

Weak selection revealed by the whole-genome comparison of the X chromosome and autosomes of human and chimpanzee

Jian Lu and Chung-I Wu*

Department of Ecology and Evolution, University of Chicago, Chicago, IL 60637

Communicated by Tomoko Ohta, National Institute of Genetics, Mishima, Japan, January 19, 2005 (received for review November 22, 2004)

The effect of weak selection driving genome evolution has attracted much attention in the last decade, but the task of measuring the strength of such selection is particularly difficult. A useful approach is to contrast the evolution of X-linked and autosomal genes in two closely related species in a whole-genome analysis. If the fitness effect of mutations is recessive, X-linked genes should evolve more rapidly than autosomal genes when the mutations are advantageous, and they should evolve more slowly than autosomal genes when the mutations are deleterious. We found synonymous substitutions on the X chromosome of human and chimpanzee to be less frequent than those on the autosomes. When calibrated against substitutions in the intergenic regions and pseudogenes to filter out the differences in the mutation rate and ancestral population size between X chromosomes and autosomes, X-linked synonymous substitutions are still 10% less frequent. At least 90% of the synonymous substitutions in human and chimpanzee are estimated to be deleterious, but the fitness effect is weaker than the effect of genetic drift. However, X-linked non-synonymous substitutions are $\approx 30\%$ more frequent than autosomal ones, suggesting the fixation of advantageous mutations that are recessive.

nearly neutral evolution | synonymous substitution | codon usage | purifying selection | positive selection

It is a central tenet of the neutral theory of molecular evolution that the fixation of neutral variations is prevalent, or even predominant, at the molecular level (1–3). The detection of natural selection, both positive and negative, has thus been the focus of many recent analyses of genomic sequences (4–7). However, there have been fewer attempts at measuring the strength of selection (5, 8–10). Ohta (2, 3, 11, 12) may have been the first to stress the importance of weak selection in molecular evolution. With genomic data, we may now be able to answer the question: “How much of the molecular divergence between species is affected by weak selection, the strength of which does not overwhelm genetic drift?”

In some cases, positive selection has been shown to be more extensive than predicted by the neutral theory (6, 7, 13), but the analyses may not always inform about the strength of selection (see refs. 14 and 15). In principle, the strength of selection can be estimated from the polymorphism data because the changes in the level and pattern of polymorphism are determined by recombination and selection (5, 16). However, this signature is short lived (17). Similarly, the strength of negative selection within populations has been estimated (5), but how much it contributes to the divergence between species is not clear. Negative selection, if sufficiently weak, does allow divergence to proceed. Some (but probably not all) synonymous substitutions in *Drosophila* are likely such cases (8). Ironically, in species with a small effective population size, such as humans, in which divergence under negative selection is even more plausible, the low codon usage bias makes the measurement of selection on synonymous changes impractical.

An alternative approach to measuring the extent and strength of selection, both positive and negative, is to contrast the evolution of X-linked and autosomal genes (18, 19). If the fitness effect of a mutation is (partially) recessive, then this effect can be more readily manifested on the X chromosome than on the autosomes (20). When the recessive mutations are still rare, they will nonetheless be expressed in the hemizygous males of the XY system. On the other hand, autosomal mutations have to become sufficiently frequent to form homozygotes to be influenced by natural selection under random mating. Therefore, if recessive mutations are common, X-linked genes will evolve more rapidly when advantageous and more slowly when deleterious. The X-linked vs. autosomal approach has its limitations, because there are other factors influencing the relative rates of evolution in the X chromosome and autosomes (21–23). To correct for these factors, DNA sequences from much of the genomes of two closely related species would be needed. Human and chimpanzee are two species providing sufficient genomic data.

Materials and Methods

DNA Sequences. The human–chimpanzee “reciprocal best” alignments (made by using the July 2003 human assembly and the Nov 13, 2003, Arachne 4X draft chimpanzee assembly from the Broad Institute, Cambridge, MA) were downloaded from <http://genome.ucsc.edu/goldenPath/hg16/vsPt0/axtBest>. The quality scores of chimpanzee sequences were downloaded from <http://hgdownload.cse.ucsc.edu/goldenPath/panTro1/bigZips>. The bases in the chimpanzee genome sequence with quality values < 20 were masked according to the position coordinate. The repetitive sequences in human–chimpanzee alignments were masked according to the annotations.

Intergenic sequences. We selected only aligned fragments that are longer than 1 kb and 50 kb away from any annotated human gene (ENSEMBL 22) as intergenic regions. The cutoff in divergence is 5%, although the level is $< 2\%$ for the majority of sequences used. In total, 71,667 autosomal and 6,115 X-linked fragments were used in this study.

Processed pseudogene sequences. Processed pseudogenes result from new retrotransposition. The sequences of 7,868 processed pseudogenes in the human genome were downloaded from www.pseudogene.org. These sequences were mapped to the corresponding human–chimpanzee alignments by using the FASTA34 program (24) and our own scripts. The mapped pseudogene alignments were then translated into protein sequences for further confirmation. In total, 277 X-linked and 4,735 autosomal pseudogene alignments were used. The divergence cutoff in human and chimpanzee is the same as for intergenic regions. (We also used another pseudogene data set from www.bork.embl-heidelberg.de/Docu/Human_Pseudogenes; the

Abbreviations: TV, transversional substitution rate; X/A, X-linked/autosomal genes.

*To whom correspondence should be addressed at: University of Chicago, 1101 East 57th Street, Chicago, IL 60637. E-mail: cwu@uchicago.edu.

© 2005 by The National Academy of Sciences of the USA

Table 1. The numbers of substitutions at synonymous sites, at intergenic sites, and in pseudogenes in human and chimpanzee

Site	$K_s \times 100$		$K_i \times 100$		$K_{\psi} \times 100$		$K_s/K_i [K_s/K_{\psi}]$	
	Unmasked	Masked	Unmasked	Masked	Unmasked	Masked	Unmasked	Masked
X	0.805 ± 0.033*	0.466 ± 0.027*	1.141 ± 0.006†	0.852 ± 0.005†	1.118 ± 0.044‡	0.799 ± 0.039‡	0.706 (0.635, 0.777) [0.720 (0.636, 0.810)]	0.547 (0.474, 0.623) [0.583 (0.496, 0.677)]
A	1.115 ± 0.009§	0.652 ± 0.007§	1.409 ± 0.002¶	1.026 ± 0.001¶	1.405 ± 0.011	0.958 ± 0.010	0.791 (0.778, 0.804) [0.793 (0.774, 0.810)]	0.635 (0.619, 0.649) [0.681 (0.658, 0.698)]
X/A	0.722 (0.650, 0.797)	0.714 (0.632, 0.798)	0.810 (0.802, 0.818)	0.830 (0.822, 0.840)	0.796 (0.740, 0.851)	0.834 (0.769, 0.903)	0.893 (P = 0.011) [0.907 (P = 0.059)]	0.861 (P = 0.010) [0.856 (P = 0.026)]

Repetitive sequences and CpG sites are shown. The 95% confidence intervals of the ratios are given in parentheses. The boldface X/A values are used in estimations. *P* is the probability that the X/A ratio is ≥ 1. X, X chromosomes; A, autosomes.

*No. of genes = 529.

†No. of fragments = 6,115.

‡No. of pseudogenes = 277.

§No. of genes = 12,779.

¶No. of fragments = 71,667.

||No. of pseudogenes = 4,735.

results are highly consistent and are presented in Table 4, which is published as supporting information on the PNAS web site).

Coding sequences. Ensemble human gene coding sequences (version 22, www.ensembl.org) were used to blast against the human sequences in the human–chimpanzee alignments. We used only the best hits, with similarity of 100% and alignable length ≥ 90 bp for further analysis. We extracted the coding sequences from the human–chimpanzee alignments by using the SIM4 program (25) and our own scripts. The extracted human and chimpanzee coding sequences were then translated into protein sequences for confirmation (by using the human ENSEMBL coding sequences as references). We discarded any alignment with stop codons or indels that can cause a frameshift in either species. For genes with alternative splicing, the transcripts with the least human–chimpanzee divergence were adopted. In the end, 12,779 autosomal and 529 X-linked genes were selected for analysis.

Analyses. The K_a (the number of nonsynonymous substitutions per nonsynonymous site) and K_s (the number of synonymous substitutions per synonymous site) values for coding sequences were computed by the method of Li (26). The divergence for intergenic sequences (K_i) and for pseudogenes (K_{ψ}) was computed by the two-parameter method of Kimura (27).

The substitution rate and pattern at the CpG sites are very different from the rest of the sequences, because C changes to T at a very high rate (28). If that had occurred, we would have observed CG ↔ TG changes or CG ↔ CA changes (CA being the reverse complement of TG). Therefore, we masked the CpG sites by removing all of the CG dinucleotides in the human–chimpanzee alignments.

The number of 4-fold degenerate sites in coding sequences and

the transversion substitutions at these sites (reported in Table 2) were counted by the divergence analyses implemented in the GCG 10.2 package.

Bootstrap. The bootstrap method (29) was used to infer 95% confidence interval estimations. In each replicate, the X-linked and/or autosomal genes (with the same sample size as the original data set) were randomly sampled with replacement from the original data set. The statistic values (mean K_a , K_s , K_{ψ} , or K_i) and the ratios [K_s/K_{ψ} , K_s/K_i , TV_{4-f}/TV_i (where TV is the transversion substitution rate), K_a/K_s , K_a/K_{ψ} , and K_a/K_i and those ratios for X-linked/autosomal genes (X/A)] were calculated based on the sampled data set(s). For each ratio estimation, the bootstrap method was replicated 10,000 times, and the 95% confidence intervals for that quantity were estimated (or the probability that the ratios were ≥ or ≤ 1; see Tables 1–3 for details).

Results

Lower K_s for X-Linked Genes than for Autosomal Genes. For coding sequence comparisons, we used 529 X-linked and 12,779 autosomal genes from the human–chimpanzee reciprocal best alignments (see *Materials and Methods*). In Table 1, the number of synonymous substitutions per 100 sites ($K_s \times 100$) between human and chimpanzee is shown to be 0.805 ± 0.033 for X-linked genes, ≈30% lower than the corresponding number for the autosomes (1.115 ± 0.009). Given the large sample sizes used, the difference is highly significant ($P < 0.0001$, Kolmogorov–Smirnov test).

Calibration for X-Linked–Autosomal Differences in Mutation Rate, Ancestral Polymorphism, and GC Content. There are two known sources that could contribute to a smaller K_s on the X chromo-

Table 2. The numbers of transversional substitutions at 4-fold degenerate sites and intergenic sites in human and chimpanzee

Site	$TV_{4-f} \times 100$		$TV_i \times 100$		Ratio TV_{4-f}/TV_i	
	Unmasked	Masked	Unmasked	Masked	Unmasked	Masked
X	0.236 ± 0.023 (48,265)	0.139 ± 0.018 (38,842)	0.411 ± 0.003 (24,886,727)	0.343 ± 0.003 (13,159,821)	0.574 ($P < 0.0001$)	0.405 ($P < 0.0001$)
A	0.325 ± 0.005 (1,530,548)	0.198 ± 0.005 (1,180,609)	0.474 ± 0.001 (739,637,558)	0.391 ± 0.001 (389,317,592)	0.686 ($P < 0.0001$)	0.506 ($P < 0.0001$)
X/A	0.726 ($P = 0.0007$)	0.702 ($P = 0.0012$)	0.867 ($P < 0.0001$)	0.877 ($P < 0.0001$)	0.837 ($P = 0.033$)	0.800 ($P = 0.029$)

Repetitive sequences and CpG sites are shown. The number of sites appears in parentheses. *P* is the probability that the X/A ratio is ≥ 1. X, X chromosomes; A, autosomes.

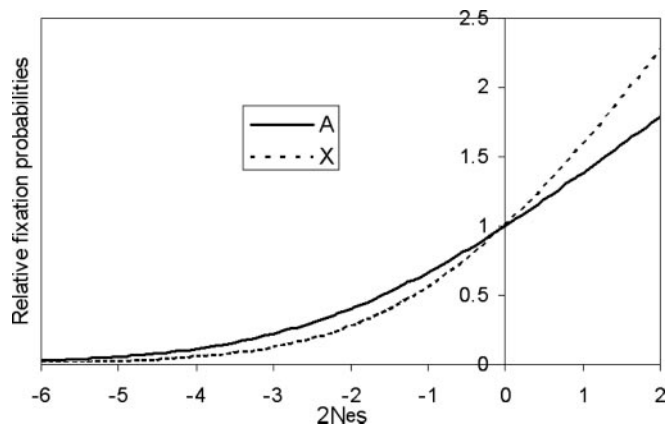


Fig. 1. The fixation probabilities for X-linked and autosomal mutations under selection, relative to the neutral ones. The fixation probabilities are functions of effective population size (N_e), selective coefficient (s), dominance coefficient (h), and the effective male-to-female ratio (m). (See *Supporting Methods*.) Data shown are for cases of $N_e = 10,000$, $h = 0.1$, and $m = 0.33$.

The above fitting ignores the possibility that a portion of synonymous changes might be advantageous. In what follows, we assume a portion of synonymous substitutions, p , to be deleterious with selection intensity s . The rest, $1-p$, are assumed to be advantageous to the same degree, as in ref. 35. (Because we are dealing with very weak selection, with $s < 1/2N_e$, we do not define a separate neutral class with $s = 0$.) For any $2N_e s$ value, there is a unique p value that would make the expected X/A ratio of Fig. 1 equal in fixation probability to an observed value. The observed X/A ratio of 0.90 (Table 1, the last two rows) was chosen to find the p value. The results are given in Fig. 2 and explained in the legends. In Fig. 2, we conclude that $P > 0.90$ and $0.5 < 2N_e s < 0.8$. Although Fig. 2 uses the same parameter values for h and m as Fig. 1, the tight clustering of curves in Fig. 3 suggests the robustness of the conclusion of large p and small $2N_e s$ values. Numerical evaluations corroborate the suggestion, and Table 6, which is published as supporting information on the PNAS web site, provides the point estimates of p and $2N_e s$ over a wide range of parameter values. In conclusion, the bulk of synonymous substitutions in human and chimpanzee are deleterious, but the selection intensity is extremely small, weaker than the effect of genetic drift.

Our estimation of p should not have been affected by any

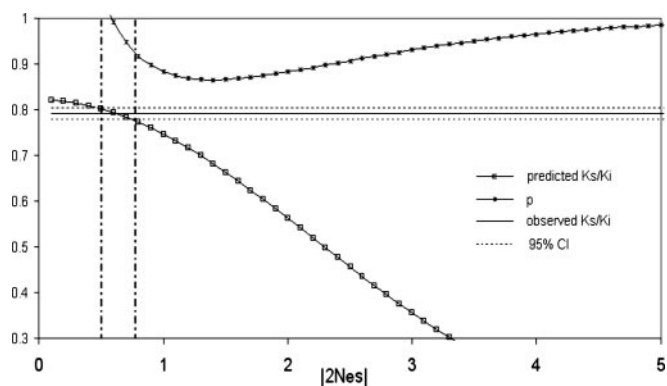


Fig. 2. Estimation of p (the proportion of synonymous substitutions under negative selection) and $2N_e s$ when the X/A ratio is 0.9 (see the last two rows of Table 1). For all values of $2N_e s$, p is always >0.85 . The K_s/K_i value is also a function of $2N_e s$; the case for $h = 0.1$, $m = 0.33$, and $N_e = 10,000$ is shown in the lower curve. The observed autosomal K_s/K_i at 0.791 further constrains $2N_e s$ to <1 and p to >0.9 . K_i is assumed to be the neutral rate.

possible difference in recombination rate between X chromosomes and autosomes. Because the X chromosome does not recombine in males, it is expected to experience less recombination than the autosomes. According to the general Hill–Robertson effect (36), negative selection against X-linked mutations should be less effective, and their K_s should be somewhat higher than the autosomal values. The observation is in the opposite direction.

Positive Selection Driving Nonsynonymous Substitutions. In contrast with synonymous substitutions, the rate of nonsynonymous substitutions is higher on the X chromosome than on the autosomes when calibrated against K_s , K_a/K_ψ , or K_i . In Table 3, the ratios of K_a/K_s , K_a/K_ψ , and K_a/K_i are given for X-linked and autosomal genes. Although K_a/K_s is often used to indicate the rate of nonsynonymous substitutions relative to the neutral rate, K_s is, in fact, not a neutral rate between human and chimpanzee. K_a/K_ψ and K_a/K_i are clearly better indicators of the selective constraints on nonsynonymous substitutions. The X/A ratios are 1.297 and 1.275, respectively, for K_a/K_i and K_a/K_ψ (both significantly >1 ; see Table 3). The results suggest that recessive advantageous mutations do leave a footprint in the genomes of human and chimpanzee, in the form of a higher average K_a for X-linked genes.

Discussion

The contrast in the evolutionary rates in the coding regions of X-linked and autosomal genes suggests that both positive and negative selection operate extensively on the genomes of human and chimpanzee. For recessive deleterious mutations, the intensity of selection is very weak on the homozygotes and even weaker on the heterozygotes. The conclusion that $>90\%$ of synonymous substitutions in human and chimpanzee are deleterious raises challenging issues that are addressed below.

The Flux Model for Synonymous Substitutions. What is the nature of negative selection against synonymous changes? In the “flux” model of synonymous changes (35), the change from a preferred codon to an unpreferred one is governed by negative selection with intensity $-s$. In the other direction, it is positive selection with intensity s . Although a gene with 100% preferred codons is assumed to be the fittest, the actual codon usage is kept at a mutation-selection equilibrium that is below the optimum. If a population is in equilibrium, the numbers of advantageous and deleterious substitutions should be equal. Sometimes, the equilibrium is shifted downward because of, say, a reduction in effective population size, and there would be a larger flux of deleterious substitutions from preferred to unpreferred codons than advantageous mutations going in the other direction. However, to account for the observation of $>90\%$ deleterious substitutions, codon usage would have to be experiencing a very drastic shift from a strongly biased pattern to a neutral one. Because mammals generally have only weak or no codon usage bias (28), the flux model cannot account for the observation.

A Model of Compensatory Mutations. We believe our observation of pervasive weak selection against synonymous changes between human and chimpanzee demands a new model for synonymous substitutions. The large number of deleterious synonymous changes must be compensated by a smaller number of advantageous mutations, each, on average, having a larger effect on fitness. The idea of compensatory mutations was proposed by Ohta (11, 37). We are hopeful that such models will be developed and tested. Here, we wish to suggest a possible outline of such a model. A fundamental assumption of the flux model is that genes with 100% preferred codons are the fittest. This assumption has not been tested empirically; in fact, contrary evidence exists (8, 38). Nor is the assumption biologically justified. It

