# Variance Components Models for Gene–Environment Interaction in Quantitative Trait Locus Linkage Analysis

Shaun Purcell and Pak Sham

*Social, Genetic and Developmental Psychiatry Research Centre, Institute of Psychiatry, King's College, London, UK*

Gene–environment interaction (G × E) is likely to be a common and important source of variation for complex behavioral traits. Gene–environment interaction, or genetic control of sensitivity to the environment, can be incorporated into variance components twin and sib-pair analyses by partitioning genetic effects into a mean part, which is independent of the environment, and a part that is a linear function of the environment. An approach described in a companion paper (Purcell, 2002) is applied to sib-pair variance components linkage analysis in two ways: allowing for quantitative trait locus by environment interaction and utilizing information on any residual interactions detected prior to analysis. As well as elucidating environmental pathways, consideration of G × E in quantitative and molecular studies will potentially direct and enhance gene-mapping efforts.

Gene–environment interaction $(G \times E)$ is most tractable when dealing with measured (as opposed to latent) genetic and environmental effects. In a companion paper dealing with latent $G \times$ measured $E$ interaction (Purcell, 2002), an approach is outlined with reference to the classical twin study. In this paper the same approach is extended to quantitative trait locus (QTL) sib-pair linkage analysis, within a variance components framework. Although analysis of $G \times E$ should eventually lead to a better understanding of the etiology of complex traits and diseases, in the context of linkage analysis (which only identifies fairly large genomic regions likely to harbor disease-causing genes) the current goal is simply to increase power to detect genes of small effect, rather than dissecting genetic-environmental architectures *per se*.

The variance components approach to sib-pair QTL linkage analysis is, in essence, only a trivial extension of the twin model (Amos, 1994; Fulker et al., 1999; Kruglyak & Lander, 1995). Assuming we have only full sibling pairs, the likelihood is parameterized in terms of three variance components: variance due to the QTL, $Q$, variance due to shared sibling effects, $S$, and variance due to nonshared sibling effects, $N$. Polygenic additive effects load onto both $S$ and $N$. The basic allele-sharing test of linkage is of the relationship between phenotypic sib-pair similarity and IBD sharing at the test locus. A "weighted-likelihood conditioning-on-trait-values" approach (Sham et al., 2000) is adopted in the following analyses, in order to provide a robust test of linkage in selected samples.
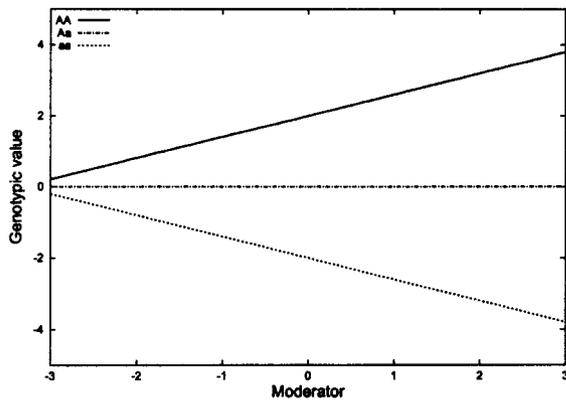
Two ways in which $G \times E$ can feature in QTL linkage analysis are considered: (1) when the actual QTL effect is moderated by a measured covariate $(Q \times M)$, and (2) when a residual variance component is moderated by a measured covariate (i.e., $S \times M$ and $N \times M$).

## $Q \times M$ in Linkage Analysis

Analogous to the modelling of a moderating effect on the additive genetic path, $a$, in the twin model (Purcell, 2002), the QTL path $q$ is simply modified to $(q+\beta_Q M)$ or even $(q+\beta_Q M+\delta_Q M^2)$ to incorporate $Q \times M$ interaction, representing linear and nonlinear interactions, respectively, between the additive genetic value at the QTL, $a_Q$, and the moderator. We assume that the presence or absence of a particular allele is unrelated to the moderator (i.e., no gene–environment correlation).

The simulations reported in Table 1 show four conditions varying in (1) QTL effect $(a_Q > 0)$ (2) $Q \times M$ interaction $(\beta_Q > 0)$ and (3) residual $G \times E$ interaction $(A \times M$ in fact, i.e., $\beta_X > 0)$. For each condition 200 replicate datasets were simulated, and a number of likelihood ratio test statistics were constructed. The base model $SN$ has no QTL effects; model $Q – SN$ allows for a simple QTL effect; model $Q – SN – X_Q$ allows for a moderator-linked QTL effect as well as a main effect. Two additional models also allow for the possibility of interaction effects between the residual variance components ($S$ and $N$) and the measured moderator variable $M$. From left to right, the three likelihood ratio tests shown in Table 1 therefore represent (1) a simple 1 degree of freedom test for an additive QTL effect (2) a 2 degree of freedom test for a QTL effect that might be moderated by the variable $M$ and (3) as for the previous test, but allowing for $S \times M$ and $N \times M$ effects under both the super- and submodel. In all cases, 1000 DZ twin pairs were simulated, with residual variance components $a = c = e = 1$ and an additive biallelic QTL with equal allele frequencies. The expected variance components associated with the $Q \times M$ corresponding to the fourth row of Table 1 are illustrated in Figure 1.

**Figure 1**

Example of a simulated $Q \times E$ interaction $\beta_Q = 0.3$ with additive genetic value $a_Q = 2$ and dominance deviation $d_Q = 0$.

**Table 1**

QTL Linkage Incorporating Q x M Interaction in DZ Twin Pairs, with and Without $A \times M$ Interaction Also

| | | | Likelihood ratio tests | | |
|---|---|---|---|---|---|
| | | | $Q - SN$ | $Q - SN - X_Q$ | $Q - SN - X_Q - X_S X_N$ |
| Simulated | | | | | |
| $a_Q$ | $\beta_Q$ | $\beta_X$ | SN | SN | $SN - X_S X_N$ |
| 0 | . | . | 0.62 | 1.93 | 1.77 |
| 0 | . | 0.2 | 0.58 | 13.74 | 1.81 |
| 2 | . | . | 48.10 | 48.27 | 48.67 |
| 2 | 0.3 | . | 45.13 | 106.87 | 49.66 |

The first two rows of Table 1 represent the case of no QTL effect. In the first row (no QTL effect, no interactions) the test statistics are all close to their expected values under the null. The second row (no QTL effect, residual interaction) shows that the combined test of a *moderated* QTL effect (second column) is highly anti-conservative in the presence of residual $A \times M$, with a highly significant $\chi_2^2$ = 13.74 (expected $\chi^2$ is 1.5, as the test of $Q$ only involves 0.5 degree of freedom). This bias is due to the greater variance at higher levels of the moderator (due to the residual interaction) which the $\beta_Q$ parameter attempts to account for. Properly modelling this residual non-additivity (i.e., by the inclusion of $S \times M$ and $N \times M$ terms, as in the third column) reduces this bias. Therefore, it is unwise to perform a simple $Q \times M$ type of analysis when conducting a linkage test when there is non-additivity in the data.

The next two rows of Table 1 represent the case of a large QTL effect ($a_Q = 2$). In the third row (QTL effect, no interactions) the likelihood ratio tests are all similar (although the first has one less degree of freedom). In the fourth row (QTL effect, QTL interacts with moderator), the "robust" linkage test (third column) gives very little extra information compared to a simple QTL test (first column). This is because the $S \times M$ and $N \times M$ components will soak up the $Q \times M$ effect (i.e., the opposite of the above effect). If one were able to be sure that there were

no significant residual interaction effects, however, then the basic test of a moderated QTL effect (second column) would in fact provide more power under the alternate.

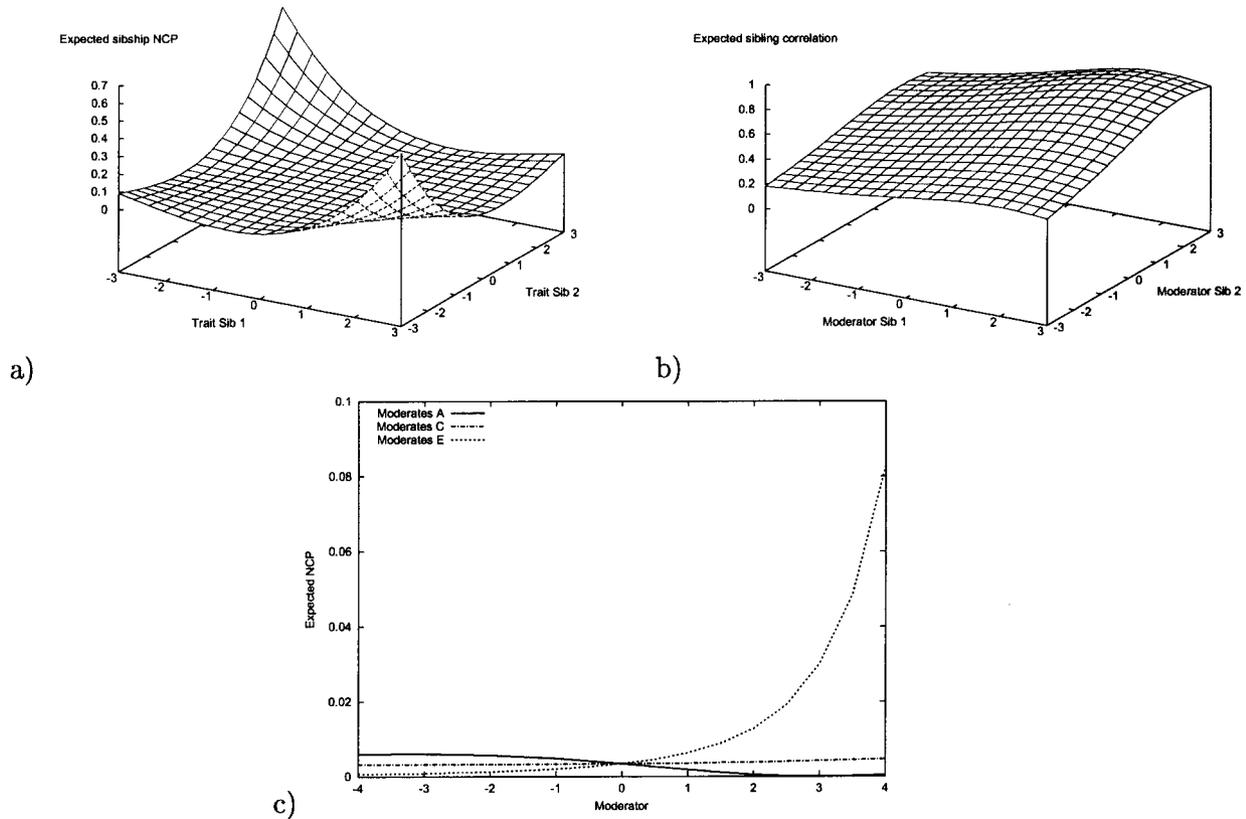## $N \times M$ and QTL Linkage in Selected Samples

The use of selective sampling schemes for linkage is highly desirable, especially when working with sibling pairs, as most pairs will yield very little information for linkage. To facilitate this, an index of potential informativeness for each pair can be calculated, conditional on their observed trait scores, the QTL allele frequency and effect size, and the residual sibling (Purcell et al., 2001; Sham & Purcell, 2001). Irrespective of QTL effect size, a higher residual correlation increases power to detect a QTL (Sham et al., 2000).

Typically, a single value for the sample residual correlation is specified when selecting or analyzing a sample for linkage. However, in the presence of $G \times E$ there will, by definition, be heterogeneity in the residual correlation across the sample. This section explores the possibility of using prior knowledge of such heterogeneity (when the relevant moderating variables have also been measured in the linkage sample) to better specify pair-specific residual correlations in order to increase power.

A correlation is a property of a number of paired observations: specifying a pair-specific correlation implies that the pair belongs to a particular subset which has that correlation. If a moderator variable $M$ interacts with either $A$, $C$ or $E$ components, then $M$ can predict which pairs will have higher residual sibling correlations. Consider, for example, an $E \times M$ interaction such that individuals scoring higher on $M$ will tend to have lower effects of $E$. In this case, pairs in which both members score high on $M$ will have a higher residual correlation. All other things being equal, it would therefore be preferable to select this pair over a pair with a lower residual correlation.

Figure 2 illustrates the relationship between sample selection and $G \times E$ in three graphs. Panel a) illustrates the relationship between trait score and expected informativeness: concordant high and low pairs and discordant pairs in particular are most informative. Panel a) assumes a constant sibling correlation across the sample however, which might not be the case. Panel b) illustrates how the residual sibling correlation might change as a function of a moderator variable, in the presence of an $E \times M$ interaction similar to that described above. It would therefore be desirable to take this information into account when selecting and analyzing pairs for linkage: panel c) shows the marked impact on the expected non-centrality parameter (via the expected residual correlation) for the linkage test in the presence of $G \times E$. The graph shows the expected non-centrality parameter (NCP) per randomly-selected sibpair as a function of sib pair moderator (assuming, in this case, that the moderator is identical between sibs and that the main effect of the moderator has been partialled out of the trait). In particular, modelling $E \times M$ interaction can greatly increase power — it seems that residual $A \times M$ and $C \times M$ do not influence the test so much (as they operate on both sibling variance and covariance, and so have less impact on the correlation).

It is interesting to note that these results are related to an observation regarding bivariate linkage analysis and the

a)

b)

c)

**Figure 2**

G × E and sample selection for linkage. See text for explanation.

source of residual cross-trait phenotypic covariance: that power increases dramatically with decreasing nonshared sources of covariance (Evans, 2002). In this sense, bivariate analysis and including a moderator variable can have a similar effect: the impact on the NCP of modelling $E \times M$, as shown above in panel c), seems to reflect a similar trend to that shown in Figure 2 of Evans (2002).

Focusing on $E \times M$, we assume that prior twin analyses have estimated a significantly nonzero value for $\beta_Z$. For a phenotyped sample of pairs also measured on $M$, this prior knowledge can be used (1) to select sibling pairs which are most informative for linkage, by calculating the residual correlation applicable to that pair conditional on measured $M$ and (2) in analysis, to use the pair-specific residual correlations. Ideally, the sample in which $\beta_Z$ was estimated will be as close as possible to the linkage sample (for example, the linkage sample could be all the DZ pairs from the twin sample). Effects of misspecifying $\beta_Z$ are explored below in the simulations.

Using a method based on the Haseman-Elston linkage test (Haseman & Elston, 1972; Sham & Purcell, 2001), the expected noncentrality parameter (NCP) for pair $i$ is

$$\frac{q^4}{16} \left[ \frac{(T_{i1} + T_{i2})^2}{(1 + r)^2} - \frac{(T_{i1} - T_{i2})^2}{(1 - r)^2} + \frac{4r}{1 - r^2} \right]^2$$

assuming complete marker informativeness, where $T_{i1}$ and $T_{i2}$ are the standardized trait scores for the pair, $r$ is the sibling correlation, and $q^2$ is the proportion of variance due to the QTL. This index can be used to rank order sibling pairs by potential informativeness. In the presence of heterogeneity, it is possible to calculate pair-specific correlations, which will more accurately model the residual variance in the sample. For pair $i$, conditional on estimated values of $a^2$, $c^2$, $e^2$ and $\beta_Z$ and measured $M_{i1}$ and $M_{i2}$, then $r_i$ can be calculated as

$$r_i = \frac{0.5 \times a^2 + c^2}{\sqrt{a^2 + c^2 + (e + \beta_Z M_{i1})^2} \sqrt{a^2 + c^2 + (e + \beta_Z M_{i2})^2}}$$

which can be substituted in the above expression. The trait score for sib $j$ of pair $i$, $T_{ij}$, also has to be standardized to unit variance conditional on the moderator. In the case of $a^2$, $c^2$, $e^2$ and $\beta_Z$ having been previously estimated

$$T'_{ij} = T_{ij} / \sqrt{a^2 + c^2 + (e + \beta_Z M_{ij})^2}$$

although the expressions for the moderator-conditional standardized scores and correlations will change depending on which models are being used to give the prior parameter estimates. Sibships, not twins, may only have been available, for example.

The formulation of the linkage model used here (Sham et al., 2000) has only a single free parameter, the QTL variance, $q^2$. The total variance and residual correlation are fixed, either to their sample values or other values estimated in previous studies (e.g., in the case of a selected sample). In the present case, the variance is fixed to unity and the residual correlation fixed to the pair-specific values, conditional on the moderator. The covariance matrices conditional on IBD sharing at the test locus are therefore

$$\begin{bmatrix} 1 & r_i - q^2/2 \\ r_i - q^2/2 & 1 \end{bmatrix}, \begin{bmatrix} 1 & r_i \\ r_i & 1 \end{bmatrix} \text{ and } \begin{bmatrix} 1 & r_i + q^2/2 \\ r_i + q^2/2 & 1 \end{bmatrix}$$

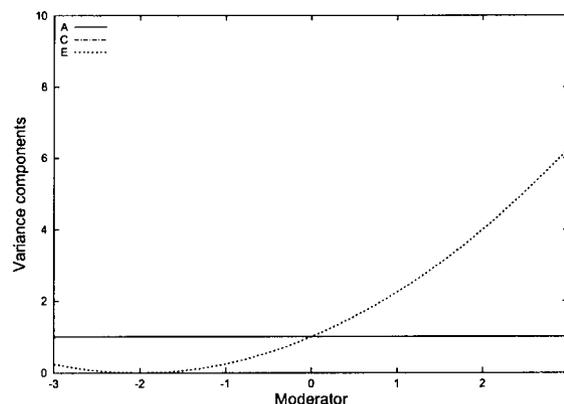for pairs sharing 0, 1 and 2 alleles IBD, respectively.

### Simulations

Simulations based on sib-pair datasets featuring a residual $E \times M$ interaction in all cases were conducted under a number of conditions: varying QTL effect, sample selection scheme and whether or not the residual interaction was included or misspecified in the analysis model (Table 2). Under each condition a dataset of 5000 DZ pairs was simulated 200 times. Selected sample analyses were based on the most informative 10% (i.e., 500 pairs). The QTL effect was specified in terms of the additive genetic value, $a_Q$, which was 0, 0.5 or 1, for a fully informative biallelic test locus with equifrequent alleles. Three final conditions concerned the residual interaction, which was simulated as $\beta_Z = 0.5$ in all cases (illustrated in Figure 3). In the first case, "w/ $E \times M$", the correct moderator variable was incorporated into the analysis with the correct estimate of $\beta_Z$ to form the pair-specific residual correlations used in selection and analysis. In the second condition, "w/ out $E \times M$", both selection and analysis were performed as usual, ignoring the moderator $M$. In the third condition, the true moderator was replaced with an unrelated random variable (i.e., which would have no moderating properties with respect to the trait) but $\beta_Z$ was still assumed to be 0.5,

representing a misspecification of the moderating effect in selection and analysis.

Under the null of no QTL effect ($a_Q = 0$), all models show average test statistics close the expected value (0.5), whether or not the moderator was included or misspecified and whether or not the analysis was performed on the whole or a selected sample. The $\hat{q}$ column gives a standardized estimate of the QTL variance, which are all close to zero under the null. For the selected and unselected samples, Table 2 also gives the % of replicates (out of 200) significant at various significance levels, which are all close to expected values.

Under the alternate hypothesis (i.e., $a_Q > 0$) it is clear that selected samples are more efficient than unselected samples (e.g., for $a_Q = 1$, in the condition not incorporating the moderator, on average 54% (11.40/21.30 = 0.535) of the signal was recovered by 10% of the sample). Incorporating the moderator results in a considerable gain in information. In terms of the average test statistic, for $a_Q = 1$ in unselected samples, there is a gain of 50%



**Figure 3**

$E \times M$ interaction with residual components $a = c = e = 1$ and $\beta_Z = 0.5$, as used in all simulations.

**Table 2**

Results of QTL Linkage Simulations Incorporating $E \times M$ Interaction

| | 10% most informative | | | | | | Unselected | | | | |
| | | | % significant at $p =$ | | | | | | % significant at $p =$ | | |
| $a_Q$ | $\hat{q}$ | LRT | 0.025 | 0.005 | 0.0005 | | $\hat{q}$ | LRT | 0.025 | 0.005 | 0.0005 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| w/ $E \times M$ | | | | | | | | | | | |
| 0 | 0.012 | 0.53 | 2 | 1 | 0 | | 0.010 | 0.56 | 2.5 | 1 | 0.5 |
| 0.5 | 0.047 | 2.73 | 27 | 9.5 | 2.5 | | 0.042 | 3.21 | 30 | 14 | 4.5 |
| 1 | 0.180 | 20.18 | 100 | 98 | 88 | | 0.184 | 31.82 | 100 | 100 | 99 |
| w/out $E \times M$ | | | | | | | | | | | |
| 0 | 0.011 | 0.43 | 1.5 | 0 | 0 | | 0.010 | 0.53 | 2.5 | 1 | 0 |
| 0.5 | 0.028 | 1.58 | 14 | 3.5 | 1 | | 0.033 | 2.35 | 24 | 9 | 1.5 |
| 1 | 0.101 | 11.40 | 90.5 | 75.5 | 45 | | 0.121 | 21.30 | 99.5 | 95.5 | 91.5 |
| w/ incorrect $E \times M$ | | | | | | | | | | | |
| 0 | 0.006 | 0.57 | 3 | 0.5 | 0 | | 0.005 | 0.43 | 2 | 0 | 0 |
| 0.5 | 0.015 | 1.29 | 9.5 | 2.5 | 0 | | 0.016 | 1.70 | 15.5 | 5 | 1 |
| 1 | 0.079 | 9.41 | 80.5 | 62 | 34 | | 0.100 | 17.53 | 97 | 91 | 78.5 |

(i.e., comparing "w/ $E \times M$" and "w/ out $E \times M$" conditions, (31.82-21.30)/21.30). For $a_Q = 1$ in selected samples, there is a gain of 77% (20.18-11.40)/11.40. In terms of the percentage significant with this sample size at a particular significance level, the gains can be great; for example, 88% are significant for $a_Q = 1$ at $p = 0.0005$ when the moderator is included compared to only 45% when it is not.

The "w/ incorrect $E \times M$" rows represent the scenario where the moderator is actually completely unrelated to the trait (i.e., the estimate of $\beta_Z$ obtained from another dataset is completely unwarranted in this one). As can be seen, this does reduce power to some extent, although the test still appears to have the correct performance under the null. In the case of $a_Q = 1$ the average test statistic drops by approximately 18% for both selected and unselected samples, the majority of the signal remains intact despite the complete misspecification.

If there is strong reason to believe that the moderating effect does exist in the linkage sample, then both selecting and analyzing incorporating the moderator seems desirable. If the effect is less certain, then it might not be advisable to select on the basis of the moderator, although it would be of interest to conduct the analyses both with and without incorporation of the putative moderator.

## Discussion

This paper has shown some of the potential gains and loses involved with modelling $G \times E$ in the QTL linkage analyses. In particular, it was shown (1) that a simple approach to $Q \times M$ interaction can lead to false positives and (2) that the benefits of bivariate linkage can potentially be harnessed in a $G \times E$ framework with a moderator that interacts with the residual nonshared component, whether or not the second trait is influenced by the QTL. Further work on gene-environment interaction will hopefully increase the chance of detecting QTL using linkage in small sibships.

### Software

Scripts to perform the above analyses using Mx (Neale, 1997) can be found at *http://statgen.iop.kcl.ac.uk/gxe/*.

## References

Amos, C. I. (1994). Robust variance-components approach for assessing genetic linkage in pedigrees. *American Journal of Human Genetics, 54,* 535–543.

Evans, D. (2002). The power of multivariate quantitative-trait loci linkage analysis is influenced by the correlation between variables. *American Journal of Human Genetics, 70,* 1599–1602.

Fulker, D., Cherney, S., Sham, P., & Hewitt, J. (1999). Combined linkage and association sib-pair analysis for quantitative traits. *American Journal of Human Genetics, 64*(1), 259–267.

Haseman, J. K., & Elston, R. C. (1972). The investigation of linkage between a quantitative trait and a marker locus. *Behavior Genetics, 2,* 3–19.

Kurglyak, L., & Lander, E. S. (1995). Complete multipoint sib-pair analysis of qualitative and quantitative traits. *American Journal of Human Genetics, 57,* 439–454.

Neale, M. (1997). *Mx: Statistical modeling.* Box 980126 VCU, Richmond VA 23298.

Purcell, S. (2002). Variance components models for gene–environment interaction in twin analysis. *Twin Research, 5,* 554–571.

Purcell, S., Cherny, S., Hewitt, J., & Sham, P. (2001). Optimal sibship selection for genotyping in QTL linkage analysis. *Human Heredity, 52,* 1–13.

Sham, P., Cherny, S., Purcell, S., & Hewitt, J. (2000). Power of linkage versus association analyis of quantitative traits, by use of variance-components models, for sibship data. *American Journal of Human Genetics, 66,* 1616–1630.

Sham, P., & Purcell, S. (2001). Equivalence between Haseman-Elston and variance components linkage analysis for sib pairs. *American Journal of Human Genetics, 68,* 1527–1532.

Sham, P. C., Zhao, J. H., Cherny, S., & Hewitt, J. (2000). Variance components QTL linkage analysis: conditioning on trait values. *Genetic Epidemiology, 19*(S1), S22–S28.