

A comparison of Bayesian methods for  
haplotype reconstruction from population  
genotype data

Matthew Stephens<sup>1</sup> and Peter Donnelly<sup>2</sup>

Address for Correspondence: Matthew Stephens  
Department of Statistics, University of Washington  
Box # 354322, Seattle WA 98195-4322  
Tel (206) 543-4302, FAX (206) 685-7419  
email: [stephens@stat.washington.edu](mailto:stephens@stat.washington.edu)

Running Title: Bayesian haplotype reconstruction

<sup>1</sup>Department of Statistics, University of Washington

<sup>2</sup>Department of Statistics, University of Oxford

## **Abstract**

In this report we compare and contrast three previously-published Bayesian methods for inferring haplotypes from genotype data in a population sample. We review the methods, emphasising the differences between them in terms of both the models (“priors”) they use, and the computational strategies they employ. We introduce a new algorithm that combines the modelling strategy of one, with the computational strategies of another. In comparisons on real and simulated data, this new algorithm outperforms all three existing methods. The new algorithm is included in the software package PHASE, version 2.0, available from <http://www.stat.washington.edu/stephens/software.html>

Current high-throughput genotyping technologies, when applied to DNA from a diploid individual, are able to determine which two alleles are present at each locus, but not the haplotype information: that is, which combinations of alleles are present on each of the two chromosomes. Knowledge of the haplotypes carried by sampled individuals would be helpful in many settings – including linkage disequilibrium mapping, and inference of population evolutionary history – essentially because genetic inheritance operates through the transmission of chromosomal segments. Experimental methods for haplotype determination exist, but are currently time-consuming and expensive. Statistical methods for inferring haplotypes are therefore of considerable interest. While in some studies data may be available on related individuals to assist in this endeavour, in general such data may be either unavailable, or only partially informative. We focus here on the problem of statistically inferring haplotypes from unphased genotype data for a sample of (“unrelated”) individuals from a population.

Several approaches to this problem have been proposed, notably Clark’s algorithm (Clark 1990), and maximum likelihood estimation of haplotype frequencies via the EM algorithm (e.g. Excoffier and Slatkin 1995). Stephens et al. (2001a) (henceforth SSD) introduced two Bayesian approaches, one (their *Algorithm 2*, or “naive Gibbs sampler”) which used a simple Dirichlet prior distribution, and a second, more sophisticated approach (their *Algorithm 3*), in which the prior approximated the coalescent. Results on simulated SNP and microsatellite data, and more limited comparisons using real data (Stephens et al. 2001b), suggested that this second approach, implemented in the software PHASE, produced consistently more accurate hap-

lotype estimates than previous methods. In addition SSD point out other advantages of a Bayesian approach to this problem, including the ability to provide accurate measures of uncertainty in statistically estimated haplotypes, which in principle could be exploited in subsequent analyses (although practical considerations mean that this has seldom been fully exploited in practice).

More recently, two other Bayesian approaches to this problem have been published: Niu et al. (2002) (henceforth NQXL) introduced an algorithm, which they refer to as PL, implemented in the software HAPLOTYPER; and Lin et al. (2002) (henceforth LCZC) also introduced a Bayesian algorithm, which they generously attribute to us, but which nevertheless differs in substantive ways from the algorithms in SSD.

Here we highlight the conceptual differences between these different Bayesian methods, some of which may be unclear from the original papers. We note that the main contribution of NQXL — the introduction of computational strategies to greatly reduce running times — can also be applied to the other algorithms, and we describe a new version of PHASE that exploits these strategies. In our comparisons on datasets considered by NQXL and LCZC, this new version of PHASE outperforms the other two methods. Our comparisons also demonstrate that the apparently inferior performance of PHASE compared to HAPLOTYPER in some of NQXL’s comparisons was not, as they suggest, due to the fact that the data sets considered in these comparisons deviated from the implicit (coalescent-based) modelling assumptions underlying PHASE. Rather, it was due to the fact that PHASE required longer runs than NQXL employed to produce reliable results.

**Components of Bayesian approaches.** Bayesian haplotype reconstruction methods treat the unknown haplotypes as random quantities and combine

- *Prior information* — beliefs about what sorts of patterns of haplotypes we would expect to observe in population samples, with
- *The likelihood* — the information in the observed data,

in order to calculate the *posterior distribution*, the conditional distribution of the unobserved haplotypes (or haplotype frequencies) given the observed genotype data. The haplotypes themselves can then be estimated from this posterior distribution, for example by choosing the *a posteriori* most likely haplotype reconstruction for each individual.

In Bayesian approaches to complicated statistical problems, it is helpful, conceptually, to distinguish two separate issues.

I The *model* or *prior distribution* for the quantities of interest, in this case for population haplotype frequencies. For a given data set, different prior assumptions will in general lead to different posterior distributions, and hence to different estimates.

II The computational algorithm used. For challenging problems, including this one, the posterior distribution cannot be calculated exactly. Instead, computational methods – typically Markov chain Monte Carlo (MCMC) – are used to approximate it. Different computational tricks, or different numbers of iterations, will change the quality of approximations to the Bayesian answer.

The three Bayesian approaches we consider here differ in both the prior *and* the computational algorithms used, as we now outline.

**Prior distributions.** SSD described two algorithms based on two different prior distributions for the haplotype frequencies: the first “naive Gibbs sampler” used a Dirichlet prior distribution; the second, implemented in PHASE, used a prior approximating the coalescent. In their comparisons the algorithm based on the approximate coalescent prior substantially outperformed the algorithm based on the Dirichlet prior. The subsequent algorithms of NQXL and LCZC are each based on the Dirichlet prior.

Interestingly, NQXL and LCZC attribute rather different properties to the Dirichlet prior. NQXL state that their method “imposes no assumptions on the population evolutionary history”. In contrast, LCZC attribute the success of their method to the fact that the “neutral coalescent model, which [it] incorporates, is a reasonable approximation of the random collection of human sequences used as test data”. The truth, we suggest, lies somewhere in between. The Dirichlet prior arises naturally in genetics models with so-called “parent-independent” mutation (Stephens and Donnelly 2000) – that is, when the genetic sequence of a mutant offspring does not depend on the progenitor sequence. This assumption on the mutation process does not apply (even approximately) to DNA sequence data, nor to data at multiple SNP or microsatellite loci. Thus the use of a Dirichlet prior can be thought of as making simple, but highly unrealistic, assumptions about the genetic processes underlying the evolution of the study population. In contrast the approximate coalescent prior used in SSD is based on the arguably more complex, but decidedly more realistic, assumption that the genetic sequence

of a mutant offspring will differ only slightly from the progenitor sequence (often by a single base change).

We can informally illustrate an important operational difference between the Dirichlet prior, and what we have called an approximate coalescent prior, as follows. Sometimes an unresolved genotype can be broken up into two haplotypes, one or both of which is already known (or assumed) to be present in the sample. Both priors will put substantial weight on this possibility. In contrast, suppose that an unresolved genotype cannot be broken up in such a way, but that it can be broken up so that both haplotypes are similar to, but not identical to, known haplotypes (where “similar to” here means that one haplotype can be formed from the other by one or a small number of single base changes). The approximate coalescent prior will put substantial weight on this reconstruction, but the Dirichlet prior will choose randomly between all possible reconstructions, giving no additional weight to the one involving haplotypes similar to those already seen. Analogous problems occur with Clark’s method and the EM algorithm. (Indeed, the maximum likelihood estimate for haplotype frequencies, which the EM algorithm aims to find, corresponds to the mode of the posterior distribution for a particular Dirichlet prior. In this sense the EM algorithm gives the same answer as a Bayesian procedure with an unrealistic prior.)

Whatever one’s view on the accuracy of the coalescent as a model for real data, it is difficult to imagine any actual population sample where guessing the haplotypes at random will be more accurate than choosing haplotypes that are similar to others in the sample. Indeed, the main innovation in LCZC can be thought of as an *ad hoc* modification of the Dirichlet prior to

avoid this undesirable “guessing-at-random” behaviour. Their modification is that, when considering whether an individual’s genotypes can be resolved into haplotypes that match other haplotypes in the sample, they look for matches *only* at positions where the individual is heterozygous, ignoring the data at positions where the individual is homozygous. As a consequence of this, the algorithm never reaches the situation considered above where no “matching” haplotypes exist, and therefore avoids choosing randomly between all possible reconstructions. This modification has certain computational advantages over the approximate coalescent prior – in particular LCZC exploited a computational trick from the naive Gibbs sampler in SSD to produce an efficient algorithm. However, as our comparisons below demonstrate, the resulting algorithm is less accurate than one based on the approximate coalescent prior. This is because the genotypes at positions where an individual is homozygous carry potentially-valuable information about the phase relationships at the other (heterozygous) positions – information that is exploited by the approximate coalescent prior when looking for close-matching haplotypes.

Whether the posterior distribution for one prior will provide better estimates than the posterior distribution for a different prior will depend on which of the priors does a better job of capturing features of the real data. We continue to believe, both on general population genetics grounds, and on the evidence of the superior performance of the PHASE algorithm on comparisons here, and in SSD and Stephens et al. (2001b), that the use of an approximate coalescent prior will lead to better estimates than the use of a Dirichlet prior (even with LCZC’s modification).

**Computational approach.** Both algorithms in SSD, and the algorithm in LCZC, used relatively unsophisticated MCMC algorithms based on Gibbs sampling. An important innovation in NQXL’s work is the introduction of two computational tricks — *prior annealing* and *partition-ligation* — for reducing the computational effort required to obtain a good approximation to the true posterior distribution. These ideas are largely independent of the prior used, and similar ideas can be applied to the approximate coalescent prior used in PHASE (and could be applied to the modified Dirichlet prior of LCZC), as we outline below. Qin et al. 2002 apply similar ideas to make the EM algorithm computationally tractable for large data sets.

Our discussion above may appear to construct a rather concrete divide between models (priors) and computation. This is deliberate: we want to emphasise the distinct role that each of these components can play in the quality of the final solution obtained. However, in practice there is often a strong interaction between these two components of a Bayesian analysis. Indeed, the algorithm implemented in PHASE was not actually developed in the conventional way of writing down a prior and likelihood, and then developing a computational method for sampling from the corresponding posterior (and neither, incidentally, was the algorithm in LCZC). Rather the posterior is defined implicitly, as the stationary distribution of a particular Markov chain, which in turn is defined via a set of (inconsistent) conditional distributions. Although defining posterior distributions in this way is not without its potential pitfalls, the algorithm in SSD was designed to circumvent these (see appendix). Further, this unconventional approach has the advantage of avoiding some of the computational difficulties of sampling from the poste-

rior corresponding to an exact coalescent prior, whilst capturing the salient features of such a prior, notably the tendency for haplotypes in a population to be similar to other haplotypes in the population. Although NQXL claim that the “pseudoposterior probabilities” that our “pseudo-Bayesian” algorithm attaches to the constructed haplotypes are difficult to interpret, simulation results in Stephens et al. (2001a) show these probabilities to be relatively well calibrated relative to a coalescent prior, even in the presence of moderate amounts of recombination.

**A modified version of PHASE.** We now outline a modified version of PHASE, that continues to make use of an approximate coalescent prior, but exploits ideas from NQXL to improve computational efficiency, and to increase the size of problem that can be handled. For convenience, in the following description we use “frequency” to refer to relative frequency.

As in NQXL, our algorithm follows a “divide and conquer” strategy, of initially estimating haplotype frequencies within short blocks of consecutive loci (SNPs), before successively combining estimates for adjacent blocks to obtain estimates of haplotypes across the whole region under consideration. Note that we are using the term “block” to refer simply to a set of consecutive loci, with no implication about the patterns of linkage disequilibrium present. Results from NQXL suggest that the way in which the block boundaries are chosen is relatively unimportant: to encourage independence of results from multiple runs of the algorithm we randomly chose the length of each block to be 6, 7, or 8 loci, with probabilities 0.3, 0.3 and 0.4 respectively. To each block we applied Algorithm 3 from SSD, with the following alterations:

- i) we updated each individual in turn, in a random order (with a different

- random order for each sweep);
- ii) when updating an individual, we updated all ambiguous loci in the block under consideration, rather than choosing 5 at random (we consider a locus to be ambiguous in a particular individual if the individual is either heterozygous, or is missing one or both alleles at that locus);
  - iii) to improve mixing during burn-in iterations, with probability  $\beta$  (whose value is specified below) we computed the probability of each haplotype pair as being proportional to the sum, rather than the product, of the appropriate conditional probabilities. This modification makes the algorithm more likely to visit configurations in which only one of the two haplotypes is similar to other haplotypes in the sample, thus improving mixing of the MCMC scheme.

The value of  $\beta$  was decreased linearly from 1.0 to 0.0 over 100 “burn-in” iterations (where one iteration means updating every individual once), and was then fixed at zero for 100 further iterations. During these further 100 iterations, for each haplotype that could possibly occur in the sample, we obtained a (Rao-Blackwellized) estimate of the posterior mean of its frequency in the sample.

The above procedure results in an estimate of the haplotype frequencies within each short block. We then apply a variant on the idea of progressive ligation from NQXL to iteratively combine consecutive blocks. NQXL suggest taking from each block the  $B$  haplotypes with the highest estimated frequencies, and forming a list,  $L$  say, of the  $B^2$  possible concatenated haplotypes. We follow this suggestion, but rather than taking a fixed value of  $B$

(as NQXL seem to suggest), we choose  $B$  separately for each block, in such a way to include all haplotypes whose estimated sample frequency is  $f/2n$  or greater, where  $n$  is the number of diploid individuals in the sample, and we took  $f = 0.001$ . (Bigger values of  $f$  result in shorter lists, and hence faster runs, at the cost of a potential decrease in the accuracy of the approximation to the posterior distribution.) Once we have formed  $L$ , we obtain new estimates for the haplotype frequencies within the newly-created block by applying the same MCMC algorithm described above (including the burn-in with linearly decreasing  $\beta$  from 1.0 to 0.0) to the new block, allowing each individual to be made up only of pairs of haplotypes in  $L$ . We then continue this ligation procedure, each time concatenating the last-formed block with the adjacent small block. When all blocks have been combined (into a single final block containing all loci), we estimate each individual’s pair of haplotypes by its posterior mode obtained from the final 100 iterations.

The algorithm above includes several variables (particularly  $f$ , and the number of iterations) whose values will affect both run times, and the reliability of the approximation to the posterior distribution. The particular values we used were chosen so that, in preliminary tests on a few of the data sets, multiple runs of the algorithm from different starting points typically gave similar haplotype estimates. To aid in the comparison of results of different runs we monitored the value of a “pseudo-likelihood” (Besag 1974), defined as

$$\prod_{i=1}^n \sum_{h_1 \in L} \sum_{h_2 \in L} \Pr(h_1, h_2 | H_{-i}) I((h_1, h_2) \text{ consistent with } G_i),$$

where  $H_{-i}$  is the set of all haplotypes in the current MCMC configuration, excluding the  $i$ th individual,  $G_i$  is the genotype of the  $i$ th individual, and

$I(\cdot)$  is the indicator function.

This pseudo-likelihood can be thought of as providing a measure of the goodness of fit of the estimated haplotypes to the underlying model. When different runs give very different values for the goodness of fit this suggests that the runs may be too short to provide reliable results. Further, among multiple runs on the same dataset we would expect those with the highest values of this pseudo-likelihood (averaged over the final 100 iterations, say) to provide the more accurate results.

Our preliminary tests suggested that, with the parameter values we used, multiple runs of the algorithm did occasionally produce results with rather different values of the pseudo-likelihood, suggesting that the algorithm sometimes converged to a local, rather than global, mode of the posterior distribution. In our first set of comparisons below, to alleviate this problem, we ran the whole algorithm on each data set 5 times independently, and chose the solution corresponding to the run that maximises a pseudo-likelihood averaged over the final 100 iterations. In our second set of comparisons, to reduce computation, we ran the algorithm on each data set only once – we would expect a multiple-run strategy to slightly improve average accuracy.

**Comparisons.** Our first set of comparisons is based on similar comparisons made by NQXL, who ran PHASE and HAPLOTYPER on several data sets, and found that HAPLOTYPER performed more accurately in many cases. We compared

1. The algorithm from LCZC (using code kindly provided by S. Lin) run at its default run length;
2. HAPLOTYPER run at its default settings;

3. PHASE v1.0 (which implements Algorithm 3 in SSD) run at its default settings; and
4. the modified version of PHASE described here;

on several data sets for which NQXL found PHASE performed poorly. We used two different criteria for assessing accuracy:

- (a) the error rate, as defined by NQXL, namely the proportion of individuals whose haplotype estimates are not completely correct.
- (b) a more stringent measure of accuracy, which measures the similarity between the estimated haplotypes and the true haplotypes. Specifically, we count how many individual nucleotides must be changed in the estimated haplotypes to make them the same as the known haplotypes, and divide this by the largest value it could possibly take (given the genotype information) to obtain a number between 0 and 1.

We describe the second measure as “more stringent” because it makes a more detailed comparison between the estimated and true haplotypes, rather than simply determining whether each estimate is correct or incorrect. This measure can thus discriminate between methods even in cases where it is unrealistic to expect a statistical method to completely determine haplotypes at every site, as may be the case for many real data sets, particularly those including low-frequency alleles, or sites/loci spread over a large genetic distance.

Table 1A gives, for each type of data and for each method, the mean individual error rate ((a) above). By this measure of accuracy, the modified

version of PHASE and HAPLOTYPYER perform similarly. PHASE v1.0 run at its default values performs considerably better than in NQXL, but perhaps slightly less well than the modified version. Note however that the apparently large difference in error rates for the ACE data (0.18 vs 0.28) actually corresponds to making an error on just one additional individual. Somewhat surprisingly, LCZC’s algorithm performed consistently less well than the other methods, most notably on the simulated data. Runs 100 times longer than the default settings produced almost identical average performance for these simulated data sets (results not shown). Nevertheless, computational problems may still be (partly) responsible for the poor performance of the algorithm in these data sets, and performance might be improved by the use of a more sophisticated computational scheme.

Table 1B summarizes the performance of each algorithm on the more stringent criterion (b) above, which we call the “single-site error rate”. By this measure of accuracy, the modified version of PHASE consistently, and substantially, outperforms both HAPLOTYPYER and LCZC’s algorithm on these data sets. This indicates that when the methods are unable to reconstruct the haplotypes completely, the PHASE estimated haplotypes tend to be much more similar to the true haplotypes – presumably because the true haplotypes conform more closely to the assumptions of the approximate coalescent prior than to those of the Dirichlet prior. In particular, it seems that the apparently inferior performance of PHASE in NQXL was not due to its sensitivity to deviations from the assumptions of the coalescent model (as NQXL suggest), but rather due to the fact that the 5,000 updates they used (compared with the default of 2,000,000 updates, which we used here),

were insufficient for the algorithm to provide a reasonable approximation to the posterior distribution.

For our second set of comparisons we examine the performance of the algorithm from LCZC, HAPLOTYPYER, and the modified version of PHASE for the data sets in LCZC. (PHASE v1.0 is omitted from these comparisons, due to its high computational demands for data sets of this size.) For ease of comparison, we use here two measures of accuracy based on those in LCZC:

- (a) the error rate, as defined by SSD, namely the proportion of *ambiguous* individuals whose haplotype estimates are not completely correct. Note that although this differs from the definition of error rate in NQXL used above, methods that perform well by one of these criterion will tend also to perform well by the other.
  
- (b) the switch error, which measures the proportion of heterozygote positions whose phase is wrongly inferred relative to the previous heterozygote position. This differs qualitatively from the single-site error rate used above in that it does not depend on the accuracy of a method in inferring longer-range phase relationships, and so is perhaps most appropriate where these longer-range phase relationships cannot be accurately inferred by statistical means (which may be the case for some of these data sets). Note that the switch error is  $(1 - \text{the switch accuracy defined by LCZC})$ . We make this change to an error rate, so that, like the other measures we use, small values indicate accurate haplotype estimates.

Because these data sets have some missing genotypes, not all the true

haplotypes are completely known, and so, strictly, it is not actually possible to compute either of these criteria. To finesse this problem LCZC scored phase calls in each individual only at sites where neither allele was missing (S. Lin, personal communication). To allow comparisons with LCZC's results we take the same approach here. Table 2 shows the performance of each of the methods by both criteria. The modified version of PHASE appreciably outperforms the other two methods, again presumably due to the greater accuracy of the approximate coalescent prior.

Finally, we note that LCZC made additional comparisons between their method, HAPLOTYPED, and the EM algorithm using only the common variants (minor allele frequency  $> 0.2$ ) on the same datasets, and found that their algorithm outperformed the others. For these data sets, which naturally contain many fewer SNPs than the full data sets, and where all three algorithms perform better in absolute terms, PHASE and LCZC's algorithm perform more similarly (results not shown).

**Discussion.** The estimation of haplotypes from population data, for the sizes of data sets currently being generated, and those likely in the context of the proposed Haplotype Map project, is a challenging problem, which requires sophisticated computational methods. It appears that Bayesian approaches have much to offer, firstly in terms of accuracy of estimation, but also in their ability (i) to incorporate, in a natural way, features such as genotyping error, missing data, or additional information, for example from pedigrees; and (ii) to provide a coherent framework in which to account for the uncertainty associated with estimates of haplotypes or haplotype frequencies in later analyses. Amongst Bayesian approaches, the comparisons

reported here and elsewhere suggest that PHASE provides the most accurate reconstructions.

The modified version of PHASE reported here very substantially reduces the computational time involved in running the method. For example, on our desktop machine with a 800MHz Pentium III processor, the implementation we used for our comparisons took roughly 30 minutes of CPU time per dataset for the largest gene used in the second set of comparisons (TRPC5), which consisted of 20 diploid individuals typed at 165 SNPs. Shorter runs taking roughly 2 minutes each produced almost identical average accuracy for this gene, suggesting that in this case our choice of run-length was conservative. Although these times exceed the 10 seconds LCZC quote for their algorithm on the same data set, or the 35 seconds it takes HAPLOTYPED (with ROUNDS set to 20) on our machine, it is clear that even much larger problems will remain well within the bounds of practicality. Further, if necessary the efficiency of our current implementation could be improved in several ways. However, in our view other aspects of the problem deserve more urgent attention. For example, all current methods, including PHASE, ignore the decay of linkage disequilibrium with distance between markers; a new version of PHASE developed by one of us (MS) allows for this, and results in still more accurate haplotype estimates (details will be published elsewhere). Further, in most applications, estimating haplotypes, or even population haplotype frequencies, will not be the ultimate goal. To fully capitalise on Bayesian methods for haplotype reconstruction it will be necessary to integrate the analysis of the haplotypes – be it testing for association with a disease phenotype, or estimating recombination rates, for example –

together with the haplotype estimation procedure, to fully allow for uncertainty in the haplotype estimates.

**Acknowledgements** This work was supported by NIH grant 1RO1HG/LM02585-01 to MS.

**Software** The algorithm described in this report is included in the software package PHASE, version 2.0, available from

`http://www.stat.washington.edu/stephens/software.html`.

## A Appendix

In the text we note that there are potential pitfalls associated with the fact that the posterior distribution sampled from by PHASE (both the version implementing the algorithm from SSD, and the new version described here) is defined implicitly, as the stationary distribution of a particular Markov chain, which in turn is defined via a set of inconsistent conditional distributions. Here “inconsistent” means that there is no joint distribution that has these conditional distributions. The fact that these conditional distributions are inconsistent is potentially problematic, as a Gibbs sampler based on inconsistent conditional distributions is not, in general, guaranteed to converge to a proper probability distribution. However, in this case convergence to a proper distribution *is* guaranteed, because the Markov chain has a finite state space (the space of all possible haplotype reconstructions) and is irreducible and aperiodic. (All such Markov chains have a stationary distribution, and converge to this stationary distribution; e.g. Theorem 7.4 in Behrens 2000).

A second potential technical problem with using inconsistent conditional distributions in Gibbs sampling is that, using the standard “fixed scan” approach to Gibbs sampling, where each individual is updated in turn in some fixed order, the stationary distribution could depend on the order used. This seems undesirable, and so to avoid this SSD used a “random scan” Gibbs sampler, where at each iteration a random individual is chosen for updating (with each individual being equally likely). In this paper we used a different, and perhaps slightly preferable, random scan strategy, where at each iteration all individuals are updated in a random order, with a different random order for each iteration. Both schemes clearly ensure that the stationary

distribution is independent of the order that the individuals are input into the algorithm.

## References

- Behrends E (2000). Introduction to Markov chains : with special emphasis on rapid mixing. Vieweg: Braunschweig/Wiesbaden
- Besag J (1974). Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society, series B* 36:192–236
- Clark AG (1990). Inference of haplotypes from PCR-amplified samples of diploid populations. *Molecular Biology and Evolution* 7(2):111–122
- Drysdale C, McGraw D, Stack C, Stephens J, Judson R, Nandabalan K, Arnold K, Ruano G, Liggett S (2000). Complex promoter and coding region  $\beta_2$ -adrenergic receptor haplotypes alter receptor expression and predict in vivo responsiveness. *Proceedings of the National Academy of Science, USA* 97:10483–10488
- Excoffier L, Slatkin M (1995). Maximum-Likelihood Estimation of Molecular Haplotype Frequencies in a Diploid Population. *Molecular Biology and Evolution* 12(5):921–927
- Kerem B, Rommens J, Buchanan J, Markiewicz D, Cox T, Chakravarti A, Buchwald M, Tsui LC (1989). Identification of the cystic fibrosis gene: genetic analysis. *Science* 245:1073–1080
- Lin S, Cutler DJ, Zwick ME, Chakravarti A (2002). Haplotype Inference in Random Population Samples. *American Journal of Human Genetics* 71:1129–1137
- Niu T, Qin ZS, Xu X, Liu JS (2002). Bayesian Haplotype Inference for Multiple Linked Single-Nucleotide Polymorphisms. *American Journal*

of Human Genetics 70:157–169

Qin ZS, Niu T, Liu JS (2002). Partial-Ligation-Expectation-Maximisation for Haplotype Inference with Single Nucleotide Polymorphisms. *American Journal of Human Genetics* 71:1242–1247

Rieder MJ, Taylor SL, Clark AG, Nickerson DA (1999). Sequence variation in the human angiotensin converting enzyme. *Nature Genetics* 22:59–62

Stephens M, Donnelly P (2000). Inference in molecular population genetics. *Journal of The Royal Statistical Society, Series B* 62:605–655

Stephens M, Smith NJ, Donnelly P (2001a). A new statistical method for haplotype reconstruction from population data. *American Journal of Human Genetics* 68:978–989

Stephens M, Smith NJ, Donnelly P (2001b). Reply to Zhang et al. *American Journal of Human Genetics* 69:912–914

	$\beta_2AR^1$	ACE <sup>2</sup>	CFTR <sup>3</sup>	Simulated data <sup>4</sup>
A) Individual Error Rates:				
LCZC	0.18	0.31	0.54	0.40
HAPLOTYPER	0.09	0.19	<b>0.40</b>	<b>0.020</b>
PHASE	<b>0.04</b>	0.28	0.46	0.068
modified PHASE	0.05	<b>0.18</b>	0.47	0.045
B) Single-Site Error Rates:				
LCZC	0.19	0.11	0.43	0.30
HAPLOTYPER	0.11	0.11	0.66	0.031
PHASE	<b>0.03</b>	0.10	<b>0.36</b>	0.047
modified PHASE	<b>0.03</b>	<b>0.03</b>	0.37	<b>0.028</b>

Table 1: Mean Individual and Single-Site error rates (see main text for definitions). Each number in the table is an average over 100 (or, for the last column, 20) data sets. The results for the best-performing method in each column are highlighted in **bold**.

## Notes for Table 1

1. These data (Drysdale et al. 2000), were used by NQXL to explore sensitivity of methods to deviations from Hardy Weinberg Equilibrium (HWE). We simulated 100 data sets, each of 15 individuals, by pairing randomly-chosen haplotypes according to NQXL's "strong heterozygote favoring" model, as described in their paper. Each of the 100 data sets contained either 0, 1 or 2 homozygotes, which were the cases where NQXL saw the poorest performance for PHASE compared with HAPLOTYPER. Although NQXL suggest that excess heterozygosity might result from a selective advantage for heterozygotes, selection will not cause deviations from HWE unless the fitness differences are *very* extreme (e.g. lethal recessives).
2. These data (Rieder et al. 1999) were used by NQXL to test the stability of algorithms. As in NQXL, we ran each algorithm 100 times on the known genotypes, each run using a different initial value for the seed of the random number generator.
3. Cystic Fibrosis data from Kerem et al. (1989). As in NQXL, we randomly permuted the subset of 57 haplotypes with no missing data 100 times, to generate 100 data sets, each containing 28 hypothetical individuals.
4. 20 datasets simulated under NQXL's bottleneck model, each containing data for 20 individuals at 20 loci, were kindly provided by Z.S. Qin. We understand (Z.S. Qin, personal communication) that NQXL's bottleneck simulations made the assumption that in the bottleneck pop-

ulation all loci were in complete linkage equilibrium. This is a non-standard assumption in this context, and seems likely to produce simulated haplotypes that exhibit very different patterns to those expected under what we would consider more plausible assumptions.

	GLRA2	MAOA	KCND1	ATR	GLA	TRPC5	BRS3	MECP2
A) Error Rate:								
LCZC	.79	.61	.54	.62	.89	<b>.58</b>	.72	.85
HAPLOTYPYPER	.89	.76	.72	.72	.79	.72	.79	<b>.64</b>
PHASE	<b>.76</b>	<b>.54</b>	<b>.46</b>	<b>.45</b>	<b>.68</b>	<b>.58</b>	<b>.67</b>	.77
B) Switch Error:								
LCZC	.14	.10	.22	.29	.22	<b>.13</b>	.14	.23
HAPLOTYPYPER	.16	.12	.27	.32	.16	.20	.15	.19
PHASE	<b>.10</b>	<b>.07</b>	<b>.13</b>	<b>.18</b>	<b>.11</b>	<b>.13</b>	<b>.10</b>	<b>.15</b>

Table 2: Individual error rate and switch error (see main text for definitions) for the data sets considered by LCZC. The results for HAPLOTYPYPER and the algorithm in LCZC are taken from LCZC’s Table 1. The results for PHASE were obtained by us on 100 data sets simulated in the same way as those used to produce Table 1 in LCZC (i.e. by randomly pairing the 40 X-chromosome haplotypes used by LCZC, kindly provided by D. Cutler). Each number in the table is based on results for 100 data sets. For example, the error rates are the total number of mistakes made across all 100 data sets, divided by the total number of ambiguous individuals in all 100 data sets. The results for the best-performing method in each column are highlighted in **bold**.