# Identify Superior Parental Lines for Biparental Crossing via Genomic Prediction: Rice as an Example

**Ping-Yuan Chung**
  National Taiwan University

**Chen-Tuo Liao** ( ✉ ctliao@ntu.edu.tw )
  Southern Taiwan University    https://orcid.org/0000-0001-9777-3701

1 **Identify Superior Parental Lines for Biparental Crossing**

2 **via Genomic Prediction: Rice as an Example**

3 Ping-Yuan Chung and Chen-Tuo Liao *

4 Affiliations: PY Chung and CT Liao, Department of Agronomy, National Taiwan

5 University, Taipei, Taiwan. *Corresponding author (ctliao@ntu.edu.tw)

6 **ABSTRACT**

7 **Background:** A set of superior parental lines is the key to high-performing recombinant

8 inbred lines (RILs) for biparental crossing in a rice breeding program. The number of

9 possible crosses in such a breeding program is often far greater than the number that

10 breeders can handle in the field. A practical parental selection method via genomic

11 prediction (GP) is therefore developed to help breeders identify a set of superior

12 parental lines from a candidate population before field trials.

13 **Results:** The parental selection via GP often involves truncation selection, selecting the

14 top fraction of accessions based on their genomic estimated breeding values (GEBVs).

15 However, the truncation selection inevitably causes a loss of genomic diversity in the

16 breeding population. To preserve genomic variation, the selection of closely related

17 accessions should be avoided. We first proposed a new index to quantify the genomic

18 diversity for a set of candidate accessions. Then, we compared the performance of three

19 classes of strategy for the parental selection, including those consider (a) GEBV only, (b)

20 genomic diversity only, and (c) both GEBV and genomic diversity. We analyzed two

21 rice (*Oryza sativa* L.) genome datasets for the comparison. The results show that the

22 strategies considering both GEBV and genomic diversity have the best or second-best

23 performance for all the traits analyzed in this study.

24 **Conclusion:** Combining GP with Monte Carlo simulation can be a useful means of

25 parental selection for rice pre-breeding programs. Different strategies can be

26 implemented to identify a set of superior parental lines from a candidate population. In

27 consequence, the strategies considering both GEBV and genomic diversity that can

28 balance the starting GEBV average with maintenance of genomic diversity should be

29      recommended for practical use.

30      **Keywords:** genomic prediction, genomic selection, mixed effects model, rice breeding.

31

32                                    **BACKGROUND**

33      Biparental crossing is a commonly used scheme in pure-line breeding for self-pollinated

34      crops such as rice, wheat (*Triticum aestivum* L.), soybean [*Glycine max* (L.) Merr.] and

35      oat (*Avena sativa* L.). Plant breeders cross two inbred parental lines to produce $F_1$

36      population, then a subset of diverse individuals of the $F_2$ population is selected to

37      produce potential RILs after several generations of selfing. Obviously, the parental lines

38      play a fundamental role in the line development and significantly affect the performance

39      of the resulting RILs. However, the identification of superior parental lines from

40      germplasm collections for creating genetic variation to maximize selection response in

41      subsequent cycles is still a challenge for plant breeders (Bernardo 2003; Witcombe et al.

42      2013). Another practical concern is that the number of possible crosses in such a

43      breeding program is often far greater than the number that breeders can handle in the

44      field. Therefore, it should be of great help to breeders if a limited number of superior

45      parents can be identified before the field trial.

46          Genomic selection based on the statistical method of GP has been used to improve

47      breeding efficiency in dairy cattle (Hayes et al. 2009) and a variety of crops (Massman

48      et al. 2013; Asoro et al. 2011; Heffner et al. 2011; Lorenz et al. 2012; Spindel et al.

49      2015). The main concept of GP is to capture all the effects of quantitative trait loci

50      (QTLs) by using dense DNA markers over the whole genome, assuming that the DNA

51      markers are in strong linkage disequilibrium with one or more QTLs (Meuwissen et al.

52      2001). The most commonly used DNA markers are single nucleotide polymorphisms

53      (SNPs). A GP model is first built using the phenotype and genotype data of a training

54      population. Then, GEBVs for the candidate individuals with known genotype data are

55      predicted through the resulting GP model. There are two kinds of mixed linear model

56      methods are widely employed to obtain the GEBVs: (i) best linear unbiased prediction

57      (BLUP) based on markers and (ii) BLUP based on a genomic relationship matrix. For

58    the BLUP of (i), the marker effects are treated as random effects and the GEBVs of

59    individuals are calculated by multiplying their marker scores by these BLUP estimates.

60    Ridge regression BLUP (rr-BLUP) method (Meuwissen et al. 2001; Piepho 2009)

61    follows this approach. For the BLUP of (ii), the genotypic values of individuals are

62    treated as random effects and estimated through a genomic relationship matrix. The

63    genomic BLUP (GBLUP) method (Habier et al. 2007; VanRaden 2008) follows this

64    approach. For more details regarding the GP models and the estimation methods used

65    for their model parameters, refer to Xavier et al. (2016).

66        Gaynor et al. (2017) proposed a two-part strategy for implementing genomic

67    selection for line development, addressing the two components: (i) a product

68    development component, to identify inbred lines either for hybrid parent development

69    or cultivar release; (ii) a population improvement, to increase the frequency of favorable

70    alleles through rapid recurrent genomic selection. Conducting a stochastic simulation,

71    they showed that programs using the two-part strategy generated up to 2.5 times more

72    genetic gain than conventional programs, and up to 1.5 times more genetic gain than the

73    best performing standard genomic selection strategy. Also, Yao et al. (2018) combined

74    GP with Monte Carlo simulation to select superior parents in wheat breeding before the

75    field trial. They used the criterion of usefulness function on a selection index,

76    incorporating yield and two quality traits, to evaluate a cross. Their usefulness function

77    took into account both the mean genetic value and genetic variance of progeny

78    populations. Yao et al. (2018) simulated the required progeny populations using the

79    R/qtl package (Broman et al. 2003), and calculated their usefulness function estimates.

80    It was concluded that the use of the usefulness function for parental selection resulted in

81    higher genetic gain than the use of mid-parent GEBV, implying that the strategy for the

82    parental selection cannot only consider GEBVs of the candidate accessions.

83        Selecting the parental lines with the highest GEBVs (truncation selection),

84    breeders hope to maximally pass favorable properties of the parental lines on to their

85    progeny populations. However, several favorable QTLs can risk being eliminated from

86    the breeding population using the truncation selection (Vanavermaete et al. 2020). We

87    therefore take both GEBV and genomic diversity into account for identifying superior

88    parents in a biparental crossing program. For a specific target trait, we construct a

89  GBLUP model to predict the GEBVs for the candidate accessions. Furthermore, we

90  propose a new index to quantify the genomic diversity for a set of candidate accessions

91  according to the GBLUP model. We simulate the genotype data for progeny populations

92  over successive generations derived from a cross between two parental lines. The

93  GEBVs of the progeny populations are then predicted by the trained GBLUP model. We

94  further make generation advancement decisions according to the resulting GEBVs.

95  Finally, we assess a set of parental lines based on their $F_{10}$ RILs which are assumed to

96  be a fixed population. Several selection strategies are evaluated within two rice genome

97  datasets.


98


## MATERIALS AND METHODS


**The Rice Genome Datasets**


101  **Dataset I**: We first used the rice genome dataset presented in Zhao et al. (2011) to

102  illustrate our proposed procedure. This dataset was originally collected for

103  genome-wide association study (GWAS). The dataset contains 44,100 SNP variants and

104  36 traits of 413 *O. sativa* accessions, and has a strong subpopulation structure

105  containing six different groups. We deleted any SNPs with a missing rate of > 0.05 and

106  a minor allele frequency of < 0.05. To reduce redundant collinearity in calculation of the

107  genomic relationship matrix, we only retained about one-third of the SNPs which are

108  evenly distributed over each chromosome. We then imputed a missing SNP marker from

109  its corresponding major homozygous alleles. The final marker matrix consists of 413

110  accessions and 11,047 SNPs. We here analyzed the six traits: brown rice seed width

111  (BRSW), florets per panicle (FPP), flowering time at Arkansas (FTAA), flowering time

112  at Faridpur (FTAF), plant height (PH), and panicle number per plant (PNPP).


113  **Dataset II**: We further analyzed the rice genome dataset presented in Spindle et al.

114  (2015), which was collected for genomic selection study. The dataset contains 73,147

115  SNP variants and 363 elite breeding lines belonging to *indica* or *indica-admixed* group.

116  The phenotype data include the four years (2009-2012), two seasons per year (dry and

4

117    wet), of grain yield (YLD), flowering time (FT), and plant height (PH). Note that the

118    PH data in 2009 wet season are not available. The adjusted means for 328 out of the 363

119    individuals and 10,772 out of the 73,147 SNP markers were used for this study. We here

120    chose one marker every 0.1cM over each chromosome.


121    **Monte Carlo Simulation for the Genotype of Progeny Populations**

122    To simulate the genotype data for progeny populations, we used Gramene Annotated

123    Nipponbare Sequence (Youens-Clark et al. 2011) to estimate recombination rates

124    between two adjacent SNPs. The Gramene Annotated Nipponbare Sequence database

125    contains both the physical and linkage distances between SNPs, which can be

126    downloaded from http://archive.gramene.org. The genetic positions of the SNPs are

127    estimated via linear interpolation between the two markers flanking each SNP. Once the

128    genetic positions were obtained, the recombination rates between adjacent SNPs were

129    estimated via Haldane's mapping function (Haldane 1919):


130    $$r_{AB} = \frac{1}{2}(1 - e^{-2X_{AB}}),$$


131    where $r_{AB}$ is the recombination rate and $X_{AB}$ is the linkage distance between SNP

132    markers A and B. Through a series of Bernoulli distributions and the estimated

133    recombination rates, the crossover of each chromosome was simulated to yield the

134    sequence of a gamete, then two gametes were paired to produce the genotype data for

135    the progeny.


136    **GBLUP Model**

137    We considered the following single-trait GBLUP model for GP:


138    $$\boldsymbol{y} = \mu\mathbf{1}_n + \boldsymbol{g} + \boldsymbol{e}, \quad [1]$$


139    where **y** denotes the vector of phenotypic values of a training population with $n$

140    individuals; $\mu$ is a constant term; $\mathbf{1}_n$ is the vector of order $n$ with all elements equal to

141    1; $\boldsymbol{g}$ stands for the vector of genotypic values and $\boldsymbol{e}$ is the vector of random errors. It

142    is assumed that $\boldsymbol{g}$ follows a multi-variate normal distribution $\text{MVN}(\mathbf{0}, \sigma_g^2\boldsymbol{K})$, where $\mathbf{0}$

143   is a zero vector; $\sigma_g^2$ is the genetic variance of additive effects and $K$ is a genomic

144   relationship matrix among the individuals. Furthermore, $e$ follows $\text{MVN}(\mathbf{0}, \sigma_e^2 I_n)$,

145   where $\sigma_e^2$ is the random error variance and $I_n$ denotes the identity matrix of order $n$.

146   Here, $g$ and $e$ are assumed to be mutually independent. In this study, we considered

147   the genomic relationship matrix $K = MM^T/p$, where $M$ is the marker score matrix

148   and $p$ is the number of SNP markers. The elements of $M$ are coded as $-1$, 0, and 1 for

149   the minor homozygous alleles ($A_1A_1$), the heterozygous alleles ($A_1A_2$), and the major

150   homozygous alleles ($A_2A_2$), respectively. The model parameters of the GBLUP model

151   can be estimated through Henderson's equations (Henderson 1984), given by:

152
$$\begin{bmatrix} n & \mathbf{1}_n^T \\ \mathbf{1}_n & I_n + \lambda K^{-1} \end{bmatrix} \begin{bmatrix} \hat{\mu} \\ \hat{g} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_n^T y \\ y \end{bmatrix}, \quad [2]$$

153   where the regularization parameter $\lambda$ is given by $\lambda = \dfrac{\sigma_e^2}{\sigma_g^2}$. We used the R function

154   mmer( ) in the R package sommer (Covarrubias-Pazaran 2016) to obtain the restricted

155   maximum likelihood estimates (REMLs) for the two variance components of $\sigma_g^2$ and

156   $\sigma_e^2$, and then plugged the resulting estimates into Eq. [2] to get $\hat{\mu}$ and $\hat{g}$.

157       Let $\hat{g}_{bp}$ be the vector of estimated genotypic values for a breeding population and

158   $K_{bp}$ be the genomic relationship matrix between the breeding population and the

159   training population. In the case, we have:

160
$$\hat{g}_{bp} = K_{bp} K^{-1} \hat{g}.$$

161   The GEBV for the breeding population is $\hat{g}_{bp}$ plus the estimate of the constant term $\hat{\mu}$.

162   **The Index to Quantify Genomic Diversity**

163   Let $g_0$ be the vector of genotypic values and $K_0$ be the genomic relationship matrix

164   for a particular set of accessions with size $n_0$. According to the GBLUP model of Eq.

165   [1], the covariance matrix for $g_0$ is given by:

166
$$\text{Var}(g_0) = \sigma_g^2 K_0.$$

167   The determinant of the covariance matrix represents the overall variability for the

168     genotypic values, which is calculated as:

169 $$|\text{Var}(\boldsymbol{g}_0)| = (\sigma_g^2)^{n_0}|\boldsymbol{K}_0|. \text{ [3]}$$

170     Clearly, the determinant of Eq. [3] is proportional to the $D$-score defined below:

171 $$D\text{-score} = |\boldsymbol{K}_0|. \text{ [4]}$$

172     The $D$-score of Eq. [4] ranges from 0 to 1. For a fixed number of $n_0$, a subset of

173     accessions chosen from a breeding population that achieves the maximal D-score will

174     have greater genomic diversity than the competing choices with size $n_0$. The concept of

175     the $D$-score is adopted from optimum experimental designs (Atkinson and Donev 1992).

176     A simple example is given to illustrate the $D$-score. Suppose that there are $n = 3$

177     accessions in the candidate set with the genomic relationship matrix:

178 $$\boldsymbol{K} = \begin{bmatrix} 1 & 0.7 & 0.5 \\ 0.7 & 1 & 0.3 \\ 0.5 & 0.3 & 1 \end{bmatrix}.$$

179     For $n_0 = 2$, the $D$-score for $g_1$ and $g_2$ is calculated as $|\boldsymbol{K}_0| = \begin{vmatrix} 1 & 0.7 \\ 0.7 & 1 \end{vmatrix} = 0.51$.

180     Similarly, the $D$-scores for $g_1$ and $g_3$, and for $g_2$ and $g_3$ are given by 0.75 and 0.91,

181     respectively. Clearly, the two accessions with $g_2$ and $g_3$ have greater genomic

182     variation (smaller genomic correlation) than the other competing choices. A set of

183     accessions with the maximal D-score can avoid the selection of closely related

184     individuals.

185     **An Algorithm to Search for Accessions with the Maximal $D$-Score**

186     We required a highly efficient algorithm to search for a subset of accessions within a

187     candidate population so that it can achieve the maximal $D$-score. We used a genetic

188     algorithm to complete this task, which is an exchange algorithm with the three different

189     operators: roulette wheel selection, crossover, and mutation (Whitley 1994). For a given

190     candidate set $S_c$ with $n_c$ accessions, we searched for an optimal subset $S_0$ with $n_0$

191     individuals from $S_c$. Our algorithm began with a set of $m$ random solutions, each of

192     which is a vector of 0 or 1 with a length equal to $n_c$. The number of values with a score

193     of 1 in the vector is equal to $n_0$, corresponding to the chosen accessions at the current

194     stage. Here, we fixed $m = n_0$. We then obtained the elite solutions from the initial $m$

195     random solutions after a large number of iterations, where each iteration repeated all the

196     three operators. We stopped the algorithm when the maximal *D*-score among the current

197     elite solutions converged.


**The Procedure for Selecting Parental Lines**

198

199     To evaluate a variety of strategies in determining parental lines, we carried out the

200     following steps.


201     Step 1: For a specific target trait, we used all of the phenotypic values available from

202     the rice genome dataset to build the corresponding GBLUP model of Eq. [1].


203     Step 2: We predicted the GEBVs for all of the accessions in the dataset through the

204     trained GBLUP model developed in Step 1. Seven strategies were used to select a subset

205     of 10 parental lines according to their GEBVs: (i) the GEBV only (GEBV-O) approach,

206     which chose the top 10 accessions (either maximal or minimal); the genomic diversity

207     only (GD-O) approaches: (ii) GD-O-30, (iii) GD-O-50, and (iv) GD-O-100, which

208     applied the genetic exchange algorithm to search for an optimal subset of 10 accessions

209     from each of the three candidate sets composed of the top 30, 50, and 100 accessions,

210     respectively, such that the chosen subset had the maximal *D*-score; and the approaches

211     (GEBV-GD) considering both GEBV and genomic diversity: (v) GEBV-GD-30, (vi)

212     GEBV-GD-50, and (vii) GEBV-GD-100, which retained the top two accessions, then

213     applied the genetic exchange algorithm to search for another eight accessions from the

214     remainder of each candidate set for GD-O-30, GD-O-50, and GD-O-100, respectively,

215     so that the resulting 10 accessions had the maximal *D*-score.


216     Step 3: For each subset of 10 accessions determined by the seven strategies, we crossed

217     any two parental lines to produce 45 $F_1$ hybrids. Here, we started to simulate the

218     genotype data for successive generations of progeny populations through the Monte

219     Carlo simulation. Each of the 45 $F_1$ hybrids produced 60 individuals of the $F_2$

220     population by self-pollination, resulting in 2700 $F_2$ individuals. After obtaining the

221     GEBVs for the 2700 $F_2$ individuals via the trained GBLUP model of Step 1, we then

222     retained the top 45 $F_2$ individuals. Again, we used these 45 $F_2$ individuals to produce

223     2700 $F_3$ individuals (each $F_2$ individual produced 60 $F_3$ individuals) and retained the top

224     45 $F_3$ individuals. We then repeated the same procedure to produce 2700 $F_{10}$ individuals

225     which are assumed to be a fixed population.

226     Step 4: For the resulting 2700 $F_{10}$ individuals generated according to each strategy, we

227     found the best $F_{10}$ RIL with the top GEBV.

228       A flowchart of the procedure is displayed in Figure 1. We repeated this analysis

229     procedure 30 times to obtain the best $F_{10}$ RILs from each repetition for each strategy.

230     The average of the GEBVs for the best $F_{10}$ RILs was then calculated and used as the

231     measure of efficiency for the strategy. Note that for the traits of BRSW, FPP, and PNPP

232     in Dataset I; and YLD in Dataset II, larger GEBVs are preferable (i.e., these traits

233     follow the rule that the larger, the better). The remaining five traits of FTAA, FTAF, and

234     PH in Dataset I; and FT, and PH in Dataset II are those for which the rule is "the smaller,

235     the better".

236     **Calculation of Genetic Gain**

237       To gain an understanding of the genetic improvement on a target trait using

238     different strategies, we estimated genetic gain as

239    
$$\text{genetic gain} = \overline{GEBV}_{F_{10}} - \overline{GEBV}_P, \quad [5]$$

240     where $\overline{GEBV}_{F_{10}}$ denotes the GEBV average among the resulting 2700 $F_{10}$ RILs and

241     $\overline{GEBV}_P$ denotes the GEBV average among the 10 selected parental lines for each

242     strategy (Rutkoski 2019). The larger absolute value of the genetic gain indicates the

243     more improvement on the target trait.

244

245                             **RESULTS**

**Strategies Comparison Based on the best $F_{10}$ RILs**

246

247    The GEBV averages of the best $F_{10}$ RILs from the 30 repetitions using each of the

248    seven strategies are displayed in Tables 1 and 2 for the two datasets. The results in the

249    tables show that the strategies considering both GEBV and genomic diversity

250    (GEBV-GD-30, -50, -100) generally have satisfactory efficiency, because they achieve

251    the best or second-best performance for all the traits. Therefore, this kind of strategies

252    could be a reliable means of determining the parental lines. On the other hand, the

253    strategies accounting for genomic diversity only (GD-O-30, -50, -100) don't have

254    satisfactory efficiency for all the traits, with the exception of GD-O-100 for YLD in

255    Dataset II. For the strategy based on GEBV only, the GEBV-O has the best or

256    second-best performance for FPP, and PH in Dataset I; and PH, and FT in Dataset II, but

257    also has the worst or second-worst performance for the remaining four traits in Dataset I

258    and YLD in Dataset II. Thus, the GEBV-O could be a high-risk strategy.

259    We also displayed the GEBV averages with the plus and minus one unit of their

260    corresponding standard deviations for the best individuals from the 30 repetitions over

261    consecutive generations in Figures 2 and 3. From the figures, the four strategies of

262    GEBV-O, GEBV-GD-30, -50, -100 selected the same best individual from the 30

263    repetitions at parental generation, and also at $F_1$ generation, so there is no standard

264    deviation shown with the corresponding GEBV averages. The GEBV averages of the

265    best selected parental lines by the strategies can be ranked as GEBV-O = GEBV-GD-30

266    = GEBV-GD-50 = GEBV-GD-100 > GD-O-30 > GD-O-50 > GD-O-100 in decreasing

267    desirability. The desirability at parental generation decreases as the degree of diversity

268    increases for the three strategies considering the genomic diversity only. Also, the

269    desirability declines from parental generation to $F_1$ generation for every strategy, due to

270    the heterogenous alleles in $F_1$ hybrids.

271    To explore the extent to which the top two accessions contribute to the subset of

272    ten parental lines determined by the four strategies of GEBV-O, GEBV-GD-30, -50,

273    -100, we compared each subset with a reduced group consisting of $F_1$ hybrids whose

274    parental lines contain at least one of the top two accessions for each subset. Every

275    reduced group consists of 17 $F_1$ hybrids. Similarly, we followed the analysis procedure

276 to obtain the GEBV averages for the best $F_{10}$ RILs from 30 repetitions based on the
277 reduced group. The results are displayed in Table 3 with the corresponding GEBV
278 averages based on the group of the original 45 $F_1$ hybrids. From the table, there is no
279 practical significant difference between these two groups for all the traits using the four
280 strategies.

281 **Genetic Gains for the Strategies**

282 The average among the genetic gains on a target trait for each strategy calculated
283 by Eq. [5] from the 30 repetitions is displayed in Tables 4 and 5 for Datasets I and II,
284 respectively. It is reasonable to compare the performance of the strategies according to
285 the endpoint of $\overline{GEBV}_{F_{10}}$. From the tables, we found that the comparison results based
286 on $\overline{GEBV}_{F_{10}}$ are consistent with the above results based on the best $F_{10}$ RILs. Also, the
287 strategies considering genomic diversity (GD-O-30, -50, -100; GEBV-GD-30, -50, -100)
288 have greater genetic gain than the GEBV-O for all the traits except PH in Dataset I
289 (Table 4). As expected, the genetic gain usually increases with the increase of the
290 genomic diversity (GD-O-100 outperforms both GD-O-50 and GD-O-30 for all the
291 traits except BRSW, and FTAF in Dataset I; GEBV-GD-100 outperforms both
292 GEBV-GD-50 and GEBV-GD-30 for all the traits). In addition, GEBV-O has the best
293 $\overline{GEBV}_P$; GEBV-GD-30 has better $\overline{GEBV}_P$ than GD-O-30; GEBV-GD-50 has better
294 $\overline{GEBV}_P$ than GD-O-50 and GEBV-GD-100 has better $\overline{GEBV}_P$ than GD-O-100 for all
295 the traits. Namely, a strategy has a relatively good starting point as it considers more
296 degree of GEBV.

297

298 **DISCUSSION**

299 From the results for comparing the proposed strategies, those considering both GEBV
300 and genomic diversity or considering GEBV only can be recommended for practical use.
301 Furthermore, from the results for exploring the extent to which the top two accessions
302 contribute to the parental lines determined by the four strategies of GEBV-O,
303 GEBV-GD-30, -50, -100, we have the conclusion: the economical strategies with 17 $F_1$

304      hybrids whose parental lines contain at least one of the top two accessions for each

305      selected subset can be a practical alternative to those with 45 $F_1$ hybrids composed of all

306      of the possible crosses.


307      From Tables 4 and 5, the strategies considering genomic diversity only (GD-O-30,

308      -50, -100) generally have greater genetic gain, mainly due to their more genomic

309      variation but less favorable $\overline{GEBV}_P$, so they have more room to improve. Also, the

310      GEBV-O has the best starting $\overline{GEBV}_P$ but the least genomic diversity in the base

311      population, so it has less potential to improve. The strategies considering both GEBV

312      and genomic diversity (GEBV-GD-30, -50, -100) could balance the tradeoff between

313      starting $\overline{GEBV}_P$ and genomic variation of the base population.


314      Dataset II was specifically collected for genomic selection. All of the available

315      accessions in the dataset belong to *indica* or *indica-admixed* group. From the results of

316      the performance based on the best $F_{10}$ RILs in Table 2, all the seven strategies seem to

317      have close performance for the three target traits. The resulting GEBV averages of the

318      best $F_{10}$ RILs range from 6472 to 6546 kg/ha for YLD, from 85.889 to 91.852 cm for

319      PH, and from 77.725 to 78.410 days for FT. This could be due to the fact that the

320      candidate accessions in Dataset II are elite breeding lines which have limited genomic

321      diversity and similar phenotypic values for the target traits. However, the two strategies

322      with greater genomic diversity, GD-O-100 and GEBV-GD-100 for YLD (their

323      corresponding GEBV averages are 6546 and 6539 kg/ha), led to larger YLD than the

324      other five strategies (their corresponding GEBV averages range from 6472 to 6506

325      kg/ha). The four strategies of GEBV-O, GEBV-GD-30, -50, -100 performed equally

326      well for PH (their corresponding GEBV averages range between 85.817 and 86.062 cm),

327      but slightly better than GD-O-30, -50, -100 (their corresponding GEBV averages are

328      87.517, 89.920, and 91.799 cm). The consistent results based on the $\overline{GEBV}_{F_{10}}$ can be

329      found in Table 5.


330      It is known that Dataset I contains more genomic diversity than Dataset II, since it

331      consists of five subpopulations and one admixed group. The more genomic diversity of

332      Dataset I could lead to a bigger difference between the strategies considering both

333      GEBV and genomic diversity, and the strategy considering GEBV only for some traits.

334    For example, the difference of the GEBV averages among the best $F_{10}$ RILs between
335    GEBV-GD-50 and GEBV-O is about -9.06 days for FTAA, and -2.55 days for FTAF in
336    Dataset I (Table 1), but the corresponding difference is just -0.09 days for FT in Dataset
337    II (Table 2). However, the flowering time is very sensitive to environments, so the
338    genomic diversity cannot solely amount to the different results between these two
339    datasets. More interestingly, the more genomic diversity of Dataset I could lead to a
340    larger genetic gain for a specific trait. From Table 4, the mean of the genetic gains using
341    the seven strategies for PH in Dataset I is given by -42.15 cm. But, from Table 5, the
342    corresponding mean in Dataset II is just -13.79 cm.

343    Daetwyler et al. (2015) and Goiffon et al. (2017) highlighted that an increase in
344    rare favorable alleles in a population can help improve selection responses. Selecting
345    only parental lines with the highest GEBVs can result in a loss of rare favorable alleles
346    for some target traits, thus missing potential RILs over future generations. From the
347    results of BRSW, FTAA, FTAF, and PNPP in Figure 2; and YLD in Figure 3, the
348    performance of GEBV-O appears to be inferior to GEBV-GD-30, -50, -100. This
349    indicates that an increase in genomic diversity in parental lines could compensate for
350    this possible deficiency, and then improve the long-term response to the target traits.
351    The greater genomic diversity could increase the possibility of containing favorable
352    alleles in parental lines, and we therefore expect that the chance of harboring the
353    favorable alleles would increase in RIL populations.

354    Apparently, the numbers of accessions fixed in the proposed strategies seem to be a
355    little arbitrary, such as those of selecting 10 parental lines, retaining the top 2 accessions,
356    and searching 10 or another 8 accessions from the three candidate sets composed of the
357    top 30, 50 and 100 accessions, respectively. A user certainly can adjust these numbers in
358    the strategies for her/his own study. Also, it was required to have historical phenotypic
359    data used to build the GP model. If the historical phenotypic data are not available, then
360    a pilot experiment is needed to phenotype a set of accessions, which can be determined
361    using an optimization algorithm (Ou and Liao 2019). An R function for performing the
362    proposed procedure of selecting parental lines is available from the authors upon
363    request.

364    As mentioned earlier, Yao et al. (2018) evaluated genetic gain using the usefulness

365    function in parental selection, and showed that their selection strategy outperformed the

366    strategy using the mid-parent GEBV. In this study, we emphasized both GEBV and

367    genomic diversity in parental selection, and we made generation advancement decisions

368    for each selection strategy according to the GEBVs of the top individuals. Finally, we

369    compared different strategies based on the performance of the best $F_{10}$ RILs, and

370    discussed genetic gain for target traits using the strategies. Moreover, Yao et al. (2018)

371    showed that applying a selection index incorporating multiple traits can simultaneously

372    improve both yield and quality in wheat than the individual trait selection. Also, Jia and

373    Jannink (2012), Hayashi and Iwata (2013), and Guo et al. (2014) highlighted that

374    multiple-trait GP models can provide better prediction accuracy than single-trait GP

375    models for those traits with low heritability. We will consider the selection index and the

376    multiple-trait GP models into the framework of the current study, so as to investigate the

377    multiple-trait situations in a future study.

378

## CONCLUSIONS

380    Combining GP with Monte Carlo simulation can be a useful means of detecting superior

381    parents for rice pre-breeding programs. Different strategies can be implemented to

382    identify a set of superior parental lines from a candidate population. The strategy

383    considering GEBV only can have a better starting GEBV average but less genomic

384    diversity in the base population. On the other hand, the strategies considering genomic

385    diversity only can have greater genomic diversity but a less favorable starting GEBV

386    average in the base population. The strategies considering both GEBV and genomic

387    diversity that can balance the starting GEBV average with maintenance of genomic

388    diversity should be recommended for practical use.

389

390    **Abbreviations:** BLUP, best linear unbiased predictor; BRSW, brown rice seed width;
391    GBLUP, genomic best linear unbiased predictor; GEBV, genomic estimated breeding

392  value; GEBV-GD, algorithms considering both GEBV and genomic diversity; GEBV-O,
393  algorithms considering GEBV only; FPP, florets per panicle; FT, flowering time; FTAA,
394  flowering time at Arkansas; FTAF, flowering time at Faridpur; GD-O, algorithms
395  considering genomic diversity only; PH, plant height; PNPP, panicle number per plant;
396  RIL, recombinant inbred line; SNP, single nucleotide polymorphism; YLD, grain yield.

397
398

399

400                                    **DECLARATIONS**

401  **Ethical Approval and Consent to participate:** Not applicable.

402  **Consent for publication:** Not applicable.

403  **Availability of supporting data:** Not applicable.

404  **Competing interests:** The authors declare that there is no conflict of interest.

405  **Funding:** Not applicable.

406  **Authors' contributions:** PY analyzed the datasets, wrote the R functions, prepared the
407  tables and figures, and drafted the manuscript. CT supervised the research, derived the
408  analysis approach, and drafted the manuscript.

411

412                                    **REFERENCES**

413

414  Asoro, F.G., M.A. Newell, W.D. Beavis, M.P. Scott, and J.L. Jannink. 2011. Accuracy
415        and training population design for genomic selection on quantitative traits in
416        elite north American oats. Plant Genome 4: 132–144.


417  Atkison, A.C., and A.N. Donev. 1992. Optimum experimental designs. New York:

418       Oxford University Press.

419    Bernardo, R. 2003. Parental selection, number of breeding populations, and size of each
420       population in inbred development. Theor. Appl. Genet. 107: 1252-1256.

421    Broman, K.W., H. Wu, S. Sen, and G.A. Churchill. 2003. R/qtl: QTL mapping in
422       experimental crosses. Bioinformatics 19: 889-890.

423    Covarrubias-Pazaran, G. 2016. Genome-assisted prediction of quantitative traits using
424       the R package sommer. PLOS One 11: e0156744.

425    Daetwyler, H.D., M.J. Hayden, G.C. Spangenberg, and B.J. Hayes. 2015. Selection on
426       optimal haploid value increases genetic gain and preserves more genetic
427       diversity relative to genomic selection. Genetics 200: 1341–1348.

428    Gaynor, R.C., G. Gorjanc, A.R. Bentley, E.S. Ober, P. Howell, R. Jackson, et al. 2017. A
429       two-part strategy for using genomic selection to develop inbred lines. Crop Sci.
430       57: 2372–2386.

431    Goiffon, M., A. Kusmec, L. Wang, G. Hu, and P.S. Schnable. 2017. Improving response
432       in genomic selection with a population-based selection strategy: Optimal
433       population value selection. Genetics 206: 1675–1682.

434    Guo, G., F. Zhao, Y. Wang, Y. Zhang, L. Du, and G. Su. 2014. Comparison of
435       single-trait and multiple-trait genomic prediction models. BMC Genetics 15: 30.

436    Habier, D., R.L. Fernando, and J.C.M. Dekkers. 2007. The impact of genetic
437       relationship information on genome-assisted breeding values. Genetics 177:
438       2389–2397.

439    Haldane, J.B.S. 1919. The combination of linkage values and the calculation of distance
440       between the loci for linked factors. Genetics 8: 299–309.

441    Hayashi, T., and H. Iwata. 2013. A Bayesian method and its variational approximation

442          for prediction of genomic breeding values in multiple traits. BMC
443          bioinformatics 14: 1–14.

444    Hayes, B.J., P.J. Bowman, A.J. Chamberlain, and M.E. Goddard. 2009. Genomic
445          selection in dairy cattle: Progress and challenges. J. Dairy Sci. 92: 433–443.

446    Heffner, E.L., J.L. Jannink, and M.E. Sorrells. 2011. Genomic selection accuracy using
447          multifamily prediction models in wheat breeding program. Plant Genome 4: 65–
448          75.

449    Henderson, C.R. 1984. Applications of linear models in animal breeding. Univ. of
450          Guelph, Guelph, Ontario.

451    Jia, Y., and J.L. Jannink. 2012. Multiple-trait genomic selection methods increase
452          genetic value prediction accuracy. Genetics 192: 1513-1522.

453    Lorenz, A.J., K.P. Smith, and J.L. Jannink. 2012. Potential and optimization of genomic
454          selection for *Fusarium* head blight resistance in six-row barley. Crop Sci. 52:
455          1609–1621.

456    Massman, J.M., H.J.G. Jung, and R. Bernardo. 2013. Genomewide selection versus
457          marker-assisted recurrent selection to improve grain yield and stover-quality
458          traits for cellulosic ethanol in maize. Crop Sci. 53: 58–66.

459    Meuwissen, T.H.E., B.J. Hayes, and M.E. Goddard. 2001. Prediction of total genetic
460          value using genome-wide dense marker maps. Genetics 157: 1819–1829.

461    Ou, J.H., and C.T. Liao. 2019. Training set determination for genomic selection. Theor.
462          Appl. Genet. 132: 2781–2792.

463    Piepho, H.P. 2009. Ridge regression and extensions for genome-wide selection in maize.
464          Crop Sci. 49: 1165–1176.

465    Rutkoski, J.E. 2019. A practical guide to genetic gain. Adv. Agron. 157: 217-249.

466    Spindel, J., H. Begum, D. Akdemir, P. Virk, B. Collard, E. Redona, et al. 2015. Genomic

467         selection and association mapping in rice (*Oryza sativa*): Effect of trait genetic

468         architecture, training population composition, marker number and statistical

469         model on accuracy of rice genomic selection in elite, tropical rice breeding lines.

470         PLOS Genetics 11: e1005350.

471    VanRaden, P.M. 2008. Efficient methods to compute genomic predictions. J. Dairy Sci.

472         91: 4414–4423.

473    Vanavermaete, D., J. Fostier, S. Maenhout, and B. De Baets. 2020. Preservation of

474         genetic variation in a breeding population for long-term genetic gain. G3:

475         Genes|Genomes|Genetics doi:10.1534/g3.120.401354.

476    Whitley, D. 1994. A genetic algorithm tutorial. Stat. Comput. 4: 65–85.

477    Witcombe, J.R., S. Gyawali, M. Subedi, D.S. Virk, and K.D. Joshi. 2013. Plant breeding

478         can be made more efficient by having fewer, better crosses. BMC Plant Biol. 13:

479         22.

480    Xavier, A., W.M. Muir, B. Craig, and K.M. Rainey. 2016. Walking through the

481         statistical black boxes of plant breeding. Theor. Appl. Genet. 129: 1933–1949.

482    Yao, J., D. Zhao, X. Chen, Y. Zhang, and J. Wang. 2018. Use of genomic selection and

483         breeding simulation in cross prediction for improvement of yield and quality in

484         wheat (*Triticum aestivum* L.). Crop J. 6: 353-365.

485    Youens-Clark, K., E. Buckler, T. Casstevens, C. Chen, G. DeClerck, P. Derwent, et al.

486         2011. Gramene database in 2010: updates and extensions. Nucleic Acids Res. 39:

487         D1085–D1094.

488    Zhao, K., C.W. Tung, G.C. Eizenga, M.H. Wright, M.L. Ali, A.H. Price, et al. 2011.

489         Genome-wide association mapping reveals a rich genetic architecture of

490         complex traits in *Oryza sativa*. Nature Communications 2: 467.

Table 1: The ranking and the GEBV average (in parentheses) for the best $F_{10}$ RILs from the 30 repetitions using the seven proposed strategies in Dataset I.

|  | BRSW | FPP | FTAA | FTAF | PH | PNPP |
|---|---|---|---|---|---|---|
| GEBV-O | 6 (3.418) | **2 (5.961)** | 6 (56.521) | 6 (61.856) | **1 (42.185)** | 6 (4.125) |
| GD-O-30 | 7 (3.408) | 5 (5.951) | 3 (51.564) | 3 (59.353) | 5 (49.337) | 3 (4.188) |
| GD-O-50 | 3 (3.576) | 6 (5.916) | 5 (53.348) | 5 (60.126) | 6 (49.801) | 5 (4.138) |
| GD-O-100 | 4 (3.496) | 7 (5.882) | 7 (56.835) | 7 (61.967) | 7 (51.788) | 7 (4.086) |
| GEBV-GD-30 | 5 (3.419) | 3 (5.954) | **1 (47.136)** | **1 (59.216)** | **2 (42.699)** | **1 (4.225)** |
| GEBV-GD-50 | **1 (3.656)** | **1 (5.964)** | **2 (47.457)** | **2 (59.304)** | 3 (43.232) | **2 (4.214)** |
| GEBV-GD-100 | **2 (3.634)** | 4 (5.953) | 4 (51.382) | 4 (59.634) | 4 (43.498) | 4 (4.171) |

(i) The best and second-best strategies are indicated in bold text, and the worst and second-worst strategies are indicated by underlining.

(ii) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs; GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen from the candidate sets composed of the top 30, 50, and 100 accessions, respectively; GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions, respectively, which have the maximal D-scores.

508 (iii) BRSW: brown rice seed width; FPP: florets per panicle; FTAA: flowering time at

509 Arkansas; FTAF: flowering time at Faridpur; PH: plant height; PNPP: panicle number

510 per plant.

511

512

513

514

515

516 Table 2: The ranking and the GEBV average (in parentheses) for the best $F_{10}$ RILs from

517 the 30 repetitions using the seven proposed strategies in Dataset II.

518

| | YLD | PH | FT |
|---|---|---|---|
| GEBV-O | 7 (6472) | **1 (85.817)** | **2 (77.818)** |
| GD-O-30 | 4 (6491) | 5 (87.517) | 7 (78.410) |
| GD-O-50 | 5 (6489) | 6 (89.920) | 5 (78.164) |
| GD-O-100 | **1 (6546)** | 7 (91.799) | 6 (78.359) |
| GEBV-GD-30 | 3 (6506) | **2 (85.976)** | 4 (77.883) |
| GEBV-GD-50 | 6 (6485) | 3 (85.917) | **1 (77.725)** |
| GEBV-GD-100 | **2 (6539)** | 4 (86.062) | 3 (77.873) |

519

520 (i) The best and second-best strategies are indicated in bold text, and the worst and

521 second-worst strategies are indicated by underlining.

522 (ii) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs;

523 GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen

524 from the candidate sets composed of the top 30, 50, and 100 accessions, respectively;

525 GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen

526 from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions,

527 respectively, which have the maximal D-scores.

528 (iii) YLD: yield; PH: plant height; FT: flowering time.

529

530

531

532

533

534 Table 3: The GEBV averages for the best $F_{10}$ RILs from the 30 repetitions based on the
535 group of the original 45 $F_1$ hybrids and the reduced group of 17 $F_1$ hybrids using the
536 four strategies of GEBV-O, GEBV-GD-30, GEBV-GD-50, and GEBV-GD-100.

537

|  | GEBV-O | | GEBV-GD-30 | | GEBV-GD-50 | | GEBV-GD-100 | |
|---|---|---|---|---|---|---|---|---|
| Dataset I | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ |
| BRSW | 3.418 | 3.423 | 3.419 | 3.418 | 3.656 | 3.652 | 3.634 | 3.650 |
| FPP | 5.961 | 5.965 | 5.954 | 5.957 | 5.964 | 5.958 | 5.953 | 5.943 |
| FTAA | 56.521 | 57.513 | 47.136 | 46.961 | 47.457 | 47.421 | 51.382 | 51.734 |
| FTAF | 61.856 | 61.850 | 59.216 | 59.123 | 59.304 | 59.232 | 59.634 | 59.713 |
| PH | 42.185 | 43.409 | 42.699 | 43.271 | 43.232 | 43.791 | 43.498 | 43.854 |
| PNPP | 4.125 | 4.129 | 4.225 | 4.226 | 4.214 | 4.204 | 4.171 | 4.161 |
| Dataset II | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ | 45 $F_1$ | 17 $F_1$ |
| YLD | 6472 | 6476 | 6506 | 6499 | 6485 | 6484 | 6539 | 6534 |
| PH | 85.817 | 85.991 | 85.976 | 85.844 | 85.917 | 86.092 | 86.062 | 86.060 |
| FT | 78.818 | 77.834 | 77.883 | 77.750 | 77.725 | 77.778 | 77.873 | 77.690 |

538

539 (i) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs;

540 GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen

541 from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions,

542 respectively, which have the maximal D-scores.

543 (ii) BRSW: brown rice seed width; FPP: florets per panicle; FTAA: flowering time at

544 Arkansas; FTAF: flowering time at Faridpur; PH: plant height; PNPP: panicle number

545 per plant.

546 (iii) YLD: yield; PH: plant height; FT: flowering time.

547

548

549

550 Table 4: The average of genetic gains from the 30 repetitions for Dataset I.

551

| | BRSW | | | FPP | | |
|---|---|---|---|---|---|---|
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 3.17 | 3.42 | 0.25 | 5.51 | 5.96 | 0.45 |
| GD-O-30 | 3.10 | 3.41 | 0.31 | 5.48 | 5.95 | 0.47 |
| GD-O-50 | 3.00 | 3.57 | 0.57 | 5.41 | 5.91 | 0.50 |
| GD-O-100 | 2.94 | 3.49 | 0.55 | 5.31 | 5.88 | 0.57 |
| GEBV-GD-30 | 3.12 | 3.42 | 0.30 | 5.48 | 5.95 | 0.47 |
| GEBV-GD-50 | 3.04 | 3.65 | 0.61 | 5.43 | 5.96 | 0.53 |
| GEBV-GD-100 | 3.00 | 3.63 | 0.63 | 5.34 | 5.95 | 0.61 |
| | FTAA | | | FTAF | | |
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 64.30 | 56.57 | -7.73 | 63.45 | 61.87 | -1.58 |
| GD-O-30 | 72.25 | 49.26 | -22.99 | 64.93 | 59.40 | -5.53 |
| GD-O-50 | 75.41 | 53.54 | -21.87 | 65.82 | 60.16 | -5.66 |
| GD-O-100 | 80.01 | 57.00 | -23.01 | 67.34 | 62.01 | -5.33 |
| GEBV-GD-30 | 71.09 | 47.31 | -23.78 | 64.68 | 59.25 | -5.43 |
| GEBV-GD-50 | 72.86 | 47.64 | -25.22 | 65.40 | 59.35 | -6.05 |
| GEBV-GD-100 | 77.16 | 51.53 | -25.63 | 66.46 | 59.68 | -6.78 |
| | PH | | | PNPP | | |
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 83.77 | 42.52 | -41.25 | 3.93 | 4.12 | 0.19 |
| GD-O-30 | 89.50 | 49.69 | -39.81 | 3.86 | 4.19 | 0.33 |
| GD-O-50 | 90.11 | 50.13 | -39.98 | 3.80 | 4.14 | 0.34 |
| GD-O-100 | 92.10 | 52.10 | -40.00 | 3.64 | 4.08 | 0.44 |
| GEBV-GD-30 | 87.26 | 42.99 | -44.27 | 3.90 | 4.22 | 0.32 |
| GEBV-GD-50 | 87.95 | 43.50 | -44.45 | 3.84 | 4.21 | 0.37 |
| GEBV-GD-100 | 89.27 | 43.95 | -45.32 | 3.70 | 4.17 | 0.47 |

552

553

554 (i) $\overline{GEBV}_P$: the GEBV average among the 10 selected parental lines. $\overline{GEBV}_{F_{10}}$ : the

555 GEBV average among the resulting 2700 $F_{10}$ RILs.

556 (ii) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs;

557 GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen

558 from the candidate sets composed of the top 30, 50, and 100 accessions, respectively;

559 GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen

560 from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions,

561 respectively, which have the maximal D-scores.

562 (iii) BRSW: brown rice seed width; FPP: florets per panicle; FTAA: flowering time at

563 Arkansas; FTAF: flowering time at Faridpur; PH: plant height; PNPP: panicle number

564 per plant.

565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585

586
587
588
589
590
591
592
593
594
595
596

597 Table 5: The average of genetic gains from the 30 repetitions for Dataset II.

598
599

| | YLD | | |
|---|---|---|---|
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 5571.61 | 6468.60 | 896.99 |
| GD-O-30 | 5452.39 | 6488.02 | 1035.63 |
| GD-O-50 | 5436.58 | 6484.58 | 1048.00 |
| GD-O-100 | 5289.74 | 6540.72 | 1250.98 |
| GEBV-GD-30 | 5538.44 | 6501.23 | 962.79 |
| GEBV-GD-50 | 5522.45 | 6482.13 | 959.68 |
| GEBV-GD-100 | 5454.37 | 6535.79 | 1081.42 |
| | PH | | |
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 97.75 | 85.89 | -11.86 |
| GD-O-30 | 102.20 | 87.59 | -14.61 |
| GD-O-50 | 103.66 | 89.99 | -13.67 |
| GD-O-100 | 106.83 | 91.85 | -14.98 |
| GEBV-GD-30 | 99.00 | 86.01 | -12.99 |
| GEBV-GD-50 | 99.39 | 85.99 | -13.40 |
| GEBV-GD-100 | 101.15 | 86.13 | -15.02 |
| | FT | | |
| | $\overline{\text{GEBV}}_P$ | $\overline{\text{GEBV}}_{F_{10}}$ | genetic gain |
| GEBV-O | 83.14 | 77.84 | -5.30 |
| GD-O-30 | 83.98 | 78.73 | -5.25 |
| GD-O-50 | 84.57 | 78.19 | -6.38 |
| GD-O-100 | 85.62 | 78.39 | -7.23 |
| GEBV-GD-30 | 83.44 | 77.90 | -5.54 |
| GEBV-GD-50 | 83.69 | 77.76 | -5.93 |
| GEBV-GD-100 | 84.16 | 77.89 | -6.27 |

600

601 (i) $\overline{GEBV}_P$: the GEBV average among the 10 selected parental lines. $\overline{GEBV}_{F_{10}}$ : the

24

602    GEBV average among the resulting 2700 $F_{10}$ RILs.

603    (ii) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs;

604    GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen

605    from the candidate sets composed of the top 30, 50, and 100 accessions, respectively;

606    GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen

607    from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions,

608    respectively, which have the maximal D-scores.

609    (iii) YLD: yield; PH: plant height; FT: flowering time.

610
611
612
613
614
615
616
617
618
619
620

Figure 1: The working flow for the Monte Carlo simulation.

GEBV: genomic estimated breeding value; GBLUP: genomic best linear unbiased predictor; RIL: recombinant inbred line.

BRSW (the larger the better)

FPP (the larger the better)

626

FTAA (the smaller the better)

FTAF (the smaller the better)

627

PH (the smaller the better)

PNPP (the larger the better)

628

27

629 Figure 2: The GEBV averages for the best individuals from the 30 repetitions at
630 consecutive generations for the six chosen traits in Dataset I.

631 (i) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs;
632 GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen
633 from the candidate sets composed of the top 30, 50, and 100 accessions, respectively;
634 GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen
635 from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions,
636 respectively, which have the maximal D-scores.

637 (ii) BRSW: brown rice seed width; FPP: florets per panicle; FTAA: flowering time at
638 Arkansas; FTAF: flowering time at Faridpur; PH: plant height; PNPP: panicle number
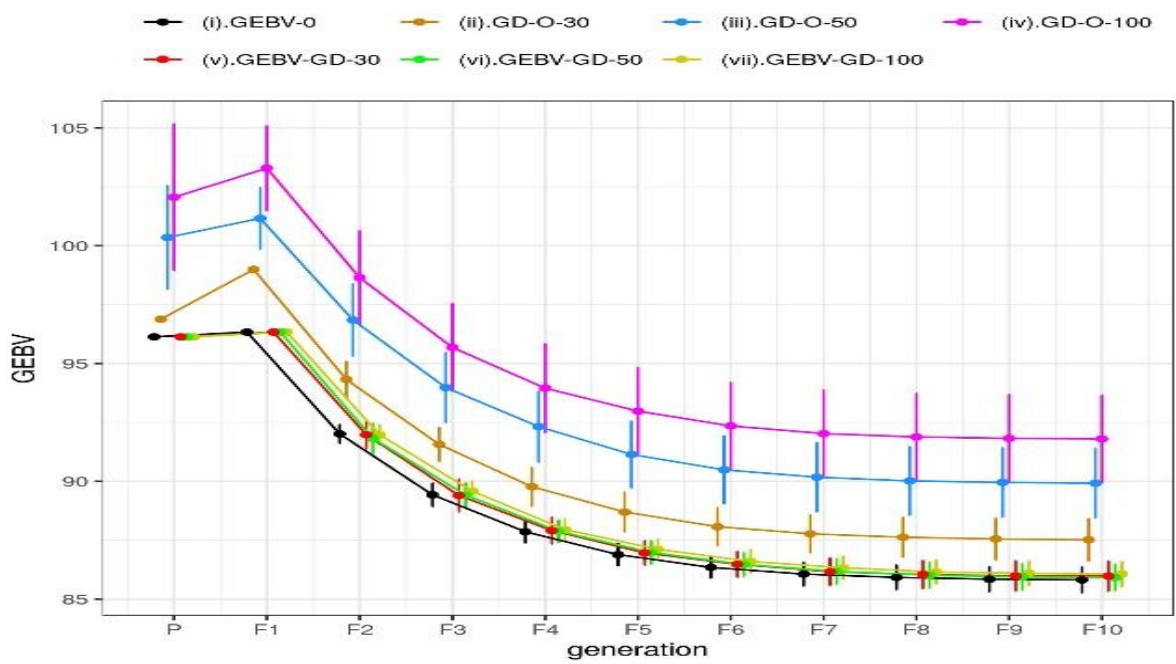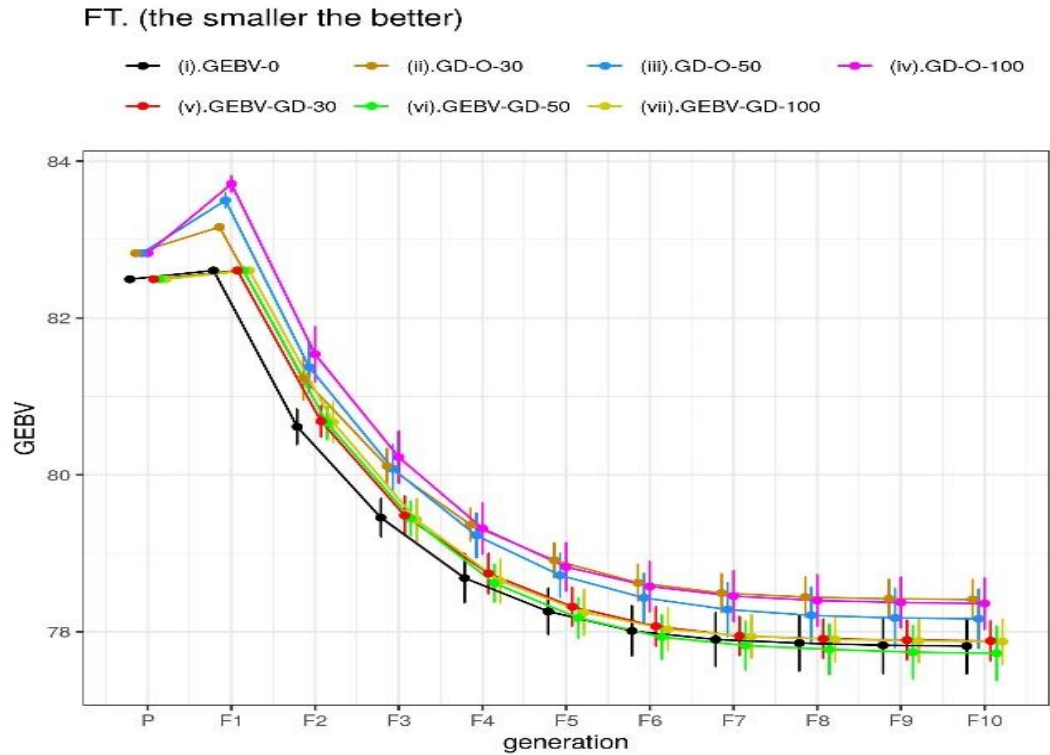639 per plant.

640

641

642

643

YLD. (the larger the better)

644

645



PH. (the smaller the better)

646

FT. (the smaller the better)

647

Figure 3: The GEBV averages for the best individuals form the 30 repetitions at consecutive generations for the three target traits in Dataset II.
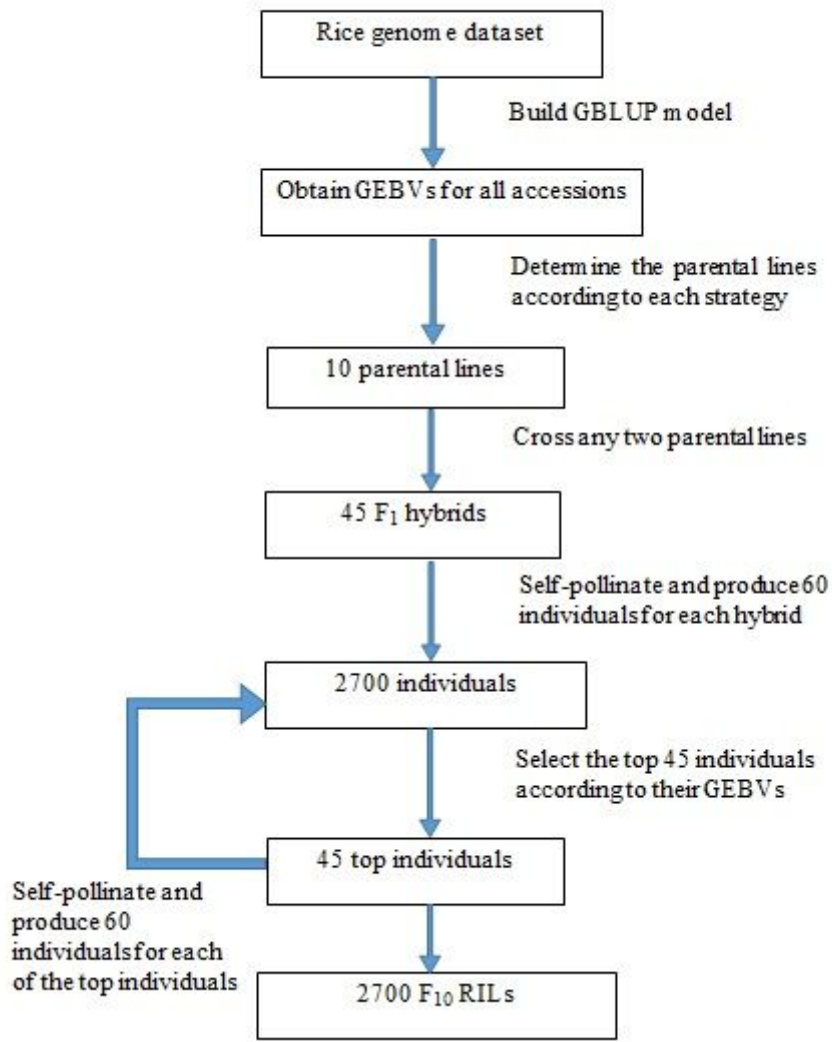
(i) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs; GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen from the candidate sets composed of the top 30, 50, and 100 accessions, respectively; GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions, respectively, which have the maximal D-scores.

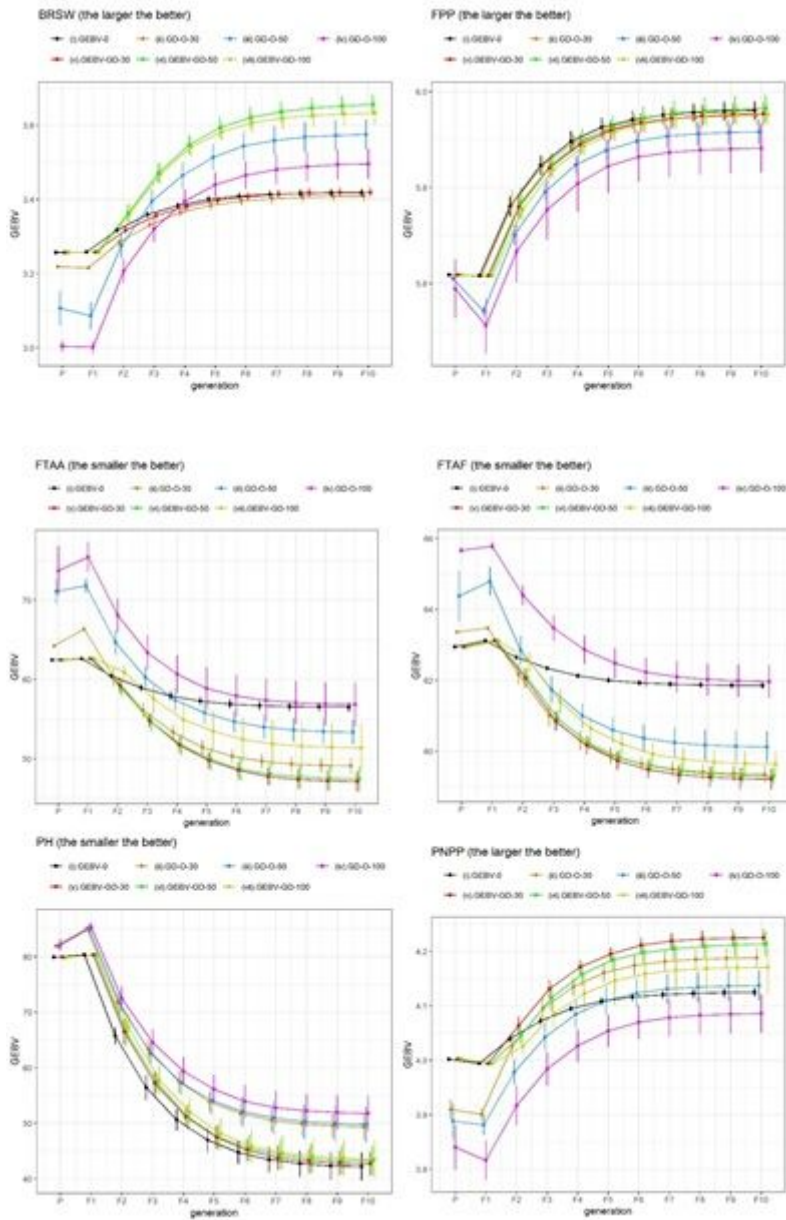(ii) YLD, yield; PH, plant height; FT, flowering time.

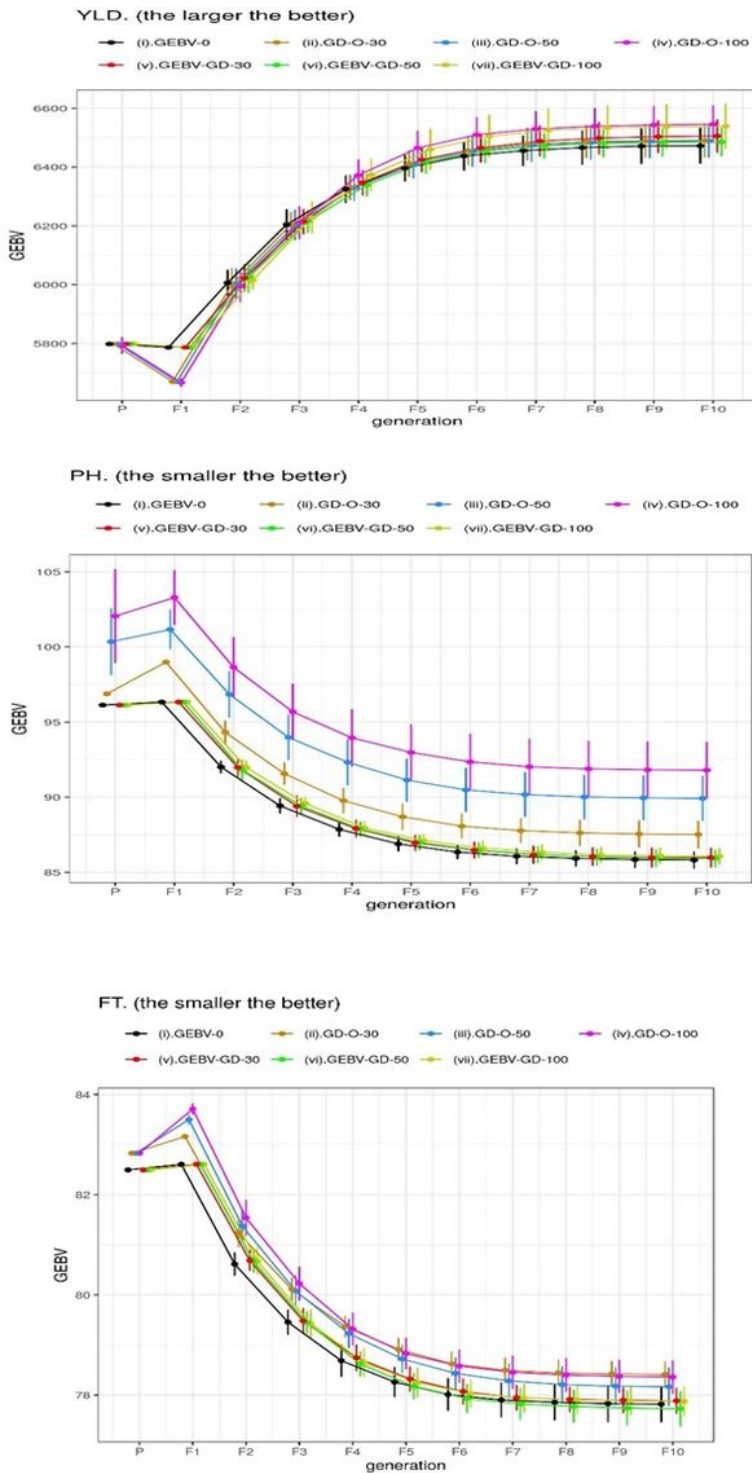657

658

# Figures



**Figure 1**

The working flow for the Monte Carlo simulation. GEBV: genomic estimated breeding value; GBLUP: genomic best linear unbiased predictor; RIL: recombinant inbred line.

**Figure 2**

The GEBV averages for the best individuals from the 30 repetitions at consecutive generations for the six chosen traits in Dataset I. (i) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs; GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen from the candidate sets composed of the top 30, 50, and 100 accessions, respectively; GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen from the remainder of the candidate sets composed of the top 30, 50, and 100 accessions, respectively, which have the maximal D-scores. (ii) BRSW: brown rice seed width; FPP: florets per panicle; FTAA: flowering time at Arkansas; FTAF: flowering time at Faridpur; PH: plant height; PNPP: panicle number per plant.

## Figure 3

The GEBV averages for the best individuals form the 30 repetitions at consecutive generations for the three target traits in Dataset II. (i) GEBV-O: the subset of the top 10 accessions with the minimal or maximal GEBVs; GD-O-30, -50, -100: the subsets of 10 accessions with the maximal D-scores chosen from the candidate sets composed of the top 30, 50, and 100 accessions, respectively; GEBV-GD-30, -50, -100: the subsets of the top 2 accessions plus 8 accessions chosen from the remainder of the candidate

sets composed of the top 30, 50, and 100 accessions, respectively, which have the maximal D-scores. (ii) YLD, yield; PH, plant height; FT, flowering time.