

# UC Berkeley

## UC Berkeley Previously Published Works

**Title**

Haplotypes versus Genotypes on Pedigrees

**Permalink**

<https://escholarship.org/uc/item/6n11t0pw>

**Journal**

Algorithms for Molecular Biology, 6(1)

**ISSN**

1748-7188

**Author**

Kirkpatrick, Bonnie B

**Publication Date**

2011-04-19

**DOI**

<http://dx.doi.org/10.1186/1748-7188-6-10>

Peer reviewed

RESEARCH

Open Access

# Haplotypes versus genotypes on pedigrees

Bonnie B Kirkpatrick<sup>1,2</sup>

## Abstract

**Background:** Genome sequencing will soon produce haplotype data for individuals. For pedigrees of related individuals, sequencing appears to be an attractive alternative to genotyping. However, methods for pedigree analysis with haplotype data have not yet been developed, and the computational complexity of such problems has been an open question. Furthermore, it is not clear in which scenarios haplotype data would provide better estimates than genotype data for quantities such as recombination rates.

**Results:** To answer these questions, a reduction is given from genotype problem instances to haplotype problem instances, and it is shown that solving the haplotype problem yields the solution to the genotype problem, up to constant factors or coefficients. The pedigree analysis problems we will consider are the likelihood, maximum probability haplotype, and minimum recombination haplotype problems.

**Conclusions:** Two algorithms are introduced: an exponential-time hidden Markov model (HMM) for haplotype data where some individuals are untyped, and a linear-time algorithm for pedigrees having haplotype data for all individuals. Recombination estimates from the general haplotype HMM algorithm are compared to recombination estimates produced by a genotype HMM. Having haplotype data on all individuals produces better estimates. However, having several untyped individuals can drastically reduce the utility of haplotype data.

Pedigree analysis, both linkage and association studies, has a long history of important contributions to genetics, including disease-gene finding and some of the first genetic maps for humans. Recent contributions include fine-scale recombination maps in humans [1], regions linked to Schizophrenia that might be missed by genome-wide association studies [2], and insights into the relationship between cystic fibrosis and fertility [3]. Algorithms for pedigree problems are of great interest to the computer science community, in part because of connections to machine learning algorithms, optimization methods, and combinatorics [4-8].

Single-molecule sequencing is an attractive alternative to genotyping and would yield haplotypes for individuals in a pedigree [9]. Such technologies are being developed and may become commercial within five to ten years. Sequencing methods would apparently yield more information from the same set of sampled individuals, which is critical due to the limited availability of individuals for sampling in multi-generational pedigrees (i.e. individuals

usually must be living at the time of sampling). There is substantial evidence that haplotypes can be more useful than genotypes for both population and family based studies when using methods such as association studies [10,11] and pedigree analysis [12,13]. While it is intuitive that haplotypes provide more information than genotypes, there are instances with family data in which there are few enough typed individuals that there is little practical difference between haplotype and genotype data. Additionally, in order to exploit the information contained in haplotype data, we need to understand the instances where diploid inheritance is computationally tractable given haplotype data.

Pedigree analysis with genotype data is well studied in terms of complexity [6,7] and algorithms [14-16]. Less is known about haplotype data on pedigrees. This paper shows that, given haplotype data on a pedigree, finding both minimum recombination and maximum probability haplotypes is as tractable as computing the same quantities for pedigrees with genotype data (i.e., these problems are NP- and #P-hard, respectively). To obtain a reduction that applies equally well to several types of pedigree calculations, we will consider a modular polynomial-time mapping from the genotype problem to the

Correspondence: [bbkirk@eecs.berkeley.edu](mailto:bbkirk@eecs.berkeley.edu)

<sup>1</sup>Electrical Engineering and Computer Sciences, University of California Berkeley, Berkeley, CA 94720-1776, USA

Full list of author information is available at the end of the article

haplotype problem. The reduction preserves the solutions to the analyses, meaning that the solution to the haplotype problem is the solution to the genotype problem after adjusting by constant factors or coefficients.

Since the reduction uses a biologically unlikely recombination scenario, we will investigate the accuracy and information of realistic examples with haplotypes and genotype data on the same pedigree. Pedigree data was simulated having a known number of recombinations. The recombination distributions were computed at a particular locus of interest and compared to the ground-truth. Since both the haplotypes and genotypes of a specific person contain the same alleles, the differences between the haplotype and genotype recombination distributions were determined by the extra information in the haplotype data. As expected, the haplotype data reliably yields greater accuracy when all the pedigree individuals are typed. However, as fewer pedigree individuals are typed, there is less practical difference between the utility of haplotype versus genotype data. The number of untyped generations that separate typed individuals influences whether haplotype data are actually more accurate than genotype data. For instance with two half-siblings, having two untyped parents results in estimates from genotype data that are nearly as accurate as the estimates computed from haplotype data.

Finally, there is an important instance where haplotype data is more computationally tractable than genotype data. When all individuals in the pedigree are typed, although unlikely from a practical perspective of obtaining genetic samples, the computational problem decomposes into conditionally independent sub-problems, and has a linear-time algorithm. This can be contrasted with the known hardness of the genotype problem even when all individuals are genotyped. The existence of this linear-time algorithm for haplotype data could facilitate useful approaches that combine population genetic and pedigree methods. For instance, if the haplotypes of the founders are drawn from a coalescent and the pedigree individuals are all haplotyped, the probability of a combined model could easily be computed for certain coalescent models.

### Introduction to Pedigree Analysis

A *pedigree* is a directed acyclic graph where the set of nodes,  $I$ , are individuals, and directed edges indicate genetic inheritance between parent and child. A diploid pedigree (i.e. for humans) necessarily has either zero or two incoming edges for each person. The set,  $F$ , of individuals without incoming edges are referred to as pedigree *founders*. An individual,  $i$ , with two parents is a *non-founder*, and we will refer to their two parents as  $m(i)$  and  $p(i)$ .

As is commonly done to accommodate inheritance of genetic information, we will extend this model to include a representation of the alleles of each individual and of the inheritance origin of each allele. More formally, we represent a single chromosome as an ordered sequence of variables,  $x_j$ , where each variable takes on an *allele* value in  $\{1, \dots, k_j\}$ . Each variable represents a *polymorphic site*,  $j$ , in the genome, where there are  $k_j$  possible sequence variants. Since diploid individuals have two copies of each chromosome, one copy inherited from each parent, we will use a superscript  $m$  and  $p$  to indicate the maternal and paternal chromosomes respectively. For a particular individual  $i$ , the information on both copies of a particular chromosome at site  $j$  is represented as  $x_{i,j}^m$  and  $x_{i,j}^p$ .

Furthermore, we assume that inheritance in the pedigree proceeds with recombination and without mutation (i.e. Mendelian inheritance at each site). This imposes consistency rules on parents and children: the allele  $x_{i,j}^m$  must appear in the mother  $m(i)$ 's genome as either the grand-maternal or grand-paternal allele,  $x_{m(i),j}^m$  or  $x_{m(i),j}^p$  and similarly for the paternal allele and the father  $p(i)$ 's genome.

Let  $x$  be a vector containing all the haplotypes  $x_i^m, x_i^p$  for all individuals  $i \in I$ , then we are interested in the probability

$$\mathbb{P}[x] = \prod_{f \in F} \mathbb{P}\left[x_f^p\right] \mathbb{P}\left[x_f^m\right] \cdot \prod_{i \in I \setminus F} \mathbb{P}\left[x_i^p \mid x_{p(i)}^p, x_{p(i)}^m\right] \mathbb{P}\left[x_i^m \mid x_{m(i)}^p, x_{m(i)}^m\right],$$

where the superscript  $m$  and  $p$  indicate maternal and paternal alleles, while the functions  $m(i)$  and  $p(i)$  indicate parents of  $i$ . The first product is over the independent founder individuals whose haplotypes are drawn from a uniform prior distribution, while the second product, over the non-founders, contains the probabilities for the children to inherit their haplotypes from their parents. The unobserved vector  $x$  is not immediately derived from observed haplotype data, since vector  $x$  contains haplotype alleles labeled with their parental origins for all the individuals. To compute this quantity, we need notation to represent the parental origins of each allele where differing origins for neighboring haplotype alleles will indicate recombination events.

For each non-founder, let us indicate the source of each maternal allele using the binary variable  $s_{i,j}^m \in \{m, p\}$ , where the value  $m$  indicates that  $x_{i,j}^m$  allele has grand-maternal origin and  $p$  indicates grand-paternal origin. Similarly, we define  $s_{i,j}^p$  for the origin of  $i$ 's paternal allele. For a particular site, these indicators for

all the individuals,  $s_j$ , is commonly referred to as the identity-by-descent (IBD) inheritance path. A recombination is observed at consecutive sites as a change in the binary value of a source vector, for instance,  $s_{i,j}^m = p$  and  $s_{i,j+1}^m = m$ . To compute the inheritance portion of the equation for  $P[x]$ , we will sum over the inheritance options  $\mathbb{P}[x] = \sum_s \mathbb{P}[x|s] \mathbb{P}[s]$  where  $\mathbb{P}[s] = 1/2^{2|I \setminus F|}$ . We can observe two kinds of data for pedigree individuals whose genetic material is available. The first, and most common, is genotype data, a tuple of alleles  $(g_{i,j}^0, g_{i,j}^1)$  that must appear in the variables  $x_{i,j}^m$  and  $x_{i,j}^p$  for each site  $j$ . Since these alleles are unlabeled for origin, we do not know which allele was inherited from which parent. The second type of data is haplotypes, where we observe two sequences of alleles  $h_i^0$  and  $h_i^1$  and each sequence represents alleles that were inherited together from the same parent. However, we do not know which sequence is maternal and which is paternal. For either type of data define a function  $C_{i,j}$  for locus  $j$  which indicates compatibility of the assigned haplotype alleles with the data and requires inheritance consistency between generations. Specifically, for genotype data  $C_{i,j} = 1$  if  $x_{i,j}^p = x_{f(i),j}^p$ ,  $x_{i,j}^m = x_{f(i),j}^m$  and  $\{x_{i,j}^m, x_{i,j}^p\} = \{g_{i,j}^0, g_{i,j}^1\}$ . Under haplotype data, the  $C_{i,j} = 1$  when the first two equalities, above, hold and  $\{x_{i,j}^m, x_{i,j}^p\} = \{h_{i,j}^0, h_{i,j}^1\}$ , which are the haplotype alleles at locus  $j$ .

Now, we write the equation for  $P[x]$  as a function of the per-site recombination probability  $\theta \leq 0.5$ . For particular values of all the haplotype alleles  $x_{i,j}^m$  and  $x_{i,j}^p$ , the haplotype probability conditional on the inheritance options and the observed data through  $C_{i,j}$  is

$$\mathbb{P}[x|s] = \prod_{f \in F} \prod_{j=1}^l C_{f,j} \mathbb{P}[x_{f,j}^p] \mathbb{P}[x_{f,j}^m] \prod_{i \in I \setminus F} C_{i,1} \cdot \prod_{j=2}^l C_{i,j} \cdot \theta^{(R_{i,j}^m + R_{i,j}^p)} \cdot (1 - \theta)^{(2 - R_{i,j}^m - R_{i,j}^p)}$$

where  $R_{i,j}^m = \mathbb{1}[s_{i,j-1}^m \neq s_{i,j}^m]$  and  $R_{i,j}^p = \mathbb{1}[s_{i,j-1}^p \neq s_{i,j}^p]$ .

**Pedigree Problem Formulations**

Given a pedigree and some observed genotype or haplotype data, there are three problem formulations that we might be interested in. The first is to compute the probability of some observed data, while the last two problems find values for the unobserved haplotypes of individuals in the pedigree.

**Likelihood**

Find the probability of the observed data by summing over all the possible unobserved haplotypes, i.e.  $\sum_s \sum_s \mathbb{P}[x|s] \mathbb{P}[s]$ .

**Maximum Probability**

Find the values of  $x_{i,j}^m$  and  $x_{i,j}^p$  that maximize the probability of the data, i.e.  $\max_x \sum_s \mathbb{P}[x|s] \mathbb{P}[s]$ .

**Minimum Recombination**

Find the values of  $x_{i,j}^m$  and  $x_{i,j}^p$  that minimize the number of required recombinations, i.e.  $\min_{x,s} \sum_i \sum_{j>2} \mathbb{1}[s_{i,j-1}^p \neq s_{i,j}^p] + \mathbb{1}[s_{i,j-1}^m \neq s_{i,j}^m]$ .

The likelihood is commonly used for estimating site-specific recombination rates, relationship testing, computing p-values for association tests, and performing linkage analysis. Haplotype and/or IBD inferences, obtained by maximizing the probability or minimizing the recombinations, are useful for non-parametric association tests, tests on haplotypes, and tests where there is disease information for unobserved genomes.

**Hardness Results**

With genotype data, the likelihood and minimum recombination problems are NP-hard, while the maximum probability problem is #P-hard. Piccolboni and Gusfield [6] proved the hardness of the likelihood and maximum probability computations by relying on a single locus sub-pedigree with half-siblings. Although their paper discussed a more elaborate setting involving a phenotype, their proof, however, applies to this setting. Li and Jiang proved the minimum recombination problem to be hard by using a two-locus sub-pedigree with half-siblings [7]. In all these proofs, half-siblings were pivotal to establishing reductions from well known NP and #P problems.

In this paper, we introduce a simple and powerful reduction that converts any genotype problem on a pedigree of  $n$  individuals into a haplotype problem on a pedigree of at most  $6n$  individuals. This reduction is simple, because it merely introduces four full-siblings and an extra parent for each genotyped individual. We do not need complicated structures involving inbreeding or half-siblings. The reduction works equally well for all three problem formulations.

**Mapping**

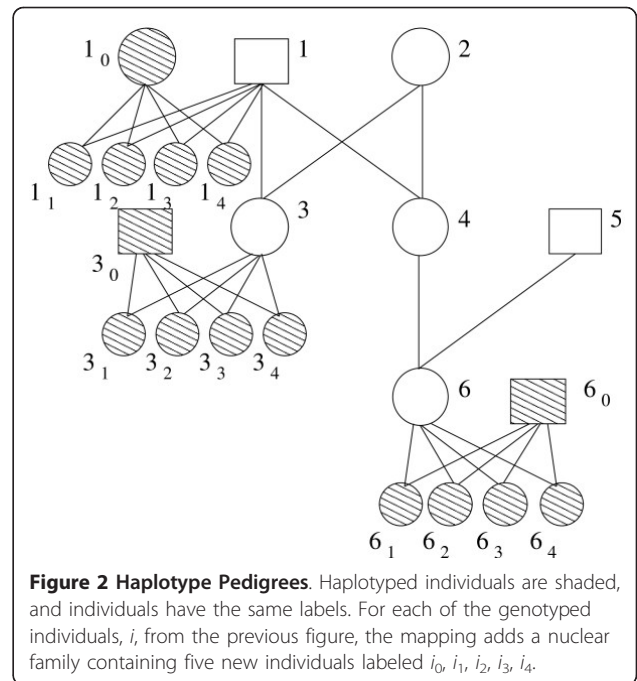
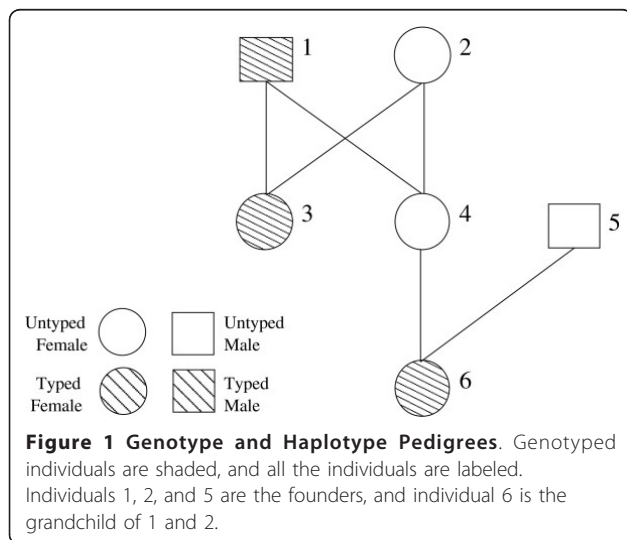
Given a pedigree with genotype data, for any of the three pedigree problems, we define a polynomial mapping to a corresponding haplotype problem with exactly  $5|G|$  individuals haplotyped. First we create the pedigree graph for the new haplotype instance, and later we construct the required haplotype observations from the genotype data.

Let  $G \subset I$  represent the set of genotyped individuals in a pedigree having individuals  $I$  and edges  $E$ . We will create a haplotype instance of the problem, with individuals  $H \cup I$  and edges  $R \cup E$ . To obtain the set  $H$ , we add five

individuals,  $i_0, i_1, i_2, i_3, i_4$ , to  $H$  for every individual  $i \in G$ . The set of new relationship edges,  $R$ , will connect individuals in sets  $H$  and  $G$ . Specifically, the edges stipulate that  $i$  and  $i_0$  are the parents of full-siblings  $i_1, i_2, i_3$ , and  $i_4$  by including the edges:  $i_0 \rightarrow i_1, i_0 \rightarrow i_2, i_0 \rightarrow i_3, i_0 \rightarrow i_4, i \rightarrow i_1, i \rightarrow i_2, i \rightarrow i_3$ , and  $i \rightarrow i_4$ . We will refer to these five individuals,  $i_0, i_1, i_2, i_3$ , and  $i_4$ , and their relationships with  $i$  as the *proxy family* for individual  $i$ . For example, the 6-individual genotype pedigree in Figure 1 becomes a 21-individual genotype pedigree in Figure 2. This produces a pedigree graph with exactly  $5|G| + |I|$  individuals and  $8|G| + |E|$  edges.

To obtain the new haplotype data from the genotype data, we type only individuals in  $H$  such that the corresponding genotyped individual in  $G$  is required, by the rules of inheritance, to have the observed genotypes. Without loss of generality, assume that the genotype alleles are sorted such that  $g_{i,j}^0 < g_{i,j}^1$ . Now we can easily constrain the parental genotype for individual  $i \in G$  by giving the spouse,  $i_0$ , homozygous haplotypes of all ones while giving child  $i_1$  the haplotypes  $\{\bar{1}, g_i^0\}$ , child  $i_2$  haplotypes  $\{\bar{1}, g_i^1\}$ . This guarantees the correct genotype, but does not ensure that the haplotypes of that genotype have the same probability or number of recombinations.

Since there is an arbitrary assorting of genotype alleles at neighboring loci into the parent haplotypes  $x_i^p$  and  $x_i^m$ , we will use the remaining two children to represent possible re-assortments of the genotyped parent's  $T_i$  heterozygous loci, indexed by  $t_j$  where  $1 \leq j \leq T_i$ . In addition to the haplotype  $\bar{1}$ , child  $i_3$ , will have haplotype consisting of  $h_{i_3,t_j} := g_{i,t_j}^{1-j \bmod 2}$ , while child  $i_4$  has the genotyped parent's complementary alleles  $h_{i_4,t_j} := g_{i,t_j}^{j \bmod 2}$ .



This results in child  $i_3$  and  $i_4$  alternating in having the smaller allele at every other heterozygous locus.

This reduction preserves the solutions to the three problems up to constant factors or constant coefficients. Specifically, the solution to the haplotype version of the problem is the solution to the genotype version with the values of the functions being related by constant factors or coefficients, depending on whether the function is a recombination count or a probability.

**Lemma 1.** Let  $r_g$  be the minimum number of recombinations in the genotype problem instance. The mapping yields a haplotype problem instance having

$$r_h = r_g + \sum_{i \in G} 2(T_i - 1) \quad (1)$$

for the minimum number of recombinations, where  $T_i$  is the number of heterozygous sites in genotype  $i$ .

To prove this result, we exploit the alternating pattern of alleles assigned to the four children. This pattern forces there to be two recombinations, among the four children, between consecutive heterozygous loci.

*Proof.* Consider the haplotype instance of the problem. Recall that set  $G$  is defined as the individuals who are genotyped in the genotype problem instance, and, by construction, they are not haplotyped in the haplotype problem instance. For each  $i \in G$  the rules of inheritance applied to  $i$ 's proxy family dictate that the set of alleles at each position are given by  $g_{i,j}^0$  and  $g_{i,j}^1$ . Therefore, the proxy family dictates the genotype of  $i$ .



Since the haplotypes for all the typed individuals are completely given, we only need to consider the assortment of the alleles from  $g_i^0$  and  $g_i^1$  into the maternal and paternal alleles of individual  $i$ . Clearly this assortment determines the number of recombinations that the proxy family contributes to Eqn. (1). However, we will use induction along the genome to show that every possible phasing of the parental genotype induces the same minimum number of recombinations among the four children, namely  $2(T_i - 1)$ .

Now we define an arbitrary assortment of the genotype alleles into two haplotypes for person  $i$ . We can think of this parental genotype for  $l$  loci as a string  $s \in \{H, T\}^l$ , where  $H$  represents a homozygous site and  $T$  a heterozygous site. Recall that  $T_i$  is the number of heterozygous sites in the genotype string, and those sites appear at indices  $t_j$  where  $1 \leq j \leq T_i$ . For this genotype there are  $2^{T_i-1}$  pairs of haplotypes that phase the given genotype. Represent each pair by setting  $T_i - 1$  binary variables

$$P_{t_j} = \begin{cases} 0, & \text{if } x_{i,t_j}^p < x_{i,t_j}^m, \\ 1, & \text{otherwise.} \end{cases}$$

Note, that we are only interested in the origin of the children's haplotypes, rather than in the origin of  $i$ 's haplotypes, so the  $p$  and  $m$  can arbitrarily label either haplotype.

Since  $\{i_1, i_2\}$  between them have the parent genotype at every locus, one of them has origin  $p$  while the other has origin  $m$ , and similarly for  $\{i_3, i_4\}$ . For each locus, indicate the paternal origin of the allele for individuals  $i_1$  and  $i_3$ , respectively with  $Q_j$  and  $S_j$ . Formally,  $Q_j = 1$  if both  $h_{i_1,j} = x_{i_1,j}^p$  and  $h_{i_3,j} = x_{i_3,j}^m$  while  $Q_j = 0$  otherwise. Similarly,  $S_j = 1$  if both  $h_{i_2,j} = x_{i_2,j}^p$  and  $h_{i_4,j} = x_{i_4,j}^m$  while  $S_j = 0$  otherwise. Define  $R_j$  as the minimum recombination count before locus  $j$ . Notice that  $P_{t_1}$  sets the origin of all the child haplotypes, therefore  $R_{t_1} = 0$ , since all preceding homozygous loci can have the same origin as locus  $t_1$ .

From  $t_j$  to  $t_{j+1}$  we have two cases:

1. If  $P_{t_j} = P_{t_{j+1}}$ , then  $Q_{t_j} = Q_{t_{j+1}}$  and  $S_{t_j} \neq S_{t_{j+1}}$  by the alternating construction of children  $i_3$  and  $i_4$  as compared with  $i_1$  and  $i_2$ .
2. Similarly, if  $P_{t_j} \neq P_{t_{j+1}}$ , then  $Q_{t_j} \neq Q_{t_{j+1}}$  and  $S_{t_j} = S_{t_{j+1}}$ .

Furthermore, regardless of the number of homozygous loci separating  $t_j$  and  $t_{j+1}$ , the number of recombinations can only be increased. Therefore, we have the recursion

$$R_{t_{j+1}} = 2 + R_{t_j},$$

proving the lemma.  $\square$

After applying the mapping, the haplotype probability turns out to have a coefficient that is independent of the haplotype assignment to the non-founding parent of the proxy family. This coefficient can be computed in linear time from the genotype data using a Markov chain. The Markov chain has 16 states and has a transition step between each pair of neighboring loci. This small Markov model can be thought of in the sum-product algorithm as an elimination of the typed individuals in the proxy family; alternatively, it is also equivalent to peeling-off the typed proxy individuals in the Elston-Stewart algorithm [14]. Once we have this coefficient, independent of the haplotype assignment, it is clear that the likelihood and maximum probability haplotype problems also have haplotype solutions related proportionally to the genotype solution.

**Lemma 2.** *The mapping yields a haplotype problem instance having haplotype probabilities proportional to the haplotype probabilities of the genotype instance. Specifically, for all  $x$ ,*

$$\mathbb{P}_h[x] = (\mathbb{P}_g[\{x_i | i \in I\}]) \cdot \prod_{i \in G} p_t(i) \prod_j \mathbb{P}[x_{i_0,j}^p = 1] \mathbb{P}[x_{i_0,j}^m = 1]$$

where the proxy family transmission probability is a function of genotype  $g_b$ , the recombination rate  $\theta \leq 0.5$ , and of the transition matrices  $P$ ,  $Q_{0110}$ , and  $Q_{1001}$ ,

$$p_t(i) = \left(\frac{1}{16}\right) \vec{1} \cdot P^{h_0} \cdot \prod_{j=0}^{T_i} (O_j Q_{0110} + (1 - O_j) Q_{1001}) \cdot P^{h_j} \cdot \vec{1}^T$$

and  $O_j$  indicates whether index  $j$  is odd,  $h_0$  is the number of homozygous loci that begin proxy parent's genotype, and  $h_j$  is the number of consecutive homozygous loci after the  $j$ 'th heterozygous locus where there are  $T_i$  heterozygous loci for proxy parent  $i$ . The transition probabilities are given by  $P_{ij} = \theta^{H(i,j)}(1 - \theta)^{4-H(i,j)}$  where  $H(i,j)$  is the Hamming distance between inheritance states  $i$  and  $j$ . Let  $Q_{0110}$  be a transition matrix having non-zero recombination probabilities only in column 0110 (i.e.  $Q_{0110, i,j} = P_{ij}$  when  $j = 0110$ ). Similarly, let  $Q_{1001}$  be a transition matrix with non-zero recombination probabilities only in column 1001.

*Proof.* Without loss of generality, assume that individuals  $i \in G$  are all fathers in their proxy family. This is simply for convenience of notation.

Let  $x$  be any fixed assignment of haplotypes to all the individuals in the pedigree. When conditioning on the assigned haplotypes for individual  $i$ , the probability of the proxy family of  $i$  is independent of the probability

for the rest of the pedigree. Since we can say this for all the proxy families, the terms in the probability for the pedigree individuals in set  $I$  (i.e. those also in the genotype pedigree) are equal to the probability on the genotype data in the genotype pedigree. Therefore, we write that

$$\mathbb{P}_h[x] = \sum_s \mathbb{P}_g[\{x_i | i \in I\} | \{s_i | i \in I\}] \mathbb{P}[\{s_i \in I\}] \cdot \prod_{i \in G} \left( \prod_j \mathbb{P}[x_{i_0,j}^p = 1] \mathbb{P}[x_{i_0,j}^m = 1] \right) \cdot \prod_k \mathbb{P}[x_{i_k}^p | x_{f(i_k)}^p, x_{f(i_k)}^m, s_{i_k}^p] \cdot \mathbb{P}[x_{i_k}^m | x_{m(i_k)}^p, x_{m(i_k)}^m, s_{i_k}^m] \mathbb{P}[s_{i_k}^p] \mathbb{P}[s_{i_k}^m].$$

The sum over vector  $s$  can be split into sums over the component pieces. The sums involving the  $s_{i_k}$  can be distributed into the product over  $k$ , since that is the only place they are used. Let  $s_{i_k} = (s_{i_k}^p, s_{i_k}^m)$ . We easily see that  $\mathbb{P}[x_{i_k}^m | x_{m(i_k)}^p, x_{m(i_k)}^m, s_{i_k}^m] \mathbb{P}[s_{i_k}^m] = 1$ , since there are two ways to inherit the 1-allele from the mother, and all of them are compatible.

$$\mathbb{P}_h[x] = \sum_{\{s_i | i \in I\}} \mathbb{P}_g[\{x_i | i \in I\} | \{s_i | i \in I\}] \mathbb{P}[\{s_i \in I\}] \cdot \prod_{i \in G} \left( \prod_j \mathbb{P}[x_{i_0,j}^p = 1] \mathbb{P}[x_{i_0,j}^m = 1] \right) \cdot \left( \prod_k \sum_{s_{i_k}} \mathbb{P}[x_{i_k}^p | x_{f(i_k)}^p, x_{f(i_k)}^m, s_{i_k}^p] \mathbb{P}[s_{i_k}^p] \right).$$

Let  $p_t(i)$  be the transmission probability for the proxy family, defined as

$$p_t(i) = \prod_k \sum_{s_{i_k}} \mathbb{P}[x_{i_k}^p | x_{f(i_k)}^p, x_{f(i_k)}^m, s_{i_k}^p] \mathbb{P}[s_{i_k}^p].$$

View this probability as a Markov chain along the genome with a state space of size  $2^4$  where each state indicates the inheritance of  $(s_{i_1}, s_{i_2}, s_{i_3}, s_{i_4})$ . The transition probabilities are given by  $P_{ij} = \theta^{H(i,j)}(1-\theta)^{4-H(i,j)}$  where  $H(i,j)$  is the Hamming distance between inheritance states  $i$  and  $j$ . By design, the transitions allowed by the data have an unusual structure dictated by the heterozygous loci of the proxy parent. Specifically, at a heterozygous locus, there is exactly one inheritance state that satisfies the children's haplotypes. At homozygous loci, all the inheritance states are allowed. So, we compute this probability using the  $l$ -state transition probabilities to determine the contribution of a particular stretch of  $l$  homozygous loci that are followed by a heterozygous locus. Notice that the heterozygous locus has, as inheritance indicators, either  $(0, 1, 1, 0)$  or  $(1, 0, 0, 1)$ , and these alternate between consecutive heterozygous loci.

Let  $Q_{0110}$  be a transition matrix having non-zero recombination probabilities only in column 0110 (i.e.  $Q_{0110,i,j} = P_{ij}$  when  $j = 0110$ ). Similarly, let  $Q_{1001}$  be a transition matrix with non-zero recombination probabilities only in column 1001. Let  $h_0$  be the number of homozygous loci that begin proxy parent's genotype and let  $h_j$  be the number of consecutive homozygous loci after the  $j$ 'th heterozygous locus where  $1 \leq j \leq T_i$  and  $T_i$  is the number of heterozygous loci for proxy parent  $i$ . Now, we can write the transmission probability in terms of matrix operations

$$p_t(i) = \left( \frac{1}{16} \right) \vec{1} \cdot P^{h_0} \cdot \prod_{j=0}^{T_i} (Z_j Q_{0110} + (1 - Z_j) Q_{1001}) P^{h_j} \cdot \vec{1}^T$$

where  $Z_j$  indicates whether the  $j$ 'th heterozygous locus has inheritance indicators  $(0, 1, 1, 0)$ . The column vector of ones at the end simply sums all final state probabilities to obtain the total probability.

Finally, notice that the two heterozygous inheritance states  $(0, 1, 1, 0)$  and  $(1, 0, 0, 1)$  are arbitrarily labeled. The main feature is that these states alternate at heterozygous loci, and it does not matter which one occurs first. So, we can write  $p_t(i)$  as in the statement of the lemma in terms of  $O_j$  which indicates the event that  $j$  is odd. Now we have a quantity that is a function of the genotype data and not dependent on the haplotypes under consideration.  $\square$

**Corollary 3.** *The mapping yields a haplotype problem instance having a likelihood and maximum probability proportional, respectively, to the likelihood and maximum probability of the genotype instance. Specifically,*

$$\sum_x \mathbb{P}_h[x] = \left( \sum_{\{x_i | i \in I\}} \mathbb{P}_g[\{x_i | i \in I\}] \right) \cdot \prod_{i \in G} p_t(i) \prod_j \mathbb{P}[x_{i_0,j}^p = 1] \mathbb{P}[x_{i_0,j}^m = 1]$$

and

$$\max_x \sum_x \mathbb{P}_h[x] = \left( \max_{\{x_i | i \in I\}} \mathbb{P}_g[\{x_i | i \in I\}] \right) \cdot \prod_{i \in G} p_t(i) \cdot \prod_j \mathbb{P}[x_{i_0,j}^p = 1] \mathbb{P}[x_{i_0,j}^m = 1]$$

where  $p_t(i)$  is proxy family  $i$ 's transmission probability as defined in Lemma 2.

*Proof.* Lemma 2 shows that  $X$  is independent of the coefficient of proportionality between the haplotype probability and the genotype probability. Therefore, this coefficient factors out of both the likelihood and the maximum probability equations.  $\square$

Although this reduction establishes the hardness of these haplotype pedigree problems, it does so by constructing children whose haplotypes require many recombinations and would be extremely unlikely to occur naturally. Accordingly, we suspect that realistic instances of these haplotyping problems may provide more information about the locations of recombinations than genotype instances.

### Algorithms and Accuracy of Estimates

One indication that the haplotype problem might be practically more tractable is the amount of information in the haplotype data relative to the genotype data. To understand this, we can consider a pedigree with a fixed set of sampled individuals. Assume that there are two input data sets available, either the haplotype or the genotype data, for all the sampled individuals. Note that the alleles observed will be identical in both the haplotype and genotype data, so we are interested in the distribution that these data impose on the inheritance probabilities. By comparing the accuracy of the recombination estimates under these two data sets, we can get an idea for how useful the respective probability distributions are.

Let  $R_j$  be a random variable representing the number of recombinations in the whole pedigree that occur between loci  $j - 1$  and  $j$ . Similar to our notation before,  $R_j = \sum_{i \in I} R_{i,j}^p + R_{i,j}^m$ . We want to compute the distribution of  $R_j$  under both the genotype and haplotype inheritance probability distributions. These two inheritance distributions are different precisely because there are haplotypes and inheritance paths that are consistent with the genotype constraints but disallowed by the haplotype constraints.

These distributions are obtained by constructing a hidden Markov model for the linkage dependencies along the genome. At each locus, the HMM considers the constraints given by either the haplotype or genotype data (i.e. the haplotype data HMM is a variation on the Lander-Green algorithm [15]). We first use the forward-backward algorithm to compute the marginal inheritance probabilities for each locus using a hidden Markov model. Once we have the marginal probabilities, we can easily obtain the distribution for  $R_j$ .

### General Haplotype and Genotype HMMs

The likelihood can be modeled using a hidden Markov model along the genome with inheritance paths as hidden states. An *inheritance path* is a graph with nodes being the alleles of individuals and directed edges between alleles that are inherited from parent to child. The transition probabilities are functions of  $\theta$  and the number of recombinations between a given pair of inheritance graphs.

Given the data, we compute the marginal inheritance path probabilities at each site by using the forward-backward algorithm for HMMs. Sobel and Lange described a method for enumerating the inheritance paths compatible with the allele data observed at each locus [16]. There are at most  $k = 2^{2^{|\Lambda F|}}$  inheritance paths when  $\Lambda F$  is the set of non-founder individuals, and both the forward and backward recursions do an  $O(k^2)$  calculation at each site.

To compute the analogous probability for haplotype data, we use a similar HMM. For haplotypes, the hidden states must consider the haplotype orientations, which specify the parental origins of all the observed haplotypes. Notice that these orientations are not equivalent to inheritance paths, since they only specify inheritance edges between haplotyped individuals and their parents. For each of the  $2^{2^{|H|}}$  haplotype orientations, where  $H$  is the set of haplotyped individuals, we enumerate the inheritance paths compatible with the haplotype alleles, their orientations, and the pedigree relationships. Alternatively, each of the inheritance paths enumerated for the genotype algorithm induces a particular orientation on the haplotypes heterozygous for that locus (i.e. parental origin of the entire haplotype). Thus, the hidden states for the haplotype HMM are the cross-product of the orientations and the inheritance paths.

The haplotype HMM has transition probabilities that are nearly identical to the genotype HMM with the exception that transitions between inheritance paths with different haplotype orientations have probability zero. Recombinations are only allowed when they do not occur between typed haplotypes.

The forward-backward algorithm is also used on the haplotype HMM. However, there are  $2^{2^{(|I|+|H|-|F|)}}$  hidden states, yielding a slightly slower calculation. Fortunately, the haplotype recursions can be run simultaneous with the genotype recursions, meaning that the inheritance paths need only be enumerated once.

### Haplotype Likelihoods in Linear Time

There is one obvious instance of the haplotyping problems where there are polynomial-time algorithms. Even though it is impractical to assume that we can sample genetic material from deceased individuals in a multi-generational pedigree, for a moment, let us consider the case where all the individuals in the pedigree are haplotyped.

The Elston-Stewart algorithm [14] for genotype data has a direct analogue for haplotype data. This algorithm calculates the likelihood via the belief propagation algorithm by eliminating individuals recursively from the bottom up. Each individual is "peeled off", after their descendants have been peeled off, by using



a forward-backward algorithm on the HMM for the mother-father-child trio.

The haplotype version of this algorithm is linear when all the individuals are haplotyped, since each elimination step is conditionally independent of all the others. Given the parents' haplotypes, regardless of which was inherited from which grand-parent, the probability of the child's haplotype is independent of all other trios. Therefore, we can take a product over the likelihoods for all the trios, and compute each trio likelihood using a 4-state HMM. Then for  $k$  non-founding individuals, and  $l$  loci, this algorithm has  $O(kl)$  running time.

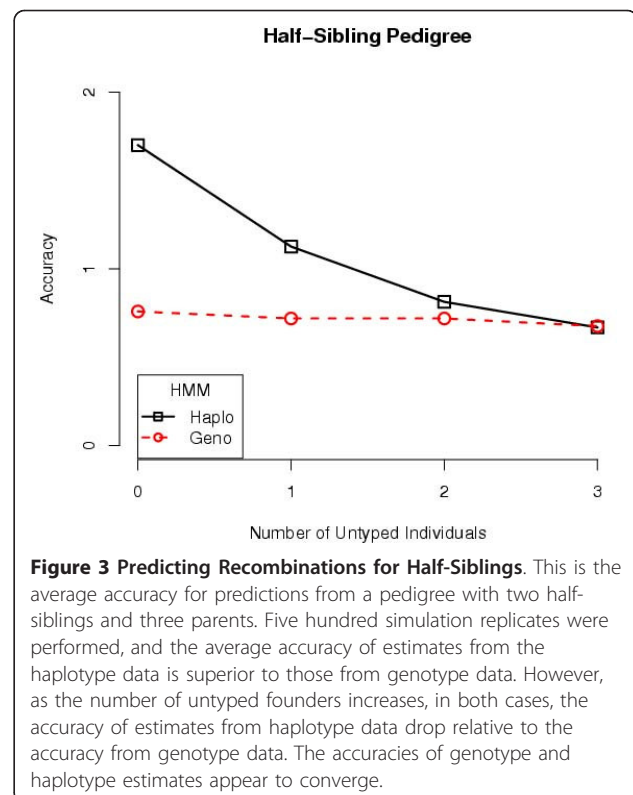
This same intuition carries through to the minimum recombination problem, and each trio can be considered independent of the others. This contrasts with the genotype minimum recombination problem which is known to be hard, even when all the individuals are genotyped [7].

## Results

To simulate realistic pedigree data, SNPs were selected from HapMap that span 100 mb on both sides of a loosely-linked pair of sites. There are 40 SNPs total, with 20 tightly linked SNPs on each side of a strong recombination breakpoint having  $\theta = 0.25$ . The haplotypes for these SNPs were selected randomly from HapMap. Pedigree haplotype and genotype data were simulated for each child by uniformly selecting one of the parental alleles for the first locus, and subsequent loci were selected on the same parental haplotype with probability  $\theta_j$  for each locus  $j$ . Inheritance was simulated for 500 simulation replicates.

The simulation yielded completely typed pedigrees. For each pedigree, we removed the genotype and haplotype information for increasing numbers of untyped individuals. For each instance of a specific number of untyped individuals, two values were computed on the estimated number of recombinations between the central pair of loci: the haplotype and genotype accuracies. Accuracy was computed as a function of the  $l_1$  distance between the deterministic number of recombinations and the calculated distribution. Specifically, accuracy was  $2 - \sum_{i \geq 0} |x_i - a_i|$ , where  $x_i$  was the estimated probability for  $i$  recombinations and  $a_i$  was the deterministic indicator of whether there were  $i$  recombinations in the data simulated on the pedigree.

In all the instances we observed a trend where the best accuracy was obtained with haplotype data where everyone in the pedigree was haplotyped. For example, a five-individual pedigree with two half-siblings is shown in Figure 3. With the three founders untyped, the haplotype data yielded similar accuracy as the genotype data. Consider a three-generation pedigree having two parents, their two children, an in-law, and a grandchild for a total of six individuals, three of them founders. This

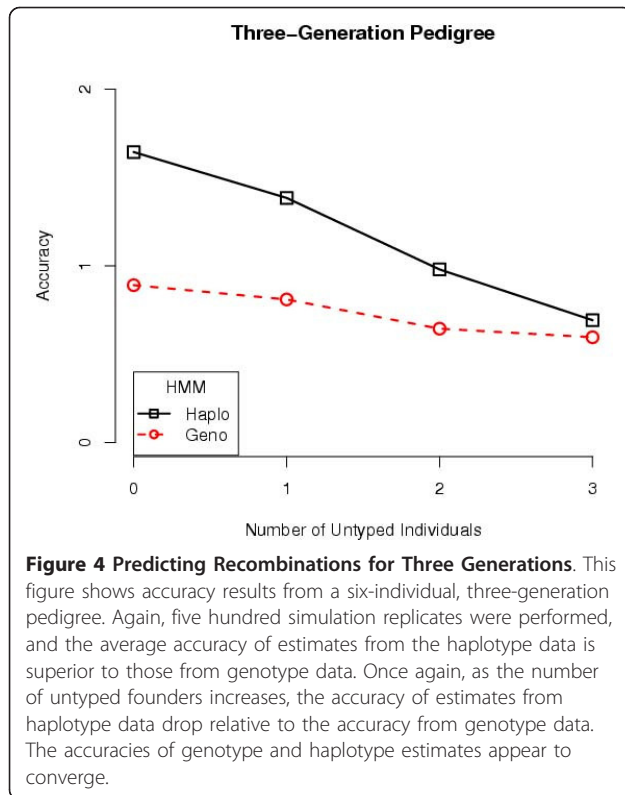


pedigree has a similar trend in accuracy as the number of untyped founders increases, Figure 4. As the number of untyped individuals increases, the accuracies of genotype and haplotype estimates appear to converge.

## Discussion

Sequencing technologies would seem to solve the phasing problem by yielding haplotype data. However, if we wish to consider diploid inheritance with recombination, the phasing problem remains, even when we are given chromosome-length haplotype data. This is demonstrated by reduction of the phasing problem for genotypes to the phased version of the same problem for three common pedigree problems. This theoretical result is due largely to the unavailability of genetic material for deceased individuals.

Three pedigree calculations were discussed: likelihood, maximum probability, and minimum recombination. Each of these calculations on haplotype data have the same computational complexity as the same computation on genotype data. In the worst case, it takes only a single generation to remove the correlation between sites in the haplotype. This worst case provided the reduction that proves the complexity results for the haplotype computations, and it worked equally well for all three pedigree computations. The worst-case is not biologically realistic, since it requires roughly  $2(m - 1)$



recombinations for  $m$  sites in 4 meioses. This is very unlikely to occur under typical models for inheritance. To investigate more likely scenarios, sequences were simulated in a region of the genome surrounding a recombination breakpoint. From haplotype and genotype data, we estimated the distribution of the number of recombinations at the breakpoint and compared the estimates to the ground-truth for accuracy.

When typing everyone in the pedigree, the estimates from haplotype data were very accurate, because the haplotype data provides enough constraints to determine where the recombinations must have occurred. With decreasing numbers of typed individuals, the accuracy of haplotype-based estimates dropped until it seemed to converge to the genotype accuracy due to a lack of constraints. From the structure of the calculations, we observed that with fewer typed individuals there were more haplotype orientations to consider, and the haplotype calculation more closely resembled the genotype calculation. However, the haplotype calculation had more constraints and lost accuracy at a slower rate.

Several interesting open problems remain. First, approximation algorithms might be a useful approach for haplotypes on pedigrees. The existence of a linear-time algorithm when all individuals are haplotyped may suggest that the general haplotype problem instance could be amenable to approximation algorithms.

Second, these proofs apply when there is no missing data in a genotyped individual (i.e. a proxy parent).

The proof requires knowing whether the proxy parent is heterozygous or homozygous at each locus, and this is unknown when there is missing data. Third, there is an interesting case of mixed haplotypes and genotypes. For this case to be interesting, the ends of haplotypes must occur at different locations in different individuals in the pedigree. Otherwise, the haplotypes that start and end at the same positions in all individuals can easily be converted into multi-allelic genotypes, with an allele for each haplotype. The mixed haplotype-genotype problem is not amenable to the proof techniques used here. However, the haplotype HMM in Section can easily be revised to handle the mixed case. This is important because the data produced by single polymer sequencing is more likely to resemble the mixed case than either the haplotype or the genotype cases.

#### Acknowledgements

I want to thank Richard M. Karp for reviewing a draft of the manuscript and the National Science Foundation for support through the Graduate Research Fellowship.

#### Author details

<sup>1</sup>Electrical Engineering and Computer Sciences, University of California Berkeley, Berkeley, CA 94720-1776, USA. <sup>2</sup>International Computer Science Institute, 1947 Center St. Suite 600, Berkeley, CA 94704, USA.

#### Authors' contributions

BK conceived of the problem, proved the results, and implemented the algorithms.

#### Competing interests

The authors declare that they have no competing interests.

Received: 10 August 2010 Accepted: 19 April 2011

Published: 19 April 2011

#### References

- Coop G, Wen X, Ober C, Pritchard J, Przeworski M: High-Resolution Mapping of Crossovers Reveals Extensive Variation in Fine-Scale Recombination Patterns Among Humans. *Science* 2008, **319**(5868):1395-1398.
- MY N, DF L, et al: Meta-analysis of 32 genome-wide linkage studies of schizophrenia. *Mol Psychiatry* 2009, **14**:774-85.
- Romero I, Ober C: CFTR mutations and reproductive outcomes in a population isolate. *Human Genet* 2008, **122**:583-588.
- Geiger D, Meek C, Wexler Y: Speeding up HMM algorithms for genetic linkage analysis via chain reductions of the state space. *Bioinformatics* 2009, **25**(12):i196.
- Xiao J, Liu L, Xia L, Jiang T: Efficient Algorithms for Reconstructing Zero-Recombinant Haplotypes on a Pedigree Based on Fast Elimination of Redundant Linear Equations. *SIAM Journal on Computing* 2009, **38**:2198.
- Piccolboni A, Gusfield D: On the Complexity of Fundamental Computational Problems in Pedigree Analysis. *Journal of Computational Biology* 2003, **10**(5):763-773.
- Li J, Jiang T: An Exact Solution for Finding Minimum Recombinant Haplotype Configurations on Pedigrees with Missing Data by Integer Linear Programming. *Proceedings of the 7th Annual International Conference on Research in Computational Molecular Biology* 2003, 101-110.
- Thattai BD: Combinatorics of Pedigrees I: Counterexamples to a Reconstruction Question. *SIAM Journal on Discrete Mathematics* 2008, **22**(3):961-970.

9. Eid J, *et al*: **Real-Time DNA Sequencing from Single Polymerase Molecules**. *Science* 2009, **323**(5910):133-138.
10. Barrett J, Hansoul S, Nicolae D, Cho J, Duerr R, Rioux J, Brant S, Silverberg M, Taylor K, Barmada M, *et al*: **Genome-wide association defines more than 30 distinct susceptibility loci for Crohn's disease**. *Nature Genetics* 2008, **40**:955-962.
11. Chen WM, Abecasis G: **Family-Based Association Tests for Genomewide Association Scans**. *American Journal of Human Genetics* 2007, **81**:913-926.
12. Burdick J, Chen W, Abecasis G, Cheung V: **In silico method for inferring genotypes in pedigrees**. *Nature Genetics* 2006, **38**:1002-1004.
13. Kirkpatrick B, Halperin E, Karp R: **Haplotype Inference in Complex Pedigrees**. *Journal of Computational Biology* 2010, **17**(3):269-280.
14. Elston R, Stewart J: **A general model for the analysis of pedigree data**. *Human Heredity* 1971, **21**:523-542.
15. Lander E, Green P: **Construction of Multilocus Genetic Linkage Maps in Humans**. *Proceedings of the National Academy of Science* 1987, **84**(5):2363-2367.
16. Sobel E, Lange K: **Descent Graphs in Pedigree Analysis: Applications to Haplotyping, Location Scores, and Marker-Sharing Statistics**. *American Journal of Human Genetics* 1996, **58**(6):1323-1337.

doi:10.1186/1748-7188-6-10

**Cite this article as:** Kirkpatrick: Haplotypes versus genotypes on pedigrees. *Algorithms for Molecular Biology* 2011 **6**:10.

**Submit your next manuscript to BioMed Central  
and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
[www.biomedcentral.com/submit](http://www.biomedcentral.com/submit)

