

The Effects of Multilocus Balancing Selection on Neutral Variability

Arcadio Navarro¹ and Nick H. Barton

Institute of Cell, Animal and Population Biology, University of Edinburgh, Edinburgh EH9 3JT, United Kingdom

Manuscript received October 22, 2001

Accepted for publication February 22, 2002

ABSTRACT

We studied the effect of multilocus balancing selection on neutral nucleotide variability at linked sites by simulating a model where diallelic polymorphisms are maintained at an arbitrary number of selected loci by means of symmetric overdominance. Different combinations of alleles define different genetic backgrounds that subdivide the population and strongly affect variability. Several multilocus fitness regimes with different degrees of epistasis and gametic disequilibrium are allowed. Analytical results based on a multilocus extension of the structured coalescent predict that the expected linked neutral diversity increases exponentially with the number of selected loci and can become extremely large. Our simulation results show that although variability increases with the number of genetic backgrounds that are maintained in the population, it is reduced by random fluctuations in the frequencies of those backgrounds and does not reach high levels even in very large populations. We also show that previous results on balancing selection in single-locus systems do not extend to the multilocus scenario in a straightforward way. Different patterns of linkage disequilibrium and of the frequency spectrum of neutral mutations are expected under different degrees of epistasis. Interestingly, the power to detect balancing selection using deviations from a neutral distribution of allele frequencies seems to be diminished under the fitness regime that leads to the largest increase of variability over the neutral case. This and other results are discussed in the light of data from the *Mhc*.

MORE than 50 complete genomes are currently accessible on public databases. Within the next few years, this quantity is expected to increase by at least an order of magnitude (BERNAL *et al.* 2001). Such an awesome increase in the availability of multilocus data has stimulated a growing interest in complex genetic systems and in their effects on neutral DNA variability. This question has been successfully addressed for simple single-locus systems by means of the structured coalescent, a method that traces back in time genealogies of sets of neutral genes that can be associated with different selected alleles (HUDSON 1990; NORDBORG 1997, 2001). Some results for multilocus systems have been obtained either by the use of approximations to the exact structured coalescent, as for purifying selection (CHARLESWORTH *et al.* 1993; CHARLESWORTH 1994; HUDSON and KAPLAN 1994, 1995), or, less frequently, by extensions of the structured coalescent (KELLY and WADE 2000; BARTON and NAVARRO 2002) that have been applied to the simpler case of balancing selection.

Both single-locus and multilocus results rely on the similarity between a genetically subdivided population and a spatially subdivided one. In the same way that natural populations can be spatially structured into local demes, they can also be genetically structured into diverse “genetic backgrounds.” Genetic backgrounds are

defined as combinations of variants from different selected sites in a chromosome. Thus, they can be thought of as, for example, haplotypes, defined by different combinations of selected alleles from different genes, or as alleles, defined by combinations of variants from different selected sites in a gene. The genetic structure produced by selected backgrounds influences diversity at neutral loci in an analogous way to spatial structure. Just as with fluctuating deme sizes (WHITLOCK and BARTON 1997) diversity at linked neutral loci is reduced if background frequencies fluctuate more than expected by drift either because selection eliminates deleterious variants (CHARLESWORTH *et al.* 1993; CHARLESWORTH 1994; HUDSON and KAPLAN 1994, 1995) or because it forces the fixation of advantageous alleles (KAPLAN *et al.* 1989; AQUADRO and BEGUN 1993; AQUADRO *et al.* 1994). In contrast, if some sort of balancing selection maintains genetic backgrounds at stable frequencies, then neutral diversity is enhanced (STROBECK 1983; HUDSON and KAPLAN 1988; KAPLAN *et al.* 1988; HEY 1991; STEPHAN *et al.* 1992; NORDBORG 1997), just as happens with stable geographical subdivision (NAGYLAKI 1982). The analogy is maintained in the case of fluctuating selection, because neutral variability can be either enhanced or reduced, depending on the timescale of the fluctuations relative to coalescence times (SVED 1983; WHITLOCK and BARTON 1997; BARTON 2000).

BARTON and NAVARRO (2002) used this analogy to set up a general analytical method that extends the coalescent to systems with an arbitrary number of se-

¹ Corresponding author: Institute of Cell, Animal and Population Biology, University of Edinburgh, W. Mains Rd., Edinburgh EH9 3JT, Scotland. E-mail: arcadi@holyrood.ed.ac.uk

lected loci. They followed the approach of MARUYAMA (1972, extended by NAGYLAKI 1982) for a spatially structured population and adapted it to genetical structure. Then, they applied their method to the case of multilocus balancing selection. They focused on a model where diallelic polymorphisms are maintained at equilibrium by strong overdominant selection. Their analytical results provide accurate predictions for systems formed by a small number of selected loci, but as the number of loci increases, an extremely large hitchhiking effect is predicted, with the probability of identity between two randomly chosen alleles dropping almost to zero (see Figure 2 in BARTON and NAVARRO 2002). To study this behavior, they carried out a preliminary simulation analysis and concluded that frequency fluctuations of each of the possible genetic backgrounds must be accounted for. They discussed several ways to do this and concluded that forward simulations are probably the most straightforward method. Here, we investigate when and why the analytical predictions of the extended multilocus coalescent become inaccurate, aiming to clarify the extent to which the method developed by BARTON and NAVARRO (2002) can be used in the case of balancing selection.

The main purpose of this work is to study the levels and patterns of neutral variability to be expected under different multilocus balancing selection regimes. Because epistasis is a key component of any multilocus system, we hope to suggest ways in which the study of neutral variability may allow us to distinguish between different kinds of interactions among selected loci. To do this, we have simulated the effects on linked neutral variability of balancing selection acting on a set of diallelic loci. Simulation results are compared with predictions obtained by the method proposed by BARTON and NAVARRO (2002). For the sake of simplicity we consider the case of symmetrical overdominance at each locus. Fitnesses across loci are allowed to combine additively, multiplicatively, or with different degrees of epistasis. The method of BARTON and NAVARRO (2002) is general and can be used to study nonsymmetrical cases and other forms of balancing selection, such as, for example, frequency-dependent selection. However, we study the simpler case of symmetric overdominance because we aim to focus on qualitative results that will not depend on the exact nature of selection or on the details of the model. An exploration of the parameter space allows us to find out which general patterns of neutral variability are to be expected under multilocus balancing selection.

METHODS OF COMPUTER SIMULATION

We simulate a Wright-Fisher diploid population of size N whose life cycle consists of drift, selection, and recombination. In every run, we consider chromosomes formed by a number of selected loci, each segregating

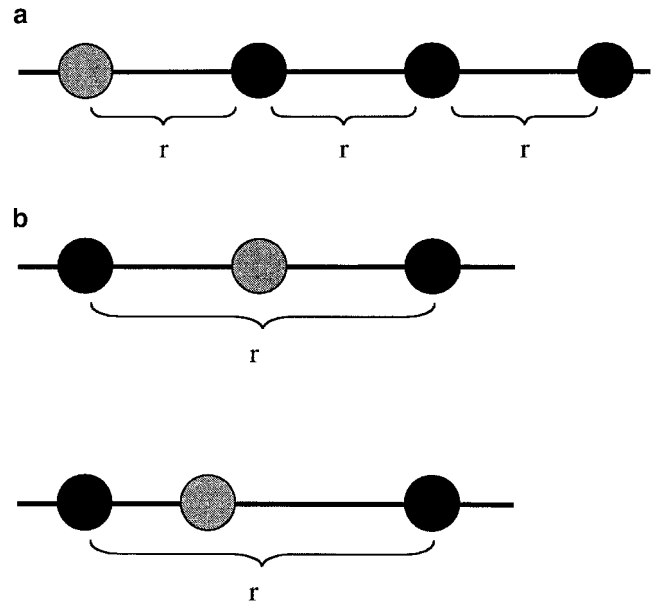


FIGURE 1.—Different possible configurations of the neutral and selected loci. (a) The neutral locus (gray) lies at the left of the set of selected loci (black). (b) The neutral locus lies at two different positions between two selected loci.

for two alleles, and a neutral locus lying among them. Genetic backgrounds are defined by combinations of alleles at different loci, so n loci produce 2^n potential backgrounds. Selected loci are assumed to be spaced at equal intervals along the genetic map, the recombination rate between any two adjacent selected loci being r (Figure 1). The neutral locus can be in any position, either at an extreme (Figure 1a) or within the set of selected loci (Figure 1b). It recombines with the selected loci, but there is no intragenic recombination. Mutations are generated at the neutral locus with rate μ immediately after recombination and according to the infinite sites model; *i.e.*, each mutation occurs at a new site. Also, it is important to stress that, all along this work, we use terms such as “locus” and “alleles” for the sake of simplicity. A genetic background is an abstract entity and can also be defined by combinations of variants within a gene, a single exon, or even a noncoding sequence.

The coalescent approach proposed by BARTON and NAVARRO (2002) considers a polymorphic equilibrium where all the possible 2^n backgrounds are present at constant frequencies and where there is no linkage disequilibrium among selected loci (hereafter, we use D , described by LEWONTIN and KOJIMA 1960, as a measure of linkage disequilibrium). Such an equilibrium can be reproduced under the n -locus symmetric viability model with negative epistasis (KARLIN and AVNI 1981; CHRISTIANSEN 1987, 1988, 2000; BARTON and SHPAK 2000). In natural populations, however, different fitness schemes will allow for different polymorphic equilibria with different background frequencies. In fact, multilocus selec-

tion and the nature of multilocus polymorphic equilibria have been the subject of much research, recently reviewed by CHRISTIANSEN (2000). The main factors influencing the degree of polymorphism and linkage disequilibrium at the selected loci are, together with recombination and population size, the strength and nature of selection. The amount of linkage disequilibrium in the system is an important issue, because it determines the actual degree of subdivision in the population. Deterministic analysis shows that all the symmetric multilocus fitness schemes that we use here have an equilibrium with $D = 0$ (where all possible 2^n haplotypes are present at equal frequencies) but in many cases this equilibrium is unstable for small r (*cf.* CHRISTIANSEN 2000). For example, when overdominant selection is operating and fitnesses are multiplicative across loci, as in Equation 2 below, selection favors linkage disequilibrium. If there is no recombination, D is maximum and the equilibrium population is formed by a pair of complementary genotypes (LEWONTIN and KOJIMA 1960; FRANKLIN and LEWONTIN 1970). In this case subdivision is simple, with only two subpopulations present. With some recombination, all possible genotypes are produced, but two predominate. Only above a threshold recombination rate does linkage equilibrium become stable (KARLIN and LIBERMAN 1979) so that all the possible backgrounds are found at even frequencies. Other multilocus fitness schemes, such as negative epistasis, allow for greater population subdivision, because they favor gametic equilibrium ($D = 0$) and, thus, the presence of a larger number of backgrounds in the population.

To investigate these scenarios, we used two different fitness functions in our simulations. First, possible interactions among loci were studied in a series of runs where fitnesses were computed according to the n -locus symmetric viability model. For simplicity, we assumed that all loci contribute equally to fitness and that selection acts only on the proportion of heterozygous loci of an individual. The fitness of an individual is given by

$$w(\alpha, h, k) = 1 + \alpha h^k, \quad (1)$$

where h is the proportion of heterozygous loci in a given individual ($0 \leq h \leq 1$) and α is the strength of selection. Epistasis enters the function by means of k , a parameter that allows for different selective regimes. Figure 2 shows three different selection schemes and two different selection strengths. If $k = 1$ (Figure 2a, straight lines), there is additive selection on heterozygosity. Selection acts on each locus individually and tends to maximize average heterozygosity. Linkage disequilibrium becomes important with epistasis, because the average fitness of the population will then depend on its variance in heterozygosity (CHRISTIANSEN 1987, 1988). If $k > 1$ (Figure 2, concave lines), there is positive epistasis on heterozygosity, so selection favors increased variance in heterozygosity and, hence, $D \neq 0$. In other words, selection favors

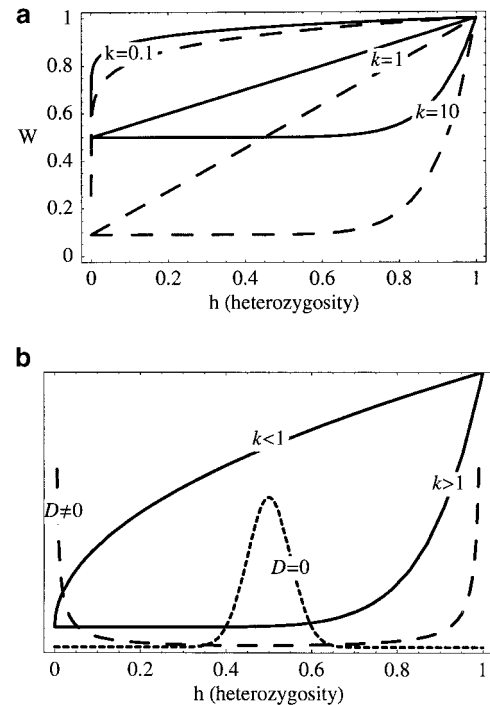


FIGURE 2.—(a) The fitness function as in Equation 1. To make values comparable, we divide by $1 + \alpha$. Solid lines, $\alpha = 1$; dashed lines, $\alpha = 10$. Straight lines, additive fitness ($k = 1$); concave lines, positive epistasis ($k > 1$); convex lines, negative epistasis ($k < 1$). (b) Two hypothetical fitness schemes and their associated degrees of linkage disequilibrium. Negative epistasis ($k < 1$) favors low variance in heterozygosity (dotted line) and, thus, low linkage disequilibrium. Positive epistasis ($k > 1$) favors high variance in heterozygosity (dashed line) and, thus, high linkage disequilibrium.

maximum linkage disequilibrium, either positive or negative, because in that case the population tends to be formed by individuals that are either completely heterozygous or completely homozygous. With $k < 1$ (Figure 2, convex lines), there is negative epistasis, so selection favors decreased variance in heterozygosity and $D = 0$ (CHRISTIANSEN 1988). Linkage disequilibrium would produce individuals with low heterozygosity, which would be eliminated by selection due to their extremely low relative fitness (see convex lines in Figure 2).

Also, to study the case of independent selective effects of each locus we ran a series of simulations in which fitnesses were multiplicative across loci (with heterozygotes having a fitness of 1 and homozygotes of $1 - s$). The fitness of an individual was given by

$$w(s, h) = (1 - s)^{i(1-h)}, \quad (2)$$

where h is again the proportion of heterozygous loci and i is the total number of loci. A multiplicative fitness scheme favors $D \neq 0$ (LEWONTIN and KOJIMA 1960; FRANKLIN and LEWONTIN 1970). For very weak selection ($s \rightarrow 0$) it converges to additive selection on heterozygosity.

For every set of parameters, an initial population was

randomly generated assuming equal allele frequencies and $D = 0$. The population was run to drift-selection equilibrium and then several diversity measures were taken for the neutral locus. The probability of identity between two randomly chosen alleles, f , was computed for the whole population. BARTON and NAVARRO (2002) used this homozygosity measure in their coalescent approach (their Equations 30 and 33) and thus it allows for straightforward comparisons between analytical and simulation results. Also, for every simulation run two more measures of genetic variability were obtained from samples of randomly chosen alleles: the mean number of nucleotide differences between pairs of sequences, d , and the number of segregating sites, S . Under a neutral model, the expectation of these two quantities is

$$d = \theta \quad (3)$$

$$S = \theta \sum_{i=1}^{n-1} \frac{1}{i} \quad (4)$$

(EWENS 1979), where $\theta = 4N\mu$ is the neutral evolution parameter and n is the sample size. The mean number of nucleotide differences is proportional to the pairwise coalescence time and identities can be used as the moment-generating function of coalescence times (HUDSON 1990). Thus, we can also compare simulated values of d with corresponding analytical predictions based on the method of BARTON and NAVARRO (2002).

Single-locus balancing selection has been shown to increase the number of heterozygotes, that is, to decrease f (STROBECK 1983; HUDSON and KAPLAN 1988; KAPLAN *et al.* 1988; HEY 1991; STEPHAN *et al.* 1992; NORDBORG 1997; SATTA *et al.* 1998; BARTON and NAVARRO 2002). This increase will affect d more than S , since d is weighted toward variants at intermediate frequencies. Such deviations in the frequency spectrum can be measured by means of the Tajima's D statistic (TAJIMA 1989; FU 1997; henceforth Tajima's D is referred to as TD to avoid confusion with linkage disequilibrium) as

$$TD = \frac{d - \theta_w}{\sqrt{\text{Var}(d - \theta_w)}}, \quad (5)$$

where

$$\theta_w = \frac{S}{\sum_{i=1}^{n-1} 1/i}. \quad (6)$$

Positive values of Tajima's D reflect a deficiency of homozygotes and are usually associated with stable population subdivision, either spatial or genetical (TAJIMA 1989; FU 1997).

Every variability measure value plotted in the figures is the average of 10 runs. The program was tested by using it to compute some well-known population genetics quantities, such as expected times for the fixation or loss of a neutral or a selected allele or patterns of

decay of gametic disequilibrium. The results obtained agreed with the literature.

RESULTS

Identities: The coalescent approach developed by BARTON and NAVARRO (2002) relies on two related assumptions: that every background is abundant and that its frequency is both known and constant. As discussed in BARTON and NAVARRO (2002), these assumptions will be valid for strong and constant selection, but unreliable predictions are made when many selected loci are considered. To find out why and under which parameter values the theory breaks down, and to explore a wider parameter space, we simulated a multilocus system as described in METHODS OF COMPUTER SIMULATION and we compared our simulation results with the analytical predictions of BARTON and NAVARRO (2002).

Figures 3 and 4 show the way in which the probability of identity changes at a neutral locus located at the extreme of a set of selected loci, as the number of selected loci increases, for different selective regimes and population sizes. Analytical predictions obtained from Equations 30 and 33 in BARTON and NAVARRO (2002) are also shown. For multilocus systems where $k < 1$ (negative epistasis), the coalescent-based predictions hold quite well when compared with the simulations (Figure 3, a and b). In this case, negative epistasis favors low variance in heterozygosity and therefore selection makes the system tend to $D = 0$. In other words, when the number of selected loci is small, the simulations meet the analytical assumptions because selection causes all the possible backgrounds to be present at even and constant frequencies; the population is as subdivided as the theory assumes. As expected, the analytical and simulation results start to diverge when the number of loci in the system becomes too large. The smaller the population, the fewer selected loci are needed for identities to diverge.

When $k = 1$ (additive selection, Figure 3, c and d) or $k > 1$ (positive epistasis, Figure 3e), multilocus balancing selection fails to boost variability beyond the effect of a single diallelic locus. A similar divergence between multilocus analytical predictions and simulation results is registered if fitnesses are multiplicative across loci (compare Figure 4, a and c, with low recombination). If selection is additive ($k = 1$) this discrepancy is due to drift. The number of possible backgrounds increases exponentially with the number of loci and, as discussed by BARTON and NAVARRO (2002), it may be too large for all of them to be simultaneously present in a finite population. However, this is certainly not true in Figure 3, where the number of possible backgrounds is between 2 (one locus) and 512 (nine loci) and the population size is either 10^3 or 10^4 . Although additive fitness precludes alleles to be lost and, indeed, maintains allelic frequencies very close to deterministic expecta-

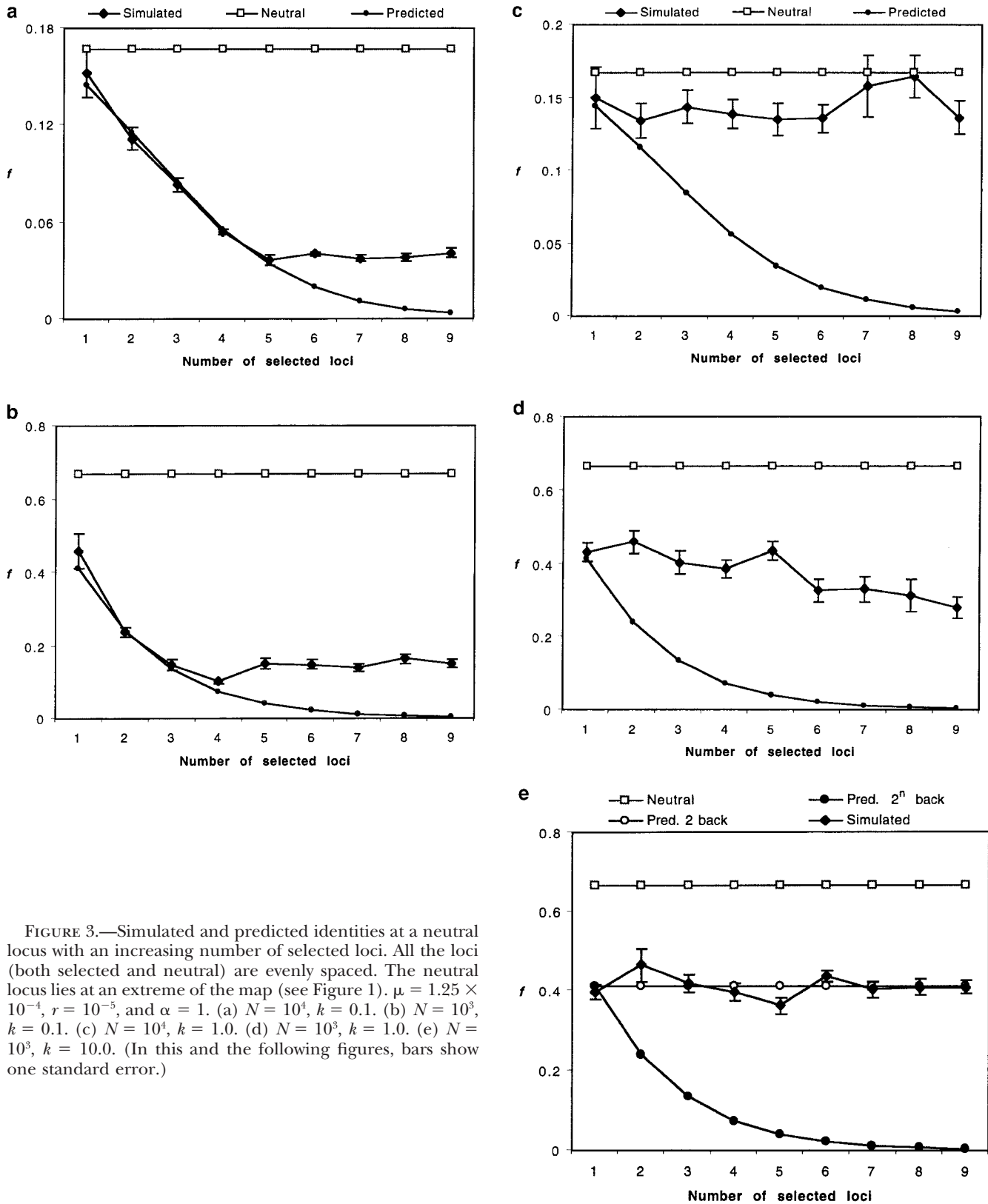


FIGURE 3.—Simulated and predicted identities at a neutral locus with an increasing number of selected loci. All the loci (both selected and neutral) are evenly spaced. The neutral locus lies at an extreme of the map (see Figure 1). $\mu = 1.25 \times 10^{-4}$, $r = 10^{-5}$, and $\alpha = 1$. (a) $N = 10^4$, $k = 0.1$. (b) $N = 10^3$, $k = 0.1$. (c) $N = 10^4$, $k = 1.0$. (d) $N = 10^3$, $k = 1.0$. (e) $N = 10^3$, $k = 10.0$. (In this and the following figures, bars show one standard error.)

tions at every single locus, selection does not act strongly on the level of linkage disequilibrium if it is $\sim D = 0$. Therefore, unless N is very large and/or selection is strong enough to make drift negligible, additive selec-

tion fails to maintain a large number of backgrounds at stable frequencies. In a finite population, drift generates random associations among selected loci and, hence, the assumptions of the extended coalescent of BARTON

and NAVARRO (2002) are not met and the population will not be as subdivided as the theory assumes. Moreover, if the number of possible backgrounds is large, drift will allow background frequencies to undergo strong fluctuations and variability will be lost. Figure 5 shows some of these fluctuations. As can be seen, weak selection and a high number of selected loci allow strong

fluctuations. Some backgrounds may be temporarily lost (Figure 5, a and b), but they can eventually be restored by recombination (Figure 5b) because all the necessary alleles are segregating in the population. In contrast, neutral variability associated with a lost background can be restored by mutation only. The simulation results fit the analytical assumptions when the number of backgrounds is small enough for them to be maintained at roughly constant frequencies (Figure 5c).

When epistasis is positive ($k > 1$, Figure 3e) or fitnesses are multiplicative (Figure 4c), there is also a discrepancy between simulated and theoretically predicted results, but this is only apparent. As we have already mentioned, under such fitness schemes selection favors linkage disequilibrium among selected loci, so that the equilibrium population is dominated by two complementary genotypes (FRANKLIN and LEWONTIN 1970). The results shown in Figures 3e and 4c (with low recombination) become clear if one considers that the population depicted there consists essentially of two complementary backgrounds, independently of the number of selected loci. Selection tends to maintain the two dominant backgrounds at even and constant frequencies. Even though the rest of the backgrounds are continuously produced by recombination, they are kept at low frequencies and eventually lost. Thus, there is no discrepancy between analytical predictions and simulations: one needs only to recalculate the coalescent predictions, taking into account the real background equilibrium frequencies (Figures 3e and 4a). When recombination is low, nondominant backgrounds have such low frequencies that most of the time they are absent from the population and, thus, can be ignored. With intermediate or high recombination one needs to take into account the exact equilibrium frequencies in the coalescent, so the helpful assumption made by BARTON and NAVARRO (2002) of no linkage disequilibrium between selected loci cannot be used. Still, if the number of selected loci is small the extended coalescent can be applied and the right identities obtained after some tedious but feasible algebra.

Identities and recombination: The results in Figure 4 suggest that, independently of the number of loci and the selection regime, neutral variability does not increase for markers at an extreme of the set of selected

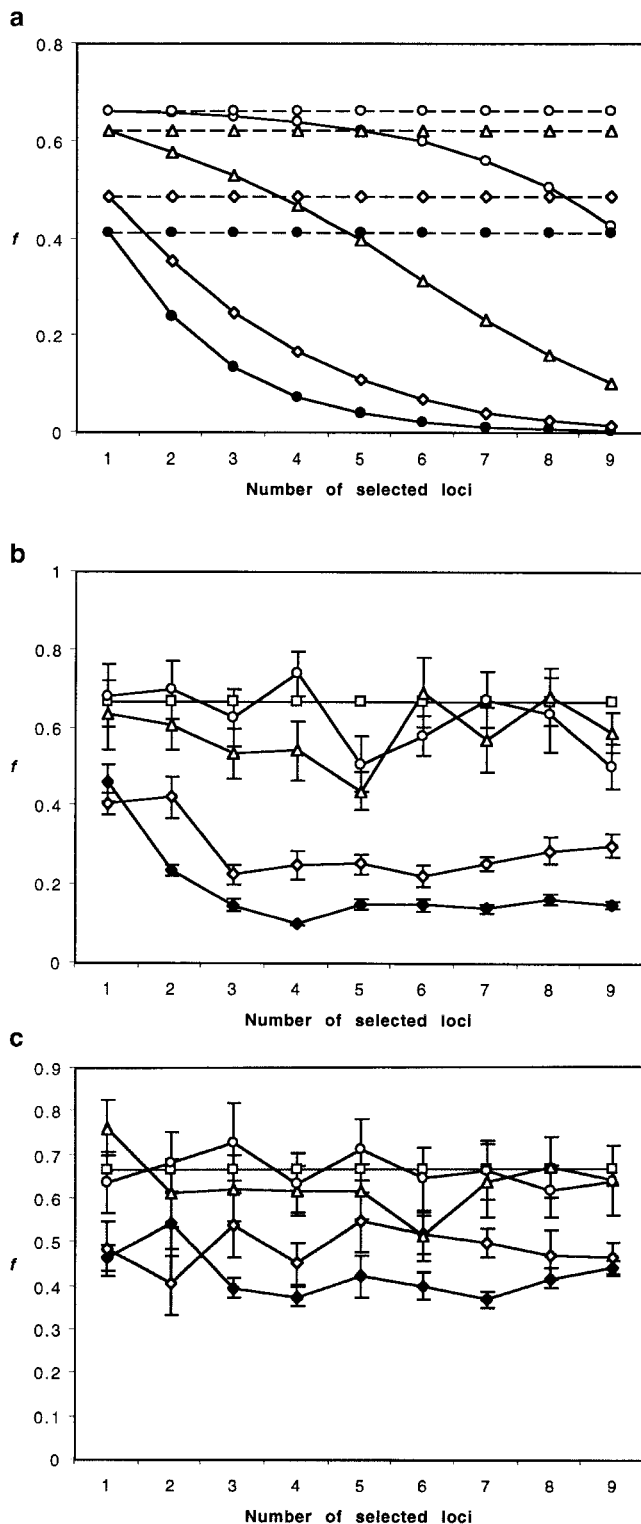


FIGURE 4.—Simulated and predicted identities at a neutral locus with an increasing number of selected loci. All the loci (both selected and neutral) are evenly spaced. The neutral locus lies at an extreme of the map (see Figure 1). $N = 10^3$, $\mu = 1.25 \times 10^{-4}$. (a) Predicted values with several backgrounds (solid lines) or only two backgrounds (dashed lines). (\bullet) $r = 10^{-5}$. (\diamond) $r = 10^{-4}$. (Δ) $r = 10^{-3}$. (\circ) $r = 10^{-2}$. (b) Negative epistasis, $\alpha = 1$ and $k = 0.1$. (\square) Neutral. (\blacklozenge) $r = 10^{-5}$. (\diamond) $r = 10^{-4}$. (Δ) $r = 10^{-3}$. (\circ) $r = 10^{-2}$. (c) Multiplicative fitness, $s = 0.01$. (\square) Neutral. (\blacklozenge) $r = 10^{-5}$. (\diamond) $r = 10^{-4}$. (Δ) $r = 10^{-3}$. (\circ) $r = 10^{-2}$.

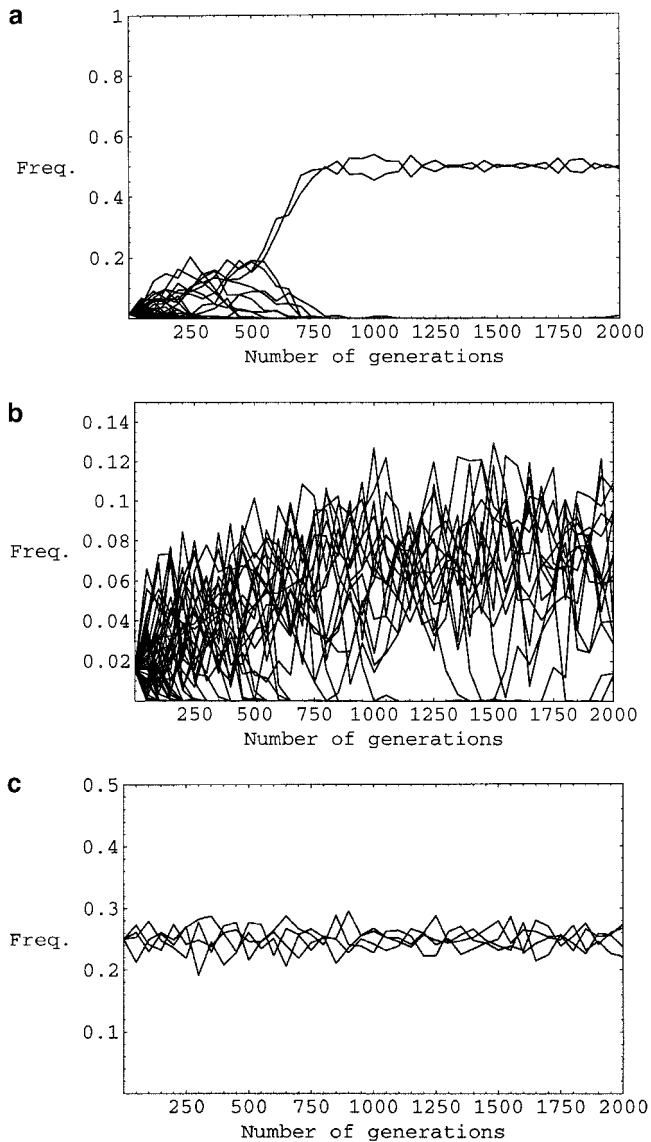


FIGURE 5.—Background frequency fluctuations under different parameter values. Each line stands for the frequency on one of the possible backgrounds. Selected loci are evenly spaced ($r = 10^{-5}$) and $\alpha = 1$. (a) $N = 10^3$, $k = 1$, six loci. (b) $N = 10^3$, $k = 0.1$, six loci. (c) $N = 10^3$, $k = 0.1$, two loci. (Note the different scales.)

loci when $Nr > 1$ ($r > 10^{-3}$ in the figure) between the neutral marker and the closest selected locus. The effect of recombination is further studied in Figures 6–8. Figure 6 plots the identities for a neutral locus lying at the extreme of a group of selected loci (either two or seven) under negative epistasis and for different recombination values. As expected, simulated and predicted results diverge when the number of loci is large, but, independently of this discrepancy, the effect of the set of selected loci on neutral variability dissipates when recombination between them and the neutral locus is $> 1/N$ ($Nr > 1$). The same threshold is found for other population sizes (not shown).

In all the results presented up to now, the neutral

locus lies at an extreme of the map. We now study the chromosomal regions between selected loci. Figures 7 and 8 show theoretically predicted and simulated identities for a neutral locus at different relative positions between two selected loci. Predicted identities for the single-selected locus case are also shown in Figure 7. Selected loci are at positions 0 and 1 and the recombination fraction between them can change. With negative epistasis and intermediate or low recombination ($r \leq 10^{-3}$, Figure 7a), maximum neutral variability (which is always found in regions closely linked to the selected loci) is increased beyond the one-locus (two backgrounds) limit, because the population is more subdivided with two loci (four backgrounds) than with a single one. Variability is increased in a wider region of the chromosome than expected for a single locus alone. The reason is simple: a higher level of variability is attained and it decays over a longer distance. Still, high recombination rates ($r > 10^{-3}$) preclude any relevant multilocus effect and the predicted values are almost identical for one as for two selected loci (Figure 7b).

Results for multiplicative fitnesses are shown in Figure 8. Results for positive epistasis (not shown) are qualitatively equivalent. As expected, when recombination is high and selection is moderate (Figure 8a with $r = 10^{-2}$) only the individual effect of each locus is relevant. If recombination is low (Figure 8a with $r = 10^{-4}$) maximum neutral variability near selected loci is still the same as with high recombination. Variability is not increased beyond the one-locus limit because the population is dominated by only two backgrounds. Moving away from the extremes of the set of selected loci, variability decays at the same rate as in the single-locus case (not shown). In segments between selected loci, however, a second kind of multilocus effect is detected. Variability is increased over a much larger region than expected for a single locus. This is due to the fact that selection generates linkage disequilibrium between selected loci. Crossing over between the two selected loci, which would allow neutral alleles to recombine away, breaks linkage disequilibrium and generates gametes that are eliminated by selection. Thus, the effective recombination rate is reduced in regions between selected loci and, if selection is strong enough, variability can be increased even when recombination is high (Figure 8b with $r = 10^{-2}$). Simulated and analytical results fit quite well because coalescent predictions can be calculated taking into account the exact haplotype frequencies at equilibrium. Just as previously shown in Figure 3e and 4a, the extended coalescent needs only information about the frequencies of the selectively relevant haplotypes and not about the kind of selection producing them.

Coalescence times and the frequency spectrum: To complement the information provided by identities, we considered two more variability measures: d and S (see METHODS OF COMPUTER SIMULATION). Analytical predic-

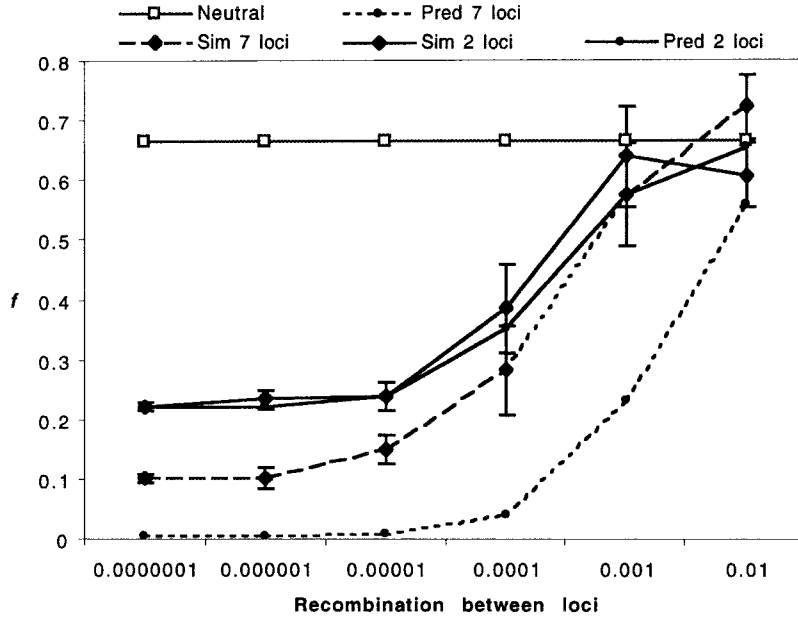


FIGURE 6.—Simulated and predicted identities under different recombination frequencies. All the loci are evenly spaced. The neutral locus lies at an extreme of the map (see Figure 1). $N = 10^3$, $\mu = 1.25 \times 10^{-4}$. Negative epistasis is shown with $\alpha = 1$, $k = 0.1$.

tions for d can be obtained by using the identities provided by BARTON and NAVARRO (2002) as the moment-generating function of the distribution of pairwise coalescence times. The average pairwise coalescence time, $E(T)$, is given by the first moment of the distribution taken on $\mu = 0$:

$$E(T) = - \left[\frac{\partial f}{\partial \mu} \right]_{\mu=0}. \tag{7}$$

BARTON and NAVARRO (2002) show that, if allele frequencies are even at each locus and if μ , r , $N^{-1} \ll 1$, the average identity between gametes chosen at random from the population, f , is

$$f = \frac{1}{4N\mu(1 + \sum_U 1/(4N\mu + 4N\rho_U))} = \frac{1}{(1 + 4N\mu + \sum_{U \neq \emptyset} \mu/(\mu + \rho_U))}, \tag{8}$$

where U are all the possible sets of loci and ρ_U are the total recombination rates between the limits of any given set (for details, see Equations 26–33 in BARTON and NAVARRO 2002). From this equation, the average pairwise coalescence times are

$$E(T) = 2N + \sum_{U \neq \emptyset} \frac{1}{2\rho_U}. \tag{9}$$

The mean number of nucleotide differences between pairs of randomly chosen alleles is given by

$$E(d) = 2\mu E(T) = 4N\mu + \sum_{U \neq \emptyset} \frac{\mu}{\rho_U}. \tag{10}$$

This shows that, when $r \rightarrow 0$, $d \rightarrow \infty$. The reason is that under an infinite sites model with no recombination, different genetic backgrounds eventually accumulate infinite differences. Variances can be calculated using the same method.

Figure 9 shows changes in the values of d and S as the number of selected loci increases, for different selective regimes. In general, these two summary statistics behave as expected. Positive epistasis ($k > 1$) increases variability, but not beyond the one-locus two-backgrounds limit, so analytical results and simulations coincide independently of the number of selected loci. With negative epistasis ($k < 1$), many backgrounds are maintained by selection and neutral variability is increased beyond the single-locus limit, even though these backgrounds undergo strong fluctuations (Figure 5). In this case simulations fit the analytical predictions until the number of selected loci is large enough for background frequency fluctuations to sweep variability away. When the number of selected loci is large, fluctuations can be so strong that the timescale of background loss and recovery becomes very small. In that case, d and S can be smaller in systems with a large number of loci than when the number of loci is intermediate (compare, for example, d values for three loci with values for eight or nine loci in Figure 9a). Note that although d and S become quite low, the corresponding identities do not undergo such a radical change (Figure 3) because, first, they can fluctuate only between 0 and 1 and, second, they converge very quickly to their equilibrium values (see below).

By comparing Figure 9a and 9b, it can be seen that d and S increase at different rates with increasing number of backgrounds. These different rates result in a change in the shape of the frequency spectrum with increasing number of loci. When there is positive epistasis ($k > 1$) Tajima's D is highly positive and remains so with increasing number of loci (Figure 10). In this case, balancing selection makes the frequency spectrum even and an excess of variants at average frequencies is detected, just as expected from previous results for the

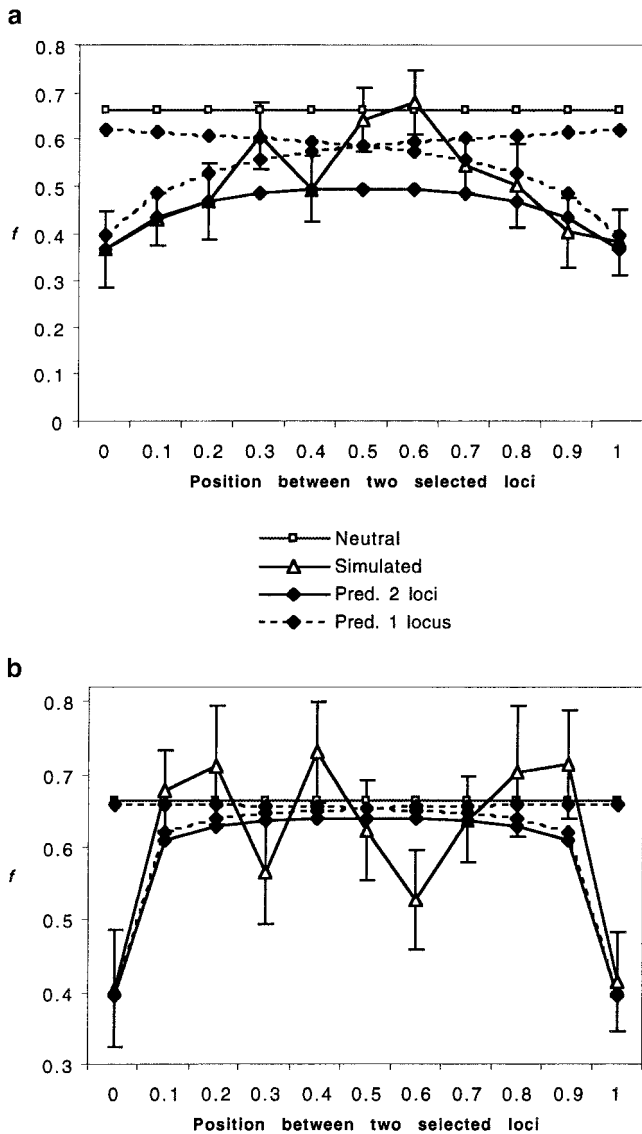


FIGURE 7.—Simulated and predicted identities at a neutral locus between two selected loci. The neutral locus lies at different relative positions between two selected loci that are at distance r from each other. $N = 10^3$, $\mu = 1.25 \times 10^{-4}$. Negative epistasis is shown with $\alpha = 1$, $k = 0.1$. (a) $r = 10^{-3}$. (b) $r = 10^{-2}$. (Note the different scales.)

one-locus two-background case. With negative epistasis ($k < 1$) the situation is different. Figure 10 shows that Tajima's D gets less positive with increasing number of loci. Paradoxically, although variability is maximized with negative epistasis, the power to detect deviations from a neutral distribution of allele frequencies is decreased. There are two reasons for this behavior. First, with so many neutral variants segregating in the population, it becomes easy to find several low frequency alleles in a sample. Second, background frequency fluctuations act as selective sweeps or, in terms of the analogy with spatially subdivided populations, as extinction-recolonization events. Such events are known to make Tajima's

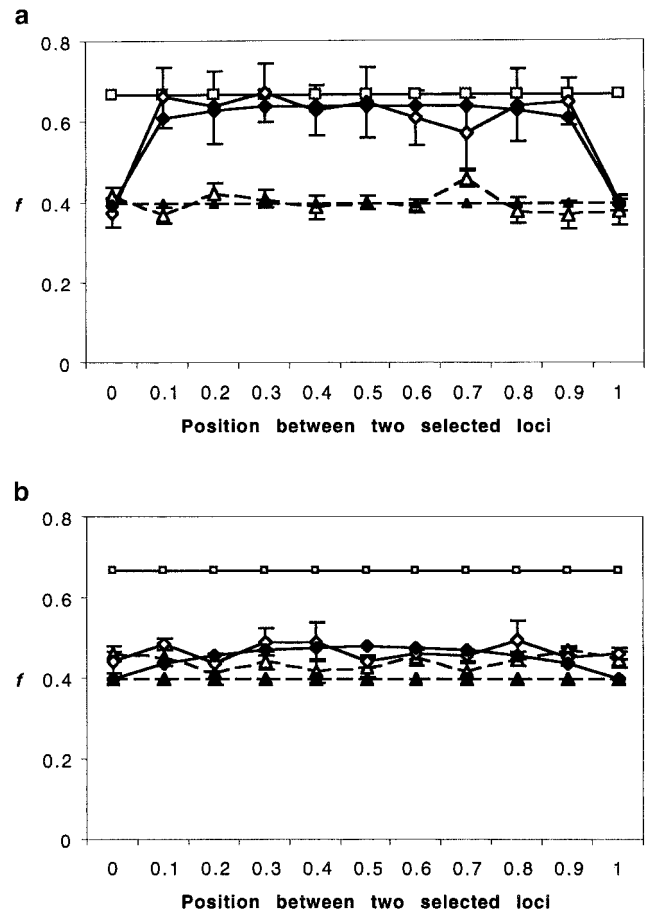


FIGURE 8.—Simulated and predicted identities at a neutral locus between two selected loci. (\square) Neutral, (Δ) simulated $r = 1^{-4}$, (\blacktriangle) predicted $r = 1^{-4}$, (\diamond) simulated $r = 1^{-2}$, (\blacklozenge) predicted $r = 1^{-2}$. The neutral locus lies at different relative positions between two selected loci that are at distance r from each other. $N = 10^3$, $\mu = 1.25 \times 10^{-4}$. Multiplicative fitness is shown with (a) $s = 0.1$ and (b) $s = 0.9$.

D negative because low frequency variants tend to accumulate while mutation restores variability in a recently lost and recovered background.

DISCUSSION

The multilocus coalescent and balancing selection: We have shown that in a multilocus system the ways in which selected loci interact are factors as important as population size, recombination, or the strength of selection, because they determine a key factor: the extent to which the population is subdivided, *i.e.*, the number of genetic backgrounds that are maintained in the population. Different kinds of epistasis allow for different degrees of population subdivision and, thus, for different degrees of diversity, both selected and neutral. In general, balancing selection acting on groups of loci generates two related kinds of multilocus effects. First, variability at sites closely linked to each of the selected

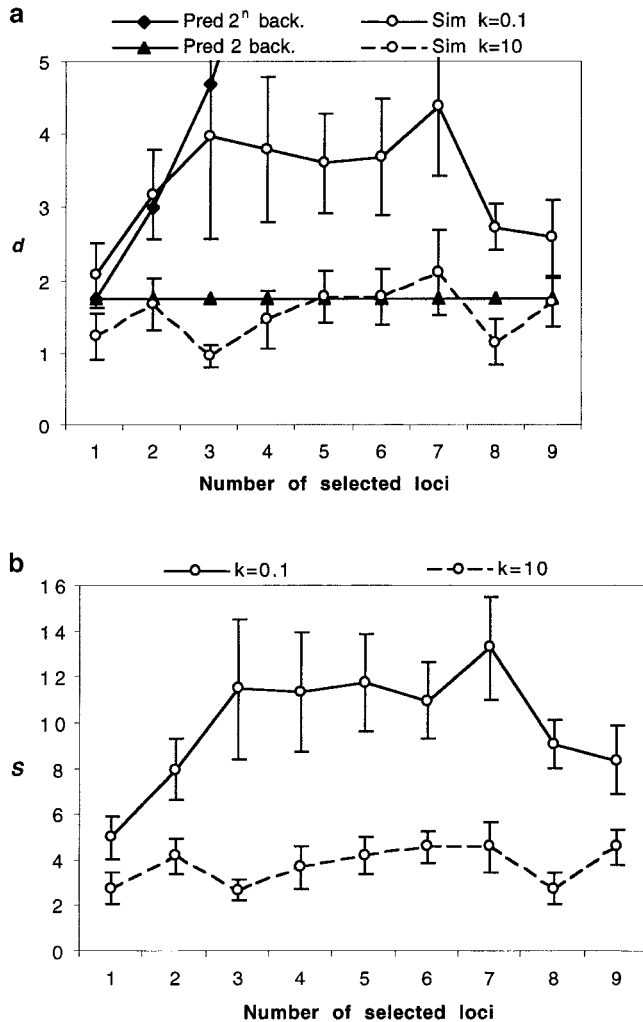


FIGURE 9.— d and S values at a neutral locus with an increasing number of selected loci and different selective regimes. All the loci (both selected and neutral) are evenly spaced. The neutral locus lies at an extreme of the map (see Figure 1). $N = 10^3$, $\mu = 1.25 \times 10^{-4}$, $r = 10^{-4}$, $\alpha = 1$, sample size $n = 10$. (a) Simulated and predicted values of the average number of pairwise differences. (b) Simulated values of the number of segregating sites.

loci can be enhanced beyond the expectations for a single-locus system. This effect is produced when the population is highly subdivided, that is, when a large number of backgrounds are maintained. In our model, this is achieved under negative epistasis (Figures 3, a and b, for example). This variability increment dissipates quickly as the neutral locus moves away from the set of selected loci and completely disappears when $Nr > 1$ (Figures 6 and 7). Second, when multiple selected loci are involved, the variability enhancement extends to a larger section of the map. With negative epistasis this extension affects all neutral variability linked to a set of selected loci and is a trivial consequence of the higher levels of variability reached near each of the selected loci (Figures 6 and 7a). In contrast, with multiplicative

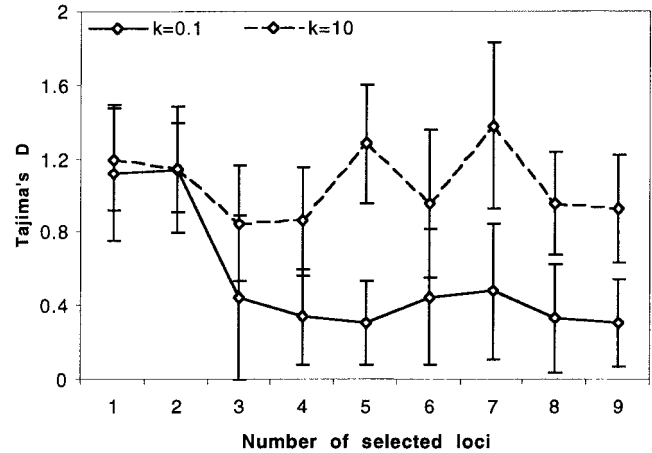


FIGURE 10.—Simulated Tajima's D values at a neutral locus with an increasing number of selected loci. All the loci (both selected and neutral) are evenly spaced. The neutral locus lies at an extreme of the map (see Figure 1). $N = 10^3$, $\mu = 1.25 \times 10^{-4}$, $\alpha = 1$, $r = 10^{-4}$, sample size $n = 10$.

fitness or positive epistasis (Figures 7b and 8), the extension is directly due to selection and affects only variability in regions between selected loci. Under the latter selective regimes, the population reaches an equilibrium in which there is maximum linkage disequilibrium and only two backgrounds dominate. As a consequence, selection opposes the homogenizing effect of recombination because crossing over between two selected loci produces unfit gametes that selection tends to eliminate. The rate of decay of molecular diversity remains the same as in the single-locus two-backgrounds case for neutral markers outside the set of selected loci. In contrast, diversity is enhanced at neutral markers located in regions among the set, even if they are a long way from the selected loci themselves (see Figure 8 with $r = 10^{-2}$, $Nr = 10$). This mechanism has been already described by KELLY and WADE (2000), who analyzed a system formed by two diallelic, epistatically interacting loci and obtained predictions for the patterns of neutral sequence variation linked to them. Their model is an instance of the positive epistasis scenario presented here.

Under negative epistasis, simulations show that, when the number of backgrounds is large relative to the population size, variability stops increasing with the addition of more selected loci (Figures 3a, 3b, and 4b). It is clear that the multilocus coalescent expectations of a huge increase in variability (BARTON and NAVARRO 2002) are not fulfilled by the simulations. As discussed in BARTON and NAVARRO (2002), the causes of deviations from the coalescent predictions are fluctuations in background frequencies (Figure 5). These fluctuations are not taken into account by the extended coalescent, which considers only drift within each of the subpopulations defined by the backgrounds. Although allele frequencies at ev-

ery locus are maintained very close to their equilibrium values, strong frequency fluctuations eliminate backgrounds and generate linkage disequilibrium between selected loci, thus violating several of the assumptions of the analytical study. The two parameters implied in the fitness scheme, α and k , together with N , r , and the number of loci in the system determine how many backgrounds can be maintained in the population and how stable their frequencies will be. As N and α increase, and as k decreases, the more subdivision and the more variability the population can harbor. On the other hand, the larger the number of loci, the less intense will be selection on each individual locus, and because the number of backgrounds grows exponentially with the number of loci, selection will be even weaker on each individual background. For example, in Figure 3, a and b, selection is still very efficient in maintaining allelic frequencies when the number of selected loci is greater than five, but it cannot keep haplotypic frequencies stable against drift and, thus, background frequencies fluctuate.

The nature of multilocus balancing selection equilibria has been studied in detail (*cf.* CHRISTIANSEN 2000), and, ideally, one could hope to find a statistic that, given α , k , N , and r , would summarize the amount of genetic subdivision in an equilibrium population and use it to predict the amount of neutral variability that this population can sustain. An obvious candidate for such a statistic is the effective number of selected backgrounds, n_e , which can be defined by analogy to the effective number of alleles, $n_e = 1/\sum p_i^2$, where p_i are the frequencies of the i backgrounds present in the population (CROW and KIMURA 1970, p. 324). We averaged n_e over the 1000 generations previous to taking the identity measures in Figure 3, b and c. In the case of five selected loci with negative epistasis (Figure 3b, $f = 0.15$), the effective number of backgrounds is ~ 13 . With the same number of selected loci but under an additive fitness scheme, n_e is only ~ 4 (Figure 3c, $f = 0.43$). Variability, therefore, seems to increase with the effective number of backgrounds. Yet, the relationship is far from simple. For example, $n_e \sim 4$ is achieved with two selected loci in a system under negative epistasis (Figure 3b), but in this case the probability of identity ($f = 0.24$) is only one-half as big as the corresponding 4-background value in the additivity case. This implies higher variability with the same effective number of backgrounds. This is due to the fact that fluctuations are much stronger with five selected loci and additivity than with two selected loci and negative epistasis. The degree of population subdivision may be similar in both cases, but n_e contains no information about fluctuations, which are what really matters. Unfortunately, as discussed in BARTON and NAVARRO (2002), it is difficult to account for them analytically. The extended coalescent approach proposed by BARTON and NAVARRO (2002) will, therefore, be a useful tool as long as the degree of genetic subdivision

of the population is known and stable. This implies knowledge of the targets of selection but not necessarily of the kind of selection acting upon them.

Variability patterns associated with multilocus balancing selection: One of the main purposes of this study was to investigate the levels and patterns of neutral variability to be expected under multilocus balancing selection and, conversely, to suggest ways in which the study of neutral variability may allow us to distinguish between different types of fitness regimes. We focused on qualitative results because of the sheer complexity of the problem and because we expected that a preliminary exploration of the parameter space would allow us to detect general patterns that would not depend on the exact nature of selection or on the details of the model. We have been able to do so and found that, although multilocus balancing selection always increases variability, very different patterns are expected under different parameter values. This challenges the traditional consensus view, on the basis of results from single-locus studies, that the footprints of balancing selection are high variability, strong linkage disequilibrium, and uniformly high allele frequencies across sites (NORDBORG 1997; KELLY and WADE 2000). Our results show that this is not always the case and that predictions from single-locus studies do not extend in a straightforward way to multilocus systems.

Under positive epistasis, variability patterns resemble the intuitions provided by single-locus studies. Large "islands" of high variability and high linkage disequilibrium are expected to be found within the region spanned by the selected loci. Variability is expected to be high and roughly constant along the intervening region, even if this region is highly recombining, and to decay for markers at the extremes of the set of selected loci. In sharp contrast, with negative epistasis the effect of multilocus balancing selection depends strongly on recombination. In that case, multilocus balancing selection tends to produce peaks of higher variability around each of the selected loci (not shown) and there is a threshold of recombination beyond which no effect on neutral variability is expected ($Nr > 1$). Actually, the multilocus threshold coincides with the single-locus one, as previously described (HUDSON and KAPLAN 1988; SATTÀ 1997; TAKAHATA and SATTÀ 1998). A particularly interesting result of our simulation study concerns frequency spectrum. With negative epistasis and an increasing number of selected loci, variability increases (up to a certain point) whereas Tajima's D gets closer to neutrality (compare Figures 9a and 10). The reasons for this behavior are clear: the amount of variability segregating in the population is so great that some low-frequency variants are represented a single time in our sample, and background frequency fluctuations act as selective sweeps. Nevertheless, this raises an unexpected paradox: the parameter values that allow balancing selection to produce a maximization of neu-

tral variability make it more difficult to detect selection out of the evenness of the frequency spectrum. This effect is even greater if the population is not yet at equilibrium (results not shown). In contrast with identities, which quickly converge to their equilibrium values, the average number of pairwise differences, d , can take a very long time to reach its equilibrium value, because mutations keep accumulating at differentiating backgrounds even when the probability of identity between them is close to zero. Equations 9 and 10 show that in a population structured by a large number of backgrounds, when recombination is low, d is very high and, thus, the time needed for mutation to generate all that variability is enormous.

Data on multilocus balancing selection: Can our results help us to understand the patterns of DNA variability found in natural populations? An initial point can be made on the issue of the overall amount of balancing selection in the genome. Even ignoring the important problems of the expected distribution of epistatic effects and of the possibility of unequal contributions to fitness of different loci, it is clear that if multilocus balancing selection were common, a considerable increase of variability beyond the neutral expectations would be predicted in regions of low recombination. In fact, the opposite correlation has been found in a variety of organisms (AQUADRO *et al.* 1994; NACHMAN 1997; DVORAK *et al.* 1998; STEPHAN and LANGLEY 1998; PRZEWORSKI *et al.* 2000), which seems to rule out the possibility of an overall effect due to widespread balancing selection. The absence of these overall effects can have several causes. First, recombination rates across the genome might be too high. Second, variability-reducing forms of selection, such as purifying or positive selection, might be dominant. Finally, balancing selection may fluctuate with time, which, depending on the timescale of the fluctuations, might produce a reduction in associated neutral variability (WHITLOCK and BARTON 1997; BARTON 2000).

Nevertheless, an increasing number of cases of multi-locus balancing selection are being reported. Some examples are self-incompatibility systems (*cf.* in CHARLESWORTH and AWADALLA 1998), meiotic drive systems (LYTTLE 1991; PALOPOLI 2000); *Mhc* (*cf.* in HUGHES and YEAGER 1997, 1998; BECK and TROWSDALE 2000; MEYER and THOMSON 2001), or R-genes in plants (*cf.* in BERGELSON *et al.* 2001). The human *Mhc* multigene family (known as *HLA* in humans) is one of the best-documented examples of the action of selection maintaining polymorphism and has become a paradigm for molecular evolutionary studies. Evidence for the action of some sort of balancing selection, recently reviewed by MEYER and THOMSON (2001), comes from different sources, such as the abundance of *trans*-specific polymorphisms or a high rate of nonsynonymous over synonymous substitutions in the exons corresponding to the peptide-binding regions.

Although the *Mhc* is usually modeled as a single-locus multiple-allele system (see, for example, TAKAHATA and NEI 1990; SLATKIN and MUIRHEAD 2000), it has features that make it an ideal candidate to be studied by means of multilocus models. First, both *Mhc* alleles and haplotypes are analogous to our genetic backgrounds: There are several targets of selection in different sites of the *Mhc* genes and alleles are defined by combinations of variants of these targets (at this level an allele would be a genetic background). Correspondingly, *Mhc* haplotypes, upon which selection is acting (*cf.* BECK and TROWSDALE 2000; MEYER and THOMSON 2001), are defined by combinations of these different alleles (and, thus, at this level different genetic backgrounds would correspond to different haplotypes). Second, gene conversion and/or intragenic recombination events have been detected within *Mhc* exons (TAKAHATA and SATTA 1998) and recombination can generate both new haplotypes and new alleles (TAKAHATA and SATTA 1998). Third, there is strong linkage disequilibrium both within and between *Mhc* genes (SANCHEZ-MAZAS *et al.* 2000; *cf.* MEYER and THOMSON 2001), which suggests the existence of interactions between them. Finally, even if the previous arguments were not convincing, it is clear that a single-locus, multiple-alleles model can be approximated by the multilocus approach presented here if recombination is assumed to be absent.

SLATKIN and MUIRHEAD (2000) estimate the intensity of overdominant selection in several human *Mhc* loci from populations all over the world. Their estimations are based on a highly symmetric, single-locus multiple-alleles model in which every heterozygous individual has a fitness of $1 + s$ relative to every homozygous individual. This model is a special case of the one presented here, provided $r = 0$ and $k \rightarrow 0$. The estimates of Ns obtained by SLATKIN and MUIRHEAD (2000) range between 10^2 and 10^3 , which agrees with previous estimates from TAKAHATA *et al.* (1992) and SATTA *et al.* (1994) and coincides with the strength of selection we have assumed in this article ($N = 10^2$ – 10^3 and $\alpha = 1$, so $N\alpha = 10^2$ – 10^3). However, the kind of interactions involved in the *Mhc* is a far more complex issue.

In the coding regions of *Mhc* genes, there is an overall enhancement of variability, sometimes accompanied by linkage disequilibrium both within and among genes (HUGHES and YEAGER 1997, 1998; MEYER and THOMSON 2001). The predominance of positive epistasis cannot be deduced from high linkage disequilibrium levels because, first, linkage disequilibrium can have other causes besides the specific kind of positive epistasis we use in this study; second, linkage disequilibrium is usually not complete, several haplotypes being present at high frequencies in *Mhc* loci; and, third, the exact targets of selection are generally unknown, making it difficult to ascertain how many relevant backgrounds are actually involved. Some insight may be gained by considering the different expectations produced by different kinds

of epistasis, because patterns of variability within the *Mhc* suggest that different loci may be involved in different kinds of interaction. Effective population size in humans is estimated to be $\sim 10^4$ (PRZEWORSKI *et al.* 2000), so the recombination threshold beyond which multilocus balancing selection with negative epistasis will have no effect on neutral variability is $r = 10^{-4}$ ($Nr \sim 1$). It is clear that intragenic recombination is below the threshold, but the whole *Mhc* region spans ~ 4 Mb on chromosome 6 (BECK and TROWSDALE 2000), which, assuming a recombination rate of 1 cM/Mb (CULLEN *et al.* 1997) means $r = 0.04$ for the whole region. Class II loci HLA-DQA1 and HLA-DQB1 are separated by 20 kb, so the recombination rate in the intervening region is $\sim 10^{-4}$. Those are the two genes for which the best evidence of linkage disequilibrium has been found (GAUDIERI *et al.* 1999, 2000; SANCHEZ-MAZAS *et al.* 2000) and, if epistatic interactions favoring linkage disequilibrium were dominant in the system, variability and disequilibrium should be high and the frequency spectrum of neutral mutations should be even in the intervening region. Class II locus HLA-DPB1 is 400 kb away from HLA-DQA1 ($r \sim 4 \times 10^{-3}$, $Nr \sim 4$), so some of the variability within these two loci is beyond the threshold of action of multilocus selection if epistasis is such that it tends to make linkage disequilibrium low. This would seem to be the case because, although HLA-DPB1 has high variability, it is the locus showing less evidence of linkage disequilibrium with any other HLA loci. Interestingly, its frequency spectrum does not present significant deviations from neutrality in the available studies (SANCHEZ-MAZAS *et al.* 2000). Further suggestive evidence comes from the presence of recombination hotspots and coldspots in the *Mhc* (*cf.* BECK and TROWSDALE 2000). It is hardly surprising that the levels of linkage disequilibrium are lower in highly recombining regions. However, selection may also have some role in shaping the distribution of recombination in the *Mhc*, because different kinds of epistasis change the effective recombination rates by selectively getting rid of recombinants (positive epistasis or multiplicative fitness) or by favoring them (negative epistasis). For example, when one detects low recombination regions, one might actually be detecting regions under strong selection. If that were the case, differences in the frequency spectrum should be expected between high variability regions with different recombination rates.

All this evidence suggests that studying regions between *Mhc* loci, and not only the *Mhc* loci themselves, is a potentially fruitful strategy to ascertain the kind of selection involved. Unfortunately, only a few such studies are available (GUILLAUMEUX *et al.* 1998; HORTON *et al.* 1998; GAUDIERI *et al.* 1999, 2000; O'HUIGIN *et al.* 2000). These studies describe the existence within the *Mhc* of "polymorphic frozen blocks" (GAUDIERI *et al.* 1999, 2000), stretches of extremely high nucleotide variability interrupted by regions of low variability. This

feature raises hopes of distinguishing between possible kinds of interaction affecting different regions of the *Mhc*. However, much work still needs to be done. It is not yet established, for example, whether the high variability in noncoding regions can be attributed exclusively to hitchhiking with HLA loci, and selection may be directly acting on those regions (GAUDIERI *et al.* 1999). Clearly, more than mere accumulation of sequence data is needed to assess this sort of question. It is also necessary to analyze available data sets with a multilocus frame of mind, focusing on how the *Mhc*-wide patterns of linkage disequilibrium and frequency spectrum correlate with variability levels.

Several simplifying assumptions underlie our model. First, although the analytical approach used by BARTON and NAVARRO (2002) is completely general, our simulations focus on symmetric overdominance acting on several equally spaced, diallelic loci that contribute equally to fitness. The patterns of neutral variability and linkage disequilibrium expected in more complex multilocus systems are difficult to predict, but certainly worth studying. For example, the removal of symmetry will allow for more backgrounds to be present in the population even with positive epistasis (a scenario that, incidentally, will be closer to *Mhc* variability, which presents a large number of haplotypes with high linkage disequilibrium). Second, due to its complexity we did not consider the effect of gene conversion on the multilocus scenario. Gene conversion may be important in the positive epistasis scenario, because, in contrast with crossing over, it provides a mechanism of background homogenization that is not opposed by selection. Gene conversion events in segments between selected loci will allow neutral variants to segregate away without breaking down linkage disequilibrium. Also, the intragenic recombination generated by gene conversion may be at the origin of selectively relevant backgrounds (HOGSTRAND and BOHME 1999). Finally, we assume that the selection-drift equilibrium has been reached, but KELLY and WADE (2000) show that the patterns of variability expected during the approach to equilibrium in a two-locus model are quite different from the ones expected at equilibrium. We have shown that in a multilocus scenario equilibrium is even more difficult to achieve, so the differences are expected to be larger. This is relevant for the study of *Mhc* and other balancing selection systems. In short distances within the human *Mhc*, for example, recombination is of the order of $r \sim 10^{-5}$ /kb (TAKAHATA and SATTA 1998). In this case, in a system formed by 10 selected sites (say, 10 amino acids within the same exon) the average pairwise coalescence time is $\sim 10^3 N$ generations. Assuming a generation time of ~ 10 years this gives an average pairwise coalescence time of $\sim 10^8$ years. This shows that the populations cannot possibly be at equilibrium. Even if a species survived to such an amount of time, turnover of alleles by mutations at the selected sites would preclude equilibrium. Some of the

patterns we have detected, particularly the ones referring to frequency spectrum, are stronger in populations that have not reached equilibrium. Deeper insight is to be gained by the study of nonequilibrium multilocus systems. Such analysis is currently under way.

We thank P. Andolfatto, P. Awadalla, B. Charlesworth, D. Charlesworth, F. Depaulis, S. Otto, J. Rozas, and three anonymous reviewers for valuable discussion and criticism. A.N. is grateful to F. Depaulis, whose comments were particularly helpful (and extremely funny), and to D. Charlesworth, whose ideas made this work readable. This work was supported by Biotechnology and Biological Sciences Research Council/Engineering and Physical Sciences Research Council.

LITERATURE CITED

- AQUADRO, C. F., and D. J. BEGUN, 1993 Evidence for and implications of genetic hitch-hiking in the *Drosophila* genome, pp. 159–178 in *Mechanisms of Molecular Evolution*, edited by N. TAKAHATA and A. G. CLARK. Sinauer Press, Sunderland, MA.
- AQUADRO, C. F., D. J. BEGUN and E. C. KINDAHL, 1994 Selection, recombination and DNA polymorphism in *Drosophila*, pp. 46–56 in *Non-Neutral Evolution*, edited by B. GOLDING. Chapman & Hall, New York.
- BARTON, N., 2000 Genetic hitchhiking. *Philos. Trans. R. Soc. Lond. B* **355**: 1–10.
- BARTON, N. H., and A. NAVARRO, 2002 Extending the coalescent to multilocus systems: the case of balancing selection. *Genet. Res.* **79**: 129–139.
- BARTON, N. H., and M. SHPAK, 2000 The stability of symmetrical solutions to polygenic models. *Theor. Popul. Biol.* **57**: 249–263.
- BECK, S., and J. TROWSDALE, 2000 The human major histocompatibility complex: lessons from the DNA sequence. *Annu. Rev. Genomics Hum. Genet.* **1**: 117–138.
- BERGELSSON, J., M. KREITMAN, E. A. STAHL and D. TIAN, 2001 Evolutionary dynamics of plant R-genes. *Science* **292**: 2281–2285.
- BERNAL, A., U. EAR and N. KYRPIDES, 2001 Genomes OnLine Database (GOLD): a monitor of genome projects world-wide. *Nucleic Acids Res.* **29**: 126–127.
- CHARLESWORTH, B., 1994 The effect of background selection against deleterious mutations on weakly selected, linked variants. *Genet. Res.* **63**: 213–228.
- CHARLESWORTH, B., M. T. MORGAN and D. CHARLESWORTH, 1993 The effect of deleterious mutations on neutral molecular variation. *Genetics* **134**: 1289–1303.
- CHARLESWORTH, D., and P. AWADALLA, 1998 Flowering plant self-incompatibility: the molecular population genetics of Brassica S-loci. *Heredity* **81**: 1–9.
- CHRISTIANSEN, F. B., 1987 The deviation from linkage equilibrium with multiple loci varying in a stepping-stone cline. *J. Genet.* **66**: 45–67.
- CHRISTIANSEN, F. B., 1988 Epistasis in the multiple locus symmetric viability model. *J. Math. Biol.* **26**: 595–618.
- CHRISTIANSEN, F. B., 2000 *Population Genetics of Multiple Loci*. John Wiley & Sons, New York.
- CROW, J. F., and M. KIMURA, 1970 *An Introduction to Population Genetics Theory*. Harper & Row, New York.
- CULLEN, M., J. NOBLE, H. ERLICH, K. THORPE, S. BECK *et al.*, 1997 Characterization of recombination in the HLA class II region. *Am. J. Hum. Genet.* **60**: 397–407.
- DVORAK, J., M. C. LUO and Z. L. YANG, 1998 Restriction fragment length polymorphism and divergence in the genomic regions of high and low recombination in self-fertilizing and cross-fertilizing aegilops species. *Genetics* **148**: 423–434.
- EWENS, W. J., 1979 *Mathematical Population Genetics*. Springer Verlag, Berlin.
- FRANKLIN, I., and R. C. LEWONTIN, 1970 Is the gene the unit of selection? *Genetics* **65**: 701–734.
- FU, Y. X., 1997 Statistical tests of neutrality of mutations against population growth, hitchhiking and background selection. *Genetics* **147**: 915–925.
- GAUDIERI, S., J. K. KULSKI, R. L. DAWKINS and T. GOJOBORI, 1999 Extensive nucleotide variability within a 370 kb sequence from the central region of the major histocompatibility complex. *Gene* **238**: 157–161.
- GAUDIERI, S., R. L. DAWKINS, K. HABARA, J. K. KULSKI and T. GOJOBORI, 2000 SNP profile within the human Major Histoincompatibility Complex reveals an extreme and interrupted level of nucleotide diversity. *Genome Res.* **10**: 1579–1586.
- GUILLAUMEUX, T., M. JANER, G. K. S. WONG, T. SPIES and D. E. GERAGHTY, 1998 The complete genomic sequence of 424,015 bp at the centromeric end of the HLA class I region: gene content and polymorphism. *Proc. Natl. Acad. Sci. USA* **95**: 9494–9499.
- HEY, J., 1991 The structure of genealogies and the distribution of fixed differences between DNA sequence samples from natural populations. *Genetics* **128**: 831–840.
- HOGSTRAND, K., and J. BOHME, 1999 Gene conversion can create new MHC alleles. *Immunol. Rev.* **167**: 305–317.
- HORTON, R., D. NIBLETT, S. MILNE, S. PALMER, B. TUBBY *et al.*, 1998 Large-scale sequence comparisons reveal unusually high levels of variation in the HLA-DQB1 locus in the class II region of the human MHC. *J. Mol. Biol.* **282**: 71–87.
- HUDSON, R., 1990 Gene genealogies and the coalescent process. *Oxf. Surv. Evol. Biol.* **7**: 1–44.
- HUDSON, R. B., and N. L. KAPLAN, 1988 The coalescent process in models with selection and recombination. *Genetics* **120**: 831–840.
- HUDSON, R. R., and N. L. KAPLAN, 1994 Gene trees with background selection, pp. 140–153 in *Non-neutral Evolution: Theories and Molecular Data*, edited by G. B. GOLDING. Chapman & Hall, London/New York.
- HUDSON, R. R., and N. L. KAPLAN, 1995 The coalescent process with background selection. *Philos. Trans. R. Soc. B* **349**: 19–23.
- HUGHES, A. L., and M. YEAGER, 1997 Molecular evolution of the vertebrate immune system. *BioEssays* **19**: 777–786.
- HUGHES, A. L., and M. YEAGER, 1998 Natural selection and the evolutionary history of major histocompatibility complex loci. *Front. Biosci.* **3**: 509–516.
- KAPLAN, N. L., T. DARDEN and R. B. HUDSON, 1988 The coalescent process in models with selection. *Genetics* **120**: 819–829.
- KAPLAN, N. L., R. R. HUDSON and C. H. LANGLEY, 1989 The hitchhiking effect revisited. *Genetics* **123**: 887–899.
- KARLIN, S., and H. AVNI, 1981 Analysis of central equilibria in multilocus systems. *Theor. Popul. Biol.* **20**: 241–280.
- KARLIN, S., and U. LIBERMAN, 1979 Central equilibria in multilocus systems. II. Bisexual generalized nonepistatic selection schemes. *Genetics* **91**: 799–816.
- KELLY, J. K., and M. J. WADE, 2000 Molecular evolution near a two-locus balanced polymorphism. *J. Theor. Biol.* **204**: 83–101.
- LEWONTIN, R. C., and K. KOJIMA, 1960 The evolutionary dynamics of complex polymorphisms. *Evolution* **14**: 450–472.
- LYTTLE, T. W., 1991 Segregation distorters. *Annu. Rev. Genet.* **25**: 511–557.
- MARUYAMA, T., 1972 Rate of decrease of genetic variability in a two-dimensional continuous population of finite size. *Genetics* **70**: 639–651.
- MEYER, D., and G. THOMSON, 2001 How selection shapes variation on the human major histoincompatibility complex: a review. *Ann. Hum. Genet.* **65**: 1–26.
- NACHMAN, M. W., 1997 Patterns of DNA variability at X-linked loci in *Mus domesticus*. *Genetics* **147**: 1303–1316.
- NAGYLAKI, T., 1982 Geographical invariance in population genetics. *J. Theor. Biol.* **99**: 159–172.
- NORDBORG, M., 1997 Structured coalescent processes on different time scales. *Genetics* **146**: 1501–1514.
- NORDBORG, M., 2001 Coalescent theory, pp. 179–212 in *Handbook of Statistical Genetics*, edited by D. BALDING, M. BISHOP and C. CANNINGS. Wiley & Sons, New York.
- O'HUIGIN, C., Y. SATTI, A. HAUSMANN, R. L. DAWKINS and J. KLEIN, 2000 The implications of intergenic polymorphism for major histocompatibility complex evolution. *Genetics* **156**: 867–877.
- PALOPOLI, M. F., 2000 Genetic partners in crime: evolution of a coadapted gene complex that specializes in sperm sabotage, pp. 113–126 in *Epistasis and the Evolutionary Process*, edited by E. D. BRODIE, M. U. WADE and J. WOLF. Oxford University Press, London.
- PRZEWORSKI, M., R. R. HUDSON and A. DI RIENZO, 2000 Adjusting the focus on human variation. *Trends Genet.* **16**: 296–302.
- SANCHEZ-MAZAS, A., S. DJOULAH, M. BUSSON, I. LE MONNIER DE GOU-

- VILLE, J. C. POIRIER *et al.*, 2000 A linkage disequilibrium map of the Mhc region based on the analysis of 14 loci haplotypes in 50 French families. *Eur. J. Hum. Genet.* **8**: 33–41.
- SATTA, Y., 1997 Effect of intra-locus recombination on HLA polymorphism. *Hereditas* **127**: 105–112.
- SATTA, Y., G. O’HUIGIN, N. TAKAHATA and J. KLEIN, 1994 Intensity of natural selection at the major histocompatibility complex loci. *Proc. Natl. Acad. Sci. USA* **91**: 7184–7188.
- SATTA, Y., Y. J. LI and N. TAKAHATA, 1998 The neutral theory and natural selection in the HLA region. *Front. Biosci.* **27**: 459–467.
- SLATKIN, M., and C. A. MUIRHEAD, 2000 A method for estimating the intensity of overdominant selection from the distribution of allele frequencies. *Genetics* **156**: 2119–2126.
- STEPHAN, W., and C. H. LANGLEY, 1998 DNA polymorphism in *Lycopodium obscurum* and crossing-over per physical length. *Genetics* **150**: 1585–1593.
- STEPHAN, W., T. H. WIEHE and M. LENZ, 1992 The effect of strongly selected substitutions on neutral polymorphism: analytical results based on diffusion theory. *Theor. Popul. Biol.* **41**: 237–254.
- STROBECK, C., 1983 Expected linkage disequilibrium for a neutral locus linked to a chromosome rearrangement. *Genetics* **103**: 545–555.
- SVED, J. A., 1983 Does natural selection increase or decrease variability at linked loci? *Genetics* **105**: 239–340.
- TAJIMA, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**: 585–595.
- TAKAHATA, N., and M. NEI, 1990 Allelic genealogy under overdominant and frequency-dependent selection and polymorphism of MHC loci. *Genetics* **124**: 967–978.
- TAKAHATA, N., and Y. SATTA, 1998 Footprints of intragenic recombination at HLA loci. *Immunogenetics* **47**: 430–441.
- TAKAHATA, N., Y. SATTA and J. KLEIN, 1992 Polymorphism and balancing selection at major histocompatibility loci. *Genetics* **130**: 925–938.
- WHITLOCK, M. C., and N. H. BARTON, 1997 The effective size of a subdivided population. *Genetics* **146**: 427–441.

Communicating editor: N. TAKAHATA

