



# A Simulation Study of Genotype Phasing Software



Dent A. Earl<sup>1</sup>, David Haussler<sup>1,2,3</sup>

1. Bioinformatics Graduate program, UCSC 2. Center for Biomolecular Science and Engineering, UCSC, 3. Howard Hughes Medical Institute

BIOINFORMATICS  
BIOMOLECULAR ENGINEERING DEPT



## Abstract

Here we report the results of a simulation study of three genotype phasing software packages. The three were tested on ten simulated data sets built using an identical pedigree structure, a high resolution recombination map, and real genotypes taken from HapMap III data for two populations: Utah residents with ancestry from northern and western Europe (CEU, 88 individuals) and Yoruba in Ibadan, Nigerian (YRI, 100 individuals).

The most surprising result shown here is that though none of the methods tested were provided with information on the underlying pedigree that gave rise to the data, they were all affected indirectly by the structure: the algorithms performed better on individuals that were more related to everyone else in the data set (their expected number of identical genotypes in the data was large).

## Introduction

Diploid organisms have, by definition, homologous pairs of chromosomes. An individual's genotype is a combination of the individual's underlying genetic haplotypes, a haplotype being the sequence of alleles that are physically connected to one another along one chromosome. Without labor, and cost intensive methods (cloning and sequencing), haplotypes are not directly observable, and instead must be inferred through a process called *phasing*.

Phasing has experienced a large increase in attention in recent years, likely due to the boom in single nucleotide polymorphism (SNP) microarray market penetration. SNP microarrays (SNP chips) allow for the observation of an individual's genotype at predefined loci. Unfortunately SNP genotypes contain no direct information about the sequence of the observed alleles, information which is crucial for disease, population genetic, and evolutionary studies.

Two endeavors rely on haplotype information: clinical assays and evolutionary reconstruction. In a clinical assay for a cis-acting genetic disease a clinician would like to determine what the chromosome-local pattern of inherited alleles is for a patient.

For evolutionary reconstruction the goal of phasing genotype data into haplotypes is to understand and record the pattern of transmission of alleles through individuals over generations.

These endeavors are not mutually exclusive, but complementary. For example, knowing the genetic ancestry of an individual could also help to inform the clinician about the potential efficacy of certain drugs. And knowing the underlying molecular genetics and epidemiological genetics of a particular region of the genome could be illuminating in reconstructing the evolutionary history of that region.

## Methods

**PEDIGREE**  
The structure of the pedigree that we simulate, shown in Figure 1, was chosen in an attempt to have a gradient of difficulty of phasing problems, from the easy when an individual has one chromosome from one population and one from another, to the more difficult where the individual's chromosomes come from the same population. With this in mind I created a symmetrical pedigree with founding haplotypes taken directly from the source populations, admixture from descendants of the different populations, and with exogamous individuals from the ancestral populations breeding into the descendant family.

**FOUNDERS**  
Haplotypes are chosen at random without replacement to make up founding members of the pedigree and subsequent exogamous individuals. Founders for the CEU side of the family are I-a and I-b and for the YRI side of the family, I-c and I-d. The subsequent exogamous individuals in the CEU side are II-a and III-c and for the YRI side are II-f and III-h.

A large pool of software with varying abilities has been released over the past decade, but there is no standard methodology for phasing and no clear gold standard algorithm at the moment.

Following the simulation of genotype data, I test the major phasing software packages as to their abilities to resolve the simulated genotype data.

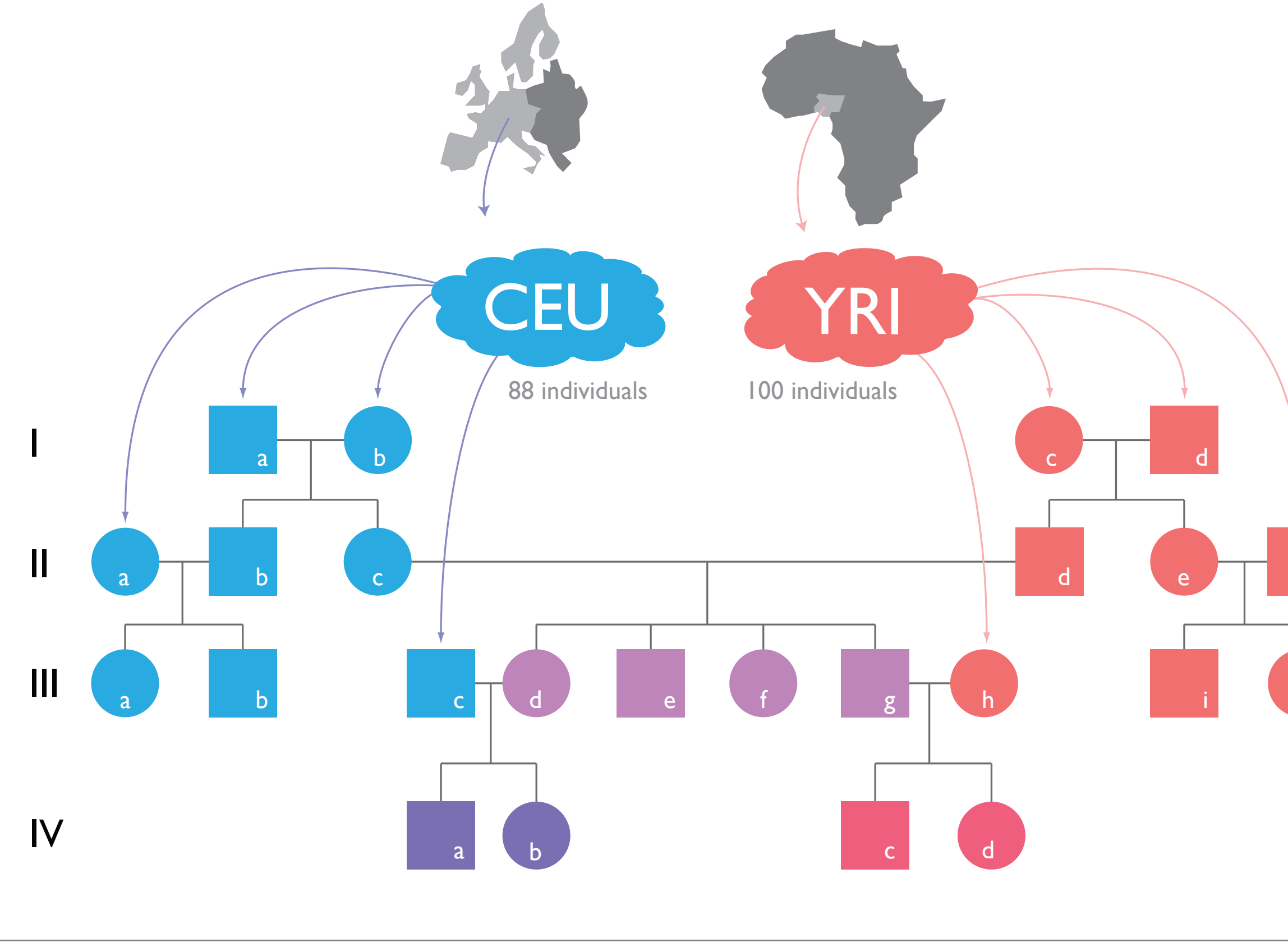
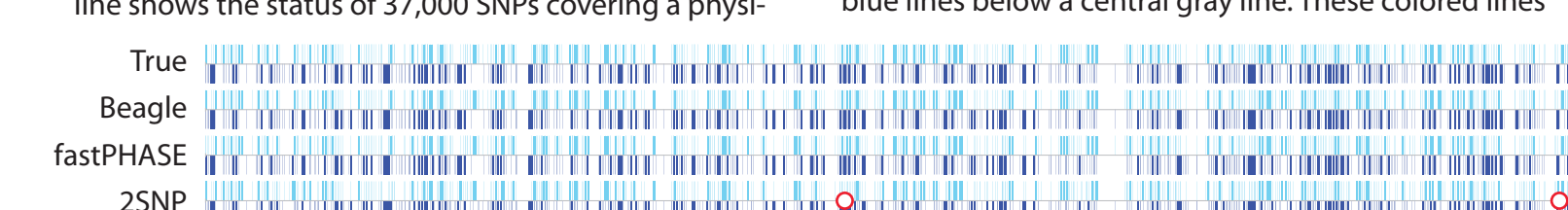
**SOFTWARE**  
fastPHASE, published in 2006, was developed as a compromise between accuracy and speed. fastPHASE does this by omitting the MCMC scheme of PHASE in favor of a hidden Markov model (HMM). As the authors note, however, the price paid here is that the HMM incorporates no information about demographics, or evolutionary processes.

Beagle, published in 2007, uses a directed acyclic graph (DAG) to model localized haplotype-clusters and then an HMM to find the most likely haplotype pair, conditional on an individual's genotype. The algorithm allows the use of known haplotypes in a panel, which is used to populate the paths of the HMM in the localized haplotype-cluster. The final step of the algorithm is to use the Viterbi algorithm to find the most-likely haplotype pair for each individual.

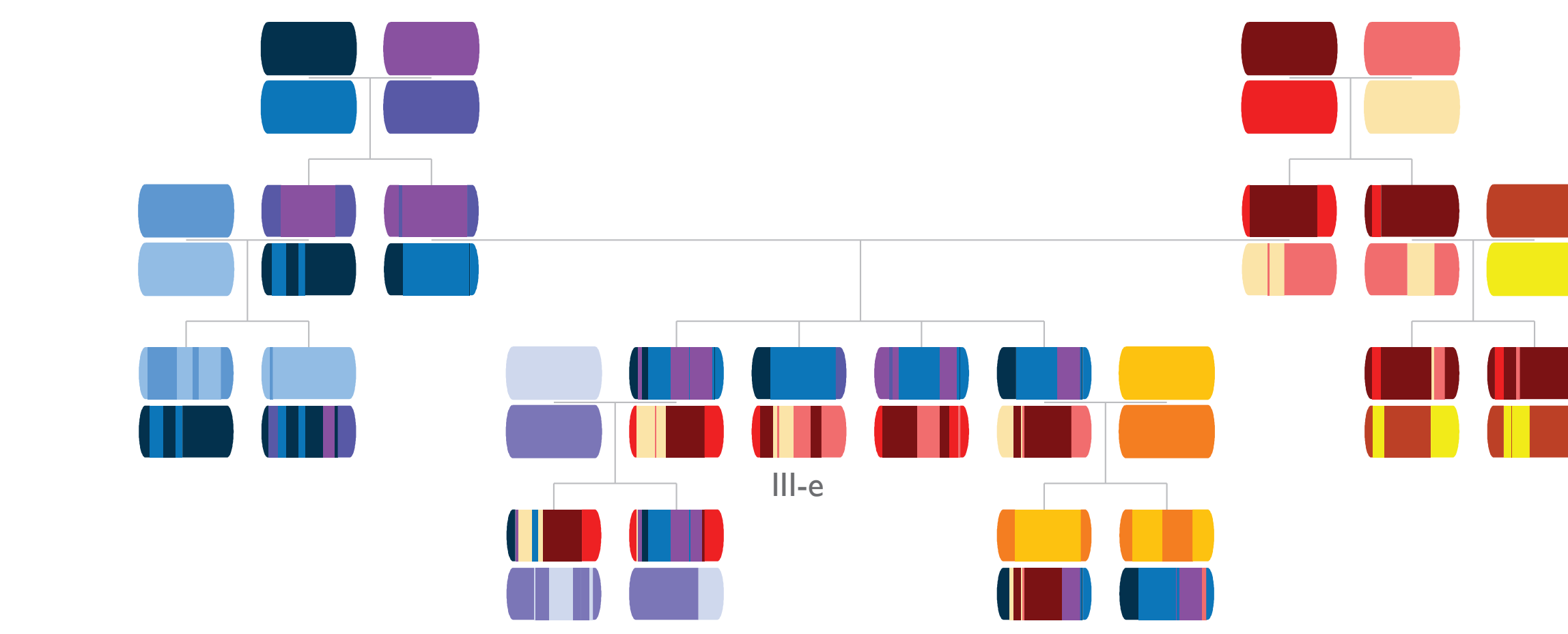
2SNP, published in 2008, when presented with trio data (father, mother and child) uses a combination of approaches to phasing. It enforces simple patterns of Mendelian inheritance where possible, and elsewhere uses a graph structure approach and a maximum spanning tree based on the population genetic metrics of linkage disequilibrium (LD) and Hardy-Weinberg equilibrium (HWE) to assign haplotypes. Essentially, the algorithm constrains the problem based on the known demographic structure and focuses solely on complex positions.

**SWITCH ERROR**  
In order to quantitatively assess the performance of the algorithms I implemented the metric switch error, which is one minus the switch accuracy of Li et al. (2004). Given the correct (True) phase for a particular segment of a chromosome, and a solution given by a phasing algorithm, the switch error is calculated to be the number of times the phase pattern of adjacent alleles is not the same in the test as in the True, divided by the total number of heterozygous positions minus one. Switch error ranges from 1 to 0.

**FIGURE 4**  
Haplotype line plots representing the true and inferred states of a simulated individual's genotype (individual III-e from the first simulation shown in Figure 2). Each line is made up of cyan lines above and blue lines below a central gray line. This color coding indicates heterozygous SNPs and their placement on the top (cyan) or bottom (blue) area corresponds to which chromosome the represented data originates. Homozygous positions are not drawn because they add no information to the plot. A "switch error" occurs when a mistake is made as to which chromosome the allele is placed on. Switch errors are marked by red circles.



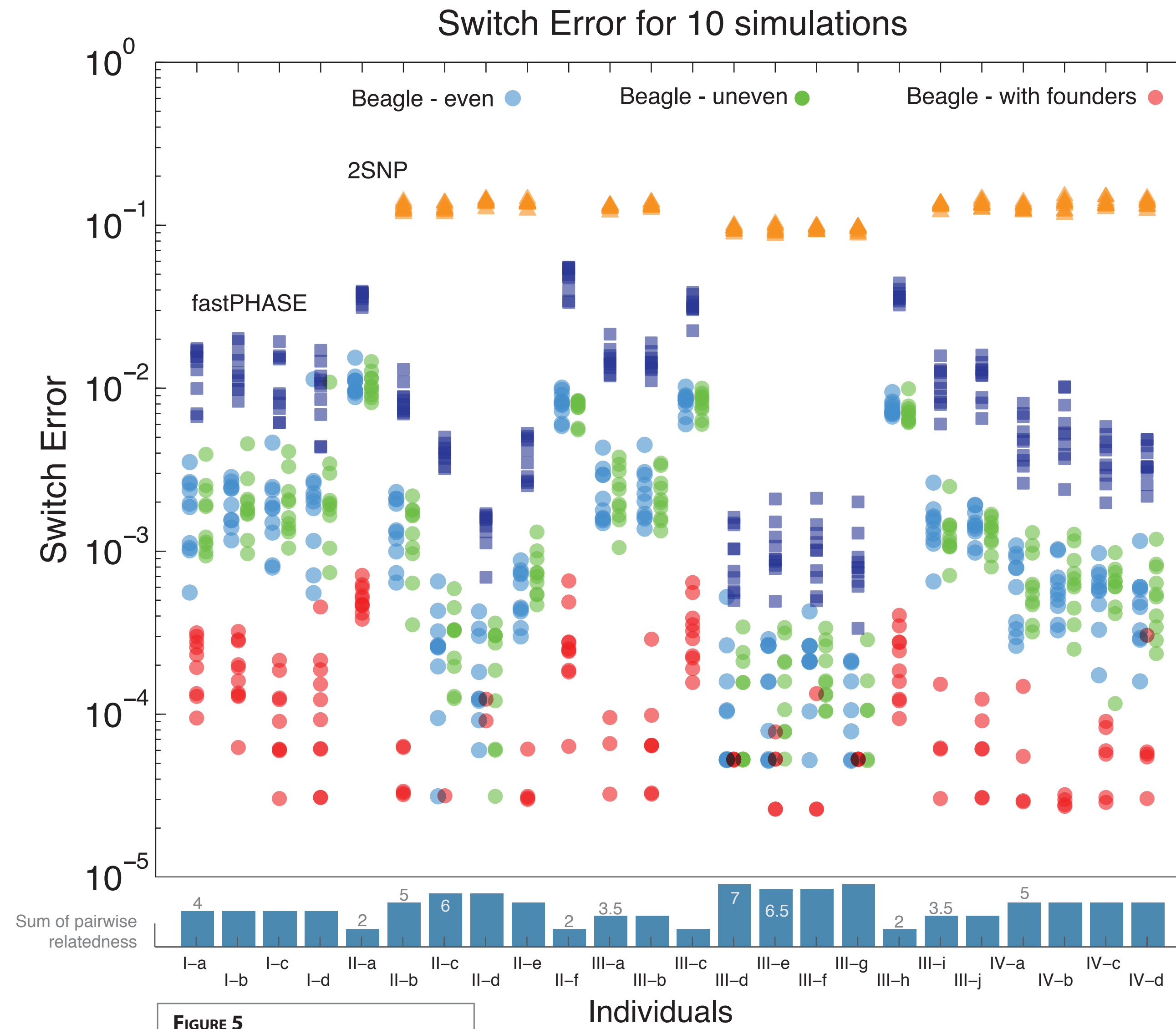
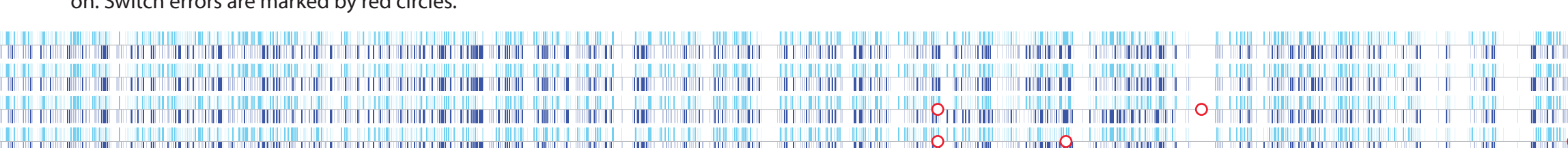
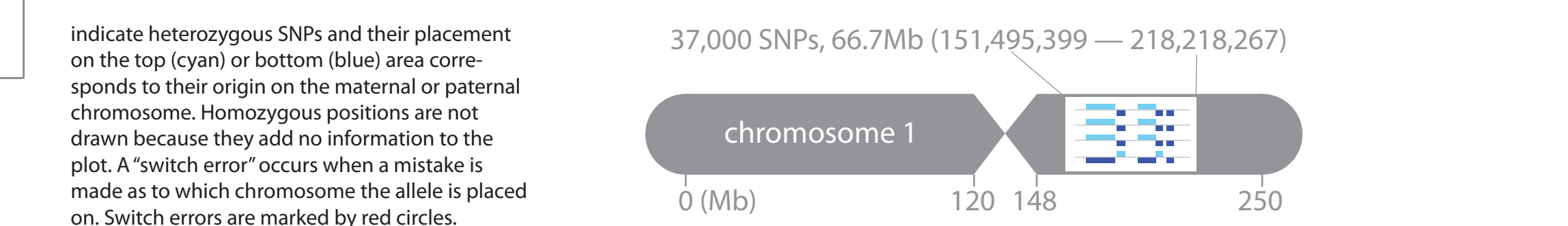
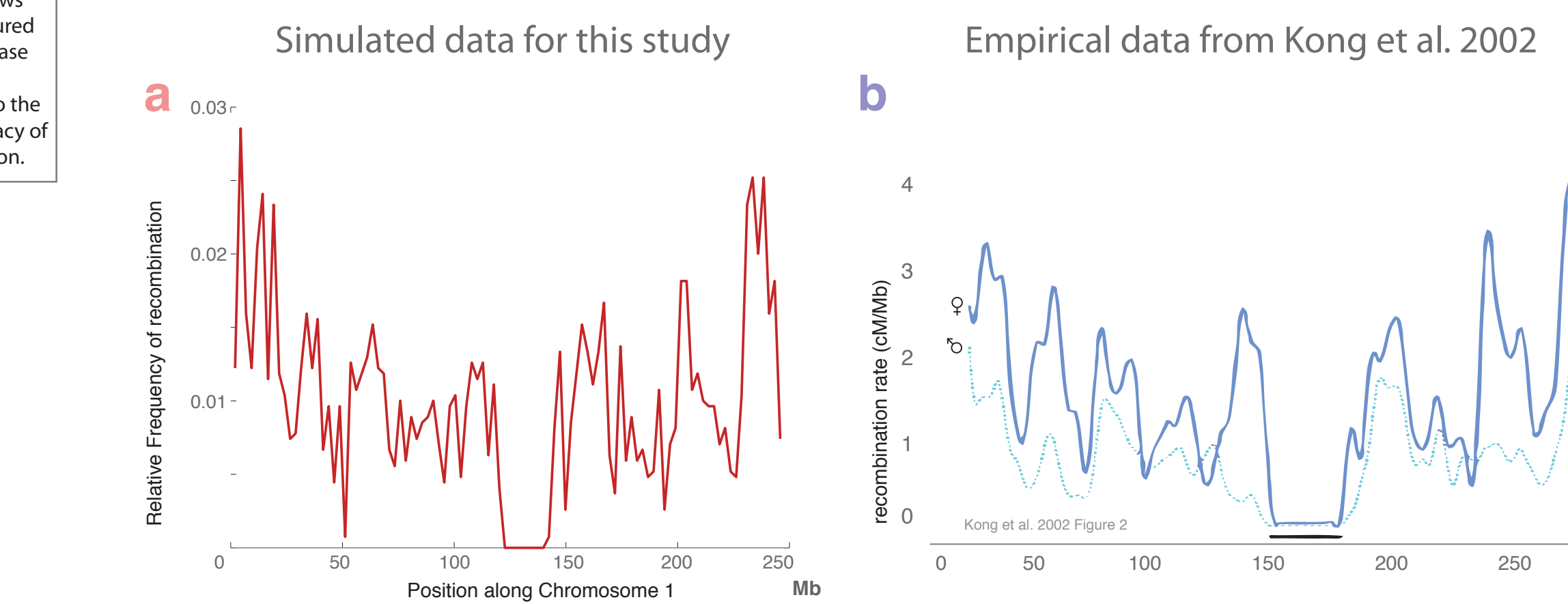
**FIGURE 1**  
Pedigree showing experimental design. Colored clouds represent source populations and colored curves are indicative of founding individuals in the pedigree. Generation number is shown in roman numerals on the left. Individuals are referred to as III-a, or IV-c.



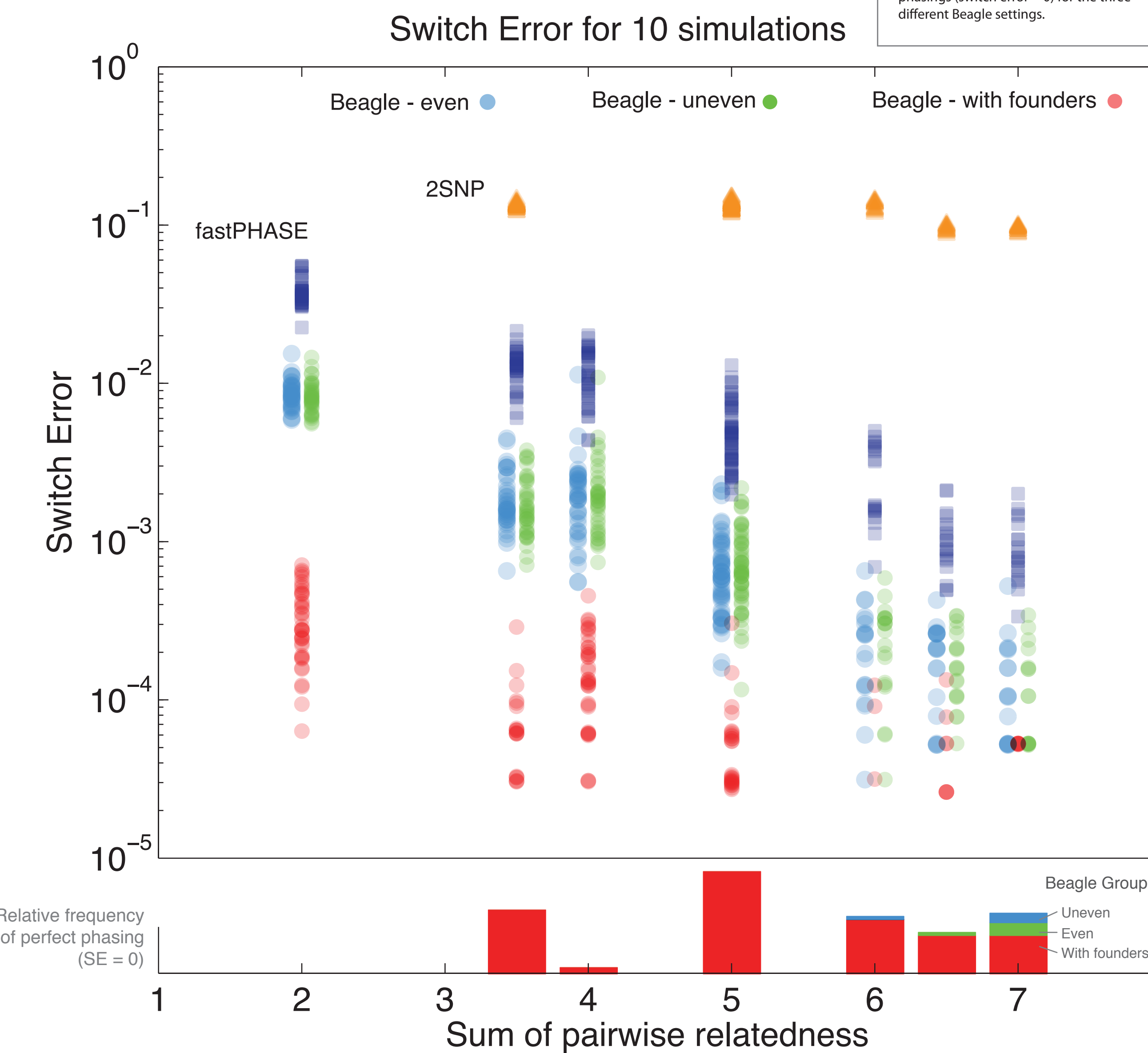
**FIGURE 2**  
Ten simulated pedigrees. Each individual is comprised of a pair of stacked colored sausage shapes (each representing a copy of chromosome 1 with 120,000 SNPs). Different colors represent different haplotype blocks with warm colors (reds, oranges) indicating YRI population haplotypes and cool colors (blues, purples) indicating CEU population haplotypes. Offspring of matings receive one recombinant chromosome from each parent. The first simulation is shown at larger size for detail.



**FIGURE 3**  
Comparison between simulated gametes and empirical data from the literature. Part (a) shows the relative frequency of recombination events along chromosome 1 for 1,000 simulated gametes. Part (b) shows Figure 2 from Kong et al. 2002, empirically measured recombination rates in centimorgans per megabase (cM/Mb) for males and females. That the simulated data appears similar to the data from the literature is indicative of the accuracy of the recombination rate data used in the simulation.



**FIGURE 5**  
Scatter plot of log of switch error for the three different software packages tested, shown as a function of individual. Individuals are coded as in the pedigree in Figure 1. Beneath the scatter plot is a histogram of the sum of pairwise relatedness of each individual to the data set. Beagle is shown with three different settings.



**FIGURE 6**  
Scatter plot of log of switch error for the three different software packages tested, shown as a function of the individual's sum of pairwise relatedness to the data set. Beneath the scatter plot is a histogram showing the relative frequency of perfect phasings (switch error = 0) for the three different Beagle settings.

## Results

**TESTS**  
All algorithms were run on the same server powered by dual, quad-core processors (eight cores total, Intel Xeon 2.83GHz) with 16 gigs of RAM.

**fastPHASE**  
fastPHASE showed consistently good performance, with switch errors ranging between 4.927e-04 and 4.76e-02. While the authors claim the software can make use of a panel of known haplotypes to improve phasing I was unable to get the feature to work without crashing due to memory allocation issues (segmentation faults). Running without panel information the algorithm took approximately seven and a half hours for each replicate to run and performed very well, easily besting the switch error scores of the 200 times faster 2SNP.

I note here that the large number of markers in the present study (approximately 120,000) does not seem to have adversely affected fastPHASE. The small individual sample size, 24, could also work against the algorithm, as it uses the presented data set in part as a template for finding the phasing solutions. However, fastPHASE performed well even with the large number of markers and relatively small number of individuals.

**BEAGLE**  
Beagle had the best performance of the algorithms, solving a few haplotypes with no errors (perfect phasing) and generally running switch errors below 1 in 10,000 to 1 in 1,000 range. Beagle does not take in any information about the pedigree and the success it enjoyed in this data set is likely due to the prevalence of similar haplotypes to the founder haplotypes being in the panel. Beagle took approximately an hour and a half to run.

**BEAGLE SUBSETS**  
In order to determine how much of a role Beagle's library played in its performance, I tested three library types: uneven, where the library was made up of 168 CEU haplotypes and 192 YRI haplotypes; even, where the library was made up of 168 CEU haplotypes and 168 (randomly chosen) YRI haplotypes; and with founders, where the library was made up of 176 CEU and 200 YRI haplotypes, including the founding haplotypes (the correct solutions).

**2SNP**  
2SNP, arguably the simplest package evaluated in this study was also the poorest performer with switch errors consistently in the range of 1 in 10,000 to 1 in 1,000 range. 2SNP is limited in that it does not attempt to predict the true phase of the parents in a trio but will simply predict the transmitted and un-transmitted chromosomes to the child. This is why there are missing data in the 2SNP column in Table switchError. 2SNP was remarkably efficient in terms of time taken, completing work after an average of 130 seconds.

## Conclusions

We presented here the results of 10 independent simulations of a pedigree of 24 individuals over four generations, tracking 120,000 SNPs over chromosome 1.

The most exciting results from this study are the instances of near perfect haplotype reconstruction carried out by Beagle, and the apparent increase in phasing accuracy in all software when analyzing individuals more related to the data set than individuals less related to the data set. It was unexpected that any of the software tested would be able to solve the phasing correctly, but Beagle returned at least one perfect phasing under all subsets of its library that were tested.

The increased accuracy relative to relatedness effect appears to have some impact on 2SNP, which is unexpected because 2SNP does not consider the entire data set at once but instead runs on trios (mother father child). A more likely explanation for the slight increase in accuracy seen in 2SNP with the increase of relatedness is that the individuals with the greatest relatedness values are the ones in the center of the pedigree with parents from two different ancestral populations. Without further testing it seems compelling that the slight increase in accuracy for 2SNP is due to the extreme difference between the mother and father haplotypes for these individuals. To test for this a second experiment would need to be run with an altered pedigree in which an individual was mated into the family from the opposite ancestral population (e.g. a new coupling between a CEU individual and III-b), or a YRI individual and III-b).

The strong downward trend of Beagle and fastPHASE with increasing relatedness seems to be best explained by the fact that both models rely on HMMs and both use the entire data set to aid in phasing individuals. Thus, more closely related individuals would have access to better templates for finding the best solution. The increased accuracy in Beagle relative to fastPHASE seems best explained by its incorporation of the haplotype library, as evidenced by the order of magnitude increase in accuracy when the original founder haplotypes are included in the haplotype library. We note that difference in switch error between the uneven and even sized library is negligible.

Another unexpected result was the sheer speed of 2SNP. The program phased the children of the trio subsets in only 134 seconds. This raw speed seems to be achieved by parsing the trios by simple Mendelian determinism (which could be implemented through boolean logic) and then to applying a maximum spanning tree to whatever is left over. This second part of the algorithm is likely the source of 2SNP's errors, and is would be an area worth exploring competing methodologies. 2SNP had a switch error ranging from 0.086 to 0.14, meaning it made mistakes on up to 14% of the data. Despite the raw speed of the algorithm, this error rate is quite high.

fastPHASE performed quite well. The program had switch errors ranging from 4.9e-04 to 4.8e-02, with the lowest number of errors occurring in the admixed offspring and the largest errors occurring in the founding individuals. The program took seven and a half hours to run on average, which was the longest amount of time.

[1] Gao, J., B. Albrecht, S. S. Chern, W. Li, D. C. C. Cookson, and L. R. Cardon. Multiscale analysis of dense genetic maps using sparse grid trees. *Nat Genet*, 38(11):1011-1016, Nov 2006.  
[2] Matthew Stephens and Peter Donnelly. A comparison of Bayesian methods for haplotype reconstruction from population genotype data. *Am J Hum Genet*, 78(3):1432-9, Nov 2005.  
[3] G. Goren and H. C. Long. High density linkage disequilibrium mapping using models of haplotype block variation. *Hum Genetics*, 20 Suppl 1:162-167, Aug 2004.  
[4] G. Goren and H. C. Long. Genetic diversity and linkage disequilibrium in the human genome. *Proc Natl Acad Sci USA*, 103(11):158-162, Jun 2006.  
[5] Paul Scheet and Matthew Stephens. A fast and flexible statistical model for large-scale population genotyping: applications to inferring missing genotypes and haplotype phase. *Am J Hum Genet*, 78(4):623-41, Apr 2006.  
[6] Tom Iwatake, Raul V. Mouton, Peter J. Conroy, Soham Ghosh, Sudeep Shrivastava, Anshul Arora, and Michael Knapp. Identification of probable genotyping errors by comparison of haplotypes. *PLoS ONE*, 1(4):e104, Apr 2006.  
[7] H. C. Long and R. L. Lu. Coherent trends, extending single- and double-dose effects of the human genome and its specific variation in recombination. *Am J Hum Genet*, 68(3):1013-1021, Sep 1998.  
[8] Iwan Franke, Florian Grotz, and Hans-Joachim Teichmann. Efficient and accurate large-scale reconstruction of haplotypes. *BMC Bioinformatics*, 7:242, Jun 2006.  
[9] Dhanraj Datta and Alexander Zaslavsky. Scalable phasing method for trios and quartets. *ACM SIGSAC*, 49(3):313-34, Jun 2008.  
[10] H. C. Long, Philip Hecatt, and Paul Forthofer. A multiresolution method for detecting and correcting recombination from gene sequences. *Genetics*, 160(3):1231-41, May 2002.  
[11] Sherry R. Brahm and Brent E. Bouvier. Rapid and accurate haplotype phasing and missing data inference by whole-genome association studies by use of localized haplotype reconstruction. *J Hum Genet*, 48(11):1097-1103, Nov 2007.  
[12] N. Hartz, P. Donnelly, and J. Maccioni. Flexible and accurate genotype imputation method for next-generation genome-wide association studies. *PLoS ONE*, 5(12):e12281, Dec 2010.  
[13] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[14] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[15] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[16] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[17] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[18] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[19] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[20] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[21] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[22] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[23] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[24] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.  
[25] J. K. Pritchard, P. Donnelly, and M. J. Rosenberg. Population structure and the analysis of genetic variation. *Journal of Genetic Epidemiology*, 37(6):133-164, Dec 2000.