1      **Efficiency of genomic prediction of non-assessed single crosses**

2      José Marcelo Soriano Viana,[*1] Helcio Duarte Pereira,[1] Gabriel Borges Mundim,[†] Hans-Peter

3      Piepho,[‡] and Fabyano Fonseca e Silva[§]

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

5      [†]Down AgroSciences Seeds and Biotechnology Brazil Ltda, Indianópolis, MG, Brazil.

6      [‡]University of Hohenheim, Institute of Crop Science, Biostatistics Unit, 70599, Stuttgart, Germany.

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

8      Reference      number      for      data      available      in      public      repository:

9      https://doi.org/10.6084/m9.figshare.5035130.v1

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

11    Running title: Genomic prediction of single crosses.

12    **KEYWORDS** genomic selection; linkage disequilibrium; general combining ability; specific

13    combining ability; doubled haploids.

14    [1]Corresponding author: José Marcelo Soriano Viana. Federal University of Viçosa, Department of

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

16    +55(31)3899-2514.

17    **ABSTRACT** The objective was to provide additional relevant information on efficiency of

18    prediction of non-assessed single crosses. We derived the genetic model for genomic prediction.

19    The SNP and QTL genotypic data for DH lines, the QTL genotypic data of SCs, and the phenotypic

20    data for DH lines and SCs were simulated assuming 10,000 SNPs, 400 QTLs, two groups of 70

21    selected DH lines, and 4,900 SCs. The heritabilities for the assessed SCs were 30, 60 and 100%.

22    The scenarios included three sampling processes of DH lines, two sampling processes of SCs for

23    testing, two SNP densities, DH lines from distinct and same populations, DH lines from populations

24    with lower LD, two genetic models, three statistical models, and three statistical approaches. The

25    efficiency of prediction of untested SCs was based on the prediction accuracy and the efficacy of

26    identification of the best 300 (7-9%) non-assessed SCs (coincidence index), computed based on the

27    true genotypic values. Concerning the prediction accuracy and coincidence, our results proved that

28    prediction of untested SCs is very efficient. The accuracies and coincidences ranged from

29    approximately 0.8 and 0.5, respectively, under low heritability, to 0.9 and 0.7, assuming high

30    heritability. Additionally, we highlighted the relevance of the overall LD and evidenced that

31    efficient prediction of untested SCs can be achieved for crops that show no heterotic pattern, for

32    reduced training set size (10%), for SNP density of 1 cM, and for distinct sampling processes of DH

33    lines, based on random choice of the SCs for testing.

34                                            **INTRODUCTION**

35      Genomic selection is very commonly used in animal breeding programs, especially for dairy

36    cattle Van Eenennaam et al. (2014). The same cannot yet be said to the same degree concerning

37    crop breeding. The main reasons for the effective application of genomic selection in livestock

38    breeding are: it is efficient, that is, the process has high prediction accuracy, the cost of phenotyping

39    (mainly progeny test) is higher than the cost of genotyping, and the process significantly shortens

40    the selection cycle (Meuwissen et al. 2013). In spite of the many field and simulation-based studies

41    with genomic selection in plant breeding, in general the cost of phenotyping is often still much

42    lower than the cost of genotyping, restricting its application in breeding programs. Jonas and de

43    Koning (2013) consider that genomic selection has the potential to improve existing plant breeding

44    schemes. However, based also on the high diversity and complexity of plant breeding methods, they

45    stated that there are great obstacles to overcome.

46        An important application of genomic selection in plant breeding is the prediction of untested

47    single crosses (genotypic value prediction) and testcrosses (general combining ability effect

48    prediction) in hybrid breeding (Zhao et al. 2015). Prediction accuracy of barley two- and three-way

49    crosses has been investigated (Philipp et al. 2016). The prediction of untested single crosses was

50    pioneered by Bernardo (1994), based on best linear unbiased prediction (BLUP). Many significant

51    studies on prediction of untested single cross and testcross performance have been published in the

52    last 23 years, focused on the assessment of the prediction accuracy. Most investigations were based

53    on empirical data and estimated the prediction accuracy using a cross-validation procedure. Very

54    few were based on simulated data (Li et al. 2017; Technow et al. 2012a). With no exception, the

55    inference was that prediction of untested single crosses and testcrosses can be an efficient,

56    depending on heritability, training set size, and number of tested inbreds in hybrid combination

57    (both, one, and none parents tested). Remarkably, this conclusion was drawn from studies differing

58    in the type of molecular marker, density of markers, number of inbreds, level of relatedness,

59    diversity, and linkage disequilibrium (LD) between inbreds, heterotic pattern, training set size,

genetic model, and statistical approach (Zhao et al. 2015). Efficient prediction of barley two- and three-way crosses has been achieved when training and validation sets include the same class of hybrid (Philipp et al. 2016).

Most papers on genomic prediction of maize single cross performance published since 2011 have employed single nucleotide polymorphism (SNP), with the number SNPs filtered ranging from 425 (Zhao et al. 2013a) to 39,627 (Technow et al. 2012a). Based on the physical length of the maize genome (approximately 2,106 megabase pairs (Mb) according to Maize genetics and genomics database), the SNP density ranged from approximately 5 to 0.05 Mb, respectively. For grain yield, the relative prediction accuracies (computed as accuracy/root square of the heritability) in these two papers ranged from 0.27 to 0.62 and from 0.65 to 0.95, respectively. The number of inbreds in each heterotic group was highly variable too, ranging from six and nine (Bernardo 1994) to 75 and 75 (Technow et al. 2012a). The relative accuracy observed by Bernardo (1994) ranged between 0.72 and 0.89. The number of testcrosses ranged between 255 (Windhausen et al. 2012) and 1,894 (Albrecht et al. 2014). The relative accuracies ranged from 0.46 to 0.52 and from 0.33 to 0.65, respectively. The level of relatedness ranged from non-related inbreds in each group (Technow et al. 2012a) to a maximum average value of 0.58 (Bernardo 1995). The relative accuracy obtained by Bernardo (1995) ranged from 0.41 to 0.80. The common heterotic groups were Stiff Stalk and non-Stiff Stalk (Kadam et al. 1916) or Dent and Flint (Technow et al. 2014). The study of Bernardo (1996a) involved nine heterotic groups and the (statistically significant from zero) relative accuracies ranged from 0.43 to 0.88. No study provided clearly greater prediction accuracy of the additive-dominance model relative to the additive model. Finally, only with testcrosses the genomic BLUP (GBLUP) approach outperformed BLUP (Albrecht et al. 2014; Albrecht et al. 2011) concerning prediction accuracy.

Technow et al. (2012a) provided the most comprehensive study on prediction of untested single cross performance. Our assessment on the efficiency of prediction of non-assessed single

4

85    crosses provides additional relevant information. Our simulation-based study is the first to provide

86    for breeders a direct measure of efficiency of identification of the best non-assessed single crosses

87    (coincidence index), additionally to the standard prediction accuracy. What is the efficiency of

88    identification of the best 300 untested single crosses if the prediction accuracy is, for example,

89    approximately 0.90? Our results evidence that the efficacy range between 0.60 and 0.70, depending

90    on the doubled haploid (DH) lines derivation process. These measures of efficiency were provided

91    for a large data set (10,000 SNPs, 400 quantitative trait loci (QTLs), 4,900 single crosses) and for

92    low (30%) to high heritability (100%), assuming scenarios not favorable to prediction of non-

93    assessed single cross performance, as low level of relatedness and a not high heterotic pattern. Low

94    heritability has been observed in some CIMMYT's global maize and wheat breeding programs

95    (Crossa et al. 2014). Additionally, we derived the genetic model for genomic prediction, supported

96    by quantitative genetics theory, highlighted the relevance of the overall LD (not only for linked

97    SNPs and QTLs), and evidenced that efficient prediction of untested single crosses can be achieved

98    for crops that show no clear heterotic pattern, as rice, wheat, and barley, for reduced training set

99    size (10%), for SNP density of 1 centiMorgan (cM), and for distinct processes of (DH) lines

100   sampling. Finally, we showed that the choice of the single crosses for testing must be based on a

101   random process, but never by sampling DH or inbreds lines for a diallel. By sampling 76% of the

102   available genotyped DH lines in each group for a diallel (Technow et al. (2012a) sampled 75% of

103   the inbreds), the prediction accuracies and coincidence indexes were 38 to 77% and 39 to 98%

104   lower, respectively, compared with random sampling of 30% of the possible single crosses for

105   testing. Thus, our objective was to provide to breeders additional relevant information on prediction

106   of non-assessed single crosses.

107                                **MATERIALS AND METHODS**

108   **Theory**

5

109    Generally, most papers on genomic selection presents only statistical aspects and the genetic

110    models are deduced from gene to SNP effects. Importantly, when there is some quantitative

111    genetics theory, the LD is completely ignored. The theory developed provides, based on

112    quantitative genetics including LD, the genetic model for genomic prediction of single crosses. The

113    model developed offers the genetic background to the models fitted in important previously papers

114    on prediction of untested single crosses and testcrosses (Massman et al. 2013; Technow et al.

115    2012a; Albrecht et al. 2011). Notice, however, that the derived model has distinct presuppositions.

116    ***LD in a group of selected DH or inbred lines***

117    Consider a group of DH or inbred lines selected from a population or heterotic group. Assume

118    also a QTL (alleles B/b) and a SNP (alleles C/c) where B and b are the alleles that increase and

119    decrease the trait expression, respectively. Define the joint genotype probabilities as

120    $P(BBCC) = f_{22}$,   $P(BBcc) = f_{20}$,   $P(bbCC) = f_{02}$,   and   $P(bbcc) = f_{00}$, where the subscript

121    indicates the number of copies of the major allele (B and C). The measure of LD between the QTL

122    and the SNP is   $\Delta_{bc} = f_{22}f_{00} - f_{20}f_{02}$   (Kempthorne 1954) and the haplotype frequencies are

123    $P(BC) = f_{22} = p_b p_c + \Delta_{bc}$,     $P(Bc) = f_{20} = p_b q_c - \Delta_{bc}$,     $P(bC) = f_{02} = q_b p_c - \Delta_{bc}$,     and

124    $P(bc) = f_{00} = q_b q_c + \Delta_{bc}$, where  p  is the frequency of the major allele (B or C) and  $q = 1 - p$  is

125    the frequency of the minor allele (b or c). Notice that  $p_b = f_{22} + f_{20}$  and  $p_c = f_{22} + f_{02}$. It is

126    important to highlight the fact that we are not assuming that the QTL and the SNP are linked and in

127    LD in the population or heterotic group, because this is not necessary for genomic prediction. But

128    we are assuming that they are in LD in the group of DH or inbred lines. Furthermore, because

129    selection, genetic drift, and inbreeding (only for inbreds and linked QTLs and SNPs), the gene and

130    genotypic frequencies and the LD values concerning the selected DH or inbred lines cannot be

131    traced to the values in the population or heterotic group.

6

132 *SNP genotypic values of DH or inbred lines*

133 The average genotypic value for a group of selected DH or inbred lines is

134 $M_{IL} = m_b + \left(p_b - q_b\right)a_b$, where $m_b$ is the mean of the genotypic values of the homozygotes and

135 $a_b$ is the deviation between the genotypic value of the homozygote of higher expression and $m_b$.

136 Thus, the average SNP genotypic values for the DH or inbred lines CC and cc are

137 $$G_{CC} = \frac{1}{f_{.2}}\left[f_{22}\left(m_b + a_b\right) + f_{02}\left(m_b - a_b\right)\right] = M_{IL} + 2q_c\alpha_{SNP} = M_{IL} + A_{CC}$$

138 $$G_{cc} = \frac{1}{f_{.0}}\left[f_{20}\left(m_b + a_b\right) + f_{00}\left(m_b - a_b\right)\right] = M_{IL} - 2p_c\alpha_{SNP} = M_{IL} + A_{cc}$$

139 where $\alpha_{SNP} = \left[\dfrac{\Delta_{bc}}{p_c q_c}\right]a_b = \kappa_{bc}a_b$ is the average effect of a SNP substitution in the group of DH

140 or inbred lines and A is the SNP additive value for a DH or inbred line. Notice that E(A) = 0.

141 Assuming two QTLs (alleles B and b, and E and e) in LD with the SNP, the average effect of

142 a SNP substitution in the selected DH or inbred lines is $\alpha_{SNP} = \kappa_{bc}a_b + \kappa_{ce}a_e$, where

143 $\kappa_{ce} = \left[\dfrac{\Delta_{ce}}{p_c q_c}\right]$. Thus, in general, the average effect of a SNP substitution (and the SNP additive

144 value) is proportional to the measure of LD and to the a deviation for each QTL that is in LD with

145 the marker.

146 *SNP genotypic values of single crosses*

147 Aiming to maximize the heterosis, maize breeders commonly assess single crosses originating

148 from selected DH or inbred lines from distinct heterotic groups. Consider $n_1$ DH or inbred lines

149 from a population or heterotic group and $n_2$ DH or inbred lines from a distinct population or

7

150     heterotic group. The average genotypic value for the single crosses derived by crossing the DH or

151     inbred lines from group 1 with the DH or inbred lines from group 2 is

152     $M_H = m_b + \left( p_{b1}p_{b2} - q_{b1}q_{b2} \right) a_b + \left( p_{b1}q_{b2} + q_{b1}p_{b2} \right) d_b$

153     where $d_b$ is the dominance deviation (the deviation between the genotypic value of the

154     heterozygote and $m_b$).

155     The average genotypic values for the single crosses derived from DH or inbred lines CC and

156     cc of the group 1 are

157     $$M_{CC1} = M_H + q_{c1}\kappa_{bc1}\left[ a_b + \left( q_{b2} - p_{b2} \right)d_b \right] = M_H + q_{c1}\kappa_{bc1}\alpha_{b2} = M_H + q_{c1}\alpha_{SNP1}$$
$$= M_H + GCA_{CC1}$$

158     $M_{cc1} = M_H - p_{c1}\kappa_{bc1}\alpha_{b2} = M_H - p_{c1}\alpha_{SNP1} = M_H + GCA_{cc1}$

159     where $\alpha_{b2}$ is the average effect of allelic substitution in the population derived by random crosses

160     between the DH or inbred lines of group 2, $\alpha_{SNP1}$ is the SNP effect of allelic substitution in the

161     hybrid population relative to a SNP derived from group 1, and GCA stands for the general

162     combining ability effect for a SNP locus. Notice that $\alpha_{SNP1}$ depends on the LD in group 1

163     ( $\kappa_{bc1} = \Delta_{bc1}/p_{c1}q_{c1}$ ) and the average effect of allelic substitution in the population derived by

164     random crosses between the DH or inbred lines of group 2. Further,

165     $E(GCA) = p_{c1}GCA_{CC1} + q_{c1}GCA_{cc1} = 0$. Concerning the single crosses derived from DH or

166     inbred lines CC and cc of the group 2 we have

167     $$M_{CC2} = M_H + q_{c2}\kappa_{bc2}\left[ a_b + \left( q_{b1} - p_{b1} \right)d_b \right] = M_H + q_{c2}\kappa_{bc2}\alpha_{b1} = M_H + q_{c2}\alpha_{SNP2}$$
$$= M_H + GCA_{CC2}$$

8

168 $$M_{cc2} = M_H - p_{c2}\kappa_{bc2}\alpha_{b1} = M_H - p_{c2}\alpha_{SNP2} = M_H + GCA_{cc2}$$

169    Notice that $E(GCA) = 0$ also. The average genotypic values for the single crosses concerning

170    the SNP locus are

171
$$M_{CC1xCC2} = M_H + q_{c1}\alpha_{SNP1} + q_{c2}\alpha_{SNP2} - 2q_{c1}q_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{CC1} + GCA_{CC2} + SCA_{CC1xCC2}$$

172
$$M_{cc1xcc2} = M_H - p_{c1}\alpha_{SNP1} - p_{c2}\alpha_{SNP2} - 2p_{c1}p_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{cc1} + GCA_{cc2} + SCA_{cc1xcc2}$$

173
$$M_{CC1xcc2} = M_H + q_{c1}\alpha_{SNP1} - p_{c2}\alpha_{SNP2} + 2q_{c1}p_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{CC1} + GCA_{cc2} + SCA_{CC1xcc2}$$

174
$$M_{cc1xCC2} = M_H - p_{c1}\alpha_{SNP1} + q_{c2}\alpha_{SNP2} + 2p_{c1}q_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{cc1} + GCA_{CC2} + SCA_{cc1xCC2}$$

175    where $\kappa_{bc1}\kappa_{bc2}d_b = d_{SNP}$ is the SNP dominance deviation in the hybrid population and SCA

176    stands for the specific combining ability effect for a SNP locus. Notice that $E(SCA) =$

177    $p_{c1}p_{c2}SCA_{CC1xCC2} + p_{c1}q_{c2}SCA_{CC1xcc2} + q_{c1}p_{c2}SCA_{cc1xCC2} + q_{c1}q_{c2}SCA_{cc1xcc2} = 0$ and

178    , for each group, $E(SCA|CC) = E(SCA|cc) = 0$. That is, the expectation of the SNP SCA effects

179    given a SNP genotype for the common DH or inbred line is also zero. Notice also that the four

180    genotypic values depends on four parameters ($M_H$, $\alpha_{SNP1}$, $\alpha_{SNP2}$, and $d_{SNP}$).

181    Assuming two QTLs (alleles B and b, and E and e) in LD with the SNP, the SNP dominance

182    deviation is $d_{SNP} = \kappa_{bc1}\kappa_{bc2}d_b + \kappa_{ce1}\kappa_{ce2}d_e$. Thus, generally, the SNP dominance deviation

183    (and the SNP SCA effect) is proportional to the product of the LD values in both groups of DH or

184    inbred lines and to the dominance deviation for each QTL that is in LD with the marker.

9

185        The previous model expressed as a function of the GCA and SCA effects is that proposed by

186    Massman et al. (2013), but these authors assumed $GCA_{CC} + GCA_{cc} = 0$ (for each heterotic group

187    and for each SNP) and $SCA_{CC1xCC2} = SCA_{cc1xcc2} = -SCA_{CC1xcc2} = -SCA_{cc1xCC2}$.

188    Technow et al. (2012b) have used a standard extension from QTL to SNP, defining the single cross

189    genotypic value for a SNP as a function of the SNP a and d deviations. That is,

190    $M = M_H + u_1 a_1 + u_2 a_2 + u_3 d$, where $u_1$ and $u_2$ equal to 1/2 or $-1/2$ if the corresponding DH or

191    inbred line is homozygous for distinct SNP alleles (CC or cc), and $u_3$ equal to 0 if the single cross

192    is homozygous or 1 if heterozygous.

193    ***SNP genotypic values of single crosses from DH or inbred lines derived from the same***

194    ***population or heterotic group***

195       Well defined heterotic groups are known for maize, but not for special maize as popcorn and

196    sweet corn and for other crops as wheat (Zhao et al. 2013b), rice (Xu et al. 2014), and barley

197    (Philipp et al. 2016). Thus, for many breeders, it is interesting to know about the efficiency of

198    genomic prediction of singles crosses when there are no heterotic groups. Assuming n DH or inbred

199    lines derived from the same population or heterotic group, the average genotypic values for the

200    single crosses concerning the SNP locus are

201    $M_{CCxCC} = M + 2q_c \alpha_{SNP} - 2q_c^2 \kappa_{bc}^2 d_b = M + 2GCA_{CC} + SCA_{CCxCC}$

202    $M_{ccxcc} = M - 2p_c \alpha_{SNP} - 2p_c^2 \kappa_{bc}^2 d_b = M + 2GCA_{cc} + SCA_{ccxcc}$

203    $M_{CCxcc} = M + 2(q_c - p_c) \alpha_{SNP} + 2p_c q_c \kappa_{bc}^2 d_b = M + GCA_{CC} + GCA_{cc} + SCA_{CCxcc}$

204    where    $M = m_b + (p_c - q_c) a_b + 2p_c q_c d_b$    is    the    hybrid    population    mean,

205    $\alpha_{SNP} = \kappa_{bc}[a_b + (q_b - p_b) d_b] = \kappa_{bc} \alpha_b$ is the average effect of a SNP substitution in the hybrid

206    population, and $d_{SNP} = \kappa_{bc}^2 d_b$ is the SNP dominance deviation. Notice that the SNP GCA effects

207 are equal to half the SNP additive value for the single crosses (A), the SNP SCA effects are the SNP

208 dominance deviations for the single crosses (D), and that the three genotypic values depends on

209 three parameters ($M$, $\alpha_{SNP}$, and $d_{SNP}$). Notice also that E(GCA) = E(A) = E(SCA) =

210 E(SCA|CC) = E(SCA|cc) = E(D) = 0.

### *Accuracy of single cross genomic prediction*

212     Assuming a QTL and a SNP in LD in the two groups of DH or inbred lines, the predictor of

213 the single cross QTL genotypic value is the single cross SNP genotypic value (because they are

214 proportional). Thus, the covariance between the predictor and the genotypic value is

$$
\begin{aligned}
\text{Cov}\left(\tilde{G}, G\right) &= f_{22}^1 f_{22}^2 \left[ M_H + GCA_{CC1} + GCA_{CC2} + SCA_{CC1xCC2} \right]\left[ M_H + GCA_{BB1} + GCA_{BB2} + SCA_{BB1xBB2} \right] + \\
&\quad + f_{22}^1 f_{20}^2 \left[ M_H + GCA_{CC1} + GCA_{cc2} + SCA_{CC1xcc2} \right]\left[ M_H + GCA_{BB1} + GCA_{BB2} + SCA_{BB1xBB2} \right] + \\
&\quad \ldots \\
&\quad + f_{00}^1 f_{00}^2 \left[ M_H + GCA_{cc1} + GCA_{cc2} + SCA_{cc1xcc2} \right]\left[ M_H + GCA_{bb1} + GCA_{bb2} + SCA_{bb1xbb2} \right] - \left( M_H \right)^2 \\
&= p_{c1}q_{c1}\left( \kappa_{bc1}\alpha_{b2} \right)^2 + p_{c2}q_{c2}\left( \kappa_{bc2}\alpha_{b1} \right)^2 + 4p_{c1}q_{c1}p_{c2}q_{c2}\left( \kappa_{bc1}\kappa_{bc2}d_b \right)^2 \\
&= p_{c1}q_{c1}\left( \alpha_{SNP1} \right)^2 + p_{c2}q_{c2}\left( \alpha_{SNP2} \right)^2 + 4p_{c1}q_{c1}p_{c2}q_{c2}\left( d_{SNP} \right)^2 \\
&= \sigma_{GCA_{SNP}}^{2(1)} + \sigma_{GCA_{SNP}}^{2(2)} + \sigma_{SCA_{SNP}}^2 = \sigma_{G(SNP)}^2
\end{aligned}
$$

216

217 where the GCA and SCA effects for the QTL are $GCA_{BB1} = q_{b1}\alpha_{b2}$, $GCA_{bb1} = -p_{b1}\alpha_{b2}$,

218 $GCA_{BB2} = q_{b2}\alpha_{b1}$, $\quad\quad\quad GCA_{bb2} = -p_{b2}\alpha_{b1}$, $\quad\quad\quad SCA_{BB1xBB2} = -2q_{b1}q_{b2}d_b$,

219 $SCA_{BB1xbb2} = 2q_{b1}p_{b2}d_b$, $\quad SCA_{bb1xBB2} = 2p_{b1}q_{b2}d_b$, and $\quad SCA_{bb1xbb2} = -2p_{b1}p_{b2}d_b$,

220 $\sigma_{GCA}^2$ and $\sigma_{SCA}^2$ are the GCA and SCA variances for the SNP locus, and $\sigma_G^2$ is the SNP

221 genotypic variance. The GCA and SCA variances for the QTL are $\sigma_{GCA}^{2(1)} = p_{b1}q_{b1}\left( \alpha_{b2} \right)^2$,

222     $\sigma_{GCA}^{2(2)} = p_{b2}q_{b2}\left(\alpha_{b1}\right)^2$, and $\sigma_{SCA}^2 = 4p_{b1}q_{b1}p_{b2}q_{b2}\left(d_b\right)^2$. The QTL genotypic variance is

223     $\sigma_G^2 = \sigma_{GCA}^{2(1)} + \sigma_{GCA}^{2(2)} + \sigma_{SCA}^2$ Thus, the single cross prediction accuracy is

224     $\rho_{\widetilde{G},G} = \sqrt{\dfrac{\sigma_{G(SNP)}^2}{\sigma_G^2}}$

225     Assuming s SNPs,

226     $\rho_{\widetilde{G},G} = \displaystyle\sum_{r=1}^{s} \sigma_{G(SNP(r))}^2 \Big/ \sqrt{\sigma_{\widetilde{G}}^2 \sigma_G^2}$

227     where $\sigma_{\widetilde{G}}^2$ is the variance of the predicted single cross genotypic values and $\sigma_G^2$ is the single cross

228     genotypic variance. Further,

229     $\alpha_{SNP(r)1} = \displaystyle\sum_{i=1}^{k'} \left[\dfrac{\Delta_{ri1}}{p_{r1}q_{r1}}\right]\alpha_{i2} = \sum_{i=1}^{k'} \kappa_{ri1}\alpha_{i2}$, where k' is the number of QTLs in LD with the SNP

230                                   r) in group 1, and

231     $d_{SNP(r)} = \displaystyle\sum_{i=1}^{k''} \left[\dfrac{\Delta_{ri1}}{p_{r1}q_{r1}}\right]\left[\dfrac{\Delta_{ri2}}{p_{r2}q_{r2}}\right]d_i = \sum_{i=1}^{k''} \kappa_{ri1}\kappa_{ri2}d_i$ where k'' is the number of QTLs in LD with

232                                 the SNP r in both groups

233     Notice that because the accuracy of genomic prediction of single crosses depends on the

234     squares of the average effects of SNP substitution and the SNP dominance deviations, it is not

235     affected by the linkage phase (coupling or repulsion), as it does not depend on linkage. But it

236     depends on the magnitude of the LD in each group of DH or inbred lines.

12

237      Assuming single crosses derived from DH or inbred lines of a single population or heterotic

238      group      we      have      $\sigma^2_{G(SNP)} = 2p_c q_c (\alpha_{SNP})^2 + (2p_c q_c d_{SNP})^2$      and

239      $\sigma^2_G = 2p_b q_b (\alpha_b)^2 + (2p_b q_b d_b)^2$.

240      **The statistical model for single cross genomic prediction**

241      Assume $n_1$ and $n_2$ (several tens) DH or inbred lines from two populations or heterotic groups

242      genotyped for s (thousands) SNPs and the experimental assessment of h (few hundred) single-

243      crosses (h much lower than $n_1.n_2$) in e (several) environments (a combination of growing seasons,

244      years, and locals). Defining y as the adjusted single cross phenotypic mean, the statistical model

245      for prediction of the average effects of SNP substitution and the SNP dominance deviations is

246      $y = M_H + \sum_{r=1}^{s} \left( z_{1_r} \alpha_{SNP1_r} + z_{2_r} \alpha_{SNP2_r} + z_{3_r} d_{SNP_r} \right) + error$

247      where $z_{1_r} = q_{r1}$, $z_{2_r} = q_{r2}$, and $z_{3_r} = -2q_{r1}q_{r2}$ if the SNP genotypes for the DH or inbred lines

248      are CC (group 1) and CC (group 2), $z_{1_r} = -p_{r1}$, $z_{2_r} = -p_{r2}$, and $z_{3_r} = -2p_{r1}p_{r2}$ if the SNP

249      genotypes for the DH or inbred lines are cc (group 1) and cc (group 2), $z_{1_r} = q_{r1}$, $z_{2_r} = -p_{r2}$, and

250      $z_{3_r} = 2q_{r1}p_{r2}$ if the SNP genotypes for the DH or inbred lines are CC (group 1) and cc (group 2),

251      and $z_{1_r} = -p_{r1}$, $z_{2_r} = q_{r2}$, and $z_{3_r} = p_{r1}q_{r2}$ if the SNP genotypes for the DH or inbred lines are

252      cc (group 1) and CC (group 2).

253      Regarding the single crosses obtained from DH or inbred lines of the same population or

254      heterotic group we have

255      $y = M + \sum_{r=1}^{s} \left( z_{1_r} \alpha_{SNP_r} + z_{2_r} d_{SNP_r} \right) + error$

256    where $z_{1_r} = 2q_r$ and $z_{2_r} = -2q_r^2$ if the SNP genotypes for the DH or inbred lines are CC and CC,

257    $z_{1_r} = -2p_r$ and $z_{2_r} = -2p_r^2$ if the SNP genotypes for the DH or inbred lines are cc and cc, and

258    $z_{1_r} = 2(q_r - p_r)$ and $z_{2_r} = 2p_r q_r$ if the SNP genotypes for the DH or inbred lines are CC and cc.

259        The statistical problem of genomic prediction when there are a very large number of

260    molecular markers and relatively few observations have been addressed thorough several

261    regularized whole-genome regression and prediction methods (Daetwyler et al. 2013; de Los

262    Campos et al. 2013). Then, the predicted effects of SNP substitution ($\tilde{\alpha}$) and SNP dominance

263    deviations ($\tilde{d}$) must be used to provide genomic prediction of non-assessed single crosses. The

264    predicted genotypic value for a non-assessed single cross of DH or inbred lines from two groups is

265    $$\widetilde{G} = \hat{M}_H + \sum_{r=1}^{s} \left( z_{1_r} \tilde{\alpha}_{SNP1_r} + z_{2_r} \tilde{\alpha}_{SNP2_r} + z_{3_r} \tilde{d}_{SNP_r} \right)$$

266        For a non-assessed single cross of DH or inbred lines from the same group, the predicted

267    genotypic value is

268    $$\widetilde{G} = \hat{M} + \sum_{r=1}^{s} \left( z_{1_r} \tilde{\alpha}_{SNP_r} + z_{2_r} \tilde{d}_{SNP_r} \right)$$

269    **Simulation**

270        The SNP and QTL genotypic data for DH lines, the QTL genotypic data of single crosses, and

271    the phenotypic data for DH lines and single crosses were simulated using the software

272    *REALbreeding*. The program has been developed by the first author using the software *REALbasic*

273    *2009* (Viana et al. 2017a; Viana et al. 2017b; Viana et al. 2016; Azevedo et al. 2015; Viana et al.

274    2013). Based on our input, the software distributed 10,000 SNPs and 400 QTLs in ten

275    chromosomes (1,000 SNPs and 40 QTLs by chromosome). The average SNP density was 0.1 cM.

276    The QTLs were distributed in the regions covered by the SNPs (approximately 100

277    cM/chromosome). Initially, *REALbreeding* sampled 700 DH lines from two non-inbred populations

278   (heterotic groups) in LD (350 from each population). The populations were composites of two

279   populations in linkage equilibrium. In a composite, there is LD only for linked SNPs and QTLs

280   (Viana et al. 2016). The number of DH lines from each $S_0$ plant was one (scenario 1) or ranged

281   from 1 to 5 (scenario 2). We also sampled 350 DH lines from each population after three

282   generations of selfing (using the single seed descent process). The number of DH lines from each $S_3$

283   plant ranged from 1 to 5 (scenario 3). For each scenario, the software then crossed 70 selected DH

284   lines from each population, using a diallel design. The heritability for the DH lines was 30%.

285       The genotypic values of the DH lines and of the single crosses were generated assuming a

286   single set of 400 QTLs and two degrees of dominance. To simulate grain yield and expansion

287   volume, a measure of popcorn quality, we defined positive dominance ($0 < (d/a)_i \leq 1.2$, $i = 1, ...,$

288   400) and bidirectional dominance ($-1.2 \leq (d/a)_i \leq 1.2$), respectively, where d/a is the degree of

289   dominance. To compute the genotypic values, *REALbreeding* used our input relative to the

290   maximum and minimum genotypic values for homozygotes. For grain yield and expansion volume,

291   we defined 140 and 30 g/plant and 55 and 15 mL/g, respectively. The phenotypic values were

292   obtained from the sum of the population mean, genotypic value, and experimental error. The error

293   variance was computed from the broad sense heritability. To avoid outliers, we defined the

294   maximum and minimum phenotypic values as 160 and 10 g/plant and 65 and 5 mL/g.

295       The heritabilities for the assessed single crosses were 30, 60, and 100%. Thus, the genotypic

296   value prediction accuracies of the assessed single crosses were 0.55, 0.77, and 1.00, respectively.

297   For each scenario were processed 50 resamplings of 30 and 10% of the single crosses (1,470 and

298   490 assessed single crosses). That is, we predicted 70 and 90% of the single crosses (3,430 and

299   4,410 non-assessed single crosses). Additionally, to assess the relevance of the number of DH lines

300   sampled, we fixed the number of DH lines to achieve the same number of assessed single crosses,

301   using a diallel. That is, we sampled 50 times 38 and 22 DH lines in each group for a diallel

302   (scenario 4), generating 1,444 and 484 single crosses for assessment, respectively. We called these

15

303  processes as sampling of single crosses (scenarios 1 to 3) and sampling of DH lines (scenario 4).

304  Other additional scenarios were: genomic prediction of single crosses from selected DH lines from

305  same heterotic group (interestingly for wheat, rice, and barley breeders, for example) (scenario 5)

306  and from selected DH lines from populations with lower LD (scenario 6), to emphasize that the

307  prediction accuracy depends on the LD in the groups of DH or inbred lines. A last scenario

308  (seventh) was genomic prediction of single crosses under an average density of one SNP each cM.

309  This lower density was obtained by random sampling of 100 SNPs per chromosome using a

310  *REALbreeding* tool (*sampler*). To investigate the single cross prediction efficiency based on our

311  model and on the models proposed by Massman et al. (2013) and Technow et al. (2012b), we used

312  another *REALbreeding* tool (*Incidence matrix*) to generate the incidence matrices for the three

313  models and for the two DH lines sampling processes. To assess the relevance of the SCA effects

314  prediction on genomic prediction of single cross performance, we also fitted the additive model

315  (including only the GCA effects). For comparison purpose, we also processed single cross

316  prediction based on GBLUP (with the observed additive and dominance relationship matrices) and

317  BLUP (with the expected additive and dominance relationship matrices).

318  **Statistical analysis**

319      The methods used for prediction were ridge regression BLUP (RR-BLUP), GBLUP and

320  BLUP. For the analyses we used the *rrBLUP* package (Endelman 2011). The accuracies of single

321  cross genotypic value prediction were obtained by the correlation between the true values of the

322  non-assessed single crosses computed by *REALbreeding* and the values predicted by RR-BLUP,

323  GBLUP, and BLUP. We also computed the efficiency of identification of the 300 non-assessed

324  single crosses of higher genotypic value (coincidence index). The parametric average coincidence

325  index was computed by ordering the average phenotypic values of the 4,900 single crosses for each

326  heritability and for each DH lines derivation process. Regarding grain yield, for heritability of 30%

327  the coincidence index was 0.2533, 0.2833, and 0.2433 assuming one DH line per $S_0$ plant, one to

328  five DH lines per $S_0$ plant, and one to five DH lines per $S_3$ plant, respectively. The corresponding

329  values for heritability of 60% were, respectively, 0.4800, 0.4900, and 0.4567. Concerning

330  expansion volume, the corresponding values for heritabilities of 30 and 60% were, respectively,

331  0.2600, 0.2833, and 0.2700, and 0.4733, 0.5100, and 0.4533. The assumed average parametric

332  coefficient index was 0.26 and 0.48 for heritabilities of 30 and 60%, respectively, for both traits.

333  For the population structure analysis we employed *Structure* (Falush et al. 2003) and fitted the no

334  admixture model with independent allelic frequencies. The number of SNPs, sample size, burn-in

335  period, and number of MCMC (Markov chain Monte Carlo) replications were 1,000 (sampled at

336  random), 140 (70 DH lines from each population), 10,000, and 40,000, respectively. The number of

337  populations assumed (*K*) ranged from 1 to 4, and the most probable *K* value was determined based

338  on the inferred plateau method (Viana et al. 2013). The LD analyses were performed with

339  *Haploview* (Barrett et al. 2005).

340  **Data availability**

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

342  https://doi.org/10.6084/m9.figshare.5035130.v1. Data citation:

343  Viana, José Marcelo Soriano; Pereira, Helcio Duarte; Mundim, Gabriel Borges; Piepho, Hans-Peter;

344  Fonseca e Silva, Fabyano (2017): Efficiency of genomic prediction of non-assessed single crosses.

345  figshare. https://doi.org/10.6084/m9.figshare.5035130.v1

346  **RESULTS**

347  The parametric mean and genotypic variance in the populations 1 and 2 were 108.5 and 87.3

348  (g/plant) and 4.7680 and 6.2580 (g/plant)$^2$. The DH lines derivation processes (one and one to five

349  per $S_0$ plant and one to five per $S_3$ plant) provided, for each population, selected DH lines with

350  similar mean (approximately 97 and 76 g/plant for populations 1 and 2), inbreeding depression

351  (approximately −10 and −13% for populations 1 and 2), and genotypic variance (approximately 6

352  and 7 (g/plant)$^2$ for populations 1 and 2) and groups of single crosses also similar for mean

17

353    (approximately 103 g/plant), heterosis (approximately 19%), and genotypic variance

354    (approximately 4 $(g/plant)^2$). Because we derived one to few DH lines from unrelated $S_0$ and $S_3$

355    plants, the average level of relatedness between the selected DH lines was very low (zero and zero,

356    0.0041 and 0.0041, and 0.0054 and 0.0074 assuming one DH line per $S_0$, one to five DH lines per

357    $S_0$, and one to five DH lines per $S_3$, for populations 1 and 2, respectively). Concerning SNP data,

358    the frequency distribution of the minor allele frequency (MAF) and the absolute value of the

359    difference between a SNP allele frequency were also similar for both groups of selected DH lines,

360    regardless of the DH line derivation process (Figure 1a, b, c). The average MAF was 0.33,

361    regardless of the population and DH line derivation process. However, the evidence obtained by the

362    population structure analysis was that the DH lines belong to two distinct subpopulations (suggested

363    $K$ equal to 2.4 by the inferred plateau method). The percentages of non-polymorphic SNPs were

364    very low (0.1 to 0.4%). No differences between allelic frequencies were observed for only 1.7 to

365    2.1% of the SNPs. For approximately 70% of the SNPs, the absolute difference between allelic

366    frequencies ranged from 0.1 to 0.6. Regarding LD, for the groups of selected DH lines the evidence

367    based on the analysis of chromosome 1 (no difference between chromosomes is expected) is that

368    LD extents for up to 35 cM, regardless of the DH lines derivation process (Figure 1c, d). Ignoring

369    the non-significant LD values (LOD score lower than 3), for 17 to 20% of the SNP pairs the $r^2$

370    values ranged from 0.2 to 0.5 (average of 0.16, regardless of the DH lines group and derivation

371    process).

372        Assuming our model, average SNP density of 0.1 cM, training set size of 30%, positive

373    dominance (grain yield), additive-dominance model, and sampling of single crosses, the prediction

374    accuracies of the non-assessed single crosses were greater than the accuracies of the assessed single

375    crosses for low (up to 46% higher) and intermediate (up to 16% higher) heritabilities (Table 1;

376    Figure 2a). As the prediction accuracy of assessed single crosses approaches 1.0, the accuracy of the

377    non-assessed single crosses approaches approximately 0.9 (up to 11% lower). Sampling one to five

378   DH lines per $S_3$ plant was only slightly superior to the other DH lines derivation processes,

379   regardless of the prediction accuracy of the assessed single crosses (up to 5% higher). Fitting the

380   additive model provided essentially the same prediction accuracies since the maximum decrease

381   was approximately 1%. No significant differences between the prediction accuracies of non-

382   assessed single crosses were also observed assuming bidirectional dominance (expansion volume).

383   The differences compared to positive dominance ranged from approximately −5 to 2%. However, a

384   striking difference was observed between the sampling processes of single crosses for testing.

385   Random sampling of single crosses provided much greater prediction accuracies of non-assessed

386   single crosses, compared to sampling DH lines for a diallel. The increases in the accuracies by

387   sampling single crosses ranged from approximately 38 to 77%, proportional to the heritability.

388   Decreasing the average SNP density to 1 cM led to a slight decrease in the prediction accuracy of

389   non-assessed single crosses of approximately −4%). Decreasing the training set size to 10%

390   decreased the prediction accuracy of non-assessed single crosses in approximately −5 to −15%,

391   inversely proportional to the heritability. To evidence that the prediction accuracy of non-assessed

392   single crosses depends on the level of (overall) LD in the groups of selected DH or inbred lines, we

393   derived DH lines from the same base populations after 10 generations of random crosses (to

394   decrease the LD). The accuracies were also high, ranging from 0.83 to 0.95, proportional to the

395   heritability. The prediction accuracies of non-assessed single crosses from DH lines of the same

396   population were equivalent to the accuracies for single crosses derived from DH lines belonging to

397   distinct heterotic groups, ranging from 0.83 to 0.91, also proportional to the heritability. Comparing

398   our statistical model with the models proposed by Massman et al. (2013) and Technow et al.

399   (2012a), we observed no differences for the prediction accuracies of non-assessed single crosses

400   (maximum difference of 1%). Finally, no significant differences between the prediction accuracies

401   for RR-BLUP, GBLUP, and BLUP occurred (maximum of 2%), excepting for one to five DH lines

402   per $S_3$ plant, where BLUP was 9 to 10% inferior, regardless of the heritability.

19

403    Concerning the coincidence index, in general the inferences are the same established from the

404    prediction accuracy analysis (Table 2; Figure 2b). There were no differences between the

405    coincidence indexes regarding our model and the models proposed by Massman et al. (2013) and

406    Technow et al. (2012a) (maximum difference of 3%), and between the RR-BLUP, GBLUP, and

407    BLUP approaches, except for one to five DH lines per $S_3$ plant, where BLUP was $-19$ to $-27\%$

408    inferior, proportional to the heritability. The coincidence indexes were also high for single crosses

409    derived from selected DH lines obtained from the base populations with lower LD (ranging from

410    0.55 to 0.76, proportional to the heritability) and from selected DH lines of the same population

411    (ranging from 0.61 to 0.76, also proportional to the heritability). Sampling single crosses for

412    assessment also provided much greater coincidence index compared to sampling DH lines for a

413    diallel (39 to 98% higher, proportional to the heritability). Decreasing the SNP density and the

414    training set size decreased the coincidence index from 5 to 10% (proportional to the heritability)

415    and from 17 to 26% (inversely proportional to the heritability), respectively. The maximum

416    difference in the coincidence index by fitting the additive-dominant and the additive models was

417    $-3\%$. Only for one DH line per $S_0$ plant the coincidence indexes assuming bidirectional dominance

418    were slightly greater than the values assuming positive dominance (9 to 14% greater). This

419    sampling process of DH lines provided the higher values of coincidence index, compared to the

420    other sampling processes (7 to 26% higher, inversely proportional to the heritability). Finally, the

421    coincidence index of the non-assessed single crosses are greater than the parametric values for all

422    assessed single crosses assuming low (up to 117% higher) and intermediate (up to 39% higher)

423    heritabilities (Table 1). However, as the parametric coincidence of assessed single crosses

424    approaches 1.0, the coincidence values of the non-assessed single crosses approach approximately

425    0.60 to 0.74 (up to 26 to 40% lower), depending on the DH line sampling process.

426                                    **DISCUSSION**

20

427     It was twenty-three years ago today, Bernardo (1994) taught the breeders to use BLUP (more

428     precisely, GBLUP) for predicting untested maize single cross performance. BLUP, as well known,

429     is the Henderson's (1974) approach for genetic assessment. Based on the prediction accuracies

430     obtained by Bernardo (1994, 1995, 1996a, 1996b, 1996c), for grain yield and other traits (distinct

431     genetic controls), a breeder should realize that the performance of untested single crosses can be

432     effectively predicted using relationship information from molecular or pedigree data, unbalanced

433     and large data set, and diverse heterotic patterns. This general inference has been confirmed with

434     maize (Zhao et al. 2015) and other important crops, as rice (Xu et al. 2014), wheat (Zhao et al.

435     2013b) and barley (Philipp et al. 2016), along the last 20 years. Why, then, there is no published

436     evidence that prediction of untested single crosses is of general use by breeders of worldwide seed

437     companies? What should be additionally proved to make prediction of untested single crosses as

438     successful as the Jenkins' (1934) method for predicting double crosses performance was? We

439     believe that this paper offers a significant contribution.

440     Our assessment on efficiency of prediction of untested single cross performance keeps some

441     similarities with few earlier studies but sharp differences for most previous investigations. This

442     study is based on simulated data set, as the study of Technow et al. (2012a), assuming 400 QTLs

443     distributed along ten chromosomes. Thus, the prediction accuracies and coincidence indexes (a

444     measure of untested single crosses selection efficiency) are for really non-assessed single crosses

445     since the values were computed based on the true genotypic values of the non-assessed single

446     crosses and not on a cross-validation procedure involving assessed single crosses. This does not

447     mean that we consider simulated data better than field data or have any criticism on the cross-

448     validation procedure. We know that simulated data, because the presuppositions, cannot integrally

449     describe the complexity of populations and genetic determination of traits (Daetwyler et al. 2013).

450     To highlight the relevance of (overall) LD, our study is based on scenarios not favorable to

451     prediction of untested single cross performance: very low level of relationship between the DH

452    lines, low and intermediate heritabilities for the assessed single crosses, and not higher heterotic

453    pattern. In the studies of Massman et al. (2013) and Bernardo (1994, 1995, 1996a) the relationship

454    among inbreds from the same heterotic group ranged from 0.11 to 0.58. Riedelsheimer et al. (2012)

455    observed high relationship only between the non-Stiff Stalk inbreds. Technow et al. (2012a)

456    assumed non-related inbreds. For most of the investigations on prediction of untested single crosses

457    and testcrosses, the grain yield heritability ranged from 0.72 to 0.88. The common heterotic patterns

458    in these previous studies are Stiff Stalk and non-Stiff Stalk, and Dent and Flint. The MAF in the

459    groups of Dent and Flint inbreds were approximately 0.10 and 0.20, respectively, and

460    approximately 20% of the SNPs showed a difference of allelic frequency of at least 0.6.

461        Concerning the prediction accuracy and the efficiency of identification of the superior 300

462    non-assessed single crosses, our results prove that prediction of untested single crosses is a very

463    efficient procedure (note that we are not saying genomic prediction), specially for low and

464    intermediate heritabilities of the assessed single crosses. The prediction accuracy of the non-

465    assessed single crosses under low (0.55 to 0.71) and intermediate (0.74 to 0.87) accuracies of

466    assessed single crosses achieved 0.85 and 0.89, respectively. It is important to highlight that these

467    are not relative accuracies. Most important, the coincidence of the non-assessed single crosses

468    under low (0.26 to 0.39) and intermediate (0.44 to 0.66) parametric coincidences of assessed single

469    crosses achieved 0.59 and 0.64, respectively. For high heritability (80 to 95%; accuracies from 0.89

470    to 0.97), as observed in most of the studies on prediction of untested single cross performance, we

471    can state (based on values predicted by fitting a quadratic regression model) that the prediction

472    accuracy of non-assessed single crosses is up to only 10% lower (0.87 to 0.92) and, most

473    impressive, the coincidence index can range from 0.61 to 0.71 (parametric coincidences between

474    0.72 to 0.93). Under maximum accuracy of assessed single crosses (1.0), the prediction accuracy

475    and coincidence of non-assessed single crosses achieved 0.93 and 0.76. Thus, assuming high

476    heritability, high density, and training set size of 30%, the accuracy can achieve 0.92 and the

22

477    efficiency of identification of the best 9% of the non-assessed single crosses can achieve 0.71. It is

478    important to highlight that this efficacy can be higher by using more related DH or inbred lines,

479    under high LD. Thus, we strong recommend that maize breeders, as well as rice, wheat, and barley

480    breeders, make widespread use of prediction of non-assessed single crosses, at least for preliminary

481    screening or prior to field testing.

482    To take advantage of genomic prediction, Kadam et al. (2016) recommend redesigning hybrid

483    breeding programs. However, because breeders are unlikely to rely solely on genomic predictions

484    when selecting superior untested hybrids, Technow et al. (2014) believe that genomic prediction

485    will be combined with field testing of the most promising experimental hybrids. For grain yield, the

486    prediction accuracies observed by Bernardo (1994, 1995, 1996a) ranged from 0.14 to 0.80,

487    proportional to the heritability (in the range 35-74%) and training set size. The non-relative

488    accuracies (relative accuracy x root square of heritability) observed in the studies of Kadam et al.

489    (2016), Technow et al. (2014), Massman et al. (2013), Technow et al. (2012a), and Riedelsheimer et

490    al. (2012) ranged between 0.20 and 0.86, also proportional to the heritability (in the range 53-98%)

491    and training set size.

492    We hope that readers of this paper have realized the importance of (overall) LD for effective

493    prediction of non-assessed single crosses, as well as genetic variability (see the parametric accuracy

494    of genomic prediction). Breeders have no control over LD and relatedness between the DH or

495    inbred lines. However, selection should always provide high level of overall LD in the groups of

496    selected DH or inbred lines. Comparison of our LD assessment with the LD analyses from other

497    studies is inadequate because we have distances in cM and not in base-pairs. But in general the level

498    of LD was high ($r^2$ of approximately 0.3) only for SNPs separated by up to 0.5 Mb (Technow et al.

499    2014; Massman et al. 2013; Technow et al. 2012a; Riedelsheimer et al. 2012). To maximize the

500    prediction accuracy and the efficiency of identification of the best non-assessed single crosses it is

501    necessary to adopt the random sampling of single crosses for testing instead of the random sampling

23

502    of DH or inbred lines for a diallel. This is because sampling 30 or even 10% of the single crosses

503    leads to single crosses for testing derived from all DH or inbred lines from each group. In our case,

504    in every resampling assuming training set size of 30 and 10% we always get groups of assessed

505    single crosses (1,470 and 490 single crosses, respectively) derived from the 70 DH lines of each

506    group. However, sampling DH lines for a diallel provided 1,440 and 484 single crosses for testing

507    derived from 38 and 22 DH lines, respectively. Thus, the sampling of single crosses provides best

508    prediction of the SNP average effects of substitution. Riedelsheimer et al. (2012) emphasized the

509    need for large genetic variability to obtain high prediction accuracies. Further, their results indicated

510    that pairs of closely related lines and population structuring only weakly contributed to the high

511    prediction accuracies. Regarding dominance, because it can be a relevant genetic effect, breeders

512    should always fit the additive-dominance model to maximize the prediction accuracy and the

513    efficiency of identification of the best non-assessed single crosses. Interestingly, in most of the

514    studies on prediction of non-assessed single crosses the prediction accuracy did not significantly

515    increase when modeling SCA in addition to GCA effects (Zhao et al. 2015).

516        Concerning SNP density and training set size, factors related with the costs of genotyping and

517    phenotyping, breeders should find a balance between efficiency and expenses, since maximizing

518    SNP density and training set size maximizes the efficiency of untested single cross prediction.

519    Based on our results, because the decreases in the prediction accuracy (approximately 4%) and

520    coincidence index (5 to 10%) by decreasing the average SNP density from 0.1 to 1 cM are of

521    reduced magnitude, we consider sufficient to employ custom genotyping to provide an average SNP

522    density of 1 cM. Decreasing the training set size from 30 to 10% of the single crosses does not

523    significantly affect the prediction accuracy under intermediate to high heritability (decrease of up to

524    9%), but the coincidence index can be reduced in up to 21%. However, considering that the

525    coincidence index will be kept in the range 0.48 to 0.61, proportional to the heritability, and that the

526    maximum values are in the range 0.48 to 0.61, we also consider sufficient to assess at least 10% of

24

527    the possible single crosses. As highlighted by Zhao et al. (2015), marker density only marginally

528    affects the prediction accuracy of untested single crosses. For biparental populations, a plateau for

529    the accuracy is reached with a few hundred markers. Technow et al. (2014) did not find an

530    improvement of prediction accuracies by using higher SNP density. Additionally, the increase in the

531    training set size led to a relative small increase in the prediction accuracy. However, the prediction

532    accuracies obtained by Riedelsheimer et al. (2012) under high density (38,019 SNPs) were

533    substantially greater than those reached with a low-density marker panel (1,152 SNPs). In the study

534    of Technow et al. (2012a), the prediction accuracies increased with SNP density and number of

535    parents tested in hybrid combination.

536         The DH lines sampling process, the heterotic pattern, and the statistical approach should not

537    be worries for breeders. However, under high heritability notice that sampling more than one DH

538    line per $S_0$ or $S_3$ plant provided the higher coincidence values and high prediction accuracy in our

539    study. For rice, wheat, and barley breeders our message is: high prediction accuracy and high

540    efficiency of identification of superior non-assessed single crosses does not depend on heterotic

541    groups but on the (overall) LD in the group or in each group of DH or inbred lines. In other words,

542    the efficiency of prediction of non-assessed single crosses derived from DH or inbred lines from the

543    same population can be as high as the efficiency of prediction of untested single crosses derived

544    from DH or inbred lines from distinct heterotic groups. This is not confirmed comparing the relative

545    prediction accuracies for grain yield of maize untested single crosses (from approximately 0.50 to

546    0.95, for most studies) with those obtained with rice, wheat, and barley untested hybrids (0.50 to

547    0.60, approximately) (Philipp et al. 2016; Xu et al. 2014; Zhao et al. 2013b). However, the lower

548    relative prediction accuracies for untested rice, wheat, and barley hybrids should be due to

549    prediction of two- and three-way crosses. Regarding the statistical approach, our model did not

550    provide an increase in the efficiency of non-assessed single cross prediction, compared to the

551    models proposed by Massman et al. (2013) and Technow et al. (2012a). It is important to highlight

25

552  that our results showed that these two models are really identical (data no shown). Thus, because

553  the simplified definition of the incidence matrices for these two previous models, it is quite safe to

554  use any of them. Finally, the choice between the statistical approaches RR-BLUP (prediction of

555  genotypic values of non-assessed single crosses based on prediction of SNP average effects of

556  substitution), GBLUP (prediction of genotypic values of non-assessed single crosses based on

557  additive and dominance genomic matrices), and BLUP (prediction of genotypic values of non-

558  assessed single crosses based on additive and dominance matrices from pedigree records) is not a

559  serious worry for breeders too. Our evidence is that there is no significant difference between RR-

560  BLUP and GBLUP regarding prediction accuracy and efficiency of identification of the best

561  untested single crosses. Further, even when the level of relatedness between the DH or inbred lines

562  in each group is low, in general BLUP is as efficient as genomic prediction, excepting when the DH

563  lines are derived from inbred population. Thus, DNA polymorphism is not essential for an efficient

564  prediction of non-assessed single cross performance. In his review on genomic selection in hybrid

565  breeding, Zhao et al. (2015) state that the choice of the biometrical model has no substantial impact

566  on the prediction accuracy of untested single crosses. Technow et al. (2014) observed that

567  prediction methods GBLUP and BayesB resulted in very similar prediction accuracies. In the study

568  of Massman et al. (2013), BLUP and RR-BLUP models did not lead to prediction accuracies that

569  differed significantly. Comparing GBLUP and BayesB, Technow et al. (2012a) concluded that the

570  latter method produced significantly higher accuracies for the additive-dominance model.

571  Our main contributions on the prediction efficiency of non-assessed single cross performance

572  are: 1) the prediction accuracy of untested single crosses ranged from approximately 0.80 to 0.90 as

573  the heritability of tested single crosses ranged from low (30%) to high (100%); however, the

574  efficacy of identification of the best 9% of the untested single crosses ranged from approximately

575  0.50 to 0.70, depending on the DH lines sampling process; 2) the prediction accuracy for crops

576  showing no defined heterotic pattern can be as efficient as with maize, for which there is well

26

577   defined heterotic groups; this is because the most important factor affecting the prediction

578   efficiency is the overall LD; 3) to maximize prediction accuracy and coincidence the choice of

579   single crosses for testing should be based on a random process; this procedure maximizes the

580   number of DH lines in hybrid combination and provides better predictions of the SNP average

581   effects of substitution and dominance deviations; 4) because non significant decreases in the

582   prediction accuracy and coincidence, the prediction of untested single crosses can be efficient

583   assuming reduced training set size (10%) and SNP density of 1 cM; 5) RR-BLUP and GBLUP

584   provides equivalent prediction efficiencies of untested single crosses; 6) excepting for DH lines

585   derived from inbred populations, BLUP is as efficient as genomic prediction of untested single

586   crosses; and 7) the theoretical accuracy shows that the prediction accuracy is not affected by the

587   linkage phase.

588                                              **ACKNOWLEDGMENTS**

592                                              **LITERATURE CITED**

593   Albrecht, T., H.-J. Auinger, V. Wimmer, J.O. Ogutu, C. Knaak *et al.*, 2014 Genome-based

594           prediction of maize hybrid performance across genetic groups, testers, locations, and years.

595           *Theoretical and Applied Genetics* 127 (6):1375-1386.

596   Albrecht, T., V. Wimmer, H.-J. Auinger, M. Erbe, C. Knaak *et al.*, 2011 Genome-based prediction

597           of testcross values in maize. *Theoretical and Applied Genetics* 123 (2):339-350.

598   Azevedo, C.F., M.D. Vilela de Resende, F. Fonseca e Silva, J.M. Soriano Viana, M.S. Ferreira

599           Valente *et al.*, 2015 Ridge, Lasso and Bayesian additive-dominance genomic models. *BMC*

600           *Genet* 16.

601    Barrett, J.C., B. Fry, J. Maller, and M.J. Daly, 2005 Haploview: analysis and visualization of LD
602           and haplotype maps. *Bioinformatics* 21 (2):263-265.

603    Crossa, J., P. Perez, J. Hickey, J. Burgueno, L. Ornella *et al.*, 2014 Genomic prediction in
604           CIMMYT maize and wheat breeding programs. *Heredity (Edinb)* 112 (1):48-60.

605    Daetwyler, H.D., M.P.L. Calus, R. Pong-Wong, G. de los Campos, and J.M. Hickey, 2013 Genomic
606           Prediction in Animals and Plants: Simulation of Data, Validation, Reporting, and
607           Benchmarking. *Genetics* 193 (2):347-+.

608    de Los Campos, G., J.M. Hickey, R. Pong-Wong, H.D. Daetwyler, and M.P. Calus, 2013 Whole-
609           genome regression and prediction methods applied to plant and animal breeding. *Genetics*
610           193 (2):327-345.

611    Endelman, J.B., 2011 Ridge Regression and Other Kernels for Genomic Selection with R Package
612           rrBLUP. *Plant Genome* 4 (3):250-255.

613    Jonas, E., and D.J. de Koning, 2013 Does genomic selection have a future in plant breeding? *Trends*
614           *in Biotechnology* 31 (9):497-504.

615    Kadam, D.C., S.M. Potts, M.O. Bohn, A.E. Lipka, and A.J. Lorenz, 2016 Genomic Prediction of
616           Single Crosses in the Early Stages of a Maize Hybrid Breeding Pipeline. *G3-Genes*
617           *Genomes Genetics* 6 (11):3443-3453.

618    Li, Z., N. Philipp, M. Spiller, G. Stiewe, J.C. Reif *et al.*, 2017 Genome-Wide Prediction of the
619           Performance of Three-Way Hybrids in Barley. *Plant Genome* 10 (1).

620    Massman, J.M., A. Gordillo, R.E. Lorenzana, and R. Bernardo, 2013 Genomewide predictions from
621           maize single-cross data. *Theor Appl Genet* 126 (1):13-22.

622    Meuwissen, T., B. Hayes, and M. Goddard, 2013 Accelerating Improvement of Livestock with
623           Genomic Selection. *Annual Review of Animal Biosciences, Vol 1* 1:221-237.

624    Philipp, N., G.Z. Liu, Y.S. Zhao, S. He, M. Spiller *et al.*, 2016 Genomic Prediction of Barley
625           Hybrid Performance. *Plant Genome* 9 (2).

626  Riedelsheimer, C., A. Czedik-Eysenberg, C. Grieder, J. Lisec, F. Technow *et al.*, 2012 Genomic
627      and metabolic prediction of complex heterotic traits in hybrid maize. *Nature Genetics* 44
628      (2):217-220.

629  Technow, F., C. Riedelsheimer, T.A. Schrag, and A.E. Melchinger, 2012a Genomic prediction of
630      hybrid performance in maize with models incorporating dominance and population specific
631      marker effects. *Theoretical and Applied Genetics* 125 (6):1181-1194.

632  Technow, F., C. Riedelsheimer, T.A. Schrag, and A.E. Melchinger, 2012b Genomic prediction of
633      hybrid performance in maize with models incorporating dominance and population specific
634      marker effects. *Theor Appl Genet* 125 (6):1181-1194.

635  Technow, F., T.A. Schrag, W. Schipprack, E. Bauer, H. Simianer *et al.*, 2014 Genome Properties
636      and Prospects of Genomic Prediction of Hybrid Performance in a Breeding Program of
637      Maize. *Genetics* 197 (4):1343-U1469.

638  Van Eenennaam, A.L., K.A. Weigel, A.E. Young, M.A. Cleveland, and J.C.M. Dekkers, 2014
639      Applied Animal Genomics: Results from the Field. *Annual Review of Animal Biosciences,*
640      *Vol 2* 2:105-139.

641  Viana, J.M.S., H.-P. Piepho, and F.F. Silva, 2016 Quantitative genetics theory for genomic
642      selection and efficiency of breeding value prediction in open-pollinated populations.
643      *Scientia Agricola* 73 (3):243-251.

644  Viana, J.M.S., H.P. Piepho, and F.F. Silva, 2017a Quantitative genetics theory for genomic
645      selection and efficiency of genotypic value prediction in open-pollinated populations.
646      *Scientia Agricola* 74 (1):41-50.

647  Viana, J.M.S., F.F. Silva, G.B. Mundim, C.F. Azevedo, and H.U. Jan, 2017b Efficiency of low
648      heritability QTL mapping under high SNP density. *Euphytica* 213 (1).

649  Viana, J.M.S., M.S.F. Valente, F.F. Silva, G.B. Mundim, and G.P. Paes, 2013 Efficacy of

650      population structure analysis with breeding populations and inbred lines. *Genetica* 141 (7-

651      9):389-399.

652  Windhausen, V.S., G.N. Atlin, J.M. Hickey, J. Crossa, J.-L. Jannink *et al.*, 2012 Effectiveness of

653      Genomic Prediction of Maize Hybrid Performance in Different Breeding Populations and

654      Environments. *G3-Genes Genomes Genetics* 2 (11):1427-1436.

655  Xu, S., D. Zhu, and Q. Zhang, 2014 Predicting hybrid performance in rice using genomic best linear

656      unbiased prediction. *Proceedings of the National Academy of Sciences of the United States

657      of America* 111 (34):12456-12461.

658  Zhao, Y., M. Gowda, W. Liu, T. Wuerschum, H.P. Maurer *et al.*, 2013a Choice of shrinkage

659      parameter and prediction of genomic breeding values in elite maize breeding populations.

660      *Plant Breeding* 132 (1):99-106.

661  Zhao, Y., M.F. Mette, and J.C. Reif, 2015 Genomic selection in hybrid breeding. *Plant Breeding*

662      134 (1):1-10.

663  Zhao, Y., J. Zeng, R. Fernando, and J.C. Reif, 2013b Genomic Prediction of Hybrid Wheat

664      Performance. *Crop Science* 53 (3):802.

665

666 **Table 1** Average prediction accuracies of non-assessed single crosses and its standard deviation,

667 assuming single crosses from selected DH lines, 30 and 10% of assessed single crosses, two traits

668 (grain yield - GY, g/plant, and expansion volume - EV, mL/g), two sampling processes of single

669 crosses, four statistical models, three DH lines sampling processes, two genetic models, and three

670 accuracies of assessed single crosses

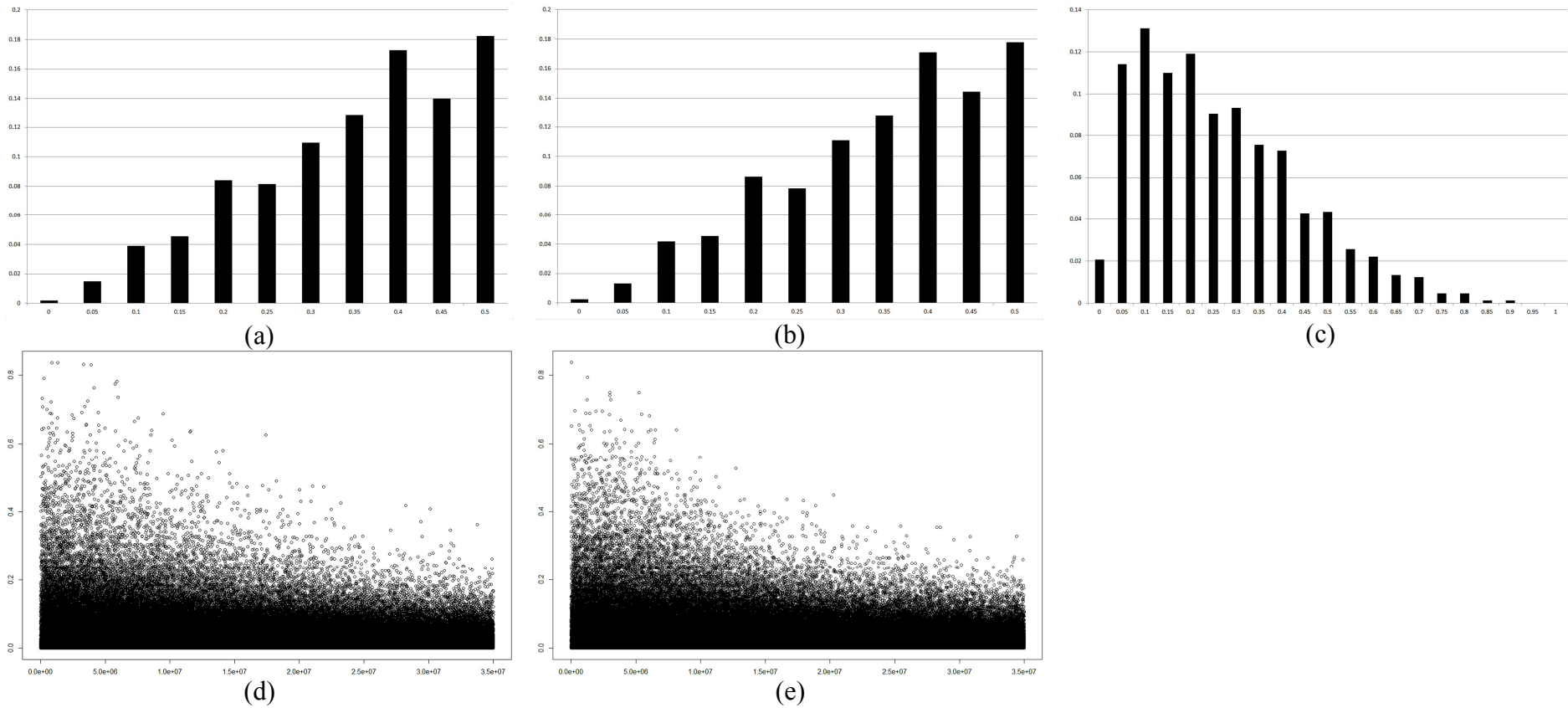| Trait | Samp. proc. | Statistical model | DH lines | Gen. mod. | Accuracy of assessed single crosses | | |
|---|---|---|---|---|---|---|---|
| | | | | | 0.55 | 0.77 | 1.00 |
| GY | SCs | Viana et al. | $1/S_0$ | AD | $0.7790 \pm 0.0124$ | $0.8447 \pm 0.0066$ | $0.8859 \pm 0.0018$ |
| | | | | A | $0.7688 \pm 0.0132$ | $0.8380 \pm 0.0067$ | $0.8821 \pm 0.0019$ |
| | | | $1\text{-}5/S_0$ | AD | $0.7947 \pm 0.0125$ | $0.8525 \pm 0.0072$ | $0.8896 \pm 0.0025$ |
| | | | | A | $0.7895 \pm 0.0126$ | $0.8465 \pm 0.0077$ | $0.8858 \pm 0.0027$ |
| | | | $1\text{-}5/S_3$ | AD | $0.8010 \pm 0.0145$ | $0.8678 \pm 0.0054$ | $0.9276 \pm 0.0025$ |
| | | | | A | $0.7954 \pm 0.0145$ | $0.8627 \pm 0.0056$ | $0.9238 \pm 0.0026$ |
| | | | $1\text{-}5/S_3$ | AD[a] | $0.7718 \pm 0.0161$ | $0.8371 \pm 0.0079$ | $0.8888 \pm 0.0043$ |
| | | | $1\text{-}5/S_3$ | AD[b] | $0.6836 \pm 0.0277$ | $0.7885 \pm 0.0139$ | $0.8817 \pm 0.0049$ |
| | | | $1/S_0$ | AD[c] | $0.8293 \pm 0.0131$ | $0.8944 \pm 0.0049$ | $0.9479 \pm 0.0017$ |
| | | | $1\text{-}5/S_3$ | AD[d] | $0.8267 \pm 0.0082$ | $0.8928 \pm 0.0043$ | $0.9083 \pm 0.0023$ |
| | | Massman et. al.[e] | $1/S_0$ | AD | $0.7874 \pm 0.0118$ | $0.8519 \pm 0.0053$ | $0.8924 \pm 0.0026$ |
| | | | $1\text{-}5/S_0$ | AD | $0.7982 \pm 0.0140$ | $0.8622 \pm 0.0055$ | $0.8973 \pm 0.0025$ |
| | | | $1\text{-}5/S_3$ | AD | $0.8074 \pm 0.0112$ | $0.8753 \pm 0.0056$ | $0.9314 \pm 0.0026$ |
| | | GBLUP | $1/S_0$ | AD | $0.7841 \pm 0.0122$ | $0.8477 \pm 0.0064$ | $0.8906 \pm 0.0019$ |
| | | | $1\text{-}5/S_0$ | AD | $0.7973 \pm 0.0124$ | $0.8574 \pm 0.0070$ | $0.8978 \pm 0.0019$ |
| | | | $1\text{-}5/S_3$ | AD | $0.7911 \pm 0.0146$ | $0.8639 \pm 0.0056$ | $0.9319 \pm 0.0023$ |
| | | BLUP | $1/S_0$ | AD | $0.7855 \pm 0.0129$ | $0.8541 \pm 0.0059$ | $0.8899 \pm 0.0019$ |
| | | | $1\text{-}5/S_0$ | AD | $0.7803 \pm 0.0143$ | $0.8435 \pm 0.0074$ | $0.8830 \pm 0.0024$ |
| | | | $1\text{-}5/S_3$ | AD | $0.7227 \pm 0.0203$ | $0.7915 \pm 0.0077$ | $0.8373 \pm 0.0048$ |
| | DHs | Viana et al. | $1/S_0$ | AD | $0.5012 \pm 0.0416$ | $0.5117 \pm 0.0467$ | $0.5343 \pm 0.0467$ |
| | | | $1\text{-}5/S_0$ | AD | $0.4827 \pm 0.0423$ | $0.5000 \pm 0.0420$ | $0.5036 \pm 0.0465$ |
| | | | $1\text{-}5/S_3$ | AD | $0.5799 \pm 0.0437$ | $0.6106 \pm 0.0413$ | $0.6357 \pm 0.0429$ |
| EV | SCs | Viana et al. | $1/S_0$ | AD | $0.7779 \pm 0.0157$ | $0.8458 \pm 0.0069$ | $0.8820 \pm 0.0024$ |
| | | | $1\text{-}5/S_0$ | AD | $0.8019 \pm 0.0155$ | $0.8656 \pm 0.0050$ | $0.9055 \pm 0.0020$ |
| | | | $1\text{-}5/S_3$ | AD | $0.7589 \pm 0.0143$ | $0.8424 \pm 0.0058$ | $0.9165 \pm 0.0027$ |

[a]density of 1 cM; [b]training set of 490 single crosses (10%); [c]after 10 generations of random crosses; [d]single crosses from DH lines of the same population; [e]and Technow et al..

671 **Table 2** Average coincidence of the best 300 predicted single crosses and its standard deviation,

672 assuming single crosses from selected DH lines, 30 and 10% of assessed single crosses, two traits

673 (grain yield - GY, g/plant, and expansion volume - EV, mL/g), two sampling processes of single

674 crosses, four statistical models, three DH lines sampling processes, two genetic models, and three
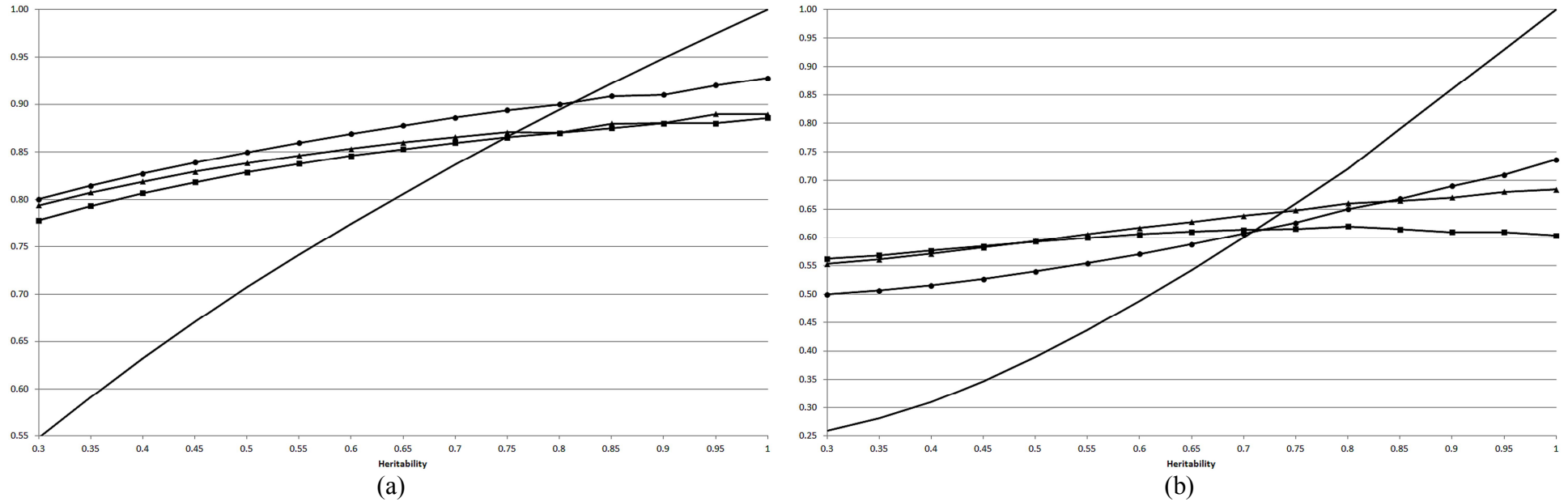
675 parametric coincidence of assessed single crosses

| Trait | Samp. proc. | Statistical model | DH lines | Gen. mod. | Coincidence of assessed single crosses | | |
|---|---|---|---|---|---|---|---|
| | | | | | 0.26 | 0.48 | 1.00 |
| GY | SCs | Viana et al. | $1/S_0$ | AD | $0.4523 \pm 0.0334$ | $0.5525 \pm 0.0190$ | $0.6037 \pm 0.0170$ |
| | | | | A | $0.4396 \pm 0.0346$ | $0.5449 \pm 0.0176$ | $0.5976 \pm 0.0172$ |
| | | | $1\text{-}5/S_0$ | AD | $0.5686 \pm 0.0273$ | $0.6369 \pm 0.0221$ | $0.6842 \pm 0.0140$ |
| | | | | A | $0.5640 \pm 0.0283$ | $0.6299 \pm 0.0221$ | $0.6816 \pm 0.0152$ |
| | | | $1\text{-}5/S_3$ | AD | $0.5129 \pm 0.0235$ | $0.6044 \pm 0.0200$ | $0.7363 \pm 0.0183$ |
| | | | | A | $0.5063 \pm 0.0225$ | $0.5993 \pm 0.0193$ | $0.7305 \pm 0.0190$ |
| | | | $1\text{-}5/S_3$ | AD[a] | $0.4881 \pm 0.0278$ | $0.5691 \pm 0.0229$ | $0.6620 \pm 0.0215$ |
| | | | $1\text{-}5/S_3$ | AD[b] | $0.3805 \pm 0.0511$ | $0.4797 \pm 0.0354$ | $0.6087 \pm 0.0233$ |
| | | | $1/S_0$ | AD[c] | $0.5528 \pm 0.0298$ | $0.6489 \pm 0.0203$ | $0.7571 \pm 0.0162$ |
| | | | $1\text{-}5/S_3$ | AD[d] | $0.6116 \pm 0.0214$ | $0.7156 \pm 0.0150$ | $0.7581 \pm 0.0166$ |
| | | Massman et. al.[e] | $1/S_0$ | AD | $0.4670 \pm 0.0346$ | $0.5663 \pm 0.0174$ | $0.6157 \pm 0.0157$ |
| | | | $1\text{-}5/S_0$ | AD | $0.5651 \pm 0.0310$ | $0.6431 \pm 0.0164$ | $0.6955 \pm 0.0144$ |
| | | | $1\text{-}5/S_3$ | AD | $0.5279 \pm 0.0291$ | $0.6139 \pm 0.0204$ | $0.7423 \pm 0.0172$ |
| | | GBLUP | $1/S_0$ | AD | $0.4622 \pm 0.0308$ | $0.5660 \pm 0.0190$ | $0.6092 \pm 0.0163$ |
| | | | $1\text{-}5/S_0$ | AD | $0.5650 \pm 0.0280$ | $0.6384 \pm 0.0204$ | $0.6849 \pm 0.0137$ |
| | | | $1\text{-}5/S_3$ | AD | $0.5010 \pm 0.0245$ | $0.5937 \pm 0.0216$ | $0.7294 \pm 0.0168$ |
| | | BLUP | $1/S_0$ | AD | $0.4641 \pm 0.0331$ | $0.5709 \pm 0.0176$ | $0.6081 \pm 0.0127$ |
| | | | $1\text{-}5/S_0$ | AD | $0.5531 \pm 0.0323$ | $0.6272 \pm 0.0194$ | $0.6699 \pm 0.0130$ |
| | | | $1\text{-}5/S_3$ | AD | $0.4172 \pm 0.0258$ | $0.4731 \pm 0.0211$ | $0.5377 \pm 0.0196$ |
| | DHs | Viana et al. | $1/S_0$ | AD | $0.2753 \pm 0.0374$ | $0.3056 \pm 0.0445$ | $0.3169 \pm 0.0401$ |
| | | | $1\text{-}5/S_0$ | AD | $0.3268 \pm 0.0642$ | $0.3400 \pm 0.0691$ | $0.3461 \pm 0.0728$ |
| | | | $1\text{-}5/S_3$ | AD | $0.3699 \pm 0.0583$ | $0.3931 \pm 0.0579$ | $0.4300 \pm 0.0633$ |
| EV | SCs | Viana et al. | $1/S_0$ | AD | $0.5156 \pm 0.0331$ | $0.6081 \pm 0.0159$ | $0.6599 \pm 0.0146$ |
| | | | $1\text{-}5/S_0$ | AD | $0.5506 \pm 0.0285$ | $0.6337 \pm 0.0203$ | $0.6944 \pm 0.0141$ |
| | | | $1\text{-}5/S_3$ | AD | $0.4746 \pm 0.0294$ | $0.5843 \pm 0.0174$ | $0.7141 \pm 0.0171$ |

[a]density of 1 cM; [b]training set of 490 single crosses (10%); [c]after 10 generations of random crosses; [d]single crosses from DH lines of the same population; [e]and Technow et al..

676    **Figure 1** Frequency distribution of the MAF in the groups of selected DH lines (a and b) and the absolute value of the difference between a SNP allele

677    frequency (c), and LD ($r^2$) in relation to distance (cM) in the two groups of selected DH lines (d and e), regarding SNPs in chromosome 1 separated by

678    zero to 35 cM, assuming one DH line per $S_0$ plant.

**Figure 2** Predicted accuracies (a) and coincidence indexes (b) for untested single crosses (square: $1/S_0$; triangle: $1-5/S_0$; circle: $1-5/S_3$), and parametric

680   accuracies and coincidence indexes for tested single crosses (continuous line), assuming our model, average SNP density of 0.1 cM, training set size of

681   30%, positive dominance (grain yield), additive-dominance model, and sampling of single crosses.