

---

## Variation in synonymous codon use and DNA polymorphism within the *Drosophila* genome

Nicolas Bierne<sup>1,2,\*</sup> and Adam Eyre-Walker<sup>1</sup>

<sup>1</sup> Centre for the Study of Evolution & School of Biological Sciences, University of Sussex, Brighton, BN1 9QG, UK

<sup>2</sup> Laboratoire Génome, Populations, Interactions, Adaptation, UMR5171 CNRS-UMII IFREMER, SMEL, 1 Quai de la daurade, 34200 Sète, France

\*: Corresponding author : N. Bierne, email address : [n-bierne@univ-montp2.fr](mailto:n-bierne@univ-montp2.fr)

---

### Abstract:

A strong negative correlation between the rate of amino-acid substitution and codon usage bias in *Drosophila* has been attributed to interference between positive selection at nonsynonymous sites and weak selection on codon usage. To further explore this possibility we have investigated polymorphism and divergence at three kinds of sites: synonymous, nonsynonymous and intronic in relation to codon bias in *D. melanogaster* and *D. simulans*. We confirmed that protein evolution is one of the main explicative parameters for interlocus codon bias variation ( $r^2 \sim 40\%$ ). However, intron or synonymous diversities, which could have been expected to be good indicators of local interference [here defined as the additional increase of drift due to selection on tightly linked sites, also called 'genetic draft' by Gillespie (2000)] did not covary significantly with codon bias or with protein evolution. Concurrently, levels of polymorphism were reduced in regions of low recombination rates whereas codon bias was not. Finally, while nonsynonymous diversities were very well correlated between species, neither synonymous nor intron diversities observed in *D. melanogaster* were correlated with those observed in *D. simulans*. All together, our results suggest that the selective constraint on the protein is a stable component of gene evolution while local interference is not. The pattern of variation in genetic draft along the genome therefore seems to be instable through evolutionary times and should therefore be considered as a minor determinant of codon bias variance. We argue that selective constraints for optimal codon usage are likely to be correlated with selective constraints on the protein, both between codons within a gene, as previously suggested, and also between genes within a genome.

**Keywords:** Codon usage, selective constraints, substitution rate, diversity, GC-content, linkage, Hill-Robertson interference, *Drosophila*

# 1 **Introduction**

2 It is now widely accepted that weak selection for codon usage is acting upon synonymous  
3 mutations in some organisms (for review see Kurland, 1991; Sharp *et al.*, 1995; Akashi &  
4 Eyre-Walker, 1998) but is relaxed in others (e.g. in some bacteria, Sharp *et al.*, 2005; in  
5 Mammals, Duret, 2002). Variation in the intensity of synonymous selection has sometimes  
6 been detected at small evolutionary time scales between closely related species: in *Drosophila*  
7 for instance, population genetics analyses suggest that selection for codon usage is currently  
8 active in *D. simulans* (Akashi & Schaeffer, 1997; Kliman, 1999; Begun, 2001) while it seems  
9 to be relaxed in *D. melanogaster* owing to a recent reduction in population size in this lineage  
10 (Akashi, 1996; McVean & Vieira, 2001). However, it has proven extremely difficult to  
11 distinguish between competing explanations that account for interlocus variance in codon bias  
12 within a genome, once selection had been validated at the scale of the whole genome (Akashi,  
13 2001; Duret, 2002). The pattern of codon usage results from a balance among three factors: (i)  
14 selection, (ii) random genetic drift and (iii) mutation (Bulmer 1991). All three factors are  
15 likely to vary within a recombining genome and have been proposed in turn to be the  
16 causative agent.

17 (i) Interlocus variation in codon usage bias may more simply be the consequence of  
18 unequal selection coefficients across genes. The fitness difference between synonymous  
19 codons most probably relies on translation efficiency. In *Drosophila*, codon usage is biased  
20 toward preferred codons that generally correspond to the most abundant cognate tRNA  
21 (Moriyama & Powell, 1997). Variation in translational selection across genes is attested by a  
22 positive correlation between codon usage bias and the level of gene expression (Duret &  
23 Mouchiroud, 1999). In addition, a variable selective regime on synonymous mutations is  
24 further suggested by a negative correlation between codon bias and synonymous substitution  
25 rates (Sharp & Li, 1989; Moriyama & Hartl, 1993; Bierne & Eyre-Walker, 2003). However, it

1 remains unclear whether codon usage primarily affects the elongation rate or the fidelity of  
2 protein synthesis (Akashi, 2001; Duret, 2002). The latter hypothesis is supported in  
3 *Drosophila* and *Caenorhabditis*, where codon bias is stronger at constrained than at  
4 substituted amino acids (Akashi, 1994; Marais & Duret, 2001). However, the translational  
5 accuracy hypothesis would predict a positive correlation between codon bias and gene length,  
6 as is observed in the Prokaryote *Escherichia coli* (Eyre-Walker, 1996), while the reverse is  
7 observed in the metazoan genomes analysed to date (Moriyama & Powell, 1998; Duret &  
8 Mouchiroud, 1999; Comeron *et al.*, 1999). How much strengths of selection for the speed or  
9 the accuracy of translation vary across genes therefore remains unclear.

10 (ii) Another factor which can cause within-genome variation in selection efficacy is the  
11 Hill-Robertson (HR) effect (Hill & Robertson, 1966). The HR effect corresponds to a  
12 decrease in the efficacy of selection acting upon a mutation due to selection on other  
13 genetically linked segregating mutations. Selection, whatever its direction, increases the  
14 variance in reproductive success and consequently inflates genetic drift (Wright, 1931). The  
15 HR effect and related models for neutral mutations (i.e. hitchhiking, Maynard Smith & Haigh,  
16 1974; background selection, Charlesworth *et al.*, 1993) can therefore be understood as local  
17 variations in genetic drift (Felsenstein, 1974). Gillespie (2000) suggested that “genetic draft”  
18 could prove a useful label for the stochastic effects induced by indirect selection that are  
19 different in their origin and their statistic properties from purely demographic random drift.  
20 We will here follow Gillespie’s terminology and will state that the magnitude of the genetic  
21 draft but not genetic drift can vary along a recombining genome. The correlation observed  
22 between local recombination rates and gene diversity in *D. melanogaster* (Begun & Aquadro,  
23 1992) has classically been attributed to higher interferences (be they caused by positive or  
24 negative selection) in genomic regions of low recombination, in accordance with this idea of  
25 within-genome variation in genetic draft. Although some confounding factors obscure the

1 correlation (Marais *et al.*, 2001), codon usage bias also correlates with local recombination  
2 rates in *D. melanogaster* (Hey & Kliman, 2002; Marais & Piganeau, 2002), as expected by  
3 the HR effect. However, the correlation is very weak, accounting for only ~1% of the codon  
4 bias variance, and is restricted to lower recombination rate values (Marais & Piganeau, 2002).  
5 HR effect, although operating on codon bias (Hey & Kliman, 2002), was therefore thought to  
6 be a minor determinant of interlocus codon bias variance (Marais & Piganeau, 2002).  
7 However, Betancourt and Presgraves (2002) have documented a very strong negative  
8 correlation between the rate of amino-acid substitution and codon usage bias in *Drosophila*.  
9 These authors proposed that a small scale HR effect accounts for the correlation: genetic draft  
10 would be more intense in fast-evolving genes that undergo a high rate of selectively driven  
11 amino-acid substitutions and would thus be unable to optimise their codon usage. The HR  
12 effect hypothesis is an indirect interpretation of the correlation that would require further  
13 evidence; however Betancourt and Presgraves (2002) discussed and refuted alternative  
14 hypotheses making the HR effect the last possible explanation (but see Marais &  
15 Charlesworth, 2003). In addition, Kim (2004) has recently shown theoretically that the model  
16 is reasonable.

17 (iii) Finally, codon bias may still vary across genes that share similar selection intensities  
18 if a correlated mutational bias is superposed on selection on synonymous codon use. In  
19 *Drosophila*, all preferred codons end in G or C which results in a very strong positive  
20 correlation ( $R^2 > 0.9$ ) between codon bias and GC content at synonymous third coding  
21 positions (hereafter GC<sub>3</sub>). Furthermore, a good correlation ( $0.1 < R^2 < 0.3$ ) is observed between  
22 GC<sub>3</sub> and intron GC content (hereafter GC<sub>i</sub>) in accordance with the hypothesis that a non  
23 negligible mutational bias is superimposed on selection on codon use in this taxa (Kliman &  
24 Hey, 1994; Akashi *et al.*, 1998). Marais *et al.* (2001) have pointed out that GC<sub>i</sub> correlates  
25 positively with local recombination rate, as does GC<sub>3</sub>, in *D. melanogaster*, suggesting that a

1 part of the mutational bias may be associated with recombination. However, the correlation is  
2 very weak ( $R^2 < 1\%$ ), sometimes not detected (Hey & Kliman, 2002), and recombination only  
3 accounts for a very small fraction of the correlation between  $GC_3$  and  $GC_i$ .

4 In short, each factor (selection for the speed or the accuracy of translation, genetic draft  
5 and mutation bias) appears to be operating concomitantly in *Drosophila*. However, their  
6 relative contribution is still unclear. Most importantly, the recent observation that codon usage  
7 bias is strongly correlated with the rate of amino-acid substitution, led to the suggestion that  
8 small scale variations in local interference within the *Drosophila* genome may have a stronger  
9 impact on interlocus codon bias variance than previously thought (Betancourt & Presgraves,  
10 2002; Kim, 2004).

11 In the present study we propose (i) to explore the effect of within-genome variation in  
12 local interference on codon usage bias using a more direct measure of genetic draft, relative  
13 levels of synonymous and intron polymorphism, (ii) to use a comparative analysis between *D.*  
14 *melanogaster* and *D. simulans* in order to investigate the relative stability through time of  
15 local interference on the one hand, and selective constraint on the protein on the other hand,  
16 and (iii) to measure the relative contributions of various factors to synonymous codon bias.  
17 Since our sample size is constrained by the availability of polymorphism data at three kinds of  
18 sites (intron, synonymous and non-synonymous sites) which could only be found for a few  
19 tens of genes at present, we could only study the factors which have a major impact on  
20 synonymous codon use.

## 1 **Data**

### 2 **Dataset 1- Intron and synonymous diversity, synonymous codon bias and within-genome** 3 **variation in local interference.**

4 One aim of the present work was to investigate intron and synonymous diversity as a  
5 correlate of local interference in order to further explore the model of small scale HR effect  
6 proposed by Betancourt and Presgraves (2002). To do this, we have compiled from the  
7 literature a first dataset composed of 38 genes in *D. melanogaster* and 34 genes in *D.*  
8 *simulans* for which polymorphism data were available at three kinds of sites –intron,  
9 synonymous and non-synonymous sites (see Bierne & Eyre-Walker, 2004). Data were  
10 available in both species for 23 genes. Unfortunately, the constraint on the data resulted in too  
11 few genes sharing a similar sampling scheme in *D. melanogaster* to conduct correlation  
12 analyses with reasonable sampling sizes. We therefore have chosen to use as many genes as  
13 possible keeping in mind the potential caveat surrounding the dataset. On the other hand, the  
14 sampling scheme of the *D. simulans* dataset was very homogeneous as the sequence data  
15 comes from a single population in California (Begun & Whitley, 2000). For each locus we  
16 have computed the non-synonymous, synonymous and intron diversities within a species as  
17 well as substitution rates between species (respectively  $\theta_{c_n}$ ,  $\theta_{c_s}$ ,  $\theta_i$  and  $D_{c_n}$ ,  $D_{c_s}$ ,  $D_i$ ) using  
18 DnaSP (Rozas & Rozas, 1999). Since there is ambiguity surrounding the definition of a site in  
19 coding sequences which can lead to a spurious correlation between rates of non-synonymous  
20 substitution and codon bias (Muse, 1996; Bierne & Eyre-Walker, 2003), we used synonymous  
21 and non-synonymous rates *per codon* instead of *per site*, while intron rates were of course *per*  
22 *site*. As a consequence, intron and synonymous evolutionary rates were not directly  
23 comparable; however this was not the purpose of this analysis where a correlative approach  
24 was used. In addition, correction for multiple hits could safely be ignored in the  
25 polymorphism as well as the divergence between the closely related species *D. melanogaster*

1 and *D. simulans*. Diversities were estimated from the number of polymorphic sites and  
2 sample size (Watterson, 1975) or the average number of nucleotide differences per site  
3 between two sequences (Nei, 1987); however the results were very similar with both  
4 estimators and we arbitrarily choose Watterson's estimator to present in the results below.  
5 Codon usage bias was first measured by the frequency of optimal codons (*Fop*, Ikemura,  
6 1985) by using the program CODONW (Peden, 1999). Other measures of codon bias such as  
7 the effective number of codons (ENC, Wright, 1990) or the codon bias index (CBI, Morton,  
8 1993) are sometimes used in the literature. These measures were very well correlated to *Fop*  
9 and gave qualitatively similar results. We chose *Fop* because optimal codons have been nicely  
10 defined from data on expression levels in *Drosophila* (Duret & Mouchiroud, 1999). As  
11 explained above however, *Fop* and other measures of codon bias are correlated to GC content  
12 in *Drosophila*, which is taken as evidence for the action of a mutational bias sometimes  
13 favouring the same codons as selection does. In order to account for variations in mutational  
14 patterns we computed the residuals of the *Fop*/GC<sub>i</sub> correlation, that we named *Fop*-GC<sub>i</sub>. *Fop*-  
15 GC<sub>i</sub> is expected to measure the single action of selection freed from the cumulative effect of  
16 mutation bias. Finally we estimated recombination rates in *D. melanogaster* (hereafter *R<sub>mel</sub>*)  
17 by using the data and standard method of Kliman and Hey (1993). We are aware of the debate  
18 surrounding the accuracy of different estimates of recombination rates in *Drosophila* (Marais  
19 *et al.*, 2003; Kliman & Hey, 2003); however we chose the same estimator as the one used by  
20 Betancourt and Presgraves (2002) as a basis for comparison, having verified that other  
21 estimates gave qualitatively similar results. The data were compiled into a spreadsheet which  
22 is provided in the Supplementary file available on the journal web site.

23

24

25

## 1 **Dataset 2- relative contribution of each factor**

2 In a second analysis we have investigated the relative contribution of each factor potentially  
3 involved in codon bias variation –e.g. protein evolution, expression level, gene length and  
4 surrounding non-coding GC content. Since we were not constrained by the need to have  
5 intron polymorphism data in this analysis we were able to compile a much larger dataset of  
6 genes. Thanks to the recent effort to produce large polymorphism datasets in *D. simulans* (e.g.  
7 Begun & Whitley, 2000; Schlenke & Begun, 2003), we were able to compute non-  
8 synonymous and synonymous diversities from 105 genes in this species. We chose to measure  
9 the level of constraint on a protein using polymorphism data, rather than divergence data,  
10 because adaptive substitutions can affect divergence estimates. However, qualitatively similar  
11 results were obtained using dN/dS. Codon usage bias was measured by *Fop*. We did not  
12 correct for local GC content in this analysis because we were interested in assessing the  
13 mutation bias effect (while in the previous dataset we wanted to remove it). Non-coding GC  
14 content (hereafter GC<sub>nc</sub>) usually was GC<sub>i</sub> but for the few genes without introns, GC<sub>nc</sub> was  
15 computed from the surrounding non-coding DNA (500 to 1000 bp on either side depending  
16 on the distance from adjacent genes) having verified that GC<sub>i</sub> and GC<sub>nc</sub> were very well  
17 correlated in the *Drosophila* genome ( $R^2=0.99$ ,  $P < 0.001$ ). Rough estimates of gene  
18 expression levels were measured by EST-counting using the procedure described in Duret and  
19 Mouchiroud (1999). Finally, gene length was also considered in the analysis as it is known to  
20 correlate with codon usage in *Drosophila* (Powell & Moriyama, 1997; Duret & Mouchiroud,  
21 1999). The data were compiled into a spreadsheet which is provided together with dataset 1 in  
22 the Supplementary file available on the journal web site.



## 1 **Analyses**

### 2 **Comparison of the effect of recombination and amino acid substitution rates on silent** 3 **diversity and codon usage bias in *D. melanogaster***

4 To begin with, we have explored the possibility that the relevant scale for variations in local  
5 interference could be better captured by amino acid substitution rates than by local  
6 recombination rates. We were also interested in verifying whether the within-genome  
7 variance in local interference was accurately captured by the data in *D. melanogaster* despite  
8 an unbalanced sampling scheme between loci. In accordance with a well known observation  
9 in *Drosophila* (Begun & Aquadro, 1992), the diversity measured at synonymous and intronic  
10 sites ( $\theta_{+s}$ ) was significantly correlated with recombination ( $R_{mel}$ ) in our dataset 1 (Fig. 1A). In  
11 order to estimate the noise introduced by heterogeneous sampling in our meta-analysis of *D.*  
12 *melanogaster* data, we have plotted, in the same graph, results from a survey which was  
13 devoid of bias in the sampling strategy (Andolfatto & Przeworski, 2001). Andolfatto and  
14 Przeworski (2001) used sequence data from a single population of Zimbabwe in Africa; few  
15 loci, though, were screened simultaneously in exons and introns as required for the present  
16 analysis. Figure 1A shows that the correlation we obtained with dataset 1 does not differ  
17 greatly from the one obtained by Andolfatto and Przeworski (2001). A Levene's test of  
18 homogeneity of variance reveals that the variance of  $\theta$  are not significantly different in the  
19 two datasets ( $F_{1,73}=2.79$ ,  $P=0.1$ ). In addition, the method of Stephan (1995) was used to fit the  
20 curve expected under a model of recurrent selective sweeps (Stephan *et al.*, 1992). The fitted  
21 curves were roughly the same with the two datasets (Fig. 1A). This comparison suggests that  
22 variations in local interference are accurately captured in dataset 1 and that the heterogeneity  
23 of this dataset introduces neither bias nor substantial statistical noise. In the same dataset  
24 however, codon usage bias ( $F_{op-GC_i}$ ) did not correlate significantly with  $R_{mel}$  (Fig. 1B) while  
25 it was strongly correlated with amino acid substitution rate,  $Dc_n$  (Fig. 1D), as previously

1 reported (Betancourt & Presgraves, 2002; Marais *et al.*, 2004). The hypothesis of a small  
2 scale HR effect responsible for low codon bias in fast evolving genes would have predicted a  
3 correlation between  $Dc_n$  and  $\theta_{i+s}$ , but this was not the case (Fig. 1C). Conflicting results were  
4 therefore obtained depending on whether amino acid substitution rates or local recombination  
5 rates were used to assess local interference.

### 6 [Figure 1]

7

### 8 **The relationship between DNA variation and codon usage bias**

9 Table 1 presents the various correlations obtained between codon bias as measured by *Fop-*  
10  $GC_i$  and DNA variation decomposed into three classes of mutations (i.e. non-synonymous,  
11 synonymous and introns) within a species (i.e. diversity) and between species (i.e.  
12 divergence). In accordance with previous results (Sharp & Li, 1989; Marais *et al.*, 2004),  
13 significant negative correlations were obtained between codon bias and non-synonymous, as  
14 well as synonymous substitution rates. In contrast, intron divergence did not correlate  
15 significantly with codon bias (Table 1), or with  $GC_i$  ( $r_s=-0.12$ , n.s.), in accordance with the  
16 neutral expectation. The correlation previously observed with non-synonymous divergence  
17 was here extended to non-synonymous diversity within both species. Non-synonymous  
18 diversity is most likely composed of neutral or nearly neutral mutations and it was unclear  
19 whether one would have expected non-synonymous diversity to be a good index of the  
20 density of selected mutations and local genetic draft. Nevertheless, non-synonymous diversity  
21 is known to be highly correlated with non-synonymous divergence in *Drosophila* because  
22 selective constraint on the protein is the main determinant of non-synonymous variation (see  
23 Bierne and Eyre-Walker 2004). Therefore, we would not take this observation as an argument  
24 against the HR effect hypothesis. Synonymous diversity, on the other hand, was not  
25 significantly correlated with codon bias which was predictable in *D. melanogaster* where

1 synonymous selection is thought to be relaxed (Akashi, 1996) but could have been expected  
2 in *D. simulans*. Indeed, the results obtained with synonymous mutations were consistent with  
3 the effect of weak selection acting on synonymous codons (Bulmer, 1991). Finally, and most  
4 importantly, intron diversity, which could have been expected to be an unbiased indicator of  
5 local interference, did not covary significantly with codon bias (Table 1).

## 6 [Table 1]

### 7 **Comparative analysis**

8 Our dataset 1 allowed us to compare diversities between different classes of mutations within  
9 a species, and for the same class of mutations, compare the diversity realised in each species.  
10 Such an analysis was conducted to further investigate how local interference could evolve  
11 between species. Both in *D. melanogaster* and *D. simulans*, synonymous and intron  
12 diversities were strongly correlated (Fig. 2A). However, neither intron nor synonymous  
13 diversity within *D. melanogaster* was correlated with intron or synonymous diversity in *D.*  
14 *simulans* (Fig. 2B). Non-synonymous diversity on the other hand did not correlate  
15 significantly with intron (Fig. 2C) or synonymous (not shown, *D. melanogaster*:  $r_s=0.29$ , n.s.;  
16 *D. simulans*:  $r_s=0.13$ , n.s.) diversity. Non-synonymous diversity within *D. melanogaster*  
17 however was strongly correlated with non-synonymous diversity within *D. simulans* (Fig.  
18 2D).

## 19 [Figure 2]

20 Taken together, these results illustrate that intron and synonymous variations are mainly  
21 driven by stochastic processes (genetic drift and draft) that are not stable components through  
22 evolutionary times, while non-synonymous variation is mainly driven by selective constraint  
23 on the protein which in contrast seems to be a stable element, at least between closely related  
24 species.

25

1 **The relative contribution of each factor**

2 We can now reconsider the relative contribution of various factors thought to be involved in  
3 the codon bias variance with our dataset 2. Correlation statistics are presented in Table 2. We  
4 chose the ratio  $\omega = \theta_{c_n} / \theta_{c_s}$  as a measure of the selective constraint on the protein (small  $\omega$   
5 indicates more constraint on the amino acid sequence). In decreasing order the parameters that  
6 appeared to explain most of the variation in codon usage were (i) the selective constraint on  
7 the protein as measured by the  $\omega$  ratio, which explains ~40% of the codon bias variance, (ii)  
8 the local mutational pattern as measured by surrounding non-coding GC content which  
9 explains ~15% of the codon bias variance and (iii) expression levels as measured by EST-  
10 counting which explains ~10% of the codon bias variance. None of these three parameters co-  
11 vary significantly with each other in this dataset suggesting that they correspond to almost  
12 independent factors. Note however that a correlation between the level of gene expression and  
13 non-synonymous evolutionary rates has been described elsewhere, although with a much  
14 larger dataset (Marais *et al.*, 2004) or in other organisms, where selection on codon usage is  
15 relaxed such as Mammals (Duret & Mouchiroud, 2000). Finally, the correlation between  
16 codon bias and gene length which has previously been reported with very large datasets  
17 (Powell & Moriyama, 1997; Duret & Mouchiroud, 1999) was not significant in our dataset 2  
18 (Table 2).

19 **[Table 2]**

## 1 **Discussion**

2 Polymorphism and divergence data at three kinds of sites –synonymous, nonsynonymous and  
3 intronic– were used to investigate the importance of within-genome variations in local  
4 interference on the evolution of codon usage in *Drosophila*. We first argue that our results  
5 suggest that fast evolving genes do not have conspicuously higher levels of genetic draft. In  
6 addition, a comparative analysis between *D. melanogaster* and *D. simulans* suggests that local  
7 interference is unlikely to be a stable component of gene evolution while selective constraint  
8 on the protein is. All together our results suggest that the correlation between synonymous  
9 codon usage and protein evolution cannot be exclusively interpreted by local interference  
10 between selection at non-synonymous and synonymous sites. We will finally discuss  
11 alternative explanations involving some connections between selection on the protein and  
12 selection for the speed or the accuracy of translation.

13

### 14 **Synonymous and intron diversities do not corroborate a more intense genetic draft in** 15 **the recent history of fast-evolving genes**

16 Since the publication of the correlation between local recombination rates and gene diversity  
17 (Begun & Aquadro, 1992), local variation in genetic draft within the *Drosophila* genome has  
18 been thoroughly investigated. Using recombination rates to assess the intensity of genetic  
19 draft, within-genome variation in local interference has been suspected to influence the  
20 efficacy of weak selection on various genomic components such as intron length (Carvalho &  
21 Clark, 1999) or codon usage (Kliman & Hey, 1993; Comeron *et al.*, 1999). More recently,  
22 other parameters thought to correlate with local interference have been investigated such as  
23 gene length, the presence/absence of introns, or the spatial situation of targeting sites in the  
24 gene (Comeron & Kreitman, 2002). Most of these correlations are minute, accounting for a  
25 minor part of the total variance and thus require very large datasets (often exhaustive genome-

1 wide datasets) to be detected. In addition, some confounding mutational biases have  
2 sometimes been identified (Marais *et al.*, 2001). The weakness of the correlation are perhaps  
3 not surprising given the relevant estimates of recombination and mutation rates (Marais &  
4 Piganeau, 2002). On the contrary, the correlation observed between the rate of protein  
5 evolution and codon bias is surprisingly strong ( $R^2 > 40\%$ ). It is so strong that it does not  
6 require a large dataset to detect; neither does it require the presence of genes with particularly  
7 high rates of amino-acid substitution. As a consequence, if the correlation was entirely due to  
8 HR effects, one could have expected a detectable effect on levels of polymorphism (McVean  
9 & Charlesworth, 2000). However, neither synonymous nor intron diversities significantly  
10 correlate with protein evolution nor do they correlate with codon bias. In the same dataset,  
11 diversities were significantly reduced in region of low recombination rates whereas codon  
12 bias was not. Therefore, it seems difficult to summarise the results obtained within a single  
13 framework, namely HR effects.

14 *Drosophila* populations are known to exhibit complex patterns of genetic diversity that are  
15 not consistent with any simple model at demographic equilibrium (Andolfatto & Przeworski,  
16 2000; Begun, 2001; Wall *et al.*, 2002). *D. melanogaster* and *D. simulans* are thought to have  
17 spread across the world from Africa after the last glaciation (David & Capy, 1988). Derived  
18 populations are known to depart from demographic equilibrium (Begun & Aquadro, 1993;  
19 Begun, 2001; Baudry *et al.*, 2004) but the situation in Africa is not straightforward either  
20 (Glinka *et al.*, 2003; Veuille *et al.*, 2004). Indeed, it is likely that natural populations never  
21 conform to the standard population genetic assumptions (Lewontin, 2002). One may therefore  
22 suspect that departures from equilibrium could introduce unpredicted stochastic variance  
23 preventing any solid interpretation of the data. However, we would argue that (i) equilibrium  
24 does not need to be assumed here as demographic processes should affect the whole genome  
25 in a similar way such that within-genome variation captured in a correlation analysis can only

1 come from non-demographic processes (*i.e.* genetic draft), (ii) the significant correlation  
2 obtained between polymorphism levels and recombination rates attests that a fraction of the  
3 within-genome variation in local interference is accurately captured in the data and (iii) for a  
4 factor to have a bearing on the long term evolution of a trait with such a minuscule phenotypic  
5 consequence as codon bias, its effect should probably surpass the stochastic variance  
6 inevitably generated in every natural population. In our dataset, two correlations –between  
7 recombination and silent diversity and between non-synonymous polymorphism and codon  
8 bias– have proved to persist despite enduring the stochastic pressure.

9       Alone, though, the apparently conflicting observations we reported are not sufficient to  
10 completely refute small scale HR effects because codon bias, polymorphism levels and  
11 recombination are parameters that evolve at different time-scales. Diversity may not be  
12 reduced in fast-evolving genes nowadays but might have been in the past. Because codon bias  
13 depends on long-term evolution (Marais *et al.*, 2004), forces acting on it should be rather  
14 stable components of gene evolution.

15

## 16 **Local interference is not a stable component of gene evolution**

17 Local interference depends on the density of selected sites, the strength of the selection acting  
18 on selected sites and local recombination rates (McVean & Charlesworth, 2000; Stephan &  
19 Kim, 2002). Evidence has recently accumulated which suggests that local recombination rates  
20 are not stable over even short timescales (e.g. Munte *et al.*, 2001; Takano-Shimizu, 2001;  
21 Meunier & Duret, 2004). For instance, Ptak *et al.* (2005) have demonstrated that the  
22 recombination landscape has markedly changed during the human/chimp divergence. These  
23 results would suggest that local interference may vary accordingly in time. However, the  
24 possibility remains that the variation in local selection (density and strength of selection)  
25 prevails over the variation in local recombination rate, as implicitly assumed in the model of

1 Betancourt and Presgraves (2002). To assess the stability of local interference, we have here  
2 conducted a comparative analysis of polymorphism levels between *D. melanogaster* and *D.*  
3 *simulans*. Neither intron nor synonymous diversity within *D. melanogaster* was correlated  
4 with intron or synonymous diversity in *D. simulans*, suggesting that local interference is not a  
5 very stable component of gene evolution. Instead of the apparent stochastic nature of silent  
6 diversity, non-synonymous diversities were very well correlated between species suggesting  
7 that selective constraints are conserved across species. In accordance with this view, Munte *et*  
8 *al.* (2001) showed that the recombinational environment of a gene strongly conditions  
9 synonymous substitution rates while it has no detectable effect on amino acid evolutionary  
10 rates in *Drosophila*.

11 All together, our results suggest that the selective constraint on the protein is a stable  
12 component of gene evolution (also see Skibinski & Ward, 2004) while local interference is  
13 not.

14

### 15 **Correlated selective constraints on synonymous and non-synonymous sites**

16 Our evidence suggests that HR effects are not a strong determinant of codon bias, but why  
17 then is there a correlation between synonymous codon usage and rates of protein evolution?  
18 The alternatives have been well discussed elsewhere (Akashi, 1994; Betancourt & Presgraves,  
19 2002; Marais *et al.*, 2004) but we reiterate them here briefly. Although attractive at first sight,  
20 non-synonymous changes that transform a preferred codon into an unpreferred codon  
21 (Lipman & Wilbur, 1985) cannot reasonably account for the correlation. Indeed, removing  
22 such codons (which represent 19% of nonsynonymous changes) has no effect on the strength  
23 of the correlation between the rate of non-synonymous substitution and codon usage bias  
24 (data not shown, see Akashi, 1994; Marais & Duret, 2001; Betancourt & Presgraves, 2002).



1 It is also easy to refute another possibility –that the correlation arises through the way in  
2 which sites are counted in the estimation of the non-synonymous substitution rate. In the  
3 method of Goldman and Yang (1994), the method used by Betancourt and Presgraves (2002),  
4 sites are counted as mutational opportunities (see Bierne and Eyre-Walker 2003), so as codon  
5 bias increases the number of synonymous sites decreases and the number of non-synonymous  
6 sites increases. This means that genes with high codon bias will tend to have lower rates of  
7 non-synonymous substitution per site (i.e. if two genes have undergone similar numbers of  
8 non-synonymous substitutions per codon, the gene with the higher level of codon bias will  
9 actually have a lower rate of non-synonymous substitution per site). However, the rate of  
10 nonsynonymous substitution *per codon* is also correlated to codon bias.

11 This leaves an idea originally proposed by Akashi (1994), that the strength of selection  
12 acting upon synonymous mutations is correlated to that acting upon non-synonymous  
13 mutations. This could be due to selection on translational accuracy –genes in which most  
14 amino acid sites need to be occupied by a particular amino acid will evolve slowly and will  
15 need to accurately translate. Betancourt and Presgraves (2002) offered several lines of  
16 evidence against this hypothesis. First they noted that the rate of synonymous substitution was  
17 positively correlated to codon bias in their analysis while it was generally accepted the  
18 correlation was negative. However, this was an artefact of the method they used, as we have  
19 discussed elsewhere (Bierne & Eyre-Walker, 2003) –the rate of synonymous substitution does  
20 correlate negatively and significantly with codon bias (Table 1) as previously reported (Sharp  
21 & Li, 1989). Second, Betancourt and Presgraves (2002) tested this hypothesis by considering  
22 the correlation between the level of codon bias in codons which had not undergone a non-  
23 synonymous substitution and the overall rate of non-synonymous substitution. They found the  
24 correlation was unchanged and concluded that there was no evidence of correlated strengths  
25 of selection. To explain the logic of their test let us consider a pair of two fold degenerate

1 codons - phenylalanine for example. Let us imagine that the average strength of selection  
2 against non-synonymous mutations is  $s_n$ . Errors during translation will have an effect on the  
3 fitness of the individual which is correlated to this average strength (the correlation will not be  
4 perfect, because while TTT to TTA mutations might be common, TTT to TTA translational  
5 errors may not be). This will manifest itself as selection on synonymous codon bias; so the  
6 strength of selection on codon bias will be correlated to the strength of selection against  
7 deleterious mutations. The average strength of selection against non-synonymous mutations  $s_s$   
8 is therefore equal to  $k s_n$ , where  $k$  is a constant. It seems likely, unless the translational error  
9 rate is very high that  $k < 1$ . Let us now think about all the phenylalanines in a gene. Some will  
10 be very important because they are critical for function and others will not be. We can divide  
11 the sites into three categories; (i) sites at which  $N_e s_s < 1$  and  $N_e s_n < 1$  – i.e. selection at both sites  
12 is ineffective; (ii) sites at which  $N_e s_s < 1$  and  $N_e s_n > 1$  – sites at which selection is effective  
13 against the non-synonymous mutations, but ineffective on synonymous codon use; and (iii)  
14 sites at which  $N_e s_s > 1$  and  $N_e s_n > 1$  – codons at which selection is effective on both non-  
15 synonymous and synonymous mutations. The rate of non-synonymous substitution, ignoring  
16 adaptive evolution is determined by the proportion of sites in category (i) relative to  
17 categories (ii) and (iii), while the level of synonymous codon use is determined by the  
18 proportion of sites in (i) and (ii) relative to (iii). Betancourt and Presgraves (2002) just looked  
19 at synonymous codon use at codons with no amino acid substitution which would be  
20 equivalent to looking at the relative number of codons in category (ii) versus (iii). It is clear  
21 that if the  $s_n$ 's in a gene are independently and randomly drawn from some distribution then  
22 there will be no correlation between the rate of nonsynonymous substitution and the level of  
23 bias in codons which have not undergone amino acid substitution – this would be equivalent  
24 to randomly allocating codons to the three categories and so there is no expectation of a  
25 correlation between (i)/(i+ii+iii) and ii/(ii+iii). However, their test is not valid if there is a

1 correlation between  $s_n$  at different sites within a gene; *i.e.* if genes with strong selection  
2 against non-synonymous mutations at one codon also tend to have strong selection at other  
3 codons. This is indeed the case – for example the two halves of a gene have correlated rates of  
4 non-synonymous substitution (Smith and Eyre-Walker 2002).

5         The strong correlation between codon bias and rates of non-synonymous substitution,  
6 or levels of non-synonymous polymorphism, and our explanation for the correlation, suggest  
7 that selection on codon usage bias is primarily driven by translational accuracy. This is  
8 supported by the fact that constrained codons tend to have higher levels of codon bias  
9 (Akashi, 1994; Marais & Duret, 2001). However, this effect was not very strong and the  
10 positive correlation between codon bias and gene length predicted by the translational  
11 accuracy hypothesis (Eyre-Walker, 1996) was negative instead (Duret & Mouchiroud, 1999).  
12 Codon bias is expected to be stronger in longer genes under translational accuracy because  
13 mistakes in longer genes will be energetically more costly. However, controlling for gene  
14 function is difficult in this type of analysis – *i.e.* it may be that longer genes tend to be poorly  
15 constrained and therefore fast evolving. If the inter-locus variance in selection regime  
16 overwhelms the gene length effect, the correlation would not necessarily be found.  
17 Furthermore, the fact that the correlation was not strong could be explained by an  
18 overabundance of non-selectively constrained codons in the non-substituted class when the  
19 comparison involved closely related species (Akashi, 1994) and an overabundance of  
20 selectively constrained codons in the substituted class (*i.e.* covarion-like evolution, Fitch,  
21 1971) when the comparison involved distant species (Marais & Duret, 2001). Finally, the  
22 variance in selective constraints for optimal protein synthesis may be more easily  
23 encapsulated between genes within a genome than it is between codons within a gene. Indeed,  
24 it is likely that very constrained proteins that play a major role in the correlation, are  
25 constrained at nearly every amino-acid, the reverse being true for fast-evolving proteins.

1

## 2 **Conclusion**

3 We found evidence against HR effects as a suitable explanation for the correlation between  
4 the rate of amino-acid substitution and codon usage bias in *Drosophila*. Although there are  
5 theoretical reasons to believe (Hill & Robertson, 1966; McVean & Charlesworth, 2000; Kim,  
6 2004) and empirical data to suggest (Hey and Kliman 2002) that HR interferences are  
7 operating on codon bias they cannot reasonably explain such a strong correlation and should  
8 be viewed as a minor determinant of interlocus codon bias variance (Marais and Piganeau  
9 2002).

10 We would therefore conclude that variation in codon usage within the *Drosophila* genome  
11 is mainly a simple consequence of unequal selection coefficients across genes. Discriminating  
12 between selection for the speed or the accuracy of protein synthesis is difficult but our  
13 analysis suggests that the fidelity of translation may be a more important component than  
14 previously thought. Usually investigated at the codon level within a gene, the effect of  
15 selection on the accuracy of translation may more markedly be seen at the gene level within a  
16 genome.

17

## 18 **Acknowledgements**

19 We are very grateful to people of the Centre for the Study of Evolution for helpful discussions  
20 on the issue of within-genome variations in effective population size and to two anonymous  
21 referees for insightful comments on the manuscript. The authors were supported by the  
22 Biotechnology and Biological Sciences Research Council and the Royal Society.

## 1 **References**

- 2 Akashi, H. 1994 Synonymous codon usage in *Drosophila melanogaster*: natural selection and  
3 translational accuracy. *Genetics* **136**: 927-935.
- 4 Akashi, H. 1996 Molecular evolution between *Drosophila melanogaster* and *D.simulans*:  
5 reduced codon bias, faster rates of amino acid substitution, and larger proteins in  
6 *D.melanogaster*. *Genetics* **144**: 1297-1307.
- 7 Akashi, H. 2001 Gene expression and molecular evolution. *Curr. Op. Genet. Dev.* **11**: 660-  
8 666.
- 9 Akashi, H., and A. Eyre-Walker 1998 Translational selection and molecular evolution. *Curr.*  
10 *Op. Genet. Dev.* **8**: 688-693.
- 11 Akashi, H., R. M. Kliman and A. Eyre-Walker 1998 Mutation pressure, natural selection, and  
12 the evolution of base composition in *Drosophila*. *Genetica* **102/103**: 49-60.
- 13 Akashi, H., and S. W. Schaeffer 1997 Natural selection and the frequency distributions of  
14 “silent” DNA polymorphism in *Drosophila*. *Genetics* **146**: 295-307.
- 15 Andolfatto, P., and M. Przeworski 2000 A genome-wide departure from the standard neutral  
16 model in natural populations of *Drosophila*. *Genetics* **156**: 257-268.
- 17 Andolfatto, P., and M. Przeworski 2001 Regions of lower crossing over harbor more rare  
18 variants in African populations of *Drosophila melanogaster*. *Genetics* **158**: 657-665.
- 19 Baudry, E., B. Viginier and M. Veuille 2004 Non-African populations of *Drosophila*  
20 *melanogaster* have a unique origin. *Mol Biol Evol* **21**: 1482-1491.
- 21 Begun, D. 2001 The frequency distribution of nucleotide variation in *Drosophila simulans*.  
22 *Mol. Biol. Evol.* **18**: 1343-1352.
- 23 Begun, D., and C. F. Aquadro 1993 African and North American populations of *Drosophila*  
24 *melanogaster* are very different at the DNA level. *Nature* **365**: 548-550.
- 25 Begun, D. J., and C. F. Aquadro 1992 Levels of naturally occurring DNA polymorphism  
26 correlate with recombination rates in *D. melanogaster*. *Nature* **356**: 519-520.
- 27 Begun, D. J., and P. Whitley 2000 Reduced X-linked nucleotide polymorphism in *Drosophila*  
28 *simulans*. *Proceedings of the National Academy of Sciences of the United States of*  
29 *America* **97**: 5960-5965.
- 30 Betancourt, A. J., and D. C. Presgraves 2002 Linkage limits the power of natural selection in  
31 *Drosophila*. *Proc Natl Acad Sci U S A* **99**: 13616-13620.
- 32 Bierne, N., and A. Eyre-Walker 2003 The problem of counting sites in the estimation of the  
33 synonymous and non-synonymous substitution rates: implications for the correlation  
34 between the synonymous substitution rate and codon usage bias. *Genetics* **165**: 1587–  
35 1597.
- 36 Bierne, N., and A. Eyre-Walker 2004 The genomic rate of adaptive amino-acid substitution in  
37 *Drosophila*. *Mol Biol Evol* **21**: 1350-1360.
- 38 Bulmer, M. 1991 The selection-mutation-drift theory of synonymous codon usage. *Genetics*  
39 **129**: 897-907.
- 40 Carvalho, A. B., and A. G. Clark 1999 Intron size and natural selection. *Nature* **401**: 344.
- 41 Charlesworth, B., M. T. Morgan and D. Charlesworth 1993 The effect of deleterious  
42 mutations on neutral molecular variation. *Genetics* **134**: 1289-1303.
- 43 Comeron, J. M., and M. Kreitman 2002 Population, evolutionary and genomic consequences  
44 of interference selection. *Genetics* **161**: 389-410.
- 45 Comeron, J. M., M. Kreitman and M. Aguade 1999 Natural selection on synonymous sites is  
46 correlated with gene length and recombination in *Drosophila*. *Genetics* **151**: 239-249.
- 47 David, J. R., and P. Capy 1988 Genetic variation of *Drosophila melanogaster* natural  
48 populations. *Trends Genet.* **4**: 106-111.

- 1 Duret, L. 2002 Evolution of synonymous codon usage in metazoans. *Curr. Opin. Genet. Dev.*  
2 **12**: 640-649.
- 3 Duret, L., and D. Mouchiroud 1999 Expression pattern and, surprisingly, gene length shape  
4 codon usage in *Caenorhabditis*, *Drosophila*, and *Arabidopsis*. *Proc. Natl. Acad. Sci.*  
5 *USA* **96**: 4482-4487.
- 6 Duret, L., and D. Mouchiroud 2000 Determinants of substitution rates in mammalian genes:  
7 expression pattern affects selection intensity but not mutation rate. *Mol Biol Evol* **17**:  
8 68-74.
- 9 Eyre-Walker, A. 1996 Synonymous codon bias is related to gene length in *Escherichia coli*:  
10 Selection for translational accuracy? *Mol. Biol. Evol.* **13**: 864-872.
- 11 Felsenstein, J. 1974 The evolutionary advantage of recombination. *Genetics* **78**: 737-756.
- 12 Fitch, W. M. 1971 Rate of change of concomitantly variable codons. *J Mol Evol* **1**: 84-96.
- 13 Gillespie, J. H. 2000 Genetic drift in an infinite population: The pseudohitchhiking model.  
14 *Genetics* **155**: 909-919.
- 15 Glinka, S., L. Ometto, S. Mousset, W. Stephan and D. De Lorenzo 2003 Demography and  
16 natural selection have shaped genetic variation in *Drosophila melanogaster*: a multi-  
17 locus approach. *Genetics* **165**: 1269-1278.
- 18 Goldman, N., and Z. Yang 1994 A codon based model of nucleotide substitution for protein-  
19 coding DNA sequences. *Mol. Biol. Evol.* **11**: 725-736.
- 20 Hey, J., and R. M. Kliman 2002 Interactions between natural selection, recombination and  
21 gene density in the genes of *Drosophila*. *Genetics* **160**: 595-608.
- 22 Hill, W. G., and A. Robertson 1966 The effect of linkage on limits to artificial selection.  
23 *Genet. Res.* **8**: 269-294.
- 24 Ikemura, T. 1985 Codon usage and tRNA content in unicellular and multicellular organisms.  
25 *Mol. Biol. Evol.* **2**: 13-34.
- 26 Kim, Y. 2004 Effect of strong directional selection on weakly selected mutations at linked  
27 sites: Implication for synonymous codon usage. *Mol Biol Evol* **21**: 286 - 294.
- 28 Kliman, R. 1999 Recent selection on synonymous codon usage in *Drosophila*. *J. Mol. Evol.*  
29 **49**: 343-351.
- 30 Kliman, R. M., and J. Hey 1993 Reduced natural selection associated with low recombination  
31 in *Drosophila melanogaster*. *Molecular Biology and Evolution* **10**: 1239-1258.
- 32 Kliman, R. M., and J. Hey 1994 The effects of mutation and natural selection on codon bias in  
33 the genes of *Drosophila*. *Genetics* **137**: 1049-1056.
- 34 Kliman, R. M., and J. Hey 2003 Hill-Robertson interference in *Drosophila melanogaster*:  
35 reply to Marais, Mouchiroud and Duret. *Genetical Research* **81**: 89-90.
- 36 Kurland, C. G. 1991 Codon bias and gene expression. *FEBS Letters* **285**: 165-169.
- 37 Lewontin, R. C. 2002 Directions in evolutionary biology. *Annu Rev Genet* **36**: 1-18.
- 38 Lipman, D. J., and W. J. Wilbur 1985 Interaction of silent and replacement changes in  
39 eukaryotic coding sequences. *J. Mol. Evol.* **21**: 161-167.
- 40 Marais, G., and B. Charlesworth 2003 Genome evolution: recombination speeds up adaptive  
41 evolution. *Curr. Biol.* **13**: R68-R70.
- 42 Marais, G., T. Domazet-Loaso, D. Tautz and B. Charlesworth 2004 Correlated evolution of  
43 synonymous and nonsynonymous sites in *Drosophila*. *J Mol Evol* **59**: 771-779.
- 44 Marais, G., and L. Duret 2001 Synonymous codon usage, accuracy of translation, and gene  
45 length in *Caenorhabditis elegans*. *J. Mol. Evol.* **52**: 275-280.
- 46 Marais, G., D. Mouchiroud and L. Duret 2001 Does recombination improve selection on  
47 codon usage? Lessons from nematode and fly complete genomes. *Proc Natl Acad Sci*  
48 *U S A* **98**: 5688-5692.
- 49 Marais, G., D. Mouchiroud and L. Duret 2003 Neutral effect of recombination on base  
50 composition in *Drosophila*. *Genetical Research* **81**: 79-87.

- 1 Marais, G., and G. Piganeau 2002 Hill-Robertson interference is a minor determinant of  
2 variation in codon bias across *Drosophila melanogaster* and *Caenorhabditis elegans*  
3 genome. *Mol. Biol. Evol.* **19**: 1399-1406.
- 4 Maynard Smith, J., and J. Haigh 1974 The hitch-hiking effect of a favourable gene. *Genet.*  
5 *Res.* **23**: 23-35.
- 6 McVean, G., and J. Vieira 2001 Inferring parameters of mutation, selection and demography  
7 from patterns of synonymous site evolution in *Drosophila*. *Genetics* **157**: 245-257.
- 8 McVean, G. A., and B. Charlesworth 2000 The effects of Hill-Robertson interference between  
9 weakly selected mutations on patterns of molecular evolution and variation. *Genetics*  
10 **155**: 929-944.
- 11 Meunier, J., and L. Duret 2004 Recombination drives the evolution of GC-content in the  
12 human genome. *Mol Biol Evol* **21**: 984-990.
- 13 Moriyama, E. N., and D. L. Hartl 1993 Codon usage bias and base composition of nuclear  
14 genes of *Drosophila*. *Genetics* **134**: 847-858.
- 15 Moriyama, E. N., and J. R. Powell 1997 Codon usage bias and tRNA abundance in  
16 *Drosophila*. *J. Mol. Evol.* **45**: 514-523.
- 17 Moriyama, E. N., and J. R. Powell 1998 Gene length and codon usage bias in *Drosophila*  
18 *melanogaster*, *Saccharomyces cerevisiae* and *Escherichia coli*. *Nucl. Acids Res.* **26**:  
19 3188-3193.
- 20 Morton, B. R. 1993 Chloroplast DNA codon use: evidence for selection at the psb A locus  
21 based on tRNA availability. *J Mol Evol* **37**: 273-280.
- 22 Munte, A., M. Aguade and C. Segarra 2001 Changes in the recombinational environment  
23 affect divergence in the yellow gene of *Drosophila*. *Mol Biol Evol* **18**: 1045-1056.
- 24 Muse, S. V. 1996 Estimating the synonymous and nonsynonymous substitution rates. *Mol*  
25 *Biol Evol* **13**: 105-114.
- 26 Nei, M. 1987 *Molecular evolutionary genetics*. Columbia University Press, New York.
- 27 Peden, J. F. 1999 Analysis of codon usage, pp. University of Nottingham, UK.
- 28 Powell, J. R., and E. N. Moriyama 1997 Evolution of codon usage bias in *Drosophila*. *Proc.*  
29 *Natl. Acad. Sci. USA* **94**: 7784-7790.
- 30 Ptak, S. E., D. A. Hinds, K. Koehler, B. Nickel, N. Patil *et al.* 2005 Fine-scale recombination  
31 patterns differ between chimpanzees and humans. *Nat. Genet.* **18**: [Epub ahead of  
32 print].
- 33 Rozas, J., and R. Rozas 1999 DnaSP version 3: an integrated program for molecular  
34 population genetics and molecular evolution analysis. *Bioinformatics* **15**: 174-175.
- 35 Schlenke, T. A., and D. J. Begun 2003 Natural selection drives drosophila immune system  
36 evolution. *Genetics* **164**: 1471-1480.
- 37 Sharp, P. M., M. Averof, A. T. Lloyd, G. Matassi and J. F. Peden 1995 DNA sequence  
38 evolution: the sounds of silence. *Philos Trans R Soc Lond B Biol Sci* **349**: 241-247.
- 39 Sharp, P. M., E. Bailes, R. J. Grocock, J. F. Peden and R. E. Sockett 2005 Variation in the  
40 strength of selected codon usage bias among bacteria. *Nucl. Acids Res.* **33**: 1141-  
41 1153.
- 42 Sharp, P. M., and W.-H. Li 1989 On the rate of DNA sequence evolution in *Drosophila*. *J.*  
43 *Mol. Evol.* **28**: 3398-3402.
- 44 Skibinski, D. O., and R. D. Ward 2004 Average allozyme heterozygosity in vertebrates  
45 correlates with Ka/Ks measured in the human-mouse lineage. *Mol Biol Evol* **21**:  
46 1753-1759.
- 47 Stephan, W. 1995 An improved method for estimating the rate of fixation of favourable  
48 mutations based on DNA polymorphism data. *Mol. Biol. Evol.* **12**: 959-962.
- 49 Stephan, W., and Y. Kim 2002 Recent applications of diffusion theory to population genetics,  
50 pp. 72-93 in *Modern developments in theoretical population genetics, the legacy of*

1 *Gustave Malecot*, edited by M. Slatkin and M. Veuille. Oxford University Press,  
2 Oxford.

3 Stephan, W., T. H. E. Wiehe and M. W. Lenz 1992 The effect of strongly selected  
4 substitutions on neutral polymorphism: analytical results based on diffusion theory.  
5 *Theor Pop Biol* **41**: 237-254.

6 Takano-Shimizu, T. 2001 Local changes in GC/AT substitution biases and in crossover  
7 frequencies on *Drosophila* chromosomes. *Mol Biol Evol* **18**: 606-619.

8 Veuille, M., E. Baudry, M. Cobb, N. Derome and E. Gravot 2004 Historicity and the  
9 population genetics of *Drosophila melanogaster* and *D. simulans*. *Genetica* **120**: 61-  
10 70.

11 Wall, J. D., P. Andolfatto and M. Przeworski 2002 Testing models of selection and  
12 demography in *Drosophila simulans*. *Genetics* **162**: 203-216.

13 Watterson, G. A. 1975 On the number of segregating sites in genetical models without  
14 recombination. *Theor. Pop. Biol.* **7**: 256-276.

15 Wright, F. 1990 The 'effective number of codons' used in a gene. *Gene* **87**: 23-29.

16 Wright, S. 1931 Evolution in Mendelian populations. *Genetics* **16**: 97-159.

17



**Table 1.** Spearman's correlation coefficients between codon bias ( $Fop-GC_i$ ) and DNA variation, and significance levels. Significant correlations are in bold. In brackets are the results obtained with the larger dataset, dataset 2, where possible.

	$Fop-GC_i$	$Dc_n$
<b>Non-synonymous</b>		
$Dc_n$	<b>-0.63<sup>***</sup></b>	
$\theta_{c_n}$ sim	<b>-0.71<sup>***</sup></b> ( <b>-0.51<sup>***</sup></b> )	<b>0.69<sup>***</sup></b> ( <b>0.62<sup>***</sup></b> )
$\theta_{c_n}$ mel	<b>-0.36<sup>*</sup></b>	<b>0.52<sup>***</sup></b>
<b>Synonymous</b>		
$Dc_s$	<b>-0.42<sup>**</sup></b>	<b>0.41<sup>**</sup></b>
$\theta_{c_s}$ sim	-0.27 <sup>NS</sup> (0.08 <sup>NS</sup> )	-0.11 <sup>NS</sup> (-0.09 <sup>NS</sup> )
$\theta_{c_s}$ mel	0.20 <sup>NS</sup>	0.20 <sup>NS</sup>
<b>Intron</b>		
$D_i$	-0.10 <sup>NS</sup>	0.00 <sup>NS</sup>
$\theta_i$ sim	0.02 <sup>NS</sup>	0.18 <sup>NS</sup>
$\theta_i$ mel	0.18 <sup>NS</sup>	-0.17 <sup>NS</sup>

<sup>NS</sup>: not significant; \* 0.01 < p < 0.05; \*\* : 0.001 < p < 0.01; \*\*\* : p < 0.001.

**Table 2.** Spearman's correlation coefficients between codon bias (*Fop*) and explicative parameters and significance levels (see the data section). Significant correlations that remained unchanged after a partial correlation analysis are in bold.

	<i>Fop</i>	GC <sub>nc</sub>	Ln(#ESTs)	Gene length
$\omega$	<b>-0.65<sup>***</sup></b>	-0.30 <sup>a</sup>	-0.07 <sup>NS</sup>	-0.01 <sup>NS</sup>
GC <sub>nc</sub>	<b>0.44<sup>***</sup></b>		0.12 <sup>NS</sup>	-0.01 <sup>NS</sup>
Ln(#ESTs)	<b>0.25<sup>**</sup></b>			-0.01 <sup>NS</sup>
Gene length	0.00 <sup>NS</sup>			

<sup>NS</sup>: not significant; \* 0.01 < p < 0.05; \*\* : 0.001 < p < 0.01; \*\*\* : p < 0.001;

<sup>a</sup> not significant after a partial correlation analysis ( $r_p = -0.05$ , n.s.).

## Figure legends

**Figure 1.** Comparison of the effect of recombination and the rate of amino-acid substitution on silent diversity and codon usage bias in *Drosophila melanogaster*. (A) Scatterplot of recombination ( $R_{mel}$ ) and diversity measured at synonymous and intronic sites ( $\theta_{i+s}$ ) in our dataset 1 (black dots) or silent diversity reported in Andolfatto and Przeworski (2001) (crosses). The solid and dotted lines are the best fit of Stephan's model (1995) to dataset 1 and data from Andolfatto and Przeworski (2001), respectively. (B) Scatterplot of recombination ( $R_{mel}$ ) and codon bias calibrated by intron GC-content ( $Fop-GC_i$ ). (C) Scatterplot of nonsynonymous substitution rate ( $Dc_n$ ) and diversity ( $\theta_{i+s}$ ). (D) Scatterplot of nonsynonymous substitution rate ( $Dc_n$ ) and codon bias ( $Fop-GC_i$ ).

**Figure 2.** Comparison of diversity between different classes of mutations within a species, and comparison of diversity between species for the same class of mutations. (A) Scatterplot of intron diversity ( $\theta_i$ ) and synonymous diversity ( $\theta_{c_s}$ ) in *D. simulans* (black dots) and *D. melanogaster* (open dots). (B) Scatterplot of diversity in *D. simulans* and *D. melanogaster* at synonymous sites (black dots) and intronic sites (open dots). (C) Scatterplot of intron diversity ( $\theta_i$ ) and nonsynonymous diversity ( $\theta_{c_n}$ ) in *D. simulans* (black dots) and *D. melanogaster* (open dots). (D) Scatterplot of diversity in *D. simulans* and *D. melanogaster* at nonsynonymous sites.

Figure 1

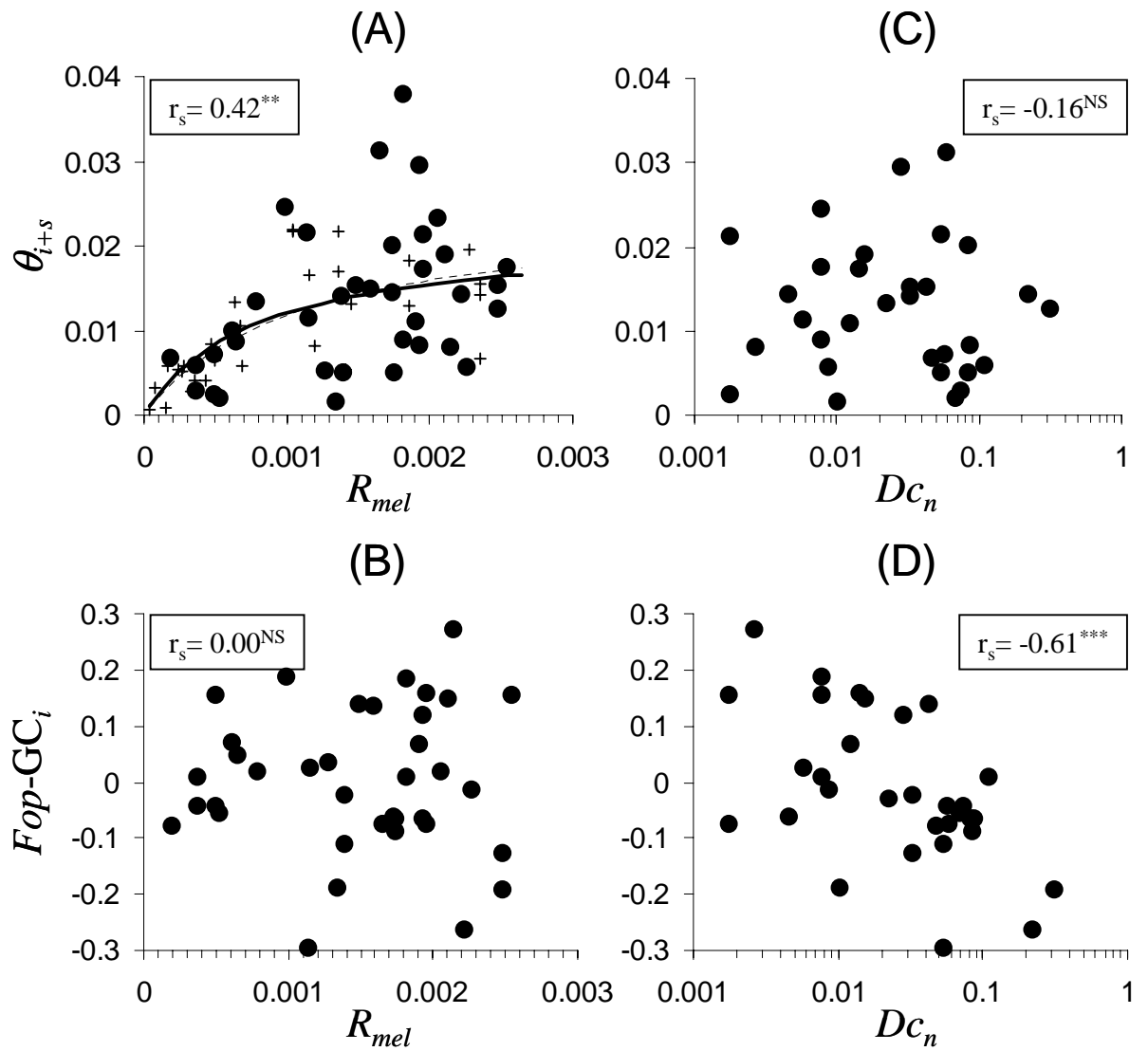


Figure 2

