

AN ASYMPTOTIC SAMPLING FORMULA FOR THE COALESCENT WITH RECOMBINATION

BY PAUL A. JENKINS* AND YUN S. SONG*,†

University of California, Berkeley

Ewens sampling formula (ESF) is a one-parameter family of probability distributions with a number of intriguing combinatorial connections. This elegant closed-form formula first arose in biology as the stationary probability distribution of a sample configuration at one locus under the infinite-alleles model of mutation. Since its discovery in the early 1970s, the ESF has been used in various biological applications, and has sparked several interesting mathematical generalizations. In the population genetics community, extending the underlying random-mating model to include recombination has received much attention in the past, but no general closed-form sampling formula is currently known even for the simplest extension, that is, a model with two loci. In this paper, we show that it is possible to obtain useful closed-form results in the case the population-scaled recombination rate ρ is large but not necessarily infinite. Specifically, we consider an asymptotic expansion of the two-locus sampling formula in inverse powers of ρ and obtain closed-form expressions for the first few terms in the expansion. Our asymptotic sampling formula applies to arbitrary sample sizes and configurations.

1. Introduction. The probability of a sample configuration provides a useful ground for analyzing genetic data. Popular applications include obtaining maximum likelihood estimates of model parameters and performing ancestral inference [see Stephens (2001)]. In principle, model-based full-likelihood analyses, such as that based on the coalescent [Kingman (1982a,b)], should be among the most powerful methods since they make full use of the data. However, in most cases it is intractable to obtain a closed-form formula for the probability of a given dataset. A well-known exception to this hurdle is the Ewens sampling formula (ESF), which describes the stationary probability distribution of a sample configuration under the one-locus infinite-alleles model in the diffusion limit [Ewens (1972)]. Notable biolog-

*Supported in part by an NIH grant R00-GM080099

†Supported in part by an Alfred P. Sloan Research Fellowship and a Packard Fellowship for Science and Engineering

AMS 2000 subject classifications: Primary 92D15; secondary 65C50,92D10

Keywords and phrases: Ewens sampling formula, coalescent process, recombination, two-locus model, infinite alleles, Golding's recursion

ical applications of this closed-form formula include the test of selective neutrality [see Slatkin (1994); Watterson (1977)]. Hoppe (1984) provided a Pólya-like urn model interpretation of the formula, and recently Griffiths and Lessard (2005) provided a new combinatorial proof of the ESF and extended the framework to obtain new results for the case with a variable population size. We refer the reader to the latter paper for a nice summary of previous works related to the ESF. Note that the ESF also arises in several interesting contexts outside biology, including random partition structures and Bayesian statistics; see Arratia et al. (2003) for examples of intricate combinatorial connections. The ESF is a special case of the two-parameter sampling formula constructed by Pitman (1992, 1995) for exchangeable random partitions.

Golding (1984) considered generalizing the infinite-alleles model to include recombination and constructed a recursion relation satisfied by the two-locus sampling probability distribution at stationarity in the diffusion limit. Ethier and Griffiths (1990) later undertook a more mathematical analysis of the two-locus model and provided several interesting results. However, to date a general closed-form formula for the two-locus sampling distribution remains unknown. Indeed, it is widely recognized that recombination adds a formidably challenging layer of complexity to population genetics analysis. Because obtaining exact analytic results in the presence of recombination is difficult, recent research has focused on developing sophisticated and computationally-intensive Monte Carlo techniques. Examples of such techniques applied to the coalescent include Monte Carlo simulations [see Hudson (1985, 2001)], importance sampling [see De Iorio and Griffiths (2004a,b); Fearnhead and Donnelly (2001); Griffiths and Marjoram (1996); Griffiths et al. (2008); Stephens and Donnelly (2000)], and Markov Chain Monte Carlo methods [see Kuhner et al. (2000); Nielsen (2000); Wang and Rannala (2008)].

Being the simplest model with recombination, the two-locus case has been extensively studied in the past [Ethier and Griffiths (1990); Golding (1984); Griffiths (1981, 1991); Hudson (1985)], and a renewed wave of interest was recently sparked by Hudson (2001), who proposed a composite likelihood method which uses two-locus sampling probabilities as building blocks. LDhat, a widely-used software package developed by McVean and colleagues, is based on this composite likelihood approach, and it has been used to produce a fine-scale map of recombination rate variation in the human genome [McVean et al. (2004); Myers et al. (2005)]. LDhat relies on the importance sampling scheme proposed by Fearnhead and Donnelly (2001) for the coalescent with recombination, to generate exhaustive lookup tables containing

two-locus probabilities for all inequivalent sample configurations and a range of relevant parameter values. This process of generating exhaustive lookup tables is very computationally expensive. A fast and accurate method of estimating two-locus probabilities would be of practical value.

In this paper, we revisit the tantalizing open question of whether a closed-form sampling formula can be found for the coalescent with recombination. We show that, at least for the two-locus infinite-alleles model with the population-scaled recombination rate ρ large but not necessarily infinite, it is possible to obtain useful closed-form analytic results. Our work generalizes previous results [Ethier and Griffiths (1990); Golding (1984)] for $\rho = \infty$, in which case the loci become independent and the two-locus sampling distribution is given by a product of one-locus ESFs. Our main results can be summarized as follows.

Main results. Consider the diffusion limit of the two-locus infinite-alleles model with population-scaled mutation rates θ_A and θ_B at the two loci. For a sample configuration \mathbf{n} (defined later in the text), we use $q(\mathbf{n}|\theta_A, \theta_B, \rho)$ to denote the probability of observing \mathbf{n} given the parameters θ_A, θ_B , and ρ . For an arbitrary \mathbf{n} , our goal is to find an asymptotic expansion of $q(\mathbf{n}|\theta_A, \theta_B, \rho)$ in inverse powers of ρ , i.e., for large values of the recombination rate ρ , our goal is to find

$$q(\mathbf{n}|\theta_A, \theta_B, \rho) = q_0(\mathbf{n}|\theta_A, \theta_B) + \frac{q_1(\mathbf{n}|\theta_A, \theta_B)}{\rho} + \frac{q_2(\mathbf{n}|\theta_A, \theta_B)}{\rho^2} + O\left(\frac{1}{\rho^3}\right),$$

where q_0, q_1 , and q_2 are independent of ρ . As mentioned before, $q_0(\mathbf{n}|\theta_A, \theta_B)$ is given by a product of one-locus ESFs. In this paper, we derive a closed-form formula for the first-order term $q_1(\mathbf{n}|\theta_A, \theta_B)$. Further, we show that the second-order term $q_2(\mathbf{n}|\theta_A, \theta_B)$ can be decomposed into two parts, one for which we obtain a closed-form formula and the other that satisfies a simple strict recursion. The latter can be easily evaluated using dynamic programming. Details of these results are described in Section 3. In a similar vein, in Section 4 we obtain a simple asymptotic formula for the joint probability distribution of the number of alleles observed at the two loci.

2. Preliminaries. In this section, we review the ESF for the one-locus infinite-alleles model, as well as Golding's (1984) recursion relation for the two-locus generalization. Our notational convention generally follows that of Ethier and Griffiths (1990).

Given a positive integer k , $[k]$ denotes the k -set $\{1, \dots, k\}$. For a non-negative real number x and a positive integer n , $(x)_n := x(x+1)\dots(x+n-1)$ denotes the n th ascending factorial of x . We use $\mathbf{0}$ to denote either

a vector or a matrix of all zeroes; it will be clear from context which is intended. Throughout, we consider the diffusion limit of a neutral haploid exchangeable model of random mating with constant population size $2N$. We refer to the haploid individuals in the population as gametes.

2.1. *Ewens sampling formula for the one-locus model.* In the one-locus model, a sample configuration is denoted by a vector of multiplicities $\mathbf{n} = (n_1, \dots, n_K)$, where n_i denotes the number of gametes with allele i at the locus and K denotes the total number of distinct allelic types observed. We use n to denote $\sum_{i=1}^K n_i$, the total sample size. Under the infinite-alleles model, any two gametes can be compared to determine whether or not they have the same allele, but it is not possible to determine how the alleles are related when they are different. Therefore, allelic label is arbitrary. The probability of a mutation event at the locus per gamete per generation is denoted by u . In the diffusion limit, $N \rightarrow \infty$ and $u \rightarrow 0$ with the population-scaled mutation rate $\theta = 4Nu$ held fixed. Each mutation gives rise to a new allele that has never been seen before in the population. For the one-locus model just described, Ewens (1972) obtained the following result:

PROPOSITION 2.1 (Ewens). *At stationarity in the diffusion limit of the one-locus infinite-alleles model with the scaled mutation parameter θ , the probability of an unordered sample configuration $\mathbf{n} = (n_1, \dots, n_K)$ is given by*

$$(2.1) \quad p(\mathbf{n} \mid \theta) = \frac{n!}{n_1 \dots n_K} \frac{1}{\alpha_1! \dots \alpha_n!} \frac{\theta^K}{(\theta)_n},$$

where α_i denotes the number of allele types represented i times, i.e., $\alpha_i := |\{k \mid n_k = i\}|$.

Let \mathcal{A}_n denote an *ordered* configuration of n sequentially sampled gametes such that the corresponding unordered configuration is given by \mathbf{n} . By exchangeability, the probability of \mathcal{A}_n is invariant under all permutations of the sampling order. Hence, we can write this probability of an ordered sample as $q(\mathbf{n})$ without ambiguity. It is given by

$$(2.2) \quad q(\mathbf{n} \mid \theta) = p(\mathbf{n} \mid \theta) \left[\frac{n!}{\prod_{i=1}^K n_i! \alpha_1! \dots \alpha_n!} \right]^{-1} = \left[\prod_{i=1}^K (n_i - 1)! \right] \frac{\theta^K}{(\theta)_n},$$

which follows from the fact that there are $\frac{n!}{\prod_{i=1}^K n_i! \alpha_1! \dots \alpha_n!}$ orderings corresponding to \mathbf{n} [Hoppe (1984)]. It is often more convenient to work with an

ordered sample than with an unordered sample. In this paper, we will work with the former; i.e., we will work with $q(\mathbf{n} \mid \theta)$ rather than $p(\mathbf{n} \mid \theta)$.

In the coalescent process going backwards in time, at each event a lineage is lost either by coalescence or mutation. By consideration of the most recent event back in time, one can show that $q(\mathbf{n} \mid \theta)$ satisfies

$$(2.3) \quad n(n-1+\theta)q(\mathbf{n} \mid \theta) = \sum_{i=1}^K n_i(n_i-1)q(\mathbf{n}-\mathbf{e}_i \mid \theta) + \theta \sum_{i=1}^K \delta_{n_i,1}q(\mathbf{n}-\mathbf{e}_i \mid \theta),$$

where $\delta_{n_i,1}$ is the Kronecker delta and \mathbf{e}_i is a unit vector with the i th entry equal to one and all other entries equal to zero. The boundary condition is $q(\mathbf{e}_i \mid \theta) = 1$ for all $i \in [K]$, and $q(\mathbf{n} \mid \theta)$ is defined to be zero if \mathbf{n} contains any negative component. It can be easily verified that the formula of $q(\mathbf{n} \mid \theta)$ shown in (2.2) satisfies the recursion (2.3).

Ewens (1972) also obtained the following result regarding the number of allelic types:

PROPOSITION 2.2 (Ewens). *Let K_n denote the number of distinct allelic types observed in a sample of size n . Then,*

$$(2.4) \quad \mathbb{P}(K_n = k \mid \theta) = \frac{s(n, k)\theta^k}{(\theta)_n},$$

where $s(n, k)$ are the unsigned Stirling numbers of the first kind. Note that $(\theta)_n = s(n, 1)\theta + s(n, 2)\theta^2 + \dots + s(n, n)\theta^n$.

It follows from (2.1) and (2.4) that K_n is a sufficient statistic for θ .

2.2. *Golding's recursion for the two-locus case.* Golding (1984) first generalized the one-locus recursion (2.3) to two loci, and Ethier and Griffiths (1990) later undertook a more mathematical study of the model. We denote the two loci by A and B , and use θ_A and θ_B to denote the respective population-scaled mutation rates. We use K and L to denote the number of distinct allelic types observed at locus A and locus B , respectively. The population-scaled recombination rate is denoted by $\rho = 4Nr$, where r is the probability of a recombination event between the two loci per gamete per generation. A key observation is that to obtain a closed system of equations, the type space must be extended to allow some gametes to be specified only at one of the two loci.

DEFINITION 2.1 (Extended sample configuration for two loci). *The two-locus sample configuration is denoted by $\mathbf{n} = (\mathbf{a}, \mathbf{b}, \mathbf{c})$, where $\mathbf{a} = (a_1, \dots, a_K)$*

with a_i being the number of gametes with allele i at locus A and unspecified alleles at locus B , $\mathbf{b} = (b_1, \dots, b_L)$ with b_j being the number of gametes with unspecified alleles at locus A and allele j at locus B , and $\mathbf{c} = (c_{ij})$ is a $K \times L$ matrix with c_{ij} being the multiplicity of gametes with allele i at locus A and allele j at locus B . Further, we define

$$\begin{aligned} a &= \sum_{i=1}^K a_i, & c_i &= \sum_{j=1}^L c_{ij}, & c &= \sum_{i=1}^K \sum_{j=1}^L c_{ij}, \\ b &= \sum_{j=1}^L b_j, & c_{.j} &= \sum_{i=1}^K c_{ij}, & n &= a + b + c. \end{aligned}$$

We use $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ to denote the sampling probability of an ordered sample with configuration $(\mathbf{a}, \mathbf{b}, \mathbf{c})$. For ease of notation, we do not show the dependence on parameters. For $0 \leq \rho < \infty$, Golding's (1984) recursion for $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ takes the following form:

$$\begin{aligned} [n(n-1) + \theta_A(a+c) + \theta_B(b+c) + \rho c]q(\mathbf{a}, \mathbf{b}, \mathbf{c}) &= \\ & \sum_{i=1}^K a_i(a_i - 1 + 2c_i)q(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{c}) + \sum_{j=1}^L b_j(b_j - 1 + 2c_{.j})q(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{c}) \\ & + \sum_{i=1}^K \sum_{j=1}^L [c_{ij}(c_{ij} - 1)q(\mathbf{a}, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij}) + 2a_i b_j q(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{c} + \mathbf{e}_{ij})] \\ & + \theta_A \sum_{i=1}^K \left[\sum_{j=1}^L \delta_{a_i+c_i,1} \delta_{c_{ij},1} q(\mathbf{a}, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}) + \delta_{a_i,1} \delta_{c_i,0} q(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{c}) \right] \\ & + \theta_B \sum_{j=1}^L \left[\sum_{i=1}^K \delta_{b_j+c_{.j},1} \delta_{c_{ij},1} q(\mathbf{a} + \mathbf{e}_i, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij}) + \delta_{b_j,1} \delta_{c_{.j},0} q(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{c}) \right] \\ & + \rho \sum_{i=1}^K \sum_{j=1}^L c_{ij} q(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}). \end{aligned} \tag{2.5}$$

Relevant boundary conditions are $q(\mathbf{e}_i, \mathbf{0}, \mathbf{0}) = q(\mathbf{0}, \mathbf{e}_j, \mathbf{0}) = 1$ for all $i \in [K]$ and $j \in [L]$. For notational convenience, we deviate from Ethier and Griffiths (1990) and allow each summation to range over all allelic types. To be consistent, we define $q(\mathbf{a}, \mathbf{b}, \mathbf{c}) = 0$ whenever any entry in \mathbf{a} , \mathbf{b} , or \mathbf{c} is negative.

For ease of discussion, we define the following terms:

DEFINITION 2.2 (Degree). *The degree of $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is defined to be $a + b + 2c$.*

DEFINITION 2.3 (Strictly recursive). *We say that a recursion relation is strictly recursive if it contains only a single term of the highest degree.*

Except in the special case $\rho = \infty$, a closed-form solution for $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is not known. Notice that the terms $q(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{c} + \mathbf{e}_{ij})$ and $q(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij})$ on the right-hand side of (2.5) have the same degree as $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ on the left-hand side. Therefore, (2.5) is not strictly recursive. For each degree, we therefore need to solve a system of coupled equations, and this system grows very rapidly with n . For example, for a sample with $a = 0, b = 0$, and $c = 40$, computing $q(\mathbf{0}, \mathbf{0}, \mathbf{c})$ requires solving a system of more than 20,000 coupled equations [Hudson (2001)]; this is around the limit of sample sizes that can be handled in a reasonable time. In the following section, we revisit the problem of obtaining a closed-form formula for $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ and obtain an asymptotic expansion for large ρ .

3. An asymptotic sampling formula for the two-locus case. For large ρ , our objective is to find an asymptotic expansion of the form

$$(3.1) \quad q(\mathbf{a}, \mathbf{b}, \mathbf{c}) = q_0(\mathbf{a}, \mathbf{b}, \mathbf{c}) + \frac{q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})}{\rho} + \frac{q_2(\mathbf{a}, \mathbf{b}, \mathbf{c})}{\rho^2} + O\left(\frac{1}{\rho^3}\right),$$

where q_0, q_1 , and q_2 are independent of ρ . Our closed-form formulas will be expressed using the following notation:

DEFINITION 3.1. *For a given multiplicity vector $\mathbf{a} = (a_1, \dots, a_K)$ with $a = \sum_{i=1}^K a_i$, we define*

$$(3.2) \quad q^A(\mathbf{a}) = \left[\prod_{i=1}^K (a_i - 1)! \right] \frac{\theta_A^K}{(\theta_A)_a}.$$

Similarly, for a given multiplicity vector $\mathbf{b} = (b_1, \dots, b_L)$ with $b = \sum_{i=1}^L b_i$, we define

$$(3.3) \quad q^B(\mathbf{b}) = \left[\prod_{j=1}^L (b_j - 1)! \right] \frac{\theta_B^L}{(\theta_B)_b}.$$

As discussed in Section 2.1, q^A (respectively, q^B) gives the probability of an ordered sample taken from locus A (respectively, B).

DEFINITION 3.2 (Marginal configuration). *We use $\mathbf{c}_A = (c_i)_{i \in [K]}$ and $\mathbf{c}_B = (c_{\cdot j})_{j \in [L]}$ to denote the marginal sample configurations of \mathbf{c} restricted to locus A and locus B, respectively.*

The leading order term $q_0(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is equal to $q(\mathbf{a}, \mathbf{b}, \mathbf{c})$ when $\rho = \infty$, in which case the two loci are independent. Theorem 2.3 of Ethier and Griffiths (1990) states that $q_0(\mathbf{0}, \mathbf{0}, \mathbf{c}) = q^A(\mathbf{c}_A)q^B(\mathbf{c}_B)$. More generally, one can obtain the following result for the leading order contribution:

PROPOSITION 3.1. *In the asymptotic expansion (3.1) of the two-locus sampling formula, the zeroth order term $q_0(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is given by*

$$(3.4) \quad q_0(\mathbf{a}, \mathbf{b}, \mathbf{c}) = q^A(\mathbf{a} + \mathbf{c}_A)q^B(\mathbf{b} + \mathbf{c}_B).$$

Although this result is intuitively obvious, in Section 5.1 we provide a detailed new proof, since it well illustrates our general strategy. One of the main results of this paper is a closed-form formula for the next order term $q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})$. The case with $\mathbf{c} = \mathbf{0}$ admits a particularly simple solution:

LEMMA 3.1. *In the asymptotic expansion (3.1) of the two-locus sampling formula, the first order term satisfies*

$$q_1(\mathbf{a}, \mathbf{b}, \mathbf{0}) = 0$$

for arbitrary \mathbf{a} and \mathbf{b} .

That $q_1(\mathbf{a}, \mathbf{b}, \mathbf{0})$ vanishes is not expected *a priori*. Below we shall see that $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) \neq 0$ in general. For an arbitrary configuration matrix \mathbf{c} of non-negative integers, we obtain the following closed-form formula for $q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})$:

THEOREM 3.1. *In the asymptotic expansion (3.1) of the two-locus sampling formula, the first order term $q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is given by*

$$(3.5) \quad q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) = \frac{1}{2} \left[\begin{aligned} & c(c-1)q^A(\mathbf{a} + \mathbf{c}_A)q^B(\mathbf{b} + \mathbf{c}_B) \\ & - q^B(\mathbf{b} + \mathbf{c}_B) \sum_{i=1}^K c_i(c_i - 1)q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i) \\ & - q^A(\mathbf{a} + \mathbf{c}_A) \sum_{j=1}^L c_j(c_j - 1)q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j) \\ & + \sum_{i=1}^K \sum_{j=1}^L c_{ij}(c_{ij} - 1)q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i)q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j) \end{aligned} \right],$$

for arbitrary configurations $\mathbf{a}, \mathbf{b}, \mathbf{c}$ of non-negative integers.

Lemma 3.1 is used in proving Theorem 3.1. A proof of Theorem 3.1 is provided in Section 5.2, while a proof of Lemma 3.1 is given in Section 5.3. In principle, similar arguments can be used to find the $(j+1)$ th order term given the j th, although a general expression does not seem to be easy to obtain. In Section 5.4, we provide a proof of the following result for $q_2(\mathbf{a}, \mathbf{b}, \mathbf{c})$:

THEOREM 3.2. *In the asymptotic expansion (3.1) of the two-locus sampling formula, the second order term $q_2(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is of the form*

$$(3.6) \quad q_2(\mathbf{a}, \mathbf{b}, \mathbf{c}) = q_2(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}) + \sigma(\mathbf{a}, \mathbf{b}, \mathbf{c}),$$

where $\sigma(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is given by the analytic formula shown in Appendix A, and $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$ satisfies the following strict recursion:

$$(3.7) \quad \begin{aligned} & [a(a + \theta_A - 1) + b(b + \theta_B - 1)]q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) = \\ & \sum_{i=1}^K a_i(a_i - 1)q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \sum_{j=1}^L b_j(b_j - 1)q_2(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\ & + \theta_A \sum_{i=1}^K \delta_{a_i,1} q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \theta_B \sum_{j=1}^L \delta_{b_j,1} q_2(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\ & + 4 \left[a\theta_A - (\theta_A + a - 1) \sum_{i=1}^K \delta_{a_i,1} \right] \left[b\theta_B - (\theta_B + b - 1) \sum_{j=1}^L \delta_{b_j,1} \right] q^A(\mathbf{a})q^B(\mathbf{b}), \end{aligned}$$

with boundary conditions $q_2(\mathbf{e}_i, \mathbf{0}, \mathbf{0}) = q_2(\mathbf{0}, \mathbf{e}_j, \mathbf{0}) = 0$ for all $i \in [K]$ and $j \in [L]$.

In contrast to $q_1(\mathbf{a}, \mathbf{b}, \mathbf{0})$ (c.f., Lemma 3.1), it turns out that $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$ does not vanish in general. We do not have an analytic solution for $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$, but note that (3.7) is strictly recursive and that it can be easily solved numerically using dynamic programming. Deriving an analytic expression for $\sigma(\mathbf{a}, \mathbf{b}, \mathbf{c})$ in (3.6) is a laborious task, as the long equation in Appendix A suggests. We have written a computer program to verify numerically that our analytic result is correct.

Numerical study (not shown) suggests that the relative contribution of $q_2(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0})$ to $q_2(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is in most cases extremely small. For very large sample sizes, solving the dynamic programming problem for $q_2(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0})$ could become prohibitive, and so we might consider simply omitting it for an analytic estimate of $q_2(\mathbf{a}, \mathbf{b}, \mathbf{c})$. For a given sample configuration it would therefore be desirable to estimate the contribution of this

term without having to calculate it directly. In Appendix B we obtain a close upper bound on $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$ which can be found using an $a \times b$ dynamic programming table. This is much simpler than the dynamic programming problem for (3.7), whose table has $[\prod_{i=1}^K a_i][\prod_{j=1}^L b_j]$ entries.

4. Joint distribution of the number of alleles at the two loci in a sample. Following the same strategy as in the previous section, we can obtain the asymptotic behavior of the joint distribution of the number of alleles observed at the two loci in a sample. To make explicit the dependence of these numbers on the sample size, write the number of alleles at locus A as $K_{a,b,c}$ and the number of alleles at locus B as $L_{a,b,c}$. Ethier and Griffiths (1990) proved that the probability $p(a, b, c; k, l) := \mathbb{P}(K_{a,b,c} = k, L_{a,b,c} = l)$ satisfies the recursion

$$\begin{aligned}
 & [n(n-1) + \theta_A(a+c) + \theta_B(b+c) + \rho c]p(a, b, c; k, l) = \\
 & \quad a(a-1+2c)p(a-1, b, c; k, l) + b(b-1+2c)p(a, b-1, c; k, l) \\
 & \quad + c(c-1)p(a, b, c-1; k, l) + 2abp(a-1, b-1, c+1; k, l) \\
 & \quad + \theta_A [ap(a-1, b, c; k-1, l) + cp(a, b+1, c-1; k-1, l)] \\
 & \quad + \theta_B [bp(a, b-1, c; k, l-1) + cp(a+1, b, c-1; k, l-1)] \\
 (4.1) \quad & + \rho cp(a+1, b+1, c-1; k, l),
 \end{aligned}$$

where $p(a, b, c; k, l) = 0$ if $a < 0, b < 0, c < 0, k < 0, l < 0, a = b = c = 0$, or $k = l = 0$. Equation (4.1) has a unique solution satisfying the initial conditions

$$p(1, 0, 0; k, l) = \delta_{k,1}\delta_{l,0}, \quad p(0, 1, 0; k, l) = \delta_{k,0}\delta_{l,1},$$

for $k, l = 0, 1, \dots, n$.

As with Golding's recursion, equation (4.1) can be solved numerically, but quickly becomes computationally intractable with growing n . The only exception is the special case of $\rho = \infty$, for which the distribution is given by the product of (2.4) for each locus. In what follows, we use the following notation in writing an asymptotic series for $p(a, b, c; k, l)$:

DEFINITION 4.1. *For loci A and B , respectively, we define the analogues of (2.4) as*

$$(4.2) \quad p^A(a; k) = \frac{s(a, k)\theta_A^k}{(\theta_A)_a},$$

and

$$(4.3) \quad p^B(b; l) = \frac{s(b, l)\theta_B^l}{(\theta_B)_b},$$

where $s(a, k)$ and $s(b, l)$ are the Stirling numbers of the first kind.

We pose the expansion

$$(4.4) \quad p(a, b, c; k, l) = p_0(a, b, c; k, l) + \frac{p_1(a, b, c; k, l)}{\rho} + O\left(\frac{1}{\rho^2}\right),$$

for large ρ . Then, in Section 5.5 we prove the following result for the zeroth order term:

PROPOSITION 4.1. *For an asymptotic expansion of the form (4.4) satisfying the recursion (4.1), $p_0(a, b, c; k, l)$ is given by*

$$(4.5) \quad p_0(a, b, c; k, l) = p^A(a + c; k)p^B(b + c; l).$$

Similar to Lemma 3.1, we obtain the following vanishing result for the first order term in the case of $c = 0$:

LEMMA 4.1. *For an asymptotic expansion of the form (4.4) satisfying the recursion (4.1), we have*

$$p_1(a, b, 0; k, l) = 0.$$

Using this lemma, it is then possible to obtain the following result for an arbitrary c :

PROPOSITION 4.2. *For an asymptotic expansion of the form (4.4) satisfying the recursion (4.1), $p_1(a, b, c; k, l)$ is given by*

$$(4.6) \quad p_1(a, b, c; k, l) = \frac{c(c-1)}{2} [p^A(a + c; k) - p^A(a + c - 1; k)] \times [p^B(b + c; l) - p^B(b + c - 1; l)].$$

Proofs of Proposition 4.2 and Lemma 4.1 are provided in Section 5.6 and Section 5.7, respectively.

5. Proofs of main results. In what follows, we provide proofs of the results mentioned in the previous two sections.

5.1. *Proof of Proposition 3.1.* First assume $c > 0$. Substitute the expansion (3.1) into Golding's recursion (2.5), divide by ρc and let $\rho \rightarrow \infty$. We are then left with

$$(5.1) \quad q_0(\mathbf{a}, \mathbf{b}, \mathbf{c}) = \sum_{i=1}^K \sum_{j=1}^L \frac{c_{ij}}{c} q_0(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}).$$

Now, applying (5.1) repeatedly gives

$$q_0(\mathbf{a}, \mathbf{b}, \mathbf{c}) = \sum_{\text{orderings}} \frac{\prod_{(i,j) \in [K] \times [L]} c_{ij}!}{c!} q_0(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}),$$

where the summation is over all distinct orderings of the c gametes with multiplicity $\mathbf{c} = (c_{ij})$. There are $\frac{c!}{\prod_{(i,j)} c_{ij}!}$ such orderings and since the summand is independent of the ordering, we conclude

$$(5.2) \quad q_0(\mathbf{a}, \mathbf{b}, \mathbf{c}) = q_0(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}).$$

Clearly, (5.2) also holds for $c = 0$. From a coalescent perspective, this equation tells us that any gamete with specified alleles (i.e., "carrying ancestral material") at both loci must undergo recombination instantaneously backwards in time.

Now, by substituting the asymptotic expansion (3.1) with $\mathbf{c} = \mathbf{0}$ into Golding's recursion (2.5) and letting $\rho \rightarrow \infty$, we obtain

$$(5.3) \quad \begin{aligned} [n(n-1) + \theta_A a + \theta_B b] q_0(\mathbf{a}, \mathbf{b}, \mathbf{0}) = & \\ & \sum_{i=1}^K a_i(a_i - 1) q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \sum_{j=1}^L b_j(b_j - 1) q_0(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\ & + 2 \sum_{i=1}^K \sum_{j=1}^L a_i b_j q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) \\ & + \theta_A \sum_{i=1}^K \delta_{a_i,1} q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \theta_B \sum_{j=1}^L \delta_{b_j,1} q_0(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}). \end{aligned}$$

Equation (5.2) implies $q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) = q_0(\mathbf{a}, \mathbf{b}, \mathbf{0})$, so with a bit of rearranging we are left with

$$(5.4) \quad \begin{aligned} [a(a + \theta_A - 1) + b(b + \theta_B - 1)] q_0(\mathbf{a}, \mathbf{b}, \mathbf{0}) = & \\ & \sum_{i=1}^K a_i(a_i - 1) q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \sum_{j=1}^L b_j(b_j - 1) q_0(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\ & + \theta_A \sum_{i=1}^K \delta_{a_i,1} q_0(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \theta_B \sum_{j=1}^L \delta_{b_j,1} q_0(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}), \end{aligned}$$

with boundary conditions $q_0(\mathbf{e}_i, \mathbf{0}, \mathbf{0}) = q_0(\mathbf{0}, \mathbf{e}_j, \mathbf{0}) = 1$ for all $i \in [K]$ and $j \in [L]$. Noting that (5.4) is the sum of two independent recursions of the form (2.3), one for each locus and each with appropriate boundary condition, we conclude that $q_0(\mathbf{a}, \mathbf{b}, \mathbf{0})$ is given by

$$(5.5) \quad q_0(\mathbf{a}, \mathbf{b}, \mathbf{0}) = q^A(\mathbf{a})q^B(\mathbf{b}),$$

a product of two (ordered) ESFs. It is straightforward to verify that (5.5) satisfies (5.4). Finally, using (5.2) and (5.5), we arrive at (3.4). \square

5.2. *Proof of Theorem 3.1.* First assume $c > 0$. Substitute the asymptotic expansion (3.1) into Golding's recursion (2.5), eliminate terms of order ρ by applying (5.1), and let $\rho \rightarrow \infty$. After applying (5.2) to the remaining terms and invoking (5.4), with some rearrangement we obtain

$$\begin{aligned} cq_1(\mathbf{a}, \mathbf{b}, c) - \sum_{i=1}^K \sum_{j=1}^L c_{ij} q_1(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, c - e_{ij}) = \\ c(c-1)q_0(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}) - \sum_{i=1}^K c_i(c_i - 1)q_0(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i, \mathbf{b} + \mathbf{c}_B, \mathbf{0}) \\ - \sum_{j=1}^L c_j(c_j - 1)q_0(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B - \mathbf{e}_j, \mathbf{0}) \\ + \sum_{i=1}^K \sum_{j=1}^L c_{ij}(c_{ij} - 1)q_0(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i, \mathbf{b} + \mathbf{c}_B - \mathbf{e}_j, \mathbf{0}). \end{aligned}$$

Now, by utilizing (5.5), this can be written in the form

$$(5.6) \quad q_1(\mathbf{a}, \mathbf{b}, c) = f(\mathbf{a}, \mathbf{b}, c) + \sum_{i=1}^K \sum_{j=1}^L \frac{c_{ij}}{c} q_1(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, c - e_{ij}),$$

where

$$\begin{aligned} f(\mathbf{a}, \mathbf{b}, c) := & (c-1)q^A(\mathbf{a} + \mathbf{c}_A)q^B(\mathbf{b} + \mathbf{c}_B) \\ & - q^B(\mathbf{b} + \mathbf{c}_B) \sum_{i=1}^K \frac{c_i(c_i - 1)}{c} q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i) \\ & - q^A(\mathbf{a} + \mathbf{c}_A) \sum_{j=1}^L \frac{c_j(c_j - 1)}{c} q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j) \\ (5.7) \quad & + \sum_{i=1}^K \sum_{j=1}^L \frac{c_{ij}(c_{ij} - 1)}{c} q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i)q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j). \end{aligned}$$

Above we assumed $c > 0$. We define $f(\mathbf{a}, \mathbf{b}, \mathbf{c}) = 0$ if $\mathbf{c} = \mathbf{0}$. Iterating the recursion (5.6), we may write $q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})$ as

$$q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) = f(\mathbf{a}, \mathbf{b}, \mathbf{c}) + \sum_{i=1}^K \sum_{j=1}^L \frac{c_{ij}}{c} \left[f(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}) \right. \\ \left. + \sum_{i'=1}^K \sum_{j'=1}^L \frac{c_{i'j'} - \delta_{ii'} \delta_{jj'}}{c-1} q_1(\mathbf{a} + \mathbf{e}_i + \mathbf{e}_{i'}, \mathbf{b} + \mathbf{e}_j + \mathbf{e}_{j'}, \mathbf{c} - \mathbf{e}_{ij} - \mathbf{e}_{i'j'}) \right].$$

Similarly, repeatedly iterating (5.6) yields

$$\begin{aligned} q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) &= q_1(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}) \\ &\quad + f(\mathbf{a}, \mathbf{b}, \mathbf{c}) \\ &\quad + \sum_{i_1 j_1} \frac{c_{i_1 j_1}}{c} f(\mathbf{a} + \mathbf{e}_{i_1}, \mathbf{b} + \mathbf{e}_{j_1}, \mathbf{c} - \mathbf{e}_{i_1 j_1}) \\ &\quad + \sum_{i_1 j_1, i_2 j_2} \frac{c_{i_1 j_1} c_{i_2 j_2} - \delta_{i_1 j_1, i_2 j_2}}{c(c-1)} \\ &\quad \quad \times f(\mathbf{a} + \mathbf{e}_{i_1} + \mathbf{e}_{i_2}, \mathbf{b} + \mathbf{e}_{j_1} + \mathbf{e}_{j_2}, \mathbf{c} - \mathbf{e}_{i_1 j_1} - \mathbf{e}_{i_2 j_2}) \\ (5.8) \quad &+ \dots + \sum_{i_1 j_1, \dots, i_c j_c} \frac{\prod_{ij} c_{ij}!}{c!} f(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}). \end{aligned}$$

The key observation is that the right-hand side of (5.8) has a nice probabilistic interpretation which allows us to obtain a closed-form formula. To be more precise, consider the first summation

$$\sum_{i_1 j_1} \frac{c_{i_1 j_1}}{c} f(\mathbf{a} + \mathbf{e}_{i_1}, \mathbf{b} + \mathbf{e}_{j_1}, \mathbf{c} - \mathbf{e}_{i_1 j_1}).$$

For a fixed sample configuration \mathbf{c} , this can be interpreted as the sum over all possible ways of throwing away a gamete at random and calculating f based on the remaining subsample, which we will denote $\mathbf{c}^{(c-1)}$. Equivalently, it is the expected value of f with respect to subsampling without replacement $c-1$ of the gametes in \mathbf{c} . Write this as

$$\mathbb{E}[f(\mathbf{A}^{(c-1)}, \mathbf{B}^{(c-1)}, \mathbf{C}^{(c-1)})],$$

where $\mathbf{C}^{(c-1)}$ is the random subsample obtained by sampling without replacement $c-1$ gametes from \mathbf{c} , and $\mathbf{A}^{(c-1)} := \mathbf{a} + \mathbf{c}_A - \mathbf{C}_A^{(c-1)}$, $\mathbf{B}^{(c-1)} := \mathbf{b} + \mathbf{c}_B - \mathbf{C}_B^{(c-1)}$. Note that once the subsample $\mathbf{c}^{(c-1)}$ is obtained, then $\mathbf{a}^{(c-1)}$

and $\mathbf{b}^{(c-1)}$ are fully specified. More generally, consider the $(c-m)$ th sum in (5.8). A particular term in the summation corresponds to an ordering of $c-m$ gametes in \mathbf{c} , which, when removed leave a subsample $\mathbf{c}^{(m)}$. With respect to this subsample, the summand is

$$\prod_{i=1}^K \prod_{j=1}^L \frac{c_{ij}!}{c_{ij}^{(m)}!} \frac{m!}{c!} f(\mathbf{a}^{(m)}, \mathbf{b}^{(m)}, \mathbf{c}^{(m)}),$$

and for each such subsample $\mathbf{c}^{(m)}$ there are $\binom{c-m}{\mathbf{c}-\mathbf{c}^{(m)}}$ distinct orderings of the remaining types in \mathbf{c} , with each ordering contributing the same amount to the sum. Here, $\binom{c-m}{\mathbf{c}-\mathbf{c}^{(m)}}$ denotes the multinomial coefficient:

$$\binom{c-m}{\mathbf{c}-\mathbf{c}^{(m)}} = \frac{(c-m)!}{\prod_{i=1}^K \prod_{j=1}^L (c_{ij} - c_{ij}^{(m)})!}.$$

Gathering identical terms, the $(c-m)$ th sum in (5.8) can therefore be written over all distinct subsamples of \mathbf{c} of size m :

$$\begin{aligned} & \sum_{\mathbf{c}^{(m)}} \binom{c-m}{\mathbf{c}-\mathbf{c}^{(m)}} \prod_{i=1}^K \prod_{j=1}^L \frac{c_{ij}!}{c_{ij}^{(m)}!} \frac{m!}{c!} f(\mathbf{a}^{(m)}, \mathbf{b}^{(m)}, \mathbf{c}^{(m)}) \\ &= \sum_{\mathbf{c}^{(m)}} \frac{1}{\binom{c}{\mathbf{c}^{(m)}}} \prod_{i=1}^K \prod_{j=1}^L \binom{c_{ij}}{c_{ij}^{(m)}} f(\mathbf{a}^{(m)}, \mathbf{b}^{(m)}, \mathbf{c}^{(m)}) \\ &= \mathbb{E}[f(\mathbf{A}^{(m)}, \mathbf{B}^{(m)}, \mathbf{C}^{(m)})], \end{aligned}$$

where, for a fixed m , $\mathbf{C}^{(m)} = (C_{ij}^{(m)})$ is a multivariate hypergeometric(c, \mathbf{c}, m) random variable; that is,

$$\mathbb{P}\left(\bigcap_{(i,j) \in [K] \times [L]} [C_{ij}^{(m)} = c_{ij}^{(m)}]\right) = \frac{1}{\binom{c}{\mathbf{c}^{(m)}}} \prod_{(i,j) \in [K] \times [L]} \binom{c_{ij}}{c_{ij}^{(m)}}.$$

Furthermore, marginally we have

$$\begin{aligned} C_{ij}^{(m)} &\sim \text{hypergeometric}(c, c_{ij}, m), \\ C_{i\cdot}^{(m)} &\sim \text{hypergeometric}(c, c_{i\cdot}, m), \\ C_{\cdot j}^{(m)} &\sim \text{hypergeometric}(c, c_{\cdot j}, m). \end{aligned}$$

In summary, (5.8) can be written as

$$(5.9) \quad q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) = q_1(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0}) + \sum_{m=1}^c \mathbb{E}[f(\mathbf{A}^{(m)}, \mathbf{B}^{(m)}, \mathbf{C}^{(m)})].$$

According to Lemma 3.1, the first term $q_1(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0})$ vanishes, so we are left with

$$(5.10) \quad q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) = \sum_{m=1}^c \mathbb{E}[f(\mathbf{A}^{(m)}, \mathbf{B}^{(m)}, \mathbf{C}^{(m)})].$$

Finally, since $\mathbf{A}^{(m)} + \mathbf{C}_A^{(m)} = \mathbf{a} + \mathbf{c}_A$ and $\mathbf{B}^{(m)} + \mathbf{C}_B^{(m)} = \mathbf{b} + \mathbf{c}_B$, (5.7) and (5.10) together imply

$$\begin{aligned} q_1(\mathbf{a}, \mathbf{b}, \mathbf{c}) = & \sum_{m=1}^c \left[(m-1)q^A(\mathbf{a} + \mathbf{c}_A)q^B(\mathbf{b} + \mathbf{c}_B) \right. \\ & - q^B(\mathbf{b} + \mathbf{c}_B) \frac{1}{m} \sum_{i=1}^K \mathbb{E}[C_{i\cdot}^{(m)}(C_{i\cdot}^{(m)} - 1)] q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i) \\ & - q^A(\mathbf{a} + \mathbf{c}_A) \frac{1}{m} \sum_{j=1}^L \mathbb{E}[C_{\cdot j}^{(m)}(C_{\cdot j}^{(m)} - 1)] q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j) \\ & \left. + \frac{1}{m} \sum_{i=1}^K \sum_{j=1}^L \mathbb{E}[C_{ij}^{(m)}(C_{ij}^{(m)} - 1)] q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i) q^B(\mathbf{b} + \mathbf{c}_B - \mathbf{e}_j) \right]. \end{aligned}$$

The moments in this equation are easy to compute and one can sum them over m to obtain the desired result (3.5). \square

5.3. Proof of Lemma 3.1. First note that for any sample $(\mathbf{a}, \mathbf{b}, \mathbf{c})$ and any subsample of the form $(\mathbf{a}^{(1)}, \mathbf{b}^{(1)}, \mathbf{c}^{(1)})$, we have $f(\mathbf{a}^{(1)}, \mathbf{b}^{(1)}, \mathbf{c}^{(1)}) = 0$, since every term on right-hand side of (5.7) has a vanishing coefficient. So, equation (5.9) implies

$$(5.11) \quad q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) = q_1(\mathbf{a}, \mathbf{b}, \mathbf{0}),$$

for any $(i, j) \in [K] \times [L]$. Now, substitute the asymptotic expansion (3.1) with $\mathbf{c} = \mathbf{0}$ into Golding's recursion (2.5). Note that terms of order ρ are absent since $\mathbf{c} = \mathbf{0}$. Eliminate terms with coefficients independent of ρ by applying (5.3), multiply both sides of the recursion by ρ , and let $\rho \rightarrow \infty$ to

obtain the following:

$$\begin{aligned}
& [n(n-1) + \theta_A a + \theta_B b] q_1(\mathbf{a}, \mathbf{b}, \mathbf{0}) = \\
& \sum_{i=1}^K a_i(a_i - 1) q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \sum_{j=1}^L b_j(b_j - 1) q_1(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\
& + 2 \sum_{i=1}^K \sum_{j=1}^L a_i b_j q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) \\
& + \theta_A \sum_{i=1}^K \delta_{a_i,1} q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \theta_B \sum_{j=1}^L \delta_{b_j,1} q_1(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}),
\end{aligned}$$

with boundary conditions $q_1(\mathbf{e}_i, \mathbf{0}, \mathbf{0}) = q_1(\mathbf{0}, \mathbf{e}_j, \mathbf{0}) = 0$ for all $i \in [K]$ and $j \in [L]$. This equation can be made strictly recursive by applying (5.11) to $q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij})$. It therefore follows from the boundary conditions (for example by induction) that $q_1(\mathbf{a}, \mathbf{b}, \mathbf{0}) = 0$. \square

5.4. *Proof of Theorem 3.2.* Here, we provide only an outline of a proof; details are similar to the proof of Theorem 3.1. Substitute the asymptotic expansion (3.1) into Golding's recursion (2.5), eliminate terms with coefficients proportional to ρ or independent of ρ . Then, multiply both sides of the recursion by ρ and let $\rho \rightarrow \infty$ to obtain

$$\begin{aligned}
& c q_2(\mathbf{a}, \mathbf{b}, \mathbf{c}) - \sum_{i=1}^K \sum_{j=1}^L c_{ij} q_2(\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}) = \\
& \sum_{i=1}^K a_i(a_i - 1 + 2c_i) q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{c}) + \sum_{j=1}^L b_j(b_j - 1 + 2c_j) q_1(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{c}) \\
& + \sum_{i=1}^K \sum_{j=1}^L [c_{ij}(c_{ij} - 1) q_1(\mathbf{a}, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij}) + 2a_i b_j q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{c} + \mathbf{e}_{ij})] \\
& + \theta_A \sum_{i=1}^K \left[\sum_{j=1}^L \delta_{a_i+c_i,1} \delta_{c_{ij},1} q_1(\mathbf{a}, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}) + \delta_{a_i,1} \delta_{c_i,0} q_1(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{c}) \right] \\
& + \theta_B \sum_{j=1}^L \left[\sum_{i=1}^K \delta_{b_j+c_j,1} \delta_{c_{ij},1} q_1(\mathbf{a} + \mathbf{e}_i, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij}) + \delta_{b_j,1} \delta_{c_j,0} q_1(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{c}) \right] \\
& - [n(n-1) + \theta_A(a+c) + \theta_B(b+c)] q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})
\end{aligned} \tag{5.12}$$

By substituting our expression (3.5) for $q_1(\mathbf{a}, \mathbf{b}, \mathbf{c})$, the right-hand side can be expressed as a function $g(\mathbf{a}, \mathbf{b}, \mathbf{c})$ which is completely known but rather

cumbersome to write down. As in the proof of Theorem 3.1, the same ‘unwrapping’ maneuver can be applied to rearrange (5.12) into the form (3.6), where

$$\sigma(\mathbf{a}, \mathbf{b}, \mathbf{c}) = \sum_{m=1}^c \mathbb{E}[g(\mathbf{A}^{(m)}, \mathbf{B}^{(m)}, \mathbf{C}^{(m)})].$$

This time $\mathbb{E}[g(\mathbf{A}^{(m)}, \mathbf{B}^{(m)}, \mathbf{C}^{(m)})]$ is a function of fourth-order moments of the multivariate hypergeometric distribution. The formula shown in Appendix A is obtained by evaluating the expectations and summing over m .

We now show that $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$ satisfies the recursion shown in (3.7). We will use the fact that for a sample $(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij})$, we have

$$\begin{aligned} g(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) &= 2(a-1)(b-1)q^A(\mathbf{a})q^B(\mathbf{b}) \\ &\quad - 2(b-1)(a_i-1)q^A(\mathbf{a} - \mathbf{e}_i)q^B(\mathbf{b}) \\ &\quad - 2(a-1)(b_j-1)q^A(\mathbf{a})q^B(\mathbf{b} - \mathbf{e}_j) \\ &\quad + 2(a_i-1)(b_j-1)q^A(\mathbf{a} - \mathbf{e}_i)q^B(\mathbf{b} - \mathbf{e}_j). \end{aligned} \tag{5.13}$$

For an arbitrary \mathbf{c} , $g(\mathbf{a}, \mathbf{b}, \mathbf{c})$ is much more complicated.

Now, one can adopt the approach used in the proof of Lemma 3.1 to obtain a strict recursion for $q_2(\mathbf{a} + \mathbf{c}_A, \mathbf{b} + \mathbf{c}_B, \mathbf{0})$. First, note that (3.6) and (5.13) imply

$$\begin{aligned} q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) &= q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) + \mathbb{E}[g((\mathbf{A} - \mathbf{e}_i)^{(1)}, (\mathbf{B} - \mathbf{e}_j)^{(1)}, \mathbf{e}_{ij}^{(1)})] \\ &= q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) + g(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) \\ &= q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) + 2(a-1)(b-1)q^A(\mathbf{a})q^B(\mathbf{b}) \\ &\quad - 2(b-1)(a_i-1)q^A(\mathbf{a} - \mathbf{e}_i)q^B(\mathbf{b}) \\ &\quad - 2(a-1)(b_j-1)q^A(\mathbf{a})q^B(\mathbf{b} - \mathbf{e}_j) \\ &\quad + 2(a_i-1)(b_j-1)q^A(\mathbf{a} - \mathbf{e}_i)q^B(\mathbf{b} - \mathbf{e}_j). \end{aligned} \tag{5.14}$$

As before, substitute the asymptotic expansion (3.1) for $\mathbf{c} = \mathbf{0}$ into Golding’s recursion (2.5), eliminate terms with coefficients independent of ρ or proportional to ρ^{-1} , and let $\rho \rightarrow \infty$ to obtain

$$\begin{aligned} [n(n-1) + \theta_A a + \theta_B b]q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) &= \\ &\sum_{i=1}^K a_i(a_i-1)q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \sum_{j=1}^L b_j(b_j-1)q_2(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}) \\ &\quad + 2 \sum_{i=1}^K \sum_{j=1}^L a_i b_j q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij}) \\ &\quad + \theta_A \sum_{i=1}^K \delta_{a_i,1} q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{0}) + \theta_B \sum_{j=1}^L \delta_{b_j,1} q_2(\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{0}), \end{aligned}$$

with boundary conditions $q_2(\mathbf{e}_i, \mathbf{0}, \mathbf{0}) = q_2(\mathbf{0}, \mathbf{e}_j, \mathbf{0}) = 0$ for all $i \in [K]$ and $j \in [L]$. This equation can be made strictly recursive by applying (5.14) to $q_2(\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{e}_{ij})$. After some simplification, this leads to the recursion (3.7). \square

5.5. *Proof of Proposition 4.1.* The proof is similar to the proof of Proposition 3.1, working with the system (4.1) rather than Golding's recursion (2.5). First assume $c > 0$. Substitute the expansion (4.4) into the recursion (4.1), divide by ρc and let $\rho \rightarrow \infty$. We are left with

$$p_0(a, b, c; k, l) = p_0(a + 1, b + 1, c - 1; k, l),$$

which implies

$$(5.15) \quad p_0(a, b, c; k, l) = p_0(a + c, b + c, 0; k, l).$$

Clearly, (5.15) also holds for $c = 0$.

Now, by substituting the asymptotic expansion (4.4) with $c = 0$ into (4.1) and letting $\rho \rightarrow \infty$, we obtain

$$(5.16) \quad \begin{aligned} & [n(n-1) + \theta_A a + \theta_B b] p_0(a, b, 0; k, l) = \\ & a(a-1)p_0(a-1, b, 0; k, l) + b(b-1)p_0(a, b-1, 0; k, l) \\ & + 2abp_0(a-1, b-1, 1; k, l) \\ & + \theta_A a p_0(a-1, b, 0; k-1, l) + \theta_B b p_0(a, b-1, 0; k, l-1). \end{aligned}$$

After invoking (5.15) on $p_0(a-1, b-1, 1; k, l)$ and rearranging, we are left with

$$(5.17) \quad \begin{aligned} & [a(a + \theta_A - 1) + b(b + \theta_B - 1)] p_0(a, b, 0; k, l) = \\ & a(a-1)p_0(a-1, b, 0; k, l) + b(b-1)p_0(a, b-1, 0; k, l) \\ & + \theta_A a p_0(a-1, b, 0; k-1, l) + \theta_B b p_0(a, b-1, 0; k, l-1), \end{aligned}$$

with boundary conditions $p_0(1, 0, 0; k, l) = \delta_{k,1}\delta_{l,0}$, and $p_0(0, 1, 0; k, l) = \delta_{k,0}\delta_{l,1}$. Equation (5.17) can be expressed as a linear sum of two independent recursions:

$$\begin{aligned} (a-1 + \theta_A)p_0^A(a; k) &= (a-1)p_0^A(a-1; k) + \theta_A p_0^A(a-1; k-1), \\ (b-1 + \theta_B)p_0^B(b; l) &= (b-1)p_0^B(b-1; l) + \theta_B p_0^B(b-1; l-1), \end{aligned}$$

with respective boundary conditions $p_0^A(1; k) = \delta_{k,1}$ and $p_0^B(1; l) = \delta_{l,1}$. These recursions are precisely those considered by Ewens (1972, eq. 21), with respective solutions (4.2) and (4.3). Hence, $p_0^A(a; k) = p^A(a; k)$ and $p_0^B(b; l) = p^B(b; l)$, and it is straightforward to verify that $p^A(a; k)p^B(b; l)$ satisfies (5.17). Substituting this solution into (5.15), we arrive at (4.5), as required. \square

5.6. *Proof of Proposition 4.2.* First assume $c > 0$. Substitute the asymptotic expansion (4.4) into the recursion (4.1), eliminate terms with coefficients linear in ρ by applying (5.15), and let $\rho \rightarrow \infty$. After applying (5.15) to the remaining terms and invoking (5.17), with some rearrangement we obtain

$$(5.18) \quad \begin{aligned} p_1(a, b, c; k, l) - p_1(a + 1, b + 1, c - 1; k, l) = \\ (c - 1)[p_0(a + c, b + c, 0; k, l) - p_0(a + c - 1, b + c, 0; k, l) \\ - p_0(a + c, b + c - 1, 0; k, l) + p_0(a + c - 1, b + c - 1, 0; k, l)]. \end{aligned}$$

Applying the recursion repeatedly, this becomes

$$(5.19) \quad \begin{aligned} p_1(a, b, c; k, l) = & p_1(a + c, b + c, 0; k, l) + [p_0(a + c, b + c, 0; k, l) \\ & - p_0(a + c - 1, b + c, 0; k, l) - p_0(a + c, b + c - 1, 0; k, l) \\ & + p_0(a + c - 1, b + c - 1, 0; k, l)] \sum_{m=0}^{c-1} (c - 1 - m). \end{aligned}$$

According to Lemma 4.1, the first term $p_1(a + c, b + c, 0; k, l)$ vanishes. Hence, since $p_0(a, b, c; k, l)$ is given by (4.5), the right-hand side of (5.19) is fully known. With some rearrangement we are left with (4.6). \square

5.7. *Proof of Lemma 4.1.* First note that (5.18) implies

$$(5.20) \quad p_1(a - 1, b - 1, 1; k, l) = p_1(a, b, 0; k, l).$$

Now, substitute the asymptotic expansion (4.4) with $c = 0$ into (4.1), eliminate leading order terms by applying (5.16), and let $\rho \rightarrow \infty$. The result is made strictly recursive by invoking (5.20), and we obtain

$$\begin{aligned} [a(a + \theta_A - 1) + b(b + \theta_B - 1)]p_1(a, b, 0; k, l) = \\ a(a - 1)p_1(a - 1, b, 0; k, l) + b(b - 1)p_1(a, b - 1, 0; k, l) \\ + \theta_A a p_1(a - 1, b, 0; k - 1, l) + \theta_B b p_1(a, b - 1, 0; k, l - 1), \end{aligned}$$

with boundary conditions $p_1(1, 0, 0; k, l) = p_1(0, 1, 0; k, l) = 0$, for $k, l = 0, \dots, n$. It therefore follows (for example by induction) that $p_1(a, b, 0; k, l) = 0$. \square

Acknowledgments. We thank Robert C. Griffiths, Charles H. Langley, and Joshua Paul for useful discussions.

APPENDIX A: EXPRESSION FOR $\sigma(\mathbf{a}, \mathbf{b}, \mathbf{c})$

We use Q^A to denote $q^A(\mathbf{a} + \mathbf{c}_A)$, Q_i^A to denote $q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i)$, Q_{ik}^A to denote $q^A(\mathbf{a} + \mathbf{c}_A - \mathbf{e}_i - \mathbf{e}_k)$, and so on. Then,

$$\begin{aligned}
\sigma(\mathbf{a}, \mathbf{b}, \mathbf{c}) = & \frac{c}{3} \left[\frac{(c-1)(c+1)(3c-2)}{8} + (c-1)(3a+3b+2c-1) + 6ab \right] Q^A Q^B \\
& - \frac{\theta_A(c-1)}{2} \sum_{i=1}^K \delta_{a_i,0} \delta_{c_i,1} Q_i^A Q^B - \frac{\theta_B(c-1)}{2} \sum_{j=1}^L \delta_{b_j,0} \delta_{c_j,1} Q^A Q_j^B \\
& + \sum_{i=1}^K \left[\frac{\theta_A - c(c-3) + 2a + 4b - 4}{4} c_i(c_i - 1) \right. \\
& \quad \left. - (2b + c - 1)c_i(a_i + c_i - 1) \right] Q_i^A Q^B \\
& + \frac{1}{2} \sum_{i=1}^K \left[\frac{\theta_A}{2} \delta_{c_i,2} + \frac{5 - 6a_i - 4c_i}{6} \right] c_i(c_i - 1) Q_{ii}^A Q^B \\
& + \sum_{j=1}^L \left[\frac{\theta_B - c(c-3) + 2b + 4a - 4}{4} c_j(c_j - 1) \right. \\
& \quad \left. - (2a + c - 1)c_j(b_j + c_j - 1) \right] Q^A Q_j^B \\
& + \frac{1}{2} \sum_{j=1}^L \left[\frac{\theta_B}{2} \delta_{c_j,2} + \frac{5 - 6b_j - 4c_j}{6} \right] c_j(c_j - 1) Q^A Q_{jj}^B \\
& + \sum_{i,k=1}^K \frac{c_i(c_i - 1)c_k(c_k - 1)}{8} Q_{ik}^A Q^B + \sum_{j,l=1}^L \frac{c_j(c_j - 1)c_l(c_l - 1)}{8} Q^A Q_{jl}^B \\
& - \frac{\theta_A + \theta_B - c(c-5) + 2a + 2b - 4}{4} \sum_{i=1}^K \sum_{j=1}^L c_{ij}(c_{ij} - 1) Q_i^A Q_j^B \\
& + \sum_{i=1}^K \sum_{j=1}^L \left[\frac{c_i(c_i - 1)c_j(c_j - 1)}{4} + \frac{c_{ij}(c_{ij} + 1 - 2c_i + 2c_i c_j - 2c_j)}{2} \right. \\
& \quad \left. + c_{ij}b_j(c_i - 1) + c_{ij}a_i(c_j - 1) + 2a_i b_j c_{ij} \right. \\
& \quad \left. + \frac{\theta_B}{2} \delta_{b_j,0} \delta_{c_j,1} \delta_{c_{ij},1} (c_i - 1) + \frac{\theta_A}{2} \delta_{a_i,0} \delta_{c_i,1} \delta_{c_{ij},1} (c_j - 1) \right] Q_i^A Q_j^B
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2} \sum_{i=1}^K \left[a_i + c_i - 1 - \frac{\theta_A}{2} \delta_{a_i+c_i,2} \right] \sum_{j=1}^L c_{ij} (c_{ij} - 1) Q_{ii}^A Q_j^B \\
& + \frac{1}{2} \sum_{j=1}^L \left[b_j + c_{.j} - 1 - \frac{\theta_B}{2} \delta_{b_j+c_{.j},2} \right] \sum_{i=1}^K c_{ij} (c_{ij} - 1) Q_i^A Q_{jj}^B \\
& - \frac{1}{4} \sum_{i=1}^K \sum_{j=1}^L \sum_{k=1}^K c_{ij} (c_{ij} - 1) c_{k.} (c_{k.} - 1) Q_{ik}^A Q_j^B \\
& - \frac{1}{4} \sum_{i=1}^K \sum_{j=1}^L \sum_{l=1}^L c_{ij} (c_{ij} - 1) c_{.l} (c_{.l} - 1) Q_i^A Q_{jl}^B \\
& + \frac{1}{8} \sum_{i=1}^K \sum_{j=1}^L \sum_{k=1}^K \sum_{l=1}^L c_{ij} (c_{ij} - 1) c_{kl} (c_{kl} - 1) Q_{ik}^A Q_{jl}^B \\
& - \frac{1}{12} \sum_{i=1}^K \sum_{j=1}^L c_{ij} (c_{ij} - 1) (2c_{ij} - 1) Q_{ii}^A Q_{jj}^B.
\end{aligned}$$

To check the correctness of the above expression, we also solved the recursion (5.12) numerically for all sample configurations of sizes $n = 10, 20$, and 30 (with $K, L \leq 2$), and confirmed that the above analytic expression agreed in all cases. We also implemented a *Mathematica* program to solve $q(a, b, c)$ exactly in the special case $K = L = 1$. The program can return series expansions in terms of ρ^{-1} as $\rho \rightarrow \infty$ which are symbolic in θ_A and θ_B . We could then compare the first three terms against q_0, q_1 , and q_2 , for various sample configurations (a, b, c) .

APPENDIX B: AN UPPER BOUND ON $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$

To obtain an upper bound on $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$, we can exploit the linearity in (3.7). Let

$$\begin{aligned}
s_1(\mathbf{a}, \mathbf{b}) & := ab\theta_A\theta_B, \\
s_2(\mathbf{a}, \mathbf{b}) & := -b\theta_B(\theta_A + a - 1) \sum_i \delta_{a_i,1}, \\
s_3(\mathbf{a}, \mathbf{b}) & := -a\theta_A(\theta_B + b - 1) \sum_j \delta_{b_j,1}, \\
\text{(B.1)} \quad s_4(\mathbf{a}, \mathbf{b}) & := (\theta_A + a - 1)(\theta_B + b - 1) \sum_i \delta_{a_i,1} \sum_j \delta_{b_j,1}.
\end{aligned}$$

Now note that if $r_k(\mathbf{a}, \mathbf{b})$ ($k = 1, \dots, 4$) satisfies

$$(B.2) \quad [a(a + \theta_A - 1) + b(b + \theta_B - 1)]r_k(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^K a_i(a_i - 1 + \theta_A \delta_{a_i,1})r_k(\mathbf{a} - \mathbf{e}_i, \mathbf{b}) + \sum_{j=1}^L b_j(b_j - 1 + \theta_B \delta_{b_j,1})r_k(\mathbf{a}, \mathbf{b} - \mathbf{e}_j) + 4q^A(\mathbf{a})q^B(\mathbf{b})s_k(\mathbf{a}, \mathbf{b}),$$

then $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) = r_1(\mathbf{a}, \mathbf{b}) + r_2(\mathbf{a}, \mathbf{b}) + r_3(\mathbf{a}, \mathbf{b}) + r_4(\mathbf{a}, \mathbf{b})$. We can then bound each $r_k(\mathbf{a}, \mathbf{b})$ individually and sum them for bounds on $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$.

By iteratively unwrapping (B.2), each term that contributes to $r_k(\mathbf{a}, \mathbf{b})$ can be mapped bijectively to a sequence of subsamples of (\mathbf{a}, \mathbf{b}) with the following constraints:

1. The first term in the sequence is (\mathbf{a}, \mathbf{b}) ,
2. The m th term in the sequence has total sample size $a + b + 1 - m$,
3. The sequence has length m such that $1 \leq m \leq a + b$.

For example, the sequence $((\mathbf{a}, \mathbf{b}))$ is associated with the contribution

$$\frac{4q^A(\mathbf{a})q^B(\mathbf{b})s_k(\mathbf{a}, \mathbf{b})}{a(a + \theta_A - 1) + b(b + \theta_B - 1)},$$

while the sequence $((\mathbf{a}, \mathbf{b}), (\mathbf{a} - \mathbf{e}_i, \mathbf{b}))$ is associated with the contribution

$$\frac{a_i(a_i - 1 + \theta_A \delta_{a_i,1})}{[a(a + \theta_A - 1) + b(b + \theta_B - 1)]} \cdot \frac{4q^A(\mathbf{a} - \mathbf{e}_i)q^B(\mathbf{b})s_k(\mathbf{a} - \mathbf{e}_i, \mathbf{b})}{(a - 1)(a + \theta_A - 2) + b(b + \theta_B - 1)},$$

and so on. Call a sequence satisfying the constraints 1–3 a *path*. Denote a subsample of (\mathbf{a}, \mathbf{b}) of size (m_A, m_B) by $(\mathbf{a}^{(m_A)}, \mathbf{b}^{(m_B)})$. To see how to obtain a bound on $r_k(\mathbf{a}, \mathbf{b})$, let us consider the contribution from a single example path,

$$\begin{aligned} & ((\mathbf{a}, \mathbf{b}), (\mathbf{a}^{(a-1)}, \mathbf{b}), (\mathbf{a}^{(a-1)}, \mathbf{b}^{(b-1)}), (\mathbf{a}^{(a-2)}, \mathbf{b}^{(b-1)})) \\ & := ((\mathbf{a}, \mathbf{b}), (\mathbf{a} - \mathbf{e}_i, \mathbf{b}), (\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j), (\mathbf{a} - 2\mathbf{e}_i, \mathbf{b} - \mathbf{e}_j)). \end{aligned}$$

If $a_i = 2$ and $b_j > 1$, the contribution from this path to $r_k(\mathbf{a}, \mathbf{b})$ is

$$(B.3) \quad \left[\frac{a_i(a_i - 1)}{D(a, b)} \cdot \frac{b_j(b_j - 1)}{D(a - 1, b)} \cdot \frac{\theta_A}{D(a - 1, b - 1)} \right] \times \frac{4q^A(\mathbf{a}^{(a-2)})q^B(\mathbf{b}^{(b-1)})}{D(a - 2, b - 1)} s_k(\mathbf{a}^{(a-2)}, \mathbf{b}^{(b-1)}),$$

where $D(a, b) := a(a-1+\theta_A) + b(b-1+\theta_B)$. Evaluating $r_k(\mathbf{a}, \mathbf{b})$ completely is problematic because of terms like

$$[D(a, b)D(a-1, b)D(a-1, b-1)D(a-2, b-1)]^{-1},$$

which differ for each path. Suppose for now we have obtained an upper bound $U_{a-2, b-1}$ on this term which depends only on the subsample size $(a-2, b-1)$ for the final term in the path, and not on the complete path. Then the contribution (B.3) can be rewritten

$$\begin{aligned} & 4 \left[\frac{a(a-1+\theta_A)}{D(a, b)} \cdot \frac{b(b-1+\theta_B)}{D(a-1, b)} \cdot \frac{(a-1)(a-2+\theta_A)}{D(a-1, b-1)} \cdot \frac{1}{D(a-2, b-1)} \right] \\ & \quad \times \left[\frac{a_i}{a} \frac{b_j}{b} \frac{1}{a-1} \right] q^A(\mathbf{a}) q^B(\mathbf{b}) s_k(\mathbf{a}^{(a-2)}, \mathbf{b}^{(b-1)}) \\ \text{(B.4)} \quad & \leq 4U_{a-2, b-1} \frac{a!}{(a-2)!} \frac{b!}{(b-1)!} \frac{(\theta_A)_a}{(\theta_A)_{a-2}} \frac{(\theta_B)_b}{(\theta_B)_{b-1}} \\ & \quad \times \left[\frac{a_i}{a} \frac{b_j}{b} \frac{1}{a-1} \right] q^A(\mathbf{a}) q^B(\mathbf{b}) s_k(\mathbf{a}^{(a-2)}, \mathbf{b}^{(b-1)}). \end{aligned}$$

The reason for this approach is that one can now sum (B.4):

- (I) First, over all paths to this subsample, then
- (II) over all paths ending at a subsample of size $(a-2, b-1)$, and finally
- (III) over all subsample sizes.

There are $\binom{3}{2} = 3$ paths to this subsample. Hence, summing the bound (B.4) over (I), a bound on the total contribution of paths to this subsample is

$$\begin{aligned} \text{(B.5)} \quad & 4U_{a-2, b-1} \frac{a!}{(a-2)!} \frac{b!}{(b-1)!} \frac{(\theta_A)_a}{(\theta_A)_{a-2}} \frac{(\theta_B)_b}{(\theta_B)_{b-1}} \\ & \times \binom{3}{2} \mathbb{P}(\mathbf{A}^{(a-2)} = \mathbf{a}^{(a-2)}) \mathbb{P}(\mathbf{B}^{(b-1)} = \mathbf{b}^{(b-1)}) q^A(\mathbf{a}) q^B(\mathbf{b}) s_k(\mathbf{a}^{(a-2)}, \mathbf{b}^{(b-1)}), \end{aligned}$$

where the random variable $\mathbf{A}^{(m_A)}$ is distributed according to a multivariate hypergeometric(a, \mathbf{a}, m_A) distribution, and the random variable $\mathbf{B}^{(m_B)}$ is distributed according to a multivariate hypergeometric(b, \mathbf{b}, m_B) distribution.

Next, sum this bound over (II). Thus, a bound on the contribution from

all paths ending at a subsample of size $(a-2, b-1)$ is

$$4U_{a-2, b-1} \frac{a!}{(a-2)!} \frac{b!}{(b-1)!} \frac{(\theta_A)_a}{(\theta_A)_{a-2}} \frac{(\theta_B)_b}{(\theta_B)_{b-1}} \times \binom{3}{2} q^A(\mathbf{a}) q^B(\mathbf{b}) \mathbb{E}[s_k(\mathbf{A}^{(a-2)}, \mathbf{B}^{(b-1)})].$$

More generally, one can follow the same argument to find a bound on the total contribution from paths ending at a subsample of size (m_A, m_B) :

$$(B.6) \quad 4U_{m_A, m_B} \binom{a+b-m_A-m_B}{a-m_A} \frac{a!}{m_A!} \frac{b!}{m_B!} \frac{(\theta_A)_a}{(\theta_A)_{m_A}} \frac{(\theta_B)_b}{(\theta_B)_{m_B}} \times q^A(\mathbf{a}) q^B(\mathbf{b}) \mathbb{E}[s_k(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})].$$

The binomial coefficient $\binom{a+b-m_A-m_B}{a-m_A}$ accounts for the number of ways of interspersing the sequences $(\mathbf{a}, \mathbf{a}^{(a-1)}, \dots, \mathbf{a}^{(m_A)})$ and $(\mathbf{b}, \mathbf{b}^{(b-1)}, \dots, \mathbf{b}^{(m_B)})$. Referring to (B.1), the expectations in (B.6) can be evaluated:

$$(B.7) \quad \begin{aligned} \mathbb{E}[s_1(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})] &= m_A m_B \theta_A \theta_B, \\ \mathbb{E}[s_2(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})] &= -m_B \theta_B (\theta_A + m_A - 1) \sum_i \frac{a_i \binom{a-a_i}{m_A-1}}{\binom{a}{m_A}}, \\ \mathbb{E}[s_3(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})] &= -m_A \theta_A (\theta_B + m_B - 1) \sum_j \frac{b_j \binom{b-b_j}{m_B-1}}{\binom{b}{m_B}}, \\ \mathbb{E}[s_4(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})] &= (\theta_A + m_A - 1)(\theta_B + m_B - 1) \\ &\quad \times \sum_i \frac{a_i \binom{a-a_i}{m_A-1}}{\binom{a}{m_A}} \sum_j \frac{b_j \binom{b-b_j}{m_B-1}}{\binom{b}{m_B}}. \end{aligned}$$

Summing (B.6) over all subsample sizes and over $k = 1, \dots, 4$, we have therefore proven the following:

PROPOSITION B.1.

$$q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) \leq 4q^A(\mathbf{a}) q^B(\mathbf{b}) \sum_{m_A=2}^a \sum_{m_B=2}^b (a+b-m_A-m_B) \binom{a}{m_A} \binom{b}{m_B} \times \frac{(\theta_A)_a}{(\theta_A)_{m_A}} \frac{(\theta_B)_b}{(\theta_B)_{m_B}} U_{m_A, m_B} \sum_{k=1}^4 \mathbb{E}[s_k(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})],$$

where $\mathbb{E}[s_k(\mathbf{A}^{(m_A)}, \mathbf{B}^{(m_B)})]$ is given by (B.7) for $k = 1, \dots, 4$.

In a similar fashion,

$$q_2(\mathbf{a}, \mathbf{b}, \mathbf{0}) \geq 4q^A(\mathbf{a})q^B(\mathbf{b}) \sum_{m_A=2}^a \sum_{m_B=2}^b (a+b-m_A-m_B)! \binom{a}{m_A} \binom{b}{m_B} \\ \times \frac{(\theta_A)_a}{(\theta_A)_{m_A}} \frac{(\theta_B)_b}{(\theta_B)_{m_B}} L_{m_A, m_B} \sum_{k=1}^4 \mathbb{E}[s_k(\mathbf{A}^{m_A}, \mathbf{B}^{m_B})],$$

where L_{m_A, m_B} is a lower bound on all coefficients of the form

$$[D(a, b) \dots D(m_A, m_B)]^{-1},$$

corresponding to paths through the recursion from (\mathbf{a}, \mathbf{b}) to any subsample $(\mathbf{a}^{(m_A)}, \mathbf{b}^{(m_B)})$ of size (m_A, m_B) .

It remains to choose U_{m_A, m_B} . The closest possible bound is given by the actual path which maximizes this term, and this can be found by constructing an $a \times b$ dynamic programming table as follows. Define

$$U_{m_A, m_B} = \begin{cases} \frac{1}{D(a, b)} & \text{if } (m_A, m_B) = (a, b), \\ \frac{U_{m_A+1, b}}{D(m_A, b)} & \text{if } m_B = b, \\ \frac{U_{a, m_B+1}}{D(a, m_B)} & \text{if } m_A = a, \\ \frac{\max\{U_{m_A+1, m_B}, U_{m_A, m_B+1}\}}{D(m_A, m_B)} & \text{if } 2 \leq m_A \leq a-1, \text{ and } 2 \leq m_B \leq b-1, \end{cases}$$

and similarly for L_{m_A, m_B} (replacing max with min). Note that constructing this single table finds all the relevant bounds. In fact, U_{m_A, m_B} provides upper bounds for $r_1(\mathbf{a}, \mathbf{b})$ and $r_4(\mathbf{a}, \mathbf{b})$; and L_{m_A, m_B} provides upper bounds for $r_2(\mathbf{a}, \mathbf{b})$ and $r_3(\mathbf{a}, \mathbf{b})$, since these terms are less than or equal to zero.

The bounds presented above require $O(ab)$ operations, which, while growing with a and b , require much less time than a complete evaluation of $q_2(\mathbf{a}, \mathbf{b}, \mathbf{0})$. Empirical testing of this upper bound suggest that it is generally very close. We also investigated some constant time upper bounds as well as defining lower bounds, but these proved to be less useful.

REFERENCES

- Arratia, A., Barbour, A. D., and Tavaré, S. *Logarithmic Combinatorial Structures: A Probabilistic Approach*. European Mathematical Society Publishing House, Switzerland, 2003.

- De Iorio, M. and Griffiths, R. C. (2004a). Importance sampling on coalescent histories I. *Adv. Appl. Prob.*, **36**, 417–433.
- De Iorio, M. and Griffiths, R. C. (2004b). Importance sampling on coalescent histories II. *Adv. Appl. Prob.*, **36**, 434–454.
- Ethier, S. N. and Griffiths, R. C. (1990). On the two-locus sampling distribution. *J. Math. Biol.*, **29**, 131–159.
- Ewens, W. J. (1972). The sampling theory of selectively neutral alleles. *Theor. Popul. Biol.*, **3**, 87–112.
- Fearnhead, P. and Donnelly, P. (2001). Estimating recombination rates from population genetic data. *Genetics*, **159**, 1299–1318.
- Golding, G. B. (1984). The sampling distribution of linkage disequilibrium. *Genetics*, **108**, 257–274.
- Griffiths, R. C. (1981). Neutral two-locus multiple allele models with recombination. *Theor. Popul. Biol.*, **19**, 169–186.
- Griffiths, R. C. The two-locus ancestral graph. In Basawa, I. V. and Taylor, R. L., editors, *Selected proceedings of the Sheffield symposium on applied probability: 18. IMS Lecture Notes—Monograph series*, volume 18, pages 100–117, 1991.
- Griffiths, R. C. and Lessard, S. (2005). Ewens’ sampling formula and related formulae: combinatorial proofs, extensions to variable population size and applications to ages of alleles. *Theor. Popul. Biol.*, **68**, 167–77.
- Griffiths, R. C. and Marjoram, P. (1996). Ancestral inference from samples of DNA sequences with recombination. *J. Comput. Biol.*, **3**, 479–502.
- Griffiths, R. C., Jenkins, P. A., and Song, Y. S. (2008). Importance sampling and the two-locus model with subdivided population structure. *Adv. Appl. Prob.*, **40**, 473–500.
- Hoppe, F. (1984). Pólya-like urns and the Ewens’ sampling formula. *J. Math. Biol.*, **20**, 91–94.
- Hudson, R. R. (1985). The sampling distribution of linkage disequilibrium under an infinite allele model without selection. *Genetics*, **109**, 611–631.
- Hudson, R. R. (2001). Two-locus sampling distributions and their application. *Genetics*, **159**, 1805–1817.
- Kingman, J. F. C. (1982a). On the genealogy of large populations. *J. Appl. Prob.*, **19**, 27–43.
- Kingman, J. F. C. (1982b). The coalescent. *Stoch. Process. Appl.*, **13**, 235–248.
- Kuhner, M. K., Yamato, J., and Felsenstein, J. (2000). Maximum likelihood estimation of recombination rates from population data. *Genetics*, **156**, 1393–1401.
- McVean, G. A. T., Myers, S., Hunt, S., Deloukas, P., Bentley, D. R., and Donnelly, P. (2004). The fine-scale structure of recombination rate variation in the human genome. *Science*, **304**, 581–584.
- Myers, S., Bottolo, L., Freeman, C., McVean, G., and Donnelly, P. (2005). A fine-scale map of recombination rates and hotspots across the human genome. *Science*, **310**, 321–324.
- Nielsen, R. (2000). Estimation of population parameters and recombination rates from single nucleotide polymorphisms. *Genetics*, **154**, 931–942.
- Pitman, J. The two-parameter generalization of Ewens’ random partition structure. Technical Report 345, Dept. Statistics, U.C. Berkeley, 1992.
- Pitman, J. (1995). Exchangeable and partially exchangeable random partitions. *Probab. Th. Rel. Fields*, **102**, 145–158.
- Slatkin, M. (1994). An exact test for neutrality based on the Ewens sampling distribution. *Genet. Res.*, **64**, 71–4.
- Stephens, M. Inference under the coalescent. In Balding, D., Bishop, M., and Cannings,

- C., editors, *Handbook of Statistical Genetics*, chapter 8. Wiley, Chichester, UK, 2001.
- Stephens, M. and Donnelly, P. (2000). Inference in molecular population genetics. *J. R. Statist. Soc. B*, **62**, 605–655.
- Wang, Y. and Rannala, B. (2008). Bayesian inference of fine-scale recombination rates using population genomic data. *Phil. Trans. R. Soc. B*, **363**, 3921–3930.
- Watterson, G. A. (1977). Heterosis or neutrality? *Genetics*, **85**, 789–814.

P. A. JENKINS
COMPUTER SCIENCE DIVISION
UNIVERSITY OF CALIFORNIA, BERKELEY
BERKELEY, CA 94720
USA
E-MAIL: pauljenk@eecs.berkeley.edu

Y. S. SONG
DEPARTMENT OF STATISTICS AND
COMPUTER SCIENCE DIVISION
UNIVERSITY OF CALIFORNIA, BERKELEY
BERKELEY, CA 94720
USA
E-MAIL: yss@stat.berkeley.edu