

Counting All Possible Ancestral Configurations of Sample Sequences in Population Genetics

Yun S. Song, Rune Lyngsø, and Jotun Hein

Abstract—Given a set D of input sequences, a genealogy for D can be constructed backward in time using such evolutionary events as mutation, coalescent, and recombination. An ancestral configuration (AC) can be regarded as the multiset of all sequences present at a particular point in time in a possible genealogy for D . The complexity of computing the likelihood of observing D depends heavily on the total number of distinct ACs of D and, therefore, it is of interest to estimate that number. For D consisting of binary sequences of finite length, we consider the problem of enumerating exactly all distinct ACs. We assume that the root sequence type is known and that the mutation process is governed by the infinite-sites model. When there is no recombination, we construct a general method of obtaining closed-form formulas for the total number of ACs. The enumeration problem becomes much more complicated when recombination is involved. In that case, we devise a method of enumeration based on counting contingency tables and construct a dynamic programming algorithm for the approach. Last, we describe a method of counting the number of ACs that can appear in genealogies with less than or equal to a given number R of recombinations. Of particular interest is the case in which R is close to the minimum number of recombinations for D .

Index Terms—Ancestral configurations, coalescent, recombination, contingency table, enumeration.

1 INTRODUCTION

ONE of the standard problems in mathematical population genetics is to compute the likelihood $P(D)$ of observing a given data set D under the assumed model of evolution. The likelihood $P(D)$ can be defined as a formal sum over all possible genealogies consistent with D . A genealogy, in turn, can be viewed as a sequence of Markov states and the likelihood $P(D)$ can be determined, in principle, by summing products of transition probabilities over all possible sequences of states in a Markov chain with given initial and final states. There are recursion relations for computing $P(D)$ exactly [1], [2], [15], but that approach quickly becomes infeasible as data size grows and one must then resort to Monte Carlo methods [3], [5], [6], [7], [11], [12], [16]. Given a set D of input sequences, a genealogy for D can be constructed backward in time using such evolutionary events as mutation, coalescent, and recombination. An ancestral configuration (AC) can be regarded as the multiset of all sequences present at a particular point in time in a possible genealogy for D . For both deterministic and stochastic approaches, the efficiency of a method of computing $P(D)$ depends heavily on the total number of distinct ACs of D .

In this paper, we address the problem of enumerating all distinct ACs for a given data set D , consisting of binary sequences of finite length. That is, we are interested in computing the total number of distinct ACs in all possible

genealogies that could have generated D . If an AC appears in two or more different genealogies, it is counted only once. We consider both the classic coalescent [9], [10], in which recombination is absent, and its extension, where recombination is allowed. When a mutation event occurs in a genealogy, then it always leads to a new AC in the genealogy, but a coalescent or a recombination event may lead to an AC that has already been encountered in the genealogy.

In the absence of recombination, genealogies can be represented by time-ordered binary trees, whereas a case with recombination requires a more complicated graphical representation, called the ancestral recombination graph (ARG) [5]. It is well known that the coalescent with recombination is considerably more difficult to study than the classic coalescent. One obvious reason for that contrast is that there are many more inequivalent ARGs than trees. As many computations of interest—computing $P(D)$, for example—involve studying a set of genealogies consistent with the observed data, perhaps it is not so surprising that including recombination in the model of evolution poses many challenges. One of our goals here is to illustrate more precisely why recombination is difficult to study by comparing the total number of ACs when there is recombination with that when it is absent. It has been known to many, if not most, people in population genetics that there can be many more ACs when there is recombination than when it is absent, but we are not aware of any work that tried to examine exactly by how much.

As mentioned above, when there are many ACs, solving recursion relations to compute $P(D)$ exactly is often infeasible. In such a case, we are interested in asking whether it would be possible to approximate the true recursion by a “collapsed” recursion in which ACs are lumped together such that we just need to consider transitions, with appropriately defined probability, between

- Y.S. Song is with the Department of Computer Science, University of California at Davis, 2063 Kemper Hall, One Shields Avenue, Davis, CA 95616. E-mail: yssong@cs.ucdavis.edu.
- R. Lyngsø and J. Hein are with the Department of Statistics, Oxford University, 1 South Parks Road, Oxford OX1 3TG, UK. E-mail: {lyngsoe, hein}@stats.ox.ac.uk.

Manuscript received 17 May 2005; revised 27 Dec. 2005; accepted 16 Jan. 2006; published online 31 July 2006.

For information on obtaining reprints of this article, please send e-mail to: tcbb@computer.org, and reference IEEECS Log Number TCBB-0048-0505.

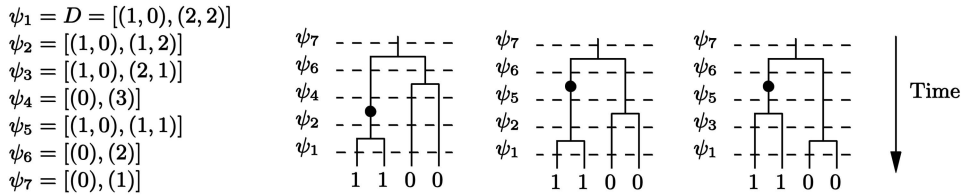


Fig. 1. Coalescent trees and ancestral configurations ψ_1, \dots, ψ_7 . A cross section, corresponding to a particular time slice, of a tree defines an AC.

lumped ACs to estimate $P(D)$. We show in this paper that this idea of lumping ACs works at least in the context of our enumeration problem. What we achieve here can be viewed as a small step toward meeting our desired goal.

In our work, we assume the infinite sites model of mutation, which means that there can be at most one mutation per site. Hence, the input data can be regarded as defining a set of binary sequences. For ease of discussion, we assume that the root sequence is known and use 0 to denote the ancestral type at the root and 1 to denote a mutant type. We remark that the techniques described in this paper can be generalized to the root-unknown case. When there is no recombination, we construct a general method of obtaining closed-form formulas for the total number of ACs; these formulas are polynomials in the multiplicity of distinct sequences in D . We have implemented this method in *Mathematica* and the program is called *aceTrees* (short for “ancestral configuration enumeration for trees”). As expected, the enumeration problem becomes much more complicated when recombination is involved. In that case, we discuss a method based on counting nonnegative integer valued tables with fixed row and column sums, commonly known as contingency tables. For two sites, we show how a closed-form formula can be obtained. For an arbitrary number of sites, we construct a dynamic programming algorithm for counting ACs. We have implemented this algorithm in C++; the software is called *aceARGs* (short for “ancestral configuration enumeration for ARGs”). Last, we discuss a method of counting the number of ACs that can appear in ARGs with less than or equal to a given number R of recombinations. Of particular interest is the case in which R is (or is close to) the minimum number of recombinations for D . We have implemented this algorithm in C; the software is called *greven*. This method interpolates between the case where there is no recombination and the case where an arbitrary number of recombinations are allowed. As it should, *greven* agrees with *aceTrees* and *aceARGs* in those two respective limits. As a further check, we have made alternative, independent implementations of all three programs in *Python*. All our programs implementing the methods discussed in this paper are available upon request.

The organization of this paper is as follows: In Section 2, we consider the problem of enumerating ACs when there is no recombination, in which case, genealogies can be represented by coalescent trees. Mitochondria data from [17] are considered there as an example. By counting the total number of ACs, we illustrate why an exact computation of the likelihood is difficult for that data. The aforementioned method of counting ACs in unconstrained recombination graphs is discussed in Section 3. We

compare the number of ACs in coalescent trees with that in ARGs and highlight their differences. In Section 4, we consider enumerating ACs that can appear in ARGs with at most R recombinations. We conclude in Section 5 with some general remarks on our work.

2 ANCESTRAL CONFIGURATIONS IN COALESCENT TREES

In this section, we focus on the case in which the evolution under the infinite-sites model can be represented by a tree. A given data set D is compatible with a tree if it passes the three gamete test: For every pair of sites, not all of the allele types 01, 10, and 11 appear in the data. If the three gamete test is passed, then there exists a unique perfect phylogeny τ for D [4], [8], with at most one mutation per site, but it is important to note that, in general, there are many coalescent trees consistent with τ . (Recall that, unlike coalescent trees, a perfect phylogeny may be nonbinary and that relative time ordering of two of its interior vertices is not defined if one vertex is not a descendant of the other.)

2.1 Definition of an Ancestral Configuration

In our algorithm for counting ACs, we partition the information contained in a configuration into two parts, one encoding sequence types and the other their multiplicity. We use x_1, \dots, x_d to denote d distinct finite binary sequences (allele types) x_i of some fixed length—for example, $x_1 = 0100$, $x_2 = 0010$, and $x_3 = 1100$, with d being 3. We define $T = (x_1, \dots, x_d)$ and $\mathbf{n} = (n_1, \dots, n_d)$, where $n_i > 0$, for all $i \in \{1, \dots, d\}$, denote their multiplicity. We say that (T, \mathbf{n}) and (T', \mathbf{n}') are equivalent, denoted $(T, \mathbf{n}) \sim (T', \mathbf{n}')$, if there exists a permutation $\sigma \in S_d$ such that $(T_\sigma, \mathbf{n}_\sigma) = (T', \mathbf{n}')$, where $T_\sigma = (x_{\sigma(1)}, \dots, x_{\sigma(d)})$ and $\mathbf{n}_\sigma = (n_{\sigma(1)}, \dots, n_{\sigma(d)})$. By a *configuration*, we mean an equivalence class $[T, \mathbf{n}] \in \{(T, \mathbf{n})\} / \sim$.

As we describe presently, every coalescent tree defines a sequence of configurations. Consider a data set $D = [T, \mathbf{n}]$, where $T = (x_1, \dots, x_d)$ and $\mathbf{n} = (n_1, \dots, n_d)$. A configuration is said to be *ancestral* to D if it “appears” in at least one possible coalescent tree that derives D . (Note that D is ancestral to itself according to this definition.) Let us illustrate this definition through a specific example. Consider the simple example $D = [(1,0), (2,2)]$, consisting of four binary sequences of length one. There are several coalescent trees that can give rise to D under the infinite-sites model of mutation. Three of them are shown in Fig. 1. A cross section, corresponding to a particular time slice, of a coalescent tree defines an AC. Going backward in time, configuration changes whenever either a mutation or a coalescent event occurs.

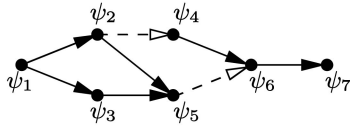


Fig. 2. Markov chain transition graph C_D for $D = [(1, 0), (2, 2)]$. Dashed lines denote mutation events, while solid lines denote coalescent events. This directed graph is acyclic and contains three distinct paths from the initial configuration $\psi_1 = D$ to the absorbing configuration ψ_7 .

For $D = [(1, 0), (2, 2)]$, the seven configurations shown in Fig. 1 form a complete set of configurations ancestral to D , i.e., any configuration appearing in a coalescent tree for D must be one of the seven configurations. The corresponding Markov chain transition graph C_D , where vertices correspond to the ACs for D and edges correspond to allowed transitions, is as shown in Fig. 2. This graph shows that the sequence of ACs arising in any coalescent tree for D must be one of the following:

$$\begin{aligned} \psi_1 &\rightarrow \psi_2 \rightarrow \psi_4 \rightarrow \psi_6 \rightarrow \psi_7, \\ \psi_1 &\rightarrow \psi_2 \rightarrow \psi_5 \rightarrow \psi_6 \rightarrow \psi_7, \\ \psi_1 &\rightarrow \psi_3 \rightarrow \psi_5 \rightarrow \psi_6 \rightarrow \psi_7. \end{aligned}$$

In general, every coalescent tree compatible with any given data set D starts in the configuration corresponding to D and ends in the configuration consisting of a single all-zero sequence (corresponding to the most recent common ancestor).

Let $P_{i,j}$ denote some appropriately defined transition probability corresponding to the transition $\psi_i \rightarrow \psi_j$. Then, the probability of the sequence $\psi_{i_1} \rightarrow \psi_{i_2} \rightarrow \dots \rightarrow \psi_{i_k}$ can be defined as $P_{i_1,i_2} P_{i_2,i_3} \dots P_{i_{k-1},i_k}$. In [1], Ethier and Griffiths showed that the probability of observing D can be obtained by summing such probabilities over all possible sequences of ACs. In the example $D = [(1, 0), (2, 2)]$ considered above,

$$\begin{aligned} P(D) &= P_{1,2} P_{2,4} P_{4,6} P_{6,7} + P_{1,2} P_{2,5} P_{5,6} P_{6,7} \\ &\quad + P_{1,3} P_{3,5} P_{5,6} P_{6,7}. \end{aligned}$$

Ethier and Griffiths [1] formulated a recursion for evaluating $P(D)$. In general, computing $P(D)$ exactly using the recursion becomes infeasible when sample size is large and one must then resort to Monte Carlo methods (see [6], [7]). In what follows, we develop a method of counting exactly the total number of inequivalent ACs of an arbitrary data set D . This kind of enumeration should prove useful for studying when exact computation of $P(D)$ becomes infeasible.

2.2 Enumeration of ACs: A General Solution

We now describe a general, efficient method of obtaining closed-form formulas for the total number of inequivalent ACs. The actual ACs themselves can easily be extracted from our method.

2.2.1 Notation

For ease of discussion, we first introduce some useful notation. We use $\mathbf{1}_k$ to denote a k -tuple of 1s, and \mathbf{e}_i to denote a vector with a 1 at the i th entry and 0s elsewhere; the length of \mathbf{e}_i will be clear from the context of its usage. For $\mathbf{n} = (n_1, \dots, n_k)$, define $\pi(\mathbf{n}) := \prod_{i=1}^k n_i$. Let \mathcal{D}_i denote a

deletion operator which, when acting on a vector of length $k \geq i$, deletes the i th entry of the vector, thus changing its length to $k - 1$. Let $\mathcal{R}_{i,z}$ denote a replacement operator which replaces the i th entry of a vector with z .

Let $\mathcal{T} = (\mathbf{y}_1, \dots, \mathbf{y}_k)$ and $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)$. As usual, $\mathbf{y}_1, \dots, \mathbf{y}_k$ are k distinct binary sequences. In \mathcal{T} , a site is called a *singlet* if there is exactly one allele type $\mathbf{y}_i \in \{\mathbf{y}_1, \dots, \mathbf{y}_k\}$ with value 1 at that site. Here note that μ_i need not be 1. In $\mathcal{T} = (100, 101, 010)$, for example, sites two and three are singlets. Let a be a singlet site and \mathbf{y}_i the unique allele type with a 1 at that site. With \mathbf{y}_i^a denoting the allele type obtained from \mathbf{y}_i after the 1 at site a mutates to a 0, define $\mathcal{M}_a(\mathcal{T}, \boldsymbol{\mu})$ as

$$\begin{cases} (\mathcal{R}_{i,\mathbf{y}_i^a}(\mathcal{T}), \mathcal{R}_{i,1}(\boldsymbol{\mu})), & \text{if } \mathbf{y}_i^a \notin \{\mathbf{y}_1, \dots, \mathbf{y}_k\}, \\ (\mathcal{D}_i(\mathcal{T}), \mathcal{D}_i(\boldsymbol{\mu} + \mathbf{e}_j)), & \text{if } \mathbf{y}_i^a = \mathbf{y}_j \in \{\mathbf{y}_1, \dots, \mathbf{y}_k\}. \end{cases}$$

Since a singlet site a uniquely determines i , we omit i in writing $\mathcal{M}_a(\mathcal{T}, \boldsymbol{\mu})$. In the example $\mathcal{T} = (100, 101, 010)$ considered above, the unique allele type associated with $a = 2$ is \mathbf{y}_3 , and $\mathbf{y}_3^2 = (000)$ is not in $\{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3\}$. For $a = 3$, the associated unique allele type is \mathbf{y}_2 and $\mathbf{y}_2^3 = (100)$, which is equal to \mathbf{y}_1 . Hence, $\mathcal{M}_2(\mathcal{T}, \boldsymbol{\mu}) = ((100, 101, 000), (\mu_1, \mu_2, 1))$ and $\mathcal{M}_3(\mathcal{T}, \boldsymbol{\mu}) = ((100, 010), (\mu_1 + 1, \mu_3))$.

2.2.2 Coalescent and Mutation Events

Consider a data set $D = [T, \mathbf{n}]$. The multiplicity n_i of an allele type \mathbf{x}_i decreases by exactly one when two sequences of that allele type coalesce. The number of ACs with allele types T is therefore equal to $\pi(\mathbf{n})$, which is equal to 1 (for $[T, \mathbf{n}]$ itself) plus the number of configurations that can be reached from $[T, \mathbf{n}]$ via coalescent events only. When every allele type has coalesced completely, we end up with the configuration $[T, \mathbf{1}_d]$.

In $T = (\mathbf{x}_1, \dots, \mathbf{x}_d)$, suppose that \mathbf{x}_i is the only allele type with a 1 at site a . Consider a configuration $[T, \mathbf{m}]$ with $m_i = 1$, reached by some coalescent events starting from $[T, \mathbf{n}]$. Now, a mutation event can occur to change the character of \mathbf{x}_i at site a from 1 to 0, making $[T, \mathbf{m}]$ jump to a different configuration $[T', \mathbf{m}']$, where $T' \neq T$. Let \mathbf{x}_i^a denote the allele type obtained after such a mutation event. There are two possible cases. First, if \mathbf{x}_i^a is equal to one of the alleles in T , say \mathbf{x}_j , then the configuration after the mutation event is $[\mathcal{D}_i(T), \mathcal{D}_i(\mathbf{m} + \mathbf{e}_j)]$. Note that the multiplicity of \mathbf{x}_j has increased by one. How many configurations are there with allele types $\mathcal{D}_i(T)$? For $k \neq i$ and $k \neq j$, the multiplicity of \mathbf{x}_k can be any integer between n_k and 1; independently of the mutation event, coalescent events within allele type \mathbf{x}_k can reduce the multiplicity from n_k to 1. Further, since the multiplicity of \mathbf{x}_j increases by 1 after the mutation event, it can range from $n_j + 1$ to 1. Hence, the total number of possible ACs with allele types $\mathcal{D}_i(T)$ is $\pi(\mathcal{D}_i(\mathbf{n} + \mathbf{e}_j))$. Second, if \mathbf{x}_i^a is not equal to any allele in T , then the configuration after the mutation event is $[\mathcal{R}_{i,\mathbf{x}_i^a}(T), \mathbf{m}]$. Note that we still have $m_i = 1$. How many configurations are there with allele types $\mathcal{R}_{i,\mathbf{x}_i^a}(T)$? Similarly to the first case, since \mathbf{m} is bounded by $\mathcal{R}_{i,1}(\mathbf{n})$ and $\mathbf{1}_d$, the number of configurations with allele types $\mathcal{R}_{i,\mathbf{x}_i^a}(T)$ is $\pi(\mathcal{R}_{i,1}(\mathbf{n}))$.

$$\begin{aligned}
v_1 &= ((100, 110, 001), (n_1, n_2, n_3)) \\
v_2 &= ((100, 001), (n_1 + 1, n_3)) \\
v_3 &= ((100, 110, 000), (n_1, n_2, 1)) \\
v_4 &= ((100, 000), (n_1 + 1, 1)) \\
v_5 &= ((000, 001), (1, n_3)) \\
v_6 &= ((000), (2))
\end{aligned}$$

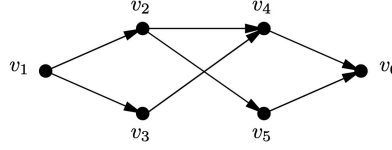


Fig. 3. Application of our enumeration algorithm to $D = [(100, 110, 001), (n_1, n_2, n_3)]$. Shown on the right-hand side is the final graph G_D constructed by our algorithm.

2.2.3 A Graph Construction Algorithm

Our method of counting ACs is to deal with coalescent events for each allele type configuration \mathcal{T} combinatorially and find all possible allele type configurations using a simple graph construction algorithm. More precisely, the general idea goes as follows: For a given input data set D , suppose that there are ACs with \mathcal{T} as allele type configuration. If the maximum multiplicity of \mathcal{T} over all such ACs is μ , then, as discussed above, a simple coalescent argument shows that the total number of ACs with allele type configuration \mathcal{T} is $\pi(\mu)$. Hence, if we know all possible allele type configurations $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_p$ and their maximum multiplicity $\mu_1, \mu_2, \dots, \mu_p$, respectively, then we can easily determine the total number of ACs of D by summing over $\pi(\mu_j)$. Our graph construction algorithm below finds such $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_p$ and $\mu_1, \mu_2, \dots, \mu_p$.

We wish to construct a directed graph G_D , where a vertex v_j is labeled by (T_j, n_j) , to be determined by the following iterative procedure:

1. Given a data set $D = [T, n]$, let G_D initially be a graph with no edges and with $v_1 = (T, n)$ as its only vertex. Here, (T, n) is an arbitrary representative of the equivalence class $[T, n]$.
2. Let V_0 denote the set of all vertices in G_D with out-degree zero. For all $v_j = (T_j, n_j) \in V_0$, determine the set S_{v_j} of singlet sites (defined above) in T_j .
3. If $S_{v_j} = \emptyset$ for all $v_j = (T_j, n_j) \in V_0$, terminate the procedure. Otherwise, arbitrarily order V_0 and sequentially carry out the following steps for v_j satisfying $S_{v_j} \neq \emptyset$: Determine $\mathcal{M}_a(T_j, n_j)$ for all $a \in S_{v_j}$. If $\mathcal{M}_a(T_j, n_j) = v_k \in G_D$, draw a directed edge from v_j to v_k . If not, then add $v_k = \mathcal{M}_a(T_j, n_j)$ to G_D and draw a directed edge from v_j to v_k .
4. Go back to Step 2.

The graph G_D compactly encodes all possible ACs. In a vertex $v_j = (T_j, n_j) \in G_D$, T_j captures the set of distinct binary strings, whereas n_j determines the range of possible multiplicity of the strings, i.e., if there are r distinct binary sequences in T_j , then the multiplicity can be anything between n_j and 1_r . Further, it is straightforward to show that if v_j and v_k are two distinct vertices in G_D , then $T_j \neq T_k$. The total number $\alpha(D)$ of configurations ancestral to D is therefore given by

$$\alpha(D) = \sum_{i \in V(G_D)} \pi(n_i), \quad (1)$$

where $V(G_D)$ is the index set of the vertices in G_D . We have implemented the above algorithm in *Mathematica*. For a given data set $D = [(x_1, \dots, x_d), (n_1, \dots, n_d)]$, our program

generates closed-form formulas for $\alpha(D)$ in terms of n_1, \dots, n_d .

Remark. We have an independent program, to be discussed later, that can exhaustively search through evolutionary histories to compute the number of ACs. That program can be used to analyze small sample sizes (usually less than 30 sequences). We have checked that it produces the same answers as does our combinatorial method described above.

2.3 A Toy Example

Consider the data set $D = [(100, 110, 001), (n_1, n_2, n_3)]$, where n_1, n_2 , and n_3 are some arbitrary positive integers. Applying the algorithm from Section 2.2 leads to the graph G_D illustrated in Fig. 3. From (1), the total number of configurations ancestral to D is thus given by

$$\begin{aligned}
\alpha(D) &= \sum_{i=1}^6 \pi(n_i) \\
&= n_1 n_2 n_3 + n_3(n_1 + 1) + n_1 n_2 + (n_1 + 1) + n_3 + 2 \\
&= n_1(n_2 + 1)(n_3 + 1) + 2n_3 + 3.
\end{aligned}$$

2.4 Mitochondria DNA Data

We now consider the data set from Ward et al. [17], which consists of sequences from the control region of mitochondria DNA (mtDNA), sampled from 63 individuals in a single Amerindian tribe. There are 18 segregating sites in the data, shown in Table 1.

Computing the likelihood of observing a given data set is important—for example, for parameter estimation—but, as mentioned in Section 2.1, doing it exactly is infeasible when the sample size is large. The same mtDNA data set was considered by Griffiths and Tavaré [6], who developed a Markov chain Monte Carlo method to analyze the data. Our goal here is to compute the total number of inequivalent ACs of the data, thus illustrating why computing the likelihood exactly using recursions is difficult.

Let $T = (x_1, \dots, x_{14})$ denote the 14 distinct haplotypes shown in Table 1, with x_i being the i th row. Applying our algorithm from Section 2.2 to $D = [T, (n_1, \dots, n_{14})]$ generates a directed graph G_D with 12,896 vertices and one can easily obtain a closed-form formula, which is too long to write down here, for the total number $\alpha(D)$ of ACs. For $n = (1, 1, \dots, 1)$, $\alpha(D) = 128,640$. For $n = (2, 2, 1, 3, 19, 1, 1, 1, 4, 8, 5, 4, 3, 1)$, which is the multiplicity of the actual data, $\alpha(D) = 413,243,616$.

TABLE 1
Mitochondria Data from [17], Consisting of
63 Sequences and 18 Segregating Sites

Haplotype	Multiplicity
100101000000010000	2
100101000100010000	2
010000000000000001	1
001000001000000000	3
000101000000010000	19
000111000000010000	1
000000000011000001	1
000000000011100001	1
0000000000000000111	4
0000000000000000101	8
0000000000000000000	5
0000000000000000001	4
0000000100000000000	3
0000001000000001000	1

There are 413, 243, 616 ACs for this data set.

3 ANCESTRAL CONFIGURATIONS IN UNCONSTRAINED ARGs

In this section, we turn to enumerating ancestral configurations in unconstrained ancestral recombination graphs (ARGs) [5]. This case is much more complicated than the classic coalescent case and, in general, it is difficult to obtain closed-form formulas for the number of ACs. Below, we translate the problem of counting ACs into counting contingency tables and provide a dynamic programming algorithm.

3.1 Definition of an AC in ARGs

In the absence of recombination, if the input data D contains n sequences, then any configuration ancestral to D contains at most n sequences. This is no longer true when recombinations are allowed. Going backward in time, when a recombination event is encountered, the lineage of a sequence breaks up into two parts, distributing its ancestral material to two different sequences which carry additional nonancestral material (denoted by “*”s). This concept is illustrated in Fig. 4, where possible coalescent events are also described. To simplify things, we assume that recombination breakpoints occur at the midpoints of consecutive sites in D . Hence, if D consists of some segregating sites in a region X , the distribution of ancestral material between two consecutive segregating sites i and $i + 1$ in X is completely determined by the configuration at sites i and $i + 1$.

Of particular interest is a class of ARGs in which, for any recombination event, both the prefix and the suffix involved in a recombination event contain some ancestral material. Such a class of ARGs was the main focus of [5] and we also limit our attention to such ARGs in this paper. Hence, for our purpose, a sequence appearing in an ARG is generally a string over $\{0, 1, *\}$, but the all-* string is not allowed. As before, 0 denotes the ancestral type at the root and 1 denotes

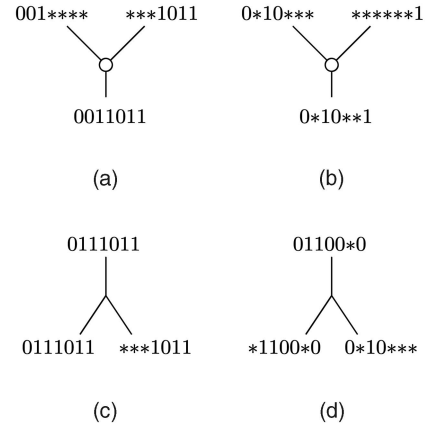


Fig. 4. Illustration of possible events. Time flows from top to bottom. (a) and (b) are examples of recombination events. Let $s[i : j]$ denote the substring of s in between i and j , inclusive. Going backward in time, when a recombination event is encountered, the lineage of a sequence $s[1 : L]$ breaks up into two parts, one corresponding to the lineage of a sequence s_1 whose prefix $s_1[1 : k]$ is identical to $s[1 : k]$ and the other corresponding to the lineage of a second sequence s_2 whose suffix $s_2[k + 1 : L]$ is identical to $s[k + 1 : L]$. The suffix $s_1[k + 1 : L]$ and the prefix $s_2[1 : k]$ carry nonancestral material, denoted by “*”s. (c) and (d) are examples of coalescent events. Two sequences s_1 and s_2 can coalesce if their characters are identical at common ancestral positions. The final sequence contains the union of the ancestral material in s_1 and s_2 .

a mutant type. A configuration is defined as in Section 2.1, except that now strings are defined over $\{0, 1, *\}$. To every ARG, there corresponds a sequence of configurations which can be viewed as generating the ARG backward in time. Given a data set D , a configuration is said to be *ancestral* to D if it appears in at least one possible sequence of configurations that corresponds to some ARG consistent with D . Shown in Fig. 5 is a sequence of ACs and its corresponding ARG. It is important to note that there are many—in fact, infinitely many—distinct ARGs that can derive the same initial data and that the total number of ACs can be immensely large. This is the main point that we wish to illustrate in this paper. For the simple example $D = [(10, 11, 01), (1, 1, 2)]$ considered in Fig. 5, there are only 220 ACs in total. We shall soon see that, as the number of sites and the number of sequences increase, the total number of ACs grows extremely fast.

3.2 A Warm-Up Example

Before we plunge into the core of our enumeration work, let us consider a simple example for which it is not too difficult to obtain a complete set of ACs by hand. The reader not too familiar with the ARG is recommended to go through this example. Consider the initial data $D = [(00, 11), (1, 1)]$.

If recombinations are not allowed, then it is clear that there are only five ACs, namely, $[(00, 11), (1, 1)]$, $[(00, 01), (1, 1)]$, $[(00, 10), (1, 1)]$, $[(00), (2)]$, and $[(00), (1)]$. In the presence of recombination, there are 30 ACs, shown in Fig. 6. These ACs can be generated iteratively by asking what happens to a given AC under a mutation, a coalescent, or a recombination event; note that some of these events may not be possible, depending on the AC. On the right-hand side of Fig. 6 is a Markov chain transition graph C_D , depicting the possible transitions between ACs. Note that the graph C_D

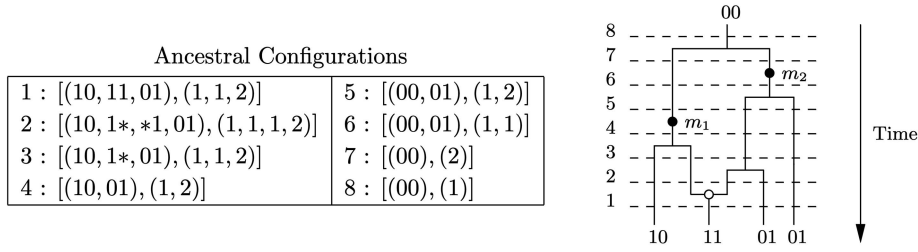


Fig. 5. A sequence of ACs for $D = [(10, 11, 01), (1, 1, 2)]$ and its corresponding ARG. Filled circles denote mutation events. An open circle denotes a recombination event with breakpoint between the first and the second sites. Mutation events at the first and the second sites are denoted by m_1 and m_2 , respectively. Note that there are other ARGs that could have generated D . In total, there are 220 ACs of D .

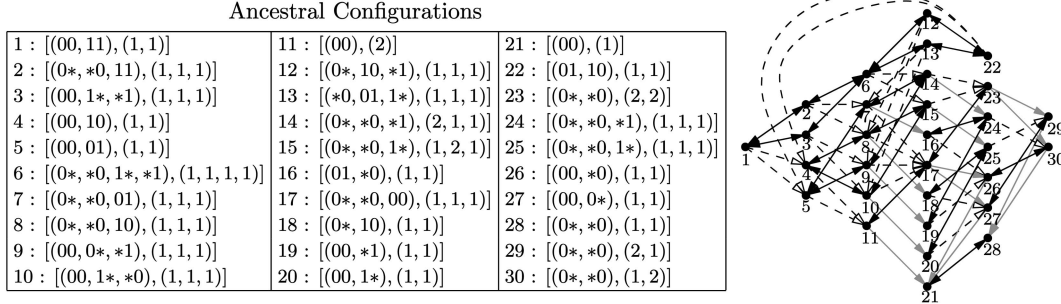


Fig. 6. Configurations ancestral to $D = [(00, 11), (1, 1)]$ and the Markov chain transition graph C_D . Mutation events are denoted by dashed arrows. Coalescent and recombination events are denoted by solid arrows. Note that some solid arrows are bidirectional, one direction for recombination and the other for coalescent. Unidirectional arrows denote coalescent events.

contains directed cycles. AC 1 is the initial configuration and AC 21 is called the grand common ancestor. To any ARG with the grand common ancestor as the root, there corresponds a unique path from AC 1 to AC 21.

3.3 Possible Number of 1s and 0s in an AC

Let n denote the total number of binary sequences in D , with each sequence being of length L . Let n_0^i (respectively, n_1^i) denote the total number of 0s (respectively, 1s) at site i . Note that $n_0^i + n_1^i = n$ for all $i \in \{1, \dots, L\}$. (For $D = [(10, 01, 11), (2, 3, 1)]$, $n_0^1 = 3$, $n_1^1 = 3$ for the first site and $n_0^2 = 2$, $n_1^2 = 4$ for the second site.) Let r_0^i and r_1^i denote the total number of 0s and of 1s, respectively, at site i in an arbitrary AC ψ . The infinite-sites model of mutation imposes constraints on possible values of (r_0^i, r_1^i) . For $n_1^i \neq 0$, a mutation event at site i can occur in AC ψ only if $r_1^i = 1$. If $r_1^i = 1$ and the mutation event occurs, then AC ψ jumps to a new AC ψ' with $(r_0^i + 1)$ 0s and no 1s.

To determine the possible range of r_0^i and r_1^i , let us first consider the $L = 1$ case. Since every sequence is of length one, it follows from the above discussion that there exists *exactly* one AC for any pair $(r_0^1, r_1^1) \in I(n_0^1, n_1^1)$, where $I(n_0, n_1)$ is defined as the set of all nonnegative integer pairs (p, q) such that if $n_0 \neq 0$ and $n_1 \neq 0$, then

$$\begin{aligned} &\text{either (i) } 1 \leq p \leq n_0 \text{ and } 1 \leq q \leq n_1, \\ &\text{or (ii) } 1 \leq p \leq n_0 + 1 \text{ and } q = 0; \end{aligned} \quad (2)$$

if $n_0 = 0$ and $n_1 = n$, then

$$\begin{aligned} &\text{either (i) } p = 0 \text{ and } 1 \leq q \leq n, \\ &\text{or (ii) } p = 1 \text{ and } q = 0; \end{aligned} \quad (3)$$

and, if $n_0 = n$ and $n_1 = 0$, then

$$1 \leq p \leq n \text{ and } q = 0. \quad (4)$$

There are $n_0(n_1 + 1) + 1$ possible values of (p, q) in the first case, $n + 1$ in the second, and n in the last. There exists no AC with r_0^1 0s and r_1^1 1s if $(r_0^1, r_1^1) \notin I(n_0^1, n_1^1)$. Now comes a crucial point. Even for $L > 1$, if an arbitrary number of recombination events are allowed, there exists *at least one* AC for any combination of $(r_0^1, r_1^1), \dots, (r_0^L, r_1^L)$ satisfying $(r_0^i, r_1^i) \in I(n_0^i, n_1^i)$, for all $i \in \{1, \dots, L\}$. There exists no AC with r_0^i 0s and r_1^i 1s at site i if $(r_0^i, r_1^i) \notin I(n_0^i, n_1^i)$. One way to see these points is as follows: Starting from the input data D , use recombination events to distribute the ancestral material of the first sequence to L new sequences, each carrying exactly one ancestral site. Do the same for all other sequences in D so that we end up with $n \cdot L$ sequences, each carrying exactly one ancestral site. Then, determining the range of r_0^i and r_1^i for each site i reduces to the $L = 1$ case. We stress that the above statement would not hold for $L > 1$ in the absence of recombination (or if too few recombination events are allowed). For example, consider $D = [(01), (2)]$. In the absence of recombination, there is no AC with $r_0^1 = 1$ and $r_1^1 = 2$.

3.4 Important Observations

We now highlight a couple of important facts. First, the set of ACs does not depend on the order of 0s and 1s within a site when arbitrary number of recombinations are allowed. For instance, $[(01, 10), (1, 1)]$ has the same set of ACs as the example $[(00, 11), (1, 1)]$ considered in Section 3.2 (see Fig. 6). This fact can easily be explained as follows: Let D and D' be two data sets of the same size such that the i th site of D is equal to the i th site of D' up to some rearrangement of elements within the site. Then, D can be

$\begin{bmatrix} 00 \\ 00 \\ 1* \\ *1 \end{bmatrix} \quad \begin{matrix} c_1 = 2 \\ c_2 = 1 \\ c_3 = 1 \end{matrix}$	$\begin{bmatrix} 00* \\ 000 \\ 1*0 \\ *10 \\ **0 \end{bmatrix} \quad \begin{array}{cccc c} 1 & 1 & 1 & 1 & 4 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ \hline 2 & 1 & 1 & 1 & \end{array}$	$\begin{bmatrix} 000 \\ 000 \\ 1** \\ *10 \\ **0 \end{bmatrix} \quad \begin{array}{cccc c} 2 & 0 & 1 & 1 & 4 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ \hline 2 & 1 & 1 & 1 & \end{array}$	$\begin{bmatrix} 000 \\ 000 \\ 1*0 \\ *1* \\ **0 \end{bmatrix} \quad \begin{array}{cccc c} 2 & 1 & 0 & 1 & 4 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ \hline 2 & 1 & 1 & 1 & \end{array}$
(a)	(b)	(c)	(d)

Fig. 7. Constructing ACs for $D = [(0001, 0000, 1101, 1110), (1, 1, 1, 1)]$ by extending strings. For clarity of illustration, ACs are represented as lists of strings inside $[\]$; a string is read horizontally. (a) An AC for the first two sites. (b)-(d) ACs and their corresponding restricted contingency tables obtained from extending the AC in (a) by an additional column with $r_0 = 4, r_1 = 0, r_* = 1$. Note that $c_* = 1$ in (b)-(d), thus satisfying $c + c_* = r_0 + r_1 + r_*$.

transformed into D' , or vice versa, using a series of appropriate recombination and coalescent events. It therefore follows that D and D' have the same set of ACs.

Second, if two input data sets D and D' differ by some permutation σ of their sites, then they have the same number of ACs. Any AC of D can be transformed into an AC of D' via σ and vice versa.

For D , a data set containing n sequences of length L , let $X(D) = \{(n_0^1, n_1^1), \dots, (n_0^L, n_1^L)\}$ denote the multiset of ordered pairs (n_0^i, n_1^i) , where n_0^i and n_1^i are the total number of 0s and 1s, respectively, at site i of D . Then, the above two facts together imply that two data sets D and D' have the same number of ACs if $X(D) = X(D')$ as multisets.

3.5 Main Idea: Counting Restricted Contingency Tables

Here we describe our method of counting ACs. Intuitively, our approach is to construct ACs moving along the sites, say from left to right. Given that we have constructed ACs for the first $k-1$ sites, we construct ACs for the first k sites by appending additional material (i.e., 0, 1, or *) to the right. We iterate this procedure until the last site is processed. As we describe below, this procedure of extending ACs can be translated into constructing contingency tables, where column sums are partially determined by the multiplicity of the strings in the AC that are being extended (i.e., an AC for the first $k-1$ sites) and row sums are partially determined by the allowed numbers of 0s and 1s at site k (c.f., Section 3.3).

3.5.1 Extending an AC

Suppose that ψ is an AC for the first $k-1$ sites. Further suppose that ψ contains d distinct strings, with multiplicity c_1, c_2, \dots, c_d satisfying $c_1 + c_2 + \dots + c_d = c$. Since ψ is an AC for the first $k-1$ sites, it does not contain any all-* string of length $k-1$. However, a string s in an AC for the first k sites may contain the all-* string of length $k-1$ as its prefix if the character at the k th site of s is either 0 or 1, but not *. For illustration, consider $D = [(0001, 0000, 1101, 1110), (1, 1, 1, 1)]$. An AC for the first two sites is shown in Fig. 7a, where we simply list strings to represent the AC. It has $c = 4$ and the multiplicity of "00" is 2, while every other sequence type has multiplicity 1. In any AC for the first two sites, the string "**" is not allowed, but "**0" or "**1" may appear in an AC for the first three sites. Shown in Fig. 7b, Fig. 7c, and Fig. 7d are three examples of ACs for the first three sites; they each contain

"**0". The all-* string "**" is not allowed in any AC for the first three sites.

In an AC for the first k sites, how many 0s and 1s can be present at site k ? We have already answered that question in Section 3.3, i.e., the number r_0 of 0s and the number r_1 of 1s satisfy $(r_0, r_1) \in I(n_0^k, n_1^k)$, where n_0^k and n_1^k , respectively, denote the number of 0s and 1s at site k in the input data D . How about the number r_* of *s at site k ? Since an all-* string of length k is not allowed, r_* is bounded from above, while the minimum value of r_* is determined by the total number of strings in ψ and the number of non-* characters at site k . A moment's thought leads to the conclusion that $\max(c - r_0 - r_1, 0) \leq r_* \leq c$. Similarly, the number c_* of all-* prefixes of length $k-1$ satisfies $\max(r_0 + r_1 - c, 0) \leq c_* \leq r_0 + r_1$. Note that c_* and r_* must satisfy $c + c_* = r_0 + r_1 + r_*$ for consistency. In Fig. 7b, Fig. 7c, and Fig. 7d, the ACs shown have $(r_0, r_1) = (4, 0) \in I(3, 1)$ and $r_* = 1$ at site three, and each AC has exactly one sequence containing "**" as a prefix of length 2 (i.e., $c_* = 1$). As required, $c + c_* = r_0 + r_1 + r_*$.

3.5.2 Restricted Contingency Tables

As above, let ψ be an AC for the first $k-1$ sites containing d distinct strings. The pairing of a character at site k with a length- $(k-1)$ string in ψ or a length- $(k-1)$ all-* string can be concisely summarized by a table of the following form:

	row sums			
	$A_{0,1}$	\cdots	$A_{0,d}$	$A_{0,*}$
	$A_{1,1}$	\cdots	$A_{1,d}$	$A_{1,*}$
	$A_{*,1}$	\cdots	$A_{*,d}$	0
column sums:	c_1	\cdots	c_d	c_*

For $j = 1, \dots, d$, a particular entry $A_{i,j}$ corresponds to the multiplicity of a length- k string obtained from appending character i to the right of a length- $(k-1)$ string of type j in ψ , whereas $A_{i,*}$ corresponds to the multiplicity of a length- k string obtained from appending character i to the right of an all-* string of length $k-1$. To avoid generating an all-* string of length k , we impose the condition $A_{*,*} = 0$. A column sum c_j , for $j = 1, \dots, d$, is given by the multiplicity of string type j in ψ and row sums r_0 and r_1 , respectively, denoting the number of 0s and 1s at site k in the corresponding new AC, satisfy $(r_0, r_1) \in I(n_0^k, n_1^k)$ (see Section 3.3). Further, as described above, the number c_* of all-* prefixes of length $k-1$ and the number r_* of *s at site k satisfy

$$\begin{aligned} \max(c - r_0 - r_1, 0) &\leq r_* \leq c, \\ \max(r_0 + r_1 - c, 0) &\leq c_* \leq r_0 + r_1, \end{aligned}$$

and $c_1 + \dots + c_d + c_* = r_0 + r_1 + r_*$. The number of nonzero entries in the above contingency table gives the number of *distinct* length- k strings in the corresponding new AC. In the next iteration, nonzero $A_{i,j}$ will appear as possible column sums when constructing contingency tables for site $k+1$. Returning to our example, contingency tables corresponding to the ACs discussed before are shown in Fig. 7b, Fig. 7c, and Fig. 7d. In Fig. 7b, all strings in the AC are distinct and the corresponding contingency table contains only 1s as nonzero entries. In Fig. 7c and Fig. 7d, the string 000 has multiplicity 2 and the entry $A_{0,1}$, which corresponds to the multiplicity of the length-3 string obtained from appending 0 to the right of 00, is 2.

In summary, there is a one-to-one correspondence between the nonnegative integer valued contingency tables described above and ACs for the first k sites that can be created by appending 0s, 1s, and *s to the right of the strings in an AC for the first $k-1$ sites. Hence, we can enumerate ACs by counting contingency tables.

Since the table described above has $A_{*,*} = 0$, we call it a *restricted* contingency table with row sums $\mathbf{r} = (r_0, r_1, r_*)$ and column sums $\mathbf{c} = (c_1, \dots, c_d, c_*)$. If the entry $A_{*,*}$ were not restricted to be zero, we would have a standard (or an *unrestricted*) contingency table with row sums \mathbf{r} and column sums \mathbf{c} . Restricted and unrestricted contingency tables are related as follows: Let $N_u(\mathbf{r}, \mathbf{c})$ denote the number of *unrestricted* contingency tables with column sums \mathbf{c} and row sums \mathbf{r} . Similarly, let $N_r(\mathbf{r}, \mathbf{c})$ denote the number of *restricted* contingency tables with row sums \mathbf{r} and column sums \mathbf{c} . For $\mathbf{r} = (r_0, r_1, r_*)$ and $\mathbf{c} = (c_1, \dots, c_d, c_*)$, define $\mathbf{r}' = (r_0, r_1, r_* - 1)$ and $\mathbf{c}' = (c_1, \dots, c_d, c_* - 1)$. Then, it is straightforward to show that

$$N_r(\mathbf{r}, \mathbf{c}) = \begin{cases} N_u(\mathbf{r}, \mathbf{c}) - N_u(\mathbf{r}', \mathbf{c}') & \text{if } r_* \geq 1 \text{ and } c_* \geq 1, \\ N_u(\mathbf{r}, \mathbf{c}) & \text{otherwise.} \end{cases} \quad (5)$$

We have not yet discussed how to obtain the total number of ACs. The case with two sites is amenable to analytic techniques. For more than two sites, we propose a solution via dynamic programming. These topics are subsequently discussed in the next two subsections.

3.6 Counting ACs for Two Sites

If there are only two sites in D , we only need to consider contingency tables of size 3×3 . This simplifies the problem considerably and it is possible to sum over all necessary 3×3 contingency tables explicitly to obtain closed-form formulas. We describe this analytical result below.

It follows from the discussion in the previous section that the number of ACs with c_0 0s and c_1 1s at the first site and r_0 0s and r_1 1s at the second site is given by the following sum of restricted contingency table numbers:

$$\sum_{j=\max(r_0+r_1, c_0+c_1)}^{r_0+r_1+c_0+c_1} N_r((r_0, r_1, j - r_0 - r_1), (c_0, c_1, j - c_0 - c_1)),$$

where the sum over j accounts for the allowed values of c_* and r_* . By using (5), this sum can be shown to be equal to

$N_u((r_0, r_1, c_0 + c_1), (c_0, c_1, r_0 + r_1))$. Hence, the total number $\beta(D)$ of ACs for the case with two sites is

$$\sum_{\substack{(c_0, c_1) \in I(n_0^1, n_1^1), \\ (r_0, r_1) \in I(n_0^2, n_1^2)}} N_u((r_0, r_1, c_0 + c_1), (c_0, c_1, r_0 + r_1)), \quad (6)$$

where, as before, $I(n_0^i, n_1^i)$ determines the allowed number of 0s and 1s at site i .

The key observation is that the expression shown in (6) can be summed explicitly. A 3×3 *unrestricted* contingency table of the form

			row sums
a	b	v	r_0
c	d	w	r_1
x	y	z	$c_0 + c_1$

column sums: $c_0 \quad c_1 \quad r_0 + r_1$

has only four degrees of freedom; once the entries a, b, c, d are chosen, the other entries get fixed by given row and column sums. Moreover, a, b, c, d satisfy the following set of constraints:

$$\begin{aligned} 0 &\leq a + b \leq r_0, \\ 0 &\leq c + d \leq r_1, \\ 0 &\leq a + c \leq c_0, \\ 0 &\leq b + d \leq c_1. \end{aligned}$$

These constraints imply

$$\begin{aligned} a \in \mathcal{A} &= \{0, \dots, \min(c_0, r_0)\}, \\ b \in \mathcal{B} &= \{0, \dots, \min(c_1, r_0 - a)\}, \\ c \in \mathcal{C} &= \{0, \dots, \min(r_1, c_0 - a)\}, \\ d \in \mathcal{D} &= \{0, \dots, \min(r_1 - c, c_1 - b)\}, \end{aligned}$$

and, therefore, (6) can be written as

$$\beta(D) = \sum_{\substack{(c_0, c_1) \in I(n_0^1, n_1^1), \\ (r_0, r_1) \in I(n_0^2, n_1^2)}} \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} \sum_{c \in \mathcal{C}} \sum_{d \in \mathcal{D}} 1.$$

This sum can be performed explicitly. For example, for $n_0^1 = n_1^1 = n_0^2 = n_1^2 = n/2$, for n even, the total number $\beta(D)$ of ACs is given by the following degree-8 polynomial:

$$\begin{aligned} \beta(D) = 2 &+ \frac{179}{84}n + \frac{2,339}{1,260}n^2 + \frac{151}{160}n^3 + \frac{3,851}{11,520}n^4 \\ &+ \frac{31}{384}n^5 + \frac{259}{23,040}n^6 + \frac{11}{13,440}n^7 + \frac{11}{430,080}n^8. \end{aligned} \quad (7)$$

Note that the degree 8 is equal to the number $3^L - 1$ of distinct strings over $\{0, 1, *\}$, except for the all-* string, of length $L = 2$.

3.7 Counting ACs for Arbitrary Number of Sites

For more than two sites, we adopt a dynamic programming approach for counting *restricted* contingency tables. As mentioned in Section 3.5, our algorithm progresses sequentially along the sites, from left to right. There also exists a more efficient, albeit somewhat more complicated, dynamic programming formulation in terms of counting *unrestricted*

contingency tables, but we will not discuss that here. Further, we remark that we considered an alternative approach to counting ACs that is based on counting lattice points bounded by polytopes. Computing the number of such lattice points is a widely studied problem in mathematics (for example, see [13] and references therein). Moreover, there is a public software called *LattE* (available at <http://www.math.ucdavis.edu/~latte/>) that enumerates lattice points, but we found that it is significantly slower than our own program for enumerating ACs via counting restricted contingency tables; it seems that $L = 2$ is the only case that *LattE* can handle.

Recall that both contingency tables shown in Fig. 7c and Fig. 7d have 2, 1, 1, 1 as nonzero entries, although the two tables have different corresponding ACs. In each AC, as the corresponding contingency table captures, there are four distinct types of length-3 strings, with exactly one type being of multiplicity 2, while every other type is of multiplicity 1. Now, consider extending the ACs in Fig. 7c and Fig. 7d to create ACs for the first four sites. For both cases, we need to find contingency tables of the following form:

					row sums
$A_{0,1}$	$A_{0,2}$	$A_{0,3}$	$A_{0,4}$	$A_{0,*}$	r_0
$A_{1,1}$	$A_{1,2}$	$A_{1,3}$	$A_{1,4}$	$A_{1,*}$	r_1
$A_{*,1}$	$A_{*,2}$	$A_{*,3}$	$A_{*,4}$	0	r_*
column sums:	2	1	1	1	c_*

where $(r_0, r_1) \in I(n_0^4, n_1^4)$. As a consequence, the same set of contingency tables will be generated, although the corresponding set of ACs will be different for the two cases. (In the ACs obtained from extending the AC in Fig. 7c, possible prefixes of length 3 are "000", "1**", "*10", "**0" and "***", whereas, in the ACs obtained from extending the AC in Fig. 7d, possible prefixes of length 3 are "000", "1*0", "*1*", "**0", and "***.") Hence, since we are only interested in *counting* ACs and not in constructing the actual ACs themselves, we do not need to consider extensions of the ACs in Fig. 7c and Fig. 7d separately. At each iteration of our algorithm, we only keep track of distinct multiplicity configurations and how often each configuration appears. Then, we construct contingency tables for each distinct multiplicity configuration. This method considerably reduces the number of contingency tables to construct.

We now describe our enumeration algorithm. Given a matrix M , let $\Lambda(M)$ denote the descending array of positive integers in M . For example, if

$$M = \begin{pmatrix} 3 & 2 & 0 & 4 \\ 1 & 0 & 2 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix},$$

then $\Lambda(M) = (4, 3, 2, 2, 1, 1)$. Suppose that M and M' are contingency tables that arise while considering site k . When constructing contingency tables for site $k+1$, we do not need to distinguish M from M' if $\Lambda(M) = \Lambda(M')$; we just need to keep track of how many tables for site k have $\Lambda(M)$ as a descending array of positive integers. (Note that $\Lambda(M)$

are used as column sums when constructing contingency tables for site $k+1$.)

Given a k -tuple $v = (v_1, v_2, \dots, v_k)$ of nonnegative integers, we define $|v| = \sum_{i=1}^k v_i$. Let $\mathcal{R}(c, r)$ be the set of all restricted contingency tables with column sums (c, c_*) and row sums (r, r_*) , where $\max(|r| - |c|, 0) \leq c_* \leq |r|$ and $\max(|c| - |r|, 0) \leq r_* \leq |c|$. As usual, (n_0^i, n_1^i) denote the number of 0s and 1s at site i in the input data D .

For each site i , we associate a set \mathcal{C}_i of column sums. To initialize the algorithm, we define

$$\mathcal{C}_1 := \{\Lambda(r_0, r_1) \mid (r_0, r_1) \in I(n_0^1, n_1^1)\},$$

where $\Lambda(r_0, r_1)$ is (r_0, r_1) if $r_0 \geq r_1$ or (r_1, r_0) if $r_0 < r_1$. For each $c^1 \in \mathcal{C}_1$, the *multiplicity* of c^1 is defined as

$$\mu(c^1) := |\{(r_0, r_1) \in I(n_0^1, n_1^1) \mid \Lambda(r_0, r_1) = c^1\}|.$$

For the k th site, where $k > 1$, \mathcal{C}_k is recursively defined as

$$\mathcal{C}_k := \left\{ \Lambda(M) \mid \begin{array}{l} M \in \mathcal{R}(c^{k-1}, r^k), \text{ where} \\ c^{k-1} \in \mathcal{C}_{k-1} \text{ and } r^k \in I(n_0^k, n_1^k) \end{array} \right\}.$$

Note that only distinct $\Lambda(r^k)$ need to be considered when constructing \mathcal{C}_k . For every $c^{k-1} \in \mathcal{C}_{k-1}$ and every $c^k \in \mathcal{C}_k$, we count the following number of restricted contingency tables M with $\Lambda(M) = c^k$:

$$w(c^{k-1}, c^k) := \sum_{r^k \in I(n_0^k, n_1^k)} |\{M \in \mathcal{R}(c^{k-1}, r^k) \mid \Lambda(M) = c^k\}|.$$

This can then be used to define the multiplicity of $c^k \in \mathcal{C}_k$ as

$$\mu_k(c^k) = \sum_{c^{k-1} \in \mathcal{C}_{k-1}} \mu_{k-1}(c^{k-1}) w(c^{k-1}, c^k).$$

In practice, \mathcal{C}_k and $\mu_k(\cdot)$ can be determined concurrently.

The total number $\beta(D)$ of ACs for D is obtained by summing the multiplicity of column sums for the last site:

$$\beta(D) = \sum_{c^L \in \mathcal{C}_L} \mu_L(c^L). \quad (8)$$

More generally, the total number of ACs for $D[1:k]$, restriction of D to the first k sites, is

$$\beta(D[1:k]) = \sum_{c^k \in \mathcal{C}_k} \mu_k(c^k).$$

We have written a C++ program that implements the dynamic programming algorithm described above. For two sites, we have checked numerically that (8) gives the same answers as does (6).

3.8 Some Explicit Enumeration of ACs

For n even, $n_0^i = n_1^i = n/2$ for every site i leads to the largest number of ACs. For n odd, the largest number of ACs is achieved if $n_0^i = \lfloor n/2 \rfloor + 1$ and $n_1^i = \lfloor n/2 \rfloor$ for every site i , where $\lfloor \cdot \rfloor$ denotes the floor function. The left-hand side of Table 2 shows such largest possible number of ACs for a small number n of sequences and a small number L of sites. It is quite striking that the number of ACs grows so fast as the data size increases. These numbers were obtained using our software that implements the dynamic programming algorithm described in Section 3.7. For $L = 2$ and n even,

TABLE 2
Total Number of ACs for Data D with, for Every Site i , $n_0^i = a_0$ and $n_1^i = a_1$,
where $a_0 = a_1 = \lfloor n/2 \rfloor$ for n Even, and $a_0 = \lfloor n/2 \rfloor + 1, a_1 = \lfloor n/2 \rfloor$ for n Odd

$\beta(D)$ (With Recombination)					$\max(\alpha(D))$ (Without Recombination)				
n	L				n	L			
	2	3	4	5		2	3	4	5
2	30	573	16 875	689 175	2	5	9	17	33
3	108	6 286	743 387	149 861 079	3	9	17	33	65
4	330	62 589	32 482 009	35 523 729 489	4	11	19	35	67
5	866	445 137	893 479 326	4 938 627 635 669	5	16	28	52	100
6	2 143	3 302 506	29 521 615 942	962 962 451 049 968	6	19	31	55	103
7	4 611	17 409 443	568 860 072 916	91 812 561 254 804 105	7	25	41	73	137

Under these conditions, the maximum value of $\alpha(D)$ is $1 + a_0(a_1 + 2^L - 1)$. These values are shown on the right-hand side. Shown on the left-hand side is the total number of ACs when an arbitrary number of recombinations are allowed.

our closed-form formula (7) agrees with the numerical values shown in Table 2.

For fixed L , the total number of ACs is bounded from above by a polynomial in n of degree $3^L - 1$. There are $3^L - 1$ nonnegative integer valued variables, m_1, \dots, m_{3^L-1} , and the maximum value that any variable can take is n . Hence, the total number of ACs is $O(n^{3^L-1})$. As this simple analysis shows, the total number of ACs depends more crucially on the number of sites than on the number of sequences. This qualitative behavior of growth is also apparent in Table 2.

In the absence of recombination, there are 2^L distinct binary sequences of length L , but the genealogy that is given by a tree puts tight restriction on which sequences can appear together in an AC. As in the above case, for every site i , let $n_0^i = a_0$ and $n_1^i = a_1$, where $a_0 = a_1 = n/2$ if n is even, and $a_0 = \lfloor n/2 \rfloor + 1, a_1 = \lfloor n/2 \rfloor$ if n is odd. Under these conditions, the number $\alpha(D)$ of ACs is maximum if the input data contains a_0 all-0 sequences and a_1 all-1 sequences. One can use the method described in Section 2.2 to show that $\max(\alpha(D)) = 1 + a_0(a_1 + 2^L - 1)$, where L is the number of sites. Some numerical values of this function are shown on the right-hand side of Table 2. These numbers are negligibly small compared to the value of $\beta(D)$ shown on the left-hand side of Table 2.

For L sites, let $B_L(k)$ denote the number of ACs with exactly k strings and let $S_L(k)$ denote the cumulative count $\sum_{i=1}^k B_L(i)$, which gives the number of ACs with at most k strings. Table 3 shows $B_L(k)$ and $S_L(k)$ for data sets with $n = 5$ and $(n_0^i, n_1^i) = (3, 2)$ for every site i . The largest value of $S_L(k)$ corresponds to $\beta(D)$ for $n = 5$ in Table 2.

4 ANCESTRAL CONFIGURATIONS IN CONSTRAINED ARGs

In this section, we enumerate ACs in ARGs with at most R recombinations. This study is motivated by the problem of approximating the likelihood $P(D)$ of the input data D when the recombination rate is low by summing over the probability of ARGs with a small number of recombinations. Further, the work described here can be used to enumerate all ARGs that can generate D using at most R recombinations under the infinite-sites

model of mutation. The constraint on the number of recombinations further complicates the problem of enumerating ACs and combinatorial analysis is more difficult to carry out. We provide a dynamic programming algorithm that enumerates ACs by explicitly tracking all possible evolutionary histories backward in time.

Cautionary Remark. Given a data set D , an AC ψ that is R recombinations away from D may not appear in any ARG for D with at most R recombinations; further recombinations may be required for ψ to reach the grand common ancestor. Hence, the number of ACs in ARGs for D with at most R recombinations is, in general, smaller than the number of ACs that can be reached from D using at most R recombinations. In what follows, we are interested in counting the former.

4.1 Examples Revisited

As mentioned in Section 3.4, when an arbitrary number of recombinations are allowed, the set of ACs does not depend on the order of 0s and 1s within a column. For example, $D = [(00, 11), (1, 1)]$ and $D' = [(01, 10), (1, 1)]$ have the same set of ACs, shown in Fig. 6, when an arbitrary number of recombinations are allowed. That no longer holds true if there is a restriction on the number of allowed recombinations. This fact is easiest to see in the absence of recombination, in which case, the set of ACs for D is $\{[(00, 11), (1, 1)], [(00, 01), (1, 1)], [(00, 10), (1, 1)], [(00), (2)], [(00), (1)], [(01, 10), (1, 1)], [(01, 00), (1, 1)], [(00, 10), (1, 1)], [(00), (2)], [(00), (1)]\}$. Similarly, when only a single recombination is allowed, it is straightforward to show that D and D' have different sets of ACs, although the two sets are both of size 19.

When there is a restriction on the number of recombinations, even the number of ACs may depend on the order of 0s and 1s within a column. Consider the two data sets $D = [(00, 11), (1, 2)]$ and $D' = [(01, 10, 11), (1, 1, 1)]$. In the absence of recombination, D has six ACs, while there exists no coalescent tree for D' (c.f., the three gamete test mentioned in Section 2). The minimum number of recombinations for D' is one and there are 21 inequivalent ACs in the set of ARGs for D' with exactly one recombination. In contrast, there are 26 ACs for D when one recombination is allowed.

TABLE 3
Counting of ACs with a Restricted Number of Strings

k	$B_2(k)$	$S_2(k)$	$B_3(k)$	$S_3(k)$	$B_4(k)$	$S_4(k)$	$B_5(k)$	$S_5(k)$
1	1	1	1	1	1	1	1	1
2	12	13	62	63	312	313	1562	1563
3	60	73	1143	1206	21720	22033	412683	414246
4	163	236	9517	10723	575921	597954	35132963	35547209
5	244	480	40006	50729	6811757	7409711	1168994740	1204541949
6	210	690	89716	140445	38651220	46060931	16622300644	17826842593
7	120	810	118047	258492	115829376	161890307	112796639915	130623482508
8	45	855	99013	357505	201958486	363848793	408212643773	538836126281
9	10	865	56395	413900	222978515	586827308	868467110362	1407303236643
10	1	866	22818	436718	167052986	753880294	1177704929588	2585008166231
11			6732	443450	89613906	843494200	1087103087864	3672111254095
12			1447	444897	35821104	879315304	719916684420	4392027938515
13			218	445115	10971625	890286929	356517388653	4748545327168
14			21	445136	2620233	892907162	136353481775	4884898808943
15			1	445137	491380	893398542	41276692230	4926175501173
16					72009	893470551	10070266029	4936245767202
17					8076	893478627	2004602973	4938250370175
18					662	893479289	327842330	4938578212505
19					36	893479325	44105460	4938622317965
20					1	893479326	4854891	4938627172856
21							431132	4938627603988
22							30060	4938627634048
23							1565	4938627635613
24							55	4938627635668
25							1	4938627635669

For L sites, $B_L(k)$ denotes the number of ACs with exactly k strings, while $S_L(k)$ denotes the number of ACs with at most k strings. This table is for data with $n = 5$ and $(n_0^i, n_1^i) = (3, 2)$ for every site i .

In Section 3.4, we observed that, when an arbitrary number of recombinations is allowed, two data sets differing by some permutation of their columns have the same number of ACs. This also is no longer true, in general, if the number of recombinations is constrained. For example, consider the two data sets $D = [(001, 110, 111), (1, 1, 1)]$ and $D' = [(010, 101, 111), (1, 1, 1)]$. While D has an ARG with only one recombination, no valid ARG with less than two recombinations exists for D' .

4.2 The Search Algorithm

For the case with a restricted number R of recombinations, the approach we take is to explicitly trace all possible evolutionary histories backward in time. For any particular AC, we try all possible events—that is, recombination, coalescent, or mutation events—that can happen to that AC. Possible recombination and coalescent events are as illustrated in Fig. 4. We do not consider recombination events that produce an all- $*$ sequence. Recall that, because we assume the infinite-sites model, a mutation of 1 to 0 at a site is possible if and only if exactly one sequence carries a 1 at that site. Going backward in time in this manner, event by event, we can exhaustively find all sequences of events

with less than or equal to the allowed number of recombinations, while keeping track of all distinct ACs encountered. To accelerate the search, we propose two ideas, one based on branch and bound and the other on dynamic programming.

We maintain a hash table H that stores ACs already encountered, together with the number of recombinations further allowed for each AC. Initially, H contains the pair (D, R) , where D is the given data set and R is the maximum number of allowed recombinations. Given an AC ψ and its associated limit c on the number of recombinations, we first compute the minimum number k of recombinations needed in the evolution of ψ ; our method of computing k is described in [14]. If $k > c$, we backtrack the search. Otherwise, we check whether ψ is present in H . Suppose that (ψ, c') with $c' \geq c$ is present in H . Then, continuing along the current search path will not lead to any AC that has not been encountered previously, so, again, it is safe to backtrack. If $c' < c$ or if ψ is not present in H , then we need to add the entry (ψ, c) to H and iteratively try all possible events backward in time starting from ψ . When either a mutation or a coalescent event is tried, the number of

TABLE 4

The Maximum Number, over All Data Sets D Satisfying $n_0^i = \lceil n/2 \rceil$ and $n_1^i = \lfloor n/2 \rfloor$ for Every Site i , of ACs that Can Be Encountered in ARGs for D with the Minimum Number $R_{\min}(D)$ of Recombinations

n	$\max_D [\gamma_{R_{\min}(D)}(D)]$			
	L			
	2	3	4	5
2	5(0)	9(0)	17(0)	33(0)
3	9(0)	17(0)	33(0)	65(0)
4	34(1)	216(2)	1524(3)	12085(4)
5	60(1)	408(2)	2986(3)	23958(4)
6	78(1)	2709(3)	17053(4)	1210197(6)

The numbers in parentheses indicate the minimum number $R_{\min}(D)$ of recombinations for D with the maximum value of $\gamma_{R_{\min}(D)}(D)$.

allowed recombinations remains c ; when a recombination event is tried, the number of allowed recombinations decreases to $c - 1$. The hash table serves a dual purpose. First, as just described, it allows us to avoid repeating a search already carried out. Second, when all branches have been explored, the number of entries in the hash table equals $\gamma_R(D)$, the number of ACs that can be encountered in ARGs deriving D with at most R recombinations. It is straightforward to produce a list of all $\gamma_R(D)$ inequivalent ACs encountered during the search.

We have written a C program, called *greven*, that implements the method described above. For small data sets, we have checked that this program agrees with *aceTrees* when no recombination is allowed and with *aceARGs* when an arbitrary number of recombinations are allowed.

4.3 Further Explicit Enumeration of ACs

For small values of n , we applied our program to all data sets with the same values of n_0^i and n_1^i as in Table 2. Let $R_{\min}(D)$ denote the minimum number of recombinations needed in the evolution of a given set D of sequences. Shown in Table 4 is $\max_D [\gamma_{R_{\min}(D)}(D)]$, i.e., the maximum number, over all D , of ACs that can be encountered in ARGs for D with $R_{\min}(D)$ recombinations. The numbers in parentheses denote the minimum number $R_{\min}(D)$ of recombinations for D with the maximum value of $\gamma_{R_{\min}(D)}(D)$. Table 4 should be compared with Table 2. As this enumeration shows, although the number of ACs encountered in minimal ARGs is noticeably larger than the number of ACs possible in the absence of recombination, it is still negligibly small compared to the total number of ACs when an arbitrary number of recombinations are allowed.

5 DISCUSSION

In this paper, we addressed the problem of enumerating all configurations (or Markov states) ancestral to a given data set D , consisting of binary sequences of finite length. Our main motivation was to understand in detail the difficulties involved in computing the likelihood $P(D)$ of observing D . We constructed a general method of obtaining closed-form

formulas (as functions of the multiplicity of distinct sequences in D) for the total number of ACs when recombination is absent. For the case in which recombination is allowed, we devised a method of enumeration based on counting contingency tables and constructed a dynamic programming algorithm. In terms of the number of ACs, one can use our methods to rigorously study how the case with recombination compares with the case without recombination. In our work, we assumed that recombination breakpoints occur at the midpoints of consecutive sites in D . If this assumption is relaxed, then the number of ACs will increase.

We also discussed a method of counting the number of ACs that can appear in ARGs with less than or equal to a given number R of recombinations. Of particular interest is the case in which R is (or is close to) the minimum number of recombinations for D . When the recombination rate is low, it is widely believed that ARGs with the minimum or “near minimum” number of recombinations should make the dominant contribution to the likelihood $P(D)$. However, we are not aware of any work that has actually tried to sum the probability over all minimal or near minimal ARGs. That problem is the main motivation for our present work. Ultimately, we would like to see how the total probability carried by ARGs with at most R recombinations changes as R increases. The goal of this paper was to understand some basic properties of the Markov state space itself and to gauge the size of data that we may be able to handle.

As discussed in Section 3.8, the total number of ACs can be inordinately large when an arbitrary number of recombinations are allowed. Not all ACs play an equally important role in computing the likelihood, however. We mentioned in Section 3.4 that if D and D' are related by a series of operations involving permutations of characters within a column and/or permutations of entire columns, then they have the same total number of ACs. But, D and D' may look radically different to anybody's eye and carry quite different information. For example, consider the following two contrived data sets:

$$D = \begin{pmatrix} 1010101 \\ 1111111 \\ 0101010 \\ 0000000 \end{pmatrix} \text{ and } D' = \begin{pmatrix} 1111111 \\ 1111111 \\ 0000000 \\ 0000000 \end{pmatrix}.$$

When an arbitrary number of recombinations are allowed, D and D' have the same number of ACs, but they are qualitatively very different; clear signals of past recombinations are present in D , while D' is compatible with a tree. The configuration space is very large, containing about 2.26×10^{17} ACs, but, for low recombination rates, a large portion of those ACs should be insignificant for computing the likelihood $P(D')$. In general, if we are to devise a deterministic method for estimating the likelihood of observing a given data set, we would need to investigate how we can systematically take advantage of our understanding of the configuration space and of the structure of the Markov chain transition graph.

ACKNOWLEDGMENTS

This work is supported by grants EIA-0220154 and IIS-0513910 (Song) from the US National Science Foundation, by grant HAMJW (Hein, Lyngsø) from EPSRC, and by grant HAMKA (Hein, Lyngsø) from MRC.

REFERENCES

- [1] S.N. Ethier and R.C. Griffiths, "The Infinitely-Many-Sites Model as a Measure Valued Diffusion," *Annals of Probability*, vol. 15, pp. 515-545, 1987.
- [2] S.N. Ethier and R.C. Griffiths, "On the Two-Locus Sampling Distribution," *J. Math. Biology*, vol. 29, pp. 131-159, 1990.
- [3] P. Fearnhead and P. Donnelly, "Estimating Recombination Rates from Population Genetic Data," *Genetics*, vol. 159, pp. 1299-1318, 2001.
- [4] R.C. Griffiths, "Genealogical-Tree Probabilities in the Infinitely-Many-Site Mode," *J. Math. Biology*, vol. 27, pp. 667-680, 1989.
- [5] R.C. Griffiths and P. Marjoram, "Ancestral Inference from Samples of DNA Sequences with Recombination," *J. Computational Biology*, vol. 3, pp. 479-502, 1996.
- [6] R.C. Griffiths and S. Tavaré, "Ancestral Inference in Population Genetics," *Statistics in Science*, vol. 9, pp. 307-319, 1994.
- [7] R.C. Griffiths and S. Tavaré, "Simulating Probability Distributions in the Coalescent," *Theoretical Population Biology*, vol. 46, pp. 131-159, 1994.
- [8] D. Gusfield, "Efficient Algorithms for Inferring Evolutionary Trees," *Networks*, vol. 21, pp. 19-28, 1991.
- [9] J.F.C. Kingman, "The Coalescent," *Stochastic Processing Applications*, vol. 13, pp. 235-248, 1982.
- [10] J.F.C. Kingman, "On the Genealogy of Large Populations," *J. Applied Probability*, vol. 19A, pp. 27-43, 1982.
- [11] M.K. Kuhner, J. Yamato, and J. Felsenstein, "Estimating Effective Population Size and Mutation Rate from Sequence Data Using Metropolis-Hastings Sampling," *Genetics*, vol. 140, pp. 1421-1430, 1995.
- [12] M.K. Kuhner, J. Yamato, and J. Felsenstein, "Maximum Likelihood Estimation of Recombination Rates from Population Data," *Genetics*, vol. 156, pp. 1393-1401, 2000.
- [13] J. De Loera, R. Hemmecke, J. Tauzer, and R. Yoshida, "Effective Lattice Point Counting in Rational Convex Polytopes," *J. Symbolic Computation*, vol. 38, pp. 1273-1302, 2004.
- [14] R. Lyngsø, Y.S. Song, and J. Hein, "Minimum Recombination Histories by Branch and Bound," *Proc. 2005 Workshop Algorithms in Bioinformatics*, pp. 239-250, 2005.
- [15] K.L. Simonsen and G.A. Churchill, "A Markov Chain Model of Coalescence with Recombination," *Theoretical Population Biology*, vol. 52, pp. 43-59, 1997.
- [16] M. Stephens and P. Donnelly, "Inference in Molecular Population Genetics," *J. Royal Statistical Soc. Series B*, vol. 62, pp. 605-655, 2000.
- [17] R.H. Ward, B.L. Frazier, K. Dew, and S. Pääbo, "Extensive Mitochondria Diversity within a Single Amerindian Tribe," *Proc. Nat'l Academy Science*, vol. 88, pp. 8720-8724, 1991.



Yun S. Song received the BS degrees in mathematics and physics from the Massachusetts Institute of Technology (MIT) and the PhD degree in physics from Stanford University. He was a postdoctoral researcher at the University of Oxford, where he became interested in mathematical population genetics and computational biology. He is now a postdoctoral researcher in the Department of Computer Science at the University of California, Davis.



Rune Lyngsø received the PhD degree in computer science from the University of Aarhus in 2000. Since then, he has been working as a postdoctoral researcher, first with Professor David Haussler at the University of California at Santa Cruz and currently with Professor Jotun Hein at the University of Oxford. His main contributions are in the areas of RNA secondary structure prediction, analysis of stochastic models of formal languages, and string algorithms.



Jotun Hein is a professor in the Department of Statistics at the University of Oxford and is a member of the Royal Danish Academy of Sciences and Letters. Before moving to Oxford in 2001, he was an associate professor at the University of Aarhus, where he served as the director of Bioinformatics Research Center. Most of his work is centered around developing methods for analyzing data that arise in molecular genetics and genomics, with a comparative/evolutionary focus. He was an associate editor for *Genetics* from 2001-2003 and is a coauthor of a book on molecular population genetics entitled *Gene Genealogies, Sequence Variation and Evolution*, published by Oxford University Press.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.