

# Genotyping marker density reduction is not an effective approach in long-term prediction-based breeding of cross-pollinated crops

**Julio Cesar DoVale** (✉ [juliodovale@ufc.br](mailto:juliodovale@ufc.br))

Federal University of Ceara: Universidade Federal do Ceara <https://orcid.org/0000-0002-3497-9793>

**Humberto Fanelli Carvalho**

Universidad Politecnica de Madrid

**Felipe Sabadin**

Virginia Tech: Virginia Polytechnic Institute and State University

**Roberto Fritsche-Neto**

International Rice Research Institute

---

## Research Article

**Keywords:** Linkage disequilibrium, genomic recurrent selection, marker subsets, additive-dominant model, SCA-GBLUP

**Posted Date:** November 12th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-1005929/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 Genotyping marker density reduction is not an effective approach in long-term prediction-based breeding of  
2 cross-pollinated crops

<sup>4</sup> Júlio César Do Vale<sup>a\*</sup>, Humberto Fanelli Carvalho<sup>b</sup>, Felipe Sabadin<sup>c</sup>, Roberto Fritsche-Neto<sup>d</sup>

5           <sup>a</sup>Department of Crop Science, Federal University of Ceará, Fortaleza, Ceará, Brazil  
6           <sup>b</sup>Universidad Politecnica de Madrid, Madrid, Spain  
7           <sup>c</sup>Virginia Tech: Virginia Polytechnic Institute and State University, Blacksburg, United States

11 \*Corresponding author

12 Júlio César DoVale  
13 email: [juliodovale@ufc.br](mailto:juliodovale@ufc.br)  
14  
15 **ORCID of the authors:**  
16 Julio César DoVale (JCD)  
17 Humberto Fanelli Carvalho  
18 Felipe Sabadin (FC): 0000-0003-1839-638X  
19 Roberto Fritzsche-Neto (RFN): 0000-0002-1349-464X

20 ABSTRACT

21

22 Reductions of genotyping marker density have been extensively evaluated as potential strategies to reduce the genotyping  
23 costs of genomic selection (GS). Low-density marker panels are appealing in GS because they entail lower  
24 multicollinearity and computational time-consumption and allow more individuals to be genotyped for the same cost.  
25 However, statistical models used in GS are usually evaluated with empirical data, using "static" training sets and  
26 populations. This may be adequate for making predictions during a breeding program's initial cycles, but not for the long  
27 term. Moreover, to the best of our knowledge, no GS models consider the effect of dominance, which is particularly  
28 important for breeding outcomes in cross-pollinated crops. Hence, dominance effects are an important and unexplored  
29 issue in GS for long-term programs involving allogamous species. To address it, we employed two approaches: analysis  
30 of empirical maize datasets and simulations of long-term breeding applying phenotypic and genomic recurrent selection  
31 (intrapopulation and reciprocal schemes). In both schemes, we simulated twenty breeding cycles and assessed the effect  
32 of marker density reduction on the population mean, the best crosses, additive variance, selective accuracy, and response  
33 to selection with models (additive, additive-dominant, general (GCA), and specific combining ability (SCA)). Our results  
34 indicate that marker reduction based on linkage disequilibrium levels provides useful predictions only within a cycle, as  
35 accuracy significantly decreases over cycles. In the long-term, high-marker density provides the best responses to  
36 selection. The model to be used depends on the breeding scheme: additive for intrapopulation and additive-dominant or  
37 SCA for reciprocal.

38

39 **Keywords:** Linkage disequilibrium; genomic recurrent selection; marker subsets; additive-dominant model; SCA-  
40 GBLUP;

41 **DECLARATIONS**

42

43 **Funding:** Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) – Process 104371/2019-6.

44

45 **Conflicts of Interest:** On behalf of all authors, the corresponding author states that there is no conflict of interest.

46

47 **Codes and datas availability:** codes and supporting data is available in:

48 <https://data.mendeley.com/datasets/96p3khznj/1>

49

50 **Ethical standards:** Not applicable.

51

52 **AUTHOR CONTRIBUTION STATEMENT**

53 JCD wrote the manuscript, discussion, contributing to ideas and graphs. HFC and FS contributed to the analysis of  
54 empirical data, in writing, mainly discussion. RFN elaborated on the hypothesis, conducted the analysis of the simulated  
55 experiments, interpreted the results, and contributed to the writing. All authors read and approved the final manuscript.

56

57 **ACKNOWLEDGMENTS**

58 To the Queiroz College of Agriculture (University of São Paulo) and the Federal University of Ceará.

59

60 **KEY MESSAGE**

61 In genomic recurrent selection, the more markers, the better because they buffer the linkage disequilibrium losses caused  
62 by recombination over cycles, and consequently, provide higher responses to selection.

63 INTRODUCTION

64

65 Genomic selection (GS), using high-density, single-nucleotide polymorphism (SNP) chips, has been widely  
66 adopted in breeding programs to improve predictive ability and response to selection (Hayes *et al.*, 2009; Crossa *et al.*,  
67 2010). However, these chips' genotyping cost is currently prohibitive, especially for low profitability species (Hou *et al.*  
68 2020), and programs in poor regions of the world, such as those in Latin America, Africa, and Southwest Asia.

69 In this context, marker density reduction has been extensively evaluated as an alternative to reduce the  
70 genotyping costs of GS. This reduction can be carried out based on different criteria, for instance, marker effects (e Sousa  
71 *et al.* 2019), trait genetic architecture (Zhang *et al.* 2015), haplotype block analysis (Ma *et al.* 2016), genomewide  
72 association studies (GWAS) (Subedi *et al.* 2013), and linkage disequilibrium (LD) (Al-Tobasei *et al.*, 2021). Of course,  
73 all strategies have advantages and disadvantages. However, the LD criterion seems to be the more suitable strategy  
74 because it eliminates markers for redundancy that are highly correlated without interfering with marker effects and  
75 thereby reduces multicollinearity (an obstacle for highly dense SNP panels). Moreover, it avoids marker effect noise (Xu  
76 2013), and reduces the probability of overfitting by having the marker reduction process be independent of the phenotype  
77 (Vallejo *et al.*, 2018).

78 The use of low-density marker panels in breeding programs is also appealing due to a reduction in computational  
79 time-consumption that allows more individuals to be genotyped for the same cost (Gorjanc *et al.* 2015). Furthermore,  
80 with larger datasets of phenotyped and genotyped individuals, GS' predictive ability may be increasingly driven by LD  
81 rather than by linkage information (Hickey *et al.* 2014). Furthermore, several studies claim that it is possible to  
82 substantially reduce the number of markers while maintaining high predictive ability (Tayeh *et al.* 2015; Ma *et al.* 2016;  
83 e Sousa *et al.* 2019; Al-Tobasei *et al.* 2020). Thus, the use of genotyping techniques with highly discriminative markers  
84 could be a viable alternative for current, low-cost SNP array strategies, with the potential to increase the fraction of the  
85 genome captured in a cost-efficient manner (Elshire *et al.* 2011).

86 The studies mentioned above made it possible to obtain an informative marker subset with either empirical or  
87 simulated data. The results of this were often satisfactory, as they considered the LD within a generation in a fixed  
88 scenario, i.e., used fixed training sets (TS) in the early stages of a breeding program to infer across multiple breeding  
89 cycles. However, both selection and recombination events that occur during breeding cycles may change LD between  
90 markers and Quantitative Trait Loci (QTLs) (in the latter case, due to cross-over occurrences), especially in more complex  
91 traits (Xu 2013). As a result, there are changes in allelic frequencies and, consequently, allele substitution effects which  
92 may even cause allele fixation (Walsh and Lynch 2018).

93 Although previous studies, like those mentioned above, evaluated the effect of marker density on the  
94 performance of genomic prediction models, there are no reports of its effect on models that incorporate additive as well

95 as dominance effects in cross-pollinated crops under recurrent selection (RS) schemes. The vast majority of published  
96 papers so far have considered only additive models for genomic predictions, which may not represent the genetic  
97 complexity of all crops. Even though only additive effects are transmitted over generations, the *per se* performance of a  
98 population and its single-crosses depends on dominance deviations (Falconer and Mackay, 1996)., This is the basis of  
99 heterosis and is crucial in the expression of agronomic traits of cross-pollinated crops. (Bernardo 2010). Additionally, as  
100 estimated by marker effects, Galli *et al.* (2020) observed negative covariance in genetic values of parental lines versus  
101 their single-crosses. This reinforces the assertion that poor predictions are made about single-crosses' performance when  
102 they are solely based on the additive effects of the parents' genome.

103 Simulations are a suitable tool to assess the effect of marker density reduction in a long-term breeding program  
104 applying GS. Several factors can be controlled to make inferences on given genetic parameters over several breeding  
105 cycles in a fast, inexpensive and consistent way (Dai *et al.* 2020). Simulations are essential for breeding programs of  
106 cross-pollinated crops. These usually adopt RS-based methods for obtaining cultivars, either by intrapopulation recurrent  
107 selection (IRS) or reciprocal recurrent selection (RRS), and hence genetic progress is achieved very slowly over breeding  
108 cycles. Several published studies address data simulation in the conducted long-term breeding programs, but with  
109 different objectives. Some sought to evaluate methods for updating of training sets (e.g., Neyhart *et al.*, 2016), others  
110 using RS to optimize GS (Muleta *et al.*, 2019) and even relate the number of parents in the persistence of accuracy and  
111 genetic gains (Muller *et al.*, 2017). However, a study involving IRS and RRS schemes, which are widely used and allow  
112 a closer approach to the daily life of cross-pollinated crops breeding programs, has not yet been carried out with  
113 information genomic.

114 The flexibility of the models and methods used in GS allows for the integration of allelic interaction effects, as  
115 the combining abilities of inbred parental lines (Werner *et al.* 2018). Furthermore, Reif *et al.* (2013) observed an increase  
116 in their model's predictive capacity when incorporating general combining ability (GCA) effects. Despite this, no studies  
117 of long-term breeding programs compare the superiority of models that incorporate specific combining ability (SCA)  
118 effects using the genomic best linear unbiased predictor (G-BLUP) model, including additive and dominance effects.  
119 Thus, our objectives were to: (i) assess the impact of marker density reduction on prediction accuracy and responses to  
120 selection in long-term breeding programs, (ii) assess whether the inclusion of dominance effects affects the main genetic  
121 parameters, and (iii) to compare genomic prediction models using combining ability, additive and additive-dominance  
122 effects over many cycles.

123 **MATERIALS AND METHODS**

124

125 To achieve our objectives, we worked with empirical data and simulation experiments.

126

127 ***Empirical data***

128

129 Overview: we used genotyping chips to perform GS for grain yield in maize and evaluated the effects of different marker  
130 reduction strategies on GS accuracy.

131

132 ***Phenotypic data***

133

134 Our study considered grain yield (GY, in Mg ha<sup>-1</sup>) from two datasets of maize single-crosses carried out by the  
135 Helix Seeds Company (HEL) and the University of São Paulo (USP). The plots were manually harvested, and GY was  
136 corrected to 13% moisture. Each dataset was analyzed independently.137 The HEL data were obtained from 452 single-crosses obtained from 111 inbred lines. The hybrids were  
138 evaluated in a randomized complete block design in five Brazilian sites during the 2015 season.139 The USP data were composed of 903 maize single-crosses, obtained from a full diallel of 49 inbred lines. The  
140 single-crosses were evaluated in 2016 and 2017 at two sites (Piracicaba and Anhumas) in São Paulo State, Brazil. Field  
141 trials were performed in an augmented block design, each block consisting of 16 unique single-crosses and two  
142 commercial hybrids as checks. Single-crosses were evaluated at each site and year under two nitrogen fertilization levels  
143 (ideal and low N). These eight treatment combinations comprise what we herein call "environments". More details on  
144 the experimental design and cultivation practices for HEL and USP datasets can be found in Sousa *et al.* (2017) and Galli  
145 *et al.* (2020), respectively.

146

147 ***Genetic-statistical models to estimate BLUEs***

148

149 Mixed model equations were used for the statistical analysis of single-cross phenotypes in each dataset. First,  
150 the joint analysis of each phenotype was performed to estimate genotype means across environments. Next, we fitted a  
151 mixed model to obtain the Best Linear Unbiased Estimator (BLUE) for each genotype's yield. We then estimated the  
152 adjusted means in each environment via the *breedR* package (Muñoz and Rodriguez, 2016) in R software (R Development  
153 Core Team, 2019). The analyses for the USP dataset were performed using this model:

154

$$\mathbf{y} = \mathbf{Ql} + \mathbf{Sb} + \mathbf{Tc} + \mathbf{Ug} + \mathbf{Vi} + \boldsymbol{\varepsilon}$$

155 where  $\mathbf{y}$  is the vector of phenotypic values of single-crosses and checks;  $\mathbf{l}$  is the vector of fixed effects of the environment  
156 (combination of site  $\times$  year  $\times$  N level);  $\mathbf{b}$  is the vector of random effect of block nested within environment, where  $\mathbf{b} \sim$   
157  $N(\mathbf{0}, \mathbf{I}\sigma_b^2)$ ;  $\mathbf{c}$  is the vector of fixed effects of checks;  $\mathbf{g}$  is the vector of fixed effects of single-crosses;  $\mathbf{i}$  is the vector of  
158 fixed effects of interaction checks  $\times$  environments;  $\boldsymbol{\varepsilon}$  is the vector of random residual effects, which are confounded in  
159 the final residual term, where  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{D}_e)$ .  $\mathbf{Q}$ ,  $\mathbf{S}$ ,  $\mathbf{T}$ ,  $\mathbf{U}$ , and  $\mathbf{V}$  are the incidence matrices for  $\mathbf{l}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{g}$ , and  $\mathbf{i}$ . We assumed  
160 an unstructured covariance matrix across environments for the residual term ( $\mathbf{D}_e$ ).

161 For the phenotypic analysis of the HEL dataset, we used the following reduced model:

$$162 \quad \mathbf{y} = \mathbf{Xr} + \mathbf{Tg} + \mathbf{Bs} + \mathbf{Hx} + \boldsymbol{\varepsilon}$$

163 where  $\mathbf{r}$  is the vector of block effect considered as random, where  $\mathbf{r} \sim N(\mathbf{0}, \mathbf{I}\sigma_r^2)$ ;  $\mathbf{s}$  is the vector of random effects of the  
164 environment (sites), where  $\mathbf{s} \sim N(\mathbf{0}, \mathbf{I}\sigma_s^2)$ ;  $\mathbf{x}$  is the vector of the random effects of single-crosses by environment  
165 interaction, where  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{I}\sigma_x^2)$ ; and  $\boldsymbol{\varepsilon}$  is the vector of error, where  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}\sigma_\varepsilon^2)$ .  $\mathbf{X}$ ,  $\mathbf{T}$ ,  $\mathbf{B}$ , and  $\mathbf{H}$  are incidence matrices  
166 for  $\mathbf{r}$ ,  $\mathbf{g}$ ,  $\mathbf{s}$ ,  $\mathbf{x}$ .

167 In both HEL and USP models, single-crosses' effects were considered random in order to obtain the variance  
168 components via the restricted maximum likelihood (REML/BLUP) procedure (Henderson 1975).

169

170 *Genotypic data*

171

172 All parental inbred lines (USP and HEL) were genotyped using the Affymetrix® Axiom Maize Genotyping array  
173 containing 616 K SNPs (Unterseer et al. 2014). As a standard quality control, markers with a low call rate (< 90%) and  
174 those found to be non-biallelic were removed. Missing data were imputed with the Beagle 5.0 algorithm (Browning et  
175 al. 2018). Next, loci showing heterozygosity in one or more individuals were removed. Single-cross genotypes were  
176 obtained *in silico* by genomic information from parental inbred lines. Finally, markers with minor allele frequency (MAF)  
177 < 0.05 were removed from the hybrid genomic matrix, resulting in 63,104 and 30,467 SNP markers for the USP and HEL  
178 datasets.

179

180 *Marker reduction*

181

182 Different SNP marker subsets were obtained via the following two approaches: (i) removing SNP markers based  
183 on the LD between them and (ii) applying the LA-GA-T genetic algorithm (Akdemir (2017) in the *STPGA R* package,  
184 which aims to select optimized training sets (OTS).

185 In the first approach, pairwise LD was calculated as the squared allele frequencies correlation ( $r^2$ ). In order to  
186 remove redundant markers, we applied the following threshold  $r^2$  values, 0.99, 0.90, 0.80, 0.70, 0.60, 0.50, 0.40, 0.30,  
187 0.20, 0.10, and 0.01, using the *SNPRelate* package (Zheng et al., 2012). After performing quality control, we used the  
188 full SNP markers (M\_Full) matrix as the benchmark scenario for downstream comparisons.

189 In the second approach, we applied a Singular Value Decomposition (SVD) using the number of eigenvalues  
190 that explained 98% of the variance, according to Pocrník *et al.* (2016). To reduce computation time and represent only  
191 the LD among markers' physical components, we applied SVD + OTS within each chromosome. This procedure resulted  
192 in thirteen SNP matrices for use in genomic prediction models.

193

194 *Genomic predictions and model comparisons*

195

196 We used the thirteen SNP matrices to build genomic relationship matrices via the *snpReady* package (Granato  
197 et al. 2018). Then, we applied additive and additive-dominance G-BLUP models to perform genomic predictions of grain  
198 yield using the following equations, respectively:

199 
$$\hat{\mathbf{g}} = \mathbf{1}\mu + \mathbf{Z}_a\mathbf{a} + \boldsymbol{\varepsilon}$$

200 
$$\hat{\mathbf{g}} = \mathbf{1}\mu + \mathbf{Z}_a\mathbf{a} + \mathbf{Z}_d\mathbf{d} + \boldsymbol{\varepsilon}$$

201 where  $\hat{\mathbf{g}}$  is the vector of adjusted environmental grain yield means of single-crosses described in the 'Phenotypic data'  
202 subsection;  $\mu$  is the mean (intercept);  $\mathbf{a}$  is the vector of additive genetic effects of the individuals, where  $\mathbf{a} \sim N(\mathbf{0}, \mathbf{G}_a\sigma_a^2)$ ;  
203  $\mathbf{d}$  is the vector of dominance effects, where  $\mathbf{d} \sim N(\mathbf{0}, \mathbf{G}_d\sigma_d^2)$ ; and  $\boldsymbol{\varepsilon}$  is the vector of random residuals  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I}\sigma_\varepsilon^2)$ .  $\mathbf{1}$   
204 is the vector of ones;  $\mathbf{Z}_a$  and  $\mathbf{Z}_d$  are incidence matrices for  $\mathbf{a}$  and  $\mathbf{d}$ .  $\sigma_a^2$  is the genomic additive variance,  $\sigma_d^2$  is the genomic  
205 dominance variance, and  $\sigma_\varepsilon^2$  is the residual variance.  $\mathbf{G}_a$  and  $\mathbf{G}_d$  are the additive and dominance genomic relationship  
206 matrices, following the equations:  $\mathbf{G}_a = \mathbf{W}_A \mathbf{W}'_A / 2 \sum_{i=1}^n p_i (1 - p_i)$  and  $\mathbf{G}_d = \mathbf{W}_D \mathbf{W}'_D / 4 \sum_{i=1}^n (p_i(1 - p_i))^2$ , where  $p_i$  is  
207 the frequency of one allele of the locus  $i$  and  $\mathbf{W}$  is the marker incidence matrix (VanRaden 2008). Genotypes in the  $\mathbf{W}_A$   
208 matrix were coded as 0 for homozygote A<sub>1</sub>A<sub>1</sub>, 1 for the heterozygote A<sub>1</sub>A<sub>2</sub>, and 2 for the homozygote A<sub>2</sub>A<sub>2</sub>. In the  $\mathbf{W}_D$   
209 matrix, genotypes were coded as 0 for both homozygotes and 1 for heterozygotes. The genomic prediction models were  
210 performed using the *sommer* package (Covarrubias-Pazaran 2016).

211 In order to compare marker reduction scenarios, we used model accuracy, effective population size ( $N_e$ ), and  
212 several metrics of the genetic relationship among individuals. We obtained estimates of  $N_e$  based on the  $\mathbf{G}_a$  matrix.  
213 Correlation among  $\mathbf{G}_a$  matrices across scenarios were assessed with Mantel tests (*vegan* R package, (Oksanen et al.  
214 2019)), as was the correlation among  $\mathbf{G}_d$  ones. We used  $\mathbf{G}_a$  and  $\mathbf{G}_d$  calculated from the M\_Full SNP marker matrix to  
215 evaluate clustering of single-crosses based on Euclidean distance. The *superheat* R package created genetic distance

216 heatmaps (Barter and Yu 2018) using M\_Full, M\_0.5, M\_0.01, and M\_OTS. Subsequently, the singles-crosses order was  
217 used to sort heatmaps M\_0.5, M\_0.01, and M\_OTS.

218       The predictive ability of the genomic prediction models was estimated as the correlation between the predicted  
219 and the observed genotypic values via a cross-validated, alpha-based design (CV- $\alpha$ , Yassue *et al.*, 2020), which is an  
220 extension of the methodology presented by Shao (1993). This model validation procedure consists of assigning  
221 observations to train or validation folds (groups) in each replication using alpha-lattice sorting premises and evaluating  
222 the model's prediction in each assignment iteration performed. CV- $\alpha$  was developed in a way that allows for the creation  
223 of validation scenarios with two, three, or four replicates per group, regardless of the number of treatments. The number  
224 of groups is determined by the percentage of training and validation sets. Our assignments of observations to groups were  
225 based on the  $\alpha(0,1)$  lattice design in order to reduce the concurrences of any two treatments in the same group across the  
226 replicates (Patterson and Williams 1976). Here we used a validation scenario with four replicates and five groups (i.e.,  
227 5-fold validation), and model predictive abilities were averaged across the different assignments of replicates to those  
228 groups.

229

230       ***Simulation experiments***

231

232       Overview: We simulated long-term breeding programs that apply intrapopulation (IRS) and reciprocal (RRS) recurrent  
233 selection schemes-based, using genomic (GRS) and phenotypic recurrent selection (PRS). Different marker densities and  
234 model types were applied.

235

236       *Characteristics of the founder population, trait genetic architecture, and genotyping chips*

237

238       We simulated a population of maize single-crosses from inbred parental lines to perform phenotypic and  
239 genomic predictions for a quantitative trait. For this, we used the *AlphaSimR* package (Gaynor *et al.* 2020). The scripts  
240 can be found in the supplementary files.

241       A historical population founded by 1,000 unique individuals was simulated stochastically (i.e., without artificial  
242 selection), with genomic information consisting of ten chromosomes pairs each using a Markovian Coalescent Simulator  
243 (MaCS) (Chen *et al.* 2009). Thirty-one thousand biallelic segregating sites were considered uniformly distributed across  
244 chromosomes, and 1,000 segregating loci were randomly sampled as quantitative trait nucleotides (QTN). The founders  
245 were diploidized but not inbred. From the historical population, base populations were obtained to perform recurrent  
246 selection (RS) schemes on, as described below.

247        The target of the simulation was a quantitative trait, such as yield. The genetic parameters obtained by Hallauer  
248        *et al.* (2010) were used. Each QTN received randomly allocated additive and dominance effects. Genetic values for each  
249        genotype were obtained by summing all additive and dominance effects for all QTN. Additive effects ( $a$ ) were sampled  
250        of a gamma distribution with scale and shape parameters equal to 1 and randomly assigned for each QTN. Similarly,  
251        dominance effects ( $d$ ) for each QTN were computed by multiplying the absolute value of its additive effect ( $|a_i|$ ) by locus-  
252        specific dominance degree ( $\delta_i$ ). Dominance degrees were sampled of a Gaussian distribution with  $\delta_i \sim N(\mu_\delta, \sigma_\delta^2)$ , where  
253         $\mu_\delta$  is the average dominance degree equal to 0.50 and  $\sigma_\delta^2$  is the dominance variance equal to 1. Dominance effects were  
254        assigned for each QTN according to the equation below:

255        
$$d_i = \begin{cases} 0, & \text{if QTN is homozygous} \\ \delta_i \times |a_i|, & \text{if QTN is heterozygous} \end{cases}$$

256        Phenotypic values were obtained by adding a random error sampled of a Gaussian distribution with mean equal  
257        to 0 and variance ( $\sigma_e^2$ ) equal to 1, which was defined by broad-sense ( $H^2 = 0.50$ ) and narrow ( $h^2 = 0.33$ ) heritabilities.

258        In order to represent the genotyping chips, we used three marker densities sampled from the segregating loci as  
259        SNP: the first was 30K SNPs, which on average is roughly equivalent to 15 markers per cM; the second was 1,5K SNPs  
260        equivalent to 1 marker per cM; and the third, 0,75K SNPs equivalent to one marker for every 2 cM.

261

262        *Breeding schemes: base population and burn-in phase*

263

264        We used IRS and RRS, two main breeding schemes RS-based, to evaluate the effect of marker reduction on  
265        predictions over twenty cycles with steps conducted independently. To create the base population for the IRS scheme,  
266        we selected 15 parents from the historical population (Fig. 1) and built a full diallel to generate the base population (Cycle  
267        0, C<sub>0</sub>). For the RRS scheme, the base population was made similarly, but in this case, 100 parents were used for the  
268        diallel. Parents of the best single-crosses were selected to build two heterotic groups (HG). Regardless of the breeding  
269        scheme, stage, or cycle, the progeny of 50 individuals was considered (equivalent to one trial with 2 plots of 25 plants).

270

271

272        **Fig. 1** Schematic representation of the recurrent selection breeding schemes applied in the simulations. After obtaining  
273        the historical population, base populations were selected from it randomly, and these served as the basis for the IRS  
274        (orange) and RRS (blue) schemes. Thus, information in black applies to both schemes.

275

276        Once the base populations were created, a burn-in stage followed. For the burn-in stage (not represented in detail  
277        in Figure 1), we used three traditional recurrent selection cycles, as described below:

278 **IRS** – 200 half-sib progenies were obtained from the base population; then, the top 30 progenies were selected, and  
279 random recombination of 1,500 self-fertilized individuals was simulated ( $S_1$ ) from those 30 best parents; using the  
280 resulting population, the process was repeated two more times, finally resulting in the  $C_0$  population.

281 **RRS** – 200 full-sib offspring were obtained through a crossing-plan with 200 unique parents from each HG. Then, the  
282 best single-crosses were selected, and we proceeded with the recombination ( $S_1$ ) of the parents within HG. At the end of  
283 each cycle, the fixation index (Fst) was calculated to assess the genetic distance between the HGs.

284

285 *Breeding cycles: genomic and phenotypic selection*

286

287 After the three cycles of the burn-in phase, 3,000 individuals were randomly selected as each scheme's  $C_0$  (Fig.  
288 1). From those 3,000 individuals, 2,000 were assigned to the training set (TS) for genomic predictions. First, marker  
289 effects were estimated using the three different SNP-chip densities described earlier. For that, genotyping and  
290 phenotyping information was used. Then, these 3,000 individuals were also used as population  $C_0$  to begin the subsequent  
291 stage, consisting of twenty cycles of RS (GRS or PRS).

292 In each cycle, for PRS simulations (in both schemes – IRS or RRS), breeding cycles were run as described in  
293 the burn-in phase. For GRS simulations, the genotype of the 200 parents allowed us to skip the generation of progeny,  
294 select directly based on genomic breeding values, and simulate planting only the selected individuals for recombination,  
295 closing one breeding cycle in just one growing season. Furthermore, based on the parents' genotype, we obtained *in*  
296 *silico* all 40,000 possible single-crosses genotypes (200 x 200), increasing the population size and the intensity of  
297 selection.

298

299 *Genomic predictions, methods, and model comparisons*

300

301 As in the equations mentioned in the 'Genomic prediction' subsection for empirical data, we used additive and  
302 additive-dominance G-BLUP models to perform genomic predictions using three SNP-chip densities in both breeding  
303 schemes. Additionally, for the RRS scheme, we modeled using combining ability effects. Thus, the genomic estimated  
304 value (GEV) of single-crosses was obtained using estimates of the general combining ability (GCA) and specific  
305 combining ability (SCA) of the parents, according to the following expression:

$$306 \quad GEV_{ij} = GCA_i + GCA_j + SCA_{i\times j}$$

307 where  $GEV_{ij}$  is the genomic estimated value of single-cross  $ij$ ,  $GCA_i$  is the general combining ability of parent  $i$ ,  $GCA_j$   
308 is the general combining ability of parent  $j$ , and  $SCA_{i\times j}$  is the specific combining ability from the cross between  $i$  and

309  $j$  parentals. Kronecker products were used to model the interaction between single-crosses' parental genomes to capture  
310 SCA effects (Acosta-Pech et al. 2017; Basnet et al. 2019).

311 We evaluated two methods at different levels for both IRS and RRS schemes. In the first, based on the GRS, we  
312 considered the three markers densities already mentioned (GRS\_30K, GRS\_1.5K, and GRS\_0.75K). In the second, based  
313 on PRS, used as benchmark strategies, we consider two selection ways: RS\_Traditional – traditional recurrent selection  
314 based on selection and recombination per cycle of the best phenotypically identified individuals, and RS\_Drift – recurrent  
315 selection also conducted with phenotypic data, but from randomly identified individuals per cycle. In this sense, we call  
316 for RS\_Traditional the phenotypic model while for RS\_Drift the random model. The methods within each breeding  
317 scheme were simulated over 20 breeding cycles and replicated 20 times.

318 The average genetic values (population mean), the performance of the best-crosses, Fst between heterotic  
319 groups, additive genetic variance, and prediction accuracy were calculated in each breeding cycle. The prediction  
320 accuracy for methods GRS-based was calculated as a Pearson correlation between true genetic values and GEV estimates.  
321 For methods PRS-based, the prediction accuracy was computed as the square root of heritability ( $h^2$ ). Furthermore, the  
322 response to selection was calculated by modified breeder's equation with other components that remained fixed between  
323 the studied methods:

$$324 \quad RS = \frac{PA}{t}$$

325 where  $RS$  is the relative response to selection to the previous cycle;  $PA$  is prediction accuracy; and  $t$  is the breeding cycle  
326 time. We consider that a PRS (RS\_Traditional or RS\_Drift) cycle takes 1.75 years on average to complete, while 0.35  
327 years is needed for GRS.

328    **RESULTS**

329

330    ***Marker density reduction effect - empirical data***

331

332        The number of markers in the SNP sub-sets decreased substantially with the decreasing LD level cut-offs applied  
333        to the empirical data (Table 1). In general, in both datasets, the reduction was steadier leading up to the subset with 10%  
334        LD between markers (M\_0.1) than after it. The matrix with all markers (M\_Full) had approximately 11.5 and 3.3 times  
335        more markers (USP and HEL datasets, respectively) than the one with only 1% LD between markers (M\_0.01). However,  
336        the subset obtained applying the algorithm to optimize training sets (M\_OTS) had drastically fewer markers than M\_0.01  
337        (3.9 and 5.4 times less, in USP and HEL, respectively).

338

339

340        **Table 1** Number of SNP-markers per chromosome as a function of the marker density reduction performed using different  
341        levels of linkage disequilibrium and the OTS algorithm. USP and HEL datasets.

342

343        Despite considerable marker reduction using LD level cut-offs, effective population sizes ( $N_e$ ) did not change  
344        in either dataset when considering the different SNP subsets (reduction scenarios) for estimating this parameter (Table  
345        2). In addition, the significant and high correlation coefficients among additive matrices ( $\mathbf{G}_a$ ) ( $> 0.98$  for both datasets)  
346        as well as among the dominance ( $\mathbf{G}_d$ ) matrices ( $> 0.73$  for USP and  $> 0.99$  for HEL) across marker densities (Tables S1  
347        and S2), indicates that the matrices with low-density markers are as representative as those with high density.  
348        Furthermore, comparing the genetic distance heatmaps across marker reduction scenarios, including the most extreme  
349        scenarios, shows that genetic relationships were not altered by marker density (Figs. S1 and S2).

350

351

352        **Table 2** Effective population size ( $N_e$ ) obtained using additive effect matrices ( $A$ ) with marker density reduction  
353        performed using different linkage disequilibrium levels and the OTS algorithm. USP and HEL datasets.

354

355        Similarly, results showed that the predictive ability remained unchanged throughout the various marker  
356        reduction scenarios, including the most extreme ones (Fig. 2). Therefore, with approximately 1,500 markers (1,408 for  
357        USP and 1,734 for HEL), it was possible to predict maize single-cross performance with the same accuracy as that  
358        achieved with all markers after the quality control process. Additionally, greater predictive abilities were observed when

359 we adopted the additive-dominance model, regardless of marker density. On average, the compound model accuracy  
360 surpassed the additive one by 9.8% and 19% for the USP and HEL datasets, respectively.

361

362

363 **Fig. 2** Predictive ability across marker density reduction scenarios based on different linkage disequilibrium levels and  
364 an OTS algorithm, using additive and additive-dominant models. **A** USP dataset **B** HEL dataset.

365

366 ***Marker density reduction effect - long-term IRS breeding simulations***

367

368 In simulations of long-term IRS breeding, RS\_Traditional surpassed that of GRS for the population mean  
369 (genetic mean), regardless of marker density (Fig. 3A). RS\_Traditional's superiority started to become evident in the third  
370 cycle. In relation to GRS, the density of 30K SNP had superior average performance compared to the other two densities,  
371 reaching almost double at the end of the twenty breeding cycles. We also observed that the additive model tended to be  
372 superior to the additive-dominance one for this parameter. The IRS scheme conducted by the drift method had the worst  
373 performance and, unlike the others, it declined over the cycles.

374

375

376 **Fig. 3** Simulation of 20 breeding cycles via PRS (Traditional and Drift) and GRS (with three marker densities) using the  
377 IRS breeding scheme with additive and additive-dominant G-BLUP models. **A** Population mean **B** Additive variance **C**  
378 Accuracy

379

380 The additive variance available for GRS increased up till the second breeding cycle and then consumed  
381 vertiginously, cycle after cycle, but at a lower rate than that observed for RS\_Traditional (Fig. 3B). Thus, the results of  
382 this genetic parameter were exactly opposite to those observed in the previous one (population mean). On the other hand,  
383 the trend of the curves for the accuracy parameter was very similar to those obtained for the population mean (Fig. 3C).  
384 The difference is that the GRS with the highest marker density started better than the RS\_Traditional. In general, after  
385 the second cycle, GRS showed abrupt reductions in reliability. As for RS\_Traditional, there was a loss in accuracy over  
386 the cycles, but it was less steep. In contrast, RS\_Drift led to accuracies tending to zero throughout the entire breeding  
387 process. Regardless of the genetic model used, GRS accuracy increased with the number of markers used. Furthermore,  
388 the additive model has slightly outperformed the additive-dominant one in accuracy for IRS long-term breeding.

389

390 ***Marker density reduction effect - long-term RRS breeding simulations***

391

392       The simulations of breeding programs with the RRS scheme showed different results regarding population mean  
393       with IRS (Fig. 4A). Here GRS outperformed PRS, regardless of marker density and model adopted. The differences  
394       started to become bigger after the fourth breeding cycle. The curves referring to the performance of the best crosses  
395       behaved similarly to the population mean parameter (Fig. 4B). The best crosses obtained via GRS also outperformed  
396       those obtained via PRS. In general, the yield achieved by applying GRS exceeded  $18 \text{ Mg ha}^{-1}$ , on average 64% and 350%  
397       higher than that obtained by using RS\_Traditional and RS\_Drift, respectively. In addition, GRS allowed for the  
398       maintenance of greater genetic divergence (as measured by fixation index, Fst) between heterotic groups (Fig. 4C) than  
399       did PRS, especially from the eighth cycle.

400

401

402       **Fig. 4** Simulation of 20 breeding cycles via PRS (Traditional and Drift) and GRS (with three marker densities) using the  
403       RRS breeding scheme with additive, additive-dominant, CGA, and SCA G-BLUP models. **A** Population mean **B** The  
404       best crosses **C** Fst between HGs **D** Additive variance in HG1 **E** Additive variance in HG2 **F** Accuracy

405

406       Unlike what was observed in the IRS scheme, additive variance (for both heterotic groups) tended to be  
407       consumed more quickly in the breeding program conducted using GRS than in the program applying PRS, regardless of  
408       marker density or model (Fig. 4D and 4E). Similar to the IRS scheme, the PRS led by drift consumed less genetic  
409       variability, and RS\_Traditional showed greater accuracy considering the whole simulation, from the first to the last  
410       breeding cycle. (Fig. 4F). Once again, the accuracy tended to be higher with a higher marker density. Overall, the models  
411       did not significantly affect parameters analyzed in the RRS scheme. However, there was a slight tendency for the additive-  
412       dominant and SGA models' accuracy to be higher than the additive and GCA models.

413

414       *Genetic gains in both breeding schemes*

415

416       Our results show that genetic gains are directly proportional to marker density regardless of the scheme used  
417       (Fig. 5) concerning the selection response to selection after seven years of breeding. In the IRS breeding scheme, the  
418       genetic gain obtained with the GRS\_30K SNP reached more than 100% than that obtained with GRS\_1,5K SNP and  
419       approximately 250% with the GRS\_0,75K SNP (Fig. 5A). Regarding RS (traditional or via drift), the genetic gain of  
420       GRS\_30K reached almost 600% than that obtained by RS\_Traditional. In this scheme, the additive model outperformed  
421       the additive-dominant one when the GRS was conducted with any mark density.

422

423

424     **Fig. 5** Responses to selection after seven years of breeding via PRS (Traditional and Drift) and GRS (with three marker  
425     densities) using different G-BLUP models. **A** IRS **B** RRS breeding schemes

426

427              Genetic gains obtained in the RRS reached levels higher than those IRS schemes (Fig. 5B). However, the  
428     differences between the GRS methods were smaller in the RRS scheme. GRS\_30K achieved 48% superiority over  
429     GRS\_1.5K and 96% over GRS\_0.75K. At the end of the seven years of breeding for GRS\_30K, the dominant additive  
430     model surpassed by almost 22, 56, and 8% of the additive, GCA, and SCA models. Concerning PRS, the genetic gain of  
431     GRS\_30K reached nearly 1,150% than that obtained by RS\_Traditional. In this scheme, the additive-dominant model  
432     surpassed all mark density and was followed by the SCA model. RS\_Drift generated genetic gains close to zero in both  
433     breeding schemes.

434

435 **DISCUSSION**

436

437       High-density genotyping chips can span the whole genome with markers, making it possible to capture the total  
438       genetic variability of a trait throughout the genome and allow predictions of individuals' breeding values without  
439       collecting their phenotypic data (Meuwissen et al. 2001).

440       Implementing low-density marker chips in breeding programs applying GS is a desirable, cost-saving strategy,  
441       given a reduction in genotyping costs. It is especially important in species with low profitability. Several studies have  
442       documented the benefits of using marker subsets in GS analyses (Gorjanc et al. 2015; Tayeh et al. 2015; Ma et al. 2016;  
443       Li et al. 2018; e Sousa et al. 2019). Thus, the finding that few markers were sufficient to achieve a predictive ability  
444       similar to that obtained with high marker density, without changing the genetic relationship estimates among individuals  
445       (Tables 1 and 2; Figs. S1 and S2) spells significant advantages for low-density chips. Their use could lead to technological  
446       improvements and diminishing costs (Heaton et al. 2005; Chessa et al. 2007), as well as reduced computational demand.

447       Nevertheless, over breeding cycles, recombination between markers and QTLs will cause LD decreases, while  
448       selection and drift will potentially act to generate new LD or to strengthen the LD between closely linked loci (Hill and  
449       Robertson, 1968; Lorenz et al., 2011). Thus, there is a need to study the long-term effects of reducing marker density,  
450       particularly given that long-term recurrent selection reduces the effectiveness of GS. Simulations of different reduction  
451       scenarios and controlled factors may do this quickly and at a lower cost.

452

453       ***Consequences of marker reduction for GRS in long-term breeding schemes***

454

455       *Simulated breeding program conducted with the IRS scheme*

456

457       When a breeding program is conducted using recurrent selection on a single population (IRS), the aim is to  
458       release open-pollinated varieties into the field. Hence, the population means, genetic variance, and accuracy are essentials  
459       parameters for evaluating IRS scheme performance. However, at the end of the breeding cycles, our results showed that  
460       GRS applied in an IRS scheme is inferior to phenotypic selection. Probably because RS\_Traditional is driven based on  
461       owner phenotypes. The GRS, on the other hand, is conducted based on the effect of the markers to predict the phenotype  
462       of supposedly superior individuals. As several recombinations are necessary over the breeding cycles, the effect of the  
463       markers must lose efficiency in predicting phenotype, which is reflected in these parameters.

464       GRS accuracy with higher markers density was superior in the first cycle but decreased abruptly from the second  
465       breeding cycle. That suggests that the model training set (TS) must be re-calibrated (updated) during the recurrent  
466       selection process to maintain predictive ability (Neyhart et al. 2017). Indeed, simulated experiments with eucalyptus

467 (Jannink 2010), barley (Denis and Bouvet 2013), and empirical data from advanced-cycle rye (Auinger et al. 2016)  
468 revealed that the accuracy of GS was improved by recalibration of the TS. However, each recurring selection cycle's  
469 recombination events change the QTL-marker LD pattern, diminishing its usefulness for GS. Furthermore, changes in  
470 allele frequencies due to selection over breeding cycles impact the estimation of the effects of allelic substitutions and,  
471 consequently, affects genomic prediction accuracy.

472 The above situation can be aggravated by very low initial marker density, compromising the TS's recalibration.  
473 However, our results show that loss of accuracy over breeding cycles was lower when higher marker densities were used  
474 (Fig. 3C). Furthermore, Müller *et al.* (2017) used simulations to show the persistence of predictive accuracy for a longer  
475 timeframe and, consequently, better selection outcomes when higher marker density was used. Thus, high marker density  
476 seems to buffer the effect of the loss of LD between markers and QTLs over the recombination cycles more effectively  
477 than low marker density, even without TS recalibration. Although, however, there is an expectation that genotyping costs  
478 for low-density chips are lower than those for high-density chips. In practice, the actual differences in terms of genotyping  
479 cost per sample among these chips are negligible. In light of this, and given our results, we stipulate that the economic  
480 gain between high- and low-density chips do not compensate for the prediction accuracy loss in the case of long-term  
481 breeding programs.

482

483 *Simulated breeding program conducted with the RRS scheme*

484

485 When a breeding program is conducted using an RRS scheme, two populations are improved simultaneously  
486 due to reciprocal crosses. At first, the frequencies of desirable alleles in each population are maximized, improving  
487 performance via additive effects. Subsequently, the populations are crossed to capitalize on the non-additive effects of  
488 heterozygosity (Hallauer *et al.* 2010). Thus, RRS is more complex than IRS. However, when conducted with GS support,  
489 it can be carried out with efforts similar to those requiring programs applying an IRS scheme. Population performance  
490 was better with RRS than IRS (Figs. 3A and 4A), likely due to more significant heterosis in the former than in the open-  
491 pollinated individuals of the latter.

492 The performance of single-crosses selected by GRS was higher than those selected using phenotypic selection  
493 (Figs. 4A and 4B). This may be explained by the fact that GRS allows *in-silico* generation of all single-cross combinations  
494 (40,000 in the present study), a process that increases the chances of identifying the best crosses but is unfeasible with  
495 PRS (traditional or via drift). Furthermore, genomic information makes it possible to distinguish between alleles  
496 belonging to different heterotic groups and hence allows for the maintenance of divergence over breeding cycles (Fig.  
497 4C). Given the selection of individuals is carried out on single-crosses, which have high heterotic effects, GRS is better  
498 equipped to capture complementarity, heterozygosity, and genetic distance elements that may contribute to an individual's

499 performance. That is important for a breeding program that aims to release single-crosses, since those are the main  
500 components of heterosis (Moll et al. 1965; Prasad and Singh 1986).

501 Long-term RS may reduce genetic variance to the point where genetic gains are limited (Jannink, 2010). Indeed,  
502 genetic variance showed a substantial decline over breeding cycles in both selection schemes and methods. However, in  
503 the RRS scheme, GRS resulted in more substantial losses of genetic diversity than PRS (Figs. 4D and 4E). Muleta *et al.*  
504 (2019) observed similar results, with higher losses in the simulation of an oligogenic trait with high heritability and lower  
505 losses in a polygenic trait with low heritability. The authors attributed this finding to higher selection intensity in GS,  
506 given more cycles per year (two and three cycles of GS per year were simulated). Essentially, selection intensity  
507 (standardized for cycles per year) will be higher in GRS than in PRS because it has a larger sample size, as previously  
508 mentioned.

509 The maintenance of prediction accuracy across selection cycles is critical for long-term genetic gain (Müller *et*  
510 *al.*, 2017). Interestingly, in our study, RS\_Traditional accuracy surpassed that of GRS throughout cycles. Mathematically,  
511 accuracy is a function of additive variance; a reduction in variance causes a decrease in prediction accuracy. Our finding  
512 may be explained by the breakdown of LD between markers and QTLs due to recombination (Jannink, 2010). Moreover,  
513 over cycles to selection and drift of alleles are not tagged by markers in TS, which cannot be targeted by GRS (Rutkoski  
514 *et al.*, 2015), culminating in low accuracy. Another suggested explanation for this result is that selection is based on  
515 genetic divergence, in which the alleles present in individuals from both heterotic groups change in direction and sense  
516 (upper and lower tail of the genetic variation curve), affecting the reliability of the GS.

517 As seen with the IRS scheme, we noted that higher marker density results in higher accuracy over breeding  
518 cycles with RRS. In addition, a higher number of markers should better capture the LD between QTLs and markers since  
519 some markers may dissociate from the QTLs due to recombination during the cycles, while others may not, keeping some  
520 level of LD.

521

## 522 *Comparison between prediction models in response to selection*

523

524 We consider in this study that a PRS cycle takes 1.75 years on average to complete while 0.35 years is needed  
525 for GRS. Regarding time in a study of this nature is important because each method has a different length. Therefore,  
526 weighting for the time factor makes comparisons fairer (realistic). Thus, over seven years of breeding, it would be possible  
527 to carry out the twenty cycles by GRS, while PRS only four.

528 Predictions obtained using GRS tended to be more accurate than those obtained using PRS, as the former  
529 considers the real genetic relatedness between genotypes rather than the average expectation. Studies have been  
530 conducted with empirical data (Al-Tobasei et al. 2020) and with simulations (Hou et al. 2020) for animal and plant

531 breeding to understand markers' effects on the accuracy of genomic predictions. However, these only consider additive  
532 genetic effects to identify marker subsets and, consequently, to make the predictions.

533 Several other studies have already shown that information about dominance deviations incorporated via kernels  
534 is essential to increase the efficiency of GS prediction models (Azevedo et al. 2015; Dias et al. 2018; Matias et al. 2018;  
535 Dai et al. 2020), especially when the species is cross-pollinated and heterosis is explored in the final product (Santos *et*  
536 *al.*, 2016). Our empirical data findings (Fig. 2) corroborate others that show the importance of incorporating additive  
537 effects and dominance deviations in the models to increase predictive ability (Technow et al. 2012; Alves et al. 2019). In  
538 the presence of dominance, the additive-dominant model is expected to be more accurate, as the additive model falsely  
539 assumes that residuals are independent and identically distributed (Duenk *et al.*, 2017). Furthermore, additive-dominant  
540 models allow for finer-grain assessment of genomic contributions to performance than models that only consider additive  
541 effects. However, when we consider only the additive kernel in the models, additive effects, dominance, and residuals  
542 are not readily distinguishable from one another (Alves *et al.*, 2019).

543 We found in long-term breeding simulations that the additive model was superior primarily in terms of accuracy  
544 to the additive-dominant model for the IRS scheme. This result is most evident in response to selection (Fig. 5) and makes  
545 sense, as the products generated by this scheme (open-pollinated varieties and inbred lines) effectively exploit additive  
546 effects (Hallauer et al., 2010). Contrasting the RRS scheme, the additive-dominant model was superior to those  
547 accompanied by the SGA model. Analogously to the previous explanation, this scheme's genetic material (single-crosses)  
548 effectively exploits additive and non-additive genetic effects. Thus, a model that better captures and separates these  
549 effects should generate the best responses to selection. Furthermore, Duenk *et al.* (2017) showed that in the presence of  
550 dominance, the mean squared error of the average effect of allelic substitution with the additive-dominant model was  
551 always smaller, especially when the heritability was low. Therefore, for cross-pollinated species in which single-crosses  
552 are commercially exploited, it can be expected that genetic models that include dominance will be more robust and have  
553 higher accuracy in estimating the average effect of allelic substitution purely additive models.

554 GS is expected to be not consistently well-suited for RRS since, in the modeling process, only one vector of  
555 additive effects and dominance deviations is considered for the entire population and their respective heterotic groups.  
556 Thus, we use two more approaches to simulate breeding cycles with RRS; G-BLUP modeled only with the GCA effects  
557 and another with SCA effects. Furthermore, based on quantitative genetics theory, with the application RS\_Traditional,  
558 individuals from both heterotic groups can be expected to be selected for breeding based on the same allele's effects,  
559 which might lead to a reduction in the genetic distance among HG over the cycles, vanishing the possibility of the  
560 materialization of heterosis (Hallauer et al. 2010).

561 Breeders that exploit heterosis to generate products (especially single-crosses) can theoretically be expected to  
562 benefit from models that can predict non-additive genetic effects, as this would allow more accurate predictions of the

563 performance of single-crosses. In essence, the two models that best capture dominance deviations are the additive-  
564 dominance model and the SCA one. These models presented similar estimates for the parameters studied, such as the  
565 population mean, the best crosses, and Fst. Werner et al. (2018) used different models to verify the increase in predictive  
566 ability in oilseed rape with the GS approach. Interestingly, that study found that incorporating SCA effects combined  
567 with GCA into the RR-BLUP method (similar to G-BLUP) was the approach that caused the most fluctuation in  
568 predictions over 200 validation cycles. Except for this, all methods (including those based on Bayesian statistics) revealed  
569 comparatively robust predictive ability and high accuracy.

570 We found different trends among the results generated with the models employed after seven years of breeding.  
571 The responses to selection were different, mainly for GRS\_30K. However, the selection process is standardized in the  
572 RRS scheme, i.e., we select the best single-crosses first and then recombine their parents. Even under these circumstances,  
573 it is important to be aware of how to model. The smallest difference between the additive-dominant and SCA models  
574 must be related to the fact that the dominance deviations are captured by SCA effects for heterosis (Falconer and Mackay,  
575 1996).

576 Additionally, it is well-known that non-additive effects are partially included in the parents breeding value,  
577 GCA, or in allele substitution effects (Falconer and Mackay, 1996; Bernardo, 2010; Werner et al., 2018). Despite this, it  
578 is noteworthy that even if, in some situations, the additive-dominance model achieves lower performance than the  
579 additive one, the additive-dominance model generates more reliable estimates (Alves et al. 2019). Reif et al. (2007)  
580 demonstrated that the dominance variance decreases the respective additive variance but increases according to the  
581 populations' divergence. Thus, dominance effects are increasingly absorbed by the population mean or become  
582 inseparable from additive effects over breeding cycles (Technow et al. 2014).

583

584 *Final considerations*

585

586 Disentangling the different factors that may contribute to heterogeneity and predictive ability failure in empirical  
587 datasets is challenging (Dai et al. 2020). However, several models and tools are available to minimize bias. Using  
588 simulations, we can control several uncontrolled factors in the field and focus on answering specific questions. Our results  
589 clearly showed that reducing the marker density is efficient only in a restricted timeframe (within a breeding cycle or a  
590 static group of parents). Hence, low-density marker sets are ineffective for GRS (either with IRS or RRS schemes),  
591 usually conducted for long-term breeding programs, and that does not aim to reach a plateau as inbreeding schemes to  
592 increase homozygosity. Overall, the higher the marker density, the higher the capacity of the marker set to retain relevant  
593 genetic information while being subject to the LD losses between markers and QTLs that result from crossing-over events  
594 during recombination cycles. This brings a direct gain (linear and ascending) for the response to selection. Therefore, this

595 study complements the results of others published recently. Although GS allows a higher genetic gain, especially when  
596 weighted by time, Muleta et al. (2019) also observed more significant accuracy losses in GRS than in PRS throughout  
597 the breeding cycles. In this sense, there is a need for more studies on how often and in what way to update the TS to  
598 maximize prediction accuracy (Rincent et al. 2012), retaining in satisfactory levels the LD between markers and QTLs  
599 (Neyhart et al., 2017), and the relationship between TS and target populations Finally, we recommend that genotyping  
600 should be carried out with high marker density, so that, in long-term breeding programs, the marker set's ability to capture  
601 LD over cycles is maintained, and does not represent a bottleneck. This may avoid the need to update the TS with each  
602 breeding cycle. The additive model is sufficient to predict breeding populations' performance in the IRS breeding scheme  
603 and the additive-dominance one when employed RRS scheme. Furthermore, regarding non-additive genetic effects, we  
604 found a slight difference between the additive-dominance and SCA G-BLUP models.

605 **REFERENCES**

606

- 607 Acosta-Pech R, Crossa J, de los Campos G, et al (2017) Genomic models with genotype  $\times$  environment interaction for  
608 predicting hybrid performance: an application in maize hybrids. *Theor Appl Genet* 130:1431–1440.  
609 <https://doi.org/10.1007/s00122-017-2898-0>
- 610 Akdemir D (2017) STPGA: Selection of training populations with a genetic algorithm. *bioRxiv*.  
611 <https://doi.org/10.1101/111989>
- 612 Al-Tobasei R, Ali A, Garcia A, et al (2020) Genomic Predictions for Muscle Yield and Fillet Firmness in Rainbow  
613 Trout using Reduced-Density SNP Panels. <https://doi.org/10.21203/rs.3.rs-36925/v1>
- 614 Alves FC, Granato ÍSC, Galli G, et al (2019) Bayesian analysis and prediction of hybrid performance. *Plant Methods*  
615 15:1–18. <https://doi.org/10.1186/s13007-019-0388-x>
- 616 Auinger HJ, Schönleben M, Lehermeier C, et al (2016) Model training across multiple breeding cycles significantly  
617 improves genomic prediction accuracy in rye (*Secale cereale* L.). *Theor Appl Genet* 129:2043–2053.  
618 <https://doi.org/10.1007/s00122-016-2756-5>
- 619 Azevedo CF, de Resende MDV, e Silva FF, et al (2015) Ridge, Lasso and Bayesian additive-dominance genomic  
620 models. *BMC Genet* 16:1–13. <https://doi.org/10.1186/s12863-015-0264-2>
- 621 Bandeira e Sousa M, Cuevas J, Couto EG de O, et al (2017) Genomic-enabled prediction in maize using kernel models  
622 with genotype  $\times$  environment interaction. *G3 Genes, Genomes, Genet* 7:1995–2014.  
623 <https://doi.org/10.1534/g3.117.042341>
- 624 Barter RL, Yu B (2018) Superheat: An R Package for Creating Beautiful and Extendable Heatmaps for Visualizing  
625 Complex Data. *J Comput Graph Stat* 27:910–922. <https://doi.org/10.1080/10618600.2018.1473780>
- 626 Basnet BR, Crossa J, Dreisigacker S, et al (2019) Hybrid Wheat Prediction Using Genomic, Pedigree, and  
627 Environmental Covariates Interaction Models. *Plant Genome* 12:180051.  
628 <https://doi.org/10.3835/plantgenome2018.07.0051>
- 629 Bernardo R (2010) Breeding for quantitative traits in plants Stemma Press. Stemma Press, Woodbury
- 630 Chen GK, Marjoram P, Wall JD (2009) Fast and flexible simulation of DNA sequence data. *Genome Res* 19:136–142.  
631 <https://doi.org/10.1101/gr.083634.108>
- 632 Chessa S, Chiatti F, Ceriotti G, et al (2007) Development of a single nucleotide polymorphism genotyping microarray  
633 platform for the identification of bovine milk protein genetic polymorphisms. *J Dairy Sci* 90:451–464.  
634 [https://doi.org/10.3168/jds.S0022-0302\(07\)72647-4](https://doi.org/10.3168/jds.S0022-0302(07)72647-4)
- 635 Covarrubias-Pazaran G (2016) Genome-Assisted prediction of quantitative traits using the r package sommer. *PLoS*  
636 One 11:1–15. <https://doi.org/10.1371/journal.pone.0156744>

- 637 Crossa J, De Los Campos G, Pérez P, et al (2010) Prediction of genetic values of quantitative traits in plant breeding  
638 using pedigree and molecular markers. *Genetics* 186:713–724. <https://doi.org/10.1534/genetics.110.118521>
- 639 Dai Z, Long N, Huang W (2020) Influence of genetic interactions on polygenic prediction. *G3 Genes, Genomes, Genet*  
640 10:109–115. <https://doi.org/10.1534/g3.119.400812>
- 641 Denis M, Bouvet JM (2013) Efficiency of genomic selection with models including dominance effect in the context of  
642 Eucalyptus breeding. *Tree Genet Genomes* 9:37–51. <https://doi.org/10.1007/s11295-012-0528-1>
- 643 Dias KODG, Gezan SA, Guimarães CT, et al (2018) Improving accuracies of genomic predictions for drought  
644 tolerance in maize by joint modeling of additive and dominance effects in multi-environment trials. *Heredity*  
645 (Edinb) 121:24–37. <https://doi.org/10.1038/s41437-018-0053-6>
- 646 Dos Santos JPR, De Castro Vasconcellos RC, Pires LPM, et al (2016) Inclusion of dominance effects in the  
647 multivariate GBLUP model. *PLoS One* 11:1–21. <https://doi.org/10.1371/journal.pone.0152045>
- 648 Duenk P, Calus MPL, Wientjes YCJ, Bijma P (2017) Benefits of dominance over additive models for the estimation of  
649 average effects in the presence of dominance. *G3 Genes, Genomes, Genet* 7:3405–3414.  
650 <https://doi.org/10.1534/g3.117.300113>
- 651 e Sousa MB, Galli G, Lyra DH, et al (2019) Increasing accuracy and reducing costs of genomic prediction by marker  
652 selection. *Euphytica* 215:. <https://doi.org/10.1007/s10681-019-2339-z>
- 653 Elshire RJ, Glaubitz JC, Sun Q, et al (2011) A robust, simple genotyping-by-sequencing (GBS) approach for high  
654 diversity species. *PLoS One* 6:1–10. <https://doi.org/10.1371/journal.pone.0019379>
- 655 Falconer D, Mackay T (1996) Introduction to quantitative genetics, Longman, 4th edn. Longmans Green, Harlow,  
656 Essex, UK.
- 657 Galli G, Alves FC, Morosini JS, Fritsche-Neto R (2020) On the usefulness of parental lines GWAS for predicting low  
658 heritability traits in tropical maize hybrids. *PLoS One* 15:1–15. <https://doi.org/10.1371/journal.pone.0228724>
- 659 Gaynor RC, Gorjanc G, Hickey JM (2020) AlphaSimR: An R-package for Breeding Program Simulations 2 3. *bioRxiv*  
660 2020.08.10.245167
- 661 Gorjanc G, Cleveland MA, Houston RD, Hickey JM (2015) Potential of genotyping-by-sequencing for genomic  
662 selection in livestock populations. *Genet Sel Evol* 47:. <https://doi.org/10.1186/s12711-015-0102-z>
- 663 Granato ISC, Galli G, de Oliveira Couto EG, et al (2018) snpReady: a tool to assist breeders in genomic analysis. *Mol*  
664 *Breed* 38:. <https://doi.org/10.1007/s11032-018-0844-8>
- 665 Hallauer A, Carena M, Filho JM (2010) Quantitative genetics in maize breeding
- 666 Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME (2009) Invited review: Genomic selection in dairy cattle:  
667 Progress and challenges. *J Dairy Sci* 92:433–443. <https://doi.org/10.3168/jds.2008-1646>
- 668 Heaton MP, Keen JE, Clawson ML, et al (2005) Use of bovine single nucleotide polymorphism markers to verify

- 669 sample tracking in beef processing. *J Am Vet Med Assoc* 226:1311–1314.
- 670 <https://doi.org/10.2460/javma.2005.226.1311>
- 671 Henderson CR (1975) Best Linear Unbiased Estimation and Prediction under a Selection Model. *Biometrics* 31:423.
- 672 <https://doi.org/10.2307/2529430>
- 673 Hickey JM, Dreisigacker S, Crossa J, et al (2014) Evaluation of genomic selection training population designs and
- 674 genotyping strategies in plant breeding programs using simulation. *Crop Sci* 54:1476–1488.
- 675 <https://doi.org/10.2135/cropsci2013.03.0195>
- 676 Hill W, Robertson A (1968) Linkage disequilibrium in finite populations. *Theor Appl Genet* 38:226–231.
- 677 <https://doi.org/10.1080/03071848408522227>
- 678 Hou L, Liang W, Xu G, et al (2020) Accuracy of genomic prediction using mixed low-density marker panels. *Anim*
- 679 *Prod Sci* 60:999–1007. <https://doi.org/10.1071/AN18503>
- 680 Jannink JL (2010) Dynamics of long-term genomic selection. *Genet Sel Evol* 42:1–11. <https://doi.org/10.1186/1297-9686-42-35>
- 681 Li B, Zhang N, Wang YG, et al (2018) Genomic prediction of breeding values using a subset of SNPs identified by
- 682 three machine learning methods. *Front Genet* 9:1–20. <https://doi.org/10.3389/fgene.2018.00237>
- 683 Lorenz AJ, Chao S, Asoro FG, et al (2011) Genomic Selection in Plant Breeding. Knowledge and Prospects., 1st edn.
- 684 Elsevier Inc.
- 685 Ma Y, Reif JC, Jiang Y, et al (2016) Potential of marker selection to increase prediction accuracy of genomic selection
- 686 in soybean (*Glycine max* L.). *Mol Breed* 36:1–10. <https://doi.org/10.1007/s11032-016-0504-9>
- 687 Matias FI, Barrios SCL, Bearari LM, et al (2018) Contribution of additive and dominance effects on agronomical and
- 688 nutritional traits, and multivariate selection on *Urochloa* spp. hybrids. *Crop Sci* 58:2444–2458.
- 689 <https://doi.org/10.2135/cropsci2018.04.0261>
- 690 Meuwissen THE, Hayes BJ, Goddard ME (2001) Prediction of total genetic value using genomewide dense marker
- 691 maps. *Genetics* 157:1819–1829. <https://doi.org/10.1093/genetics/157.4.1819>
- 692 Moll R, Lonnquist J, Vélez Fortuno J, Johnson E (1965) The relationship of heterosis and genetic divergence in maize.
- 693 *Genetics* 52:139–144
- 694 Muleta KT, Pressoir G, Morris GP (2019) Optimizing genomic selection for a sorghum breeding program in Haiti: A
- 695 simulation study. *G3 Genes, Genomes, Genet* 9:391–401. <https://doi.org/10.1534/g3.118.200932>
- 696 Müller D, Schopp P, Melchinger AE (2017) Persistency of prediction accuracy and genetic gain in synthetic
- 697 populations under recurrent genomic selection. *G3 Genes, Genomes, Genet* 7:801–811.
- 698 <https://doi.org/10.1534/g3.116.036582>
- 699 Muñoz F, Rodriguez LS (2016) breedR : An open statistical package to analyse genetic data
- 700

- 701 Neyhart JL, Tiede T, Lorenz AJ, Smith KP (2017) Evaluating methods of updating training data in long-term  
702 genomewide selection. *G3 Genes, Genomes, Genet* 7:1499–1510. <https://doi.org/10.1534/g3.117.040550>
- 703 Oksanen J, Blanchet FG, Friendly M, et al (2019) Package 'vegan' Title Community Ecology Package. *Community*  
704 *Ecol Packag* 2:1–297
- 705 Patterson HD, Williams ER (1976) A new class of resolvable incomplete block designs. *Biometrika* 63:83–92.  
706 <https://doi.org/10.1093/biomet/63.1.83>
- 707 Pocrnic I, Lourenco DAL, Masuda Y, Misztal I (2016) Dimensionality of genomic information and performance of the  
708 Algorithm for Proven and Young for different livestock species. *Genet Sel Evol* 48:1–9.  
709 <https://doi.org/10.1186/s12711-016-0261-6>
- 710 Prasad SK, Singh TP (1986) Heterosis in relation to genetic divergence in maize (*Zea mays* L.). *Euphytica* 35:919–924.  
711 <https://doi.org/10.1007/BF00028600>
- 712 Reif JC, Gumpert FM, Fischer S, Melchinger AE (2007) Impact of interpopulation divergence on additive and  
713 dominance variance in hybrid populations. *Genetics* 176:1931–1934. <https://doi.org/10.1534/genetics.107.074146>
- 714 Rincent R, Laloë D, Nicolas S, et al (2012) Maximizing the reliability of genomic selection by optimizing the  
715 calibration set of reference individuals: Comparison of methods in two diverse groups of maize inbreds (*Zea*  
716 *mays* L.). *Genetics* 192:715–728. <https://doi.org/10.1534/genetics.112.141473>
- 717 Rutkoski J, Singh RP, Huerta-Espino J, et al (2015) Efficient Use of Historical Data for Genomic Selection: A Case  
718 Study of Stem Rust Resistance in Wheat. *Plant Genome* 8:1–10.  
719 <https://doi.org/10.3835/plantgenome2014.09.0046>
- 720 Shao J (1993) Linear model selection by cross-validation. *J Am Stat Assoc* 88:486–494.  
721 <https://doi.org/10.1016/j.jspi.2003.10.004>
- 722 Subedi S, Feng Z, Deardon R, Schenkel FS (2013) SNP selection for predicting a quantitative trait. *J Appl Stat* 40:600–  
723 613. <https://doi.org/10.1080/02664763.2012.750282>
- 724 Tayeh N, Klein A, Le Paslier MC, et al (2015) Genomic prediction in pea: Effect of marker density and training  
725 population size and composition on prediction accuracy. *Front Plant Sci* 6:1–11.  
726 <https://doi.org/10.3389/fpls.2015.00941>
- 727 Technow F, Riedelsheimer C, Schrag TA, Melchinger AE (2012) Genomic prediction of hybrid performance in maize  
728 with models incorporating dominance and population specific marker effects. *Theor Appl Genet* 125:1181–1194.  
729 <https://doi.org/10.1007/s00122-012-1905-8>
- 730 Technow F, Schrag TA, Schipprack W, et al (2014) Genome properties and prospects of genomic prediction of hybrid  
731 performance in a breeding program of maize. *Genetics* 197:1343–1355.  
732 <https://doi.org/10.1534/genetics.114.165860>

- 733 Unterseer S, Bauer E, Haberer G, et al (2014) A powerful tool for genome analysis in maize: Development and  
734 evaluation of the high density 600 k SNP genotyping array. *BMC Genomics* 15:1–15.  
735 <https://doi.org/10.1186/1471-2164-15-823>
- 736 VanRaden PM (2008) Efficient methods to compute genomic predictions. *J Dairy Sci* 91:4414–4423.  
737 <https://doi.org/10.3168/jds.2007-0980>
- 738 Walsh B, Lynch M (2018) Evolution and Selection of Quantitative Traits
- 739 Werner CR, Qian L, Voss-Fels KP, et al (2018) Genome-wide regression models considering general and specific  
740 combining ability predict hybrid performance in oilseed rape with similar accuracy regardless of trait  
741 architecture. *Theor Appl Genet* 131:299–317. <https://doi.org/10.1007/s00122-017-3002-5>
- 742 Xu S (2013) Genetic mapping and genomic selection using recombination breakpoint data. *Genetics* 195:1103–1115.  
743 <https://doi.org/10.1534/genetics.113.155309>
- 744 Yassue RM, Sabadin JFG, Galli G, et al (2020) CV- $\alpha$ : designing validations sets to increase the precision and enable  
745 multiple comparison tests in genomic prediction
- 746 Zhang Z, Erbe M, He J, et al (2015) Accuracy of whole-genome prediction using a genetic architecture-enhanced  
747 variance-covariance matrix. *G3 Genes, Genomes, Genet* 5:615–627. <https://doi.org/10.1534/g3.114.016261>

# Figures

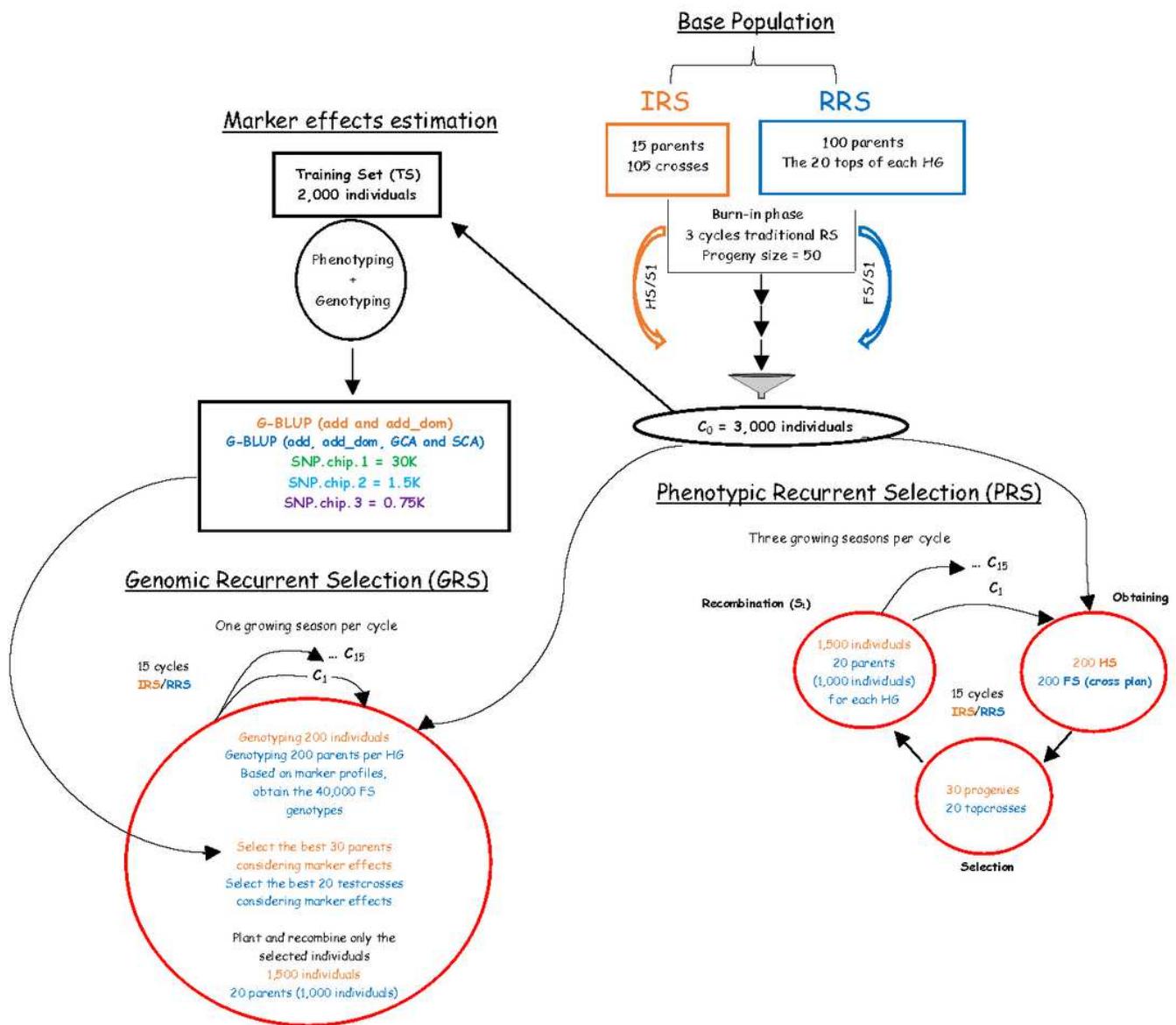
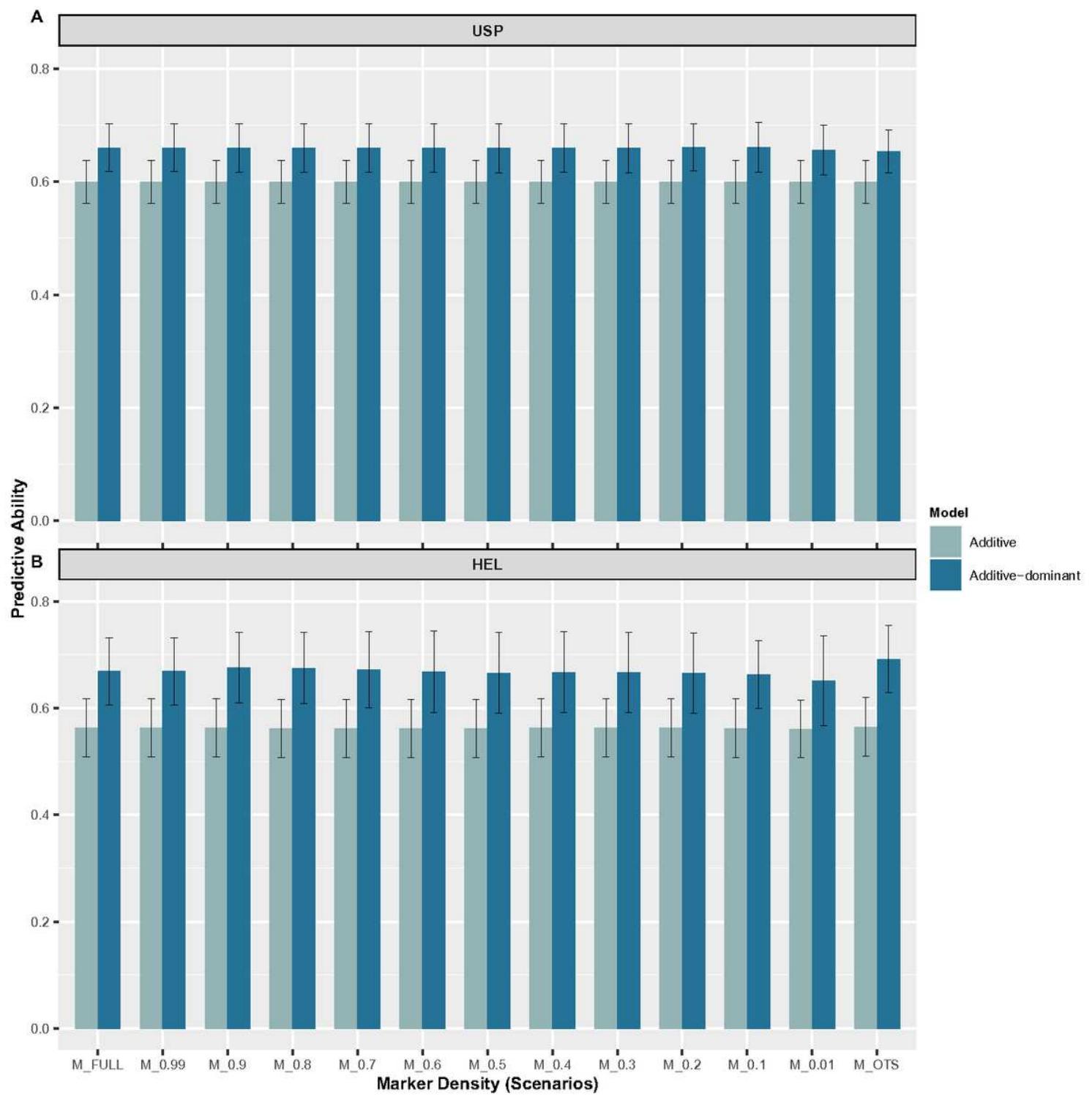


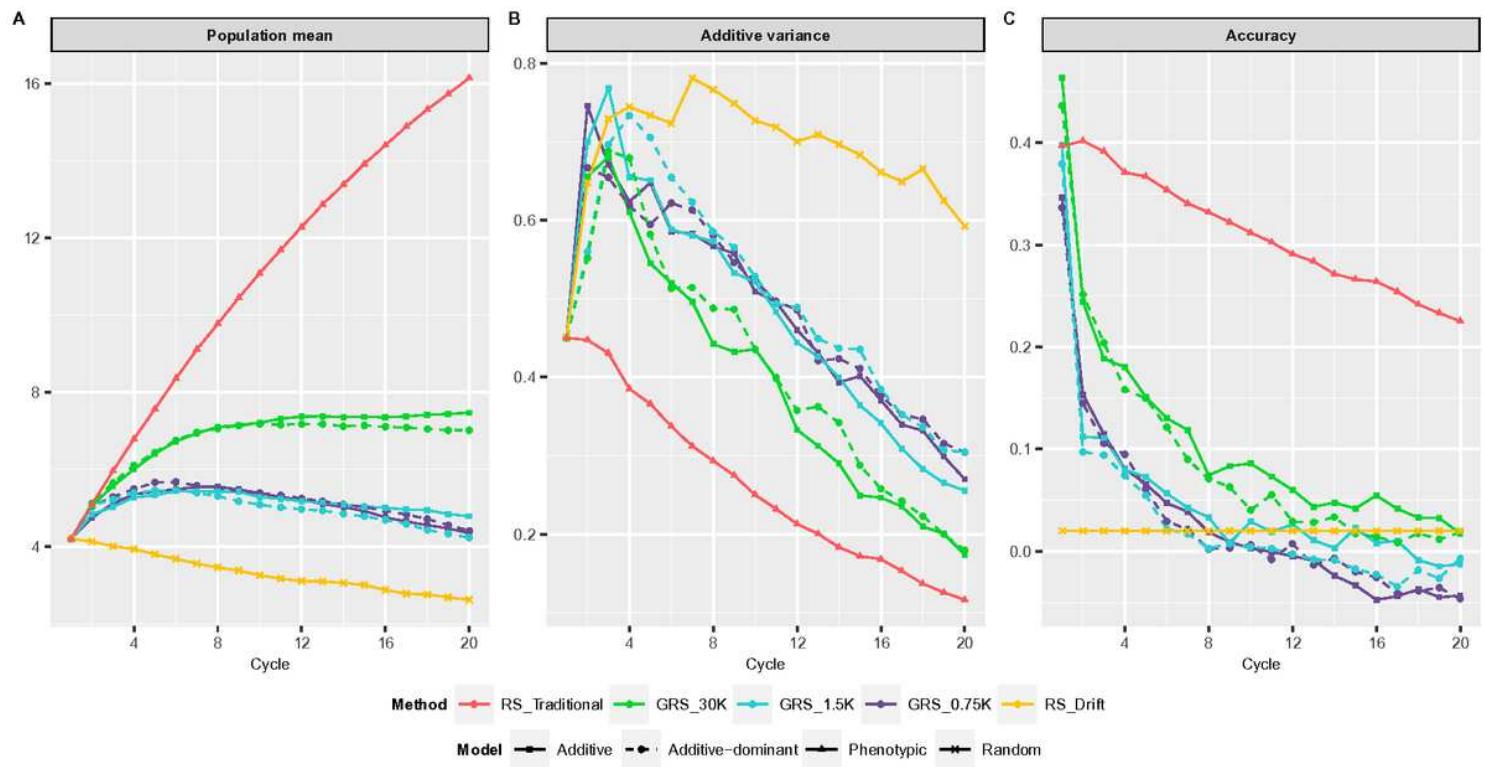
Figure 1

Schematic representation of the recurrent selection breeding schemes applied in the simulations. After obtaining the historical population, base populations were selected from it randomly, and these served as the basis for the IRS (orange) and RRS (blue) schemes. Thus, information in black applies to both schemes.



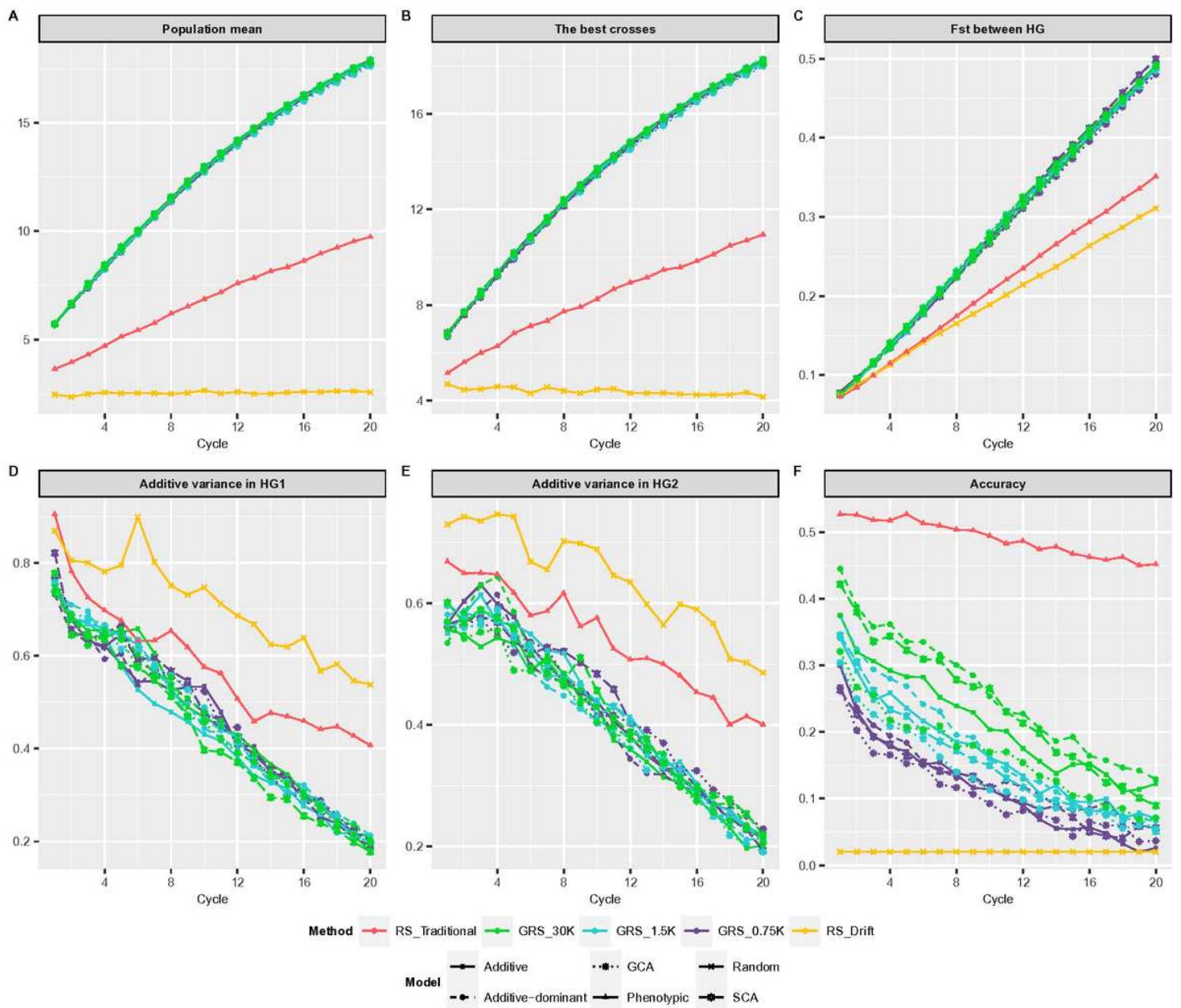
**Figure 2**

Predictive ability across marker density reduction scenarios based on different linkage disequilibrium levels and 363 an OTS algorithm, using additive and additive-dominant models. A USP dataset B HEL dataset.



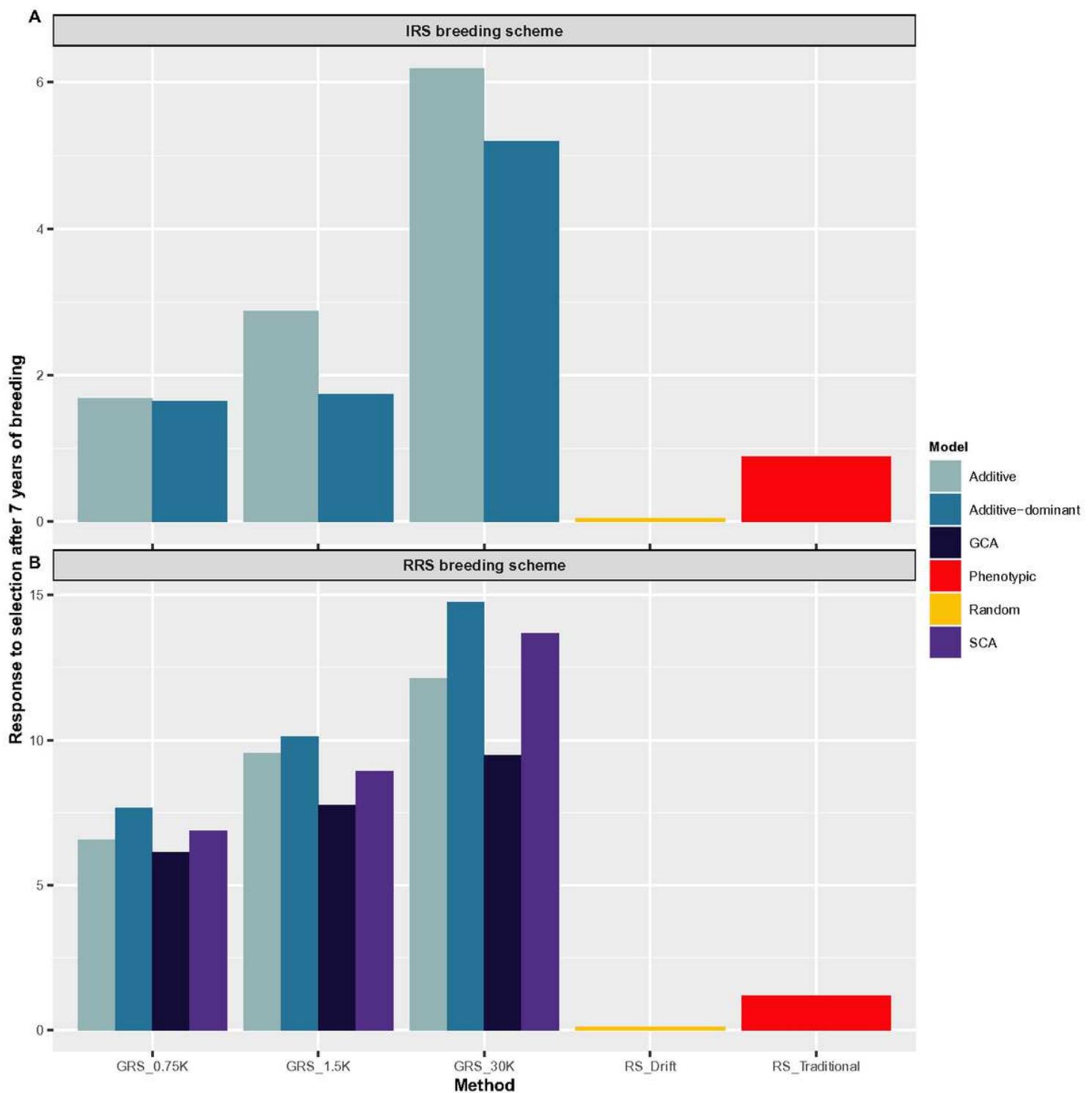
**Figure 3**

Simulation of 20 breeding cycles via PRS (Traditional and Drift) and GRS (with three marker densities) using the IRS breeding scheme with additive and additive-dominant G-BLUP models. A Population mean B Additive variance C Accuracy



**Figure 4**

Simulation of 20 breeding cycles via PRS (Traditional and Drift) and GRS (with three marker densities) using the RRS breeding scheme with additive, additive-dominant, CGA, and SCA G-BLUP models. A Population mean B The best crosses C Fst between HGs D Additive variance in HG1 E Additive variance in HG2 F Accuracy



**Figure 5**

Responses to selection after seven years of breeding via PRS (Traditional and Drift) and GRS (with three marker densities) using different G-BLUP models. An IRS B RRS breeding schemes

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- Table1.xlsx
- Table2.xlsx
- SUPPLEMENTARYMATERIAL.pdf