

Interrogating multiple aspects of variation in a full resequencing data set to infer human population size changes

Benjamin F. Voight^{†,‡}, Alison M. Adams^{†,§}, Linda A. Frisse^{†,¶}, Yudong Qian[†], Richard R. Hudson^{||}, and Anna Di Rienzo^{†,††}

[§]Committee on Genetics and Departments of [†]Human Genetics and ^{||}Ecology and Evolution, University of Chicago, Chicago, IL 60637

Edited by Daniel L. Hartl, Harvard University, Cambridge, MA, and approved October 28, 2005 (received for review August 22, 2005)

We present an expanded data set of 50 unlinked autosomal noncoding regions, resequenced in samples of Hausa from Cameroon, Italians, and Chinese. We use these data to make inferences about human demographic history by using a technique that combines multiple aspects of genetic data, including levels of polymorphism, the allele frequency spectrum, and linkage disequilibrium. We explore an extensive range of demographic parameters and demonstrate that our method of combining multiple aspects of the data results in a significant reduction of the compatible parameter space. In agreement with previous reports, we find that the Hausa data are compatible with demographic equilibrium as well as a set of recent population expansion models. In contrast to the Hausa, when multiple aspects of the data are considered jointly, the non-Africans depart from an equilibrium model of constant population size and are compatible with a range of simple bottleneck models, including a 50–90% reduction in effective population size occurring some time after the appearance of modern humans in Africa 160,000–120,000 years ago.

bottlenecks | combining *P* values | human demographic inference | population growth

Elucidating how and when populations change in size is an important element in reconstructing evolutionary history because these changes often reflect crucial events in the history of a species, such as range expansions, environmental changes, and mixture between groups (1). In addition, making inferences based on population variation data typically requires the specification of a demographic model. Such applications include detecting the signature of natural selection or estimating recombination rates from patterns of linkage disequilibrium (LD) (2–5). Finally, better knowledge of demographic histories in human populations is particularly important for whole-genome, LD-based association studies (6, 7).

Motivated by the excess of rare variants observed in mitochondrial DNA data, attention initially focused on models of ancient population growth and on the idea that population expansions may have accompanied the dispersal out of Africa or the emergence of new tool technology in the Upper Paleolithic (8–13). However, the accumulation of nuclear sequence variation surveys showed that this simple growth model was consistent with the observed frequency spectrum only for a subset of the loci (14–16). Likewise, LD surveys revealed marked differences in the rate of LD decay in African populations compared with that in non-African populations (17–19). These results together with the higher levels of sequence variation in African populations compared with that in non-African populations led to the proposal that population size reduction, such as bottlenecks, account for patterns of variation and LD in non-African populations (15, 18, 19). This bottleneck was hypothesized to correspond with the dispersal of modern humans out of Africa (18).

However, the investigation of formal bottleneck models has typically used a single aspect of genetic variation data, either the allele frequency spectrum (15, 20) or patterns of LD (18, 21), raising the question of whether such models were indeed consistent with

the data when multiple aspects of genetic variation were considered simultaneously (22–24). Specifically, it is not known whether simple bottleneck models can generate the marked differences in LD levels seen between Africans and non-Africans with only a limited reduction in polymorphism levels outside Africa. Although previous work suggested that variation in recombination rate may explain the decay in LD observed in a multiethnic sample (7), it is not obvious that it could also explain the differences between Africans and non-Africans.

Ideally, making inferences about population history should be based on data from a large number of unlinked and neutrally evolving loci and on statistical methodology that makes efficient use of all or most of the information in the data. Full resequencing studies, in which the sequence of the surveyed segments is determined for every individual in the sample, represent one scheme for generating data sets in which multiple aspects of sequence variation are characterized. With regard to data analysis, full likelihood methods have been successfully applied to nonrecombining data (Y chromosome or mitochondrial DNA) to reconstruct population histories (25–27). However, for regions with recombination, the currently available methods are computationally infeasible. As a result, a variety of statistics, each summarizing different aspects of genetic variation data, may be used (13, 15, 28), with the subsequent reduction in information content traded for computational tractability. It is still desirable to combine the results of tests based on individual statistics because the joint distributions of multiple summaries of the data should contain more information than the marginal distributions of multiple single summaries considered separately.

We previously developed a full resequencing scheme in which pairs of tightly but not completely linked segments, referred to as “locus pairs,” are surveyed (19). This study design aims to maximize the information content for a given amount of sequencing effort because, by skipping the intervening segment, many more independent loci can be surveyed. Using this scheme, we previously surveyed 10 noncoding regions in three human population samples: Hausa of Cameroon, Italians, and Chinese. Here, we survey an additional 40 locus pairs in the same samples. This data set allows the simultaneous characterization of polymorphism levels, allele frequency spectrum, and LD in each sample; in addition, it obviates the need to correct for ascertainment bias with its associated uncertainties and possible loss of information (29–31). In choosing only noncoding regions distant from genes, we limit the possibility that our analysis of demographic history will be confounded by the

Conflict of interest statement: No conflicts declared.

This paper was submitted directly (Track II) to the PNAS office.

Freely available online through the PNAS open access option.

Abbreviations: LD, linkage disequilibrium; ML, maximum likelihood.

[‡]B.F.V., A.M.A., and L.A.F. contributed equally to this work.

[¶]Present address: Computercraft Corporation, McLean, VA 22102.

^{††}To whom correspondence should be addressed at: 920 East 58th Street, Chicago, IL 60637.

E-mail: dirienzo@genetics.uchicago.edu.

© 2005 by The National Academy of Sciences of the USA

effects of natural selection. To analyze these data, we implement an approach to determine P values associated with several observed summaries of genetic data considered jointly over a grid of demographic parameter values. These summaries include the average Tajima's D (\bar{D}) and the variance of Tajima's D across loci $\text{Var}[D]$ (32), the average number of segregating sites across loci (\bar{S}), and an overall composite likelihood estimator of the population cross-over rate parameter ($\hat{\rho}$) (33). By combining P values obtained from these individual statistics into a single statistical test, we greatly improve the power to reject demographic scenarios incompatible with the data. Although it is well established that other demographic features apply to these populations (e.g., population subdivision and gene flow) (34, 35), we chose to focus solely on population size changes to reduce modeling complexity. We explore an extensive grid of the demographic parameter space that revealed a confidence set of relatively simple bottleneck models that explain the patterns of variation in the non-African samples. Our results combine aspects of genetic variation from allele frequency spectrum, LD, and polymorphism levels within noncoding autosomal regions to infer the history of human populations. Because our data set was collected without ascertainment, it may be useful for validating the results of SNP genotyping surveys.

Materials and Methods

DNA Samples. Sequence variation was surveyed in DNA samples from the same three human populations investigated in Frisse *et al.* (2001): 15 Hausa samples from Yaounde, Cameroon; 15 individuals from central Italy; and 15 Han Chinese from Taiwan. In addition, one common chimpanzee DNA sample was also sequenced at each region. This study was approved by the Institutional Review Board of the University of Chicago.

Resequencing Data Collection. We selected 40 unlinked genomic regions for resequencing using the locus pair approach (19): For each unlinked region, we sequenced two segments of ≈ 1 kb separated by ≈ 8 kb. The selection of genomic targets was aimed at regions that did not contain nor were tightly linked to known or strongly predicted coding regions. Most surveyed segments also did not contain and were not tightly linked to noncoding regions strongly conserved between human and mouse (as determined by inspection of the VISTA genome browser). These regions were selected as described in ref. 19 except that here we deliberately included regions with a broader range of cross-over rates and %G+C content. The local cM:Mb (Mb-megabase) ratio was obtained based on the interval defined by the two closest flanking markers on the DeCode Genetics (Reykjavik, Iceland) genetic map (36). The average and variance of the cM:Mb ratio across the 50 segments (i.e., 40 locus pairs from this study and the 10 given in ref. 19) are 1.31 and 0.83, respectively. The average and variance of %G+C across the 50 locus pairs are 38.3 and 46.6, respectively. Detailed information on each surveyed segment is provided in Table 2, which is published as supporting information on the PNAS web site. PCR and sequencing was performed as described in refs. 19 and 37. All sequencing reactions were run on automated capillary sequencers (ABI3100 and ABI3700). Sequence reads were scored by using POLYPHRED (38); all putative polymorphisms and software-derived genotype calls were visually inspected and individually confirmed.

Testing Demographic Models. For each demographic model of interest, we performed a separate test for each summary statistic of genetic variation. In addition, for some of the models (equilibrium and bottleneck), we also calculated a test statistic, C , which combines the P values of multiple summary statistics as follows:

$$C = -2 \sum_{i=1}^k \ln(p_i), \quad [1]$$

where P_i is the estimated P value of the i th summary statistic of k summary statistics.

For models defined by more than one demographic parameter (i.e., simple growth and bottleneck models), these tests were performed over a grid of parameter values. The combinations of parameter values that are compatible with the observed values of the test statistic(s) constitute the accepted portion of the parameter space for each model. For simple growth models, the test was based on Fu and Li's D^* (39), whereas for bottleneck models, the test was based on combining P values from multiple summary statistics, as discussed below. The P values, P_i , for each individual summary statistic were estimated from Monte Carlo simulations using a modification of the program MS (40), as follows. We used coalescent simulations to generate 50,000 replicates, each consisting of 50 independent locus pairs, for each combination of parameter values; mutation and recombination rates were allowed to vary across locus pairs as described below. Samples of sequences 10 kb in length were generated in which the intervening 8 kb were ignored to mimic the locus pair data. The probability, P , of observing a value greater than that found in the data were estimated by simulations and converted to a two-tailed P value by applying the formula $1 - 2 \cdot |0.5 - P|$.

The P values for the combined test statistic C were estimated by using the empirical distribution of the statistic from simulations. For each combination of parameter values, we recorded the values of each summary statistic in each replicate and generate the distribution of these simulated values. For each replicate, we treated the value of each summary statistic as the "observed" value and determined its P value relative to the empirical distribution from the remaining 49,999 replicates. For each replicate, we combined these P values to calculate a value of C . By following this procedure for each of the 50,000 replicates (for a single demographic scenario of interest), we obtained a distribution of the combined statistic. This distribution can be used to estimate a one-tailed P value for the observed value of C .

Mutation Rate Model. We assumed an infinite sites model, where we modeled the variation in mutation rate across locus pairs by using a gamma(12.46, 2.11×10^{-9}) distribution. The mean and variance for this distribution matched the observed mean and variance for the mutation rates estimated based on human-chimpanzee sequence divergence in our locus pair data (assuming 6 million years since divergence and a generation time of 25 years). The 90% central interval of this distribution is (1.54×10^{-8} , 3.96×10^{-8}) with $E\mu = 2.63 \times 10^{-8}$.

Recombination Rate Model. We modeled the variation in the crossing-over rate, c , across locus pairs using a lognormal $[-18.148, (0.5802)^2]$ distribution; cross-over rate was assumed to be homogeneous within each locus pair. The 90% central interval of this distribution is (0.51×10^{-8} , 3.41×10^{-8}). The median of this distribution matched the overall recombination rate for the Hausa data (1.31×10^{-8}) based on the composite likelihood estimator, $\hat{\rho}$, of Hudson (33). Because we cannot accurately estimate the variance in recombination rate across surveyed segments as short as 10 kb, we matched the variance of the lognormal distribution to the variance of cM:Mb values estimated from the Marshfield genetic map for the interval containing each locus pair (41). We acknowledge that this model may capture some but not all of the recombination rate variation estimated across the human genome (42).

Summary Statistics. We summarize the locus pair data in terms of the average Tajima's D (\bar{D}), the variance of Tajima's D ($\text{Var}[D]$), the average Fu and Li's D^* (\bar{D}^*), the average number of segregating sites (\bar{S}), and the average nucleotide diversity across the 50 locus pairs ($\bar{\pi}$), as well as $\hat{\rho}$, an overall estimate of the population crossing-over parameter ($4Nc$) as obtained by composite likelihood (33). Because there is not enough information in our data to accurately estimate $\hat{\rho}$ and the gene conversion parameters (43), we

Table 1. Observed summary statistics

Population	\bar{D}	$\widehat{\text{Var}}[D]$	\bar{D}^*	\bar{S}	$\bar{\pi}$, %	$\hat{\rho}$
Hausa	-0.20	0.55	-0.17	11.1	0.110	0.0006
Italian	0.28*	1.19**	0.18	7.1	0.085	0.0003
Chinese	0.18	1.08*	0.05	6.9	0.079	0.0002*

Observed summary statistics of polymorphism data for 50 locus pairs. *, $P < 0.05$; **, $P < 0.01$; under an equilibrium model.

assumed a model of gene conversion with rate (f) twice that of cross-over and tract lengths exponentially distributed with mean (L) 500 bp and estimate $\hat{\rho}$. Alternative models of gene conversion ($f = 10, L = 55$ bp) based on sperm-typing data (44) yielded qualitatively similar results (data not shown).

Results

Summary of Sequence Variation and Tests of the Equilibrium Model.

We resequenced 40 unlinked locus pairs in 15 individuals from each of three population samples: Hausa, Italians, and Chinese. The results of this survey are analyzed together with data for an additional 10 unlinked locus pairs previously resequenced in the same population samples (19), for a total of 50 unlinked locus pairs. The average surveyed length per locus pair was 2,365 bp (for a total of 118,259 bp surveyed in each individual), and the average unsurveyed intervening segment was 7,921 bp long.

The values of summary statistics used for demographic testing are shown in Table 1, with a synopsis of the summary statistics for the 40 new locus pairs presented in Table 3, which is published as supporting information on the PNAS web site. The allele frequency spectrum was summarized by the average and variance of Tajima's D and Fu and Li's D^* across loci, polymorphism levels are summarized by the average number of polymorphic sites (\bar{S}) across loci, and LD decay was summarized in terms of an overall composite likelihood estimator of the population cross-over rate parameter $\hat{\rho}$ (33). The results of this expanded data set are in qualitative agreement with those from our previous survey (19) and with other similar data sets (2, 5, 15, 16). With regard to the allele frequency spectrum, the Hausa show a skew toward rare variants and a low variance across loci, whereas both non-African samples have an excess of intermediate frequency variants and high variance across loci. In addition, polymorphism levels and LD decay are higher in the Hausa compared with both non-African samples, but this difference is greater for LD decay (1.9- to 3.2-fold) than polymorphism levels (1.6-fold).

To determine whether the levels of LD decay and the frequency spectrum were consistent with a model of constant population size, we conducted coalescent simulations under equilibrium to determine the P values of the observed summary statistics. We obtained the effective population size, denoted N_A , for each population by using an estimator of the population mutation rate parameter ($4N_A\mu$) based on the number of polymorphic sites and sample size (45), and an estimate of μ based on sequence divergence between human and chimpanzee for the 50 locus pairs. Each summary statistic for the Hausa data are consistent with the equilibrium model (Table 1). However, for the non-African populations, the skew toward intermediate frequency variants, and the elevated LD are incompatible with a simple equilibrium model; a combined statistic based on \bar{D} , $\widehat{\text{Var}}[D]$, and $\hat{\rho}$, obtained by using Eq. 1, is significant for the Italian ($P \leq 0.0148$) and Chinese data ($P \leq 0.0052$).

Estimating the Ancestral Population Size Under a Growth Model. Even though a model of constant population size could not be rejected for the Hausa, human populations certainly experienced rapid growth recently and, perhaps, in more ancient times. Thus, the negative but nonsignificant values of Tajima's D and Fu and Li's

D^* in the Hausa may simply reflect limited power and suggest that some expansion models are appropriate for this population. By following the approach in ref. 28, we considered a model in which an ancestral population at equilibrium size N_A grows exponentially beginning t_{onset} generations in the past at rate α , such that the present population size is $N_A e^{\alpha t_{\text{onset}}}$ (8). To test this model, we fixed the ancestral population size for each combination of demographic parameter values, such that the expected number of segregating sites matched the average number observed in the Hausa sample (28).

Unlike in ref. 28, we estimated the best-fit growth parameters for the Hausa data, α and t_{onset} , along with the associated point estimate of N_A , via approximate maximum likelihood (ML) based on the summary statistic, Fu and Li's D^* . We focused on the average D^* across the 50 locus pairs, denoted \bar{D}_{obs}^* , because it was previously shown to be the most informative for discriminating between equilibrium and growth models (28). For each demographic growth model, we obtained distributions of \bar{D}^* by simulation and estimated the probability that $|\bar{D}^* - \bar{D}_{\text{obs}}^*| < 0.001$ and then chose the model for which this probability was highest. This procedure returns the approximate ML estimate of the growth parameters, α and t_{onset} , compatible with the Hausa data based on \bar{D}_{obs}^* . Note that we refer to this as approximate ML on a summary statistic because we do not use the full data and because we approximate rather than obtain the probabilities exactly. We found that the model with the highest overall probability was at an α of 0.75×10^{-3} and t_{onset} of 1,000 generations, which corresponds to a model with ≈ 2 -fold growth starting 25,000 years ago, assuming a generation time of 25 years, from an ancestral population size of 10,659. We present confidence sets of α and t_{onset} for which \bar{D}_{obs}^* are consistent with the observed Hausa data in Fig. 3, which is published as supporting information on the PNAS web site. The span of acceptable models is consistent with previous reports (28), with a slight reduction in confidence set due to the inclusion of additional data.

To assess the uncertainty in N_A , we obtain a range of N_A consistent with the ML estimate of $\hat{\alpha} = 0.75 \times 10^{-3}$ and $\hat{t}_{\text{onset}} = 1,000$ as follows. We performed additional coalescent simulations as described earlier, where we used the ML parameters for the demographic history and gradually lowered or raised the value of N_A until \bar{S} was incompatible with the observed data at the 5% level. We found these high and low values of N_A to be 9,450 and 12,300, respectively. Later, we will use this information to assess the effect of our choice of N_A in testing bottleneck models.

Testing Bottleneck Models in the Non-African Data. The positive \bar{D} values and large $\widehat{\text{Var}}[D]$ along with the low polymorphism and high LD levels observed in the non-African populations (Table 1) suggest that models including a reduction in population size may be compatible with the data. We considered one family of bottleneck models for these data, where a population of constant size N_A instantaneously shrinks in size to $b \cdot N_A$ at time t_{start} generations before the present. The population remains at that size for t_{dur} generations and then instantaneously recovers to its original size (Fig. 4, which is published as supporting information on the PNAS web site).

Under the assumption that non-African populations originated from an ancestral population in sub-Saharan Africa, we set the ancestral population size in the bottleneck simulations to the values of N_A obtained by ML based on the Hausa data and the simple growth model ($N_A = 10,659$). This assumption has important implications for our subsequent inferences about compatible bottleneck scenarios. We then used coalescent simulations to estimate the P values for each summary statistic, point on a grid of bottleneck severities (b), bottleneck duration (t_{dur}), and time since the beginning of the bottleneck (t_{start}). This procedure allows defining the portion of the multidimensional parameter space that is compatible with the data.

By combining P values of different summaries as described by Eq.

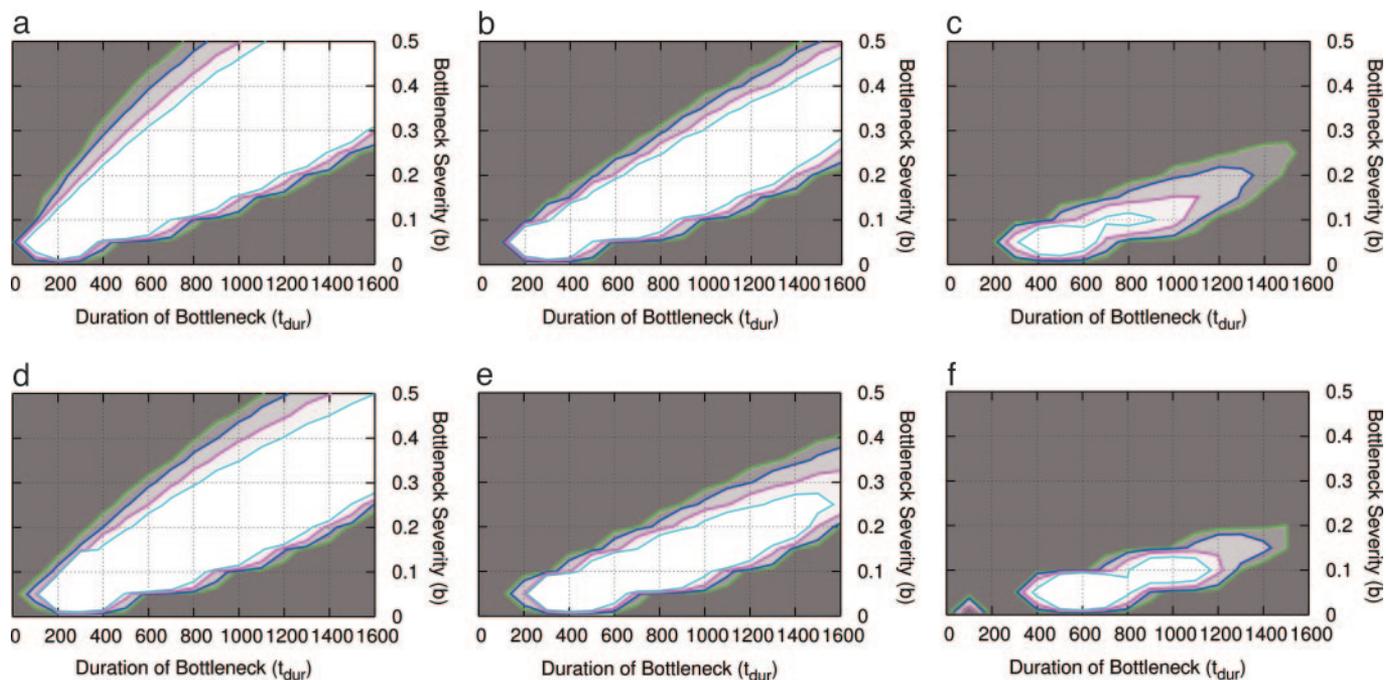


Fig. 2. Confidence sets for a bottleneck with t_{start} of 40,000 years. Results are shown for the Italian (a–c) and Chinese (d–f) data sets for N_A values of 9,450 (a and d), 10,659 (b and e), and 12,300 (c and f). The combined statistics are $\bar{D}-\bar{S}-\hat{p}$. The contours represent the confidence region of parameter space with P values of 0.1, 0.05, 0.02, and 0.01 from innermost to outermost, with darker shading indicating lower P values.

severity), there may be a range of values for the other parameters (e.g., time of onset and duration of bottleneck) that are equally consistent with the data. It is particularly important in this context to make efficient use of the information in the data. Although it may be useful to generate point estimates of the demographic parameters, it is even more important to obtain the multidimensional confidence set if specific hypotheses about human evolution are to be tested.

The present study represents an important step toward improving our inferences about human demography. Although the present data set is not as large as other resequencing surveys (2, 5), it was specifically designed for demographic inference and will provide a useful reference for analyses of gene regions, because, in an attempt to select neutrally evolving regions, we focused on segments that neither contain nor are tightly linked to coding regions. In addition, most of these segments neither contain nor are tightly linked to noncoding sequences conserved between human and mouse. Our scheme for data collection aimed at maximizing the information content of the data so that multiple aspects of genetic variation could be analyzed for the same set of independent loci. Owing to the use of ethnically identified samples, we could provide evidence for different demographic histories in different populations.

Our analytical approach also improves on previous studies of human demography. First, it provides a full characterization of the uncertainty around the best-fitting model by identifying the portion of the multidimensional parameter space that is consistent with genetic variation data in each population. The inclusion of multiple aspects of genetic variation by combining the P values for different summary statistics provides greater power than any single summary alone, allowing us to reduce substantially the accepted space for each model. Our study is based on an extensive exploration of the demographic parameter space including onset, duration, and severity of the bottleneck. It is important to note that the reduction in bottleneck parameter space was greatly aided by our inference about N_A based on the Hausa data. Because the N_A is restricted, the range of compatible values for summary statistics that depend on N_A (i.e., \hat{p} and \bar{S}) is also constrained.

An important limitation of our analysis is that we considered only models of randomly mating populations. Although this is a common assumption in modeling studies of population size change, it is unlikely to be satisfied by human populations, even if geographically defined (34, 51). In fact, it is possible that population structure alone could account for the observed patterns of human variation (2, 5, 15, 35). Interestingly, the addition of $\text{Var}[D]$ into the bottleneck analysis results in a further reduction of the accepted parameter space (Figs. 8–11, which are published as supporting information on the PNAS web site), although combining this statistic with \bar{D} , \bar{S} , and \hat{p} reduces the power to reject the constant size model (Fig. 1). This observation suggests that additional features, such as population structure, are required to produce $\text{Var}[D]$ values that are more consistent with our data. Although it is desirable and certainly more realistic to include elements of population structure in models of human demography (52), there is insufficient data to indicate the most plausible family of such models. For these reasons, testing simple growth and bottleneck models is a reasonable first step toward developing more complex and realistic models. Obviously, if changes in population size and population structure were considered jointly rather than separately, the accepted range of values for the growth and bottleneck parameters is likely to be different.

A main conclusion of this study is that simple bottleneck models can explain the non-African data even when multiple aspects of genetic variation are considered simultaneously. Several previous studies of human sequence variation had modeled specific bottleneck scenarios on the basis of either frequency spectrum information (2, 5, 15, 48), LD decay (18), or polymorphism levels (21). Wall and Przeworski (15) analyzed full resequencing data and proposed that a bottleneck and selective sweeps at some loci could explain the frequency spectrum observed in non-Africans but did not provide information regarding the likely parameter values. The frequency spectrum was used also by Marth *et al.* (20) to estimate a best-fit bottleneck model for Europeans and East Asians. We used our simulation scheme to estimate the probability of the Italian and Chinese data for the corresponding best-fit models of Marth *et al.* (20). In our parameterization, the best fit model for the Asian

