

Genetic Structure of Human Populations

Noah A. Rosenberg,^{1*} Jonathan K. Pritchard,² James L. Weber,³
Howard M. Cann,⁴ Kenneth K. Kidd,⁵ Lev A. Zhivotovsky,⁶
Marcus W. Feldman⁷

We studied human population structure using genotypes at 377 autosomal microsatellite loci in 1056 individuals from 52 populations. Within-population differences among individuals account for 93 to 95% of genetic variation; differences among major groups constitute only 3 to 5%. Nevertheless, without using prior information about the origins of individuals, we identified six main genetic clusters, five of which correspond to major geographic regions, and subclusters that often correspond to individual populations. General agreement of genetic and predefined populations suggests that self-reported ancestry can facilitate assessments of epidemiological risks but does not obviate the need to use genetic information in genetic association studies.

Most studies of human variation begin by sampling from predefined “populations.” These populations are usually defined on the basis of culture or geography and might not reflect underlying genetic relationships (*1*). Because knowledge about genetic structure of modern human populations can aid in inference of human evolutionary history, we used the HGDP-CEPH Human Genome Diversity Cell Line Panel (*2, 3*) to test the correspondence of predefined groups with those inferred from individual multilocus genotypes (supporting online text).

The average proportion of genetic differences between individuals from different human populations only slightly exceeds that

between unrelated individuals from a single population (*4–9*). That is, the within-population component of genetic variation, estimated here as 93 to 95% (Table 1), accounts for most of human genetic diversity. Perhaps as a result of differences in sampling schemes (*10*), our estimate is higher than previous estimates from studies of comparable geographic coverage (*4–6, 9*), one of which also used microsatellite markers (*6*). This overall similarity of human populations is also evident in the geographically widespread nature of most alleles (fig. S1). Of 4199 alleles present more than once in the sample, 46.7% appeared in all major regions represented: Africa, Europe, the Middle East, Central/

REPORTS

South Asia, East Asia, Oceania, and America. Only 7.4% of these 4199 alleles were exclusive to one region; region-specific alleles were usually rare, with a median relative frequency of 1.0% in their region of occurrence (11).

Despite small among-population variance components and the rarity of "private" alleles, analysis of multilocus genotypes allows inference of genetic ancestry without relying on information about sampling locations of individuals (12–14). We applied a model-based clustering algorithm that, loosely speaking, identifies subgroups that have distinctive allele frequencies. This procedure, implemented in the computer program *structure* (14), places individuals into K clusters, where K is chosen in advance but can be varied across independent runs of the algorithm. Individuals can have membership in multiple clusters, with membership coefficients summing to 1 across clusters.

In the worldwide sample, individuals from the same predefined population nearly always shared similar membership coefficients in inferred clusters (Fig. 1). At $K = 2$

the clusters were anchored by Africa and America, regions separated by a relatively large genetic distance (table S1). Each increase in K split one of the clusters obtained with the previous value. At $K = 5$, clusters corresponded largely to major geographic regions. However, the next cluster at $K = 6$ did not match a major region but consisted largely of individuals of the isolated Kalash group, who speak an Indo-European language and live in northwest Pakistan (Fig. 1 and table S2). In several populations, individuals had partial membership in multiple clusters, with similar membership coefficients for most individuals. These populations might reflect continuous gradations in allele frequencies across regions or admixture of neighboring

groups. Unlike other populations from Pakistan, Kalash showed no membership in East Asia at $K = 5$, consistent with their suggested European or Middle Eastern origin (15).

In America and Oceania, regions with low heterozygosity (table S3), inferred clusters corresponded closely to predefined populations (Fig. 2). These regions had the largest among-population variance components, and they required the fewest loci to obtain the clusters observed with the full data. Inferred clusters for Africa and the Middle East were also consistent across runs but did not all correspond to predefined groups. For the other samples, among-population variance components were below 2%, and independent *structure* runs were less consistent. For $K \geq 3$, similarity coefficients for pairs of runs

Table 1. Analysis of molecular variance (AMOVA). Eurasia, which encompasses Europe, the Middle East, and Central/South Asia, is treated as one region in the five-region AMOVA but is subdivided in the seven-region design. The World-B97 sample mimics a previous study (6).

Sample	Number of regions	Number of populations	Variance components and 95% confidence intervals (%)		
			Within populations	Among populations within regions	Among regions
World	1	52	94.6 (94.3, 94.8)	5.4 (5.2, 5.7)	
World	5	52	93.2 (92.9, 93.5)	2.5 (2.4, 2.6)	4.3 (4.0, 4.7)
World	7	52	94.1 (93.8, 94.3)	2.4 (2.3, 2.5)	3.6 (3.3, 3.9)
World-B97	5	14	89.8 (89.3, 90.2)	5.0 (4.8, 5.3)	5.2 (4.7, 5.7)
Africa	1	6	96.9 (96.7, 97.1)	3.1 (2.9, 3.3)	
Eurasia	1	21	98.5 (98.4, 98.6)	1.5 (1.4, 1.6)	
Eurasia	3	21	98.3 (98.2, 98.4)	1.2 (1.1, 1.3)	0.5 (0.4, 0.6)
Europe	1	8	99.3 (99.1, 99.4)	0.7 (0.6, 0.9)	
Middle East	1	4	98.7 (98.6, 98.8)	1.3 (1.2, 1.4)	
Central/South Asia	1	9	98.6 (98.5, 98.8)	1.4 (1.2, 1.5)	
East Asia	1	18	98.7 (98.6, 98.9)	1.3 (1.1, 1.4)	
Oceania	1	2	93.6 (92.8, 94.3)	6.4 (5.7, 7.2)	
America	1	5	88.4 (87.7, 89.0)	11.6 (11.0, 12.3)	

¹Molecular and Computational Biology, 1042 West 36th Place DRB 289, University of Southern California, Los Angeles, CA 90089, USA. ²Department of Human Genetics, University of Chicago, 920 East 58th Street, Chicago, IL 60637, USA. ³Center for Medical Genetics, Marshfield Medical Research Foundation, Marshfield, WI 54449, USA. ⁴Foundation Jean Dausset-Centre d'Etude du Polymorphisme Humain (CEPH), 27 rue Juliette Dodu, 75010 Paris, France. ⁵Department of Genetics, Yale University School of Medicine, 333 Cedar Street, New Haven, CT 06520, USA. ⁶Vavilov Institute of General Genetics, Russian Academy of Sciences, 3 Gubkin Street, Moscow 117809, Russia. ⁷Department of Biological Sciences, Stanford University, Stanford, CA 94305, USA.

*To whom correspondence should be addressed. E-mail: noahr@usc.edu

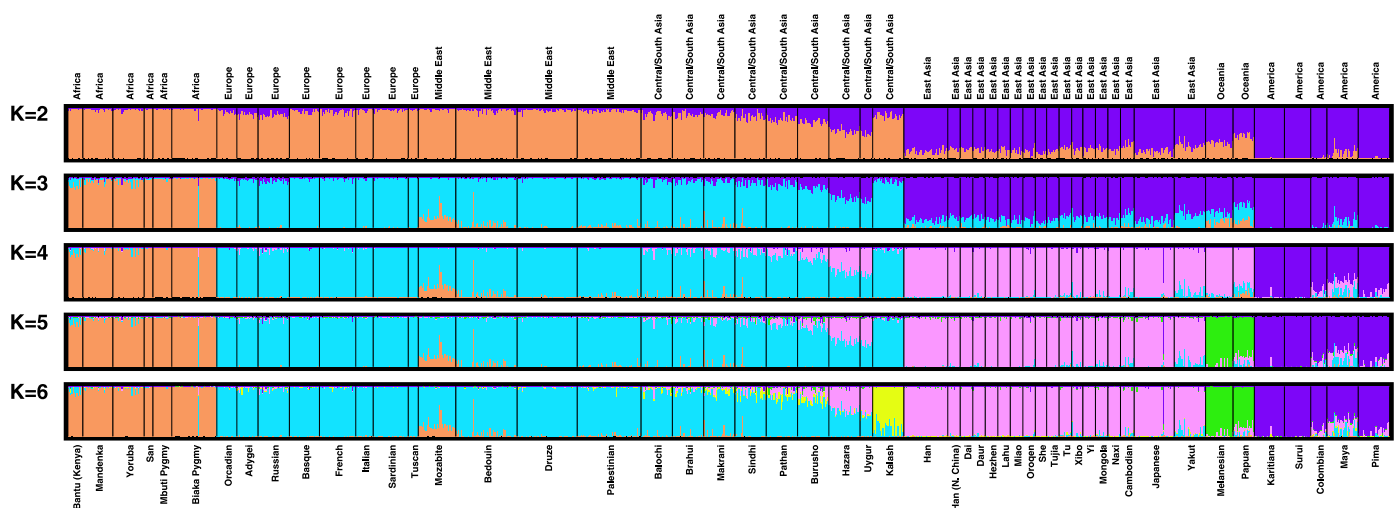


Fig. 1. Estimated population structure. Each individual is represented by a thin vertical line, which is partitioned into K colored segments that represent the individual's estimated membership fractions in K clusters. Black lines separate individuals of different populations. Populations are labeled below the figure, with their regional affiliations above it. Ten *structure* runs at each

K produced nearly identical individual membership coefficients, having pairwise similarity coefficients above 0.97, with the exceptions of comparisons involving four runs at $K = 3$ that separated East Asia instead of Eurasia, and one run at $K = 6$ that separated Karitiana instead of Kalash. The figure shown for a given K is based on the highest probability run at that K .

REPORTS

were typically moderate (0.1 to 0.85), rather than large (0.85 to 1.0). However, various patterns were observed across runs.

In East Asia, Yakut, whose language is Altaic, and Japanese, whose language is often classified as Altaic, were usually identified as distinctive. Other speakers of Altaic languages, including Daur, Hezhen, Mongola, Oroqen, and Xibo, all from northern China, shared a greater degree of membership with Japanese and Yakut than with more southerly groups from other language families, such as Cambodian, Dai, Han, Miao, Naxi, She, Tujia, and Yi. However, Tu, who speak an Altaic language and live in north-central China, largely grouped with the southern populations. Lahu, who speak a Sino-Tibetan language and were the least heterozygous pop-

ulation in the region, frequently separated despite their proximity with other groups sampled from southern China (16).

Eurasia frequently separated into its component regions, along with Kalash. Adygei, from the Caucasus, shared membership in Europe and Central/South Asia. Within Central/South Asia, Burusho of northern Pakistan, a linguistic isolate, largely separated from other groups, although less clearly than the genetic isolate, Kalash. Perhaps as a result of shared Mongol ancestry (15, 16), Hazara of Pakistan and Uyghur of northwestern China, whose languages are Indo-European and Altaic, respectively, clustered together. For Balochi, Makrani, Pathan, and Sindhi, all of whose languages are Indo-European, and less so for Dravidian-speaking Brahui, multiple

clusters were found, with individuals from many populations having membership in each cluster.

Europe, with the smallest among-population variance component (0.7%), was the most difficult region in which to detect population structure. The highest-likelihood run for $K = 3$ found no structure; in other runs, Basque and Sardinian were identified as distinctive. Russians variously grouped with Adygei and Orcadians; Russian-Orcadian similarity might derive from shared Viking contributions (17). French, Italians, and Tuscans showed mixed membership in clusters that contained other populations.

Because genetic drift occurs rapidly in small populations, particularly in those that are also isolated, these groups quickly accu-

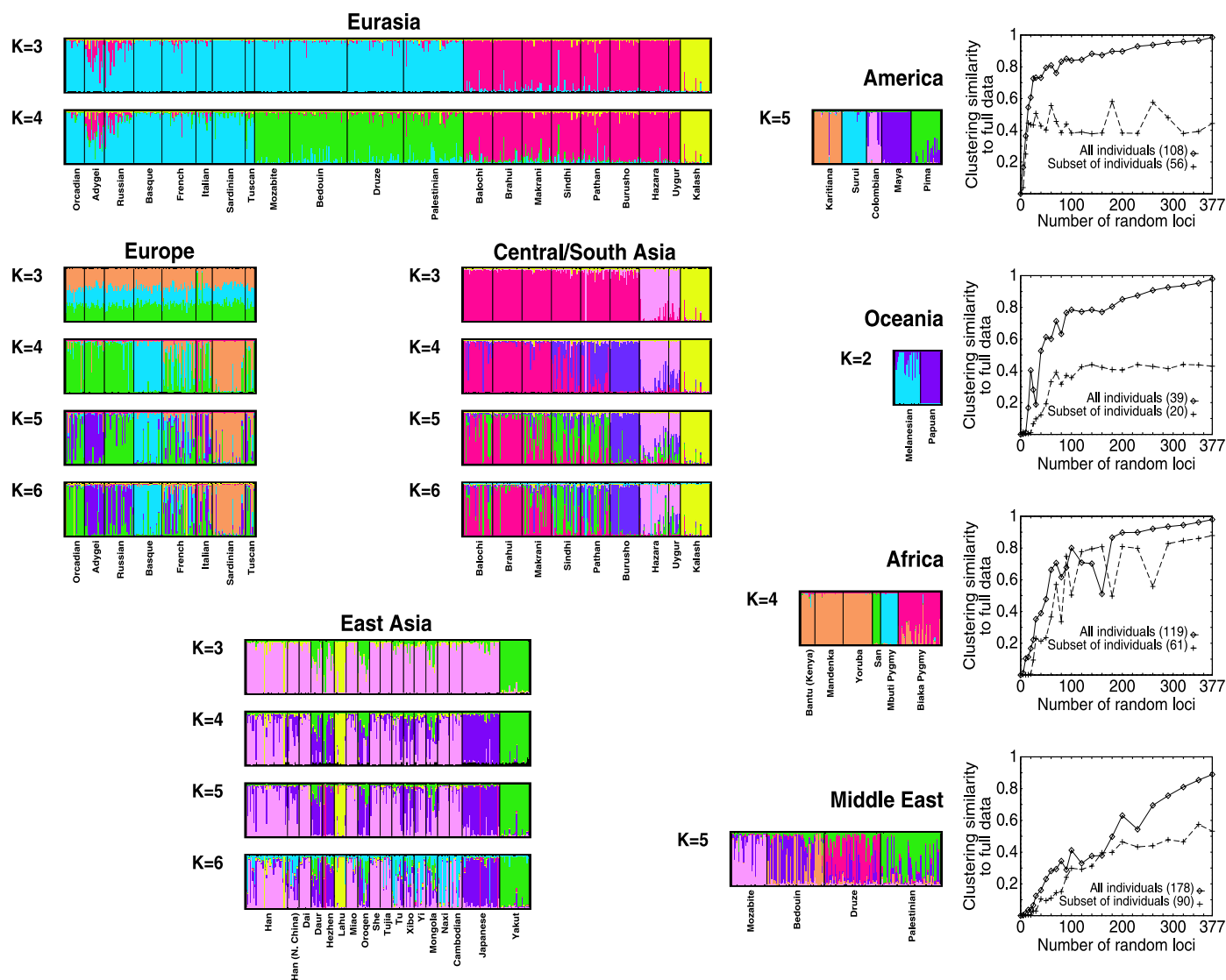


Fig. 2. Estimated population structure for regions. For America, Oceania, Africa, and the Middle East, solutions were consistent across 10 runs (all similarity coefficients above 0.97, 0.93, 0.97, and 0.86, respectively, except those involving one run with Africa that assigned many Biaka individuals partial membership with San). Values of K shown for these samples are the highest values for which this was true, and the highest

probability runs are shown. For remaining regions, solutions were more variable across runs, and the highest probability runs for various values of K are displayed. Graphs for America, Oceania, Africa, and the Middle East display median similarity coefficients between runs based on the full data and runs based on subsets of the data. Correspondence of colors across figures for different regions is not meaningful.

mulate distinctive allele frequencies. Thus, *structure* efficiently detects isolated and relatively homogeneous groups, even if the times since their divergences or exchanges with other groups are short (18). This phenomenon may explain the inferred distinctiveness of groups with low heterozygosity, such as Lahu and American groups, and those that are small and isolated, such as Kalash. Groups with larger sample sizes are also more easily separated; thus, the difficulty of clustering in East Asia was exacerbated by small sample sizes. Because sampling was population based, the sample likely produced clusters that were more distinct than would have been found in a sample with random worldwide representation. However, world-level boundaries between major clusters mostly corresponded to major physical barriers (oceans, Himalayas, Sahara).

The amount of among-group variation affects the number of loci required to produce clusters similar to those obtained with the full data. For the Middle East, with an among-population variance component of 1.3%, nearly all the loci were required to achieve a similarity of 0.8 to the clustering on the basis of full data, and use of more loci would likely produce more consistent clustering. For Oceania and Africa, only ~200 loci were needed; for the world sample, ~150 were needed (fig. S2), and ~100 were sufficient for America. Fewer loci would probably suffice for larger samples (18); conversely, accuracy decreased considerably when only half the sample was used (Fig. 2). The number of loci required would also decrease if extremely informative markers, such as those with particularly high heterozygosity (table S4), were genotyped (18). The loci here form a panel intended for use primarily in individuals of European descent (19). Although 10 of the loci had heterozygosity less than 0.5 in East Asia, none had similarly low European heterozygosities; thus, inference of subclusters using "random" markers might be more difficult than observed here, especially in Europe. However, the effect of excluding markers with low European heterozygosity is likely minimal, because generally high microsatellite heterozygosities ensure that relatively few loci are discarded on these grounds (20). The fact that regional heterozygosities here (table S3) follow the same relative order as and have nearly equal values to those of loci that were ascertained in a geographically diverse panel (12) provides further evidence that the ascertainment effect on heterozygosity estimates and on statistics derived from these estimates, such as genetic variance components (21), is small.

Genetic clusters often corresponded closely to predefined regional or population groups or to collections of geographically and linguistically similar populations. Among ex-

ceptions, linguistic similarity did not provide a general explanation for genetic groupings of populations that were relatively distant geographically, such as Hazara and Uyghur or Tu and populations from southern China. Our finer clustering results compared with other multilocus studies derive from our use of more data. General correspondence between regional affiliation and genetic ancestry has been reported (12–14), with clearer correspondence in studies that used more loci (13) than in those that used fewer loci (9, 22); we have further identified correspondence between genetic structure and population affiliation in regions with among-population variance components larger than 2 to 3%.

The structure of human populations is relevant in various epidemiological contexts. As a result of variation in frequencies of both genetic and nongenetic risk factors, rates of disease and of such phenotypes as adverse drug response vary across populations (22, 23). Further, information about a patient's population of origin might provide health-care practitioners with information about risk when direct causes of disease are unknown (23). Recent articles have considered whether it is preferable to use self-reported population ancestry or genetically inferred ancestry in such situations (22–25). We have found that predefined labels were highly informative about membership in genetic clusters, even for intermediate populations, in which most individuals had similar membership coefficients across clusters. Sizable variation in ancestry within predefined populations was detected only rarely, such as among geographically proximate Middle Eastern groups.

Thus, for many applications in epidemiology, as well as for assessing individual disease risks, self-reported population ancestry likely provides a suitable proxy for genetic ancestry. Self-reported ancestry can be obtained less intrusively than genetic ancestry, and if self-reported ancestry subdivides a genetic cluster into multiple groups, it may provide useful information about unknown environmental risk factors (23, 25). One exception to these general comments may arise in recently admixed populations, in which genetic ancestry varies substantially among individuals; this variation might correlate with risk as a result of genetic or cultural factors (24). In some contexts, however, use of genetic clusters is more appropriate than use of self-reported ancestry. In genetic case-control association studies, false positives can be obtained if disease risk is correlated with genetic ancestry (24, 26). Basing analyses on self-reported ancestry reduces the proportion of false positives considerably (25). However, association studies are usually analyzed by significance testing, in which slight differences in genetic ancestry between

cases and controls can produce statistically significant false-positive associations in large samples. Thus, errors incurred by using self-reported rather than genetic ancestry might cause serious problems in large studies that will be required for identifying susceptibility loci with small effects (26). Genetic clustering is also more appropriate for some types of population genetic studies, because unrecognized genetic structure can produce false positives in statistical tests for population growth or natural selection (27).

The challenge of genetic studies of human history is to use the small amount of genetic differentiation among populations to infer the history of human migrations. Because most alleles are widespread, genetic differences among human populations derive mainly from gradations in allele frequencies rather than from distinctive "diagnostic" genotypes. Indeed, it was only in the accumulation of small allele-frequency differences across many loci that population structure was identified. Patterns of modern human population structure discussed here can be used to guide construction of historical models of migration and admixture that will be useful in inferential studies of human genetic history.

References and Notes

1. M. W. Foster, R. R. Sharp, *Genome Res.* **12**, 844 (2002).
2. H. M. Cann *et al.*, *Science* **296**, 261 (2002).
3. Genotypes from this study are available at <http://research.marshfieldclinic.org/genetics/Freq/FreqInfo.htm>.
4. R. C. Lewontin, *Evol. Biol.* **6**, 381 (1972).
5. B. D. H. Latter, *Am. Nat.* **116**, 220 (1980).
6. G. Barbujani, A. Magagni, E. Minch, L. L. Cavalli-Sforza, *Proc. Natl. Acad. Sci. U.S.A.* **94**, 4516 (1997).
7. L. B. Jorde *et al.*, *Am. J. Hum. Genet.* **66**, 979 (2000).
8. R. A. Brown, G. J. Armelagos, *Evol. Anthropol.* **10**, 34 (2001).
9. C. Romualdi *et al.*, *Genome Res.* **12**, 602 (2002).
10. Smaller within-population variance components of comparable studies may result from their use of isolated and geographically well-separated populations to construct samples. Such a scheme might exaggerate among-group differences compared with those in the present sample, which had a smaller proportion of such populations. Indeed, when we restricted analysis to a set of populations that approximated a previous data set (6), we obtained a larger among-region component. Variance components also depend on sample sizes and on marker properties (7–9). Differential natural selection on protein variants across geographic regions might exaggerate among-group differences. Conversely, for a fixed level of within-group diversity, recurrent microsatellite mutations reduce among-group differences in comparison with those observed at markers for which each mutation produces a novel allele (28).
11. Recurrent mutation might be expected to influence allelic distributions considerably. However, widespread distributions of most alleles and the paucity of alleles found only in two disconnected regions suggest that recurrent mutations are only rarely followed by independent drift to sizable frequencies in multiple regions (29).
12. A. M. Bowcock *et al.*, *Nature* **368**, 455 (1994).
13. J. L. Mountain, L. L. Cavalli-Sforza, *Am. J. Hum. Genet.* **61**, 705 (1997).
14. J. K. Pritchard, M. Stephens, P. Donnelly, *Genetics* **155**, 945 (2000).
15. R. Qamar *et al.*, *Am. J. Hum. Genet.* **70**, 1107 (2002).
16. R. Du, V. F. Yip, *Ethnic Groups in China* (Lubrecht and Cramer, Port Jervis, NY, 1996).

REPORTS

17. J. Haywood, *The Penguin Historical Atlas of the Vikings* (Penguin Books, London, 1995).
18. N. A. Rosenberg *et al.*, *Genetics* **159**, 699 (2001).
19. J. L. Weber, K. W. Broman, *Adv. Genet.* **42**, 77 (2001).
20. A. R. Rogers, L. B. Jorde, *Am. J. Hum. Genet.* **58**, 1033 (1996).
21. M. Urbanek, D. Goldman, J. C. Long, *Mol. Biol. Evol.* **13**, 943 (1996).
22. J. F. Wilson *et al.*, *Nature Genet.* **29**, 265 (2001).
23. N. Risch, E. Burchard, E. Ziv, H. Tang, *Genome Biol.* **3**, comment2007.1 (2002).
24. D. C. Thomas, J. S. Witte, *Cancer Epidemiol. Biomark. Prev.* **11**, 505 (2002).
25. S. Wacholder, N. Rothman, N. Caporaso, *Cancer Epidemiol. Biomark. Prev.* **11**, 513 (2002).
26. J. K. Pritchard, P. Donnelly, *Theor. Popul. Biol.* **60**, 227 (2001).
27. S. E. Ptak, M. Przeworski, *Trends Genet.* **18**, 559 (2002).
28. L. Jin, R. Chakraborty, *Heredity* **74**, 274 (1995).
29. F. Calafell *et al.*, *Eur. J. Hum. Genet.* **6**, 38 (1998).
30. D. Altshuler, M. Cho, D. Falush, H. Innan, L. Kurina, J. Mountain, D. Nettle, M. Nordborg, M. Przeworski, N. Risch, D. Rosenberg, M. Stephens, D. Thomas, and E. Ziv provided helpful comments. The Mammalian Genotyping Service is supported by the National Heart, Lung, and Blood Institute. This work was supported by an NSF Biological Informatics Postdoctoral Fellowship (N.A.R.), a Burroughs-Wellcome Fund Hitchings Elion grant (J.K.P.), and NIH GM28428 (M.W.F.).

Supporting Online Material

www.sciencemag.org/cgi/content/full/298/5602/2381/DC1

Materials and Methods

Supporting Text

Figs. S1 and S2

Tables S1 to S4

References

19 June 2002; accepted 30 October 2002