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, Indianápolis, 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.

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 a definitive proof that prediction of non-assessed single

18    crosses (SCs) is efficient. We provided a new genetic model for genomic prediction. The SNP and

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

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

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

22    included three sampling processes of DH lines, two sampling processes of SCs for testing, two SNP

23    densities, DH lines from the same population, DH lines from populations with lower LD, two

24    genetic models, three statistical models, and three statistical approaches. The efficiency of

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

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

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

28    untested SCs is very efficient. The accuracies and coincidences ranged from approximately 0.80

29    and 0.50, respectively, under low heritability, to 0.90 and 0.7, assuming high heritability.

30    Additionally, we highlighted the relevance of the overall LD and evidenced that efficient prediction

31    of untested SCs can be achieved for crops that show no heterotic pattern, for reduced training size

32    set (10%), for SNP density of 1 cM, and for distinct sampling processes of DH lines, based on

33    random choice of the SCs for testing.

## INTRODUCTION

35    The genomic selection is a reality in animal breeding, especially for dairy cattle (Van

36    Eenennaam et al. 2014). The same cannot yet be said concerning crop breeding, with exceptions.

37  The main reasons for the effective application of genomic selection in livestock breeding are: it is

38  efficient, that is, the process has high prediction accuracy, the cost of phenotyping (mainly progeny

39  test) is higher than the cost of genotyping, and the process significantly shorten the selection cycle

40  (Meuwissen et al. 2013). It is worth to remember that prediction of breeding and genotypic values is

41  not exclusive for genomic selection, having been pioneered by the best linear unbiased prediction

42  method (BLUP) (Henderson 1974). In spite of the many field and simulation-based studies with

43  genomic selection in plant breeding, in general the cost of phenotyping is much lower than the cost

44  of genotyping, restricting its application in breeding programs. Jonas and de Koning (2013)

45  consider that genomic selection has the potential to improve existing plant breeding schemes.

46  However, based also on the high diversity and complexity of plant breeding methods, they stated

47  that there are great obstacles to overcome.

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

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

50  prediction) in hybrid breeding (Zhao et al. 2015). The prediction of untested single crosses was

51  pioneered by Bernardo (1994), also based on BLUP. Many significant studies on prediction of

52  untested single cross and testcross performance have been published in the last 23 years, focused on

53  the assessment of the prediction accuracy. Most investigations were based on empirical data and

54  estimated the prediction accuracy using a cross-validation procedure. Very few were based on

55  simulated data (Li et al. 2017; Technow et al. 2012a). With no exception, the inference was that

56  prediction of untested single crosses and testcrosses is an efficient process, proportional to

57  heritability, training set size, and number of tested inbreds in hybrid combination (both, one, and

58  none parents tested). It is impressive that this inference have been stated from studies differing for

59  molecular markers, density of markers, number of inbreds, level of relatedness, diversity and

60  linkage disequilibrium (LD) between inbreds, heterotic patterns, training set size, genetic model,

61  and statistical approach (Zhao et al. 2015).

62      Most papers on genomic prediction of maize single cross performance published since 2011

63      have employed single nucleotide polymorphism (SNP), with the SNP number in the range 425

64      (Zhao et al. 2013a) to 39,627 (Technow et al. 2012a). Based on the physical length of the maize

65      genome (approximately 2,000 megabase pairs (Mb) according to Maize genetics and genomics

66      database), the density ranged from approximately 5 to 0.05 Mb, respectively. For grain yield, the

67      relative prediction accuracies (accuracy/root square of the heritability) in these two papers ranged

68      from 0.27 to 0.62 and from 0.65 to 0.95, respectively. The number of inbreds in each heterotic

69      group was highly variable too, ranging from six and nine (Bernardo 1994) to 75 and 75 (Technow et

70      al. 2012a). The relative accuracy observed by Bernardo (1994) ranged between 0.72 and 0.89. The

71      number of testcrosses ranged between 255 (Windhausen et al. 2012) and 1,894 (Albrecht et al.

72      2014). The relative accuracies ranged from 0.46 to 0.52 and from 0.33 to 0.65, respectively. The

73      level of relatedness ranged from non-related inbreds in each group (Technow et al. 2012a) to an

74      maximum average value of 0.58 (Bernardo 1995). The relative accuracy obtained by Bernardo

75      (1995) ranged from 0.41 to 0.80. The common heterotic groups were Stiff Stalk and non-Stiff Stalk

76      (Kadam et al. 1916) or Dent and Flint (Technow et al. 2014). The study of Bernardo (1996a)

77      involved nine heterotic groups and the (statistically significant) relative accuracies ranged from 0.43

78      to 0.88. No study provided clearly greater prediction accuracy of the additive-dominance model

79      relative to the additive model. Finally, only with testcrosses the genomic BLUP (GBLUP) approach

80      outperformed BLUP (Albrecht et al. 2014; Albrecht et al. 2011) concerning prediction accuracy.

81      After so many years of research on prediction of untested single crosses, with consistent

82      results from reduced and large data sets, it is was a challenge to plan a study that could provide a

83      new and significant contribution on efficiency of prediction of untested single cross performance.

84      We believe have achieved our purpose. For the first time, our simulation study has provided for

85      breeders a direct measure of efficiency of identification of the best 300 of the really non-assessed

86      single crosses, additionally to the standard prediction accuracy (coincidence index). These measures

87      of efficacy were provided for a large data set (4,900 single crosses) and for low (30%) to high

88   heritability (100%), assuming scenarios not favorable to prediction of non-assessed single cross

89   performance, as low level of relatedness and a not high heterotic pattern. Additionally, we provided

90   a new genetic model for genomic prediction, supported by quantitative genetics theory, highlighted

91   the relevance of the overall LD (not only for linked SNPs and QTLs), and evidenced that efficient

92   prediction of untested single crosses can be achieved for crops that show no clear heterotic pattern,

93   as rice, wheat, and barley, for reduced training size set (10%), for SNP density of 1 cM, and for

94   distinct processes of doubled haploid (DH) lines sampling. Finally, we showed that the choice of

95   the single crosses for testing must be based on a random process, but never by sampling DH or

96   inbreds lines for a diallel. Thus, our objective was to provide to breeders a definitive proof that

97   prediction of non-assessed single crosses can be efficient and that they should make widespread use

98   of this procedure for identification of the best hybrids, prior to field testing.

## MATERIALS AND METHODS

100  **Theory**

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

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

103  also a quantitative trait locus (QTL) (alleles B/b) and a SNP (alleles C/c) where B and b are the

104  alleles that increase and decrease the trait expression, respectively. Define the joint genotype

105  probabilities (equal to the joint haplotype probabilities) as $P(BBCC) = f_{22}$, $P(BBcc) = f_{20}$,

106  $P(bbCC) = f_{02}$, and $P(bbcc) = f_{00}$, where the subscript indicates the number of copies of the

107  major allele (B and C). The measure of LD between the QTL and the SNP is

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

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

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

111  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

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

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

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

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

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

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

118 ***SNP genotypic values of DH or inbred lines***

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

120 $M_{IL} = m_b + \left(p_b - q_b\right)a_b$, where $m_b$ is the mean of the genotypic values of the homozygotes and

121 $a_b$ is the deviation between the genotypic value of the homozygote of higher expression and $m_b$.

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

123 $$G_{CC} = \frac{1}{f_{.2}}\left[f_{22}\left(m_b + a_b\right) + f_{02}\left(m_b - a_b\right)\right] = M_{IL} + 2q_c\alpha_{SNP} = M_{IL} + A_{CC}$$

124 $$G_{cc} = \frac{1}{f_{.0}}\left[f_{20}\left(m_b + a_b\right) + f_{00}\left(m_b - a_b\right)\right] = M_{IL} - 2p_c\alpha_{SNP} = M_{IL} + A_{cc}$$

125 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

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

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

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

129 $\kappa_{ce} = \left[\dfrac{\Delta_{ce}}{p_c q_c}\right]$. Thus, in general, the average effect of a SNP substitution (and the SNP additive

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

131 the marker.

132  *SNP genotypic values of single crosses*

133      Aiming to maximize the heterosis, maize breeders commonly assess single crosses originated

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

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

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

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

138
$$M_H = m_b + \left( p_{b1}p_{b2} - q_{b1}q_{b2} \right) a_b + \left( p_{b1}q_{b2} + q_{b1}p_{b2} \right) d_b$$

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

140  heterozygote and $m_b$ ).

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

142  cc of the group 1 are

143
$$M_{CC1} = M_H + q_{c1}\kappa_{bc1}\left[ a_b + \left( q_{b2} - p_{b2} \right) d_b \right] = M_H + q_{c1}\kappa_{bc1}\alpha_{b2} = M_H + q_{c1}\alpha_{SNP1}$$
$$= M_H + GCA_{CC1}$$

144
$$M_{cc1} = M_H - p_{c1}\kappa_{bc1}\alpha_{b2} = M_H - p_{c1}\alpha_{SNP1} = M_H + GCA_{cc1}$$

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

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

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

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

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

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

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

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

153
$$M_{CC2} = M_H + q_{c2}\kappa_{bc2}\left[a_b + \left(q_{b1} - p_{b1}\right)d_b\right] = M_H + q_{c2}\kappa_{bc2}\alpha_{b1} = M_H + q_{c2}\alpha_{SNP2}$$
$$= M_H + GCA_{CC2}$$

154
$$M_{cc2} = M_H - p_{c2}\kappa_{bc2}\alpha_{b1} = M_H - p_{c2}\alpha_{SNP2} = M_H + GCA_{cc2}$$

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

156     the SNP locus are

157
$$M_{CC1xCC2} = M_H + q_{c1}\alpha_{SNP1} + q_{c2}\alpha_{SNP2} - 2q_{c1}q_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{CC1} + GCA_{CC2} + SCA_{CC1xCC2}$$

158
$$M_{cc1xcc2} = M_H - p_{c1}\alpha_{SNP1} - p_{c2}\alpha_{SNP2} - 2p_{c1}p_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{cc1} + GCA_{cc2} + SCA_{cc1xcc2}$$

159
$$M_{CC1xcc2} = M_H + q_{c1}\alpha_{SNP1} - p_{c2}\alpha_{SNP2} + 2q_{c1}p_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{CC1} + GCA_{cc2} + SCA_{CC1xcc2}$$

160
$$M_{cc1xCC2} = M_H - p_{c1}\alpha_{SNP1} + q_{c2}\alpha_{SNP2} + 2p_{c1}q_{c2}\kappa_{bc1}\kappa_{bc2}d_b$$
$$= M_H + GCA_{cc1} + GCA_{CC2} + SCA_{cc1xCC2}$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

176 $M = M_H + u_1 a_1 + u_2 a_2 + u_3 d$, where $u_1$ and $u_2$ equal to 1/2 or −1/2 if the corresponding DH or

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

178 is homozygous or 1 if heterozygous.

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

180 ***population or heterotic group***

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

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

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

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

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

186 single crosses concerning the SNP locus are

187 $M_{CCxCC} = M + 2q_c\alpha_{SNP} - 2q_c^2\kappa_{bc}^2 d_b = M + 2GCA_{CC} + SCA_{CCxCC}$

188 $M_{ccxcc} = M - 2p_c\alpha_{SNP} - 2p_c^2\kappa_{bc}^2 d_b = M + 2GCA_{cc} + SCA_{ccxcc}$

189 $M_{CCxcc} = M + 2(q_c - p_c)\alpha_{SNP} + 2p_c q_c\kappa_{bc}^2 d_b = M + GCA_{CC} + GCA_{cc} + SCA_{CCxcc}$

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

191 $\alpha_{SNP} = \kappa_{bc}[a_b + (q_b - p_b)d_b] = \kappa_{bc}\alpha_b$ is the average effect of a SNP substitution in the hybrid

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

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

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

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

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

### *Accuracy of single cross genomic prediction*

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

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

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

201

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

202

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

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

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

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

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

208 $\qquad \sigma_{GCA}^{2(2)} = p_{b2}q_{b2}\left(\alpha_{b1}\right)^2$, and $\sigma_{SCA}^2 = 4p_{b1}q_{b1}p_{b2}q_{b2}\left(d_b\right)^2$. The QTL genotypic variance is

209 $\qquad \sigma_G^2 = \sigma_{GCA}^{2(1)} + \sigma_{GCA}^{2(2)} + \sigma_{SCA}^2$ Thus, the single cross prediction accuracy is

210 $\qquad \rho_{\widetilde{G},G} = \sqrt{\dfrac{\sigma_{G(SNP)}^2}{\sigma_G^2}}$

211 $\qquad$ Assuming s SNPs,

212 $\qquad \rho_{\widetilde{G},G} = \sum\limits_{r=1}^{s} \sigma_{G(SNP(r))}^2 \Big/ \sqrt{\sigma_{\widetilde{G}}^2 \sigma_G^2}$

213 $\qquad$ where $\sigma_{\widetilde{G}}^2$ is the variance of the predicted single cross genotypic values and $\sigma_G^2$ is the single cross

214 $\qquad$ genotypic variance. Further,

215 $\qquad \alpha_{SNP(r)1} = \sum\limits_{i=1}^{k'} \left[\dfrac{\Delta_{ri1}}{p_{r1}q_{r1}}\right]\alpha_{i2} = \sum\limits_{i=1}^{k'} \kappa_{ri1}\alpha_{i2}$, where k' is the number of QTLs in LD with the SNP

216 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ r) in group 1, and

217 $\qquad d_{SNP(r)} = \sum\limits_{i=1}^{k''} \left[\dfrac{\Delta_{ri1}}{p_{r1}q_{r1}}\right]\left[\dfrac{\Delta_{ri2}}{p_{r2}q_{r2}}\right]d_i = \sum\limits_{i=1}^{k''} \kappa_{ri1}\kappa_{ri2}d_i$ where k" is the number of QTLs in LD with

218 $\qquad\qquad\qquad\qquad\qquad\qquad$ the SNP r in both groups

219 $\qquad$ Notice that because the accuracy of genomic prediction of single crosses depends on the

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

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

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

223 $\qquad$ Assuming single crosses derived from DH or inbred lines of a single population or heterotic

224 group $\qquad$ we $\qquad$ have $\qquad \sigma_{G(SNP)}^2 = 2p_c q_c \left(\alpha_{SNP}\right)^2 + \left(2p_c q_c d_{SNP}\right)^2 \qquad$ and

225 $\sigma_G^2 = 2p_b q_b \left(\alpha_b\right)^2 + \left(2p_b q_b d_b\right)^2$. Therefore, the prediction accuracy of single crosses derived

226 from DH or inbred lines from two distinct populations or heterotic groups differ from the prediction

227 accuracy of single crosses resulting from DH or inbred lines obtained from each population or

228 heterotic group.

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

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

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

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

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

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

235 $$y = M_H + \sum_{r=1}^{s} \left( z_{1_r} \alpha_{SNP1_r} + z_{2_r} \alpha_{SNP2_r} + z_{3_r} d_{SNP_r} \right) + error$$

236 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

237 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

238 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

239 $z_{3_r} = 2q_{r1}p_{r2}$ if the SNP genotypes for the DH or inbred lines are CC (group 1) and cc (group 2),

240 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

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

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

243 heterotic group we have

244 $$y = M + \sum_{r=1}^{s} \left( z_{1_r} \alpha_{SNP_r} + z_{2_r} d_{SNP_r} \right) + error$$

245    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,

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

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

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

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

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

251    Campos et al. 2013). Then, the predicted effects of SNP substitution and SNP dominance deviations

252    must be used to provide genomic prediction of non-assessed single crosses. The predicted genotypic

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

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

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

256    genotypic value is

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

258    **Simulation**

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

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

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

262    *2009* (Viana et al. 2017a; Viana et al. 2017b; Viana et al. 2016; Azevedo et al. 2015; Viana et al.

263    2013). Based on our input, the software distributed 10,000 SNPs and 400 QTLs in ten

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

265    centiMorgan (cM). The QTLs were distributed in the regions covered by the SNPs (approximately

266    100 cM/chromosome). Initially, *REALbreeding* sampled 700 DH lines from two non-inbred

267    populations (heterotic groups) in LD (350 from each population). The populations were composites

268    of two populations in linkage equilibrium. In a composite, there is LD only for linked SNPs and

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

287    490 assessed single crosses). That is, we predicted 70 and 90% of the single crosses (3,430 and

288    4,410 non-assessed single crosses). Additionally, to assess the relevance of the number of DH lines

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

304    (including only the GCA effects). We also processed single cross prediction based on GBLUP and

305    BLUP.

**Statistical analysis**

307        The methods used for prediction were ridge regression BLUP (RR-BLUP), GBLUP (with the

308    observed additive and dominance relationship matrices) and BLUP (with the expected additive and

309    dominance relationship matrices). For the analyses we used the *rrBLUP* package (Endelman 2011).

310    The accuracies of single cross genotypic value prediction were obtained by the correlation between

311    the true values of the non-assessed single crosses computed by *REALbreeding* and the values

312    predicted by RR-BLUP, GBLUP, and BLUP. We also computed the efficiency of identification of

313    the 300 non-assessed single crosses of higher genotypic value (coincidence index). The parametric

314    average coincidence index was computed by ordering the average phenotypic values of the 4,900

315    single crosses for each heritability and for each DH lines derivation process. Regarding grain yield,

316    for heritability of 30% the coincidence index was 0.2533, 0.2833, and 0.2433 assuming one DH line

317    per $S_0$ plant, one to five DH lines per $S_0$ plant, and one to five DH lines per $S_3$ plant, respectively.

318    The corresponding values for heritability of 60% were, respectively, 0.4800, 0.4900, and 0.4567.

319    Concerning expansion volume, the corresponding values for heritabilities of 30 and 60% were,

320 respectively, 0.2600, 0.2833, and 0.2700, and 0.4733, 0.5100, and 0.4533. The assumed average

321 parametric coefficient index was 0.26 and 0.48 for heritabilities of 30 and 60%, respectively, for

322 both traits. For the population structure analysis we employed *Structure* (Falush et al. 2003) and

323 fitted the no admixture model with independent allelic frequencies. The number of SNPs, sample

324 size, burn-in period, and number of MCMC (Markov chain Monte Carlo) replications were 1,000

325 (sampled at random), 140 (70 DH lines from each population), 10,000, and 40,000, respectively.

326 The number of populations assumed ($K$) ranged from 1 to 4, and the most probable $K$ value was

327 determined based on the inferred plateau method (Viana et al. 2013). The LD analyses were

328 performed with *Haploview* (Barrett et al. 2005).

329 **Data availability**

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

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

332 Viana, José Marcelo Soriano; Pereira, Hélcio Duarte; Mundim, Gabriel Borges; Piepho, Hans-Peter;

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

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

335             **RESULTS**

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

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

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

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

340 (approximately $-10$ and $-13\%$ for populations 1 and 2), and genotypic variance (approximately 6

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

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

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

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

345 0.0041 and 0.0041, and 0.0054 and 0.0074 assuming one DH line per $S_0$, one to five DH lines per

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

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

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

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

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

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

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

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

354    2.1% of the SNPs. For approximately 70% of the SNPs, the absolute difference between allelic

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

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

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

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

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

360    process).

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

377    Decreasing the average SNP density to 1 cM led to a slightly decrease in the prediction accuracy of

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

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

380    proportional to the heritability. To evidence that the prediction accuracy of non-assessed single

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

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

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

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

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

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

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

388    (2012a), we observed no differences for the prediction accuracies of non-assessed single crosses

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

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

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

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

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

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

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

396    BLUP approaches, except for one to five DH lines per $S_3$ plant, where BLUP was −19 to −27%

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

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

399    0.55 to 0.76, proportional to the heritability) and from selected DH lines of the same population

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

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

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

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

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

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

406    −3%. Only for one DH line per $S_0$ plant the coincidence indexes assuming bidirectional dominance

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

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

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

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

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

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

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

414    0.60 to 0.74 (up to 26 to 40% lower), depending on the DH line sampling process.

## DISCUSSION

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

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

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

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

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

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

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

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

424 2013b) and barley (Philipp et al. 2016), along the last 20 years. Why, then, we did not find

425 published information that prediction of untested single crosses is of general use by breeders of

426 worldwide seed companies? What the scientific investigation should additionally prove to make

427 prediction of untested single crosses as successful as the Jenkins' (1934) method for predicting

428 double crosses performance was? We believe that this paper offers the final proof.

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

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

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

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

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

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

435 crosses and not on a cross-validation procedure involving assessed single crosses. This not means

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

437 procedure. We know that simulated data, because the presuppositions, cannot integrally describe the

438 complexity of populations and genetic determination of traits (Daetwyler et al. 2013). To highlight

439 the relevance of (overall) LD, our study is based on scenarios not favorable to prediction of untested

440 single cross performance: very low level of relationship between the DH lines, low and intermediate

441 heritabilities for the assessed single crosses, and not higher heterotic pattern. In the studies of

442 Massman et al. (2013) and Bernardo (1994, 1995, 1996a) the relationship among inbreds from the

443 same heterotic group ranged from 0.11 to 0.58. Riedelsheimer et al. (2012) observed high

444 relationships only within the non-Stiff Stalk inbreds. Technow et al. (2012a) assumed non-related

445 inbreds. For most of the investigations on prediction of untested single crosses and testcrosses, the

446 grain yield heritability ranged from 0.72 to 0.88. The common heterotic patterns in these previous

447 studies are Stiff Stalk and non-Stiff Stalk, and Dent and Flint. The MAF in the groups of Dent and

448 Flint inbreds were approximately 0.10 and 0.20, respectively, and approximately 20% of the SNPs

449 showed a difference of allelic frequency of at least 0.6.

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

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

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

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

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

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

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

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

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

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

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

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

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

463      0.72 to 0.93). Under maximum accuracy of assessed single crosses (1.0), the prediction accuracy

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

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

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

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

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

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

470      screening or prior to field testing.

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

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

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

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

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

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

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

478     (2016), Technow et al. (2014), Massman et al. (2013), Technow et al. (2012a), and Riedelsheimer et

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

480     and training set size.

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

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

483     of genomic prediction). Although breeders do not have control on LD and relatedness between the

484     DH or inbred lines, because selection they should always expect high level of overall LD in the

485     groups of selected DH or inbred lines. Comparison of our LD assessment with the LD analyses

486     from other studies is inadequate because we have distances in cM and not in base-pairs. But in

487     general the level of LD was high ($r^2$ of approximately 0.3) only for SNPs separated by up to 0.5 Mb

488     (Technow et al. 2014; Massman et al. 2013; Technow et al. 2012a; Riedelsheimer et al. 2012). To

489     maximize the prediction accuracy and the efficiency of identification of the best non-assessed single

490     crosses it is necessary to adopt the random sampling of single crosses for testing instead of the

491     random sampling of DH or inbred lines for a diallel. This is because sampling 30 or even 10% of

492     the single crosses leads to single crosses for testing derived from all DH or inbred lines from each

493     group. In our case, in every resampling assuming training set size of 30 and 10% we always get

494     groups of assessed single crosses (1,470 and 490 single crosses, respectively) derived from the 70

495     DH lines of each group. However, sampling DH lines for a diallel provided 1,440 and 484 single

496     crosses for testing derived from 38 and 22 DH lines, respectively. Thus, the sampling of single

497     crosses provides best prediction of the SNP average effects of substitution. Riedelsheimer et al.

498     (2012) emphasized the need for large genetic variability to obtain high prediction accuracies.

499     Further, their results indicated that pairs of closely related lines and population structuring only

500     weakly contributed to the high prediction accuracies. Regarding dominance, because it can be a

501     relevant genetic effect, breeders should always fit the additive-dominance model to maximize the

502   prediction accuracy and the efficiency of identification of the best non-assessed single crosses.

503   Interestingly, in most of the studies on prediction of non-assessed single crosses the prediction

504   accuracy did not significantly increase when modeling SCA in addition to GCA effects (Zhao et al.

505   2015).

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

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

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

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

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

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

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

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

514   9%), but the coincidence index can be reduced in up to 21%. However, considering that the

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

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

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

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

519   the accuracy is reached with a few hundred markers. Technow et al. (2014) did not improved

520   prediction accuracies by using higher SNP density. Additionally, the increase in the training set size

521   led to a relative small increase in the prediction accuracy. However, the prediction accuracies

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

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

524   et al. (2012a), the prediction accuracies increased with SNP density and number of parents tested in

525   hybrid combination.

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

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

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

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

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

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

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

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

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

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

536   0.95, for most studies) with those obtained with rice, wheat, and barley untested hybrids (0.50 to

537   0.60, approximately) (Philipp et al. 2016; Xu et al. 2014; Zhao et al. 2013b). However, the lower

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

539   LD level. Regarding the statistical approach, our model did not provide an increase in the efficiency

540   of non-assessed single cross prediction, compared to the models proposed by Massman et al. (2013)

541   and Technow et al. (2012a). It is important to highlight that our results showed that these two

542   models are really identical (data no shown). Thus, because the simplified definition of the incidence

543   matrices for these two previous models, it is quite safe to use any of them. Finally, the choice

544   between the statistical approaches RR-BLUP (prediction of genotypic values of non-assessed single

545   crosses based on prediction of SNP average effects of substitution), GBLUP (prediction of

546   genotypic values of non-assessed single crosses based on additive and dominance genomic

547   matrices), and BLUP (prediction of genotypic values of non-assessed single crosses based on

548   additive and dominance matrices from pedigree records) is not a serious worry for breeders too. Our

549   evidence is that there is no significant difference between RR-BLUP and GBLUP regarding

550   prediction accuracy and efficiency of identification of the best untested single crosses. Further, even

551   when the level of relatedness between the DH or inbred lines in each group is low, in general BLUP

552   is as efficient as genomic prediction, excepting when the DH lines are derived from inbred

553   population. Thus, DNA polymorphism is not essential for an efficient prediction of non-assessed

554    single cross performance. In his review on genomic selection in hybrid breeding, Zhao et al. (2015)

555    state that the choice of the biometrical model has no substantial impact on the prediction accuracy

556    of untested single crosses. Technow et al. (2014) observed that prediction methods GBLUP and

557    BayesB resulted in very similar prediction accuracies. In the study of Massman et al. (2013), BLUP

558    and RR-BLUP models did not lead to prediction accuracies that differed significantly. Comparing

559    GBLUP and BayesB, Technow et al. (2012a) concluded that the latter method produced

560    significantly higher accuracies for the additive-dominance models.

**LITERATURE CITED**

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

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

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

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

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

571    Bernardo, R., 1994 Prediction of maize single-cross performance using RFLPs and information

572      form related hybrids. *Crop Science* 34: 20-25.

573    Bernardo, R., 1995 Genetic models for predicting maize single-cross performance in unbalanced

574      yield trial data. *Crop Science* 35: 141-147.

575    Bernardo, R., 1996a Best linear unbiased prediction of maize single-cross performance. Crop Sci

576      36: 50-56.

577    Bernardo, R., 1996b Best linear unbiased prediction of maize single-cross performance given

578      erroneous inbred relationships. *Crop Science* 36: 862-866.

579 Bernardo, R., 1996c Best linear unbiased prediction of the performance of crosses between untested

580      maize inbreds. *Crop Science* 36: 872-876.

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

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

583      *Genet* 16.

584 Barrett, J.C., B. Fry, J. Maller, and M.J. Daly, 2005 Haploview: analysis and visualization of LD

585      and haplotype maps. *Bioinformatics* 21 (2):263-265.

586 Daetwyler, H.D., M.P.L. Calus, R. Pong-Wong, G. de los Campos, and J.M. Hickey, 2013 Genomic

587      Prediction in Animals and Plants: Simulation of Data, Validation, Reporting, and

588      Benchmarking. *Genetics* 193 (2):347-+.

589 de Los Campos, G., J.M. Hickey, R. Pong-Wong, H.D. Daetwyler, and M.P. Calus, 2013 Whole-

590      genome regression and prediction methods applied to plant and animal breeding. *Genetics* 193

591      (2):327-345.

592 Endelman, J.B., 2011 Ridge Regression and Other Kernels for Genomic Selection with R Package

593      rrBLUP. *Plant Genome* 4 (3):250-255.

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

595      genotype data: linked loci and correlated allele frequencies. *Genetics* 164:1567–1587

596 Henderson, C.R., 1974 General flexibility of linear model techniques for sire evaluation. *Journal of*

597      *Dairy Science* 57:963–972.

598 Jenkins, M.T, 1934 Methods of estimating the performance of double crosses in corn. *Journal of the*

599      *American Society of Agronomy* 26:199–204.

600 Jonas, E., and D.J. de Koning, 2013 Does genomic selection have a future in plant breeding? *Trends*

601      *in Biotechnology* 31 (9):497-504.

602 Kadam, D.C., S.M. Potts, M.O. Bohn, A.E. Lipka, and A.J. Lorenz, 2016 Genomic Prediction of

603      Single Crosses in the Early Stages of a Maize Hybrid Breeding Pipeline. *G3-Genes Genomes*

604      *Genetics* 6 (11):3443-3453.

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

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

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

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

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

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

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

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

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

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

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

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

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

637    Viana, J.M.S., M.S.F. Valente, F.F. Silva, G.B. Mundim, and G.P. Paes, 2013 Efficacy of
638        population structure analysis with breeding populations and inbred lines. *Genetica* 141 (7-
639        9):389-399.

640    Windhausen, V.S., G.N. Atlin, J.M. Hickey, J. Crossa, J.-L. Jannink *et al.*, 2012 Effectiveness of
641        Genomic Prediction of Maize Hybrid Performance in Different Breeding Populations and
642        Environments. *G3-Genes Genomes Genetics* 2 (11):1427-1436.

643    Xu, S., D. Zhu, and Q. Zhang, 2014 Predicting hybrid performance in rice using genomic best linear
644        unbiased prediction. *Proceedings of the National Academy of Sciences of the United States of*
645        *America* 111 (34):12456-12461.

646    Zhao, Y., M. Gowda, W. Liu, T. Wuerschum, H.P. Maurer *et al.*, 2013a Choice of shrinkage
647        parameter and prediction of genomic breeding values in elite maize breeding populations. *Plant*
648        *Breeding* 132 (1):99-106.

649    Zhao, Y., M.F. Mette, and J.C. Reif, 2015 Genomic selection in hybrid breeding. *Plant Breeding*
650        134 (1):1-10.

651    Zhao, Y., J. Zeng, R. Fernando, and J.C. Reif, 2013b Genomic Prediction of Hybrid Wheat
652        Performance. *Crop Science* 53 (3):802.

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

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

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

656 crosses, four statistical models, three DH lines sampling processes, two genetic models, and three
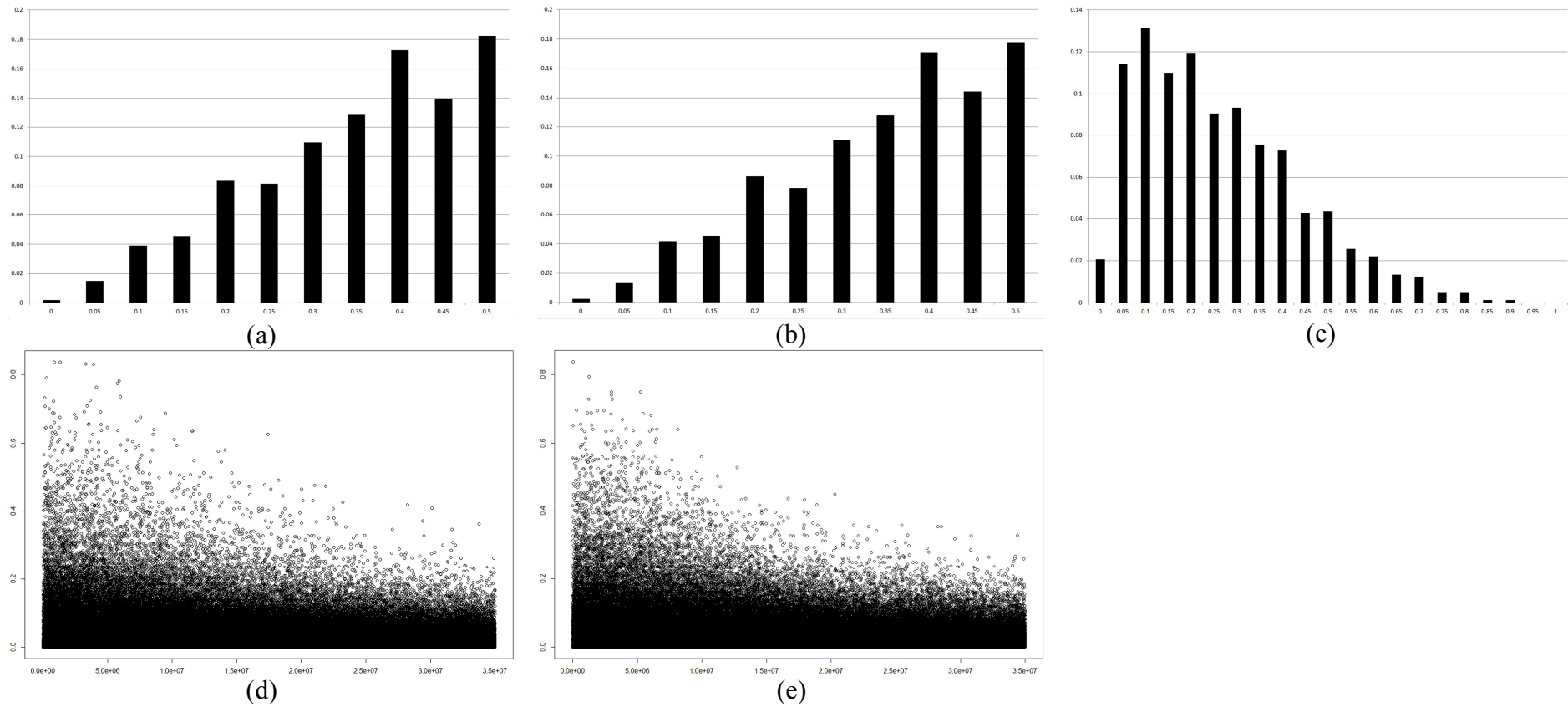
657 accuracies of assessed single crosses

| 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. | $1/S_0$ | 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.
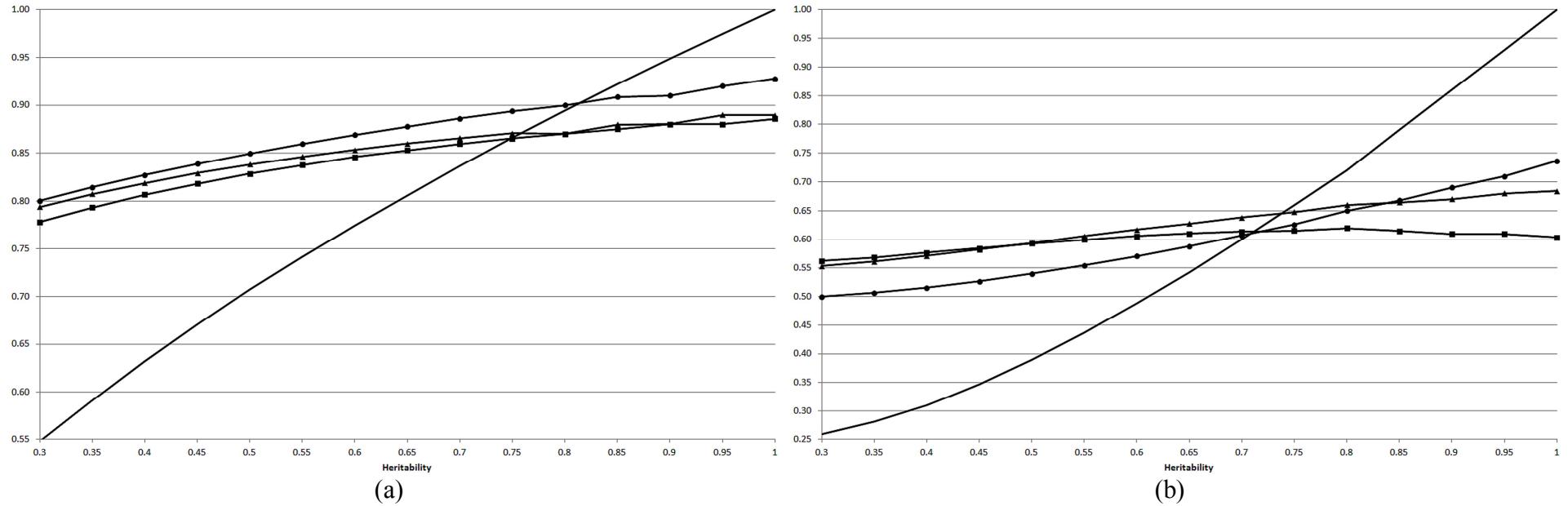
658 **Table 2** Average coincidence of the best 300 predicted single crosses and its standard deviation,

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

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

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

662 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. | $1/S_0$ | 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.

**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 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 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

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

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