

Genetics and population analysis

## Mapping the genetic architecture of complex traits in experimental populations

Jian Yang<sup>1</sup>, Jun Zhu<sup>1,\*</sup> and Robert W. Williams<sup>2</sup>

<sup>1</sup>Institute of Bioinformatics, College of Agriculture and Biotechnology, Zhejiang University, Hangzhou, 310029, P. R. China and <sup>2</sup>University of Tennessee Health Science Center, Memphis, Tennessee, 38163, USA

Received on November 2, 2006; revised and accepted on April 7, 2007

Advance Access publication April 25, 2007

Associate Editor: Martin Bishop

### ABSTRACT

**Summary:** Understanding how interactions among set of genes affect diverse phenotypes is having a greater impact on biomedical research, agriculture and evolutionary biology. Mapping and characterizing the isolated effects of single quantitative trait locus (QTL) is a first step, but we also need to assemble networks of QTLs and define non-additive interactions (epistasis) together with a host of potential environmental modulators. In this article, we present a full-QTL model with which to explore the genetic architecture of complex trait in multiple environments. Our model includes the effects of multiple QTLs, epistasis, QTL-by-environment interactions and epistasis-by-environment interactions. A new mapping strategy, including marker interval selection, detection of marker interval interactions and genome scans, is used to evaluate putative locations of multiple QTLs and their interactions. All the mapping procedures are performed in the framework of mixed linear model that are flexible to model environmental factors regardless of fix or random effects being assumed. An *F*-statistic based on Henderson method III is used for hypothesis tests. This method is less computationally greedy than corresponding likelihood ratio test. In each of the mapping procedures, permutation testing is exploited to control for genome-wide false positive rate, and model selection is used to reduce ghost peaks in *F*-statistic profile. Parameters of the full-QTL model are estimated using a Bayesian method via Gibbs sampling. Monte Carlo simulations help define the reliability and efficiency of the method. Two real-world phenotypes (BXD mouse olfactory bulb weight data and rice yield data) are used as exemplars to demonstrate our methods.

**Availability:** A software package is freely available at <http://ibi.zju.edu.cn/software/qtlnetwork>

**Contact:** [jzhu@zju.edu.cn](mailto:jzhu@zju.edu.cn)

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

### 1 INTRODUCTION

In contrast to Mendelian traits controlled by individual genes, the phenotypic variations of complex traits result from the segregation of alleles at multiple quantitative trait loci (QTLs) with effects sensitive to genetic, sexual, parental and environmental factors (Mackay, 2001). Understanding the genetic

architecture of complex traits is a major challenge in the post-genomic era, especially for QTL-by-QTL interactions (epistasis), QTL-by-sex (QS) interactions, QTL-by-environment (QE) interactions, epistasis-by-sex interactions, epistasis-by-environment interactions and more complex higher order interactions.

The term epistasis was primarily coined to describe the distortions of Mendelian segregation ratios that were due to one gene masking the effects of another in classical genetic studies on qualitative variations such as coat color (i.e. the albino allele of tyrosinase masks the phenotypes of other loci such as brown and agouti; Carlborg and Haley, 2004). Intensive work on quantitative variation also provided evidence of epistatic interactions. For example, susceptibility to lung cancer in mouse is significantly influenced by interactions among QTLs (Fijneman *et al.*, 1996). Molecular dissection of bristle number in *Drosophila* has also revealed a substantial interaction between two QTLs: the combination of these loci had much larger effect than that predicted by the sum of their individual effects (Gurganus *et al.*, 1999; Long *et al.*, 1995). Strong interactions between QTLs have also been observed in maize (Lukens and Doebley, 1999) and soybean (Lark *et al.*, 1995). The genotypic effect of one locus on phenotype might depend on the genotype at several or many other loci, such as the dependence of mutant phenotype on modifier genes in mouse (Gerlai, 1996) and fruit fly (de Belle and Heisenberg, 1996). In addition, QTL with minor or no individual effect can also be involved in epistatic interaction, a finding that is well documented for a number of physiological traits in *Drosophila melanogaster* (Montooth *et al.*, 2003).

With the advance of molecular marker techniques, fine-scale genetic and physical chromosomal maps of various organisms are now available. Using these maps, methods of mapping QTLs in experimental populations have become pervasive if not even high throughput (Haley and Knott, 1992; Jansen, 1994; Lander and Botstein, 1989; Zeng, 1994). Methods have been developed for detecting QE interactions (Piepho, 2000) and for detecting epistasis among QTLs (Kao *et al.*, 1999; Ljungberg *et al.*, 2004; Sen and Churchill, 2001; Yi *et al.*, 2003). Some of these methods are based on simultaneous scans to detect epistatic QTLs that do not have individual effects (Ljungberg *et al.*, 2004; Sen and Churchill, 2001). However, none of these methods has integrated QE interactions and epistasis into one

\*To whom correspondence should be addressed.

mapping system. Wang *et al.* (1999) proposed a two-locus model for data from multi-environment trials (METs), which could simultaneously analyze QTL main effects, epistatic effects as well as their interactions with environments. But their approach has some drawbacks in cofactor selection, false positive rate control and computational tractability. Moreover, parameter estimation of this method is conducted using a two-locus model, which does not take the whole genetic architecture into account and might contribute to biased estimation of genetic effects.

In the present study, a full-QTL model is proposed for modeling the genetic architecture of complex trait, which integrates the effects of multiple QTLs, epistasis and QE interactions into one mapping system. A Bayesian method implemented with Gibbs sampling is used to estimate genetic parameters in the full-QTL model. A systematic mapping strategy is developed to search for QTLs and their interactions by the *F*-test based on Henderson method III, which requires less computation than the likelihood ratio test. Monte Carlo simulation studies and real data sets in mouse and rice are used to demonstrate the utility of the method.

## 2 METHODS

### 2.1 Modeling the genetic architecture of complex trait from multi-environment trials (METs)

Consider a population consisting of  $N$  recombinant inbred lines (RILs) or doubled-haploid lines (DHLs) derived from a cross between two homozygous inbred lines ( $P_1$  and  $P_2$ ). The experiments are conducted in  $p$  different environments. Suppose there are  $s$  segregating QTLs ( $Q_1, Q_2, \dots, Q_s$ ) each with two genotypes  $QQ$  and  $qq$ , in which  $t$  pairs of QTLs are involved in epistatic interactions. Let a random variable  $\xi_{ki}$  be the genotype of  $Q_k$  from the  $i$ -th line taking  $\xi_{ki} = 1$  if the genotype of  $Q_k$  is  $Q_kQ_k$  and  $\xi_{ki} = -1$  if the genotype of  $Q_k$  is  $q_kq_k$ . Let  $x_{ki} = E(\xi_{ki} | \text{the genotypes of flanking markers})$  which can be achieved by a general algorithm proposed by Jiang and Zeng (1997) with dominant, codominant or missing markers. Regarding the environmental effects as random effects, the phenotypic value of the  $i$ -th line in the  $j$ -th environment ( $y_{ij}$ ) can be expressed by the following mixed linear model

$$y_{ij} = \mu + \sum_k^s a_k x_{ki} + \sum_{k,h \in \{1,2,\dots,s\}, k \neq h}^t aa_{kh} x_{ki} x_{hi} + e_j + \sum_k^s ae_{kj} x_{ki} + \sum_{k,h \in \{1,2,\dots,s\}, k \neq h}^t aae_{khj} x_{ki} x_{hi} + \varepsilon_{ij} \quad (1)$$

where  $\mu$  is the population mean;  $a_k$  is the additive effect of  $Q_k$ , fixed effect;  $aa_{kh}$  is the additive-additive epistatic effect between  $Q_k$  and  $Q_h$ , fixed effect;  $e_j$  is the main effect of the  $j$ -th environment, random effect;  $ae_{kj}$  is additive-environment interaction effect between  $Q_k$  and environment  $j$ , random effect;  $aae_{khj}$  is the interaction effect between  $aa_{kh}$  and environment  $j$ , random effect;  $\varepsilon_{ij}$  is the residual effect. In such a mixed linear model, there is no constraint that any pair of effects involved in interaction must have their individual effects.

We can express Equation (1) in matrix form as

$$\begin{aligned} \mathbf{y} &= \mathbf{1}\mu + \mathbf{X}_A \mathbf{b}_A + \mathbf{X}_{AA} \mathbf{b}_{AA} + \mathbf{U}_E \mathbf{e}_E + \sum_{k=1}^s \mathbf{U}_{A_k E} \mathbf{e}_{A_k E} + \sum_{h=1}^t \mathbf{U}_{AA_h E} \mathbf{e}_{AA_h E} + \mathbf{e}_\varepsilon \\ &= [\mathbf{1} : \mathbf{X}_A : \mathbf{X}_{AA}] [\mu : \mathbf{b}'_A : \mathbf{b}'_{AA}]' + \sum_{u=1}^r \mathbf{U}_u \mathbf{e}_u + \mathbf{I}_\varepsilon \mathbf{e}_\varepsilon \\ &= \mathbf{X} \mathbf{b} + \sum_{u=1}^{r+1} \mathbf{U}_u \mathbf{e}_u \end{aligned} \quad (2)$$

where  $\mathbf{y}$  is an  $n \times 1$  vector of phenotypic values;  $n$  is the total number of observations;  $\mathbf{1}$  is an  $n \times 1$  vector with all the elements = 1.  $\mathbf{b}_A = [a_1 \ a_2 \ \dots \ a_s]'$  and  $\mathbf{b}_{AA} = [aa_1 \ aa_2 \ \dots \ aa_t]'$  with the coefficient matrix  $\mathbf{X}_A$  and  $\mathbf{X}_{AA}$ ;  $\mathbf{e}_E = [e_1 \ e_1 \ \dots \ e_p]'$   $\sim N(\mathbf{0}, \mathbf{I}\sigma_E^2)$ ,  $\mathbf{e}_{A_k E} = [ae_{k1} \ ae_{k2} \ \dots \ ae_{kp}]' \sim N(\mathbf{0}, \mathbf{I}\sigma_{A_k E}^2)$  and  $\mathbf{e}_{AA_h E} = [aae_{h1} \ aae_{h2} \ \dots \ aae_{hp}]' \sim N(\mathbf{0}, \mathbf{I}\sigma_{AA_h E}^2)$  with the coefficient matrices  $\mathbf{U}_E$ ,  $\mathbf{U}_{A_k E}$  and  $\mathbf{U}_{AA_h E}$ , respectively;  $\mathbf{e}_\varepsilon$  is an  $n \times 1$  vector of residual effects,  $\mathbf{e}_\varepsilon \sim N(\mathbf{0}, \mathbf{I}\sigma_\varepsilon^2)$ ;  $\mathbf{I}(\mathbf{U}_{r+1})$  is an  $n \times n$  identity matrix.

### 2.2 Scanning genome for QTLs and epistasis

In Equation (1), we assume that the locations of all the QTLs and the epistatic interactions among them are known. In reality, however, such information is actually unavailable before mapping. In the following sections, we will introduce a systematic mapping strategy to search for QTLs with and/or without epistatic effects in Equation (1).

**2.2.1 Mapping QTLs by 1D genome scan** To address the problem of multi-dimensional searches for the multiple loci in the whole genome, Zeng (1994) proposed an approach called composite interval mapping (CIM), which could simplify the process of mapping multiple QTLs from multiple dimensional to 1D search problem. It can be accomplished by testing for a QTL in a particular genomic region conditioned on the selected markers controlling the genetic variance of other QTLs located outside the genomic region to be tested. Denote by  $I_1, I_2, \dots, I_c$  the marker intervals selected as cofactors when testing a genomic region for putative QTL. Let  $(M_l^-, M_l^+)$  be interval  $I_l$  and  $M_l^-$  and  $M_l^+$  be its flanking markers. Let  $M_l M_l$  and  $m_l m_l$  be the two genotypes of  $M_l$ . The QTL mapping model for testing a locus  $k$  within a particular genomic region can be written as

$$y_{ij} = \mu_j + a_{kj} x_{ki} + \sum_{l=1}^c (\alpha_{jl}^- \xi_{il}^- + \alpha_{jl}^+ \xi_{il}^+) + \varepsilon_{ij} \quad (3)$$

where  $\mu_j$  is the population mean in the  $j$ -th environment;  $\alpha_{jl}$  is the effect of the  $l$ -th marker in the  $j$ -th environment;  $\xi_{il}$  takes the value of 1 or  $-1$  depending on whether the genotype of  $M_l$  is  $M_l M_l$  or  $m_l m_l$ ; the remaining variables and parameters have the same definition as those in Equation (1).

Equation (3) can be written in matrix form as

$$\mathbf{y} = \mathbf{W}_Q \mathbf{b}_Q + \mathbf{W}_M \mathbf{b}_M + \boldsymbol{\varepsilon} \quad (4)$$

where  $\mathbf{b}_Q = [a_1 \ a_2 \ \dots \ a_L]'$ ;  $\mathbf{b}_M = [\mu' \ \alpha'_1 \ \alpha'_2 \ \dots \ \alpha'_L]'$  with  $\mu = [\mu_1 \ \mu_2 \ \dots \ \mu_L]'$  and  $\boldsymbol{\alpha}_j = [\alpha_{j1}^- \ \alpha_{j1}^+ \ \alpha_{j2}^- \ \alpha_{j2}^+ \ \dots \ \alpha_{jc}^- \ \alpha_{jc}^+]'$ ;  $\mathbf{W}_Q$  and  $\mathbf{W}_M$  are the coefficient matrices corresponding to  $\mathbf{b}_Q$  and  $\mathbf{b}_M$ , respectively;  $\mathbf{y}$  and  $\boldsymbol{\varepsilon}$  have the same definitions as those in Equation (2). Without having to know which elements of  $\mathbf{b}_Q$  and  $\mathbf{b}_M$  are fixed or random, a general equation for the expected reduction sum of square of  $\mathbf{b}_Q$  can be obtained by Henderson method III (Searle, 1992) as

$$\begin{aligned} E[\text{SSR}(\mathbf{b}_Q | \mathbf{b}_M)] &= E(\mathbf{y}' \mathbf{W} \mathbf{W}^+ \mathbf{y} - \mathbf{y}' \mathbf{W}_M \mathbf{W}_M^+ \mathbf{y}) \\ &= \text{tr}[\mathbf{W}'_Q \mathbf{A}_M \mathbf{W}_Q E(\mathbf{b}_Q \mathbf{b}'_Q)] + \sigma_\varepsilon^2 (r_W - r_{W_M}) \end{aligned} \quad (5)$$

where  $\mathbf{W} = (\mathbf{W}_Q : \mathbf{W}_M)$ ,  $\mathbf{A}_M = \mathbf{I} - \mathbf{W}_M (\mathbf{W}'_M \mathbf{W}_M)^{-1} \mathbf{W}'_M$ ,  $r_W$  is the rank of  $\mathbf{W}$  and  $r_{W_M}$  is the rank of  $\mathbf{W}_M$ . Accordingly, we can have the following *F*-statistic under the null hypothesis  $H_0: a_1 = a_2 = \dots = a_p = 0$

$$F = \frac{\text{SSR}(\mathbf{b}_Q | \mathbf{b}_M) / (r_W - r_{W_M})}{\text{SSE} / (n - r_W)} \quad (6)$$

The *F*-test can be conducted along the whole genome step by step (1D genome scan). When the *F*-values for a region exceed a pre-defined critical threshold, a QTL is indicated at that position with the regional maximum *F*-value.

Before scanning the genome for putative QTLs, marker interval analysis is required to select the candidate marker intervals as cofactors for Equation (3). Piepho and Gauch (2001) proposed a marker pair selection (MPS) approach to select cofactors for CIM. The MPS

approach has two advantages: (1) markers enter the model in adjacent pairs, which reduces the number of models to be considered, thus alleviating the problem of overfitting and increasing the chances of detecting QTLs; (2) an exhaustive search for all the marker pairs per chromosome is used instead of simple forward selection, which maximizes the chance of finding the best-fitting models. For a trait value of the  $i$ -th line in the  $j$ -th environment, a single-interval model of the  $l$ -th marker interval is written as

$$y_{ij} = \mu_j + \alpha_{jl}^- \xi_{il}^- + \alpha_{jl}^+ \xi_{il}^+ + \varepsilon_{ij} \quad (7)$$

where, all the parameters and variables have same definition as those in Equation (3). We can conduct the  $F$ -test based on the aforementioned Henderson method III for all the marker intervals, and plot the  $F$ -values along the whole genome. When the  $F$ -values at a region exceed the threshold value, a candidate marker interval is selected at the position with the regional maximum  $F$ -value.

### 2.2.2 Mapping epistasis by two-dimensional (2D) genome scan

Suppose that  $s$  QTLs have been mapped by the 1D genome scan. In order to find all possible epistasis, we adopt the 2D genome scan procedure conditional on the effects of the  $s$  QTLs mapped by the 1D genome scan as well as a group of marker interval pairs selected by marker interval interaction analysis. For any pair of marker intervals ( $I^A$  and  $I^B$ ), we can use the following two-interval model to test the interaction effect between them

$$y_{ij} = \mu_j + \alpha\alpha_j^{A^-B^-} \xi_i^{A^-} \xi_i^{B^-} + \alpha\alpha_j^{A^+B^+} \xi_i^{A^+} \xi_i^{B^+} + \sum_{k=1}^c (\alpha_{jk}^- \xi_{ik}^- + \alpha_{jk}^+ \xi_{ik}^+) + \varepsilon_{ij} \quad (8)$$

where  $\alpha\alpha_j^{A^-B^-}$  and  $\alpha\alpha_j^{A^+B^+}$  are the interaction effects between the flanking markers of  $I^A$  and  $I^B$ , ( $M^{A^-}, M^{B^-}$ ) and ( $M^{A^+}, M^{B^+}$ ). Under the null hypothesis  $H_0: \alpha\alpha_1^{A^-B^-} = \alpha\alpha_1^{A^+B^+} = \alpha\alpha_2^{A^-B^-} = \alpha\alpha_2^{A^+B^+} = \dots = \alpha\alpha_p^{A^-B^-} = \alpha\alpha_p^{A^+B^+} = 0$ , the  $F$ -test based on Henderson method III can be performed for all possible pair-wise marker intervals. When the  $F$ -values for a region exceed the pre-defined threshold value, a candidate marker interval interaction is selected at that position with the regional maximum  $F$ -value.

Let  $(I_1^A, I_1^B), (I_2^A, I_2^B), \dots, (I_f^A, I_f^B)$  be the candidate marker interval interactions selected above. The model for testing the significance of epistatic interaction between loci  $k$  and  $h$  can be written as

$$y_{ij} = \mu_j + aa_{k,hj} x_{ki} x_{hi} + \sum_s^p a_{sj} x_{si} + \sum_l^f [\alpha\alpha_{jl}^{A^-B^-} \xi_{il}^{A^-} \xi_{il}^{B^-} + \alpha\alpha_{jl}^{A^+B^+} \xi_{il}^{A^+} \xi_{il}^{B^+}] + \varepsilon_{ij} \quad (9)$$

where, all the parameters and variables have similar definitions as those described above. Under the null hypothesis  $H_0: aa_{k,h1} = aa_{k,h2} = \dots = aa_{k,hp} = 0$ , the  $F$ -statistic can be used to test any pair of loci in the genome (2D whole-genome scan). However, this 2D whole-genome scan strategy is very time-consuming, especially for the data from METs. To reduce the computational complexity, we skip the genomic regions, in which no significant marker interval interaction is detected by the interval interaction analysis.

### 2.2.3 Threshold determination and model selection

As described above, the present mapping strategy consists of the procedures of marker interval selection, detection of marker interval interaction, 1D and 2D genome scans. In each of these procedures, multiple hypothesis tests are performed across the entire genome. Thus, it is necessary to adjust for the critical threshold value of the  $F$ -statistic to control the experiment-wise false positive rate. The permutation testing (Doerge and Churchill, 1996) is employed to determine an empirical threshold value of the  $F$ -statistic. However, Equations (3, 7–9) are complex models, which contain not only the variables to be tested, but also the variables for background variance control. If the

observation values are directly permuted, the relationship between trait and background control variables will be destroyed, and thus permutation testing will give artificially low threshold. Therefore, some adjustments are needed to apply the permutation testing for complex models. Without loss of generality, we express the Equations (3, 7–9) in a general matrix form as

$$\mathbf{y} = \mathbf{W}_T \mathbf{b}_T + \mathbf{W}_C \mathbf{b}_C + \varepsilon \quad (10)$$

where,  $\mathbf{b}_T$  represents the effects to be tested with coefficient matrix  $\mathbf{W}_T$ ;  $\mathbf{b}_C$  represents the effects for background variance control with coefficient matrix  $\mathbf{W}_C$ . For each permutation, we randomly shuffle the line order of the matrix  $\mathbf{W}_T$  to destroy the relationship between the variables to be tested and the trait values, but keep the matrix  $\mathbf{W}_C$  unchanged. This method is distribution free, applicable in different population structures and especially simple to be conducted. However, a significant disadvantage is the computational burden. At least 1000 and 10000 permutations are suggested to obtain reasonably accurate estimates of the threshold value for type I error rates of 0.05 and 0.01, respectively.

Furthermore, at the end of each mapping procedure, the peaks which exceed the critical  $F$ -value calculated by permutation testing are selected from the  $F$ -statistic profile. However, some of the peaks are ghost peaks due to the high correlation of closely linked markers and random noise, etc. Thus, we perform a stepwise selection on all the significant peaks selected from the  $F$ -statistic profile after each mapping procedure using the  $F$ -statistic as criterion. The stepwise selection strategy is comprised of two steps, i.e. the forward selection step and backward selection step. Firstly, all of the selected peaks are ranked by  $F$ -values, and the one with the maximum  $F$ -value is picked up into Equation (10) as the effects for background control ( $\mathbf{b}_C$ ). In forward selection step, the remaining peaks are put into Equation (10) as tested effects ( $\mathbf{b}_T$ ) for  $F$ -test one by one, and the most significant one is retained in the model. In backward selection step, the peaks retained in the model are tested by  $F$ -statistic to examine if some retained peaks will return to non-significance due to the new peak selected by forward selection. These two selection steps are iteratively performed until the forward selection step is unable to select any peak into the model and the backward selection step is unable to drop any retained peak out of the model.

## 2.3 Parameter estimation using a Bayesian method via Gibbs sampling

After we obtain the locations of the putative QTLs and the epistatic interactions among them, we can estimate all the parameters of Equation (1) by a Bayesian method via Gibbs sampling (Wang *et al.*, 1994). In this section, we are not intending to describe this Bayesian method in detail, but providing a brief introduction of it. For a general mixed model equation, the estimates of all the effects can be obtained by solving the following normal equation,

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{U}_1 & \mathbf{X}'\mathbf{U}_2 & \dots & \mathbf{X}'\mathbf{U}_r \\ \mathbf{U}'_1\mathbf{X} & \mathbf{U}'_1\mathbf{U}_1 + \mathbf{R}_1^{-1} \frac{\sigma_{r+1}^2}{\sigma_1^2} & \mathbf{U}'_1\mathbf{U}_2 & \dots & \mathbf{U}'_1\mathbf{U}_r \\ \mathbf{U}'_2\mathbf{X} & \mathbf{U}'_2\mathbf{U}_1 & \mathbf{U}'_2\mathbf{U}_2 + \mathbf{R}_2^{-1} \frac{\sigma_{r+1}^2}{\sigma_2^2} & \dots & \mathbf{U}'_2\mathbf{U}_r \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{U}'_r\mathbf{X} & \mathbf{U}'_r\mathbf{U}_1 & \mathbf{U}'_r\mathbf{U}_2 & \dots & \mathbf{U}'_r\mathbf{U}_r + \mathbf{R}_r^{-1} \frac{\sigma_{r+1}^2}{\sigma_r^2} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{e}}_1 \\ \hat{\mathbf{e}}_2 \\ \vdots \\ \hat{\mathbf{e}}_r \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{y} \\ \mathbf{U}'_1\mathbf{y} \\ \mathbf{U}'_2\mathbf{y} \\ \vdots \\ \mathbf{U}'_r\mathbf{y} \end{bmatrix} \quad (11)$$

where  $\mathbf{R}_r$  is the Wright's relationship matrix.

In Bayesian analysis of mixed linear model, the prior distributions of all the unknown parameters are given as

$$p(\mathbf{b}) \propto \text{constant} \quad (12)$$

$$\mathbf{e}_u | \mathbf{R}_u, \sigma_u^2 \sim N_{q_u}(\mathbf{0}, \mathbf{R}_u \sigma_u^2) \quad (13)$$

$$p(\sigma_u^2) \propto (\sigma_u^2)^{-v_u/2-1} \exp\left(-\frac{1}{2} v_u s_u^2 / \sigma_u^2\right) \quad (14)$$

where  $q_u$  is the rank of  $\mathbf{U}_u$ ;  $v_u$  is a degree of belief; and  $s_u^2$  is a prior value of  $\sigma_u^2$ .

The joint posterior distribution density functions of all the parameters are

$$p(\boldsymbol{\theta}, \boldsymbol{\sigma} | \mathbf{y}, \mathbf{s}, \mathbf{v}) \propto \prod_{i=1}^{r+1} \left[ (\sigma_i^2)^{-(q_i+v_i+2)/2} \exp\left(-\frac{\lambda_i}{2\sigma_i^2}\right) \right] \quad (15)$$

where  $\boldsymbol{\theta} = [\mathbf{b}' \mathbf{e}'_1 \cdots \mathbf{e}'_r]'$ ,  $\boldsymbol{\sigma} = [\sigma_1^2 \sigma_2^2 \cdots \sigma_{r+1}^2]'$ ,  $\mathbf{v} = [v_1 \ v_2 \ \dots \ v_{r+1}]'$  and  $\mathbf{s} = [s_1^2 \ s_2^2 \ \dots \ s_{r+1}^2]'$ ;  $\lambda_i = \mathbf{e}'_i \mathbf{R}^{-1} \mathbf{e}_i$  if  $i \leq r$ , and  $\lambda_{r+1} = (\mathbf{y} - \mathbf{X}\mathbf{b} - \sum_{j=1}^r \mathbf{U}_j \mathbf{e}_j)' (\mathbf{y} - \mathbf{X}\mathbf{b} - \sum_{j=1}^r \mathbf{U}_j \mathbf{e}_j) + v_{r+1} s_{r+1}^2$ .

The full conditional posterior distributions of the parameters are

$$\theta_i | \mathbf{y}, \boldsymbol{\theta}_{-i}, \boldsymbol{\sigma} \sim N(\tilde{\theta}_i, \tilde{v}_i) \quad (16)$$

$$\sigma_i^2 | \mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\sigma}_{-i}, \mathbf{s}, \mathbf{v} \sim \tilde{v}_i \tilde{s}_i^2 \chi_{\tilde{v}_i}^{-2} \quad (17)$$

where  $\boldsymbol{\theta}_{-i}$  and  $\boldsymbol{\sigma}_{-i}$  denote  $\boldsymbol{\theta}$  and  $\boldsymbol{\sigma}$  without the  $i$ -th element;  $\tilde{\theta}_i = (\pi_i - \sum_{j=1, j \neq i}^N \omega_{ij} \theta_j) / \omega_{ii}$ ,  $\tilde{v}_i = \sigma_{r+1}^2 / \omega_{ii}$  and  $\omega_{ij}$  is the  $ij$ -th element of the first matrix of Equation (11);  $\tilde{s}_i^2 = \lambda_i / \tilde{v}_i$ .

The full conditional posterior distributions (Equations 16 and 17) are called the Gibbs samplers. Our objective is to generate random samples of parameters from the joint posterior distribution, by updating and drawing samples from the Gibbs samplers. The Gibbs sampling procedure is formally performed as

- (1) set arbitrary initial values for  $\boldsymbol{\theta}$  and  $\boldsymbol{\sigma}$ ;
- (2) generate  $\boldsymbol{\theta}$  from Equation (16) and update  $\boldsymbol{\theta}$ ;
- (3) generate  $\boldsymbol{\sigma}$  from Equation (17) and update  $\boldsymbol{\sigma}$ ;
- (4) repeat (2) and (3) for  $k$  times.

When  $k \rightarrow \infty$ , a Markov chain with an equilibrium distribution is created with Equation (15) as its density. The initial iterations are usually not collected as samples for that the chain may not have reached the equilibrium distribution yet. We can check the convergence of the iterations by running the chains under different specifications (initial values, chain length and number of samplers saved). If these different specifications result in similar results, it is assumed that the Gibbs samplers have converged to equilibrium distribution. This procedure is called burn-in. After burn-in, we run a chain of length  $L$ , and collect samples from every  $d$ -th iteration cycle (thinning interval). In practice, a conservative burn-in period of 20 000 cycles, a chain length ( $L$ ) of 200 000 and a thinning interval ( $d$ ) of 10 cycles are used for all the parameters. Parameter estimations and statistical inferences of the parameters are conducted by summarizing the Gibbs samplers.

### 3 RESULTS

#### 3.1 Monte Carlo simulation

The simulation study was conducted under a range of scenarios with different heritabilities and sample sizes. To make the

simulation as close to the reality as possible, we use a real genetic map of mice, which consists of 1095 markers covering 2037.6 cM with an average spacing of 1.86 cM. Suppose that a complex trait is controlled by a genetic architecture with five QTLs. The QTLs are designated as Q1, Q2, ..., Q5. Four of the five QTLs are involved in three pairs of epistatic interactions designated as EQ1, EQ2 and EQ3. Detailed information on the positions and genetic effects of these QTLs is presented in Tables 1 and 2. Two mapping populations with 100 and 200 RILs were generated according to the real linkage map and the hypothesized genetic architecture. Denote I and II as sample sizes of 100 and 200, respectively. Two levels of heritability (20 and 40%) are used to generate the phenotypic values. All the individuals are assumed to be investigated in three different environments. Two hundred simulations were run for each case and the average estimates and their SEs were computed.

The estimated positions and effects of QTLs and epistasis and their SEs were presented in Tables 1 and 2. The support interval (SI) calculated by the odd ratio reduced by a factor 10 (Lander and Botstein, 1989) was averaged for each of the QTLs. In general, our model can provide reasonably accurate estimates of the parameters and have acceptable performance in false positive rate control. At both of the two heritability levels with increasing sample size, the powers of detecting QTLs as well as the precision of parameter estimation increased, support intervals narrowed and the false discovery rates decreased. Comparing the accuracy of parameter estimation, powers of detecting QTLs and epistasis as well as the false discovery rates of QTLs and epistasis (Tables 1 and 2), it was found that the method performance was slightly improved with the increase of heritability. For the QTLs with RCs larger than 4.0%, the powers of detecting them were all higher than 90% even for the sample size of 100 at the heritability level of 20%. In all of the cases, the power of detecting epistasis was lower than that of QTL. Taking case I at the heritability level of 40% as an instance, although the RC of Q5 (2.31%) was lower than that of EQ1 (2.74%), the power of Q1 (52.5%) was distinctly higher than that of EQ1 (29.0%). Thus, a relatively large sample size is required for efficiently detecting epistasis.

#### 3.2 Analysis of mouse data

A data set from mouse BXD population consisting of 358 animals belonging to 35 BXD recombinant inbred strains (Williams *et al.*, 2001), was re-analyzed by the present method. These strains were generated by crossing C57BL/6J (B6) and DBA/2J (D2) parental strains in the 1970s (BXD1 through 32) and 1990s (BXD33 through 42) (Taylor *et al.*, 1999). Genotypic data of these BXD strains was obtained from the following URL: <http://www.genenet-work.org/dbdoc/BXDGeno.html>, which included 3795 markers covering 19 autosomes and the sex chromosome. Because the population size is relatively small, some of the adjacent markers have identical strain distribution patterns. We screened out 1095 haplotype markers and constructed a genetic linkage map covering 2037.6 cM with an average spacing of 1.86 cM. All the 358 animals were measured for olfactory bulb weight (OBW), brain weight (BrW) and body weight (BoW), and the trait data were downloaded from the following URL: <http://www.nervenet.org/m-ain/databases.html>

**Table 1.** Summarized simulation results for mapping QTLs with individual effects

Heritability	QTL	20%					40%				
		Q1	Q2	Q3	Q4	Q5	Q1	Q2	Q3	Q4	Q5
<sup>a</sup> RC (%)		2.42	1.88	4.20	3.43	1.15	4.84	3.76	8.40	6.86	2.31
Chromosome		3	5	6	8	10	3	5	6	8	10
Position (cM)		31.98	60.31	42.64	52.27	65.50	31.98	60.31	42.64	52.27	65.50
Estimates (SE)	I	32.22 (1.74)	60.43 (1.87)	42.70 (1.11)	52.74 (1.90)	65.10 (2.89)	32.07 (1.47)	60.68 (1.11)	42.60 (0.94)	52.58 (1.61)	65.54 (2.55)
	II	32.04 (0.95)	61.04 (3.63)	42.63 (0.51)	52.59 (0.68)	65.81 (1.54)	32.04 (0.79)	60.72 (0.57)	42.59 (0.44)	52.57 (0.58)	65.74 (1.31)
<sup>b</sup> SI Length (cM)	I	4.06	4.97	4.10	3.18	3.66	3.54	4.71	3.58	2.73	3.63
	II	3.52	4.80	4.04	3.05	4.23	2.64	3.58	3.11	2.37	3.01
<i>a</i>		-2.78	0.00	2.90	-3.31	1.92	-2.78	0.00	2.90	-3.31	1.92
Estimates (SE)	I	-2.94 (0.50)	-0.17 (0.81)	2.92 (0.62)	-3.35 (0.56)	2.37 (0.65)	-2.79 (0.39)	-0.04 (0.47)	2.90 (0.45)	-3.26 (0.45)	1.98 (0.36)
	II	-2.80 (0.26)	-0.04 (0.32)	2.93 (0.26)	-3.28 (0.27)	1.91 (0.23)	-2.78 (0.18)	-0.02 (0.21)	2.92 (0.19)	-3.28 (0.19)	1.88 (0.18)
<i>ae</i> <sub>1</sub>		0.00	-2.04	2.58	0.00	0.00	0.00	-2.04	2.58	0.00	0.00
Estimates (SE)	I	0.04 (0.22)	-2.17 (0.61)	2.38 (0.66)	-0.02 (0.31)	-0.01 (0.24)	0.02 (0.11)	-2.04 (0.36)	2.47 (0.43)	-0.01 (0.18)	-0.00 (0.18)
	II	-0.00 (0.15)	-2.04 (0.38)	2.51 (0.37)	0.01 (0.15)	0.00 (0.16)	0.00 (0.11)	-2.05 (0.26)	2.55 (0.26)	0.01 (0.11)	0.00 (0.10)
<i>ae</i> <sub>2</sub>		0.00	-0.67	-1.31	0.00	0.00	0.00	-0.67	-1.31	0.00	0.00
Estimates (SE)	I	0.00 (0.20)	-0.72 (0.56)	-1.19 (0.51)	0.01 (0.21)	0.05 (0.30)	0.00 (0.11)	-0.68 (0.37)	-1.26 (0.36)	0.01 (0.15)	0.01 (0.18)
	II	0.02 (0.16)	-0.68 (0.32)	-1.29 (0.33)	-0.03 (0.15)	0.00 (0.15)	0.01 (0.12)	-0.67 (0.23)	-1.31 (0.23)	-0.02 (0.11)	-0.00 (0.10)
<i>ae</i> <sub>3</sub>		0.00	2.72	-1.27	0.00	0.00	0.00	2.72	-1.27	0.00	0.00
Estimates (SE)	I	-0.05 (0.23)	2.89 (0.55)	-1.19 (0.52)	0.00 (0.26)	-0.03 (0.28)	-0.02 (0.11)	2.70 (0.34)	-1.23 (0.36)	-0.00 (0.15)	-0.01 (0.17)
	II	0.02 (0.14)	2.65 (0.38)	-1.23 (0.34)	0.01 (0.16)	-0.00 (0.15)	-0.01 (0.10)	2.67 (0.25)	-1.25 (0.23)	0.01 (0.12)	-0.00 (0.10)
Power (%)	I	76.5	41.5	95.5	95.0	28.5	92.0	83.5	99.5	99.5	52.5
	II	100.0	99.5	100.0	100.0	89.5	100.0	100.0	100.0	100.0	99.0

<sup>a</sup>RC, relative contribution.<sup>b</sup>SI Length, average length of support interval of 200 simulation replicates.

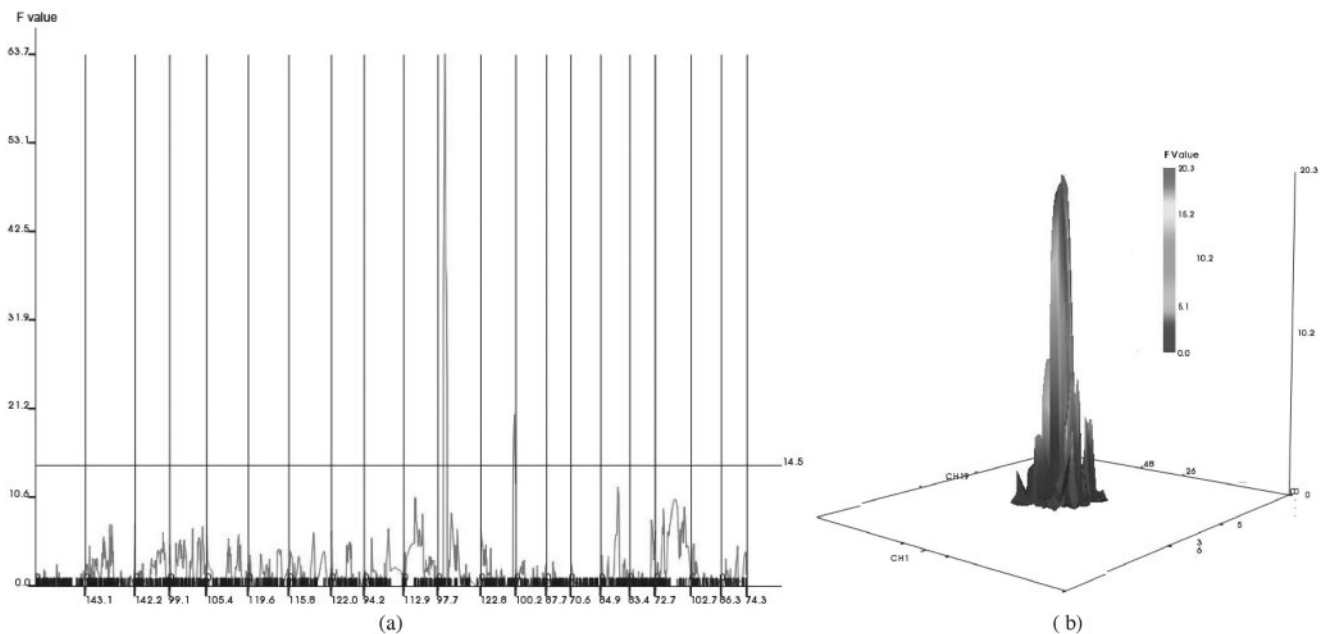
I and II represent the population sizes of 100 and 200, respectively. Each estimate presented in this table is an average of the estimates from 200 runs with the standard error (SE) in the parentheses. The false discovery rates of QTLs for cases I and II are 0.0413 and 0.0412 at the heritability level of 20%, and are 0.0584 and 0.0413 at the heritability level of 40%.

**Table 2.** Summarized simulation results for mapping epistasis

Heritability		20%			40%		
Epistasis		EQ1 (Q1–Q4)	EQ2 (Q2–Q4)	EQ3 (Q2–Q5)	EQ1 (Q1–Q4)	EQ2 (Q2–Q4)	EQ3 (Q2–Q5)
RC <sup>a</sup> (%)		1.37	3.33	2.22	2.74	6.67	4.43
aa		–2.09	2.58	–2.66	–2.09	2.58	–2.66
Estimates (SE)	I	–2.35 (0.41)	2.29 (0.53)	–2.44 (0.44)	–2.03 (0.34)	2.37 (0.42)	–2.45 (0.41)
	II	–2.02 (0.27)	2.48 (0.29)	–2.50 (0.34)	–2.04 (0.50)	2.51 (0.22)	–2.53 (0.27)
aae <sub>1</sub>		0.00	–1.79	0.00	0.00	–1.79	0.00
Estimates (SE)	I	0.01 (0.19)	–1.41 (0.63)	–0.03 (0.27)	0.01 (0.12)	–1.57 (0.42)	0.02 (0.12)
	II	0.02 (0.15)	–1.69 (0.32)	–0.00 (0.13)	0.02 (0.11)	–1.74 (0.23)	–0.00 (0.09)
aae <sub>2</sub>		0.00	2.16	0.00	0.00	2.16	0.00
Estimates (SE)	I	–0.04 (0.21)	1.86 (0.64)	–0.00 (0.21)	–0.02 (0.12)	1.94 (0.46)	0.00 (0.13)
	II	–0.02 (0.15)	2.06 (0.37)	0.00 (0.12)	–0.01 (0.11)	2.10 (0.25)	0.00 (0.08)
aae <sub>3</sub>		0.00	–0.37	0.00	0.00	–0.37	0.00
Estimates (SE)	I	0.03 (0.20)	–0.43 (0.47)	0.03 (0.12)	0.01 (0.12)	–0.38 (0.37)	–0.02 (0.15)
	II	–0.01 (0.14)	–0.39 (0.31)	–0.00 (0.12)	–0.01 (0.11)	–0.38 (0.22)	0.00 (0.09)
Power (%)	I	16.5	85.0	39.5	29.0	93.0	67.0
	II	78.0	99.0	94.5	87.5	100.0	98.0

<sup>a</sup>RC, relative contribution.

I and II represent the population sizes of 100 and 200, respectively. Each estimate presented in this table is an average of the estimates from 200 runs with the SE in the parentheses. The names of QTLs involved in each epistasis are given in the parenthesis after the name of epistasis. The false discovery rates of epistasis for the cases I and II are 0.1627 and 0.0512 at the heritability level of 20%, and are 0.1538 and 0.0473 at the heritability of 40%.



**Fig. 1.** *F*-statistic plots from (a) 1D genome scan for QTLs with individual effects and (b) 2D genome scan for epistasis of OBW in mouse. (a) Two peaks exceed the threshold *F*-value (14.5) calculated by permutation tests on chromosome 11 and 12, respectively. (b) *F*-statistic profile obtained by 2D genome scan between the regions from 26 to 48 cM on chromosome 1 and from the 5 to 36 cM on chromosome 19. A significant peak is detected which much larger than the threshold *F*-value (10.2). Colour version of this figure is available as Supplementary material online.

With the objective of identifying the OBW-specific QTLs, Williams *et al.* (2001) excluded the contributions of BrW–OBW (BrW minus OBW), BoW, age and sex to OBW by regression analysis, and used the regression-corrected values (residues) for QTL analysis. We used the original BrW but not BrW–OBW for regression analysis. Our method detected two significant QTLs

with sharp and narrow *F*-statistic peaks located on chromosome 11 and 12 by 1D genome scan (Fig. 1a). Moreover, one significant epistatic interaction was detected between two QTLs on chromosome 1 and chromosome 19 by 2D genome scan (Fig. 1b). The critical *F*-values for both 1D and 2D scan were calculated by permutation testing at genome-wide 0.05

**Table 3.** QTLs and epistasis for OBW in mouse

QTL	OBW11–17	OBW12–55	Epistasis	(OBW1–28, OBW19–13)
Chromosome	11	12	Chromosome	(1, 19)
Position (cM)	20.3	95.7	Position (cM)	(39.9, 15.2)
SI <sup>a</sup> (cM)	18.4–22.1	92.9–98.2	<sup>a</sup> SI (cM)	(36.0–40.9, 13.9–21.7)
<i>a</i> ( <i>P</i> -value)	–0.75 (0.0000)	–0.32 (0.0004)	<i>aa</i> ( <i>P</i> -value)	0.43 (0.0000)

<sup>a</sup>SI, support interval.

The QTLs are designated as ‘OBW’ with the serial numbers of chromosome and marker interval.

significance level. Genetic effects and support intervals of all the QTLs are presented in Table 3. The QTLs are designated as ‘OBW’ with the serial numbers of chromosome and marker interval. The additive effect of QTL OBW11–12 is –0.75, indicating that alleles from D2 could increase OBW by ~0.75 mg above population mean. In the study of Williams *et al.* (2001), they mapped a QTL on chromosome 11 between markers D11Mit51 and D11Mit20 (18–20 cM) with an additive effect of –0.7, which was consistent with the QTL OBW11–12 in the present study. The estimated RC of OBW11–12 was 14.6%, suggesting that ~14.6% of OBW variation is generated by this QTL. In addition, none of the two epistatic QTLs could be detected with a significant individual effect. We have strong confidence in this epistasis, not only because the *F*-statistic peak of this epistatic interaction was especially sharp, but also because there was a significant marker interaction in these two genomic regions in the previous pair-wise marker interval interaction analysis.

### 3.3 Analysis of rice data

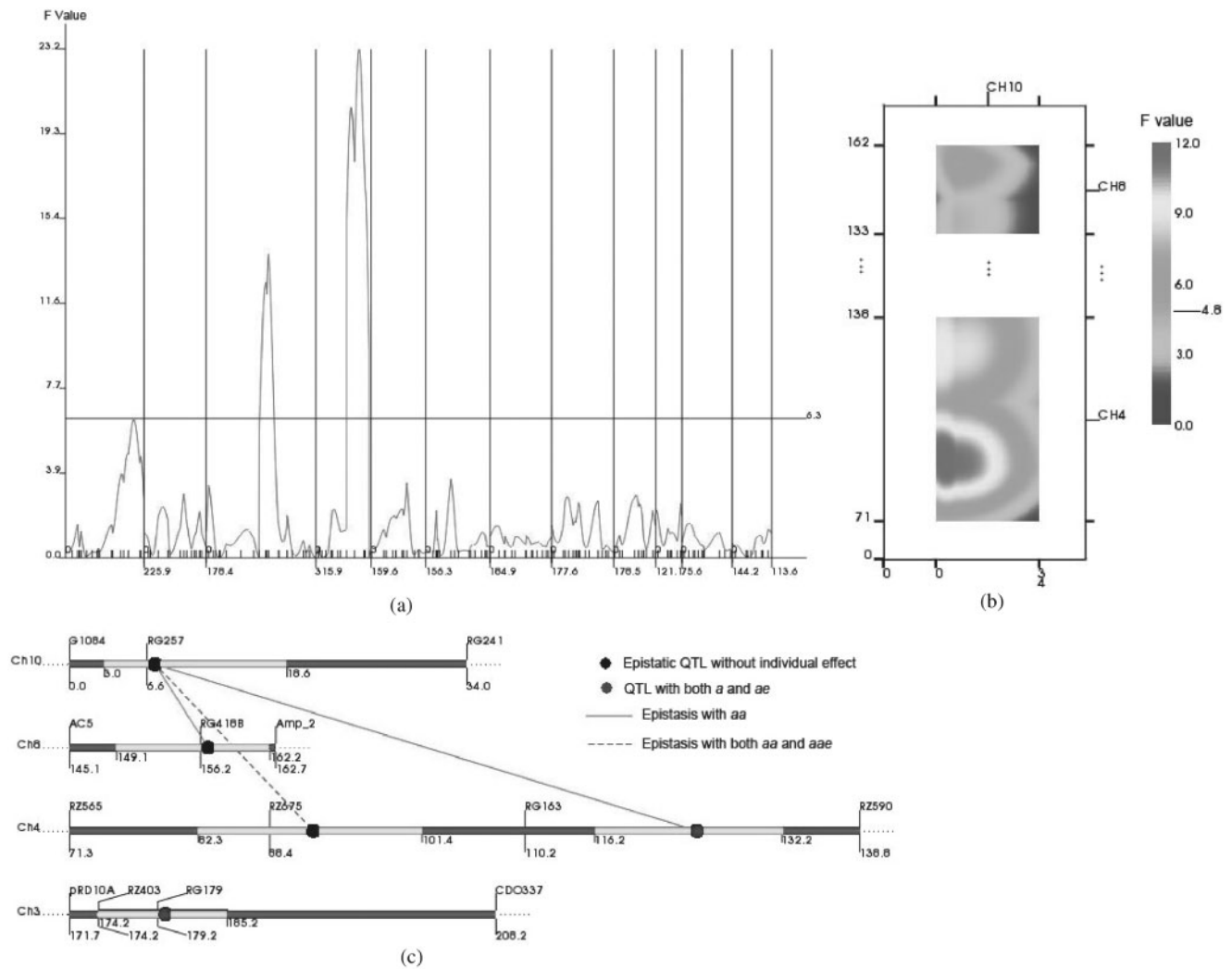
A rice DH population consisting of 123 lines derived from two inbred parents, IR64 and Azucena, was grown in a randomized block design with two replications at Hainan (18° N) in 1995 and Hangzhou (32° N) in 1996 and 1998. A total of 175 polymorphic markers covering 2005 cM of the genome along 12 chromosomes were used in this analysis (Huang *et al.*, 1997). Yield (Yd) was investigated in a series of three-year field experiments.

A total of six QTLs with three pair of epistatic interactions were detected for yield at genome-wise significance level of 0.05 by permutation testing. The *F*-statistic profiles obtained from the 1D and partial 2D genome scan procedures are shown in Figure 2a and b, and the whole genetic architecture information is summarized in an informative QTL network map (Fig. 2c). The QTLs are designated as ‘Yd’ with the serial numbers of chromosome and marker interval. Of these six QTLs, only two QTLs (Yd3–13 and Yd4–10) had additive effects, and both of them were sensitive to the environments (Table 4). Except for Yd4–10, all the epistatic QTLs had no individual effect. Special attention should be paid to the QTL Yd10–2, which was involved in all the three epistasis (Fig. 2b and c) but without individual effect according to the *F*-statistic profile obtained from 1D genome scan (Fig. 2a). It might be implied that there may be one or several modifier gene(s) that play(s) an important role in regulating the reproductive growth located in this locus. In addition, it was revealed that nearly 31.8% of

the genetic variance is attributed to those epistatic interactions between QTLs without individual effects.

## 4 DISCUSSION

In the present study, we propose a full-QTL model to integrate the effects of QTLs, epistasis and QE interaction into one mapping system, and developed a systematic mapping strategy to search for QTLs and epistasis among them. Other two popular QTL mapping methods, the multiple interval mapping (MIM; Kao *et al.*, 1999) in Windows QTL Cartographer software and the multiple imputation method (Sen and Churchill, 2001) in R/qtl software, can also handle QTLs with individual effects and epistasis. Under the assumption that the QTLs without individual effects will not be involved in epistasis, MIM method uses model selection technique to detect the epistatic interactions among the QTLs detected by CIM method. The multiple imputation method was developed based on a Monte Carlo algorithm to implement Bayesian QTL analysis. It uses 2D genome scan to search for multiple interacting QTLs. Both of these two programs cannot analyze QE interaction effects. When using these two programs to analyze data from METs, users have to analyze the data in separate environments, and compare the results from different environments to indicate the QE interaction. However, the differences in mapping results among different environments could not precisely indicate the existence of QE interaction (Jansen *et al.*, 1995). Wang *et al.* (1999) proposed a two-locus model that could also map epistasis and QE interaction for data from METs. However, it has several disadvantages especially in the mapping strategy. (1) It selects the cofactors separately in each of the environments by strategy of stepwise model selection, without controlling the genome-wise false positive rate; (2) in the genome scan procedure, only the adjacent marker intervals of the candidate marker pair are searched by the two-locus model. In the present method, marker interval analysis is conducted by data from all the environments with permutation testing to control the genome-wise false positive rate. It is much more conservative in cofactor selection than the separate analysis without genome-wise false positive rate control. With respect to the pair-wise genomic regions for 2D scan, we extend scanning regions around the candidate marker pair until no pair of marker interval has significant interaction; (3) comparison-wise significant level is used to test the significance of putative QTL and epistasis, which could result in high rate of false positive epistatic interactions. In the



**Fig. 2.** *F*-statistic plots from (a) 1D genome scan for QTLs with individual effects and (b) 2D genome scan for epistasis and (c) the predicted genetic architecture of yield in rice. (a) Two peaks exceed the threshold *F*-value (6.3) calculated by permutation tests on chromosome 3 and 4, respectively. (b) 2D genome scan is performed between the region from 0 to 32 cM on chromosome 10 and two regions on chromosome 8 and 14 for epistasis. Three peaks have been detected exceeding the threshold *F*-value (4.8). (c) The blue ball represents QTL with both additive effect and additive–environment interaction effect. The black ball represents epistatic QTL without individual effect. Chromosome region in yellow indicates the support interval of a QTL. Colour version of this figure is available as Supplementary material online.

**Table 4.** QTLs and epistasis for yield in rice

QTL	Yd3–13	Yd4–10	Epistasis	(Yd4–9, Yd10–2)	(Yd4–10, Yd10–2)	(Yd8–15, Yd10–2)
Chromosome	3	4	Chromosome	(4, 10)	(4, 10)	(8, 10)
Position	179.2	124.2	Position	(91.4, 6.6)	(124.2, 6.6)	(156.2, 6.6)
SI	174.2–185.2	116.2–132.2	SI	(82.3–101.4, 3.0–18.6)	(11(116.2–132.2, 3.0–18.6)	(1(149.1–162.2, 3.0–18.6)
<i>a</i> ( <i>P</i> -value)	–0.88 (0.0002)	–1.77 (0.0000)	<i>aa</i> ( <i>P</i> -value)	0.65 (0.0149)	0.82 (0.005)	–0.88 (0.0002)
<i>ae</i> <sub>1</sub> ( <i>P</i> -value)	1.11 (0.0236)	1.56 (0.0097)	<i>aae</i> <sub>1</sub> ( <i>P</i> -value)	–0.68 (0.1923)	0.00 (0.9971)	0.46 (0.2539)
<i>ae</i> <sub>2</sub> ( <i>P</i> -value)	–1.21 (0.0013)	–1.77 (0.0001)	<i>aae</i> <sub>2</sub> ( <i>P</i> -value)	1.07 (0.0091)	0.00 (0.9964)	–0.52 (0.1169)
<i>ae</i> <sub>3</sub> ( <i>P</i> -value)	0.12 (0.7337)	0.18 (0.6686)	<i>aae</i> <sub>3</sub> ( <i>P</i> -value)	–0.43 (0.2510)	0.00 (0.9993)	0.07 (0.8243)

<sup>a</sup>SI, support interval.

The QTLs are designated as ‘Yd’ with the serial numbers of chromosome and marker interval.



present method, the threshold  $F$ -values determined by permutation testing are used to control the genome-wide false positive rate and (4) it uses likelihood ratio test to search for the putative QTLs and epistasis, and uses Jackknife re-sampling technique for significance test. Both methods require calculation of the inverse of an  $n \times n$   $\mathbf{V}$  matrix ( $n$  is the total number of observations). We use the  $F$ -test based on Henderson method III for hypothesis tests throughout the whole mapping procedure. It completely avoids inverting the  $\mathbf{V}$  matrix. Moreover, Bayesian method is used for parameter estimation and statistical inference of the full-QTL model, which can avoid calculation of the inverse of the  $n \times n$  matrix, and also provide reasonably unbiased estimates of all the genetic effects (Tables 1 and 2). In our software, we provide an option to choose the conventional mixed model approach instead of the Bayesian method. The conventional mixed model approach consists of the REML (restricted maximum likelihood) method for variance component estimation, GLS (generalized least square) for fixed effect estimation, BLUP (best linear unbiased prediction) for random effect prediction and Jackknife re-sampling technique for significance tests of the parameters.

In the worked examples of the present study, we found some epistatic QTLs which were not detectable with individual effects by 1D genome scan. This suggested that epistatic interactions from modifier loci could be a common type of epistasis. According to Greenspan's network model to describe the flexible genome, we believe that these types of epistatic interactions are important genetic buffer which can provide much functional redundancy for species to survive perturbations, and also can generate much more phenotypic polymorphism in response to natural and artificial selection (Greenspan, 2001). However, we should bear it in mind that any strong conclusion on epistasis should be based on a relatively large population, at least more than 200 for RILs. The simulation study in the present study reveals that for the RIL population with a size of 100, the false positive rate of epistasis is higher than 0.15 and the power of detecting epistasis is relatively low (Table 2). Moreover, further experimental techniques such as genetic complementation and co-isogenic mutation analysis are required to identify genes involved in epistatic interactions to gain insights into the genetic network of evolutionarily and agriculturally important complex traits. The permanent genetic resources, such as RILs and DHLs, need to be genotyped only once and can be investigated in multiple environments, which cannot only benefit the analysis of traits with low heritabilities, but also allow an examination of the QE interactions. The method proposed in the present study provides a way to identify the QE interaction, which can play a significant role in genetic improvement of crops to obtain location-specific or broadly adapted elite varieties by marker-assisted selection. Another significant source of phenotypic variance is the sex effect, especially for the behavioral traits. By treating male and female as two environments, the present method can also analyze the sexual dimorphism of QTL and epistasis. In addition, researches in genetic epidemiology have implied that particular genes and particular genetic variation-associated disease risk in one population can differ from another because the spectrum of environments is different between populations. Identification of the interaction between

environment and gene or a group of genes would play a key role in prevention and treatment of common diseases. Thus, a more flexible method is required to handle the data from natural populations.

Based on the models and methods proposed in the present study, the computer software QTLNetwork was developed in the C++ programming language (<http://ibi.zju.edu.cn/software/qtlnetwork>). This software can be run on most of the commonly used operation systems including Microsoft Windows, Linux and UNIX. The Graphic User Interface (GUI) and graphic visualization was developed based on Microsoft Foundation Class (MFC) and Visualization Toolkit (VTK), respectively. Various kinds of populations (DH, RI, backcross,  $F_2$  and other arbitrary designs) can be handled by this software for QTL mapping.

## ACKNOWLEDGEMENTS

We are grateful to William G. Hill and Gurdev S. Khush for insightful comments during their visiting Zhejiang University. We thank Xiangyang Lou and Yousaf Hayat for their valuable suggestions on this manuscript. We also thank two anonymous reviewers for useful comments and suggestions on the earlier version of the manuscript. This research was partially supported by the National Basic Research Program of China (973 Program), the National High Technology Research and Development Program of China (863 Program) and the 111 Project (B06014).

*Conflict of Interest:* none declared.

## REFERENCES

- Carlborg, O. and Haley, C.S. (2004) Epistasis: too often neglected in complex trait studies? *Nat. Rev. Genet.*, **5**, 618–625.
- de Belle, J.S. and Heisenberg, M. (1996) Expression of *Drosophila* mushroom body mutations in alternative genetic backgrounds: a case study of the mushroom body miniature gene (*mbm*). *Proc. Natl. Acad. Sci. USA*, **93**, 9875–9880.
- Doerge, R.W. and Churchill, G.A. (1996) Permutation tests for multiple loci affecting a quantitative character. *Genetics*, **142**, 285–294.
- Fijneman, R.J. *et al.* (1996) Complex interactions of new quantitative trait loci, *Sluc1*, *Sluc2*, *Sluc3*, and *Sluc4*, that influence the susceptibility to lung cancer in the mouse. *Nat. Genet.*, **14**, 465–467.
- Gerlai, R. (1996) Gene-targeting studies of mammalian behavior: is it the mutation or the background genotype? *Trends Neurosci.*, **19**, 177–181.
- Greenspan, R.J. (2001) The flexible genome. *Nat. Rev. Genet.*, **2**, 383–387.
- Gurganus, M.C. *et al.* (1999) High-resolution mapping of quantitative trait loci for sternopleural bristle number in *Drosophila melanogaster*. *Genetics*, **152**, 1585–1604.
- Haley, C.S. and Knott, S.A. (1992) A simple regression method for mapping quantitative trait loci in line crosses using flanking markers. *Heredity*, **69**, 315–324.
- Huang, N. *et al.* (1997) RFLP mapping of isozymes, RAPD and QTLs for grain shape, brown planthopper resistance in a doubled haploid rice population. *Mol. Breed.*, **3**, 105–113.
- Jansen, R.C. (1994) Controlling the type I and type II errors in mapping quantitative trait loci. *Genetics*, **138**, 871–881.
- Jansen, R.C. *et al.* (1995) Genotype-by-environment interaction in genetic mapping of multiple quantitative trait loci. *Theor. Appl. Genet.*, **91**, 33–37.
- Jiang, C. and Zeng, Z.B. (1997) Mapping quantitative trait loci with dominant and missing markers in various crosses from two inbred lines. *Genetica*, **101**, 47–58.
- Kao, C.H. *et al.* (1999) Multiple interval mapping for quantitative trait loci. *Genetics*, **152**, 1203–1216.

- Lander, E.S. and Botstein, D. (1989) Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics*, **121**, 185–199.
- Lark, K.G. et al. (1995) Interactions between quantitative trait loci in soybean in which trait variation at one locus is conditional upon a specific allele at another. *Proc. Natl. Acad. Sci. USA*, **92**, 4656–4660.
- Ljungberg, K. et al. (2004) Simultaneous search for multiple QTL using the global optimization algorithm DIRECT. *Bioinformatics*, **20**, 1887–1895.
- Long, A.D. et al. (1995) High resolution mapping of genetic factors affecting abdominal bristle number in *Drosophila melanogaster*. *Genetics*, **139**, 1273–1291.
- Lukens, L.N. and Doebley, J. (1999) Epistatic and environmental interactions for quantitative trait loci involved in maize evolution. *Genet. Res.*, **74**, 291–302.
- Mackay, T.F. (2001) The genetic architecture of quantitative traits. *Annu. Rev. Genet.*, **35**, 303–339.
- Montooth, K.L. et al. (2003) Mapping determinants of variation in energy metabolism, respiration and flight in *Drosophila*. *Genetics*, **165**, 623–635.
- Piepho, H.P. (2000) A mixed-model approach to mapping quantitative trait loci in barley on the basis of multiple environment data. *Genetics*, **156**, 2043–2050.
- Piepho, H.P. and Gauch, H.G. Jr (2001) Marker pair selection for mapping quantitative trait loci. *Genetics*, **157**, 433–444.
- Searle, S.R. et al. (1992) *Variance Components*. John Wiley & Sons, New York.
- Sen, S. and Churchill, G.A. (2001) A statistical framework for quantitative trait mapping. *Genetics*, **159**, 371–387.
- Taylor, B.A. et al. (1999) Genotyping new BXD recombinant inbred mouse strains and comparison of BXD and consensus maps. *Mamm. Genome*, **10**, 335–348.
- Wang, C.S. et al. (1994) Bayesian analysis of mixed linear models via Gibbs sampling with an application to litter size in Iberian pigs. *Genet. Sel. Evol.*, **26**, 91–115.
- Wang, D.L. et al. (1999) Mapping QTLs with epistatic effects and QTL environment interactions by mixed linear model approaches. *Theor. Appl. Genet.*, **99**, 1255–1264.
- Williams, R.W. et al. (2001) Genetic dissection of the olfactory bulbs of mice: QTLs on four chromosomes modulate bulb size. *Behav. Genet.*, **31**, 61–77.
- Yi, N. et al. (2003) Bayesian model choice and search strategies for mapping interacting quantitative trait Loci. *Genetics*, **165**, 867–883.
- Zeng, Z.B. (1994) Precision mapping of quantitative trait loci. *Genetics*, **136**, 1457–1468.