

# Statistical Tests for Detection of Misspecified Relationships by Use of Genome-Screen Data

Mary Sara McPeck and Lei Sun

Department of Statistics, University of Chicago, Chicago

Misspecified relationships can have serious consequences for linkage studies, resulting in either reduced power or false-positive evidence for linkage. If some individuals in the pedigree are untyped, then Mendelian errors may not be observed. Previous approaches to detection of misspecified relationships by use of genotype data were developed for sib and half-sib pairs. We extend the likelihood calculations of Göring and Ott and Boehnke and Cox to more-general relative pairs, for which identity-by-descent (IBD) status is no longer a Markov chain, and we propose a likelihood-ratio test. We also extend the identity-by-state (IBS)-based test of Ehm and Wagner to nonsib relative pairs. The likelihood-ratio test has high power, but its drawbacks include the need to construct and apply a separate Markov chain for each possible alternative relationship and the need for simulation to assess significance. The IBS-based test is simpler but has lower power. We propose two new test statistics—conditional expected IBD (EIBD) and adjusted IBS (AIBS)—designed to retain the simplicity of IBS while increasing power by taking into account chance sharing. In simulations, the power of EIBD is generally close to that of the likelihood-ratio test. The power of AIBS is higher than that of IBS, in all cases considered. We suggest a strategy of initial screening by use of EIBD and AIBS, followed by application of the likelihood-ratio test to only a subset of relative pairs, identified by use of EIBD and AIBS. We apply the methods to a Genetic Analysis Workshop 11 data set from the Collaborative Study on the Genetics of Alcoholism.

## Introduction

Pedigree error, or misclassification of the relationship between individuals, can have potentially serious consequences for linkage studies. It can lead to false positive evidence for linkage or reduce the power of linkage detection. In some cases, misspecified relationships can be identified through discovery of Mendelian errors. However, if some individuals in the pedigree are untyped, then Mendelian errors may not result, so other methods are needed to test for deviation from the reported relationship. Genome-screen data collected for mapping studies have the potential to be highly informative for verifying relationships among individuals.

Statistical methods for detecting misspecified relationships based on genotype data have been developed specifically for the cases of sib pairs and half-sib pairs. Göring and Ott (1997) and Boehnke and Cox (1997) compute the likelihood of the observed genotype data for the pair in these cases, under the assumption of no interference. In the case when the sib pair has a single

typed parent, Göring and Ott (1997) also compute the likelihood conditional on the genotype of the parent. For each putative sib pair, Göring and Ott (1997) assign prior probabilities to the relationships sib, half-sib, and unrelated, and compute the posterior probabilities of the relationships given the data. Boehnke and Cox (1997) use a hidden Markov model to calculate the likelihood of the data for each of a set of possible relationships for the pair, namely, sib, half-sib, unrelated, and MZ twin, and they select the one that maximizes the likelihood. Ehm and Wagner (1998) propose an approximately normally distributed test statistic based on the number of alleles shared identical by state (IBS) by a sib pair, summed over a large number of genetic markers.

We consider the problem of relationship testing for more general relative pairs. We extend both the IBS-based test of Ehm and Wagner (1998) and the likelihood calculation of Göring and Ott (1997) and Boehnke and Cox (1997) to more general outbred relative pairs. In their likelihood calculations, both Göring and Ott (1997) and Boehnke and Cox (1997) assume that the identity-by-descent (IBD) process for a relative pair is Markov. However, as noted by Donnelly (1983) and Feingold (1993), the Markov assumption fails to hold for all but the simplest relative pairs. For instance, it does not hold for avuncular and first-cousin relationships, even under the assumption of no interference (see

Received January 25, 1999; accepted December 29, 1999; electronically published March 9, 2000.

Dr. Mary Sara McPeck, Department of Statistics, University of Chicago, 5734 South University Avenue, Chicago, IL 60637. E-mail: mcpeek@galton.uchicago.edu

© 2000 by The American Society of Human Genetics. All rights reserved.  
0002-9297/2000/6603-0028\$02.00

both the *IBD Process for a Pair of Relatives Often Non-Markov* subsection, in the Methods section, and Appendix A). Donnelly (1983) showed how to construct an augmented IBD process that is Markov under no interference. We apply a hidden Markov method to the augmented IBD process to calculate the likelihood. In Appendix B, we describe an extension allowing for interference. In Appendix C, we describe an extension to inbred relative pairs. Using the above likelihood calculation, we propose a likelihood-ratio test for misspecified relationships. In addition, we propose two new test statistics that do not require specification of an alternative relationship, yet do take into account chance sharing, thus combining some of the strong points of both the IBS and likelihood-ratio tests.

**Methods**

*IBD Process for a Pair of Relatives Often Non-Markov*

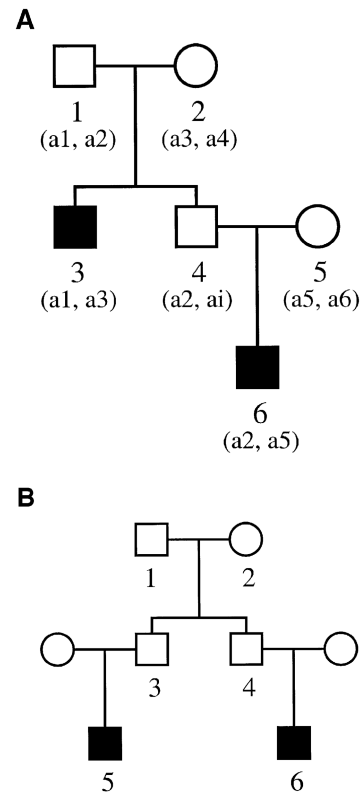
A set of alleles is said to be IBD if the alleles are copies inherited from the same ancestral allele. IBD is to be distinguished from IBS; a set of alleles being IBS simply means that the set of alleles is of the same observed allelic type. A set of alleles being IBD implies that they are IBS (ignoring the possibilities of genotyping error and mutation), but not vice versa.

Define the IBD process  $\{D\}$  for a pair of outbred relatives, call them individuals 1 and 2, by letting  $D_m$  equal the number of alleles shared IBD by the pair at marker  $m$ ,

$$D_m = 1_{g_{11}=g_{21}} + 1_{g_{11}=g_{22}} + 1_{g_{12}=g_{21}} + 1_{g_{12}=g_{22}} ,$$

where  $1_{g_{1i}=g_{2j}}$  is the indicator of the event that allele  $i$  of individual 1 and allele  $j$  of individual 2 at marker  $m$  are IBD, with arbitrary labeling of the two alleles of an individual. For outbred relative pairs,  $D_m$  takes values in  $\{0, 1, 2\}$  for each  $m$ . In order to calculate the likelihood in the cases of sib pairs and half-sib pairs, Goring and Ott (1997) and Boehnke and Cox (1997) make the assumption that the IBD process  $\{D\}$  is Markov. First, this assumption does not hold in the presence of interference, so these works make an implicit assumption of no interference. Second, although the Markov assumption holds, under the assumption of no interference, in the cases of full sibs, half-sibs, parent-child, and grandparent-grandchild (hereafter denoted as “grand-PC”), it does not hold for general relationships even if no interference is assumed. In particular, it fails in the cases of avuncular and first-cousin relationships (Donnelly 1983; Feingold 1993).

To understand why this is so, first consider the avuncular case. Let the individuals be labeled 1-6 as in fig. 1A, with individuals 3 and 6 forming an avuncular pair.



**Figure 1** A, Pedigree for avuncular pair. B, Pedigree for first-cousin pair.

The Markov property requires that conditional on the IBD value  $D_A$  for the avuncular pair at a locus  $A$ , the IBD values at loci to the right of locus  $A$  are independent of the IBD values at loci to the left of locus  $A$ . The violation of the Markov property in the avuncular case arises as follows: conditional on the number of alleles shared IBD by individuals 3 and 6 at locus  $A$ , if the  $A$  allele not transmitted from individual 4 to individual 6 is shared IBD by individuals 3 and 4 (call this event  $S_A$ ), then the chance is increased that individuals 3 and 6 share an allele IBD at any other locus linked to  $A$ . This induces a positive correlation in sharing at loci linked to  $A$ , conditional on IBD sharing at  $A$  (see Appendix A). By conditioning on the event  $S_A$  or its complement, we show in Appendix A that if locus  $B$  is to the right of locus  $A$  and locus  $C$  is to the left of locus  $A$ , both linked to  $A$ , then for the avuncular pair 3 and 6,

$$P(D_C = 1 | D_A = j, D_B = 1) > P(D_C = 1 | D_A = j) , \quad (1)$$

violating the Markov property.

The first-cousin case is shown in figure 1B. Let  $D_A$  be the number of alleles shared IBD by the cousin pair 5 and 6 at locus  $A$ , let  $I_A$  be the indicator of the event that the cousins’ paternally inherited alleles at  $A$  came from

the same grandparent, and let  $D'_A$  be the number of alleles shared IBD by individuals 3 and 4. Conditional on  $D_A$ , either holding  $D'_A$  fixed and increasing  $I_A$  or holding  $I_A$  fixed and increasing  $D'_A$  increases the chance that the cousin pair shares one IBD at any locus linked to  $A$ . As in the avuncular case, this induces a positive correlation in sharing at loci linked to  $A$ , conditional on the number of alleles shared IBD at  $A$ , resulting in a similar violation of the Markov property.

*Likelihood Calculation for an Outbred Relative Pair Under No Interference*

Let  $\{D\}$  be the IBD process for a pair of outbred relatives. To extend the likelihood calculation of Goring and Ott (1997) and Boehnke and Cox (1997) to the case when  $\{D\}$  is no longer Markov, we construct an augmented IBD process  $\{A\}$  that is Markov under the assumption of no interference and that contains all the information of the process  $\{D\}$ . For the avuncular and first-cousin cases, we give the state spaces for such augmented IBD Markov processes  $\{A\}$  in, respectively, table 1A and B, with transition matrices given in table 2A and B, and with the probability distributions of the next state entered, conditional on the current state, and the leaving rate for each state given in table 3A and B. Note that in general, under no interference, the augmented Markov process could be chosen to be  $\{A'\}$ , where  $A'_m$  is the inheritance vector at marker  $m$  defined by Lander and Green (1987), or  $\{A''\}$ , used by Feingold (1993) in the avuncular and first-cousin cases, where  $A''_m$  is the equivalence class of inheritance vectors at  $m$  obtained by identifying inheritance vectors that differ only by interchanges of maternal and paternal haplotypes within founders (Kruglyak et al. 1996). The use of  $\{A'\}$  or  $\{A''\}$  provides an automatic way to construct the augmented Markov chain, although these processes contain unnecessary information. For instance, in the avuncular case, the state space of the process  $\{A'\}$  is of size 64, and that for  $\{A''\}$  is of size 8, whereas our augmented process  $\{A\}$  requires only 4 states. Similarly, in the first-cousin case, the state space of the process  $\{A'\}$  is of size 256, that for  $\{A''\}$  is of size 16, and that for our augmented process  $\{A\}$  is 7. In both cases, our augmented chain  $\{A\}$  has the minimal number of states needed to both contain all the information of the IBD process  $\{D\}$  and satisfy the Markov property under no interference. Donnelly (1983) constructs similar minimal-state augmented chains for a number of relationships such as “ $m$ th generation descendant” and “ $s$ th cousin  $t$  times removed.” The problem of implementing an automated method for generating a minimal-state augmented Markov chain for any given pairwise relationship is still an open one.

In the cases of full sibs and half-sibs, for which the IBD process  $\{D\}$  is Markov under no interference,

**Table 1**

**State Spaces of Augmented Markov Chains for Avuncular and First-Cousin Pairs**

STATE LABEL	VALUE THAT DEFINES STATE <sup>a</sup>		
	IBD(A3,A4)	IBD(A3,A6)	
A. Avuncular chain:			
AV1	0	0	
AV2	1	0	
AV3	1	1	
AV4	2	1	
	IBD(B3,B4)	IBD(B5,B6)	G(B5,B6)
B. First-cousin chain:			
FC1	0	0	0
FC2	0	0	1
FC3	1	0	0
FC4	1	0	1
FC5	1	1	1
FC6	2	0	0
FC7	2	1	1

<sup>a</sup> “IBD( $X_i, X_j$ )” is the number of alleles shared IBD by individuals  $i$  and  $j$ , where individuals are as labeled in figure 1X, where “ $X$ ” denotes the figure panel (A or B); “G(B5,B6)” is the indicator of event that allele inherited by individual 5 from individual 3 and the allele inherited by individual 6 from individual 4 are both descended from either individual 1 or individual 2—that is, that they are both from the same grandparent. Individuals are labeled as in figure 1B).

Boehnke and Cox (1997) use a hidden Markov method to calculate the probability  $P_R(G_1, G_2, \dots, G_{n_c})$ , where  $n_c$  is the number of markers on the  $c$ th chromosome,  $G_m$  is the genotype data for the pair at marker  $m$ , and the subscript  $R$  denotes calculation of the probability under the assumed relationship  $R$ . When the IBD status of a pair is not Markov, the likelihood can still be calculated by use of a hidden Markov method, applying the algorithm of Baum (1972). However the hidden Markov method is applied to the augmented Markov chain  $\{A\}$  instead of to the IBD process  $\{D\}$ . The calculation, using Baum’s (1972) forward probabilities, is summarized as follows: for any given relationship  $R$ , we define  $\alpha_1(j) = P_R(A_1 = j)$  and  $\alpha_k(j) = P_R(G_1, G_2, \dots, G_{k-1}, A_k = j)$ , for  $k > 1$ . Note that  $\alpha_1(j)$  is just the stationary distribution of the augmented Markov chain  $\{A\}$  for relationship  $R$ . The stationary distribution for our augmented avuncular Markov chain is  $\pi_1 = \pi_2 = \pi_3 = \pi_4 = \frac{1}{4}$ , and the stationary distribution for our augmented first-cousin Markov chain is  $\pi_1 = \pi_2 = \pi_4 = \pi_5 = \pi_6 = \pi_7 = \frac{1}{8}$ ,  $\pi_3 = \frac{1}{4}$ . Then the recursion formula, similar to that in Boehnke and Cox (1997), is  $\alpha_{k+1}(j) = \sum_i \alpha_k(i) P_R(A_{k+1} = j | A_k = i) P(G_k | A_k = i)$ , where  $P_R(A_{k+1} = j | A_k = i)$  is the transition probability of the augmented Markov chain, which, for the cases of avuncular and first-cousin pairs, is given in, respectively, table 2A and B. Since the augmented Markov chain  $\{A\}$  contains all the information of the IBD process  $\{D\}$ , and IBD status is sufficient to calculate the conditional prob-

**Table 2**

**Transition Matrices of Augmented Markov Chains for Avuncular and First-Cousin Pairs, Where  $\psi = \theta^2 + (1 - \theta)^2$  and  $\phi = 1 - \psi = 2\theta(1 - \theta)$**

CURRENT STATE	STATE AT RECOMBINATION FRACTION $\theta$ FROM CURRENT STATE						
	AV1	AV2	AV3	AV4			
<b>A. Avuncular chain:<sup>a</sup></b>							
AV1	$\psi^2$	$\psi\phi$	$\psi\phi$	$\phi^2$			
AV2	$\psi\phi$	$(1 - \theta)\psi^2 + \theta\phi^2$	$\theta\psi^2 + (1 - \theta)\phi^2$	$\psi\phi$			
AV3	$\psi\phi$	$\theta\psi^2 + (1 - \theta)\phi^2$	$(1 - \theta)\psi^2 + \theta\phi^2$	$\psi\phi$			
AV4	$\phi^2$	$\psi\phi$	$\psi\phi$	$\psi^2$			
<b>B. First cousin:<sup>b</sup></b>							
FC1	$\psi^3$	$\psi^2\phi$	$2\psi^2\phi$	$\psi\phi^2$	$\psi\phi^2$	$\psi\phi^2$	$\phi^3$
FC2	$\psi^2\phi$	$\psi^3$	$2\psi\phi^2$	$\psi^2\phi$	$\psi^2\phi$	$\phi^3$	$\psi\phi^2$
FC3	$\psi^2\phi$	$\psi\phi^2$	$\psi^3 + \psi\phi^2$	$\theta(1 - \theta)(\psi^2 + \phi^2)$	$\theta(1 - \theta)(\psi^2 + \phi^2)$	$\psi^2\phi$	$\psi\phi^2$
FC4	$\psi\phi^2$	$\psi^2\phi$	$\psi^2\phi + \phi^3$	$(1 - \theta)^2\psi^2 + \theta^2\phi^2$	$\theta^2\psi^2 + (1 - \theta)^2\phi^2$	$\psi\phi^2$	$\psi^2\phi$
FC5	$\psi\phi^2$	$\psi^2\phi$	$\psi^2\phi + \phi^3$	$\theta^2\psi^2 + (1 - \theta)^2\phi^2$	$(1 - \theta)^2\psi^2 + \theta^2\phi^2$	$\psi\phi^2$	$\psi^2\phi$
FC6	$\psi\phi^2$	$\phi^3$	$2\psi^2\phi$	$\psi\phi^2$	$\psi\phi^2$	$\psi^3$	$\psi^2\phi$
FC7	$\phi^3$	$\psi\phi^2$	$2\psi\phi^2$	$\psi^2\phi$	$\psi^2\phi$	$\psi^2\phi$	$\psi^3$

<sup>a</sup> States are as labeled in table 1A.

<sup>b</sup> States are as labeled in table 1B.

ability of genotype data, we obtain  $P(G_k|A_k = i) = P(G_k|D_k = \text{IBD status associated with state } i \text{ of } A)$ , where the IBD status associated with state  $i$  of  $A$  is given in table 1A and B, for the cases of avuncular and first-cousin pairs, respectively. Thus, for the outbred case, computation of the probabilities  $P(G_k|A_k = i)$  can be reduced to computation of the probabilities  $P(G_k|D_k = j)$  that appear in Thompson (1975). The summation  $\sum_j \alpha_{n_c}(j)$  gives  $P_R(G_1, G_2, \dots, G_{n_c})$  for the  $c$ th chromosome, and the  $P_R(G_1, G_2, \dots, G_{n_c})$  are multiplied over all chromosomes  $c$  to complete calculation of the likelihood of the genotype data throughout the genome.

In Appendix B, we discuss an extension of the likelihood calculation to take into account interference, using the chi-square model. In Appendix C, we discuss extension to the case of an inbred relative pair.

*Likelihood-Ratio Test (LRT) and Maximized Likelihood-Ratio test (MLRT)*

In order to test the null relationship against a specific alternative relationship, we calculate the likelihoods  $L_O$  and  $L_A$  for the data under the null and alternative relationships, respectively, using the above method. The test statistic for the likelihood-ratio test (LRT) is the log-likelihood ratio  $\text{LLR} = \log(L_A) - \log(L_O)$ . We assess significance by simulation under the null hypothesis. In our simulations, we find that in fact the normal distribution provides a very close approximation to the distribution of LLR under the null (or alternative) relationship when genome-screen data are used. However, the null mean and variance of LLR would still need to be obtained by

simulation in order to use the normal approximation for assessing significance.

In practice, one often does not have a specific alternative relationship in mind. In that case, one can maximize the likelihood over a set of alternatives  $\mathcal{A}$  to obtain  $\hat{L}_{\mathcal{A}}$ , and then consider the maximized log-likelihood ratio  $\text{MLLR} = \log(\hat{L}_{\mathcal{A}}) - \log(L_O)$ , where this statistic depends on the particular set of alternatives  $\mathcal{A}$  considered. In our data analysis and simulations,  $\mathcal{A}$  typically consists of sib, half-sib, grand-PC, avuncular, unrelated and first-cousin relationships (excluding the null relationship). Unlike LLR, MLLR often has a rather skewed distribution, so the normal approximation is not appropriate. Thus we require simulation to assess significance for the test by calculation of empirical  $p$  values. Note that one may detect a misspecified relationship by use of the maximum likelihood-ratio test (MLRT), although the true alternative relationship may not be in the class  $\mathcal{A}$ . For a data set in which thousands of pairwise relationships are examined, use of a large alternative class  $\mathcal{A}$  may not be computationally feasible. Judicious choice of  $\mathcal{A}$  may give power against a wide range of alternatives, including many not in  $\mathcal{A}$ . If significant deviation from the null is detected in a few cases, a larger class  $\mathcal{A}$  may be considered for relationship estimation in those cases.

*Markov Approximation to LRT and MLRT*

One strategy, which would eliminate the need to construct augmented Markov chains for a wide variety of relationships, and which would reduce the computational burden of the MLRT or LRT, would be to ap-

**Table 3**  
**Probability Distributions for Next State Entered, Conditional on Current State, with Leaving Rates (in Terms of Genetic Distance)**

CURRENT STATE	PROBABILITY THAT NEXT STATE ENTERED IS							LEAVING RATE				
	AV1	AV2	AV3	AV4	FC1	FC2	FC3		FC4	FC5	FC6	FC7
Avuncular chain of table 2A:												
AV1	0	1/2	1/2	0								4
AV2	2/5	0	1/5	2/5								5
AV3	2/5	1/5	0	2/5								5
AV4	0	1/2	1/2	0								4
First-cousin chain of table 2B:												
FC1	0	1/3	2/3	0	0	0	0	0	0	0	0	6
FC2	1/3	0	0	1/3	1/3	0	0	0	0	0	0	6
FC3	1/3	0	0	1/6	1/6	1/3	0	0	0	0	0	6
FC4	0	1/3	1/3	0	0	0	1/3	0	0	0	0	6
FC5	0	1/3	1/3	0	0	0	1/3	0	0	0	0	6
FC6	0	0	2/3	0	0	0	1/3	0	0	0	0	6
FC7	0	0	0	1/3	1/3	1/3	0	0	0	0	0	6

proximate the IBD process  $\{D\}$  by a Markov process  $\{B\}$ , with the correct conditional probabilities  $P(D_{m2} = j | D_{m1} = i)$  used as transition probabilities for  $\{B\}$ , where  $m1$  and  $m2$  label adjacent markers; that is,  $\{B\}$  is a Markov chain on the set of markers, with  $P(B_{m1} = i) = P(D_{m1} = i)$  for all markers  $m1$ , and  $P(B_{m2} = j | B_{m1} = i) = P(D_{m2} = j | D_{m1} = i)$  for all pairs of adjacent markers  $m1$  and  $m2$ . We approximate the likelihood of the data by letting  $\{B\}$  represent the IBD process; that is, we treat the IBD process as Markov, although it is not. The likelihood for the data would then be calculated with  $\{B\}$  as the hidden Markov chain, as in Boehnke and Cox (1997). Algorithms for calculation of  $P(D_{m2} = j | D_{m1} = i)$  for outbred relationships are discussed by Denniston (1975), Thompson (1988), and Tiwari and Elston (1999). Explicit formulae are given by Bishop and Williamson (1990) for the relationships sib, half-sib, parent-child, grand-PC, avuncular, and first cousin. In our data analysis and simulation, we use the correct likelihood, not the approximation. However, we also compare the correct and approximate likelihoods for the avuncular and first-cousin cases in the Results section.

*Test Based on IBS*

Ehm and Wagner (1998) propose an approximately normally distributed test statistic, which we call  $S'$ , based on half the number of alleles shared IBS summed up over a large number of markers:  $S' = \sum_m S_m / 2$ , where  $S_m$  is the number of alleles shared IBS by the pair at locus  $m$ . Letting  $(g_{11}, g_{12})$  and  $(g_{21}, g_{22})$  be the genotypes of individuals 1 and 2 at locus  $m$ , and letting  $g_{1i} \approx g_{2j}$  denote the event that allele  $i$  of individual 1 and allele  $j$  of individual 2 at locus  $m$  are IBS, then  $S_m$  is defined as follows:

$$S_m = 2 \text{ if and only if } (g_{11} \approx g_{21} \text{ and } g_{12} \approx g_{22})$$

$$\text{or } (g_{11} \approx g_{22} \text{ and } g_{12} \approx g_{21}) ;$$

$$S_m = 0 \text{ if and only if } g_{11} \not\approx g_{21} \text{ and } g_{11} \not\approx g_{22}$$

$$\text{and } g_{12} \not\approx g_{21} \text{ and } g_{12} \not\approx g_{22} ;$$

$$S_m = 1 \text{ otherwise .}$$

In the case of sib pair, Ehm and Wagner (1998) describe the calculation of the mean and variance of  $S'$  and apply the normal approximation to assess significance. Note that Ehm and Wagner (1998) consider a one-sided hypothesis test, whereas all of our hypotheses are two sided. A two-sided test is appropriate even for the case in which the null relationship is sib pair, because excess sharing over the null could indicate that the sib pair is inbred.

For more general relative pairs, we consider the statistic  $S = \frac{1}{n} \sum_{m=1}^n S_m$ , where  $n$  is the total number of markers. We verify by simulation that the normal approximation works well for assessing significance when data are from a genome screen. For each relationship examined, the null mean and variance,  $E_O(S)$  and  $Var_O(S)$ , are needed in order to apply the normal approximation. Let  $f_1, f_2, \dots, f_i$  be the allele frequencies at marker  $m$ . Then in addition to the null IBD probabilities  $p_i = P_O(D_m = i)$ ,  $i = 0, 1, 2$ , the calculation of  $E_O(S)$  and  $Var_O(S)$  requires the following probabilities, valid for outbred pairs:

$$\begin{aligned}
 P(S_m = 2|D_m = 2) &= 1, \\
 P(S_m = 2|D_m = 1) &= \sum_i f_i^2, \\
 P(S_m = 1|D_m = 1) &= 1 - \sum_i f_i^2, \\
 P(S_m = 2|D_m = 0) &= \sum_i f_i^4 + 2 \sum_i \sum_{j \neq i} f_i^2 f_j^2, \\
 P(S_m = 1|D_m = 0) &= 4 \sum_i \sum_{j \neq i} f_i^3 f_j \\
 &\quad + 4 \sum_i \sum_{j \neq i} \sum_{k \neq i, j} f_i^2 f_j f_k, \\
 P(S_{m1}, S_{m2} | D_{m1} = i, D_{m2} = j) \\
 &= P(S_{m1} | D_{m1} = i) P(S_{m2} | D_{m2} = j),
 \end{aligned}$$

for  $m1 \neq m2$ . Finally, the null probabilities  $P_O(D_{m2} = i2 | D_{m1} = i1)$  for  $i1, i2 \in \{0, 1, 2\}$  are also required, and are given in table 1 of Bishop and Williamson (1990) for the relationships sib, half-sib, parent-child, grand-PC, avuncular, and first cousin. More generally, they can be determined from the augmented Markov chain  $\{A\}$ , if it is known, or, in some cases, by the algorithms of Denniston (1975), Thompson (1988), and Tiwari and Elston (1999). Then we have

$$E_O(S) = \frac{1}{n} \sum_{m=1}^n \sum_{i=0}^2 \sum_{j=0}^2 j p_i P\{S_m = j | D_m = i\}$$

and

$$\begin{aligned}
 \text{Var}_O(S) &= \left[ \frac{1}{n^2} \sum_{m1=1}^n \sum_{m2=1}^n \sum_{i1=0}^2 \sum_{i2=0}^2 \sum_{j1=0}^2 \sum_{j2=0}^2 \right. \\
 &\quad j_1 j_2 p_{i1} P_O\{D_{m2} = i2 | D_{m1} = i1\} \\
 &\quad \times P\{S_{m1} = j_1 | D_{m1} = i1\} \\
 &\quad \left. \times (P\{S_{m2} = j_2 | D_{m2} = i2\})^{I\{m1 \neq m2\}} - E_O(S)^2 \right],
 \end{aligned}$$

where  $I\{m1 \neq m2\}$  is the indicator of the event  $\{m1 \neq m2\}$ .

*Test Based on Conditional Expected IBD (EIBD)*

Although the MLRT has high power, its drawbacks include the need to construct augmented Markov chains for the null and for each alternative relationship, the need to implement a separate hidden Markov chain calculation for each possible alternative relationship, and the need for simulation to assess significance. In principle it should be possible to perform simulations for each class of relatives in the data set, e.g. for avuncular pairs, and then use the same simulations to evaluate significance of MLLR for all avuncular pairs. However, it is often the case that not all individuals are typed at the same markers, so separate simulations may be needed for each pair of individuals. Thus, the LRT or MLRT

may be very cumbersome to use as a diagnostic tool unless the set of relationships to consider is very limited, as in the sib pair case described in Goring and Ott (1997) and Boehnke and Cox (1997), or the number of relative pairs in the data is small. The IBS-based test is much simpler computationally, but it loses power by not explicitly considering chance sharing. We propose two alternative test statistics designed to retain the simplicity of IBS, but to increase power by taking into account chance sharing.

The first test statistic, denoted EIBD, is the average of the conditional expected number of alleles shared IBD at each marker, conditional on the data for that marker, the null relationship, and the allele frequencies; that is,  $EIBD = \frac{1}{n} \sum_{m=1}^n E_O(D_m | G_m)$ , where  $D_m$  is the number of alleles shared IBD at marker  $m$ ,  $G_m$  is the genotype information for the pair at marker  $m$ ,  $n$  is the number of markers, and the subscript  $O$  indicates that the expectation is calculated under the null relationship. Note that the expectation at a locus is taken conditional only on the data for that locus. Letting  $p_i = P_O(D = i)$ ,  $i = 0, 1, 2$ , we have

$$E_O(D_m | G_m) = \frac{2P(G_m | D_m = 2)p_2 + P(G_m | D_m = 1)p_1}{\sum_{i=0,1,2} P(G_m | D_m = i)p_i},$$

where the probabilities  $P(G_m | D_m = i)$  are given in Thompson (1975) for outbred relative pairs.

We find that the normal distribution gives a close approximation to the sampling distribution of EIBD when applied to genome screen data. Thus, to assess significance, one need only calculate the mean and variance of the statistic under the null hypothesis. Note that  $E_O(EIBD) = E_O[E_O(D_m | G_m)] = E_O(D_m) = 2p_2 + p_1$ , the mean number of alleles shared IBD under the null relationship. The calculation of the null variance is very similar to that for IBS. We can think of  $E_O(D_m | G_m)$  as a function of  $G_m$ . Then to calculate the variance of EIBD, we need the probabilities  $P_O(D_{m2} = i2 | D_{m1} = i1)$  as in the IBS case and the probabilities  $P(G_m | D_m = i)$  instead of the probabilities  $P(S_m | D_m = i)$  used in the IBS case.

*Test Based on Adjusted IBS (AIBS)*

One possible drawback of the EIBD statistic is that if the null relationship has  $p_2 = 0$ , then  $E_O(D_m | G_m)$  is restricted to be within the range 0–1. This may give less than optimal power if the alternative relationship has moderate  $p_2$ . To avoid this problem, we also propose an adjusted IBS statistic, which is an average over all markers  $m$  of  $A_m$ , where  $A_m$  is a sum over each shared allele of its null conditional probability of being shared IBD given that the allele is shared IBS and under the assumption that the shared alleles result from a random draw of one allele from each of the individuals; that is,

AIBS =  $\frac{1}{n} \sum_{m=1}^n A_m$ , where  $A_m = 0$  if no alleles are shared IBS,  $A_m = \Phi_o / [\Phi_o + (1 - \Phi_o)f_i]$  if one allele is shared IBS and it is allele  $i$ , and  $A_m = \Phi_o / [\Phi_o + (1 - \Phi_o)f_i] + \Phi_o / [\Phi_o + (1 - \Phi_o)f_j]$  if two alleles are shared IBS and they are alleles  $i$  and  $j$ . Here  $\Phi_o$  is the kinship coefficient under the null relationship,  $\Phi_o = p_1/4 + p_2/2$ . If an allele were drawn at random from each individual's genotype at a given locus, the quantity  $\Phi_o / [\Phi_o + (1 - \Phi_o)f_i]$  would represent the probability that the two alleles are shared IBD given that they are shared IBS for allele  $i$ , conditional on the null relationship.

For AIBS, we find that the normal approximation is quite satisfactory for assessing significance. Calculation of the null mean and variance for AIBS is very similar to that for IBS and EIBD.  $A_m$  can be thought of as a function of  $G_m$ , so the probabilities  $P(D_{m2} = i2 | D_{m1} = i1)$  and  $P(G_m | D_m = i)$  are used to find the null mean and variance, as for EIBD.

#### Preliminary Screening with EIBD and AIBS, Followed by Application of MLRT

In a reasonably large data set, such as the Genetic Analysis Workshop (GAW) 11 Committee on the Genetics of Alcoholism (COGA) data set (Begleiter et al. 1999) analyzed in the *Application to GAW 11 COGA Data* subsection of the Results section, there may be thousands of relationship pairs to be tested. To use the MLRT, one would apply a computationally demanding calculation to each pair, with significance assessed by simulation in which that calculation is repeated  $10^5$  times. This large number of replicates would be needed when significance level .001 is used to reduce the number of false positive detections when screening a large number of pairs. If not all pairs are typed on the same markers, as in the GAW 11 COGA data set, then separate simulations would be performed for each pair, even among those that have the same null relationship. To make this more practical, we suggest a strategy of preliminary screening using EIBD and AIBS to determine a subset of pairs for which it is worthwhile to perform the MLRT. We establish a somewhat arbitrary cutoff of .2 for the smaller of the two  $p$  values obtained by use of EIBD and AIBS, in order for a relative pair to be chosen for the MLRT. The appropriate assessment of significance for this two-step procedure is easily obtained by simulation, at virtually no additional cost over the simulations that are required to assess significance for MLRT. In the Results section, we show that virtually no power is lost by doing this instead of applying the MLRT to all relative pairs, while the savings in computing time is substantial.

#### Relationship Estimation

If the null relationship is rejected in a hypothesis test, it is of interest to know what relationship is suggested by the data. One strategy would be to formulate the augmented Markov chain for each of a large number of relationships, calculate the likelihood under each, and compare them. However, in practice, the need to specify and implement an augmented Markov chain for every relationship considered would involve a substantial investment of time for each alternative considered. We propose a simpler preliminary strategy, involving estimation of  $\underline{p} = (p_0, p_1, p_2)$ , the probabilities of sharing zero, one, or two alleles IBD. This simple method could be used to suggest some likely alternative relationships, whose likelihoods could then be compared.

We estimate  $\underline{p}$  by maximizing  $\sum_{m=1}^n \log [L(G_m; \underline{p})]$ , where

$$L(G_m; \underline{p}) = p_0 P(G_m | D_m = 0) \\ + p_1 P(G_m | D_m = 1) + p_2 P(G_m | D_m = 2)$$

is the likelihood of the genotype data at marker  $m$  in terms of  $\underline{p}$  and the allele frequencies. If the markers were unlinked, this estimate of  $\underline{p}$  would be the maximum likelihood estimate derived in Thompson (1975). However, we apply this procedure here to linked markers. The quantity  $\sum_{m=1}^n \log [L(G_m; \underline{p})]$  can be quickly maximized by a simple application of the EM algorithm. Where the current estimate of  $\underline{p}$  is  $\underline{p}^{(k)} = [p_0^{(k)}, p_1^{(k)}, p_2^{(k)}]$ , we obtain the updated estimate by the formula

$$p_i^{(k+1)} = \frac{p_i^{(k)}}{n} \sum_{m=1}^n \frac{P(G_m | D_m = i)}{L[G_m; \underline{p}^{(k)}]} .$$

We investigate the properties of this estimator when applied to genome screen data below.

As shown in Thompson (1986), under the assumption of no inbreeding, the constraint  $p_1^2 \geq 4p_0p_2$  must be satisfied. The quantity  $\sum_{m=1}^n \log [L(G_m; \underline{p})]$  could be maximized subject to this constraint by first finding the maximizing  $\underline{p}$  in the unconstrained case. If the constraint is violated, then the condition  $p_1^2 = 4p_0p_2$  is imposed and a one-dimensional search algorithm is used. If the true relationship is not known, however, one may want to use the unconstrained estimate to allow for the possibility of inbreeding.

## Results

### Simulation Studies

We perform simulations to compare the power of the four test statistics, MLLR, EIBD, AIBS, and IBS, to detect misspecified pairwise relationships. We also include

in the comparison the LRT based on the correct alternative relationship, which sets a benchmark of close to optimal power that is not realistically achievable in practice when the correct alternative relationship is unknown. (Power of the LRT used here is slightly suboptimal because of the presence of interference, but we expect the effect to be almost negligible. See Appendix B for the extension of the LRT to the case of interference.) In our initial simulations, we consider the following five relationships: sib, half-sib, grand-PC, avuncular, and first cousin. For this set of simulations, we take the MLLR to be the maximum of the log-likelihood over the four possible alternative relationships from this set (excluding the null relationship), minus the null log-likelihood. We simulate marker data from an autosomal genome screen for which we vary the allele frequencies and marker resolution. Our simulated scenarios include panels of microsatellite markers equally spaced at recombination fractions of .07, .15, and .25, with sex-averaged chromosome lengths taken from Broman et al. (1998), and with all markers having allele frequencies .40, .20, .20, .05, .05, .05, and .05. We also simulate SNPs equally spaced at recombination fractions .01 and .07, with allele frequencies .7 and .3. The allele frequencies for these simulated SNP and microsatellite panels were chosen so that the markers would be somewhat less informative than the ideal, but within the range of what might be typical. Our results show that the conclusions about power comparisons across the statistics depend very little on the assumptions about the allele frequency distributions. Our final simulated marker panel is based on the markers actually typed in the GAW 11 COGA data set. This panel is more realistic because marker spacings are unequal with an average intermarker recombination fraction of .13, allele frequency distributions differ across markers, and some marker data are missing. We consider the power of the hypothesis tests at significance levels of .01 and .001. The significance level of .01 would be appropriate for a single hypothesis test, whereas we use the level of .001 in our screening of 2,810 relative pairs in the COGA data set in order to reduce the number of false positives that would be expected to occur in screening a large number of pairs. All simulations are performed by use of the chi-square model for crossovers with interference (Cobbs 1978; Stam 1979; Foss et al. 1993; McPeck and Speed 1995; Zhao et al. 1995) with parameter  $m = 4$  for humans, corresponding to a gamma shape parameter of 5, as suggested by the results of Lin and Speed (1996). Although it is convenient to assume no interference in the development and implementation of the testing methods, the actual data do contain interference. Thus, in order to give as close an indication as possible of the performance of the methods on real data, we simulate the data with interference. The number of replications

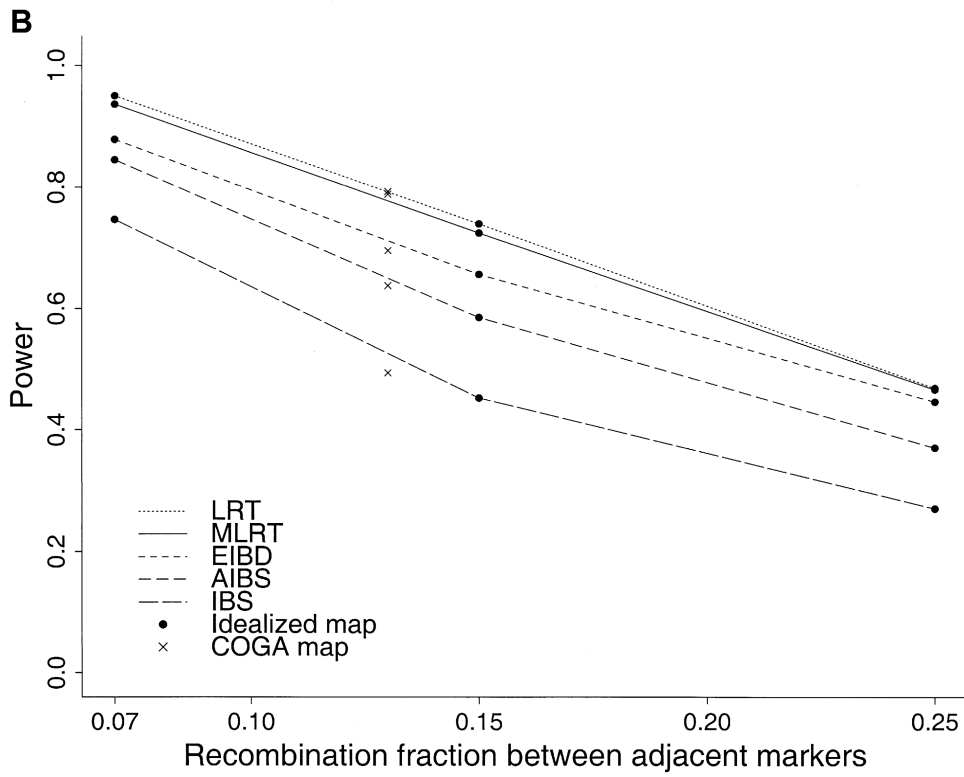
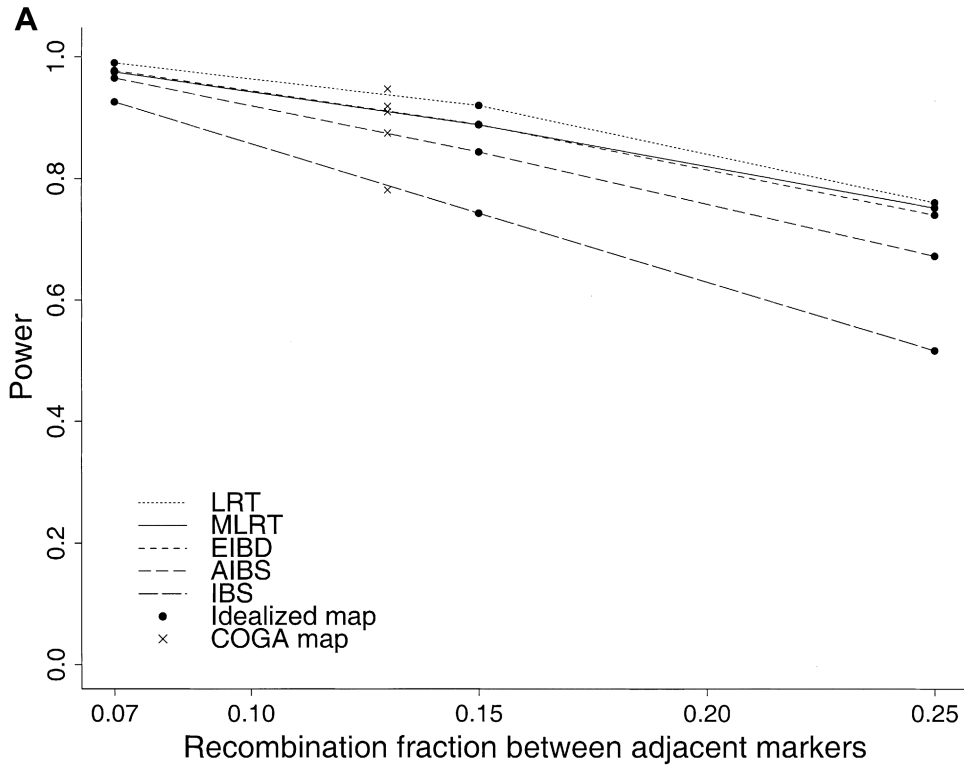
in each simulation is 1 million, and the five testing methods are analyzed on the same 1 million data sets in each case, minimizing any effects of sampling variability.

Figure 2 and tables 4–7 give the results of the power studies. In each case, we give power obtained by comparison of each of 1 million data sets simulated under the alternative relationship to an empirical null distribution obtained from 1 million data sets simulated under the null relationship. This procedure leads to very accurate power comparisons. However, in practice, for the tests based on EIBD, AIBS, and IBS, the normal approximation is adequate to assess significance. This does not hold for the MLRT, for which simulation is required to assess significance, but 100,000 simulated realizations may be adequate if significance level .001 is used or 10,000 simulated realizations if significance level .01 is used. For EIBD, AIBS, and IBS, we compared the  $p$  values calculated from the normal approximation to those calculated from the empirical null distribution and found them to be very close (results not shown). In our simulations, we also found that the distributions of these test statistics were approximately normal when the null relationship was not the true relationship.

Figure 2A gives the power for testing the null of half-sib against the alternative of first cousin at significance level .01 and figure 2B gives the same at significance level .001, both for microsatellite markers with frequencies .40, .20, .20, .05, .05, .05, and .05. Power results based on the GAW 11 COGA data are also shown. From the plot it is clear that the LRT has the highest power, followed by MLRT, EIBD, and AIBS, with IBS having substantially lower power than the others. Note that at significance level .01, MLRT and EIBD have roughly equal power to detect the first-cousin alternative, whereas at significance level .001 the power of MLRT is higher. Since the normal approximation does not hold for MLRT, we do not necessarily expect it to follow a similar pattern to the other statistics across the two plots. The results are similar to those in figure 2 if the null of first cousin is tested against the alternative of half-sib or avuncular (see table 4), or the null of avuncular is tested against the alternative of first cousin (results not shown). Table 4 also gives power for testing a null of first cousin against the alternative of grand-PC for significance levels .01 and .001. The ordering of the five statistics in terms of their power is the same as above.

The three relationships half-sib, avuncular, and grand-PC are similar in that they have the same probabilities of IBD. However, the transition rate of the grand-PC IBD process is only 1/2 that of half-sibs and 2/5 that of avuncular. (In each case, we define the transition rate to be the rate of transition to IBD value  $1 - i$  conditional on current IBD value  $i$  for the stationary IBD process.) The LRT and MLRT are the only ones among the five tests considered that take into account the information





**Figure 2** Power versus genome-screen resolution for the four test statistics with 1 million replications, microsatellite marker frequencies .40, .20, .20, .05, .05, .05, and .05. A, Null relationship half-sib, alternative relationship first-cousin, significance level .01. B, Null relationship half-sib, alternative relationship first-cousin, significance level .001.

**Table 4**  
**Power of Tests Based on LLR, MLLR, EIBD, AIBS, and IBS, against the Alternative of Half-Sib, Grand-PC, or Avuncular Relationship, When Null Relationship Is First Cousin.**

SIGNIFICANCE LEVEL, MARKER TYPE ( $\theta$ ), AND TEST STATISTIC <sup>a</sup>	POWER WHEN ALTERNATIVE (TRUE) RELATIONSHIP IS <sup>b</sup>		
	Half-Sib	Grand-PC	Avuncular
Significance level .01:			
SNP (.01):			
LRT	1.00	1.00	1.00
MLRT	1.00	1.00	1.00
EIBD	.99	.97	.99
AIBS	.97	.95	.98
IBS	.97	.95	.97
Microsatellite (.07):			
LRT	.99	1.00	.99
MLRT	.99	1.00	.99
EIBD	.98	.96	.98
AIBS	.96	.94	.97
IBS	.94	.91	.94
COGA map (average .13):			
LRT	.95	.98	.95
MLRT	.95	.97	.94
EIBD	.92	.89	.93
AIBS	.87	.85	.88
IBS	.81	.79	.81
Microsatellite (.15)			
LRT	.93	.96	.92
MLRT	.92	.95	.92
EIBD	.90	.88	.91
AIBS	.84	.82	.84
IBS	.74	.73	.75
Microsatellite (.25):			
LRT	.77	.81	.76
MLRT	.76	.80	.75
EIBD	.75	.73	.75
AIBS	.65	.64	.65
IBS	.52	.52	.53
SNP (.07):			
LRT	.69	.75	.68
MLRT	.69	.74	.67
EIBD	.65	.64	.65
AIBS	.45	.45	.45
IBS	.42	.42	.42
Significance level .001:			
SNP (.01):			
LRT	.99	1.00	.98
MLRT	.98	.99	.98
EIBD	.95	.91	.95
AIBS	.89	.85	.89
IBS	.88	.85	.89
Microsatellite (.07):			
LRT	.96	.99	.95
MLRT	.96	.98	.95
EIBD	.91	.87	.92
AIBS	.86	.82	.87
IBS	.79	.76	.79
COGA map (average .13):			
LRT	.82	.90	.81
MLRT	.82	.88	.80
EIBD	.75	.73	.76
AIBS	.65	.63	.65

IBS	.54	.54	.54
Microsatellite (.15):			
LRT	.77	.86	.75
MLRT	.76	.84	.74
EIBD	.71	.69	.72
AIBS	.59	.58	.59
IBS	.45	.46	.45
Microsatellite (.25):			
LRT	.49	.57	.48
MLRT	.48	.56	.47
EIBD	.47	.47	.47
AIBS	.35	.35	.34
IBS	.23	.24	.23
SNP (.07):			
LRT	.40	.48	.39
MLRT	.39	.47	.37
EIBD	.36	.37	.36
AIBS	.19	.20	.18
IBS	.18	.19	.18

<sup>a</sup> Microsatellite markers have allele frequencies .40, .20, .20, .05, .05, .05, and .05, SNPs have allele frequencies .7 and .3, and markers are equally spaced, with given recombination fraction  $\theta$  between adjacent pairs. GAW 11 COGA map has average marker spacing of 13.6 cM, which would correspond to  $\theta \approx .13$  when the Kosambi map function is used.

<sup>b</sup> Power is based on  $10^6$  simulated realizations.

in the data on the transition rate of the process. For each of the statistics EIBD, AIBS, and IBS, its mean value does not vary among the three relationships, although its variance does vary. Thus, these statistics have almost no power to distinguish among these three relationships. The LRT and MLRT have some power to distinguish among them based mainly on the different transition rates of the IBD processes. However, table 5 shows that the power of the LRT (which is always higher than that of MLRT) is generally very low in these cases. Note that the power is higher when one in the pair of relationships (null or alternative) is grand-PC than when neither relationship is grand-PC. This is explained by the fact that the transition rate for the grand-PC IBD process is very different from those for the avuncular and half-sib IBD processes.

We find that for EIBD, AIBS, and IBS, when the null relationship is first cousin, the grand-PC relationship tends to be slightly more difficult to detect as an alternative than are half-sib and avuncular (table 4) and grand-PC is more difficult to correctly reject as a null than are half-sib and avuncular relationships (results not shown). The fact that the transition rate is much smaller for the grand-PC IBD process than for the half-sib and avuncular IBD processes causes the variances of EIBD, AIBS, and IBS to be higher under the grand-PC model than under either the half-sib or avuncular model. This explains the increased difficulty in detecting grand-PC as an alternative or correctly rejecting it as a null for those statistics. Given two relationship pairs, the power

**Table 5**  
**Power of LRT to Distinguish Half-Sib, Grand-PC, and Avuncular Relationships, at Significance Level .001**

MARKER TYPE ( $\theta$ ) AND ALTERNATIVE (TRUE) RELATIONSHIP	POWER OF LRT WHEN NULL (FALSE) RELATIONSHIP IS <sup>a</sup>		
	Half-Sib	Grand-PC	Avuncular
SNP (.01):			
Half-sib	NA	.47	.01
Grand-PC	.40	NA	.73
Avuncular	.02	.78	NA
Microsatellite (.07):			
Half-sib	NA	.23	.01
Grand-PC	.20	NA	.41
Avuncular	.01	.47	NA
COGA map (average .13):			
Half-sib	NA	.06	.03
Grand-PC	.05	NA	.10
Avuncular	.00	.13	NA
Microsatellite (.15):			
Half-sib	NA	.04	.00
Grand-PC	.04	NA	.07
Avuncular	.00	.09	NA
Microsatellite (.25):			
Half-sib	NA	.01	.00
Grand-PC	.01	NA	.02
Avuncular	.00	.02	NA
SNP (.07):			
Half-sib	NA	.01	.00
Grand-PC	.01	NA	.01
Avuncular	.00	.02	NA

NOTE.—For details, see table 4.

<sup>a</sup> NA = not applicable.

to distinguish between them depends on the choice of the null relationship. For example, it is much easier to distinguish grand-PC (alternative) from first-cousin (null) than the other way around.

We find that the full sib relationship is relatively easy to distinguish from the other relationships, either as a null or as an alternative. Table 6A and B, respectively, show power when full sib is either the null or the alternative relationship, for significance level .001, for microsatellites with recombination fraction between adjacent markers of .25 and for SNPs with recombination fraction between adjacent markers of .07. For lower significance levels or increased marker density, power is nearly perfect for all methods when full sibs is either the null or the alternative, so the results are not shown. Note that when the alternative relationship is full sibs while the null relationship has  $p_2 = 0$ , EIBD performs worse than the other statistics. This is because the conditional expected number of alleles shared IBD can never be  $>1$  in that case. Even so, the power of EIBD in this case is very high,  $\geq 88\%$  for all cases simulated.

In all cases, the LRT is the most powerful of the five methods. Note that interference is present in the simulated data, whereas the likelihood ratio is calculated under the assumption of no interference for the sake of

computational simplicity. Thus, it is not the true likelihood-ratio test for the model used in the simulations, although, in practice, we expect this to have a negligible effect on power. Figure 2A and B and table 4 show that the power achieved by MLRT and EIBD is close to that of LRT, with MLRT usually having somewhat higher power than EIBD. When the true alternative relationship is not contained in the set  $\mathcal{A}$  of relationships considered for the MLRT, the power of EIBD and that of MLRT tend to be very close for the cases we have simulated (results not shown), including double first cousin, siblings from a first-cousin mating, and the relationship, which we call half-sib plus first cousin, that is depicted in figure 3, and which our analysis suggests may occur in the GAW 11 COGA data set. EIBD has higher power than both AIBS and IBS in all cases simulated except when the alternative relationship is full sibs while the

**Table 6**  
**Power, at Significance Level .001, of Tests Based on LLR, MLLR, EIBD, AIBS, and IBS, against (A) the Alternative of Full Sib, When Null Relationship Is Half-Sib, Grand-PC, Avuncular, or First Cousin, and (B) the Alternative of Half-Sib, Grand-PC, Avuncular, or First Cousin, When Null Relationship Is Full Sib**

MARKER TYPE ( $\theta$ ) AND TEST STATISTIC	POWER OF TEST			
	Half-Sib	Grand-PC	Avuncular	First Cousin
A. Alternative of Full Sib, When Null Relationship Is Half-Sib, Grand-PC, Avuncular, or First Cousin				
Microsatellite (.25):				
LRT	1.00	1.00	1.00	1.00
MLRT	.99	.99	1.00	1.00
EIBD	.92	.88	.93	1.00
AIBS	1.00	.99	1.00	1.00
IBS	1.00	.99	1.00	1.00
SNP (.07):				
LRT	1.00	1.00	1.00	1.00
MLRT	1.00	1.00	1.00	1.00
EIBD	.95	.92	.96	1.00
AIBS	.99	.99	1.00	1.00
IBS	.99	.99	1.00	1.00
B. Alternative of Half-Sib, Grand-PC, Avuncular, or First Cousin, When Null Relationship Is Full Sib				
Microsatellite (.25):				
LRT	1.00	1.00	1.00	1.00
MLRT	1.00	1.00	1.00	1.00
EIBD	1.00	1.00	1.00	1.00
AIBS	1.00	1.00	1.00	1.00
IBS	1.00	.99	1.00	1.00
SNP (.07):				
LRT	1.00	1.00	1.00	1.00
MLRT	1.00	1.00	1.00	1.00
EIBD	1.00	1.00	1.00	1.00
AIBS	1.00	.99	1.00	1.00
IBS	.99	.99	.99	1.00

NOTE.—For details, see table 4.

**Table 7**

**Power of Suggested Strategy of Application of MLRT to Those Pairs for Which at Least One of EIBD and AIBS Has  $p < .2$ , Divided by Power of MLRT, for Significance Level .001 in Both Cases**

Null Relationship	True Relationship	Power Relative to MLRT
First cousin	Half-sib	1.000
First cousin	Grand-PC	1.000
First cousin	Avuncular	.999
Half-sib	First cousin	1.000
Grand-PC	First cousin	.999
Avuncular	First cousin	1.000
Sibling	All others	1.000
All others	Sibling	1.000

NOTE.—The number of simulated realizations is  $10^6$ . The GAW 11 COGA map is used, with average intermarker distance of 13.6 cM.

null relationship has  $p_2 = 0$ . In that case, however, all five statistics have very high power. The power of AIBS is always at least as high as, and sometimes substantially higher than, IBS, as shown in figure 2 and tables 4 and 6.

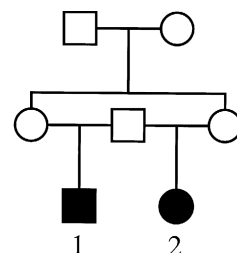
The type of map used, SNP map, microsatellite map, or GAW 11 COGA map, did not have a substantial impact on the power comparisons among the statistics. The power results for SNPs with allele frequencies .7 and .3 are similar to results for microsatellites at a lower density. Using SNPs at recombination fraction .01 generally gives a test with slightly more power than that using microsatellites at recombination fraction .07, whereas using SNPs at recombination fraction .07 generally gives a test with slightly less power than that using microsatellites at recombination fraction .25. When SNPs with allele frequencies .5 and .5 are used, power is slightly higher, but the increase is fairly small (results not shown). The results for the GAW 11 COGA map, which has average intermarker recombination fraction of .13 and different allele frequency distributions across markers, were quite similar to those for the idealized microsatellite map with intermarker recombination fraction of .15.

In practice, for a data set in which there may be thousands of pairwise relationships to consider, we have suggested a strategy of preliminary screening with EIBD and AIBS, using a  $p \leq .2$  on at least one of these tests as a prerequisite for applying the MLRT. The rationale is that simulation to assess significance for the MLRT is computationally intensive, so could be reserved for those pairs that are most likely to yield an unusual result. We find that significance levels, obtained by simulation, for this two-step procedure are not detectably different from significance levels, obtained by simulation, for MLRT, even when the number of simulations used is

$10^6$ ; that is, under the null relationship, the chance that a relative pair will have  $p \geq .2$  for both AIBS and EIBD conditional on having  $p < .001$  for MLRT is so small that the resulting difference in simulated significance levels for the two procedures (straight MLRT vs. screening followed by MLRT) is smaller than the amount of chance variation in simulated significance level for either procedure, even when the number of realizations is  $10^6$ . Table 7 shows that virtually no power is lost by use of this strategy compared with use of the MLRT on all pairs, at least for the relationships we have considered.

To use the MLRT for a wider range of pairwise relationships, one would need to construct an augmented Markov chain for each relationship considered as either a null or an alternative. Another approach is the Markov approximation suggested in the Methods section. For the cases of avuncular and first-cousin relationships, we have performed simulations to compare the likelihood calculated by use of the augmented Markov chain (correct likelihood) with the approximate likelihood obtained by use of the Markov approximation. Our simulations consisted of 100,000 replicates of the given relationship, with genotypes simulated on the basis of the GAW 11 COGA map. For the avuncular relationship, out of 100,000 simulations, the maximum relative error of the likelihood for genome screen data when the Markov approximation was used was .000124. For the first-cousin relationship, out of 100,000 simulations, the maximum relative error of the likelihood for genome screen data when the Markov approximation was used was .000168. When the Markov approximation was used to calculate the likelihoods for the LRT, with first cousin as the alternative and avuncular as the null, the power was very close to that when the correct likelihood was used, similarly with avuncular as the null and first cousin as the alternative (results not shown). This suggests that, at least for these two relationships, the Markov approximation to the likelihood is adequate.

We performed simulations to assess the estimation of relationships, using the method described in Methods above. The results are shown in table 8. We give the mean and standard deviation of the estimated IBD shar-



**Figure 3** Half-sib plus first-cousin pedigree: individuals 1 and 2 are half-sibs through their father and first cousins through their mother.

**Table 8**

**Relationship Estimation, Where  $p_i$  Is the Probability of  $i$  Alleles Shared IBD, for  $i = 0, 1, 2$**

MARKER TYPE ( $\theta$ ) AND RELATIONSHIP <sup>a</sup>	MEAN (SD) OF ESTIMATED IBD SHARING PROBABILITY <sup>b</sup>		
	$p_0$	$p_1$	$p_2$
SNP (.01):			
Full sib	.251 (.045)	.499 (.050)	.250 (.042)
Half-sib	.504 (.057)	.488 (.060)	.008 (.011)
Grand-PC	.505 (.068)	.488 (.070)	.008 (.011)
Avuncular	.505 (.055)	.488 (.057)	.007 (.011)
First cousin	.755 (.051)	.238 (.054)	.007 (.010)
Microsatellite (.07):			
Full sib	.249 (.047)	.500 (.052)	.251 (.044)
Half-sib	.502 (.060)	.490 (.061)	.007 (.011)
Grand-PC	.502 (.070)	.491 (.071)	.007 (.011)
Avuncular	.503 (.057)	.490 (.058)	.007 (.011)
First cousin	.752 (.054)	.242 (.056)	.006 (.009)
COGA map (average .13):			
Full sib	.249 (.053)	.501 (.062)	.250 (.048)
Half-sib	.503 (.068)	.487 (.071)	.009 (.014)
Grand-PC	.504 (.077)	.487 (.079)	.009 (.014)
Avuncular	.503 (.066)	.487 (.069)	.009 (.014)
First cousin	.754 (.065)	.238 (.067)	.008 (.012)
Microsatellite (.15):			
Full sib	.250 (.055)	.500 (.064)	.250 (.048)
Half-sib	.503 (.070)	.486 (.073)	.010 (.016)
Grand-PC	.503 (.079)	.486 (.081)	.010 (.016)
Avuncular	.504 (.068)	.485 (.071)	.011 (.016)
First cousin	.756 (.069)	.235 (.071)	.009 (.014)
Microsatellite (.25):			
Full sib	.250 (.065)	.500 (.080)	.251 (.055)
Half-sib	.503 (.084)	.483 (.088)	.013 (.021)
Grand-PC	.504 (.090)	.482 (.094)	.014 (.021)
Avuncular	.505 (.083)	.481 (.087)	.014 (.021)
First cousin	.758 (.082)	.231 (.086)	.012 (.018)
SNP (.07):			
Full sib	.250 (.074)	.500 (.097)	.250 (.059)
Half-sib	.510 (.093)	.470 (.102)	.019 (.029)
Grand-PC	.513 (.101)	.468 (.110)	.020 (.029)
Avuncular	.513 (.093)	.468 (.102)	.019 (.028)
First cousin	.766 (.094)	.215 (.103)	.019 (.027)

<sup>a</sup> Microsatellite markers have allele frequencies .40, .20, .20, .05, .05, .05, and .05; SNPs have allele frequencies .7 and .3, and markers are equally spaced, with given recombination fraction  $\theta$  between adjacent pairs.

<sup>b</sup> Each mean and SD is based on  $10^4$  simulated realizations.

ing probabilities when the data are simulated under various relationships, where the number of replicates in each simulation is 10,000. For the relationships with  $p_2 = 0$ , there is a slight bias in the estimates, amounting to no more than ~5%. This bias is expected because one can estimate  $p_2$  only at or above its actual value, never below, in those cases. For microsatellite markers at recombination fraction .07, the bias is quite small. The standard deviations of the estimates tend to be rather large at the marker resolutions considered. Thus, the procedure may give only a rough idea of the true relationship. For instance, double first cousins has  $p_2 =$

$1/16$ ,  $p_1 = 3/8$ ,  $p_2 = 9/16$ , and quadruple half first cousins has  $p_2 = 1/32$ ,  $p_1 = 7/16$ ,  $p_0 = 17/32$ . With the level of standard deviation seen in the simulations, such relationships may be difficult to distinguish from each other and from half-sib/grand-PC/avuncular ( $p_2 = 0$ ,  $p_1 = 1/2$ ,  $p_0 = 1/2$ ), by use of the estimation procedure. We explored the use of the maximized value of  $\sum_{m=1}^n \log [L(G_m; \underline{p})]$  as a test statistic to detect misspecified relationships, but its power was lower than AIBS, perhaps because there is not enough information in these data to accurately estimate relationships by this method. Still, this quick approximate method can be used to suggest a class of relationships whose likelihoods could then be compared.

#### Application to GAW 11 COGA Data

The GAW 11 COGA data were collected for the purpose of mapping genes for susceptibility to alcohol dependence and related phenotypes (Begleiter et al. 1999). The data consist of 105 pedigrees, generally 3- or 4-generation, with 1,214 individuals, 992 of whom are genotyped. The genome screen includes 296 markers with average heterozygosity of .73 and average intermarker distance of 13.6 cM. In our analysis, we consider only the autosomal markers, of which there are 285 with average intermarker distance of 13.5 cM. Allele frequencies, estimated with the USER M13 program (Boehnke 1991), were distributed with the data, as were marker order and distances estimated with the CRIMAP program (Lander and Green 1987). In the present analysis for detection of misspecified relationships, we consider only five null relationships: full sib, half sib, grand-PC, avuncular and first cousin. Among the typed individuals, we use the pedigree information to identify 2,810 such relative pairs. The majority (2,625 pairs) have >200 typed markers in common. The minimum number of shared typed markers for any of the 2,810 pairs is 25.

For each pair, performance of the MLRT, including 100,000 simulated realizations to assess significance at the .001 level, takes ~4 min on a Sun Ultra II with 360-MHz processor and .5 GB RAM. Thus, for 2,810 pairs, the approximate time to perform the MLRT would be 7.8 days. We instead prescreen, using AIBS and EIBD, performing the MLRT only for those pairs for which at least one of the AIBS and EIBD-based tests has a  $p < .2$ . This results in a time of ~2.6 d to complete the calculations, a savings of ~ $\frac{2}{3}$ , or >5 days of CPU time.

For each of the 2,810 pairs, we calculate EIBD and AIBS and perform the corresponding hypothesis tests for relationship misspecification. We identify 949 of the 2,810 pairs that have  $p < .2$  for either the EIBD or AIBS test. For each of these 949 pairs, we calculate the MLLR statistic where the set of alternative relationships  $\mathcal{A}$ , over

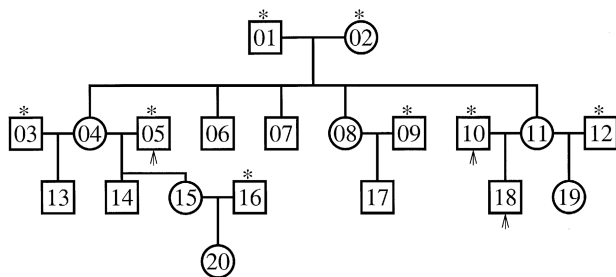
which the statistic is maximized, consists of full sib, half sib, grand-PC, avuncular, first cousin, and unrelated. To assess significance of deviation from the null relationship, for a particular pair among the 949, we simulate 100,000 realizations of the genotype data for that pair under the null relationship, with the same markers typed as in the data for that pair. For each realization, we can perform the two-step procedure of first screening with EIBD or AIBS and then calculating MLLR for those pairs having  $p < .2$  for EIBD or AIBS. In this data set, we find that the significance level obtained this way is virtually indistinguishable from that obtained for MLRT alone. Among the 949 pairs, there are 26 that are significant at level .001, a level at which we would expect approximately three false positives if all tests were independent, and these 26 pairs occur in 11 pedigrees.

In pedigree data, when one individual has, say, mis-specified paternity, or if there is a switched sample, this may be expected to create multiple incorrect pairwise relationships, some of which may be observable, depending on how many close relatives have genotype data available. When a particular pair is observed to have significant deviation from its null relationship, the pattern among the other relative pairs in the pedigree can often confirm and point to a likely explanation for the finding. We consider some examples in the GAW 11 COGA data. For instance, there is a sibship of size 4 in which one particular sib has significant relationship misfit with each of his three putative sibs. There is a similar sibship of size 3 in which a particular sib has relationship misfit with each of his two putative sibs. There is another sibship of size 3 in which this occurs and furthermore, the particular sib who shows misfit also shows significant relationship misfit with each of his three nieces and nephews. For the first and third of these families, there are no genotype data available on either of the parents, whereas for the second family, there are genotype data available on the mother only. In the second family, the possibility that the individual is a half-sib to each of his other putative sibs is compatible with the data. For the families with no parental genotype data available, for each pairing of the problematic individual with his sibs, the relationship maximizing the likelihood, among those in  $\mathcal{A}$ , is first cousin. In these two cases, the relationship of half-sibs is not compatible with the data, nor is the relationship of unrelated. There is an additional family containing a sibship of size 4 in which one sib has significant relationship misspecification with two of his three putative sibs, with father's genotype missing. In this case, the possibility that the one sib is actually a half-sib to all three of his putative siblings is compatible with the data. There is a sib pair that shows significant relationship misspecification, where the father is typed, but not the mother, and there are no other siblings. Neither sib shows any relationship misfit with any of his

avuncular or first-cousin relatives. Among the relationships in  $\mathcal{A}$ , the one that fits best for this pair is half-sib; however, for rejection of half-sib as a null,  $p = .04$ , and the estimated values of  $p_0$ ,  $p_1$ , and  $p_2$  suggest the possibility of a closer relationship such as half-sib plus first cousin as shown in figure 3. There is a sibship of size 9 among which there are four pairs showing significant misfit. These pairs appear to be much more closely related than sibs, but are clearly not MZ twins, suggesting an inbred relationship. In another family, there is a sib pair with significant deviation that also appears to be much more closely related than sibs, but less than MZ twins, again suggesting inbreeding.

In some cases, when a particular pair is observed to have significant deviation from its null relationship, the pattern among other relatives does not confirm this. For instance, there is a family containing a sibship of size 4 in which one sib has significant relationship misspecification with two of this three putative sibs, with father's genotype missing. In this case, the possibility that the one sib is a half-sib to each of his three putative sibs is strongly rejected by the data, and it is difficult to come up with a consistent explanation for the results. The one sib appears to be a half-sib to two of the three putative sibs, and a full sib to the other, but the remaining three all appear to be full sibs, an impossibility. In another family, there is a cousin pair showing significant relationship misspecification. Both cousins' parents are typed, and both have several typed siblings, all with data at  $>170$  markers, but there are no other significant pairs. With so many individuals typed, Mendelian errors would arise under most of the possible alternative scenarios. The GAW 11 COGA data have already been cleaned of Mendelian errors. One might presume that if the level of Mendelian errors had been above that expected from genotyping error, then the problematic individuals' genotypes would have been removed from the data set at that point. However, one should ideally have full information on the Mendelian errors detected in order to better assess relationship misspecification. Given that there are no other significant pairs in the pedigree, and under the assumption that most relatives are typed and none show excessive Mendelian errors, then it is reasonable to speculate that such a finding may be a false positive. There are two other similar cases, one involving an uncle-niece pair and one involving a grand-PC pair.

Perhaps the most interesting cases, from the point of view of the methods developed here, are those in which the apparently misspecified relationships could not have been detected based only on analysis of sibs and half-sibs, but require consideration of other relationships such as avuncular and first cousin. Consider the family shown in figure 4. Here we have removed extraneous individuals from the pedigree and changed the sexes of



**Figure 4** A pedigree from the GAW 11 COGA data set, with some extraneous individuals removed and sexes of some individuals changed. The starred individuals are untyped, and all other individuals are typed for  $\geq 250$  markers. Arrows indicate individuals discussed in the text.

some individuals to provide an extra level of confidentiality for the family. In the pedigree, individuals 1, 2, 3, 5, 9, 10, 12, and 16 are untyped, whereas all other individuals are typed for  $\geq 250$  markers. In this pedigree, significant relationship misfit was detected for the cousin pair 18 and 14 and the cousin pair 18 and 15. For all four testing methods, EIBD, AIBS, IBS, and MLRT,  $p < 10^{-5}$  for rejection of the first-cousin relationship for both of these pairs. The value of  $(p_0, p_1, p_2)$  for the first-cousin relationship is  $(.75, .25, 0)$ , but the estimated values for these two pairs are  $(.28, .56, .16)$  and  $(.27, .57, .16)$ , respectively, which is between half and full sibs. There is no excess sharing between individuals 4 and 11, who are the mothers of the cousins. There is no misfit between individual 18 and his half-sib 19, or between individuals 14, 15, or 18 and their avuncular relatives. One possible explanation consistent with the data is that individuals 5 and 10 are the same person. In that case, individuals 18 and 14 would have the relationship shown in figure 3, as would individuals 18 and 15, whereas all other relationships in figure 4 would be preserved. In another case, which would also not be detected by use of full or half sibs, a pair of putative first cousins shows significant relationship misfit, with the relationship of double first cousins actually being consistent with the data. Neither father is typed, and the cousin relationship is through the mothers.

## Discussion

In order to extend the likelihood calculations of Goring and Ott (1997) and Boehnke and Cox (1997) to more general pairwise relationships for which IBD status is no longer a Markov chain, we define an augmented Markov chain that contains the information of IBD status. A generally applicable way to define such an augmented Markov chain is to consider the process whose state space is the equivalence relation on inheritance vectors

obtained by identifying inheritance vectors that differ only by interchanges of maternal and paternal haplotypes within founders, although this will not give the most parsimonious chain. For the avuncular and first-cousin cases, we give augmented Markov chains containing the minimal amount of information needed beyond IBD status in order to make the process Markov. Similar chains for some additional relationships are given by Donnelly (1983). Using these augmented Markov chains as the basis of our likelihood calculations, we describe and implement a MLRT, for which significance is assessed by simulation. Extensions of the likelihood calculation that take into account interference and inbreeding are described in Appendices A and B, respectively.

We extend the IBS-based test of Ehm and Wagner (1998) to more general relative pairs. The implementation of the IBS-based test is much simpler computationally than that of MLRT, but the IBS-based test loses power by not taking into account chance sharing. The tests based on EIBD and AIBS are a compromise between the two, with higher power than IBS, but with the desirable features of the IBS-based test, namely computational simplicity and no need to specify a particular alternative relationship. In simulations, EIBD has power close to that of the MLRT and has higher power than both the AIBS and IBS-based tests, except in the case when the null relationship has  $p_2 = 0$  and the alternative is full sib. However, the power of EIBD is still quite high in this case, generally over 90% in our simulations. The statistic AIBS outperforms IBS in every simulation, but in most simulations does not perform as well as EIBD.

Among the four statistics considered, MLLR, EIBD, AIBS, and IBS, the test based on MLLR has the highest power to detect misspecified relationships. However, although the tests based on EIBD, AIBS, and IBS do not require specification of any alternative relationship, the MLRT requires one to construct augmented Markov chains for a number of alternative relationships as well as for the null relationship. Furthermore, assessment of significance for the MLRT requires computationally intensive simulation, whereas a normal approximation can be used for EIBD, AIBS, and IBS. When a large number of relative pairs are to be considered, as in the GAW 11 COGA data, we recommend preliminary screening using EIBD and AIBS, with MLRT applied only to pairs having  $p \leq .2$  for at least one of EIBD and AIBS. Our results show that the significance level is hardly detectably changed, and the reduction in power is minuscule. In the GAW 11 COGA data set, this resulted in the MLRT being performed for 1/3 of the pairs considered, which took  $\sim 2.6$  d of computer time on a SUN Ultra II with 360-MHz processor and .5 GB RAM,

instead of the 7.8 days it would be expected to take to perform MLRT for all pairs.

To assess significance for the EIBD, AIBS, or IBS-based tests using the normal approximation, the null mean and variance must be computed. In all three cases, this requires calculation of the two-locus null conditional IBD probabilities  $P_O(D_{m2} = i2 | D_{m1} = i1)$  for  $i1, i2 \in \{0, 1, 2\}$  in addition to the single-locus null IBD probabilities  $p_0, p_1,$  and  $p_2$  and allele frequencies at each locus. In particular, map information in the form of all pairwise recombination fractions between markers and allele frequencies at all loci are used to assess significance. For the MLRT, recombination fractions between all pairs of adjacent markers and allele frequencies are used both to compute the statistic and for simulation to assess significance.

In most of our simulations, we have assumed that allele frequencies were known, whereas in practice, one would generally need to estimate these from the data. It is of interest to know how robust the various tests are to misspecification of allele frequencies. Allele frequencies are involved in calculation of all statistics besides IBS, but even for that statistic, allele frequencies are needed in order to assess significance. Our preliminary simulation studies suggest that MLRT and EIBD are comparably robust to misspecified allele frequencies, with AIBS and IBS apparently somewhat less robust (results not shown).

Another aspect of the simulations that is somewhat unrealistic is that there is no genotyping error. With genome screen data, a low rate of genotyping error would not be expected to have serious consequences for relationship error detection by use of the MLRT, except in the cases of parent-offspring and MZ twin relationships, in which genotyping errors might result in an outcome that has likelihood zero under that relationship. In contrast, the other statistics, EIBD, AIBS, and IBS would be expected to be more robust. In the cases of parent-offspring and MZ twin relationships, the likelihood approach could be slightly modified to include a low rate of random genotyping error as suggested by Broman and Weber (1998). The GAW 11 COGA data set was apparently already cleaned of Mendelian errors before being distributed. Ideally, detection of misspecified relationships and identification of Mendelian errors should be performed simultaneously. Although some Mendelian errors will occur because of genotyping errors, their presence or absence and overall level can provide important clues to the detection of relationship errors and to the understanding of their likely cause.

One way to make the application of MLRT feasible for a wider class of relationships is to use the Markov approximation to the likelihood proposed in the Methods section. In that case, rather than construct an augmented Markov chain for each null and alternative re-

lationship considered, one need only calculate the one-step conditional IBD probabilities  $\Pr(D_{m2} = j | D_{m1} = i)$ , where  $D_m$  is the number of alleles shared IBD by a given pair at locus  $m$ . Our results indicate that at least for the avuncular and first-cousin relationships, this approximation is adequate for relationship testing.

In the GAW 11 COGA data set, a number of pairwise relationships showing significant misfit were detected, including some that could not have been detected had only full and half-sibs been considered. To reduce the occurrence of false positive detection of relationship error among the 2,810 relative pairs in the data set, we set the significance level to .001. Our simulations indicate that power of the MLRT is still quite high at this significance level. When an apparent error is detected, one may be able to distinguish a true relationship error from chance rejection of the null by consideration of the pattern of results among multiple pairs from the same pedigree. In the case of a true error, there may be a pattern of results among close relatives all pointing to a particular alternative explanation for the data. Ideally, the methods presented here should be extended to simultaneous inference on a number of relatives. For instance, an entire sibship could be considered simultaneously in a single likelihood analysis, rather than separate consideration of relative pairs.

When the null relationship is rejected by a hypothesis test, it is natural to consider the problem of estimation of the correct relationship. We describe a simple approach to estimation of  $p_0, p_1,$  and  $p_2$ , the null probabilities of IBD. The estimates are rough, but can be used to suggest candidate relationships whose likelihoods could then be calculated and compared.

For the test statistics EIBD, AIBS, and IBS, one could consider trying to increase power by weighting markers differently depending on their location in the genome, giving isolated markers more weight than those in densely mapped regions, because correlation between markers is a decreasing function of distance. Our preliminary work on optimal weights indicates that, at least in the complete data case, the increase in power tends to be small (results not shown).

## Acknowledgments

We would like to thank Dr. Theodore Reich and Dr. Henri Begleiter of the Collaborative Study on the Genetics of Alcoholism and Dr. Jean MacCluer of the Genetics Analysis Workshop for permission to use the GAW11 COGA data. We would like to thank Dr. Nancy Cox for helpful discussions. This work is supported by National Institutes of Health/National Human Genome Research Institute grant HG01645. The GAW is supported by National Institutes of Health grant GM31575.



## Appendix A

### More Detailed Explanation of the Violation of the Markov Property by the Avuncular and First-Cousin IBD Processes

To see that the avuncular and first-cousin IBD processes are not Markov, it would suffice to provide a counterexample in each case. However, in order to understand the augmented Markov processes we introduce, it is necessary to understand the nature of the violation of the Markov property, which we now describe in more detail.

First consider the avuncular case. Suppose the genotypes at locus  $A$  are as given in figure 1A, where the maternally inherited allele of individual 4,  $a_i$ , is either  $a_3$  or  $a_4$ . Here, the avuncular pair of individuals 3 and 6 share no alleles IBD at locus  $A$ . Consider a nearby locus  $B$  to the right of locus  $A$ , linked to  $A$ . We are interested in the distribution of the number of alleles shared IBD by individuals 3 and 6 at locus  $B$  conditional on the genotype information at locus  $A$  given in figure 1A. We make the relatively weak assumption that the crossover process is a regular, stationary point process, with chiasma interference permitted, but with no chromatid interference, i.e. the choices of chromatid strands for different crossovers are independent and uniform. In the example shown in figure 1A, if  $a_i$  is  $a_3$ , i.e. if the allele not transmitted from 4 to 6 is shared IBD between individuals 3 and 4, then IBD sharing of 1 for the avuncular pair at locus  $B$  may be achieved by a single crossover in any one of three meioses: the meioses involving transmission of genetic material from individual 1 to individual 3, 1 to 4, or 4 to 6. If  $a_i$  is  $a_4$ , i.e. if the allele not transmitted from 4 to 6 is not shared IBD between individuals 3 and 4, then IBD sharing of 1 for the avuncular pair at locus  $B$  may be achieved by a single crossover in either of two meioses: 1 to 3, or 1 to 4. Thus, the instantaneous rate of transition at  $A$  from IBD 0 to IBD 1 for the avuncular pair is 3 if  $a_i = a_3$  and 2 if  $a_i = a_4$ . This suggests that in the example shown in figure 1A, if we condition on all the allele information at locus  $A$ , then the distribution of the number of alleles shared IBD by the avuncular pair at  $B$  depends on  $a_i$ , i.e. it depends on whether the  $A$  allele not transmitted from individual 4 to individual 6 is shared IBD between individuals 3 and 4. In fact, letting  $D_A$  and  $D_B$  be the number of alleles shared IBD by individuals 3 and 6 at loci  $A$  and  $B$ , respectively, letting  $S_A$  denote the event that the  $A$  allele not transmitted from individual 4 to individual 6 is shared IBD between individuals 3 and 4, and letting  $S_A^c$  denote the complementary event to  $S_A$ , we have that  $P(D_B = 1|D_A = j, S_A) > P(D_B = 1|D_A = j, S_A^c)$  when  $A$  and  $B$  are linked. This inequality can be deduced

from the transition probabilities given in table 2A, in light of the fact that  $2\theta(1 - \theta) < \theta^2 + (1 - \theta)^2$  for  $0 \leq \theta < .5$ . Immediate consequences of this inequality are (i)  $P(D_B = 1|D_A = j, S_A) > P(D_B = 1|D_A = j)$ , and (ii)  $P(D_B = 0|D_A = j, S_A) < P(D_B = 0|D_A = j, S_A^c)$ . From (ii), we have (iii)  $P(D_B = 0|D_A = j) > P(D_B = 0|D_A = j, S_A)$ . Using (i) and (iii), we have that

$$\begin{aligned} P(S_A|D_B = 1, D_A = j) &= P(S_A, D_B = 1|D_A = j) / \\ P(D_B = 1|D_A = j) &= P(D_B = 1|D_A = j, S_A) \\ &\quad \times P(S_A|D_A = j) / P(D_B = 1|D_A = j) \\ &> P(D_B = 0|D_A = j, S_A) P(S_A|D_A = j) / \\ P(D_B = 0|D_A = j) &= P(S_A|D_B = 0, D_A = j) , \end{aligned}$$

which implies (iv)  $P(S_A|D_A = j, D_B = 1) > P(S_A|D_A = j)$ . Thus, conditional on the number of alleles shared IBD by the avuncular pair at locus  $A$ , the probability that  $S_A$  occurs is increased if the IBD sharing by the avuncular pair is 1 at a nearby locus  $B$ . Suppose  $B$  is to the right of  $A$ , and let  $C$  be another nearby locus to the left of  $A$ . Then, using (iv) and the fact that  $(S, D)$  is Markov, we have that

$$\begin{aligned} P(D_C = 1|D_A = j, D_B = 1) &= P(D_C = 1|D_A = j, S_A, D_B = 1) P(S_A|D_A = j, D_B = 1) \\ &\quad + P(D_C = 1|D_A = j, S_A^c, D_B = 1) P(S_A^c|D_A = j, D_B = 1) \\ &= P(D_C = 1|D_A = j, S_A) P(S_A|D_A = j, D_B = 1) \\ &\quad + P(D_C = 1|D_A = j, S_A^c) [1 - P(S_A|D_A = j, D_B = 1)] \\ &= P(D_C = 1|D_A = j, S_A^c) + [P(D_C = 1|D_A = j, S_A) \\ &\quad - P(D_C = 1|D_A = j, S_A^c)] P(S_A|D_A = j, D_B = 1) \\ &> P(D_C = 1|D_A = j, S_A^c) + [P(D_C = 1|D_A = j, S_A) \\ &\quad - P(D_C = 1|D_A = j, S_A^c)] P(S_A|D_A = j) \\ &= P(D_C = 1|D_A = j) . \end{aligned}$$

This shows inequality 1, in violation of the Markov property. The first-cousin case follows by use of a generalization of this argument.

## Appendix B

### Calculation of the Likelihood Under the $\chi^2$ Interference Model

The chi-square interference model can be viewed as a hidden Markov model. By applying the Baum algorithm,

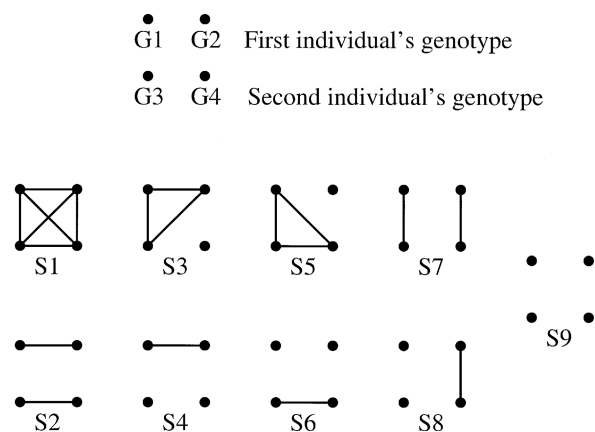
the likelihood can be obtained. To see this, note that for the chi-square model with parameter  $m$ , the crossover process on four strands can be obtained by first constructing a Poisson process with rate  $2(m + 1)$  in terms of genetic distance. Start at one end of the chromosome and label the first point of the Poisson process  $X_1$  where  $X_1$  is chosen uniformly at random from among the integers  $\{0, 1, 2, \dots, m\}$ . For all  $i > 1$ , label the  $i$ th point of the process (counting in order from the end of the chromosome)  $X_i = X_1 + i \pmod{m + 1}$ . Then every point with label 0 is a crossover point for the four-strand process. To obtain the single-strand crossover process, independently keep each point of the four-strand crossover process with chance  $1/2$  or eliminate it with chance  $1/2$ . Consider a single strand inherited by an offspring from its parent. For a given chromosomal location  $t$ , define  $Z(t) = [X(t), Y(t)]$ , where  $X(t) = X_i$  for  $C_i \leq t < C_{i+1}$ ,  $C_i$  being the  $i$ th point of the original Poisson  $(2(m + 1))$  process, and  $Y(t) = 1$  if the offspring inherited the parent's paternal DNA and 0 if the offspring inherited the parent's maternal DNA at location  $t$ . Then  $Z(t) = [X(t), Y(t)]$  is a Markov process with  $P[\text{Next state is } (x_2, y_2) | \text{Current state is } (x_1, y_1)] = P_{(x_1, y_1), (x_2, y_2)}$ , where  $P_{(i,j), (i+1, j)} = 1$  for  $i = 0, 1, \dots, m - 1, j = 0, 1, P_{(m, j), (0, k)} = 1/2$  for  $j = 0, 1, k = 0, 1$ , and all other entries are 0. The leaving rate for each state is  $2(m + 1)$ . The observed data give partial information on  $Y(t)$  only. Thus, the Markov chain  $Z(t) = [X(t), Y(t)]$  is hidden. Now consider a pair of individuals in a pedigree. Define such a hidden Markov chain for each meiosis in the pedigree, and consider the product Markov chain  $[Z_1(t), Z_2(t), \dots, Z_n(t)] = [X_1(t), Y_1(t), X_2(t), Y_2(t), \dots, X_n(t), Y_n(t)]$ , where  $Z_i(t) = [X_i(t), Y_i(t)]$  is the Markov chain for the  $i$ th meiosis in the pedigree, and there are  $n$  meioses in total. Since the meioses are independent, the transition matrix is the  $n$ -fold Kronecker product of the transition matrix for a single meiosis. If genotype data for the pair of individuals is observed, then in principle the Baum algorithm can be applied to the product Markov chain to calculate the likelihood. As in Kruglyak et al. (1996), a reduction in dimensionality could be achieved by identifying states that differ by one or more interchanges of founders' paternally and maternally inherited haplotypes.

### Appendix C

#### Inbred Relative Pair

In the case of an inbred relative pair, instead of three IBD states (0, 1, or 2 alleles shared IBD), there are now nine IBD states (Jacquard 1974). Let the (unordered) genotype of individual 1 be  $(G_1, G_2)$  and let the (unordered) genotype of individual 2 be  $(G_3, G_4)$ . Jacquard (1974) depicts each of the nine IBD states by a graph

with four nodes, each representing one of  $G_1, G_2, G_3, G_4$ , with an edge present between  $G_i$  and  $G_j$  if and only if  $G_i$  and  $G_j$  are IBD. See figure C1. States  $S_7, S_8$ , and  $S_9$  correspond to outbred states of 2, 1, and 0 alleles shared IBD, respectively. The other six IBD states involve inbreeding in one or both individuals. Thompson (1975) gives the distribution of genotype given IBD state for an outbred relative pair, which is needed in order to use each of the four testing methods described in the current work. However, this is not sufficient to determine the distribution of genotype given IBD state for inbred relative pairs. The relevant distributions for inbred relative pairs appear in Jacquard (1974). In principle, an augmented Markov chain could be derived for an inbred relative pair, for instance the chain  $\{A''\}$  described in Methods above could be used. Then the likelihood of the data could be computed by use of the distribution of genotype given IBD state that appears in Jacquard (1974). This would allow construction of a likelihood-ratio test for an inbred relative pair. In order to extend the definition of EIBD to inbreds, there is more than one reasonable approach. One could define the of number of alleles shared IBD for inbred relative pairs by defining states  $S_1, S_2, \dots, S_9$  to have 4, 0, 2, 0, 2, 0, 2, 1, and 0 alleles shared IBD, respectively. Alternatively, one might prefer to define states  $S_1, S_2, \dots, S_9$  to have 2, 0, 1, 0, 1, 0, 2, 1, and 0 alleles shared IBD, respectively, as for outbreds. To extend the definition of IBS to inbreds, one could think of 9 IBS states analogous to the nine IBD states. Again, to apply the IBS score statistic to inbreds, one could define the number of alleles shared IBS for inbred relative pairs by defining IBS states  $S_1, S_2, \dots, S_9$  to have 4, 0, 2, 0, 2, 0, 2, 1, and 0 alleles shared IBS, respectively. However, if the level of inbreeding is low, one may have more power by defining IBS



**Figure C1** Graphic representation of the nine IBD states of Jacquard (1974). An edge is present between a given pair of nodes if and only if they are identical by descent.

states  $S_1, S_2, \dots, S_9$  to have 2, 0, 1, 0, 1, 0, 2, 1, and 0 alleles shared IBS, respectively, as for outbreeds.

## References

- Baum LE (1972) An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes. *Inequalities* 3:1–8
- Begleiter H, Reich T, Nurnberger J Jr, Li TK, Conneally PM, Edenberg H, Crowe R, et al (1999) Description of the Genetic Analysis Workshop 11 Collaborative Study on the Genetics of Alcoholism. *Genet Epidemiol* 17 Suppl 1:S25–S30
- Bishop DT, Williamson JA (1990) The power of identity-by-state methods for linkage analysis. *Am J Hum Genet* 46:254–265
- Boehnke M (1991) Allele frequency estimation from data on relatives. *Am J Hum Genet* 48:22–25
- Boehnke M, Cox NJ (1997) Accurate inference of relationships in sib-pair linkage studies. *Am J Hum Genet* 61:423–429
- Broman KW, Murray JC, Sheffield VC, White RL, Weber JL (1998) Comprehensive human genetic maps: individual and sex-specific variation in recombination. *Am J Hum Genet* 63:861–869
- Broman KW, Weber JL (1998) Estimation of pairwise relationships in the presence of genotyping errors. *Am J Hum Genet* 63:1563–1564
- Cobbs G (1978) Renewal process approach to the theory of genetic linkage: case of no chromatid interference. *Genetics* 89:563–581
- Denniston C (1975) Probability and genetic relationship: two loci. *Ann Hum Genet* 39:89–104
- Donnelly KP (1983) The probability that related individuals share some section of genome identical by descent. *Theor Popul Biol* 23:34–63
- Ehm MG, Wagner M (1998) A test statistic to detect errors in sib-pair relationships. *Am J Hum Genet* 62:181–188
- Feingold E (1993) Markov processes for modeling and analyzing a new genetic mapping method. *J Appl Prob* 30:766–779
- Foss E, Lande R, Stahl FW, Steinberg CM (1993) Chiasma interference as a function of genetic distance. *Genetics* 133:681–691
- Göring HHH, Ott J (1997) Relationship estimation in affected sib pair analysis of late-onset diseases. *Eur J Hum Genet* 5:69–77
- Jacquard A (1974) *The genetic structure of populations*. Springer-Verlag, New York
- Kruglyak L, Daly MJ, Reeve-Daly MP, Lander ES (1996) Parametric and non-parametric linkage analysis: a unified multipoint approach. *Am J Hum Genet* 58:1347–1363
- Lander, Green (1987) Construction of multilocus genetic linkage maps in humans. *Proc Natl Acad Sci USA* 84:2363–2367
- Lin SL, Speed TP (1996) Incorporating crossover interference into pedigree analysis using the chi(2) model. *Hum Hered* 46:315–322
- McPeck MS, Speed TP (1995) Modeling interference in genetic recombination. *Genetics* 139:1031–1044
- Stam P (1979) Interference in genetic crossing over and chromosome mapping. *Genetics* 92:573–594
- Thompson EA (1975) The estimation of pairwise relationships. *Ann Hum Genet* 39:173–188
- (1986) *Pedigree analysis in human genetics*. The Johns Hopkins University Press, Baltimore
- (1988) Two-locus and three-locus gene identity by descent in pedigrees. *IMA J Math Appl Med Biol* 5:261–280
- Tiwari H, Elston R (1999) A new explicit algorithm to calculate identity by descent probabilities for pairs of relatives with respect to two linked loci. *Genet Epidemiol* 17:216
- Zhao H, Speed TP, McPeck MS (1995) Statistical analysis of crossover interference using the chi-square model. *Genetics* 139:1045–1056