

Worldwide Human Relationships Inferred from Genome-Wide Patterns of Variation

Jun Z. Li,^{1,2,*†} Devin M. Absher,^{1,2,*} Hua Tang,¹ Audrey M. Southwick,^{1,2} Amanda M. Casto,¹ Sohini Ramachandran,⁴ Howard M. Cann,⁵ Gregory S. Barsh,^{1,3} Marcus Feldman,^{4‡} Luigi L. Cavalli-Sforza,^{1‡} Richard M. Myers^{1,2‡}

Human genetic diversity is shaped by both demographic and biological factors and has fundamental implications for understanding the genetic basis of diseases. We studied 938 unrelated individuals from 51 populations of the Human Genome Diversity Panel at 650,000 common single-nucleotide polymorphism loci. Individual ancestry and population substructure were detectable with very high resolution. The relationship between haplotype heterozygosity and geography was consistent with the hypothesis of a serial founder effect with a single origin in sub-Saharan Africa. In addition, we observed a pattern of ancestral allele frequency distributions that reflects variation in population dynamics among geographic regions. This data set allows the most comprehensive characterization to date of human genetic variation.

In the past 30 years, the ability to study DNA sequence variation has dramatically increased our knowledge of the relationships among and history of human populations. Analyses of mitochondrial, Y chromosomal, and autosomal markers have revealed geographical structuring of human populations at the continental level (1–3) and suggest that a small group of individuals migrated out of eastern Africa and their descendants subsequently expanded into most of today's populations (3–6). Despite this progress, these studies were limited to a small fraction of the genome, to

limited populations, or both, and yield an incomplete picture of the relative importance of mutation, recombination, migration, demography, selection, and random drift (7–10). To substantially increase the genomic and population coverage of past studies (e.g., the HapMap Project), we have examined more than 650,000 single-nucleotide polymorphisms (SNPs) in samples from the Human Genome Diversity Panel (HGDP-CEPH), which represents 1064 fully consenting individuals from 51 populations from sub-Saharan Africa, North Africa,

Europe, the Middle East, South/Central Asia, East Asia, Oceania, and the Americas (11). This data set is freely available (12) and allows a detailed characterization of worldwide genetic variation.

We first studied genetic ancestry of each individual without using his/her population identity. This analysis considers each person's genome as having originated from K ancestral but unobserved populations whose contributions are described by K coefficients that sum to 1 for each individual. To increase computational efficiency, we developed new software, *frappe*, that implements a maximum likelihood method (13) to analyze all 642,690 autosomal SNPs in 938 unrelated and successfully genotyped HGDP-CEPH individuals (14). Figure 1A shows the results for $K = 7$; those for $K = 2$ through 6 are in fig. S1. At $K = 5$, the 938 individuals segregate into five continental groups, similar to those re-

¹Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305–5120, USA. ²Stanford Human Genome Center, Stanford University School of Medicine, Stanford, CA 94305–5120, USA. ³Department of Pediatrics, Stanford University School of Medicine, Stanford, CA 94305–5120, USA. ⁴Department of Biological Sciences, Stanford University, Stanford, CA 94305–5120, USA. ⁵Foundation Jean Dausset-Centre d'Etude du Polymorphisme Humain (CEPH), 75010 Paris, France.

*These authors contributed equally to this work.

†Present address: Department of Human Genetics, University of Michigan, 5789A MS II, Ann Arbor, MI 48109–5618, USA.

‡To whom correspondence should be addressed. E-mail: marc@charles.stanford.edu (M.F.); cavalli@stanford.edu (L.L.C.S.); myers@shgc.stanford.edu (R.M.M.)

ported in a microsatellite-based study of the same panel (3). At $K = 6$, the new component accounts for a major portion of ancestry for individuals from South/Central Asia, separating this region from the Middle East and Europe. This result differs from that in (3), where the sixth component contained the Kalash individuals, but South/Central Asia, the Middle East, and Europe were not clearly distinguished unless analyzed separately from the rest of the world. At $K = 7$, the new component occurs at highest proportions in the Middle Eastern populations, separating them from European populations. In many populations, ancestry is derived predominantly from

one of the inferred components, whereas in others, especially those in the Middle East and South/Central Asia, there are multiple sources of ancestry. For example, Palestinians, Druze, and Bedouins have contributions from the Middle East, Europe, and South/Central Asia. Burusho, Pathan, and Sindhi have an East Asian contribution. Hazara and Uyghur share a similar profile of combined South/Central Asian, East Asian, and European ancestry. In East Asia, only the Yakuts share ancestry with both Europe and America, although these contributions are small. Although much of sub-Saharan Africa, Europe, and East Asia appears to be homogeneous in Fig. 1A, finer

substructures can be detected when individual regions are analyzed separately. For example, we identified two components that separate the 16 East Asian populations and correspond to a north-south genetic gradient (fig. S2A). Han Chinese can be divided into a southern and a northern group. A similar analysis for South/Central Asia is shown in fig. S2B.

Mixed ancestries inferred from genetic data can often be interpreted as arising from recent admixture among multiple founder populations. In the current setting, however, the estimated mixed ancestry can be due either to recent admixture or to shared ancestry before the diver-

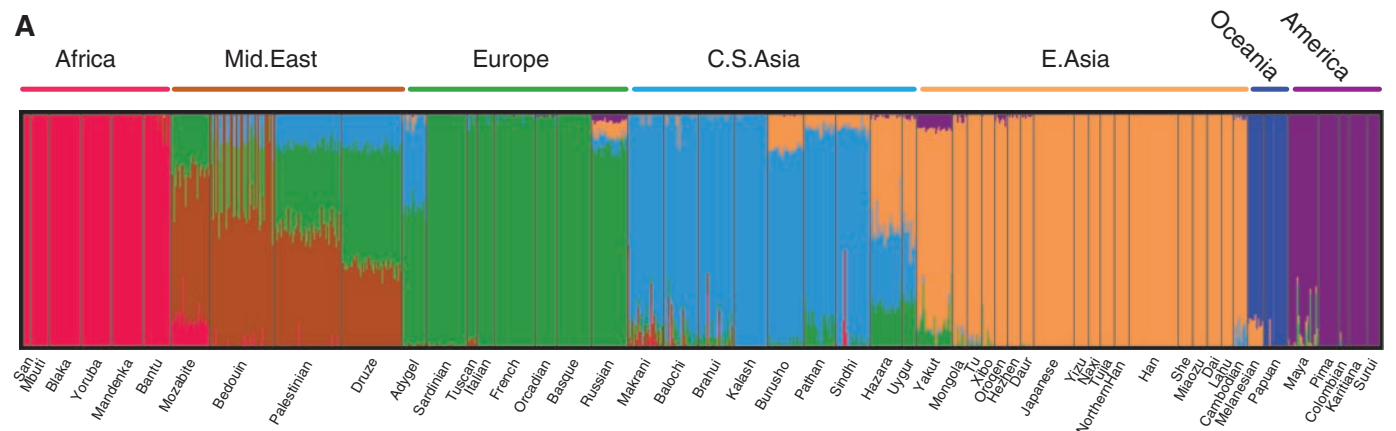
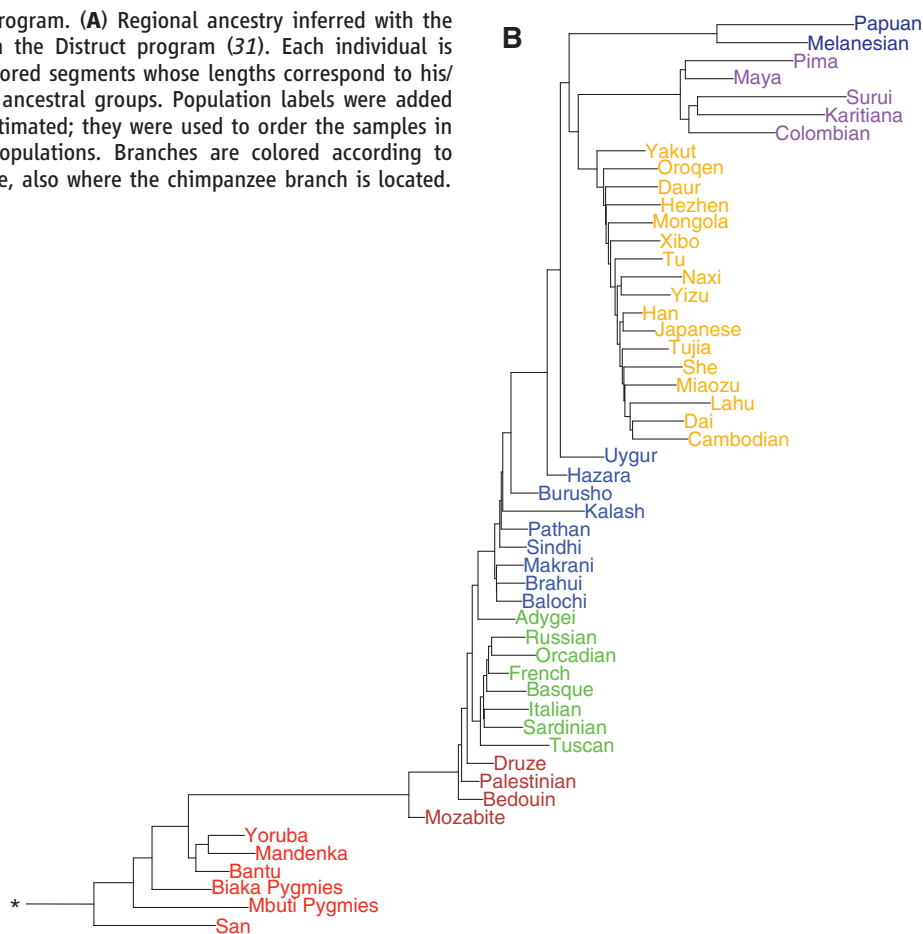


Fig. 1. Individual ancestry and population dendrogram. **(A)** Regional ancestry inferred with the *frappe* program at $K = 7$ (13) and plotted with the *Distruct* program (31). Each individual is represented by a vertical line partitioned into colored segments whose lengths correspond to his/her ancestry coefficients in up to seven inferred ancestral groups. Population labels were added only after each individual's ancestry had been estimated; they were used to order the samples in plotting. **(B)** Maximum likelihood tree of 51 populations. Branches are colored according to continents/regions. * indicates the root of the tree, also where the chimpanzee branch is located.



gence of two populations but without subsequent gene flow between them. For example, the European and Asian ancestries seen in Uyghur and Hazara populations are likely due to relatively recent admixture, whereas the inferred Native American ancestry in Yakuts and Russians likely reflects shared ancestry before the predecessors of the Native Americans crossed the Bering Strait. The Middle Eastern populations may have experienced both continuous gene flow and shared ancestry with the rest of Eurasia.

Because individuals belonging to the same recognized population almost always show similar ancestry proportions (Fig. 1A and fig. S2), it is meaningful to statistically evaluate the genetic relationships among populations. We calculated the Wright's fixation indices F_{st} s among the 51 populations from the population allele frequencies across all autosomal SNPs (15) and constructed a phylogenetic tree by the maximum likelihood method (16), using orthologous chimpanzee alleles as the outgroup. The sub-Saharan African populations are located nearest to the root of the tree (Fig. 1B), outward from which are branches that correspond, sequentially, to pop-

ulations from North Africa, the Middle East, Europe, South/Central Asia, Oceania, America, and East Asia. This population tree shows not only major splits between different continents but also sublineages within continents (14) consistent with the ancestry analysis shown above as well as with results from microsatellite markers (17). The branching pattern largely agrees with the approximate order of human expansion (2) and supports the "out of Africa" model of human origin.

We performed principal component analyses (PCA) on the F_{st} matrix to capture a major portion of genetic variability. The first and second PCs explain 59% and 26% of the F_{st} variation, respectively (fig. S3A) and separate the 51 populations into the known continental groups, with the first PC primarily describing the contrast between sub-Saharan Africans and non-Africans and the second driven by the East-West difference in Eurasia. The third and fourth PCs distinguish the Native American and Oceanian populations, respectively (figs. S3, C and D). The regional clusters are more clearly separated than was possible with 782 microsatellites (16). A PCA plot of the 938 unrelated individuals (fig. S3B) is

similar to the 51-population plot and illustrates the regional clusteredness at the individual level.

The PCA for individual continents/regions clearly delineates fine-scale population structure. In Fig. 2A, the eight European populations, including the central populations (Orcadian, French, Northern Italian from Bergamo, and Tuscan) which were previously indistinguishable with fewer markers, can be separated (3). In Fig. 2B, the four populations from the Middle East are distinguished; the Bedouins can be divided into two subgroups, one of which is similar to the Palestinians. The PC1-PC2 plots for four other continental groups and descriptions and interpretations are in (14) and figs. S4 and S5. These individual-level results, along with ancestry analyses in Fig. 1A and fig. S2, indicate that although some populations have limited sample size (<10), the population structuring appears robust.

We carried out an analysis of molecular variance (AMOVA) (18, 19) to partition overall genetic variation into three components: within-population (WP), among-population-within-group (i.e., geographical region) (AP/WG), and among geographical region (AG). The 51 populations are assigned to the seven geographical regions shown in Fig. 1A. The results are similar among autosomal chromosomes: the WP, AP/WG, and AG components explain $88.9 \pm 0.3\%$, $2.1 \pm 0.05\%$, and $9.0 \pm 0.3\%$ (mean \pm SD across 22 chromosomes) of the variance, respectively (Fig. 3A). For comparison, the WP, AP/WG, and AG components for 783 microsatellite markers are 94.0%, 2.3%, and 3.7%, respectively (3, 5). The difference between the SNP-based estimates (this study) and the microsatellite-based results can be partly explained by higher mutation rates of microsatellites, which are more driven by shorter-term evolutionary processes. For X chromosome (ChrX) SNPs, the WP, AP/WG, and AG components are 84.7%, 2.4%, and 12.9% (Fig. 3A), consistent with estimates based on ChrX microsatellites (20). The greater AG component for ChrX than autosomes is discussed in (14) and fig. S6. Together, these results reaffirm that within-population variation accounts for most of the genetic diversity in humans. However, the between-population variance is sufficient to reveal consistent population structure because subtle but nonrandom differences between populations accumulate over a large number of loci and yield principal components that can account for a major portion of the variation (21).

We compared SNP haplotype heterozygosity across populations and found, consistent with earlier reports (22), that it is highest in sub-Saharan Africa and decreases steadily with distance from this region (Fig. 3B). The mean heterozygosity across autosomal haplotypes (using 295 haplotype blocks in Chr16) (14) is negatively correlated with distance from Addis Ababa, Ethiopia (5, 23), with a correlation coefficient r of -0.91 and a slope of -1.1×10^{-5} per km (Fig. 3B). This trend is consistent with a serial founder effect, a scenario in which population

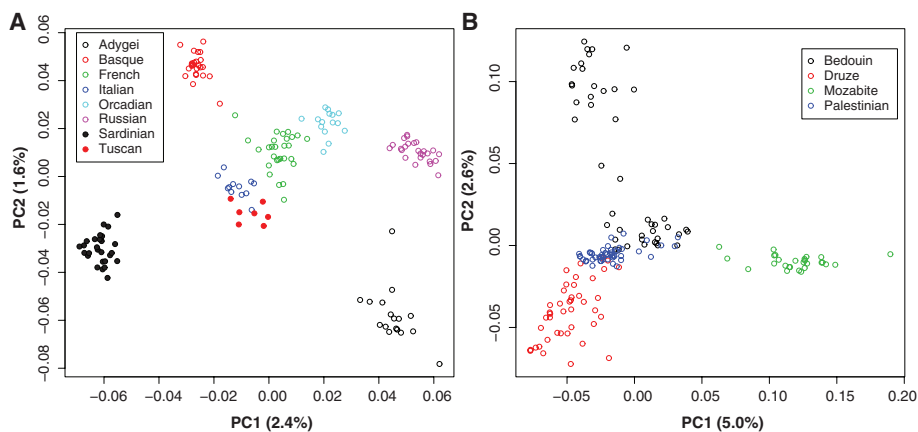


Fig. 2. Fine-scale population structure principal component analyses in two geographic regions, using all autosomal SNPs. (A) Europe. (B) The Middle East.

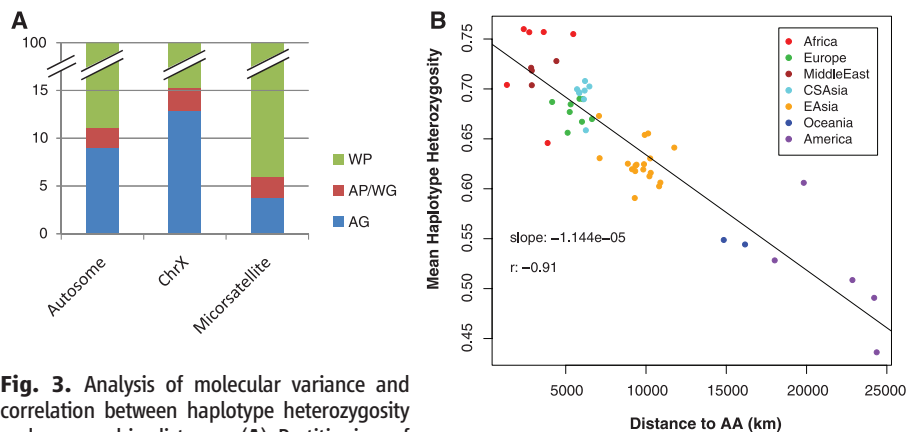


Fig. 3. Analysis of molecular variance and correlation between haplotype heterozygosity and geographic distance. (A) Partitioning of genetic variance into three components (18): Within-Population (WP), Among-Population-Within-Region (AP/WG), and Among-Region (AG), by using autosomal SNPs, microsatellite markers, and ChrX SNPs, respectively. (B) SNP haplotype heterozygosity versus geographic distances from Addis Ababa (AA), Ethiopia. The linear regression slope is indicated along with the Pearson correlation r .

expansion involves successive migration of a small fraction of individuals out of the previous location, starting from a single origin in sub-Saharan Africa. For ChrX haplotypes (using 453 haplotype blocks), the correlation and slope are -0.85 and -1.3×10^{-5} per km, respectively; the slightly higher geographic gradient (i.e., steeper slope) than for autosomes agrees with the higher ChrX Fsts (fig. S6). These values are similar to those reported for microsatellites (5): $r = -0.87$, slope = -6.52×10^{-6} per km, and to SNP-based heterozygosities: $r = -0.81$, slope = -3.8×10^{-6} ($r = -0.93$ when only non-African populations are considered). SNP-based heterozygosities depend on allele frequencies and are affected by ascertainment bias, whereas the haplotype and microsatellite heterozygosities are less affected as a result of their greater polymorphism (22).

By genotyping two chimpanzee samples, we were able to define the putative ancestral allele for $\sim 95.5\%$ of the SNPs in the 650 K panel. We compared the distribution of these ancestral allele frequencies (AAFs) among the 51 populations. Figure 4A shows four examples of the AAF spectrum for Yoruba, French, Chinese, and Japanese populations. Yorubans and other sub-Saharan Africans have more SNPs with high AAFs (>0.6 , on the right of the distribution) and fewer with low AAFs, producing a steeper slope of SNP counts in the midrange of AAF spectrum. The slopes

within 20 to 80% AAF are plotted in Fig. 4B for all 51 populations, showing a progressive reduction moving away from Africa, from ~ 0.04 in sub-Saharan Africa, to ~ 0.03 in Eurasia, ~ 0.02 in East Asia, and ~ 0.01 in Oceania and the Americas. This steady flattening of the AAF distribution may be related to the SNP panel used (which primarily includes common SNPs in Europe, East Asia, and sub-Saharan Africa), but the ascertainment scheme alone cannot explain the entire trend, as the SNPs analyzed by the International HapMap Project show a similar phenomenon (24). In particular, the AAF spectra of ENCODE (ENCyclopedia of DNA Elements) regions, where genotyped SNPs were discovered by re-sequencing, follow a similar pattern (fig. S7), where the 20 to 80% slopes in HapMap Chinese and Japanese populations are about half of that in the Yoruban population.

The flattening of the AAF spectrum reflects the interplay of multiple demographic forces and may yield clues to the history of individual populations. Generally, populations that had a small effective population size and/or experienced a severe bottleneck would have more pronounced genetic drift, resulting in a more rapid increase in derived allele frequencies. Populations that maintained a large size or experienced expansion would tend to preserve the ancestral states of the variant loci. Theoretical work (25), empirical data,

and simulation (26) have shown that demography plays a major role in the change of the AAF spectrum over time. Our result is consistent with the serial founder model, in which non-African populations form a sequential chain of colonies. Those that are more peripheral and younger have relatively smaller effective sizes, and perhaps experienced greater selective pressure. For example, the European and Asian spectra can be explained by a reduction of population size followed by recent recovery (a bottleneck), whereas the spectrum for an African-American population suggests a history of moderate but uninterrupted expansion (26, 27).

Compared to the HapMap panel, HGDP-CEPH includes Oceanian and American populations, as well as a dense collection from the Middle East and South/Central Asia. Characterization of the added populations is important for studying evolution and disease processes not only in these populations but also in those that share common ancestry with them due to recent migration (e.g., U.S. Latino populations). HGDP-CEPH is not a random sample of the world's populations: Some parts of the world (e.g., China and Pakistan) are more densely covered than others (Africa, the Americas, and Oceania). Many populations have been isolated from each other by geography or custom. The observation that they are genetically distinguishable suggests that

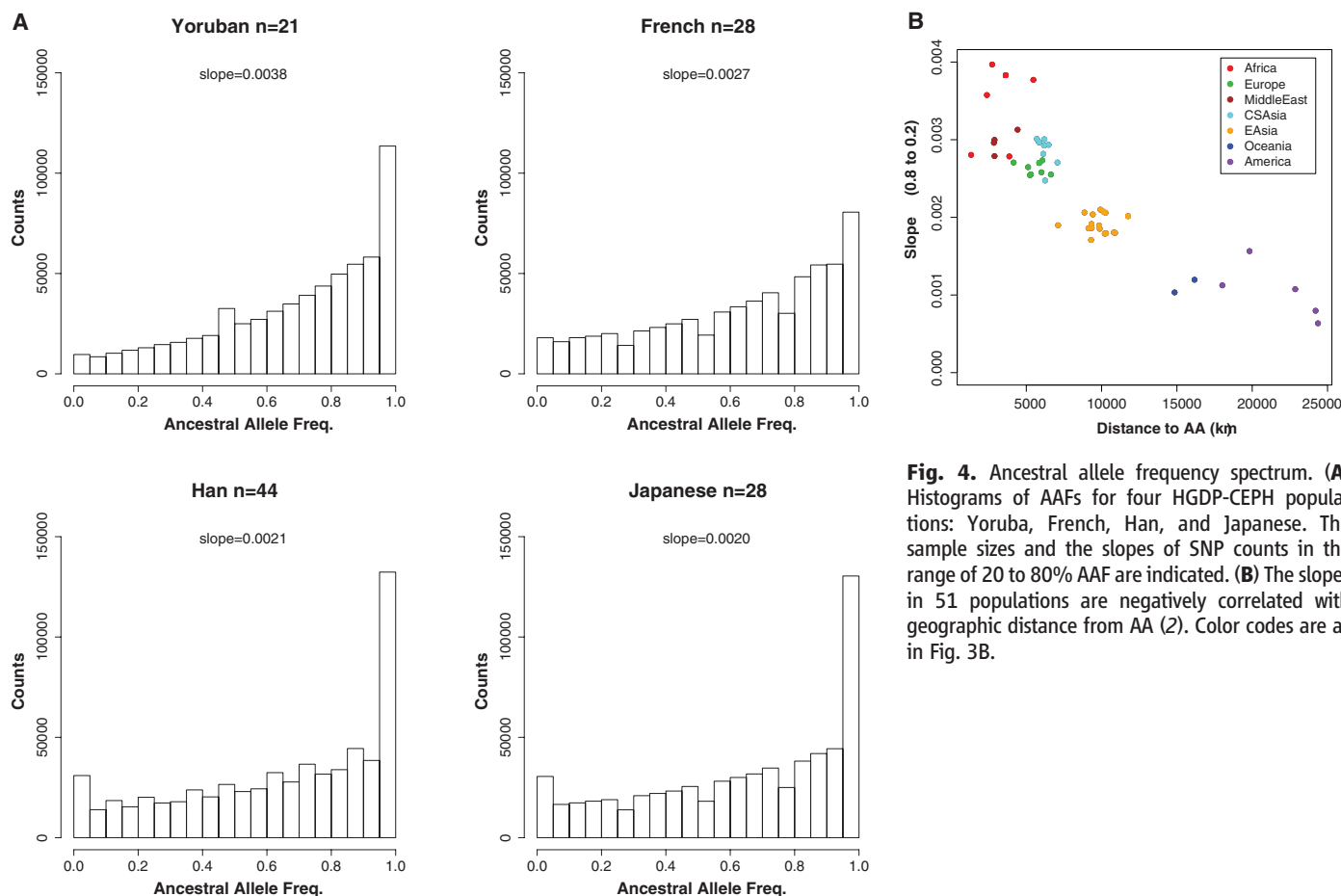


Fig. 4. Ancestral allele frequency spectrum. **(A)** Histograms of AAFs for four HGDP-CEPH populations: Yoruba, French, Han, and Japanese. The sample sizes and the slopes of SNP counts in the range of 20 to 80% AAF are indicated. **(B)** The slopes in 51 populations are negatively correlated with geographic distance from AA (2). Color codes are as in Fig. 3B.

self-reported ancestry is sufficiently accurate for assessing population stratification in disease studies, except for those involving recent admixture (3, 28). These results, however, say nothing about the origin and distribution of human phenotypic variation. The observed population structure can be largely explained by random drift at neutral loci. Nevertheless, some regions of the genome may have experienced accelerated divergence due to local selection (9, 24, 29) as anatomically modern humans spread around the globe during the past 100,000 years, adapting to a wide range of habitats and climates. The population richness of HGDP-CEPH makes it possible to detect correlation between genomic variation and local environmental and/or phenotypic variation (30), thus leading to more detailed understandings of selective forces acting in different regions of the world.

References and Notes

- M. Bamshad *et al.*, *Nat. Rev. Genet.* **5**, 598 (2004).
- L. L. Cavalli-Sforza *et al.*, *Nat. Genet.* **33** (suppl.), 266 (2003).
- N. A. Rosenberg *et al.*, *Science* **298**, 2381 (2002).
- A. M. Bowcock *et al.*, *Nature* **368**, 455 (1994).
- S. Ramachandran *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 15942 (2005).
- N. A. Rosenberg *et al.*, *PLoS Genet.* **1**, e70 (2005).
- C. D. Bustamante *et al.*, *Nature* **437**, 1153 (2005).
- L. L. Cavalli-Sforza, *Annu. Rev. Genomics Hum. Genet.* **8**, 1 (2007).
- B. F. Voight *et al.*, *PLoS Biol.* **4**, e72 (2006).
- S. H. Williamson *et al.*, *PLoS Genet.* **3**, e90 (2007).
- H. M. Cann *et al.*, *Science* **296**, 261 (2002).
- Human Genome Diversity Panel, <http://shgc.stanford.edu/hgdp> or ftp://ftp.cephb.fr/hgdp_suppl1.
- H. Tang *et al.*, *Genet. Epidemiol.* **28**, 289 (2005).
- Materials and methods are available as supporting material on *Science* Online.
- J. Reynolds *et al.*, *Genetics* **105**, 767 (1983).
- J. Felsenstein, *Am. J. Hum. Genet.* **25**, 471 (1973).
- L. A. Zhivotovsky *et al.*, *Am. J. Hum. Genet.* **72**, 1171 (2003).
- L. Excoffier *et al.*, *Genetics* **131**, 479 (1992).
- B. Weir, *Genetic Data Analysis II* (Sinauer, Sunderland, MA, 1996).
- S. Ramachandran *et al.*, *Hum. Genomics* **1**, 87 (2004).
- A. W. Edwards, *Bioessays* **25**, 798 (2003).
- D. F. Conrad *et al.*, *Nat. Genet.* **38**, 1251 (2006).
- F. Prugnolle *et al.*, *Curr. Biol.* **15**, R159 (2005).
- International HapMap Consortium, *Nature* **437**, 1299 (2005).
- M. Kimura *et al.*, *Genetics* **75**, 199 (1973).
- G. T. Marth *et al.*, *Genetics* **166**, 351 (2004).
- A. Keinan *et al.*, *Nat. Genet.* **39**, 1251 (2007).
- H. Tang *et al.*, *Am. J. Hum. Genet.* **76**, 268 (2005).
- S. A. Tishkoff *et al.*, *Nat. Genet.* **39**, 31 (2007).
- A. Manica *et al.*, *Nature* **448**, 346 (2007).
- N. A. Rosenberg, *Mol. Ecol. Notes* **4**, 137 (2004).
- We thank H. Greely for support and helpful discussions, A. Aggarwal and S. Brady for technical assistance, and the Biological Resource Center of the Foundation Jean Dausset-CEPH for preparing and formatting HGDP-CEPH diversity panel DNAs. The BeadArrays used for genotyping the HGDP-CEPH samples were provided by 23andMe and Illumina, Inc. All authors declare no competing financial interest. M.W.F. and H.T. were supported in part by NIH grants GM28016 and GM073059, respectively.

Supporting Online Material

www.sciencemag.org/cgi/content/full/319/5866/1100/DC1
 Materials and Methods
 SOM Text
 Figs. S1 to S7
 Table S1
 References

3 December 2007; accepted 17 January 2008
 10.1126/science.1153717