

A Unified Markov Chain Monte Carlo Framework for Mapping Multiple Quantitative Trait Loci

Nengjun Yi¹

Section on Statistical Genetics, Department of Biostatistics, University of Alabama, Birmingham, Alabama 35294-0022

Manuscript received January 6, 2004

Accepted for publication February 25, 2004

ABSTRACT

In this article, a unified Markov chain Monte Carlo (MCMC) framework is proposed to identify multiple quantitative trait loci (QTL) for complex traits in experimental designs, based on a composite space representation of the problem that has fixed dimension. The proposed unified approach includes the existing Bayesian QTL mapping methods using reversible jump MCMC algorithm as special cases. We also show that a variety of Bayesian variable selection methods using Gibbs sampling can be applied to the composite model space for mapping multiple QTL. The unified framework not only results in some new algorithms, but also gives useful insight into some of the important factors governing the performance of Gibbs sampling and reversible jump for mapping multiple QTL. Finally, we develop strategies to improve the performance of MCMC algorithms.

MANY complex traits are controlled by multiple genetic [quantitative trait loci (QTL)] and environmental factors. Mapping QTL is the process of estimating the number of QTL, their genomic positions, and genetic effects conditional on the observed phenotypic data and marker data. This is essentially a problem of model selection (*e.g.*, BROMAN and SPEED 2002; SILLANPÄÄ and CORANDER 2002). QTL mapping is complicated by the fact that the number of QTL and hence the dimensionality of the parameter space are unknown. Recently, the Bayesian methods and Markov chain Monte Carlo (MCMC) algorithms have been applied to jointly infer the number of QTL, their genomic positions, and genetic effects. The reversible-jump MCMC algorithm introduced by GREEN (1995) can move between models of different dimension and has become an almost routine tool in Bayesian QTL mapping (HOESCHELE 2001). Using the reversible-jump MCMC method, we can in principle jointly infer the genetic model of a complex trait and the associated genetic parameters, including the number, positions, and genetic effects of the identified QTL. Recently, a variety of reversible-jump algorithms have been conducted to map QTL in both experimental designs (SATAGOPAN and YANDELL 1996; SILLANPÄÄ and ARJAS 1998, 1999; STEPHENS and FISCH 1998; YI and XU 2000; GAFFNEY 2001) and pedigrees (HEATH 1997; UIMARI and HOESCHELE 1997; XU and YI 2000; UIMARI and SILLANPÄÄ 2001; YI and XU 2001).

The reversible-jump MCMC algorithm is a very gen-

eral and widely applicable technique (GREEN 1995, 2003). It appears to be suited for implementing model selection procedures across a wide range of possible genetic architectures. However, this flexible method has been deemed somewhat “difficult” to understand, cumbersome to conduct, and difficult to tune. It also has been noted that the reversible-jump MCMC is usually subject to poor mixing and slow convergence. Therefore, there seems to be a need for further methodological work on improving the efficiency of reversible jump. The improved frameworks have been established recently for conventional statistical models (GODSILL 2001; BROOKS *et al.* 2003; GREEN 2003). It is clear that Bayesian QTL mapping can benefit from renewed research efforts.

For conventional linear models, a variety of MCMC methods have been proposed for variable selection, including the variable selection algorithms of SMITH and KOHN (1996) and KUO and MALLICK (1998), the MCMC model combination (MC³) technique of RAFTERY *et al.* (1997), the Gibbs variable selection of DELLAPORTAS *et al.* (2002), and the stochastic search variable selection of GEORGE and McCULLOCH (1993). For certain situations, these different methods have their own advantages. To date, however, they have been rarely applied to the area of mapping QTL (but see BROMAN and SPEED 2002; YI *et al.* 2003b). The variable selection methods, originally derived from diverse procedures, have been recently shown to relate closely to Green’s reversible-jump MCMC (CLYDE 1999; NTZOUFRAS 1999; GODSILL 2001; DELLAPORTAS *et al.* 2002). GODSILL (2001) recently introduced a composite model space framework that embraces not only all of the above variable selection methods, but also the reversible jump. The composite space method,

¹Address for correspondence: Department of Biostatistics, Ryals Public Health Bldg., 1665 University Blvd., University of Alabama, Birmingham, Alabama 35294-0022. E-mail: nyi@ms.soph.uab.edu

which is a modification of the product space used by CARLIN and CHIB (1995), provides an interesting viewpoint on model selection, since it allows MCMC simulation to be performed, at least conceptually, on a fixed dimension space. Under the composite space representation, the Bayesian variable selection methods described above and the reversible-jump algorithm can be shown to derive straightforwardly from a general framework. These relationships between the methods can aid our understanding of MCMC model selection procedures and may assist in the development of improved procedures.

In this study, we propose a composite space presentation for the multiple-QTL model and develop a unified MCMC framework for exploring the posterior of the composite space. The proposed unified approach includes the existing Bayesian QTL mapping methods using reversible-jump MCMC algorithm as special cases. We also show that a variety of Bayesian variable selection methods using Gibbs sampling can be applied to map multiple QTL. The unified framework sheds some light upon the important factors governing the performance of Gibbs sampling and reversible jump for mapping multiple QTL. We also develop strategies to improve the performance of MCMC algorithms.

THE MULTIPLE-QTL MODEL

We consider a mapping population derived from two or multiple inbred lines. Suppose that the quantitative trait under investigation is affected by l loci (QTL). If no epistasis is assumed, the observed phenotypic value of individual i , y_i , can be described by the linear model

$$y_i = \mu + \sum_{j=1}^l \mathbf{x}_{ij} \boldsymbol{\beta}_j + e_i, \tag{1}$$

where μ is the population mean, \mathbf{x}_{ij} denotes the genotype indicator of the j th QTL for individual i , $\boldsymbol{\beta}_j$ is the vector of genetic effects associated with the j th QTL, and e_i is the residual error assumed to follow $N(0, \sigma^2)$. The definitions of \mathbf{x}_{ij} and $\boldsymbol{\beta}_j$ depend on the experimental design. For an F_2 cross, for example, we have that

$$\mathbf{x}_{ij} = \begin{cases} (1 \quad -0.5)^T & \text{if the genotype is QQ} \\ (0 \quad 0.5)^T & \text{if the genotype is Qq and } \boldsymbol{\beta}_j = (a_j, d_j)^T \\ (-1 \quad -0.5)^T & \text{if the genotype is qq,} \end{cases}$$

where a_j and d_j are the additive and dominance effects of the j th QTL, respectively.

The above model is typically used in Bayesian mapping implemented via the reversible-jump MCMC algorithm. In this model, the number of QTL is treated as a random variable, and thus the total number of possible effects is unknown. In practical implementation of reversible-jump MCMC, we usually assume that the random variable l has an upper bound K . Model (1) can be rewritten as

$$y_i = \mu + \sum_{j=1}^K \gamma_j \mathbf{x}_{ij} \boldsymbol{\beta}_j + e_i, \tag{2}$$

where γ_j is an indicator variable that denotes that the j th QTL is included ($\gamma_j = 1$) in the model or excluded from the model ($\gamma_j = 0$). Note that the number of QTL does not explicitly appear in model (2). This parameter equals the number of 1's in $\boldsymbol{\gamma} = \{\gamma_j\}_{j=1}^K$.

Model (2) is similar to that used in Bayesian variable selection for the linear regression model (e.g., KUO and MALLICK 1998). The idea of adding the indicator variable in the model facilitates setting up MCMC algorithms. As in the linear regression model, we treat K as known and thus in model (2) the total number of possible effects is fixed.

The choice of the constant K depends on the method and the aim in the analysis. In marker analysis, each marker is treated as a potential QTL and thus K equals the number of markers (BALL 2001; BROMAN and SPEED 2002; XU 2003; Yi *et al.* 2003b). In QTL mapping, one does not know *a priori* how many QTL to expect for a given trait. We here propose two methods for choosing K : (1) As in almost all existing Bayesian mapping methods, we assume that there are at most K QTL in the entire genome, and (2) we assume that there are at most K_c QTL on the c th chromosome. Then we have $K = \sum_c K_c$. As an extreme case, we could assume that each marker interval is associated with a QTL and thus K_c is identical to the number of marker intervals on the c th chromosome. The assumption that there is at most one QTL on a marker interval is not a fundamental requirement for the proposed method. Generally, the value of K can be smaller than the number of marker intervals. The value of K should account for the data information and the previous results obtained by using other QTL mapping methods. As is seen later, alternatively, we can use particular prior distributions on the parameters in the model to relieve the influence of K on the performance of the proposed algorithms. In particular, these prior distributions account for the sample size n , the marker information, and the upper bound of QTL K .

In the following sections, we first propose a composite space representation of the problem for mapping multiple QTL based on model (2). We then discuss the specifications of the prior distributions on the unknowns. Finally, we develop a unified MCMC framework for exploring the composite model space.

THE COMPOSITE SPACE FOR THE MULTIPLE-QTL MODEL

In QTL studies, we observe the phenotypic trait and a set of marker genotypes. Assume that marker linkage maps have been developed on the basis of the observed marker data so that the locations of the markers on each chromosome are known *a priori*. Our aim is to

jointly infer the number of QTL, their genomic positions, and genetic effects. This can be viewed essentially as a problem of model selection. In model (2), the number of QTL is determined by the vector of indicator variables $\boldsymbol{\gamma} = \{\gamma_j\}_{j=1}^K$. Hereafter, we call the vector $\boldsymbol{\gamma}$ the model index, which indicates which QTL are present in the model.

In marker analysis, model (2) is essentially a usual linear regression model in that each coefficient x_{ij} is observed. In QTL mapping, the coefficients in the model, $\mathbf{x} = \{x_{ij}\}_{i=1, j=1}^{n, K}$, are unobservable, and the locations of K QTL, $\boldsymbol{\lambda} = \{\lambda_j\}_{j=1}^K$, are also unknown. Denote $\boldsymbol{\beta} = \{\beta_j\}_{j=1}^K$ and $\boldsymbol{\theta} = (\boldsymbol{\beta}, \mu, \sigma^2)$. We partition $(\boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$ into $(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)$ and $(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma})$, representing the unknowns included ($\gamma_j = 1$) or excluded ($\gamma_j = 0$) from the model, respectively, where $\boldsymbol{\theta}_\gamma = (\boldsymbol{\beta}_\gamma, \mu, \sigma^2)$ and $\boldsymbol{\theta}_{-\gamma} = \boldsymbol{\beta}_{-\gamma}$. Hereafter, we call $(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$ the composite space for the multiple-QTL model. For the detailed description about the composite space for model uncertainty problems, the reader is referred to GODSILL (2001, 2003).

Under model (2), the likelihood function for a particular $\boldsymbol{\gamma}$ depends only upon the parameters $(\mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)$ used by that model, *i.e.*,

$$p(\mathbf{y}|\boldsymbol{\gamma}, \mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{y}|\boldsymbol{\gamma}, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma). \quad (3)$$

We assume that the prior distribution of $(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$ can be partitioned as

$$\begin{aligned} p(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}) &= p(\boldsymbol{\gamma})p(\boldsymbol{\theta}, \boldsymbol{\lambda}, \mathbf{x}|\boldsymbol{\gamma}) \\ &= p(\boldsymbol{\gamma})p(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma|\boldsymbol{\gamma})p(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma}, \boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma). \end{aligned} \quad (4)$$

The full posterior distribution of the composite model space $(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$ can now be expressed as

$$\begin{aligned} p(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) &\propto p(\mathbf{y}|\boldsymbol{\gamma}, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)p(\boldsymbol{\gamma})p(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma|\boldsymbol{\gamma}) \\ &\cdot p(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma}, \boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma). \end{aligned} \quad (5)$$

Note that here we have suppressed the notation for conditional on the observed marker data.

In the above posterior distribution, $p(\boldsymbol{\gamma})$ is the prior distribution of the model index. $p(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma|\boldsymbol{\gamma})$ is the joint prior distribution of the used unknowns, which can be partitioned into three components:

$$p(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma|\boldsymbol{\gamma}) = p(\boldsymbol{\lambda}_\gamma|\boldsymbol{\gamma})p(\mathbf{x}_\gamma|\boldsymbol{\lambda}_\gamma)p(\boldsymbol{\theta}_\gamma|\boldsymbol{\gamma}, \mathbf{x}_\gamma). \quad (6)$$

The prior for the unused unknowns, $p(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma}, \boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)$, may be called ‘‘pseudo-prior.’’ It is reasonable to assume that $(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma})$ is *a priori* independent of $(\boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)$. The pseudo-prior can be factorized into three components:

$$p(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma}, \boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma) = p(\boldsymbol{\lambda}_{-\gamma}|\boldsymbol{\gamma})p(\mathbf{x}_{-\gamma}|\boldsymbol{\lambda}_{-\gamma})p(\boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma}). \quad (7)$$

In Equations 6 and 7, $p(\boldsymbol{\lambda}_\gamma|\boldsymbol{\gamma})$ and $p(\boldsymbol{\lambda}_{-\gamma}|\boldsymbol{\gamma})$ are the prior distributions of the locations of QTL. $p(\mathbf{x}_\gamma|\boldsymbol{\lambda}_\gamma)$ and $p(\mathbf{x}_{-\gamma}|\boldsymbol{\lambda}_{-\gamma})$ are the probability distributions of QTL

genotype indicators, which can be calculated using multipoint methods (JIANG and ZENG 1997). $p(\boldsymbol{\theta}_\gamma|\boldsymbol{\gamma}, \mathbf{x}_\gamma)$ is the prior distribution of the used parameters, which may depend on \mathbf{x}_γ . $p(\boldsymbol{\theta}_{-\gamma}|\boldsymbol{\gamma})$ is the prior distribution of the unused genetic effects.

The key feature of the composite model space is that the dimension remains fixed even when the model index $\boldsymbol{\gamma}$ or the number of QTL changes. This remarkable feat is achieved by augmenting the varying dimensional space $(\boldsymbol{\gamma}, \boldsymbol{\lambda}_\gamma, \mathbf{x}_\gamma, \boldsymbol{\theta}_\gamma)$ to the fixed dimensional space $(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$. Simulation of $p(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ then can be addressed via standard MCMC algorithms for a distribution of fixed dimension (GODSILL 2001). Thus convergence properties of these algorithms are inherited from standard MCMC theory. Furthermore, the composite space approach provides a method to use the important parameters for models other than the current model for efficient proposal design (GODSILL 2003; GREEN 2003).

PRIOR SPECIFICATIONS

The statistical properties of the Bayesian approach rest squarely on the specification of the prior distributions on the unknowns. This is especially true in mapping multiple QTL across the entire genome. In this section, we discuss the prior distribution of the composite model space for the multiple-QTL model.

For the specification of the model index, most Bayesian variable selection implementations have used independence priors of the form

$$p(\boldsymbol{\gamma}) = \prod_{j=1}^K w_j^{\gamma_j} (1 - w_j)^{1-\gamma_j}. \quad (8)$$

Under this prior, each QTL enters the model independently of the other QTL, with probability $p(\gamma_j = 1) = 1 - p(\gamma_j = 0) = w_j$. In QTL mapping, a reasonable reduction may be to set $w_j \equiv w$, yielding

$$p(\boldsymbol{\gamma}) = w^l (1 - w)^{K-l}, \quad (9)$$

where l is the number of QTL equal to the number of 1's in $\boldsymbol{\gamma}$. The hyperparameter w is the prior expected proportion of QTL included in the model. In particular, setting $w = 1/2$ yields the popular uniform prior

$$p(\boldsymbol{\gamma}) = 1/2^K, \quad (10)$$

which gives the same prior weight to all models and is widely used as noninformative prior in variable selection problems (NTZOUFRAS 1999). However, this prior actually puts most of its weight near models with $K/2$ QTL (CHIPMAN *et al.* 2001). Alternatively, we could put a Poisson prior with a predetermined mean L on the number of QTL or the number of 1's in $\boldsymbol{\gamma}$, *i.e.*,

$$p(\boldsymbol{\gamma}) = \frac{L^l}{l!} e^{-L}. \quad (11)$$

The locations of QTL are assumed to be independent

a priori and uniformly distributed across the entire genome or the corresponding regions. If we suppose that there are at most K_c QTL on the c th chromosome, K_c QTL could be uniformly distributed on this chromosome. Since it is usually difficult to distinguish multiple QTL on a marker interval (e.g., LYNCH and WALSH 1998), it may be reasonable to assume that there is at most one QTL on a marker interval although this is not a fundamental requirement for our method. Furthermore, this assumption can limit the model space.

The prior for the overall mean μ is normally distributed with mean η_0 and variance τ_0^2 . We could choose $\eta_0 = 0$ or $\bar{y} = (1/n)\sum_{i=1}^n y_i$ and $\tau_0^2 = s_y^2 = (1/(n-1))\sum_{i=1}^n (y_i - \bar{y})^2$. We choose an inverse $\Gamma(a, b)$ as the prior of σ^2 . GAFFNEY (2001) suggested $a = 3$ and $b = s_y^2$, which has prior mean and variance equal to $s_y^2/2$. Alternatively, we could take $p(\sigma^2) \propto 1/\sigma^2$ or $p(\sigma^2) \propto 1$.

We could use three types of prior distributions for the genetic effect β . First, we could use a normal prior for each vector of genetic effects, i.e., $\beta_j \sim N(0, \Sigma)$, $j = 1, \dots, K$, where the prior mean of zero reflects indifference between positive and negative values, and Σ is the prior covariance matrix. The covariance matrix Σ could be chosen to be diagonal. In this prior specification, the prior distribution for each QTL is identical and is independent of γ_j . Most Bayesian mapping methods have used this type of prior distribution. This prior has been used by KUO and MALLICK (1998) in Bayesian variable selection for a linear regression model. Second, we could use the prior $p(\beta_j | \gamma_j) = (1 - \gamma_j)N(0, \Sigma) + \gamma_j N(0, c^2 \Sigma)$, where c^2 is a predetermined constant. This prior has been used by GEORGE and MCCULLOCH (1993) and DELLAPORTAS *et al.* (2002) for a linear regression model. Third, we could use $p(\beta_j | \gamma_j, \mathbf{x}_j) \sim N(0, c^2(\mathbf{x}_j' \mathbf{x}_j)^{-1} \sigma^2)$, where c^2 is a hyperparameter. This prior has been used extensively in Bayesian variable selection for the conventional linear model.

The prior distributions on the model index γ and the genetic effects β may be the most critical factors influencing the performance of the algorithms and thus deserve careful attention. The hyperparameter w or L in the prior of γ controls the expected proportion of genetic effects and the number of QTL included in the model. The prior covariance matrix or the hyperparameter c in the prior of β controls the expected size of genetic effects included in the model. Small w or L and large prior variance or c would concentrate the prior on parsimonious models with large effects, and large w or L and small prior variance or c would concentrate on saturated models with small effects. The reasonable choices of c and w would account for the sample size n , the marker information, and the upper bound of QTL K . For the conventional linear model, FERNANDEZ *et al.* (2001) recommended $c = \max\{n, K^2\}$. GEORGE and FOSTER (2000) proposed treating c and w as unknown parameters and using empirical Bayes estimates of c and w based on the data. Finally, we could consider

hyperprior distributions on w or L and Σ or c (e.g., CHIPMAN *et al.* 2001; GAFFNEY 2001).

POSTERIOR CALCULATION AND EXPLORATION

In this section, we develop a unified MCMC framework for simulating from the posterior distribution, $p(\gamma, \lambda, \mathbf{x}, \theta | \mathbf{y})$, which includes the existing reversible-jump algorithms as special cases and also provides some new methods for mapping multiple QTL. It is seen that standard Gibbs samplers applied to the composite model space produce several well-known Bayesian variable selection methods developed for the linear regression model, while a more sophisticated Metropolis-Hastings (M-H) approach produces a version of the reversible-jump algorithm. We also propose strategies to improve efficiency of Bayesian mapping.

The full conditional posterior distributions for θ_γ and $\theta_{-\gamma}$ are given by

$$p(\theta_\gamma | \gamma, \lambda, \mathbf{x}, \theta_{-\gamma}, \mathbf{y}) \propto p(\mathbf{y} | \gamma, \mathbf{x}_\gamma, \theta_\gamma) p(\theta_\gamma | \gamma, \mathbf{x}_\gamma) \quad (12)$$

$$p(\theta_{-\gamma} | \gamma, \lambda, \mathbf{x}, \theta_\gamma, \mathbf{y}) \propto p(\theta_{-\gamma} | \gamma). \quad (13)$$

The full conditional posterior distributions for elements of θ_γ , e.g., μ , β_γ , and σ^2 , can be easily derived from Equation 12. These posteriors have standard forms and thus can be easily sampled (e.g., GELMAN *et al.* 1995). It can be seen that the parameters unused in the model do not influence the posterior of θ_γ . Since the unused parameters do not contribute to the likelihood, the posterior of $\theta_{-\gamma}$ is identical to its prior. For some algorithms developed later, the values of $\theta_{-\gamma}$ are used to update the model index and thus $\theta_{-\gamma}$ needs to be generated from the pseudo-prior. The pseudo-prior is likely to influence the performance of these algorithms and hence it should be specified with caution.

Since the position λ_j is highly dependent on \mathbf{x}_j , we jointly update λ_j and \mathbf{x}_j . The joint full conditional posterior distribution for the location and genotype indicator of the j th QTL is

$$p(\lambda_j, \mathbf{x}_j | \gamma, \lambda_{-j}, \mathbf{x}_{-j}, \theta, \mathbf{y}) \propto \begin{cases} p(\mathbf{y} | \gamma, \mathbf{x}_j, \theta_\gamma) p(\lambda_j | \lambda_{-j}) p(\mathbf{x}_j | \lambda_j) p(\theta_j | \gamma, \mathbf{x}_j) & \text{if } \gamma_j = 1 \\ p(\lambda_j | \lambda_{-j}) p(\mathbf{x}_j | \lambda_j) & \text{if } \gamma_j = 0, \end{cases} \quad (14)$$

where λ_{-j} (\mathbf{x}_{-j}) represents all elements of λ (\mathbf{x}) except λ_j (\mathbf{x}_j). This posterior is not a standard distribution, and thus the M-H algorithm is needed to update λ_j and \mathbf{x}_j jointly. We first propose a new location λ_j^* from $q(\lambda_j^*; \lambda_j)$, and then generate genotype indicator \mathbf{x}_j^* at this new location for all individuals from the posterior $q(\mathbf{x}_j^*) = p(\mathbf{x}_j | \gamma, \lambda, \mathbf{x}_{-j}, \theta, \mathbf{y})$. The proposals for the new location and the genotype indicator are then accepted or rejected simultaneously with probability

$$\min\left(1, \frac{p(\lambda_j^*, \mathbf{x}_j^* | \gamma, \lambda_{-j}, \mathbf{x}_{-j}, \theta, \mathbf{y}) q(\lambda_j; \lambda_j^*) q(\mathbf{x}_j)}{p(\lambda_j, \mathbf{x}_j | \gamma, \lambda_{-j}, \mathbf{x}_{-j}, \theta, \mathbf{y}) q(\lambda_j^*; \lambda_j) q(\mathbf{x}_j^*)}\right) \quad (15)$$

(Yi and Xu 2001). Two schemes can be employed to propose a new location, λ_j^* : (1) local move, propose λ_j^* from Uniform($\lambda_j - d, \lambda_j + d$), where d is a predetermined tuning parameter; and (2) long range move, propose λ_j^* uniformly from the corresponding region (GAFFNEY 2001). Note that under the conjugate prior, the parameter θ can be integrated out from the posterior distribution (14). Therefore, the joint full conditional posterior distribution becomes

$$p(\lambda_j, \mathbf{x}_j | \boldsymbol{\gamma}, \boldsymbol{\lambda}_{-j}, \mathbf{x}_{-j}, \mathbf{y}) \propto \begin{cases} p(\mathbf{y} | \boldsymbol{\gamma}, \mathbf{x}_j) p(\lambda_j | \boldsymbol{\lambda}_{-j}) p(\mathbf{x}_j | \lambda_j) & \text{if } \gamma_j = 1 \\ p(\lambda_j | \boldsymbol{\lambda}_{-j}) p(\mathbf{x}_j | \lambda_j) & \text{if } \gamma_j = 0, \end{cases} \quad (16)$$

which is independent of θ .

The full conditional posterior distribution of x_{ij} is given by

$$p(x_{ij} | \boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}_{-ij}, \boldsymbol{\theta}, \mathbf{y}) \propto \begin{cases} p(y_i | \boldsymbol{\gamma}, \mathbf{x}_j, \boldsymbol{\theta}_\gamma) p(x_{ij} | \lambda_j) p(\boldsymbol{\theta}_\gamma | \boldsymbol{\gamma}, \mathbf{x}_j) & \text{if } \gamma_j = 1 \\ p(x_{ij} | \lambda_j) & \text{if } \gamma_j = 0, \end{cases} \quad (17)$$

where \mathbf{x}_{-ij} represents all elements of \mathbf{x} except x_{ij} . This posterior is a discrete distribution and thus easily sampled.

Note that when the j th QTL is not included in the model, the posteriors of the genetic effects, location, and genotype indicators for this QTL are identical to the corresponding priors. The values of these unknowns are required to update the indicator variable of the QTL. Therefore, we first describe the methods for updating the model index on the basis of the values of these unknowns sampled from their priors. However, sampling from the priors does not update or make use of our current knowledge about these unknowns and hence cannot possibly produce an optimal sampler.

The standard MCMC procedures, Gibbs sampler and M-H algorithm, can be applied to update the model index $\boldsymbol{\gamma}$. Several different methods can be developed as follows.

Method I: The full conditional posterior distribution of the indicator variable γ_j is given by

$$p(\gamma_j = s | \boldsymbol{\gamma}_{-j}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}, \mathbf{y}) \propto p(\mathbf{y} | \gamma_j = s, \boldsymbol{\gamma}_{-j}, \mathbf{x}_j, \boldsymbol{\theta}_\gamma) p(\gamma_j = s, \boldsymbol{\gamma}_{-j}) p(\boldsymbol{\beta}_j | \gamma_j = s), \quad j = 1, \dots, K, \quad (18)$$

where $\boldsymbol{\gamma}_{-j}$ represents all elements of $\boldsymbol{\gamma}$ except γ_j . This posterior is a Bernoulli distribution and thus easily sampled. The sampling can be implemented sequentially or in random order.

This Gibbs sampler includes several Bayesian variable selection methods as special cases, depending on the prior specifications of $\boldsymbol{\beta}_j$ (NTZOUFRAS 1999; DELLAPORTAS *et al.* 2002). KUO and MALLICK (1998) use a prior distribution $p(\boldsymbol{\beta}_j)$, which is independent of γ_j so that $p(\boldsymbol{\beta}_j | \gamma_j = 1) = p(\boldsymbol{\beta}_j | \gamma_j = 0)$. Then, the third term on the right-hand side of (18) can be omitted. Similar to GEORGE and McCULLOCH (1993), DELLAPORTAS *et al.* (2002) use a mixture of normal distribution for model

parameters, *i.e.*, $p(\boldsymbol{\beta}_j | \gamma_j) = (1 - \gamma_j)N(0, \Sigma) + \gamma_j N(0, c^2 \Sigma)$, so that $p(\boldsymbol{\beta}_j | \gamma_j = 1) = N(0, c^2 \Sigma)$ and $p(\boldsymbol{\beta}_j | \gamma_j = 0) = N(0, \Sigma)$. Then, the third term should remain in (18).

We also can apply the M-H algorithm to update the model index $\boldsymbol{\gamma}$ conditional on other unknowns. The M-H algorithm is not based on sampling directly from the full conditional, but on a proposal for a move from $\boldsymbol{\gamma}$ to $\boldsymbol{\gamma}'$, followed by acceptance or rejection of this proposal. Although the M-H sampler can in principle update multiple components of $\boldsymbol{\gamma}$ simultaneously, we discuss only the simplest strategy where only one component in $\boldsymbol{\gamma}$ is proposed; thus at each iteration we actually propose to add or delete one QTL. We assume that the j th element of $\boldsymbol{\gamma}$ is proposed with probability $q(\boldsymbol{\gamma}'; \boldsymbol{\gamma})$; then the acceptance probability, using the standard M-H algorithm, is given by $\min(1, r)$, where the acceptance ratio r is

$$r = \frac{p(\gamma_j = s | \boldsymbol{\gamma}_{-j}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{p(\gamma_j = 1 - s | \boldsymbol{\gamma}_{-j}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}, \mathbf{y})} \cdot \frac{q(\boldsymbol{\gamma}; \boldsymbol{\gamma}')}{q(\boldsymbol{\gamma}'; \boldsymbol{\gamma})} \\ = \frac{p(\mathbf{y} | \boldsymbol{\gamma}', \mathbf{x}_j, \boldsymbol{\theta}_{\boldsymbol{\gamma}'}) p(\boldsymbol{\gamma}') p(\boldsymbol{\beta}_j | \gamma_j = s)}{p(\mathbf{y} | \boldsymbol{\gamma}, \mathbf{x}_j, \boldsymbol{\theta}_\gamma) p(\boldsymbol{\gamma}) p(\boldsymbol{\beta}_j | \gamma_j = 1 - s)} \cdot \frac{q(\boldsymbol{\gamma}; \boldsymbol{\gamma}')}{q(\boldsymbol{\gamma}'; \boldsymbol{\gamma})}, \quad (19)$$

where $\boldsymbol{\gamma} = (\gamma_j = 1 - s, \boldsymbol{\gamma}_{-j})$, $\boldsymbol{\gamma}' = (\gamma_j = s, \boldsymbol{\gamma}_{-j})$, and $s = 1$ or 0 corresponding to adding or deleting one QTL, respectively.

The proposals $q(\boldsymbol{\gamma}'; \boldsymbol{\gamma})$ can be set to p_a or $p_a / (l + 1)$, depending on $s = 1$ or 0 , where l is the number of 1's in $\boldsymbol{\gamma}$, which equals the current number of QTL, and p_a and p_a are constants satisfying $p_a + p_a = 1$. This proposal scheme is equivalent to that commonly used in Bayesian QTL mapping. Alternatively, we can set $[p(\boldsymbol{\gamma}) q(\boldsymbol{\gamma}; \boldsymbol{\gamma}')] / [p(\boldsymbol{\gamma}') q(\boldsymbol{\gamma}'; \boldsymbol{\gamma})] = 1$ (GAFFNEY 2001). Under our composite model space, however, two new schemes can be developed, borrowing the idea of variable selection: (1) We pick one of the K variables (QTL) at random and either delete or add it if it is currently or not, respectively, in the model; thus we have that $q_1(\boldsymbol{\gamma}'; \boldsymbol{\gamma}) = q_1(\boldsymbol{\gamma}; \boldsymbol{\gamma}') = 1/K$, or (2) we can update γ_j for all $j = 1, \dots, K$ sequentially or in random order; thus we have that $q_1(\boldsymbol{\gamma}'; \boldsymbol{\gamma}) = q_1(\boldsymbol{\gamma}; \boldsymbol{\gamma}') = 1$. Under both these schemes, the move proposal probability cancels from the acceptance probability ratio (19).

The above M-H algorithm is equivalent to a reversible-jump algorithm with reflecting boundaries at 0 and K QTL. To describe this relationship, we assume $s = 1$, which corresponds to adding one QTL into the model. The reversible jump can proceed to generate a new location and the genotype indicators at the new location from the priors and the associated effects $\boldsymbol{\beta}_j$ from $p(\boldsymbol{\beta}_j | \gamma_j = 0)$. Then, the acceptance ratio is given, using the reversible-jump algorithm of GREEN (1995, 2003), by (19). This reversible-jump algorithm has been widely used in Bayesian QTL mapping (*e.g.*, HEATH 1997; SIL-LANPÄÄ and ARJAS 1998, 1999; STEPHENS and FISCH 1998; Yi and Xu 2000).

Method II: The above algorithm is conditional on the

value of β_j sampled from the pseudo-prior when β_j is proposed to add into the model. More efficient MCMC algorithms by using blocking strategies can be devised to yield improved performance. Under the linear model (2), we can choose the pseudo-prior for β_j to be the conditional posterior for β_j with $\gamma_j = 1$; that is, set

$$p(\beta_j | \gamma_j = 0) = p(\beta_j | \gamma_j = 1, \gamma_{-j}, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y}),$$

where $\theta_{-\beta_j}$ means all elements of θ except β_j . The sampling step for γ_j then reduces to

$$p(\gamma_j = s | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y}) = \int_{\beta_j} p(\gamma_j = s, \beta_j | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y}) d\beta_j \quad (20)$$

(GODSILL 2001). This approach is equivalent to a sampling scheme, which first draws γ_j from $p(\gamma_j | \gamma_{-j}, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})$ and then draws β_j from $p(\beta_j | \gamma_j = s, \gamma_{-j}, \lambda, \mathbf{x}, \theta, \mathbf{y})$. This scheme actually draws jointly for (γ_j, β_j) . This blocking procedure can be viewed as equivalent to that used by GEWEKE (1996).

A Metropolis-Hastings version of the above blocking procedure can be easily designed. Assume that the j th element of γ is proposed with probability $q(\gamma'; \gamma)$. The acceptance ratio is given by

$$r = \frac{p(\gamma_j = s | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})}{p(\gamma_j = 1 - s | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})} \cdot \frac{q(\gamma; \gamma')}{q(\gamma'; \gamma)}, \quad (21)$$

where $\gamma = (\gamma_j = 1 - s, \gamma_{-j})$, $\gamma' = (\gamma_j = s, \gamma_{-j})$, and $s = 1$ or 0 corresponding to adding or deleting one QTL, respectively.

Using the identity $p(\gamma_j = s | \gamma_{-j}, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y}) = p(\gamma_j = s, \beta_j | \gamma_{-j}, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y}) / p(\beta_j | \gamma_j = s, \gamma_{-j}, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})$, the acceptance ratio (21) for $s = 1$ then becomes

$$\begin{aligned} r &= \frac{p(\gamma_j = 1, \beta_j | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})}{p(\gamma_j = 0, \beta_j | \gamma_{-j}, \lambda, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})} \cdot \frac{q(\gamma; \gamma') p(\beta_j | \gamma, \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})}{q(\gamma'; \gamma) p(\beta_j | \gamma', \mathbf{x}, \theta_{-\beta_j}, \mathbf{y})} \\ &= \frac{p(\mathbf{y} | \gamma', \mathbf{x}_{\gamma'}, \theta_{\gamma'}) p(\gamma') p(\beta_j | \gamma')}{p(\mathbf{y} | \gamma, \mathbf{x}_{\gamma}, \theta_{\gamma}) p(\gamma) p(\beta_j | \gamma)} \cdot \frac{q(\gamma; \gamma') p(\beta_j | \gamma, \mathbf{x}, \theta_{\gamma}, \mathbf{y})}{q(\gamma'; \gamma) p(\beta_j | \gamma', \mathbf{x}, \theta_{\gamma}, \mathbf{y})} \\ &= \frac{p(\mathbf{y} | \gamma', \mathbf{x}_{\gamma'}, \theta_{\gamma'}) p(\gamma') p(\beta_j | \gamma')}{p(\mathbf{y} | \gamma, \mathbf{x}_{\gamma}, \theta_{\gamma}) p(\gamma)} \cdot \frac{q(\gamma; \gamma')}{q(\gamma'; \gamma) p(\beta_j | \gamma', \mathbf{x}, \theta_{\gamma}, \mathbf{y})}, \end{aligned} \quad (22)$$

where $\gamma = (\gamma_j = 0, \gamma_{-j})$ and $\gamma' = (\gamma_j = 1, \gamma_{-j})$.

This M-H algorithm is equivalent to the reversible-jump algorithm, which proceeds to generate a new location and the genotype indicators at the new location from the priors and the associated effects β_j from the full conditional posterior, $p(\beta_j | \gamma', \mathbf{x}, \theta_{\gamma}, \mathbf{y})$. This reversible-jump algorithm is similar to that developed by YI and XU (2001, 2002).

Method III: Under the linear model (2), in fact, all parameters θ can be integrated out from the conditional posterior distribution (18), *i.e.*,

$$p(\gamma_j = s | \gamma_{-j}, \mathbf{x}, \mathbf{y}) \propto p(\mathbf{y} | \gamma_j = s, \gamma_{-j}, \mathbf{x}_{\gamma}) p(\gamma_j = s, \gamma_{-j}). \quad (23)$$

Therefore, γ can be updated independent of θ . This method is equivalent to that by SMITH and KOHN (1996) for the conventional linear regression model. BROMAN and SPEED (2002) have applied this method to marker selection in a backcross design. This approach is equivalent to a blocking scheme, which first draws θ_{γ} from the full conditional posterior $p(\theta_{\gamma} | \gamma_j = 1, \gamma_{-j}, \mathbf{x}, \mathbf{y})$ and then draws γ_j from $p(\gamma_j = s | \gamma_{-j}, \lambda, \mathbf{x}, \theta, \mathbf{y})$. This equivalence can be seen more explicitly from the following M-H version.

Assume that a proposal for a move from $\gamma = (\gamma_j = 1 - s, \gamma_{-j})$ to $\gamma' = (\gamma_j = s, \gamma_{-j})$ with probability $q(\gamma'; \gamma)$; the acceptance ratio is given, using the standard M-H procedure, by

$$\begin{aligned} r &= \frac{p(\gamma_j = s | \gamma_{-j}, \lambda, \mathbf{x}, \mathbf{y})}{p(\gamma_j = 1 - s | \gamma_{-j}, \lambda, \mathbf{x}, \mathbf{y})} \cdot \frac{q(\gamma; \gamma')}{q(\gamma'; \gamma)} \\ &= \frac{p(\mathbf{y} | \gamma', \mathbf{x}_{\gamma'}) p(\gamma')}{p(\mathbf{y} | \gamma, \mathbf{x}_{\gamma}) p(\gamma)} \cdot \frac{q(\gamma; \gamma')}{q(\gamma'; \gamma)}. \end{aligned} \quad (24)$$

Using the identity $p(\gamma_j = s | \gamma_{-j}, \mathbf{x}, \mathbf{y}) = p(\gamma_j = s, \theta_{\gamma} | \gamma_{-j}, \mathbf{x}, \mathbf{y}) / p(\theta_{\gamma} | \gamma, \mathbf{x}, \mathbf{y})$, the acceptance ratio then becomes

$$\begin{aligned} r &= \frac{p(\gamma_j = s, \theta_{\gamma} | \gamma_{-j}, \lambda, \mathbf{x}, \mathbf{y})}{p(\gamma_j = 1 - s, \theta_{\gamma'} | \gamma_{-j}, \lambda, \mathbf{x}, \mathbf{y})} \cdot \frac{q(\gamma; \gamma') p(\theta_{\gamma} | \gamma, \mathbf{x}_{\gamma}, \mathbf{y})}{q(\gamma'; \gamma) p(\theta_{\gamma'} | \gamma', \mathbf{x}_{\gamma'}, \mathbf{y})} \\ &= \frac{p(\mathbf{y} | \gamma', \mathbf{x}_{\gamma'}, \theta_{\gamma'}) p(\gamma') p(\theta_{\gamma'} | \gamma')}{p(\mathbf{y} | \gamma, \mathbf{x}_{\gamma}, \theta_{\gamma}) p(\gamma) p(\theta_{\gamma} | \gamma)} \cdot \frac{q(\gamma; \gamma') p(\theta_{\gamma} | \gamma, \mathbf{x}_{\gamma}, \mathbf{y})}{q(\gamma'; \gamma) p(\theta_{\gamma'} | \gamma', \mathbf{x}_{\gamma'}, \mathbf{y})}. \end{aligned} \quad (25)$$

This is exactly the acceptance ratio for the reversible-jump sampler, which proceeds to generate a new location and the genotype indicators at the new location from the priors and all the associated parameters θ_{γ} from the full conditional posterior, $p(\theta_{\gamma} | \gamma, \mathbf{x}, \mathbf{y})$. Such a scheme can be viewed as equivalent to the MC³ method of RAFTERY *et al.* (1997) for linear regressions. Note that the method developed by GAFFNEY (2001), in which all the associated effects β_{γ} are sampled from $p(\beta_{\gamma} | \gamma, \mathbf{x}, \mu, \sigma^2, \mathbf{y})$, is close to the above approach.

The major difference among these three methods is the proposal distribution on the genetic effects. Proposing the new genetic effects from the prior does not seem to be the best choice. It places an extraordinary burden on prior specification for the effects (GAFFNEY 2001). In methods II and III, the corresponding parameters are integrated out from the posteriors, or equivalently the blocking strategies are used. Since the acceptance probabilities are independent of parameter values, the samplers would lead to excellent exploration of model space (GODSILL 2001). This advantage results from the use of the full conditional posterior as reversible-jump proposals. This would suggest that reversible-jump proposals should be designed to approximate as close as possible the full conditionals.

General formula and improved strategies: We can derive a general formula that includes the algorithms discussed above as special cases. Consider a proposal

from the current state of the composite space $(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$ to a new state $(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}')$ with the proposal distribution $q(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}'; \boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})$. Using the standard M-H procedure, the acceptance probability for this proposal is given by

$$\min\left(1, \frac{p(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}' | \mathbf{y}) q(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}; \boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}')}{p(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) q(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}'; \boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})}\right). \quad (26)$$

The proposal can be split into three components:

$$q(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}'; \boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}) = q_1(\boldsymbol{\gamma}'; \boldsymbol{\gamma}) q_2(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma}) p(\boldsymbol{\lambda}'_{-\gamma'}, \mathbf{x}'_{-\gamma'}, \boldsymbol{\theta}'_{-\gamma'} | \boldsymbol{\gamma}').$$

The first component q_1 proposes a move to a new model index $\boldsymbol{\gamma}'$. The second term q_2 is the proposal for the unknowns used by model $\boldsymbol{\gamma}'$. The third term is the proposal probability for the remaining unused unknowns, which is chosen to equal the pseudo-prior $p(\boldsymbol{\lambda}'_{-\gamma'}, \mathbf{x}'_{-\gamma'}, \boldsymbol{\theta}'_{-\gamma'} | \boldsymbol{\gamma}')$. The acceptance ratio then reduces to

$$\begin{aligned} r &= \frac{p(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}' | \mathbf{y}) q(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta}; \boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}')}{p(\boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) q(\boldsymbol{\gamma}', \boldsymbol{\lambda}', \mathbf{x}', \boldsymbol{\theta}'; \boldsymbol{\gamma}, \boldsymbol{\lambda}, \mathbf{x}, \boldsymbol{\theta})} \\ &= \frac{p(\boldsymbol{\gamma}', \boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'} | \mathbf{y}) p(\boldsymbol{\lambda}'_{-\gamma'}, \mathbf{x}'_{-\gamma'}, \boldsymbol{\theta}'_{-\gamma'} | \boldsymbol{\gamma}')}{p(\boldsymbol{\gamma}, \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma} | \mathbf{y}) p(\boldsymbol{\lambda}_{-\gamma}, \mathbf{x}_{-\gamma}, \boldsymbol{\theta}_{-\gamma} | \boldsymbol{\gamma})} \\ &\quad \cdot \frac{q_1(\boldsymbol{\gamma}; \boldsymbol{\gamma}') q_2(\boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma}; \boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}) p(\boldsymbol{\lambda}_{-\gamma}, \boldsymbol{\theta}_{-\gamma}, \mathbf{x}_{-\gamma} | \boldsymbol{\gamma})}{q_1(\boldsymbol{\gamma}'; \boldsymbol{\gamma}) q_2(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma}) p(\boldsymbol{\lambda}'_{-\gamma'}, \boldsymbol{\theta}'_{-\gamma'}, \mathbf{x}'_{-\gamma'} | \boldsymbol{\gamma}')} \\ &= \frac{p(\boldsymbol{\gamma}', \boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'} | \mathbf{y})}{p(\boldsymbol{\gamma}, \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma} | \mathbf{y})} \cdot \frac{q_1(\boldsymbol{\gamma}; \boldsymbol{\gamma}') q_2(\boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma}; \boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'})}{q_1(\boldsymbol{\gamma}'; \boldsymbol{\gamma}) q_2(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma})}. \end{aligned} \quad (27)$$

This is exactly the acceptance ratio for the reversible-jump sampler with the proposal distribution factored into two components, $q_1(\cdot)$ and $q_2(\cdot)$ (GREEN 1995, 2003; GODSILL 2001, 2003). This derivation of reversible jump is obtained purely from an application of the standard M-H method to fixed-dimensional composite model space. We see that the acceptance probability is independent of the value of any parameters that are unused by both models k and k' . Hence sampling of these unused unknowns is only a “conceptual” step, which need not be performed in practice. The aim of including these unused parameters is to build a fixed-dimensional model space.

The performance of the above M-H sampler is determined by the proposal distributions $q_1(\cdot)$ and $q_2(\cdot)$. The optimal choice of proposal $q_2(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma})$ should be the full conditional $p(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'} | \boldsymbol{\gamma}', \mathbf{y})$. This scheme produces an M-H sampler with the posterior model $p(\boldsymbol{\gamma} | \mathbf{y})$ as the target distribution and thus leads to excellent exploration of model space (GODSILL 2001). Unfortunately, this full conditional is not available analytically. We have to design a proposal that approximates as closely as possible the full conditional. As in all existing Bayesian mapping methods, we use three sequential

steps to propose the values of unknowns $(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'})$ and then have the factorization

$$q_2(\boldsymbol{\lambda}'_{\gamma'}, \mathbf{x}'_{\gamma'}, \boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma}, \boldsymbol{\theta}_{\gamma}) = q_{21}(\boldsymbol{\lambda}'_{\gamma'} | \boldsymbol{\gamma}') q_{22}(\mathbf{x}'_{\gamma'} | \boldsymbol{\lambda}'_{\gamma'}) q_{23}(\boldsymbol{\theta}'_{\gamma'}; \boldsymbol{\theta}_{\gamma}). \quad (28)$$

Conditional on $(\boldsymbol{\gamma}, \boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma})$, model (2) is a conventional linear model and thus q_{23} can be taken to be the full conditional posterior $p(\boldsymbol{\theta}_{\gamma} | \boldsymbol{\gamma}, \mathbf{x}, \mathbf{y})$, which results in the acceptance probability independent of $\boldsymbol{\theta}_{\gamma}$. Sampling $(\boldsymbol{\lambda}_{\gamma}, \mathbf{x}_{\gamma})$ is a special problem in QTL mapping. Therefore, performance of MCMC mapping procedures should depend highly on the specifications of q_{21} and q_{22} . In all the previous algorithms, the location λ'_j and the genotypes \mathbf{x}'_j are proposed from their priors. This sampling scheme may be suboptimal since each locus is chosen with equal probability no matter which one has weak or strong linkage evidence. Sampling λ'_j from the prior also means that the information about the j th QTL is totally lost as soon as we delete this QTL from the model; this usually causes low acceptance probability and greatly influences the mixing behavior. To improve performance of reversible jump, it may be desirable to choose a location with stronger linkage evidence. The proposal $q_{21}(\lambda'_j | \boldsymbol{\gamma}')$ then has unequal probability over the genome. LEE and THOMAS (2000) developed a method to propose a location by scanning the unoccupied regions of the entire genome for evidence of linkage of the trait residuals. Although the method of LEE and THOMAS (2000) has greatly improved the acceptance ratio, it largely increases computational load. With the composite model space approach, we are able to design an algorithm in which the values for any locus can be retained until this locus is next visited. An efficient scheme could be designed as follows: If the j th QTL was ever included in the model, the last location of this QTL is directly taken; otherwise, a new location is sampled from the prior. Clearly this easy-to-use method makes use of our current knowledge about the QTL locations and thus should improve the performance of MCMC algorithms.

DISCUSSION

Mapping multiple QTL can be viewed essentially as a problem of model selection (*e.g.*, BROMAN and SPEED 2002; SILLANPÄÄ and CORANDER 2002). A variety of Bayesian model selection procedures have been developed for conventional statistical models (see CHIPMAN *et al.* 2001; GODSILL 2001; DELLAPORTAS *et al.* 2002). Although some of these procedures, *e.g.*, reversible-jump algorithm, have been applied to map multiple QTL, others have not yet. To date, most applications of reversible jump have conducted proposals on an *ad hoc* basis. Therefore, there is a need for further methodological work on improving the reversible-jump algorithms for mapping QTL. This article presents a unified MCMC framework for mapping multiple QTL in experi-

mental designs, based on a composite space representation of the QTL model. We show that various Bayesian model selection procedures can be modified to map multiple QTL. We also demonstrate that the composite space approach leads directly to the reversible-jump algorithm. The results add to the overall understanding of the reversible-jump and the Bayesian model selection procedures for QTL mapping and lead to new classes of Bayesian mapping algorithms that combine the benefits of several different schemes within the composite model space.

The main difficulty with the existing reversible-jump algorithms is that the acceptance probability is too low. Another major challenge remains to ascertain convergence of the reversible-jump sampler and obtain a rapidly converging sampler. This is especially true in searching for multiple QTL across the entire genome. With the existing reversible-jump algorithms, the information about a QTL is totally lost as soon as we delete this QTL from the model; this greatly influences the acceptance probability and the mixing behavior. The key to the proposed composite space approach is that transdimensional problems can be considered by looking only at a fixed dimension space. As GODSILL (2003) notes, therefore, convergence properties of the reversible-jump algorithm are inherited from the standard M-H scheme on the composite model space. A further advantage is that in principle the parameters for models other than the current model can be stored and then used for an efficient proposal design when a model is revisited. In this study, we develop strategies to achieve this advantage, which may improve the MCMC procedures for mapping QTL.

The proposed unified procedure includes various Bayesian model and variable selection methods. These methods are derived from a single framework, and thus may be incorporated into a unified computer program. Future work includes full-scale simulation studies and real data analyses, which can verify the mathematical derivations involved in the theory and test the efficiency of the proposed methods. The statistical properties of the Bayesian approach rest squarely on the specification of the prior distributions on the unknowns. Substantial effort will be devoted to prior selection and investigating the robustness of the proposed methods.

In this study, we ignored gene-by-environment interactions and gene-by-gene interactions (epistatic effects). A growing number of experiments provide strong evidence of the presence of gene-by-environment interactions and epistasis for many complex traits. Thus it is important to include potential interactions into the proposed methods. Recently, Yi and XU (2002) and Yi *et al.* (2003a) developed reversible-jump algorithms for searching for complex epistatic QTL across the entire genome. The composite space sampler proposed can in principle be extended to include interactions and

thus may improve efficiency of detecting complex interacting QTL.

This work was supported by the National Institutes of Health (NIH) (NIH RO1ES09912, NIH RO1 DK056366, and NIH P30DK056336) and an Obesity-Related Pilot/Feasibility Studies grant at University of Alabama at Birmingham (528176).

LITERATURE CITED

- BALL, R. D., 2001 Bayesian methods for quantitative trait loci mapping based on model selection: approximate analysis using the Bayesian information criterion. *Genetics* **159**: 1351–1364.
- BROMAN, K. W., and T. P. SPEED, 2002 A model selection approach for identification of quantitative trait loci in experimental crosses. *J. R. Stat. Soc. B* **64**: 641–656.
- BROOKS, S. P., P. GIUDICI and G. O. ROBERTS, 2003 Efficient construction of reversible MCMC proposal distributions. *J. R. Stat. Soc. B* **65**: 3–56.
- CARLIN, B. P., and S. CHIB, 1995 Bayesian model choice via Markov chain Monte Carlo. *J. Am. Stat. Assoc.* **88**: 881–889.
- CHIPMAN, H., E. I. EDWARDS and R. E. MCCULLOCH, 2001 The practical implementation of Bayesian model selection, pp. 65–116 in *Model Selection*, edited by P. LAHIRI. Institute of Mathematical Statistics, Beachwood, OH.
- CLYDE, M. A., 1999 Bayesian model averaging and model search strategies, pp. 157–185 in *Bayesian Statistics 6*, edited by J. M. BERNARDO, J. O. BERGER, A. P. DAWID and A. F. M. SMITH. Oxford University Press, London/New York/Oxford.
- DELLAPORTAS, P., J. J. FORSTER and I. NTZOUFRAS, 2002 On Bayesian model and variable selection using MCMC. *Stat. Comput.* **12**: 27–36.
- FERNANDEZ, C., E. LEY and M. F. J. STEEL, 2001 Benchmark priors for Bayesian model averaging. *J. Econom.* **100**: 381–427.
- GAFFNEY, P. J., 2001 An efficient reversible jump Markov chain Monte Carlo approach to detect multiple loci and their effects in inbred crosses. Ph.D. Dissertation, Department of Statistics, University of Wisconsin, Madison, WI.
- GELMAN, A., J. CARLIN, H. STERN and D. RUBIN, 1995 *Bayesian Data Analysis*. Chapman & Hall, London.
- GEORGE, E. I., and D. P. FOSTER, 2000 Calibration and empirical Bayes variable selection. *Biometrika* **87**: 731–747.
- GEORGE, E. I., and R. E. MCCULLOCH, 1993 Variable selection via Gibbs sampling. *J. Am. Stat. Assoc.* **88**: 881–889.
- GEWEKE, J., 1996 Variable selection and comparison in regression, pp. 609–620 in *Bayesian Statistics 5*, edited by J. M. BERNARDO, J. O. BERGER, A. P. DAWID and A. F. M. SMITH. Oxford University Press, London/New York/Oxford.
- GODSILL, S. J., 2001 On the relationship between MCMC model uncertainty methods. *J. Comput. Graph. Stat.* **10**: 230–248.
- GODSILL, S. J., 2003 Proposal densities, and product space methods, pp. 199–203 in *Highly Structured Stochastic System*, edited by P. J. GREEN, N. L. HJORT and S. RICHARDSON. Oxford University Press, London/New York/Oxford.
- GREEN, P. J., 1995 Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**: 711–732.
- GREEN, P. J., 2003 Trans-dimensional Markov chain Monte Carlo, pp. 179–198 in *Highly Structured Stochastic System*, edited by P. J. GREEN, N. L. HJORT and S. RICHARDSON. Oxford University Press, London/New York/Oxford.
- HEATH, S. C., 1997 Markov chain Monte Carlo segregation and linkage analysis for oligogenic models. *Am. J. Hum. Genet.* **61**: 748–760.
- HOESCHELE, I., 2001 Mapping quantitative trait loci in outbred pedigrees, pp. 599–644 in *Handbook of Statistical Genetics*, edited by D. J. BALDING, M. BISHOP and C. CANNINGS. John Wiley & Sons, New York.
- JIANG, C., and Z-B. ZENG, 1997 Mapping quantitative trait loci with dominant and missing markers in various crosses from two inbred lines. *Genetica* **101**: 47–58.
- KUO, L., and B. MALLICK, 1998 Variable selection for regression models. *Sankhya. Ser. B* **60**: 65–81.

- LEE, J. K., and D. C. THOMAS, 2000 Performance of Markov chain Monte Carlo approaches for mapping genes in oligogenic models with an unknown number of loci. *Am. J. Hum. Genet.* **67**: 1232–1250.
- LYNCH, M., and B. WALSH, 1998 *Genetics and Analysis of Quantitative Traits*. Sinauer Associates, Sunderland, MA.
- NTZOUFRAS I., 1999 Aspects of Bayesian model and variable selection using MCMC. Ph.D. Dissertation, Department of Statistics, Athens University of Economics and Business, Athens, Greece.
- RAFTERY, A. E., D. MADIGAN and J. A. HOETING, 1997 Bayesian model averaging for linear regression models. *J. Am. Stat. Assoc.* **92**: 179–191.
- SATAGOPAN, J. M., and B. S. YANDELL, 1996 Estimating the number of quantitative trait loci via Bayesian model determination. Special Contributed Paper Session on Genetic Analysis of Quantitative Traits and Complex Disease, Biometric Section, Joint Statistical Meeting, Chicago.
- SILLANPÄÄ, M. J., and E. ARJAS, 1998 Bayesian mapping of multiple quantitative trait loci from incomplete inbred line cross data. *Genetics* **148**: 1373–1388.
- SILLANPÄÄ, M. J., and E. ARJAS, 1999 Bayesian mapping of multiple quantitative trait loci from incomplete outbred offspring data. *Genetics* **151**: 1605–1619.
- SILLANPÄÄ, M. J., and J. CORANDER, 2002 Model choice in gene mapping: what and why. *Trends Genet.* **18**: 301–307.
- SMITH, M., and R. KOHN, 1996 Nonparametric regression using Bayesian variable selection. *J. Econom.* **75**: 317–344.
- STEPHENS, D. A., and R. D. FISCH, 1998 Bayesian analysis of quantitative trait locus data using reversible jump Markov chain Monte Carlo. *Biometrics* **54**: 1334–1347.
- UIMARI, P., and I. HOESCHELE, 1997 Mapping linked quantitative trait loci using Bayesian method analysis and Markov chain Monte Carlo algorithms. *Genetics* **146**: 735–743.
- UIMARI, P., and M. J. SILLANPÄÄ, 2001 Bayesian oligogenic analysis of quantitative and qualitative traits in general pedigrees. *Genet. Epidemiol.* **21**: 224–242.
- XU, S., 2003 Estimating polygenic effects using markers of the entire genome. *Genetics* **163**: 789–801.
- XU, S., and N. YI, 2000 Mixed model analysis of quantitative trait loci. *Proc. Natl. Acad. Sci. USA* **97**: 14542–14547.
- YI, N., and S. XU, 2000 Bayesian mapping of quantitative trait loci for complex binary traits. *Genetics* **155**: 1391–1403.
- YI, N., and S. XU, 2001 Bayesian mapping of quantitative trait loci under complicated mating designs. *Genetics* **157**: 1759–1771.
- YI, N., and S. XU, 2002 Mapping quantitative trait loci with epistatic effects. *Genet. Res.* **79**: 185–198.
- YI, N., D. B. ALLISON and S. XU, 2003a Bayesian model choice and search strategies for mapping multiple epistatic quantitative trait loci. *Genetics* **165**: 867–883.
- YI, N., V. GEORGE and D. B. ALLISON, 2003b Stochastic search variable selection for identifying multiple quantitative trait loci. *Genetics* **164**: 1129–1138.

Communicating editor: J. B. WALSH