

---

# Genealogical properties of subsamples in highly fecund populations

Bjarki Eldon  
Museum für Naturkunde  
43 Invalidenstraße  
10115 Berlin  
Germany

Fabian Freund  
University of Hohenheim  
Institute 350b  
Fruwirthstraße 21  
D-70599 Stuttgart, Germany

Alison M. Etheridge  
University of Oxford  
Department of Statistics  
24–29 St Giles  
OX1 3LB Oxford, UK

July 17, 2017

**Abstract** We consider some genealogical properties of nested samples. The complete sample is assumed to have been drawn from a natural population characterised by high fecundity and sweepstakes reproduction (HFSR). The random gene genealogies of the samples are modeled by random trees which allow for multiple mergers of ancestral lineages looking back in time. In contrast, the classical Kingman coalescent only admits asynchronous pairwise mergers of ancestral lineages. The pattern of genetic diversity observed in samples from HFSR populations differs strongly from expectations based on Kingman's  $n$ -coalescent. Among the genealogical properties we consider are the probability that the complete sample and the nested subsample share the most recent common ancestor; we also compare lengths of 'internal' branches of nested genealogies. The results indicate how 'informative' a subsample is about the properties of the larger complete sample from which the subsample is drawn, and by implication how much information is gained by increasing the sample size.

keywords: coalescent; high fecundity; nested samples; multiple mergers; time to most recent common ancestor

AMS subject classification: 92D15, 60J28

## Contents

1	Introduction . . . . .	2
2	Sharing the MRCA . . . . .	6
3	Relative times and lengths . . . . .	13
4	Proofs . . . . .	23
5	Conclusion and open questions . . . . .	31
A1	Coalescent processes . . . . .	36
A2	Goldschmidt and Martin's construction of the Bolthausen-Sznitman $n$ -coalescent . . . . .	38

## 1 Introduction

The study of the evolutionary history of natural populations usually proceeds by drawing inference from a random sample of DNA sequences. To this end the coalescent approach initiated by [55, 57, 56, 82, 51] - i.e. the probabilistic modeling of the random ancestral relations of the sampled DNA sequences - has proved to be very useful [84, cf.]. Inference based on the coalescent relies on the key assumption, as in standard statistical inference, that the evolutionary history of the (finite) sample approximates, or is informative about, the evolutionary history of the population from which the sample is drawn. We would like to know how much some basic genealogical sample-based statistics tell us about the population in a multiple-merger coalescent framework. Does the 'informativeness' of the various genealogical statistics depend on the underlying coalescent process? A more practical approach to this question is, instead of comparing a sample with the population, to ask how much of the genetic information of a sample is already contained in a subsample, i.e. what is gained by enlarging the sample? A related question concerns the size of the sample; i.e. how large does our sample need to be for a reliable inference? We approach these problems by studying some genealogical properties of nested samples, by which we mean where a sample (a subsample)

is drawn (uniformly at random without replacement) from a larger sample (the complete sample). By way of an example, [73] consider nested samples whose ancestries are governed by the Kingman  $n$ -coalescent [55, 57, 56]. One of the results of [73] concerns the probability that a subsample shares its most recent common ancestor (MRCA) with the complete sample. If the subsample shares the MRCA with the complete sample, it also, with high probability, shares the most ancient genealogical information. Thus, the effects on the genetic structure of the complete sample of this most ancient part of the genealogy are also present in the subsample. In addition, the complete sample and the subsample have had exactly the same timespan to collect mutations. [73] show that the probability that a subsample of a fixed size  $m$  shares the MRCA with the complete sample of arbitrarily large size  $n$  ( $n \rightarrow \infty$ ) converges to  $(m-1)/(m+1)$ . Even a subsample of size 2 shares the MRCA with probability  $1/3$ , while a sample of size 19 already shares with probability 0.9. This shows that by this measure (the probability of sharing the MRCA) even a rather small subsample drawn from a large complete sample whose ancestry is governed by the Kingman coalescent captures properties of the complete sample quite well. We will return to this example in Sec. 2.3.

Subsampling and the time to the MRCA (TMRCA) have biological applications. [49] uses subsampling to infer species trees. [6] compare methods for estimating TMRCA. Divergence times of populations can be inferred from estimates of TMRCA [88], especially if one has data from several populations [68]. [81] derives the distribution of the number of descendants of the MRCA of a subsample. [54] applies a Bayesian approach on data from unlinked loci to estimate coalescent times. Estimates of TMRCA are used to infer the time at which populations became ‘established’ [42]. The age of genes coding for functional objects such as proteins is also correlated with gene function, and therefore associations with diseases [21]. However, we will only be concerned with neutrally evolving non-coding and non-recombining segments of the genome.

A universal mechanism among all biological populations is reproduction and inheritance. Reproduction refers to the generation of offspring, and inheritance refers to the transmission of information necessary for viability and reproduction. Mendel’s laws on independent segregation of chromosomes into gametes describe the transmission of information from a parent to an offspring in a diploid population. For our purposes, however, it suffices to think of haploid populations where one can think of an individual as a single gene copy. By tracing gene copies as they are passed on from one generation to the next one automatically stores two sets of information. On the one hand one stores how frequencies of genetic types change going forwards in time; on the other hand one keeps track of the ancestral, or genealogical, relations among the different copies. This duality has been successfully exploited for example in modeling selection [34, 35]. To model genetic variation in natural populations one requires a mathematically tractable model of how genetic information is passed from parents to offspring. In the Wright-Fisher model offspring choose their parents independently and uniformly at random. Suppose we are tracing the ancestry of  $n \geq 2$  gene copies in a haploid Wright-Fisher population of  $N$  gene copies in total. For any pair, the chance that they have a common ancestor in the previous generation is  $1/N$ . Informally, we trace the genealogy of our gene copies on the order of  $\mathcal{O}(N)$  generations until we see the first merger, i.e. when at least 2 gene copies (or their ancestral lines) find a common ancestor.

If  $n$  is small relative to  $N$ , when a merger occurs, with probability  $1 - \mathcal{O}(1/N)$  it involves just two ancestral lineages. This means that if we measure time in units of  $N$  generations, and assume  $N$  is very large, the random ancestral relations of our sampled gene copies can be described by a continuous-time Markov chain in which each pair of ancestral lines merges at rate 1 and no other mergers are possible. We have, in an informal way, arrived at the Kingman-coalescent [55,57,56]. One can derive the Kingman-coalescent not just from the Wright-Fisher model but from any population model which satisfies certain assumptions on the offspring distribution [60,70,63]. These assumptions mainly dictate that higher moments of the offspring number distribution are small relative to (an appropriate power of) the population size. The Kingman-coalescent, and its various extensions, are used almost universally as the ‘null model’ for a gene genealogy in population genetics. The Kingman-coalescent is a remarkably good model for populations characterised by low fecundity, i.e. whose individuals have small numbers of offspring relative to the population size.

The classical Kingman-coalescent is derived from a population model in which the population size is constant between generations. Extensions to stochastically varying population size, in which the population size does not vary ‘too much’ between generations, have been made [53]; the result is a time-changed Kingman-coalescent. Probably the most commonly applied model of deterministically changing population size is the model of exponential population growth (see eg. [25,41,30]). In each generation the population size is multiplied by a factor  $(1 + \beta/N)$ , where  $\beta > 0$ . Therefore, the population size in generation  $k$  going forward in time is given by  $N_k = N_0(1 + \beta/N)^k$  where  $N_0$  is taken as the ‘initial’ population size. It follows that the population size  $[Nt]$  generations ago is  $Ne^{-\beta t}$ . [30] show that exponential population growth can be distinguished from multiple-merger coalescents (in which at least three ancestral lineages can merge), derived from population models of high fecundity and sweepstakes reproduction, using population genetic data from a single locus, provided that sample size and number of mutations (segregating sites) are not too small.

A diverse group of natural populations, including some marine organisms [45], fungi [1,79,50], and viruses [83] are highly fecund. By way of example, individual Atlantic codfish [59,67] and Pacific oysters [58] can lay millions of eggs. This high fecundity counteracts the high mortality rate among the larvae (juveniles) of these populations (Type III survivorship). The term ‘sweepstakes reproduction’ has been proposed to describe the reproduction mode of highly fecund populations with Type III survivorship [44]. Population models which admit high fecundity and sweepstakes reproduction (HFSR) through skewed or heavy-tailed offspring number distributions have been developed [63,64,78,31,72,52]. In the haploid model of [78], each individual independently contributes a random number  $X$  of juveniles where  $(C, \alpha > 0)$

$$\mathbb{P}(X \geq k) \sim \frac{C}{k^\alpha}, \quad k \rightarrow \infty, \quad (1)$$

and  $x_n \sim y_n$  means  $x_n/y_n \rightarrow 1$  as  $n \rightarrow \infty$ . The constant  $C > 0$  is a normalising constant, and the constant  $\alpha$  determines the skewness of the distribution. The next generation of individuals is then formed by sampling (uniformly without replacement) from the pool of juveniles. In the case  $\alpha < 2$  the random ancestral relations

of gene copies can be described by specific forms of multiple-merger coalescent processes [71].

Coalescent processes derived from population models of HFSR (see (1) for an example) admit multiple mergers of ancestral lineages [24, 69, 70, 75, 64, 71, 62]. Mathematically, we consider exchangeable  $n$ -coalescent processes, which are Markovian processes  $(\Pi_t^{(n)})_{t \geq 0}$  on the set of partitions of  $[n] := \{1, 2, \dots, n\}$  whose transitions are mergers of partition blocks (a ‘block’ is a subset of  $[n]$ , see Sec. A1) with rates specified in Sec. A1. The blocks of  $\Pi_t^{(n)}$  show which individuals in  $[n]$  share a common ancestor at time  $t$  measured from the time of sampling. Thus, the blocks of  $\Pi_t^{(n)}$  can be interpreted as ancestral lineages. The specific structure of the transition rates allows to treat a multiple-merger  $n$ -coalescent as the restriction of an exchangeable Markovian process  $(\Pi_t)_{t \geq 0}$  on the set of partitions of  $\mathbb{N}$ , which is called a multiple-merger coalescent (MMC) process. MMC processes are referred to as  $\Lambda$ -coalescents ( $\Lambda$  a finite measure on  $[0, 1]$ ) [24, 69, 70] if any number of ancestral lineages can merge at any given time, but only one such merger occurs at a time. By way of an example, if  $1 \leq \alpha < 2$  in (1) one obtains a so-called Beta( $2 - \alpha, \alpha$ )-coalescent [71] (Beta-coalescent, see Eq. (A37)). Processes which admit at least two (multiple) mergers at a time are referred to as  $\Xi$ -coalescents ( $\Xi$  a finite measure on the infinite simplex  $\Delta$ ) [75, 63, 64]. See Sec. A1 for details. Specific examples of these MMC processes have been shown to give a better fit to genetic data sampled from Atlantic cod [12, 18, 2, 16, 19] and Japanese sardines [66] than the classical Kingman-coalescent. See e.g. [29] for an overview of inference methods for MMC processes. [45] review the evidence for sweepstakes reproduction among marine populations and conclude ‘that it plays a major role in shaping marine biodiversity’.

MMC models arise in contexts other than high fecundity. [17] show that repeated strong bottlenecks in a Wright-Fisher population lead to time-changed Kingman-coalescents which look like  $\Xi$ -coalescents. [27, 28] show that the genealogy of a locus subjected to repeated beneficial mutations is well approximated by a  $\Xi$ -coalescent. [74] provides rigorous justification of the claims of [65, 22] that the genealogy of a population subject to repeated beneficial mutations can be described by the Beta-coalescent with  $\alpha = 1$  (also referred to as the Bolthausen-Sznitman coalescent [20]). These examples show that MMC processes are relevant for biology.

Overviews of mathematical population genetics can be found in a number of books (see e.g. [26, 36]), monographs, or papers. We refer the interested reader to e.g. [10, 25, 5, 33, 9, 13] for a more detailed background on coalescent theory. Our aim is to study how well some genealogical properties of a subsample capture those of the sample from which the subsample was drawn. We compare these properties between the classical Kingman-coalescent, exponential population growth, and  $\Lambda$ - and  $\Xi$ -coalescents using theoretical derivations and simulations. A precise description of MMC processes is given in Appendix (Sec. A1). For ease of reference we include a table (Table 1) of notation and terminology.

**Table 1** Notation and terminology.

symbol	explanation
HFSR	high fecundity and sweepstakes reproduction
MMC	multiple-merger coalescent
MRCA	most recent common ancestor
TMRC	time to MRCA
leaves	special kind of vertices in a random graph (genealogy); correspond to sampled DNA sequences
$n$ -coalescent	a coalescent process started from $n$ leaves
$\mathbb{N}$	the set of the natural numbers $\mathbb{N} := \{1, 2, \dots\}$
$[n]$	$[n] := \{1, 2, \dots, n\}$ , $n \in \mathbb{N}$
$[n]_a$	$[n]_a := \{a, a+1, \dots, n\}$ for $n, a \in \{0\} \cup \mathbb{N}$ , $a \leq n$
$\mathbb{1}_{(A)}$	$\mathbb{1}_{(A)} = 1$ if $A$ holds, and zero otherwise
$x \wedge y$	$\min\{x, y\}$
$T_{\text{MRCA}}^{(\infty)}$	the random TMRC of the population current at some stated time
$T_{\text{MRCA}}^{(n)}$	the random TMRC of a sample of size $n \in [2, \infty)$
$T_{\text{MRCA}}^{(m;n)}$	the random TMRC of a subsample of size $m$ taken from a complete sample of size $n > m$
$T_{\text{MRCA}}^{(M)}$	the random TMRC of a finite sample $M \subset \mathbb{N}$
$\Pi$	coalescent process; $\Pi \equiv \Pi_t := \{\Pi(t), t \geq 0\}$
$\Pi^{(n)}$	$\Pi$ restricted to $[n]$
$\Pi^{(\Lambda)}$	$\Lambda$ -coalescent
$\Pi^{(\Xi)}$	$\Xi$ -coalescent
$\mathbb{P}^{(\Pi)}(A)$	probability of event $A$ under $\Pi$
$p_{n,m}^{(\Pi)}$	$p_{n,m}^{(\Pi)} := \mathbb{P}^{(\Pi)}(T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)})$ ; the probability that subsample and complete sample share the MRCA
$\Delta$	the infinite simplex $\Delta := \{(x_1, x_2, \dots)   x_i \in [0, 1], \sum_{i \in \mathbb{N}} x_i \leq 1\}$
$\rho_T^{(m;n)}$	the ratio $T_{\text{MRCA}}^{(m;n)} / T_{\text{MRCA}}^{(n)}$ ; see Sec. 3.1
$\rho_l^{(m;n)}$	the ratio of ‘internal’ edge lengths between subsample and complete sample; see Sec. 3.1

## 2 Sharing the MRCA

We consider a  $\Xi$ - or  $\Lambda$ - $n$ -coalescent with a starting partition  $\pi = \{\{1\}, \dots, \{n\}\}$ , i.e. initially all the blocks  $\pi_i \in \pi$  are singleton blocks. We refer to the elements of the starting partition as ‘leaves’. A common ancestor of a set  $A \subset \mathbb{N}$  of leaves is any block containing  $A$ . A set  $A$  of leaves has a common ancestor if and only if the coalescent passes through a partition with a block containing  $A$ . This allows us to identify the common ancestor with blocks of the partition-valued states of the coalescent. The MRCA of a set  $A$  of leaves is the smallest block which contains  $A$  (whenever that block appears). Given that we start from a finite set  $[n]$  of leaves ( $n < \infty$ ) we will eventually (i.e. in finite time almost surely) observe the partition  $\{[n]\}$  containing only the block  $[n]$ . Let  $T_{\text{MRCA}}^{(n)}$  denote the TMRC of the set  $[n]$  of leaves, i.e. we define

$$T_{\text{MRCA}}^{(n)} := \inf \left\{ t \geq 0 : \Pi_t^{(n)} = \{[n]\} \right\}. \quad (2)$$

The  $T_{\text{MRCA}}^{(n)}$  is therefore the first time  $\Pi_t^{(n)}$  arrives at partition  $\{[n]\} \in \mathcal{P}_n$ . Write  $T_{\text{MRCA}}^{(\infty)}$  for the TMRCA of the whole population. By a finite sample we mean a finite set  $A$  of leaves.

A subsample is a subset of a given sample (a given set of leaves). We let  $m$  denote the size of the subsample. For convenience and w.l.o.g. we assume leaves 1 to  $m$  are the leaves of the subsample, and we assume block  $\pi_1$  in any partition always contains element 1. A common ancestor of the subsample is any block containing  $[m]$ ; the MRCA of the subsample is the smallest block containing  $[m]$  (whenever it appears). We define the time to the MRCA of a subsample of size  $m$  of a sample of size  $n \geq m$  as

$$T_{\text{MRCA}}^{(m;n)} := \inf \left\{ t \geq 0 : [m] \subseteq \pi_1 \in \Pi_t^{(n)} \right\}; \quad (3)$$

i.e.  $T_{\text{MRCA}}^{(m;n)}$  is the time of first occurrence of the subset  $[m]$  in block  $\pi_1$  in a partition of  $\Pi^{(n)}$ . The sample and the subsample share the MRCA if the smallest block containing  $[m]$  ever observed in  $\Pi^{(n)}$  is  $[n]$ ; this happens almost surely if  $T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)}$ .

Our main mathematical results concern the probability

$$p_{n,m}^{(\Pi)} := \mathbb{P}^{(\Pi)} \left( T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \right), \quad (4)$$

which is the probability that the sample (of size  $n$ ) and the nested subsample (of size  $m < n$ ) share their MRCA under the coalescent process  $\Pi$ . From now on it should be understood that we always look at nested samples. We are able to obtain representations of  $p_{n,m}^{(\Pi)}$  both for finite  $n$  and  $m$  and also for the limit  $\lim_{n \rightarrow \infty} p_{n,m}^{(\Pi)}$ ,  $m$  fixed, for some MMC processes. We will let  $p_{n,m}^{(\Xi)}$  denote  $p_{n,m}^{(\Pi)}$  in (4) when  $\Pi$  is a  $\Xi$ -coalescent, and  $p_{n,m}^{(\Lambda)}$  denote  $p_{n,m}^{(\Pi)}$  when  $\Pi$  is a  $\Lambda$ -coalescent.

## 2.1 Finite $n$

Our main focus is to compare genealogical properties of nested samples between different coalescent processes in order to learn what is gained by enlarging the sample size. In this context, a natural question to address is which  $n$ -coalescent  $\Pi$  maximises  $p_{n,m}^{(\Pi)}$  for a given finite sample size  $n$  and subsample size  $m$ ? Trivially this is the ‘star-shaped’ coalescent with  $\Lambda$ -measure  $\Lambda(dx) = \delta_1(x)dx$  so that  $\Lambda(\{1\}) = 1$ , all  $n$  blocks merge after an exponential waiting time, and  $p_{n,m}^{(\delta_1)} = 1$ . We now compare  $p_{n,m}^{(\text{Kingman})}$  (meaning  $p_{n,m}^{(\Pi)}$  when  $\Pi$  is the Kingman-coalescent) to all  $p_{n,m}^{(\Lambda)}$  with  $\Lambda(\{1\}) = 0$ . We can show the following (see Sec. 4.1 for a proof).

**Proposition 1** *For any given sample size  $n$  and subsample size  $m < n$  there is a  $\Lambda'$  with  $\Lambda'(\{1\}) = 0$  which fulfills  $p_{n,m}^{(\Lambda')} > p_{n,m}^{(\text{Kingman})}$ .*

One can think of  $\Lambda'$  as given by  $\Lambda = \delta_\psi$  for some fixed  $\psi \in (0, 1)$  and very close to 1. Prop. 1 holds for any finite sample size  $n$  and subsample size  $m$ . Regarding



the limit  $p_m^{(\Pi)} = \lim_{n \rightarrow \infty} p_{n,m}^{(\Pi)}$  with  $m$  fixed we conjecture that  $p_m^{(\text{Kingman})} > p_m^{(\Lambda)}$  for every  $\Lambda$ -coalescent with  $\Lambda(\{1\}) = 0$ . Should our conjecture be true, the limits compare in the opposite way to the comparison of the non-limit probabilities given in Prop. 1.

The result in Prop. 1 holds for a very special  $\Lambda$ -coalescent. One can numerically evaluate  $p_{n,m}^{(\Lambda)}$  with a recursion (see Sec. 4.7.1 for a proof), and thus compare  $p_{n,m}^{(\Lambda)}$  for different  $\Lambda$ -coalescents. Let  $\lambda(n)$  (see Eq. (A36)) denote the total rate of mergers given  $n$  blocks, and  $\lambda_k(n)$  (see Eq. (A35)) denote the rate at which any  $k$  of  $n$  blocks merge. Write  $\beta(n, n-k+1) := \lambda_k(n)/\lambda(n)$  for the probability of a single merger of  $k$  blocks (a  $k$ -merger) given  $n$  blocks ( $2 \leq k \leq n$ ). Then

$$p_{n,m}^{(\Lambda)} = \sum_{k=2}^n \beta(n, n-k+1) \sum_{\ell=0}^{k \wedge m} \frac{\binom{n-m}{k-\ell} \binom{m}{\ell}}{\binom{n}{k}} p_{n-k+1, m'}^{(\Lambda)} \quad (5)$$

where  $\binom{n-m}{k-\ell} := 0$  if  $n-m < k-\ell$  and  $m' = (m-\ell+1)\mathbb{1}_{(\ell>1)} + m\mathbb{1}_{(\ell \leq 1)}$ . In the case  $m=2$  recursion (5) simplifies to

$$p_{n,2}^{(\Lambda)} = \sum_{k=2}^{n-2} \beta(n, n-k+1) \frac{(n-k)(n+k-1)}{n(n-1)} p_{n-k+1,2}^{(\Lambda)} + \beta(n,2) \frac{2}{n} + \beta(n,1). \quad (6)$$

Recursion (5) further simplifies in the case of the Kingman coalescent, since then  $\beta(n, n-1) = 1$  for  $n \geq 2$ . [73] obtain, with  $\Pi$  the Kingman-coalescent,

$$p_{n,m}^{(\Pi)} = \frac{m-1}{m+1} \frac{n+1}{n-1}. \quad (7)$$

Since the representation (7) only depends on which mergers are possible, the result (7) holds for a time-changed Kingman-coalescent as derived for example in [53] from a population model of ‘modest’ changes in population size.

As remarked in the introduction, the Beta-coalescent (see Eq. (A37)) with coalescent parameter  $\alpha \in [1, 2)$ , is an example of a  $\Lambda$ -coalescent (see Eq. (A33)) and can be derived from population model (1). Figure 7 shows graphs of  $p_{n,m}^{(\Pi)}$  when  $\Pi$  is the Beta-coalescent (see Eq. (A37)) with coalescent parameter  $\alpha \in [1, 2)$  as a function of  $\alpha$ ; the results indicate that  $p_{n,m}^{(\text{Beta-coal})} < p_{n,m}^{(\text{Kingman})}$  for  $n$  large enough and any  $m$ . This shows that one needs a larger subsample under the Beta-coalescent than under the Kingman-coalescent for a given sample size to have the same value of  $p_{n,m}^{(\Pi)}$ . By implication, one gains more information by enlarging the sample under the Beta-coalescent than under the Kingman-coalescent.

We conclude this subsection with two closed-form representations of  $p_{n,m}^{(\Pi)}$ . To prepare for the first one we recall the concept of ‘coming down from infinity’. This property is defined as follows. If a  $\Xi$ - $n$ -coalescent  $(\Pi_t^{(\Xi)})_{t \geq 0}$  comes down from infinity then, with probability 1, the number of blocks is finite for any  $t > 0$ , which is equivalent to  $\lim_{n \rightarrow \infty} T_{\text{MRCA}}^{(n)} < \infty$  a.s. If  $\Pi_t^{(\Xi)}$ , for  $t > 0$ , has infinitely many blocks with probability 1, we say that the coalescent ‘stays infinite’. Conditions for  $\Xi$  to fall into one of these two classes are available, see e.g. [77, 76, 48]. If



$\Xi(\{\mathbf{x} \in \Delta \mid \sum_{i=1}^k x_i = 1 \text{ for } k \in \mathbb{N}\}) > 0$ , the  $\Xi$ -coalescent does not stay infinite [76, p.39], but does not necessarily come down from infinity. In fact, there is a.s. a finite (random) time  $T \geq 0$  so that the number of blocks is finite for all  $t > T$  (see [75, p. 39]). This means that for such a coalescent,  $\lim_{n \rightarrow \infty} T_{\text{MRCA}}^{(n)}$  is finite almost surely. For processes that stay infinite ( $\tilde{\Pi}$ ),  $\lim_{n \rightarrow \infty} p_{n,m}^{(\tilde{\Pi})} = 0$  since the MRCA of the set  $\mathbb{N}$  of leaves in the starting partition  $\{\{1\}, \{2\}, \dots\}$  is never reached.

We have a representation of  $p_{n,m}^{(\Xi)}$  (see Sec. 4.2 for a proof).

**Proposition 2** For any finite measure  $\Xi$  on  $\Delta$ , we have

$$p_{n,m}^{(\Xi)} = 1 - \mathbb{E} \left[ \sum_{i \in \mathbb{N}} \prod_{\ell=0}^{m-1} \frac{B_{[i]}^{(n)} - \ell}{n - \ell} \right] > 0, \quad (8)$$

where  $B_{[1]}^{(n)}, B_{[2]}^{(n)}, \dots$  are the sizes of the blocks of  $\Pi_{T_{\text{MRCA}}^{(n)}}^{(n)}$ , ordered by size from biggest to smallest where the sequence  $B_{[1]}^{(n)}, B_{[2]}^{(n)}, \dots$  is extended to an infinite sequence by taking  $B_{[i]}^{(n)} = 0$  for  $i > \#\Pi_{T_{\text{MRCA}}^{(n)}}^{(n)}$ . If the  $\Xi$ -coalescent comes down from infinity, we have

$$p_{n,m}^{(\Xi)} \rightarrow 1 - \mathbb{E} \left[ \sum_{i \in \mathbb{N}} P_{[i]}^m \right] = 1 - \mathbb{E} [X^{m-1}] = 1 - \frac{\mathbb{E} [Y^m]}{\mathbb{E} [Y]} > 0 \quad (9)$$

for fixed  $m$  and  $n \rightarrow \infty$ , where  $P_{[i]} := \lim_{n \rightarrow \infty} B_{[i]}^{(n)} / n$  is the (almost surely existing) asymptotic frequency of the  $i$ th biggest block of  $\Pi_{T_{\text{MRCA}}^{(\infty)}}^{(\infty)}$ ,  $X$  is the asymptotic frequency of a size-biased pick from the blocks of  $\Pi_{T_{\text{MRCA}}^{(\infty)}}^{(\infty)}$ , while  $Y$  is the asymptotic frequency of a block picked uniformly at random from  $\Pi_{T_{\text{MRCA}}^{(\infty)}}^{(\infty)}$ .

In the case of the Bolthausen-Sznitman (BS-coal)  $n$ -coalescent [20], which is a  $\Lambda$ - $n$ -coalescent with  $\Lambda(dx) = dx$  (see Eq. (A34)), i.e. the density associated with the uniform distribution on  $[0, 1]$ , we can give a characterisation of  $p_{n,m}^{(\Pi)}$  in terms of independent Bernoulli r.v.'s (see Sec. 4.3 for a proof).

**Proposition 3** Let  $B_1, \dots, B_{n-1}$  be independent Bernoulli random variables with  $\mathbb{P}(B_i = 1) = 1/i$ . Let  $\Pi$  denote the Bolthausen-Sznitman  $n$ -coalescent. For  $2 \leq m < n$ ,

$$p_{n,m}^{(\Pi)} = \mathbb{E} \left[ \frac{B_1 + \dots + B_{m-1}}{B_1 + \dots + B_{n-1}} \right]. \quad (10)$$

Moreover,  $\log n p_{n,m}^{(\Pi)} \rightarrow \sum_{i=1}^{m-1} i^{-1}$  for  $n \rightarrow \infty$  and  $m$  fixed.

## 2.2 Two variants of $p_{n,m}^{(\Pi)}$

The probability  $p_{n,m}^{(\Pi)}$  (see Eq. (5)) is an indication of how likely it is that the ‘oldest’ genealogical branches, or the edges connected directly to the MRCA, are (partially) shared between subsample and sample. We remark that the sample and the subsample may share the MRCA without sharing any of the internal edges if the associated coalescent admits multiple mergers (see Fig. 1C for an example). Such events are highly unlikely though for  $n$  large enough if the  $\Lambda$ -coalescent comes down from infinity, see Corollary 1. If the sample and the subsample share the MRCA then the subsample is more likely to include the oldest allele of the complete sample. To derive the actual probability of the event that the subsample carries the oldest allele of the complete sample one needs to include mutation. Consider a  $\Lambda$ - $n$ -coalescent with neutral mutation. Mutations are modelled by a homogeneous Poisson point process on the branches of the  $\Lambda$ - $n$ -coalescent with (scaled) mutation rate  $\theta > 0$ . We assume the infinitely-many-alleles model. This means that the allelic type of each individual is seen by tracing its ancestral line back to the first mutation on it. The ancestral line shares the type of the MRCA if there is no mutation on the line before the MRCA is reached. We are interested in the event that the oldest allele from the sample is also found in the subsample. This can be expressed by using the concept of ‘frozen’ and ‘active’ ancestral lines in a  $n$ -coalescent with mutation [23]. At a given time  $t$ , an ancestral lineage is called frozen if there has been a mutation on it, otherwise it is called active. The age of a sampled allele ( $i$ , say) is the waiting time  $\tau_i$  until its’ ancestral lineage is frozen. For consistency we prolong the  $n$ -coalescent after reaching the MRCA (at time  $T_{\text{MRCA}}^{(n)}$ ) by a single ancestral line. The first mutation on the prolonged line is seen after an additional  $\text{Exp}(\theta/2)$  time which freezes the line. Thus, the oldest allele of a sample is given by the ancestral lineage which is frozen last (active the longest), and this age is  $\max\{\tau_i : i \in [n]\}$  for the sample and  $\max\{\tau_i : i \in [m]\}$  for the subsample. Let  $A^{(n)}(t)$  denote the count of active ancestral lineages in the sample at time  $t$ . We write

$$p_{n,m}^{(\Pi,\theta)} := \mathbb{P}^{(\Pi,\theta)} \left( A^{(n)}(\max\{\tau_i : i \in [m]\}) = 0 \right). \quad (11)$$

for the probability that the subsample includes the oldest allele of the sample.

We consider  $p_{n,m}^{(\Lambda,\theta)} \equiv p_{n,m}^{(\Pi^{(\Lambda)},\theta)}$  for  $n, m \in \mathbb{N}_0$ ,  $\theta > 0$ . The case  $n = m = 1$  (or  $n > m = 1$ ) means we trace back a single lineage until it is hit by a mutation (either in the sample and/or subsample). The boundary conditions are  $p_{n,n}^{(\Lambda,\theta)} = 1$  and  $p_{n,0}^{(\Lambda,\theta)} = 0$  for  $n > 0$ . The recursion for  $p_{n,m}^{(\Lambda,\theta)}$  is

$$\begin{aligned} p_{n,m}^{(\Lambda,\theta)} &= \frac{\theta m}{2\lambda(n) + \theta n} p_{n-1,m-1}^{(\Lambda,\theta)} + \frac{\theta(n-m)}{2\lambda(n) + \theta n} p_{n-1,m}^{(\Lambda,\theta)} \\ &\quad + \frac{2\lambda(n)}{2\lambda(n) + \theta n} \sum_{k=2}^n \beta(n, n-k+1) \sum_{\ell=0}^{k \wedge m} \frac{\binom{n-m}{k-\ell} \binom{m}{\ell}}{\binom{n}{k}} p_{n-k+1,m'}^{(\Lambda,\theta)}, \end{aligned} \quad (12)$$

where  $m' = (m - \ell + 1)\mathbb{1}_{(\ell > 1)} + m\mathbb{1}_{(\ell \leq 1)}$  (see Sec. 4.7.3 for a proof),

The probability  $p_{n,m}^{(\Lambda,\theta)}$  is a function of the scaled mutation rate  $\theta$ . Here, and in most models in population genetics which include mutation,  $\theta := \mu_N/c_N$  where  $\mu_N$  is the rate of mutation per locus per generation, and  $c_N$  is the pairwise coalescence probability, or the probability that 2 distinct individuals sampled at the same time from a population of size  $N$  have the same parent. Since (usually)  $c_N \rightarrow 0$  as  $N \rightarrow \infty$  to ensure convergence to a continuous-time limit [70, 63], and since  $\theta$  is usually assumed to be of order  $\mathcal{O}(1)$ , we let  $\mu_N$  depend on  $N$ . The key point here is that  $\theta$  depends on  $c_N$ . By way of an example,  $c_N = 1/N$  for the haploid Wright-Fisher model, while  $c_N = \mathcal{O}(N^{1-\alpha})$  for the Beta( $2 - \alpha, \alpha$ )-coalescent,  $1 < \alpha < 2$  [71]. This means that the scaled mutation rates ( $\theta$ ) are not directly comparable between different coalescent processes; this again means that expressions  $(p_{n,m}^{(\Lambda,\theta)})$ , defined in (11), for example) which depend on the mutation rate cannot be directly compared between different coalescent processes that may have different timescales. We further remark that we must define  $\theta$  to be proportional to  $c_N$  since the branch lengths on which the mutation process runs are in units of  $1/c_N$ ; thus if we don't rescale the mutation rate  $\mu_N$  with  $1/c_N$  we would never see any mutations. It is therefore the mutation rate  $\mu_N$ , which must be determined from molecular (or DNA sequence) data, which sets the timescale; the  $c_N$  comes from the model.

The probability  $p_{n,m}^{(\Pi)}$  is also the probability that the MRCA of the subsample (of size  $m$ ) subtends all the  $n$  leaves ( $[n]$  is the smallest block containing  $[m]$ ). A related more general question is to ask about the distribution of the size (number of elements) of the smallest block which contains  $[m]$ . This is the same as asking about the distribution of the number of leaves subtended by the MRCA of the subsample. For Kingman's  $n$ -coalescent, the distribution is computed in [81, Thm. 1]. The probability of the event that the MRCA of the subsample subtends only the leaves of the subsample is especially interesting, see also e.g. [87, p. 184, Eq. 2], where this probability is also described recursively. This recursion can be easily extended to  $\Lambda$ -coalescents. Define  $T_{\text{MRCA}}^{(\Lambda)}$  to be the first time that  $A$  is completely contained within a block of  $\Pi_t$ . Write

$$q_{n,m}^{(\Lambda)} := \mathbb{P}^{(\Pi^{(\Lambda)})} \left( T_{\text{MRCA}}^{([m] \cup \{i\})} > T_{\text{MRCA}}^{([m])} \forall i \in \{m+1, \dots, n\} \right) \quad (13)$$

for the probability that the MRCA of the subsample subtends only the leaves of the subsample. Let  $\beta(n, n-k+1)$  be the probability of a  $k$ -merger ( $2 \leq k \leq n$ ) given  $n$  active lines. The recursion for  $q_{n,m}^{(\Lambda)}$  is (see Sec. 4.7.2 for a proof)

$$q_{n,m}^{(\Lambda)} = \sum_{k=2}^{(n-m) \wedge m} \frac{\beta(n, n-k+1)}{\binom{n}{k}} \left( \binom{m}{k} q_{n-k+1, m-k+1}^{(\Lambda)} + \binom{n-m}{k} q_{n-k+1, m}^{(\Lambda)} \right) \quad (14)$$

with boundary conditions  $q_{n,n}^{(\Lambda)} = q_{n,1}^{(\Lambda)} = 1$  for  $n \in \mathbb{N}$ . One may use  $q_{n,m}^{(\Lambda)}$  to calculate a  $p$ -value in a test for observing block  $[m]$  under a  $\Lambda$ -coalescent. As one might expect (see Sec. 4.3 for a proof), for  $m$  fixed and for any  $\Lambda$ -coalescent,

$$\lim_{n \rightarrow \infty} q_{n,m}^{(\Lambda)} = 0. \quad (15)$$

In the case of the Bolthausen-Sznitman (BS-coal)  $n$ -coalescent we obtain an exact representation of  $q_{n,m}^{(\text{BS-coal})}$  (see Sec. 4.4 for a proof).

**Proposition 4** Let  $B_1, \dots, B_{n-1}, B'_1, \dots, B'_{n-m}$  be independent Bernoulli variables with  $\mathbb{P}(B_i = 1) = \mathbb{P}(B'_i = 1) = i^{-1}$ . For the Bolthausen-Sznitman  $n$ -coalescent we have, for  $2 \leq m < n$ ,

$$q_{n,m}^{(BS-coal)} = \binom{n-1}{m-1}^{-1} \mathbb{E} \left[ \left( \frac{\sum_{i \in [m-1]} B_i + \sum_{i \in [n-m]} B'_i}{\sum_{i \in [m-1]} B_i} \right)^{-1} \right]. \quad (16)$$

### 2.3 The limit $\lim_{n \rightarrow \infty} p_{n,m}^{(\Pi)}$

In this subsection we discuss the limit  $\lim_{n \rightarrow \infty} p_{n,m}^{(\Pi)}$  with  $m$  fixed. For a fixed  $m \in \mathbb{N}$ , write

$$p_m^{(\Pi)} := \lim_{n \rightarrow \infty} \mathbb{P}^{(\Pi)} \left( T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \right) \quad (17)$$

for the probability, under coalescent  $\Pi$ , that a subsample of size  $m$  shares the MRCA with an arbitrarily large sample. The limit  $\lim_{n \rightarrow \infty} T_{\text{MRCA}}^{(n)}$  is a valid limit for any coalescent (even if it diverges) and therefore (17) is well defined. For any  $\Xi$ -coalescent  $(p_{n,m}^{(\Xi)})_{n > m}$  is monotonically decreasing as  $n$  increases. The limit  $p_m^{(\Xi)}$  is derived under the assumption that the same  $\Xi$ -coalescent is obtained for arbitrarily large sample size. This assumption may not hold when one wants to relate to finite real populations. The quantity  $p_m^{(\Pi)}$  should therefore only be regarded as a limit. See further discussion on this point in Sec. 5.

For the Kingman-coalescent we have the following result, first obtained in [73] by solving a recursion,

$$p_m^{(\text{Kingman})} = \frac{m-1}{m+1}. \quad (18)$$

To see (18) without solving a recursion, we consider the process forwards in time from the MRCA. Label the two ancestral lines generated by the first split (of the MRCA) as  $a_1$  and  $a_2$ . The fraction of the population that is a descendant of  $a_1$  is distributed as a uniform random variable on the unit interval, see e.g. the remark after Thm. 1.2 in [8]. Therefore, with  $U$  a uniform r.v. on  $[0, 1]$ , and any finite  $m \in \mathbb{N}$ ,

$$p_m^{(\text{Kingman})} = 1 - \mathbb{E}[U^m] - \mathbb{E}[(1-U)^m] = 1 - 2 \int_0^1 x^m dx = \frac{m-1}{m+1}. \quad (19)$$

For the Bolthausen-Sznitman coalescent (BS-coal)  $\lim_{n \rightarrow \infty} p_{n,m}^{(\text{BS-coal})} = 0$  for  $m$  fixed. We remark in this context that the Bolthausen-Sznitman coalescent does not come down from infinity.

Result (18) is also an indication that the statistic  $T_{\text{MRCA}}^{(n)}$  is a good statistic for capturing a property of the population with a small sample, at least under the Kingman coalescent. We remark that the Kingman coalescent comes down from infinity. Result (18) (and (19)) is the ‘spark’ for the current work.

Our main mathematical result, Thm. 1, is a representation of  $p_m^{(\text{Beta-coal})}$ , i.e.  $p_m^{(\Pi^{(\Lambda)})}$  when  $\Pi^{(\Lambda)}$  is the Beta( $2 - \alpha, \alpha$ )-coalescent [78] (Beta-coalescent; see Eq.

(A38)). The Beta-coalescent is well-studied, there are connections to superprocesses, continuous-state branching processes (CSBP) and continuous stable random trees as described e.g. in [14] and [7]. For  $\alpha \in (1, 2)$ , the Beta-coalescent comes down from infinity. The representation of  $p_m^{(\text{Beta-coal})}$  given in Thm. 1 can be directly derived from [8, Thm. 1.2], which is a result based on the connection between the Beta-coalescent and a CSBP (see Sec. 4.5 for a proof).

**Theorem 1** Define  $p_m^{(\text{Beta-coal})} \equiv p_m^{(\Pi)}$  (see Eq. (17)) when  $\Pi$  is the Beta-coalescent for  $\alpha \in (1, 2)$ . Let  $K$  denote the random number of blocks involved in the merger upon which the MRCA of  $[n]$  is reached;  $K$  has generating function  $\mathbb{E}[u^K] = \alpha u \int_0^1 (1-x)^{1-\alpha} ((1-ux)^{\alpha-1} - 1) dx$  for  $u \in [0, 1]$  [47, Thm. 3.5]. Let  $(Y_i)_{i \in \mathbb{N}}$  be a sequence of i.i.d. r.v. with Slack's distribution on  $[0, \infty)$ , i.e.  $Y_1$  has Laplace transform  $\mathbb{E}[e^{-\lambda Y_1}] = 1 - (1 + \lambda)^{-1/(\alpha-1)}$  [80]. We have the representation

$$p_m^{(\text{Beta-coal})} = 1 - \sum_{k \in \mathbb{N}} k \mathbb{E}[(Y_1 + \dots + Y_k)^{1-\alpha}]^{-1} \mathbb{E}\left[\frac{Y_1^m}{(Y_1 + \dots + Y_k)^{\alpha+m-1}}\right] \mathbb{P}(K = k). \quad (20)$$

A comparison of  $p_m^{(\Pi)}$  between coalescent processes that come down from infinity is complicated by at least two things. First, it is not clear what such a comparison would mean in terms of inference for real populations. Second, the representation (20) is highly non-trivial to evaluate. However, the result in Thm. 1 is of mathematical interest in its own right.

We close this subsection with a consideration of the limit  $\lim_{n \rightarrow \infty} p_{n,m}^{(\Pi)}$  when  $\Pi$  is a  $\Xi$ -coalescent (A32). We give a criterion for when  $p_m^{(\Xi)} > 0$ . This question is closely related to the question of coming down from infinity for a  $\Xi$ -coalescent. We have the following result (see Sec. 4.6 for a proof).

**Proposition 5** Consider any  $\Xi$ -coalescent. For any fixed  $m \in \mathbb{N}$ ,  $m \geq 2$ ,  $p_m^{(\Xi)}$  exists. If the coalescent comes down from infinity or  $\Xi(\{x \in \Delta \mid \sum_{i=1}^k x_i = 1 \text{ for } k \in \mathbb{N}\}) > 0$ ,  $p_m^{(\Xi)} > 0$ . If it stays infinite,  $p_m^{(\Xi)} = 0$ .

### 3 Relative times and lengths

So far we have considered how well the MRCAs match between a sample and a nested subsample when both are finite; we have also discussed some limit results when the sample size tends to infinity. We would also like to have some understanding of how the distributions of the various genealogical statistics compare between different coalescent processes. In this section we use simulations to estimate the distribution of two genealogical statistics and compare them between the Kingman-coalescent, exponential population growth, and the Beta-coalescent.

#### 3.1 Simulation method

To generate realisations of our statistics we simulate genealogies by drawing waiting times between mergers, and merger sizes, governed by the transition rates of the corresponding  $n$ -coalescent.

Denote by  $T_j$  the random waiting time for the first merger of the  $j$ -coalescent. The coalescent process under exponential population growth is a time-changed Kingman-coalescent (see e.g. [25,30]). [41] give a way of sampling  $T_j$  under exponential growth. Let  $\beta > 0$  denote the growth rate under exponential growth. Write  $S_j = T_n + \dots + T_j$  for  $2 \leq j \leq n$ , with  $S_{n+1} = 0$  a.s. If  $\{U_j : 2 \leq j \leq n\}$  denotes a collection of i.i.d. uniform  $(0, 1]$  random variables, then [41]

$$S_j = T_j + S_{j+1} = \frac{1}{\beta} \log \left( \exp(\beta S_{j+1}) - \frac{2\beta}{j(j-1)} \log(U_j) \right), \quad 2 \leq j \leq n. \quad (21)$$

Eq. (21) tells us that if  $\beta$  is very large, the time intervals  $T_j$  near the MRCA become quite small. The time intervals near the leaves are much less affected. We choose the grid of values for  $\beta$  as

$$\beta \in \{0.1, 0.5, 1, 10, 50, 100, 500, 1000, 5000, 10000\}.$$

Recall in this context the growth model  $N_k = N_0(1 + \beta/N)^k$  for the population size in generation  $k \geq 0$  going forward in time, and where  $N_0$  is the population size at the start of the growth. Our choice of grid values for  $\beta$  should reflect the range of growth from weak ( $\beta = 0.1$ ) to very strong ( $\beta = 10^4$ ) and most estimates of  $\beta$  obtained for natural populations should fall within this range.

Under the Beta-coalescent  $T_j$  is an exponential with rate  $\lambda(j) = \lambda_2(j) + \dots + \lambda_j(j)$  where  $\lambda_i(j)$  is given in Eq. (A38).

Let  $n$  and  $m$  denote the current number of sample and subsample lines; since the subsample is nested within the sample a subsample line is necessarily also a sample line. Let  $M \in \{2, \dots, n\}$  be the size of the first merger of the corresponding  $n$ -coalescent. Under a  $\Lambda$ - $n$ -coalescent  $M = k$  with probability  $\mathbb{P}(M = k) = \lambda_k(n)/\lambda(n)$  (see Eq. (A35) and (A36)) for  $2 \leq k \leq n$ ; under a (time-changed) Kingman- $n$ -coalescent  $M = 2$  a.s. We partition the sample and subsample lines as follows. Let  $m_{\text{ext}}$  resp.  $\tilde{m}_{\text{ext}}$  denote the current number of ‘external’ lines of the subsample, resp. number of ‘external’ lines not belonging to the subsample. Let  $m_{\text{int}}$  resp.  $\tilde{m}_{\text{int}}$  denote the current number of ‘internal’ lines of the subsample, resp. current number of ‘internal’ lines not belonging to the subsample. External lines are subtended by exactly 1 leaf, while internal lines are subtended by at least 2 leaves.

The distinction between the lines will now be further explained. Given that we start with  $n$  sample lines of which  $m$  belong to the subsample then initially  $m_{\text{ext}} = m$ ,  $\tilde{m}_{\text{ext}} = n - m$ , and  $m_{\text{int}} = \tilde{m}_{\text{int}} = 0$ . Now suppose we have drawn  $k$  lines to merge, of which  $x_1 \leq m_{\text{ext}}$  were drawn from the external lines of the subsample,  $x_2 \leq \tilde{m}_{\text{ext}}$  from the external lines not belonging to the subsample,  $x_3 \leq m_{\text{int}}$  from the internal lines of the subsample, and  $x_4 \leq \tilde{m}_{\text{int}}$  from the internal lines not belonging to the subsample. The following transitions of the lines then occur:

$$\begin{aligned} m_{\text{ext}} &\rightarrow m_{\text{ext}} - x_1, \\ \tilde{m}_{\text{ext}} &\rightarrow \tilde{m}_{\text{ext}} - x_2, \\ m_{\text{int}} &\rightarrow m_{\text{int}} - x_3 + \mathbb{1}_{(x_1+x_3 \geq 1)}, \\ \tilde{m}_{\text{int}} &\rightarrow \tilde{m}_{\text{int}} - x_4 + \mathbb{1}_{(x_2+x_4=k)}. \end{aligned} \quad (22)$$

The transitions reflect our assumption that if at least 1 subsample line is involved in a given merger, the continuing ancestral line is considered to belong to the

subsample; mutations that arise on the continuing line will then be carried by the subsample, and visible in the subsample unless all the subsample lines were involved in the merger ( $x_1 = m_{\text{ext}}$  and  $x_3 = m_{\text{int}}$  and  $x_1 + x_3 \geq 1$ ). Therefore, if a single external subsample line, and no other subsample line, is involved in a merger ( $x_1 = 1, x_3 = 0$ ) we regard the continuing line as an ‘internal’ line of the subsample. An external line of the subsample therefore remains so only until it is involved in a merger. By way of example, the continuing line of the first merger in Fig. 1D counts as an internal line of the subsample.

Given the merger size  $k \in [2, m_{\text{ext}} + \tilde{m}_{\text{ext}} + m_{\text{int}} + \tilde{m}_{\text{int}}]$ , due to the exchangeability of the  $n$ -coalescent, we draw  $k$  lines from a multivariate hypergeometric by

$$\mathbb{P}(X = x) = \frac{\binom{m_{\text{ext}}}{x_1} \binom{\tilde{m}_{\text{ext}}}{x_2} \binom{m_{\text{int}}}{x_3} \binom{\tilde{m}_{\text{int}}}{x_4}}{\binom{m_{\text{ext}} + \tilde{m}_{\text{ext}} + m_{\text{int}} + \tilde{m}_{\text{int}}}{k}}, \quad x_1 + \dots + x_4 = k, \quad (23)$$

where the  $x_i$  denote the number of lines drawn from each of the four groups.

The statistics we consider are the relative times  $\rho_T^{(m;n)} := T_{\text{MRCA}}^{(m;n)} / T_{\text{MRCA}}^{(n)}$  and the relative lengths  $\rho_I^{(m;n)} := L_{\text{int}}^{(m;n)} / L_{\text{int}}^{(n)}$  where  $L_{\text{int}}^{(m;n)}$  is the sum of the lengths of the internal edges associated with the subsample and  $L_{\text{int}}^{(n)}$  is the sum of the lengths of internal edges of the complete sample. The ratio  $\rho_I^{(m;n)}$  indicates how much of the ‘ancestral variation’, or mutations present in at least 2 copies in the sample, are captured by the subsample. The ratio  $\rho_T^{(m;n)}$  indicates how likely we are to capture with the subsample the ancestral variation in the complete sample.

A realisation of  $\rho_I^{(m;n)}$  is obtained as follows. Given  $j = m_{\text{ext}} + m_{\text{int}} + \tilde{m}_{\text{ext}} + \tilde{m}_{\text{int}}$  current sample lines, let  $t_j$  denote a realisation of  $T_j$ , the random time during which there are  $j$  lines of the complete sample. We update the total lengths  $\ell_{\text{int}}^{(m;n)}$  of internal subsample lines, and  $\ell_{\text{int}}^{(n)}$  of internal lines of the complete sample, as

$$\begin{aligned} \ell_{\text{int}}^{(m;n)} &\rightarrow \ell_{\text{int}}^{(m;n)} + \mathbb{1}_{(m_{\text{ext}} + m_{\text{int}} > 1)} m_{\text{int}} t_j, \\ \ell_{\text{int}}^{(n)} &\rightarrow \ell_{\text{int}}^{(n)} + \mathbb{1}_{(j > 1)} (m_{\text{int}} + \tilde{m}_{\text{int}}) t_j. \end{aligned} \quad (24)$$

The updating rule for  $\ell_{\text{int}}^{(m;n)}$  in Eq. (24) reflects the fact that mutations on the common ancestor line of the subsample, for example the continuing line after the merger of all 3 subsample lines in Fig. 1D, are not visible in the subsample. The updating rule for  $\ell_{\text{int}}^{(n)}$  in Eq. (24) similarly reflects the fact that mutations on the continuing line of the MRCA of the complete sample are not visible in the sample; but once the MRCA of the complete sample is reached we stop the process.

A realisation of  $\rho_I^{(m;n)}$  is then recorded as  $r_I^{(m;n)} := \ell_{\text{int}}^{(m;n)} / \ell_{\text{int}}^{(n)}$ . By way of example, the edges marked with a black dot in Fig. 1B are internal edges of the complete sample while the edges marked with a circle in Fig. 1A are internal edges associated with the subsample as well as the complete sample, and we have  $\rho_I^{(m;n)} = (T_5 + T_4 + T_3) / (T_5 + 2T_4 + T_3)$  for the genealogy in Fig. 1A. There are no internal edges associated with the subsample in Fig. 1B and 1C; therefore  $\rho_I^{(m;n)} = 0$  for the genealogies in Fig. 1B and 1C. The sample and the subsam-



ple share all the internal edges in the genealogy shown in Fig. 1D and therefore  
 $\rho_I^{(m;n)} = 1$ .

Realisations of  $T_{\text{MRCA}}^{(m;n)}(t^{(m;n)})$  and  $T_{\text{MRCA}}^{(n)}(t^{(n)})$  are recorded as

$$\begin{aligned} t^{(m;n)} &= \inf\{t \geq 0 : m_{\text{ext}} + m_{\text{int}} = 1\}, \\ t^{(n)} &= \inf\{t \geq 0 : m_{\text{ext}} + m_{\text{int}} + \tilde{m}_{\text{ext}} + \tilde{m}_{\text{int}} = 1\}, \end{aligned} \quad (25)$$

by adding up the realised waiting times  $t_j$  of  $T_j$ . We record a realisation of  $\rho_T^{(m;n)}$   
as  $r_T^{(m;n)} := t^{(m;n)} / t^{(n)}$ .

### 3.2 Simulation results

Figures 2 and 3 show estimates, in the form of boxplots (see a description in the  
captions), of the distributions of  $\rho_T^{(m;n)}$  (left column) and  $\rho_I^{(m;n)}$  (right column);  
under the Beta-coalescent as a function of  $\alpha$  (Figure 2) and under exponential  
growth as a function of  $\beta$  (Figure 3). In some of the boxplots of  $\rho_T^{(m;n)}$  in fact all  
of them under exponential growth (Fig. 3) the interquartile range (the difference  
between the 75th and 25th percentile) is zero. The estimates shown in Fig. 3 of  
the distribution of  $\rho_T^{(m;n)}$  indicate that  $\rho_T^{(m;n)}$  becomes more concentrated at 1 as  $\beta$   
increases. Recall in this context that  $p_{n,m}^{(\text{exp. growth})} = (m-1)(n+1)/((m+1)(n-1))$  since exp. growth results in a time-changed Kingman-coalescent.

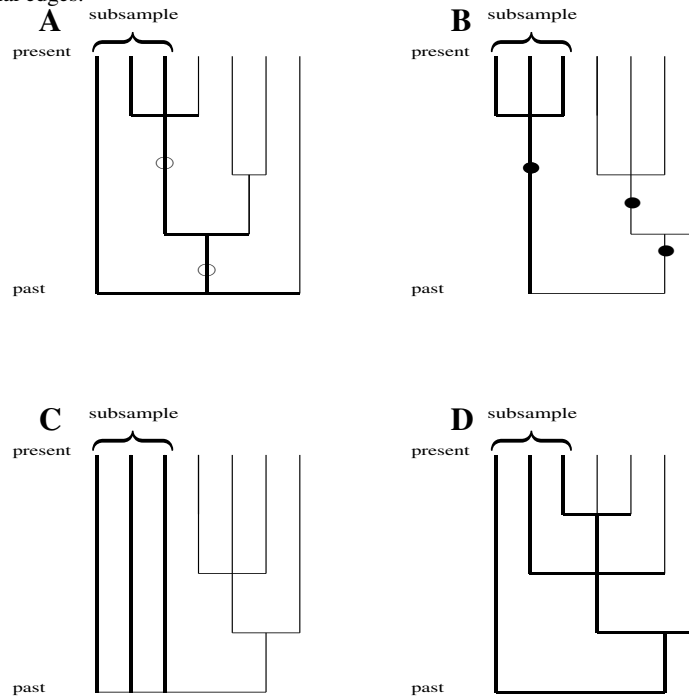
In Figure 2 we see a gradual shift in the distribution of  $\rho_T^{(m;n)}$  as subsample  
size increases; from being skewed to the right (ie. towards higher values) to being  
skewed to the left (ie. towards smaller values). This is in sharp contrast to the  
distribution under exponential growth (Figure 3) where the distribution of  $\rho_T^{(m;n)}$   
is always skewed to the left. This indicates that under a MMC process a subsample  
is much less informative about the larger complete sample than under exponential  
growth. In contrast, under exponential growth, even a small subsample can be very  
informative about the larger complete sample, especially in a strongly growing  
(large  $\beta$ ) population. Estimates of the means  $\mathbb{E}^{(\Pi)}[\rho_T^{(m;n)}]$ , shown in Figure 4  
(circles) for the Beta-coalescent, and in Figure 5 (circles) for exponential growth,  
further strengthen our conclusion.

The distribution of  $\rho_I^{(m;n)}$ , the relative lengths of internal edges, also behaves  
differently between the Beta-coalescent and exponential growth (Figures 2 and 3,  
middle columns). The distribution of  $\rho_I^{(m;n)}$  becomes more concentrated around  
smaller values as growth becomes stronger ( $\beta$  increases) while it stays highly  
variable as  $\alpha$  tends to 1, although the median decreases as skewness increases ( $\alpha$   
tends to 1). This indicates that we capture less and less of the ‘ancestral variation’  
(mutations observed in at least 2 copies in the sample) in the larger sample as  
growth or skewness increase. By implication our sample captures less of the an-  
cestral variation in a strongly growing population, or in a population with highly  
skewed reproduction. Estimates of  $\mathbb{E}^{(\Pi)}[\rho_I^{(m;n)}]$  (Figures 4 and 5, ‘+’) also indi-

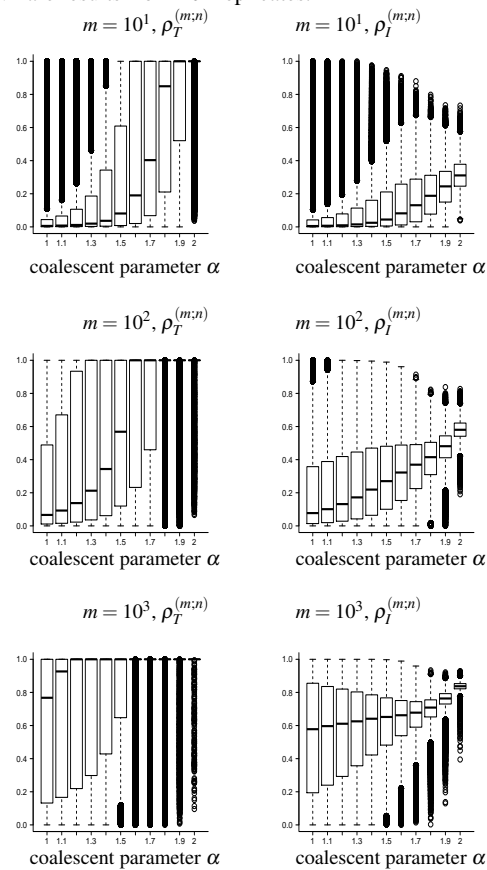
cate that one would need a large sample to capture at least half of the ancestral variation if growth or skewness is high.

To conclude,  $\rho_T^{(m;n)}$  and  $\rho_I^{(m;n)}$  seem to converge to opposite values under exponential growth;  $\rho_T^{(m;n)}$  to 1 and  $\rho_I^{(m;n)}$  to small values, as  $\beta$  increases. Thus, even if we are sharing the MRCA with higher probability as  $\beta$  increases (recall that the samples are nested), we are capturing less and less of the ancestral variation. Essentially the opposite trend is seen for both  $\rho_T^{(m;n)}$  and  $\rho_I^{(m;n)}$  under the Beta-coalescent; the distributions of both statistics stay highly variable as  $\alpha \rightarrow 1$ .

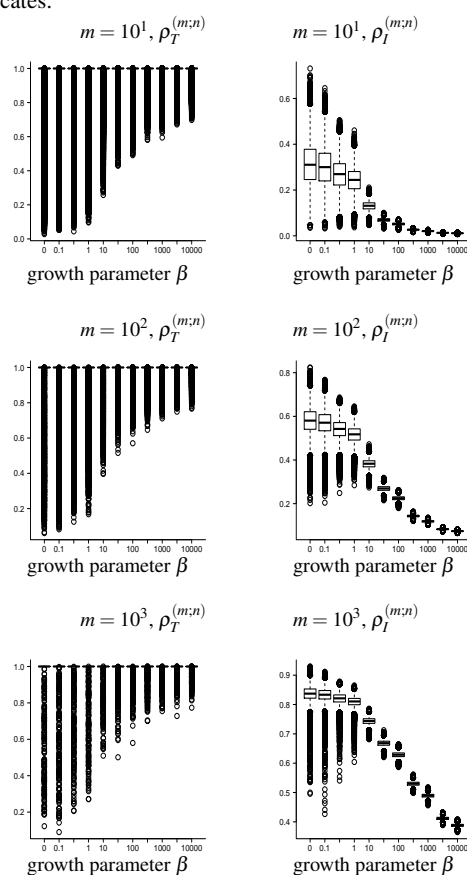
**Fig. 1** Examples of genealogies. Thick edges denote lineages ancestral to the subsample of size  $m = 3$ ; sample size  $n = 7$ . The marked edges in **A** denote internal ancestral lineages to both the subsample and the whole sample; the marked edges in **B** denote lineages internal only to the whole sample. In **C** the sample and subsample share the MRCA without sharing any of the internal edges. The genealogies are shown from the time of sampling (present) until the MRCA of the whole sample is reached (past). In **C** the sample and the subsample share the MRCA without sharing any internal edges. In **D** the sample and the subsample share the MRCA and all the internal edges.



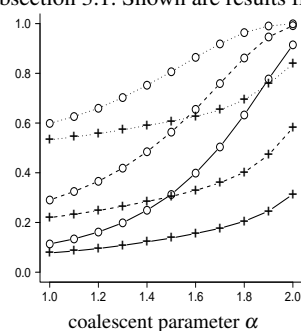
**Fig. 2** Estimates, shown in the form of boxplots, of the distributions of  $\rho_T^{(m;n)}$  and of  $\rho_I^{(m;n)}$  as functions of the coalescent parameter  $\alpha$  of the Beta( $2 - \alpha, \alpha$ )-coalescent for values of sample size  $n = 10^4$  and subsample size  $m$  as shown. The coalescent process at  $\alpha = 2$  is the Kingman-coalescent. Each box shows the interquartile range ( $IQR$ ; the difference between the 75th and 25th percentile), and the whiskers extend to  $\pm 1.5 \times IQR$ ; values outside  $\pm 1.5 \times IQR$  are shown as circles. The thick line within each box shows the median. For explanation of symbols see Subsection 3.1. Shown are results from  $10^5$  replicates.



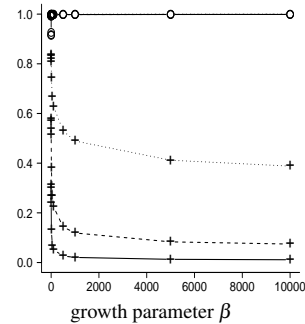
**Fig. 3** Estimates, shown in the form of boxplots, of the distributions of  $\rho_T^{(m;n)}$  and of  $\rho_I^{(m;n)}$  as functions of the exponential growth parameter  $\beta$  for values of sample size  $n = 10^4$  and subsample size  $m$  as shown. Each box shows the interquartile range ( $IQR$ ; the difference between the 75th and 25th percentile), and the whiskers extend to  $\pm 1.5 \times IQR$ ; values outside  $\pm 1.5 \times IQR$  are shown as circles. The thick line within each box shows the median. The coalescent process at  $\beta = 0$  is the Kingman-coalescent. For explanation of symbols see Subsection 3.1. The grid of values of  $\beta$  is  $\{0.1, 0.5, 1.0, 10.0, 50.0, 100.0, 500.0, 1000.0, 5000.0, 10000.0\}$ . Shown are results from  $10^5$  replicates.



**Fig. 4** Estimates of  $\mathbb{E}[\rho_T^{(m;n)}]$  (○), and of  $\mathbb{E}[\rho_I^{(m;n)}]$  (+) as functions of the coalescent parameter  $\alpha$  for values of sample size  $n = 10^4$  and subsample size  $m = 10^1$  (solid lines);  $m = 10^2$  (dashed lines);  $m = 10^3$  (dotted lines). The coalescent process at  $\alpha = 2$  is the Kingman-coalescent. For explanation of symbols see Subsection 3.1. Shown are results from  $10^5$  replicates.



**Fig. 5** Estimates of  $\mathbb{E}[\rho_T^{(m;n)}]$  ( $\circ$ ) and of  $\mathbb{E}[\rho_I^{(m;n)}]$  ( $+$ ) as functions of the exponential growth parameter  $\beta$  for values of sample size  $n = 10^4$  and subsample size  $m = 10^1$  (solid lines);  $m = 10^2$  (dashed lines);  $m = 10^3$  (dotted lines). The coalescent process at  $\beta = 0$  is the Kingman-coalescent. For explanation of symbols see Subsection 3.1. The grid of values of  $\beta$  is  $\{0.1, 0.5, 1.0, 10.0, 50.0, 100.0, 500.0, 1000.0, 5000.0, 10000.0\}$ . Shown are results from  $10^5$  replicates.





## 4 Proofs

### 4.1 Proof of Prop. 1

*Proof* Let  $\Lambda = \delta_p$  for  $p \in (0, 1)$  which fulfills  $\Lambda(\{0\}) = 0$ . The probability that the associated  $\Lambda$ - $n$ -coalescent is star-shaped, i.e. all blocks merge at the first (and then only) collision, is

$$\frac{p^{n-2}}{p^{-2}(1 - (1-p)^n - np(1-p)^{n-1})} = \frac{p^n}{\sum_{i=2}^n \binom{n}{i} p^i (1-p)^{n-i}} > p^n.$$

For any star-shaped path of a  $n$ -coalescent, we have  $T_{\text{MRCA}}^{(n)} = T_{\text{MRCA}}^{(m)}$  for any  $m < n$ . Thus, we can choose  $\Lambda'$  s.t.  $\Lambda' = \delta_p$  with

$$p = \left( \mathbb{P}^{(\delta_0)} \left( T_{\text{MRCA}}^{(m)} = T_{\text{MRCA}}^{(n)} \right) \right)^{\frac{1}{n}}.$$

□

### 4.2 Proof of Prop. 2

*Proof* Assume  $\Pi$  is a  $\Xi$ -coalescent. The event  $\left\{ T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \right\}$  is the complement of the event

$$A_{m,n} := \left\{ [m] \subseteq B, B \text{ is a block of } \Pi_{T_{\text{MRCA}}^{(n)} -}^{(n)} \right\}. \quad (26)$$

Due to the exchangeability of the  $\Xi$ -coalescent,  $\Pi_{T_{\text{MRCA}}^{(n)} -}^{(n)}$  is an exchangeable partition of  $[n]$ . Given the (ordered) block sizes  $\left( B_{[i]}^{(n)} \right)_{i \in \mathbb{N}}$ , the probability that all individuals are in the same block  $i$  is given by drawing without replacement, i.e.

$$\mathbb{P} \left( [m] \subseteq \text{block } i \mid B_{[i]}^{(n)} \right) = \prod_{\ell=0}^{m-1} \frac{B_{[i]}^{(n)} - \ell}{n - \ell}.$$

Summing this up over all blocks and taking the expectation yields  $\mathbb{P}(A_{m,n}) = 1 - p_{n,m}^{(\Xi)}$ , thus establishing Eq. (8) (by definition there is more than one block at time  $T_{\text{MRCA}}^{(n)}$  so  $p_{n,m}^{(\Xi)} > 0$ .)

To show the convergence in Eq. (9) we first establish that all objects are well defined. Assume now that the  $\Xi$ -coalescent comes down from infinity, so at any time  $t > 0$ , there are only finitely many blocks in the partition  $\Pi_t$  almost surely. For  $n \rightarrow \infty$ , Kingman's correspondence [55, Thm. 2] ensures that the asymptotic frequencies of the blocks in the partition  $\Pi_t$  of  $\mathbb{N}$  exist almost surely and are limits of the block frequencies in the  $n$ -coalescent as written in the proposition. Pick an arbitrarily small  $t > 0$ . Then, consider only paths where  $\Pi_t$  has more than one block. Since the number of blocks of  $\Pi_t$  is finite a.s. we can find  $n_0 \in \mathbb{N}$  so that

$\Pi_t^{(n_0)}$  has at least one individual in any block of  $\Pi_t$  (thus has the same number of blocks). By construction, from time  $t$  onwards, the  $\Xi$ - $n_0$ -coalescent merges the blocks in exactly the same (Markovian) manner as the  $\Xi$ -coalescent. So if  $\Pi_t$  has more than one block,  $T_{\text{MRCA}}^{(n)} = T_{\text{MRCA}}^{(\infty)}$  for  $n \geq n_0$  and the asymptotic frequencies at  $T_{\text{MRCA}}^{(\infty)}$  exist (since their corresponding blocks are a specific merger of the blocks of  $\Pi_t$  whose block frequencies exist). Now  $T_{\text{MRCA}}^{(2)} \leq T_{\text{MRCA}}^{(\infty)}$  almost surely and  $T_{\text{MRCA}}^{(2)}$  is  $\text{Exp}(\Xi(\Delta))$ -distributed. Therefore, for almost every path, we can choose  $t < T_{\text{MRCA}}^{(2)}$  so that  $\Pi_t$  has more than one block.

We have established that all objects are well defined; now we show the actual convergence in (9). For  $\mathbf{x} \in \Delta$  let

$$f_{n,m}(\mathbf{x}) := \sum_{i \in \mathbb{N}} \prod_{\ell=0}^{m-1} \frac{nx_i - \ell}{n - \ell}$$

and  $f_m(\mathbf{x}) := \sum_{i \in \mathbb{N}} x_i^m$ . We have  $f_{n,m} \rightarrow f_m$  uniformly on  $\Delta$  and that  $f_m$  is continuous on  $\Delta$  in the  $\ell^1$ -norm with  $0 \leq f_m \leq 1$ . We can rewrite, using Eq. (8),

$$p_{n,m}^{(\Xi)} = 1 - \mathbb{E} \left[ f_{n,m} \left( \left( \frac{1}{n} B_{[i]}^{(n)} \right)_{i \in \mathbb{N}} \right) \right].$$

For any  $\varepsilon > 0$  we find  $n_0$  so that for  $n \geq n_0$

$$\begin{aligned} & \left| \mathbb{E} \left[ f_{n,m} \left( \left( \frac{1}{n} B_{[i]} \right)_{i \in \mathbb{N}} \right) \right] - \mathbb{E} \left[ f_m \left( P_{[i]} \right)_{i \in \mathbb{N}} \right] \right| \\ & \leq \left| \mathbb{E} \left[ f_{n,m} \left( \left( \frac{1}{n} B_{[i]} \right)_{i \in \mathbb{N}} \right) \right] - \mathbb{E} \left[ f_m \left( \left( \frac{1}{n} B_{[i]} \right)_{i \in \mathbb{N}} \right) \right] \right| \\ & \quad + \left| \mathbb{E} \left[ f_m \left( \left( \frac{1}{n} B_{[i]} \right)_{i \in \mathbb{N}} \right) \right] - \mathbb{E} \left[ f_m \left( P_{[i]} \right)_{i \in \mathbb{N}} \right] \right| \leq 2\varepsilon. \end{aligned}$$

We have used uniform convergence of  $f_{n,m}$  to control the first difference and the convergence (in law) of  $(n^{-1} B_{[i]}^{(n)})_{i \in \mathbb{N}}$  to  $(P_{[i]})_{i \in \mathbb{N}}$  to control the second.

The representation of the limit in Eq. (9) in terms of  $X$  and  $Y$  follows directly from the properties of exchangeable partitions (c.f. for example [9, 10]). The first equality is [9, Eq. (1.4)], while the second equality uses the correspondence between the distribution of a size-biased and a uniform pick of a block, see [9, Eq. (1.2)]. By definition  $\Pi_{T_{\text{MRCA}}^{(\infty)}}$  has more than one block almost surely so the limit in Eq. (9) is  $> 0$ .

□

*Remark 1* Reordering the block frequencies, e.g. in order of least elements of blocks, does not change Eq. (9).

### 4.3 Proof of Prop. 3

*Proof* We use the construction of [39] in which the Bolthausen-Sznitman coalescent is obtained by cutting a random recursive tree  $\mathbb{T}_n$  with  $n$  nodes at independent  $\text{Exp}(1)$  times, see Sec. A2. Consider the last merger in the Bolthausen-Sznitman  $n$ -coalescent. In terms of cutting edges of  $\mathbb{T}_n$ , it is reached when the last edge

connected to the root of  $\mathbb{T}_n$  is cut. Let  $E_n$  be the number of such edges in  $\mathbb{T}_n$ . For  $T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)}$ , we need that not all  $i \in [m]$  are in a single block of the  $n$ -coalescent before the last merger (see proof of Prop. 2).

By construction, for any node with label in  $[m]$ , on the path to the node labelled 1 (root) in  $\mathbb{T}_n$ , the last node passed before reaching the root must also have a label from  $[m]$ . Thus, any node connected to the root of  $\mathbb{T}_n$  that is labelled from  $[n]_{m+1}$  cannot root a subtree that includes any nodes labelled from  $[m]$ .

Now, we consider the last edge of  $\mathbb{T}_n$  cut in the construction of the Bolthausen-Sznitman  $n$ -coalescent, which causes the MRCA of the  $n$ -coalescent to be reached. It has to be connected to the root. Consider the two subtrees on both sides of the edge cut last. One subtree contains the root, thus includes at least the label 1 from  $[m]$ . If the other subtree is rooted in a node labelled from  $[m]$ , we have

$T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)}$ , since both subtrees contain labels of  $[m]$ , thus not all  $i \in [m]$  are in a single block of the  $n$ -coalescent before the last merger. If the subtree not containing the root has a root labelled from  $[n]_{m+1}$ , as argued above, it contains no labels from  $[m]$ . Additionally, since we are at the last cut, all other edges connected to the root of  $\mathbb{T}_n$  have already been cut and all labels in the subtrees rooted by them joined with label 1. Thus, all labels in  $[m]$  are labelling the root before the last cut, which corresponds to  $[m]$  being a subset of a block of the  $n$ -coalescent before the last merger, hence  $T_{\text{MRCA}}^{(m;n)} \neq T_{\text{MRCA}}^{(n)}$ .

This shows  $T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)}$  if and only if the last edge cut is an edge connecting a node labelled from  $[m]$  with the root. Let  $E_m$  be the number of edges of  $\mathbb{T}_n$  connected to the root labelled from  $[m]$  and  $E_n$  be the total number of edges connected to the root. Then,

$$\mathbb{P}\left(T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)}\right) = \mathbb{E}\left[\frac{E_m}{E_n}\right], \quad (27)$$

because given  $\mathbb{T}_n$ ,  $E_m/E_n$  is the probability that the edge cut last is connected to a node with a label from  $[m]$ ; edges are cut at i.i.d. times, so the edge cut last is uniformly distributed among all edges connected to the root.

As we see from the sequential construction of  $\mathbb{T}_n$ ,  $E_m$  is the number of edges connected to 1 when the first  $m$  nodes are set, the resulting tree is a random recursive tree  $\mathbb{T}_m$  with  $n$  leaves. The numbers  $E_n$  and  $E_m$  can be described in terms of a Chinese restaurant process (CRP), see [39, p. 724]: The number of edges connected to node 1 is distributed as the number of tables in a CRP with  $n$  (resp.  $m$ ) customers.

This distribution is  $E_i \stackrel{d}{=} B_1 + \dots + B_i$  ( $i \in \{m, n\}$ ), where  $B_1, \dots, B_i$  are independent Bernoulli variables with  $P(B_j = 1) = j^{-1}$ , see e.g. [3, p. 10]. The sequential construction of the random recursive trees (and the connected CRPs) ensures that the  $B_1, \dots, B_m$  are identical for  $E_m$  and  $E_n$ . This establishes the equality of Eq.'s (27) and (10).

From the proof of [38, Lemma 3], we have  $\frac{\log(n)}{E_n} \rightarrow 1$  in  $L^1$  for  $n \rightarrow \infty$ . The sequence  $(E_m/E_n)_{n \in \mathbb{N}}$  is bounded a.s. Thus, bounded convergence ensures

$$\lim_{n \rightarrow \infty} \mathbb{E}\left[\log(n) \frac{E_m}{E_n}\right] = \mathbb{E}\left[\lim_{n \rightarrow \infty} \frac{\log(n) E_m}{E_n}\right] = \mathbb{E}[E_m] = 1 + \frac{1}{2} + \dots + \frac{1}{m-1}.$$

□

#### 4.4 Proof of Prop. 4

*Proof* As in Subsection 3 we use the construction of the Bolthausen-Sznitman  $n$ -coalescent described in [39]. We wish to establish the probability that the MRCA of a subsample of size  $m$  from a sample of size  $n$  subtends only the subsample in the  $n$ -coalescent. We will also use the Bernoulli variables  $B_i$ ,  $i \in [n]$  of  $\mathbb{T}_n$  as in Prop. 3, where  $B_i = 1$  if the node labelled  $i$  is directly connected to the root (node labelled 1). If we look at the cutting procedure which constructs the Bolthausen-Sznitman  $n$ -coalescent from  $\mathbb{T}_n$ , we observe that no path of  $\mathbb{T}_n$  can contribute positive probability to  $q_{n,m}^{(\text{BS-coal})}$  that attaches any node labelled from  $[n]_{m+1}$  to a node labelled from  $[m]_2$ . If we do attach a node labelled  $i \in [n]_{m+1}$  to a node labelled from  $[m]_2$ , when constructing the Bolthausen-Sznitman  $n$ -coalescent we will cut an edge on the path from the node labelled  $i$  to the root labelled 1 before the MRCA of  $[m]$  is reached, thus  $i$  would be subtended by the MRCA of  $[m]$ . The probability that a node labelled  $i \in [n]_{m+1}$  is not connected to a node labelled from  $[m]_2$  in  $\mathbb{T}_n$  is

$$\prod_{i=1}^{n-m} \frac{i}{m+i-1} = \binom{n-1}{m-1}^{-1}.$$

Even when there is no edge connecting a node labelled from  $[n]_{m+1}$  directly with a node labelled from  $[m]_2$ , not all such paths of  $\mathbb{T}_n$  will contribute to  $q_{n,m}^{(\text{BS-coal})}$ . To contribute, we need that the cutting procedure does not lead to any  $i \in [n]_{m+1}$  being subtended by the MRCA of  $[m]$ . For the mentioned paths, this happens if and only if we cut all edges connecting nodes labelled from  $[m]_2$  to 1 before cutting any edge connecting 1 to nodes labelled from  $[n]_{m+1}$ . We have  $\sum_{i \in [m-1]} B_i$  edges adjacent to node 1, see the proof of Prop. 3. With the constraint that no edge connects a node labelled from  $[n]_{m+1}$  directly with a node labelled from  $[m]_2$ , the sequential construction yields that, after relabelling, the nodes labelled with  $\{1\} \cup [n]_{m+1}$  form a  $\mathbb{T}_{n-m+1}$  and thus there are  $\sum_{i \in [n-m]} B'_i$  edges adjacent to the root of  $\mathbb{T}_n$  connecting to the nodes labelled with  $\{1\} \cup [n]_{m+1}$ , where  $B'_i \stackrel{d}{=} B_i$  for independent  $B'_i$ . All edges adjacent to the node labelled 1 need to be cut before the MRCA of  $[n]$  is reached and they are cut at independent  $\text{Exp}(1)$  times. This means that the probability of cutting all edges connecting 1 to nodes labelled from  $[m]_2$  first is just drawing  $\sum_{i \in [m-1]} B_i$  times without replacement from  $\sum_{i \in [m-1]} B_i + \sum_{i \in [n-m]} B'_i$  edges, where all  $\sum_{i \in [m-1]} B_i$  edges connecting nodes labelled from  $[m]_2$  have to be drawn. This probability equals

$$\left( \frac{\sum_{i \in [m-1]} B_i + \sum_{i \in [n-m]} B'_i}{\sum_{i \in [m-1]} B_i} \right)^{-1}.$$

Integrating over all contributing paths of  $\mathbb{T}_n$  with the cutting constraint described above finishes the proof.  $\square$

#### 4.5 Proof of Thm. 1

*Proof* We track the asymptotic frequencies  $(P_{[i]}(t))_{t \geq 0}$  of the  $i$ th biggest block for all  $t > 0$  and  $i \in \mathbb{N}$ . Consider a non-negative and measurable  $[0, \infty)$ -valued function

667  $g$  on the  $k$ -dimensional simplex

$$\Delta_k := \{(x_1, \dots, x_k) : x_1 \geq x_2 \geq \dots \geq x_k \geq 0, \sum_{i \in [k]} x_i = 1\}$$

668 that is invariant under permutations  $(x_1, \dots, x_k) \mapsto (x_{\sigma(1)}, \dots, x_{\sigma(k)})$ . [8, Thm. 1.2]  
669 shows that

$$\begin{aligned} & \mathbb{E} [g((P_{[i]}(T_k))_{i \in \mathbb{N}}) | N(T_k) = k] \\ &= \mathbb{E} [(Y_1 + \dots + Y_k)^{1-\alpha}]^{-1} \mathbb{E} \left[ (Y_1 + \dots + Y_k)^{1-\alpha} g \left( \left( \frac{Y_i}{Y_1 + \dots + Y_k} \right)_{i \in [k]} \right) \right], \end{aligned}$$

670 where  $T_k$  is the waiting time until a state with  $\leq k$  blocks is hit by the Beta-  
671 coalescent and  $N(t)$  is the number of blocks of  $\Pi_t$ , thus we condition on the coa-  
672 lescent to hit a state with exactly  $k$  blocks.

673 We can apply this formula to compute  $\mathbb{E} [\sum_{i \in [K]} P_i^m]$  from Eq. (9), where  $K$  is the  
674 number of blocks at the last collision of the Beta-coalescent. For this, condition on  
675  $K = k$ . With  $\{K = k\} = \{N(T_k) = k\} \cap \{\text{all blocks of } \Pi_{T_k} \text{ merge at the next merger}\}$ ,  
676 the strong Markov property shows that the block frequencies at  $T_k$  are independent  
677 of them merging at the next collision. However, these frequencies are, conditioned  
678 on  $K$ , just  $(P_i)_{i \in [K]}$ . For  $\underline{x} \in \Delta_k$  we set  $g_m(\underline{x}) = \sum_{i=1}^k x_i^m$  (which fulfills all necessary  
679 conditions to apply [8, Thm. 1.2]) and compute

$$\begin{aligned} & \mathbb{E} \left[ \sum_{i \in [K]} P_i^m \right] \\ &= \sum_{k \in \mathbb{N}} \mathbb{E} \left[ \sum_{i \in [K]} P_i^m | K = k \right] \mathbb{P}(K = k) \\ &= \sum_{k \in \mathbb{N}} \mathbb{E} [g((P_i(T_k))_{i \in \mathbb{N}}) | N(T_k) = k] \mathbb{P}(K = k) \\ &= \sum_{k \in \mathbb{N}} \mathbb{E} [(Y_1 + \dots + Y_k)^{1-\alpha}]^{-1} \mathbb{E} \left[ \sum_{i=1}^k \frac{Y_i^m}{(Y_1 + \dots + Y_k)^{\alpha+m-1}} \right] \mathbb{P}(K = k). \end{aligned}$$

680 The distribution of  $K$  for the Beta-coalescent is known from [47, Thm. 3.5]. Using  
681 that

$$\left( \frac{Y_i^m}{(Y_1 + \dots + Y_k)^{\alpha+m-1}} \right)_{i \in [k]}$$

are identically distributed and  $p_m^{(\text{Beta-coal})} = 1 - \mathbb{E} [\sum_{i \in [K]} P_i^m]$  completes the proof.  $\square$

## 682 4.6 Proof of Prop. 5

683 *Proof* Consider any  $\Xi$ -coalescent (and its restrictions to  $[n]$ ,  $n \in \mathbb{N}$ ). Since, for  
684 nested samples,  $T_{\text{MRCA}}^{(m;n)} \leq T_{\text{MRCA}}^{(n)}$  almost surely for any  $n \geq m$ , we have

$$\{T_{\text{MRCA}}^{(m;m+i)} = T_{\text{MRCA}}^{(m+i)}\} \supseteq \{T_{\text{MRCA}}^{(m;m+i+1)} = T_{\text{MRCA}}^{(m+i+1)}\}$$

for any  $i \in \mathbb{N}$ . Thus,  $p_m^{(\mathcal{E})} = \lim_{n \rightarrow \infty} p_{n,m}^{(\mathcal{E})} = \mathbb{P}^{(\mathcal{E})} \left( T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \forall n > m \right)$  exists.

Suppose first that the  $\mathcal{E}$ -coalescent comes down from infinity. Then, Eq. (9) shows  $p_m^{(\mathcal{E})} > 0$ .

If the  $\mathcal{E}$ -coalescent stays infinite,  $\tau_m$  is almost surely finite, while  $T_{\text{MRCA}}^{(n)} \rightarrow \infty$  almost surely. Thus,  $p_m^{(\mathcal{E})} = 0$ .

Consider a  $\mathcal{E}$ -coalescent that neither comes down from infinity nor stays infinite. Then,  $\mathbb{E}(\{\mathbf{x} \in \Delta \mid \sum_{i=1}^k x_i = 1 \text{ for } k \in \mathbb{N}\}) > 0$ . As stated in the introduction, in this case there is an almost surely finite waiting time  $T$  with  $\#\Pi_T < \infty$  almost surely. Let  $n_T$  be the finite number of blocks at time  $T$ . Again, exchangeability ensures, as in proving Eq. (8), that there is a positive probability that not all  $i \in [m]$  are in the same block of  $\Pi_T$  (so in particular, with positive probability,  $T_{\text{MRCA}}^{(m)} > T$ ). The strong Markov property of the  $\mathcal{E}$ -coalescent ensures that, given  $n_T$ ,  $\Pi_T$  evolves like a  $\mathcal{E}$ - $n_T$ -coalescent, which can have at most  $n_T$  mergers. In summary, with positive probability, more than one of the  $n_T$  blocks at time  $T$  includes individuals from the subset  $[m]$  and the  $n_T$  blocks are merged following a  $\mathcal{E}$ - $n_T$ -coalescent. Then, Eq. (8) shows that with positive probability, conditioned on the event that  $k > 1$  blocks of  $\Pi_T$  contain individuals from  $m$ , also more than one block of the  $\mathcal{E}$ -coalescent at its last collision contains individuals of  $[m]$ .  $\square$

Prop. 5 shows that  $\mathbb{P}^{(\mathcal{E})} \left( T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \right) \rightarrow 0$  for fixed  $m$  and  $n \rightarrow \infty$  if the  $\mathcal{E}$ -coalescent stays infinite. The Bolthausen-Sznitman coalescent stays infinite [77, Example 15]; however, convergence to 0 is only of order  $\mathcal{O}(1/\log(n))$ .

#### 4.7 Proof of recursions (5), (14), and (12)

The strong Markov property of a  $\Lambda$ -coalescent together with a natural coupling which we will introduce below allows us to describe many functionals of multiple-merger  $n$ -coalescents recursively by conditioning on their first jump, e.g. see [40] or [61]. We use this to prove recursions (5) and (12).

##### 4.7.1 Proof of Eq. (5)

Consider the probability  $p_{n,m}^{(\Lambda)}$  (see Eq. (5)) that a sample of size  $n$  shares the MRCA with a subsample of size  $m \in [n-1]_2$ . The boundary conditions  $p_{m,m} = 1$  and  $p_{n,1} = 0$  for  $n > 1$  follow directly from the definition. We record how many individuals are merged at the first jump of the  $n$ -coalescent. Suppose a  $k$ -merger occurs which happens with probability  $\beta(n, n-k+1)$ . Conditional on a  $k$ -merger,  $\ell \leq m$  of individuals that merge are taken from the subsample and  $n-\ell$  are not with probability  $\binom{m}{\ell} \binom{n-m}{k-\ell} / \binom{n}{k}$ , since the individuals that merge are picked uniformly at random without replacement. For  $p_{n,m}^{(\Lambda)} > 0$ , we need that not all  $m$  individuals are merged unless all  $n$  individuals are merged, thus  $\ell < m$  or  $k = n$ . Writing  $C(k, \ell)$  for the event that exactly  $\ell$  lineages from the subsample are merged (with  $\ell < m$  or

705  $k = n$ ), the strong Markov property shows that

$$\mathbb{P}^{(\Pi^{(A)})} \left( T_{\text{MRCA}}^{(m;n)} = T_{\text{MRCA}}^{(n)} \mid C(k, \ell) \right) = \mathbb{P}^{(\Pi^A)} \left( T_{\text{MRCA}}^{(m'; n-k+1)} = T_{\text{MRCA}}^{(n-k+1)} \right)$$

706 with  $m' = (m - \ell + 1)\mathbb{1}_{(\ell > 1)} + m\mathbb{1}_{(\ell \leq 1)}$ , since among the ancestral lines (blocks) af-  
707 ter the first collision,  $m'$  subtend the subsample. Summing over all possible values  
708 (recall the boundary conditions) yields recursion (5).  $\square$

#### 709 4.7.2 Proof of Eq. (14)

710 We again condition on the event that  $k$  blocks are merged at the first jump. Only  $k$ -  
711 mergers where either all merged individuals are picked from the subsample  $[m]$  or  
712 none is sampled from  $[m]$  contribute positive probability to  $q_{n,m}^{(A)}$ . After the jump,  
713 we thus have  $n - k + 1$  ancestral lineages present, from which either  $m - k + 1$   
714 or  $m$  are connected to the subsample. The strong Markov property and sampling  
715 without replacement for the  $k$ -merger then yields Eq. (14).

#### 716 4.7.3 Proof of Eq. (12)

717 Recall the natural coupling: if we restrict an  $n$ -coalescent with mutation rate  $\theta$  to  
718 any  $\ell$ -sized subset  $L \subseteq [n]$ , the restriction is an  $\ell$ -coalescent with mutation with the  
719 same rate  $\theta$ . To prove recursion (12) we partition over three possible outcomes of  
720 the first event: it is a mutation on a lineage subtending the subsample ( $E_1$ ), it is a  
721 mutation on a lineage not subtending the subsample ( $E_2$ ), or it is a merger ( $E_3$ ).  
722 Naturally, before any mutation occurs, all edges are active.

723 We recall a few elementary facts. The time to the first mutation on any lineage  
724 is  $\text{Exp}(\theta/2)$ -distributed (mutations on different/disjoint lineages are independent)  
725 and independent of the waiting time for the first merger. The minimum of indepen-  
726 dent exponential r.v.'s  $X_1, \dots, X_i$  with parameters  $\alpha_1, \dots, \alpha_i$  is again exponentially  
727 distributed with parameter  $\sum_{j=1}^i \alpha_j$ . Finally,  $\mathbb{P}(X_1 \leq X_2) = \alpha_1/(\alpha_1 + \alpha_2)$ .  
728

729 The waiting times  $X_i$  for events  $E_i$  for  $1 \leq i \leq 3$  are all exponential; the one for  
730  $E_1$  with rate  $\theta m/2$ , for  $E_2$  with rate  $\theta(n - m)/2$ , and for  $E_3$  with rate  $\lambda(n)$ . The  
731 probability of event  $E_1$  is  $\mathbb{P}(E_1) = \theta m/(2\lambda(n) + \theta n)$  and, conditional on  $E_1$ ,  
732

$$\left\{ A^{(n)}(\max\{\tau_i : i \in [m]\}) = 0 \right\}$$

733 is determined by the  $n - 1$  active lineages after the event. The memoryless property  
734 of the exponential and natural coupling imply that after the first event, conditional  
735 on that event being  $E_1$ , the remaining  $n - 1$  lineages, of which  $(m - 1)$  subtend the  
736 subsample, follow an  $(n - 1)$ -coalescent with mutation rate  $\theta$ . Thus,

$$\mathbb{P} \left( A^{(n)}(\max\{\tau_i : i \in [m]\}) = 0 \mid E_1 \right) = p_{n-1, m-1}^{(\Pi^{(A)})}.$$



Analogously, the probability of event  $E_2$  is  $\mathbb{P}(E_2) = \theta(n-m)/(2\lambda(n) + \theta n)$ . Given  $E_2$ , we need to follow the coalescent of  $n-1$  lineages, of which  $m$  are from the subsample, which gives

$$\mathbb{P}\left(A^{(n)}(\max\{\tau_i : i \in [m]\}) = 0 | E_2\right) = p_{n,m-1}^{(\Pi^{(\Lambda)}, \theta)}.$$

We have  $P(E_3) = 1 - P(E_1) - P(E_2) = 2\lambda(n)/(2\lambda(n) + \theta n)$ . To compute

$$\mathbb{P}\left(A^{(n)}(\max\{\tau_i : i \in [m]\}) = 0 | E_3\right),$$

proceed exactly as in the proof of recursion (5) by partitioning over the number of mutant lineages involved in the merger, but with changed boundary conditions since  $p_{i,1}^{(\Pi^{(\Lambda)}, \theta)} > 0$ , while  $p_{i,1}^{(\Pi^{(\Lambda)})} = 0$  for  $i > 1$ . Summing over  $E_1, E_2, E_3$  yields Eq. (12).  $\square$

#### 4.8 Proof of Eq. (15)

Recall our assumption that block  $\pi_1$  always contains element 1. To see (15), we will show that, for  $n$  large enough,

$$\mathbb{P}^{(\Pi^{(\Lambda)})}\left(T_{\text{MRCA}}^{(m;n)} \geq \inf\{t \geq 0 : \pi_1 \cap [n]_{m+1} \neq \emptyset, \pi_1 \in \Pi_t\}\right) = 1. \quad (28)$$

In words, the smallest block containing  $[m]$  appearing in the  $n$ -coalescent will always contain at least  $m+1$  elements; block  $[m]$  will almost never be observed.

Hence,  $\lim_{n \rightarrow \infty} q_{n,m}^{(\Pi^{(\Lambda)})} = 0$ .

Consider first  $\Lambda$  with  $\int_{[0,1]} x^{-1} \Lambda(dx) = \infty$ , which makes the  $\Lambda$ -coalescent dust-free (no singleton blocks almost surely for  $t > 0$ ) - see the proof of [69, Lemma 25]. For  $t > 0$ , [69, Prop. 30] shows that the partition block  $\pi_1 \in \Pi_t^{(n, \Lambda)}$  containing individual 1 at time  $t$  in the  $\Lambda$ - $n$ -coalescent  $\{\Pi_t^{(n, \Lambda)}, t \geq 0\}$  fulfills  $\lim_{n \rightarrow \infty} \#\pi_1/n > 0$  almost surely. Thus, individual 1 has already merged before any time  $t > 0$  if  $n > N'$ , where  $N'$  is a random variable on  $\mathbb{N}$  almost surely. However, within the subsample of fixed size  $m$ , we wait an exponential time with rate  $\lambda(m)$  for any merger of individuals in  $[m]$ . Thus, for  $n$  large enough individual 1 has almost surely already merged with individuals of  $[n]_{m+1}$  before merging with another individual in the subsample. Consider now  $\Lambda$  with  $\int_{[0,1]} x^{-1} \Lambda(dx) < \infty$ , which shows that the coalescent has dust, i.e. there is a positive probability that there is a positive fraction of singleton blocks at any time  $t$ , see [69, Prop. 26]. In this case [37, Corollary 2.3] shows that at its first merger, for  $n \rightarrow \infty$ , individual 1 merges with a positive fraction of all individuals  $\mathbb{N}$  almost surely, which has to include individuals in  $[n]_{m+1}$ . Since this is the earliest merger where the MRCA of  $[m]$  can be reached, the proof is complete.  $\square$

Analogously we have the following:

**Corollary 1** Consider any  $\Xi$ -coalescent which comes down from infinity and its restrictions to  $[n]$ ,  $n \in \mathbb{N}$ . Fix subsample size  $m \in \mathbb{N}$ . Let  $\tilde{T}_m^{(n)}$  be the first time that any  $i \in [m]$  is involved in a merger in the  $\Xi$ - $n$ -coalescent for  $n \geq m$ . We have

$$\lim_{n \rightarrow \infty} \mathbb{P} \left( \tilde{T}_m^{(n)} = T_{MRCA}^{(n)} \right) = 0.$$

*Proof* If the  $\Xi$ -coalescent comes down from infinity, it fulfills  $\int_{[0,1]} x^{-1} \Lambda(dx) = \infty$ , since it has to be dust-free. As above, we see that individual 1 has already merged before  $T_{MRCA}^{(n)}$  for  $n \rightarrow \infty$ , which establishes the corollary.  $\square$

## 5 Conclusion and open questions

By studying properties of nested samples we have aimed at understanding how much information about the evolutionary history of a population can be extracted from a sample, i.e. how the genealogical information increases if we enlarge the sample. In particular, we have focussed on multiple-merger coalescent processes derived from population models of high fecundity and sweepstakes reproduction. In comparison with the Kingman-coalescent the general conclusion, at least for the statistics we consider, is that a subsample represents less well the ‘population’ or the complete sample from which the subsample was drawn when the underlying coalescent mechanism admits multiple mergers. The subsample reaches its MRCA sooner and shares less of the ancestral genetic variants (internal branches) with the complete sample under a MMC process than under the Kingman-coalescent. A similar conclusion can be broadly reached in comparison with exponential population growth. This seems to imply that one would need a larger sample for inference under a multiple-merger coalescent than under a (time-changed) Kingman-coalescent. Large sample size has been shown to impact inference under the Wright-Fisher model [11], in particular if the sample size exceeds the effective size [85]. The main effect is that when sample size is large enough, one starts to notice multiple and/or simultaneous mergers in the trees. The implication is that if one sampled a whole finite population with Wright-Fisher reproduction, the genealogy of the whole population is not well approximated by the Kingman-coalescent. One would also expect an impact of large sample size on inference under MMC. The effective size in HFSR populations can be much smaller than in a Wright-Fisher population with the same census size [78, 46, 86]. The implication is that for almost any finite population, the genealogy of the whole population is not the genealogy one derives under the assumption of a small sample size compared to population size. This therefore leaves the question of what one is making an inference about when one applies a coalescent-based inference method.

Our sample size considerations raise another point. Our main result on a representation for  $p_m^{(\text{Beta-coal})}$  rests on the assumption that the Beta-coalescent holds for an infinitely large population. The limit behavior of  $p_{n,m}^{(\Xi)}$  as  $n \rightarrow \infty$  described in Prop. 5 also depends on the ability of  $\Xi$ -coalescents to be applicable to an infinitely large sample size. The Beta-coalescent was initially derived as an approximation of the random tree describing the ancestral relations of a finite number of gene copies; the same holds for the Kingman-coalescent, and in general

for any coalescent process derived from a population model. Large sample size should impact the coalescent process derived for a HFSR model just as it does for the Wright-Fisher model. Yet the Beta-coalescent, and any coalescent that satisfies certain assumptions (see [77]) comes down from infinity. There is a curious disconnect between the concept of coming down from infinity and the impact of large sample size which has not been bridged. At least,  $p_m^{(\Xi)}$  gives a lower bound on  $p_{n,m}^{(\Xi)}$  for any  $n$  until the coalescent approximation breaks down. This interpretation makes sense for real populations.

All our results are applicable to a single non-recombining locus. A natural question to ask is if and how our results might change if we considered multiple unlinked loci. How would the statistics we consider, averaged over many unlinked loci, behave under MMC in comparison with a (time-changed) Kingman-coalescent? DNA sequencing technology has advanced to the degree that sequencing whole genomes is now almost routine (see eg. [43,4]). One could ask how large a sample from a HFSR population does one need to be confident to have sampled a significant fraction of the genome-wide ancestral variation? In this context, let  $T_{\text{MRCA}}^{(n,\ell)}$  denote the TMRCA of the complete sample of size  $n$  at a non-recombining locus  $\ell \in [L]$ , and  $T_{\text{MRCA}}^{(m;n,\ell)}$  the TMRCA of a nested subsample of size  $m$  at same locus. Then we would like to compare the probability

$$\mathbb{P}(\Pi) \left( \bigcap_{\ell \in [L]} \left\{ T_{\text{MRCA}}^{(m;n,\ell)} = T_{\text{MRCA}}^{(n,\ell)} \right\} \right)$$

between different coalescent processes. And in fact, the independence of the genealogies at unlinked loci under the Kingman-coalescent, and Eq. (7), gives

$$\mathbb{P}(\text{Kingman}) \left( \bigcap_{\ell \in [L]} \left\{ T_{\text{MRCA}}^{(m;n,\ell)} = T_{\text{MRCA}}^{(n,\ell)} \right\} \right) = \left( \frac{(m-1)(n+1)}{(m+1)(n-1)} \right)^L.$$

Under a MMC process the genealogies at unlinked loci are not independent (see e.g. [32, 15]).

We compared results from single-locus MMC models with a time-changed Kingman-coalescent derived from a single-locus model of exponential population growth. Naturally one would like to compare results between genomic (multi-locus) models of HFSR with population growth to genomic models of HFSR without growth, and to genomic models of growth without HFSR. Some mathematical handle on the distributions of the quantities we simulated would (obviously) also be nice. However, these will have to remain important open tasks.

**Acknowledgements** BE was funded by DFG grant STE 325/17-1 to Wolfgang Stephan through Priority Programme SPP1819: Rapid Evolutionary Adaptation. FF was funded by DFG grant FR 3633/2-1 through Priority Program 1590: Probabilistic Structures in Evolution.

## References

1. Agrios, G.: Plant pathology. Academic Press, Amsterdam (2005)

2. Árnason, E., Halldórsdóttir, K.: Nucleotide variation and balancing selection at the *Ckma* gene in Atlantic cod: analysis with multiple merger coalescent models. *PeerJ* **3**, e786 (2015). DOI 10.7717/peerj.786. URL <http://dx.doi.org/10.7717/peerj.786>
3. Arratia, R., Barbour, A.D., Tavaré, S.: Logarithmic Combinatorial Structures: A Probabilistic Approach. European Mathematical Society (EMS), Zürich (2003)
4. Barney, B.T., Munkholm, C., Walt, D.R., Palumbi, S.R.: Highly localized divergence within supergenes in atlantic cod (*gadus morhua*) within the gulf of maine. *BMC Genomics* **18**(1) (2017). DOI 10.1186/s12864-017-3660-3. URL <https://doi.org/10.1186/s12864-017-3660-3>
5. Barton, N.H., Etheridge, A.M., Véber, A.: Modelling evolution in a spatial continuum. *Journal of Statistical Mechanics: Theory and Experiment* **2013**(01), P01,002 (2013). URL <http://stacks.iop.org/1742-5468/2013/i=01/a=P01002>
6. Basu, A., Majumder, P.P.: A comparison of two popular statistical methods for estimating the time to most recent common ancestor (tmrca) from a sample of DNA sequences. *Journal of genetics* **82**(1-2), 7–12 (2003)
7. Berestycki, J., Berestycki, N., Schweinsberg, J.: Beta-coalescents and continuous stable random trees. *Ann Probab* **35**, 1835–1887 (2007)
8. Berestycki, J., Berestycki, N., Schweinsberg, J.: Small-time behavior of beta coalescents. *Ann Inst H Poincaré Probab Statist* **44**, 214–238 (2008)
9. Berestycki, N.: Recent progress in coalescent theory. *Ensaio Matemáticos* **16**, 1–193 (2009)
10. Bertoin, J.: Exchangeable coalescents. *Cours d'école doctorale* pp. 20–24 (2010)
11. Bhaskar, A., Clark, A., Song, Y.: Distortion of genealogical properties when the sample size is very large. *PNAS* **111**, 2385–2390 (2014)
12. Birkner, M., Blath, J.: Computing likelihoods for coalescents with multiple collisions in the infinitely many sites model. *J Math Biol* **57**, 435–465 (2008)
13. Birkner, M., Blath, J.: coalescents and population genetic inference. *Trends in stochastic analysis* (353), 329 (2009)
14. Birkner, M., Blath, J., Capaldo, M., Etheridge, A.M., Möhle, M., Schweinsberg, J., Wakolbinger, A.: Alpha-stable branching and beta-coalescents. *Electron. J. Probab* **10**, 303–325 (2005)
15. Birkner, M., Blath, J., Eldon, B.: An ancestral recombination graph for diploid populations with skewed offspring distribution. *Genetics* **193**, 255–290 (2013)
16. Birkner, M., Blath, J., Eldon, B.: Statistical properties of the site-frequency spectrum associated with  $\Lambda$ -coalescents. *Genetics* **195**, 1037–1053 (2013)
17. Birkner, M., Blath, J., Möhle, M., Steinrück, M., Tams, J.: A modified lookdown construction for the Xi-Fleming-Viot process with mutation and populations with recurrent bottlenecks. *ALEA Lat. Am. J. Probab. Math. Stat.* **6**, 25–61 (2009)
18. Birkner, M., Blath, J., Steinrück, M.: Analysis of DNA sequence variation within marine species using Beta-coalescents. *Theor Popul Biol* **87**, 15–24 (2013)
19. Blath, J., Cronjäger, M.C., Eldon, B., Hammer, M.: The site-frequency spectrum associated with  $\Xi$ -coalescents. *Theoretical Population Biology* **110**, 36–50 (2016). DOI 10.1016/j.tpb.2016.04.002
20. Bolthausen, E., Sznitman, A.: On Ruelle's probability cascades and an abstract cavity method. *Comm Math Phys* **197**, 247–276 (1998)
21. Capra, J.A., Stolzer, M., Durand, D., Pollard, K.S.: How old is my gene? *Trends in Genetics* **29**(11), 659–668 (2013)
22. Desai, M.M., Walczak, A.M., Fisher, D.S.: Genetic diversity and the structure of genealogies in rapidly adapting populations. *Genetics* **193**(2), 565–585 (2013)
23. Dong, R., Gnedin, A., Pitman, J.: Exchangeable partitions derived from markovian coalescents. *The Annals of Applied Probability* pp. 1172–1201 (2007)
24. Donnelly, P., Kurtz, T.G.: Particle representations for measure-valued population models. *Ann Probab* **27**, 166–205 (1999)
25. Donnelly, P., Tavaré, S.: Coalescents and genealogical structure under neutrality. *Annual review of genetics* **29**(1), 401–421 (1995)
26. Durrett, R.: Probability models for DNA sequence evolution, 2nd edn. Springer, New York (2008)
27. Durrett, R., Schweinsberg, J.: Approximating selective sweeps. *Theor Popul Biol* **66**, 129–138 (2004)

28. Durrett, R., Schweinsberg, J.: A coalescent model for the effect of advantageous mutations on the genealogy of a population. *Stoch Proc Appl* **115**, 1628–1657 (2005)
29. Eldon, B.: Inference methods for multiple merger coalescents. In: P. Pontarotti (ed.) *Evolutionary Biology: convergent evolution, evolution of complex traits, concepts and methods*, pp. 347–371. Springer (2016)
30. Eldon, B., Birkner, M., Blath, J., Freund, F.: Can the site-frequency spectrum distinguish exponential population growth from multiple-merger coalescents. *Genetics* **199**, 841–856 (2015)
31. Eldon, B., Wakeley, J.: Coalescent processes when the distribution of offspring number among individuals is highly skewed. *Genetics* **172**, 2621–2633 (2006)
32. Eldon, B., Wakeley, J.: Linkage disequilibrium under skewed offspring distribution among individuals in a population. *Genetics* **178**, 1517–1532 (2008)
33. Etheridge, A.: *Some Mathematical Models from Population Genetics*. Springer Berlin Heidelberg (2011). DOI 10.1007/978-3-642-16632-7. URL <http://dx.doi.org/10.1007/978-3-642-16632-7>
34. Etheridge, A., Griffiths, R.: A coalescent dual process in a Moran model with genic selection. *Theor Popul Biol* **75**, 320–330 (2009)
35. Etheridge, A.M., Griffiths, R.C., Taylor, J.E.: A coalescent dual process in a Moran model with genic selection, and the Lambda coalescent limit. *Theor Popul Biol* **78**, 77–92 (2010)
36. Ewens, W.J.: *Mathematical population genetics 1: theoretical introduction*, vol. 27. Springer Science & Business Media (2012)
37. Freund, F., Möhle, M.: On the size of the block of 1 for  $\Xi$ -coalescents with dust. *ArXiv e-prints* (2017)
38. Freund, F., Siri-Jégousse, A.: Minimal clade size in the bolthausen-sznitman coalescent. *Journal of Applied Probability* **51**(3), 657–668 (2014)
39. Goldschmidt, C., Martin, J.B.: Random recursive trees and the bolthausen-sznitman coalescent. *Electron. J. Probab* **10**(21), 718–745 (2005)
40. Griffiths, R.C., Tavaré, S.: Monte carlo inference methods in population genetics. *Mathematical and computer modelling* **23**(8-9), 141–158 (1996)
41. Griffiths, R.C., Tavaré, S.: The age of a mutation in a general coalescent tree. *Comm Statistic Stoch Models* **14**, 273–295 (1998)
42. Griswold, C.K., Baker, A.J.: Time to the most recent common ancestor and divergence times of populations of common chaffinches (*Fringilla coelebs*) in Europe and North Africa: insights into Pleistocene refugia and current levels of migration. *Evolution* **56**(1), 143–153 (2002)
43. Halldórsdóttir, K., Árnason, E.: Whole-genome sequencing uncovers cryptic and hybrid species among Atlantic and Pacific cod-fish (2015). DOI 10.1101/034926. [Http://dx.doi.org/10.1101/034926](http://dx.doi.org/10.1101/034926)
44. Hedgecock, D.: Does variance in reproductive success limit effective population sizes of marine organisms? In: A. Beaumont (ed.) *Genetics and evolution of Aquatic Organisms*, pp. 1222–1344. Chapman and Hall, London (1994)
45. Hedgecock, D., Pudovkin, A.I.: Sweepstakes reproductive success in highly fecund marine fish and shellfish: a review and commentary. *Bull Marine Science* **87**, 971–1002 (2011)
46. Hedrick, P.: Large variance in reproductive success and the  $N_e/N$  ratio. *Evolution* **59**(7), 1596 (2005). DOI 10.1554/05-009
47. Hénard, O.: The fixation line in the  $\Lambda$ -coalescent. *The Annals of Applied Probability* **25**(5), 3007–3032 (2015)
48. Herriger, P., Möhle, M.: Conditions for exchangeable coalescents to come down from infinity. *Alea* **9**(2), 637–665 (2012)
49. Hird, S., Kubatko, L., Carstens, B.: Rapid and accurate species tree estimation for phylogeographic investigations using replicated subsampling. *Molecular Phylogenetics and Evolution* **57**(2), 888–898 (2010)
50. Hovmøller, M.S., Sørensen, C.K., Walter, S., Justesen, A.F.: Diversity of *Puccinia striiformis* on cereals and grasses. *Annual review of phytopathology* **49**, 197–217 (2011)
51. Hudson, R.R.: Properties of a neutral allele model with intragenic recombination. *Theor Popul Biol* **23**, 183–201 (1983)
52. Huillet, T., Möhle, M.: On the extended Moran model and its relation to coalescents with multiple collisions. *Theor Popul Biol* **87**, 5–14 (2013)
53. Kaj, I., Krone, S.M.: The coalescent process in a population with stochastically varying size. *Journal of Applied Probability* **40**(01), 33–48 (2003)

54. King, L., Wakeley, J.: Empirical bayes estimation of coalescence times from nucleotide sequence data. *Genetics* **204**(1), 249–257 (2016). DOI 10.1534/genetics.115.185751
55. Kingman, J.F.C.: The coalescent. *Stoch Proc Appl* **13**, 235–248 (1982)
56. Kingman, J.F.C.: Exchangeability and the evolution of large populations. In: G. Koch, F. Spizzichino (eds.) *Exchangeability in Probability and Statistics*, pp. 97–112. North-Holland, Amsterdam (1982)
57. Kingman, J.F.C.: On the genealogy of large populations. *J App Probab* **19A**, 27–43 (1982)
58. Li, G., Hedgecock, D.: Genetic heterogeneity, detected by PCR-SSCP, among samples of larval Pacific oysters (*Crassostrea gigas*) supports the hypothesis of large variance in reproductive success. *Can. J. Fish. Aquat. Sci.* **55**(4), 1025–1033 (1998). DOI 10.1139/f97-312
59. May, A.W.: Fecundity of Atlantic cod. *J Fish Res Brd Can* **24**, 1531–1551 (1967)
60. Möhle, M.: Robustness results for the coalescent. *Journal of Applied Probability* **35**(02), 438–447 (1998)
61. Möhle, M.: On sampling distributions for coalescent processes with simultaneous multiple collisions. *Bernoulli* **12**(1), 35–53 (2006)
62. Möhle, M.: Coalescent processes derived from some compound Poisson population models. *Elect Comm Probab* **16**, 567–582 (2011)
63. Möhle, M., Sagitov, S.: A classification of coalescent processes for haploid exchangeable population models. *Ann Probab* **29**, 1547–1562 (2001)
64. Möhle, M., Sagitov, S.: Coalescent patterns in diploid exchangeable population models. *J Math Biol* **47**, 337–352 (2003)
65. Neher, R.A., Hallatschek, O.: Genealogies of rapidly adapting populations. *Proceedings of the National Academy of Sciences* **110**(2), 437–442 (2013)
66. Niwa, H.S., Nashida, K., Yanagimoto, T.: Reproductive skew in Japanese sardine inferred from DNA sequences. *ICES Journal of Marine Science: Journal du Conseil* **73**(9), 2181–2189 (2016). DOI 10.1093/icesjms/fsw070. URL <http://dx.doi.org/10.1093/icesjms/fsw070>
67. Oosthuizen, E., Daan, N.: Egg fecundity and maturity of North Sea cod, *Gadus morhua*. *Netherlands Journal of Sea Research* **8**(4), 378–397 (1974)
68. Pettengill, J.B.: The time to most recent common ancestor does not (usually) approximate the date of divergence. *PloS one* **10**(8), e0128407 (2015)
69. Pitman, J.: Coalescents with multiple collisions. *Ann Probab* **27**, 1870–1902 (1999)
70. Sagitov, S.: The general coalescent with asynchronous mergers of ancestral lines. *J Appl Probab* **36**, 1116–1125 (1999)
71. Sagitov, S.: Convergence to the coalescent with simultaneous mergers. *J Appl Probab* **40**, 839–854 (2003)
72. Sargsyan, O., Wakeley, J.: A coalescent process with simultaneous multiple mergers for approximating the gene genealogies of many marine organisms. *Theor Pop Biol* **74**, 104–114 (2008)
73. Saunders, I.W., Tavaré, S., Watterson, G.A.: On the genealogy of nested subsamples from a haploid population. *Advances in Applied Probability* **16**(3), 471 (1984). DOI 10.2307/1427285
74. Schweinsberg, J.: Rigorous results for a population model with selection II: genealogy of the population. *ArXiv:1507.00394*
75. Schweinsberg, J.: Coalescents with simultaneous multiple collisions. *Electron J Probab* **5**, 1–50 (2000)
76. Schweinsberg, J.: Coalescents with simultaneous multiple collisions. *Electronic Journal of Probability* **5**, 1–50 (2000)
77. Schweinsberg, J.: A necessary and sufficient condition for the-coalescent to come down from the infinity. *Electronic Communications in Probability* [electronic only] **5**, 1–11 (2000)
78. Schweinsberg, J.: Coalescent processes obtained from supercritical Galton-Watson processes. *Stoch Proc Appl* **106**, 107–139 (2003)
79. Simon, M., Cordo, C.: Inheritance of partial resistance to *Septoria tritici* in wheat (*Triticum aestivum*): limitation of pycnidia and spore production. *Agronomie* **17**(6-7), 343–347 (1997)
80. Slack, R.: A branching process with mean one and possibly infinite variance. *Probability Theory and Related Fields* **9**(2), 139–145 (1968)
81. Spouge, J.L.: Within a sample from a population, the distribution of the number of descendants of a subsample's most recent common ancestor. *Theoretical population biology* **92**, 51–54 (2014)



- 1021 82. Tajima, F.: Evolutionary relationships of DNA sequences in finite populations. *Genetics*  
1022 **105**, 437–460 (1983)
- 1023 83. Timm, A., Yin, J.: Kinetics of virus production from single cells. *Virology* **424**(1), 11–17  
1024 (2012)
- 1025 84. Wakeley, J.: *Coalescent theory*. Roberts & Co (2007)
- 1026 85. Wakeley, J., Takahashi, T.: Gene genealogies when the sample size exceeds the effective  
1027 size of the population. *Mol Biol Evol* **20**, 208–213 (2003)
- 1028 86. Waples, R.S.: Tiny estimates of the  $N_e/N$  ratio in marine fishes: Are they real? *Journal of*  
1029 *Fish Biology* **89**(6), 2479–2504 (2016). DOI 10.1111/jfb.13143
- 1030 87. Wiuf, C., Donnelly, P.: Conditional genealogies and the age of a neutral mutant. *Theoretical*  
1031 *Population Biology* **56**(2), 183 – 201 (1999). DOI <http://dx.doi.org/10.1006/tpbi.1998.1411>.  
1032 URL <http://www.sciencedirect.com/science/article/pii/S0040580998914113>
- 1033 88. Zhou, J., Teo, Y.Y.: Estimating time to the most recent common ancestor (tmrca): compari-  
1034 son and application of eight methods. *European Journal of Human Genetics* (2015)

## 1035 A1 Coalescent processes

1036 To keep our presentation self-contained a precise definition of the coalescent pro-  
1037 cesses we will need will now be given. We follow the description of [19]. A coa-  
1038 lescent process  $\Pi$  is a continuous-time Markov chain on the partitions of  $\mathbb{N}$ . Let  
1039  $\Pi^{(n)}$  denote the restriction to  $[n]$ , and write  $\mathcal{P}_n$  for the space of partitions of  $[n]$ .  
1040 A partition  $\pi = \{\pi_1, \dots, \pi_{\#\pi}\} \in \mathcal{P}_n$  has  $\#\pi$  blocks which are disjoint subsets of  
1041  $[n]$ . We assume the blocks  $\pi_i$  are ordered by their smallest element; therefore we  
1042 always have  $1 \in \pi_1$ . In general a merging event can involve  $r$  distinct groups of  
1043 blocks merging simultaneously. We write  $\underline{k} = (k_1, \dots, k_r)$  where  $k_i \geq 2$  denotes the  
1044 number of blocks merging in group  $i$ . Here  $r \in [\lfloor \#\pi/2 \rfloor]$ ,  $k_1 + \dots + k_r \in [\#\pi]_2$  and  
1045  $i_1^{(a)}, \dots, i_{k_a}^{(a)}$  will denote the indices of the blocks in the  $a$ th group. By  $\pi' \prec_{\#\pi, \underline{k}} \pi$   
1046 we denote a transition from  $\pi$  to  $\pi' = A \cup B$  where

$$\begin{aligned} A &= \left\{ \pi_\ell : \ell \in [\#\pi], \ell \notin \bigcup_{a=1}^r \{i_1^{(a)}, \dots, i_{k_a}^{(a)}\} \right\}, \\ B &= \bigcup_{b=1}^r \left\{ \pi_{i_1^{(b)}}, \dots, \pi_{i_{k_b}^{(b)}} \right\}. \end{aligned} \quad (\text{A29})$$

1047 In (A29), set  $A$  (possibly empty) contains the blocks not involved in a merger,  
1048 and  $B$  lists the blocks involved in each of the  $r$  mergers. By  $\pi' \prec_{\#\pi, k} \pi$  we denote  
1049 the transition in a  $\Lambda$ -coalescent where  $k \in [\#\pi]_2$  merge in a single merger and  
1050  $\pi'$  is given as in (A29) with  $r = 1$ ; ie. only one group of blocks merges in each  
1051 transition. By  $\pi' \prec_{\#\pi} \pi$  we denote a transition in the Kingman-coalescent where  
1052  $r = 1$  and 2 blocks merge in each transition.

1053 Now that we have specified the possible transitions, we can state the rates of  
1054 the transitions. Let  $\Delta$  denote the infinite simplex  $\Delta = \{(x_1, x_2, \dots) : x_1 \geq x_2 \geq$   
1055  $\dots \geq 0, \sum_i x_i \leq 1\}$ ; let  $\mathbf{x}$  denote an element of  $\Delta$ . Define the functions  $f(\mathbf{x}; \#\pi, \underline{k})$



1056 and  $g(\mathbf{x}; \# \pi, \underline{k})$  on  $\Delta_0 := \Delta \setminus \{(0, 0, \dots)\}$  where  $(\prod_{m=1}^0 x_{i_r+m} := 1)$ , by

$$\begin{aligned} f(\mathbf{x}; \# \pi, \underline{k}) &= \frac{1}{\sum_j x_j^2} \sum_{\ell=0}^s \sum_{i_1 \neq \dots \neq i_{r+\ell}} \binom{s}{\ell} x_{i_1}^{k_1} \cdots x_{i_r}^{k_r} \prod_{m=1}^{\ell} x_{i_{r+m}} \left(1 - \sum_j x_j\right)^{s-\ell}, \\ g(\mathbf{x}; n) &= \frac{1 - \sum_{\ell=0}^n \sum_{i_1 \neq \dots \neq i_{\ell}} \binom{n}{\ell} x_{i_1} \cdots x_{i_{\ell}} (1 - \sum_j x_j)^{n-\ell}}{\sum_j x_j^2}. \end{aligned} \quad (\text{A30})$$

1057 where  $x_{i_0} := 1$ . Write  $\Xi_0$  for a finite measure on  $\Delta_0$ , and define [75], for some  
1058  $a \geq 0$ ,

$$\begin{aligned} \lambda_{n, \underline{k}} &:= \int_{\Delta_0} f(\mathbf{x}, n, \underline{k}) \Xi_0 d\mathbf{x} + a \mathbb{1}_{(r=1, k_1=2)}, \\ \lambda_n &:= \int_{\Delta_0} g(\mathbf{x}, n) \Xi_0 d\mathbf{x} + a \binom{n}{2}. \end{aligned} \quad (\text{A31})$$

1059 A  $\Xi$ -coalescent [75] is a continuous-time  $\mathcal{P}_n$ -valued Markov chain with tran-  
1060 sitions  $q_{\pi, \pi'}$  given by, where  $\lambda_{n, \underline{k}}$  and  $\lambda_n$  are given in (A31),

$$q_{\pi, \pi'} = \begin{cases} \lambda_{n, \underline{k}} & \text{if } \pi' \prec_{\# \pi, \underline{k}} \pi, \# \pi = n, \\ -\lambda_n & \text{if } \pi' = \pi \text{ and } n = \# \pi, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{A32})$$

1061 A  $\Lambda$ -coalescent [24, 69, 70] is a specific case of a  $\Xi$ -coalescent where we restrict  
1062 to the subset  $\Delta_0 := \Delta_0 \cap \{(x_1, x_2, \dots) : x_1 \in (0, 1], x_{1+i} = 0 \forall i \in \mathbb{N}\}$  [75] and the  
1063 transition rate of  $\pi' \prec_{\# \pi, k} \pi$  becomes, where  $\# \pi = n, 2 \leq k \leq n$ ,

$$\lambda_{n, k} = \int_{\Delta_0} x^k (1-x)^{n-k} x^{-2} \Xi_0 + a \mathbb{1}_{(k=2)}. \quad (\text{A33})$$

1064 The Kingman-coalescent is obviously obtained in the case  $\Xi_0(\Delta_0) = 0$  and  $a = 1$ .

1065 When we refer to a  $\Lambda$ -coalescent we will refer to the measure  $\Lambda = \Xi_0 + a \delta_0$   
1066 with  $\Xi_0$  restricted to  $\Delta_0$ . For  $\Lambda$  a finite measure on  $[0, 1]$  one can also represent the  
1067 coalescent rate  $\lambda_{n, k}$  of a  $\Lambda$ -coalescent as

$$\lambda_{n, k} = \int_0^1 x^{k-2} (1-x)^{n-k} \Lambda(dx), \quad 2 \leq k \leq n. \quad (\text{A34})$$

1068 The total rate of  $k$ -mergers in a  $\Lambda$ -coalescent is given by

$$\lambda_k(n) = \binom{n}{k} \int_0^1 x^k (1-x)^{n-k} x^{-2} d\Lambda(x), \quad 2 \leq k \leq n; \quad (\text{A35})$$

1069 and the total rate of mergers given  $n \geq 2$  active blocks is

$$\lambda(n) = \lambda_2(n) + \cdots + \lambda_n(n). \quad (\text{A36})$$

An important example of a  $\Lambda$ -coalescent is the Beta( $2 - \alpha, \alpha$ )-coalescent [78] where the  $\Lambda$  measure is associated with the beta density, where  $B(\cdot, \cdot)$  is the beta function,

$$\Lambda(dx) = \frac{x^{1-\alpha}(1-x)^{\alpha-1}}{B(2-\alpha, \alpha)} dx, \quad 1 \leq \alpha < 2. \quad (\text{A37})$$

The total rate of a  $k$ -merger (A35) is then given by, for  $2 \leq k \leq n$ ,

$$\lambda_k(n) = \binom{n}{k} \frac{B(k-\alpha, n-k+\alpha)}{B(2-\alpha, \alpha)}, \quad 1 \leq \alpha < 2. \quad (\text{A38})$$

When  $\alpha = 1$  the Beta( $2 - \alpha, \alpha$ )-coalescent corresponds to the Bolthausen-Sznitman coalescent [20, 39].

## A2 Goldschmidt and Martin's construction of the Bolthausen-Sznitman $n$ -coalescent

From [39], we recall the construction of the Bolthausen-Sznitman  $n$ -coalescent by cutting the edges of a random recursive tree. Let  $\mathbb{T}_n$  be a random recursive tree with  $n$  nodes. We can construct  $\mathbb{T}_n$  sequentially as follows

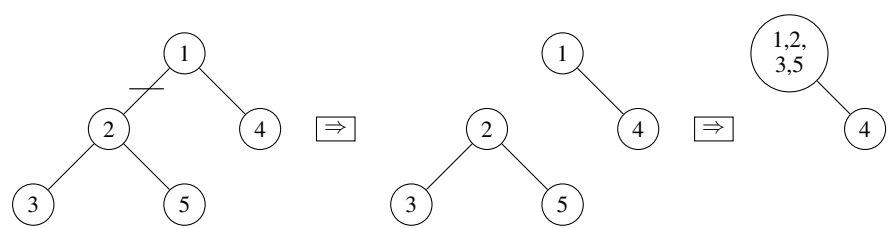
- (i) Start with node 1 (the root) and no edges,
- (ii) If  $i < n$  nodes are present, add node  $i + 1$  and one edge connecting it to a node in  $[i]$  picked uniformly,
- (iii) stop if  $n$  nodes are present.

The object  $\mathbb{T}_n$  is a labelled tree, each node has a single label.

We consider a realisation of  $\mathbb{T}_n$  and transform this tree over time into labelled trees with fewer nodes with nodes amassing multiple labels.

- (i) Each edge of  $\mathbb{T}_n$  is linked to an exponential clock. Clocks are i.i.d.  $\text{Exp}(1)$ -distributed.
- (ii) We wait for the first clock to ring. At this time, we cut/remove the edge whose clock rang first. The tree is thus split in two trees from which one includes label 1. We denote the tree with label 1 by  $\mathbb{T}^{(1)}$ , the other tree by  $\mathbb{T}^{(2)}$ . Let  $e_1$  be the node of  $\mathbb{T}^{(1)}$  that was connected by the removed edge
- (iii) All labels of  $\mathbb{T}^{(2)}$  are added to the set of labels of  $e_1$ . Remove  $\mathbb{T}^{(2)}$  including its clocks.
- (iv) Repeat from (ii), using  $\mathbb{T}^{(1)}$  labelled as in (iii) with the (remaining) clocks from (i). Stop when  $\mathbb{T}^{(1)}$  in step (iii) consists of only a single node and no edges.
- (v) For any time  $t$ , label sets at the nodes of  $\mathbb{T}^{(1)}$  ( $\mathbb{T}_n$  before the first clock has rang) give a partition  $\Pi_t^{(n)}$  of  $[n]$ .  $(\Pi_t^{(n)})_{t \geq 0}$  is a Bolthausen-Sznitman  $n$ -coalescent (set  $\Pi_t^{(n)} = [n]$  if  $t$  is bigger than the time at which we stopped the cutting procedure).

Figure A2 shows an illustration of steps (i)-(iii) for a realisation of  $\mathbb{T}_5$ .



**Fig. 6** Example for the first cutting and relabelling step (ii), (iii) for the construction from [39].

**Fig. 7** Graphs of  $p_{n,m}^{(B)}$  (see Eq. (5)) as a function of  $\alpha$  for  $(n, m) = (10^2, 10^1)$  (circles);  $(10^3, 10^1)$  (—);  $(10^3, 10^2)$  (+). The corresponding results for the Kingman-coalescent,  $p_{n,m}^{(Kingman)}$  (A) and  $p_m^{(Kingman)}$  (B) are shown as lines.

