

Mammalian Gene Evolution: Nucleotide Sequence Divergence Between Mouse and Rat

Kenneth H. Wolfe, Paul M. Sharp

Department of Genetics, University of Dublin, Trinity College, Dublin 2, Ireland

Received: 14 July 1992 / Revised: 15 September 1992

Abstract. As a paradigm of mammalian gene evolution, the nature and extent of DNA sequence divergence between homologous protein-coding genes from mouse and rat have been investigated. The data set examined includes 363 genes totalling 411 kilobases, making this by far the largest comparison conducted between a single pair of species. Mouse and rat genes are on average 93.4% identical in nucleotide sequence and 93.9% identical in amino acid sequence. Individual genes vary substantially in the extent of nonsynonymous nucleotide substitution, as expected from protein evolution studies; here the variation is characterized. The extent of synonymous (or silent) substitution also varies considerably among genes, though the coefficient of variation is about four times smaller than for nonsynonymous substitutions. A small number of genes mapped to the X-chromosome have a slower rate of molecular evolution than average, as predicted if molecular evolution is "male-driven." Base composition at silent sites varies from 33% to 95% G + C in different genes; mouse and rat homologues differ on average by only 1.7% in silent-site G + C, but it is shown that this is not necessarily due to any selective constraint on their base composition. Synonymous substitution rates and silent site base composition appear to be related (genes at intermediate G + C have on average higher rates), but the relationship is not as strong as in our earlier analyses. Rates of synonymous and nonsynonymous substitution are correlated, apparently because of an excess of substitutions involv-

ing adjacent pairs of nucleotides. Several factors suggest that synonymous codon usage in rodent genes is not subject to selection.

Key words: Molecular clocks — Rodents — Genome evolution — G + C content — Codon usage — Dinucleotide mutation effects

As increasingly large amounts of DNA sequence data accumulate, our understanding of the pattern and dynamics of gene sequence divergence is growing clearer. The molecular clock, first postulated by Zuckerkandl and Pauling (1962) and later championed by Allan Wilson (see, for example, Wilson et al. 1977, 1987), is a description of the observed regularity with which substitutions accrue in nucleotide and amino acid sequences. Under the neutral theory (Kimura 1983) the rate of molecular evolution (k) of a sequence is determined by the product of the mutation rate (u_T) and the fraction of sites in the sequence (or the average fraction of mutations at those sites) that are selectively neutral (f_0):

$$k = u_T f_0$$

(Kimura 1977). Thus, the rate of sequence change in different genes or proteins will be similar only if these two factors are the same.

From the first studies of protein evolution it was apparent that the rate of amino acid replacement varies enormously among proteins (see, e.g., Dickerson 1971) and that this can be attributed to differences in the fraction of possible neutral amino acid replacements (f_0) between, say, fibrinopeptides on

the one hand and histone H4 on the other. (See Table 4.4 of Nei 1987.) While these extreme examples of fast and slow protein sequence evolution have now become the stuff of textbooks, few data have been gathered on the distribution of rates among the less-exceptional proteins. An understanding of the shape of the underlying distribution of rates would contribute toward an understanding of protein evolution in general. However, documentation of rates of evolution of a large number of proteins has been hampered in part by the necessity to combine data from sequence comparisons across different taxonomic groups, thus requiring the use of fossil-based divergence dates along with the ensuing uncertainties (e.g., Li et al. 1985).

Here we overcome this problem by limiting our comparisons to a single pair of species (mouse and rat) so that the relative extents of change seen in different proteins are a direct measure of their relative rates of evolution, regardless of the data of divergence between the two species. (In fact, the actual date of the mouse-rat speciation event is a matter of considerable controversy—see Catzeflis et al. 1992—and will not be discussed in much detail.) The massive popularity of DNA sequencing as a biochemical tool, and of both mouse and rat as model organisms, has resulted in the sequencing of several hundred homologous gene pairs from these species in the decade following the first comparative study (Jagodzinski et al. 1981). By comparing nucleotide as well as amino acid sequences we have also been able to characterize the variation among genes in the rate of silent (synonymous) nucleotide substitution and show that, as with proteins, there are some genes with exceptionally fast or exceptionally slow rates. Since almost all silent sites in mammalian genes are likely to be free to accept nucleotide substitutions, these deviations from a “silent molecular clock” may be due to local variation in the mutation rate (μ_T); the possible basis for this variation is discussed. A seemingly paradoxical correlation between the rates of synonymous and nonsynonymous substitution in genes is also investigated.

Data and Methods

Sequences were collected between 1986 and 1992 and were either taken from releases 59–71 of GenBank or obtained from the EMBL and GenBank database electronic mail servers; in a few cases sequences were obtained directly from the literature. Genomic sequences were used in preference to cDNAs where possible at the time of collection. Homologous mouse and rat sequence pairs in GenBank were identified initially by browsing through the sequence definitions with the aid of a text editor or by keyword searches using the retrieval system ACNUC (Gouy et al. 1985). Other sequence pairs were identified by a semiautomated procedure whereby data on polypeptide length and

amino acid composition were extracted from each of the 5,568 mouse and 2,543 rat protein-coding sequences in GenBank (release 71) using ACNUC, after which each mouse sequence was compared to each rat sequence and candidate matches of similar length and amino acid composition were flagged for further investigation. This method has the potential to identify all homologous pairs very rapidly (<30 min on a shared VAX computer to make all 14 million possible pairwise comparisons of murid genes) but is biased in favor of highly conserved sequences and cannot find matches between incomplete sequences and full-length homologues. In practice, this method found only about 70% of the pairs that had already been identified manually, but also led to the discovery of 35 additional pairs of homologues that were not evident from the annotation of their database entries.

Coding sequence pairs were then extracted from GenBank using ACNUC, and the deduced protein sequences were aligned using a rapid Needleman-Wunsch method (CLUSTAL; Higgins and Sharp 1989). The DNA sequences were then aligned using the protein alignments as templates. All alignments were inspected by eye for possible frameshift sequencing errors and a number of putative errors were corrected in such a way as to maximize the ratio of synonymous to nonsynonymous nucleotide substitutions in the frameshifted region. Three genes were found where the reported mouse and rat sequences were essentially identical; we suspect that the source species is incorrectly identified in database entries M84361 (“rat” CSF-1), X61479 (“rat” *c-fms*), and M36660 (“mouse” NADPH: menadione oxidoreductase). We assume that these result from laboratory or database confusion and have not included these sequences in the analysis.

The extents of synonymous and nonsynonymous divergence (in terms of nucleotide substitutions per site) were calculated by the method of Li et al. (1985). In this method, sites in codons are classified as 0-fold, 2-fold, or 4-fold degenerate, according to how many alternative nucleotides at the site encode the same amino acid. A correction for superimposed substitutions at single sites is made by the two-parameter method of Kimura (1980), which allows for different rates of transitions and transversions. The numbers of synonymous and nonsynonymous substitutions per site are then calculated as appropriately weighted averages of these values. (See Li et al. 1985 for more details, including the estimation of standard errors.)

A few sequence pairs were found to give unusually high K_S values (Fig. 2), so we considered the possibility that these might be paralogous. In the case of β_2 -microglobulin ($K_S = 0.69$) the mouse and rat genes have each been sequenced by at least two independent groups. Multiple independent sequences have also been determined for rat SVS IV ($K_S = 0.78$). Mouse SVS IV has been sequenced by only one laboratory (Chen et al. 1987), but this group obtained the mouse cDNA by using a sequenced rat SVS IV cDNA probe and also directly sequenced most of the mouse protein. Furthermore, Dietrich et al. (1992) were able to map SVS IV (*svp-4*) to mouse chromosome 2 by polymerase chain reaction amplification using primers corresponding to the published mouse sequence. There is thus no reason to doubt that the mouse and rat SVS IV sequences are orthologues. The gene with the third-greatest K_S value (0.51) is SPOT-1, which is discussed by Dickinson et al. (1989).

Nucleotide and Amino Acid Sequence Divergence Between Mouse and Rat

The nucleotide sequences of the coding regions of 363 genes for which DNA sequences have been de-

Table 1. Nucleotide and amino acid sequence divergence, and silent-site G+C content, in 363 homologous genes of mouse and rat^a

	Mean (\pm SD)	Range	L ^b
Amino acid identity (%)	93.9 (\pm 8.1)	56.0–100.0	136,729
Nucleotide identity (%)	93.4 (\pm 4.1)	69.6–99.0	411,300
Nonsynonymous substitutions (K_A)	0.032 (\pm 0.049)	0.000–0.372	318,873
Synonymous substitutions (K_S)	0.224 (\pm 0.084)	0.041–0.780	91,315
Silent site G+C content (%)	62.0 (\pm 11.1)	32.8–95.4	145,357

^a The 363 genes, and their individual values, are detailed in Appendix 1. Means and standard deviations are weighted by the number of sites in each gene. Silent site G+C content (GC_S) is the G+C content at 2-fold and 4-fold degenerate sites in codons.

^b The total number of residues in each category

terminated from both mouse and rat were aligned and compared. The genes, and some statistics describing their molecular evolution, are listed in Appendix 1. The results are summarised in Table 1. The total length of aligned sequences (excluding gaps) is 411 kb. The genes range in size from 135 bp (thymosin β 4) to 8,247 bp (inositol-1,4,5-triphosphate receptor), with an average length of 1,133 bp. As expected, the deduced protein sequences range in degree of divergence. Twenty-five proteins (7%), including some actins, ion-channel proteins, and ribosomal proteins, are identical in mouse and rat. The most divergent proteins are the salivary SPOT-1 protein (56% identity; discussed by Dickinson et al. 1989) and interleukin-3 (59%; see Cohen et al. 1986). Nucleotide sequence identity varies from 70% in the SPOT-1 gene to 99% in the gene encoding the Y-box binding transcription factor. The mean level of sequence identity, weighted by gene length, is 93.4% (standard deviation 4.1%) for nucleotides and 93.9% (\pm 8.1%) for amino acids.

The extents of amino acid and nucleotide sequence divergence are (necessarily) correlated, but the relationship between these two measures is interesting. In genes encoding very highly conserved proteins, amino acid sequence identity exceeds nucleotide identity (Fig. 1), because silent (synonymous) nucleotide substitutions are permitted in these genes. However, in genes encoding less-conserved proteins, nucleotide similarity exceeds protein similarity. This presumably reflects the degeneracy of the genetic code: for example, a single codon with nucleotide substitutions at positions 2 and 3 (the former causing an amino acid replacement, the latter probably silent) exhibits 33% nucleotide sequence identity but 0% amino acid sequence identity. For the mouse-vs-rat data, the point at which nucleotide and amino acid sequence identity levels cross (i.e., are similar) is approximately 93% (Fig. 1b); the mean values happen to be close to this point. This picture is markedly different from that seen in a study of bacterial sequence divergence (67 genes compared between *Escherichia coli* and *Salmonella typhimurium*; Sharp 1991). In the bacterial data, protein sequences are invariably more similar

than nucleotide sequences over a range of 76 to 100% amino acid sequence identity, and extrapolation of the relationship suggests a crossover point at around 65% identity (P.M.S., unpublished results). This difference is probably related to a much higher mean ratio of synonymous-to-nonsynonymous divergence in the bacterial genes than in the rodent genes, which is discussed below.

Nucleotide substitutions between two species can be classified as either nonsynonymous (replacement) or synonymous (silent), depending on whether or not they alter the protein sequence. The estimated numbers of nonsynonymous substitutions per site (K_A) and synonymous substitutions per site (K_S) were calculated by the method of Li et al. (1985), which corrects for multiple hits (Table 1, Appendix 1). These values, and the relationship between them, will be discussed in turn.

Nonsynonymous Nucleotide Substitution and Protein Evolution

The estimates of the extents of nonsynonymous nucleotide substitution (K_A) among genes reflect the diversity of levels of protein sequence conservation and range from zero in the 25 genes mentioned earlier (and $K_A < 0.01$ substitution per site in a further 103, or 28% of genes) to $K_A = 0.372$ (\pm 0.053) nonsynonymous substitutions per site in the salivary SPOT-1 gene; the weighted mean K_A is 0.032 (with standard deviation 0.049).

The K_A values for the 363 genes (Fig. 2a) form a broad, largely unskewed, distribution when plotted on a histogram with a semilogarithmic scale (as originally used by Ochman and Wilson 1987). There is, however, a distinct peak due to a number of genes with extremely low values of K_A (i.e., encoding almost invariant proteins). This result was also apparent in the small number of genes analyzed by Li et al. (1985; see histograms in Ochman and Wilson 1987; and in Hartl and Clark 1989:369). A similar result is obtained in comparisons of sequences (about 700 genes) between humans and rodents (K.H.W., unpublished results).

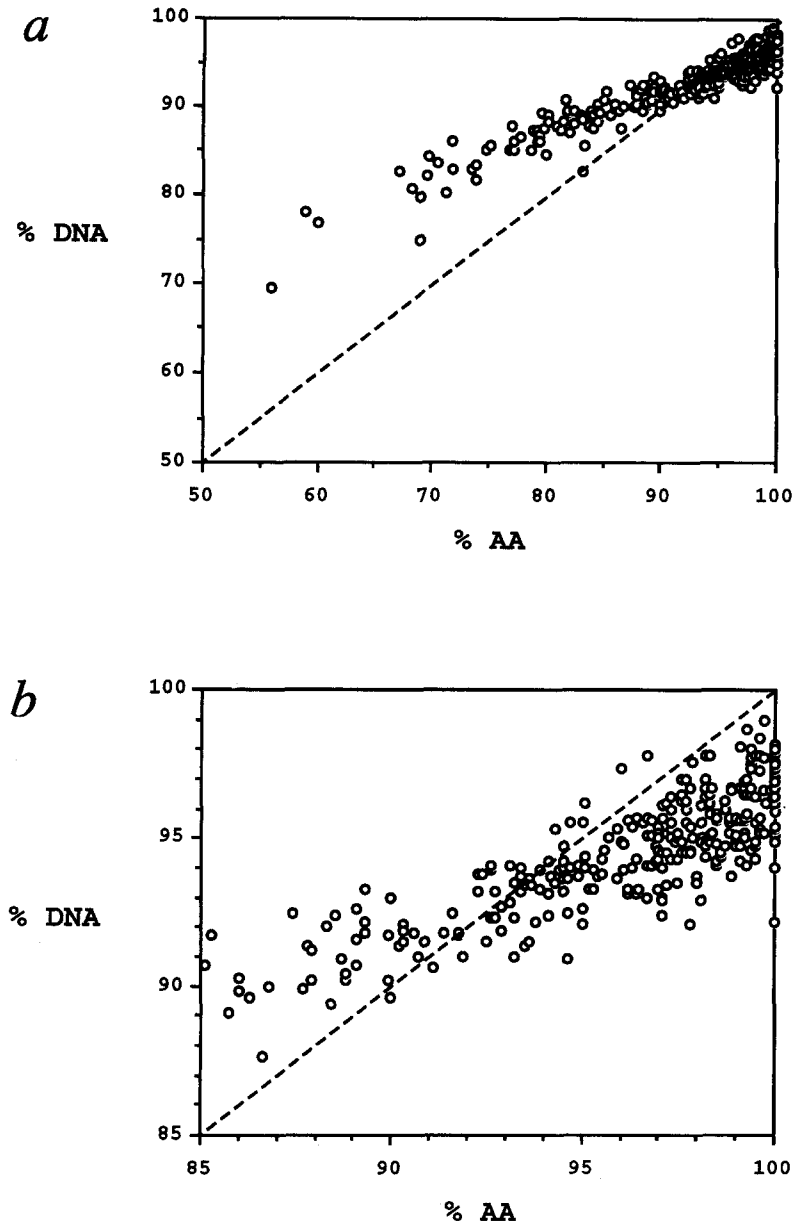


Fig. 1. Relationship between nucleotide (% DNA) and amino acid (% AA) sequence identity among 363 genes compared between mouse and rat. The region between 85% and 100% identity is enlarged in b. The dashed lines indicate equal levels of amino acid and nucleotide sequence identity.

The extent of synonymous substitution also varies between genes, although to a much lesser extent. (See below.) Setting that variation aside for the moment, it is interesting to note that the ratio between the mean synonymous and mean nonsynonymous divergence in the rodent genes is 7.1. This value is a little higher than that reported earlier (about five, among 35 genes; Li et al. 1985) in comparisons among mammalian orders. More significantly, this value is much lower than that seen in bacterial genes (for which the ratio is 24; Sharp 1991). Ochman and Wilson (1987; see also Lawrence et al. 1991) have proposed that this difference is due, at least in part, to the enormous difference in effective population size between bacteria (specifically *E. coli*) and mammals: selection against slightly deleterious mutant protein sequences may be far more efficient in bacteria. The

effect may also be partly due to genes of different functions being included in the two data sets.

If synonymous substitutions are essentially neutral (as discussed below), the ratio of K_S to K_A provides a measure of the degree of selective constraint on a protein sequence, independent of other factors (such as local mutation rate differences) that may cause the underlying nucleotide substitution rate to be different in different genes. Values of the K_S/K_A ratio for individual genes vary from infinity (for the set of 25 genes encoding invariant proteins) to 1.22 for the Blast-1 antigen and 0.89 for interleukin-3. The last value may be remarkable because the nonsynonymous rate is expected to exceed the synonymous rate only in exceptional cases of positive natural selection for amino acid sequence divergence. (See, e.g., Hughes and Nei 1988.) However, interleukin-3 is the only gene (out of 363) for

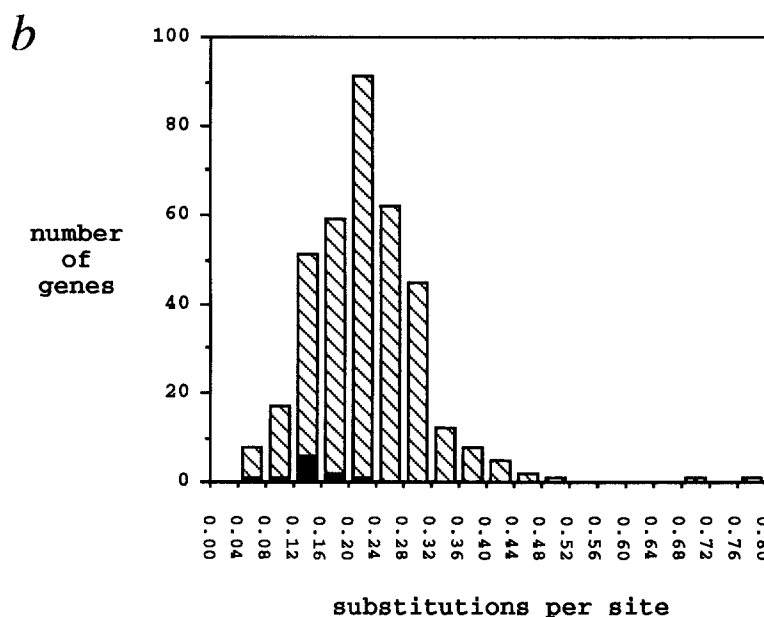
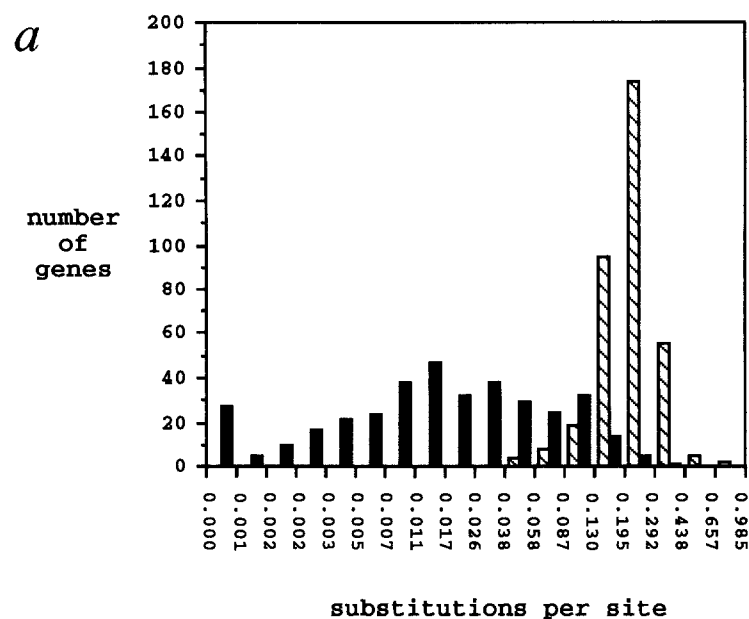


Fig. 2. Extent of nucleotide sequence divergence in 363 genes compared between mouse and rat. **a** Distribution of values of nonsynonymous (K_A , *solid*) and synonymous (K_S , *hatched*) substitutions per site. Note the semilog scale on the X-axis. **b** Distribution of values of synonymous substitutions per site (K_S), with X-linked genes highlighted.

which $K_S < K_A$, and the sequence is sufficiently short that the ratio may not be significantly less than one. Furthermore, it has not been established whether the estimation method of Li et al. (1985) will actually result in a K_S/K_A ratio of precisely 1.0 for a neutral DNA sequence; it is possible that the method suffers from a small bias, particularly in its treatment of 2-fold degenerate sites, which may overestimate K_S . Many of the other genes with low K_S/K_A ratios are either members of the immunoglobulin superfamily (Blast-1, Ig-C-delta, Ig-V-lambda, CD4, CD8 α , and CD43 antigens) or secretory proteins (SPOT-1, beta- and kappa-caseins, whey acidic protein). It is perhaps noteworthy that

interleukin-3 is the only gene for which $K_S < K_A$ despite the considerable overlap between the ranges of K_S and K_A (Fig. 2a; there are 99 genes for which K_A is greater than the lowest K_S seen in any gene and 350 genes for which K_S is less than the highest K_A seen).

Silent Sites

The weighted mean extent of synonymous substitution (K_S) between mouse and rat genes is 0.224 substitutions per site (standard deviation 0.084). (Li et al. 1987 reported a mean value of 0.237 for 24

genes.) The values for individual genes (Table 1 and Fig. 2) range 19-fold between 0.041 (± 0.014) in the Y-box transcription factor gene and 0.780 (± 0.179) in the gene for seminal vesicle secretory protein IV. [We have consulted the original publications of the sequences at either end of the range of K_S values (Fig. 2) to verify that these are bona fide homologous single-copy genes from mouse and rat; see also Data and Methods, above.]

The K_S values form quite a broad, roughly symmetric distribution about the mean (Fig. 2). Is this variation among genes in K_S greater than that expected by chance due to sampling effects? This can be tested statistically: if the variation in K_S is due to sampling error, the standard deviation of values around the mean K_S should be that expected under a binomial distribution. From equations (1)–(10) in Li et al. (1985) (using mean numbers of transitions and transversions per site in all genes, and the harmonic mean number of sites per gene), the standard error of the mean K_S in mouse-rat genes should be 0.043. This is only half as great as the observed standard deviation of K_S (0.084; Table 1), implying that there is genuine variation in synonymous substitution rates among genes. However, the synonymous substitution rates are considerably less variable than nonsynonymous rates: the coefficient of variation (the ratio of the standard deviation to the mean) of K_A is about four times that of K_S (Table 1).

The mean K_S for mouse-vs-rat genes is very close to the figure of 0.231 (± 0.009) calculated by Catzefflis et al. (1987, 1989) as the overall extent of sequence divergence between mouse and rat genomes, obtained by DNA-DNA hybridization of single copy genomic DNA (with correction for multiple-hit kinetics). We note, however, that K_S is generally larger (perhaps for the reason mentioned above) than the substitution rate observed at 4-fold degenerate sites alone (Li et al. 1985), even though both purport to estimate the rate of silent nucleotide substitution. The similarity in substitution rates at silent codon positions and in the genome as a whole (which is largely noncoding) strongly suggests that codon usage in mammalian genes is not constrained by natural selection, as we (Sharp 1989; Wolfe et al. 1989) and others (Eyre-Walker 1991) have argued elsewhere (also, see below).

To estimate the absolute rate of synonymous substitution (per unit time) it is necessary to know when mouse and rat last shared a common ancestor, and estimates of that date are quite controversial. Palaeontological evidence has been interpreted as indicating the date to be about 10 million years ago (see Catzefflis et al. 1992), suggesting a mean absolute rate of 11.2×10^{-9} substitutions per site per year per lineage (or a divergence rate down two lineages of 2.24% per million years). This is approx-

imately two times the average rate in mammals (including rodents) previously expressed as 4.7×10^{-9} substitutions per site per year (Li et al. 1985) or about 1% per million years down two lineages (Wilson et al. 1987). However, Wilson et al. (1987) have argued that the murid fossil record is difficult to interpret, and that the common ancestry of mouse and rat may date from as much as 35 million years ago. In that case, obviously, the rate estimates would be 3.5 times smaller.

Sex-linked Genes

Some of the heterogeneity in silent substitution rates among mammalian genes could be due to the location of some genes on the sex chromosomes. Miyata et al. (1987a,b) proposed that, because of the vastly greater number of cell generations that occur on average per year in spermatocytes as compared to oocytes, most nucleotide substitutions in mammalian genes arise in the male germline (so-called ‘‘male-driven’’ molecular evolution). As a consequence, the average silent substitution rates for autosomal, X-linked, and Y-linked genes should be in the ratio 1 : 0.67 : 2, reflecting the proportions of time (50%, 33%, and 100%, respectively) that each type of chromosome spends in males when averaged over many generations. More recently, Miyata et al. (1990) compared sequences between humans and rodent and found that the mean synonymous substitution rate for six X-linked genes was 58% of that in 35 autosomal genes. Some studies have, however, reported individual X-linked genes with high rates of synonymous substitution (e.g., Iizasa et al., 1989).

The data in Appendix 1 includes 11 genes known to be X-linked, and their mean K_S value is 0.143 (± 0.035). This is 64% of the mean value for the entire data set (which can be assumed to be largely autosomal), or 61% of the mean K_S (0.236 ± 0.063) for the 179 genes that have actually been mapped to an autosome in one or both species. The result is thus in close agreement with the predictions from the hypothesis of Miyata et al. (1987a,b). [The mean K_A for the X-linked genes is also low, 0.009 (± 0.009), as compared to 0.032 (± 0.049) for the whole data set.] The X-linked genes do not form a distinct group of slowly evolving sequences, but rather they lie within but toward the lower end of the distribution of K_S values seen in other genes (Fig. 2b): in order of increasing K_S the X-linked genes are ranked at positions 6 (*a-raf*), 10 (PLP), 37 (5-HT-1c receptor), 38 (androgen receptor), 40 (OTC), 51 (HPRT), 63 (NCAM-L1), 76 (PGK1), 84 (connexin-32), 104 (factor IX), and 220 (RPS4X) out of 363. Thus, there are heterogeneous rates for both

X-linked and autosomal genes, but there may well be a systematic (mutation rate) effect such that the former evolve at about two-thirds the rate of the latter. At present, there are no genes in the mouse-rat data set that are known to be Y-linked: such sequences would provide a valuable test of the "male-driven" molecular evolution hypothesis.

Variation in Silent-Site G + C Content

Synonymous codon usage is highly heterogeneous among mammalian genes, with the principal variation being in the base composition (G + C content) at silent sites (Ikemura 1985; Sharp et al. 1988). This variation seems to reflect local chromosomal base composition, since the G + C content at silent sites of individual mammalian genes is highly correlated with the G + C content of their introns and flanking sequences (Aota and Ikemura 1986; Ikemura and Aota 1988). This appears to be related to the organization of the mammalian genome into "isochores," i.e., domains of several hundred kilobases each having a relatively homogeneous base composition internally but being different from neighboring isochores (Bernardi et al. 1985). As expected, the mean G + C content at silent codon positions (designated GC_S and defined as the G + C content at 2-fold and 4-fold degenerate sites) in the mouse-rat genes studied here varies substantially from 32.8% to 95.4% (in the genes for nucleolar protein B23 and AGP/EBP transcription factor, respectively). The standard deviation of GC_S in the 363 genes is 11.1%, which is 3.7 times greater than that expected due to sampling error (under a binomial distribution around the mean GC_S of 62.0%).

The values of GC_S in homologous mouse and rat genes are very strongly correlated ($r = 0.973$). The greatest difference in GC_S between the species in a single gene is 10.2% in SVS IV (41.7% in mouse, 51.8% in rat), but the mean absolute difference in GC_S (ΔGC_S) is only 1.69%, with no significant net bias in either direction (61.8% mean GC_S in mouse; 62.2% in rat). Furthermore, there is no significant correlation between ΔGC_S and the mean GC_S in mouse and rat. This contrasts with the divergence in base composition that has occurred between human and rodent genes: human genes tend to have more extreme base compositions than their rodent homologues, a phenomenon that has been termed the "minor shift" (Mouchiroud et al. 1988; Bernardi et al. 1988).

Bernardi et al. (1988) have cited the strong correlation between silent-site G + C contents in homologous mouse and rat genes as evidence that their codon usage (or overall base composition) is being constrained by purifying selection. In fact,

the high correlation coefficient may be entirely attributable to the common ancestry of the sequences, since the mean K_S value of 0.224 substitutions per site implies that about 80% of the silent codon positions are identical by descent. We have investigated this by conducting a simple computer simulation: 363 ancestral sequences were generated (corresponding in length and G + C content to the silent sites in the rodent genes) and nucleotide substitutions (corresponding in number to the product of K_S and the number of silent sites in each gene) were then made at random assuming that all types of substitution are equally probable. Even with this simplifying assumption (which causes the G + C contents of the daughter sequences to converge toward 50%) the mean correlation coefficient from 1,000 replications was 0.950 (range 0.931–0.964), which is close to the correlation coefficient obtained with the real data. Indeed, if the simulation is conducted considering only 4-fold degenerate sites (instead of 2-fold and 4-fold sites combined), the mean correlation coefficient in simulations (0.958) exceeds that from the real data (0.947). There is thus little or no need to invoke causes other than common ancestry to explain the similarity of GC_S in the two species.

Excessive Variation in K_S and GC_S

The variation among genes in rates of nonsynonymous substitution can be interpreted in terms of differing extents of selective constraints on amino acid sequences (i.e., differences in the proportion of the sequence that is critical to the function of the protein). In contrast, the variation among genes described above in both silent substitution rates and silent-site G + C contents could be explained either by selective constraints on codon usage or by variation in the underlying mutation patterns among loci.

We and others have reported a relationship between the substitution rate and base composition at silent sites in mammalian genes, with genes of high GC_S (and possibly also those of low GC_S) having lower substitution rates than those of intermediate base compositions (Filipski 1988; Wolfe et al. 1989; Ticher and Graur 1989; Bulmer et al. 1991). When this trend is reinvestigated using the present data set (representing a fivefold increase in data over our 1989 study), the relationship between K_S and GC_S is less striking than originally reported (Fig. 3a). In particular, there are several large genes with intermediate GC_S contents (around 60%) that have low K_S values. Nevertheless, there is still a paucity of genes having both a high silent substitution rate and extreme values of GC_S .

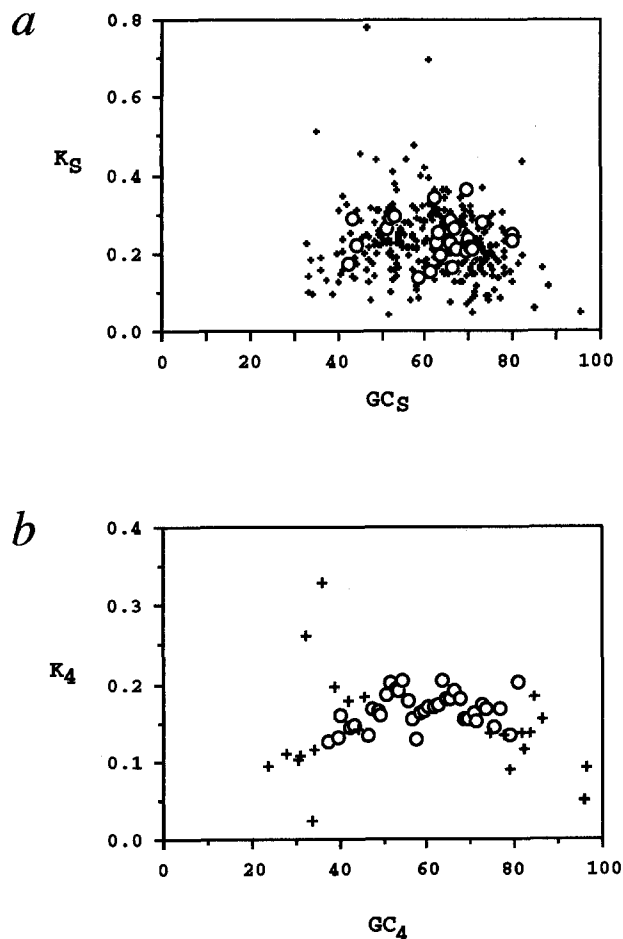


Fig. 3. Relationship between silent site divergence (between mouse and rat) and base composition (G + C content). **a** Number of synonymous substitutions per site (K_S) vs G + C content at silent sites (GC_S) for 363 genes. Circles represent genes with ≥ 500 silent sites. **b** Number of substitutions (K_4) vs G + C content (GC_4) at 4-fold degenerate sites, considering only codons where no more than one nucleotide substitution has occurred and correcting for multiple hit kinetics by the method of Tajima and Nei (1984). (Cf. Wolfe et al. 1989.) Each point represents the pooled result for all genes within a particular 1% interval of GC_4 . Circles represent intervals for which the number of sites compared is ≥ 500 .

For clarity, in our earlier study (Wolfe et al. 1989) we calculated substitution rates (K_4 , using the multiple-hits correction formula of Tajima and Nei 1984), and base composition (GC_4) only at 4-fold degenerate sites, and only in codons where there were no additional nucleotide substitutions. These modifications eliminate the correlation between synonymous and nonsynonymous substitution rates (see below), and also eliminate a possible artifact resulting from multiple-hits corrections on sequences of biased base compositions. Applying these criteria to the expanded data set does not result in a significant relationship between K_4 and GC_4 for individual genes. Nevertheless, an effect is apparent when the silent substitution rate is aver-

aged over all genes within each 1% interval of silent-site G + C content (Fig 3b). With the exception of a few intervals containing only one or a small number of genes, K_4 appears to peak at approximately 60% GC_4 , similar to what has been reported in other studies (Wolfe et al. 1989; Bulmer et al. 1991).

Thus, this study has documented the excessive variability of both substitution rates and base composition at silent sites in codons, but the relationship between the two remains unclear. In a recent analysis of sequences from three orders of mammal (Bulmer et al. 1991), we suggested that unknown factor(s) in addition to silent-site G + C content systematically affect the silent substitution rate (even after correcting for rate differences among lineages) in a manner that is consistent for a particular gene in different species. The discrepancy between the results of the present analysis and those based on the sequence data available four years ago (Filipski 1988; Wolfe et al. 1989; Ticher and Graur 1989) illustrates the need for very large data sets (and hence large numbers of nucleotide substitutions) before generalizations about substitution patterns or biases can be made with confidence. The data studied here comprise perhaps 0.5% of all genes in the mammalian genome and a total of 27,348 nucleotide differences between mouse and rat and so may constitute a representative sample.

The wide range of silent substitution rates and base compositions seen in groups of genes with similar functions (Appendix 1) also points to a lack of gene-expression-related constraint on codon usage. It has been suggested that mammalian codon usage is related to the major tissue of gene expression (Newgard et al. 1986), but this does not seem to be a general observation. Among the mouse-rat genes examined here, apolipoproteins A-II and A-IV are both expressed in liver, but have GC_S values of 60% and 78%, respectively (and serum albumin, also expressed in liver, has a GC_S value of 53%). There is also a wide variation in GC_S values among genes expressed in the testis (41% in a cytochrome *c* isoform, 78% in HSP70.2), brain (47% in calmodulin I; 79% in creatine kinase B), and pancreas (36% in α_2 -amylase, 73% in ribonuclease), for example. Alternatively, there are 13 ribosomal protein genes listed in Appendix 1 with K_S values that range between 0.172 and 0.359, and GC_S values between 48% and 70%. Because they are highly expressed, ribosomal protein genes generally provide the clearest examples of selective codon usage in organisms where translational efficiency is at a premium (e.g., *E. coli*, Ikemura 1985; *Saccharomyces cerevisiae*, Sharp et al. 1986; or *Drosophila melanogaster*, Shields et al. 1988).

Table 2. Numbers of substituted nucleotides at adjacent codon positions 2 and 3^a

		Position 3		Ratio ^b	chi ²
		Identical	Different		
A. All sites:					
Position 2	{ Identical	52,060	9,225	1.67	63.1
	{ Different	600	202		
B. Excluding C _p G dinucleotides:					
Position 2	{ Identical	50,276	7,154	2.09	107.5
	{ Different	484	170		

^a Only codons in which (in both species) position 3 is 4-fold degenerate are considered. Position 2 is always 0-fold degenerate.

^b Ratio of [D:D/(D:D+D:I)]/[I:D/(I:D+I:I)], where I and D represent identical and different nucleotides, respectively, and a colon separates the two bases of a dinucleotide. For example, I:D indicates a dinucleotide where the first base (at codon position 2) is identical in the two species, and the second base (at position 3) is different

Correlation Between Rates of Synonymous and Nonsynonymous Substitution—Evidence for Doublet Mutations

In the first extensive investigation of nucleotide substitution rates at synonymous and nonsynonymous sites in mammalian genes (Li et al. 1985), there was a positive correlation across genes between the two rates (Graur 1985). A similar correlation was reported by Miyata et al. (1987a,b), and by Ticher and Graur (1989), in comparisons of 39 and 42 genes, respectively, between humans and rodents. In the present mouse-rat data set the correlation coefficient (weighting each gene by its length) between K_S and K_A is $r = 0.45 (\pm 0.05)$, which is highly significant. This correlation is surprising because the two rates are expected to reflect different causes: the nonsynonymous rate should largely reflect protein sequence conservation, while the synonymous rate reflects the local mutation rate and any possible codon selection. The correlation could exist either because the two rates are similar over the whole gene (such that conserved proteins have low divergence at silent sites, for some reason) or because substitutions at adjacent nucleotide positions have occurred at a frequency greater than would be predicted from either substitution rate alone. Inspection of the 363 mouse-rat sequences suggests that the latter effect is occurring. For example, if the synonymous substitution rate is recalculated ignoring those codons where the species differ by more than a single nucleotide substitution, the resultant modified value of K_S is not significantly correlated with K_A ($r = 0.10 \pm 0.05$).

To look at this tandem substitution effect in more detail we identified homologous codons in the mouse and rat sequences where the third position is completely (4-fold) degenerate in both species. (Note that second positions are always 0-fold degenerate, i.e., nondegenerate.) These codons were then subdivided according to whether none, one, or both of these positions differ between the two species. The results reveal a highly significant excess of tandem nucleotide substitutions: where position 2 is unchanged, position 3 is changed in 15.1% of codons, but where a substitution has occurred at position 2, the fraction of differences at position 3 is 25.2% (Table 2A). It is well known (see, for example, Giannelli et al. 1991) that C_pG dinucleotides have a high mutation rate due to methylation of the cytosine (on either strand) followed by deamination, resulting in either a T_pG or a C_pA sequence. As a consequence, C_pG sites contribute a disproportionate number of single-base mutations, where an ancestral codon NCG is replaced by NTG or NCA in a descendant. If sites containing C_pG dinucleotides (in either species) are excluded from the analysis, the excess of tandem nucleotide substitutions becomes even more pronounced (Table 2B).

An excess of multiply-substituted codons was first commented upon by Fitch (1980), in a comparison of β -globin sequences from three orders of mammals. The question arises as to whether the excess of tandem substitutions arises from mutational events involving simultaneous substitution of two adjacent nucleotides or from two separate (consecutive) events. Lipman and Wilbur (1985) suggested that a nonsynonymous substitution may lead not only to an amino acid replacement but also to the replacement of an optimal codon for one amino acid by a nonoptimal codon for the new amino acid. Then, in the wake of the nonsynonymous substitution, a mutation leading to a synonymous change in the same codon may be positively selected as a way of generating an optimal codon for the newly specified amino acid. This seems unlikely because selection among synonymous codons is not thought to operate in many (if any) mammalian genes. (See above, and Sharp 1989.) Two other points should be noted. First, the selection coefficients driving sequences toward "re-optimized" codon usage could only be of the same order of magnitude as the nonsynonymous substitution rate. Second, if nonoptimal codons generated by a nonsynonymous mutation were selectively disadvantageous, the nonsynonymous change would be expected to be selected against even if the amino acid replacement itself were neutral.

To test whether the tandem differences reflect synonymous codon selection, we have examined

Table 3. Numbers of substituted nucleotides at adjacent codon positions 3 and 1^a

		Position 3		Ratio ^b	chi ²
		Identical	Different		
A. All sites:					
Position 1	{ Identical	45,657	7,968	1.62	109.5
	{ Different	1,304	413		
B. Excluding C _p G dinucleotides:					
Position 1	{ Identical	44,097	6,167	2.07	222.5
	{ Different	1,096	373		

^a Only codons in which (in both species) position 3 is 4-fold degenerate, and position 1 of the following codon is nondegenerate, are considered.

^b Ratio of $[D:D/(D:D+I:D)]/[D:I/(D:I+I:I)]$, where I and D represent identical and different nucleotides, respectively, and a colon separates the two bases of a dinucleotide. For example, I:D indicates a dinucleotide where the first base (at codon position 3) is identical in the two species, and the second base (at position 1 of the next codon) is different

the frequency of tandem differences at bases in neighboring codons (i.e., the 3:1 position involving the third base of a codon and the first base of the next codon). An excess of tandem substitutions, very similar in magnitude to that at 2:3 positions, is seen at the 3:1 position, both when all sites are considered (Table 3A) and when C_pG sites are excluded (Table 3B). This result *might* arise if codon pairs are under selection (see, e.g., Gouy 1987) but such selection should be secondary to selection within a codon, *if* that exists. Thus, from the magnitude of the excess of tandem substitutions at both 2:3 and 3:1 positions we conclude that this phenomenon most probably reflects mutational events simultaneously replacing both nucleotides. In support of this, we note that a significant excess of dinucleotide substitutions is seen in noncoding sequences from the primate eta-globin pseudogene region (our unpublished analyses of the data of Goodman et al. (1989)), where codon selection cannot be a factor. Furthermore, we have recently proposed that a significant excess of switches between the TCN and AGY groups of serine codons at highly conserved sites (in a range of proteins from a wide range of species) must also be due to doublet mutations (M. Averhof, K.H.W. and P.M.S., manuscript submitted).

Synonymous and nonsynonymous substitution rates have also been found to be correlated in studies on bacterial (Sharp and Li 1987; Sharp 1991) and chloroplast genes (Wolfe and Sharp 1988). In at least the case of bacteria, selection among synonymous codons can be a powerful force, and it will be interesting to investigate the possible causes of the K_S-K_A correlation in those species.

Conclusions

Data from individual genes have, in the past, pointed to the diversity of patterns of molecular evolution in mammalian genes. For example, the slow rate of evolution of actins (Alonso et al. 1986), the rapid rate of silent substitution in β_2 -microglobulin (Li et al. 1985), and the extremely high G + C content of the metallothionein-I gene (Durnam et al. 1980) have all been reported. In this study, we have attempted to put these observations into a clearer perspective by investigating the mean and range of several measures of the molecular evolutionary process, calculated from a large number of genes compared between a single pair of species (and so having a single divergence time). It is realistic to hope that the large number of genes studied here may be truly representative of those in the mammalian genome.

We have examined the evolution of silent sites in these genes in some detail. Silent sites are of particular interest because they were initially expected to be neutral (King and Jukes 1969). However, while early investigations focused on the comparative homogeneity (relative to nonsynonymous substitution rates) of synonymous substitution rates among mammalian genes (Miyata et al. 1980; Kimura 1981), it is now clear that these rates vary quite substantially. It has also become apparent that alternative synonymous codons are not neutral in many genes in the genomes of bacteria, fungi, and even insects (reviewed in Ikemura 1985; Sharp et al. 1988; Andersson and Kurland 1990), and that different intensities of synonymous codon selection lead to variation in synonymous substitution rates among genes (Sharp and Li 1987, 1989; Sharp 1991). Since it is also clear that codon usage in mammalian genes is highly nonrandom, it might seem reasonable to speculate that codon selection accounts for the silent substitution rate variation; however, several lines of argument and evidence presented here suggest that this is not the case. First, different patterns of codon usage in different mammalian genes are not related to any obvious expression differences among the genes, and codon selection would not be expected to overcome random genetic drift in mammals due to their small effective population sizes. Second a slow rate of synonymous substitution in X-linked genes, and the correlation of nonsynonymous and synonymous rates among genes can each be attributed to mutational causes. Several hypotheses, based on models of variation in DNA polymerase or DNA repair activities, have been put forward to explain why mutation rates and patterns might vary among different regions of the mammalian genome (Filipski 1987, 1988; Sueoka 1988;

Wolfe et al. 1989; Wolfe 1991). However, whether mutational causes alone are responsible for all of the variation seen in GC_S and K_S remains to be established. More detailed analysis of the data is in progress and should provide further insight into mammalian gene evolution.

Acknowledgments. This work utilized the facilities of the Irish National Centre for Bioinformatics. We are grateful to Manolo Gouy, Des Higgins, Liz Cowe, and Andrew Lloyd for help with database access, to Anne-Marie Murphy for some preliminary data analysis, and to Michael Bulmer and Adam Eyre-Walker for discussion. We also thank the managers and sponsors of the electronic mail servers at EMBL and GenBank. This study was supported in part by EOLAS grant SC/91/603.

References

- Alonso S, Minty A, Bourlet Y, Buckingham M (1986) Comparison of three actin-coding sequences in the mouse; evolutionary relationships between the actin genes of warm-blooded vertebrates. *J Mol Evol* 23:11–22
- Andersson SGE and Kurland CG (1990) Codon preferences in free-living microorganisms. *Microbiol Rev* 54:198–210
- Aota S, Ikemura T (1986) Diversity in G + C content at the third position of codons in vertebrate genes and its cause. *Nucleic Acids Res* 14:6345–6355, and correction 14:8702
- Bernardi G, Olofsson B, Filipinski J, Zerial M, Salinas J, Cuny G, Meunier-Rotival M, Rodier F (1985) The mosaic genome of warm-blooded vertebrates. *Science* 228:953–958
- Bernardi G, Mouchiroud D, Gautier C, Bernardi G (1988) Compositional patterns in vertebrate genomes: conservation and change in evolution. *J Mol Evol* 28:7–18
- Bulmer M, Wolfe KH, Sharp PM (1991) Synonymous nucleotide substitution rates in mammalian genes: implications for the molecular clock and the relationship of mammalian orders. *Proc Natl Acad Sci USA* 88:5974–5978
- Catzefflis FM, Sheldon FH, Ahlquist JE, Sibley CG (1987) DNA-DNA hybridization evidence of the rapid rate of murid rodent DNA evolution. *Mol Biol Evol* 4:242–253
- Catzefflis FM, Nevo E, Ahlquist JE, Sibley CG (1989) Relationship of the chromosomal species in the Eurasian mole rats of the *Spalax ehrenbergi* group as determined by DNA-DNA hybridization, and an estimate of the spalacid-murid divergence time. *J Mol Evol* 29:223–232
- Catzefflis FM, Aguilar J-P, Jaeger J-J (1992) Murid rodents: phylogeny and evolution. *Trends Ecol Evol* 7:122–126
- Chen YH, Pentecost BT, McLachlan JA, Teng CT (1987) The androgen-dependent mouse seminal vesicle secretory protein IV: characterization and complementary deoxyribonucleic acid cloning. *Mol Endocrinol* 1:707–716
- Cohen DR, Hapel AJ, Young IG (1986) Cloning and expression of the rat interleukin-3 gene. *Nucleic Acids Res* 14:3641–3658
- Dickerson RE (1971) The structure of cytochrome *c* and the rates of molecular evolution. *J Mol Evol* 1:26–45
- Dickinson DP, Mirels L, Tabak LA, Gross KW (1989) Rapid evolution of variants in a rodent multigene family encoding salivary proteins. *Mol Biol Evol* 6:80–102
- Dietrich W, Katz H, Lincoln SE, Shin H-S, Friedman J, Dracopoli NC, Lander ES (1992) A genetic map of the mouse suitable for typing intraspecific crosses. *Genetics* 131:423–447
- Durnam DM, Perrin F, Gannon F, Palmiter RD (1980) Isolation and characterization of the mouse metallothionein-I gene. *Proc Natl Acad Sci USA* 77:6511–6515
- Eyre-Walker A (1991) An analysis of codon usage in mammals: selection or mutation bias? *J Mol Evol* 33:442–449
- Filipinski J (1987) Correlation between molecular clock ticking, codon usage, fidelity of DNA repair, chromosome banding and chromatin compactness in germline cells. *FEBS Lett* 217:184–186
- Filipinski J (1988) Why the rate of silent codon substitutions is variable within a vertebrate's genome. *J Theor Biol* 134:159–164
- Fitch WM (1980) Estimating the total number of nucleotide substitutions since the common ancestor of a pair of genes: comparison of several methods and three beta hemoglobin messenger RNA's. *J Mol Evol* 16:153–209
- Giannelli F, Green PM, High KA, Sommer S, Lillicrap DP, Ludwig M, Olek K, Reitsma PH, Goossens M, Yoshioka A, Brownlee GG (1991) Haemophilia B: database of point mutations and short additions and deletions—second edition. *Nucleic Acids Res* 19 (suppl.):2193–2219
- Goodman M, Koop BF, Czelusniak J, Fitch DHA, Tagle DA, Slightom JL (1989) Molecular phylogeny of the family of apes and humans. *Genome* 31:316–335
- Gouy M (1987) Codon contexts in enterobacterial and coliphage genes. *Mol Biol Evol* 4:426–444
- Gouy M, Gautier C, Attimonelli M, Lanave C, di Paola G (1985) ACNUC—a portable retrieval system for nucleic acid sequence databases: logical and physical designs and usage. *Comp Appl Biosci* 1:167–172
- Graur D (1985) Amino acid composition and the evolutionary rates of protein-coding genes. *J Mol Evol* 22:53–62
- Hartl DL, Clark AG (1989) Principles of population genetics. Sinauer Associates, Sunderland, MA
- Higgins DG, Sharp PM (1989) Fast and sensitive multiple sequence alignments on a microcomputer. *Comp Appl Biosci* 5:151–153
- Hughes AL, Nei M (1988) Pattern of nucleotide substitution at major histocompatibility complex loci reveals overdominant selection. *Nature* 335:167–170
- Iizasa T, Taira M, Shimada H, Ishijima S, Tatibana M (1989) Molecular cloning and sequencing of human cDNA for phosphoriboxyl pyrophosphate synthetase subunit II. *FEBS Lett* 244:47–50
- Ikemura T (1985) Codon usage and tRNA content in unicellular and multicellular organisms. *Mol Biol Evol* 2:13–34
- Ikemura T, Aota S (1988) Global variation in G + C content along vertebrate genome DNA. *J Mol Biol* 203:1–13
- Jagodzinski LL, Sargent TD, Yang M, Glackin C, Bonner J (1981) Sequence homology between RNAs encoding rat alpha-fetoprotein and rat serum albumin. *Proc Natl Acad Sci USA* 78:3521–3525
- Kimura M (1977) Preponderance of synonymous changes as evidence for the neutral theory of molecular evolution. *Nature* 267:275–276
- Kimura M (1980) A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *J Mol Evol* 16:111–120
- Kimura M (1981) Estimation of evolutionary distances between homologous nucleotide sequences. *Proc Natl Acad Sci USA* 78:454–458
- Kimura M (1983) The neutral theory of molecular evolution. Cambridge University Press, Cambridge, England
- King JL, Jukes TH (1969) Non-Darwinian evolution. *Science* 164:788–798
- Lawrence JG, Hartl DL, Ochman H (1991) Molecular considerations in the evolution of bacterial genes. *J Mol Evol* 33:241–250

- Lipman DJ, Wilbur WJ (1985) Interaction of silent and replacement changes in eukaryotic coding sequences. *J Mol Evol* 21:161-167
- Li W-H, Wu C-I, Luo C-C (1985) A new method for estimating synonymous and nonsynonymous rates of nucleotide substitution considering the relative likelihood of nucleotide and codon changes. *Mol Biol Evol* 2:150-174
- Li W-H, Tanimura M, Sharp PM (1987) An evaluation of the molecular clock hypothesis using mammalian DNA sequences. *J Mol Evol* 25:330-342
- Miyata T, Yasunaga T, Nishida T (1980) Nucleotide sequence divergence and functional constraint in mRNA evolution. *Proc Natl Acad Sci USA* 77:7328-7332
- Miyata T, Hayashida H, Kuma K, Yasunaga T (1987a) Male-driven molecular evolution demonstrated by different rates of silent substitutions between autosome- and sex chromosome-linked genes. *Proc Jpn Acad Ser B* 63:327-331
- Miyata T, Hayashida H, Kuma K, Mitsuyasu K, Yasunaga T (1987b) Male-driven molecular evolution: a model and nucleotide sequence analysis. *Cold Spring Harbor Symp Quant Biol* 52:863-867
- Miyata T, Kuma K, Iwabe N, Hayashida H, Yasunaga T (1990) Different rates of evolution of autosome-, X chromosome- and Y chromosome-linked genes: hypothesis of male-driven molecular evolution. In: Takahata N, Crow JF (eds) *Population biology of genes and molecules*. Baifukan, Tokyo, Japan, pp 341-357
- Mouchiroud D, Gautier C, Bernardi G (1988) The compositional distribution of coding sequences and DNA molecules in humans and murids. *J Mol Evol* 27:311-320
- Nei M (1987) *Molecular evolutionary genetics*. Columbia University Press, New York
- Newgard CB, Nakano K, Hwang PK, Fletterick RJ (1986) Sequence analysis of the cDNA encoding human liver glycogen phosphorylase reveals tissue-specific codon usage. *Proc Natl Acad Sci USA* 83:8132-8136
- Ochman H, Wilson AC (1987) Evolution in bacteria: evidence for a universal substitution rate in cellular genomes. *J Mol Evol* 26:74-86
- Sharp PM (1989) Evolution at 'silent' sites in DNA. In: Hill WG, Mackay TFC (eds) *Evolution and animal breeding; reviews on molecular and quantitative approaches in honour of Alan Robertson*. C.A.B. International, Wallingford, UK, pp 24-32
- Sharp PM (1991) Determinants of DNA sequence divergence between *Escherichia coli* and *Salmonella typhimurium*: codon usage, map position, and concerted evolution. *J Mol Evol* 33:23-33
- Sharp PM, Li W-H (1987) The rate of synonymous substitution in enterobacterial genes is inversely related to codon usage bias. *Mol Biol Evol* 4:222-230
- Sharp PM, Li W-H (1989) On the rate of DNA sequence evolution in *Drosophila*. *J Mol Evol* 28:398-402
- Sharp PM, Tuohy TMF, Mosurski KR (1986) Codon usage in yeast: cluster analysis clearly differentiates between highly and lowly expressed genes. *Nucleic Acids Res* 14: 5125-5143
- Sharp PM, Cowe E, Higgins DG, Shields DC, Wolfe KH, Wright F (1988) Codon usage patterns in *Escherichia coli*, *Bacillus subtilis*, *Saccharomyces cerevisiae*, *Schizosaccharomyces pombe*, *Drosophila melanogaster* and *Homo sapiens*; a review of the considerable within-species diversity. *Nucleic Acids Res* 16:8207-8211
- Shields DC, Sharp PM, Higgins DG, Wright F (1988) "Silent" sites in *Drosophila* genes are not neutral: evidence of selection among synonymous codons. *Mol Biol Evol* 5:704-716
- Sueoka N (1988) Directional mutation pressure and neutral molecular evolution. *Proc Natl Acad Sci USA* 85:2653-2657
- Tajima F, Nei M (1984) Estimation of evolutionary distances between nucleotide sequences. *Mol Biol Evol* 1:269-285
- Ticher A, Graur D (1989) Nucleic acid composition, codon usage, and the rate of synonymous substitution in protein-coding genes. *J Mol Evol* 28:286-298
- Wilson AC, Carlson SS, White TJ (1977) Biochemical evolution. *Ann Rev Biochem* 46:573-639
- Wilson AC, Ochman H, Prager EM (1987) Molecular time scale for evolution. *Trends Genet* 3:241-247
- Wolfe KH (1991) Mammalian DNA replication: mutation biases and the mutation rate. *J Theor Biol* 149:441-451
- Wolfe KH, Sharp PM (1988) Identification of functional open reading frames in chloroplast genomes. *Gene* 66:215-222
- Wolfe KH, Sharp PM, Li W-H (1989) Mutation rates differ among regions of the mammalian genome. *Nature* 337:283-285
- Zuckerkindl E, Pauling L (1962) Molecular disease, evolution, and genic heterogeneity. In: Kasha M, Pullman B (eds) *Horizons in biochemistry*. Academic Press, New York, pp 189-225

Appendix 1. Molecular evolution and base composition of 363 nuclear genes compared between mouse and rat

Gene	Identity (%)		Nonsynonymous		Synonymous		GC ₃	AGC ₃	Accession number	
	DNA	prot.	K _A x 100	L _A	K _S x 100	L _S			(%)	(%)
7B2 neuroendocrine protein	95.7	97.1	1.2	493	17.1	137	53.3	-1.7	X15830	M63901
alpha-1 protease inhibitor (A1PI)	83.0	73.5	15.9	488	31.9	145	62.0	-2.9	X00945	D00675
nicotinic acetylcholine receptor epsilon (nAChRE)	94.8	96.1	1.8	1128	18.3	348	60.4	0.5	J04698	X13252
nicotinic acetylcholine receptor gamma (nAChRG)	92.9	97.1	1.3	236	30.3	76	77.7	3.5	X03818	X06364
lysosomal acid phosphatase	95.7	96.4	1.7	975	14.1	288	60.3	-0.3	X57199	M27893
alpha-actin (vascular smooth muscle)	94.7	99.5	0.2	886	28.1	246	64.7	-0.4	X07935	X06801
alpha-actin (cardiac)	95.4	100.0	0.0	221	25.8	58	63.2	-1.1	M15501	X00306
alpha-actin (skeletal muscle)	97.3	100.0	0.0	882	13.6	249	76.8	0.3	M12347	J00692
beta-actin (cytoplasmic)	96.6	99.7	0.1	875	16.2	251	72.8	-0.7	X03672	J00691
gamma-actin (smooth muscle)	95.1	100.0	0.0	882	25.7	246	63.2	-4.6	M26689	M23233
alpha-c-adaptin (subunit of clathrin AP-2 complex)	95.7	99.6	0.2	2167	19.5	647	63.5	2.7	X14972	X53773
alcohol dehydrogenase (ADH-1; liver)	90.2	89.9	5.0	865	31.9	260	57.0	-0.9	M11307	M15327
alpha-fetoprotein (AFP)	87.0	82.0	9.2	1439	36.3	376	53.4	-2.0	M16381	J00695
AGP/EBP transcription factor (LAP, IL6DBP, SFB)	98.7	99.3	0.3	668	4.7	220	95.4	0.3	M61007	X54626
alpha-lactalbumin	87.6	86.6	6.8	346	47.8	80	57.7	-6.5	M80909	X00461
delta-aminolevulinatase dehydratase (ALAD)	95.5	97.3	1.3	749	16.0	241	62.5	0.4	X13752	M14479
serum albumin (ALB)	89.6	90.0	4.9	980	37.7	274	52.8	-5.8	M16111	J00698
aldolase A	96.5	99.2	0.4	833	14.6	259	65.3	2.1	Y00516	X04261
aldolase C	94.6	97.8	1.0	523	24.0	155	62.3	0.4	X03796	X06984
murinoglobulin (alpha(1)-inhibitor III)	87.7	81.0	10.1	3391	26.4	959	51.0	-0.5	M65238	X52984
alpha-2-amylase (pancreatic)	94.0	93.4	3.3	1200	18.9	309	36.1	0.5	J00360	J00703
amyloid beta protein	94.7	99.0	0.4	1635	26.5	450	65.3	3.1	M18373	X07648
androgen receptor (Tfm locus)	96.3	97.6	1.1	2085	13.8	612	58.6	0.0	M37890	M20133
atrial natriuretic factor (ANF, ANP)	93.2	93.4	3.5	347	19.9	109	65.6	-2.9	K02781	K02062
angiotensinogen	90.0	86.8	6.4	1080	26.1	351	73.9	-1.4	J03046	L00091
apolipoprotein A-II	82.2	69.6	15.6	239	41.7	67	60.0	4.9	X04119	X03468

Appendix 1. Continued

Gene	Identity (%)		Nonsynonymous		Synonymous		GC _S (%)	ΔGC _S (%)	Accession number	
	DNA	prot.	K _A x 100	L _A	K _S x 100	L _S			mouse	rat
TIS11 gene (rat clone cMG1; begins MetThrThr....)	97.7	99.7	0.4	770	8.6	244	77.0	-0.6	M58566	X52590
TIS21 gene / PC3 (NGF-inducible protein)	94.3	97.5	1.1	370	25.3	104	78.2	-5.8	M64292	M60921
TIS7 gene / PC4 (interferon-related protein)	96.1	97.1	1.3	1050	13.9	297	43.3	-3.7	X17400	J04511
tumor necrosis factor (TNF) alpha	93.6	94.5	2.6	542	21.8	163	72.6	-3.2	Y00467	D00475
TNF receptor (Goodwin TNFR2; Lewis TNFR1)	89.4	83.9	8.5	1043	22.0	316	66.4	1.2	M59377	M63122
transition protein 1 (TP1)	96.4	98.2	1.6	128	8.6	37	58.5	0.0	X12521	M17096
transition protein 2 (TP2)	89.0	84.2	8.4	265	25.9	77	67.0	6.8	J03494	X14776
tissue plasminogen activator (tPA)	91.0	91.9	3.9	1312	33.6	365	61.3	-2.0	J03520	M23697
transin-1 (pTR1)	91.7	89.9	4.9	1118	24.7	307	46.5	-0.5	X63162	X02601
transthyretin (prealbumin)	91.0	93.2	4.0	332	29.8	109	61.3	-1.6	X04191	K03252
trkB oncogene (tyrosine protein kinase; gp145 form)	94.2	98.5	0.6	1926	28.2	537	65.8	-0.2	X17647	M55291
alpha-tropomyosin TM2 isoform 2 (fibroblast)	97.3	99.6	0.1	680	13.8	172	61.5	-0.6	M22479	M16432
beta-tropomyosin (skeletal muscle)	97.4	100.0	0.0	682	13.9	170	69.4	1.5	X12650	L00372
trypsin Ta (clone pMPT9)	91.4	87.8	6.0	568	20.2	170	62.7	-2.3	X04574	J00778
thyrotropin beta subunit (TSHB)	91.6	89.1	4.9	324	25.5	90	54.4	-4.6	M20536	M13897
alpha-tubulin (mouse M-alpha-1)	97.1	100.0	0.0	1052	14.4	301	64.3	0.6	M13445	J00797
UBF1 transcription factor	94.7	98.7	0.5	1843	29.2	449	69.1	3.3	X60831	M61726
UDP-glucuronosyltransferase (UDPGT; 17-beta hydroxysteroid)	90.2	87.9	5.9	1262	30.6	328	40.1	-1.5	X06358	Y00156
urate oxidase (UOX)	93.9	94.4	2.6	720	22.8	189	63.9	-2.7	M27695	M63593
vascular cell adhesion molecule (VCAM-1)	89.6	86.3	6.9	1722	27.7	495	48.2	-0.9	M84487	M84488
vitamin D binding protein (Gc globulin)	91.0	90.7	4.7	1110	30.3	306	51.0	-1.1	M55413	M60205
whey acidic protein (WAP)	80.8	68.2	19.6	314	31.8	82	59.2	4.6	X01157	X01153
tryptophan hydroxylase (Tph)	93.6	95.9	1.7	1046	27.4	286	56.7	-1.4	J04758	X53501
mKBP transcription factor (CRE-BP2; RARF2; C/EBP-related)	96.7	98.4	0.7	719	12.7	220	40.1	-1.3	M31629	M65148
Y-box binding protein 1 / enhancer factor I subunit A	99.0	99.7	0.1	744	4.1	222	51.5	-0.9	M60419	M57299
tyrosine hydroxylase	94.5	97.2	1.3	1137	21.9	357	66.9	-0.8	M69200	M10244
cytoplasmic protein-tyrosine phosphatase (PTP-S)	95.1	96.7	1.7	863	19.7	226	39.9	-1.3	M81477	X58828
ypt-1 (rab-1) oncogene (ras-related)	95.6	99.0	0.4	486	21.8	129	48.3	-7.7	Y00094	J02998

Notes: K_S and K_A are the numbers of synonymous and nonsynonymous substitutions per site, respectively (Li et al. 1985). L_S and L_A are the numbers of synonymous and nonsynonymous sites, respectively. GC_S is the mean silent-site G + C content in the mouse and rat genes. ΔGC_S is GC_S of the mouse gene minus that of the rat gene. In some cases the accession number listed is one of several from which the sequence used was constructed. A dash indicates that the sequence was taken directly from the literature.