

Positive and Negative Selection on Noncoding DNA in *Drosophila simulans*

Penelope R. Haddrill,* Doris Bachtrog,† and Peter Andolfatto‡

*Institute of Evolutionary Biology, School of Biological Sciences, University of Edinburgh, Edinburgh, United Kingdom; †Section of Ecology, Behavior and Evolution, Division of Biological Sciences, University of California, San Diego; and ‡Department of Ecology and Evolutionary Biology and Lewis-Sigler Institute for Integrative Genomics, Princeton University

There is now a wealth of evidence that some of the most important regions of the genome are found outside those that encode proteins, and noncoding regions of the genome have been shown to be subject to substantial levels of selective constraint, particularly in *Drosophila*. Recent work has suggested that these regions may also have been subject to the action of positive selection, with large fractions of noncoding divergence having been driven to fixation by adaptive evolution. However, this work has focused on *Drosophila melanogaster*, which is thought to have experienced a reduction in effective population size (N_e), and thus a reduction in the efficacy of selection, compared with its closest relative *Drosophila simulans*. Here, we examine patterns of evolution at several classes of noncoding DNA in *D. simulans* and find that all noncoding DNA is subject to the action of negative selection, indicated by reduced levels of polymorphism and divergence and a skew in the frequency spectrum toward rare variants. We find that the signature of negative selection on noncoding DNA and nonsynonymous sites is obscured to some extent by purifying selection acting on preferred to unpreferred synonymous codon mutations. We investigate the extent to which divergence in noncoding DNA is inferred to be the product of positive selection and to what extent these inferences depend on selection on synonymous sites and demography. Based on patterns of polymorphism and divergence for different classes of synonymous substitution, we find the divergence excess inferred in noncoding DNA and nonsynonymous sites in the *D. simulans* lineage difficult to reconcile with demographic explanations.

Introduction

Noncoding DNA makes up a large fraction of most eukaryotic genomes, and yet relatively little is known about the functional importance of sequences that are not translated into proteins. Recently, a number of multilocus and comparative genomics studies have examined patterns of molecular evolution in noncoding DNA in various *Drosophila* species, and these have revealed slower rates of evolution and higher levels of selective constraint in long (>80 base pairs [bp]) introns and intergenic sequences, when compared with synonymous sites in coding regions (Bergman and Kreitman 2001; Halligan et al. 2004; Kohn et al. 2004; Andolfatto 2005; Haddrill, Charlesworth, et al. 2005; Bachtrog and Andolfatto 2006; Halligan and Keightley 2006). This is assumed to be due to the presence of *cis*-regulatory elements (Casillas et al. 2007) or conserved RNA secondary structures, for which there is some direct experimental evidence (e.g., Stephan and Kirby 1993; Kirby et al. 1995; Leicht et al. 1995; Carlini et al. 2001; Chen and Stephan 2003; Bergman et al. 2005; Gallo et al. 2006).

Although these comparative genomics approaches show that divergence between species is reduced in some classes of noncoding DNA, indicating that these sequences are functionally constrained and thus subject to negative selection, this conclusion cannot be firmly supported without evidence from within-species patterns of variability to rule out mutation rate variation. Similarly, if noncoding sequences are functionally important, they are likely to be subject to the action of positive selection as well as negative selection, and the signatures of these different types of selection can only be distinguished using an approach that combines

within-species polymorphism data with between-species measures of divergence (McDonald and Kreitman 1991). This type of approach has previously provided evidence that a considerable proportion of the protein-coding sequence divergence between species has been driven to fixation by positive selection (Fay et al. 2002; Smith and Eyre-Walker 2002; Sawyer et al. 2003).

Andolfatto (2005), expanding upon previous findings by Jenkins et al. (1995) and Kohn et al. (2004), examined patterns of molecular evolution in several classes of noncoding DNA (long introns, untranslated transcribed regions [UTRs], and intergenic regions), using within-species polymorphism data in *Drosophila melanogaster* and between-species divergence to the closely related *Drosophila simulans*. For all classes of noncoding sequence (compared with synonymous sites), Andolfatto (2005) found reduced levels of polymorphism and divergence, high selective constraint (ca. 40–70%), and a skew in the frequency spectrum of mutations toward rare variants, all of which indicate the action of negative selection. However, he also found a significant excess of between-species divergence relative to polymorphism (again compared with synonymous sites) for almost all classes of noncoding sequence, which is a signature of adaptive evolution. Using an extension of the McDonald–Kreitman approach (McDonald and Kreitman 1991; Fay et al. 2001), Andolfatto (2005) estimated that a substantial fraction of the divergence in noncoding regions between *D. melanogaster* and *D. simulans* was driven to fixation by positive selection (ca. 20% for intronic and intergenic sequences and 60% for UTRs). These results indicate that noncoding regions of the *D. melanogaster* genome are functionally significant and have been subject to the action of both positive and negative selection.

However, *D. melanogaster* may not be the most appropriate species for studying these types of patterns of molecular evolution because it may be unusual compared with other *Drosophila* species. Based on measures of nucleotide diversity at synonymous sites, *D. melanogaster* is thought to have experienced a reduction in effective population size (N_e) compared with its closest relative *D. simulans*

Key words: *Drosophila simulans*, noncoding DNA, natural selection, adaptive evolution, McDonald–Kreitman test, codon usage bias.

E-mail: p.haddrill@ed.ac.uk.

Mol. Biol. Evol. 25(9):1825–1834, 2008

doi:10.1093/molbev/msn125

Advance Access publication May 29, 2008

(Aquadro et al. 1988; Akashi 1995, 1996; Moriyama and Powell 1996; Andolfatto 2001; Eyre-Walker et al. 2002). This difference in N_e predicts a reduction in the efficacy of selection in *D. melanogaster* compared with *D. simulans* and thus lower levels of adaptive evolution and higher rates of fixation of mildly deleterious mutations in the *D. melanogaster* lineage (Hill and Robertson 1966; Kimura 1983).

There is some evidence for a reduction in the efficacy of selection and thus a lower N_e in the *D. melanogaster* lineage compared with the *D. simulans* lineage. In particular, levels of codon usage bias are reduced, and proteins are longer in *D. melanogaster* relative to *D. simulans* (Akashi 1995, 1996). In addition, levels of amino acid (relative to synonymous) polymorphism are higher in *D. melanogaster*, consistent with a smaller N_e (Choudhary and Singh 1987; Aquadro et al. 1988; Moriyama and Powell 1996; Andolfatto 2001).

N_e is likely to be particularly important when selection is weak, such as selection for codon usage (Akashi 1995; Akashi and Schaeffer 1997; McVean and Vieira 2001) because a mutation with mild fitness effects in a large population could be effectively neutral in a smaller population. Given that selection on noncoding sites may be fairly weak compared with nonsynonymous sites (Haddrill et al. 2007), the difference in effective population size in *D. simulans* may have a substantial impact on patterns of evolution at noncoding sites. Of particular concern is that a lower population size and thus a lower efficacy of selection, in *D. melanogaster* relative to *D. simulans*, may have led to inaccurate estimates of adaptive evolution parameters in previous studies. With polymorphism data from *D. simulans*, we can ask whether previous inferences are robust to the choice of species in which polymorphism was surveyed. Further, with information from a suitable outgroup species (*Drosophila yakuba*), we can investigate lineage-specific changes and distinguish between divergence that has accumulated due to a relaxation in selection from that accumulated due to adaptive evolution.

A recent study has examined genome-wide polymorphism and divergence data for *D. simulans*, finding evidence for the effects of both negative and positive selection (Begun et al. 2007). The scale of this analysis is impressive. However, the level of coverage in this study was, on average, only 3.9 individuals per locus. This small sample size prohibits analyses based on the frequencies of polymorphisms, which can be extremely useful in identifying signatures of negative and positive selection (Akashi 1999; Nielsen 2005). In addition, *D. simulans* populations are structured (Hamblin and Veuille 1999), and some are likely to have experienced recent bottleneck events associated with the expansion of the species out of Africa (Hamblin and Veuille 1999; Andolfatto 2001; Wall et al. 2002; Baudry et al. 2006). Thus, mixed population samples, or specifically those focusing on non-African populations, may not be ideal for making inferences about selection, given that patterns of variation may be dominated by demographic factors.

Here, we analyze polymorphism and divergence data for 67 loci (~33 kb per individual in total). We surveyed polymorphism in a sample of 20 lines of *D. simulans* from a Madagascan population (the proposed geographic origin of the species; Dean and Ballard 2004), yielding considerable information about polymorphism frequencies. The loci are

largely a subset of the loci examined by Andolfatto (2005) in *D. melanogaster* and include coding DNA and both intronic and UTR noncoding DNA. We use these data to look for signatures of both negative and positive selection in noncoding DNA in order to examine the significance of these processes in shaping patterns of molecular evolution in *D. simulans* and also to assess whether previous results for *D. melanogaster* are typical of the *melanogaster* subgroup.

Materials and Methods

Data Collection

We collected data for 21 coding regions (average surveyed length 675 bp), 22 UTRs (average surveyed length 382 bp), and 24 introns (average surveyed length 449 bp). A subset of these were surveyed by Andolfatto (2005) in *D. melanogaster* (19 coding, 20 UTRs, and 9 introns). All loci were X-linked genomic fragments and were surveyed in a sample of 20 *D. simulans* individuals from a Madagascan population (Dean and Ballard 2004) and reside in regions of high recombination in *D. simulans* (Wall et al. 2002). Further information about all 67 surveyed loci can be found in supplementary tables 1 and 2 (Supplementary Material online).

Briefly, sequence data were collected as follows. A single male fly was selected from each line and genomic DNA extracted using the Puregene DNA extraction kit. Polymerase chain reaction was used to amplify the appropriate genomic DNA fragment and then primers and unincorporated nucleotides were removed using exonuclease I and shrimp alkaline phosphatase. Fragments were directly sequenced on both strands using the Big Dye version 3.0 cycle sequencing kit (Applied Biosystems, Foster City, CA) and run on an ABI 3730 capillary sequencer. Sequence trace files were edited using Sequencher (Gene Codes Corporation, Ann Arbor, MI). The orthologous regions from *D. melanogaster* and *D. yakuba* were added to each alignment, using sequences downloaded from FlyBase (<http://flybase.org/>, Release 4.2) and the *D. yakuba* genome project (<http://insects.eugenics.org/species/blast/>), respectively, and aligned using MUSCLE (<http://www.drive5.com/muscle>) with adjustments to preserve reading frames. In some cases, regions that were particularly difficult to align were masked. This disproportionately affected introns and is expected to bias estimates of divergence downward. Details of these regions are given in supplementary table 3 (Supplementary Material online). The sequence data from this study have been submitted to GenBank under accession numbers EU744978–EU746317.

Analysis

The estimated number of synonymous sites, nonsynonymous sites, average pairwise diversity (π), average pairwise divergence to *D. melanogaster* (D_{xy}), as well as counts of the number of polymorphisms (S) were performed with a library of Perl scripts (Polymorphorama) written by P.A. and D.B. The number of nonsynonymous and synonymous sites was estimated using the Nei and Gojobori (1986) method. Average pairwise

divergence (D_{xy}) estimates were either corrected for multiple hits using a Jukes–Cantor correction (Jukes and Cantor 1969) or, in the case of synonymous sites, the Kimura (1980) 2-parameter model. Multiply hit sites were included in all analyses, but insertion–deletion polymorphisms and polymorphic sites overlapping alignment gaps were excluded.

For lineage-specific estimates of divergence, we reconstructed an ancestor of *D. melanogaster*–*D. simulans* (ANC) sequence, using *D. yakuba* as an outgroup, by maximum likelihood as implemented in the “codeml” (for coding regions, free ratio model [model = 1]) and “baseml” (for noncoding regions) programs of PAML (Yang 1997). We assigned codon usage states (with the most likely states given probabilities of 1), unpreferred (U) or preferred (P), to each codon according to the codon preference table of Andolfatto (2007), which is based on a genome-wide analysis of codon usage in *D. melanogaster*.

For analyses comparing polymorphism and divergence, we pooled site classes across loci. Each of the putatively selected noncoding site classes was compared with a subset of synonymous sites (see below) as a putatively neutral standard. Using the McDonald–Kreitman approach (McDonald and Kreitman 1991), we compared the ratio of putatively neutral with putatively selected polymorphic and divergent sites using a Fisher’s exact test. Following Andolfatto (2005), we used an extension of this approach (Fay et al. 2001) to estimate the proportion of divergence driven to fixation by positive selection (α) as $\alpha = 1 - (D_S P_X / D_X P_S)$, where the subscripts *S* and *X* denote the putatively neutral and putatively selected sequence classes, respectively, and $D = \sum_{i=1}^n D_i$ and $P = \sum_{i=1}^n P_i$, where D_i and P_i are the number of divergent and polymorphic variants at locus *i*, respectively, and *n* is the number of loci in a particular sequence class. The number of divergent sites at a locus (*D*) was corrected for multiple hits using a Jukes–Cantor correction (Jukes and Cantor 1969). Confidence intervals (CIs) for α were estimated using a nonparametric bootstrap procedure, with resampling by site. This method has been shown to only slightly underestimate CIs, given estimates of recombination and nucleotide diversity in *Drosophila* (Andolfatto 2005a).

For comparison, we also carried out these analyses using nonsynonymous sites as the putatively selected site class. In addition, for nonsynonymous sites, it was possible to compare estimates of α calculated using the summing across loci method above to those calculated using the methods of Smith and Eyre-Walker (2002), Bierne and Eyre-Walker (2004), and Welch (2006) to ensure that estimates from summing across loci were not markedly different from other methods (see supplementary results, Supplementary Material online).

Results

Reduced Polymorphism and Divergence in Noncoding Regions

We surveyed a total of ~14 kb of coding sequence and ~19 kb of noncoding sequence per individual. Polymorphism and divergence summaries for each sequence class

are shown in table 1, using both *D. melanogaster* and ANC as an outgroup. Data for individual loci are presented in supplementary table 4 (Supplementary Material online).

This data set represents a subset of the loci examined by Andolfatto (2005) in *D. melanogaster* and, consistent with that study, we found levels of divergence are similarly reduced in all noncoding DNA classes relative to synonymous sites. Like *D. melanogaster*, nonsynonymous and noncoding DNA in *D. simulans* show reduced levels of polymorphism compared with synonymous sites. Mean synonymous site diversity is 3.02%, significantly higher than at nonsynonymous sites (0.19%, Wilcoxon 2-sample test, $P < 10^{-4}$) and noncoding sites (1.13%, $P < 10^{-4}$). In fact, average levels of polymorphism for all classes of sites in *D. simulans* are surprisingly similar to those in *D. melanogaster* (see table 1 of Andolfatto [2005]). This is in contrast to some previous studies, which have reported estimates of nucleotide diversity at synonymous sites as being significantly lower in *D. melanogaster* than *D. simulans*, consistent with a reduced effective population size (N_e) in *D. melanogaster* relative to *D. simulans* (Aquadro et al. 1988; Akashi 1995, 1996; Eyre-Walker et al. 2002). However, it should be noted that there is a large X-autosome component to the reported difference in levels of diversity between the 2 species (Andolfatto 2001). Because loci surveyed in this study are all X linked, similar levels of variability are expected based on previous findings (Andolfatto 2001; Nolte and Schlötterer 2008).

A Genome-Wide Negative Skew in the Distribution of Polymorphism Frequencies

Reduced levels of polymorphism and divergence indicate that nonsynonymous and noncoding variants segregating in the population may be subject to selective constraints, resulting in slower rates of evolution (as is typically assumed for nonsynonymous sites). If this is the case, the frequency of polymorphisms in these sequence classes will be skewed toward rare variants because the action of negative selection will keep such variants at lower frequencies than those that are neutral (Tajima 1989; Akashi 1999; Nielsen and Weinreich 1999; Andolfatto 2005; Bachtrog and Andolfatto 2006). Alternatively, reduced polymorphism and divergence at noncoding sites relative to synonymous sites could reflect lower mutation rates, but this would not be expected to result in a skew in the polymorphism frequency spectrum.

Figure 1 shows mean Tajima’s *D* values for all sequence classes. This analysis shows that the distribution of polymorphism frequencies is negatively skewed for both nonsynonymous sites and all noncoding sequence classes, as indicated by negative mean Tajima’s *D* values, consistent with the hypothesis that these polymorphisms are subject to purifying selection.

However, in contrast to the patterns in *D. melanogaster* (Andolfatto 2005), mean Tajima’s *D* for synonymous sites is also strongly negative, and the distributions of Tajima’s *D* values are not significantly different between any of the sequence classes examined (Wilcoxon 2-sample tests, all $P > 0.16$). The strong skew toward rare polymorphisms

Table 1
Polymorphism and Divergence in Coding and Noncoding DNA of *Drosophila simulans*

Sequence Class	No. Regions	Mean π^a	Mean D_{xy}^b	D^c	$P^{all(d)}$	$P^{all(e)}$	$P^{l(d,f)}$	$P^{l(e,f)}$
<i>melanogaster</i> outgroup								
Synonymous	21	3.02	13.87	59 ^g	110 ^g	—	65 ^g	—
Nonsynonymous	21	0.19	1.24	123	158	0.025	50	$<10^{-4}$
Noncoding	46	1.13	5.44	809	1,212	0.078	463	$<10^{-4}$
Introns	24	1.29	6.13	505	804	0.188	312	0.001
UTRs	22	0.95	4.69	304	408	0.021	151	$<10^{-4}$
5' UTRs	10	0.89	4.95	166	189	0.003	78	$<10^{-4}$
3' UTRs	12	0.99	4.47	138	219	0.218	73	$<10^{-4}$
ANC outgroup ^h								
Synonymous	21	3.02	6.10	29 ^g	110 ^g	—	63 ^g	—
Nonsynonymous	21	0.18	0.68	57	142	0.128	43	$<10^{-4}$
Noncoding	46	1.01	2.31	233	983	0.650	357	0.167
Introns	24	1.22	2.33	115	678	0.074	247	1.000
UTRs	22	0.79	2.29	118	305	0.119	110	0.001
5' UTRs	10	0.87	2.45	67	180	0.180	71	0.014
3' UTRs	12	0.71	2.14	51	125	0.118	39	0.001

^a π is the average pairwise divergence per nucleotide site between alleles (%).
^b D_{xy} is the average Jukes–Cantor corrected divergence from the outgroup (%).
^c D is the number of divergent sites (Jukes–Cantor corrected).
^d P^{all}/P^l are the number of polymorphic sites including all polymorphisms/excluding those at a frequency of less than 5%.
^e P^{all}/P^l are probabilities from McDonald–Kreitman tests including all polymorphisms/excluding those at a frequency of less than 5%.
^f Excluding polymorphisms present at a frequency of less than 10% and 15% did not alter the conclusions.
^g The synonymous site class counts for the McDonald–Kreitman tests include only P → P and U → U changes (see text for details).
^h ANC outgroup is the reconstructed ancestor of *D. simulans* and *Drosophila melanogaster* used as the outgroup.

at synonymous sites may suggest that, despite having high levels of polymorphism and divergence compared with other sequence classes, these sites may be subject to purifying selection. It has previously been suggested that selection on synonymous sites is stronger in *D. simulans* than *D. melanogaster* (Akashi 1995; Akashi and Schaeffer 1997; Nielsen et al. 2007). This is also consistent with the observation of reduced divergence at synonymous sites in the *D. simulans* lineage compared with the *D. melanogaster* lineage (Akashi 1996; Begun et al. 2007).

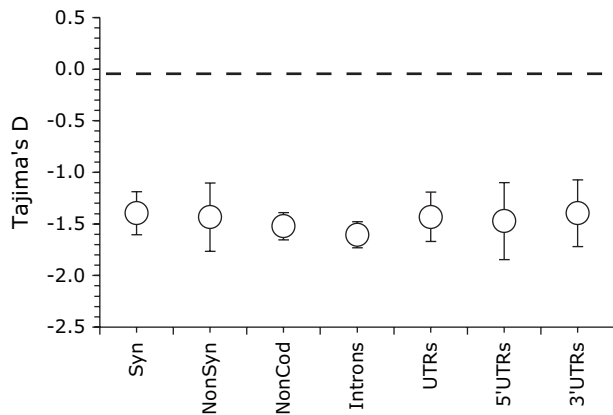


FIG. 1.—Mean Tajima's D values for coding and noncoding DNA in *Drosophila simulans*, using the *Drosophila melanogaster* outgroup. The dashed line indicates the mean expected value of Tajima's D under a model of neutral evolution for a sample of 20 individuals with similar levels of variability to those found here (no recombination, see table 1 of Tajima [1989]). Error bars indicate 2 standard errors. Syn: synonymous sites, NonSyn: nonsynonymous sites, and NonCod: pooled noncoding sites.

Evidence for Purifying Selection on Synonymous, Nonsynonymous, and Noncoding DNA Sites in *D. simulans*

To investigate the evidence for selection on synonymous sites in more detail, we divided all synonymous changes into 3 categories, based on predictions under a model of selection for codon usage bias, using preferred (P)/unpreferred (U) codon classifications from *D. melanogaster* (Andolfatto 2007). The standard model of codon usage bias (Bulmer 1991; Akashi 1995) assumes that on average: 1) preferred to preferred (P → P) and unpreferred to unpreferred (U → U) changes are neutral, 2) preferred to unpreferred (P → U) changes are deleterious, and 3) unpreferred to preferred (U → P) changes are advantageous. We examined the frequency distributions of polymorphisms in each of these categories and each of the noncoding DNA classes.

Figure 2a shows that selection to maintain optimal codon usage is likely to account for a considerable fraction of the overall skew toward low-frequency variants in the frequency spectrum of synonymous polymorphisms. Given the codon selection model envisioned, we propose that the P → P and U → U classes (accounting for ~22% of our polymorphisms) should be used as a neutral frame of reference. P → U changes (accounting for ~53% of our polymorphisms) are putatively negatively selected and, consistent with this expectation, this category shows the strongest skew toward rare variants and is significantly more negatively skewed than the P → P/U → U class (Wilcoxon 2-sample test, $P < 10^{-3}$). In contrast, a greater proportion of U → P changes, the putatively positively selected class and ~25% of our polymorphisms, are seen at intermediate or high frequency. The frequency spectrum for U → P changes is significantly different from that for P → U changes

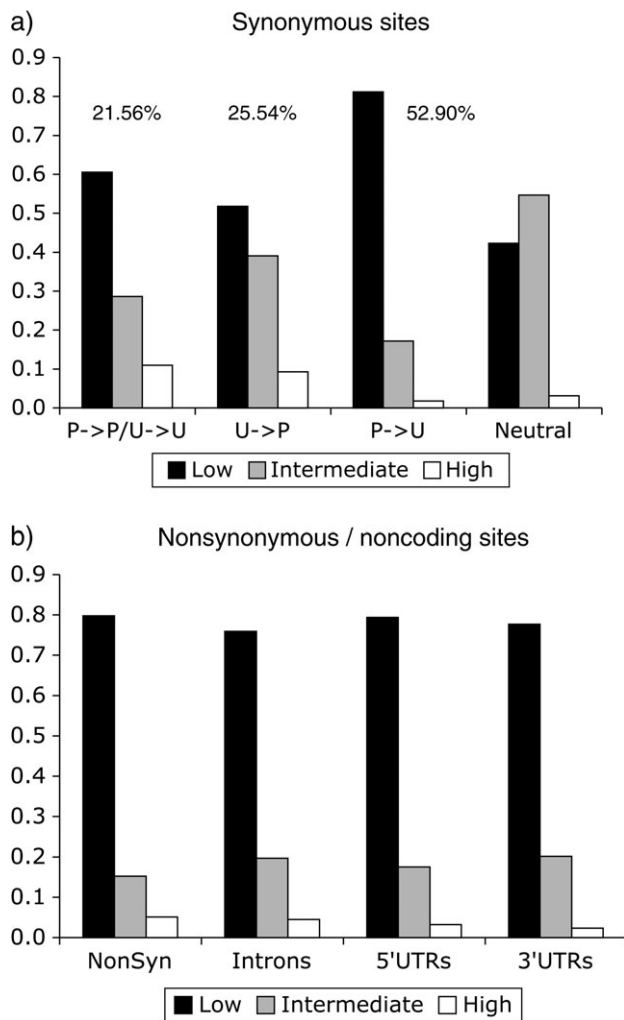


FIG. 2.—The distribution of frequency classes of polymorphisms for different types of synonymous site changes, nonsynonymous site changes, and noncoding DNA changes. Within synonymous site changes, $P \rightarrow P/U \rightarrow U$ includes preferred to preferred and unpreferred to unpreferred changes, $P \rightarrow U$ are preferred to unpreferred changes, and $U \rightarrow P$ are unpreferred to preferred changes. “Neutral” indicates the neutral expectation. The low-frequency class includes polymorphisms at a frequency of 1–2/20, the intermediate frequency class 3–17/20, and the high-frequency class 18–19/20. The numbers above each type of synonymous site change indicate the percentage of the total synonymous polymorphisms of each type.

($P < 10^{-3}$) but not significantly different from the putatively neutral class (i.e., vs. $P \rightarrow P/U \rightarrow U$, $P > 0.05$).

Although the shift toward higher frequencies for $U \rightarrow P$ changes is suggestive of positive selection, there is no statistical support for this claim. It is apparent in figure 2a that the frequency spectrum for $U \rightarrow P$ changes is still negatively skewed relative to the neutral model (as is that for $P \rightarrow P/U \rightarrow U$), suggesting either that these polymorphisms are also deleterious (as has been suggested for within-class variants, Begun 2001), on average, or that other population genetic factors are also influencing synonymous polymorphism frequencies, such as genetic hitchhiking (Braverman et al. 1995; Andolfatto 2007), background selection (Gordo et al. 2002), or population expansion (Dean and Ballard 2004). Given the significantly

negative skew in the frequency spectrum of $P \rightarrow U$ changes and possible positive selection on $U \rightarrow P$ polymorphisms, we propose to use $P \rightarrow P$ and $U \rightarrow U$ changes in comparisons to nonsynonymous and noncoding sites below.

Using our proposed neutral class of synonymous sites ($P \rightarrow P$ and $U \rightarrow U$) as a reference, we detect a significant shift toward lower polymorphism frequencies at nonsynonymous sites and in every class of noncoding DNA (fig. 2b; Wilcoxon 2-sample test, all $P < 10^{-3}$). This pattern is consistent with the findings of Andolfatto (2005) and Bachtrog and Andolfatto (2006) and suggests stronger selective constraint on nonsynonymous and noncoding DNA sites than on $P \rightarrow P$ and $U \rightarrow U$ synonymous changes. We failed to detect a difference in the frequency spectra of $P \rightarrow U$ changes, nonsynonymous sites, and all classes of noncoding DNA (all $P > 0.05$).

Quantifying Positive Selection

In a study of 188 X-linked loci in *D. melanogaster*, Andolfatto (2005) found that as well as being subject to the action of purifying selection, a large fraction of the nonsynonymous and noncoding divergence between *D. melanogaster* and *D. simulans* had been driven to fixation by positive selection. Several studies have also identified substantial levels of adaptive evolution at nonsynonymous sites using polymorphism data from *D. melanogaster* and *D. simulans* (reviewed in Eyre-Walker 2006), but noncoding DNA has not been as widely investigated. We can look for the effects of positive selection in our data using the McDonald–Kreitman approach (McDonald and Kreitman 1991), which distinguishes negative and positive selection from neutrality by comparing levels of polymorphism and divergence for a putatively selected class of sequence with that of a neutral standard. If both types of sequence are evolving neutrally, the ratio of polymorphism and divergence will be the same for the 2 classes. In contrast, positive selection will increase the amount of divergence relative to polymorphism, whereas negative selection will reduce it.

In the past, this approach has been used using all synonymous sites as a neutral reference. However, theoretical results show that selection on synonymous sites can inflate levels of polymorphism relative to divergence relative to neutral expectations, biasing tests based on the McDonald–Kreitman approach toward showing evidence for positive selection (Akashi 1995, 1999; McVean and Charlesworth 1999; Eyre-Walker 2002; Andolfatto 2005). Andolfatto (2005) showed that the McDonald–Kreitman approach will not lead to spurious evidence for positive selection unless the neutral reference class used is under stronger selective constraint than the putatively selected class.

In our analysis of the frequency spectrum of polymorphisms (above), we found that all classes of synonymous substitution exhibit a skew toward rare variants relative to neutral expectations, including $P \rightarrow P/U \rightarrow U$ polymorphisms, consistent with purifying selection (Begun 2001), among other possible factors. An alternative, highly conservative, approach is to use $U \rightarrow P$ synonymous changes as the reference class (Akashi 1995, 1999; Begun et al. 2007). However,

we find that ratios of polymorphism to divergence are not significantly different for the $P \rightarrow P/U \rightarrow U$ and $U \rightarrow P$ classes of synonymous substitution (supplementary table 5, Supplementary Material online; Fisher's exact test, $P > 0.05$), and using $U \rightarrow P$ synonymous changes does not significantly alter estimates of α reported below for the *D. simulans* lineage (see Discussion). We thus conclude that our lineage-specific estimates of α in *D. simulans* are not likely to be noticeably upwardly biased by weak purifying selection on $P \rightarrow P/U \rightarrow U$ polymorphisms. Further, the use of $U \rightarrow P$ polymorphisms is limited to lineage-specific comparisons and comes with a loss of statistical power (because only a small fraction of synonymous substitutions fall into this class). We thus propose to use $P \rightarrow P/U \rightarrow U$ polymorphisms as our neutral reference class.

The inclusion of all observed polymorphisms in the putatively selected class can make the McDonald–Kreitman test conservative because the presence of mildly deleterious variants that are kept at low frequency by selection can mask the signature of positive selection (Fay et al. 2001). Because these mutations are subject to weak negative selection, they will contribute to polymorphism but not to divergence and thus will obscure an excess of divergence due to positive selection. It is therefore possible to increase the power of this test by excluding variants that are present at low frequency, as long as the putatively neutral and selected sites are treated equally (Templeton 1996). There is no standard cutoff frequency for removing rare variants, thus we examined the effect of removing variants up to a frequency of 5% (singletons only), 10% (singletons and doubletons), and 15% (singletons, doubletons, and tripletons) on the results of the McDonald–Kreitman test (see table 1).

For all cutoff frequency criteria, we observe a significant excess of divergence relative to synonymous sites ($P \rightarrow P$ and $U \rightarrow U$ sites only) for nonsynonymous sites and all classes of noncoding DNA (Fisher's exact test, nonsynonymous sites, $P < 10^{-4}$; introns, $P \leq 0.001$; all UTR sites, $P < 10^{-4}$; 5' UTR sites, $P < 10^{-4}$; and 3' UTR sites, $P \leq 0.002$; see table 1). These results are consistent with the hypothesis that a significant fraction of divergence at nonsynonymous and all noncoding sites has been driven to fixation by positive selection and is in agreement with Andolfatto (2005), who also found an excess of divergence at these classes of sequence using polymorphism data from *D. melanogaster*.

The McDonald–Kreitman test assumes constancy in levels of selective constraint, and changes in population size can violate this assumption (see Discussion). Given the presumed difference in the efficacy of selection in the *D. melanogaster* and *D. simulans* lineages, we repeated the above analyses using ANC as the outgroup and looking specifically at divergence along the *D. simulans* lineage. Using all polymorphisms, and at all cutoff frequencies, we find an excess of lineage-specific divergence relative to the neutral class for nonsynonymous sites (Fisher's exact test, $P \leq 0.002$), all UTR sites ($P = 0.001$), 5' UTR sites ($P \leq 0.014$), and 3' UTR sites ($P \leq 0.004$). Despite evidence for a marginally significant excess of polymorphism at intronic sites (suggesting the action of purifying selection, $P = 0.074$), we observe no signature of positive selection in introns along the *D. simulans* lineage.

Because positive selection appears to have been an important force in the evolution of several classes of noncoding DNA, we also attempt to quantify the fraction of the divergence at these sites that has been driven to fixation by positive selection (α), using an extension of the McDonald–Kreitman test. Using the approach of Fay et al. (2001), we calculate α by summing the number of divergent and polymorphic variants for each locus within the putatively neutral and selected classes (see Materials and Methods). Again we use only the $P \rightarrow P$ and $U \rightarrow U$ synonymous changes as the neutral reference class and compare these to the nonsynonymous and noncoding (putatively selected) sequence classes.

Figure 3 shows estimates of α for nonsynonymous sites and all types of noncoding sequence, using both *D. melanogaster* and ANC as outgroups. Again we examine the effect of excluding rare variants of different frequencies because the estimation of α will also be affected by the presence of weakly deleterious mutations, as discussed above. Figure 3 shows that excluding singleton polymorphisms generally results in a large increase in α , but as the cutoff frequency increases from singletons to tripletons (from 5% to 15%), the corresponding increase in α is small and tends to level off.

When divergence is measured to *D. melanogaster* (fig. 3a), we estimate the fraction of divergence driven to fixation by positive selection to be ~45% in introns, ~50–60% for 3' UTRs, and ~60–70% at nonsynonymous sites and 5' UTRs. Using polymorphism data from *D. melanogaster*, Andolfatto (2005) reported similar estimates of α for nonsynonymous and UTR sites (~60%) but lower estimates for intronic sites (~20%). When we consider α along the *D. simulans* lineage only (fig. 3b), we find very similar estimates of the fraction of divergence driven to fixation by positive selection for nonsynonymous sites and 5' UTRs (~60–70%) and slightly higher estimates for 3' UTRs (~65%), compared with estimates using the *D. melanogaster* outgroup. However, a clear signature of positive selection at intron sites is no longer evident with values of α actually being significantly negative when all polymorphisms are included. Although removing low-frequency polymorphisms results in positive estimates of α in introns, these are not significantly greater than zero.

However, it should be noted that these estimates of α have large CIs, so are actually compatible with a fairly wide range of values, including zero. It is possible that there is greater variation in α values for introns than UTRs on the *D. simulans* lineage, making the signal of adaptive evolution at these sites less clear. It is of note that if we use only the 9 intron loci surveyed both here and by Andolfatto (2005) in *D. melanogaster* (those surveyed by Haddrill, Thornton, et al. [2005]), the estimate of α in introns along the *D. simulans* lineage is ~25% (excluding singletons). Although this value is still not significantly greater than zero (90% CIs = -0.24 to 0.55), it is not inconsistent with Andolfatto's (2005) estimate of α in introns in *D. melanogaster* (~20%). However, the estimate of α in this subset of 9 introns using divergence to *D. melanogaster*, although overlapping with the *D. simulans* lineage estimate, is still higher (~50% excluding singletons) and is significantly greater than zero (90% CIs = 0.31–0.66).

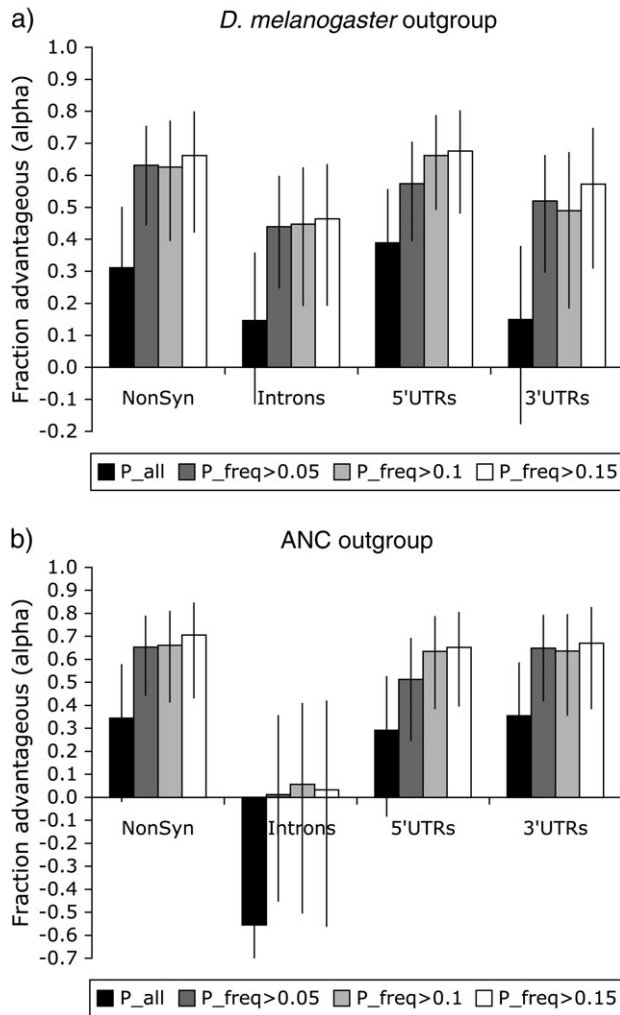


FIG. 3.—Estimates of the fraction of adaptively driven nucleotide substitutions (α) in coding and noncoding DNA (a) between *Drosophila simulans* and *Drosophila melanogaster* and (b) inferred along the *D. simulans* lineage using a reconstructed ancestor. NonSyn: nonsynonymous sites. For each class of sequence, 4 estimates are presented based on the subset of polymorphisms that were included in the analysis: P_all, all polymorphisms included; P_freq > 0.05, polymorphisms at a frequency of greater than 5% (singletons) included; P_freq > 0.1, polymorphisms at a frequency of greater than 10% (doubletons) included; and P_freq > 0.15, polymorphisms at a frequency of greater than 15% (tripletons) included. Error bars indicate 90% CIs.

Discussion

Overall, these results support recent claims that a large fraction of noncoding DNA in *Drosophila* is functionally important (Bergman and Kreitman 2001; Andolfatto 2005; Haddrill, Charlesworth, et al. 2005; Bachtrog and Andolfatto 2006; Halligan and Keightley 2006). In agreement with previous studies in *D. melanogaster* (Kohn et al. 2004; Andolfatto 2005), UTR sequences appear to have been subject to both negative and positive selection in *D. simulans*, suggesting that they contain a great deal of functionally important sequence. Whether along the lineage separating *D. simulans* and *D. melanogaster* or along the *D. simulans* lineage only, we conclude that a large fraction of nucleotide divergence in both 5' UTRs and 3' UTRs has

been driven to fixation by positive selection. In fact, the UTR sequences tend to show similar fractions of adaptively driven divergence as nonsynonymous sites, and these values for α are also similar to those estimated for nonsynonymous sites in other comparisons (reviewed in Eyre-Walker 2006). Evolution at UTR sites is therefore likely to have had a substantial impact on adaptive divergence between *D. simulans* and *D. melanogaster*, and the influence of both negative and positive selection suggests that a large fraction of UTR sequence is functionally important.

Similar to UTRs, intronic sites show clear signatures of negative selection in *D. simulans*. However, in contrast to UTRs, the signature of positive selection is clear only when divergence between *D. simulans* and *D. melanogaster* is considered and not when considering the *D. simulans* lineage only. There are several possible explanations for this discrepancy. First, intron sequences may not have evolved adaptively along the *D. simulans* lineage, and our results could indicate that the positively selected divergence between these 2 species at intronic sites may have occurred primarily along the *D. melanogaster* lineage. Second, the pattern could reflect the accumulation of deleterious fixations in introns along the *D. melanogaster* lineage, which would be expected if *D. melanogaster* has indeed experienced a period of population size reduction (Fay and Wu 2001), and deleterious mutations at introns are subject to weaker purifying selection, on average, than at nonsynonymous or UTR sites.

Finally, biases in our analyses could be contributing to this result. For example, biases in the reconstruction of the ancestor of *D. simulans* and *D. melanogaster* could cause divergence (and thus α) to be underestimated, although this seems unlikely because we do not see the same effect in UTRs. Another potential bias, as mentioned previously (see Materials and Methods), results from the fact that some difficult-to-align regions were masked from alignments, and this disproportionately affected introns, thus potentially underestimating their divergence. We repeated the analyses with the masked regions included in the alignments and found that estimates of α in introns increased only slightly (e.g., from 5.6% to 7.1%, excluding singletons). This suggests that the exclusion of these regions has not resulted in a major bias in our results. However, there may still be a bias toward underestimating divergence in introns. Although some of the originally masked regions are likely to reflect insertion/deletion divergence between species, some may have been incorrectly aligned as insertions/deletions. Because regions of insertion/deletion are excluded from the analyses, this would bias estimates of divergence downward and could therefore be contributing to the lack of evidence for adaptive evolution at intronic sites.

To investigate this further, we calculated lineage-specific estimates of α along the *D. melanogaster* lineage, using data from Andolfatto (2005), for all loci that overlap with this study (19 coding regions, 9 introns, ten 5' UTRs, and ten 3' UTRs). Table 2 also shows lineage-specific estimates of the total number of divergent sites for each sequence class. For nonsynonymous and UTR sites, the sum of the 2 lineage-specific estimates of divergence is close to the estimates of divergence along the entire *D. simulans*–*D. melanogaster* lineage. However, this is not the case for

Table 2

Estimates of the Number of Divergent Sites and the Fraction of Divergence Driven to Fixation by Positive Selection, α (with 90% CIs) for Coding and Noncoding DNA in *Drosophila simulans* and *Drosophila melanogaster*, Using Loci Overlapping between This Study and Andolfatto (2005)

Lineage	NonSyn ^a	Introns	5' UTRs	3' UTRs
Divergence				
sim-mel ^b	123	187	161	138
sim ^c	57	53	67	51
mel-sim ^d	132	227	138	151
mel ^e	62	74	65	81
α^f				
sim-mel ^b	0.63 (0.39–0.77)	0.51 (0.25–0.67)	0.66 (0.49–0.79)	0.49 (0.18–0.67)
sim ^c	0.66 (0.41–0.81)	0.22 (–0.27 to 0.53)	0.63 (0.38–0.79)	0.64 (0.35–0.80)
mel-sim ^d	0.68 (0.54–0.81)	0.04 (–0.23 to 0.24)	0.43 (0.22–0.60)	0.75 (0.64–0.85)
mel ^e	0.71 (0.56–0.82)	–0.29 (–0.76 to 0.06)	0.41 (0.14–0.60)	0.79 (0.67–0.88)

^a Nonsynonymous sites.

^b Estimates using *D. simulans* polymorphism and divergence to *D. melanogaster*.

^c *Drosophila simulans* lineage-specific estimates.

^d Estimates using *D. melanogaster* polymorphism and divergence to *D. simulans*.

^e *Drosophila melanogaster* lineage-specific estimates.

^f Estimates of α were calculated using only P \rightarrow P/U \rightarrow U synonymous sites in *D. simulans* but using all synonymous sites in *D. melanogaster*. Polymorphisms at a frequency of up to 10% were excluded in all cases.

introns, where the sum of the 2 lineage-specific estimates of divergence is substantially lower than the estimates along the entire divergence time between species. This suggests that the addition of a more distant species (*D. yakuba*) to the alignments, for use in reconstructing the ANC outgroup, has resulted in the exclusion of large unalignable regions specifically in introns, thus considerably underestimating lineage-specific divergences. The suggestion that introns have not evolved adaptively along the *D. simulans* lineage may therefore be unreliable. However, this does not seem to be the case for coding regions and UTRs, so we can be confident in our lineage-specific estimates of α in those classes.

One other potential bias in our data could result if *D. simulans* has experienced a population size expansion in the recent past, as has been suggested previously for the Madagascar population (Dean and Ballard 2004). Such a demographic event can lead to artifactual evidence of adaptive evolution in the putatively selected class because mildly deleterious mutations that are currently removed from the population would have become fixed in the past, resulting in higher levels of divergence than would be expected given current patterns of polymorphism (Ohta 1993; Fay and Wu 2001; Eyre-Walker 2002). The negative skew in the frequency spectrum of P \rightarrow P/U \rightarrow U polymorphisms compared with the neutral expectation may indicate a population size expansion in *D. simulans*, and this could be affecting our estimates of α at nonsynonymous and noncoding sites. We can determine whether we see the predicted effects of a population size expansion by examining the divergence to polymorphism ratio for mildly deleterious P \rightarrow U synonymous sites. Under a population size expansion, these sites would be expected to have accumulated deleterious fixations in the past, when the population size was smaller, and will thus have higher levels of divergence than would be expected, given current patterns of polymorphism. The ratio for P \rightarrow U sites is 0.17 compared with 0.26 for our neutral sites, the P \rightarrow P and U \rightarrow U synonymous sites (see supplementary table 5, Supplementary Material online). This difference is in the opposite direction

to that predicted under a weak selection—population size expansion model.

In addition to this, as mentioned above, the use of U \rightarrow P synonymous changes as the neutral reference class does not significantly alter estimates of α along the *D. simulans* lineage (estimates of α [with 90% CIs]: nonsynonymous sites = 0.74 [0.58–0.86]; introns = 0.28 [–0.05 to 0.54]; 5' UTRs = 0.72 [0.56–0.84]; and 3' UTRs = 0.72 [0.54–0.85]). The use of these changes can be viewed as being robust to a population size expansion because the effect of such a demographic event would be to enhance adaptive evolution at U \rightarrow P sites. This makes the use of these sites as the neutral reference class even more conservative than in the absence of an expansion. These results suggest that, although we cannot rule out the possibility that a population size expansion along the *D. simulans* lineage has had an influence on our results, it cannot alone explain patterns of adaptive evolution in our putatively selected classes. Similarly, although a population size expansion may be contributing to the overall negative skew in polymorphism frequencies (although other possible explanations are listed above), the fact that P \rightarrow U synonymous, nonsynonymous, and noncoding polymorphisms are significantly more negatively skewed than P \rightarrow P/U \rightarrow U and U \rightarrow P synonymous changes indicates that a population size expansion cannot alone account for signatures of negative selection, even if differences in mutation rate are invoked.

We also compared patterns of evolution at synonymous sites in *D. simulans* with Andolfatto's (2005) results for *D. melanogaster*. In contrast to previous studies (Akashi 1995, 1996; Eyre-Walker et al. 2002), levels of polymorphism and divergence on the X chromosome are remarkably similar in the 2 species, suggesting that they have similar effective population sizes. A difference in population size between the 2 species is suggested by substantially higher levels of autosomal diversity in *D. simulans*, as well as several consistent patterns of polymorphism and divergence (Aquadro et al. 1988; Akashi 1995, 1996; Akashi and

Schaeffer 1997; Andolfatto 2001; Eyre-Walker et al. 2002). In contrast to *D. melanogaster* (Andolfatto 2005), we found signatures of purifying selection at synonymous sites in *D. simulans* and this seems to be at least partly explained by selection for codon usage bias. This is in agreement with previous studies that have reported a relaxation of selection for optimal codon usage in *D. melanogaster* due to a reduction in effective population size along the *D. melanogaster* lineage (Akashi 1996).

In summary, we have examined patterns of evolution at several classes of noncoding DNA in *D. simulans* and find that all noncoding DNA is subject to the action of negative selection, indicated by low levels of polymorphism and divergence and a skew in the frequency spectrum toward rare variants. Although the signature of negative selection is the only consistent pattern we see in introns, evidence for an excess of divergence relative to polymorphism indicates that a large fraction of nucleotide divergence in UTRs has been driven to fixation by positive selection in the *D. simulans* lineage. These findings add to the increasing wealth of evidence that some of the most important regions of the genome are found outside those that encode proteins, emphasizing the importance of evolution in regulatory regions.

Supplementary Material

Supplementary results and tables 1–5 are available at *Molecular Biology and Evolution* online (<http://www.mbe.oxfordjournals.org/>).

Acknowledgments

We thank Hiroshi Akashi, Brian Charlesworth, Daniel Halligan, and 1 anonymous reviewer for comments on the manuscript. We thank John Welch for comments on the manuscript, for providing the “MKtest” program, and for extremely valuable discussions and advice regarding the calculation of α statistics. We thank Adam Eyre-Walker for providing the “Distribution of Fitness Effects” program and for advice on its use. P.A. was supported in part by an A. P. Sloan fellowship in Molecular and Computational Biology. D.B. was supported in part by a National Institutes of Health Grant (GM076007) and an A. P. Sloan fellowship in Molecular and Computational Biology.

Literature Cited

- 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. Detecting the ‘footprint’ of natural selection in within and between species DNA sequence data. *Gene*. 238:39–51.
- Akashi H, Schaeffer SW. 1997. Natural selection and the frequency distributions of “silent” DNA polymorphism in *Drosophila*. *Genetics*. 146:295–307.
- Andolfatto P. 2001. Contrasting patterns of X-linked and autosomal nucleotide variation in *Drosophila melanogaster* and *Drosophila simulans*. *Mol Biol Evol*. 18:279–290.
- Andolfatto P. 2005a. Adaptive evolution of non-coding DNA in *Drosophila*. *Nature*. 437:1149–1152.
- Andolfatto P. 2007. Hitchhiking effects of recurrent beneficial amino acid substitutions in the *Drosophila melanogaster* genome. *Genome Res*. 17:1755–1762.
- Andolfatto P. 2005b. Adaptive evolution of non-coding DNA in *Drosophila*. *Nature*. 437:S2.
- Aquadro CF, Lado KM, Noon WA. 1988. The *rosy* region of *Drosophila melanogaster* and *Drosophila simulans*. I. Contrasting levels of naturally occurring DNA restriction map variation and divergence. *Genetics*. 119:875–888.
- Bachtrog D, Andolfatto P. 2006. Selection, recombination and demographic history in *Drosophila miranda*. *Genetics*. 174:2045–2059.
- Baudry E, Derome N, Huet M, Veuille M. 2006. Contrasted polymorphism patterns in a large sample of populations from the evolutionary genetics model *Drosophila simulans*. *Genetics*. 173:759–767.
- Begun DJ. 2001. The frequency distribution of nucleotide variation in *Drosophila simulans*. *Mol Biol Evol*. 18:1343–1352.
- Begun DJ, Holloway AK, Stevens K, et al. (13 co-authors). 2007. Population genomics: whole-genome analysis of polymorphism and divergence in *Drosophila simulans*. *PLoS Biol*. 5: 2534–2559.
- Bergman CM, Carlson JW, Celniker SE. 2005. *Drosophila* DNase I footprint database: a systematic annotation of transcription factor binding sites in the fruitfly, *Drosophila melanogaster*. *Bioinformatics*. 21:1747–1749.
- Bergman CM, Kreitman M. 2001. Analysis of conserved non-coding DNA in *Drosophila* reveals similar constraints in intergenic and intronic sequences. *Genome Res*. 11:1335–1345.
- Bierne N, Eyre-Walker A. 2004. The genomic rate of adaptive amino acid substitution on *Drosophila*. *Mol Biol Evol*. 21:1350–1360.
- Braverman JM, Hudson RR, Kaplan NL, Langley CH, Stephan W. 1995. The hitchhiking effect on the site frequency spectrum of DNA polymorphisms. *Genetics*. 140:783–796.
- Bulmer M. 1991. The selection-mutation-drift theory of synonymous codon usage. *Genetics*. 129:897–907.
- Carlini DB, Chen Y, Stephan W. 2001. The relationship between third-codon position nucleotide content, codon bias, mRNA secondary structure and gene expression in the *Drosophilid* alcohol dehydrogenase genes *Adh* and *Adhr*. *Genetics*. 139:625–633.
- Casillas S, Barbadilla A, Bergman CM. 2007. Purifying selection maintains highly conserved noncoding sequences in *Drosophila*. *Mol Biol Evol*. 24:2222–2234.
- Chen Y, Stephan W. 2003. Compensatory evolution of a precursor messenger RNA secondary structure in the *Drosophila melanogaster* *Adh* gene. *Proc Natl Acad Sci USA*. 100:11499–11504.
- Choudhary M, Singh RS. 1987. A comprehensive study of genic variation in natural populations of *Drosophila melanogaster*. III. Variations in genetic structure and their causes between *Drosophila melanogaster* and its sibling species *Drosophila simulans*. *Genetics*. 117:697–710.
- Dean MD, Ballard JW. 2004. Linking phylogenetics with population genetics to reconstruct the geographic origin of a species. *Mol Phylogenet Evol*. 32:998–1009.
- Eyre-Walker A. 2002. Changing effective population size and the McDonald-Kreitman test. *Genetics*. 162:2017–2024.
- Eyre-Walker A. 2006. The genomic rate of adaptive evolution. *Trends Ecol Evol*. 21:569–575.
- Eyre-Walker A, Keightley PD, Smith NGC, Gaffney D. 2002. Quantifying the slightly deleterious mutation model of molecular evolution. *Mol Biol Evol*. 19:2142–2149.

- Fay JC, Wu C-I. 2001. The neutral theory in the genomic era. *Curr Opin Genet Dev*. 11:642–646.
- Fay JC, Wyckoff GJ, Wu C-I. 2001. Positive and negative selection on the human genome. *Genetics*. 158:1227–1234.
- Fay JC, Wyckoff GJ, Wu C-I. 2002. Testing the neutral theory of molecular evolution with genomic data from *Drosophila*. *Nature*. 415:1024–1026.
- Gallo SM, Li L, Hu Z, Halfon MS. 2006. REDfly: a regulatory element database for *Drosophila*. *Bioinformatics*. 22:381–383.
- Gordo I, Navarro A, Charlesworth B. 2002. Muller's ratchet and the pattern of variation at a neutral locus. *Genetics*. 161:835–848.
- Haddrill PR, Charlesworth B, Halligan DL, Andolfatto P. 2005. Patterns of intron sequence evolution in *Drosophila* are dependent upon length and GC content. *Genome Biol*. 6:R67.
- Haddrill PR, Thornton KR, Charlesworth B, Andolfatto P. 2005. Multilocus patterns of nucleotide variability and the demographic and selection history of *Drosophila melanogaster* populations. *Genome Res*. 15:790–799.
- Haddrill PR, Halligan DL, Tomaras D, Charlesworth B. 2007. Reduced efficacy of selection in regions of the *Drosophila* genome that lack crossing over. *Genome Biol*. 8:R18.
- Halligan DL, Eyre-Walker A, Andolfatto P, Keightley PD. 2004. Patterns of evolutionary constraints in intronic and intergenic DNA of *Drosophila*. *Genome Res*. 14:273–279.
- Halligan DL, Keightley PD. 2006. Ubiquitous selective constraints in the *Drosophila* genome revealed by a genome-wide interspecies comparison. *Genome Res*. 16:875–884.
- Hamblin MT, Veuille M. 1999. Population structure among African and derived populations of *Drosophila simulans*: evidence for ancient subdivision and recent admixture. *Genetics*. 153:305–317.
- Hill WG, Robertson A. 1966. The effect of linkage on the limits of artificial selection. *Genet Res*. 8:269–294.
- Jenkins DL, Ortori CA, Brookfield JF. 1995. A test for adaptive change in DNA sequences controlling transcription. *Proc R Soc B*. 261:203–207.
- Jukes TH, Cantor CR. 1969. Evolution of protein molecules. In: Munro HN, editor. *Mammalian protein metabolism III*. New York: Academic Press. p. 21–132.
- Kimura M. 1980. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide substitutions. *J Mol Evol*. 16:111–120.
- Kimura M. 1983. *The neutral theory of molecular evolution*. Cambridge: Cambridge University Press.
- Kirby DA, Muse SV, Stephan W. 1995. Maintenance of pre-mRNA secondary structure by epistatic selection. *Proc Natl Acad Sci USA*. 92:9047–9051.
- Kohn MH, Fang S, Wu C-I. 2004. Inference of positive and negative selection on the 5' regulatory regions of *Drosophila* genes. *Mol Biol Evol*. 21:374–383.
- Leicht BG, Muse SV, Hanczyc M, Clark AG. 1995. Constraints on intron evolution in the gene encoding the myosin alkali light chain in *Drosophila*. *Genetics*. 139:299–308.
- McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the *Adh* locus in *Drosophila*. *Nature*. 351:652–654.
- McVean GAT, Charlesworth B. 1999. A population genetic model for the evolution of synonymous codon usage: patterns and predictions. *Genet Res*. 74:145–158.
- McVean GAT, Vieira J. 2001. Inferring parameters of mutation, selection and demography from patterns of synonymous site evolution in *Drosophila*. *Genetics*. 157:245–257.
- Moriyama EN, Powell JR. 1996. Intraspecific nuclear DNA variation in *Drosophila*. *Mol Biol Evol*. 13:261–277.
- Nei M, Gojobori T. 1986. Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Mol Biol Evol*. 3:418–426.
- Nielsen R. 2005. Molecular signatures of natural selection. *Annu Rev Genet*. 39:197–218.
- Nielsen R, Bauer DuMont VL, Hubisz MJ, Aquadro CF. 2007. Maximum likelihood estimation of ancestral codon usage bias parameters in *Drosophila*. *Mol Biol Evol*. 24:228–235.
- Nielsen R, Weinreich DM. 1999. The age of nonsynonymous and synonymous mutations in mtDNA and implications for the mildly deleterious theory. *Genetics*. 153:497–506.
- Nolte V, Schlötterer C. 2008. African *Drosophila melanogaster* and *D. simulans* populations have similar levels of sequence variability, suggesting comparable effective population sizes. *Genetics*. 178:405–412.
- Ohta T. 1993. Amino acid substitution at the *Adh* locus of *Drosophila* is facilitated by small population size. *Proc Natl Acad Sci USA*. 90:4548–4551.
- Sawyer SA, Kulathinal RJ, Bustamante CD, Hartl DL. 2003. Bayesian analysis suggests that most amino acid replacements in *Drosophila* are driven by positive selection. *J Mol Evol*. 57:S154–S164.
- Smith NGC, Eyre-Walker A. 2002. Adaptive protein evolution in *Drosophila*. *Nature*. 415:1022–1024.
- Stephan W, Kirby DA. 1993. RNA folding in *Drosophila* shows a distance effect for compensatory fitness interactions. *Genetics*. 135:97–103.
- Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics*. 123:585–595.
- Templeton A. 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 the hominoid primates. *Genetics*. 144:1263–1270.
- Wall JD, Andolfatto P, Przeworski M. 2002. Testing models of selection and demography in *Drosophila simulans*. *Genetics*. 162:203–216.
- Welch JJ. 2006. Estimating the genomewide rate of adaptive protein evolution in *Drosophila*. *Genetics*. 173:821–837.
- Yang Z. 1997. PAML: a program package for phylogenetic analysis by maximum likelihood. *Comput Appl Biosci*. 13:555–556.

John H. McDonald, Associate Editor

Accepted May 20, 2008