

Analysis of genomic diversity in Mexican Mestizo populations to develop genomic medicine in Mexico

Irma Silva-Zolezzi¹, Alfredo Hidalgo-Miranda¹, Jesus Estrada-Gil¹, Juan Carlos Fernandez-Lopez, Laura Uribe-Figueroa, Alejandra Contreras, Eros Balam-Ortiz, Laura del Bosque-Plata, David Velazquez-Fernandez, Cesar Lara, Rodrigo Goya, Enrique Hernandez-Lemus, Carlos Davila, Eduardo Barrientos, Santiago March, and Gerardo Jimenez-Sanchez²

National Institute of Genomic Medicine (INMEGEN), Periferico Sur No. 4124, Torre Zafiro II, 6to. Piso, Col. Jardines del Pedregal, Mexico D.F. 01900, Mexico

Communicated by Eric S. Lander, The Broad Institute, Cambridge, MA, March 23, 2009 (received for review March 23, 2008)

Mexico is developing the basis for genomic medicine to improve healthcare of its population. The extensive study of genetic diversity and linkage disequilibrium structure of different populations has made it possible to develop tagging and imputation strategies to comprehensively analyze common genetic variation in association studies of complex diseases. We assessed the benefit of a Mexican haplotype map to improve identification of genes related to common diseases in the Mexican population. We evaluated genetic diversity, linkage disequilibrium patterns, and extent of haplotype sharing using genomewide data from Mexican Mestizos from regions with different histories of admixture and particular population dynamics. Ancestry was evaluated by including 1 Mexican Amerindian group and data from the HapMap. Our results provide evidence of genetic differences between Mexican subpopulations that should be considered in the design and analysis of association studies of complex diseases. In addition, these results support the notion that a haplotype map of the Mexican Mestizo population can reduce the number of tag SNPs required to characterize common genetic variation in this population. This is one of the first genomewide genotyping efforts of a recently admixed population in Latin America.

admixture | genetic variation | population genetics | SNP tagging

More than 560 million people live in Latin American countries, and according to U.S. Census Bureau estimates the Latino population reached ≈ 45.5 million in 2007, representing the largest and fastest-growing minority group in the United States. Mexican Mestizos, as other Latino populations, are a recently admixed population composed of Amerindian, European, and, to a lesser extent, African ancestries. Although the diversity of Latino populations poses several challenges for genetic studies (1), it makes them a powerful resource for analyzing the genetic bases of complex diseases (2). In the past 5 years, Mexico has been committed to develop a human and technological infrastructure for genomics with special emphasis on the development of a national platform of genomic medicine to improve healthcare of Mexicans (3–6). This effort, together with a population of ≈ 105 million inhabitants including 60 Amerindian groups and a complex history of admixture, makes Mexico an ideal country in which to perform genomic analysis of common complex diseases.

Two current approaches to identify genes influencing complex diseases are genomewide association studies (GWAS) and admixture mapping (AM). GWAS depend on efficient SNP tagging (7, 8), and AM on the availability of panels of genomewide markers with frequency differences between parental populations (9, 10). For populations not comprehensively represented in the HapMap (11), such as Latinos, limitations exist for an efficient tagging and imputation, because of the need of a higher number of markers to achieve the same relative power compared to that for Asians and Europeans (12) and the lack of knowledge about population-specific linkage disequilibrium (LD) patterns (13). In addition, false positives because of population structure are minimized in GWAS by excluding individuals with ancestry differences (7). This is not practical in studies including Latinos such as Mexicans, where $>80\%$ of the population consists of Mestizos with known differ-

ences in ancestral proportions (2). As for AM, there are a few SNP panels developed for Latino populations (14–16); however, detailed genomewide information from Mestizo and Amerindian populations remains limited (17, 18). Recent studies of Latin American populations have shown differential ancestral contribution patterns between and within groups that correlate with pre-Columbian native population density and with patterns of recent demographic growth (2). These differences should be considered to improve AM panels for Latin American populations.

Historically, admixture patterns throughout Mexico have been influenced by differences in parental population densities and demographic growth (19–21). Genetic heterogeneity between and within Mestizos from different regions has been documented (22–29). However, no genomewide comparison of different Mestizo and Amerindian populations in Mexico is currently available in the public domain. To analyze genomic diversity and LD patterns in Mexicans, we developed the Mexican Genome Diversity Project (MGDP). This resource will be useful to develop strategies for the genetic analysis of Mexican and related admixed populations, such as marker selection for optimal coverage of common genetic variation in GWA and targeted association studies, and also for the adequate application of tagging and imputation approaches (30, 31) and for AM (10) in Mexicans and other Latino populations. Our study is one of the first extensive genomewide genotyping efforts performed in Latin America. The MGDP will contribute to the development of genomic medicine in Mexico and the rest of Latin America.

Results

We analyzed data from 300 nonrelated self-identified Mestizo individuals from 6 states located in geographically distant regions in Mexico: Sonora (SON) and Zacatecas (ZAC) in the north, Guanajuato (GUA) in the center, Guerrero (GUE) in the center–Pacific, Veracruz (VER) in the center–Gulf, and Yucatan (YUC) in the southeast. Considering that Zapotecos have been shown as a good ancestral population for predicting Amerindian (AMI) ancestry in Mexican Mestizos (16), we included 30 Zapotecos (ZAP) from the southwestern state of Oaxaca (Fig. 1). For comparative purposes, we included similar data sets from HapMap populations: northern Europeans (CEU), Africans (YRI), and East Asians (EA), including Chinese (CHB) and Japanese (JPT). A HapMap-like database with SNP frequencies in Mexicans and HapMap populations was generated (<http://diversity.inmegen.gob.mx>).

Author contributions: I.S.-Z., A.H.-M., J.E.-G., C.L., and G.J.-S. designed research; I.S.-Z., A.H.-M., J.E.-G., L.U.-F., A.C., E.B.-O., L.d.B.-P., D.V.-F., C.L., E.B., S.M., and G.J.-S. performed research; J.E.-G. and C.D. contributed new reagents/analytic tools; I.S.-Z., A.H.-M., J.E.-G., J.C.F.-L., L.U.-F., R.G., E.H.-L., C.D., and G.J.-S. analyzed data; and I.S.-Z., A.H.-M., J.E.-G., and G.J.-S. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

¹I.S.-Z., A.H.-M., and J.E.-G. contributed equally to this work.

²To whom correspondence should be addressed. E-mail: gjimenez@inmegen.gob.mx.

This article contains supporting information online at www.pnas.org/cgi/content/full/0903045106/DCSupplemental.

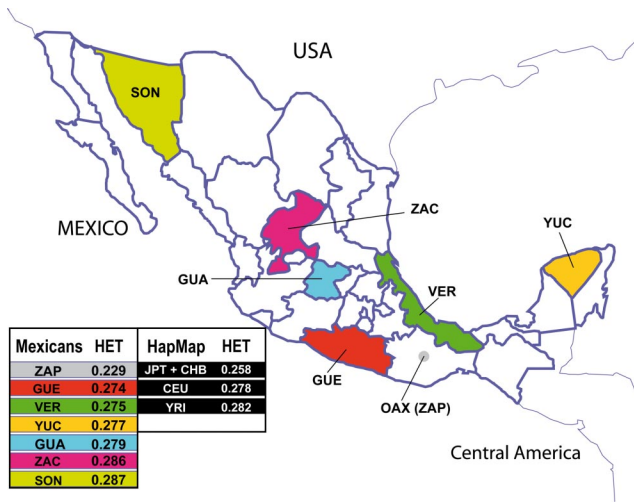


Fig. 1. Genetic diversity measured by heterozygosity (HET) in Mexican and HapMap populations. Northern, central, central-Gulf, central-Pacific, and southern regions in Mexico were included. Average HET values are shown for Amerindian Zapotecos (ZAP), 6 Mexican Mestizo subpopulations (GUA, GUE, SON, VER, YUC, and ZAC), and HapMap populations (YRI, CEU, and JPT + CHB).

Analysis of Genetic Diversity in Mexicans. We measured heterozygosity (HET), performed principal components analysis (PCA) (32), and calculated F_{ST} statistics using data sets obtained for Mexican and HapMap populations. Mexican Mestizo subpopulations had HET values between 0.274 in GUE and 0.287 in SON. Among HapMap populations, YRI displayed the highest genetic diversity (HET = 0.282) and JPT + CHB the lowest (HET = 0.258), as previously reported (33). Among Mexicans, northern subpopulations (SON and ZAC) had the highest HET values, suggesting more genetic diversity, and the ZAP Amerindian samples had the lowest (HET = 0.229), as expected for an isolated population. For PCA analysis, we used different combinations of data sets and conditions. In all scenarios the 2 most informative eigenvectors for each data set are displayed (Fig. 2 A–D). When included, the HapMap and ZAP populations formed defined clusters, while the Mexican Mestizo subpopulations were widely distributed between the CEU and ZAP samples (Fig. 2 A and B). The ZAP population clustering in the PCA plot suggests the absence of recent admixture in this Amerindian group. As expected, when all groups were analyzed (Fig. 2 A), the largest genetic distance exists between the YRI population and the rest of the groups. In the second axis, the ZAP cluster is located between CEU and EA and, in both the first and the second axes, all Mexican Mestizos are spread between CEU and ZAP (Fig. 2 A and B). To better display the distribution of Mexican Mestizos, we generated 2 additional data sets, one leaving out YRI samples (Fig. 2 B) and another including only CEU and ZAP. These analyses gave evidence of genetic diversity between and within Mexican Mestizo populations. In addition, a PCA including only CEU, ZAP, and the 2 Mestizo groups with the largest HET difference (SON and GUE) showed that samples from SON were closer to the CEU, and those from GUE were closer to the ZAP (Fig. 2 C and D). In both plots, some individuals were displaced along eigenvector 2, reflecting additional ancestral contributions in Mestizos. To evaluate whether this effect is related to African (AFR) ancestry, we analyzed an additional data set including YRI [supporting information (SI) Fig. S1 A and B]. The distribution of Mestizos in eigenvector 3 (Fig. S1 B) indicates that the spread observed in eigenvector 2 (Fig. 2 C and D) reflects AFR ancestral contribution. Interestingly, Mestizos did not organize in a straight line between CEU and ZAP (Fig. 2 C and D). This is most probably because those 2 groups of samples do not fully

represent the genetic variability of European and Amerindian ancestral origin present in these Mestizos (2).

To measure genetic distances between Mexican subpopulations, and between these populations and those from the HapMap, we performed a pairwise F_{ST} statistical analysis (Table 1). Of all Mexican groups, the Amerindian ZAP population showed the highest F_{ST} values when compared to all HapMap populations. As expected, the highest value was observed when compared to YRI (23.9), followed by CEU (15.4), JPT (11.9), and CHB (12.0). F_{ST} values between ZAP and each Mestizo subpopulation (Table 1) were consistent with their distribution in the PCA plot (Fig. 2 C), with GUE and VER closest to the ZAP cluster (F_{ST} values 3.2 and 3.8, respectively) and SON at the other end of the distribution (F_{ST} of 8.2). Pairwise comparisons between Mexican groups showed that SON when compared to all other Mestizo subpopulations had higher F_{ST} values than that observed between CHB and JPT. Moreover, the F_{ST} value between SON and ZAP (8.2) was higher than that of any other comparison between any Mestizo subpopulation and non-African HapMap group (Table 1). These results support the presence of considerable genetic heterogeneity between Mexican Mestizo subpopulations and suggest that this diversity is mainly related to a differential distribution of AMI and EUR ancestral components.

To assess genetic ancestry in Mexicans, we determined individual and population average ancestral proportions using STRUCTURE (34, 35). For this, we used 1,814 ancestry informative markers (AIMs) selected using different criteria to ensure genome-wide distribution and minimize LD between SNPs (see *Materials and Methods*). We used HapMap data and the ZAP population as EUR, AFR, EA, and AMI ancestral sources in the analyses. Our results were most consistent with 4 population groups ($K = 4$), explaining the major substructure in this set of Mexican Mestizos (Fig. 3 A and B). In this model, their mean ancestries (\pm SD) were 0.552 ± 0.154 for AMI, 0.418 ± 0.155 for EUR, 0.018 ± 0.035 for AFR, and 0.012 ± 0.018 for EA (Table S1). We observed differences within and between Mestizo subpopulations, mainly in EUR and AMI ancestries (Fig. 3 A and B). The highest and lowest estimates of mean EUR ancestry were 0.616 ± 0.085 for SON and 0.285 ± 0.120 for GUE. Most Mestizo subpopulations displayed statistically significant differences in mean EUR ancestral contribution, and both SON and GUE showed differences when compared to any other Mestizo subpopulation (Table S2). Mestizo groups with similar mean EUR ancestry were those from central and central-coastal regions (VER, YUC, and GUA). In contrast, most Mestizo subpopulations had a similar average AMI ancestral contribution—GUE the highest (0.660 ± 0.138) and SON the lowest (0.362 ± 0.089) (Fig. 3 B)—and only subpopulations in the northern states (SON and ZAC) showed statistically significant differences compared with all other Mestizo groups (Table S2). The other 2 ancestries analyzed, AFR and EA, were smaller and almost homogenous among all Mestizo subpopulations. Significant differences in AFR ancestry were observed for SON and ZAC against VER and YUC (Table S2). To evaluate the contribution of ancestry differences to the overall regional genetic diversity between Mestizo subpopulations, we calculated Pearson correlation coefficients between pairwise F_{ST} values and differences in AMI, EUR, and AFR ancestral proportions. This analysis revealed a high correlation between overall genetic diversity (F_{ST}) and EUR ($r = 0.937$) and AMI ($r = 0.944$) ancestry differences. To estimate the size of this effect, we calculated genetic distance between Mexican subpopulations, specifically attributable to differences in the 2 main continental ancestry proportions (Table S3). This analysis revealed that for most pairwise comparisons between Mestizo subpopulations (10 of 15), 50% of the genetic distance between them is attributable to differences in continental ancestry. Interestingly, most comparisons with low contribution of continental ancestry differences to overall genetic distance included the subpopulation of YUC. These

approach assumes similar genomewide LD patterns between the analyzed samples and the reference panel (30). Tagging or imputation using HapMap information is not as efficient in Mexicans and other Latinos as it is in other populations because of the presence of a genetic component not captured by HapMap data (13). The MGDGP data set will be of great value to test the accuracy of the imputation paradigm in Mexicans and to improve imputation approaches by the inclusion of adequate estimates of individual and local ancestry. The MGDGP data will also be useful to optimize existing sets of AIMs (14–16, 29) to perform AM studies in traits and diseases showing ethnicity-based differences in prevalence in Mexicans, such as HDL cholesterol levels (42), gall bladder disease (43), and type 2 diabetes (44).

We are currently increasing the SNP density to ≈ 1.5 million SNPs per genome using a combination of microarray platforms. Here we present one of the first public genomewide data sets for Mexican Mestizo and Amerindian populations. This effort will contribute to the design of better strategies aimed at characterizing the genetic factors underlying common complex diseases in Mexicans. In addition, this information will increase our knowledge of genomic variability in Latino populations. The scientific and technological infrastructure derived from this project will significantly contribute to the development of genomic medicine in Mexico and Latin America (3, 6).

Materials and Methods

Anonymous blood samples from 300 nonrelated and self-defined Mestizos and 30 Amerindian Zapotecos were collected in 7 states in Mexico: Guanajuato, Guerrero, Sonora, Veracruz, Yucatan, Zacatecas, and Oaxaca (ZAP). The Scientific,

Ethics, and Bio-Security Review Boards from the National Institute of Genomic Medicine (INMEGEN) approved this study. An ad hoc process for community consultation and engagement was implemented. Genomic DNA was extracted from blood (QIAGEN). Genotyping was performed according to the Affymetrix 100K SNP array protocol and 99,953 SNPs passed quality control in all populations. Phasing was performed with fastPhase v1.1.4 (45). All genotypes and raw signal intensity files are available (<ftp://ftp.inmegen.gob.mx>). Average HET was calculated with PLINK (<http://pngu.mgh.harvard.edu/purcell/plink/>) (46). The PCA was done with EIGENSTRAT (32), and F_{ST} with EIGENSOFT (39). Ancestral contributions were assessed with Mann–Whitney U tests, Pearson correlations, box-plot distributions, and their coefficients of variation. For ancestry analysis 1,814 AIMs were used to run STRUCTURE v.2.1 (34, 35). Scripts for informativeness for assignment were kindly provided by N. Rosenberg (37). Alleles private to the Mexican population had a MAF > 0.05 in any of the Mexican subpopulations, but were absent in all HapMap populations. Alleles private to any particular Mexican subpopulation had a MAF > 0.05 in 1 Mexican group and were absent in the other 6. LD calculations, long-range haplotype diversity, and H5 analysis were done with Haploview and special-purpose code, as previously described (47, 48). All data analyses were performed at INMEGEN in Mexico City. (see *SI Materials and Methods*).

ACKNOWLEDGMENTS. We thank the Federal Government of Mexico, particularly the Ministry of Health for valuable support throughout the project. Participation of the governments and universities of the states of Guanajuato, Guerrero, Oaxaca, Sonora, Veracruz, Yucatan, and Zacatecas contributed significantly to this work. We thank all volunteers in the study and the National Institute of Genomic Medicine (INMEGEN)'s personnel for important support; Alejandro López, José Bedolla, Alejandro Rodríguez, and Lucía Orozco for their major contributions to the thorough communication strategy; and Blanca Gonzalez-Sobrinio for helpful advice on Mexican ethnohistory. This work was supported by funds from the Federal Government of Mexico to the National Institute of Genomic Medicine and by infrastructure donated by the Mexican Health Foundation (FUNSALUD) and the Gonzalo Río Arronte Foundation.

- Gonzalez Burchard E, et al. (2005) Latino populations: A unique opportunity for the study of race, genetics, and social environment in epidemiological research. *Am J Public Health* 95:2161–2168.
- Wang S, et al. (2008) Geographic patterns of genome admixture in Latin American Mestizos. *PLoS Genet* 4:e1000037.
- Jimenez-Sanchez G (2003) Developing a platform for genomic medicine in Mexico. *Science* 300:295–296.
- Hardy BJ, et al. (2008) The next steps for genomic medicine: Challenges and opportunities for the developing world. *Nat Rev Genet* 9(Suppl 1):S23–S27.
- Seguin B, Hardy BJ, Singer PA, Daar AS (2008) Genomics, public health and developing countries: The case of the Mexican National Institute of Genomic Medicine (INMEGEN). *Nat Rev Genet* 9(Suppl 1):S5–S9.
- Jimenez-Sanchez G, Silva-Zolezzi I, Hidalgo A, March S (2008) Genomic medicine in Mexico: Initial steps and the road ahead. *Genome Res* 18:1191–1198.
- The Wellcome Trust Case Control Consortium (2007) Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* 447:661–678.
- McCarthy MI, et al. (2008) Genome-wide association studies for complex traits: Consensus, uncertainty and challenges. *Nat Rev Genet* 9:356–369.
- Smith MW, O'Brien SJ (2005) Mapping by admixture linkage disequilibrium: Advances, limitations and guidelines. *Nat Rev Genet* 6:623–632.
- Seldin MF (2007) Admixture mapping as a tool in gene discovery. *Curr Opin Genet Dev* 17:177–181.
- The International HapMap Consortium (2005) A haplotype map of the human genome. *Nature* 437:1299–1320.
- de Bakker PI, et al. (2006) Transferability of tag SNPs in genetic association studies in multiple populations. *Nat Genet* 38:1298–1303.
- Huang L, et al. (2009) Genotype-imputation accuracy across worldwide human populations. *Am J Hum Genet* 84:235–250.
- Mao X, et al. (2007) A genomewide admixture mapping panel for Hispanic/Latino populations. *Am J Hum Genet* 80:1171–1178.
- Tian C, et al. (2007) A genomewide single-nucleotide-polymorphism panel for Mexican American admixture mapping. *Am J Hum Genet* 80:1014–1023.
- Price AL, et al. (2007) A genomewide admixture map for Latino populations. *Am J Hum Genet* 80:1024–1036.
- Li JZ, et al. (2008) Worldwide human relationships inferred from genome-wide patterns of variation. *Science* 319:1100–1104.
- Jakobsson M, et al. (2008) Genotype, haplotype and copy-number variation in world-wide human populations. *Nature* 451:998–1003.
- Gerhard P (1986) *Historical Geography of New Spain, 1519–1821* (Spanish) (Universidad Nacional Autónoma de México, Mexico City).
- Gerhard P (1991) *La Frontera Sureste de la Nueva España* (Universidad Nacional Autónoma de México, Mexico City) (in Spanish).
- Gerhard P (1996) *La Frontera Norte de la Nueva España* (Universidad Nacional Autónoma de México, Mexico City) (in Spanish).
- Buentello-Malo L, Penalzoza-Espinosa RI, Loeza F, Salamanca-Gomez F, Cerda-Flores RM (2003) Genetic structure of seven Mexican indigenous populations based on five polymarker loci. *Am J Hum Biol* 15:23–28.
- Cerda-Flores RM, et al. (1992) Gene diversity and estimation of genetic admixture among Mexican-Americans of Starr County, Texas. *Ann Hum Biol* 19:347–360.
- Cerda-Flores RM, et al. (2002) Genetic admixture in three Mexican Mestizo populations based on D1S80 and HLA-DQA1 loci. *Am J Hum Biol* 14:257–263.
- De Leo C, et al. (1997) HLA class I and class II alleles and haplotypes in Mexican Mestizos established from serological typing of 50 families. *Hum Biol* 69:809–818.
- Gorodetzky C, et al. (2001) The genetic structure of Mexican Mestizos of different locations: Tracking back their origins through MHC genes, blood group systems, and microsatellites. *Hum Immunol* 62:979–991.
- Lisker R, et al. (1986) Gene frequencies and admixture estimates in a Mexico City population. *Am J Phys Anthropol* 71:203–207.
- Lisker R, Ramirez E, Briceno RP, Granados J, Babinsky V (1990) Gene frequencies and admixture estimates in four Mexican urban centers. *Hum Biol* 62:791–801.
- Martinez-Marignac VL, et al. (2007) Admixture in Mexico City: Implications for admixture mapping of type 2 diabetes genetic risk factors. *Hum Genet* 120:807–819.
- Marchini J, Howie B, Myers S, McVean G, Donnelly P (2007) A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet* 39:906–913.
- Zeggini E, et al. (2008) Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nat Genet* 40:638–645.
- Patterson N, Price AL, Reich D (2006) Population structure and eigenanalysis. *PLoS Genet* 2:e190.
- Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R (2005) Ascertainment bias in studies of human genome-wide polymorphism. *Genome Res* 15:1496–1502.
- Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics* 164:1567–1587.
- Falush D, Stephens M, Pritchard JK (2007) Inference of population structure using multilocus genotype data: Dominant markers and null alleles. *Mol Ecol Notes* 7:574–578.
- Ramachandran S, et al. (2005) Support from the relationship of genetic and geographic distance in human populations for a serial founder effect originating in Africa. *Proc Natl Acad Sci USA* 102:15942–15947.
- Rosenberg NA, Li LM, Ward R, Pritchard JK (2003) Informativeness of genetic markers for inference of ancestry. *Am J Hum Genet* 73:1402–1422.
- Rangel-Villalobos H, et al. (2008) Genetic admixture, relatedness, and structure patterns among Mexican populations revealed by the Y-chromosome. *Am J Phys Anthropol* 135:448–461.
- Aguirre-Beltran G, ed (1972) *La Población Negra de México: Estudio Etnográfico* (Fondo de Cultura Económica, Mexico City) (in Spanish).
- Matsuzaki H, et al. (2004) Genotyping over 100,000 SNPs on a pair of oligonucleotide arrays. *Nat Methods* 1:109–111.
- Conrad DF, et al. (2006) A worldwide survey of haplotype variation and linkage disequilibrium in the human genome. *Nat Genet* 38:1251–1260.
- Cosrow N, Falkner B (2004) Race/ethnic issues in obesity and obesity-related comorbidities. *J Clin Endocrinol Metab* 89:2590–2594.
- Everhart JE, et al. (2002) Prevalence of gallbladder disease in American Indian populations: Findings from the Strong Heart Study. *Hepatology* 35:1507–1512.
- Hamman RF, et al. (1989) Methods and prevalence of non-insulin-dependent diabetes mellitus in a biethnic Colorado population. The San Luis Valley Diabetes Study. *Am J Epidemiol* 129:295–311.
- Scheet P, Stephens M (2006) A fast and flexible statistical model for large-scale population genotype data: Applications to inferring missing genotypes and haplotypic phase. *Am J Hum Genet* 78:629–644.
- Purcell S, et al. (2007) PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 81:559–575.
- Bonnen PE, et al. (2006) Evaluating potential for whole-genome studies in Kosrae, an isolated population in Micronesia. *Nat Genet* 38:214–217.
- Barrett JC, Fry B, Maller J, Daly MJ (2005) Haploview: Analysis and visualization of LD and haplotype maps. *Bioinformatics* 21:263–265.