

## ORIGINAL ARTICLE

# Bayesian QTL mapping for multiple families derived from crossing a set of inbred lines to a reference line

T Hayashi<sup>1</sup> and H Iwata<sup>2</sup><sup>1</sup>Division of Animal Sciences, National Institute of Agrobiological Sciences, 2 Ikenodai, Tsukuba, Ibaraki, Japan and <sup>2</sup>Data Mining and Grid Research Team, National Agricultural Research Center, Kannondai, Tsukuba, Ibaraki, Japan

In some crop species, germplasm collections consisting of a large number of accessions that include traditional landraces, modern cultivars and wild species have recently been established. Such collections are regarded as useful stocks of genes for breeding programs. However, to efficiently utilize these collections for plant breeding, understanding genetic variation in agronomic traits at the QTL level between the accessions is indispensable. One effective way to extract the actual QTL information included in these collections is to perform QTL analysis jointly for multiple families derived from crossing some accessions of the collection with a single reference line such as a standard commercial variety. We

developed a Bayesian method for jointly analyzing QTL in such interconnected multiple families derived from a set of inbred lines crossed to the reference line, to detect QTL segregating between any of the inbred lines and the reference line. In this study, we considered multiple recombinant inbred lines, each of which was derived from crossing each of the inbred lines to the reference line. The method was evaluated through the use of simulated data sets for its efficiency in detecting QTL and identifying families segregating at each QTL.

*Heredity* (2009) **102**, 497–505; doi:10.1038/hdy.2009.6; published online 11 February 2009

**Keywords:** QTL mapping; MCMC estimation; multiple families; recombinant inbred lines

## Introduction

In some crop species such as maize, rice and wheat, germplasm collections that include accessions consisting of traditional landraces, modern cultivars and wild species have recently been established (Flint-Garcia *et al.*, 2005; Kojima *et al.*, 2005; Crossa *et al.*, 2007), and used for evaluation of the genetic diversity present in a species. Such collections are also useful as stocks of genes for breeding programs. Understanding the genetic variation in agronomic traits at the QTL level in collections is required to utilize these collections as breeding materials. Although association studies using some accessions sampled from a collection are a straightforward way to evaluate QTL diversity within the collection, whole genome association analysis requires the development of high-density markers that cover the whole genome and is generally prohibited by the enormous cost of developing and genotyping a large number of markers.

One effective way to extract the QTL information in a crop collection would be to utilize the segregating multiple families derived from crossing some accessions sampled from the collection to a single reference line such as a standard commercial variety for QTL mapping.

This mating design was recently adopted by Yu *et al.* (2008) to reinforce the association study in founder lines of maize, where the segregating multiple families of recombinant inbred lines (RILs) were derived from crosses between 25 diverse founders and a reference founder line. They showed that, by projecting high-density marker information from the founder lines to the RILs, more accurate association mapping was made possible in a cost-effective way using the RILs with a moderate number of the selected markers genotyped. They also showed that the effect of population structure present in the founder lines on association mapping, causing the frequent false positives, was minimized by the multiple RILs because of reshuffling of genomes between two parental lines.

This mating design is also useful for linkage-based QTL mapping to investigate the diversity of QTL affecting the agricultural traits in the germplasm collections of the crops for which whole genome association studies are unrealistic at present due to the limited availabilities of a sufficient number of SNP markers and/or high throughput genotyping systems. For the future association studies of such crops, the targets to be analyzed can be specified by QTL mapping in the multiple families and development and genotyping of SNPs can effectively be confined to the specified regions, not on a whole genome. Linkage-based QTL mapping using the multiple families derived from the founder lines can accurately identify the QTL regions with lower false-positive rate than association mapping using the founder lines in which the unknown population structure might be present although the specified regions are

Correspondence: Dr T Hayashi, Division of Animal Sciences, National Institute of Agrobiological Sciences, 2 Ikenodai, Tsukuba, Ibaraki 305-0901, Japan.

E-mail: hayatk@affrc.go.jp

Received 18 August 2008; revised 9 December 2008; accepted 23 December 2008; published online 11 February 2009

relatively broader. Moreover, with this mating design, we can detect the QTL at which any accession possesses a different allele from that of a reference line, which would provide some useful information for breeding of the crop. However, statistical methods of QTL mapping to effectively analyze a large number of multiple families with a common parental line remain to be developed.

In this paper, we develop a Bayesian method for jointly analyzing QTL for such interconnected multiple families with a common parental reference line; the families are derived by crossing a set of inbred lines (referred to as the 'tested lines' hereafter) sampled from a collection with a common parent that serves as the reference line, to detect QTL segregating between any of the tested lines and the reference line. It would be desirable to analyze as many tested lines as possible for understanding of QTL diversity in the collections. Accordingly, a large number of families, each of which is derived from each tested line crossed to a reference line, should be treated. As the number of families increases, however, each family would necessarily be confined to smaller size owing to limitations on the available space and cost, and this might decrease the accuracy in the estimation of the effect of the specific QTL allele derived from each tested line. Therefore, we treated the effects of alleles from the tested lines as random effects, but treated the effect of the allele from the reference line as a fixed effect. Here we discriminate 'random effect' and 'fixed effect' from a frequentist stand point although all effects included in the model are random in a Bayesian framework. When an effect is treated as random in a frequentist framework, a probability distribution is considered for the effect by a frequentist, which can be regarded as a prior distribution by a Bayesian. An effect with such a probability distribution provided by frequentist consideration is termed as 'random effect' whereas an effect to which no probability distribution is assigned by a frequentist is as 'fixed effect' in this Bayesian study.

Information about accessions that possess QTL alleles different from that of the reference line will be very useful in future breeding programs. We therefore incorporated a variable that indicates a segregation of each QTL in each family into the statistical model to infer which of the tested lines possess QTL alleles different from that of the reference line.

Our consideration was confined to multiple families of RILs derived from crosses between a considerable number of tested lines and a common reference line. However, the statistical model would easily be applicable to other families such as  $F_2$  or backcross with slight modification. The method was then evaluated for its efficiency in detecting QTL and identifying families that segregate for each QTL using simulated data sets.

## Materials and methods

### Analyzed families

We consider multiple families of RILs derived from crosses between a considerable number of tested lines and a common reference line, where the tested lines are crossed to the reference line, followed by selfing, to generate segregating  $F_2$  populations, from each of which RILs are derived through single-seed descent with repeated cycles of selfing. The number of families, equal

to the number of tested lines used for crossing with the reference line, is  $m$  and the size of the  $i$ th RIL family derived from the  $i$ th tested line is  $n_i$ .

### Statistical model

We assume that observations of the phenotype of a trait are available for individuals in the multiple families of RILs, as is marker information, including genotypic data at markers for the tested lines, the reference line and individuals in the multiple families of RILs and a linkage map of the markers, where all individuals in the RILs are assumed to be homozygous at all QTL and markers. We assumed that there is no epistatic interaction between QTL in this study although this assumption can be relaxed without difficulty. The phenotypic value of the  $j$ th individual in the  $i$ th RIL family is denoted by  $y_{ij}$  ( $i = 1, 2, \dots, m; j = 1, 2, \dots, n_i$ ), for which we can apply the following linear model,

$$y_{ij} = \mu + \sum_{l=1}^N \{u_{lij}a_l + (1 - u_{lij})b_l\} + e_{ij}. \quad (1)$$

In this model,  $\mu$  is the intercept of the model, which is a mean of the genotypic values obtained by omitting segregating QTL in the multiple families and  $N$  is the number of QTL affecting the phenotypic value. The variable  $u_{lij}$  indicates the genotype of the individual at the  $l$ th QTL, where the alleles at the QTL are denoted by  $Q_l$  and  $q_{li}$  for the reference line and the  $i$ th tested line, respectively, and  $u_{lij} = 1$  for the genotype  $Q_lQ_l$  and 0 for  $q_{li}q_{li}$ . The genotypic contributions of the QTL corresponding to  $Q_lQ_l$  and  $q_{li}q_{li}$  are denoted by  $a_l$  and  $b_{li}$ , respectively, and  $e_{ij}$  is the residual error following a normal distribution with mean 0 and variance  $\sigma_e^2$ . In multiple families that share a single parental line (the reference line), the effects of QTL alleles derived from the reference line are well estimated by a large number of degrees of freedom allocated for the estimation, but the instability in the estimation of the effects of the alleles from each tested line might be caused by the limited size of each family. For the QTL effects, we thus treated  $a_l$  as a fixed effect and  $b_{li}$  as a random effect sampled from a normal distribution with mean 0 and variance  $\sigma_{b_{li}}^2$ . It is noted that the variance of  $b_{li}$  is indexed by ' $l$ ' because the QTL effect has a specific distribution for each QTL.

Moreover, we incorporate a variable,  $s_{li}$ , that indicates whether each QTL is segregating or not in each family, where  $s_{li} = 1$  if the  $l$ th QTL is segregating in the  $i$ th family and  $s_{li} = 0$  otherwise. When  $s_{li} = 0$ , the  $i$ th tested line has the same allele at the  $l$ th QTL as the reference line; accordingly, the genotypic values at the  $l$ th QTL are expressed as  $a_l$  for all individuals in the  $i$ th family. Denoting the genotypic contribution from the  $l$ th QTL to the phenotypic value by  $\Delta_{lij}$ , we can write  $\Delta_{lij} = s_{li}\{u_{lij}a_l + (1 - u_{lij})b_{li}\} + (1 - s_{li})a_l$ . Therefore, considering a segregation variable  $s_{li}$ , model (1) can be modified as

$$y_{ij} = \mu + \sum_{l=1}^N \Delta_{lij} + e_{ij}. \quad (2)$$

### Prior and posterior distributions of parameters and variables

The parameters and variables included in model (2) and the locations of  $N$  QTL, denoted as  $\lambda_1, \lambda_2, \dots, \lambda_N$ , are

collectively written as  $\theta$  and are referred to as unobservables. The observed phenotypic values are denoted by  $y = \{y_{ij}\}$  for  $i = 1, 2, \dots, m$  and  $j = 1, 2, \dots, n_i$  for each  $i$ . The likelihood is written as

$$f(\mathbf{y}|\theta) = \prod_{i=1}^m \prod_{j=1}^{n_i} f(y_{ij}|\theta) \propto (\sigma_e^2)^{-n_T/2} \times \exp\left\{-\frac{1}{2\sigma_e^2} \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \mu - \sum_{l=1}^N \Delta_{lij})^2\right\}$$

where  $n_T = \sum_{i=1}^m n_i$  is total number of individuals in the whole families. Denoting the joint prior distribution and the joint posterior distribution of  $\theta$  with  $p(\theta)$  and a  $p(\theta|\mathbf{y})$ , respectively, we can write

$$p(\theta|\mathbf{y}) \propto f(\mathbf{y}|\theta)p(\theta) = f(\mathbf{y}|\theta)p(\mu)p(\sigma_e^2) \times \prod_{l=1}^N p(a_l) \times \prod_{l=1}^N \prod_{i=1}^m p(b_{li}|\sigma_{bl}^2, s_{li}) \times \prod_{l=1}^N p(\sigma_{bl}^2|s_{l1}, s_{l2}, \dots, s_{lm}) \times \prod_{l=1}^N \prod_{i=1}^m p(s_{li}) \times \prod_{l=1}^N \prod_{i=1}^m \prod_{j=1}^{n_i} p(u_{lij}|\lambda_i) \times \prod_{l=1}^N p(\lambda_l) \times p(N)$$

where  $p(\mu)$ ,  $p(\sigma_e^2)$ ,  $p(a_l)$ ,  $p(b_{li}|\sigma_{bl}^2, s_{li})$ ,  $p(\sigma_{bl}^2|s_{l1}, s_{l2}, \dots, s_{lm})$ ,  $p(s_{li})$ ,  $p(u_{lij}|\lambda_i)$ ,  $p(\lambda_i)$  and  $p(N)$  are the priors of components of  $\theta$ . For  $\mu$ ,  $\sigma_e^2$  and  $a_l$ , we chose the following prior distributions,  $p(\mu) \propto 1$ ,  $p(\sigma_e^2) \propto 1/\sigma_e^2$  and  $p(a_l) \propto 1$ .

It should be noted that  $b_{li}$  is not included in the likelihood  $f(\mathbf{y}|\theta)$  when  $s_{li} = 0$ , meaning that the  $l$ th QTL is not segregating in the  $i$ th family, whereas  $b_{li}$  is included in the likelihood  $f(\mathbf{y}|\theta)$  when  $s_{li} = 1$ . Therefore, the full conditional posterior distribution of  $b_{li}$  is independent of the data  $\mathbf{y}$  and equated to the prior  $p(b_{li}|\sigma_{bl}^2, s_{li})$  when  $s_{li} = 0$ . Such priors as  $p(b_{li}|\sigma_{bl}^2, s_{li} = 0)$  were referred to as 'pseudo-priors' by Carlin and Chib (1995) in the context of Bayesian model choice. We assumed that  $p(b_{li}|\sigma_{bl}^2, s_{li}) = \phi(b_{li}|c, d)$  for both  $s_{li} = 1$  and  $s_{li} = 0$ , where  $\phi(y|c, d)$  denotes the normal density function with mean  $c$  and variance  $d$ . For  $p(\sigma_{bl}^2|s_{li})$  which is also a pseudo-prior, we assumed  $p(\sigma_{bl}^2|s_{li}) \propto 1/\sigma_{bl}^2$  (Xu, 2003) for both  $s_{li} = 0$  and  $s_{li} = 1$  although this form of a prior of  $\sigma_{bl}^2$  leads to the improper posteriors of  $\sigma_{bl}^2$  and  $b_{li}$  (ter Braak *et al.*, 2005). We would give some consideration to the problem of improper posteriors in Discussion.

As the prior distribution of  $u_{lij}$ , we adopted the conditional probability of a QTL genotype given linked marker genotypes near the QTL location as described by Jiang and Zeng (1997) for a biparental  $F_1$  population. The prior probabilities of  $s_{li} = 0$  and 1 were given as 0.5 for QTL segregation. The prior distribution of  $\lambda_i$  is assumed uniform across the whole chromosomal region. The prior probability of  $N$  was a Poisson distribution with a pre-specified mean  $\delta$ . In the following simulation experiments, we assumed that  $\delta = 2$ .

We estimate  $\theta$  by using a Markov chain Monte Carlo (MCMC) algorithm. After the initial values are given to  $\theta$ , MCMC cycles are repeated for updating the values. A Gibbs sampling scheme is applied to the update of  $\theta$  except for  $N$  and  $\lambda_i$  ( $i = 1, 2, \dots, N$ ), which are updated based on Metropolis-Hastings algorithm (Metropolis *et al.*, 1953; Hastings, 1970) including a reversible-jump

MCMC (RJ-MCMC) sampling (Green, 1995) for  $N$ . Details of the updating process for  $\theta$  are given in the Appendix A.

### Simulation experiments

**Simulation settings:** We evaluated the proposed Bayesian method for the efficiency of detecting QTL segregating in any family and identifying families segregating at each QTL with the analyses of simulated data sets. We considered multiple  $F_8$  families, where a set of tested lines were crossed to a reference line to generate multiple  $F_2$  families from each of which  $F_8$  families were derived through single-seed descent with repeated cycles of selfing. In our simulation, we assumed that the family size, denoted by  $n$ , is equal for all families (that is,  $n = n_1 = n_2 = \dots = n_m$ ). We assumed three combinations for the number of families,  $m$ , and family size,  $n$ , as  $(m, n) = (50, 40)$ ,  $(100, 20)$  and  $(200, 10)$  with total number of individuals in all families fixed as 2000.

The simulated genome consisted of four chromosomes, Chr1, Chr2, Chr3 and Chr4, each of length 100 cM, on which 21 markers per chromosome were located every 5 cM. We assumed that there were five alleles with equal frequencies at each marker in the founder generation. Accordingly, each allele was randomly allocated to each marker of the reference line and tested lines with probability 0.2 in our simulations. We generated three QTL, QTL1, QTL2 and QTL3, located at 23 cM on Chr1, 72 cM on Chr2 and 12 cM on Chr3; Chr4 harbored no QTL, and was used to investigate the false-positive rate (FPR), which is described in more detail in the next section. The numbers of QTL alleles existing in all tested lines were three for QTL1, two for QTL2 and five for QTL3. We denoted the  $k$ th allele at the  $l$ th QTL as  $A_{lk}$ . We assumed that the reference line had the first allele at each QTL (that is,  $A_{11}$ ,  $A_{21}$  and  $A_{31}$ ). The allele frequency of  $A_{lk}$  was denoted by  $f_{lk}$  and the QTL effect of the homozygote with  $A_{lk}$  was denoted by  $\alpha_{lk}$ , which is referred to hereafter as the allelic effect of  $A_{lk}$ . These frequencies and effects were set to the values shown in Table 1 for our simulations. The proportions of families segregating at each QTL, which were derived from the tested lines possessing the alleles other than  $A_{l1}$  at the  $l$ th QTL ( $l = 1, 2, 3$ ) were assumed as 0.3, 0.2 and 0.8 for QTL1, QTL2 and QTL3, respectively, as shown in Table 1. When generating each data set, the QTL alleles were randomly allocated to  $m$  tested lines such that the allele frequencies were those given in Table 1, where the allele allocation in the tested lines was recorded and used for summarizing the results of simulation analyses. In addition, we considered 10 unlinked biallelic additive QTL, each with equal frequency of two homozygous in founder lines and with effects of 0.1 and  $-0.1$  for two homozygous, to include polygenic effects whose variances were summed to be 0.1. The phenotypic values of individuals in the  $F_8$  generation were determined by the sum of genotypic effects corresponding to the genotypes at the three QTL and 10 unlinked QTL and environmental effects sampled from a normal distribution with mean 0 and variance 1. The proportion of phenotypic variance explained by each QTL (referred to as PVQ) was also shown in Table 1, as this might affect the power of detecting each QTL.

We generated 100 data sets for each of the three settings for  $(m, n)$ . The power of the QTL detection and the accuracy

**Table 1** Summary of QTL and parameter values used to generate simulated data sets of multiple families

QTL	Location	Allele number	Allele frequencies <sup>a</sup>	Effect	PVQ <sup>b</sup>
QTL1	23 cM on Chr1	3	$f_{11} = 0.7, f_{12} = 0.1, f_{13} = 0.2$	$\alpha_{11} = 0.0, \alpha_{12} = -0.6, \alpha_{13} = 0.3$	0.022
QTL2	72 cM on Chr2	2	$f_{21} = 0.8, f_{22} = 0.2$	$\alpha_{21} = -0.8, \alpha_{22} = 0.0$	0.046
QTL3	12 cM on Chr3	5	$f_{31} = f_{32} = f_{33} = f_{34} = f_{35} = 0.2$	$\alpha_{31} = -0.2, \alpha_{32} = 0.2, \alpha_{33} = -0.6, \alpha_{34} = 0.0, \alpha_{35} = 0.4$	0.052

At each QTL, the first allele is assumed to be derived from a reference line such that RILs from tested lines possessing the first allele are not segregating at the QTL.

<sup>a</sup>Frequencies of alleles at the QTL in the collection of tested lines.

<sup>b</sup>Proportion of phenotypic variance explained by the QTL in the whole families. The variance explained by a QTL,  $Var(QTL)$ , is calculated as:  $Var(QTL) = E_{family}(Var(QTL | family)) + Var_{family}(E(QTL | family))$ , where  $Var(QTL | family)$  and  $E(QTL | family)$  are the conditional variance and expectation of QTL effect given the alleles existing in each family and  $E_{family}(\cdot)$  and  $Var_{family}(\cdot)$  indicate the expectation and variance taken over all families.

We can calculate  $\sigma_{bi}^2$ , which is the variance of the effect of the allele derived from a tested line, for QTL1, QTL2 and QTL3 as 0.054, 0.1024 and 0.1184, respectively, with the allele frequencies in the collection of tested lines and allelic effects.

in identifying the families that were segregating at each QTL were evaluated through analyses of the 100 data sets for each setting of  $(m, n)$ . For comparison, the same data sets were also analyzed using a method based on interval mapping for multiple families proposed by Xu (1998), referred to as IM, treating QTL effect as a random effect due to the large number of families (that is,  $m \geq 50$ ). Moreover, to evaluate the incremental efficiency obtained by incorporating a segregation variable,  $s_{ij}$ , we applied an additional Bayesian method based on model (1), without consideration of the segregation variable, to the analyses of simulated data sets. Hereafter, the Bayesian methods based on models (1) and (2) are referred to as Bayes1 and Bayes2, respectively. For each of the Bayesian methods, we performed 50 000 cycles of MCMC and sampled the values of the unobservables every 20 cycles during the last 40 000 cycles with the first 10 000 cycles discarded as burn-in.

In our Bayesian methods, the posterior QTL intensity (Sillanpää and Arjas, 1998) for each small interval with 1 cM length on the genome was calculated for QTL detection. We obtained a summed QTL intensity, referred to as *SQI* (Hayashi and Awata, 2008), by summing the posterior QTL intensity over all intervals on each chromosome, and used *SQI* as a test statistic for detecting QTL on a chromosome. Thresholds of *SQI* were determined from the empirical null distributions of the maximum of *SQI* over all chromosomes obtained by analyses of 100 null data sets that were generated on the assumption of no QTL in each setting of  $(m, n)$ . The empirical null distributions of the maximum of *SQI* over all chromosomes were established by analyzing 100 null data sets for Bayes1 and Bayes2. The values of maximized *SQI* corresponding to 5% significant level of the empirical null distributions were regarded as the thresholds for *SQI*. When *SQI* exceeded the thresholds for any chromosome, detection of a QTL on the chromosome is declared. The Bayesian estimates of the positions and effects of the detected QTL were given in the analysis of each data set as described in Hayashi and Awata (2008), where the positions and effects of the QTL fitted in the model were averaged over the chromosome, with the QTL intensity of intervals that harbored the QTL used as a weight. Such a weighted average for the posterior probabilities of QTL segregation in each family (that is,  $s_{ij} = 1$ ) was also considered to identify the families that were segregating at QTL in Bayes2.

In IM, the likelihood-ratio test statistic (*LRT*) was adopted for QTL detection. Thresholds for *LRT* were

determined similarly to the approach used for *SQI*. In IM, the position of the peak of *LRT* was regarded as the estimated QTL position.

**Results of simulation experiments:** Table 2 shows the powers of QTL detection and the estimates of the QTL position and effect of allele from a reference line at each QTL for Bayes1 and Bayes2 as well as IM in which the estimated of QTL effects were not given as variances of QTL effects were treated in IM with a random effect model (Xu, 1998). The averages and s.d. for the estimated QTL positions and QTL effects were calculated over the repetitions that successfully detected the QTL. In the simulation, Chr4, which harbored no QTL, was used to evaluate FPR, for QTL detection, where FPR was defined as the number of repetitions that falsely detected a QTL on Chr4 in the analyses of 100 data sets. For  $(m, n) = (50, 40)$ ,  $(100, 20)$  and  $(200, 10)$ , the respective FPRs were 1, 2 and 2 in IM; 1, 2 and 3 in Bayes1; and 2, 2 and 2 in Bayes2. Therefore, the thresholds corresponding to the genome-wide 5% significance level empirically determined by the analyses of 100 null data sets appropriately controlled the FPR for all three methods, such that the powers of these methods were suitably compared.

The powers of QTL detection were decreased as the number of families ( $m$ ) was increased with family size ( $n$ ) decreased in all three methods. At a given  $(m, n)$ , the Bayesian methods showed higher powers of detecting QTL than IM whereas the powers were comparable between Bayes1 and Bayes2. The powers of detection for QTL1 were much lower than those for QTL2 and QTL3, which were 38 and 41% for  $(m, n) = (50, 40)$  and decreased to 26% for  $(m, n) = (100, 20)$  and to 6 and 14% for  $(m, n) = (200, 10)$  with Bayes1 and Bayes2, respectively. As shown in Table 1, *PVQ* of QTL1 was considerably smaller than that of the other QTL, and this was responsible for the poor powers for QTL1. For QTL2 and QTL3 with moderate *PVQ* values, both Bayesian methods showed higher powers than IM; powers were higher than 80% at  $(m, n) = (200, 10)$  and increased to about 95% at  $(m, n) = (50, 40)$  with the Bayesian methods.

The estimates of the positions were slightly biased for QTL2 and QTL3, but were noticeably biased for QTL1 at  $(m, n) = (200, 10)$  in the Bayesian methods. The estimates of the effects of the alleles from the reference line obtained with Bayesian methods were considerably biased for QTL2 and QTL3. For example, the simulated effect of QTL2 was  $-0.8$  (Table 1), but the estimates were

**Table 2** Results of simulation experiments

	Methods <sup>a</sup>								
	IM			Bayes1			Bayes2		
	QTL1	QTL2	QTL3	QTL1	QTL2	QTL3	QTL1	QTL2	QTL3
<i>m</i> = 50, <i>n</i> = 40									
Power <sup>b</sup>	20	80	84	38	95	97	41	97	96
Position <sup>c</sup>	21.9 (7.1)	73.2 (5.6)	12.0 (6.5)	30.3 (15.9)	71.3 (4.8)	12.9 (5.7)	28.3 (12.8)	69.3 (5.6)	14.9 (6.6)
Effect <sup>d</sup>	—	—	—	0.01 (0.05)	−0.17 (0.05)	−0.17 (0.05)	0.05 (0.10)	−0.42 (0.12)	−0.32 (0.15)
<i>m</i> = 100, <i>n</i> = 20									
Power	12	57	67	26	88	95	26	95	99
Position	25.3 (8.6)	73.6 (8.0)	12.4 (5.3)	20.1 (9.7)	71.2 (6.5)	13.1 (5.3)	22.2 (7.0)	70.3 (5.3)	15.3 (5.6)
Effect	—	—	—	−0.02 (0.06)	−0.16 (0.05)	−0.17 (0.05)	0.01 (0.13)	−0.36 (0.11)	−0.35 (0.16)
<i>m</i> = 200, <i>n</i> = 10									
Power	10	35	50	9	82	90	14	88	87
Position	26.5 (13.0)	71.6 (9.2)	11.7 (4.4)	32.8 (22.4)	70.6 (6.5)	12.8 (5.1)	32.8 (15.5)	69.4 (6.2)	15.4 (8.5)
Effect	—	—	—	−0.01 (0.04)	0.01 (0.33)	−0.19 (0.31)	−0.03 (0.11)	−0.12 (0.46)	−0.26 (0.50)

In each of the settings of (*m,n*) = (50,40), (100,20) and (200,10), 100 data sets were simulated and analyzed.

<sup>a</sup>IM, interval mapping with random effect model; Bayes1, Bayesian method based on model (1) without variables indicating QTL segregation in each family; Bayes2, Bayesian method based on model (2) taking variables for QTL segregation into consideration.

<sup>b</sup>Number of replicates that successfully detected each QTL with a genome-wide significance level of *P* = 0.05 in a total of 100 replicates.

<sup>c</sup>Average of the estimated QTL positions calculated over the replicates that successfully detected the QTL, with s.d. given in the parenthesis.

<sup>d</sup>Average of the estimated QTL effects of alleles derived from a reference line over the replicates that successfully detected the QTL, with s.d. given in the parenthesis. These estimates were not obtained in IM as the variance of QTL effects were treated by IM with a random effect model.

shrunk towards zero (Table 2). This shrinkage was less in Bayes2 than in Bayes1; that is, in Bayes2, the estimated values were closer to the simulated values given in Table 1. Bayes2, however, provided biased estimates of the effects of QTL3 for (*m,n*) = (50,40) and (100,20), where the simulated effect was given as −0.2, but the respective estimates were inflated to −0.32 and −0.37, respectively (Table 2).

In Table 3, we have summarized the inferences about QTL segregation in the families for the analyses with Bayes2. In the analysis of each simulated data set, we obtained a posterior probability of QTL segregation (i.e., *s<sub>li</sub>* = 1) for each family and averaged the probabilities over the families derived from tested lines that possessed identical alleles at each QTL. We further averaged the probabilities over the repetitions with successful detection of the QTL in each setting of (*m,n*) and the results are listed for each allele in the rows labeled ‘Probability of segregation’ of Table 3.

In addition, to evaluate the ability of Bayes2 to identify tested lines with alleles that differ from that of the reference line, resulting in QTL segregation in the families derived from the tested lines crossed with the reference line, we investigated the proportions of the tested lines with the posterior probabilities of QTL segregation exceeding two pre-determined values 0.6 and 0.9 for each QTL allele in the replications with successful QTL detection (Table 3). For example, consider QTL1 at a setting of (*m,n*) = (50,40). There were 35 lines with allele A<sub>11</sub>, five lines with A<sub>12</sub> and 10 lines with A<sub>13</sub> in each replication, given allele frequencies 0.7, 0.1 and 0.2 for A<sub>11</sub>, A<sub>12</sub> and A<sub>13</sub>, respectively, as given in Table 1. Therefore, the total numbers of the tested lines with A<sub>11</sub>, A<sub>12</sub> and A<sub>13</sub> investigated in 41 replications that successfully detected QTL1 were 1435, 205 and 410, respectively, from which we obtained the numbers of tested lines with posterior probability for QTL segregation exceeding 0.6

as 215, 125 and 98, with proportions 15, 61 and 23%, respectively, as listed in Table 3. Similarly, the proportions of tested lines with posterior probability of QTL segregation exceeding 0.9 were 1, 20 and 3% for the lines with alleles A<sub>11</sub>, A<sub>12</sub> and A<sub>13</sub>, respectively. As the first alleles at three QTL (A<sub>11</sub>, A<sub>21</sub> and A<sub>31</sub>) were allocated to the reference line in our simulations, the proportions of the tested lines possessing these QTL alleles with posterior probabilities of QTL segregation greater than 0.6 or 0.9 were regarded as the false discovery rates for QTL segregation in the non-segregating families derived from the tested lines. For the tested lines with QTL alleles that differed from those of the reference line, the proportions indicated the capability of correct identification for the families, derived from the tested lines, which segregated for the QTL. The accuracies of inference for segregating tested lines were enhanced as the effects of alleles or the family size increased (Table 3).

## Discussion

### The efficiencies of the Bayesian methods in analyzing simulated data sets

As shown in Table 2, the powers of QTL detection were greater for both Bayesian methods than IM, indicating the possibility that information on QTL that distinguishes a reference line relative to the tested lines might be effectively elucidated by the Bayesian methods using the experimental design adopted in this study. Especially, in a setting of (*m,n*) = (200,10), Bayes2 method showed noticeably higher powers (≥87%) for the detection of QTL2 and QTL3 than IM (35 and 50%).

However, the estimates of the QTL effects of alleles from the reference line were biased in the Bayesian methods (Table 2). This might have been caused by inaccuracies in the inference about QTL segregation in

Table 3 Inference for QTL segregation in each family using Bayes2 in simulation experiments

	QTL1			QTL2			QTL3				
	$A_{11}^a$	$A_{12}$	$A_{13}$	$A_{21}$	$A_{22}$	$A_{31}$	$A_{32}$	$A_{33}$	$A_{34}$	$A_{35}$	
$m = 50, n = 40$											
Probability of segregation <sup>b</sup>	0.46 (0.02)	0.69 (0.04)	0.49 (0.03)	0.35 (0.04)	0.85 (0.02)	0.38 (0.03)	0.58 (0.04)	0.30 (0.03)	0.47 (0.05)	0.69 (0.04)	
Proportion of segregating lines (0.6) <sup>c</sup>	0.15	0.61	0.23	0.13	0.92	0.14	0.50	0.06	0.28	0.72	
Proportion of segregating lines (0.9) <sup>d</sup>	0.01	0.20	0.03	0.02	0.55	0.00	0.06	0.00	0.03	0.16	
$m = 100, n = 20$											
Probability of segregation	0.48 (0.01)	0.61 (0.03)	0.51 (0.02)	0.42 (0.03)	0.75 (0.03)	0.40 (0.03)	0.56 (0.04)	0.33 (0.03)	0.48 (0.03)	0.65 (0.03)	
Proportion of segregating lines (0.6)	0.13	0.45	0.21	0.14	0.78	0.13	0.40	0.07	0.26	0.61	
Proportion of segregating lines (0.9)	0.01	0.10	0.02	0.01	0.25	0.00	0.03	0.00	0.01	0.07	
$m = 200, n = 10$											
Probability of segregation	0.49 (0.01)	0.54 (0.02)	0.51 (0.01)	0.48 (0.03)	0.57 (0.05)	0.48 (0.04)	0.53 (0.04)	0.44 (0.06)	0.51 (0.04)	0.56 (0.05)	
Proportion of segregating lines (0.6)	0.10	0.24	0.17	0.22	0.49	0.27	0.36	0.25	0.31	0.43	
Proportion of segregating lines (0.9)	0.00	0.04	0.00	0.02	0.04	0.02	0.04	0.03	0.02	0.07	

<sup>a</sup>QTL allele is denoted by  $A_{ik}$  for the  $k$ th allele at the  $l$ th QTL.

<sup>b</sup>Averaged posterior probabilities for QTL segregation in the families from tested lines possessing each allele, which are further averaged over the repetitions that successfully detected the QTL and listed with standard deviations given in the parenthesis. As the first alleles at QTL,  $A_{11}$ , were possessed by a reference line, these probabilities indicate the rate of false discovery of QTL segregation for the non-segregating families from tested lines that possessed  $A_{11}$ .

<sup>c</sup>Proportion of the tested lines with each QTL allele inferred to be segregating at the QTL in the families derived from the tested lines on the criterion of the posterior probability of QTL segregation  $> 0.6$  in the repetitions with successful QTL detection.

<sup>d</sup>Proportion of the tested lines with each QTL allele inferred to be segregating at QTL in the families derived from the tested lines on the criterion of the posterior probability of segregation  $> 0.9$  in the repetitions with successful QTL detection.

each family. In Bayes1, as the inference about QTL segregation in each family was not incorporated to the analyses, the estimates of alleles from the reference line were shrunk to zero. In Bayes2 which could infer QTL segregation in each family using a variable indicating QTL segregation, the accuracies in the inference were varied depending on the effects of QTL alleles and the combinations of the number of families and each family size (Table 3), as was the accuracies in the estimation of QTL effects. For example, at QTL3, the QTL segregation in the families from the tested lines with allele  $A_{33}$  were frequently undiscoverable, where the posterior probabilities of segregation were only 0.30, 0.33 and 0.44 for  $(m,n) = (50,40)$ ,  $(100,20)$  and  $(200,10)$ , respectively. Accordingly, allele  $A_{33}$  was frequently misidentified as the allele from the reference line,  $A_{31}$ , especially, in  $(m,n) = (50,40)$  and  $(100,20)$ . Therefore, the effect of  $A_{33}$  ( $\alpha_{33} = -0.6$ , Table 1) was confounded with the effect of  $A_{31}$  ( $\alpha_{31} = -0.2$ , Table 1) causing considerable downward bias in the estimates of  $\alpha_{31}$  in  $(m,n) = (50,40)$  and  $(100,20)$ , as shown in Table 2. The accuracies in the inference about QTL segregation decreased as the number of families increased and each family size decreased owing to sampling error in segregation caused by small family size. Taking QTL2 as an example, the power of identifying QTL segregation in the families with allele  $A_{22}$  reduced as  $m$ , consequently, the estimate of  $\alpha_{21}$  was increasingly biased with increasing  $m$  (Tables 2 and 3).

The posterior probability of segregation in each family at each QTL obtained with Bayes2 can be used to identify the tested lines that have QTL alleles different from that of the reference line. As shown in Table 3, at  $(m,n) = (50,40)$ , tested lines that had QTL alleles with effects greatly different from the QTL alleles of the reference line were efficiently identified. For example, the power of correctly identifying the segregation at QTL3 in the families derived from the tested lines having alleles  $A_{35}$  was 72% based on the criterion of the posterior probability of segregation greater than 0.6. In this criterion, however, the false discovery of segregation at QTL3 in non-segregating families, which were derived from tested lines with allele  $A_{31}$ , occurred at a rate of 14%. Increasing the threshold for the posterior probability of segregation to 0.9 decreased the power of correct identification of segregating families to 16% for the allele  $A_{35}$  at QTL3, but the rate of false discovery of QTL segregation for non-segregating families was negligible (Table 3). Using the threshold of 0.9 for the posterior probability of QTL segregation in  $(m,n) = (50,40)$ , tested lines with  $A_{22}$  at QTL2 were still correctly identified with 55% as having a different allele from that of the reference line. Therefore, Bayes2 showed a practical capability to identify tested lines with QTL alleles different from that of a reference line in  $(m,n) = (50,40)$  although the rates of successful identification for the segregating families were lower at settings of  $(m,n) = (100,20)$  and  $(200,10)$ , as shown in Table 3.

#### MCMC algorithm in Bayesian model selection for multiple families

The dimensionality of the parameters in the models for QTL mapping changes depending on the number of QTL included. Although effective sampling schemes based on Gibbs sampling, such as stochastic search variable

selection (SSVS) (Yi *et al.*, 2003) and Bayesian shrinkage estimation (Xu, 2003), have recently been proposed, we adopted RJ-MCMC for the inference of the number of QTL,  $N$ , in this study. A Gibbs sampling scheme for model selection can only be performed over a composite model that is a product space of candidate models and their parameters (Godsill, 2001). The model for multiple families considered in the present study is determined not only by the number of QTL but also by the configurations of QTL alleles in the tested lines in contrast to a model for the QTL analysis in a biparental cross family, where the model is simply determined by the QTL number. Accordingly, the number of possible models becomes intractably large in the analysis of multiple families designed in this study, in which the composite model space is difficult to be dealt with for Gibbs sampling schemes. Therefore, we chose the RJ-MCMC sampling for estimation of the QTL number and a random model approach was introduced to cope with the enormous number of possible configurations of alleles in the tested lines at each QTL for the Bayesian estimation.

In the Bayesian method proposed by Xu (2003), the effects of QTL assumed at each position in a genome were treated as random effects, the priors of which were normal distributions with mean zero and different variances for different QTL. In our study, the priors of the effects of QTL alleles from tested lines were also assumed as normal distributions with mean zero and different variances for different QTL. In Bayes2, in addition, we incorporated a binary variable indicating QTL segregation in each family at each QTL, which can be regarded as analogous to the indicator variable for the presence of a QTL at each genome position used by Yi *et al.* (2003) for SSVS. Although Jannink and Wu (2003) applied RJ-MCMC for the inference about allele configurations in multiple interconnected families, accurate estimation of the allele configurations in a large number of families with moderate to small family sizes would be difficult as the number of possible configurations becomes enormously large owing to the increase in the number of potential alleles. Moreover, the difference in alleles between the tested lines can only be indirectly inferred in the multiple families, considered in the present study, through a single reference line shared by the families, which would make suitable configuration of alleles in the families more difficult.

As shown in the simulation experiments (Tables 2 and 3), Bayes2 might be a practical method to detect QTL segregating between a reference line and tested lines and to allow the inference about QTL segregation in each family unless the family size is too small, as in the setting of  $(m,n) = (50,40)$  in simulations. Slower convergence and a poorer mixing property of RJ-MCMC compared with Gibbs sampling would be compensated for to some extent by increasing the iterations, which is possible for the high-performance computers that are now available without requiring excessive computational time.

For the prior of the variance  $\sigma_{b_i}^2$  of QTL effects  $b_i$  of the alleles from the tested lines, we adopted  $p(\sigma_{b_i}^2 | s_i) \propto 1/\sigma_{b_i}^2$ . As shown by ter Braak *et al.* (2005), this form of a prior for  $\sigma_{b_i}^2$  yielded the improper posteriors for  $\sigma_{b_i}^2$  and  $b_i$ , which had infinite mass near zero, thus, if the Markov chain truly converged, the values of  $\sigma_{b_i}^2$  and  $b_i$  should be

fixed at zero. Hobert and Casella (1996) discussed that the MCMC procedure with improper posterior cannot converge. They found, however, that the posterior sample from the MCMC could show nice-looking behavior despite the improper posterior. In our method (Bayes2), our main concern is to detect the segregation of QTL in each family, indicated by a binary variable  $s_i$ , rather than to estimate  $\sigma_{b_i}^2$  and  $b_i$ . Posterior samples of  $s_i$  might be robust to the impropriety of the posteriors of  $\sigma_{b_i}^2$  and  $b_i$ . Therefore, in the present analyses, we daringly used the prior  $p(\sigma_{b_i}^2 | s_i) \propto 1/\sigma_{b_i}^2$  and the posterior samples of  $\sigma_{b_i}^2$  and  $b_i$  seemingly behaved well along with  $s_i$  while this problem of the improper posterior requires further consideration.

One might be interested in the influence of a prior mean of QTL number,  $N$ , on the power of QTL detection for Bayesian methods. We, thus, applied additional analyses for the same data sets used in simulations assuming the prior means of  $N$  equal to 1 and 10. As the results of these additional analyses, we obtained almost the same powers as the original analyses with the prior mean of  $N$  being 2 (results not shown).

#### Utility of multiple families derived from germplasm collections

The germplasm collections that have been recently established for some crops are useful for association mapping of traits of economic importance. Some statistical methods, including mixed linear model and Bayesian method, have been devised for whole genome association studies in such collections (Yu *et al.*, 2006, 2008; Iwata *et al.*, 2007). A whole genome scan with association mapping requires a considerable number of markers that cover the entire genome at high density, making it both expensive and time-consuming. Therefore, multiple populations of segregating families derived by crossing some accessions in the collection with a reference line such as a popular commercial variety, as described here, would be valuable for obtaining preliminary QTL information, including the number of QTL and their positions for a subsequent association study, in which the target regions can be confined to the QTL regions estimated from the preliminary linkage QTL analysis. The population structure present in the original collections, which decreases the efficiencies in association mapping, is also minimized by reshuffling the genomes of two parents in each family to construct the multiple families (Yu *et al.*, 2008). Therefore, adopting a linkage mapping strategy in the multiple families derived from germplasm collections will improve the power of QTL detection although the mapping resolution is inferior to association mapping approach. In addition, in the analysis of the multiple families described in the present study, we can select the tested lines that will be useful for the future breeding programs. The results of our study show that the Bayesian method developed for analyzing such families can play a practical role in QTL analysis in germplasm collections.

The program (written with Fortran 77) used in the simulation experiment of this study can be applied to actual data of multiple RILs derived from crossing a reference line to several tested lines and a Windows executable version of the program is available on request to the authors.

## References

- Carlin BP, Chib S (1995). Bayesian model choice via Markov chain Monte Carlo. *J R Stat Soc Ser B* **57**: 473–484.
- Crossa J, Burgueño J, Dreisigacker S, Vargas M, Herrera-Foessel SA, Lillemo M *et al.* (2007). Association analysis of historical bread wheat germplasm using additive genetic covariance of relatives and population structure. *Genetics* **177**: 1889–1913.
- Flint-Garcia SA, Thuillet A, Yu J, Pressoir G, Romero SM, Mitchell SE *et al.* (2005). Maize association population: a high resolution platform for QTL dissection. *Plant J* **44**: 1054–1064.
- Green PJ (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**: 711–732.
- Godsill SJ (2001). On the relationship between MCMC model uncertainty methods. *J Comput Graph Stat* **10**: 230–248.
- Hastings WK (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**: 98–109.
- Hayashi T, Awata T (2008). A Bayesian method for simultaneously detecting mendelian and imprinted quantitative trait loci in experimental crosses of outbred species. *Genetics* **178**: 527–538.
- Hobert JP, Casella G (1996). The effect of improper priors on Gibbs sampling in hierarchical linear mixed models. *J Am Stat Assoc* **91**: 1461–1473.
- Iwata H, Uga Y, Yoshioka Y, Ebana K, Hayashi T (2007). Bayesian association mapping of multiple quantitative trait loci and its application to the analysis of genetic variation among *Oryza sativa* L. germplasms. *Theor Appl Genet* **114**: 1437–1449.
- Jannink JL, Fernando RL (2004). On the Metropolis-Hastings acceptance probability to add or drop a quantitative trait locus in Markov chain Monte Carlo-based Bayesian analysis. *Genetics* **166**: 641–643.
- Jannink JL, Wu XL (2003). Estimating allelic number and identity in state of QTLs in interconnected families. *Genet Res* **81**: 133–144.
- Jiang C, Zeng ZB (1997). Mapping quantitative trait loci with dominant and missing markers in various crosses from two inbred lines. *Genetica* **101**: 47–58.
- Kojima Y, Ebana K, Fukuoka S, Nagamine T, Kawase M (2005). Development of an RFLP-based rice diversity research set of germplasm. *Bred Sci* **55**: 431–440.
- Metropolis N, Rosenbluth AW, Rosenbluth MN, Teller AH, Teller E (1953). Equation of state calculations by fast computing machines. *J Chem Phys* **21**: 1087–1092.
- Sillanpää MJ, Arjas E (1998). Bayesian mapping of multiple quantitative trait loci from incomplete inbred line cross data. *Genetics* **148**: 1373–1388.
- Sillanpää MJ, Gasbarra D, Arjas E (2004). Comment on ‘on the Metropolis-Hastings acceptance probability to add or drop a quantitative trait locus in Markov Chain Monte Carlo-based Bayesian analyses’. *Genetics* **167**: 1037.
- ter Braak CJF, Boer MP, Bink MAM (2005). Extending Xu’s Bayesian model for estimating polygenic effects using markers of the entire genome. *Genetics* **170**: 1435–1438.
- Xu S (1998). Mapping quantitative trait loci using multiple families of line crosses. *Genetics* **148**: 517–524.
- Xu S (2003). Estimating polygenic effects using markers of entire genome. *Genetics* **163**: 789–801.
- Yi N, George V, Allison DB (2003). Stochastic search variable selection for identifying multiple quantitative trait loci. *Genetics* **164**: 1129–1138.
- Yu J, Pressoir G, Briggs WH, Bi IV, Yamasaki M, Doebley JF *et al.* (2006). A unified mixed-model method for association mapping accounting for multiple levels of relatedness. *Nat Genet* **38**: 203–208.
- Yu J, Holland JB, McMullen MD, Buckler ES (2008). Genetic design and statistical power of nested association mapping in maize. *Genetics* **178**: 539–551.

## Appendix A

## MCMC sampling

The MCMC cycle to estimate each element of  $\theta$  consists of the following steps:

- Updating the effects of allele derived from a reference line at each QTL,  $a_l$  ( $l = 1, 2, \dots, N$ ).
- Updating the effects of alleles derived from the tested lines at each QTL,  $b_{li}$  ( $l = 1, 2, \dots, N; i = 1, 2, \dots, m$ );
- Updating the variance of the effects of alleles from the tested lines at each QTL,  $\sigma_{bl}^2$  ( $l = 1, 2, \dots, N$ );
- Updating the intercept  $\mu$  and the residual variance  $\sigma_e^2$ ;
- Updating the variables indicating QTL genotypes for each of the individuals in the whole families,  $u_{lij}$  ( $l = 1, 2, \dots, N; j = 1, 2, \dots, n_i$  for  $i = 1, 2, \dots, m$ );
- Updating the QTL locations,  $\lambda_l$  ( $l = 1, 2, \dots, N$ );
- Updating the variables indicating segregation at each QTL in each family,  $s_{li}$  ( $l = 1, 2, \dots, N; i = 1, 2, \dots, m$ );
- Adding one new QTL to the model or removing one existing QTL from the model.

Steps a, b, c, d, e and g are performed by means of Gibbs sampling whereas steps f and h are performed using Metropolis–Hastings algorithm. The full conditional posterior distributions of some unobservables, from which updating values are sampled in Gibbs sampling algorithm, can be constructed from the likelihood function of the phenotype  $y_{ij}$ , which is presented as a normal distribution with mean  $\mu + \sum_{i=1}^m \Delta_{lij}$  and variance  $\sigma_e^2$ , denoted by  $\phi(y_{ij} | \mu + \sum_{i=1}^m \Delta_{lij}, \sigma_e^2)$ , and the prior distribution of the unobservables. Here we explain how to perform each of the MCMC steps.

**Update of effects of alleles from a reference line:** Assuming the prior distribution of  $a_l$  as  $p(a_l) \propto 1$ , the full conditional posterior distribution of  $a_l$  is a normal distribution with mean

$$\frac{\sum_{i=1}^m \sum_{j=1}^{n_i} \{s_{li} u_{lij} + (1 - s_{li})\} \left( y_{ij} - \mu - \sum_{k \neq l}^N \Delta_{kij} \right)}{\sum_{i=1}^m \sum_{j=1}^{n_i} \{s_{li} u_{lij} + (1 - s_{li})\}^2}$$

and variance

$$\frac{\sigma_e^2}{\sum_{i=1}^m \sum_{j=1}^{n_i} \{s_{li} u_{lij} + (1 - s_{li})\}^2}$$

from which  $a_l$  is sampled.

**Update of effects of alleles from the tested lines:** For both  $s_{li} = 1$  and  $s_{li} = 0$ , we assumed the prior distribution of  $b_{li}$  as  $p(b_{li} | \sigma_{bl}^2, s_{li}) = \phi(b_{li} | 0, \sigma_{bl}^2)$ . When  $s_{li} = 1$ , the full conditional posterior distribution of  $b_{li}$  is a normal distribution with mean

$$\frac{\sum_{j=1}^{n_i} (1 - u_{lij}) \left( y_{ij} - \mu - \sum_{k \neq l}^N \Delta_{kij} \right)}{\sum_{j=1}^{n_i} (1 - u_{lij})^2 + \sigma_e^2 / \sigma_{bl}^2}$$



and variance

$$\frac{1}{\left\{ (1/\sigma_e^2) \sum_{j=1}^{n_i} (1 - u_{ij})^2 + (1/\sigma_{bl}^2) \right\}}$$

When  $s_{li} = 0$ , the full conditional posterior distribution of  $b_{li}$  is independent of the data  $y$  and given as the prior distribution  $p(b_{li} | \sigma_{bl}^2, s_{li})$ .

**Update of variances of effects of alleles from the tested lines:** Assuming that  $p(\sigma_{bl}^2) \propto 1/\sigma_{bl}^2$ , the full conditional posterior distribution of  $\sigma_{bl}^2$  is a scaled inverted  $\chi^2$  distribution. Updated value of  $\sigma_{bl}^2$  is given as  $\sum_{i=1}^m b_{li}^2 / \chi_m^2$ , where  $\chi_m^2$  is a random number sampled from a  $\chi^2$  distribution with  $m$  d.f.

**Update of intercept and residual variance:** Assuming that  $p(\mu) \propto 1$  and  $p(\sigma_e^2) \propto 1/\sigma_e^2$ , the full conditional posterior distribution of  $\mu$  is a normal distribution with mean  $\sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \sum_{l=1}^N \Delta_{lij}) / n_T$  and variance  $\sigma_e^2 / n_T$  and that of  $\sigma_e^2$  is a scaled inverted  $\chi^2$  distribution. The updated value of  $\sigma_e^2$  is obtained as  $\sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \mu - \sum_{l=1}^N \Delta_{lij})^2 / \chi_{n_T}^2$ , where  $\chi_{n_T}^2$  is a  $\chi^2$  variable with  $n_T$  degrees of freedom.

**Update of variable indicating QTL genotype:** Given the prior probabilities for  $u_{ij} = 1$  and  $0$ , the full conditional posterior probability of  $u_{ij} = k$  ( $k = 0$  or  $1$ ) is written as

$$\frac{p(u_{ij} = k) \phi\left(y_{ij} | \mu + s_{li} \{ka_i + (1 - k)b_{li}\} + (1 - s_{li})a_i + \sum_{h \neq l}^N \Delta_{hil}, \sigma_e^2\right)}{p(u_{ij} = 1) \phi\left(y_{ij} | \mu + a_i + \sum_{h \neq l}^N \Delta_{hil}, \sigma_e^2\right) + p(u_{ij} = 0) \phi\left(y_{ij} | \mu + s_{li}b_{li} + (1 - s_{li})a_i + \sum_{h \neq l}^N \Delta_{hil}, \sigma_e^2\right)}$$

When  $s_{li} = 0$ , this probability is equal to the prior probability  $p(u_{ij} = k)$ .

$$\min \left\{ 1, \frac{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + \sum_{l=1}^N \Delta_{lij} + s_{N^*i} \{u_{N^*ij} a_{N^*} + (1 - u_{N^*ij}) b_{N^*i}\} + (1 - s_{N^*i}) a_{N^*}, \sigma_e^2\right)}{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + \sum_{l=1}^N \Delta_{lij}, \sigma_e^2\right)} \frac{\delta p_d}{p_a(N + 1)} \right\},$$

**Update of QTL location:** For updating the present location of the  $l$ th QTL  $\lambda_l$ , a new location  $\lambda_l^*$  is proposed by sampling a value from a uniform distribution over a small interval including  $\lambda_l$  and the new genotypes  $u_{ij}^*$  ( $i = 1, 2, \dots, m; j = 1, 2, \dots, n_i$ ) of all individuals in RILs are proposed from  $p(u_{ij}^* | \lambda_l)$  corresponding to the new QTL location to the present QTL genotypes  $u_{ij}$ . The proposed QTL location is accepted with probability  $\gamma$ , which is written as

$$\gamma = \min \left\{ 1, \frac{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + s_{li} \{u_{ij}^* a_i + (1 - u_{ij}^*) b_{li}\} + (1 - s_{li}) a_i + \sum_{k \neq l}^N \Delta_{kij}, \sigma_e^2\right)}{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + s_{li} \{u_{ij} a_i + (1 - u_{ij}) b_{li}\} + (1 - s_{li}) a_i + \sum_{k \neq l}^N \Delta_{kij}, \sigma_e^2\right)} \right\}$$

Then  $\lambda_l^*$  is the updated QTL location and  $u_{ij}^*$  is the new QTL genotypes. If the proposed QTL is rejected with probability  $1 - \gamma$ , the location and genotypes at the QTL remain  $\lambda_l$  and  $u_{ij}$ .

**Update of variables indicating QTL segregation:** Assuming the prior probabilities of  $s_{li} = 1$  and  $s_{li} = 0$ , which were set at  $0.5$  for both values of  $s_{li}$  in this study, the full conditional posterior probability of  $s_{li} = k$  ( $k = 1$  or  $0$ ) is expressed as

$$\frac{\prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu_i + k \{u_{ij} a_i + (1 - u_{ij}) b_{li}\} + (1 - k) a_i + \sum_{h \neq l}^N \Delta_{hij}, \sigma_e^2\right)}{\prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu_i + a_i + \sum_{h \neq l}^N \Delta_{hij}, \sigma_e^2\right) + \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu_i + u_{ij} a_i + (1 - u_{ij}) b_{li} + \sum_{h \neq l}^N \Delta_{hij}, \sigma_e^2\right)}$$

**Update of QTL number:** The QTL number is updated with RJ-MCMC algorithm. The number of QTL,  $N$ , is updated by adding one new QTL to the model with probability  $p_a$  or deleting one existing QTL from the model with probability  $p_d$  in the way described by Jannink and Fernando (2004) and Sillanpää *et al.* (2004). For a proposed QTL number,  $N^*$ , there are three possible values;  $N^* = N + 1$ ,  $N^* = N - 1$  and  $N^* = N$  with probabilities  $p_a$ ,  $p_d$  and  $1 - p_a - p_d$ .

When attempting to add one new QTL, firstly the location of the QTL  $\lambda_{N^*}$  is sampled from a uniform distribution over a whole genome region,  $p(\lambda_{N^*})$ . Then the segregation of the additional QTL in each family,  $s_{N^*i}$  ( $i = 1, 2, \dots, m$ ), and the QTL genotypes of each individual,  $u_{N^*ij}$  ( $i = 1, 2, \dots, m; j = 1, 2, \dots, n_i$ ), are determined by sampling from  $p(s_{N^*i})$  and  $p(u_{N^*ij} | \lambda_{N^*})$ , respectively. Moreover, the QTL effects of allele from the reference line and tested lines are sampled from  $p(a_{N^*})$  and  $p(b_{N^*i})$  ( $i = 1, 2, \dots, m$ ), respectively. The new QTL is accepted with probability

where  $\delta$  is a mean of Poisson distribution used as the prior for QTL number.

For deleting one existing QTL, a random choice is made among the existing QTL. The chosen QTL is then proposed to be deleted from the model. If the  $l$ th QTL is proposed to be deleted, the probability for accepting the proposal is

$$\min \left\{ 1, \frac{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + \sum_{k \neq l}^N \Delta_{kij}, \sigma_e^2\right)}{\prod_{i=1}^m \prod_{j=1}^{n_i} \phi\left(y_{ij} | \mu + \sum_{k=1}^N \Delta_{kij}, \sigma_e^2\right)} \frac{p_a N}{\delta p_d} \right\}$$