

1 Genomic prediction and training set optimization in a 2 structured Mediterranean oat population

3 Simon Rio* · Luis Gallego-Sánchez* ·
4 Gracia Montilla-Bascón · Francisco J.
5 Canales · Julio Isidro y Sánchez · Elena
6 Prats

7
8 Received: date / Accepted: date

9 * These authors contributed equally to this study

10 **Abstract Key message: The strong genetic structure observed in Mediter-**
11 **anean oats affects the predictive ability of genomic prediction as well**
12 **as the performance of training set optimization methods.** In this study,
13 we investigated the efficiency of genomic prediction and training set optimization in
14 a highly structured population of cultivars and landraces of cultivated oat (*Avena*
15 *sativa*) from the Mediterranean basin, including white (subsp. *sativa*) and red
16 (subsp. *byzantina*) oats, genotyped using genotype-by-sequencing markers, and
17 evaluated for agronomic traits in Southern Spain. For most traits, the predictive
18 abilities were moderate to high with little differences between models, except for
19 biomass for which Bayes-B showed a substantial gain compared to other models.
20 The consistency between the structure of the training population and the popula-
21 tion to be predicted was key to the predictive ability of genomic predictions. The
22 predictive ability of inter-subspecies predictions was indeed much lower than that
23 of intra-subspecies predictions for all traits. Regarding training set optimization,
24 the linear mixed model optimization criteria (PEVmean and CDmean) performed
25 better than the heuristic approach “partitioning around medoids”, even under high
26 population structure. The superiority of CDmean and PEVmean could be explained
27 by their ability to adapt the representation of each genetic group according to
28 those represented in the population to be predicted. These results represent an
29 important step towards the implementation of genomic prediction in oat breeding
30 programs and address important issues faced by the genomic prediction community
31 regarding population structure and training set optimization.

Simon Rio · Julio Isidro y Sánchez

Centro de Biotecnología y Genómica de Plantas (CBGP, UPM-INIA) Universidad Politécnica
de Madrid (UPM) - Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria
(INIA) Campus de Montegancedo-UPM 28223-Pozuelo de Alarcón, (Madrid), Spain
E-mail: j.isidro@upm.es

Luis Gallego-Sánchez · Gracia Montilla-Bascón · Francisco J. Canales · Elena Prats
Institute for Sustainable Agriculture, Spanish Research Council (CSIC), Córdoba, Spain
E-mail: elena.prats@ias.csic.es

32 **Keywords** Oat · *Avena sativa* · Genomic prediction · Training set optimization ·
33 Genetic structure · Environmental adaptation

34 Introduction

35 The cultivated oat (*Avena sativa* L.) is an important and versatile crop that can be
36 grown for grain, forage, feed and straw (Welch, 2012). In the last twenty years, the
37 oat cultivation area within the Mediterranean region increased with an approximate
38 rate of 7,500 Ha per year (FAO., 2017). However, the yields achieved in Northern
39 Europe remains much greater than that of the Mediterranean area, as a result of
40 the limited adaptation of oats to the environmental conditions of Southern Europe
41 (e.g., water availability or temperature). Thus, there is a critical need for selecting
42 oat genotypes that are better adapted to the Mediterranean environment by taking
43 advantage of the existing diversity in the region: the white (subsp. *sativa*) and red
44 (subsp. *byzantina*) oats.

45 The advent of molecular markers has revolutionized the methodologies for
46 selecting individuals in animal and plant breeding programs. Among these new
47 methods, genomic prediction (GP) was proposed by Meuwissen et al. (2001) and
48 stands today as one of the most promising tool. In its simplest application, a set
49 of individuals is evaluated for a given trait and genotyped using single nucleotide
50 polymorphisms (SNPs). A statistical model is trained on this data set, referred to
51 as the training set (TRS), and is used to predict the breeding value of individuals
52 for whom only SNP information is available, referred to as the test set (TS) (Isidro
53 et al., 2016). Several methods have been proposed in the literature including
54 models making different hypotheses on the distribution of the effects of quantitative
55 trait loci (QTL) like GBLUP, BayesA, BayesB or BayesC π (Meuwissen et al.,
56 2001; Habier et al., 2011; Heslot et al., 2012), semi-parametric methods like the
57 reproducing kernel Hilbert space (RKHS) (Gianola et al., 2006; Gianola and van
58 Kaam, 2008), or tree-based methods like random forests (Breiman, 2001; Chen and
59 Ishwaran, 2012).

60 One of the most critical steps in GP is the selection of the TRS since it is critical
61 to the predictive ability of the models. In the last few years, several studies have
62 investigated different approaches to optimize TRSs (Rincent et al., 2012; Hickey
63 et al., 2014; Akdemir et al., 2015; Isidro et al., 2015; Lorenz and Smith, 2015;
64 Tayeh et al., 2015; Akdemir, 2017; Rincent et al., 2017; Brandariz and Bernardo,
65 2018; Norman et al., 2018; Akdemir and Isidro-Sánchez, 2019; Berro et al., 2019;
66 Edwards et al., 2019; Guo et al., 2019; Mangin et al., 2019; Ou and Liao, 2019;
67 Sarinelli et al., 2019; Alvarenga et al., 2020; Olatoye et al., 2020; Roth et al.,
68 2020). Among the optimization criteria, some approaches have been the subject of
69 particular consideration including the mean coefficient of determination (CDmean)
70 and the mean prediction error variance (PEVmean) initially presented for contrasts
71 between individuals (Laloë, 1993; Rincent et al., 2012), or clustering methods
72 such as stratified sampling (Isidro et al., 2015; Akdemir and Isidro-Sánchez, 2019;
73 Guo et al., 2019) and partitioning around medoids (PAM) (Guo et al., 2019). The
74 CDmean, PEVmean and PAM criteria are now routinely used for TRS optimization,
75 especially when the TRS size is small (Akdemir and Isidro-Sánchez, 2019). As part
76 of the optimization process, the population structure plays a key role as it impacts

77 both the performance of the optimization methods (Isidro et al., 2015) and the
78 GP predictive ability.

79 In a population stratified into genetics groups, when the same genetic groups are
80 found within the TRS and the TS, the differences in means between groups are often
81 implicitly taken into account by the model and contribute to the predictive ability
82 (Guo et al., 2014; Rio et al., 2019). Conversely, when targeting a group-specific
83 TS, training a model on a different group can dramatically limit the predictive
84 ability, as shown in dairy and beef cattle (Olson et al., 2012; Chen et al., 2013)
85 or maize (Technow et al., 2013; Lehermeier et al., 2014). As proposed by de Roos
86 et al. (2009), genetic groups can also be combined into generic multi-group TRSs
87 that show a good predictive ability regardless of the target population, as shown
88 in dairy cattle (Brøndum et al., 2011; Pryce et al., 2011; Zhou et al., 2013), maize
89 (Technow et al., 2013; Rio et al., 2019) or soybean (Duhnen et al., 2017). Several
90 models have been proposed that explicitly account for genetic structure such as
91 modeling genetic covariances between individuals from different groups by adapting
92 multi-trait models (Karoui et al., 2012; Lehermeier et al., 2015).

93 The use of genomics in oat breeding is rather limited compared to other cereals
94 like maize, wheat or rice, due to the scarcity of available tools. It can be explained
95 by the complexity of its allo-hexaploid genome ($2n=6x=42$) with high content in
96 repetitive sequences (Yan et al., 2016). Nevertheless, thanks to the efforts of the
97 oat community over the past few years, several genome tools have been developed
98 such as the Illumina 6K gene chip (Tinker et al., 2014), genotyping-by-sequencing
99 (GBS) (Huang et al., 2014; Bekele et al., 2018), and a consensus map (Chaffin
100 et al., 2016). Those tools enabled many genetic studies and breeding applications
101 (Esvelt Klos et al., 2016; Tumino et al., 2016; Yan et al., 2016; Bjørnstad et al.,
102 2017; Tumino et al., 2017; Carlson et al., 2019; Kebede et al., 2019; Sunstrum
103 et al., 2019; Isidro-Sánchez et al., 2020a,b; Yan et al., 2020). More recently, the
104 draft of the hexaploid *Avena sativa* genome sequence: OT3098 v1 - PepsiCo¹ and
105 the sequence of two diploid oat genomes: *Avena Atlantica* and *Avena Eriantha*
106 (Maughan et al., 2019) have been released. They will open a new frontier for the
107 study of the oat genome and for the development of genomics-assisted breeding.

108 A few studies have focused on the application of GP and genomic selection
109 (GS) in oat (Asoro et al., 2011, 2013; Bekele et al., 2018; Mellers et al., 2020;
110 Haikka et al., 2020a,b). In these studies, the objectives included the comparison of
111 GS to traditional phenotypic and marker-assisted selection for β -glucans (Asoro
112 et al., 2013), GP of heading date using SNPs and tag-levels haplotype markers
113 (Bekele et al., 2018), GP of agronomic traits and Fusarium head blight in an oat
114 commercial breeding program (Haikka et al., 2020b,a), and the implementation of
115 GP within a winter oat biparental cross (Mellers et al., 2020). These empirical GS
116 applications have demonstrated the effective use of GS within breeding populations
117 to accelerate oat breeding. Nevertheless, there is a lack of experimental studies
118 focusing on the efficiency of GP and TRS optimization in highly structured oat
119 populations.

120 In this paper, a structured Mediterranean oat population, including both white
121 and red oat inbred lines, was evaluated for agronomic traits. The objectives were
122 to (i) evaluate the predictive ability of different GP models, (ii) optimize TRSs

¹ https://wheat.pw.usda.gov/GG3/graingenes_downloads/oat-ot3098-pepsico

123 using different methods, and (iii) evaluate the impact of genetic structure on both
124 the GP predictive ability and the performance of optimization methods.

125 **Materials and Methods**

126 Genetic material and genotypic data

127 Genetic material consists of a collection of 709 cultivated oat (*Avena sativa*) inbred
128 lines, including landraces from the Mediterranean area, provided by the “Centro
129 de Recursos Filogenéticos” of the “Instituto Nacional de Investigación y Tecnología
130 Agraria y Alimentaria” (INIA, Madrid, Spain) and the United States Department
131 of Agriculture (Washington, USA), along with cultivars and breeding lines provided
132 by different institutions, as presented in Sánchez-Martín et al. (2014) and Canales
133 et al. (2021a). All individuals were genotyped using GBS-SNP markers as detailed
134 in Canales et al. (2021a). Genotypes were pooled into libraries of 96 genotypes at
135 the genomic platform of the McGill University (Canada), following the PstI-MspI
136 method (Huang et al., 2014). Each GBS library was sequenced on a single line of a
137 HiSeq 2500 at the “Plateforme d’Analyses Génomiques of the Institut de Biologie
138 Intégrative et des Systèmes” of the “Université Laval” (Quebec City, Canada). Raw
139 FASTQ sequences were processed using the Haplotag pipeline (Tinker et al., 2016).
140 After a filtering on the minor allele frequency ($> 5\%$), the heterozygosity rate
141 ($< 20\%$) and the percentage of missing values ($< 50\%$), a total of 17,288 bi-allelic
142 SNP markers corresponding to 12,418 tags were obtained. Ten individuals were also
143 discarded due to a large heterozygosity rate and/or a large percentage of missing
144 values, leaving 609 individuals for subsequent analyses. Missing values at SNPs were
145 then imputed using the multivariate normal expectation maximization algorithm
146 (Poland and Rife, 2012) implemented in the R package rrBLUP (Endelman, 2011).
147 The marker dataset is available at Dryad Data (Canales et al., 2021b) and sequence
148 read data are available from NCBI SRA archive as BioProject ID PRJNA693576
149 (<http://www.ncbi.nlm.nih.gov/bioproject/693576>).

150 Structure analysis

151 A structure analysis was performed using the STRUCTURE software (Pritchard
152 et al., 2000) for a number of genetic groups Q ranging from 2 to 5 and using the
153 admixture model with correlated allele frequencies between groups (Falush et al.,
154 2003). Each analysis consisted of 10,000 MCMC iterations and a burn-in of 1,000
155 iterations. Admixture barplots are presented in Fig. 1A for $Q = 3$ and in Online
156 Resource Fig. S1 for other values of Q . The population could be separated into
157 two groups corresponding to the two oat subspecies forming the population (*sativa*
158 and *byzantina*), further referred to as Byzantina and Sativa. As the Sativa group
159 was mainly structured into two sub-groups, we considered three genetic groups for
160 further analyses by assigning individuals using their maximal admixture coefficient:
161 Byzantina (257 lines), Sativa_A (243 lines) and Sativa_B (199 lines). Global position
162 system coordinates were available for most individual accessions (Canales et al.,
163 2021a) and revealed a relationship between the site where individuals were collected
164 and the genetic group to which they were assigned (Fig. 1B). The Evanno method

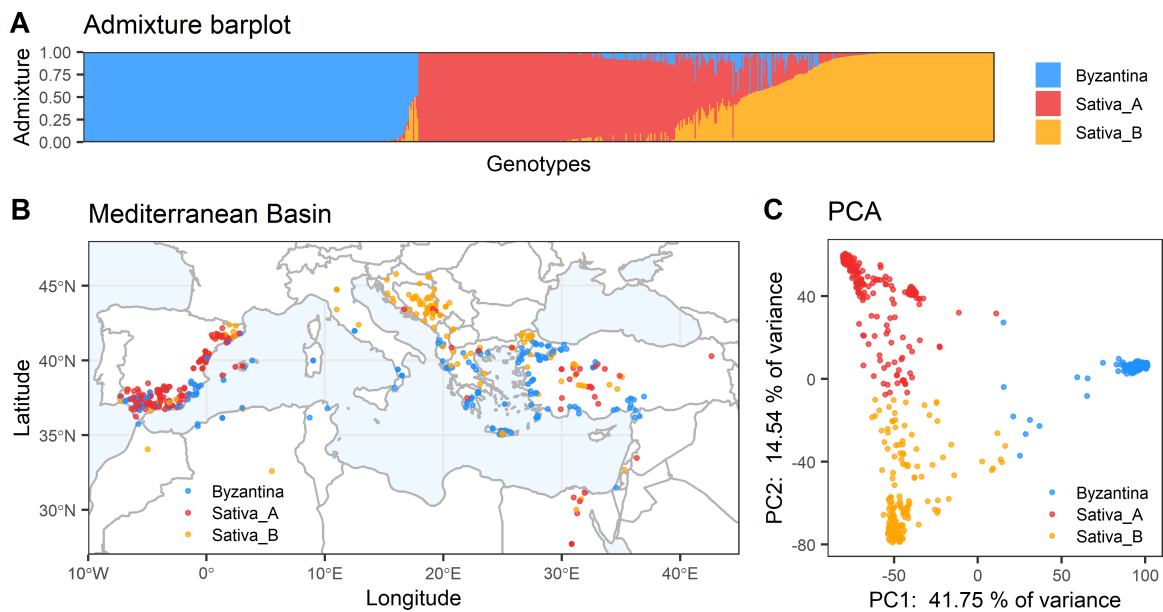


Fig. 1 Graphs illustrating the structure analysis performed on the oat dataset using STRUC-TURE and showing the existence of three genetic groups: Byzantina (257 individuals), Sativa_A (243 individuals) and Sativa_B (199 individuals). Graphs include (A) an admixture barplot showing the admixture proportions of each individual, (B) a map of the Mediterranean Basin with dots indicating the location where individuals have been collected (Canales et al., 2021a), and (C) the two principal components of a principal component analysis (PCA) performed on genomic data. For graphs B and C, dots were colored according to the group to which each individual was assigned based on its maximal admixture coefficient.

165 was applied was applied and supported the use of $Q = 2$ groups due to the high 166 genetic divergence between Byzantina and Sativa subspecies (Online Resource 167 Figure S2). Like in Canales et al. (2021a), where a similar structure analysis 168 was performed on this same dataset with very consistent results, and a previous 169 study based on SSR markers (Montilla-Bascón et al., 2013), we subdivided the 170 Sativa group into two groups, which was supported by a substantial proportion of 171 variance explained by the axis differentiating these groups (14.54%) for a principal 172 component analysis (PCA) performed on the genomic data (Fig 1C). Note that, 173 unlike in Canales et al. (2021a), we did not define additional categories for admixed 174 individuals.

175 Phenotypic analysis

176 The oat collection was evaluated in 2017 and 2018 in Cordoba (Spain) with an 177 altitude of 90 m and light clay calcic cambisol soil, and in 2018 in Santaella (Spain) 178 with an altitude of 238 m, and a light clay eutric gleysol soil, forming three distinct 179 trials (Co17, Co18, and Sa18). Each trial was an alpha lattice square design with 180 three replicates using the cultivar Patones as check. Four agronomic traits were 181 evaluated: heading time (Heading) in growing degree-days, plant height (Height)

182 in cm, vegetative biomass (Biomass) in t/ha and grain yield (Yield) in t/ha, see
 183 Canales (2019) and Canales et al. (2021a) for more details. The contribution of
 184 genotype-by-environment (GxE) interactions to the phenotypic variance and the
 185 broad-sense heritability were investigated using the following model:

$$Y_{ijk_r} = \mu + \alpha_k + \beta_j + G_{ik} + (G \times \beta)_{ijk} + \gamma_{rj} + E_{ijk_r},$$

186 where Y_{ijk_r} is the phenotype of individual i from group k in block r of trial
 187 j , μ is the global intercept, α_k is the effect of group k with $k \in \{By, SaA, SaB\}$
 188 (By : Byzantina, SaA : Sativa_A and SaB : Sativa_B), β_j is the effect of trial j , G_{ik}
 189 is the random genotypic effect of individual i from group k with $G_{ik} \sim \mathcal{N}(0, \sigma_{G_k}^2)$
 190 independent, $(G \times \beta)_{ijk}$ is the random genotype-by-environment (GxE) effect of
 191 individual i from group k in trial j with $(G \times \beta)_{ijk} \sim \mathcal{N}(0, \sigma_{G \times \beta_k}^2)$ independent,
 192 γ_{rj} is the effect of block r in trial j , and E_{ijk_r} is the error of individual i from
 193 group k in block r of trial j with $E_{ijk_r} \sim \mathcal{N}(0, \sigma_E^2)$ independent and identically
 194 distributed. All random effects are assumed to be independent of each other.
 195 Model parameters were estimated using the R package ‘‘MM4LMM’’ (Laporte and
 196 Mary-Huard, 2020). The group-specific means were calculated as following:

$$\mu_k = \mu + \alpha_k + \frac{1}{J} \sum_{j=1}^J \beta_j + \frac{1}{JR} \sum_{j=1}^J \sum_{r=1}^R \gamma_{rj},$$

197 where $J = 3$ is the number of trials and $R = 3$ is the number of blocks in each
 198 trial. The group-specific broad-sense heritabilities were calculated as following:

$$H_k^2 = \frac{\sigma_{G_k}^2}{\sigma_{G_k}^2 + \frac{1}{J} \sigma_{G \times \beta_k}^2 + \frac{1}{JR} \sigma_E^2},$$

199 Least-square means (LS-means) of each individual (Y_{ik}^*) were computed based
 200 on the same model with G_{ik} and $(G \times \beta)_{ijk}$ as fixed effects using:

$$Y_{ik}^* = \mu_k + G_{ik} + \frac{1}{J} \sum_{j=1}^J (G \times \beta)_{ijk},$$

201 and were further referred to as phenotypes for GP analyses.

202 Genomic prediction models

203 In this study, the first four GP models (GBLUP, MGBLUP, Bayes-B and RKHS)
 204 can be written as:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{g} + \mathbf{e}, \quad (1)$$

205 where \mathbf{y} is the vector of reference phenotypes (i.e., the LSmeans), \mathbf{X} is the design
 206 matrix for fixed effects, $\boldsymbol{\beta}$ is the vector of fixed effects, \mathbf{Z} is the incidence matrix
 207 linking phenotypes to breeding values, \mathbf{g} is the vector of breeding values and \mathbf{e} is
 208 the vector of errors. All models assume independence between \mathbf{g} and \mathbf{e} .

209 *GBLUP* A standard additive GBLUP model was applied using the R package
 210 “rrBLUP” (Endelman, 2011) with the following assumptions: $\beta = \mu$ is the overall
 211 mean, $\mathbf{g} \sim \mathcal{N}(0, \mathbf{K}\sigma_G^2)$, \mathbf{K} is the kinship matrix, σ_G^2 is the genetic variance,
 212 $\mathbf{e} \sim \mathcal{N}(0, \mathbf{I}\sigma_E^2)$, \mathbf{I} is the identity matrix and σ_E^2 is the error variance. The kinship
 213 between individuals i and j ($K_{i,j}$) was computed following VanRaden (2008):

$$K_{i,j} = \frac{\sum_{m=1}^M (W_{im} - f_m)(W_{jm} - f_m)}{\sum_{m=1}^M f_m(1 - f_m)}, \quad (2)$$

214 where M is the number of markers, W_{im} is the genotypic score of individual i at
 215 locus m (coded 0; 0.5; 1) and f_m is the frequency of allele “1” at locus m estimated
 216 on the whole dataset.

217 *MGBLUP* A multi-group GBLUP (MGBLUP) model (Lehermeier et al., 2015) was
 218 applied using the R package “MTM”². This model consists of adapting a multi-trait
 219 model to the analysis of one trait in different groups, which allows for group-specific
 220 genetic variances and specific genetic covariances between groups. In this model,

221 $\beta = (\mu_{By}, \mu_{SaA}, \mu_{SaB})^T$ is the vector of group-specific means, $\mathbf{g} = \begin{bmatrix} \mathbf{g}_{By}^* \\ \mathbf{g}_{SaA}^* \\ \mathbf{g}_{SaB}^* \end{bmatrix}$ is the

222 expanded vector of breeding values of each individual in each group of size of $3N$,
 223 N being the number of individuals, where:

224 $\begin{bmatrix} \mathbf{g}_{By}^* \\ \mathbf{g}_{SaA}^* \\ \mathbf{g}_{SaB}^* \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} \sigma_{G_{By}}^2 & \sigma_{G_{By},B} & \sigma_{G_{By},SaB} \\ \sigma_{G_{By},SaA} & \sigma_{G_{SaA}}^2 & \sigma_{G_{SaA},SaB} \\ \sigma_{G_{By},SaB} & \sigma_{G_{SaA},SaB} & \sigma_{G_{SaB}}^2 \end{bmatrix} \otimes \mathbf{K}\right)$, with $\sigma_{G_{X,Y}}$ being the

225 genetic covariance between groups X and Y (the letters X, Y are further used as
 226 group names when not specifically designating given groups), \mathbf{K} the same kinship

227 matrix computed following Eq. (2), and $\mathbf{e} = \begin{bmatrix} \mathbf{e}_{By} \\ \mathbf{e}_{SaA} \\ \mathbf{e}_{SaB} \end{bmatrix}$ is the vector of errors of size

228 N where: $\begin{bmatrix} \mathbf{e}_A \\ \mathbf{e}_B \\ \mathbf{e}_C \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} \mathbf{I}_{By}\sigma_{E_{By}}^2 & 0 & 0 \\ 0 & \mathbf{I}_{SaA}\sigma_{E_{SaA}}^2 & 0 \\ 0 & 0 & \mathbf{I}_{SaB}\sigma_{E_{SaB}}^2 \end{bmatrix}\right)$, \mathbf{I}_X is the identity

229 matrix with dimensions equal to the number of observations from group X and
 230 $\sigma_{E_X}^2$ is the error variance in group X . The choice of hyper-parameters for the
 231 inference was done following Lehermeier et al. (2015).

232 *Bayes-B* The Bayesian shrinkage regression Bayes-B proposed by Meuwissen et al.
 233 (2001) was applied using the R package “BGLR” (Pérez and de los Campos, 2014).

234 In Bayes-B, only a proportion of a markers can have a non-zero effect with a
 235 variance specific to each marker. This modeling represents genetic architectures
 236 for which some SNPs are not associated to any QTL while others are associated
 237 to QTL with potentially large effects. In this model, $\mathbf{g} = \mathbf{W}\mathbf{u}$, \mathbf{W} is the centered
 238 genotyping matrix and \mathbf{u} is the vector of marker effects where the prior distribution
 239 of each u_m is the following mixture distribution:

$$P(u_m|\pi) = \begin{cases} 0 & \text{with probability } \pi, \\ t(0, \nu, S^2) & \text{otherwise,} \end{cases}$$

² available at <https://github.com/QuantGen/MTM>

240 where π is the proportion of marker with null effect, $t(0, \nu, S^2)$ is the scaled-t
 241 distribution with ν and S^2 being the number of degrees of freedom and the scale
 242 parameter, respectively. Other terms are identical to those of GBLUP.

243 *RKHS* The reproducing kernel Hilbert space (RKHS) semiparametric approach
 244 for genomic prediction (Gianola et al., 2006; Gianola and van Kaam, 2008) was
 245 applied using R package “BGLR” (Pérez and de los Campos, 2014). This approach
 246 combines a classical additive genetic model with a kernel function which converts
 247 predictor variables into set of distances among observations. The RKHS model
 248 based on a Gaussian kernel has been demonstrated to capture epistatic effects
 249 between markers (Jiang and Reif, 2015). In this model, $\mathbf{g} = \mathbf{K}_h \boldsymbol{\alpha}$, \mathbf{K}_h is the matrix
 250 of kernel entries and $\boldsymbol{\alpha}$ is the vector of individual effects. Other terms are identical
 251 to those of GBLUP. The kernel function implemented here was a Gaussian kernel:

$$\mathbf{K}_h(\mathbf{W}_i, \mathbf{W}_j) = e^{-hd_{i,j}},$$

252 where h is a smoothing parameter and $d_{i,j}$ is the marker-based Euclidean distance
 253 between individuals i and j . The value of the smoothing parameter was chosen
 254 following the kernel averaging method proposed in Pérez and de los Campos (2014).

255 *Random Forest* The tree-based machine learning approach called Random Forest
 256 (Breiman, 2001; Chen and Ishwaran, 2012) was applied using the R package
 257 “randomForest” (Liaw and Wiener, 2007). The grouping property of trees enables
 258 the Random Forest to adequately deal with correlations and interactions between
 259 predictor variables (Chen and Ishwaran, 2012). In this approach, the vector \mathbf{y}
 260 of phenotypes was used as the vector of response variable while the centered
 261 genotyping matrix \mathbf{W} was used as the matrix of predictor variables.

262 For the GP models based on a Bayesian inference (MGBLUP, Bayes-B and
 263 RKHS), 10,000 MCMC iterations were considered with a burn-in of 1,000 iterations
 264 and a thinning of 3 (i.e., one out of three samples were conserved to compute
 265 posterior means).

266 Evaluation of the predictive ability of genomic prediction

267 The precision of the models was evaluated using three different cross-validation
 268 (CV) procedures where the predictive ability was calculated by correlating the
 269 predictions of breeding values of the TS to the reference phenotypes (i.e., the
 270 LSmeans).

271 The first CV procedure, referred to as holdout cross-validation (HO-CV), was
 272 performed by repeatedly splitting (x 100) the oat population into a TRS and a
 273 TS with proportions being 4/5 and 1/5, respectively. This CV procedure makes it
 274 possible to study the level of precision that can be obtained when neglecting the
 275 role of genetic structure.

276 The second CV procedure, referred to as leave-one-out cross-validation (LOO-
 277 CV), was performed by predicting the breeding value of each individual using a
 278 model trained on all the remaining individuals. It allowed a joint graphic repre-
 279 sentation of the quality of prediction of all individuals depending on their genetic
 280 group.

281 The third CV procedure, referred to as structured-holdout cross validation
 282 (SHO-CV), allowed to study the impact of genetic structure on the predictive ability,
 283 as presented in Rio et al. (2019). In SHO-CV, group-specific TSs of 49 individuals
 284 were predicted using a model calibrated on 150 other individuals. Depending on
 285 the scenario, the training set included either members of a single group (e.g.,
 286 150 individuals from the Byzantina group), or of the three groups in balanced
 287 proportions (i.e., 50 individuals from each group). The sampling was repeated 100
 288 times for each scenario.

289 Training set optimization

290 For TRS optimization, the oat population was repeatedly split (x 30) into a
 291 candidate set (CS) of N_{CS} individuals and a TS of N_{TS} individuals (Fig. 2).
 292 Both the CS and the TS were defined differently depending on the optimization
 293 scenario (see below). The TRS individuals were then selected among CS individuals
 294 based on PEVmean, CDmean or PAM (see below) and using a genetic algorithm
 295 implemented in the R package “STPGA” (Akdemir, 2017), using 500 iterations.
 296 For each TS, 30 random TRS were sampled as a benchmark. For each trait, the
 297 breeding values of TS individuals were then predicted using a GBLUP model
 298 trained on TRS individuals, and the predictive ability was calculated to compare
 299 optimization methods. Three different criteria were considered:

300 *PEVmean and CDmean* The PEVmean and CDmean optimization criteria are
 301 derived from the GBLUP linear mixed model. The PEVmean criterion is defined
 302 as the mean of the predictor error variance over a set of individuals, where each
 303 PEV_i can be computed as:

$$PEV_i = \text{Var}(\mathbf{g}_i - \hat{\mathbf{g}}_i) = \left(\mathbf{ZM}\mathbf{Z}^T + \mathbf{K}^{-1}\lambda \right)_{i,i}^{-1} \times \sigma_E^2,$$

304 where \mathbf{g}_i is the breeding value of i , $\hat{\mathbf{g}}_i$ is the best linear unbiased prediction (BLUP)
 305 of \mathbf{g}_i , $\mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-}\mathbf{X}^T$ is an orthogonal projector on the subspace spanned
 306 by the columns of \mathbf{X} where $(\mathbf{X}^T\mathbf{X})^{-}$ is a generalized inverse of $\mathbf{X}^T\mathbf{X}$ (Laloë, 1993),
 307 and $\lambda = \frac{\sigma_E^2}{\sigma_G^2}$. All other terms correspond to those described in the GBLUP model.

308 The CDmean criterion is defined as the mean of the coefficient of determination
 309 (i.e., the square correlation between the breeding value of an individual and its
 310 corresponding prediction) where each individual CD can be computed as:

$$CD_i = \text{cor}(\mathbf{g}_i, \hat{\mathbf{g}}_i)^2 = \frac{\left(\mathbf{K} - \lambda (\mathbf{ZM}\mathbf{Z}^T + \mathbf{K}^{-1}\lambda)^{-1} \right)_{i,i}}{\mathbf{K}_{i,i}}$$

311 In this study, both PEVmean and CDmean optimizations were “targeted” (Akdemir
 312 and Isidro-Sánchez, 2019), meaning that the criteria were computed directly over
 313 the TS individuals. Note that finding the best TRS is be done by minimizing
 314 PEVmean while it is done by maximizing CDmean.

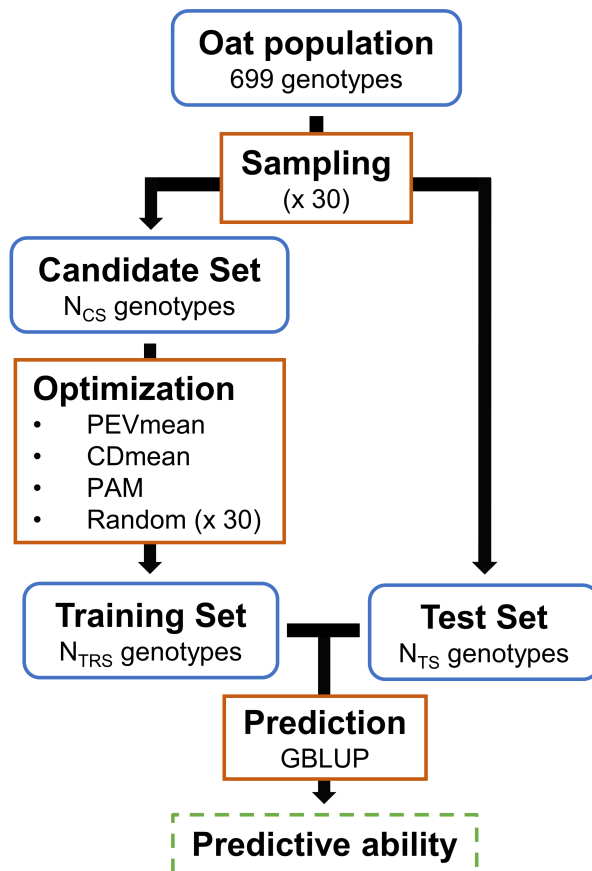


Fig. 2 Diagram illustrating the training set (TRS) optimization procedure where the oat population was repeatedly split (x 30) into a candidate set (CS) of N_{CS} individuals and a test set (TS) of N_{TS} individuals. Training set individuals were selected among CS individuals using a genetic algorithm (STPGA) and different methods: (i) the mean of the prediction error variance (PEVmean) of the TS, (ii) the mean of the coefficient of determination (CDmean) of the TS and (iii) partitioning around medoids (PAM). For each TS, 30 random TRS were sampled as a benchmark. The breeding values of TS individuals were then predicted using a GBLUP model trained on TRS individuals for each trait, before calculating the predictive ability for validation.

315 *PAM* Partitioning around medoids is a clustering method that classifies individuals
 316 into clusters by minimizing the sum of dissimilarities between the individuals of each
 317 cluster, and designating a central individual, or medoid of that cluster (Kaufman
 318 and Rousseeuw, 1987). The application of PAM to TRS optimization and genomic
 319 data was first presented by Guo et al. (2019).

320 Two optimization scenarios were considered in this study:

321 *Standard optimization* Test sets of 99 individuals were randomly sampled (i.e.,
322 by neglecting genetic structure). The CS consisted of all the remaining individ-
323 uals. Optimized TRSs were selected for a gradient of TRS sizes (i.e., $N_{TRS} \in$
324 $\{10, 30, 60, 100, 200, 300, 600\}$). Note that when $N_{TRS} = 600$, all CS individuals
325 are included in the TRS.

326 *Structure-based optimization.* We sampled group-specific TSs of 49 individuals
327 (e.g., 49 individuals from the Byzantina group). The CS consisted of 150 remaining
328 individuals from each genetic group. Optimized TRSs were selected for a gradient
329 of TRS sizes (i.e., $N_{TRS} \in \{9, 15, 30, 60, 90, 120, 150\}$). Note that $N_{TRS} = 150$
330 corresponds to the limit for which only individuals from the same groups as TS
331 individuals can be selected during the optimization procedure.

332 Results

333 Phenotypic characterization of the population

334 Our oat population has been evaluated for four agronomic traits: Heading, Height, 335
Biomass and Yield. As expected, agronomic performance of the collection varied 336
overall in the different environments. In addition to the different altitudes and 337 soil
structure of the different sites already stated, Co17 was characterised by a 338 mean
maximum T^a of 21.49°C, a mean minimum of 7.51°C and a rainfall of 415 339 mm during
the growing season. Co18 was slightly warmer with a mean maximum 340 T^a of 21.50°C, a
mean minimum of 8.46°C and a rain of 497 mm during the 341 growing season. Sa18, was
slightly colder and rainier than any of the environments 342 of Cordoba, with a mean
maximum T^a of 20.07 °C, a mean minimum of 7.64°C 343 and a rainfall of 513 mm during
the growing season (Online Resource Table S1). 344 In these environments, mean
Heading values ranged from 174 days to heading 345 at Co18 to 150 days at Co17, with a
minimum of 118 days and a maximum of 346 200 days to heading. Mean Height values
ranged from 115 cm at Co18 to 139 cm 347 at Sa18 with the shortest accession reaching 68
cm and the longest reaching 191 348 cm. Regarding Biomass, mean values ranged from
6832 kg/ha at Co17 to 2893 at 349 Co18 with minimum values of 190 kg/ha and
maximum of 11766 Kg/ha. Mean 350 yield ranged from 2326 Kg/ha at Sa18 to 983 kg/ha
at Co18, with the highest 351 yielding accessions reaching 5326 Kg/ha and the lowest one
yielding 180 Kg /ha. All 352 evaluated traits showed moderate to high broad-sense
heritabilities, with variations 353 depending on the genetic group (Table 1). For instance, a
broad-sense heritability 354 of 0.57, 0.74 and 0.55 were estimated for Biomass in the
Byzantina, Sativa A and 355 Sativa B groups, respectively. Those variations could be
connected to the differences 356 in genetic variances observed between groups (i.e., the
larger the group-specific

357 genetic variance $\sigma_{G_X}^2$, the higher the group-specific broad-sense heritability H_X^2).
358 The GxE variances were not large for most traits and also variable depending on
359 the genetic group. In terms of means, the performances of each genetic group were
360 comparable for most traits, but the Byzantina group outperformed both Sativa
361 groups for Yield and Biomass.

Table 1 Group-specific means, variances and broad-sense heritabilities estimated in the phenotypic analysis for Heading (in 100 growing degree-days (GDD)), Height (in cm), Biomass (in t/ha) and Yield (in t/ha).

| | Heading | Height | Biomass | Yield |
|----------------------------------|---------|--------|---------|-------|
| μ_{By} | 9.62 | 131.33 | 5.91 | 1.85 |
| μ_{SaA} | 10.71 | 124.63 | 5.19 | 1.43 |
| μ_{SaB} | 10.60 | 133.54 | 5.25 | 1.25 |
| σ_{GBy}^2 | 0.76 | 76.43 | 0.95 | 0.14 |
| σ_{GSaA}^2 | 2.01 | 113.13 | 2.38 | 0.46 |
| σ_{GSaB}^2 | 1.04 | 290.85 | 1.38 | 0.14 |
| $\sigma_{(G \times \beta)By}^2$ | 0.04 | 0.00 | 0.52 | 0.11 |
| $\sigma_{(G \times \beta)SaA}^2$ | 0.28 | 79.98 | 0.84 | 0.12 |
| $\sigma_{(G \times \beta)SaB}^2$ | 0.18 | 11.44 | 1.72 | 0.06 |
| σ_E^2 | 0.27 | 189.33 | 4.95 | 0.32 |
| H_{By}^2 | 0.95 | 0.78 | 0.57 | 0.65 |
| H_{SaA}^2 | 0.94 | 0.70 | 0.74 | 0.86 |
| H_{SaB}^2 | 0.92 | 0.92 | 0.55 | 0.72 |

362 Predictive ability of genomic prediction models

363 The GP predictive ability was evaluated using HO-CV for different GP models:
 364 GBLUP, Bayes-B, RKHS, Random-Forest and MGBLUP (Fig. 3). The mean
 365 predictive abilities averaged over all models were 0.87, 0.72, 0.63 and 0.81 for
 366 Heading, Height, Biomass and Yield, respectively. In terms of model comparison,
 367 GBLUP, RKHS and MGBLUP showed similar predictive abilities for all traits.
 368 Bayes-B showed similar predictive ability to GBLUP (our reference model) for
 369 Heading, Height and Yield, and higher predictions than GBLUP for Biomass (0.68
 370 vs. 0.60).

371 Training set optimization targeting random test sets

372 Three TRS optimization methods (PEVmean, CDmean and PAM) were compared
 373 to random sampling for a gradient of TRS sizes and validated using the four traits
 374 (Fig. 4). For all traits, the predictive capability increased with the size of the TRS,
 375 but at a rate that decreased as the size of the TRS increased. Optimizations of
 376 TRS based on PEVmean and CDmean performed very similarly and generally
 377 allowed for higher gains compared to the optimization based on PAM. For Height
 378 and $N_{TRS} = 100$, TRSs selected by CDmean and PEVmean showed a mean
 379 predictive ability of 0.70 compared to 0.61 and 0.60 for PAM and random sampling,
 380 respectively. The gains obtained with the optimization based on PAM were very
 381 variable depending on the trait, and even led to a substantial loss in predictive
 382 ability for Biomass (e.g., for $N_{TRS} = 100$, TRSs selected by PAM showed a
 383 mean predictive ability of 0.35 compared to 0.43 for random sampling). The mean
 384 group proportions within selected TRSs did not reveal major differences between
 385 optimization methods (Online Resource Fig. S3).

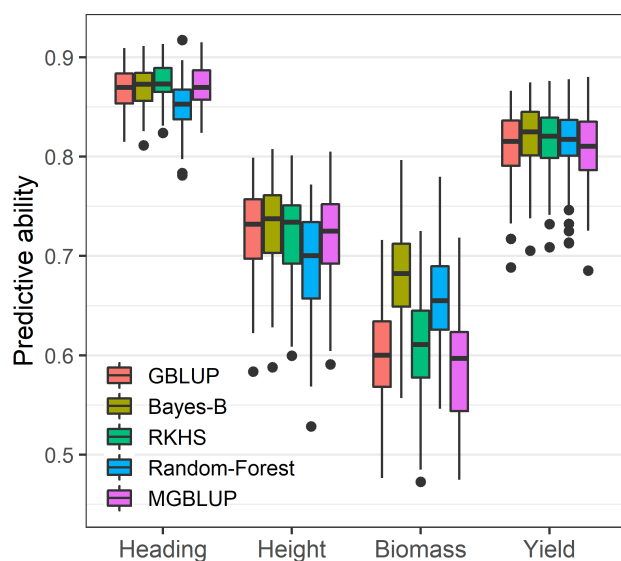


Fig. 3 Boxplots of predictive abilities obtained for each trait using holdout cross-validations (x 100 replicates) and five GP models: GBLUP, Bayes-B, RKHS, Random-Forest and multi-group GBLUP (MGBLUP).

386 Genetic structure and predictive ability

387 The impact of genetic structure on the predictive ability was first investigated
 388 graphically using LOO-CV (Online Resource Fig. S4). This approach confirmed
 389 the ability of the GBLUP model to predict differences in mean between groups and
 390 suggested an ability of the model to predict beyond this simple effect of genetic
 391 structure. We then investigated the impact of genetic structure using SHO-CV
 392 to explore within/across/multi-group scenarios (Fig. 5). In general, to predict
 393 a group-specific TS, the best strategy was to train a model using individuals
 394 from the same genetic group. The worst predictive abilities were obtained by
 395 training a model using the Byzantina group to predict any of the Sativa groups, or
 396 vice versa. Interestingly, a negative mean correlation of -0.55 was obtained when
 397 predicting the Sativa group using the Byzantina group for Heading. In general, the
 398 predictive abilities obtained with across-group scenarios between the Sativa_A and
 399 Sativa_B groups were moderate but not always symmetric. For instance, with Yield,
 400 a predictive ability of 0.58 was achieved when training a model on the Sativa_B
 401 group to predict the Sativa_A group, while a predictive ability of 0.22 was obtained
 402 for the opposite scenario. Multi-group TRSs always allowed for moderate to high
 403 predictive abilities regardless of the targeted TS.

404 Training set optimization targeting group-specific test sets

405 The impact of genetic structure on TRS optimization was investigated using
 406 structure-based optimization scenarios (Fig. 6). Two optimization methods (CD-

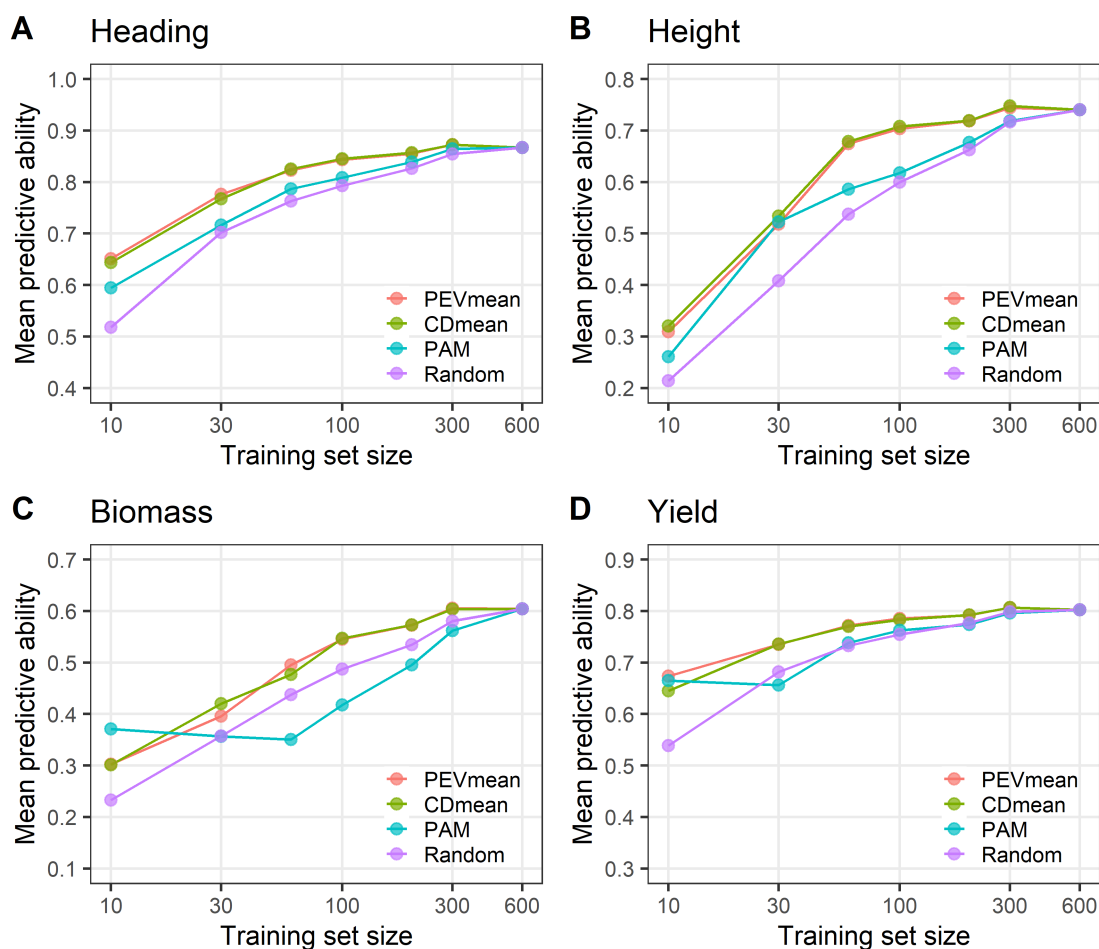


Fig. 4 Plots of mean GBLUP predictive abilities over random test sets (TSs) of 99 individuals (x 30 replicates) according to the size of the training set (TRS) for **(A)** Heading, **(B)** Height, **(C)** Biomass and **(D)** Yield. Different optimization methods were compared: PEVmean, CDmean and PAM, along with random sampling as a benchmark (x 30 replicates for each TS).

407 mean and PAM) were compared to random sampling for selecting TRSs that best
 408 predict group-specific TSs. The optimization based on CDmean always led to higher
 409 gains compared to PAM. The differences between optimization methods could
 410 largely be explained by the ability of CDmean to preferentially select individuals
 411 from the same genetic group as the one in the TS, unlike PAM or random sam-
 412 pling (Online Resource Fig. S5). However, the observed gains were very variable
 413 depending on the group-specific TS and the trait. For $N_{TRS} = 30$ and Yield, the
 414 gains in predictive ability from selecting TRSs using CDmean compared to random
 415 was +0.35 for Byzantina TSs compared to +0.03 for Sativa_A TSs (Fig. 6).

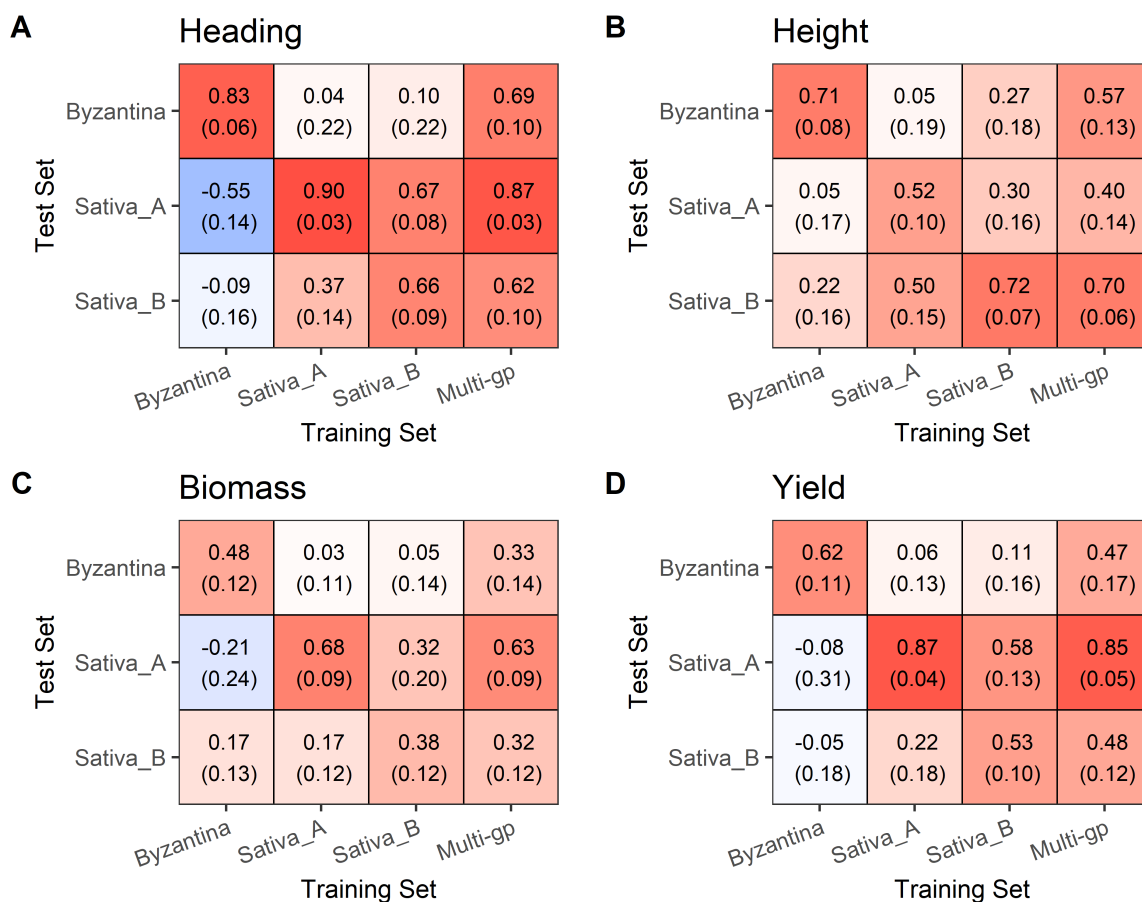


Fig. 5 Heatmaps of mean GBLUP predictive abilities obtained using the structured holdout cross-validations (x 100 replicates) for **(A)** Heading, **(B)** Height, **(C)** Biomass and **(D)** Yield. Group-specific training sets (TRSs) and multi-group TRSs of 150 individuals are indicated on the x axis while group-specific test sets of 49 individuals are indicated on the y axis. Standard deviations are shown between brackets.

416 Discussion

417 Oat genomics-assisted breeding in the Mediterranean basin

418 The GP predictive ability of oat agronomic and phenological traits has only been
 419 the subject of a few studies in the last past years (Asoro et al., 2011; Bekele et al.,
 2018; Haikka et al., 2020b,a; Mellers et al., 2020). The moderate to high predictive
 421 abilities obtained in our study are comparable to the ones obtained in the latter
 422 studies. It confirms the value of GBS-SNP markers as a genotyping technology for
 423 implementing GS in oat breeding programs, as proposed by Huang et al. (2014).
 424 The similar performances achieved by the different GP models tested in this study
 425 is a common feature in the GP literature (Heslot et al., 2012; de los Campos
 426 et al., 2013). The large gain in predictive ability obtained by applying Bayes-B for

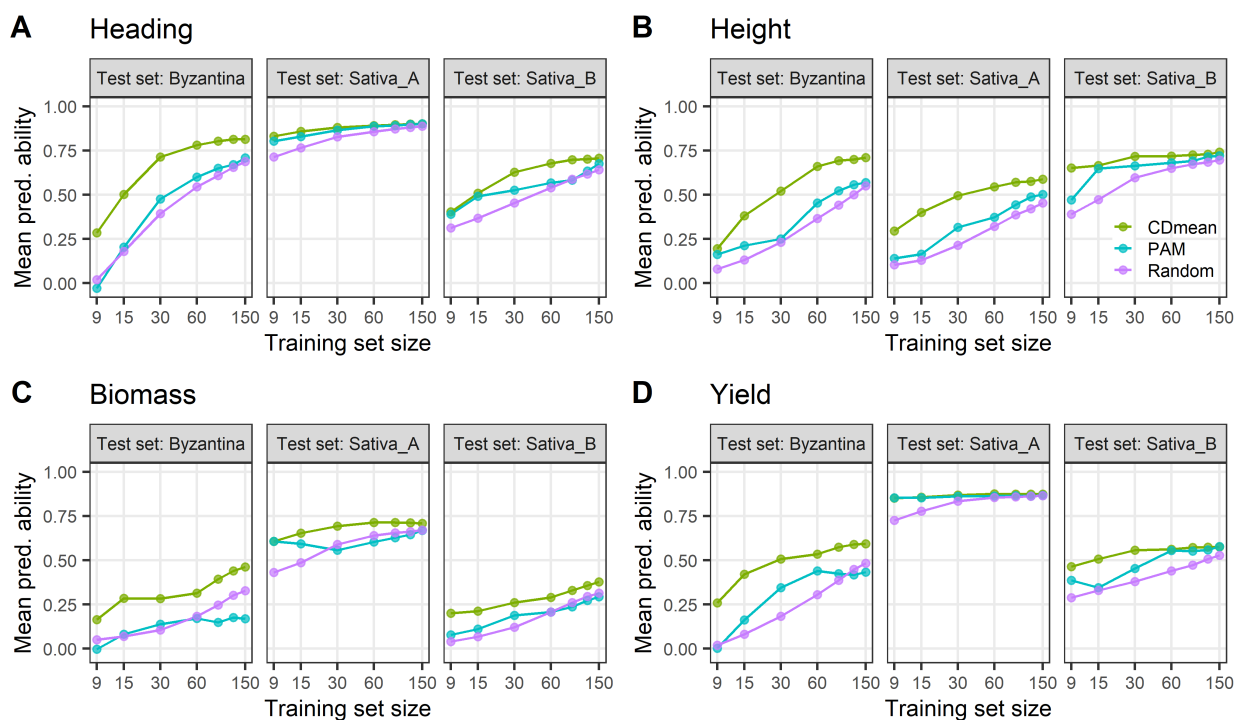


Fig. 6 Plots of mean GBLUP predictive abilities (Mean pred. ability) over group-specific test sets (TS) of 49 individuals (x 30 replicates) according to the size of the training set (TRS) for (A) Heading, (B) Height, (C) Biomass and (D) Yield. Two optimization methods were compared: CDmean and PAM, along with random sampling with equal group contributions in the TRS as a benchmark (x 30 replicates for each TS). Each candidate set consisted of 450 individuals, including 150 from each of the three genetic groups.

427 Biomass may result from the existence of a QTL with large effect segregating in
 428 the population that can better be accounted for using Bayes-B (Meuwissen et al.,
 429 2001; Pérez and de los Campos, 2014).

430 To improve the productivity of oats in the hot and dry environment of the
 431 Mediterranean basin, breeding programs must not only be based on the use of
 432 efficient tools like GP, but also on the introgression of favorable alleles from various
 433 sources of genetic diversity. Landraces harbour a great genetic potential for oat
 434 improvement as they are endowed with a higher genetic variability compared to high-
 435 yielding cultivars (Montilla-Bascón et al., 2013; Sánchez-Martín et al., 2016; Winkler
 436 et al., 2016). The population evaluated in our study includes landraces and cultivars
 437 of both white and red oats. A restricted set of this broad germplasm has already been
 438 characterized in field trials in different countries of the Mediterranean rim (Sánchez-
 439 Martín et al., 2014) and showed potential to detect QTL for powdery mildew and
 440 crown rust resistance (Montilla-Bascón et al., 2015), as well as agronomic traits
 441 (Rispaill et al., 2018). In our study, all genetic groups showed a substantial genetic
 442 variance and comparable means for all traits. We can reasonably assume that
 443 different QTL are involved in the trait genetic variability depending on the genetic
 444 group. Since population structure results from differences in allele frequencies

445 between groups, a QTL may indeed be polymorphic in one group and contribute
446 to its genetic variance, while being fixed in another group. This suggests a great
447 potential of both white and red oat landraces for harnessing new alleles to be
448 introgressed into elite germplasm using pre-breeding methods based on GP (Gorjanc
449 et al., 2016; Allier et al., 2020).

450 Genomic prediction in a highly structured population

451 When assessing the predictive ability of GP within structured populations, the
452 observed predictive ability may result mostly from the ability of the model to
453 predict differences in mean between groups (Guo et al., 2014; Rio et al., 2019). For
454 all traits, the SHO-CV showed that the predictive ability was also moderate to
455 high when considering group-specific TSs, meaning that it would be possible to
456 identify the best individuals within each genetic group (Fig. 5). Our results were
457 in concordance with previous studies on GP in structured populations (Rio et al.,
458 2019). Thus, (i) a given group-specific TS was generally best predicted using a TRS
459 including only individuals from the same genetic group, (ii) applying across-group
460 predictions could highly depreciate the GP predictive ability, and (iii) a multi-group
461 TRS showed a high predictive ability, regardless of the targeted TS.

462 The genetic dissimilarity between the two oat subspecies *sativa* and *byzantina*
463 was illustrated by the difficulty to predict one of them using a model trained
464 on the other (Fig. 5). The substantial negative predictive ability obtained when
465 predicting the *Sativa_A* group using a model trained on the *Byzantina* group for
466 Heading is uncommon in GP. It may result from the existence of QTL segregating
467 in the population with effects of opposite signs depending on the genetic group.
468 Interestingly, previous results based on the same data (Canales et al., 2021a)
469 identified a marker associated with heading date (*avgbs_cluster_1918.1*) located
470 in chromosome 1D in the hexaploid reference genome. Comparison of the region
471 around the marker between the *byzantina* and *sativa* preliminary genome assemblies
472 have identified the *GAS4*-like gibberellin responsive gene, and a gene with homology
473 to *miR172* as potential candidates genes (Canales et al., 2021a). Conversely to what
474 was observed between *Byzantina* and *Sativa* groups, both *Sativa_A* and *Sativa_B*
475 were able to predict each other with moderate predictive abilities. This difference
476 can be explained by a greater genetic similarity between the *Sativa* groups compared
477 to that observed between the *Sativa* and *Byzantina* groups (Fig. 1C). These results
478 are in concordance with the existence of moderate population structure in white
479 oat, as observed in previous studies (Asoro et al., 2011; Newell et al., 2011, 2012;
480 Huang et al., 2014; Tumino et al., 2016; Esvelt Klos et al., 2016; Winkler et al., 2016;
481 Bjørnstad et al., 2017; Haikka et al., 2020b,a; Isidro-Sánchez et al., 2020a). They
482 also illustrate the need for the evaluation of population structure when applying
483 GP to a broad diversity (Isidro et al., 2015; Guo et al., 2019). Defining a TRS must
484 be done by selecting individuals from the same genetic groups as those represented
485 in the target population to maximize the predictive ability of GPs. If the target
486 population is not clearly identified, the best strategy is to define a generic TRS
487 that include all genetic groups (de Roos et al., 2009; Rio et al., 2019).

488 Training set optimization in a highly structured population

489 Training set optimization methods like CDmean, PEVmean and PAM can be used
490 to select a subset of individuals to be evaluated when budget limitations limit
491 the possibility to evaluate all possible individuals through extensive field trials
492 (see Akdemir and Isidro-Sánchez (2019) and references herein). Unlike previous
493 TRS optimization results in highly structured population presented by Isidro et al.
494 (2015), in our results CDmean and PEVmean allowed for substantial gains compared
495 to random sampling for all traits. This might be explained by several differences
496 regarding the methodology: (i) the criteria were computed directly on the TS using
497 the targeted optimization recommended by Akdemir and Isidro-Sánchez (2019)
498 rather than on the remaining candidates, and (ii) the optimization algorithm was
499 the genetic algorithm implemented in the “STPGA” R package (Akdemir, 2017)
500 rather than an exchange algorithm. Note also that the CD and the PEV computed
501 were those associated with the prediction of each breeding value and not with a
502 contrast between a set of breeding values. A possible extension of the CDmean and
503 PEVmean optimization criteria could be to compute the CD and PEV associated
504 with contrasts between the breeding value of each TS individual and the mean of
505 the TS, as recommended by Rincent et al. (2017).

506 The performance of the optimization based on PAM proposed by Guo et al.
507 (2019) was highly variable and could even lead to predictive abilities that were worse
508 than those obtained by random sampling. Our results illustrate the superiority
509 of using criteria that are directly connected to the quantity of interest (e.g., CD
510 corresponds to the model-based square correlation between the breeding value
511 of an individual and its prediction) rather than using heuristic approaches like
512 PAM. The aim of the optimization based on PAM is indeed to maximize the GP
513 predictive ability by maximizing the genetic distances between TRS individuals.
514 The differences between PAM and CDmean/PEVmean optimizations show that a
515 targeted optimization based on CDmean/PEVmean not only aims at maximizing
516 the distances between TRS individuals, but also implicitly accounts for other
517 criteria such as the minimization of the genetic distances between TRS and TS
518 individuals (Pszczola et al., 2012; Albrecht et al., 2011; Clark et al., 2012). The
519 superiority of the targeted CDmean compared to PAM was also shown by its
520 ability to select the TRS according to the nature of the target. For instance, if
521 the target consists only of Byzantina individuals, then a straightforward strategy
522 would be to preferentially select Byzantina individuals to form the TRS. While
523 this strategy is implicitly applied using a targeted CDmean optimization, the PAM
524 optimization does not account for any information of the TS, and may lead to
525 selecting individuals that are poorly connected to the targeted population.

526 **Conclusion**

527 This manuscript presents results on the GP of key agronomic traits in a diverse
528 populations of Mediterranean oat cultivars and landraces. The consistency between
529 the structure of the training population and the population to be predicted was
530 key to the predictive ability of genomic predictions. Regarding TRS optimization,
531 the superiority of CDmean and PEVmean compared to PAM was illustrated by
532 their ability to adapt the representation of each genetic group according to those

533 represented in the population to be predicted. Our findings are useful for future
534 studies that aims to implement genomics-assisted breeding tools in presence of
535 high population structure in oat and other species.

536 **Funding**

537 This work was supported by the Spanish Ministry of Science and Innovation
538 [PID2019-104518RB-100], (AEI/FEDER, UE) and the regional government through
539 the AGR-253 group and the European Regional and Social Development Funds.
540 LGS is holder of a FPI fellowship from the Spanish Ministry of Economy and
541 Competitiveness [BES-2017-080152]. This project has received funding from the
542 European Union's Horizon 2020 research and innovation program under grant
543 agreement No 818144, and also the Severo Ochoa Program for Centres of Excellence
544 in R&D. JIS was supported by the Beatriz Galindo Program (BEAGAL18/00115)
545 from the Ministerio de Educación y Formación Profesional of Spain and the Severo
546 Ochoa Program for Centres of Excellence in R&D from the Agencia Estatal de
547 Investigación of Spain, grant SEV-2016-0672 (2017-2021) to the CBGP.

548 **Author contribution statement**

549 SR, JIS and EP conceived the study. FJC and GMB participated in the genotyping
550 and phenotyping. SR and LGS performed statistical analyses. SR, LGS, JIS and
551 EP drafted the manuscript. All authors discussed the results and reviewed the
552 manuscript.

553 **Conflict of interest**

554 The authors declare that they have no conflict of interest.

555 **Ethical standards**

556 The authors declare that the experiments comply with the current laws of the
557 countries in which the experiments were performed.

558 **References**

- 559 Akdemir, D. (2017). Stpga: Selection of training populations with a genetic
560 algorithm. *bioRxiv*.
- 561 Akdemir, D. and Isidro-Sánchez, J. (2019). Design of training populations for
562 selective phenotyping in genomic prediction. *Scientific Reports*, 9(1):1446.
- 563 Akdemir, D., Sanchez, J. I., and Jannink, J.-L. (2015). Optimization of genomic
564 selection training populations with a genetic algorithm. *Genetics Selection
565 Evolution*, 47(1):38.

- 566 Albrecht, T., Wimmer, V., Auinger, H.-J., Erbe, M., Knaak, C., Ouzunova, M.,
567 Simianer, H., and Schön, C.-C. (2011). Genome-based prediction of testcross
568 values in maize. *Theoretical and Applied Genetics*, 123(2):339.
- 569 Allier, A., Teyssèdre, S., Lehermeier, C., Moreau, L., and Charcosset, A. (2020).
570 Optimized breeding strategies to harness genetic resources with different perfor-
571 mance levels. *BMC Genomics*, 21(1).
- 572 Alvarenga, A. B., Veroneze, R., Oliveira, H. R., Marques, D. B., Lopes, P. S.,
573 Silva, F. F., and Brito, L. F. (2020). Comparing alternative single-step gblup
574 approaches and training population designs for genomic evaluation of crossbred
575 animals. *Frontiers in genetics*, 11:263.
- 576 Asoro, F., Newell, M., Beavis, W., Scott, P., Tinker, N., and Jannink, J.-L. (2013).
577 Genomic, marker-assisted, and pedigree-blup selection methods for beta-glucan
578 concentration in elite oat. *Crop Sci.*, 53:1894–1906.
- 579 Asoro, F. G., Newell, M. A., Beavis, W. D., Scott, M. P., and Jannink, J.-L. (2011).
580 Accuracy and training population design for genomic selection on quantitative
581 traits in elite north american oats. *The Plant Genome*, 4(2):132–144.
- 582 Bekele, W. A., Wight, C. P., Chao, S., Howarth, C. J., and Tinker, N. A. (2018).
583 Haplotype-based genotyping-by-sequencing in oat genome research. *Plant*
584 *Biotechnology Journal*, 16(8):1452–1463.
- 585 Berro, I., Lado, B., Nalin, R. S., Quincke, M., and Gutiérrez, L. (2019). Training
586 population optimization for genomic selection. *The Plant Genome*, 12(3):1–14.
- 587 Bjørnstad, Å., He, X., Tekle, S., Klos, K., Huang, Y.-F., Tinker, N. A., Dong, Y.,
588 and Skinnnes, H. (2017). Genetic variation and associations involving fusarium
589 head blight and deoxynivalenol accumulation in cultivated oat (*avena sativa* l.).
590 *Plant Breeding*, 136(5):620–636.
- 591 Brandariz, S. P. and Bernardo, R. (2018). Maintaining the accuracy of genomewide
592 predictions when selection has occurred in the training population. *Crop Science*,
593 58(3):1226–1231.
- 594 Breiman, L. (2001). Random forests. *Mach. Learn.*, 45(1):5–32.
- 595 Brøndum, R., Rius-Vilarrasa, E., Strandén, I., Su, G., Guldbrandtsen, B., Fikse,
596 W., and Lund, M. (2011). Reliabilities of genomic prediction using combined
597 reference data of the nordic red dairy cattle populations. *Journal of Dairy*
598 *Science*.
- 599 Canales, F. J. (2019). *Improving oat for adaptation to Mediterranean environments*.
600 PhD thesis, Universidad de Córdoba.
- 601 Canales, F. J., Montilla-Bascón, G., Bekele, W. A., Howarth, C. J., Langdon,
602 T., Risipail, N., Tinker, N. A., and Prats, E. (2021a). Population genomics of
603 mediterranean oat (*a. sativa*) reveals high genetic diversity and three loci for
604 heading date. *Theoretical and Applied Genetics*.
- 605 Canales, F. J., Montilla-Bascón, G., Bekele, W. A., Howarth, C., Langdon, T.,
606 Risipail, N., Tinker, N., and Prats, E. (2021b). Data set from: Population
607 genomics of mediterranean oat (*a. sativa*) reveals high genetic diversity and three
608 loci for heading date. *Dryad, Dataset*.
- 609 Carlson, M. O., Montilla-Bascon, G., Hoekenga, O. A., Tinker, N. A., Poland,
610 J., Baseggio, M., Sorrells, M. E., Jannink, J.-L., Gore, M. A., and Yeats, T. H.
611 (2019). Multivariate genome-wide association analyses reveal the genetic basis
612 of seed fatty acid composition in oat (*avena sativa* l.). *G3: Genes, Genomes,*
613 *Genetics*, 9(9):2963–2975.

- 614 Chaffin, A. S., Huang, Y.-F., Smith, S., Bekele, W. A., Babiker, E., Gnanesh, B. N.,
615 Foresman, B. J., Blanchard, S. G., Jay, J. J., Reid, R. W., et al. (2016). A
616 consensus map in cultivated hexaploid oat reveals conserved grass synteny with
617 substantial subgenome rearrangement. *The plant genome*, 9(2):1–21.
- 618 Chen, L., Schenkel, F., Vinsky, M., Crews, D. H., and Li, C. (2013). Accuracy of
619 predicting genomic breeding values for residual feed intake in angus and charolais
620 beef cattle. *Journal of Animal Science*, 91:4669–4678.
- 621 Chen, X. and Ishwaran, H. (2012). Random forests for genomic data analysis.
622 *Genomics*, 99(6):323 – 329.
- 623 Clark, S. A., Hickey, J. M., Daetwyler, H. D., and van der Werf, J. H. (2012). The
624 importance of information on relatives for the prediction of genomic breeding
625 values and the implications for the makeup of reference data sets in livestock
626 breeding schemes. *Genetics Selection Evolution*, 44(1):4.
- 627 de los Campos, G., Hickey, J. M., Pong-Wong, R., Daetwyler, H. D., and Calus,
628 M. P. (2013). Whole-genome regression and prediction methods applied to plant
629 and animal breeding. *Genetics*, 193(2):327–345.
- 630 de Roos, A. P. W., Hayes, B. J., and Goddard, M. E. (2009). Reliability of genomic
631 predictions across multiple populations. *Genetics*, 183(4):1545–1553.
- 632 Duhnen, A., Gras, A., Teyssèdre, S., Romestant, M., Claustres, B., Daydé, J.,
633 and Mangin, B. (2017). Genomic selection for yield and seed protein content
634 in soybean: A study of breeding program data and assessment of prediction
635 accuracy. *Crop Sci.*, 57:1–13.
- 636 Edwards, S. M., Buntjer, J. B., Jackson, R., Bentley, A. R., Lage, J., Byrne, E.,
637 Burt, C., Jack, P., Berry, S., Flatman, E., et al. (2019). The effects of training
638 population design on genomic prediction accuracy in wheat. *Theoretical and
639 Applied Genetics*, 132(7):1943–1952.
- 640 Endelman, J. B. (2011). Ridge regression and other kernels for genomic selection
641 with r package rrblup. *The Plant Genome*, 4(3):250–255.
- 642 Esvelt Klos, K., Huang, Y.-F., Bekele, W. A., Obert, D. E., Babiker, E., Beattie,
643 A. D., Bjørnstad, Å., Bonman, J. M., Carson, M. L., Chao, S., et al. (2016).
644 Population genomics related to adaptation in elite oat germplasm. *The Plant
645 Genome*, 9(2):1–12.
- 646 Falush, D., Stephens, M., and Pritchard, J. K. (2003). Inference of population
647 structure using multilocus genotype data: Linked loci and correlated allele
648 frequencies. *Genetics*, 164(4):1567–1587.
- 649 FAO. (2017). *FAO statistical yearbook*. FAO.
- 650 Gianola, D., Fernando, R. L., and Stella, A. (2006). Genomic-assisted prediction
651 of genetic value with semiparametric procedures. *Genetics*, 173(3):1761–1776.
- 652 Gianola, D. and van Kaam, J. B. C. H. M. (2008). Reproducing kernel hilbert
653 spaces regression methods for genomic assisted prediction of quantitative traits.
654 *Genetics*, 178(4):2289–2303.
- 655 Gorjanc, G., Jenko, J., Hearne, S. J., and Hickey, J. M. (2016). Initiating maize
656 pre-breeding programs using genomic selection to harness polygenic variation
657 from landrace populations. *BMC Genomics*, 17(1).
- 658 Guo, T., Yu, X., Li, X., Zhang, H., Zhu, C., Flint-Garcia, S., McMullen, M. D.,
659 Holland, J. B., Szalma, S. J., Wissler, R. J., and Yu, J. (2019). Optimal designs
660 for genomic selection in hybrid crops. *Molecular Plant*, 12(3):390–401.
- 661 Guo, Z., Tucker, D. M., Basten, C. J., Gandhi, H., Ersoz, E., Guo, B., Xu, Z., Wang,
662 D., and Gay, G. (2014). The impact of population structure on genomic prediction

- 663 in stratified populations. *Theoretical and Applied Genetics*, 127(3):749–762.
- 664 Habier, D., Fernando, R. L., Kizilkaya, K., and Garrick, D. J. (2011). Extension of
665 the bayesian alphabet for genomic selection. *BMC Bioinformatics*, 12(1):186.
- 666 Haikka, H., Knurr, T., Manninen, O., Pietila, L., Isolahti, M., Teperi, E., Mantysaari,
667 E. A., and Strandén, I. (2020a). Genomic prediction of grain yield in commercial
668 Finnish oat (*Avena sativa*) and barley (*Hordeum vulgare*) breeding programmes.
669 *PLANT BREEDING*, 139(3):550–561.
- 670 Haikka, H., Manninen, O., Hautsalo, J., Pietila, L., Jalli, M., and Veteläinen, M.
671 (2020b). Genome-wide Association Study and Genomic Prediction for Fusarium
672 graminearum Resistance Traits in Nordic Oat (*Avena sativa* L.). *AGRONOMY-
673 BASEL*, 10(2).
- 674 Heslot, N., Yang, H., Sorrells, M. E., and Jannink, J. (2012). Genomic selection in
675 plant breeding: A comparison of models. *Crop Sci.*, 52:146–160.
- 676 Hickey, J. M., Dreisigacker, S., Crossa, J., Hearne, S., Babu, R., Prasanna, B. M.,
677 Grondona, M., Zambelli, A., Windhausen, V. S., Mathews, K., et al. (2014).
678 Evaluation of genomic selection training population designs and genotyping
679 strategies in plant breeding programs using simulation. *Crop Science*, 54(4):1476–
680 1488.
- 681 Huang, Y.-F., Poland, J. A., Wight, C. P., Jackson, E. W., and Tinker, N. A.
682 (2014). Using genotyping-by-sequencing (gbs) for genomic discovery in cultivated
683 oat. *PLOS ONE*, 9(7):1–16.
- 684 Isidro, J., Akdemir, D., and Burke, J. (2016). Genomic selection. In William, A.,
685 Alain, B., and Maarten, V. G., editors, *The world wheat book: a history of wheat
686 breeding*, volume 3, chapter 32, pages 1001–1023. Lavoisier, Paris.
- 687 Isidro, J., Jannink, J.-L., Akdemir, D., Poland, J., Heslot, N., and Sorrells, M. E.
688 (2015). Training set optimization under population structure in genomic selection.
689 *Theoretical and Applied Genetics*, 128(1):145–158.
- 690 Isidro-Sánchez, J., D’Arcy Cusack, K., Verheecke-Vaessen, C., Kahla, A., Bekele,
691 W., Doohan, F., Magan, N., and Medina, A. (2020a). Genome-wide association
692 mapping of fusarium langsethiae infection and mycotoxin accumulation in oat
693 (*avena sativa* l.). *The Plant Genome*, page e20023.
- 694 Isidro-Sánchez, J., Prats, E., Howarth, C., Langdon, T., and Montilla-Bascón, G.
695 (2020b). Genomic approaches for climate resilience breeding in oats. In *Genomic
696 Designing of Climate-Smart Cereal Crops*, pages 133–169. Springer.
- 697 Jiang, Y. and Reif, J. C. (2015). Modeling epistasis in genomic selection. *Genetics*,
698 201(2):759–768.
- 699 Karoui, S., Carabaño, M. J., Díaz, C., and Legarra, A. (2012). Joint genomic
700 evaluation of french dairy cattle breeds using multiple-trait models. *Genetics
701 Selection Evolution*, 44(1):39.
- 702 Kaufman, L. and Rousseeuw, P. (1987). *Clustering by Means of Medoids*. Delft
703 University of Technology : reports of the Faculty of Technical Mathematics and
704 Informatics. Faculty of Mathematics and Informatics.
- 705 Kebede, A. Z., Friesen-Enns, J., Gnanesh, B. N., Menzies, J. G., Fetch, J. W. M.,
706 Chong, J., Beattie, A. D., Paczos-Grzeda, E., and McCartney, C. A. (2019).
707 Mapping oat crown rust resistance gene pc45 confirms association with pckm.
708 *G3: Genes, Genomes, Genetics*, 9(2):505–511.
- 709 Laloë, D. (1993). Precision and information in linear models of genetic evaluation.
710 *Genetics Selection Evolution*, 25(6):557.

- 711 Laporte, F. and Mary-Huard, T. (2020). *MM4LMM: Inference of Linear Mixed*
712 *Models Through MM Algorithm*. R package version 2.0.2.
- 713 Lehermeier, C., Krämer, N., Bauer, E., Bauland, C., Camisan, C., Campo, L., Fla-
714 ment, P., Melchinger, A. E., Menz, M., Meyer, N., Moreau, L., Moreno-González,
715 J., Ouzunova, M., Pausch, H., Ranc, N., Schipprack, W., Schönleben, M., Walter,
716 H., Charcosset, A., and Schön, C.-C. (2014). Usefulness of multiparental popula-
717 tions of maize (*zea mays* l.) for genome-based prediction. *Genetics*, 198(1):3–16.
- 718 Lehermeier, C., Schön, C.-C., and de los Campos, G. (2015). Assessment of genetic
719 heterogeneity in structured plant populations using multivariate whole-genome
720 regression models. *Genetics*, 201(1):323–337.
- 721 Lorenz, A. J. and Smith, K. P. (2015). Adding genetically distant individuals to
722 training populations reduces genomic prediction accuracy in barley. *Crop science*,
723 55(6):2657–2667.
- 724 Mangin, B., Rincant, R., Rabier, C.-E., Moreau, L., and Goudemand-Dugue, E.
725 (2019). Training set optimization of genomic prediction by means of ethacc.
726 *PLOS ONE*, 14(2):1–21.
- 727 Maughan, P. J., Lee, R., Walstead, R., Vickerstaff, R. J., Fogarty, M. C., Brouwer,
728 C. R., Reid, R. R., Jay, J. J., Bekele, W. A., Jackson, E. W., Tinker, N. A.,
729 Langdon, T., Schlueter, J. A., and Jellen, E. N. (2019). Genomic insights from
730 the first chromosome-scale assemblies of oat (*avena* spp.) diploid species. *BMC*
731 *Biology*, 17(1).
- 732 Mellers, G., Mackay, I., Cowan, S., Griffiths, I., Martinez-Martin, P., Poland,
733 J. A., Bekele, W., Tinker, N. A., Bentley, A. R., and Howarth, C. J. (2020).
734 Implementing within-cross genomic prediction to reduce oat breeding costs. *The*
735 *Plant Genome*, 13(1):e20004.
- 736 Meuwissen, T. H. E., Hayes, B. J., and Goddard, M. E. (2001). Prediction of total
737 genetic value using genome-wide dense marker maps. *Genetics*, 157(4):1819–1829.
- 738 Montilla-Bascón, G., Rispail, N., Sánchez-Martín, J., Rubiales, D., Mur, L. A. J.,
739 Langdon, T., Howarth, C. J., and Prats, E. (2015). Genome-wide association
740 study for crown rust (*puccinia coronata* f. sp. *avenae*) and powdery mildew
741 (*blumeria graminis* f. sp. *avenae*) resistance in an oat (*avena sativa*) collection of
742 commercial varieties and landraces. *Frontiers in Plant Science*, 6.
- 743 Montilla-Bascón, G., Sánchez-Martín, J., Rispail, N., Rubiales, D., Mur, L., Lang-
744 don, T., Griffiths, I., Howarth, C., and Prats, E. (2013). Genetic diversity and
745 population structure among oat cultivars and landraces. *Plant molecular biology*
746 *reporter*, 31(6):1305–1314.
- 747 Newell, M., Cook, D., Tinker, N., and Jannink, J.-L. (2011). Population structure
748 and linkage disequilibrium in oat (*avena sativa* l.): implications for genome-wide
749 association studies. *Theoretical and Applied Genetics*, 122(3):623–632.
- 750 Newell, M. A., Asoro, F. G., Scott, M. P., White, P. J., Beavis, W. D., and Jannink,
751 J.-L. (2012). Genome-wide association study for oat (*avena sativa* l.) beta-glucan
752 concentration using germplasm of worldwide origin. *Theoretical and Applied*
753 *Genetics*, 125(8):1687–1696.
- 754 Norman, A., Taylor, J., Edwards, J., and Kuchel, H. (2018). Optimising genomic
755 selection in wheat: Effect of marker density, population size and population
756 structure on prediction accuracy. *G3: Genes, Genomes, Genetics*, 8(9):2889–
757 2899.
- 758 Olatoye, M. O., Clark, L. V., Labonte, N. R., Dong, H., Dwiyantri, M. S., Anzoua,
759 K. G., Brummer, J. E., Ghimire, B. K., Dzyubenko, E., Dzyubenko, N., et al.

- (2020). Training population optimization for genomic selection in miscanthus. *G3: Genes, Genomes, Genetics*, 10(7):2465–2476.
- Olson, K. M., Van Raden, P. M., and Tooker, M. E. (2012). Multibreed genomic evaluations using purebred holsteins, jersey, and brown swiss. *Journal of Dairy Science*, 95(9):5378–5383.
- Ou, J.-H. and Liao, C.-T. (2019). Training set determination for genomic selection. *Theoretical and Applied Genetics*, 132(10):2781–2792.
- Pérez, P. and de los Campos, G. (2014). Genome-wide regression and prediction with the bglr statistical package. *Genetics*, 198(2):483–495.
- Poland, J. A. and Rife, T. W. (2012). Genotyping-by-sequencing for plant breeding and genetics. *The Plant Genome*, 5(3):92–102.
- Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959.
- Pryce, J. E., Gredler, B., Bolormaa, S., Bowman, P. J., Egger-Danner, C., Fuerst, C., Emmerling, R., Solkner, J., Goddard, M. E., and Hayes, B. J. (2011). Short communication: Genomic selection using a multi-breed, across-country reference population. *Journal of Dairy Science*, 94(5):2625–2630.
- Pszczola, M., Strabel, T., Mulder, H., and Calus, M. (2012). Reliability of direct genomic values for animals with different relationships within and to the reference population. *Journal of dairy science*, 95(1):389–400.
- Rincent, R., Charcosset, A., and Moreau, L. (2017). Predicting genomic selection efficiency to optimize calibration set and to assess prediction accuracy in highly structured populations. *Theoretical and Applied Genetics*, 130(11):2231–2247.
- Rincent, R., Laloë, D., Nicolas, S., Altmann, T., Brunel, D., Revilla, P., Rodríguez, V., Moreno-Gonzalez, J., Melchinger, A., Bauer, E., Schoen, C.-C., Meyer, N., Giauffret, C., Bauland, C., Jamin, P., Laborde, J., Monod, H., Flament, P., Charcosset, A., and Moreau, L. (2012). Maximizing the reliability of genomic selection by optimizing the calibration set of reference individuals: Comparison of methods in two diverse groups of maize inbreds (*zea mays* l.). *Genetics*, 192(2):715–728.
- Rio, S., Mary-Huard, T., Moreau, L., and Charcosset, A. (2019). Genomic selection efficiency and a priori estimation of accuracy in a structured dent maize panel. *Theoretical and Applied Genetics*, 132(1):81–96.
- Rispail, N., Montilla-Bascón, G., Sánchez-Martín, J., Flores, F., Howarth, C., Langdon, T., Rubiales, D., and Prats, E. (2018). Multi-environmental trials reveal genetic plasticity of oat agronomic traits associated with climate variable changes. *Frontiers in Plant Science*, 9:1358.
- Roth, M., Muranty, H., Di Guardo, M., Guerra, W., Patocchi, A., and Costa, F. (2020). Genomic prediction of fruit texture and training population optimization towards the application of genomic selection in apple. *Horticulture research*, 7(1):1–14.
- Sánchez-Martín, J., Rispail, N., Flores, F., Emeran, A. A., Sillero, J. C., Rubiales, D., and Prats, E. (2016). Higher rust resistance and similar yield of oat landraces versus cultivars under high temperature and drought. *Agronomy for Sustainable Development*, 37(1).
- Sarinelli, J. M., Murphy, J. P., Tyagi, P., Holland, J. B., Johnson, J. W., Mergoum, M., Mason, R. E., Babar, A., Harrison, S., Sutton, R., et al. (2019). Training population selection and use of fixed effects to optimize genomic predictions in a historical usa winter wheat panel. *Theoretical and Applied Genetics*, 132(4):1247–

1261.

- 809 Sunstrum, F. G., Bekele, W. A., Wight, C. P., Yan, W., Chen, Y., and Tinker,
810 N. A. (2019). A genetic linkage map in southern-by-spring oat identifies multiple
811 quantitative trait loci for adaptation and rust resistance. *Plant breeding*, 138(1):82–
812 94.
- 814 Sánchez-Martín, J., Rubiales, D., Flores, F., Emeran, A., Shtaya, M., Sillero, J.,
815 Allagui, M., and Prats, E. (2014). Adaptation of oat (*avena sativa*) cultivars to
816 autumn sowings in mediterranean environments. *Field Crops Research*, 156:111 –
817 122.
- 818 Tayeh, N., Klein, A., Le Paslier, M.-C., Jacquin, F., Houtin, H., Rond, C., Chabert-
819 Martinello, M., Magnin-Robert, J.-B., Marget, P., Aubert, G., et al. (2015).
820 Genomic prediction in pea: effect of marker density and training population size
821 and composition on prediction accuracy. *Frontiers in plant science*, 6:941.
- 822 Technow, F., Burger, A., and Melchinger, A. E. (2013). Genomic prediction of
823 northern corn leaf blight resistance in maize with combined or separated training
824 sets for heterotic groups. *G3: Genes—Genomes—Genetics*, 3(2):197–203.
- 825 Tinker, N. A., Bekele, W. A., and Hattori, J. (2016). Haplotag: Software for
826 haplotype-based genotyping-by-sequencing analysis. *G3: Genes, Genomes, Ge-
827 netics*, 6(4):857–863.
- 828 Tinker, N. A., Chao, S., Lazo, G. R., Oliver, R. E., Huang, Y.-F., Poland, J. A.,
829 Jellen, E. N., Maughan, P. J., Kilian, A., and Jackson, E. W. (2014). A snp
830 genotyping array for hexaploid oat. *The Plant Genome*, 7(3):1–8.
- 831 Tumino, G., Voorrips, R. E., Morcia, C., Ghizzoni, R., Germeier, C. U., Paulo, M.-
832 J., Terzi, V., and Smulders, M. J. (2017). Genome-wide association analysis for
833 lodging tolerance and plant height in a diverse european hexaploid oat collection.
834 *Euphytica*, 213(8):163.
- 835 Tumino, G., Voorrips, R. E., Rizza, F., Badeck, F. W., Morcia, C., Ghizzoni, R.,
836 Germeier, C. U., Paulo, M.-J., Terzi, V., and Smulders, M. J. (2016). Population
837 structure and genome-wide association analysis for frost tolerance in oat using
838 continuous snp array signal intensity ratios. *Theoretical and Applied Genetics*,
839 129(9):1711–1724.
- 840 Welch, R. W. (2012). *The oat crop: production and utilization*. Springer Science &
841 Business Media.
- 842 Winkler, L. R., Michael Bonman, J., Chao, S., Admassu Yimer, B., Bockelman,
843 H., and Esvelt Klos, K. (2016). Population structure and genotype–phenotype
844 associations in a collection of oat landraces and historic cultivars. *Frontiers in
845 plant science*, 7:1077.
- 846 Yan, H., Bekele, W. A., Wight, C. P., Peng, Y., Langdon, T., Latta, R. G., Fu,
847 Y.-B., Diederichsen, A., Howarth, C. J., Jellen, E. N., et al. (2016). High-density
848 marker profiling confirms ancestral genomes of *avena* species and identifies
849 d-genome chromosomes of hexaploid oat. *Theoretical and Applied Genetics*,
850 129(11):2133–2149.
- 851 Yan, H., Zhou, P., Peng, Y., Bekele, W. A., Ren, C., Tinker, N. A., and Peng, Y.
852 (2020). Genetic diversity and genome-wide association analysis in chinese hulless
853 oat germplasm. *Theoretical and Applied Genetics*, pages 1–16.
- 854 Zhou, L., Ding, X., Zhang, Q., Wang, Y., Lund, M. S., and Su, G. (2013). Consis-
855 tency of linkage disequilibrium between chinese and nordic holsteins and genomic
856 prediction for chinese holsteins using a joint reference population. *Genetics
857 Selection Evolution*, 45(1):7.

Supplementary Materials

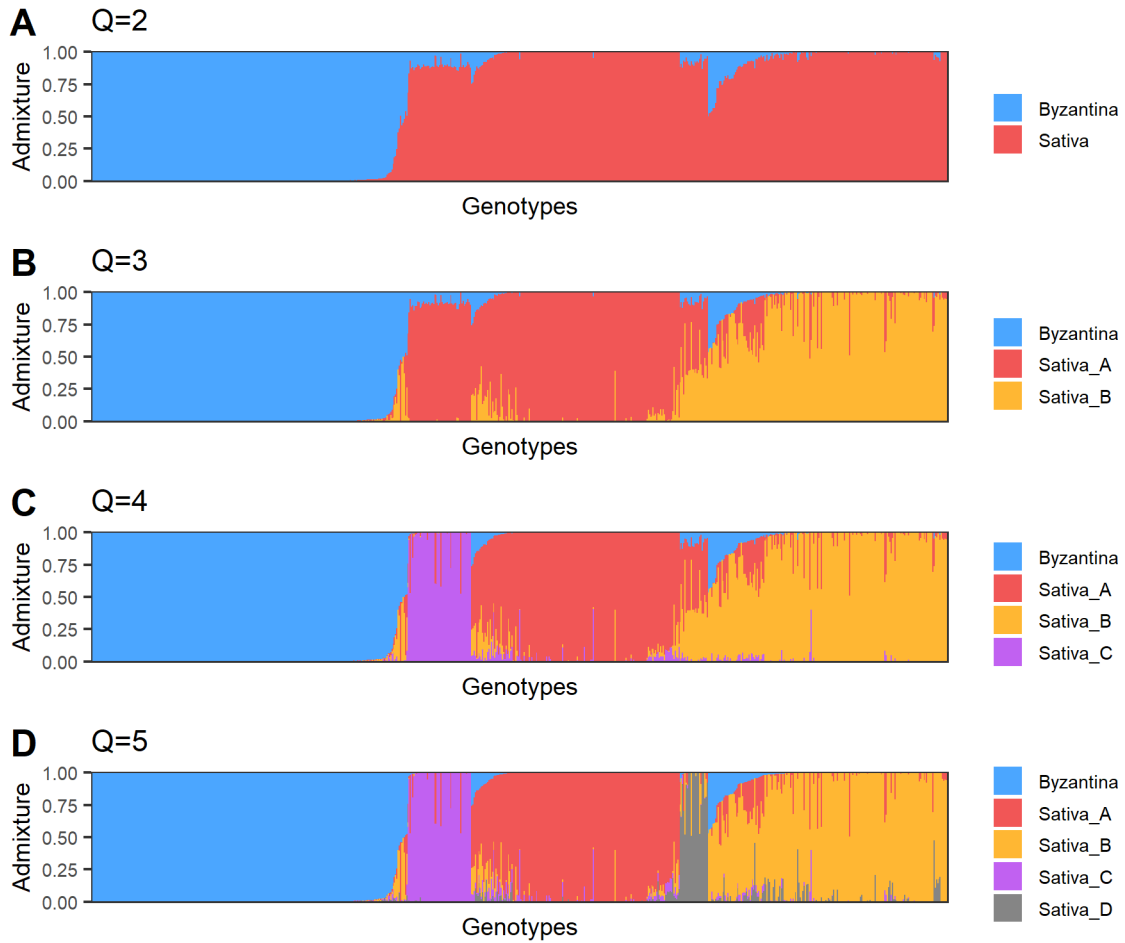
Genomic prediction and training set optimization in a structured Mediterranean oat population

Simon Rio^{1*}, Luis Gallego-Sánchez^{2*}, Gracia Montilla-Bascón², Francisco J. Canales², Julio Isidro y Sánchez¹, and Elena Prats²

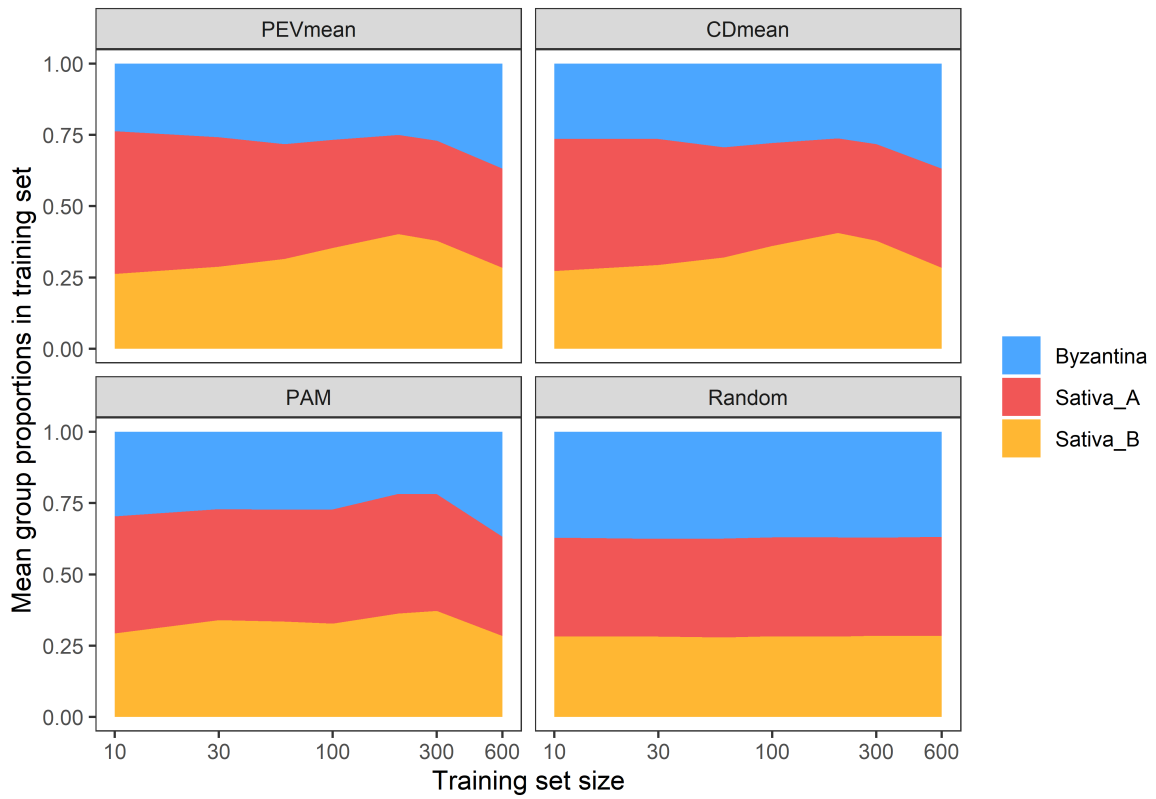
¹Centro de Biotecnología y Genómica de Plantas (CBGP, UPM-INIA) Universidad Politécnica de Madrid (UPM) - Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria (INIA) Campus de Montegancedo-UPM 28223-Pozuelo de Alarcón, (Madrid), Spain

²Institute for Sustainable Agriculture, Spanish Research Council (CSIC), Córdoba, Spain

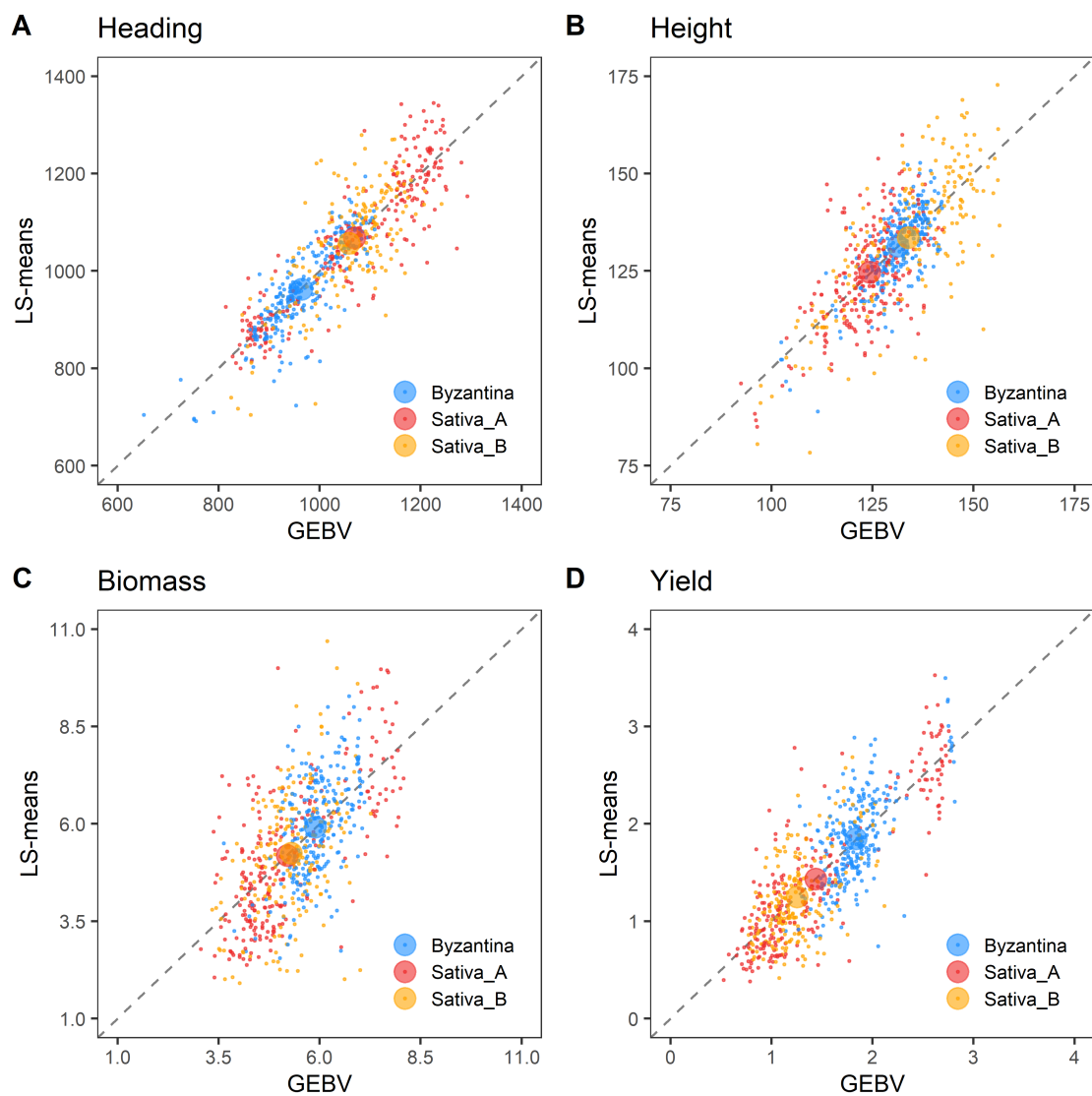
*These authors contributed equally to this study



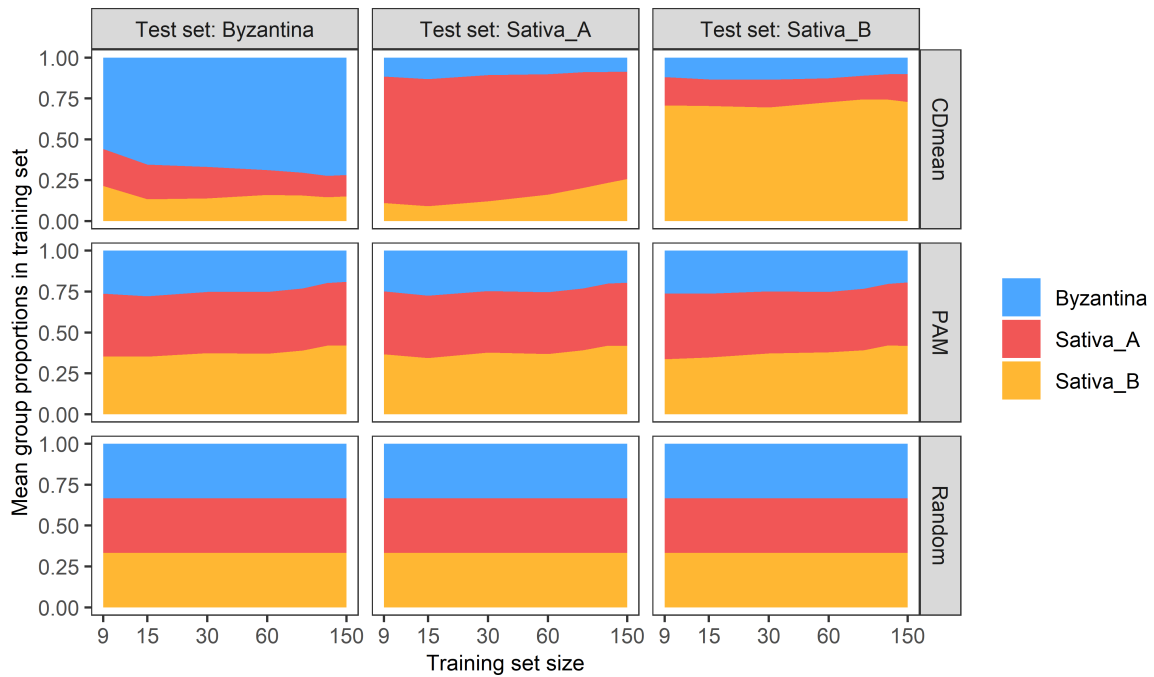
Supplementary Figure S1. Admixture barplots showing the admixture proportions of each individual and obtained using STRUCTURE for different numbers of genetic groups: (A) $Q = 2$, (B) $Q = 3$, (C) $Q = 4$, and (D) $Q = 5$.



Supplementary Figure S2. Mean group proportions in the training set (TRS) according to TRS size obtained using different optimization methods (over 30 random test sets (TSs)): PEVmean, CDmean and PAM, along with random sampling as a benchmark (x 30 replicates for each TS)



Supplementary Figure S3. Plots of genomic estimated breeding values (GEBVs) against LS-means obtained by leave-one-out (LOO) cross-validation using GBLUP for (A) Heading (in growing degree-days (GDD)), (B) Height (in cm), (C) Biomass (in t/ha) and (D) Yield (in t/ha). Each dot represents one individual and was colored according to its genetic group. Big dots represent the mean of GEBVs and LS-means for each group.



Supplementary Figure S4. Mean group proportions in the training set (TRS) according to TRS size obtained using different optimization methods (over 30 random test sets (TSs)): CDmean and PAM, along with random sampling with equal group contributions in the TRS as a benchmark (x 30 replicates for each TS). Each candidate set consisted in a of 450 individuals, including 150 of each of the three genetic groups.