

Whole genome identity-by-descent determination

Hadi Sabaa, Zhipeng Cai*, Yining Wang, Randy Goebel

*Department of Computing Science, University of Alberta
Edmonton, Alberta T6G 2E8, Canada*

Stephen Moore†

*Department of Agriculture, Food, and Nutritional Science, University of Alberta
Edmonton, Alberta T6G 2P5, Canada*

Guohui Lin‡

*Department of Computing Science, University of Alberta
Edmonton, Alberta T6G 2E8, Canada*

Emails: {sabaa, zcai, yining, rgoebel, stephen.moore, guohui}@ualberta.ca

High-throughput single nucleotide polymorphism genotyping assays conveniently produce genotype data for genome-wide genetic linkage and association studies. For pedigree datasets, the unphased genotype data is used to infer the haplotypes for individuals, according to Mendelian inheritance rules. Linkage studies can then locate putative chromosomal regions based on the haplotype allele sharing among the pedigree members and their disease status. Most existing haplotyping programs require rather strict pedigree structures and return a single inferred solution for downstream analysis. In this research, we relax the pedigree structure to contain ungenotyped founders, and present an cubic time whole genome haplotyping algorithm to minimize the number of zero-recombination haplotype blocks. With or without explicitly enumerating all the haplotyping solutions, the algorithm determines all distinct haplotype allele identity-by-descent (IBD) sharings among the pedigree members, in linear time in the total number of haplotyping solutions. Our algorithm is implemented as a computer program *iBDD*. Extensive simulation experiments using two sets of sixteen pedigree structures from previous studies showed that, in general, there are trillions of haplotyping solutions, but only up to a few thousand distinct haplotype allele IBD sharings. *iBDD* is able to return all these sharings for downstream genome-wide linkage and association studies.

*Currently with Department of Computer Science, Georgia State University. Atlanta, GA 30303, USA.

†Currently with Centre for Animal Science, Queensland Alliance for Agriculture & Food Innovation, University of Queensland. Brisbane St Lucia, QLD 4072, Australia.

‡Corresponding Author. Phone: +1 780 4923737; Fax: +1 780 4921071

Keywords: Haplotyping, identical by descent, haplotype allele sharing, linkage analysis

Biographical Notes

Dr. Hadi Sabaa received his PhD degree from the Department of Computing Science at the University Alberta. His research interests lie in bioinformatics and genomics, typically in SNP haplotyping and GWAS.

Dr. Zhipeng Cai received his BS degree in computer science from Beijing Institute of Technology, China, and his MS and PhD degrees in Computing Science from University of Alberta, Canada. He is currently an Assistant Professor in the Department of Computer Science at the Georgia State University, USA. His research interests include bioinformatics, wireless networks and optimization theory. Dr. Cai received a number of awards and honors, including the NSERC Postdoc Fellowship and the Ph.D Outstanding Research Achievement Award. He also served as a guest editor for *Algorithmica* and *Journal of Combinatorial Optimization*, and as an Editorial Board Member for *International Journal of Sensor Networks*.

Yining Wang is a PhD student in the Department of Computing Science at the University Alberta. His research interests lie in bioinformatics and applications of machine learning to genomics.

Professor Randy Goebel is vice president of the innovates Centre of Research Excellence (iCORE) at Alberta Innovates Technology Futures (www.albertainnovates.ca), chair of the Alberta Innovates Academy, professor of Computing Science at the University of Alberta (www.cs.ualberta.ca), and principle investigator in the Alberta Innovates Centre for Machine Learning (www.aicml.ca). Randy has research interests that include applications of machine learning to systems biology, visualization, and web mining, as well as work on natural language processing, web semantics, and belief revision.

Professor Stephen Moore is a highly regarded research scientist who prior to his appointment with QAAFI worked for the cattle industry in Alberta, Canada. He has more than 20 years' experience in bovine genomics, including his role as Chair in Bovine Genomics at the University of Alberta since 1999. In his role as CEO Livestock Gentec, University of Alberta, Professor Moore led many successful projects to identify genes that underlie production and quality traits in cattle. The expertise he brings to QAAFI's Centre for Animal Science reflects the centre's capacity to embark on research to help Australian and international livestock industries flourish now and into the future. Professor Moore took up his appointment as Director of the Centre For Animal Science at QAAFI in September 2011.

Dr. Guohui Lin is an Associate Professor of Computing Science at the University of Alberta. He received his PhD in Theoretical Computer Science from the Chinese Academy of Sciences in 1998. His research interests include Bioinformatics and Computational Biology, and the recent work focuses on whole genome phylogenetic

analysis for viruses, cancer bioinformatics, and bovine genomics. He is a member of ACM and a member of IEEE Computer Society.

1. Background

A *single nucleotide polymorphism* (SNP) is a DNA sequence variation occurring when a single nucleotide in the genome differs between individuals, or between paired homologous chromosomes in an individual of a diploid species (*e.g.*, humans). SNPs have been used as genetic markers in linkage analysis and association studies, where the sharing status of the alleles among members is used to draw inferences about the inheritable properties. In this research, we have developed a program called *iBDD* which determines all distinct haplotype allele identity-by-descent (IBD) sharings *in one whole genome scan*, for the most complex pedigree genotype datasets.

SNPs are believed to contribute to the most genetic variations in human populations [1]. The rapid development of genotyping technology has led to the identification of thousands to millions of SNPs for various species; for humans, these *common variants* provide the foundation for *genome-wide association studies* (GWAS) under the *common disease – common variant* (CDCV) hypothesis. Recently, GWAS have achieved a great deal of success [2], but genetic fine-mapping for complex diseases such as cancer and mental illness is still a great challenge.

Nevertheless, the unphased genotype data is recognized as a fundamental bottleneck in general genetic linkage and association studies, particularly for rare diseases. For diploid species including humans, at each biallelic SNP locus, the unphased genotype data contains two alleles (nucleotides), without specifying their parental origins. A haplotypes is *phased* genotype which, at each SNP locus, contains the allele of the same parental origin.

In the mapping of disease-susceptible genes in genetic linkage and association studies, one important assumption is that such disease-susceptible genes are in linkage disequilibrium to certain SNPs, so that these SNP markers can be the anchors of disease-susceptible genes. Given its biallelic nature, genetic linkage and association studies based on SNP genotype in general requires a large number of samples, positive samples in particular, so that the association study results are statistically significant [2]. On the other hand, if haplotypes for a large non-recombinant chromosomal region can be determined for all the samples, then the disease-susceptible haplotype alleles may be easily identified, since the haplotype allele sharing can be better determined. This approach has been particularly successful for simple Mendelian disease genetic linkage on pedigree data. For general population-based GWAS, haplotypes clearly contain more inheritance information than unphased genotype, and the ideal case is to build the dense SNP haplotype map: this provides more detailed and deterministic haplotype allelic information for the mapping of disease-susceptible genes. It should be emphasized that, as SNPs act as anchors, it is the allele sharing that is used for linkage inference, and thousands to millions of

haplotype allele configurations give the identical sharing. Our *i*BDD program from this research is designed to determine all distinct haplotype allele sharings.

Haplotypes can be experimentally determined, but it is very expensive to do so [3]. In practice, a less costly alternative is to collect genotype data. Therefore, efficient and accurate computational methods for the inference of haplotypes from genotype data are of considerable value. There is a rich and growing literature on haplotype inference from genotype data, also commonly referred to as *phasing* or *haplotyping*. Research that focuses on unrelated individuals — population data — is reviewed in [4] and [5], with its recent representative fastPHASE by Scheet and Stephens [6]; Research on related individuals — pedigree data — is reviewed in [7], including (exact and approximate) likelihood-based methods [8–16] and genetic rule-based strategies [17–22], with its representative PedPhase by Li and Jiang [17]. The likelihood-based methods usually work for low-density SNP data but not high-density data; neither can they handle large (in many cases, even moderately large) datasets because of the extensive computations required. Additional information and assumptions, such as Hardy-Weinberg equilibrium and marker recombination rates, are generally required to calculate the likelihoods.

Rule-based methods for haplotyping exploit the Mendelian laws of inheritance to minimize the total number of *crossover events* (also called *recombination events* or *breakpoints*) in all pedigree members [7] needed for explaining the observed genotype data. They generally run faster than likelihood-based methods. Nevertheless, this computational minimization problem is NP-hard [17] in general, indicating that there is unlikely a fast algorithm that reconstructs such optimal haplotyping solutions (also called *configurations*). When no recombination events are allowed and the pedigree structure is full (*i.e.* every non-founder pedigree member has both parents genotyped), the problem is then to infer zero-recombination haplotypes from given pedigree genotypes, which is called the *zero-recombination haplotype configuration* (ZRHC) problem. ZRHC turns out to be polynomial time solvable. Li and Jiang [17] presented an $O(m^3n^3)$ time algorithm and a computer program PedPhase which solves ZRHC, by reducing the problem to solving a system of linear equations over the cyclic group Z_2 , where m is the number of loci and n is the number of members in the pedigree. Doan and Evans [23] presented another $O(2^{m^2}m^3n^2)$ time algorithm for ZRHC, a consequence of a fixed-parameter tractable algorithm for the general minimization problem.

It should be noted that the PedPhase and related rule-based methods, as well as those based on likelihood (*e.g.*, PhyloPed [16]), require (either specifically or *conventionally*) the full pedigree structure — any non-founder pedigree member must have both parents genotyped. As an exception, *i*Linker is rule-based, does not require a full pedigree, but the pedigree can have only a couple or a single

founder [20]. In summary, most existing rule-based haplotyping methods have rather strict pedigree structure requirements, while a practical pedigree often contains multiple founders yet some of them may have passed away and their genotype data can no longer be collected. For instance, none of the pedigrees used in seven previous case studies [20, 24–29] is full.

In our research, we relax the past pedigree structure requirement in our haplotyping algorithm to allow for ungenotyped founders, as long as the pedigree stays connected and every ungenotyped member appears in exactly one nuclear family. We note that, while we deal with ungenotyped pedigree members, every genotyped member must have complete genotype data. There is existing work that deals with missing genotype data in genotyped members [30]. (When there is a large portion of missing data in a population dataset, there is another line of work that does haplotyping and imputation using the inferred haplotypes [6, 31–33].) Under the zero-recombination assumption, our novel rule-based haplotyping algorithm has thus the largest pedigree coverage; it can produce all haplotyping solutions in $O(m^3 n^3)$ time, where m is the number of loci, n is the number of (genotyped) members in the pedigree; and it is extended into a whole genome haplotyping algorithm to minimize the number of zero-recombination chromosomal regions (*i.e.* haplotype blocks). Moreover, note that genetic linkage and association studies use allele sharing information amongst the individuals to make inferences, but not the detailed haplotype alleles of each individual [20, 34–38]. We therefore take advantage of all the haplotyping solutions (produced implicitly or explicitly) to determine all distinct genome-wide haplotype allele identity-by-descent (IBD) and *identity-by-state* (IBS) [34] sharings among all pedigree members, together with their associated numbers of haplotyping solutions. It is important to point out that the use of sharing properties of haplotypes adequately addresses the issue of haplotype ambiguity (particularly, the founder haplotype ambiguity) in linkage analysis and haplotype-based association studies. All these functions are coded in a Perl program *iBDD*, which is available upon request.

2. Materials and Methods

A *pedigree* describes the parent-offspring relationship among individuals. In our discussion of a pedigree, members are present only if they are genotyped. Those members that have no parents are the *founders* of the pedigree. A *nuclear family* in a pedigree consists of the parent(s) and all the children. In this paper, connected complex pedigrees are considered, in which a non-present parent of a nuclear family is an ungenotyped founder. Our pedigrees are *connected*, such that each nuclear family has at least one parent and an ungenotyped founder appears in exactly one nuclear family. The genotype at a (biallelic) SNP locus is *homozygous* if the two

alleles are the same (*i.e.*, AA or BB), or is *heterozygous* if the two alleles are different (*i.e.*, AB).

Our haplotyping algorithm is based on the Mendelian laws of inheritance, which states that, at each locus of a pair of homologous autosomes of a child, one allele is inherited from her father (the *paternal* allele) and the other from mother (the *maternal* allele). A child does not inherit a complete autosome from each parent, since *crossover* (also called *recombination*) events occur. That is, during the meiosis process, the two homologous autosomes of a parent may be shuffled and four chromatids are generated, each of which is a shuffled copy of the two homologous autosomes of the parent. One of these chromatids is passed on to the child. Between any two consecutive SNP loci along an autosome of the parent, if recombination occurs then there is a *breakpoint site* between these two loci; this creates a *breakpoint* between the two loci on the corresponding chromatid passed on to the child.

2.1. Zero-recombination haplotyping

Given a pedigree and the unphased genotype data of its members, Li and Jiang [17] presented a constraint-based haplotyping algorithm, PedPhase, under the recombination-free assumption. PedPhase defines a binary parental source (PS) variable z_i for each child z at every SNP locus i , to record the parental origin of the two alleles: $z_i = 1$ if and only if SNP i is heterozygous and child z inherits allele A from her mother. Adopting the cyclic group rule “ $1 + 1 = 0$ ”, PedPhase produces a system of linear equations as haplotyping constraints, where subsets of equations are collected trio by trio. For example, let x denote the mother, y denote the father, and z be their child; let i and j denote two consecutive SNP loci at which x is heterozygous (*i.e.*, x is homozygous at loci $i + 1, i + 2, \dots, j - 1$); assume z is heterozygous as well at these two loci (but z is not necessarily heterozygous at any loci of $i + 1, i + 2, \dots, j - 1$). Depending on the genotype of y at these two loci, PedPhase writes down the following set of constraints:

$$\begin{cases} \text{either } x_i = x_j \text{ or } x_i + x_j = 1, \\ \text{either } y_i = y_j \text{ or } y_i + y_j = 1, \\ x_i + x_j = z_i + z_j \text{ and/or } y_i + y_j = z_i + z_j. \end{cases}$$

It has been proven that all the linear equations are satisfied if and only if the extracted haplotyping solution conforms with the input genotype data. Because the equations are written down for trios, PedPhase requires a full pedigree structure, *i.e.*, every non-founder member must have both parents genotyped.

We re-write these haplotyping constraints on a trio (x, y, z) to separate them into two independent subsets of constraints, for the two parent-child pairs (x, z)

and (y, z) respectively. For example, for the pair (x, z) , only one linear constraint is written down, disregarding what the genotype y has at loci i and j :

$$x_i + x_j + z_i + z_j = 0.$$

In this way, we relax the pedigree structure to allow for one ungenotyped parent per nuclear family — this factoring provides the advantage of the constraint re-writing scheme. That is, in the above case, parent y can be ungenotyped. Furthermore, if indeed y is ungenotyped, then besides the equations for all (x, child) pairs, additional equations can be expressed to constrain the number of haplotypes for y to be at most two. It should be noted that in the case where the pedigree structure is full, our re-writing scheme becomes exactly the same as the PedPhase original scheme. Similarly, it can be shown that these constraints in the re-writing scheme are all satisfied if and only if the corresponding haplotyping solution is feasible. For the detailed re-writing scheme that handles dozens of distinct scenarios and its mathematical proof of correctness, interested readers may refer to Li and Jiang [17] and Cheng *et al.* [21]. Moreover, the total number of linear equations in our new system is proven to be no more than $7mn$, resulting in an $O(m^3n^3)$ -time haplotyping algorithm, where m is the total number of SNPs along the chromosome and n is the size of the pedigree.

2.2. Parsimonious whole genome haplotyping

For small chromosomal regions, the zero-recombination assumption can be reasonable and the genetic linkage and association studies based on zero-recombination haplotyping solutions are meaningful. For genotype datasets on whole chromosomes, we need to relax the zero-recombination assumption but take advantage of the zero-recombination haplotyping algorithm to greedily and optimally determine the maximal zero-recombination haplotype blocks, and their associated haplotyping solutions.

The following outlines at a high level on how this is accomplished in *i*BDD. Starting with $i = 2$, the zero-recombination haplotyping algorithm is run to check whether solutions exist for chromosomal region $[1, i]$, which contains loci $1, 2, \dots, i$. If affirmative, i is incremented by 1; and the checking process is repeated until at some point, there is no solution to the linear system written for chromosomal region $[1, i]$. This gives a maximal zero-recombination chromosomal region $[1, i - 1]$, and all haplotyping solutions are recovered for the region (either implicitly by storing those determined PS variables or explicitly by listing all haplotype configurations). Next, chromosomal region $[1, i - 1]$ is chopped, and the haplotyping process moves on to repeatedly determine the other maximal zero-recombination chromosomal regions starting SNP locus i .

Note that in this incremental haplotyping process, haplotyping solutions for chromosomal region $[1, j]$ are fully used in the zero-recombination haplotyping algorithm to compute the haplotyping solutions for chromosomal region $[1, j + 1]$. Our implementation ensures that the total number of linear equations is confined in $O(mn)$, to guarantee the $O(m^3n^3)$ running time for the whole genome scan.

2.3. *The identity-by-descent determination*

Given an explicit haplotyping solution for a zero-recombination haplotype block, or given an implicit haplotyping solution represented as a solution to the linear system (*i.e.*, the PS values), we can use genotype data to trace the inheritance of each haplotype allele of a genotyped founder member within the entire block. Essentially, we identify for each child's two haplotype alleles their parental source. After this is done, all the pedigree descendants who share the same founder haplotype allele form a cluster. The haplotype allele identical-by-descent (IBD) sharing is characterized as a collection of such clusters labeled with the founder identity. Note that this process does not require (and thus allows the user to choose to not create) explicit haplotyping solutions. The process runs in $O(kn)$ time, where k is the total number of maximal zero-recombination haplotype blocks.

3. Results

For our experiments, we use a real dataset from Wirtenberger *et al.* [39] to generate simulation datasets. This real dataset contains 877 SNPs on chromosome 1 for independent individuals genotyped by GeneChip Human Mapping 10K Xba array. To demonstrate the performance of our *i*BDD, six real pedigree structures are employed, which have been used in previous genetic linkage and association studies (Figure 1 in Lin *et al.* [20], Figure 1 in Martin *et al.* [25], Figure 2 in Sinsheimer *et al.* [29], Figure 2 in Lin *et al.* [24], Figure 1 in Hauser *et al.* [27], and Figure 1 in Howell *et al.* [28], respectively). These pedigrees vary wildly in size, the number of nuclear families, the number of founders, and the number of ungenotyped founders. A summary of their structural characteristics is included in Table 1. We followed a trio simulation process which generates the two whole chromosomal haplotypes for a child from her parents' haplotypes according to the $\chi^2(m)$ -model for crossover events with $m = 4$ [40–42]. Using this trio generation process, the pedigree genotype datasets are simulated, 100 datasets for each pedigree.

3.1. *Breakpoint recovery*

One important performance measure of a haplotyping algorithm is whether a true/simulated breakpoint can be correctly recovered. In the pedigree genotype

Table 1. Properties of the six pedigrees we used in the simulation study.

Pedigree No.	1	2	3	4	5	6
#Members	16	19	17	24	10	20
#Generations	3	3	3	5	3	3
#Nuclear families	4	3	4	12	2	4
#Founders	5	4	5	9	3	5
#Ungenotyped	3	2	2	3	1	4

dataset simulation process, a simulated parental breakpoint site could arise in between two consecutive homozygous SNP loci, thus it cannot be precisely recovered by any computational approach. We adopted the criteria mentioned in previous work [16, 20, 42] to do the mapping between the simulated breakpoint sites and the breakpoint sites computed by *i*BDD: a simulated breakpoint site is classified as *recovered* if there is a computed breakpoint site on the same parent such that the SNPs in between these two locations, if any, are all homozygous; otherwise, the breakpoint is classified as *missed*.

The breakpoint recovery *precision* is defined as the ratio of the number of simulated breakpoint sites that are recovered (true positives) over the number of breakpoint sites computed by *i*BDD (true and false positives). The breakpoint recovery *recall* is defined as the ratio of the number of true positives over the total number of simulated breakpoint sites. We calculated the precision and recall on each simulated dataset as the averages over all the haplotyping solutions. The mean precision and recall values over all 100 datasets associated with the pedigree are reported in Table 2. The general conclusion on breakpoint recovery is that *i*BDD generates slightly fewer breakpoints than the simulated truth, achieving average precision of 79.3%, to recover most of the true breakpoints (average recall 79.0%). Fig. 1 plots the mean recall (*y*-axis) versus precision (*x*-axis) values on the 100 datasets simulated for pedigree No. 1. One can see from this plot and Table 2 that the breakpoint recovery is quite stable across the 100 datasets for each pedigree, suggesting its variation depends mostly on the pedigree structure.

3.2. Haplotype allele sharing recovery

We use the *F*-score from information retrieval to measure the accuracy of haplotype allele sharing recovery, compared with the simulated sharing. More specifically, we determine all the zero-recombination chromosomal regions for the simulated haplotype configuration. These simulated zero-recombination chromosomal regions intersect with the *i*BDD zero-recombination chromosomal regions to produce a set

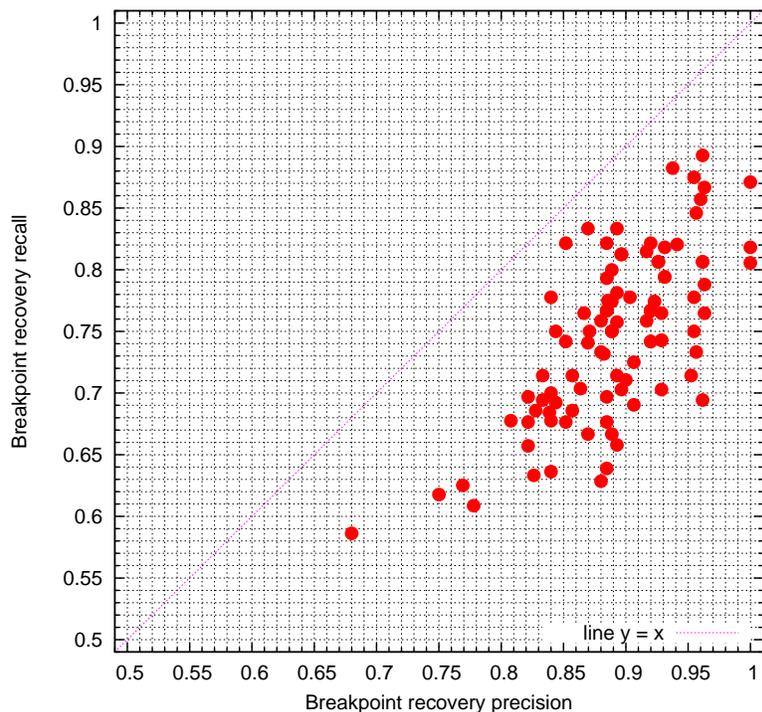


Fig. 1. The mean recall versus precision values on the 100 datasets simulated for pedigree No. 1.

Table 2. The mean breakpoint recovery precision and recall by *i*BDD on each of the six pedigrees, over the averages of 100 simulated datasets, and their standard deviations.

Pedigree No.	Precision	Recall
1	0.893 ± 0.053	0.748 ± 0.066
2	0.684 ± 0.064	0.879 ± 0.061
3	0.862 ± 0.057	0.793 ± 0.060
4	0.915 ± 0.043	0.683 ± 0.053
5	0.729 ± 0.082	0.886 ± 0.065
6	0.678 ± 0.075	0.751 ± 0.074
Average	0.793	0.790

of *common* zero-recombination chromosomal regions to the simulated sharing and the *i*BDD sharing. Note that on each of these zero-recombination chromosomal regions, the simulated sharing is also a collection of pedigree member clusters labeled with the founder identity. We first arbitrarily name the two clusters associated with

the same founder *paternal* and *maternal*, respectively. The two clusters associated with the same founder in the *i*BDD sharing are accordingly named, *paternal* and *maternal* respectively, such that the F -score of the mapping between the simulated and the *i*BDD sharings is maximized (over two possible choices). Upon completing all the founders, the (weighted) F -score for this chromosomal region is computed between the two collections of labeled clusters. The chromosomal F -score between the simulated and the *i*BDD haplotype allele IBD sharings is taken as the weighted average F -score over all maximal zero-recombination chromosomal regions, where the weight of a region is taken as the number of SNPs in that region.

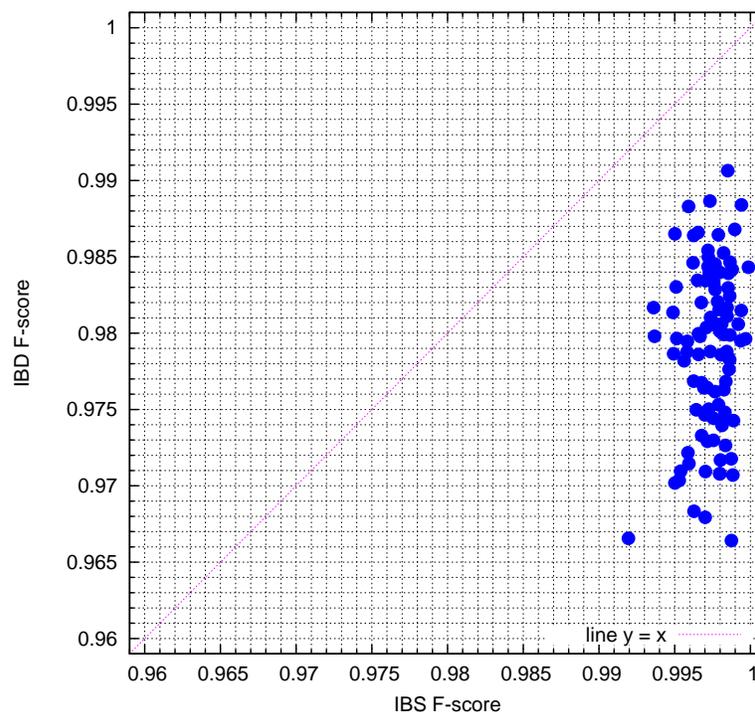


Fig. 2. The mean F -scores of the haplotype allele IBD and IBS sharings compared with the simulated sharings, respectively, on the 100 datasets simulated for pedigree No. 1.

In addition to the haplotype allele IBD sharing, we have also collected all distinct haplotype allele IBS sharings produced by *i*BDD, and calculated the F -score between the simulated haplotype allele IBS sharing and each *i*BDD haplotype allele IBS sharing. That is, on each simulated dataset, we calculated the F -score of each IBD (IBS, respectively) sharing against the simulated IBD (IBS, respectively) sharing to demonstrate the performance of haplotype allele sharing recovery by *i*BDD. The average F -score over all distinct IBD (IBS, respectively) sharings by

*i*BDD against the simulated IBD (IBS, respectively) sharing, across all 100 datasets simulated for each of the six pedigrees, is reported in Table 3. Note that given a haplotyping solution, the associated haplotype allele IBD sharing is a refinement of the associated haplotype allele IBS sharing. The number of distinct haplotype allele IBS sharings produced by *i*BDD is about four times the number of distinct haplotype allele IBD sharings, for all the 600 datasets in our simulation studies. Fig. 2 plots the average recovery F -scores (*i.e.* against simulation) on the 100 datasets simulated for pedigree No. 1. One can see from Fig. 2 and Table 3 that most of the haplotype allele sharings produced by *i*BDD are close to the simulated truth, yet some could be a bit far away. This suggests that using only one sharing in linkage and association studies might not work well, as the sharing could significantly deviate from the truth.

Table 3. The mean haplotype allele sharing F -scores by *i*BDD against the simulated ones on each of the six pedigrees, over the averages of 100 simulated datasets, and their standard deviations.

Pedigree No.	F -score	
	IBD	IBS
1	0.978±0.005	0.996±0.002
2	0.968±0.007	0.998±0.001
3	0.980±0.003	0.996±0.002
4	0.984±0.003	0.999±0.001
5	0.986±0.003	0.998±0.002
6	0.975±0.004	0.996±0.002
Average	0.978	0.997

4. Discussion

Guaranteed by its worst-case cubic running time, *i*BDD runs fast to produce all haplotyping solutions, from which it takes only a fraction of time to determine all the distinct haplotype allele IBD and IBS sharings together with their associated numbers of haplotyping solutions. Among the six pedigrees used in our study, pedigree No. 1 is moderately complex (Table 1). *i*BDD was able to terminate in about two minutes on average, using an Intel E6850 3.0GHz processor with 4GB RAM.

4.1. *The huge number of haplotyping solutions versus a handful distinct haplotype allele sharings*

Previously, the uncertainty in the inferred multiple haplotyping solutions was conjectured to be a potential problem for haplotype-based association studies. By using the haplotype allele sharing status, we believe that haplotype ambiguities in multiple haplotyping solutions can be eliminated.

From the simulation study, we observed that for each simulated dataset, the number of haplotyping solutions by *i*BDD is always large, in the trillions. Indeed, for a pedigree of n non-founders and a total number of H heterozygous loci across all genomes of all individuals, there could be 2^{H-n} haplotyping solutions. Nevertheless, there are only a few to dozens of, or up to a few thousand distinct haplotype allele IBD sharings (*i.e.*, up to 2^n), each associated with thousands to millions of haplotyping solutions. For example, across the 100 simulated datasets for pedigree No. 1 (plotted in Fig. 3) the average number of distinct haplotype allele IBS sharings is 48.24, the average number of distinct haplotype allele IBD sharings is 235.78, and their associated numbers of haplotyping solutions are in the range of 2^{60} 's.

The observation suggests that using only one or a few inferred haplotyping solutions in genetic linkage and association studies [16, 20, 26] is not appropriate. Given that trying all haplotyping solutions is computationally prohibitive, our *i*BDD becomes helpful in feasibly enumerating all distinct haplotype allele IBD (or IBS) sharings.

4.2. *Comparison to existing haplotyping methods*

Our *i*BDD is derived from PedPhase [17], and therefore it is reasonable to make comparison to PedPhase. Unfortunately, since PedPhase requires full pedigree structure (*i.e.*, every non-founder member must have both parents genotyped), it does not run on any of the six pedigrees we used earlier. Moreover, PedPhase is a zero-recombination haplotyping algorithm only, which has to be extended into a general haplotyping algorithm for comparison purposes. We thus adopted the parsimony rule in the same way as in *i*BDD to extend it to xPedPhase, which determines all maximal zero-recombination chromosomal regions together with all haplotyping solutions. We also made a comparison to *i*Linker [20], another rule-based haplotyping algorithm which returns only one haplotyping solution. For fair comparison to *i*Linker, we adopted the scheme to use the first haplotyping solutions by *i*BDD and xPedPhase, respectively, and made comparison among these three using breakpoint recovery and haplotype allele IBD sharing recovery.

Note that *i*Linker accepts only pedigrees with a couple founders or a single founder, while xPedPhase accepts only full pedigrees. We therefore used another

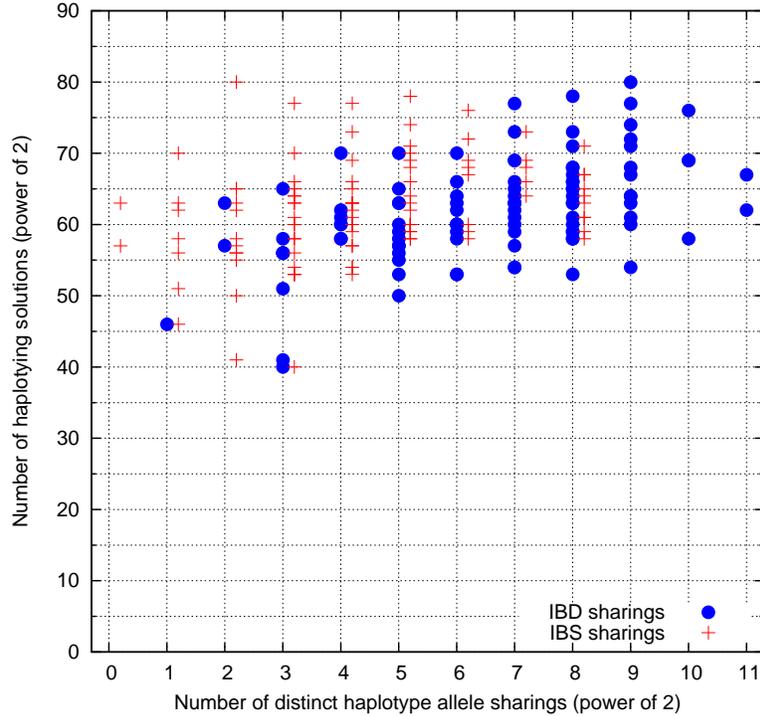


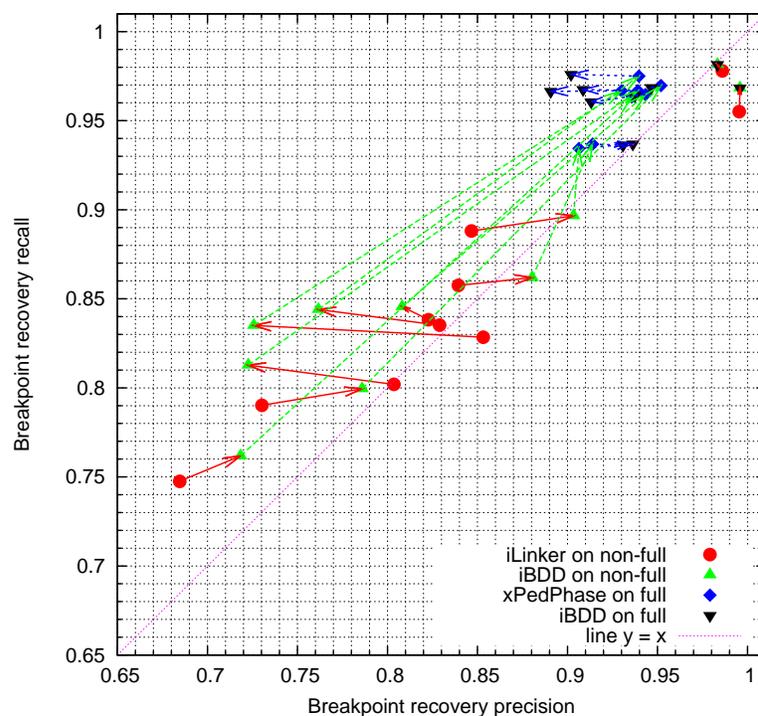
Fig. 3. The numbers of haplotyping solutions, distinct haplotype allele IBD and IBS sharings on the 100 datasets simulated for pedigree No. 1. It is interesting to note that 2 datasets have a unique IBS sharing — 2 crosses in column 0.

set of ten pedigrees in this comparative study, summarized in Table 4. This, however, makes the direct comparison between *i*Linker and xPedPhase impossible; so the only comparison is via *i*BDD. That is, these ten pedigrees all have a couple founders, and thus we can run both *i*Linker and *i*BDD. For each of them, missing founders are added to make the pedigree full, and subsequently simulation datasets are generated for running xPedPhase and *i*BDD. Again, 100 full-datasets are simulated for each full pedigree, from which the genotype data for the missing founders are dropped to form the nonfull-datasets for the corresponding non-full pedigree. The average breakpoint precision and recall, and haplotype allele IBD sharing F -score are collected, for xPedPhase on the full-datasets only, for *i*BDD on both the full-datasets and the nonfull-datasets, and for *i*Linker on the nonfull-datasets only.

Fig. 4 plots the average recall against average precision on the ten pedigrees. In the plot, there are four dots associated with each pedigree, except for the first two pedigrees on which xPedPhase regularly failed. These four dots correspond to the four pairs of (precision, recall) by *i*Linker on the non-full datasets, by *i*BDD on

Table 4. Properties of pedigrees used to compare the performance of xPedPhase, *i*Linker, and *i*BDD.

Pedigree No.	1	2	3	4	5	6	7	8	9	10
#Members	4	5	7	9	10	13	11	13	15	16
#Generations	2	2	3	3	3	3	3	3	3	3
#Nuclear families	1	1	2	3	3	4	3	3	4	4
#Founders	2	2	3	4	4	5	4	4	5	5
#Ungenotyped	0	0	1	2	2	3	2	2	3	3

Fig. 4. The mean recall versus precision values of breakpoint recovery for *i*Linker (red dots) and *i*BDD (green triangles) on ten non-full pedigrees, and xPedPhase (blue triangles) and *i*BDD (black diamonds) on the corresponding full pedigrees, respectively; each pedigree is associated with 100 simulated genotype datasets.

the non-full datasets, by xPedPhase on the full datasets, and by *i*BDD on the full datasets. They are connected sequentially using red, blue, and green arrows. From these arrows, one can see the general tendency of improved performance. In terms of breakpoint recovery, *i*BDD performance was not significantly different from xPe-

dPhase on the full datasets — slightly better recall versus slightly worse precision, *i*BDD performed better than *i*Linker on the non-full datasets, and non-surprisingly breakpoint recoveries on the full datasets are better than on the corresponding non-full datasets. Note that on the full datasets, *i*BDD and xPedPhase are expected to perform no differently, since they essentially use the same haplotyping algorithm; yet the slight difference we have seen is perhaps due to different implementation.

Besides rule-based haplotyping methods, we have also tried to make comparisons with likelihood-based haplotyping algorithms including Haploview [43], fastPHASE [6], and SuperLink [15]. Unfortunately, for the most recent versions of all three programs, Haploview and SuperLink did not produce the complete haplotype for all pedigree members, while fastPHASE did not seem to follow the pedigree structure, resulting in haplotyping solutions that break the Mendelian rules of inheritance. Consequently we do not have results on this comparison to present here.

4.3. *Possible reasons for low breakpoint recovery*

The Results section shows that the *i*BDD haplotype allele sharing recovery is almost perfect, but the breakpoint recovery is not. Fig. 4 shows that for each pair of a full and the corresponding non-full pedigrees, *i*BDD performed much better on the datasets associated with the full pedigree; the difference on average is 0.1324 (or 16.0%) in precision and 0.1274 (or 14.8%) in recall. This suggests that the non-full pedigree structure seems a major reason for low breakpoint recovery. Theoretically, in the face of non-full pedigree structure where some founder members are not genotyped, the constraints derived from them are only a small subset (less than 50%) of the constraints derived from that member if the member is genotyped. Though there are no parental source (PS) variables defined for ungenotyped founders, this could lead to a much larger solution space compared to the solution space associated with the corresponding full pedigree; furthermore, many of these additional solutions may be farther away from the true haplotyping solution.

4.4. *High haplotype allele sharing recovery*

While pedigree structure has a major impact on breakpoint recovery, fortunately, as we showed in the Results section, it has relatively minor effects to the haplotype allele sharing. Our comparative study shows that the IBD/IBS sharing recovery difference is only 0.0069/0.0020, or 0.7%/0.2%. One possible reason is that the haplotype allele sharing is the same for many distinct haplotyping solutions, and consequently more robust against the uncertainties arising from phase inference. This phenomenon also suggests that the use of this more robust haplotype allele sharing directly in genetic linkage and association studies may resolve the issue of

phase ambiguities, which has been previously observed [8, 10, 16, 35, 44, 45].

4.5. Use of haplotype allele sharing in genetic linkage

Our simulation study shows that for each simulated dataset, trillions of haplotyping solutions give rise to only a few close, though distinct, haplotype allele IBD (and IBS) sharings, most of which are very close to the simulation sharing.

We have access to the real genotype data associated with pedigree No. 1, which was used for the linkage analysis of a family disease [20] (data not to be released). We ran *i*BDD on this dataset for all haplotyping solutions and for all distinct haplotype allele IBD sharings. Since member M is the diseased founder, we are interested in the two clusters of pedigree members who share a haplotype allele with M. There are 4096 distinct haplotype allele IBD sharings found (for 2^{58} haplotyping solutions). These sharings differ slightly towards the tail of the chromosome, and they all point to two chromosomal regions — 19.62–20.87Mbps (7 SNPs) and 22.48–29.57Mbps (68 SNPs) — that are shared exclusively by *all 9 diseased members* including the diseased founder M. Surprisingly, the LOD score approach [13] does not detect significant linkage. In this sense, the beauty of our *i*BDD program — its ability to determine the explicit sharing of each haplotype allele — is confirmed, while the traditional LOD score approaches for linkage analysis is not really effective, in particular when the pedigree is small.

5. Conclusions

We presented an efficient $O(m^3n^3)$ -time whole genome haplotyping algorithm for a pedigree genotype dataset to minimize the number of zero-recombination haplotype blocks, where n is the size of the pedigree and m is the number of SNPs, and the pedigree can contain ungenotyped founders. With or without explicitly enumerating all the haplotyping solutions, our *i*BDD program determines all distinct haplotype allele IBD sharings among the pedigree members. Extensive simulation experiments supported that *i*BDD is able to return all these sharings for downstream genome-wide linkage and association studies. It would be interesting to develop a similar score as the LOD score in Merlin, such that the score incorporates all distinct haplotype allele sharings and their associated number of haplotype configurations. The new score landscape will be more useful in linkage analysis.

6. Acknowledgments

This research is supported in part by AARI, AICML, ALIDF, iCORE, and NSERC.

References

1. International Human Genome Sequencing Consortium. Initial sequencing and analysis of the human genome. *Nature*, 409:860–921, 2001.
2. D. Altshuler, M. J. Daly, and E. S. Lander. Genetic mapping in human disease. *Science*, 322:881–888, 2008.
3. S. Lin, A. Chakravarti, and D. J. Cutler. Haplotype and missing data inference in nuclear families. *Genome Research*, 14:1624–1632, 2004.
4. T. Niu. Algorithms for inferring haplotypes. *Genetic Epidemiology*, 27:334–347, 2004.
5. R. M. Salem, J. Wessel, and N. J. Schork. A comprehensive literature review of haplotyping software and methods for use with unrelated individuals. *Human Genomics*, 2:39–66, 2005.
6. P. Scheet and M. Stephens. A fast and flexible statistical model for large-scale population genotype data: applications to inferring missing genotypes and haplotypic phase. *American Journal of Human Genetics*, 78:629–644, 2006.
7. G. Gao, D. B. Allison, and I. Hoeschele. Haplotyping methods for pedigrees. *Human Heredity*, 67:248–266, 2009.
8. E. S. Lander and P. Green. Construction of multilocus genetic linkage maps in human. *Proceedings of National Academy of Sciences of USA*, 84:2363–2367, 1987.
9. E. Sobel, K. Lange, J. R. O’Connell, and D. E. Weeks. Haplotype algorithms. In T. P. Speed and M. S. Waterman, editors, *Genetic Mapping and DNA Sequencing. IMA Volumes in Mathematics and Its Applications*, pages 89–110. Springer, New York, 1995.
10. L. Kruglyak, M. J. Daly, M. P. Reeve-Daly, and E. S. Lander. Parametric and nonparametric linkage analysis: a unified multipoint approach. *American Journal of Human Genetics*, 58:1347–1363, 1996.
11. S. Lin and T. P. Speed. An algorithm for haplotype analysis. *Journal of Computational Biology*, 4:535–546, 1997.
12. D. F. Gudbjartsson, K. Jonasson, M. L. Frigge, and A. Kong. Allegro, a new computer program for multipoint linkage analysis. *Nature Genetics*, 25:12–13, 2000.
13. G. R. Abecasis, S. S. Cherny, W. O. Cookson, and L. R. Cardon. Merlin – rapid analysis of dense genetic maps using sparse gene flow trees. *Nature Genetics*, 30:97–101, 2002.
14. G. Gao, I. Hoeschele, P. Sorensen, and F. X. Du. Conditional probability methods for haplotyping in pedigrees. *Genetics*, 167:2055–2065, 2004.
15. M. Fishelson, N. Dovgolevsky, and D. Geiger. Maximum likelihood haplotyping for general pedigrees. *Human Heredity*, 59:41–60, 2005.
16. B. Kirkpatrick, E. Halperin, and R. M. Karp. Haplotype inference in complex pedigrees. *Journal of Computational Biology*, 17:269–280, 2010.
17. J. Li and T. Jiang. Efficient rule-based haplotyping algorithms for pedigree data. In *Proceedings of the 7th Annual Conference on Research in Computational Molecular Biology (RECOMB’03)*, pages 197–206, 2003.
18. K. Zhang, F. Sun, and H. Zhao. Haplore: a program for haplotype reconstruction in general pedigrees without recombination. *Bioinformatics*, 21:90–103, 2005.
19. J. Xiao, L. Liu, L. Xia, and T. Jiang. Fast elimination of redundant linear equations and reconstruction of recombination-free Mendelian inheritance on a pedigree. In *Proceedings of the 18th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA’07)*, pages 655–664, 2007.

20. G. Lin, Z. Wang, L. Wang, Y.-L. Lau, and W. Yang. Identification of linked regions using high-density SNP genotype data for linkage analyses. *Bioinformatics*, 24:86–93, 2008.
21. Y. Cheng, H. Sabaa, Z. Cai, R. Goebel, and G. Lin. Efficient haplotype inference algorithms in one whole genome scan for pedigree data with non-genotyped founders. *Acta Mathematicae Applicatae Sinica (English Series)*, 25:477–488, 2009.
22. L. Liu and T. Jiang. A linear-time algorithm for reconstructing zero-recombinant haplotype configuration on pedigrees without mating loops. *Journal of Combinatorial Optimization*, 19:217–240, 2010.
23. D. D. Doan and P. A. Evans. An FPT haplotyping algorithm on pedigrees with a small number of sites. *Algorithms for Molecular Biology*, 6:8, 2011.
24. S. Lin, E. Thompson, and E. Wijsman. Achieving irreducibility of the Markov chain Monte Carlo method applied to pedigree data. *Mathematical Medicine and Biology*, 10:1–17, 1993.
25. E. R. Martin, S. A. Monks, L. L. Warren, and N. L. Kaplan. A test for linkage and association in general pedigrees: the pedigree disequilibrium test. *American Journal of Human Genetics*, 67:146–154, 2000.
26. P. Tapadar, S. Ghosh, and P. M. Majumder. Haplotyping in pedigrees via a genetic algorithm. *Human Heredity*, pages 43–56, 2000.
27. M. A. Hauser, C. B. Conde, V. Kowaljow, G. Zeppa, A. L. Taratuto, U. M. Torian, J. Vance, M. A. Pericak-Vance, M. C. Speer, and A. L. Rosa. Myotilin mutation found in second pedigree with LGMD1A. *The American Journal of Human Genetics*, 71:1428–1432, 2002.
28. N. Howell, C. B. Smejkal, D. A. Mackey, P. F. Chinnery, D. M. Turnbull, and C. Herrstadt. The pedigree rate of sequence divergence in the human mitochondrial genome: there is a difference between phylogenetic and pedigree rates. *The American Journal of Human Genetics*, 72:659–670, 2003.
29. J. S. Sinsheimer, C. L. Plaisier, A. Huertas-Vazquez, C. Aguilar-Salinas, T. Tusie-Luna, P. Pajukanta, and K. Lange. Estimating ethnic admixture from pedigree data. *The American Journal of Human Genetics*, 82:748–755, 2008.
30. X. Li and J. Li. Efficient haplotype inference from pedigrees with missing data using linear systems with disjoint-set data structures. In *Proceedings of the Seventh Computational Systems Bioinformatics Conference*, pages 297–308, 2008.
31. S. R. Browning and B. L. Browning. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *American Journal of Human Genetics*, 81:1084–1097, 2007.
32. J. Marchini, B. Howie, S. Myers, G. McVean, and P. Donnelly. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nature Genetics*, 39:906–913, 2007.
33. Y. Li, C. J. Willer, J. Ding, P. Scheet, and G. R. Abecasis. MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genetic Epidemiology*, 34:816–834, 2010.
34. C. W. Cotterman. *A calculus for statistical genetics*. PhD thesis, Ohio State University, Columbus, 1940.
35. S. W. Guo. Proportion of genome shared identical by descent by relatives: concept, computation, and applications. *American Journal of Human Genetics*, 56:1468–1476, 1995.
36. D. J. Schaid, C. M. Rowland, D. E. Tines, R. M. Jacobson, and G. A. Poland. Score

- tests for association between traits and haplotypes when linkage phase is ambiguous. *American Journal of Human Genetics*, 70:425–434, 2002.
37. J. Y. Tzeng, B. Devlin, L. Wasserman, and K. Roeder. On the identification of disease mutations by the analysis of haplotype similarity and goodness of fit. *American Journal of Human Genetics*, 72:891–902, 2003.
 38. D. Y. Lin and D. Zeng. Likelihood-based inference on haplotype effects in genetic association studies. *American Journal of Human Genetics*, 101:89–104, 2006.
 39. M. Wirtenberger, K. Hemminki, B. Chen, and B. Burwinkel. SNP microarray analysis for genome-wide detection of crossover regions. *Human Genetics*, 117:389–397, 2005.
 40. H. Zhao, T. P. Speed, and M. S. McPeck. Statistical analysis of crossover interference using the chi-square model. *Genetics*, 139:1045–1056, 1995.
 41. K. W. Broman and J. L. Weber. Characterization of human crossover interference. *American Journal of Human Genetics*, 66:1911–1926, 2000.
 42. Z. Cai, H. Sabaa, Y. Wang, R. Goebel, Z. Wang, J. Xu, P. Stothard, and G. Lin. Most parsimonious haplotype allele sharing determination. *BMC Bioinformatics*, 10:115, 2009.
 43. J. C. Barrett, B. Fry, J. Maller, and M. J. Daly. Haploview: analysis and visualization of LD and haplotype maps. *Bioinformatics*, 21:263–265, 2005.
 44. M. E. Hawley and K. K. Kidd. HAPLO: a program using the EM algorithm to estimate the frequencies of multi-site haplotypes. *Journal of Heredity*, 86:409–411, 1995.
 45. D. E. Weeks, E. Sobel, J. R. O’Connell, and K. Lange. Computer programs for multilocus haplotyping of general pedigrees. *American Journal of Human Genetics*, 56:1506–1507, 1995.