

Expected patterns of local ancestry in a hybrid zone

Joel Smith^{1*}, Bret Payseur², John Novembre^{1,3}

1 Department of Ecology and Evolution, University of Chicago, Chicago, IL

2 Laboratory of Genetics, University of Wisconsin-Madison, Madison, WI

3 Department of Human Genetics, University of Chicago, Chicago, IL

*joelsmith@uchicago.edu

1 **1 Abstract**

2 The initial drivers of reproductive isolation between species are poorly characterized. In cases where
3 partial reproductive isolation exists, genomic patterns of variation in hybrid zones may provide clues
4 about the barriers to gene flow which arose first during the early stages of speciation. Purifying
5 selection against incompatible substitutions that reduce hybrid fitness has the potential to distort
6 local patterns of ancestry relative to background patterns across the genome. The magnitude and
7 qualitative properties of this pattern are dependent on several factors including migration history and
8 the relative fitnesses for different combinations of incompatible alleles. We present a model which
9 may account for these factors and highlight the potential for its use in verifying the action of natural
10 selection on candidate loci implicated in reducing hybrid fitness.

11 **2 Introduction**

12 A large fraction of research aiming to describe the process of speciation involves mapping genetic
13 variants responsible for reproductive isolation. Despite its difficulty, this task has nevertheless been
14 carried out for a number of cases in which the link between a reproductive isolating mechanism mapped
15 in a laboratory setting and its effect on an individual's fitness in nature is demonstrated [Schluter,
16 2009]. However, in many of these cases, reproductive isolation is already complete such that the initial
17 cause of speciation cannot be attributed to any one locus or set of loci due to a lack of information
18 regarding the order in which these isolating barriers arose [Turelli et al., 2014]. Hybrid zones present
19 a convenient situation where reproductive isolation is incomplete. In these cases, the mechanisms of
20 reproductive isolation are both fewer and more recently derived. Relative to scenarios with complete
21 reproductive isolation, systems with ongoing hybridization may provide a more narrow set of candidate
22 loci to consider as the initial drivers of speciation.

23 The next task would be to describe the mechanism by which the incompatible substitutions were
24 fixed. Functional annotations for the implicated loci can yield some clues about the ecological context
25 or genetic causes that resulted in these substitutions. A rigorously tested explanation would require
26 that field experiments be carried out to establish their effect on fitness in nature [Schemske, 2000,
27 Schemske and Bradshaw, 1999]. However, patterns of genomic variation can provide a complementary
28 source of evidence for the action of natural selection on genetic variants which are relevant to a
29 phenotype of interest [Tiffin and Ross-Ibarra, 2014]. The robustness of any given metric or model for
30 the signature of natural selection depends on well-conceived theory that describes both the conditions
31 under which the signature is detectable as well as any non-selective processes that can explain the
32 pattern. This observational approach has been a driver of both theoretical and empirical research which
33 aims to implicate loci responsible for genetic incompatibilities that decrease fitness among hybrids in
34 nature [Barton, 1979, Barton and Hewitt, 1985, Endler, 1973, White, 1968].

35 Hybrid zones are thought to present a useful situation where the interaction between gene flow and
36 natural selection can leave identifiable patterns associated with genetic incompatibilities in genomic
37 data [Harrison and Larson, 2016, Payseur, 2010, Payseur and Rieseberg, 2016]. Historically, most
38 work on this problem has relied on using differences in allele frequencies across the hybrid zone while
39 ignoring patterns of linkage disequilibrium among neighboring sites [Barton and Hewitt, 1985]. More
40 recently, increased access to sequencing technology has prompted the use of methods which can infer
41 local ancestry across the genomes of admixed individuals [Gompert and Buerkle, 2013]. In this regard,
42 population genetic inference has made a significant shift toward developing models which leverage this
43 information for a variety of purposes. Several models aim to infer the migration history between
44 genetically distinct populations using the length of ancestry tract lengths among admixed individuals
45 [Gravel, 2012, Harris and Nielsen, 2013, Hellenthal et al., 2014, Liang and Nielsen, 2014, Loh et al.,
46 2013, Patterson et al., 2012, Pool and Nielsen, 2009, Price et al., 2009, Sedghifar et al., 2015]. As
47 the primary intention of these approaches has been to focus on populations within a species, there
48 is a lack of work which aims to describe the effect of genetic incompatibilities which commonly arise
49 between species after a prolonged period of geographic isolation.

50 Theory with formal treatment of genetic incompatibilities and ancestry tracts has been slow to
51 accumulate, in large part due to the large parameter space of both migration histories and genetic
52 architectures that may contribute to reduced fitness in hybrid individuals. As a result, forward sim-
53 ulations of whole chromosomes under differing migration and selection regimes have been used to
54 describe some general patterns [Gompert et al., 2012, Hvala et al., 2018, Lindtke and Buerkle, 2015,
55 Schumer and Brandvain, 2016]. In a few of these cases, the primary goal is to describe the conditions

56 which may account for the heterogeneous patterns of genomic differentiation which have been widely
57 observed across hybrid zones [Harrison and Larson, 2016]. For example, Gompert et al. [2012] focus
58 on describing differences in both the number of contributing loci and the mechanism of their effect
59 through either underdominance at single loci or two-locus epistasis. They also introduce a formalized
60 approach to identify outlier loci responsible for reduced hybrid fitness using allele frequency clines
61 across the genome. Lindtke and Buerkle [2015] pay particular attention to two-locus models of genetic
62 incompatibilities and compare the relative efficiency with which different kinds of epistatic interactions
63 can maintain genomic differentiation in a hybrid zone under both high and low migration.

64 In an effort to make use of ancestry tract lengths rather than allele frequencies at individual loci,
65 Sedghifar et al. [2015] derive a null expectation for the length of ancestry tracts in a geographic
66 context where distance from the contact zone of two genetically distinct populations is explicitly
67 modeled. They then provide a likelihood function which they use to infer the age of the contact
68 zone, or time at which admixture between the populations began. Sedghifar et al. [2016] extends this
69 spatially-explicit framework further to model the mean ancestry tract length which is contiguous with
70 an under-dominant locus.

71 Another approach that uses local ancestry inference to identify genetic incompatibilities relies on
72 computing correlations in ancestry among pairs of loci in a hybrid zone [Schumer et al., 2014]. Schumer
73 and Brandvain [2016] use simulation to demonstrate how selection against incompatible alleles at two
74 loci can lead to a positive correlation in species ancestry at those loci. They find good power to
75 identify these associations for genetic architectures that feature ubiquitous selection (see Figure 1d).
76 The intuition for this pattern is that genotypes with the same ancestry at both loci are the only
77 genotypes with high fitness, such that an over-representation of ancestry at those loci relative to
78 background levels of linkage disequilibrium (LD) should lead to an identifiable signal. For genetic
79 architectures that only feature strong selection against derived allele combinations, they find much
80 less power to identify significant pairs.

81 The variety of approaches and data available to study this problem have prompted a few questions
82 of where to proceed next. We first describe a few of the well-studied genetic architectures for two-locus
83 genetic incompatibilities as well as others that have received less attention but which have also been
84 identified in nature. We then present a model to compute the expected distribution of ancestry tract
85 lengths around incompatibility loci.

86 2.1 Two-Locus Genetic Incompatibilities

87 The two-locus fitness matrix provides a useful representation of different genetic architectures which
 88 might contribute to genetic incompatibility between species (Figure 1 and Table 1). In addition to
 89 theoretical arguments and simulations, much of our current understanding for how relevant any of
 90 these genetic architectures might be in nature has been driven by genetic dissection of reproductive
 91 barrier phenotypes in the lab [White et al., 2011]. There are a number of empirical examples in
 92 a variety of species which have hinted at the potential importance of meiotic drive and neutral (or
 93 nearly neutral) causes for the fixation of incompatible substitutions [Maheshwari and Barbash, 2011,
 94 Presgraves, 2010, Sweigart and Willis, 2012]. While the precise combination of evolutionary forces
 95 which are responsible for incompatibility formation remain unknown, the evolution of incompatibili-
 96 ties in hybrid populations can be reasonably approximated with simple epistasis [Schumer et al., 2014].

97

		bb ($\mathbf{B}_1\mathbf{B}_1$)	Bb ($\mathbf{B}_1\mathbf{B}_2$)	BB ($\mathbf{B}_2\mathbf{B}_2$)
aa	($\mathbf{A}_1\mathbf{A}_1$)	1	$1 - s_a h_a$	$1 - s_a$
Aa	($\mathbf{A}_1\mathbf{A}_2$)	$1 - s_e h_1$	$1 - s_e h_0$	$1 - s_a h_a$
AA	($\mathbf{A}_2\mathbf{A}_2$)	$1 - s_e$	$1 - s_e h_1$	1

Table 1. Genotype fitnesses for the DMI and symmetric incompatibility models. The first pairs of bold letters are DMI model genotypes and the genotypes in parentheses indicate the symmetric model. s_a and s_e denote the selection coefficient against the ancestral and incompatible alleles, respectively. h_a , h_0 and h_1 denote the dominance effects of ancestral, double-heterozygotes and single-heterozygotes, respectively.

98 The most well-known model is described in Dobzhansky [1937] in which alleles fix at two interacting
 99 loci among populations that are geographically isolated. The top row in Figure 1 shows a range of
 100 possible fitness matrices that might result from this scenario, also known as the Dobzhansky-Muller
 101 incompatibility model (DMI). If we denote the ancestral genotype as **aaBB** in all of these cases,
 102 then the derived genotypes before coming into secondary contact are **aabb** and **AABB**. We chose
 103 these example matrices to emphasize the diversity of fitness configurations that might result from this
 104 model. The fitness matrix in Figure 1a is an example where the the derived substitutions were fixed
 105 by positive selection, such that the ancestral genotype suffers a fitness cost. Figures 1a and 1b are
 106 examples where the derived alleles interact dominantly; whereas in Figure 1c, derived alleles interact
 107 recessively.

108 Lindtke and Buerkle [2015] draw attention to a different model of genetic incompatibility where
 109 allele substitutions occur at two loci in both populations leading to a symmetric pattern of fitnesses

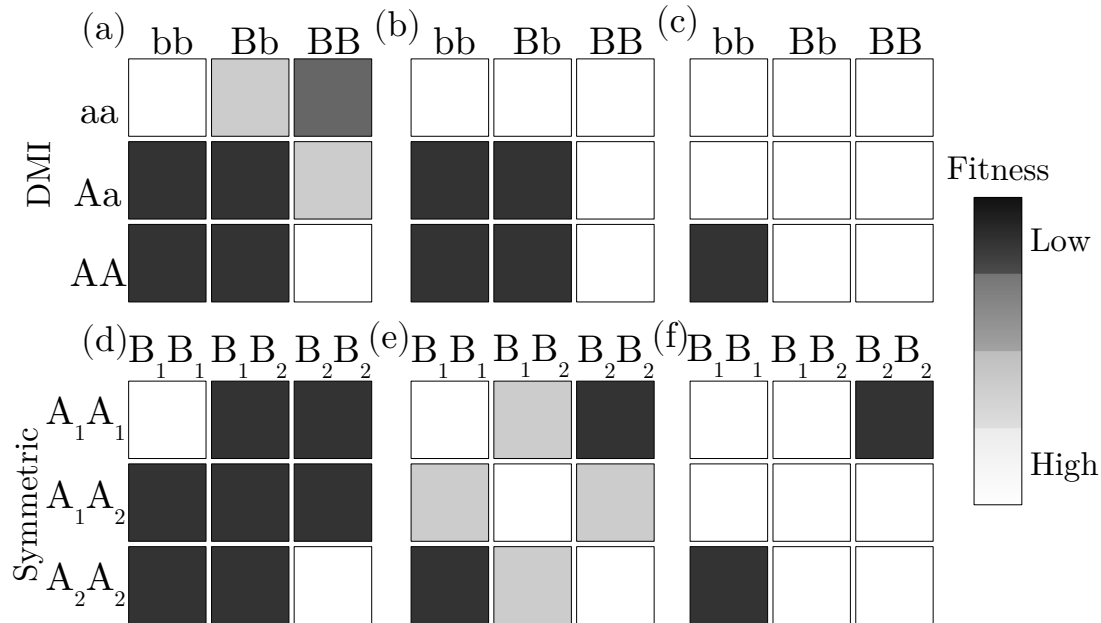


Fig 1. Two-locus fitness matrices for six models of genetic incompatibility. Each matrix includes the fitnesses of all possible two-locus genotypes where each locus is biallelic. Shaded boxes represent genotypes with a fitness cost that varies positively with the amount of shading. The top row of matrices are variations of the DMI model with the **aaBB** genotype representing the ancestral state and the bottom row shows variations of a symmetric incompatibility model. For both rows, the dominance effect of derived substitutions decreases from left to right.

110 between the two derived genotypes **A₁A₁B₁B₁** and **A₂A₂B₂B₂** (Figure 1d, 1e, 1f). Their results
 111 suggest that this mechanism could provide a better explanation for the observed patterns of genetic
 112 differentiation that occur at extended genomic distances between species that hybridize [Harrison
 113 and Larson, 2016]. Regulatory interactions between a transcription factor encoded at one locus and
 114 the corresponding binding site at a second locus would be one scenario consistent with this model.
 115 Seehausen et al. [2014] note that this model could also be common in meiotic drive scenarios where a
 116 substitution that promotes biased transmission of a selfish genetic element at one locus is counteracted
 117 by a substitution at a second locus which restores unbiased inheritance. The bottom row in Figure 1
 118 shows a range of possible fitness matrices under this model, where the left-most matrix results from
 119 dominant substitutions which interfere between haplotypes, and the right most matrix results from
 120 recessive substitutions. Simulated data in Figure 2 (using the software dfuse from Lindtke and Buerkle
 121 [2015]) illustrates the effect of the DMI model in Figure 1b where selection against derived alleles leads

122 to a bias towards the ancestral genotype (**aaBB**) of recombined ancestries.

123 In the following section, we first review an approach taken by Gravel [2012] to model the distri-
124 bution of ancestry tract lengths across the genomes of an admixed population. We then describe the
125 framework for our own extension to this approach which aims to model the distribution of ancestry
126 tract lengths that are contiguous with a locus undergoing epistatic interactions according to any of
127 the incompatibility scenarios outlined above.

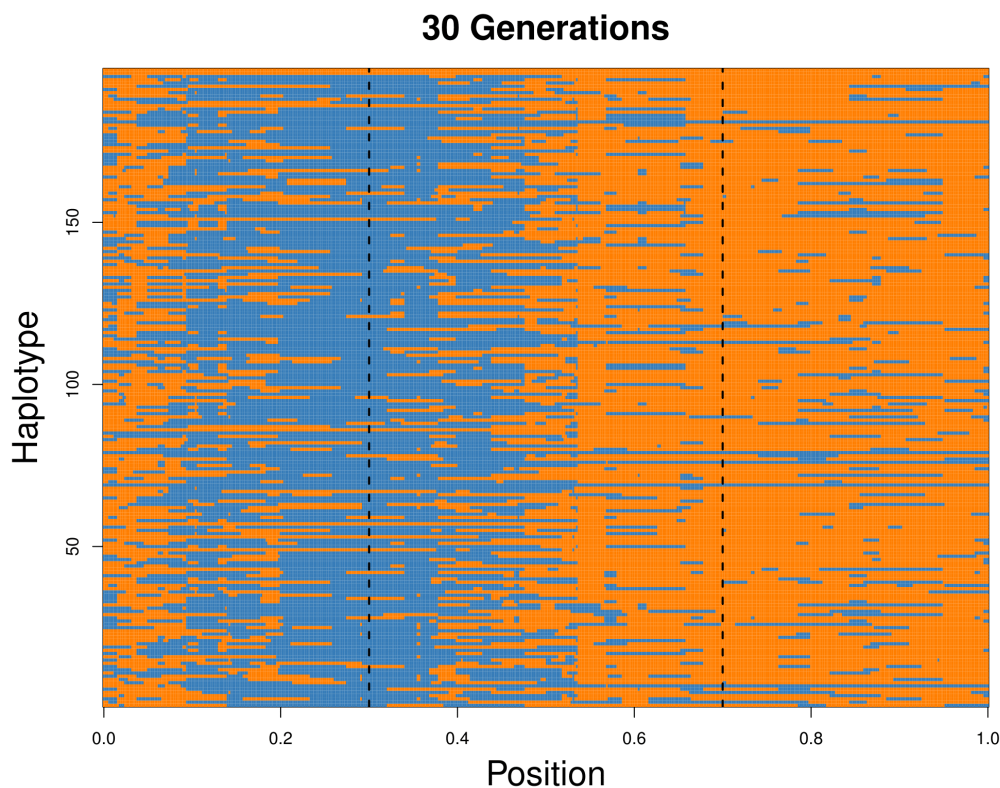


Fig 2. Haplotype data simulated using the software *dfuse* with the fitness matrix in Figure 1b. The forward-in-time simulation begins with two infinite source populations contributing equal fractions of ancestry (0.5) to a target population of 100 individuals 30 generations in the past. Each generation to the present follows a Wright-Fisher model, whereby both source populations contribute a fraction of individuals m to the target population. In this case $m = 0.1$. Recombination occurs uniformly along the chromosome at rate 1 crossover per chromosome per generation. After recombination, individuals are removed from the population according to a specified fitness matrix. The parameter values defined in Table 1 take the following values: $s_a = 0$, $s_e = 0.9$, $h_1 = 1$, $h_0 = 1$, and $h_a = 0$. The incompatibility loci are indicated by the vertical dotted lines.

128 3 Model Description

129 3.1 Tract Length Distributions Under Neutral Admixture

130 Gravel [2012] defines a Markov chain along a chromosome with transition rates between both an ances-
 131 try state variable, p , and the time, t , at which ancestry p arrives in a hybrid population. Consider the
 132 demography of a sample up to the first hybridization event T generations ago, where each generation
 133 is labeled $s \in \{0, 1, 2, \dots, T - 1\}$. Let $m_p(t)$ denote the fraction of individuals in the target population
 134 replaced by individuals from source population p at time t . $m(t)$ is the total fraction of individuals
 135 in the target population replaced by migrants in generation t where $\sum_p m_p(t) \leq 1$. Moving along a
 136 chromosome from any point, the probability of encountering state (p, t) after a recombination event
 137 that occurred at generation τ is

$$P(p, t | \tau) = m_p(t) \prod_{t'=\tau+1}^{t-1} (1 - m(t')). \quad (1)$$

138 τ is uniformly distributed on $(1, t - 1)$, so the discrete transition probabilities can be expressed as

$$R(p, t \rightarrow p', t') = \sum_{\tau=1}^{\min(t, t')-1} \frac{P(p', t' | \tau)}{(t - 1)} \quad (2)$$

139 To get the continuous transition rate, one can multiply the discrete transition rate by the continuous
 140 overall transition rate $t - 1$. This follows from the fact that a recombination event occurs at each gen-
 141 eration such that probability of observing an ancestry junction depends on the number of generations
 142 since admixture:

$$Q(p, t \rightarrow p', t') = m_{p'}(t') \sum_{\tau=1}^{\min(t, t')-1} \prod_{s=\tau+1}^{t'-1} (1 - m(s)). \quad (3)$$

143 Using Q , one can compute the tract length distribution for a given ancestry. Q is first uniformized
 144 to adjust self-transition probabilities such that the total transition rate from each state is equal to
 145 the rate of the state with the highest transition rate, Q_0 [Stewart, 1994]. One can then compute the
 146 distribution of the number of steps spent in a particular ancestry, $\{b_n\}_{n=1, \dots, \Lambda}$, up to a cutoff Λ , where
 147 $\sum_{i=1}^{\Lambda} b_i \approx 1$. $\{b_n\}_{n=1, \dots, \Lambda}$ is computed by multiplying the state vector with the transition matrix for
 148 Λ iterations while recording the amount of probability absorbed by the non- p ancestries at each step.
 149 The Erlang distribution models the length of a trajectory, l , with k steps as:

$$\mathbb{E}_{k, Q_0}(l) = \frac{Q_0^k l^{k-1} e^{-Q_0 l}}{(k-1)!} \quad (4)$$

150 This leads to the tract length distribution:

$$\phi(l) = \sum_{k=1}^{\Lambda+1} b_k \mathbb{E}_{k, Q_0}(l) \quad (5)$$

151 **3.2 A Locus-Specific Tract Length Distribution With Selection**

152 Equation 5 describes the length of tracts in a way that is not locus specific. We are interested in how
153 the effects of purifying selection against alleles at two loci under negative selection, according to the
154 incompatibility models described above, may skew the tract length distribution. More specifically, we
155 want to model the distribution of ancestry tracts lengths that are contiguous with a negatively selected
156 allele on a chromosome. In this case, the probability of observing a transition, or recombination event,
157 depends on its recombination distance from the incompatibility loci of interest.

158 We define the number of basepairs between loci A and B to be $v + w = L$, where v is the number
159 of basepairs from the A locus to the v th position and w is the number of basepairs from position $v + 1$
160 to L (Figure 3). We extend the transition matrix Q in equation (3) such that each value of v denotes
161 a new Q_v by multiplying each transition rate by the probability, Ψ_v^τ , that an ancestry junction which
162 arises at time τ at position v survives to the present:

$$Q_v(p, t \rightarrow p', t') = m_{p'}(t') \sum_{\tau=1}^{\min(t, t')-1} \Psi_v^\tau \prod_{s=\tau+1}^{t'-1} (1 - m(s)). \quad (6)$$

163 Equation 6 is computed as a function of the sequence of genotypic backgrounds the junction encounters
164 each generation to the present. Using a two-allele model, let **A** and **a** refer to alternative alleles at the
165 locus of interest, and alleles **B** and **b** refer to the second locus located at some distance away from the
166 **A** locus. We can define a state space, S , of two-locus genotypes in which the junction can exist:

$$S = \begin{bmatrix} \mathbf{AB}|ab \\ \mathbf{AB}|Ab \\ \mathbf{AB}|aB \\ \mathbf{AB}|AB \\ \mathbf{Ab}|ab \\ \mathbf{Ab}|Ab \\ \mathbf{Ab}|aB \\ \mathbf{Ab}|AB \\ \mathbf{aB}|ab \\ \mathbf{aB}|Ab \\ \mathbf{aB}|aB \\ \mathbf{aB}|AB \\ \mathbf{ab}|ab \\ \mathbf{ab}|Ab \\ \mathbf{ab}|aB \\ \mathbf{ab}|AB \\ \epsilon \end{bmatrix}$$

167 where the bold pair of alleles refers to the chromosome on which the junction resides. In cases where
 168 the interacting loci are on different chromosomes, the bold alleles refer to the genomic complement
 169 from which the junction is inherited.

170 Let $\mathbf{P}_v^{t,t-1}$ be a symmetric 17×17 transition matrix among the states in S from time t to $t-1$ for a
 171 junction at the v th position, where $\mathbf{P}_{v,i,j}^{t,t-1}$ refers to the transition from state i to j . The first row in this
 172 matrix is shown below in Equation 7 (see Appendix A for rows 1-17). The transition probabilities in
 173 $\mathbf{P}_v^{t,t-1}$ depend on the fitness of genotypes carrying the junction, ω , the recombination rate between the
 174 interacting loci, r , and the frequency of possible gametes with which to pair in the hybrid population
 175 at time $t-1$: $x_1^{t-1}, x_2^{t-1}, x_3^{t-1}, x_4^{t-1}$. Let x_1, x_2, x_3, x_4 refer to the frequencies of gametes $\mathbf{AB}, \mathbf{Ab},$
 176 \mathbf{aB} and \mathbf{ab} , respectively. Gamete frequencies are computed numerically by simulation [Gavrilets 1997,
 177 Appendix A.5]. Let ω_i denote the marginal fitness of gamete i where $\omega_1, \omega_2, \omega_3, \omega_4$ refer to gametes
 178 $\mathbf{AB}, \mathbf{Ab}, \mathbf{aB}$ and \mathbf{ab} , respectively. Let ω_{ij} refer to the fitness of an individual with gametes i and j .
 179 Figure 3 provides some intuition for how the following transition probabilities in $\mathbf{P}_v^{t,t-1}$ are computed.

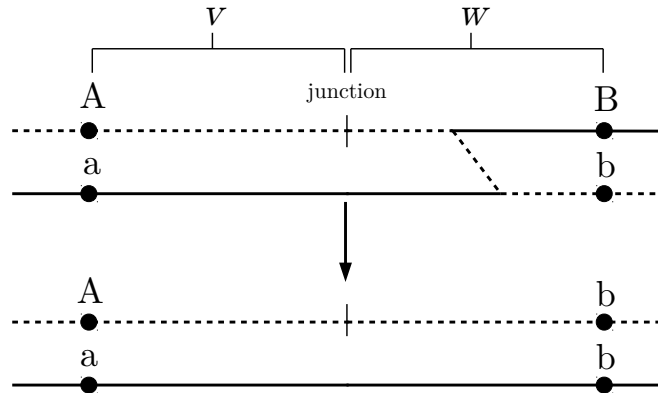


Fig 3. A visual description of the transition probability $\mathbf{P}_{v,1,5}^{t,t-1}$. For the first state in S , $\mathbf{AB}|ab$, the transition probability to state $\mathbf{Ab}|ab$, is a product of the probability that the bold haplotype (\mathbf{AB}) is chosen (0.5), a recombination event occurs between the junction and locus B, $r\frac{w}{v+w}$, the recombined gamete gets paired with gamete x_4 at time $t - 1$, and the individual with genotype $\mathbf{Ab}|ab$ survives, ω_{14} .

$$\mathbf{P}_{v,1,j}^{t,t-1} = \begin{cases} .5\omega_{14}(1-r)x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{14}(1-r)x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{14}(1-r)x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{14}(1-r)x_1^{t-1} & \text{if } j = 4; \\ .5\omega_{14}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 5; \\ .5\omega_{14}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 6; \\ .5\omega_{14}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 7; \\ .5\omega_{14}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 8; \\ .5\omega_{14}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{14}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{14}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{14}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{1,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (7)$$

183 We can define the initial probabilities, π_0^τ , of a junction in each state when it occurs at a partic-
 184 ular time τ . These probabilities will vary depending on the ancestry of interest for the tract length
 185 distribution. Conditional on a recombination event occurring between the two loci, the probability
 186 that the junction occurs at any particular position is uniform ($1/L$). If the ancestry of interest is that
 187 of the A allele, then

$$\pi_0^\tau = \begin{bmatrix} 2p_A^\tau p_B^\tau x_1^\tau \\ 2p_A^\tau p_B^\tau x_2^\tau \\ 2p_A^\tau p_B^\tau x_3^\tau \\ 2p_A^\tau p_B^\tau x_4^\tau \\ 2p_A^\tau p_b^\tau x_1^\tau \\ 2p_A^\tau p_b^\tau x_2^\tau \\ 2p_A^\tau p_b^\tau x_3^\tau \\ 2p_A^\tau p_b^\tau x_4^\tau \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

188 where $p_A^\tau, p_a^\tau, p_B^\tau, p_b^\tau$ are the allele frequencies at time τ . The probability that the junction resides
 189 among each of the states after its origination at time τ to the present is

$$\pi_v^\tau = \pi_0^\tau \prod_{t=0}^{\tau} \mathbf{P}_v^{t,t-1}. \quad (8)$$

190 After defining the vector $\eta = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]$, the survival probability of the
 191 junction is

$$\Psi_v^\tau = 1 - \pi_v^\tau \eta. \quad (9)$$

192 The transition matrix Q_v can now be computed using Equation 6 for all values of v where $v \in$
 193 $\{1 \dots L\}$. In contrast to the transition matrix Q defined in Equation 3, the set of transition matrices

194 Q_v are inhomogeneous over positions v . As a result, the uniformization technique outlined in Stewart
195 [1994] does not apply. However, Andreychenko [2010] describes an approach to uniformize a time-
196 inhomogeneous Markov chain which relies on partitioning the transition matrix into time-dependent
197 and time-independent components. Whereas the time-homogeneous case relies on uniformizing by the
198 constant transition rate of the state with the largest value, the time-inhomogeneous case relies on using
199 the average rate of the state with the largest transition rate value. As before, the distribution for the
200 number of steps in a trajectory, $\{b_n\}_{n=1,\dots,\Lambda}$, can be computed and used with Equations 4 and 5 to
201 calculate the tract length distribution.

202 4 Discussion

203 The model presented above describes an approach which may prove useful in verifying the role of
204 purifying selection against incompatible alleles in a hybrid zone. If shown to be robust under a
205 reasonable set of demographic scenarios and genetic architectures for incompatibility, this model would
206 provide an additional tool for testing the effects of selection on candidate loci which have been identified
207 by QTL mapping of hybrid sterility or inviability traits [White et al., 2011]. This model could also be
208 used to develop an independent test of loci identified by steep clines in allele frequency across a hybrid
209 zone relative to the genomic background [Gompert et al., 2012]. We would like to emphasize the
210 novelty and potential power of considering a genetic architecture of reproductive isolation for which
211 there is strong empirical support (The Dobzhansky-Muller model) in the context of ancestry.

212 While any formal statements regarding the expected tract length distribution would require a
213 full implementation of this model, we can make a few intuitive statements which follow from previous
214 theory and simulation [Hvala et al., 2018, Lindtke and Buerkle, 2015]. In particular, Hvala et al. [2018]
215 show that the number and density of ancestry junctions scales negatively with selection strength at
216 incompatibility loci. This signature was further influenced by the genetic distance of junctions from
217 the loci, the form of selection and dominance.

218 There are several challenges that remain before computing expected tract length distributions
219 and performing inference on parameters of interest. In particular, computing Q_v for a large set of
220 positions may be difficult considering the repeated summation over products in Equation 6, the matrix
221 multiplication required both for Equation 8, and computing $\{b_n\}_{n=1,\dots,\Lambda}$. While Gravel [2012] intended
222 to model admixture events which occurred relatively recently, many hybrid zones of interest are likely
223 to have formed more than 100 generations ago, which produces more computational burden given that
224 the state space of Q_v is $2Tv$. However, it is likely that differences in the junction survival probability,

225 Ψ_v^τ , beyond some value of τ become negligible. The simplified two-locus, two-allele model that we
226 consider is another effort to reduce the parameter space of genotype fitnesses that might result from
227 higher-order epistasis of 3 or more loci.

228 Because Ψ_v^τ is dependent on hybrid zone gamete frequencies in a linear stepping-stone model,
229 deviation from this simplifying assumption will most likely affect the results. The linear stepping-stone
230 model which we borrow from Gavrilets [1997] can be generalized to any number of demes between the
231 two infinite source populations. By implementing our model with this population structure, one could
232 compute tract length distributions as a function of distance from the hybrid zone in a similar spirit
233 to the more geography-explicit approach of Sedghifar et al. [2015, 2016]. The assumption of large
234 population size is another assumption which could also be relaxed (see Appendix 1 in Gravel [2012]),
235 as this will likely influence the rate at which ancestry junctions are fixed or lost from the population.

236 Aside from the challenges of model misspecification, performing inference will be particularly
237 difficult considering the computational burden of computing the tract length distribution for a set
238 of migration rates and fitness matrix parameters. Gravel [2012] uses a maximum-likelihood scheme to
239 identify the set of parameters that best describe the magnitude and timing of migration events from
240 a source into a target population. Given that our primary interest is to infer the effects of purifying
241 selection, it may be more efficient to treat the migration history as a latent variable to be marginalized
242 over using Markov chain Monte Carlo.

243 Despite these challenges, our framework for computing statistical properties of haplotypes in a
244 hybrid zone represents one of only a few recent efforts which aim to exploit the combination of whole-
245 genome sequencing and dense genotyping approaches that have emerged for non-model systems. In
246 particular, this model is the only example that we know of for deriving locus-specific haplotype patterns
247 under epistasis. Given the complexity of this problem, an alternative option may be to use simulation-
248 based classification in a machine learning framework [Chan et al., 2018, Schrider and Kern, 2018,
249 Sheehan and Song, 2016]. Rather than focusing on any one summary statistic, several summary
250 statistics with potential relevance to purifying selection against genetic incompatibilities could be
251 used simultaneously. Alternatively, Chan et al. [2018] describe another machine learning approach
252 which could instead use genotype data directly.

253 Regardless of the methods used to identify genomic patterns of purifying selection against incom-
254 patibility loci, this effort represents one facet of the many lines of evidence necessary to identify and
255 describe the causes of reproductive isolation between species.

256 Acknowledgements

257 We thank Megan Frayer and John Hvala at The University of Wisconsin, Madison for their helpful
 258 discussions and insight. We also thank Yaniv Brandvain for additional advice and encouragement.
 259 Members of the Novembre, Stephens, and He labs provided useful feedback at an early stage. This
 260 work was funded by NSF grant DEB-1353737 to Bret Payseur, DEB-1353737 to Bret Payseur and John
 261 Novembre as well as NSF Graduate Research Fellowship and National Institute Of General Medical
 262 Sciences of the National Institutes of Health under award numbers DGE-1144082 and T32GM007197
 263 to Joel Smith.

264 Appendix A

265 The full transition matrix used to compute the junction survival probabilities in Equation 9:

$$\mathbf{P}_{v,1,j}^{t,t-1} = \begin{cases} .5\omega_{14}(1-r)x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{14}(1-r)x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{14}(1-r)x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{14}(1-r)x_1^{t-1} & \text{if } j = 4; \\ .5\omega_{14}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 5; \\ .5\omega_{14}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 6; \\ .5\omega_{14}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 7; \\ .5\omega_{14}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 8 \\ .5\omega_{14}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{14}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{14}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{14}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 12 \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{1,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (10)$$

$$\mathbf{P}_{v,2,j}^{t,t-1} = \begin{cases} .5\omega_{12}(((1-r)x_4^{t-1}) + (r\frac{v}{v+w}x_4^{t-1})) & \text{if } j = 1; \\ .5\omega_{12}(((1-r)x_2^{t-1}) + (r\frac{v}{v+w}x_2^{t-1})) & \text{if } j = 2; \\ .5\omega_{12}(((1-r)x_3^{t-1}) + (r\frac{v}{v+w}x_3^{t-1})) & \text{if } j = 3; \\ .5\omega_{12}(((1-r)x_1^{t-1}) + (r\frac{v}{v+w}x_1^{t-1})) & \text{if } j = 4; \\ .5\omega_{12}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 5; \\ .5\omega_{12}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 6; \\ .5\omega_{12}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 7; \\ .5\omega_{12}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 8 \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12 \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{2,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (11)$$

$$\mathbf{P}_{v,3,j}^{t,t-1} = \begin{cases} .5\omega_{13}(((1-r)x_4^{t-1}) + (r\frac{w}{v+w}x_4^{t-1})) & \text{if } j = 1; \\ .5\omega_{13}(((1-r)x_2^{t-1}) + (r\frac{w}{v+w}x_2^{t-1})) & \text{if } j = 2; \\ .5\omega_{13}(((1-r)x_3^{t-1}) + (r\frac{w}{v+w}x_3^{t-1})) & \text{if } j = 3; \\ .5\omega_{13}(((1-r)x_1^{t-1}) + (r\frac{w}{v+w}x_1^{t-1})) & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{13}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{13}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{13}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{13}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 12 \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{3,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (12)$$

$$\mathbf{P}_{v,4,j}^{t,t-1} = \begin{cases} .5\omega_{11}(((1-r)x_4^{t-1}) + (rx_4^{t-1})) & \text{if } j = 1; \\ .5\omega_{11}(((1-r)x_2^{t-1}) + (rx_2^{t-1})) & \text{if } j = 2; \\ .5\omega_{11}(((1-r)x_3^{t-1}) + (rx_3^{t-1})) & \text{if } j = 3; \\ .5\omega_{11}(((1-r)x_1^{t-1}) + (rx_1^{t-1})) & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{4,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (13)$$

$$\mathbf{P}_{v,5,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ .5\omega_{24}(((1-r)x_4^{t-1}) + (r\frac{w}{v+w}x_4^{t-1})) & \text{if } j = 5; \\ .5\omega_{24}(((1-r)x_2^{t-1}) + (r\frac{w}{v+w}x_2^{t-1})) & \text{if } j = 6; \\ .5\omega_{24}(((1-r)x_3^{t-1}) + (r\frac{w}{v+w}x_3^{t-1})) & \text{if } j = 7; \\ .5\omega_{24}(((1-r)x_1^{t-1}) + (r\frac{w}{v+w}x_1^{t-1})) & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ .5\omega_{24}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 13; \\ .5\omega_{24}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 14; \\ .5\omega_{24}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 15; \\ .5\omega_{24}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{5,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (14)$$

$$\mathbf{P}_{v,6,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ .5\omega_{22}(((1-r)x_4^{t-1}) + (rx_4^{t-1})) & \text{if } j = 5; \\ .5\omega_{22}(((1-r)x_2^{t-1}) + (rx_2^{t-1})) & \text{if } j = 6; \\ .5\omega_{22}(((1-r)x_3^{t-1}) + (rx_3^{t-1})) & \text{if } j = 7; \\ .5\omega_{22}(((1-r)x_1^{t-1}) + (rx_1^{t-1})) & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{6,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (15)$$

$$\mathbf{P}_{v,7,j}^{t,t-1} = \begin{cases} .5\omega_{23}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{23}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{23}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{23}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 4; \\ .5\omega_{23}(1-r)x_4^{t-1} & \text{if } j = 5; \\ .5\omega_{23}(1-r)x_2^{t-1} & \text{if } j = 6; \\ .5\omega_{23}(1-r)x_3^{t-1} & \text{if } j = 7; \\ .5\omega_{23}(1-r)x_1^{t-1} & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ .5\omega_{23}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 13; \\ .5\omega_{23}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 14; \\ .5\omega_{23}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 15; \\ .5\omega_{23}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{7,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (16)$$

$$\mathbf{P}_{v,8,j}^{t,t-1} = \begin{cases} .5\omega_{12}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{12}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{12}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{12}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 4; \\ .5\omega_{12}(((1-r)x_4^{t-1}) + (r\frac{v}{v+w}x_4^{t-1})) & \text{if } j = 5; \\ .5\omega_{12}(((1-r)x_2^{t-1}) + (r\frac{v}{v+w}x_2^{t-1})) & \text{if } j = 6; \\ .5\omega_{12}(((1-r)x_3^{t-1}) + (r\frac{v}{v+w}x_3^{t-1})) & \text{if } j = 7; \\ .5\omega_{12}(((1-r)x_1^{t-1}) + (r\frac{v}{v+w}x_1^{t-1})) & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{8,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (17)$$

$$\mathbf{P}_{v,9,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{34}(((1-r)x_4^{t-1}) + (r\frac{v}{v+w}x_4^{t-1})) & \text{if } j = 9; \\ .5\omega_{34}(((1-r)x_2^{t-1}) + (r\frac{v}{v+w}x_2^{t-1})) & \text{if } j = 10; \\ .5\omega_{34}(((1-r)x_3^{t-1}) + (r\frac{v}{v+w}x_3^{t-1})) & \text{if } j = 11; \\ .5\omega_{34}(((1-r)x_1^{t-1}) + (r\frac{v}{v+w}x_1^{t-1})) & \text{if } j = 12; \\ .5\omega_{34}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 13; \\ .5\omega_{34}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 14; \\ .5\omega_{34}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 15; \\ .5\omega_{34}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{9,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (18)$$

$$\mathbf{P}_{v,10,j}^{t,t-1} = \begin{cases} .5\omega_{23}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{23}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{23}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{23}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{23}(1-r)x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{23}(1-r)x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{23}(1-r)x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{23}(1-r)x_1^{t-1} & \text{if } j = 12; \\ .5\omega_{23}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 13; \\ .5\omega_{23}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 14; \\ .5\omega_{23}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 15; \\ .5\omega_{23}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{10,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (19)$$

$$\mathbf{P}_{v,11,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{33}(((1-r)x_4^{t-1}) + (rx_4^{t-1})) & \text{if } j = 9; \\ .5\omega_{33}(((1-r)x_2^{t-1}) + (rx_2^{t-1})) & \text{if } j = 10; \\ .5\omega_{33}(((1-r)x_3^{t-1}) + (rx_3^{t-1})) & \text{if } j = 11; \\ .5\omega_{33}(((1-r)x_1^{t-1}) + (rx_1^{t-1})) & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{11,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (20)$$

$$\mathbf{P}_{v,12,j}^{t,t-1} = \begin{cases} .5\omega_{13}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 1; \\ .5\omega_{13}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 2; \\ .5\omega_{13}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 3; \\ .5\omega_{13}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 4 \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{13}(((1-r)x_4^{t-1}) + (r\frac{w}{v+w}x_4^{t-1})) & \text{if } j = 9; \\ .5\omega_{13}(((1-r)x_2^{t-1}) + (r\frac{w}{v+w}x_2^{t-1})) & \text{if } j = 10; \\ .5\omega_{13}(((1-r)x_3^{t-1}) + (r\frac{w}{v+w}x_3^{t-1})) & \text{if } j = 11; \\ .5\omega_{13}(((1-r)x_1^{t-1}) + (r\frac{w}{v+w}x_1^{t-1})) & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16 \\ 1 - \sum_j^{16} \mathbf{P}_{12,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (21)$$

$$\mathbf{P}_{v,13,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ .5\omega_{44}(((1-r)x_4^{t-1}) + (rx_4^{t-1})) & \text{if } j = 13; \\ .5\omega_{44}(((1-r)x_2^{t-1}) + (rx_2^{t-1})) & \text{if } j = 14; \\ .5\omega_{44}(((1-r)x_3^{t-1}) + (rx_3^{t-1})) & \text{if } j = 15; \\ .5\omega_{44}(((1-r)x_1^{t-1}) + (rx_1^{t-1})) & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{13,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (22)$$

$$\mathbf{P}_{v,14,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ .5\omega_{24}r\frac{v}{v+w}x_4^{t-1} & \text{if } j = 5; \\ .5\omega_{24}r\frac{v}{v+w}x_2^{t-1} & \text{if } j = 6; \\ .5\omega_{24}r\frac{v}{v+w}x_3^{t-1} & \text{if } j = 7; \\ .5\omega_{24}r\frac{v}{v+w}x_1^{t-1} & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ .5\omega_{24}(((1-r)x_4^{t-1}) + (r\frac{w}{v+w}x_4^{t-1})) & \text{if } j = 13; \\ .5\omega_{24}(((1-r)x_2^{t-1}) + (r\frac{w}{v+w}x_2^{t-1})) & \text{if } j = 14; \\ .5\omega_{24}(((1-r)x_3^{t-1}) + (r\frac{w}{v+w}x_3^{t-1})) & \text{if } j = 15; \\ .5\omega_{24}(((1-r)x_1^{t-1}) + (r\frac{w}{v+w}x_1^{t-1})) & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{14,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (23)$$

$$\mathbf{P}_{v,15,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ .5\omega_{34}r\frac{w}{v+w}x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{34}r\frac{w}{v+w}x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{34}r\frac{w}{v+w}x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{34}r\frac{w}{v+w}x_1^{t-1} & \text{if } j = 12; \\ .5\omega_{34}(((1-r)x_4^{t-1}) + (r\frac{v}{v+w}x_4^{t-1})) & \text{if } j = 13; \\ .5\omega_{34}(((1-r)x_2^{t-1}) + (r\frac{v}{v+w}x_2^{t-1})) & \text{if } j = 14; \\ .5\omega_{34}(((1-r)x_3^{t-1}) + (r\frac{v}{v+w}x_3^{t-1})) & \text{if } j = 15; \\ .5\omega_{34}(((1-r)x_1^{t-1}) + (r\frac{v}{v+w}x_1^{t-1})) & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{15,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (24)$$

$$\mathbf{P}_{v,16,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ .5\omega_{44}r \frac{v}{v+w} x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{44}r \frac{v}{v+w} x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{44}r \frac{v}{v+w} x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{44}r \frac{v}{v+w} x_1^{t-1} & \text{if } j = 12 \\ .5\omega_{44}r \frac{w}{v+w} x_4^{t-1} & \text{if } j = 9; \\ .5\omega_{44}r \frac{w}{v+w} x_2^{t-1} & \text{if } j = 10; \\ .5\omega_{44}r \frac{w}{v+w} x_3^{t-1} & \text{if } j = 11; \\ .5\omega_{44}r \frac{w}{v+w} x_1^{t-1} & \text{if } j = 12 \\ .5\omega_{44}(1-r)x_4^{t-1} & \text{if } j = 13; \\ .5\omega_{44}(1-r)x_2^{t-1} & \text{if } j = 14; \\ .5\omega_{44}(1-r)x_3^{t-1} & \text{if } j = 15; \\ .5\omega_{44}(1-r)x_1^{t-1} & \text{if } j = 16; \\ 1 - \sum_j^{16} \mathbf{P}_{16,j}^{t,t-1} & \text{if } j = 17 \end{cases} \quad (25)$$

$$\mathbf{P}_{v,17,j}^{t,t-1} = \begin{cases} 0 & \text{if } j = 1; \\ 0 & \text{if } j = 2; \\ 0 & \text{if } j = 3; \\ 0 & \text{if } j = 4; \\ 0 & \text{if } j = 5; \\ 0 & \text{if } j = 6; \\ 0 & \text{if } j = 7; \\ 0 & \text{if } j = 8; \\ 0 & \text{if } j = 9; \\ 0 & \text{if } j = 10; \\ 0 & \text{if } j = 11; \\ 0 & \text{if } j = 12; \\ 0 & \text{if } j = 13; \\ 0 & \text{if } j = 14; \\ 0 & \text{if } j = 15; \\ 0 & \text{if } j = 16; \\ 1 & \text{if } j = 17 \end{cases} \quad (26)$$

268 References

- 269 Aleksandr Andreychenko. *Uniformization for time-inhomogeneous Markov population models*. PhD
270 thesis, Saarland University, 2010.
- 271 NH Barton. The dynamics of hybrid zones. *Heredity*, 43(3):341, 1979.
- 272 Nicholas H Barton and Godfrey M Hewitt. Analysis of hybrid zones. *Annual Review of Ecology and*
273 *systematics*, 16(1):113–148, 1985.
- 274 Jeffrey Chan, Valerio Perrone, Jeffrey P Spence, Paul A Jenkins, Sara Mathieson, and
275 Yun S Song. A likelihood-free inference framework for population genetic data us-
276 ing exchangeable neural networks. *bioRxiv*, 2018. doi: 10.1101/267211. URL
277 <https://www.biorxiv.org/content/early/2018/02/18/267211>.
- 278 Theodosius Dobzhansky. *Genetics and the Origin of Species*. Columbia University Press, 1937.
- 279 John A Endler. Gene flow and population differentiation: studies of clines suggest that differentiation
280 along environmental gradients may be independent of gene flow. *Science*, 179(4070):243–250, 1973.
- 281 Sergey Gavrilets. Hybrid zones with dobzhansky-type epistatic selection. *Evolution*, 51(4):1027–1035,
282 1997.
- 283 Zachariah Gompert and C Alex Buerkle. Analyses of genetic ancestry enable key insights for molecular
284 ecology. *Molecular Ecology*, 22(21):5278–5294, 2013.
- 285 Zachariah Gompert, Thomas L Parchman, and C Alex Buerkle. Genomics of isolation in hybrids.
286 *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 367(1587):439–
287 450, 2012.
- 288 Simon Gravel. Population genetics models of local ancestry. *Genetics*, 191(2):607–619, 2012.
- 289 Kelley Harris and Rasmus Nielsen. Inferring demographic history from a spectrum of shared haplotype
290 lengths. *PLoS Genetics*, 9(6):e1003521, 2013.
- 291 Richard G Harrison and Erica L Larson. Heterogeneous genome divergence, differential introgression,
292 and the origin and structure of hybrid zones. *Molecular Ecology*, 25(11):2454–2466, 2016.
- 293 Garrett Hellenthal, George BJ Busby, Gavin Band, James F Wilson, Cristian Capelli, Daniel Falush,
294 and Simon Myers. A genetic atlas of human admixture history. *Science*, 343(6172):747–751, 2014.

- 295 John A Hvala, Megan E Frayer, and Bret A Payseur. Signatures of hybridization and speciation in
296 genomic patterns of ancestry. *Evolution*, 2018.
- 297 Mason Liang and Rasmus Nielsen. The lengths of admixture tracts. *Genetics*, 197(3):953–967, 2014.
- 298 Dorothea Lindtke and C Alex Buerkle. The genetic architecture of hybrid incompatibilities and their
299 effect on barriers to introgression in secondary contact. *Evolution*, 69(8):1987–2004, 2015.
- 300 Po-Ru Loh, Mark Lipson, Nick Patterson, Priya Moorjani, Joseph K Pickrell, David Reich, and Bonnie
301 Berger. Inferring admixture histories of human populations using linkage disequilibrium. *Genetics*,
302 193(4):1233–1254, 2013.
- 303 Shamoni Maheshwari and Daniel A Barbash. The genetics of hybrid incompatibilities. *Annual Review*
304 *of Genetics*, 45:331–355, 2011.
- 305 Nick Patterson, Priya Moorjani, Yontao Luo, Swapan Mallick, Nadin Rohland, Yiping Zhan, Teri
306 Genschoreck, Teresa Webster, and David Reich. Ancient admixture in human history. *Genetics*,
307 192(3):1065–1093, 2012.
- 308 Bret A Payseur. Using differential introgression in hybrid zones to identify genomic regions involved
309 in speciation. *Molecular Ecology Resources*, 10(5):806–820, 2010.
- 310 Bret A Payseur and Loren H Rieseberg. A genomic perspective on hybridization and speciation.
311 *Molecular Ecology*, 25(11):2337–2360, 2016.
- 312 John E Pool and Rasmus Nielsen. Inference of historical changes in migration rate from the lengths
313 of migrant tracts. *Genetics*, 181(2):711–719, 2009.
- 314 Daven C Presgraves. The molecular evolutionary basis of species formation. *Nature Reviews Genetics*,
315 11(3):175, 2010.
- 316 Alkes L Price, Arti Tandon, Nick Patterson, Kathleen C Barnes, Nicholas Rafaels, Ingo Ruczinski,
317 Terri H Beaty, Rasika Mathias, David Reich, and Simon Myers. Sensitive detection of chromosomal
318 segments of distinct ancestry in admixed populations. *PLoS Genetics*, 5(6):e1000519, 2009.
- 319 Douglas W Schemske. Understanding the origin of species. *Evolution*, 54(3):1069–1073, 2000.
- 320 Douglas W Schemske and HD Bradshaw. Pollinator preference and the evolution of floral traits in
321 monkeyflowers (*Mimulus*). *Proceedings of the National Academy of Sciences*, 96(21):11910–11915,
322 1999.

- 323 Dolph Schluter. Evidence for ecological speciation and its alternative. *Science*, 323(5915):737–741,
324 2009.
- 325 Daniel R Schrider and Andrew D Kern. Supervised machine learning for population genetics: A new
326 paradigm. *Trends in Genetics*, 2018.
- 327 Molly Schumer and Yaniv Brandvain. Determining epistatic selection in admixed populations. *Molec-*
328 *ular Ecology*, 25(11):2577–2591, 2016.
- 329 Molly Schumer, Rongfeng Cui, Daniel L Powell, Rebecca Dresner, Gil G Rosenthal, and Peter An-
330 dolfatto. High-resolution mapping reveals hundreds of genetic incompatibilities in hybridizing fish
331 species. *Elife*, 3, 2014.
- 332 Alisa Sedghifar, Yaniv Brandvain, Peter Ralph, and Graham Coop. The spatial mixing of genomes in
333 secondary contact zones. *Genetics*, 201(1):243–261, 2015.
- 334 Alisa Sedghifar, Yaniv Brandvain, and Peter Ralph. Beyond clines: lineages and haplotype blocks in
335 hybrid zones. *Molecular Ecology*, 25(11):2559–2576, 2016.
- 336 Ole Seehausen, Roger K Butlin, Irene Keller, Catherine E Wagner, Janette W Boughman, Paul A Ho-
337 henlohe, Catherine L Peichel, Glenn-Peter Saetre, Claudia Bank, Åke Brännström, et al. Genomics
338 and the origin of species. *Nature Reviews Genetics*, 15(3):176, 2014.
- 339 Sara Sheehan and Yun S Song. Deep learning for population genetic inference. *PLoS Computational*
340 *Biology*, 12(3):e1004845, 2016.
- 341 William J Stewart. *Introduction to the numerical solution of Markov chains*. Princeton University
342 Press, 1994.
- 343 Andrea L Sweigart and John H Willis. Molecular evolution and genetics of postzygotic reproductive
344 isolation in plants. *F1000 Biology Reports*, 4, 2012.
- 345 Peter Tiffin and Jeffrey Ross-Ibarra. Advances and limits of using population genetics to understand
346 local adaptation. *Trends in Ecology & Evolution*, 29(12):673–680, 2014.
- 347 Michael Turelli, Jeremy R Lipkowitz, and Yaniv Brandvain. On the Coyne and Orr-igin of species:
348 effects of intrinsic postzygotic isolation, ecological differentiation, X chromosome size, and sympatry
349 on *Drosophila* speciation. *Evolution*, 68(4):1176–1187, 2014.

- 350 Michael A White, Brian Steffy, Tim Wiltshire, and Bret A Payseur. Genetic dissection of a key
351 reproductive barrier between nascent species of house mice. *Genetics*, 189(1):289–304, 2011.
- 352 MJD White. Models of speciation. *Science*, 159(3819):1065–1070, 1968.