# Maximum Likelihood Implementation of an Isolation-with-Migration Model for Three Species

DANIEL DALQUEN[1*], TIANQI ZHU[2*] and ZIHENG YANG[1, 2, 3]

[1] *Department of Genetics, Evolution and Environment, University College London, Darwin Building, Gower Street, London WC1E 6BT, UK*

[2] *Center for Computational Genomics, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, China*

**Running head**: IMPLEMENTATION OF IM MODEL FOR 3 SPECIES

* Those authors contributed equally to the study.

**Key words:** Multispecies coalescent, maximum likelihood, speciation, IM model, migration.

[3]*Correspondence to*:    Ziheng Yang
            *Address*:    Department of Genetics, Evolution and Environment
                    University College London
                    Darwin Building
                    Gower Street
                    London WC1E 6BT
                    England
            *Email*:    z.yang@ucl.ac.uk
            *Phone:*    +44 (20) 7679 4379

*Abstract.*— We develop a maximum likelihood (ML) method for estimating migration rates between species using genomic sequence data. A species tree is used to accommodate the phylogenetic relationships among three species, allowing for migration between the two sister species, while the third species is used as an outgroup. A Markov chain characterization of the genealogical process of coalescence and migration is used to integrate out the migration histories at each locus analytically, while Gaussian quadrature is used to integrate over the coalescent times on each genealogical tree numerically. This is an extension of our early implementation of the symmetrical isolation-with-migration model for three species to accommodate arbitrary loci with two or three sequences per locus and to allow asymmetrical migration rates. Our implementation can accommodate tens of thousands of loci, making it feasible to analyze genome-scale datasets to test for gene flow. We calculate the posterior probabilities of gene trees at individual loci to identify genomic regions that are likely to have been transferred between species due to gene flow. We conduct a simulation study to examine the statistical properties of the likelihood ratio test for gene flow between the two ingroup species and of the maximum likelihood estimates of model parameters such as the migration rate. Inclusion of data from a third outgroup species is found to increase dramatically the power of the test and the precision of parameter estimation. We compiled and analyzed several genomic datasets from the *Drosophila* fruit flies. Our analyses suggest no migration from *D. melanogaster* to *D. simulans*, and a significant amount of gene flow from *D. simulans* to *D. melanogaster*, at the rate of ~0.02 migrant individuals per generation. We discuss the utility of the multispecies coalescent model for species tree estimation, accounting for incomplete lineage sorting and migration.

Migration or gene flow is an important biological process that affects our interpretation of genetic data from both within and between species (e.g., Patterson et al., 2006; Innan and Watanabe, 2006; Yamamichi et al., 2012; Leaché et al., 2013; Mallet et al., 2016). For example, different models of speciation make different predictions about the presence or absence of gene flow at the time of species formation. There is a rich body of literature in population genetics concerning models of population subdivision and migration, starting from Wright (1931; 1943). For example, in the finite-island model, any population can exchange migrants with any other (Wright, 1943), while in the stepping-stone model, only neighboring populations can exchange migrants (Kimura and Weiss, 1964). The standard single-population coalescent theory (Kingman, 1982) has been extended to deal with such models of population structure and migration, in the so-called *structured coalescent* (e.g., Li, 1976; Strobeck, 1987; Takahata, 1988; Notohara, 1990; Nath and Griffiths, 1993; Wilkinson-Herbots, 1998). Models of population structure have been implemented in computer programs such as GENETREE (Bahlo and Griffiths, 2000) and MIGRATE (Beerli and Felsenstein, 1999; 2001; Beerli, 2006), which allow joint estimation of population sizes and migration rates from genetic data.

However, population structure models ignore the phylogenetic relationships among the populations and their divergence times. The isolation-with-migration (IM) model is attractive as it

73 incorporates the population/species phylogeny in a model of migration. They allow us to estimate the
74 migration rates and other parameters such as the species divergence times and population sizes under
75 more realistic models (Nielsen and Wakeley, 2001; Hey and Nielsen, 2004; Wilkinson-Herbots, 2008;
76 2012). Another yet unexplored use of the IM model is species tree estimation under the multispecies
77 coalescent model with migration, accounting for both incomplete lineage sorting and introgression.
78 Coalescent-based phylogenetic inference, which accommodate gene tree-species tree discordance due
79 to incomplete lineage sorting, has been heralded as a paradigm shift in molecular phylogenetics
80 (Edwards, 2009). Recent analyses of genomic datasets have found widespread conflicts among
81 nuclear gene trees and between the mitochondrial gene tree and the nuclear species tree, for example,
82 in mosquitos (Fontaine et al., 2015), butterflies (Martin et al., 2013), frogs (Zhou et al., 2012), birds
83 (Ellegren et al., 2012), hares (Melo-Ferreira et al., 2012), bears (Liu et al., 2014; Kutschera et al.,
84 2014), and gibbons (Chan et al., 2013). Hybridization both between sister species and between non-
85 sister species is commonly observed between modern species, so it is natural to expect it to have
86 occurred in ancestral species as well, especially during adaptive radiations (Mallet, 2005; Mallet et al.,
87 2016). Many empirical studies have highlighted incomplete lineage sorting (or rapid radiation) and
88 gene flow (introgression) as the two major challenges to species tree estimation when the species are
89 closely related. While the multispecies coalescent model with gene flow should accommodate both
90 factors naturally, full likelihood methods of species tree estimation under the model are currently
91 lacking.

92 Full likelihood implementation of the IM model for the analysis of genetic sequence data is
93 challenging because calculation of the likelihood function has to average over the genealogical history
94 at every locus, which includes the gene tree topology, the branch lengths (the coalescent times), and
95 the whole migration trajectory (the number, directions and times of all migration events). The IM
96 programs (Nielsen and Wakeley, 2001; Hey and Nielsen, 2004; Hey, 2010), for example, are not
97 practical for analyzing datasets with a few hundred loci (Hey, 2010). Approximations are often
98 necessary to analyze genome-scale data with many loci (Gronau et al., 2011).

99 When there are only a few sequences at a locus, it is possible to integrate out the migration history
100 either numerically or analytically (Wang and Hey, 2010; Lohse et al., 2011; Zhu and Yang, 2012;
101 Andersen et al., 2014). It is then feasible to analyze tens of thousands of loci even though only a few
102 sequences are sampled at each locus. Here loci may be defined as loosely linked short genomic
103 segments that are far apart from each other, so that recombination within a locus is unlikely to affect
104 the gene tree distribution, while different loci are nearly independent due to recombination events
105 (Burgess and Yang, 2008; Lohse et al., 2011). Wang and Hey (2010) used numerical integration and
106 special functions to integrate out the migration history under the IM model for two species when the
107 data at every locus consist of two sequences, with one from each species. A more efficient approach
108 is to integrate out the migration trajectory analytically by using the Markov chain characterization of
109 the coalescent process with migration developed in the structured coalescent framework (Notohara,

110   1990; Nath and Griffiths, 1993; Hobolth et al., 2011; Zhu and Yang, 2012; Andersen et al., 2014).

111   For example, with only two sequences at a locus, the probability of the sequence data at any locus

112   depends on the sequence divergence time $t$ only, and not on the number and times of the migration

113   events. The density for $t$ can be calculated analytically (Hobolth et al., 2011; see also Nath and

114   Griffiths, 1993; Wilkinson-Herbots, 2008). Lohse *et al.* (2011) derived probabilistic distributions of

115   gene trees using generating functions and symbolic algebra in Mathematica. The implementation

116   allows more than two sequences at each locus, thus increasing the power of the analysis (Lohse et al.,

117   2011).

118       Zhu and Yang (2012) implemented the IM model for three species, assuming symmetry in the

119   migration rates and population sizes between species 1 and 2 (with $M_{12} = M_{21} = M$, and $\theta_1 = \theta_2$), while

120   a third species (species 3) is used as the outgroup. They constructed a likelihood ratio test (LRT) by

121   comparing this model, M2 (gene flow), with a null model of no migration with $M = 0$ (M0: no gene

122   flow). In their implementation, the data at every locus are assumed to consist of three sequences, with

123   one sequence from each species (this data configuration is referred to in this paper as '123'). This

124   restriction on data leads to reduced power of the test and to an unusual case of unidentifiability (Zhu

125   and Yang, 2012). Recently, Andersen *et al.* (2014) have considered the IM model in a general setting,

126   in which one ancestral species splits into an arbitrary number of populations at a time in the past (so

127   that the populations are related by a star phylogeny), allowing for migration between any two

128   populations. The authors developed a strategy for 'lumping' states in the Markov chain to alleviate

129   the problem of state-space explosion. Their implementation, for the case of two diploid individuals

130   from two species (four sequences per locus), assumed free recombination between any two sites

131   (alignment columns). Under this assumption, the data at different sites are independent (conditional

132   on the species phylogeny and parameters in the model) so that the sequence dataset can be

133   summarized as counts of $4^4$ possible site patterns (nucleotide combinations), and the authors were able

134   to integrate out the coalescent times in the gene trees for each site analytically (Andersen et al., 2014,

135   sections 5 and 8.4).

136       In this study we extend the implementation of Zhu and Yang (2012). Like many previous studies

137   such as Takahata et al. (1995), Wang and Hey (2010), and Lohse et al. (2011), we work under the

138   assumption of complete linkage within a locus and free recombination between loci. We note that

139   both free recombination and complete linkage within a locus are extreme assumptions, and their

140   impact on the inference is not yet well understood (but see Burgess and Yang, 2008; Zhu and Yang,

141   2012). We accommodate loci of two or three sequences of arbitrary configurations, including '11'

142   (two sequences from species 1), '112' (two sequences from species 1 and one sequence from species

143   2), and so on. Extension to arbitrary loci (with two or three sequences per locus) improves the power

144   of the likelihood ratio test of gene flow and makes it possible to estimate the migration rates, which

145   are unidentifiable with '123' loci alone (Zhu and Yang, 2012). We focus on migration between

146   species 1 and 2, and include species 3 as an outgroup to improve the power of the analysis. As nicely

147     discussed by Lohse *et al.* (2011), the outgroup may be informative about the gene tree topology as

148     well as the branch lengths and about the ancestral nucleotide states in the common ancestor of species

149     1 and 2.  Inclusion of the outgroup may also make the inference more robust to mutation rate variation

150     among loci (Yang, 2002).  We remove the symmetry assumption of the model, so that the inference

151     can be conducted under a more realistic model.  We develop an empirical Bayes approach to

152     calculating the posterior probabilities of gene tree topologies at individual loci, which may be

153     informative about whether the locus has been transferred between species due to gene flow.  We

154     conduct a simulation study to examine the false positive rate and power of the LRT of gene flow as

155     well as the bias and variance of maximum likelihood estimates of model parameters.  We use the

156     genome sequences of *Drosophila melanogaster, D. simulans*, and *D. yakuba* to construct multi-locus

157     datasets and apply our new method to infer the pattern and rate of migration between those fruit-fly

158     species.

159

160     ## THEORY AND METHODS

161     ### *Model and Data*

162     The terms species and population are used interchangeably in this paper.  The species tree is $((1, 2),$

163     $3)$, with 4 and 5 to be the ancestral species (Fig. 1a).  The two divergence events on the species tree

164     define three time epochs: $E_1$: $(0, \tau_1)$, $E_2$: $(\tau_1, \tau_0)$ and $E_3$: $(\tau_0, \infty)$ (Fig. 1a).  We consider two models.

165     M0 (no gene flow) assumes no gene flow and is the multispecies coalescent model for three species

166     (Takahata et al., 1995; Yang, 2002; Rannala and Yang, 2003).  Model M2 (gene flow) allows

167     migration between species 1 and 2 (during time epoch $E_1$), but not from or to species 3.

168          There are nine parameters in the general IM model for three species, including two species

169     divergence times ($\tau_0$ and $\tau_1$), five effective population sizes ($\theta_1$, $\theta_2$, $\theta_3$, $\theta_4$, $\theta_5$), and two migration rates

170     ($M_{12}$ and $M_{21}$).  Here $\tau_0$ and $\tau_1$ are scaled by the mutation rate and are measured by the expected

171     number of mutations per site, and $\theta_i = 4N_i\mu$ $(i = 1, …, 5)$ are the population size parameters for the

172     five species, with $N_i$ being the (effective) population size of species $i$ and $\mu$ the mutation rate per site

173     per generation.  The migration rate is $M_{ij} = N_j m_{ij}$, where $m_{ij}$ is the proportion of individuals in

174     population $j$ that are immigrants from population $i$.  We define parameters by referring to the real-

175     world process with time running forward (rather than the coalescent view with time running

176     backward) so that $M_{ij}$ is the expected number of migrant individuals from populations $i$ to $j$ per

177     generation.  The parameters under M2 (gene flow) are $\Theta_2 = \{\tau_0, \tau_1, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5, M_{12}, M_{21}\}$.

178     Model 0 (no gene flow) is a special case of M2 with $M_{12} = M_{21} = 0$, with parameters $\Theta_0 = \{\tau_0, \tau_1, \theta_1,$

179     $\theta_2, \theta_3, \theta_4, \theta_5\}$.  Note that the symmetrical versions of M0 and M2 assume $\theta_1 = \theta_2$ and $M_{12} = M_{21}$ (Zhu

180     and Yang, 2012).

181          The data consist of multiple neutral loci.  At each locus, two or three sequences are sampled, each

182    from any of the three species. We focus mainly on the case of three sequences at a locus. The case of

183    two sequences is much simpler and will be described briefly. Let the three sequences at a locus be $a$,

184    $b$, and $c$. Each sequence will also be labelled by the population it is sampled from. For example, the

185    initial state for a locus with data configuration '123' (with one sequence from each of the three

186    species) is recorded as $1_a 2_b 3_c$. The Markov chain runs backwards in time, describing the change of

187    states due to coalescent and migration. For example a locus with initial state $1_a 2_b 3_c$ may enter the

188    state $2_{ab} 3_c$, which means that sequences $a$ and $b$ have coalesced so that only two sequences remain in

189    the sample and the ancestor of sequences $a$ and $b$ is in population 2 while sequence $c$ is in population

190    3. There are six gene tree shapes for three sequences: $G_1$-$G_6$ (Fig. 1b-g), depending on the time

191    epochs during which the two coalescent events occur. When we keep track of both the sequence IDs

192    ($a$, $b$, $c$) and the population IDs (1, 2, 3), each gene tree shape may correspond to three distinct gene

193    trees (Fig. 2). For example, tree shape $G_6$ corresponds to three gene trees: $G_{6c}$: $((a, b), c)$; $G_{6a}$: $((b, c),$

194    $a)$; and $G_{6b}$: $((c, a), b)$, where the subscript is the more distantly related sequence in the gene tree.

195    However, depending on the initial data configuration, some of the gene trees may not be possible (for

196    example, for a '123' locus, only gene trees $G_{3c}$, $G_{5c}$, $G_{6c}$, $G_{6a}$, $G_{6b}$ are possible under M2), and

197    furthermore some of the gene trees have the same probability distribution under the model (such as

198    $G_{6c}$, $G_{6a}$, and $G_{6b}$). To avoid excessive notation we make a distinction between gene tree shapes and

199    gene trees only if there is a risk of confusion.

### *Likelihood Function for Three Sequences at a Locus*

201    We assume that the sequences at each locus are already aligned, with alignment gaps and ambiguity

202    nucleotides removed. We use the JC69 mutation model (Jukes and Cantor, 1969) to correct for

203    multiple substitutions. The different loci are assumed to have the same mutation rate, although

204    relative rates for the loci can be incorporated in the likelihood calculation (if available, for example,

205    through comparison with an outgroup species, Yang, 2002). The sequence alignment at any locus $i$

206    with three sequences can be summarized as the counts, $D_i = (n_0, n_1, n_2, n_3, n_4)$, of sites with five

207    different site patterns: *xxx*, *xxy*, *yxx*, *xyx*, and *xyz*, where $x$, $y$ and $z$ are any distinct nucleotides. The

208    probability of the data given the gene tree topology ($G$) and branch lengths ($b_0$, $b_1$) (Fig. 2), $P(D_i|G,$

209    $b_0, b_1)$, is thus given by the multinomial distribution, with the probabilities of the five site patterns

210    calculated efficiently under the JC69 model (Saitou, 1988; Yang, 1994). Conveniently, $P(D_i|G, b_0, b_1)$

211    depends on the gene tree topology and branch lengths, but not on which time epoch each coalescent

212    event occurs in (Yang, 2002; 2010).

213        The probability of data at locus $i$ is an average over the gene tree topologies and coalescent times

$$f(D_i \mid \Theta) = \sum_k \int_{l_0}^{u_0} \int_{l_1}^{u_1} P(D_i \mid G_k, b_0, b_1) f(G_k, t_0, t_1 \mid \Theta)\, \mathrm{d}t_1 \mathrm{d}t_0 , \tag{1}$$

215    where the sum is over all possible gene trees for the locus, while the integrals are over the coalescent

216    times $t_0$ and $t_1$, with the integral limits $t_0 \in (l_0, u_0)$ and $t_1 \in (l_1, u_1)$ given below. Note the branch

217  lengths $b_0$ and $b_1$ in the gene tree are simple linear functions of $t_0$ and $t_1$ (Figs. 1 and 2 and Table 1).

218  The probability of the genealogy, $f(G_k, t_0, t_1|\Theta)$, depends on the model ($M_0$ or $M_2$) and will be

219  described in the next section.  For data configurations with three sequences, there are up to $6 \times 3 = 18$

220  gene trees to average over.

221      Finally, the log likelihood of the data at all $L$ loci, $D = \{D_i\}$, is a sum over the $L$ loci

222
$$\ell(\Theta; D) = \sum_{i=1}^{L} \log f(D_i \mid \Theta). \tag{2}$$

223      Note that our model assumes that the $n$ sites in the sequence at the locus share the same

224  genealogical tree (topology and coalescent times).  This contrasts with the implementation of

225  Andersen *et al.* (2014), which assumes that the different sites have independent histories.

## *Implementation of Model $M_0$ (No Gene Flow)*

227  We first discuss our ML implementation of model M0, which assumes no migration between any two

228  populations.  The implementation of Yang (2002) considered '123' loci only so that the model

229  involve only four parameters: $\Theta_0 = \{\tau_0, \tau_1, \theta_4, \theta_5\}$.  Here we allow arbitrary loci of two or three

230  sequences, with up to seven parameters in the model: $\Theta_0 = \{\tau_0, \tau_1, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5\}$.  Note that the

231  population size parameter for a modern species ($\theta_1$, $\theta_2$, or $\theta_3$) exists in the model only if two or more

232  sequences are sampled from that species at least at one locus.

233      Consider a locus with three sequences.  In general, the probability density of the gene tree has the

234  form

235
$$f(G_k, t_0, t_1) = \text{rates} \times e^{-T} = \frac{2}{\theta_i} \frac{2}{\theta_j} e^{-T}, \tag{3}$$

236  where parameters $\theta_i$ and $\theta_j$ are for the populations in which the two coalescent events occur and the

237  exponential term $e^{-T}$ is the probability that no coalescent event occurs in the rest of the gene tree, with

238  $T$ being the *total per-lineage-pair coalescent waiting time* of Yang (2014, p.336).  Note that the

239  coalescent rate for a pair of sequences in a population with population size parameter $\theta$ is $2/\theta$: for

240  very small $\Delta t$, the probability that the pair will coalesce during the time interval $(t, t + \Delta t)$ is $\frac{2}{\theta}\Delta t$ .

241      Take, for example, configuration '111', with the initial state $1_a 1_b 1_c$.  The probability of data for

242  the locus (Eq. 1) is an average over $6 \times 3$ gene trees.  For example, in the case of gene tree $G_{1c}$: $((a, b),$

243  $c)$, the probability density of the gene tree (with coalescent times) is

244
$$f(G_{1c}, t_0, t_1) = \frac{2}{\theta_1} \frac{2}{\theta_1} e^{-T} = \frac{2}{\theta_1} \frac{2}{\theta_1} e^{-\frac{6}{\theta_1} t_1 - \frac{2}{\theta_1} t_0}, \quad t_0 > 0, \ t_1 > 0, \ t_0 + t_1 < \tau_1, \tag{4}$$

245  where $\frac{2}{\theta_1}$ and $\frac{2}{\theta_1}$ are the rates for the two coalescent events, both occurring in species 1.  Because of

246  the symmetry of the '111' locus, the density is the same for the three gene trees: $G_{1c}$, $G_{1a}$, and $G_{1b}$.

247  The densities and rates for all data configurations and gene trees are summarized in Table S1.  Note

248  that some gene trees are not possible for certain configurations of loci (e.g., gene trees $G_{1c}$, $G_{1a}$, and

249  $G_{1b}$ for '112' loci).

7

250    To compute the integrals of equation (1) numerically, we apply a linear transform. Let $x_0 = \frac{2}{\theta_i} t_0$

251    and $x_1 = \frac{2}{\theta_j} t_1$ be the coalescent times measured in generations, where $\theta$s are for the populations in

252    which the coalescent events occur. Each integral in equation (1) then becomes

253    $$\int_{l_0}^{u_0} \int_{l_1}^{u_1} P(D_i \mid G_k, b_0, b_1) f(G_k, t_0, t_1) \, dt_1 dt_0 = \int_{l_0'}^{u_0'} \int_{l_1'}^{u_1'} P(D_i \mid G_k, b_0, b_1) f(G_k, x_0, x_1) \left| \frac{\partial(t_0, t_1)}{\partial(x_0, x_1)} \right| dx_1 dx_0 . \qquad (5)$$

254    In several cases (gene tree shapes $G_1$ and $G_4$ for initial state '111'; $G_4$ for '112'; and $G_1$, $G_2$ and $G_4$

255    for '333'), the integration region is a triangle (for instance, the region for $G_1$ is given by $t_0 > 0$, $t_1 > 0$,

256    $t_0 + t_1 < \tau_1$; see Fig. 1). As we calculate the 2-D integral of equation (5) by calculating two 1-D

257    integrals using Gaussian quadrature (the so-called product rule), the integral region has to be a

258    rectangle. We thus apply a transform to achieve this. For example, in the case of $G_1$ for the initial

259    state '111', we use $x_0 = \frac{2}{\theta_1}(t_0 + t_1)$, $x_1 = \frac{t_1}{t_0 + t_1}$, so that $t_0 = \frac{\theta_1}{2} x_0 (1 - x_1)$, $t_1 = \frac{\theta_1}{2} x_0 x_1$. The new limits

260    are $0 < x_0 < \frac{2}{\theta_1} \tau_1$, $0 < x_1 < 1$, and the Jacobi of the transform is $\left| \frac{\partial(t_0, t_1)}{\partial(x_0, x_1)} \right| = \frac{\theta_1}{2} \frac{\theta_1}{2} x_0$. Then

261    $$\int_0^{\tau_1} \int_0^{\tau_1 - t_0} P(D_i \mid G_{1k}, b_0, b_1) \times \frac{2}{\theta_1} \frac{2}{\theta_1} e^{-\frac{6}{\theta_1} t_1 - \frac{2}{\theta_1} t_0} \, dt_1 dt_0 = \int_0^{\frac{2}{\theta_1} \tau_1} \int_0^1 P(D_i \mid G_{1k}, b_0, b_1) \times x_0 e^{-2x_0 x_1 - x_0} \, dx_1 dx_0 , \qquad (6)$$

262    where $b_0 = t_0$ and $b_1 = t_1$ in the integral on the left-hand side, and $b_0 = \frac{\theta_1}{2} x_0 (1 - x_1)$ and $b_1 = \frac{\theta_1}{2} x_0 x_1$ in

263    the integral on the right-hand side.

264    The transforms from $(t_0, t_1)$ to $(x_0, x_1)$ are summarized in Table S2. We use Gaussian quadrature

265    to calculate the 2-D integrals of equations (5) or (6). Except where stated otherwise, we used $K = 16$

266    points in the quadrature. See Yang (2010) for details. It is necessary to apply scaling to avoid

267    underflows as the probabilities of equation (1) may be too small to represent in the computer.

268    *The case of two sequences.* In the case of two sequences at a locus, the possible initial states are

269    11, 12, 22, 13, 23, and 33, depending on which populations the two sequences are sampled from. The

270    simple gene tree has two branches, which have the same length $t$, with density $f(t|\Theta)$ (Table 1). For

271    instance, with the initial state 11 (two sequences from species 1), $f(t|\Theta)$ is a piecewise continuous

272    function because the population size and thus the coalescent rate may differ in the three time epochs.

273    The sequence data at the locus are summarized as $d_i$ differences out of $n_i$ sites. Then the probability

274    of observing $d_i$ differences at $n_i$ sites given that the two sequences separated time $t$ ago is

275    $$f(d_i|t) = \left( \frac{3}{4} - \frac{3}{4} e^{-8t/3} \right)^{d_i} \left( \frac{1}{4} + \frac{3}{4} e^{-8t/3} \right)^{n_i - d_i} . \qquad (7)$$

276    The (unconditional) probability of observing the data at the locus is an average over the coalescent

277    time

278    $$f(d_i|\Theta) = \int_0^\infty f(t \mid \Theta) f(d_i \mid t) \, dt . \qquad (8)$$

279    Gaussian quadrature is used to calculate the 1-D integral, with the transform $x = \frac{2}{\theta_j} t$ (Table 1).

*Implementation of Model M2 (gene flow)*

Under model M2 (gene flow), the likelihood is given by equation (1) as before, and the probability of the data at each locus $P(D_i|G_k, b_0, b_1)$ remains the same. However, the probability density for the gene trees, $f(G_k, t_0, t_1)$, depends on the migration rates and differs from that under model M0. Our aim in this section is thus to describe the gene-tree density. We use a Markov chain to characterize the process of coalescent and migration when we trace the gene genealogy backwards in time. In the general case, the states of the Markov chain will include both the population IDs and sequence IDs. Because of our assumption of no migration involving species 3, the coalescent process during time epochs $E_2$ and $E_3$ are essentially the standard single-population coalescent. Thus, we focus on epoch $E_1$. While it is possible to use one Markov chain for all initial states, we use different Markov chains depending on the initial states to increase computational efficiency (Table 2). The Markov chain characterization allows one to calculate the probability density for the gene tree topology and coalescent times, $f(G_k, t_0, t_1)$, with the migration history integrated out analytically (Hobolth et al., 2011; Zhu and Yang, 2012; Andersen et al., 2014). We do not use the idea of Andersen et al. (2014) for lumping states in the Markov chain because it would add much complexity to the algorithm with no or little gain for the cases of two or three sequences per locus. For the general migration case with three species, lumping actually increases the number of states from 12 to 15 for two sequences, and from 57 to 70 for three sequences (Andersen et al., 2014, table 2). We note that for four or more sequences per locus, Andersen et al.'s algorithm may lead to considerable reduction of the state space.

We illustrate the theory using gene tree $G_{1c}$: $((a, b), c)$ and initial state $s =$ '111'. We take advantage of the symmetry of the initial state and consider a reduced Markov chain with eight states, dropping the sequence IDs: {111, 112, 122, 222, 11, 12, 22, 1|2} (Table 2). Here the state '1|2' means one sequence in either population 1 or 2. When both coalescent events have occurred and there is only one sequence in the sample, there will be no need to keep track of the population ID, so that states 1 and 2 can be lumped into one artificial absorbing state (Andersen et al., 2014). The rate matrix is given in Table 3. For gene tree shape $G_1$, we have $f(G_{1c}, t_0, t_1) = f(G_{1a}, t_0, t_1) = f(G_{1b}, t_0, t_1) = \frac{1}{3} f(G_1, t_0, t_1)$, with

$$f(G_1, t_0, t_1) = 3 \frac{2}{\theta_1} P_{s,111}(t_1) \left( \frac{2}{\theta_1} P_{11,11}(t_0) + \frac{2}{\theta_2} P_{11,22}(t_0) \right) + \frac{2}{\theta_1} P_{s,112}(t_1) \left( \frac{2}{\theta_1} P_{12,11}(t_0) + \frac{2}{\theta_2} P_{12,22}(t_0) \right)$$
$$+ \frac{2}{\theta_2} P_{s,122}(t_1) \left( \frac{2}{\theta_1} P_{12,11}(t_0) + \frac{2}{\theta_2} P_{12,22}(t_0) \right) + 3 \frac{2}{\theta_2} P_{s,222}(t_1) \left( \frac{2}{\theta_1} P_{22,11}(t_0) + \frac{2}{\theta_2} P_{22,22}(t_0) \right). \tag{9}$$

Note that the probability density function here has the interpretation that $f(G_1, t_0, t_1) \Delta t_0 \Delta t_1$, for very small $\Delta t_0$ and $\Delta t_1$, is the probability that the gene tree topology is $G_1$ (that is, $t_0 + t_1 < \tau_1$), that the first coalescent occurs during the time interval $(t_1, t_1 + \Delta t_1)$, and that the second coalescent occurs during the time interval $(t_1 + t_0, t_1 + t_0 + \Delta t_0)$ (see Fig. 1). Equation (9) gives this probability as the sum of four terms. The first term is for the case where the Markov chain is in state 111 right before $t_1$, with probability $P_{s,111}(t_1)$; the first coalescent occurs in species 1 during $(t_1, t_1 + \Delta t_1)$, with probability

314     $3 \times \frac{2}{\theta_1} \Delta t_1$, the factor 3 due to there being 3 possible pairs for coalescent with the state 111; and then the

315     second coalescent occurs during $(t_1 + t_0, t_1 + t_0 + \Delta t_0)$ either in population 1, with probability

316     $P_{11,11}(t_0) \times \left(\frac{2}{\theta_1} \Delta t_0\right)$, or in population 2, with probability $P_{11,22}(t_0) \times \left(\frac{2}{\theta_2} \Delta t_0\right)$. Note that in this scenario,

317     the first coalescence changes the state of the chain from 111 to 11. Similarly the $2^{nd}$, $3^{rd}$, and $4^{th}$ terms

318     in equation (9) are for the cases where the state right before the first coalescent at time $t_1$ is 112, 122,

319     and 222, respectively, with the second coalescent occurring either in population 1 or in population 2.

320       The densities for the other gene trees and for the other initial states are presented in Appendix A

321     and summarized in Tables S3 and S4.

322       This Markov chain characterization of the genealogical process of coalescent and migration also

323     allows easy calculation of the probabilities of gene tree topologies, integrating over the coalescent

324     times. For example with the initial state '123', the transition probability $P_{123,\,13|23}(\tau_1)$ calculated from

325     the Markov chain of Table 2 (case III) is the probability that sequences 1 and 2 have coalesced by

326     time $\tau_1$. This then gives the probabilities for the five gene trees for the initial state '123' as $P(G_{3c}) =$

327     $P_{123,\,13|23}(\tau_1)$, $P(G_{6c}) = P(G_{6a}) = P(G_{6b}) = \frac{1}{3}\left(1 - P_{123,13|23}(\tau_1)\right) \times e^{-2/\theta_5(\tau_0-\tau_1)}$, and $P(G_{5c}) = 1 - P(G_{3c}) -$

328     $3P(G_{6c})$ (Fig. 1). Here $e^{-2/\theta_5(\tau_0-\tau_1)}$ is the probability that sequences 1 and 2 do not coalesce in epoch

329     $E_2$.

330       In the case of two sequences at a locus, the likelihood calculation given the branch length $t$ is

331     given by equations (7) and (8). The probability density of the genealogy $f(t)$ under M2 (gene flow) is

332     the same as under $M_0$ for the initial states 13, 23, or 33 (Table 1). For initial states $s = 11$, 12, or 22,

333     the two sequences can coalesce in any of the three time intervals: $(0, \tau_1)$, $(\tau_1, \tau_0)$, and $(\tau_0, \infty)$, so that

334     the density is given as

335
$$
f(t) = \begin{cases}
\frac{2}{\theta_1} P_{s,11}(t) + \frac{2}{\theta_2} P_{s,22}(t), & t < \tau_1, \\[2ex]
\sum_{j \in B_2} P_{s,j}(\tau_1) \times \frac{2}{\theta_5} e^{-\frac{2}{\theta_5}(t-\tau_1)}, & \tau_1 < t < \tau_0, \\[2ex]
\sum_{j \in B_2} P_{s,j}(\tau_1) e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)} \times \frac{2}{\theta_4} e^{-\frac{2}{\theta_4}(t-\tau_0)}, & t > \tau_0.
\end{cases}
\tag{10}
$$

336     where $B_2 = \{11, 12, 22\}$ is the set of states with two sequences. The transition probability $P_{s,j}(t)$ is

337     calculated using a Markov chain with four states 11, 12, 22, and 1|2. See Hobolth *et al.* (2011).

338

339     *Likelihood Ratio Test Comparing Models M0 (No Gene Flow) and M2 (gene flow)*

340     As M0 is a special case of M2, we use an LRT to compare them. However, we note that the large-

341     sample $\chi^2$ approximation is not valid and the null distribution (that is, the distribution of the test

342     statistic $2\Delta\ell = 2[\ell_2 - \ell_0]$ when the null hypothesis M0 is true) depends on the data configurations at

343     the loci.

344       As discussed by Zhu and Yang (2012), if the data consist of loci of configuration 123 only, the

345     symmetric version of model M2 has two more parameters than M0: $\theta_1$ (=$\theta_2$) and $M$. However, for

346     two reasons, the large-sample $\chi_2^2$ approximation to the test statistic is not valid. First, the null

347     hypothesis M0 corresponds to the alternative hypothesis M2 with $M = 0$, but this parameter value is at

348     the boundary of the parameter space. Second, when $M = 0$, parameter $\theta_1$ (=$\theta_2$) in model M2 becomes

349     unidentifiable. As a result of the violations of the regularity conditions for the $\chi^2$ approximation, the

350     true null distribution is unknown. Furthermore, analysis of data of configuration '123' under M2

351     leads to an unusual unidentifiability problem: two sets of $\theta_1$ (=$\theta_2$) and $M$ values always give the same

352     log likelihood value.

353       It is easy to see that this unidentifiability problem exists for the symmetric model if the data

354     consist of a mixture of loci with configurations 12 and 123, or if the 12 and 123 loci are supplemented

355     with an arbitrary mixture of loci of configurations 33, 13, 23, 333, 133, and 233, without any loci of

356     configurations 11, 22, 112, 122, 111, 222, 113, and 223. All such datasets will show the

357     unidentifiability problem under M2 and the two violations of the regularity conditions for the $\chi_2^2$

358     asymptotics. In this study, we follow Zhu and Yang (2012) and use $\chi_2^2$ as the null distribution to

359     conduct the test and consider the test to be significant if $2\Delta\ell > 5.99$. For data of a mixture of loci with

360     configurations 11, 22, and 12, or of a mixture of 113, 223, and 123, parameter $\theta_1$ (=$\theta_2$) is identifiable

361     in both models M0 and M2. While we still have the problem with the parameter value $M = 0$ at the

362     boundary, the problem is an instance of case 5 in Self and Liang (1987). As a result, the null

363     distribution is known to be the 50:50 mixture of 0 and $\chi_1^2$, with the 5% critical value to be 2.71. The

364     critical values for different mixtures of two initial states under the symmetric model are given in

365     supplementary Table S5.

366       A similar unidentifiability problem exists under the asymmetrical model for certain combinations

367     of loci. Let $U_1 = \{11, 111, 112, 113\}$ and $U_2 = \{22, 122, 222, 223\}$. If a dataset consists of at least

368     one of the states in $U_1$ and one of the states in $U_2$, then M2 is identifiable. In this case, M2 has two

369     more free parameters ($M_{12}$ and $M_{21}$) than M0 and a 50:50 mixture of 0 and $\chi_2^2$ is the null distribution,

370     with the significance value $2\Delta\ell = 4.61$. If a dataset consists of at least one state in $U_1$ but none in $U_2$

371     or at least one state in $U_2$ but none in $U_1$, the model is unidentifiable. In this case the null distribution

372     is unknown and we use $\chi_3^2$ to conduct the test, with critical value 7.82. If a dataset contains none of

373     the states in either $U_1$ or $U_2$, we use $\chi_4^2$ to conduct the test, with the critical value 9.49, since M0 and

374     M2 differ by four parameters. The critical values for the likelihood ratio test under the asymmetric

375     model for different mixtures of loci are given in Table S6.

### *Posterior probabilities of gene tree topologies*

When there is gene flow, it may be of interest to know which loci are most likely to have been transferred between species, and to further examine whether the transferred genes share a particular function or are located in the same chromosomal region. Our formulation of the IM model does not allow us to address this question in a straightforward manner. However, we can use an Empirical Bayes approach to calculate the posterior probabilities of the 18 gene tree topologies for each locus, which may be informative about whether the locus is involved in cross-species gene flow. For example, for a '123' locus, the possible gene trees are $G_{3c}$, $G_{5c}$, $G_{6c}$, $G_{6a}$, and $G_{6b}$, with $G_{3c}$ being possible only if the locus is transferred between species 1 and 2 (Fig. 1). Similarly for a '112' locus, gene tree shape $G_1$ is possible only with gene flow. We note that loci of certain configurations, such as '113' or '223', may not provide such information about gene flow.

The probability of data at a locus, $f(D_i|\Theta)$, is a sum over the 18 gene trees (equation 1). The posterior probabilities of the gene trees can be calculated by rescaling those 18 terms so that they sum to 1.

$$f(G_k|D_i,\Theta) = \frac{f(G_k|\Theta)f(D_i|G_k,\Theta)}{f(D_i|\Theta)} = \frac{\int_{l_0}^{u_0}\int_{l_1}^{u_1} P(D_i|G_k,b_0,b_1)f(G_k,t_0,t_1|\Theta)\,\mathrm{d}t_1\mathrm{d}t_0}{\sum_k \int_{l_0}^{u_0}\int_{l_1}^{u_1} P(D_i|G_k,b_0,b_1)f(G_k,t_0,t_1|\Theta)\,\mathrm{d}t_1\mathrm{d}t_0}. \quad (11)$$

We replace the parameters ($\Theta$) by their MLEs $(\hat{\Theta})$, and the method is known as Empirical Bayes (EB). The EB procedure does not account for sampling errors in the MLEs, which may be a concern if the dataset is small and the MLEs involve considerable sampling errors. This is the same EB procedure as used in reconstructing ancestral sequences in molecular phylogenetics (Yang et al., 1995) and in detecting positively selected sites in a protein-coding gene (Nielsen and Yang, 1998).

We conducted a small simulation to examine the reliability of the calculation using equation (11). We simulated datasets using the parameter values: $\tau_0 = 0.0243$, $\tau_1 = 0.0136$, $\theta_4 = 0.0400$, $\theta_5 = 0.0106$, $\theta_1 = 0.0052$, $\theta_2 = 0.0127$, $M_{12} = 0$ and $M_{21} = 0.0183$, which are the MLEs under M2 from the *Drosophila* dataset D1 (auto), to be described and analyzed later (Tables 4 and 9). We simulated two replicate datasets, each of the same size and configurations as the real data. The results are very similar between the two datasets so we discuss only those for the first dataset. The MLEs from the simulated dataset are $\hat{\tau}_0 = 0.0242$, $\hat{\tau}_1 = 0.0137$, $\hat{\theta}_4 = 0.0402$, $\hat{\theta}_5 = 0.0104$, $\hat{\theta}_1 = 0.0058$, $\hat{\theta}_2 = 0.0126$, $\hat{M}_{12} = 0.0018$ and $\hat{M}_{21} = 0.0196$, very close to the true values. The calculated posterior probabilities for gene tree topologies for the '123' loci (Fig. 3a) are accurate in the sense that a posterior probability of 90% is for a correct gene tree about 90% of the time (Fig. 3b). However, the power may not be very high. While the posterior for gene trees $G_{6a}$ and $G_{6b}$ may reach high values, that for $G_{6c}$ is seldom very high (Fig. 3c). It may be hard to distinguish among gene trees $G_{3c}$, $G_{5c}$, and $G_{6c}$. Lastly, approximately equal proportions of loci are inferred to have gene trees $G_{6c}$, $G_{6a}$ and

409 $G_{6b}$ (Fig. 3a), and they are also close to the expected proportions. Overall the results indicate a well-
410 behaved method.

## *Program Implementation, Validation, and Availability*

412 While the general theory of the gene-tree distribution under the Markov chain characterization of the
413 genealogical process under the IM model is straightforward (Zhu and Yang, 2012; Andersen et al.,
414 2014), development of a computer program that can analyze tens of thousands of loci with an
415 arbitrary mixture of loci of different configurations is challenging. Note that under both models M0
416 (no gene flow) and M2 (gene flow), the number of possible gene trees, the probability density of each
417 gene tree and its coalescent times, and the integration limits for the integrals over the coalescent times
418 all depend on the data configuration at the locus. This dependence makes the programming effort
419 rather tedious and error-prone. Thus we decided to tabulate the necessary results, in Tables S1 and S2
420 for M0 and similarly in Tables S3 and S4 for M2.

421    We conducted extensive tests to validate our implementation. The Mᴄᴄᴏᴀʟ program, which is
422 part of the ʙᴘᴘ package (Yang and Rannala, 2010; Zhang et al., 2011), was used to simulate sequence
423 data under models M0 and M2 for different data configurations and parameter values. We ensured
424 consistency of the MLEs: when the same model is used to generate the data and to analyze them, the
425 MLEs should converge to the true parameter values when the size of the dataset (the number of loci)
426 increases. We also confirmed that the likelihood stabilizes when the number of points in the Gaussian
427 quadrature is increased. We simulated $10^6$ (true) gene trees under M2 to confirm that the observed
428 frequencies of gene tree topologies match their probabilities calculated from the Markov chain
429 characterization.

430    Both models M0 and M2 are implemented in the program 3ꜱ. We identified two bottlenecks in
431 calculating the likelihood and improved performance in both areas. First, for most initial states, the
432 transition probability matrix $P(t)$ needs to be calculated numerically, involving an expensive matrix
433 exponential. We use the GNU Scientific Library (GSL) (Galassi et al., 2013) to optimize this step.
434 Second, the likelihood calculation is proportional to the number of loci in the data, as it is dominated
435 by the computation of the probability of data at each locus, $f(D_i|\Theta)$. We take advantage of the
436 independence among loci and use OpenMP to parallelize the computation (Dagum and Menon, 1998).
437 While both optimizations are optional, they offer significant speed-ups on genome-scale datasets (Fig.
438 S1). The program, with instructions on how to compile and run it with and without GSL and
439 OpenMP, is available at http://abacus.gene.ucl.ac.uk/software/3s.html.

## Drosophila *genomic datasets*

441 We compiled multi-locus datasets for three *Drosophila* species, *D. melanogaster* (M), *D. simulans* (S)
442 and *D. yakuba* (Y). We used Flybase FB2016_01 (Attrill et al., 2016) genome releases of *D.*
443 *melanogaster* (r6.09, January 2016), *D. simulans* (r2.01, Hu et al., 2013) and *D. yakuba* (r1.05,
444 January 2016), as well as the assembly of *D. simulans* strain M252 (Palmieri et al., 2014). We treated

445    the two *D. simulans* genomes (r2.01 from North American and M252 from Madagascar) as two

446    random samples from the same species.   Five datasets of MSSY loci were constructed (Table 4): D1

447    (auto) for autosomes 2 and 3, D2 (noncoding) for intergenic regions and introns from chromosomes 2

448    and 3, D3 (chrX) for the X chromosome, D4 (exons complete) and D5 (exons split).  D4 (exons

449    complete) was compiled using non-overlapping complete exons on chromosomes 2 and 3.  When

450    exons were overlapping, only the longest was kept.  For all datasets except D4 (exons complete),

451    sequences were split into chunks between 100 and 500bp that were separated by at least 2kb.  These

452    criteria were from Wang and Hey (2010), based on previous estimates of recombination rates for

453    *Drosophila* (Hey and Nielsen, 2004).  To construct each of datasets D1-D4, we extracted the loci from

454    the *D. melanogaster* genome as a starting point and then ran NCBI BLAST (Camacho et al., 2009)

455    with default settings to find matching sequences in the other genomes. We discarded short matches

456    (<40% of the query sequence), and removed loci where the two longest matches differed in length by

457    less than 10% to avoid paralogues. The remaining loci were aligned using MAFFT, using default

458    settings (Katoh and Standley, 2013).  We reduced each of the MSSY loci to either MSY or SSY by

459    randomly removing either the *D. melanogaster* or one of the *D. simulans* sequences.  Dataset D5

460    (exons split) was constructed by splitting the alignments of D4 (exons complete) into loci of between

461    100 and 500bp and removing chunks that did not fulfill the 2kb-separation criterion.  Thus all loci in

462    D5 are also in D4, but the alignments of the same loci in D5 may be shorter.  Some loci in D4 (374 of

463    them) were longer than 2600bp, and were split into more than one locus in D5.  Finally, we added the

464    378 MMY loci from Hutter et al. (2007) to all datasets except D2 (chrX) after updating their

465    coordinates to the current *D. melanogaster* release and confirming that they do not overlap with the

466    MSSY loci we compiled.

467        Note that D2 (noncoding) includes both intergenic regions and introns: these were found to

468    produce very similar estimates in a preliminary analysis and were thus merged into one dataset.  D1

469    (auto) and D3 (chrX) include both noncoding regions and exons.  The loci in D2 (noncoding), D4

470    (exons complete), and D5 (exons split) may not be included in D1 (auto).

471        The five datasets were analyzed using the program 3s under models M0 and M2 to estimate

472    parameters and to test for gene flow.  Fitting the two models to each dataset took about 20 minutes on

473    a single core and ~1 minute using 32 cores on a Sun Fire X4600M2 server (with 32 Opteron AMD

474    cores at 2.7GHz).  We also calculated the posterior probabilities of gene tree topologies under M2 to

475    identify the gene loci that are most likely to have been transferred across species barriers during

476    introgression (Eq. 11).

477    **RESULTS**

478    ***Computer Simulation to Examine the Statistical Properties of the new model***

479    We conducted computer simulations to examine the false positive rate and the power of the LRT

480    comparing models M0 (no gene flow) and M2 (gene flow) to test for migration between species 1 and

2. We also examined the biases and variances of MLEs of parameters under M2. Our simulation design largely follows that of Zhu and Yang (2012).

To examine the false positive rate of the test, we simulated replicate datasets under the symmetrical version of M0 and analyzed them under both M0 and M2, assuming symmetry (Table 5). We used four sets of parameter values (Zhu and Yang, 2012: table 1). The first two sets are based roughly on parameter estimates from the hominoids (Burgess and Yang, 2008) and the mangroves (Zhou et al., 2007). Sets 3 and 4 have larger parameter values and also different values for the three $\theta$s. The number of loci was fixed at $L = 10, 100, 1000$, and $15{,}000$, with each locus having 500 sites. Gene trees with branch lengths (coalescent times) were generated from the multispecies coalescent model (Rannala and Yang, 2003) using the program MCCOAL, which is part of the BPP pacakge (Rannala and Yang, 2003; Yang and Rannala, 2010). Given the gene tree, the sequences were allowed to evolve along the branches of the tree, under the JC69 mutation model (Jukes and Cantor, 1969). The resulting sequences at the tips of the tree constituted the data. Each replicate dataset thus consisted of $L$ sequence alignments, with 500 base pairs at each locus. We considered three kinds of data: (a) all loci of configuration 123, (b) a mixture of loci of configurations 11 and 12 in equal proportions, and (c) a mixture of loci of configurations 113 and 123 in equal proportions. The number of replicates was 1000.

Overall, the use of the $\chi_2^2$ distribution for data of configuration (a) 123 made the test conservative, as the false positive rate was always <1%, while an error rate of 5% was allowed (Table 5). For the 'pairs' data (configuration b, 11&12), we observed false positive rates of up to10% for parameter sets 2 and 3. The analysis seemed to suffer from a lack of information when only two sequences were available at each locus. In theory the false positive rate should converge to 5% when the number of loci increases, so it appears that more loci are needed for the asymptotics to be reliable for the 'pairs' data than for the 'triplet' data (c: 113&123). Adding an outgroup sequence increased the information content in the data, reducing the false positive rate to below 5%.

We examined the power of the test by simulating sequence alignments under the symmetrical version of M2 (gene flow). We used parameter values of Set 1 (hominoid) and Set 2 (mangroves), with $M_{12} = M_{21} = 1$ (Table 6). The test has virtually no power with $L = 10$ loci. With $L = 100$ or 1000, there are large performance differences between the two sets of parameter values. This is because the sequences are far more divergent and thus more informative for the mangroves set than for the hominoid set. Power is quite high with 1000 loci, when three sequences are used at each locus. Power is similar for the '123' data and for the '113&123' data. There is dramatic difference in power between the 'pairs' data (b, 11&12) and the 'triplet' data (c, 113&123). The use of the outgroup species improves the power of the test dramatically. This is consistent with Lohse *et al.* (2011), who suggested that triplet samples provide qualitatively new information about historical parameters in the joint distribution of topologies and branch lengths.

517     Table 7 lists the means and standard deviations of the MLEs of parameters under model M2 for

518     the same data analyzed in Table 6.  Datasets with '123' loci only suffer from the problem of

519     unidentifiability and do not allow the estimation of the migration rate.  Inclusion of the '113' loci

520     allows the model to estimate $\theta_1$ (=$\theta_2$) and $M$ and the unidentifiability problem disappears, leading to

521     better parameter estimation.  Furthermore the 'triplet' data provided much better parameter estimates

522     than the 'pair' data.

523     We also simulated data under the general (asymmetrical) model M2 (gene flow) to examine the

524     estimation of migration rates.  Given that the estimation was poor for the 'pair' data even under the

525     symmetrical model (Table 7) and that the asymmetrical model involves even more parameters, we

526     focus on the 'triplet' data only, with three sequences per locus.  We used the mangrove set of

527     parameters, with the migration rates set at $M_{12} = 0.1$ and $M_{21} = 1$ migrant individuals per generation.

528     We explored two different data configurations, with each dataset consisting of (a) '223' and '123' loci

529     in equal proportions, and (b) '113', '223', and '123' loci in equal proportions (Table 8).  The results

530     suggest that 100 loci may be too few to obtain reliable parameter estimates.  In particular, the lack of

531     polymorphism data for species 1 in the 223&123 configuration led to large fluctuations in the

532     estimates of $\theta_5$, $\theta_1$ and $M_{21}$.  Even with 1000 loci, we encountered several datasets in which the MLEs

533     of parameters hit the boundary set in the program (with $M_{12} = M_{21} = 0$), or the MLEs imply a star tree

534     (with $\tau_0 \approx \tau_1$ and $\theta_5 \approx 0$ or $\infty$).  With 15000 loci, the estimates are close to the true values.  Estimates

535     of migration rates are seen to involve a positive bias, but the bias is small with 15000 loci.  To fit the

536     asymmetrical IM model, it appears important to include thousands of loci, and to include population

537     data for both species 1 and 2 (such as '113' and '223' loci), as well as the '123' loci.

538

### *Analysis of* Drosophila *genomic datasets*

540     For each of the five datasets (Table 4), we performed three runs of 3S and used the results from the

541     run with the highest log likelihood.  Integration over coalescent times in the gene trees used Gaussian

542     quadrature with $K = 16$ points.  We used both the symmetrical and asymmetrical versions of models

543     M0 and M2, but here we focus on the asymmetrical models as they fit the data much better (Table 9).

544     We describe some general features of the results before discussing results specific to individual

545     datasets.  In every dataset, the LRT comparing M0 and M2 is significant.  Furthermore, the parameter

546     estimates under M2 suggest no migration from *D. melanogaster* to *D. simulans*, and about 0.016 to

547     0.044 immigrants per generation from *D. simulans* to *D. melanogaster*.  The consistency among the

548     datasets suggests that this pattern of unidirectional migration may be real.  Estimates of $\tau$ and $\theta$

549     parameters have very small standard errors because of the large size of the datasets.  Parameter

550     estimates are nearly identical between datasets D1 (auto) and D2 (noncoding), and between D4 (exons

551     complete) and D5 (exons split), suggesting that with such large genomic datasets, how extensively the

552     genomes were sampled to compile the datasets did not matter much.  Note that the autosomal dataset

553   D1 (auto) is dominated by noncoding DNA, even though different noncoding loci may be included in

554   D1 and D2, and that loci in D5 (exons split) are a subset of those in D4 (exons complete). While

555   model M0 did not fit the data as well as M2, it produced stable and reasonable estimates of $\theta$ and $\tau$

556   parameters, which were also similar to estimates from M2. (The exon datasets D4 and D5 are

557   exceptions to this pattern, to be discussed later.) For example, in datasets D1 (auto) and D2

558   (noncoding), both M0 and M2 estimates suggest that $\theta_S$ ($\approx 0.013$) is more than twice as large as $\theta_M$

559   ($\approx 0.005\text{-}0.006$), consistent with previous studies which suggest that *D. simulans* has a larger effective

560   population size than *D. melanogaster* (e.g., Langley et al., 2012; Wang and Hey, 2010). Also from

561   datasets D1 (auto) and D2 (noncoding) we obtained $\hat{\tau}_{MS} = 0.011$ and $\hat{\theta}_{MS} = 0.013\text{-}0.014$ under M0,

562   and $\hat{\tau}_{MS} = 0.012\text{-}0.014$ and $\hat{\theta}_{MS} = 0.011\text{-}0.012$ under M2 (Table 9). The slightly smaller estimates of

563   $\tau_{MS}$ and larger estimates of $\theta_M$ under M0 than under M2 may be expected because a more recent

564   divergence between *D. melanogaster* and *D. simulans* and a larger population size for *D.*

565   *melanogaster* may help M0 (which does not allow gene flow) to explain the genetic variation

566   introduced by immigrants from *D. simulans*.

567       Dataset D3 (chrX) for the X chromosome showed very different patterns from the autosomal

568   datasets D1 (auto) and D2 (noncoding), with a smaller estimate of $\theta_S$, and slightly larger estimates of

569   the other $\theta$ parameters. The estimated migration rate $M_{SM}$ was much higher for the X than for the

570   autosomes. By the simple model of random mating and neutral evolution, and assuming the same

571   mutation rate for the X and the autosomes, one would expect the effective population size for the X

572   chromosome to be ¾ that for the autosome, so that $\theta_S$ for X should be ¾ times as large as $\theta_S$ for the

573   autosomes, while the $\tau$s and $M$s should be identical. The parameter estimates suggested that this

574   simplistic model may not fit the data well. However the estimates of $\theta_M$ and $M_{SM}$ from D3 (chrX)

575   were associated with large sampling errors. Indeed D3 (chrX) does not include any MMY loci, so

576   that the data contain only very weak information concerning $\theta_M$ even though the model is identifiable.

577   The correlation between estimates of $\theta_M$ and $M_{SM}$ means that estimation of $M_{SM}$ may be affected as

578   well. We thus reran M2 under the constraint that $\theta_M = \frac{1}{2}\theta_S$ or $\theta_M = \theta_S$, obtaining estimates of $M_{SM}$ to

579   be 0.016 and 0.008 (Table 9). Thus there was no evidence for a large $M_{SM}$ for the X than for the

580   autosomes. The large changes to $\theta_M$ and $M_{SM}$ caused virtually no change to the log likelihood or to

581   estimates of other parameters, suggesting that the data are uninformative about $\theta_M$ and $M_{SM}$ while the

582   other parameters were well estimated. We leave it to future investigations, perhaps by including some

583   MMY or MMM loci with polymorphism for *D. melanogaster*, to generate more reliable parameter

584   estimates for the X and to understand possible differences in the evolutionary process between the X

585   chromosome and the autosomes.

586       The two exon datasets, D4 (exons complete) and D5 (exons split), are exceptional to the general

587   pattern of high similarity of parameter estimates between M0 and M2. For those two datasets,

588    estimates of $\tau_{MS}$ under M2 are much larger than those under M0. However those M2 estimates are

589    unreliable, because ML optimization under M2 converged to a star tree with $\tau_{MSY} \approx \tau_{MS}$ and $\theta_{MS} \approx 0$

590    (Table 9). We were unable to determine the reasons for this behavior. We note that the same

591    behavior was encountered in a few simulated datasets, as mentioned earlier, and that the problem did

592    not occur for dataset D1 (auto), which includes both coding and non-coding loci. The estimates of $\theta_M$

593    and $\theta_S$ from D4 (exons complete) and D5 (exons split) were smaller than those from D1 (auto) or D2

594    (noncoding), which can be explained by the reduced neutral mutation rate in the exons due to

595    selective constraint on nonsynonymous mutations. Again, the estimates suggest no migration from *D.*

596    *melanogaster* to *D. simulans*, but the migration rate from *D. simulans* to *D. melanogaster* is much

597    higher than for the autosome. We note that estimates of $\tau$ and $\theta$ parameters under M0 from those

598    exon datasets were similar to the M0 estimates from D1 (auto) and D2 (non-coding), and that the

599    estimates of $\tau_{MSY}$ were very similar between M0 and M2 for the same dataset. Thus we reran the M2

600    analysis of the two exon datasets, with $\tau_{MSY} = 0.020$ and $\tau_{MS} = 0.013$ fixed, to estimate the other

601    parameters. The results appear much more reasonable (Table 9). Both datasets D4 and D5 suggested

602    no migration from *D. melanogaster* to *D. simulans*, but the estimates of $M_{SM}$, at ~0.02 immigrants

603    from *D. simulans* to *D. melanogaster* per generation, were very similar to those from D1 (auto) and

604    D2 (noncoding).

605        To examine the robustness of our estimates of migration rates and to explore the impact of the

606    correlation between population sizes and migration rates, we re-analyzed the datasets under M2 (gene

607    flow) assuming asymmetrical migration rates (with $M_{MS} \neq M_{SM}$) but symmetrical population sizes ($\theta_M$

608    $= \theta_S$) (Table S7). Again the LRT is significant in every dataset, and parameter estimates suggested

609    unidirectional migration, with $\hat{M}_{MS} = 0$ in every dataset. However, estimates of $M_{SM}$ were much

610    larger than those of Table 9 in every dataset except for D3 (chrX), which has been discussed above.

611    For example, $\hat{M}_{SM} = 0.036\text{-}0.041$ from D1 (auto) and D2 (noncoding) under the constraint $\theta_M = \theta_S$

612    (Table S7), in comparison with 0.016-0.018 without the constraint (Table 9). We note that, except for

613    $\theta_M$ and $M_{SM}$, the parameter estimates were virtually identical with and without the constraint $\theta_M = \theta_S$

614    (compare Tables S7 and 9). There are far more SSY than MMY loci in those datasets (Table 4), so

615    that the estimates of $\theta_M = \theta_S$, at 0.012 (Table S7), were dominated by the *D. simulans* polymorphism

616    data, and were too large for *D. melanogaster*. This has lead to overestimates of $M_{SM}$, apparently

617    because a large $M_{SM}$ is more compatible with the (unrealistically assumed) large $\theta_M$. Thus the

618    assumption $\theta_M = \theta_S$ has caused serious biases in the estimation of migration rates, highlighting the

619    importance of the asymmetrical model. Note that the data contain strong evidence against the

620    assumption $\theta_M = \theta_S$; for example, relaxing the assumption improves the log likelihood by 66-82 units

621    in datasets D1 (auto) and D2 (noncoding). D3 (chrX) does not include any MMY loci. As a result,

622    $\theta_M$ is unidentifiable under M0 (so that the log likelihood is the same with and without the constraint

623      $\theta_M = \theta_S$), while under M2, $\theta_M$ is identifiable but very poorly estimated (so that the log likelihoods are

624      distinct but extremely similar with and without the constraint) (Tables 9 and S7).

625      We used equation (11) to calculate the posterior probabilities for gene trees for the MSY loci in

626      the five datasets (Table 4). Here we discuss the results for D5 (exons split) (Fig. 4), and those for D1

627      (auto) and D3 (chrX) are presented in Figs. S2 and S3. At the MLEs under M2 (Table 9, with $\tau_{MSY} =$

628      0.020 and $\tau_{MS} = 0.013$ fixed), the expected gene tree probabilities for any MSY locus are $P(G_{3c}) =$

629      0.1324, $P(G_{5c}) = 0.7368$, and $P(G_{6c}) = P(G_{6a}) = P(G_{6b}) = 0.0436$, with the gene tree-species tree

630      mismatch probability $P(G_{6a}) + P(G_{6b}) = 0.0872$. Most loci have gene tree $G_{5c}$ (Fig. 4), because the

631      migration rate is low, so that $G_{3c}$ is uncommon and because the outgroup species is quite distant so

632      that there is not much gene tree-species tree discordance. A small proportion of loci very likely have

633      the gene tree $G_{3c}$, and are likely to have been transferred across species (from *D. simulans* to *D.*

634      *melanogaster* since $M_{MS} \approx 0$). The top 41 loci, with $P(G_{3c}) > 95\%$, are listed in Table S8. More than

635      half of those loci were also inferred to have $P(G_{3c}) > 95\%$ in the analysis of dataset D4 (exons

636      complete) (Table S8), suggesting that this inference was not very sensitive to the different filtering

637      procedures applied to compile the datasets.

638      An intriguing feature in Fig. 4 (and also in Figs. S2 and S3 for datasets D1 and D3) is that many

639      more loci seem to support gene tree $G_{6c}$ than $G_{6a}$ or $G_{6b}$, while the model predicts equal proportions

640      for those three gene trees. This is in contrast to the simulated dataset, in which the three gene trees

641      are inferred to occur with similar proportions, as expected under the model (Fig. 3A). The reasons for

642      this pattern are unknown, but are likely to be some kind of model violation.

643      To explore the potential of the IM model for species tree estimation under the multispecies

644      coalescent with migration, we applied model M2 to dataset D1 (auto), assuming alternative species

645      trees for M, S, and Y. The MLEs and log likelihood values are shown in Table 10. The ((MS)Y) tree

646      has a much greater log likelihood value than the two alternative trees (by about 20,000 units). Indeed,

647      both alternative trees converge to the star tree with $\tau_0 = \tau_1$. Migration is detected only in the direction

648      of S→M when the assumed tree is ((MS)Y). Note that our model assumes migration between the two

649      ingroup species only. In theory a stratified bootstrap resampling procedure can be used to assess the

650      significance of the ML species tree, sampling loci and then sampling sites for each sampled locus.

651      This is not pursued here since there does not seem to be any uncertainty about the species phylogeny

652      in this case (Russo et al., 1995; Obbard et al., 2012).

653

## DISCUSSION

655 *Utilities and limitations of our implementation*

656 In this paper, we have extended our previous implementation of the IM model (Zhu and Yang, 2012)

657 in several important ways. First, we have relaxed the symmetry assumption, so that the test of gene

658    flow and estimation of migration rates and population size parameters can be conducted under more

659    realistic models.  For the *Drosophila* datasets, our analyses suggest that gene flow is indeed

660    asymmetrical, the population sizes of *D. melanogaster* and *D. simulans* are very different, and

661    accounting for such asymmetries in the model is important to accurate estimation of the migration

662    rates.  Second, we have extended the implementation so that a locus can have 2 or 3 sequences of

663    arbitrary configurations.  This removes the unidentifiability problem that we encountered when '123'

664    loci alone were used, making it possible to estimate the migration rates.  It also improves the power of

665    the LRT of gene flow because the null distribution becomes known.  The extension to arbitrary loci

666    also paves the way for implementing more complex models of migration.

667        We envisage that a major future use of the IM model is to infer species phylogenies under the

668    multispecies coalescent model with migration, accommodating two major factors that thwart species

669    tree estimation, especially for species formed during radiative speciations: incomplete lineage sorting

670    (ILS) and gene flow (Mallet et al., 2016).  Heuristic methods based on the model that treat estimated

671    gene tree topologies as observed data are being developed (Wen et al., 2016), but full likelihood

672    methods have the advantage of accommodating the different sources of uncertainties appropriately.

673    However the functionality of 3S in this regard is limited.  The assumption of gene flow between sister

674    species only may be too restrictive and gene flow between non-sister species needs to be allowed as

675    well (Mallet et al., 2016).  Furthermore, our implementation is restricted to three species, with two or

676    three sequences per locus.  This limitation is mainly due to our use of numerical integration (Gaussian

677    quadrature) to integrate over the coalescent times, with the dimension of the integrals to be one less

678    than the number of sequences at the locus.  With four or more sequences per locus, this calculation

679    may not be feasible.  Furthermore, the number of states in the Markov chain used to characterize the

680    genealogical process also increases explosively with the increase of the number of sequences per

681    locus (Andersen et al., 2014).  We suggest that to analyze genomic datasets involving more than three

682    species and more than three sequences per locus, a subsampling procedure may be useful, similarly to

683    our analysis of the *Drosophila* datasets (see also Wang and Hey, 2010).  Suppose there are $s > 3$

684    species.  We specify a 'master' species tree including all $s$ species and use it to define the parameters:

685    the $(s-1)$ species divergence times ($\tau$s) and up to $(2s-1)$ population size parameters ($\theta$s).  At every

686    locus, we sample three sequences, which may be from different species, so that the data

687    configurations may be 123, 114, 255, etc.  The species tree for the sequences of any particular locus

688    can be constructed from the master species tree by pruning off branches for species not sampled in the

689    data at the locus.  The theory developed in Zhu and Yang (2012) and in this paper will then be

690    applicable with the only complication that the coalescent rate (the population size) and the migration

691    rate may change along the same branch on the species subtree at the speciation events in the master

692    species tree.  Such rate changes are relatively straightforward to accommodate.  This strategy involves

693    filtering of data but the information loss may not be very serious for such large genomic datasets.

694    Note that given the data, this strategy calculates the likelihood correctly.

695       In the future, we also hope to implement models of nonhomogeneous migration rates over time.

696      Gene flow may be common at the early stage of species formation and decrease until the two

697      populations achieve complete isolation. A simple model may assume a constant migration rate $M$

698      since species divergence until a time point $T$ $(0 < T < \tau_1)$ when gene flow ceases. In this model of

699      *isolation with initial migration*, both the migration rate $M$ and the time point $T$ will be parameters to

700      be estimated from the sequence data (Wilkinson-Herbots, 2012). The same Markov chain

701      characterization as used here can be used to derive the density of gene trees by breaking the time

702      epoch $E_1$ into two segments: $E_{1a}$: $0 < t < T$ and $E_{1b}$: $T < t < \tau_1$. Alternatively, one may use a

703      deterministic mathematical function such as an exponential decay to describe the changing migration

704      rate over time. The initial migration rate and the exponential decay rate will be parameters to be

705      estimated. If reproductive isolation builds up gradually after species split, such nonhomogeneous

706      migration models may be more realistic than the usual IM model with a constant migration rate after

707      species divergence.

708       Similarly, introgression or hybridisation may be modelled in the same framework (Twyford and

709      Ennos, 2011). Recent introgression or contamination may be modelled by assuming that a proportion

710      of individuals sampled from species 1 are in fact from species 2. Introgression can then be tested

711      using a likelihood ratio test. As the model naturally accommodates ancestral polymorphism and

712      incomplete lineage sorting (ILS), the test will distinguish introgression from ILS. Note that

713      introgression affects all loci of the introgressed individual, while with ILS, caused by the coalescent

714      process, the different genomic loci have independent histories.

### 715 *Asymmetrical Migration in* Drosophila *fruit flies*

716 Wang and Hey (2010: Table 7) compiled and analyzed a *Drosophila* dataset similar to our dataset D1

717 (auto), consisting of 30,323 autosomal loci but including only two sequences for each locus, of

718 configurations SS, MS, and MM. Under the asymmetrical model, their estimates of population size

719 parameters are $\theta_M = 0.0055$ and $\theta_S = 0.01352$, which are close to our estimates from D1 (auto). The

720 ancestral population size $\theta_{MS}$ estimated by Wang and Hey ranges from 0.007 to 0.010, whereas our

721 estimates are larger, at $\theta_{MS} = 0.011$ and $\theta_{MSY} = 0.040$. The M-S divergence time parameter is $\tau_{MS} =$

722 0.017 by Wang and Hey and 0.0136 in our analysis. A strong negative correlation between $\tau_{MS}$ and

723 $\theta_{MS}$ is expected in such analyses (Yang, 2002). Wang and Hey (2010) estimated the migration rate (in

724 our notation) to be $M_{MS} = N_S m_{MS} = 0$ from *D. melanogaster* to D. *simulans* and $M_{SM} = N_M m_{SM} =$

725 $4.846 \times 0.00552/4 = 0.0067$ from *simulans* to *melanogaster*. Our estimates under M2 are $M_{MS} = 0$ as

726 well and $M_{SM} = 0.0183$, which is much larger.

727       The data of Wang and Hey (2010) were also analyzed by Lohse *et al.* (2011, Table 1), who

728 compared parameter estimates from two datasets which have either two or three sequences per locus

729 for the same set of loci. The authors found that the estimate of the migration rate from the 'triplet'

730 data was nearly twice as large as that for the 'pair' data. This is consistent with our finding.

731    We note that our datasets are based on updated genome sequences, relative to the data analyzed

732    by Wang and Hey (2010) and Lohse *et al.* (2011).  Also different filters were applied and different

733    loci were included in those datasets.  Furthermore, Wang and Hey (2010) removed loci at which the

734    pairwise sequence distances indicated gene tree-species tree conflict.  We did not apply this filtering

735    because such loci are informative about the gene tree distribution and about the parameters in our

736    analysis of loci of three sequences.  Lohse *et al.* (2011) removed highly variable loci and highly

737    variable sites so that the data could be analyzed under the infinite-sites model.  Given the multiple

738    differences among the datasets, we conclude that the estimates obtained from those studies are largely

739    consistent.

740    Different from Wang and Hey (2010), we also compiled and analyzed a dataset for the X

741    chromosome (D3 chrX) as well as two exon datasets: D4 (exons complete) and D5 (exons split).  The

742    use of multiple datasets, even though some of them are overlapping, allows us to confirm the

743    robustness of our analyses, as processes such as migration are expected to have genome-wide effects,

744    and to discover similarities and differences in the evolutionary process among different parts of the

745    genome.  Indeed all five datasets we analyzed support a model of unidirectional gene flow, from *D.*

746    *simulans* to *D. melanogaster*, at the rate of ~0.02 migrant individuals per generation.  We included the

747    two exon datasets even though we do not expect exons to be evolving neutrally.  Note that the

748    multispecies coalescent model implemented in 3S assumes neutral evolution of the gene sequences,

749    such that mutations in the sequences do not affect the genealogical process or the gene tree

750    distribution.  Nevertheless, most proteins appear to perform the same conserved function in closely

751    related species and their coding genes are under similar purifying selection in the different species.

752    The main effect of the selective constraint may then be a reduction of the neutral mutation rate.

753    Species-specific natural selection such as positive selection would be more problematic but loci

754    undergoing positive selection or responsible for between-species incompatibilities are expected to be

755    rare.  Similar points have been made by Ebersberger et al. (2007;  see also Yang, 2015) in their

756    analysis of hominoid genomic sequence data.

757

765

766

22

## References

Andersen, L. N., T. Mailund, and A. Hobolth. 2014. Efficient computation in the IM model. J. Math. Biol. 68:1423-1451.

Attrill, H., K. Falls, J. L. Goodman, G. H. Millburn, G. Antonazzo, A. J. Rey, and S. J. Marygold. 2016. FlyBase: establishing a Gene Group resource for *Drosophila melanogaster*. Nucleic Acids Res. 44:D786-792.

Bahlo, M., and R. C. Griffiths. 2000. Inference from gene trees in a subdivided population. Theor. Popul. Biol. 57:79-95.

Beerli, P. 2006. Comparison of Bayesian and maximum-likelihood inference of population genetic parameters. Bioinformatics 22:341-345.

Beerli, P., and J. Felsenstein. 1999. Maximum-likelihood estimation of migration rates and effective population numbers in two populations using a coalescent approach. Genetics 152:763-773.

Beerli, P., and J. Felsenstein. 2001. Maximum likelihood estimation of a migration matrix and effective population sizes in *n* subpopulations by using a coalescent approach. Proc. Natl. Acad. Sci. U.S.A. 98:4563-4568.

Burgess, R., and Z. Yang. 2008. Estimation of hominoid ancestral population sizes under Bayesian coalescent models incorporating mutation rate variation and sequencing errors. Mol. Biol. Evol. 25:1979-1994.

Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, and T. L. Madden. 2009. BLAST+: architecture and applications. BMC Bioinformatics 10:421.

Chan, Y. C., C. Roos, M. Inoue-Murayama, E. Inoue, C. C. Shih, K. J. Pei, and L. Vigilant. 2013. Inferring the evolutionary histories of divergences in Hylobates and Nomascus gibbons through multilocus sequence data. BMC Evol. Biol. 13:82.

Dagum, L., and R. Menon. 1998. OpenMP: an industry standard API for shared-memory programming. Computational Science & Engineering, IEEE5 1:46-55.

Ebersberger, I., P. Galgoczy, S. Taudien, S. Taenzer, M. Platzer, and A. von Haeseler. 2007. Mapping human genetic ancestry. Mol. Biol. Evol. 24:2266-2276.

Edwards, S. V. 2009. Is a new and general theory of molecular systematics emerging? Evolution 63:1–19.

Ellegren, H., L. Smeds, R. Burri, P. I. Olason, N. Backstrom, T. Kawakami, A. Kunstner, H. Makinen, K. Nadachowska-Brzyska, A. Qvarnstrom, S. Uebbing, and J. B. W. Wolf. 2012. The genomic landscape of species divergence in Ficedula flycatchers. Nature 491:756-760.

Fontaine, M. C., J. B. Pease, A. Steele, R. M. Waterhouse, D. E. Neafsey, I. V. Sharakhov, X. Jiang, A. B. Hall, F. Catteruccia, E. Kakani, S. N. Mitchell, Y. C. Wu, H. A. Smith, R. R. Love, M. K. Lawniczak, M. A. Slotman, S. J. Emrich, M. W. Hahn, and N. J. Besansky. 2015. Mosquito genomics. Extensive introgression in a malaria vector species complex revealed by phylogenomics. Science 347:1258524.

Galassi, M., J. Davies, J. Theiler, B. Gough, R. Priedhorsky, G. Jungman, and M. Booth. 2013. *GNU Scientific Library Reference Manual*. The GSL Project.

Gronau, I., M. J. Hubisz, B. Gulko, C. G. Danko, and A. Siepel. 2011. Bayesian inference of ancient human demography from individual genome sequences. Nature Genet. 43:1031-1034.

Hey, J. 2010. Isolation with migration models for more than two populations. Mol. Biol. Evol. 27:905-920.

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 coalescence time density in an isolation-with-migration model with few samples. Genetics 187:1241-1243.

Hu, T. T., M. B. Eisen, K. R. Thornton, and P. Andolfatto. 2013. A second-generation assembly of the *Drosophila simulans* genome provides new insights into patterns of lineage-specific divergence. Genome Res. 23:89-98.

Hutter, S., H. Li, S. Beisswanger, D. De Lorenzo, and W. Stephan. 2007. Distinctly different sex ratios in African and European populations of *Drosophila melanogaster* inferred from

820    chromosomewide single nucleotide polymorphism data. Genetics 177:469-480.
821    Innan, H., and H. Watanabe. 2006. The effect of gene flow on the coalescent time in the human-
822        chimpanzee ancestral population. Mol. Biol. Evol. 23:1040-1047.
823    Jukes, T. H., and C. R. Cantor. 1969. Evolution of protein molecules. Pages 21-123 *in* Mammalian
824        Protein Metabolism (H. N. Munro, ed.) Academic Press, New York.
825    Katoh, K., and D. M. Standley. 2013. MAFFT multiple sequence alignment software version 7:
826        improvements in performance and usability. Mol. Biol. Evol. 30:772-780.
827    Kimura, M., and W. H. Weiss. 1964. The stepping stone model of genetic structure and the decrease
828        of genetic correlation with distance. Genetics 49:561-576.
829    Kingman, J. F. C. 1982. The coalescent. Stochastic Process Appl. 13:235-248.
830    Kutschera, V. E., T. Bidon, F. Hailer, J. L. Rodi, S. R. Fain, and A. Janke. 2014. Bears in a forest of
831        gene trees: Phylogenetic inference is complicated by incomplete lineage sorting and gene
832        flow. Mol. Biol. Evol. 31:2004-2017.
833    Langley, C. H., K. Stevens, C. Cardeno, Y. C. Lee, D. R. Schrider, J. E. Pool, S. A. Langley, C.
834        Suarez, R. B. Corbett-Detig, B. Kolaczkowski, S. Fang, P. M. Nista, A. K. Holloway, A. D.
835        Kern, C. N. Dewey, Y. S. Song, M. W. Hahn, and D. J. Begun. 2012. Genomic variation in
836        natural populations of *Drosophila melanogaster*. Genetics 192:533-598.
837    Leaché, A. D., R. B. Harris, M. E. Maliska, and C. W. Linkem. 2013. Comparative species divergence
838        across eight triplets of spiny lizards (Sceloporus) using genomic sequence data. Genome Biol.
839        Evol. 5:2410-2419.
840    Li, W.-H. 1976. Distribution of nucleotide differences between two randomly chosen cistrons in a
841        subdivided population: the finite island model. Theor Popul Biol 10:303-308.
842    Liu, S., E. D. Lorenzen, M. Fumagalli, B. Li, K. Harris, Z. Xiong, L. Zhou, T. S. Korneliussen, M.
843        Somel, C. Babbitt, G. Wray, J. Li, W. He, Z. Wang, W. Fu, X. Xiang, C. C. Morgan, A.
844        Doherty, M. J. O'Connell, J. O. McInerney, E. W. Born, L. Dalen, R. Dietz, L. Orlando, C.
845        Sonne, G. Zhang, R. Nielsen, E. Willerslev, and J. Wang. 2014. Population genomics reveal
846        recent speciation and rapid evolutionary adaptation in polar bears. Cell 157:785-794.
847    Lohse, K., R. J. Harrison, and N. H. Barton. 2011. A general method for calculating likelihoods under
848        the coalescent process. Genetics 189:977-987.
849    Mallet, J. 2005. Hybridization as an invasion of the genome. Trends Ecol. Evol. 20:229-237.
850    Mallet, J., N. Besansky, and M. W. Hahn. 2016. How reticulated are species? BioEssays.
851    Martin, S. H., K. K. Dasmahapatra, N. J. Nadeau, C. Salazar, J. R. Walters, F. Simpson, M. Blaxter,
852        A. Manica, J. Mallet, and C. D. Jiggins. 2013. Genome-wide evidence for speciation with
853        gene flow in Heliconius butterflies. Genome Res 23:1817-1828.
854    Melo-Ferreira, J., P. Boursot, M. Carneiro, P. J. Esteves, L. Farelo, and P. C. Alves. 2012. Recurrent
855        introgression of mitochondrial DNA among hares (Lepus spp.) revealed by species-tree
856        inference and coalescent simulations. Syst. Biol. 61:367-381.
857    Nath, H. B., and R. C. Griffiths. 1993. The coalescent in two colonies with symmetric migration. J.
858        Math. Biol. 31:841–852.
859    Nielsen, R., and J. Wakeley. 2001. Distinguishing migration from isolation: a Markov chain Monte
860        Carlo approach. Genetics 158:885-896.
861    Nielsen, R., and Z. Yang. 1998. Likelihood models for detecting positively selected amino acid sites
862        and applications to the HIV-1 envelope gene. Genetics 148:929-936.
863    Notohara, M. 1990. The coalescent and the genealogical process in geographically structured
864        populations. J. Math. Biol. 29:59-75.
865    Obbard, D. J., J. Maclennan, K. W. Kim, A. Rambaut, P. M. O'Grady, and F. M. Jiggins. 2012.
866        Estimating divergence dates and substitution rates in the *Drosophila* phylogeny. Mol. Biol.
867        Evol. 29:3459-3473.
868    Palmieri, N., V. Nolte, J. Chen, and C. Schlotterer. 2014. Genome assembly and annotation of a
869        *Drosophila simulans* strain from Madagascar. Mol. Ecol. Resour. 15:372-381.
870    Patterson, N., D. J. Richter, S. Gnerre, E. S. Lander, and D. Reich. 2006. Genetic evidence for
871        complex speciation of humans and chimpanzees. Nature 441:1103-1108.
872    Rannala, B., and Z. Yang. 2003. Bayes estimation of species divergence times and ancestral
873        population sizes using DNA sequences from multiple loci. Genetics 164:1645-1656.
874    Russo, C. A., N. Takezaki, and M. Nei. 1995. Molecular phylogeny and divergence times of

875           Drosophilid species. Mol. Biol. Evol. 12:391-404.

876 Saitou, N. 1988. Property and efficiency of the maximum likelihood method for molecular phylogeny.
877           J. Mol. Evol. 27:261-273.

878 Self, S. G., and K.-Y. Liang. 1987. Asymptotic properties of maximum likelihood estimators and
879           likelihood ratio tests under nonstandard conditions. J. Am. Stat. Assoc. 82:605-610.

880 Strobeck, K. 1987. Average number of nucleotide differences in a sample from a single
881           subpopulation: A test for population subdivision. Genetics 117:149-153.

882 Takahata, N. 1988. The coalescent in two partially isolated diffusion populations. Genet. Res.
883           (Camb.) 52:213-222.

884 Takahata, N., Y. Satta, and J. Klein. 1995. Divergence time and population size in the lineage leading
885           to modern humans. Theor. Popul. Biol. 48:198-221.

886 Twyford, A. D., and R. A. Ennos. 2011. Next-generation hybridization and introgression. Heredity
887           108:179-189.

888 Wang, Y., and J. Hey. 2010. Estimating divergence parameters with small samples from a large
889           number of loci. Genetics 184:363-379.

890 Wen, D., Y. Yu, M. W. Hahn, and L. Nakhleh. 2016. Reticulate evolutionary history and extensive
891           introgression in mosquito species revealed by phylogenetic network analysis. Mol. Ecol.

892 Wilkinson-Herbots, H. M. 1998. Genealogy and subpopulation differentiation under various models
893           of population structure. J. Math. Biol. 37:535-585.

894 Wilkinson-Herbots, H. M. 2008. The distribution of the coalescence time and the number of pairwise
895           nucleotide differences in the "isolation with migration" model. Theor. Popul. Biol. 73:277-
896           288.

897 Wilkinson-Herbots, H. M. 2012. The distribution of the coalescence time and the number of pairwise
898           nucleotide differences in a model of population divergence or speciation with an initial period
899           of gene flow. Theor. Popul. Biol. 82:92-108.

900 Wright, S. 1931. Evolution in Mendelian populations. Genetics 16:97-159.

901 Wright, S. 1943. Isolation by distance. Genetics 28:114-138.

902 Yamamichi, M., J. Gojobori, and H. Innan. 2012. An autosomal analysis gives no genetic evidence for
903           complex speciation of humans and chimpanzees. Mol. Biol. Evol. 29:145-156.

904 Yang, Z. 1994. Statistical properties of the maximum likelihood method of phylogenetic estimation
905           and comparison with distance matrix methods. Syst. Biol. 43:329-342.

906 Yang, Z. 2002. Likelihood and Bayes estimation of ancestral population sizes in Hominoids using
907           data from multiple loci. Genetics 162:1811-1823.

908 Yang, Z. 2010. A likelihood ratio test of speciation with gene flow using genomic sequence data.
909           Genom. Biol. Evol. 2:200-211.

910 Yang, Z. 2014. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford,
911           England.

912 Yang, Z. 2015. The BPP program for species tree estimation and species delimitation. Curr. Zool.
913           61:854-865.

914 Yang, Z., S. Kumar, and M. Nei. 1995. A new method of inference of ancestral nucleotide and amino
915           acid sequences. Genetics 141:1641-1650.

916 Yang, Z., and B. Rannala. 2010. Bayesian species delimitation using multilocus sequence data. Proc.
917           Natl. Acad. Sci. U.S.A. 107:9264-9269.

918 Zhang, C., D.-X. Zhang, T. Zhu, and Z. Yang. 2011. Evaluation of a Bayesian coalescent method of
919           species delimitation. Syst. Biol. 60:747-761.

920 Zhou, R., K. Zeng, W. Wu, X. Chen, Z. Yang, S. Shi, and C.-I. Wu. 2007. Population genetics of
921           speciation in nonmodel organisms: I. ancestral polymorphism in mangroves. Mol. Biol. Evol.
922           24:2746-2754.

923 Zhou, W. W., Y. Wen, J. Fu, Y. B. Xu, J. Q. Jin, L. Ding, M. S. Min, J. Che, and Y. P. Zhang. 2012.
924           Speciation in the *Rana chensinensis* species complex and its relationship to the uplift of the
925           Qinghai-Tibetan Plateau. Mol. Ecol. 21:960-973.

926 Zhu, T., and Z. Yang. 2012. Maximum likelihood implementation of an isolation-with-migration
927           model with three species for testing speciation with gene flow. Mol. Biol. Evol. 29:3131-
928           3142.

929

930

## APPENDIX A.

## DISTRIBUTION OF GENE TREES FOR THREE SEQUENCES UNDER M2 (GENE FLOW)

### Case I: Initial states 111 and 222

With the initial state $s = 111$ or $222$, all three sequences at the locus are from the same species (1 or 2). Due to the symmetry, the densities of the three gene trees of the same shape (such as $G_{1c}$, $G_{1a}$, and $G_{1b}$) are identical. There is thus no need to keep track of the sequence IDs, even though the likelihood averages over all 18 gene trees (Table S1). Thus we consider a Markov chain with 8 states: 111, 112, 122, 222, 11, 12, 22, 1|2, with '1|2' to be an artificial state formed by merging states 1 and 2. The rate matrix is given in Table 3. The density for gene tree shape $G_1$ is given in equation (9). By a similar argument we obtain the densities for tree shapes $G_2$-$G_6$, as follows.

$$f(G_2, t_0, t_1) = \frac{2}{\theta_5} e^{-\frac{2}{\theta_5} t_0} \times$$
$$\sum_{j \in S_2} \left[ 3 \frac{2}{\theta_1} P_{s,111}(t_1) P_{11,j}(\tau_1 - t_1) + \frac{2}{\theta_1} P_{s,112}(t_1) P_{12,j}(\tau_1 - t_1) + \frac{2}{\theta_2} P_{s,221}(t_1) P_{12,j}(\tau_1 - t_1) + 3 \frac{2}{\theta_2} P_{s,222}(t_1) P_{22,j}(\tau_1 - t_1) \right],$$

$$f(G_3, t_0, t_1) = \frac{2}{\theta_4} e^{-\frac{2}{\theta_5}(\tau_0 - \tau_1)} e^{-\frac{2}{\theta_4} t_0} \times$$
$$\sum_{j \in S_2} \left[ 3 \frac{2}{\theta_1} P_{s,111}(t_1) P_{11,j}(\tau_1 - t_1) + \frac{2}{\theta_1} P_{s,112}(t_1) P_{12,j}(\tau_1 - t_1) + \frac{2}{\theta_2} P_{s,122}(t_1) P_{12,j}(\tau_1 - t_1) + 3 \frac{2}{\theta_2} P_{s,222}(t_1) P_{22,j}(\tau_1 - t_1) \right],$$

$$f(G_4, t_0, t_1) = \frac{6}{\theta_5} e^{-\frac{6}{\theta_5} t_1} \frac{2}{\theta_5} e^{-\frac{2}{\theta_5} t_0} \times \sum_{j \in S_3} P_{s,j}(\tau_1), \qquad 0 < t_1 + t_0 < \tau_0 - \tau_1,$$

$$f(G_5, t_0, t_1) = \frac{6}{\theta_5} e^{-\frac{6}{\theta_5} t_1} e^{-\frac{2}{\theta_5}(\tau_0 - \tau_1 - t_1)} \frac{2}{\theta_4} e^{-\frac{2}{\theta_4} t_0} \times \sum_{j \in S_3} P_{s,j}(\tau_1), \qquad 0 < t_1 < \tau_0 - \tau_1, \ 0 < t_0 < \infty,$$

$$f(G_6, t_0, t_1) = e^{-\frac{6}{\theta_5}(\tau_0 - \tau_1)} \frac{6}{\theta_4} e^{-\frac{6}{\theta_4} t_1} \frac{2}{\theta_4} e^{-\frac{2}{\theta_4} t_0} \times \sum_{j \in S_3} P_{s,j}(\tau_1), \qquad 0 < t_0, t_1 < \infty,$$

$$(12)$$

where $S_2$ and $S_3$ are the sets of states with two and three sequences, respectively, that can be reached by the initial state (Table 2). Again each density for a tree shape should be divided by 3 to give the density for the gene tree: e.g., $f(G_{2a}, t_0, t_1) = f(G_2, t_0, t_1)/3$.

### Case II: Initial states 112 and 122

For initial state $s = 112$ or $122$, the likelihood calculation at each locus averages over all 18 gene trees (Table S1). This is the only case in this study where it is necessary to keep track of both the sequence IDs and the population IDs in our Markov chain characterization of the process of coalescent with migration. The initial states are thus $1_a 1_b 2_c$ or $1_a 2_b 2_c$. However, for states of three sequences, we always arrange the sequence IDs in the order $a$, $b$, and $c$ to simplify the notation and thus the subscripts are dropped. Thus $1_a 1_b 1_c$, $1_a 1_b 2_c$ and $1_a 2_b 2_c$ are written as 111, 112 and 122, respectively. There are 21 states in the chain: 111, 112, 121, 122, 211, 212, 221, 222, $1_{bc} 1_a$, $1_{ca} 1_b$, $1_{ab} 1_c$, $1_{bc} 2_a$, $1_{ca} 2_b$, $1_{ab} 2_c$, $1_a 2_{bc}$, $1_b 2_{ca}$, $1_c 2_{ab}$, $2_{bc} 2_a$, $2_{ca} 2_b$, $2_{ab} 2_c$, and 1|2. The states of two sequences have the subscripts to

26

957      indicate the sequence IDs. For example, $1_{bc}2_a$ means that sequences $b$ and $c$ have coalesced and their

958      ancestor is in population 1 while sequence $a$ is in population 2.

959      For gene tree $G_{1c}$, with $0 < t_0 + t_1 < \tau_1$, we have

$$
f(G_{1c}, t_0, t_1)
$$

960

$$
= \tfrac{2}{\theta_1} P_{s,111}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{ab}1_c,1_{ab}1_c}(t_0) + \tfrac{2}{\theta_2} P_{1_{ab}1_c,2_{ab}2_c}(t_0)\right) + \tfrac{2}{\theta_1} P_{s,112}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{ab}2_c,1_{ab}1_c}(t_0) + \tfrac{2}{\theta_2} P_{1_{ab}2_c,2_{ab}2_c}(t_0)\right) \tag{13}
$$

$$
+ \tfrac{2}{\theta_2} P_{s,221}(t_1)\left(\tfrac{2}{\theta_1} P_{1_c2_{ab},1_{ab}1_c}(t_0) + \tfrac{2}{\theta_2} P_{1_c2_{ab},2_{ab}2_c}(t_0)\right) + \tfrac{2}{\theta_2} P_{s,222}(t_1)\left(\tfrac{2}{\theta_1} P_{2_{ab}2_c,1_{ab}1_c}(t_0) + \tfrac{2}{\theta_2} P_{2_{ab}2_c,2_{ab}2_c}(t_0)\right),
$$

961 The densities for gene trees $G_{1b}$ and $G_{1a}$ are similar.

$$
f(G_{1b}, t_0, t_1)
$$

962

$$
= \tfrac{2}{\theta_1} P_{s,111}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{ca}1_b,1_{ca}1_b}(t_0) + \tfrac{2}{\theta_2} P_{1_{ca}1_b,2_{ca}2_b}(t_0)\right) + \tfrac{2}{\theta_2} P_{s,212}(t_1)\left(\tfrac{2}{\theta_1} P_{1_b2_{ca},1_{ca}1_b}(t_0) + \tfrac{2}{\theta_2} P_{1_b2_{ca},2_{ca}2_b}(t_0)\right) \tag{14}
$$

$$
+ \tfrac{2}{\theta_1} P_{s,121}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{ca}2_b,1_{ca}1_b}(t_0) + \tfrac{2}{\theta_2} P_{1_{ca}2_b,2_{ca}2_b}(t_0)\right) + \tfrac{2}{\theta_2} P_{s,222}(t_1)\left(\tfrac{2}{\theta_1} P_{2_{ca}2_b,1_{ca}1_b}(t_0) + \tfrac{2}{\theta_2} P_{2_{ca}2_b,2_{ca}2_b}(t_0)\right),
$$

$$
f(G_{1a}, t_0, t_1)
$$

963

$$
= \tfrac{2}{\theta_1} P_{s,111}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{bc}1_a,1_{bc}1_a}(t_0) + \tfrac{2}{\theta_2} P_{1_{bc}1_a,2_{bc}2_a}(t_0)\right) + \tfrac{2}{\theta_2} P_{s,122}(t_1)\left(\tfrac{2}{\theta_1} P_{1_a2_{bc},1_{bc}1_a}(t_0) + \tfrac{2}{\theta_2} P_{1_a2_{bc},2_{bc}2_a}(t_0)\right)
$$

$$
+ \tfrac{2}{\theta_1} P_{s,211}(t_1)\left(\tfrac{2}{\theta_1} P_{1_{bc}2_a,1_{bc}1_a}(t_0) + \tfrac{2}{\theta_2} P_{1_{bc}2_a,2_{bc}2_a}(t_0)\right) + \tfrac{2}{\theta_2} P_{s,222}(t_1)\left(\tfrac{2}{\theta_1} P_{2_{bc}2_a,1_{bc}1_a}(t_{01}) + \tfrac{2}{\theta_2} P_{2_{bc}2_a,2_{bc}2_a}(t_0)\right),
$$

964      For gene tree $G_2$, we have $t_1 < \tau_1$, $t_0 < \tau_0 - \tau_1$, and

965

$$
f(G_{2c}, t_0, t_1) = \tfrac{2}{\theta_5} \mathrm{e}^{-\tfrac{2}{\theta_5}t_0} \times
$$

$$
\sum_{j \in S_2}\left[\tfrac{2}{\theta_1} P_{s,111}(t_1) P_{1_{ab}1_c,j}(\tau_1 - t_1) + \tfrac{2}{\theta_1} P_{s,112}(t_1) P_{1_{ab}2_c,j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,221}(t_1) P_{1_c2_{ab},j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,222}(t_1) P_{2_{ab}2_c,j}(\tau_1 - t_1)\right],
$$

966

$$
f(G_{2b}, t_0, t_1) = \tfrac{2}{\theta_5} \mathrm{e}^{-\tfrac{2}{\theta_5}t_0} \times
$$

$$
\sum_{j \in S_2}\left[\tfrac{2}{\theta_1} P_{s,111}(t_1) P_{1_{ca}1_b,j}(\tau_1 - t_1) + \tfrac{2}{\theta_1} P_{s,121}(t_1) P_{1_{ca}2_b,j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,212}(t_1) P_{1_b2_{ca},j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,222}(t_1) P_{2_{ca}2_b,j}(\tau_1 - t_1)\right],
$$

$$
f(G_{2a}, t_0, t_1) = \tfrac{2}{\theta_{12}} \tfrac{2}{\theta_5} \mathrm{e}^{-\tfrac{2}{\theta_5}t_0} \times
$$

$$
\sum_{j \in S_2}\left[\tfrac{2}{\theta_1} P_{s,111}(t_1) P_{1_{bc}1_a,j}(\tau_1 - t_1) + \tfrac{2}{\theta_1} P_{s,211}(t_1) P_{1_{bc}2_a,j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,122}(t_1) P_{1_a2_{bc},j}(\tau_1 - t_1) + \tfrac{2}{\theta_2} P_{s,222}(t_1) P_{2_{bc}2_a,j}(\tau_1 - t_1)\right].
$$

967                        (15)

968      For gene tree $G_3$, with $t_1 < \tau_1 < \tau_0 < t_0$, we have

$$f(G_{3c},t_0,t_1) = e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0} \times$$
$$\sum_{j\in S_2}\left[\frac{2}{\theta_1}P_{s,111}(t_1)P_{1_{ab}1_c,j}(\tau_1-t_1)+\frac{2}{\theta_1}P_{s,112}(t_1)P_{1_{ab}2_c,j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,221}(t_1)P_{1_c2_{ab},j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,222}(t_1)P_{2_{ab}2_c,j}(\tau_1-t_1)\right],$$

$$f(G_{3b},t_0,t_1) = \frac{2}{\theta_{12}}e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0} \times$$
$$\sum_{j\in S_2}\left[\frac{2}{\theta_1}P_{s,111}(t_1)P_{1_{ca}1_b,j}(\tau_1-t_1)+\frac{2}{\theta_1}P_{s,121}(t_1)P_{1_{ca}2_b,j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,212}(t_1)P_{1_b2_{ca},j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,222}(t_1)P_{2_{ca}2_b,j}(\tau_1-t_1)\right]$$

$$f(G_{3a},t_0,t_1) = \frac{2}{\theta_{12}}e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0} \times$$
$$\sum_{j\in S_2}\left[\frac{2}{\theta_1}P_{s,111}(t_1)P_{1_{bc}1_a,j}(\tau_1-t_1)+\frac{2}{\theta_1}P_{s,211}(t_1)P_{1_{bc}2_a,j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,122}(t_1)P_{1_a2_{bc},j}(\tau_1-t_1)+\frac{2}{\theta_2}P_{s,222}(t_1)P_{2_{bc}2_a,j}(\tau_1-t_1)\right]\right].$$

$$(16)$$

For gene trees $G_4$, $G_5$, and $G_6$, the probability density does not depend on the sequence IDs.

$$f(G_{4k},t_0,t_1) = \frac{2}{\theta_5}e^{-\frac{6}{\theta_5}t_1}\frac{2}{\theta_5}e^{-\frac{2}{\theta_5}t_0} \times \sum_{j\in S_3}P_{s,j}(\tau_1), \qquad 0<t_1+t_0<\tau_0-\tau_1,$$

$$f(G_{5k},t_0,t_1) = \frac{2}{\theta_5}e^{-\frac{6}{\theta_5}t_1}e^{-\frac{2}{\theta_5}(\tau_0-\tau_1-t_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0} \times \sum_{j\in S_3}P_{s,j}(\tau_1), \qquad 0<t_1<\tau_0-\tau_1, 0<t_0<\infty, \qquad (17)$$

$$f(G_{6k},t_0,t_1) = e^{-\frac{6}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{6}{\theta_4}t_1}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0} \times \sum_{j\in S_3}P_{s,j}(\tau_1), \qquad 0<t_1,t_0<\infty,$$

where $k = c$, $a$, and $b$.

### *Case III: Initial states 113, 123, and 223*

For initial state $s = 113$, $123$, or $223$, only three gene tree shapes are possible: $G_3$, $G_5$, and $G_6$ (Table

S1). For tree shapes $G_3$ and $G_5$, the only gene tree possible is $G_{3c}$ or $G_{5c}$: $((a, b), c)$, while for the tree

shape $G_6$, the three gene trees $G_{6c}$: $((a, b), c)$; $G_{6a}$: $((b, c), a)$; and $G_{6b}$: $((c, a), b)$ have the same prior

density. Thus there is no need to trace the sequence IDs. There are four states in the chain: 113, 123,

223, 13|23, with the rate matrix given as follows.

|       | 113 | 123 | 223 | 13\|23 |
|-------|-----|-----|-----|--------|
| 113   | $-(2w_{21}+c_1)$ | $2w_{21}$ | $0$ | $c_1$ |
| 123   | $w_{12}$ | $-(w_{12}+w_{21})$ | $w_{21}$ | $0$ |
| 223   | $0$ | $2w_{12}$ | $-(2w_{12}+c_2)$ | $c_2$ |
| 13\|23 | $0$ | $0$ | $0$ | $0$ |

$$(18)$$

For tree shapes $G_3$ and $G_5$, only one gene tree is possible, so that

$$f(G_{3c},t_0,t_1) = \frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0}\times\left[\frac{2}{\theta_1}P_{s,113}(t_1)+\frac{2}{\theta_2}P_{s,223}(t_1)\right],$$

$$f(G_{5c},t_0,t_1) = \frac{2}{\theta_5}\frac{2}{\theta_4}e^{-\frac{2}{\theta_5}t_1}e^{-\frac{2}{\theta_4}t_0}\times\sum_{j\in S_3}P_{s,j}(\tau_1).$$

$$(19)$$

For tree shape $G_6$, the three gene trees have the same density.

$$f(G_{6k}, t_0, t_1) = \frac{2}{\theta_4}e^{-\frac{6}{\theta_4}t_1}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}t_0}\times\sum_{j\in S_3}P_{s,j}(\tau_1)e^{-\frac{2}{\theta_5}\tau_0}, \qquad (20)$$

986    where $k = c$, $a$, and $b$.

### *Case IV: Initial states 133, 233, and 333*

988    For initial state $s = 133$, 233, or 333, there is no need to trace the sequence IDs.  We first discuss the

989    initial state 333.  The genealogical process is the single-population coalescent, with different

990    population size parameters: $\theta_3$ for $t < \tau_0$ or $\theta_4$ for $t > \tau_0$.  There is no need to distinguish among $G_1$, $G_2$,

991    and $G_4$, or between $G_3$ and $G_5$, so we consider only $G_1$ and $G_3$, but with the range of the coalescent

992    times modified accordingly.  There are thus three tree shapes: $G_1$, $G_3$, and $G_6$.  For each one, we sum

993    over three gene trees.  Thus with initial state $s = 333$, we have

994
$$f(G_k, t_0, t_1) = \begin{cases} \frac{2}{\theta_3}\frac{2}{\theta_3}e^{-\frac{6}{\theta_3}t_1}e^{-\frac{2}{\theta_3}t_0}, & 0 < t_1 + t_0 < \tau_0, \text{ for } k = 1c, 1a, 1b, \\[2mm] \frac{2}{\theta_3}\frac{2}{\theta_4}e^{-\frac{6}{\theta_3}t_1}e^{-\frac{2}{\theta_3}(\tau_0 - t_1)}e^{-\frac{2}{\theta_4}t_0}, & t_1 < \tau_0, \quad\quad \text{ for } k = 3c, 3a, 3b, \\[2mm] \frac{2}{\theta_4}\frac{2}{\theta_4}e^{-\frac{6}{\theta_4}\tau_0}e^{-\frac{6}{\theta_4}t_1}e^{-\frac{2}{\theta_4}t_0}, & 0 < t_1, t_0 < \infty, \text{ for } k = 6c, 6a, 6b. \end{cases} \tag{21}$$

995    Similarly, for initial state $s = 133$ or 233, we consider two tree shapes $G_3$ and $G_6$.

996
$$f(G_k, t_0, t_1) = \begin{cases} \frac{2}{\theta_3}\frac{2}{\theta_4}e^{-\frac{2}{\theta_3}t_1}e^{-\frac{2}{\theta_4}t_0}, & t_1 < \tau_0, \quad\quad \text{ for } k = 3, \\[2mm] \frac{2}{\theta_4}\frac{2}{\theta_4}e^{-\frac{2}{\theta_3}\tau_0}e^{-\frac{6}{\theta_4}t_1}e^{-\frac{2}{\theta_4}t_0}, & 0 < t_1, t_0 < \infty, \text{ for } k = 6c, 6a, 6b. \end{cases} \tag{22}$$

997

998

**FIGURE LEGENDS**

1000

1001   FIGURE 1.  (**a**) Species tree illustrating parameters in model M2 (gene flow) for three species (1, 2,

1002   and 3) and (**b**)-(**g**) possible gene tree shapes for a locus with three sequences ($a$, $b$, and $c$).  With

1003   certain initial states (data configurations at the locus), we have to keep track of the sequence IDs ($a$, $b$,

1004   and $c$) as well as the population IDs, so that each gene tree shape may correspond to three distinct

1005   gene trees.  For example, with the data configuration (initial state) $1_a2_b3_c$, the tree shape $G_6$ represents

1006   three distinct gene trees: $G_{6c}$: $((a, b), c)$; $G_{6a}$: $((b, c), a)$; and $G_{6b}$: $((c, a), b)$.

1007

1008

1009   FIGURE 2.  The three gene trees with branch lengths for three sequences $a$, $b$, and $c$.  Branch lengths $b_0$

1010   and $b_1$ are simple linear functions of coalescent times $t_0$ and $t_1$ in the gene trees of Fig. 1.  For

1011   example, for the tree $G_1$ of Fig. 1, $b_0 = t_0$ and $b_1 = t_1$, while for $G_2$, $b_0 = t_0 + \tau_1 - t_1$ and $b_1 = t_1$.

1012

1013

1014   FIGURE 3.  Posterior probabilities of the six possible gene trees ($G_{3c}$, $G_{5c}$, $G_{6c}$, $G_{6a}$, and $G_{6b}$) for the

1015   '123' loci in a dataset simulated using the MLEs of parameters for the *Drosophila* dataset D1 (auto).

1016

1017

1018   FIGURE 4.  Posterior probabilities of gene trees for the MSY loci for dataset D5 (exons split).  The red

1019   lines for gene tree $G_{3c}$ indicated loci that are likely to have been transferred across species, with

1020   $P(G_{3c}) > 95\%$.

1021

1022

1023

1024

**TABLE 1. Summary of the density for coalescent time for two sequences under M0 (no gene flow)**

| State | $f(t)$ before transform | $t$ limits | $f(x)$ after transform | $x$ limits | $b$ |
|---|---|---|---|---|---|
| 11 | $\frac{2}{\theta_1}e^{-\frac{2}{\theta_1}t}$ | $(0, \tau_1)$ | $e^{-x}$ | $(0,\frac{2}{\theta_1}\tau_1)$ | $\frac{\theta_1}{2}x$ |
| | $e^{-\frac{2}{\theta_1}\tau_1}\frac{2}{\theta_5}e^{-\frac{2}{\theta_5}(t-\tau_1)}$ | $(\tau_1, \tau_0)$ | $e^{-\frac{2}{\theta_1}\tau_1}e^{-x}$ | $(0,\frac{2}{\theta_5}(\tau_0-\tau_1))$ | $\tau_1+\frac{\theta_5}{2}x$ |
| | $e^{-\frac{2}{\theta_1}\tau_1}e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}(t-\tau_0)}$ | $(\tau_0, \infty)$ | $e^{-\frac{2}{\theta_1}\tau_1}e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}e^{-x}$ | $(0, \infty)$ | $\tau_0+\frac{\theta_4}{2}x$ |
| 22 | As for 11 above, with $\theta_1$ replaced by $\theta_2$ | | | | |
| 12 | $\frac{2}{\theta_5}e^{-\frac{2}{\theta_5}(t-\tau_1)}$ | $(\tau_1, \tau_0)$ | $e^{-x}$ | $(0,\frac{2}{\theta_5}(\tau_0-\tau_1))$ | $\tau_1+\frac{\theta_5}{2}x$ |
| | $e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}(t-\tau_0)}$ | $(\tau_0, \infty)$ | $e^{-\frac{2}{\theta_5}(\tau_0-\tau_1)}e^{-x}$ | $(0, \infty)$ | $\tau_0+\frac{\theta_4}{2}x$ |
| 13/23 | $\frac{2}{\theta_4}e^{-\frac{2}{\theta_4}(t-\tau_0)}$ | $(\tau_0, \infty)$ | $e^{-x}$ | $(0, \infty)$ | $\tau_0+\frac{\theta_4}{2}x$ |
| 33 | $\frac{2}{\theta_3}e^{-\frac{2}{\theta_3}t}$ | $(0, \tau_0)$ | $e^{-x}$ | $(0,\frac{2}{\theta_3}\tau_0)$ | $\frac{\theta_3}{2}x$ |
| | $\frac{2}{\theta_4}e^{-\frac{2}{\theta_3}\tau_0}e^{-\frac{2}{\theta_4}(t-\tau_0)}$ | $(\tau_0, \infty)$ | $e^{-x}e^{-\frac{2}{\theta_3}\tau_0}$ | $(0, \infty)$ | $\tau_0+\frac{\theta_4}{2}x$ |

**TABLE 2. Markov chains and their states for characterizing the genealogical process of epoch $E_1$ in model M2 (gene flow)**

| Case | Initial states | States in chain | Calculation of $P(t)$ |
|---|---|---|---|
| | Loci with 3 sequences | | |
| I | {111, 222} | {111, 112, 122, 222, 11, 12, 22, 1\|2} | Numerical |
| | | 8 states | |
| II | {112, 122} | {111, 112, 121, 122, 211, 212, 221, 222, $1_{bc}1_a$, $1_{ca}1_b$, $1_{ab}1_c$, $1_{bc}2_a$, $1_{ca}2_b$, $1_{ab}2_c$, $1_a2_{bc}$, $1_b2_{ca}$, $1_c2_{ab}$, $2_{bc}2_a$, $2_{ca}2_b$, $2_{ab}2_c$, 1\|2} | Numerical |
| | | 21 states | |
| III | {113, 123, 223} | {113, 123, 223, 13\|23} | Numerical |
| IV | {133, 233, 333} | {133, 233, 13, 23, 33, 3} | Analytical |
| | | | |
| | Loci with 2 sequences | | |
| V | {11, 12, 22} | {11, 12, 22, 1\|2} | Numerical |
| VI | {13, 23, 33} | {13, 23, 33, 3} | Analytical |

Note.— In case II (with initial states 112 or 122), it is necessary to keep track of both the population ID (1, 2, 3) and the sequence ID ($a$, $b$, $c$), so that state $1_{ab}2_c$ means two lineages in the sample, with the common ancestor of $a$ and $b$ in population 1, and sequence $c$ in population 2.

**TABLE 3. Rate matrix $Q$ for the Markov chain for initial states 111 and 222 under model M2**

| | 111 | 112 | 122 | 222 | 11 | 12 | 22 | 1\|2 |
|---|---|---|---|---|---|---|---|---|
| 111 | · | $3\times 4M_{21}/\theta_1$ | | | $3\times 2/\theta_1$ | | | |
| 112 | $4M_{12}/\theta_2$ | · | $2\times 4M_{21}/\theta_1$ | | | $2/\theta_2$ | | |
| 122 | | $2\times 4M_{12}/\theta_2$ | · | $4M_{21}/\theta_1$ | | $2/\theta_1$ | | |
| 222 | | | $3\times 4M_{12}/\theta_2$ | · | | | $3\times 2/\theta_2$ | |
| 11 | | | | | · | $2\times 4M_{21}/\theta_1$ | | $2/\theta_1$ |
| 12 | | | | | $4M_{12}/\theta_2$ | · | $4M_{21}/\theta_1$ | |
| 22 | | | | | | $2\times 4M_{12}/\theta_2$ | · | $2/\theta_2$ |
| 1\|2 | | | | | | | | · |

Note.— We define parameters using the real-world process (with time running forward), so that the migration rate $M_{ij} = N_j m_{ij}$ is the expected number of migrant individuals from populations $i$ to $j$ per generation (in the real world) and $m_{ij}$ is the proportion of individuals in population $j$ that are immigrants from population $i$. The Markov chain is then used to describe the process of coalescent with migration, with time running backwards. For example $Q_{111,\,112}$ is the rate for the transition from state 111 to state 112, which in the real world means one of the three sequences in population 1 is an immigrant from population 2, which has the rate $3m_{21}$ per generation. Since time is measured by the mutational distance and one time unit is the expected time to accumulate one mutation per site (that is, one time unit is $1/\mu$ generations), the rate per time unit is $Q_{111,\,112} = 3m_{21}\times 1/\mu = 3\times 4N_1 m_{21}/(4N_1\mu)$ $=3\times 4M_{21}/\theta_1$, as in the table. Given the rate matrix $Q = \{Q_{ij}\}$, the transition probability matrix over time $t$ is given as $P(t) = \{P_{ij}(t)\} = e^{Qt}$. This is the same calculation as in the Markov chain models for nucleotide substitution such as Jukes and Cantor (Jukes and Cantor, 1969).

**TABLE 4. Five *Drosophila* datasets analyzed in this paper**

| Dataset | #MMY loci | #MSY loci | #SSY loci | Total |
|---|---|---|---|---|
| D1 auto | 378 | 19,224 | 9,425 | 29,027 |
| D2 noncoding | 378 | 14,498 | 7,211 | 22,087 |
| D3 chrX | 0 | 4,381 | 2,105 | 6,486 |
| D4 exons complete | 378 | 27,200 | 13,500 | 41,078 |
| D5 exons split | 378 | 10,979 | 5,342 | 16,699 |

**TABLE 5. False positive rate, percentage of zeros, and 95% quantile of the null distribution of the LRT statistic ($2\Delta\ell$) comparing the symmetrical versions of models M0 (no gene flow) and M2 (gene flow)**

| Data | $L=10$ | 100 | 1000 | 15,000 |
|---|---|---|---|---|
| Set 1 (hominoid): $\theta_4 = \theta_5 = \theta_{12} = 0.005$, $\tau_0 = 0.006$, $\tau_1 = 0.004$ | | | | |
| (a) 123 | 0.000 0.829 0.034 | 0.001 0.641 2.217 | 0.005 0.528 2.708 | 0.004 0.506 2.443 |
| (b) 11&12 | 0.003 0.851 0.578 | 0.019 0.680 1.528 | 0.045 0.504 2.542 | 0.084 0.479 3.492 |
| (c) 113&123 | 0.002 0.848 0.307 | 0.027 0.674 2.073 | 0.037 0.576 2.161 | 0.035 0.507 2.329 |
| Set 2 (mangroves): $\theta_4 = \theta_5 = \theta_{12} = 0.01$, $\tau_0 = 0.02$, $\tau_1 = 0.01$ | | | | |
| (a) 123 | 0.001 0.883 0.616 | 0.006 0.798 1.330 | 0.009 0.709 2.060 | 0.004 0.345 1.772 |
| (b) 11&12 | 0.009 0.881 0.454 | 0.020 0.741 1.542 | 0.100 0.439 3.872 | 0.078 0.570 3.481 |
| (c) 113&123 | 0.010 0.906 0.418 | 0.035 0.791 1.983 | 0.031 0.712 2.013 | 0.039 0.722 2.136 |
| Set 3: $\theta_4 = \theta_{12} = 0.02$, $\theta_5 = 0.03$, $\tau_0 = 0.06$, $\tau_1 = 0.04$ | | | | |
| (a) 123 | 0.000 0.957 0.000 | 0.002 0.904 0.501 | 0.001 0.896 0.424 | 0.006 0.884 0.975 |
| (b) 11&12 | 0.007 0.864 0.796 | 0.032 0.727 1.979 | 0.035 0.713 1.814 | 0.009 0.839 0.422 |
| (c) 113&123 | 0.003 0.945 0.017 | 0.008 0.902 0.535 | 0.007 0.895 0.589 | 0.008 0.910 0.198 |
| Set 4: $\theta_4 = \theta_{12} = 0.02$, $\theta_5 = 0.01$, $\tau_0 = 0.02$, $\tau_1 = 0.01$ | | | | |
| (a) 123 | 0.000 0.854 1.137 | 0.003 0.782 1.469 | 0.001 0.717 0.841 | 0.002 0.685 2.003 |
| (b) 11&12 | 0.008 0.823 0.479 | 0.032 0.757 1.707 | 0.047 0.625 2.470 | 0.049 0.656 2.687 |
| (c) 113&123 | 0.013 0.823 1.056 | 0.040 0.775 2.069 | 0.034 0.719 1.782 | 0.030 0.666 2.136 |

Note.— In each cell, the three numbers are the false positive rate, the proportion of replicates in which the test statistic is $2\Delta\ell = 0$, and the estimated 95% critical value. The critical value used for the test is $\chi^2_{2,5\%} = 5.99$ for (a) configuration 123, and is 2.71 for (b) 11&12 and (c) 113&123.


**TABLE 6. Power of the LRT comparing the symmetrical versions of models M0 (no gene flow) and M2 (gene flow)**

| Data | $L=10$ | 100 | 1000 | 15,000 |
|---|---|---|---|---|
| Set 1 (hominoid): $\theta_4 = \theta_5 = \theta_{12} = 0.005$, $\tau_0 = 0.006$, $\tau_1 = 0.004$, $M=1$ | | | | |
| (a) 123 | 0.6% | 5.3% | 81.6% | 100% |
| (b) 11&12 | 4.6% | 7.0% | 16.1% | 65.7% |
| (c) 113&123 | 3.3 % | 17.9% | 88.3% | 100% |
| Set 2 (mangroves): $\theta_4 = \theta_5 = \theta_{12} = 0.01$, $\tau_0 = 0.02$, $\tau_1 = 0.01$, $M=1$ | | | | |
| (a) 123 | 3.0% | 52.1% | 100% | 100% |
| (b) 11&12 | 8.0% | 27.3% | 32.0% | 89.3% |
| (c) 113&123 | 13.8% | 69.3% | 100% | 100% |

Note.— The critical value used is 5.99 for (a) 123, and is 2.71 for (b) 11&12 and (c) 113&123.

**Table 7. Means and SDs of MLEs from datasets simulated under the symmetrical model M2 (gene flow)**

| Data | (a) 11&12 | | | | | | (b) 113&123 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\theta_4$ | $\theta_5$ | $\tau_0$ | $\tau_1$ | $\theta_{12}$ | $M$ | $\theta_4$ | $\theta_5$ | $\tau_0$ | $\tau_1$ | $\theta_{12}$ | $M$ |
| Set 1 (hominoid): $\theta_4 = \theta_5 = \theta_{12} = 0.005$, $\tau_0 = 0.006$, $\tau_1 = 0.004$, $M = 1$ | | | | | | | | | | | | |
| Truth | 5 | 5 | 6 | 4 | 5 | 1 | 5 | 5 | 6 | 4 | 5 | 1 |
| $L = 100$ | $6.7 \pm 4.1$ | $33.7 \pm 191.0$ | $6.7 \pm 3.0$ | $3.4 \pm 2.3$ | $9.3 \pm 64.0$ | $1.4 \pm 1.7$ | $4.9 \pm 1.0$ | $10.8 \pm 90.2$ | $6.0 \pm 0.4$ | $3.6 \pm 1.9$ | $6.6 \pm 8.1$ | $1.3 \pm 1.4$ |
| $L = 1000$ | $5.5 \pm 2.5$ | $20.0 \pm 152.5$ | $7.4 \pm 3.6$ | $3.4 \pm 1.9$ | $6.9 \pm 56.9$ | $1.1 \pm 0.7$ | $5.0 \pm 0.3$ | $4.7 \pm 2.0$ | $6.0 \pm 0.1$ | $4.0 \pm 1.2$ | $5.1 \pm 0.6$ | $1.1 \pm 0.6$ |
| $L = 15000$ | $5.1 \pm 1.0$ | $14.1 \pm 98.3$ | $7.4 \pm 4.1$ | $3.5 \pm 1.3$ | $5.0 \pm 0.4$ | $0.9 \pm 0.2$ | $5.0 \pm 0.1$ | $5.0 \pm 0.6$ | $6.0 \pm 0.0$ | $4.0 \pm 0.3$ | $5.0 \pm 0.1$ | $1.0 \pm 0.1$ |
| Set 2 (mangroves): $\theta_4 = \theta_5 = \theta_{12} = 0.01$, $\tau_0 = 0.02$, $\tau_1 = 0.01$, $M = 1$ | | | | | | | | | | | | |
| Truth | 10 | 10 | 20 | 10 | 10 | 1 | 10 | 10 | 20 | 10 | 10 | 1 |
| $L = 100$ | $13.1 \pm 7.5$ | $17.8 \pm 87.2$ | $18.6 \pm 7.5$ | $8.8 \pm 5.0$ | $10.9 \pm 7.3$ | $1.5 \pm 1.7$ | $9.9 \pm 1.9$ | $9.6 \pm 3.9$ | $20.1 \pm 0.9$ | $9.9 \pm 4.2$ | $14.0 \pm 70.0$ | $1.4 \pm 1.4$ |
| $L = 1000$ | $10.9 \pm 4.3$ | $13.4 \pm 64.5$ | $18.6 \pm 7.7$ | $8.6 \pm 4.0$ | $10.0 \pm 1.7$ | $1.1 \pm 0.5$ | $10.0 \pm 0.6$ | $9.9 \pm 1.2$ | $20.0 \pm 0.3$ | $10.0 \pm 0.2$ | $10.1 \pm 0.6$ | $1.1 \pm 0.4$ |
| $L = 15000$ | $10.1 \pm 2.2$ | $16.9 \pm 103.4$ | $20.8 \pm 7.8$ | $9.5 \pm 2.0$ | $10.0 \pm 0.2$ | $1.0 \pm 0.2$ | $10.0 \pm 0.2$ | $10.0 \pm 0.3$ | $20.0 \pm 0.1$ | $10.0 \pm 0.3$ | $10.0 \pm 0.1$ | $1.0 \pm 0.1$ |

Note.— Estimates of $\theta$s and $\tau$s are multiplied by 1000. For $L = 100$ or 1000, some estimates are very large ($\infty$) in certain datasets, causing the mean and SD to be very large. See table 5 for the power of the LRT from the same data.

TABLE **8**. Means and SDs of MLEs from datasets simulated under the asymmetrical IM model M2 (gene flow)

| Data | Parameters (true values in parentheses) | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\theta_4$ (10) | $\theta_5$ (10) | $\tau_0$ (20) | $\tau_1$ (10) | $\theta_1$ (5) | $\theta_2$ (10) | $M_{12}$ (0.1) | $M_{21}$ (1) |
| (a) 223&123 | | | | | | | | |
| $L = 100$ | $9.9 \pm 2.0$ | $16.8 \pm 63.1$ | $20.1 \pm 0.9$ | $10.4 \pm 5.0$ | $9.7 \pm 19.3$ | $9.4 \pm 5.9$ | $0.2 \pm 0.5$ | $1.2 \pm 0.8$ |
| $L = 1000$ | $10.0 \pm 0.6$ | $12.6 \pm 38.9$ | $20.0 \pm 0.3$ | $10.0 \pm 4.9$ | $9.5 \pm 22.0$ | $9.6 \pm 1.6$ | $0.2 \pm 0.2$ | $1.6 \pm 2.6$ |
| $L = 15000$ | $10.0 \pm 0.2$ | $9.7 \pm 1.2$ | $20.0 \pm 0.1$ | $10.3 \pm 2.9$ | $5.4 \pm 3.5$ | $10.0 \pm 0.4$ | $0.1 \pm 0.0$ | $1.1 \pm 0.7$ |
| (b) 113&223&123 | | | | | | | | |
| $L = 99$ | $9.8 \pm 2.0$ | $10.9 \pm 26.9$ | $20.1 \pm 1.0$ | $10.2 \pm 5.0$ | $7.5 \pm 5.8$ | $9.3 \pm 6.1$ | $0.4 \pm 1.0$ | $1.4 \pm 1.5$ |
| $L = 999$ | $10.0 \pm 0.6$ | $11.8 \pm 37.6$ | $20.0 \pm 0.3$ | $10.1 \pm 4.7$ | $5.4 \pm 1.3$ | $9.5 \pm 2.1$ | $0.2 \pm 0.2$ | $1.0 \pm 0.3$ |
| $L = 15000$ | $10.0 \pm 0.1$ | $9.7 \pm 1.3$ | $20.0 \pm 0.1$ | $10.1 \pm 2.8$ | $5.0 \pm 0.3$ | $9.9 \pm 0.5$ | $0.1 \pm 0.1$ | $1.0 \pm 0.1$ |

Note.— Estimates of $\theta$s and $\tau$s are multiplied by 1000. For $L \leq 1000$, several datasets produced large estimates of $\theta_5$ at the upper bound set by the program. The means and SDs were calculated by excluding those estimates.

**TABLE 9. MLEs and standard errors from the five *Drosophila* datasets of Table 4**

| Data & model | $\tau_{MSY}$ | $\tau_{MS}$ | $\theta_{MSY}$ | $\theta_{MS}$ | $\theta_M$ | $\theta_S$ | $M_{MS}$ | $M_{SM}$ | $\ell$ | $2\Delta\ell$ |
|---|---|---|---|---|---|---|---|---|---|---|
| D1 auto | | | | | | | | | | |
| M0 | 24.6±0.1 | 11.3±0.1 | 39.4±0.3 | 13.3±0.2 | 6.0±0.4 | 12.8±0.2 | | | −4,763,806.0 | |
| M2 | 24.3±0.1 | 13.6±0.2 | 40.0±0.3 | 10.6±0.3 | 5.2±0.6 | 12.7±0.2 | 0.0 | 18.3±3.1 | −4,763,452.5 | 707.0 |
| D2 noncoding | | | | | | | | | | |
| M0 | 24.5±0.1 | 10.8±0.1 | 41.6±0.4 | 13.9±0.2 | 6.0±0.4 | 13.1±0.2 | | | −3,326,330.8 | |
| M2 | 24.3±0.1 | 12.6±0.2 | 42.1±0.4 | 12.0±0.2 | 5.3±0.4 | 13.0±0.2 | 0.0 | 16.2±2.5 | −3,326,145.1 | 371.2 |
| D3 chrX | | | | | | | | | | |
| M0 | 28.0±0.2 | 12.3±0.2 | 41.1±0.6 | 15.3±0.4 | NA | 8.2±0.2 | | | −1,027,233.4 | |
| M2 | 27.8±0.2 | 14.2±0.3 | 41.6±0.6 | 13.0±0.5 | 20.9±9.4 | 8.3±0.2 | 0.0 | 40.2±16.9 | −1,027,161.6 | 143.5 |
| M2 ($\theta_M = \theta_S/2$) | 27.8±0.2 | 14.2±0.3 | 41.6±0.6 | 13.0±0.5 | 4.1±NA | 8.3±0.2 | 0.0 | 8.0±1.1 | −1,027,161.7 | 143.5 |
| M2 ($\theta_M = \theta_S$) | 27.8±0.2 | 14.2±0.3 | 41.6±0.6 | 13.0±0.5 | 8.3±0.2 | | 0.0 | 15.9±NA | −1,027,161.7 | 143.5 |
| D4 exons complete | | | | | | | | | | |
| M0 | 20.2±0.1 | 10.9±0.1 | 33.7±0.2 | 9.9±0.1 | 5.9±0.4 | 10.7±0.1 | | | −7,853,901.6 | |
| M2 | 18.3±0.1 | 18.3±0.1 | 38.2±0.2 | 0.0±0.0 | 4.5±0.5 | 10.7±0.1 | 0.0 | 43.6±4.0 | −7,853,313.7 | 1175.8 |
| M2 ($\tau_{MSY} = 0.020$, $\tau_{MS} = 0.013$) | 20 | 13 | 34.3±0.2 | 7.4±0.0 | 5.1±NA | 10.6±0.1 | 0.0 | 20.7±NA | −7,853,425.1 | 952.9 |
| D5 exons split (subset of D4) | | | | | | | | | | |
| M0 | 19.6±0.1 | 10.9±0.1 | 38.9±0.3 | 9.4±0.2 | 5.9±0.4 | 10.2±0.2 | | | −2,139,639.5 | |
| M2 | 18.0±0.1 | 18.0±0.1 | 42.6±0.4 | 0.0±0.0 | 4.2±0.3 | 10.2±0.2 | 0.0 | 37.8±2.9 | −2,139,182.0 | 915.1 |
| M2 ($\tau_{MSY} = 0.020$, $\tau_{MS} = 0.013$) | 20 | 13 | 38.5±0.3 | 7.4±0.4 | 4.7±0.4 | 10.1±0.2 | 0.0 | 20.4±3.3 | −2,139,414.4 | 450.2 |

Note.— Estimates of $\tau$, $\theta$, and $M$ are multiplied by 1000. See Table 4 for information about the datasets.

**TABLE 10. MLEs and log likelihood values under M2 assuming different species trees for dataset D1 (auto) of Table 4**

| Species tree | $\tau_{MSY}$ | $\tau_1$ | $\theta_{MSY}$ | $\theta_5$ | $\theta_M$ | $\theta_S$ | $\theta_Y$ | $M_{12}$ | $M_{21}$ | $\ell$ |
|---|---|---|---|---|---|---|---|---|---|---|
| ((MS)Y) | 24.3±0.1 | 13.6±0.2 ($\tau_{MS}$) | 40.0±0.3 | 10.6±0.3 ($\theta_{MS}$) | 5.2±0.6 | 12.7±0.2 | NA | 0.0 ($M_{MS}$) | 18.3±3.1 ($M_{SM}$) | −4,763,452.5 |
| ((MY)S) | 10.7±0.1 | 10.7±1.0 ($\tau_{MY}$) | 53.5±0.3 | ∞ ($\theta_{MY}$) | 5.7±0.4 | ∞ | 8.2±0.1 | 0.0 ($M_{MY}$) | 0.0 ($M_{YM}$) | −4,780,884.0 |
| ((SY)M) | 11.4±0.1 | 11.4±0.1 ($\tau_{SY}$) | 52.8±0.3 | ∞ ($\theta_{SY}$) | 11.3±0.1 | ∞ | 4.2±0.3 | 0.0 ($M_{SY}$) | 0.0 ($M_{YS}$) | −4,783,156.2 |

Note.— Estimates of $\tau$, $\theta$, and $M$ are multiplied by 1000. Estimates of $\theta_5$ and $\theta_S$ hit the upper bound set in the program for trees ((MY)S) and ((SY)M).

$\tau_0$    $\theta_4$    $\theta_5$    $M_{12}$    $M_{21}$    $\tau_1$    $\theta_1$    $\theta_2$    $\theta_3$

$E_3$    $E_2$    $E_1$

$t_0$    $t_1$

1    2    3

(**a**) Species tree    (**b**) $G_1$    (**c**) $G_2$    (**d**) $G_3$    (**e**) $G_4$    (**f**) $G_5$    (**g**) $G_6$

**(a)** $T_c$    **(b)** $T_a$    **(c)** $T_b$