

SUPPLEMENTARY INFORMATION - 1

Modeling selection against introgression.

Ivan Juric, Simon Aeschbacher, Graham Coop

1 Introduction

Here, we describe several models of a single pulse of admixture between Neanderthal and modern humans, and derive approximations for the present-day frequency of a neutral introgressed Neanderthal allele linked to one or multiple sites under purifying selection in humans. We then demonstrate the accuracy of these approximations by comparing them to numerically iterated recursion equations and individual-based simulations. Lastly, we consider models of single and multiple waves of continuous introgression and show that one cannot distinguish between these models and a single-pulse admixture model using the present-day frequency of introgressed alleles as the only source of information.

2 Single-pulse introgression models

In the following, we derive deterministic approximations to the frequency of a neutral allele linked to locus under purifying selection after a single pulse of admixture. We consider a neutral polymorphism on an autosome and on the X chromosome, and in both cases we allow for one or multiple linked loci under selection.

2.1 A single autosomal locus under selection

We model the allele-frequency dynamics at a neutral locus linked to a selected locus following a single pulse of admixture (introgression) from the Neanderthal population that happened t generations ago. Let N_1 and N_2 be the Neanderthal and human alleles at the neutral locus, and S_1 and S_2 the two alleles at the locus under selection. We denote the recombination rate between the two loci by r , and assume that the Neanderthal-derived allele S_1 is deleterious in humans. Specifically, the viability of a human individual heterozygous at the locus under selection is $w_{12} = w_{21} = 1 - s$, where $0 < s < 1$, and the viability of a human S_2 homozygote is $w_{22} = 1$. We further assume that the frequency of S_1 in humans is low, so that the deleterious S_1S_1 homozygotes are very rare and can be ignored. Numerical and individual-based simulations (see below) show that this assumption is reasonable for the parameter values we consider. We assume that, prior to admixture, Neanderthals and humans were fixed for allele N_1 and N_2 , respectively. At the time of introgression, the frequency of N_1 rises instantaneously from 0 to p_0 in the human population.

In the following, we describe how p_t , the frequency of N_1 t generations after introgression, depends on its initial frequency p_0 , the heterozygote selection coefficient s , and the recombination fraction r . Let x_0 and y_0 denote the frequency of haplotypes S_1N_1 and S_2N_1 in humans at the time of admixture. Hence, the frequency of allele N_1 immediately after admixture can be written as

$$p_0 = x_0 + y_0. \tag{1}$$

After recombination and random union of gametes, the haplotype frequencies in zygotes can be approximated by

$$x_0^* = x_0(1 - r), \tag{2a}$$

$$y_0^* = x_0r + y_0, \tag{2b}$$

still assuming that S_1 is initially rare in humans.

After viability selection, the haplotype frequencies in the next generation of adult humans become

$$x_1 = x_0^*(1 - s) = x_0(1 - s)(1 - r), \quad (3a)$$

$$y_1 = y_0^* = x_0r + y_0. \quad (3b)$$

From here, it is straightforward to obtain the explicit equations for the haplotype frequencies at generation t as

$$x_t = x_0[(1 - s)(1 - r)]^t, \quad (4a)$$

$$y_t = r \sum_{i=0}^{t-1} x_i + y_0. \quad (4b)$$

This can be simplified to

$$x_t = x_0[(1 - s)(1 - r)]^t, \quad (5a)$$

$$y_t = x_0r \frac{1 - [(1 - s)(1 - r)]^t}{1 - (1 - s)(1 - r)} + y_0. \quad (5b)$$

Lastly, we plug equation (5) into equation (1), which yields

$$p_t = x_0f(s, r, t) + y_0, \quad (6)$$

where

$$f(s, r, t) = \frac{[(1 - s)(1 - r)]^t [1 - r - (1 - s)(1 - r)] + r}{1 - (1 - s)(1 - r)}. \quad (7)$$

As time t goes to infinity, $f(s, r, t)$ approaches Bengtsson's Bengtsson [1985] gene flow factor for the case when selection happens before migration (cf. equation A3 in Charlesworth et al. [1997]):

$$p_\infty = x_0 \frac{r}{1 - (1 - s)(1 - r)} + y_0. \quad (8)$$

If s and r are small, at equilibrium $f(s, r, t)$ further simplifies to $r/(r+s)$, which is equal to the approximation to the effective rate of gene flow found by Petry Petry [1983] based on a diffusion approximation. If $y_0 = 0$, we can replace x_0 by p_0 in the equations above.

2.2 A single X-chromosomal locus under selection

The non-pseudoautosomal X chromosome is a special case, because of the differences in transmission and the fact that recombination only happens in females. We take this into account by modifying equation (2) to

$$x_{X,0}^* = \frac{2}{3}x_{X,0}(1-r) + \left(1 - \frac{2}{3}\right)x_{X,0} = x_{X,0} \left(1 - \frac{2}{3}r\right) \quad (9a)$$

$$y_{X,0}^* = \frac{2}{3}x_{X,0}r + y_{X,0} \quad (9b)$$

The factor of $2/3$ accounts for the fact that two thirds of all X chromosomes are found in females.

As above, we assume that selection acts on viability and after recombination and random mating, but we now allow for selection to be sex specific. We denote by s_f and s_m the strength of selection against female and male carriers of a single S_1 allele, respectively. Recalling that, at the time of selection, a proportion of $2/3$ of the X chromosomes are found in females, and $1/3$ in males, we obtain the haplotype frequencies in adults of the next generation as

$$\begin{aligned} x_{X,1} &= \frac{2}{3}x_{X,0}^*(1-s_f) + \frac{1}{3}x_{X,0}^*(1-s_m) \\ &= x_{X,0} \left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right), \end{aligned} \quad (10a)$$

$$\begin{aligned} y_{X,1} &= y_{X,0}^* \\ &= \frac{2}{3}x_{X,0}r + y_{X,0}. \end{aligned} \quad (10b)$$

Iteration of equation (10) yields the explicit haplotype frequencies at time t ,

$$x_{X,t} = x_{X,0} \left[\left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right) \right]^t, \quad (11a)$$

$$y_{X,t} = rx_{X,0} \frac{2}{3} \frac{1 - \left[\left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right)\right]^t}{1 - \left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right)} + y_{X,0}. \quad (11b)$$

The frequency of allele N_1 at time t , $p_{X,t} = x_{X,t} + y_{X,t}$, can therefore be written as

$$p_{X,t} = x_{X,0} \frac{\left(1 - \frac{2}{3}r\right)^{t+1} \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right)^t \left(\frac{1}{3}s_m + \frac{2}{3}s_f\right) + \frac{2}{3}r}{1 - \left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right)} + y_{X,0}. \quad (12)$$

By setting $s_f = s_m = s$ and allowing recombination to happen in the entire population, i.e. by replacing $1/3$ and $2/3$ by 0 and 1 , respectively, we recover equation (6), as expected. At equilibrium ($t = \infty$), equation (12) becomes

$$p_{x,\infty} = x_{X,0} \frac{\frac{2}{3}r}{1 - \left(1 - \frac{2}{3}r\right) \left(1 - \frac{1}{3}s_m - \frac{2}{3}s_f\right)} + y_{X,0}, \quad (13)$$

which can be approximated as

$$\begin{aligned} p_{x,\infty} &\approx x_{X,0} \frac{\frac{2}{3}r}{\frac{1}{3}s_m + \frac{2}{3}s_f + \frac{2}{3}r} + y_{X,0} \\ &= x_{X,0} \left(1 + \frac{\frac{1}{2}s_m + s_f}{r}\right)^{-1} + y_{X,0} \end{aligned} \quad (14)$$

if both selection and recombination are weak ($s_m, s_f, r \ll 1$).

2.3 Multiple autosomal loci under selection

In the following, we consider a neutral locus embedded in a suite of multiple autosomal loci under purifying selection. Due to the complexities of multilocus models, we make the following simplifications. First, we assume that, immediately before introgression, all deleterious alleles were fixed in Neanderthals, but absent in humans. Second, we directly consider the situation where the frequency of the neutral Neanderthal-derived alleles has reached equilibrium in

the human population (i.e. $t = \infty$). This can be justified if purifying selection is strong relative to the time since introgression, the admixture proportion, and recombination. Third, we directly draw on the above-mentioned analogy between the surviving proportion of a cohort of introgressed alleles and the effective rate of gene flow at a neutral locus experiencing linkage-mediated selection (cf. equation 7 and subsequent text). We can therefore approximate the equilibrium frequency of the introgressed neutral allele N_1 by modifying the multilocus version of the effective migration rate in equation (24) of reference Aeschbacher and Bürger [2014].

Specifically, let there be I and J loci under selection on the left- and right-hand side of the focal neutral site, and denote them by \mathbf{A}_i and \mathbf{B}_j , respectively. The equilibrium frequency of the introgressed neutral allele N_1 can then be approximated as

$$p_{\infty, IJ} \approx p_0 \left(1 + \frac{a_1}{r_{\mathbf{A}_1}}\right)^{-1} \times \prod_{i=2}^I \left(1 + \frac{a_i}{\sum_{k=1}^{i-1} a_k + r_{\mathbf{A}_i}}\right)^{-1} \times \left(1 + \frac{b_1}{r_{\mathbf{B}_1}}\right)^{-1} \times \prod_{j=2}^J \left(1 + \frac{b_j}{\sum_{k=1}^{j-1} b_k + r_{\mathbf{B}_j}}\right)^{-1}, \quad (15)$$

where a_i and b_i are the selection coefficients against the deleterious mutations at locus \mathbf{A}_i and \mathbf{B}_j , respectively, and $r_{\mathbf{A}_i}$ and $r_{\mathbf{B}_j}$ are the recombination fractions between the neutral locus and the respective locus under selection. Equation (15) assumes that both selection and recombination are weak.

If we set the selection coefficient at all loci under selection to the same value s , equation (15) can be simplified to

$$p_{\infty, IJ} \approx p_0 \prod_{i=1}^I \left[1 + \frac{s}{s(i-1) + r_i}\right]^{-1} \prod_{j=1}^J \left[1 + \frac{s}{s(j-1) + r_j}\right]^{-1}, \quad (16)$$

where r_i and r_j are short cuts for $r_{\mathbf{A}_i}$ and $r_{\mathbf{B}_j}$, respectively.

To assess the accuracy of equation (15), we derived discrete-time recursion equations for a model with one neutral and two selected loci. As before, we

assumed that the admixture proportion is small, so that homozygous carriers of Neanderthal-derived alleles are very rare in the human population and the dynamics of the full diploid model can be approximated by considering a haploid model with four haplotypes. In addition, we assumed that the mean fitness of the human population was not affected by the few carriers of deleterious introgressed mutations. To simplify our notation, we denote the two loci under selection by A and B, and use A_1 (A_2) and B_1 (B_2) for the deleterious (advantageous) alleles at locus A and B, respectively. We consider the following two configurations: one in which the neutral locus N is flanked by one locus under selection on each side (A–N–B), and another where the neutral locus is flanked by two selected loci on one side (A–B–N). We denote by r_{XY} the recombination rate between locus X and Y, where $r_{XY} = r_{YX}$, and we assume that recombination distances accumulate additively across loci.

For configuration A–N–B, the four focal haplotypes are $A_1N_1B_1$, $A_1N_1B_2$, $A_2N_1B_1$, and $A_2N_1B_2$. We denote their frequencies among adults of generation t after viability selection but before recombination and random mating by $x_1(t)$, $x_2(t)$, $x_3(t)$, and $x_4(t)$, respectively. After recombination, random mating, and viability selection, the haplotype frequencies among adults of the next generation are

$$\begin{aligned}
 x_1(t+1) &= w_1 (1 - r_{AN}) (1 - r_{BN}) x_1(t), \\
 x_2(t+1) &= w_2 [(1 - r_{AN}) r_{BN} x_1(t) + (1 - r_{AN} + r_{AN} r_{BN}) x_2(t)], \\
 x_3(t+1) &= w_3 [r_{AN} (1 - r_{BN}) x_1(t) + (1 - r_{BN} + r_{AN} r_{BN}) x_3(t)], \\
 x_4(t+1) &= w_4 [r_{AN} r_{BN} x_1(t) + r_{AN} (1 - r_{BN}) x_2(t) \\
 &\quad + (1 - r_{AN}) r_{BN} x_3(t) + x_4(t)],
 \end{aligned} \tag{17}$$

where w_i denotes the relative fitness of haplotype i .

For configuration A–B–N, the four focal haplotypes are $A_1B_1N_1$, $A_1B_2N_1$, $A_2B_1N_1$, and $A_2B_2N_1$. In this case, the haplotype frequencies follow the fol-

lowing recursions:

$$\begin{aligned}
x_1(t+1) &= w_1 (1 - r_{\text{AB}}) (1 - r_{\text{BN}}) x_1(t), \\
x_2(t+1) &= w_2 (1 - r_{\text{AB}} + r_{\text{AB}} r_{\text{BN}}) x_2(t), \\
x_3(t+1) &= w_3 [r_{\text{AB}} (1 - r_{\text{BN}}) x_1(t) + (1 - r_{\text{BN}}) x_3(t)], \\
x_4(t+1) &= w_4 [r_{\text{BN}} x_1(t) + r_{\text{AB}} (1 - r_{\text{BN}}) x_2(t) + r_{\text{BN}} x_3(t) + x_4(t)].
\end{aligned}
\tag{18}$$

For both configurations, the frequency of the introgressed neutral allele N_1 at time t can be approximated by

$$p_t = p(t) = \sum_{i=1}^4 x_i(t), \tag{19}$$

where the $x_i(t)$ evolve according to equations (17) and (18), depending on the configuration.

We assume that fitness is additive across loci, and parametrize it as

$$\begin{aligned}
w_1 &= 1 - a - b, \\
w_2 &= 1 - a, \\
w_3 &= 1 - b, \\
w_4 &= 1,
\end{aligned}
\tag{20}$$

where $0 \leq a, b \leq 1$.

We numerically iterated equations (17) and (18) and computed p_t according to equation (19) at each step until an equilibrium was reached. Specifically, we terminated the process when the absolute difference between consecutive values p_t and p_{t+1} became smaller than 10^{-9} . We also iterated equations (17) and (18) over a fixed number of $t = 2000$ generations and computed p_t . The approximation in equation (15) performs very well if the underlying assumptions are met and an equilibrium has been reached (Fig S1.6). However, if an equilibrium has not been reached, the approximation in (15) should not be used (compare upward black triangles to respective black curves in Fig S1.6). Moreover, if the

assumption of recombination being weak relative to selection is violated, equation (15) tends to underestimate the actual equilibrium frequency of the neutral introgressed allele, as expected. This effect is particularly strong if genetic distances between consecutive linked loci under selection are highly asymmetric (compare blue triangles and respective blue curves in Fig S1.6B for weak r_{AB}).

2.4 Multiple X-chromosomal loci under selection

Finally, we turn to the case of a neutral locus linked to multiple loci under selection on the X chromosome. Our autosomal results indicate the agreement of parameters fit under the single and multiple locus equilibrium models to empirical data. Therefore, given that the X chromosome results are in a similar region of parameter space we do not fit this final model to data, but we include it here for completeness. Let there be I and J loci under selection on the left- and right-hand side of the focal neutral site, and denote them by A_i and B_j , respectively, where $i = 1, \dots, I$ and $j = 1, \dots, J$. Together, our equation (13) and equation (24) from reference Aeschbacher and Bürger [2014] suggest that the equilibrium frequency of the introgressed neutral allele N_1 may be approximated as

$$p_{\infty, X, IJ} \approx p_{0, X} \left(1 + \frac{\frac{1}{2}a_{1,m} + a_{1,f}}{r_{A_1}} \right)^{-1} \times \prod_{i=2}^I \left[1 + \frac{\frac{1}{2}a_{i,m} + a_{i,f}}{\sum_{k=1}^{i-1} (\frac{1}{2}a_{k,m} + a_{k,f}) + r_{A_i}} \right]^{-1} \\ \times \left(1 + \frac{\frac{1}{2}b_{1,m} + b_{1,f}}{r_{B_1}} \right)^{-1} \times \prod_{j=2}^J \left[1 + \frac{\frac{1}{2}b_{j,m} + b_{j,f}}{\sum_{k=1}^{j-1} (\frac{1}{2}b_{k,m} + b_{k,f}) + r_{B_j}} \right]^{-1}, \quad (21)$$

where $a_{i,f}$ ($a_{i,m}$) and $b_{j,f}$ ($b_{j,m}$) are the coefficients of selection against heterozygous carriers of the deleterious mutation at locus A_i and B_j in females (males). Moreover, r_{A_i} and r_{B_j} are the recombination rates between the neutral locus and locus A_i and B_j , respectively. The factor of $1/2$ accounts for the fact that a given X-chromosomal haplotype spends half of its time in males relative to the time spent in females.

If we fix the selection coefficients in females and males across all loci to s_f and s_m as above, equation (21) can be simplified to

$$p_{\infty, X, IJ} \approx p_{0, X} \prod_{i=1}^I \left[1 + \frac{\frac{1}{2}s_m + s_f}{\left(\frac{1}{2}s_m + s_f\right)(i-1) + \frac{2}{3}r_i} \right]^{-1} \times \prod_{j=1}^J \left[1 + \frac{\frac{1}{2}s_m + s_f}{\left(\frac{1}{2}s_m + s_f\right)(j-1) + \frac{2}{3}r_j} \right]^{-1}, \quad (22)$$

where r_i and r_j are short cuts for r_{A_i} and r_{B_j} , respectively.

To test our conjecture in equation (21), we derived discrete-time recursions for a model with one neutral and two selected loci analogous to those in equations (17) and (18), but for the case of the X chromosome. As before, we assumed that the admixture proportion is small, so that the diploid model can be approximated by a haploid model with four haplotypes, and the mean fitness of the human population was not affected by introgression of deleterious mutations. We again simplify our notation by denoting the two loci under selection by A and B, and we use A_1 (A_2) and B_1 (B_2) for the deleterious (advantageous) alleles at locus A and B, respectively. As above, we consider the two configurations A–N–B and A–B–N, and assume that recombination fractions accumulate additively across loci.

For configuration A–N–B, the four haplotypes of interest are $A_1N_1B_1$, $A_1N_1B_2$, $A_2N_1B_1$, and $A_2N_1B_2$, and their frequencies after recombination in generation t are given by

$$\begin{aligned} x_1^*(t) &= \frac{1}{3}x_1(t) + \frac{2}{3}(1 - r_{AN})(1 - r_{BN})x_1(t), \\ x_2^*(t) &= \frac{1}{3}x_2(t) + \frac{2}{3}[(1 - r_{AN})r_{BN}x_1(t) + (1 - r_{AN} + r_{AN}r_{BN})x_2(t)], \\ x_3^*(t) &= \frac{1}{3}x_3(t) + \frac{2}{3}[r_{AN}(1 - r_{BN})x_1(t) + (1 - r_{BN} + r_{AN}r_{BN})x_3(t)], \\ x_4^*(t) &= x_4(t) \\ &\quad + \frac{2}{3}[r_{AN}r_{BN}x_1(t) + r_{AN}(1 - r_{BN})x_2(t) + (1 - r_{AN})r_{BN}x_3(t)], \end{aligned} \quad (23)$$

respectively. After random mating and viability selection, the frequency of

haplotype i among adults in the next generation is

$$x_i(t+1) = \frac{1}{3}m_i x_i^*(t) + \frac{2}{3}f_i x_i^*(t), \quad (24)$$

where m_i and f_i are the relative fitnesses of haplotype i in males and females, respectively. The frequency of the introgressed neutral allele N_1 at time t is then obtained from equation (19), where the $x_i(t)$ behave as described in equation (24).

For configuration A–B–N, the four haplotypes of interest are $A_1B_1N_1$, $A_1B_2N_1$, $A_2B_1N_1$, and $A_2B_2N_1$. Their frequencies after recombination in generation t are

$$\begin{aligned} x_1^*(t) &= \frac{1}{3}x_1(t) + \frac{2}{3}(1-r_{AB})(1-r_{BN})x_1(t), \\ x_2^*(t) &= \frac{1}{3}x_2(t) + \frac{2}{3}(1-r_{AB}+r_{AB}r_{BN})x_2(t), \\ x_3^*(t) &= \frac{1}{3}x_3(t) + \frac{2}{3}[r_{AB}(1-r_{BN})x_1(t) + (1-r_{BN})x_3(t)], \\ x_4^*(t) &= x_4(t) + \frac{2}{3}[r_{BN}x_1(t) + r_{AB}(1-r_{BN})x_2(t) + r_{BN}x_3(t)], \end{aligned} \quad (25)$$

respectively. Note that $r_{AN} = r_{AB} + r_{BN}$ by assumption. Equations (24) and (19) remain unchanged.

As above, we assume that fitness is additive across loci. For females, we parametrize fitnesses as

$$\begin{aligned} f_1 &= 1 - a_f - b_f, \\ f_2 &= 1 - a_f, \\ f_3 &= 1 - b_f, \\ f_4 &= 1, \end{aligned} \quad (26)$$

and for males, we set

$$\begin{aligned}
m_1 &= 1 - a_m - b_m, \\
m_2 &= 1 - a_m, \\
m_3 &= 1 - b_m, \\
m_4 &= 1,
\end{aligned}
\tag{27}$$

where $0 \leq a_f, a_m, b_f, b_m \leq 1$.

We numerically iterated equation (24) and computed p_t according to equation (19) at each step until an equilibrium was reached. As above, we terminated the process when the absolute difference between consecutive values p_t and p_{t+1} became smaller than 10^{-9} . We also iterated equation (24) over a fixed number of $t = 2000$ generations and computed p_t . As expected, our conjecture in equation (21) provides a very good approximation if the underlying assumptions are met and an equilibrium has been reached (Fig S1.7). However, as in the autosomal case discussed above, if an equilibrium has not been reached, or if the assumption of recombination being weak relative to selection is violated, the approximation in (21) tends to underestimate the actual frequency of the neutral introgressed allele (Fig S1.7).

2.5 Evaluating the accuracy of approximations

We evaluate the accuracy of equation (6) and (12) in two ways. For each test we pick the range of values of r based on the size of windows used in our analysis of human genome. We pick s such that it incorporates the values inferred by our inference procedure (see S2 for details).

In the first test, we numerically solve the recursion for a model in which two loci are under selection (Lewontin and Kojima [1960] or equation 2.9 on page 45 in Rice [2004] for autosomal chromosomes, and the corrected equation 10 from Connallon and Clark [2013] for the X chromosome) by setting the selection coefficient of second selected locus to zero. We compare the recursion results to

equations (6) and (12). This test the effects of ignoring homozygous individuals. To obtain the equilibrium results using recursions, we iterate until the difference in allele frequency between two generations is less than 10^{-9} . The results of this test show that ignoring homozygous individuals does not affect p_t substantially (Figures S1.1 - S1.4).

Our RSS inference method (S2 Text) works with expectations based on deterministic expressions for the Neanderthal allele frequency, and we view drift as ‘noise’ around these expectations. In the second test, we therefore test whether genetic drift has a substantial influence on the average frequency of neutral allele. We do this by performing individual based simulations and calculating the difference between the frequency of neutral alleles at the end of simulation, generation 2000, ($p_{t,sim}$) and that obtained by equation (6). For individual-based simulations, we assume a single constant population size of $N = 10000$ diploid individuals of which $2Np_0$ are double heterozygous (N_2S_2/N_1S_1) at the start of simulation. Generations are non-overlapping. We model soft selection (roulette wheel selection) and assume a multiplicative fitness scheme. This means that individuals homozygous for deleterious allele have fitness $(1 - s)^2$. We plot the $p_{t,sim}$ against p_t (Figure S1.5) along with approximate 95% CI for $p_{t,sim}$. The plot shows a good agreement between average observed and expected deterministic results (Figure S1.5), thus our approach of using deterministic equations is valid in our parameter regime.

In conclusion, our tests suggest that for the parameter range of interest, our approximations describe the expected frequency of neutral alleles well.

2.6 Models of waves of introgression

The interbreeding with Neanderthals likely happened over many generations, rather than in one generation, as we assume with the single pulse admixture model. Furthermore, there is evidence of at least two waves of Neanderthal introgression into the East Asian population [Wall et al., 2013, Vernot and Akey, 2014, 2015, Kim and Lohmueller, 2015]. In this section we derive the expres-

sion for present-day frequency of a neutral Neanderthal allele linked to selected allele under models of continuous and dual waves of introgression. We start by considering how haplotype frequencies change over time during continuous introgression. We then use this result to derive the equivalent expressions for single and dual wave introgression models. Lastly, we show that a single pulse model can produce the same present-day frequencies as the wave models.

Haplotype frequencies during continuous introgression

In the simplest continuous introgression model, each generation a constant fraction of N_1S_1 and N_1S_2 haplotypes introgress into the human population. We consider a slightly more complex model where the fraction of haplotypes that introgresses into humans during the first generation is x_0 and y_0 and in all following generations mx_0 and ny_0 , where $m > 0$ and $n > 0$. We consider this, more complicated, model because it, rather than the simplest model, is used for deriving the dual wave model.

First generation

The following events happen each generation: 1) admixture (inflow of Neanderthal alleles), 2) recombination and mating 3) selection.

Before admixture the frequency of N_1S_1 and N_1S_2 haplotypes in humans is zero. In other words, $x_0 = y_0 = 0$.

After the admixture we have

$$\begin{aligned}x_0^m &= x_0, \\y_0^m &= y_0.\end{aligned}\tag{28}$$

Recombination changes haplotype frequency in the same way as in single pulse model:

$$\begin{aligned}x_0^* &= x_0^m(1 - r), \\y_0^* &= x_0^m r + y_0^m.\end{aligned}\tag{29}$$

After selection the haplotype frequencies at the end of the first generation (or beginning of the second generation) are

$$\begin{aligned}x_1 &= x_0^*(1 - s), \\y_1 &= y_0^*.\end{aligned}\tag{30}$$

Second and later generations

The life cycle in the second and following generations are the same as during the first generation, except now admixture increases the Neanderthal haplotype frequency in humans by mx_0 and nx_0 .

Therefore, for $t = 1, 2, 3, \dots$:

After admixture

$$\begin{aligned}x_t^m &= x_t + x_0m, \\y_t^m &= y_t + y_0n.\end{aligned}\tag{31}$$

After recombination

$$\begin{aligned}x_t^* &= x_t^m(1 - r), \\y_t^* &= x_t^m r + y_t^m.\end{aligned}\tag{32}$$

After selection

$$\begin{aligned}x_{t+1} &= x_t^*(1 - s), \\y_{t+1} &= y_t^*.\end{aligned}\tag{33}$$

We obtain the following explicit expression for the x_t and y_t by using same methods as in single pulse model

$$\begin{aligned}x_t &= x_0 \left(G^t + m \frac{G^t - G}{G - 1} \right), \\y_t &= x_0 r \left(1 + \frac{G^t - G}{G - 1} + m \frac{G^t - tG + t - 1}{(G - 1)^2} \right) + y_0 [1 + (t - 1)n].\end{aligned}\tag{34}$$

where

$$G = (1 - s)(1 - r) \quad (35)$$

Single wave introgression model

We now consider a single wave introgression model where introgression occurs continuously for first τ generations. Therefore, before generation τ the haplotype frequency change in the single wave model is the same as in the continuous introgression model (equation 34). After τ , haplotype frequencies change according to the single pulse model with initial frequencies x_τ and y_τ (section 1.1, equation 5). Therefore we can write

$$x_t = \begin{cases} x_0 \left(G^t + m \frac{G^t - G}{G - 1} \right), & \text{if } t < \tau. \\ x_\tau G^{t - \tau}, & \text{if } \tau < t, \end{cases} \quad (36)$$

and

$$y_t = \begin{cases} rx_0 \left(1 + \frac{G^t - G}{G - 1} + m \frac{G^t - tG + t - 1}{(G - 1)^2} \right) + y_0 [1 + (t - 1)n], & \text{if } t < \tau, \\ rx_\tau \left(\frac{1 - G^{t - \tau}}{1 - G} \right) + y_\tau, & \text{if } \tau < t. \end{cases} \quad (37)$$

Where

$$\begin{aligned} x_\tau &= x_0 \left(G^\tau + m \frac{G^\tau - G}{G - 1} \right) \\ y_\tau &= x_0 r \left(1 + \frac{G^\tau - G}{G - 1} + m \frac{G^\tau - tG + \tau - 1}{(G - 1)^2} \right) + y_0 [1 + (\tau - 1)n] \\ m &= n = 1. \end{aligned} \quad (38)$$

Dual wave introgression model

In this model, introgression occurs first τ_1 generations, after which it stops until generation τ_2 , when the second wave of introgression starts. The second wave ends at generation τ_3 . Up until τ_2 this model is equal to single wave

introgression model. Therefore, the expression for x_t and y_t for $t < \tau_2$ are given by equations 36, 37 and 38 and replacing τ by τ_1 . At generation τ_2 the introgression starts again, but the initial haplotype frequencies are x_{τ_2} and y_{τ_2} rather than 0. This situation is mathematically equivalent to continuous introgression model starting from $t = 1$ (rather than $t = 0$) and with x_0 and y_0 replaced by x_{τ_2} and y_{τ_2} .

Then, the expressions for x_t and y_t can be found to be

$$x_t = \begin{cases} x_0 \left(G^t + m \frac{G^t - G}{G-1} \right), & \text{if } t < \tau_1, \\ x_{\tau_1} G^{t-\tau_1}, & \text{if } \tau_1 < t < \tau_2, \\ x_{\tau_2} \left(G^{t-\tau_2+1} + m_{\tau_2} \frac{G^{t-\tau_2+1} - G}{G-1} \right), & \text{if } \tau_2 < t < \tau_3, \\ x_{\tau_3} G^{t-\tau_3}, & \text{if } \tau_3 < t \end{cases} \quad (39)$$

and

$$y_t = \begin{cases} rx_0 \left(1 + \frac{G^t - G}{G-1} + m \frac{G^t - tG + t - 1}{(G-1)^2} \right) + y_0 [1 + (t-1)n], & \text{if } t < \tau_1, \\ rx_{\tau_1} \left(\frac{1 - G^{t-\tau_1}}{1-G} \right) + y_{\tau_1}, & \text{if } \tau_1 < t < \tau_2, \\ rx_{\tau_2} \left(1 + \frac{G^{t-\tau_2+1} - G}{G-1} + m_{\tau_2} \frac{G^{t-\tau_2+1} - (t-\tau_2+1)G + (t-\tau_2+1) - 1}{(G-1)^2} \right) + y_{\tau_2} [1 + (t-\tau_2)n_{\tau_2}], & \text{if } \tau_2 < t < \tau_3, \\ rx_{\tau_3} \left(\frac{1 - G^{t-\tau_3}}{1-G} \right) + y_{\tau_3}, & \text{if } \tau_3 < t. \end{cases} \quad (40)$$

Where:

$$\begin{aligned}
x_{\tau_1} &= x_0 \left(G^{\tau_1} + m \frac{G^{\tau_1} - G}{G - 1} \right) \\
x_{\tau_2} &= x_{\tau_1} G^{\tau_2 - \tau_1} \\
x_{\tau_3} &= x_{\tau_2} \left(G^{\tau_3 - \tau_2 + 1} + m_{\tau_2} \frac{G^{\tau_3 - \tau_2 + 1} - G}{G - 1} \right) \\
y_{\tau_1} &= r x_0 \left(1 + \frac{G^{\tau_1} - G}{G - 1} + m \frac{G^{\tau_1} - \tau_1 G + \tau_1 - 1}{(G - 1)^2} \right) + y_0 [1 + (\tau_1 - 1)n] \\
y_{\tau_2} &= r x_{\tau_1} \left(\frac{1 - G^{\tau_2 - \tau_1}}{1 - G} \right) + y_{\tau_1} \\
y_{\tau_3} &= r x_{\tau_2} \left(1 + \frac{G^{\tau_3 - \tau_2 + 1} - G}{G - 1} + m_{\tau_2} \frac{G^{\tau_3 - \tau_2 + 1} - (t - \tau_2 + 1)G + (\tau_3 - \tau_2 + 1) - 1}{(G - 1)^2} \right) \\
&\quad + y_{\tau_2} [1 + (\tau_3 - \tau_2)n_{\tau_2}] \\
m &= n = 1 \\
m_{\tau_2} &= x_0 / x_{\tau_2} \\
n_{\tau_2} &= y_0 / y_{\tau_2}
\end{aligned} \tag{41}$$

In both the pulse and wave models haplotype frequencies change at the same rate once introgression stops, and this change is determined by r and s . The difference is that the haplotype frequencies at the time when introgression stops can be different under different models. However, if we do not know the duration of introgression or the time between two introgression waves, we can always represent wave models as the single pulse model in which pulse happened more towards present. To approximate the single wave model, we replace t by $t - (\tau - \log(x_\tau/x_0)/\log(G))$ in equation 5, while to approximate the dual wave model, we replace t by $t - (\tau_3 - \log(x_{\tau_3}/x_0)/\log(G))$. The effect of this approximation is that after the time of the last introgression, the single pulse and wave models are indistinguishable (figure S1.8). That is, we have shown that if one does not know the details of introgression it is impossible to distinguish between the single pulse and wave models based only on the present day frequency of

Neanderthal alleles. Therefore, the single pulse model is a good model for our analysis, despite the fact that it is obviously an approximation.

References

- B. O. Bengtsson. The flow of genes through a genetic barrier. In P. J. Greenwood, P.H. Harvey, and M. Slatkin, editors, *Evolution – Essays in honour of John Maynard Smith*, volume 1, chapter 3, pages 31–42. Cambridge University Press, New York, NY, 1985.
- Brian Charlesworth, Magnus Nordborg, and Deborah Charlesworth. The effects of local selection, balanced polymorphism and background selection on equilibrium patterns of genetic diversity in subdivided populations. *Genet Res*, 70(2):155–174, 10 1997. ISSN 1469-5073. doi: null.
- D Petry. The effect on neutral gene flow of selection at a linked locus. *Theoretical population biology*, 23(3):300–313, 1983. ISSN 0040-5809. doi: 10.1016/0040-5809(83)90020-5.
- Simon Aeschbacher and Reinhard Bürger. The effect of linkage on establishment and survival of locally beneficial mutations. *Genetics*, 197(1):317–336, May 2014. doi: 10.1534/genetics.114.163477.
- RC Lewontin and Ken-ichi Kojima. The evolutionary dynamics of complex polymorphisms. *Evolution*, pages 458–472, 1960.
- Sean H Rice. *Evolutionary theory: mathematical and conceptual foundations*. Sinauer Associates Sunderland, 2004.
- Tim Connallon and Andrew G Clark. Antagonistic versus nonantagonistic models of balancing selection: characterizing the relative timescales and hitchhiking effects of partial selective sweeps. *Evolution*, 67(3):908–917, 2013.
- J. D. Wall, M. A. Yang, F. Jay, S. K. Kim, E. Y. Durand, L. S. Stevison, C. Gignoux, A. Woerner, M. F. Hammer, and M. Slatkin. Higher levels of

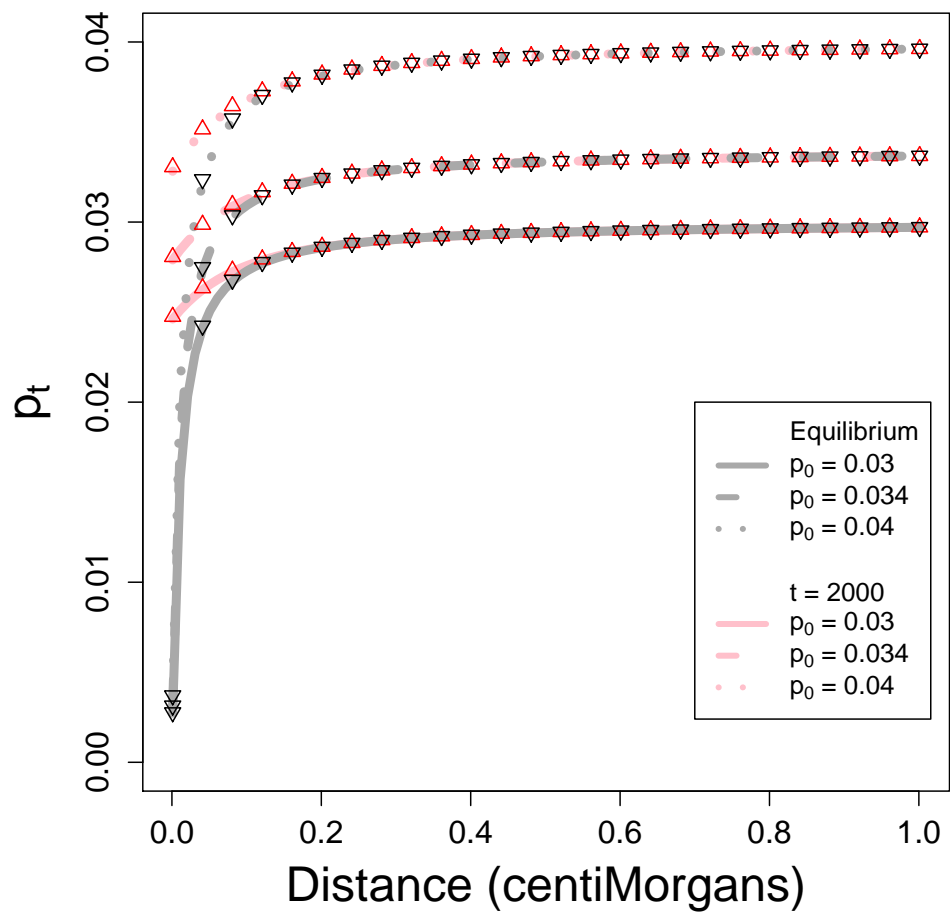


Figure S1.1: Approximate p_t as a function of the recombinational distance r . Lines represent the equation (6) for $t = 2000$ (red) and equilibrium equation (8) (grey). Numerical solution of the recursion are represented by black upward and downward facing triangles. Other parameters: $s = 0.0001$ and $y_0 = 0$ for all lines, $p_0 = 0.04$ (dotted), 0.034 (dashed) and 0.03 (full line).

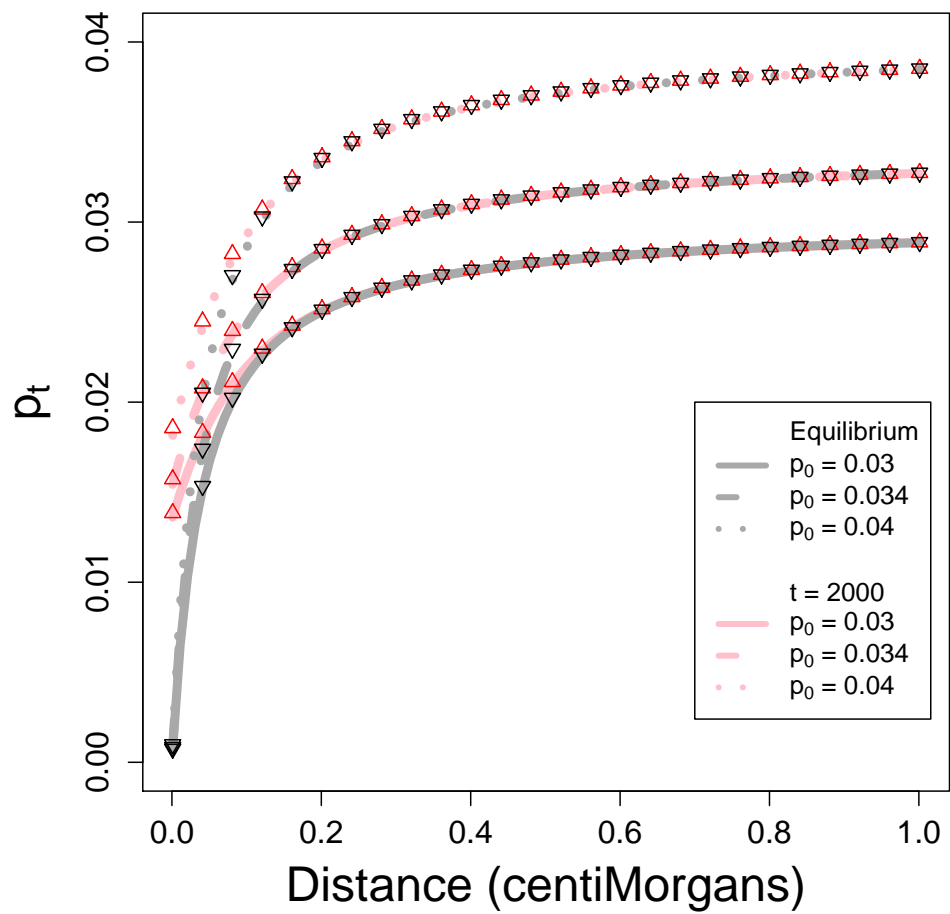


Figure S1.2: Approximate p_t as a function of the recombinational distance r . Lines represent the equation (6) for $t = 2000$ (red) and equilibrium equation (8) (grey). Numerical solution of the recursion are represented by black upward and downward facing triangles. Other parameters: $s = 0.0004$ and $y_0 = 0$ for all lines, $p_0 = 0.04$ (dotted), 0.034 (dashed) and 0.03 (full line).

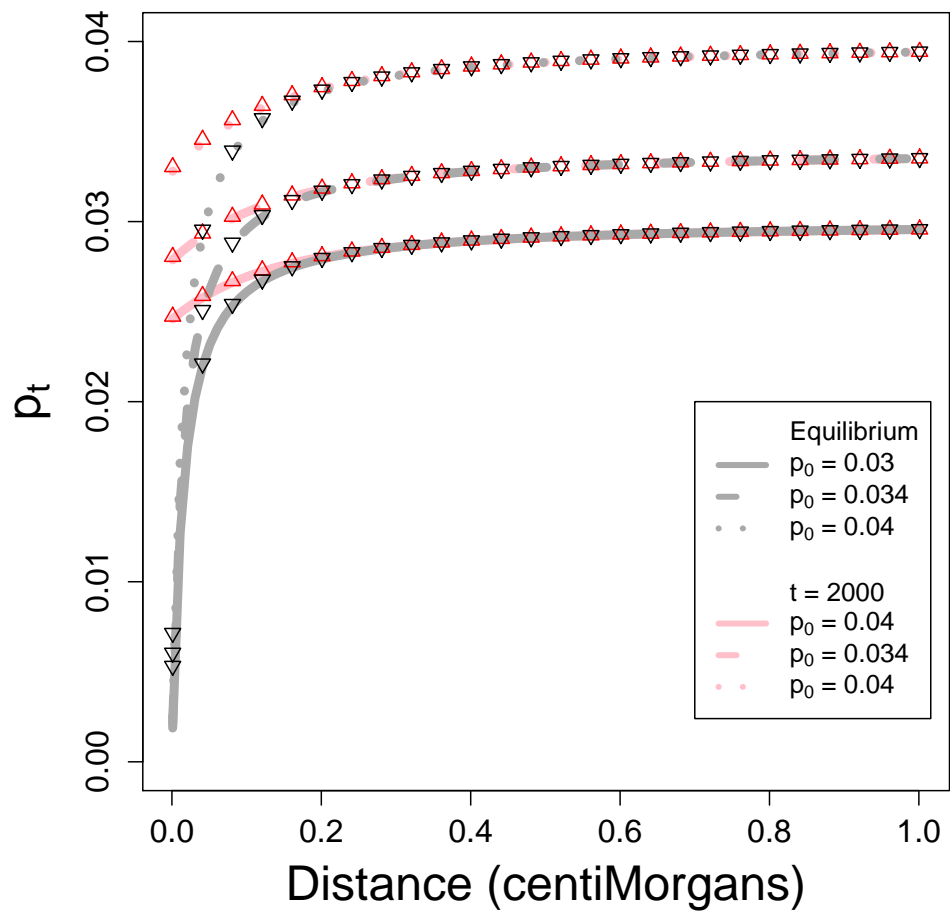


Figure S1.3: Approximate p_t as a function of the recombinational distance r . Lines represent the equation (12) for $t = 2000$ (red) and equilibrium equation (13) (grey). Numerical solution of the recursion are represented by black upward and downward facing triangles. Other parameters: $s_f = s_m = 0.0001$ and $y_{X,0} = 0$ for all lines, $p_0 = 0.04$ (dotted), 0.034 (dashed) and 0.03 (full line).

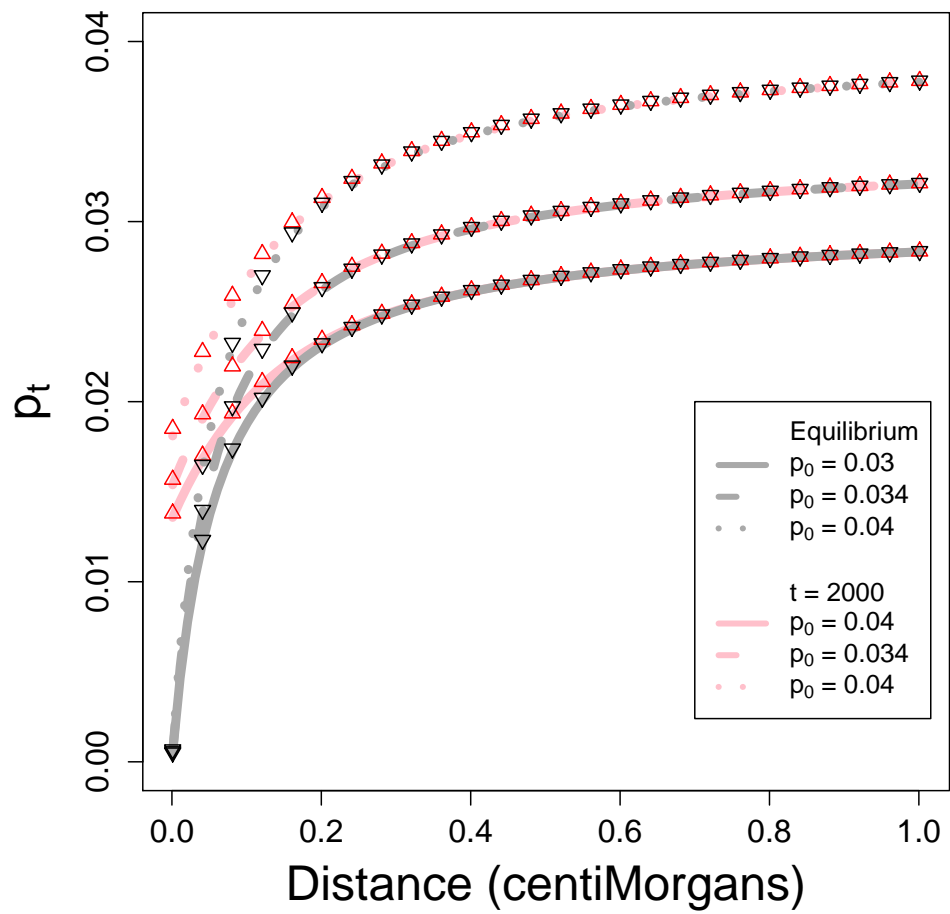


Figure S1.4: Approximate p_t as a function of the recombinational distance r . Lines represent the equation (12) for $t = 2000$ (red) and equilibrium equation (13) (grey). Numerical solution of the recursion are represented by black upward and downward facing triangles. Other parameters: $s_f = s_m = 0.0004$ and $y_{X,0} = 0$ for all lines, $p_0 = 0.04$ (dotted), 0.034 (dashed) and 0.03 (full line).

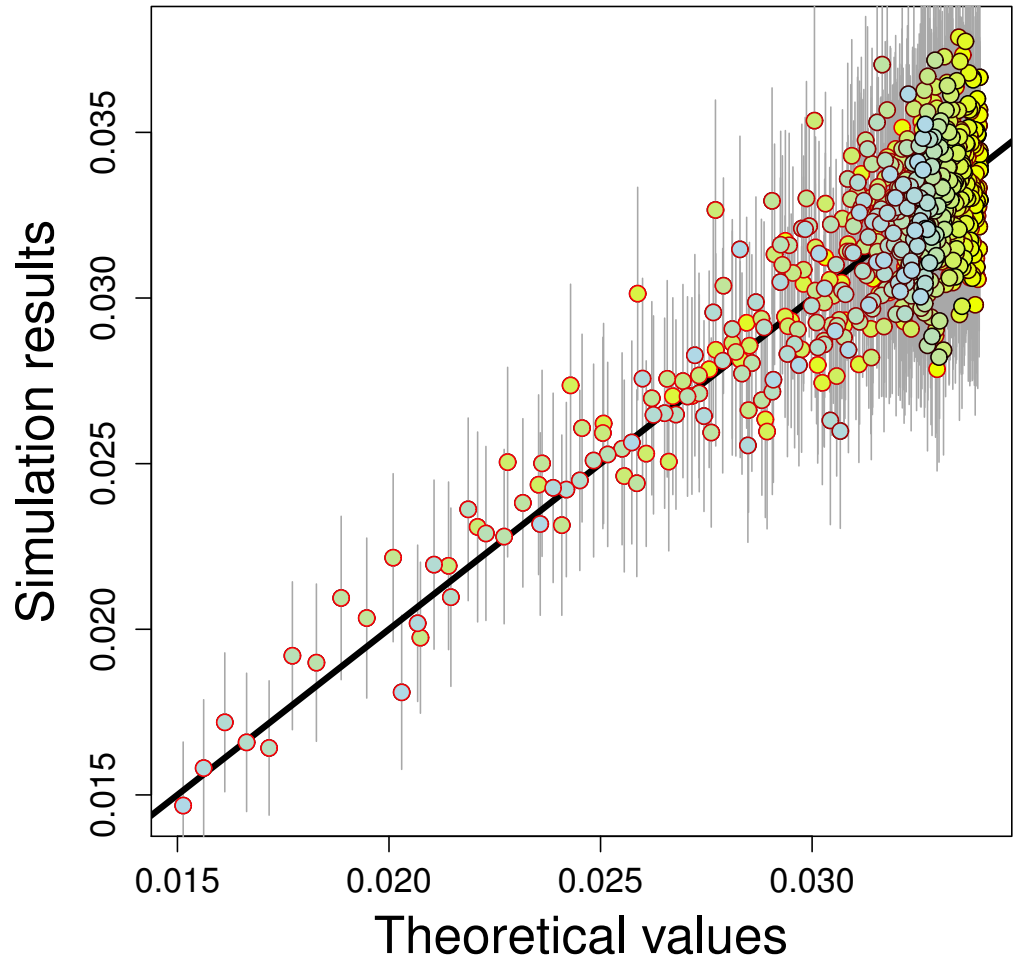


Figure S1.5: Comparison between mean p_t obtained from individual based simulations and those obtained from equation (6). Plot contain 676 circles representing different combinations of r and s . The r ranges from 1×10^{-5} (red circumference) to 1×10^{-2} (black circumference), s ranges from 1×10^{-5} (yellow circle area) to 4×10^{-4} (light blue circle area). For each parameter combination p_t was calculated based on 1000 independent runs. $t = 2000$. Grey lines represent approximate 95% confidence intervals for simulation results (mean $\pm 1.96 \cdot \text{stderror}$) Diagonal is shown with dotted black line.

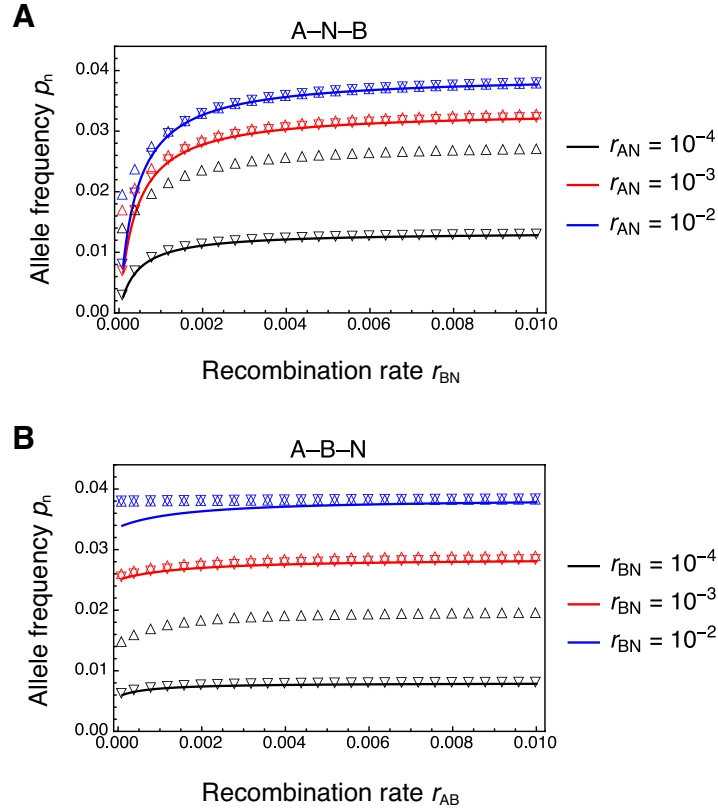


Figure S1.6: Accuracy of approximation to the frequency of a neutral allele N_1 linked to multiple autosomal loci under purifying selection. Curves show $p_{\infty, IJ}$ from equation (15) for various recombination distances between the focal neutral locus N and the two loci under selection, A and B . Upward and downward facing triangles give values obtained after iterating deterministic recursions over $t = 2000$ generations and until the equilibrium is reached, respectively. A: The neutral locus is flanked by one locus under selection on each side, and recursions followed equation (17). B: The neutral locus is flanked by two selected loci on one side and recursions followed equation (18). A, B: Selection coefficients against introgressed deleterious mutations at locus A and B are $a = 0.0002$ and $b = 0.0004$, respectively. The initial frequency of N_1 is $p_0 = 0.04$.

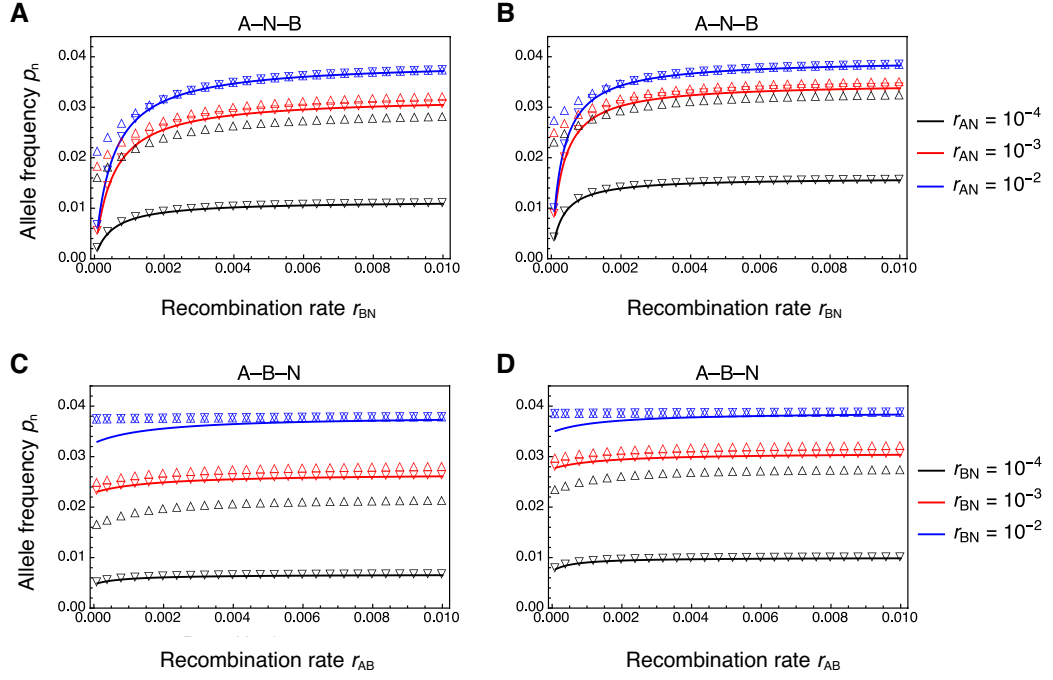


Figure S1.7: Accuracy of approximation to the frequency of a neutral allele N_1 linked to multiple X-chromosomal loci under purifying selection. Curves show $p_{\infty, X, I, J}$ from equation (21) for various recombination distances between the focal neutral locus N and the two loci under selection, A and B. Upward and downward facing triangles give values obtained after iterating equation (24) over $t = 2000$ generations and until the equilibrium is reached, respectively. A, B: The neutral locus is flanked by one locus under selection on each side. C, D: The neutral locus is flanked by two selected loci on one side. A, C: Selection coefficients against introgressed deleterious mutations at locus A and B in females (males) are $a_f = 0.0001$ ($a_m = 0.0003$) and $b_f = 0.0002$ ($b_m = 0.0006$), respectively. B, D: Selection coefficients are identical in the two sexes; $a_f = a_m = 0.0001$ and $b_f = b_m = 0.0002$. In all panels, the initial frequency of N_1 is $p_{0, X} = 0.04$.

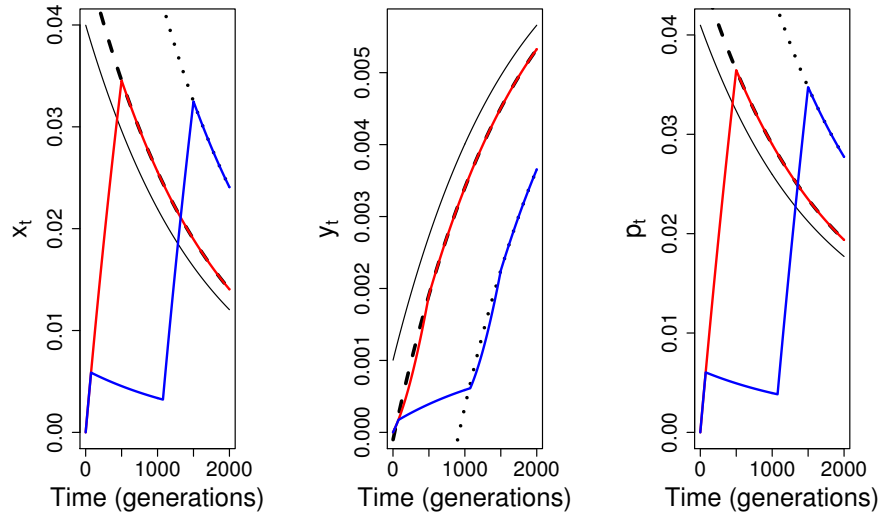


Figure S1.8: Approximating single (red line) and dual (blue line) wave models with the single pulse model. By changing time in the single pulse model (dashed and dotted black lines) we can approximate haplotype frequencies in wave models. Parameters: $r = 10^{-4}$, $s = 5 \times 10^{-4}$, $x_0 = 0.04$, $y_0 = 0.001$. Duration of admixture in single wave model, $\tau = 500$. Additional parameters for dual wave model: $\tau_1 = 75$, $\tau_2 = 1075$, $\tau_3 = 1500$. Black full line represents a single pulse model without changing time.

neanderthal ancestry in East Asians than in Europeans. *Genetics*, 194(1):199–209, May 2013.

Benjamin Vernot and Joshua M. Akey. Resurrecting Surviving Neandertal Lineages from Modern Human Genomes. *SCIENCE*, 343(6174):1017–1021, FEB 28 2014. ISSN 0036-8075. doi: 10.1126/science.1245938.

B. Vernot and J. M. Akey. Complex history of admixture between modern humans and Neandertals. *Am. J. Hum. Genet.*, 96(3):448–453, Mar 2015.

B. Y. Kim and K. E. Lohmueller. Selection and reduced population size cannot explain higher amounts of Neandertal ancestry in East Asian than in European human populations. *Am. J. Hum. Genet.*, 96(3):454–461, Mar 2015.