

1 Using Bayesian multilevel whole-genome regression models for
2 partial pooling of estimation sets in genomic prediction

3 Frank Technow and L. Radu Totir

4 December 19, 2014

5 DuPont Pioneer, Johnston, Iowa 50131

6

Running Head: Partial pooling of estimation sets

7

Key Words: genomic prediction, data pooling, Bayesian whole genome regression, multilevel

8

models

9

Corresponding Author:

10

Frank Technow

11

DuPont Pioneer

12

8305 NW 62nd Ave

13

Johnston, Iowa 50131

14

`Frank.Technow@pioneer.com`

15

Tel.: +1 (515) 535-8317

16

Fax.: +1 (515) 535-0226

17

Abstract

18 Estimation set size is an important determinant of genomic prediction accuracy. Plant breeding
19 programs are characterized by a high degree of structuring, particularly into populations. This
20 hampers establishment of large estimation sets for each population. Pooling populations increases
21 estimation set size but ignores unique genetic characteristics of each. A possible solution is partial
22 pooling with multilevel models, which allows estimating population specific marker effects
23 while still leveraging information across populations. We developed a Bayesian multilevel whole-
24 genome regression model and compared its performance to that of the popular BayesA model
25 applied to each population separately (no pooling) and to the joined data set (complete pooling).
26 As example we analyzed a wide array of traits from the nested association mapping maize pop-
27 ulation. There we show that for small population sizes (e.g., < 50), partial pooling increased
28 prediction accuracy over no or complete pooling for populations represented in the estimation set.
29 No pooling was superior however when populations were large. In another example data set of
30 interconnected biparental maize populations either partial or complete pooling were superior, de-
31 pending on the trait. A simulation showed that no pooling is superior when differences in genetic
32 effects among populations are large and partial pooling when they are intermediate. With small
33 differences, partial and complete pooling achieved equally high accuracy. For prediction of new
34 populations, partial and complete pooling had very similar accuracy in all cases. We conclude that
35 partial pooling with multilevel models can maximize the potential of pooling by making optimal
36 use of information in pooled estimation sets.

INTRODUCTION

37

38 Genomic selection (Meuwissen *et al.* 2001) in animal and plant breeding rests on the accurate
39 prediction of genomic breeding values (GEBV). An important determinant of prediction accuracy
40 is the size of the estimation set (Daetwyler *et al.* 2010). In animal breeding, assembling large
41 estimation sets is relatively straight forward for large dairy breeds like Holstein Friesian, where
42 genomic selection is applied most successfully to date (Hayes *et al.* 2009). For smaller dairy cattle
43 breeds and in particular for beef cattle breeds, however, assembling sufficiently large estimation
44 sets within each breed is often not possible (Weber *et al.* 2012). Creation of multi-population esti-
45 mation sets by pooling several breeds is therefore of great interest and subject of current research
46 (Lund *et al.* 2014).

47 A similar situation exists in plant breeding, which is characterized by a high degree of struc-
48 turing (Albrecht *et al.* 2014). This structuring results from the importance of keeping distinct
49 heterotic groups for maximum exploitation of heterosis (Melchinger and Gumber 1998), from the
50 predominance of distinct biparental populations (Riedelsheimer *et al.* 2013) and the need for spe-
51 cialized breeding programs targeting specific traits or environments (Windhausen *et al.* 2012). This
52 requires that the phenotyping and genotyping resources available to a breeding program have to
53 be allocated to multiple populations, which prevents the creation of sufficiently large estimation
54 sets for each population. Several studies therefore investigated the merit of pooled estimation sets
55 combining populations (Asoro *et al.* 2011; Heffner *et al.* 2011; Lorenz *et al.* 2012; Riedelsheimer
56 *et al.* 2013; Lehermeier *et al.* 2014) or even heterotic groups (Technow *et al.* 2013; Lehermeier
57 *et al.* 2014).

58 However, pooling estimation sets is complicated by genetic differences among populations,
59 such as in linkage disequilibrium, allele frequencies or relationship structure (Windhausen *et al.*
60 2012; Weber *et al.* 2012; Riedelsheimer *et al.* 2013; Technow *et al.* 2014). This might be the reason
61 why using pooled estimation sets failed to increase prediction accuracy in some applications in
62 plant (Desta and Ortiz 2014) and animal breeding (Lund *et al.* 2014).

63 Therefore, Brøndum *et al.* (2012) proposed to use separate estimation sets for each population

64 but to derive genome position specific priors from estimation results in the other population. In
65 this way, unique genome properties of each population could be accounted for while still using
66 information from other populations. A similar, but perhaps more formal approach is “partial pool-
67 ing”, facilitated by Bayesian multilevel models (Gelman and Hill 2006; Gelman and Pardoe 2006;
68 Gelman 2006a). In multilevel models, parameters (e.g., marker effects) are estimated specific for
69 each population but are “shrunk” towards an overall marker effect. Both the specific and overall
70 marker effects are estimated simultaneously from the data, thereby allowing that the former are
71 still informed by data from the other populations. Partial pooling thus strikes a middle ground be-
72 tween “no pooling” (specific marker effects estimated from data of specific population only) and
73 “complete pooling” (unspecific marker effects estimated from pooled estimation set).

74 Our objectives were to (i) demonstrate the use of Bayesian multilevel whole-genome regression
75 models for genomic prediction and (ii) determine in which scenarios partial pooling might be
76 superior over no or complete pooling of estimation sets. Our investigations were based on two
77 publicly available maize breeding data sets and supported by a simulation study.

78 MATERIALS AND METHODS

79 **Multilevel whole genome regression model** The model fitted to the data was

$$y_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_e^2) \quad (1)$$
$$\mu_{ij} = \beta_0 + \sum_k z_{ijk} u_{jk},$$

80 where y_{ij} was the observed phenotypic value of the i^{th} individual from the j^{th} population and
81 μ_{ij} its linear predictor. The phenotypic data y_{ij} was centered to mean zero and scaled to unit
82 variance. The Normal density function, which was used as likelihood, was denoted as \mathcal{N} with σ_e^2
83 denoting the residual variance component. The common intercept was β_0 . Finally, u_{kj} denoted the
84 additive effect of the k^{th} biallelic single nucleotide polymorphism (SNP) marker in population j .
85 The genotype of individual i from population j at marker k was represented by z_{ijk} , which was the

86 number of reference alleles, centered by twice the reference allele frequency. Which of the alleles
87 was chosen as reference allele depended on the data set and is described below. Effects u_{kj} were
88 only estimated when the corresponding marker was polymorphic in population j . Otherwise it was
89 set to 0 and treated as a constant.

90 [Figure 1 about here.]

91 The hierarchical prior distribution setup will be explained next. A graphical display is shown
92 in Figure 1A. The prior of u_{kj} was

$$u_{jk} \sim \mathcal{N}(u_k, \gamma_k^2), \quad (2)$$

93 where u_k was the overall effect of the k^{th} marker and variance parameter γ_k^2 quantified the devi-
94 ations of the specific effects u_{kj} from u_k . Note that all else equal, the shrinkage toward u_k is the
95 stronger the smaller γ_k^2 .

96 Both parameters were associated with prior distributions themselves and estimated from the
97 data. For u_k this was $u_k \sim \mathcal{N}(0, \sigma_k^2)$. Here, the variance parameter σ_k^2 controls the amount of
98 shrinkage towards 0. It was associated with a scaled inverse Chi-square prior with 4.001 degree of
99 freedom and scale parameter S^2 . The prior for u_k thus corresponded to the well known “BayesA”
100 prior (Meuwissen *et al.* 2001).

101 For the variance parameter γ_k^2 , we specified

$$\gamma_k \sim \mathcal{N}(m, d^2, 0 < a, b = \infty) \quad (3)$$

102 which is a Normal distribution prior on γ_k with mean parameter m and standard deviation d , left
103 truncated at zero. Note that the mean of the truncated distribution $\mathcal{N}(m, d^2, 0 < a, b = \infty)$, which
104 is a function of m , d and the truncation points, can be interpreted as the “typical” deviation of the
105 specific marker effects u_{kj} from u_k . Higher values of this mean indicate larger deviations and vice
106 verse. This parameter might therefore be used to quantify population divergence.

107 An uniform prior $Uni(0.001, 0.5)$ was used for the hyperparameters S^2 , m and d . The prior for
108 the intercept β_0 was a Normal distribution with mean 0 and a very large variance. For the residual

109 variance σ_e^2 we specified a uniform distribution prior over the interval $[0, 1]$ on σ_e , which agrees
110 with recommendations for uninformative priors on variance components (Gelman 2006b).

111 Samples from the posterior distribution were drawn with Gibbs sampling, implemented in the
112 JAGS Gibbs sampling environment (Plummer 2003). The total number of samples was 1000,
113 drawn from a single chain with burn in of 10000 and thinning intervals of 500. These settings
114 ensured convergence and an effective sample size (ESS) of > 100 for all parameters (ESS of u_k
115 and u_{jk} were typically > 500).

116 The ESS was calculated with the R (R Core Team 2013) package CODA (Plummer *et al.* 2006),
117 which was also used to monitor convergence using diagnostic plots.

118 **Conventional whole genome regression model** We used the popular Bayesian whole genome
119 regression method “BayesA” (Meuwissen *et al.* 2001), with the modifications of Yang and Tem-
120 pelman (2012) pertaining to the hyperparameter S^2 (see Figure 1B for a graphical representation).
121 The linear model was

$$y_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_e^2) \quad (4)$$
$$\mu_{ij} = \beta_0 + \sum_k z_{ijk} u_k,$$

122 which is principally the same as in (1), with the difference that the population index j was dropped.
123 For no pooling, the model was applied to each population in turn, for complete pooling to the joint
124 data set. For σ_e^2 we used an improper scaled inverse Chi-square prior with -1 degrees of freedom
125 and scale equal to zero. This is equivalent to a uniform prior on σ_e (Gelman 2006b), as was used
126 for the multilevel model, but exploits conjugacy.

127 The BayesA Gibbs Sampler was implemented as a C routine compatible with the R statistical
128 software environment. Again we drew a total number of 1000 samples from a single chain with
129 burn in of 10000 and thinning of 500.

130 **Estimation, prediction and testing procedure** Let Π denote the set of P populations represented
131 in the estimation set and the set of N_p individuals from a population in Π as Λ_p , where p indexes the
132 population in Π . A graphical representation is presented in Figure 2. Further, let those individuals
133 from a population in Π that are not in Λ_p be denoted as $\bar{\Lambda}_p$ and the set of populations not in Π as
134 $\bar{\Pi}$. Populations in $\bar{\Pi}$ will be referred to as “new” populations. The estimation set thus comprised
135 all individuals belonging to Λ_p , for $p \in \Pi$. The test set used for calculating prediction accuracy,
136 comprised individuals in $\bar{\Lambda}_p$ from populations in Π and all individuals from populations in $\bar{\Pi}$.
137 The phenotypic observations of test individuals were masked in the estimation procedure. The
138 separation of populations into Π and $\bar{\Pi}$ and of individuals within a population into Λ_p and $\bar{\Lambda}_p$ was
139 done at random.

140 Within each population, prediction accuracy was computed as the correlation between GEBVs
141 and observed phenotypic values of individuals in the testing set. The within population prediction
142 accuracies were subsequently averaged for populations in Π and $\bar{\Pi}$. These average within popula-
143 tion prediction accuracies will henceforth be denoted as r_{Π} and $r_{\bar{\Pi}}$. Thus, r_{Π} and $r_{\bar{\Pi}}$ correspond
144 to the prediction accuracy for populations represented and not represented in the estimation set,
145 respectively.

146 When using partial pooling, GEBVs of individuals in $\bar{\Lambda}_p$ were predicted using the posterior
147 means of the marker effects estimated for the corresponding population (*i.e.*, u_{jk}). GEBVs of indi-
148 viduals from populations in $\bar{\Pi}$ were predicted using the posterior means of the overall (unspecific)
149 marker effects u_k .

150 When using complete pooling, GEBVs of all individuals in the test set were predicted from the
151 posterior means of marker effects u_k estimated from the joint data set with model (4).

152 Finally, when using no pooling, GEBVs of individuals in $\bar{\Lambda}_p$ were predicted using the posterior
153 means of the marker effects u_k obtained after applying model (4) to the estimation data from the
154 corresponding set Λ_p . The no pooling approach does not provide a direct way of predicting GEBVs
155 of individuals from populations in $\bar{\Pi}$. Thus, $r_{\bar{\Pi}}$ was not evaluated for the no pooling approach.

156 [Figure 2 about here.]

157 **Application to nested association mapping (NAM) maize populations** The NAM data set was
158 obtained from <http://www.panzea.org>. It comprised 4699 recombinant inbred lines (RILs)
159 from 25 biparental crosses between a genetically diverse set of maize inbred lines and line B73
160 as common parent (McMullen *et al.* 2009). The average population size was 188. The RILs
161 were genotyped with 1106 polymorphic SNP markers covering the whole genome. The non-B73
162 allele was defined as the reference allele. We confirmed that all SNP were biallelic and thereby
163 that the reference allele corresponded to the same nucleotide in all 25 populations. To facilitate
164 computations, we used a thinned set of 285 markers, chosen in such a way that there was one
165 marker per 5 cM interval, on average. A previous study showed that a density of one marker
166 per 10 cM interval is sufficient for genomic prediction in the NAM population (Guo *et al.* 2012).
167 We analyzed the traits days to silking (DS), ear height (EH), ear length (EL), southern leaf blight
168 resistance (SLB), near-infrared starch measurements (NS) and upper leaf angle (ULA), which were
169 phenotyped in multi-environment field trials. The phenotypic records used for fitting the models
170 were averages over the single environment phenotypes. The number of environments were 10, 11,
171 8, 3, 7 and 9 for DS, EH, EL, SLB, NS and ULA, respectively. The traits chosen represent the
172 major trait categories available: yield component (EL), agronomic (EH), disease resistance (SLB),
173 flowering (DS), quality (NS) and morphology (ULA).

174 To investigate the effect of total number of lines N , number of populations P and number of
175 lines per population N_p in the estimation set on prediction accuracy and the relative performance
176 of the pooling approaches, the following combinations of P and N_p were considered: $P = 5$ and
177 $N_p = 50$ and 100, $P = 10$ and $N_p = 25$, 50 and 100, $P = 20$ and $N_p = 12.5$, 25, and 50. For $P =$
178 20 and $N_p = 12.5$, we sampled 19 populations with 12 individuals and one with 22, which results
179 in an average N_p of 12.5. The P and N_p combinations thus gave rise to N of either 250, 500 or
180 1000. For each combination of trait, P and N_p , 50 estimation-testing data sets were generated by
181 repeating the sampling of Π and Λ_p as described above. Throughout, the three pooling approaches
182 were applied to the same data sets. The sampling variation between different data sets thus does
183 not enter the comparisons among pooling approaches.

184 **Application to interconnected biparental (IB) maize populations** This data set was obtained
185 from the supplement of Riedelsheimer *et al.* (2013). It comprised 635 doubled haploid (DH) lines
186 from five biparental populations with average size of 127. The populations were derived from
187 crosses between four European flint inbred lines. For all DH lines 16741 SNP markers polymor-
188 phic across populations were available. We replaced missing marker genotypes with twice the
189 frequency of the reference allele, which was the allele with the lower frequency. When analyzing
190 the data we used a thinned set of 285 markers. Because the data set did not include a map of the
191 markers, the markers were chosen randomly.

192 The DH lines were phenotyped in multi-environment field trials for Giberella ear rot (GER)
193 severity, a fungal disease caused by *Fusarium graminearum*, deoxynivalenol (DON) content (ma-
194 jor mycotoxin produced by the fungus), ear length (EL), kernel rows (KR) and kernels per row
195 (KpR). A more detailed description of this data set can be found in Riedelsheimer *et al.* (2013) and
196 Martin *et al.* (2012).

197 As described above, populations were randomly split into Λ_p and $\bar{\Lambda}_p$. However, because there
198 were only five populations in total, we did not exclude any populations from Π . Set $\bar{\Pi}$ was thus
199 empty and we did not evaluate $r_{\bar{\Pi}}$.

200 The sets Λ_p comprised 25%, 50% and 75% of the lines in each population, which corresponded
201 to an average N_p of 31, 63 and 95, respectively. For each trait and percentage value of estimation
202 individuals, 100 estimation-testing data sets generated, each time resampling the subset of 285
203 markers too.

204 **Application to simulated data set** We conducted a simulation study to specifically investigate
205 the performance of the pooling approaches under increasing levels of differences in QTL effects
206 among populations. The basis for the simulation were the marker genotypes of the lines in the
207 NAM populations. To simulate genetic values, we first randomly chose 20 marker loci as QTL,
208 which were subsequently removed from the set of observed markers. We drew additive overall
209 effects a_q from a standard normal distribution. Then population specific QTL effects a_{jq} were

210 sampled from $\mathcal{N}(a_q, \tau_q^2)$. The variance parameter τ^2 was chosen such that the relative standard
211 deviation (rSD), *i.e.*, τ_q/a_q , was equal to 2, 1, 0.5, 0.25 and 0.0. The greater rSD, the less similar the
212 population specific QTL effects are. True genetic values were obtained by summing QTL effects
213 a_{jq} according the QTL genotypes of each individual. Finally phenotypic values were simulated by
214 adding a normally distributed noise variable to the true genetic values. The variance of the noise
215 variable was chosen such that the heritability across populations was equal to 0.70. The average
216 within family heritability necessarily increased with decreasing rSD, and was 0.53, 0.58, 0.64, 0.68
217 and 0.70 at rSD 2, 1, 0.5, 0.25 and 0.0, respectively.

218 Set II comprised $P = 10$ populations and sets Λ_p had size $N_p = 25$. For each rSD value
219 50 estimation-testing data sets were generated. The QTL positions and effects were randomly
220 generated anew for each data set. Also in this case we used a thinned set of 285 markers. Because
221 the true genetic values were known, r_{Π} and $r_{\bar{\Pi}}$ were computed as the correlation between true
222 genetic values and GEBVs.

223 RESULTS

224 **NAM maize populations** Trends typically held across traits. The results presented and discussed
225 therefore apply to all traits, unless otherwise mentioned.

226 Increasing N_p while keeping N constant (*i.e.*, having fewer but larger populations in the esti-
227 mation set) generally increased r_{Π} and decreased $r_{\bar{\Pi}}$ (Table 1). However, the increase in r_{Π} was
228 much more pronounced than the decrease in $r_{\bar{\Pi}}$.

229 When increasing N_p with constant P or when increasing P with constant N_p , both r_{Π} and $r_{\bar{\Pi}}$
230 increased (Table 1). However, while in the first case, r_{Π} and $r_{\bar{\Pi}}$ increased in similar magnitudes,
231 the increase in r_{Π} was much smaller than the increase in $r_{\bar{\Pi}}$ in the second case, in particular when
232 N_p was high. Per definition, the accuracy of no pooling is not expected to change as long as N_p
233 remains constant.

234 For low P and high N_p , *e.g.*, $P = 5$ and $N_p = 100$, no pooling achieved the highest r_{Π} and
235 complete pooling the lowest (Table 1). For high P and low N_p , *e.g.*, $P = 20$ and $N_p = 25$, partial

236 pooling achieved the highest r_{II} . Here no pooling resulted in the lowest r_{II} . The only exception to
 237 this was trait DS, where no pooling had a r_{II} equal or higher to partial and complete pooling also
 238 for low N_p .

239 Partial and complete pooling achieved virtually identical prediction accuracies $r_{\overline{II}}$ for new pop-
 240 ulations (Table 1). In general, $r_{\overline{II}}$ of a particular pooling approach was considerably lower than the
 241 corresponding r_{II} . The differences between r_{II} and $r_{\overline{II}}$ tended to be larger for high N_p .

TABLE 1: Average within population prediction accuracies in NAM maize populations

P	N_p	trait	r_{II}			$r_{\overline{II}}$	
			no	partial	complete	partial	complete
5	50	DS	0.41	0.34	0.26	0.19	0.19
		EH	0.47	0.44	0.39	0.31	0.32
		EL	0.39	0.37	0.28	0.19	0.19
		NS	0.39	0.37	0.32	0.25	0.26
		SLB	0.49	0.49	0.45	0.37	0.37
		ULA	0.50	0.48	0.44	0.36	0.36
	100	DS	0.52	0.41	0.28	0.21	0.20
		EH	0.57	0.51	0.43	0.34	0.34
		EL	0.49	0.46	0.35	0.23	0.23
		NS	0.47	0.44	0.36	0.29	0.29
		SLB	0.58	0.58	0.50	0.41	0.41
		ULA	0.58	0.54	0.47	0.40	0.40
10	25	DS	0.32	0.28	0.22	0.18	0.17
		EH	0.38	0.38	0.35	0.30	0.31
		EL	0.31	0.31	0.25	0.21	0.21
		NS	0.30	0.33	0.30	0.26	0.27

Continued on next page

		SLB	0.40	0.46	0.43	0.38	0.39
		ULA	0.39	0.44	0.41	0.36	0.37
50		DS	0.42	0.35	0.26	0.22	0.22
		EH	0.47	0.45	0.40	0.36	0.36
		EL	0.40	0.39	0.29	0.23	0.23
		NS	0.38	0.40	0.35	0.30	0.30
		SLB	0.49	0.52	0.46	0.42	0.43
		ULA	0.48	0.50	0.45	0.41	0.41
100		DS	0.51	0.42	0.30	0.25	0.25
		EH	0.57	0.53	0.44	0.39	0.39
		EL	0.48	0.46	0.33	0.27	0.27
		NS	0.48	0.46	0.38	0.33	0.33
		SLB	0.57	0.57	0.49	0.45	0.45
		ULA	0.59	0.56	0.48	0.45	0.44
20	12.5	DS	0.23	0.23	0.21	0.17	0.17
		EH	0.28	0.34	0.33	0.30	0.31
		EL	0.22	0.27	0.23	0.19	0.19
		NS	0.21	0.30	0.29	0.27	0.28
		SLB	0.31	0.43	0.42	0.38	0.39
		ULA	0.28	0.40	0.39	0.35	0.36
25		DS	0.32	0.30	0.24	0.22	0.23
		EH	0.38	0.42	0.39	0.36	0.37
		EL	0.31	0.34	0.28	0.22	0.22
		NS	0.30	0.36	0.33	0.30	0.31
		SLB	0.39	0.48	0.45	0.42	0.43

	ULA	0.38	0.46	0.44	0.42	0.42
50	DS	0.42	0.37	0.29	0.26	0.26
	EH	0.48	0.49	0.42	0.40	0.40
	EL	0.39	0.40	0.30	0.28	0.29
	NS	0.38	0.41	0.36	0.34	0.34
	SLB	0.49	0.54	0.48	0.46	0.47
	ULA	0.49	0.52	0.47	0.46	0.46

Values shown are average within population prediction accuracies for test individuals, averaged over 50 random estimation-test data splits. The standard errors were < 0.013 . P gives the size of set Π , i.e., the number of populations represented in the estimation set, column N_p gives the number of individuals from each population in Π that were used for estimation, i.e., the sizes of sets Λ_p . The traits were: days to silking (DS), ear height (EH), ear length (EL), southern leaf blight resistance (SLB), near-infrared starch measurements (NS) and upper leaf angle (ULA).

242 **IB maize populations** The prediction accuracy r_{Π} increased with increasing N_p , for all traits and
 243 pooling approaches (Table 2). Averaged over traits, the increase was largest for no pooling, where
 244 the accuracy increased from an average of 0.35 at $N_p = 31$ to 0.48 at $N_p = 95$. The accuracies
 245 for the partial and complete pooling approaches increased from 0.39 and 0.38, respectively, at
 246 $N_p = 31$ to 0.48 at $N_p = 95$.

247 At $N_p = 31$, partial pooling had the highest r_{Π} for traits EL, KpR, complete pooling for traits
 248 DON and KR. For GER both had the same accuracy. The no pooling approach had the lowest r_{Π} ,
 249 except for EL and KpR, where it had the same accuracy as complete pooling. For the highest N_p
 250 of 95, the accuracy differences among the pooling approaches decreased. Partial pooling still had
 251 the highest accuracy for EL and KpR and the same as complete pooling for DON and GER. While
 252 never better than partial pooling, no pooling had higher prediction accuracy than complete pooling
 253 for EL and KpR.

254

[Table 1 about here.]

255 **Simulated maize populations** For all pooling approaches, r_{II} increased with decreasing rSD
256 (Table 3). The increase for no pooling, however, was comparatively small and a result of the
257 increasing within family heritability with decreasing rSD. The relative performance of the pooling
258 approaches also depended on rSD. For the highest rSD value considered, no pooling had the highest
259 r_{II} , for the intermediate rSD value of 1.0 partial pooling. For the lower rSD values complete and
260 partial pooling achieved similarly high r_{II} .

261 Also $r_{\bar{II}}$ for both partial and complete pooling increased strongly with decreasing rSD and the
262 differences to r_{II} decreased (Table 3). Partial and complete pooling achieved almost identical $r_{\bar{II}}$.

263 The mean of the truncated Normal distribution prior $\mathcal{N}(m, d^2, 0 < a, b = \infty)$ for parameter γ_k
264 increased with increasing rSD. Its average values were 0.0111, 0.0153, 0.0190, 0.0269 and 0.0296
265 for rSD of 0.0, 0.25, 0.5, 1.0 and 2.0, respectively.

266 [Table 2 about here.]

267 DISCUSSION

268 **Comparison of pooling approaches** Partial pooling allows estimation of population specific
269 marker effects while still facilitating “borrowing” of information across populations. It is therefore
270 a compromise between no pooling, which models unique characteristics of each population but
271 ignores shared information, and complete pooling, in which the opposite is the case.

272 When population sizes N_p are sufficiently large, borrowing information from other populations
273 is not required for achieving high prediction accuracy of new individuals from the same population
274 (r_{II}). Further enlarging estimation sets by pooling with other populations might then even be detri-
275 mental (Riedelsheimer *et al.* 2013). This explains why no pooling was the most accurate approach
276 when N_p was large (e.g., ≥ 50), particularly in the NAM population, and why it profited most
277 from increases in N_p . Therefore, pooling of estimation sets is most promising if N_p is small due to
278 budget of other constraints. We indeed observed that pooling was more accurate than no pooling
279 when N_p was small (e.g., < 50). The superiority of either pooling approach over no pooling also
280 increased with increasing P , because information from more populations was available, which is

281 not used in no pooling. Thus, pooling is expected to most advantageous when P is relatively high
282 and N_p low. Whether partial or complete pooling is the better approach will then also depend on
283 the similarity of the pooled populations. The greater the similarity, the relatively better complete
284 pooling is expected to perform, because the ability to estimate population specific marker effects
285 becomes less important. In this situation partial pooling might even be of disadvantage, because it
286 requires estimation of many more effects which might lead to problems associated with noniden-
287 tifiability (Gelfand and Sahu 1999). The parents of the IB populations are from the same breeding
288 program (Riedelsheimer *et al.* 2013), whereas the non-common parents of the NAM populations
289 were chosen to be maximally diverse and comprise temperate, tropical and specialty (sweet and
290 popcorn) maize germplasm (McMullen *et al.* 2009). Accommodating for unique characteristics of
291 the populations is therefore more important in NAM than in IB, which might explain why complete
292 pooling was always inferior to partial pooling in the former but often equal or even superior in the
293 latter and also why no pooling never achieved the highest prediction accuracy in IB, even for large
294 N_p .

295 The relative performance of the pooling approaches was very stable across traits in the NAM
296 data set, with the exception of DS. For this trait the no pooling approach was generally superior,
297 even at high P and low N_p . Buckler *et al.* (2009) found evidence for an allelic series at the QTL
298 identified for DS in the NAM population. Thus, while the positions of the QTL are conserved
299 across populations, their effects differ. Possible reasons are presence of multiple alleles or QTL
300 by genetic background interaction. In this situation, pooling of data is not expected to have an
301 advantage over no pooling. This example also shows that decisions about whether to pool data
302 or not have to be made on a by trait basis and should incorporate prior knowledge about genetic
303 architecture, if available.

304 The dependence of the relative performance of the pooling approaches on the similarity of
305 populations was also reinforced by the results from our simulation study. There we also observed
306 that the mean of $\mathcal{N}(m, d^2, 0 < a, b = \infty)$, the prior distribution of γ_k^2 , which quantifies the devia-
307 tions of specific marker effects u_{jk} from the overall effect u_k , increased with increasing simulated

308 differences among population specific QTL effects. This was expected, but demonstrates that the
309 data was informative for the highlevel hyperparameters. Averaged over P and N_p , this mean was
310 largest for DS and ULA in NAM (results not shown). This might reflect the noted differences
311 between population specific QTL effects for DS. Trait ULA, however, did not diverge from the
312 pattern observed for the remainder of traits and there does not seem to be any strong indication of
313 an allelic series as in DS (reference tba). There was also no obvious relation between the mean of
314 $\mathcal{N}(m, d^2, 0 < a, b = \infty)$ and performance of the pooling approaches in IB (results not shown).

315 Modeling unique characteristics of populations requires that these populations are represented
316 in the estimation set. Prediction of individuals from new populations in $\bar{\Pi}$ therefore has to rely on
317 the overall, unspecific marker effects u_k , in both partial and complete pooling. It was thus expected
318 that both achieved very similar prediction accuracies $r_{\bar{\Pi}}$ for new populations.

319 Our results demonstrate that partial pooling is able to model unique characteristics of popu-
320 lations within the estimation set without compromising on the ability of prediction of individuals
321 from new populations. This is one reason why Gelman (2006a) see the the greatest potential of
322 partial pooling with multilevel models in predictive applications.

323 We exemplified the use of multilevel models for partial pooling in the context of multiple
324 populations, a scenario of high relevance for plant (Lehermeier *et al.* 2014) and animal (Lund
325 *et al.* 2014) breeding. However, the concept is readily applicable in a wide array of scenarios.
326 Examples are pooling data across multiple top-cross testers or environments, as is of particular
327 relevance in plant breeding (Albrecht *et al.* 2014). Extending the models to more than two levels is
328 straightforward, too, for example for pooling multiple populations from multiple heterotic groups
329 or breeding programs.

330 **Alternative approaches to partial pooling** There are alternatives to multilevel models for partial
331 pooling. Brøndum *et al.* (2012) leveraged information across populations by using results obtained
332 from one population to derive genome position specific priors for the analysis of another. For ex-
333 ample, when there were two populations A and B, then A was analyzed first and result so obtained

334 used as prior information when analyzing B. One disadvantage of their approach is that because
335 analyses are done sequentially, information is not shared simultaneously among populations. In
336 the example above, information from A is used for B but not vice versa. To use information from
337 B for A, the analyses had to be repeated in reverse order. It is also not obvious how the approach
338 of Brøndum *et al.* (2012) can be generalized to more than two populations or to prediction of
339 individuals from new populations. Another potential source of concern is that the priors derived
340 from population A are too informative to allow substantial Bayesian learning, especially when
341 population B is small (Gelfand and Sahu 1999; Gianola 2013).

342 Lund *et al.* (2014) proposed to consider phenotypic observations from different populations
343 as different traits and to analyze pooled data sets with multi-trait models. This would facilitate
344 simultaneous sharing of information across populations through covariances. When the number of
345 populations becomes large this might prove challenging, however, because of the need of estimat-
346 ing large unstructured covariance matrices. The problem is exacerbated when unique covariance
347 matrices are estimated for each marker, as would be necessary to accommodate for varying link-
348 age phases between markers and QTL among populations (Lund *et al.* 2014). In this case too,
349 prediction of individuals from new populations would not be possible directly.

350 Schulz-Streeck *et al.* (2012) proposed a model that simultaneously fits main and population
351 specific marker effects (u_{snp} and u_{psnp} in their notation). The principal difference to our approach
352 is that both effects are on the same hierarchical level, such that the genetic value of an individual
353 is modeled as the sum of u_{snp} and u_{psnp} . As a consequence, both sets of marker terms “compete”
354 for the same underlying information. This might compromise the ability of prediction in new
355 populations which has to be based on u_{snp} . Prediction targeting individuals from new populations
356 was not attempted by the authors, however.

357 **Composition of estimation set** Increasing the number of individuals from a population in the
358 estimation set (N_p) always increased prediction accuracy for untested individuals from the same
359 population (r_{Π}), regardless if the estimation set was further enlarged by individuals from other

360 populations (partial and complete pooling) or not (no pooling).

361 However, because plant breeding programs have to operate under budget constraints, optimum
362 allocation of resources is of great importance for maximizing the potential of genomic selection
363 (Lorenz 2013; Riedelsheimer and Melchinger 2013). With a fixed budget for phenotyping that is
364 proportional to N , the number of populations P and the number of individuals per population
365 N_p have to be optimized under the constraint that $N = P \cdot N_p$. Such an optimization could be
366 accomplished using basic theory about response to selection (Falconer and Mackay 1996) and
367 accounting for the different prediction accuracy for populations represented and not represented
368 in the estimation set (r_{II} and $r_{\bar{II}}$, respectively), as exemplified by Technow *et al.* (2013). A key
369 point hereby is that r_{II} will increase with increasing N_p but it will apply to fewer populations
370 because of the decrease in P . This is exacerbated by the decrease in $r_{\bar{II}}$ that we observed was
371 associated with decreasing P . Thus, if the total number of populations is large, as is typically the
372 case in plant breeding programs, having very low P is likely to be undesirable. In the context of
373 plant breeding this and other studies, most recently Lehermeier *et al.* (2014), showed that pooling
374 data across populations can at least partly compensate for low N_p if populations are related and
375 there is evidence for the merit of pooling very divergent germplasm too (Technow *et al.* 2013).
376 Using pooled estimation sets therefore has the potential to allow for high P without compromising
377 too much on $r_{\bar{II}}$. We showed that partial pooling with multilevel models can further enhance this
378 potential by making optimal use of the information in pooled estimation sets.

379

LITERATURE CITED

- 380 Albrecht, T., H.-J. Auinger, V. Wimmer, J. Ogutu, C. Knaak, et al., 2014 Genome-based prediction
381 of maize hybrid performance across genetic groups, testers, locations, and years. *Theor Appl*
382 *Genet* 127: 1375–1386.
- 383 Asoro, F. G., M. Newell, M. Beavis, M. Scott, and J.-L. Jannink, 2011 Accuracy and training
384 population design for genomic selection on quantitative traits in elite North American oats.
385 *Plant Gen.* 4: 132–144.

- 386 Brøndum, R. F., G. Su, M. S. Lund, P. J. Bowman, M. E. Goddard, et al., 2012 Genome position
387 specific priors for genomic prediction. *BMC genomics* 13: 543.
- 388 Buckler, E. S., J. B. Holland, P. J. Bradbury, C. B. Acharya, P. J. Brown, et al., 2009 The genetic
389 architecture of maize flowering time. *Science* 325: 714–718.
- 390 Daetwyler, H. D., R. Pong-Wong, B. Villanueva, and J. A. Woolliams, 2010 The impact of genetic
391 architecture on genome-wide evaluation methods. *Genetics* 185: 1021–1031.
- 392 Desta, Z. A. and R. Ortiz, 2014 Genomic selection: genome-wide prediction in plant improvement.
393 *Trends Plant Sci.* 19: 592 – 601.
- 394 Falconer, D. S. and T. F. C. Mackay, 1996 *Introduction to Quantitative Genetics* (4 ed.), Chapter
395 Response to selection, pp. 184–193. Addison Wesley Longman Limited, Harlow.
- 396 Gelfand, A. E. and S. K. Sahu, 1999 Identifiability, improper priors and gibbs sampling for gener-
397 alized linear models. *J Am Stat Assoc* 94: 247–253.
- 398 Gelman, A., 2006a Multilevel (hierarchical) modeling: what it can and cannot do. *Technomet-*
399 *rics* 48: 432 – 435.
- 400 Gelman, A., 2006b Prior distributions for variance parameters in hierarchical models. *Bayesian*
401 *Analysis* 1: 515 – 533.
- 402 Gelman, A. and J. Hill, 2006 *Data analysis using regression and multilevel/hierarchical models*.
403 Cambridge University Press.
- 404 Gelman, A. and I. Pardoe, 2006 Bayesian measures of explained variance and pooling in multilevel
405 (hierarchical) models. *Technometrics* 48: 241 – 251.
- 406 Gianola, D., 2013 Priors in Whole-Genome Regression: The Bayesian Alphabet Returns. *Genet-*
407 *ics* 194: 573–596.
- 408 Guo, Z., D. Tucker, J. Lu, V. Kishore, and G. Gay, 2012 Evaluation of genome-wide selection
409 efficiency in maize nested association mapping populations. *Theor. Appl. Genet.* 124: 261–
410 275.

- 411 Hayes, B. J., P. J. Bowman, A. J. Chamberlain, and M. E. Goddard, 2009 Invited review: genomic
412 selection in dairy cattle: progress and challenges. *J. Dairy Sci.* 92: 433–443.
- 413 Heffner, E. L., J.-L. Jannink, and M. E. Sorrells, 2011 Genomic selection accuracy using multi-
414 family prediction models in a wheat breeding program. *Plant Gen.* 4: 65–75.
- 415 Lehermeier, C., N. Krmer, E. Bauer, C. Bauland, C. Camisan, et al., 2014 Usefulness of Mul-
416 tiparental Populations of Maize (*Zea mays* L.) for Genome-Based Prediction. *Genetics* 198:
417 3–16.
- 418 Lorenz, A. J., 2013 Resource Allocation for Maximizing Prediction Accuracy and Genetic Gain
419 of Genomic Selection in Plant Breeding: A Simulation Experiment. *G3* 3: 481–491.
- 420 Lorenz, A. J., K. P. Smith, and J.-L. Jannink, 2012 Potential and optimization of genomic selection
421 for fusarium head blight resistance in six-row barley. *Crop Sci.* 52: 1609–1621.
- 422 Lund, M. S., G. Su, L. Janss, B. Guldbbrandtsen, and R. F. Brøndum, 2014 Invited review: Genomic
423 evaluation of cattle in a multi-breed context. *Livest. Sci.* 166: 101–110.
- 424 Martin, M., B. S. Dhillon, T. Miedaner, and A. E. Melchinger, 2012 Inheritance of resis-
425 tance to *Gibberella* ear rot and deoxynivalenol contamination in five flint maize crosses. *Plant*
426 *Breed.* 131: 28–32.
- 427 McMullen, M. D., S. Kresovich, H. S. Villeda, P. Bradbury, H. Li, et al., 2009 Genetic Properties
428 of the Maize Nested Association Mapping Population. *Science* 325: 737–740.
- 429 Melchinger, A. E. and R. K. Gumber, 1998 Overview of heterosis and heterotic groups in agro-
430 nomic crops, pp. 29–44 in *Concepts and Breeding of Heterosis in Crop Plants*, edited by K. R.
431 Lamkey and J. E. Staub. CSSA, Madison, WI.
- 432 Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard, 2001 Prediction of Total Genetic Value
433 Using Genome-Wide Dense Marker Maps. *Genetics* 157: 1819–1829.
- 434 Plummer, M., 2003 JAGS: a program for analysis of Bayesian graphical models using Gibbs
435 sampling.

- 436 Plummer, M., N. Best, K. Cowles, and K. Vines, 2006 CODA: convergence diagnosis and output
437 analysis for MCMC. *R News* 6: 7–11.
- 438 R Core Team, 2013 *R: A Language and Environment for Statistical Computing*.
- 439 Riedelsheimer, C., J. B. Endelman, M. Stange, M. E. Sorrells, J.-L. Jannink, et al., 2013 Genomic
440 Predictability of Interconnected Biparental Maize Populations. *Genetics* 194: 493–503.
- 441 Riedelsheimer, C. and A. E. Melchinger, 2013 Optimizing the allocation of resources for genomic
442 selection in one breeding cycle. *Theor Appl Genet* 126: 2835–2848.
- 443 Schulz-Streeck, T., J. Ogutu, Z. Karaman, C. Knaak, and H. Piepho, 2012 Genomic selection
444 using multiple populations. *Crop Sci* 52: 2453–2461.
- 445 Technow, F., A. Bürger, and A. E. Melchinger, 2013 Genomic prediction of northern corn leaf
446 blight resistance in maize with combined or separated training sets for heterotic groups. *G3* 3:
447 197–203.
- 448 Technow, F., T. A. Schrag, W. Schipprack, E. Bauer, H. Simianer, et al., 2014 Genome properties
449 and prospects of genomic prediction of hybrid performance in a breeding program of maize.
450 *Genetics* 197: 1343–1355.
- 451 Weber, K. L., R. M. Thallman, J. W. Keele, W. M. Snelling, G. L. Bennett, et al., 2012 Accuracy
452 of genomic breeding values in multibreed beef cattle populations derived from deregressed
453 breeding values and phenotypes. *J. Anim. Sci.* 90: 4177–4190.
- 454 Windhausen, V. S., G. N. Atlin, J. M. Hickey, J. Crossa, J.-L. Jannink, et al., 2012 Effectiveness
455 of Genomic Prediction of Maize Hybrid Performance in Different Breeding Populations and
456 Environments. *G3* 2: 1427–1436.
- 457 Yang, W. and R. J. Tempelman, 2012 A Bayesian Antedependence Model for Whole Genome
458 Prediction. *Genetics* 190: 1491–1501.

List of Figures

459

460

461

462

463

464

465

466

- 1 Graphical visualization of the multilevel model (A) and the conventional BayesA model (B). 24
- 2 Graphical visualization of the testing strategy for evaluating prediction accuracy. The estimation set comprises Λ_1 and Λ_2 from populations P_1 and P_2 (set Π). The prediction accuracy of lines from populations represented in estimation set (r_{Π}) was computed from $\bar{\Lambda}_1$ and $\bar{\Lambda}_2$, the prediction accuracy of lines from populations not represented in estimation set from lines in P_3 and P_4 (set $\bar{\Pi}$). 25

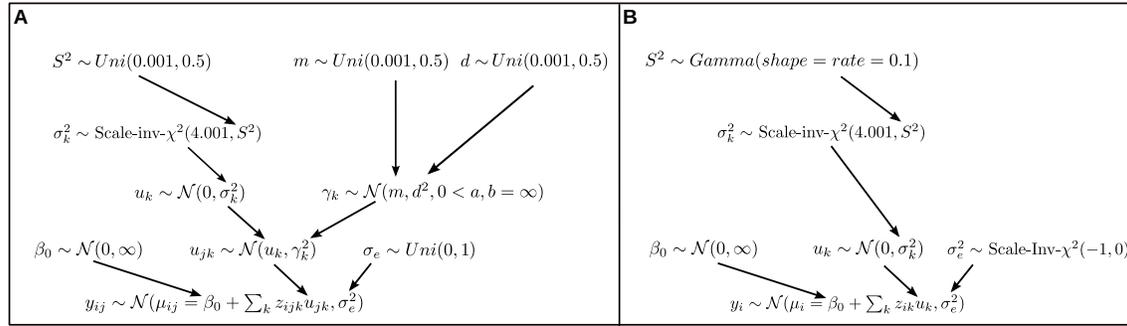


FIGURE 1: Graphical visualization of the multilevel model (A) and the conventional BayesA model (B).

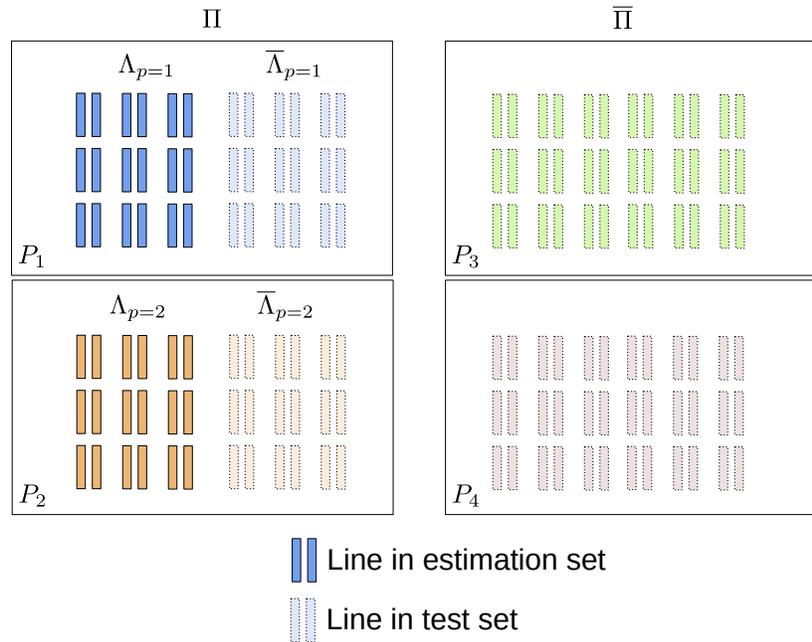


FIGURE 2: Graphical visualization of the testing strategy for evaluating prediction accuracy. The estimation set comprises Λ_1 and Λ_2 from populations P_1 and P_2 (set II). The prediction accuracy of lines from populations represented in estimation set (r_{II}) was computed from $\bar{\Lambda}_1$ and $\bar{\Lambda}_2$, the prediction accuracy of lines from populations not represented in estimation set from lines in P_3 and P_4 (set \bar{II}).

List of Tables

467			
468	1	Average within population prediction accuracies in NAM maize populations	12
469	2	Average within population prediction accuracies in interconnected biparental maize	
470		populations	27
471	3	Average prediction accuracies for simulated maize populations	28

TABLE 2: Average within population prediction accuracies in interconnected biparental maize populations

N_p	Trait	Pooling		
		no	partial	complete
31	EL	0.31	0.33	0.31
	DON	0.38	0.44	0.46
	GER	0.38	0.43	0.43
	KR	0.46	0.50	0.52
	KpR	0.21	0.23	0.21
62	EL	0.40	0.41	0.39
	DON	0.47	0.51	0.51
	GER	0.47	0.50	0.49
	KR	0.53	0.56	0.58
	KpR	0.28	0.29	0.27
95	EL	0.44	0.46	0.43
	DON	0.51	0.53	0.53
	GER	0.51	0.53	0.53
	KR	0.56	0.58	0.59
	KpR	0.31	0.32	0.30

Values shown are average within population prediction accuracies for test individuals, averaged over 100 random estimation-test data splits. Standard errors were < 0.01 . N_p denotes the average number of individuals per population in the estimation set. The traits were ear length (EL), deoxynivalenol content (DON), Gibberella ear rot severity (GER) kernel rows (KR) and kernels per row (KpR)

TABLE 3: Average prediction accuracies for simulated maize populations

rSD	r_{Π}			$r_{\bar{\Pi}}$	
	no	partial	complete	partial	complete
0.0	0.54	0.89	0.89	0.89	0.89
0.25	0.51	0.84	0.85	0.84	0.84
0.5	0.50	0.76	0.76	0.73	0.73
1.0	0.48	0.57	0.53	0.48	0.49
2.0	0.44	0.41	0.30	0.20	0.21

Values shown are average within population prediction accuracies for test individuals, averaged over 50 random estimation-test data splits. Standard errors were < 0.015 . rSD is the relative standard deviation of simulated population specific QTL effects.