

Evolution of compensatory substitutions through G·U intermediate state in *Drosophila* rRNA

(compensatory mutation models/rRNA phylogeny/deleterious substitutions/rates of evolution/expansion segments)

FRANÇOIS ROUSSET, MICHEL PÉLANDAKIS, AND MICHEL SOLIGNAC

Laboratoire de Biologie et Génétique Évolutives, Centre National de la Recherche Scientifique, 91198 Gif-sur-Yvette, France

Communicated by Francisco J. Ayala, August 12, 1991 (received for review March 13, 1991)

ABSTRACT It has often been suggested that the frequently observed Watson–Crick base-pair compensatory substitutions in RNA helical structures occur mainly through a slightly deleterious G·U intermediate state. We have scored base substitutions in a set of 82 related *Drosophila* species for the D1 and D2 variable domains of the large rRNA subunit. In all locations where a G-C ↔ A-U compensatory base change occurred, a G·U pair has been observed in one or several species. As this dominant process implies two transitions, their rate was far higher in paired regions (92%) than in unpaired regions (47%). The other types of compensation were rarer and no intermediate states were observed. Most of the G·U base pairs observed in a species are not slightly deleterious. The rate of evolution of compensatory substitution is close to that predicted by a simple model of compensatory substitution through slightly deleterious or slightly advantageous G·U pairs, although some exceptions are presented.

The secondary structure of rRNAs is remarkably uniform across taxa (1–3). This conservation is ensured by a special pattern of base change known as compensatory mutation (although what is observed is a compensatory substitution): when a substitution has occurred at a given site, the corresponding site, located vis-à-vis in the helical structure formed by the folding of the single RNA strand, also exhibits a change that restores the Watson–Crick base complementarity.

This observation is so general for the “stable” RNAs that the most efficient method used to confirm a secondary structure inferred from a sequence is based on the observation in various species of compensatory substitutions in the putative helices. Biochemical studies (4) or functional studies of double mutants (5, 6) have confirmed results obtained with the comparative method.

A simple model assumes that A-U* and G-C are optimal and stable states and that A-U ↔ G-C compensatory substitutions occur mainly through a slightly deleterious intermediate G·U state that is somewhat less stable but retains the helical structure. This slightly deleterious intermediate would be short-lived and, therefore, rarely observed. The low frequency of G·U in RNA sequences is generally explained in this way (7, 8). However, when only distantly related species are compared, as usual, there is no evidence that the G·U pairs effectively observed are deleterious or fugacious (9). In fact, some of these pairs may be deleterious states whereas others may be conserved over more or less prolonged times. The terms fugacious and stable states have a temporal significance and the time reference will be the average life span of a neutral pair (see *Results*).

The recognition of intermediate states *per se* must be achieved through comparison of numerous and related se-

quences that have evolved during a short time. This requires the study in related species of a region of the molecule variable enough to allow the observation of a sufficient number of substitutions. For this purpose, we have focused our study on the sequence of the divergent (or variable) domains D1 and D2 (10) of the 28S rRNA of 82 species of *Drosophila* and related drosophilids. This situation has several advantages. (i) The regions studied offer a high number of mutations in these two divergent domains. (ii) The phylogenetic analysis of the substitutions allows the events to be polarized. (iii) The fact that species are very related allows us to reconstruct the substitution pattern step by step, with a good approximation. This is because, in species clusters, the fixation of compensated stable states may not be achieved in all species and one or several sequences may still carry remnants of the deleterious states.

MATERIALS AND METHODS

Species. We have used the D1 and D2 sequences of the large rRNA subunits from 82 *Drosophila* and related species (Fig. 1), obtained in this laboratory. D1 and D2 domains of *Drosophila* (*Scaptodrosophila*) *dimorpha*, *Chymomyza bicolor*, and 26 cyclorrhaphous Diptera have also been sequenced, and 20 additional D2 sequences have been taken from the literature (11). Sequences of these 48 additional species have only been used to design secondary structures but they were not taken into account to establish the pattern of substitutions.

Sequences. The divergent domains D1 and D2 [sensu (10)] of the large rRNA subunit (28S), totaling about 545 nucleotides have been sequenced by the direct method, using the RNA as template (12). Total RNA was extracted by the lithium chloride method (13).

Sequence Alignment and Secondary Structure. Multiple alignment was carried out manually as well as automatically with the Clustal program (14). A very short segment (2–5 bases, positions 218–221 in Fig. 2) in a loop of the D2 domain has been excluded from the analysis because a high rate of insertion and deletion in these positions prevented safe reconstitution of substitution pathways. Secondary structures were determined using Zuker's algorithm (15). These structures have been confirmed or corrected through evidence of compensatory base changes from the sequences of the 130 species listed above.

Phylogeny and Pattern of Substitution. The general phylogeny of Drosophilidae has been analyzed by various parsimony and distance methods (M.P. and M.S., unpublished data). A more detailed analysis of the phylogeny of the

Abbreviation: Myr, million year(s).

*Standard IUB nomenclature for nucleic acids (33) is not used in this report. Rather, the following nomenclature is used. A hyphen is used for Watson–Crick base pairs in rRNA helices: G-C and A-U. A dot is used for G·U base pairs. No punctuation is used for unpaired nucleotides: A C, A A, and U U.

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked “advertisement” in accordance with 18 U.S.C. §1734 solely to indicate this fact.

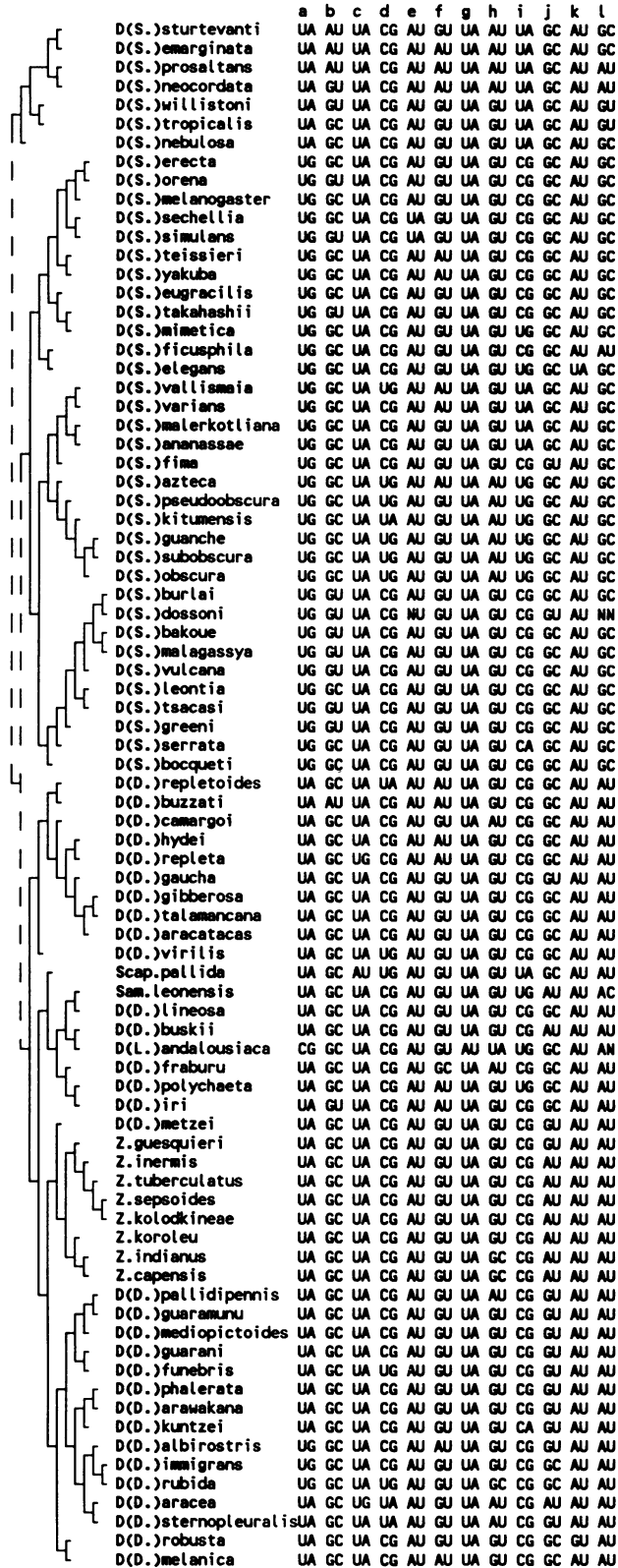


FIG. 1. Working phylogeny and sequences of the 12 compensating base pairs. The tree was obtained using the neighbor-joining method and rooted with *Leucophenga maculata*. Dotted lines correspond to the deepest branches of the phylogeny, which were not taken into account in the analysis. D., *Drosophila*; S., *Sophophora*; L., *Lordiphosa*; Scap., *Scaptomyza*; Sam., *Samoia*. Letters a to l at the top identify the base pairs on Fig. 2.

subgenus *Sophophora* is presented elsewhere (16). The tree has been rooted using outgroups progressively closer to the

set of species analyzed. A reliable and closely related outgroup, *Leucophenga maculata* (Drosophilidae, Steganinae), has thus been finally retained. The phylogenetic tree presently used (Fig. 1) has been established using the neighbor-joining distance method (17). The genera *Zaprionus*, *Scaptomyza*, and *Samoia* are clustered within the subgenus *Drosophila*, in agreement with Throckmorton (18).

We have plotted on the tree, for each site, the substitution pathways compatible with the general topology, retaining the most parsimonious one(s). Substitutions were polarized (e.g., A → U and U → A substitutions were distinguished). When alternative substitutions were possible, we did not attempt to choose between them. So, about 10% of substitutions were not included in the analysis (they did not appear to be of a particular type). Moreover, no substitutions were inferred for the deepest branches of the phylogeny (dotted lines in Fig. 1) because such data would have been unreliable. The variability has been estimated for each site by the minimum number of substitutions (including the 10% mentioned above not used for substitution nature) required to explain the distribution of nucleotide states across species.

RESULTS

Secondary Structure. The secondary structures of the D1 and D2 domains are presented in Fig. 2. These structures illustrate a specific type of interaction between bases, the one that usually emerges from compensatory mutation studies. They are neither a record of all possible interactions between bases nor do they show the folding of naked rRNA deduced from a specific sequence.

Time Scale. The total length of the tree, estimated by the sum of the length of its internal and terminal branches (the deepest ones have been discarded), has been converted into absolute time. We have used divergence time estimates available in the literature for some species pairs [subgenus *Drosophila*-subgenus *Sophophora*, 50 million years (Myr) (19); *Drosophila willistoni*-*Drosophila melanogaster*, 50 Myr (19); *Drosophila obscura*-*D. melanogaster*, 45 Myr (19, 20); *D. melanogaster*-*Drosophila orena*, 6 Myr (21)] and their corresponding patristic distances (data not shown) to extrapolate time calculations to the total set of species used. An estimate of 1-2 billion years was obtained for the whole tree. The most variable sites experienced 14 (3 sites) and 16 (1 site) substitutions within the unpaired regions. So, the corresponding evolutionary rate (10^{-8} per site per year if the total time is 1.5 billion years) of these rDNA sites is slightly under the *Drosophila* neutral rate equal to 1.8×10^{-8} (20).

Substitutions. Substitutions were inferred separately for nonpaired and paired positions and, for the latter, were also classified as a function of the vis-à-vis nucleotide. Their location and number are plotted on secondary structures in Fig. 2. The overall substitution rate is higher in unpaired regions (412 substitutions for 275 nucleotide positions) than in paired ones (144 substitutions for 133 pairs, i.e., 266 nucleotide positions). Sequences of pairs where compensation occurs are given in Fig. 1.

If we consider now the nature of substitutions (they were not determined for about 10% of the substitutions), we found for nonpaired regions, a total of 176 transitions and 201 transversions, 140 of the latter being A ↔ U. This bias persists in the relative rates, when the number of substitutions is weighted by the frequency of the nucleotide considered. By contrast, there are only 11 transversions for 123 transitions in paired regions (Fig. 3).

Evidence for Compensatory Substitutions and Intermediate States. The pathways of polarized substitutions from one type of pair to others inferred in helices are given in Fig. 3. By far, the most common pathway is G-C ↔ G-U ↔ A-U; although not oriented 5' → 3' in the figure, these compensations

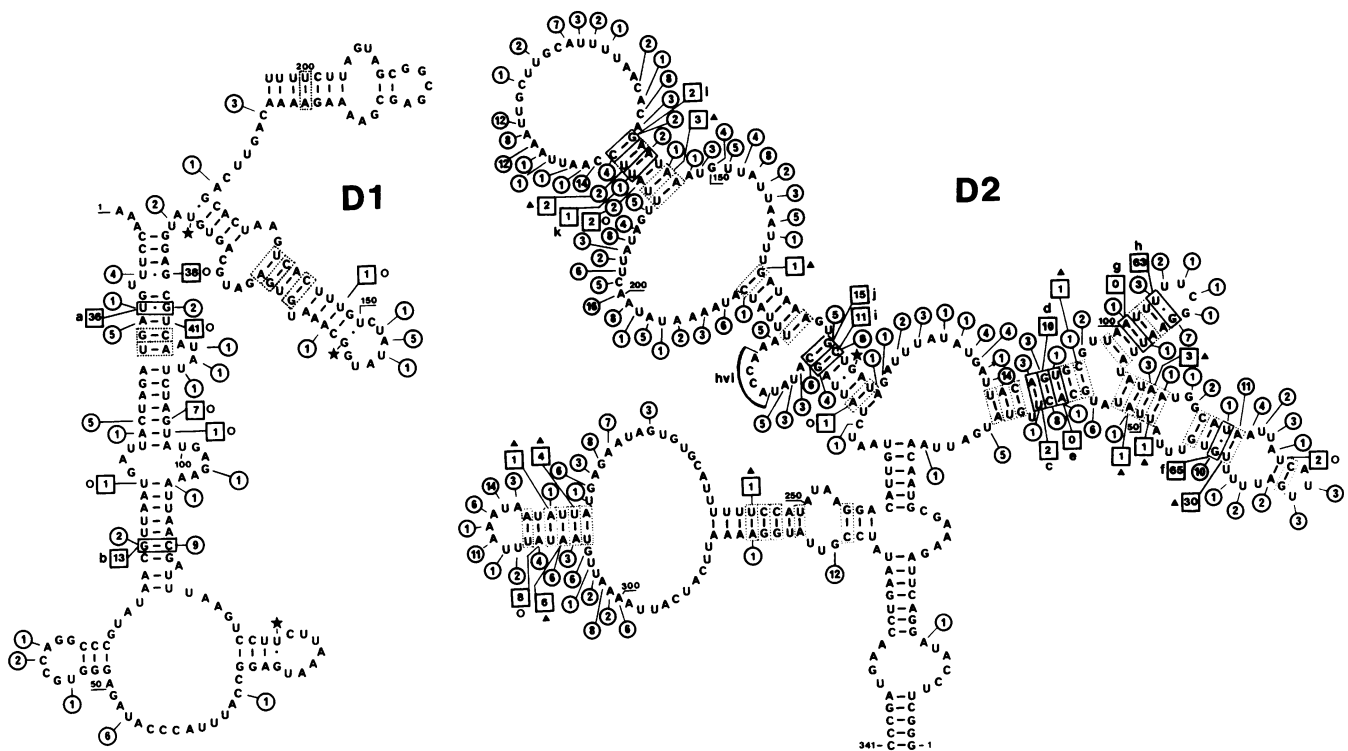


FIG. 2. Secondary structure and number of substitutions in the D1 and D2 regions. The sequence is that of *Drosophila melanogaster*, and positions are numbered from 3'(RNA) to 5'. Substitutions have not been counted in the hypervariable loop (hvl). Symbols: circled numbers, number of substitutions inferred at the specified site; boxed numbers, number of G-U pairs observed at the specified base pair in the 82 analyzed drosophilids; boxed pairs, pairs where compensation was observed within the drosophilids; dotted boxes, pairs of sites for which compensatory substitutions were observed only in other dipteran species. Base pairs from categories A to E of Table 1 are represented by the following symbols: B, Δ ; C, \circ ; D, $*$; A and E, indexed by letters a to l as in Fig. 1.

always appeared as the result of two transitions, $G \leftrightarrow A$ and $C \leftrightarrow U$, never of two transversions, $G \leftrightarrow U$ and $C \leftrightarrow A$. The distribution of G-U among base pairs and among species is detailed in Table 1. The potential intermediate pair C-A, which also allows compensation through two transitions, was very scarce in the sequences and was observed only three times (Fig. 1).

The only other compensation observed was $A-U \leftrightarrow U-A$ (four observations, Fig. 3) but no intermediate states (A A or U U) are known in these positions.

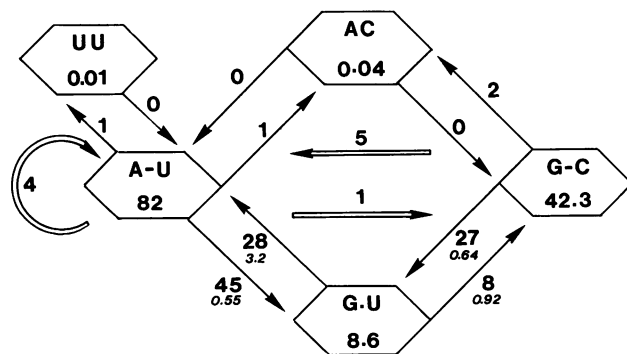


FIG. 3. Substitution pathways in paired regions. In the hexagonal boxes are indicated the average numbers of the types of pairs observed in the sequences of the 82 *Drosophila* species over 133 positions. The total number of polarized substitutions observed from one pair to another is indicated on the arrows; the relative substitution rate (number of substitutions divided by the number of pairs) is indicated for the four higher rates by smaller italic numbers. Open-shafted arrows correspond to the pathways requiring two substitutions. In addition to the ways indicated, a change, $A C \rightarrow U-A$ (two transversions, no intermediate), was observed once.

A Substitution Model. The more deleterious a G-U is, the rarer it is because, when generated by a mutation, it rarely reaches fixation and if it does, it is rapidly replaced by an advantageous mutant. So, we may expect some relationship between the frequency of G-U among species and the substitution rate. The relation between these two parameters can be established from a Markovian model—i.e., by assuming that probabilities of change between different pair states are constant over evolutionary time, for the main pathway of compensatory substitution:

$$A-U \xrightleftharpoons[r_2]{r_1} G-U \xrightleftharpoons[r_4]{r_3} G-C$$

where the r_i values are the substitution rates (all substitutions in this scheme are transitions) between corresponding pairs. These substitution rates are defined as the product of the mutation rate (in fact transition rate u) by the population size $2N$ and the probability of fixation of a mutant $k(N, s) = 2s/[1$

Table 1. Distribution of G-U pairs among base pairs and species

Type of position	No. of positions	G-U in indicated type of position, average no. per species
A	8	2.61
B	12	0.66
C	10	1.24
D	4	4
E	4	0.04
Total	38	8.55

A, positions where $A-U \leftrightarrow G-C$ substitutions occur in the analyzed drosophilids; B, positions where $A-U \leftrightarrow G-C$ substitutions occur among other dipteran species; C, variable positions where no compensatory substitutions are known; D, G-U present in all analyzed drosophilids; E, position where only $A-U \leftrightarrow U-A$ compensatory substitutions are observed among the analyzed drosophilids.

– $\exp(-4Ns)$], dependent on the selective pressure s and on N (22). So $r_i = 2Nuk(N, s)$. From basic Markov chain theory (23), the equilibrium frequency of G-U at a base pair (f_{G-U}) in the sequences is:

$$f_{G-U} = [1 + (r_2/r_1) + (r_3/r_4)]^{-1} \quad [1]$$

Homozygotes for A-U and G-C are considered to be selectively equivalent, homozygotes for G-U have a fitness of $1 + 2s$ and we suppose semidominance. If we assume that all transition rates are identical, then

$$f_{G-U} = [1 + 2(1 - e^{-4Ns})/(e^{4Ns} - 1)]^{-1} \\ = (1 - 1/X)/(2X - 1/X - 1), \quad [2]$$

where $X = e^{-4Ns}$. Conversely,

$$X = (1 - f_{G-U})/(2f_{G-U}). \quad [3]$$

So Ns can easily be estimated from f_{G-U} .

For a given paired position in helices, the average substitution rate, denoted Sr and corresponding to the sum for the two bases of the pair, is the sum of the r_i weighted by the frequency of the starting state. If equilibrium is assumed, then its value in such a three-state model is necessarily twice the rate of substitution of the obligate intermediate state toward the two others:

$$Sr = 2f_{G-U}(r_2 + r_3) = 4f_{G-U}[2Nu2s/(1 - e^{-4Ns})] \\ = 4uf_{G-U}(1 - f_{G-U})\ln[1 - f_{G-U}]/2f_{G-U}/(1 - 3f_{G-U}) \\ \text{(from Eq. 3).} \quad [4]$$

This gives the shape of the theoretical curve (Fig. 4). To draw this curve, we need the total time (1.5 billion years, see above) and the transition rate u . The value of u has been

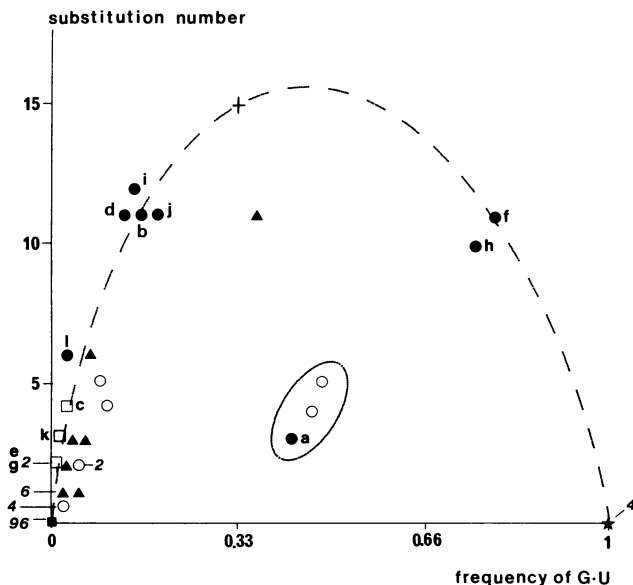


FIG. 4. Relationship between the observed frequency of G-U at a given pair position in helices among species and the inferred number of substitutions for the two paired sites, as compared to the Markovian model (see Results). Others symbols are as in Fig. 2, except ●, positions where A-U ↔ G-C substitutions occur; ■, constant Watson-Crick base pairs; and □, positions where A-U ↔ U-A, but no A-U ↔ G-C, substitutions occur (attributable to two transversions, given only for indication). Italic numbers are the number of points. Letters a to l are as in Figs. 1 and 2. The curve was drawn as described in Results. +, Neutrality point; *, category D in Table 1.

estimated from the highest synonymous substitution rate K_S of Sharp and Li (20) (1.8×10^{-8}), the approximation $K_S = 3(u + v)/2$ [from equation 7 of Li *et al.* (24)], and $u/(2v) = 176/201$ as inferred from unpaired positions of the D1 and D2 domains, where v is the transversion rate. From these relationships, $u = 0.75 \times 10^{-8}$. In this case, a neutral transition occurs every $1/u$ and the life span of a neutral G-U will be half this time, $1/(2u) = 66.5$ Myr (a transition on either nucleotide restoring a Watson-Crick pair). This time can be taken as an objective boundary between fugacious (deleterious) and stable (advantageous) states.

In such a model, if $f_{G-U} = 1/6$, then s is approximately equal to $-0.9/(4N)$ (from Eq. 3), and the mean life span of G-U pairs will be $[2u2Nk(N, -s)]^{-1} = 43.5$ Myr. This is also the time during which the base pair was in the state G-U, divided by half the number of substitutions: so it is inversely proportional to the slope of the line joining the origin to the corresponding point on the curve. The value of N is unknown but s is obviously very small.

Points for individual base pairs are plotted on the same figure, their coordinates being observed G-U frequency and inferred substitution numbers. The curve closely approximates the points for the most variable positions.

DISCUSSION

Reliability of the Substitution Record. The secondary structure of the D1 domain is the one found for the mouse (10) and for various eukaryotes, including *D. melanogaster* (25), by thermodynamic, comparative, and chemical analyses. The general outline of the D2 domain, with the relative position of its three arms, is similar to the one established for several eukaryotes by Michot and Bachellerie (26). The comparison of 130 dipteran sequences rendered the secondary structures established unambiguous.

The estimations of number of substitutions appear robust with regard to other possible phylogenies. The record of substitutions would not be significantly modified if, in some branches, an alternative topology had been used: the substitutions inferred often correspond to nucleotides singular to a given species or common to a group of species; whatever the real local topology, the nucleotide change inferred would have been the same. This can be checked for positions where compensatory substitutions occur from Fig. 1.

The high density of closely related species used allows us to detect most of the intermediate states. However, by application of the parsimony criterion, there is a risk of bias in the substitution record ensuing from the fact that undetected substitutions belong preferentially to the types occurring faster (positively selected). This may explain why the number of substitutions from G-U to Watson-Crick is apparently lower than the reverse (Fig. 3).

The Multiple Nature of G-U States. The frequency of G-U pairs is low in the sequences (8.6/133). This is expected if G-U pairs are deleterious. However, when a G-U pair is conserved in all species, it is clear evidence of selection in favor of G-U in that position. We have observed four such pairs in our sequences (three in D1 and one in D2, this last one being well conserved in other Diptera). Such situations are also known elsewhere: in 16S rRNA (27), 5S rRNA (28), 23S rRNA (29, 30), self-splicing introns (31), and many other RNAs.

Positions where compensation occurs correspond to pairs with high G-U frequency as well as pairs with low G-U frequency. In some positions, G-U pairs are more frequent than expected under neutrality (Fig. 4) and thus these long-lived pairs are, at least most of the time, at selective advantage. Only in the remaining base pairs, where G-U pairs are rare and thus contribute only a small part of the G-U pairs

encountered in a sequence, may these G-U pairs be effectively deleterious.

Relation Between G-U Frequency and Nucleotide Variability.

The overall frequency of G-U in sequences (8.6/133) is not appropriate information to elucidate the constraints on the evolution of compensatory substitutions, since the nature of G-U pairs is variable from position to position. More generally, it is unsafe to conclude from data pooled from various positions under different constraints (such as Fig. 3). It is better to investigate each base pair separately, as in Fig. 4. Although the frequency of G-U pairs is an observation from sequences that are not independent, the comparison of data to expectations from the Markovian model yields some clear results.

The rarest G-U pairs are found in sites where the substitution rate is lowest. An exception is the isolated group of three points (circled in Fig. 4). It corresponds to three sites located near the stem of the D1 domain that exhibit parallel patterns of substitution: two G-U pairs were conserved in the *Sophophora* subgenus and were variable in the other species. The reverse is true for the third pair. This observation is best explained if the selective value of these G-U pairs was changed in the course of evolution within the genus *Drosophila*, perhaps through the exchange of the positions where G-U occur. Their intermediate position on Fig. 4 should reflect this change.

The absence of compensation in some positions where the frequency of G-U is relatively high and substitutions are numerous suggests that in these positions only one paired state is admissible within the set of drosophilid species considered. The large excess of G-U ↔ A-U substitutions (Fig. 3) comes from these positions. It is however remarkable that compensatory substitutions were observed for most of these positions in the D2 domain of some distantly related cyclorrhaphous Diptera. Some change of selective constraints may also be involved here.

The few compensatory substitutions of the type A-U ↔ U-A we have observed, always without intermediate states, may have evolved through the compensation pathway suggested by Kimura (32): if the intermediate state is fairly deleterious, it does not reach fixation but it can be maintained at low frequency by the rather high A ↔ U transversion rate, "waiting" for a compensatory mutation, the latter then being fixed.

A direct extrapolation of the results obtained from rapidly evolving regions in related species to slowly evolving regions in distant species is assuredly incorrect. The D2 domain is highly variable; in more strongly constrained regions, G-U pairs may be more deleterious, but many G-U pairs may also be advantageous. The empirical foundations of the current model for compensatory substitution (7, 8) were based upon the occasional observation of G-U pairs in tRNAs. Paradoxically, most are probably not deleterious states, but the model of compensation through these pairs is probably yet correct.

We thank Dominique Vautrin for excellent technical assistance, the Bowling Green University for providing strains, Henriette Rutt-

kay for providing sequences of the group *obscura*, and Y. d'Aubenton Carafa, J.-M. Cornuet, and J.-C. Mounolou for comments on the manuscript.

1. Brimacombe, R. (1984) *Trends Biochem. Sci.* **5**, 273-277.
2. Woese, C. R. (1987) *Microbiol. Rev.* **51**, 221-271.
3. Raué, H. A., Klootwijk, J. & Musters, W. (1988) *Prog. Biophys. Mol. Biol.* **51**, 77-129.
4. Moazed, D., Stern, S. & Noller, H. F. (1986) *J. Mol. Biol.* **187**, 399-416.
5. Ozeki, H., Inokuchi, H., Yamao, F., Kodaira, M., Sakano, H., Ikemura, T. & Shimura, Y. (1980) in *tRNA: Biological Aspects*, eds. Söll, D., Abelson, J. & Schimmel, P. (Cold Spring Harbor Lab., Cold Spring Harbor, NY), pp. 341-362.
6. Waring, R. B., Towner, P., Minter, S. J. & Davies, R. W. (1986) *Nature (London)* **321**, 133-139.
7. Holmquist, R., Jukes, T. H. & Pangburn, S. (1973) *J. Mol. Biol.* **78**, 91-116.
8. Ohta, T. (1974) *Nature (London)* **252**, 351-354.
9. Hancock, J. M., Tautz, D. & Dover, G. A. (1988) *Mol. Biol. Evol.* **5**, 393-414.
10. Hassouna, N., Michot, B. & Bachellerie, J. P. (1984) *Nucleic Acids Res.* **12**, 4259-4279.
11. Vossbrinck, C. R. & Friedman, S. (1989) *Syst. Entomol.* **14**, 417-431.
12. Qu, L.-H., Michot, B. & Bachellerie, J.-P. (1983) *Nucleic Acids Res.* **11**, 5903-5920.
13. Maccacchini, M. L., Rudin, Y., Blobel, G. & Schatz, G. (1979) *Proc. Natl. Acad. Sci. USA* **76**, 343-347.
14. Higgins, D. G. & Sharp, P. M. (1989) *Cabios* **5**, 151-153.
15. Zuker, M. & Stiegler, P. (1981) *Nucleic Acids Res.* **9**, 133-148.
16. Pélandakis, M., Higgins, D. G. & Solignac, M. (1991) *Genetica* **84**, 87-94.
17. Saitou, N. & Nei, M. (1987) *Mol. Biol. Evol.* **4**, 406-425.
18. Throckmorton, L. H. (1975) in *Handbook of Genetics*, ed. King, R. (Plenum, New York), Vol. 3, pp. 421-469.
19. Beverley, S. M. & Wilson, A. C. (1984) *J. Mol. Evol.* **21**, 1-13.
20. Sharp, P. M. & Li, W.-H. (1989) *J. Mol. Evol.* **28**, 398-402.
21. Lachaise, D., Cariou, M.-L., David, J. R., Lemeunier, F., Tsacas, L. & Ashburner, M. (1988) *Evol. Biol.* **22**, 159-225.
22. Kimura, M. (1962) *Genetics* **47**, 713-719.
23. Feller, W. (1968) *An Introduction to Probability Theory and Its Applications* (Wiley, New York), 3rd Ed.
24. Li, W.-H., Wu, C.-H. & Luo, C.-C. (1985) *Mol. Biol. Evol.* **2**, 150-174.
25. Qu, L. H. (1986) Dissertation (Univ. of Toulouse, France).
26. Michot, B. & Bachellerie, J.-P. (1987) *Biochimie* **69**, 11-23.
27. Woese, C. R., Gutell, R., Gupta, R. & Noller, H. F. (1983) *Microbiol. Rev.* **47**, 621-669.
28. Westhof, E., Romby, P., Romaniuk, P. J., Ebel, J.-P., Ehresmann, C. & Ehresmann, B. (1989) *J. Mol. Biol.* **207**, 417-431.
29. Egebjerg, J., Douthwaite, S. R., Liljas, A. & Garrett, R. A. (1990) *J. Mol. Biol.* **213**, 275-288.
30. Michot, B., Qu, L.-H. & Bachellerie, J.-P. (1990) *Eur. J. Biochem.* **188**, 219-229.
31. Inoue, T., Sullivan, F. X. & Cech, T. R. (1986) *J. Mol. Biol.* **189**, 143-165.
32. Kimura, M. (1985) *J. Genet.* **64**, 7-19.
33. IUPAC-IUB Commission on Biochemical Nomenclature (CBN) (1970) *Eur. J. Biochem.* **15**, 203-208.