1

# Efficiency of genome-wide association study in open-pollinated populations

2   José Marcelo Soriano Viana,[*1] Gabriel Borges Mundim,[*] Fabyano Fonseca e Silva,[†] and Antonio

3   Augusto Franco Garcia[‡]

4   [*]Federal University of Viçosa, Department of General Biology, 36570-900, Viçosa, MG, Brazil.

5   [†]Federal University of Viçosa, Department of Animal Science, 36570-900, Viçosa, MG, Brazil.

6   [‡]ESALQ - University of São Paulo, Department of Genetics, 13418-900, Piracicaba, SP, Brazil.

7   Reference     number     for     data     available     in     public     repository:

8   https://dx.doi.org/10.6084/m9.figshare.3201838.v1.

9   *REALbreeding* private link: https://figshare.com/s/618bee7accd410464232.

10

1

11    Running title: GWAS in open-pollinated populations.

12    **KEYWORDS** association mapping; GWAS; linkage disequilibrium; inbred lines panel; RIL.

13    [1]Corresponding author: José Marcelo Soriano Viana. Federal University of Viçosa, Department of

14    General Biology, 36570-900, Viçosa, MG, Brazil. E-mail: jmsviana@ufv.br. Telephone:

15    +55(31)3899-2514.

16    **ABSTRACT** Genome-wide association studies (GWAS) with plant species have employed inbred

17    lines panels. Thus, to our knowledge, no information is available on theory and efficiency of

18    GWAS in open-pollinated populations. Our objectives are to present quantitative genetics theory for

19    GWAS, evaluate the relative efficiency of GWAS in non-inbred and inbred populations and in an

20    inbred lines panel, and assess factors affecting GWAS, such as linkage disequilibrium (LD), sample

21    size, and quantitative trait locus (QTL) heritability. Fifty samples of 400 individuals from

22    populations with LD were simulated. Individuals were genotyped for 10,000 single nucleotide

23    polymorphisms (SNPs) and phenotyped for traits with different degrees of dominance controlled by

24    10 QTLs and 90 minor genes. The average SNP density was 0.1 centiMorgan and the trait

25    heritabilities were 0.4 and 0.8. We assessed GWAS efficiency based on the power of QTL

26    detection, number of false-positive associations, bias in the estimated QTL position, and range of

27    the significant SNPs for the same QTL. When the LD between a QTL and one or more SNPs is

28    restricted to markers very close to or within the QTL, GWAS in open-pollinated populations can be

29    highly efficient, depending mainly on QTL heritability and sample size. GWAS achieved the

30    highest power of QTL detection, the smallest number of false-positive associations, and the lowest

31    bias in the estimated QTL position for the inbred lines panel correcting for population structure.

32    Under low QTL heritability and reduced sample size, GWAS is ineffective for non-inbred and

33    inbred populations and for inbred lines panel.

**INTRODUCTION**

34

35      Association mapping is a high-resolution method for mapping quantitative trait locus (QTL)

36      based on linkage disequilibrium (LD) (Yu and Buckler 2006). Linkage disequilibrium is commonly

37      defined as the non-random association of alleles at two loci carried on the same gamete, caused by

38      their shared history of mutation and recombination (Weir 2008; Flint-Garcia *et al.* 2003).

39      Association mapping has been successful in detecting genes controlling human diseases and

40      quantitative traits in plant and animal species (Pearson and Manolio 2008; Zhu *et al.* 2008;

41      Barendse *et al.* 2007). There are two main association mapping strategies: the candidate gene

42      approach, which focuses on polymorphisms in specific genes that control the traits of interest, and

43      the genome-wide association study (GWAS), which surveys the entire genome for polymorphisms

44      associated with complex traits (Rafalski 2010).

45      With the advent of high-throughput genotyping and sequencing technologies, breeders have

46      used GWAS to identify genes underlying quantitative trait variation. Compared to QTL mapping,

47      which has statistical precision in the range of 10 to 30 centiMorgans (cM) (confidence interval or

48      highest probability density interval), the main advantage of GWAS is a more precise identification

49      of candidate genes (Zhu *et al.* 2008). Another advantage is the use of a breeding population instead

50      of one derived by crossing two inbred or pure lines (Flint-Garcia *et al.* 2005). However, as

51      highlighted by Weir (2010), the efficiency of GWAS is considerably affected by relatedness and

52      population structure, which can generate spurious association between unlinked marker and QTL.

53      Rafalski (2010) emphasized that the choices of population (due to the degree of LD and genotypic

54      variation), marker density, and sample size are crucial decisions for achieving greater power of

55      QTL detection. Ingvarsson and Street (2011) discussed the influence of population size, extent of

56      LD, trait heritability (precision of phenotyping), and population structure on GWAS efficiency,

57      highlighting that studies with plant species should greatly increase population size to detect QTLs

58      with lower effect (heritability of 1−2%).

59    Yu *et al.* (2006) proposed a mixed model approach for GWAS analysis called the Q + K (or

60    QK) method, where Q and K are the population structure and kinship matrices, respectively. This

61    method has provided the best results and greatly has improved the control of both type I and type II

62    error rates compared with other methods. Stich and Melchinger (2009) and Yang *et al.* (2010)

63    compared GWAS methods based on simulated and field data. Based on type I error control and

64    power of QTL detection, they concluded that the mixed model approach using only the kinship

65    matrix (K model) to correct for relatedness was more efficient than the approaches controlling only

66    population structure (Q model) and both population structure and relatedness because the spurious

67    associations could not be completely controlled by population structure. Based on simulated inbred

68    lines panel, Bernardo (2013) demonstrated that his models G and QG were superior to the K and

69    QK, respectively, where G indicates a model that uses genome-wide markers to account for QTLs

70    on background chromosomes. The new approach showed a better balance between power of QTL

71    detection and false discovery rate (FDR).

72    Recently, many instances of GWAS have been published with plant species, including barley,

73    sorghum, wheat, rice, sugarcane, soybean and particularly maize (Ingvarsson and Street 2011). Pace

74    *et al.* (2015) carried out a GWAS with 384 maize inbred lines evaluated for 22 seedling root

75    architecture traits and genotyped with 681,257 single nucleotide polymorphisms (SNPs). They

76    identified 268 marker-trait associations. Some of these SNPs were located within or near (less than

77    one kilo base pairs) to candidate genes involved in root development at the seedling stage.

78    Thirunavukkarasu *et al.* (2014) evaluated 240 elite inbred lines of subtropical maize under water

79    stress and used a set of 29,619 high-quality SNPs. The GWAS identified 50 SNPs consistently

80    associated with agronomic traits related to functional traits that could lead to drought tolerance.

81    Thirty-one of the SNPs detected were situated near drought-tolerance genes. Schaefer and Bernardo

82    (2013) used GWAS on a collection of 284 historical maize inbred lines and 39,166 SNPs and

83    identified 19 QTLs for flowering time, 13 for kernel composition, and 22 for disease resistance.

84    However, only two candidate genes were suggested: one regulating days to anthesis and one

85    regulating oil concentration. Additionally, several QTL hot spots (chromosome regions with

86    previously mapped QTLs) were also identified, affecting days to anthesis (four), days to silking

87    (two), and resistance to northern corn leaf blight (four) and Goss's wilt and blight (one).

88        Genome-wide association studies with plant species have employed inbred lines panels. Thus,

89    to our knowledge, no information is available on theory and efficiency of GWAS in open-pollinated

90    populations. Our objectives are to present quantitative genetics theory for GWAS, to evaluate the

91    relative efficiency of GWAS in non-inbred and inbred populations and in an inbred lines panel, and

92    to assess factors that affect GWAS, such as LD, sample size, and QTL heritability.

93                              **MATERIALS AND METHODS**

94    **Quantitative genetics theory for GWAS in open-pollinated populations**

95        Consider a biallelic QTL (alleles **B**/**b**) and a SNP (alleles **C**/**c**) located in the same

96    chromosome, and a population (generation 0) of an open-pollinated species. Assuming LD, the joint

97    gamete and joint genotype probabilities in the population are presented by Viana et al. (2016). The

98    QTL genotypic values are $G_{\mathbf{BB}} = m_b + a_b$, $G_{\mathbf{Bb}} = m_b + d_b$, and $G_{\mathbf{bb}} = m_b - a_b$, where $m_b$ is

99    the mean of the genotypic values of the homozygotes, $a_b$ is the deviation between the genotypic

100   value of the homozygote of higher expression and $m_b$, and $d_b$ is the dominance deviation (the

101   deviation between the genotypic value of the heterozygote and $m_b$ ). The average genotypic values

102   of individuals with the genotypes **CC**, **Cc**, and **cc** are

103
$$G_{\mathbf{CC}} = \frac{1}{p_c^2}\left( f_{22}G_{\mathbf{BBCC}} + f_{12}G_{\mathbf{BbCC}} + f_{02}G_{\mathbf{bbCC}} \right)$$
$$= M + 2q_c\kappa_{bc}\alpha_b + \left( -2q_c^2\kappa_{bc}^2 d_b \right) = M + 2\alpha_{\mathbf{C}} + D_{\mathbf{CC}} = M + A_{\mathbf{CC}} + D_{\mathbf{CC}} = m_c + a_c$$

104
$$G_{\mathbf{Cc}} = \frac{1}{2p_c q_c}\left( f_{21}G_{\mathbf{BBCc}} + f_{11}G_{\mathbf{BbCc}} + f_{01}G_{\mathbf{bbCc}} \right)$$
$$= M + (q_c - p_c)\kappa_{bc}\alpha_b + 2p_c q_c\kappa_{bc}^2 d_b = M + (\alpha_{\mathbf{C}} + \alpha_{\mathbf{c}}) + D_{\mathbf{Cc}} = M + A_{\mathbf{Cc}} + D_{\mathbf{Cc}} = m_c + d_c$$

$$G_{\mathbf{cc}} = \frac{1}{q_c^2}\left(f_{20}G_{\mathbf{BBcc}} + f_{10}G_{\mathbf{Bbcc}} + f_{00}G_{\mathbf{bbcc}}\right)$$

$$= M + \left(-2p_c\kappa_{bc}\alpha_b\right) + \left(-2p_c^2\kappa_{bc}^2 d_b\right) = M + 2\alpha_{\mathbf{c}} + D_{\mathbf{cc}} = M + A_{\mathbf{cc}} + D_{\mathbf{cc}} = m_c - a_c$$

106     where     $p$     is     the     frequency     of     the     major     allele

107     (**B** or **C**), $q = 1 - p$ is the frequency of the minor allele (**b** or **c**), $f_{ij}$ is the probability of the

108     individual with i and j copies of the allele **B** of the QTL and allele **C** of the SNP (i, j = 2, 1, or 0)

109     (for simplicity, we omitted the superscript (0) - for generation 0 - in all parameters that depend on

110     the LD measure of generation $-1$), $M = m_b + \left(p_b - q_b\right)a_b + 2p_bq_bd_b$ is the population mean,

111     $\kappa_{bc} = \left[\dfrac{\Delta_{bc}^{(-1)}}{p_c q_c}\right]$,     $\alpha_b = a_b + \left(q_b - p_b\right)d_b$   is   the   average   effect   of   a   gene   substitution,

112     $\alpha_{\mathbf{C}} = q_c\kappa_{bc}\alpha_b$ and $\alpha_{\mathbf{c}} = -p_c\kappa_{bc}\alpha_b$ are the average effects of the SNP alleles, and A and D are

113     the SNP additive and dominance values. $\Delta_{bc}^{(-1)} = P_{\mathbf{BC}}^{(-1)}P_{\mathbf{bc}}^{(-1)} - P_{\mathbf{Bc}}^{(-1)}P_{\mathbf{bC}}^{(-1)}$ is the measure of LD in

114     the gametic pool of generation $-1$ (Kempthorne 1957), where $P^{(-1)}$ indicates a joint gamete

115     probability. Another common measure of LD is the square of the correlation between the values of

116     the alleles at the two loci ($r_{bc}^{(-1)}$) in the gametic pool of generation $-1$ (Hill and Robertson 1968).

117     Note that $\Delta_{bc}^{(-1)} = r_{bc}^{(-1)}\sqrt{p_bq_bp_cq_c}$ . The average effect of substituting the allele **C** for **c** is

118     $\alpha_{SNP} = \alpha_{\mathbf{C}} - \alpha_{\mathbf{c}} = \kappa_{bc}\alpha_b$ . The dominance deviation for the SNP is $d_{SNP} = \kappa_{bc}^2 d_b$ . The other

119     SNP parameters are $m_c = M + \left(q_c - p_c\right)\alpha_{SNP} - \left(1 - 2p_cq_c\right)d_{SNP}$, $a_c = \alpha_{SNP} - \left(q_c - p_c\right)d_{SNP}$,

120     and $d_c = d_{SNP}$.

121     Assuming no QTL in LD with the SNP, $G_{\mathbf{CC}} = G_{\mathbf{Cc}} = G_{\mathbf{cc}} = M$. Thus, the identification of

122     the QTL can be based on testing the hypothesis that there is no difference between these genotypic

6

123 means (based on analysis of variance). Assuming thousands of SNPs, it is necessary to employ a

124 Bonferroni-type procedure to control the type I error when there are multiple-comparisons, as that

125 proposed by Benjamini and Hochberg (1995). Note that $\alpha_{SNP} = a_c + (q_c - p_c)d_{SNP}$, where

126 $a_c = G_{CC} - m_c$, $d_{SNP} = G_{Cc} - m_c$, and $m_c = (G_{CC} + G_{cc})/2$.

127 Alternatively, the QTL identification can be done by testing that there is no relationship

128 between the genotypic values for the individuals **CC**, **Cc**, and **cc** with the number of copies of one

129 SNP allele. The parameters of the additive-dominance model can be derived by fitting the model

130 $G = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$ (x = 2, 1, or 0), where G is the QTL genotypic value. The model can be

131 expressed as $y_{(9 \times 1)} = X_{(9 \times 3)}.\beta_{(3 \times 1)}$ + error vector$_{(9 \times 1)}$, where y is the vector of QTL genotypic

132 values, conditional on the SNP genotype, X is the incidence matrix, and β is the parameter vector.

133 The matrix of genotype probabilities is $P_{(9 \times 9)}$ = diagonal$\{ f_{ij} \}$. Thus, for the complete model or a

134 reduced model, $\beta = (X'PX)^{-1}(X'Py)$. The parameters for the complete model are

135 $\beta_0 = M - 2p_c\alpha_{SNP} - 2p_c^2 d_{SNP}$

136 $\beta_1 = \alpha_{SNP} + (1 + 2p_c)d_{SNP}$

137 $\beta_2 = -d_{SNP}$

138 The alternative regression model is $G = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$, where $x_1$ = 1, 0, or −1 if the

139 individual is **CC**, **Cc**, or **cc**, and $x_2$ = 0 or 1 if the individual is homozygous or heterozygous,

140 respectively. Fitting the complete model, $\beta_0 = m_c$, $\beta_1 = a_c$, and $\beta_2 = d_{SNP}$. Assuming no QTL in

141 LD with the SNP, $\beta_1 = \beta_2 = 0$ and $\beta_0 = M$, regardless of the model. Fitting the additive model,

142 $G = \beta_0 + \beta_1 x + \varepsilon$ or $G = \beta_0 + \beta_1 x_1 + \varepsilon$ (no dominance), $\beta_1 = \alpha_{SNP}$.

143 If there are two **QTLs** (alleles **B/b** and **E/e**) in LD with the SNP (alleles **C/c**), it can be

144 demonstrated that (Viana *et al.* 2016)

145    $\alpha_{SNP} = \kappa_{bc}\alpha_b + \kappa_{ce}\alpha_e$

146    $d_{SNP} = \kappa_{bc}^2 d_b + \kappa_{ce}^2 d_e$

147    where $\kappa_{ce} = \left[\dfrac{\Delta_{ce}^{(-1)}}{p_c q_c}\right]$. Thus, the average effect of a SNP substitution (and the SNP additive value)

148    is proportional to the measure of LD and to the average effect of a gene substitution for each QTL

149    that is in LD with the marker, and the SNP dominance deviation (and the SNP dominance value) is

150    proportional to the square of the LD value and to the dominance deviation for each QTL that is in

151    LD with the marker.

152         If there is population structure, this must be corrected in the GWAS to avoid spurious

153    associations due to admixture LD. For simplicity, consider two subpopulations in Hardy-Weinberg

154    equilibrium and a SNP (alleles **C**/**c**) and a QTL (alleles **B**/**b**) unlinked, in linkage equilibrium in

155    both subpopulations. Assuming that p and q are the allelic frequencies in one subpopulation and r

156    and s are the allelic frequencies in the other subpopulation, the average genotypic value of

157    individuals **CC**, **Cc**, and **cc** are

158    $$G_{\mathbf{CC}} = m_b + \left(\frac{1}{u_1 p_c^2 + u_2 r_c^2}\right)\left\{\left[u_1(p_b - q_b)p_c^2 + u_2(r_b - s_b)r_c^2\right]a_b \right.$$
$$\left. + \left(u_1 2p_b q_b p_c^2 + u_2 2r_b s_b r_c^2\right)d_b\right\}$$

159    $$G_{\mathbf{Cc}} = m_b + \left(\frac{1}{u_1 2p_c q_c + u_2 2r_c s_c}\right)\left\{\left[u_1(p_b - q_b)2p_c q_c + u_2(r_b - s_b)2r_c s_c\right]a_b \right.$$
$$\left. + \left(u_1 2p_b q_b 2p_c q_c + u_2 2r_b s_b 2r_c s_c\right)d_b\right\}$$

160    $$G_{\mathbf{cc}} = m_b + \left(\frac{1}{u_1 q_c^2 + u_2 s_c^2}\right)\left\{\left[u_1(p_b - q_b)q_c^2 + u_2(r_b - s_b)s_c^2\right]a_b \right.$$
$$\left. + \left(u_1 2p_b q_b q_c^2 + u_2 2r_b s_b s_c^2\right)d_b\right\}$$

8

161   where $u_1$ and $u_2$ are the proportions of individuals from subpopulations 1 and 2 (probabilities of

162   an individual belongs to subpopulations 1 and 2). Only if there is no population structure ($u_1 = 1$ or

163   0), $G_{CC} = G_{Cc} = G_{cc} = M$ (and $\beta_1 = \beta_2 = 0$ and $\beta_0 = M$).

**Quantitative genetics theory for GWAS with inbred lines panel**

165       In general, the inbred lines in a panel represent the genetic variability for the traits being

166   assessed. Therefore, an inbred lines panel includes inbreds from distinct populations or heterotic

167   groups. Consider again a QTL (alleles **B/b**) and a SNP (alleles **C/c**) located in the same

168   chromosome, and that they are in LD in a population (generation 0). Assuming n ($n \rightarrow \infty$)

169   generations of selfing, the (limits of the) probabilities of the inbreds (recombinant inbred lines;

170   RILs) are (for simplicity, we omitted again the superscript (0) - for generation 0 - in all parameters

171   that depend on the LD measure of generation $-1$)

172   $$f_{22}^{(n)} = f_{22} + \frac{1}{2}\left(f_{21} + f_{12}\right) + \frac{1}{4}f_{11} + \frac{1}{2}\left(\frac{1-2\theta_{bc}}{1+2\theta_{bc}}\right)\Delta_{bc}^{(-1)}$$

173   $$f_{20}^{(n)} = f_{20} + \frac{1}{2}\left(f_{21} + f_{10}\right) + \frac{1}{4}f_{11} - \frac{1}{2}\left(\frac{1-2\theta_{bc}}{1+2\theta_{bc}}\right)\Delta_{bc}^{(-1)}$$

174   $$f_{02}^{(n)} = f_{02} + \frac{1}{2}\left(f_{01} + f_{12}\right) + \frac{1}{4}f_{11} - \frac{1}{2}\left(\frac{1-2\theta_{bc}}{1+2\theta_{bc}}\right)\Delta_{bc}^{(-1)}$$

175   $$f_{00}^{(n)} = f_{00} + \frac{1}{2}\left(f_{01} + f_{10}\right) + \frac{1}{4}f_{11} + \frac{1}{2}\left(\frac{1-2\theta_{bc}}{1+2\theta_{bc}}\right)\Delta_{bc}^{(-1)}$$

176   where $\theta_{bc}$ is the frequency of recombinant gametes. The haplotypes are $P_{BC}^{(n)} = p_b p_c + \Delta_{bc}^{(n)}$,

177   $P_{Bc}^{(n)} = p_b q_c - \Delta_{bc}^{(n)}$,        $P_{bC}^{(n)} = q_b p_c - \Delta_{bc}^{(n)}$,        and        $P_{bc}^{(n)} = q_b q_c + \Delta_{bc}^{(n)}$,        where

178   $\Delta_{bc}^{(n)} = \left(\frac{1}{1+2\theta_{bc}}\right)\Delta_{bc}^{(-1)}$. Thus, if there is crossing over ($0 < \theta_{bc} \leq 0.5$), the LD in this inbred

179     population is lower than the LD in generation $-1$. If the SNP and QTL are completely linked ($\theta_{bc}$

180     $= 0$), the LD in the inbred population is the same LD in generation $-1$. The maximum decrease is

181     50%, achieved with $\theta_{bc} = 0.5$. Compared with the LD in generation 0, the LD in generation n is

182     $\Delta_{bc}^{(n)} = \left[ \dfrac{1}{\left(1 + 2\theta_{bc}\right)\left(1 - \theta_{bc}\right)} \right] \Delta_{bc}^{(0)}$. Thus, the maximum decrease is 12%, achieved with $\theta_{bc} = 0.25$.

183     In contrast, after n generations of random crosses $\Delta_{bc}^{(n)} = \left(1 - \theta_{bc}\right)^{n+1} \Delta_{bc}^{(-1)} = \left(1 - \theta_{bc}\right)^n \Delta_{bc}^{(0)}$.

184     Thus, if $0 < \theta_{bc} \leq 0.5$, the maximum decrease is 100% since $\lim\limits_{n \to \infty} \Delta_{bc}^{(n)} = 0$.

185     For the inbreds sampled from a population, we have

186     $G_{CC}^{(n)} = \dfrac{1}{f_{.2}^{(n)}}\left[ f_{22}^{(n)}\left(m_b + a_b\right) + f_{02}^{(n)}\left(m_b - a_b\right) \right] = M_{IL} + 2q_c \alpha_{SNP}^{(n)} = M_{IL} + A_{CC}^{(n)}$

187     $G_{cc}^{(n)} = \dfrac{1}{f_{.0}^{(n)}}\left[ f_{20}^{(n)}\left(m_b + a_b\right) + f_{00}^{(n)}\left(m_b - a_b\right) \right] = M_{IL} - 2p_c \alpha_{SNP}^{(n)} = M_{IL} + A_{cc}^{(n)}$

188     where $M_{IL} = m_b + \left(p_b - q_b\right)a_b$ is the inbred population mean, $\alpha_{SNP}^{(n)} = \left(\dfrac{1}{2 + 4\theta_{bc}}\right)\kappa_{bc}a_b$ is the

189     SNP average effect of allele substitution in the inbred population, and A is the SNP additive value

190     for an inbred line. Assuming no QTL in LD with the SNP, $G_{CC}^{(n)} = G_{cc}^{(n)} = M_{IL}$. Note that

191     $\alpha_{SNP}^{(n)} = \left( G_{CC}^{(n)} - G_{cc}^{(n)} \right)/2$.

192     The haplotypes of an inbred lines panel including inbreds from N populations are

193     $P_{BC}^{(n)\prime} = \bar{p}_b\bar{p}_c + \Delta_{bc}^{(n)\prime}$,    $P_{Bc}^{(n)\prime} = \bar{p}_b\bar{q}_c - \Delta_{bc}^{(n)\prime}$,    $P_{bC}^{(n)\prime} = \bar{q}_b\bar{p}_c - \Delta_{bc}^{(n)\prime}$,   and   $P_{bc}^{(n)\prime} = \bar{q}_b\bar{q}_c + \Delta_{bc}^{(n)\prime}$,

194     where $\Delta_{bc}^{(n)\prime} = \sum\limits_{i=1}^{N} u_i\left[ \Delta_{bc_i}^{(n)} + p_{b_i}p_{c_i} \right] - \left( \sum\limits_{i=1}^{N} u_i p_{b_i} \right)\left( \sum\limits_{i=1}^{N} u_i p_{c_i} \right) = \overline{\Delta}_{bc}^{(n)} + \overline{p_b p_c} - \bar{p}_b\bar{p}_c$ and $u_i$ is the

195     probability of an inbred line belonging to population i. Because this function is too complex to

196     interpret, the analysis of the LD value in an inbred lines panel relative to the LD in the inbreds from

197     each population is presented using the simulated data.

198         Due to population structure, associations involving unlinked SNP and QTL in linkage

199     equilibrium in the non-inbred populations can be declared. For simplicity, assume an inbred lines

200     panel with inbreds from two populations where an unlinked pair of SNP (alleles **C**/**c**) and QTL

201     (alleles **B**/**b**) is in linkage equilibrium. Let $u_1$ and $u_2$ be the proportions of inbreds from these

202     populations. Assuming that p and q are the allelic frequencies in one population, that r and s are the

203     allelic frequencies in the other population, and that $p \neq q$ or $r \neq s$,

204     $$G_{CC}^{(n)} = m_b + \left( \frac{1}{u_1 p_c + u_2 r_c} \right) [u_1 (p_b - q_b) p_c + u_2 (r_b - s_b) r_c] a_b$$

205     $$G_{cc}^{(n)} = m_b + \left( \frac{1}{u_1 q_c + u_2 s_c} \right) [u_1 (p_b - q_b) q_c + u_2 (r_b - s_b) s_c] a_b$$

206         If there is no population structure ($u_1 = 1$ or $0$), $G_{CC}^{(n)} = G_{cc}^{(n)} = M_{IL}$.

207     **Simulation**

208         We simulated 50 samples of populations with LD using the software *REALbreeding* (Viana *et*

209     *al.* 2016, 2013; Azevedo *et al.* 2015). This software has been developed by the first author using the

210     program *REALbasic 2009*. Population 1, generation 10r, is the advanced generation of a composite

211     of two populations in linkage equilibrium (population 1, generation 0) obtained after 10 generations

212     of random crosses, assuming a sample size of 400 individuals. Population 1, generations 10s and

213     10r10s, were obtained from Population 1, generation 0, assuming 10 generations of selfing and 10

214     generations of random crosses followed by 10 generations of selfing, respectively, assuming sample

215     sizes of 100 and 400, respectively. Populations 2, 3, and 4, generation 10s, are also inbred

216     populations (10 generations of selfing) derived from composites of two populations, also assuming

217     a sample size of 100. The parents of populations 2 and 3 were assumed to be non-improved and

218     improved populations, respectively. An improved population was defined as having frequencies of

11

219  favorable genes greater than 0.5, while a non-improved population was defined as having

220  frequencies less than 0.5. A composite is a Hardy-Weinberg equilibrium population with LD for

221  only linked markers and genes. In the case of a composite of two populations in linkage

222  equilibrium, $\Delta_{bc}^{(-1)} = \left(\dfrac{1 - 2\theta_{bc}}{4}\right)\left(p_b^1 - p_b^2\right)\left(p_c^1 - p_c^2\right)$, where the indices 1 and 2 refer to the

223  parental populations.

224  Based on our input, *REALbreeding* randomly distributed 10,000 SNPs, 10 QTLs (of higher

225  effect) and 90 minor genes (QTLs of lower effect) in 10 chromosomes (1,000 SNPs and 10 genes

226  by chromosome). The average SNP density was 0.1 cM. The genes were distributed in the regions

227  covered by the SNPs. Four, three, two, and one QTLs were inserted in chromosomes 1, 5, 9, and 10,

228  respectively. We also specified one SNP within each QTL and a minimum distance between linked

229  QTLs of 10 cM. To allow *REALbreeding* to compute the phenotypic value for each genotyped

230  individual, we informed the minimum and maximum genotypic values for homozygotes, proportion

231  between the parameter *a* for a QTL and the parameter *a* for a minor gene ($a_{QTL}/a_{mg}$), degree of

232  dominance (($d/a$)$_i$, i = 1, ..., 100), direction of dominance, and broad sense heritability.

233  *REALbreeding* saves two main files, one with the marker genotypes and another with the additive,

234  dominance, and phenotypic values (non-inbred populations) or the genotypic and phenotypic values

235  (inbred populations). The true additive and dominance genetic values or genotypic values are

236  computed from the population gene frequencies (random values), LD values, average effects of

237  gene substitution or *a* deviations, and dominance deviations. The phenotypic values are computed

238  from the true population mean, additive and dominance values or genotypic values, and from error

239  effects sampled from a normal distribution. The error variance is computed from the broad sense

240  heritability.

241  We simulated three popcorn traits. The minimum and maximum genotypic values of

242  homozygotes for grain yield, expansion volume, and days to maturity were 30 and 180 g per plant,

243  15 and 65 mL/g, and 100 and 170 days, respectively. We defined positive dominance for grain yield

12

244      $(0 < (d/a)_i \leq 1.2)$, bidirectional dominance for expansion volume $(-1.2 \leq (d/a)_i \leq 1.2)$, and no

245      dominance for days to maturity $((d/a)_i = 0)$. The broad sense heritabilities were 0.4 and 0.8. These

246      values can be associated with individual and progeny assessment, respectively. Assuming $a_{QTL}/a_{mg}$

247      = 10, each QTL explained approximately 4 and 8% of the phenotypic variance for heritabilities of

248      0.4 and 0.8. The GWAS was performed in population 1, generations 10r and 10r10s, and in the

249      inbred lines panel obtained from inbreds of the populations 1 through 4, generation 0 (generations

250      10s). To assess the influence of the sample size on the GWAS efficiency, we considered sample

251      sizes of 400 and 200. Thus, we used 100 or 50 inbreds from populations 1 through 4 to generate the

252      inbred lines panel. To assess the influence of the QTL heritability on the GWAS efficiency, we

253      converted four QTLs (QTLs 3, 7, 8, and 10 on chromosomes 1, 5, 9, and 10, respectively) to minor

254      genes and assumed QTL heritability of 12% (for trait heritability of 0.7). Then, the GWAS was

255      performed on population 1, generation 10r.

256      **Statistical analyses**

257      The analyses of LD and association were performed with the software *PowerMarker* (Liu and

258      Muse 2005) for open-pollinated populations, and *Tassel* (Bradbury *et al.* 2007) for the inbred lines

259      panel and RILs. Because there is no relationship between the inbred lines, the GWAS with the

260      inbred lines panel was based on the general linear model, correcting for population structure (Q

261      model). For the population structure analysis, we used *Structure* software (Falush *et al.* 2003) and

262      fitted the admixture model with correlated allelic frequencies and the no admixture model with

263      independent allelic frequencies. The number of SNPs, sample size, burn-in period, and number of

264      MCMC (Monte Carlo Markov chain) replications were 100 (10 random SNPs by chromosome),

265      400 (simulation 1), 10,000, and 40,000, respectively. The number of populations assumed (K)

266      ranged from 1 to 7, and the most probable K value was determined based on the inferred plateau

267      method (Viana *et al.* 2013). We used Benjamini-Hochberg FDR of 5 and 1% to control the type I

268      error (Benjamini and Hochberg 1995).

13

269       To classify each significant association as true or false, we used a program developed in

270       *REALbasic 2009* by the first author. The classification criterion was based on the difference

271       between the position of the SNP and the position of a true QTL (candidate gene). If the difference

272       was less than or equal to 2.5 cM (Yu *et al.* 2008), the association was classified as true. The GWAS

273       efficiency was assessed based on the power of QTL detection (probability of rejecting $H_0$ when $H_0$

274       is false; control of type II error), number of false-positive associations (control of type I error), bias

275       in the estimated QTL position (precision of mapping), and range of the significant SNPs for the

276       same QTL (Li *et al.* 2010).

277       **Data availability**

278       *REALbreeding* is available upon request. The data set is available at

279       https://dx.doi.org/10.6084/m9.figshare.3201838.v1. Supplemental file S1 contains detailed

280       description of all data files (SNP and QTL positions, SNP genotypes, and phenotypic values). Data

281       citation:

282       Viana, José Marcelo; Mundim, Gabriel Borges; Fonseca e Silva, Fabyano; Augusto Franco Garcia,

283       Antonio (2016): Efficiency of genome-wide association study in open-pollinated populations.

284       figshare. https://dx.doi.org/10.6084/m9.figshare.3201838.v1

285       **RESULTS**

286       The results for assessing the efficiency of GWAS in open-pollinated populations refer to

287       population 1, generation 10r. In generation 0, the degree of LD is so high that several significant

288       associations are observed along the length of a chromosome with one or more QTLs or in one or

289       more large chromosome regions (Figure 1). These several significant associations are not false-

290       positive (at least most of them). This is due to the degree of LD and presence of QTL. Even

291       assuming a FDR of 1%, it is worthless for the identification of candidate genes to infer that there

292       are one or more QTLs in a chromosome region spanning 20 cM. When the LD between a QTL and

293       one or more markers is restricted to SNPs very close to or within the QTL, the analysis can be

294    highly efficient, depending mainly on the QTL effect and sample size. Assuming a QTL heritability

295    of 8% and sample size 400 (simulation 1), the significant associations for expansion volume

296    observed in chromosome 1 evidenced five QTLs with a FDR of 5% or four QTLs with a FDR of

297    1% (Figure 1). This implies in a QTL detection power of 100%. Three of the four true QTLs

298    (candidate genes) were identified by SNPs located within the QTL while one was identified by five

299    or four SNPs in a region spanning approximately 2.0 or 1.7 cM, respectively, depending on the

300    FDR. The significant associations at a FDR of 5 or 1% for SNPs 223 (at position 21.7 cM), 243 (at

301    position 23.3), 245 (at position 23.4 cM), and 252 (at position 23.7 cM) are attributable to their LD

302    with QTL 2. The absolute LD values are 0.1488, 0.1494, 0.1747, and 0.1416, respectively (with

303    highly significant P values according to the chi-square test). The significant association at a FDR of

304    5% for SNP 627 (at position 61.8 cM) is not a false-positive, since it is in LD with QTL 2 ($|\Delta| =$

305    0.0366, chi-square test P value = 3.22E-6) and QTL 3 ($|\Delta| = 0.0302$, chi-square test P value =

306    7.55E-5). Then, the result is interpreted as a fifth QTL.

307        Only for intermediate to high QTL heritability (8 and 12%) and greater sample size were the

308    results from GWAS clearly different between days to maturity and the other two traits, except for

309    the power of QTL detection (Table 1). The number of significant associations, number of false-

310    positives, bias in QTL position, and average range of chromosome regions with one or more QTLs

311    were greater in the absence of dominance. With a FDR of 5%, the power of QTL detection ranged

312    from 88 to 100% but was associated with a high number of significant associations in chromosomes

313    with one to four QTLs. On average, each true QTL was identified based on two to three (for days to

314    maturity) SNPs, in chromosome regions spanning 0.8 to 2.2 cM. The bias in QTL position ranged

315    from 0.5 to 0.8 cM. Increasing the control of the type I error provided better results and greatly

316    reduced the number of false-positive associations. The power of QTL detection ranged from 75 to

317    100% and each QTL was identified based on one to two SNPs in chromosome regions spanning 0.4

318    to 1.1 cM. The bias in QTL position ranged from 0.3 to 0.6 cM.

15

319    Assuming a QTL heritability of 8% and sample size of 200 or a QTL heritability of 4% and

320    sample size of 400, it is better to assume a FDR of 5% to ensure greater power of QTL detection

321    and fewer false-positive associations. However, the power of detection ranged from 33 to 39%,

322    particularly due to the lower QTL effect (Table 1). With lower QTL heritability and reduced sample

323    size, GWAS is ineffective, showing an average power of QTL detection less than or equal to 5%.

324    This scenario does not improve when increasing the FDR to 10% (data not shown). Increasing the

325    QTL heritability to 12% resulted in an increase in the power of QTL detection, particularly when

326    assuming a sample size of 200 individuals (Table 1). There were also increases in bias in the QTL

327    position, range of chromosome regions with an identified QTL, number of false-positives, and

328    number of significant associations in chromosomes with one to two QTLs, mainly with greater

329    sample size. When assuming 200 individuals, the power of QTL detection reached 70–75%,

330    regardless of the trait.

331    We also provided results for comparing GWAS in open-pollinated population and in an

332    inbred lines panel. An impressive result from GWAS with an inbred lines panel is the efficacy of

333    discarding spurious associations due to population structure (Figure 2). The number of spurious

334    associations in chromosome 3 (no QTL) were reduced from 477 to zero in the analysis of expansion

335    volume assuming a FDR of 1%, QTL heritability of 8%, and sample size 400 (simulation 1).

336    Correcting for population structure decreased the number of significant associations in chromosome

337    1 (four QTLs) from 464 to 9. This implies a QTL detection power of 100% but with three to five

338    false-positive associations. The population structure analysis evidenced four subpopulations (Figure

339    3). In general, the efficiency of GWAS was greater with the inbred lines panel (Table 2). The power

340    of QTL detection was higher, and the number of false-positive associations was lower. Furthermore,

341    only SNPs within QTL showed significant associations in general. Also, no differences were

342    observed between the traits and similarly for open-pollinated populations the analysis is ineffective

343    when assuming lower QTL heritability and sample size.

16

344    The following were indicated by analysis of the parametric LD in the populations and in the

345    inbred lines panel based on a random 10 cM segment of chromosome 1 (100 SNPs): higher LD in

346    population 1, generation 0 (average absolute $\Delta = 0.0403$; 627 values greater than 0.1), lower LD in

347    population 3, generation 0 (average absolute $\Delta = 0.0203$; 48 values greater than 0.1), a slight

348    decrease in the LD with selfing (5−6%), and the lowest LD in the inbred lines panel (average

349    absolute $\Delta = 0.0249$; 8 values greater than 0.1) (Figures 4 and 5). The LD decay due to 10

350    generations of random crosses was approximately 25%, regardless of the population. For example,

351    the number of absolute LD values greater than 0.1 decreased 60% in population 1, generation 10r.

352    Compared to GWAS in population 1, generation 10r, at a FDR of 5%, the GWAS with RILs

353    from population 1, generation 10r (lowest parametric LD among the non-inbred populations), at a

354    FDR of 1%, showed the same power of QTL detection and a high number of significant

355    associations along the length of one or more chromosomes with one to four QTLs (Table 2). As

356    explained, this makes the GWAS ineffective for identifying candidate genes. Compared to GWAS

357    in generation 10r, the lower efficiency of GWAS with RILs for identifying candidate genes (due to

358    a greater number of significant associations in chromosomes with one to four QTLs) can be

359    attributed to higher heritability, due to increase in the genotypic variance for the same error

360    variance, and higher estimated LD. Based on simulation 1, the estimated QTL heritability with RILs

361    was approximately 9% for the three traits, assuming QTL heritability of 8% and 400 individuals

362    assessed in generation 10r (12.5% greater than the heritability at generation 10r). Due to sampling,

363    the estimated LD was greater with RILs than with non-inbred plants in generation 10r (Figure 6).

364    Based on simulation 1, the average estimated $\Delta$ and $r^2$ values were 0.0252 and 0.0241 for the RILs

365    and 0.0235 and 0.0225 for generation 10r, respectively. Although these average values are

366    equivalent, the estimated $\Delta$ values with RILs were four times greater on average than the estimated

367    $\Delta$ values in generation 10r. Once again, the GWAS was ineffective when assuming low heritability

368    and reduced sample size.

17

369 **DISCUSSION**

370       The presented theory proves that a significant association from a GWAS in a non-inbred or

371 inbred open-pollinated population and in an inbred lines panel, while controlling the type I error

372 rate and correcting for population structure and relatedness, is due to LD between the SNP and one

373 or more linked QTLs. The theory also shows that GWAS provides estimation of the average effect

374 of a SNP substitution (and consequently the estimation of SNP effects). Schaefer and Bernardo

375 (2013) estimated SNP effects for days to anthesis, days to silking, oil and starch concentration, and

376 measures of disease resistance using a maize inbred lines panel. We showed that only if there is a

377 single QTL in LD with a significant SNP it is adequate to test dominance for the QTL loci. It is

378 important to highlight that only if there is a single QTL in LD with a significant SNP, if the SNP is

379 within the QTL, and if QTL and SNP alleles have the same frequency it is adequate to consider the

380 SNP average effect of substitution as the QTL average effect of substitution. Furthermore, we also

381 proved that a significant association due to admixture LD (population structure) does not depend on

382 linkage disequilibrium between the SNP and a linked QTL.

383       To our knowledge, this is the first study on GWAS efficiency in open-pollinated population.

384 The results are very encouraging and show that the process can be highly efficient, depending on

385 LD, sample size, and QTL effect. In an open-pollinated population, the LD measure depends also

386 on the SNP and QTL allele frequencies. Thus, significant associations involving several SNPs with

387 the same QTL can be observed, including SNPs that are tens of mega base pairs (or centiMorgans)

388 from the QTL. In reality, a closely linked QTL and SNP can have a lower LD value compared to a

389 more distant QTL and SNP pair. In populations with low levels of LD, significant associations are

390 expected to occur for only SNPs within the QTL or located very close to the QTL (within a few

391 hundred base pairs), which favors the identification of a candidate gene for the QTL. In this

392 scenario, a QTL would be declared based on one to a small number of significant associations

393 spanning a chromosome region of a few kilo base pairs (not mega base pairs or centiMorgans).

18

394       A genome-wide association study is ineffective for lower sample size (200 individuals) and

395    QTL heritability (4%), regardless of the population, i.e., including inbred lines panel and RILs. This

396    scenario does not improve when increasing the FDR. Thus, we recommend that breeders employ

397    larger sample size (400 individuals) and achieve high trait heritability (70−80%) (such as by

398    genotyping parents and phenotyping replicated progeny). With intermediate (8%) to high (12%)

399    QTL heritability and larger sample size, it is important to define a FDR of 1% to decrease the

400    number of false-positive associations (note that some associations cannot really be false-positives).

401    According to Larsson *et al.* (2013), false-positive associations can arise from markers that are in

402    long-range LD with causative polymorphisms, which despite being rare are typically unaccounted

403    for in association studies.

404       Our results are comparable to those obtained by Yu *et al.* (2008) in a simulation study

405    investigating the genetic and statistical properties (power of QTL detection and FDR) of the nested

406    association mapping (NAM) design. With 5,000 genotypes, they achieved an average power of

407    QTL detection of 57% (with a range of 30 to 85%) when considering two trait heritabilities (0.4 and

408    0.7) and two different numbers of QTL controlling the trait (20 and 50). They also observed that a

409    higher heritability always gave higher QTL detection power, particularly for QTL with moderate to

410    small effect. Hung *et al.* (2012) assessed the maize NAM population and achieved heritabilities

411    greater than 0.8 for traits related to flowering time and plant architecture, resulting in a good power

412    to detect QTL. In contrast, traits with lower heritabilities (up to 0.6) and stronger sensitivity to

413    environmental variation allow only a reasonable power of QTL detection. Also using the NAM

414    population, Kump *et al.* (2011) evaluated resistance to southern leaf blight (SLB) disease and

415    obtained a heritability of 87%. They identified 32 QTLs with predominantly small and additive

416    effects on SLB resistance and many of the SNPs within and outside of QTL intervals were also

417    within or near to genes previously shown to be involved in plant disease resistance.

418          Field results have demonstrated that GWAS are best carried out with a large sample size (Yu

419     and Buckler 2006). According to Flint-Garcia *et al.* (2005), increasing the population size increases

420     the number of individuals with rare alleles, thus improving the power to test the association between

421     these rare alleles and the trait of interest. Yu *et al.* (2008) showed that the gain in efficiency by

422     increasing sample size was evidenced by increased power of QTL detection and smaller FDR,

423     mainly with heritability of 0.7 in comparison with a heritability of 0.4. Based on a simulation study,

424     Long and Langley (1999) demonstrated that approximately 500 individuals should be genotyped for

425     20 SNP loci within the candidate gene region to detect marker-trait associations for QTLs that

426     account for as little as 5% of the phenotypic variation. They observed that more power was

427     achieved by increasing the population size than by increasing the SNP density within the candidate

428     gene.

429          Compared to QTL mapping, GWAS is much more precise for mapping QTLs and identifying

430     candidate genes. In QTL mapping studies based on simulated data, the bias in QTL position ranged

431     from 2.0 to 6.0 cM depending on sample size, heritability, and marker density (Li *et al.* 2010). The

432     bias with GWAS should be much lower because the significant SNPs are frequently within or very

433     close to the candidate genes.

434          Compared to GWAS in an inbred lines panel, GWAS in open-pollinated population was less

435     efficient, i.e., showed slightly lower power of QTL detection, higher number of false-positive

436     associations, slightly higher bias in QTL position, and higher number of significant associations for

437     the same QTL. The increase in efficiency by using the inbred lines panel was due to the lower

438     degree of LD achieved by mixing groups of inbreds with positive and negative LD values. This is

439     probably the main advantage of GWAS based on inbred lines panel. In contrast, when fixing the

440     FDR, GWAS in a non-inbred population tends to be more efficient than GWAS with RILs from the

441     same population. According to Flint-Garcia *et al.* (2005), the inbred lines panel exploits the rapid

442     breakdown of LD in diverse maize lines, enabling very high resolution for QTL mapping.

443     Population structure results from constructing a panel with inbreds from various breeding programs

444     and distinct heterotic groups, which can cause false-positive marker-trait associations if the data is

445     not corrected (Yan *et al.* 2009). The lowest parametric LD values for the inbred lines panel occurred

446     in published studies (Yan *et al.* 2009, Remington *et al.* 2001). Moreover, with the inbred lines

447     panel, generally, only SNP loci within the QTL showed significant association, which is a

448     highlighted result from GWAS that can serve as a basis for a fine mapping strategy for marker-

449     assisted selection and map-based cloning genes (Gupta *et al.* 2005).

450     GWAS in plant breeding has been effective for identifying candidate genes for quantitative

451     traits such as plant architecture, kernel composition, root development, flowering time, drought

452     tolerance, pathogen resistance, and metabolic processes (Zhu *et al.* 2008). Based on our evidence,

453     breeders can employ non-inbred and inbred breeding populations while taking into account that the

454     level of LD should be low, the sample size should be higher than that necessary for QTL mapping,

455     and the QTL heritability should be intermediate to high to achieve greater power of QTL detection

456     and precise mapping of candidate genes.

457    

461     **LITERATURE CITED**

462     Azevedo, C. F., M. D. V. Resende, F. F. Silva, J. M. S. Viana, M. S. F. Valente *et al.*, 2015 Ridge,

463        Lasso and Bayesian additive-dominance genomic models. BMC Genet. 16: 105-118.

464     Barendse, W., A. Reverter, R. J. Bunch, B. E. Harrison, W. Barris *et al.*, 2007 A validated whole-

465        genome association study of efficient food conversion in cattle. Genetics 176: 1893-905.

466     Benjamini, Y., and Y. Hochberg, 1995 Controlling the false discovery rate: a practical and powerful

467        approach to multiple testing. J. R. Stat. Soc. 57: 289-300.

468  Bernardo, R., 2013 Genomewide markers for controlling background variation in association

469      mapping. Plant Genome 6.

470  Bradbury, P. J., Z. Zhang, D. E. Kroon, T. M. Casstevens, Y. Ramdoss *et al.*, 2007 TASSEL:

471      software for association mapping of complex traits in diverse samples. Bioinformatics 23: 2633-

472      2635.

473  Falush, D., M. Stephens, and J. K. Pritchard, 2003 Inference of population structure using

474      multilocus genotype data: linked loci and correlated allele frequencies. Genetics 164: 1567-1587.

475  Flint-Garcia, S. A., J. M. Thornsberry, and E. S. Buckler, 2003 Structure of linkage disequilibrium

476      in plants. Annu. Rev. Plant Biol. 54: 357-374.

477  Flint-Garcia, S. A., A. C. Thuillet, J. Yu, G. Pressoir, S. M. Romero *et al.*, 2005 Maize association

478      population: a high-resolution platform for quantitative trait locus dissection. Plant J. 44: 1054-

479      1064.

480  Gupta, P. K., S. Rustgi, and P. L. Kulwal, 2005 Linkage disequilibrium and association studies in

481      higher plants: present status and future prospects. Plant Mol. Biol. 57: 461-485.

482  Hill, W. G., and A. Robertson, 1968 Linkage disequilibrium in finite populations. Theor. Appl.

483      Genet. 38: 226-231.

484  Hung, H. Y., C. Browne, K. Guill, N. Coles, M. Eller *et al.*, 2012 The relationship between parental

485      genetic or phenotypic divergence and progeny variation in the maize nested association mapping

486      population. Heredity 108: 490-499.

487  Ingvarsson, P. K., and N. R. Street, 2011 Association genetics of complex traits in plants. New

488      Phytol. 189: 909-922.

489  Kempthorne, O., 1957 *An Introduction to Genetic Statistics*. John Wiley and Sons Inc., New York.

490  Kump, K. L., P. J. Bradbury, R. J. Wisser, E. S. Buckler, A. R. Belcher *et al.*, 2011 Genome-wide

491      association study of quantitative resistance to southern leaf blight in the maize nested association

492      mapping population. Nat. Genet. 43: 163-169.

493    Larsson, S. J., A. E. Lipka, and E. S. Buckler, 2013 Lessons from Dwarf8 on the strengths and

494        weaknesses of structured association mapping. Plos Genet. 9: e1003246.

495    Li, H., S. Hearne, M. Banziger, Z. Li, and J. Wang, 2010 Statistical properties of QTL linkage

496        mapping in biparental genetic populations. Heredity 105: 257-267.

497    Liu, K., and S. V. Muse, 2005 PowerMarker: integrated analysis environment for genetic marker

498        data. Bioinformatics 21: 2128-2129.

499    Long, A. D., and C. H. Langley, 1999 The power of association studies to detect the contribution of

500        candidate genetic loci to variation in complex traits. Genome Res. 9: 720-731.

501    Pace, J., C. Gardner, C. Romay, B. Ganapathysubramanian, and T. Lübberstedt, 2015 Genome-wide

502        association analysis of seedling root development in maize (*Zea mays* L.). BMC Genomics 16:

503        47-58.

504    Pearson, T. A., and T. A. Manolio, 2008 How to interpret a genome-wide association study. J. Am.

505        Med. Assoc. 299: 1335-1344.

506    Rafalski, J. A., 2010 Association genetics in crop improvement. Curr. Opin. Plant Biol. 13: 174-

507        180.

508    Remington, D. L., J. M. Thornsberry, Y. Matsuoka, L. M. Wilson, S. R. Whitt *et al.*, 2001 Structure

509        of linkage disequilibrium and phenotypic associations in the maize genome. PNAS 98: 11479-

510        11484.

511    Schaefer, C. M., and R. Bernardo, 2013 Genome-wide association mapping of flowering time,

512        kernel composition, and disease resistance in historical Minnesota maize inbreds. Crop Sci. 53:

513        2518-2529.

514    Stich, B., and A. E. Melchinger, 2009 Comparison of mixed-model approaches for association

515        mapping in rapeseed, potato, sugar beet, maize, and Arabidopsis. BMC Genomics 10: e94.

516   Thirunavukkarasu, N., F. Hossain, K. Arora, R. Sharma, K. Shiriga *et al.*, 2014 Functional
517       mechanisms of drought tolerance in subtropical maize (*Zea mays* L.) identified using genome-
518       wide association mapping. BMC Genomics 15: 1182-1193.

519   Viana, J. M. S., H.-P. Piepho, and F. F. Silva, 2016 Quantitative genetics theory for genomic
520       selection and efficiency of breeding value prediction in open-pollinated populations. Sci. Agric.
521       73: 243-251.

522   Viana, J. M. S., M. S. F. Valente, F. F. Silva, G. B. Mundim, and G. P. Paes, 2013 Efficacy of
523       population structure analysis with breeding populations and inbred lines. Genetica 141: 389-399.

524   Weir, B. S., 2008 Linkage disequilibrium and association mapping. Ann. Rev. Genomics Hum.
525       Genet. 9: 129-142.

526   Weir, B., 2010 Statistical genetic issues for genome-wide association studies. Genome 53: 869-875.

527   Yan, J. B., T. Shah, M. Warburton, E. S. Buckler, M. D. McMullen *et al.*, 2009 Genetic
528       characterization of a global maize collection using SNP markers. Plos One 4: e8451.

529   Yang, X., J. Yan, T. Shah, M. L. Warbuton, Q. Li *et al.*, 2010 Genetic analysis and characterization
530       of a new maize association mapping panel for quantitative trait loci dissection. Theor. Appl.
531       Genet. 121: 417-431.

532   Yu, J., and E. S. Buckler, 2006 Genetic association mapping and genome organization of maize.
533       Curr. Opin. Biotechnol. 17: 1-6.

534   Yu, J., J. B. Holland, M. D. McMullen, and E. S. Buckler, 2008 Genetic design and statistical
535       power of nested association mapping in maize. Genetics 178: 539-551.

536   Yu, J., G. Pressoir, W. H. Briggs, I. Vroh Bi, M. Yamasaki *et al.*, 2006 A unified mixed-model
537       method for association mapping that accounts for multiple levels of relatedness. Nat. Genet. 38:
538       203-208.

539   Zhu, C., M. Gore, E. S. Buckler, and J. Yu, 2008 Status and prospects of association mapping in
540       plants. Plant Genome 1: 5-20.

**Table 1** Average number of significant associations with a FDR of 5 or 1%, power of QTL detection (%), number of false-positive associations in chromosomes with no QTL and one to four QTLs, bias in the QTL position (cM), and average range for the regions with identified QTL, regarding population 1, generation 10r (random cross), three traits (expansion volume (EV; mL/g), grain yield (GY; g per plant), and days to maturity (DM)), two sample sizes, and three QTL heritabilities[a]
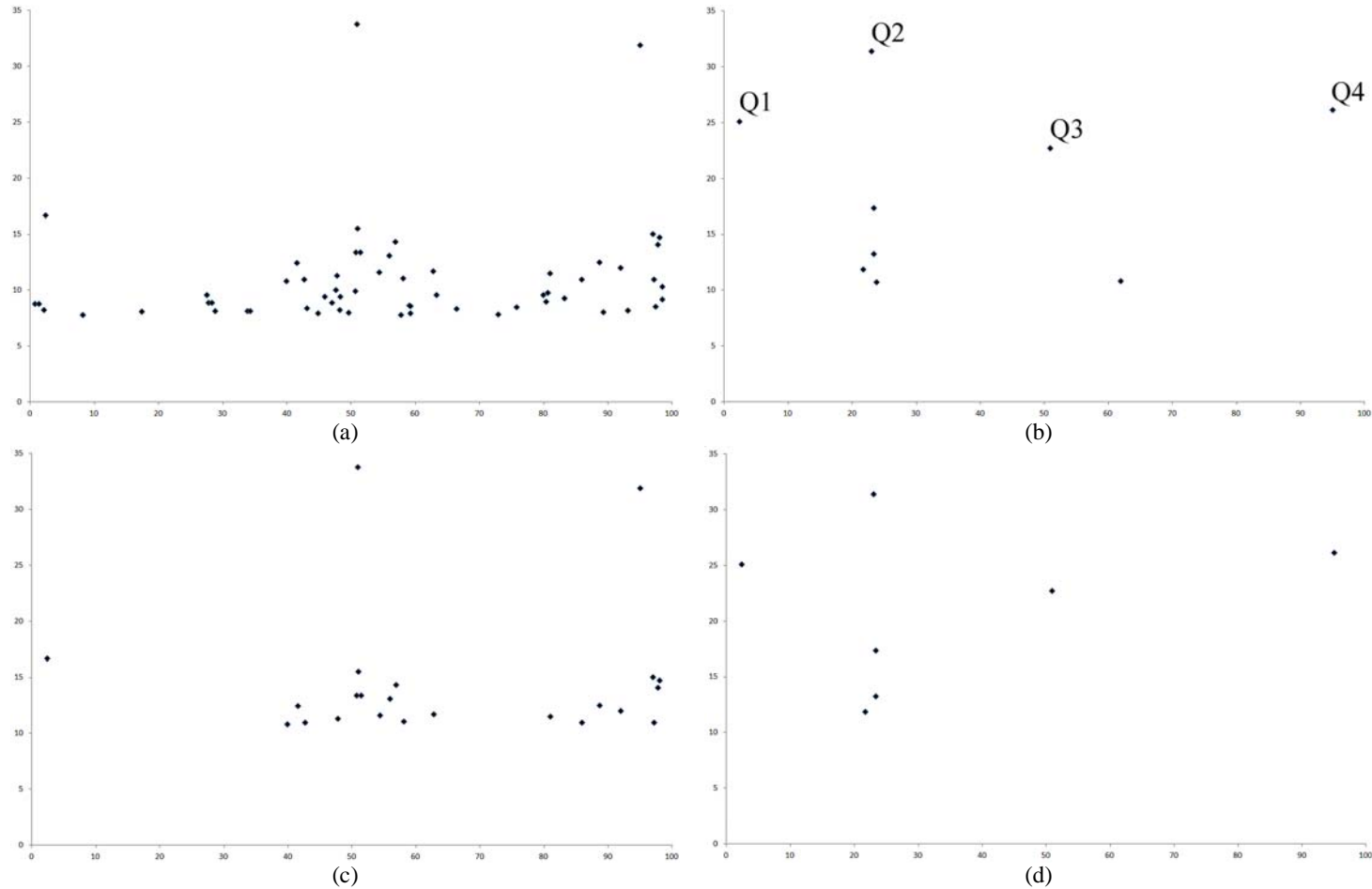
| FDR | Trait | Sample | $h^2$ | Sig. Assoc. | Power | False+0 | False+1-4 | Bias | Av. range |
|---|---|---|---|---|---|---|---|---|---|
| 5% | EV | 400 | 4% | 6.7 (0; 22) | 37.2 (0; 80) | 0.6 (0; 5) | 0.9 (0; 8) | 0.21 (0.00; 0.98) | 0.28 (0.00; 2.45) |
| | | | 8% | 32.1 (13; 73) | 88.6 (60; 100) | 3.3 (0; 17) | 7.9 (1; 28) | 0.52 (0.09; 0.83) | 0.84 (0.11; 2.18) |
| | | | 12% | 26.8 (7; 63) | 99.2 (60; 100) | 2.9 (0; 11) | 8.1 (0; 30) | 0.59 (0.00; 0.92) | 1.27 (0.00; 2.74) |
| | | 200 | 4% | 0.8 (0; 5) | 5.2 (0; 40) | 0.2 (0; 2) | 0.3 (0; 1) | 0.22 (0.00; 1.76) | 0.14 (0.00; 1.31) |
| | | | 8% | 6.2 (0; 33) | 37.0 (0; 70) | 0.5 (0; 3) | 1.1 (0; 12) | 0.16 (0.00; 0.88) | 0.17 (0.00; 1.42) |
| | | | 12% | 6.0 (2; 20) | 70.8 (20; 100) | 0.5 (0; 4) | 0.8 (0; 5) | 0.18 (0.00; 0.87) | 0.25 (0.00; 1.58) |
| | GY | 400 | 4% | 5.7 (0; 25) | 34.4 (0; 90) | 0.4 (0; 2) | 0.8 (0; 7) | 0.20 (0.00; 0.85) | 0.18 (0.00; 1.03) |
| | | | 8% | 31.3 (10; 82) | 87.8 (70; 100) | 3.3 (0; 12) | 7.5 (0; 28) | 0.51 (0.00; 0.97) | 0.77 (0.00; 1.54) |
| | | | 12% | 31.6 (8; 74) | 99.6 (80; 100) | 3.6 (0; 16) | 9.8 (0; 43) | 0.68 (0.14; 1.01) | 1.39 (0.12; 2.60) |
| | | 200 | 4% | 0.8 (0; 16) | 3.8 (0; 30) | 0.6 (0; 4) | 0.4 (0; 6) | 0.15 (0.00; 1.39) | 0.06 (0.00; 0.94) |
| | | | 8% | 5.9 (0; 18) | 32.6 (10; 80) | 0.8 (0; 8) | 1.1 (0; 7) | 0.16 (0.00; 0.97) | 0.16 (0.00; 1.75) |
| | | | 12% | 7.2 (3; 19) | 74.8 (20; 100) | 0.7 (0; 7) | 1.2 (0; 6) | 0.21 (0.00; 0.75) | 0.30 (0.00; 1.49) |
| | DM | 400 | 4% | 8.6 (0; 33) | 39.0 (0; 100) | 1.1 (0; 5) | 1.4 (0; 9) | 0.29 (0.00; 0.70) | 0.38 (0.00; 1.64) |
| | | | 8% | 50.7 (12; 119) | 92.8 (70; 100) | 6.5 (0; 21) | 14.9 (3; 48) | 0.65 (0.06; 1.02) | 1.22 (0.07; 2.82) |
| | | | 12% | 50.2 (20; 110) | 100.0 (100; 100) | 6.5 (0; 23) | 18.4 (3; 50) | 0.77 (0.49; 1.05) | 2.02 (0.76; 3.49) |
| | | 200 | 4% | 1.0 (0; 7) | 5.0 (0; 30) | 0.5 (0; 2) | 0.5 (0; 3) | 0.17 (0.00; 0.82) | 0.15 (0.00; 2.45) |
| | | | 8% | 6.6 (0; 41) | 34.0 (0; 70) | 0.7 (0; 3) | 1.2 (0; 12) | 0.23 (0.00; 0.91) | 0.27 (0.00; 2.00) |
| | | | 12% | 8.8 (2; 31) | 75.2 (40; 100) | 1.0 (0; 9) | 1.9 (0; 12) | 0.26 (0.00; 0.96) | 0.35 (0.00; 1.38) |
| 1% | EV | 400 | 8% | 15.0 (7; 39) | 76.6 (50; 100) | 0.4 (0; 3) | 2.0 (0; 13) | 0.32 (0.00; 0.75) | 0.41 (0.00; 1.70) |
| | | | 12% | 13.3 (4; 38) | 98.4 (60; 100) | 0.5 (0; 2) | 2.7 (0; 16) | 0.40 (0.00; 0.79) | 0.65 (0.00; 1.92) |
| | GY | 400 | 8% | 14.8 (6; 39) | 75.4 (60; 100) | 0.3 (0; 3) | 2.1 (0; 10) | 0.33 (0.00; 0.80) | 0.43 (0.00; 1.32) |
| | | | 12% | 15.6 (5; 43) | 98.4 (80; 100) | 0.6 (0; 7) | 3.2 (0; 21) | 0.53 (0.03; 0.81) | 0.82 (0.03; 1.63) |
| | DM | 400 | 8% | 20.6 (6; 51) | 80.0 (40; 100) | 0.6 (0; 3) | 4.1 (0; 19) | 0.44 (0.00; 0.93) | 0.61 (0.00; 1.99) |
| | | | 12% | 22.3 (8; 49) | 100.0 (100; 100) | 0.9 (0; 6) | 6.5 (0; 21) | 0.58 (0.02; 0.99) | 1.14 (0.03; 2.98) |

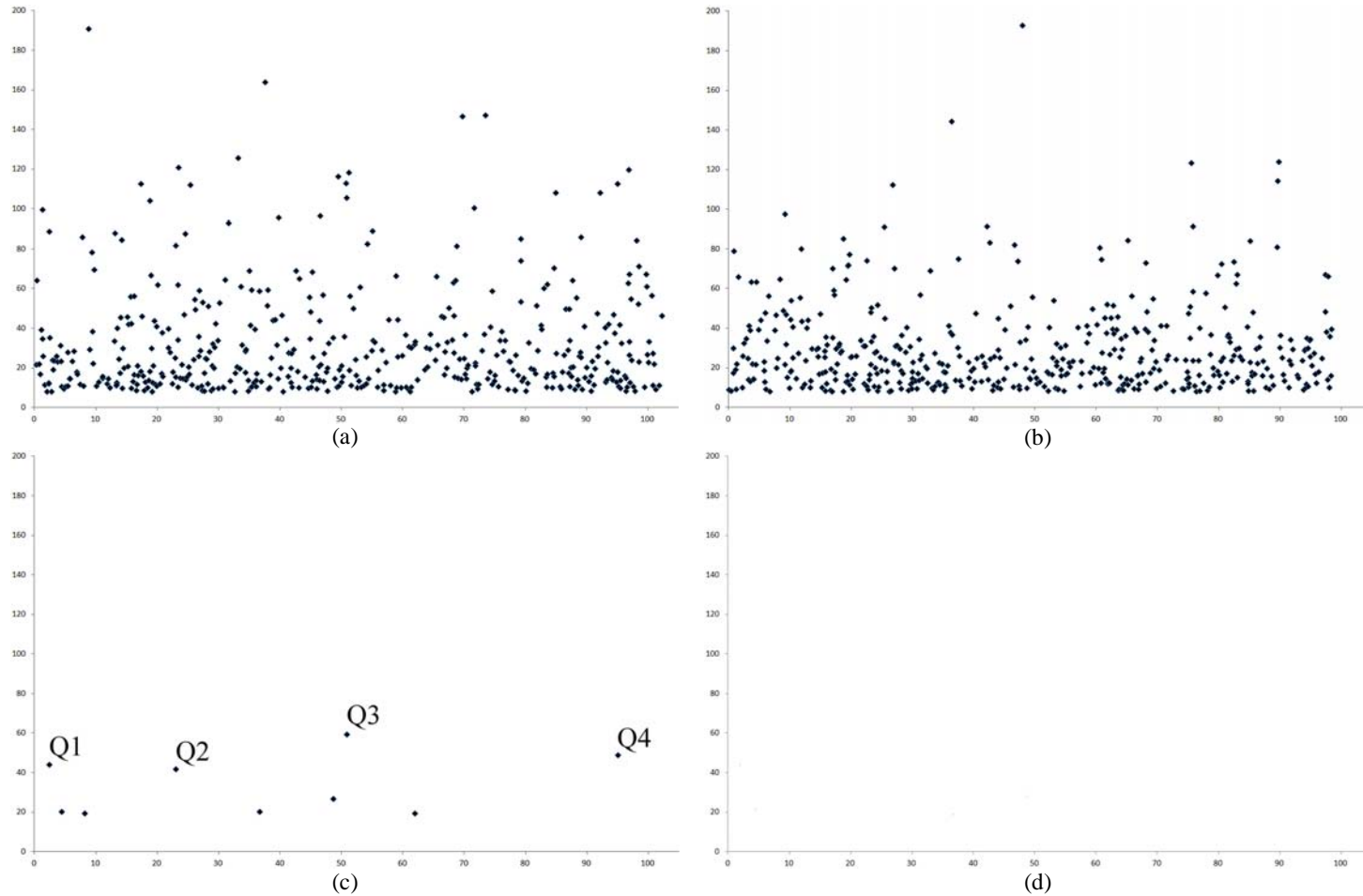[a]the values between parentheses are the minimum and maximum.

25

**Table 2** Average number of significant associations with a FDR of 5 or 1%, power of QTL detection (%), number of false-positive associations in chromosomes with no QTL and one to four QTLs, bias in the QTL position (cM), and average range for the regions with identified QTL, regarding an inbred lines panel and RILs from population 1, generation 10r (random cross), three traits (expansion volume (EV; mL/g), grain yield (GY; g per plant), and days to maturity (DM)), two sample sizes, and two QTL heritabilities[a]

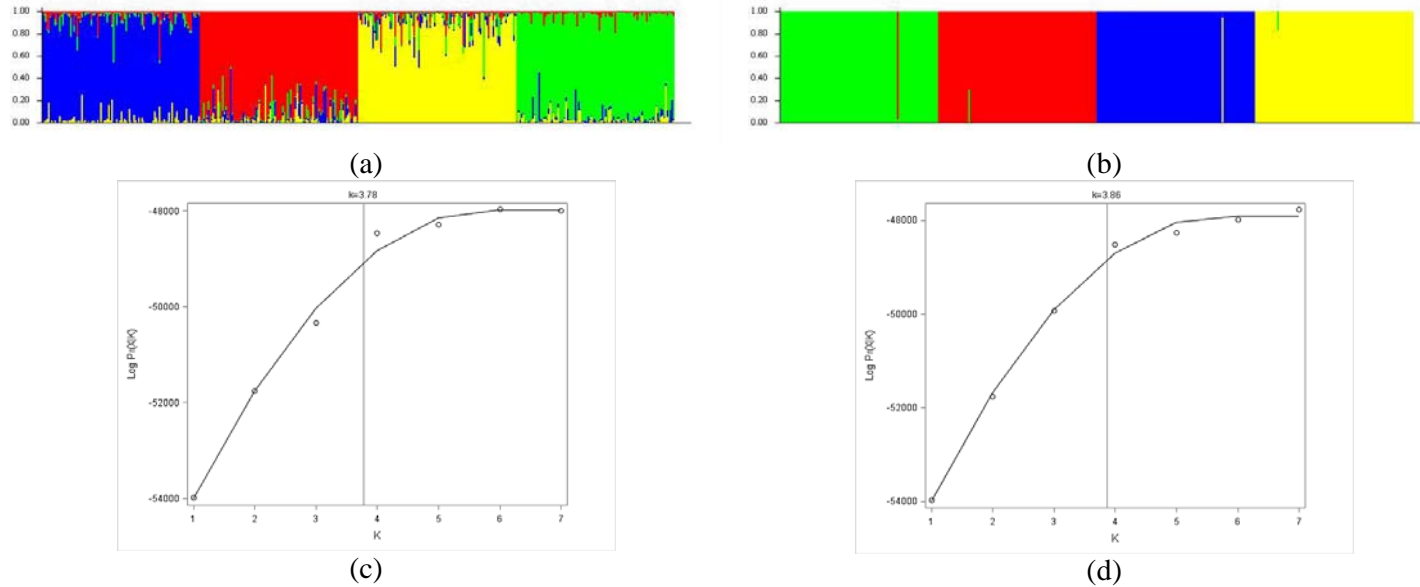| Population | FDR | Trait | Sample | h2 | Sig. Assoc. | Power | False+0 | False+1-4 | Bias | Av. range |
|---|---|---|---|---|---|---|---|---|---|---|
| Inbred lines panel | 5% | EV | 400 | 4% | 6.9 (2; 14) | 58.0 (20; 100) | 0.3 (0; 2) | 0.6 (0; 4) | 0.04 (0.00; 0.35) | 0.04 (0.00; 0.42) |
| | | | | 8% | 14.6 (9; 27) | 96.0 (90; 100) | 0.5 (0; 2) | 3.1 (0; 13) | 0.12 (0.00; 0.70) | 0.14 (0.00; 0.58) |
| | | | 200 | 4% | 0.7 (0; 5) | 5.4 (0; 40) | 0.1 (0; 1) | 0.4 (0; 1) | 0.09 (0.00; 0.92) | 0.06 (0.00; 1.12) |
| | | | | 8% | 5.4 (0; 10) | 45.2 (0; 80) | 0.1 (0; 1) | 0.6 (0; 4) | 0.04 (0.00; 0.69) | 0.05 (0.00; 1.04) |
| | | GY | 400 | 4% | 7.9 (2; 17) | 61.4 (20; 100) | 0.4 (0; 3) | 1.0 (0; 5) | 0.05 (0.00; 0.28) | 0.05 (0.00; 0.33) |
| | | | | 8% | 13.9 (7; 23) | 96.2 (70; 100) | 0.3 (0; 2) | 2.8 (0; 11) | 0.11 (0.00; 0.41) | 0.12 (0.00; 0.50) |
| | | | 200 | 4% | 1.0 (0; 5) | 8.2 (0; 30) | 0.2 (0; 2) | 0.2 (0; 2) | 0.00 (0.00; 0.00) | 0.00 (0.00; 0.00) |
| | | | | 8% | 5.1 (0; 13) | 42.6 (0; 70) | 0.2 (0; 2) | 0.5 (0; 5) | 0.05 (0.00; 0.72) | 0.04 (0.00; 1.04) |
| | | DM | 400 | 4% | 9.1 (4; 15) | 70.8 (40; 90) | 0.4 (0; 2) | 1.2 (0; 5) | 0.06 (0.00; 0.62) | 0.05 (0.00; 0.52) |
| | | | | 8% | 15.4 (8; 29) | 96.0 (80; 100) | 0.5 (0; 3) | 3.7 (0; 17) | 0.14 (0.00; 0.43) | 0.17 (0.00; 0.57) |
| | | | 200 | 4% | 1.1 (0; 6) | 9.6 (0; 40) | 0.1 (0; 1) | 0.1 (0; 1) | 0.01 (0.00; 0.23) | 0.00 (0.00; 0.05) |
| | | | | 8% | 5.6 (0; 12) | 46.8 (0; 80) | 0.1 (0; 2) | 0.6 (0; 4) | 0.05 (0.00; 0.53) | 0.04 (0.00; 0.54) |
| | 1% | EV | 400 | 8% | 10.8 (6; 21) | 91.6 (60; 100) | 0.1 (0; 1) | 1.0 (0; 10) | 0.05 (0.00; 0.63) | 0.06 (0.00; 0.61) |
| | | GY | 400 | 8% | 10.2 (7; 15) | 91.2 (70; 100) | 0.0 (0; 1) | 0.7 (0; 6) | 0.03 (0.00; 0.23) | 0.03 (0.00; 0.26) |
| | | DM | 400 | 8% | 10.7 (7; 16) | 91.6 (70; 100) | 0.0 (0; 1) | 1.0 (0; 5) | 0.07 (0.00; 0.39) | 0.07 (0.00; 0.47) |
| RILs | 1% | EV | 400 | 4% | 7.3 (1; 31) | 39.0 (10; 70) | 0.1 (0; 2) | 1.7 (0; 12) | 0.17 (0.00; 0.73) | 0.25 (0.00; 1.25) |
| | | | | 8% | 34.5 (4; 122) | 87.0 (40; 100) | 0.3 (0; 2) | 12.7 (0; 68) | 0.61 (0.00; 1.05) | 0.90 (0.00; 2.43) |
| | | | 200 | 8% | 4.9 (0; 24) | 27.0 (0; 70) | 0.1 (0; 2) | 1.1 (0; 9) | 0.15 (0.00; 1.15) | 0.20 (0.00; 1.42) |
| | | GY | 400 | 4% | 9.5 (2; 42) | 43.0 (10; 70) | 0.1 (0; 1) | 2.8 (0; 23) | 0.30 (0.00; 1.09) | 0.41 (0.00; 2.00) |
| | | | | 8% | 34.1 (5; 123) | 86.0 (50; 100) | 0.2 (0; 2) | 12.9 (0; 73) | 0.60 (0.00; 1.05) | 0.88 (0.00; 2.29) |
| | | | 200 | 8% | 6.3 (0; 27) | 30.0 (0; 60) | 0.1 (0; 4) | 1.8 (0; 16) | 0.22 (0.00; 1.29) | 0.30 (0.00; 2.71) |
| | | DM | 400 | 4% | 16.6 (4; 65) | 61.0 (30; 100) | 0.2 (0; 2) | 5.2 (0; 44) | 0.47 (0.00; 1.09) | 0.62 (0.00; 1.94) |
| | | | | 8% | 40.4 (11; 142) | 89.0 (70; 100) | 0.4 (0; 3) | 15.8 (0; 81) | 0.66 (0.12; 1.07) | 1.01 (0.13; 2.62) |
| | | | 200 | 8% | 7.4 (0; 42) | 32.0 (0; 80) | 0.2 (0; 3) | 2.1 (0; 15) | 0.26 (0.00; 1.21) | 0.36 (0.00; 2.11) |
| | 5% | EV | 200 | 4% | 1.8 (0; 34) | 7.0 (0; 30) | 0.2 (0; 1) | 1.2 (0; 23) | 0.19 (0.00; 1.12) | 0.33 (0.00; 2.47) |
| | | GY | 200 | 4% | 2.9 (0; 19) | 11.0 (0; 30) | 0.3 (0; 1) | 1.5 (0; 9) | 0.41 (0.00; 1.77) | 0.67 (0.00; 3.59) |
| | | DM | 200 | 4% | 4.6 (0; 33) | 17.0 (0; 50) | 0.3 (0; 2) | 2.0 (0; 19) | 0.32 (0.00; 1.26) | 0.37 (0.00; 1.81) |

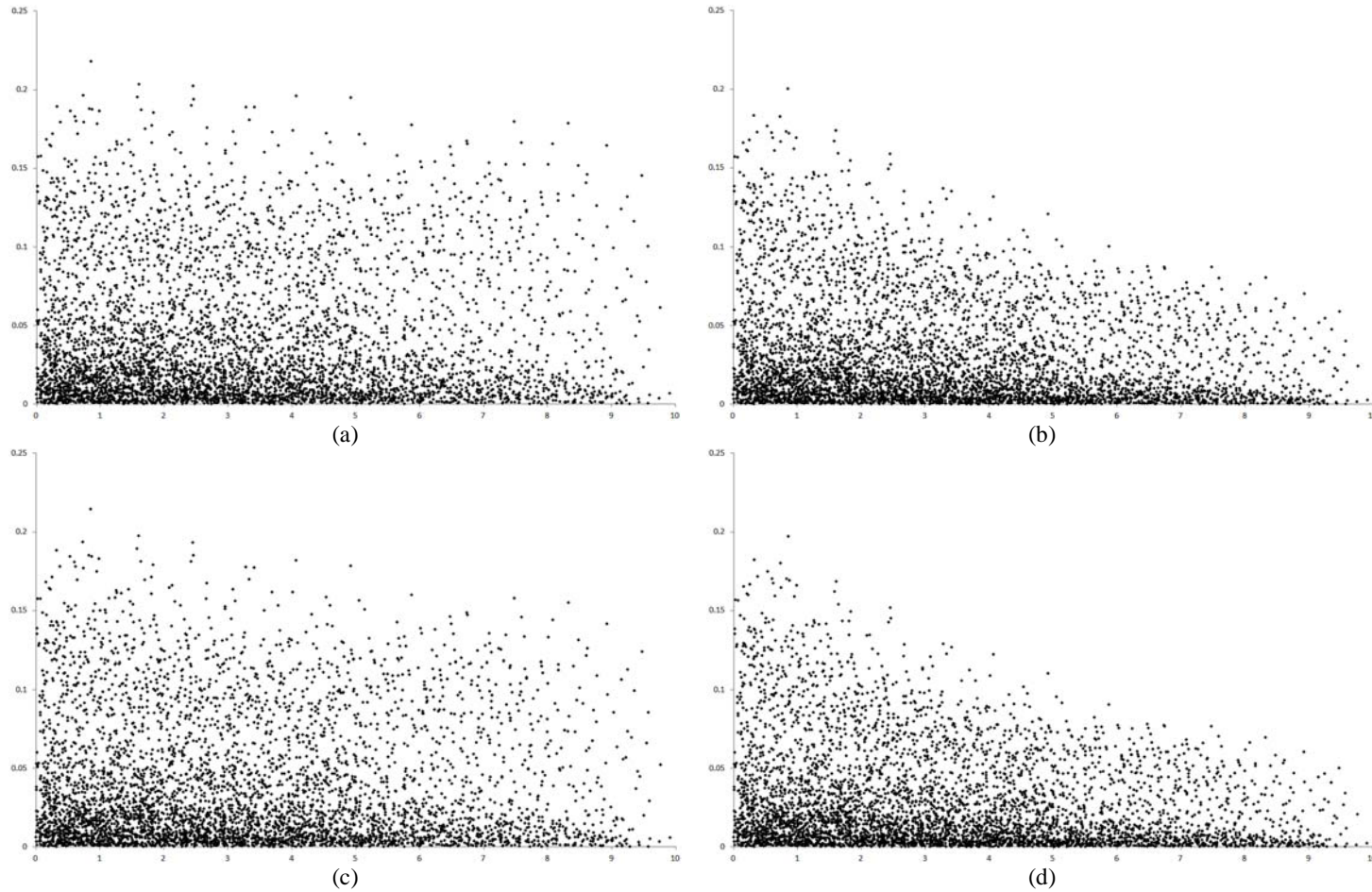[a]the values between parentheses are the minimum and maximum.

26

550 **Figure 1** Significant associations at a FDR of 5 (a and b) or 1% (c and d) (F test; Y axe) in chromosome 1 (SNP position (cM); X axe), from the
551 GWAS in population 1, generations 0 (a and c) and 10r (random cross) (b and d), regarding expansion volume, QTL heritability of 8%, and sample
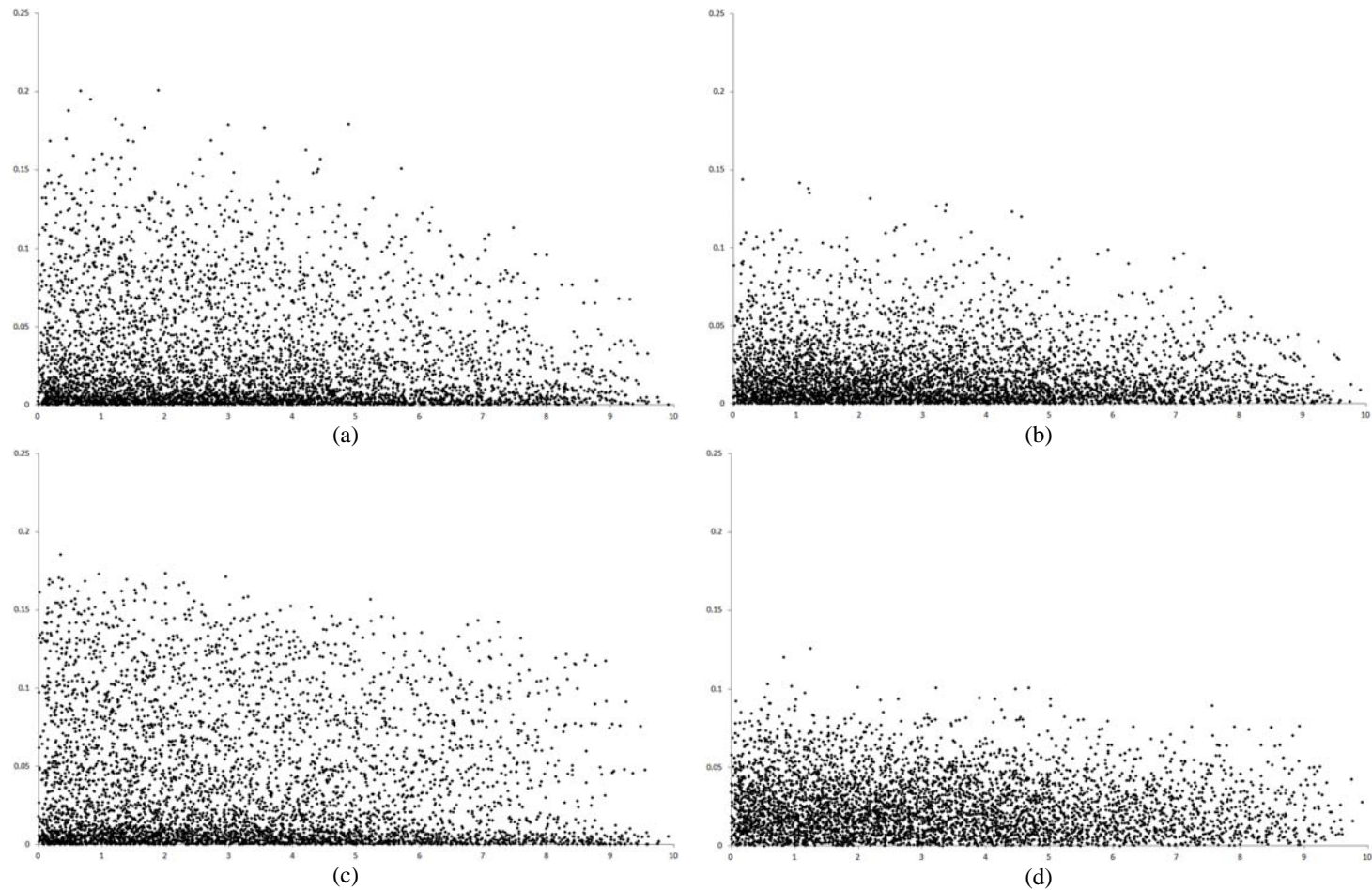552 size 400 (simulation 1) (Q = QTL).

553 **Figure 2** Significant associations at a FDR of 1% (F test; Y axe) in chromosomes 1 and 3 (SNP position (cM); X axe) ignoring (a and b, respectively)
554 and correcting for the population structure (c and d, respectively), from the GWAS in an inbred lines panel regarding expansion volume, QTL
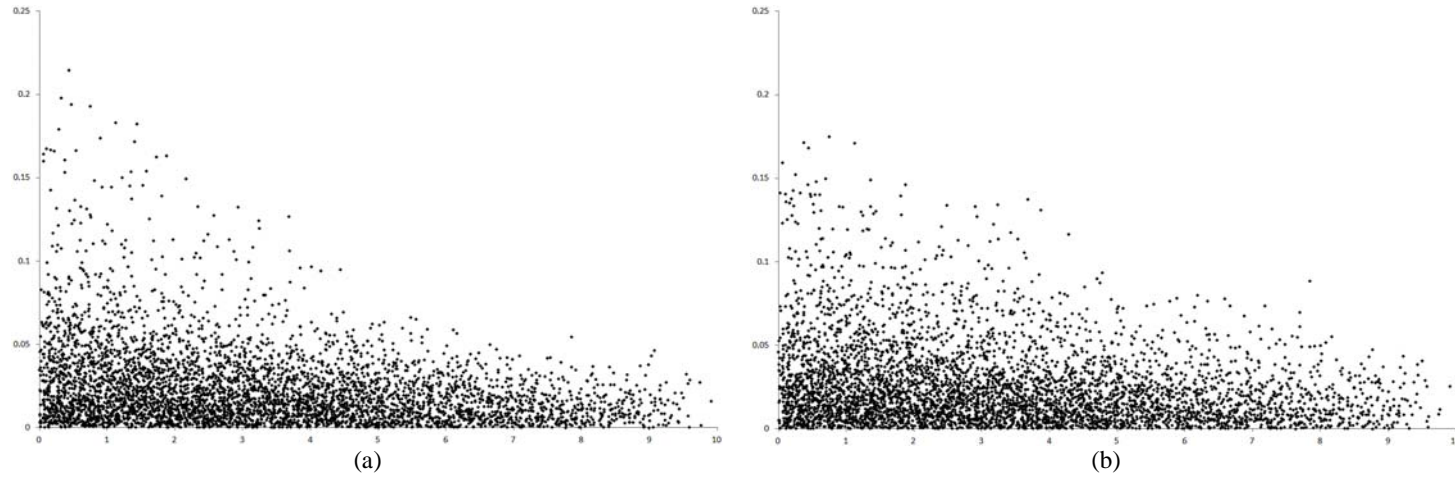555 heritability of 8%, and sample size 400 (simulation 1) (Q = QTL).

(a)

(b)

(c)

(d)

**Figure 3** Results from the population structure analysis and the inferred plateau method, based on the admixture model with correlated allelic frequencies (a and c) and the no admixture model with independent allelic frequencies (b and d).

558 **Figure 4** Relationship between the parametric LD value (absolute value; Y axe) and distance (cM; X axe) in population 1, generations 0 (a), 10r
559 (random cross) (b), 10s (selfing) (c), and 10r10s (d), assuming a segment of 10 cM of chromosome 1 (centered on QTL 3).

**Figure 5** Relationship between the parametric LD value (absolute value; Y axe) and distance (cM; X axe) in populations 2 (a), 3 (b), and 4 (c), generation 10s (selfing), and in the inbred lines panel (d), assuming a segment of 10 cM of chromosome 1 (centered on QTL 3).

(a)
(b)

562  **Figure 6** Relationship between the estimated LD value (absolute value; Y axe) and distance (cM; X axe) in population 1, generations 10r (random
563  cross) (a) and 10r10s (random cross and selfing) (b), simulation 1, assuming a segment of 10 cM of chromosome 1 (centered on QTL 3).