \tilde{\beta}^{k-1-\ell}}{\beta - \tilde{\beta}} \right) z^\ell = \sum_{\ell=0}^{99} \left( \sum_{k=\ell+2}^{100} u_k \sum_{j=0}^{k-2-\ell} \beta^j \tilde{\beta}^{k-2-\ell-j} \right) z^\ell.$$

Note that on  $\mathcal{B}$ ,

$$\begin{aligned} |R(z)| &\leq \sum_{\ell=0}^{99} \sum_{k=\ell+2}^{100} \sum_{j=0}^{k-2-\ell} |u_k| |\beta|^j |\tilde{\beta}|^{k-2-\ell-j} |z|^\ell \\ (54) \quad &\leq \sum_{\ell=0}^{99} \sum_{k=\ell+2}^{100} \sum_{j=0}^{k-2-\ell} \left(\frac{9}{5}\right)^k \left(\frac{1}{2}\right)^{k-2} \approx 3234.224489, \end{aligned}$$

where we used the bound for  $|u_n|$  from Lemma A.3 and the fact that  $\beta, \tilde{\beta}, |z| \leq \frac{1}{2}$ .

Next, let us consider the function

$$S(r, \theta) \equiv \sum_{\ell=0}^{99} a_\ell r^\ell \cos(\ell\theta)$$

defined over the rectangle  $(r, \theta) \in [0, \frac{1}{2}] \times [0, \pi]$ , where  $S(r, \theta) = \Re(Q(z))$  if  $z = r[\cos(\pm\theta) + i \sin(\pm\theta)] \in \mathcal{B}$ . We need the following bound for the gradient of  $S$ :

$$\begin{aligned} |\nabla S| &= \left| \left( \sum_{\ell=0}^{99} \ell a_\ell r^{\ell-1} \cos(\ell\theta), \sum_{\ell=0}^{99} -\ell a_\ell r^\ell \sin(\ell\theta) \right) \right| \\ &= \left| \sum_{\ell=0}^{99} (\ell a_\ell r^{\ell-1} \cos(\ell\theta), -\ell a_\ell r^\ell \sin(\ell\theta)) \right| \\ (55) \quad &= \left| \sum_{\ell=0}^{99} \ell a_\ell r^{\ell-1} (\cos(\ell\theta), -r \sin(\ell\theta)) \right| \leq \sum_{\ell=0}^{99} \ell |a_\ell| |r|^{\ell-1} |(\cos(\ell\theta), -r \sin(\ell\theta))| \\ &\leq \sum_{\ell=0}^{99} \ell |a_\ell| |r|^{\ell-1} \leq \sum_{\ell=0}^{99} \ell |a_\ell| \left(\frac{1}{2}\right)^{\ell-1} \approx 89.628949. \end{aligned}$$

Here, we have made use of  $|r| < \frac{1}{2}$  and for  $|r| < 1$ ,  $\sqrt{\cos^2 x + r^2 \sin^2 x} \leq \sqrt{\cos^2 x + \sin^2 x} = 1$ .

A numerical calculation shows that over the grid  $\mathcal{I} \equiv \{(\frac{k}{2000}, \frac{j\pi}{1000}) : (k, j) \in \mathbb{Z}^2, 0 \leq k, j \leq 1000\}$ , we have

$$(56) \quad \min_{(r, \theta) \in \mathcal{I}} |S(r, \theta)| = \left| S\left(\frac{1}{2}, \frac{502\pi}{1000}\right) \right| \approx 0.9518894218.$$

We now show—with a similar method to that used to prove Lemma A.4—that

$$(57) \quad \min_{(r, \theta) \in [0, \frac{1}{2}] \times [0, \pi]} |S(r, \theta)| \geq \frac{3235}{2^{50}}.$$

Suppose for contradiction that there exists  $z_0 = (r_0, \theta_0) \in [0, \frac{1}{2}] \times [0, \pi]$  such that  $|S(r_0, \theta_0)| < 3235/2^{50}$ . Then let us take  $z_1 = (r_1, \theta_1) \in \mathcal{I}$  such that

$$(58) \quad |z_1 - z_0| < \sqrt{\frac{1}{16} + \frac{\pi^2}{4}} \left(\frac{1}{1000}\right) \leq \frac{1}{500}.$$



By the mean value theorem, there exists a point  $(r, \theta)$  on the line segment from  $(r_0, \theta_0)$  to  $(r_1, \theta_1)$  such that

$$\nabla S(r, \theta) \cdot (z_1 - z_0) = S(r_1, \theta_1) - S(r_0, \theta_0),$$

where  $\cdot$  is the inner product of  $\mathbb{R}^2$ . By using the Cauchy–Schwarz inequality together with (55), (56) and (58), the assumption  $|S(r_0, \theta_0)| < 3235/2^{50}$  would thus give

$$\begin{aligned} \frac{90}{500} &\geq |\nabla S(r, \theta)| |z_1 - z_0| \geq |\nabla S(r, \theta) \cdot (z_1 - z_0)| = |S(r_1, \theta_1) - S(r_0, \theta_0)| \\ &\geq |S(r_1, \theta_1)| - |S(r_0, \theta_0)| \geq \frac{9}{10} - \frac{3235}{2^{50}} > 0.89, \end{aligned}$$

which is a contradiction. Hence, equation (57) holds.

Finally, because for  $z \in \mathcal{B}$  we have

$$|Q(z)| \geq |\Re(Q(z))| \geq \min_{(r, \theta) \in [0, \frac{1}{2}] \times [0, \pi]} |S(r, \theta)|,$$

by using equations (53), (54) and (57), it follows that in  $\mathcal{B}$ ,

$$\begin{aligned} \left| \frac{U_1(z)}{z - \beta} \right| &= |Q(z) + (\beta - \tilde{\beta})R(z)| \geq ||Q(z)| - |(\tilde{\beta} - \beta)R(z)|| \geq \frac{3235}{2^{50}} - |(\tilde{\beta} - \beta)||R(z)| \\ &\geq \frac{3235}{2^{50}} - \frac{|R(z)|}{2^{50}} > \frac{3235}{2^{50}} - \frac{3234.224489 \dots}{2^{50}} > 0. \end{aligned}$$

This concludes the proof.  $\square$

Combining Lemmas A.4 and A.5 with the inequality in equation (47), we obtain the following proposition.

**PROPOSITION A.6.** *The function  $U(z)$  has a unique (simple) root  $\alpha$  inside  $\mathcal{B}$ , where  $\alpha \approx 0.4889986317$ .*

**PROOF.** For the decomposition  $U(z) = U_1(z) + U_2(z)$ , equation (47) together with Lemma A.4 gives for  $z \in \partial\mathcal{B}$ ,

$$|U_1(z)| \geq \frac{3}{1000} > 0.00025 > |U_2(z)|.$$

Hence, from Rouché’s theorem, inside  $\mathcal{B}$  the function  $U(z)$  has the same number of roots (considered with multiplicity) as polynomial  $U_1(z)$ . From Lemma A.5, we know that  $U_1(z)$  has one (simple) root inside  $\mathcal{B}$ .

The only remaining step is the numerical computation of  $\alpha$ , whose first ten digits turn out to coincide with the constant  $\beta$  found in Lemma A.5 as the root of  $U_1(z)$  inside  $\mathcal{B}$ . We again decompose  $U(z)$ :

$$U(z) = \sum_{k=0}^{\infty} u_k z^k = \sum_{k=0}^{500} u_k z^k + \sum_{k=501}^{\infty} u_k z^k = \tilde{U}_1(z) + \tilde{U}_2(z).$$

Note that from our bound for  $|u_k|$  (Lemma A.3), for each  $z \in \mathcal{B}$  we have

$$(59) \quad |\tilde{U}_2(z)| \leq \sum_{k=501}^{\infty} |u_k| |z|^k \leq \sum_{k=501}^{\infty} \left(\frac{9}{5}\right)^k \left(\frac{1}{2}\right)^k = 10 \left(\frac{9}{10}\right)^{501} \leq 10^{-21}.$$

Let us now consider

$$\alpha' = \frac{550,563,513,910,285}{1,125,899,906,842,624} \approx 0.48899863172938484723,$$

$$\alpha'' = \frac{1,101,127,027,820,571}{2,251,799,813,685,248} \approx 0.48899863172938529132.$$

These values were chosen using the bisection method such that

$$\tilde{U}_1(\alpha') = 2.708185805 \dots \cdot 10^{-16} \quad \text{and} \quad \tilde{U}_1(\alpha'') = -4.953373282 \dots \cdot 10^{-15}.$$

From the bound of  $|\tilde{U}_2(z)|$  in equation (59), it is clear that  $U(\alpha') > 0$  and  $U(\alpha'') < 0$ . Let  $\alpha$  be the unique root of  $U(z)$  in  $\mathcal{B}$ , which by the intermediate value theorem must be a real root in  $(\alpha', \alpha'')$ , and let  $\epsilon \equiv \alpha - \alpha' \leq 10^{-14}$ . Note that

$$\frac{1}{\alpha'} - \frac{1}{\alpha} = \frac{\epsilon}{\alpha'(\alpha' + \epsilon)} \leq \frac{\epsilon}{(\alpha')^2} \leq 5 \cdot 10^{-14}.$$

Thus, we can use

$$\alpha' = 0.48899863172938484723,$$

$$(\alpha')^{-1} = 2.0449954971518340953$$

to approximate  $\alpha$  and  $\alpha^{-1}$ , respectively.  $\square$

**Acknowledgments.** This work developed from discussions at the Banff International Research Station.

**Funding.** Support was provided by a Rita Levi-Montalcini grant from the Ministero dell'Istruzione, dell'Università e della Ricerca (FD), grants MOST-104-2923-M-009-006-MY3 and MOST-107-2115-M-009-010-MY2 (MF, ARP) and National Institutes of Health grant R01 GM131404 (NAR).

## REFERENCES

- [1] ALDOUS, D. (1996). Probability distributions on cladograms. In *Random Discrete Structures (Minneapolis, MN, 1993)*. IMA Vol. Math. Appl. **76** 1–18. Springer, New York. MR1395604 [https://doi.org/10.1007/978-1-4612-0719-1\\_1](https://doi.org/10.1007/978-1-4612-0719-1_1)
- [2] ALIMPIEV, E. and ROSENBERG, N. A. (2021). Enumeration of coalescent histories for caterpillar species trees and  $p$ -pseudocaterpillar gene trees. *Adv. in Appl. Math.* **131** Paper No. 102265. MR4305882 <https://doi.org/10.1016/j.aam.2021.102265>
- [3] BERGERON, F., FLAJOLET, P. and SALVY, B. (1992). Varieties of increasing trees. In *CAAP '92 (Rennes, 1992)*. Lecture Notes in Computer Science **581** 24–48. Springer, Berlin. MR1251994 [https://doi.org/10.1007/3-540-55251-0\\_2](https://doi.org/10.1007/3-540-55251-0_2)
- [4] BLUM, M. G. B., FRANÇOIS, O. and JANSON, S. (2006). The mean, variance and limiting distribution of two statistics sensitive to phylogenetic tree balance. *Ann. Appl. Probab.* **16** 2195–2214. MR2288718 <https://doi.org/10.1214/105051606000000547>
- [5] BODINI, O., COURTIÉL, J., DOVGAL, S. and HWANG, H.-K. (2018). Asymptotic distribution of parameters in random maps. In *29th International Conference on Probabilistic, Combinatorial and Asymptotic Methods for the Analysis of Algorithms. LIPIcs. Leibniz Int. Proc. Inform.* **110** Art. No. 13, 12. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern. MR3826132
- [6] BROWN, J. K. M. (1994). Probabilities of evolutionary trees. *Syst. Biol.* **43** 78–91.
- [7] CHANG, H. and FUCHS, M. (2010). Limit theorems for patterns in phylogenetic trees. *J. Math. Biol.* **60** 481–512. MR2587587 <https://doi.org/10.1007/s00285-009-0275-6>
- [8] DEGNAN, J. H. (2005). *Gene Tree Distributions Under the Coalescent Process*. ProQuest LLC, Ann Arbor, MI. Thesis (Ph.D.)—The University of New Mexico. MR2707346
- [9] DEGNAN, J. H., ROSENBERG, N. A. and STADLER, T. (2012). A characterization of the set of species trees that produce anomalous ranked gene trees. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **9** 1558–1568.

- [10] DEGNAN, J. H., ROSENBERG, N. A. and STADLER, T. (2012). The probability distribution of ranked gene trees on a species tree. *Math. Biosci.* **235** 45–55. MR2901026 <https://doi.org/10.1016/j.mbs.2011.10.006>
- [11] DISANTO, F. and MUNARINI, E. (2019). Local height in weighted Dyck models of random walks and the variability of the number of coalescent histories for caterpillar-shaped gene trees and species trees. *SN Appl. Sci.* **1** 578.
- [12] DISANTO, F. and ROSENBERG, N. A. (2015). Coalescent histories for lodgepole species trees. *J. Comput. Biol.* **22** 918–929. MR3407775 <https://doi.org/10.1089/cmb.2015.0015>
- [13] DISANTO, F. and ROSENBERG, N. A. (2016). Asymptotic properties of the number of matching coalescent histories for caterpillar-like families of species trees. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **13** 913–925.
- [14] DISANTO, F. and ROSENBERG, N. A. (2017). Enumeration of ancestral configurations for matching gene trees and species trees. *J. Comput. Biol.* **24** 831–850. MR3705080 <https://doi.org/10.1089/cmb.2016.0159>
- [15] DISANTO, F. and ROSENBERG, N. A. (2019). On the number of non-equivalent ancestral configurations for matching gene trees and species trees. *Bull. Math. Biol.* **81** 384–407. MR3902906 <https://doi.org/10.1007/s11538-017-0342-x>
- [16] DISANTO, F. and ROSENBERG, N. A. (2019). Enumeration of compact coalescent histories for matching gene trees and species trees. *J. Math. Biol.* **78** 155–188. MR3932643 <https://doi.org/10.1007/s00285-018-1271-5>
- [17] DISANTO, F., SCHLIZIO, A. and WIEHE, T. (2013). Yule-generated trees constrained by node imbalance. *Math. Biosci.* **246** 139–147. MR3122195 <https://doi.org/10.1016/j.mbs.2013.08.008>
- [18] DISANTO, F. and WIEHE, T. (2013). Exact enumeration of cherries and pitchforks in ranked trees under the coalescent model. *Math. Biosci.* **242** 195–200. MR3068684 <https://doi.org/10.1016/j.mbs.2013.01.010>
- [19] FELSENSTEIN, J. (1978). The number of evolutionary trees. *Syst. Zool.* **27** 27–33.
- [20] FILL, J. A. (1996). On the distribution of binary search trees under the random permutation model. *Random Structures Algorithms* **8** 1–25. MR1368848 [https://doi.org/10.1002/\(SICI\)1098-2418\(199601\)8:1<1::AID-RSA1>3.0.CO;2-1](https://doi.org/10.1002/(SICI)1098-2418(199601)8:1<1::AID-RSA1>3.0.CO;2-1)
- [21] FILL, J. A. and KAPUR, N. (2004). Limiting distributions for additive functionals on Catalan trees. *Theoret. Comput. Sci.* **326** 69–102. MR2094243 <https://doi.org/10.1016/j.tcs.2004.05.010>
- [22] FLAJOLET, P., GOURDON, X. and MARTÍNEZ, C. (1997). Patterns in random binary search trees. *Random Structures Algorithms* **11** 223–244. MR1609509 [https://doi.org/10.1002/\(SICI\)1098-2418\(199710\)11:3<223::AID-RSA2>3.0.CO;2-2](https://doi.org/10.1002/(SICI)1098-2418(199710)11:3<223::AID-RSA2>3.0.CO;2-2)
- [23] FLAJOLET, P. and SEDGEWICK, R. (2009). *Analytic Combinatorics*. Cambridge Univ. Press, Cambridge. MR2483235 <https://doi.org/10.1017/CBO9780511801655>
- [24] FUCHS, M., HOLMGREN, C., MITSCHKE, D. and NEININGER, R. (2021). A note on the independence number, domination number and related parameters of random binary search trees and random recursive trees. *Discrete Appl. Math.* **292** 64–71. MR4203518 <https://doi.org/10.1016/j.dam.2020.12.013>
- [25] HARDING, E. F. (1971). The probabilities of rooted tree-shapes generated by random bifurcation. *Adv. in Appl. Probab.* **3** 44–77. MR0282451 <https://doi.org/10.2307/1426329>
- [26] HIMWICH, Z. M. and ROSENBERG, N. A. (2020). Roadblocked monotonic paths and the enumeration of coalescent histories for non-matching caterpillar gene trees and species trees. *Adv. in Appl. Math.* **113** 101939, 33. MR4025683 <https://doi.org/10.1016/j.aam.2019.101939>
- [27] HOLMGREN, C. and JANSON, S. (2015). Limit laws for functions of fringe trees for binary search trees and random recursive trees. *Electron. J. Probab.* **20** no. 4, 51. MR3311217 <https://doi.org/10.1214/EJP.v20-3627>
- [28] LAMBERT, A. and STADLER, T. (2013). Birth-death models and coalescent point processes: The shape and probability of reconstructed phylogenies. *Theor. Popul. Biol.* **90** 113–128.
- [29] MADDISON, W. P. (1997). Gene trees in species trees. *Syst. Biol.* **46** 523–536.
- [30] MADDISON, W. P. and KNOWLES, L. L. (2006). Inferring phylogeny despite incomplete lineage sorting. *Syst. Biol.* **55** 21–30. <https://doi.org/10.1080/10635150500354928>
- [31] MCKENZIE, A. and STEEL, M. (2000). Distributions of cherries for two models of trees. *Math. Biosci.* **164** 81–92. MR1747204 [https://doi.org/10.1016/S0025-5564\(99\)00060-7](https://doi.org/10.1016/S0025-5564(99)00060-7)
- [32] ROSENBERG, N. A. (2006). The mean and variance of the numbers of  $r$ -pronged nodes and  $r$ -caterpillars in Yule-generated genealogical trees. *Ann. Comb.* **10** 129–146. MR2233885 <https://doi.org/10.1007/s00026-006-0278-6>
- [33] ROSENBERG, N. A. (2007). Counting coalescent histories. *J. Comput. Biol.* **14** 360–377. MR2312126 <https://doi.org/10.1089/cmb.2006.0109>
- [34] ROSENBERG, N. A. (2013). Coalescent histories for caterpillar-like families. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **10** 1253–1262.

- [35] ROSENBERG, N. A. (2019). Enumeration of lonely pairs of gene trees and species trees by means of antipodal cherries. *Adv. in Appl. Math.* **102** 1–17. MR3866986 <https://doi.org/10.1016/j.aam.2018.09.001>
- [36] ROSENBERG, N. A. and TAO, R. (2008). Discordance of species trees with their most likely gene trees: The case of five taxa. *Syst. Biol.* **57** 131–140.
- [37] STADLER, T. and DEGNAN, J. H. (2012). A polynomial time algorithm for calculating the probability of a ranked gene tree given a species tree. *Algorithms Mol. Biol.* **7** 7.
- [38] STANLEY, R. P. (1999). *Enumerative Combinatorics. Vol. 2. Cambridge Studies in Advanced Mathematics* **62**. Cambridge Univ. Press, Cambridge. With a foreword by Gian-Carlo Rota and appendix 1 by Sergey Fomin. MR1676282 <https://doi.org/10.1017/CBO9780511609589>
- [39] STEEL, M. (2016). *Phylogeny—Discrete and Random Processes in Evolution. CBMS-NSF Regional Conference Series in Applied Mathematics* **89**. SIAM, Philadelphia, PA. MR3601108 <https://doi.org/10.1137/1.9781611974485.ch1>
- [40] STEEL, M. and MCKENZIE, A. (2001). Properties of phylogenetic trees generated by Yule-type speciation models. *Math. Biosci.* **170** 91–112. MR1818785 [https://doi.org/10.1016/S0025-5564\(00\)00061-4](https://doi.org/10.1016/S0025-5564(00)00061-4)
- [41] THAN, C. and NAKHLEH, L. (2009). Species tree inference by minimizing deep coalescences. *PLoS Comput. Biol.* **5** e1000501. MR2559310 <https://doi.org/10.1371/journal.pcbi.1000501>
- [42] THAN, C., RUTHS, D., INNAN, H. and NAKHLEH, L. (2007). Confounding factors in HGT detection: Statistical error, coalescent effects, and multiple solutions. *J. Comput. Biol.* **14** 517–535. MR2318615 <https://doi.org/10.1089/cmb.2007.A010>
- [43] THAN, C. V. and ROSENBERG, N. A. (2013). Mathematical properties of the deep coalescence cost. *IEEE/ACM Trans. Comput. Biol. Bioinform.* **10** 61–72.
- [44] THAN, C. V. and ROSENBERG, N. A. (2014). Mean deep coalescence cost under exchangeable probability distributions. *Discrete Appl. Math.* **174** 11–26. MR3215453 <https://doi.org/10.1016/j.dam.2014.02.010>
- [45] WAGNER, S. (2015). Central limit theorems for additive tree parameters with small toll functions. *Combin. Probab. Comput.* **24** 329–353. MR3318048 <https://doi.org/10.1017/S0963548314000443>
- [46] WU, Y. (2012). Coalescent-based species tree inference from gene tree topologies under incomplete lineage sorting by maximum likelihood. *Evolution* **66** 763–775. <https://doi.org/10.1111/j.1558-5646.2011.01476.x>
- [47] WU, Y. (2016). An algorithm for computing the gene tree probability under the multispecies coalescent and its application in the inference of population tree. *Bioinformatics* **32** i225–i233. <https://doi.org/10.1093/bioinformatics/btw261>
- [48] YULE, G. U. (1925). A mathematical theory of evolution based on the conclusions of Dr. J. C. Willis, F. R. S. *Philos. Trans. R. Soc. Lond. B* **213** 21–87.