

# Long-Range LD Can Confound Genome Scans in Admixed Populations

*To the Editor:* In the September 2007 issue of *The Journal*, Tang et al. analyzed data from 192 Puerto Ricans genotyped at 112,584 autosomal markers and identified three regions with a deficiency in the proportion of European ancestry. They concluded that recent selection occurred at these regions after the admixture of European, African, and Native American ancestors.<sup>1</sup> These signals of selection are very strong: We estimate that they each correspond to selection coefficients of  $>0.08$  per generation, which if confirmed would represent the three most powerful selective adaptations discovered to date in humans. Here, we demonstrate that on the basis of the method the authors applied, these signals of selection could be explained as artifacts of the unusual long-range linkage disequilibrium (LD) that occurs at these regions and that is not specific to Puerto Ricans. We failed to replicate the signal of selection in an independent and larger study of 364 Puerto Rican samples, when we applied a method that is not susceptible to this confounder. Our results highlight a complexity in the analysis of dense genotype data from recently admixed populations; this complexity needs to be taken into account not only in genome-wide screens for selection but also in genome-wide association studies to ensure that false-positive signals are avoided.

The signals of selection were identified with methods described in Tang et al.,<sup>2</sup> which uses an extension of a Hidden Markov Model (HMM) to infer segments of ancestry from dense genotype data. The authors note that the assumptions of an HMM “are violated when the marker map is dense and linkage disequilibrium (LD) exists within an ancestral population”; they partially address this confounder by modeling the LD between consecutive pairs of markers but describe this approach as a “compromise” because they do not account for higher order LD.<sup>2</sup> In light of the phenomenon that nearby sites in a region may be in weak LD, whereas more distant sites may be in much stronger LD, the approach of modeling only LD between consecutive markers is potentially inadequate.<sup>3</sup> As we demonstrate below, local-ancestry estimates in regions where LD is not fully modeled will not only be overconfident but will also be systematically biased, thereby leading to false-positive deficiencies in the population contributing majority ancestry.

In a separate analysis focusing on long-range LD in European populations, we applied principal components analysis (PCA) to several genome-wide data sets and identified 24 autosomal long-range LD regions, each spanning  $>2$  megabases (Mb) (Table 1). The functional basis for these regions is currently being explored. The 24 PCA regions were identified by running the EIGENSOFT software<sup>4,5</sup> on a data set of 327 European Americans genotyped on

the Illumina 550K array and identifying all regions where there was significant long-range LD extending  $>2$  Mb that explained one of the top eigenvectors. The regions were independently replicated in 1593 European Americans from the Illumina iControl data set genotyped on the Illumina 550K array and in 1504 + 1500 British samples from the Wellcome Trust Case Control Consortium (1958 Birth Cohort and National Blood Service Cohorts, genotyped on the Affymetrix 500K array), confirming that these regions genuinely harbor long-range LD in European populations.

Strikingly, all three of the signals of selection reported by Tang et al.<sup>1</sup> lie in one of the PCA regions (Table 1). Because the PCA regions comprise  $<4.7\%$  of the autosomal genome, the hypothesis that the regions discussed in Tang et al.<sup>1</sup> and the PCA regions are independent is violated with a  $p$  value of  $(0.047)^3 = 0.0001$ . As we will show, the presence of long-range LD in populations ancestral to Puerto Ricans could explain both the signals from Tang et al.<sup>1</sup> and the PCA results.

Long-range LD can arise for reasons unrelated to selection. For example, inversions are known to suppress viable recombination, and a known inversion polymorphism at position 8–12 Mb on chromosome 8 has previously been shown to be the cause of long-range LD<sup>6</sup> (also see Table 1). (Interestingly, this inversion polymorphism appears to produce a signal of unusual ancestry in Figure 1 of Tang et al.,<sup>1</sup> in addition to the three regions highlighted in the same paper.) It is important for studies inferring the action of selection to rule out alternative explanations for the observed data. For the regions identified by Tang et al.,<sup>1</sup> long-range LD that arose because of inversion polymorphism or other reasons provides a plausible alternative explanation.

LD that is not properly modeled impacts not only the uncertainty in local-ancestry estimates but also the expected value of these estimates, leading to large systematic biases in regions of long-range LD. To demonstrate this, we consider a hypothetical admixed population with ancestry  $\alpha_1 = 80\%$  from ancestral population 1 and  $\alpha_2 = 20\%$  from ancestral population 2. We then consider an A/C marker in which the A allele has frequency  $p_1 = 25\%$  in population 1 and  $p_2 = 75\%$  in population 2, so that its frequency in the admixed population is  $p = \alpha_1 p_1 + \alpha_2 p_2 = 35\%$ . Let  $q_1 = 75\%$ ,  $q_2 = 25\%$ , and  $q = 65\%$  denote the corresponding frequencies of the C allele. If local ancestry on a single-haploid chromosome is inferred with only information from that marker, we obtain  $P(\text{population 1} | A) = \alpha_1 p_1 / (\alpha_1 p_1 + \alpha_2 p_2) = 0.57$  and  $P(\text{population 1} | C) = \alpha_1 q_1 / (\alpha_1 q_1 + \alpha_2 q_2) = 0.92$ , so that the expected value of the ancestry estimate is  $E(P(\text{population 1})) = p P(\text{population 1} | A) + q P(\text{population 1} | C) = 0.80$ , which is an unbiased estimate of  $\alpha_1$ . Now, we consider a second marker that has identical allele frequencies and that is in perfect LD with the first and suppose that the two markers are used to infer local ancestry, treating them as if they were unlinked (this could happen with the method of Tang et al.<sup>2</sup> if the markers

**Table 1. Correspondence between Regions from Tang et al. and Regions of Extended LD in European Populations**

Chromosome	SNP at Region Peak, from Tang et al. <sup>1</sup>	SNP Position from PCA Analysis	Extended LD Region, Mb
6	rs169679	29.0 Mb	25.5–33.5 Mb
8	rs896760	113.5 Mb	112–115 Mb
11	rs637249	56.0 Mb	46–57 Mb

For each region reported to be under selection, we list the SNP defining the peak of this region as described in Tang et al.,<sup>1</sup> the physical position of the SNP, and the physical position of the corresponding region of extended LD from PCA analysis. The other autosomal long-range LD regions identified by PCA analysis were chromosome 1: 48–52 Mb, 2: 86–100.5 Mb, 2: 134.5–138 Mb, 2: 183–190 Mb, 3: 47.5–50 Mb, 3: 83.5–87 Mb, 3: 89–97.5 Mb, 5: 44.5–50.5 Mb, 5: 98–100.5 Mb, 5: 129–132 Mb, 5: 135.5–138.5 Mb, 6: 57–64 Mb, 6: 140–142.5 Mb, 7: 55–66 Mb, 8: 8–12 Mb, 8: 43–50 Mb, 10: 37–43 Mb, 11: 87.5–90.5 Mb, 12: 33–40 Mb, 12: 109.5–112 Mb, and 20: 32–34.5 Mb.

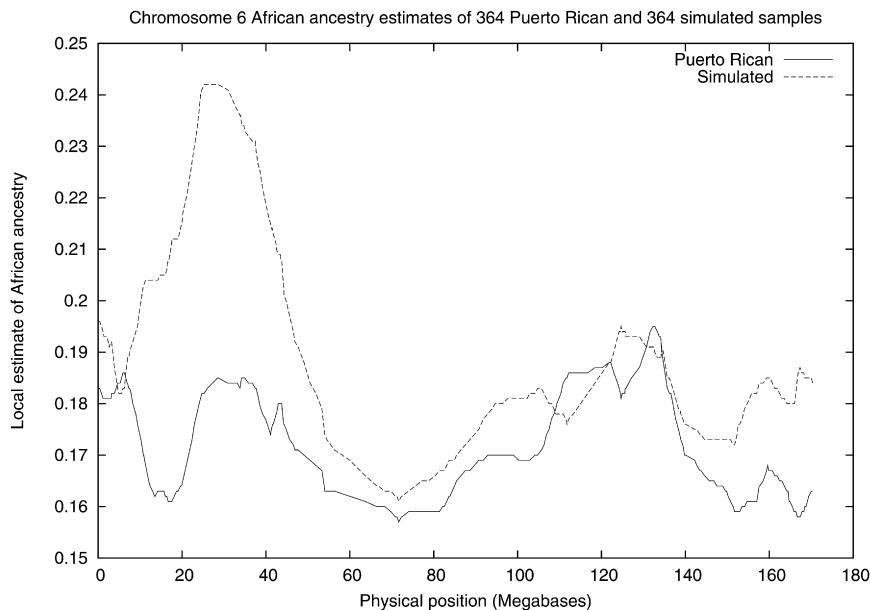
are nonconsecutive). The resulting local-ancestry estimates are  $P(\text{population 1|AA}) = \alpha_1 p_1^2 / (\alpha_1 p_1^2 + \alpha_2 p_2^2) = 0.31$  and  $P(\text{population 1|CC}) = \alpha_1 q_1^2 / (\alpha_1 q_1^2 + \alpha_2 q_2^2) = 0.97$ , so that the expected value of the ancestry estimate is  $E(P(\text{population 1})) = p P(\text{population 1|AA}) + q P(\text{population 1|CC}) = 0.74$ , a downwardly biased estimate of  $\alpha_1$ . More generally, when  $n$  perfectly linked markers are used to infer ancestry and are treated as unlinked, for large  $n$  (e.g.,  $n \geq 5$ ), the evidence of ancestry associated to a particular allele becomes overwhelming, and the estimated ancestry proportion will equal the allele frequency:  $P(\text{population 1|A}^n) = \alpha_1 p_1^n / (\alpha_1 p_1^n + \alpha_2 p_2^n) \approx 0$  and  $P(\text{population 1|C}^n) = \alpha_1 q_1^n / (\alpha_1 q_1^n + \alpha_2 q_2^n) \approx 1$ , so that  $E(P(\text{population 1})) = p P(\text{population 1|A}^n) + q P(\text{population 1|C}^n) = q = 0.65$ . The deficiency of 15% local ancestry, compared to genome-wide ancestry of 80%, shows that the bias could produce effects as large as the 14% deficiencies in European ancestry reported by Tang et al.<sup>1</sup>; such deficiencies will persist when local-ancestry estimates are incorporated into an HMM. In a data set of 112,584 markers, the regions of long-range LD listed in Table 1 would be expected to contain at least 100 markers each. As in our example, unmodeled LD could bias ancestry estimates in the direction of allele frequencies, thereby favoring a deficiency of the population contributing majority ancestry—just as reported in Tang et al.<sup>1</sup>

In addition to their analysis of 112,584 markers, Tang et al.<sup>1</sup> report evidence of selection in analyses of individual HLA markers (Table 1 of their paper). These single-marker analyses are immune to the effects of long-range LD but may be affected by their use of inaccurate ancestral populations to model Puerto Rican ancestry. In particular, the Native American ancestry of Puerto Ricans derives from the Taino, a Native South American population that is likely to be highly genetically diverged from the Native North American populations such as the Pima and Maya used by Tang et al.<sup>1</sup> to model Native American ancestry.<sup>7</sup> Frequency differences among Native American populations could explain why Table 1 of Tang et al.<sup>1</sup> reports

a 13% increase in Native American ancestry based on allele frequencies of individual markers at the HLA locus, whereas Figure 1 of Tang et al.<sup>1</sup> reports no deviation in Native American ancestry at the same locus when flanking genomic data were used.<sup>2</sup> We note that if single-marker analyses are affected by the use of inaccurate ancestral populations, analyses of individual markers in new samples from the same populations would not provide an independent replication because the genetic drift underlying the inaccuracy occurs at the population level, not at the individual level.

As an independent test for selection at the chromosome 6 locus, we analyzed 364 new Puerto Rican samples, consisting of 170 individuals with Crohn's disease and 194 matched controls recruited at the University of Puerto Rico School of Medicine. We genotyped these samples at 2459 autosomal markers from our published admixture map that were powerful for distinguishing African from non-African ancestry.<sup>8</sup> (Most markers in the map have relatively similar frequencies in Europeans and Native Americans, with very different frequencies in Africans.) Genotyping was performed with the Illumina Golden Gate technology, and standard quality filters were applied.<sup>9</sup> After additional filtering to exclude markers that were highly differentiated between Europeans and Native Americans (so as to ensure an effective two-way African versus non-African admixture analysis in a three-way admixed population<sup>7</sup>) and disallow LD between markers in the ancestral populations,<sup>10</sup> we retained 1438 markers for downstream analysis. We found that these markers were sufficient to generate useful ancestry estimates: Our calculations indicate that we capture 61% of maximum information about African versus non-African ancestry at the chromosome 6 region, so our effective sample size is  $(0.61)(364) = 223$ , which is larger than the sample size of 192 in Tang et al.<sup>1</sup>

By using the ANCESTRYMAP software<sup>1</sup> to obtain local-ancestry estimates, we failed to replicate the finding of Tang et al.<sup>1</sup> of an increase in African ancestry at chromosome 6 (Figure 1) and did not observe an unusual deviation in ancestry at any region of the genome. (These results do not shed light on selection signals at the chromosome 8 and 11 regions because Tang et al.<sup>1</sup> reported deviations in European and Native American ancestry at these loci, whereas our 1,438 markers only distinguish African versus non-African ancestry.) To test whether our negative result could be a consequence of low power, we simulated a data set of 364 samples from an admixed population that has 18% African ancestry genome wide but 32% at the chromosome 6 region.<sup>1</sup> In detail, we simulated samples by generating ancestry segments and genotypes at the same set of 1438 markers (with the same pattern of missing data as our Puerto Rican samples) assuming 18% African ancestry, 82% European ancestry, and an average of nine generations since admixture (This quantity was inferred from the Puerto Rican data and is similar to values for other Latino populations.<sup>7</sup>). We preferentially selected samples



**Figure 1. A Replication Study in 364 Puerto Ricans Finds No Significant Rise in African Ancestry at the Chromosome 6 Locus**

Local estimates of percent African ancestry on chromosome 6 for 364 Puerto Rican samples and the same number of samples from a hypothetical admixed population simulated to have unusually high African ancestry at the chromosome 6 region centered at position 29.0 Mb as reported in Tang et al.<sup>1</sup> Local ancestry was estimated by ANCESTRYMAP with unlinked markers. We note that the Puerto Rican samples in our study show a slight peak at this region, but this is not significant because there are 41 larger peaks of African ancestry elsewhere in the genome. In contrast, the simulated samples show an excess of African ancestry at this locus, and this is more than twice as large as is observed anywhere else in the genome.

with African ancestry at marker rs451774 (position 28.6 Mb on chromosome 6) so as to achieve 32% African ancestry at this locus. By running ANCESTRYMAP on 364 simulated samples, we detected a large rise in African ancestry at the chromosome 6 region (Figure 1). Although the local estimate of 24% African ancestry at this region is less than the value of 32% used to simulate the data (because ANCESTRYMAP assumes the null model of no unusual deviation in local ancestry and thus imposes a strong prior of 18% African ancestry), the excess of African ancestry is more than twice what is observed anywhere else in the genome. Thus, our failure to identify a rise in African ancestry in Puerto Rican samples on chromosome 6 is not due to a lack of power.

To test the robustness of our negative result, we reran our analysis of the 364 Puerto Rican samples with marker sets chosen to have different thresholds for maximum differentiation between Europeans and Native Americans and reran with all African and European allele-frequency data omitted to ensure that our results were not affected by inaccurate ancestral populations. We also reran with the control individuals only, to ensure that our results were not influenced by the inclusion of Crohn's disease cases. In none of these runs did we observe a signal of a rise in African ancestry at the chromosome 6 locus. The above runs used markers that are not in LD in ancestral populations, as required by ANCESTRYMAP. However, as a demonstration of the pitfalls of not accounting for LD between markers, we reran ANCESTRYMAP on a larger set of 1852 markers in which no constraint was applied to disallow LD in ancestral populations. African-ancestry estimates across the genome varied wildly from 15% to 54%, corresponding to large deficiencies in European ancestry analogous to the signals from Tang et al.<sup>1</sup>

Our analysis demonstrates that the signals of recent selection reported by Tang et al.<sup>1</sup> could theoretically be explained as artifacts caused by regions of long-range LD (with which they strikingly coincide) and inaccurate ancestral populations. Furthermore, we empirically failed to replicate the finding of an unusual deviation in African ancestry at the chromosome 6 region in our analysis of a larger Puerto Rican sample set. We believe that the hypothesis of selection since admixture should therefore be viewed with caution. We note that in a joint analysis of more than 10,000 African American samples that we have scanned in admixture-mapping studies, we have not yet found a single locus at which there is signal of a local-ancestry deviation that is not specific to disease cases. We consider it unlikely that recent selection events could lead to three distinct local-ancestry deviations that are large enough to be detected with only 192 Puerto Rican samples, when we failed to detect any such effect in African Americans using >50-fold more samples.

These results also have methodological significance for genome-wide association studies in admixed populations such as Latinos and African Americans. To have maximum power, such studies need to take advantage of admixture association signals (deviations in local ancestry in disease cases compared to their genome-wide average) as well as case-control association signals. The method of Tang et al.<sup>2</sup> has been shown to accurately infer ancestry in simulated data sets, but our results suggest that it may produce false-positive admixture association signals in regions of long-range LD in admixed populations. In association studies, such errors can be controlled by computation of local-ancestry estimates in both cases and controls. However, case-only admixture association analyses are known to provide higher statistical power.<sup>11</sup> Thus, carrying out robust, fully powered genome-wide association studies in admixed

populations will require methods that rigorously account for the confounding effects of long-range LD.

Alkes L. Price,<sup>1,2</sup> Michael E. Weale,<sup>3</sup> Nick Patterson,<sup>2</sup> Simon R. Myers,<sup>2,4</sup> Anna C. Need,<sup>3</sup> Kevin V. Shianna,<sup>3</sup> Dongliang Ge,<sup>3</sup> Jerome I. Rotter,<sup>5</sup> Esther Torres,<sup>6</sup> Kent D. Taylor,<sup>5</sup> David B. Goldstein,<sup>3</sup> and David Reich<sup>1,2,\*</sup>

<sup>1</sup>Department of Genetics, Harvard Medical School, Boston, MA 02115, USA; <sup>2</sup>Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA; <sup>3</sup>Institute for Genome Sciences and Policy, Duke University, Durham, NC 27710, USA; <sup>4</sup>Department of Statistics, University of Oxford, Oxford OX1 3TG, UK; <sup>5</sup>Cedars-Sinai Medical Center, University of California Los Angeles, Los Angeles, CA 90048, USA; <sup>6</sup>University of Puerto Rico School of Medicine San Juan, PR 00936, USA

\*Correspondence: [reich@genetics.med.harvard.edu](mailto:reich@genetics.med.harvard.edu)

## Acknowledgments

A.L.P. is supported by a Ruth Kirschstein National Research Service Award from the NIH. N.P. is supported by a K-01 career development award from the NIH. D.R. is supported by a Burroughs Wellcome Career Development Award in the Biomedical Sciences. This research was also supported by U-01 award HG004168 from the NIH (D.R. and N.P.), by NIDDK grant PO1DK46763 (J.I.R.), and by the Board of Governor's Chair in Medical Genetics at Cedars-Sinai Medical Center (J.I.R.). Genotyping of the Puerto Rican samples was supported in part by grant M01-RR00425 to the Cedars-Sinai GCRC genotyping core (K.D.T.) and by NIH grant DK62413 (K.D.T.).

## References

1. Tang, H., Choudhry, S., Mei, R., Morgan, M., Rodriguez-Clin-ton, W., Burchard, E.G., and Risch, N.J. (2007). Recent genetic selection in the ancestral admixture of Puerto Ricans. *Am. J. Hum. Genet.* 81, 626–633.

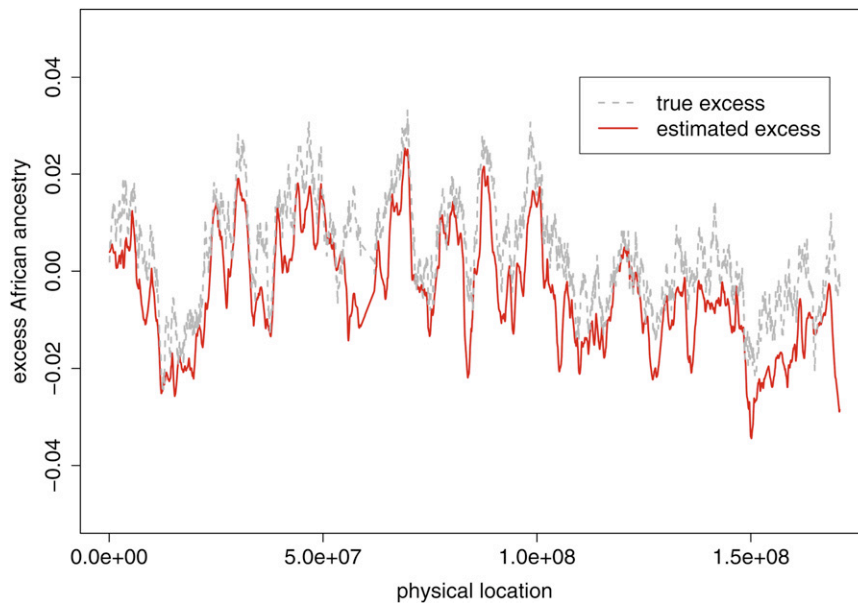
2. Tang, H., Coram, M., Wang, P., Zhu, X., and Risch, N. (2006). Reconstructing genetic ancestry blocks in admixed individuals. *Am. J. Hum. Genet.* 79, 1–12.
3. Wall, J.D., and Pritchard, J.K. (2003). Linkage disequilibrium in the human genome. *Nat. Rev. Genet.* 4, 587–597.
4. Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A., and Reich, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* 38, 904–909.
5. Patterson, N., Price, A.L., and Reich, D. (2006). Population structure and eigenanalysis. *PLoS Genet* 2, e190.
6. Tian, C., Plenge, R.M., Ransom, M., Lee, A., Villoslada, P., Selmi, C., Klareskog, L., Pulver, A.E., Qi, L., Gregersen, P.K., et al. (2008). Analysis and application of European genetic substructure using 300 K SNP information. *PLoS Genet* 4, e4.
7. Price, A.L., Patterson, N., Yu, F., Cox, D.R., Waliszewska, A., McDonald, G.J., Tandon, A., Schirmer, C., Neubauer, J., Bed-oya, G., et al. (2007). A genomewide admixture map for Latino populations. *Am. J. Hum. Genet.* 80, 1024–1036.
8. Smith, M.W., Patterson, N., Lautenberger, J.A., Truelove, A.L., McDonald, G.J., Waliszewska, A., Kessing, B.D., Malasky, M.J., Scafe, C., Le, E., et al. (2004). A high-density admixture map for disease gene discovery in African Americans. *Am. J. Hum. Genet.* 74, 1001–1013.
9. Fan, J.B., Oliphant, A., Shen, R., Kermani, B.G., Garcia, F., Gunderson, K.L., Hansen, M., Steemers, F., Butler, S.L., Deloukas, P., et al. (2003). Highly parallel SNP genotyping. *Cold Spring Harb. Symp. Quant. Biol.* 68, 69–78.
10. Reich, D., Patterson, N., De Jager, P.L., McDonald, G.J., Waliszewska, A., Tandon, A., Lincoln, R.R., DeLoa, C., Fruhan, S.A., Cabre, P., et al. (2005). A whole-genome admixture scan finds a candidate gene for multiple sclerosis susceptibility. *Nat. Genet.* 37, 1113–1118.
11. Patterson, N., Hattangadi, N., Lane, B., Lohmueller, K.E., Hafler, D.A., Oksenberg, J.R., Hauser, S.L., Smith, M.W., O'Brien, S.J., Altshuler, D., et al. (2004). Methods for high-density admixture mapping of disease genes. *Am. J. Hum. Genet.* 74, 979–1000.

DOI 10.1016/j.ajhg.2008.06.005. ©2008 by The American Society of Human Genetics. All rights reserved.

## Response to Price et al.

*To the Editor:* In 2006, Tang and colleagues<sup>1</sup> presented a novel statistical method for genetic admixture analysis based on high-density SNP arrays rather than conventional ancestry informative markers (AIMs). The chromosomes of an admixed individual represent a consecutive patchwork of ancestry blocks representing the ancestral populations contributing to the admixed individual. Their approach<sup>1</sup> is based on the probabilistic reconstruction of those chromosomal ancestry blocks within single individuals. From the block reconstructions, estimates of ancestry at any location in the genome can be derived. The authors recognized that high-density SNP arrays could include nearby markers that are in linkage disequilibrium (LD) in the ancestral population and that such LD could contrib-

ute noise to the block reconstructions and subsequent locus-specific ancestry estimation. Therefore, they proposed a Markov-Hidden Markov Model (MHMM) that allowed for pairwise dependency between adjacent markers in the ancestral populations in the estimation process and developed a computer program (SABER) to perform these calculations. They showed, through extensive simulations with data derived from the HapMap project,<sup>2</sup> that the method was robust in reconstructing ancestry blocks, even for very dense sets of markers and for an individual with three ancestral components, and when some of the model parameters were misspecified.<sup>1</sup> Subsequently, Tang et al.<sup>3</sup> used the MHMM to reconstruct ancestry blocks from Affymetrix 100K data in a sample of 192 Puerto Ricans from the Genetics of Asthma in Latino Americans (GALA) study<sup>4</sup> and examined the genome-wide distribution of African, European, and Native American ancestry in this sample.



**Figure 1. Comparison of Estimated and True Excess African Ancestry on Chromosome 6p**

markers in perfect LD, SABER would effectively treat this collection as a single marker only, because no additional information is provided after accounting for the background LD, and produce an accurate ancestry estimate. The MHMM method in SABER also uses a great deal more information than just single SNP genotypes because it uses the empirical distribution of ancestry-block sizes in determining for a given individual the ancestry state at a specific location. In fact, in their extensive simu-

The authors found strong evidence for statistical deviation in ancestry at three chromosomal locations (chromosome 6p, 8q, and 11q), allowing for both statistical variation due to sample size and for ancestry genetic drift, which creates random ancestry variation around the genome.<sup>5</sup> In particular, the location on chromosome 6p overlaps with the HLA cluster of loci, and the authors replicated an observed excess of African ancestry and deficit of European ancestry in an independent sample of Puerto Ricans from the literature, by using also published HLA allele frequencies.

Price et al. now raise a number of concerns regarding both the accuracy and unbiased nature of our ancestry estimation with the MHMM method,<sup>1</sup> as well as our conclusion regarding historic selection as the cause for the significant local ancestry deviations we observed.<sup>3</sup> Their primary concern regarding ancestry estimation is that inclusion of markers that are not in linkage equilibrium (LE) in the ancestral populations can lead to both increased noise and bias. They provide an example of  $n$  consecutive SNPs that are in perfect LD (with identical allele frequencies) and show that one can get distorted ancestry estimates if the loci are assumed to be independent. They also provide an example from simulated data in which the inclusion of markers in LD in a set of 1852 markers leads to excess noise and bias in the ancestry estimates. However, all these analyses were performed with the program ANCESTRYMAP and the theory described therein.<sup>6</sup> As the authors have stated, ANCESTRYMAP requires the use of statistically independent markers (i.e., no LD) and furthermore only allows for two ancestral populations.<sup>6</sup> We agree that these requirements may create problems for high-density array data, or more generally for data with markers that are in LD, or for populations with three ancestral components. However, the examples they present are unrealistic and have little relevance for analyses with SABER.<sup>1</sup> SABER allows for LD between adjoining markers in ancestral populations. Therefore, for the example of  $n$  consecutive

simulations with SABER, Tang et al.<sup>1</sup> clearly showed using real data (from HapMap) that the Markovian assumption of pairwise dependency in the ancestral populations was critical to obtain accurate ancestry-block reconstruction and locus-specific estimates. These simulations were performed with markers with an average spacing of 30 kb, 6 kb, and 3 kb (corresponding to the density of a 100K, 500K, and 1000K chip, respectively). Although the ancestry estimation became somewhat noisier with a higher density of SNPs when markers were assumed to be independent, the authors clearly showed robust reconstruction, even at the highest SNP density, when the Markov assumption of pairwise dependency was used via SABER.<sup>1</sup> For Price et al. to imply that their examples have relevance for SABER is incorrect and misleading.

Because Price et al. were particularly concerned about our results on chromosome 6p because of putative long-range linkage disequilibrium in this region in Europeans, we specifically re-examined the results on chromosome 6 from Simulation 2 in Tang et al.<sup>3</sup> According to those authors, “our simulated data incorporates a realistic level of high-order dependency among linked markers, and we have the opportunity to examine whether the MHMM is adequate.”<sup>3</sup> Thus, LD between nearby but nonconsecutive SNPs in the real data in this region is featured in the simulated data as well. Figure 1 compares the estimated excess of African ancestry (that is, the estimated locus-specific African ancestry subtracting out the genome-wide average African ancestry) with the true excess African ancestry along chromosome 6. The red line provides the estimated values, and the gray line provides the true values. Overall, the estimated excess of African ancestry is within 2% of the true values, and in fact there is no evidence of any systematic bias near the MHC region located between 26.0 and 34.0 Mb. These results provide additional reassurance from real data that the methods employed in SABER provide unbiased results in the presence of possible background

LD in the ancestral populations on chromosome 6p. Furthermore, in our original paper,<sup>3</sup> we studied the even and odd subsets of markers and found comparable deviations in both subsets in all three regions reported. Thus, the concern raised by Price et al. of systematic distortion in our local ancestry estimates appears to be unwarranted.

Price et al. also find fault in our analysis of HLA data in Puerto Ricans,<sup>1</sup> specifically regarding the appropriateness of populations we used to represent Native Americans in that analysis (Pima and Mayan). Although it is true that there is some genetic variation among Native American groups, and the Taino Indians were the Puerto Rican ancestors, the methods of ancestry estimation that we used and that were based on maximum likelihood (FRAPPE)<sup>7</sup> including the admixed subjects in the estimation of ancestral allele frequencies on the basis of the admixed subjects and not just the ancestral-population surrogates. We have shown previously that by allowing for re-estimation, we can accurately recapture the correct ancestral allele frequencies even when the surrogate-population allele frequencies are somewhat different.<sup>7</sup> Furthermore, a far more serious concern of bias in this type of analysis would arise from assuming that the Native American ancestry component in Puerto Ricans is 0, as Price et al. have done.

The admixture analysis of Price et al. of an independent sample of Puerto Rican Crohn's disease patients and controls, for which they claim no replication of our observed excess African and decreased European ancestry on chromosome 6p, also deserves comment. As they've stated, they reduced an initial marker set of 2459 SNPs to 1438 to eliminate markers that had allele frequency differences between Europeans and Native Americans as well as to "disallow LD between markers in the ancestral populations." This is because the ANCESTRYMAP program is not robust to background LD and also does not allow for more than two ancestral populations. Although this marker density corresponds to approximately one marker for every 2 Mb, because chromosome 6p putatively has an extended region of LD from 25.5 to 33.5 Mb,<sup>8</sup> we assume they allowed very few markers in this region, perhaps only one (rs451774 at 28.6 Mb). If so, the claim that "61% of maximum information about African vs. non-African ancestry at the chromosome 6p region" was obtained is difficult to imagine, especially because the allele frequency difference between Africans and Europeans for that marker is only approximately 0.40. Furthermore, the lack of allowance for Native American ancestry in their analysis makes their results difficult to interpret. They also show that their method is highly conservative because a simulated ancestry excess of .14 was reduced by more than a factor of two upon estimation. Despite the low power of their analysis, they still observed a modest increase in African ancestry at chromosome 6p and might have observed a greater increase with greater marker density and information.

Of course, we agree that all initial genetic observations, be they disease associations or arguments for ancestral selection, require independent replication. We therefore

also conducted an independent replication study, this time with AIMs rather than high-density chip data. We examined a new sample of 383 Puerto Rican subjects from the GALA study,<sup>4</sup> approximately double in size of our original sample. We typed 104 AIMs from around the genome and obtained a genome-wide estimate of African, European, and Native American ancestry for each individual using FRAPPE.<sup>7</sup> For comparison, we estimated ancestry on chromosome 6p using five ancestry informative markers: rs393228, rs7773913, rs853693, rs6456883, and rs847851. These markers span from 25.07 to 35.01 Mb on chromosome 6p. The estimated average African ancestry outside of chromosome 6p was 25.5%; by contrast, the estimated African ancestry based on the five markers on chromosome 6p was 40.0%, an excess of 14.5%, comparable to the difference we observed in our original study.<sup>3</sup> To assess statistical significance of this difference, we estimated the African ancestry at chromosome 6p for each individual. To do this, we first performed a single-marker analysis, in which we computed the posterior probability that an allele is derived from an African ancestor, given the ancestral allele frequencies and the individual's genome-wide ancestry:

$$\hat{z} = P(\text{African} | (t_{\text{afr}}, t_{\text{eur}}, t_{\text{amr}}), (p_{\text{afr}}, p_{\text{eur}}, p_{\text{amr}})) \\ = \frac{t_{\text{afr}} p_{\text{afr}}}{t_{\text{afr}} p_{\text{afr}} + t_{\text{eur}} p_{\text{eur}} + t_{\text{amr}} p_{\text{amr}}},$$

where  $(p_{\text{afr}}, p_{\text{eur}}, p_{\text{amr}})$  are the allele frequencies in the three ancestral populations, respectively, and  $(t_{\text{afr}}, t_{\text{eur}}, t_{\text{amr}})$  denote the genome-wide ancestry proportions for the individual. We then computed the location-specific ancestry of an individual by averaging over the five SNPs. This analysis is quite conservative because the ancestry estimate at chromosome 6p is shrunken significantly back toward the individual's genome-wide estimate by the Bayesian calculation. Thus, in this case we observed an average of 30.1% African ancestry at 6p, still greater than the 25.5% genome-wide estimate. We then calculated, for each individual, the difference between the estimated African ancestry at chromosome 6p (as derived above) and the genome-wide African ancestry. The mean of this difference was .0525, with a standard error of .0066. A t test to determine whether the mean is significantly different from 0 yielded a t value of 8.4,  $p < 10^{-15}$ . Thus, the conclusion of excess African ancestry on 6p compared with the rest of the genome in this sample is unequivocal and confirms our original observation.

The specter of bias in our analysis was probably raised by Price et al. due to the fact that the three locations we identified as sites of ancestral selection mapped into three regions with long-range LD, as they have described in Table 1 of their letter. Ironically, long-range LD has been cited as evidence for historical selection, not by us but by others, including some of the authors of the current letter.<sup>8</sup> In fact, long-range LD was used as an argument for historical

selection on the lactase persistence SNP on chromosome 2q.<sup>8</sup> Interestingly, this region of chromosome 2q (134.5 to 138.0 Mb) is also on the list of extended LD in Table 1 of Price et al. It is puzzling that on the one hand long-range LD has been used as evidence for selection in one analysis<sup>8</sup> and on the other as evidence for bias and against selection in the current letter.

Furthermore, our initial distribution of genome-wide excess African ancestry was quite symmetric and fit a simulated null distribution quite well, with the exception of a very small number of outlier loci (Figure 2 in Tang et al.<sup>3</sup>). These outlier loci were on chromosome 6p. If polymorphic inversions in the European population and associated regions of extended LD were an important source of bias in our analyses, as suggested by Price et al., we would have expected to see more outlier points in this distribution, specifically at locations corresponding to the 24 regions identified in Table 1 of Price et al. Aside from the three regions already mentioned, and possibly another region at 8p, none of the remaining 20 regions showed any deviation of ancestry from background levels.<sup>3</sup>

Extended regions of LD in the human genome have been previously described. Huttley et al.<sup>9</sup> studied 5048 autosomal microsatellite markers in Europeans and identified ten regions with putative evidence of long-range LD. Price et al. have now extended these findings by examining 550K SNP markers. Although the two studies identified some overlapping regions (chromosomes 2p, 6p, and 7p), many other regions are distinct. Whereas numerous authors, including Huttley et al.<sup>9</sup> and Bersaglieri et al.,<sup>8</sup> have suggested these regions represent targets of historical selection, Price et al. now propose that regions of long-range LD they identified are due to polymorphic inversions but have provided no evidence to support this contention. We believe other evidence argues against this conclusion. Jorgenson et al.<sup>10</sup> compared genetic maps across four major racial-ethnic groups in a very large sample of sibships. They noted that polymorphic inversions impact genetic map distances, and when the frequencies of these inversions differ across groups, map distances between markers in and around the inversion will consequently also differ significantly between groups. They identified two regions, one on chromosome 8p and another on 12q, that displayed ethnic-specific map differences. The region on chromosome 8p coincided with a previously described polymorphic inversion.<sup>11</sup> However, they found no other genomic region (aside from 12q) with significant ethnic-specific map differences (including on chromosomes 6p, 8q, and 11q); in particular, none of the regions in Table 1 of Price et al. aside from chromosome 8p showed evidence of map differences. Thus, the suggestion that the regions of long-range LD identified by Price et al. (aside from 8p) are due to polymorphic inversions appears highly speculative, at best.

We do agree that the fact that our three regions on chromosomes 6p, 8q, and 11q coincided with three regions of extended LD in Table 1 of Price et al. is unlikely to be due to chance. However, it seems inconsistent to argue that long-

range LD provides evidence of historical selection in one population but when similar evidence is found in a population derived from it, selection is deemed unlikely and artifact is invoked. For example, Price et al. and others have argued that the HLA region on chromosome 6p is particularly interesting because of its broad impact on disease. Again, it seems contradictory to argue that the HLA region on chromosome 6p has been a target of selection in Europeans and other populations but could not have been in Puerto Ricans, leading to a differential ancestry distribution. The region on chromosome 8p harboring a polymorphic inversion, which showed a suggestive but not significant ancestry deviation in our analysis, harbors an olfactory gene cluster and has shown phenotypic effects in other studies.<sup>11</sup> Furthermore, two of the other regions we identified (6p and 8q) also harbor olfactory gene clusters,<sup>3</sup> an observation that seems unlikely to be due purely to chance.

Price et al. argue that because they observed no evidence of ancestry distortions in African Americans, there must not be any in Puerto Ricans either. We do not see the relevance of this observation because these are populations with distinct and nonoverlapping social, demographic, and genetic histories.

In summary, we have shown that the MHMM approach, as implemented in the program SABER, is robust to putative regions of extended LD in real data. This method should be particularly useful for investigators studying admixed populations with high-density chips. Furthermore, we have shown a convincing replication of our prior results of excess African ancestry on chromosome 6p in Puerto Ricans.

Hua Tang,<sup>1</sup> Shweta Choudhry,<sup>2</sup> Rui Mei,<sup>5</sup>  
Martin Morgan,<sup>6</sup> William Rodriguez-Cintron,<sup>7,8</sup>  
Esteban Gonzalez Burchard,<sup>2,4</sup> and Neil J. Risch<sup>3,4,9,\*</sup>

<sup>1</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305, USA; <sup>2</sup>Departments of Biopharmaceutical Sciences and Medicine, <sup>3</sup>Department of Epidemiology and Biostatistics, <sup>4</sup>Institute for Human Genetics, University of California San Francisco, San Francisco, CA 94143, USA; <sup>5</sup>Affymetrix, Santa Clara, CA 95051, USA; <sup>6</sup>Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA; <sup>7</sup>Veterans Caribbean Health Care System, San Juan, PR 00921, USA; <sup>8</sup>University of Puerto Rico School of Medicine San Juan, PR 00936-5067, USA; <sup>9</sup>Division of Research, Kaiser Permanente, Oakland, CA 94611, USA

\*Correspondence: [rischn@humgen.ucsf.edu](mailto:rischn@humgen.ucsf.edu)

## References

1. Tang, H., Coram, M., Wang, P., Zhu, X., and Risch, N. (2006). Reconstructing genetic ancestry blocks in admixed individuals. *Am. J. Hum. Genet.* 79, 1–12.
2. International HapMap Consortium (2005). A haplotype map of the human genome. *Nature* 437, 1299–1320.
3. Tang, H., Choudhry, S., Mei, R., Morgan, M., Rodriguez-Cintron, W., Burchard, E.G., and Risch, N.J. (2007). Recent genetic

selection in the ancestral admixture of Puerto Ricans. *Am. J. Hum. Genet.* 81, 626–633.

4. Burchard, E.G., Avila, P.C., Nazario, S., Casal, J., Torres, A., Rodriguez-Santana, J.R., Syliva, J.S., Fagan, J.K., Salas, J., Lilly, C.M., et al. (2004). Lower bronchodilator responsiveness in Puerto Rican than in Mexican asthmatic subjects. *Am. J. Respir. Crit. Care Med.* 169, 386–392.

5. Long, J. (1991). The genetic structure of admixed populations. *Genetics* 127, 417–428.

6. Patterson, N., Hattangadi, N., Lane, B., Lohmueller, K., Hafler, D., Oksenberg, J., Hauser, S., Smith, M., O'Brien, S., Altschuler, D., et al. (2004). Methods for high-density admixture mapping of disease genes. *Am. J. Hum. Genet.* 74, 979–1000.

7. Tang, H., Peng, J., Wang, P., and Risch, N. (2005). Estimation of individual admixture: Analytical and study design considerations. *Genet. Epidemiol.* 28, 289–301.

8. Bersaglieri, T., Sabeti, P.C., Patterson, N., Vanderploeg, T., Schaffner, S.F., Drake, J.A., Rhodes, M., Reich, D.E., and Hirsch-

horn, J.N. (2004). Genetic signatures of strong recent positive selection at the lactase gene. *Am. J. Hum. Genet.* 74, 1111–1120.

9. Huttley, G.A., Smith, M.W., Carrington, M., and O'Brien, S.J. (1999). A scan for linkage disequilibrium across the human genome. *Genetics* 152, 1711–1722.

10. Jorgenson, E., Tang, H., Gadde, M., Province, M., Leppert, M., Kardia, S., Schork, N., Cooper, R., Rao, D.C., Boerwinkle, E., and Risch, N. (2005). Ethnicity and human genetic linkage maps. *Am. J. Hum. Genet.* 76, 276–290.

11. Giglio, S., Broman, K.W., Matsumoto, N., Calvari, V., Gimelli, G., Neumann, T., Ohashi, H., Voullaire, L., Larizza, D., Giorda, R., et al. (2001). Olfactory receptor-gene clusters, genomic-inversion polymorphisms, and common chromosome rearrangements. *Am. J. Hum. Genet.* 68, 874–883.

DOI 10.1016/j.ajhg.2008.06.009. ©2008 by The American Society of Human Genetics. All rights reserved.

## East Asian and Melanesian Ancestry in Polynesians

*To the Editor:* Kayser et al.<sup>1</sup> estimated the ancestry of Polynesians by using 377 autosomal microsatellite loci and concluded that 0.79 of the ancestry was from East Asians (95% CI, 0.76–0.84) and 0.21 from Melanesians. In contrast, maternally inherited mtDNA ancestry was previously estimated to be 0.94 East Asian and 0.06 Melanesian and paternally inherited Y chromosome ancestry was estimated to be 0.28 East Asian and 0.66 Melanesian.<sup>2</sup> One might guess that the East Asian autosomal ancestry would be approximately the arithmetic average of the mtDNA and Y ancestry, 0.61, but the estimated autosomal ancestry of 0.79 is substantially higher. To account for this difference and the different estimates in ancestry from different chromosomes, strong sex differences in gene flow, occurring in a particular chronological order, are necessary. Here I present a simple two-phase scenario to explain the different observed ancestries for autosomal, mtDNA, and Y markers and then discuss how this scenario could be modified and still result in the observed patterns.

First, assume that a population of East Asian ancestry, which eventually became the Polynesians, settled in Melanesia and that subsequently there was male gene flow from Melanesians into this population. This pattern is consistent with both matrilocality and matrilinearity in this population.<sup>1</sup> The effect of this male gene flow at a rate of  $m_m$  per generation over  $t$  generations on Y ancestry can be given<sup>3</sup> as

$$q_t = (1 - m_m)^t q_0 + [1 - (1 - m_m)^t] q_{Mel}$$

where  $q_0$  and  $q_t$  are the initial and  $t$  generation East Asian ancestry in the population and  $q_{Mel}$  is the East Asian ances-

try in the Melanesian migrants. Assuming that  $q_{Mel} = 0$ ,  $q_0 = 1$ , and  $q_t = 0.28$ , then

$$0.28 = (1 - m_m)^t \text{ or } m_m = 1 - e^{\ln(0.280)/t}.$$

For example, if  $t = 50$ , then  $m_m = 0.0251$ .

For autosomal loci in this population, the East Asian ancestry is

$$q_t = \left[1 - \frac{1}{2}(m_f + m_m)\right]^t q_0 + \left\{1 - \left[1 - \frac{1}{2}(m_f + m_m)\right]^t\right\} q_{Mel}$$

where  $m_f$  is the per-generation rate of female gene flow. Again, assume that  $q_{Mel} = 0$ ,  $q_0 = 1$ ,  $m_f = 0$ , and with the estimated  $m_m$  of 0.0251 used,  $q_t = (0.987)^t$ . For example, if  $t = 50$ , then  $q_t = 0.532$ .

Second, assume that subsequently there was female gene flow from the East Asians into this population for  $x$  generations so that the autosomal East Asian ancestry can be expressed as

$$q_{t+x} = \left[1 - \frac{1}{2}(m_f + m_m)\right]^x q_t + \left\{1 - \left[1 - \frac{1}{2}(m_f + m_m)\right]^x\right\} q_{EA}$$

where  $q_{EA}$  is the East Asian ancestry in the East Asian female migrants. Assuming that  $q_{EA} = 1$ ,  $q_t = 0.532$ ,  $q_{t+x} = 0.79$ , and  $m_m = 0$ ,

$$0.79 = 1 - 0.468 \left(1 - \frac{1}{2}m_f\right)^x \text{ or } m_f = 2(1 - e^{\ln(0.449)/x}).$$

For example, if  $x = 50$ , then  $m_f = 0.0318$ .

This two-phase scenario is presented in **Figure 1**, which shows a decline in Y and autosomal East Asian