

# Effect of population size and unbalanced data sets on QTL detection using genome-wide association mapping in barley breeding germplasm

Hongyun Wang · Kevin P. Smith · Emily Combs ·  
Tom Blake · Richard D. Horsley · Gary J. Muehlbauer

Received: 29 April 2011 / Accepted: 17 August 2011 / Published online: 7 September 2011  
© Springer-Verlag 2011

**Abstract** Over the past two decades many quantitative trait loci (QTL) have been detected; however, very few have been incorporated into breeding programs. The recent development of genome-wide association studies (GWAS) in plants provides the opportunity to detect QTL in germplasm collections such as unstructured populations from breeding programs. The overall goal of the barley Coordinated Agricultural Project was to conduct GWAS with the intent to couple QTL detection and breeding. The basic idea is that breeding programs generate a vast amount of phenotypic data and combined with cheap genotyping it should be possible to use GWAS to detect QTL that would be immediately accessible and used by breeding programs. There are several constraints to using breeding program-derived phenotype data for conducting GWAS namely: limited population size and unbalanced data sets. We chose the highly heritable trait heading date to study these two variables. We examined 766 spring barley breeding lines (panel #1) grown in balanced trials and a subset of 384

spring barley breeding lines (panel #2) grown in balanced and unbalanced trials. In panel #1, we detected three major QTL for heading date that have been detected in previous bi-parental mapping studies. Simulation studies showed that population sizes greater than 384 individuals are required to consistently detect QTL. We also showed that unbalanced data sets from panel #2 can be used to detect the three major QTL. However, unbalanced data sets resulted in an increase in the false-positive rate. Interestingly, one-step analysis performed better than two-step analysis in reducing the false-positive rate. The results of this work show that it is possible to use phenotypic data from breeding programs to detect QTL, but that careful consideration of population size and experimental design are required.

## Introduction

Quantitative trait loci (QTL) mapping has been primarily conducted in bi-parental populations segregating for the trait of interest. In many crop species, hundreds of QTL have been detected using bi-parental populations. However, there are limitations to this approach including the need to create a population segregating for each trait to be mapped, the ability to assess only two alleles per locus, and the limited number of meioses within a single population. An alternative approach is to conduct association mapping in a large germplasm collection (Thornsberry et al. 2001). Association mapping is based on linkage disequilibrium (LD) or the non-independence of alleles in a population (Gaut and Long 2003). This approach circumvents the need for constructing genetic mapping populations for each trait of interest and instead utilizes all of the recombination events that have occurred throughout the evolutionary

Communicated by R. Waugh.

**Electronic supplementary material** The online version of this article (doi:10.1007/s00122-011-1691-8) contains supplementary material, which is available to authorized users.

H. Wang · K. P. Smith · E. Combs · G. J. Muehlbauer (✉)  
Department of Agronomy and Plant Genetics,  
University of Minnesota, St. Paul, MN 55108, USA  
e-mail: muehl003@umn.edu

T. Blake  
Department of Plant Sciences and Plant Pathology,  
Montana State University, Bozeman, MT 59717, USA

R. D. Horsley  
Department of Plant Sciences, North Dakota State University,  
Fargo, ND 58108, USA

history of the plant. Thus, association mapping considers many more informative meioses than traditional bi-parental mapping.

Association mapping in plants of candidate genes and genome-wide association studies (GWAS) have successfully identified associations of marker alleles with traits (e.g., Thornsberry et al. 2001; Zhao et al. 2007; Beló et al. 2008; Wang et al. 2008; Breseghello and Sorrells 2006; Pressoir et al. 2009; Wilson et al. 2004; Cockram et al. 2010). In barley, GWAS have been conducted to detect QTL in elite germplasm for yield and agronomic traits (Kraakman et al. 2004, 2006), Fusarium head blight resistance (Massman et al. 2011), winter hardiness (von Zitzewitz et al. 2011), and growth habit and inflorescence type (Cuesta-Marcos et al. 2010). GWAS were conducted to detect QTL for spot blotch resistance in a wild barley collection (Roy et al. 2010). In addition, GWAS were conducted as the starting point to identify markers associated with lateral floret fertility, which directly led to the isolation of the *INTERMEDIUM-C* gene (Ramsay et al. 2011).

One goal of QTL mapping is to identify marker-trait associations that can be used for marker-assisted selection (MAS) within breeding programs. Although thousands of QTL have been detected in crop species, the vast majority have not been deployed in breeding programs (Bernardo 2008). The lack of QTL deployment has been due to several issues including: (1) the lack of high-throughput marker technology; (2) the beneficial QTL allele is already fixed in the breeding germplasm; (3) linkage drag; (4) the marker does not have predictive value in the breeding germplasm; and (5) the trait is easier to score phenotypically.

Coupling QTL mapping within breeding program activities would greatly advance the utility of QTL detection. Those marker–QTL associations detected within a breeding program would be immediately useful for MAS approaches. Breeding programs collect a vast amount of phenotypic data on breeding lines on a yearly basis. Most of the phenotypic data are used to select a small number of breeding lines that are advanced to the next stage of evaluation. Typically, none of these data are used to gain a better understanding of the underlying genetics of traits. In the era of cheap genotyping there exists the possibility of leveraging the phenotyping capacity of breeding programs for QTL identification. One approach to leveraging this information is to conduct GWAS for QTL detection using breeding trial data. The barley CAP is a large collaborative project with a goal to conduct GWAS within barley breeding programs (Waugh et al. 2009). To achieve this goal, each year 96 lines were contributed from each of 10 US breeding programs for phenotypic evaluation in various environments and each of the lines was genotyped with 3,072 single nucleotide polymorphism (SNP) markers

(Close et al. 2009; Rostoks et al. 2006). One attractive feature of breeding program data is that large populations of breeding lines can be assembled for mapping. A disadvantage is that in order to assemble these large populations it is necessary to combine phenotypic data from multiple breeding trials. To assess the utility of GWAS within breeding programs, an understanding of the constraints of population size and experimental design are required.

Generally, larger population sizes will lead to higher power and greater precision in QTL detection. In bi-parental QTL mapping analysis, the standard number of lines used has been about 100 (Vales et al. 2005). Small population sizes in bi-parental mapping populations are common and have resulted in underestimates of the number of QTL controlling a trait, inflated QTL effects, and the inability to detect QTL  $\times$  QTL interactions (Vales et al. 2005). In addition, population size affects the estimation of LD and LD estimates are a function of sample size, inter-marker distance and marker heterozygosity (Khatkar et al. 2008; Slate and Pemberton 2007). Therefore, the primary questions for GWAS are how many lines are required to identify QTL and to obtain an accurate estimate of the number of the QTL.

GWAS provide the opportunity to take advantage of large data sets generated by breeding programs. However, breeding programs will often evaluate different sets of breeding lines in different experiments because they differ in stage of testing or breeding objectives. In some cases, the use of checks that are common to all experiments can aid in making comparisons among lines. In addition, several unbalanced experimental designs, e.g. lattice designs and cross-over designs, can generate unbalanced data sets (Spilke et al. 2005). Moreover, unbalanced data sets may also come from balanced designs with one or more observations missing. Linear mixed model analysis has been commonly used to obtain best linear unbiased predictions (BLUPs) or breeding values for the lines tested in an unbalanced design (Smith et al. 2005), which can then be used for subsequent analysis. Thus, an important question is, for GWAS what is the effect of unbalanced designs resulting from combining multiple breeding data sets?

To assess the impact of population size and unbalanced data sets on GWAS, we chose to focus on the trait heading date obtained as part of the barley CAP. Heading date is a highly heritable agronomic trait, is easy to score, and has been the subject of many prior QTL studies; thus it is a useful trait to study the impacts of population size and unbalanced data sets on GWAS. The objectives of this study were to: (1) identify QTL associated with heading date in different environments through GWAS within breeding program germplasm; (2) investigate the effect of population size to identify robust QTL; and (3) examine the utility of unbalanced data sets for GWAS.

## Materials and methods

### Mapping panels and experimental design

Two mapping panels were used to conduct GWAS. All lines in the mapping panels were from spring barley breeding programs and were inbred to at least the  $F_4$  via single seed descent. The first mapping panel (panel #1) contained 766 lines derived from eight spring barley breeding programs: USDA-ARS, Aberdeen, ID (AB), Busch Agricultural Resources Inc. (BA), University of Minnesota (MN), Montana State University (MT), North Dakota State 2-Row (N2), North Dakota State 6-Row (N6), Utah State University (UT), and Washington State University (WA). Each program contributed 96 lines to the panel, except two lines from the AB program that were not included. This mapping panel was planted in 2007 in a randomized complete block design with two replications in a dryland trial in Huntley, MT and an irrigated trial in Bozeman, MT. The checks used in the Montana trials were: Craft, Eslick, Haxby, Hockett, Geraldine, Tradition, Drummond, AC\_Metcalf, Harrington, Baronesse, and Robust. The second mapping panel (panel #2) consisted of 384 lines from four breeding programs in the Upper Midwest: BA, MN, N2, and N6. These lines were a subset of the lines in the larger panel of 766 lines. These 384 lines were grown in balanced and unbalanced experiments. In the balanced experiment, all 384 lines were evaluated in a Fusarium head blight resistance experiment in Fargo, ND and Langdon, ND in 2006 and in Crookston, MN in 2007 (Massman et al. 2011). This experiment utilized a randomized complete block design with two replications in each environment. The checks used in these balanced experiments were Robust, Stander, and MNBrite. These same 384 lines were also evaluated in standard breeding trials conducted by each of the individual breeding programs. In these unbalanced data sets, all 384 lines were evaluated in 9 experiments (see Supplemental Table 1). Each experiment consisted of more than one trial, each of which contained part of the 384 lines. Each trial utilized a randomized complete block design with one, two or three replications. At least three of the following genotypes were used as checks for each trial: Conrad, Merit, Legacy, Tradition, Baronesse, Harrington, Lacey, MNBrite, Robust, Stander, and Stellar.

### Genotype and phenotype data

Heading date was measured as the number of days after planting in which at least 50% of the spikes in a row were emerged at least half way from the boot. The entry mean based broad-sense heritability of heading date was estimated using PROC MIXED in SAS v9.2 (Holland et al. 2003).

Leaf tissue was sampled from a single plant from each of the 766 lines. DNA was extracted using the method of Slotta et al. (2008). SNP genotyping was conducted using two Illumina oligo pool assays (OPAs) containing 1,536 SNPs and referred to as barley OPA1 and barley OPA2 (Close et al. 2009). The Illumina BeadStation was used to genotype each line utilizing the Golden Gate assay, as described in Fan et al. (2006). All marker and phenotypic data used in this study are available at The Hordeum Toolbox (<http://thehordeumtoolbox.org>).

The genotypic data sets were downloaded from The Hordeum Toolbox after setting the maximum missing data criterion at 20%. Using TASSEL v3.0 (Bradbury et al. 2007), the heterozygous SNPs, the monomorphic SNPs and those with minor allele frequencies (MAF) <0.001 were removed from the data sets. After these steps, 2,653 SNPs and 2,538 SNPs including unmapped SNPs met our criteria and were used in panel #1 and panel #2, respectively.

### Population structure in mapping panels

The SNP data were used to determine the population structure of the two mapping panels. Principal component analysis (PCA) was performed in TASSEL v3.0 after imputing missing data using the 3-nearest neighborhood method measured by Manhattan distances. STRUCTURE v2.3.1 was also used to infer the population structure in the two mapping panels (Pritchard et al. 2000). When performing the structure analysis, the number of populations ( $k$ ) was set from 1 to 12 and both burn-in time and MCMC (Markov Chain Monte Carlo) iteration numbers were set to 10,000, and a model for admixture and correlated allele frequencies was used. Five independent replicate runs were performed for this parameter setting. The best value of  $k$  was determined by  $\ln P(d)$  (log posterior probability of data) using the method proposed by Evanno et al. (2005).

### Population size simulation

Based on the 766 lines in panel #1, the mapping population size was reduced by random resampling without replacement from each breeding program. After random resampling of 12, 24, 36, 48, 60, 72, and 84 lines from each breeding program, mapping populations with 96, 192, 288, 384, 480, 576, and 672 lines were formed. For each resampling, five repeats were conducted. The arithmetic mean of two environments (dryland and irrigated) for each line was used. For each new population, PCA was used to determine population structure. GWAS were conducted for each new mapping population generated from each repeat of resampling (see below).

## Balanced and unbalanced data sets

In the 384 lines that composed panel #2, balanced data sets were obtained from three environments. An arithmetic mean of three environments for each line was used for the GWAS. The unbalanced data sets were obtained from the 384 lines in panel #2, which were tested in nine experiments comprising 3–7 environments per experiment connected by 3–9 checks (Supplemental Table 1). For the unbalanced data set, one-step and two-step analysis were used to find marker-trait associations. In the two-step analysis, first BLUPs for each line were calculated using SAS PROC MIXED. BLUPs were calculated using the equation:  $y = X\beta + Zu + e$  and solved using the mixed model equations.  $y$  is a vector of observed phenotypes.  $\beta$  is the effect of experiments, which is the mean of trials and treated as fixed effect (Bernardo 2010).  $u$  is the polygene background effects and treated as random effects. The polygene background effects were assumed to be independently and identically distributed.  $X$  and  $Z$  are incidence matrices for the relationship of  $y$  to  $\beta$  and  $u$ , respectively.  $e$  is an error term and treated as a random effect. Two random variables  $u$  and  $e$  have variance:  $\text{Var}(u) = IV_g$ , and  $\text{Var}(e) = RV_R$ . Where  $I$  is an identity matrix and  $R$  is a diagonal matrix in which the off-diagonal elements are zero and the diagonal elements are the reciprocal of the number of trials in each experiment.  $V_g$  is the genetic variance and  $V_R$  is the residual variance. The restricted maximum likelihood (REML) estimation method was used to estimate  $V_g$  and  $V_R$ . Common checks among different trials were used to adjust for environmental effects. Then the BLUPs calculated in the first step were used as a phenotypic trait in the second step for the GWAS (see methods below). In the one-step analysis, the environments, marker effects, populations and polygenic effects were fitted in one linear mixed model.

A simulation method of subsetting the balanced data sets to obtain unbalanced data sets was also used to compare the GWAS results using balanced and unbalanced data sets. The full balanced data sets from panel #1 evaluated in the dryland and irrigated environments were used for the simulation. To create unbalanced data sets from the balanced data sets, half of the lines from each breeding program were randomly chosen for one environment and the remaining half of the lines were chosen for another environment. Thus, each of the 766 lines has a single heading date value in a single environment. In contrast, in the original balanced data sets each line has two heading date values in both environments. Five repeats were implemented for the simulation. For each new population, PCA was used to determine population structure. GWAS were conducted with both one- and two-step analyses using the simulated unbalanced data sets.

## Association mapping analyses and comparisons

For the GWAS of balanced data sets and two-step analysis of unbalanced data sets, the compressed MLM in TASSEL v3.0 was used (Bradbury et al. 2007; Zhang et al. 2010; Yu et al. 2006). The  $R^2$  for each marker was reported from the TASSEL output. For the one-step analysis of unbalanced data sets, SAS PROC MIXED was used to identify marker-trait associations. In the MLM, the first three principal components were used to construct the  $Q$  matrix for population structure correction. The relative kinship matrix  $K$  was calculated from the percentage of shared alleles between paired lines. The threshold for detecting QTL considering multiple testing was determined by setting the false discovery rate at 0.05 (Benjamini and Hochberg 1995). The marker  $p$  values using balanced and unbalanced data sets (empirical and simulated) were compared by quantile–quantile plot of  $-\log_{10}(p)$ . Observed  $p$  values were from marker-trait association analysis. The expected distribution of marker  $p$  values was uniform, which assumes random markers are unlinked to the polymorphisms controlling heading date.

## Results

### Phenotypic data

In the two mapping panels, the heading date mean from each trial was used for all of the analyses except those using simulated unbalanced data sets. In mapping panel #1 composed of 766 lines, there were two environments (trials) and heading date had a range from 49.5 to 57 days (Table 1). The Pearson's correlation coefficient between two environments in the first mapping panel of 766 lines was 0.75. In mapping panel #2 containing 384 lines, there were three environments (trials) for obtaining the balanced data sets and nine experiments (each with 3–7 environments) for obtaining the unbalanced data sets. As expected, the broad-sense heritabilities for heading date based on entry means were relatively high, ranging from 0.52 to 0.96. The Pearson's correlation coefficients in panel #2 for the balanced data sets were moderate among the three environments and ranged from 0.39 to 0.50 (Table 1). For the unbalanced data sets derived from panel #2, the Pearson's correlation coefficients ranged from 0.16 to 0.93 (Table 1).

### Population structure

The population structure of panel #1 was estimated using STRUCTURE and PCA based on 2,653 SNPs after filtering missing data and removing monomorphic markers. The

**Table 1** Descriptive statistics and entry mean based heritability of heading date in the two mapping panels

Experiment <sup>a</sup>	No. of lines	No. of env. <sup>b</sup>	Range of Pearson's $r^c$	Range <sup>d</sup>	Mean $\pm$ SD	$H^2$ <sup>e</sup>
Mapping panel #1						
Drought experiment	766	2	0.75	49.5–57	51.8 $\pm$ 2	0.89
Mapping panel #2						
Unbalanced data sets						
2401	33	4	0.16–0.70	63.8–66.3	65 $\pm$ 0.7	0.63
6401	32	4	0.18–0.66	59.3–64.3	62 $\pm$ 0.9	0.52
INT-IM	29	2	0.47	64.5–70	68.3 $\pm$ 1.3	0.7
CAPBA	89	1	na <sup>f</sup>	61–74	67.8 $\pm$ 2.8	0.95
PYT1	96	3	0.44–0.63	56.1–60.3	57.8 $\pm$ 0.8	0.73
Expt12	27	5	0.48–0.89	54.7–56.8	54.7 $\pm$ 1.2	0.91
Expt13	69	5	0.71–0.86	53.4–59	55.5 $\pm$ 1.1	0.71
Expt3	74	7	0.81–0.93	52.7–60.3	56.1 $\pm$ 1.7	0.96
Expt5	22	4	0.38–0.79	54.3–58	56.4 $\pm$ 0.9	0.79
Balanced data set						
FHBN	384	3	0.39–0.50	48–57.5	53.1 $\pm$ 3.1	0.69

<sup>a</sup> Descriptive names of experiments: *Drought Experiment*, all CAP I lines from eight spring barley breeding programs, *2401* two-row elite yield trial in the BA breeding program, *6401* six-row elite yield trial in the BA breeding program, *INT-IM* international two-row yield trial in the BA breeding program, *CAPBA* in the BA breeding program, *Expt12* advanced two-row yield trial in the N2 breeding program, *Expt13* intermediate two-row yield trial in the N2 breeding program, *Expt3* intermediate malting barley yield trial in the N6 breeding program, *Expt5* intermediate low protein malting barley yield trial in the N6 breeding program, *PYT1* preliminary yield trial 1 in the MN breeding program, *FHBN* Fusarium head blight resistance trial using lines from the BA, N2, N6, MN breeding programs. MN, University of Minnesota breeding program; N2, North Dakota State University 2-row breeding program; N6, North Dakota 6-row breeding program; BA, Busch Agricultural Resources Inc. breeding program. For more details on the location of each trial and the number of replications at each location see Supplementary Table 1

<sup>b</sup> The number of environments in which the breeding lines were evaluated

<sup>c</sup> Range of Pearson's correlation coefficients among environments within experiment

<sup>d</sup> Range of heading date calculated as the arithmetic means of multiple environments

<sup>e</sup> Entry mean based broad-sense heritability of heading date

<sup>f</sup> Not available

best value of  $k = 6$  was determined by  $\ln P(d)$  (log posterior probability of data) using STRUCTURE as described by Evanno et al. (2005). Thus, the lines from the eight breeding programs were divided into six subpopulations, approximately corresponding to the breeding programs (Fig. 1a). Three of those subpopulations were six-row barley and another three were two-row barley. Previously, Hamblin et al. (2010) showed a similar population structure for these same lines except that the Utah program was not defined as a separate population. PCA of panel #1 gave similar results as the STRUCTURE analysis. The first principal component explained 26.8% of the total variation and separated two-row and six-row subpopulations. The second principal component explained 6.3% of the total variation and further separated breeding programs in two-row and six-row subpopulations. For example, the N2 breeding program was clearly different from other two-row subpopulations, and the UT breeding program was clearly different from the other six-row subpopulations (Fig. 1b).

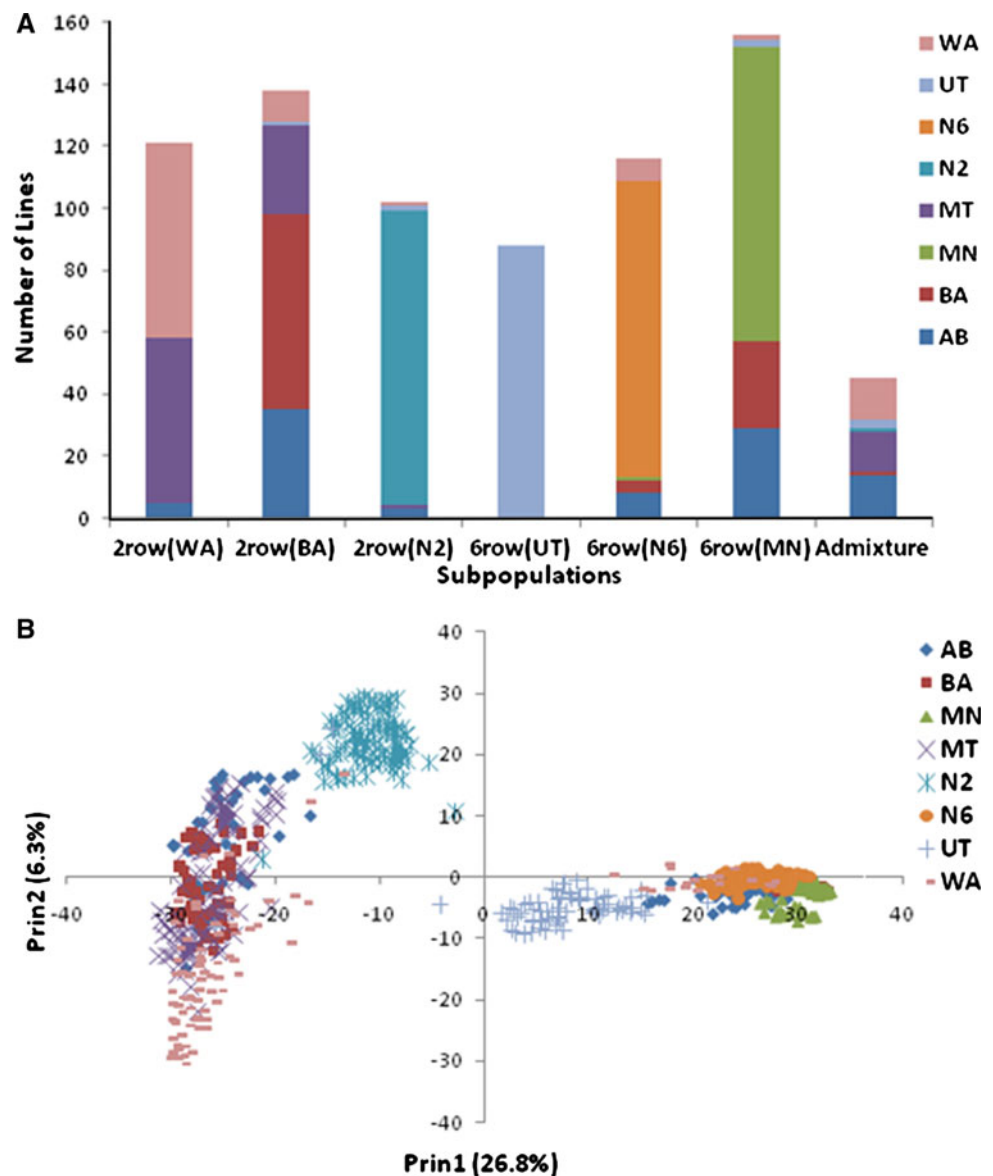
Population structure was also estimated for panel #2 using STRUCTURE and PCA with 2,538 SNPs left after

filtering based on missing data and removing monomorphic markers (Supplemental Fig. 1C). STRUCTURE analysis revealed that there are four subpopulations, which was consistent with population structure analysis using the same mapping panel and 1,536 SNPs (Massman et al. 2011). The four subpopulations corresponded to each breeding population. The BA breeding program had both six-row and two-row barley and 64 two-row barley lines were in a subpopulation with one N2 line, 28 six-row barley lines were in a subpopulation with MN six-row barley, and four of the six-row barley lines were in a subpopulation with N6 six-row barley. In the PCA of panel #2, the first three principal components explained 41.4, 7.5 and 6% of the total variation (Supplemental Fig. 1A, B). The first principal component separated two-row with six-row barley, the second principal component separated two-row barley in the BA and N2 breeding programs, and the third principal component separated six-row barley in N6 and a combination of the MN and BA breeding programs.

In cultivated barley, the difference in inflorescence row type (two or six row) contributes to population structure



**Fig. 1** Population structure of panel #1. **a** Model based population structure analysis (Pritchard et al. 2000) of the 766 member mapping panel from eight breeding programs. Results shown here are inferred after  $k$  is set to six populations. A model for admixture and correlated allele frequencies was assumed. The eight breeding programs are University of Idaho (AB), Busch Agricultural Resources Inc (BA), University of Minnesota (MN), Montana State University (MT), North Dakota State (NDSU 2-Row) (N2), North Dakota State (NDSU 6-Row) (N6), Utah State University (UT), and Washington State University (WA). The six subpopulations are labeled according to the breeding program that contributed the most lines to a subpopulation. **b** Principal component analysis of panel #1 and plot of the first two principal components. The first and second principal components explained 26.8 and 6.3% of total variation, respectively



within breeding germplasm (e.g. Comadran et al. 2011). Thus, Hamblin et al. (2010) advised dividing the population into two-row and six-row subpopulations before conducting association mapping. However, according to our population structure analysis of panel #1 and panel #2, both STRUCTURE and PCA methods separated two-row and six-row types efficiently. Therefore, we combined the lines across inflorescence row types for the GWAS using a mixed linear model in which population structure was included as a fixed effect.

#### GWAS for heading date using panel #1

To account for population structure, we used measures of both fine ( $K$ ) and gross ( $Q$ ) population structure. The fine

level population structure ( $K$ ) was measured using relative kinship. Instead of calculating  $K$  using relative kinship coefficients estimated using SPAGeDi, a matrix of relative kinship coefficients was calculated based on the proportion of shared SNPs. Because SNP markers have a very low mutation rate, shared SNPs imply identity by descent (Zhao et al. 2007). This relative kinship estimate is both simpler to calculate and more accurately reflects the genetic structure of a breeding population where all lines are expected to be related to some extent (Myles et al. 2009). Gross population structure,  $Q$ , was accounted for using PCA as suggested by Price et al. (2006). Using PCA to account for population structure was faster and less computationally intensive than STRUCTURE as has been previously shown (Zhao et al. 2007). We used both the

**Table 2** Association mapping analysis for heading date from breeding lines

QTL <sup>a</sup>	766 lines <sup>b</sup>			384 lines <sup>c</sup>			References <sup>d</sup>
				Balanced data sets	Unbalanced data sets		
	Dryland	Irrigated	Mean	Mean	Two step	One step	
QHd2H.0	0.02 <sup>e</sup>	ND	ND	ND	ND	ND	QHD.umn-2H.1 (EBmac0615, Bmac0093), Mesfin et al. (2003)
QHd2H.64	0.02–0.03	0.01–0.03	0.02–0.03	0.03	0.03–0.04	** <sup>f</sup>	
QHd2H.128	0.02	ND	ND	ND	0.03–0.04	**	
QHd3H.100	ND	ND	ND	ND	0.03	ND	QHD.umn-2H.2 (ABC252), Mesfin et al. (2003)
QHd3H.126-127	0.01–0.02	0.02–0.03	0.02	0.02	0.04	**	
QHd4H.0	ND	ND	ND	ND	0.03	**	
QHd4H.48-50	ND	ND	ND	ND	0.03–0.04	ND	QHD.BIE2-3H.1(ABG377), Rostoks et al. (2005)
QHd4H.63	ND	ND	ND	0.03	ND	ND	
QHd5H.6	ND	ND	ND	0.04	ND	ND	
QHd5H.48	ND	ND	ND	ND	0.03	ND	QHD.IgDa-7H.1 (MWG527), Rostoks et al. (2005)
QHd5H.129	ND	ND	ND	ND	0.03	ND	
QHd6H.44	ND	0.02	ND	ND	ND	ND	
QHd6H.56	ND	ND	ND	0.03	ND	ND	QHD.IgDa-7H.1 (MWG527), Rostoks et al. (2005)
QHd6H.71	ND	ND	ND	ND	ND	**	
QHd6H.127	0.02	ND	ND	ND	ND	ND	
QHd7H.37-41	0.01–0.03	0.02–0.03	0.02–0.03	0.03	0.04	**	QHD.IgDa-7H.1 (MWG527), Rostoks et al. (2005)
QHd7H.50	ND	ND	ND	ND	0.03	ND	

ND not detected

<sup>a</sup> Heading date QTL are named by trait and chromosome position

<sup>b</sup> A drought experiment tested in Huntley (dryland) and Bozeman (irrigated), MT in 2007. Mean is the arithmetic mean of dryland and irrigated environments

<sup>c</sup> Four breeding programs in unbalanced and balanced design. Mean is the arithmetic mean of three environments in the balanced design. Two-step analysis and one-step analysis are the two methods used for detecting QTL in unbalanced data sets

<sup>d</sup> Previously identified heading date QTL and the associated markers that are in coincident locations with the heading date QTL detected in this study

<sup>e</sup>  $R^2$  range for significant markers.  $R^2$  value was output from TASSEL v3.0

<sup>f</sup> \*\* QTL detected, calculated from SAS, no  $R^2$

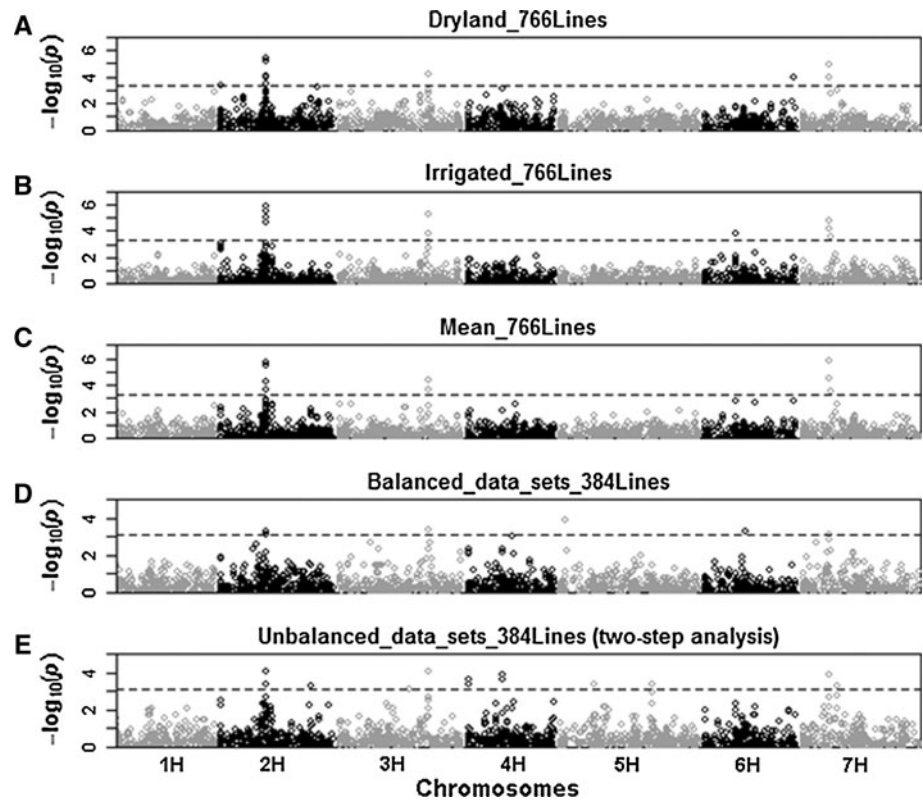
$K$  only model and full model ( $Q + K$ ) for the GWAS for panel #1 and did not find significant differences between the two methods (data not shown). Therefore, all further analyses used the full model ( $Q + K$ ). To further reduce the false-positive rate, the compressed MLM was used instead of the GLM model (Zhang et al. 2010). After correcting for population structure using PCA ( $Q$ ) and using percentage of shared alleles between each pair of lines to account for relative kinship ( $K$ ), the compressed MLM in TASSEL detected six heading date QTL in the dryland environment, four heading date QTL in the irrigated environment, and three heading date QTL using the arithmetic mean of the two environments (Table 2; Fig. 2). Thus, three QTL, QHd2H.64, QHd3H.126-127 and QHd7H.37-41, were identified in both environments (dryland and irrigated) and from the arithmetic mean of the two

environments. The percentage of phenotypic variation explained by significant SNP markers, as assessed by  $R^2$ , had a range from 0.01 to 0.03.

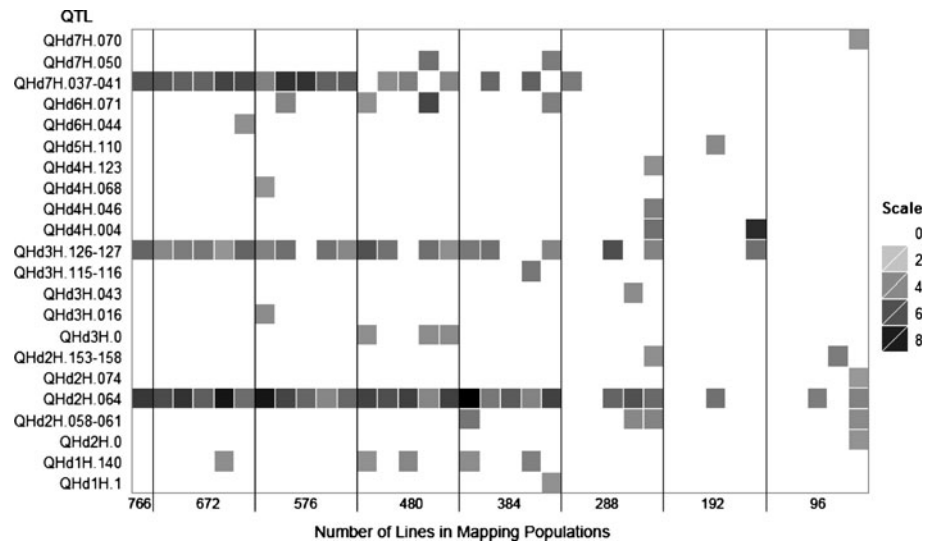
#### Effect of population size

To examine the impact of population size on the ability to detect QTL using GWAS, we randomly resampled the 766 individuals from panel #1 such that we created new mapping populations with 96, 192, 288, 384, 480, 576, and 672 individuals (Fig. 3). PCA of each of the smaller populations showed that the first principal component explained about 25 to 26% of the total variation and the second principal component explained about 5–6% of the total variation, which were similar to those when using 766 lines from panel #1. When the population size was

**Fig. 2** Genome wide scan for significant markers associated with heading data QTL in two mapping panels. **a** 766 breeding lines (panel #1) evaluated in the dryland environment. **b** Panel #1 evaluated in the irrigated environment. **c** The arithmetic mean of panel #1 evaluated in the dryland and irrigated environments. **d** 384 breeding lines (panel #2) evaluated in a balanced design. **e** Panel #2 evaluated in an unbalanced design and two-step analysis applied



**Fig. 3** Effects of population size on identifying common QTL. Heading date QTL are named by trait and chromosome position. Resampling without replacement of 84, 72, 60, 48, 36, 24, 12 lines from each breeding program was conducted. For each resampling, five repeats were done. The squares are  $-\log_{10}(p)$  which uses the smallest marker  $p$  value for each of the QTL



reduced from 766 to 672, 576, 480 and 384, QHd2H.064 was identified in all five repeats of each resampling. For those same four population sizes, QHd3H.127-127 and QHd7H.37-41 were identified in at least two repeats of each resampling. Reducing the population size below 384 further reduced the number of repeats in which those three QTL were detected and we observed at least one repeat of each resampling with no QTL identified. In addition, when population size was reduced, there were new QTL detected in addition to the three QTL identified

using the whole mapping population, indicating an increase of falsely identified QTL. At the smallest population size (96) only QHd2H.064 was detected in two out of the five repeats.

Comparison of results from balanced and unbalanced data sets using panel #2

To assess the utility of balanced and unbalanced data sets for GWAS, we examined the 384 breeding lines in panel

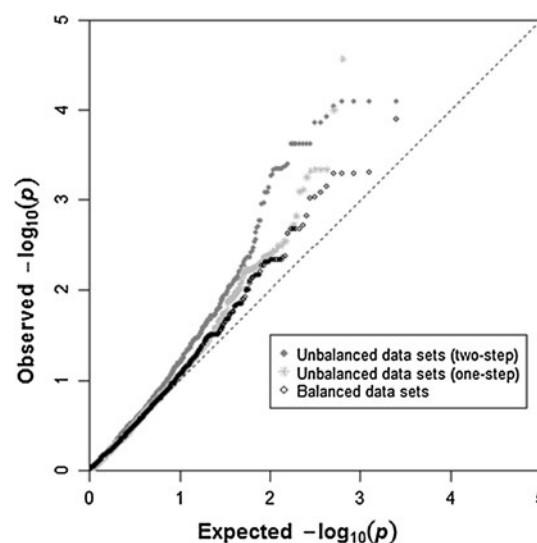


#2. When we analyzed the balanced data sets, we identified six heading date QTL (Fig. 2; Table 2), three of which were the same as the three major QTL identified using the 766 member mapping panel. To analyze the unbalanced data sets, we evaluated the same set of 384 lines and used both one-step and two-step approaches. In the two-step approach, we calculated BLUPs for each line and found that the Pearson's correlation coefficient between the BLUPs and the arithmetic mean of three environments in the balanced design was 0.675. Using these BLUPs, we conducted the second step of GWAS in TASSEL and detected ten heading date QTL (Table 2; Fig. 2). When one-step analysis was conducted in the unbalanced data sets, six heading date QTL were identified using SAS PROC MIXED (Table 2). The three heading date QTL (QHd2H.64, QHd3H.126-127 and QHd7H.37-41) identified in panel #1 were identified in both one- and two-step analyses of the unbalanced data sets. Five QTL (QHd3H.100, QHd4H.48-50, QHd5H.48, QHd5H.129, and QHd7H.50) were identified in the two-step analysis of the unbalanced data sets but not in the one-step analysis or the balanced data sets of panel #1 and #2. Only one locus (QHd6H.71) was identified in the one-step analysis of the unbalanced data sets but not in the two-step analysis or the balanced data sets from panel #1 and #2. The percentage of phenotypic variation explained by the significant SNP markers, as assessed by  $R^2$ , had a range from 0.02 to 0.04.

The quantile–quantile plots indicate that balanced datasets are better than unbalanced datasets for reducing the false-positive rate and that the one-step analysis is better than two-step analysis in reducing the false-positive rate for GWAS using the unbalanced data sets (Fig. 4). The distributions of observed and estimated  $-\log_{10}(p)$  from the GWAS obtained from balanced data sets is slightly less deviated from the distribution of expected  $-\log_{10}(p)$  than that obtained from unbalanced data sets. In the unbalanced data sets analysis, the distribution of observed  $-\log_{10}(p)$  obtained from one-step analysis is slightly less deviated from the distribution of expected  $-\log_{10}(p)$  than that obtained from two-step analysis.

#### Simulation study of unbalanced data sets

To further explore the effect of unbalanced data on GWAS, we conducted a simulation study. Unbalanced data sets were simulated from the full balanced data sets comprising 766 lines (panel #1). In the balanced data sets, the arithmetic mean of the dryland and irrigated environments in panel #1 were used as phenotypic data for the GWAS. In the five repeated simulations of unbalanced data sets, heading date values of lines in either dryland or irrigated environment were used as phenotypic data in the

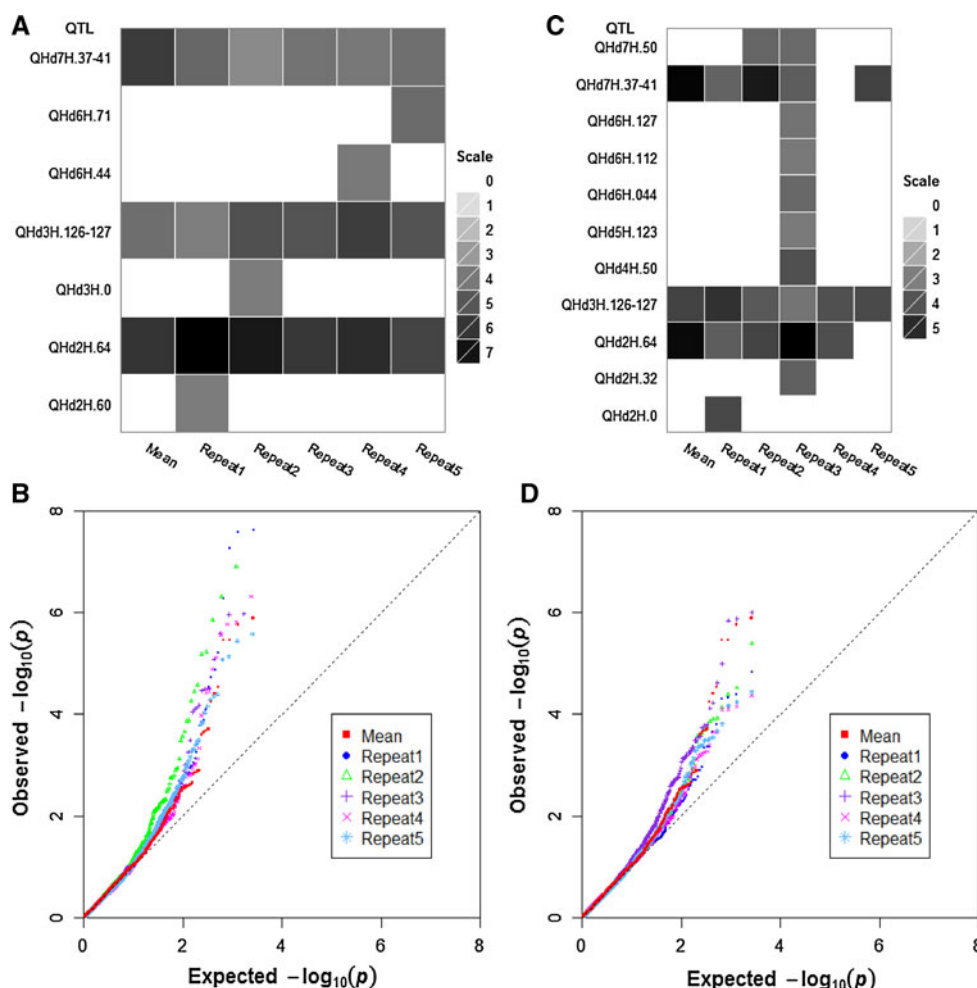


**Fig. 4** Quantile–quantile plot of distribution of estimated and observed  $-\log_{10}(p)$  from the GWAS. Mixed linear model ( $Q + K$ ) method was used for the analysis. Results were from the empirical data analysis using balanced and unbalanced data sets in four breeding programs (384 lines) in the Upper Midwest. One- and two-step analyses methods were conducted for unbalanced data sets

one-step analysis and BLUPs of lines were used as phenotypic data in the two-step analysis. In the one-step analysis, the three heading date QTL, QHd2H.64, QHd3H.126-127 and QHd7H.37-41, detected using the balanced data sets were also detected using the simulated unbalanced data sets (Fig. 5a). In addition to those three QTL, four other QTL were detected. As illustrated in the quantile–quantile plot of distributions of observed and estimated  $-\log_{10}(p)$  from the GWAS, all five repeats of the unbalanced data set simulation were as good as the balanced data sets at reducing the false-positive rate (Fig. 5b).

In the two-step analysis, line BLUPs were calculated in the first step. In all five repeats of the simulation, the Pearson's correlation coefficients between the calculated BLUPs and the arithmetic mean of two environments in the balanced data sets ranged from 0.87 to 0.88. The three major heading date QTL, QHd2H.64, QHd3H.126-127 and QHd7H.37-41, were detected in at least four of the five repeats (Fig. 5c). In addition to those three QTL, eight other QTL were also detected. Seven of these eight were detected in a single repeat. As illustrated in the quantile–quantile plot of distributions of observed and estimated  $-\log_{10}(p)$  from the GWA analysis, four out of five repeats were as good as the balanced data sets at reducing the false-positive rate; the one exception, repeat 3, was slightly more deviated from the diagonal line than the arithmetic mean (Fig. 5c, d).

**Fig. 5** Comparison of balanced data sets and *simulated* unbalanced data sets in the GWAS. One- and two-step analyses were conducted using unbalanced data sets *simulated* from the balanced data sets of 766 lines consisting of 8 breeding programs. The *mean* is the arithmetic mean of two environments (dryland and irrigated) for each line in the balanced data sets of 766 lines. *Repeat1* to *Repeat5* indicates five repeats of simulation of the balanced data sets to unbalanced data sets for the one-step and two-step analysis. **a** Heatmap of heading date QTL detected in *one-step analysis*. The squares are  $-\log_{10}(p)$ , which uses the smallest marker  $p$  value for each QTL. **b** Quantile–quantile plot of distribution of estimated and observed  $-\log_{10}(p)$  in *one-step analysis*. **c** Heatmap of heading date QTL detected in *two-step analysis*. The squares are  $-\log_{10}(p)$ , which uses the smallest marker  $p$  value for each of the QTL. **d** Quantile–quantile plot of distribution of estimated and observed  $-\log_{10}(p)$  in *two-step analysis*



## Discussion

### Heading date QTL detection in breeding program germplasm

In this study, GWAS were conducted using barley lines from multiple breeding programs in the United States. Three heading date QTL, QHd2H.64, QHd3H.126-127 and QHd7H.37-41, were consistently detected in panel #1 containing 766 lines from 8 breeding programs and panel #2 containing 384 lines from four breeding programs, which was a subset of the 766 lines in panel #1. All three of these QTL have been detected in previous mapping studies. QHd2H.64 was detected by SNP markers located on the long arm of chromosome 2H near the centromere. Lorenz et al. (2010) used both single SNP and haplotype analysis to identify heading date QTL in the same region. In addition, an early maturity gene, *Eam6*, was located in this region (Cuesta-Marcos et al. 2008; Franckowiak and Konishi 2002; Horsley et al. 2006). In a recent study, SNP marker 11\_20438 linked to heading date QTL was in the same chromosomal location as *Eam6* (Comadran et al.

2011). The five SNPs associated with QHd2H.64 in our study, 11\_10191, 11\_10685, 11\_20585, 11\_21399 and 12\_30265, were co-located with SNP 11\_20438 on the consensus map (Close et al. 2009). Moreover, all five SNP markers were co-located with 11\_10632, which is tightly linked to microsatellite markers EBmac0615 and Bmac0093 that were found to be associated with heading date (Mesfin et al. 2003; Nduulu et al. 2007). Similarly, one of the significant SNP markers associated with QHd3H.126-127, 12\_30096, is tightly linked to RFLP marker ABG377, which is associated with the heading date QTL QHd.BIE2-3H.1 (Rostoks et al. 2005). For the heading date QTL identified on Chromosome 7, one of the significant SNP markers associated with QHd7H.37-41, 12\_30893, is tightly linked to RFLP marker MWG527, which is associated with QHd.IgDa-7H.1 (Rostoks et al. 2005). Furthermore, the SNP 12\_30893 associated with QHd7H.37-41 is within the *VRN-H3* gene. *VRN-H3* controls the vernalization response and is homologous to the Arabidopsis Flowering Locus T (Yan et al. 2006). In addition to those three QTL, QHd2H.128 was detected in the one- and two-step analysis of panel #2. The significant

markers identifying this QTL are tightly linked to RFLP marker ABC252 that was found to be associated with heading date in the Fredrickson/Stander mapping population (Mesfin et al. 2003).

The significant SNP markers identifying each of the heading date QTL explained 1–4% of the phenotypic variation. This percentage is much lower than that reported in previous heading date QTL mapping studies using biparental populations. For example, heading date QTL explained 2.2–57.8% in several barley bi-parental mapping studies developed from wide crosses (i.e., Mesfin et al. 2003). In addition, it is lower than recent GWAS using rice landraces, in which all identified QTL explained about 36% of the phenotypic variance on average for each trait (Huang et al. 2010). Several possibilities exist to explain the lower percentage of explained phenotypic variation by QTL in our study. One, many major effect heading date QTL are likely fixed in elite populations whereas they would be segregating in a geographically diverse set of landraces or wide biparental crosses (Massman et al. 2011). Two, the range of phenotypic variation is much smaller in elite breeding germplasm. For example, the range of heading date was about 90 days for rice landraces (Huang et al. 2010) compared to 13 days in our GWAS. Three, heading date in elite breeding germplasm may be controlled by epistasis or other non-additive effects; quantitative traits in germplasm that have experienced selection have been shown to be controlled by epistasis (Dudley and Johnson 2009; Ordonez et al. 2010). Four, if there is a strong association between heading date and a specific breeding program, the correction for population structure may reduce the ability to detect heading date QTL (Myles et al. 2009). For example, the QHd2H.64 locus was fixed in the MN and MT breeding programs and was nearly fixed in the N2 and N6 breeding programs, and the QHd7H.37-41 was fixed in MN and N6 breeding programs, indicating that the number of lines segregating for these loci were low and associated with specific breeding programs (Supplemental Table 2). Five, the SNP-based marker coverage used in this study may be insufficient to have a SNP in high LD with each of the heading date QTL segregating in the population (Yu et al. 2006; Lorenz et al. 2010).

#### The effect of population size on QTL detection

In our simulation study of population size, we found that population sizes greater than 384 individuals are required to consistently detect the three major heading date QTL. Small population size may introduce bias when estimating population structure and familial relatedness compared to larger populations. Considering correction for population structure and relatedness is not always effective using a

mixed linear model and therefore could affect the results of GWAS, the correction is expected to be even less accurate for populations of decreasing size. However, in our simulation study, the first and second principal components captured similar percentages of the total variation when comparing the small mapping populations with the full mapping population. This could be explained by the fact that the individual lines were resampled from each breeding program. Therefore, the population structure was similar in the small mapping populations as that in the full mapping population. Furthermore, in a small population, the QTL must have a very large effect and must be in high LD with tested markers in order to be effectively detected (Zhu et al. 2008). Lastly, small population sizes will be inadequate to identify QTL with low minor allele frequencies. As noted above, some heading date QTL were fixed or nearly fixed in several breeding programs, indicating that they may be segregating at a low frequency in the population. All of these factors could explain why in smaller mapping populations the three major QTL were not consistently detected.

In addition to population size, the power to detect QTL is a function of heritability and the number of QTL controlling the trait. Bradbury et al. (2011) used the marker genotypes from a set of breeding lines from the Barley CAP, which included the lines that were used in this study, to investigate the effects of population size, heritability, and number of QTL on QTL detection in GWAS. They showed that population sizes of 300 were sufficient to detect QTL for traits with relatively high heritability (0.75–1.0); however as the number of QTL increased, the larger effect QTL were detected but not the smaller effect QTL. These results appear to agree with our empirical study wherein population sizes of 384 were sufficient to detect several QTL consistently for a trait with moderate to high heritability (Fig. 3).

#### The impact of unbalanced data sets on QTL detection

A comparison of unbalanced data sets and balanced data sets was conducted through empirical data and simulation analysis. From both the empirical data and simulation analyses, both balanced and unbalanced data sets detected three heading date QTL, QHd2H.64, QHd3H.126-127 and QHd7H.37-41. More heading date QTL were detected using the unbalanced data sets compared to the balanced data sets. This could indicate that the use of balanced data sets was better at reducing the number of false-positive associations (Fig. 5a, c; Table 2).

In GWAS using multiple-environmental and unbalanced data sets, an intuitive method is to include both phenotypic and genotypic data in one step, which is called one-step analysis (Cullis et al. 1998; Smith et al. 2001). To reduce

the computational burden in this one-step approach, a two-step analysis can be conducted. To conduct the two-step analysis, the entry means, adjusted means or BLUPs are calculated across multiple environments in the first step and then used for GWAS in the second step (Lorenz et al. 2010; Stich et al. 2008). Both one- and two-step analyses were used in our GWAS using empirical and simulated unbalanced data sets. Our GWAS results indicated that the two-step analysis was slightly less effective at reducing the false-positive rate compared to the one-step approach. We used SAS and TASSEL to conduct the one- and two-step analyses, respectively. Although using different software packages could confound the comparison between one- and two-step approaches because different criteria for convergence and maximum number of iterations were applied, our results were consistent with Stich et al. (2008) in which they found that the empirical type I error rates based on adjusted entry means calculated by a two-step analysis were slightly higher than that of one-step analysis. When using unbalanced data sets, the one-step approach appears to be more effective at reducing false discovery while maintaining power to detect QTL.

#### Integrating breeding and genetic mapping

From a practical standpoint, the ability to map QTL using unbalanced data sets from breeding trials will be of considerable utility. Breeders routinely obtain phenotypic data on a large number of breeding lines within their program. However, much of these data are scattered in separate experiments and must be combined across years and/or locations to generate the large data sets that can reap the benefits of mapping using larger population sizes. Using somewhat casual data (e.g. visual 1–5 scale) from breeding trials has been shown to be useful in bi-parental mapping (e.g., Spaner et al. 1998). Given the exponential decrease in genotyping costs it may make sense for breeders to replace casual measurements with more careful measurements to exploit the potential for GWAS in addition to selection. Strategically assembling mapping panels that are aligned with breeding objectives should provide more useful information for MAS. For example, mapping panels that include germplasm from different breeding programs or regions could be used to design strategies that would take advantage of genetic crosses between programs or regions. Our results demonstrate that breeders can use their data to detect QTL without specialized balanced experimental designs. While this study shows promise using a highly heritable trait, future studies should investigate traits with more modest heritabilities.

Several kinds of valuable information can be gained from conducting GWAS from breeding trial data. First,

generating large population sizes should increase the power of QTL detection and better characterize the genetic architecture of the traits of interest. More simply inherited traits can be managed using traditional marker-assisted selection while highly multigenic traits can be handled more effectively using genomic selection. Second, it should be possible to determine the extent to which previously published QTL are contributing to genetic variation in breeding populations. As previously noted, few QTL discovered in bi-parental mapping populations have been effectively used for MAS. In part, this is due to the fact that favorable alleles discovered in wide cross mapping populations may already be fixed in more elite germplasm (Condon et al. 2008). Alternatively, beneficial alleles tightly linked to deleterious alleles are not selected by breeders. Third, using the same marker platform, different breeders conducting GWAS in their breeding programs will be able to assess the extent to which different genes are affecting their traits of interest. This will aid in the design of crosses of parents between breeding programs and facilitate germplasm exchange. As the cost of genotyping continues to decline, it will be increasingly attractive for breeders to genotype more lines in their program and leverage large amounts of breeding data to conduct GWAS.

To utilize these types of data sets, breeders will need to carefully assess several variables when conducting QTL analysis including the heritability and genetic architecture of the trait, the population size and structure, and experimental design. Among all these variables, population size and experimental design are controlled by the breeder and the response to selection is maximized for various criteria including field space allocation and cost. The number and composition of lines in a mapping population could affect the population structure and also determine whether the traits are correlated with the population structure. If the traits are highly correlated with the population structure, it may be difficult to detect QTL for those traits using GWAS. As genotyping costs decrease, making larger mapping populations more feasible, breeders' attention will shift to experimental designs that can generate the most useful phenotypic data (Myles et al. 2009). Regardless of improvements in both phenotyping and genotyping, it will be highly desirable to combine large breeding data sets to characterize the genetic architecture of important traits in breeding germplasm and then use the appropriate marker-based breeding method for trait improvement.

**Acknowledgments** This research was supported by USDA-CSREES-NRI Grant No. 2006-55606-16722 and USDA-NIFA Grant No. 2009-85606-05701, "Barley Coordinated Agricultural Project: Leveraging Genomics, Genetics, and Breeding for Gene Discovery and Barley Improvement".

## References

- Beló A, Zheng P, Luck S, Shen B, Meyer DJ, Li B, Tingey S, Rafalski A (2008) Whole genome scan detects an allelic variant of *fad2* associated with increased oleic acid levels in maize. *Mol Genet Genomics* 279:1–10
- Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser B-Stat Methodol* 57:289–300
- Bernardo R (2008) Molecular markers and selection for complex traits in plants: learning from the last 20 years. *Crop Sci* 48:1649–1664
- Bernardo R (2010) Breeding for quantitative traits in plants, 2nd edn. Stemma Press Woodburn, MN
- Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES (2007) TASSEL: software for association mapping of complex traits in diverse samples. *Bioinformatics* 23:2633–2635
- Bradbury P, Parker T, Hamblin MT, Jannink JL (2011) Assessment of power and false discovery in genome-wide association studies using the BarleyCAP germplasm. *Crop Sci* 51:52–59
- Breseghello F, Sorrells ME (2006) Association mapping of kernel size and milling quality in wheat (*Triticum aestivum* L.) cultivars. *Genetics* 172:1165–1177
- Close TJ, Bhat PR, Lonardi S, Wu Y, Rostoks N, Ramsay L, Druka A, Stein N, Svensson JT, Wanamaker S, Bozdog S, Roose ML, Moscou MJ, Chao S, Varshney RK, Szucs P, Sato K, Hayes PM, Matthews DE, Kleinhofs A, Muehlbauer GJ, DeYoung J, Marshall DF, Madishetty K, Fenton RD, Condamine P, Graner A, Waugh R (2009) Development and implementation of high-throughput SNP genotyping in barley. *BMC Genomics* 10:582
- Cockram J, White J, Zuluaga DL, Smith D, Comadran J, Macaulay M, Luo Z, Kearsey MJ, Werner P, Harrap D, Tapsell C, Liu H, Hedley PE, Stein N, Schulte D, Steuernagel B, Marshall DF, Thomas WTB, Ramsay L, Mackay I, Balding DJ, Waugh R, O'Sullivan DM, AGOUEB Consortium (2010) Genome-wide association mapping to candidate polymorphism resolution in the unsequenced barley genome. *Proc Natl Acad Sci USA* 107:21611–21616
- Comadran J, Russell JR, Booth A, Pswarayi A, Ceccarelli S, Grando S, Stanca AM, Pecchioni N, Akar T, Al-Yassin A, Benbelkacem A, Ouabbou H, Bort J, van Eeuwijk FA, Thomas WT, Romagosa I (2011) Mixed model association scans of multi-environmental trial data reveal major loci controlling yield and yield related traits in *Hordeum vulgare* in Mediterranean environments. *Theor Appl Genet* 122:1363–1373
- Condon F, Gustus C, Rasmusson DC, Smith KP (2008) Effect of advanced cycle breeding on genetic diversity in barley breeding germplasm. *Crop Sci* 48:1027–1036
- Cuesta-Marcos A, Igartua E, Ciudad FJ, Codesal P, Russell JR, Molinacano JL, Moralejo M, Szucs P, Gracia MP, Lasa JM, Casas AM (2008) Heading date QTL in a spring x winter barley cross evaluated in Mediterranean environments. *Mol Breed* 21:455–471
- Cuesta-Marcos A, Szucs P, Close TJ, Filichkin T, Muehlbauer GJ, Smith KP, Hayes PM (2010) Genome-wide SNPs and resequencing of growth habit and inflorescence genes in barley: implications for association mapping in germplasm arrays varying in size and structure. *BMC Genomics* 11:707
- Cullis B, Gogel B, Verbyla A, Thompson R (1998) Spatial analysis of multi-environment early generation variety trials. *Biometrics* 54:1–18
- Dudley JW, Johnson GR (2009) Epistatic models improve prediction of performance in corn. *Crop Sci* 49:763–770
- Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study. *Mol Ecol* 14:2611–2620
- Fan JB, Chee MS, Gunderson KL (2006) Highly parallel genomic assays. *Nat Rev Genet* 7:632–644
- Franckowiak JD, Konishi T (2002) Early maturity 6, Eam6. *Barley Genet Newsl* 32:86–87
- Gaut BS, Long AD (2003) The lowdown on linkage disequilibrium. *Plant Cell* 15:1502–1506
- Hamblin MT, Close TJ, Bhat PR, Chao S, Kling JG, Abraham KJ, Blake T, Brooks WS, Cooper B, Griffey CA, Hayes PM, Hole DJ, Horsley RD, Obert DE, Smith KP, Ullrich SE, Muehlbauer GJ, Jannink J-L (2010) Population structure and linkage disequilibrium in U.S. barley germplasm: implications for association mapping. *Crop Sci* 50:556–566
- Holland JB, Nyquist WE, Cervantes-Martínez CT (2003) Estimating and interpreting heritability for plant breeding: an update. *Plant Breed Rev* 22:9–112
- Horsley RD, Schmierer D, Maier C, Kudrna D, Urrea CA, Steffenson BJ, Schwarz PB, Franckowiak JD, Green MJ, Zhang B, Kleinhofs A (2006) Identification of QTLs associated with *Fusarium* head blight resistance in barley accession CIho 4196. *Crop Sci* 46:145–156
- Huang X, Wei X, Sang T, Zhao Q, Feng Q, Zhao Y, Li C, Zhu C, Lu T, Zhang Z, Li M, Fan D, Guo Y, Wang A, Wang L, Deng L, Li W, Lu Y, Weng Q, Liu K, Huang T, Zhou T, Jing Y, Li W, Lin Z, Buckler ES, Qian Q, Zhang Q, Li J, Han B (2010) Genome-wide association studies of 14 agronomic traits in rice landraces. *Nat Genet* 42:961–967
- Khatkar MS, Nicholas FW, Collins AR, Zenger KR, Cavanagh JA, Barris W, Schnabel RD, Taylor JF, Raadsma HW (2008) Extent of genome-wide linkage disequilibrium in Australian Holstein-Friesian cattle based on a high-density SNP panel. *BMC Genomics* 9:187
- Kraakman ATW, Niks RE, Van den Berg PMMM, Stam P, Van Eeuwijk FA (2004) Linkage disequilibrium mapping of yield and yield stability in modern spring barley cultivars. *Genetics* 168:435–446
- Kraakman ATW, Martínez F, Mussiraliev B, Van Eeuwijk F, Niks R (2006) Linkage disequilibrium mapping of morphological, resistance, and other agronomically relevant traits in modern spring barley cultivars. *Mol Breed* 17:41–58
- Lorenz AJ, Hamblin MT, Jannink JL, Baxter I (2010) Performance of single nucleotide polymorphisms versus haplotypes for genome-wide association analysis in barley. *PLoS One* 5:e14079
- Massman J, Cooper B, Horsley R, Neate S, Dill-Macky R, Chao S, Dong Y, Schwarz P, Muehlbauer GJ, Smith KP (2011) Genome-wide association mapping of *Fusarium* head blight resistance in contemporary barley breeding germplasm. *Mol Breed* 27:439–454
- Mesfin A, Smith KP, Dill-Macky R, Evans CK, Waugh R, Gustus CD, Muehlbauer GJ (2003) Quantitative trait loci for *Fusarium* head blight resistance in barley detected in a two-rowed by six-rowed population. *Crop Sci* 43:307–318
- Myles S, Peiffer J, Brown PJ, Ersoz ES, Zhang Z, Costich DE, Buckler ES (2009) Association mapping: critical considerations shift from genotyping to experimental design. *Plant Cell* 21:2194–2202
- Nduulu LM, Mesfin A, Muehlbauer GJ, Smith KP (2007) Analysis of the chromosome 2 (2H) region of barley associated with the correlated traits *Fusarium* head blight resistance and heading date. *Theor Appl Genet* 115:561–570
- Ordóñez SA, Silva J, Oard JH (2010) Association mapping of grain quality and flowering in elite japonica rice germplasm. *J Cereal Sci* 51:337–343
- Pressoir G, Brown PJ, Zhu W, Upadaya N, Rocheford T, Buckler ES, Kresovich S (2009) Natural variation in maize architecture is mediated by allelic differences at the PINOID co-ortholog *barren inflorescence2*. *Plant J* 58:618–628



- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D (2006) Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* 38:904–909
- Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus genotype data. *Genetics* 155:945–959
- Ramsay L, Comadran J, Druka A, Marshall DF, Thomas WTB, Macaulay M, MacKenzie K, Simpson C, Fuller J, Bonar N, Hayes PM, Lundqvist U, Franckowiak JD, Close TJ, Muehlbauer GJ, Waugh R (2011) *INTERMEDIUM-C*, a modifier of lateral spikelet fertility in barley, is an ortholog of the maize domestication gene *TEOSINTE BRANCHED 1*. *Nat Genet* 43:169–172
- Rostoks N, Mudie S, Cardle L, Russell J, Ramsay L, Booth A, Svensson JT, Wanamaker SI, Walia H, Rodriguez EM (2005) Genome-wide SNP discovery and linkage analysis in barley based on genes responsive to abiotic stress. *Mol Genet Genomics* 274:515–527
- Rostoks N, Ramsay L, MacKenzie K, Cardle L, Bhat PR, Roose ML, Svensson JT, Stein N, Varshney RK, Marshall DF (2006) Recent history of artificial outcrossing facilitates whole-genome association mapping in elite inbred crop varieties. *Proc Natl Acad Sci USA* 103:18656–18661
- Roy JK, Smith KP, Muehlbauer GJ, Chao S, Close TJ, Steffenson BJ (2010) Association mapping of spot blotch resistance in wild barley. *Mol Breed* 26:243–256
- Slate J, Pemberton JM (2007) Admixture and patterns of linkage disequilibrium in a free-living vertebrate population. *J Evol Biol* 20:1415–1427
- Slotta TA, Brady L, Chao S (2008) High throughput tissue preparation for large-scale genotyping experiments. *Mol Ecol Resour* 8:83–87
- Smith A, Cullis B, Thompson R (2001) Analyzing variety by environment data using multiplicative mixed models and adjustments for spatial field trend. *Biometrics* 57:1138–1147
- Smith AB, Cullis BR, Thompson R (2005) The analysis of crop cultivar breeding and evaluation trials: an overview of current mixed model approaches. *J Agric Sci* 143:449–462
- Spaner D, Shugar LP, Choo TM, Falak I, Briggs KG, Legge WG, Falk DE, Ullrich SE, Tinker NA, Steffenson BJ, Mather DE (1998) Mapping of disease resistance loci in barley on the basis of visual assessment of naturally occurring symptoms. *Crop Sci* 38:843–850
- Spilke J, Piepho HP, Hu X (2005) Analysis of unbalanced data by mixed linear models using the MIXED procedure of the SAS system. *J Agron Crop Sci* 191:47–54
- Stich B, Möhring J, Piepho HP, Heckenberger M, Buckler ES, Melchinger AE (2008) Comparison of mixed-model approaches for association mapping. *Genetics* 178:1745–1754
- Thornsberry JM, Goodman MM, Doebley J, Kresovich S, Nielsen D, Buckler ES (2001) *Dwarf8* polymorphisms associate with variation in flowering time. *Nat Genet* 28:286–289
- Vales MI, Schön CC, Capetini F, Chen XM, Corey AE, Mather DE, Mundt CC, Richardson KL, Sandoval-Islas JS, Utz HF, Hayes PM (2005) Effect of population size on the estimation of QTL: a test using resistance to barley stripe rust. *Theor Appl Genet* 111:1260–1270
- von Zitzewitz J, Condon F, Corey A, Cuesta-Marcos A, Filichkina T, Haggard K, Fisk SP, Smith KP, Muehlbauer GJ, Karsai I, Hayes PM (2011) The genetics of winterhardiness in barley: perspectives from genome-wide association mapping. *The Plant Genome* 4:76–91
- Wang J, McClean PE, Lee R, Goos RJ, Helms T (2008) Association mapping of iron deficiency chlorosis loci in soybean (*Glycine max* L. Merr.) advanced breeding lines. *Theor Appl Genet* 116:777–787
- Waugh R, Jannink JL, Muehlbauer GJ, Ramsay L (2009) The emergence of whole genome association scans in barley. *Curr Opin Plant Biol* 12:218–222
- Wilson LM, Whitt SR, Ibanez AM, Rocheford TR, Goodman MM, ES Buckler IV (2004) Dissection of maize kernel composition and starch production by candidate gene association. *Plant Cell* 16:2719–2733
- Yan L, Fu D, Li C, Blechl A, Tranquilli G, Bonafede M, Sanchez A, Valarik M, Yasuda S, Dubcovsky J (2006) The wheat and barley vernalization gene *VRN3* is an orthologue of *FT*. *Proc Natl Acad Sci USA* 103:19581–19586
- Yu J, Pressoir G, Briggs WH, Bi IV, Yamasaki M, Doebley JF, McMullen MD, Gaut BS, Nielsen DM, Holland JB (2006) A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat Genet* 38:203–208
- Zhang Z, Ersoz E, Lai CQ, Todhunter RJ, Tiwari HK, Gore MA, Bradbury PJ, Yu J, Arnett DK, Ordovas JM, Buckler ES (2010) Mixed linear model approach adapted for genome-wide association studies. *Nat Genet* 42:355–360
- Zhao K, Aranzana MJ, Kim S, Lister C, Shindo C, Tang C, Toomajian C, Zheng H, Dean C, Marjoram P (2007) An *Arabidopsis* example of association mapping in structured samples. *PLoS Genet* 3:e4
- Zhu C, Gore M, Buckler ES, Yu J (2008) Status and prospects of association mapping in plants. *Plant Genome* 1:5–20