



# Clades, clans, and reciprocal monophyly under neutral evolutionary models

Sha Zhu, James H. Degnan, Mike Steel\*

Biomathematics Research Centre, University of Canterbury, Private Bag 4800, Christchurch 8140, New Zealand

## ARTICLE INFO

### Article history:

Received 16 December 2010

Available online 21 March 2011

### Keywords:

Yule tree  
Coalescent model  
Clade  
Clan

## ABSTRACT

The Yule model and the coalescent model are two neutral stochastic models for generating trees in phylogenetics and population genetics, respectively. Although these models are quite different, they lead to identical distributions concerning the probability that pre-specified groups of taxa form monophyletic groups (clades) in the tree. We extend earlier work to derive exact formulae for the probability of finding one or more groups of taxa as clades in a rooted tree, or as ‘clans’ in an unrooted tree. Our findings are relevant for calculating the statistical significance of observed monophyly and reciprocal monophyly in phylogenetics.

© 2011 Elsevier Inc. All rights reserved.

## 1. Introduction

When gene trees are estimated from multiple lineages taken from two or more populations, there is an increased chance that lineages within each population form monophyletic groups compared to sampling multiple lineages from a single population. In asking whether a particular group of lineages came from a taxonomically distinct population (Cummins et al., 2008; Rosenberg, 2007), this observation has led to the adoption of a null hypothesis that a set of lineages belongs to a single population or taxonomic group. Statistical tests for reciprocal monophyly between two sister taxa can then be developed to test against this null hypothesis (Hudson and Coyne, 2002; Rosenberg, 2003). Here, ‘reciprocal monophyly’ is the condition that, as one looks back in time, lineages coalesce within each of the two taxa, before any coalescence events take place between the two taxa.

Reciprocal monophyly is central to the genealogical species concept. According to this concept, two groups come from different species if they form distinct monophyletic groups (DeQuiroz, 2007; Hudson and Coyne, 2002). Gene trees from lineages sampled from one or more populations are typically estimated, and monophyly (or lack of monophyly) of these groups can be observed from the clades of the gene tree. Statistical tests for whether observed levels of monophyly provide sufficient evidence to conclude that a group is taxonomically distinct can be performed, given a probabilistic model for the clades on a tree (Rosenberg, 2007).

Two neutral models – involving different evolutionary scales – are useful in this context. The Yule (pure birth, or birth–death) model describes the speciation (and extinction) of lineages at the

species level as one moves forward in time, while Kingman’s coalescent process is a population genetic model which traces the ancestry of individual lineages back in time (and which thereby forms a tree). These are two quite different processes, and they lead to different branch lengths on trees; remarkably, however, they generate identical distributions of tree topologies (Aldous, 1995). Thus, while the coalescent process is a natural model for trees in single populations, the equivalence of the Yule and coalescent models for tree topologies means that results for the Yule model can be exploited in studying probabilities of clades for coalescent trees in single populations.

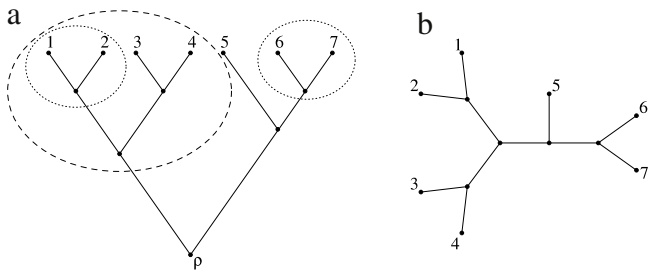
Although there has been an emphasis on testing for the taxonomic distinctiveness of one group of lineages, joint probabilities of clades could be used to examine whether the observed monophyly of several groups is statistically significant using a single test. Such an omnibus test of the null hypothesis that all groups come from one population might be more powerful than testing several groups one at a time.

In this note, we derive exact formulae for the joint probabilities of  $k$  clades for a random Yule/coalescent gene tree under the conditions that the  $k$  clades are mutually exclusive (they have no leaves of the gene tree in common), and are either exhaustive (all leaves of the gene tree occur in one of the  $k$  clades), or form only a subset of the leaves of the gene tree. These results generalize results from Rosenberg (2003), which provided an explicit formula for the probability that two mutually exclusive and exhaustive sets of leaves formed clades on a Yule/coalescent gene tree.

In addition, we extend the results to unrooted trees by giving the probabilities of ‘clans’ (sets of leaves that are all on one side of a split Wilkinson et al., 2007), as well as the joint probability of  $k > 1$  clans, on Yule/coalescent trees which have been unrooted. This extension is relevant when only unrooted trees can be estimated, which is particularly common in microbial evolution (Lapointe et al., 2010).

\* Corresponding author.

E-mail address: [mathmomike@gmail.com](mailto:mathmomike@gmail.com) (M. Steel).



**Fig. 1.** (a) This rooted tree has 13 ‘clades’, including the three sets circled ( $\{1, 2\}$ ,  $\{1, 2, 3, 4\}$ ,  $\{6, 7\}$ ). In this tree,  $\{1, 2\}$  and  $\{3, 4\}$  are sister clades, but  $\{1, 2\}$  and  $\{6, 7\}$  are not. (b) The unrooted tree  $T^{-\rho}$  obtained from the tree  $T$  in (a) by suppressing the root vertex  $\rho$ . This tree has  $\{3, 4, 5, 6, 7\}$  as a ‘clan’, even though this set is not a clade of  $T$ .

Our arguments throughout rely on just a few generic properties of neutral phylogenetic processes, such as the Yule and coalescent models. One of these properties is ‘exchangeability’ which, roughly speaking, requires that the probability of any rooted phylogenetic tree on a given leaf set depends just on the tree’s ‘shape’ and not on how its leaves are labelled. Of course, for trees constructed from biological data, the posterior probability of different phylogenetic trees will depend very much on how the leaves are labelled – with species that share high sequence similarity tending to be clustered together in the tree. Thus the ‘exchangeability’ property is one that should be viewed as appropriate for a ‘prior’ distribution on trees – that is, before one has considered the data. This is particularly relevant to a statistical test we describe in the final section, which relies solely on the exchangeability property.

**2. Clades**

Throughout this paper, we will let  $X_n$  (or, more briefly,  $X$ ) denote a set of taxa of size  $n$ . Given a rooted phylogenetic  $X$ -tree  $T_X$  (more briefly  $T$ ), with leaf set  $X = X_n$ , a *clade of  $T$*  is either an element of  $X$  or a subset of  $X$  that corresponds to the set of leaves that are descended from any internal vertex. For example, in Fig. 1(a), the sets  $\{3\}$ ,  $\{3, 4\}$  and  $\{1, 2, 3, 4\}$  are three clades. Throughout, we use  $A, B \subseteq X$ , etc., to denote clades, and  $\mathcal{E}, \mathcal{F}$ , etc., to denote events. Any two clades  $A$  and  $B$  of  $T$  satisfy the following *compatibility condition*:

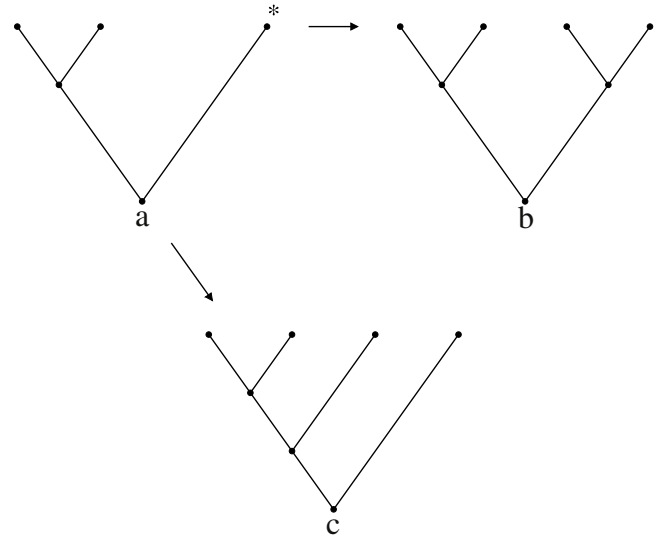
$$A \cap B \in \{A, B, \emptyset\}. \tag{1}$$

This is equivalent to requiring that  $A = B$ , one set is a strict subset of the other, or the two sets are disjoint.

We will let  $c(T)$  denote the set of clades of  $T$ , and say that a clade is *proper* if it is a strict subset of  $X$ . Notice that a rooted phylogenetic  $X$ -tree has at most  $2n - 1$  clades, and it has precisely this number if and only the tree is *binary*, that is, if each non-leaf vertex has two descendant vertices.

**3. The Yule–Harding–Kingman process**

Consider the probability distribution on binary phylogenetic  $X$ -trees described by a model that grows a tree by selecting a leaf uniformly at random and ‘splitting’ it into two new leaves, as illustrated in Fig. 2. Since we are ignoring branch lengths in this paper and concentrating just on tree topologies, the resulting probability distribution on rooted binary tree topologies is the same as that given by any (stationary or non-stationary) birth–death process on trees in which birth (speciation) and death (extinction) events apply exchangeably to all the species extant at any given moment (see Aldous, 1995 for further details). This is useful, since the rates of speciation and extinction throughout time



**Fig. 2.** From a rooted binary tree on three leaves (a), splitting the right leaf (\*) leads to a ‘balanced’ tree shape (b), while splitting either of the other two leaves produces an unbalanced tree (c). Thus the balanced tree shape has probability 1/3, and as there are three distinct ways to label the leaves, each of these rooted binary phylogenetic trees has probability 1/9 under the YHK process. For a phylogenetic tree of shape (c), the probability is 1/18.

may be both time dependent and variable according to the number of taxa present (Rabosky and Lovette, 2008).

The study of such pure-birth trees was initiated in Yule’s 1925 paper (Yule, 1925), and the probability distribution on tree topologies (without reference to branch lengths) was further studied by Harding (1971). Moreover, this probability distribution on trees is precisely the same as that given by a quite different process, namely Kingman’s coalescent process (Kingman, 1982) in population genetics, which starts at the leaves and successively combines pairs of elements, provided that, once again, we ignore branch lengths (Aldous, 1995).

To emphasize this equivalence between a model in macro-evolution (speciation and extinction) and micro-evolution (population genetics), we will refer to it as the *Yule–Harding–Kingman (YHK) process* for generating tree topologies.

We will also refer to a random binary phylogenetic  $X$ -tree produced by any of these stochastically equivalent processes as  $\mathcal{T}_X$  (or often just  $\mathcal{T}$  if  $X$  is clear), and so  $\mathbb{P}(\mathcal{T}_X = T)$  is the probability that  $T$  is the actual phylogenetic  $X$ -tree produced by the process. The process, viewed as a pure-birth model, is illustrated in Fig. 2.

In this paper, we exploit two important properties of the process that generates  $\mathcal{T}$ . First, we recall some notation that will be used throughout: for any phylogenetic  $X$ -tree and any non-empty subset  $Y$  of  $X$ , let  $T_{X|Y}$  be the phylogenetic tree induced by restricting the leaf set to  $Y$  (as in Semple and Steel (2003)). The two properties that the YHK process enjoys, and which we will exploit throughout this paper, are the following.

(EP) If  $T'$  is obtained from  $T$  by permuting its leaves, then

$$\mathbb{P}(\mathcal{T} = T') = \mathbb{P}(\mathcal{T} = T).$$

(GE) For any proper (and non-empty) subset  $A$  of  $X$ , and any rooted binary phylogenetic tree  $T$  with leaf set  $X - A$ ,

$$\mathbb{P}(\mathcal{T}_{X|(X-A)} = T \mid A \in c(\mathcal{T})) = \mathbb{P}(\mathcal{T}_{X-A} = T).$$

Property (EP) is the *exchangeability* property (Aldous, 1995), which requires that the probability of a particular phylogenetic tree depends just on its shape and not on how its leaves are labeled (it is called ‘label-invariance’ in Steel and Penny (1993)). Property (GE) is the *group elimination* property from Aldous (1995); it states

that, conditional on  $A$  forming a clade in the tree, the tree structure on the remaining taxa is also described by the YHK process. In turn, (GE) implies the following *sampling consistency* property (Aldous, 1995). For any rooted binary tree  $T$  with leaf set  $A \subseteq X$ , we have

$$(SC) \mathbb{P}(\mathcal{T}_{X|A} = T) = \mathbb{P}(\mathcal{T}_A = T).$$

To see that (GE) implies (SC), one sequentially deletes leaves that are not in  $A$ , noting that each leaf is, trivially, a clade in any tree.

**4. Clade probabilities under the YHK process**

The following result is stated and proved in Rosenberg (2006, Theorem 4.4) (it also appears as Proposition 2 of Blum and Francois, 2005, and in the appendix of Heard, 1992). A further proof of this result is also possible based on induction on  $n$  and using the well-known property of the YHK model that the number of leaves in one of the (randomly selected) maximal subtrees of  $\mathcal{T}_X$  is uniformly distributed between 1 and  $n - 1$ .

**Lemma 4.1.** *Let  $X_n(a)$  be the number of proper clades of size  $a$  in  $\mathcal{T}_X$ , where  $n = |X|$ . Then*

$$\mathbb{E}[X_n(a)] = \frac{2n}{a(a+1)}, \quad 1 \leq a \leq n-1. \quad \square$$

For a subset  $A$  of  $X$ , let  $p_n(A)$  be the probability that  $A$  is a proper clade of  $\mathcal{T}_X$ . From the exchangeability property (EP) it is clear that this probability depends only on  $a = |A|$  and  $n$ , and so we can write  $p_n(a)$  for this probability. From Rosenberg (2003), we have

**Lemma 4.2.**

$$p_n(a) = \begin{cases} \frac{2n}{a(a+1)} \binom{n}{a}^{-1}, & \text{if } 1 \leq a \leq n-1; \\ 0, & \text{otherwise.} \end{cases}$$

The proof of this result from Rosenberg (2003) relies on a combinatorial identity to sum a series. Here, we point out how Lemma 4.2 follows very directly from Lemma 4.1.

**Proof of Lemma 4.2.** For  $1 \leq a \leq n - 1$ , the exchangeability property (EP) implies that

$$p_n(A) = \sum_{k \geq 0} \mathbb{P}(\mathcal{T} \text{ has } k \text{ clades of size } a) \cdot \frac{k}{\binom{n}{a}} = \mathbb{E}[X_n(a)] \binom{n}{a}^{-1},$$

where  $X_n(a)$  is as defined in Lemma 4.1. This completes the proof.  $\square$

4.1. Pairs of clades

For a pair  $A, B$  of disjoint subsets of  $X$ , let  $\hat{p}_n(A, B)$  be the probability that  $A$  and  $B$  are *sister clades* of  $\mathcal{T}_X$  (i.e.,  $A, B$  and  $A \cup B$  are clades of  $\mathcal{T}_X$ ). By exchangeability (EP), this probability depends on  $a = |A|, b = |B|$  and  $n$  only, and so we will denote it  $\hat{p}_n(a, b)$ .

Consider first the special case where  $n = a + b$ ; that is,  $A$  and  $X - A$  are sister clades, which is equivalent to saying that  $A$  is a maximal proper clade. From Brown (1994, Eq. 6) (see also Rosenberg, 2003), the probability of this event is given as follows.

**Lemma 4.3.** *For  $1 \leq a \leq n$ , we have*

$$\hat{p}_n(a, n - a) = \frac{2}{n - 1} \binom{n}{a}^{-1}.$$

We generalize this slightly, as follows.

**Lemma 4.4.** *Let  $k = a + b \leq n$ . Then*

$$\hat{p}_n(a, b) = \frac{4a!b!(n - k)!}{(n - 1)!k(k^2 - 1)}.$$

**Proof.**

$$\hat{p}_n(A, B) = \mathbb{P}(A \in c(\mathcal{T}_{X|A \cup B}) \mid A \cup B \in c(\mathcal{T}_X)) \cdot \mathbb{P}(A \cup B \in c(\mathcal{T}_X)).$$

Applying Lemma 4.2 to the first term, and property (SC) and Lemma 4.3 to the second term, we have

$$\hat{p}_n(A, B) = \frac{2}{a + b - 1} \binom{a + b}{a}^{-1} \cdot \frac{2n}{(a + b)(a + b + 1)} \binom{n}{a + b}^{-1},$$

from which the result follows.  $\square$

Now, for any two arbitrary subsets  $A \subseteq X_n$  and  $B \subseteq X_n$ , let  $p_n(A, B)$  be the probability that a YHK tree  $\mathcal{T}$  on  $X_n$  has  $A$  and  $B$  as proper clades. As usual, let  $a = |A|$  and  $b = |B|$ .

**Theorem 4.5.**

$$p_n(A, B) = \begin{cases} p_n(a) & \text{if } A = B \text{ [Case 1];} \\ R_n(a, b), & \text{if } A \subsetneq B \text{ [Case 2];} \\ R_n(b, a), & \text{if } B \subsetneq A \text{ [Case 3];} \\ \hat{p}_n(a, n - a), & \text{if } A \cap B = \emptyset, A \cup B = X_n \text{ [Case 4];} \\ r_n(a, b), & \text{if } A \cap B = \emptyset, A \cup B \subsetneq X_n \text{ [Case 5];} \\ 0, & \text{otherwise [Case 6];} \end{cases}$$

where

$p_n(a)$ , and  $\hat{p}_n(a, n - a)$  are given by Lemmas 4.2 and 4.3,

$$R_n(a, b) := \frac{4n}{a(a + 1)(b + 1)} \binom{n}{b}^{-1} \binom{b}{a}^{-1},$$

$$r_n(a, b) := \frac{4a!b!(n - a - b)!}{(n - 1)!} G_n(a, b), \quad \text{and where}$$

$$G_n(a, b) := \frac{n}{ab(a + 1)(b + 1) - \frac{a(a + 1) + b(b + 1) + ab}{ab(a + 1)(b + 1)(a + b + 1)} + \frac{1}{(a + b)((a + b)^2 - 1)}}.$$

**Proof.** Cases 1 and 4 are given by Lemmas 4.2 and 4.3, respectively. For the second case ( $A \subsetneq B$ ), we have

$$p_n(A, B) = \mathbb{P}(A \in c(\mathcal{T}_X) \mid B \in c(\mathcal{T}_X)) \cdot \mathbb{P}(B \in c(\mathcal{T}_X)).$$

Since  $A \subsetneq B$ , we can apply property (SC) and Lemma 4.2 to deduce that the first term in this product is  $\frac{2b}{a(a+1)} \binom{b}{a}^{-1}$ , while the second term is  $\frac{2n}{b(b+1)} \binom{n}{b}^{-1}$ , from which the result follows. Case 3 follows by an analogous argument. For Case 5, consider the following two pairs of events:

- $\mathcal{E}_1 : A, B \in c(\mathcal{T}_X)$ ,
- $\mathcal{E}_2 : A \cup B, B \in c(\mathcal{T}_X)$ ,
- $\mathcal{F}_1 : A \in c(\mathcal{T}_{X|(X-B)})$ ,
- $\mathcal{F}_2 : B \in c(\mathcal{T}_X)$ .

We are interested in computing  $\mathbb{P}(\mathcal{E}_1)$ , since this is  $p_n(A, B)$ , and by the principle of inclusion and exclusion we have

$$\mathbb{P}(\mathcal{E}_1) = \mathbb{P}(\mathcal{E}_1 \cup \mathcal{E}_2) + \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) - \mathbb{P}(\mathcal{E}_2). \tag{2}$$

Now,  $\mathcal{E}_1 \cup \mathcal{E}_2$  occurs precisely if  $\mathcal{F}_1 \cap \mathcal{F}_2$  occurs (since  $\mathcal{E}_1 \cup \mathcal{E}_2$  is the event that  $B \in c(\mathcal{T}_X)$  and either  $A \in c(\mathcal{T}_X)$  or  $A \cup B \in c(\mathcal{T}_X)$ ). Thus

$$\mathbb{P}(\mathcal{E}_1 \cup \mathcal{E}_2) = \mathbb{P}(\mathcal{F}_1 \mid \mathcal{F}_2) \cdot \mathbb{P}(\mathcal{F}_2).$$

Combining this equation with (2), and noting that  $\mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) = \hat{p}_n(A, B)$  and  $p_n(A, B) = \mathbb{P}(\mathcal{E}_1)$ , we obtain

$$p_n(A, B) = \mathbb{P}(\mathcal{E}_1) = \mathbb{P}(\mathcal{F}_1 \mid \mathcal{F}_2) \cdot \mathbb{P}(\mathcal{F}_2) - \mathbb{P}(\mathcal{E}_2) + \hat{p}_n(A, B). \quad (3)$$

Now, by (GE),

$$\mathbb{P}(\mathcal{F}_1 \mid \mathcal{F}_2) = \mathbb{P}(A \in c(\mathcal{T}_{X-B})) = p_{n-b}(a), \quad (4)$$

and

$$\begin{aligned} \mathbb{P}(\mathcal{E}_2) &= \mathbb{P}(B \in c(\mathcal{T}_X) \mid A \cup B \in c(\mathcal{T}_X)) \cdot \mathbb{P}(A \cup B \in c(\mathcal{T}_X)) \\ &= p_{a+b}(b) \cdot p_n(a+b). \end{aligned} \quad (5)$$

Thus, substituting (4) and (5) and the equality  $\mathbb{P}(\mathcal{F}_2) = p_n(b)$  into (3), we obtain

$$p_n(A, B) = p_{n-b}(a) \cdot p_n(b) - p_{a+b}(b) \cdot p_n(a+b) + \hat{p}_n(a, b).$$

Case 5 now follows from Lemmas 4.2 and 4.4. Case 6 follows from the compatibility condition (1) for clades.  $\square$

We now ask whether the events ‘A is a clade’ and ‘B is a clade’ are positively or negatively correlated under the YHK process. Let  $\eta_A$  (respectively  $\eta_B$ ) be the Bernoulli (0, 1) random variables that take the value 1 if A (respectively B) is a clade of a YHK tree  $\mathcal{T}$  on  $X_n$ , and let  $\rho_n(A, B)$  denote the correlation coefficient of  $\eta_A$  and  $\eta_B$ , which is given by

$$\rho_n(A, B) = \frac{p_n(A, B) - p_n(A)p_n(B)}{\sqrt{p_n(A)(1 - p_n(A))p_n(B)(1 - p_n(B))}}.$$

**Corollary 4.6.** For any two strict subsets A, B of X, the correlation  $\rho_n(A, B)$  is

- strictly negative, if A, B are not compatible, and undefined if  $|A| = 1$  or  $|B| = 1$ , and
- strictly positive, otherwise.

**Proof.** If A and B are not compatible, then  $p_n(A, B) = 0$ , but both  $p_n(A)$  and  $p_n(B)$  are greater than zero, and so  $\rho_n(A, B) < 0$ . If  $|A| = 1$ , then  $p_n(A) = 1$  and  $p_n(A, B) = p_n(B)$  (regardless of whether A is a subset of B or is disjoint from B). Thus the numerator and denominator of  $\rho_n(A, B)$  are both zero. A similar argument holds if  $|B| = 1$ .

In the remaining cases, we consider the ratio  $p_n(A, B)/(p_n(A)p_n(B))$ . For example, in Case 2, we have

$$\frac{p_n(A, B)}{p_n(A) \cdot p_n(B)} = \frac{(n-1) \cdots (n-a+1)}{(b-1) \cdots (b-a+1)}.$$

This is strictly  $> 1$  since  $\frac{n-1}{b-1} > 1, \dots, \frac{n-a+1}{b-a+1} > 1$ . Similar arguments apply in the other cases; however, Case 5 requires some detailed algebraic manipulation.  $\square$

Fig. 3 illustrates the correlation coefficient  $\rho_n(A, B)$  for  $n = 25$  in Cases 2, 4 and 5. Notice that the correlation is typically much smaller in Cases 2 and 5 than in Case 4.

### 5. Extension to partitions of X

Suppose that the collection of sets  $A_1, A_2, \dots, A_k$  forms a partition of X, and let  $a_i = |A_i|$ , for  $i = 1, \dots, k$ , so that  $n = |X| = \sum_{i=1}^k a_i$ . For a rooted YHK tree  $\mathcal{T}$ , let  $p(a_1, \dots, a_k)$  be the probability that  $A_1, A_2, \dots, A_k$  are clades of  $\mathcal{T}$  (this probability depends only on the cardinality of the sets by the exchangeability property). For example,  $p(2, 2, 2) = 2/225$ , and from Lemma 4.3 we have  $p(a_1, a_2) = \frac{2}{a_1+a_2-1} \binom{a_1+a_2}{a_1}^{-1}$ . Our aim in this section is to generalize this to larger values of k. In order to do so, we describe a new result for the Yule model, which requires a further definition.

For a rooted YHK tree  $\mathcal{T}$ , and a rooted phylogenetic tree  $T_k$  with leaf set  $\{1, \dots, k\}$ , let  $p(a_1, \dots, a_k; T_k)$  be the probability that  $A_1, A_2, \dots, A_k$  are clades of  $\mathcal{T}$  and that  $T_k$  is the tree obtained from  $\mathcal{T}$  by replacing each clade  $A_i$  by a single leaf labelled  $i$ , for  $i = 1, \dots, k$ . Let  $\mathcal{I}(T_k)$  denote the set of interior vertices of  $T_k$ .

**Theorem 5.1.** For  $k > 1$ , we have

$$(i) \quad p(a_1, \dots, a_k; T_k) = \frac{2^{k-1} \prod_{i=1}^k a_i!}{n!} \prod_{v \in \mathcal{I}(T_k)} \left( \frac{1}{\sum_{i=1}^k a_i I_v(A_i) - 1} \right),$$

where  $I_v(A_i)$  is the indicator variable that takes the value 1 if  $i$  is a descendant of  $v$  in  $T_k$  and 0 otherwise;

$$(ii) \quad p(a_1, \dots, a_k) = \sum_{T_k} p(a_1, \dots, a_k; T_k),$$

where the summation is over all distinct rooted binary phylogenetic trees on leaf set  $\{1, \dots, k\}$ .

**Proof.** We prove the result by induction on k. For  $k = 2$ , Lemma 4.3 gives  $p(a_1, a_2; T_2) = \hat{p}_n(a_1, a_2) = \frac{2}{n-1} \binom{n}{a}^{-1}$ , where  $n = a_1 + a_2$ , which agrees with the expression given in part (i) with  $k = 2$ .

Now, suppose that part (i) holds whenever k is less than or equal to  $m \geq 2$ ; we will show that it also holds when  $k = m + 1$ . Thus, suppose that we have a collection  $C = \{A_1, \dots, A_{m+1}\}$  that partitions X, and also have a rooted binary phylogenetic tree  $T_{m+1}$  on leaf set  $\{1, \dots, m+1\}$ . Then  $T_{m+1}$  has a cherry (two leaves adjacent to the same vertex). Without loss of generality (by reordering the sets if necessary), we may suppose that these two leaves are  $m$  and  $m+1$ . Consider the collection of  $m$  sets obtained from C by replacing  $A_m$  and  $A_{m+1}$  by their union, and let  $T'$  be the tree obtained from  $T_{m+1}$  by deleting the leaves  $m$  and  $m+1$  along with their incident edges and labelling the exposed vertex by  $m$ . Notice that  $T'$  is a rooted binary phylogenetic tree that has leaf set  $\{1, \dots, m\}$ . By the exchangeability and group elimination (via sampling consistency) properties, we have, for  $a'_m := a_m + a_{m+1}$ , the following identity:

$$p(a_1, \dots, a_{m+1}; T_{m+1}) = p(a_1, \dots, a'_m; T') \cdot \hat{p}_{a'_m}(a_m, a_{m+1}),$$

where  $\hat{p}_{a'_m}(a_m, a_{m+1})$  is the probability that a YHK tree on leaf set  $A_m \cup A_{m+1}$  has  $A_m$  and  $A_{m+1}$  as sister (and thus maximal) clades. Applying the induction hypothesis for the first term on the right-hand side of this equation, namely  $p(a_1, \dots, a'_m; T')$ , and applying Lemma 4.3 for the second term, and collecting terms, leads to the expression in part (i) for  $k = m + 1$ , and thereby justifies the induction step.

Part (ii) follows by observing that each tree  $\mathcal{T}$  that has  $A_1, \dots, A_k$  as clades has one (and only one) associated tree  $T_k$ , and so these trees provide a partition of the event for which the probability is given by  $p(a_1, \dots, a_k)$ .  $\square$

As an illustration of Theorem 5.1, we have the following result for  $k = 3$ :

$$p(a_1, a_2, a_3) = \frac{4a_1!a_2!a_3!}{n!(n-1)} \left[ \sum_{i=1}^3 \frac{1}{n-a_i-1} \right],$$

where  $n = a_1 + a_2 + a_3$ .

We note that, as well as being a generalization of Lemma 4.3 to  $k > 2$ , Theorem 5.1(i) also generalizes the classic result that the probability that a YHK tree  $\mathcal{T}$  has a given tree topology  $T_k$  is  $\frac{2^{n-1}}{k!} \prod_{v \in \mathcal{I}(T_k)} \binom{1}{n_v-1}$ , where  $n_v$  is the number of leaves of  $T_k$  below  $v$  (see Brown, 1994 or Semple and Steel, 2003). This can be seen by setting  $a_1 = a_2 = \dots = a_n = 1$  in Theorem 5.1(i).

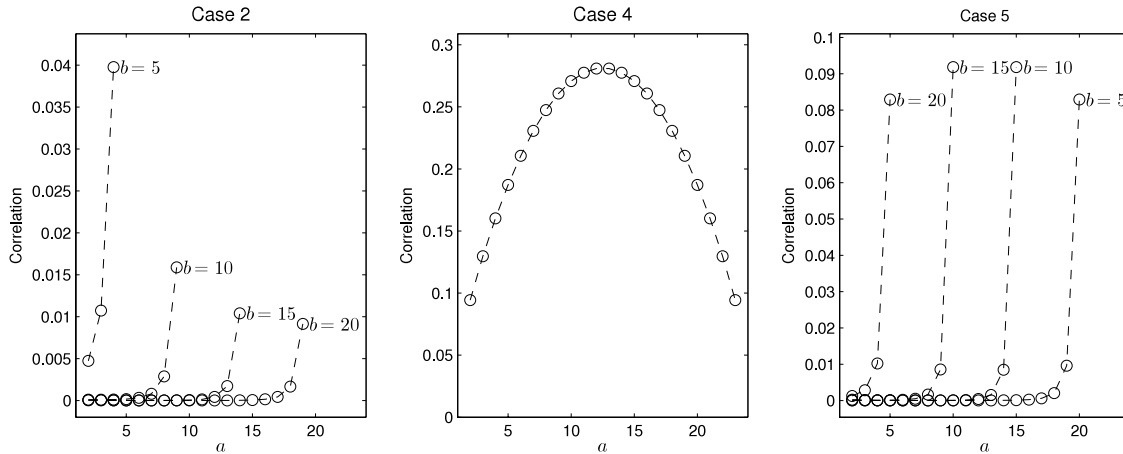


Fig. 3. Graphs of  $\rho_n(A, B)$  for  $n = 25$ , in Cases 2, 4 and 5, with  $a = |A|$  and  $b = |B|$ .

6. Extension to unrooted trees

If we suppress the root  $\rho$  of a rooted binary phylogenetic  $X$ -tree  $T$ , we obtain an unrooted binary phylogenetic  $X$ -tree, which we will denote as  $T^{-\rho}$  (as shown in Fig. 1(b)). Following (Wilkinson et al., 2007) (see also Lapointe et al., 2010), we say that a subset  $A$  of  $X$  is a *clan* of an unrooted phylogenetic  $X$ -tree  $T'$  if  $A|(X - A)$  is a split of  $T'$ . Note that any clade of the rooted tree  $T$  becomes a clan of  $T^{-\rho}$ . However, this latter tree also has additional clans that do not correspond to a clade of  $T$ . The precise relationship is given as follows.

**Lemma 6.1.** Given a rooted binary  $X$ -tree,  $T$ , a set  $A$  is a clan of  $T^{-\rho}$  if and only if either  $A$  is a clade of  $T$  or  $X - A$  is a clade of  $T$ . □

Now, suppose that the rooted phylogenetic tree  $T$  is generated under the YHK process. Then we obtain an induced probability for the unrooted tree  $T^{-\rho}$ . Note that the same unrooted tree can arise from different rootings. This probability distribution on unrooted phylogenetic trees can also be described directly as a Yule-type process on unrooted trees in which, at each stage, a leaf is selected uniformly at random and a new leaf (with a random label) is attached to its incident edge (see, e.g., Steel and Penny, 1993). Fig. 4 illustrates how different leaf choices in this process lead to different shapes of unrooted trees.

For a strict non-empty subset  $A$  of  $X_n$ , let  $q_n(A)$  be the probability that  $A$  is a clan of the unrooted YHK tree on leaf set  $X_n$ ; by the exchangeability property (EP), this depends only on  $a = |A|$  and  $n$ , so we will also write it as  $q_n(a)$ .

**Lemma 6.2.**

$$q_n(a) = 2n \left[ \frac{1}{a(a+1)} + \frac{1}{b(b+1)} - \frac{1}{(n-1)n} \right] \binom{n}{a}^{-1},$$

where  $a = |A|, b = n - a$ .

**Proof.** By Lemma 6.1, we have

$$q_n(A) = p_n(A) + p_n(X - A) - p_n(A, X - A).$$

Applying Lemmas 4.2 and 4.3, and noting that  $p_n(A, X - A) = \hat{p}_n(A, X - A)$ , leads to the claimed equation. □

Now, consider two disjoint subsets  $A$  and  $B$  of  $X$ , and let  $q_n(A, B)$  be the probability that  $A$  and  $B$  are both clans of the unrooted YHK tree on leaf set  $X_n$ . By exchangeability (EP), this probability depends only on  $a = |A|, b = |B|$ , and  $n$ , and so we will denote it as  $q_n(a, b)$ . As an example, we have

$$q_6(2, 2) = 7/225.$$

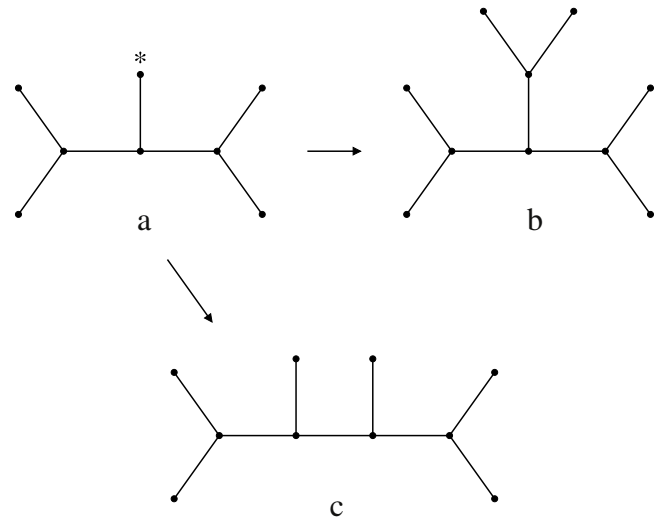


Fig. 4. Only one unrooted binary tree shape is possible with five leaves (a), but two are possible with six leaves (b, c). If the ‘central’ leaf (\*) of tree (a) is split to form two leaves, then we obtain tree shape (b), while splitting any one of the remaining four leaves produces tree shape (c). Thus, tree shape (b) has probability 1/5. Since there are  $6!/3!2^3 = 15$  distinct ways to label its leaves, each of the resulting phylogenetic trees has probability 1/75. By contrast, any phylogenetic tree of shape (c) has probability  $4/5 \times 1/90 = 2/225$ .

To see this, observe that, if we take (say)  $A = \{1, 2\}, B = \{3, 4\}$ , then, referring to Fig. 4, there is just one tree of shape (b) and two of shape (c) that has both clans  $A$  and  $B$ . Thus,  $q_6(2, 2) = 1 \times \frac{1}{75} + 2 \times \frac{2}{225}$ . We now give an exact analytical formula for  $q_n(a, b)$ .

**Theorem 6.3.**

(i) If  $a + b = n$ , then

$$q_n(a, b) = q_{a+b}(A) = \frac{2a!b!}{(a+b-1)!} \left[ \frac{1}{a(a+1)} + \frac{1}{b(b+1)} - \frac{1}{(a+b)(a+b-1)} \right].$$

(ii) If  $a + b < n$  then

$$q_n(a, b) = r_n(a, b) + R_n(a, n - b) + R_n(b, n - a) - \hat{p}_n(b, n - b)p_{n-b}(a) - \hat{p}_n(a, n - a)p_{n-a}(b),$$

where the first three quantities are as given in Theorem 4.5 (Cases 2, 3 and 5), while the last two terms are given by Lemmas 4.2 and 4.3.

**Proof.** Part (i) follows from Lemma 6.2, noting that  $n = a + b$ . For part (ii), Lemma 6.1 implies that  $A$  and  $B$  are clans of  $\mathcal{T}^{-\rho}$  precisely if one of the following three events occurs:

- (a)  $A$  and  $B$  are clades of  $T$ ;
- (b)  $A$  and  $X - B$  are clades of  $T$ , but  $B$  is not a clade of  $T$ ;
- (c)  $B$  and  $X - A$  are clades of  $T$ , but  $A$  is not a clade of  $T$ .

(Note that  $X - A$  and  $X - B$  cannot both be clades of  $T$ , by the compatibility condition (1), since  $(X - A) \cap (X - B) \neq \emptyset$  by the assumption that  $a + b < n$ , and since  $X - A$  neither contains nor is contained in  $X - B$ .) Moreover, the three events (a), (b), and (c) are mutually exclusive, by virtue of the assumption that  $A$  and  $B$  are disjoint and their union is a strict subset of  $X$ . The probability of Event (a) is  $r_n(a, b)$ , while the probability of Event (b) is  $R_n(a, n - b) - \hat{p}_n(b, n - b)p_{n-b}(a)$ , since the first term is the probability that  $A$  and  $X - B$  are clades of  $\mathcal{T}$ , and  $\hat{p}_n(b, n - b)p_{n-b}(a)$  is the probability that  $A, X - B$ , and  $B$  are clades of  $\mathcal{T}$ . Similarly,  $R_n(b, n - a) - \hat{p}_n(a, n - a)p_{n-a}(b)$  is the probability of Event (c). The result now follows by adding the probabilities of these three mutually exclusive events.  $\square$

6.1. Extensions of the clan condition (I)

For a pair  $A, B$  of disjoint subsets of  $X$ , a weaker condition than requiring that  $A$  and  $B$  are both clans of  $\mathcal{T}^{-\rho}$  is simply to require that at least one edge of this tree separates  $A$  from  $B$ . Let  $Q_n(A, B)$  be the probability of this event for an unrooted YHK tree on the leaf set  $X_n$ . Note that the sampling consistency property (SC) for  $\mathcal{T}$  implies an analogous property for  $\mathcal{T}^{-\rho}$ . Namely, for any subset  $W$  of  $X$  (the leaf set of  $\mathcal{T}$ ), the unrooted tree obtained from  $\mathcal{T}^{-\rho}$  by restricting the leaf set to  $W$  has the same distribution as  $(\mathcal{T}_W)^{-\rho}$ . Accordingly, if we apply this with  $W = A \cup B$ , we obtain the following result:

$$Q_n(A, B) = q_{a+b}(A), \tag{6}$$

where  $q_{a+b}(A)$  is given by Theorem 6.3(i).

6.2. Extensions of the clan condition (II)

We now describe a second extension. Suppose that  $A_1, A_2, \dots, A_k$  partition  $X$ , and, as usual, let  $a_i = |A_i|$ . For an unrooted YHK tree  $\mathcal{T}$ , let  $q(a_1, \dots, a_k)$  be the probability that  $A_1, A_2, \dots, A_k$  are clans of  $\mathcal{T}$ , and let  $q'(a_1, \dots, a_k)$  be the probability that  $A_1, A_2, \dots, A_k$  are convex on  $\mathcal{T}$  (that is, the minimal subtree connecting the leaves in  $A_i$  is vertex disjoint from the minimal subtree connecting the leaves in  $A_j$  for all pairs  $i, j$ ; see Semple and Steel (2003) for further details and the biological significance of convexity).

We have calculated  $q$  when  $k = 2$  above (and  $q' = q$  in this case). We turn now to the next case of interest,  $k = 3$ , where, for example, we have

$$q(2, 2, 2) = 1/75, \quad \text{and} \quad q'(2, 2, 2) = 1/15.$$

The following result provides an exact formula for these two quantities for arbitrary  $(a_1, a_2, a_3)$ .

**Theorem 6.4.** Let  $n = a_1 + a_2 + a_3$ . Then

- (i)  $q(a_1, a_2, a_3) = \frac{4a_1!a_2!a_3!}{(n-1)!} \left[ \sum_{i=1}^3 \frac{1}{(n-a_i)((n-a_i)^2-1)} \right]$ .
- (ii)  $q'(a_1, a_2, a_3) = q_n(a_1, a_2) + q_n(a_1, a_3) + q_n(a_2, a_3) - 2q(a_1, a_2, a_3)$ , where  $q_n(a_i, a_j)$  is given in Theorem 6.3(ii), and  $q(a_1, a_2, a_3)$  is from part (i).

**Proof.** For part (i), the event that  $A_1, A_2$ , and  $A_3$  (which partition  $X$ ) are clans of  $\mathcal{T}^{-\rho}$  is the union of three disjoint events  $E_{jk}$  over the three choices of  $\{j, k\} \in \{\{1, 2\}, \{1, 3\}, \{2, 3\}\}$ , where  $E_{jk}$  is the event that the union of two of the sets – say  $A_j$  and  $A_k$  – must be a clade of  $\mathcal{T}$ , and that this clade has maximal clades  $A_j$  and  $A_k$ . The exchangeability and group elimination conditions then give

$$q(a_1, a_2, a_3) = \mathbb{P}(E_{12}) + \mathbb{P}(E_{13}) + \mathbb{P}(E_{23}) \\ = \sum_{i=1}^3 p_n(n - a_i) \cdot \hat{p}_{a_j+a_k}(a_j, a_k),$$

where  $\{a_i, a_j, a_k\} = \{1, 2, 3\}$  in the term on the right-hand side of this last equation. By Lemmas 4.2 and 4.3, this gives

$$q(a_1, a_2, a_3) \\ = \sum_{i=1}^3 \frac{2n}{(n - a_i)(n - a_i + 1)} \frac{(n - a_i)!a_i!}{n!} \cdot \frac{2}{(n - a_i - 1)} \frac{a_j!a_k!}{(n - a_i)!},$$

which simplifies to the expression given in (ii).

For part (ii), the event that  $A_1, A_2$ , and  $A_3$  are convex on  $\mathcal{T}^{-\rho}$  is the union of three (non-disjoint!) events  $E'_{jk}$  over the three choices of  $\{j, k\} \in \{\{1, 2\}, \{1, 3\}, \{2, 3\}\}$ , where  $E'_{jk}$  is the event that two of the sets – say  $A_j$  and  $A_k$  – are clans of  $\mathcal{T}^{-\rho}$ . Note that the intersection of any two (or three) of these three events is simply the event that all three sets are clans of  $\mathcal{T}$ , which was dealt with in part (i). Thus, by the principle of inclusion and exclusion, we have

$$q'(a_1, a_2, a_3) = \mathbb{P}(E'_{12}) + \mathbb{P}(E'_{13}) + \mathbb{P}(E'_{23}) - 2q(a_1, a_2, a_3),$$

and the result in part (iii) now follows.  $\square$

Deriving explicit formulae for  $q(a_1, \dots, a_k)$  and  $q'(a_1, \dots, a_k)$  for  $k > 3$  is, in principle, possible, but the formulae quickly become complex.

6.3. Extensions of the clan condition (III)

A third extension is to consider the probability  $Q_n(A_1, A_2)$  that two sets  $A_1$  and  $A_2$  are clans of a YHK tree on  $n$  leaves when these two sets are *not* disjoint. For this setting, we have the following result.

**Proposition 6.5.** Suppose that  $A_1$  and  $A_2$  are non-disjoint subsets of  $X$ , and that  $a_i = |A_i|$ .

- (i) If  $A_1 \subset A_2$ , then

$$Q_n(A_1, A_2) = q_n(a_1, n - a_2),$$

where  $q_n(*, *)$  is given by Theorem 6.3. Similarly, if  $A_2 \subset A_1$ , then  $Q_n(A_1, A_2) = q_n(n - a_1, a_2)$ .

- (ii) Otherwise, if neither set  $A_1, A_2$  is a subset of the other, then

$$Q_n(A_1, A_2) = \begin{cases} q_n(a_1 - a_{12}, a_2 - a_{12}), & \text{if } A_1 \cup A_2 = X, \\ 0, & \text{otherwise,} \end{cases}$$

where  $a_{12} = |A_1 \cap A_2|$ .

**Proof.** First, observe that, if  $A_1 \subset A_2$ , then  $A_1$  and  $A_2$  are clans of an unrooted phylogenetic  $X$ -tree  $T$  if and only if  $A_1$  and  $X - A_2$  are clans of  $T$ . Noting that these are disjoint sets, the first part of Proposition 6.5 follows from Theorem 6.3. For the second case, where neither set  $A_1$  nor  $A_2$  is a subset of the other, first observe that in order for  $A_1$  and  $A_2$  to be clans of the same unrooted phylogenetic  $X$ -tree  $T$  a necessary condition is that  $A_1 \cup A_2 = X$ . Moreover, under this condition,  $A_1$  and  $A_2$  are clans of  $T$  if and only if  $A_1 - A_1 \cap A_2$  and  $A_2 - A_1 \cap A_2$  are clans of  $T$ ; as these are disjoint sets, the second part of Proposition 6.5 follows from Theorem 6.3.  $\square$

## 7. Discussion

The arguments we have used in our analysis have primarily relied on repeated application of the properties of exchangeability (EP) and group elimination (GE) (or its corollary, sampling consistency (SC)) for the YHK model, together with Lemmas 4.2 and 4.3. However, other natural models for trees can also satisfy some of these properties. Indeed, the distribution that assigns each rooted binary phylogenetic tree on  $X_n$  the same probability (sometimes known as the proportional to distinguishable arrangements model, or the PDA model) satisfies both exchangeability and group elimination (Aldous, 1995). This suggests that by finding and applying the corresponding results to Lemmas 4.2 and 4.3 for the PDA model, one could develop a parallel line of results for the PDA model to most of the analysis we have provided in this paper for the YHK model.

Unfortunately, only one other model, apart from PDA and YHK, is known to satisfy both exchangeability and group elimination, and this model is not of biological interest, as it only generates pectinate (comb-like) tree shapes. Aldous (1995) has conjectured that these are the *only* three distributions on rooted binary phylogenetic trees that satisfy both exchangeability and group elimination. Nonetheless, it may be of interest to explore models that satisfy weakened assumptions—for example, the exchangeability property (EP) and (SC), or just exchangeability by itself.

Even with exchangeability alone, one can devise meaningful statistical significance tests. For example, suppose that  $N$  taxa include one or more particular (disjoint) subsets (different ‘types’ of taxa)  $A_1, A_2, \dots, A_k$ , where  $k \geq 1$ . Consider any model for generating a rooted binary tree that satisfies the exchangeability property (EP), and let  $p_n$  be the probability that a tree on this set of taxa as leaves, generated under this model, has at least one clade of size at least  $n$  consisting of just one type (i.e., all leaves in the clade are a subset of one of the sets  $A_1, \dots, A_k$ ). Then we have the following result, the proof of which is given in the Appendix.

**Proposition 7.1.** *For any probability distribution on rooted binary trees satisfying the exchangeability property (EP), we have*

$$p_n \leq \sum_{i=1}^k \sum_{m=n}^{a_i} \frac{\binom{a_i}{m}}{\binom{N-1}{m-1}},$$

where  $a_i = |A_i|$ .

As a simple example, suppose that we have  $N = 40$  taxa, including two disjoint groups, each containing six taxa. For a tree generated under any model that satisfies the exchangeability property, the probability that this tree would contain a clade of size four or larger consisting entirely of taxa from one of the two groups is, at most,

$$2 \cdot \left( \frac{\binom{6}{4}}{\binom{39}{3}} + \frac{\binom{6}{5}}{\binom{39}{4}} + \frac{\binom{6}{6}}{\binom{39}{5}} \right) < 0.005.$$

We end with a caution. In applying Proposition 7.1 as a part of significance test using a given phylogenetic tree, it is important that the groups  $A_1, A_2, \dots$  are specified *a priori*, and not identified based on the data used to construct the given tree (or, indeed, the tree itself).

## Acknowledgments

We thank the Royal Society of NZ (Marsden Fund and James Cook Fellowship) and the Allan Wilson Centre for Molecular Ecology and Evolution for funding this work. We also thank Arne Mooers and Tyler Kuhn for comments that motivated the development of Proposition 7.1, and Erick Matsen and two anonymous reviewers for helpful comments on an earlier version of this manuscript.

## Appendix. Proof of Proposition 7.1

Let  $X_{m,i}$  be the number of clades of size  $m$  in the randomly generated tree that has the property that the taxa are all of type  $A_i$ , and let  $X := \sum_{i=1}^k \sum_{m=n}^{a_i} X_{m,i}$ . Then  $p_n = \mathbb{P}(X > 0)$ . Since  $X$  is a non-negative integer random variable, we have

$$\mathbb{P}(X > 0) \leq \mathbb{E}[X]. \quad (7)$$

By linearity of expectation, we have

$$\mathbb{E}[X] = \sum_{i=1}^k \sum_{m=n}^{a_i} \mathbb{E}[X_{m,i}]. \quad (8)$$

Moreover,

$$\mathbb{E}[X_{m,i}] = \sum_t \mathbb{E}[X_{m,i}|t] \mathbb{P}(t), \quad (9)$$

where the summation is over all binary tree shapes on the given leaf set of size  $N$ ,  $\mathbb{E}[X_{m,i}|t]$  is the conditional expectation of  $X_{m,i}$  given that  $t$  is the tree shape generated by the random speciation process, and  $\mathbb{P}(t)$  is the probability of generating tree shape  $t$ . For any given the tree shape  $t$ ,

$$\mathbb{E}[X_{m,i}|t] = \sum_{v:n_v=m} \mathbb{E}[I_{v,i}|t], \quad (10)$$

where the summation is over all the interior vertices of  $t$  for which the number of leaves below  $v$  ( $n_v$ ) is  $m$ , and where  $I_{v,i}$  is the binary random variable that takes the value 1 precisely if all the leaves below  $v$  are of type  $A_i$ , and  $I_{v,i} = 0$  otherwise. Now, by exchangeability, we have the following identity for any vertex  $v$  of  $t$  with  $n_v = m$ :

$$\mathbb{E}[I_{v,i}|t] = \mathbb{P}(I_{v,i} = 1|t) = \frac{\binom{a_i}{m}}{\binom{N}{m}}. \quad (11)$$

Now any tree shape on  $N$  leaves has, at most,  $N/m$  vertices  $v$  for which  $n_v = m$ , and so we obtain, from (10) and (11),  $\mathbb{E}[X_{m,i}|t] \leq \frac{N}{m} \cdot \frac{\binom{a_i}{m}}{\binom{N}{m}} = \frac{\binom{a_i}{m}}{\binom{N-1}{m-1}}$ . Since this inequality holds for all tree shapes

$t$ , Eq. (9) implies that  $\mathbb{E}[X_{m,i}] \leq \frac{\binom{a_i}{m}}{\binom{N-1}{m-1}}$ . The expression for  $p_n$  now follows from Eqs. (7) and (8).  $\square$

## References

- Aldous, D., 1995. Probability distributions on cladograms. In: Aldous, D., Pemantle, R. (Eds.), *Random Discrete Structures*. In: IMA Volumes in Mathematics and its Applications, vol. 76. Springer, pp. 1–18.
- Blum, M.G.B., Francois, O., 2005. External branch length and minimal clade size under the neutral coalescent. *Adv. in Appl. Probab.* 37, 647–662.
- Brown, J.K.M., 1994. Probabilities of evolutionary trees. *Syst. Biol.* 43, 8–91.
- Cummings, M.P., Neel, M.C., Shaw, K.L., 2008. A genealogical approach to quantifying lineage divergence. *Evolution* 62, 2411–2422.
- DeQuieroz, K., 2007. Species concepts and species delimitation. *Syst. Biol.* 56, 879–886.
- Harding, E.F., 1971. The probabilities of rooted tree shapes generated by random bifurcation. *Adv. in Appl. Probab.* 3, 44–77.
- Heard, S.B., 1992. Patterns in tree balance among cladistic, phenetic, and randomly generated phylogenetic trees. *Evolution* 46, 1818–1826.
- Hudson, R.R., Coyne, J.A., 2002. Mathematical consequences of the genealogical species concept. *Evolution* 56, 1557–1565.
- Kingman, J.F.C., 1982. On the genealogy of large populations. *J. Appl. Probab.* 19A, 27–43.
- Lapointe, F.-J., Lopez, P., Boucher, Y., Koenig, J., Baptiste, E., 2010. Cladistics: a multi-level perspective for harvesting unrooted gene trees. *Trends Microbiol.* 18 (8), 341–347.
- Rabosky, D.L., Lovette, I.J., 2008. Explosive evolutionary radiations: decreasing speciation or increasing extinction through time? *Evolution* 62, 1866–1875.
- Rosenberg, N.A., 2003. The shapes of neutral gene genealogies in two species: probabilities of monophyly, paraphyly, and polyphyly in a coalescent model. *Evolution* 57, 1465–1477.
- 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.

- Rosenberg, N.A., 2007. Statistical tests for taxonomic distinctiveness from observations of monophyly. *Evolution* 61, 317–323.
- Semple, C., Steel, M., 2003. *Phylogenetics*. Oxford University Press.
- Steel, M.A., Penny, D., 1993. Distributions of tree comparison metrics—some new results. *Syst. Biol.* 42, 126–141.
- Wilkinson, M., McInerney, J.O., Hirt, R.P., Foster, P.G., Embley, T.M., 2007. Of clades and clans: terms for phylogenetic relationships in unrooted trees. *Trends Ecol. Evol.* 22, 114–115.
- Yule, G.U., 1925. A mathematical theory of evolution. Based on the Conclusions of Dr. J.C. Willis, F.R.S.. *Phil. Trans. R. Soc.* 213, 21–87.