

Can the Site-Frequency Spectrum Distinguish Exponential Population Growth from Multiple-Merger Coalescents?

Bjarki Eldon,^{*,1} Matthias Birkner,[†] Jochen Blath,^{*} and Fabian Freund^{*}

^{*}TU Berlin, Institut für Mathematik, 10623 Berlin, Germany, [†]JGU Mainz, Institut für Mathematik, 55099 Mainz, Germany, and

[‡]University of Hohenheim, Institute of Plant Breeding, Seed Science, and Population Genetics, 70599 Stuttgart, Germany

ABSTRACT The ability of the site-frequency spectrum (SFS) to reflect the particularities of gene genealogies exhibiting multiple mergers of ancestral lines as opposed to those obtained in the presence of population growth is our focus. An excess of singletons is a well-known characteristic of both population growth and multiple mergers. Other aspects of the SFS, in particular, the weight of the right tail, are, however, affected in specific ways by the two model classes. Using an approximate likelihood method and minimum-distance statistics, our estimates of statistical power indicate that exponential and algebraic growth can indeed be distinguished from multiple-merger coalescents, even for moderate sample sizes, if the number of segregating sites is high enough. A normalized version of the SFS (nSFS) is also used as a summary statistic in an approximate Bayesian computation (ABC) approach. The results give further positive evidence as to the general eligibility of the SFS to distinguish between the different histories.

KEYWORDS coalescent; multiple mergers; population growth; approximate maximum likelihood test; approximate Bayesian computation; site-frequency spectrum

THE site-frequency spectrum (SFS) at a given locus is one of the most important and popular statistics based on genetic data sampled from a natural population. In combination with the postulation of the assumptions of the infinitely-many-sites mutation model (Watterson, 1975) and a suitable underlying coalescent framework, the SFS allows one to draw inferences about evolutionary parameters, such as coalescent parameters associated with multiple-merger coalescents or population-growth models.

The Kingman coalescent, developed by Kingman (1982a, b,c), Hudson (1983a,b), and Tajima (1983), describing the random ancestral relations among DNA sequences drawn from natural populations, is a prominent and widely used model from which one can make predictions about genetic diversity. Many quantities of interest, such as the expected values and covariances of the SFS associated with the Kingman coalescent, are easily computed thanks to results by Fu

(1995). The robustness of the Kingman coalescent is quite remarkable; indeed, a large number of genealogy models can be shown to have the Kingman coalescent or a variant thereof as their limit process (*cf.*, *e.g.*, Möhle 1998). A large volume of work is thus devoted to inference methods based on the Kingman coalescent [see, *e.g.*, Donnelly and Tavaré (1995), Hudson (1990), Nordborg (2001), Hein *et al.* (2005), and Wakeley (2007) for reviews].

However, many evolutionary histories can lead to significant deviations from the Kingman coalescent model. Such deviations can be detected using a variety of statistical tools, such as Tajima's *D* (Tajima 1989a), Fu and Li's *D* (Fu and Li 1993), and Fay and Wu's *H* (Fay and Wu 2000), which are all functions of the SFS. However, they do not always allow one to identify the actual evolutionary mechanisms leading to such deviations. Developing statistical tools that allow one to distinguish between different evolutionary histories is therefore of fundamental importance.

This work focuses on properties of the (folded and unfolded) SFS in the infinitely-many-sites model for three population histories: (1) classical Kingman coalescent, (2) population growth, in particular, exponential population growth, and (3) high fecundity coupled with skewed offspring distributions (HFSODs), resulting in gene genealogies being

Copyright © 2015 by the Genetics Society of America

doi: 10.1534/genetics.114.173807

Manuscript received December 16, 2014; accepted for publication January 6, 2015; published Early Online January 9, 2015.

Supporting information is available online at <http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.173807/-/DC1>

¹Corresponding author: TU Berlin, Institut für Mathematik, Straße des 17. Juni 136, 10623 Berlin, Germany. E-mail: eldon@math.tu-berlin.de

described by so-called Lambda-coalescents (Sagitov 1999; Pitman 1999; Donnelly and Kurtz 1999). Briefly, multiple-merger coalescents may be more appropriate for organisms exhibiting HFSODs than the Kingman coalescent (*cf.*, *e.g.*, Beckenbach 1994; Árnason 2004; Eldon and Wakeley 2006; Sargsyan and Wakeley 2008; Hedgcock and Pudovkin 2011) [see also the review by Tellier and Lemaire (2014)].

Recent population growth as well as multiple-merger coalescents may lead to an excess of singletons in the SFS compared with the classical Kingman coalescent-based SFS, which *e.g.*, contributes to shifting Tajima's D values (Tajima 1989b) to the negative. Indeed, Durrett and Schweinsberg (2005) proved that Tajima's D will be negative, at least for large sample size, under fairly general multiple-merger coalescents.

The associated genealogical trees are, however, qualitatively different. While moderate fluctuations in population size lead to a time change of the Kingman coalescent (Kaj and Krone 2003), multiple-merger coalescents by definition change the topology of the genealogical tree. There is thus hope that each demographic effect leaves specific signatures in the resulting SFS not only with respect to an excess of singletons but also, *e.g.*, with respect to its right tail.

Indeed, one observes that the Kingman coalescent will not be a good match to genetic data containing a large fraction of singleton polymorphisms (relative to the total number of polymorphisms) because of a lack of free (coalescent) parameters as opposed to multiple-merger and population-growth models, both of which can predict an excess of singletons. Encouragingly, multiple-merger and population-growth models exhibit noticeable differences in the bulk of the site-frequency spectrum, in particular, in the lumped tail (Figure 1; see also Figures 4 and S2 in Neher and Hallatschek 2013. (Neher and Hallatschek 2013). In Figure 1, the normalized expected spectrum $\varphi_i^{(n,\Pi)}$ [see Equation (2)] for a given coalescent Π , *i.e.*, the expected spectrum scaled by the expected total number of segregating sites, is compared for different multiple-merger coalescents [B = beta-coalescents (Schweinsberg, 2003) or D = Dirac coalescents (Eldon and Wakeley, 2006)] and exponential (E) and algebraic (A) growth models leading to time-changed Kingman coalescents for sample size (number of leaves) n as shown. Details for these coalescent models are given at the beginning of File S1. The first five classes (representing relative length of external branches, two-leaf branches, etc.) are shown, with classes from six onward collected together (labeled 5+). In Figure 1, the relative external branch lengths were matched between the different coalescent processes. Even though the relative external branch lengths and, by implication, the number of singletons relative to the total number of segregating sites can be matched between the different processes, the collapsed tail (group 5+ in Figure 1) differs noticeably between the multiple-merger coalescents and the growth models. One also observes that the parameters have been chosen to match $\varphi_1^{(n,\Pi)}$ for $\Pi \in \{A, D, E\}$ with $\varphi_1^{(n,B)}$ when $\alpha = 1$, where α is the coalescent parameter associated with B. Thus, $\varphi_1^{(n,B)}$ is

maximized for the given n (because $\alpha \in [1, 2]$), but $\varphi_1^{(n,D)}$, $\varphi_1^{(n,E)}$, and $\varphi_1^{(n,A)}$ all can increase by increasing the relevant parameters (ψ , β , or γ).

Matching the relative external branch lengths $\varphi_1^{(n,\Pi)}$ [see Equation 2] and observing how the rest of the normalized expected spectrum behaves, as illustrated in Figure 1, give hope that multiple-merger processes may be distinguished from (at least) particular population-growth models with adequate statistical power. In the limit of large n , for the Kingman coalescent, $\varphi_1^{(n,K)} = O[1/\log(n)]$.

Inference methods for distinguishing population growth from the usual Kingman coalescent have been studied extensively (see, *e.g.*, Tajima 1989a; Slatkin and Hudson 1991; Rogers and Harpending 1992; Kaj and Krone 2003; Sano and Tachida 2005). Simulation-based work includes Ramírez-Soriano *et al.* (2008), who considered the statistical power of several tests under population size increase and decrease and the impact of recombination. Ramos-Onsins and Rozas (2002) considered the statistical power of statistics based on the site-frequency spectrum to distinguish deterministic population growth from the Kingman coalescent. On the theoretical side, Myers *et al.* (2008), Bhaskar and Song (2014), and Kim *et al.* (2014) considered principal questions of identifiability of demographic histories. In particular, Bhaskar and Song (2014) showed theoretically that complete knowledge of the SFS for large sample sizes carries enough information to fully recover demographic history under mild assumptions on the possible fluctuations of the demography.

Detecting multiple-merger coalescents in populations deviating from the Kingman coalescent assumptions is a relatively new direction of research. Indeed, deriving inference methods based on multiple-merger coalescents has only just begun (Eldon and Wakeley 2006; Birkner and Blath 2008; Eldon 2011; Birkner *et al.* 2011, 2013a, b; Steinrücken *et al.* 2013; Rödelberger *et al.* 2014; Koskela *et al.* 2015). In particular, Birkner *et al.* (2013b) obtained recursions for the expected site-frequency spectrum associated with Lambda-coalescents. In this work, we address the issue of distinguishing multiple-merger coalescents from exponential population growth by proposing statistical tests based on the (normalized) SFS, estimating statistical power for interval hypotheses via simulation. Because we can only work with approximate likelihood functions and our methods, in particular, the so-called fixed- s method, can be sensitive to an (unknown) true coalescent mutation rate $\theta/2$, we complement our analysis by an approximate Bayesian computation approach (ABC) (Rubin 1984; Tavaré *et al.* 1997; Pritchard *et al.* 1999; Cucala and Marin 2013; Baragatti and Pudlo 2014).

Materials and Methods

Basic properties of the site-frequency spectrum

Consider a sample of n DNA sequences taken at a given genetic locus, and assume that we can distinguish between

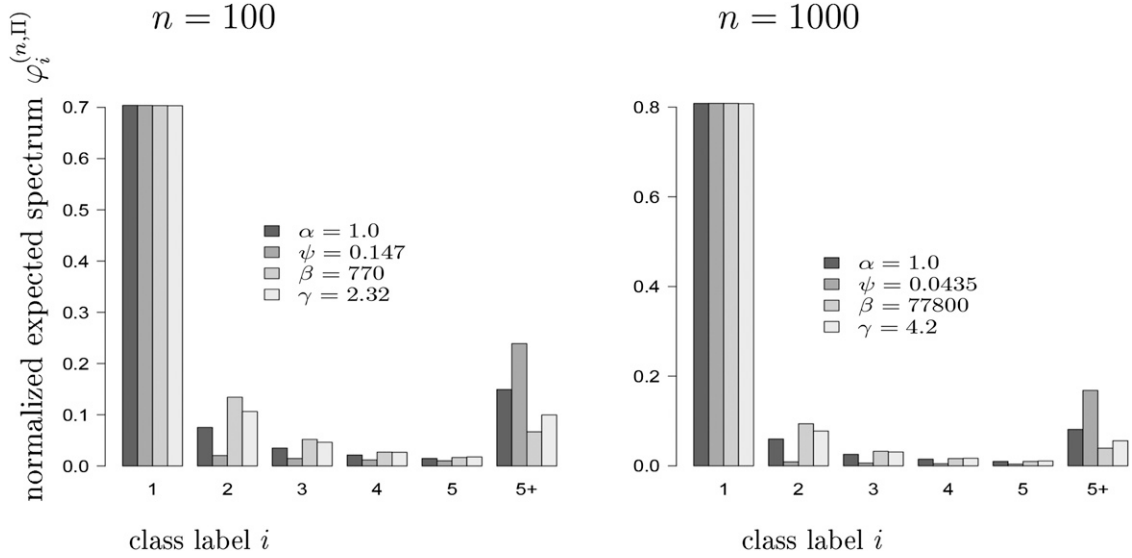


Figure 1 Matching $\varphi_i^{(n,\Pi)}$ [see Equation 2] for the different coalescent processes $\Pi \in \{A, B, D, E\}$ with number of leaves n as shown. Expected values were computed exactly. The processes and their associated parameters are algebraic growth (A, γ), beta($2 - \alpha$)-coalescent (B, α), Dirac coalescent (D, ψ), and exponential growth (E, β). The values with label 5+ represent the collapsed tail $\sum_{i>5} \varphi_i^{(n,\Pi)}$.

derived (new mutations) and ancestral states. For $n \in \mathbb{N}$, let $n := \{1, \dots, n\} \setminus \{1, \dots, n\}$. We denote by $\xi_i^{(n)}$ the total number of sites at which the mutant base appears $i \in [n - 1]$ times. Then

$$\underline{\xi}^{(n)} := \left(\xi_1^{(n)}, \dots, \xi_{n-1}^{(n)} \right)$$

is referred to as the *unfolded* site-frequency spectrum based on the n DNA sequences. If mutant and wild type cannot be distinguished, one often considers the *folded* spectrum $\underline{\eta}^{(n)} := (\eta_1^{(n)}, \dots, \eta_{\lfloor n/2 \rfloor}^{(n)})$, where ancestral and derived states are not distinguished and hence

$$\eta_i^{(n)} := \frac{\xi_i^{(n)} + \xi_{n-i}^{(n)}}{1 + \delta_{i,n-i}}, \quad 1 \leq i \leq \lfloor n/2 \rfloor$$

(Fu 1995), where $\delta_{i,j} = 1$ if $i = j$ and 0 otherwise. In this study, we will mostly be concerned with the unfolded site-frequency spectrum. Define $\zeta_i^{(n)} := \xi_i^{(n)} / |\xi^{(n)}|$, where $|\xi^{(n)}| := \xi_1^{(n)} + \dots + \xi_{n-1}^{(n)}$ denotes the total number of segregating sites. Thus, $\underline{\zeta}^{(n)} = (\zeta_1^{(n)}, \dots, \zeta_{n-1}^{(n)})$ is the “normalized” unfolded SFS (nSFS), with the convention that $\zeta_i^{(n)} = 0$ in the trivial case of complete absence of segregating sites ($|\xi^{(n)}| = 0$).

In order to compute expected values, variances, and covariances of the SFS, an explicit underlying probabilistic model is needed. In the following, we assume that the genealogy of a sample can be described by a coalescent process, more precisely by either (a time change of) the Kingman coalescent or a multiple-merger coalescent. In addition, the infinitely-many-sites mutation model (Watterson 1975) is assumed, and mutations are modeled by a Poisson process on the coalescent branches with rate $\theta/2$. With this parameterization, the expected number of segregating

sites in a sample of size 2, and hence the expected number of pairwise differences in a sample from the population, equals θ .

Closed-form expressions for the expected values and (co-)variances of $\underline{\xi}^{(n)}$ have been determined by Fu (1995) when associated with the Kingman coalescent. One can represent the expected values of $\underline{\xi}^{(n)}$ in a unified way using the results of Griffiths and Tavaré (1998), Kaj and Krone (2003), and Birkner *et al.* (2013b), which allow one to treat the expected values (and covariances) of the SFS for all coalescent models in question.

Let $\Pi^{(n)} = (\Pi_t^{(n)}, t \geq 0)$ be a (partition-valued exchangeable) coalescent process started from n leaves (partition blocks) corresponding to the random genealogy of a sample of size n . By discussing leaves rather than DNA sequences, we are emphasizing our viewpoint of the genealogy as a random graph, where the leaves are a particular kind of vertex. Our interest is in the topology of the genealogy and how it is reflected in the associated site-frequency spectrum.

If the initial number of leaves is not specified, we simply speak of Π . One may think of Π as the Kingman coalescent, but the point is that the following result will also stay true for externally time-changed Kingman coalescents as well as multiple-merger coalescents (a.k.a. Lambda- or Xi-coalescents in the mathematical literature) and even externally time-changed multiple-merger coalescents.

Given n and a coalescent model Π , let $(Y_t^{(n)})_{t \geq 0}$ be the block-counting process of the underlying coalescent $\Pi^{(n)}$ started from n lineages; *i.e.*, $Y_t^{(n)}$ gives the number of ancestral lines (blocks) present/active at (backward) time t . For $2 \leq k < n$, let $T_k^{(n)}$ be the random amount of time that $(Y_t^{(n)})_{t \geq 0}$ spends in state k . Given a coalescent $\Pi^{(n)}$ started from n (unlabeled) lineages, denote by $p^{(n),\Pi}[k, i]$ as the probability that *conditional* on the event that $Y_t^{(n)} = k$ for

some time point t a given one of the k blocks subtends exactly $i \in [n - 1]$ leaves. A general representation of $\mathbb{E}^{\Pi, \theta}[\xi_i^{(n)}]$ is then

$$\mathbb{E}^{\Pi, \theta}[\xi_i^{(n)}] = \frac{\theta}{2} \sum_{k=2}^{n-i+1} p^{(n), \Pi}[k, i] \cdot k \cdot \mathbb{E}^{\Pi}[T_k^{(n)}], \quad i \in [n - 1] \quad (1)$$

The normalized expected SFS $\varphi_i^{(n, \Pi)}$ for $i \in [n - 1]$ is defined as

$$\varphi_i^{(n, \Pi)} = \frac{\sum_{k=2}^{n-i+1} p^{(n), \Pi}[k, i] \cdot k \cdot \mathbb{E}^{\Pi}[T_k^{(n)}]}{\sum_{\ell=2}^n \ell \mathbb{E}^{\Pi}[T_{\ell}^{(n)}]} \quad (2)$$

where the denominator in Equation 2 is the expected total tree length when starting from n leaves. One can interpret the quantity $\varphi_i^{(n, \Pi)}$ as the probability that a mutation, under the infinitely-many-sites assumption and the coalescent model Π , with known ancestral types, appears i times in a sample of size n . Importantly, $\varphi_i^{(n, \Pi)}$ is not a function of the mutation rate, unlike $\mathbb{E}^{\Pi, \theta}[\xi_i^{(n)}]$. One also can view $\varphi_i^{(n, \Pi)}$ as a first-order approximation of the expected value $\mathbb{E}^{\Pi, \theta}[\zeta_i^{(n)}]$ of the nSFS.

As examples for Π , we will consider the classical Kingman coalescent (K), exponential (E) and algebraic (A) growth models, and the beta($2 - \alpha$, α) (B) and Dirac (D) multiple-merger coalescents, as shown in File S1. Simulations suggest that $\varphi_i^{(n, B)}$ is a good approximation of $\mathbb{E}^{B, \theta}[\zeta_i^{(n)}]$ when α is not too close to 1 and n and θ are not too small (Birkner *et al.* 2013b). Similar conclusions hold in the case of exponential and algebraic growth (results not shown).

One can use the recursive formulas obtained by Birkner *et al.* (2013b) to compute $\mathbb{E}^{\Pi, \theta}[\xi_i^{(n)}]$ associated with Lambda-coalescents. To compute $\varphi_i^{(n, \Pi)}$ associated with growth models, we use the results of Polanski and Kimmel (2003), whose recursions are given in File S1.

A comparison of the observed $\zeta_i^{(n)}$ (instead of $\xi_i^{(n)}$) with an expected value $\mathbb{E}^{\Pi, \theta}[\zeta_i^{(n)}]$ —obtained under a particular coalescent model Π —enables one to do inference without having to jointly estimate the mutation rate θ using, *e.g.*, a minimum-distance statistic. Indeed, it appears that under any coalescent model Π , $\mathbb{E}^{\Pi, \theta}[\zeta_i^{(n)}]$ is almost constant as a function of the mutation rate θ (unless θ is very small); we provide some evidence for this in Equation S20 in File S1. Unfortunately, there seems to be no explicit way of representing $\mathbb{E}^{\Pi, \theta}[\zeta_i^{(n)}]$ as a simple function of the associated coalescent parameters and sample size n . As mentioned earlier, one may instead work with $\varphi_i^{(n, \Pi)}$.

Time scales, segregating sites, and mutation rates

The choice of a multiple-merger coalescent model (*i.e.*, demographic history) Π and its underlying parameters strongly affects classical estimates of the coalescent mutation rate $\theta/2$ (*i.e.*, the Poisson rate at which mutations appear on coalescent branches). Assume without loss of generality for all multiple-merger coalescents in question that the underlying

coalescent measure Λ is always a probability measure: this normalization fixes the coalescent time unit as the expected time to the most recent common ancestor of two individuals sampled uniformly from the population.

Given an observed number of segregating sites S in a sample of size n , a common estimate $\hat{\theta}^{\Pi}$ of the scaled mutation rate θ associated with coalescent model Π is the Watterson estimate, *i.e.*,

$$\hat{\theta}^{\Pi} := \frac{2S}{\mathbb{E}^{\Pi}[B^{(n)}]} \quad (3)$$

where $\mathbb{E}^{\Pi}[B^{(n)}]$ is the expectation of the total tree length $B^{(n)}$ of an (n -)coalescent model Π . One can, of course, also estimate θ as a (different) linear combination of the site-frequency spectrum [*cf.* Achaz (2009) in the case of the Kingman coalescent]. Using the recursions for $\mathbb{E}^{\Pi, \theta}[\xi_i^{(n)}]$ obtained by Birkner *et al.* (2013b), one also can estimate θ using either Equation 3 or a linear combination of the expected SFS in the case of a Lambda-coalescent.

Given an estimate $\hat{\theta}^{\Pi}$ and knowledge of the mutation/substitution rate $\hat{\mu}$ per year at the locus under consideration, one can find a real-time embedding of the coalescent history via the approximate identity

$$\text{Coalescent time unit} \times \frac{\hat{\theta}^{\Pi}}{2} \approx \text{year} \times \hat{\mu} \quad (4)$$

(see Steinrücken *et al.* 2013, Section 4.2), which, of course, depends on Π .

If one has additional information on the specific reproductive mechanisms of an approximating population model, this can even enable one to estimate the model census population size. For example, given a Cannings population model (Cannings, 1974, 1975) of fixed size N , let c_N be the probability that two gene copies, drawn uniformly at random and without replacement from a population of size N , derive from a common parental gene copy in the previous generation. While for the usual haploid Wright-Fisher model $c_N = 1/N$, in a class of population models studied by Schweinsberg (2003) leading to the beta($2 - \alpha$, α)-coalescent, c_N is proportional to $1/N^{\alpha-1}$ for $\alpha \in (1, 2]$ (but note that the proportionality constant depends on finer details of the particular model). By a limit theorem for Cannings models by Möhle and Sagitov (2001), one coalescent time unit corresponds to approximately $1/c_N$ generations in the original model with population size N . Thus the mutation rate $\hat{\mu}$ at the locus under consideration per individual per generation must be scaled with $1/c_N$ [as noted *e.g.*, in Eldon and Wakeley (2006)], and the relation between $\hat{\mu}$, the coalescent mutation rate $\theta^{\Pi}/2$, and c_N is then given by the (approximate) identity

$$c_N \approx \frac{2\hat{\mu}}{\theta^{\Pi}} \quad (5)$$

In particular, if the Cannings model class (and thus c_N as a function of N) is known, N can be estimated via Equation 5.

In this context, it is important to note that different population models on very different time scales still can have the Kingman coalescent as their ancestral limit process; two examples are the Wright-Fisher [$O(c_N^{-1}) = N$] and Moran [$O(c_N^{-1}) = N^2$] models. This is certainly also the case for multiple-merger coalescents. In particular, c_N is *a priori* not a function of the limiting coalescent model (this appears to be a rather frequent misperception).

Distinguishing real and coalescent time scales is important because nonlinear scaling otherwise may easily lead to confusion: for example, the expected total tree length $\mathbb{E}^B[B^{(n)}]$ (measured in coalescent time units) *decreases* as a function of $\alpha \in (1, 2]$, while the corresponding quantity (measured in real-time generations) $\mathbb{E}^B[B^{(n)}]/c_N$ *increases* in $\alpha \in (1, 2]$ (cf. Figure S1 in File S1).

The time scaling applied to a classical Wright-Fisher model with *fluctuating* population size [as in Kaj and Krone (2003)] in order to obtain a (time-changed) Kingman coalescent is shown in particular in Equation (S10) in File S1. Again, the estimate (Equation 3) of θ depends on the growth model and growth parameter.

Approximate likelihood-ratio tests for the SFS

Our aim is to construct a statistical test to distinguish among the model classes E, A, D, and B (which intersect exactly in the Kingman coalescent K). In order to distinguish, say, E from B based on an observed site-frequency spectrum $\underline{\xi}^{(n)}$ with sample size n and segregating sites $S = |\underline{\xi}^{(n)}|$, a natural approach is to construct a likelihood-ratio test.

Recall that we think of our observed spectrum as a realization of a coalescent tree with n leaves obtained from a coalescent model Π , with mutations distributed on the tree according to an independent Poisson ($\theta/2$) process. For each model Π from classes {E, A, D, B}, the coalescent will be uniquely determined by a single coalescent parameter $\beta \in [0, \infty)$ (for E), $\gamma \in [0, \infty)$ (for A), $\psi \in [0, 1]$ (for D), and $\alpha \in [1, 2]$ (for B). [Note that the beta-coalescent is well defined for $\alpha \in (0, 2]$, but we restrict to a smaller parameter range corresponding to the population model in Schweinsberg (2003).]

Suppose that our null hypothesis H_0 is the presence of recent exponential population growth (E) with (unknown) parameter $\beta \in [0, \infty)$, and we wish to test it against the alternative H_1 hypothesis of a multiple-merger coalescent, say, the beta($2 - \alpha, \alpha$)-coalescent (B) for (unknown) $\alpha \in [1, 2]$, where $\beta = 0$ and $\alpha = 2$ correspond to the Kingman coalescent. In this framework, the coalescent mutation rate θ is not directly observable but plays the role of a nuisance parameter. In particular, it is the interplay of the coalescent model Π and the mutation-rate parameter θ that governs the law of the observed number of mutations (see the discussion in the preceding section). To take θ explicitly into account, one could test

$$H_0 : (\Pi, \theta) \in \Theta^E := \{(\beta, \theta) : \beta \in [0, \infty), \theta \in (0, \infty)\}$$

(exponential growth) against

$$H_1 : (\Pi, \theta) \in \Theta^B := \{(\alpha, \theta) : \alpha \in [1, 2], \theta \in (0, \infty)\}$$

if the beta-coalescent family is the alternative [by slight abuse of notation, we identify the coalescent model Π with the corresponding coalescent parameter β (resp. α) in each model class when appropriate]. The underlying parameter ranges are two-dimensional, and although an explicit likelihood-ratio test based on methods described in Simonsen *et al.* (1995) can be constructed, it will likely pose computational challenges.

Instead, given an observed number of segregating sites $S = s$, we simplify our framework by employing the fixed- s method discussed, e.g., in Depaulis and Veuille (1998) and Ramos-Onsins and Rozas (2002). Here we treat the observed number of segregating sites as a *fixed parameter* $s \in \mathbb{N}$, not as (observation of a) random variable S . We will thus obtain the empirical distributional quantities of our test by Monte Carlo simulations, placing uniformly at random s mutations along the branches of the simulated tree.

The fixed- s method is different from generating samples for a given θ by conditioning on $S = s$, yet the fixed- s method usually leads to reasonable tests when the true θ is close to the Watterson estimate $\hat{\theta}(\Pi, s)$ based on s (Wall and Hudson 2001). However, it can lead to substantial deviations from the conditional distribution if θ is extreme [see, e.g., Markovtsova *et al.* (2001), who show that for a test in a related framework, the probability of rejection can be substantially different from 5%]. Regarding this caveat, Depaulis *et al.* (2001) took a Bayesian viewpoint and showed that the values of θ that lead to unreliable tests are highly unlikely given s . We address this issue by using rejection sampling to check the robustness of the fixed- s method against varying θ (see File S1). Our analysis is complemented with an ABC approach using the normalized frequency spectrum $\underline{z}^{(n)}$, which should be insensitive to the actual value of θ as long as θ is not too small [cf. Equation (S20) in File S1].

By fixing $S = s$ and treating it as a parameter of our test, we may consider the new pair of hypotheses

$$H_0^s : \Pi \in \Theta_s^E := \{\beta : \beta \in [0, \infty)\}$$

and

$$H_1^s : \Pi \in \Theta_s^B := \{\alpha : \alpha \in [1, 2]\}$$

Define a likelihood function $L(\Pi, \underline{k}^{(n)}, s)$ for the observed frequency spectrum $\underline{k}^{(n)} = (k_1^{(n)}, \dots, k_{n-1}^{(n)})$ with fixed $|\underline{k}^{(n)}| = s$ under the coalescent model $\Pi \in \Theta_s^E$ (resp. $\Pi \in \Theta_s^B$) by

$$\begin{aligned} L(\Pi, \underline{k}^{(n)}, s) &= \mathbb{P}^{\Pi, s} \left\{ \xi_i^{(n)} = k_i^{(n)}, i \in [n-1] \right\} \\ &= \mathbb{E}^{\Pi} \left[\frac{s!}{k_1^{(n)}! \cdots k_{n-1}^{(n)}!} \prod_{i=1}^{n-1} \binom{B_i^{(n)}}{B^{(n)}}^{k_i^{(n)}} \right] \end{aligned} \quad (6)$$

where $B_i^{(n)}$ are the random lengths of branches subtending $i \in [n - 1]$ leaves, and $B^{(n)}$ is the total branch length of the coalescent under Π . The fixed- s paradigm thus leads to a mixture of multinomial distributions where the parameters are given by the respective relative branch lengths. The hope is that the location of the maximum of $L(\Pi, \underline{k}^{(n)}, s)$ is typically not far from the location of the corresponding coordinate of the maximizer in the full two-dimensional explicit- θ model, in which one can additionally maximize over all $\theta \in [0, \infty)$.

Now we can construct a likelihood-ratio test based on $L(\Pi, \underline{k}^{(n)}, s)$ via the likelihood-ratio function

$$\varrho_{(E,B;s)}(\underline{\xi}^{(n)}) := \frac{\sup\{L(\Pi, \underline{k}^{(n)}, s), \Pi \in \Theta_s^E\}}{\sup\{L(\Pi, \underline{k}^{(n)}, s), \Pi \in \Theta_s^B\}} \quad (7)$$

Given a significance level $a \in (0, 1)$ (say, $a = 0.05$), let $\varrho_{(E,B;s)}^*(a)$ be the a -quantile of $\varrho_{(E,B;s)}(\underline{\xi}^{(n)})$ under E , chosen as the largest values so that

$$\sup_{\Pi \in \Theta_s^E} \mathbb{P}^{\Pi,s}\{\varrho_{(E,B;s)}(\underline{\xi}^{(n)}) \leq \varrho_{(E,B;s)}^*(a)\} \leq a \quad (8)$$

The decision rule that constitutes the fixed- s likelihood-ratio test, given s and sample size n , is

$$\text{Reject } H_0^s \Leftrightarrow \varrho_{(E,B;s)}(\underline{\xi}^{(n)}) \leq \varrho_{(E,B;s)}^*(a)$$

This formulation is free of the nuisance parameter θ . To assess the justification for the fixed- s assumption, we investigate how close this is to the corresponding quantiles for different values of θ , including the Watterson estimator

$$\hat{\theta} = \hat{\theta}(\Pi; s) = \frac{2s}{\mathbb{E}^\Pi[B^{(n)}]} \quad (9)$$

[cf. (Equation 3)], for selected choices of Π . The agreement appears reasonably good and seems to increase with sample size (see [File S1](#)).

The corresponding power function of the test, *i.e.*, the probability of rejecting a false null hypothesis, is given by

$$G_{(E,B;s)}(\Pi) = \mathbb{P}^\Pi\{\varrho_{(E,B;s)}(\underline{\xi}^{(n)}) \leq \varrho_{(E,B;s)}^*(a, S)\}, \quad \Pi \in \Theta_s^B \quad (10)$$

The likelihood (Equation 6) cannot be represented as a simple formula involving the coalescent parameters; one can approximate (Equation 6) via a Monte Carlo approach, but this is computationally expensive. An approximation is

$$L(\Pi, \underline{k}^{(n)}, s) \approx \frac{s!}{k_1^{(n)}! \cdots k_{n-1}^{(n)}!} \prod_{i=1}^{n-1} (\varphi_i^{(n,\Pi)})^{k_i^{(n)}} \quad (11)$$

where we replaced the random quantities $B_i^{(n)}/B^{(n)}$ in Equation 6 by $\varphi_i^{(n,\Pi)} = \mathbb{E}^\Pi[B_i^{(n)}]/\mathbb{E}^\Pi[B^{(n)}]$ (Equation 2).

Interestingly, an approximate maximum likelihood method based on Equation 11 is equivalent to the following approach: consider a family of (approximate) likelihood functions

$$\tilde{L}(\Pi, \underline{\xi}^{(n)}, s) = \prod_{i=1}^{n-1} e^{-[\hat{\theta}(\Pi,s)/2]} \mathbb{E}^\Pi[B^{(n)}] \varphi_i^{(n,\Pi)} \xi_i^{(n)} \times \frac{\left[\frac{\hat{\theta}(\Pi,s)}{2} \mathbb{E}^\Pi[B^{(n)}] \varphi_i^{(n,\Pi)}\right]^{\xi_i^{(n)}}}{\xi_i^{(n)}!} \quad (12)$$

where $\hat{\theta}(\Pi, s) = 2s/\mathbb{E}^\Pi[B^{(n)}]$ is the Watterson estimator for the mutation rate under a Π -coalescent with n leaves when $S = s$ segregating sites are observed [recall Equation 3]. In Equation 12, \tilde{L} is well defined even if $|\underline{\xi}^{(n)}| \neq s$.

The rationale behind Equation 12 is simple: it pretends that the classes are approximately independent and Poisson distributed (this is, of course, not literally true but encouraged by the fact that the off-diagonal entries of the covariance matrix of $\underline{\xi}^{(n)}$ are small compared with the diagonal terms) (see Birkner *et al.* 2013b). Equation 12 is indeed equal to the one obtained from the Poisson random field (PRF) of Sawyer and Hartl (1992), which considers unlinked sites. Within the PRF framework, Equation 12 is an exact likelihood function. In our model of completely linked sites at a single locus, the assumption of independence is merely a convenient computational tool. An analogous approximation of likelihood functions is considered by Bhaskar *et al.* (2015) in the context of varying population sizes; these authors also provide a detailed discussion of the intermediate situation when there is some but not too much recombination between sites at a given locus but free recombination between loci.

For fixed s , we can view $(\Pi, \hat{\theta}(\Pi, s))$ as parameterizing a one-dimensional curve in the full two-dimensional space $H_0 \cup H_1$ defined by the requirement that $\mathbb{E}^{\Pi, \hat{\theta}(\Pi, s)}[S] = s$. The two approximate maximum likelihood approaches based on Equations 11 and 12 are equivalent. Indeed,

$$\tilde{L}(\Pi, \underline{k}^{(n)}, s) = \prod_{i=1}^{n-1} e^{-s\varphi_i^{(n,\Pi)}} \frac{(s\varphi_i^{(n,\Pi)})^{k_i^{(n)}}}{k_i^{(n)}!} = e^{-s} L(\Pi, \underline{k}^{(n)}, s) \quad (13)$$

because the $\varphi_i^{(n,\Pi)}$ sum to 1. Hence, both likelihood functions differ only by the fixed prefactor e^{-s} , so they attain their maximum at the same position.

Thus now we consider the statistic

$$\tilde{Q}_{(E,B)}(\underline{\xi}^{(n)}) := \frac{\sup\{\tilde{L}(\Pi, \underline{k}^{(n)}, |\underline{k}^{(n)}|), \Pi \in \Theta^E\}}{\sup\{\tilde{L}(\Pi, \underline{k}^{(n)}, |\underline{k}^{(n)}|), \Pi \in \Theta^B\}} \quad (14)$$

[where Θ^E and Θ^B refer to the projection of H_0 (resp. H_1) on the coalescent parameter]. For a given value of s , we can then (by simulations using the fixed- s approach) determine approximate quantiles $\tilde{Q}_{(E,B;s)}^*(a)$ associated with a significance level a as in Equation 6 and base our test on the criterion $\tilde{Q}_{(E,B)}(\underline{\xi}^{(n)}) \leq \tilde{Q}_{(E,B;s)}^*(a)$. Similarly, the (approximate) power function $\tilde{G}_{(E,B;s)}$ can be estimated using simulations.

An alternative approach to the (approximate) likelihood-based tests would be rejection rules based on minimal-distance statistics, *i.e.*,

$$Q_{(E,B)}^{(d)}(\underline{\xi}^{(n)}) := \frac{\inf\{d(\varphi^{(n,\Pi)}, \underline{\xi}^{(n)}), \Pi \in \Theta_s^E\}}{\inf\{d(\varphi^{(n,\Pi)}, \underline{\xi}^{(n)}), \Pi \in \Theta_s^B\}} \quad (15)$$

for some suitable distance measure d (*e.g.*, the ℓ_p distance with $p = 2$) with corresponding power function $G_{(E,B;s)}^{(d)}$. We will not discuss the theoretical justification for this method. However, we will use the ℓ_2 distance between normalized expected spectra under various coalescent models to produce three-dimensional heat maps that give some intuitive insight into how a pair of different models out of $\{E, B, A, D\}$ relates to each other depending on the underlying pair of coalescent parameters (*cf. Results*).

We conclude this section with a remark on lumping. One often observes $k_i^{(n)} = 0$ for most i greater than some (small) number m in observed data, in particular, for large n . It thus seems natural to consider (approximate) likelihood functions for *lumped spectra* (*e.g.*, collapsing all entries in classes to the right of some number m into one class m^+), as we have done, *e.g.*, in Figure 1. Another natural type of lumping may be to collect together classes so that $\sum_i \varphi_i^{(n)} \geq x$ for some $x \in (0, 1/2]$. This may not always be feasible, though, if the individual $\varphi_i^{(n)}$ quickly become quite small, and we will refrain from going into a more detailed theoretical discussion of optimal lumpings. However, we will see in our subsequent ABC analysis that adequate lumping can improve the reliability of our model-selection procedure.

Approximate Bayes factors and model selection

In view of the approach and notation of the preceding section, an analogous method of model selection could be based on a *Bayes factor* of the form

$$Q_{(B,E)}^{\mathcal{B}}(\underline{\zeta}^{(n)}) := \frac{\int_{\Theta_s^B} \tilde{L}(\Pi, \underline{\zeta}^{(n)}; s) d\pi_B(\alpha)}{\int_{\Theta_s^E} \tilde{L}(\Pi, \underline{\zeta}^{(n)}; s) d\pi_E(\beta)} \quad (16)$$

(and similar for all other combinations of classes A, D, E, and B) given a pair of priors π_B, π_E on Θ_s^B, Θ_s^E . While the approach

will also work in principle for the two-dimensional prior ranges Θ^B and Θ^E , we will present the (approximate) Bayesian methods with one-dimensional prior ranges (where $S = s$ is treated as a fixed parameter, motivated from the fixed- s approach) so that they complement our previous methods.

Our simulations will be obtained using the rationale behind Equation 12; *i.e.*, after simulating a tree according to a given coalescent parameter, say, α from π_B , mutations are placed on the tree according to a Poisson process, with mutation rate $\hat{\theta}(\Pi; s)$ estimated using Equation 3. However, in Equation 16, we use the normalized site-frequency spectrum $\underline{\zeta}^{(n)}$ as observed statistics because it should be more insensitive to the true coalescent mutation rate θ [potentially deviating from $\hat{\theta}(\Pi, s)$], as argued in the corresponding section in File S1, and thus yield more robust results. \tilde{L} in Equation 16 thus denotes the likelihood function of the observed nSFS under the chosen coalescent model, with the mutation parameter given by Watterson's estimator based on s . Because we estimate the mutation rate based on s , the information loss of using the nSFS instead of the SFS should be only slight. We also experimented with ABC based on simulations using the fixed- s method, *i.e.*, distributing a fixed number of mutations uniformly on the simulated tree, and generally found higher misclassification probabilities (results not shown).

To overcome the problem of exact computation of $\tilde{L}(\Pi, \underline{\zeta}^{(n)}; s)$, which appears infeasible in practice, we employ approximate Bayesian methods (see, *e.g.*, Beaumont 2010) based on the ℓ_2 distance between observed and simulated nSFS. Bayes factors based on further (lumped) distances d and/or the folded nSFS may, of course, also be considered. In line with classical Bayes factor philosophy (*cf.*, *e.g.*, Kass and Raftery 1995), one interprets an observed value of $Q_{(B,E)}^{\mathcal{B}} \gg 1$ as evidence in favor of Θ_s^B over Θ_s^E .

For the ABC analyses, we consider as before on exponential growth (E), algebraic growth (A), and beta- and Dirac coalescents (B and D). Given sample size n and number of segregating sites s , again the coalescent model classes can be parameterized by a single parameter each, which are the exponential growth rate $\beta \in [0, \infty)$, algebraic growth rate $\gamma \in [0, \infty)$, beta-coalescent parameter $\alpha \in [1, 2]$, and mass point location $\psi \in (0, 1]$ for the Dirac coalescent.

For convenience, we employ a simple rejection-based ABC scheme to approximate the Bayes factor for the model (class) comparison given an observed nSFS (resp. folded and/or lumped versions, which can be treated analogously). First, select a number of models (out of $\{E, B, A, D\}$) that should be compared, say, E and B, and choose the corresponding prior distributions on the coalescent parameter ranges. To simulate, say, n_r independent samples of the nSFS from each model, say, from E, independently generate n_r coalescent parameters from the prior π_E and a corresponding coalescent tree Π for each generated coalescent parameter. Distribute independently Poisson mutations with parameter $\hat{\theta}(\Pi; s)/2$ on each such tree, and record the corresponding normalized site-frequency spectra. This should be done independently for all models.

Then fix a tolerance level $x \in (0, 1)$ and count the number of simulations N_E, N_B from each model that are among the $100 \cdot x\%$ best fits with respect to the ℓ_2 distance to the observed nSFS $\underline{z}^{(n)}$ (the “accepted” simulations). Here we use an additional scaling by dividing each class (resp. lumped class) in the nSFS by the median (if nonzero) within this class observed in all simulations as implemented in the R package *abc* (Csilléry *et al.* 2012). The Bayes factor for model E vs. B then can be approximated by

$$Q_{(E,B)}^{\mathcal{B}}(\underline{z}^{(n)}) \approx \frac{N_E}{N_B}$$

To assess how well this ABC approach allows one to distinguish, say, E from B (or, more generally, simultaneously among $\{E, B, A, D\}$), we use two approaches from the R package *abc*. Both are based on leave-one-out cross-validation. More precisely, we pick n_{cv} simulations at random from each model, treat them as the observed value of the nSFS, and then run the ABC approach with the same parameters and simulations as earlier. For each cross-validation sample, say, $\underline{z}_E^{(n)}(i), i \in [n_{cv}]$ from model E, we record the counts of accepted simulations $N_A[\underline{z}_E^{(n)}(i)], N_B[\underline{z}_E^{(n)}(i)], N_D[\underline{z}_E^{(n)}(i)],$ and $N_E[\underline{z}_E^{(n)}(i)]$ from the model classes A, B, D, and E (recall that the chosen cross-validation sample is left out). As measures for the distinction ability of this approach, we record for each model class, borrowing notation from Stoehr *et al.* (2014):

The (estimated) mean posterior probabilities π for model B given the observed nSFS under the true model E, say,

$$\mathbb{E}^E \left[\pi \left(B | \underline{z}^{(n)} \right) \right] \approx \frac{1}{n_{cv}} \sum_{i=1}^{n_{cv}} \frac{N_B[\underline{z}_E^{(n)}(i)]}{N_a}$$

where $N_a = N_A[\underline{z}_E^{(n)}(i)] + N_B[\underline{z}_E^{(n)}(i)] + N_E[\underline{z}_E^{(n)}(i)] + N_D[\underline{z}_E^{(n)}(i)]$ is the number of accepted simulations.

The (estimated) mean misclassification probabilities

$$\begin{aligned} & \mathbb{E}^E \left[\pi \left(\min_{Y \neq B} Q_{(B,Y)}^{\mathcal{B}} \geq 1 \mid \underline{z}^{(n)} \right) \right] \\ & \approx \frac{1}{n_{cv}} \sum_{i=1}^{n_{cv}} \mathbb{1} \left\{ N_B[\underline{z}_E^{(n)}(i)] \geq N_Y[\underline{z}_E^{(n)}(i)] \forall Y \neq B \right\} \end{aligned}$$

for $Y \in \{A, E, D\}$. To ease the notation, we will from now on omit n in the formulas.

In practice, we need to efficiently generate samples of the nSFS under the different models that can be achieved by backward-in-time coalescent simulations. For the exponential growth models (E), we use Hudson’s *ms* (Hudson 2002), as implemented in the R package (R Core Team 2012) *phyclus* (Chen 2011). For algebraic growth models (A), the beta-coalescents (B), and the Dirac coalescents (D), we use custom R and C scripts to generate samples of the nSFS (available

at: <http://page.math.tu-berlin.de/~eldon/programs.html>). To conduct the actual ABC analysis including cross-validation techniques, we employed the R package *abc* (Csilléry *et al.* 2012). Additionally, because we use Watterson’s estimator to set the mutation rate within each model, we compute the mean total length of each coalescent model as described in File S1.

Results

Power estimates of approximate likelihood-ratio tests

To assess the sensitivity of our approximate likelihood-ratio test associated with the likelihood-ratio function (Equation 7), we estimate its power $\tilde{G}_{(E,B;s)}$ from the analog of Equation 10 based on the approximate likelihood from Equation 12 as a function of α (Figure 2A) with $H_0^s = \Theta_s^E$ and $H_1^s = \Theta_s^B$ and estimate $\tilde{G}_{(B,E;s)}$ as a function of β with $H_0^s = \Theta_s^B$ and $H_1^s = \Theta_s^E$ (Figure 2B).

As shown in Figure 2, reasonably high power is obtained to reject Θ_s^E for $n = 500$ and even for a smaller sample size $n = 100$, but the power also depends, as one would expect, on the size of the test. As a side note, we remark that the power estimates $\tilde{G}_{(E,B;s)}$, as a function of α , are right at the size of each corresponding test when $\alpha = 2$ (the Kingman case) as required.

The mitochondrial DNA (mtDNA)–genome analysis of Carr and Marshall (2008), who scanned whole mitochondrial genomes (15,655 bp) of the highly fecund Atlantic cod (*Gadus morhua*), prompted us to briefly investigate the power (Figures S3 and S4 in File S1) with the number of segregating sites $s = 300$. This is nearly the total number of polymorphisms (298) observed among the 32 mtDNA genomes sampled by Carr and Marshall (2008). Our results show that while we may not quite have enough power when $n = 30$ and $s = 300$ (Figure S4A in File S1), we would be in good shape for $n = 100$ (Figures S3 and S4B in File S1). It would be very interesting to analyze such a sample, once available, because it appears to be an open debate whether beta-coalescents should be favored over classical models (including recent population growth) in HFSOD populations (*cf.*, *e.g.*, Steinrücken *et al.* 2013).

Another quite striking observation is that the power of our test is apparently nonmonotone as a function of β when $H_0^s = \Theta_s^B$, in particular, for a smaller type I error. We will present a possible heuristic explanation for this in the *Discussion* section. A rather high power in general is obtained when comparing Θ_s^E and Θ_s^D associated with the Dirac coalescent (Figure S2 in File S1) for $(n, s) = (100, 50)$. For further combinations, we refer to File S1, where Θ_s^A , associated with algebraic growth, is compared with Θ_s^B in Figure S6 and with Θ_s^D in Figure S5. The power functions $\tilde{G}_{(A,D;s)}(\psi)$ (Figure S5 in File S1) are decidedly nonmonotone, as is $\tilde{G}_{(A,B;s)}(\psi)$ (Figure S6A in File S1).

We conclude with a short remark on the sensitivity of our results on lumping of classes in the observed spectrum.

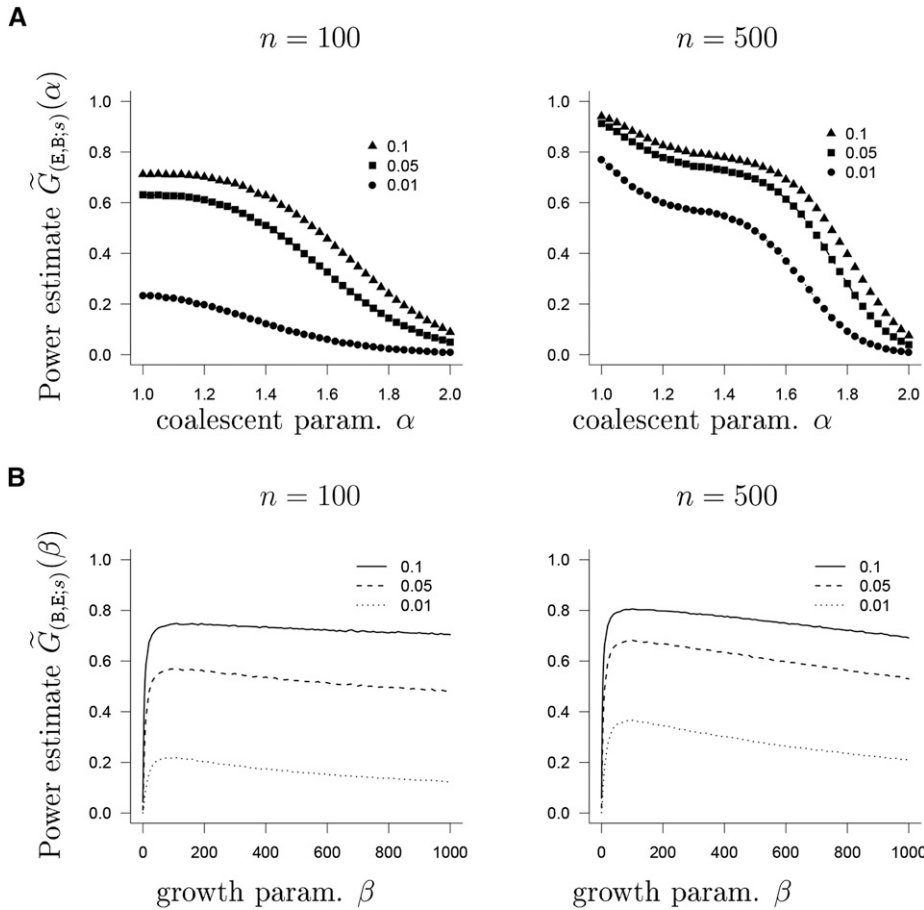


Figure 2 (A) Estimate of $\tilde{G}_{(E,B;S)}$ from Equation 10 based on the approximate likelihood from Equation 12 as a function of α (no lumping) with number of leaves n as shown and $s = 50$. (B) Estimate of $\tilde{G}_{(E,B;S)}$ from Equation 10 based on the approximate likelihood from Equation 12 as a function of β (no lumping) with number of leaves n as shown and $s = 50$. The symbols denote the size of the test, as shown in the legend. The interval hypotheses are discretized to $\Theta_s^E = \{\beta : \beta \in \{0, 1, 2, \dots, 10, 20, \dots, 1000\}\}$ and $\Theta_s^B = \{\alpha : \alpha \in \{1, 1.025, \dots, 2\}\}$. In A, the beta(2 - α , α)-coalescent is the alternative; in B, exponential growth is the alternative.

Indeed, our power estimates suggest that keeping at least the first five classes of the SFS intact and collecting the rest into one other class have little effect on the power of the test (results not shown). Keeping only the singleton ($\xi_1^{(n)}$) class intact and collecting all the rest into one class, however, significantly diminish power (results not shown). C code (cf. Kernighan and Ritchie 1988) written for estimating the power of our tests, where use was made of the GNU Scientific Library (Galassi *et al.* 2013), is available at <http://page.math.tu-berlin.de/~eldon/programs.html>.

Mean-squared distance landscapes for the normalized expected SFS under different growth and coalescent models

Given the potential ability to distinguish between growth and multiple-merger coalescent models, the following questions arise: how does the distance between $\varphi_i^{(n), \Pi_1}$ and $\varphi_i^{(n), \Pi_2}$ behave as a function of the underlying coalescent and growth parameters? Is it possible to visibly identify a one-dimensional curve given by coalescent parameter pairs corresponding to (Π_1, Π_2) along which minimal distance is achieved? Figure 3 and Figure 4 are a brief effort to understand the relation between the expected nSFS for the models in question by graphing the ℓ_2 distance $d_2^{(n)}(X, Y) = \left[\sum_{i=1}^{n-1} (\varphi_i^{(n,X)} - \varphi_i^{(n,Y)})^2 \right]^{1/2}$ as a function of the

coalescent and growth parameters associated with X and Y. In Figure 3, E is compared with B and D. In Figure 4, A is compared with B and D. In Figure 3 and Figure 4, the upper panels show the distance as the respective growth-parameter ranges from 0 to 1000, while the lower panels zoom in on the range from 0 to 10.

Figure 3 indicates the presence of a region, essentially a curve in the two-dimensional (α, β) parameter space, along which the lowest ℓ_2 distance is reached. However, one should be aware that this curve shifts in space when sample size n is increased (data not shown).

Figure 3 suggests that we should have good power to distinguish between algebraic growth and beta-coalescents. However, this seems not to be the case for distinguishing algebraic growth from Dirac coalescents: extreme growth (large γ) seems to produce an almost star-shaped genealogy—consequently, the distance to a Dirac coalescent with ψ close to 1 becomes very small (recall that $\psi = 1$ exactly corresponds to the star-shaped coalescent).

Mean misclassification and posterior probabilities for the ABC approach

In this section we analyze how far an ABC approach using the nSFS (resp. the folded nSFS, abbreviated as nfSFS) and the lumped variants as summary statistics supports our claim that one can distinguish between exponential growth

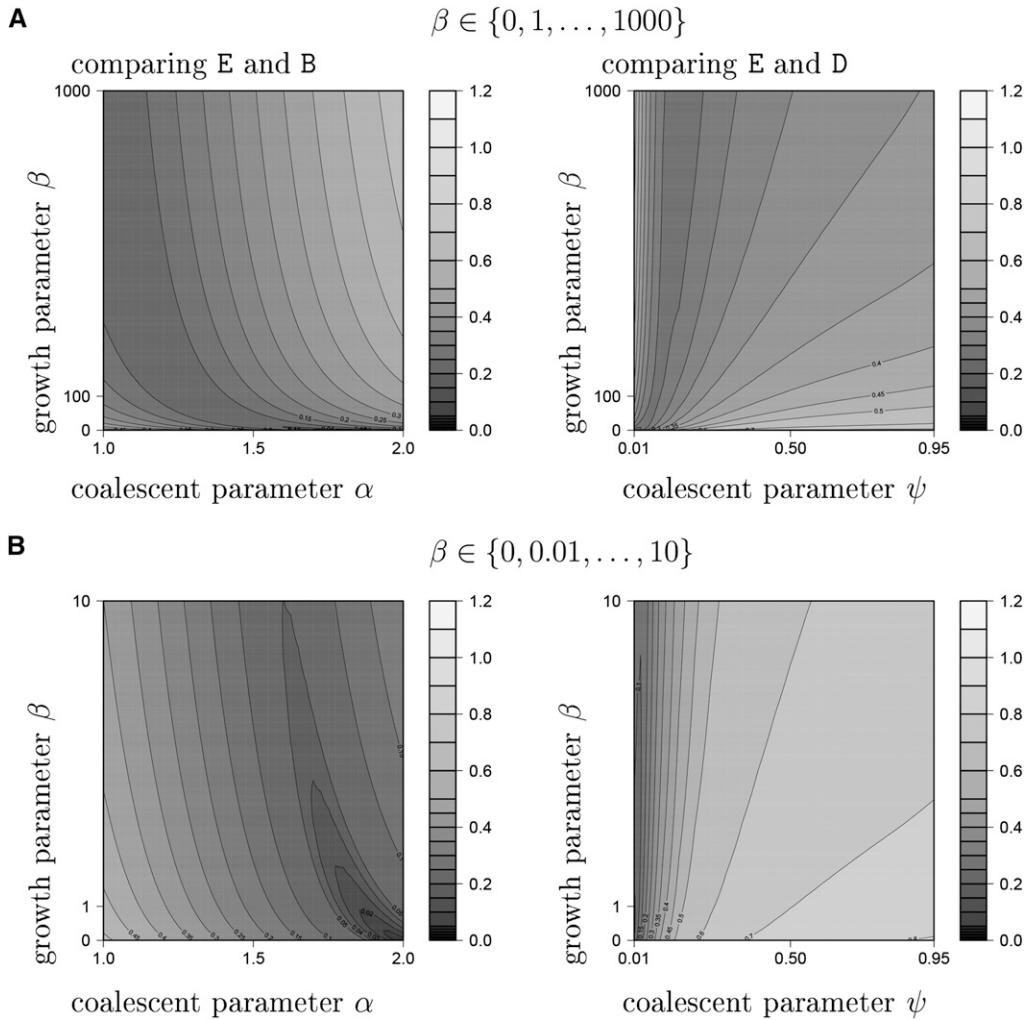


Figure 3 The ℓ_2 distance $d_2^{(n)}(E, X)$ for $X \in \{B, D\}$ of the normalized expected spectra $\phi_i^{(n,E)}$ [see Equation 2] and $\phi_i^{(n,D)}$ as a function of α ($X = B$) [resp. ψ ($X = D$)] and β (E) for number of leaves $n = 100$. Expected values were computed exactly. The grid points are $\alpha \in \{1, 1.025, \dots, 2\}$ and $\psi \in \{0.01, 0.02, \dots, 0.1, 0.15, 0.2, \dots, 0.95\}$; for β as shown.

E and the beta($2 - \alpha, \alpha$)-coalescent B as well as between E, B, and the Dirac coalescent D or between B, D, and algebraic growth A. The distinction ability of the ABC model comparison is assessed based on the simulation procedure and notation described in *Materials and Methods*. Priors were uniform over the full range of coalescent parameters in models B and D and uniform until a maximal cutoff for models E (on $[0, \beta_{\max}]$) and A (on $[0, \gamma_{\max}]$). We discretized the parameter range for the growth models by using increments of 1 or 10 for exponential growth (the first used in all multiple-model comparisons, the latter in the pairwise comparisons between E and B) and increments of 1 for algebraic growth. If not specified otherwise, we used $\beta_{\max} = \gamma_{\max} = 1000$. We fixed a sample size $n = 200$. The number of replications was set to $n_r = 2 \times 10^5$. See Table 1, Table 2, and Table 3 and Tables S4–S8 in File S1 for the estimates of posterior probabilities and misclassification probabilities (some with one replication) with various degrees of lumping and various parameter settings for $n = 200$. For an example with higher sample size $n = 1278$, see Table S9 in File S1.

The estimated error probabilities range from moderate to low values. Mean posterior probabilities $\mathbb{E}^B[\pi(E|\underline{\xi})] \approx 30\%$

indicate a correct classification probability $\mathbb{E}^B[\pi(B|\underline{\xi})] \approx 70\%$, which shows that our method has good distinguishing ability. As expected, lower tolerance generally leads to smaller errors, as do larger mutation rates, while using the folded nSFS increases them. Appropriate lumping seems to decrease the error probabilities on many occasions; see, e.g., Table 1, where a positive effect for strong lumping is observed for $s = 15$ segregating sites, whereas for $s = 75$ in Table S4 in File S1, moderate lumping seems to be more appropriate (both tables show the comparison of models B and E). Not surprisingly, exponential growth rates closer to zero are harder to distinguish from the beta($2 - \alpha, \alpha$)-coalescent models than higher growth rates (see Tables S4 and S6 in File S1). The ABC model comparison distinguishes especially well, even for $s = 15$, between exponential growth from Dirac coalescents and algebraic growth from beta($2 - \alpha, \alpha$)-coalescents (see Table 2, Table 3, and Tables S7 and S8 in File S1). For a relatively low number of segregating sites ($s = 15$), some comparisons (e.g., algebraic growth with Dirac coalescents and beta-coalescents with Dirac coalescents) can lead to common misclassification, but this effect vanishes for larger s . For $s = 75$,

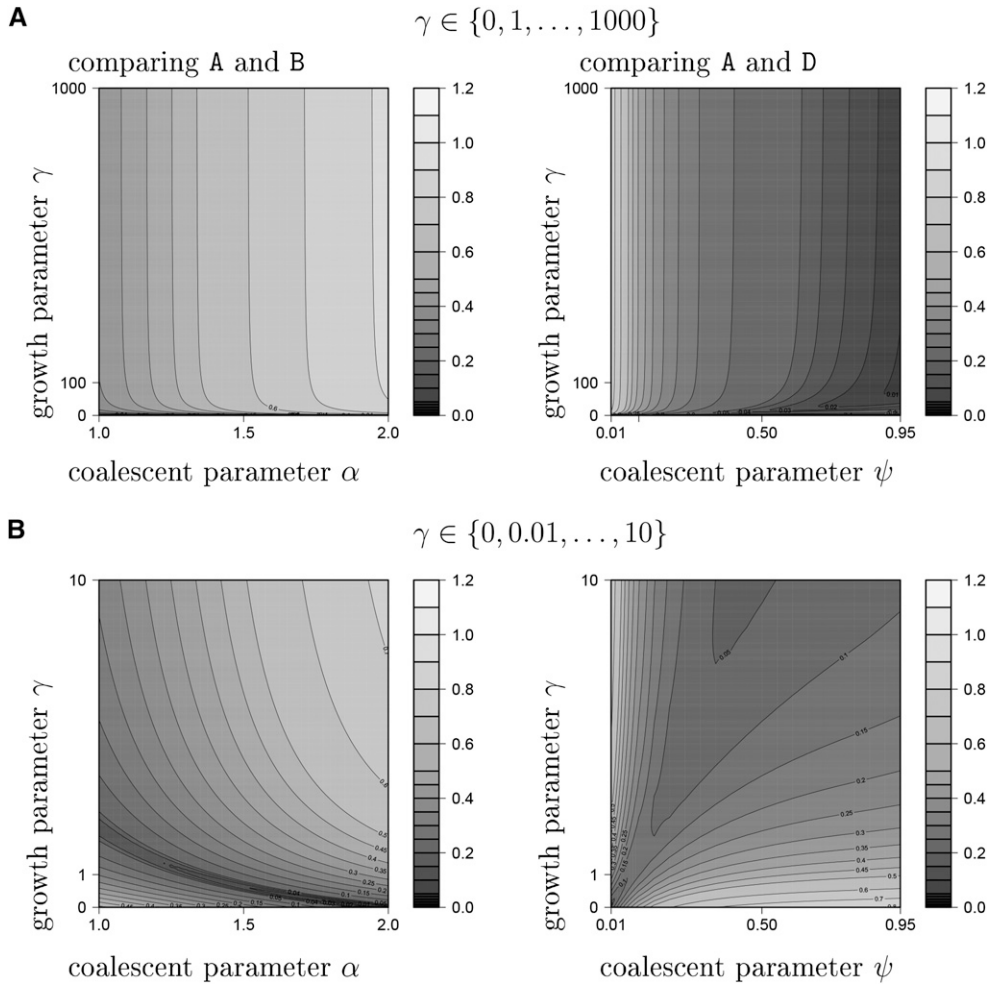


Figure 4 The ℓ_2 distance $d_2^{(n)}(A, X)$ for $X \in \{B, D\}$ of the normalized expected spectra $\phi_i^{(n,A)}$ [see Equation 2] and $\phi_i^{(n,X)}$ as a function of α ($X = B$) [resp. ψ ($X = D$)] and γ (A) for number of leaves $n = 100$. Expected values were computed exactly. The grid points are $\psi \in \{0.01, 0.02, \dots, 0.1, 0.15, 0.2, \dots, 0.95\}$ and $\alpha \in \{1, 1.025, \dots, 2\}$; for γ as shown.

Dirac coalescents can be distinguished relatively well from beta($2 - \alpha, \alpha$)-coalescents (Table 2, Table 3, and Tables S7 and S8 in File S1).

Discussion

The development of methods to distinguish between different (time-changed) coalescent scenarios for the underlying genealogy of a population on the basis of observed data is an important task, in particular, because the choice of an underlying coalescent model affects the estimated coalescent mutation rate $\hat{\theta}$ via (3) and a potential real-time embedding of the genealogy based on (5). Identification of an appropriate coalescent model also may give hints about the underlying reproductive mechanisms present in a population. By way of example, multiple-merger coalescents may indicate the presence of HFSODs in the population.

While inference methods for distinguishing population growth from the usual Kingman coalescent have been studied extensively (see, e.g., Tajima 1989a; Slatkin and Hudson 1991; Rogers and Harpending 1992; Kaj and Krone 2003; Sano and Tachida 2005) and sophisticated theoretical results on the question of identifiability of demographic histories have been

obtained (cf., e.g., Myers *et al.* 2008; Bhaskar and Song 2014; Kim *et al.* 2014), none of these studies has addressed multiple-merger coalescents. In fact, only a few results, e.g., on the statistical properties of the SFS in multiple-merger coalescents (see Birkner *et al.* 2013b) are available.

For the particular case of distinguishing multiple-merger coalescents from population-growth scenarios, this decision problem is complicated by the fact that the patterns of genetic variation produced by the two demographic effects and summarized in the SFS are expected to be similar: both lead to an excess of singletons compared with a classical Kingman coalescent-based genealogy. However, while it is usually possible to match the predicted number of singletons with the observed number in various special cases for both models, the bulk and tail of the spectrum typically will differ (cf. Figure 1 for some examples).

This paper thus is aimed at exploiting and quantifying these differences. However, for feasibility, we had to restrict both the scope and employed methods of our analysis. The first (restrictive) decision in the design of our analysis was the selection of certain subfamilies of Lambda-coalescents and demographic growth scenarios that we deemed suitable for investigation. The reason for restricting to subclasses of

Table 1 Approximations of the mean posterior probabilities and misclassification probabilities for the ABC model comparison between models E and B for tolerance $x = 0.01$, $s = 15$ segregating sites, and using either the nSFS or the nSFS as summary statistics

Fold	Lump	n_{cv}	$\mathbb{E}^B[\pi(E \underline{\zeta})]$	$\mathbb{E}^E[\pi(B \underline{\zeta})]$	$\mathbb{E}^B[\pi(q_{(E,B)}^{\mathcal{B}} \geq 1 \underline{\zeta})]$	$\mathbb{E}^E[\pi(q_{(E,B)}^{\mathcal{B}} \leq 1 \underline{\zeta})]$
No	10+	24,000	0.30	0.25	0.26	0.13
No	50+	12,000	0.32	0.29	0.26	0.17
No	100+	1,200	0.33	0.29	0.28	0.17
Yes	10+	24,000	0.32	0.25	0.28	0.12
Yes	50+	12,000	0.34	0.29	0.28	0.16
Yes	No	12,000	0.34	0.29	0.29	0.16

n_{cv} denotes the number of cross-validations; lump indicates which mutation classes are lumped into one class.

Lambda-coalescents is that the full class of multiple-merger coalescents is in one-to-one relation with the uncountable and nonparametric set of finite measures Λ on $[0, 1]$, which drastically complicates statistical questions, while most of these coalescents do not appear to have a clear biological motivation in terms of a natural underlying population model. Considering the whole Lambda-coalescent class also would raise theoretical questions concerning the unique identifiability of multiple-merger coalescents on the basis of the SFS related to Myers *et al.* (2008) and Bhaskar and Song (2014), and this is mathematically challenging and outside the scope of this work.

Hence, in case of the multiple-merger coalescents, we restricted our attention to the class of beta-coalescents (B) and the Dirac coalescents (D). These classes appeared particularly interesting to us because they both interpolate in a parametric way between the boundary points of the Kingman coalescent (K) via the Bolthausen-Sznitman coalescent (B at $\alpha = 1$) to the star-shaped coalescent (B, $\alpha = 0$; D, $\psi = 1$)—where the whole genealogy collapses to a single line in a single large merger event—among the multiple-merger coalescents. Beta-coalescents have been studied frequently in the literature (see, *e.g.*, Birkner *et al.* 2005; Bertoin and Le Gall 2003; Hallatschek and Neher 2013; Steinrücken *et al.* 2013; Birkner *et al.* 2013b) and are related to a population model with HFSODs (Schweinsberg 2003). Dirac coalescents have been chosen for their simplicity from a mathematical standpoint and also have been investigated by Eldon and Wakeley (2006). The parameter of the Dirac coalescent has a clear interpretation as the fraction of the population that is replaced in each single HFSOD reproductive event.

For similar reasons, we restricted demographic scenarios to two basic parametric growth models. Exponential growth (E) is certainly a natural model in the presence of a supercritical branching population model without geographic or resource restrictions. Our second choice, the algebraic growth model (A), appears perhaps less natural but can reflect situations in which there are spatial or resource limitations and has been analyzed in the mathematical literature (*e.g.*, Schweinsberg 2010). We refrain from more complicated scenarios, such as models with different epochs of exponential growth and recent models including superexponential growth (Reppell *et al.* 2014), which indeed could be investigated with similar methods.

Regarding statistical methodology, one could construct likelihood-based tests on the full two-dimensional parameter spaces for Π and θ given by Θ^E and Θ^B , as outlined in *Materials and Methods*, but this likely would yield considerable computational challenges. Instead, we opted to employ approximate likelihood methods based on the fixed- s method as, *e.g.*, done by Ramos-Onsins and Rozas (2002), reducing our test to a one-dimensional situation, where Equation 14 does not depend on θ at all.

Based on this method, we derive an approximate likelihood-ratio test based on a Poissonization of the SFS via Equation 12 for interval hypotheses, including large ranges of parameters such as the growth parameter β in model E and the coalescent parameter α in model B. By considering the power of our test, a key result in this setup is that even for moderate sample sizes, B and E can be distinguished reasonably well for substantial parts of the parameter space of α and β .

A well-known criticism of this method is its sensitivity on the true yet unknown coalescent mutation rate θ (*cf.* Markovtsova *et al.* 2001). We checked by rejection sampling (*cf.* File S1), conditioning on $S = s$, that for various fixed values of (Π, θ) the rejection probability in our proposed test would be reasonably close to the true rejection probability as long as the true θ is close enough to the Watterson estimate $2s/\mathbb{E}^{\Pi}[B^{(n)}]$, in line with similar observations (for different test statistics) made by Wall and Hudson (2001).

Additional information about the exact coalescent and growth parameters could lead one to test the point hypotheses (*e.g.*, B with fixed α vs. E with fixed β). Indeed, in this case, higher power can be achieved, even for relatively small numbers of segregating sites ($s = 20$), as expected (data not shown).

Distance plots (*e.g.*, Figure 3) over two-dimensional parameter ranges indicate a one-dimensional curve along which the minimal distance is reached. Note that both approaches [maximum (approximate) likelihood and minimum ℓ_2 distance] could be linked if asymptotic normality of our estimators could be established—this is a theoretical question for future work.

Finally, we consider decision rules for the normalized spectrum $\underline{\zeta}^{(n)}$ associated with models A, B, E, and D based on a simple rejection-based ABC analysis. More sophisticated techniques are available [see Beaumont (2010) for an overview] that may improve the prediction accuracy. Empirical misclassification probabilities show, for a reasonable sample

Table 2 Approximations of the mean posterior probabilities for the ABC model comparison among models E, B, and D for tolerance $x = 0.005$, sample size $n = 200$, and $s = 15$ or 75

s	Lump	n_{cv}	$\mathbb{E}^B[\pi(E \xi)]$	$\mathbb{E}^B[\pi(D \xi)]$	$\mathbb{E}^E[\pi(B \xi)]$	$\mathbb{E}^E[\pi(D \xi)]$	$\mathbb{E}^D[\pi(B \xi)]$	$\mathbb{E}^D[\pi(E \xi)]$
15	10+	24,000	0.29	0.12	0.24	0.02	0.55	0.03
15	50+	12,000	0.39	0.10	0.22	0.02	0.54	0.05
15	no	12,000	0.42	0.10	0.21	0.02	0.55	0.07
75	10+	24,000	0.22	0.09	0.12	0.00	0.11	0.01
75	50+	12,000	0.26	0.10	0.12	0.00	0.11	0.01

The nfSFS was used as summary statistics.

size of $n = 200$, at least moderate success in distinguishing among the four model classes even for as few as $s = 15$ segregating sites. Note, though, that depending on the model class comparison to be performed, reasonable error probabilities may be achieved only at higher mutation rates (a higher number s). This indicates that the genealogies produced by the different model classes (at least for suitable sample sizes) are different enough to be distinguished but that mutation rates have to be high enough that these differences are mirrored in the SFS.

In practice, our results could be used to design studies that allow one to distinguish between different conjectured scenarios with suitable power. For example, in marine species, such as Atlantic cod (*cf.*, *e.g.*, Birkner *et al.* 2013b) and Pacific oysters (*cf.* Sargsyan and Wakeley 2008), it has been suggested that certain multiple-merger coalescents could be more appropriate to describe underlying genealogies, and a reproductive mechanism (HFSOD) for population models has been proposed. We have put this to a test by performing an ABC model comparison among our four model classes for the Atlantic cod data of Árnason (2004). The model comparison clearly rejected both Dirac coalescents and algebraic growth as potential models.

While our ABC analysis indicates that exponential growth is slightly favored over the beta(2 - α , α)-coalescent, evidence is not really strong enough to rule out the latter model class. This may indicate that the SFS information of the Árnason (2004) data does not have enough polymorphic sites to distinguish between E and B clearly (our posterior predictive checks revealed that likely neither model class explains the data completely). However, rejection of the D and A model classes suggests that models that predict star-shaped genealogies do not fit the data well. Árnason (2004) used a maximum likelihood estimation method (Kuhner *et al.* 1998) and standard tests of neutrality (Tajima, 1989b; Fu, 1997) to rule out exponential population growth.

At this point, we would like to point out that while our methods and results are exemplified in certain special coalescent and growth models, they could be modified to cover different frameworks.

Before ending the discussion, we wish to comment on a few interesting side issues that appeared during the analysis.

Nonmonotonicity of the power function: At first glance, the observed nonmonotonicity of the power function in the

exponential growth parameter β when compared with certain multiple-merger coalescents (*cf.*, *e.g.*, Figure 2) may appear strange. However, the following example may suggest a heuristic way to understand such behavior in a relatively simple special case. Suppose that one wants to distinguish between an exponential growth model and a multiple-merger coalescent with a substantial Kingman component and a small weight on large multiple mergers (*e.g.*, as in the Dirac coalescent with ψ close to 1). This means that most of the time the multiple-merger coalescent will behave like a Kingman coalescent (producing frequent binary mergers), but with a small rate, comprehensive multiple mergers may occur. Certainly, when the growth parameter β is small, the exponential growth model will yield a pattern of variability close to a Kingman coalescent, and hence the power of a test to distinguish between both will be small if the Kingman component has a weight close to 1. As β increases, the power to distinguish from a Kingman coalescent will increase, in line with intuition. However, as β becomes very large, lineages will coalesce after a *very* short time in the exponential growth model. Such a scenario is certainly different from a Kingman coalescent but could produce patterns of variability closer to a multiple-merger coalescent with a drastic merger after a very short time. Seeing such a merger in the very recent past has some cost (according to the weight of the Dirac component near 1) but appears more likely than observing large amounts of Kingman-like mergers within an unnaturally short time interval, thus leading to a relative decrease in power of associated test. This last effect is nicely illustrated by the upper-right scenario in Figure 4 in the case of a large (algebraic) growth parameter.

Effect of lumping: It is intriguing to see that using the complete nSFS as summary statistics in the ABC approach can yield higher errors than using intermediate (resp. strong) lumpings of the nSFS. A possible explanation is as follows: consider the approximate likelihood function (Equation 12). Assume that the distribution of the SFS is approximately composed of independent Poisson distributions with parameter $(\theta/2)\mathbb{E}^{\Pi}[B_i^{(n)}]$ for $i \in [n - 1]$. For a Poisson-distributed random variable X with parameter κ , we have $\sqrt{\text{Var}(X)}/\mathbb{E}(X) = 1/\sqrt{\kappa}$, thus showing that smaller Poisson parameters yield a higher amount of variation relative to their expected value. Hence classes in the SFS with small

Table 3 Approximations of the mean posterior probabilities for the ABC model comparison among models A, B, and D for tolerance $x = 0.005$, sample size $n = 200$, and $s = 15$ or 75

s	Lump	n_{cv}	$\mathbb{E}^B[\pi(A \xi)]$	$\mathbb{E}^B[\pi(D \xi)]$	$\mathbb{E}^A[\pi(B \xi)]$	$\mathbb{E}^A[\pi(D \xi)]$	$\mathbb{E}^D[\pi(B \xi)]$	$\mathbb{E}^D[\pi(A \xi)]$
15	10+	24,000	0.02	0.12	0.01	0.03	0.16	0.55
15	50+	12,000	0.02	0.11	0.01	0.03	0.24	0.57
75	10+	24,000	0.01	0.09	0.01	0.05	0.10	0.29
75	50+	12,000	0.01	0.08	0.01	0.05	0.15	0.32

The nfSFS was used as summary statistics.

underlying branch lengths (which tend to be in the right tail of the SFS) and/or a low mutation rate show relatively more variation compared with their contribution to the total number of mutations than those with longer branches or if the mutation rate is higher. Lumping such classes together, under Equation 12, yields again a Poisson-distributed lumped class but with the Poisson parameter being the sum of parameters from the classes lumped together. Thus, the variation within this class relative to its contribution to the total number of mutations is reduced by lumping. If different coalescent models show different mean behavior of (lumped) classes, lumping reduces noise and thus increases the chance to correctly identify the underlying model. Naturally, this effect is weakened by higher mutation rates and/or higher sample size n [e.g., consider the limit results for the SFS in Berestycki *et al.* (2014) and Kersting and Stanciu (2015)].

Thus, using an appropriate weighing of the variables in the nSFS (resp. SFS) should improve the power to distinguish between model classes. It also would be a worthwhile future study to see whether a one-dimensional summary of the SFS similar to Tajima's D or Fay and Wu's H , as described in Achaz (2009), could yield a similar or even higher power to distinguish between the model classes than the complete (possibly reweighted) nSFS.

Acknowledgments

F. Freund thanks Luca Ferretti and Guillaume Achaz (SMILE, Collège de France, Paris) for discussions about the site-frequency spectrum. The authors thank two anonymous referees whose insightful comments and constructive criticism helped to improve the presentation. J. Blath and B. Eldon were supported by Deutsche Forschungsgemeinschaft (DFG) grant BL 1105/3-1 and M. Birkner by DFG grant BI 1058/2-1 as part of the SPP Priority Programme 1590.

Literature Cited

Achaz, G., 2009 Frequency spectrum neutrality tests: one for all and all for one. *Genetics* 183: 249–258.
 Árnason, E., 2004 Mitochondrial cytochrome *b* variation in the high-fecundity Atlantic cod: trans-Atlantic clines and shallow gene genealogy. *Genetics* 166: 1871–1885.
 Baragatti, M., and P. Pudlo, 2014 An overview on approximate Bayesian computation. *ESAIM Proc.* 44: 291–299.
 Beaumont, M. A., 2010 Approximate Bayesian computation in evolution and ecology. *Annu. Rev. Ecol. Syst.* 41: 379–406.

Beckenbach, A. T., 1994 Mitochondrial haplotype frequencies in oysters: neutral alternatives to selection models, pp. 188–198 in *Non-Neutral Evolution*, edited by B. Golding. Chapman & Hall, New York.
 Berestycki, J., N. Berestycki, and V. Limic, 2014 Asymptotic sampling formulae for lambda-coalescents. *Ann. Inst. H. Poincaré Probab. Statist.* 50: 715–731.
 Bertoin, J., and J.-F. Le Gall, 2003 Stochastic flows associated to coalescent processes. *Probab. Theory Relat. Fields* 126: 261–288.
 Bhaskar, A., and Y. Song, 2014 Descartes' rule of signs and the identifiability of population demographic models from genomic variation data. *Ann. Stat.* 42: 2469–2493.
 Bhaskar, A., Y. X. R. Wang, and Y. S. Song, 2015 Efficient inference of population size histories and locus-specific mutation rates from large-sample genomic variation data. *Genome Res.* (in press).
 Birkner, M., and J. Blath, 2008 Computing likelihoods for coalescents with multiple collisions in the infinitely many sites model. *J. Math. Biol.* 57: 435–465.
 Birkner, M., J. Blath, M. Capaldo, A. Etheridge, M. Möhle *et al.*, 2005 Alpha-stable branching and beta-coalescents. *Electron. J. Probab.* 10: 303–325.
 Birkner, M., J. Blath, and B. Eldon, 2013a An ancestral recombination graph for diploid populations with skewed offspring distribution. *Genetics* 193: 255–290.
 Birkner, M., J. Blath, and B. Eldon, 2013b Statistical properties of the site-frequency spectrum associated with lambda-coalescents. *Genetics* 195: 1037–1053.
 Birkner, M., J. Blath, and M. Steinrücken, 2011 Importance sampling for lambda-coalescents in the infinitely many sites model. *Theor. Popul. Biol.* 79: 155–173.
 Cannings, C., 1974 The latent roots of certain Markov chains arising in genetics: a new approach. I. Haploid models. *Adv. Appl. Probab.* 6: 260–290.
 Cannings, C., 1975 The latent roots of certain Markov chains arising in genetics: a new approach. II. Further haploid models. *Adv. Appl. Probab.* 7: 264–282.
 Carr, S. M., and H. D. Marshall, 2008 Intraspecific phylogeographic genomics from multiple complete mtDNA genomics in Atlantic cod (*Gadus morhua*): origins of “codmother,” transatlantic vicariance, and midglacial population expansion. *Genetics* 180: 381–389.
 Chen, W.-C., 2011 Overlapping codon model, phylogenetic clustering, and alternative partial expectation conditional maximization algorithm. Ph.D. Thesis, Iowa State University, Ames, IA. Available at: <http://gradworks.umi.com/34/73/3473002.html>.
 Csilléry, K., O. François, and M. G. B. Blum, 2012 ABC: an R package for approximate Bayesian computation (ABC). *Methods Ecol. Evol.* 3: 475–479.
 Cucala, L., and J. Marin, 2013 Bayesian inference on a mixture model with spatial dependence. *J. Comput. Graph. Stat.* 22: 584–597.
 Depaulis, F., S. Mousset, and M. Veuille, 2001 Haplotype tests using coalescent simulations conditional on the number of segregating sites. *Mol. Biol. Evol.* 18: 1136–1138.

- Depaulis, F., and M. Veuille, 1998 Neutrality tests based on the distribution of haplotypes under an infinite-site model. *Mol. Biol. Evol.* 15: 1788–1790.
- Donnelly, P., and T. G. Kurtz, 1999 Particle representations for measure-valued population models. *Ann. Probab.* 27: 166–205.
- Donnelly, P., and S. Tavaré, 1995 Coalescents and genealogical structure under neutrality. *Annu. Rev. Genet.* 29: 401–421.
- Durrett, R., and J. Schweinsberg, 2005 A coalescent model for the effect of advantageous mutations on the genealogy of a population. *Stochastic Processes Appl.* 115: 1628–1657.
- Eldon, B., 2011 Estimation of parameters in large offspring number models and ratios of coalescence times. *Theor. Popul. Biol.* 80: 16–28.
- Eldon, B., and J. Wakeley, 2006 Coalescent processes when the distribution of offspring number among individuals is highly skewed. *Genetics* 172: 2621–2633.
- Fay, J. C., and C.-I. Wu, 2000 Hitchhiking under positive darwinian selection. *Genetics* 155: 1405–1413.
- Fu, Y. X., 1995 Statistical properties of segregating sites. *Theor. Popul. Biol.* 48: 172–197.
- Fu, Y. X., 1997 Statistical tests of neutrality against population growth, hitchhiking, and background selection. *Genetics* 147: 915–925.
- Fu, Y. X., and W. H. Li, 1993 Statistical tests of neutrality of mutations. *Genetics* 133: 693–709.
- Galassi, M., J. Davies, J. Theiler, B. Gough, G. Jungman *et al.*, 2013 *GNU Scientific Library Reference Manual*, Ed. 3. Free Software Foundation, Boston.
- Griffiths, R. C., and S. Tavaré, 1998 The age of a mutation in a general coalescent tree. *Stoch. Models* 14: 273–295.
- Hallatschek, O., and R.-H. Neher, 2013 Genealogies of rapidly adapting populations. *Proc. Natl. Acad. Sci. USA* 110: 437–442.
- Hedgecock, D., and A. I. Pudovkin, 2011 Sweepstakes reproductive success in highly fecund marine fish and shellfish: a review and commentary. *Bull. Mar. Sci.* 87: 971–1002.
- Hein, J., M. H. Schierup, and C. Wiuf, 2005 *Gene Genealogies, Variation and Evolution*. Oxford University Press, Oxford, UK.
- Hudson, R. R., 1983a Properties of a neutral allele model with intragenic recombination. *Theor. Popul. Biol.* 23: 183–201.
- Hudson, R. R., 1983b Testing the constant-rate neutral allele model with protein sequence data. *Evolution* 37: 203–217.
- Hudson, R. R., 1990 Gene genealogies and the coalescent process, pp. 1–44 in *Oxford Surveys in Evolutionary Biology*, Vol. 7, edited by D. J. Futuyma, and J. Antonovics. Oxford University Press, Oxford, UK.
- Hudson, R. R., 2002 Generating samples under a wright-fisher neutral model. *Bioinformatics* 18: 337–338.
- Kaj, I., and S. Krone, 2003 The coalescent process in a population with stochastically varying size. *J. Appl. Probab.* 40: 33–48.
- Kass, R. E., and A. E. Raftery, 1995 Bayes factors. *J. Am. Stat. Assoc.* 90: 773–795.
- Kernighan, B. W., and D. M. Ritchie, 1988 *The C Programming Language*, Ed. 2. Prentice-Hall, Englewood Cliffs, NJ.
- Kersting, G., and I. Stanciu, 2015 The internal branch lengths of the Kingman coalescent. *Ann. Appl. Probab.* (in press).
- Kim, J., E. Mossel, M. Z. Rácz, and N. Ross, 2015 Can one hear the shape of a population history? *Theor. Popul. Biol.* 100: 26–38.
- Kingman, J. F. C., 1982a The coalescent. *Stoch. Processes Appl.* 13: 235–248.
- Kingman, J. F. C., 1982b Exchangeability and the evolution of large populations, pp. 97–112 in *Exchangeability in Probability and Statistics*, edited by G. Koch, and F. Spizzichino. North-Holland, Amsterdam.
- Kingman, J. F. C., 1982c On the genealogy of large populations. *J. Appl. Probab.* 19: 27–43.
- Koskela, J., P. Jenkins, and D. Spanò, 2015 Computational inference beyond Kingman’s coalescent. *J. Appl. Probab.* (in press).
- Kuhner, M. K., J. Yamato, and J. Felsenstein, 1998 Maximum likelihood estimation of population growth rates based on the coalescent. *Genetics* 149: 429–434.
- Markovtsova, L., P. Marjoram, and S. Tavaré, 2001 On a test of Depaulis and Veuille. *Mol. Biol. Evol.* 18: 1132–1133.
- Möhle, M., 1998 Robustness results for the coalescent. *J. Appl. Probab.* 35: 438–447.
- Möhle, M., and S. Sagitov, 2001 Classification of coalescent processes for haploid exchangeable coalescent processes. *Ann. Probab.* 29: 1547–1562.
- Myers, S., C. Fefferman, and N. Patterson, 2008 Can one learn history from the allelic spectrum? *Theor. Popul. Biol.* 73: 342–348.
- Neher, R. A., and O. Hallatschek, 2013 Genealogies of rapidly adapting populations. *Proc. Natl. Acad. Sci. USA* 110: 437–442.
- Nordborg, M., 2001 Coalescent theory, pp. 179–212 in *Handbook of Statistical Genetics*, Ed. 2, edited by D. J. Balding, M. J. Bishop, and C. Cannings. John Wiley & Sons, Chichester, UK.
- Pitman, J., 1999 Coalescents with multiple collisions. *Ann. Probab.* 27: 1870–1902.
- Polanski, A., and M. Kimmel, 2003 New explicit expressions for relative frequencies of single-nucleotide polymorphisms with application to statistical inference on population growth. *Genetics* 165: 427–436.
- Pritchard, J. K., M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman, 1999 Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Mol. Biol. Evol.* 16: 1791–1798.
- Ramírez-Soriano, A., S. E. Ramos-Onsins, J. Rozas, F. Calafell, and A. Navarro, 2008 Statistical power analysis of neutrality tests under demographic expansions, contractions and bottlenecks with recombination. *Genetics* 179: 555–567.
- Ramos-Onsins, S. E., and J. Rozas, 2002 Statistical properties of new neutrality tests against population growth. *Mol. Biol. Evol.* 19: 2092–2100.
- R Core Team, 2012 *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria.
- Reppell, M., M. Boehnke, and S. Zöllner, 2014 The impact of accelerating faster than exponential population growth on genetic variation. *Genetics* 196: 819–828.
- Rödelsperger, C., R. A. Neher, A. M. Weller, G. Eberhardt, H. Witte *et al.*, 2014 Characterization of genetic diversity in the nematode *pristionchus pacificus* from population-scale resequencing data. *Genetics* 196: 1153–1165.
- Rogers, A. R., and H. C. Harpending, 1992 Population growth makes waves in the distribution of pairwise genetic differences. *Mol. Biol. Evol.* 9: 552–569.
- Rubin, D. B., 1984 Bayesian justifiable and relevant frequency calculations for the applied statistician. *Ann. Stat.* 12: 1151–1172.
- Sagitov, S., 1999 The general coalescent with asynchronous mergers of ancestral lines. *J. Appl. Probab.* 36: 1116–1125.
- Sano, A., and H. Tachida, 2005 Gene genealogy and properties of test statistics of neutrality under population growth. *Genetics* 169: 1687–1697.
- Sargsyan, O., and J. Wakeley, 2008 A coalescent process with simultaneous multiple mergers for approximating the gene genealogies of many marine organisms. *Theor. Popul. Biol.* 74: 104–114.
- Sawyer, S. A., and D. L. Hartl, 1992 Population genetics of polymorphism and divergence. *Genetics* 132: 1161–1176.
- Schweinsberg, J., 2003 Coalescent processes obtained from supercritical Galton-Watson processes. *Stochastic Processes Appl.* 106: 107–139.
- Schweinsberg, J., 2010 The number of small blocks in exchangeable random partitions. *ALEA Lat. Am. J. Probab. Math. Stat.* 7: 217–242.

- Simonsen, K. L., G. A. Churchill, and C. F. Aquadro, 1995 Properties of statistical tests of neutrality for DNA polymorphism data. *Genetics* 141: 413–429.
- Slatkin, M., and R. R. Hudson, 1991 Pairwise comparisons of mitochondrial DNA sequences in stable and exponentially growing populations. *Genetics* 129: 555–562.
- Steinrücken, M., M. Birkner, and J. Blath, 2013 Analysis of DNA sequence variation within marine species using beta-coalescents. *Theor. Popul. Biol.* 87: 15–24.
- Stoehr, J., P. Pudlo, and L. Cucala, 2014 Geometric summary statistics for ABC model choice between hidden Gibbs random fields. [arXiv:1402.1380](https://arxiv.org/abs/1402.1380) [Math. ST].
- Tajima, F., 1983 Evolutionary relationship of DNA sequences in finite populations. *Genetics* 105: 437–460.
- Tajima, F., 1989a The effect of change in population size on DNA polymorphism. *Genetics* 123: 597–601.
- Tajima, F., 1989b Statistical methods for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* 123: 585–595.
- Tavaré, S., D. J. Balding, R. C. Griffiths, and P. Donnelly, 1997 Inferring coalescence times from DNA sequence data. *Genetics* 145: 505–518.
- Tellier, A., and C. Lemaire, 2014 Coalescence 2.0: a multiple branching of recent theoretical developments and their applications. *Mol. Ecol.* 23: 2637–2652.
- Wakeley, J., 2007 *Coalescent Theory*, Roberts & Company, Greenwood Village, CO.
- Wall, J. D., and R. R. Hudson, 2001 Coalescent simulations and statistical tests of neutrality. *Mol. Biol. Evol.* 18: 1134–1135.
- Watterson, G. A., 1975 On the number of segregating sites in genetical models without recombination. *Theor. Popul. Biol.* 7: 1539–1546.

Communicating editor: Y. S. Song

GENETICS

Supporting Information

<http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.173807/-/DC1>

Can the Site-Frequency Spectrum Distinguish Exponential Population Growth from Multiple-Merger Coalescents?

Bjarki Eldon, Matthias Birkner, Jochen Blath, and Fabian Freund

SUPPORTING INFORMATION

CAN THE SITE-FREQUENCY SPECTRUM DISTINGUISH EXPONENTIAL POPULATION
GROWTH FROM MULTIPLE-MERGER COALESCENTS?

B. ELDON, M. BIRKNER, J. BLATH, F. FREUND

List of equations from the main text

For ease of reference, we list below the equations from the main text used in Supporting Information. Refer to the main text for explanation of symbols.

Equation (1) from the main text:

$$\mathbb{E}^{\Pi, \theta} \left[\xi_i^{(n)} \right] = \frac{\theta}{2} \sum_{k=2}^{n-i+1} p^{(n), \Pi}[k, i] \cdot k \cdot \mathbb{E}^{\Pi} \left[T_k^{(n)} \right], \quad i \in [n-1]. \quad (\text{S1})$$

Equation (3) from the main text:

$$\hat{\theta}^{\Pi} := \frac{2S}{\mathbb{E}^{\Pi}[B^{(n)}]}, \quad (\text{S2})$$

Equation (5) from the main text:

$$c_N \approx \frac{2\tilde{\mu}}{\theta^{\Pi}}. \quad (\text{S3})$$

Equation (7) from the main text:

$$\varrho_{(\mathbf{E}, \mathbf{B}; s)}(\underline{\xi}^{(n)}) := \frac{\sup \{L(\Pi, \underline{k}^{(n)}, s), \Pi \in \Theta_s^{\mathbf{E}}\}}{\sup \{L(\Pi, \underline{k}^{(n)}, s), \Pi \in \Theta_s^{\mathbf{B}}\}}. \quad (\text{S4})$$

Equation (8) from the main text:

$$\sup_{\Pi \in \Theta_s^{\mathbf{E}}} \mathbb{P}^{\Pi, s} \{ \varrho_{(\mathbf{E}, \mathbf{B}; s)}(\underline{\xi}^{(n)}) \leq \varrho_{(\mathbf{E}, \mathbf{B}; s)}^*(a) \} \leq a. \quad (\text{S5})$$

Equation (9) from the main text:

$$\hat{\theta} = \hat{\theta}(\Pi; s) = \frac{2s}{\mathbb{E}^{\Pi}[B^{(n)}]} \quad (\text{S6})$$

Equation (10) from the main text:

$$G_{(\mathbf{E}, \mathbf{B}; s)}(\Pi) = \mathbb{P}^{\Pi} \{ \varrho_{(\mathbf{E}, \mathbf{B}; s)}(\underline{\xi}^{(n)}) \leq \varrho_{(\mathbf{E}, \mathbf{B}; s)}^*(a, S) \}, \quad \Pi \in \Theta_s^{\mathbf{B}}. \quad (\text{S7})$$

Equation (12) from the main text:

$$\tilde{L}(\Pi, \underline{\xi}^{(n)}, s) = \prod_{i=1}^{n-1} e^{-\frac{\hat{\theta}(\Pi, s)}{2} \mathbb{E}^{\Pi}[B^{(n)}] \varphi_i^{(n, \Pi)}} \frac{\left(\frac{\hat{\theta}(\Pi, s)}{2} \mathbb{E}^{\Pi}[B^{(n)}] \varphi_i^{(n, \Pi)}\right)^{\xi_i^{(n)}}}{\xi_i^{(n)}!} \quad (\text{S8})$$

Multiple merger-coalescents and the model classes K, B and D

A multiple merger- or Lambda-coalescent, formally introduced by PITMAN (1999), SAGITOV (1999), and DONNELLY and KURTZ (1999), is a partition-valued exchangeable coalescent process determined by a finite measure Λ on $[0, 1]$ which governs the dynamics of the process: If there are currently b blocks in the partition (i.e. b active ancestral lineages), k out of them merge at rate

$$\lambda_{b,k} = \int_{[0,1]} x^{k-2} (1-x)^{b-k} \Lambda(dx), \quad k = 2, \dots, b. \quad (\text{S9})$$

For an overview of the theory see e.g. BERESTYCKI (2009) or, with a biological perspective, TELLIER and LEMAIRE (2014). When Λ is associated with the beta-distribution with parameters $2 - \alpha$ and α for $1 \leq \alpha < 2$ (SCHWEINSBERG, 2003), these rates can be given explicitly by

$$\lambda_{b,k} = \frac{B(k - \alpha, b - k + \alpha)}{B(2 - \alpha, \alpha)},$$

where $B(\cdot, \cdot)$ is the classical Beta-function. Such coalescents will be called beta-coalescents, and constitute the model class B.

When Λ is associated with the Dirac coalescent (ELDON and WAKELEY, 2006), that is, $\Lambda(dx) = \delta_{\{\psi\}}(dx)$, for $\psi \in [0, 1]$, we are in class D. Here, for $\psi \in (0, 1]$, the rates are given by

$$\lambda_{b,k} = \frac{\psi^k (1 - \psi)^{b-k}}{\psi^2}.$$

Both classes intersect in the Kingman coalescent (model K), which corresponds to $\alpha = 2$ and

$\psi = 0$, and of course has coalescence rates

$$\lambda_{b,k} = \begin{cases} 1 & \text{if } k = 2, \\ 0 & \text{else,} \end{cases}$$

ie. only binary mergers are allowed. The Beta- and the Dirac coalescent each introduce a *coalescent* parameter (α, ψ) , which can be estimated from genetic data (ELDON, 2011; BIRKNER *et al.*, 2013; BIRKNER and BLATH, 2008; STEINRÜCKEN *et al.*, 2013).

Population models leading to coalescent classes K, B and D

It is well-known that the classical Wright-Fisher and the Moran model have scaling limits whose genealogy is described by a Kingman coalescent. For the more general Lambda-coalescents, MÖHLE and SAGITOV (2001) give a full classification of all Cannings models that lead to any given Lambda-coalescent. The relevant time-scaling is determined by c_N , the probability that in a population of size N , two distinct ancestral lineages merge in the previous generation. It is important to keep in mind that many different population models can lead to the same limiting coalescent, and also that the timescale, determined by c_N , may vary between different models having the same limit. For the Kingman coalescent, the classical Wright Fisher model converges on the time-scale $c_N = 1/N$, whereas for the Moran model, it is of order $1/N^2$.

A popular model that leads to the Beta($2 - \alpha, \alpha$)-coalescent has been introduced by SCHWEINSBERG (2003). For this model, the relevant time-scale is of order $1/N^{\alpha-1}$. Here, single individuals can produce positive fractions of the next generation in a single reproductive event (an instance of ‘HFSOD’) that can be related to stable branching processes, cf. BIRKNER *et al.* (2005). The size of the reproductive event is random and governed by the Beta-distribution. For details we refer to SCHWEINSBERG (2003), and for a discussion of its biological relevance eg. to STEINRÜCKEN *et al.* (2013).

The Dirac coalescent has been investigated in ELDON and WAKELEY (2006). It has a particularly simple interpretation: Given the coalescent parameter $\psi \in (0, 1]$, in each

‘substantial’ reproductive event, a fraction of $100 \cdot \psi\%$ of the generation die and are replaced by the offspring of a single parent (there can be other, ‘non-substantial’ reproductive events which, though potentially frequent, become invisible in the limit). This is an extreme case of HFSOD, and biologically it seems difficult to justify why the fraction ψ should always be the same. However, it is mathematically simple and interpolates between the Kingman coalescent $\psi = 0$ and the star-shaped coalescent $\psi = 1$, thus we included it in our study. For details see ELDON and WAKELEY (2006).

Population with varying population size and the classes E and A

In KAJ and KRONE (2003), a time-changed n -coalescent under a general model of variable population size is derived. More precisely, the authors consider a haploid Wright-Fisher model with population size N at generation $r = 0$ and consider a population size process $M_N(r), r \in \mathbb{Z}$ of the form $M_N(r) = NX_N(r), r \in \mathbb{Z}$, that is, $X_N(r)$ describes the ‘relative population size’ at generation r . Under the assumption that $X_N(\lfloor Nt \rfloor), t \in \mathbb{R}$ converges to something non-degenerate (ie. bounded away from 0 and ∞), they get the well-known limiting result that a time-changed Kingman coalescent describes the genealogy, where the infinitesimal coalescence rates are given by $1/\nu(s)$, with

$$\nu(s) = \lim_{N \rightarrow \infty} X_N(\lfloor Ns \rfloor). \quad (\text{S10})$$

Our exponential growth model E corresponds to a Kingman-coalescent with exponentially growing coalescence rates $\nu(s) = e^{\beta s}$, for $\beta \geq 0$, and can be obtained from a growth rate of β/N per generation in the pre-limiting model, ie. $N_k = N(1 + \beta/N)^k$. Indeed,

$$\nu(t) = \lim_{N \rightarrow \infty} X_N(\lfloor Nt \rfloor) = \lim_{N \rightarrow \infty} \left(1 + \frac{\beta}{N}\right)^{Nt} = e^{\beta t}.$$

Thus, the size Nt generations ago is approximately $Ne^{-\beta t}$.

The model class A is given by Kingman coalescents with algebraically growing coalescence rates, ie. $\nu(s) = s^\gamma$, for $\gamma \geq 0$. Note that if $\gamma = 0$ or $\beta = 0$, we recover the Kingman coalescent

and are back in class K.

A population model for algebraic growth was considered in (SCHWEINSBERG, 2010, Section 1.4): Fix a population size N at the present generation 0, and for notational convenience also for generation -1 (this short period of constant population size will become irrelevant after time-rescaling). For a fixed growth parameter $\gamma > 0$, the population size at the k -th generation before the present (for $k \in \mathbb{N}$) is assumed to be $\lceil Nk^{-\gamma} \rceil$. Measuring time in units of size $N^{\frac{1}{1+\gamma}}$ yields the limiting infinitesimal coalescence rate

$$\nu(t) = \lim_{N \rightarrow \infty} N^{\frac{1}{1+\gamma}} c_N(t, \gamma) = \lim_{N \rightarrow \infty} N^{\frac{1}{1+\gamma}} \frac{(N^{\frac{1}{1+\gamma}} t)^\gamma}{N} = t^\gamma,$$

where $c_N(t, \gamma)$ is the probability that two individuals in generation $N^{\frac{1}{1+\gamma}} t$ choose the same ancestor (uniformly out of the $N(N^{\frac{1}{1+\gamma}} t)^{-\gamma}$ individuals alive in that generation). Consider the time-change (for the scaling limit as $N \rightarrow \infty$)

$$T_t := \frac{t^{\gamma+1}}{\gamma+1} = \int_0^t s^\gamma ds.$$

Then, the genealogy of the algebraic growth model at previous generation t equals in law the state of a classical Kingman coalescent at time T_t . See SCHWEINSBERG (2010) for details.

The expected SFS under variable population size

The effect of fluctuations in population size on the SFS has been investigated in various articles, see eg. GRIFFITHS and TAVARÉ (1998), who derive an analog of (S1), and KAJ and KRONE (2003) who link the Wright-Fisher approximation (with fluctuating population size) with the limiting genealogy.

Recursions for the expected values and covariances of the site-frequency spectrum associated with moderate fluctuations in population size will now be briefly discussed. We will in particular consider numerically tractable recursions for the model classes E and A, based on work by POLANSKI *et al.* (2003) and POLANSKI and KIMMEL (2003).

Consider a time-inhomogeneous Kingman coalescent, started in n lineages, where each

pair of lines present at time $t \geq 0$ merges at a rate $\nu(t)$. Then, the expected frequency spectrum $\mathbb{E}[\xi_i^{(n),\nu}]$, $i \in [n-1]$, is again of the form (S1), and the time-change ν enters only in the distribution of the $T_k^{(n)} = T_k^{(n),\nu}$, $2 \leq k \leq n$, that is, the distribution of the lengths of the time intervals of the block-counting process $Y_t^{(n),\nu}$ during which there are exactly k lineages.

To evaluate $\mathbb{E}[\xi_i^{(n),\nu}]$ one needs information about $\mathbb{E}[T_k^{(n),\nu}]$. Define

$$S_j^{(n),\nu} := T_n^{(n),\nu} + T_{n-1}^{(n),\nu} + \cdots + T_j^{(n),\nu}, \quad j = n, \dots, 2 \quad (\text{S11})$$

to be the time at which the block counting process $Y^{(n),\nu}$ jumps from j to $j-1$ lineages (with the convention $S_{n+1}^{(n),\nu} := 0$). Abbreviate, for $t \geq 0$ and $j \in 2, \dots, n$,

$$F(t) := \int_0^t \nu(u) du \quad \text{and} \quad a_j^{(\vartheta)} := \int_0^\infty e^{-(j)F(s)} ds, \quad (\text{S12})$$

assuming that the first integral in (S12) is finite. It is possible to compute the marginal density of $S_m^{(n),\nu}$ using the well-known fact that the density of a convolution of exponentials with different rates can be written as a linear combination of exponential densities,

$$\mathbb{E}[S_m^{(n),\nu}] = \sum_{j=m}^n c_m^{(j,n)} a_j^{(\vartheta)}, \quad (\text{S13})$$

where

$$c_m^{(j,n)} := \prod_{\substack{m \leq i \leq n \\ i \neq j}} \frac{\binom{i}{2}}{\binom{i}{2} - \binom{j}{2}} = (-1)^{j-m} \frac{(2j-1)m}{j(j-1)} \frac{\binom{n}{j} \binom{j+m-2}{j} \binom{j}{m}}{\binom{n+j-1}{j}}, \quad (\text{S14})$$

(put $c_m^{(j,n)} = 0$ for $j < m$).

POLANSKI and KIMMEL (2003) obtain numerically stable and efficient recursions to compute $\mathbb{E}^\Pi [B_i^{(n)}]$ associated with any time-changed Kingman coalescent Π as follows. For ϑ

denoting the growth parameter associated with process Π ,

$$\mathbb{E}^\Pi \left[B_i^{(n)} \right] = \sum_{j=2}^n W_{i,j}^{(n)} a_j^{(\vartheta)} \quad (\text{S15})$$

where the constants $W_{i,j}^{(n)}$ can be computed recursively (POLANSKI and KIMMEL, 2003);

$$\begin{aligned} W_{i,2}^{(n)} &= \frac{6}{n+1}, \\ W_{i,3}^{(n)} &= \frac{30(n-2i)}{(n+1)(n+2)}, \\ W_{i,j+2}^{(n)} &= \frac{(3+2j)(n-2i)}{j(n+j+1)} W_{i,j+1}^{(n)} - \frac{(1+j)(3+2j)(n-j)}{j(2j-1)(n+j+1)} W_{i,j}^{(n)}. \end{aligned} \quad (\text{S16})$$

We now specify the main ingredient $a_j^{(\vartheta)}$ (depending on $F(t), t \geq 0$ and hence $\nu(t), t \geq 0$) explicitly for two important special cases:

a) Exponential growth. In the case of an exponentially growing population with growth parameter β , that is, $\nu(t) = e^{\beta t}$, we have

$$a_j^{(\beta)} = \frac{1}{\beta} \exp\left(\beta^{-1} \binom{j}{2}\right) E_1\left(\beta^{-1} \binom{j}{2}\right), \quad (\text{S17})$$

where

$$E_1(t) := \int_t^\infty \frac{e^{-x}}{x} dx = \int_1^\infty \frac{e^{-tx}}{x} dx \quad (\text{S18})$$

is an exponential integral function, c.f. e.g. (ABRAMOWITZ and STEGUN, 1964, 5.1.1). One can use numerical integration schemes to compute $E_1(t)$ for smaller values of t (eg. $t < 50$). For larger values of t , one can use the approximation

$$E_1(t) = t^{-1} e^{-t} \sum_{k=0}^{K-1} k! (-t)^{-k}$$

(MILGRAM, 1985), which has error of order $O(K!t^{-K})$.

b) Algebraic ('power law') growth. In the case of algebraic growth of the form $\nu(t) = t^\gamma$ for some $\gamma > 0$, we have

$$a_j^{(\gamma)} = \frac{\Gamma(1/(\gamma + 1))}{(1 + \gamma)^{\gamma/(\gamma+1)}} \binom{j}{2}^{-1/(\gamma+1)}. \quad (\text{S19})$$

Based on Equation (23) in FU (1995), it is also possible to compute the *variance* and the *covariances* of the SFS based on expressions for $\mathbb{E}_\nu[T_k^{(n),\nu} T_l^{(n),\nu}]$, $2 \leq k, l \leq n$, which in turn can be obtained from

$$\mathbb{E}_\nu[T_k^{(n),\nu} T_l^{(n),\nu}] = \mathbb{E}_\nu[S_k^{(n),\nu} S_l^{(n),\nu}] - \mathbb{E}_\nu[S_{k-1}^{(n),\nu} S_l^{(n),\nu}] - \mathbb{E}_\nu[S_k^{(n),\nu} S_{l-1}^{(n),\nu}] + \mathbb{E}_\nu[S_{k-1}^{(n),\nu} S_{l-1}^{(n),\nu}],$$

noting that, in the above notation,

$$\mathbb{E}[(S_m^{(n),\nu})^2] = \int_0^\infty s_m^2 \sum_{j=m}^n c_m^{(j,n)} \nu(s_m) \binom{j}{2} e^{-(j)F(s_m)} ds_m,$$

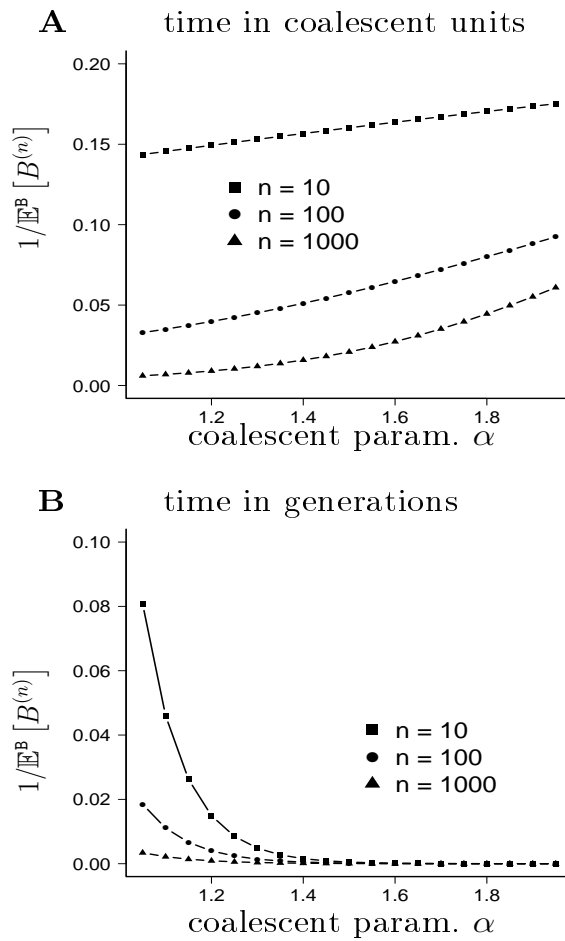
and

$$\mathbb{E}[S_m^{(n),\nu} S_k^{(n),\nu}] = \mathbb{E}[\mathbb{E}[S_m^{(n),\nu} | S_k^{(n),\nu}] S_k^{(n),\nu}],$$

where $\mathbb{E}[S_m^{(n),\nu} | S_k^{(n),\nu} = s_k]$ can be computed (it is the expectation under a regular conditional probability) as in (S13) replacing ν by $\tilde{\nu}(\cdot) := \nu(\cdot + s_k)$, $c_m^{(j,n)}$ by $\tilde{c}_m^{(j)} := c_m^{(j,k)}$ and F by $\tilde{F}(\cdot) = F(s_k + \cdot) - F(s_k)$.

The estimate $1/\mathbb{E}^{\mathbf{B}} [B^{(n)}]$ as a function of α

Figure S1: Graphs of $1/\mathbb{E}^{\mathbf{B}} [B^{(n)}]$, the estimated value of $\theta/2$ per observed mutation when using the Watterson estimator (S2) as a function of α (**A**), compare with (S2); and the estimated value of μ per observed mutation (**B**), using (S3) together with (S2), and assuming the timescale $c_N = N^{1-\alpha}$. The number of leaves n are as shown. In **B**, time is converted into generations by multiplying $\mathbb{E}^{\mathbf{B}} [B^{(n)}]$ with $N^{\alpha-1}$, when $N = 10^5$.



Approximate robustness of the expected normalized SFS w.r.t. θ

In this section, we argue that for a random genealogical tree \mathcal{T} with n leaves whose law is governed by a given coalescent mechanism Π , the expected nSFS $\mathbb{E}^{\Pi, \theta} \left[(\zeta_1^{(n)}, \dots, \zeta_{n-1}^{(n)}) \right]$ when the coalescent mutation rate is $\theta > 0$ is approximately constant as a function of θ . This is useful because it in a sense allows to “factor out” (i.e, ignore) the mutation rate parameter from a test problem when comparing different Π 's. This also means that – at least when the observed number $|\xi^{(n)}|$ of segregating sites is reasonably large – the exact observed value $|\xi^{(n)}|$ does not add much additional information for tests based on the SFS.

Indeed, we can compute

$$\begin{aligned}
 \mathbb{E}^{\Pi, \theta} \left[\zeta_i^{(n)} \right] &= \mathbb{E}^{\Pi, \theta} \left[\zeta_i^{(n)} \mathbf{1}_{\{|\xi^{(n)}| > 0\}} \right] = \mathbb{E}^{\Pi, \theta} \left[\mathbb{E}^{\Pi, \theta} \left[\zeta_i^{(n)} \mathbf{1}_{\{|\xi^{(n)}| > 0\}} \mid \mathcal{T} \right] \right] \\
 &= \mathbb{E}^{\Pi, \theta} \left[\mathbb{P}^{\Pi, \theta} (|\xi^{(n)}| > 0 \mid \mathcal{T}) \frac{\mathbb{E}^{\Pi, \theta} \left[\zeta_i^{(n)} \mathbf{1}_{\{|\xi^{(n)}| > 0\}} \mid \mathcal{T} \right]}{\mathbb{P}^{\Pi, \theta} (|\xi^{(n)}| > 0 \mid \mathcal{T})} \right] \\
 &= \mathbb{E}^{\Pi} \left[\left(1 - e^{-\frac{\theta}{2} \sum_{i=1}^{n-1} B_i^{(n)}} \right) \frac{\frac{\theta}{2} \cdot B_i^{(n)}}{\frac{\theta}{2} \sum_{i=1}^n B_i^{(n)}} \right] \\
 &= \mathbb{E}^{\Pi} \left[\frac{B_i^{(n)}}{\sum_{i=1}^n B_i^{(n)}} \right] - \mathbb{E}^{\Pi} \left[e^{-\frac{\theta}{2} \sum_{i=1}^{n-1} B_i^{(n)}} \frac{B_i^{(n)}}{\sum_{i=1}^n B_i^{(n)}} \right]. \tag{S20}
 \end{aligned}$$

Here, $B_i^{(n)}$ denotes the total length of all branches in \mathcal{T} which subtend i leaves for $i = 1, \dots, n-1$ and in the third line we used Lemma S1.1 below together with the fact that given \mathcal{T} and θ , $\xi_i^{(n)}$, $i = 1, \dots, n-1$ are independent and each $\xi_i^{(n)}$ is Poisson distributed with mean $\frac{\theta}{2} B_i^{(n)}$. Note that the first term in (S20) is independent of θ and the “correction” term is small unless θ is very small or $L_n := \sum_{i=1}^{n-1} B_i^{(n)}$, the total length of \mathcal{T} , is small under Π with substantial probability. Note that for each of the coalescent processes we consider in this investigation, it does hold that $L_n \rightarrow \infty$ as $n \rightarrow \infty$. Simulations also indicate that the distribution (not only the mean) of $\zeta_i^{(n)}$ does not depend much on θ (data not shown).

Lemma S1.1. *Let X_1, X_2 be independent Poisson-distributed variables with parameters a*

and b . Then,

$$\mathbb{E} \left[\frac{X_1}{X_1 + X_2} \mid (X_1 + X_2) > 0 \right] = \frac{a}{a + b}.$$

Proof. $X_1 + X_2$ as a sum of independent Poisson distributed random variables is again Poisson distributed with parameter $a + b$. We have

$$P(X_1 = k, X_2 = m - k \mid (X_1 + X_2) > 0) = \frac{P(X_1 = k)P(X_2 = m - k)}{P(X_1 + X_2 > 0)} = \frac{a^k b^{m-k}}{k!(m-k)!} \frac{e^{-(a+b)}}{1 - e^{-(a+b)}}$$

for $k \in \mathbb{N}_0$, $m \in \mathbb{N}$ with $k \leq m$. We compute

$$\begin{aligned} \mathbb{E} \left[\frac{X_1}{X_1 + X_2} \mid (X_1 + X_2) > 0 \right] &= \sum_{m=1}^{\infty} \sum_{k=0}^m \frac{k}{m} \frac{a^k b^{m-k}}{k!(m-k)!} \frac{e^{-(a+b)}}{1 - e^{-(a+b)}} \\ &= \frac{e^{-(a+b)}}{1 - e^{-(a+b)}} \sum_{m=1}^{\infty} \frac{a}{m(m-1)!} \sum_{k=1}^m \frac{(m-1)!}{(k-1)!((m-1)-(k-1))!} a^{k-1} b^{(m-1)-(k-1)} \\ &= \frac{e^{-(a+b)}}{1 - e^{-(a+b)}} \frac{a}{a+b} \sum_{m=1}^{\infty} \frac{(a+b)^m}{m!} = \frac{e^{-(a+b)}}{1 - e^{-(a+b)}} (e^{(a+b)} - 1) \frac{a}{a+b} = \frac{a}{a+b}. \end{aligned}$$

□

Robustness of the fixed- s -method w.r.t. θ

To check the the robustness of our fixed- s -method against varying θ under rejection sampling (cf. e.g. MARKOVTSOVA *et al.* (2001), WALL and HUDSON (2001)), we applied the following exact rejection sampling approach to simulate a coalescent tree conditional on a given number of observed segregating sites s . As input, the algorithm takes sample size n , number of segregating sites s , a coalescent model Π , and mutation rate θ , and returns a realisation of $\underline{\xi}^{(n)}$ with $|\underline{\xi}^{(n)}| = s$.

Rejection sampling algorithm :

- (i) generate a coalescent tree according to Π , read off branch lengths $B_i^{(n)}$,
- (ii) draw a total number of mutations S as realization of a Poisson random variable with parameter $(\theta/2) \sum_i B_i^{(n)}$,
- (iii) if $S = s$ the required fixed number of segregating sites, keep the $B_i^{(n)}$, otherwise discard and draw again,
- (iv) throw uniformly s mutations on the tree with branch lengths $B_i^{(n)}$, so that the probability of a mutation falling into class i is $B_i^{(n)} / (\sum_i B_i^{(n)})$.

We then computed (approximately via rejection-sampling) the size of a conditional distribution based test if one employs quantiles of the fixed- s -method derived from (S5). Of course, the hope is that both are reasonably close to each other, and this seems to hold relatively well if θ is close to the Watterson estimate $\hat{\theta}(\Pi, s)$ (S6). In particular, the results (Tables (S1)–(S3)) show that the method is particularly robust against varying θ when exponential growth is taken as null model.

Table S1: Checking size of test given fixed- s quantiles associated with size $x\%$ and $\alpha \in \{1, 1.5\}$ with Beta($2 - \alpha, \alpha$)-coalescent as null model, and exponential growth as alternative, using rejection sampling with mutation rate θ as shown. Sample size $n = 100$, segregating sites $s = 50$. The estimate $(\theta_W(\alpha))$ is obtained from (S6). All estimates from 10^5 iterates.

α	$x\%$	θ ($\theta_W(\alpha)$)	size of test	
1.0	10%	3.082453 ($\theta_W(1)$)	0.10	
		2.0	0.13	
		3.0	0.10	
		5.0	0.07	
		7.0	0.06	
	5%	3.082453 ($\theta_W(1)$)	0.05	
		2.0	0.07	
		5.0	0.03	
		7.0	0.02	
	1%	3.082453 ($\theta_W(1)$)	0.01	
		2.0	0.02	
		5.0	0.01	
		7.0	0.002	
	1.5	10%	5.7638 ($\theta_W(1.5)$)	0.11
			3.0	0.03
5.0			0.09	
7.0			0.11	
10.0			0.13	
5%		5.7638 ($\theta_W(1.5)$)	0.05	
		3.0	0.01	
		5.0	0.04	
		7.0	0.06	
1%		5.7638 ($\theta_W(1.5)$)	0.01	
		3.0	0.001	
		5.0	0.01	
		7.0	0.01	
			10.0	0.02

Table S2: Checking size of test given fixed- s quantiles associated with size $x\%$ as shown using rejection sampling with mutation rate θ as shown for exponential growth as null model, and Beta($2 - \alpha, \alpha$)-coalescent as alternative. Sample size $n = 50$, segregating sites $s = 25$. The estimate $(\theta_W(\beta))$ is obtained from (S6). All estimates from 10^5 iterates.

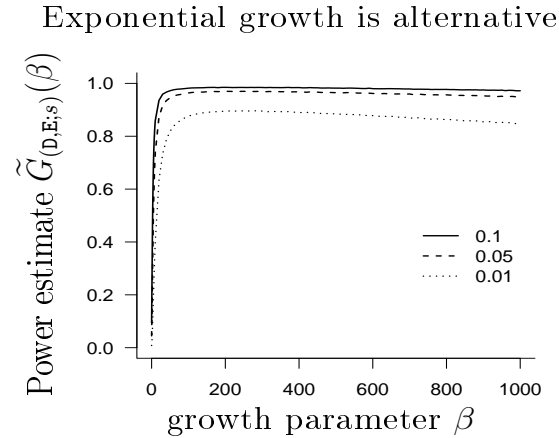
β	$x\%$	θ ($\theta_W(\beta)$)	test size
1	10%	7.895425 $\theta_W(1)$	0.10
		5.0	0.10
		7.0	0.10
		9.0	0.10
		11.0	0.10
	5%	7.895425 $\theta_W(1)$	0.05
		5.0	0.05
		7.0	0.05
		9.0	0.05
		11.0	0.05
	1%	7.895425 $\theta_W(1)$	0.01
		5.0	0.01
		7.0	0.01
		9.0	0.01
		11.0	0.01
10	10%	16.33632 $\theta_W(10)$	0.10
		12.0	0.13
		14.0	0.12
		18.0	0.10
		20.0	0.10
	5%	16.33632 $\theta_W(10)$	0.05
		12.0	0.05
		14.0	0.05
		18.0	0.05
		20.0	0.05
	1%	16.33632 $\theta_W(10)$	0.01
		12.0	0.01
		14.0	0.01
		18.0	0.01
		20.0	0.01

Table S3: Checking size of test given fixed- s quantiles associated with size $x\%$ as shown using rejection sampling with mutation rate θ as shown for exponential growth as null model ($\beta = 1000$), and Beta($2 - \alpha, \alpha$)-coalescent as alternative. The estimate ($\theta_W(\beta)$) is obtained from (S6). Sample size $n = 50$, segregating sites $s = 25$. All estimates from 10^5 iterates.

β	$x\%$	θ ($\theta_W(\beta)$)	test size
1000	10%	263.1798 $\theta_W(10^3)$	0.10
		259	0.10
		261	0.10
		265	0.10
		267	0.10
	5%	263.1798 $\theta_W(10^3)$	0.05
		259	0.05
		261	0.05
		265	0.05
		267	0.05
	1%	263.1798 $\theta_W(10^3)$	0.01
		259	0.01
		261	0.01
		265	0.01
		267	0.01

Estimation of power for $\Theta_0 = \Theta_s^D$, $\Theta_1 = \Theta_s^E$

Figure S2: Estimate of $\tilde{G}_{(D,E;s)}$ from (S7) based on the approximate likelihood (S8) as a function of ψ (no lumping) with number of leaves $n = 100$ and $s = 50$. The line types denote the size of the test as shown in the legend. The interval hypotheses are discretized to $\Theta_s^E = \{\beta : \beta \in \{0, 1, 2, \dots, 10, 20, \dots, 1000\}\}$ and $\Theta_s^D = \{\psi : \psi \in \{0, 0.01, 0.02, \dots, 0.1, 0.15, 0.2, \dots, 0.95\}\}$. Reverting the hypotheses yield very similar results (not shown).



Estimation of power for $s = 300$

Figure S3: Estimate $\tilde{G}_{(B,E;s)}(\beta)$ of power as a function of β for **(A)** $\beta \in \{0, 10, \dots, 1000\}$; **(B)** $\beta \in \{0, 1, 2, \dots, 9, 10, 20, \dots, 1000\}$ when the Beta($2 - \alpha, \alpha$)-coalescent is the null hypothesis, and the test statistic is $\sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^B\} - \sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^E\}$ (S4), with $\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s)$ the log of the Poisson likelihood function (S8) (no lumping). Values at $\beta = 0$ correspond to the Kingman coalescent. A total of 10^6 replicates for both quantiles and power estimates.

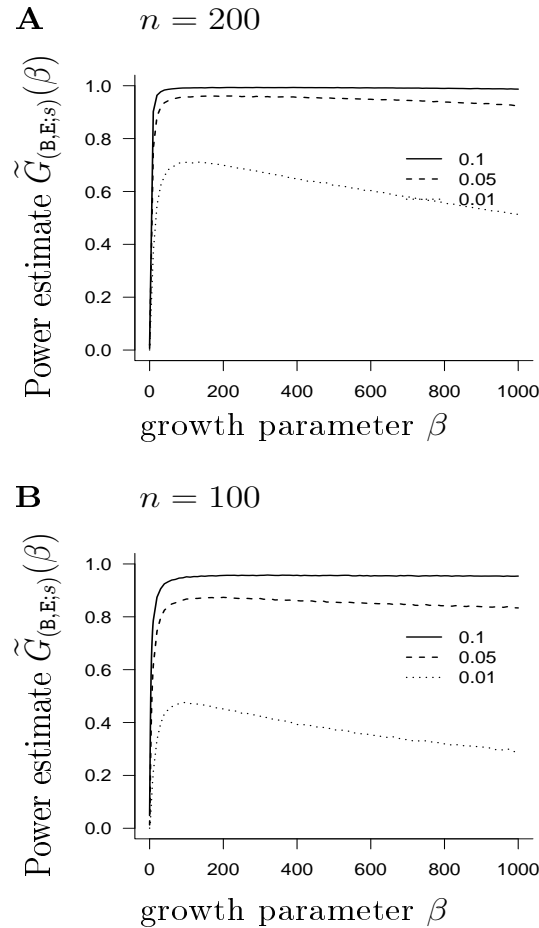
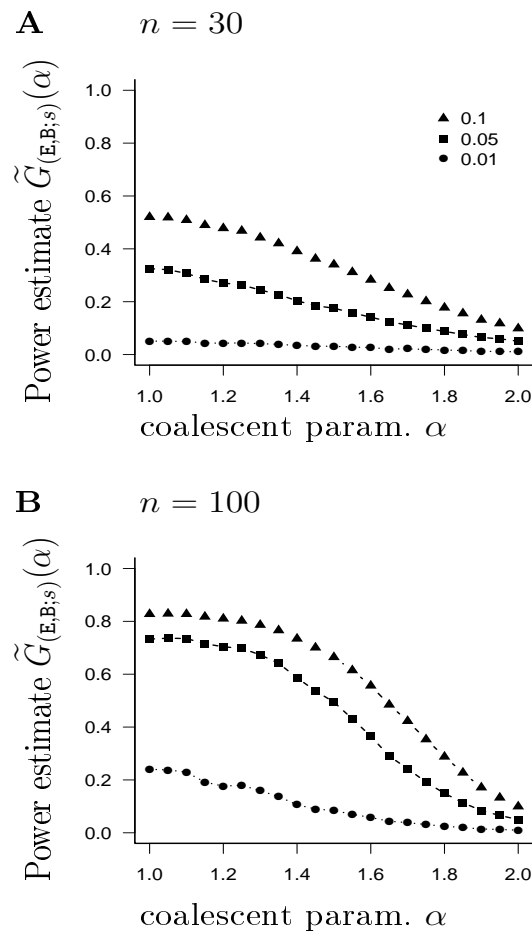
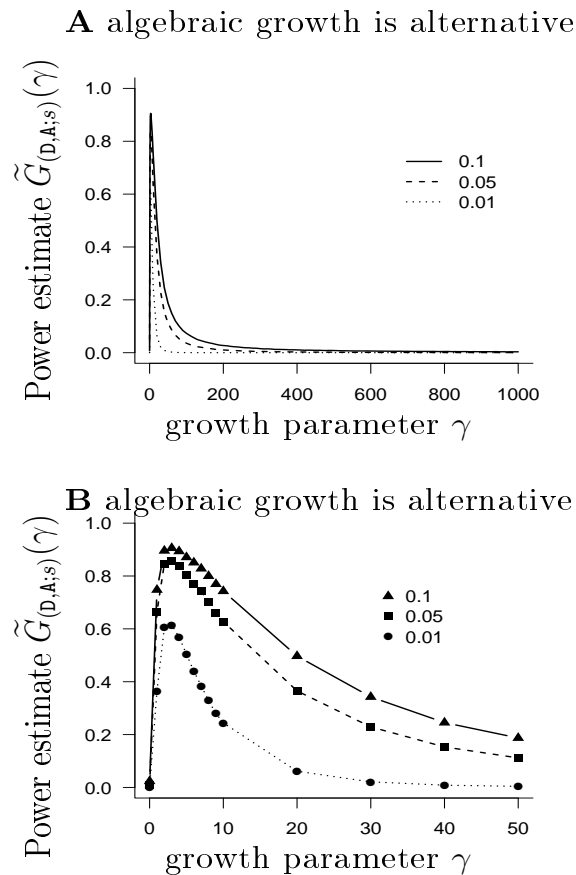


Figure S4: Estimate $\tilde{G}_{(\mathbf{E},\mathbf{B};s)}(\alpha)$ of power as a function of α for $\alpha \in [1, 2]$ when exponential growth (**E**) is the null hypothesis, Beta($2 - \alpha, \alpha$)-coalescent (**B**) is the alternative, and the test statistic is $\sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^{\mathbf{E}}\} - \sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^{\mathbf{B}}\}$ (S4), with $\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s)$ the log of the Poisson likelihood function (S8) (no lumping). Values at $\alpha = 2$ correspond to the Kingman coalescent; number of segregating sites $s = 300$; 10^6 replicates for quantiles and power estimates.



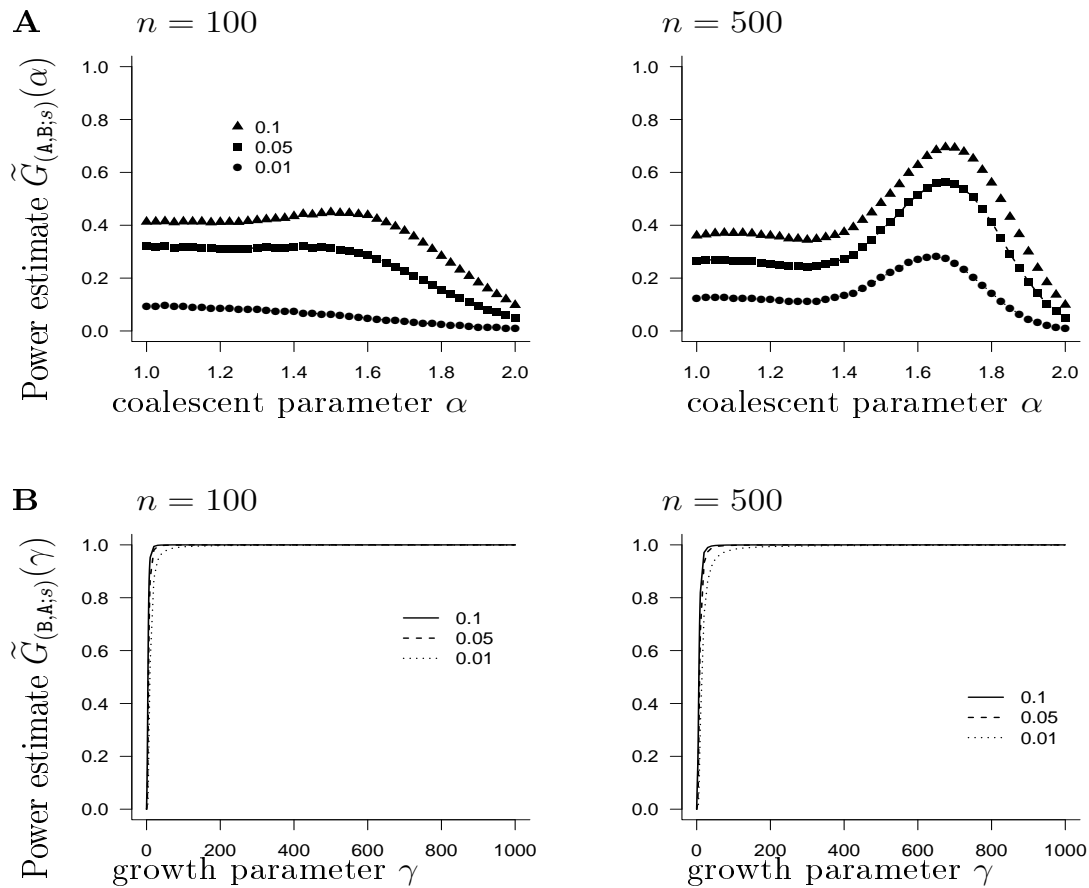
Estimate of power comparing Θ_s^A and Θ_s^D

Figure S5: Estimate $\tilde{G}_{(D,A;s)}(\gamma)$ of power (S7) between algebraic growth and the Dirac Lambda-coalescent when the test statistic is $\sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \vartheta \in \Theta_s^D\} - \sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^A\}$ (S4), with $\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s)$ the log of the Poisson likelihood function (S8) (no lumping); with $n = 100$ and number of segregating sites $s = 50$. The test sizes are as shown in the legend. The interval hypotheses are $\Theta_s^A \equiv \{\gamma : \gamma \in \{0, 1, 2, \dots, 10, 20, 30, \dots, 1000\}\}$ and $\Theta_s^D \equiv \{\psi : \psi \in \{0.01, 0.02, \dots, 0.1, 0.15, 0.2, \dots, 0.95\}\}$. Values at $\gamma = 0$ correspond to the Kingman coalescent. Expected values were computed exactly, and quantiles and power estimated from 10^5 replicates. Reverting the hypotheses shows a very similar pattern (results not shown). In **B**, we ‘zoom in’ on the range $0 \leq \gamma \leq 50$.



Estimate of power comparing Θ_s^A and Θ_s^B

Figure S6: Estimate $\tilde{G}_{(A,B;s)}(\alpha)$ (**A**) and $\tilde{G}_{(B,A;s)}(\gamma)$ (**B**) of power (S7) between algebraic growth and the Beta($2 - \alpha, \alpha$)-coalescent when the test statistic is $\sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^{\Pi_0}\} - \sup\{\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s), \Pi \in \Theta_s^{\Pi_1}\}$ (S4), with $\tilde{\ell}(\Pi, \underline{\xi}^{(n)}, s)$ the log of the Poisson likelihood function (S8) (no lumping); with number of leaves n as shown and number of segregating sites $s = 50$. The test sizes are as shown in the legend. The interval hypotheses are $\Theta_s^A \equiv \{\gamma : \gamma \in \{0, 1, 2, \dots, 10, 20, 30, \dots, 1000\}\}$ and $\Theta_s^B \equiv \{\alpha : \alpha \in \{1, 1.025, \dots, 2\}\}$. Values at $\gamma = 0$ and $\alpha = 2$ correspond to the Kingman coalescent. Expected values were computed exactly, and quantiles and power estimated from 10^5 replicates. In **A**, the Beta($2 - \alpha, \alpha$)-coalescent is the alternative hypothesis; in **B**, algebraic growth is the alternative.



Mean misclassification probabilities and posterior probabilities for ABC approach - alternative parameter choices

Table S4: Approximations of mean posterior probabilities and misclassification probabilities for the ABC model comparison between E and B for different growth parameter ranges or tolerance rates. The nSFS is used as summary statistics. β_{\max} denotes the maximal growth rate used in the growth model, n_{cv} denotes the number of cross-validations; ‘lump’ indicates which mutation classes are lumped into one class. An expected number $s = 75$ of mutations are assumed.

β_{\max}	lump	n_{cv}	tolerance	$\mathbb{E}^{\text{B}} [\pi(\text{E} \underline{\zeta})]$	$\mathbb{E}^{\text{E}} [\pi(\text{B} \underline{\zeta})]$	$\mathbb{E}^{\text{B}} [\pi(e_{(\text{E},\text{B})}^{\text{B}} \geq 1 \underline{\zeta})]$	$\mathbb{E}^{\text{E}} [\pi(e_{(\text{E},\text{B})}^{\text{B}} \leq 1 \underline{\zeta})]$
10^3	10+	24000	0.01	0.24	0.11	0.18	0.04
"	"	"	"	0.24	0.11	0.18	0.04
10^3	50+	12000	0.01	0.22	0.09	0.18	0.03
"	"	"	"	0.23	0.09	0.19	0.03
10^3	100+	1200	0.01	0.22	0.09	0.19	0.03
"	"	12000	"	0.22	0.08	0.20	0.02
10^3	no	12000	0.01	0.30	0.14	0.23	0.04
"	"	"	"	0.30	0.14	0.23	0.04
500	10+	24000	0.01	0.26	0.13	0.20	0.05
500	50+	12000	0.01	0.24	0.10	0.20	0.04
500	100+	1200	0.01	0.26	0.09	0.22	0.03
100	10+	24000	0.01	0.31	0.21	0.23	0.12
100	50+	12000	0.01	0.27	0.18	0.20	0.10
10^3	10+	24000	0.0025	0.20	0.11	0.15	0.05
10^3	50+	12000	0.0025	0.19	0.08	0.15	0.03
10^3	100+	1200	0.0025	0.18	0.08	0.16	0.03
10^3	no	1200	0.0025	0.25	0.13	0.20	0.05

Table S5: Approximations of mean posterior probabilities and misclassification probabilities for the ABC model comparison between E and B for tolerance $x = 0.0025$ and sample size $n = 200$ and assumed expected number $s = 15$ of mutations. The nSFS is used as summary statistics. n_{cv} denotes the number of cross-validations ‘lumped’ indicates which mutation classes are lumped into one class.

lump	n_{cv}	$\mathbb{E}^B [\pi(\mathbf{E} \underline{\zeta})]$	$\mathbb{E}^E [\pi(\mathbf{B} \underline{\zeta})]$	$\mathbb{E}^B [\pi(\varrho_{(\mathbf{E},\mathbf{B})}^B \geq 1 \underline{\zeta})]$	$\mathbb{E}^E [\pi(\varrho_{(\mathbf{E},\mathbf{B})}^B \leq 1 \underline{\zeta})]$
10	24000	0.28	0.24	0.23	0.14
50	12000	0.31	0.26	0.25	0.14
100	12000	0.33	0.27	0.28	0.15
no	12000	0.34	0.26	0.29	0.15

Table S6: Approximations of mean posterior probabilities and misclassification probabilities for the ABC model comparison between E and B for tolerance $x = 0.001$ and sample size $n = 200$, assumed expected number $s = 15$ of mutations and alternative prior ranges and distributions. The nSFS is used as summary statistics. n_{cv} denotes the number of cross-validations ‘lumped’ indicates which mutation classes are lumped into one class. For growth rate β , the prior is uniformly distributed on $\{\beta_{\min}, \beta_{\min} + 10, \dots, \beta_{\max}\}$. For coalescent parameter α , the prior is uniformly distributed on $[\alpha_{\min}, \alpha_{\max}]$

lump	n_{cv}	$\beta_{\min}, \beta_{\max}$	$\alpha_{\min}, \alpha_{\max}$	$\mathbb{E}^B [\pi(\mathbf{E} \underline{\zeta})]$	$\mathbb{E}^E [\pi(\mathbf{B} \underline{\zeta})]$	$\mathbb{E}^B [\pi(\varrho_{(\mathbf{E},\mathbf{B})}^B \geq 1 \underline{\zeta})]$	$\mathbb{E}^E [\pi(\varrho_{(\mathbf{E},\mathbf{B})}^B \leq 1 \underline{\zeta})]$
10	24000	0,100	1.5,2	0.39	0.34	0.30	0.23
50	12000	0,100	1.5,2	0.38	0.31	0.31	0.18
10	24000	100,1000	1,1.5	0.33	0.28	0.29	0.14
50	12000	100,1000	1,1.5	0.36	0.32	0.31	0.18

Table S7: Approximations of the misclassification probabilities for the ABC model comparison between models E, B, D for tolerance $x = 0.005$, sample size $n = 200$ and $s \in \{15, 75\}$. The folded nSFS was used as summary statistics. We use the abbreviation $mc(\Pi_1|\Pi_2) := \mathbb{E}^{\Pi_2} [\pi(\min_{\Pi \neq \Pi_1} \varrho_{(\Pi_1, \Pi)}^B \geq 1|\underline{\zeta}^{(n)})]$, $\Pi_1, \Pi_2 \in \{\mathbf{E}, \mathbf{B}, \mathbf{D}\}$.

s	lump	n_{cv}	$mc(\mathbf{E} \mathbf{B})$	$mc(\mathbf{D} \mathbf{B})$	$mc(\mathbf{B} \mathbf{E})$	$mc(\mathbf{D} \mathbf{E})$	$mc(\mathbf{B} \mathbf{D})$	$mc(\mathbf{E} \mathbf{D})$
15	10+	24000	0.27	0.07	0.12	0.01	0.62	0.01
15	50+	12000	0.39	0.06	0.08	0.01	0.60	0.03
15	no	12000	0.42	0.07	0.08	0.01	0.64	0.04
75	10+	24000	0.19	0.04	0.05	0.00	0.09	0.00
75	50+	12000	0.24	0.04	0.04	0.00	0.09	0.00

Table S8: Approximations of the misclassification probabilities for the ABC model comparison between models A, B, D for tolerance $x = 0.005$, sample size $n = 200$ and $s \in \{15, 75\}$. The folded nSFS was used as summary statistics. We use the abbreviation $mc(\Pi_1|\Pi_2) := \mathbb{E}^{\Pi_2} \left[\pi(\min_{\Pi \neq \Pi_1} \varrho_{(\Pi_1, \Pi)}^{\mathcal{B}} \geq 1 | \underline{\zeta}^{(n)}) \right]$, $\Pi_1, \Pi_2 \in \{\text{A, B, D}\}$.

s	lump	n_{cv}	$mc(\text{A} \text{B})$	$mc(\text{D} \text{B})$	$mc(\text{B} \text{A})$	$mc(\text{D} \text{A})$	$mc(\text{B} \text{D})$	$mc(\text{A} \text{D})$
15	10+	24000	0.01	0.06	0.01	0.04	0.15	0.53
15	50+	12000	0.01	0.06	0.01	0.04	0.18	0.52
75	10+	24000	0.00	0.03	0.01	0.06	0.09	0.25
75	50+	12000	0.00	0.03	0.01	0.05	0.14	0.27

ABC analysis of the cytochrome *b* mtDNA data of ÁRNASON (2004)

To investigate which model class (exponential growth **E**, algebraic growth **A**, Beta($2 - \alpha, \alpha$)-coalescents **B**, Dirac coalescents **D**) fits better to the data, we use the ABC model comparison approach given the (lumped) nfSFS of the observed mitochondrial locus. The exponential growth model class is specified by an uniform prior for growth parameter β on $\{0, 1, 2, \dots, 1000\}$, the algebraic growth class by an uniform prior for growth parameter γ on $\{0, 1, 2, \dots, 1000\}$. The class of Beta n -coalescents is specified by an uniform prior on $\{1, 1.01, \dots, 2\}$ for the coalescent parameter α , the class of Dirac coalescents by an uniform prior on $\{0.01, 0.02, \dots, 0.99\}$ for the coalescent parameter ψ (we omit the star-shaped coalescent $\psi = 1$ because the observed SFS has not only singleton mutations, thus directly violating this model). We used two tolerance levels of 0.005 and 0.00125 and perform $n_r = 200,000$ simulations for each model class. See Table **S9** for the approximated Bayes factors $\varrho_{(\mathbf{E}, \mathbf{B})}^{\mathbf{B}}$ for the model comparison of the growth model and the Beta n -coalescent model using different lumps of the nfSFS as summary statistics. The Bayes factors $\varrho_{(\mathbf{A}, \Pi)}^{\mathbf{B}}, \varrho_{(\mathbf{D}, \Pi)}^{\mathbf{B}}$ for $\Pi \in \{\mathbf{E}, \mathbf{B}\}$ have maximal values of $\approx 0.01, 0.001$ under all lumpings and both tolerances. The observed data fits slightly better to the growth model than to the Beta coalescent class, but not so much better that we could discard the Beta n -coalescents as possible genealogy models for this locus. The latter point is also highlighted by results for an ABC model comparison between only model classes **E** and **B** where all lumpings but 100+ again (slightly) favour the growth model, but for 100+ lumping this is reversed ($\varrho_{(\mathbf{E}, \mathbf{B})}^{\mathbf{B}} = 0.69$ for tolerance 0.005). The Dirac coalescents and the algebraic growth model show neglectible support for all lumpings and thus we discard them as potential models.

Table S9: Approximated Bayes factor $\varrho_{(\mathbf{E}, \mathbf{B})}^{\mathbf{B}}$ given the Atlantic cod mtDNA data

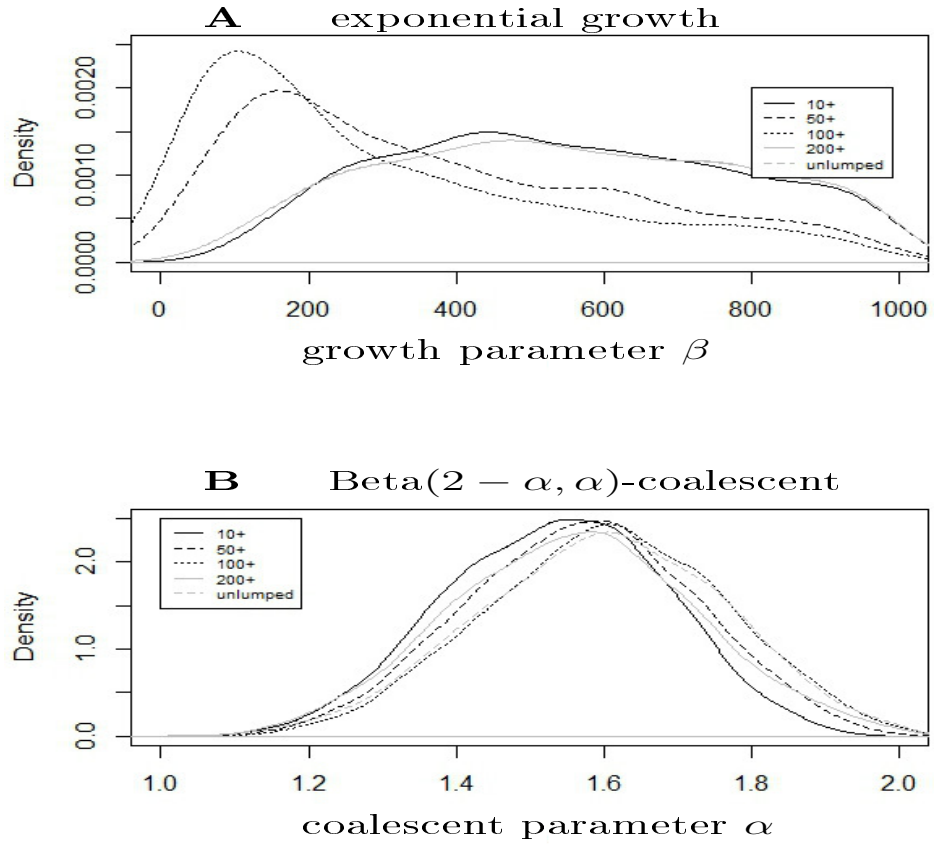
lumping number	10+	50+	100+	200+	no
tolerance 0.005	7.79	2.23	2.09	2.97	2.98
tolerance 0.00125	10.35	2.97	2.23	6.87	7.13

JEFFREYS (1961) suggested interpreting Bayes factors according to the \log_{10} scale. Lump-
ing at 10 (Table **S9**) then gives at least ‘substantial’ ($1/2 < \log_{10}(\varrho_{(\mathbf{E}, \mathbf{B})}^{\mathbf{B}}) \leq 1$) evidence

against the Beta($2 - \alpha, \alpha$)-coalescent in favor of exponential growth. Using KASS and RAFTERY (1995) suggestion of considering Bayes factors on $2 \log_e$ scale gives ‘positive’ ($2 < 2 \log_e(\varrho_{(\mathbf{E}, \mathbf{B})}^{\mathbf{B}}) < 6$) evidence in favor of exponential growth, based on lumping at 10.

Additionally to the ABC model comparison, we also evaluate which parameters fit best to the observed nfSFS at the mitochondrial locus. We omit the class of Dirac coalescents and algebraic growth models from further analysis since the observed frequency spectrum clearly does not fit to this model class. For each other model class used, we record the prior parameters from the 0.5% of the $n_r = 200,000$ simulations that have the smallest ℓ^2 distance to the observed nfSFS (summary statistics). This gives an approximate sample of the posterior distribution of $\pi(\alpha | \text{observed } \underline{\zeta}^{(n)})$ resp. $\pi(\beta | \text{observed } \underline{\zeta}^{(n)})$. Again, we used the lumped nSFS as summary statistics. Figure **S7** shows the posterior distributions for different lumping numbers.

Figure S7: Approximate posterior density of the coalescent parameter from ABC fitting of the (A) growth, and (B) Beta n -coalescent model classes to the observed nfSFS in the Atlantic cod data. Denote by α the Beta n -coalescent parameter, β the growth rate. Priors were uniform on both sets.



ABC quality control for the ÁRNASON (2004) data

We follow the recommendation from the R package `abc` (CSILLÉRY *et al.*, 2012) and perform three checks of quality for the presented ABC approach. We focus on the lumping which gives the clearest distinction, namely the lumping of all classes with mutation counts 10 or higher (class 10+). All checks are performed using the R package `abc`

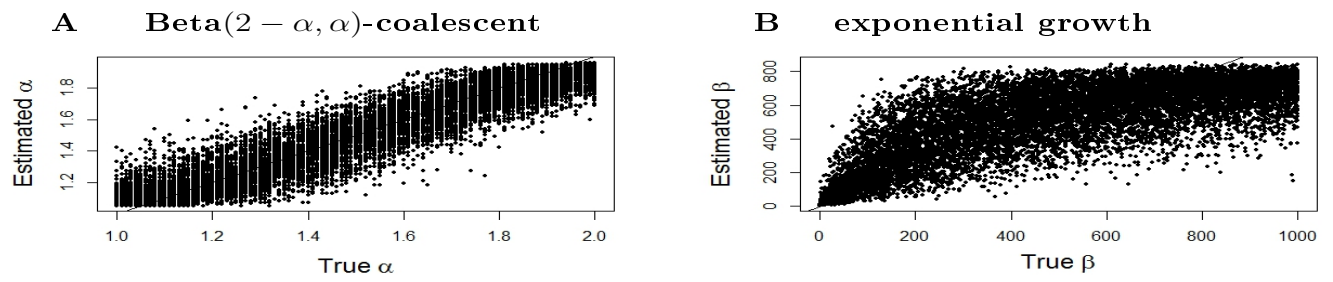
To assess the general ability to distinguish between the two model classes in the setting (i.e., number of observed mutations and sample size) given by the Atlantic cod mtDNA data from ÁRNASON (2004), we again employ a leave-one-out cross-validation as described in Methods. See Table **S10** for the results.

Table S10: Approximations of the mean posterior probabilities and misclassification probabilities (based on $n_{cv} = 12,000$ cross-validations) for the ABC model comparison between models E, A, B, D for tolerance $x = 0.005$, sample size $n = 1278$ and mutation rate estimated via Watterson’s estimator from $s = 39$ observed mutations. The lumped nfSFS (10+) was used as summary statistics. The entries are listed as $\mathbb{E}^{\Pi_{\text{row}}} [\pi(\Pi_{\text{col}}|\underline{\zeta})] / \mathbb{E}^{\Pi_{\text{row}}} [\pi(\min_{\Pi \neq \Pi_{\text{col}}} \varrho_{\Pi_{\text{col}}, \Pi}^{\mathcal{B}} \geq 1 | \underline{\zeta}^{(n)})]$.

	E	A	B	D
E	0.79/0.88	0.00/0.00	0.21/0.12	0.00/0.00
A	0.00/0.00	0.24/0.25	0.06/0.01	0.70/0.74
B	0.24/0.18	0.01/0.00	0.71/0.79	0.03/0.02
D	0.00/0.00	0.03/0.03	0.08/0.03	0.90/0.94

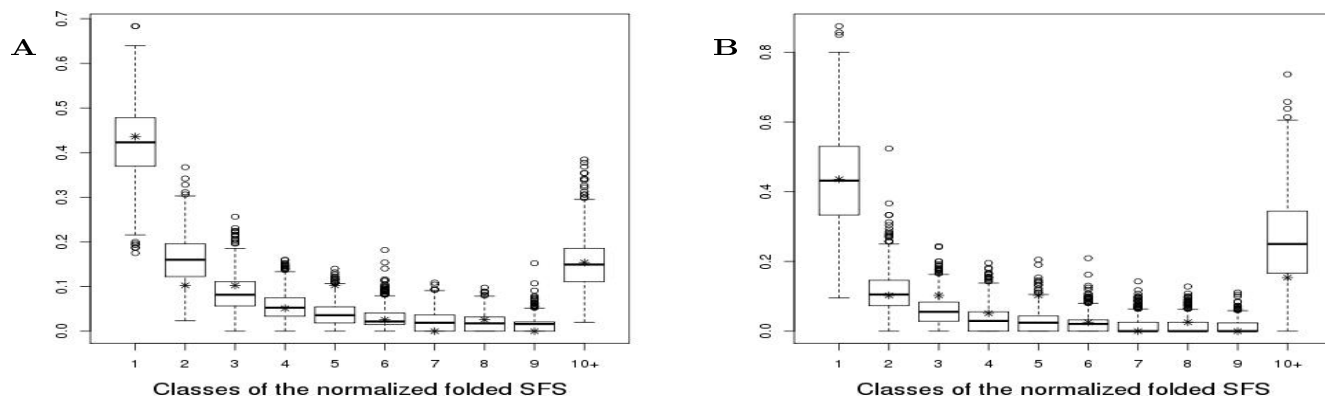
To assess the quality to distinguish the parameters within one model class, we again use leave-one-out cross-validations ($n_{cv} = 12,000$). The parameter of each simulation chosen for cross-validation is estimated as the median of the 0.5% of simulations with the smallest ℓ^2 distance to the chosen simulation. Figure **S8** shows the resulting scatter plots of the parameters of the chosen simulations and the corresponding estimations.

Figure S8: Scatter plots of estimated vs. true parameters of $n_{cv} = 12000$ cross-validated simulations in the (A) Beta coalescent; (B) exponential growth.



To see whether the posterior distributions given the cod mtDNA data from ÁRNASON (2004) define models under which the observed data is reproducible, we performed posterior predictive checks by simulating the 10+ lumped nfSFS under the posterior distribution (i.e., simulating once from each parameter set of each of the 1,000 accepted simulations) for each model class and compare these with the nfSFS observed. See Figure S9 for the results within each nfSFS class. To assess the minimal l^2 distance of the simulations using the posterior parameter distributions from the observed nfSFS, we simulated 5 replications under the posteriors. The minimal l^2 distance was 0.04 under the posterior growth model and 0.06 under the posterior Beta coalescent model.

Figure S9: Posterior predictive checks with 1,000 simulations of the nfSFS under the approximate posterior distributions given the cod data from ÁRNASON (2004) for the (A) Beta coalescent model class; (B) growth model class. Asterisks denote the observed values in the data.



The quality checks reveal that we can not distinguish well within the model classes of exponential growth and of Beta coalescents, but moderately between them. Additionally, the ABC approach distinguishes well between these two classes on one hand and the (non-fitting) other two classes A, D. The posterior predictive checks reveal that both model classes can produce the observed values in each class of the nfSFS, but do not match well in l^2 to

the actual observed nfSFS. Neither model class thus captures the observed nfSFS well.

References

- ABRAMOWITZ, M., and I. A. STEGUN, editors, 1964 *Handbook of mathematical functions with formulas, graphs, and mathematical tables*. Number 55 in Applied Mathematics Series. National Bureau of Standards, Washington, D.C.
- ÁRNASON, E., 2004 Mitochondrial cytochrome *b* variation in the high-fecundity Atlantic cod: trans-Atlantic clines and shallow gene genealogy. *Genetics* **166**: 1871–1885.
- BERESTYCKI, N., 2009 Recent progress in coalescent theory. *Ensaaios Matemáticos* **16**: 1–193.
- BIRKNER, M., and J. BLATH, 2008 Computing likelihoods for coalescents with multiple collisions in the infinitely many sites model. *J Math Biol* **57**: 435–465.
- BIRKNER, M., J. BLATH, M. CAPALDO, A. ETHERIDGE, M. MÖHLE, *et al.*, 2005 Alpha-stable branching and beta-coalescents. *Electron. J. Probab.* **10**: no. 9, 303–325 (electronic).
- BIRKNER, M., J. BLATH, and B. ELDON, 2013 Statistical properties of the site-frequency spectrum associated with lambda-coalescents. *Genetics* **195**: 1037–1053.
- CSILLÉRY, K., O. FRANÇOIS, and M. G. B. BLUM, 2012 ABC: an R package for approximate bayesian computation (ABC). *Methods in Ecology and Evolution* **3**: 475–479.
- DONNELLY, P., and T. G. KURTZ, 1999 Particle representations for measure-valued population models. *Ann Probab* **27**: 166–205.
- ELDON, B., 2011 Estimation of parameters in large offspring number models and ratios of coalescence times. *Theor Popul Biol* **80**: 16–28.
- ELDON, B., and J. WAKELEY, 2006 Coalescent processes when the distribution of offspring number among individuals is highly skewed. *Genetics* **172**: 2621–2633.
- FU, Y. X., 1995 Statistical properties of segregating sites. *Theor Popul Biol* **48**: 172–197.

- GRIFFITHS, R. C., and S. TAVARÉ, 1998 The age of a mutation in a general coalescent tree. *Comm Statistic Stoch Models* **14**: 273–295.
- JEFFREYS, H., 1961 *Theory of Probability*. Oxford University Press, Oxford, UK, 3rd edition.
- KAJ, I., and S. KRONE, 2003 The coalescent process in a population with stochastically varying size. *J Appl Probab* **40**: 33–48.
- KASS, R. E., and A. E. RAFTERY, 1995 Bayes factors. *Journal of the American Statistical Association* **90**: 773–795.
- MARKOVTSOVA, L., P. MARJORAM, and S. TAVARÉ, 2001 On a test of Depaulis and Veuille. *Molecular biology and evolution* **18**: 1132–1133.
- MILGRAM, M. S., 1985 The generalized integro-exponential function. *Math Comp* **44**: 443–458.
- MÖHLE, M., and S. SAGITOV, 2001 Classification of coalescent processes for haploid exchangeable coalescent processes. *Ann Probab* **29**: 1547–1562.
- PITMAN, J., 1999 Coalescents with multiple collisions. *Ann Probab* **27**: 1870–1902.
- POLANSKI, A., A. BOBROWSKI, and M. KIMMEL, 2003 A note on distribution of times to coalescence, under time-dependent population size. *Theor Popul Biol* **63**: 33–40.
- POLANSKI, A., and M. KIMMEL, 2003 New explicit expressions for relative frequencies of single-nucleotide polymorphisms with application to statistical inference on population growth. *Genetics* **165**: 427–436.
- SAGITOV, S., 1999 The general coalescent with asynchronous mergers of ancestral lines. *J Appl Probab* **36**: 1116–1125.
- SCHWEINSBERG, J., 2003 Coalescent processes obtained from supercritical Galton-Watson processes. *Stoch Proc Appl* **106**: 107–139.

- SCHWEINSBERG, J., 2010 The number of small blocks in exchangeable random partitions. *ALEA Lat. Am. J. Probab. Math. Stat.* **7**: 217–242.
- STEINRÜCKEN, M., M. BIRKNER, and J. BLATH, 2013 Analysis of DNA sequence variation within marine species using beta-coalescents. *Theor Popul Biol* **87**: 15–24.
- TELLIER, A., and C. LEMAIRE, 2014 Coalescence 2.0: a multiple branching of recent theoretical developments and their applications. *Mol Ecol* **23**: 2637–2652.
- WALL, J. D., and R. R. HUDSON, 2001 Coalescent simulations and statistical tests of neutrality. *Molecular biology and evolution* **18**: 1134–1135.