

Supporting Information

Nam et al. 10.1073/pnas.1419306112

SI Text

Background Selection. Background selection (BGS) reduces diversity locally in a genome by a process in which deleterious mutations are continuously pruned from the population, effectively reducing the number of genes from which future generations can be sampled. The strength of BGS is determined by the rate at which deleterious mutations enter the population, U , the recombination rate, R , and the strength of selection, s . If π_0 denotes the neutral diversity in the genome and if π is the diversity in a locus experiencing BGS, then the reduction in diversity is given by

$$\frac{\pi}{\pi_0} = \exp\left(-\frac{U}{s+R}\right)$$

(equation 6.24 in ref. 1).

We will assume a constant per-nucleotide deleterious mutation rate and recombination. Then, the rates U and R are both functions of the locus lengths, $U = uL$ and $R = rL$, where u is the per-nucleotide deleterious rate and r denotes the per-nucleotide pair recombination rate.

Comparing diversity inside with outside the regions of low diversity, we are interested in the relative reduction, which can be caused by changes to the selection rate, mutation rate, or recombination rate as factors of f_U , f_S , and f_R , respectively. The relative reduction can, thus, be expressed as

$$\frac{\exp\left(-\frac{f_U \times U}{f_S \times s + f_R \times R}\right)}{\exp\left(-\frac{U}{s+R}\right)}$$

Although all of the parameters in these equations are unknown, we do have some knowledge of their general order of magnitude or can choose conservative values to increase the relative reduction in diversity explainable by BGS.

BGS is strongest when selection is weak, but when s is very small, we expect the evolution to be nearly neutral. After s approaches $1/Ne$, we do not expect any BGS effects at all. Because Ne is on the order of 10,000–100,000, a lower limit on s is 10^{-4} – 10^{-5} . We consider both $s = 10^{-5}$ and $s = 10^{-4}$ and allow the selection inside the low-diversity regions to be one-tenth of the outside to make BGS stronger there.

We do not have recombination maps for most of the species considered, but from the human map, we know that the mean recombination rate on the X chromosome in the regions not showing reduced diversity is around 1.5 cm/Mb and that difference between the recombination rate inside and outside the reduced regions is less than a factor of two; therefore, we use $f_R = 0.5$.

We have very little information about what the deleterious mutation rate is, but we can use the mean human mutation rate, estimated to be 1.2×10^{-8} per generation (2), to explore various possibilities. If we use 1.2×10^{-8} , it would amount to assuming that 100% of mutations are under weak negative selection. By multiplying it with a number d between zero and one, we can interpret this number as the fraction of the loci that we believe are under selection. In Fig. S6, the columns correspond to different choices of d from 1% to 10% combined with different choices for s .

The rows in Fig. S6 correspond to different choices of f_U varying from 1 to 10 (therefore, the combination of $d = 0.1$ and $f_U = 10$ at the bottom right of Fig. S6 amounts to assuming that all sites within the low-diversity regions are under selection) combined with different choices for f_S , with $f_S = 1.0$ for no difference in the selection strength inside and outside of the low-diversity regions and $f_S = 0.1$ setting the selection inside the regions an order of magnitude lower. Fig. S6A shows the reduction in diversity compared with a neutral π_0 and the relative diversity of inside to outside of the low-diversity regions. In Fig. S6B, the dashed red line indicates 20%.

In the most extreme cases, we see a reduction in diversity of about 20% of the diversity within the low-diversity regions compared with outside but only in the cases where 100% of the nucleotides within regions are under selection. In the cases where 50% of the nucleotides are under selection within the regions, compared with 5% outside, the regions still retain about 50% of the diversity seen outside the regions.

Simulation of Sweeps. To assess the potential effect of hard and soft sweeps on diversity on the X chromosome, we performed a large number of simulations of a Wright–Fisher model exploring combinations of selection coefficients (s), effective population sizes (N), and frequencies of the selected variant at the onset of selection (f). We compute the time to the most recent common ancestor (TMRCA) along two recombining sequences and use this as a proxy for nucleotide diversity.

To simulate a selective sweep, we first sample frequency trajectories of a variant selected by s . We do this using rejection sampling (rejecting trajectories where the selected variant does not go to fixation). Trajectories for hard sweeps begin at 1 and proceed to $2N \times 3/4$ by repeated binomial sampling with probability parameter

$$\frac{N_{mut}}{N_{mut} + (N - N_{mut}) \times (1 - s)}$$

where N_{mut} is the number of selected variants in the previous generation. Trajectories for soft sweeps begin with an initial frequency f of the selected variant and are prepended with a trajectory from 1 to $f \times 2N \times 3/4$ representing variant frequency before the onset of the selection.

For each trajectory, we then consider a sample of two sequences representing 10 cm in length (equivalent to 10 Mb assuming a recombination rate of 1 cm/Mb). Because the effect of a sweep on flanking diversity is expected to be symmetric, we put the selected variant at the 5'-end position. To compute the TMRCA along the two sequences, we simulate backward the coalescence with recombination in discrete generations, allowing multiple mergers but only one recombination event per lineage in each generation.

Given a recombination event, the sequence downstream of the recombination point will become unlinked from the sweep with a probability equal to the frequency of chromosomes in the population not linked to the variant. A recombination event may similarly cause unlinked sequence fragments to again become linked to the selected variant with a probability equal to the frequency of chromosomes carrying the variant. Lineages that carry the selected variant can only share an immediate ancestor with other lineages that also carry this variant.

The simulation proceeds until all sequence segments separated by recombination events have found a most recent common

ancestor. For each combination of parameters s , N , and f , we perform 1,000 simulations, and the mean TMRCA along the 10 Mb is computed in bins of 10 kb (Fig. S7).

Complete List of the Great Ape Genome Project. Javier Prado-Martinez^a, Peter H. Sudmant^b, Jeffrey M. Kidd^{c,d}, Heng Li^e, Joanna L. Kelley^d, Belen Lorente-Galdos^a, Krishna R. Veeramah^f, August E. Woerner^f, Timothy D. O'Connor^b, Gabriel Santpere^a, Alexander Cagan^g, Christoph Theunert^g, Ferran Casals^a, Hafid Laayouni^a, Kasper Munch^h, Asger Hobolth^h, Anders E. Halager^h, Maika Malig^b, Jessica Hernandez-Rodriguez^a, Irene Hernandez-Herraez^a, Kay Prüfer^g, Marc Pybus^a, Laurel Johnstone^f, Michael Lachmann^g, Can Alkanⁱ, Dorina Twigg^c, Natalia Petit^a, Carl Baker^b, Fereydoun Hormozdiani^b, Marcos Fernandez-Callejo^a, Marc Dabad^a, Michael L. Wilson^j, Laurie Stevison^k, Cristina Camprubi^d, Tiago Carvalho^a, Aurora Ruiz-Herrera^{l,m}, Laura Vives^b, Marta Mele^a, Teresa Abelloⁿ, Ivanela Kondova^o, Ronald E. Bontrop^o, Anne Pusey^p, Felix Lankester^{q,r}, John A. Kiyang^q, Richard A. Bergl^r, Elizabeth Lonsdorf^r, Simon Myers^s, Mario Ventura^v, Pascal Gagneux^w, David Comas^a, Hans Siegismund^x, Julie Blanc^y, Lidia Agueda-Calpena^y, Marta Gut^y, Lucinda Fulton^z, Sarah A. Tishkoff^{aa}, James C. Mullikin^{bb}, Richard K. Wilson^z, Ivo G. Gut^y, Mary Katherine Gonder^{cc}, Oliver A. Ryder^{dd}, Beatrice H. Hahn^{cc}, Arcadi Navarro^{a,ff,gg}, Joshua M. Akey^b, Jaume Bertranpetit^a, David Reich^c, Thomas Mailund^h, Mikkel H. Schierup^{h,hh}, Christina Hvilsom^{x,ii}, Aida M. Andrés^g, Jeffrey D. Wall^k, Carlos D. Bustamante^d, Michael F. Hammer^f, Evan E. Eichler^{b,ji}, and Tomas Marques-Bonet^{a,gg}

^aInstitut de Biologia Evolutiva (Consejo Superior de Investigaciones Científicas-Universitat Pompeu Fabra), Parc de Recerca Biomèdica de Barcelona, Barcelona, Catalonia 08003, Spain; ^bDepartment of Genome Sciences, University of Washington, Seattle, WA 98195; ^cDepartment of Human Genetics, University of Michigan, Ann Arbor, MI 48109; ^dDepartment of Genetics, Stanford University, Stanford, CA 94305; ^eDepartment of Genetics, Harvard Medical School, Boston, MA 02115; ^fArizona Research Laboratories, Division of Biotechnology, University of Arizona, Tucson, AZ 85721; ^gDepartment of Evolutionary Genetics, Max Planck Institute for Evolutionary Anthropology, Leipzig 04103, Germany; ^hBioinformatics Research Centre and

^{hh}Department of Bioscience, Aarhus University, Aarhus C DK-8000, Denmark; ⁱFaculty of Engineering, Bilkent University, Ankara 06800, Turkey; ^jDepartment of Anthropology, University of Minnesota, Minneapolis, MN 55455; ^kInstitute for Human Genetics, University of California, San Francisco, CA 94143; ^lDepartament de Biologia Cel·lular, Fisiologia i Immunologia and ^mInstitut de Biotecnologia i de Biomedicina, Universitat Autònoma de Barcelona, Cerdanyola del Valles, Catalonia 08193, Spain; ⁿParc Zoològic de Barcelona, Barcelona, Catalonia 08003, Spain; ^oBiomedical Primate Research Centre, 2280 GH Rijswijk, The Netherlands; ^pDepartment of Evolutionary Anthropology, Duke University, Durham, NC 27708; ^qLimbe Wildlife Centre, BP 878 Limbe, Cameroon; ^rPaul G. Allen School for Global Animal Health, Washington State University, WA 99164; ^sNorth Carolina Zoological Park, Asheboro, NC 27205; ^tDepartment of Psychology, Franklin and Marshall College, Lancaster, PA 17604; ^uDepartment of Statistics, Oxford University, Oxford OX1 3TG, United Kingdom; ^vDepartment of Genetics and Microbiology, University of Bari, Bari 70126, Italy; ^wDepartment of Cellular and Molecular Medicine, University of California at San Diego, La Jolla, CA 92093; ^xDepartment of Biology, Bioinformatics, University of Copenhagen, Copenhagen 2200, Denmark; ^yCentro Nacional de Analisis Genómico, Parc Científic de Barcelona, Barcelona, Catalonia 08028, Spain; ^zGenome Sequencing Center, Washington University School of Medicine, St. Louis, MO 63108; ^{aa}Department of Biology and Genetics and ^{cc}Departments of Medicine and Microbiology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104; ^{bb}National Institutes of Health Intramural Sequencing Center, Bethesda, MD 20892; ^{cc}Biological Sciences, University at Albany, State University of New York, Albany, NY 12222; ^{dd}Genetics Division, San Diego Zoo's Institute for Conservation Research, Escondido, CA 92027; ^{ff}Instituto Nacional de Bioinformática, UPF, Barcelona, Catalonia 08003, Spain; ^{gg}Institut Catalana de Recerca i Estudis Avançats, Barcelona, Catalonia 08010, Spain; ⁱⁱCopenhagen Zoo, Frederiksberg DK 2000, Denmark; and ^{jj}Howard Hughes Medical Institute, Seattle, WA 98195

1. Durrett R (2008) *Probability Models for DNA Sequence Evolution* (Springer, Berlin).

2. Kong A, et al. (2012) Rate of de novo mutations and the importance of father's age to disease risk. *Nature* 488(7412):471–475.

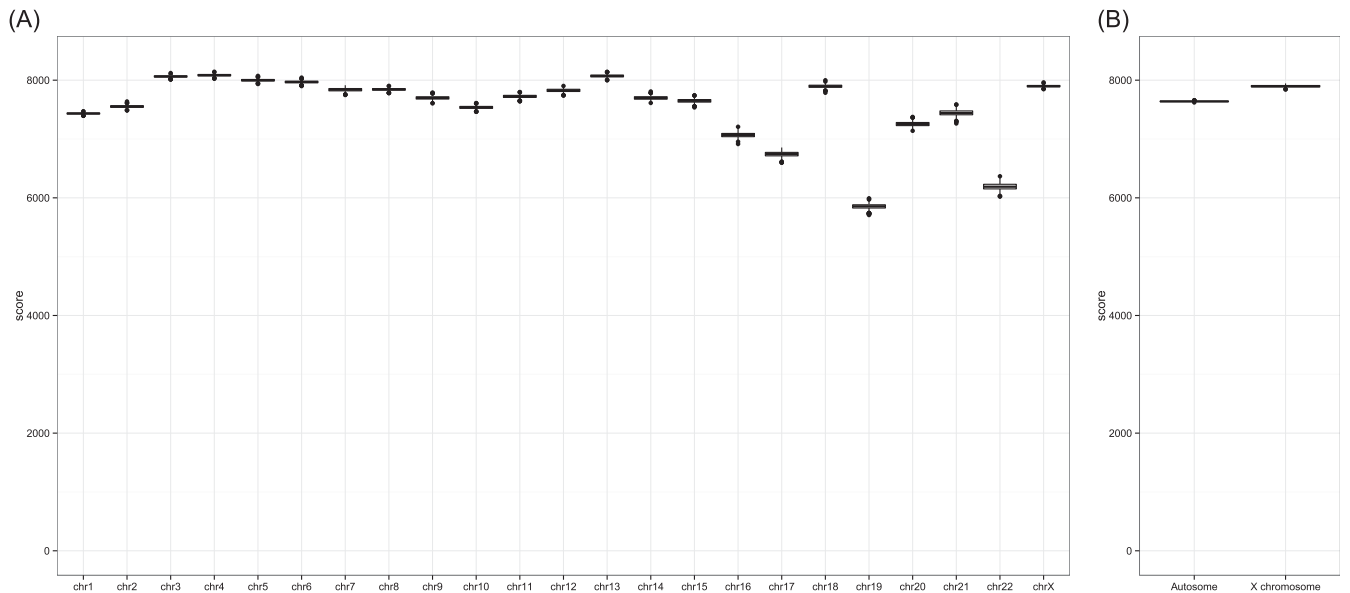


Fig. S1. The SNP quality score. (A) The SNP score in the variant call format files of each chromosome. (B) The SNP score of total autosomes and X chromosomes. The boxplots show the 95% confidence intervals calculated with 1,000 times bootstrapping replicates sampled from 1-Mb windows.

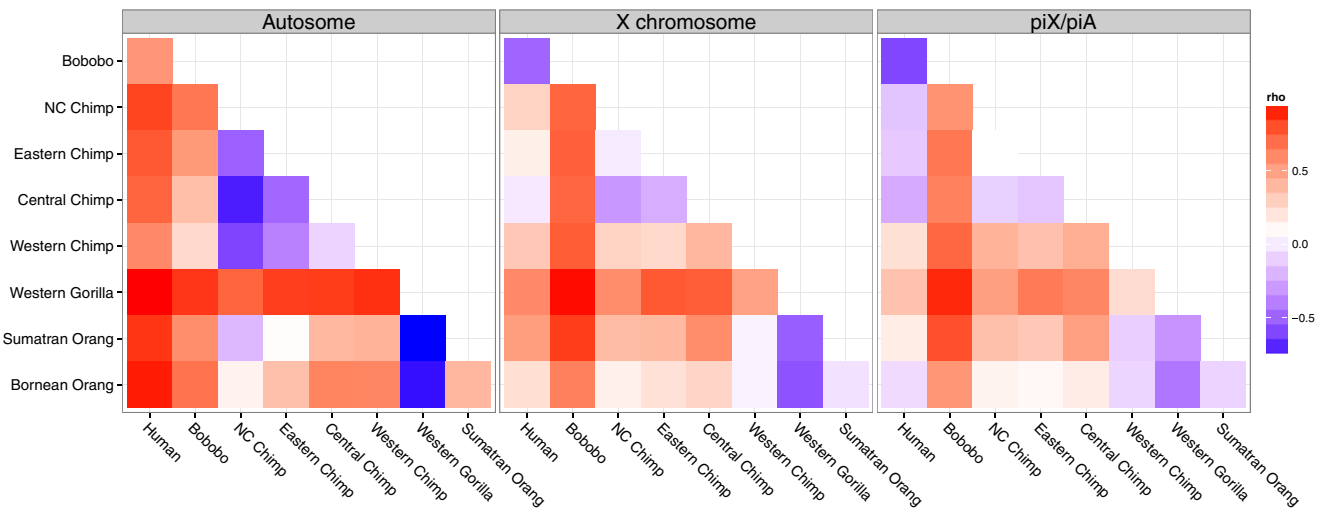


Fig. S2. Comparison of reduction in diversity between pairs of species. Heat maps of the correlation of the diversity ratio between two species and the physical distance from genes for the autosomes, the X chromosome, and the diversity ratio of X chromosomes to autosomes. The red color suggests a steeper relationship, with distance for the species indicated on the y axis; a blue color suggests a steeper relationship for the species on the x axis. NC, Nigeria-Cameroon.

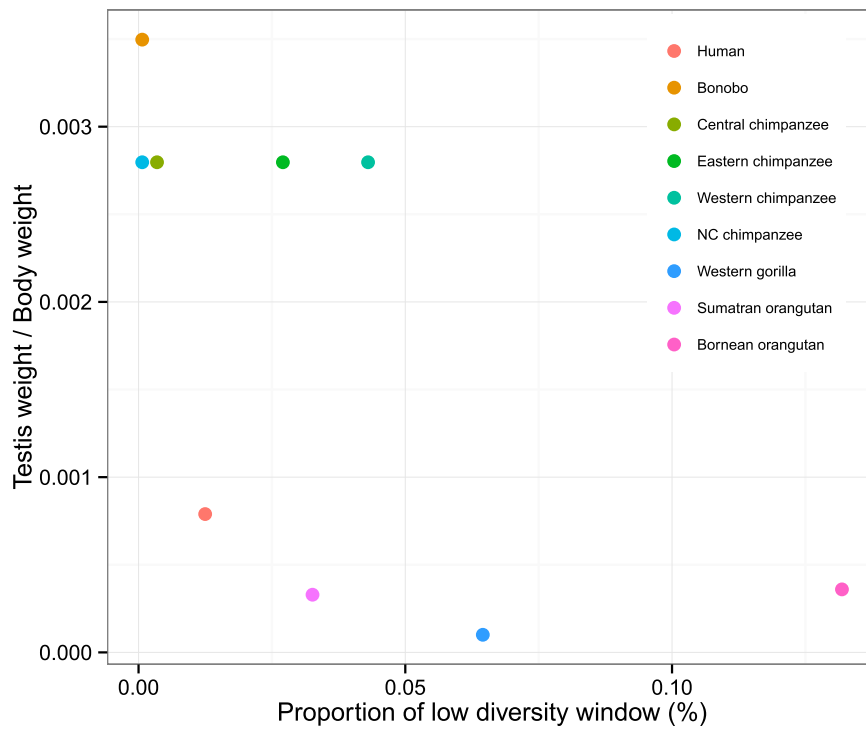


Fig. S9. Relationship between testis size and strength of sweeps. The y axis is the ratio of testicle to body weight, and the x axis is the proportion of windows that has π less than 20% of chromosomal average. NC, Nigeria–Cameroon.

Table S1. The information on taxa used in this study

Common name	Scientific name	No. of males	No. of females
Human	<i>Homo sapiens</i>	9	0
Bonobo	<i>Pan paniscus</i>	2	11
Central chimpanzee	<i>Pan troglodytes troglodytes</i>	1	3
Eastern chimpanzee	<i>Pan troglodytes schweinfurthii</i>	2	4
Western chimpanzee	<i>Pan troglodytes verus</i>	4	1
Nigeria–Cameroon chimpanzee	<i>Pan troglodytes ellioti</i>	4	6
Eastern lowland gorilla	<i>Gorilla beringei graueri</i>	2	1
Western lowland gorilla	<i>Gorilla gorilla gorilla</i>	4	23
Sumatran orangutan	<i>Pongo abelii</i>	1	4
Bornean orangutan	<i>Pongo pygmaeus</i>	1	4

