

## File S1

## SUPPORTING METHODS AND RESULTS

**MCMC algorithm** We developed a Metropolis-Hastings MCMC algorithm (GAMERMAN and HEDIBERT 2006) to obtain samples from the joint posterior probability distribution for all model parameters. Haplotype frequencies were estimated using independence or random-walk chains. When independence chains were used, proposal values for haplotype frequencies (a vector  $p_{ij}$  containing values for each locus and population) were sampled from Dirichlet distributions that were independent of  $\mathbf{p}$  from the previous time-step and similar in form to the expected posterior distribution for these parameters. This proposal distribution is very efficient when dealing with few haplotypes and intermediate haplotype frequencies. Random-walk chains were used when these criteria were not met, which involved sampling haplotype frequencies from Dirichlet distributions that were proportional to the vector  $p_{ij}$  from the previous MCMC step. At least one of these two proposal algorithms generally worked well with each data set, however, more complicated, alternative proposal distributions might be considered when a very large number of haplotypes are analyzed. The  $\alpha$  and  $\beta$  parameters associated with the conditional prior on haplotype frequencies were estimated using random-walk chains. Specifically, new values for each  $\alpha$  and  $\beta$  pair were proposed from bivariate Gaussian distributions centered on the previous parameter values with user adjusted variance and covariance. Specification of a high covariance between proposal values of  $\alpha$  and  $\beta$  was imposed to increase chain mixing. The MCMC algorithm was written in C++ using the GNU Scientific Library (GALASSI *et al.* 2009) and is available from the authors at <http://www.uwyo.edu/buerkle/software/> as the stand-alone software *bamova*.

**Simulations: estimation of  $\phi$  statistics** We conducted a series of simulations to determine whether the proposed model provided reasonable estimates of genome-level  $\phi$ -statistics. For these simulations we were solely concerned with genetic differentiation among popula-

26 tions (rather than also considering differentiation among groups of populations). For each  
27 of our three likelihood models we simulated sequence data using an infinite sites coalescent  
28 model (using R. Hudson’s software *ms*; HUDSON 2002). One group of data consisted of  
29 sequences from 25 genetic regions, whereas the second group consisted of sequences from  
30 500 genetic regions. All simulations assumed five populations split from a common ancestor  
31  $\tau$  generations in the past, where  $\tau$  has units of  $4N_e$ . We varied  $\tau$  from 0 to 1 in steps of  
32 0.05 to produce 21 data sets each for 25 and 500 loci. The ancestral population and all five  
33 descendant populations were assigned population mutation rates  $\theta = 4N_e\mu$  of 0.5, where  $\mu$   
34 is the per locus mutation rate. We assumed no migration following population subdivision.  
35 Forty gene copies were sampled from each of the five populations. For the *known haplotype*  
36 *model* analyses we treated the simulated sequences directly as the sampled data. For *NGS-*  
37 *individual model* and *NGS-population model* analyses we re-sampled the simulated sequence  
38 data sets such that coverage for each sequence was Poisson distributed ( $\lambda = 2$ ). For the  
39 *NGS-individual model* analyses we retained information on which individual each sequence  
40 came from, whereas we only retained population identification for *NGS-population model*  
41 analyses. Each data set was analyzed using our *bamova* software, with MCMC details as  
42 described in the main document.

43 MCMC implementation of the proposed Bayesian models accurately quantified genetic  
44 structure among five simulated populations with sequence data from 25 or 500 genetic re-  
45 gions (Figure S1). In general, estimates of mean genome-level  $\phi_{ST}$  ( $\mu_{ST}$ ) increased with  
46 the time since divergence of the five populations ( $\tau$ ). Credible intervals for genome-level  
47 parameters were relatively narrow, particularly when estimates were based on 500 genetic  
48 regions (Figure S1, S2). Moreover, credible intervals, and thus the uncertainty in genome-  
49 level parameters, were similar for all three first-level likelihood models (*known haplotype*  
50 *model*, *NGS-individual model*, and *NGS-population model*). We detected considerable varia-  
51 tion in the extent of population structure among genetic regions (and hence non-zero  $\sigma_{ST}$  for  
52 genome-level  $\phi_{ST}$ ), except when the population divergence time was very low (Fig. S2). Pos-

terior probability estimates for  $\mu_{ST}$  were similar to the empirical mean of the locus-specific  $\phi$  statistics calculated directly from the raw data; however, the estimates of  $\sigma_{ST}$  were generally lower than the empirical standard deviation of  $\phi_{ST}$  from the raw data.

In the analyses of simulated data sets,  $\phi_{ST}$  increased reliably and as expected with time since population divergence. Moreover, estimates of genome-level  $\phi_{ST}$  using the *known haplotypes model* were very similar to non-Bayesian point estimates of mean  $\phi_{ST}$  (Figure S1). Additionally the estimates of genome-level  $\phi_{ST}$  for the *known haplotypes model*, the *NGS-individual model*, and the *NGS-population model* were similar. This similarity in results among models suggest that high-coverage NGS data can provide parameter estimates with precision and accuracy equivalent to Sanger sequencing. Furthermore, the estimates of genome-level  $\phi_{ST}$  for the SeattleSNPs human sequence data and chromosome-level  $\phi_{ST}$  for the worldwide human SNP data (0.080–0.139) were similar to mean levels of genetic differentiation among human populations based on  $F_{ST}$  (e.g.,  $F_{ST} = 0.09$ – $0.14$  for Yoruba, European, Han Chinese and Japanese populations; WEIR *et al.* 2005; BARREIRO *et al.* 2008). An important attribute of the model is that it also provides an accurate estimate of the uncertainty in the parameter estimates. This is an attribute not necessarily shared by non-Bayesian methods of parameter estimation, particularly when hierarchical or derived parameters are involved (LINK and BAKER 2009).

**Human SeattleSNP data: alternative data subsets** In addition to analysing the SeattleSNPs data set based on the first five SNPs in each gene we analysed four additional subsets of these data: 1) sequences based on the middle five SNPs in each gene, 2) sequences based on the last five SNPs in each gene, 3) sequences based on five SNPs spaced evenly across each gene, 4) and sequences based on every 12<sup>th</sup> SNP in each gene (mean number of SNPs = 5.24, sd = 0.423). Analyses of these data sets were as described in the main text for the first five SNPs data set.

We classified four genes as high  $\phi_{ST}$  outliers (using  $a = 0.5$ ) in two or more of the data

79 subsets (Figs. 1, S3). Three of these genes, HSD11B2, FOXA2, and POLG2 were classified  
80 as  $\phi_{ST}$  outliers based on the 'first five SNPs' data subset, and are described in the main  
81 document. Other outlier gene identified in more than one data subset was CPSF4, which  
82 encodes the cleavage polyadenylation specificity factor subunit 4 protein and is an essential  
83 component of pre-mRNA 3' processing in mammals (BARABINO *et al.* 1997). Estimates  
84 of  $\phi_{ST}$  for CPSF4 were as high as 0.382 (95% ETPI 0.262–0.496; 'last five SNPs' data  
85 subset, Fig. S3). Four additional genes were identified as high  $\phi_{ST}$  outliers in single subsets  
86 of the data: FUT2, IL1F6, EPPB9, and IKBKB. When classified as outliers these genes  
87 had  $\phi_{ST}$  estimates similar to the genes detected as outliers more than once (Figs. 1, S3).  
88 Interestingly, FUT2 was classified as a candidate gene experiencing balancing selection in  
89 European Americans based on levels of polymorphism and intermediate-frequency alleles by  
90 Andres *et al.* (ANDRÉS *et al.* 2009) and is generally regarded as a well-established target  
91 of balancing selection (contrary to our findings). Variation among data subsets in whether  
92 genes were detected as outliers depended both on the distribution of divergent nucleotides  
93 along each gene and the extent of divergence at each of these nucleotides (Fig. 2). No genes  
94 were identified as low  $\phi_{ST}$  outliers, nor were any genes identified as high  $\phi_{ST}$  outliers using  
95  $a = 0.95$ .

## LITERATURE CITED

- 96  
97 ANDRÉS, A. M., M. J. HUBISZ, A. INDAP, D. G. TORGERSON, J. D. DEGENHARDT,  
98 A. R. BOYKO, R. N. GUTENKUNST, T. J. WHITE, E. D. GREEN, C. D. BUSTAMANTE,  
99 A. G. CLARK, and R. NIELSEN, 2009 Targets of balancing selection in the human  
100 genome. *Molecular Biology and Evolution* **26**: 2755–2764.
- 101 BARABINO, S. M. L., W. HUBNER, A. JENNY, L. MINVIELLESEBASTIA, and  
102 W. KELLER, 1997 The 30-kD subunit of mammalian cleavage and polyadenylation speci-  
103 ficity factor and its yeast homolog are RNA-binding zinc finger proteins. *Genes and De-*  
104 *velopment* **11**: 1703–1716.

- 105 BARREIRO, L. B., G. LAVAL, H. QUACH, E. PATIN, and L. QUINTANA-MURCI, 2008  
106 Natural selection has driven population differentiation in modern humans. *Nature Genetics*  
107 **40**: 340–345.
- 108 GALASSI, M., J. DAVIES, J. THEILER, B. GOUGH, G. JUNGMAN, M. BOOTH, and  
109 F. ROSSI, 2009 *GNU Scientific Library: Reference Manual*. Network Theory Ltd.
- 110 GAMERMAN, D. and F. L. HEDIBERT, 2006 *Markov Chain Monte Carlo: Stochastic*  
111 *Simulation for Bayesian Inference*. Chapman and Hall, New York.
- 112 HUDSON, R. R., 2002 Generating samples under a Wright-Fisher neutral model of genetic  
113 variation. *Bioinformatics* **18**: 337–338.
- 114 JAKOBSSON, M., S. W. SCHOLZ, P. SCHEET, J. R. GIBBS, J. M. VANLIERE, H.-C.  
115 FUNG, Z. A. SZPIECH, J. H. DEGNAN, K. WANG, R. GUERREIRO, J. M. BRAS, J. C.  
116 SCHYMICK, D. G. HERNANDEZ, B. J. TRAYNOR, J. SIMON-SANCHEZ, M. MATARIN,  
117 A. BRITTON, J. VAN DE LEEMPUT, I. RAFFERTY, M. BUCAN, H. M. CANN, J. A.  
118 HARDY, N. A. ROSENBERG, and A. B. SINGLETON, 2008 Genotype, haplotype and  
119 copy-number variation in worldwide human populations. *Nature* **451**: 998–1003.
- 120 LINK, W. A. and R. J. BAKER, 2009 *Bayesian Inference: with Ecological Applications*.  
121 Academic Press, Maryland Heights, MO, USA.
- 122 WEIR, B. S., L. R. CARDON, A. D. ANDERSON, D. NIELSEN, and W. HILL, 2005 Mea-  
123 sures of human population structure show heterogeneity among genomic regions. *Genome*  
124 *Research* **15**: 1468–1476.