# Contrasting Patterns of Nonneutral Evolution in Proteins Encoded in Nuclear and Mitochondrial Genomes

## Daniel M. Weinreich and David M. Rand

*Department of Ecology and Evolutionary Biology, Brown University, Providence, Rhode Island 02912*

## ABSTRACT

We report that patterns of nonneutral DNA sequence evolution among published nuclear and mitochondrially encoded protein-coding loci differ significantly in animals. Whereas an apparent excess of amino acid polymorphism is seen in most (25/31) mitochondrial genes, this pattern is seen in fewer than half (15/36) of the nuclear data sets. This differentiation is even greater among data sets with significant departures from neutrality (14/15 *vs.* 1/6). Using forward simulations, we examined patterns of nonneutral evolution using parameters chosen to mimic the differences between mitochondrial and nuclear genetics (we varied recombination rate, population size, mutation rate, selective dominance, and intensity of germ line bottleneck). Patterns of evolution were correlated only with effective population size and strength of selection, and no single genetic factor explains the empirical contrast in patterns. We further report that in *Arabidopsis thaliana*, a highly self-fertilizing plant with effectively low recombination, five of six published nuclear data sets also exhibit an excess of amino acid polymorphism. We suggest that the contrast between nuclear and mitochondrial nonneutrality in animals stems from differences in rates of recombination in conjunction with a distribution of selective effects. If the majority of mutations segregating in populations are deleterious, high linkage may hinder the spread of the occasional beneficial mutation.

S INCE the introduction of DNA sequencing technology to population genetics (Kreitman 1983), many protein-coding loci have been examined in many species. A wide range of patterns of polymorphism and divergence, consistent with a variety of selective processes, have been described in nuclear genes (Brookfield and Sharp 1994; Kreitman and Akashi 1995). The elimination of strongly deleterious mutations by the action of purifying selection appears to be a ubiquitous selective force (Kimura 1983; Kreitman 1983), although examples of balancing selection (*e.g.*, Hughes and Nei 1988; Kreitman and Hudson 1991) and directional selection (*e.g.*, Long and Langley 1993; Messier and Stewart 1997) acting on amino acid replacement mutations have also been reported.

To date, protein-coding genes on mitochondrial DNA (mtDNA) in animals have not been found to exhibit the diversity of polymorphism and divergence patterns seen in nuclear genes. On the contrary, nearly every sequencing study testing the neutrality of animal protein-coding genes in mtDNA reveals the same pattern: an excess of amino acid replacement mutations segregating within species, relative to fixed amino acid replacement mutations (Ballard and Kreitman 1994; Nachman *et al.* 1994, 1996; Rand *et al.* 1994; Rand and Kann 1996; Wise *et al.* 1998). Moreover, surveys of previously published animal mtDNA sequences have extended these observations (Hasegawa *et al.* 1998; Nachman 1998; Rand and Kann 1998).

Ohta and Kimura (1971) pointed out that slightly deleterious mutations may reach fixation in populations of finite size in spite of the pressure from purifying selection, as a consequence of genetic drift. Because the probability of fixation for deleterious mutations is an inverse function of the product of $N$ (effective population size) and $s$ (selective coefficient; Ohta 1972), a deleterious mutation with a given $s$ ($s < 0$) will be more likely to reach fixation in a small population than in a large one. Kimura (1983, Figure 3.7) made a second prediction about slightly deleterious mutations: for any negative value of $Ns$, the reduction in fixation probability relative to the strictly neutral expectation will be greater than the reduction in heterozygosity. This is a consequence of the fact that even those slightly deleterious mutations destined for loss may nevertheless persist in the population for a time due to drift and will therefore contribute to heterozygosity. For example, Akashi (1995) has argued that a large effective population size in *Drosophila simulans* is responsible for the observation of significant excess of selectively "unpreferred" codons segregating in that species, relative to fixed "unpreferred" codons. In contrast, *D. melanogaster*, which is thought to have a smaller effective population size, shows no such excess segregation of putatively mildly deleterious synonymous mutations.

Thus the observation of a relative excess number of

*Corresponding author:* Daniel M. Weinreich, Department of Biology, Muir Bldg., University of California, 9500 Gilman Dr., San Diego, CA 92093. E-mail: dmw@ucsd.edu

segregating amino acid replacement mutations in mtDNA-encoded loci is consistent with the assumption that many segregating amino acid replacement mutations are slightly deleterious. RAND and KANN (1996) partitioned segregating mutations at the (mitochondrial) ND5 locus in *D. melanogaster* into synonymous and amino acid replacement sites and applied Tajima's *D* statistic (TAJIMA 1989) to these two classes of sites independently. They found that the hypothesis of neutral evolution could be rejected only for amino acid replacement mutations, which showed a deviation consistent with weak purifying selection acting on these sites. Recently, NIELSEN and WEINREICH (1999) have shown that in models of recurrent mutation and genetic drift, mildly deleterious mutations will on average be younger than neutral mutations, even in the absence of recombination. They further showed that the mean age of segregating amino acid replacement mutations tends to be less than the mean age of segregating synonymous mutations in animal mtDNA, consistent with the view that such mutations are being weakly selected against.

Here we explore two questions suggested by these observations. First, is there significantly more diversity in the patterns of polymorphism and divergence of nuclear-encoded genes than of mitochondrially encoded genes? To assess this question, we have performed a careful survey of the literature for data sets of nuclear and mitochondrial polymorphism and divergence. Second, animal mitochondrial and nuclear DNA exhibit five gross genetic differences: mtDNA apparently lacks recombination (MORITZ *et al.* 1987; but see LUNT and HYMAN 1997; AWADALLA *et al.* 1999; EYRE-WALKER *et al.* 1999), has a smaller effective population size (BIRKY *et al.* 1983) as a consequence of maternal inheritance and an extreme population bottleneck during the course of oogenesis (BENDALL *et al.* 1996; PARSONS *et al.* 1997), and a higher mutation rate in at least some lineages (AVISE 1991), and is haploid (HAUSWIRTH and LAIPIS 1982; JENUTH *et al.* 1996). If we assume that the molecular evolution of nuclear- and mitochondrially encoded loci is driven by common selective forces, then can these genetic differences account for observed differences in the patterns of evolution? Classical diffusion-derived expressions for polymorphism (KIMURA 1969) and divergence (KIMURA 1957) are known, assuming genetic and selective independence among sites, but could not easily be extended to the present case. We therefore performed a series of computer simulations in which each of these genetic factors was independently varied under models of positive, negative, and no selection. These simulations allowed us to examine the evolutionary and sampling behavior of genes under "nuclear" and "mitochondrial" conditions. While the manifestations of selection in simulated nuclear and mitochondrial genes differ dramatically, no single genetic factor is sufficient to explain the empirical differences seen.

## MATERIALS AND METHODS

**Published DNA sequences:** Data sets consisting of DNA sequence polymorphism and divergence for 39 nuclear loci from *Drosophila* spp. and 31 data sets for 7 mtDNA-encoded loci from diverse animal species were compiled from the literature. Many of the nuclear data sets are those used in MORIYAMA and POWELL (1996); most of the mtDNA data sets are pooled from NACHMAN (1998), RAND and KANN (1998), and NIELSEN and WEINREICH (1999). Two studies of each of three genes in *D. melanogaster* are included in Table 1 (*Acp26Aa, Acp26Ab,* and *Est-6*). For the purposes of this study, a random data set for each of these genes was discarded, leaving 36 nuclear data sets. Gene name, protein length, species, sample size, fixed and polymorphic synonymous and amino acid replacement site counts, neutrality index (N.I., defined below), *P* value of the test statistic from the associated McDONALD and KREITMAN (1991) test, and citations appear in Tables 1 (nuclear encoded) and 2 (mtDNA encoded).

A data set consisting of DNA sequence polymorphism and divergence for six nuclear loci from the plant *Arabidopsis thaliana* was similarly compiled from the literature. Gene name, protein length, sample size, fixed and polymorphic synonymous and amino acid replacement site counts, N.I., *P* value of the test statistic from the associated McDONALD and KREITMAN (1991) test, and citations appear in Table 3.

**Computer simulations:** Computer simulations were written in "C" and compiled to run under UNIX. Simulations were parameterized in eight dimensions, shown in Table 4. Simulations follow *N* chromosomes, each represented by the interval (0, 1), which undergo repeated cycles (generations) of mutation, recombination, random mating and selection, and sampling. All statistics are calculated after recombination but before random mating and selection (WATTERSON 1975; R. N. NIELSEN, personal communication).

Mutations are of two sorts, selected and neutral, and in each generation the number of each sort in the population is determined by an independent Poisson-distributed deviate with mean $N\mu/2$. Chromosomes to be mutated are chosen at random and mutated "sites" are located as uniformly distributed real numbers on the interval (0, 1). Thus our simulations adhere to the infinite sites model (KIMURA 1969).

In any given generation, the number of recombination events is Poisson distributed with mean *Nc*. Pairs of "parental" chromosomes to be recombined are chosen randomly and the location of the crossover site is chosen as a uniformly distributed real number on the interval (0, 1). Each recombination event generates two novel chromosomes consisting, respectively, of all sites present on the first parental whose locations are numerically less than the crossover site together with all sites present on the second parental whose locations are numerically greater than the crossover site, and all sites present on the second parental chromosome whose locations are numerically less than the crossover site together with all sites present on the first parental whose locations are numerically greater than the crossover site.

Relative fitness is assessed for diploid genotypes. Genotype frequencies are calculated by the Hardy-Weinberg equation using allele frequencies before selection, which is equivalent to assuming random mating. Thus in these simulations *N* is both the census and effective population size. Under a multiplicative fitness model, the fitness of the *i-j*th genotype is given by

$$w_{i,j} = \frac{(1 + 2s)^{m_{i,j}}(1 + hs)^{n_{i,j}}}{\overline{w}}, \qquad (1)$$

where *s* is the selection coefficient acting on selected sites, *h* is the degree of dominance, $m_{i,j}$ is the number of selected sites

## TABLE 1

**Locus, number of codons sequenced, species, sample size, mutation class counts, N.I., statistical significance of the McDonald and Kreitman (1991) test, and citation for nuclear-encoded DNA sequence surveys used in this study**

| Locus | Codons | Species[a] | $n$ | FR[b] | FS[b] | PR[b] | PS[b] | N.I. | $P$ value[c] | Citation |
|---|---|---|---|---|---|---|---|---|---|---|
| *6-pgd* | 481 | *mel* | 13 | 7 | 45 | 1 | 3 | 2.14 | 0.60 | Begun and Aquadro (1994) |
| *Acp26Aa* | 258 | *mel* | 49 | 75 | 21 | 22 | 16 | 0.39 | 0.02 | Tsaur *et al.* (1998) |
| *Acp26Aa* | 236 | *mel* | 62 | 75 | 20 | 23 | 25 | 0.25 | 0.0002 | Aguadé (1998) |
| *Acp26Ab* | 90 | *mel* | 49 | 1 | 2 | 4 | 6 | 1.33 | 0.83 | Tsaur *et al.* (1998) |
| *Acp26Ab* | 90 | *mel* | 62 | 3 | 2 | 5 | 12 | 0.28 | 0.22 | Aguadé (1998) |
| *Acp29AB* | 234 | *mel* | 39 | 36 | 33 | 6 | 15 | 0.37 | 0.06 | Aguadé (1999) |
| *Acph-1*[d] | 447 | *sub* | 41 | 0[e] | 3 | 27 | 78 | 1.04 | 0.18 | Navarro-Sabaté *et al.* (1999) |
| *Adh* | 256 | *mel + sim* | 15 | 2 | 2 | 3 | 27 | 0.11 | 0.16 | Moriyama and Powell (1996) |
| *Adh* | 290 | *pse* | 10 | 24 | 74 | 1 | 19 | 0.16 | 0.28 | Schaeffer and Miller (1992) |
| *Adh-dup* | 208 | *pse* | 10 | 21 | 109 | 7 | 25 | 1.45 | 0.45 | Schaeffer and Miller (1992) |
| *Anon1A3* | 310 | *mel + sim* | 38 | 26 | 19 | 22 | 10 | 1.61 | 0.33 | Schmid *et al.* (1999) |
| *Anon1E9* | 577 | *mel + sim* | 23 | 44 | 22 | 37 | 34 | 0.54 | 0.083 | Schmid *et al.* (1999) |
| *Anon1G5* | 253 | *mel + sim* | 30 | 22 | 10 | 24 | 22 | 0.50 | 0.14 | Schmid *et al.* (1999) |
| *ase* | 356 | *mel + sim* | 6 | 11 | 13 | 3 | 3 | 1.18 | 0.86 | Moriyama and Powell (1996) |
| *boss* | 522 | *mel + sim* | 5 | 3 | 16 | 4 | 52 | 0.41 | 0.78 | Moriyama and Powell (1996) |
| *CecA1* | 93 | *mel* | 9 | 2 | 9 | 1 | 1 | 4.50 | 0.36 | Ramos-Onsins and Aguadé (1998) |
| *CecA2* | 91 | *mel* | 9 | 2 | 18 | 2 | 6 | 3.00 | 0.33 | Ramos-Onsins and Aguadé (1998) |
| *ci* | 319 | *mel + sim* | 10 | 21 | 23 | 1 | 0[e] | 1.10 | 0.95 | Moriyama and Powell (1996) |
| *cta* | 322 | *mel + sim* | 18 | 8 | 13 | 0 | 1 | 0.0 | 0.33 | Wayne and Kreitman (1996) |
| *Dhc-Yh3* | 573 | *sim* | 10 | 1 | 48 | 0 | 1 | 0.00 | 0.84 | Zuroycova and Eanes (1999) |
| *Est-5A* | 549 | *pse* | 8 | 9 | 13 | 17 | 27 | 0.91 | 0.86 | King (1998) |
| *Est-5B* | 546 | *pse* | 8 | 3 | 15 | 25 | 35 | 3.57 | 0.053 | King (1998) |
| *Est-5C* | 546 | *pse* | 8 | 4 | 15 | 12 | 36 | 1.25 | 0.73 | King (1998) |
| *Est-6* | 544 | *mel + sim* | 13 | 16 | 19 | 27 | 78 | 0.411 | 0.26 | Moriyama and Powell (1996) |
| *Est-6* | 545 | *mel + sim* | 15 | 19 | 25 | 18 | 69 | 0.34 | 0.008 | Balakirev *et al.* (1999) |
| *G6pd* | 516 | *mel* | 32 | 21 | 26 | 2 | 36 | 0.07 | 0.0009 | Eanes *et al.* (1993) |
| *gld* | 323 | *sim* | 11 | 5 | 49 | 1 | 26 | 0.38 | 0.34 | Hamblin and Aquadro (1996) |
| *jgw* | 255 | *tei + yak* | 30 | 21 | 16 | 4 | 27 | 0.11 | 0.0001 | Long and Langley (1993) |
| *per* | 561 | *mel + sim* | 6 | 1 | 25 | 7 | 58 | 3.02 | 0.84 | Moriyama and Powell (1996) |
| *per* | 401 | *wil* | 18 | 16 | 16 | 11 | 44 | 0.25 | 0.004 | Gleason and Powell (1997) |
| *Pgi* | 558 | *mel + sim* | 11 | 1 | 23 | 4 | 18 | 5.11 | 0.03 | Moriyama and Powell (1996) |
| *pn* | 391 | *mel + sim* | 8 | 5 | 22 | 7 | 8 | 3.85 | 0.07 | Moriyama and Powell (1996) |
| *ref(2)p* | 599 | *mel* | 10 | 34 | 28 | 13 | 4 | 2.68 | 0.60 | Wayne *et al.* (1996) |
| *Rh3* | 383 | *mel + sim* | 5 | 3 | 14 | 0 | 32 | 0.00 | 0.39 | Moriyama and Powell (1996) |
| *runt* | 511 | *mel* | 11 | 5 | 16 | 3 | 20 | 0.48 | 0.35 | Labate *et al.* (1999) |
| *runt* | 511 | *sim* | 11 | 5 | 16 | 1 | 8 | 0.40 | 0.41 | Labate *et al.* (1999) |
| *Tpi* | 250 | *mel* | 25 | 1 | 35 | 1 | 37 | 0.95 | 0.97 | Hasson *et al.* (1999) |
| *Yp2* | 349 | *mel + sim* | 6 | 5 | 11 | 1 | 10 | 0.22 | 0.24 | Moriyama and Powell (1996) |
| *z* | 268 | *mel + sim* | 6 | 2 | 15 | 0 | 16 | 0.00 | 0.030 | Moriyama and Powell (1996) |

[a] All genes from *Drosophila* spp. *mel*, *D. melanogaster*; *pse*, *D. pseudoobscura*; *sim*, *D. simulans*; *sub*, *D. subobscura*; *mel + sim*, polymorphism in *D. melanogaster* and *D. simulans* pooled; *tei + yak*, polymorphism in *D. teissieri* and *D. yakuba* pooled; *wil*, *D. willistoni*.

[b] Mutation classes: FR, fixed amino acid replacement; FS, fixed synonymous; PR, polymorphic amino acid replacement; PS, polymorphic synonymous.

[c] McDonald and Kreitman (1991) test statistic significance.

[d] $O_{3+4}$ and $O_{st}$ polymorphism pooled.

[e] Zero replaced with 1 for the purposes of calculating N.I. only. See materials and methods.

chromosomes $i$ and $j$ have in common, and $n_{i,j}$ is the number of selected sites appearing on exactly one of chromosomes $i$ and $j$. $\overline{w}$ is the population mean fitness and is given by

$$\overline{w} = \sum_i \sum_j p_i p_j w_{i,j}.$$

Under an additive fitness model, the fitness of the $i$-$j$th genotype is given by

$$w_{i,j} = \frac{1 + 2sm_{i,j} + hsn_{i,j}}{\overline{w}}, \tag{2}$$

although $w_{i,j}$ is set to 0 if $s < -1/(2m_{i,j} + hn_{i,j})$. $\overline{w}$, $m_{i,j}$, and $n_{i,j}$ are as above.

Finally, Wright-Fisher sampling is performed according to Gillespie (1993) and all chromosomes in the population are compared to identify sites newly fixed in the population. Whenever such a site is found, the corresponding fixation-

**TABLE 2**

**Locus, number of codons sequenced, species, sample size, mutation class counts, N.I., statistical significance of the McDonald and Kreitman (1991) test, and citation for mtDNA-encoded sequence surveys used in this study**

| Locus | Codons | Species[a] | $n$ | FR[b] | FS[b] | PR[b] | PS[b] | N.I. | P value[c] | Reference or citation in which N.I. is first reported[d] |
|---|---|---|---|---|---|---|---|---|---|---|
| ATPase 6 | 225 | D. mel. | 4 | 4 | 18 | 4 | 1 | 18.0 | 0.017 | Nachman (1998) |
| COI + COI | 740 | Heloconius erato | 52 | 2 | 24 | 16 | 70 | 2.74 | 0.235 | Nachman (1998) |
| COI | 188 | Alpheus spp. | 28 | 0[e] | 27 | 2 | 131 | 0.41 | 0.55 | Rand and Kann (1998) |
| COII | 228 | G. gorilla, H. sapiens and P. troglodytes | 21 | 8 | 113 | 14 | 42 | 4.70 | 0.001 | Templeton (1996) |
| cyt b | 80 | Ambystoma spp. | 16 | 0[e] | 22 | 3 | 4 | 16.5 | 0.001 | Rand and Kann (1998) |
| cyt b | 374 | Brachyramphus spp. | 19 | 0[e] | 50 | 3 | 24 | 6.25 | 0.02 | Rand and Kann (1998) |
| cyt b | 97 | Dendrobates pumilio | 12 | 4 | 6 | 0 | 6 | 0.0 | 0.234 | Nachman (1998) |
| cyt b | 291 | D. pse. | 13 | 5 | 32 | 1 | 8 | 0.80 | 0.86 | Rand and Kann (1998) |
| cyt b | 379 | Drosophila spp. | 34 | 10 | 97 | 6 | 12 | 4.85 | 0.01 | Ballard and Kreitman (1994) |
| cyt b | 260 | Emoia impar | 8 | 12 | 81 | 3 | 4 | 5.06 | 0.067 | Nachman (1998) |
| cyt b | 228 | Ensatina eschscholtzii | 24 | 21 | 27 | 38 | 171 | 0.29 | 0.3 | Nachman (1998) |
| cyt b | 100 | Gadus morhua | 41 | 0[e] | 10 | 3 | 22 | 1.36 | 0.25 | Rand and Kann (1998) |
| cyt b | 380 | Grus spp. | 13 | 2 | 49 | 10 | 25 | 9.80 | 0.002 | Rand and Kann (1998) |
| cyt b | 266 | Isothrix bistriata | 10 | 4 | 33 | 15 | 103 | 1.20 | 0.76 | Rand and Kann (1998) |
| cyt b | 144 | Melospiza melodia | 11 | 10 | 26 | 5 | 2 | 6.50 | 0.039 | Nachman (1998) |
| cyt b | 266 | Mesomys spp. | 31 | 0[e] | 14 | 30 | 126 | 3.33 | 0.075 | Rand and Kann (1998) |
| cyt b | 381 | Microtus spp. | 24 | 6 | 49 | 18 | 29 | 5.07 | 0.002 | Rand and Kann (1998) |
| cyt b | 144 | Passerella iliaca | 19 | 10 | 26 | 5 | 10 | 1.30 | 0.74 | Nachman (1998) |
| cyt b | 97 | Phyllobates lugubris | 8 | 0[e] | 19 | 11 | 59 | 3.54 | 0.06 | Nachman (1998) |
| cyt b | 94 | Pomatostomus temporalis | 35 | 8 | 18 | 0 | 17 | 0.0 | 0.014 | Nachman (1998) |
| cyt b | 380 | Scurius aberti | 20 | 18 | 146 | 12 | 38 | 2.56 | 0.025 | Rand and Kann (1998) |
| cyt b | 379 | Ursus arctos | 28 | 15 | 81 | 11 | 44 | 1.35 | 0.50 | Rand and Kann (1998) |
| NADH 2 | 342 | H. sapiens | 20 | 10 | 82 | 10 | 11 | 7.45 | 0.0005 | Wise et al. (1998) |
| NADH 2 | 342 | P. troglodytes | 20 | 10 | 82 | 7 | 32 | 1.79 | 0.40 | Wise et al. (1998) |
| NADH 3 | 178 | D. mel. + D. sim. | 37 | 2 | 13 | 1 | 5 | 1.30 | 0.86 | Rand and Kann (1996) |
| NADH 3 | 115 | Mus domesticus | 56 | 2 | 23 | 11 | 13 | 9.73 | 0.05 | Nachman et al. (1994) |
| NADH 3 | 115 | H. sapiens | 61 | 4 | 31 | 4 | 7 | 4.43 | 0.10 | Nachman et al. (1996) |
| NADH 3 | 115 | P. troglodytes | 5 | 4 | 31 | 4 | 3 | 10.3 | 0.025 | Nachman et al. (1996) |
| NADH 5 | 574 | D. mel. | 59 | 15 | 52 | 11 | 17 | 2.24 | 0.05 | Rand and Kann (1996, 1998) |
| NADH 5 | 466 | D. pse. | 22 | 9 | 25 | 2 | 19 | 0.29 | 0.59 | Rand and Kann (1998) |
| MtDNA[f] | 112[g] | H. sapiens | 147 | 179 | 915 | 23 | 33 | 3.56 | 0.00005 | Nachman et al. (1996) |

[a] Drosophila species names as in Table 1.

[b] Mutation classes: FR, fixed amino acid replacement; FS, fixed synonymous; PR, polymorphic replacement; PS, polymorphic synonymous.

[c] McDonald and Kreitman (1991) test statistic significance.

[d] Primary reference given therein.

[e] Zero replaced with 1 for the purposes of calculating N.I. only. See materials and methods.

[f] RFLP survey of entire mtDNA.

[g] Nucleotide substitutions inferred from 56 six-cutter restriction enzyme site gains, which is equivalent to sequencing 112 codons.

**TABLE 3**

**Locus, number of codons sequenced, sample size, mutation class counts, N.I., statistical significance of McDonald and Kreitman (1991) test, and citation for nuclear-encoded sequence surveys of *Arabidopsis thaliana* used in this study**

| Locus | Codons | $n$ | FR[a] | FS[a] | PR[a] | PS[a] | N.I. | $P$ value[b] | Citation |
|-------|--------|-----|-----|-----|-----|-----|------|---------|----------|
| *Adh*  | 379 | 17 | 14 | 59 | 7  | 13 | 2.02 | 0.15  | Miyashita *et al.* (1998) |
| *AP3*  | 232 | 18 | 5  | 14 | 20 | 8  | 5.83 | 0.002 | Purugganan and Suddith (1999) |
| *CAL*  | 191 | 16 | 13 | 15 | 16 | 5  | 3.69 | 0.03  | Purugganan and Suddith (1998) |
| *ChiA* | 600 | 17 | 15 | 31 | 19 | 25 | 1.57 | 0.30  | Kamabe *et al.* (1997) |
| *ChiB* | 336 | 17 | 18 | 27 | 4  | 17 | 0.35 | 0.08  | Kamabe and Miyashita (1999) |
| *PI*   | 208 | 16 | 11 | 16 | 12 | 4  | 4.36 | 0.03  | Purugganan and Suddith (1999) |

[a] Mutation classes: FR, fixed amino acid replacement; FS, fixed synonymous; PR, polymorphic amino acid replacement; PS, fixed synonymous.

[b] McDonald and Kreitman (1991) test statistic significance.

event counter is incremented [neutral ($c_{neut}$) or selected ($c_{sel}$)], and the site is removed from all chromosomes. Since sites reach fixation only in the sampling phase of the simulation, $m_{i,j}$ and $n_{i,j}$ in Equations 1 and 2 include only segregating selected sites, and the fitness of a chromosome is independent of the number of selected site fixations that have previously occurred in the simulation.

Uniform deviates on (0, 1) were generated with the UNIX library random number function (drand48()), seeded with the program's unique process identifier (getpid()). Poisson and binomial deviates were generated as described in Press *et al.* (1992).

Intralineal population bottlenecks were implemented as described in Bergstrom and Pritchard (1998) and occur after mutation but before selection. The parameters $N$ and $s$ were varied in this group of simulations, but $Nc$ was set to 0.0, $h$ was set to 1.0, and $\mu$ was set to $1/2N$. The additional parameters $M$ (the number of intralineal chromosomes before bottleneck, $M \leq N$) and $B$ (the size of the intralineal bottleneck, $B \leq M$) were employed as follows. In all cases, the intralineal bottleneck size ($B$) was set to 1 and the number of lineages was held at 1000, so that $N = 1000 \cdot M$. Thus, in these simulations, $N$ is not necessarily equal to the effective population size. The intensity of the bottleneck was parameterized by $M$, which assumed values of 10 (moderate intensity) and 100 (high intensity). When $M$ is set to 1, the Bergstrom and Pritch-

**TABLE 4**

**Definitions of parameters used in computer simulations**

| Parameter | Definition |
|-----------|------------|
| $N$ | Census number of chromosomes in the population (fixed at $1000 \cdot M$ for this study) |
| $\mu$ | Per generation per chromosome mutation rate |
| $s$ | Selection coefficient acting on selected mutations |
| $h$ | Dominance factor |
| $Nc$ | Per generation population recombination rate |
| $M$ | Intralineal chromosomal census size |
| $B$ | Intralineal bottleneck size (fixed at 1 for this study) |
| $T_{div}$ | Number of generations since species divergence (Fixed at $30 \cdot N$ for this study) |

ard model degenerates to the no-bottleneck model described above.

Simulations were performed at steady-state as previously described (Nielsen and Weinreich 1999). Briefly, a population fixed for a chromosome carrying no mutations is initiated for some point in parameter space and run for $100 \cdot N$ generations to reach quasi-equilibrium. At $2 \cdot T_{div}$ generation intervals thereafter, all the chromosomes in the population and the neutral and selected site fixation counters ($c_{neut}$ and $c_{sel}$) are recorded in a unique computer data file created for that point in parameter space. $2 \cdot T_{div}$ generations of simulation correspond to $T_{div}$ generations of divergence occurring simultaneously in two species. $c_{neut}$ and $c_{sel}$ were set to zero after the $100 \cdot N$ generation initialization and after each $2 \cdot T_{div}$ generation interval.

Since forward simulations are time intensive, these data files represented archived results, which could be reanalyzed as needed. Additionally, random chromosome samples of size $n < N$ were drawn from archived population replicates to examine the consequence of sampling on statistics of interest. Finally, recording replicate results into data files allowed us to make our simulation reentrant, thereby permitting us to utilize QUAHOG (http://www.cs.brown.edu/software/quahog/), a UNIX-based job management facility with access to >100 ULTRASparc1 workstations within the Brown University Computer Science Department. Simulations for each point in parameter space were run until 1000 replicates had accumulated in the data file for that point, unless otherwise noted.

The correctness of the simulations was verified by comparison with expectations from analytic results (Kimura 1957, 1969; Charlesworth *et al.* 1993) where possible.

**Statistics:** Published DNA sequence data sets were tested for deviation from neutral expectation with the McDonald and Kreitman (1991) test. In this test, all polymorphic sites are classified either as synonymous or as causing an amino acid replacement, and all fixed differences are similarly classified. No attempt was made to correct for multiple mutations at a nucleotide in counting fixed differences. Additionally, the neutrality index (N.I.; Rand and Kann 1996) was calculated for each McDonald/Kreitman table as

$$\text{N.I.} = \frac{(\text{no. polymorphic replacement sites/no. fixed replacement sites*})}{(\text{no. polymorphic synonymous sites*/no. fixed synonymous sites})}. \quad (3)$$

As defined by Rand and Kann (1996), N.I. values range from 0 to ∞, and under strict neutrality the ratios in the numerator and denominator are expected to be equal (McDonald and

Kreitman 1991; but see Maynard Smith 1994), giving an N.I. of 1.0. However, in those data sets in which the number of polymorphic synonymous or fixed replacement sites equals zero, we substituted 1 for the purposes of calculating N.I. (Rand and Kann 1998) to avoid division by zero. We denote this "no-division-by-zero" protocol with asterisks.

The following statistics were tabulated from the computer simulations: the number of neutral and segregating sites in the entire population in the $i$th replicate ($S^i_{N,\text{neut}}$ and $S^i_{N,\text{sel}}$, respectively) and the number of neutral and selected site fixation events in the $i$th replicate ($c^i_{\text{neut}}$ and $c^i_{\text{sel}}$, respectively). To explore the consequences of sampling from whole populations, 10 independent random samples of $x$ chromosomes each were drawn from each evolutionary replicate. We denote the number of neutral and selected sites segregating in the $j$th such sample drawn from the $i$th replicate as $S^{i,j}_{n=x,\text{neut}}$ and $S^{i,j}_{n=x,\text{sel}}$, respectively.

N.I.$_N$, the mean neutrality index for the entire population, was calculated as

$$\text{N.I.}_N = \frac{1}{r}\sum_{i=1}^{r}\text{N.I.}^i_N, \tag{4a}$$

where $r$ is the number of evolutionary replicates performed and N.I.$^i_N$, given by

$$\text{N.I.}^i_N = \frac{S^i_{n,\text{sel}}/c^i_{\text{sel}}}{S^i_{N,\text{neut}}/c^i_{\text{neut}}}, \tag{4b}$$

represents the neutrality index in the entire population in the $i$th simulated replicate. N.I.$_{n=x}$, the mean neutrality index for a sample of $x$ chromosomes drawn from the population, was calculated as

$$\text{N.I.}_{n=x} = \frac{1}{r}\sum_{i=1}^{r}\left[\frac{1}{10}\sum_{j=1}^{10}\text{N.I.}^{i,j}_{n=x}\right], \tag{5a}$$

where N.I.$^{i,j}_{n=x}$, given by

$$\text{N.I.}^{i,j}_{n=x} = \frac{S^{i,j}_{n=x,\text{sel}}/c^i_{\text{sel}}}{S^{i,j}_{n=x,\text{neut}}/c^i_{\text{neut}}}, \tag{5b}$$

represents the neutrality index in the $j$th subsample of size $x$ drawn from the $i$th simulated replicate. Values of N.I.$_{n=10}$ and N.I.$_{n=30}$ were calculated. We extended our no-division-by-zero protocol to these simulated data, substituting a 1 for $S^i_{N,\text{neut}}$, $S^{i,j}_{n=x,\text{neut}}$, or $c^i_{\text{sel}}$, in any case in which a zero was observed.

If one assumes that amino acid replacement mutations are selected and that synonymous mutations are neutral, then Equations 4 and 5 are seen to be equal to Equation 3. Though selection is known to act on some synonymous mutations in both genomes (Ballard and Kreitman 1994; Akashi 1995; Rand and Kann 1998), few would dispute that on average, selection is stronger on segregating amino acid replacement mutations, and so we do not feel that this assumption undermines our approach (see discussion).

**Power analysis:** The statistical power of the McDonald and Kreitman (1991) test to detect selection under the present model was measured as previously described (Nielsen and Weinreich 1999). For each set of parameter values simulated, power was estimated for three cases: using all polymorphisms segregating in the population, and using only that polymorphism segregating in random samples of size $n = 10$ and 30 chromosomes drawn from the population. In the latter cases, McDonald/Kreitman tests were performed on 10 replicate samples drawn from each replicate population. In all cases, the proportion of replicates that gave a test statistic significant at the 5% level was tabulated.
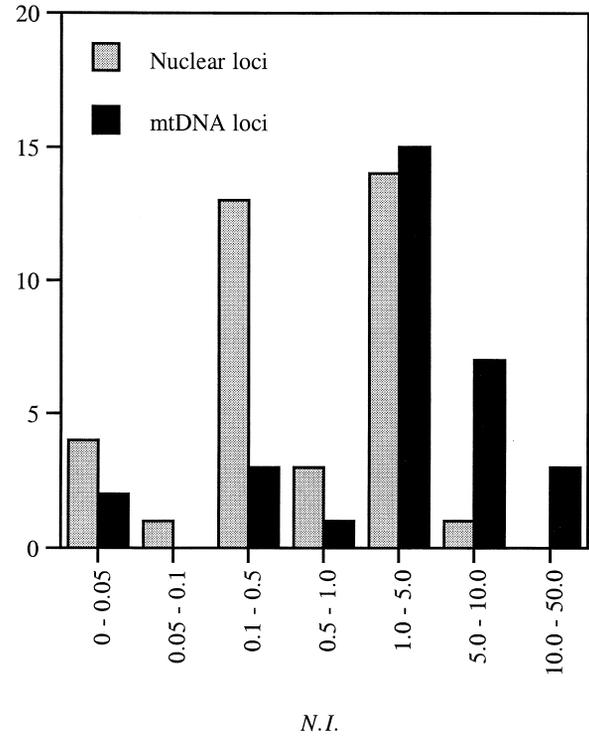


Figure 1.—Frequency histogram of observed values of N.I. (Equation 3) for 36 nuclear-encoded (shaded bars) and 31 mtDNA-encoded (solid bars) loci shown in Tables 1 and 2, respectively.

## RESULTS

**Published nuclear- and mtDNA-encoded DNA sequence analysis:** N.I. values for 39 nuclear- and 31 mtDNA-encoded loci are shown in Tables 1 and 2. After randomly removing one of each of the three duplicate nuclear data sets (see materials and methods), the mean nuclear-encoded N.I. value ($\pm$SD) is $1.21 \pm 1.38$ while the mtDNA-encoded mean N.I. is $4.41 \pm 4.52$, which differ significantly ($t = 3.82$, d.f. $= 66$, $P = 0.0002$). Figure 1 is a frequency histogram of N.I. values for nuclear- and mtDNA-encoded genes, and partitioning N.I. values into the classes shown in Figure 1 reveals a highly significant association with genome ($G = 18.56$, d.f. $= 6$, $P = 0.005$).

Moreover, we may restrict ourselves to those data sets with significant (defined as $P < 0.05$) McDonald/Kreitman test results: there were 6 such nuclear-encoded loci (*Acp26Aa*, *G6pd*, *jgw*, *per*, *Pgi*, and *z*), of which 5 have N.I. values <1.0. In contrast, 15 mtDNA-encoded loci have significant test results (*ATPase* 6 in *D. melanogaster*; *CO* II in hominoids; *Cyt b* in *Ambystoma* spp., *Brachyramphus* spp., *Drosophila* spp., *Grus* spp., *Melospiza melodia*, *Microtus* spp., *Pomatostomus temporalis*, and *Sciurus aberti*; *NADH* 2 in *Homo sapiens*; *NADH* 3 in *Mus domesticus* and *Pan troglodyte*; *NADH* 5 in *D. melanogaster*; and restriction fragment length polymorphism (RFLP) survey in *H. sapiens*), only one of which has N.I. values <1.0. These two observations jointly have a $P$ value of 0.0004 against

a null hypothesis of no difference in genome-specific bias in direction of significant deviation ($G = 12.37$, 1 d.f.).

**Published sequence analysis for nuclear DNA of *A. thaliana*:** N.I. values for six nuclear genes from *A. thaliana* are shown in Table 3. Five of the six genes have N.I. values >1.0, and the mean ($\pm$SD) N.I. value for these genes is $2.97 \pm 2.02$. Three of the genes exhibit significant McDonald and Kreitman (1991) test statistics; all of these have N.I. values >1.0.

**Diffusion approximation provides a lower bound for N.I. as a function of *Ns*:** By assuming selective and genetic site independence, Kimura (1957) developed analytic expectations for the probability of selected and neutral site fixation ($u_{sel}$ and $u_{neut}$, respectively), as well as for the number of segregating selected and neutral sites ($S_{sel}$ and $S_{neut}$, respectively; Kimura 1969). Under these assumptions, a lower bound for N.I. is given by

$$E(\text{N.I.}) = E\left(\frac{S_{sel}/c_{sel}^*}{S_{neut}^*/c_{neut}}\right) = E\left(\frac{S_{sel}/2 \cdot u_{sel}^* \cdot \mu \cdot T_{div}}{S_{neut}^*/2 \cdot u_{neut} \cdot \mu \cdot T_{div}}\right) \geq \frac{E(S_{sel})/E(u_{sel}^*)}{E(S_{neut}^*)/E(u_{neut})},$$

(6)

where $\mu$ is the per-chromosome mutation rate. The asterisks again denote our no-division-by-zero protocol. Thus $E(u_{sel}^*)$ is given by the greater of Equation 5.6 of Kimura (1957) and $1/(2 \cdot \mu \cdot T_{div})$, and $E(S_{neut}^*)$ is given by the greater of Equation 29 of Kimura (1969) and 1. The right-hand quantity in Equation 6 represents a lower bound on E(N.I.) because we have substituted the ratio of ratios of expectations for the expectation of a ratio of ratios. By Jensen's inequality, variance in either denominator will inflate the left-hand side of the equation by more than it will the right-hand side.

**Simulation of N.I. as a function of *Ns*:** In Figure 2, we present mean simulated whole-population neutrality index values (N.I.$_N$, Equation 4) under the multiplicative (Equation 1, open circles) and additive (Equation 2, solid circles) fitness schemes, for $-10/N = -0.01 \leq s \leq 10/N = 0.01$ when $N = 1000$, $\mu = 1/(2 \cdot N) = 0.0005$, $h = 1$, $Nc = 0$, $T_{div} = 30N = 30,000$ generations, and $M = B = 1$. Figure 2 also shows the diffusion-derived expression (Equation 6, solid line), which is exceeded by simulated values (under both models) for all values of *Ns*, as expected. As previously noted (Nachman 1998), N.I. is inversely related to *Ns*. The most striking pattern in the figure is the existence of a maximum N.I., the direct consequence of our no-division-by-zero protocol, which comes to dominate selected fixation values ($c_{sel}^i$) in Equation 4 and $u_{sel}$ in Equation 6. The maximum in N.I. is not the consequence of replacing the neutral segregating site count ($S_{N,neut}^i$ in Equation 4) with 1, because in our simulations of both fitness models only a very small proportion of replicates have $S_{N,neut}^i$ values equal to 0 when $Ns < 0$, and this proportion is insensitive to *Ns* [not shown, though recall that $S_{n,eut}$ is relatively insensitive to background selection (Charlesworth *et al.* 1993)]. Likewise, E($S_{neut}$) in Equa-
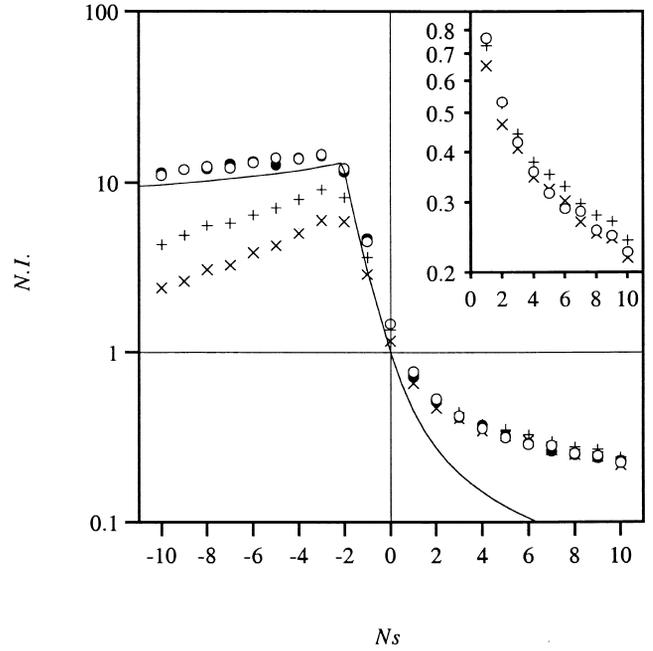


FIGURE 2.—Simulation results for values of N.I.$_N$ (Equation 4) under multiplicative (Equation 1, ○) and additive (Equation 2, ●) fitness functions for *Ns* from $-10.0$ to $10.0$. Diffusion-derived lower bound for N.I. (Equation 6, —) is shown. N.I.$_{n=10}$ (×) and N.I.$_{n=30}$ (+) (Equation 5) under multiplicative fitness function are also shown. Other parameter values: $N = 1000$, $\mu = 5 \times 10^{-4}$, $h = 1$, $Nc = 0$, $T_{div} = 30N$, and $M = B = 1$.

tion 6 is independent of *Ns* (Kimura 1969) and thus cannot be responsible for any change in slope. The location of the maximum in Figure 2 is approximately the point at which $2 \cdot u_{sel} \cdot \mu \cdot T_{div} < 1$, or equivalently, when $u_{sel} < 1/(2 \cdot \mu \cdot T_{div})$, which means that our no-division-by-zero protocol will cause N.I. to become increasingly insensitive to purifying selection as selection strength increases (driving down $u_{sel}$), as mutation rate goes down, and as divergence time decreases. The maximum empirical value of N.I. is also dependent on these parameters, and numeric values of N.I. larger than those shown in Figure 2 are possible with larger values of $\mu \cdot T_{div}$. We present results for a range of values of *s* and $\mu$, but have restricted ourselves to $T_{div} = 30N$, which we judge to be a biologically realistic number [*e.g.*, *D. melanogaster*-*D. simulans* divergence time is ~3 million years (Hey and Kliman 1993) × 10 generations/year ÷ $10^6$ effective population size (Kreitman 1983) = 30N; *H. sapiens*-*P. paniscus* divergence time is ~6 million years (Sibley 1992) ÷ 20 years/generation ÷ $10^4$ effective population size (Takahata 1993) = 30N].

Mean sample neutrality index values (N.I.$_{n=x}$, Equation 5) are also shown in Figure 2 for $x = 10$ (×) and 30 (+) drawn from populations under multiplicative fitness. Note first that the location of the maximum is unaffected by sampling. When *s* is negative, the number of segregating neutral sites in samples ($S_{n=x,neut}^{i,j}$) is again

largely independent of both sample size and strength of purifying selection (not shown), as was the number of segregating neutral sites in the whole population. Thus the location of the maximum is driven by the behavior of $c_{sel}$, which contributes equally to Equations 4 and 5. However, the value of N.I.$_{n=x}$ is conservative when $s$ is negative. Selection keeps deleterious mutations at low frequency (TAJIMA 1989), and so such sites will tend to be underrepresented in small samples. But since selection also generally prevents deleterious mutations from fixing, sampling has the effect of reducing the deviation in Equation 5 relative to neutral expectation. Because positively selected sites segregate at high frequency (TAJIMA 1989), sampling has a much smaller effect. Curiously, some underrepresentation of selected sites in very small samples ($n = 10$) occurs when $s$ is positive but small ($\leq 4/N$), thereby biasing N.I.$_{n=10}$ downward (Figure 2, inset). Thus, under weak positive selection, small sample estimates of N.I. can overstate the true population deviation from the neutral expectation, although this effect is modest.

Behavior under the additive fitness model (Equation 2) does not differ qualitatively for any parameter values examined, and no further results under this model are presented.

**Consequences of variation in $Nc$, $N$, $\mu$, $h$, and $M$ on values of N.I.:** As noted in the Introduction, the genetics of mtDNA- and nuclear-encoded genes exhibit five gross differences: recombination rate, effective population size, mutation rate, degree of selective dominance, and intralineal bottlenecks. These aspects were modeled in our simulations by the parameters $Nc$, $N$, $\mu$, $h$, and $M$, respectively, which were varied independently. Mean N.I. values from these simulations are shown in Table 5. Mean N.I. is monotonic when $Ns > -3$ (Figure 2), and selection coefficients acting on amino acid replacement mutations in mtDNA-encoded proteins have been estimated to lie in the range $-3 \leq Ns \leq 0$ (NACHMAN 1998; NIELSEN and WEINREICH 1999). In the interest of representational clarity, we now restrict ourselves to three selection coefficients, $s = -0.003$, 0.0, and 0.003, corresponding to $Ns$ of $-3.0$, 0.0, and 3.0 when $N$ is 1000. Entries in Table 5 are grouped to indicate variation of orthogonal parameters. Each entry represents a point in parameter space, and on each line results are presented in three columns, corresponding to values of $s$ equal to 0.003, 0.0, and $-0.003$. Within columns, the whole-population neutrality index (N.I.$_N$) and neutrality index values for samples of size $n = 10$ (N.I.$_{n=10}$) and 30 (N.I.$_{n=30}$) are shown to left, center, and right, respectively.

Three patterns seen in Figure 2 are also manifest in Table 5. First, in almost all cases, the inverse relationship between N.I. and $Ns$ is preserved, so that weak positive selection (represented in the left column) gives N.I. values $<1.0$ and weak purifying selection (right column) gives N.I. values $>1.0$. Second, sample neutrality index

values deviate from 1.0 less than whole-population values. And finally, sample size generally has only a modest effect on N.I.$_{n=x}$. Several additional conclusions are apparent. Most surprising to us was the general insensitivity of N.I. to recombination. In contrast, N.I.$_N$ is very sensitive to the population size (seen when $N = 10,000$ and when $M = 10$ and 100, both of which reduce the influence of genetic drift), although this sensitivity is greatly attenuated when the neutrality index is calculated for realistically sized samples. It should also be noted that small values of $N \cdot \mu$ cause a jump in the proportion of replicates in which zero segregating sites are observed (e.g., WATTERSON 1975). These zeros bias mean N.I. values down, accounting for the results shown when $N = 100$ and $\mu = 5 \times 10^{-5}$. N.I. was found to be largely insensitive to dominance, although simulations of overdominance ($s > 0$ and $h \geq 2$ or $s < 0$ and $h < 0$) could not be completed because under these conditions the number of segregating sites grew impractically large. Finally, computation time per generation of simulation increased with the number of chromosomes in the population, and the number of generations simulated increased with population size (since $T_{div} = 30 \cdot N$). Thus, $<1000$ replicates were completed for large values of $Nc$, $N$, $\mu$, and $M$.

**Consequences of variation in $Nc$, $N$, $\mu$, $h$, and $M$ on McDonald/Kreitman power:** The MCDONALD and KREITMAN (1991) test compares the ratio of segregating synonymous to amino acid replacement mutations with the ratio of fixed synonymous to amino acid replacement mutations. We performed McDonald/Kreitman tests on the simulated populations presented here as well as random samples thereof, under the assumption outlined above that synonymous mutations are selectively neutral while amino acid replacement mutations are selected. We focused on the frequency of simulated data sets in which a significant deviation is observed while $s$ is nonzero, which represents a measure of the test's statistical power to detect the action of natural selection. Since we have independently varied each of five genetic characteristics in our simulations, these data can be employed to estimate the sensitivity of power of the McDonald/Kreitman test to variations in these factors.

The proportion of replicates that give a significant McDonald/Kreitman test statistic while recombination rate, population size, mutation rate, dominance, and bottleneck size are independently varied is shown in Table 6, which has the same format as Table 5. The test was found to be more sensitive to negative selection than positive selection (AKASHI 1999), and the test's power to detect both positive and negative selection is seen to be sensitive mainly to increases in $N$ and $\mu$. (Recall that $N = 1000 \cdot M$ under the Bergstrom/Pritchard model, so increasing $M$ necessarily increases $N$.) Increasing either $N$ or $\mu$ increases the number of mutations segregating in the population (and in samples

**TABLE 5**

**Simulated neutrality index values for population and samples**

| Parameters used[a] | $s = 0.003$ | $s = 0.0$ | $s = -0.003$ |
|---|---|---|---|
| $Nc = 0.0$ | 0.436/0.408/0.447 | 1.51/1.10/1.30 | 17.28/7.51/10.97 |
| $N = 1000$ | | | |
| $\mu = 5 \times 10^{-4}$ | | | |
| $h = 1.0$ | | | |
| $M = 1$ | | | |
| $Nc = 1.0$ | 0.392/0.369/0.402 | 1.40/1.11/1.30 | 18.64/7.96/11.70 |
| $Nc = 10.0$ | 0.335/0.332/0.351[b] | 1.40/1.17/1.36 | 17.16/7.62/10.96 |
| $Nc = 100.0$ | 0.336/0.338/0.345[b] | 1.50/1.24/1.44[b] | 17.54/8.75/11.78[b] |
| $N = 100$ | 0.305/0.173/0.232 | 0.297/0.180/0.224 | 0.391/0.207/0.280 |
| $N = 10,000$ | 0.160/0.544/0.407[b] | 0.95/0.97/0.97[b] | 85.51/12.31/20.02[b] |
| $\mu = 5 \times 10^{-5}$ | 0.087/0.048/0.059 | 0.471/0.191/0.263 | 0.388/0.067/0.126 |
| $\mu = 5 \times 10^{-3}$ | 0.447/0.667/0.567[b] | 1.01/1.10/1.06[b] | 6.83/4.88/5.45[b] |
| $h = -1$ | 1.07/0.532/0.727 | 1.45/1.10/1.29 | —[c] |
| $h = 0$ | 0.562/0.401/0.479 | 1.48/1.13/1.32 | 4.94/3.59/4.36[b] |
| $h = 2$ | —[c] | 1.50/1.13/1.34 | 15.14/4.50/7.84 |
| $h = 3$ | —[c] | 1.45/1.08/1.28 | 12.27/3.02/5.50 |
| $M = 10, N = 10,000$ | 0.342/0.040/0.041 | 1.23/0.107/0.114 | 50.85/2.15/2.65 |
| $M = 100, N = 100,000$ | 0.256/0.035/0.039[d] | 1.04/0.128/0.129[d] | 39.2/2.10/3.12[d] |

Values of N.I.$_N$, N.I.$_{n=10}$, and N.I.$_{n=30}$ shown to left, center, and right, respectively. Data from 1000 evolutionary replicates except where noted.

[a] Values as on first line of table and in Figure 2 except as noted.

[b] Data from 100 evolutionary replicates.

[c] Simulation computationally impossible due to heterozygote advantage; see text.

[d] Data from 50 evolutionary replicates.

thereof), and thus by increasing the numeric values entering the $2 \times 2$ table these changes naturally increase the power of the test. Recombination increases the test's power to detect purifying selection only very slightly.

## DISCUSSION

**Patterns of polymorphism and divergence in nuclear- and mtDNA-encoded proteins differ significantly:** Mitochondrially encoded proteins exhibit a consistent pattern of excess amino acid replacement mutations segregating within species, as measured by N.I. (BALLARD and KREITMAN 1994; NACHMAN *et al.* 1994, 1996; RAND *et al.* 1994; RAND and KANN 1996, 1998; HASEGAWA *et al.* 1998; NACHMAN 1998; WISE *et al.* 1998). In contrast, nuclear-encoded proteins exhibit a variety of patterns consistent with purifying selection, neutrality, and positive selection (*e.g.*, BROOKFIELD and SHARP 1994; KREITMAN and AKASHI 1995). Our observations comparing published nuclear- and mtDNA-encoded polymorphism and divergence for protein-coding loci (Figure 1) confirm our qualitative intuition that the molecular evolution of loci in animals differs significantly as a function of the encoding genome. Whereas N.I. values for nuclear-encoded loci are evenly distributed about 1.0, the vast majority of N.I. values for mtDNA-encoded loci are >1.0.

To be fair, we did have some *a priori* expectation of this pattern (*e.g.*, MORIYAMA and POWELL 1996) before we tabulated the data in Tables 1 and 2, which may artificially inflate the significance of the values presented. But regardless of the historical contingencies by which we became aware of the signal seen in Figure 1, given the extremely small associated *P* values, patterns of nonneutral evolution clearly differ considerably between nuclear- and mtDNA-encoded proteins.

In theory, this could be a numerical artifact. Because the neutrality index (Equation 3) is a ratio of ratios, estimates of its value will be inflated by large sampling variance in either denominator (fixed amino acid replacement site or polymorphic synonymous site counts). Thus shorter genes or smaller population samples will bias the N.I. upward. And indeed, the mean gene lengths are significantly less in the mtDNA-encoded data set (mean number of amino acids ± SD encoded in nuclear data set, 381.09 ± 152.19; in mtDNA data set, 260.65 ± 156.67, $t = 3.18$, d.f. = 65, $P = 0.0011$). However, four of six mtDNA-encoded loci with N.I. values <1.0 are shorter than the sample mean, and among both data sets, N.I. is uncorrelated with gene length (not shown). Moreover, no significant difference in length exists between the 15 longest mtDNA-encoded genes and the entire nuclear data set ($t = 0.167$, d.f. = 47,

TABLE 6

**Proportion of simulation replicates with significant ($P < 0.05$)
McDonald and Kreitman (1991)
test results**

| Parameters used[a] | $s = 0.003$ | $s = 0.0$ | $s = -0.003$ |
|---|---|---|---|
| $Nc = 0.0$ | 33.2/11.5/16.1 | 7.8/4.9/6.2 | 71.7/18.7/28.1 |
| $N = 1000$ | | | |
| $\mu = 5 \times 10^{-4}$ | | | |
| $h = 1.0$ | | | |
| $M = 1$ | | | |
| $Nc = 1.0$ | 36.2/12.6/18.3 | 6.1/6.1/6.7 | 73.1/19.0/28.5 |
| $Nc = 10.0$ | 44.5/17.6/24.2[b] | 7.0/6.8/7.9 | 75.2/17.8/28.7 |
| $Nc = 100.0$ | 43.8/17.5/24.5[b] | 3.8/6.1/7.2[b] | 79.0/21.1/30.6[b] |
| $N = 100$ | 1.1/0.6/0.7 | 0.6/0.2/0.1 | 1.0/0.2/0.1 |
| $N = 10,000$ | 100.0/23.7/52.6[b] | 9.1/1.8/5.9[b] | 100.0/29.8/63.8[b] |
| $\mu = 5 \times 10^{-5}$ | 10.9/4.2/5.8 | 2.1/0.6/0.8 | 4.0/0/0 |
| $\mu = 5 \times 10^{-3}$ | 89.2/24.5/42.7[b] | 4.0/3.7/4.1[b] | 100.0/62.6/87.2[b] |
| $h = -1$ | 7.8/6.4/8.3 | 7.3/6.1/7.1 | —[c] |
| $h = 0$ | 21.4/11.7/14.5 | 7.2/5.8/7.7 | 22.5/16.1/18.8[b] |
| $h = 2$ | —[c] | 6.9/5.8/6.8 | 73.4/11.8/19.6 |
| $h = 3$ | —[c] | 9.0/6.7/8.4 | 72.8/7.2/13.5 |
| $M = 10, N = 10,000$ | 56.9/10.4/12.2 | 5.0/6.3/7.5 | 94.8/33.6/40.8 |
| $M = 100, N = 100,000$ | 100/15.0/18.0[d] | 3.9/8.6/9.2[d] | 100.3/9.5/51.9[d] |

Values for McDonald/Kreitman tables using polymorphic site counts in the population and in samples of size 10 and 30 shown to left, center, and to right, respectively; proportions expressed in percentages. Data from 1000 evolutionary replicates except where noted, and from 10 sampling replicates per evolutionary replicate.

[a] Values as on first line of table and in Figure 2 except where noted.
[b] Data from 100 evolutionary replicates.
[c] Simulation computationally impossible due to heterozygote advantage; see text.
[d] Data from 50 evolutionary replicates.

$P = 0.43$) although a highly significant association between N.I. values $>1.0$ and genome persists (partitioning N.I. values as greater than and $<1.0$: $G = 11.06$, d.f. $= 1, P = 0.0009$). Similarly, no significant difference in length exists between the 18 shortest nuclear-encoded genes and the entire mitochondrial data set ($t = 0.270$, d.f. $= 49, P = 0.39$), but again a highly significant association between N.I. and genome is detected ($G = 10.36$, d.f. $= 1, P = 0.0013$). Thus the pattern in Figure 1 seems not to be driven by any bias in gene lengths. Sample sizes also differ significantly ($n \pm SD$ for nuclear data sets, $17.69 \pm 14.98$; for mtDNA data sets, $28.99 \pm 26.91$, $t = 1.99$, d.f. $= 65, P = 0.025$); however, the larger average mtDNA data set should reduce variance in those estimates and bias estimates of N.I. downward, suggesting that the reported genome-specific difference in N.I. may be conservative. Thus differences in sampling variance (either in gene length or sample size) cannot account for the pattern in Figure 1.

The McDonald/Kreitman test is insensitive to populations not at equilibrium, to recombination, and to variation in nucleotide mutation rates (McDonald and

Kreitman 1991), and we believe the neutrality index is similarly robust. Furthermore, although different sampling strategies may underlie the nuclear and mitochondrial data sets, we do not believe that these differences bias our analysis because N.I. is based on segregating site counts rather than on site frequencies. For example, several of our nuclear data sets (*e.g.*, *Adh* and *Est*-6) represent explicitly stratified samples intended to include previously described allozyme classes. However, because truly random samples of even modest size would be expected to include allozymes segregating at moderate frequency, stratification will have only a marginal effect on segregating site counts. Additionally, since allozyme classes are the consequence of amino acid replacement mutations, one would expect that the intentional addition of very rare allozyme classes to one of our data sets would inflate segregating amino acid replacement site counts, biasing N.I. upward. Since nuclear data sets appear on average to suffer a deficit of segregating amino acid replacement sites, this effect would seem to make estimates of N.I. conservative. Several of our mitochondrial data sets may include some geographic stratification, but this is also unlikely to affect our analy-

sis. The local fixation of standing variation will not inflate N.I. as long as the haplotypes are fixed at random with respect to their number of segregating replacement sites. Although multiple-niche models of balancing selection are possible (*e.g.*, LEVINE 1953), in cases where this effect was explicitly tested for in mtDNA, no support for this model was found (FRY and ZINK 1998; BROWN *et al.* 2000). Finally, mutations under balancing selection are expected to persist in the population, but in the mtDNA data sets employed here, no such evidence exists (NIELSEN and WEINREICH 1999).

Alternatively, natural selection acting on synonymous mutations could be responsible for the pattern seen in Figure 1. For example, in *D. simulans*, segregating unpreferred synonymous mutations are overrepresented relative to fixations (AKASHI 1996), which will bias N.I. values downward. However, natural selection is similarly known to act on synonymous mutations in mtDNA-encoded genes in several species of Drosophila (BALLARD and KREITMAN 1994; RAND and KANN 1998). Thus we do not believe that natural selection acting on codon bias is a major factor contributing to the nuclear-mitochondrial contrast we report. Another possible explanation for the pattern in Figure 1 is the much broader species representation among the mitochondrial data sets. However, confining ourselves to data sets from *D. melanogaster* and *D. simulans*, a highly significant association between genome and N.I. value still exists ($G = 7.48$; d.f. $= 1$; $P = 0.006$). Although there are two drosophilid mitochondrial data sets with N.I. $<1.0$, both come from *D. pseudoobscura* and may reflect a recent population expansion in that species (RAND and KANN 1998; HAMBLIN and AQUADRO 1999). Thus the much broader species representation among mtDNA data sets cannot explain the pattern we describe.

There are only 13 mtDNA-encoded loci in metazoans (WOLSTENHOLME 1992), and thus the 31 loci in Table 2 necessarily include multiple data sets for single loci, although all such duplicates are from different species. This suggests the possibility that the statistical significance seen in Figure 1 could be the consequence of repeatedly sampling from correlated evolutionary processes. However, the association between genomes and N.I. (partitioned as N.I. $> 1.0$ and N.I. $< 1.0$) is preserved when a single data set for each locus is randomly chosen from Table 2 ($G = 6.79$, d.f. $= 1$, $P = 0.0092$). More generally, all mtDNA-encoded proteins are constituents of the enzymes responsible for oxidative phosphorylation (OXPHOS), whereas none of the nuclear-encoded proteins in Table 1 are. This common functionality among mtDNA-encoded loci might cause an evolutionary correlation, driving the observed patterns of polymorphism and divergence. For example, nearly all the nuclear-encoded proteins in Table 1 are soluble whereas OXPHOS enzymes all reside in the inner mitochondrial membrane and are extremely hydrophobic (GILLHAM 1994). It is known that strong

purifying selection acts to eliminate hydrophilic amino acids from mtDNA-encoded proteins (see NAYLOR *et al.* 1995). However, among nuclear-encoded proteins, no correlation exists between N.I. and hydrophobicity (scored by the method of KYTE and DOOLITTLE 1982; not shown). Nevertheless, the possibility that unique selective forces are acting on (mtDNA-encoded) OXPHOS proteins cannot be dismissed. An obvious approach to this question is to examine the polymorphism and divergence patterns in some of the nuclear-encoded OXPHOS proteins, an avenue that we are currently pursuing.

**Genetic factors alone seem unable to account for empirical patterns in neutrality index values:** Our simulations (Figure 2) repeat the observation that N.I. is quite sensitive to *Ns*, the strength and direction of selection (NACHMAN 1998). Thus progressively stronger positive selection increases the selected site fixation count and reduces N.I. monotonically for all parameter values examined. Progressively stronger purifying selection depresses the selected site fixation count and increases N.I., although under the no-division-by-zero protocol employed, N.I. is not permitted to climb to infinity. While we find some effect on N.I. for most of the genetic factors examined (Table 5), no single factor breaks this roughly inverse relationship between *s* and *N.I* for biologically realistic parameter values. Population size, which had the strongest influence, predicts *smaller* deviations from 1.0 in mitochondrial N.I. since mtDNA have smaller effective populations, whereas larger effects are empirically observed. Moreover, population size had its effect greatly attenuated when sample N.I. means were measured. Thus, if we wish to regard our samples of nuclear- and mtDNA-encoded genes as multiple realizations of a single evolutionary process, we are at present unable to appeal to genetic differences between genomes to account for the pattern seen in Figure 1. We acknowledge that we have not explored the effects of interaction among these factors due to the prohibitive amount of computation time required for a thorough exploration of parameter space.

As noted, only 17% (6 of 36) of the nuclear data sets in Table 1 show a significant deviation from neutral expectation by the MCDONALD and KREITMAN (1991) test, while 48% (15 of 31) of the mitochondrial data sets in Table 2 do. Inasmuch as many of the nuclear samples were constructed with an *a priori* intuition that selection might be working, while many of the mtDNA data sets were constructed to explore questions of phylogeography, this contrast is perhaps conservative. However, the results of our McDonald/Kreitman power analysis may account for this pattern. Although the *effective* population size of mtDNA is less than that of nuclear DNA (BIRKY *et al.* 1983), the *census* number of mtDNA molecules in a population is much larger (GILLHAM 1994). Furthermore, in at least some animal species, mtDNA mutation rates are higher than those of nuclear

DNA (Avise 1991). Both of these factors increase the statistical power of the test, particularly for sample sizes used in this study irrespective of the sign of *s* (Table 6). Thus if one regards the data in Tables 1 and 2 as repeated samples drawn from a single evolutionary process, on the basis of our power analysis one would predict a greater incidence of significant McDonald/ Kreitman test statistics among the mitochondrially encoded data sets. However, these results shed no light on the cause of the highly significant genome-specific differentiation in the *direction* of deviation among these data sets.

**The biological importance of the frequency distribution of *s*:** At present, there is little support for the hypothesis that unique selective forces acting on mitochondrially encoded OXPHOS proteins explain the pattern shown in Figure 1, although we cannot rule out this possibility. And no single genetic difference between nuclear and mtDNA genetics examined appears sufficient to explain this pattern. However, both our simulations and analytic expectations assume that *s* is equal for all selected mutations entering the population (we have not included deleterious mutations of large effect since such mutations contribute very little to polymorphism or divergence). This assumption of a single fixed *s* is clearly simplistic; indeed, it is theoretically problematic (Gillespie 1995). However, very little is known about the true frequency distribution of *s*. Although it is likely that the majority of amino acid replacement mutations are slightly deleterious (Ohta 1973), it seems equally likely that some are also advantageous (Gillespie 1995). Nevertheless, we believe the interplay between recombination and a common distribution of mutational *s* could at least in part be responsible for the striking contrast seen in Figure 1. If the majority of mutations that contribute to polymorphism are indeed slightly deleterious, then we reason that the patterns of polymorphism and divergence for a tightly linked chromosome will be dominated by those mutations, and that when the occasional mutation with a positive *s* occurs, it will be able to contribute to the process only if it happens to land on a relatively "unloaded" copy, an unlikely event. (Or similarly unlikely, an advantageous mutation would reach fixation in the absence of recombination only if it were sufficiently strongly selected to offset the cumulative effect of the deleterious mutations to which it was linked.) This may account for the pattern seen in mtDNA-encoded proteins, where N.I. values >1.0 predominate. In contrast, the same small fraction of advantageous mutations landing on a recombining chromosome will be more likely to get onto an unloaded segment of the chromosome before being lost. Once on an unloaded chromosomal segment, such a mutation will have its advantage expressed, resulting in differential reproductive output, and therefore will have its frequency increased by selection. We believe this could account for the pattern seen

in nuclear-encoded proteins, where N.I. values appear to be evenly distributed ~1.0. [It should be noted that although indirect evidence of recombination in animal mtDNA has recently accumulated (Lunt and Hyman 1997; Awadalla *et al.* 1999; Eyre-Walker *et al.* 1999), very low levels of mitochondrial recombination are suggested (Eyre-Walker *et al.* 1999).]

Our hypothesis predicts that empirical neutrality index values should be inversely correlated to recombination rate, although among the nuclear genes in Table 1 for which we were able to find published estimates of recombination rate, no correlation exists. Moreover, in a cursory exploration of selective frequency distribution space we were unable to find parameter values in which this effect was observed. Recently, Gillespie (1999) explored the behavior of N.I. ($\mathcal{MK}$ in his notation) under several more sophisticated frequency distributions of *s*. His simulations compared N.I. under free recombination and complete linkage as a function of population size, but, like us, he was unable to find a case in which linkage carried N.I. from ~1.0 to considerably larger values.

However, suggestive comparisons emerge from DNA polymorphism and divergence data recently accumulated from the plant *A. thaliana* (Table 3). *A. thaliana* is almost exclusively self-fertilizing, and its effective recombination rate is consequently very low (Kamabe and Miyashita 1999). Consistent with our hypothesis, N.I. is >1.0 for five of the six nuclear genes in *A. thaliana* in Table 3. Of course there are many other biological differences between Arabidopsis and Drosophila evolution that may be responsible for this observation.

Another intriguing system is the $O_{st}/O_{3+4}$ chromosomal inversion in *D. subobscura*. *Acph-1* lies very near one of the inversion breakpoints (Segarra *et al.* 1996), and since recombination between inversion haplotypes is greatly suppressed near breakpoints, under random mating the recombination rate at *Acph-1* within karyotype will be proportional to $p^2$, where *p* is the frequency of the karyotype in question. Thus our hypothesis predicts that neutrality index values calculated at *Acph-1* within karyotypes should be correlated with the square of karyotype frequency. And indeed, in a sample of *D. subobscura* taken from a population in which the frequency of $O_{3+4}$ was estimated as 0.767 and of $O_{st}$ as 0.147 (Navarro-Sabaté *et al.* 1999), N.I. values within the former karyotype are much lower (1.70) than within the latter (9.0; N.I. calculated from data in Navarro-Sabaté *et al.* 1999). Since these karyotypes exhibit a latitudinal cline (Navarro-Sabaté *et al.* 1999), this system offers the possibility of varying effective recombination rate while holding gene function constant by sampling from different points along the cline and calculating N.I. within the karyotype.

Finally, several groups (Braverman *et al.* 1997; Schug *et al.* 1998; Jensen *et al.* 1999) are exploring the interaction of recombination rate and levels of putatively silent polymorphism. Surprisingly, there are few Drosophila

data sets for protein-coding loci in regions of lowest recombination. We are now beginning to collect additional polymorphism and divergence data from such loci located on the tip of the X and on the fourth chromosome of *D. melanogaster*, regions of low recombination, to test the joint predictions that our hypothesis makes.

While our simulations revealed no single genetic factor to account for the marked difference in patterns of nonneutral evolution seen in nuclear- and mtDNA-encoded proteins (Figure 1), we suggest two (nonexclusive) hypotheses. Figure 2 and Table 5 demonstrate that N.I. is inversely related to *Ns*, so that if the selective histories of the genes in Tables 1 and 2 are distinct, N.I. will be affected. Thus, if the fraction of mildly deleterious amino acid replacement mutations entering OXPHOS genes is larger than the corresponding fraction for nuclear loci (or equivalently if the opportunities for positive selection are greater for nuclear-encoded loci), mtDNA-encoded N.I. values will be biased upward. Additionally, we speculate that genetic linkage in mtDNA results in patterns of polymorphism and divergence that are dominated by the largest class of mutations entering the population. If the frequency distribution of selection coefficients is such that a majority of mutations that contribute to polymorphism and divergence are mildly deleterious, values of N.I. >1 may result in regions of low recombination. Both hypotheses are open to experimental attack.

*Note added in proof*: Polymorphic and fixed site counts in Tables 1–3 were taken from the citations given therein. Subsequent analysis of the GenBank sequence data for some of these genes revealed small differences in some counts (C. BUSTAMANTE and B. CEZAIRLIAYN, personal communication) and concomitant differences in values of N.I. However, these differences do not materially affect the results or conclusions of this article.

## LITERATURE CITED

AGUADÉ, M., 1998 Different forces drive the evolution of the Acp26Aa and Acp26Ab accessory gland genes in *Drosophila melanogaster* species complex. Genetics **150:** 1079–1089.

AGUADÉ, M., 1999 Positive selection drives evolution of the Acp29AB accessory gland protein locus in Drosophila. Genetics **152:** 543–551.

AKASHI, H., 1995 Inferring weak selection from patterns of polymorphism and divergence at "silent" sites in Drosophila DNA. Genetics **139:** 1067–1076.

AKASHI, H., 1996 Molecular evolution between *Drosophila melanogaster* and *D. simulans*: reduced codon bias, faster rates of amino acid substitution, and larger proteins in *D. melanogaster*. Genetics **144:** 1297–1307.

AKASHI, H., 1999 Inferring the fitness effects of DNA mutations from polymorphism and divergence data: statistical power to detect directional selection under stationarity and free recombination. Genetics **151:** 221–238.

AVISE, J. C., 1991 Ten unorthodox perspectives on evolution prompted by comparative population genetic findings on mitochondrial DNA. Annu. Rev. Genet. **25:** 45–69.

AWADALLA, P., A. EYRE-WALKER and J. MAYNARD SMITH, 1999 Linkage disequilibrium and recombination in hominid mitochondrial DNA. Science **286:** 2524–2525.

BALAKIREV, E. S., E. I. BALAKIREV, F. RODRÍGUEZ-TRELLES and F. J. AYALA, 1999 Molecular evolution of two linked genes, *Est-6* and *Sod*, in *Drosophila melanogaster*. Genetics **153:** 1357–1369.

BALLARD, J. W. O., and M. KREITMAN, 1994 Unraveling selection in the mitochondrial genome of Drosophila. Genetics **138:** 757–772.

BEGUN, D. J., and C. F. AQUADRO, 1994 Evolutionary inferences from DNA variation at the 6-phosphogluconate dehydrogenase locus in natural populations of Drosophila: selection and geographic differentiation. Genetics **136:** 155–171.

BENDALL, K. E., V. A. MACAULAY, J. R. BAKER and B. C. SYKES, 1996 Heteroplasmic point mutations in the human mtDNA control region. Am. J. Hum. Genet. **59:** 1276–1287.

BERGSTROM, C. T., and J. PRITCHARD, 1998 Germline bottlenecks and the evolutionary maintenance of mitochondrial genomes. Genetics **149:** 2135–2146.

BIRKY, C. W., JR., T. MARUYAMA and P. FUERST, 1983 An approach to population and evolutionary genetic theory for genes in mitochondria and chloroplasts, and some results. Genetics **103:** 513–527.

BRAVERMAN, J., M. AGUADÉ and C. LANGLEY, 1997 Reduced level of DNA sequence variation at the *erect wing* locus of *D. melanogaster* and *D. simulans*, p. 234A in *Proceedings of the 38th Annual Drosophila Research Conference, Chicago, April 1997*. Genetics Society of America, Bethesda, MD.

BROOKFIELD, J. F. Y., and P. M. SHARP, 1994 Neutralism and selectionism face up to DNA data. Trends Genet. **10:** 109–111.

BROWN, A. F., L. M. KANN and D. M. RAND, 2000 Gene flow versus local adaptation in the Northern acorn barnacle *Semibalanus balanoides*: insights from mtDNA control region polymorphism. Evolution (in press).

CHARLESWORTH, B., M. T. MORGAN and D. CHARLESWORTH, 1993 The effect of deleterious mutations on neutral molecular variation. Genetics **134:** 1289–1303.

EANES, W. F., M. KIRCHNER and J. YOON, 1993 Evidence for adaptive evolution of the g6pd gene in *Drosophila melanogaster* and *Drosophila simulans* lineages. Proc. Natl. Acad. Sci. USA **90:** 7475–7479.

EYRE-WALKER, A., N. H. SMITH and J. MAYNARD SMITH, 1999 How clonal are human mitochondria? Proc. R. Soc. Lond. Ser. B **266:** 477–483.

FRY, A. J., and R. M. ZINK, 1998 Geographic analysis of nucleotide diversity and song sparrow (Aves: Emberizidae) population history. Mol. Ecol. **7:** 1303–1313.

GILLESPIE, J. H., 1993 Substitutional processes in molecular evolution. I. Uniform and clustered substitutions in a haploid model. Genetics **134:** 971–981.

GILLESPIE, J. H., 1995 On Otha's hypothesis: most amino acid substitutions are deleterious. J. Mol. Evol. **40:** 64–69.

GILLESPIE, J. H., 1999 The role of population size in molecular evolution. Theor. Popul. Biol. **55:** 145–156.

GILLHAM, N. W., 1994 *Organelle Genes and Genomes*. Oxford University Press, New York.

GLEASON, J. M., and J. R. POWELL, 1997 Interspecific and intraspecific comparisons of the *period* locus in the *Drosophila willistoni* sibling species. Mol. Biol. Evol. **14:** 741–753.

HAMBLIN, M. T., and C. F. AQUADRO, 1996 High nucleotide sequence variation in a region of low recombination in *Drosophila simulans* is consistent with the background selection model. Mol. Biol. Evol. **13:** 1133–1140.

HAMBLIN, M. T., and C. F. AQUADRO, 1999 DNA sequence variation and the recombinational landscape in *Drosophila pseudoobscura*: a study of the second chromosome. Genetics **153:** 859–869.

HASEGAWA, M., Y. CAO and Z. YANG, 1998 Preponderance of slightly deleterious polymorphism in mitochondrial DNA: nonsynonymous/synonymous rate ratio is much higher within species than between species. Mol. Biol. Evol. **15:** 1499–1505.

HASSON, E., I.-N. WANG, L.-W. ZENG, M. KREITMAN and W. EANES, 1999 Nucleotide variation in the Triosephosphate Isomerase (*Tpi*) locus of *Drosophila melanogaster* and *Drosophila simulans*. Mol. Biol. Evol. **15:** 756–768.

HAUSWIRTH, W. W., and P. J. LAIPIS, 1982 Mitochondrial DNA poly-

morphism in a maternal lineage of Holstein cows. Proc. Natl. Acad. Sci. USA **79:** 4686–4690.

HEY, J., and R. M. KLIMAN, 1993 Population genetics and phylogenetics of DNA sequence variation at multiple loci within the *Drosophila melanogaster* species complex. Mol. Biol. Evol. **10:** 804–822.

HUGHES, A. L., and M. NEI, 1988 Patterns of nucleotide substitution at major histocompatibility complex class I loci reveal overdominant selection. Nature **335:** 167–170.

JENSEN, M., M. KREITMAN and B. CHARLESWORTH, 1999 Background selection predominates on the fourth chromosome of *D. melanogaster*, p. a223 in *Proceedings of the 40th Annual Drosophila Research Conference.* Bellevue, WA. Genetics Society of America.

JENUTH, J. P., A. C. PETERSON, K. FU and E. A. SHOUBRIDGE, 1996 Random genetic drift in female germline explains the rapid segregation of mammalian mitochondrial DNA. Nat. Genet. **14:** 146–151.

KAMABE, A., and N. T. MIYASHITA, 1999 DNA variation in the basic Chitinase locus (*ChiB*) region of the wild plant *Arabidopsis thaliana*. Genetics **153:** 1445–1453.

KAMABE, A., H. INNAN, R. TERAUCHI and N. MIYASHITA, 1997 Nucleotide polymorphism in the Acidic Chitinase locus (*ChiA*) region of the wild plant *Arabidopsis thaliana*. Mol. Biol. Evol. **14:** 1301–1315.

KIMURA, M., 1957 Some problems of stochastic processes in genetics. Ann. Math. Stat. **28:** 882–901.

KIMURA, M., 1969 The number of heterozygous nucleotide sites maintained in a finite population due to steady flux of mutations. Genetics **61:** 893–903.

KIMURA, M., 1983 *The Neutral Theory of Molecular Evolution.* Cambridge University Press, Cambridge, United Kingdom.

KING, L. M., 1998 The role of gene conversion in determining sequence variation and divergence in the *Est-5* gene family in *Drosophila pseudoobscura*. Genetics **148:** 305–315.

KREITMAN, M., 1983 Nucleotide polymorphism at the alcohol dehydrogenase locus of Drosophila melanogaster. Nature **304:** 412–417.

KREITMAN, M., and H. AKASHI, 1995 Molecular evidence for natural selection. Annu. Rev. Ecol. Syst. **26:** 403–422.

KREITMAN, M., and R. R. HUDSON, 1991 Inferring the evolutionary histories of the Adh and Adh-dup loci in Drosophila melanogaster from patterns of polymorphism and divergence. Genetics **127:** 565–582.

KYTE, J., and R. F. DOOLITTLE, 1982 A simple method for displaying the hydropathic character of a protein. J. Mol. Biol. **157:** 105–132.

LABATE, J. A., C. H. BIERMANN and W. F. EANES, 1999 Nucleotide variation at the *runt* locus in *Drosophila melanogaster* and *Drosophila simulans*. Mol. Biol. Evol. **16:** 724–731.

LEVINE, H., 1953 Genetic equilibrium when more than one ecological niche is available. Am. Nat. **87:** 331–333.

LONG, M., and C. H. LANGLEY, 1993 Natural selection and the origin of *jingwei*, a chimeric processed functional gene in *Drosophila*. Science **260:** 91–95.

LUNT, D. H., and B. C. HYMAN, 1997 Animal mitochondrial DNA recombination. Nature **387:** 247.

MAYNARD SMITH, J., 1994 Estimating selection by comparing synonymous and substitutional changes. J. Mol. Evol. **39:** 123–128.

MCDONALD, J. H., and M. KREITMAN, 1991 Adaptive protein evolution at the *Adh* locus in *Drosophila*. Nature **351:** 652–654.

MESSIER, W., and C.-B. STEWART, 1997 Episodic adaptive evolution of primate lysozymes. Nature **385:** 151–154.

MIYASHITA, N. T., A. KAWABE, H. INNAN and R. TERAUCHI, 1998 Intra- and interspecific DNA variation and codon bias of the alcohol dehydrogenase (*Adh*) Locus in *Arabis* and *Arabidopsis thaliana*. Mol. Biol. Evol. **15:** 1420–1429.

MORITZ, C., T. E. DOWLING and W. M. BROWN, 1987 Evolution of animal mitochondrial DNA: relevance for population biology and systematics. Annu. Rev. Ecol. Syst. **18:** 269–292.

MORIYAMA, E. N., and J. R. POWELL, 1996 Intraspecific nuclear DNA variation in *Drosophila*. Mol. Biol. Evol. **13:** 261–277.

NACHMAN, M. W., 1998 Deleterious mutations in animal mitochondrial DNA. Genetica **102/103:** 61–69.

NACHMAN, M. W., S. N. BOYER and C. F. AQUADRO, 1994 Nonneutral evolution at the mitochondrial NADH dehydrogenase subunit 3 gene in mice. Proc. Natl. Acad. Sci. USA **91:** 6364–6368.

NACHMAN, M. W., W. M. BROWN, M. STONEKING and C. F. AQUADRO,

1996 Nonneutral mitochondrial DNA variation in humans and chimpanzees. Genetics **142:** 953–963.

NAVARRO-SABATÉ, A. M., M. AGUADÉ and C. SEGARRA, 1999 The relationship between allozyme and chromosomal polymorphism inferred from nucleotide variation at the *Acp-1* gene region of *Drosophila subobscura*. Genetics **153:** 871–889.

NAYLOR, G. J., T. M. COLLINS and W. M. BROWN, 1995 Hydrophobicity and phylogeny. Nature **373:** 565–566.

NIELSEN, R., and D. WEINREICH, 1999 The age of nonsynonymous and synonymous mutations in animal mtDNA and implications for the mildly deleterious theory. Genetics **153:** 497–506.

OHTA, T., 1972 Evolutionary rate of cistrons and DNA divergence. J. Mol. Evol. **1:** 150–157.

OHTA, T., 1973 Slightly deleterious mutant substitutions in evolution. Nature **246:** 96–98.

OHTA, T., and M. KIMURA, 1971 On the constancy of the evolutionary rate of cistrons. J. Mol. Evol. **1:** 18–25.

PARSONS, T. J., D. S. MUNIEC, K. SULLIVAN, N. WOODYATT, R. ALLISTON-GREINER *et al.*, 1997 A high observed substitution rate in the human mitochondrial DNA control region. Nat. Genet. **15:** 363–368.

PRESS, W. H., S. A. TEUKOLSKY, W. T. VETTERLING and B. P. FLANNERY, 1992 *Numerical Recipes in C.* Cambridge University Press, Cambridge, United Kingdom.

PURUGGANAN, M. D., and J. I. SUDDITH, 1998 Molecular population genetics of the *Arabidopsis CAULIFLOWER* regulatory gene: nonneutral evolution and naturally occurring variation in floral homeotic function. Proc. Natl. Acad. Sci. USA **95:** 8130–8134.

PURUGGANAN, M. D., and J. I. SUDDITH, 1999 Molecular population genetics of floral homeotic loci: departures from the equilibrium-neutral model at the *APETALA3* and *PISTILLATA* genes of *Arabidopsis thaliana*. Genetics **151:** 839–848.

RAMOS-ONSINS, S., and M. AGUADÉ, 1998 Molecular evolution of the *Cecropin* multigene family in Drosophila: functional genes *vs.* pseudogenes. Genetics **150:** 157–159.

RAND, D. M., and L. M. KANN, 1996 Excess amino acid polymorphism in mitochondrial DNA: contrasts among genes from Drosophila, mice and humans. Mol. Biol. Evol. **13:** 735–748.

RAND, D. M., and L. M. KANN, 1998 Mutation and selection at silent and replacement sites in the evolution of animal mitochondrial DNA. Genetica **102/103:** 393–407.

RAND, D. M., M. DORFSMAN and L. M. KANN, 1994 Neutral and nonneutral evolution of Drosophila mitochondrial DNA. Genetics **138:** 741–756.

SCHAEFFER, S. W., and E. L. MILLER, 1992 Molecular population genetics of an electrophoretically monomorphic protein in the alcohol dehydrogenase region of *Drosophila pseudoobscura*. Genetics **132:** 163–178.

SCHMID, K. J., L. NIGRO, C. F. AQUADRO and D. TAUTZ, 1999 Large number of replacement polymorphisms in rapidly evolving genes of Drosophila: implications for genome wide surveys. Genetics **153:** 1717–1729.

SCHUG, M. D., C. M. HUTTER, M. A. F. NOOR and C. F. AQUADRO, 1998 Mutation and evolution of microsatellites in *Drosophila melanogaster*. Genetica **102/103:** 359–367.

SEGARRA, C., G. RIBÒ and M. AGUADÉ, 1996 Differentiation of Muller's chromosomal elements *D* and *E* in the Obscura group of Drosophila. Genetics **144:** 139–146.

SIBLEY, C. G., 1992 DNA-DNA hybridization and the study of primate evolution, pp. 313–316 in *The Cambridge Encyclopedia of Human Evolution,* edited by S. JONES, R. MARTIN and D. PILBEAM. Cambridge University Press, Cambridge, UK.

TAJIMA, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. Genetics **123:** 585–595.

TAKAHATA, N., 1993 Allelic genealogy and human evolution. Mol. Biol. Evol. **10:** 2–22.

TEMPLETON, A. R., 1996 Contingency tests of neutrality using intra/interspecific gene trees: the rejection of neutrality for the evolution of the mitochondrial cytochrome oxidase ii gene in hominoid primates. Genetics **144:** 1263–1270.

TSAUR, S.-C., C.-T. TING and C.-I. WU, 1998 Positive selection driving the evolution of a gene of male reproduction, *Acp26Aa*, of Drosophila: II. Divergence versus polymorphism. Mol. Biol. Evol. **15:** 1040–1046.

WATTERSON, G. A., 1975 On the number of segregating sites in

genetical models without recombination. Theor. Popul. Biol. **7:** 256–276.

Wayne, M. L., and M. Kreitman, 1996   Reduced variation at *concertina*, a heterochromatic locus in *Drosophila*. Genet. Res. **68:** 101–108.

Wayne, M. L., D. Contamine and M. Kreitman, 1996   Molecular population genetics of *ref(2)P*, a locus which confers viral resistance in *Drosophila*. Mol. Biol. Evol. **13:** 191–199.

Wise, C. A., M. Sram and S. Easteal, 1998   Departure from neutrality at the mitochondrial NADH dehydrogenase subunit 2 gene in humans, but not in chimpanzees. Genetics **148:** 409–421.

Wolstenholme, D. R., 1992   Animal mitochondrial DNA: structure and evolution. Int. Rev. Cytol. **141:** 173–216.

Zurovcova, M., and W. F. Eanes, 1999   Lack of nucleotide polymorphism in the *Y*-linked sperm flagellar dynein gene *Dhc-Yh3* of *Drosophila melanogaster* and *D. simulans*. Genetics **153:** 1709–1715.

Communicating editor: A. G. Clark