

# Analyses of X-linked and autosomal genetic variation in population-scale whole genome sequencing

Srikanth Gottipati<sup>1,3</sup>, Leonardo Arbiza<sup>1,3</sup>, Adam Siepel<sup>1</sup>, Andrew G Clark<sup>1,2</sup> & Alon Keinan<sup>1</sup>

**The ratio of genetic diversity on chromosome X to that on the autosomes is sensitive to both natural selection and demography. On the basis of whole-genome sequences of 69 females, we report that whereas this ratio increases with genetic distance from genes across populations, it is lower in Europeans than in West Africans independent of proximity to genes. This relative reduction is most parsimoniously explained by differences in demographic history without the need to invoke natural selection.**

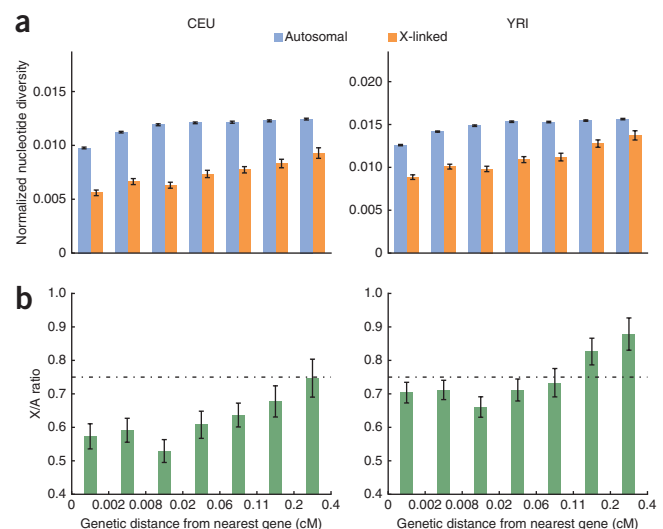
The genetic diversity of chromosome X is expected, under equilibrium assumptions, to be three-quarters of that of the autosomes in a population where the two sexes have an identical distribution of offspring numbers. However, deviations from this ratio can result from at least four forces known to have been prevalent in human history: (i) sex-biased demographic events leading to different effective population sizes of males and females; (ii) changes in population size over time (as chromosome X is proportionally more sensitive to recent epochs, owing to its reduced effective population size<sup>1</sup>); (iii) natural selection, which also affects chromosome X differently; and (iv) differences in mutation rates between sexes or between chromosome X and the autosomes. The possible effect of these forces on human genetic variation has received recent attention: Hammer *et al.*

reported that nucleotide diversity is higher than expected on chromosome X, with a mean X-to-autosome diversity ratio (X/A) of 0.9 across six populations and with no significant differences between populations<sup>2,3</sup>. Another recent study reported a significantly reduced X/A ratio in non-African populations relative to West Africans, beyond the reduction expected from known historical changes in population size<sup>4</sup>, with similar conclusions having been drawn from analyses of interpopulation allele frequency differences and the distribution of allele frequencies within populations<sup>4-7</sup>.

Estimates of the absolute X/A ratio are sensitive to details of the methods used to obtain them, including normalization by divergence from an outgroup<sup>8</sup> and differences in SNP ascertainment biases between chromosome X and the autosomes. To eliminate factors of this kind, we examined the relative X/A ratio between different populations. To compare the diversity of chromosome X and the autosomes in different populations, we considered intergenic SNPs from whole-genome sequences of 36 West African (YRI) and 33 European (CEU) females from the 1000 Genomes Project<sup>9</sup>, following rigorous quality control (**Supplementary Methods**). We normalized estimates of nucleotide diversity by divergence from a primate outgroup to correct for differences in mutation rates. Genome-wide X/A ratio estimates were  $0.73 \pm 0.016$  in YRI and  $0.61 \pm 0.018$  in CEU (**Supplementary Table 1**; normalization by divergence from rhesus macaque), which are consistent with previous estimates<sup>4</sup> and support a reduced ratio in non-Africans relative to Africans.

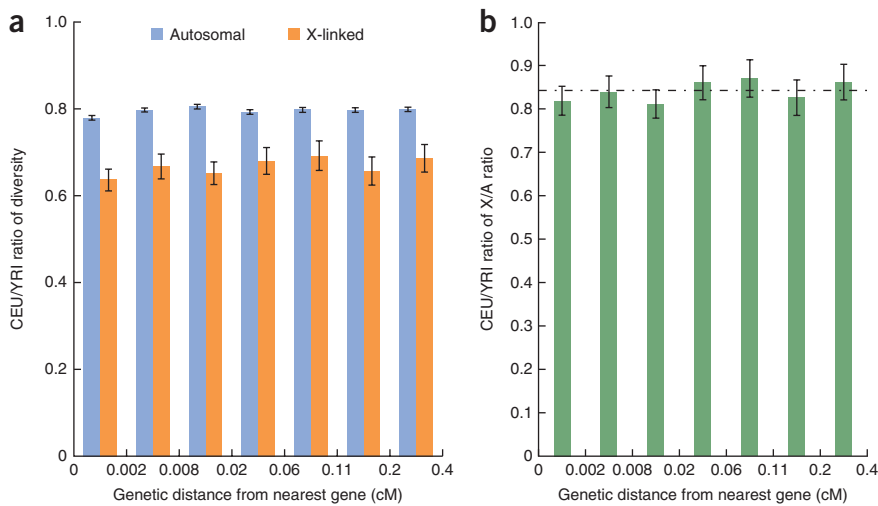
To examine the effect of natural selection, we partitioned the data by genetic distance from the nearest gene. Both X-linked and autosomal diversity increase with distance from genes (**Fig. 1a**;

**Figure 1** Autosomal, X-linked and absolute X/A diversity increase with genetic distance from the nearest gene. **(a)** Nucleotide diversity normalized by genetic divergence from rhesus macaque for a partition of the genome by distance from the nearest gene (**Supplementary Methods**). There are different scales of the y axis for the two populations (CEU and YRI), which are proportional to autosomal normalized diversity. **(b)** X/A ratios corresponding to the estimates from **a** (horizontal dashed line represents the expectation of three-quarters). In all panels, x-axis labels represent the boundaries between partitions, which were selected such that each partition encompassed an equal fraction of chromosome X (**Supplementary Fig. 2**). Error bars denote standard error estimated by a block bootstrap approach (**Supplementary Methods**). We obtained similar results when we used divergence from orangutan for normalization (**Supplementary Fig. 3**) and observed similar trends when we considered only levels of human nucleotide diversity, without any normalization by divergence (**Supplementary Fig. 4**).



<sup>1</sup>Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, New York, USA. <sup>2</sup>Department of Molecular Biology and Genetics, Cornell University, Ithaca, New York, USA. <sup>3</sup>These authors contributed equally to this work. Correspondence should be addressed to A.K. (ak735@cornell.edu).

Received 3 March; accepted 7 June; published online 24 July 2011; doi:10.1038/ng.877



**Figure 2** Relative autosomal, X-linked and X/A diversity are not correlated with genetic distance from the nearest gene. **(a,b)** For a partition of the genome as in **Figure 1**, X-linked and autosomal nucleotide diversity in CEU divided by the corresponding in YRI **(a)** and X/A ratio in CEU divided by X/A ratio in YRI **(b)**. Estimates <1 in **a** reflect the reduced diversity in non-Africans, most notably due to the out-of-Africa population bottleneck. Estimates <1 in **b** indicate a reduction in X-linked diversity compared to autosomal diversity that is specific to non-Africans (the horizontal dashed line denotes the estimate based on pooled, genome-wide intergenic data). Error bars denote standard error estimated by a block bootstrap approach (**Supplementary Methods**). These results are independent of normalization by divergence, as normalizing diversity in both populations by the same divergence estimates would have canceled out in the CEU-to-YRI ratio.

$P = 0.002$  and  $P = 0.077$  for CEU,  $P = 0.0008$  and  $P = 0.070$  for YRI). This increase in diversity with distance from genes closely matches predictions of the model of McVicker *et al.*<sup>10</sup> for both the autosomes and chromosome X (**Supplementary Fig. 1**;  $P < 0.01$  for all four cases), consistent with a diversity-reducing effect of selection on linked sites, through purifying selection (background selection), positive selection (genetic hitchhiking) or both. In addition, the site frequency spectrum is skewed toward lower-frequency alleles closer to genes, as expected from the action of natural selection (**Supplementary Note**). The increase in diversity with distance from genes is greater for chromosome X than for the autosomes (**Fig. 1a**), suggesting that diversity reduction due to selection at linked sites has been a more powerful force on chromosome X. As a result, the X/A ratio increases with distance from genes (**Fig. 1b**;  $P < 0.001$  for both CEU and YRI), consistent with recent results from six individuals of European descent<sup>3</sup> and in line with the observation that the increase in interpopulation allele frequency differentiation as recombination rate decreases is greater for chromosome X than for the autosomes<sup>11</sup>. The high X/A ratio observed in the loci sequenced by Hammer *et al.*<sup>2,3</sup> is in accordance with the large distance from genes and high local recombination rate of these loci (**Fig. 1b**).

So far, we have shown that the absolute X/A ratio is likely to have been strongly influenced by natural selection. To test whether the observed differences between Africans and non-Africans are also due to differential selective forces, we studied the relative levels of diversity between populations, considering the CEU-to-YRI ratio of nucleotide diversity (relative diversity) and the CEU-to-YRI ratio of the X/A ratio (relative X/A). Neither X-linked nor autosomal relative diversity is sensitive to distance from genes (**Fig. 2a**;  $P = 0.28$  and  $P = 0.53$  in a test of correlation), and the levels of relative diversity are consistently lower for chromosome X than for the autosomes (**Fig. 2a**). As a consequence, the relative X/A ratio remains nearly constant across all distances from genes (**Fig. 2b**;  $P = 0.42$  in a test of correlation) and is

always consistent with the genome-wide estimate of  $0.84 \pm 0.03$ , despite the pronounced dependence of selective effects on proximity to genes. Kienan *et al.* also observed no clear relationship between relative X/A ratio and distance from genes<sup>4</sup>, and the improved methodology and richer data set used here enable us to more definitively establish that relative X/A ratio is not sensitive to proximity to genes.

The lack of correlation between relative X/A and distance from genes suggests that the difference in X/A ratio between populations cannot be attributed to the effects of diversity-reducing selection acting on genes. In contrast, several plausible demographic explanations have been offered for the observed differences between populations, including the increased effect of recent history on chromosome X<sup>1,4</sup> and sex-biased demographic events<sup>7,12</sup>. One such sex-biased event has been highlighted in a recent simulation study: waves of primarily male migration during the dispersal of modern humans out of Africa<sup>12</sup>. Another recent modeling study supports that for a demographic event to explain the observed differences, it would have to coincide with the timing of the out-of-Africa event<sup>7</sup>. Our

results indicate that the difference in X/A ratio between African and non-African populations primarily derives from demographic forces such as those explored in these studies. It would require a very specific, consistent and highly improbable form of population-specific natural selection to drive the observed pattern.

In principle, our results could be influenced by ascertainment biases stemming from differences in sequencing coverage and in the sample size of individual chromosomes between chromosome X and the autosomes. However, three features of our analysis minimized the effect of such biases. First, to equalize sample size and coverage, we considered only females in all analyses. Second, differential ascertainment biases are not likely to correlate with genetic distance from genes. Third, such biases are not likely to affect estimates of relative diversity and relative X/A because ascertainment is similar for the two population samples we compared.

In conclusion, we have shown that there is a positive correlation between X/A ratio and distance from genes, indicating that diversity on chromosome X has been shaped by selection at linked sites more than has diversity on the autosomes, probably due in large part to X-linked recessive variants being exposed in males<sup>13–15</sup>. We have also shown that the reduced X/A ratio in non-Africans relative to Africans remains essentially constant across a wide range of genetic distances from genes. It has been argued that demographic history is best studied by focusing on ‘neutral’ loci that are located as far as possible from known functional elements<sup>3</sup>. Our results lead us to propose a complementary approach of analyzing ratios of diversity between different populations, which is not sensitive to the effects of natural selection if these are similar on a genome-wide average across populations. Contrasting populations allows focus—with increased resolution—on events that occurred after their split, excluding their shared history. This is much in the same spirit as studying X/A ratio on the basis of interpopulation allele frequency differentiation<sup>4–7</sup>, which considers changes in allele frequencies accumulated after the populations split. In contrast to considering putatively neutral regions in a single

population, the approach of contrasting statistics between populations is also not sensitive to unannotated functional elements confounding the inference of 'neutral' loci, normalization by genetic divergence (Fig. 2) or differential ascertainment biases between the X chromosome and the autosomes. Finally, our approach allows us to include orders of magnitude more data, thereby providing increased statistical power. Here analysis of genome-wide sequences, in conjunction with an approach focusing on more recent epochs, allowed us to identify a reduction in X/A ratio in non-Africans that probably results from demographic events associated with the human dispersal out of Africa.

Note: Supplementary information is available on the Nature Genetics website.

#### ACKNOWLEDGMENTS

We thank R. Nielsen and T. Korneliussen for sharing their software (Supplementary Note); D. Reich and E. Zhong for advice; D. Chang and E. Gazave for comments on earlier versions of this manuscript; and the 1000 Genomes Project. This work was supported in part by NIH grant U01-HG005715 and by an Alfred P. Sloan Research Fellowship (A.K.).

#### AUTHOR CONTRIBUTIONS

A.K. conceived of and designed the study. S.G. and L.A. performed the experiments. S.G., L.A. and A.K. analyzed the results and performed statistical

analysis. A.S., A.G.C. and A.K. contributed analysis tools. A.K. wrote the paper, with review and contributions by all authors.

#### COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturegenetics/>.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. Pool, J.E. & Nielsen, R. *Evolution Int. J. Org. Evolution* **61**, 3001–3006 (2007).
2. Hammer, M.F., Mendez, F.L., Cox, M.P., Woerner, A.E. & Wall, J.D. *PLoS Genet.* **4**, e1000202 (2008).
3. Hammer, M.F. *et al. Nat. Genet.* **42**, 830–831 (2010).
4. Keinan, A., Mullikin, J.C., Patterson, N. & Reich, D. *Nat. Genet.* **41**, 66–70 (2008).
5. Amato, R. *et al. PLoS ONE* **4**, e7927 (2009).
6. Casto, A.M. *et al. Genome Biol.* **11**, R10 (2010).
7. Emery, L.S., Felsenstein, J. & Akey, J.M. *Am. J. Hum. Genet.* **87**, 848–856 (2010).
8. Bustamante, C.D. & Ramachandran, S. *Nat. Genet.* **41**, 8–10 (2009).
9. Durbin, R.M. *et al. Nature* **467**, 1061–1073 (2010).
10. McVicker, G., Gordon, D., Davis, C. & Green, P. *PLoS Genet.* **5**, e1000471 (2009).
11. Keinan, A. & Reich, D. *PLoS Genet.* **6**, e1000886 (2010).
12. Keinan, A. & Reich, D. *Mol. Biol. Evol.* **27**, 2312–2321 (2010).
13. Charlesworth, B., Morgan, M.T. & Charlesworth, D. *Genetics* **134**, 1289–1303 (1993).
14. Vicoso, B. & Charlesworth, B. *Evolution* **63**, 2413–2426 (2009).
15. Betancourt, A.J., Kim, Y. & Orr, H.A. *Genetics* **168**, 2261–2269 (2004).