

Efficient strategies for calculating blockwise likelihoods under the coalescent

Konrad Lohse^{*,1}, Martin Chmelik[†], Simon H. Martin[‡] and Nicholas H. Barton[§]

^{*}Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH93FL, UK, [†][§]Institute of Science and Technology, Am Campus 1, A-3400 Klosterneuburg, Austria, [‡]Zoology Department, University of Cambridge, UK

ABSTRACT The inference of demographic history from genome data is hindered by a lack of efficient computational approaches. In particular, it has proven difficult to exploit the information contained in the distribution of genealogies across the genome. We have previously shown that the generating function (GF) of genealogies can be used to analytically compute likelihoods of demographic models from configurations of mutations in short sequence blocks (Lohse *et al.* 2011). Although the GF has a simple, recursive form, the size of such likelihood calculations explodes quickly with the number of individuals and applications of this framework have so far been mainly limited to small samples (pairs and triplets) for which the GF can be written down by hand. Here we investigate several strategies for exploiting the inherent symmetries of the coalescent. In particular, we show that the GF of genealogies can be decomposed into a set of equivalence classes which allows likelihood calculations from non-trivial samples. Using this strategy, we automated blockwise likelihood calculations for a general set of demographic scenarios in *Mathematica*. These histories may involve population size changes, continuous migration, discrete divergence and admixture between multiple populations. To give a concrete example, we calculate the likelihood for a model of isolation with migration (IM), assuming two diploid samples without phase and outgroup information. We demonstrate the new inference scheme with an analysis of two individual butterfly genomes from the sister species *Heliconius melpomene rosina* and *Heliconius cydno*.

KEYWORDS Maximum likelihood, population divergence, gene flow, structured coalescent, generating function

Author Affiliations

Genomes contain a wealth of information about the demographic and selective history of populations. However, efficiently extracting this information to fit explicit models of population history remains a considerable computational challenge. It is currently not feasible to base demographic inference on a complete description of the ancestral process of coalescence and recombination, and so inference methods generally rely on making simplifying assumptions about recombination. In the most extreme case of methods based on the site frequency spectrum (SFS), information contained in the physical linkage

of sites is ignored altogether (Gutenkunst *et al.* 2009; Excoffier *et al.* 2013). Because the SFS is a function only of the expected length of genealogical branches (Griffiths and Tavaré 1998; Chen 2012), this greatly simplifies likelihood computations. However, it also sacrifices much of the information about past demography (Terhorst and Song 2015). Other methods approximate recombination along the genome as a Markov process (Li and Durbin 2011; Harris and Nielsen 2013; Rasmussen *et al.* 2014). However, this approach is computationally intensive, limited to simple models (Schiffels and Durbin 2014) and/or pairwise samples (Li and Durbin 2011; Mailund *et al.* 2012) and requires phase information and well assembled genomes which are still only available for a handful of species.

A different class of methods assumes that recombination can be ignored within sufficiently short blocks of sequence (Hey and Nielsen 2004; Yang 2002). The benefit of this "multi-locus assumption" is that it gives a tractable framework for analysing linked sites, and so captures the information contained in the

Copyright © 2015 by the Genetics Society of America

doi: 10.1534/genetics.XXX.XXXXXX

Manuscript compiled: Monday 21st December, 2015%

¹Institute of Evolutionary Biology, University of Edinburgh, King's Bldgs., Charlotte Auerbach Road, Edinburgh, EH9 3FL, United Kingdom. E-mail: konrad.lohse@gmail.com.

distribution of genealogical branches. Multi-locus methods are also attractive in practice because they naturally apply to RAD data or partially assembled genomes that can now be generated for any species (e.g. Davey and Blaxter 2011; Hearn *et al.* 2014).

For small samples, the probability of seeing a particular configuration of mutations at a locus can be obtained analytically. For example, Wilkinson-Herbots (2008) and Wang and Hey (2010) have derived the distribution of pairwise differences under a model of isolation with migration (IM) and Wilkinson-Herbots (2012) has extended this to a history where migration is limited to an initial period. Yang (2002) derives the probability of mutational configurations under a divergence model for three populations and a single sample from each and Zhu and Yang (2012) have included migration between the most recently diverged pair of populations in this model. However, all of these particular cases can be calculated using a general procedure based on the generating function for the genealogy (Lohse *et al.* 2011). Here we explain how the GF and – from it – model likelihoods can be efficiently computed for larger samples than has hitherto been possible.

The generating function of genealogies

Assuming an infinite sites mutation model and an outgroup to polarize mutations, the information in a non-recombining block of sequence can be summarized as a vector \underline{k} of counts of mutations on all possible genealogical branches \underline{t} . Both \underline{t} and \underline{k} are labelled by the individuals that descend from them. We have previously shown that the probability of seeing a particular configuration of mutations \underline{k} can be calculated directly from the Laplace Transform or generating function (GF) of genealogical branches (Lohse *et al.* 2011). Given a large sample of unlinked blocks, this gives a framework for computing likelihoods under any demographic model and sampling scheme. Full details are given in Lohse *et al.* (2011). Briefly, the GF is defined as $\psi[\underline{\omega}] = E[e^{-\underline{\omega} \cdot \underline{t}}]$, where $\underline{\omega}$ is a vector of dummy variables corresponding to \underline{t} . Setting the $\underline{\omega}$ to zero necessarily gives one, the total probability; differentiating with respect to ω_i and setting the $\underline{\omega}$ to zero gives (minus) the expected coalescence time. If we assume an infinite sites mutation model, the probability of seeing k_s mutations on a particular branch s is (Lohse *et al.* 2011, eq. 1):

$$P[k_s] = E \left[e^{-\mu t_s} \frac{(\mu t_s)^{k_s}}{k_s!} \right] = \frac{(-\mu)^{k_s}}{k_s!} \left(\frac{\partial^{k_s} \psi}{\partial \omega_s^{k_s}} \right)_{\omega_s = \mu} \quad (1)$$

This calculation extends to the joint probability of mutations $P[\underline{k}]$. Using the GF rather than the distribution of branches itself to compute $P[\underline{k}]$ is convenient because we avoid the Felsenstein (1988) integral and because the GF has a very simple form: going backwards in time, the GF is a recursion over successive events in the history of the sample (Lohse *et al.* 2011, eq. 4):

$$\psi[\Omega] = \frac{\sum_i \lambda_i \psi[\Omega_i]}{\left(\sum_i \lambda_i + \sum_{|S|=1} \omega_S \right)} \quad (2)$$

where, going backwards in time, Ω denotes the sampling configuration (i.e. the location and state of lineages) before some event i and Ω_i the sampling configuration afterwards. Events during this interval occur with a total rate $\sum_i \lambda_i$. The numerator is a sum over all the possible events i each weighted by its rate λ_i . Equation 2 applies to any history that consists of independently occurring events. As outlined by Lohse *et al.* (2011), the GF for

models involving discrete events (population splits, bottlenecks) can be found by inverting the GF of the analogous continuous model. In other words, if we know the GF for a model that assumes an exponential rate of events at rate Λ , then taking the inverse Laplace Transform with respect to Λ gives the GF for any fixed time of the event.

In principle, the GF recursion applies to any sample size and model and can be automated using symbolic software (such as *Mathematica*). In practice however, likelihood calculations based on the GF have so far been limited to pairs and triplets: Lohse *et al.* (2011) computed likelihoods for an IM model with unidirectional migration for three sampled genomes and Lohse *et al.* (2012) and Hearn *et al.* (2014) derived likelihoods for a range of divergence histories for a single genome from each of three populations with instantaneous admixture, including the model used by Green *et al.* (2010) to infer Neandertal admixture into modern humans (Lohse and Frantz 2014).

There are several serious challenges in applying the GF framework to larger samples of individuals. First, the number of sample configurations (and hence GF equations) grows super-exponentially with sample size. Thus, the task of solving the GF and differentiating it to tabulate probabilities for all possible mutational configurations quickly becomes computationally prohibitive. Second, models involving reversible state transitions, such as two-way migration or recombination between loci, include a potentially infinite number of events. Solving the GF for such cases involves matrix inversions (Hobolth *et al.* 2011; Lohse *et al.* 2011). Third, while assuming infinite sites mutations may be convenient mathematically and realistic for closely related sequences, this assumption becomes problematic for more distantly related outgroups that are used to polarise mutations in practice. Finally, being able to uniquely map mutations onto genealogical branches assumes phased data that are rarely available for diploid organisms, given the limitations of current sequencing technologies.

In the first part of this paper, we discuss each of these problems in turn and introduce several strategies to remedy the explosion of terms and computation time. These arguments apply generally, irrespective of the peculiarities of particular demographic models and sampling schemes, and suggest a computational "pipeline" for likelihood calculations for non-trivial samples of individuals (up to $n = 6$). The accompanying *Mathematica* notebook implements this scheme for a wide range of demographic histories that may involve arbitrary divergence, admixture and migration between multiple populations, as well as population size changes. As a concrete example, we describe maximum likelihood calculations for a model of isolation with continuous migration (IM) between two populations for unphased and unpolarized data from two diploid individuals. We compare the power of this scheme to that of minimal samples of a single haploid sequence per population. Finally, to illustrate the new method, we estimate divergence and migration between the butterfly species *Heliconius melpomene* and *H. cydno* (Martin *et al.* 2013).

Models and Methods

Partitioning the GF into equivalence classes

Since the GF is defined in terms of genealogical branches and each topology is specified by a unique set of branches, an intuitive strategy for computing likelihoods is to partition the GF into contributions from different topologies. To condition on a certain topology, we simply set GF terms that are incompatible

with it to 0 (Lohse *et al.* 2011). Importantly however, such incompatible events still contribute to the total rate $\sum_i \lambda_i$ of events in the denominator of equation 2. Then, setting all ω in the topology-conditioned GF to zero gives the probability of that particular topology. Although conditioning on a particular topology gives a GF with a manageable number of terms, it is clearly not practical to do this for all possible topologies, given their sheer number even for moderate n (Table 1).

In the following, we will distinguish between ranked and unranked topologies. The GF is a sum over all possible sequences of events in the history of a sample; Edwards (1970) called them "labelled histories". Considering only coalescence events, each labelled history corresponds to a ranked topology, i.e. a genealogy with unique leaf labels and a known order of nodes. A fundamental property of the standard coalescent, which follows directly from the exchangeability of genes sampled from the same population, is that all ranked topologies are equally likely (Hudson 1983; Kingman 1982). In other words, if we could somehow assign each mutation to a particular coalescence (i.e. internode) interval, we could use a much simpler GF, defined in terms of the $(n - 1)$ coalescence intervals rather than the $2(n - 1)$ branches for inference. This logic underlies demographic methods that use the branch length information contained in well-resolved genealogies (e.g. Nee *et al.* 1995; Pybus *et al.* 2002) and coalescent-based derivations of the site frequency spectrum (Griffiths and Tavaré 1998; Chen 2012).

Unfortunately however, when analysing sequence data from sexual organisms, we are generally limited by the number of mutations on any one genealogical branch and so often cannot resolve nodes or their order. Although unranked topologies are not equiprobable, even under the standard coalescent, their leaf labels are still exchangeable. Therefore, each unranked, unlabelled topology, or "tree shape" *sensu* Felsenstein (1978, 2003), is an equivalence class that defines a set of identically distributed genealogies (Fig. 1). This means that we only need to work out the GF for one representative (random labeling) per equivalence class. The full GF can then be written as a weighted sum of the GFs for such class representatives:

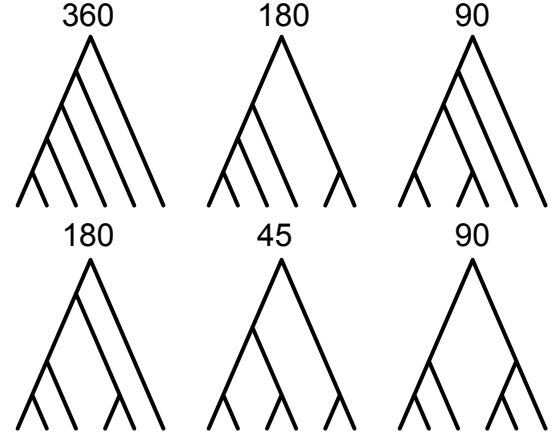
$$\psi[\underline{\omega}] = \sum_h n_h \psi[\underline{\omega}_h] \quad (3)$$

where, n_h denotes the size of equivalence class h and $\underline{\omega}_h \subset \underline{\omega}$ is the set of dummy variables that corresponds to the branches of a single class representative in h . There are necessarily many fewer equivalence classes than labelled topologies (Table 1). For example, given a sample of size $n = 6$ from a single population, there are 945 unranked topologies, but only six equivalence classes (Fig. 1).

Crucially, the idea of tree shapes as equivalence classes extends to any demographic model and sampling scheme. For samples from multiple populations, the equivalence classes are just the permutations of population labels on (unlabelled) tree shapes. It is straightforward to generate and enumerate the equivalence classes (Felsenstein 2003) for any sample. For example, for a sample of $n = 6$ from each of two populations (three per population), there are 49 equivalence classes (partially labelled shapes), which can be found by permuting the two population labels on the unlabelled tree shapes in figure 1.

In general, the size of each equivalence class n_h is a function of the number of permutations of individuals on population labels. For n_i individuals from population i , there are $n_i!$ permutations. Since the orientation of nodes is irrelevant, each

Figure 1 Unranked, unlabelled topologies define equivalence classes of genealogies. For a sample of $n = 6$ from a single population there are six equivalence classes. Their size, i.e. the number of labelled genealogies in each class (n_h) is shown above.



symmetric node in the equivalence class halves the number of unique permutations. Symmetric nodes are connected to identical subclades, that is, there exists an isomorphism ensuring that they have the same topology and the same population labels at the leaves (see Fig. 1).

$$n_h = 1/2^{n_s} \prod_i n_i! \quad (4)$$

where n_s is the number of symmetric nodes. Note also that $\sum_h n_h = (2n - 3)!!$, the total number of unranked topologies.

Any tree shape contains at least one further symmetry: there is at least one node which connects to two leaves. Because the branches descending from that node have the same length by definition, we can combine mutations (and hence ω terms) falling on them: E.g. for a triplet genealogy with topology $(a, (b, c))$, we can combine mutations on branch b and c without loss of information. The joint probability of seeing a configuration with k_b and k_c mutations can be retrieved from $P[k_b + k_c]$:

$$P[k_b, k_c] = \frac{1}{2}^{k_b+k_c} \binom{k_b+k_c}{k_b} P[k_b+k_c] \quad (5)$$

We have previously made use of this in implementing likelihood calculations for triple samples (Lohse *et al.* 2011). Although in principle, this combinatorial argument extends to arbitrary genealogies, one can show that, for larger samples, computing $P[k]$ from mutational configurations defined in terms of internode intervals is computationally wasteful compared to the direct calculation (see File S1).

Approximating models with reversible events

Migration and recombination events are fundamentally different from coalescence and population divergence. Going backwards in time, they do not lead to simpler sample configurations. Thus, the GF for models involving migration and/or recombination is

Table 1 Fundamental quantities of genealogies for small samples (n).

n	branches	ranked topologies	unranked topologies	EC ^a 1 pop.	EC 2 pop.	# of config ^b
	$2^n - 2$	$\frac{(n!(n-1)!)}{2^{(n-1)}}$	$(2n - 3)!!$	(Felsenstein 2003)		$(2 + k_m)^{2(n-1)}$
3	6	3	3	1	2	625
4	14	18	15	2	6	15625
6	62	2,700	945	6	49	9765625
8	254	1,587,600	135,135	23	560	6103515625
10	1022	2,571,912,000	34,469,425	98	7,139	3814697265625

^a the number of equivalence classes

^b the number of mutational configurations for a sample from 2 populations with up to $k_m = 3$ per mutations per branch.

a system of coupled equations, the solution of which involves matrix inversion and higher order polynomials and quickly becomes infeasible for large n (Hobolth *et al.* 2011). As an example, we consider two populations connected by symmetric migration at rate $M = 4Nm$. Since we are often interested in histories with low or moderate migration in practice, it seems reasonable to consider an approximate model in which the number of migration events is limited. Using a Taylor series expansion, the full GF can be decomposed into histories with $1, 2, \dots, n$ migration events (Lohse *et al.* 2011). The same argument applies to recombination between discrete loci and can be used to derive the GF for the sequential Markov coalescent (McVean and Cardin 2005). It is crucial to distinguish between M terms in the numerator and denominator. In other words, even if we stop including sampling configurations involving multiple migration events, M still contributes to the total rate $\sum_i \lambda_i$ in the denominator. To see how this works, we consider the simplest case of a pair of genes a and b sampled from two populations connected by symmetric migration. Following Lohse *et al.* (2011), $a \setminus b$ denotes the sampling configuration where both genes are in different populations and $a, b \setminus \emptyset$ where they are in the same population. We modify the GF (Lohse *et al.* 2011, eq. 9) to include an indicator variable γ that counts the number of migration events:

$$\begin{aligned} \psi^*[a \setminus b] &= \frac{\gamma M}{(M + \omega_a + \omega_b)} \psi^*[a, b \setminus \emptyset] \\ \psi^*[a, b \setminus \emptyset] &= \frac{1}{(1 + M + \omega_a + \omega_b)} (1 + \gamma M \psi^*[a \setminus b]) \end{aligned} \quad (6)$$

Expanding ψ^* in γ , the coefficients of $\gamma, \gamma^2, \dots, \gamma^{M_{max}}$ correspond to histories with $1, 2, \dots, M_{max}$ migration events. This is analogous to conditioning on a particular topology: the truncated GF does not sum to one (if we set the ω to zero), but rather gives the total probability of seeing no more than M_{max} events. This is convenient because it immediately gives an estimate of the accuracy of the approximation. Expanding the solution of equation 6 around $\gamma = 0$ gives:

$$\psi^*[a \setminus b] = \sum_i \frac{M^i}{((M + \omega_a + \omega_b)(1 + M + \omega_a + \omega_b))^{(i+1)/2}} \quad (7)$$

The GF conditional on there being at most one migration event is

$$\psi^*[a \setminus b | M_{max} = 1] = \frac{M}{(M + \omega_a + \omega_b)(1 + M + \omega_a + \omega_b)} \quad (8)$$

The error of this approximation is:

$$1 - \psi^*[a \setminus b | M_{max}] = 1_{\omega \rightarrow 0} = \frac{M}{M + 1} \quad (9)$$

which is just the chance that a migration event occurs before coalescence (see Fig. 2). An analogous expansion for the pairwise GF for the IM model (Lohse *et al.* 2011, eq. 13) gives:

$$\psi^*[a \setminus b | M_{max}] = \frac{1}{2} \left(2Me^{-MT} + \frac{2}{1 + M} - \frac{2e^{-(M+1)TM^2}}{1 + M} \right) \quad (10)$$

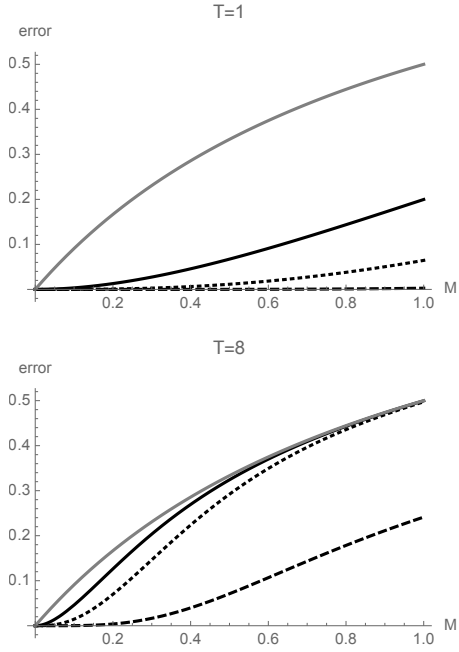
Expressions for the GF conditional on a maximum of $2, 3, \dots, n$ migration events and for larger samples can be found by automating the GF recursion. While these do not appear to have a simple form, plotting the error against M and T (Fig. 2), shows that for recent divergence ($T < 1$) and moderate gene flow ($M < 0.5$), histories involving more than two migration events are extremely unlikely ($p < 0.01$) and can be ignored to a good approximation. Considering that for large n , coalescence (at rate $n(n-1)/2$) becomes much more likely than migration (at rate Mn), this approximation should be relatively robust to sample size.

Unknown phase and root

There are at least two further complications for blockwise likelihood computations in practice: First, we have so far assumed that mutations can be polarized without error, i.e. that the infinite sites mutation model holds between in and outgroup, which is often unrealistic in practice. Second, given the current limitations of short read sequencing technology, genomic data are often unphased and one would ideally like to incorporate phase ambiguity explicitly rather than ignore it (e.g. Lohse and Frantz 2014) or rely on computational phasing.

Both unknown phase and root can be incorporated via a simple relabeling of branches. In generating the GF, we have labelled branches and corresponding ω variables by the tips (leaf-nodes) they are connected to. Crucially, the full GF expressed as a sum over equivalence class representatives (eq. 9) still has unique labels for all individuals. That is, we distinguish genes sampled from the same population. To incorporate unknown phase, we simply label leaf nodes by the population they were sampled from (Fig. 3). Because each genealogical branches are labelled by the set of leaf nodes they are connected to, this relabeling of leaf nodes defines branch types that correspond to

Figure 2 The error (eqs. 9) in limiting the number of migration events to $M_{max}=1$ (solid), 2 (dashed) and 4 (dotted) for a pairwise sample in the IM model plotted against M for different divergence times T . The results for a model of equilibrium migration without divergence is shown for comparison (grey).



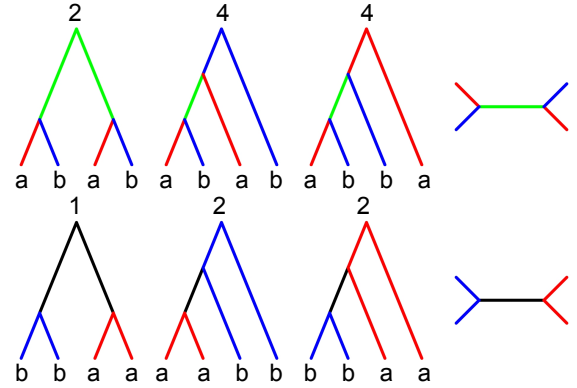
categories of the (joint) site frequency spectrum (SFS). In other words, in the absence of phase information branch types are defined by the number of descendants in each population. To see how this works, consider for example two genes from each of two populations. There are six equivalence classes of rooted genealogies (Fig. 3). Combining branches with the same population labels gives seven ω variables that correspond to site types: $\omega_a, \omega_b, \omega_{ab}, \omega_{aa}, \omega_{bb}, \omega_{aab}, \omega_{abb}$. In the absence of root information, we further combine the two branches on either side of the root. Denoting ω variables for unrooted branches by $*$ and the two sets of individuals they are connected to we have: $\omega_{\{a,abb\}}^*$, $\omega_{\{b,abb\}}^*$, $\omega_{\{ab,ab\}}^*$, $\omega_{\{aa,bb\}}^*$. The rooted branches contributing to each unrooted branch are indicated in colour in figure 3. The ω^* terms correspond to the four types of variable sites defined by the folded SFS for two populations: $k_{\{a,abb\}}^*$ (heterozygous sites unique to a), $k_{\{b,abb\}}^*$ (heterozygous sites unique to b), $k_{\{ab,ab\}}^*$ (heterozygous sites shared by both) and $k_{\{aa,bb\}}^*$ (fixed differences between a and b). Note also that without the root, the six equivalence classes collapse to two unrooted equivalence classes (defined by branches $t_{\{aa,bb\}}^*$ and $t_{\{ab,ab\}}^*$) (Fig. 3).

The combinatorial arguments outlined above extend to arbitrary sample sizes and numbers of populations. The GF for unphased data is given by combining ω variables with the same number of descendants in each population. We modify eq. 3 to write the GF of an unrooted genealogy $\psi[\underline{\omega}^*]$ as a sum over unrooted equivalence classes (denoted h^*), each of which is in turn a sum over rooted equivalence classes:

$$\psi[\underline{\omega}^*] = \sum_{h^*} \sum_{h \in h^*} n_h \psi[\underline{\omega}_h \rightarrow \underline{\omega}_h^*] \quad (11)$$

We can use this simplified GF and equation 1 to compute the

Figure 3 For a sample of two sequences (a diploid genome) from each of two populations (a and b), there are six classes of equivalent, rooted genealogies (left); their sizes n_h are shown above. Without root information, these collapse to two unrooted genealogies (right). Without phase information, there are four mutation types that map to specific branches in the rooted genealogy: heterozygous sites unique to one sample ($t_{\{a,abb\}}^*$ and $t_{\{b,abb\}}^*$, red and blue respectively), shared heterozygous sites ($t_{\{ab,ab\}}^*$, green) and fixed, homozygous differences ($t_{\{aa,bb\}}^*$, black).



probability of blockwise counts of mutation types defined by the joint SFS. We will refer to this extension of the joint SFS to blockwise data as the blockwise site frequency spectrum (bSFS), following [Bunnfeld et al. \(2015\)](#) who have used the bSFS to fit bottleneck histories in a single population.

Limiting the total number of mutational configurations

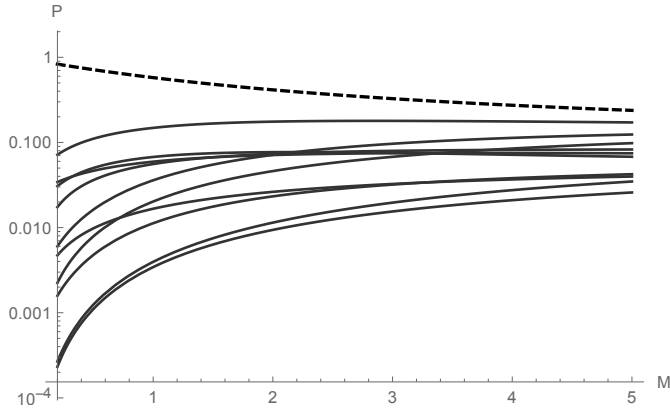
In principle, we can compute the probability of seeing arbitrarily many mutations on a particular branch from equation 1. In practice however, the extra information gained by explicitly distinguishing configurations with large numbers of mutations (which are very unlikely for short blocks) is limited, while the computational cost increases. An obvious strategy is to tabulate exact probabilities only up to a certain maximum number of mutations k_m per branch and combine residual probabilities for configurations involving more than k_m mutations on one or multiple branches. As described by [Lohse et al. \(2011\)](#) and [Lohse et al. \(2012\)](#), the residual probability of seeing more than k_m mutations on a particular branch s is given by

$$P[k_s \geq k_m] = \psi[\underline{\omega}]|_{\omega_s \rightarrow 0} - \sum_{i=0}^{k_m} P[k_s = i]$$

i.e. we subtract the sum of exact probabilities for configurations involving up to k_m mutations from the marginal probability of seeing branch s .

Assuming that we want to distinguish between all $2(n-1)$ branches in a given equivalence class and use a global k_m for all branches, there are (k_m+2) possible mutation counts per branch (including those with no mutations or more than k_m mutations on a branch) which gives $(k_m+2)^{2(n-1)}$ mutational configurations in total. For example, for $n=6$ and $k_m=3$

Figure 4 The topology spectrum for a sample of $n = 6$ from a two population IM model with asymmetric migration and $T = 1.5$. The probabilities of all 11 unrooted topologies are plotted against M . The probability of the most likely topology of reciprocal monophyly $((a, (a, a)), (b, (b, b)))$ is shown as a dashed line.



there are 9,765,625 mutational configurations per equivalence class (Table 1). Although this may seem daunting, most of these configurations are extremely unlikely, so a substantial computational saving can be made by choosing branch-specific k_m . We have implemented functions in *Mathematica* to tabulate $P[k]$ for an arbitrary vector of k_m (File S1).

The bSFS with $k_m = 0$ constitutes an interesting special case: it defines mutational configurations by the joint presence and/or absence of SFS types in a block, irrespective of their number. In the limit of very long blocks, i.e. if we assume an unlimited supply of mutations, this converges to the topological probabilities of equivalence classes which can be obtained directly from the partitioned GF by setting all $\omega \rightarrow 0$. We can think of this set of probabilities as the "topology spectrum". For a sample of 3 genes from each of 2 populations this consists of 49 equivalence classes which reduce to 11 unrooted topologies (Fig. 3). Under the IM model with unidirectional migration, the GF of each class is solvable using *Mathematica* (File S2). The most likely topology is reciprocal monophyly, i.e. $((a, (a, a)), ((b, b), b))$. As expected, its probability decreases with M and increases with T .

Data Availability

File S1 is a *Mathematica* notebook that contains the code to generate the GF and tabulate likelihoods under arbitrary demographic models. File S2 contains the code used for the analyses of the IM model, including the analyses of the *Heliconius* data and the power test. The processed input data for *Heliconius* and python scripts used are available from www.datadryad.com doi:XXX; raw sequence data are published by [Martin et al. \(2013\)](#) and available from www.datadryad.com doi:10.5061/dryad.dk712.

Results

The various combinatorial strategies for simplifying likelihood calculations based on the GF outlined above suggest a general "pipeline", each component of which can be automated:

1. Generate all equivalence classes h and enumerate their sizes n_h for a given sampling scheme.

2. Generate and solve the GF conditional on one representative within each h .
3. Take the Inverse Laplace Transform with respect to the time parameters of discrete events (e.g. divergence, admixture, bottlenecks). These processes are initially modelled as occurring with a continuous rate.
4. Re-label ω variables to combine branches and equivalence classes that are indistinguishable in the absence of root and/or phase information.
5. Find a sensible k_m for each mutation type from the data.
6. Tabulate probabilities for all mutational configurations in each equivalence class.

In the accompanying *Mathematica* notebook we have implemented this pipeline as a set of general functions. These can be used to automatically generate, solve and simplify the GF (step 1–3), and – from this – tabulate $P[k]$, the likelihood of a large range of demographic models (involving population divergence, admixture and bottlenecks) (step 6). In principle, this automation works for arbitrary sample sizes. In practice however, the inversion step (step 3) and the tabulation of probabilities (step 6) become infeasible for $n > 6$.

To give a concrete example, we derive the GF for a model of isolation at time T (scaled in $2N_e$ generations) with migration (at rate $M = 4N_e m$ migrants per generation) (IM) between two populations (labelled a and b). We further assume that migration is unidirectional, i.e. from a to b forwards in time and that both populations and their common ancestral population are of the same effective size (we later relax this assumption when analysing data). As above, we consider the special case of a single diploid sample per population without root and phase information. We first derive some basic properties of unrooted genealogies under this model. We then investigate the power of likelihood calculations based on the bSFS. Finally, we apply this likelihood calculation to genome data from two species of *Heliconius* butterflies.

The distribution of unrooted branches under the IM model

We can find the expected length of any branch (or combination of branches) s from the GF as: $E[t_s] = -\partial\psi[\omega]/\partial\omega_s|_{\omega \rightarrow 0}$ ([Lohse et al. 2011](#)). The expressions for the expected lengths of rooted branches are cumbersome (File S2). Surprisingly however, the expected lengths of the four unrooted branches $t_{\{aa,bb\}}^*$, $t_{\{ab,ab\}}^*$, $t_{\{a,abb\}}^*$ and $t_{\{b,aab\}}^*$, each of which is a sum over the underlying rooted branches (Fig. 3), have a relatively simple form (Fig. 5):

$$E[t_{\{aa,bb\}}^*] = \frac{e^{-(2+M)T}(-6e^T M^2 - 24e^{\frac{1}{2}(4+M)T}(1+M) + 2(1+M) + e^{(2+M)T} + (24 + 24M + 7M^2 + M^3))}{3M(1+M)(2+M)}$$

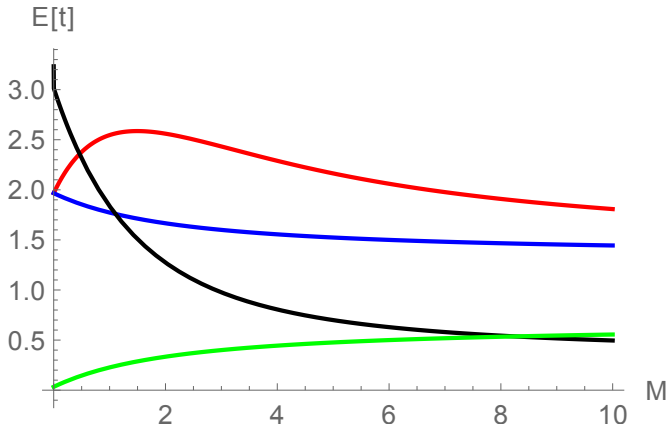
$$E[t_{\{ab,ab\}}^*] = \frac{2(2e^{-(2+M)T} + M)}{3(2+M)}$$

$$E[t_{\{a,abb\}}^*] = \frac{4e^{-(2+M)T}(3e^T M - 1 - M - 6e^{\frac{1}{2}(4+M)T}(1+M) + e^{(2+M)T}(9 + 7M + 7M^2))}{3M(1+M)(2+M)}$$

$$E[t_{\{b,aab\}}^*] = \frac{4(3 - e^{-(2+M)T} + M)}{3(2+M)} \quad (12)$$

Similarly, the probability of the two unrooted topologies reduces to:

Figure 5 The expected length of unrooted genealogical branches (eq. 12) for a sample of $n = 4$ under the IM model of two populations (a and b) with asymmetric migration and population divergence time $T = 1.5$ ($\times 2N_e$ generations). Colours correspond to those in figure 3.



$$p[t_{\{aa,bb\}}^*] = \frac{4e^{(2+M)T} + 2M}{3(2+M)} \quad (13)$$

$$p[t_{\{ab,ab\}}^*] = 1 - p[t_{\{aa,bb\}}^*]$$

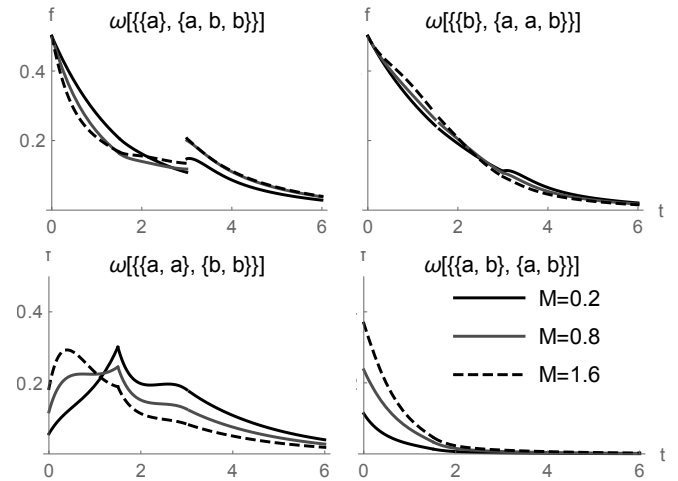
We can recover the full distribution of rooted branches from the GF by taking the Inverse Laplace Transform (using *Mathematica*) with respect to the corresponding ω^* . While this does not yield simple expressions (File S2), examining figure 6 illustrates that much of the information about population history is contained in the shape of the branch length distribution rather than its expectation (Fig. 5). For example, the length of branches carrying fixed differences $t_{\{aa,bb\}}^*$ has a multi-modal distribution with discontinuities at T and the relative size of the first mode strongly dependent on M .

Power analysis

We compared the power to detect post-divergence gene flow between two different blockwise likelihood calculations: the bSFS for a diploid genome per population ($n = 4$) and a minimal sample of a single haploid sequence ($n = 2$) per population. As a proxy for power, we computed the expected difference in support ($E[\Delta \ln L]$) between the IM model and a null model of strict divergence without gene flow and arbitrarily assumed datasets of 100 blocks. However, since we are assuming that blocks are unlinked, i.e. statistically independent, $E[\Delta \ln L]$ scales linearly with the number of blocks.

Figure 7 shows the power to detect gene flow for a relatively old split ($T = 1.5$) and sampling blocks with an average of $\theta = 4N_e\mu = 1.5$ heterozygous sites within each species. Without gene flow, this corresponds to a total number of 5.2 mutations per block on average (using eq. 12 and $E[S_T] = \theta/2 \sum E[t^*]$). Unsurprisingly, sampling a diploid sequence from each population gives greater power to detect gene flow than pairwise samples (compare black and blue lines in figure 7). However, contrasting this with the power of a simpler likelihood calculation for $n = 4$ which is based only on the total number of mutations S_T in each block (grey line in figure 7), illustrates that the additional information does not stem from the increase in sample size *per se*, but

Figure 6 The length distribution of unrooted genealogical branches for a sample of $n = 4$ under the IM model of two populations (a and b) with asymmetric migration and population divergence at $T = 1.5$ (in $2N_e$ generations).



rather the addition of topology information. Perhaps counterintuitively, we find that there is less information in the distribution of S_T for larger samples than in pairwise samples. This clearly shows that most information about post-divergence gene flow is contained in the topology, i.e. being able to assign mutations to specific branches. Similarly, adding root information almost doubles power (green lines in Fig. 7).

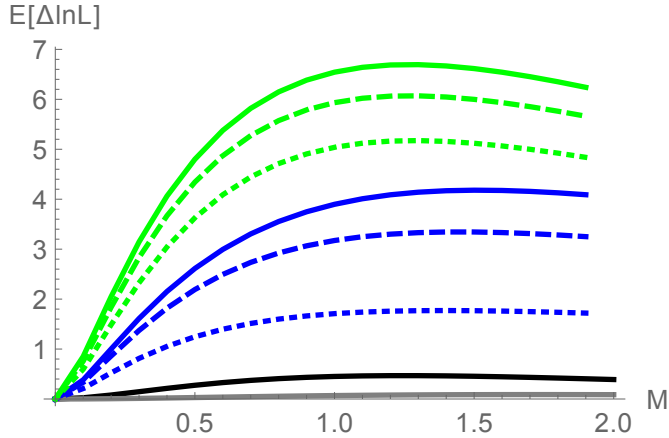
In comparison, the threshold k_m has relatively little effect on power. In other words, for realistically short blocks, most of the information is contained in the joint presence and absence of mutation types (regardless of their number).

Heliconius analysis

To illustrate likelihood calculation based on the bSFS, we estimated divergence and gene flow between two species of *Heliconius* butterflies. The sister species *H. cydno* and *H. melpomene rosina* occur in sympatry in parts of Central and South America, are known to hybridise in the wild at a low rate (Mallet *et al.* 2007), and have previously been shown to have experienced post-divergence gene flow (Martin *et al.* 2013). We sampled 225 bp blocks of intergenic, autosomal sequence for one individual genome of each species from the area of sympatry in Panama (chi565 and ro2071). These data are part of a larger resequencing study involving high coverage genomes for four individuals of each *H. cydno* and *H. m. rosina* as well as an allopatric population of *H. melpomene* from French Guiana (Martin *et al.* 2013). We excluded CpG islands and sites with low quality ($GQ < 30$ and $MQ < 30$), excessively low (< 10) or high (> 200) coverage and only considered sites that passed these filtering criteria in all individuals.

We partitioned the intergenic sequence into blocks of 225bp length. To allow some sites to violate our filtering criteria in each block whilst keeping the block length post-filtering fixed, we sampled the first 150 bases passing filtering in each block (blocks with fewer remaining sites were excluded from the analysis). 6.3% of blocks violated the 4-gamete criterion (i.e. contained both fixed differences and shared heterozygous sites) and were removed. This sampling and filtering strategy yielded 161,726 blocks with an average per site heterozygosity of 0.017 and 0.015

Figure 7 The power ($E[\Delta \ln L]$) to distinguish between an IM model and a null model of strict divergence ($T = 1.5$) from 100 unlinked blocks of length $\theta = 1.5$ for different sample sizes and data summaries: the total number of mutations in a sample of $n = 2$ (black) and $n = 4$ (grey), the bSFS for unphased data for two diploids ($n = 4$) with root (green) and without root (blue). Dotted, dashed and solid lines correspond to different maximum numbers of mutations per branch type, $k_m = 0, 1$ and 3 respectively.



in *H. m. rosina* and *H. cydno* respectively (Fig. 8). Summarizing the data by counting the four mutation types in each block gave a total of 2,337 unique mutational configurations, 1,743 of which occurred more than once.

We initially used all blocks (regardless of linkage) to obtain point estimates of parameters under three models: i) strict isolation without migration ii) isolation with migration from *H. cydno* into *H. m. rosina* ($IM_{c \rightarrow m}$) and iii) isolation with migration from *H. m. rosina* into *H. cydno* ($IM_{m \rightarrow c}$). In all cases, we assumed that the common ancestral population shared its N_e with one descendant species while the other descendant was assumed to have a different N_e . To keep computation times manageable, we did not consider more complex histories involving bi-directional gene flow or three N_e parameters.

We maximise $\ln L$ under each model using Nelder-Mead simplex optimisation implemented in the *Mathematica* function *NMaximize*. Confidence intervals (CI) for parameter estimates were obtained from 100 parametric bootstrap replicates. We used *ms* (Hudson 2002) to simulate 0.3Mb of contiguous sequence for each of the 20 *Heliconius* autosomes assuming a per site recombination rate of 1.8×10^{-9} (Jiggins et al. 2005) and the best fitting IM history. We partitioned each simulated dataset into 150bp blocks and estimated 95 % (CI) as two standard deviations of estimates across bootstrap replicates (see Discussion).

We find strong support for a model of isolation with migration from *H. cydno* into *H. m. rosina* ($IM_{c \rightarrow m}$) (Table 2) with a larger N_e in *H. cydno*. This model fits significantly better than simpler nested models of strict divergence or an IM model with a single N_e (Table 2). Our results agree with earlier genomic analyses of these species that showed support for post-divergence gene flow based on D-statistics (Martin et al. 2013), IMA analyses based on smaller numbers of loci (Kronforst et al. 2013) and genome wide SNP frequencies analysed using approximate Bayesian computation. Asymmetrical migration from *H. cydno* into *H. m. rosina* has also been reported previously, and could be

explained by the fact that F1 hybrids resemble *H. m. rosina* more closely due to dominance relationships among wing patterning alleles, possibly making F1s more attractive to *H. m. rosina* (Kronforst et al. 2006; Martin et al. 2015).

We applied a recent direct, genome-wide estimate of the mutation rate for *H. melpomene* of 2.9×10^{-9} per site and generation (Keightley et al. 2015) to convert parameter estimates into absolute values. Assuming that synonymous coding sites are evolving neutrally, we used the ratio of divergence between *H. m. rosina* and the more distantly-related 'silvaniform' clade of *Heliconius* at synonymous coding sites and the intergenic sites our analysis is based on to estimate the selective constraint on the latter. This gives an "effective mutation" rate of $\mu = 1.9 \times 10^{-9}$ (Martin et al. 2015). Applying this corrected rate to our estimate of θ and assuming four generations per year, we obtain an N_e estimate of 1.10×10^6 for *H. m. rosina* and the common ancestral population and 2.85×10^6 for *H. cydno*. We estimate species divergence to have occurred roughly 0.91 – 1.18 MY ago. Note that this is more recent than previous estimates of 1.5 million years which were obtained using approximate Bayesian computation and a different calibration based on mitochondrial genealogies (Kronforst et al. 2013; Martin et al. 2015).

Discussion

We have shown how the probabilities of genealogies, and hence of mutational configurations, can be calculated for a wide variety of demographic models. This gives an efficient way to infer demography from whole genome data. Irrespective of any particular demographic history, the possible genealogies of a sample can be partitioned into a set of equivalence classes, which are given by permuting population labels on tree shapes. We show how this fundamental symmetry of the coalescent can be exploited when computing likelihoods from blockwise mutational configurations. We have implemented this combinatorial partitioning in *Mathematica* to automatically generate and solve the generating function (GF) of the genealogy and, from this, compute likelihoods for a wide range of demographic models. Given a particular sample of genomes, we first generate a set of equivalence classes of genealogies and condition the recursion for the GF (Lohse et al. 2011) on a single representative from each class. This combinatorial strategy brings a huge computational saving. Importantly, it does not sacrifice any information. In contrast, approximating the GF for models that include reversible events in particular migration, involves a trade-off between computational efficiency and accuracy. For example, given our high estimates for unidirectional M for *Heliconius* (Table 2), it would have been unrealistic to fit a history of bi-directional migration to these data without allowing for multiple migration events in each genealogy (Figure 2).

Although these approaches make it possible to solve the GF for surprisingly large samples and biologically interesting models, the number of mutational configurations (which explodes with the number of sampled genomes) remains a fundamental limitation of such likelihood calculations in practice. Given outgroup and phase, the full information is contained in a vast table of mutational configurations which are defined in terms of the $2(n-1)$ branches of each equivalence class. For samples from two populations, the number of mutational configurations we need to calculate is the product of the last two columns of Table 1. For example, given a sample of three haploid genomes per populations and allowing for up to $k_m = 3$ mutations per branch, there are $49 \times 9,765,625 = 478,515,625$ possible muta-

Table 2 Maximum likelihood estimates of divergence and migration between *H. m. rosina* and *H. cydno*.

	θ (N_e)	θ_C (N_e)	T (τ)	M
IM estimates ^a	1.25	3.24	1.90	1.50
Scaled IM estimates ^b	1.10 (1.02 – 1.18)	2.85 (2.55 – 3.23)	1.04 (0.91 – 1.18) MY	(1.32 – 1.68)
Expected estimates ^c	1.22	3.53	1.97	1.40

^a under the best model $IM_{c \rightarrow m}$. θ_C is the scaled mutation rate in *H. cydno*

^b N_e in $\times 10^6$ individuals, τ in MY, 95 % CI in brackets

^c Mean across parametric bootstrap replicates

tional configurations, an unrealistic number of probabilities to calculate.

The blockwise site frequency spectrum

Our initial motivation for studying the bSFS was to deal with unphased data in practice. The GF of the bSFS can be obtained from the full GF simply by combining branches with equivalent leaf labels. As well as being a lossless summary of blockwise data (in the absence of phase information), the bSFS is a promising summary in general for several reasons. First, it is extremely compact compared to the full set of (phased) mutational configurations. Unlike the latter, the size of the bSFS does not depend on the number of equivalence classes (which explodes with n , Table 1), but only on n . Given a sample of n_i individuals from population i and assuming a global maximum number of mutations k_m for all mutation types, the (unfolded) bSFS comprises of a maximum of $((\prod_i (n_i + 1)) - 2)^{(k_m + 2)}$ mutational configurations. For a sample of 3 haploid genomes from each of two populations and $k_m = 3$, the bSFS has $7^5 = 16,807$ entries. Second, because equivalence classes of genealogies are defined by the presence and absence of SFS types, much of the topology information contained in the full data should still be captured in the bSFS. Finally, and perhaps surprisingly, at least for the IM model the expressions for the total length of branches contributing to unphased and unpolarized mutation types (eq. 12 & 13) are much simpler than those of the underlying rooted branches, which suggests that it may be possible to find general results.

Despite the strategies developed here, it is clear that full likelihood calculations will rarely be feasible for samples > 6 given the rapid increase in the number of equivalence classes. However, a separation of time-scales exists for many models of geographic and genetic structure (Wakeley 1998, 2009), and so full likelihood solutions for moderate ($n < 6$) samples may be sufficient for computing likelihoods for much larger samples if these contribute mainly very short branches with no mutations in the initial scattering phase during which lineages from the same population either coalesce or trace back to unsampled demes.

Dealing with linkage

A key assumption of the blockwise likelihood calculations is that there is no recombination within sequence blocks, and that different blocks are independent of each other. This latter assumption is especially problematic when we analyse whole-genome data. If we divide the genome into blocks that are small enough for recombination within them to be negligible, our method gives an unbiased estimate of the likelihood of a demographic model. However, the accuracy of the model fit will be grossly overestimated if we simply multiply likelihoods across blocks, because

adjacent blocks are strongly correlated. Ignoring this correlation amounts to a composite likelihood calculation.

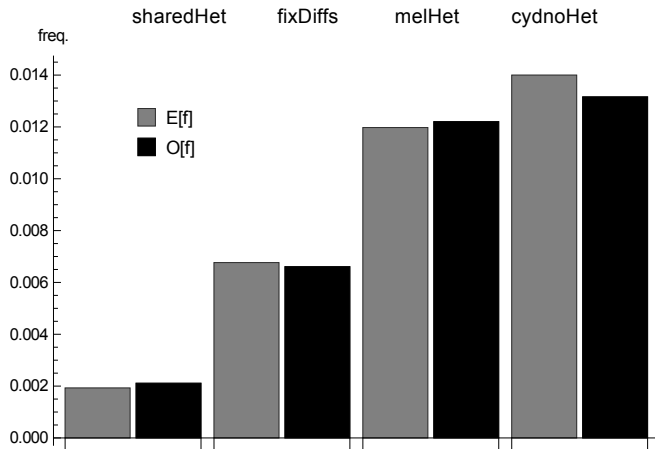
A common practice (e.g. Wang and Hey 2010; Excoffier et al. 2013; Lohse and Frantz 2014) is to assume a "safe distance" at which blocks (or SNPs) are statistically independent. This amounts to a rescaling of the $\ln L$: suppose that we multiply likelihoods across every k th block, k being chosen large enough that blocks are uncorrelated. This procedure is valid starting at any block, and so can be repeated k times, such that the whole genome is included in the analysis. Taking the average across all k analyses is equivalent to simply multiplying the likelihoods across all blocks, and then dividing the total $\ln L$ by k . However, because it is not clear how to choose k , this procedure is quite arbitrary. On the one hand, successive blocks or SNPs are not completely correlated, suggesting that this considerably underestimates the accuracy of estimates. On the other hand, however, there may be weak, long-range correlations, due to a small fraction of long regions that coalesced recently, and these may increase the variance of parameter estimates.

The safest way to account for LD is via a parametric bootstrap. Although computationally intensive, this has the added benefit that it also checks whether parameter estimates are biased (due to the assumption of no recombination within blocks). It is reassuring that in the case of the *Heliconius* data, we find that the biases in parameter estimates are very small indeed (last row in Table 2). It is important to note that the confidence intervals for the *Heliconius* estimates we derived from the bootstrap are conservative given the current limitations of coalescent simulators. Given of the limited length of continuous recombining sequence that can be simulated, the simulated datasets were over four times smaller than the data. An interesting alternative to full parametric bootstrap, which we hope to implement in the future, is to use the variance of the gradient of $\ln L$ across bootstrap replicates to adjust the Fisher Information matrix (Godambe 1960; Coffman et al. 2015).

An advantage of direct likelihood calculations is that one can easily check the absolute fit of the data to a model by asking how well the observed frequency of mutational configurations or some summary such as the SFS is predicted by the model. For example, the IM history we estimated for the two *Heliconius* species fits the observed genome-wide SFS reasonably well (Fig. 8). The fact that we slightly underestimate the heterozygosity in *H. cydno* may suggest that some process (e.g. demographic change after divergence or admixture from an unsampled ghost population/species) is not captured by our model.

In general, the GF framework makes it possible to derive the distribution of any summary statistic that can be defined as a combination of genealogical branches and understand its properties under simple demographic models and small n . Although explicit calculations based on such summaries are not feasible

Figure 8 The folded SFS has four site types: i) heterozygous sites unique to either *H. m. melpomene* or ii) *H. cydno* iii) shared by both species and iv) fixed differences. The observed genome-wide SFS is shown in black, the expectation under the IM history estimated from the bSFS (Table 2) (eq. 9) in grey.



for large n , summary statistics such as the bSFS may still have wide applicability for fitting complex models and larger samples of individuals, for example using approximate likelihood methods, or simply as a way to visualize how genealogies vary along the genome.

Acknowledgements

This work was supported by funding from the UK Natural Environment Research Council to KL (NE/I020288/1) and a grant from the European Research Council (250152) to NB. We thank Lynsey Bunnefeld for discussions throughout the project and Joshua Schraiber and one anonymous reviewer for constructive comments on an earlier version of this manuscript.

Literature Cited

- Bunnefeld, L., L. A. F. Frantz, and K. Lohse, 2015 Inferring bottlenecks from genome-wide samples of short sequence blocks. *Genetics* **201**: 1157–1169.
- Chen, H., 2012 The joint allele frequency spectrum of multiple populations: A coalescent theory approach. *Theoretical Population Biology* **81**: 179–195.
- Coffman, A. J., P. Hsieh, S. Gravel, and R. N. Gutenkunst, 2015 Computationally efficient composite likelihood statistics for demographic inference. *Molecular Biology and Evolution*.
- Davey, J. W. and M. L. Blaxter, 2011 RADseq: next-generation population genetics. *Briefings in Functional Genomics* **9**: 416–423.
- Edwards, A. W. F., 1970 Estimation of the branch points of a branching diffusion process (with discussion). *J. R. Stat. Soc. B* **32**: 155–174.
- Excoffier, L., I. Dupanloup, E. Huerta-Sánchez, V. C. Sousa, and M. Foll, 2013 Robust demographic inference from genomic and snp data. *PLoS Genet* **9**: e1003905.
- Felsenstein, J., 1978 The number of evolutionary trees. *Molecular Phylogenetics and Evolution* **27**: 27–33.
- Felsenstein, J., 1988 Phylogenies from molecular sequences: Inference and reliability. *Annu Rev Genet* **22**: 521–565.

- Felsenstein, J., 2003 *Inferring phylogenies*. Sinauer Associates, Sunderland, Massachusetts.
- Godambe, 1960 An optimum property of regular maximum likelihood estimation. *Ann. Math. Stat* **31**: 1208–1211.
- Green, R. E., J. Krause, A. W. Briggs, T. Maricic, U. Stenzel, M. Kircher, N. Patterson, H. Li, W. Zhai, M. H.-Y. Fritz, N. F. Hansen, E. Y. Durand, A.-S. Malaspina, J. D. Jensen, T. Marques-Bonet, C. Alkan, K. Prufer, M. Meyer, H. A. Burbano, J. M. Good, R. Schultz, A. Aximu-Petri, A. Butthof, B. Hober, B. Hoffner, M. Siegemund, A. Weihmann, C. Nusbaum, E. S. Lander, C. Russ, N. Novod, J. Affourtit, M. Egholm, C. Verna, P. Rudan, D. Brajkovic, Z. Kucan, I. Gusic, V. B. Doronichev, L. V. Golovanova, C. Lalueza-Fox, M. de la Rasilla, J. Fortea, A. Rosas, R. W. Schmitz, P. L. F. Johnson, E. E. Eichler, D. Falush, E. Birney, J. C. Mullikin, M. Slatkin, R. Nielsen, J. Kelso, M. Lachmann, D. Reich, and S. Pääbo, 2010 A draft sequence of the Neanderthal genome. *Science* **328**: 710–722.
- Griffiths, R. and S. Tavaré, 1998 The age of a mutation in a general coalescent tree. *Communications in Statistics. Stochastic Models* **14**: 273–295.
- Gutenkunst, R. N., R. D. Hernandez, S. H. Williamson, and B. C. D., 2009 Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genetics* **5**: e1000695.
- Harris, K. and R. Nielsen, 2013 Inferring demographic history from a spectrum of shared haplotype lengths. *PLoS Genet* **9**: e1003521.
- Hearn, J., G. N. Stone, L. Bunnefeld, J. A. Nicholls, N. H. Barton, and K. Lohse, 2014 Likelihood-based inference of population history from low-coverage de novo genome assemblies. *Molecular Ecology* **23**: 198–211.
- Hey, J. and R. Nielsen, 2004 Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*. *Genetics* **167**: 747–760.
- Hobolth, A., L. N. Andersen, and T. Mailund, 2011 On computing the coalescent time density in an isolation-with-migration model with few samples. *Genetics* **187**: 1241–1243.
- Hudson, R. R., 1983 Testing the constant-rate neutral allele model with protein sequence data. *Evolution* **37**: 203–217.
- Hudson, R. R., 2002 Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics* **18**: 337–338.
- Jiggins, C. D., J. Mavarez, M. Beltrán, W. O. McMillan, J. S. Johnston, and E. Bermingham, 2005 A genetic linkage map of the mimetic butterfly *Heliconius melpomene*. *Genetics* **171**: 557–570.
- Keightley, P. D., A. Pinharanda, R. W. Ness, F. Simpson, K. K. Dasmahapatra, J. Mallet, J. W. Davey, and C. D. Jiggins, 2015 Estimation of the spontaneous mutation rate in *Heliconius melpomene*. *Molecular Biology and Evolution* **32**: 239–243.
- Kingman, J. F. C., 1982 The coalescent. *Stochastic Processes and their Applications* **13**: 235–248.
- Kronforst, M., M. Hansen, N. Crawford, J. Gallant, W. Zhang, R. Kulathinal, D. Kapan, and S. Mullen, 2013 Hybridization reveals the evolving genomic architecture of speciation. *Cell Reports* **5**: 666–677.
- Kronforst, M. R., L. G. Young, L. M. Blume, and L. E. Gilbert, 2006 Multilocus analyses of admixture and introgression among hybridizing *Heliconius* butterflies. *Evolution* **60**: 1254–1268.
- Li, H. and R. Durbin, 2011 Inference of human population history from individual whole-genome sequences. *Nature* **475**: 493–6.

- Lohse, K., N. H. Barton, N. Melika, and G. N. Stone, 2012 A likelihood-based comparison of population histories in a parasitoid guild. *Molecular Ecology* **49**: 832–842.
- Lohse, K. and L. A. F. Frantz, 2014 Neandertal admixture in eurasia confirmed by maximum-likelihood analysis of three genomes. *Genetics* **196**: 1241–1251.
- Lohse, K., R. J. Harrison, and N. H. Barton, 2011 A general method for calculating likelihoods under the coalescent process. *Genetics* **58**: 977–987.
- Mailund, T., A. E. Halager, M. Westergaard, J. Y. Dutheil, K. Munch, L. N. Andersen, G. Lunter, K. Püfer, A. Scally, A. Hobolth, and M. H. Schierup, 2012 A new isolation with migration model along complete genomes infers very different divergence processes among closely related great ape species. *PLoS Genetics* **8**: e1003125.
- Mallet, J., M. Beltran, W. Neukirchen, and M. Linares, 2007 Natural hybridization in heliconiine butterflies: the species boundary as a continuum. *BMC Evolutionary Biology* **7**: 28.
- Martin, S. H., K. K. Dasmahapatra, N. J. Nadeau, C. Salazar, J. R. Walters, F. Simpson, M. Blaxter, A. Manica, J. Mallet, and C. D. Jiggins, 2013 Genome-wide evidence for speciation with gene flow in *Heliconius* butterflies. *Genome Research* .
- Martin, S. H., A. Eriksson, K. M. Kozak, A. Manica, and C. D. Jiggins, 2015 Speciation in *heliconius* butterflies: Minimal contact followed by millions of generations of hybridisation. *bioRxiv* .
- McVean, G. A. and N. J. Cardin, 2005 Approximating the coalescent with recombination. *Philos Trans R Soc Lond B Biol Sci* **360**: 1387–1393.
- Nee, S., E. C. Holmes, A. Rambaut, and P. H. Harvey, 1995 Inferring population history from molecular phylogenies. *Philosophical Transactions of the Royal Society of London Series B* **349**.
- Pybus, O. G., A. Rambaut, E. C. Holmes, and P. H. Harvey, 2002 New inferences from tree shape: numbers of missing taxa and population growth rates. *Systematic Biology* **51**: 881–888.
- Rasmussen, M. D., M. J. Hubisz, I. Gronau, and A. Siepel, 2014 Genome-wide inference of ancestral recombination graphs. *PLoS Genet* **10**: e1004342.
- Schiffels, S. and R. Durbin, 2014 Inferring human population size and separation history from multiple genome sequences. *Nature Genetics* **46**: 919 – 925.
- Terhorst, J. and Y. S. Song, 2015 Fundamental limits on the accuracy of demographic inference based on the sample frequency spectrum. *Proceedings of the National Academy of Sciences* **112**: 7677–7682.
- Wakeley, J., 1998 Segregating sites in Wright's island model. *Theoretical Population Biology* **53**: 166–174.
- Wakeley, J., 2009 *Coalescent theory*. Roberts & Company Publishers, Greenwood Village, Colorado.
- Wang, Y. and J. Hey, 2010 Estimating divergence parameters with small samples from a large number of loci. *Genetics* **184**: 363–373.
- Wilkinson-Herbots, H., 2012 The distribution of the coalescence time and the number of pairwise nucleotide differences in a model of population divergence or speciation with an initial period of gene flow. *Theoretical Population Biology* **82**: 92–108.
- Wilkinson-Herbots, H. M., 2008 The distribution of the coalescence time and the number of pairwise nucleotide differences in the "isolation with migration" model. *Theoretical Population Biology* **73**: 277–288.
- Yang, Z., 2002 Likelihood and Bayes estimation of ancestral population sizes in hominoids using data from multiple loci. *Genetics* **162**: 1811–1823.
- Zhu, T. and Z. Yang, 2012 Maximum likelihood implementation of an isolation-with-migration model with three species for testing speciation with gene flow. *Molecular Biology and Evolution* **49**: 832–842.